

# Noise in Data

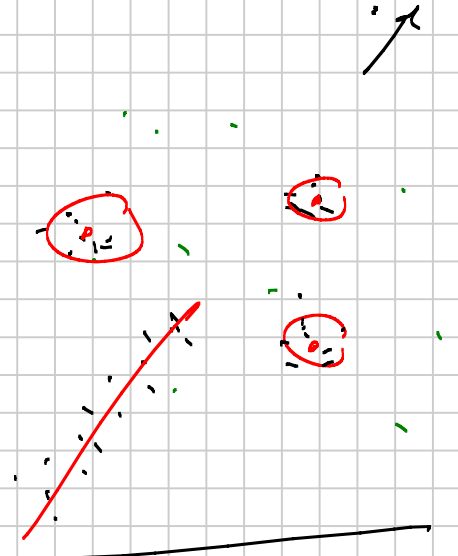
Note Title

3/30/2016

- Spurious Readings (outliers)

- Measurement Error  
small full-dimensional noise

- Background Data  
missing data



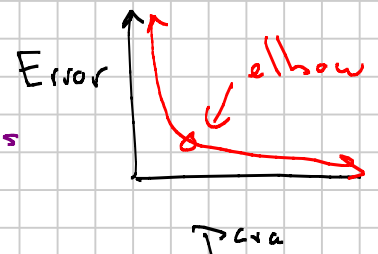
- 
- Cross-Validation
  - Outliers
  - Uncertain Data
  - Robust Estimators

---

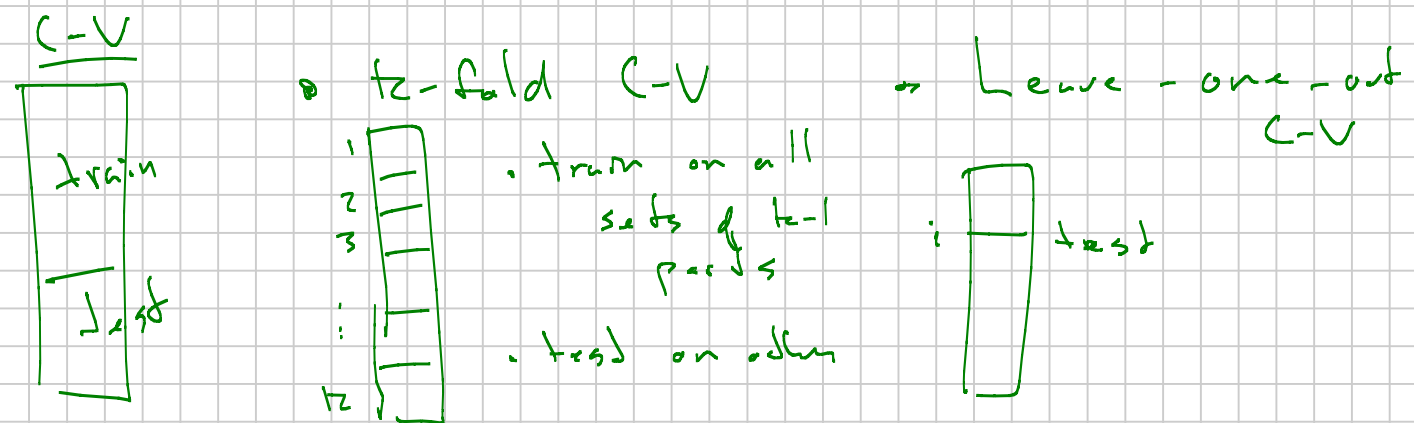
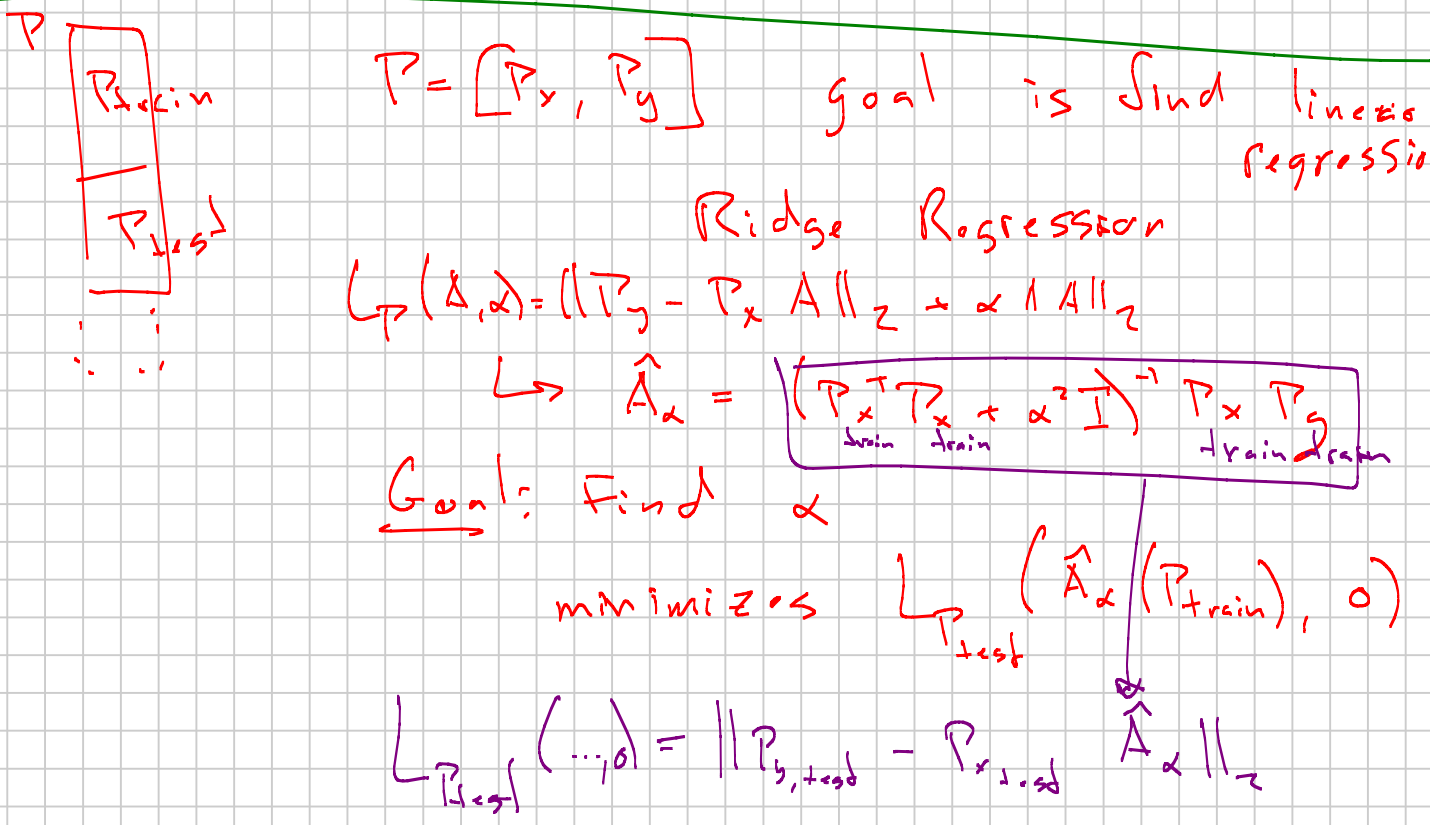
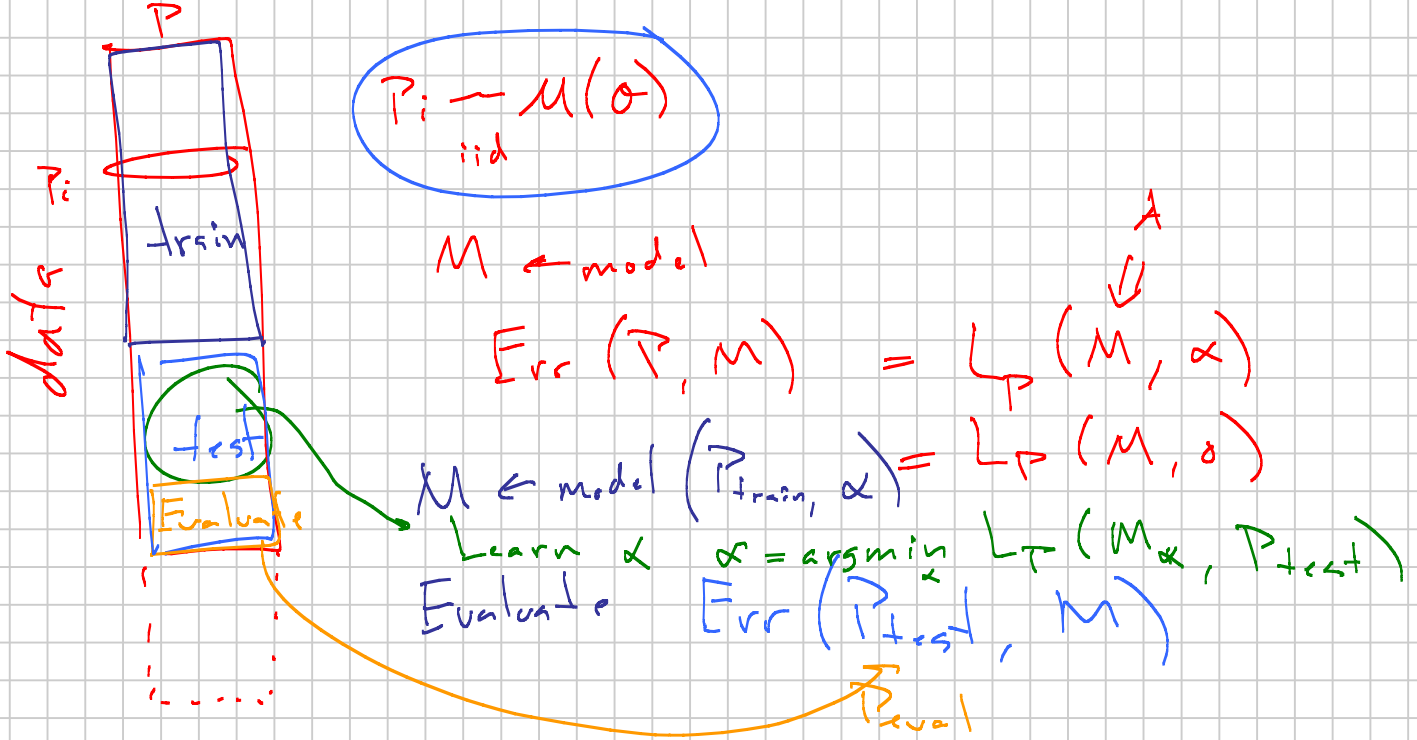
## Cross-Validation

Algo deals with noise, has a parameter

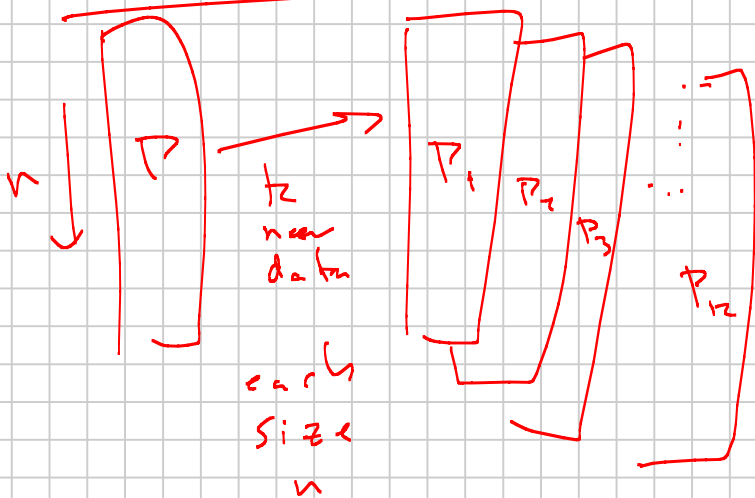
- Clustering  $\rightarrow k$
- PCA  $\rightarrow \#$  principal components
- Regression
  - degree polynomial
  - Regularization parameter



$$L_p(A, \alpha) = \|P_y - P_x A\|_2 + \alpha \|A\|_2$$



# Bootstrapping



• Sample with replacement

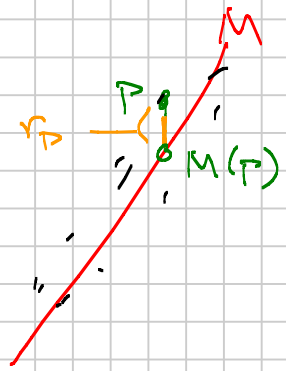
$T_i \leftarrow n$  items from  $P$  i.i.d.

$M_i \leftarrow \text{model}(T_i)$

$\{M_1, M_2, \dots, M_k\}$

# Outliers

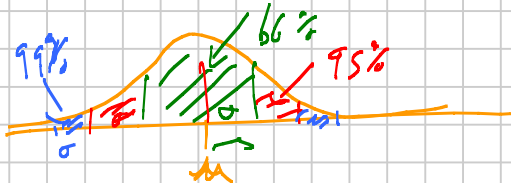
- Find + Remove Approach
- Density-based Approach



1. Build model  $M$  from input  $P$ .
2. For all  $p \in P$ , calc  $r_p = d(M(p), p)$
3. If  $r_p$  is "too large" it's an outlier!
4. Throw outliers away  $P$

iterate

Normal distribution

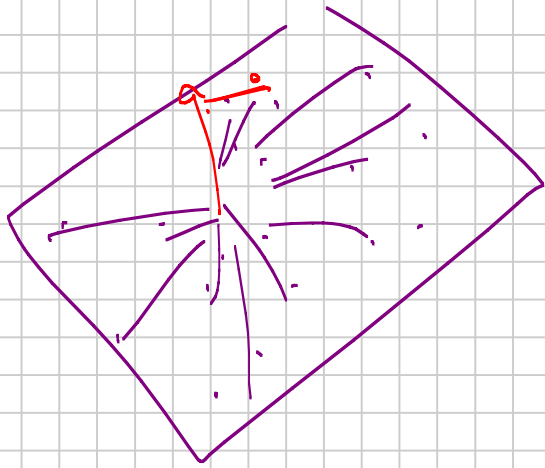


$$P(x) = e^{-\| \mu - x \|^2}$$

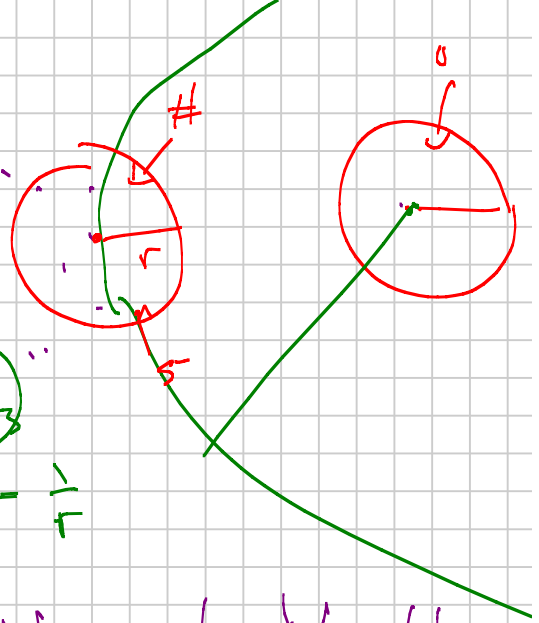
- 5% largest / 50 data points  $\rightarrow$  run out
- Pick threshold on  $r_p$  "scale"

# Density-based Approach

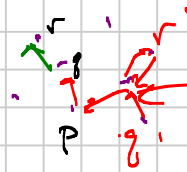
P



density  $\propto \frac{1}{r^2}$



## Reverse Nearest Neighbor



$$d(p, \phi_X(p)=q) \approx d(q, \phi_X(q)=p)$$

## Heavy-Tail Distributions

Zipf  $P(i) = c \cdot \left(\frac{1}{i}\right)$

the 7%  
of 3.5%  
and 2.8%

