# CSC 4020 Fundamental of Machine Learning: Bias-Variance Tradeoff

Baoyuan Wu
School of Data Science, CUHK-SZ

January 25/27, 2021

- We are provided by a training dataset $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$, which is drawn i.i.d. from some distribution $P(\mathcal{X}, \mathcal{Y})$

# Bias-variance tradeoff

- We are provided by a training dataset $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$, which is drawn i.i.d. from some distribution $P(\mathcal{X}, \mathcal{Y})$
- The relationship between the input features $\boldsymbol{x}$ and the output $y$ is

$$y = h(\boldsymbol{x}) + e, \ e \sim \mathcal{N}(0, \sigma^2), \tag{1}$$

$$p(y|\boldsymbol{x}) = \mathcal{N}(h(\boldsymbol{x}), \sigma^2 \mathbf{I}), \tag{2}$$

where $h(\boldsymbol{x})$ can be seen as the unknown target function and the mean of $p(y|\boldsymbol{x})$.

# Bias-variance tradeoff

$D \sim p(x, y)$

- We are provided by a training dataset $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^{n}$, which is drawn i.i.d. from some distribution $P(\mathcal{X}, \mathcal{Y})$

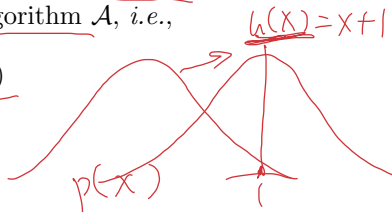- The relationship between the input features $\boldsymbol{x}$ and the output $y$ is

$$y = h(\boldsymbol{x}) + e, \ e \sim \mathcal{N}(0, \sigma^2), \tag{1}$$

$$p(y|\boldsymbol{x}) = \mathcal{N}(h(\boldsymbol{x}), \sigma^2 \mathbf{I}), \tag{2}$$

where $h(\boldsymbol{x})$ can be seen as the unknown target function and the mean of $p(y|\boldsymbol{x})$.

- The goal of machine learning is to learn a hypothesis function based on the training dataset $D$ using some learning algorithm $\mathcal{A}$, $i.e.$,

$$h_D = \mathcal{A}(D)$$

$h(x) = x + 1$

$p(x)$

# Bias-variance tradeoff

- **Expected hypothesis function** (given $\mathcal{A}$):

$$\bar{h} = E_{D \sim P^n}[h_D] = \int_D h_D \, p(D) \, dD$$

# Bias-variance tradeoff

- **Expected hypothesis function** (given $\mathcal{A}$):

$$\bar{h} = E_{D \sim P^n}[h_D] = \int_D h_D p(D) dD$$

- Given a test pair $(\boldsymbol{x}, y) \sim P(\mathcal{X}, \mathcal{Y})$ and $h_D$, the **expected test error** is defined as

$$E_{(\boldsymbol{x},y) \sim P}\left[(h_D(\boldsymbol{x}) - y)^2\right] = \int_{\boldsymbol{x}} \int_y (h_D(\boldsymbol{x}) - y)^2 p(x, y) d\boldsymbol{x} dy$$

# Bias-variance tradeoff

- **Expected hypothesis function** (given $\mathcal{A}$):

$$\bar{h} = E_{D \sim P^n}[h_D] = \int_D h_D p(D) dD$$

- Given a test pair $(\boldsymbol{x}, y) \sim P(\mathcal{X}, \mathcal{Y})$ and $h_D$, the **expected test error** is defined as

$$E_{(\boldsymbol{x}, y) \sim P}\big[(h_D(\boldsymbol{x}) - y)^2\big] = \int_{\boldsymbol{x}} \int_y (h_D(\boldsymbol{x}) - y)^2 p(x, y) d\boldsymbol{x} dy$$

- Given a test pair $(\boldsymbol{x}, y) \sim P(\mathcal{X}, \mathcal{Y})$ and $\mathcal{A}$, the **expected test error** is defined as

$$E_{(\boldsymbol{x}, y) \sim P, D \sim P^n}\big[(h_D(\boldsymbol{x}) - y)^2\big] = \int_D \int_{\boldsymbol{x}} \int_y (h_D(\boldsymbol{x}) - y)^2 p(x, y) d\boldsymbol{x} dy dD$$

# Bias-variance tradeoff

- **Expected hypothesis function** (given $\mathcal{A}$):

$$\bar{h} = E_{D \sim P^n}[h_D] = \int_D h_D p(D) dD$$

- Given a test pair $(\boldsymbol{x}, y) \sim P(\mathcal{X}, \mathcal{Y})$ and $h_D$, the **expected test error** is defined as

$$E_{(\boldsymbol{x},y) \sim P}\big[(h_D(\boldsymbol{x}) - y)^2\big] = \int_{\boldsymbol{x}} \int_y (h_D(\boldsymbol{x}) - y)^2 p(x, y) d\boldsymbol{x} dy$$

- Given a test pair $(\boldsymbol{x}, y) \sim P(\mathcal{X}, \mathcal{Y})$ and $\mathcal{A}$, the **expected test error** is defined as

$$E_{(\boldsymbol{x},y) \sim P, D \sim P^n}\big[(h_D(\boldsymbol{x}) - y)^2\big] = \int_D \int_{\boldsymbol{x}} \int_y (h_D(\boldsymbol{x}) - y)^2 p(x, y) d\boldsymbol{x} dy dD$$

- We are interested in evaluating the quality of a machine learning algorithm $\mathcal{A}$ with respect to a data distribution $P(\mathcal{X}, \mathcal{Y})$. In the following we will show that this expression decomposes into three meaningful terms.

- The **expected test error** can be decomposed as follows

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big] = E_{(\boldsymbol{x},y),D}\big[[h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x})) + (\bar{h}(\boldsymbol{x}) - y)]^2\big]$$
$$= E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + 2E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x}) - y)\big]$$
$$+ E_{(\boldsymbol{x},y),D}\big[(\bar{h}(\boldsymbol{x}) - y)^2\big]$$

$$\left( \int_D h_D(x)\, dD - \bar{h} \right)$$

$$0$$

# Bias-variance tradeoff

- The **expected test error** can be decomposed as follows

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big] = E_{(\boldsymbol{x},y),D}\big[[h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x})) + (\bar{h}(\boldsymbol{x}) - y)]^2\big]$$
$$= E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + 2E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x}) - y)\big]$$
$$+ E_{(\boldsymbol{x},y),D}\big[(\bar{h}(\boldsymbol{x}) - y)^2\big]$$

- We have

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x}) - y)\big]$$
$$= E_{(\boldsymbol{x},y)}\big[E_D[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x}) - y)]\big]$$
$$= E_{(\boldsymbol{x},y)}\big[(E_D[h_D(\boldsymbol{x})] - \bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x}) - y)) = 0$$

# Bias-variance tradeoff

- The **expected test error** can be decomposed as follows

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x})-y)^2\big] = E_{(\boldsymbol{x},y),D}\big[[h_D(\boldsymbol{x})-\bar{h}(\boldsymbol{x}))+(\bar{h}(\boldsymbol{x})-y)]^2\big]$$
$$=E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x})-\bar{h}(\boldsymbol{x}))^2\big] + 2E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x})-\bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x})-y)\big]$$
$$+ E_{(\boldsymbol{x},y),D}\big[(\bar{h}(\boldsymbol{x})-y)^2\big]$$

- We have

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x})-\bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x})-y)\big]$$
$$=E_{(\boldsymbol{x},y)}\big[E_D[(h_D(\boldsymbol{x})-\bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x})-y)]\big]$$
$$=E_{(\boldsymbol{x},y)}\big[(E_D[h_D(\boldsymbol{x})]-\bar{h}(\boldsymbol{x}))(\bar{h}(\boldsymbol{x})-y))=0$$

- Then, we have

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x})-y)^2\big] = E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x})-\bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y),D}\big[(\bar{h}(\boldsymbol{x})-y)^2\big]$$

- We also have

$y = h(x) + e$

$P(y | x) = N(h(x), \cdot)$

$$E_{(\boldsymbol{x},y),D}\left[(\bar{h}(\boldsymbol{x}) - y)^2\right] = E_{(\boldsymbol{x},y),D}\left[[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x})) + (h(x) - y)]^2\right]$$

$$=E_{\boldsymbol{x},y}\left[(h(\boldsymbol{x}) - y)^2\right] + E_{\boldsymbol{x},y}\left[\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\right] + 2E_{\boldsymbol{x},y}\left[(h(\boldsymbol{x}) - y)(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))\right]$$

$$=E_{\boldsymbol{x},y}\left[(h(\boldsymbol{x}) - y)^2\right] + E_{\boldsymbol{x},y}\left[\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\right] \tag{3}$$

$\int_{y,x} y \, P(y|x) \, dy$

$= h(x)$

# Bias-variance tradeoff

- We also have

$$E_{(\boldsymbol{x},y),D}\big[(\bar{h}(\boldsymbol{x}) - y)^2\big] = E_{(\boldsymbol{x},y),D}\big[[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x})) + (h(x) - y)]^2\big]$$
$$= E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big] + E_{\boldsymbol{x},y}\big[\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + 2E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))\big]$$
$$= E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big] + E_{\boldsymbol{x},y}\big[\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] \tag{3}$$

- Finally, we have

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big]$$
$$= E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y)}\big[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big]$$

# Bias-variance tradeoff

- We also have

$$E_{(\boldsymbol{x},y),D}\big[(\bar{h}(\boldsymbol{x}) - y)^2\big] = E_{(\boldsymbol{x},y),D}\big[[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x})) + (h(x) - y)]^2\big]$$
$$=E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big] + E_{\boldsymbol{x},y}\big[\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + 2E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))\big]$$
$$=E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big] + E_{\boldsymbol{x},y}\big[\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] \qquad (3)$$

- Finally, we have

$$y = h(x) + \underbrace{e}_{\sigma}$$

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big]$$
$$=E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y)}\big[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big]$$

$$\sigma^2$$

- Above three terms are **variance, bias, noise**, respectively.

# Bias-variance tradeoff

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big]$$
$$= E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y)}\big[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big]$$

# Bias-variance tradeoff

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big]$$
$$=E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y)}\big[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big]$$

- **variance**: Captures how much your classifier changes if you train on a different training set. How "over-specialized" is your classifier to a particular training set (overfitting)? If we have the best possible model for our training data, how far off are we from the average classifier?

# Bias-variance tradeoff

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big]$$
$$= E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y)}\big[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big]$$
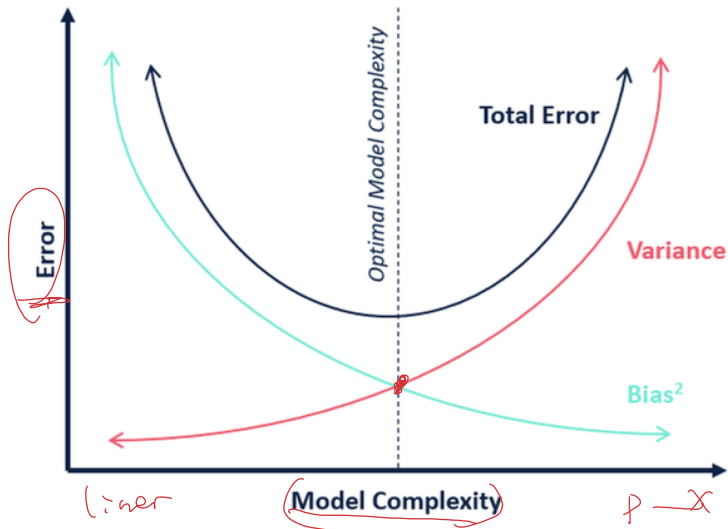
- **variance**: Captures how much your classifier changes if you train on a different training set. How "over-specialized" is your classifier to a particular training set (overfitting)? If we have the best possible model for our training data, how far off are we from the average classifier?
- **Bias**: What is the inherent error that you obtain from your classifier even with infinite training data? This is due to your classifier being "biased" to a particular kind of solution (*e.g.*, linear classifier). In other words, bias is inherent to your model.

$$E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - y)^2\big]$$
$$=E_{(\boldsymbol{x},y),D}\big[(h_D(\boldsymbol{x}) - \bar{h}(\boldsymbol{x}))^2\big] + E_{(\boldsymbol{x},y)}\big[(\bar{h}(\boldsymbol{x}) - h(\boldsymbol{x}))^2\big] + E_{\boldsymbol{x},y}\big[(h(\boldsymbol{x}) - y)^2\big]$$

- **variance**: Captures how much your classifier changes if you train on a different training set. How "over-specialized" is your classifier to a particular training set (overfitting)? If we have the best possible model for our training data, how far off are we from the average classifier?

- **Bias**: What is the inherent error that you obtain from your classifier even with infinite training data? This is due to your classifier being "biased" to a particular kind of solution (*e.g.*, linear classifier). In other words, bias is inherent to your model.

- **Noise**: How big is the data-intrinsic noise? This error measures ambiguity due to your data distribution and feature representation. You can never beat this, it is an aspect of the data.
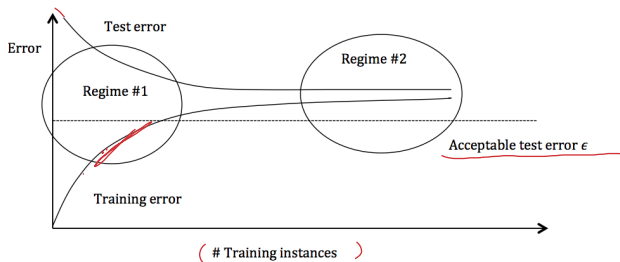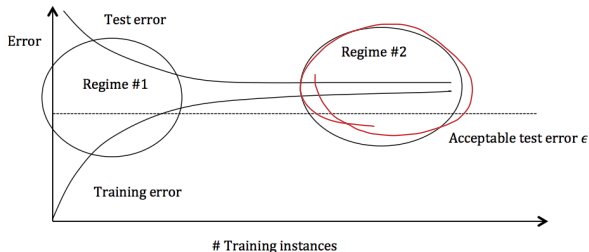
$$e \sim N(0, \delta^2)$$

# Bias-variance tradeoff

# Bias-variance tradeoff

# Bias-variance tradeoff



## Regime 1 (High Variance)

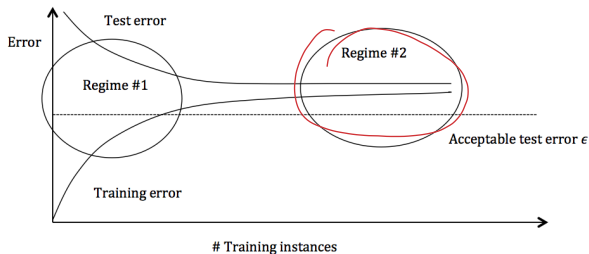In the first regime, the cause of the poor performance is high variance.

**Symptoms**:

1. Training error is much lower than test error
2. Training error is lower than $\epsilon$
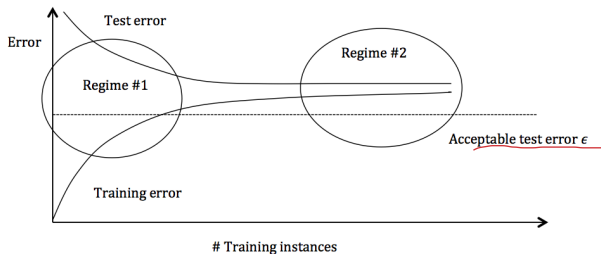3. Test error is above $\epsilon$

**Remedies**:

- Add more training data
- Reduce model complexity -- complex models are prone to high variance

# Bias-variance tradeoff

# Bias-variance tradeoff



## Regime 2 (High Bias)

Unlike the first regime, the second regime indicates high bias: the model being used is not robust enough to produce an accurate prediction.
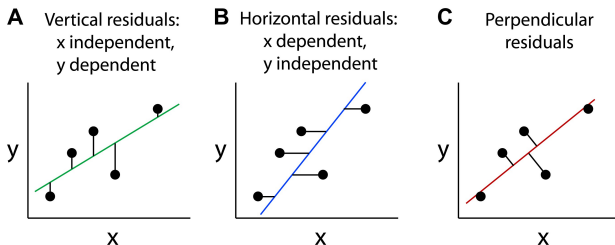
**Symptoms**:

1. Training error is higher than $\epsilon$

**Remedies**:

- Use more complex model (e.g. kernelize, use non-linear models)
- Add features

# Bias-variance tradeoff

More details can be found at `https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote12.html`

## Quiz

Q1: In the cost function of least squares estimation for linear regression, which residual we use? ( )



Q2: Suppose you have fitted a complex regression model on a dataset. Now, you are using Ridge regression with the penalty $\lambda$, *i.e.*, $\min(\boldsymbol{\theta}^\top \boldsymbol{x} - y)^2 + \lambda\|\boldsymbol{\theta}\|_2^2$. Choose the option which describes bias in best manner.

- A. In case of very large $\lambda$, bias is low
- B. In case of very large $\lambda$, bias is high
- C. We can't say about bias
- D. None of these