





Noise

Data Science Club

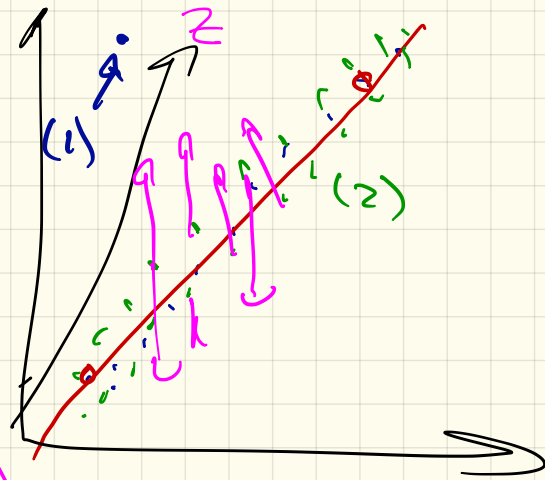
Tomorrow Tue, Apr 3 @ 5pm

Abode + Pizzas

WEB 1230

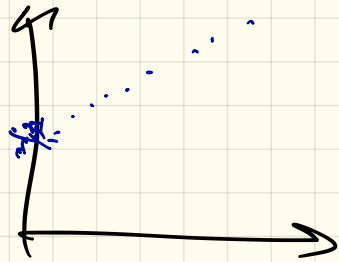
Noise in Data

- (1) Spurious Readings (outliers)
- (2) Measurement Error
hopefully unbiased
- (3) Background / Missing Data

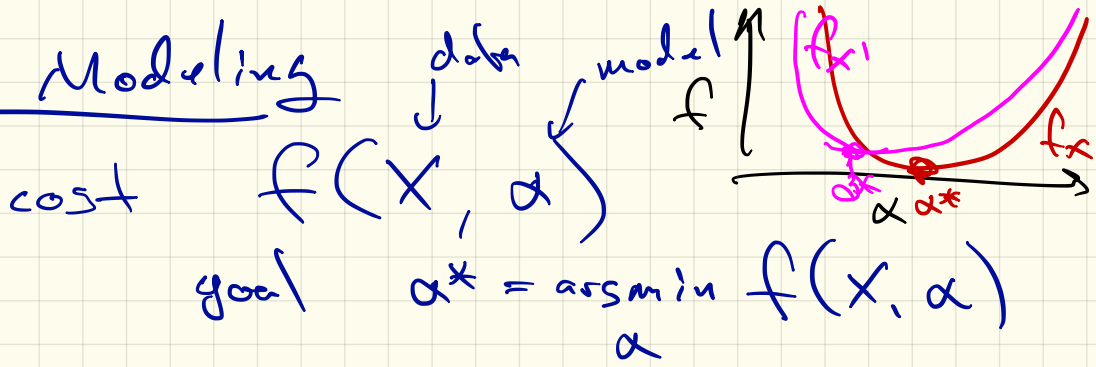


-
- Cross-Validation
 - Explicitly handling outliers

- Regularization
- Robust Statistics



Most Modeling

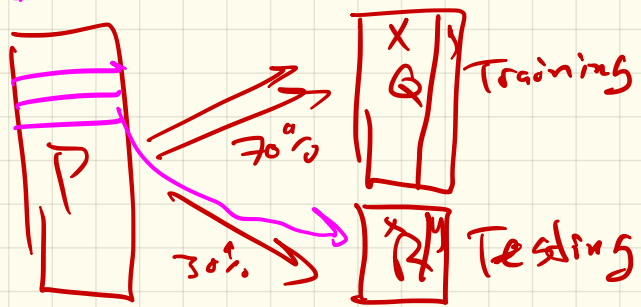


$$f(P, (\alpha, \gamma)) = \|P_y - P_x \alpha\|_2^2 + \gamma \|\alpha\|_2^2 \quad (\text{Ridge})$$

\downarrow all data observed iid

$\alpha_{x, \gamma}^* = g(x, \gamma)$

if γ fixed \rightarrow solve α^* closed form



for all γ

$$\alpha_{Q, \gamma}^* = g(Q, \gamma)$$

Evaluate $\gamma^* = \underset{\gamma}{\operatorname{argmin}} f(R, (\alpha_{Q, \gamma}^*, \gamma))$

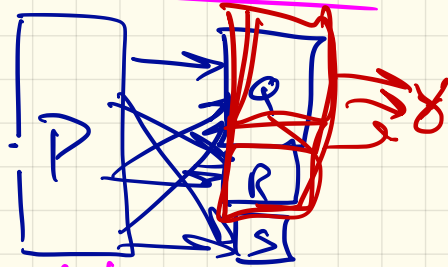
Use α_{Q, γ^*}^*

Evaluate $\gamma^* = \underset{\gamma}{\operatorname{arg\,min}} f(R, (x_{(Q, \gamma)}^*, \gamma))$

$g(Q, \gamma)$

Cross-Validation

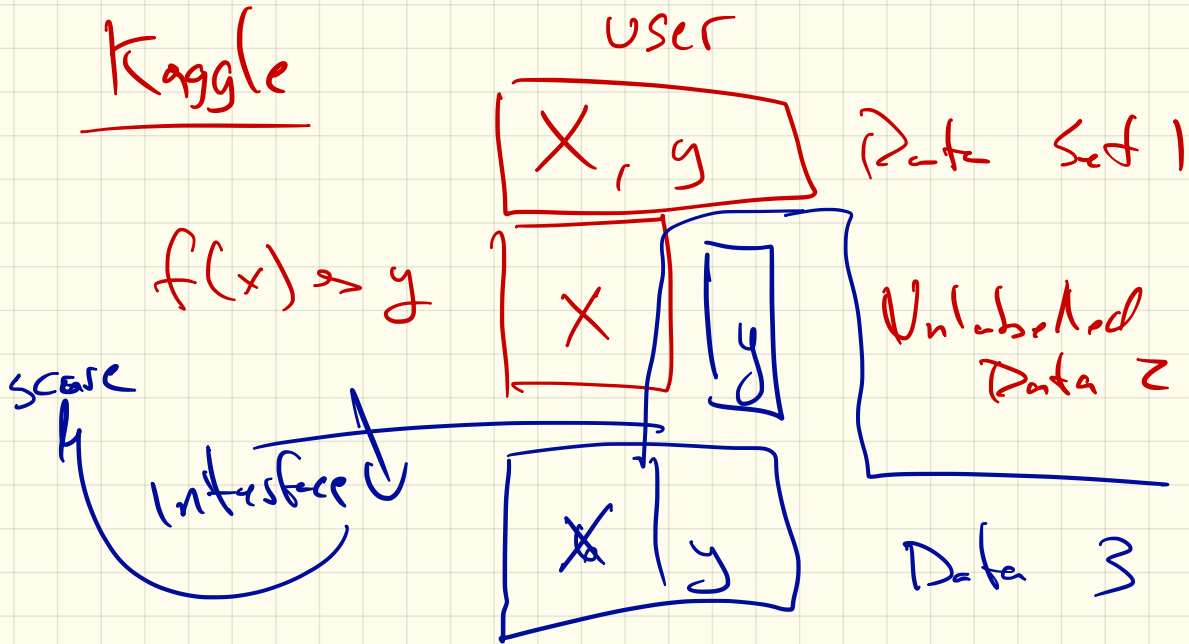
- (1) Choose parameters
- (2) Predict Generalization



Train on Q
Choose Param on R
Generalize on S .

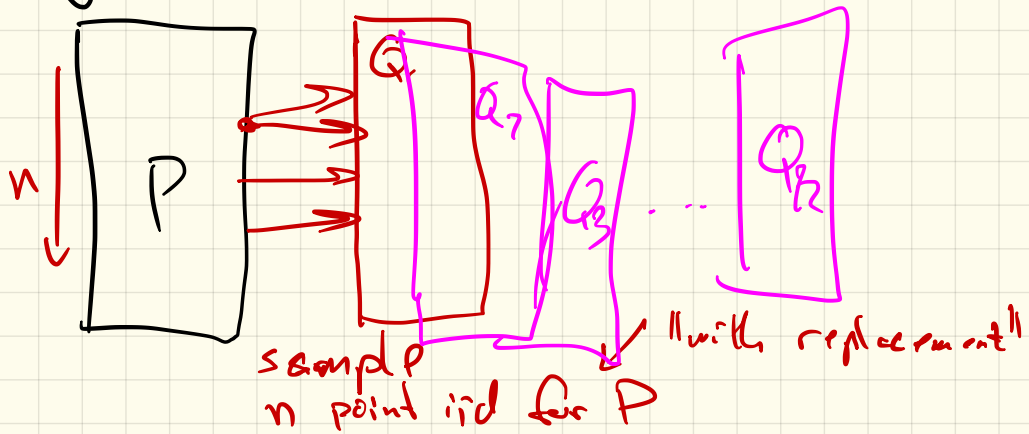
if we see new data?
how well will our model do?

Kaggle



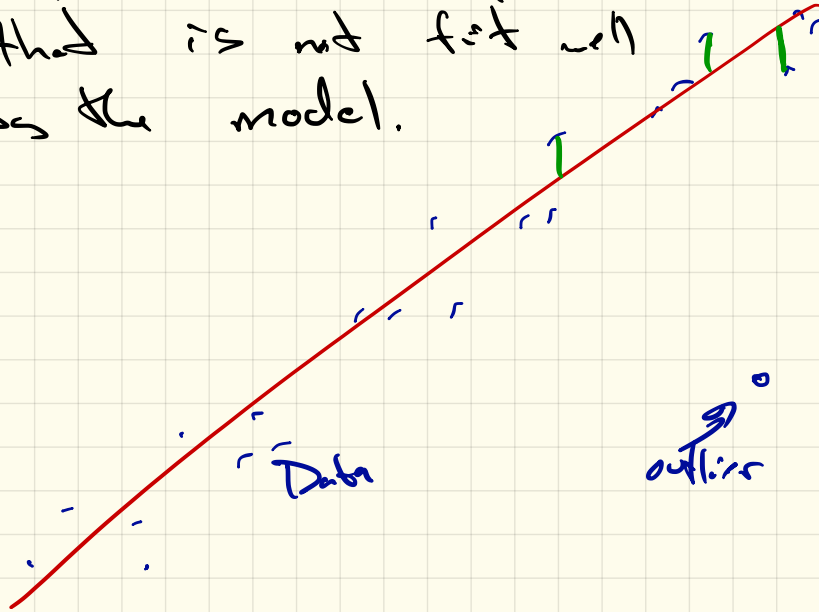
Bootstrapping

(Efron 1979)



Outliers

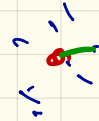
Outlier is a datapoint that is not fit well by the model.



Model

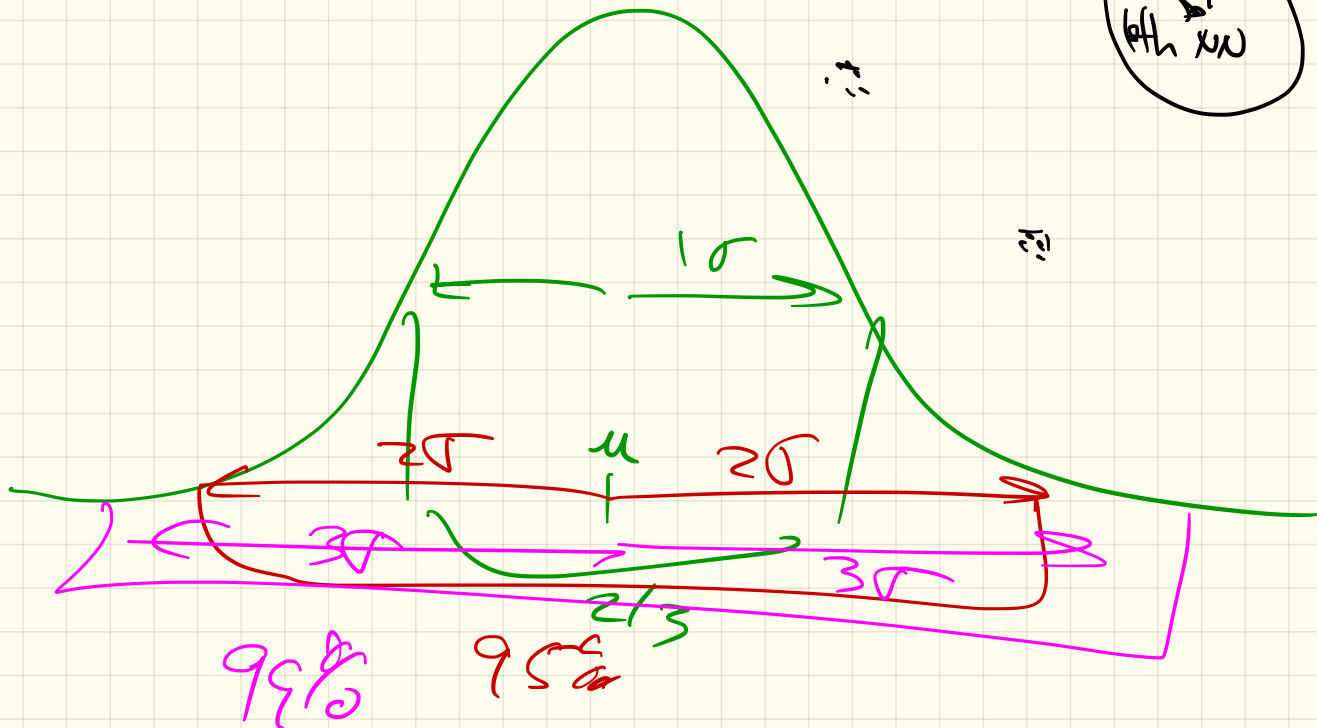
Also for outlier removal

1. Fit Model
2. Calc Resid
3. Remove x_i if resid too big.



- Density Data
- Reverse NN

$$e^{-\|x - \mu\|^2}$$



Heavy-Tailed Distrib.

Zipf's $f_i = c \cdot \left(\frac{1}{i}\right)^\alpha$

the 7%
of 3.5%
and 7.8%