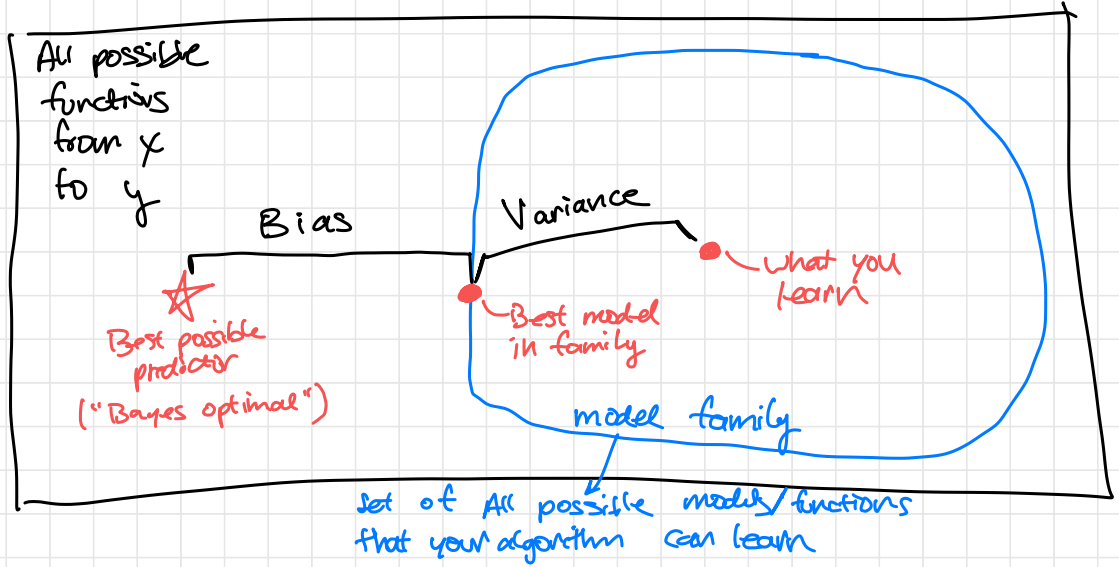


9/7/2023: Bias/Variance, MAP, Normal Equations



Bias: Error because assumptions of ML method used don't exactly match reality

Variance: Error because what you learn \neq best possible model which is because training data is incomplete

Total error can be decomposed into bias & variance

Reduce Bias

Make fewer assumptions



Make model family bigger

Reduce Variance

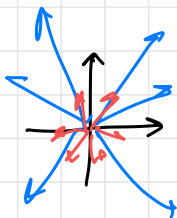
make it easier to find best model in family



make model family smaller



Regularization does this



without reg: w can have very large norm
with reg: w constrained to Ball of small norm

w can have very large norm
 w constrained to Ball of small norm

Maximum a Posteriori (MAP) estimation (extension of maximum likelihood estimation)

Idea: Think about learning as Bayesian inference problem

Think of everything as a random variable

Bayesian probabilistic story:

- ① Exists prior distribution over w , $p(w)$
- ② Some w gets sampled by nature
- ③ Dataset D is generated conditioned on w , called $p(D|w)$

Our goal: Infer most likely value of w

i.e. maximize w.r.t w $P(w|D)$

most likely w given data
observed

$$P(w|D) = \frac{P(w) P(D|w)}{P(D)}$$

$P(D)$ = Likelihood = what we maximize during MLE
doesn't depend on w so ignore it

New

Different choices of $p(w)$
give different regularization terms

For example: Let $p(w) = \prod_{j=1}^d \frac{1}{\sigma \sqrt{2\pi}} e^{-w_j^2 / 2\sigma^2}$

Gaussian with mean 0 and variance σ

$$\max_w P(w|D) = \max_w P(w) P(D|w)$$

$$= \min_w \underbrace{-\log P(w)}_{\text{regularization}} - \underbrace{\log P(D|w)}_{\text{original MLE loss}}$$

$$\begin{aligned}
 -\log P(w) &= -\sum_{j=1}^d \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \sum_{j=1}^d \frac{w_j^2}{2\sigma^2} \\
 &= \text{Constant} + \frac{1}{2\sigma^2} \sum_{j=1}^d w_j^2 \\
 &= \text{constant} + \frac{1}{2\sigma^2} \|w\|^2 \\
 &= \lambda \quad \begin{array}{l} \text{small } \sigma \Rightarrow \text{large } \lambda \\ \text{large } \sigma \Rightarrow \text{small } \lambda \end{array}
 \end{aligned}$$

Closed form solution for linear regression ("Normal Equations")

$$L(w) = \frac{1}{n} \sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2$$

$$\nabla_w L(w) = \frac{1}{n} \sum_{i=1}^n 2(w^T x^{(i)} - y^{(i)}) x^{(i)} = 0$$

$$\sum_{i=1}^n (w^T x^{(i)}) x^{(i)} = \sum_{i=1}^n y^{(i)} x^{(i)}$$

Ex: $X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}$
 "design matrix"

$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$
 n-dim vector

nxd matrix

$$\begin{array}{c}
 \begin{matrix} d \times n & n \times d \\ \downarrow & \uparrow \\ X^T X & w \end{matrix} \\
 \begin{matrix} d \times d & \in \mathbb{R}^d \\ \uparrow & \\ \end{matrix} \\
 = X^T X w = X^T y \\
 \begin{matrix} d \times n & n \\ \downarrow & \\ \end{matrix} \\
 \begin{matrix} \in \mathbb{R}^d & \\ \downarrow & \\ \end{matrix}
 \end{array}$$

So: $w = (X^T X)^{-1} X^T y$ Closed form solution!

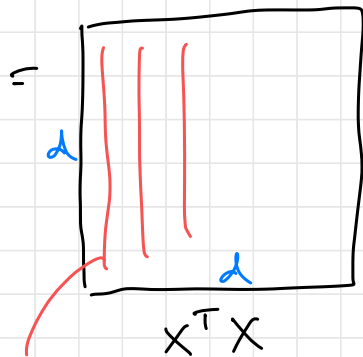
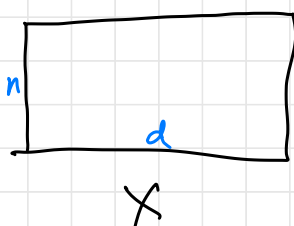
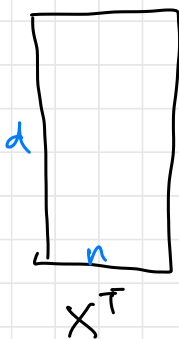
Question: What if $X^T X$ is not invertible?

Scenario 1: $n < d$

train examples

dimension of X 's
= # of features

(we have fewer examples than features)



Each column is linear combination of columns of X^T

every column of $X^T X$ lies in $\leq n$ -dim subspace of \mathbb{R}^d
 $\text{rank}(X^T X) \leq n < d$
 $\Rightarrow (X^T X)^{-1}$ does not exist!

Implication:

$$X^T X w = X^T y \text{ has } \underline{\text{many solutions}}$$

with many features & few examples,
easy to find w to get 0 train loss
many

we have high variance

→ Based on training data many equally good w 's

↪ Some will be better than others on test data

Common rule of thumb: "Always" have more examples than features

In practice, use pseudo inverse A^+

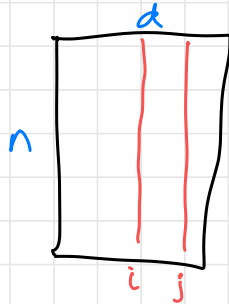
• $A^+ = A^{-1}$ when A^{-1} exists

• For $Ax = b$,

$x = A^+b$ will be a solution

So: Implement $w = (X^T X)^+ X^T y$

Scenario 2: Duplicated features



Suppose columns i & j are equal

then $X^T X$ not invertible

Intuitively: also high variance

$$w = [w_1 \dots w_i, \dots w_j, \dots]$$

+100 -100
-500 +500

many equally good w 's
hard to find best one

Rule of thumb: Avoid highly correlated/duplicate features