

Validation

Sewoong Oh

CSE446

University of Washington

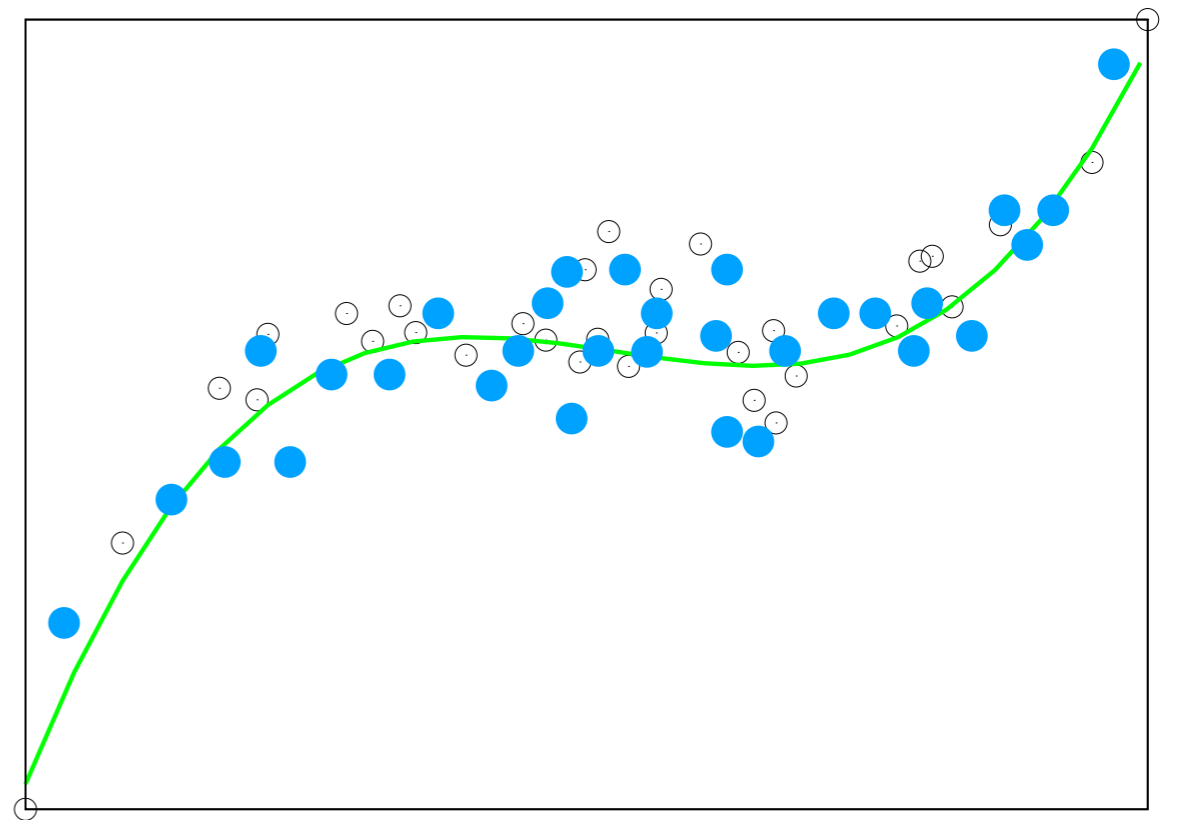
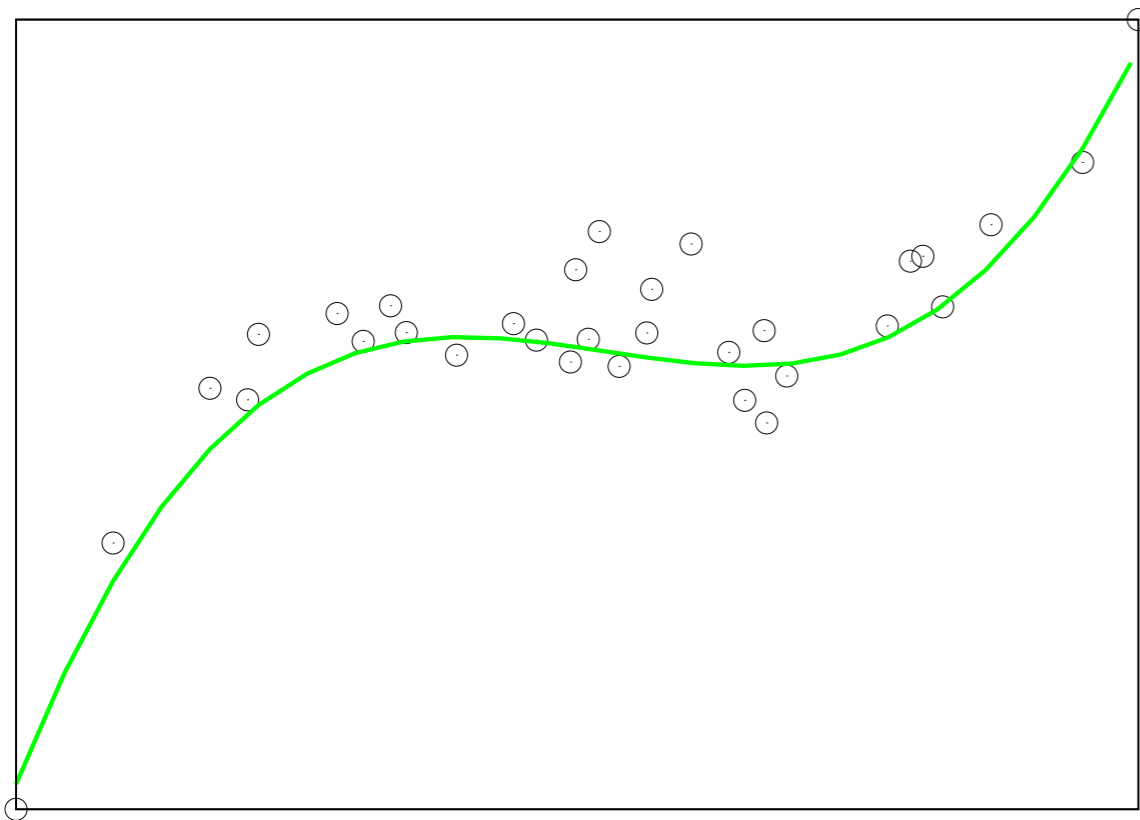
Generalization:
**how do we know which model is better in
predicting unseen data?**

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

In-sample and out-of-sample data

- the data used to train 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) white circles: MSE 174
- **predict** y for 30 (out-of-sample) blue circles: MSE 192
- conclude this predictor generalizes, as in-sample MSE \simeq out-of-sample MSE

Out-of-sample 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** used to evaluate the predictor
- we assume that test set is similar to unseen data
- test set should never be used in training

Out-of-sample Validation

- given a single dataset $S = \{(x_i, y_i)\}_{i=1}^n$
- we split the dataset into two: training set and test set
- selection of data train/test should be done randomly (80/20 or 90/10 are common)
- we use **training error** (i.e. empirical risk on training dataset) for optimization (or finding the model)

$$\text{minimize } \mathcal{L}_{\text{train}}(w) = \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} \ell(f(x_i), y_i)$$

- we use **test error** (i.e. empirical risk on test dataset) for **validation**, checking if the model behaves as expected

$$\mathcal{L}_{\text{test}}(w) = \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} \ell(f(x_i), y_i)$$

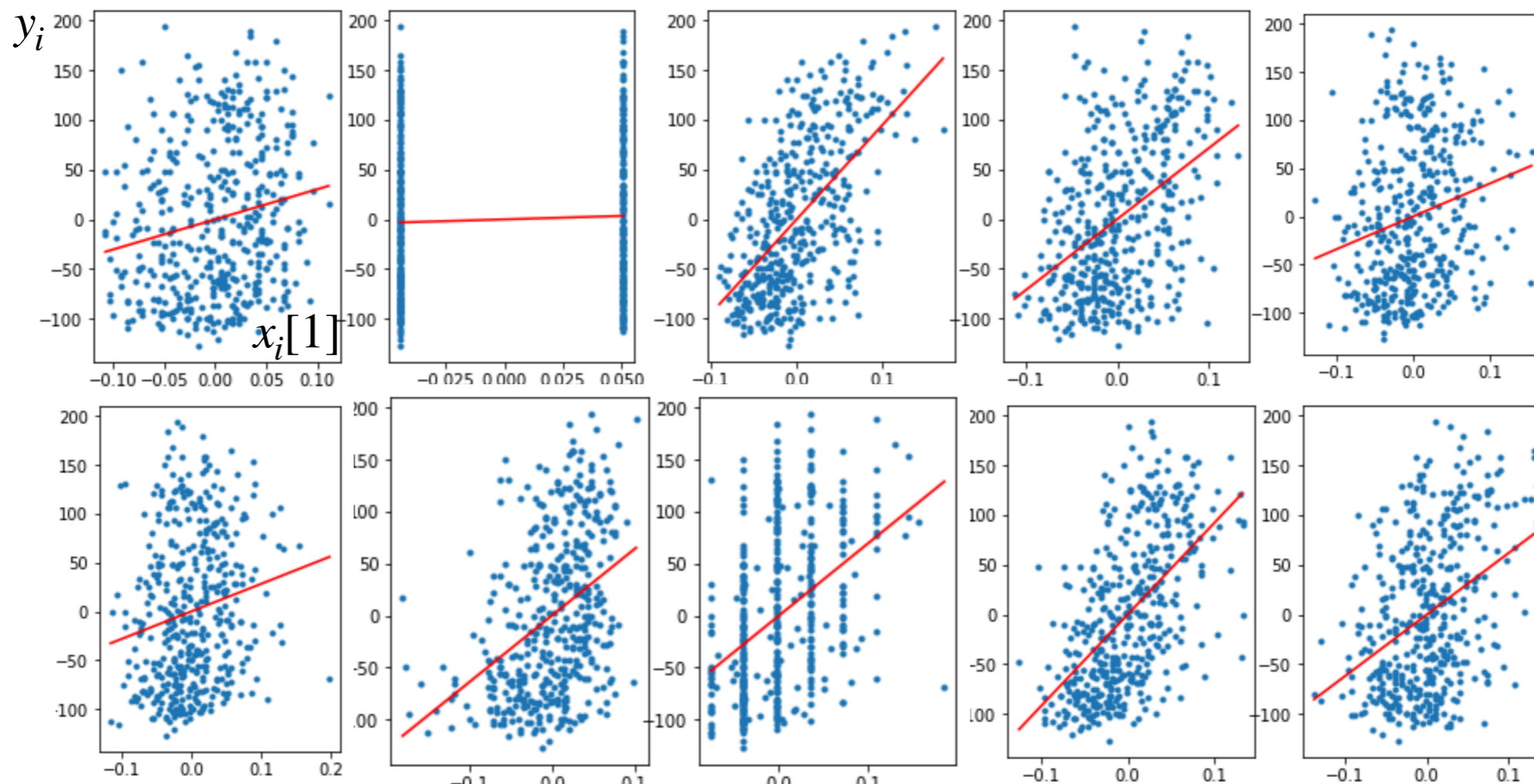
- we say a model or predictor is **overfit** if

$$\mathcal{L}_{\text{train}} \ll \mathcal{L}_{\text{test}}$$

	small training error	large training error
small test error	generalizes well performs well	possible, but unlikely
large test error	fails to generalize	generalizes well performs poorly

Choosing a predictor

- validation is useful in choosing a predictor
- typically, one trains multiple candidate predictors and chooses the predictor that has the smallest test error
- Example: Diabetes
 - 10 explanatory variables
 - from 442 patients
 - we use half for train and half for validation



Example: Diabetes

Features	Train MSE	Test MSE
All	2640	3224
S5 and BMI	3004	3453
S5	3869	4227
BMI	3540	4277
S4 and S3	4251	5302
S4	4278	5409
S3	4607	5419
None	5524	6352

- **test MSE is the primary criteria for model selection**
- Using only 2 features (S5 and BMI), one can get very close to the prediction performance of using all features
- Combining S3 and S4 does not give any performance gain

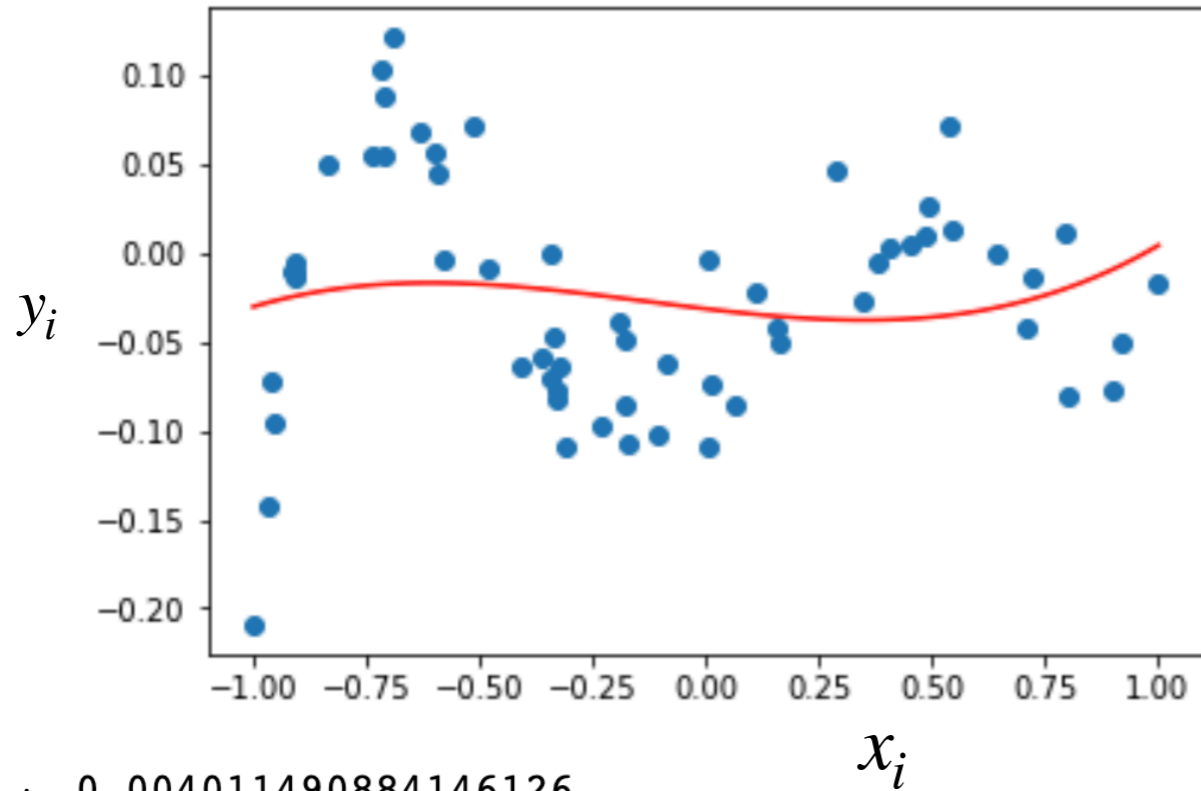
Overfitting

- 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

$$\hat{y} = f(x) = w_0 + w_1x + w_2x^2 + \dots + w_px^p$$

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

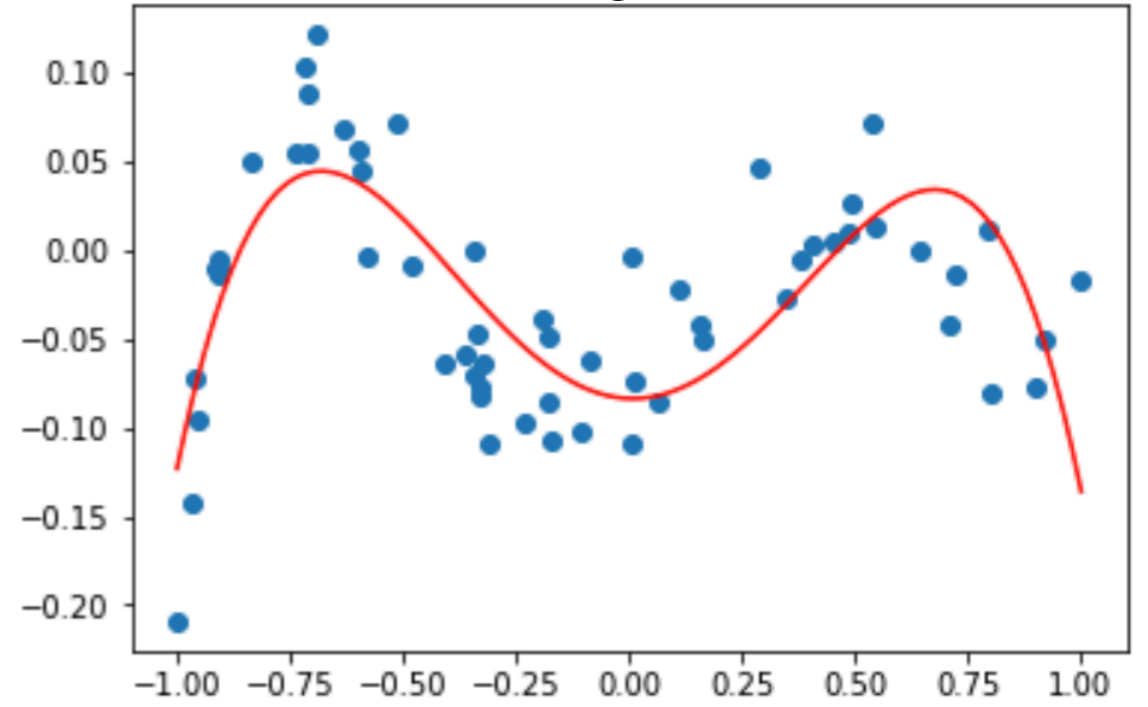
degree 3



MSE_{train} 0.004011490884146126

MSE_{test} 0.003831010290504173

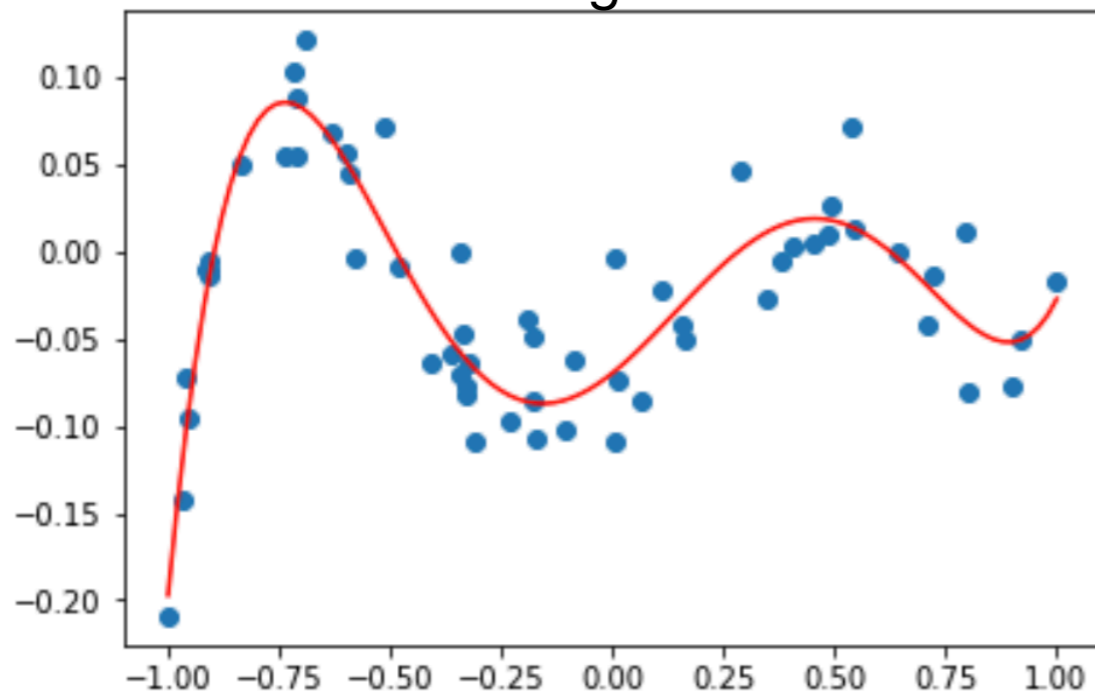
degree 4



0.0019177229761741974

0.002200492720447942

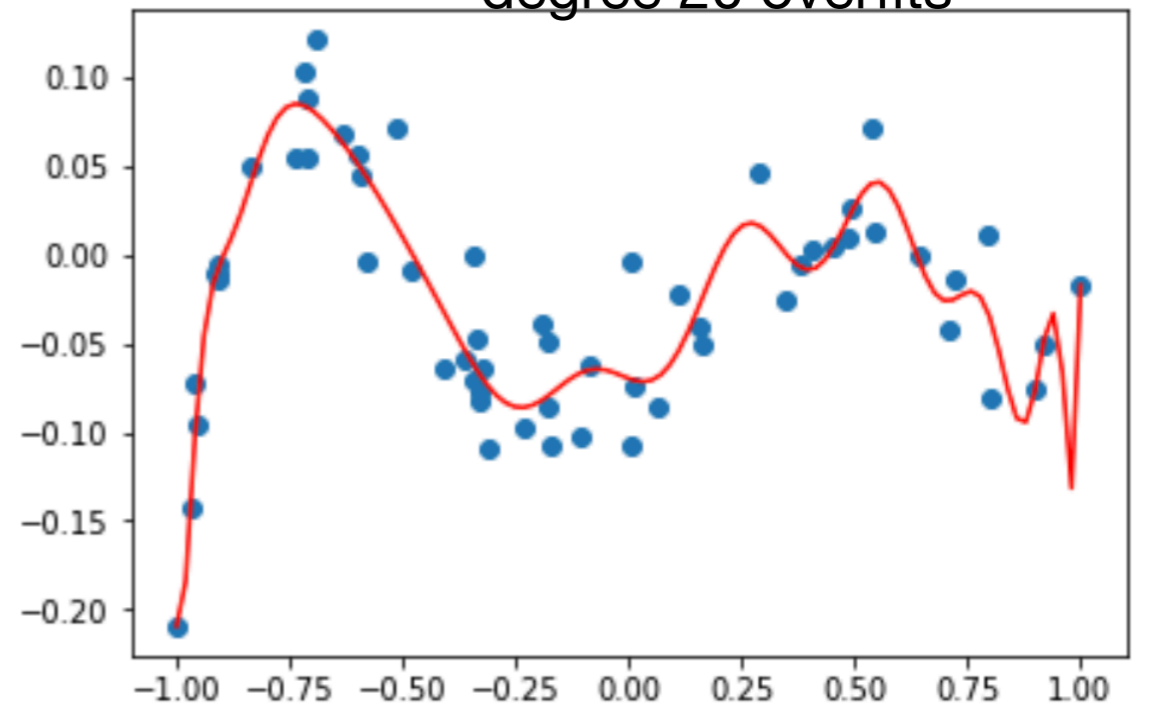
degree 5



0.0007426853089970962

0.0010239915404334238

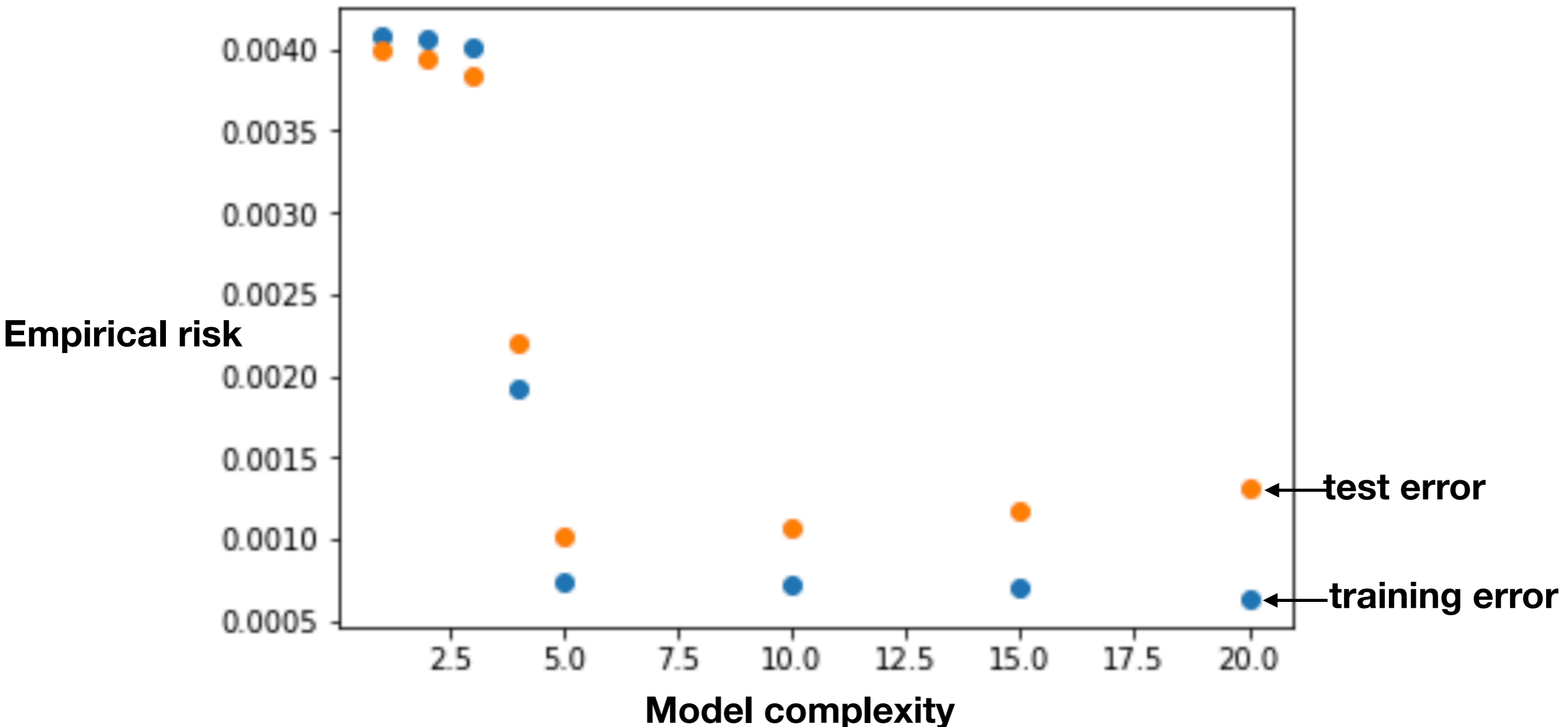
degree 20 overfits



0.0006350545819906561

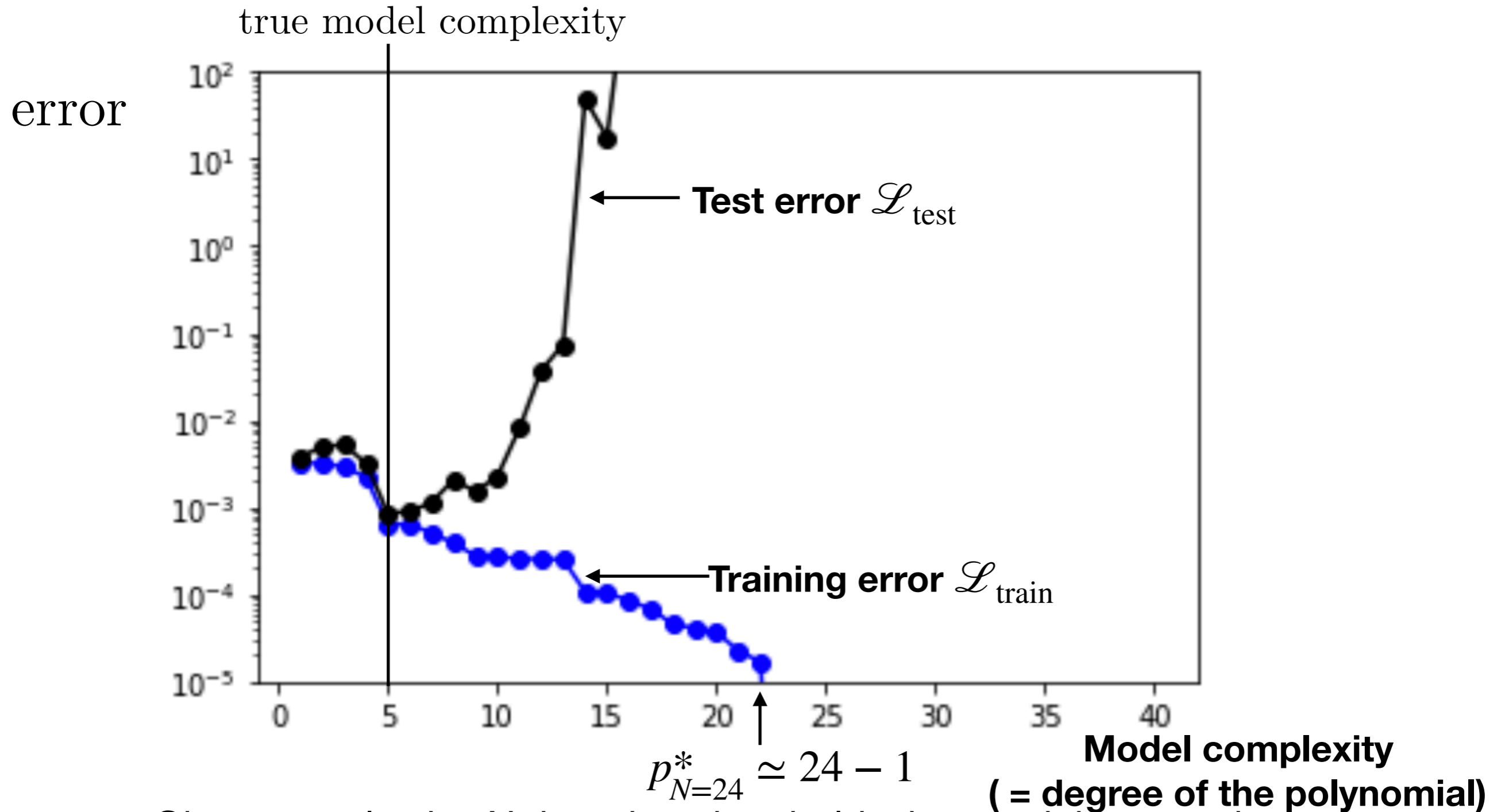
0.0013118370515986446

How does one choose which model to use?



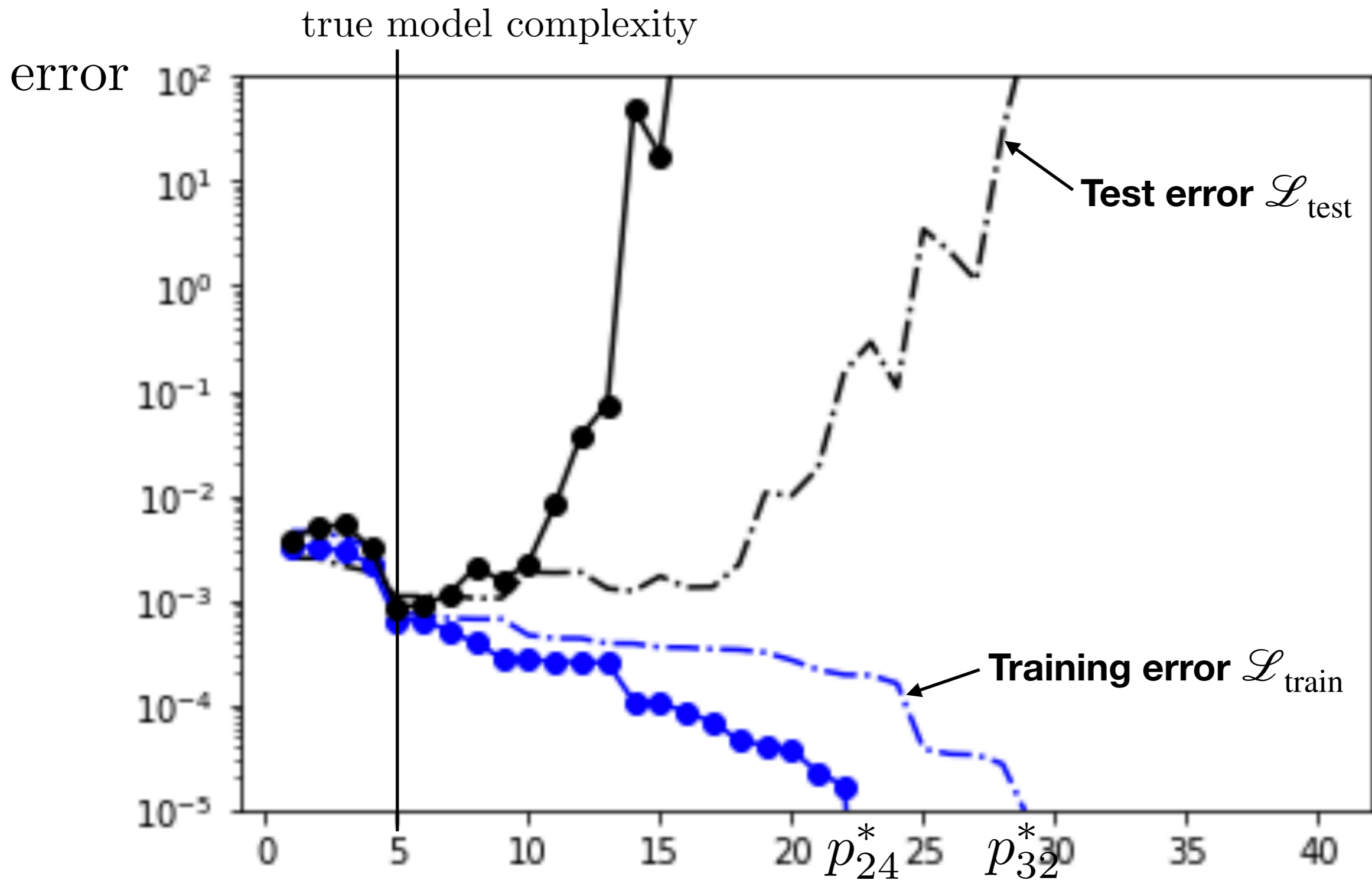
- first use 60 data points to train and 60 data points to test
- then choose degree 5 as per the above test error
- now re-train on all 120 data points with degree 5 polynomial model

- let us first fix sample size $N=30$, collect one dataset of size N , randomly shuffle the dataset, and fix one training set S_{train} and test set S_{test} via 80/20 split
- then we run multiple validations and plot the computed MSEs for all values of p that we are interested in



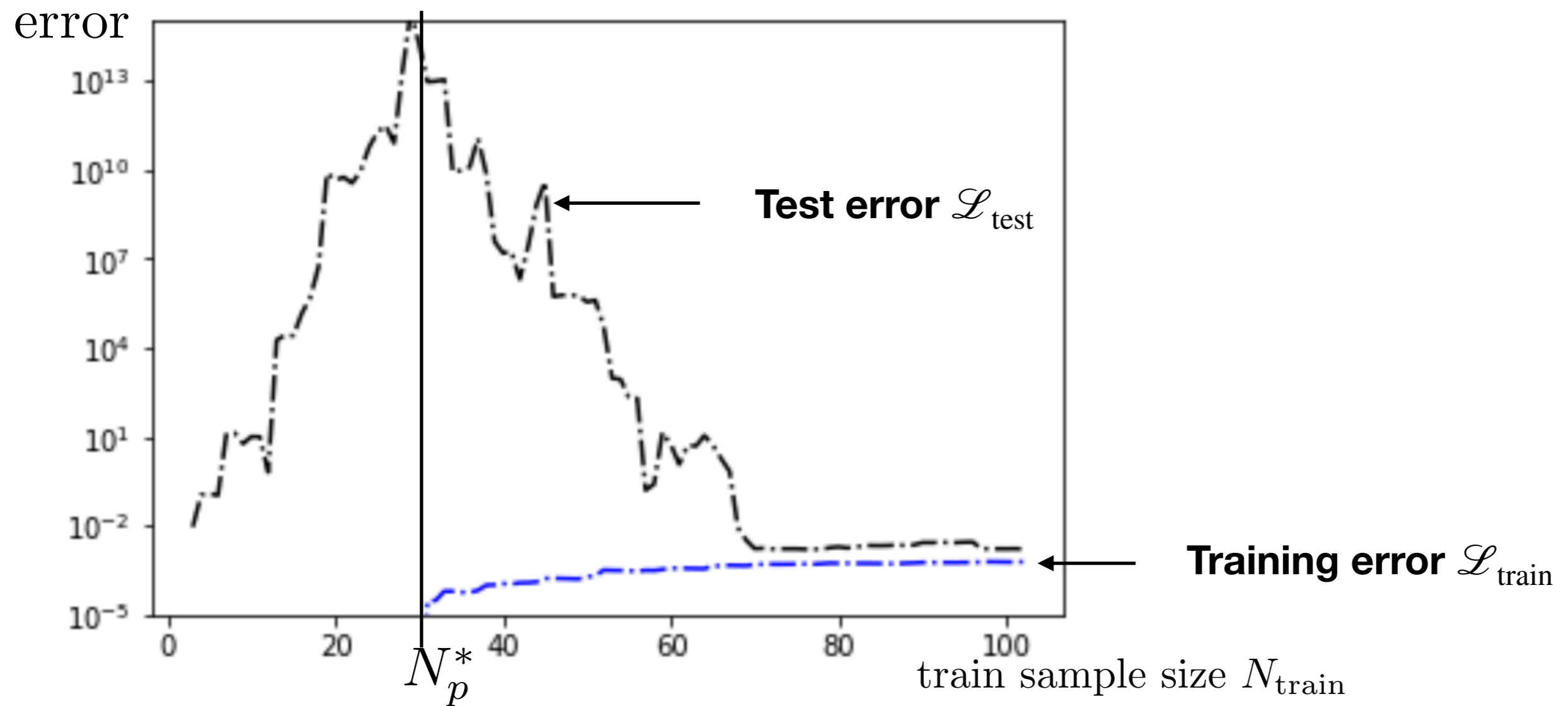
- Given sample size N there is a threshold where training error is zero
- Training error is **always** monotonically non-increasing
- Test error has a trend of going down and then up, but fluctuates

- let us now repeat the process changing the sample size to **N=40**, and see how the curves change



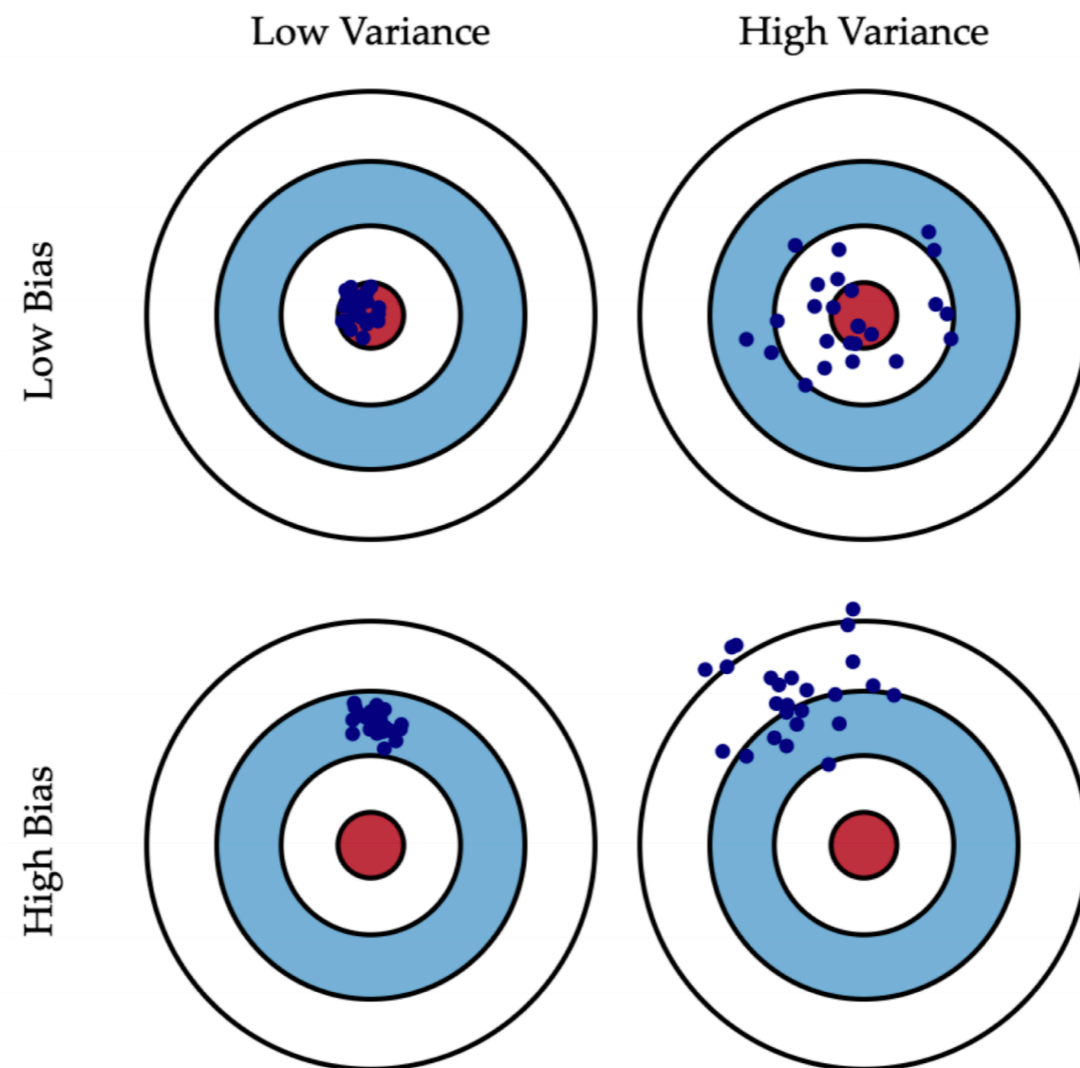
- The threshold moves right
- Training error tends to increase: more points need to fit
- Test error tends to decrease: overfitting happens later

- let us now fix predictor model complexity $p=30$, collect multiple datasets by starting with 3 samples and adding one sample at a time to the training set, but keeping a large enough test set fixed
- then we run multiple validations and plot the computed MSEs for all values of train sample size N_{train} that we are interested in



- There is a threshold below which training error is zero (extreme overfit)
 - Below this threshold, test error is meaningless, as there are multiple predictors with zero training error
 - Test error tends to decrease
 - Training error tends to increase
- 15
- why do they meet?

From practice to theory



Courtesy of Scott Fortmann-Roe

Notations

- the model is specified by the distribution of data $p_{x,y}$ for the paired examples (x_i, y_i)
- we denote our predictor by $f_{S_{\text{train}}}(x)$ to emphasize that our predictor depends on the training data $S_{\text{train}} = \{(x_i, y_i)\}_{i=1}^n$ of size n , each coming i.i.d. from distribution $p_{x,y}$
- we denote the test data by $S_{\text{test}} = \{(x_i, y_i)\}_{i=n+1}^{n+m}$ of size m , also i.i.d. from $p_{x,y}$
- when we take expectation of a function of random variables, we use the following notation to indicate what we are taking expectation of:

$$\mathbb{E}_{(x,y) \sim p_{x,y}} [F(x, y)]$$

indicating that the pair (x, y) is drawn from $p_{x,y}$

- we will simplify the subscript, whenever it is clear from the context, for example we might write

$$\mathbb{E}_{p_{x,y}} [F(x, y)]$$

or even

$$\mathbb{E}[F(x, y)]$$

Expected test error

- goal of training a predictor is to get test error small, defined as the empirical risk on the test set

$$\mathcal{L}_{\text{test}} = \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} (f_{S_{\text{train}}}(x_i) - y_i)^2$$

because this is a surrogate of the expected error on (randomly chosen) unseen data

- the expected mean squared error (true error) on a new data (x, y) is defined as

$$\begin{aligned} \mathcal{L}_{\text{true}} &= \mathbb{E}_{S_{\text{test}} \sim p_{x,y}^m, S_{\text{train}} \sim p_{x,y}^n} [\mathcal{L}_{\text{test}}] \\ &= \mathbb{E}_{S_{\text{test}} \sim p_{x,y}^m, S_{\text{train}} \sim p_{x,y}^n} \left[\frac{1}{m} \sum_{i=n+1}^{n+m} (f_{S_{\text{train}}}(x_i) - y_i)^2 \right] \\ &= \mathbb{E}_{(x,y) \sim p_{x,y}, S_{\text{train}} \sim p_{x,y}^n} [(f_{S_{\text{train}}}(x) - y)^2] \end{aligned}$$

where the last line follows from the i.i.d. assumption, and this true error is what we really care about, and hope to minimize

- for simplicity, we will write

$$\mathcal{L}_{\text{true}} = \mathbb{E}_{p_{x,y}, S_{\text{train}}} [(f_{S_{\text{train}}}(x) - y)^2]$$

- we will decompose this expected test error, to identify three sources of error

Canonical model

- recall the law of total expectation (or tower rule):

$$\mathbb{E}_{p_{x,y}} [F(x, y)] = \mathbb{E}_{p_x} \left[\mathbb{E}_{p_{y|x}} [F(x, y) | x] \right],$$

for any function $F(x, y)$ and any joint distribution $p_{x,y}$

- we focus on analyzing the conditional expectation

$$\mathbb{E}_{p_{y|x}, S_{\text{train}}} [(f_{S_{\text{train}}}(x) - y)^2 | x]$$

as by the law of total expectation (or tower rule), we have

$$\mathbb{E}_{p_{x,y}, S_{\text{train}}} [(f_{S_{\text{train}}}(x) - y)^2] = \mathbb{E}_{p_x} \left[\mathbb{E}_{p_{y|x}, S_{\text{train}}} [(f_{S_{\text{train}}}(x) - y)^2 | x] \right]$$

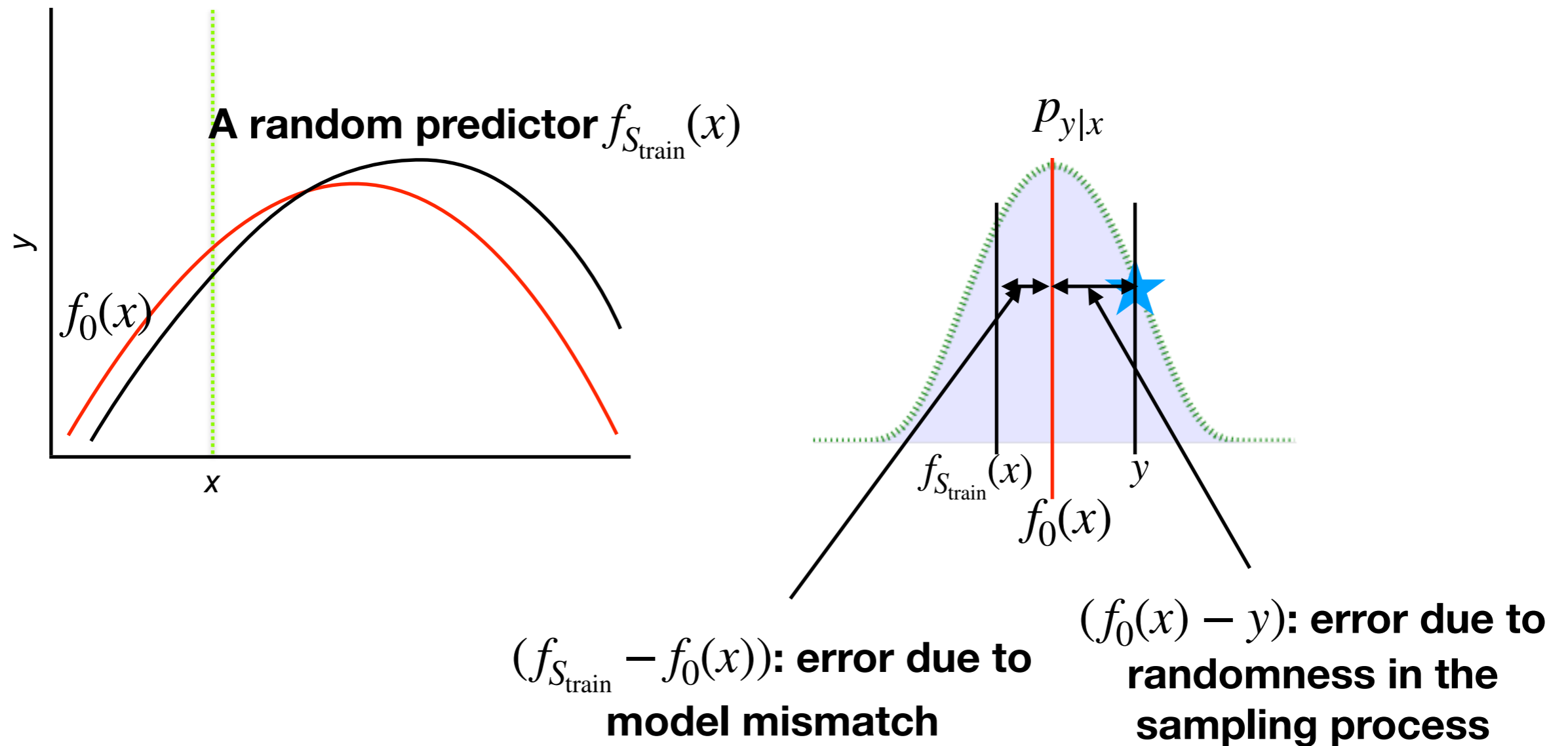
- the bias-variance tradeoff we show for the conditional expectation, will imply a similar result on the joint expectation $p_{x,y}$ by simply taking the expectation with respect to p_x to the resulting formula (we do this in equation (1) on slide 21)
- this implies that we only need to specify the conditional distribution $p_{y|x}$ to proceed with the analysis, and we will focus on the **canonical model** where

$$y = f_0(x) + \varepsilon,$$

where ε is drawn from $N(0, \sigma^2)$, zero mean Gaussian with variance σ^2

- note that this use of canonical model is without loss of generality, as it the same analysis can be used to capture bias-variance tradeoff for any $p_{x,y}$ but with heavier notations

- model: $y_i = f_0(x_i) + \varepsilon_i$ for both training and test samples



- we will analyze the conditional expectation of the true error

$$\mathbb{E}_{p_{y|x}, S_{\text{train}}} [(f_{S_{\text{train}}}(x) - y)^2 | x]$$

Bias-variance tradeoff

- the conditional true error can be written as

$$\begin{aligned}\mathbb{E}_{p_{y|x}, S_{\text{train}}}[(f_{S_{\text{train}}}(x) - y)^2 | x] &= \mathbb{E}\left[\underbrace{(f_{S_{\text{train}}}(x) - f_0(x))}_A - \underbrace{(y - f_0(x))}_B \right]^2 | x \\ &= \underbrace{\mathbb{E}_{S_{\text{train}}}[(f_{S_{\text{train}}}(x) - f_0(x))^2 | x]}_{\text{learning error} \geq 0} + \underbrace{\mathbb{E}_{p_{y|x}}[(f_0(x) - y)^2 | x]}_{\text{irreducible error} = \sigma^2}\end{aligned}$$

this follows from the fact that $\mathbb{E}[(A - B)^2] = \mathbb{E}[A^2] - 2\mathbb{E}[AB] + \mathbb{E}[B^2]$ and the noise is zero mean, i.e.

$$\mathbb{E}[AB] = \mathbb{E}[(f_{S_{\text{train}}}(x) - f_0(x))(y - f_0(x)) | x] = \mathbb{E}[f_{S_{\text{train}}}(x) - f_0(x) | x] \underbrace{\mathbb{E}[\varepsilon | x]}_{=0} = 0$$

- irreducible error**

- is due to the inherent noise in the samples, and is impossible to get rid of
- does not depend on our predictor $f(x)$
- Is a lower bound on achievable expected test error

- learning error**

- is due to the randomness (and limited sample size) in the training data
- Further decomposed into **bias** and **variance**

- the learning error can be further decomposed (with a similar trick) as

$$\begin{aligned} \mathbb{E}_{S_{\text{train}}} [(f_{S_{\text{train}}}(x) - f_0(x))^2 | x] &= \mathbb{E} \left[\left((f_{S_{\text{train}}}(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x]) - (f_0(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x]) \right)^2 | x \right] \\ &= \underbrace{\mathbb{E} \left[(f_{S_{\text{train}}}(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x])^2 | x \right]}_{\text{Variance}} + \underbrace{(f_0(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x])^2}_{\text{Bias}^2} \end{aligned}$$

this follows from $\mathbb{E}[(f_{S_{\text{train}}}(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x]) | x] = 0$

- this theoretical analysis explains the behavior of true error $\mathcal{L}_{\text{true}}$

$$\mathcal{L}_{\text{true}} = \underbrace{\sigma^2}_{\text{irreducible error}} + \underbrace{\mathbb{E}[(f_{S_{\text{train}}}(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x])^2]}_{\text{(expected) Variance}} + \underbrace{\mathbb{E}_{p_x}[(f_0(x) - \mathbb{E}[f_{S_{\text{train}}}(x) | x])^2]}_{\text{(expected) Bias}^2} \quad (1)$$

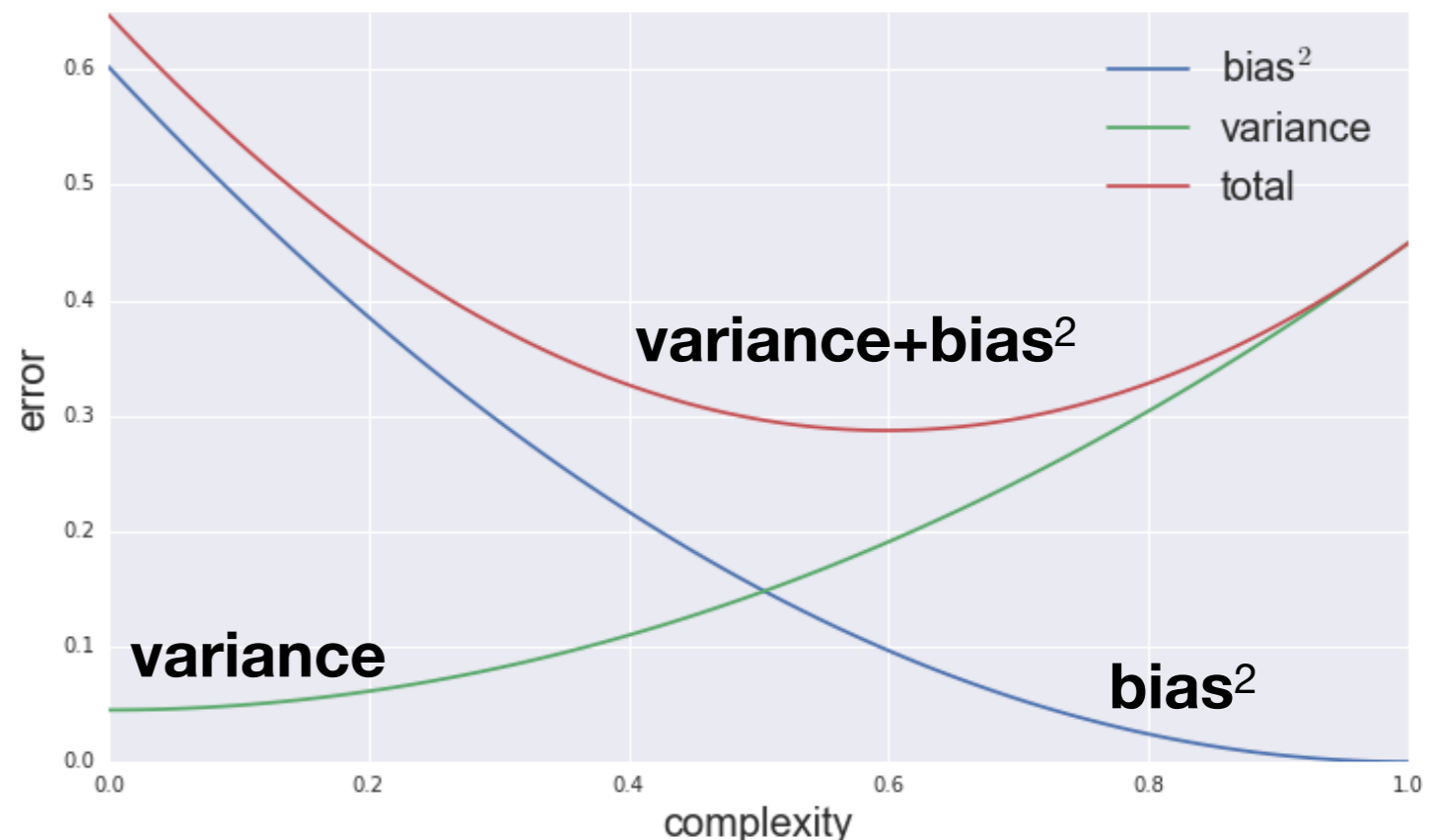
- Whether we condition on x or not when referring to **variance** and **bias** should be clear from the context

- bias**

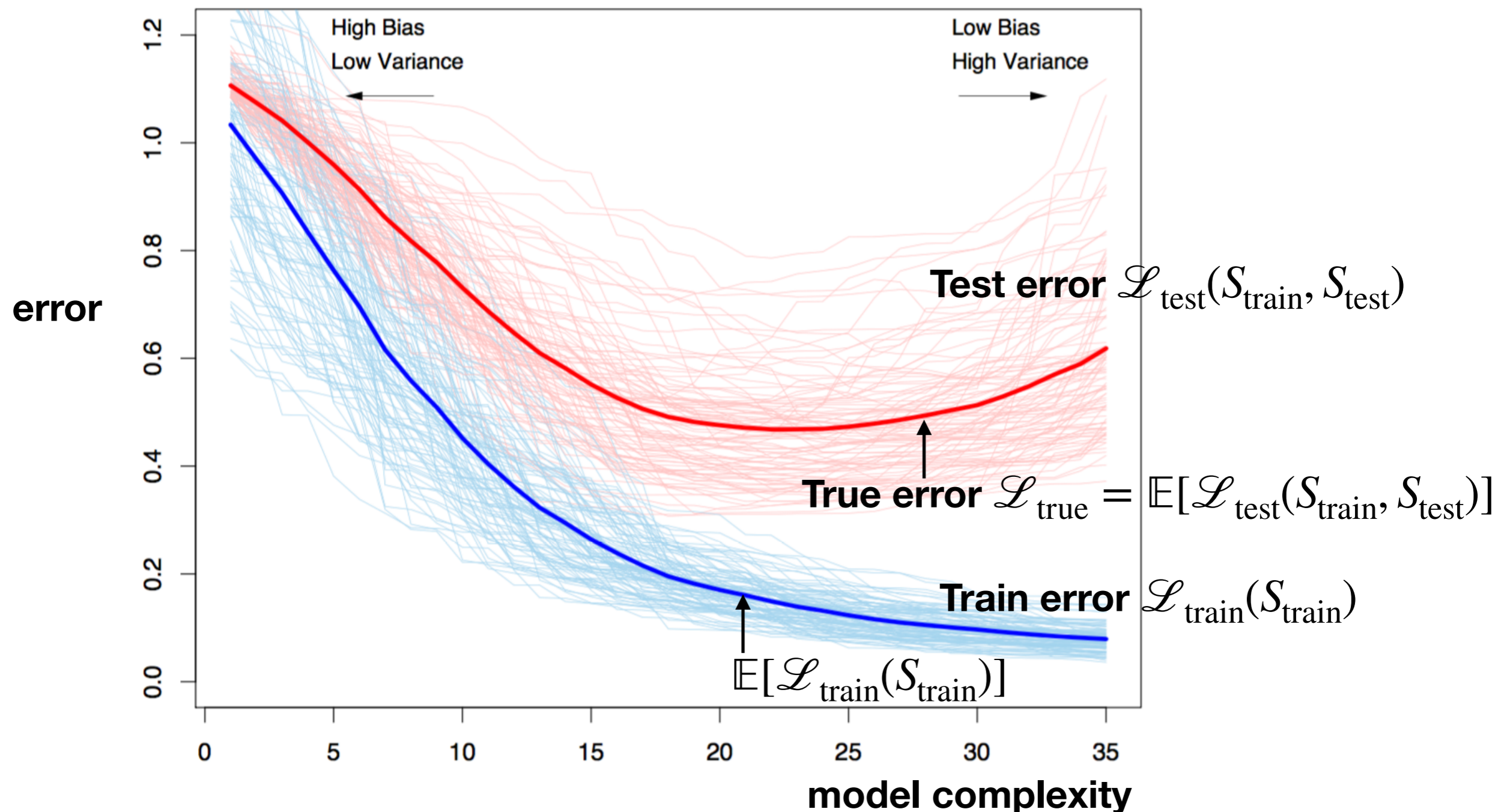
- measures how the predictor is mismatched with the true model in expectation

- variance**

- measures how the predictor varies each time with a new training datasets



- recall the test error is an unbiased estimator of the true error, i.e. $\mathcal{L}_{\text{true}} = \mathbb{E}[\mathcal{L}_{\text{test}}]$
- and theory explains **true error**, and hence expected behavior of the (random) **test error**



Simple Gaussian example

- model: $y_i = w_1 x_i[1] + w_2 x_i[2] + 0 \cdot x_i[3] + \varepsilon_i$
- $x_i[1], x_i[2], x_i[3], \varepsilon_i \sim$ i.i.d. Gaussian $N(0, \sigma^2)$
- training data $\{(x_i, y_i)\}_{i=1}^n$

- let data matrix be $\mathbf{X} = \begin{bmatrix} x_1[1] & x_1[2] & x_1[3] \\ \vdots & \vdots & \vdots \\ x_n[1] & x_n[2] & x_n[3] \end{bmatrix} \in \mathbb{R}^{n \times 3}$
- $\mathbf{X}[:, \mathbf{1}]$ denotes the first column of \mathbf{X}

Simple Gaussian example

- model: $y_i = w_1 x_i[1] + w_2 x_i[2] + 0 \cdot x_i[3] + \varepsilon_i$

- **Example 1 (simple predictor):**

- consider a **simple** model of fitting only the first feature

$$\hat{y} = \hat{w}_1 x[1]$$

- linear least squares gives

$$\hat{w}_1 = (\mathbf{X}[:, \mathbf{1}]^T \mathbf{X}[:, \mathbf{1}])^{-1} \mathbf{X}[:, \mathbf{1}]^T \mathbf{y}$$

with $\mathbf{X}[:, \mathbf{1}] = \begin{bmatrix} x_1[1] \\ \vdots \\ x_n[1] \end{bmatrix}$, and $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = w_1 \mathbf{X}[:, \mathbf{1}] + w_2 \mathbf{X}[:, \mathbf{2}] + \underbrace{\begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}}_{\triangleq \varepsilon}$

cancels each other for the first term

- plugging in \mathbf{y} in the equation, we get

$$\begin{aligned} \hat{w}_1 &= (\mathbf{X}[:, \mathbf{1}]^T \mathbf{X}[:, \mathbf{1}])^{-1} \mathbf{X}[:, \mathbf{1}]^T (\mathbf{X}[:, \mathbf{1}] w_1 + \mathbf{X}[:, \mathbf{2}] w_2 + \varepsilon) \\ &= w_1 + (\mathbf{X}[:, \mathbf{1}]^T \mathbf{X}[:, \mathbf{1}])^{-1} \mathbf{X}[:, \mathbf{1}]^T (w_2 \mathbf{X}[:, \mathbf{2}] + \varepsilon) \end{aligned}$$

Simple Gaussian example

- for large enough n ,

$$\mathbf{X}[:, \mathbf{1}]^T \mathbf{X}[:, \mathbf{1}] = \sum_{i=1}^n x_i[1]^2 \simeq n\sigma^2,$$

by law of large numbers, and we will use this approximation to simplify the formula:

$$\hat{w}_1 = w_1 + \frac{1}{n\sigma^2} \mathbf{X}[:, \mathbf{1}]^T (w_2 \mathbf{X}[:, \mathbf{2}] + \varepsilon) \quad \text{and}$$

$$\mathbb{E}[\hat{w}_1] = w_1$$

where we used the fact that $\mathbf{X}[:, \mathbf{1}]$, $\mathbf{X}[:, \mathbf{2}]$, $\varepsilon \in \mathbb{R}^n$ are independent zero-mean vectors

- we are ready to compute the **bias** and **variance**

- first, conditioned on $x = (x[1], x[2], x[3])$

$$\text{bias}^2 = (\mathbb{E}[f(x)] - f_0(x))^2 = (w_1 x[1] - (w_1 x[1] + w_2 x[2]))^2 = (w_2)^2 x[2]^2,$$

and taking expectation over p_x

$$\mathbb{E}_{p_x}[\text{bias}^2] = (w_2)^2 \sigma^2$$

- note that

- this does not decrease with (training) sample size n
- this is due to not including $x[2]$ in our prediction, in other words, using a too simple predictor

Simple Gaussian example

- Now for the variance, conditioned on $x = (x[1], x[2], x[3])$,

since we have $\hat{w}_1 = w_1 + \frac{1}{n\sigma^2} \mathbf{X}[:, \mathbf{1}]^T (w_2 \mathbf{X}[:, \mathbf{2}] + \varepsilon)$

and $\mathbb{E}[\hat{w}_1] = w_1$

- variance = $\mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2 | x]$

$$= \mathbb{E} \left[\left(\frac{1}{n\sigma^2} \mathbf{X}[:, \mathbf{1}]^T (w_2 \mathbf{X}[:, \mathbf{2}] + \varepsilon) x[1] \right)^2 | x \right]$$

$$= \frac{x[1]^2}{n^2\sigma^4} \left(\mathbb{E} \left[\left(\sum_{i=1}^n w_2 x_i[1] x_i[2] + x_i[1] \varepsilon_i \right)^2 \right] \right)$$

$$= \frac{x[1]^2}{n^2\sigma^4} \left(\mathbb{E} \left[\left(\sum_{i=1}^n (w_2 x_i[1] x_i[2])^2 + \sum_{i=1}^n (x_i[1] \varepsilon_i)^2 + 2 \sum_{i=1}^n (w_2 x_i[1]^2 x_i[2]) + \sum_{i \neq j=1}^n (w_2 x_i[1] x_i[2] + x_i[1] \varepsilon_i)(w_2 x_j[1] x_j[2] + x_j[1] \varepsilon_j) \right) \right] \right)$$

zero mean

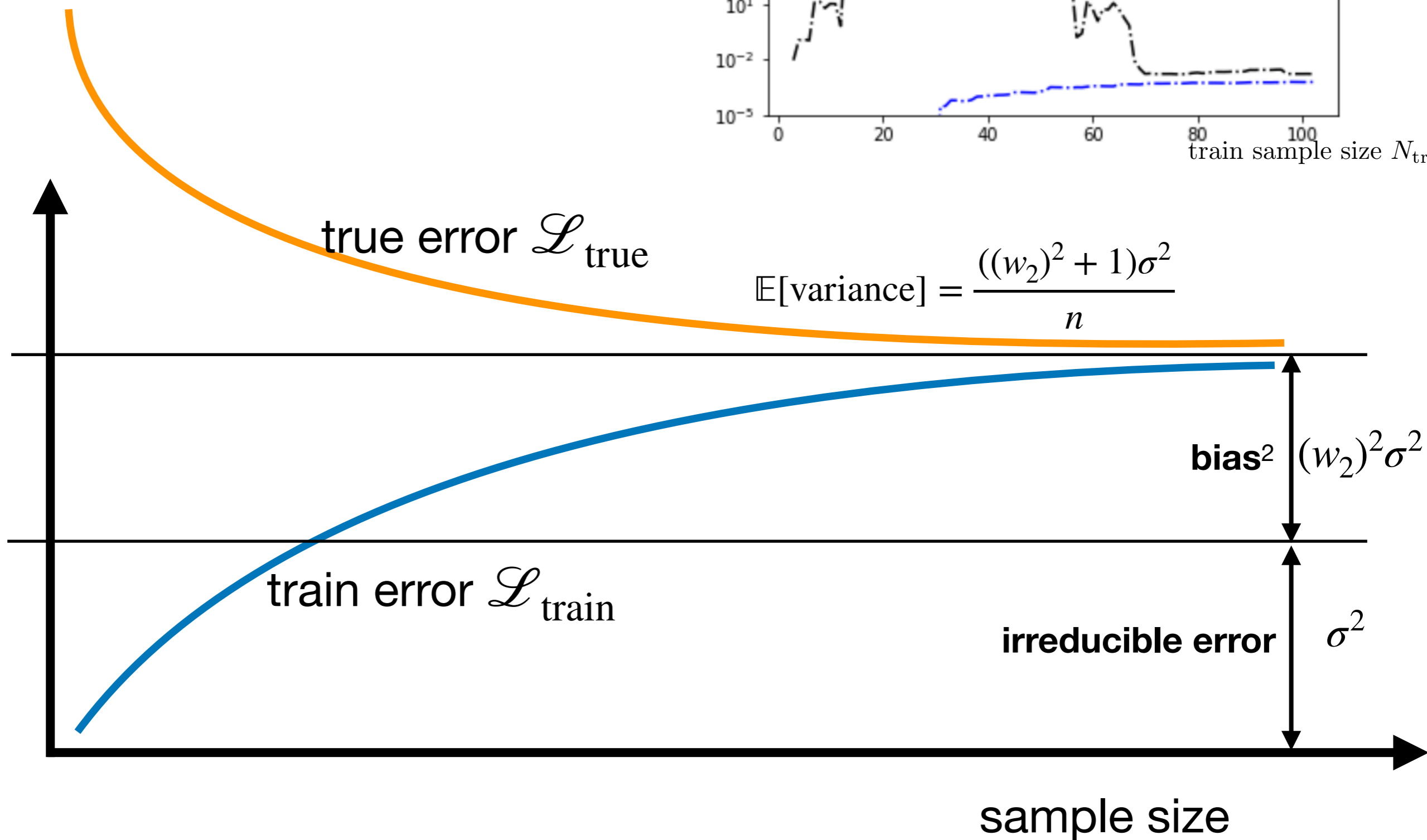
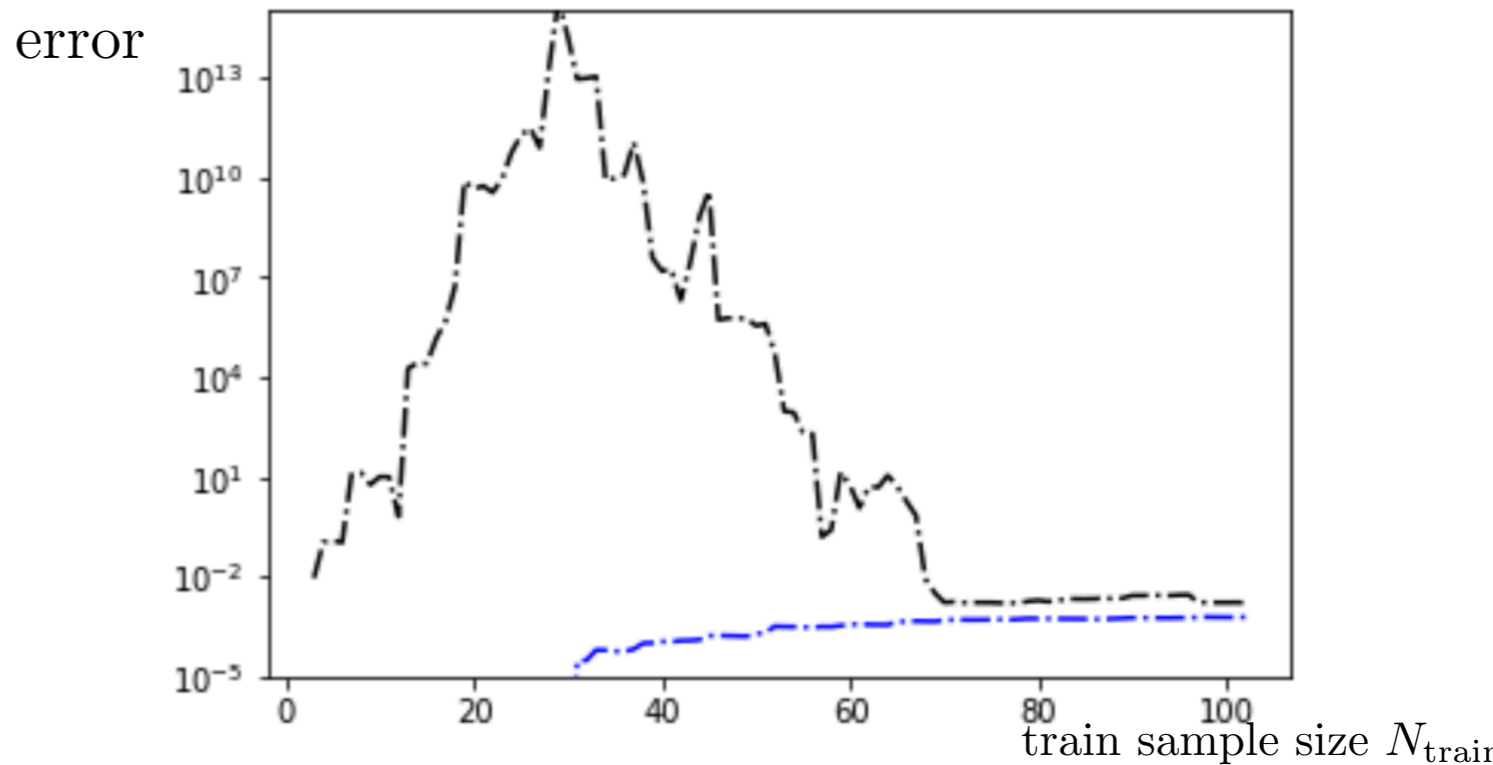
zero mean

$$= \frac{x[1]^2}{n^2\sigma^4} n \left((w_2)^2 \sigma^4 + \sigma^4 \right) = \frac{((w_2)^2 + 1)x[1]^2}{n}$$

- and $\mathbb{E}[\text{variance}] = \frac{((w_2)^2 + 1)\sigma^2}{n}$

- this decrease with (training) sample size n

- **Analysis** explains the empirical observation on the right, on error vs. sample size



Simple Gaussian example

- **Example 2 (moderate predictor):**

- consider a **moderate** model of fitting the first two features

$$\hat{y} = \hat{w}_1 x[1] + \hat{w}_2 x[2]$$

- then, we will show that bias is smaller (in fact zero) and variance is larger, i.e.

- $\mathbb{E}_{p_x}[\text{bias}^2] = 0$

- $\mathbb{E}[\text{variance}] = \frac{2\sigma^2}{n}$

- linear least squares gives

$$\begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix} = (\mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \mathbf{X}[:, \mathbf{1} : \mathbf{2}])^{-1} \mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \mathbf{y}$$

$$\text{with } \mathbf{X}[:, \mathbf{1} : \mathbf{2}] = \begin{bmatrix} x_1[1] & x_1[2] \\ \vdots & \vdots \\ x_n[1] & x_n[2] \end{bmatrix}, \text{ and } \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = w_1 \mathbf{X}[:, \mathbf{1}] + w_2 \mathbf{X}[:, \mathbf{2}] + \underbrace{\begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}}_{\triangleq \varepsilon}$$

- plugging in \mathbf{y} in the equation, we get

$$\begin{aligned} \begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix} &= (\mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \mathbf{X}[:, \mathbf{1} : \mathbf{2}])^{-1} \mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T (\mathbf{X}[:, \mathbf{1} : \mathbf{2}] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \varepsilon) \\ &= \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + (\mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \mathbf{X}[:, \mathbf{1} : \mathbf{2}])^{-1} \mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \varepsilon \end{aligned}$$

Simple Gaussian example

- for large enough n ,

$$\mathbf{X}[:, \mathbf{1}]^T \mathbf{X}[:, \mathbf{1}] = \begin{bmatrix} \sum_{i=1}^n x_i[1]^2 & \sum_{i=1}^n x_i[1]x_i[2] \\ \sum_{i=1}^n x_i[1]x_i[2] & \sum_{i=1}^n x_i[2]^2 \end{bmatrix} \simeq n\sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix},$$

by law of large numbers, and we will use this approximation to simplify the formula:

$$\begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \frac{1}{n\sigma^2} \mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \varepsilon \quad \text{and}$$

$$\mathbb{E} \begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

where we used the fact that $\mathbf{X}[:, \mathbf{1}]$, $\mathbf{X}[:, \mathbf{2}]$, $\varepsilon \in \mathbb{R}^n$ are independent zero-mean vectors

- we are ready to compute the **bias** and **variance**
- first, conditioned on $x = (x[1], x[2], x[3])$

$$\text{bias}^2 = (\mathbb{E}[f(x)] - f_0(x))^2 = (w_1x[1] + w_2x[2] - (w_1x[1] + w_2x[2]))^2 = 0,$$

- note that

- this is an unbiased predictor

Simple Gaussian example

- Now for the variance, conditioned on $x = (x[1], x[2], x[3])$,

since we have
$$\begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \frac{1}{n\sigma^2} \mathbf{X}[:, \mathbf{1} : \mathbf{2}]^T \boldsymbol{\varepsilon}$$

and

$$\mathbb{E} \begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

- variance
$$\begin{aligned} &= \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2 | x] \\ &= \mathbb{E} \left[\left(\frac{1}{n\sigma^2} \boldsymbol{\varepsilon}^T \mathbf{X}[:, \mathbf{1} : \mathbf{2}] \begin{bmatrix} x[1] \\ x[2] \end{bmatrix} \right)^2 | x \right] \\ &= \frac{1}{n^2\sigma^4} \left(\mathbb{E} \left[\left(\sum_{i=1}^n (x_i[2]x[2] + x_i[1]x[1])\varepsilon_i \right)^2 \right] \right) \\ &= \frac{(x[1]^2 + x[2]^2)n\sigma^4}{n^2\sigma^4} = \frac{x[1]^2 + x[2]^2}{n} \end{aligned}$$
- and $\mathbb{E}[\text{variance}] = \frac{2\sigma^2}{n}$
- this decrease with (training) sample size n

Simple Gaussian example

- **Example 3 (complex predictor):**

- consider a **complex** model of fitting

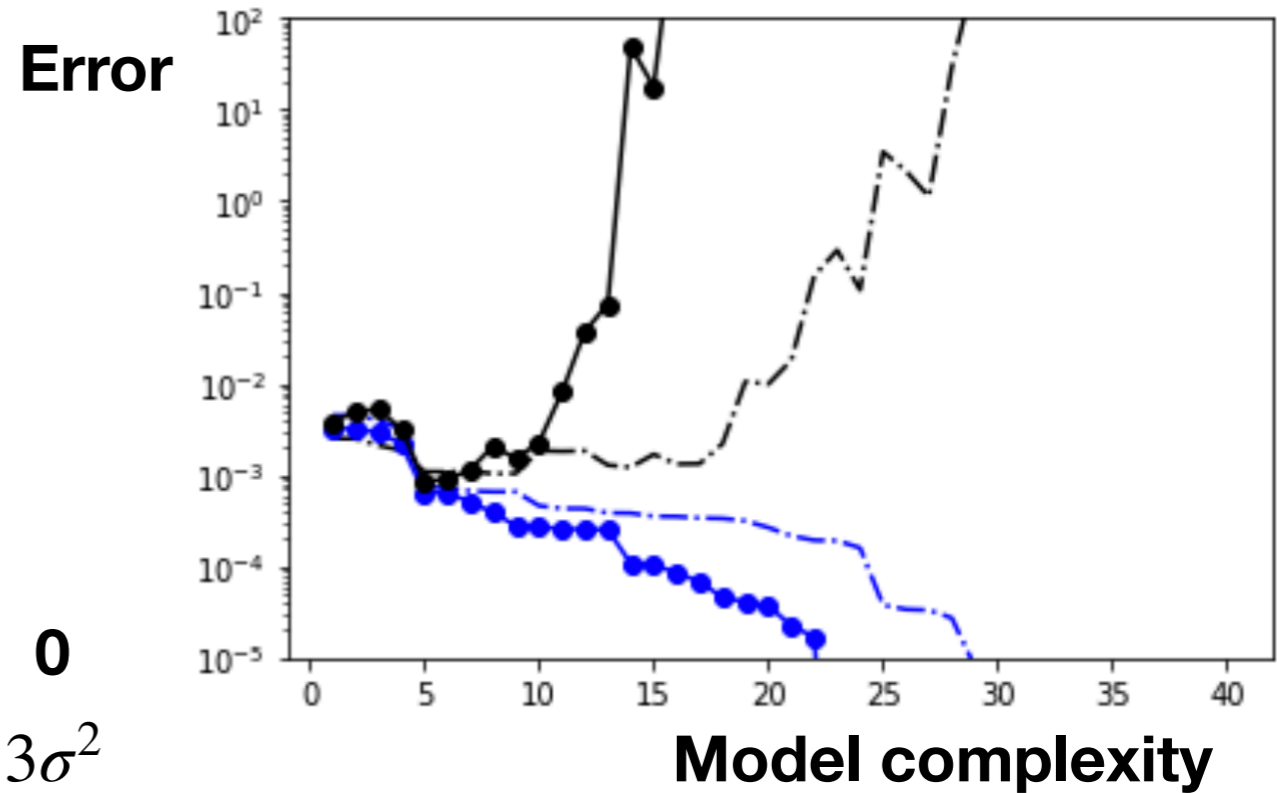
$$\hat{y} = \hat{w}_1 x[1] + \hat{w}_2 x[2] + \hat{w}_2 x[3]$$

- $\mathbb{E}_{p_x}[\text{bias}^2] = 0$

- $\mathbb{E}[\text{variance}] = \frac{3\sigma^2}{n}$

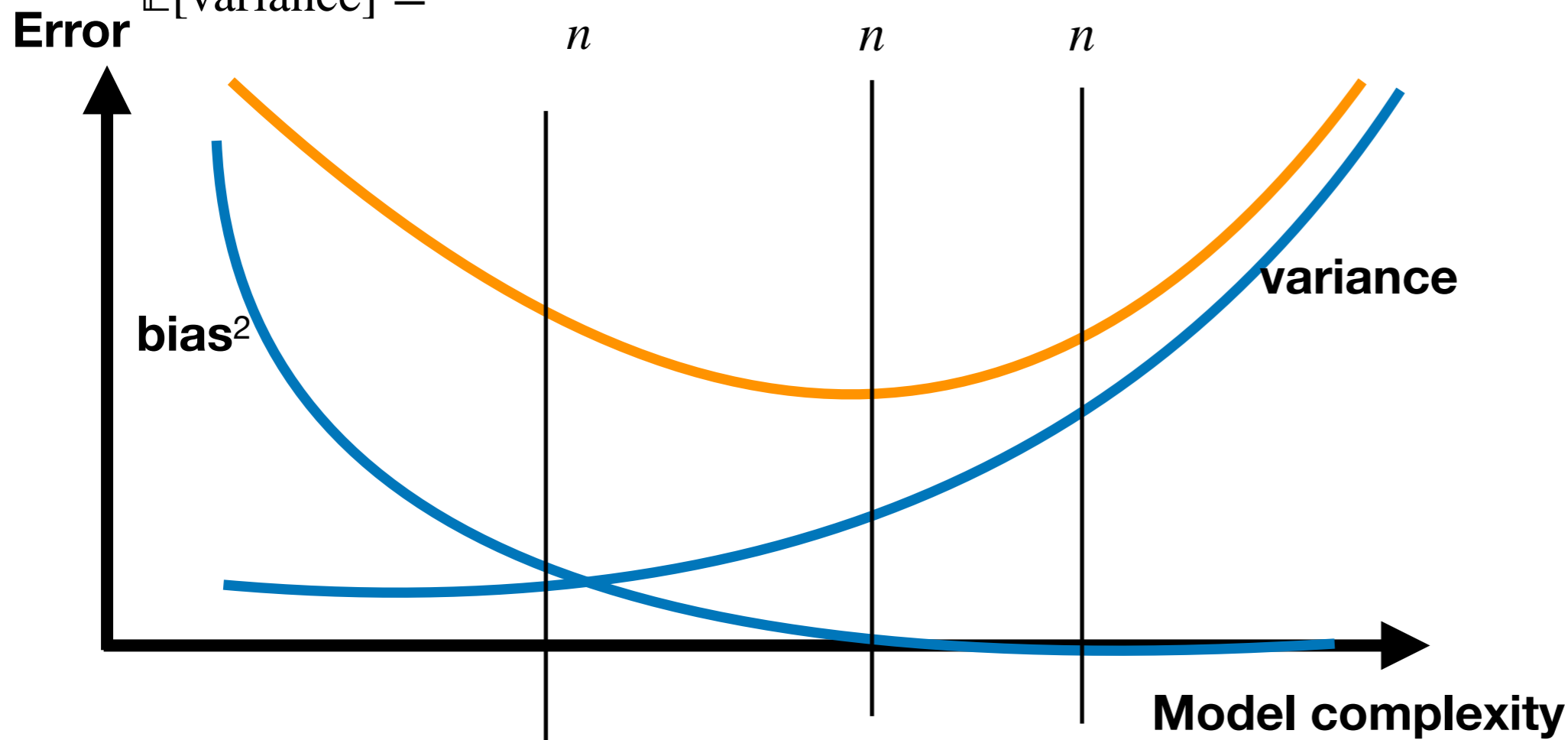
- we skip the detailed proof here, as it is almost identical to the previous one

- this explains the observation on the bias-variance tradeoff



$$\mathbb{E}_{p_x}[\text{bias}^2] = (w_2)^2 \sigma^2$$

$\mathbb{E}[\text{variance}] = \frac{((w_2)^2 + 1)\sigma^2}{n}$	0	0
	$\frac{2\sigma^2}{n}$	$\frac{3\sigma^2}{n}$



* Notations

Model: $P_{xy} \sim (X_i, Y_i)$

example: $y_i = f_0(x_i) + \varepsilon_i$

Sample: $S_{\text{train}} = \{(X_i, Y_i)\}_{i=1}^n$

$|S_{\text{train}}| = n$

$S_{\text{test}} = \{(X_i, Y_i)\}_{i=n+1}^{n+m}$

$|S_{\text{test}}| = m$

Expectation: $\mathbb{E}_{(x,y) \sim P_{xy}} [F(x,y)] = \mathbb{E}_{P_{xy}} [F(x,y)] = \mathbb{E} [F(x,y)]$

$\mathbb{E}_{y \sim P_{y|x}} [F(x,y)|x] = \mathbb{E} [F(x,y)|x]$

* Goal

(Expected) test error: $L_{\text{test}} = \frac{1}{m} \sum_{i=n+1}^{n+m} (f_{S_{\text{train}}}(x_i) - y_i)^2$

true error: $L_{\text{true}} = \mathbb{E}_{S_{\text{test}} \sim P_{xy}^m, S_{\text{train}} \sim P_{xy}^n} [L_{\text{test}}]$

$= \mathbb{E} \left[\frac{1}{m} \sum_{i=n+1}^{n+m} (f_{S_{\text{train}}}(x_i) - y_i)^2 \right]$

$= \frac{1}{m} \sum_{i=n+1}^{n+m} \mathbb{E}_{\substack{(x_i, y_i) \sim P_{xy}, \\ S_{\text{train}}}} \left[(f_{S_{\text{train}}}(x_i) - y_i)^2 \right]$

$$\mathbb{E}_{P_{xy, \text{Serain}}} \left[(f_{\text{Serain}}(x) - y)^2 \right]$$

$$= \mathbb{E}_{P_x} \left[\underbrace{\mathbb{E}_{P_{y|x, \text{Serain}}} \left[(f_{\text{Serain}}(x) - y)^2 \mid x \right]}_{\text{Variance} + \text{Bias}^2 + \sigma^2} \right]$$

Variance + Bias² + σ^2

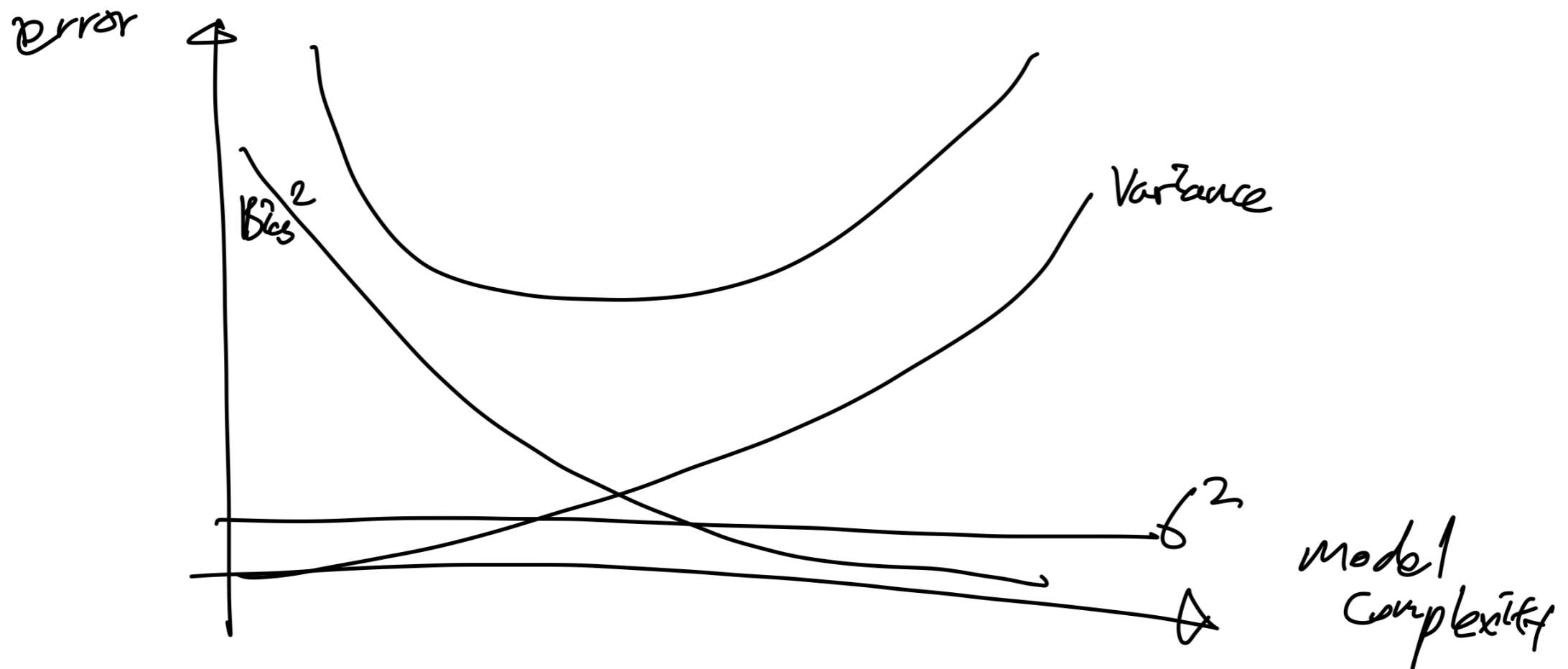
model: $y = f_0(x) + \varepsilon$

(conditional) true error $\Rightarrow \mathbb{E}_{P_{y|x, \text{Serain}}} \left[\underbrace{(f_{\text{Serain}}(x) - f_0(x))}_A \underbrace{- \varepsilon}_B \right]^2 \mid x$

$$= \underbrace{\mathbb{E} \left[\underbrace{(f_{\text{Serain}}(x) - f_0(x))}_A^2 \mid x \right]}_{\text{learning error} \geq 0} + \underbrace{\mathbb{E} \left[\underbrace{\varepsilon^2}_B \mid x \right]}_{\text{irreducible error}} + \underbrace{\mathbb{E} \left[-2AB \right]}_{\square \times \mathbb{E}[\varepsilon] = 0}$$

$$\begin{aligned}
 \text{learning error} &= \mathbb{E} \left[(f_{\text{train}}(x) - f_0(x))^2 \mid x \right] \\
 &= \mathbb{E} \left[\underbrace{(f_{\text{train}}(x) - \mathbb{E}[f_{\text{train}}(x) \mid x])^2}_A + \underbrace{(\mathbb{E}[f_{\text{train}}(x) \mid x] - f_0(x))^2}_B \mid x \right] \\
 &= \underbrace{\mathbb{E} \left[(f_{\text{train}}(x) - \hat{f}(x))^2 \mid x \right]}_{\text{Variance}} + \underbrace{(\hat{f}(x) - f_0(x))^2}_{\text{Bias}^2}
 \end{aligned}$$

$$\mathcal{L}_{\text{true}} = \sigma^2 + \mathbb{E}_{P_x} [\text{Variance}] + \mathbb{E}_{P_x} [\text{Bias}^2]$$



*Simple Gaussian Example

$$\text{model: } y_i = w_1 x_i[1] + w_2 x_i[2] + 0 \cdot x_i[3] + \epsilon_i$$

$$x_i \in \mathbb{R}^3, x_i[1], x_i[2], x_i[3], \epsilon_i \sim N(0, \sigma^2)$$

$$\text{Set of } n: \{(x_i, y_i)\}_{i=1}^n$$

Data Matrix $X \in \mathbb{R}^{n \times 3}$

$$X = \begin{bmatrix} x_1[1] & x_1[2] & x_1[3] \\ \vdots & \vdots & \vdots \\ x_n[1] & x_n[2] & x_n[3] \end{bmatrix}$$

$X[:,1] \uparrow$

*Example 1: Simple Predictor

$$\hat{y} = \hat{w}_1 \cdot X[1]$$

$$\hat{w}_1 = \underbrace{(X[:,1]^T X[:,1])^{-1} X[:,1]^T \cdot y}_{\text{linear least squares predictor}}$$

$$= w_1 + (X[:,1]^T X[:,1])^{-1} X[:,1]^T (X[:,2] \cdot w_2 + \epsilon)$$

$$\mathbb{E}[\hat{w}_1] = w_1$$

* Bias: $\mathbb{E}_x [(f(x) - f_0(x))^2]$

$$= \mathbb{E}_x [(\underbrace{\mathbb{E}[\hat{w}_1]}_{w_1} \cdot x[1] - (w_1 x[1] + w_2 x[2]))^2]$$

does not decrease with n
 Price we pay for using simple model

$$= \mathbb{E}_x [w_2^2 x[2]^2 + \epsilon^2] = w_2^2 \sigma^2$$

* Variance: $\hat{w}_1 = w_1 + \underbrace{(x[1]^\top x[1])^{-1}}_{\sum_{i=1}^n x_i[1]^2} x[1]^\top (w_2 x[2] + \epsilon)$

$$\sum_{i=1}^n x_i[1]^2 \approx n \cdot \mathbb{E}[x_i[1]^2] = n \cdot \sigma^2$$

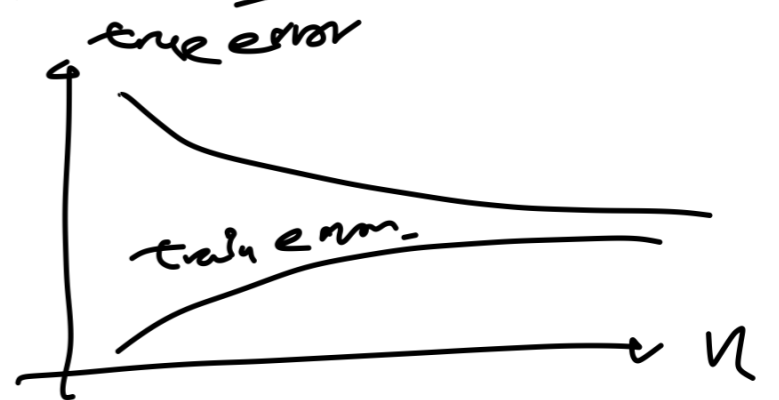
$$= \mathbb{E} [(w_1 x[1] - \hat{w}_1 x[1])^2]$$

$$= \mathbb{E} [\frac{1}{(n\sigma^2)^2} (x[1]^\top x[2] \cdot w_2 + x[1]^\top \epsilon)^2]$$

$$= \frac{\sigma^2}{n^2 \sigma^4} \left\{ w_2^2 \mathbb{E} \left[\sum_{i=1}^n x_i[1] x_i[2] \right]^2 \right\} + \mathbb{E} \left[\left(\sum_{i=1}^n x_i[1] \epsilon \right)^2 \right] \right\}$$

$$= \frac{(w_2^2 + 1) \sigma^2}{n}$$

goes down with n .
 much smaller.



* Example 2: Moderate model.

$$\hat{y} = \hat{w}_1 X[1] + \hat{w}_2 X[2]$$

claim: $\mathbb{E}[\text{bias}^2] = 0$
 $\mathbb{E}[\text{variance}] = \frac{2 \cdot 6^2}{n}$

* Example 3:

$$\hat{y} = \hat{w}_1 X[1] + \hat{w}_2 X[2] + \hat{w}_3 X[3]$$

claim: $\mathbb{E}[\text{bias}^2] = 0$
 $\mathbb{E}[\text{variance}] = \frac{36^3}{n}$

