

"Linear" Regression

Feb 26, 2018

• Clustering Hw Due Wednesday

• Data Science Club : Stack Overflow (PCA)
Tomorrow, Tue 5-6 WETS 1230

• Start of Regression / Dim-Reduction

See more in Introduction to Data Analysis (M4DA) book!

Input: $P = (p_1, p_2, \dots, p_n)$

$p_i \in \mathbb{R}^{d+1}$ today $p_i \in \mathbb{R}^2$

$$p_i = (x_i, y_i)$$

$\mathbb{R}^d \rightarrow \mathbb{R}$

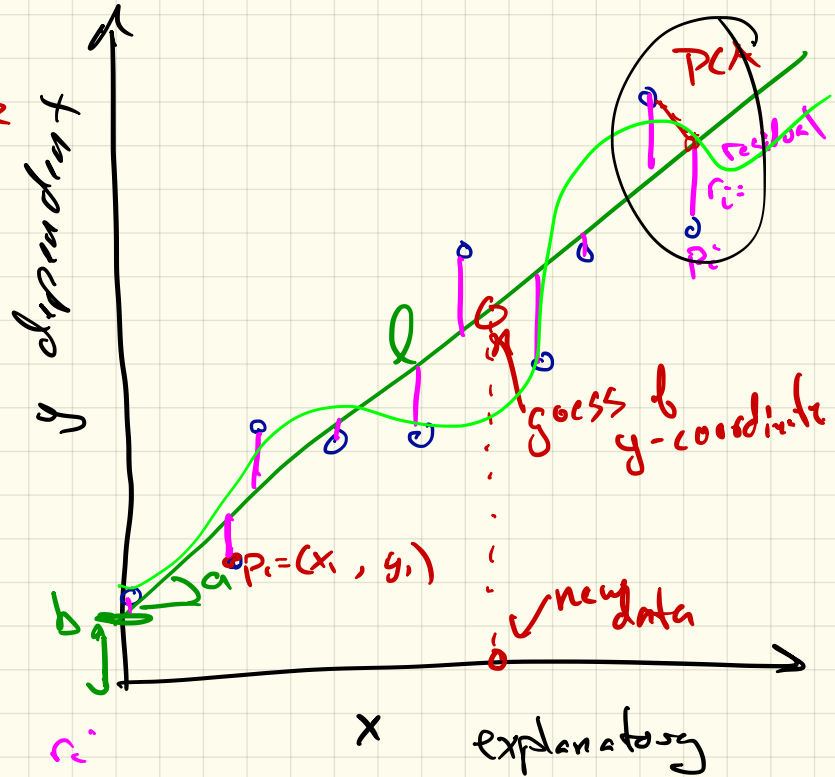
Goal fit line

$$l: l(x) = ax + b$$

$$\arg \min_l \sum_{p_i \in P} (y_i - l(x_i))^2$$

residual r_i

$$= \arg \min_{a, b} \sum_{p_i \in P} (y_i - ax_i - b)^2$$



← ordinary least squares regression

1. Closed Form Solution

$$|P| = n$$

$$\bar{P}_x = \text{average} \{x_i\} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{P}_y = \text{average} \{y_i\}$$

$$\text{Cov}[P_x, P_y] = \frac{1}{n} \sum_{P_i \in P} (P_{x_i} - \bar{P}_x) (P_{y_i} - \bar{P}_y)$$

$$\text{Var}[P_x] = \text{Cov}[P_x, P_x]$$

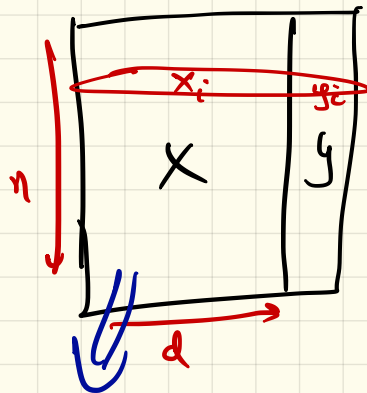
$$a = \frac{\text{Cov}[P_x, P_y]}{\text{Var}[P_x]} = \left(\frac{\langle P_x, P_y \rangle}{\|P_x\|^2} \right) \quad \text{"centered"}$$

$$b = \bar{P}_y - a \bar{P}_x$$

minimize
least squares

2. Extends to $x \in \mathbb{R}^d$

Input $P = (X, y)$ $X \subset \mathbb{R}^{n \times d}$ $y \in \mathbb{R}^n$



$P_i: x_i \in \mathbb{R}^d$ $x_i = (x_{i1}, x_{i2}, \dots, x_{id})$

$$f(x) = (b = a_0) + \sum_{j=1}^d x_{ij} a_j$$

$$\tilde{X} = \begin{matrix} \begin{matrix} x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & \vdots & & \vdots \\ x_{i1} & x_{i2} & \dots & x_{id} \\ \vdots & \vdots & & \vdots \end{matrix} \\ \downarrow n \\ \text{dim} \end{matrix}$$

Goal: $\arg \min_{a = (a_0, \dots, a_d)} \sum_{P_i \in P} (y_i - f_{a_i}(x_i))^2$

$$\Rightarrow a = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$$

↑ minimize least squares formulation

3. Polynomial Fits

Input $P \in \mathbb{R}^2$ $P_i = (x_i, y_i) \in \mathbb{R}^2$

Goal : $g_a: y = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots + a_t x^t$
 $= \sum_{j=0}^t a_j x^j$

$$\underset{a \in \mathbb{R}^t}{\text{argmin}} \sum_{P \in P} (y_i - g_a(x_i))^2$$

$$X = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 & \dots & x_1^t \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 & \dots & x_n^t \end{pmatrix}$$

4. Gauss - Markov Theorem

↳ "minimum variance solution"

Assuming

- unbias solution, expected error 0
 - all errors ϵ_i are not known to correlate
-

① minimize projection instead of vertical.?
PCA

② $(X^T X)^{-1}$ expensive (streaming, SGD)

③ Minimize something other than $\sum \epsilon_i^2$
→ Robust Estimator (Theil-Sen Est)

④ biased solutions : regularization
lasso

Theil-Sen Estimator

Breakdown point of estimator is largest fraction of data points which can be arbitrary outliers, and the estimator is still "OK".



median has a breakdown point of $1/2$.

median in \mathbb{R}^2

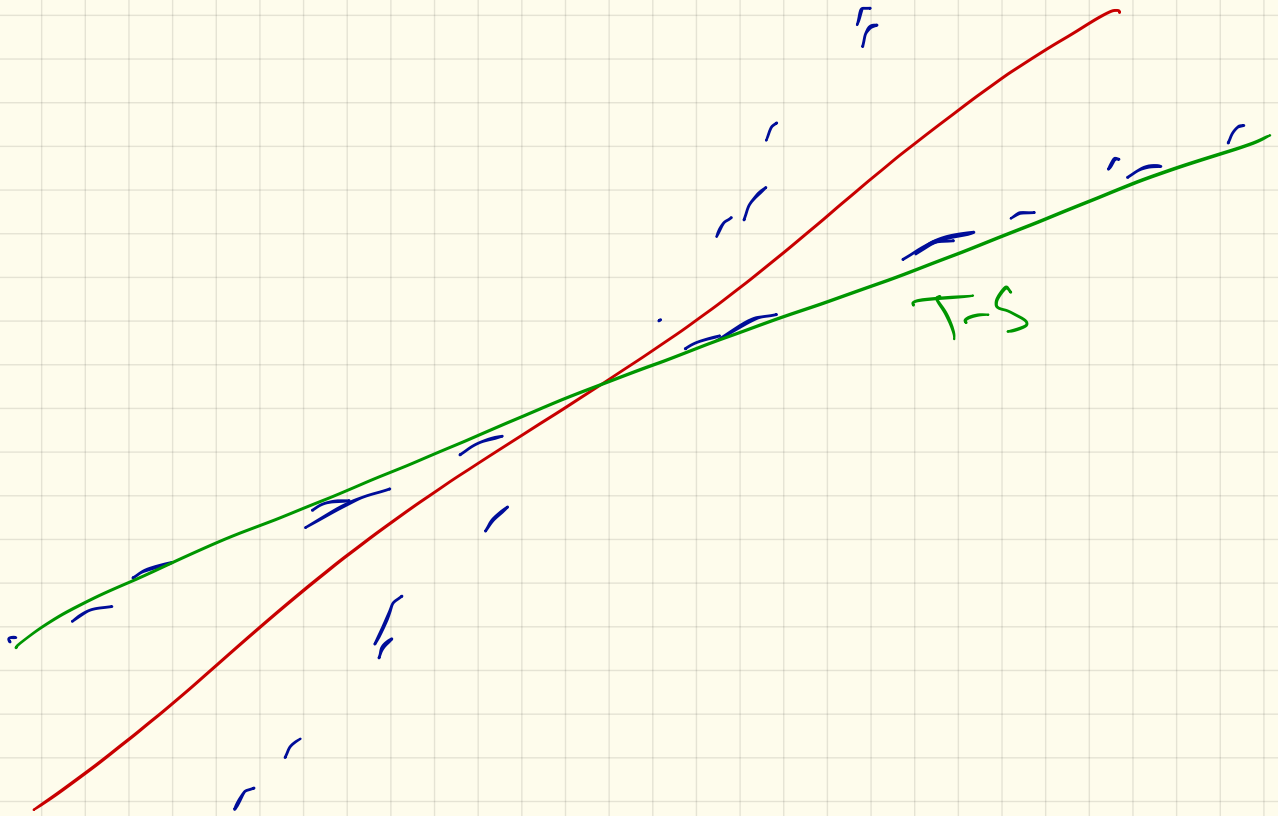
$$S = \{s_{ij} \mid \frac{y_j - y_i}{x_j - x_i}, x_i < x_j\}$$

$$a = \text{median} \{s_{ij} \in S\}$$

$$b = \text{median} \{y_i - a x_i\}$$

breakdown pt
0.293

- half points left, half right
- $\text{argmin}_{M \in \mathbb{R}^2} \sum_{P \in \mathcal{P}} |P - M|$



Tikhonov Regularization (Ridge Regression)

$$\text{Cost } L_{2,s}(P, a) = \sum_{i \in P} (y_i - ax_i)^2 + s \|a\|_2^2$$

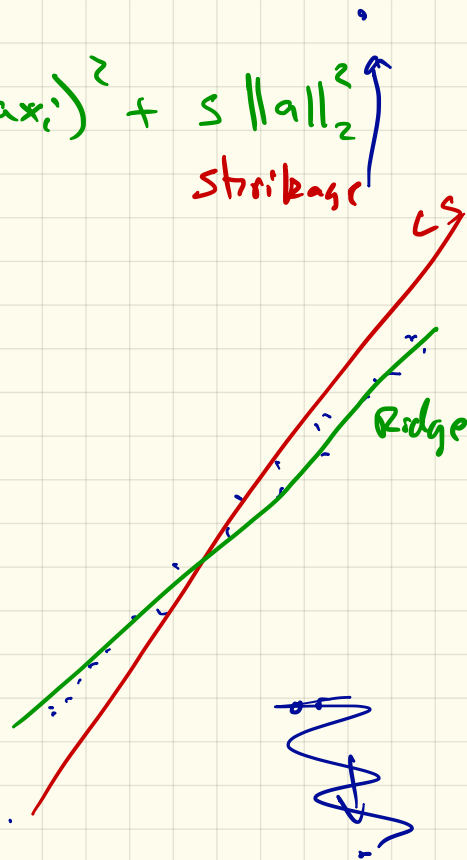
↑
shrinkage

closed form soln

$$a = (\tilde{X}^T \tilde{X} + s^2 I)^{-1} \tilde{X}^T y$$

$$= \frac{\langle P_x, P_y \rangle}{\langle P_x, P_x \rangle + s^2}$$

- Choosing s
 - + cross-validation
 - + \exists variance at least as small as LS.



Lasso (basis pursuit)

$$\text{Cost: } L_{1,s}(P, a) = \sum_{i \in P} (y_i - a x_i)^2 + s \|a\|_1$$

→ no "simple" closed form soln.

When $a \in \mathbb{R}^d$ d large
then biases toward "sparse" a .

↳ many $a_i = 0$

Even works when $d \gg n$