# Validation

Sewoong Oh
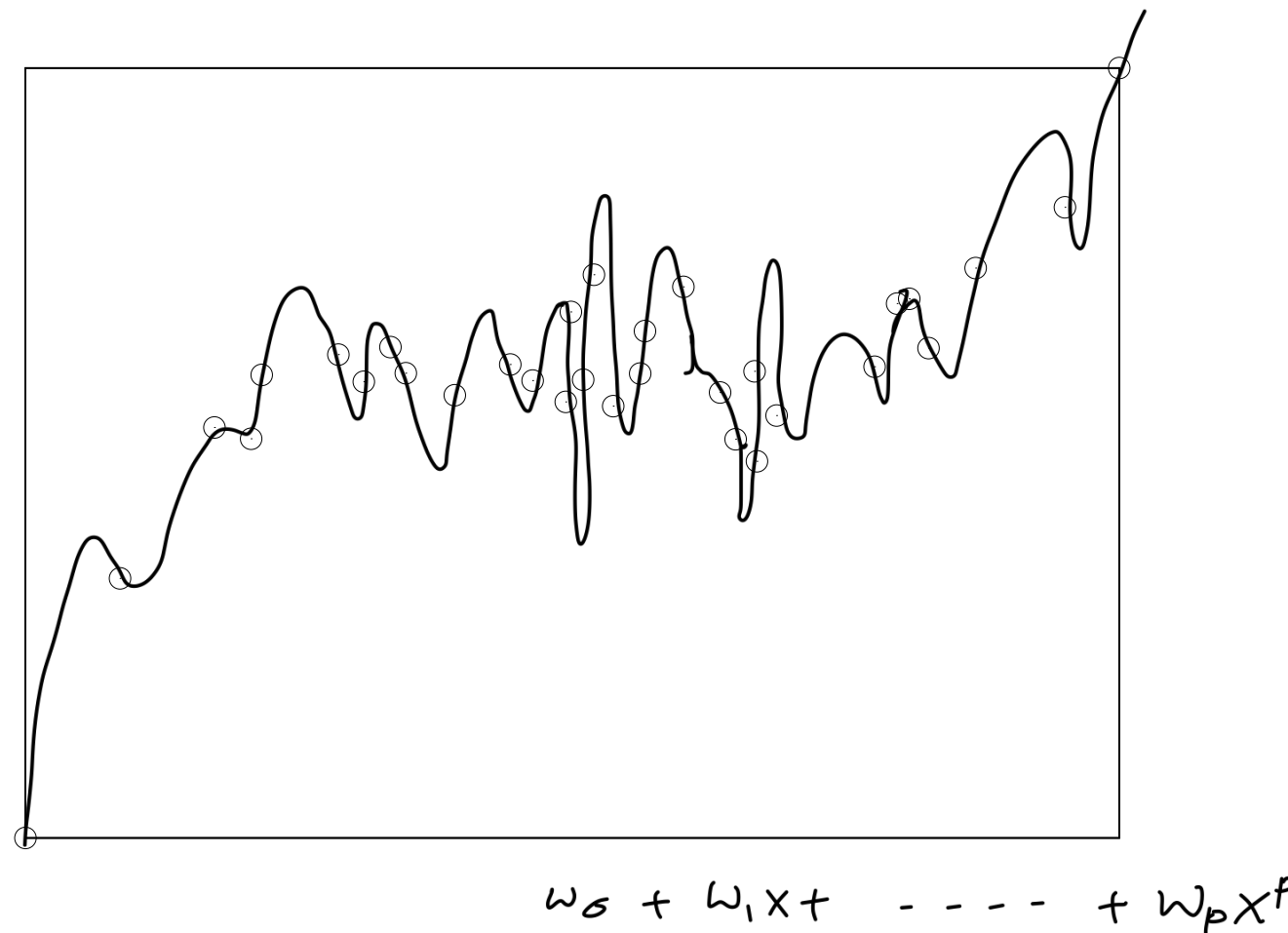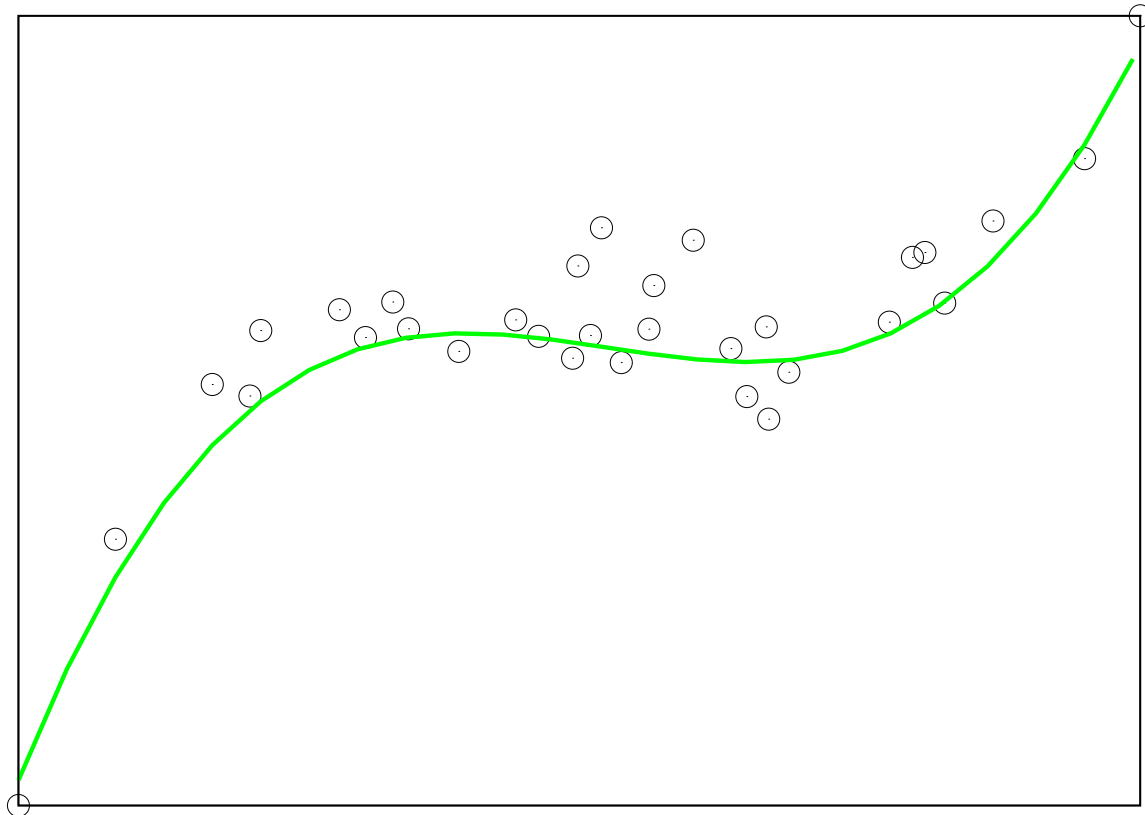
CSE/STAT 416
University of Washington

# Generalization:
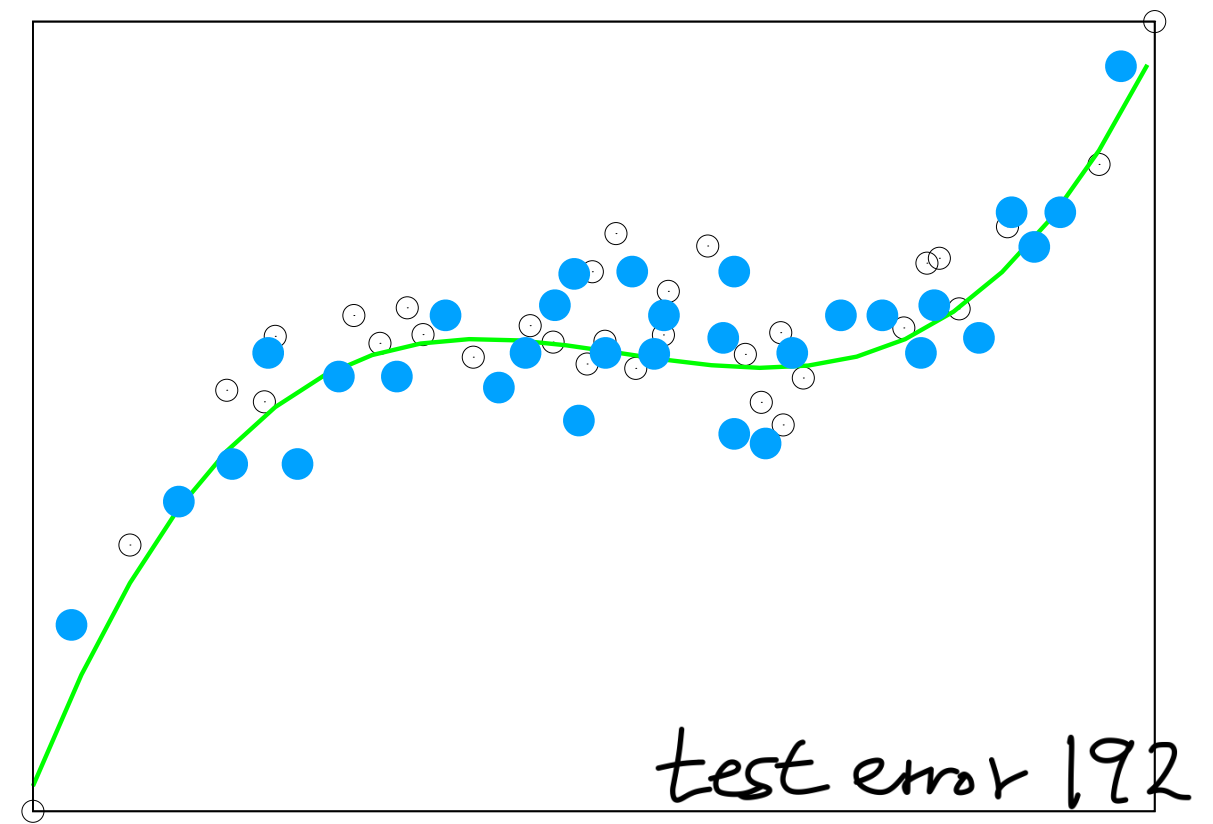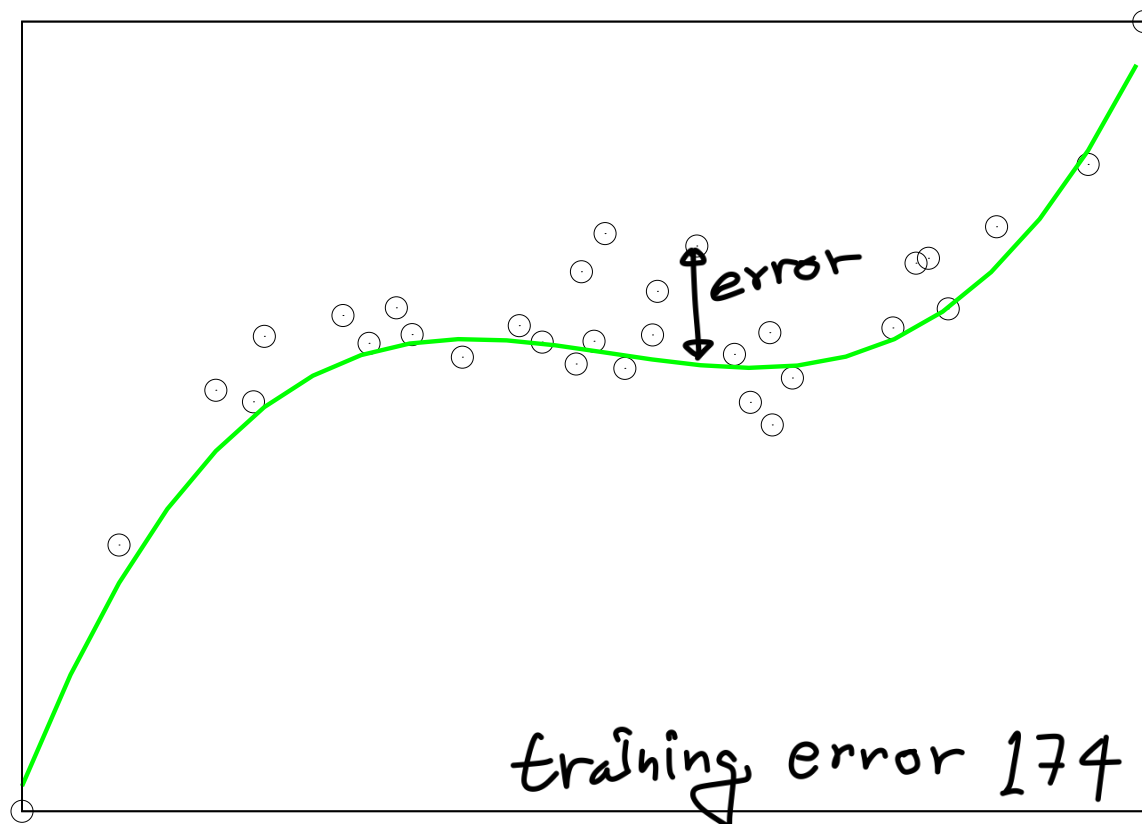# how do we validate which model is better?

# Generalization

- we say a predictor **generalizes** if it performs well on unseen data
- formal mathematical definition involves probabilistic assumptions
- first, we study practical methods for assessing generalization

$$w_0 + w_1 x + \; - - - - \; + w_p x^p$$
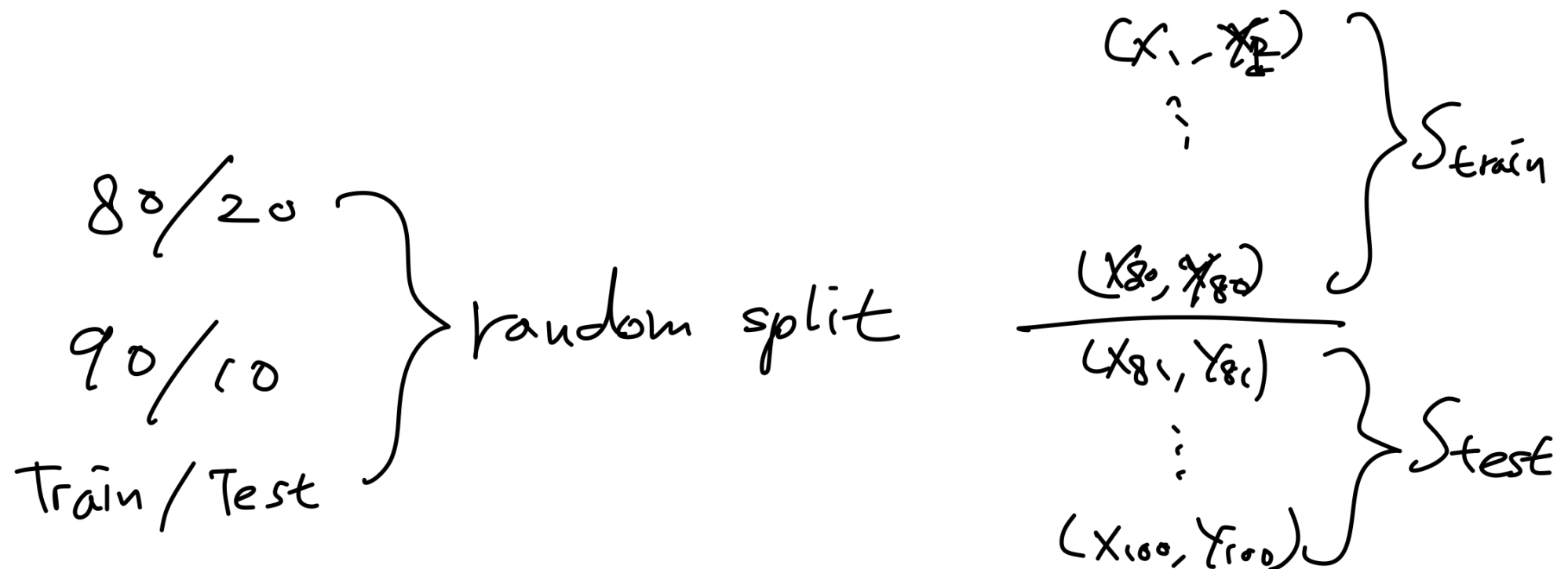
# In-sample and out-of-sample data

- the data used to construct a predictor is **training data** or **in-sample data**

- we want the predictor to work on **out-of-sample data**

- we say a predictor **fails to generalize** if it does not perform well on out-of-sample data



- **train** a cubic predictor on 32 (in-sample) black circles: MSE 174

- **predict** $y$ for 30 (out-of-sample) blue points: MSE 192

- conclude this predictor generalizes: in-sample MSE $\approx$ out-of-sample MSE

# Validation

- a way to mimic how the predictor performs on unseen data
- key idea: divide the data into two set for **training** and **testing**

- **training set** used to construct ("train") the predictor
- **test set** or **validation set** used to evaluate the predictor

- based on the assumption that test set is similar to unseen data

$$(X_1, Y_1)$$
$$\vdots$$
$$\Big\}\ S_{train}$$

$$80/20$$
$$90/10$$
$$Train / Test$$

$$\Big\}\ random\ split$$

$$(X_{80}, Y_{80})$$

$$\overline{(X_{81}, Y_{81})}$$
$$\vdots$$
$$(X_{100}, Y_{100})$$
$$\Big\}\ S_{test}$$

# Validation

- we use **training error** for optimization (or finding the model)

$$\text{MSE}_{\text{train}} \ = \ \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} (f(x_i) - y_i)^2$$
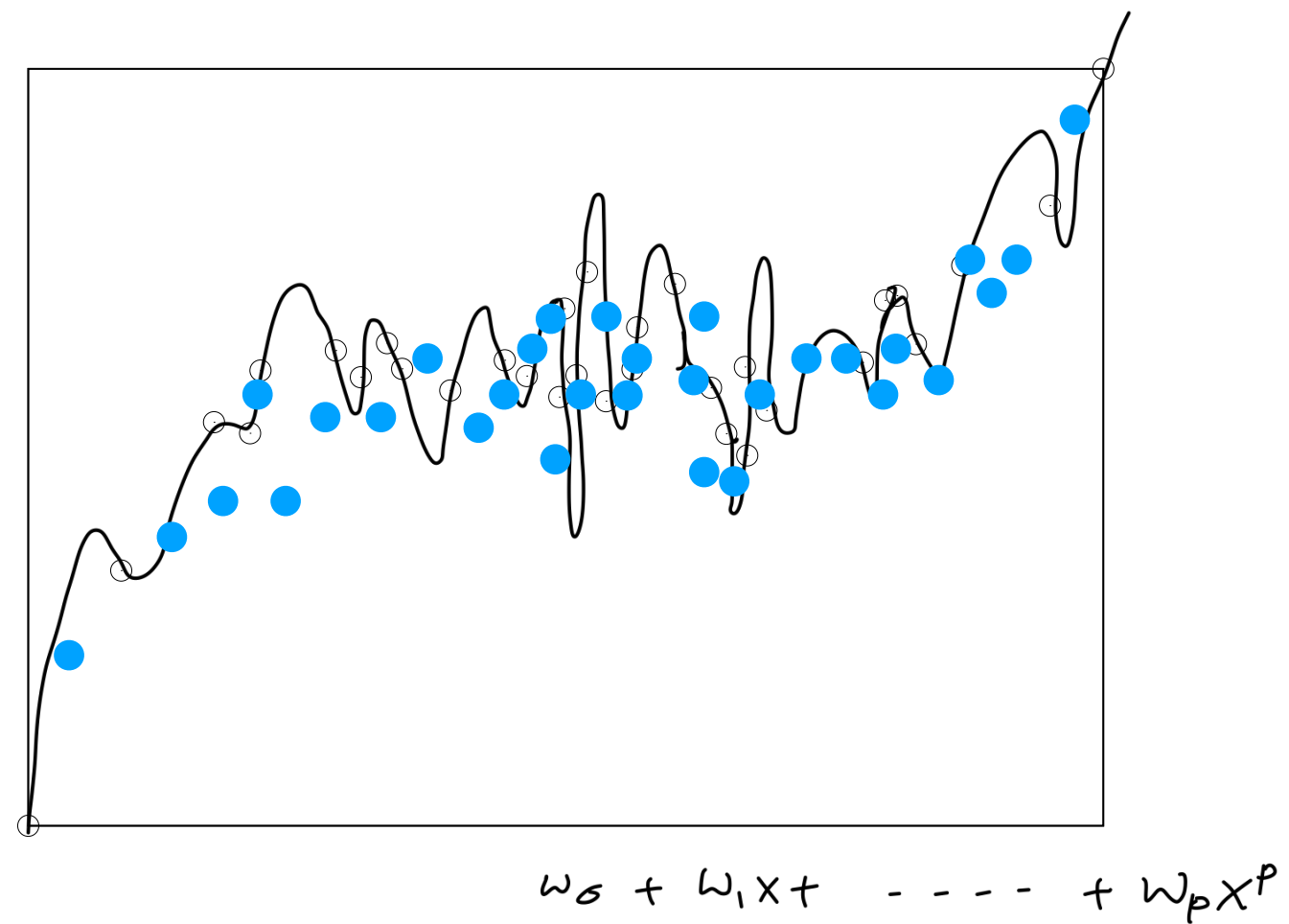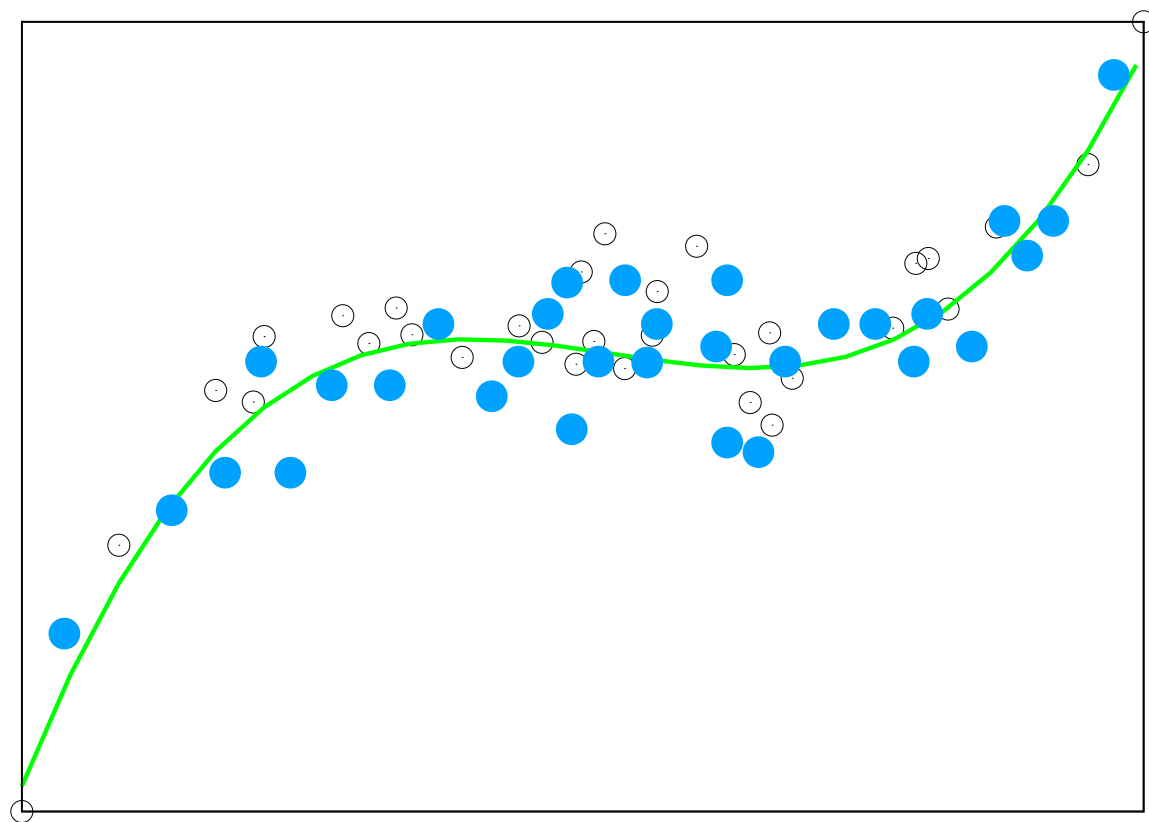
- we use **test error** for validation

$$\text{MSE}_{\text{test}} \ = \ \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} (f(x_i) - y_i)^2$$

- selecting train/test sets should be **random** (80/20 or 90/10 are common)

- we say a model or predictor is **overfit** if

$$\text{MSE}_{\text{test}} \ \gg \ \text{MSE}_{\text{train}}$$

|  | **small training error** | **large training error** |
|---|---|---|
| **small test error** | generalizes<br>perform well | possible, but lucky |
| **large test error** | fails to generalize | generalizes<br>perform bad |

$$\omega_0 + \omega_1 x + \ \text{-----} \ + \omega_p x^p$$

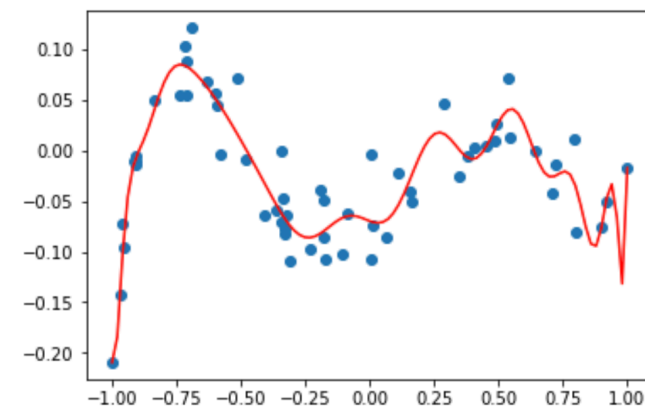- between two models, the one with smaller test error should be chosen
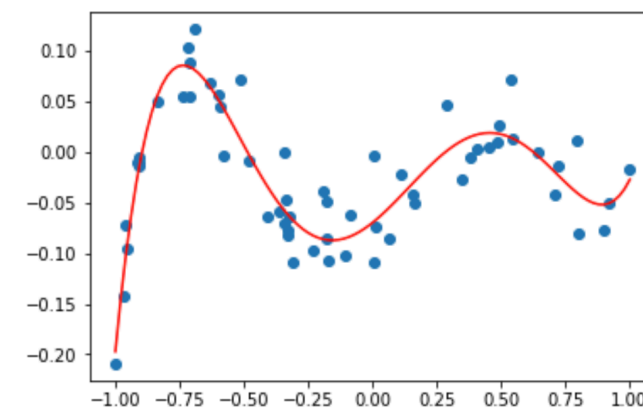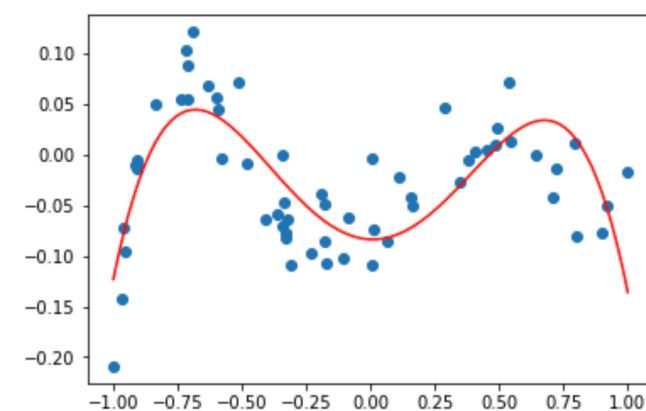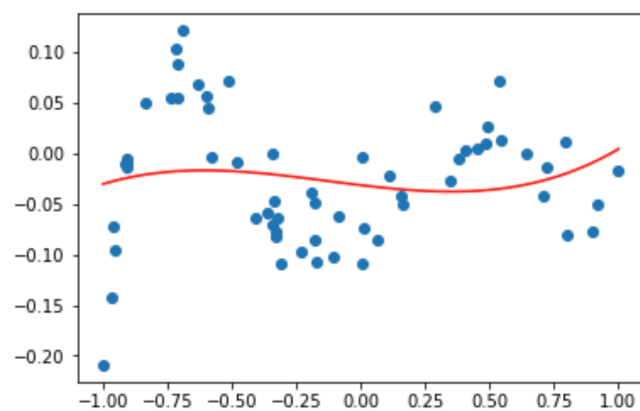
8

# Overfitting

- a model that fits the training data well but performs poorly on test data suffers from **overfitting**

- overfitting happens if we use a model with high **model complexity**

- for example, for linear regression with polynomial features

$$\hat{y} = f(x) = w_0 + w_1 x + w_2 x^2 + \cdots + w_p x^p$$

- N = 60 data points, and $p \in \{3, 4, 5, 20\}$



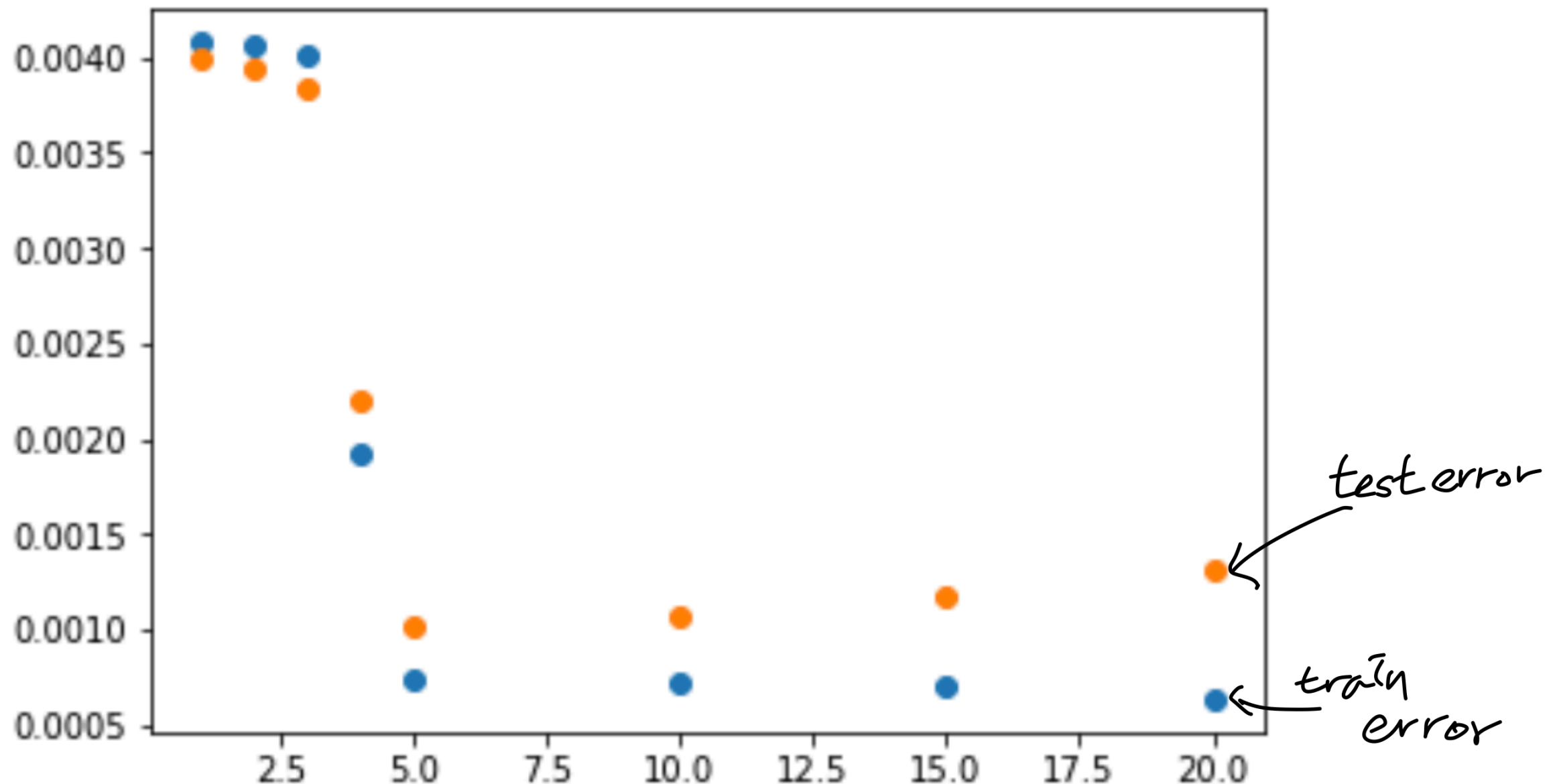| | | | |
|---|---|---|---|
| 0.004011490884146126 | 0.0019177229776174974 | 0.0007426853089970962 | 0.0006350545819906561 $\mathrm{MSE}_{\mathrm{train}}$ |
| 0.003831010290504173 | 0.0022004927204447942 | 0.0010239915404334238 | 0.0013118370515986446 $\mathrm{MSE}_{\mathrm{test}}$ |

      degree 3            degree 4            degree 5            degree 20

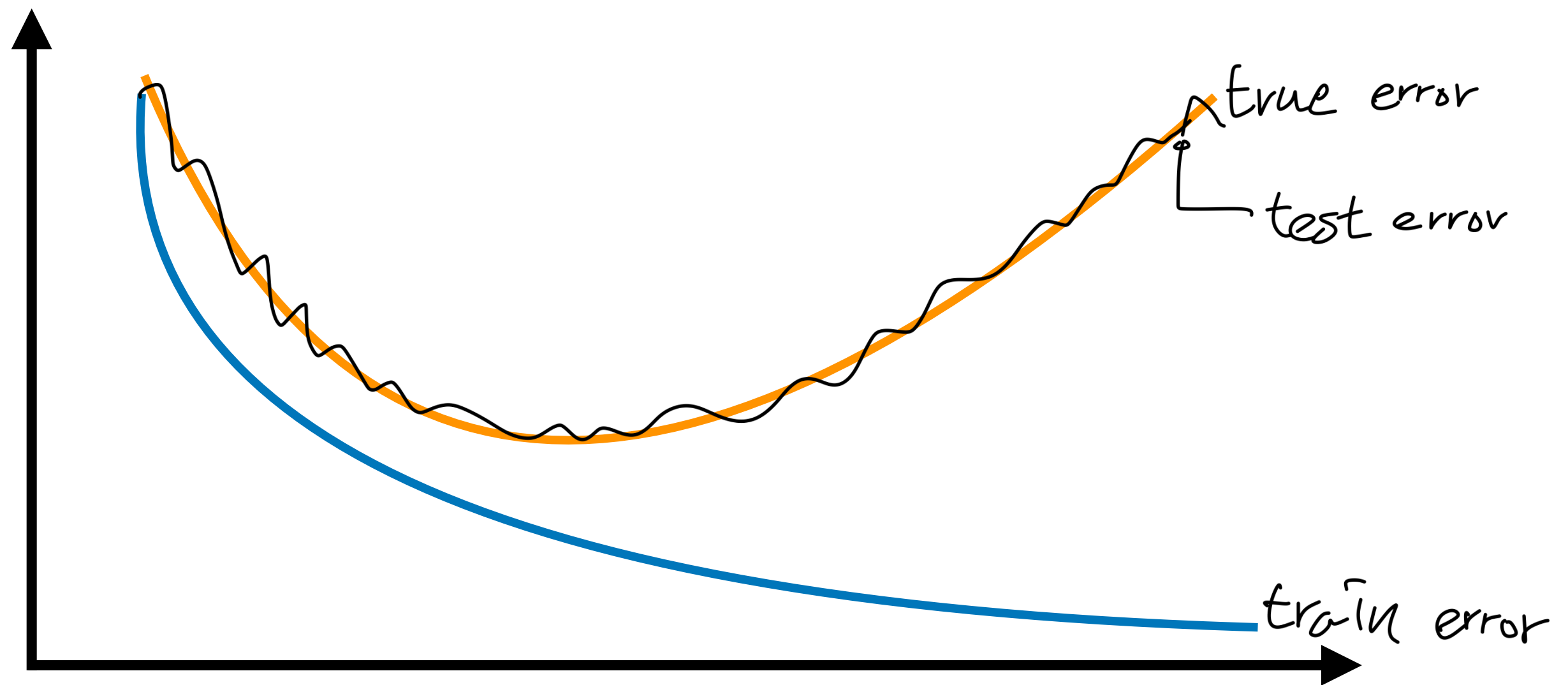# How does one choose which model to use?



- first use 60 data points to train and 60 data points to test
- then choose degree 5 as per the above test error
- now re-train on all 120 data points with degree 5 polynomial model

# Cross validation

- systematic method for out-of-sample validation
  - divide the data into *k* **folds**
  - for each i, fit predictor on all data but fold i
  - compute test error on fold i
  - average the test error across the folds


- gives some idea of the variability (or variance) of the test error


- we can estimate **stability** of the ML pipeline, by comparing the model parameters found in each fold

- test error gives an approximation of the (unknown) true error
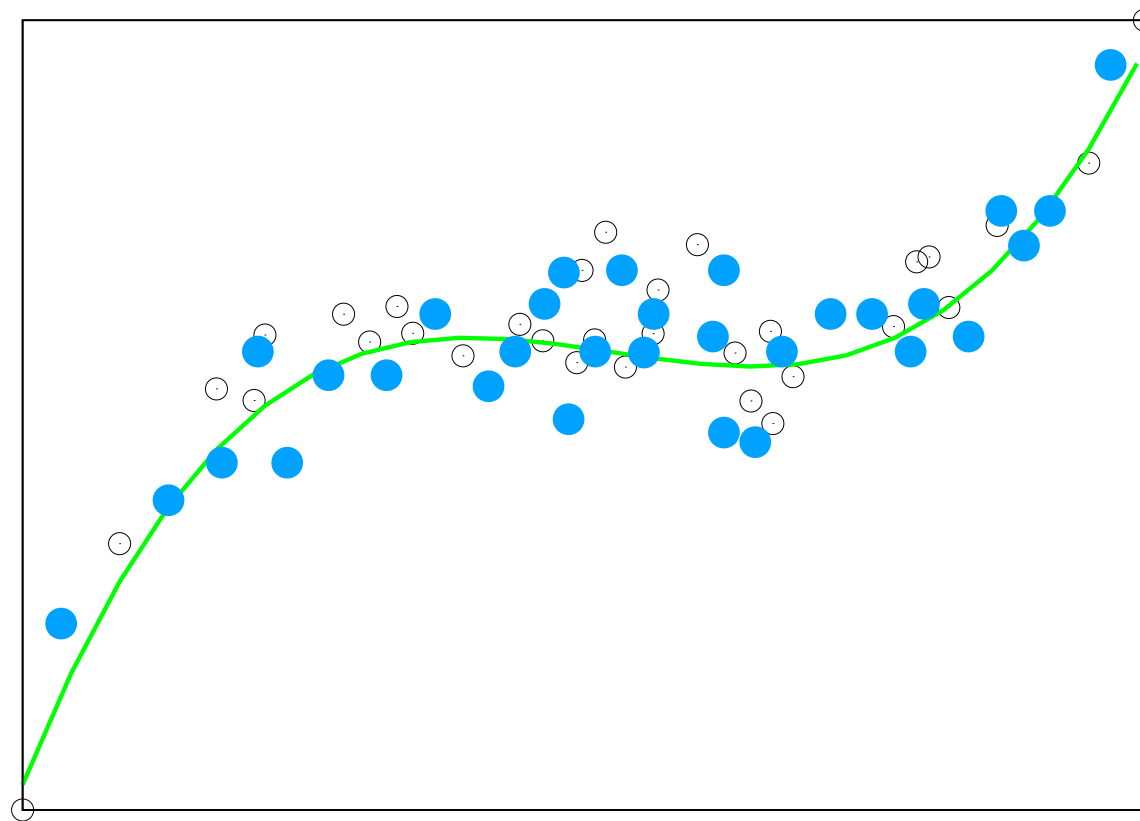- train error goes down monotonically w.r.t. model complexity

# Why does true error go down and up?

- three sources of error: **noise**, **bias**, and **variance**

$$\mathrm{MSE_{true}} = \mathrm{MSE_{noise}} + \mathrm{MSE_{bias}} + \mathrm{MSE_{variance}}$$

- error from **noise** in the data cannot be reduced



$$y = f_0(x) + \varepsilon \quad , \quad \hat{y} = f_0(x)$$

$$\mathrm{MSE_{noise}} = \mathbb{E}[\varepsilon^2] \quad \leftarrow \quad y - \hat{y} = \varepsilon$$

# Low complexity models

- suppose we train a constant function, many times each with *N* samples from $y = f_0(x) + \varepsilon$
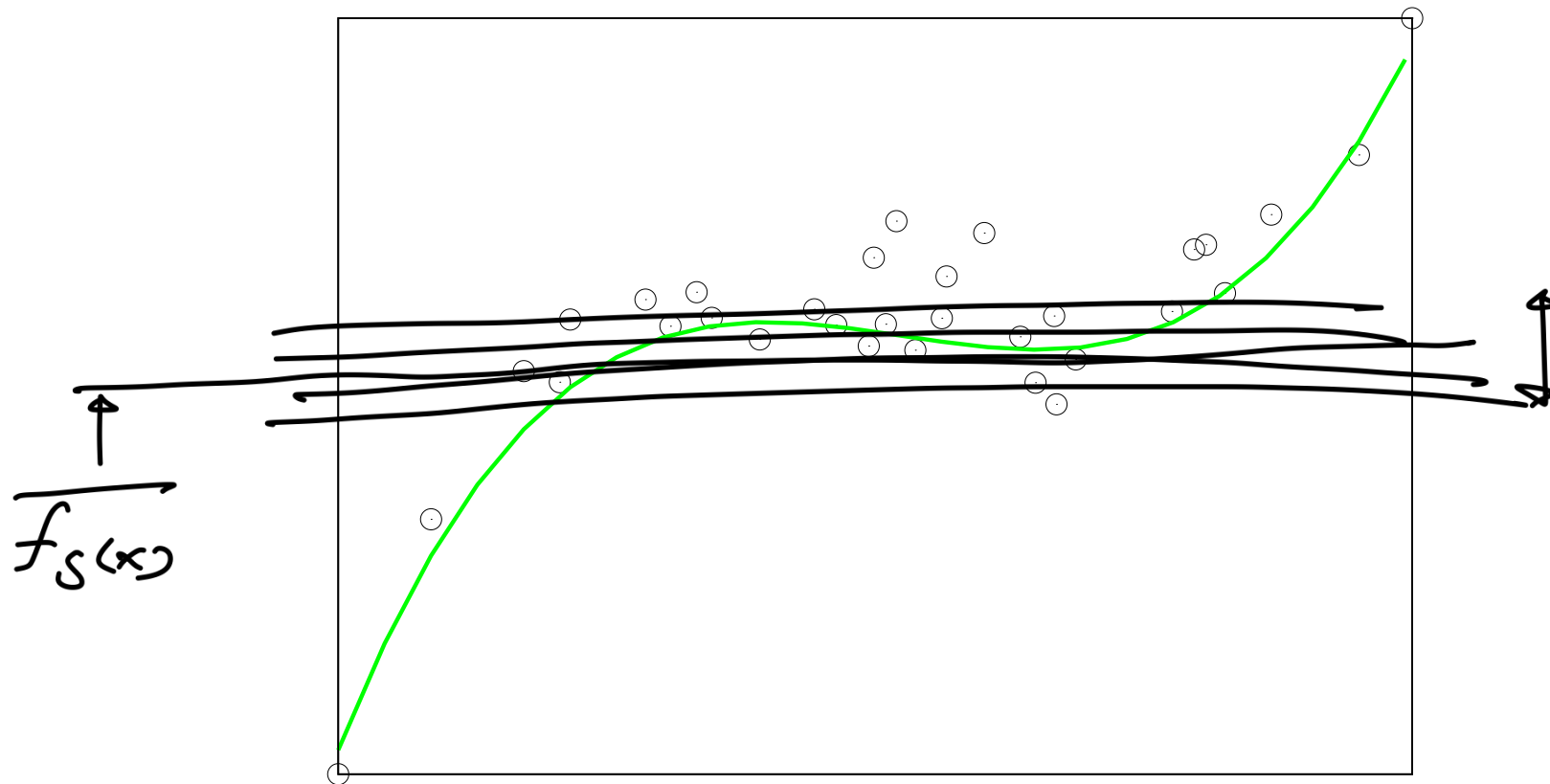


$$\text{Bias} = \mathbb{E}\big[\big|f_0(x) - \overline{f_S(x)}\big|\big]$$

$$\geq \tfrac{1}{4}\big(f_{S_1} + f_{S_2} + f_{S_3} + f_{S_4}\big)$$

- low complexity model has a large **bias**

if we change N or model class, bias changes

# Low complexity models

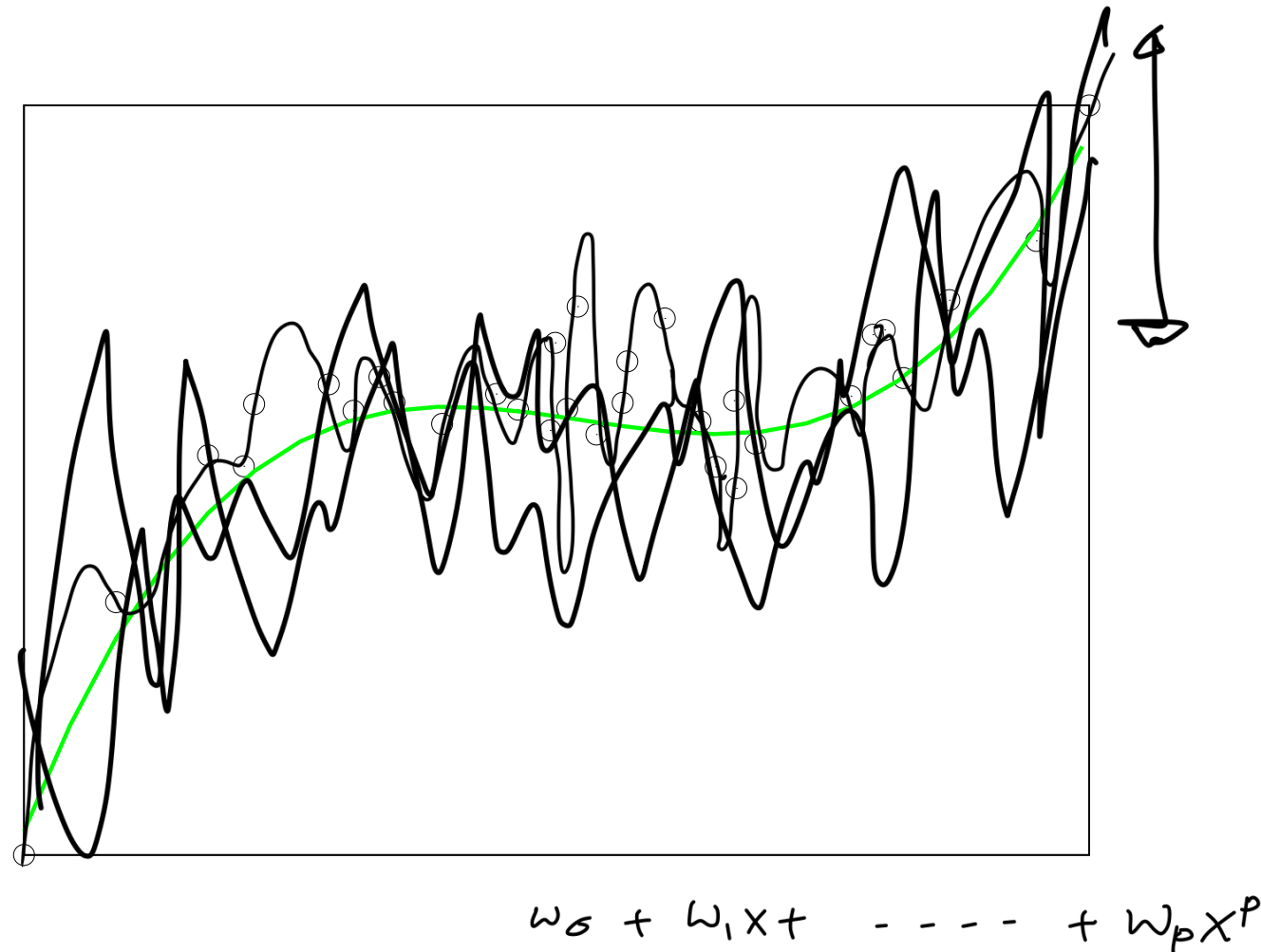- suppose we train a constant function, many times each with *N* samples from $y = f_0(x) + \varepsilon$



$$\text{Variance} = \mathbb{E}\big[\, (f_{S_i}(x) - \overline{f_S(x)})^2 \,\big]$$

- low complexity model has a small **variance**

# High complexity models

- suppose we train a high degree polynomial function, many times each with *N* samples from $\quad y = f_0(x) + \varepsilon$



$$\omega_0 + \omega_1 x + \quad \text{- - - -} \quad + \omega_p x^p$$

$$\text{Variance} = \mathbb{E}\big[ (f_S(x) - \overline{f_S(x)})^2 \big]$$

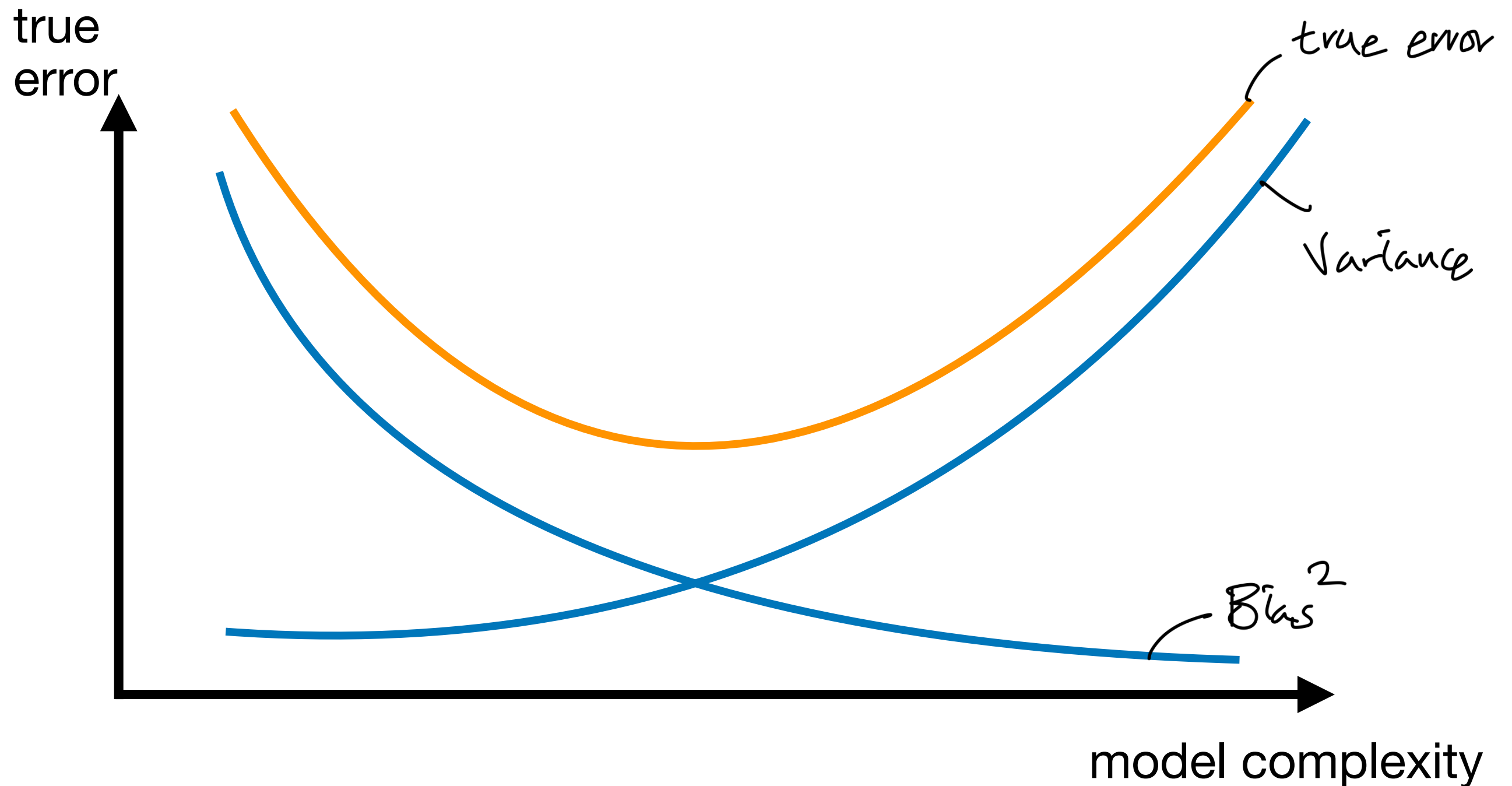- high complexity model has a large **variance** but small **bias**

# Bias-Variance tradeoff

- for fixed sample size $N$,

$$\mathrm{MSE_{true}} = \mathrm{MSE_{noise}} + \mathrm{Bias}^2 + \mathrm{Variance}$$

# For fixed model complexity

- suppose we fix model complexity such that

$$f_0(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + w_5 x^5 + \varepsilon \quad \leftarrow \text{True Model}$$

$$f(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 \quad \leftarrow \text{Predictor}$$



Variance large

true error

Bias due to model mismatch

$\text{MSE}_{\text{noise}}$

train error

easy to fit few samples

sample size