# General Guidance

Hung-yi Lee 李宏毅

# Framework of ML

Training data: $\left\{\left(\boldsymbol{x^1}, \hat{y}^1\right), \left(\boldsymbol{x^2}, \hat{y}^2\right), \dots, \left(\boldsymbol{x^N}, \hat{y}^N\right)\right\}$

Testing data: $\left\{\boldsymbol{x^{N+1}}, \boldsymbol{x^{N+2}}, \dots, \boldsymbol{x^{N+M}}\right\}$

*Speech Recognition*

$\boldsymbol{x}$:   $\hat{y}$: phoneme

*Speaker Recognition*

$\boldsymbol{x}$:   $\hat{y}$: John

(speaker)

*Image Recognition*

$\boldsymbol{x}$:   $\hat{y}$: soup

*Machine Translation*

$\boldsymbol{x}$: 痛みを知れ

$\hat{y}$: 了解痛苦吧

# Framework of ML

Training data: $\{(\boldsymbol{x^1}, \hat{y}^1), (\boldsymbol{x^2}, \hat{y}^2), \dots, (\boldsymbol{x^N}, \hat{y}^N)\}$

Training:



Step 1: function with unknown

Step 2: define loss from training data

Step 3: optimization

$$y = f_{\boldsymbol{\theta}}(\boldsymbol{x})$$

$$L(\boldsymbol{\theta})$$

$$\boldsymbol{\theta}^* = arg\min_{\boldsymbol{\theta}} L$$

Testing data: $\{\boldsymbol{x^{N+1}}, \boldsymbol{x^{N+2}}, \dots, \boldsymbol{x^{N+M}}\}$

Use $y = f_{\boldsymbol{\theta}^*}(\boldsymbol{x})$ to label the testing data

$\{y^{N+1}, y^{N+2}, \dots, y^{N+M}\}$ ⟶ Upload to Kaggle

**General Guide**

loss on training data

large — model bias → make your model complex

large — optimization → Next Lecture

small — loss on testing data

large — overfitting → more training data (not in HWs) data augmentation / make your model simpler

large — mismatch → Not in HWs, except HW 11

small → 😊

trade-off

Split your training data into training set and validation set for model selection

# _Model Bias_

- The model is too simple.

$$f_{\theta^1}(x) \qquad y = f_{\theta}(x)$$

$$f_{\theta^2}(x)$$

find a needle in a haystack …

… but there is no needle

$$f_{\theta^*}(x)$$
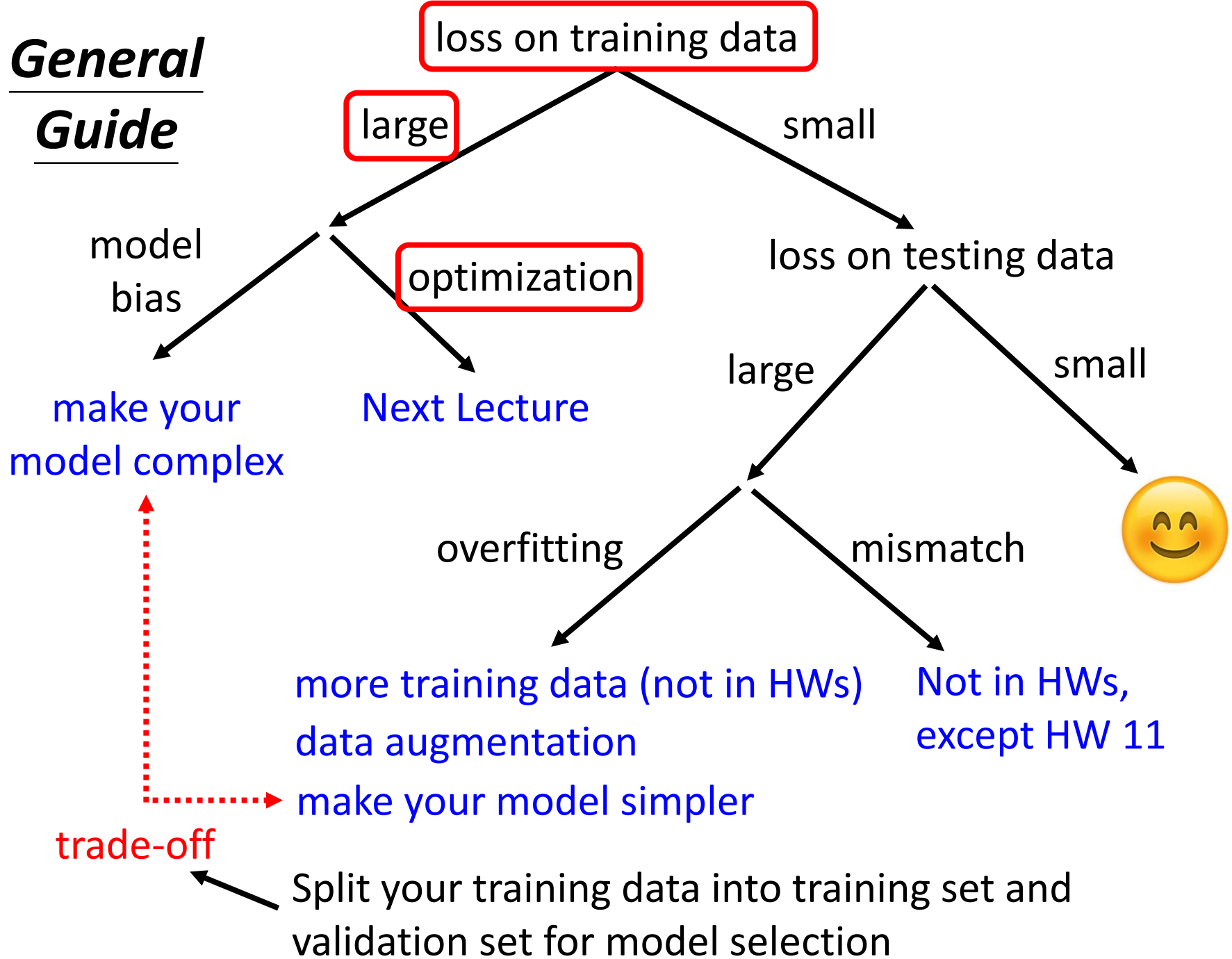
too small …  $f^*(x)$ small loss

- Solution: redesign your model to make it more flexible

$$y = b + wx_1 \xrightarrow{\text{More features}} y = b + \sum_{j=1}^{56} w_j\, x_j$$
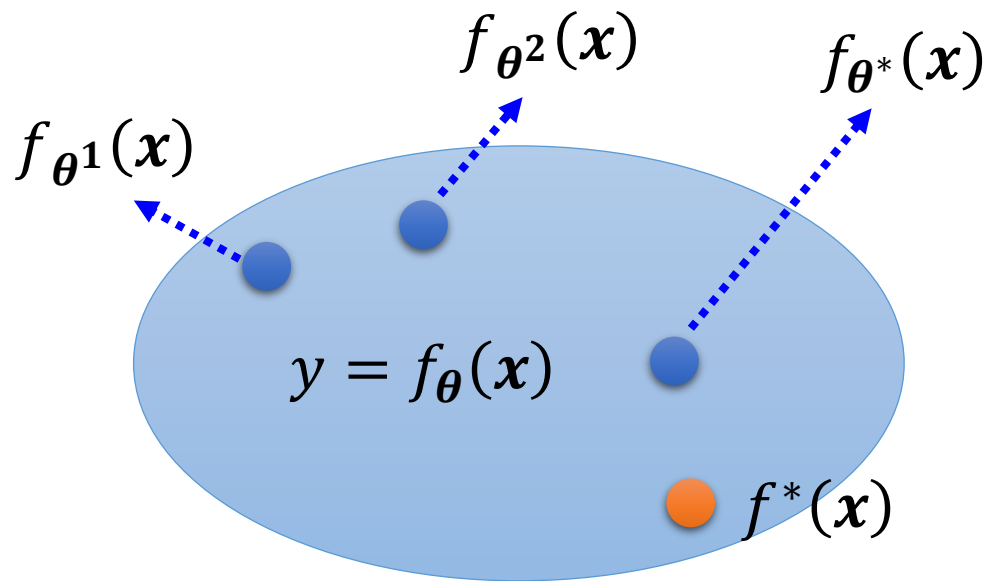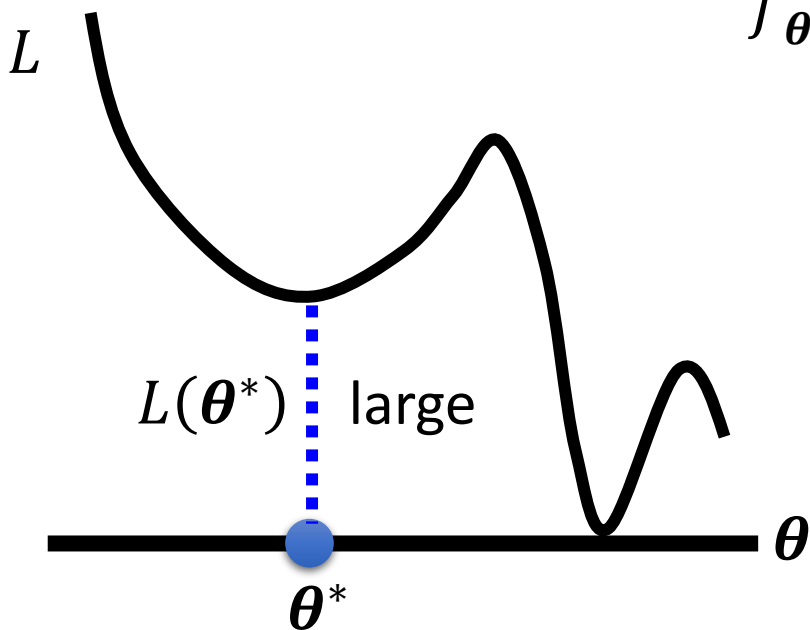
Deep Learning
(more neurons, layers)

$$y = b + \sum_i c_i \; sigmoid\left( b_i + \sum_j w_{ij} x_j \right)$$

**General Guide**

loss on training data

- large → model bias → make your model complex
- large → optimization → Next Lecture
- small → loss on testing data
  - large
    - overfitting → more training data (not in HWs) data augmentation, make your model simpler
    - mismatch → Not in HWs, except HW 11
  - small → 😊

trade-off

Split your training data into training set and validation set for model selection

# Optimization Issue

- Large loss not always imply model bias. There is another possibility ...



$L$

$L(\boldsymbol{\theta}^*)$ large

$\boldsymbol{\theta}^*$

$\boldsymbol{\theta}$

$f_{\boldsymbol{\theta^1}}(\boldsymbol{x})$

$f_{\boldsymbol{\theta^2}}(\boldsymbol{x})$

$f_{\boldsymbol{\theta^*}}(\boldsymbol{x})$

$y = f_{\boldsymbol{\theta}}(\boldsymbol{x})$

$f^*(\boldsymbol{x})$

A needle is in a haystack ...

... Just cannot find it.

## *Model Bias*

find a needle in a haystack …

   … but there is no needle

$f_{\theta^1}(x)$

$f_{\theta^2}(x)$

$f_{\theta^*}(x)$

$f^*(x)$

too small …

small loss

**Which one???**

## *Optimization Issue*

A needle is in a haystack …

   … Just cannot find it.

$f_{\theta^1}(x)$

$f_{\theta^2}(x)$

$f_{\theta^*}(x)$

$y = f_{\theta}(x)$

$f^*(x)$

# Model Bias v.s. Optimization Issue

- Gaining the insights from comparison



**Testing Data**

**Training Data**

# Optimization Issue

- Gaining the insights from comparison

- Start from shallower networks (or other models), which are easier to optimize.

- If deeper networks do not obtain smaller loss on **training data**, then there is optimization issue.
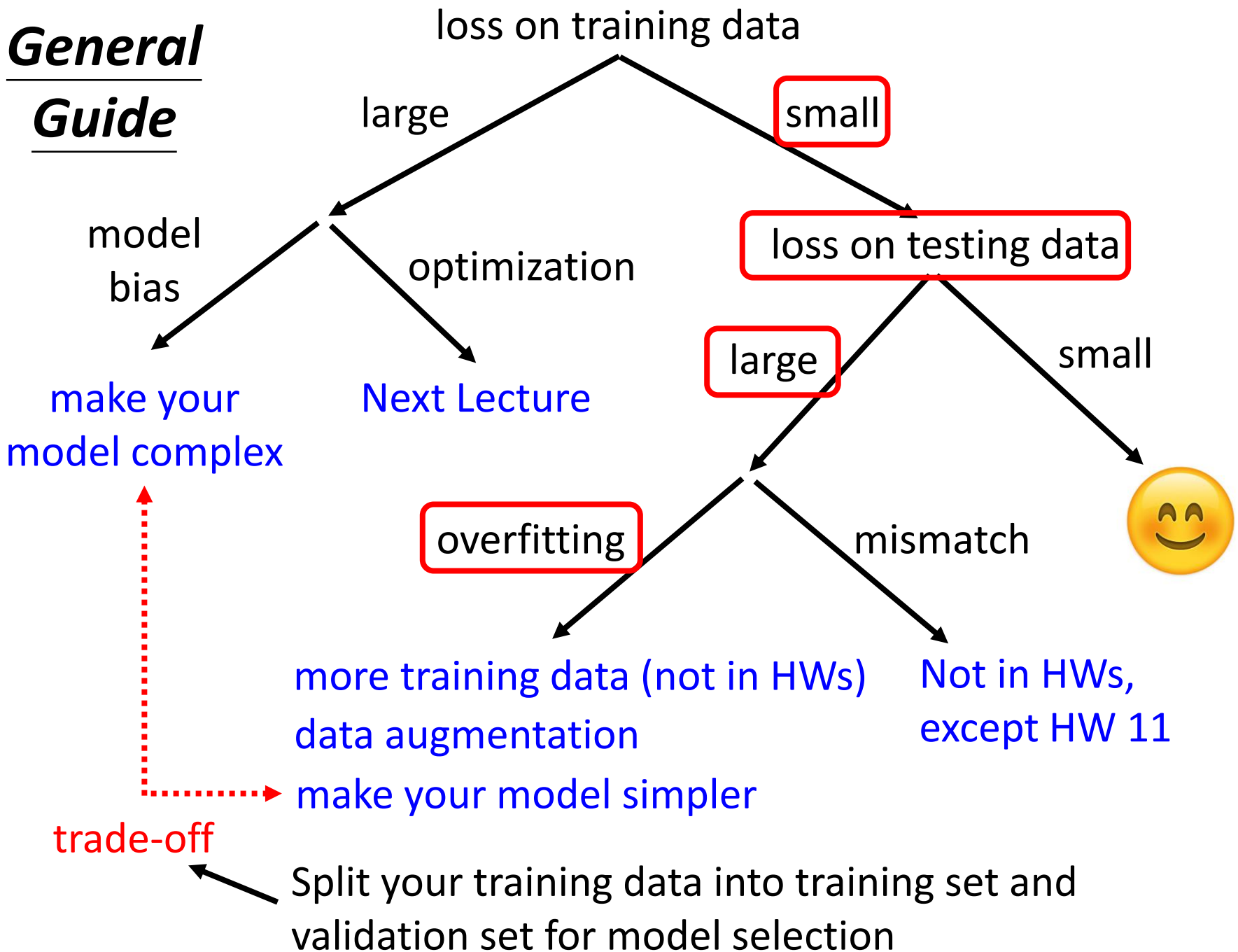
|  | 1 layer | 2 layer | 3 layer | 4 layer | 5 layer |
|---|---|---|---|---|---|
| 2017 – 2020 | 0.28k | 0.18k | 0.14k | 0.10k | 0.34k |

- Solution: More powerful optimization technology (next lecture)

# General Guide

loss on training data

large — model bias → **make your model complex**

large — optimization → **Next Lecture**

small → loss on testing data

loss on testing data — large:
- overfitting → **more training data (not in HWs) data augmentation** / **make your model simpler**
- mismatch → **Not in HWs, except HW 11**

loss on testing data — small → 😊

trade-off

Split your training data into training set and validation set for model selection

# *General Guide*

loss on training data

large → model bias → make your model complex

large → optimization → Next Lecture

small → loss on testing data

loss on testing data → large → overfitting

loss on testing data → small → 😊

overfitting → more training data (not in HWs) data augmentation

overfitting → make your model simpler

mismatch → Not in HWs, except HW 11

trade-off

Split your training data into training set and validation set for model selection

# Overfitting

- Small loss on training data, large loss on testing data. Why?

**_An extreme example_**

Training data: $\left\{ \left( \boldsymbol{x^1}, \hat{y}^1 \right), \left( \boldsymbol{x^2}, \hat{y}^2 \right), \dots, \left( \boldsymbol{x^N}, \hat{y}^N \right) \right\}$

$$f(\boldsymbol{x}) = \begin{cases} \hat{y}^i & \exists \boldsymbol{x^i} = \boldsymbol{x} \\ random & otherwise \end{cases}$$  <span style="color:red">Less than useless …</span>

This function obtains **zero training loss**, but **large testing loss**.

# Overfitting



Flexible model

"freestyle"

Large loss

**· · · ·** Real data distribution (not observable)

🔵 Training data

🟠 Testing data

# Overfitting



Flexible model →
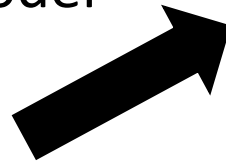
**More training data**
(cannot do it in HWs)

**Data augmentation** (you can do that in HWs)
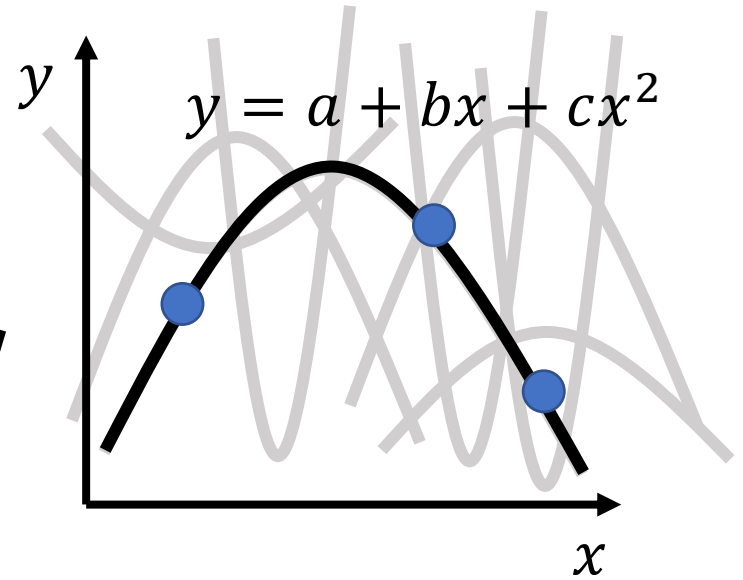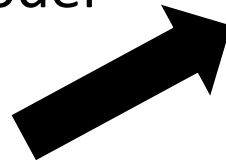
# Overfitting



**constrained** model

$$y = a + bx + cx^2$$

Real data distribution (not observable)

Training data

Testing data

# Overfitting
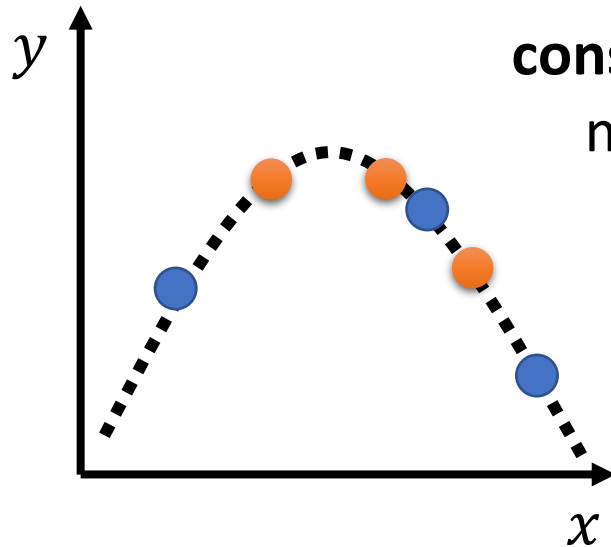


**constrained** model

$$y = a + bx + cx^2$$

- ▪▪▪▪ Real data distribution (not observable)
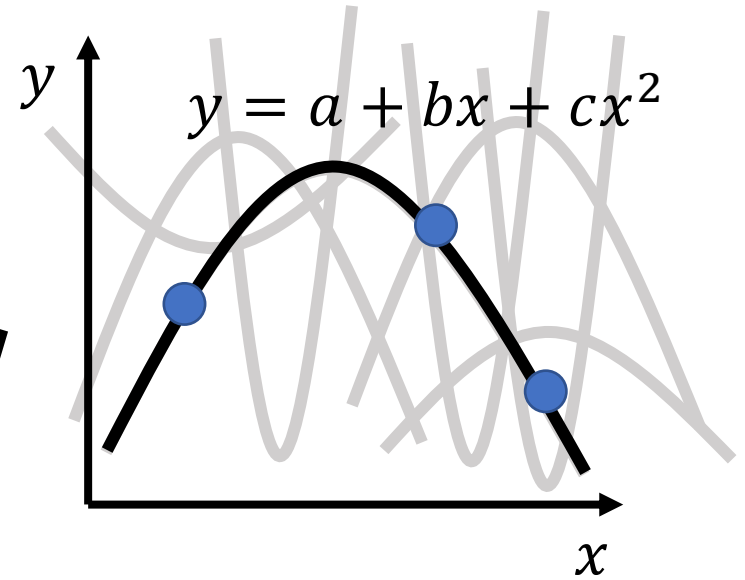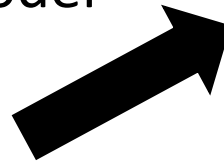- 🔵 Training data
- 🟠 Testing data

# Overfitting



**constrained** model

$$y = a + bx + cx^2$$

- Less parameters, sharing parameters
- Less features
- Early stopping
- Regularization
- Dropout

Fully-connected

CNN

# Overfitting

constrain
too much

$y = a + bx$

Back to model bias …

**⋯⋯** Real data distribution
(not observable)

🔵 Training data

🟠 Testing data

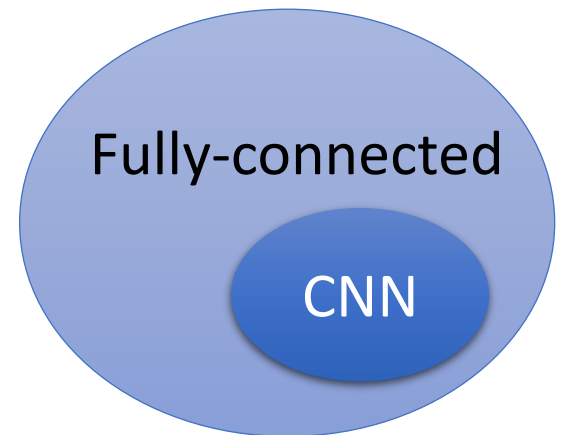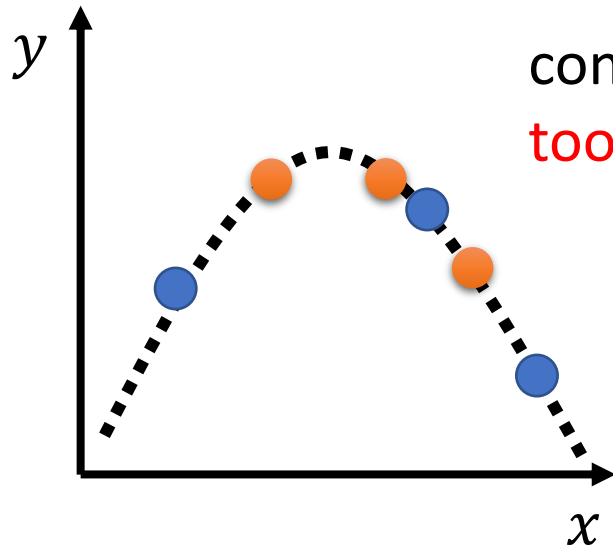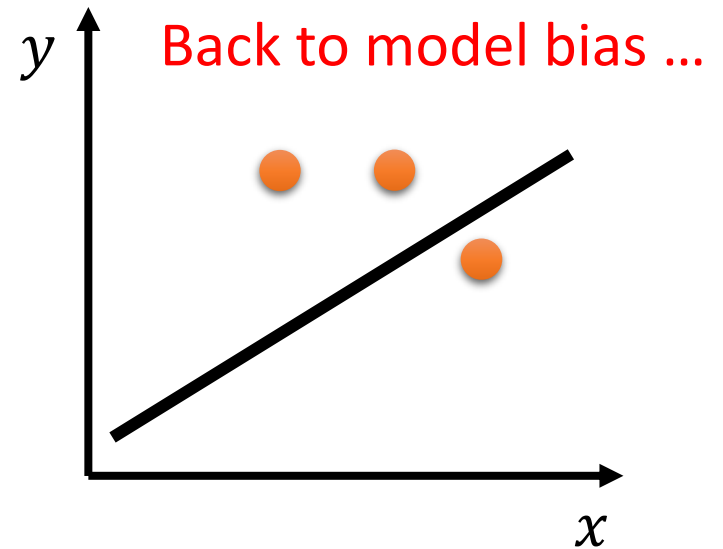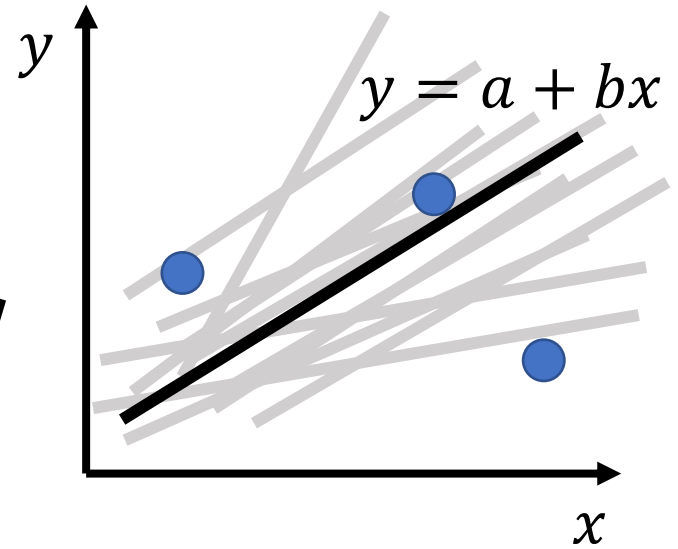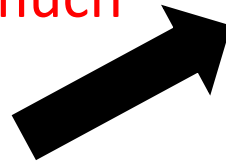# Bias-Complexity Trade-off

## *Homework*

Training Set    public    private
Testing Set    Testing Set

Model 1 ⟶ mse = 0.9

Model 2 ⟶ mse = 0.7

Model 3 ⟶ mse = 0.5 ⟶ mse > 0.5

Pick this one!    May be poor …

## *The extreme example again*

$$f_k(x) = \begin{cases} \hat{y}^i & \exists x^i = x \\ random & otherwise \end{cases}$$

$k$: 1 - 10000000000000000000

It is possible that $f_{56789}(x)$ **happens** to get good performance on public testing set.

So you select $f_{56789}(x)$ ……    Random on private testing set

# *Homework*

Training Set

Testing Set

Testing Set

Why?

Model 1 $\longrightarrow$ mse = 0.9

Model 2 $\longrightarrow$ mse = 0.7

Model 3 $\longrightarrow$ mse = 0.5 $\longrightarrow$ mse > 0.5

Pick this one!

May be poor …

What will happen?

http://www.chioka.in/how-to-select-your-final-models-in-a-kaggle-competitio/

This explains why machine usually beats human on benchmark corpora. ☺



TOP 10 IN PUBLIC LEADERBOARD

RANKED 3XX IN PRIVATE LEADERBOARD

imgflip.com

# Cross Validation

How to split?

| Training Set |
| :---: |

public

| Testing Set |
| :---: |

private

| Testing Set |
| :---: |

Training Set → Training Set / Validation set

Using the results of public testing data to select your model

You are making public set better than private set.

Model 1 ──→ mse = 0.9

Model 2 ──→ mse = 0.7

Model 3 ──→ mse = 0.5 ──→ mse > 0.5 ──→ mse > 0.5

Not recommend

# General Guide

loss on training data

— large → model bias → make your model complex

— large → optimization → Next Lecture

— small → loss on testing data

- large → overfitting → more training data (not in HWs) data augmentation
- large → overfitting → make your model simpler
- large → mismatch → Not in HWs, except HW 11
- small → 😊

trade-off

Split your training data into training set and validation set for model selection

# Let's predict no. of views of 2/26!

| | 1 layer | 2 layer | 3 layer | 4 layer |
|---|---|---|---|---|
| 2017 – 2020 | | | | 0k |
| 2021 | | | | 4k |

Red: real, Blue: predicted

2/26

e = 2.58k

loss on training data

large / small

model bias

optimization

loss on testing data

make your model complex

Next Lecture

large / small

overfitting

mismatch

😊

more training data (not in HWs) data augmentation

make your model simpler

Not in HWs, except HW 11

trade-off

Split your training data into training set and validation set for model selection

# Mismatch

- Your training and testing data have different distributions.  Be aware of how data is generated.

**_Most HWs do not have this problem, except HW11_**

**_Training Data_**



horse    bed    clock    apple    cat    plane    television    dog    dolphin    spider

Simply increasing the training data will not help.

**_Testing Data_**

loss on training data

large — small

model bias

optimization

loss on testing data

make your model complex

Next Lecture

large — small

overfitting

mismatch

😊

more training data (not in HWs) data augmentation

Not in HWs, except HW 11

make your model simpler

trade-off

Split your training data into training set and validation set for model selection