
AN INTRODUCTION TO DATA ANALYSIS
THROUGH
A GEOMETRIC LENS

JEFF M. PHILLIPS

2017

Preface

This book is meant for use with a self-contained course that introduces many basic principals and techniques needed for modern data analysis. In particular, this course was designed as preparation for students planning to take rigorous Machine Learning and Data Mining courses. With this goal in mind, it both introduces key conceptual tools which are often helpful to see multiple times, and the most basic techniques to begin a basic familiarity with the backbone of modern data analysis.

Interaction with other courses. It is recommended that students taking this class have calculus and a familiarity with programming and algorithms. They should have also taken some probability and/or linear algebra; but we also review key concepts in these areas, so as to keep the book more self-contained. Thus, it may be appropriate for students to take these classes simultaneously. If appropriately planned for, it is the hope that this course could be taken at the sophomore level so that more rigorous and advanced data analysis classes can already be taken during the junior year.

Although we touch on Bayesian Inference, we do not cover most of classical statistics; neither frequentist hypothesis testing or the similar Bayesian perspectives. Most universities have well-developed classes on these topics which while also very useful, provide a complimentary view of data analysis. Classical statistical modeling approaches are often essential when a practitioner needs to provide some modeling assumptions to harness maximum power from limited data. But in the era of big data this is not always necessary. Rather, the topics in this course provide tools for using some of the data to help choose the model.

Scope and topics. Vital concepts introduced include concentration of measure and PAC bounds, cross-validation, gradient descent, and principal component analysis. These ideas are essential for modern data analysis, but not often taught in other introductory mathematics classes in a computer science or math department. Or if these concepts are taught, they are also presented in a very different context.

We also survey basic techniques in supervised (regression and classification) and unsupervised (principal component analysis and clustering) learning. We make an effort to keep the presentation and concepts on these topics simple. We stick to those which attempt to minimize sum of squared errors, and do not go much into regularization. We stick to classic but magical algorithms like Lloyd's algorithm for k -means, the power method for eigenvectors, and perceptron for linear classification. For many students (especially those in a computer science program), these are the first iterative, non-discrete algorithms they will have encountered.

On data. While this course is mainly focused on a mathematical preparation, what would data analysis be without data? As such we provide discussion on how to use these tools and techniques on actual data, with examples given in python. We choose python since it has many powerful libraries with extremely efficient backends in C. So for most data sets, this provides the proper interface for working with these tools.

But arguably more important than writing the code itself is a discussion on when and when-not to use techniques from the immense toolbox available. This is one of the main ongoing questions a data scientist must ask. And so, the text attempts to introduce the readers to this ongoing discussion.

Jeff M. Phillips
Salt Lake City, December 2016

1 Probability Review

Probability is a critical tool for modern data analysis. It arises in dealing with uncertainty, in randomized algorithms, and in Bayesian analysis. To understand any of these concepts correctly, it is paramount to have a solid and rigorous statistical foundation. Here we review some key definitions.

1.1 Sample Spaces

We define probability through set theory, starting with a *sample space* Ω . This represents the space of all things that might happen in the setting we consider. One such potential outcome $\omega \in \Omega$ is a *sample outcome*, it is an element of the space Ω . We are usually interested in an *event* that is a subset $A \subseteq \Omega$ of the sample space.

Example: Discrete Sample Space

Consider rolling a single fair, 6-sided die. Then $\Omega = \{1, 2, 3, 4, 5, 6\}$. One roll may produce an outcome $\omega = 3$, rolling a 3. An event might be $A = \{1, 3, 5\}$, any odd number. The probability of rolling an odd number is then $\Pr(A) = |\{1, 3, 5\}|/|\{1, 2, 3, 4, 5, 6\}| = 1/2$.

A *random variable* $X : \Omega \rightarrow S$ is a **function** from the sample space Ω to a domain S . Many times $S \subseteq \mathbb{R}$, where \mathbb{R} is the space of real numbers.

Example: Random Variable

Consider flipping a fair coin with $\Omega = \{H, T\}$. If I get a head H , then I get 1 point, and if I get a T , then I get 4 points. This describes the random variable X , defined $X(H) = 1$ and $X(T) = 4$.

The *probability* of an event $\Pr(A)$ satisfies the following properties:

- $0 \leq \Pr(A) \leq 1$ for any A ,
- $\Pr(\Omega) = 1$, and
- The probability of the union of disjoint events is equivalent to the sum of their individual probabilities. Formally, for any sequence A_1, A_2, \dots where for all $i \neq j$ that $A_i \cap A_j = \emptyset$, then

$$\Pr\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \Pr(A_i).$$

Sample spaces Ω can also be continuous, representing some quantity like water, time, or land mass which does not have discrete quantities. All of the above definitions hold for this setting.

Example: Continuous Sample Space

Assume you are riding a Swiss train that is always on time, but its departure is only specified to the minute (specifically, 1:37 pm). The true departure is then in the state space $\Omega = [1:37:00, 1:38:00)$. A continuous event may be $A = [1:37:00 - 1:37:40)$, the first 40 seconds of that minute. Perhaps the train operators are risk averse, so $\Pr(A) = 0.80$. That indicates that 0.8 fraction of trains depart in the first 2/3 of that minute (less than the 0.666 expected from a uniform distribution).

1.2 Conditional Probability and Independence

Now consider two events A and B . The *conditional probability* of A given B is written $\Pr(A | B)$, and can be interpreted as the probability of A , restricted to the setting where we know B is true. It is defined in simpler terms as $\Pr(A | B) = \frac{\Pr(A \cap B)}{\Pr(B)}$, that is the probability A and B are both true, divided by (normalized by) the probability B is true.

Two **events** A and B are *independent* of each other if and only if

$$\Pr(A | B) = \Pr(A).$$

Equivalently they are independent if and only if $\Pr(B | A) = \Pr(B)$ or $\Pr(A \cap B) = \Pr(A)\Pr(B)$. By algebraic manipulation, it is not hard to see these are all equivalent properties. This implies that knowledge about B has no effect on the probability of A (and vice versa from A to B).

Example: Conditional Probability

Consider the two random variables. T is 1 if a test for cancer is positive, and 0 otherwise. Variable C is 1 if a patient has cancer, and 0 otherwise. The joint probability of the events is captured in the following table:

	$C = 1$	$C = 0$
$T = 1$	0.1	0.02
$T = 0$	0.05	0.83

Note that the sum of all cells (the joint sample space Ω) is 1. The conditional probability of having cancer, given a positive test is $\Pr(C = 1 | T = 1) = \frac{0.1}{0.1+0.02} = 0.8333$. The probability of cancer (ignoring the test) is $\Pr(C = 1) = 0.1 + 0.05 = 0.15$. Since $\Pr(C = 1 | T = 1) \neq \Pr(C = 1)$, then events $T = 1$ and $C = 1$ are not independent.

Two **random variables** X and Y are *independent* if and only if, for *all* possible events $A \subseteq \Omega_X$ and $B \subseteq \Omega_Y$ that A and B are independent: $\Pr(A \cap B) = \Pr(A)\Pr(B)$.

1.3 Density Functions

Discrete random variables can often be defined through tables (as in the above cancer example). Or we can define a function $f_X(k)$ as the probability that random variable X is equal to k . For continuous random variables we need to be more careful; we use calculus. We will next develop probability density functions and cumulative density functions for continuous random variables; the same constructions are sometimes useful for discrete random variables as well, which basically just replace a integral with a sum.

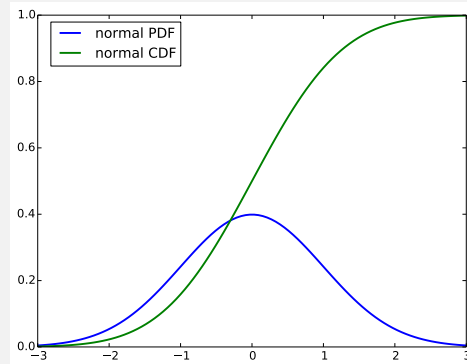
We consider a continuous sample space Ω , and a random variable X defined on that sample space. The probability density function of a random variable X is written f_X . It is defined with respect to any event A so that $\Pr(X \in A) = \int_{\omega \in A} f_X(\omega) d\omega$. The value $f_X(\omega)$ is *not equal to* $\Pr(X = \omega)$ in general, since for continuous functions $\Pr(X = \omega) = 0$ for any single value $\omega \in \Omega$. Yet, we can interpret f_X as a *likelihood* function; its value has no units, but they can be compared and larger ones are more likely.

Next we will defined the *cumulative density function* $F_X(t)$; it is the probability that X takes on a value of t or smaller. Here it is typical to have $\Omega = \mathbb{R}$, the set of real numbers. Now define $F_X(t) = \int_{\omega=-\infty}^t f_X(\omega) d\omega$.

We can also define a pdf in terms of a cdf as $f_X(\omega) = \frac{dF_X(\omega)}{d\omega}$.

Example: Normal Random Variable

A *normal* random variable X is a very common distribution to model noise. It has domain $\Omega = \mathbb{R}$. Its pdf is defined $f_X(\omega) = \frac{1}{\sqrt{2\pi}} \exp(-\omega^2/2) = \frac{1}{\sqrt{2\pi}} e^{-\omega^2/2}$, and its cdf has no closed form solution. We have plotted the cdf and pdf in the range $[-3, 3]$ where most of the mass lies:



```
import matplotlib as mpl
mpl.use('PDF')
import matplotlib.pyplot as plt
from scipy.stats import norm
import numpy as np
import math

mu = 0
variance = 1
sigma = math.sqrt(variance)
x = np.linspace(-3, 3, 201)

plt.plot(x, norm.pdf((x-mu)/sigma), linewidth=2.0, label='normal_PDF')
plt.plot(x, norm.cdf((x-mu)/sigma), linewidth=2.0, label='normal_CDF')
plt.legend(bbox_to_anchor=(.35, 1))

plt.savefig('Gaussian.pdf', bbox_inches='tight')
```

1.4 Expected Value

The expected value of a random variable X in a domain Ω is a very important constant, basically a weighted average of Ω , weighted by the range of X . For a discrete random variable X it is defined

$$\mathbf{E}[X] = \sum_{\omega \in \Omega} \omega \cdot \mathbf{Pr}[X = \omega].$$

For a continuous random variable X it is defined

$$\mathbf{E}[X] = \int_{\omega \in \Omega} \omega f_X(\omega) d\omega.$$

Linearity of Expectation: An important property of expectation is that it is a linear operation. That means for two random variables X and Y we have $\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y]$. For a scalar value α , we also $\mathbf{E}[\alpha X] = \alpha \mathbf{E}[X]$.

Example: Expectation

Let H be the random variable of the height of a man in meters without shoes. Let the pdf f_H of H be a normal distribution with expected value $\mu = 1.755\text{m}$ and with standard deviation 0.1m . Let S be the random variable of the height added by wearing a pair of shoes in centimeters (1 meter is 100 centimeters), its pdf is given by the following table:

$S = 1$	$S = 2$	$S = 3$	$S = 4$
0.1	0.1	0.5	0.3

Then the expected height of someone wearing shoes in centimeters is

$$\mathbf{E}[100 \cdot H + S] = 100 \cdot \mathbf{E}[H] + \mathbf{E}[S] = 100 \cdot 1.755 + (0.1 \cdot 1 + 0.1 \cdot 2 + 0.5 \cdot 3 + 0.3 \cdot 4) = 175.5 + 3 = 178.5$$

Note how the linearity of expectation allowed us to decompose the expression $100 \cdot H + S$ into its components, and take the expectation of each one individually. This trick is immensely powerful when analyzing complex scenarios with many factors.

1.5 Variance

The *variance* of a random variable X describes how spread out it is. It is defined

$$\mathbf{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2] - \mathbf{E}[X]^2.$$

The equivalence of those two common forms above uses that $\mathbf{E}[X]$ is a fixed scalar:

$$\mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2 - 2X\mathbf{E}[X] + \mathbf{E}[X]^2] = \mathbf{E}[X^2] - 2\mathbf{E}[X]\mathbf{E}[X] + \mathbf{E}[X]^2 = \mathbf{E}[X^2] - \mathbf{E}[X]^2.$$

For any scalar $\alpha \in \mathbb{R}$, then $\mathbf{Var}[\alpha X] = \alpha^2 \mathbf{Var}[X]$.

Note that the variance does not have the same units as the random variable or the expectation, it is that unit squared. As such, we also often discuss the *standard deviation* $\sigma_X = \sqrt{\mathbf{Var}[X]}$.

Example: Variance

Consider again the random variable S for height added by a shoe:

$S = 1$	$S = 2$	$S = 3$	$S = 4$
0.1	0.1	0.5	0.3

Its expected value is $\mathbf{E}[S] = 3$ (a fixed scalar), and its variance is

$$\begin{aligned} \mathbf{Var}[S] &= 0.1 \cdot (1 - 3)^2 + 0.1 \cdot (2 - 3)^2 + 0.5 \cdot (3 - 3)^2 + 0.3 \cdot (4 - 3)^2 \\ &= 0.1 \cdot (-2)^2 + 0.1 \cdot (-1)^2 + 0 + 0.3(1)^2 = 0.4 + 0.1 + 0.3 = 0.8. \end{aligned}$$

Then the standard deviation is $\sigma_S = \sqrt{0.8} \approx 0.894$.

The *covariance* of two random variables X and Y is defined $\mathbf{Cov}[X, Y] = \mathbf{E}[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])]$. It measures how much these random variables vary in accordance with each other; that is, if both are consistently away from the mean at the same time (in the same direction), then the covariance is high.

1.6 Joint, Marginal, and Conditional Distributions

We now extend some of these concepts to more than one random variable. Consider two random variables X and Y . Their *joint pdf* is defined $f_{X,Y} : \Omega_X \times \Omega_Y \rightarrow [0, \infty]$ where for discrete random variables this is defined by the probability $f_{X,Y}(x, y) = \mathbf{Pr}(X = x, Y = y)$. In this case, the domain of $f_{X,Y}$ is restricted so $f_{X,Y} \in [0, 1]$ and so $\sum_{x,y \in X \times Y} f_{X,Y}(x, y) = 1$.

Similarly, when $\Omega_X = \Omega_Y = \mathbb{R}$, the *joint cdf* is defined $F_{X,Y}(x, y) = \mathbf{Pr}(X \leq x, Y \leq y)$. The *marginal cumulative distribution functions* of $F_{X,Y}$ are defined as $F_X(x) = \lim_{y \rightarrow \infty} F_{X,Y}(x, y)dy$ and $F_Y(y) = \lim_{x \rightarrow \infty} F_{X,Y}(x, y)dx$.

Similarly, when Y is discrete, the *marginal pdf* is defined $f_X(x) = \sum_{y \in \Omega_Y} f_{X,Y}(x, y) = \sum_{y \in \Omega_Y} \mathbf{Pr}(X = x, Y = y)$. When the random variables are continuous, we define $f_{X,Y}(x, y) = \frac{d^2 F_{X,Y}(x, y)}{dx dy}$. And then the marginal pdf of X (when $\Omega_Y = \mathbb{R}$) is defined $f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y)dy$. Marginalizing removes the effect of a random variable (Y in the above definitions).

Now we can say random variables X and Y are independent if and only if $f_{X,Y}(x, y) = f_X(x) \cdot f_Y(y)$ for all x and y .

Then a *conditional distribution* of X given $Y = y$ is defined $f_{X|Y}(x | y) = f_{X,Y}(x, y)/f_Y(y)$ (given that $f_Y(y) \neq 0$).

Example: Marginal Distributions

Consider someone who randomly chooses his pants and shirt every day (a friend of mine actually did this in college – all clothes were in a pile, clean or dirty). Let P be a random variable for the color of pants, and S a random variable for the color of the shirt. Their joint probability is described by this table:

	$S=\text{green}$	$S=\text{red}$	$S=\text{blue}$
$P=\text{blue}$	0.3	0.1	0.2
$P=\text{white}$	0.05	0.2	0.15

Adding up along columns, the marginal distribution f_P for the color of the shirt is described by the following table:

$S=\text{green}$	$S=\text{red}$	$S=\text{blue}$
0.35	0.3	0.35

Isolating and renormalizing the middle “ $S=\text{red}$ ” column, the conditional distribution $f_{P|S}(\cdot | S=\text{red})$ is described by the following table:

$P=\text{blue}$	$P=\text{white}$
$\frac{0.1}{0.3} = 0.3333$	$\frac{0.2}{0.3} = 0.6666$

Example: Gaussian Distribution

The Gaussian distribution is a d -variate distribution $N_d : \mathbb{R}^d \rightarrow \mathbb{R}$ that generalizes the one-dimensional normal distribution. The definition of the symmetric version (we will generalize to non-trivial covariance later on) depends on a mean $\mu \in \mathbb{R}^d$ and a variance σ^2 . For any vector \mathbb{R}^d , it is defined

$$N_d(v) = \frac{1}{\sigma^d \sqrt{2\pi^d}} \exp(-\|v - \mu\|^2 / \sigma^2).$$

For the 2-dimensional case where $v = (v_x, v_y)$ and $\mu = (\mu_x, \mu_y)$, then this is defined

$$N_2(v) = \frac{1}{\sigma^2 \pi \sqrt{2}} \exp(-((v_x - \mu_x)^2 - (v_y - \mu_y)^2) / \sigma^2).$$

A *magical* property about the Gaussian distribution is that all conditional versions of it are also Gaussian, of a lower dimension. For instance, in the two dimensional case $N_2(v_x \mid v_y = 1)$ is a 1-dimensional Gaussian, or a normal distribution. There are many other essential properties of the Gaussian that we will see throughout this text, including that it is invariant under all basis transformations and that it is the limiting distribution for central limit theorem bounds.

2 Bayes' Rule

This topic is on Bayes' Rule and Bayesian Reasoning. Bayes' Rule is the key component in how to build likelihood functions, which are key to evaluating models based on data. Bayesian Reasoning is surprisingly different, much more about modeling uncertainty.

2.1 Bayes' Rule

Given two events M and D , then Bayes' Rule states

$$\Pr(M | D) = \frac{\Pr(D | M) \cdot \Pr(M)}{\Pr(D)}.$$

This assumes nothing about the independence of M and D (otherwise its pretty uninteresting). To derive this we use

$$\Pr(M \cap D) = \Pr(M | D)\Pr(D)$$

and also

$$\Pr(M \cap D) = \Pr(D \cap M) = \Pr(D | M)\Pr(M),$$

combined to get $\Pr(M | D)\Pr(D) = \Pr(D | M)\Pr(M)$, from which we can solve for $\Pr(M | D)$.

Example: Checking Bayes' Rule

Consider two events M and D with the following joint probability table:

	$M = 1$	$M = 0$
$D = 1$	0.25	0.5
$D = 0$	0.2	0.05

We can observe that indeed $\Pr(M | D) = \Pr(M \cap D) / \Pr(D) = \frac{0.25}{0.75} = \frac{1}{3}$, which is equal to

$$\frac{\Pr(D | M)\Pr(M)}{\Pr(D)} = \frac{.25(.2 + .25)}{.25 + .5} = \frac{.25}{.75} = \frac{1}{3}.$$

But Bayes' rule is not very interesting in the above example. In that example, it is actually *more* complicated to calculate the right side of Bayes' rule than it is the left side.

Example: Cracked Windshield

Consider you bought a new car and its windshield was cracked, the event W . If the car was assembled at one of three factories A , B or C , you would like to know which factory was the most likely point of origin.

Assume that in Utah 50% of cars are from factory A (that is $\Pr(A) = 0.5$) and 30% are from factory B ($\Pr(B) = 0.3$), and 20% are from factory C ($\Pr(C) = 0.2$).

Then you look up statistics online, and find the following rates of cracked windshields for each factory – apparently this is a problem! In factory A , only 1% are cracked, in factory B 10% are cracked, and in factory C 2% are cracked. That is $\Pr(W | A) = 0.01$, $\Pr(W | B) = 0.1$ and $\Pr(W | C) = 0.02$.

We can now calculate the probability the car came from each factory:

- $\Pr(A | W) = \Pr(W | A) \cdot \Pr(A) / \Pr(W) = 0.01 \cdot 0.5 / \Pr(W) = 0.005 / \Pr(W)$.
- $\Pr(B | W) = \Pr(W | B) \cdot \Pr(B) / \Pr(W) = 0.1 \cdot 0.3 / \Pr(W) = 0.03 / \Pr(W)$.
- $\Pr(C | W) = \Pr(W | C) \cdot \Pr(C) / \Pr(W) = 0.02 \cdot 0.2 / \Pr(W) = 0.004 / \Pr(W)$.

We did not calculate $\Pr(W)$, but it must be the same for all factory events, so to find the highest probability factory we can ignore it. The probability $\Pr(B | W) = 0.03 / \Pr(W)$ is the largest, and B is the most likely factory.

2.1.1 Model Given Data

In data analysis, M represents a ‘model’ and D as ‘data.’ Then $\Pr(M | D)$ is interpreted as the probability of model M given that we have observed D . A *maximum a posteriori* (or MAP) estimate is the model $M \in \Omega_M$ that maximizes $\Pr(M | D)$. That is

$$M^* = \arg \max_{M \in \Omega_M} \Pr(M | D) = \arg \max_{M \in \Omega_M} \frac{\Pr(D | M) \Pr(M)}{\Pr(D)} = \arg \max_{M \in \Omega_M} \Pr(D | M) \Pr(M).$$

Thus, to use Bayes’ Rule, we can maximize $\Pr(M | D)$ using $\Pr(M)$ and $\Pr(D | M)$. We do not need $\Pr(D)$ since our data is given to us and fixed for all models.

In some settings we may also ignore $\Pr(M)$, as we may assume all possible models are equally likely. This is not always the case, and we’ll come back to this. Thus we just need to calculate $\Pr(D | M)$. Then, in this setting $L(M) = \Pr(D | M)$ is called the *likelihood* of model M .

So what is a ‘model’ and what is ‘data?’ A *model* is usually a simple pattern which we think data is generated from, but then observed with some noise. Examples:

- The model M is a single point in \mathbb{R}^d ; the data is a set of points in \mathbb{R}^d near M .
- **linear regression:** The model M is a line in \mathbb{R}^2 ; the data is a set of points such that for each x -coordinate, the y -coordinate is the value of the line at that x -coordinate with some added noise in the y -value.
- **clustering:** The model M is a small set of points in \mathbb{R}^d ; the data is a large set of points in \mathbb{R}^d , where each point is near one of the points in M .
- **PCA:** The model M is a k -dimensional subspace in \mathbb{R}^d (for $k \ll d$); the data is a set of points in \mathbb{R}^d , where each point is near M .
- **linear classification:** The model M is a halfspace in \mathbb{R}^d ; the data is a set of labeled points (with labels + or –), so the + points are mostly in M , and the – points are mainly not in M .

Example: Gaussian MLE

Let the data D be a set of points in \mathbb{R}^1 : $\{1, 3, 12, 5, 9\}$. Let Ω_M be \mathbb{R} so that the model is a point $M \in \mathbb{R}$. If we assume that each data point is observed with independent Gaussian noise (with $\sigma = 2$, so its pdf is described as $g(x) = \frac{1}{\sqrt{8\pi}} \exp(-\frac{1}{8}(M - x)^2)$). Then

$$\Pr(D | M) = \prod_{x \in D} g(x) = \prod_{x \in D} \left(\frac{1}{\sqrt{8\pi}} \exp(-\frac{1}{8}(M - x)^2) \right).$$

Recall that we can take the product $\prod_{x \in D} g(x)$ since we assume independence of $x \in D$! To find $M^* = \arg \max_M \Pr(D | M)$ is equivalent to $\arg \max_M \ln(\Pr(D | M))$, the *log-likelihood* which is

$$\ln(\Pr(D | M)) = \ln \left(\prod_{x \in D} \left(\frac{1}{\sqrt{8\pi}} \exp(-\frac{1}{8}(M - x)^2) \right) \right) = \sum_{x \in D} \left(-\frac{1}{8}(M - x)^2 \right) + |D| \ln \left(\frac{1}{\sqrt{8\pi}} \right).$$

We can ignore the last term since it is independent of M . The first term is maximized when $\sum_{x \in D} (M - x)^2$ is minimized, which occurs precisely as $\mathbf{E}[D] = \frac{1}{|D|} \sum_{x \in D} x$, the mean of the data set D . That is, the maximum likelihood model is exactly the mean of the data D , and is quite easy to calculate.

2.2 Bayesian Inference

Bayesian inference focuses on a simplified version of Bayes's Rule:

$$\Pr(M | D) \propto \Pr(D | M) \cdot \Pr(M).$$

The symbol \propto means *proportional to*; that is there is a fixed (but possibly unknown) constant factor c multiplied on the right (in this case $c = 1/\Pr(D)$) to make them equal: $\Pr(M | D) = c \cdot \Pr(D | M) \cdot \Pr(M)$.

However, we may want to use continuous random variables, so then strictly using probability \Pr at a single point is not always correct. So we can replace each of these with pdfs

$$p(M | D) \propto f(D | M) \cdot \pi(M).$$

Each of these terms have common names. As above, the conditional probability or pdf $\Pr(D | M) \propto f(D | M)$ is called the *likelihood*. The probability or pdf of the model $\Pr(M) \propto \pi(M)$ is called the *prior*. And the left hand side $\Pr(M | D) \propto p(M | D)$ is called the *posterior*.

Again it is common to be in a situation where, given a fixed model M , it is possible to calculate the likelihood $f(D | M)$. And again, the goal is to be able to compute $p(M | D)$, as this allows us to evaluate potential models M , given the data we have seen D .

The main difference is a careful analysis of $\pi(M)$, the prior – which is not necessarily assumed uniform or “flat”. The prior allows us to encode our assumptions.

Example: Average Height

Lets estimate the height H of a typical U of U student. We can construct a data set $D = \{x_1, \dots, x_n\}$ by measuring the height of everyone in this class in inches. There may be error in the measurement, and we are an incomplete set, so we don't entirely trust the data.

So we introduce a prior $\pi(M)$. Consider we read that the average height of an full grown person is $\mu_M = 66$ inches, with a standard deviation of $\sigma = 6$ inches. So we assume

$$\pi(M) = N(66, 6) = \frac{1}{\sqrt{\pi 72}} \exp(-(\mu_M - 66)^2 / (2 \cdot 6^2)),$$

is normally distributed around 66 inches.

Now, given this knowledge we adjust the MLE example from last subsection using this prior.

- *What if our MLE estimate without the prior (e.g. $\frac{1}{|D|} \sum_{x \in D} x$) provides a value of 5.5?*
That means the data is very far from the prior. Usually this means something is wrong. We could find $\arg \max_M p(M | D)$ using this information, but that may give us an estimate of say 20 (that does not seem correct). A more likely explanation is a mistake somewhere: probably we measured in feet instead of inches!

Another vestige of Bayesian inference is that we not only can calculate the maximum likelihood model M^* , but we can also provide a posterior value for any model! This value is not an absolute probability (its not normalized, and regardless it may be of measure 0), but it is powerful in other ways:

- We can say (under our model assumptions, which are now clearly stated) that one model M_1 is twice as likely as another M_2 , if $p(M_1 | D) / p(M_2 | D) = 2$.
- We can define a range of parameter values (with more work and under our model assumptions) that likely contains the true model.
- We can now use more than one model for prediction of a value. Given a new data point x' we may want to map it onto our model as $M(x')$, or assign it a score of fit. Instead of doing this for just one "best" model M^* , we can take a weighted average of all models, weighted by their posterior; this is "marginalization."

Weight for Prior. So how important is the prior? In the average height example, it will turn out to be worth only (1/9)th of one student's measurement. But we can give it more weight.

Example: Weighted Prior for Height

Lets continue the example about the height of an average U of U student, and assume (as in the MLE estimator example) the data is generated independently from a model M with Gaussian noise with $\sigma = 2$. Thus the likelihood of the model, given the data is

$$f(D | M) = \prod_{x \in D} g(x) = \prod_{x \in D} \left(\frac{1}{\sqrt{8\pi}} \exp\left(-\frac{1}{8}(\mu_M - x)^2\right) \right).$$

Now using that the prior of the model is $\pi(M) = \frac{1}{\sqrt{\pi 72}} \exp(-(\mu_M - 66)^2/72)$, the posterior is given by

$$p(M | D) \propto f(D | M) \cdot \frac{1}{\sqrt{\pi 72}} \exp(-(\mu_M - 66)^2/72).$$

It is again easier to work with the log-posterior which is monotonic with the posterior, using some unspecified constant C (which can be effectively ignored):

$$\begin{aligned} \ln(p(M | D)) &\propto \ln(f(D | M)) + \ln(\pi(M)) + C \\ &\propto \sum_{x \in D} \left(-\frac{1}{8}(\mu_M - x)^2 \right) - \frac{1}{72}(\mu_M - 66)^2 + C \\ &\propto - \sum_{x \in D} 9(\mu_M - x)^2 + (\mu_M - 66)^2 + C \end{aligned}$$

So the maximum likelihood estimator occurs at the average of 66 along with 9 copies of the student data.

Why is student measurement data worth so much more?

We assume the standard deviation of the measurement error is 2, where as we assumed that the standard deviation of the full population was 6. In other words, our measurements had variance $2^2 = 4$, and the population had variance $6^2 = 36$ (technically, this is best to interpret at the variance when adapted to various subpopulations, e.g., U of U students): that is 9 times as much.

If instead we assumed that the standard deviation of our prior is 0.1, with variance 0.01, then this is 400 times smaller than our class measurement error variance. If we were to redo the above calculations with this smaller variance, we would find that this assumption weights the prior 400 times the effect of each student measurement in the MLE.

So what happens with more data?

Lets say, this class gets really popular, and next year 1000 students sign up! Then again the student data is overall worth more than the prior data. So with any prior, if we get enough data, it no longer becomes important. But with a small amount of data, it can have a large influence on our model.

3 Convergence

This topic will overview a variety of extremely powerful analysis results that span statistics, estimation theorem, and big data. It provides a framework to think about how to aggregate more and more data to get better and better estimates. It will cover the *Central Limit Theorem* (CLT), Chernoff-Hoeffding bounds, and Probably Approximately Correct (PAC) algorithms.

3.1 Sampling and Estimation

Most data analysis starts with some data set; we will call this data set P . It will be composed of a set of n data points $P = \{p_1, p_2, \dots, p_n\}$.

But underlying this data is almost always a very powerful assumption, that this data comes iid from a fixed, but usually unknown pdf, call this f . Lets unpack this: What does “iid” mean: Identically and Independently Distributed. The “identically” means each data point was drawn from the same f . The “independently” means that the first points have no bearing on the value of the next point.

Example: Polling

Consider a poll of $n = 1000$ likely voters in an upcoming election. If we assume each polled person is chosen iid, then we can use this to understand something about the underlying distribution f , for instance the distribution of all likely voters.

More generally, f could represent the outcome of a process, whether that is a randomized algorithm, a noisy sensing methodology, or the common behavior of a species of animals. In each of these cases, we essentially “poll” the process (algorithm, measurement, thorough observation) having it provide a sample, and repeat many times over.

Here we will talk about estimating the mean of f . To discuss this, we will now introduce a random variable $X \sim f$; a hypothetical new data point. The *mean* of f is the expected value of X : $\mathbf{E}[X]$.

We will estimate the mean of f using the *sample mean*, defined $\bar{P} = \frac{1}{n} \sum_{i=1}^n p_i$. The following diagram represents this common process: from a unknown process f , we consider n iid random variables $\{X_i\}$ corresponding to a set of n independent observations $\{p_i\}$, and take their average $\bar{P} = \frac{1}{n} \sum_{i=1}^n p_i$ to estimate the mean of f .

$$\bar{P} = \frac{1}{n} \sum \{p_i\} \xleftarrow{\text{realize}} \{X_i\} \underset{\text{iid}}{\sim} f$$

Central Limit Theorem. The central limit theorem is about how well the sample mean approximates the true mean. But to discuss the sample mean \bar{P} (which is a fixed value) we need to discuss random variables $\{X_1, X_2, \dots, X_n\}$, and their mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Note that again \bar{X} is a random variable. If we are to draw a new iid data set P' and calculate a new sample mean \bar{P}' it will likely not be exactly the same as \bar{P} ; however, the distribution of where this \bar{P}' is likely to be, is precisely \bar{X} . Arguably, this distribution is more important than \bar{P} itself.

There are many formal variants of the central limit theorem, but the basic form is as follows:

Central Limit Theorem: Consider n iid random variables X_1, X_2, \dots, X_n , where each $X_i \sim f$ for any fixed distribution f with mean μ and bounded variance σ^2 . Then $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ converges to the normal distribution with mean $\mu = \mathbf{E}[X_i]$ and variance σ^2/n .

Lets highlight some important consequences:

- For any f (that is not too crazy, since σ^2 is not infinite), then \bar{X} looks like a normal distribution.
- The mean of the normal distribution, which is the expected value of \bar{X} satisfies $\mathbf{E}[\bar{X}] = \mu$, the mean of f . This implies that we can use \bar{X} (and then also \bar{P}) as a guess for μ .
- As n gets larger (we have more data points) then the variance of \bar{X} (our estimator) decreases. So keeping in mind that although \bar{X} has the right expected value it also has some error, this error is decreasing as n increases.

```
# adapted from: https://github.com/mattnedrich/CentralLimitTheoremDemo
import random
import matplotlib as mpl
mpl.use('PDF')
import matplotlib.pyplot as plt

def plot_distribution(distribution, file, title, bin_min, bin_max, num_bins):
    bin_size = (bin_max - bin_min) / num_bins
    manual_bins = range(bin_min, bin_max + bin_size, bin_size)
    [n, bins, patches] = plt.hist(distribution, bins = manual_bins)
    plt.title(title)
    plt.xlim(bin_min, bin_max)
    plt.ylim(0, max(n) + 2)
    plt.ylabel("Frequency")
    plt.xlabel("Observation")
    plt.savefig(file, bbox_inches='tight')
    plt.clf()
    plt.cla()

minbin = 0
maxbin = 100
numbins = 50
nTrials = 1000

def create_uniform_sample_distribution():
    return range(maxbin)
sampleDistribution = create_uniform_sample_distribution()

# Plot the original population distribution
plot_distribution(sampleDistribution, 'output/SampleDistribution.pdf',
    "Population_Distribution", minbin, maxbin, numbins)

# Plot a sampling distribution for values of N = 2, 3, 10, and 30
n_vals = [2, 3, 10, 30]
for N in n_vals:
    means = []
    for j in range(nTrials):
        sampleSum = 0;
        for i in range(N):
            sampleSum += random.choice(sampleDistribution)
```



```

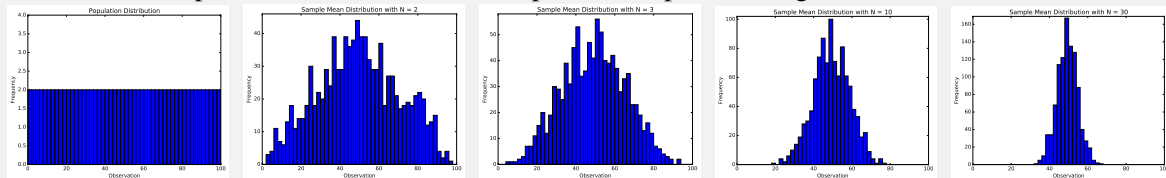
means.append(float(sampleSum) / float(N))

title = "Sample_Mean_Distribution_with_N=%s" % N
file = "output/CLT-demo-%s.pdf" % N
plot_distribution(means, file, title, minbin, maxbin, numbins)

```

Example: Central Limit Theorem

Consider f as a uniform distribution over $[0, 100]$. If we create n samples $\{p_1, \dots, p_n\}$ and their mean \bar{P} , then repeat this 1000 times, we can plot the output in histograms:



We see that starting at $n = 2$, the distributions look vaguely normal (in the technical sense of a normal distribution), and that their standard deviations narrow as n increases.

Remaining Mysteries. There should still be at least a few aspects of this not clear yet: (1) What does “convergence” mean? (2) How do we formalize or talk about this notion of error? (3) What does this say about our data \bar{P} ?

First, *convergence* refers to what happens as some parameter increases, in this case n . As the number of data points increase, as n “goes to infinity” then the above statement (\bar{X} looks like a normal distribution) becomes more and more true. For small n , the distribution may not quite look like a normal, it may be more bumpy, maybe even multi-modal. The statistical definitions of converge are varied, and we will not go into them here, we will instead replace it with more useful phrasing in explaining aspects (2) and (3).

Second, the error now has two components. We cannot simply say that \bar{P} is at most some distance ε from μ . Something crazy might have happened (the sample is random after all). And it is not useful to try to write the probability that $\bar{P} = \mu$; for equality in continuous distributions, this probability is indeed 0. But we can combine these notions. We can say the distance between \bar{P} and μ is more than ε , with probability at most δ . This is called “probably approximately correct” or PAC.

Third, we want to generate some sort of PAC bound (which is far more useful than “ \bar{X} looks kind of like a normal distribution”). Whereas a frequentist may be happy with a confidence interval and a Bayesian a normal posterior, these two options are not directly available since again, \bar{X} is not exactly a normal. So we will discuss some very common *concentration of measure* tools. These don’t exactly capture the shape of the normal distribution, but provide upper bounds for its tails, and will allow us to state PAC bounds.

3.2 Probably Approximately Correct (PAC)

We will introduce shortly the three most common concentration of measure bounds, which provide increasingly strong bounds on the tails of distributions, but require more and more information about the underlying distribution f . Each provides a PAC bound of the following form:

$$\Pr[|\bar{X} - \mathbf{E}[\bar{X}]| \geq \varepsilon] \leq \delta.$$

That is, the probability that \bar{X} (which is some random variable, often a sum of iid random variables) is further than ε to its expected value (which is μ , the expected value of f where $X_i \sim f$), is at most δ . Note we do not try to say this probability is *exactly* δ , this is often too hard. In practice there are a variety of tools, and a user may try each one, and see which ones gives the best bound.

It is useful to think of ε as the *error tolerance* and δ as the *probability of failure* i.e., that we exceed the error tolerance. However, often these bounds will allow us to write the required sample size n in terms of ε and δ . This allows us to trade these two terms off for any fixed known n ; we can allow a smaller error tolerance if we are willing to allow more probability of failure, and vice-versa.

3.3 Concentration of Measure

We will formally describe these bounds, and give some intuition of why they are true (but not proofs). But what will be the most important is what they imply. If you just know the distance of the expectation from the minimal value, you can get a very weak bound. If you know the variance of the data, you can get a stronger bound. If you know that the distribution f has a small and bounded range, then you can make the probability of failure (the δ in PAC bounds) very very small.

Markov Inequality. Let X be a random variable such that $X \geq 0$, that is it cannot take on negative values. Then for any parameter $\alpha > 0$

$$\Pr[X > \alpha] \leq \frac{\mathbf{E}[X]}{\alpha}.$$

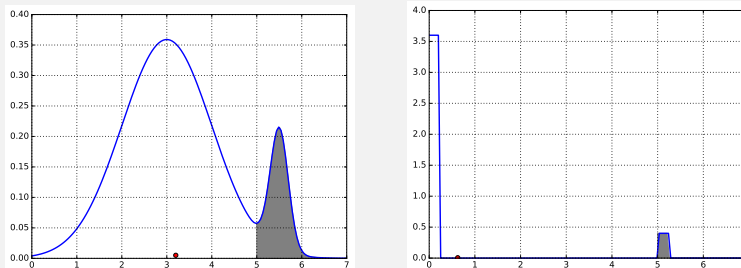
Note this is a PAC bound with $\varepsilon = \alpha - \mathbf{E}[X]$ and $\delta = \mathbf{E}[X]/\alpha$, or we can rephrase this bound as follows: $\Pr[X - \mathbf{E}[X] > \varepsilon] \leq \delta = \mathbf{E}[X]/(\varepsilon + \mathbf{E}[X])$.

This bound can be seen true by considering balancing the pdf of X at $\mathbf{E}[X]$ (think of a waiter holding a tray of food with weight of the food at various locations described by the probability distribution; the waiter's hand should be under the point $\mathbf{E}[X]$). Then, more than α fraction of its mass cannot be greater than $\mathbf{E}[X]/\alpha$ since the rest is at least 0, and it cannot balance.

If we instead know that $X \geq b$ for some constant b (instead of $X \geq 0$), then we state more generally $\Pr[X > \alpha] \leq (\mathbf{E}[X] - b)/(\alpha - b)$.

Example: Markov Inequality

Consider the pdf f drawn in blue in the following figures, with $\mathbf{E}[X]$ for $X \sim f$ marked as a red dot. The probability that X is greater than 5 (e.g. $\Pr[X \geq 5]$) is the shaded area.



Notice that in both cases that $\Pr[X \geq 5]$ is about 0.1. This is the quantity we want to bound by above by δ . But since $\mathbf{E}[X]$ is much larger in the first case (about 2.25), then the bound $\delta = \mathbf{E}[X]/\alpha$ is much larger, about 0.45. In the second case, $\mathbf{E}[X]$ is much smaller (about 0.6) so we get a much better bound of $\delta = 0.12$.

Chebyshev Inequality. Now let X be a random variable where we know $\mathbf{Var}[X]$, and $\mathbf{E}[X]$. Then for any parameter $\varepsilon > 0$

$$\Pr[|X - \mathbf{E}[X]| \geq \varepsilon] \leq \frac{\mathbf{Var}[X]}{\varepsilon^2}.$$

Again, this clearly is a PAC bound with $\delta = \mathbf{Var}[X]/\varepsilon^2$. This bound is typically stronger than the Markov one since δ decreases quadratically in ε instead of linearly.

Example: Chebyshev for IID Samples

Recall that for an average of random variables $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$, where the X_i s are iid, and have variance σ^2 , then $\mathbf{Var}[\bar{X}] = \sigma^2/n$. Hence

$$\Pr[|\bar{X} - \mathbf{E}[X_i]| \geq \varepsilon] \leq \frac{\sigma^2}{n\varepsilon^2}.$$

Consider now that we have input parameters ε and δ , our desired error tolerance and probability of failure. If can draw $X_i \sim f$ (iid) for an unknown f (with known expected value and variance σ), then we can solve for how large n needs to be: $n = \sigma^2/(\varepsilon^2\delta)$.

Chernoff-Hoeffding Inequality. Following the above example, we can consider a set of n iid random variables X_1, X_2, \dots, X_n where $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Now assume we know that each X_i lies in a bounded domain $[b, t]$, and let $\Delta = t - b$. Then for any parameter $\varepsilon > 0$

$$\Pr[|\bar{X} - \mathbf{E}[\bar{X}]| > \varepsilon] \leq 2 \exp\left(\frac{-2\varepsilon^2 n}{\Delta^2}\right).$$

Again this is a PAC bound, now with $\delta = 2 \exp(-2\varepsilon^2 n/\Delta^2)$. For a desired error tolerance ε and failure probability δ , we can set $n = (\Delta^2/(2\varepsilon^2)) \ln(2/\delta)$. Note that this has a similar relationship with ε as the Chebyshev bound, but the dependence of n on δ is exponentially less for this bound.

Relating this all back to the Gaussian distribution in the CLT, the Chebyshev bound only uses the variance information about the Gaussian, but the Chernoff-Hoeffding bound uses all of the “moments”: that it decays exponentially.

These are the most basic and common PAC concentration of measure bounds, but are by no means exhaustive.

Example: Uniform Distribution

Consider a random variable $X \sim f$ where $f(x) = \{\frac{1}{2} \text{ if } x \in [0, 2] \text{ and } 0 \text{ otherwise.}\}$, i.e, the Uniform distribution on $[0, 2]$. We know $\mathbf{E}[X] = 1$ and $\mathbf{Var}[X] = \frac{1}{3}$.

- Using the Markov Inequality, we can say $\mathbf{Pr}[X > 1.5] \leq 1/(1.5) \approx 0.6666$ and $\mathbf{Pr}[X > 3] \leq 1/3 \approx 0.33333$.
or $\mathbf{Pr}[X - \mu > 0.5] \leq \frac{2}{3}$ and $\mathbf{Pr}[X - \mu > 2] \leq \frac{1}{3}$.
- Using the Chebyshev Inequality, we can say that $\mathbf{Pr}[|X - \mu| > 0.5] \leq (1/3)/0.5^2 = \frac{4}{3}$ (which is meaningless). But $\mathbf{Pr}[|X - \mu| > 2] \leq (1/3)/(2^2) = \frac{1}{12} \approx 0.08333$.

Now consider a set of $n = 100$ random variables X_1, X_2, \dots, X_n all drawn iid from the same pdf f as above. Now we can examine the random variable $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. We know that $\mu_n = \mathbf{E}[\bar{X}] = \mu$ and that $\sigma_n^2 = \mathbf{Var}[\bar{X}] = \sigma^2/n = 1/(3n) = 1/300$.

- Using the Chebyshev Inequality, we can say that $\mathbf{Pr}[|\bar{X} - \mu| > 0.5] \leq \sigma_n^2/(0.5)^2 = \frac{1}{75} \approx 0.01333$, and $\mathbf{Pr}[|\bar{X} - \mu| > 2] \leq \sigma_n^2/2^2 = \frac{1}{1200} \approx 0.0008333$.
- Using the Chernoff-Hoeffding bound, we can say that $\mathbf{Pr}[|\bar{X} - \mu| > 0.5] \leq 2 \exp(-2(0.5)^2 n / \Delta^2) = 2 \exp(-100/8) \approx 0.0000074533$, and $\mathbf{Pr}[|\bar{X} - \mu| > 2] \leq 2 \exp(-2(2)^2 n / \Delta^2) = 2 \exp(-200) \approx 2.76 \cdot 10^{-87}$.

4 Linear Algebra Review

For this topic we quickly review many key aspects of linear algebra that will be necessary for the remainder of the text.

4.1 Vectors and Matrices

For the context of data analysis, the critical part of linear algebra deals with vectors and matrices of real numbers.

In this context, a *vector* $v = (v_1, v_2, \dots, v_d)$ is equivalent to a point in \mathbb{R}^d . By default a vector will be a column of d numbers (where d is context specific)

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

but in some cases we will assume the vector is a row

$$v^T = [v_1 \ v_2 \ \dots \ v_n].$$

An $n \times d$ matrix A is then an ordered set of n row vectors a_1, a_2, \dots, a_n

$$A = [a_1 \ a_2 \ \dots \ a_n] = \begin{bmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{bmatrix} = \begin{bmatrix} A_{1,1} & A_{1,2} & \dots & A_{1,d} \\ A_{2,1} & A_{2,2} & \dots & A_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n,1} & A_{n,2} & \dots & A_{n,d} \end{bmatrix},$$

where vector $a_i = [A_{i,1}, A_{i,2}, \dots, A_{i,d}]$, and $A_{i,j}$ is the element of the matrix in the i th row and j th column. We can write $A \in \mathbb{R}^{n \times d}$ when it is defined on the reals.

A *transpose* operation $(\cdot)^T$ reverses the roles of the rows and columns, as seen above with vector v . For a matrix, we can write:

$$A^T = \begin{bmatrix} | & | & & | \\ a_1 & a_2 & \dots & a_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} A_{1,1} & A_{2,1} & \dots & A_{n,1} \\ A_{1,2} & A_{2,2} & \dots & A_{n,2} \\ \vdots & \vdots & \ddots & \vdots \\ A_{1,n} & A_{2,d} & \dots & A_{n,d} \end{bmatrix}.$$

Example: Linear Equations

A simple place these objects arise is in linear equations. For instance

$$\begin{array}{rcl} 3x_1 & -7x_2 & +2x_3 = -2 \\ -1x_1 & +2x_2 & -5x_3 = 6 \end{array}$$

is a system of $n = 2$ linear equations, each with $d = 3$ variables. We can represent this system in matrix-vector notation as

$$Ax = b$$

where

$$b = \begin{bmatrix} -2 \\ 6 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} 3 & -7 & 2 \\ -1 & 2 & -5 \end{bmatrix}.$$

4.2 Addition

We can add together two vectors or two matrices only if they have the same dimensions. For vectors $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ and $y = (y_1, y_2, \dots, y_d) \in \mathbb{R}^d$, then vector

$$z = x + y = (x_1 + y_1, x_2 + y_2, \dots, x_d + y_d) \in \mathbb{R}^d.$$

Similarly for two matrices $A, B \in \mathbb{R}^{n \times d}$, then $C = A + B$ is defined where $C_{i,j} = A_{i,j} + B_{i,j}$ for all i, j .

4.3 Multiplication

Multiplication only requires alignment along one dimension. For two matrices $A \in \mathbb{R}^{n \times d}$ and $B \in \mathbb{R}^{d \times m}$ we can obtain a new matrix $C = AB \in \mathbb{R}^{n \times m}$ where $C_{i,j}$, the element in the i th row and j th column of C is defined

$$C_{i,j} = \sum_{k=1}^d A_{i,k} B_{k,j}.$$

To multiply A times B (where A is to the left of B , the order matters!) then we require the row dimension d of A to match the column dimension d of B . If $n \neq m$, then we *cannot* multiply BA . Keep in mind:

- Matrix multiplication is *associative* $(AB)C = A(BC)$.
- Matrix multiplication is *distributive* $A(B + C) = AB + AC$.
- Matrix multiplication is **not commutative** $AB \neq BA$.

We can also multiply a matrix A by a scalar α . In this setting $\alpha A = A\alpha$ and is defined by a new matrix B where $B_{i,j} = \alpha A_{i,j}$.

vector-vector products. There are two types of vector-vector products, and their definitions follow directly from that of matrix-matrix multiplication (since a vector is a matrix where one of the dimensions is 1). But it is worth highlighting these.

Given two column vectors $x, y \in \mathbb{R}^d$, the *inner product* or *dot product* is written

$$x^T y = x \cdot y = \langle x, y \rangle = [x_1 \ x_2 \ \dots \ x_d] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix} = \sum_{i=1}^d x_i y_i,$$

where x_i is the i th element of x and similar for y_i . This text will prefer the last notation $\langle x, y \rangle$ since the same can be used for row vectors, and there is no confusion with multiplication in using \cdot ; whether a vector is a row or a column is often arbitrary.

Note that this operation produces a single scalar value. The dot product is a linear operator. So this means for any scalar value α and three vectors $x, y, z \in \mathbb{R}^d$ we have

$$\langle \alpha x, y + z \rangle = \alpha \langle x, y + z \rangle = \alpha (\langle x, y \rangle + \langle x, z \rangle).$$

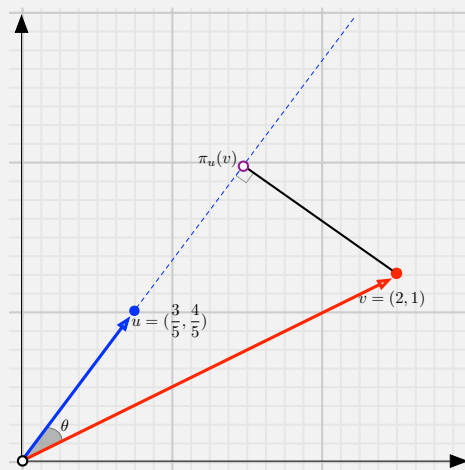
Example: Geometry of Dot Product

A dot product is one of my favorite mathematical operations! It encodes a lot of geometry. Consider two vectors $u = (\frac{3}{5}, \frac{4}{5})$ and $v = (2, 1)$, with an angle θ between them. Then it holds

$$\langle u, v \rangle = \text{length}(u) \cdot \text{length}(v) \cdot \cos(\theta).$$

Here $\text{length}(\cdot)$ measures the distance from the origin. We'll see how to measure length with a "norm" $\|\cdot\|$ soon.

Moreover, since the $\|u\| = \text{length}(u) = 1$, then we can also interpret $\langle u, v \rangle$ as the length of v projected onto the line through u . That is, let $\pi_u(v)$ be the closest point to v on the line through u (the line through u and the line segment from v to $\pi_u(v)$ make a right angle). Then $\langle u, v \rangle = \text{length}(\pi_u(v)) = \|\pi_u(v)\|$.



For two column vectors $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^d$, the *outer product* is written

$$y^T x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} [y_1 \ y_2 \ \dots \ y_d] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \dots & x_1 y_d \\ x_2 y_1 & x_2 y_2 & \dots & x_2 y_d \\ \vdots & \vdots & \ddots & \vdots \\ x_n y_1 & x_n y_2 & \dots & x_n y_d \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

Note that the result here is a matrix, not a scalar.

matrix-vector products. Another important and common operation is a matrix-vector product. Given a matrix $A \in \mathbb{R}^{n \times d}$ and a vector $x \in \mathbb{R}^d$, their product $y = Ax \in \mathbb{R}^n$.

When A is composed of row vectors $[a_1; a_2; \dots; a_n]$, then it is useful to imagine this as transposing x (which should be a column vector here, so a row vector after transposing), and taking the dot product with

each row of A .

$$y = Ax = \begin{bmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{bmatrix} x = \begin{bmatrix} \langle a_1, x \rangle \\ \langle a_2, x \rangle \\ \vdots \\ \langle a_n, x \rangle \end{bmatrix}.$$

4.4 Norms

The standard *Euclidean norm* of a vector $v = (v_1, v_2, \dots, v_d) \in \mathbb{R}^d$ is defined

$$\|v\| = \sqrt{\sum_{i=1}^d v_i^2} = \sqrt{\langle v, v \rangle}.$$

This measures the “straight-line” distance from the origin to the point at v . A vector v with norm $\|v\| = 1$ is said to be a *unit vector*; sometimes a vector x with $\|x\| = 1$ is said to be *normalized*.

However, a “norm” is a more generally concept. A class called L_p norms are well-defined for any parameter $p \in [1, \infty)$ as

$$\|v\|_p = \left(\sum_{i=1}^d |v_i|^p \right)^{1/p}.$$

Thus, when no p is specified, it is assumed to be $p = 2$. It is also common to denote $\|v\|_\infty = \max_{i=1}^d |v_i|$.

Because subtraction is well-defined between vectors $v, u \in \mathbb{R}^d$ of the same dimension, then we can also take the norm of $\|v - u\|_p$. While this is technically the norm of the vector resulting from the subtraction of u from v ; it also provides a distance between u and v . In the case of $p = 2$, then

$$\|u - v\|_2 = \sqrt{\sum_{i=1}^d (u_i - v_i)^2}$$

is precisely the straight-line (Euclidean) distance between u and v .

Moreover, all L_p norms define a distance $D_p(u, v) = \|u - v\|_p$, which satisfies a set of special properties required for a distance to be a *metric*. This include:

- **Symmetry:** For any $u, v \in \mathbb{R}^d$ we have $D(u, v) = D(v, u)$.
- **Non-negativity:** For any $u, v \in \mathbb{R}^d$ we have $D(u, v) \geq 0$, and $D(u, v) = 0$ if and only if $u = v$.
- **Triangle Inequality:** For any $u, v, w \in \mathbb{R}^d$ we have $D(u, w) + D(w, v) \geq D(u, v)$.

We can also define norms for matrices A . These take on slightly different notational conventions. The two most common are the spectral norm $\|A\| = \|A\|_2$ and the Frobenius norm $\|A\|_F$. The *Frobenius norm* is the most natural extension of the $p = 2$ norm for vectors, but uses a subscript F instead. It is defined for matrix $A \in \mathbb{R}^{n \times d}$ as

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d A_{i,j}^2} = \sqrt{\sum_{i=1}^n \|a_i\|^2},$$

where $A_{i,j}$ is the element in the i th row and j th column of A , and where a_i is the i th row vector of A . The *spectral norm* is defined for a matrix $A \in \mathbb{R}^{n \times d}$ as

$$\|A\| = \|A\|_2 = \max_{x \in \mathbb{R}^d} \|Ax\| / \|x\| = \max_{y \in \mathbb{R}^n} \|yA\| / \|y\|.$$

Its useful to think of these x and y vectors as being unit vectors, then the denominator can be ignored (as they are 1). Then we see that x and y only contain “directional” information, and the $\arg \max$ vectors point in the directions that maximize the norm.

4.5 Linear Independence

Consider a set of k vectors $x_1, x_2, \dots, x_k \in \mathbb{R}^d$, and a set of k scalars $\alpha_1, \alpha_2, \dots, \alpha_k \in \mathbb{R}$. Then because of linearity of vectors, we can write a new vector in \mathbb{R}^d as

$$z = \sum_{i=1}^k \alpha_i x_i.$$

For a set of vectors $X = \{x_1, x_2, \dots, x_k\}$, for any vector z where there exists a set of scalars α where z can be written as the above summation, then we say z is *linearly dependent* on X . If z **cannot** be written with any choice of α_i s, then we say z is *linearly independent* of X . All vectors $z \in \mathbb{R}^d$ which are linearly dependent on X are said to be in its *span*.

$$\text{span}(X) = \left\{ z \mid z = \sum_{i=1}^k \alpha_i x_i, \alpha_i \in \mathbb{R} \right\}.$$

If $\text{span}(X) = \mathbb{R}^d$ (that is for vectors $X = x_1, x_2, \dots, x_k \in \mathbb{R}^d$ all vectors are in the span), then we say X forms a *basis*.

Example: Linear Independence

Consider input vectors in a set X as

$$x_1 = \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix} \quad x_2 = \begin{bmatrix} 2 \\ 4 \\ 1 \end{bmatrix}$$

And two other vectors

$$z_1 = \begin{bmatrix} -3 \\ -5 \\ 2 \end{bmatrix} \quad z_2 = \begin{bmatrix} 3 \\ 7 \\ 1 \end{bmatrix}$$

Note that z_1 is linearly dependent on X since it can be written as $z_1 = x_1 - 2x_2$ (here $\alpha_1 = 1$ and $\alpha_2 = -2$). However z_2 is linearly independent from X since there are no scalars α_1 and α_2 so that $z_2 = \alpha_1 x_1 + \alpha_2 x_2$ (we need $\alpha_1 = \alpha_2 = 1$ so the first two coordinates align, but then the third coordinate cannot).

Also the set X is linearly independent, since there is no way to write $x_2 = \alpha_1 x_1$.

A set of vectors $X = \{x_1, x_2, \dots, x_n\}$ is *linearly independent* if there is no way to write any vector $x_i \in X$ in the set with scalars $\{\alpha_1, \dots, \alpha_{i-1}, \alpha_{i+1}, \dots, \alpha_n\}$ as the sum

$$x_i = \sum_{\substack{j=1 \\ j \neq i}}^n \alpha_j x_j$$

of the other vectors in the set.

4.6 Rank

The *rank* of a set of vectors $X = \{x_1, \dots, x_n\}$ is the size of the largest subset $X' \subset X$ which are linearly independent. Usually we report $\text{rank}(A)$ as the rank of a matrix A . It is defined as the rank of the rows of the matrix, or the rank of its columns; it turns out these quantities are always the same.

If $A \in \mathbb{R}^{n \times d}$, then $\text{rank}(A) \leq \min\{n, d\}$. If $\text{rank}(A) = \min\{n, d\}$, then A is said to be *full rank*. For instance, if $d < n$, then using the rows of $A = [a_1; a_2; \dots; a_n]$, we can describe *any* vector $z \in \mathbb{R}^d$ as the linear combination of these rows: $z = \sum_{i=1}^n \alpha_i a_i$ for some set $\{\alpha_1, \dots, \alpha_n\}$; in fact, we can set all but d of these scalars to 0.

4.7 Inverse

A matrix A is said to be *square* if it has the same number of column as it has rows. A square matrix $A \in \mathbb{R}^{n \times n}$ may have an *inverse* denoted A^{-1} . If it exists, it is a unique matrix which satisfies:

$$A^{-1}A = I = AA^{-1}$$

where I is the $n \times n$ identity matrix

$$I = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix} = \text{diag}(1, 1, \dots, 1).$$

Note that I serves the purpose of 1 in scalar algebra, so for a scalar a then using $a^{-1} = \frac{1}{a}$ we have $aa^{-1} = 1 = a^{-1}a$.

A matrix is said to be *invertable* if it has an inverse. Only square, full-rank matrices are invertable; and a matrix is always invertable if it is square and full rank. If a matrix is not square, the inverse is not defined. If a matrix is not full rank, then it does not have an inverse.

4.8 Orthogonality

Two vectors $x, y \in \mathbb{R}^d$ are *orthogonal* if $\langle x, y \rangle = 0$. This means those vectors are at a right angle to each other.

Example: Orthogonality

Consider two vectors $x = (2, -3, 4, -1, 6)$ and $y = (4, 5, 3, -7, -2)$. They are orthogonal since

$$\langle x, y \rangle = (2 \cdot 4) + (-3 \cdot 5) + (4 \cdot 3) + (-1 \cdot -7) + (6 \cdot -2) = 8 - 15 + 12 + 7 - 12 = 0.$$

A square matrix $U \in \mathbb{R}^{n \times n}$ is *orthogonal* if all of its columns $[u_1, u_2, \dots, u_n]$ are normalized and are all orthogonal with each other. It follows that

$$U^T U = I = U U^T$$

since for any normalized vector u that $\langle u, u \rangle = \|u\| = 1$.

A set of columns (for instance those of an orthogonal U) which are normalized and all orthogonal to each other are said to be *orthonormal*. If $U \in \mathbb{R}^{n \times d}$ and has orthonormal columns, then $U^T U = I$ (here I is $d \times d$) but $U U^T \neq I$.

Orthogonal matrices are norm preserving under multiplication. That means for an orthogonal matrix $U \in \mathbb{R}^{n \times n}$ and any vector $x \in \mathbb{R}^n$, then $\|Ux\| = \|x\|$.

Moreover, the columns $[u_1, u_2, \dots, u_n]$ of an orthogonal matrix $U \in \mathbb{R}^{n \times n}$ form an *basis* for \mathbb{R}^n . This means that for any vector $x \in \mathbb{R}^n$, there exists a set of scalars $\alpha_1, \dots, \alpha_n$ such that $x = \sum_{i=1}^n \alpha_i u_i$. More interestingly, we also have $\|x\|^2 = \sum_{i=1}^n \alpha_i^2$.

This can be interpreted as U describing a *rotation* (with possible mirror flips) to a new set of coordinates. That is the old coordinates of x are (x_1, x_2, \dots, x_n) and the coordinates in the new orthogonal basis $[u_1, u_2, \dots, u_n]$ are $(\alpha_1, \alpha_2, \dots, \alpha_n)$.

4.9 Python numpy Example

Python provides an excellent library called `numpy` (pronounced ‘num-pie’) for handling arrays and matrices, and performing linear basic algebra.

```
import numpy as np
from numpy import linalg as LA

#create an array, a row vector
v = np.array([1,2,7,5])
print v
#[1 2 7 5]
print v[2]
#7

#create a n=2 x d=3 matrix
A = np.array([[3,4,3], [1,6,7]])
print A
#[[3 4 3]
# [1 6 7]]
print A[1,2]
#7
print A[:, 1:3]
#[[4 3]
# [6 7]]

#adding and multiplying vectors
u = np.array([3,4,2,2])
#elementwise add
print v+u
#[4 6 9 7]
#elementwise multiply
print v*u
#[ 3  8 14 10]
# dot product
print v.dot(u)
# 35
print np.dot(u,v)
# 35

#matrix multiplication
B = np.array([[1,2], [6,5], [3,4]])
print A.dot(B)
#[[36 38]
# [58 60]]
x = np.array([3,4])
print B.dot(x)
#[11 38 25]

#norms
print LA.norm(v)
#8.88819441732
print LA.norm(v,1)
#15.0
```

```

print LA.norm(v,np.inf)
#7.0
print LA.norm(A, 'fro')
#10.9544511501
print LA.norm(A,2)
#10.704642743

#transpose
print A.T
#[[3 1]
# [4 6]
# [3 7]]
print x.T
#[3 4] (always prints in row format)

print LA.matrix_rank(A)
#2
C = np.array([[1,2],[3,5]])
print LA.inv(C)
#[[-5.  2.]
# [ 3. -1.]]
print C.dot(LA.inv(C))
#[[ 1.00000000e+00  2.22044605e-16] (nearly [[1 0]
# [ 0.00000000e+00  1.00000000e+00]] [0 1]] )

```

5 Linear Regression

We introduce the basic model of linear regression. It builds a linear model to predict one variable from one other variable or from a set of other variables. We will demonstrate how this simple technique can extend to building potentially much more complex polynomial models. Then we will introduce the central and extremely powerful idea of cross-validation. This method fundamentally changes the statistical goal of validating a model, to characterizing the data.

5.1 Simple Linear Regression

We will begin with the simplest form of linear regression. The input is a set of n 2-dimensional data points $P = (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. The ultimate goal will be to predict the y values using only the x -values. In this case x is the *explanatory variable* and y is the *dependent variable*.

In order to do this, we will “fit” a line through the data of the form

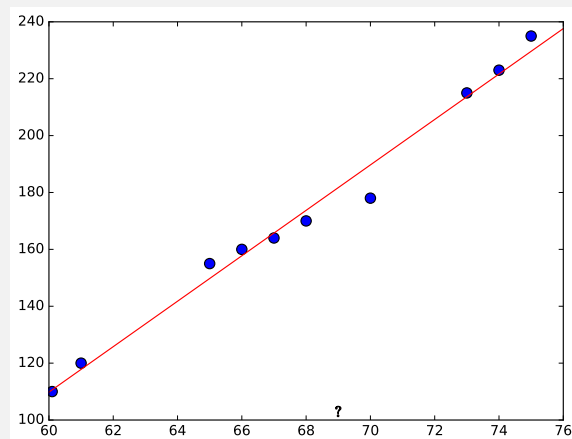
$$y = \ell(x) = ax + b,$$

where a (the slope) and b (the intercept) are parameters of this line. The line ℓ is our “model” for this input data.

Example: Fitting a line to height and weight

Consider the following data set that describes a set of heights and weights.

height (in)	weight (lbs)
66	160
68	170
60	110
70	178
65	155
61	120
74	223
73	215
75	235
67	164
69	?



Note that in the last entry, we have a height of 69, but we do not have a weight. If we were to guess the weight in the last column, how should we do this?

We can draw a line (the red one) through the data points. Then we can guess the weight for a data point with height 69, by the value of the line at height 69 inches: about 182 pounds.

Measuring error. The purpose of this line is not just to be close to all of the data (for this we will have to wait for PCA and dimensionality reduction). Rather, its goal is prediction; specifically, using the explanatory variable x to predict the dependent variable y .

In particular, for every value $x \in \mathbb{R}$, we can predict a value $\hat{y} = \ell(x)$. Then on our dataset, we can examine for each x_i how close \hat{y}_i is to y_i . This difference is called a *residual*:

$$r_i = |y_i - \hat{y}_i| = |y_i - \ell(x_i)|.$$

Note that this residual is not the distance from y_i to the line ℓ , but the distance from y_i to the corresponding point with the same x -value. Again, this is because our only goal is prediction of y . And this will be important as it allows all of the techniques to be immune to the choice of units (e.g., inches or feet, pounds or kilograms)

So the residual measures the error of a single data point, but how should we measure the overall error of the entire data set? The common approach is the sum of squared errors:

$$\text{SSE}(P, \ell) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \ell(x_i))^2.$$

Why is this the most common measure? Here are 3 explanations?

- The sum of squared errors was the optimal result for a single point estimator under Gaussian noise using Bayesian reasoning, when there was assumed Gaussian noise (See T2). In that case the answer was simply the mean of the data.
- If you treat the residuals as a vector $r = (r_1, r_2, \dots, r_n)$, then the standard way to measure total size of a vector r is through its norm $\|r\|$, which is most commonly its 2-norm $\|r\| = \|r\|_2 = \sqrt{\sum_{i=1}^n r_i^2}$. The square root part is not so important (it does not change which line ℓ minimizes this error), so removing this square root, we are left with SSE.
- For this specific formulation, there is a simple closed form solution (which we will see next) for ℓ . And in fact, this solution will generalize to many more complex scenarios.

There are many other formulations of how best to measure error for the fit of a line (and other models), but we will not cover them in this class.

Solving for ℓ . To solve for the line which minimizes $\text{SSE}(P, \ell)$ there is a very simply solution, in two steps. Calculate averages $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$, and create centered n -dimension vectors $\bar{P}_x = (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x})$ for all x -coordinates and $P_y = (y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_n - \bar{y})$ for all y -coordinates.

1. Set $a = \langle \bar{P}_y, \bar{P}_x \rangle / \|\bar{P}_x\|^2$
2. Set $b = \bar{y} - a\bar{x}$

This defines $\ell(x) = ax + b$.

We will not give the full proof of why this is optimal solution (it can be shown by expanding out the SSE expression, taking the derivative, and solving for 0), but we will give a couple points of intuition.

First lets examine the intercept

$$b = \frac{1}{n} \sum_{i=1}^n (y_i - ax_i) = \bar{y} - a\bar{x}$$

This setting of b ensures that the line $y = \ell(x) = ax + b$ goes through the point (\bar{x}, \bar{y}) at the center of the data set since $\bar{y} = \ell(\bar{x}) = a\bar{x} + b$.

Second, to understand how the slope a is chosen, it is illustrative to reexamine the dot product as

$$a = \frac{\langle \bar{P}_y, \bar{P}_x \rangle}{\|x\|^2} = \frac{\|\bar{P}_y\| \cdot \|\bar{P}_x\| \cdot \cos \theta}{\|\bar{P}_x\|^2} = \frac{\|\bar{P}_y\|}{\|\bar{P}_x\|} \cos \theta,$$

where θ is the angle between the n -dimensional vectors y and x . Now in this expression, the $\|\bar{P}_y\|/\|\bar{P}_x\|$ captures how much on (root-squared) average y increases as x does (the rise-over-run interpretation of slope). However, we may want this to be negative if there is a negative correlation between \bar{P}_x and \bar{P}_y , or really this does not matter much if there is no correlation. So the $\cos \theta$ term captures the correlation after normalizing the units of x and y . Again, this is not a formal proof, but hopefully provides some insight.

In python. This is quite easy to demonstrate in python.

```
import numpy as np

x = np.array([66, 68, 65, 70, 65, 62, 74, 70, 71, 67])
y = np.array([160, 170, 159, 188, 150, 120, 233, 198, 201, 164])

ave_x = np.average(x)
ave_y = np.average(y)

#first center the data points
xc = x - ave_x
yc = y - ave_x

a = xc.dot(yc)/xc.dot(xc)
b = ave_y - a*ave_x
print a, b

#or with scipy
from scipy import polyfit
(a,b)=polyfit(x,y,1)
print a, b

#predict weight at x=69
w=a*69+b
```

5.2 Linear Regression with Multiple Explanatory Variables

Magically, using linear algebra, everything extends gracefully to using more than one explanatory variables. Now consider a set of d explanatory variables x_1, x_2, \dots, x_d , and one dependent variable y . We would now like to use all of these variables at once to make a single (linear) prediction about the variable y . That is, we would like to create a model

$$\begin{aligned} y = M(x_1, x_2, \dots, x_d) &= \alpha_0 + \sum_{i=1}^d \alpha_i x_i \\ &= \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_d x_d. \end{aligned}$$

In this notation α_0 serves the purpose of the intercept b , and all of the α_i s replace the single coefficient a in the simple linear regression. Now we have described a more complex linear model M .

Example: Predicting customer value

A website specializing in dongles (dongles-r-us.com) wants to predict the total dollar amount that visitors will spend on their site. It has installed some software that can track three variables:

- **time** (the amount of time on the page in seconds): x_1 ,
 - **jiggle** (the amount of mouse movement in cm): x_2 , and
 - **scroll** (how far they scroll the page down in cm): x_3 .
- Also, for a set of past customers they have recorded the
- **sales** (how much they spend on dongles in cents): y .

We see a portion of their data set here with $n = 11$ customers:

time: x_1	jiggle: x_2	scroll: x_3	sales: y
232	33	402	2201
10	22	160	0
6437	343	231	7650
512	101	17	5599
441	212	55	8900
453	53	99	1742
2	2	10	0
332	79	154	1215
182	20	89	699
123	223	12	2101
424	32	15	8789

To build a model, we recast the data as an 11×4 matrix $X = [\mathbf{1}, x_1, x_2, x_3]$. We let y be the 11-dimensional column vector.

$$X = \begin{bmatrix} 1 & 232 & 33 & 402 \\ 1 & 10 & 22 & 160 \\ 1 & 6437 & 343 & 231 \\ 1 & 512 & 101 & 17 \\ 1 & 441 & 212 & 55 \\ 1 & 453 & 53 & 99 \\ 1 & 2 & 2 & 10 \\ 1 & 332 & 79 & 154 \\ 1 & 182 & 20 & 89 \\ 1 & 123 & 223 & 12 \\ 1 & 424 & 32 & 15 \end{bmatrix} \quad y = \begin{bmatrix} 2201 \\ 0 \\ 7650 \\ 5599 \\ 8900 \\ 1742 \\ 0 \\ 1215 \\ 699 \\ 2101 \\ 8789 \end{bmatrix}$$

The goal is to learn the 4-dimensional column vector $\alpha = [\alpha_0; \alpha_1; \alpha_2; \alpha_3]$ so

$$y \approx X\alpha.$$

Setting $\alpha = (X^T X)^{-1} X^T y$ obtains (roughly) $\alpha_0 = 262$, $\alpha_1 = 0.42$, $\alpha_2 = 12.72$, and $\alpha_3 = -6.50$. This implies an average customer with no interaction on the site generates $\alpha_0 = \$2.62$. That **time** does not have a strong effect here (only a coefficient α_1 at only 0.42), but **jiggle** has a strong correlation (with coefficient $\alpha_2 = 12.72$, this indicates 12 cents for every centimeter of mouse movement). Meanwhile **scroll** has a negative effect (with coefficient $\alpha_3 = -6.5$); this means that the more they scroll, the less likely they are to spend (just browsing dongles!).

Given a data point $x = (x_1, x_2, \dots, x_d)$, we can again evaluate our prediction $\hat{y} = M(x)$ using the residual value $r_i = |y_i - \hat{y}_i| = |y_i - M(x_i)|$. And to evaluate a set of n data points, it is standard to consider the sum of squared error as

$$\text{SSE}(X, y, M) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - M(x_i))^2.$$

To obtain the coefficients which minimize this error, we can now do so with very simple linear algebra.

First we construct a $n \times (d + 1)$ data matrix $X = [\mathbf{1}, x_1, x_2, \dots, x_d]$, where the first column $\mathbf{1}$ is the all ones column vector $[1; 1; \dots; 1]$. The next d columns describe for the i th row, the data values of the explanatory variables x_1 through x_d for the i th data point. Then we let y be a n -dimensional column vector with all data for the dependent variable. Now we can simply calculate the $(d + 1)$ -dimensional column vector $\alpha = [\alpha_0; \alpha_1; \dots; \alpha_d]$ as

$$\alpha = (X^T X)^{-1} X^T y.$$

Let us compare to the simple case where we have 1 explanatory variable. The $(X^T X)^{-1}$ term replaces the $\frac{1}{\|\bar{P}_x\|^2}$ term. The $X^T y$ replaces the dot product $\langle \bar{P}_y, \bar{P}_x \rangle$. And we do not need to separately solve for the intercept b , since we have created a new column in X of all 1s. Now for any dependent data values, we multiply the found coefficient α_0 by an imaginary 1 data value.

In python: We can directly write this in python as

```
import numpy as np
from numpy import linalg as LA

# directly
alpha = np.dot(np.dot(LA.inv(np.dot(X.T, X)), X.T), y.T)

# or with LA.lstsq
alpha = LA.lstsq(X, y)[0]
```

Or in many other ways.

5.3 Polynomial Regression

Sometimes linear relations are not sufficient to capture the true pattern going on in the data with even a single dependent variable x . Instead we would like to build a model of the form:

$$\hat{y} = M_2(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2$$

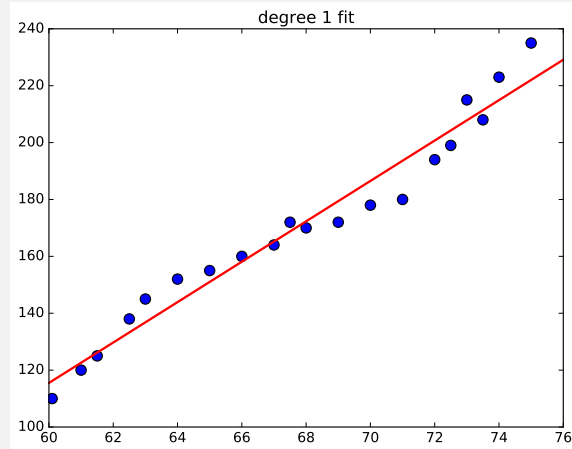
or more generally for some polynomial of degree p

$$\begin{aligned} \hat{y} = M_p(x) &= \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_p x^p \\ &= \alpha_0 + \sum_{i=1}^p \alpha_i x^i. \end{aligned}$$

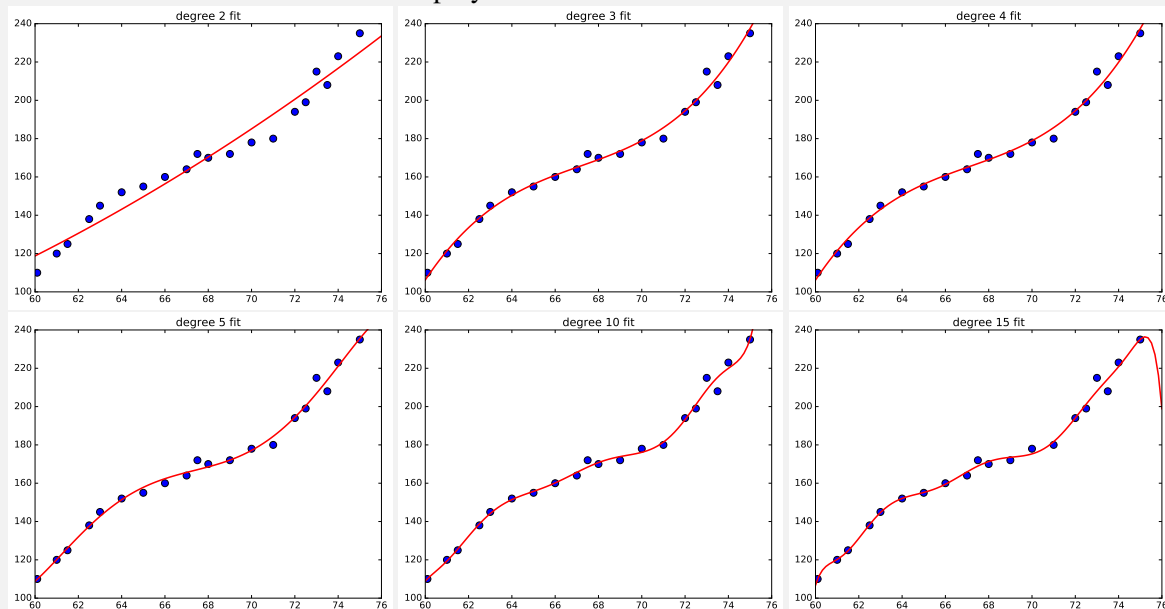
Example: Predicting Height and Weight with Polynomials

We found more height and weight data, in addition to the ones in the height-weight example above.

height (in)	weight (lbs)
61.5	125
73.5	208
62.5	138
63	145
64	152
71	180
69	172
72.5	199
72	194
67.5	172



But can we do better if we fit with a polynomial?



Again we can measure error for a single data point $p_i = (x_i, y_i)$ as a residual as $r_i = |\hat{y} - y_i| = |M_p(x_i) - y_i|$ and the error on n data points as the sum of squared residuals

$$\text{SSE}(P, M_p) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (M_p(x_i) - y_i)^2.$$

Under this error measure, it turns out we can again find a simple solution for the residuals $\alpha = [\alpha_0, \alpha_1, \dots, \alpha_p]$. For each dependent variable data value x we create a $(p + 1)$ -dimensional vector

$$v = (1, x, x^2, \dots, x^p).$$

And then for n data points $(x_1, y_1), \dots, (x_n, y_n)$ we can create an $n \times (p + 1)$ data matrix

$$X_p = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

Then we can solve the same way as if each data value raised to a different power was a different dependent variable. That is we can solve for the coefficients $\alpha = [\alpha_0; \alpha_1; \alpha_2; \dots; \alpha_n]$ as

$$\alpha = (X_p^T X_p)^{-1} X_p^T y.$$

5.4 Cross Validation

So how do we choose the correct value of p , the degree of the polynomial fit?

A (very basic) statistical (hypothesis testing) approach may be choose a model of the data (the best fit curve for some polynomial degree p , and assume Gaussian noise), then calculate the probability that the data fell outside the error bounds of that model. But maybe many different polynomials are a good fit?

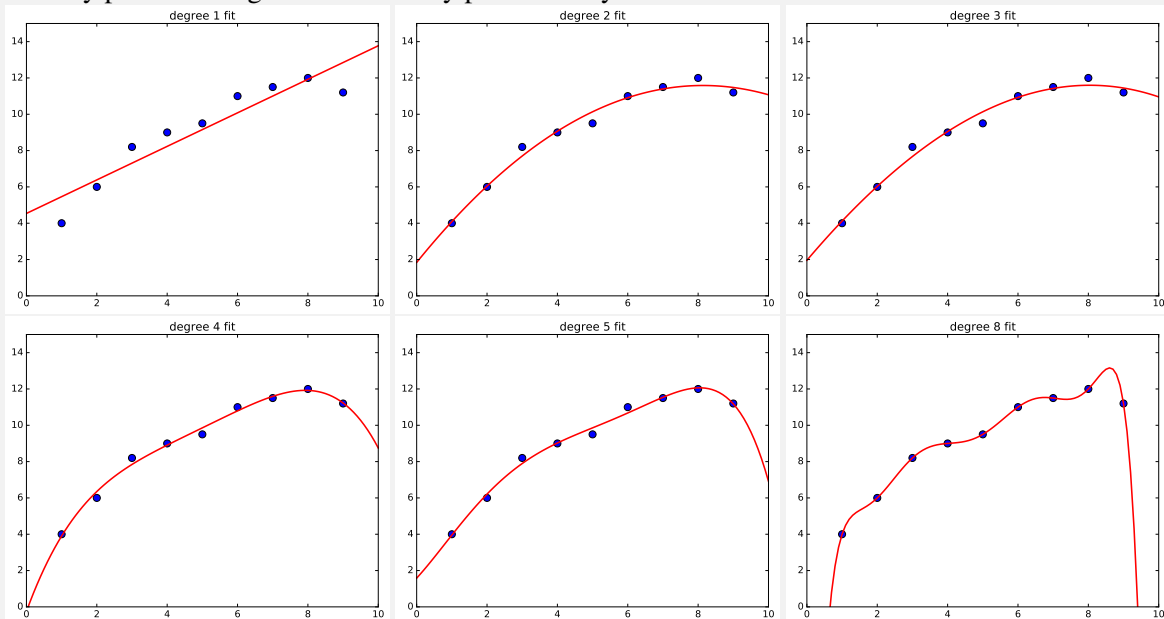
In fact, if we choose p as $n - 1$ or greater, then the curve will *polynomially interpolate* all of the points. That is, it will pass through all points, so all points have a residual of exactly 0 (up to numerical precision). This is the basis of a lot of geometric modeling (e.g., for CAD), but bad for data modeling.

Example: Simple polynomial example

Consider the simple data set of 9 points

x	1	2	3	4	5	6	7	8	9
y	4	6	8.2	9	9.5	11	11.5	12	11.2

With the following polynomial fits for $p = \{1, 2, 3, 4, 5, 8\}$. Believe your eyes, for $p = 8$, the curve actually passes through each and every point exactly.



Recall, our goal was for a new data point with only an x value to predict its y -value. Which do you think does the best job?

Generalization and Cross-Validation. Our ultimate goal in regression is *generalization* (how well do we do on new data), not SSE! Using some error measure (SSE) to fit a line or curve, is a good proxy for what we want, but in many cases (as with polynomial regression), it can be abused. We want to know how our model will generalize to new data. How would we measure this without new data?

The solution is *cross-validation*. In the simplest form, we **randomly** split our data into training data (on which we build a model) and testing data (on which we evaluate our model). The testing serves to estimate how well we would do on future data which we do not have.

- *Why randomly?:* Because you do not want to bias the model to do better on some parts than other in how you choose the split. Also, since we assume the data elements come iid from some underlying distribution, then the test data is also iid if you chose it randomly.
- *How large should the test data be?:* It depends on the data set. Both 10% and 33% are common.

Let (X, y) be the full data set (with n rows of data), and we split it into data sets $(X_{\text{train}}, y_{\text{train}})$ and $(X_{\text{test}}, y_{\text{test}})$ with n_{train} and n_{test} rows, respectively. With $n = n_{\text{train}} + n_{\text{test}}$. Next we build a model with the training data, e.g.,

$$\alpha = (X_{\text{train}}^T X_{\text{train}})^{-1} X_{\text{train}}^T y_{\text{train}}.$$

Then we evaluate the model M_α on the test data X_{test} , often using $\text{SSE}(X_{\text{test}}, y_{\text{test}}, M_\alpha)$ as

$$\text{SSE}(X_{\text{test}}, y_{\text{test}}, M_\alpha) = \sum_{(x_i, y_i) \in (X_{\text{test}}, y_{\text{test}})} (y_i - M_\alpha(x_i))^2 = \sum_{(x_i, y_i) \in (X_{\text{test}}, y_{\text{test}})} (y_i - \langle (x_i; 1), \alpha \rangle)^2.$$

We can use the testing data for two purposes:

- To estimate how well our model would perform on new data, yet unseen. That is the predicted residual of a new data point is precisely $\text{SSE}(X_{\text{test}}, y_{\text{test}}, M_\alpha)/n_{\text{test}}$.
- To choose the correct parameter for a model (which p to use)?

Its important to keep in mind that we should not use the same $(X_{\text{test}}, y_{\text{test}})$ to do both tasks. If we choose a model with $(X_{\text{test}}, y_{\text{test}})$, then we should reserve even more data for predicting the generalization error. When using the test data to choose a model parameter, the it is being used to build the model; thus evaluating generalization with this same data can suffer the same fate as testing and training with the same data.

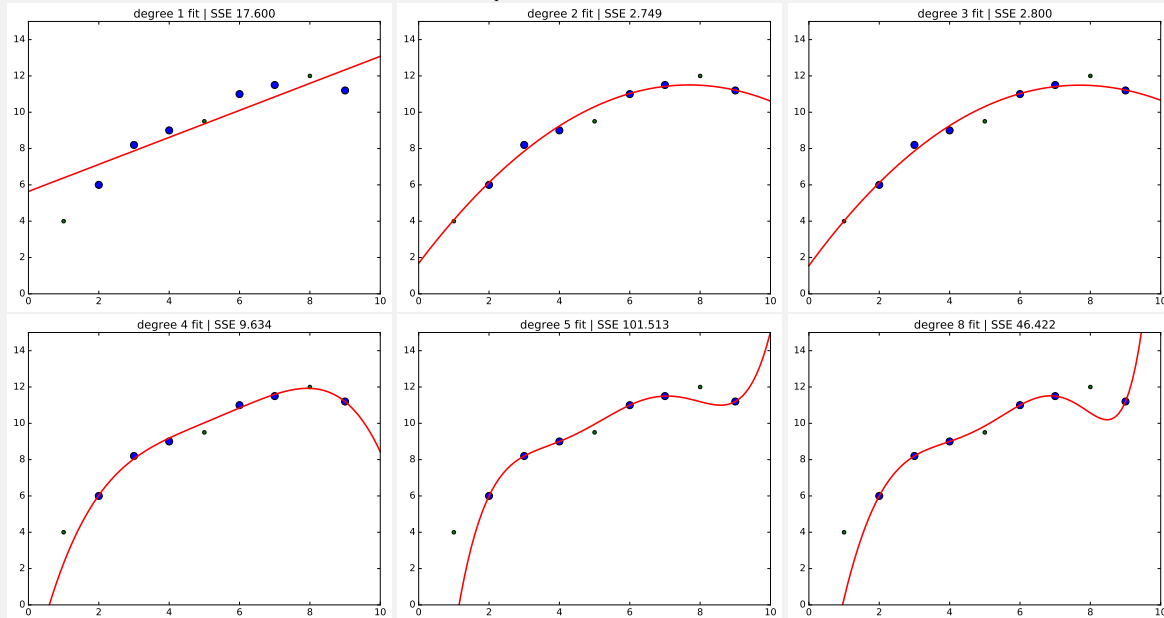
So how should we choose the best p ? We calculate models M_{α_p} for each value p on the same training data. Then calculate the model error $\text{SSE}(X_{\text{test}}, y_{\text{test}}, M_{\alpha_p})$ for each p , and see which has the smallest value. That is we train on $(X_{\text{train}}, y_{\text{train}})$ and test(evaluate) on $(X_{\text{test}}, y_{\text{test}})$.

Example: Simple polynomial example with Cross Validation

Now split our data sets into a train set and a test set:

$$\text{train: } \begin{array}{c|cccccc} x & 2 & 3 & 4 & 6 & 7 & 8 \\ \hline y & 6 & 8.2 & 9 & 11 & 11.5 & 12 \end{array} \quad \text{test: } \begin{array}{c|ccc} x & 1 & 5 & 9 \\ \hline y & 4 & 9.5 & 11.2 \end{array}$$

With the following polynomial fits for $p = \{1, 2, 3, 4, 5, 8\}$ generating model M_{α_p} on the test data. We then calculate the $\text{SSE}(x_{\text{test}}, y_{\text{test}}, M_{\alpha_p})$ score for each (as shown):



And the polynomial model with degree $p = 2$ has the lowest SSE score of 2.749. It is also the simplest model that does a very good job by the “eye-ball” test. So we would choose this as our model.

Leave-one-out Cross Validation. But, not training on the test data means that you use less data, and your model is worse! You don’t want to waste this data!

If your data is very large, then leaving out 10% is not a big deal. But if you only have 9 data points it can be. The smallest the test set could be is 1 point. But then it is not a very good representation of the full data set.

The alternative is to create n different training sets, each of size $n-1$ ($X_{1,\text{train}}, X_{2,\text{train}}, \dots, X_{n,\text{train}}$), where $X_{i,\text{train}}$ contains all points except for x_i , which is a one-point test set. Then we build n different models M_1, M_2, \dots, M_n , evaluate each model M_i on the one test point x_i to get an error $E_i = (y_i - M_i(x_i))^2$, and average their errors $E = \frac{1}{n} \sum_{i=1}^n E_i$. Again, the parameter with the smallest associated average error E is deemed the best. This allows you to build a model on as much data as possible, while still using all of the data to test.

However, this requires roughly n times as long to compute as the other techniques, so is often too slow for really big data sets.

```

import matplotlib as mpl
mpl.use('PDF')
import matplotlib.pyplot as plt
import scipy as sp
import numpy as np
import math
from numpy import linalg as LA

def plot_poly(x,y,xE,yE,p):
    plt.scatter(x,y, s=80, c="blue")
    plt.scatter(xE,yE, s=20, c="green")
    plt.axis([0,10,0,15])

    s=sp.linspace(0,10,101)

    coefs=sp.polyfit(x,y,p)
    ffit = np.polyld(coefs)
    plt.plot(s,ffit(s), 'r-', linewidth=2.0)

    #evaluate on xE, yE
    resid = ffit(xE)
    RMSE = LA.norm(resid-yE)
    SSE = xE.size * RMSE * RMSE

    title = "degree_%s_fit_|_SSE_%0.3f" % (p, SSE)
    plt.title(title)
    file = "CVpolyReg%s.pdf" % p
    plt.savefig(file, bbox_inches='tight')
    plt.clf()
    plt.cla()

# train data
xT = np.array([2, 3, 4, 6, 7, 9])
yT = np.array([6, 8.2, 9, 11, 11.5, 11.2])

#test data
xE = np.array([1, 5, 8])
yE = np.array([4, 9.5, 12])

p_vals = [1,2,3,4,5,8]
for i in p_vals:
    plot_poly(xT,yT,xE,yE,i)

```

6 Gradient Descent

In this topic we will discuss optimizing over general functions f . Typically the function is defined $f : \mathbb{R}^d \rightarrow \mathbb{R}$; that is its domain is multi-dimensional (in this case d -dimensional) and output is a real scalar (\mathbb{R}). This often arises to describe the “cost” of a model which has d parameters which describe the model (e.g., degree $(d - 1)$ -polynomial regression) and the goal is to find the parameters with minimum cost. Although there are special cases where we can solve for these optimal parameters exactly, there are many cases where we cannot. What remains in these cases is to analyze the function f , and try to find its minimum point. The most common solution for this is gradient descent where we try to “walk” in a direction so the function decreases until we no-longer can.

6.1 Functions

We review some basic properties of a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Again, the goal will be to unveil abstract tools that are often easy to imagine in low dimensions, but automatically generalize to high-dimensional data. We will first provide definitions without any calculus.

Let $B_r(x)$ define a Euclidean ball around a point x of radius r . That is, it includes all points $\{y \in \mathbb{R}^d \mid \|x - y\| \leq r\}$, within a Euclidean distance of r from x . We will use $B_r(x)$ to define a *local neighborhood* around a point x . The idea of “local” is quite flexible, and we can use *any* value of $r > 0$, basically it can be as small as we need it to be, as long as it is strictly greater than 0.

Minima and maxima. A *local maximum* of f is a point $x \in \mathbb{R}^d$ so for some neighborhood $B_r(x)$, all points $y \in B_r(x)$ have smaller (or equal) function value than at x : $f(y) \leq f(x)$. A *local minimum* of f is a point $x \in \mathbb{R}^d$ so for some neighborhood $B_r(x)$, all points $y \in B_r(x)$ have larger (or equal) function value than at x : $f(y) \geq f(x)$. If we remove the “or equal” condition for both definitions for $y \in B_r(x), y \neq x$, we say the maximum or minimum points are *strict*.

A point $x \in \mathbb{R}^d$ is a *global maximum* of f if for all $y \in \mathbb{R}^d$, then $f(y) \leq f(x)$. Likewise, a point $x \in \mathbb{R}^d$ is a *global minimum* if for all $y \in \mathbb{R}^d$, then $f(y) \geq f(x)$. There may be multiple global minimum and maximum. If there is exactly one point $x \in \mathbb{R}^d$ that is a global minimum or global maximum, we again say it is *strict*.

When we just use the term *minimum* or *maximum* (without local or global) it implies a local minimum or maximum.

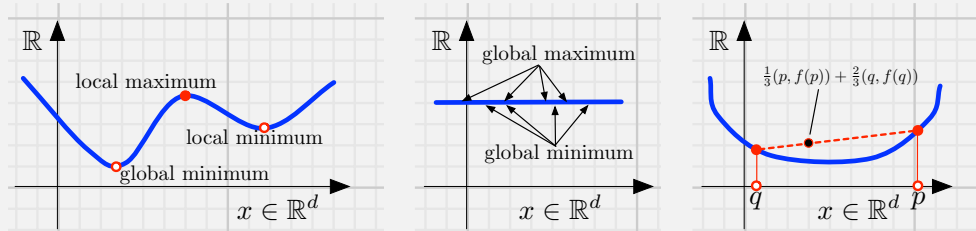
Focusing on a function restricted to a closed and bounded subset $S \subset \mathbb{R}^d$, if the function is continuous (we won’t formally define this, but it likely means what you think it does), then the function must have a global minimum and a global maximum. It may occur on the boundary of S .

A *saddle point* is a type of point $x \in \mathbb{R}^d$ so that within any neighborhood $B_r(x)$, it has points $y \in B_r(x)$ with $f(y) < f(x)$ (the lower points) and $y' \in B_r(x)$ with $f(y') > f(x)$ (the upper points). In particular, it is a saddle if there are disconnected regions of upper points (and of lower points). For $d = 1$, then there can be no saddle points. If these regions are connected, and it is not a minimum or maximum, then it is a *regular point*.

For an arbitrary (or randomly) chosen point x , it is usually a regular point (except for examples you are unlikely to encounter, the set of minimum, maximum, and saddle points are finite, while the set of regular points is infinite).

Example: Continuous Functions

Here we show some example functions, where the x -axis represents a d -dimensional space. The first function has local minimum and maximum. The second function has a constant value, so every point is a global minimum and a global maximum. The third function f is convex, which is demonstrated with the line segment between points $(p, f(p))$ and $(q, f(q))$ is always above the function f .



Convex functions. In many cases we will assume (or at least desire) that our function is convex.

To define this it will be useful to define a line $\ell \subset \mathbb{R}^d$ as follows with any two points $p, q \in \mathbb{R}^d$. Then for any scalar $\alpha \in \mathbb{R}$, a line ℓ is the set of points

$$\ell = \{x = \alpha p + (1 - \alpha)q \mid \alpha \in \mathbb{R} \text{ and } p, q \in \mathbb{R}^d\}.$$

When $\alpha \in [0, 1]$, then this defines the line segment between p and q .

A function is *convex* if for any two points $p, q \in \mathbb{R}^d$, on the line segment between them has value less than (or equal) to the values at the weighted average of p and q . That is, it is convex if

$$\text{For all } p, q \in \mathbb{R}^d \text{ and for all } \alpha \in [0, 1] \quad f(\alpha p + (1 - \alpha)q) \leq \alpha f(p) + (1 - \alpha)f(q).$$

Removing the (or equal) condition, the function becomes *strictly convex*.

There are many very cool properties of convex functions. For instance, for two convex functions f and g , then $h(x) = f(x) + g(x)$ is convex and so is $h(x) = \max\{f(x), g(x)\}$. But one will be most important for us:

- Any local minimum of a convex function will also be a global minimum. A strictly convex function will have at most a single minimum: the global minimum.

This means if we find a minimum, then we must have also found a global minimum (our goal).

6.2 Gradients

For a function $f(x) = f(x_1, x_2, \dots, x_d)$, and a unit vector $u = (u_1, u_2, \dots, u_d)$, then the *directional derivative* is defined

$$\nabla_u f(x) = \lim_{h \rightarrow 0} \frac{f(x + hu) - f(x)}{h}.$$

We are interested in functions f which are *differentiable*; this implies that $\nabla_u f(x)$ is well-defined for all x and u . The converse is not necessarily true.

Let $e_1, e_2, \dots, e_d \in \mathbb{R}^d$ be a specific set of unit vectors so that $e_i = (0, 0, \dots, 0, 1, 0, \dots, 0)$ where for e_i the 1 is in the i th coordinate.

Then define

$$\nabla_i f(x) = \nabla_{e_i} f(x) = \frac{d}{dx_i} f(x).$$

It is the derivative in the i th coordinate, treating all other coordinates as constants.

We can now, for a differentiable function f , define the *gradient of f* as

$$\nabla f = \frac{df}{dx_1}e_1 + \frac{df}{dx_2}e_2 + \dots + \frac{df}{dx_d}e_d = \left(\frac{df}{dx_1}, \frac{df}{dx_2}, \dots, \frac{df}{dx_d} \right).$$

Note that ∇f is a function from $\mathbb{R}^d \rightarrow \mathbb{R}^d$, which we can evaluate at any point $x \in \mathbb{R}^d$.

Example: Gradient

Consider the function $f(x, y, z) = 3x^2 - 2y^3 - 2xe^z$. Then $\nabla f = (6x - 2e^z, -6y^2, -2xe^z)$ and $\nabla f(3, -2, 1) = (18 - 2e, 24, -6e)$.

Linear approximation. From the gradient we can easily recover the directional derivative of f at point x , for any direction (unit vector) u as

$$\nabla_u f(x) = \langle \nabla f(x), u \rangle.$$

This implies the gradient describes the linear approximation of f at a point x . The slope of the tangent plane of f at x in any direction u is provided by $\nabla_u f(x)$.

Hence, the direction which f is increasing the most at a point x is the unit vector u where $\nabla_u f(x) = \langle \nabla f(x), u \rangle$ is the largest. This occurs at $\frac{\nabla f(x)}{\|\nabla f(x)\|}$, the normalized gradient vector.

To find the minimum of a function f , we then typically want to move from any point x in the direction $-\nabla f(x)$; this is the direction of steepest descent.

6.3 Gradient Descent

Gradient descent is a family of techniques that, for a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, try to identify either

$$\min_{x \in \mathbb{R}^d} f(x) \quad \text{and/or} \quad x^* = \arg \min_{x \in \mathbb{R}^d} f(x).$$

This is effective when f is convex and we do not have a “closed form” solution x^* . The algorithm is iterative, in that it may never reach the completely optimal x^* , but it keeps getting closer and closer.

Algorithm 6.3.1 Gradient Descent(f, x_{start})

initialize $x^{(0)} = x_{\text{start}} \in \mathbb{R}^d$.

repeat

$$x^{(k+1)} := x^{(k)} - \gamma_k \nabla f(x^{(k)})$$

until ($\|\nabla f(x^{(k)})\| \leq \tau$)

return $x^{(k)}$

Basically, for any starting point $x^{(0)}$ the algorithm moves to another point in the direction opposite to the gradient – in the direction that locally decreases f the fastest.

Stopping condition. The parameter τ is the tolerance of the algorithm. If we assume the function is differentiable, then at the minimum x^* , we must have that $\nabla f(x) = (0, 0, \dots, 0)$. So close to the minimum, it should also have a small norm. The algorithm may never reach the true minimum (and we don’t know what it is, so we cannot directly compare against the function value). So we use $\|\nabla f\|$ as a proxy.

In other settings, we may run for a fixed number T steps.

6.3.1 Learning Rate

The most critical parameter of gradient descent is γ , the *learning rate*. In many cases the algorithm will keep $\gamma_k = \gamma$ fixed for all k . It controls how fast the algorithm works. But if it is too large, when we approach the minimum, then the algorithm may go too far, and overshoot it. How should we choose γ ?

Lipschitz bound. We say a function $g : \mathbb{R}^d \rightarrow \mathbb{R}^k$ is *L-Lipschitz* if for all $x, y \in \mathbb{R}^d$ that

$$\|g(x) - g(y)\| \leq L\|x - y\|.$$

If ∇f is *L-Lipschitz*, and we set $\gamma \leq \frac{1}{L}$, then gradient descent will converge to a stationary point. Moreover, if f is convex with global minimum x^* , then after $k = O(1/\varepsilon)$ steps we can guarantee that

$$f(x^{(k)}) - f(x^*) \leq \varepsilon.$$

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *η -strongly convex* with parameter $\eta > 0$ if for all $x \in \mathbb{R}^d$ and any unit vector $u \in \mathbb{R}^d$ (that is $\|u\| = 1$) then

$$\langle \nabla f(x), u \rangle \geq \eta.$$

For an η -strongly convex function f with global minimum x^* , gradient descent with learning rate $\gamma \leq 2/(\eta + L)$ after only $k = O(\log(1/\varepsilon))$ steps will achieve

$$f(x^{(k)}) - f(x^*) \leq \varepsilon.$$

The constant in $k = O(\log(1/\varepsilon))$ depends on the condition number L/η ; recall that for any unit vector u we have $L \geq \langle \nabla f(x), u \rangle \geq \eta$. When an algorithm converges at such a rate, it is known as *linear convergence* since the log-error $\log(f(x^{(k)}) - f(x^*))$ looks like a linear function of k .

Line search. In other cases, we can search for γ_k at each step. Once we have computed the gradient $\nabla f(x^{(k)})$ then we have reduced the high-dimensional minimization problem to a one-dimensional problem. Note if f is convex, then f restricted to this one-dimensional search is also convex. We still need to find the minimum of an unknown function, but we can perform some procedure akin to binary search. We first find a value γ' such that

$$f\left(x^{(k)} - \gamma' \nabla f(x^{(k)})\right) > f(x^{(k)})$$

then we keep subdividing the region $[0, \gamma']$ into pieces which must contain the minimum (for instance the *golden section search* keeps dividing by the golden ratio).

In other situations, we can solve for the optimal γ_k exactly at each step. This is the case if we can again analytically take the derivative $\frac{d}{d\gamma} (f(x^{(k)}) - \gamma \nabla f(x^{(k)}))$ and solve for the γ where it is equal to 0.

Adjustable rate. In practice, line search is often slow. Also, we may not have a Lipschitz bound. It is often better to try a few fixed γ values, probably being a bit conservative. As long as $f(x^{(k)})$ keep decreasing, it works well. This also may alert us if there there is more than one local minimum if the algorithm converges to different locations.

An algorithm called “backtracking line search” automatically tunes the parameter γ . It uses a fixed parameter $\beta \in (0, 1)$ (preferably in $(0.1, 0.8)$; for instance use $\beta = 3/4$). Start with a large step size γ (e.g., $\gamma = 1$). Then at each step of gradient descent at location x , if

$$f(x - \gamma \nabla f(x)) > f(x) - \frac{\gamma}{2} \|\nabla f(x)\|^2$$

then update $\gamma = \beta\gamma$. This shrinks γ over the course of the algorithm to ensure it will satisfy the condition for linear convergence.

Example: Gradient Descent with Fixed Learning Rate

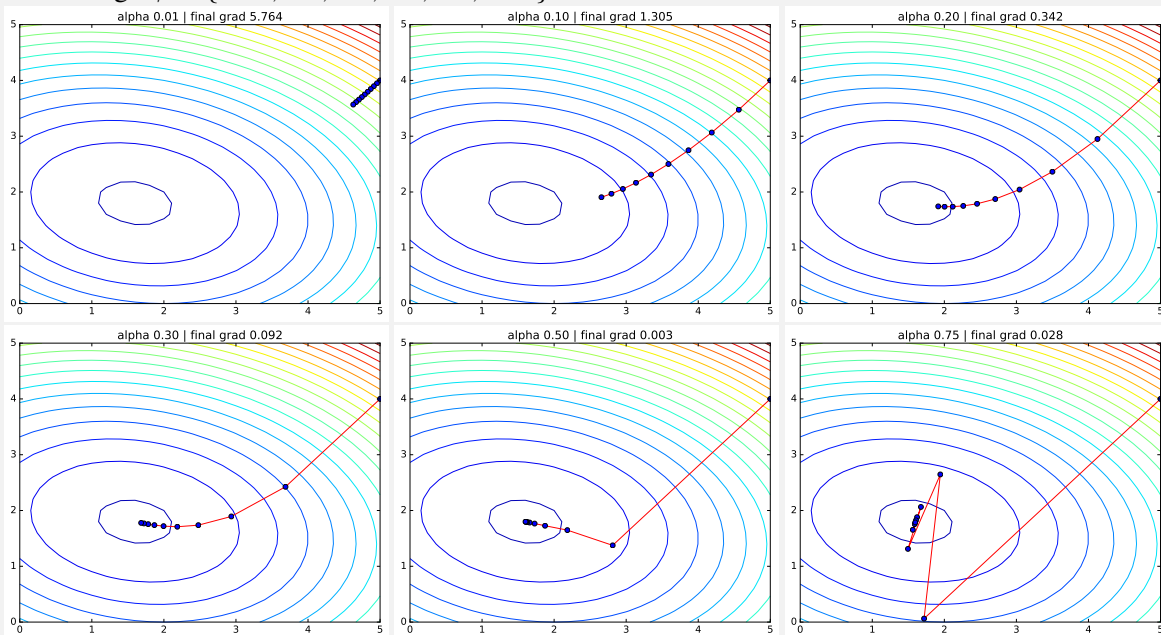
Consider the function

$$f(x, y) = \left(\frac{3}{4}x - \frac{3}{2}\right)^2 + (y - 2)^2 + \frac{1}{4}xy$$

with gradient

$$\nabla f(x, y) = \left(\frac{9}{8}x - \frac{9}{4} + \frac{1}{4}y, 2y - 4 + \frac{1}{4}x\right)$$

We run gradient descent for 10 iterations within initial position (5, 4), while varying the learning rate in the range $\gamma = \{0.01, 0.1, 0.2, 0.3, 0.5, 0.75\}$.



We see that with γ very small, the algorithm does not get close to the minimum. When γ is too large, then the algorithm jumps around a lot, and is in danger of not converging. But at a learning rate of $\gamma = 0.5$ it converges fairly smoothly and reaches a point where $\|\nabla f(x, y)\|$ is very small. Using $\gamma = 0.5$ almost overshoots in the first step; $\gamma = 0.3$ is smoother, and its probably best to use a curve that looks smooth like that one, but with a few more iterations.

```
import matplotlib as mpl
mpl.use('PDF')
import numpy as np
import matplotlib.pyplot as plt
from numpy import linalg as LA

def func(x, y):
    return (0.75*x-1.5)**2 + (y-2.0)**2 + 0.25*x*y

def func_grad(vx,vy):
    dfdx = 1.125*vx - 2.25 + 0.25*vy
    dfdy = 2.0*vy - 4.0 + 0.25*vx
    return np.array([dfdx,dfdy])

#prepare for contour plot
xlist = np.linspace(0, 5, 26)
ylist = np.linspace(0, 5, 26)
```

```

x, y = np.meshgrid(xlist, ylist)
z = func(x,y)
lev = np.linspace(0,20,21)

#iterate location
v_init = np.array([5,4])
num_iter = 10
values = np.zeros([num_iter,2])

for alpha in [0.01, 0.1, 0.2, 0.3, 0.5, 0.75]:
    values[0,:] = v_init
    v = v_init

    # actual gradient descent algorithm
    for i in range(1,num_iter):
        v = v - alpha * func_grad(v[0],v[1])
        values[i,:] = v

#plotting
plt.contour(x,y,z,levels=lev)
plt.plot(values[:,0],values[:,1], 'r-')
plt.plot(values[:,0],values[:,1], 'bo')
grad_norm = LA.norm(func_grad(v[0],v[1]))
title = "alpha_%.2f_|_final_grad_%.3f" % (alpha,grad_norm)
plt.title(title)
file = "gd-%2.0f.pdf" % (alpha*100)
plt.savefig(file, bbox_inches='tight')
plt.clf()
plt.cla()

```

6.4 Fitting a Model to Data

For data analysis, the most common use of gradient descent is to fit a model to data. In this setting we have a data set $P = \{p_1, p_2, \dots, p_n\}$ and a family of models \mathcal{M} so each possible model M_α is defined by a d -dimensional vector $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_d\}$ for p parameters.

Next we define a *loss function* $L(P, M_\alpha)$ which measures the difference between what the model predicts and what the data values are. To choose which parameters generate the best model, we let $f(\alpha) : \mathbb{R}^d \rightarrow \mathbb{R}$ be our function of interest, defined $f(\alpha) = L(P, M_\alpha)$. Then we can run gradient descent to find our model M_{α^*} . For instance we can set

$$f(\alpha) = L(P, M_\alpha) = \text{SSE}(P, M_\alpha) = \sum_{p \in P} (p_y - M_\alpha(p_x))^2. \quad (6.1)$$

This is used for examples including maximum likelihood (or maximum log-likelihood) estimators from Bayesian inference. This includes finding a single point estimator with Gaussian (where we had a closed-form solution), but also many other variants (often where there is no known closed-form solution). It also includes least squares regression and its many variants; we will see this in much more detail next. And will include other topics (including clustering, PCA, classification) we will see later in class.

6.4.1 Least Mean Squares Updates for Linear Regression

Now we will work through how to use gradient descent for simple quadratic regression. It should be straightforward to generalize to linear regression, multiple-explanatory variable linear regression, or general polynomial regression from here. This will specify the function $f(\alpha)$ in equation (6.1) to where

$d = 3$, $\alpha = (\alpha_0, \alpha_1, \alpha_2)$, for each $p \in P$ we have $p = (p_x, p_y) \in \mathbb{R}^2$, but we consider a vector $q = (q_0 = p_x^0, q_1 = p_x^1, q_2 = p_x^2)$ as the set of explanatory variables. Finally,

$$M_\alpha(p_x) = \alpha_0 + \alpha_1 p_x + \alpha_2 p_x^2 = \alpha_0 q_0 + \alpha_1 q_1 + \alpha_2 q_2.$$

To specify the gradient descent step:

$$\alpha := \alpha - \gamma \nabla f(\alpha)$$

we need to define $\nabla f(\alpha)$.

We will first show this for the case where $n = 1$, that is when there is a single data point $p = (p_x, p_y) = (x_1, y_1)$. It should now be easy to verify that the cost function $f_1(\alpha) = (\alpha_0 + \alpha_1 x_1 + \alpha_2 x_1^2 - y_1)^2$ is convex. Next, derive

$$\begin{aligned} \frac{d}{d\alpha_j} f(\alpha) &= \frac{d}{d\alpha_j} (M_\alpha(x_1) - y_1)^2 \\ &= 2(M_\alpha(x_1) - y_1) \frac{d}{d\alpha_j} (M_\alpha(x_1) - y_1) \\ &= 2(M_\alpha(x_1) - y_1) \frac{d}{d\alpha_j} \left(\sum_{j=0}^2 \alpha_j x_1^j - y_1 \right) \\ &= 2(M_\alpha(x_1) - y_1) x_1^j \end{aligned}$$

Thus, we define

$$\begin{aligned} \nabla f(\alpha) &= \left(\frac{d}{d\alpha_0} f(\alpha), \frac{d}{d\alpha_1} f(\alpha), \frac{d}{d\alpha_2} f(\alpha) \right) \\ &= 2 \left((M_\alpha(x_1) - y_1), (M_\alpha(x_1) - y_1) x_1, (M_\alpha(x_1) - y_1) x_1^2 \right) \end{aligned}$$

Applying $\alpha := \alpha - \gamma \nabla f(\alpha)$ according to this specification is known as the *LMS (least mean squares)* update rule or the *Widrow-Huff learning rule*. Quite intuitively, the magnitude the update is proportional to the residual norm $(M_\alpha(x_1) - y_1)$. So if we have a lot of error in our guess of α , then we take a large step; if we do not have a lot of error, we take a small step.

To generalize this to multiple data points ($n > 1$), there are two standard ways. Both of these take strong advantage of the cost function $f(\alpha)$ being *decomposable*. That is, we can write

$$f(\alpha) = \sum_{i=1}^n f_i(\alpha),$$

where each f_i depends only on the i th data point $p_i \in P$. In particular, where $p_i = (x_i, y_i)$, then

$$f_i(\alpha) = (M_\alpha(x_i) - y_i)^2 = (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - y_i)^2.$$

First notice that since f is the sum of f_i s, where each is convex, then f must also be convex; in fact the sum of these usually becomes strongly convex. Also two approaches towards gradient descent will take advantage of this decomposition in slightly different ways. This decomposable property holds for most loss functions for fitting a model to data.

Batch gradient descent. The first technique, called *batch gradient descent*, simply extends the definition of $\nabla f(\alpha)$ to the case with multiple data points. Since f is decomposable, then we use the linearity of the derivative to define

$$\frac{d}{d\alpha_j} f(\alpha) = \sum_{i=1}^n 2(M_\alpha(x_i) - y_i)x_i^j$$

and thus

$$\nabla f(\alpha) = \sum_{i=1}^n (2(M_\alpha(x_i) - y_i), 2(M_\alpha(x_i) - y_i)x_i, 2(M_\alpha(x_i) - y_i)x_i^2).$$

That is, the step is now just the sum of the terms from each data point. Since f is convex, then we can apply all of the nice convergence results discussed about (strongly) convex f before. However, computing $\nabla f(\alpha)$ each step takes $O(n)$ time, which can be slow.

Stochastic gradient descent. The second technique is called *incremental gradient descent*. It avoids computing the full gradient each step, and only computes it for f_i for a single data point $p_i \in P$; see algorithm 6.4.1.

Algorithm 6.4.1 Incremental Gradient Descent(f, x_{start})

initialize $x^{(0)} = x_{\text{start}} \in \mathbb{R}^d; i = 1.$

repeat

$x^{(k+1)} := x^{(k)} - \gamma_k \nabla f_i(x^{(k)})$

$i := (i + 1) \bmod n$

until ($\|\nabla f(x^{(k)})\| \leq \tau$) (or perhaps average of a few steps)

return $x^{(k)}$

A more common variant of this is called *stochastic gradient descent*. Instead of choosing the data points in order, it selects a data point p_i at random each iteration (the term “stochastic” refers to this randomness).

These algorithms are often much faster than the batch version since each iteration now takes $O(1)$ time. However, it does not automatically inherit all of the nice convergence results from what is known about (strongly) convex functions. Yet, since in most settings we are interested in, there is an abundance of data points from the same model. This implies they should have a roughly similar effect. In practice when one is far from the optimal model, these steps converge about as well as the batch version (but much much faster). When one is close to the optimal model, then the incremental / stochastic variants may not exactly converge. However, if one is willing to reach a point that is not too optimal, there are some randomized (PAC-style) guarantees possible for the stochastic variant. And in fact, for very large data sets (n is very big) they typically converge before the algorithm even uses all (or even most) of the data points.

7 Principal Component Analysis

This topic will build a series of techniques to deal with high-dimensional data. Unlike regression problems, our goal is not to predict a value (the y -coordinate), it is to understand the “shape” of the data, for instance a low-dimensional representation that captures most of meaning of the high-dimensional data. This is sometimes referred to as *unsupervised learning* (as opposed to regression and classification, where the data has labels, known as supervised learning). Like most unsupervised settings, it can be a lot of fun, but its easy to get yourself into trouble if you are not careful.

We will cover many interconnected tools, including the singular value decomposition (SVD), eigenvectors and eigenvalues, the power method, principal component analysis, and multidimensional scaling.

7.1 Data Matrices

We will start with data in a matrix $A \in \mathbb{R}^{n \times d}$, and will call upon linear algebra to rescue us. It is useful to think of each row a_i of A as a data point in \mathbb{R}^d , so there are n data points. Each dimension $j \in 1, 2, \dots, d$ corresponds with an attribute of the data points.

Example: Data Matrices

There are many situations where data matrices arise.

- Consider a set of n weather stations reporting temperature over d points in time. Then each row a_i corresponds to a single weather station, and each coordinate $A_{i,j}$ of that row is the temperature at station i at time j .
- In movie ratings, we may consider n users who have rated each of d movies on a score of 1 – 5. Then each row a_i represents a user, and the j th entry of that user is the score given to the j movie.
- Consider the price of a stock measured over time (say the closing price each day). Many time-series models consider some number of days (d days, for instance 25 days) to capture the pattern of the stock at any given time. So for a given closing day, we consider the d previous days. If we have data on the stock for 4 years (about 1000 days the stock market is open), then we can create a d -dimensional data points (the previous $d = 25$ days) for each day (except the first 25 or so). The data matrix is then comprised of n data points a_i , where each corresponds to the closing day, and the previous d days. The j th entry is the value on $(j - 1)$ days before the closing day i .
- Finally consider a series of pictures of a shape (say the Utah teapot). The camera position is fixed as is the background, but we vary two things: the rotation of the teapot, and the amount of light. Here each pictures is a set of say d pixels (say 10,000 if it is 100×100), and there are n pictures. Each picture is a row of length d , and each pixel corresponds to a column of the matrix. Similar, but more complicated scenarios frequently occur with pictures of a persons face, or $3d$ -imaging of an organ.

In each of these scenarios, there are many (n) data points, each with d attributes. The following will be very important:

- **all coordinates have the same units!**

If this “same units” property does not hold, then when we measure a distance between data points in \mathbb{R}^d , usually using the L_2 -norm, then the distance is nonsensical.

The next goal is to uncover a pattern, or a model M . In this case, the model will be a low-dimensional subspace F . It will represent a k -dimensional space, where $k \ll d$. For instance in the example with images, there are only two parameters which are changing (rotation, and lighting), so despite having $d = 10,000$ dimensions of data, 2 should be enough to represent everything.

7.1.1 Projections

Different than in linear regression this family of techniques will measure error as a projection from $a_i \in \mathbb{R}^d$ to the closest point $\pi_F(a_i)$ on F . To define this we will use linear algebra.

First recall, that given a unit vector $u \in \mathbb{R}^d$ and any data point $p \in \mathbb{R}^d$, then the dot product

$$\langle u, p \rangle$$

is the norm of p projected onto the line through u . If we multiply this scalar by u then

$$\pi_u(p) = \langle u, p \rangle u,$$

and it results in the point on the line through u that is closest to data point p . This is a *projection of p onto u* .

To understand this for a subspace F , we will need to define a basis. *For now we will assume that F contains the origin $(0, 0, 0, \dots, 0)$ (as did the line through u).* Then if F is k -dimensional, then this means there is a k -dimensional basis $U_F = \{u_1, u_2, \dots, u_k\}$ so that

- For each $u_i \in U_F$ we have $\|u_i\| = 1$, that is u_i is a unit vector.
- For each pair $u_i, u_j \in U_F$ we have $\langle u_i, u_j \rangle = 0$; the pairs are orthogonal.
- For any point $x \in F$ we can write $x = \sum_{i=1}^k \alpha_i u_i$; in particular $\alpha_i = \langle x, u_i \rangle$.

Given such a basis, then the projection on to F of a point $p \in \mathbb{R}^d$ is simply

$$\pi_F(p) = \sum_{i=1}^k \langle u_i, p \rangle u_i.$$

Thus if p happens to be exactly in F , then this recovers p exactly.

The other powerful part of the basis U_F is the it defines a *new coordinate system*. Instead of using the d original coordinates, we can use new coordinates $(\alpha_1(p), \alpha_2(p), \dots, \alpha_k(p))$ where $\alpha_i(p) = \langle u_i, p \rangle$. To be clear $\pi_F(p)$ is still in \mathbb{R}^d , but there is a k -dimensional representation if we restrict to F .

When F is d -dimensional, this operation can still be interesting. The basis we choose $U_F = \{u_1, u_2, \dots, u_d\}$ could be the same as the original coordinate axis, that is we could have $u_i = e_i = (0, 0, \dots, 0, 1, 0, \dots, 0)$ where only the i th coordinate is 1. But if it is another basis, then this acts as a rotation (with possibility of also a mirror flip). The first coordinate is rotated to be along u_1 ; the second along u_2 ; and so on. In $\pi_F(p)$, the point p does not change, just its representation.

7.1.2 SSE Goal

As usual our goal will be to minimize the sum of squared errors. In this case we define this as

$$\text{SSE}(A, F) = \sum_{a_i \in A} \|a_i - \pi_F(a_i)\|^2,$$

and our desired k -dimensional subspace F is

$$F^* = \arg \min_F \text{SSE}(A, F)$$

As compared to linear regression, this is much less a “proxy goal” where the true goal was prediction. Now we have no labels (the y_i values), so we simply try to fit a model through all of the data.

How do we solve for this?

- Linear regression does not work, its cost function is different.
- It is not obvious how to use gradient descent. The restriction that each $u_i \in U_F$ is a unit vector puts in a constraint, in fact a non-convex one. There are ways to deal with this, but we have not discussed these yet.
- ... linear algebra will come back to the rescue, now in the form of the SVD.

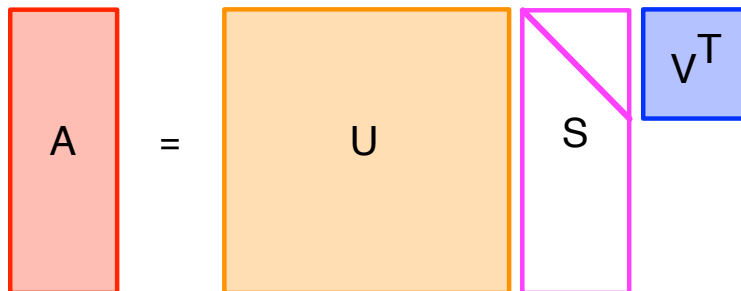
7.2 Singular Value Decomposition

A really powerful and useful linear algebra operation is called the *singular value decomposition*. It extracts an enormous amount of information about a matrix A . This section will define it and discuss many of its uses. Then we will describe one algorithm how to construct it. But in general, one simply calls the procedure in your favorite programming language and it calls the same highly-optimized back-end from the Fortran LAPACK library.

```
from scipy import linalg as LA
U, s, Vt = LA.svd(A)
```

The SVD takes in a matrix $A \in \mathbb{R}^{n \times d}$ and outputs three matrices $U \in \mathbb{R}^{n \times n}$, $S \in \mathbb{R}^{n \times d}$ and $V \in \mathbb{R}^{d \times d}$, so that $A = USV^T$.

$$[U, S, V] = \text{svd}(A)$$



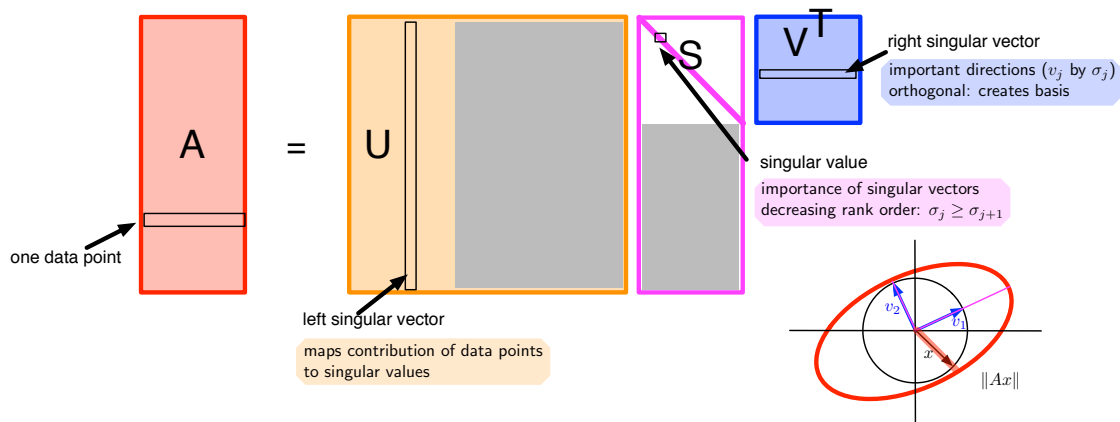
The structure that lurks beneath. The matrix S only has non-zero elements along its diagonal. So $S_{i,j} = 0$ if $i \neq j$. The remaining values $\sigma_1 = S_{1,1}$, $\sigma_2 = S_{2,2}$, ..., $\sigma_r = S_{r,r}$ are known as the *singular values*. They have the property that

$$\sigma_1 \geq \sigma_2 \geq \dots \sigma_r \geq 0$$

where $r \leq \min\{n, d\}$ is the rank of the matrix A . So the number of non-zero singular values reports the rank (this is a numerical way of computing the rank of a matrix).

The matrices U and V are orthogonal. Thus, their columns are all unit vectors and orthogonal to each other (within each matrix). The columns of U , written u_1, u_2, \dots, u_d , are called the *left singular vectors*; and the columns of V , written v_1, v_2, \dots, v_n , are called the *right singular vectors*.

This means for any vector $x \in \mathbb{R}^d$, the columns of V (the right singular vectors) provide a basis. That is, we can write $x = \sum_{i=1}^d \alpha_i v_i$ for $\alpha_i = \langle x, v_i \rangle$. Similarly for any vector $y \in \mathbb{R}^n$, the columns of U (the left singular vectors) provide a basis. This also implies that $\|x\| = \|V^T x\|$ and $\|y\| = \|yU\|$.



Tracing the path of a vector. To illustrate what this decomposition demonstrates, a useful exercise is to trace what happens to a vector $x \in \mathbb{R}^d$ as it is left-multiplied by A , that is $Ax = USV^T x$.

First $V^T x$ produces a new vector $\xi \in \mathbb{R}^d$. It essentially changes no information, just changes the basis to that described by the right singular values. For instance the new i coordinate $\xi_i = \langle v_i, x \rangle$.

Next $\eta \in \mathbb{R}^n$ is the result of $SV^T x = S\xi$. It scales ξ by the singular values of S . Note that if $d < n$ (the case we will focus on), then the last $n - d$ coordinates of η are 0. In fact, for $j > r$ (where $r = \text{rank}(A)$) then $\eta_j = 0$. For $j \leq r$, then the vector η is stretched longer in the first coordinates since these have larger values.

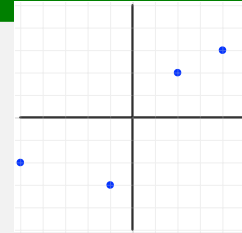
The final result is a vector $y \in \mathbb{R}^n$, the result of $Ax = USV^T x = U\eta$. This again just changes the basis of η so that it aligns with the left singular vectors. In the setting $n > d$, the last $n - d$ left singular vectors are meaningless since the corresponding entries in η are 0.

Working backwards ... this final U matrix can be thought of mapping the effect of η onto each of the data points of A . The η vector, in turn, can be thought of as scaling by the content of the data matrix A (the U and V^T matrices contain no scaling information). And the ξ vector arises via the special rotation matrix V^T that puts the starting point x into the right basis to do the scaling (from the original d -dimensional coordinates to one that suits the data better).

Example: Tracing through the SVD

Consider a matrix

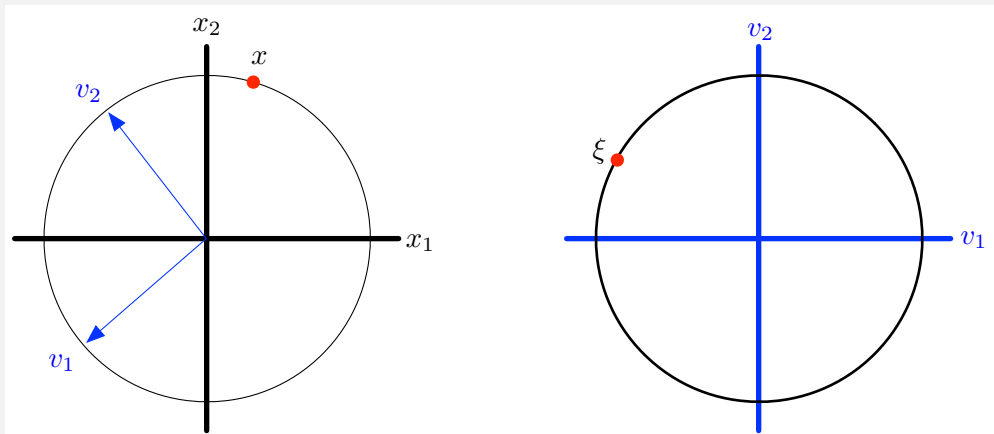
$$A = \begin{pmatrix} 4 & 3 \\ 2 & 2 \\ -1 & -3 \\ -5 & -2 \end{pmatrix},$$



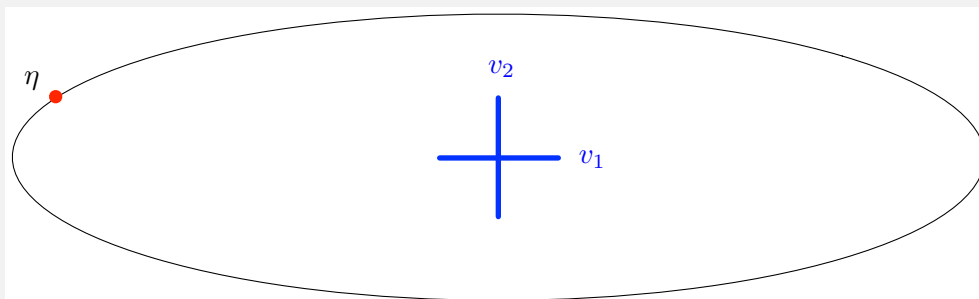
and its SVD $[U, S, V] = \text{svd}(A)$:

$$U = \begin{pmatrix} -0.6122 & 0.0523 & 0.0642 & 0.7864 \\ -0.3415 & 0.2026 & 0.8489 & -0.3487 \\ 0.3130 & -0.8070 & 0.4264 & 0.2625 \\ 0.6408 & 0.5522 & 0.3057 & 0.4371 \end{pmatrix}, \quad S = \begin{pmatrix} 8.1655 & 0 \\ 0 & 2.3074 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}, \quad V = \begin{pmatrix} -0.8142 & -0.5805 \\ -0.5805 & 0.8142 \end{pmatrix}.$$

Now consider a vector $x = (0.243, 0.97)$ (scaled very slightly so it is a unit vector, $\|x\| = 1$). Multiplying by V^T rotates (and flips) x to $\xi = V^T x$; still $\|\xi\| = 1$



Next multiplying by S scales ξ to $\eta = S\xi$. Notice there are an imaginary third and fourth coordinates now; they are both coming out of the page! Don't worry, they won't poke you since their magnitude is 0.



Finally, $y = U\eta = Ax$ is again another rotation of η in this four dimensional space.

```

import scipy as sp
import numpy as np
from scipy import linalg as LA

A = np.array([[4.0,3.0], [2.0,2.0], [-1.0,-3.0], [-5.0,-2.0]])

U, s, Vt = LA.svd(A, full_matrices=False)

print U
#[[-0.61215255 -0.05228813]
# [-0.34162337 -0.2025832 ]
# [ 0.31300005  0.80704816]
# [ 0.64077586 -0.55217683]]
print s
#[ 8.16552039  2.30743942]
print Vt
#[[-0.81424526 -0.58052102]
# [ 0.58052102 -0.81424526]]

x = np.array([0.243,0.97])
x = x/LA.norm(x)

xi = Vt.dot(x)
print xi
#[-0.7609864 -0.64876784]

S = LA.diagsvd(s,2,2)
eta = S.dot(xi)
print eta
#[-6.21384993 -1.49699248]

y = U.dot(eta)
print y
#[ 3.88209899  2.42606187 -3.1530804 -3.15508046]

print A.dot(x)
#[ 3.88209899  2.42606187 -3.1530804 -3.15508046]

```

7.2.1 Best Rank- k Approximation

So how does this help solve the initial problem of finding F^* , which minimized the SSE? The singular values hold the key.

It turns out that there is a *unique* singular value decomposition, up to ties in the singular values. This means, there is exactly one (up to singular value ties) set of right singular values which rotate into a basis so that $\|Ax\| = \|SV^T x\|$ for all $x \in \mathbb{R}^d$ (recall that U is orthogonal, so it does not change the norm, $\|U\eta\| = \|\eta\|$).

Next we realize that the singular values come in sorted order $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$. In fact, they are defined so that we choose v_1 so it maximizes $\|Av_1\|$, then we find the next singular vector v_2 which is orthogonal to v_1 and maximizes $\|Av_2\|$, and so on. Then $\sigma_i = \|Av_i\|$.

If we define F with the basis $U_F = \{v_1, v_2, \dots, v_k\}$, then

$$\|x - \pi_F(x)\|^2 = \left\| \sum_{i=1}^d v_i \langle x, v_i \rangle - \sum_{i=1}^k v_i \langle x, v_i \rangle \right\|^2 = \sum_{i=k+1}^d \langle x, v_i \rangle^2.$$

so the projection error is that part of x in the last $(d - k)$ right singular vectors.

But we are not trying to directly predict new data here (like in regression). Rather, we are trying to approximate the data we have. We want to minimize $\sum_i \|a_i - \pi_F(a_i)\|^2$. But for any unit vector u , we recall now that

$$\|Au\|^2 = \sum_{i=1}^n \langle a_i, u \rangle.$$

Thus the projection error can be measured with a set of orthonormal vectors w_1, w_2, \dots, w_{d-k} which are each orthogonal to F , as $\sum_{j=1}^{n-k} \|Aw_j\|^2$. When defining F as the first k right singular values, then these orthogonal vectors are the remaining $(n - k)$ right singular vectors, so the projection error is

$$\sum_{i=1}^n \|a_i - \pi_F(a_i)\|^2 = \sum_{j=k+1}^d \|Av_j\|^2 = \sum_{j=k+1}^d \sigma_j^2.$$

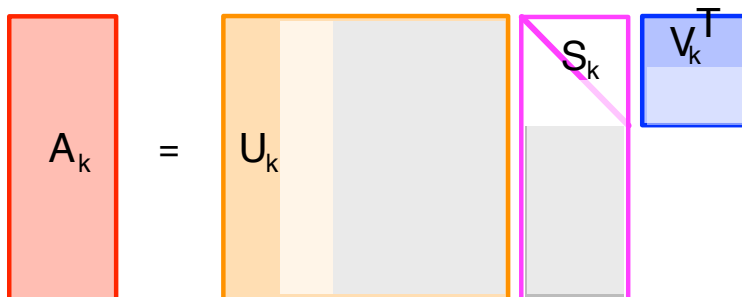
And thus by how the right singular vectors are defined, this expression is minimized when F is defined as the span of the first k singular values.

Best rank- k approximation. A similar goal is to find the *best rank- k approximation of A* . That is a matrix $A_k \in \mathbb{R}^{n \times d}$ so that $\text{rank}(A_k) = k$ and it minimizes both

$$\|A - A_k\|_2 \quad \text{and} \quad \|A - A_k\|_F.$$

Note that $\|A - A_k\|_2 = \sigma_{k+1}$ and $\|A - A_k\|_F^2 = \sum_{j=k+1}^d \sigma_j^2$.

Remarkably, this A_k matrix also comes from the SVD. If we set S_k as the matrix S in the decomposition so that all but the first k singular values are 0, then it has rank k . Hence $A_k = US_kV^T$ also has rank k and is our solution. But we can notice that when we set most of S_k to 0, then the last $(d - k)$ columns of V are meaningless since they are only multiplied by 0s in US_kV^T , so we can also set those to all 0s, or remove them entirely (along with the last $(d - k)$ columns of S_k). Similar we can make 0 or remove the last $(n - k)$ columns of U . These matrices are referred to as V_k and U_k respectively, and also $A_k = U_kS_kV_k^T$.



7.3 Eigenvalues and Eigenvectors

A related matrix decomposition to SVD is the eigendecomposition. This is only defined for a square matrix $B \in \mathbb{R}^{n \times n}$.

An *eigenvector* of B is a vector v such that there is some scalar λ that

$$Bv = \lambda v.$$

That is, multiplying B by v results in a scaled version of v . The associated value λ is called the *eigenvalue*. As a convention, we typically normalize v so $\|v\| = 1$.

In general, a square matrix $B \in \mathbb{R}^{n \times n}$ may have up to n eigenvectors (a matrix $V \in \mathbb{R}^{n \times n}$) and values (a vector $\lambda \in \mathbb{R}^n$). Some of the eigenvalues may be complex numbers (even when all of its entries are real!).

```
from scipy import linalg as LA
l, V = LA.eig(B)
```

For this reason, we will focus on positive semidefinite matrices. A *positive definite matrix* $B \in \mathbb{R}^{n \times n}$ is a symmetric matrix with all positive eigenvalues. Another characterization is for every vector $x \in \mathbb{R}^n$ then $x^T B x$ is positive. A *positive semidefinite matrix* $B \in \mathbb{R}^{n \times n}$ may have some eigenvalues at 0 and are otherwise positive; equivalently for any vector $x \in \mathbb{R}^n$, then $x^T B x$ may be zero or positive.

How do we get positive semi-definite matrices? Lets start with a data matrix $A \in \mathbb{R}^{n \times d}$. Then we can construct two positive semidefinite matrices

$$B_R = A^T A \quad \text{and} \quad B_L = A A^T.$$

Matrix B_R is $d \times d$ and B_L is $n \times n$. If the rank of A is d , then B_R is positive definite. If the rank of A is n , then B_L is positive definite.

Eigenvectors and eigenvalues relation to SVD. Next consider the SVD of A so that $[U, S, V] = \text{svd}(A)$. Then we can write

$$B_R V = A^T A V = (V S U^T)(U S V^T) V = V S^2.$$

Note that the last step follows because for orthogonal matrices U and V , then $U^T U = I$ and $V^T V = I$, where I is the identity matrix, which has no effect. The matrix S is a diagonal square¹ matrix $S = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_d)$. Then $S^2 = S S$ (the product of S with S) is again diagonal with entries $S^2 = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2)$.

Now consider a single column v_i of V (which is the i th right singular vector of A). Then extracting this column's role in the linear system $B_R V = V S^2$ we obtain

$$B_R v_i = v_i \sigma_i^2.$$

This means that i th right singular vector of A is an eigenvector (in fact the i th eigenvector) of $B_R = A^T A$. Moreover, the i th eigenvalue λ_i of B_R is the i th singular value of A squared: $\lambda_i = \sigma_i^2$.

Similarly we can derive

$$B_L U = A A^T U = (U S V^T)(V S U^T) U = U S^2,$$

and hence the left singular vectors of A are the eigenvectors of $B_L = A A^T$ and the eigenvalues of B_L are the squared singular values of A .

Eigendecomposition. In general, the eigenvectors provide a basis for a matrix $B \in \mathbb{R}^{n \times n}$ in the same way that the right V or left singular vectors U provide a basis for matrix $A \in \mathbb{R}^{n \times d}$. In fact, it is again a very special basis, and is unique up to the multiplicity of eigenvalues. This implies that all eigenvectors are orthogonal to each other.

Let $V = [v_1, v_2, \dots, v_n]$ be the eigenvectors of the matrix $B \in \mathbb{R}^{n \times n}$, as columns in the matrix V . Also let $L = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$ be the eigenvalues of B stored on the diagonal of matrix L . Then we can decompose B as

$$B = V L V^{-1}.$$

¹Technically, $S \in \mathbb{R}^{n \times d}$. To make this simple argument work, lets first assume w.l.o.g. (without loss of generality) that $d \leq n$. Then the bottom $n - d$ rows of S are all zeros, which mean the right $n - d$ rows of U do not matter. So we can ignore both these $n - d$ rows and columns. Then S is square. This makes U no longer orthogonal, so $U^T U$ is then a projection, not identity; but it turns out this is a project to the span of A , so the argument still works.

Note that the inverse of L is $L^{-1} = \text{diag}(1/\lambda_1, 1/\lambda_2, \dots, 1/\lambda_n)$. Hence we can write

$$B^{-1} = VL^{-1}V^{-1}.$$

When B is positive definite, it has n positive eigenvectors and eigenvalues; hence V is orthogonal, so $V^{-1} = V^T$. Thus in this situation, given the eigendecomposition, we now have a way to compute the inverse

$$B^{-1} = VL^{-1}V^T,$$

which was required in our almost closed-form solution for linear regression. Now we just need to compute the eigendecomposition, which we will discuss next.

7.4 The Power Method

The *power method* refers to what is probably the simplest algorithm to compute the first eigenvector and value of a matrix. By factoring out the effect of the first eigenvector, we can then recursively repeat the process on the remainder until we have found all eigenvectors and values. Moreover, this implies we can also reconstruct the singular value decomposition as well.

We will consider $B \in \mathbb{R}^{n \times n}$, a positive semidefinite matrix: $B = A^T A$.

Algorithm 7.4.1 PowerMethod(B, q)

initialize $u^{(0)}$ as a random unit vector.

for $i = 1$ **to** q **do**

$u^{(i)} := Bu^{(i-1)}$

return $v := u^{(q)} / \|u^{(q)}\|$

We can unroll the for loop to reveal another interpretation. We can directly set $v^{(q)} = B^q v^{(0)}$, so all iterations are incorporated into one matrix-vector multiplication. Recall that $B^q = B \cdot B \cdot B \cdot \dots \cdot B$, for q times. However, these q matrix multiplications are much more expensive than q matrix-vector multiplications.

Alternatively we are provided only the matrix A (where $B = A^T A$) then we can run the algorithm without explicitly constructing B (since for instance if $d > n$ and $A \in \mathbb{R}^{n \times d}$, then the size of B (d^2) may be much larger than A (nd)). Then we simply replace the inside of the for-loop with

$$u^{(i)} := A^T(Au^{(i-1)})$$

where we first multiply $\tilde{u} = Au^{(i-1)}$ and then complete $u^{(i)} = A^T \tilde{u}$.

Recovering all eigenvalues. The output of PowerMethod($B = A^T A, q$) is a single unit vector v , which we will argue is arbitrarily close to the first eigenvector v_1 . Clearly we can recover the first eigenvalue as $\lambda_1 = \|Bv_1\|$. Since we know the eigenvectors form a basis for B , they are orthogonal. Hence, after we have constructed the first eigenvector v_1 , we can factor it out from B as follows:

$$A_1 := A - Av_1v_1^T$$

$$B_1 := A_1^T A_1$$

Then we run PowerMethod($B_1 = A_1^T A_1, q$) to recover v_2 , and λ_2 ; factor them out of B_1 to obtain B_2 , and iterate.

Why does the power method work? To understand why the power method works, assume we know the eigenvectors v_1, v_2, \dots, v_n and eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ of $B \in \mathbb{R}^{n \times n}$.

Since the eigenvectors form a basis for B , and assuming it is full rank, then also for all of \mathbb{R}^n (if not, then it does not have n eigenvalues, and we can fill out the rest of the basis of \mathbb{R}^n arbitrarily). Hence, for any vector, including the initialization random vector $u^{(0)}$ can be written as

$$u^{(0)} = \sum_{j=1}^n \alpha_j v_j.$$

Recall that $\alpha_j = \langle u^{(0)}, v_j \rangle$, and since it is random, it is possible to claim that with probability at least $1/2$ that for any α_j we have that $|\alpha_j| \geq \frac{1}{2}\sqrt{n}^{-2}$. We will now assume that this holds for $j = 1$, so $\alpha_1 > 1/2\sqrt{n}$.

Next since we can interpret that algorithm as $v = B^q u^{(0)}$, then lets analyze B^q . If B has j th eigenvector v_j and eigenvalue λ_j , that is, $Bv_j = \lambda_j v_j$, then B^q has j th eigenvalue λ_j^q since

$$B^q v_j = B \cdot B \cdot \dots \cdot B v_j = B^{q-1}(v_j \lambda) = B^{q-2}(v_j \lambda) \lambda = v_j \lambda^q.$$

This holds for each eigenvalue of B^q . Hence we can rewrite output by summing over the terms in the eigenbasis as

$$v = \frac{\sum_{j=1}^n \alpha_j \lambda_j^q v_j}{\sqrt{\sum_{j=1}^n (\alpha_j \lambda_j^q)^2}}.$$

Finally, we would like to show our output v is close to the first eigenvector v_1 . We can measure closeness with the dot product (actually we will need to use its absolute value since we might find something close to $-v_1$).

$$\begin{aligned} |\langle B^q u^{(0)}, v_1 \rangle| &= \frac{\alpha_1 \lambda_1^q}{\sqrt{\sum_{j=1}^n (\alpha_j \lambda_j^q)^2}} \\ &\geq \frac{\alpha_1 \lambda_1^q}{\sqrt{\alpha_1^2 \lambda_1^{2q} + n \lambda_2^{2q}}} \geq \frac{\alpha_1 \lambda_1^q}{\alpha_1 \lambda_1^q + \lambda_2^q \sqrt{n}} = 1 - \frac{\lambda_2^q \sqrt{n}}{\alpha_1 \lambda_1^q + \lambda_2^q \sqrt{n}} \\ &\geq 1 - 2\sqrt{n} \left(\frac{\lambda_2}{\lambda_1} \right)^q. \end{aligned}$$

The first inequality holds because $\lambda_1 \geq \lambda_2 \geq \lambda_j$ for all $j > 2$. The third inequality (going to third line) holds by dropping the $\lambda_2^q \sqrt{n}$ term in the denominator, and since $\alpha_1 > 1/2\sqrt{n}$.

Thus if there is “gap” between the first two eigenvalues (λ_1/λ_2 is large), then this algorithm converges quickly to where $|\langle v, v_1 \rangle| = 1$.

7.5 Principal Component Analysis

Recall that the original goal of this topic was to find the k -dimensional subspace F to minimize

$$\|A - \pi_F(A)\|_F^2 = \sum_{a_i \in A} \|a_i - \pi_F(a_i)\|^2.$$

²Since $u^{(0)}$ is a unit vector, its norm is 1, and because $\{v_1, \dots, v_n\}$ is a basis, then $1 = \|u^{(0)}\|^2 = \sum_{j=1}^n \alpha_j^2$. Since it is random, then $\mathbf{E}[\alpha_j^2] = 1/n$ for each j . Applying a concentration of measure (almost a Markov Inequality, but need to be more careful), we can argue that with probability $1/2$ any $\alpha_j^2 > (1/4) \cdot (1/n)$, and hence $\alpha_j > (1/2) \cdot (1/\sqrt{n})$.

We have not actually solved this yet. The top k right singular values V_k of A only provided this bound assuming that F contains the origin: $(0, 0, \dots, 0)$. However, this might not be the case!

Principal Component Analysis (PCA) is an extension of the SVD when we do not restrict that the subspace V_k must go through the origin. It turns out, like with simple linear regression, that the optimal F must go through the mean of all of the data. So we can still use the SVD, after a simple preprocessing step called centering to shift the data matrix so its mean is exactly at the origin.

Specifically, *centering* is adjusting the original input data matrix $A \in \mathbb{R}^{n \times d}$ so that each column (each dimension) has an average value of 0. This is easier than it seems. Define $\bar{a}_j = \frac{1}{n} \sum_{i=1}^n A_{i,j}$ (the average of each column j). Then set each $\tilde{A}_{i,j} = A_{i,j} - \bar{a}_j$ to represent the entry in the i th row and j th column of centered matrix \tilde{A} .

There is a *centering matrix* $C_n = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$ where I_n is the $n \times n$ identity matrix, $\mathbf{1}$ is the all-ones column vector (of length n) and thus $\mathbf{1}\mathbf{1}^T$ is the all-ones $n \times n$ matrix. Then we can also just write $\tilde{A} = C_n A$.

Now to perform PCA on a data set A , we compute $[U, S, V] = \text{svd}(C_n A) = \text{svd}(\tilde{A})$.

Then the resulting singular values $\text{diag}(S) = \{\sigma_1, \sigma_2, \dots, \sigma_r\}$ are known as the *principal values*, and the top k right singular vectors $V_k = [v_1 \ v_2 \ \dots \ v_k]$ are known as the top- k *principal directions*.

This often gives a better fitting to the data than just SVD. The SVD finds the best rank- k approximation of A , which is the best k -dimensional subspace (up to Frobenius and spectral norms) **which passes through the origin**. If all of the data is far from the origin, this can essentially “waste” a dimension to pass through the origin. However, we also need to store the shift from the origin, a vector $\tilde{c} = (\tilde{a}_1, \tilde{a}_2, \dots, \tilde{a}_d) \in \mathbb{R}^d$.

7.6 Multidimensional Scaling

Dimensionality reduction is an abstract problem with input of a high-dimensional data set $P \subset \mathbb{R}^d$ and a goal of finding a corresponding lower dimensional data set $Q \subset \mathbb{R}^k$, where $k \ll d$, and properties of P are preserved in Q . Both low-rank approximations through direct SVD and through PCA are examples of this: $Q = \pi_{V_k}(P)$. However, these techniques require an explicit representation of P to start with. In some cases, we are only presented P more abstractly. There two common situations:

- We are provided a set of n objects X , and a bivariate function $d : X \times X \rightarrow \mathbb{R}$ that returns a distance between them. For instance, we can put two cities into an airline website, and it may return a dollar amount for the cheapest flight between those two cities. This dollar amount is our “distance.”
- We are simply provided a matrix $D \in \mathbb{R}^{n \times n}$, where each entry $D_{i,j}$ is the distance between the i th and j th point. In the first scenario, we can calculate such a matrix D .

Multi-Dimensional Scaling (MDS) has the goal of taking such a distance matrix D for n points and giving low-dimensional (typically) Euclidean coordinates to these points so that the embedded points have similar spatial relations to that described in D . If we had some original data set A which resulted in D , we could just apply PCA to find the embedding. It is important to note, in the setting of MDS we are typically just given D , and *not* the original data A . However, as we will show next, we can derive a matrix that will act like AA^T using only D .

A *similarity matrix* M is an $n \times n$ matrix where entry $M_{i,j}$ is the similarity between the i th and the j th data point. The similarity often associated with Euclidean distance $\|a_i - a_j\|$ is the standard inner (or dot product) $\langle a_i, a_j \rangle$. We can write

$$\|a_i - a_j\|^2 = \|a_i\|^2 + \|a_j\|^2 - 2\langle a_i, a_j \rangle,$$

and hence

$$\langle a_i, a_j \rangle = \frac{1}{2} (\|a_i\|^2 + \|a_j\|^2 - \|a_i - a_j\|^2). \quad (7.1)$$

Next we observe that for the $n \times n$ matrix AA^T the entry $[AA^T]_{i,j} = \langle a_i, a_j \rangle$. So it seems hopeful we can derive AA^T from D using equation (7.1). That is we can set $\|a_i - a_j\|^2 = D_{i,j}^2$. However, we need also need values for $\|a_i\|^2$ and $\|a_j\|^2$.

Since the embedding has an arbitrary shift to it (if we add a shift vector s to *all* embedding points, then no distances change), then we can arbitrarily choose a_1 to be at the origin. Then $\|a_1\|^2 = 0$ and $\|a_j\|^2 = \|a_1 - a_j\|^2 = D_{1,j}^2$. Using this assumption and equation (7.1), we can then derive the similarity matrix AA^T . Then we can run the eigen-decomposition on AA^T and use the coordinates of each point along the first k eigenvectors to get an embedding. This is known as *classical MDS*.

It is often used for k as 2 or 3 so the data can be easily visualized.

There are several other forms that try to preserve the distance more directly, where as this approach is essentially just minimizing the squared residuals of the projection from some unknown original (high-dimensional embedding). One can see that we recover the distances with no error if we use all n eigenvectors – if they exist. However, as mentioned, there may be less than n eigenvectors, or they may be associated with complex eigenvalues. So if our goal is an embedding into $k = 3$ or $k = 10$, there is no guarantee that this will work, or even what guarantees this will have. But MDS is used a lot nonetheless.

8 Clustering

This topic will focus on automatically grouping data points into subsets of similar points. There are numerous ways to define this problem, and most of them are quite messy. And many techniques for clustering actually lack a mathematical formulation. We will focus on what is probably the cleanest and most used formulation: k -means clustering. But, for background, we will begin with a mathematical detour in Voronoi diagrams.

8.1 Voronoi Diagrams

Consider a set $S = \{s_1, s_2, \dots, s_k\} \subset \mathbb{R}^d$ of sites. We would like to understand how these points carve up the space \mathbb{R}^d .

We can think of this more formally as the *post office problem*. Let these sites define the locations of a post office. For all points in \mathbb{R}^d (e.g., a point on the map for points in \mathbb{R}^2), we would like to assign it to the closest post office. For a fixed point, we can just check the distance to each post office:

$$\phi_S(x) = \arg \min_{s_i \in S} \|x - s_i\|.$$

However, this may be slow (naively take $O(k)$ time for each point x), and does not provide a general representation or understanding for all points. The “correct” solution to this problem is the Voronoi diagram.

The Voronoi diagram decomposes \mathbb{R}^d into k regions (a *Voronoi cell*), one for each site. The region for site s_i is defined.

$$R_i = \{x \in \mathbb{R}^d \mid \phi_S(x) = s_i\}$$

If we have these regions nicely defined, this solves the post office problem. For any point x , we just need to determine which region it lies in (for instance in \mathbb{R}^2 , once we have defined these regions, through an extension of binary search, we can locate the region containing any $x \in \mathbb{R}^2$ in only $O(\log k)$ time). But what do these regions look like, and what properties do they have.

Voronoi edges and vertices. We will start our discussion in \mathbb{R}^2 . Further, we will assume that the sites S are in *general position*: in this setting, it means that no set of three points lie on a common line, and that no set of four points lie on a common circle.

The boundary between two regions R_i and R_j , called a *Voronoi edge*, is a line or line segment. This edge $e_{i,j}$ is defined

$$e_{i,j} = \{x \in \mathbb{R}^2 \mid \|x - s_i\| = \|x - s_j\| \leq \|x - s_\ell\| \text{ for all } \ell \neq i, j\}$$

as the set of all points equal distance to s_i and s_j , and not closer to any other point s_ℓ .

Why is this set a line segment? If we only have two points in S , then it is the bisector between them. Draw a circle centered at any point x on this bisector, and if it intersects one of s_i or s_j , it will also intersect the other. This is true since we can decompose the squared distance from x to s_i along orthogonal components: along the edge, and perpendicular to the edge from s_i to $\pi_{e_{i,j}}(s_i)$.

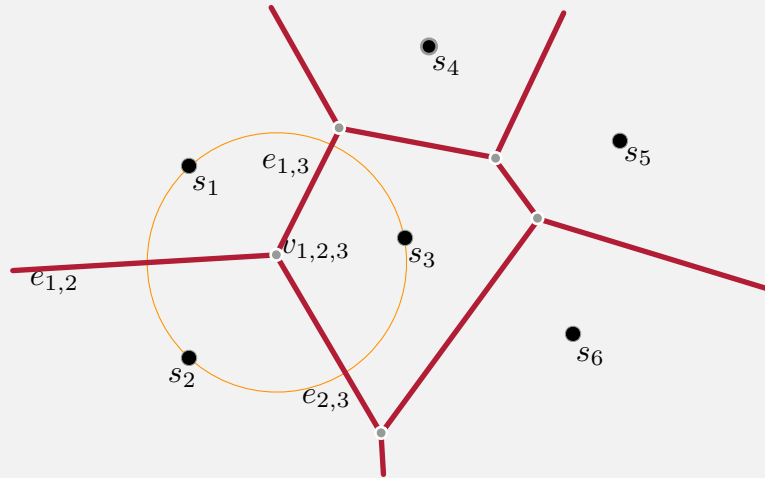
Similarly, a Voronoi vertex $v_{i,j,\ell}$ is a point where three sites s_i , s_j , and s_ℓ are all equidistance, and no other points are closer:

$$v_{i,j,k} = \{x \in \mathbb{R}^2 \mid \|x - s_i\| = \|x - s_j\| = \|x - s_\ell\| \leq \|x - s_k\| \text{ for all } k \neq i, j, \ell\}.$$

This vertex is the intersection (and end point) of three Voronoi edges $e_{i,i}$, $e_{i,\ell}$, and $e_{j,\ell}$. Think of sliding a point x along an edge $e_{i,j}$ and maintaining the circle centered at x and touching s_i and s_j . When this circle grows to where it also touches s_ℓ , then $e_{i,j}$ stops.

Example: Voronoi Diagram

See the following example with $k = 6$ sites in \mathbb{R}^2 . Notice the following properties: edges may be unbounded, and the same with regions. The circle centered at $v_{1,2,3}$ passes through s_1 , s_2 , and s_3 . Also, Voronoi cell R_3 has $5 = k - 1$ vertices and edges.



Size complexity. So how complicated can these Voronoi diagrams get? A single Voronoi cell can have $k - 1$ vertices and edges. So can the entire complex be of size $O(k^2)$ (each of k regions requiring complexity $O(k)$)? No. The Voronoi vertices and edges describe a planar graph. Some cool results from graph theory says that planar graphs have asymptotically the same number of edges, faces, and vertices. Euler's Formula for a planar graph with n vertices, m edges, and k faces is that $k + n - m = 2$. And Kuratowski's criteria says for $n \geq 3$, then $m \leq 3n - 6$. Hence, $k \leq 2n - 4$ for $n \geq 3$. The duality construction to Delaunay triangulations (discussed below) will complete the argument. Since there are k faces (the k Voronoi cells, one for each site), then there are also $O(k)$ edges and $O(k)$ vertices. In particular, there will be precisely $2n - 5$ vertices and $3k - 6$ edges.

However, this does not hold in \mathbb{R}^3 . In particular, for \mathbb{R}^3 and \mathbb{R}^4 , the complexity (number of cells, vertices, edges, faces, etc) is $O(k^2)$. This means, there could be roughly as many edges as their are pairs of vertices!

But it can get much worse. In \mathbb{R}^d (for general d) then the complexity is $O(k^{\lceil d/2 \rceil})$. This is a lot. Hence, this structure is impractical to construct in high dimensions.

: *The curse of dimensionality! ooooh*

Moreover, since this structure is explicitly tied to the post office problem, and the nearest neighbor function ϕ_S , it indicates that (a) in \mathbb{R}^2 this function is nicely behaved, but (b) in high dimensions, it is quite complicated.

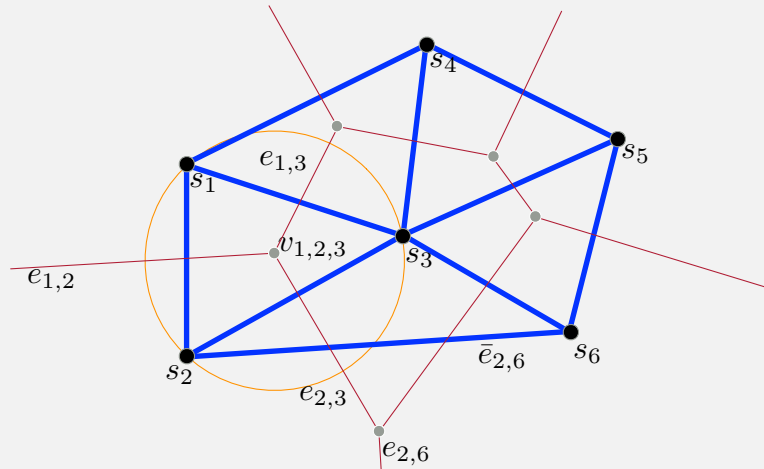
8.1.1 Delaunay Triangulation

A fascinating aspect of the Voronoi diagram is that it can be converted into a very special graph where the sites S are vertices, one called the *Delaunay triangulation*. This is the *dual* of the Voronoi diagram.

- Each face R_i of the Voronoi diagram maps to a vertex s_i in the Delaunay triangulation.
- Each vertex $v_{i,j,\ell}$ in the Voronoi diagram maps to a triangular face $f_{i,j,\ell}$ in the Delaunay triangulation.
- Each edge $e_{i,j}$ in the Voronoi diagram maps to an edge $\bar{e}_{i,j}$ in the Delaunay triangulation.

Example: Delaunay Triangulation

See the following example with 6 sites in \mathbb{R}^2 . Notice that every edge, face, and vertex in the Delaunay triangulation corresponds to a edge, vertex, and face in the Voronoi diagram. Interestingly, the associated edges may not intersect; see $e_{2,6}$ and $\bar{e}_{2,6}$.



Because of the duality between the Voronoi diagram and the Delaunay triangulation, their complexities are the same. That means the Voronoi diagram is of size $O(k)$ for k sites in \mathbb{R}^2 , but more generally is of size $O(k^{\lceil d/2 \rceil})$ in \mathbb{R}^d .

The existence of the Delaunay triangulation shows that there always exist a triangulation: A graph with vertices of a given set of points $S \subset \mathbb{R}^2$ so that all edges are straightline segments between the vertices, and each face is a triangle. In fact, there are many possible triangulations: one can always simply construct some triangulation greedily, draw any possible edges that does not cross other edges until no more can be drawn.

The Delaunay triangulation, however, is quite special. This is the triangulation that maximizes the smallest angle over all triangles; for meshing applications in graphics and simulation, skinny triangles are very hard to deal with, so this is very useful.

In circle property. Another cool way to define the Delaunay triangulation is through the *in circle* property. For any three points, the smallest enclosing ball either has all three points on the boundary, or has two points on the boundary and they are antipodal to each other. Any circle with two points antipodal on the boundary s_i and s_j , and contains no other points, then the edge $e_{i,j}$ is in the Delaunay triangulation. This is a subset of the Delaunay triangulation called the *Gabriel graph*.

Any circle with three points on its boundary s_i , s_j , and s_ℓ , and no points in its interior, then the face $f_{i,j,\ell}$ is in the Delaunay triangulation, as well as its three edges $e_{i,j}$, $e_{i,\ell}$ and $e_{j,\ell}$. But does not imply those edges satisfy the Gabriel property.

For instance, on a quick inspection, (in the example above) it may not be clear if edge $e_{3,5}$ or $e_{4,6}$ should be in the Delaunay triangulation. Clearly it can not be both since they cross. But the ball with boundary through s_3 , s_4 , and s_6 would contain s_5 , so the face $f_{3,4,6}$ cannot be in the Delaunay triangulation. On the other hand the ball with boundary through s_3 , s_6 , and s_5 does not contain s_4 or any other points in S , so the face $f_{3,5,6}$ is in the Delaunay triangulation.

8.1.2 Connection to Clustering

So what is the connection to clustering? Given a large set $X \subset \mathbb{R}^d$ of size n , we would like to find a set of k sites S (post office locations) so that each point $x \in X$ is near some post office. This is a proxy problem. So given a set of sites S , determining for each $x \in X$ which site is closest is exactly determined by the Voronoi diagram.

8.2 k -Means Clustering

Probably the most famous clustering formulation is k -means. The term “ k -means” refers to a problem formulation, *not an algorithm*. There are many algorithms with aim of solving the k -means problem formulation, exactly or approximately. We will mainly focus on the most common: Lloyd’s algorithm. Unfortunately, it is common in the literature to see “the k -means algorithm,” this typically should be replaced with Lloyd’s algorithm.

k -Means is in the family of *assignment-based clustering*. Each cluster is represented by a single point, to which all other points in the cluster are “assigned.” The input is a data set X , and the output is a set of centers $S = \{s_1, s_2, \dots, s_k\}$. This implicitly defines a set of clusters where $\phi_S(x) = \arg \min_{s \in S} \|x - s\|$ (the same as in the post office problem). Then the *k -means clustering problem* is to find the set S of k clusters to minimize

$$\text{cost}(X, S) = \sum_{x \in X} \|\phi_S(x) - x\|^2.$$

So we want every point assigned to the closest site, and want to minimize the sum of the squared distance of all such assignments.

8.3 Lloyd’s Algorithm

When people think of the k -means problem, they usually think of the following algorithm. It is usually attributed to Stuart P. Lloyd from a document in 1957, although it was not published until 1982.¹

Algorithm 8.3.1 Lloyd’s Algorithm for k -Means Clustering

Choose k points $S \subset X$	<i>arbitrarily?</i>
repeat	
for all $x \in X$: assign x to X_i for $s_i = \phi_S(x)$	<i>the closest site $s \in S$ to x</i>
for all $s_i \in S$: set $s_i = \frac{1}{ X_i } \sum_{x \in X_i} x$	<i>the average of $X_i = \{x \in X \mid \phi_S(x) = s_i\}$</i>
until (the set S is unchanged)	

Convergence. If the main loop has R rounds, then this takes roughly Rnk steps (and can be made closer to $Rn \log k$ with faster nearest neighbor search in some cases). But how large is R ; that is, how many iterations do we need?

First, we can argue that the number of steps is finite. This is true since the cost function $\text{cost}(X, S)$ will always decrease. To see this, writing it as a sum over S .

$$\text{cost}(X, S) = \sum_{x \in X} \|\phi_S(x) - x\|^2 = \sum_{s_i \in S} \sum_{x \in X_i} \|s_i - x\|^2.$$

¹Apparently, the IBM 650 computer Lloyd was using in 1957 did not have enough computational power to run the (very simple) experiments he had planned. This was replaced by the IBM 701, but it did not have quite the same “quantization” functionality as the IBM 650, and the work was forgotten. Lloyd was also worried about some issues regarding the k -means problem not having a unique minimum.

Then in each step of the **repeat-until** loop, thus must decrease. The first step holds since it moves each $x \in X$ to a subset X_i with the corresponding center s_i closer to (or the same distance to) x than before. So for each x the term $\|x - s_i\|$ is reduced (or the same). The second step holds since for each inner sum $\sum_{x \in X_i} \|s_i - x\|^2$, the single point s_i which minimizes this cost is precisely the average of X_i . So reassigning s_i as described also decreases the cost (or keeps it the same).

Importantly, if the cost decreases each step, then it cannot have the same set of centers S on two different steps, since that would imply the assignment sets $\{X_i\}$ would also be the same. Thus, in order for this to happen, the cost would need to decrease after the first occurrence, and then increase to obtain the second occurrence, which is not possible.

Since, there are finite ways each set of points can be assigned to different clusters, then, the algorithm terminates in a finite number of steps.

... but in practice usually we may run for $R = 10$, or maybe $R = 20$ steps. Or check if the change in cost function is below some sufficiently small threshold.

On clusterability: When data is easily clusterable, most clustering algorithms work quickly and well. When is not easily clusterable, then no algorithm will find good clusters.

Sometimes there is a good k -means clustering, but it is not found by Lloyd's algorithm. Then we can choose new centers again (with randomness), and try again.

Initialization. The initial paper by Lloyd advocates to choose the initial partition of X into disjoint subsets X_1, X_2, \dots, X_k *arbitrarily*. However, some choices will not be very good. For instance, if we randomly place each $x \in X$ into some x_i , then (by the central limit theorem) we expect all $s_i = \frac{1}{|X_i|} \sum_{x \in X_i} x$ to all be close to the mean of the full data set $\frac{1}{|X|} \sum_{x \in X} x$.

A bit safer way to initialize the data is to choose a set $S \subset X$ at random. Since each s_i is chosen separately (not as an average of data points), there is no centering phenomenon. However, even with this initialization, we may run Lloyd's algorithm to completion, and find a sub-optimal solution (*a local minimum!*). Thus, it is usually safer to randomly re-initialize the algorithm several times (say 3-5 times) and rerun Lloyd's algorithm for each. the probability all random restarts results in a local minimum is more rare.

A more principled way to choose an initial set S is to use an algorithm like Gonzalez (for k steps, iteratively chose the point $x \in X$ furthest from all sites chosen so far), or k -means++ (for k steps, iteratively choose a point at random proportional to the squared distance to its nearest already chosen site). These are beyond the scope of this class, but offer much stronger error guarantees of various forms. But, for k -means++, it is still suggested to try random restarts.

Number of clusters So what is the right value of k ? Like with PCA, there is no perfect answer towards choosing how many dimensions the subspace should be. When k is not given to you, typically, you would run with many different values of k . Then create a plot of $\text{cost}(S, X)$ as a function of k . This **cost** will always decrease with larger k ; but of course $k = n$ is of no use. At some point, the **cost** will not decrease much between values (this implies that probably two centers are used in the same grouping of data, so the squared distance to either is similar). Then this is a good place to choose k .

8.3.1 Extensions to the k -Means Problem and Lloyd's Algorithm

Like many algorithms in this class, Lloyd's depends on the use of a SSE cost function. In particular, it works because for $X \in \mathbb{R}^d$, then

$$\frac{1}{|X|} \sum_{x \in X} x = \text{average}(X) = \arg \min_{s \in \mathbb{R}^d} \sum_{x \in X} \|s - x\|^2.$$

There are not similar properties for other costs functions, or when X is not in \mathbb{R}^d . For instance, one may want to solve the k -medians problem where one just minimizes the sum of (non-squared) distances. In particular, this has no closed form solution for $X \in \mathbb{R}^d$ for $d > 1$.

An alternative to the averaging step is to choose

$$s_i = \arg \min_{s \in X_i} d(x, s)$$

where $d(x, s)$ is an arbitrary measure (like non-squared distance) between x and s . That is, we choose an s from the set X_i . This is particularly useful when X is in a non-Euclidean metric space where averaging may not be well-defined. For the specific case where $d(x, s) = \|x - s\|$ (for the k -median problem), then this variant of the algorithm is called *k-medoids*.

Soft clustering. Sometimes it is not desirable to assign each point to exactly one cluster. Instead, we may split a point between one or more clusters, assigning a fractional value to each. This is known as *soft clustering* whereas the original formulation is known as *hard clustering*.

There are many ways to achieve a soft clustering. For instance, consider the following Voronoi diagram-based approach called based on natural neighbor interpolation (NNI). Let $V(S)$ be the Voronoi diagram of the sites S (which decomposes \mathbb{R}^d). Then construct $V(S \cup x)$ for a particular data point x ; the Voronoi diagram of the sites S with the addition of one point x . For the region R_x defined by the point x in $V(S \cup x)$, overlay it on the original Voronoi diagram $V(S)$. This region R_x will overlap with regions R_i in the original Voronoi diagram; compute the volume v_i for the overlap with each such region. Then the fractional weight for x into each site s_i is defined $w_i(x) = v_i / \sum_{i=1}^n v_i$.

We can plug any such step into Lloyd's algorithm, and then recalculate s_i as the weighted average of all points partially assigned to the i th cluster.

8.4 Mixture of Gaussians

The k -means formulation tends to define clusters of roughly equal size. The squared cost discourages points far from any center. It also, does not adapt much to the density of individual centers.

An extension is to fit each cluster X_i with a Gaussian distribution $\mathcal{G}(\mu_i, \Sigma_i)$, defined by a mean μ_i and a covariance matrix Σ_i . Recall that the pdf of a d -dimensional Gaussian distribution is defined

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

where $|\Sigma|$ is the determinant of Σ . For instance, for $d = 2$, and the standard deviation in the x -direction of X is σ_x , and in the y -direction is σ_y , and their correlation is ρ , then

$$\Sigma = \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}.$$

Now the goal is, given a parameter k , find a set of k pdfs $F = \{f_1, f_2, \dots, f_k\}$ where $f_i = f_{\mu_i, \Sigma_i}$ to maximize

$$\prod_{x \in X} \max_{f_i \in F} f_i(x),$$

or equivalently to minimize

$$\sum_{x \in X} \min_{f_i \in F} -\log(f_i(x)).$$

For the special case where when we restrict that $\Sigma_i = I$ (the identity matrix) for each mixture, then one can check that the second formulation (the log-likelihood version) is equivalent to the k -means problem.

This hints that we can adapt Lloyd's algorithm towards this problem as well. To replace the first step of the inner loop, we assign each $x \in X$ to the Gaussian which maximizes $f_i(x)$:

$$\text{for all } x \in X: \text{ assign } x \text{ to } X_i \text{ so } i = \arg \max_{i \in 1 \dots k} f_i(x).$$

But for the second step, we need to replace a simple average with an estimation of the best fitting Gaussian to a data set X_i . This is also simple. First, calculate the mean as $\mu_i = \frac{1}{|X_i|} \sum_{x \in X_i} x$. Then calculate the covariance matrix Σ_i of X_i as the sum of outer products

$$\Sigma_i = \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T.$$

This can be interpreted as calling PCA. Calculating μ_i , and subtracting from each $x \in X_i$ is the centering step. Letting $\bar{X}_i = \{x \in \mu_i \mid x \in X_i\}$, then $\Sigma_i = VS^2V^T$ where $[U, S, V] = \text{svd}(\bar{X}_i)$.

8.4.1 Expectation-Maximization

The standard way to fit a mixture of Gaussians actually uses a soft-clustering.

Each point $x \in X$ is given a weight $w_i = f_i(x) / \sum_i f_i(x)$ for its assignment to each cluster. Then the mean and covariance matrix is estimated using weight averages.

Algorithm 8.4.1 EM Algorithm for Mixture of Gaussians

Choose k points $S \subset X$	<i>arbitrarily?</i>
for all $x \in X$: set $w_i(x)$ for $s_i = \phi_S(x)$, and $w_i(x) = 0$ otherwise	
repeat	
for $i \in [1 \dots k]$ do	
Calculate $W_i = \sum_{x \in X} w_i(x)$	<i>the total weight for cluster i</i>
Set $\mu_i = \frac{1}{W_i} \sum_{x \in X} w_i(x)x$	<i>the weighted average</i>
Set $\Sigma_i = \frac{1}{W_i} \sum_{x \in X} w_i(x - \mu_i)(x - \mu_i)^T$	<i>the weighted covariance</i>
for $x \in X$ do	
for all $i \in [1 \dots k]$: set $w_i(x) = f_i(x) / \sum_i f_i(x)$	<i>partial assignments using $f_i = f_{\mu_i, \Sigma_i}$</i>
until $(\sum_{x \in X} \sum_{i=1}^k -\log(w_i(x) \cdot f_i(x)))$ has small change	

This procedure is the classic example of a framework called *expectation-maximization*. This is an alternate optimization procedure, which alternates between maximizing the probability of some model (the partial assignment step) and calculating the most likely model using expectation (the average, covariance estimating step).

But this is a much more general framework. It is particularly useful in situations (like this one) where the true optimization criteria is messy and complex, often non-convex; but it can be broken into two or more steps where each step can be solved with a (near) closed form. Or if there is no closed form, but each part is individually convex, the gradient descent can be invoked.

8.5 Mean Shift Clustering

Now for something completely different. Clustering is a very very broad field with no settled upon approach. To demonstrate this, we will quickly review an algorithm called *mean shift clustering*. This algorithm shifts

each data point individually to its weighted center of mass. It terminates when all points converge to isolated sets.

First begin with a bivariate kernel function $K : X \times X \rightarrow \mathbb{R}$ such as the (unnormalized) Gaussian kernel

$$K(x, p) = \exp(-\|x - p\|^2 / \sigma^2)$$

for some given bandwidth parameter σ . The weighted center of mass around each point $p \in X$ is then defined as

$$\mu(p) = \frac{\sum_{x \in X} K(x, p)x}{\sum_{x \in X} K(x, p)}.$$

The algorithm just shifts each point to its center of mass: $p \leftarrow \mu(p)$.

Algorithm 8.5.1 Mean Shift

repeat

for all $p \in X$: calculate $\mu(p) = \frac{\sum_{x \in X} K(x, p)x}{\sum_{x \in X} K(x, p)}$.

for all $p \in X$: set $p \leftarrow \mu(p)$.

until (the average change $\|p - \mu(p)\|$ is small)

This algorithm does not require a parameter k . However, it has other parameters, most notably the choice of kernel K and its bandwidth σ . With the Gaussian kernel (since it has infinite support, $K(x, p) > 0$ for all x, p), it will only stop when all x are at the same point. Thus the termination condition is also important. Alternatively, a different kernel with bounded support may terminate automatically (without a specific condition); for this reason (and for speed) truncated Gaussians are often used.

This algorithm not only clusters the data, but also is a key technique for *de-noising* data. This is a process that not just removes noise (as often thought of as outliers), but attempts to return point to where they should have been before being perturbed by noise – similar to mapping a point to its cluster center.

9 Classification

This topic returns to prediction. Unlike linear regression where we were predicting a numeric value, in this case we are predicting a class: winner or loser, yes or no, rich or poor, positive or negative. Ideas from linear regression can be applied here, but we will instead overview a different, but still beautiful family of techniques based on linear classification.

This is perhaps *the* central problem in data analysis. For instance, you may want to predict:

- will a sports team win a game?
- will a politician be elected?
- will someone like a movie?
- will someone click on an ad?
- will I get a job? (If you can build a good classifier, then probably yes!)

Each of these is typically solved by building a general purpose classifier (about sports or movies etc), then applying it to the person in question.

9.1 Linear Classifiers

Our input here is a point set $X \subset \mathbb{R}^d$, where each element $x_i \in X$ also has an associated label y_i . And $y_i \in \{-1, +1\}$.

Like in regression, our goal is prediction and generalization. We assume each $(x_i, y_i) \sim \mu$; that is, each data point pair, is drawn iid from some fixed but unknown distribution. Then our goal is a function $g: \mathbb{R}^d \rightarrow \mathbb{R}$, so that if $y_i = +1$, then $g(x_i) \geq 0$ and if $y_i = -1$, then $g(x_i) \leq 0$.

We will restrict that g is linear. For a data point $x \in \mathbb{R}^d$, written $x = (x^{(1)}, x^{(2)}, \dots, x^{(d)})$ we enforce that

$$g(x) = \alpha_0 + x^{(1)}\alpha_1 + x^{(2)}\alpha_2 + \dots + x^{(d)}\alpha_d = \alpha_0 + \sum_{j=1}^d x^{(j)}\alpha_j,$$

for some set of scalar parameters $\alpha = (\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_d)$. Typically, different notation is used: we set $b = \alpha_0$ and $w = (w_1, w_2, \dots, w_d) = (\alpha_1, \alpha_2, \dots, \alpha_d) \in \mathbb{R}^d$. Then we write

$$g(x) = b + x^{(1)}w_1 + x^{(2)}w_2 + \dots + x^{(d)}w_d = \langle w, x \rangle + b.$$

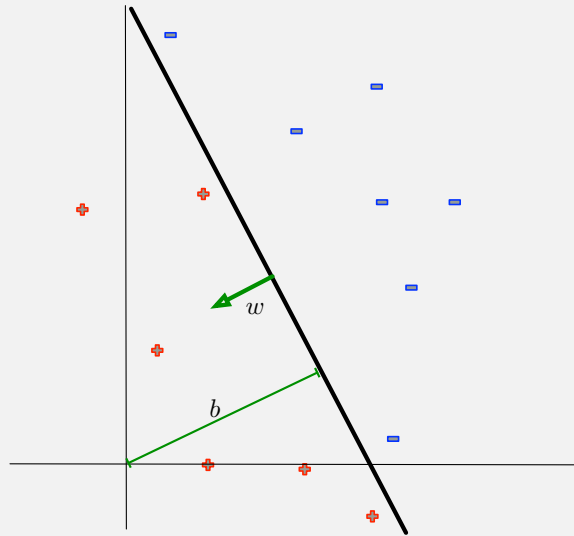
We can now interpret (w, b) as defining a halfspace in \mathbb{R}^d . Here w is the normal of that halfspace boundary (the single direction orthogonal to it) and b is the distance from the origin $\mathbf{0} = (0, 0, \dots, 0)$ to the halfspace boundary in the direction $w/\|w\|$. Because w is normal to the halfspace boundary, b is also distance from the closest point on the halfspace boundary to the origin (in any direction).

We typically ultimately use w as a unit vector, but it is not important since this can be adjusted by changing b . Let w, b be the desired halfspace with $\|w\| = 1$. Now assume we have another w', b' with $\|w'\| = \beta \neq 1$ and $w = w'/\|w'\|$, so they point in the same direction, and b' set so that they define the same halfspace. This implies $b' = b/\beta$. So the normalization of w can simply be done post-hoc without changing any structure.

Recall, our goal is $g(x) \geq 0$ if $y = +1$ and $g(x) \leq 0$ if $y = -1$. So if x lies directly on the halfspace then $g(x) = 0$.

Example: Linear Separator in \mathbb{R}^2

Here we show a set $X \in \mathbb{R}^2$ of 13 points with 6 labeled + and 7 labeled -. A linear classifier perfectly separates them. It is defined with a normal direction w (pointing towards the positive points) and an offset b .



Using techniques we have already learned, we can immediately apply two approaches towards this problem.

Linear classification via linear regression. For each data points $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, we can immediately represent x_i as the value of d explanatory variables, and y_i as the single dependent variable. Then we can set up a $n \times (d + 1)$ matrix M , where the i th row is $(1, x_i)$; that is the first coordinate is 1, and the next d coordinates come from the vector x_i . Then with a $y \in \mathbb{R}^n$ vector, we can solve for

$$\alpha = (M^T M)^{-1} M^T y$$

we have a set of $d + 1$ coefficients $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_d)$ describing a linear function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ defined

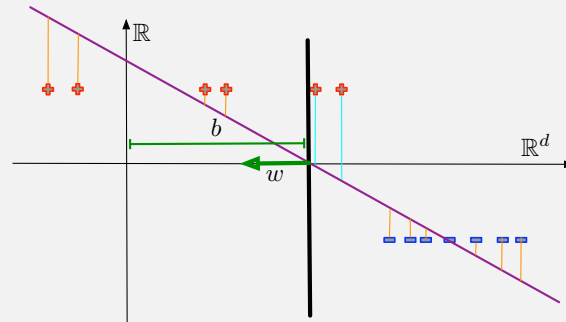
$$g(x) = \langle \alpha, (1, x) \rangle.$$

Hence $b = \alpha_0$ and $w = (\alpha_1, \alpha_2, \dots, \alpha_d)$. For x such that $g(x) > 0$, we predict $y = +1$ and for $g(x) < 0$, we predict $y = -1$.

However, this approach is optimizing the wrong problem. It is minimizing how close our predictions $g(x)$ is to -1 or $+1$, by minimizing the sum of squared errors. But our goal is to minimize the number of mispredicted values, not the numerical value.

Example: Linear Regression for Classification

We show 6 positive points and 7 negative points in \mathbb{R}^d mapped to \mathbb{R}^{d+1} . All of the d -coordinates are mapped to the x -axis. The last coordinate is mapped to the y -axis and is either $+1$ (a positive points) or -1 (a negative points). Then the best linear regression fit is shown, and the points where it has y -coordinate 0 defines the boundary of the halfspace. Note, despite there being a linear separator, this method misclassifies two points because it is optimizing the wrong measure.



Linear classification via gradient descent. Since, the linear regression SSE cost function is not the correct one, what is the correct one? We might define a cost function Δ

$$\Delta(g, (X, y)) = \sum_{i=1}^n (1 - \mathbb{1}(\text{sign}(y_i) = \text{sign}(g(x_i))))$$

which uses the *identity function* $\mathbb{1}$ (defined $\mathbb{1}(\text{TRUE}) = 1$ and $\mathbb{1}(\text{FALSE}) = 0$) to represent the *number* of misclassified points. This is what we would like to minimize.

Unfortunately, this function is discrete, so it does not have a useful (or well-defined) derivative. And, it is also not convex. Thus, encoding g as a $(d + 1)$ -dimensional parameter vector $(b, w) = \alpha$ and running gradient descent is not feasible.

However, most classification algorithms run some variant of gradient descent. To do so we will use a different cost function as a proxy for Δ , called a *loss function*. We explain this next.

9.1.1 Loss Functions

To use gradient descent for classifier learning, we will use a proxy for Δ called a *loss functions* \mathcal{L} . These are sometimes implied to be convex, and their goal is to approximate Δ . And in most cases, they are *decomposable*, so we can write

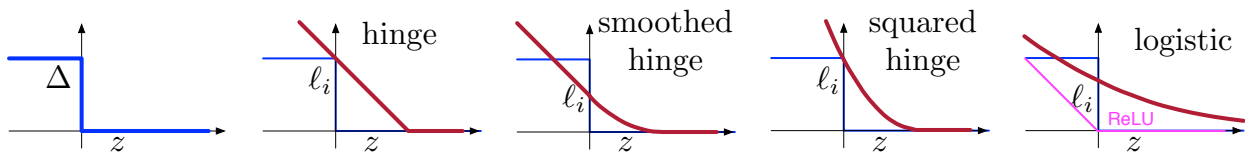
$$\begin{aligned} \mathcal{L}(g, (X, y)) &= \sum_{i=1}^n \ell_i(g, (x_i, y_i)) \\ &= \sum_{i=1}^n \ell_i(z_i) \text{ where } z_i = y_i g(x_i). \end{aligned}$$

Note that the clever expression $z_i = y_i g(x_i)$ handles when the function $g(x_i)$ correctly predicts the positive or negative example in the same way. If $y_i = +1$, and correctly $g(x_i) > 0$, then $z_i > 0$. On the other hand, if $y_i = -1$, and correctly $g(x_i) < 0$, then also $z_i > 0$. For instance, the desired cost function, Δ is written

$$\Delta(z) = \begin{cases} 0 & \text{if } z \geq 0 \\ 1 & \text{if } z < 0. \end{cases}$$

Most *loss functions* $\ell_i(z)$ which are convex proxies for Δ mainly focus on how to deal with the case $z_i < 0$ (or $z_i < 1$). The most common ones include:

- hinge loss: $\ell_i = \max(0, 1 - z)$
- smoothed hinge loss: $\ell_i = \begin{cases} 0 & \text{if } z \geq 1 \\ (1 - z)^2/2 & \text{if } 0 < z < 1 \\ \frac{1}{2} - z & \text{if } z \leq 0 \end{cases}$
- squared hinge loss: $\ell_i = \max(0, 1 - z)^2$
- logistic loss: $\ell_i = \ln(1 + \exp(-z))$



The hinge loss is the closest convex function to Δ ; in fact it strictly upper bounds Δ . However, it is non-differentiable at the “hinge-point,” (at $z = 1$) so it takes some care to use it in gradient descent. The smoothed hinge loss and squared hinge loss are approximations to this which are differentiable everywhere. The squared hinge loss is quite sensitive to outliers (similar to SSE). The smoothed hinge loss (related to the Huber loss) is a nice combination of these two loss functions.

The logistic loss can be seen as a continuous approximation to the ReLU (rectified linear unit) loss function, which is the hinge loss shifted to have the hinge point at $z = 0$. The logistic loss also has easy-to-take derivatives (does not require case analysis) and is smooth everywhere. Minimizing this loss for classification is called *logistic regression*.

9.1.2 Cross Validation and Regularization

Ultimately, in running gradient descent for classification, one typically defines the overall cost function f also using a regularization term $r(\alpha)$. For instance $r(\alpha) = \|\alpha\|^2$ is easy to use (has nice derivatives) and $r(\alpha) = \|\alpha\|_1$ (the L_1 norm) induces sparsity (for reasons not covered in this class). In general, the regularizer typically penalizes larger values of α , resulting in some bias, but less over-fitting of the data.

The regularizer $r(\alpha)$ is combined with a loss function $\mathcal{L}(g_\alpha, (X, y)) = \sum_{i=1}^n \ell_i(g_\alpha, (x_i, y_i))$ as

$$f(\alpha) = \mathcal{L}(g_\alpha, (X, y)) + \eta r(\alpha),$$

where $\eta \in \mathbb{R}$ is a *regularization parameter* that controls how drastically to regularize the solution.

Note that this function $f(\alpha)$ is still decomposable, so one can use batch, incremental, or most commonly stochastic gradient descent.

Cross-validation. Backing up a bit, the *true* goal is not minimizing f or \mathcal{L} , but predicting the class for new data points. For this, we again assume all data is drawn iid from some fixed but unknown distribution. To evaluate how well our results generalize, we can use cross-validation (holding out some data from the training, and calculating the expected error of Δ on these help out “testing” data points).

We can also choose the regularization parameter η by choosing the one that results in the best generalization on the test data after training using each on some training data.

9.2 Perceptron Algorithm

Of the above algorithms, generic linear regression is not solving the correct problem, and gradient descent methods do not really use any structure of the problem. In fact, we could have replaced the linear function $g_\alpha(x) = \langle \alpha, (1, x) \rangle$ with any function g (even non-linear ones) as long as we can take the gradient.

Now we will introduce the *perceptron algorithm* which explicitly uses the linear structure of the problem. (Technically, it only uses the fact that there is an inner product – which we will exploit in generalizations.)

Simplifications: For simplicity, we will make several assumptions about the data. First we will assume that the best linear classifier (w^*, b^*) defines a halfspace whose boundary passes through the origin. This implies $b^* = 0$, and we can ignore it. This is basically equivalent to (for data point $(x'_i, y_i) \in \mathbb{R}^{d'} \times \mathbb{R}$ using $x_i = (1, x'_i) \in \mathbb{R}^d$ where $d' + 1 = d$).

Second, we assume that for all data points (x_i, y_i) that $\|x_i\| \leq 1$ (e.g., all data point live in a unit ball). This can be done by choosing the point $x_{\max} \in X$ with largest norm $\|x_{\max}\|$, and dividing all data points by $\|x_{\max}\|$ so that point has norm 1, and all other points have smaller norms.

Finally, we assume that there exists a perfect linear classifier. One that classifies each data point to the correct class. There are variants to deal with the cases without perfect classifiers, but which are beyond the scope of this class.

The algorithm. Now to run the algorithm, we start with some normal direction w (initialized as any positive point), and then add mis-classified points to w one at a time.

Algorithm 9.2.1 Perceptron(X)

Initialize $w = y_i x_i$ for any $(x_i, y_i) \in (X, y)$

repeat

For any (x_i, y_i) such that $y_i \langle x_i, w \rangle < 0$ (is mis-classified) : update $w \leftarrow w + y_i x_i$

until (no mis-classified points or T steps)

return $w \leftarrow w / \|w\|$

Basically, if we find a mis-classified point (x_i, y_i) and $y_i = +1$, then we set $w = w + x_i$. This makes w “point” more in the direction of x_i , but also makes it longer. Having w more in the direction of w , tends to make it have dot-product (with a normalized version of w) closer to 1.

Similar, if we find a mis-classified point (x_i, y_i) with $y_i = -1$, then we set $w = w - x_i$; this points the negative of w more towards x_i , and thus w more away from x_i , and thus its dot product more likely to be negative.

The margin. To understand why the perceptron works, we need to introduce the concept of a margin. Given a classifier (w, b) , the margin is

$$\gamma = \min_{(x_i, y_i) \in X} y_i (\langle w, x_i \rangle + b).$$

It's the minimum distance of any data point x_i to the boundary of the halfspace. In this sense the optimal classifier (or the maximum margin classifier) (w^*, b^*) is the one that maximizes the margin

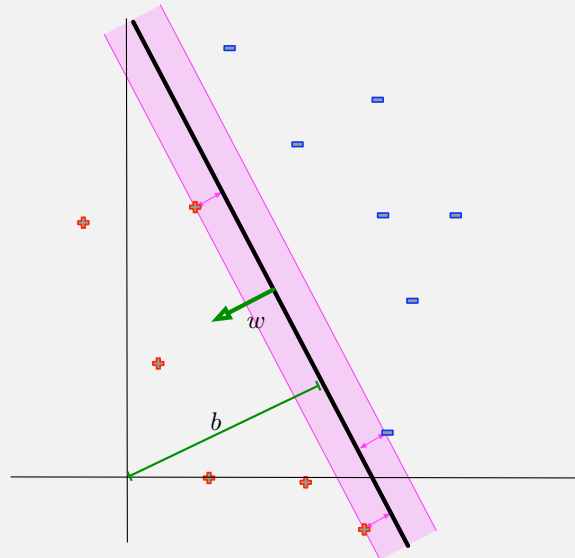
$$(w^*, b^*) = \arg \max_{(w, b)} \min_{(x_i, y_i) \in X} y_i (\langle w, x_i \rangle + b)$$
$$\gamma^* = \min_{(x_i, y_i) \in X} y_i (\langle w^*, x_i \rangle + b^*).$$

A max-margin classifier, is one that not just classifies all points correctly, but does so with the most “margin for error.” That is, if we perturbed any data point or the classifier itself, this is the classifier which

can account for the most perturbation and still predict all points correctly. It also tends to *generalize* (in the cross-validation sense) to new data better than other perfect classifiers.

Example: Margin of Linear Classifier

For a set X of 13 points in \mathbb{R}^2 , and a linear classifier defined with (w, b) . We illustrate the margin in pink. The margin $\gamma = \min_{(x_i, y_i)} y_i(\langle w, x_i \rangle + b)$. The margin is drawn with an \leftrightarrow for each support point.



The maximum margin classifier (w^*, b^*) for $X \subset \mathbb{R}^d$ can always be defined uniquely by $d + 1$ points (at least one negative, and at least one positive). These points $S \subset X$ are such that for all $(x_i, y_i) \in S$

$$\gamma^* = y_i \langle w^*, x_i \rangle + b.$$

These are known as the *support points*, since they “support” the margin strip around the classifier boundary.

Why perceptron works. We claim that after at most $T = (1/\gamma^*)^2$ steps (where γ^* is the margin of the maximum margin classifier), then there can be no more mis-classified points.

To show this we will bound two terms as a function of t , the number of mistakes found: $\langle w, w^* \rangle$ and $\|w\|^2 = \langle w, w \rangle$, before we ultimately normalize w in the **return** step.

First we can argue that $\|w\|^2 \leq t$, since each step increases $\|w\|^2$ by at most 1:

$$\langle w + y_i x_i, w + y_i x_i \rangle = \langle w, w \rangle + (y_i)^2 \langle x_i, x_i \rangle + 2y_i \langle w, x_i \rangle \leq \langle w, w \rangle + 1 + 0.$$

This is true since if x_i is mis-classified, then $y_i \langle w, x_i \rangle$ is negative.

Second, we can argue that $\langle w, w^* \rangle \geq t\gamma^*$ since each step increases it by at least γ^* . Recall that $\|w^*\| = 1$

$$\langle w + y_i x_i, w^* \rangle = \langle w, w^* \rangle + (y_i) \langle x_i, w^* \rangle \geq \langle w, w^* \rangle + \gamma^*.$$

The inequality follows from the margin of each point being at least γ^* with respect to the max-margin classifier w^* .

Combining these facts together we obtain

$$t\gamma^* \leq \langle w, w^* \rangle \leq \langle w, \frac{w}{\|w\|} \rangle = \|w\| \leq \sqrt{t}.$$

Solving for t yields $t \leq (1/\gamma^*)^2$ as desired.

9.3 Kernels

It turns out all we need to get any of the above machinery to work is a well-defined (generalized) inner-product. For two vectors $p = (p_1, \dots, p_d), q = (q_1, \dots, q_d) \in \mathbb{R}^d$, we have always used as the inner product:

$$\langle p, q \rangle = \sum_{i=1}^d p_i \cdot q_i.$$

However, we can define inner products more generally as a kernel $K(p, q)$. For instance, we can use

- $K(p, q) = \exp(-\|p - q\|^2/\sigma^2)$ for the Gaussian kernel, with bandwidth σ ,
- $K(p, q) = \exp(-\|p - q\|/\sigma)$ for the Laplace kernel, with bandwidth σ , and
- $K(p, q) = (\langle p, q \rangle + c)^r$ for the polynomial kernel of power r , with control parameter $c > 0$.

Then we define our classification function

$$g(x) = K(x, w) + b.$$

If $g(x) > 0$, we classify x as positive, and if $g(x) < 0$, we classify x as negative.

For gradient decent, again $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_d) = (b, w)$. We just need to take the gradient for each term of the loss function $\ell_i(z_i)$ for $z_i = y_i g(x_i)$ as before.

And for perceptron, we check $y_i K(x_i, w) > 0$ to see if (x_i, y_i) is mis-classified, and still simply add $w \leftarrow w + y_i x_i$ as before.

Are these linear classifiers? No. In fact, this is how you model various forms of non-linear classifiers. The “decision boundary” is no longer described by the boundary of a halfspace. For the polynomial kernel, the boundary must now be a polynomial surface of degree r . For the Gaussian and Laplace kernel it can be even more complex; the σ parameter essentially bounds the curvature of the boundary.

9.4 k NN Classifiers

Now for something completely different. There are many ways to define a classifier, and we have just touched on some of them. These include decision trees (which basically just ask a series of yes/no questions and are very interpretable) to deep neural networks (which are more complex, far less interpretable, but can achieve more accuracy). We will describe one more simple classifier.

The k -NN classifier (or *k-nearest neighbors classifier*) works as follows. Choose a scalar parameter k (it will be far simpler to choose k as an odd number, say $k = 5$). Next define a *majority* function $\text{maj} : \{-1, +1\}^k \rightarrow \{-1, +1\}$. For a set $Y = (y_1, y_2, \dots, y_k) \in \{-1, +1\}^k$ it is defined

$$\text{maj}(Y) = \begin{cases} +1 & \text{if more than } k/2 \text{ elements of } Y \text{ are } +1 \\ -1 & \text{if more than } k/2 \text{ elements of } Y \text{ are } -1. \end{cases}$$

Then for a data set X where each element $x_i \in X$ has an associated label $y_i \in \{-1, +1\}$, define a k -nearest neighbor function $\phi_{X,k}(z)$ that returns the k points in X which are closest to a query point z . Next let sign report y_i for any input point x_i ; for a set of inputs x_i , it returns the set of values y_i .

Finally, the k -NN classifier is

$$g(z) = \text{maj}(\text{sign}(\phi_{X,k}(z))).$$

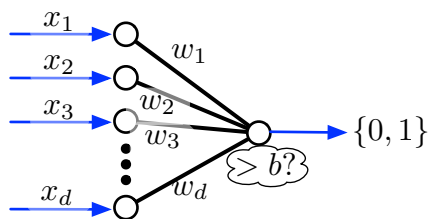
That is, it finds the k -nearest neighbors of query point z , and considers all of the class labels of those points, and returns the majority vote of those labels.

A query point z near many other positive points will almost surely return $+1$, and symmetrically for negative points. This classifier works surprisingly well for many problems but relies on a good choice of distance function to define $\phi_{X,k}$.

Unfortunately, the model for the classifier depends on all of X . So it may take a long time to evaluate on a large data set X . In contrast the functions g for all methods above take $O(d)$ time to evaluate for points in \mathbb{R}^d , and thus are very efficient.

9.5 Neural Networks

A *neural network* is a learning algorithm intuitively based on how a neuron works in the brain. A neuron takes in a set of inputs $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$, weights each input by a corresponding scalar $w = (w_1, w_2, \dots, w_d)$ and “fires” a signal if the total weight $\sum_{i=1}^d w_i x_i$ is greater than some threshold b .



$$\sum_{j=1}^d w_j \cdot x_j - b = \langle x, w \rangle - b > 0?$$

A neural network, is then just a network or graph of these neurons. Typically, these are arranged in layers. In the first layer, there may be d input values x_1, x_2, \dots, x_d . These may provide the input to t neurons (each neuron might use fewer than all inputs). Each neuron produces an output y_1, y_2, \dots, y_t . These outputs then serve as the input to the second layer, and so on.

Once the connections are determined, then the goal is to learn the weights on each neuron so that for a given input, a final neuron fires if the input satisfies some pattern (e.g., the input are pixels to a picture, and it fires if the picture contains a car). This is theorized to be “loosely” how the human brain works.

Given a data set X with labeled data points $(x, y) \in X$ (with $x \in \mathbb{R}^d$ and $y \in \{-1, +1\}$), we already know how to train a single neuron so for input x it tends to fire if $y = 1$ and not fire if $y = -1$. It is just a linear classifier. So, we can use the perceptron algorithm, or gradient descent with a well-chosen loss function!

However, for neural networks to attain more power than simple linear classifiers, they need to be at least two layers. Many amazing advances have come from so-called “deep neural networks” or “deep learning” or “deep nets” which are neural networks with many layers (say 20 or more). For these networks, the perceptron algorithm no longer works since it does not properly propagate across layers. However, a version of gradient descent called back-propagation can be used. Getting deep nets to work can be quite finicky. Their optimization function is not convex, and without various training tricks, it can be very difficult to find a good global set of weights. (The full details of this approach is well beyond the scope of this class.)