# Overfitting and Generalization

# Best Model Selection

- **We can build various models**
  - K-NN: We may choose various values for k
  - Additive Linear Model: We may user polynomials of various order

- **But.. I want to build the best model**

# Best Model Selection

- **Various models for one input**

$$f(x) = w_0$$

$$f(x) = w_1 x + w_0$$

$$f(x) = w_2 x^2 + w_1 x + w_0$$

$$f(x) = w_3 x^3 + w_2 x^2 + w_1 x + w_0$$

...

Model Complexity increases

# Best Model Selection

- **Various Models for two inputs**

$$f(x_1, x_2) = w_0$$

$$f(x_1, x_2) = w_2 x_2 + w_1 x_1 + w_0$$

$$f(x_1, x_2) = w_5 x_2^2 + w_4 x_1^2 + w_3 x_1 x_2 + w_2 x_2 + w_1 x_1 + w_0$$

$$f(x_1, x_2) = w_9 x_2^3 + w_8 x_1^3 + w_7 x_2^2 x_1 + w_6 x_2 x_1^2 + w_5 x_2^2 + w_4 x_1^2 + w_3 x_1 x_2 + w_2 x_2 + w_1 x_1 + w_0$$
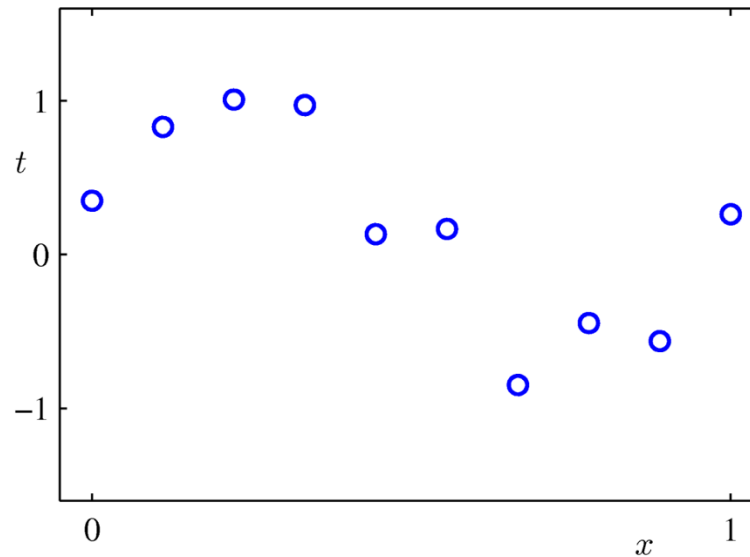
...

Model Complexity increases

증가한다!

# Best Model Selection

- **Example: Which model will be best for the data?**
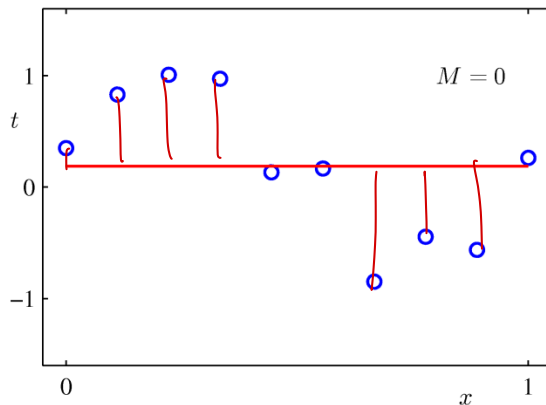
  *best의 의미?*

  – Zeroth order polynomial? First order polynomial?, or ..
  – Hmm.. Why don't we try all the possible models

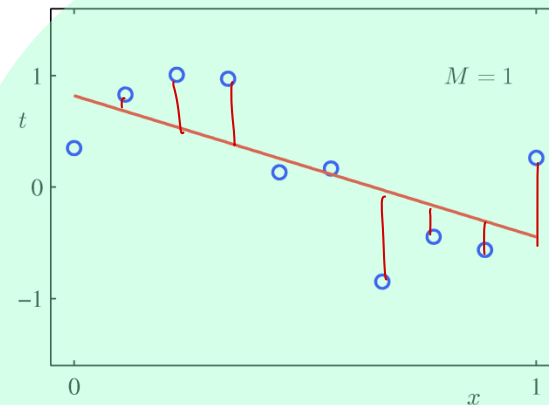## Training set

# Overfitting vs Generalization

- **Which model will be best for the data?**
  - The model which has the least error as much as possible
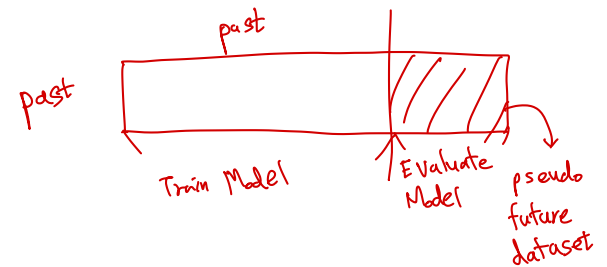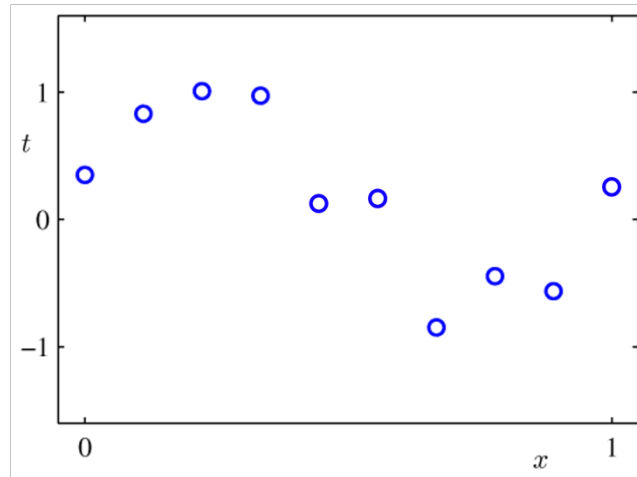


0<sup>th</sup> order polynomial regression

1<sup>st</sup> order polynomial regression
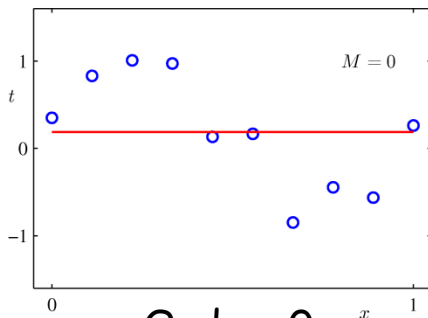
This is better
because it has less error

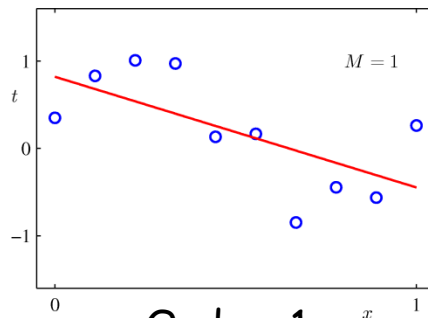# Overfitting vs Generalization

Training Data

future dataset (unknown) → Test Data

past dataset

Build Model

Evaluate model

past

## ■ Which model will be best for the data?

past

past

Train Model    Evaluate Model    pseudo future dataset

**Training Error**

No Error!

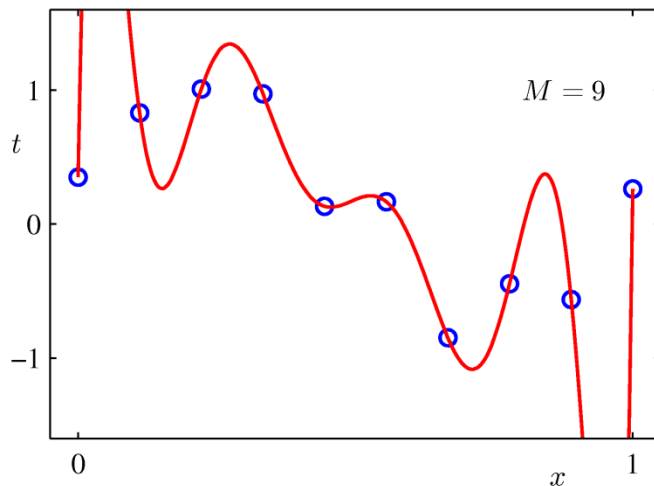| Very Large | Large | Modest | Zero |
|---|---|---|---|
| $M = 0$ | $M = 1$ | $M = 3$ | $M = 9$ |
| Order=0 | Order=1 | Order=3 | Order=9 |

# Overfitting vs Generalization

- **Which model will be best for the data?**
  - What about this?



9th order polynomial regression

*training*

  - This may be the BEST because the error is ZERO!!

Do you agree with this?

# Overfitting vs Generalization

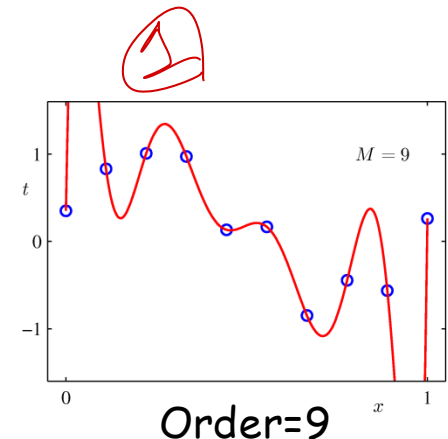- **What is the purpose of Machine Learning?**
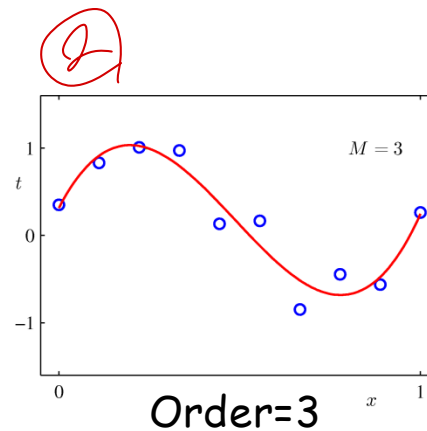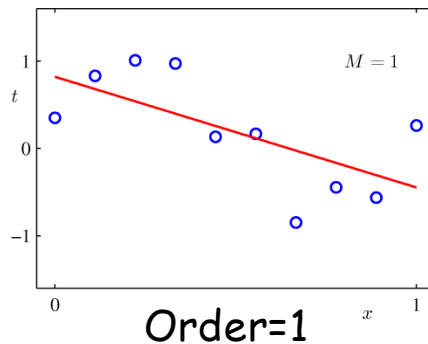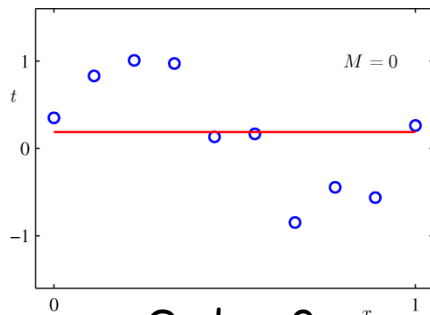
choice ① ✗

Learning the given data as exactly as possible

vs

choice ② ✓

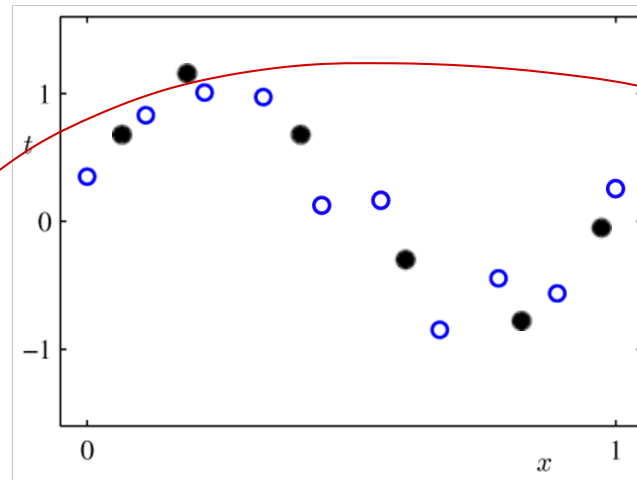Predict the unknown data as exactly as possible based on the given data

= ML의 목적



Order=0

Order=1

Order=3

Order=9

# Overfitting vs Generalization

- **As the complexity of model increases,**

New samples

Test Error

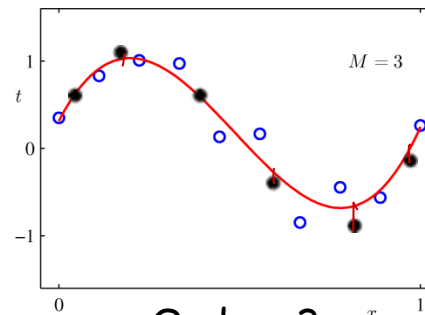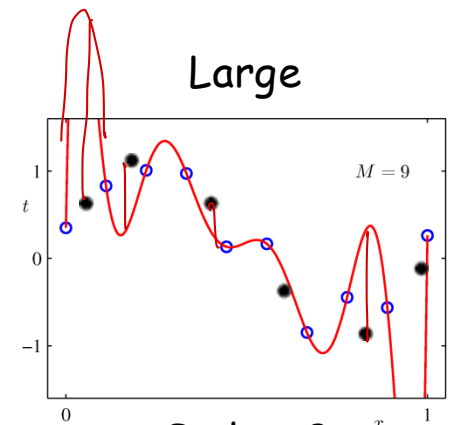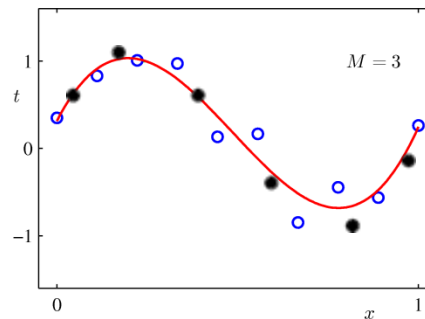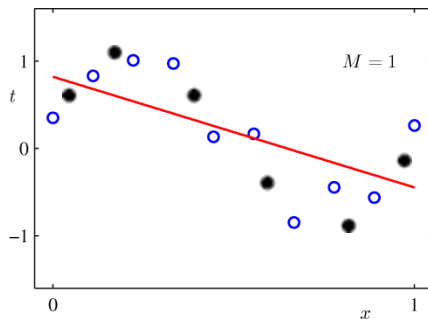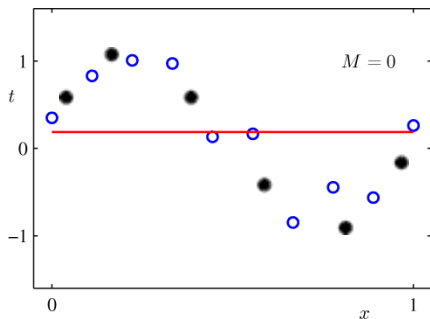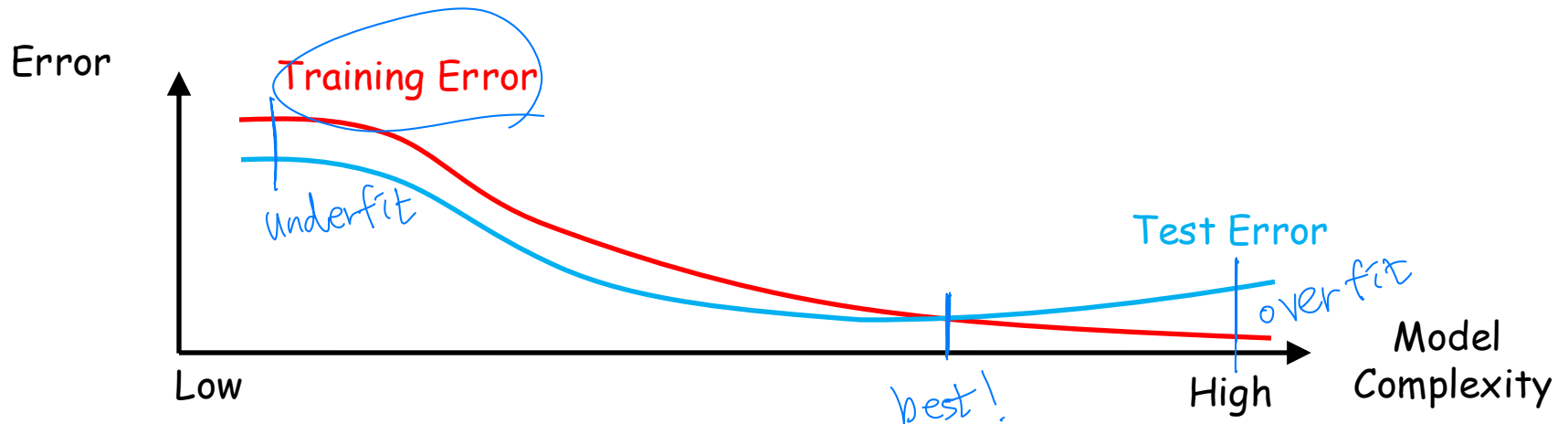| Very Large | Large | Small | Large |
|---|---|---|---|
| Order=0 | Order=1 | Order=3 | Order=9 |

$M = 0$  $M = 1$  $M = 3$  $M = 9$

# Overfitting vs Generalization

*=best model*

- **As the complexity of model increases,**
  - The model can more exactly learn the given data
  - However, the prediction accuracy (for the test set) does not necessarily increase

Test
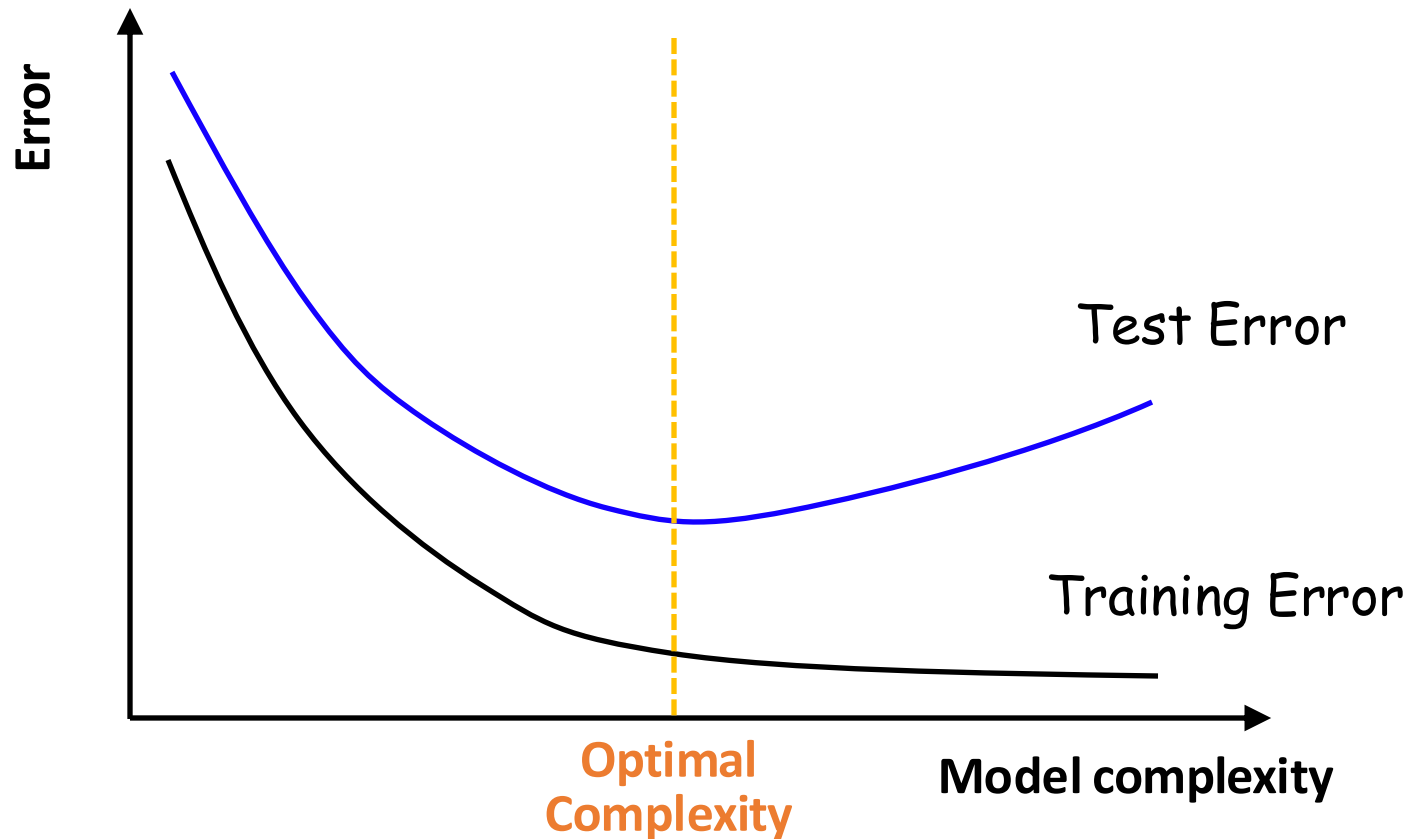Error Large

Underfit ←————————————→ Overfit

L          S          L

기저항

(= TestE, TrainE 둘다 L)

= (TrainE S,
TestE L)

best!!
Good Generalization!



Low Complexity ←————————————→ High Complexity

Training Error
Large                    Small

# Overfitting vs Generalization

- **Overfitting**



Error

Training Error

underfit

Test Error

overfit

Model Complexity

Low                    best!                    High

M = 0    M = 1    M = 3    M = 9
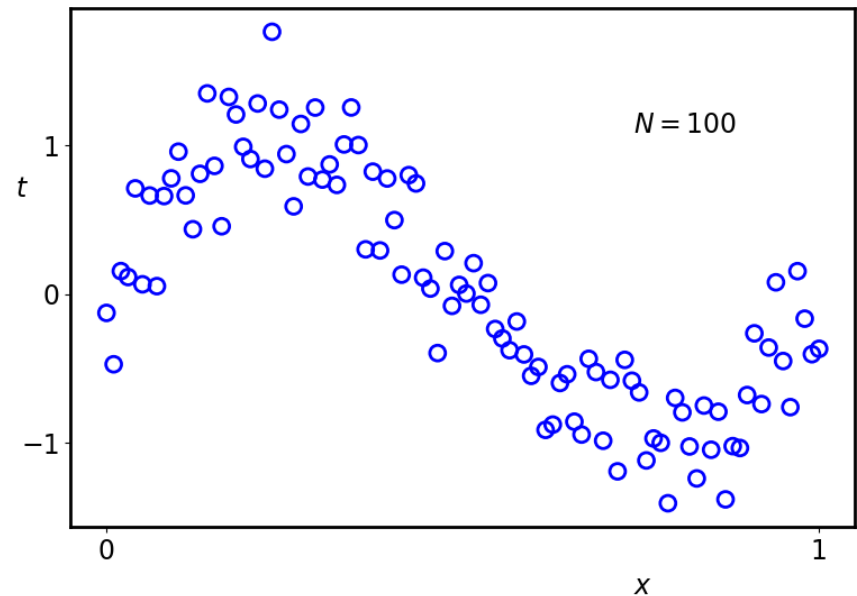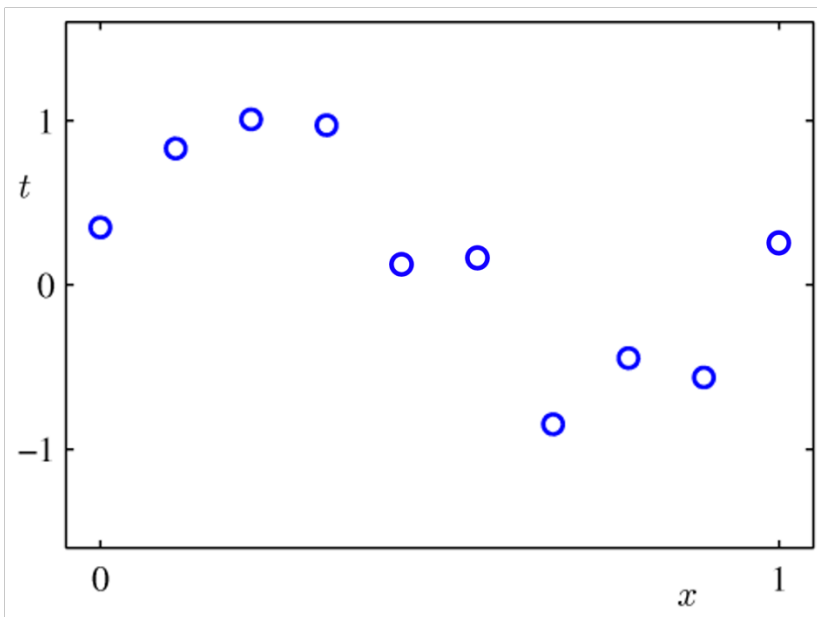
# Overfitting vs Generalization

- **Hmm.. How can I choose the optimal?**
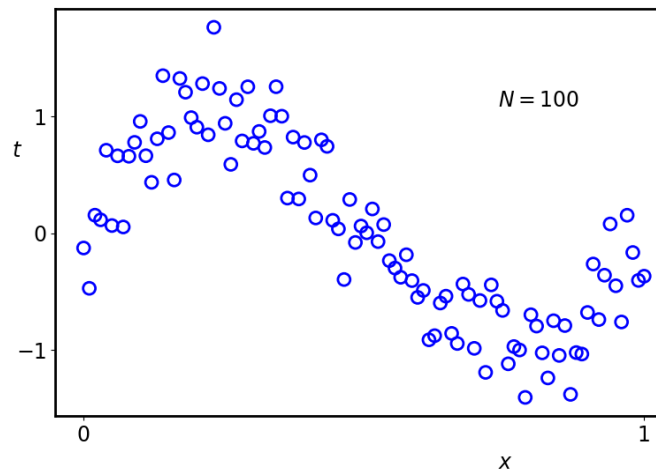
# How to Achieve Good Generalization?

- **One of easiest way is.. Collecting More Data**

많은 데이터 모아라!

# How to Achieve Good Generalization?
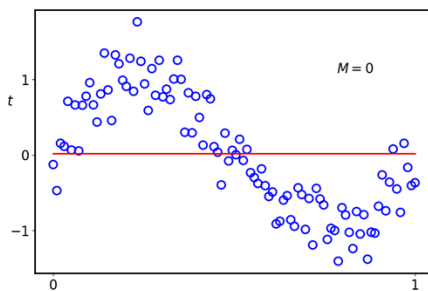
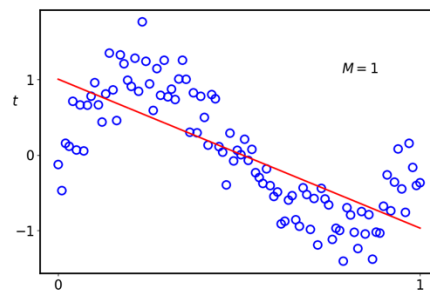- **Model building with a large dataset**



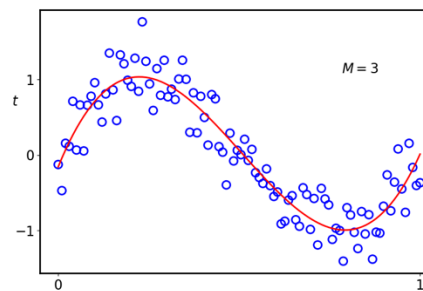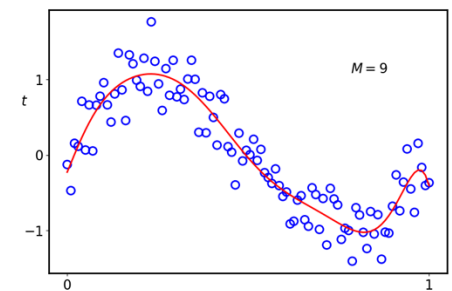**Training Error**

| Very Large | Large | Modest | Modest |
|:---:|:---:|:---:|:---:|



| Order=0 | Order=1 | Order=3 | Order=9 |
|:---:|:---:|:---:|:---:|

# How to Achieve Good Generalization?

- **Model building with a large dataset**



$N = 100$

**Model is not overfitted!!**

**Test Error**

| Very Large | Large | Modest | Modest |
|:---:|:---:|:---:|:---:|



$M = 0$    $M = 1$    $M = 3$    $M = 9$

Order=0$^x$    Order=1$^x$    Order=3$^x$    Order=9$^x$

# How to Achieve Good Generalization?

- **How many samples to avoid overfitting?**
  - Let's assume that we need <u>100</u> samples per an input feature
  - If we have two features, we need $10^4$
  - If we have three features, we need $10^6$
  - If we have four features, we need $10^8$
  - If we have ten features, we need $10^{20}$

- **Though you have a substantial amount of data, you are uncertain whether the model is overfitted or not**
  - Still, it may not be sufficient...

# How to Achieve Good Generalization?

- **Build many models and choose the best**
  - Train many models
  - Evaluate them with Cross-Validation or Hold-out method  *ch.5*
  - Choose the best  *Build & Selection*

- **Use Regularization Method**
  - There are many regularization methods  *ch.6*
  - You may train a model with the regularization method