

L20: Noise + Outliers

Mar 31, 2025
Data Mining



Jeff M. Phillips

Types of Noise

Data $X = \{x_1, \dots, x_n\}$

• Measurement Noise

$$x_i = \text{true}(x_i) + \epsilon_i$$

common noise $N(0, \Sigma) = \exp\left(-\frac{1}{2} \frac{x \cdot x}{\Sigma}\right) \epsilon_i \sim \text{Noise}$

↳ loss function (opt - criteria)

$$\sum_i \| \hat{u} - x_i \|^2$$

• Spurious Readings

↳ outliers



outliers
↳ swan

• Background data



Dealing w/ outliers

- Goals:
1. Estimate model of inliers
 2. Identify outliers, since interesting
-

Requires assumption ^{inliers} data is fit by
some model $M \in \mathcal{M}$

Family of Models \mathcal{M}

↑ family of models

- think $X \sim \mathcal{N}_d(\mu, \Sigma)$ $X \in \mathbb{R}^d$
↑ parameters
- k-means clustering / Mix Gaussians
- low-rank + $\mathcal{N}(0,1)$ noise \rightarrow PCA
- linear regression

Outlier Removal

1. Fit best model

$$M^* = \underset{M \in \mathcal{M}}{\operatorname{argmax}} \text{Likelihood}(M; X)$$

my goal!

2. For all $x_i \in X$

residual

calculate

$$r_i = D(x_i, M)$$

$$\text{e.g. } r_i = -\log(\text{Likelihood}(M, x_i))$$

\leftarrow means

e.g. $M(x_i) = \phi_s(x_i) \leftarrow$ closest cluster center

$$D(x_i, M(x_i)) = \|x_i - \phi_s(x_i)\|$$

3. Remove $x_i \in X$

s.t.

$r_i \geq$

to big ?

needs threshold

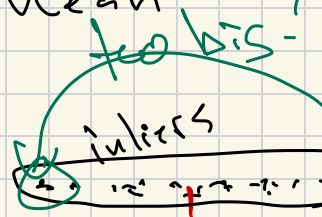
user defined

4. Repeat (Goto 1)

on

$X - \{\text{outliers}\}$.

Mean $\bar{x} \in \mathbb{R}$

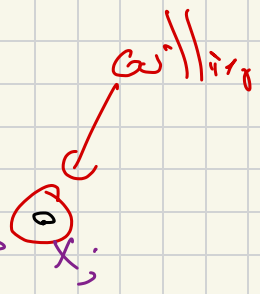


sample mean $\bar{x} \Rightarrow M^* = N(\bar{x}, I)$

x_1, x_2, \dots, x_n
 true mean

residual r_i
 is very large

Goal!



Repeat after filtering outliers
 $X_i = X \setminus \{\text{outliers}\}$

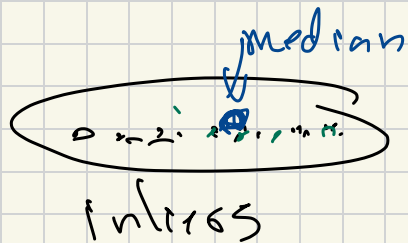


$\bar{x}_i \Rightarrow M_i^* = N(\bar{x}_i, I)$

Might not work
 if M^* is
 too poor an
 estimate

Robust Estimators

$$X \subset \mathbb{R}$$



outliers

Median = sort data, take point at 50%.

has breakdown point of 0.5

point which minimizes
$$\bar{\mu} = \arg \min_{\mu \in \mathbb{R}} \sum_{x_i \in X} |x_i - \mu|$$

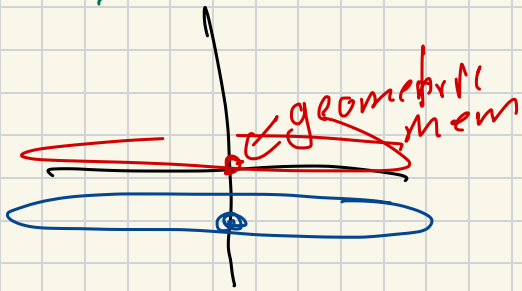
breakdown point = # points needed to move $\bar{\mu}$ "far" from inliers so outside inliers

How to compute median in \mathbb{R}^d $d > 1$?
as any robust estimator in \mathbb{R}^d .

• Geometric mean

$$\tilde{\mu} = \underset{\mu \in \mathbb{R}^d}{\text{argmin}} \sum_{x_i \in X} \|x_i - \mu\|$$

has a least breakdown point



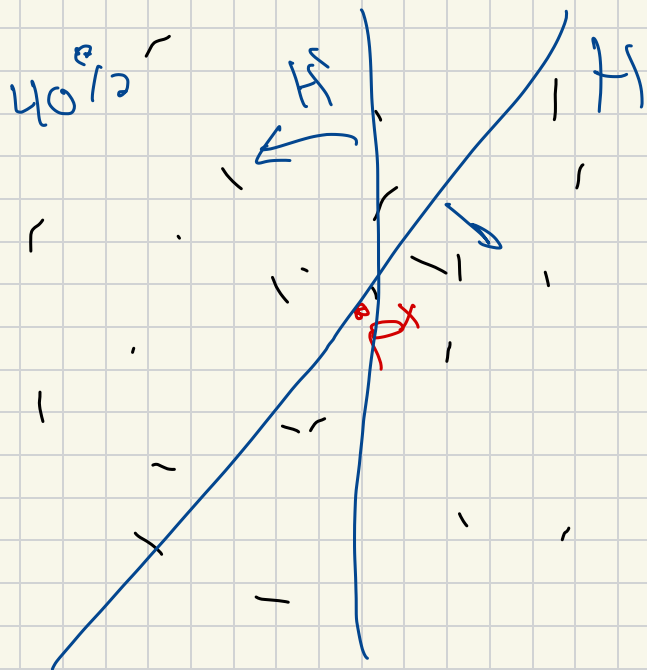
•



Tukey Median

point P^* = max A min H

half space (containing P)
 $\cap X$



gives good estimator

but hard to compute

large d

$\Theta(nd)$

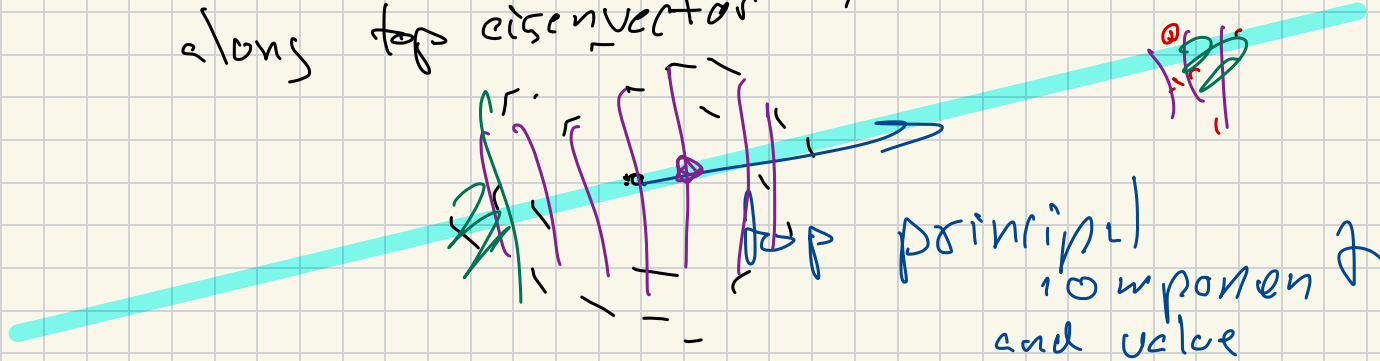
Robust Mean Estimator

$$X \sim \mathcal{N}(\mu, \Sigma)$$

1. Estimate sample mean $\hat{\mu} = \text{mean}(X)$

2. Estimate sample covariance $\hat{\Sigma} = \frac{1}{N} \sum_{x_i \in X} (x_i - \hat{\mu})(x_i - \hat{\mu})^T$

3. If $\|\hat{\Sigma}\|_2$ too big, filter points along top eigenvector



Median of Means

0. choose parameter

k (eg. $k=5$)

1. Randomly split X into k equal sets

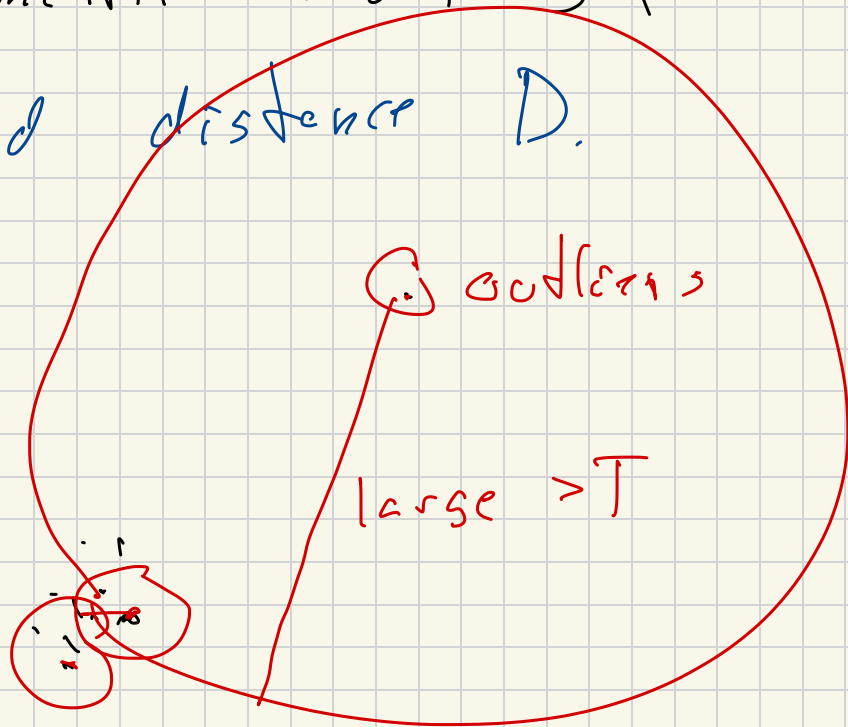
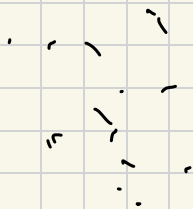
X_1, X_2, \dots, X_k

2. Calculate mean of each $X_i \rightarrow \text{mean}(X_i) = \mu_i$

3. Coordinate-wise median of $\{\mu_1, \mu_2, \dots, \mu_k\}$

Density-based Outliers

When no parametric model μ_f
Still need good distance D .



- DB Scan (clustering)

all points not in cluster \rightarrow outliers

- Parameter k (# points to define neighborhood)

threshold T

For each x_i distance D_k to k th

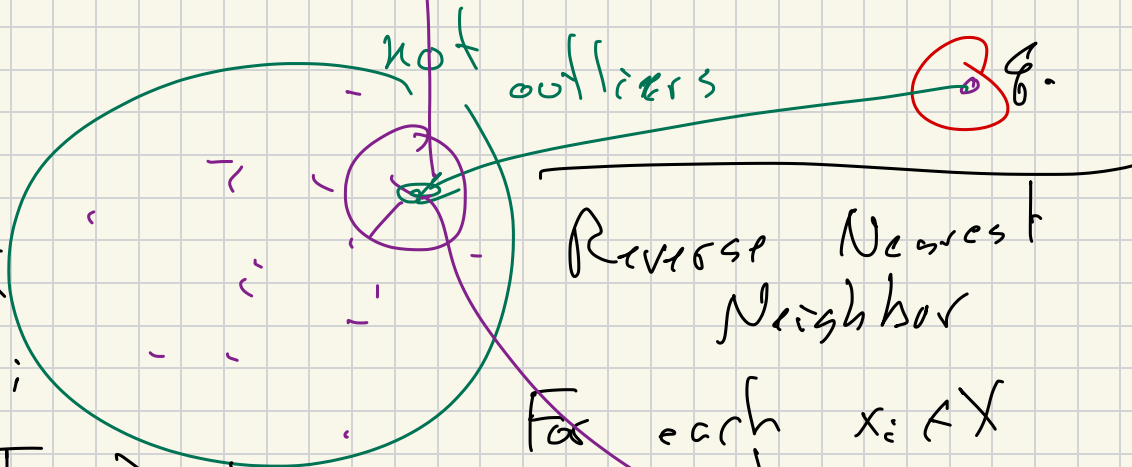
closest neighbor. if $> T \Rightarrow$ outlier

- Kernel Density Estimate

$$K(x_i, g) = \exp(-D(x_i, g)^2)$$

$$KDE_X(g) = \frac{1}{|X|} \sum_{x_i \in X} K(x_i, g) \quad \text{if } \leq T \Rightarrow g \text{ outlier}$$

All assume density is uniform
in inliers,



Let $N(q_i) = \{x_j \in X$
Nearest to q_i
if $\frac{D(q_i)}{D_i(x_i)} > T \Rightarrow$ outlier

Reverse Nearest Neighbor

For each $x_i \in X$
compute

$$D_i = \min_{x_j \in X, x_j \neq x_i} \|x_i - x_j\|$$

Outliers