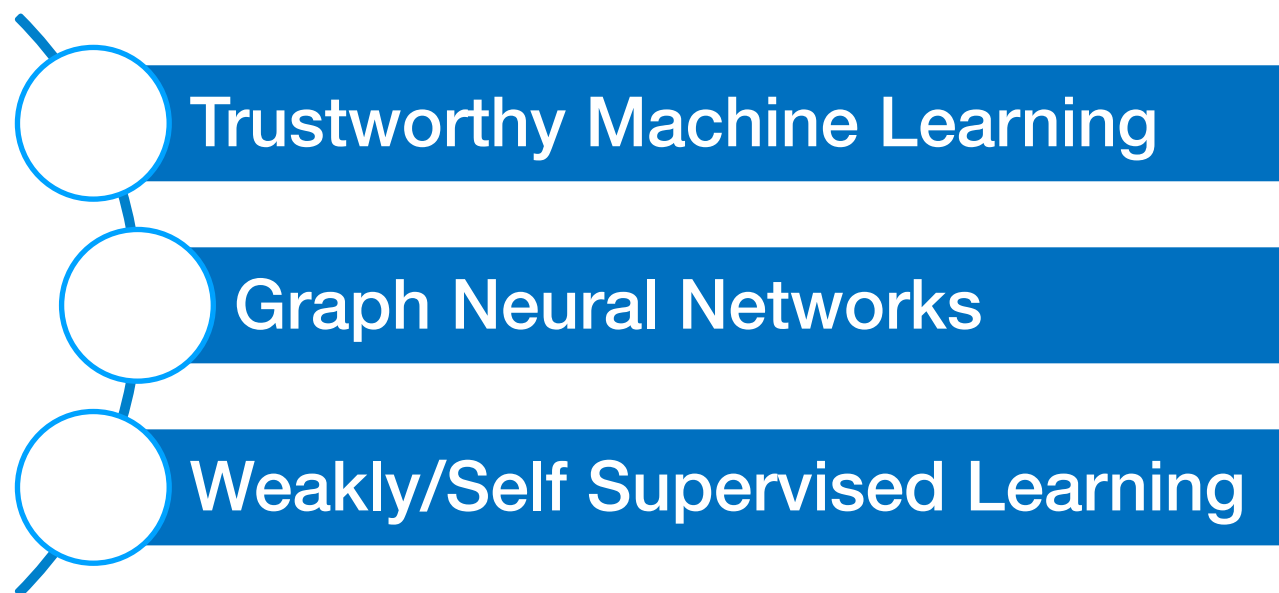# 机器学习介绍
# Introduction of Machine Learning

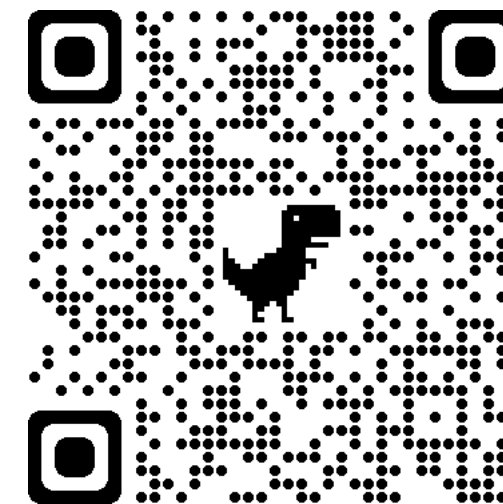**AIPKU**

主讲人：王奕森

- **Assistant Professor, Ph.D. Advisor**
- **Research Interests: Machine Learning**

Trustworthy Machine Learning

Graph Neural Networks

Weakly/Self Supervised Learning

**Homepage**

**https://yisenwang.github.io/**

中国计算机学会推荐国际学术会议

（人工智能）

- **Welcome interns!**

WE WANT YOU!

一、A 类

| 序号 | 会议简称 | 会议全称 | 出版社 | 网址 |
|---|---|---|---|---|
| 1 | AAAI | AAAI Conference on Artificial Intelligence | AAAI | http://dblp.uni-trier.de/db/conf/aaai/ |
| 2 | NeurIPS | Annual Conference on Neural Information Processing Systems | MIT Press | http://dblp.uni-trier.de/db/conf/nips/ |
| 3 | ACL | Annual Meeting of the Association for Computational Linguistics | ACL | http://dblp.uni-trier.de/db/conf/acl/ |
| 4 | CVPR | IEEE Conference on Computer Vision and Pattern Recognition | IEEE | http://dblp.uni-trier.de/db/conf/cvpr/ |
| 5 | ICCV | International Conference on Computer Vision | IEEE | http://dblp.uni-trier.de/db/conf/iccv/ |
| 6 | ICML | International Conference on Machine Learning | ACM | http://dblp.uni-trier.de/db/conf/icml/ |
| 7 | IJCAI | International Joint Conference on Artificial Intelligence | Morgan Kaufmann | http://dblp.uni-trier.de/db/conf/ijcai/ |

**ICLR, International Conference on Learning Representations**

- Machine Learning is for designing algorithms that can learn and build computable models from training data and then perform desired tasks on new data.

A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measure by P, **improves** with experience E.

So ML is not rule based and in principle it can enable machines to **evolve**.

Tom Mitchell
(Member of NAE & AAAS,
Prof. @ CMU)

## Machine Learning ≈ Looking for Function

- Speech Recognition

$$f\left( \text{[waveform]} \right) = \text{"How are you"}$$
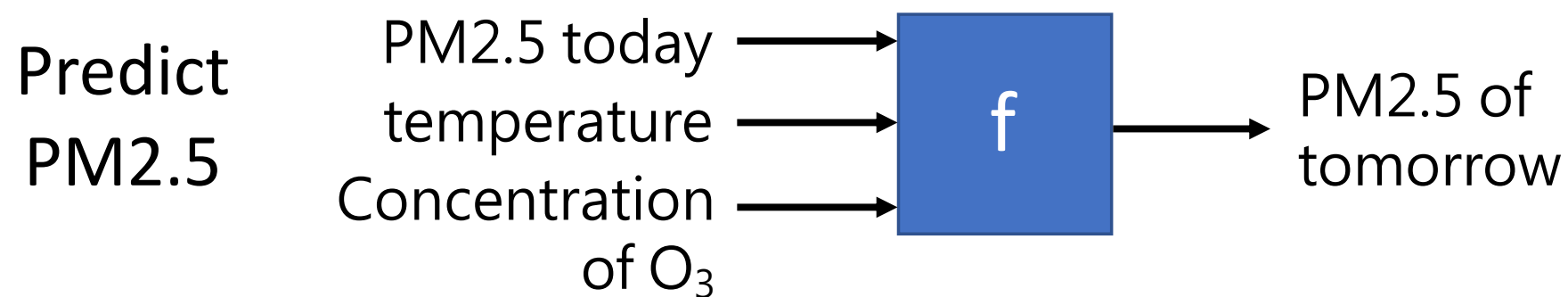
- Image Recognition

$$f\left( \text{[cat image]} \right) = \text{"Cat"}$$

- Playing Go

$$f\left( \text{[Go board]} \right) = \text{"5-5"} \text{ (next move)}$$

**Regression:** The function outputs a scalar.

Predict PM2.5

PM2.5 today $\longrightarrow$

temperature $\longrightarrow$ **f** $\longrightarrow$ PM2.5 of tomorrow

Concentration of $O_3$ $\longrightarrow$
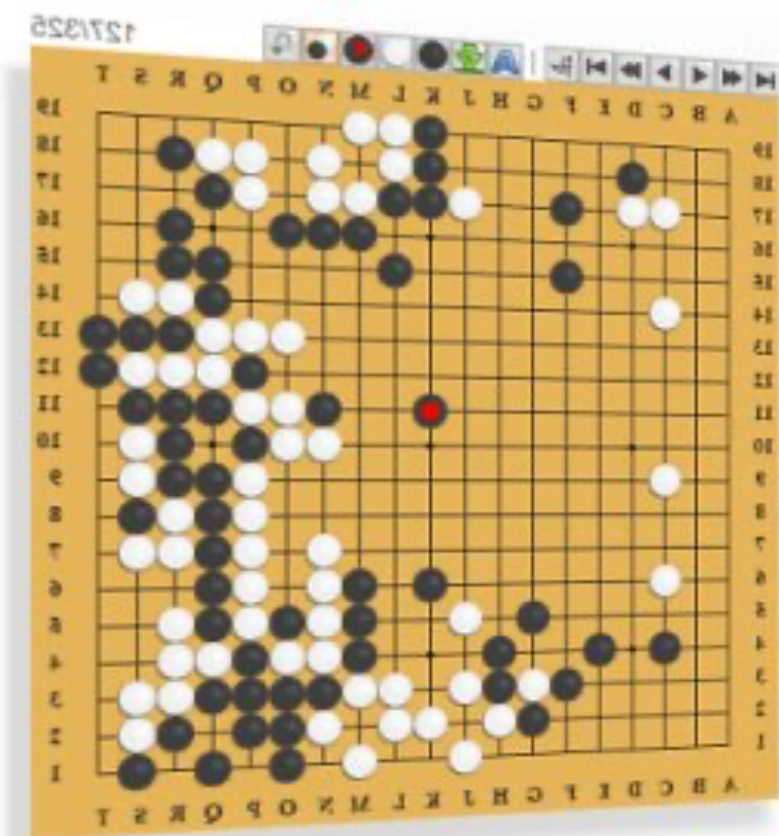
**Classification:** Given options (**classes**), the function outputs the correct one.

Spam filtering

EMAIL $\longrightarrow$ **f** $\longrightarrow$ Yes/No

**Classification:** Given options (**classes**), the function outputs the correct one.

Each position is a class

(19 x 19 classes)

Function

a position on the board

**Next move**

*Playing GO*

**Taking examples (some slides) from Prof. Hungyi Lee@NTU**



$$y = f(\quad\quad)$$

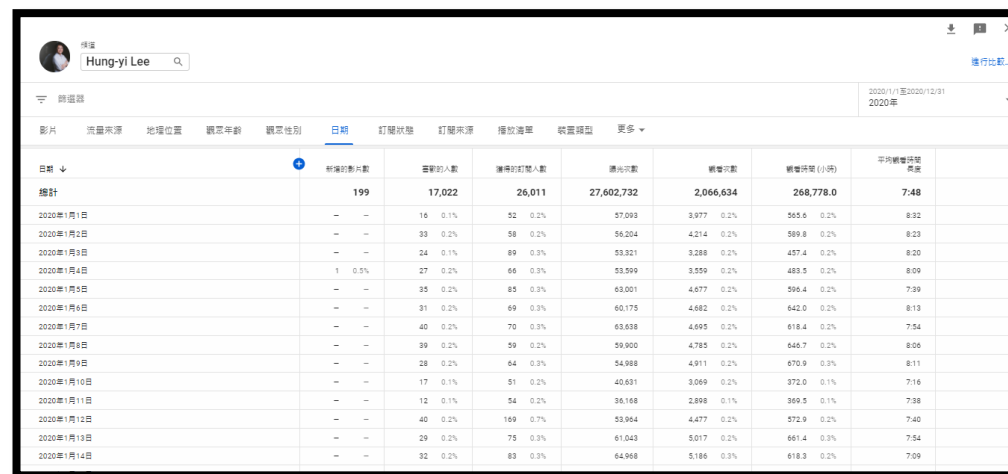no. of views on 2/26

$$y = f\left( \quad \right)$$

**Model** $\quad y = b + wx_1 \quad$ based on domain knowledge

**feature**

$y$: no. of views on 2/26, $x_1$: no. of views on 2/25

$w$ and $b$ are unknown parameters (learned from data)

**weight**    **bias**

# Defining Loss from Training Data

> Loss is a function of parameters $L(b, w)$
> Loss: how good a set of values is.

$L(0.5k, 1)$    $y = b + wx_1$ $\longrightarrow$ $y = 0.5k + 1x_1$    How good it is?

Data from 2017/01/01 − 2020/12/31

| 2017/01/01 | 01/02 | 01/03 | ...... | 2020/12/30 | 12/31 |

4.8k    4.9k    7.5k    3.4k    9.8k

$0.5k + 1x_1 = y$    5.3k

$e_1 = |y - \hat{y}| = 0.4k$
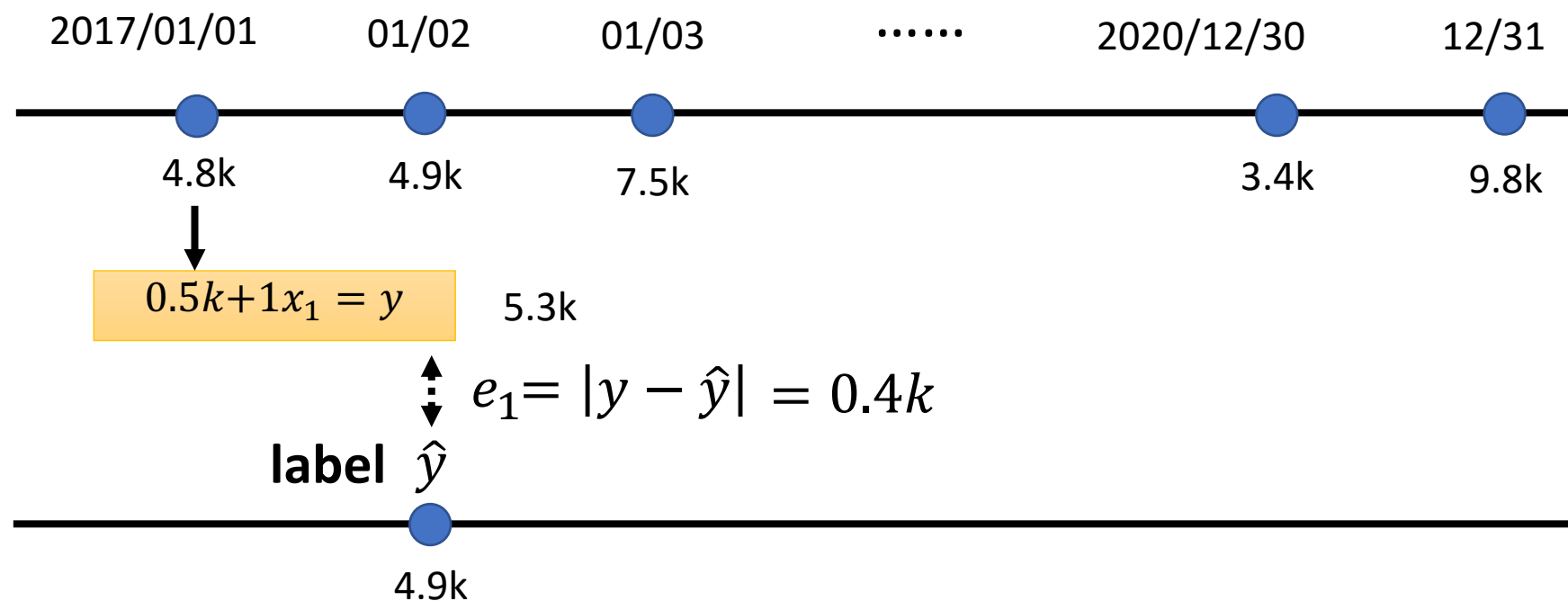
**label** $\hat{y}$

4.9k

Defining **Loss** from Training Data

➤ Loss is a function of parameters $L(b, w)$

➤ Loss: how good a set of values is.

$L(0.5k, 1)$ $\qquad y = b + wx_1 \longrightarrow y = 0.5k + 1x_1$ How good it is?

Data from 2017/01/01 − 2020/12/31

2017/01/01 01/02 01/03 ...... 2020/12/30 12/31

4.8k 4.9k 7.5k 3.4k 9.8k

$0.5k+1x_1 = y$ 5.4k $\qquad\qquad 0.5k+1x_1 = y$

$e_2 = |y - \hat{y}| = 2.1k$ $\qquad e_N$

$\hat{y}$ $\qquad\qquad\qquad\qquad\qquad \hat{y}$

4.9k 7.5k 9.8k
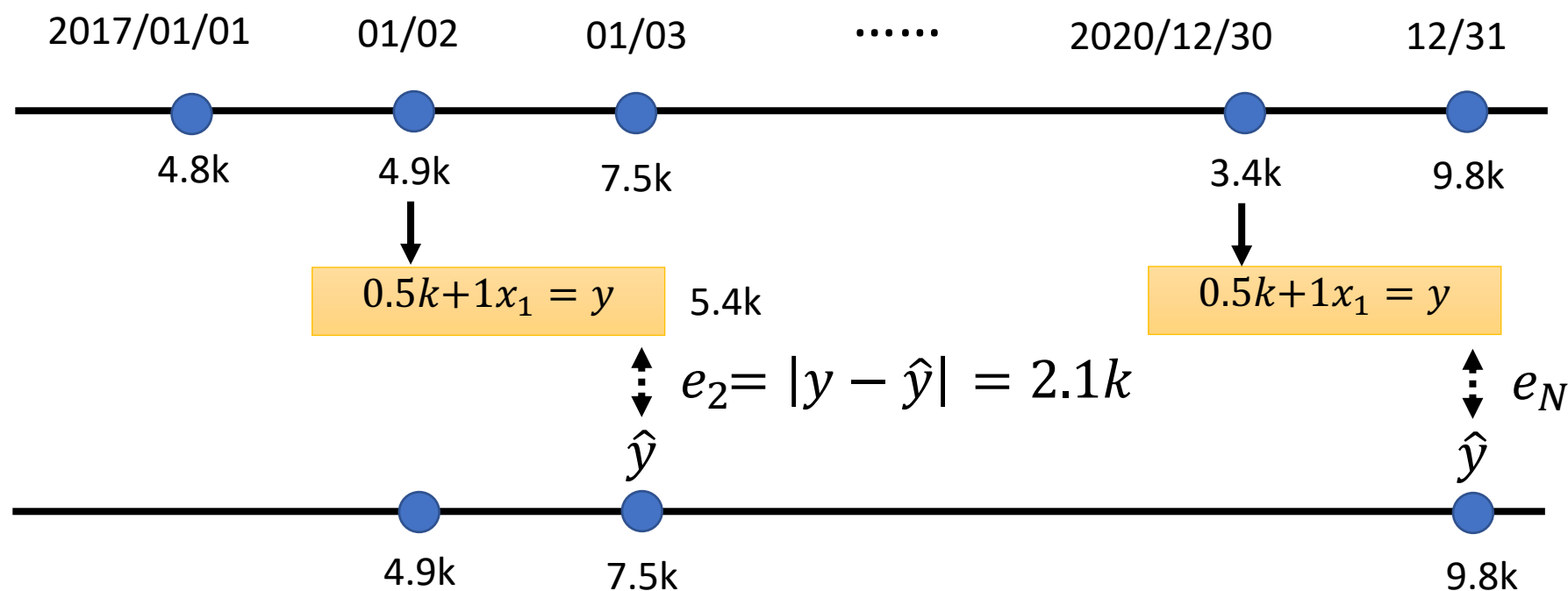
## Defining Loss from Training Data

➤ Loss is a function of parameters $L(b, w)$

➤ Loss: how good a set of values is.

4.8k      4.9k

$b + wx_1 = y$
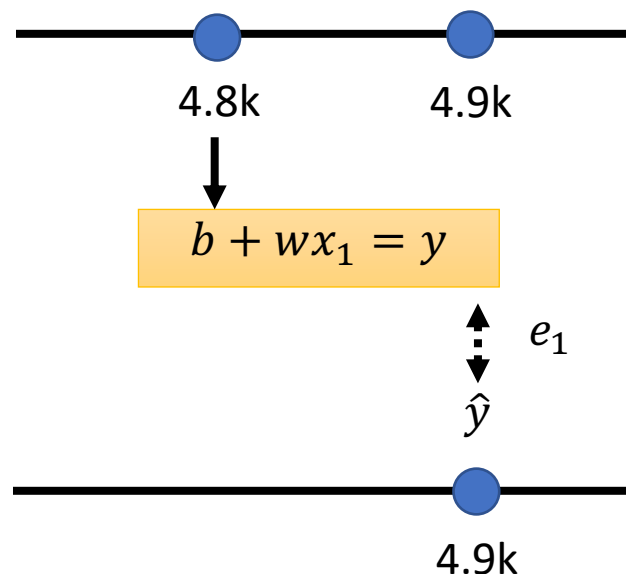
$e_1$

$\hat{y}$

4.9k

Loss:   $L = \dfrac{1}{N} \sum_n e_n$

$e = |y - \hat{y}|$     $L$ is mean absolute error (MAE)

$e = (y - \hat{y})^2$     $L$ is mean square error (MSE)

If $y$ and $\hat{y}$ are both probability distributions ➡ Cross-entropy

$e = y \log \hat{y}$

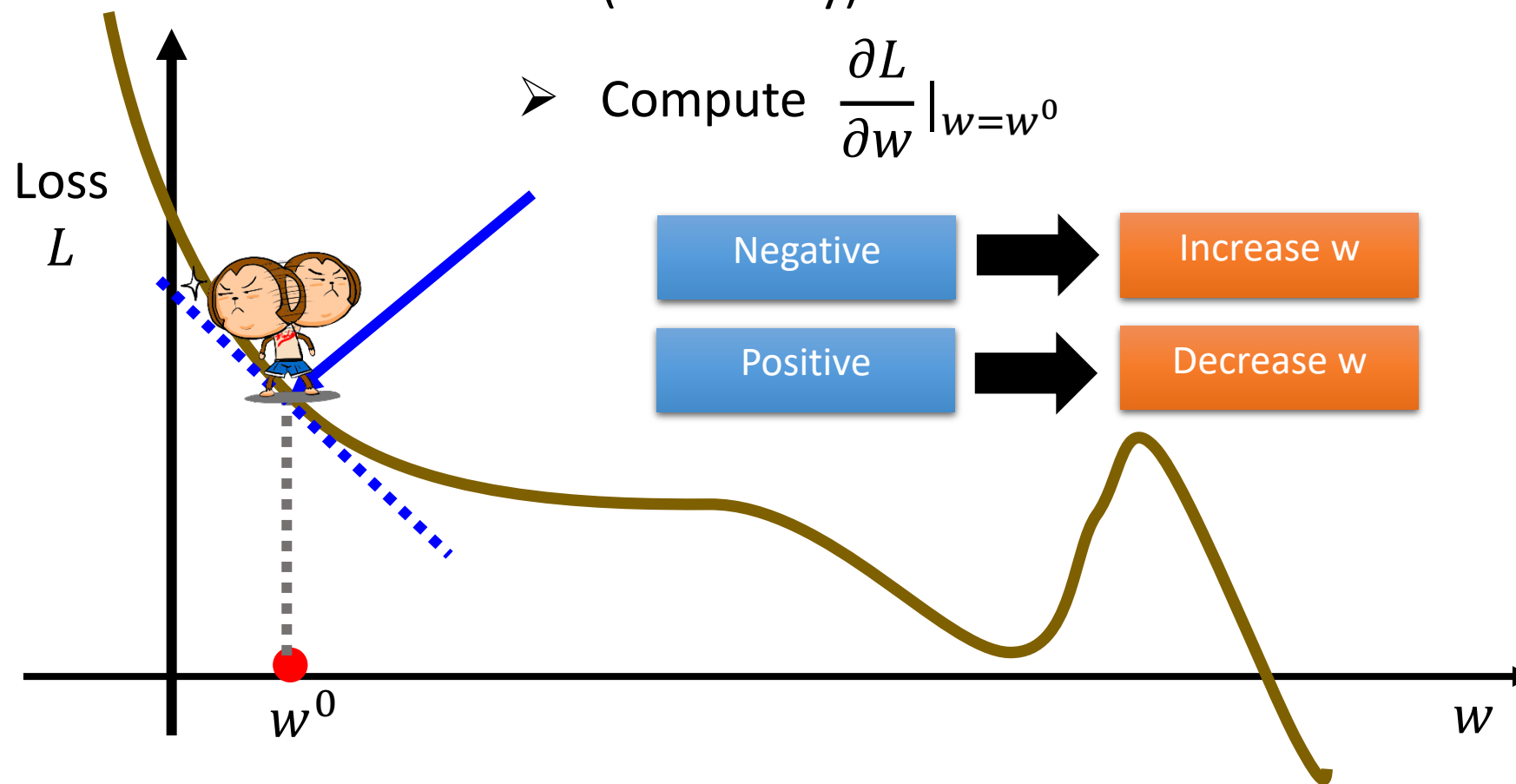$$w^* \; \blacksquare = arg\min_{w} L \; \blacksquare$$

### *Gradient Descent*

➢ (Randomly) Pick an initial value $w^0$

➢ Compute $\dfrac{\partial L}{\partial w}\Big|_{w=w^0}$

| Negative | ➡ | Increase w |
| Positive | ➡ | Decrease w |

Loss
$L$

$w^0$

$w$

$$w^* \blacksquare = arg \min_{w \blacksquare} L$$

### *Gradient Descent*

➢ (Randomly) Pick an initial value $w^0$

➢ Compute $\frac{\partial L}{\partial w}|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w}|_{w=w^0}$$

$\eta \frac{\partial L}{\partial w}|_{w=w^0}$    $\eta$: learning rate
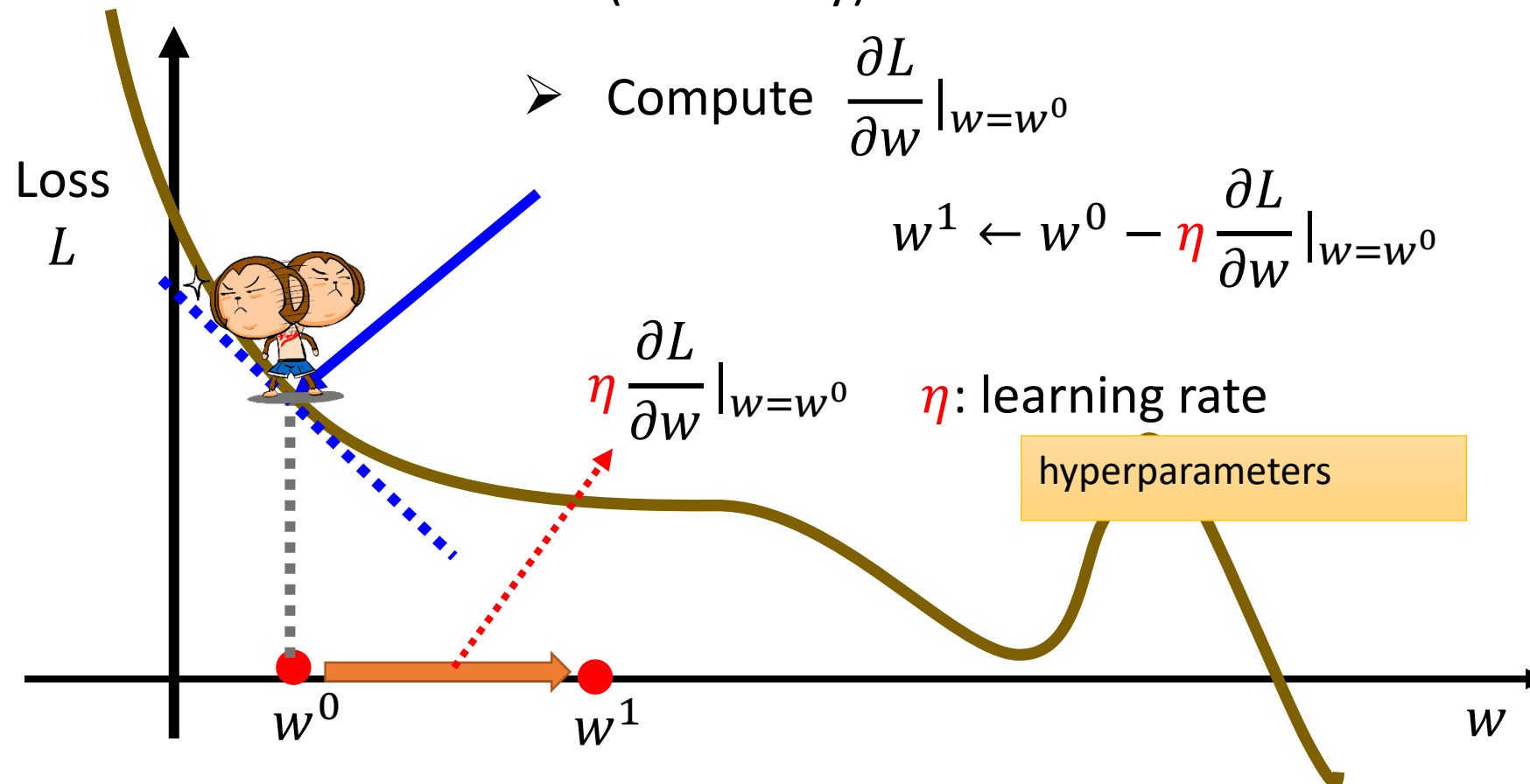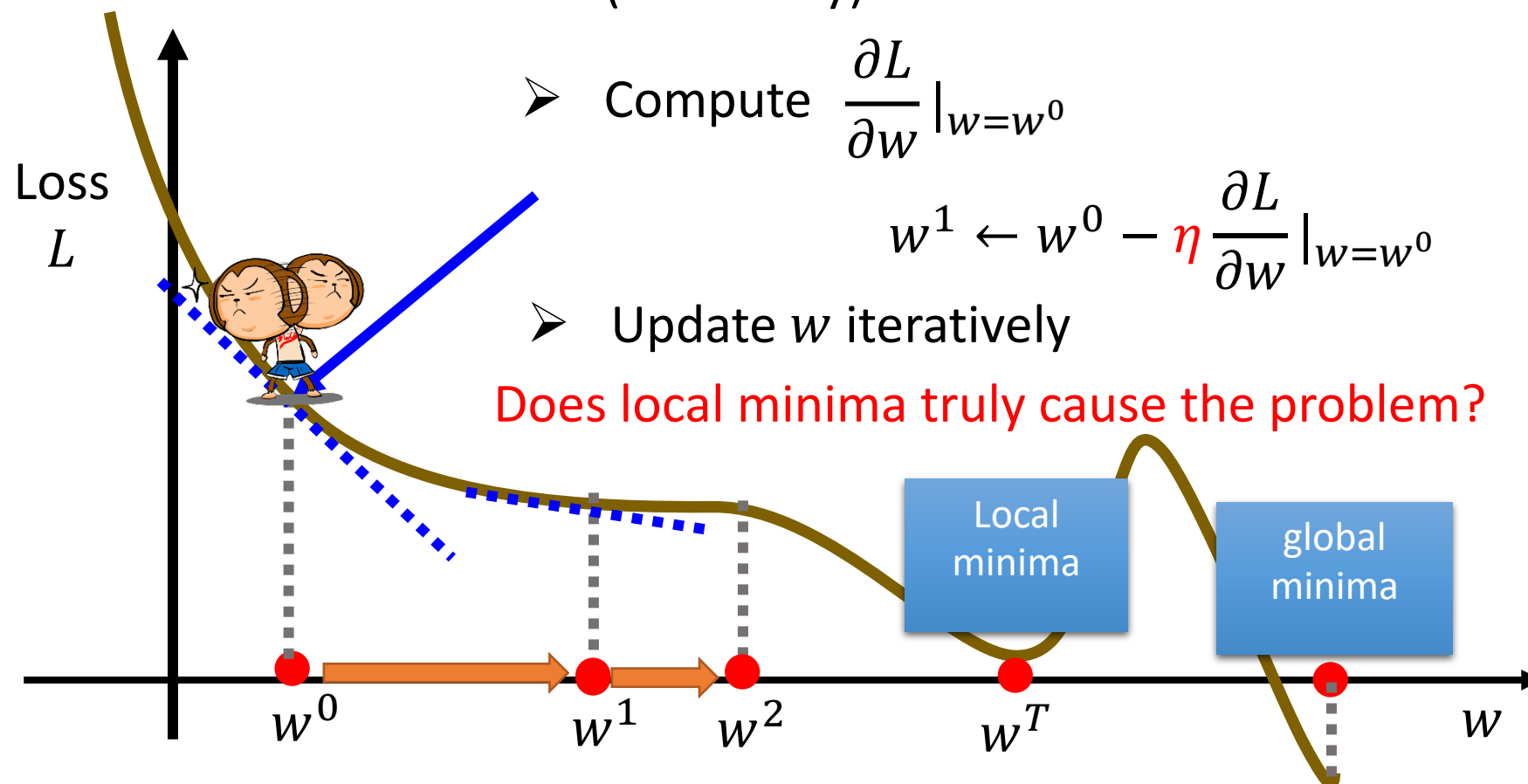
hyperparameters

Loss $L$

$w^0$    $w^1$    $w$

$$w^* \quad = arg \min_{w} L$$

### *Gradient Descent*

➢ (Randomly) Pick an initial value $w^0$

➢ Compute $\dfrac{\partial L}{\partial w}|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \dfrac{\partial L}{\partial w}|_{w=w^0}$$

➢ Update $w$ iteratively

Does local minima truly cause the problem?

Loss
$L$

Local minima

global minima

$w^0$    $w^1$   $w^2$    $w^T$     $w$

$$w^*, b^* = arg \min_{w,b} L$$

➢ (Randomly) Pick initial values $w^0$, $b^0$

➢ Compute

$$\frac{\partial L}{\partial w}\big|_{w=w^0, b=b^0} \qquad w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w}\big|_{w=w^0, b=b^0}$$

$$\frac{\partial L}{\partial b}\big|_{w=w^0, b=b^0} \qquad b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b}\big|_{w=w^0, b=b^0}$$

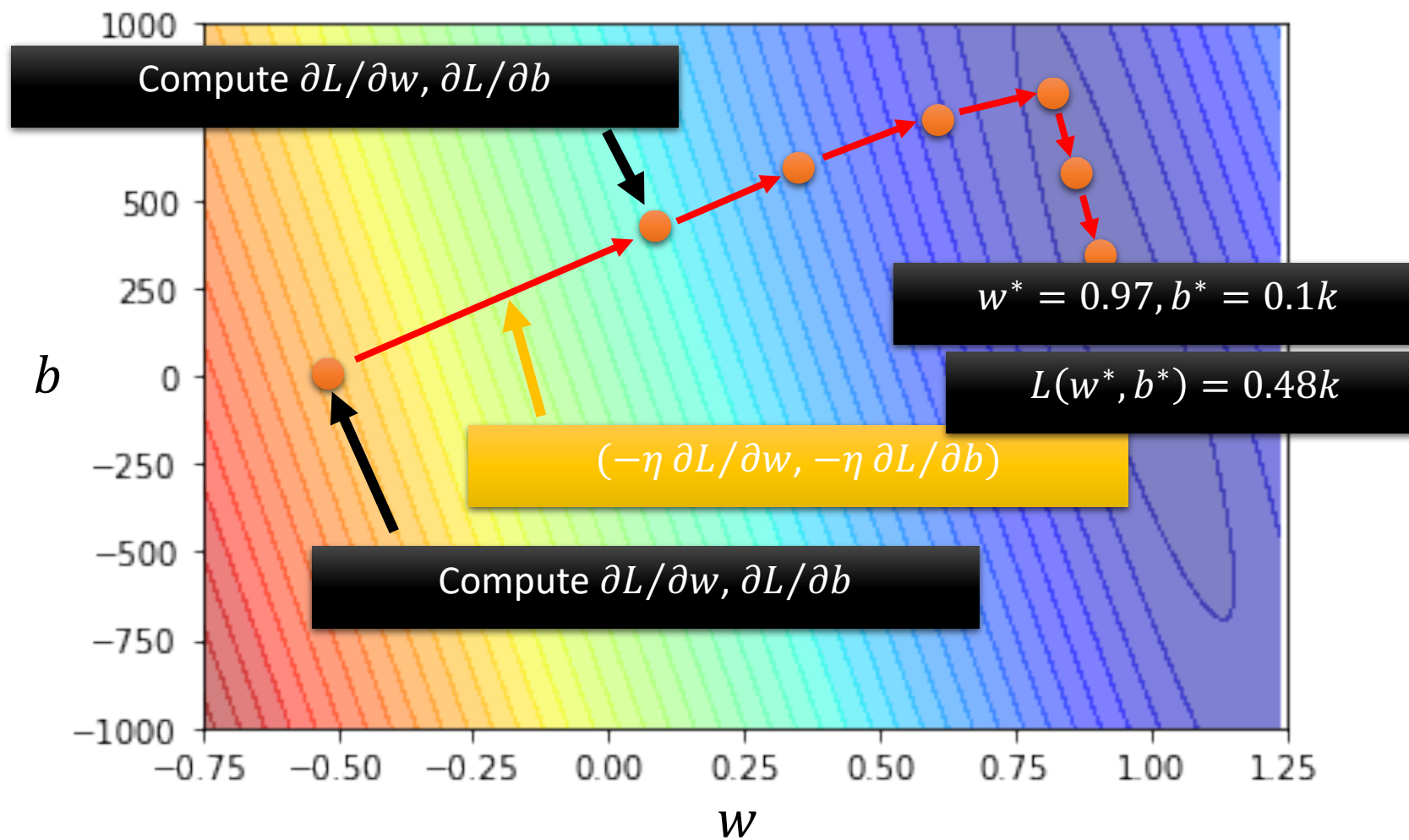Can be done in one line in most deep learning frameworks

➢ Update $w$ and $b$ interatively

**Model** $y = b + wx_1$

$$w^*, b^* = arg \min_{w,b} L$$



Compute $\partial L/\partial w, \partial L/\partial b$

$w^* = 0.97, b^* = 0.1k$

$L(w^*, b^*) = 0.48k$

$(-\eta \ \partial L/\partial w, -\eta \ \partial L/\partial b)$

Compute $\partial L/\partial w, \partial L/\partial b$

$$y = b + wx_1$$

| Step 1: function with unknown | → | Step 2: define loss from training data | → | Step 3: optimization |

$$w^* = 0.97, b^* = 0.1k$$
$$L(w^*, b^*) = 0.48k$$

$$w^* = 0.97, b^* = 0.1k$$
$$L(w^*, b^*) = 0.48k$$

$$y = b + wx_1$$

| Step 1: function with unknown | → | Step 2: define loss from training data | → | Step 3: optimization |
| --- | --- | --- | --- | --- |

***Training***

$y = 0.1k + 0.97x_1$ achieves the smallest loss $L = 0.48k$ on data of $2017 - 2020$ (**training data**)

How about data of 2021 (**unseen during training**)?

$$L' = 0.58k$$

$$y = b + wx_1$$

2017 - 2020 | 2021
$L = 0.48k$ | $L' = 0.58k$

$$y = b + \sum_{j=1}^{7} w_j x_j$$

2017 - 2020 | 2021
$L = 0.38k$ | $L' = 0.49k$

| $b$ | $w_1^*$ | $w_2^*$ | $w_3^*$ | $w_4^*$ | $w_5^*$ | $w_6^*$ | $w_7^*$ |
|---|---|---|---|---|---|---|---|
| 0.05k | 0.79 | -0.31 | 0.12 | -0.01 | -0.10 | 0.30 | 0.18 |

$$y = b + \sum_{j=1}^{28} w_j x_j$$

2017 - 2020 | 2021
$L = 0.33k$ | $L' = 0.46k$

$$y = b + \sum_{j=1}^{56} w_j x_j$$
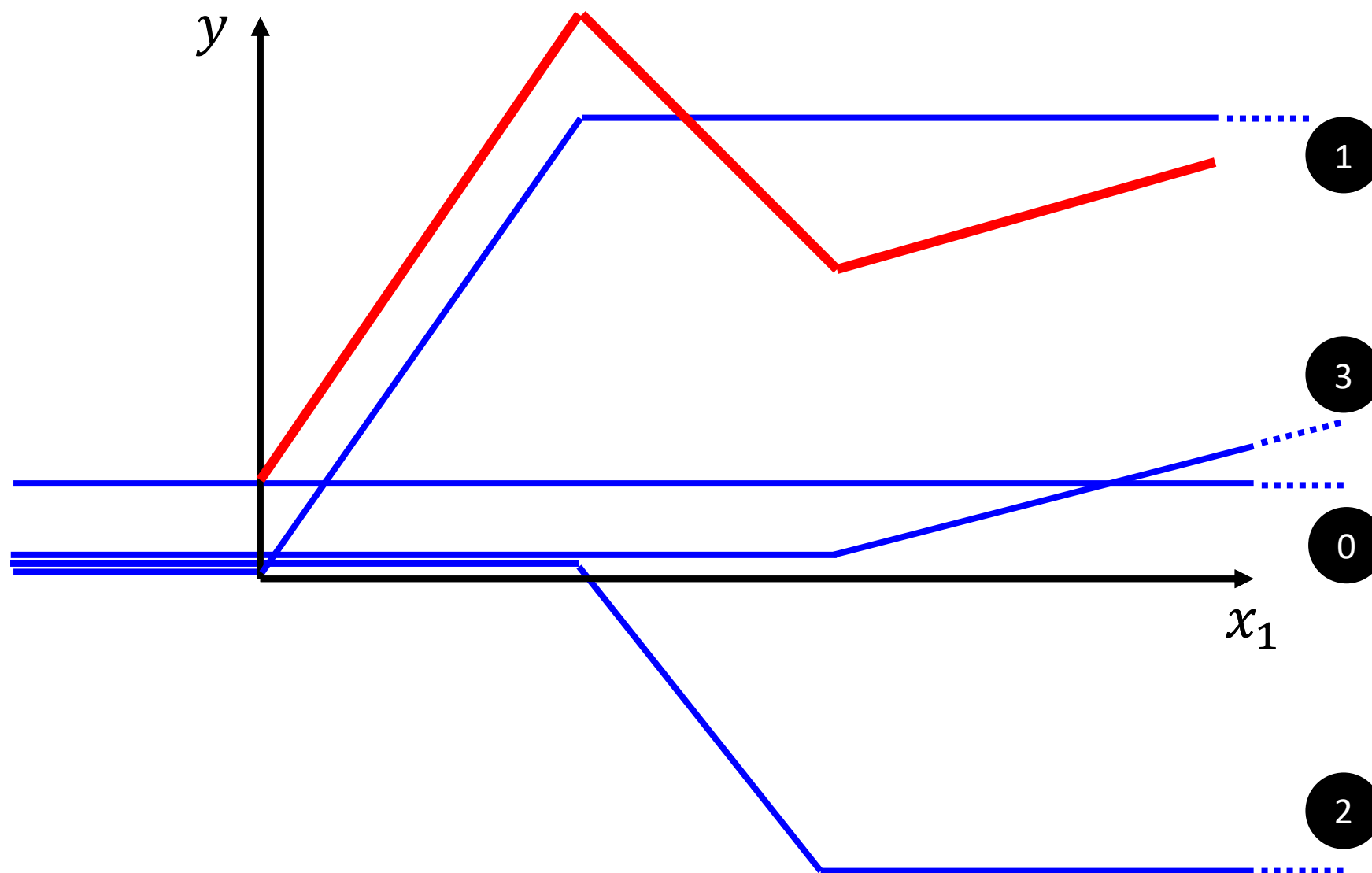
2017 - 2020 | 2021
$L = 0.32k$ | $L' = 0.46k$

***Linear models***
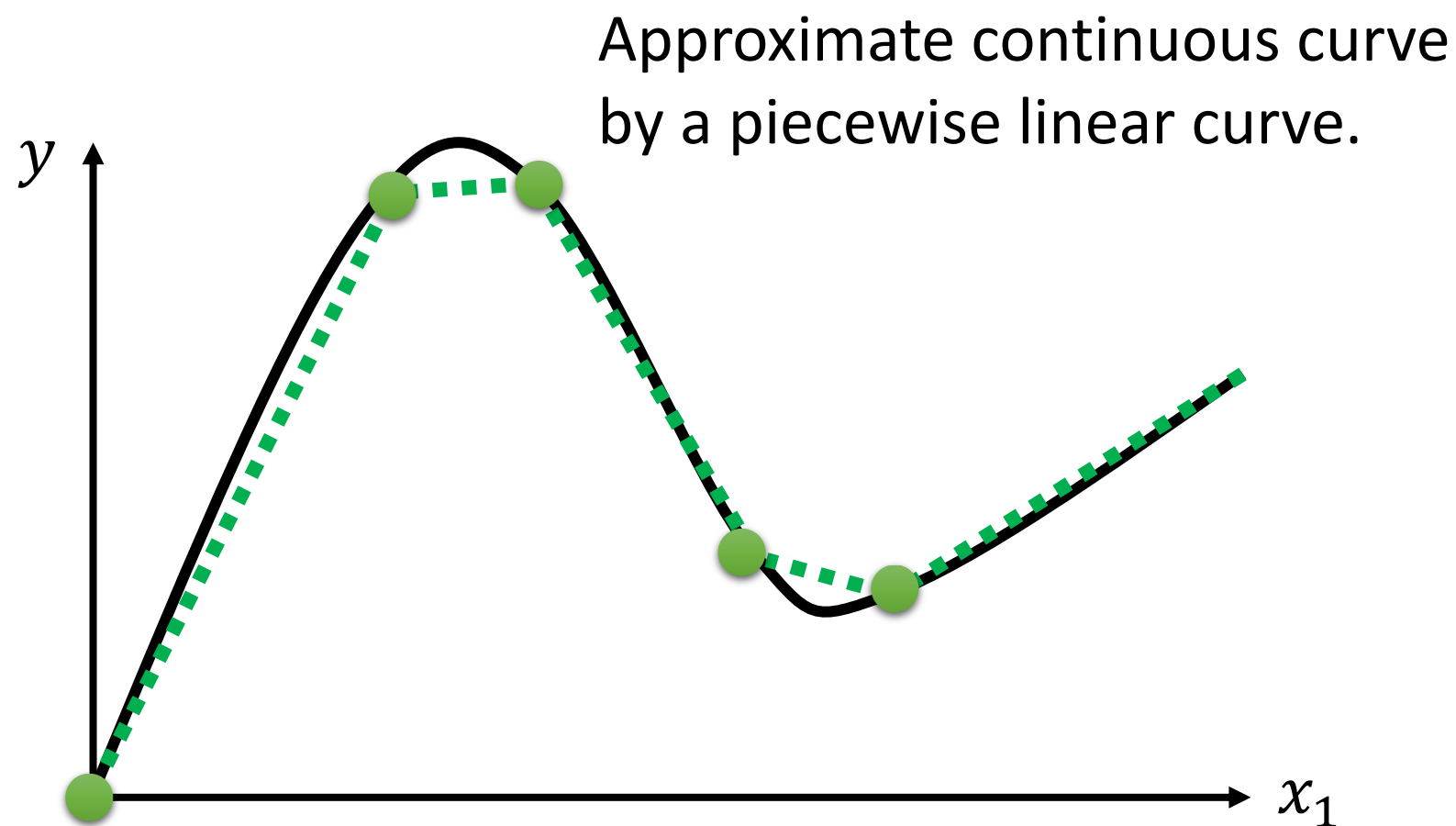
Linear models have severe limitation. **_Model Bias_**

We need a more flexible model!

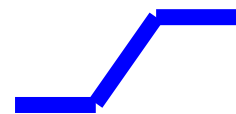Approximate continuous curve by a piecewise linear curve.

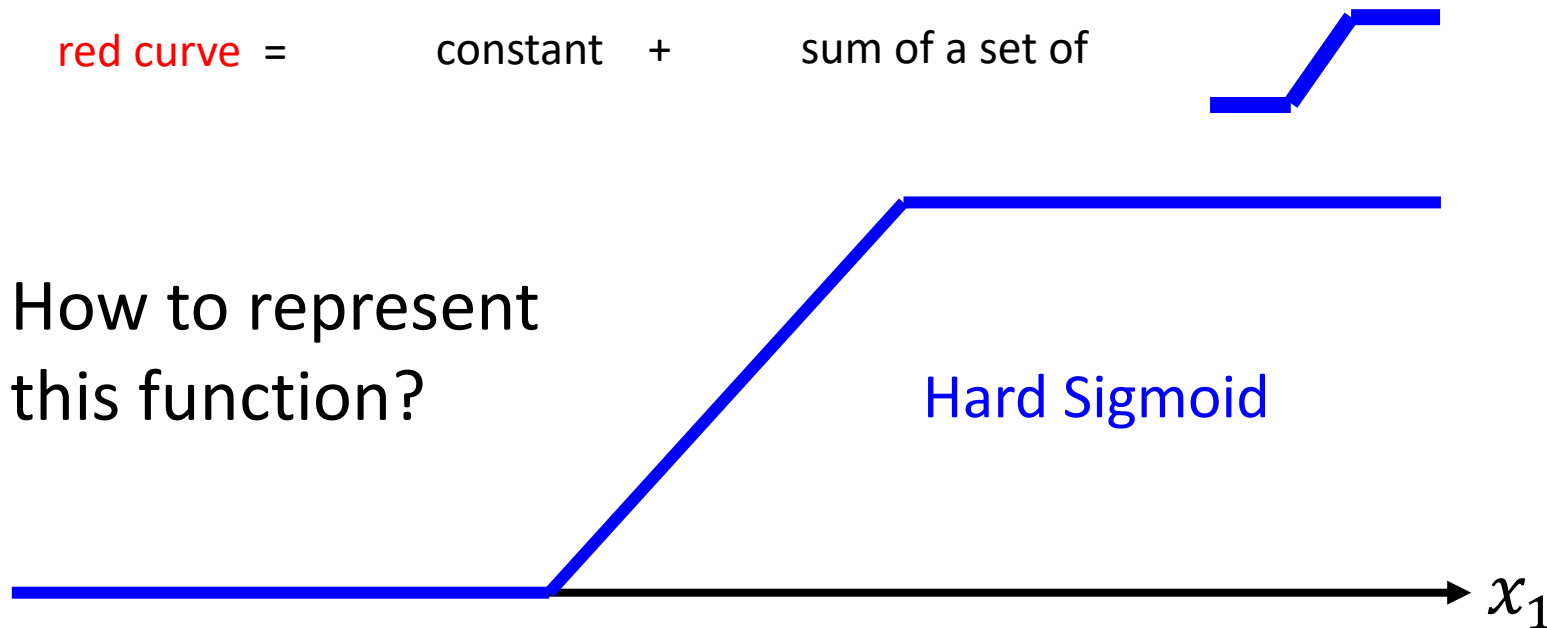To have good approximation, we need sufficient pieces.

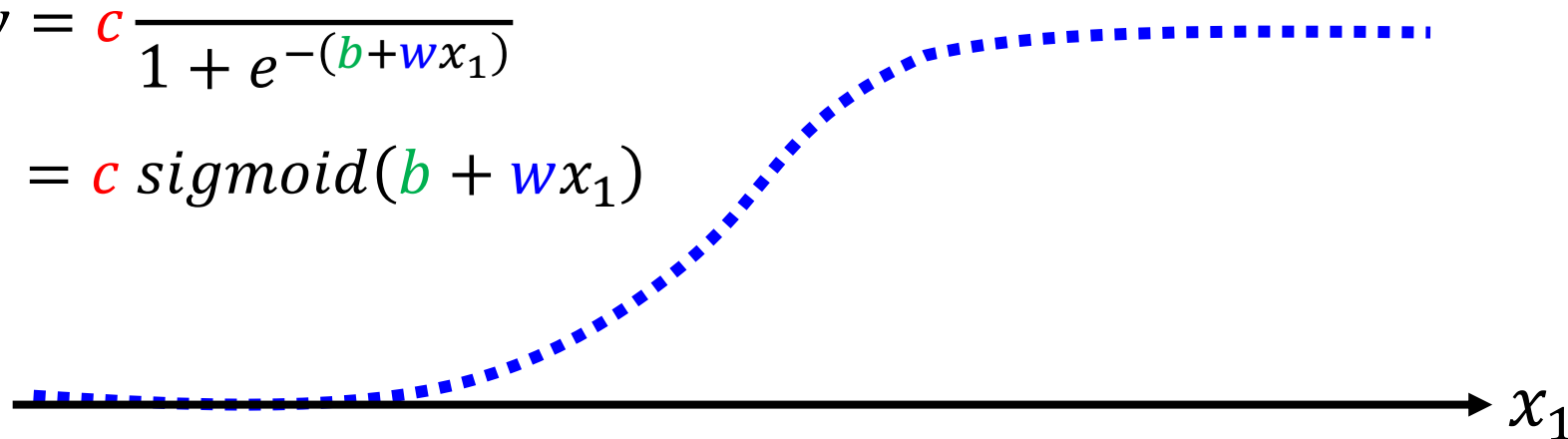red curve =     constant +     sum of a set of
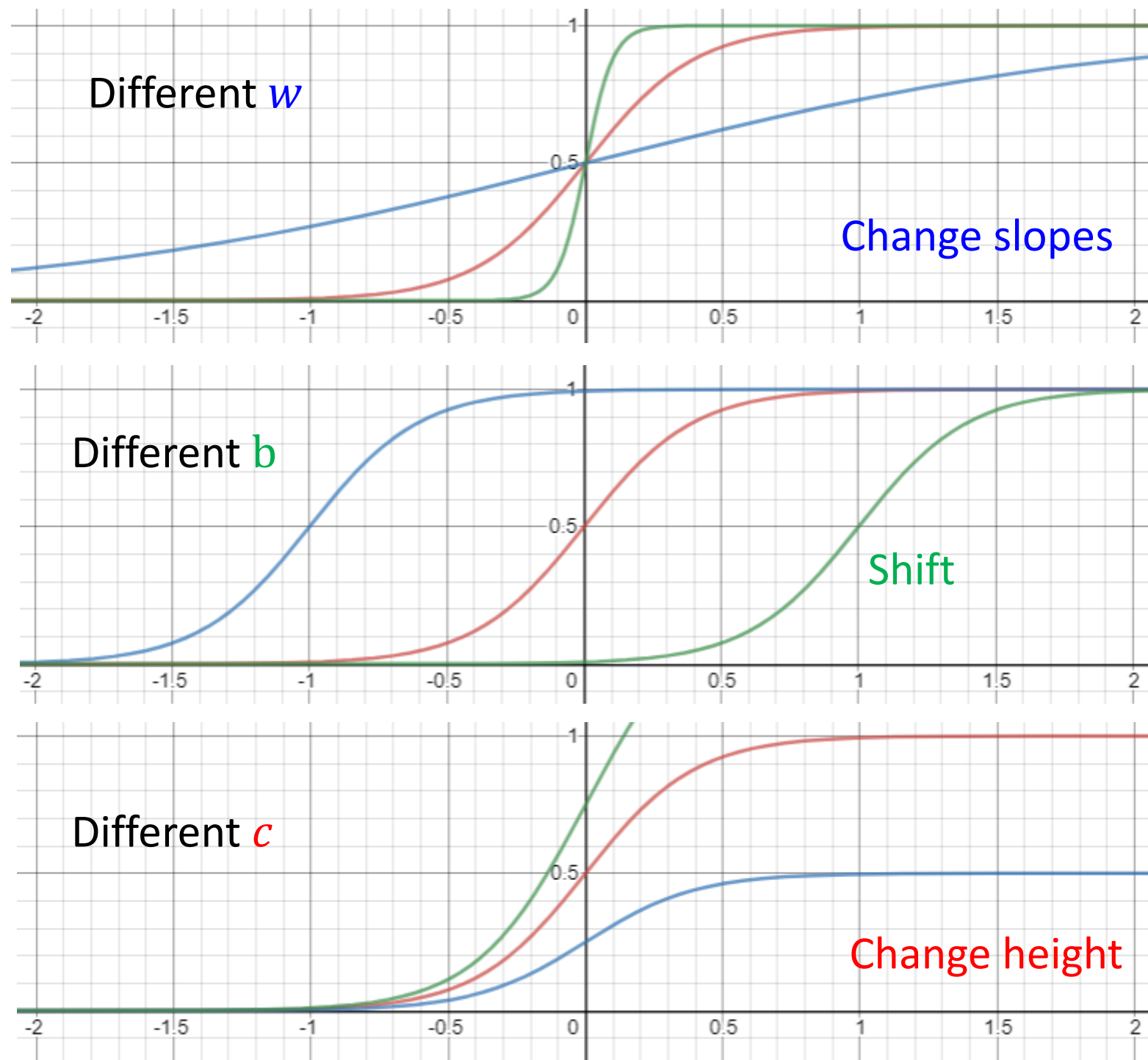
How to represent
this function?

Hard Sigmoid

$x_1$

**Sigmoid Function**

$$y = c \frac{1}{1 + e^{-(b + wx_1)}}$$
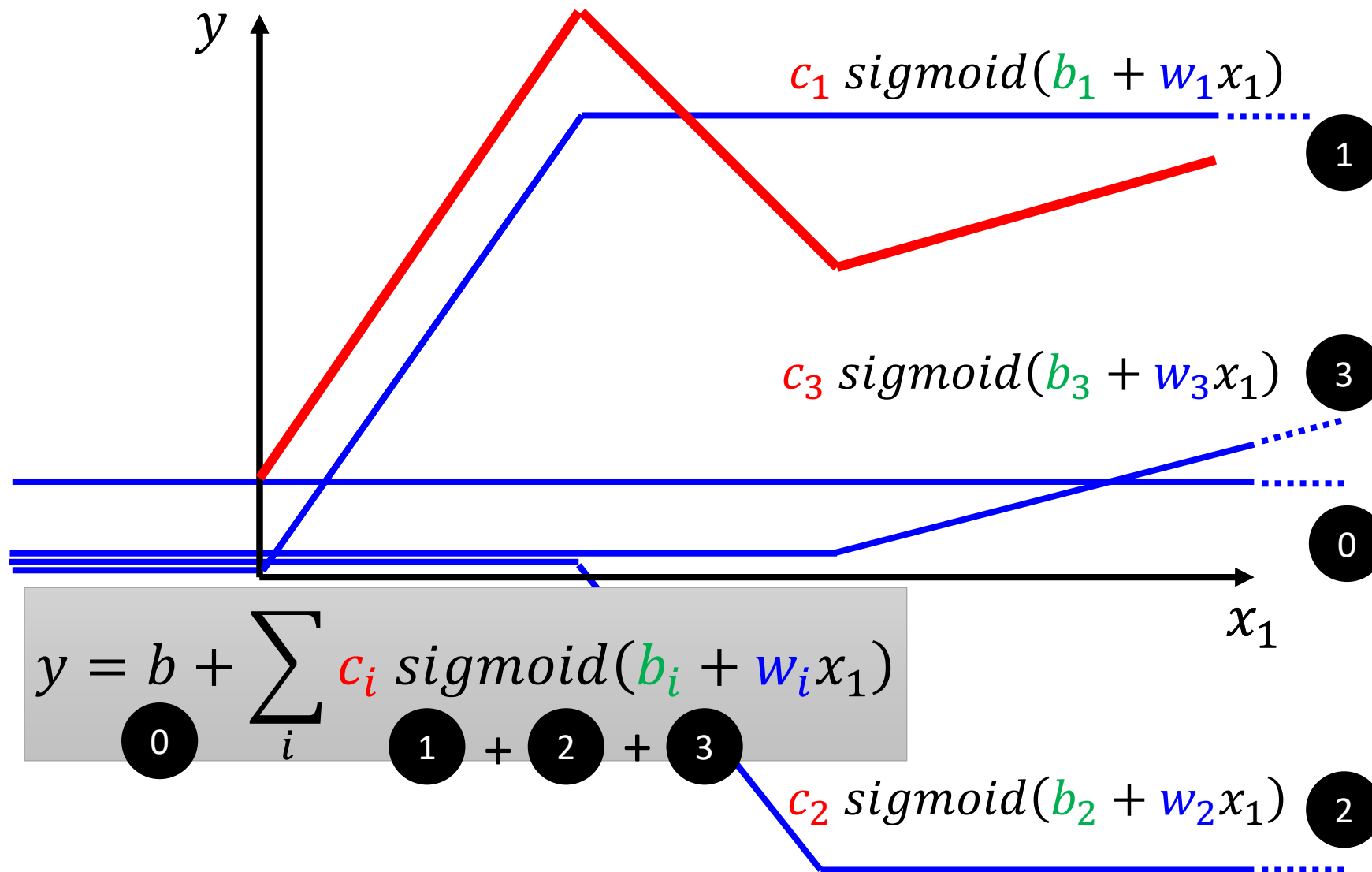
$$= c \; sigmoid(b + wx_1)$$

$x_1$

# Property of sigmoid



Different *w* — Change slopes

Different b — Shift

Different *c* — Change height

red curve = sum of a set of _____ + constant

$y$

$c_1\ sigmoid(b_1 + w_1 x_1)$

**1**

$c_3\ sigmoid(b_3 + w_3 x_1)$   **3**

**0**

$x_1$

$$y = b + \sum_i c_i\ sigmoid(b_i + w_i x_1)$$

**0**   $i$   **1** + **2** + **3**

$c_2\ sigmoid(b_2 + w_2 x_1)$   **2**

$$y = \underline{b + wx_1}$$

$$y = b + \sum_i {\color{red}c_i}\, sigmoid(\underline{{\color{green}b_i} + {\color{blue}w_i}x_1})$$

$$y = \underline{b + \sum_j w_j x_j}$$

$$y = b + \sum_i {\color{red}c_i}\, sigmoid\left({\color{green}b_i} + \underline{\sum_j {\color{blue}w_{ij}}x_i}\right)$$

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

$j: 1,2,3$   no. of features

$i: 1,2,3$   no. of sigmoid

$r_1 = b_1 + w_{11}x_1 + w_{12}x_2 + w_{13}x_3$

$w_{ij}$: weight for $x_j$ for i-th sigmoid

$r_2 = b_2 + w_{21}x_1 + w_{22}x_2 + w_{23}x_3$

$r_3 = b_3 + w_{31}x_1 + w_{32}x_2 + w_{33}x_3$

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

$i: 1,2,3$
$j: 1,2,3$

$$r_1 = b_1 + w_{11}x_1 + w_{12}x_2 + w_{13}x_3$$

$$r_2 = b_2 + w_{21}x_1 + w_{22}x_2 + w_{23}x_3$$

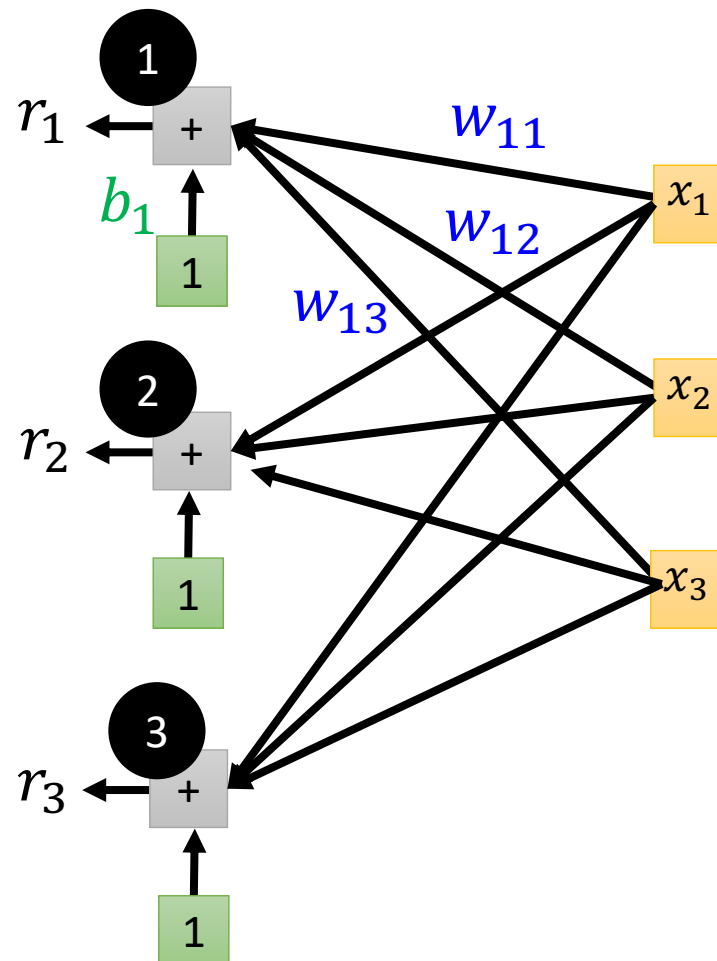$$r_3 = b_3 + w_{31}x_1 + w_{32}x_2 + w_{33}x_3$$

$$\begin{bmatrix} r_1 \\ r_2 \\ r_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} + \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\boldsymbol{r} = \boldsymbol{b} + W \quad \boldsymbol{x}$$

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

$i: 1,2,3$
$j: 1,2,3$

$$r = b + W \; x$$

$$y = b + \sum_i \textcolor{red}{c_i} \boxed{sigmoid \left( \textcolor{green}{b_i} + \sum_j \textcolor{blue}{w_{ij}} x_j \right)}$$

$i: 1,2,3$

$j: 1,2,3$



$a_1 = sigmoid(r_1) = \dfrac{1}{1 + e^{-r_1}}$

$\boxed{\boldsymbol{a}} = \sigma( \boxed{\boldsymbol{r}} )$

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

$i: 1,2,3$
$j: 1,2,3$



$$\boxed{y} = \boxed{b} + \boxed{\boldsymbol{c}^T} \; \boxed{\boldsymbol{a}}$$

$$y = b + \boldsymbol{c}^T \boldsymbol{a}$$

$$\boldsymbol{a} = \sigma(\boldsymbol{r}) \qquad \boldsymbol{r} = \boldsymbol{b} + W\boldsymbol{x}$$

$$y = \boxed{b} + \boxed{\boldsymbol{c}^T}\ \sigma(\ \boxed{\boldsymbol{b}} + \boxed{W}\ \boxed{\boldsymbol{x}}\ )$$

# Function with unknown parameters

$$y = b + \boldsymbol{c}^T \sigma( \boldsymbol{b} + W \boldsymbol{x} )$$

$\boldsymbol{x}$ feature

**Unknown parameters**

$W$  $\boldsymbol{b}$

$\boldsymbol{c}^T$  $b$

Rows of $W$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \end{bmatrix}$$

Step 1: function with unknown → Step 2: define loss from training data → Step 3: optimization

$$y = b + \boldsymbol{c}^T \sigma(\boldsymbol{b} + W\boldsymbol{x})$$

# Loss

➢ Loss is a function of parameters $L(\theta)$

➢ Loss means how good a set of values is.

**feature**

$y = b + c^T \sigma(b + Wx)$

$e$

**label** $\hat{y}$

Given a set of values

Loss: $L = \dfrac{1}{N}\sum_{n} e_n$

Step 1: function with unknown → Step 2: define loss from training data → **Step 3: optimization**

$$y = b + \boldsymbol{c}^T \sigma( \boldsymbol{b} + W \boldsymbol{x} )$$

# Optimization of New Model

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \end{bmatrix}$$

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

➢ (Randomly) Pick initial values $\boldsymbol{\theta}^0$

$$\boldsymbol{g} = \begin{bmatrix} \dfrac{\partial L}{\partial \theta_1}|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \dfrac{\partial L}{\partial \theta_2}|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix}$$

**gradient**

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \\ \vdots \end{bmatrix} - \begin{bmatrix} \textcolor{red}{\eta}\dfrac{\partial L}{\partial \theta_1}|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \textcolor{red}{\eta}\dfrac{\partial L}{\partial \theta_2}|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix}$$

$$\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^0)$$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \textcolor{red}{\eta}\boldsymbol{g}$$

# Optimization of New Model

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

➢ (Randomly) Pick initial values $\boldsymbol{\theta}^0$

➢ Compute gradient $\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^0)$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \boldsymbol{g}$$

➢ Compute gradient $\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^1)$

$$\boldsymbol{\theta}^2 \leftarrow \boldsymbol{\theta}^1 - \eta \boldsymbol{g}$$

➢ Compute gradient $\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^2)$

$$\boldsymbol{\theta}^3 \leftarrow \boldsymbol{\theta}^2 - \eta \boldsymbol{g}$$

# Optimization of New Model

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

➤ (Randomly) Pick initial values $\boldsymbol{\theta}^0$

➤ Compute gradient $\boldsymbol{g} = \nabla L^1(\boldsymbol{\theta}^0)$   $L^1$

  **update** $\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \boldsymbol{g}$

➤ Compute gradient $\boldsymbol{g} = \nabla L^2(\boldsymbol{\theta}^1)$   $L^2$

  **update** $\boldsymbol{\theta}^2 \leftarrow \boldsymbol{\theta}^1 - \eta \boldsymbol{g}$

➤ Compute gradient $\boldsymbol{g} = \nabla L^3(\boldsymbol{\theta}^2)$   $L^3$

  **update** $\boldsymbol{\theta}^3 \leftarrow \boldsymbol{\theta}^2 - \eta \boldsymbol{g}$

1 **epoch** = see all the batches once

B

batch

batch

batch

batch

$L$

N

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

Training:

| Step 1: function with unknown | Step 2: define loss from training data | Step 3: optimization |

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

$$L(\theta)$$

$$\theta^* = arg \min_\theta L$$

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

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

$$\{y^{N+1}, y^{N+2}, \dots, y^{N+M}\}$$

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

*An extreme example*

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

$$f(\boldsymbol{x}) = \begin{cases} \hat{y}^i & \exists \boldsymbol{x^i} = \boldsymbol{x} \\ random & otherwise \end{cases}$$   Less than useless ...

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**

**Data augmentation**

# Overfitting



**constrained** model

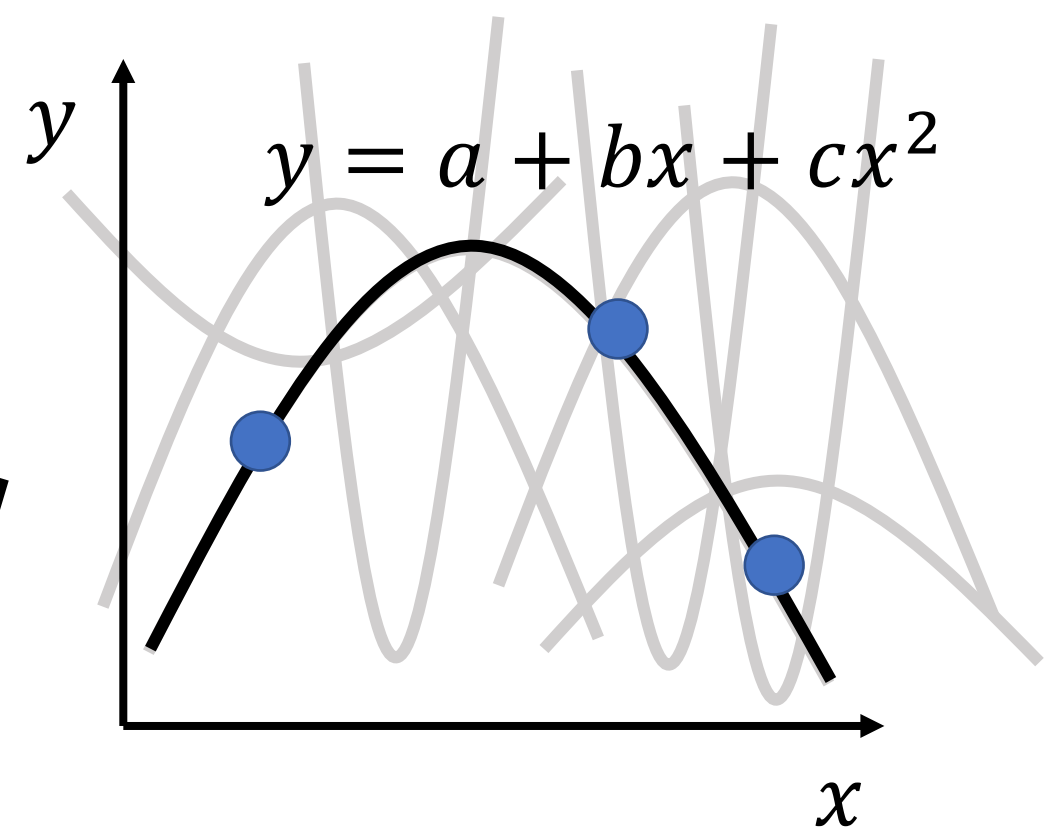$$y = a + bx + cx^2$$

•••••• Real data distribution (not observable)

● Training data
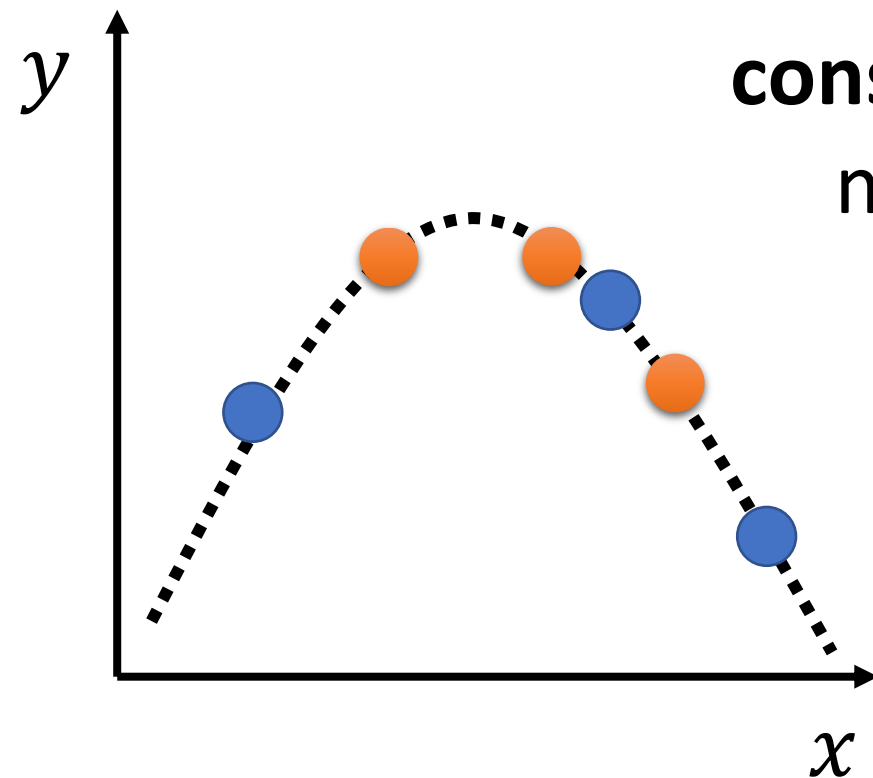
● Testing data

# Overfitting

**constrained**
model

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

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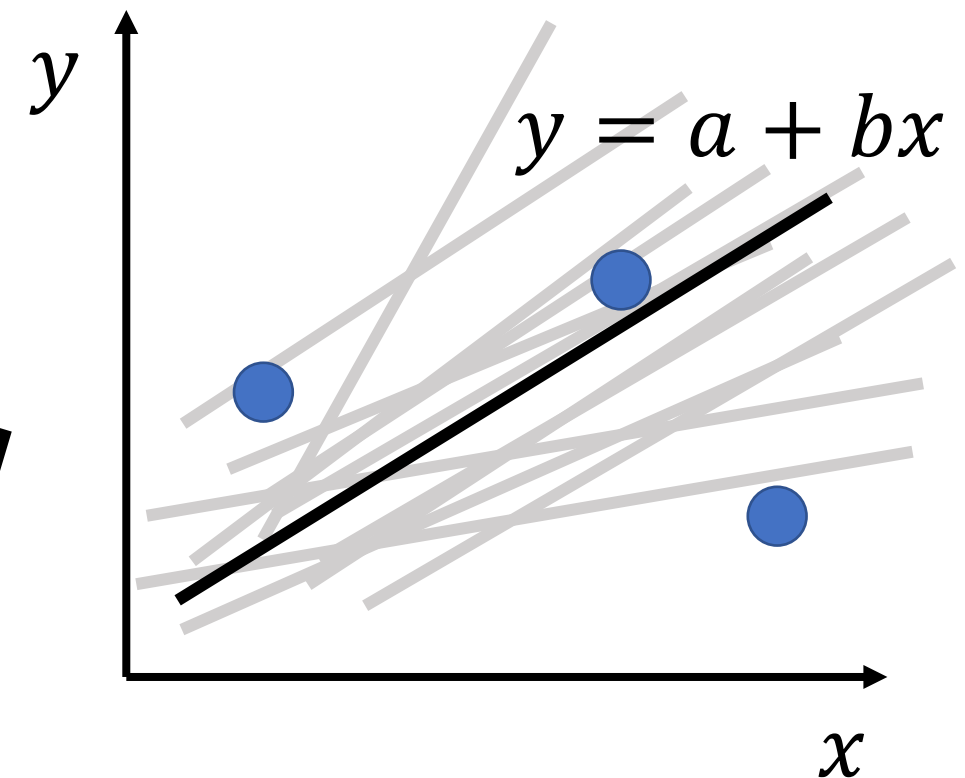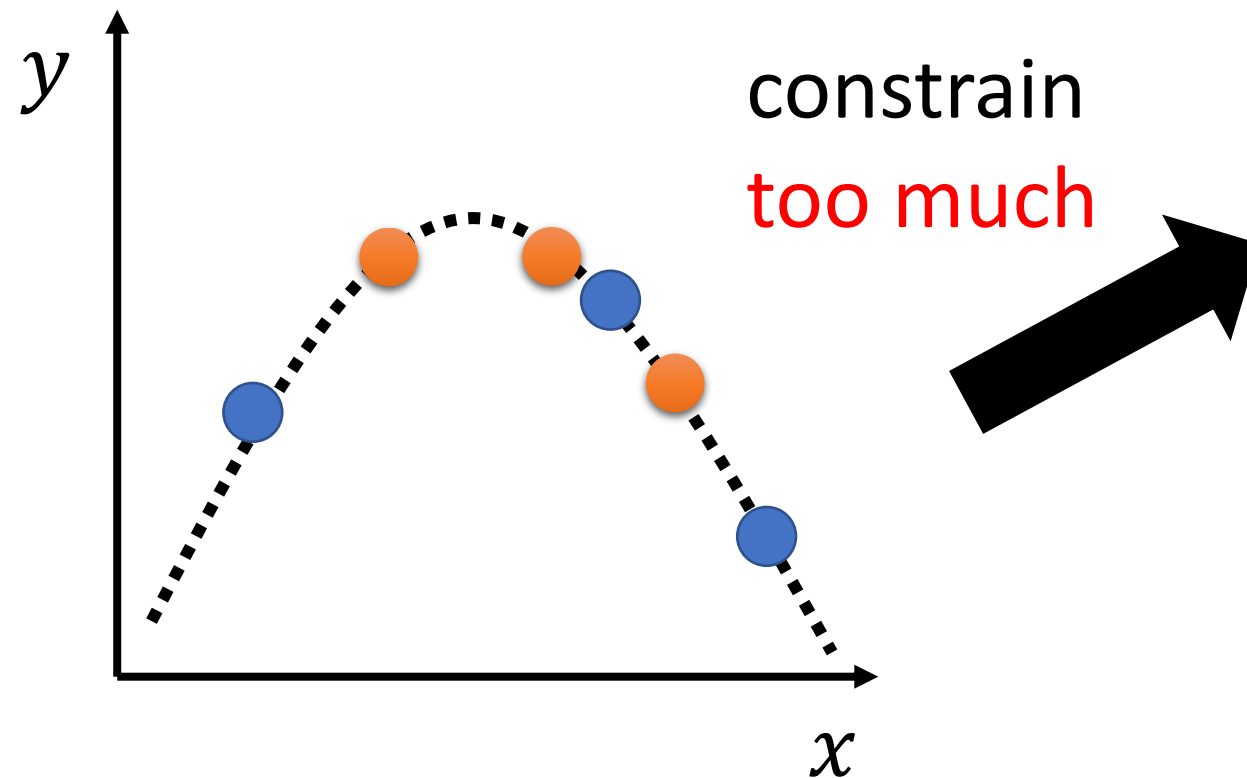- 🔵 Training data
- 🟠 Testing data

# Overfitting

**constrained** model

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

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

# Overfitting



constrain
too much

$y = a + bx$

Back to model bias …

- - - - - Real data distribution
(not observable)
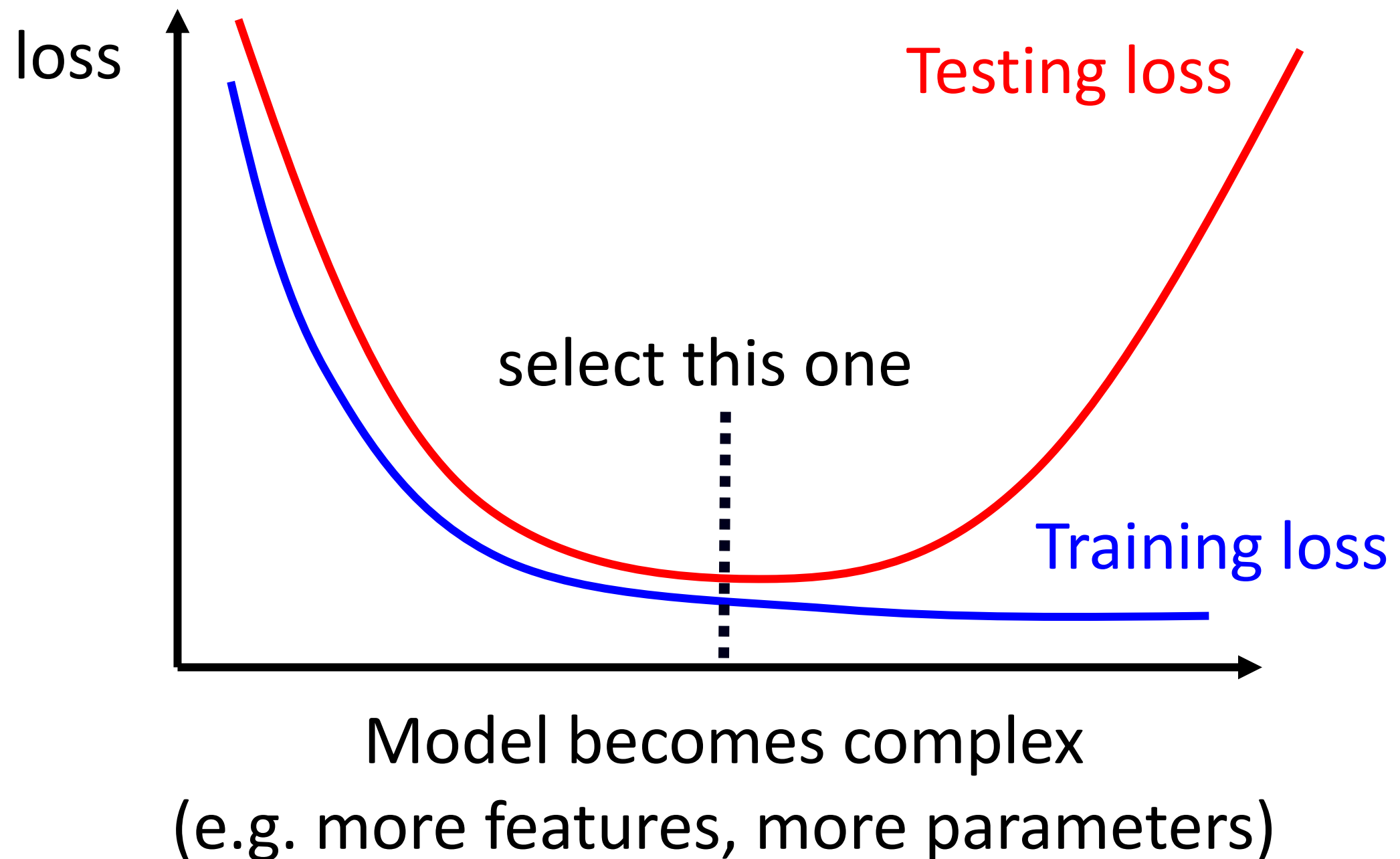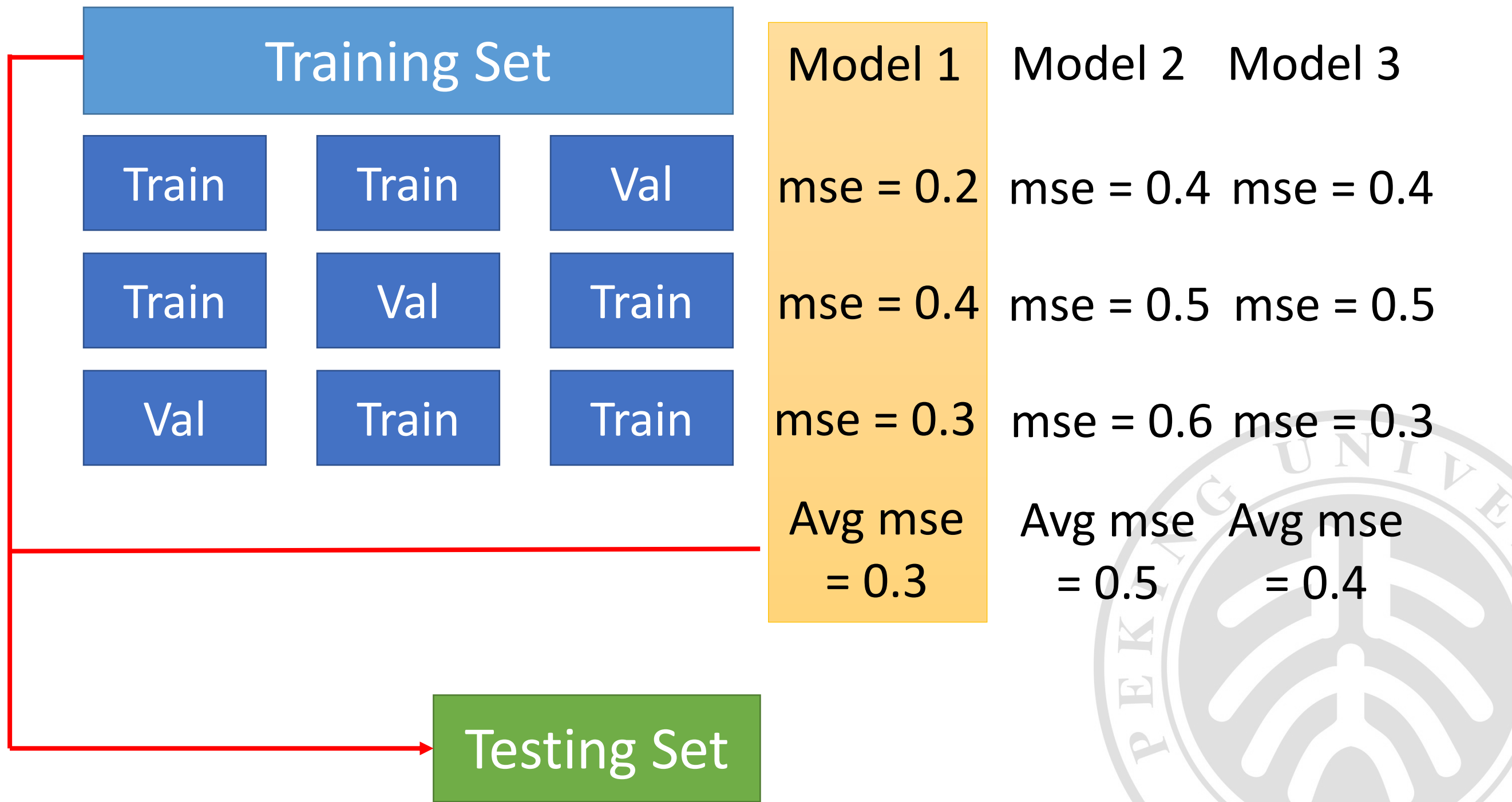
● Training data

● Testing data

# Bias-Complexity Trade-off

- What information is available for learning?
  - What does the data look like?
  - How is it annotated?

- What output is desired?
  - What should the algorithm produce?
  - How will it be used?

- Learning with a teacher
  - Explicit feedback in the form of labeled examples
  - Goal: make prediction
  - Pros: Good performance
  - Cons: Labeled data is difficult to find
  - "Classical" (labeled data simply there; do your best)
  - Query-based (can ask for labeled examples)
- Examples
  - Classification
    - Is an email spam or not?
  - Regression
    - What is the expected rate of return on a specific investment?



22-11 = ?

- Learning without labels
  - Only observed unlabeled examples
  - Goal: uncover structure in data
  - Pros: Easy to find lots of data
  - Cons: what are we looking for?

- Examples
  - Clustering
    - Group emails by topic
  - Manifold learning
    - Find a low dimensional data representation

- Learn a behavior policy by interacting with the world
  - How to navigate in a world
  - Success measured by rewards received for actions taken
  - Maximize sum of rewards
- Examples
  - Chess (and checkers)
  - Robot control
  - Piloting an airplane