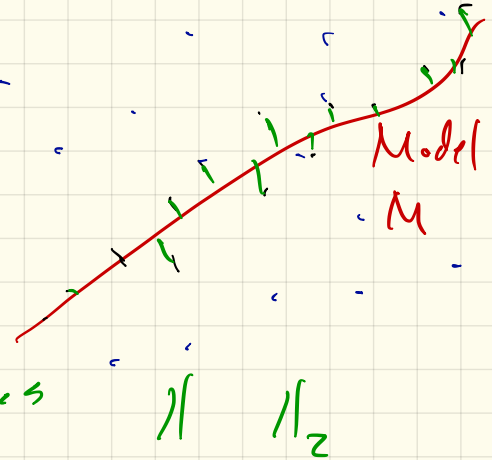


Noise in Data

- Spurious Readings
outliers

- Measurement Error

Gaussian \approx sum of squares error



- Background Data
- missing data

- Adversarial Data

Cross-Validation + Regularization

Model $M(x) = \underset{M}{\operatorname{argmin}} \left(\|M(x) - x\|_2^2 + \lambda \|M\|_1^2 \right)$

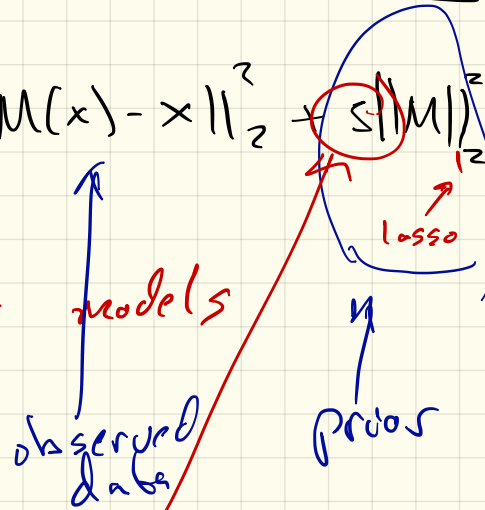
- biases towards simpler models

choose s w/ C-V
testing data



training data

generalization

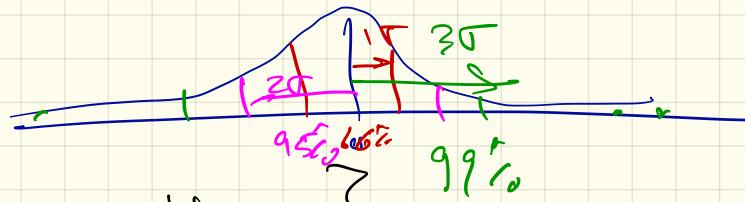


regularization parameter

Bootstrapping

for $i=1$ to n
draw n points $X_i \leftarrow X$ with replacement
measure model $M(x_1), M(x_2), \dots, M(x_n)$

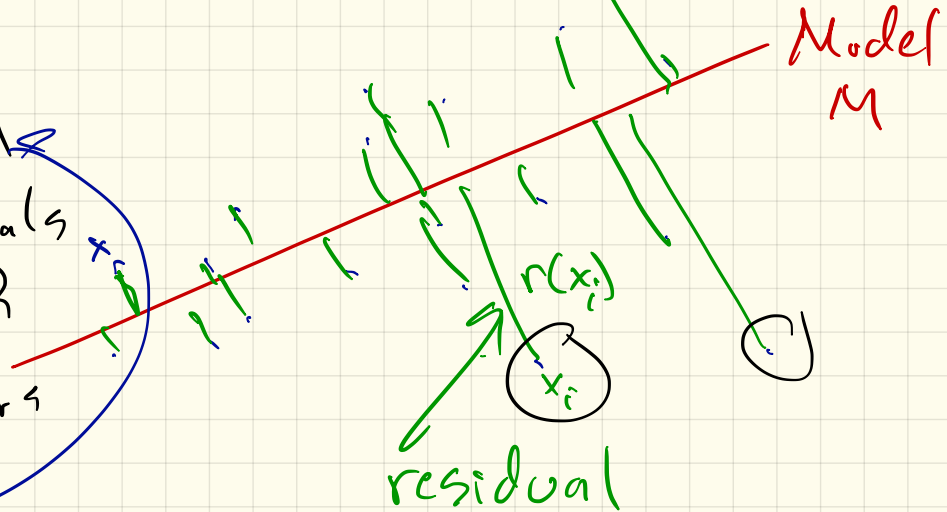
Outliers



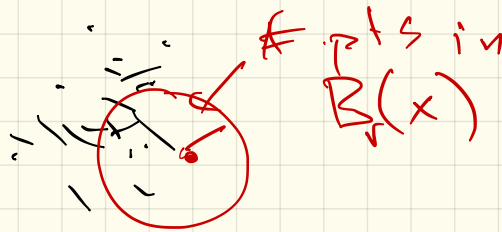
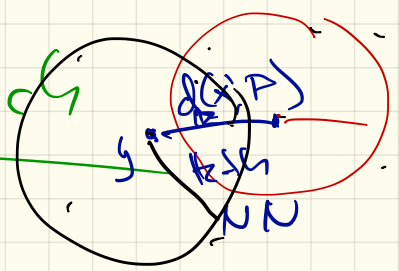
• How to remove them
Should I

1. Build Model M
 2. Calculate residuals $\{r(x) \mid x \in X\}$
 3. Remove outliers $r(x) > \tau$
4. Go to I

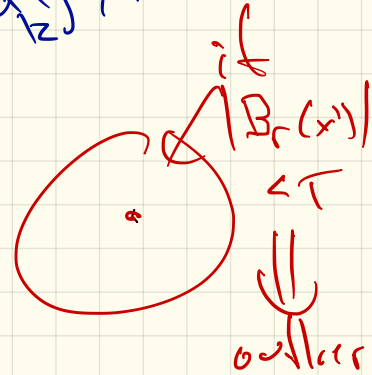
outliers = $x \in X$ w/
largest residual $r(x_i)$



Density-based Approach



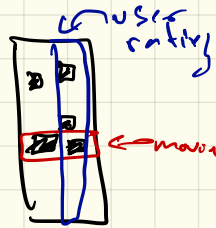
$d(x, A)$



reverse-nearest neighbors

Missing Data

Data set $P \in \mathbb{R}^{n \times d}$ matrix



$$\Omega = \{ (i, j) \in [n] \times [d] \mid P_{ij} \neq \emptyset \}$$

$\mathcal{I}_2(P)$ ← elements w/ scores

$\mathcal{I}_2^T(P)$ ← elements w/ out scores

$$P^* = \underset{X}{\operatorname{arg\,min}} \left\| \mathcal{I}_2(P - X) \right\|_F^2 + \lambda \|X\|_*$$

nuclear norm
= sum sing. values

Matrix Completion (P, Ω, λ)

0. Initialize X : $X_{ij} \leftarrow$ average of entries of row / column
or $\Pi_{\Omega}(P)$

repeat

1. $U S V^T \leftarrow \text{svd}(X)$

2. $\hat{X} \leftarrow U \boxed{\phi_{\lambda}(S)} V^T$

3. $X \leftarrow \Pi_{\Omega}(P) + \Pi_{\Omega}^{\perp}(\hat{X})$

until ("converges")

Return \hat{X}

$$\phi_{\lambda}(S) = \text{diag} \left(\begin{array}{l} (S_{11} - \lambda)_+ \\ (S_{22} - \lambda)_+ \\ \vdots \\ (S_{dd} - \lambda)_+ \end{array} \right)$$

$$(x - \lambda)_+ = \max\{0, x - \lambda\}$$