

Stochastic gradient descent

GD minimizes:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{e(h(\mathbf{x}_n), y_n)}_{\ln(1+e^{-y_n \mathbf{w}^\top \mathbf{x}_n})} \leftarrow \text{in logistic regression}$$

by iterative steps along $-\nabla E_{\text{in}}$:

$$\Delta \mathbf{w} = -\eta \nabla E_{\text{in}}(\mathbf{w})$$

∇E_{in} is based on all examples (\mathbf{x}_n, y_n)

“batch” GD

The stochastic aspect

Pick one (\mathbf{x}_n, y_n) at a time. Apply GD to $\mathbf{e}(h(\mathbf{x}_n), y_n)$

“Average” direction:

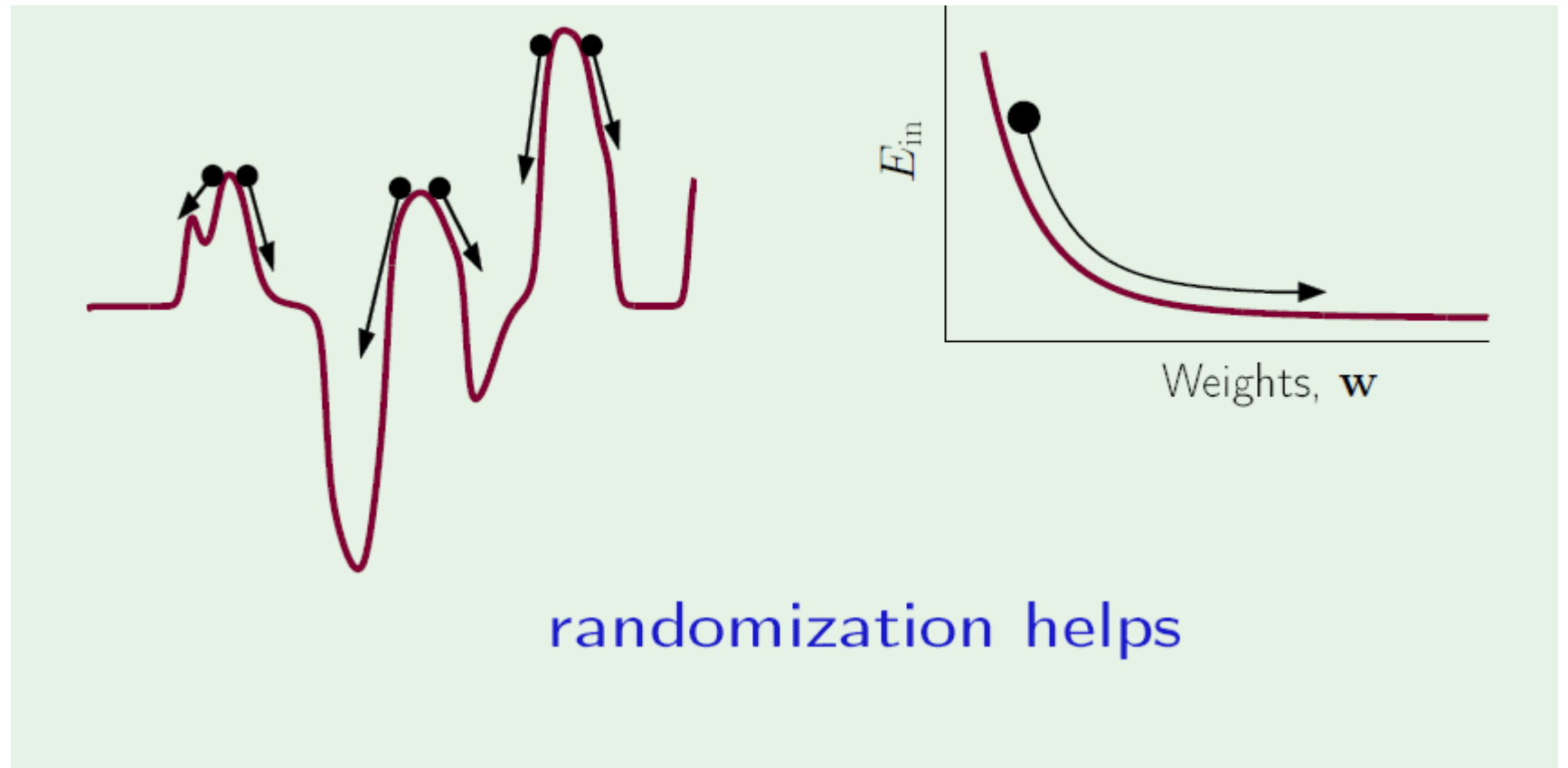
$$\begin{aligned}\mathbb{E}_n [-\nabla \mathbf{e}(h(\mathbf{x}_n), y_n)] &= \frac{1}{N} \sum_{n=1}^N -\nabla \mathbf{e}(h(\mathbf{x}_n), y_n) \\ &= -\nabla E_{\text{in}}\end{aligned}$$

randomized version of GD

stochastic gradient descent (SGD)

Benefits of SGD

1. cheaper computation
2. randomization
3. simple



Linear Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: **real value**)

- Linear model

$$y(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$e(w) = \frac{1}{N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Goal: minimize $e(w)$

Closed-form solution

Directly calculate the pseudo-inverse

Logistic Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: **0 or 1**)

- Linear model

$$y(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M)$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$e(w) = \frac{1}{N} \sum_{i=1}^N \ln \left[1 + e^{-y_n w^T x_n} \right]$$

- Goal: minimize $e(w)$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Linear Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: **real value**)

- Linear model

$$y(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$e(w) = \frac{1}{N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Goal: minimize $e(w)$

Closed-form solution

Directly calculate the pseudo-inverse

Logistic Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: **0 or 1**)

- Linear model

$$y(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M)$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$e(w) = \frac{1}{N} \sum_{i=1}^N \ln [1 + e^{-y_n w^T x_n}]$$

- Goal: minimize $e(w)$

Steps:

- Initialize w (e.g., randomly)
- Repeatedly update w based on the gradient

$$w = w - \epsilon \nabla_w \ell(w)$$

where ϵ is the learning rate.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Key Concepts in Supervised Learning

- Loss function **(measure error, or judge the fit)**
- Optimization **(how to find a good fit)**
- Generalization **(fit to unseen test data)**
- Regularization **(avoid overfitting)**

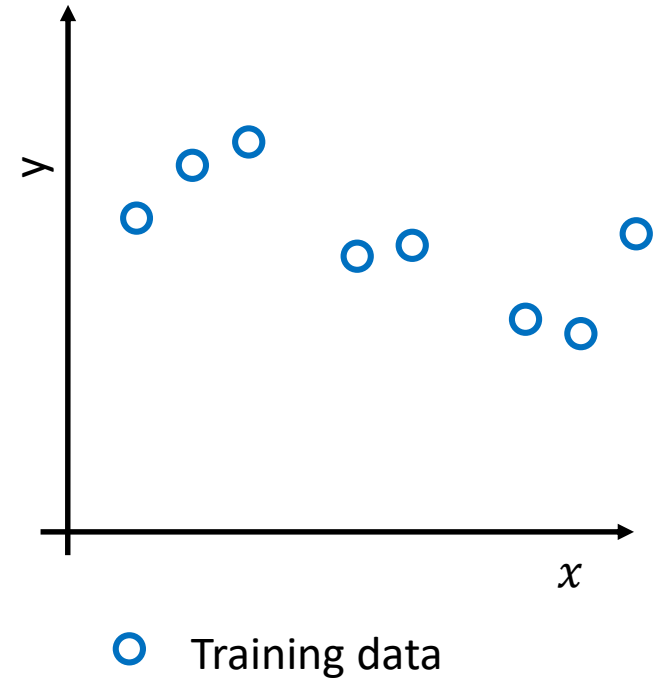
(very very ... very important!!!)

Linear Models for Regression

- **Let's discuss those key questions in the following cases**
 - **Generalization in Supervised Learning (one key concept)**
 - Linear regression using polynomial fitting

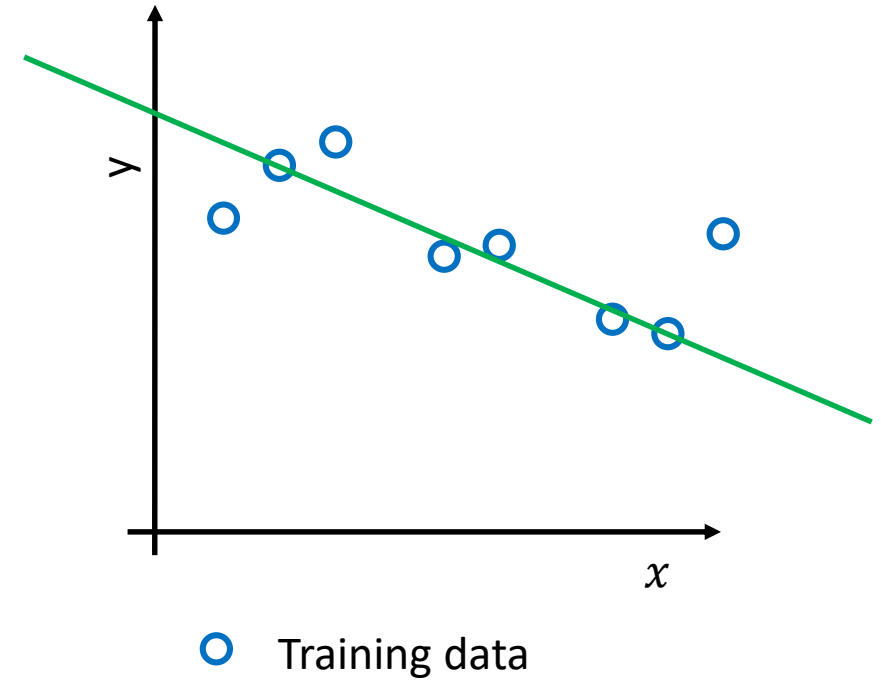
Generalization

- Generalization
 - Model's ability to predict new data



Generalization

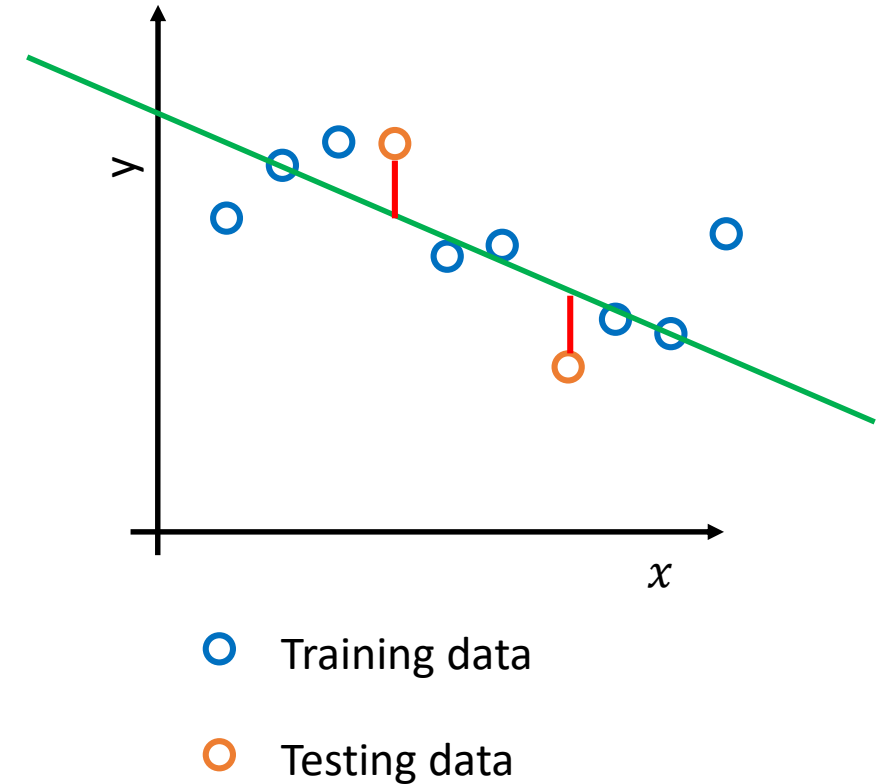
- Generalization
 - Model's ability to predict new data



Generalization

- Generalization
 - Model's ability to predict new data

In fact, what we really care about is the error on new data (in Testing datasets)



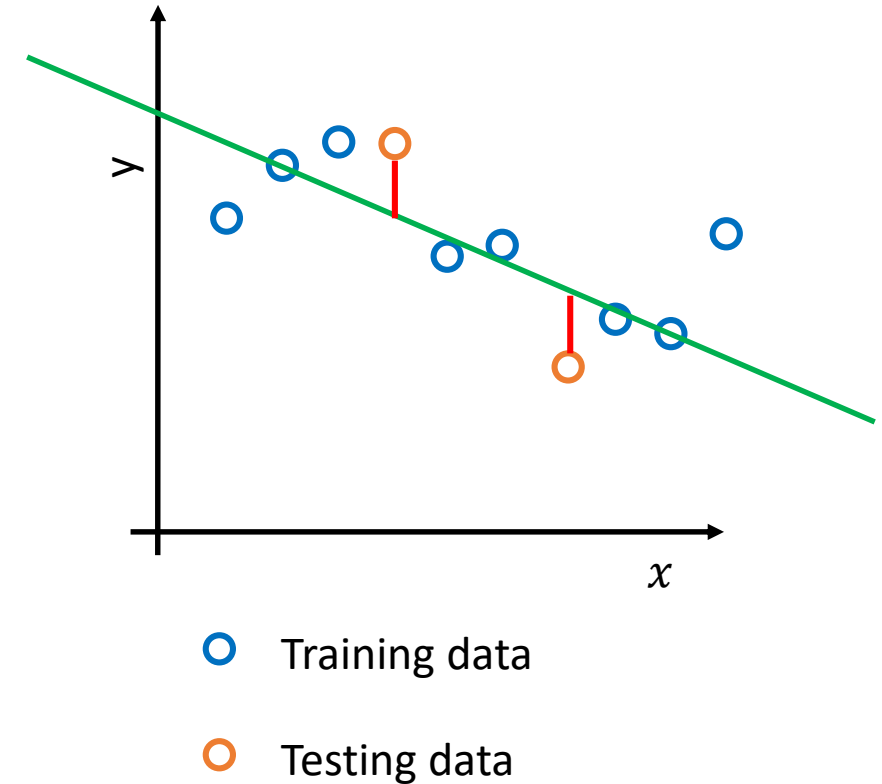
Generalization

- Generalization
 - Model's ability to predict new data

In fact, what we really care about is the error on new data (in Testing dataset)

Datasets in Machine Learning

- Training dataset
- Validation Dataset (important! we will discuss it soon.)
- Testing dataset

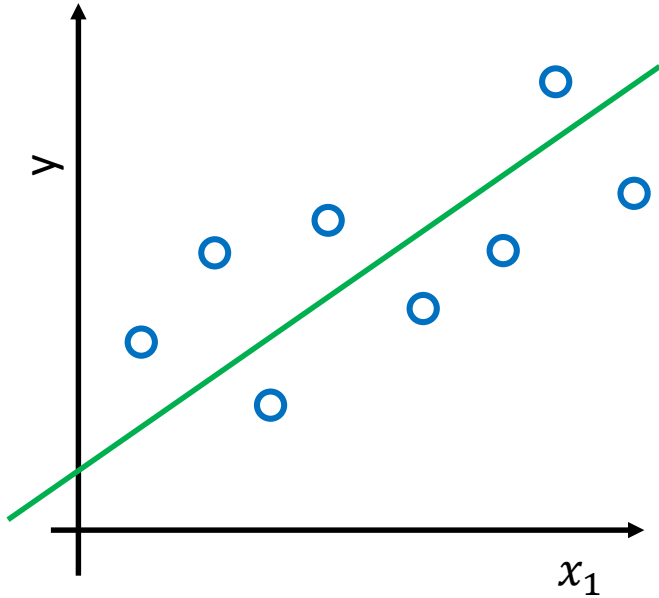


Hyperplane

Linear regression models we learned:

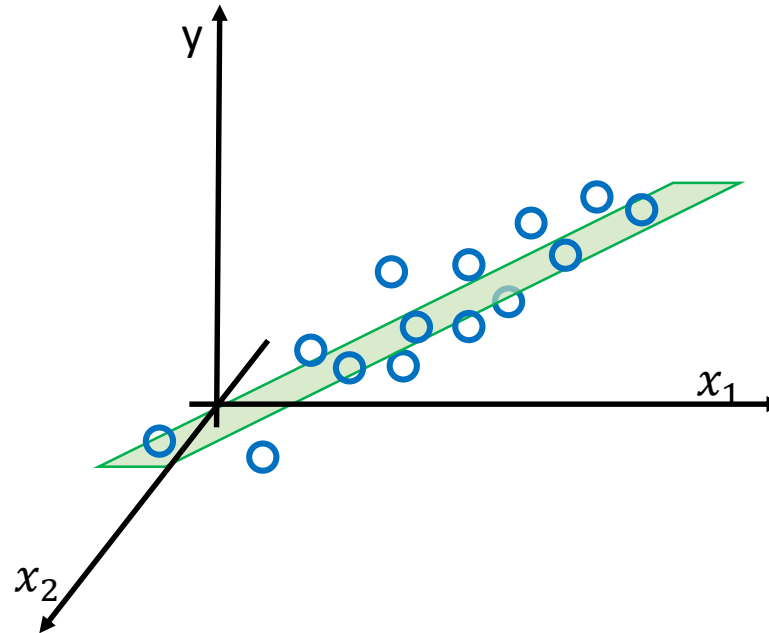
- One variable

$$y(x) = w_0 + w_1x_1$$



- Two variables

$$y(x) = w_0 + w_1x_1 + w_2x_2$$



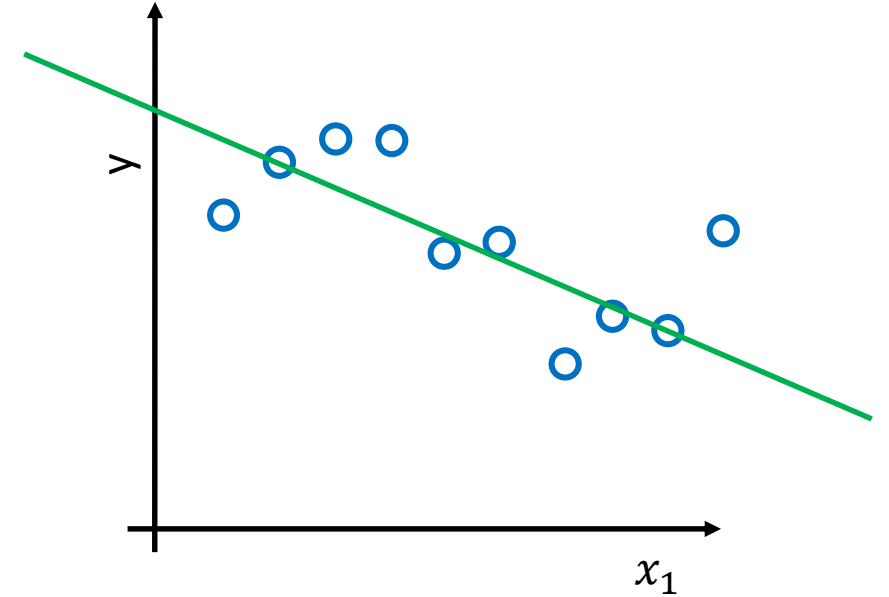
- More variables

$$y(x) = w_0 + w_1x_1 + \cdots + w_Nx_N$$

?

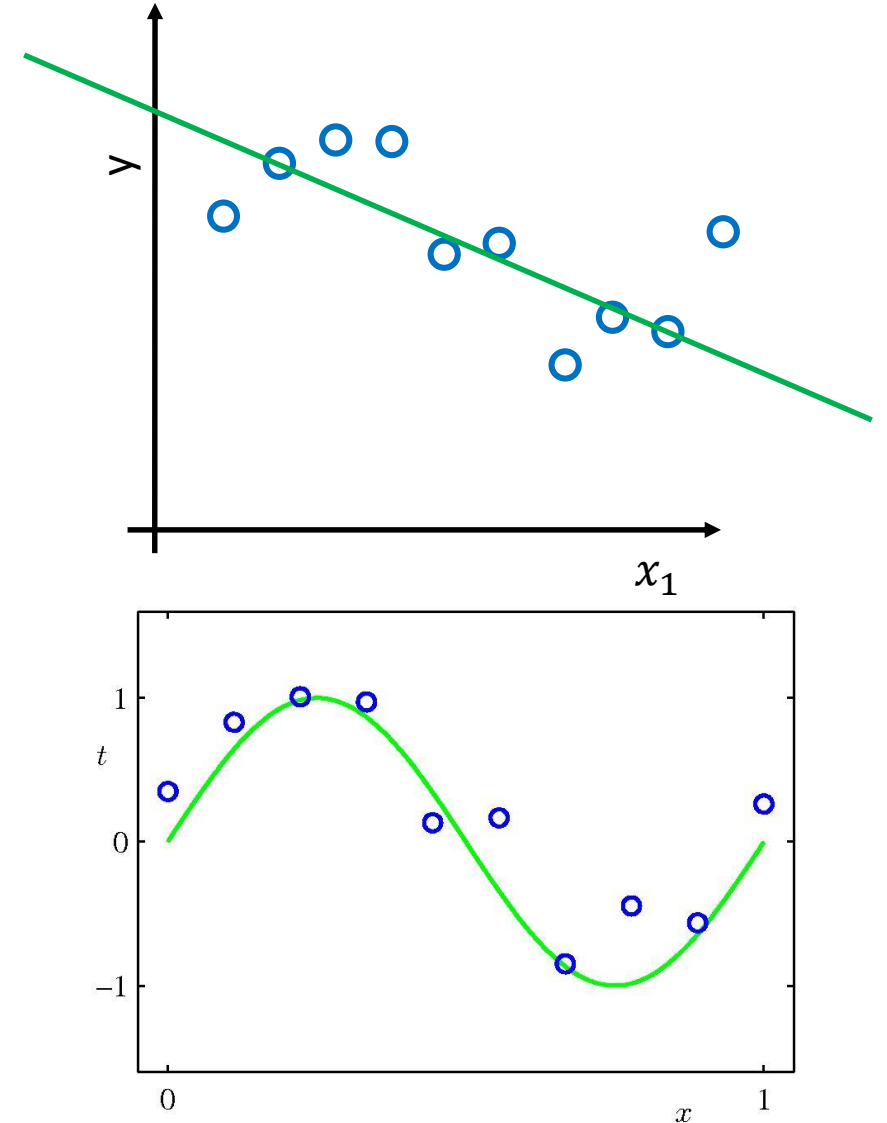
Generalization

- What if our linear model is not good? (for the data to right)



Generalization

- What if our linear model is not good? (for the data to right)
- We can use a more complicated model (polynomial)



Linear Models for Regression

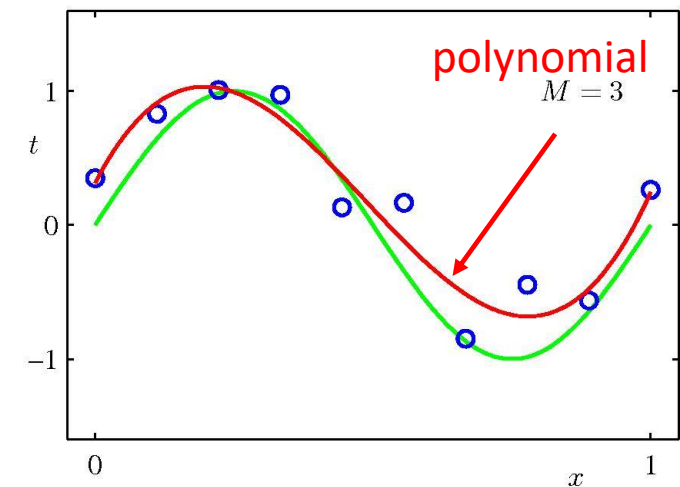
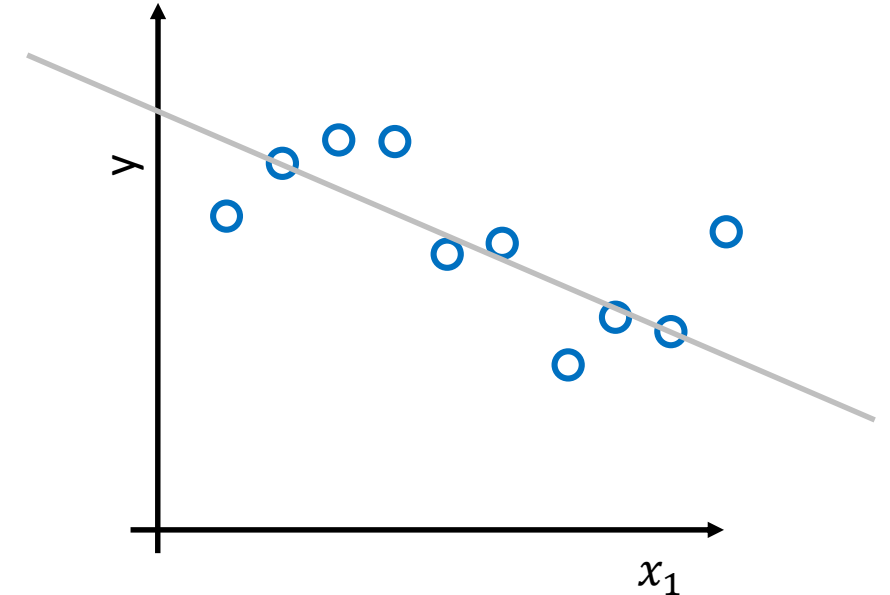
- **Let's discuss those key questions in the following cases**
 - Generalization in Supervised Learning (one key concept)
 - **Linear regression using polynomial fitting**

Fitting a Polynomial

- Example: an M -th order polynomial function of one dimensional feature x :

$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$

where x^j is the j -th power of x .



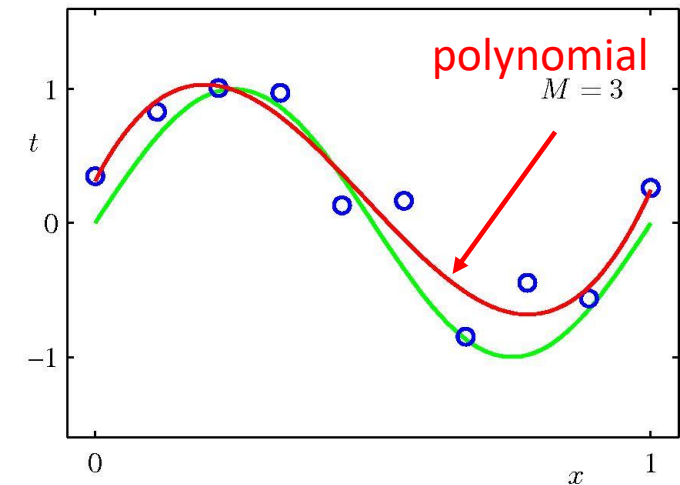
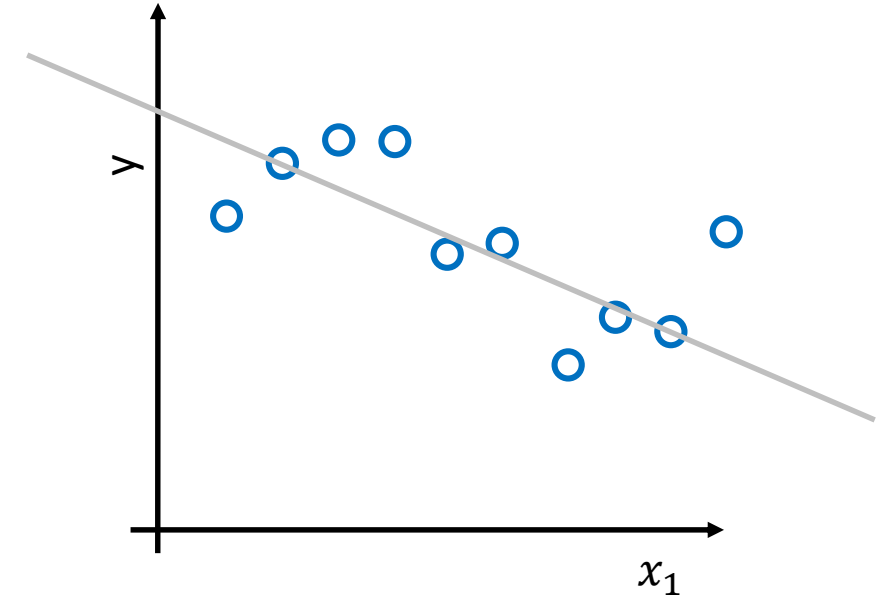
Fitting a Polynomial

- Example: an M -th order polynomial function of one dimensional feature x :

$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$

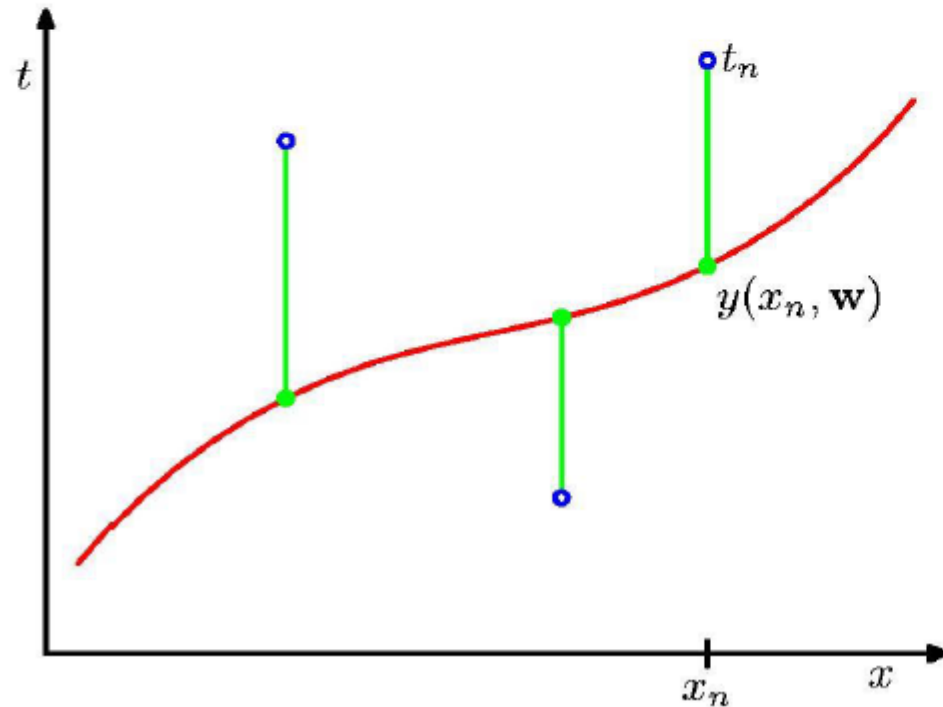
where x^j is the j -th power of x .

- Note: We can optimize for the weights w by using the same approach as we did for previous linear model.



Sum-of-Squares Error Function

The values of the coefficients will be determined by fitting the polynomial to the training data. This can be done by minimizing an error function that measures the misfit between the function $y(x, \mathbf{w})$, for any given value of \mathbf{w} , and the training set data points.



The sum of the squares error function:

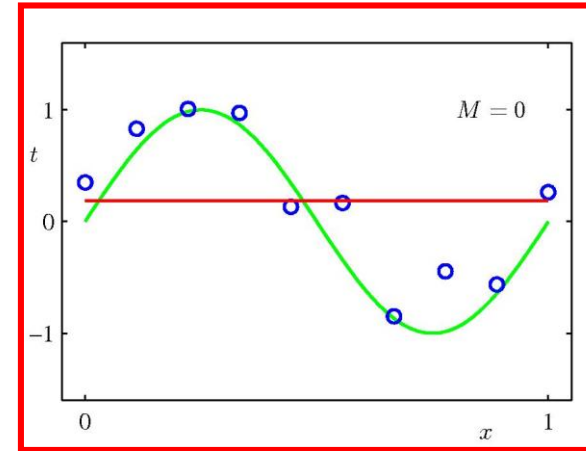
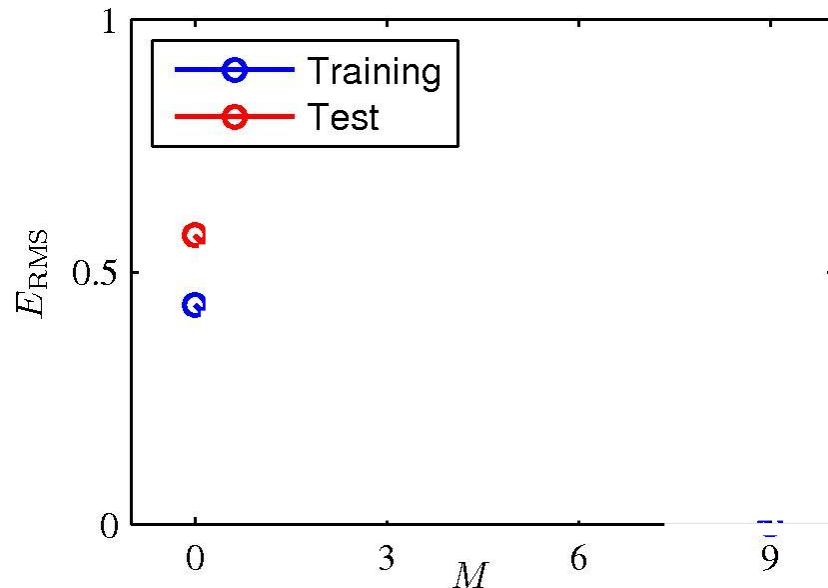
$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

Overfitting

Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

- Let's use polynomial model to explain overfitting

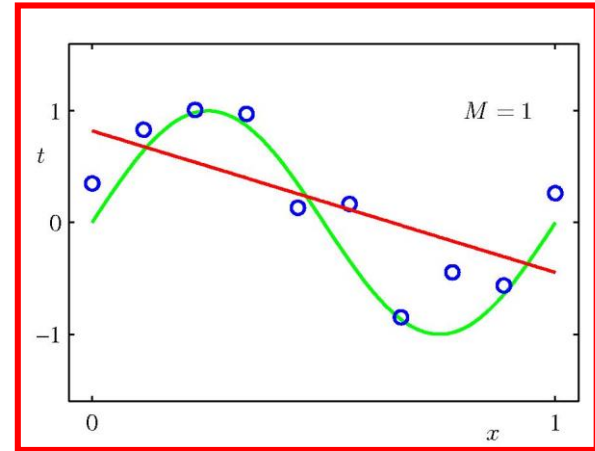
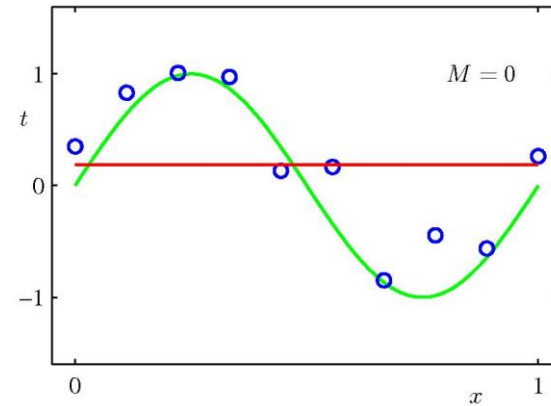
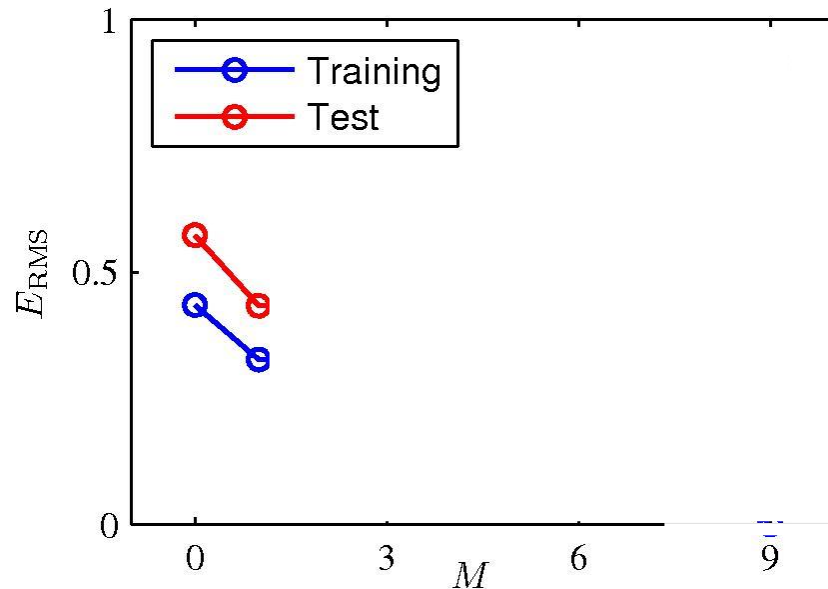
$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$



Overfitting

- Let's use polynomial model to explain overfitting

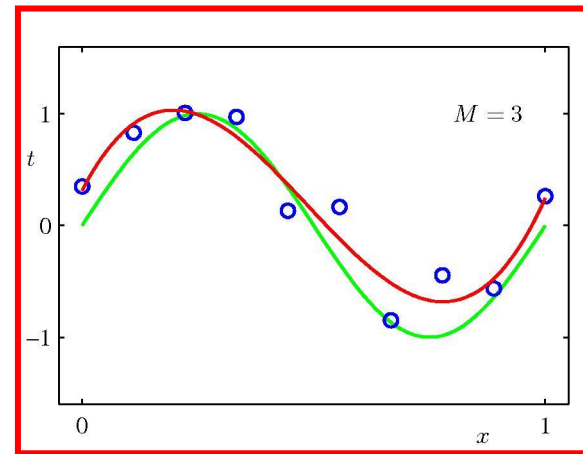
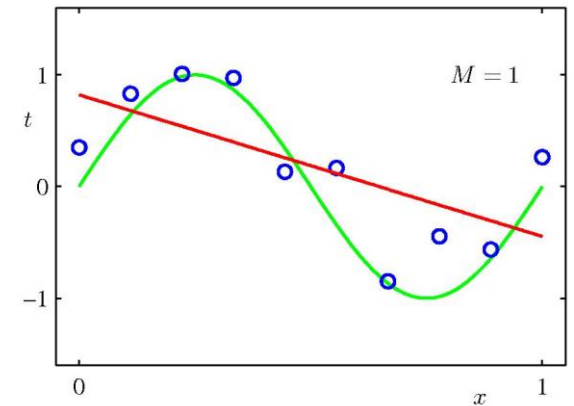
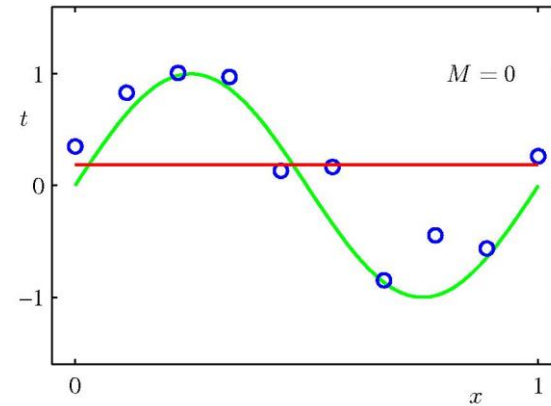
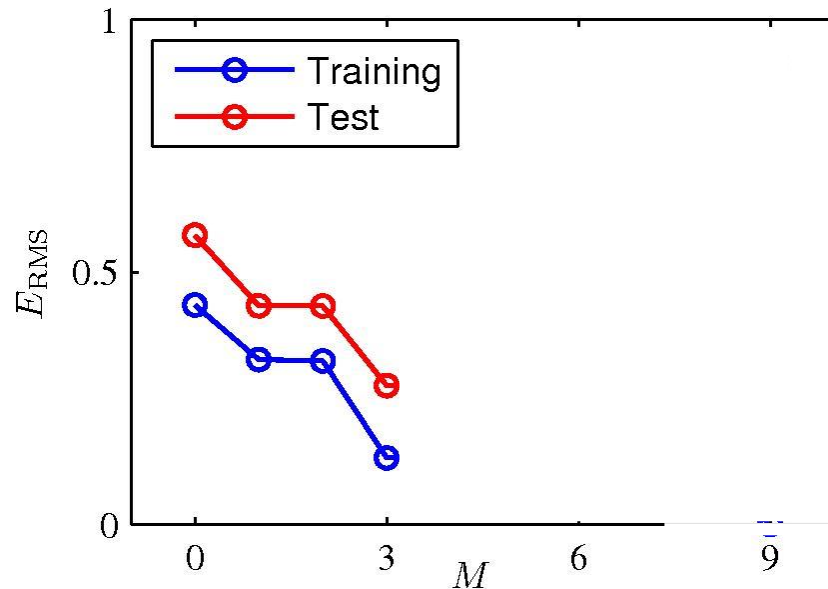
$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$



Overfitting

- Let's use polynomial model to explain overfitting

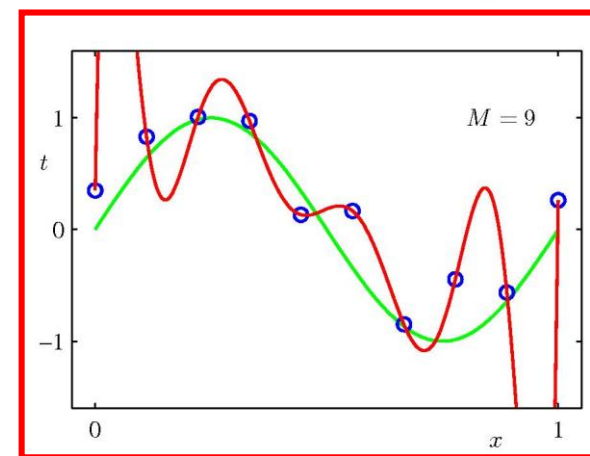
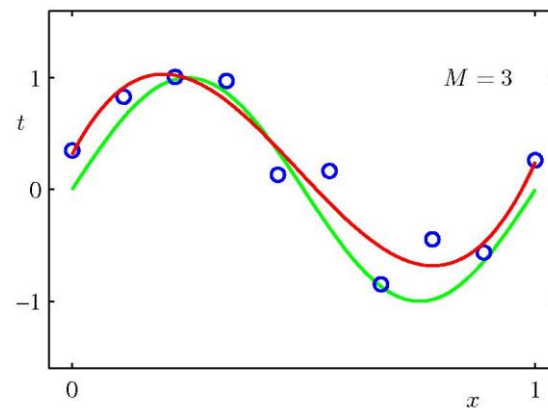
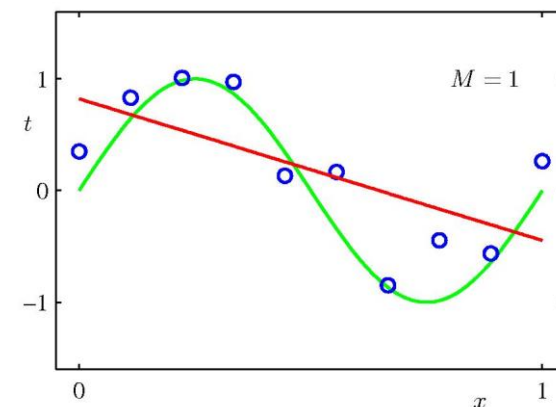
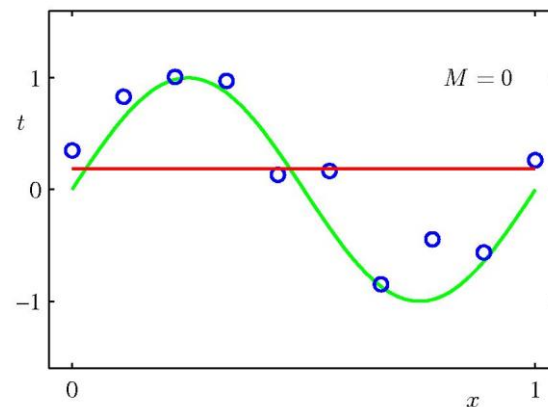
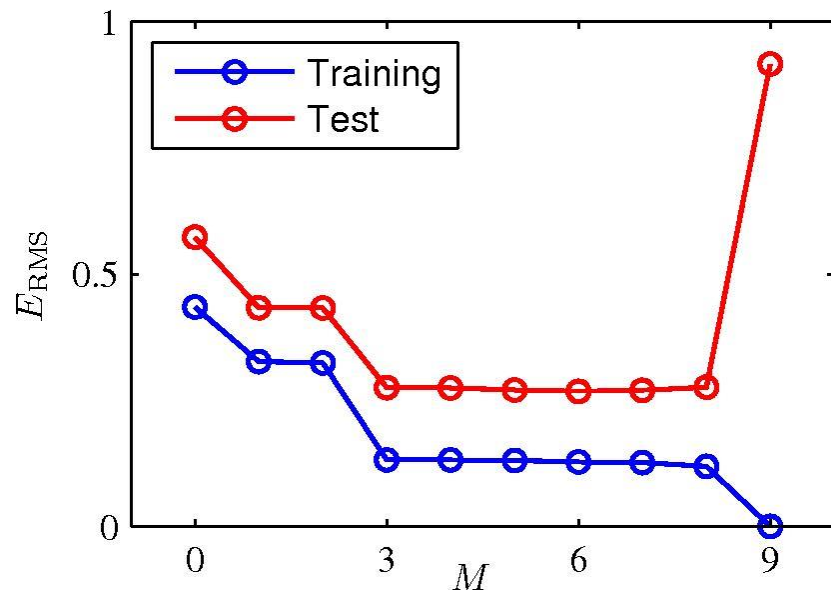
$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$



Overfitting

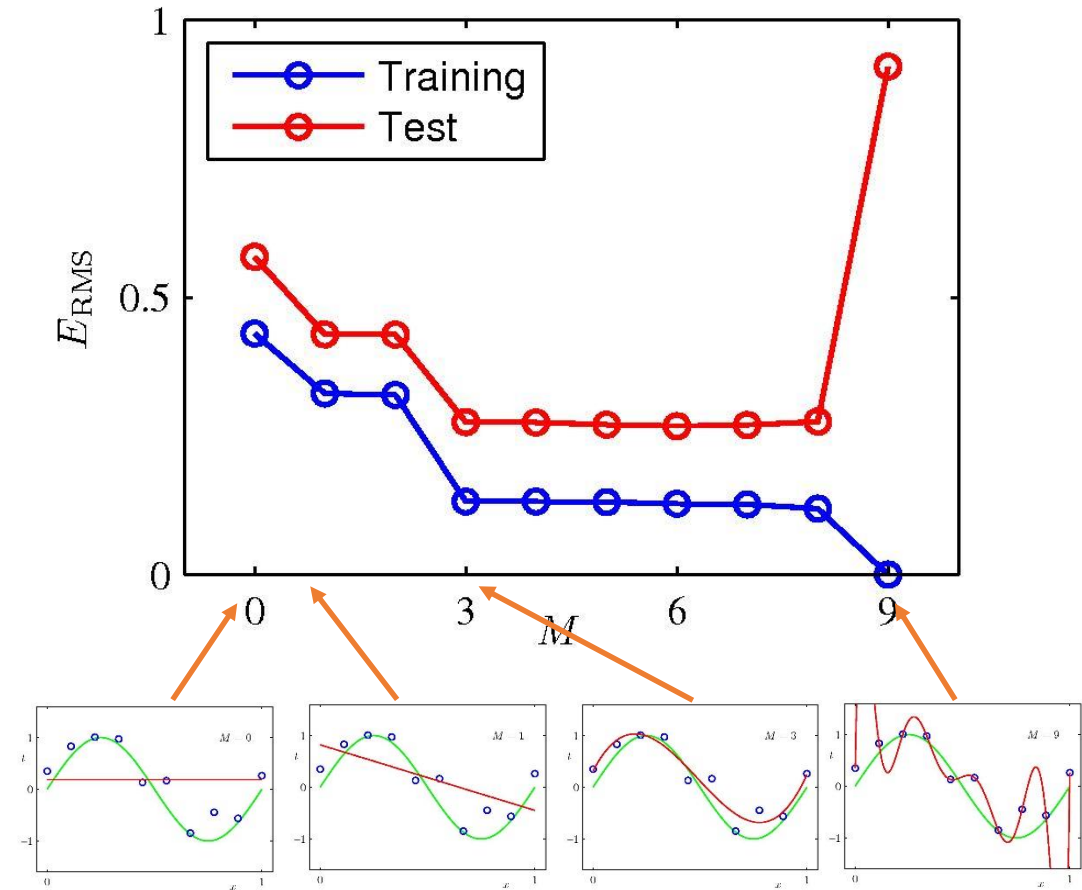
- Let's use polynomial model to explain overfitting

$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$



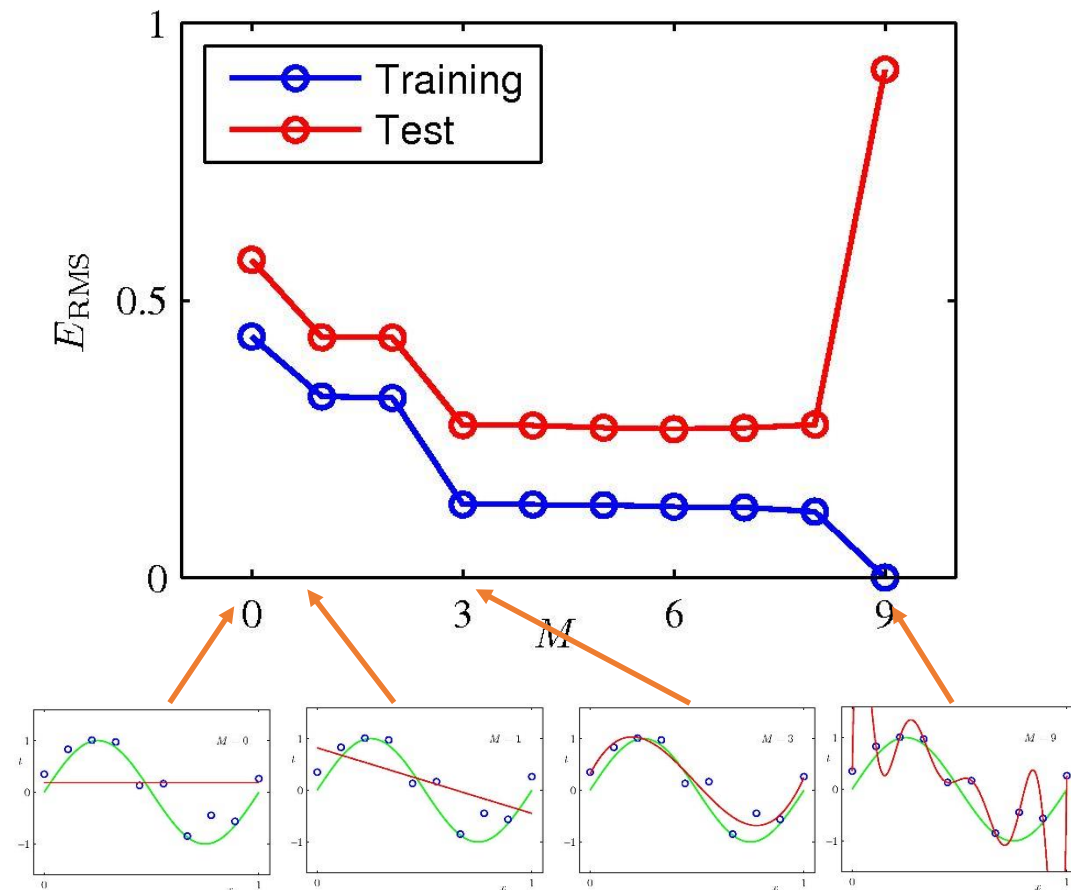
Overfitting

- Observations
 - A more complex model yields lower error on **training data**. (If we choose truly find the best function, the error on training data may go to zero)
 - A more complex model may perform very badly on **testing data** (Our model with $M = 9$ **overfits** the data)



Question?

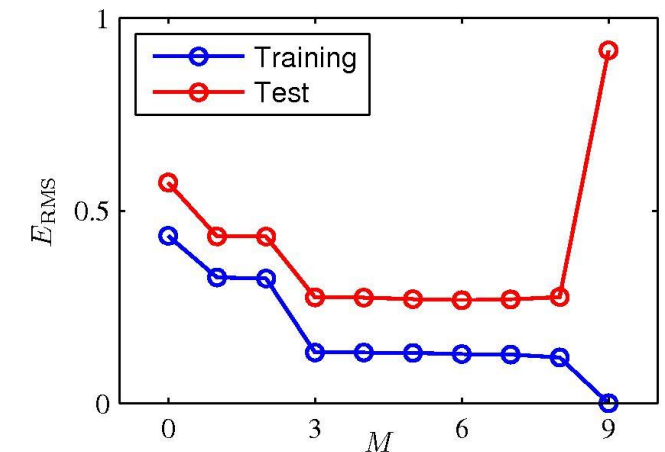
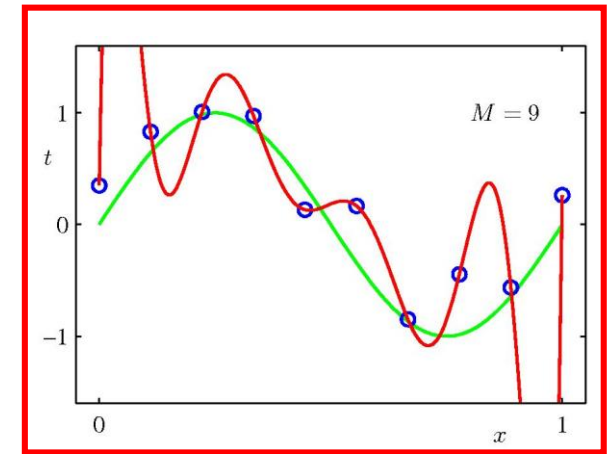
- Consider $y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$
- Why a more complex model can yield lower error on **training data**?
- Why a more complex model can perform very badly on **testing data (overfitting)**?



Overfitting

- Let's look at the estimated weights for various M

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



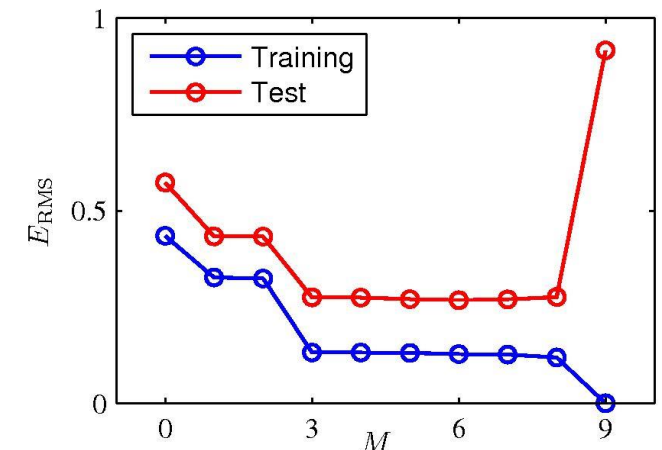
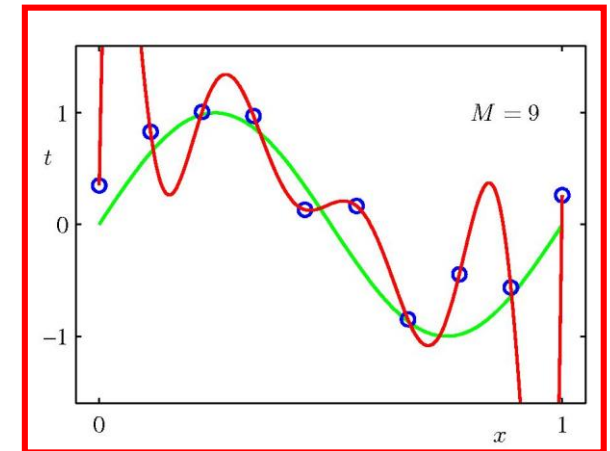
Overfitting

- Let's look at the estimated weights for various M

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

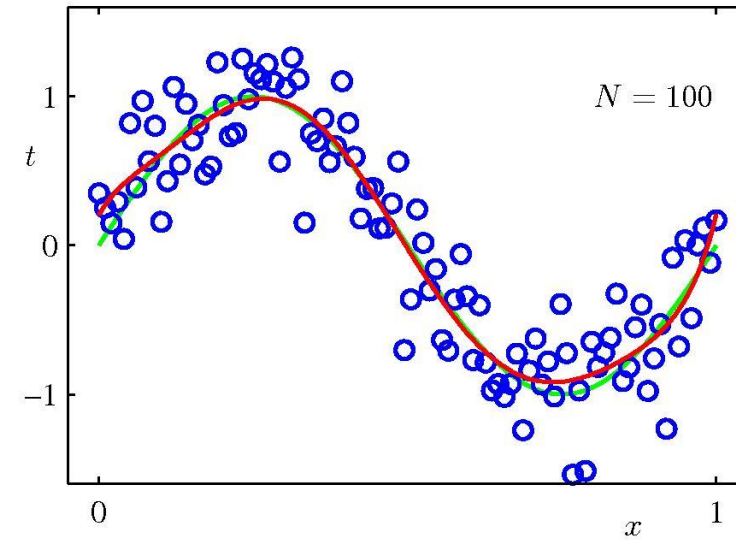
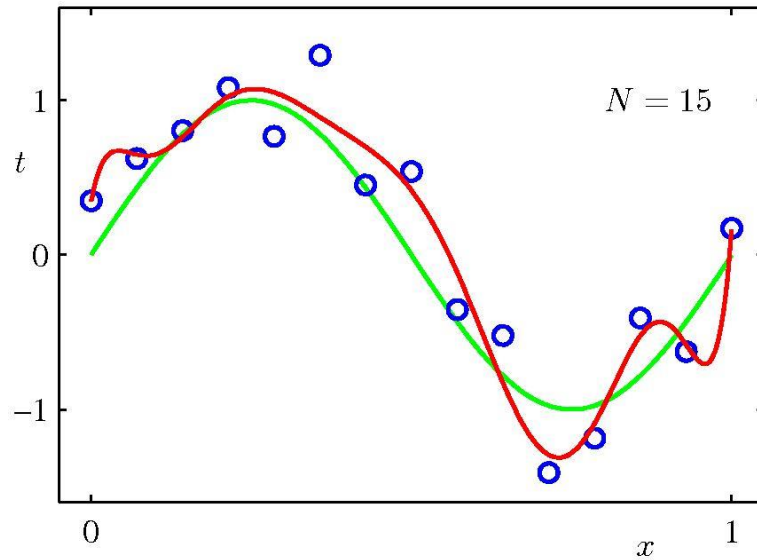
$$y(x, w) + \text{large error} \Leftarrow w_0 + \sum_{j=1}^M w_j (x + \text{small noise})^j$$

The weights are becoming huge to compensate for the noise.
(a small noise on x will have much fluctuation on prediction y)



Overfitting

- Possible solutions
 - One workaround: Use more data



Overfitting

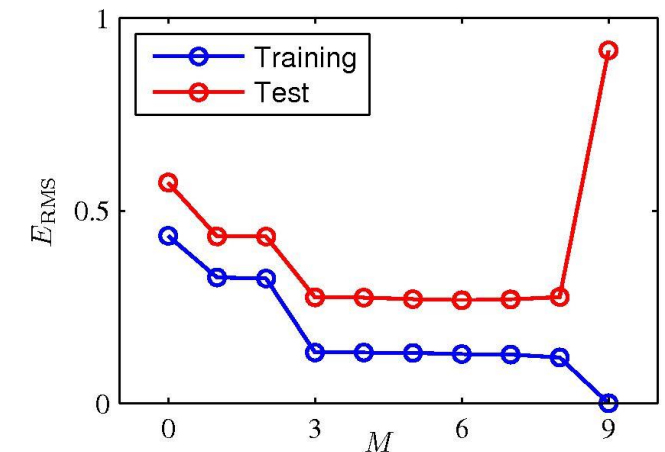
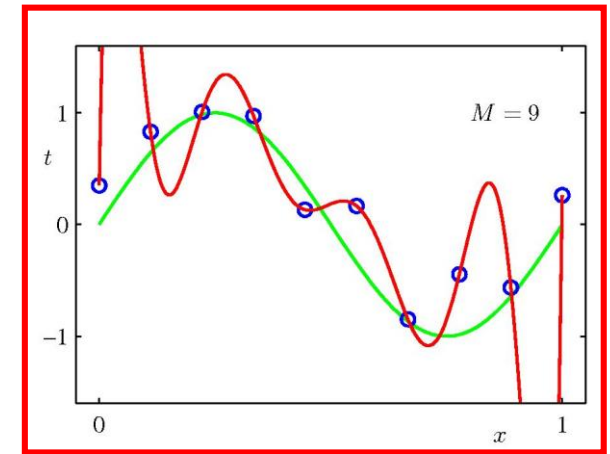
- Possible solutions
 - Second workaround: Use **regularization** (Very important!!!)
- Regularization
 - Redesign Loss function by introducing **regularization term**

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

Overfitting

- Let's look at the estimated weights for various M

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



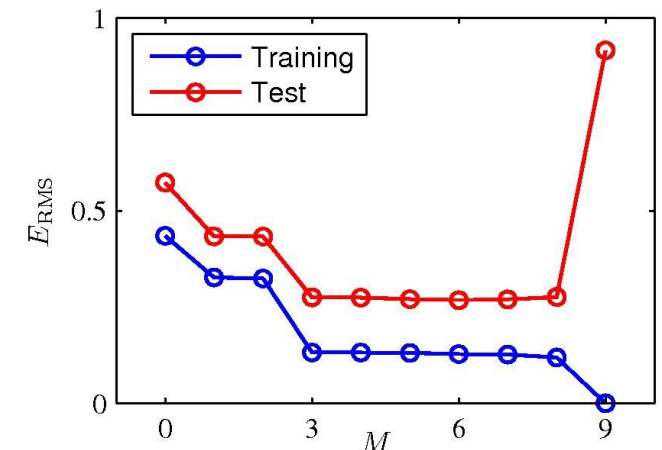
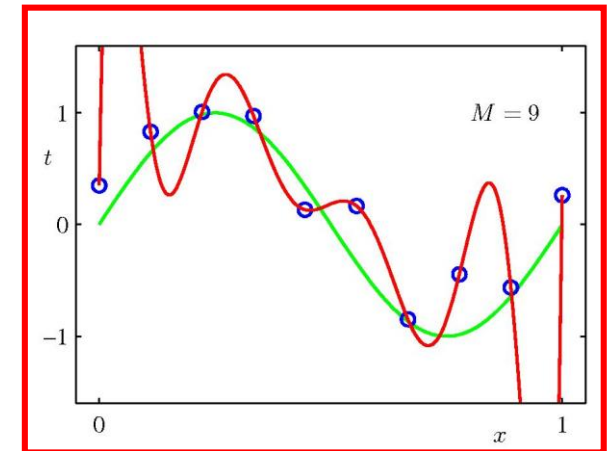
Overfitting

- Let's look at the estimated weights for various M

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

$$y(x, w) + \text{large error} \Leftarrow w_0 + \sum_{j=1}^M w_j (x + \text{small noise})^j$$

The weights are becoming huge to compensate for the noise.
(a small noise on x will have much fluctuation on prediction y)



Overfitting

- One way of dealing with this is to encourage the weights to be small. This is called regularization.

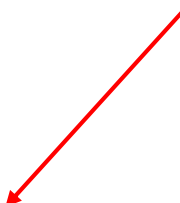
- **Standard approach**

- New loss function

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

- The penalty on the squared weights is known as ridge regression (Statistics)

Overfitting

- When w_i 's are small, prediction y will be not sensitive to small change of x
 - Smooth functions are preferred.
- 

- One way of dealing with this is to encourage the weights to be small. This is called regularization.

- **Standard approach**

- New loss function

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

- The penalty on the squared weights is known as ridge regression (Statistics)

Overfitting

- One way of dealing with this is to encourage the weights to be small. This is called regularization.

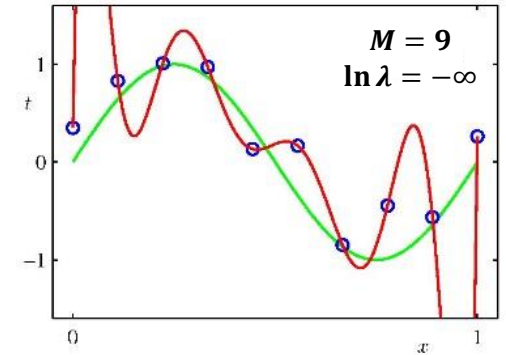
- **Standard approach**

- New loss function

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

- The penalty on the squared weights is known as ridge regression (Statistics)

- When w_i 's are small, prediction y will be not sensitive to small change of x
- Smooth functions are preferred.



Overfitting

- One way of dealing with this is to encourage the weights to be small. This is called regularization.

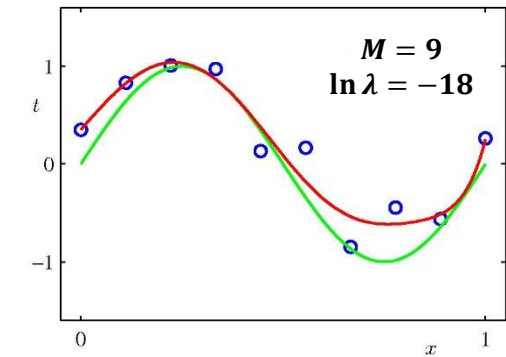
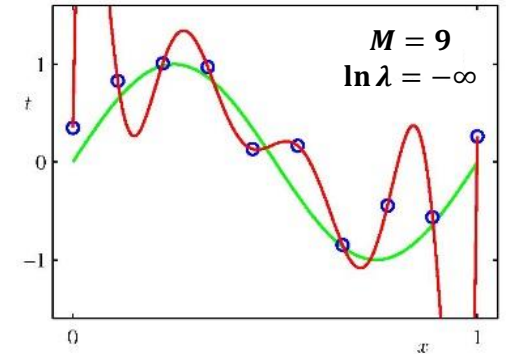
- **Standard approach**

- New loss function

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

- The penalty on the squared weights is known as **ridge regression** (Statistics)

- When w_i 's are small, prediction y will be not sensitive to small change of x
- Smooth functions are preferred.



Overfitting

- One way of dealing with this is to encourage the weights to be small. This is called regularization.

- **Standard approach**

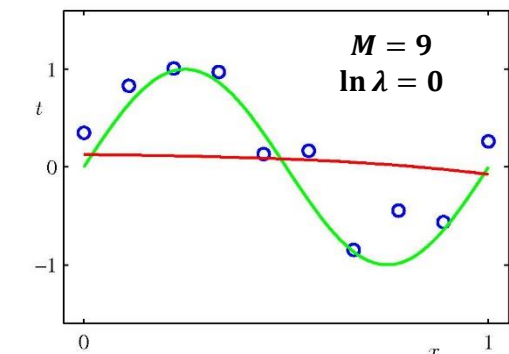
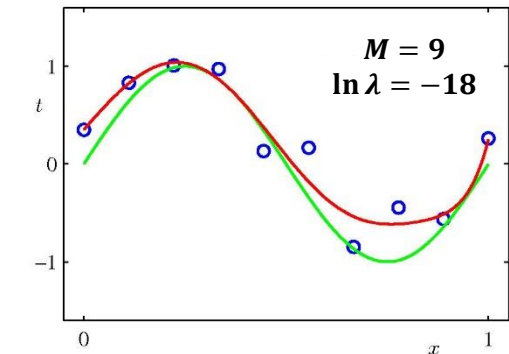
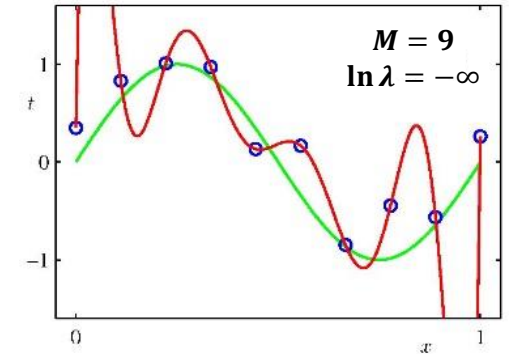
- New loss function

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

- The penalty on the squared weights is known as **ridge regression** (Statistics)

- When w_i 's are small, prediction y will be not sensitive to small change of x
- Smooth functions are preferred.

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01



Overfitting

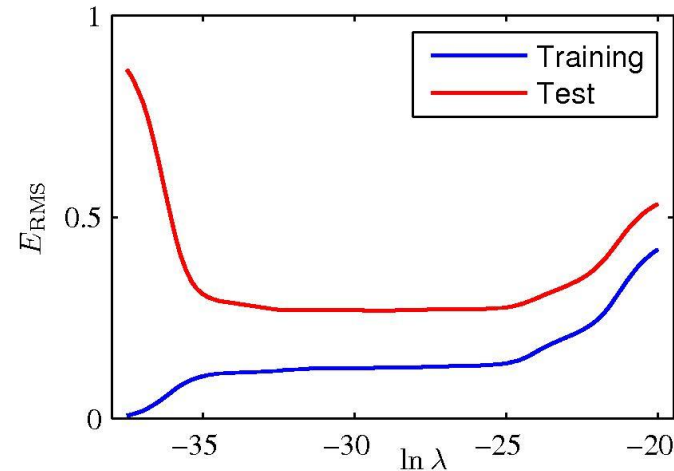
- One way of dealing with this is to encourage the weights to be small. This is called **regularization**.

- **Standard approach**

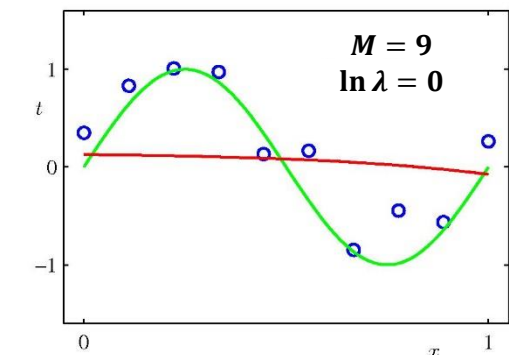
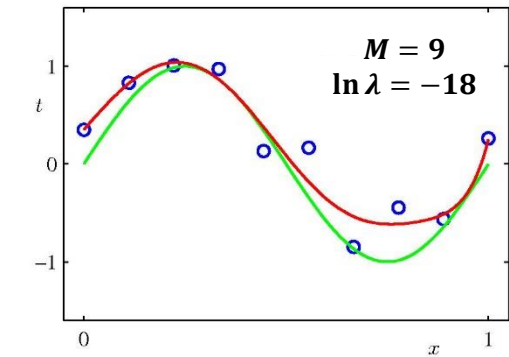
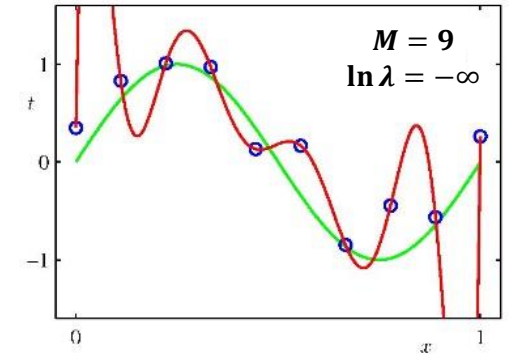
- New loss function

$$\ell(w) = \frac{1}{2N} \left[\sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2 + \lambda \sum_{i=1}^M w_i^2 \right]$$

- The penalty on the squared weights is known as ridge regression (Statistics)



However, choose value λ
carefully
(we prefer smooth function,
but don't be too smooth)



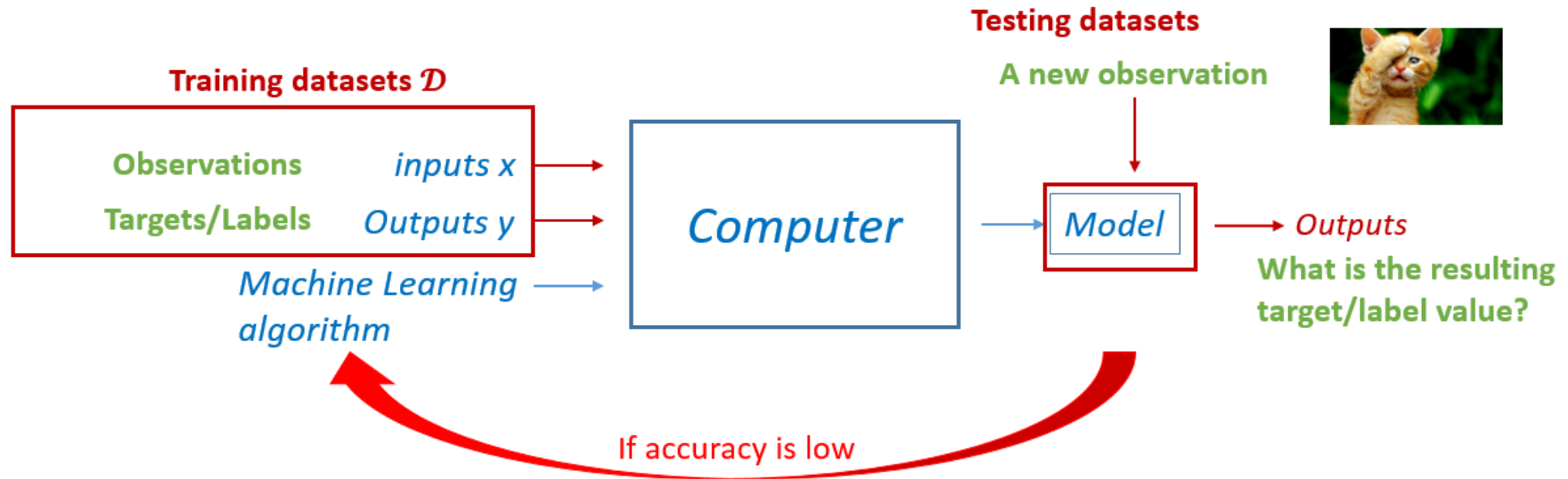
Key Concepts

Key concepts in Supervised Learning,
not just for linear models

- To find a good model
 - Loss function (measure error, or judge the fit)
 - Optimization (how to find a good fit)
 - Generalization (fit to unseen test data)
 - Regularization (avoid overfitting)
- Prediction error: validation

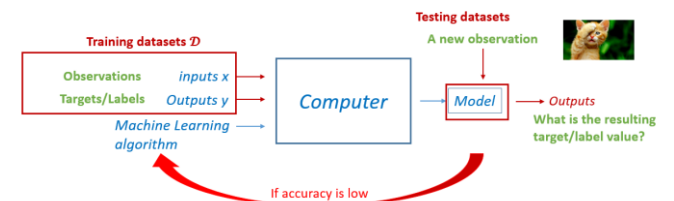
What we really care about is the
prediction error on new data

Validation



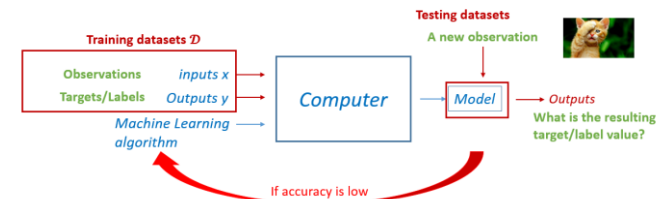
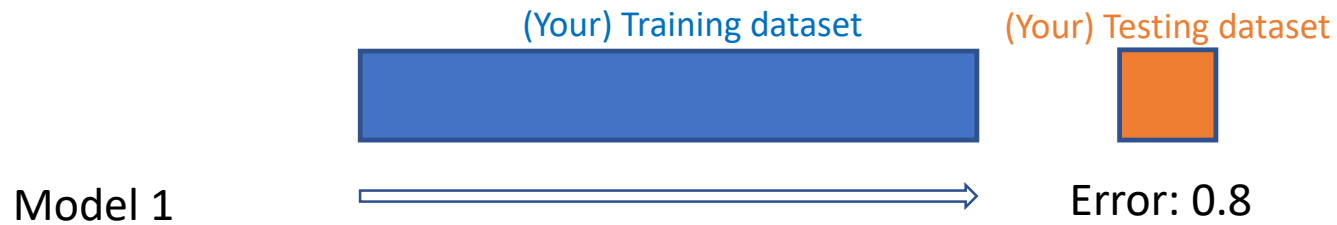
Validation

- What you should **NOT** do in Machine learning



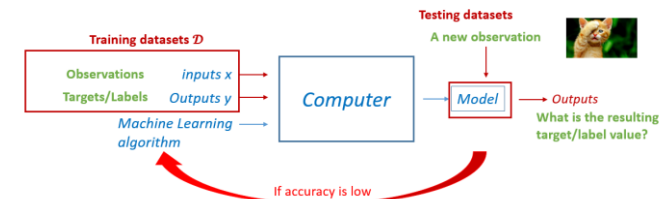
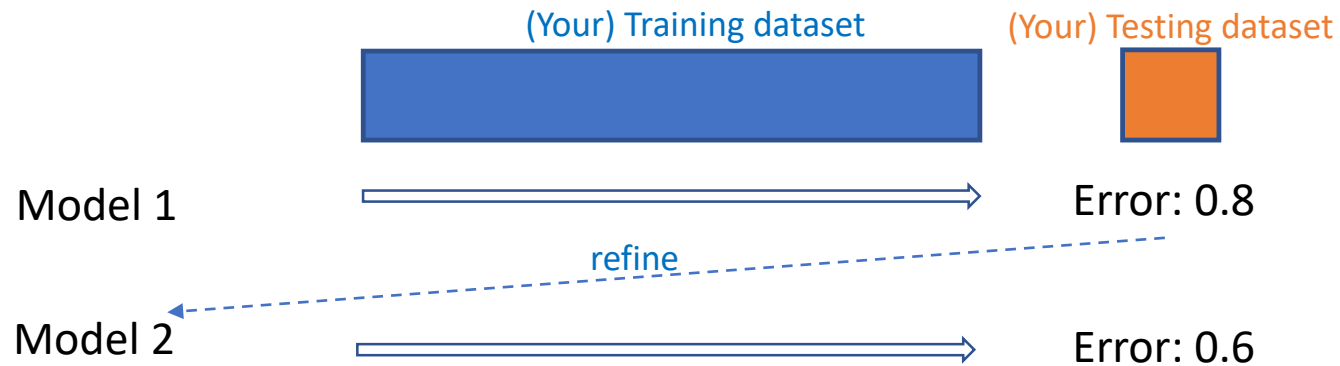
Validation

- What you should **NOT** do in Machine learning



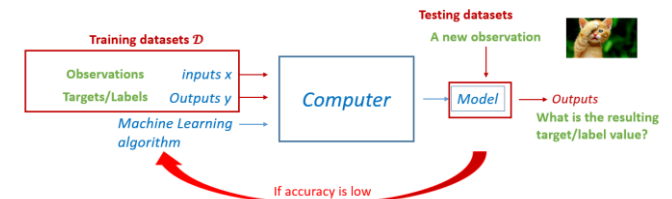
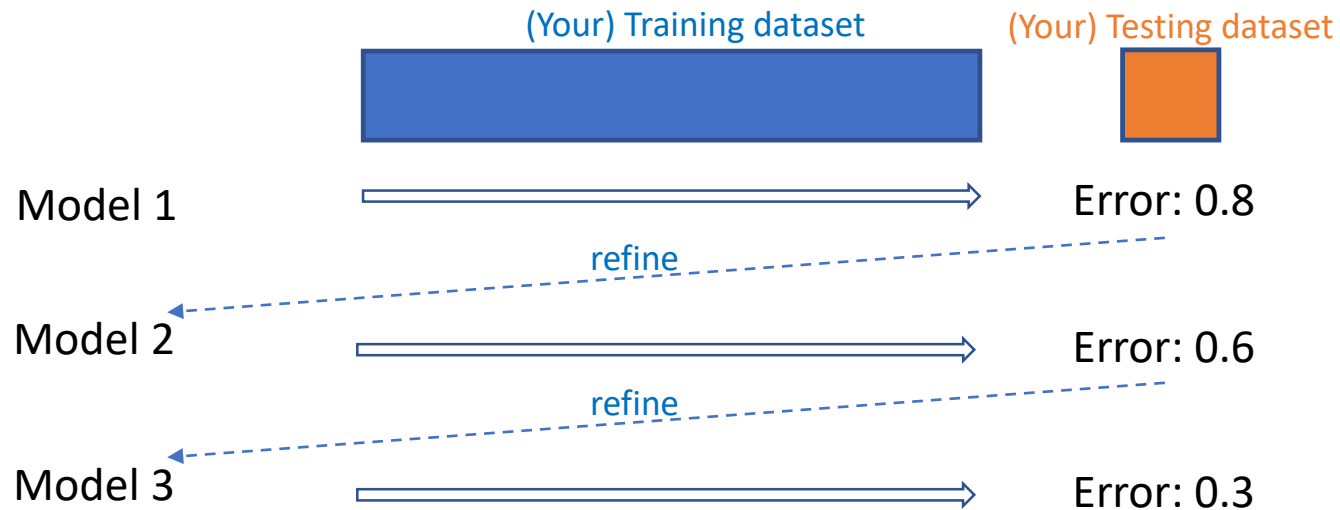
Validation

- What you should **NOT** do in Machine learning



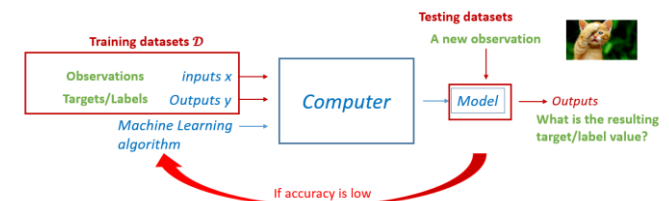
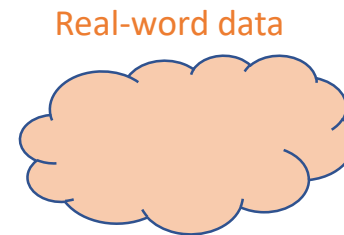
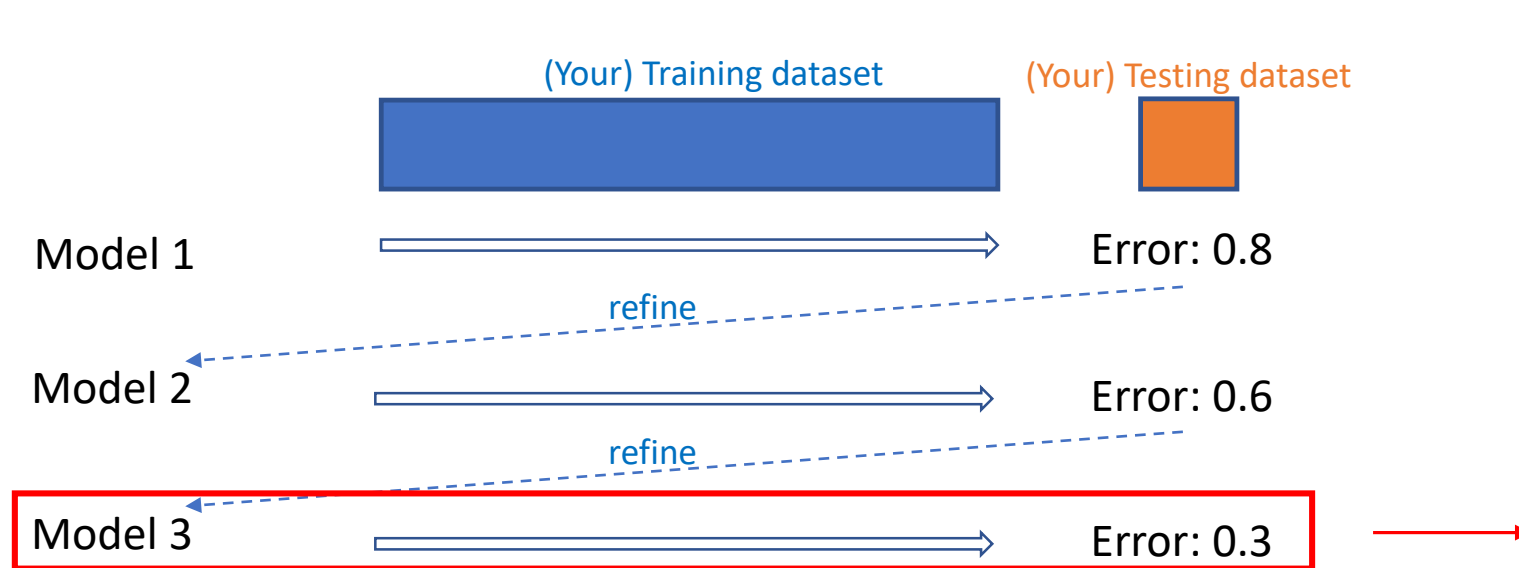
Validation

- What you should **NOT** do in Machine learning



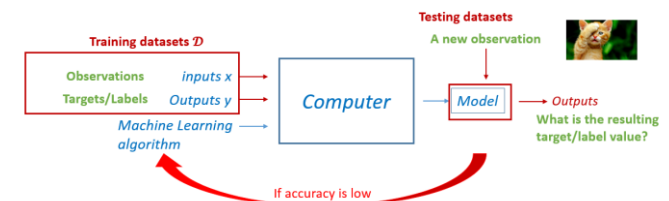
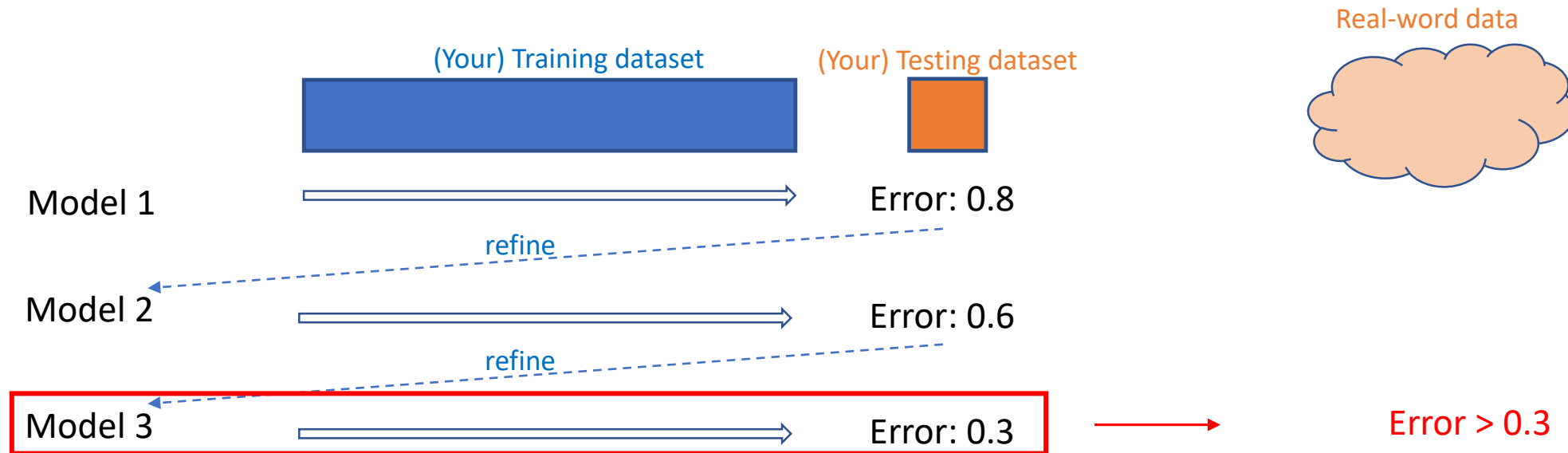
Validation

- What you should **NOT** do in Machine learning



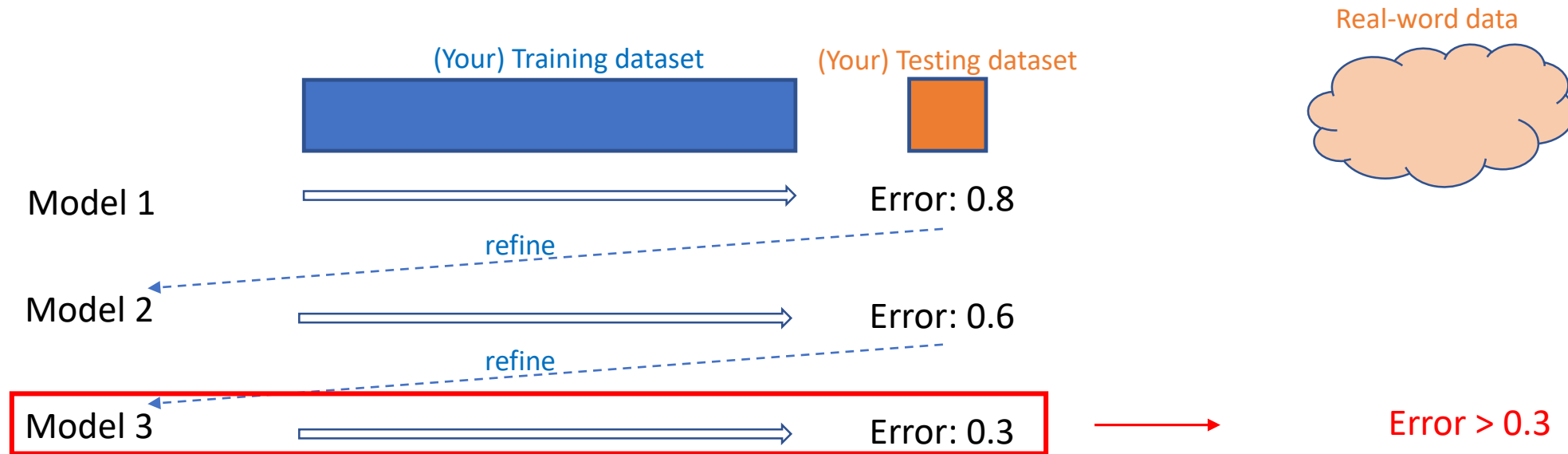
Validation

- What you should **NOT** do in Machine learning



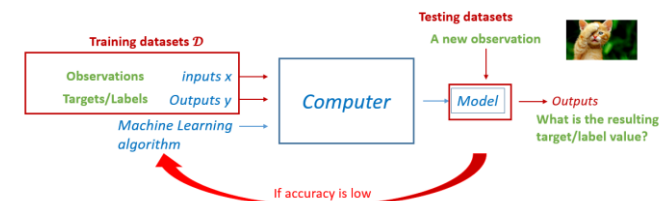
Validation

- What you should **NOT** do in Machine learning



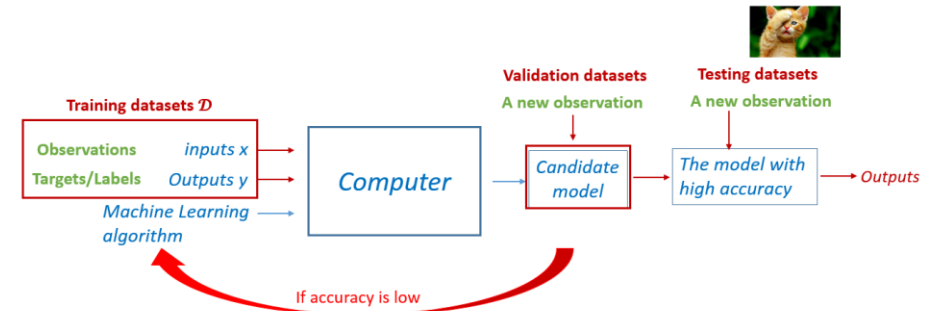
Testing data is usually used to judge the goodness of a fully-trained model. But here, it is used to refine (tune) the model.

- The testing data here is not independent of training process and may not reflect the prediction on real world data.



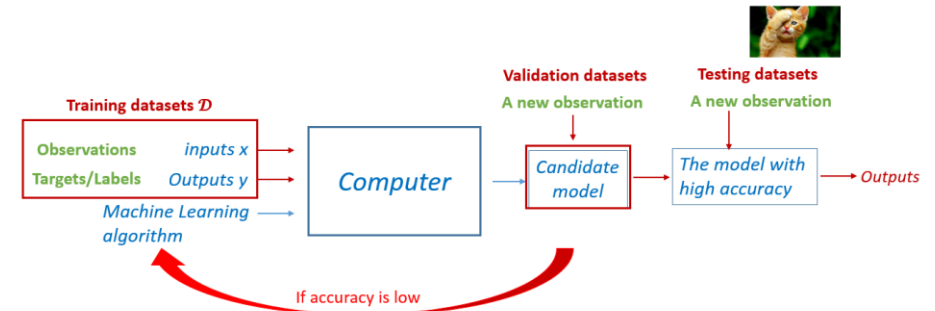
Validation

- What you **can do** in Machine learning



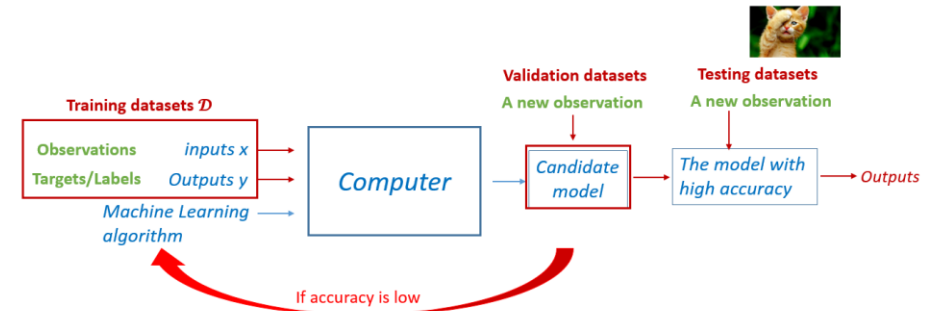
Validation

- What you **can do** in Machine learning



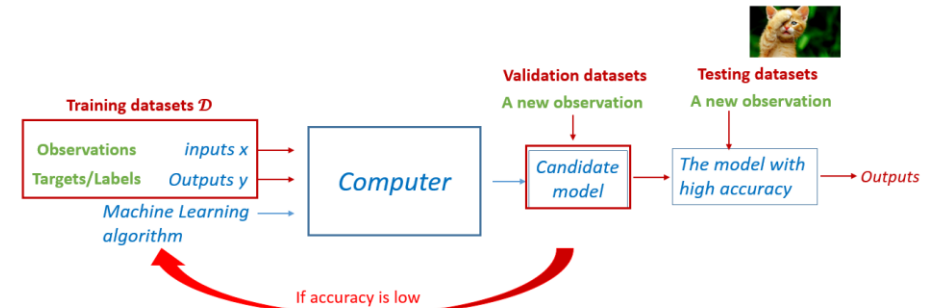
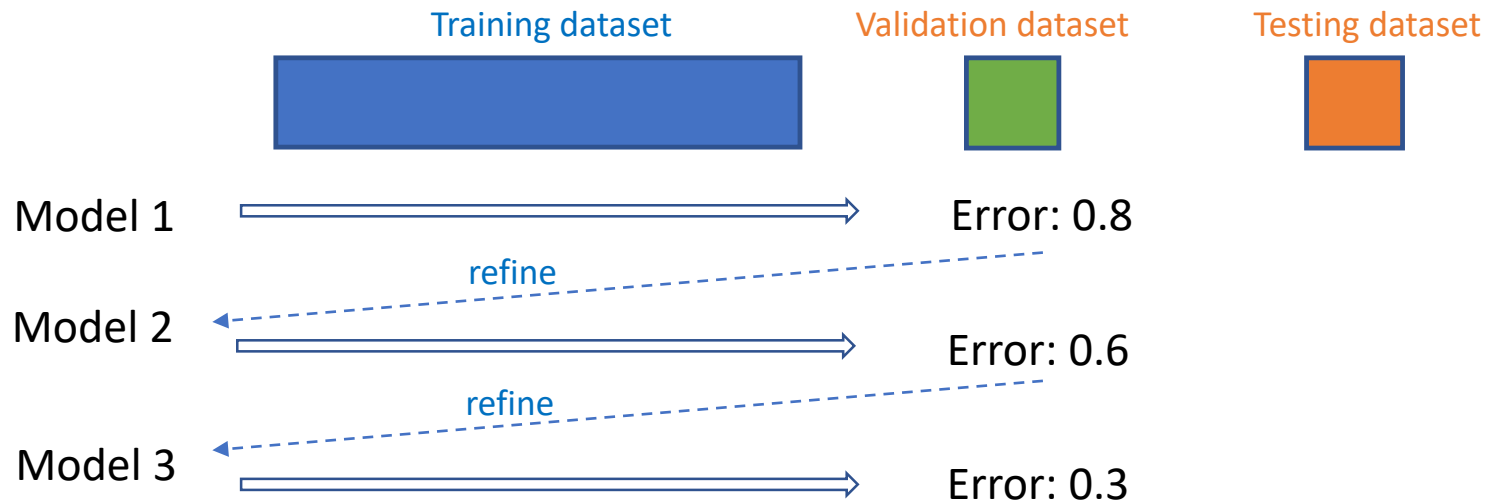
Validation

- What you **can do** in Machine learning



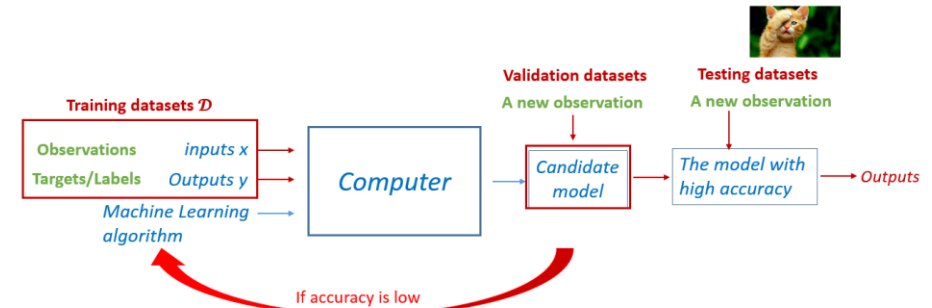
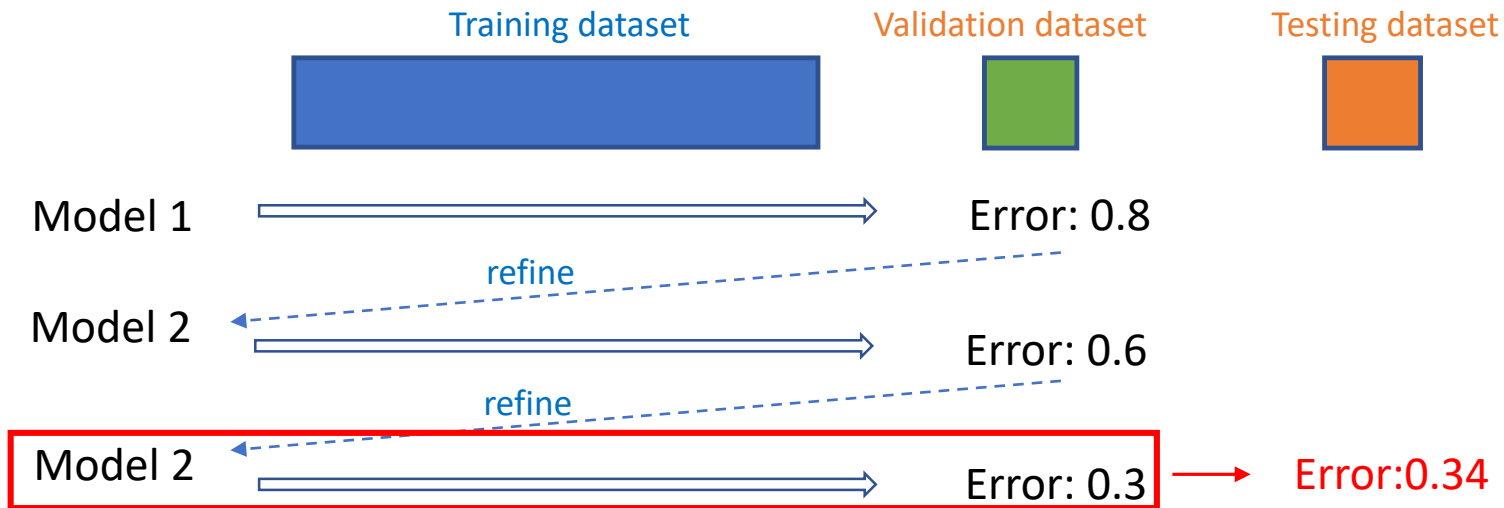
Validation

- What you **can do** in Machine learning



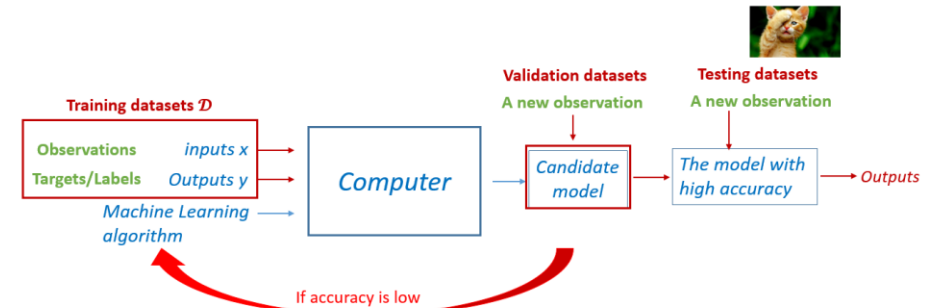
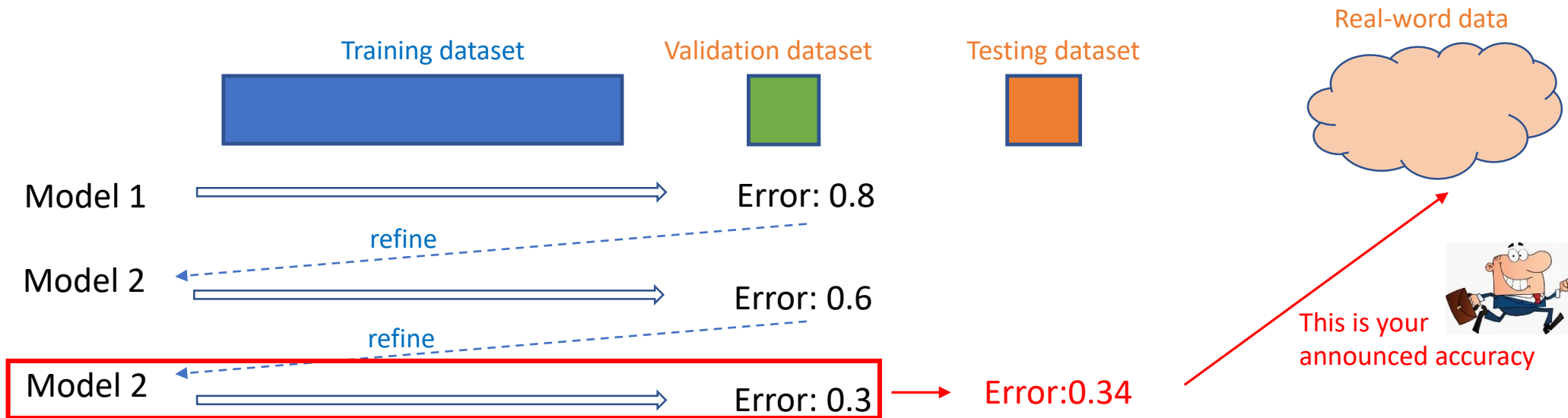
Validation

- What you **can do** in Machine learning



Validation

- What you **can do** in Machine learning



Validation

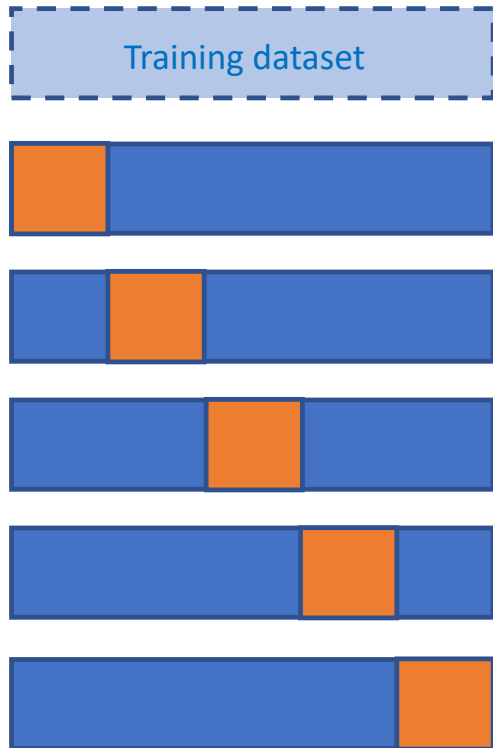
- In many applications, the supply of data for training and testing may be limited. We can use **Cross-Validation (CV)**.



Training dataset

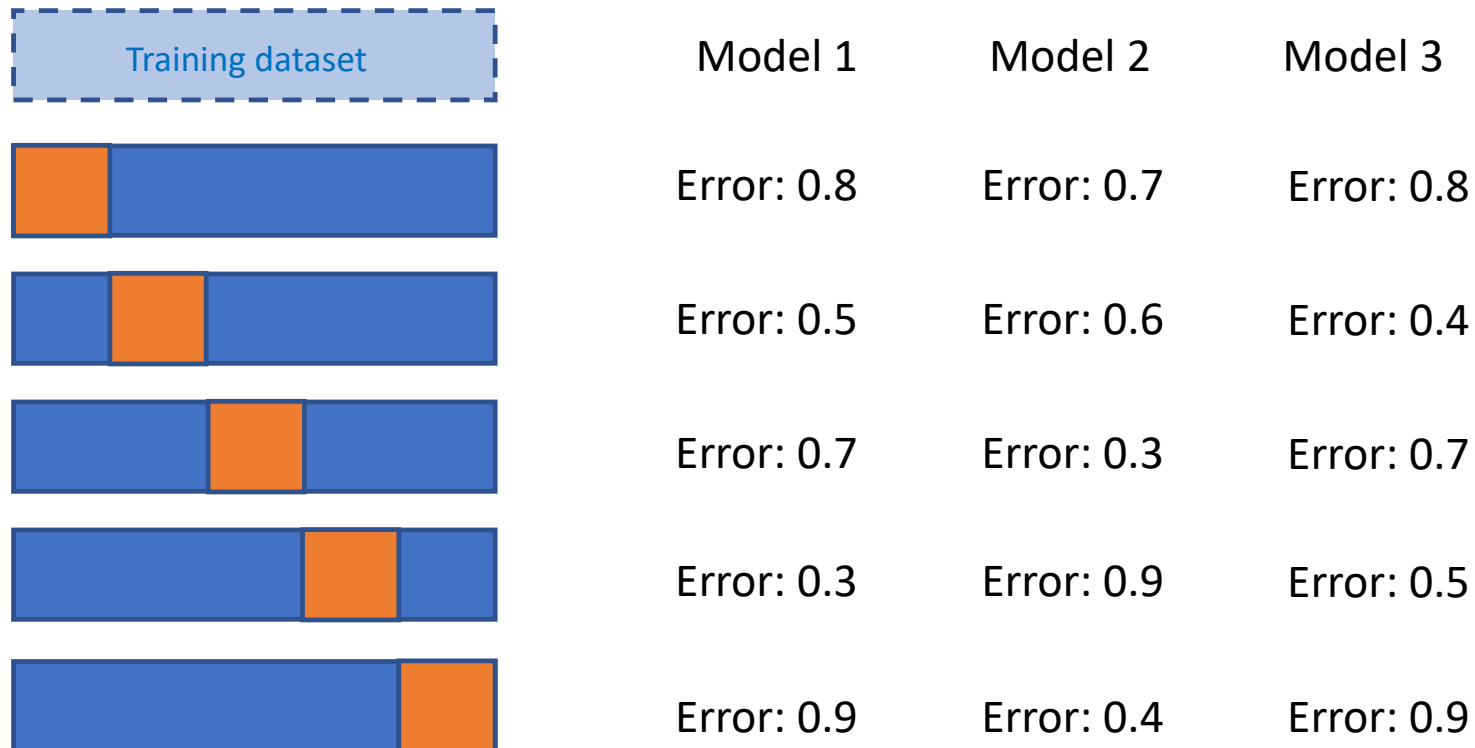
Validation

- In many applications, the supply of data for training and testing may be limited. We can use **Cross-Validation (CV)**.



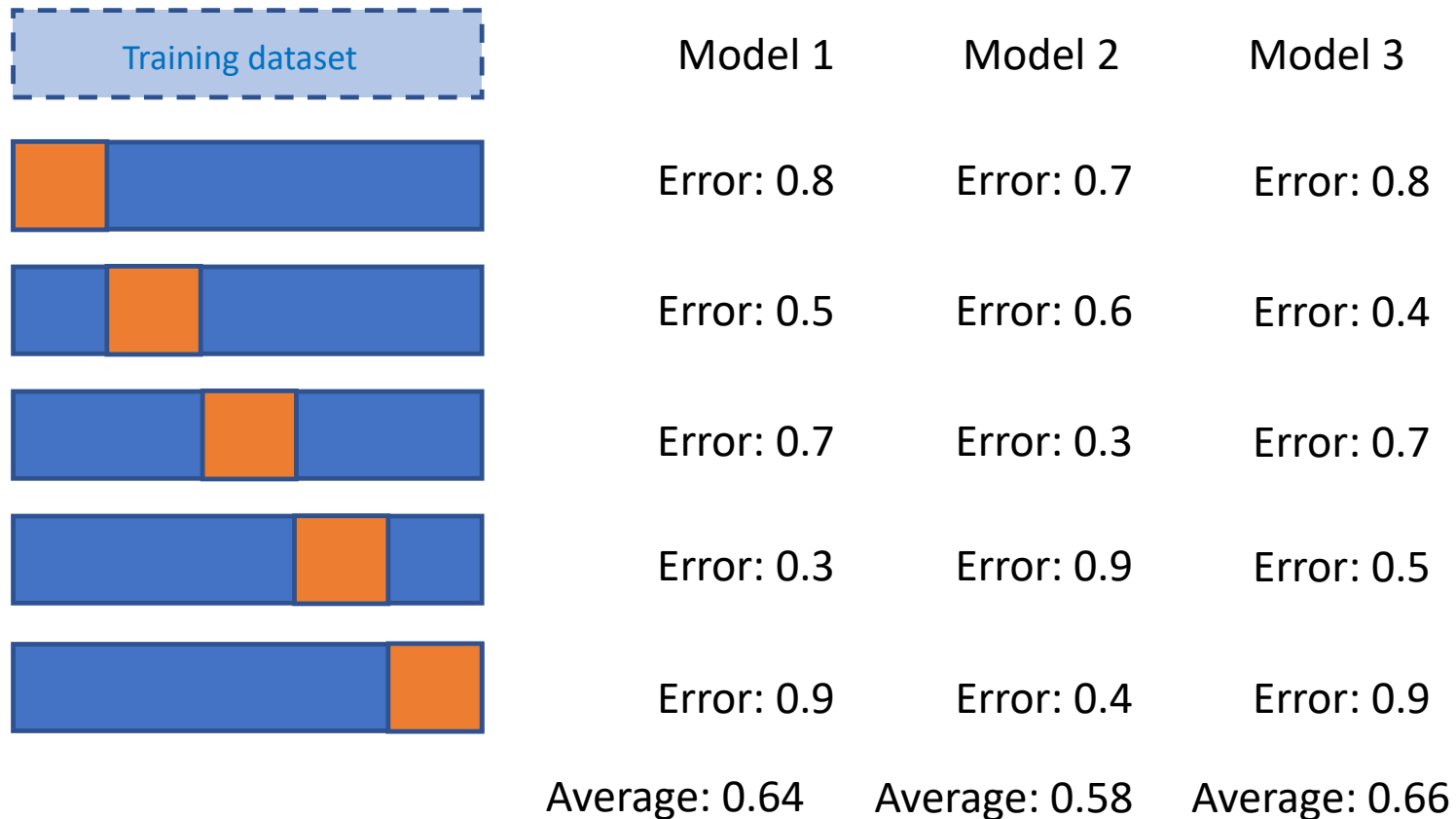
Validation

- In many applications, the supply of data for training and testing may be limited. We can use **Cross-Validation (CV)**.



Validation

- In many applications, the supply of data for training and testing may be limited. We can use **Cross-Validation (CV)**.



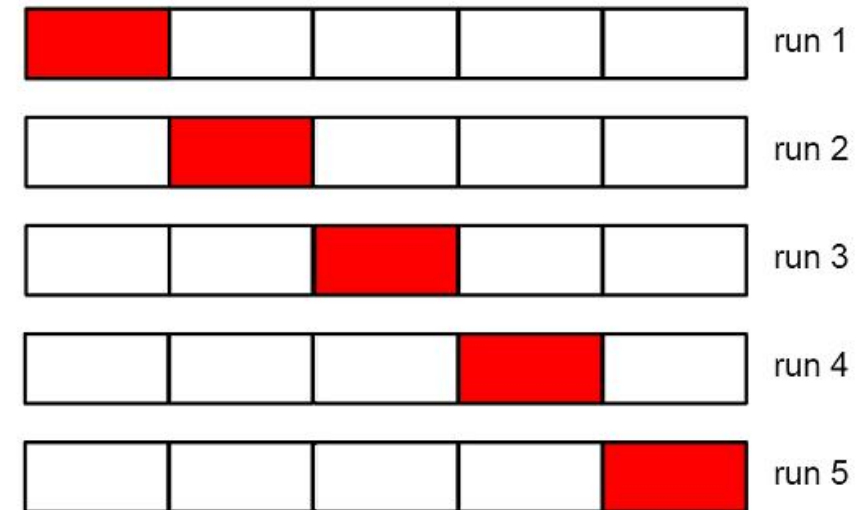
Validation

- In many applications, the supply of data for training and testing may be limited. We can use **Cross-Validation (CV)**.



Validation

- Cross Validation
 - We split the training data into K **folds**; then, for each fold $k \in \{1, \dots, K\}$, we train on all the folds but the k 'th, and test on the k 'th, in a round-robin fashion.
 - We then compute the error averaged over all the folds, and use this as a proxy for the test error.
(Note that each point gets predicted only once, although it will be used for training $K-1$ times.)
- It is common to use $K = 5$; this is called 5-fold CV.



Reading

- LFD: Chapter 4;
- Pattern Recognition and Machine Learning: 1.1 (optional)