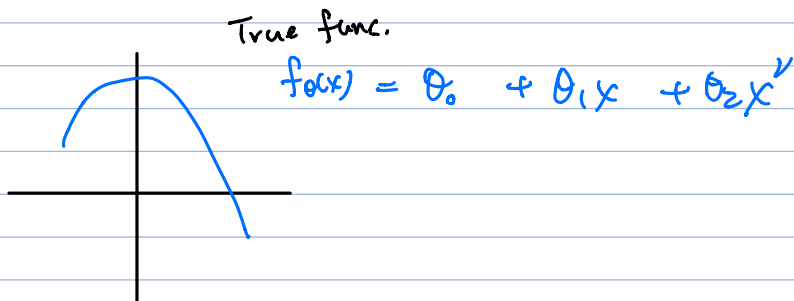
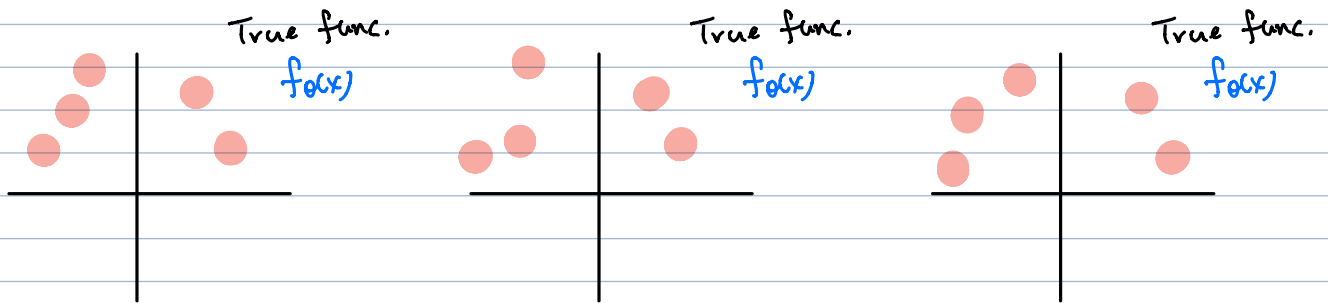


Overview

- Bias and Variance

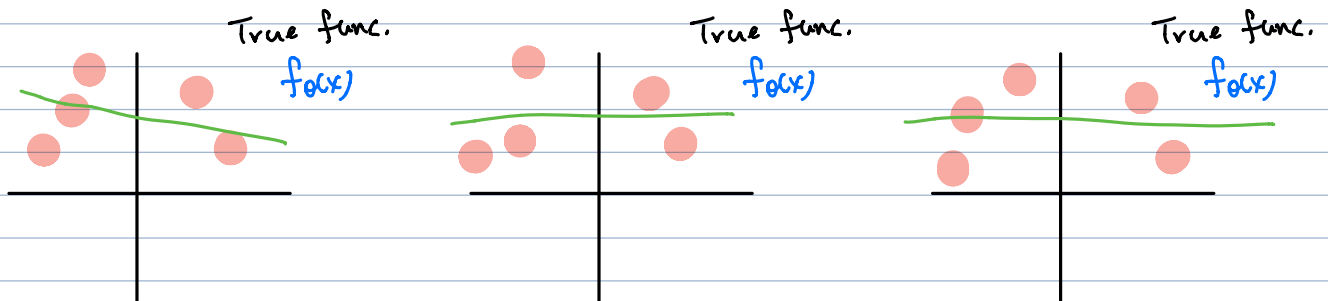


Samples



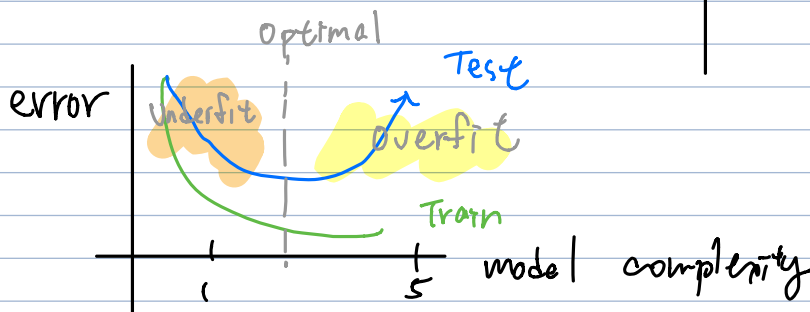
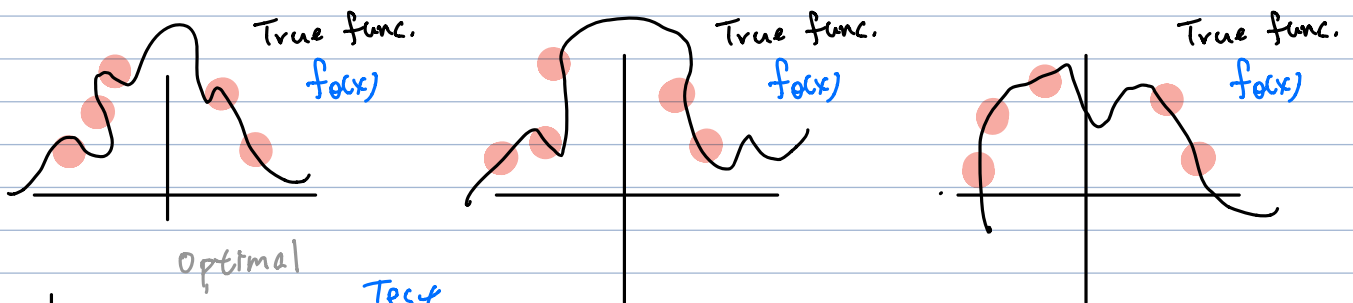
Fit a linear line

Bias, Underfit



Fit a degree 5 polynomial

High Variance, overfit



Classical Bias-Variance

Bias - Variance

$$y = x\theta + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

$\begin{matrix} \nearrow \mathbb{R}^d \text{ features} \\ \searrow \end{matrix}$

$\hookrightarrow \epsilon \in \mathbb{R} \quad \hookrightarrow \text{model } \epsilon \in \mathbb{R}^d$

Procedure

1. Draw n points $(x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$ as set S

$$y^{(i)} = x^{(i)}\theta_* + \epsilon^{(i)} \quad \epsilon^{(i)} \sim N(0, \sigma^2)$$

2. Train a linear model on $S = h_S: \mathbb{R}^d \rightarrow \mathbb{R}$

3. Draw a test sample (x, y)

$$y = h_{\theta_*}(x) + \epsilon$$

4. Measure $\mathbb{E}_{\epsilon, S} [(h_S(x) - y)^2]$

Goal: Decomposition of Error

$$\mathbb{E}_{\epsilon, S} [(h_S(x) - (h_{\theta_*}(x) + \epsilon))^2]$$

$$= \underbrace{\mathbb{E}[\epsilon^2]}_{\sigma^2} + \underbrace{\mathbb{E}[(h_S(x) - h_{\theta_*}(x))^2]}_{\text{Depend on } S, \text{ (training set)}} + \underbrace{\mathbb{E}[\epsilon(h_S(x) - h_{\theta_*}(x))]}_{\text{independent}}$$

$$\mathbb{E}[\epsilon] = 0$$

Unavailable error

Depend on S , (training set)

independent

$$h_{\text{Avg}}(x) \triangleq \mathbb{E}_S [h_S(x)] \quad \text{"long run average training err on } S \text{"}$$

(select S many times)

$$\mathbb{E} [(h_S(x) - h_{\text{Avg}}(x) + h_{\text{Avg}}(x) - h_{\theta_*}(x))^2]$$

$$= \underbrace{\mathbb{E}[(h_S(x) - h_{\text{Avg}}(x))^2]}_{\text{Variance}} + \underbrace{\mathbb{E}_S[(h_{\text{Avg}}(x) - h_{\theta_*}(x))^2]}_{\text{Bias}^2} + \underbrace{\mathbb{E}_S[(h_S(x) - h_{\theta_*}(x))(h_{\text{Avg}}(x) - h_{\theta_*}(x))]}_{0}$$

Variance

Bias²

VAR(S)

Does not depend on S

Examples	Bias	Variance
linear	high	low
degree 5	0	high

Reduce Variance that depends on the training set
Regularization

Most Classical case (linear)

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (x^{(i)T} \theta - y^{(i)})^2 + \underbrace{\frac{\lambda}{2} \|\theta\|_2^2}_{\text{Penalty of complex}}$$

Regularization Parameter

Solution $x \in \mathbb{R}^{n \times d}$

$$x^T x \theta - x^T y + \lambda \theta = 0$$

$$\Rightarrow \theta = (x^T x + \lambda I)^{-1} x^T y$$

(Normal eq)

What if $x^T x$ is not full rank? $n < d \Rightarrow \text{Rank}(x^T x) < d$
 if we set $\lambda = 0$

$$x^T x \theta = x^T y \quad v \in \text{Null}(x)$$

$$x^T x (\theta + v) = x^T x \theta = x^T y \quad (\text{Not unique})$$

if $\lambda > 0$. $x^T x + \lambda I$ is full rank

eigenvalues $(x^T x) = \sigma_1^2 \geq \sigma_2^2 \geq \dots \geq 0$

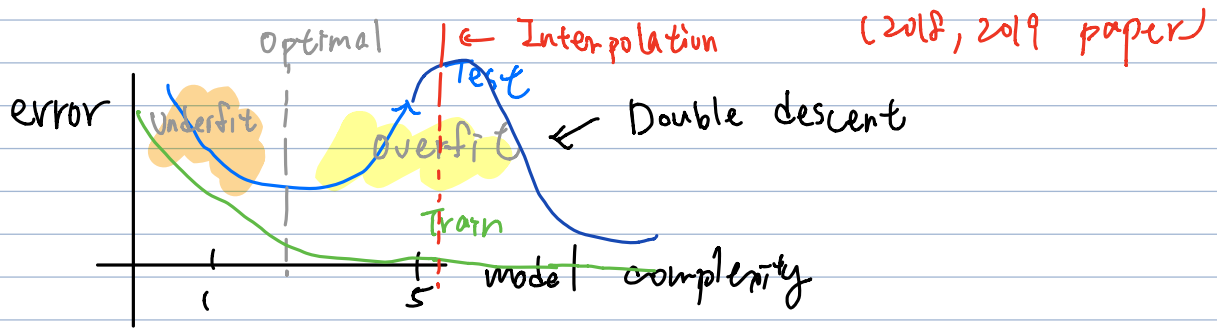
$$(x^T x + \lambda I) = \sigma_1^2 + \lambda, \sigma_2^2 + \lambda, \dots, \sigma_n^2 + \lambda \geq 0$$

$$\theta_\lambda = (x^T x + \lambda I)^{-1} x^T y$$

$$\text{Var}(s) = \mathbb{E}_s [(\theta_\lambda x - \mathbb{E}_s [\theta_\lambda x])^2]$$

$$\approx \mathbb{E}[(\theta_\lambda - \mathbb{E}[\theta_\lambda])^2]$$

As λ increases, variance decreases



Picking hyperparameters

Three set of data

daily

Train set \rightarrow fit parameters

Dev set \rightarrow tune parameters

Test set \rightarrow blind

for degree $d \in \{1, \dots, m\}$

[Train model w/ parameters d on the train set
Score on the dev set

\Rightarrow Pick best score

How to pick λ ? $\lambda \in \{0, 10^{-4}, 10^{-3}, \dots, 10^{-1}\}$

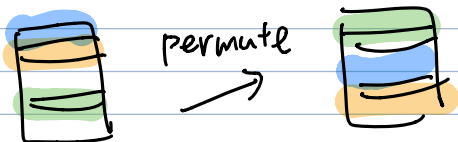
Improvements

Data Efficiency: k -fold cross validation

Compute Efficiency: Many hyperparameters
"modern ML"

[Data] k -fold Cross Validation

1.



2.

Train	Score
$S_1 S_2 S_3$	S_4
$S_1 S_3 S_4$	S_2
\vdots	\vdots

3. Combine Score (Average)

Computational

Motivation: Dropout, Regularization

More Advanced (Amesim 15') * hyperband

Run all (5, 6, 7) \Rightarrow just for a few steps
score all of them

Pick $N/2$

then they run $N/2^k$ number of models