

# Bias - Variance & Linear Regression

Xuan Zhang  
Sep. 9, 2022

I. Empirical Risk Minimization (ERM)

II. Bias-Variance decomposition

III. Linear Regression

VI. MLE & MAP

I.E.R.M

# Empirical Risk Minimization

Population Risk / Expected Risk

$$R(f) = \mathbb{E}_{XY} [\text{Loss}(Y, f(X))]$$

Bayes Optimal Rule

$$f^* = \underset{f}{\operatorname{argmin}} R(f)$$

Empirical Risk

$$R(f, D) = \frac{1}{N} \sum_{n=1}^N \text{Loss}(y_n, f(x_n))$$

Empirical Risk Minimization

$$\hat{f}_{\text{ERM}} = \underset{f}{\operatorname{argmin}} R(f, D)$$



# Why does this approximation work?

Law of large numbers: sample average converges to the expected value as sample size approaches  $\infty$ .

$$\bar{X}_N \rightarrow \mu \text{ as } N \rightarrow \infty.$$

$$\frac{1}{N} \sum_{n=1}^N \text{loss}(y_n, f(x_n)) \xrightarrow[\text{numbers}]{\text{Law of Large}} \mathbb{E}_{X,Y} [\text{loss}(Y, f(X))] = \int \text{loss}(Y, f(x)) p(x, y) dx, y$$

# Excess Error Decomposition

$H$ : our hypothesis space

$R(f)$ , <sup>expected risk /</sup> population risk

$R(f, D)$ : empirical risk

$f^{**} = \operatorname{argmin}_f R(f)$ : the function that achieves the minimal possible population risk.

$f^* = \operatorname{argmin}_{f \in H} R(f)$ : the function that achieves the minimal possible population risk in our hypothesis space.

$f_N^* = \operatorname{argmin}_{f \in H} R(f, D)$ : the function that achieves the minimal empirical risk in our hypothesis space.

excess error:

$$\mathbb{E}[R(f_N^*) - R(f^{**})] = \underbrace{\mathbb{E}[R(f^*) - R(f^{**})]}_{\text{approximation error}} + \underbrace{\mathbb{E}[R(f_N^*) - R(f^*)]}_{\text{estimation error / generation error}}$$

# Excess Error Decomposition

$$\mathbb{E}[R(f_N^*) - R(f^{**})] = \underbrace{\mathbb{E}[R(f^*) - R(f^{**})]}_{\text{approximation error}} + \underbrace{\mathbb{E}[R(f_N^*) - R(f^*)]}_{\text{estimation error / generation error}}$$

**approximation error:** (determined by the capacity of  $\mathcal{H}$ )  
measures how closely our hypothesis space  $\mathcal{H}$  can model the true optimal function  $f^{**}$ .

**generation error:** (determined by  $N$  & the capacity of  $\mathcal{H}$ )  
measures the difference in estimated risk due to having a finite training set.

II. Bias - Variance

Decomposition

# Expected Squared Loss Decomposition

Within a hypothesis space  $\mathcal{H}$ , e.g. the space of regression functions,

$f(x)$ : prediction function

$y$ : provided label (noisy)

$h(x)$ : true prediction function, true label

$$h(x) = \mathbb{E}[y|x] = \int y p(y|x) dy$$

expected squared loss:

$$\mathbb{E}[L] = \mathbb{E}[(f(x) - y)^2]$$

error from

model

$$= \mathbb{E}[(f(x) - h(x))^2] + \mathbb{E}[(h(x) - y)^2]$$

depends on the choice for  $f(x)$ .  
we want to find a  $f(x)$   
to minimize this term.

data

arises from the intrinsic  
noise on the data

# Expected Squared Loss Decomposition

(Derivation)

$$\begin{aligned} E[L] &= E[(f(x) - y)^2] \\ &= E[(f(x) - h(x) + h(x) - y)^2] \\ &= E[(f(x) - h(x))^2] + E[(h(x) - y)^2] \\ &\quad + 2E[(f(x) - h(x))(h(x) - y)] \quad \textcircled{1} \end{aligned}$$

$$\begin{aligned} \textcircled{1} &= E_{x,y} [(f(x) - E[y|x]) (E[y|x] - y)] \\ &= E_x [(f(x) - E[y|x]) \underline{E_{y|x} [E[y|x] - y]}] \quad \textcircled{2} \end{aligned}$$

$$\begin{aligned} \textcircled{2} &= E_{y|x} [y|x] - E_{y|x} [y] = 0 \\ &\quad \downarrow \qquad \qquad \downarrow \\ &\quad \int y p(y|x) dy \quad \stackrel{\text{same}}{=} \quad \int y p(y|x) dy \end{aligned}$$

# Bias - Variance Decomposition

expected squared loss:

$$\mathbb{E}[L] = \mathbb{E}[(f(x) - y)^2]$$

$$= \mathbb{E}[(f(x) - h(x))^2] + \mathbb{E}[(h(x) - y)^2]$$

error from

model

data

depends on the choice for  $f(x)$ .  
we want to find a  $f(x)$   
to minimize this term.

arises from the intrinsic  
noise on the data

first item depends on dataset  $D$ :

$$\mathbb{E}_D[(f(x; D) - h(x))^2] = (\mathbb{E}_D[f(x; D)] - h(x))^2 + \mathbb{E}_D[(f(x; D) - \mathbb{E}_D[f(x; D)])^2]$$

model risk (bias)<sup>2</sup> variance

$$\text{expected squared loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

model risk

# Bias - Variance Decomposition (Derivation)

model risk

$$\mathbb{E}_D [(f(x; D) - h(x))^2]$$

subtract and add  $\mathbb{E}_D [f(x; D)]$

$$= \mathbb{E}_D [(f(x; D) - \mathbb{E}_D [f(x; D)] + \mathbb{E}_D [f(x; D)] - h(x))^2]$$

variance

$$= \mathbb{E}_D [(f(x; D) - \mathbb{E}_D [f(x; D)])^2] + \mathbb{E}_D [(\mathbb{E}_D [f(x; D)] - h(x))^2]$$

$$+ \underbrace{2 \mathbb{E}_D [(f(x; D) - \mathbb{E}_D [f(x; D)]) (\mathbb{E}_D [f(x; D)] - h(x))]}_{=0}$$

$$\textcircled{1} = (\mathbb{E}_D [f(x; D)] - h(x))^2$$

(bias)<sup>2</sup>



# Bias - Variance Tradeoff

$$\underbrace{\mathbb{E}_D [(f(x; D) - h(x))^2]}_{\text{model risk}} = \underbrace{(\mathbb{E}_D [f(x; D)] - h(x))^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_D [(f(x; D) - \mathbb{E}_D [f(x; D)])^2]}_{\text{variance}}$$

**bias**: represents the extent to which the average prediction over all data sets differs from the best prediction function.

**variance**: measures the extent to which the solutions for individual data sets vary around the average.  
i.e. measures the extent to which the function  $f(x; D)$  is sensitive to the particular choice of data set.

# Bias - Variance Tradeoff



The world is  
black and white

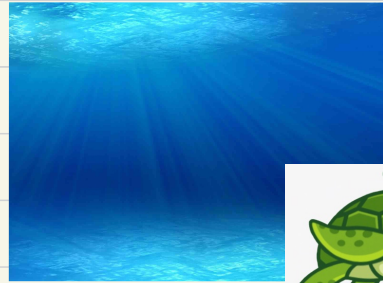


Skate

high bias (underfitting)



The world  
is green.



The world  
is blue.



high variance (overfitting)

\* Skate is the only animal that has been confirmed to see only in black and white.

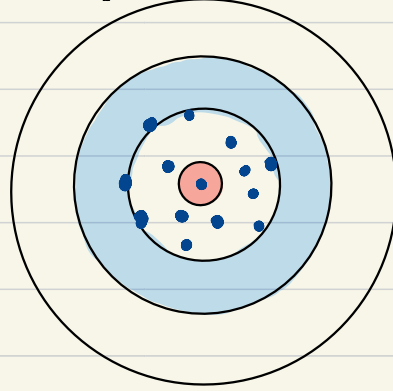
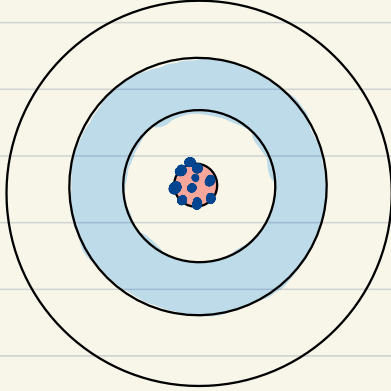
\* That turtle is purely made up.

# Bias - Variance Tradeoff

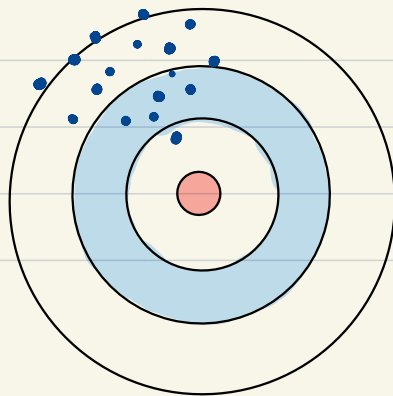
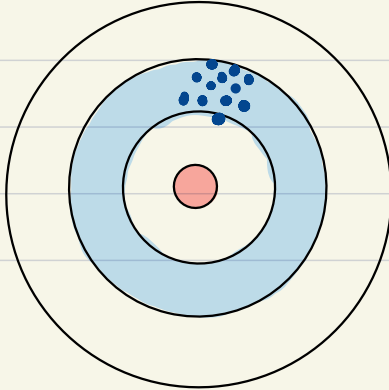
low variance

high variance

low bias



high bias



# III. Linear Regression

# Linear Regression

$$p(y|x, \theta) = \mathcal{N}(y | w^T x, \sigma^2)$$

$$x_i \in \mathbb{R}^D$$

$$y \in \mathbb{R}^N$$

$$w \in \mathbb{R}^D$$

Likelihood  $\prod_{i=1}^N \mathcal{N}(y_i | w^T x_i, \sigma^2)$

log likelihood  $\sum_{i=1}^N \log \left[ \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left( -\frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right) \right]$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^T x_i)^2 - \frac{N}{2} \log(2\pi\sigma^2)$$

maximize log likelihood  $\Leftrightarrow$  minimize  $\frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2$

# Least Squares Estimation

residual sum of squares (RSS) :

$$X \in \mathbb{R}^{N \times D}$$

$$RSS(w) = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2$$

$$= \frac{1}{2} \|Xw - y\|_2^2$$

$$= \frac{1}{2} (Xw - y)^T (Xw - y)$$

optimize :

$$\nabla_w RSS(w) = X^T X w - X^T y = 0$$

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

least squares solution:

$$w = (X^T X)^{-1} X^T y$$

Recall 2-norm (Euclidean norm)

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad \text{note that } \|x\|_2^2 = x^T x$$

# Least Squares Estimation (Derivation)

$$RSS(w) = \frac{1}{2} (Xw - y)^T (Xw - y)$$

$$= \frac{1}{2} w^T X^T X w - \frac{1}{2} y^T X w - \frac{1}{2} w^T X^T y + \frac{1}{2} y^T y$$

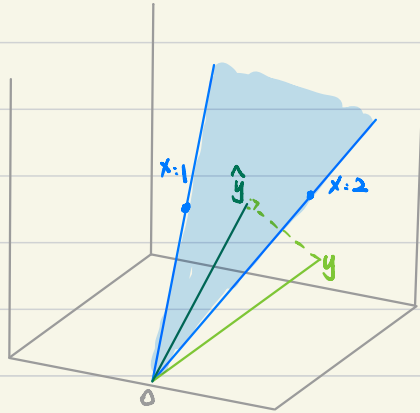
$$\nabla_w RSS(w) = \frac{1}{2} \frac{\partial w^T X^T X w}{\partial w} \textcircled{1} - \frac{1}{2} \frac{\partial y^T X w}{\partial w} \textcircled{2} - \frac{1}{2} \frac{\partial w^T X^T y}{\partial w} \textcircled{3}$$

$$= \frac{1}{2} (X^T X + X^T X) w - \frac{1}{2} X^T y - \frac{1}{2} X^T y$$

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

Vector derivatives:  $\textcircled{1} \frac{\partial x^T A x}{\partial x} = (A + A^T) x$     $\textcircled{2} \frac{\partial a^T x}{\partial x} = a$     $\textcircled{3} \frac{\partial x^T a}{\partial x} = a$

# Geometric interpretation of least squares



$\hat{y}$  lies in the linear subspace spanned by  $X$ :  $\hat{y} = Xw$

we want to find

$$\hat{y}^* = \underset{\hat{y} \in \text{span}\{x_{:1}, \dots, x_{:d}\}}{\text{argmin}} \|y - \hat{y}\|_2 \quad x_{:,d}: d\text{th column of } X$$

To minimize  $\|y - \hat{y}\|_2$ , we want  $(y - \hat{y})$  to be orthogonal to every column of  $X$ :

$$X^T(y - Xw) = 0 \quad \Rightarrow \quad w = (X^T X)^{-1} X^T y$$



VIMLE & MAP

# MLE

maximum likelihood estimation (MLE):

Pick the parameters that assign the highest probability to data.

$$\hat{\theta}_{\text{MLE}} \triangleq \underset{\theta}{\operatorname{argmax}} p(D|\theta)$$

likelihood  $p(D|\theta) = \prod_{i=1}^N p(y_i|x_i, \theta)$  iid assumption

log likelihood  $LL(\theta) \triangleq \log p(D|\theta) = \sum_{i=1}^N \log p(y_i|x_i, \theta)$

negative log likelihood  $NLL(\theta) \triangleq -\log p(D|\theta) = -\sum_{i=1}^N \log p(y_i|x_i, \theta)$

# MAP

maximum a posterior estimation (MAP):

$$\hat{\theta}_{\text{map}} = \operatorname{argmax}_{\theta} \log p(\theta | D)$$

$$= \operatorname{argmax}_{\theta} [\log p(D | \theta) + \underbrace{\log p(\theta)}_{\text{prior}} - \text{const}]$$

\* Will learn more about MAP when we learn Regularization.



Questions?