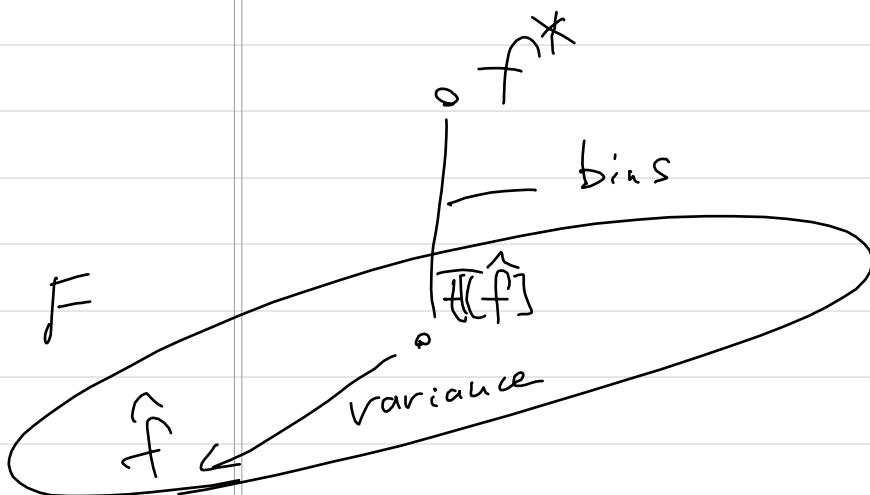


(a) Underfitting : condition that statistical model cannot capture feature and rules

Overfitting : condition that statistical model fits to not only feature but also noise

Bias : gap between expectation of estimate and true model



$R(f)$ is risk function

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

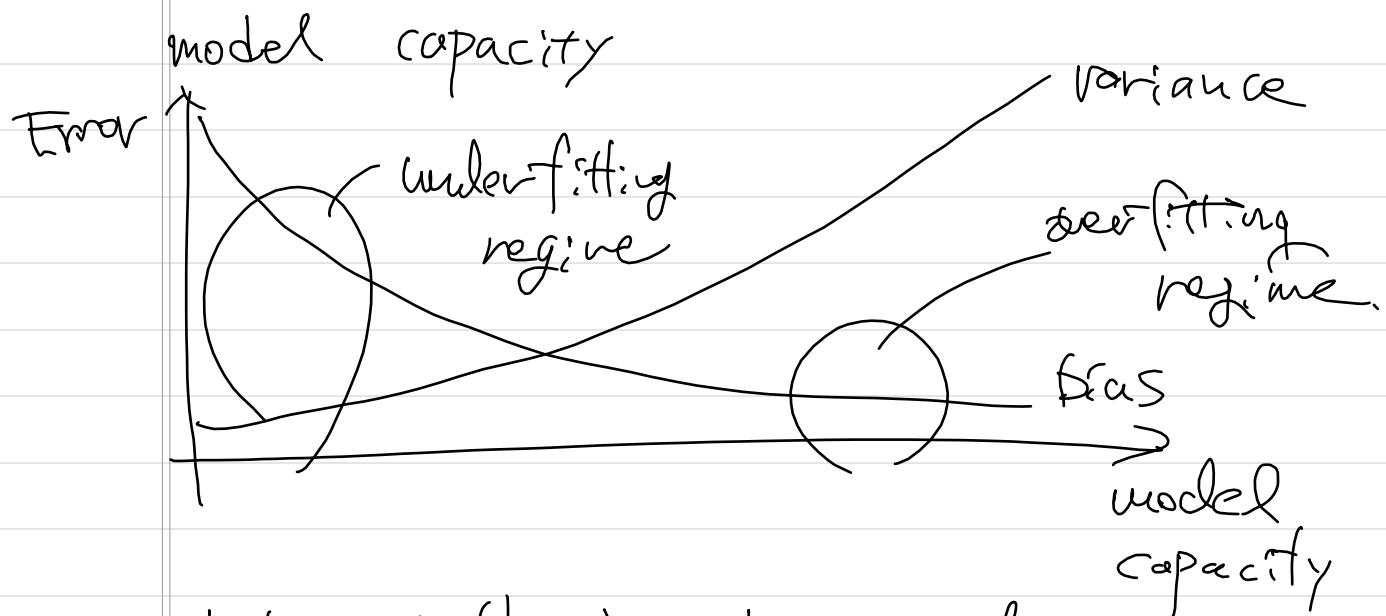
$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)$$

$$\text{Bias}[\hat{f}] = \mathbb{E}[\hat{f}] - f^*$$

$$\text{Variance}[\hat{f}] = \mathbb{E}[\hat{f}^2] - (\mathbb{E}[\hat{f}])^2$$

Bias = an error from erroneous assumption in the learning algorithm, high bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting)

Variance: an error from sensitivity to small fluctuations in the training set. High variance may result from an algorithm modeling the random noise in the training data (overfitting)



that controls the bias-and-variance trade-off,

21A-Q1

3

(cc) # Train Sample \rightarrow variance \downarrow ?

2 (a) Diff of param hyperparam

updated by training data set

choose by evaluation on validation data as model selection

which determine training behaviour and efficient capacity

b) (1) Linear Regression

$$\theta = [w + b]$$

$$\text{obj} = l(f(x), y) = \frac{1}{2} \|f(x) - y\|^2 + \lambda \|w^T w\|$$

$$f(x) = w^T x + b$$

$$w \in \mathbb{R}^{d \times d}, y \in \mathbb{R}^d, x \in \mathbb{R}^d, b \in \mathbb{R}^d.$$

↓ param

↓ hyperparam

(2-

λ : regularization

tion

loss func

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} l(\hat{y}, y)$$

η = step size

(2) MLP

$$\hat{y} = f(x; \theta)$$

$\theta \in \mathbb{R}^d$ model param

2(A-Q1)

4

$$\begin{array}{l} \theta \sim \text{Beta}(\alpha, \beta) \leftarrow \text{hyperparam} \\ X_i \sim \text{Bern}(\theta) \leftarrow \text{param} \end{array}$$

$$y \sim N(x\beta, \sigma^2 I)$$

$$\hat{\beta} = (X^T X + kI)^{-1} X^T y$$

β , θ \leftarrow param

k \leftarrow hyperparam

3

1. Predict likelihood

input given score

output likelihood.

① regression (logistic)

② MLP w/ sigmoid

2. Predict height

① linear regression
regression tree

2A - Q1

5

3. Unsupervised classification (text)

non-hierarchical
k-means (clustering)
hierarchical clustering (Ward method)

4

Cross Entropy and ML \bar{E} (relationship)

if we want to maximize log likelihood

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(y|x; \theta)$$

$$\Leftrightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log P(y|x; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} L(\theta)$$

we focus
one train sample

$$= \underset{\theta}{\operatorname{argmax}} \sum_{k=1}^K y_k \log p_k$$

$$= \underset{\theta}{\operatorname{argmin}} - \sum_{k=1}^K y_k \log p_k$$

$$\frac{1}{K} \prod_{k=1}^K p(y_i|x_i; \theta)^{y_i}$$

$$\frac{1}{K} \prod_{k=1}^K$$

it is equivalent $CLE(y_k, p_k)$

2(A-Q1)

6

15 $\mathcal{N} = \{(x_1, y_1), \dots, (x_n, y_n)\} \sim D$
 $S \subset \mathbb{R}^d \times \mathbb{R}_{-1, 1}$

(a) $\boxed{A | B} \leftarrow \text{data set}$

step 1 use A as train data

then evaluate $\hat{\theta}_A$ on data set B

Step 2 use B as vice-versa

estimated error = $\frac{1}{2} (l(B, \hat{\theta}_A) + l(A, \hat{\theta}_B))$

(b) CV error of learned classifier

an biased estimate of its true error?