# 1/25/2024: Bias/variance, MAP, Normal Equations

All possible functions from X to y

Bias

Variance

Model Family = Set of functions that our algorithm can learn

★ Best possible predictor ("Bayes optimal")

Best model in family

what you actually learn

**Bias**: Error because assumptions of ML method don't exactly match real world

**Variance**: Error because what you learn is not best possible model in model family
↳ Because training data is always incomplete, never covers all possible cases

Bias + Variance = total error of model

## Reduce Bias
- Make fewer assumptions
  ⟺
  Make model family bigger

## Reduce Variance
Make it easier to find best model in family
① Make model family smaller
   (Regularization does this)
② Add training data

# What is the probabilistic story behind regularization?

**Idea:** Think about learning as usage of Bayes Rule.

## Bayesian probabilistic story:

① Exists a prior distribution over $w$ called $p(w)$

② $w$ gets sampled from $p(w)$

③ Dataset gets generated conditioned on $w$ from distribution $p(D|w)$

**Learning goal:** Infer most likely value of $w$

ie ⟶ maximize $p(w|D)$

called MAP
maximum
a posteriori

unknown
must be learned/inferred

observed

**By Bayes Rule:** $p(w|D) = \dfrac{\boxed{p(w)}\;\boxed{p(D|w)}}{\boxed{p(D)}}$

New!
Different choices of $p(w)$
give different types
of regularization

Does not depend
on $w$,
can be ignored

likelihood of
the data
ie. what we
maximize for MLE

For example: Let's assume each $w_j$ is Gaussian centered at 0

in particular: $p(w_j) = \dfrac{1}{\sigma \sqrt{2\pi}}\, e^{-w_j^2/2\sigma^2}$

Assume
mean is 0

constant
variance

Overall: $p(w)$ is just $\displaystyle\prod_{j=1}^{d} p(w_j)$

$$\max_{w} \; p(w \mid D) = \max_{w} \; p(w) \, p(D \mid w)$$

$$= \max_{w} \; \log p(w) + \log p(D \mid w)$$

$$= \max_{w} \; \sum_{j=1}^{d} \log p(w_j) + \log p(D \mid w)$$

$$= \max_{w} \; \text{constant} + \boxed{\sum_{j=1}^{d} \frac{-w_j^2}{2\sigma^2}} + \log p(D \mid w)$$

$$= \max_{w} \; -\frac{1}{2\sigma^2} \|w\|^2 + \log p(D \mid w)$$

$$= \min_{w} \; \boxed{\frac{1}{2\sigma^2} \|w\|^2} \; \boxed{- \log p(D \mid w)}$$

$\underbrace{\phantom{xxxxxxxx}}$ L2 regularization
with strength
$\lambda = \dfrac{1}{2\sigma^2}$

$\underbrace{\phantom{xxxxxxxx}}$ MLE objective

→ If $\sigma^2$ small, $\lambda$ large
→ If $\sigma^2$ large, $\lambda$ small

---

## Closed form Solution for Linear Regression
### ("Normal Equations")

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

Solve for $w$ where $\frac{d}{dw_j} = 0$ for all $j$

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

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

define

$$X = \begin{bmatrix} - x^{(1)} - \\ \vdots \\ - x^{(n)} - \end{bmatrix} \qquad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$n \times d$ matrix

$n$-dim vector

$$X^T X w = X^T y$$

$d \times n$    $n \times d$    $d \times n$    $n\text{-dim}$

$= d \times d$

$d\text{-dim vector}$

**Solution:** $\boxed{w = (X^T X)^{-1} X^T y}$    Closed-form Solution!

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

**Scenario 1:**    $n < d$      (we have too few examples compared to # of features)

# train examples     # features



$d$   $X^T$      $X$    $n$   =    $d$

$n$        $d$        $d$

Each column of $X^T X$ is result of $X^T \cdot$ some vector, which is a linear combo of columns of $X^T$

But $X^T$ only has $n$ columns, so all of $X^T X$'s columns lie in $n$-dimensional subspace

i.e. $\text{rank}(X^T X) \leq n < d$

So $X^T X$ is not invertible

Implication: $X^T X w = X^T y$ has many solutions

there's many $w$'s that perfectly fit training data
we don't know which one is actually best!
ie. we have high variance

Rule of thumb: Want to have $n > d$
more training examples than features

In practice: we can use pseudoinverse of $X^T X$

The pseudoinverse of a matrix $A$, denoted $A^+$:
- $A^+ = A^{-1}$ if $A$ is invertible
- For any equation $Ax = b$,
$$x = A^+ b \quad \text{is a solution}$$

## Scenario 2: Duplicated features



Suppose features $i$ & $j$ are identical
then, $X^T X$ is not invertible!

Intuitively: $w$ is under-determined

$$w = [w_1, \ldots, w_i, \ldots, w_j, \ldots, w_d]$$

$+100$    $-100$
$-500$    $+500$

] all $w$'s with identical performance

Another case of high variance!
Rule of thumb: Avoid (near-) duplicate features