

Linear Regression, continued



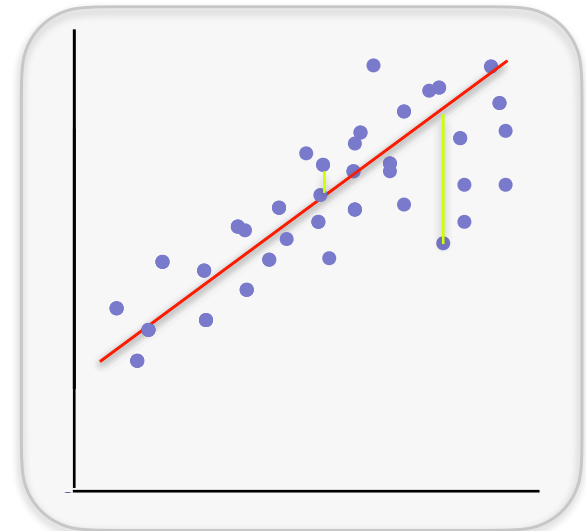
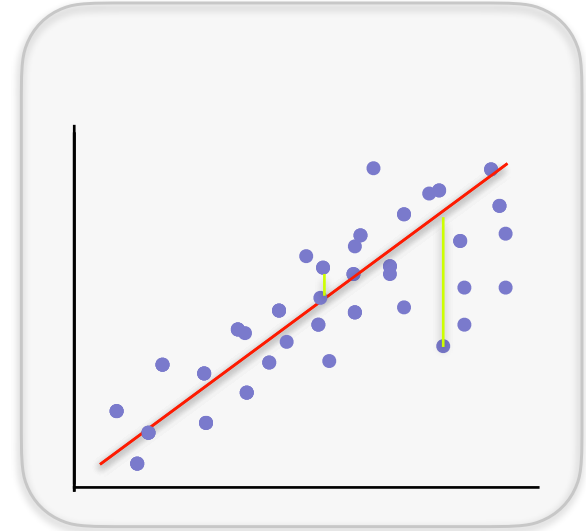
The regression problem in matrix notation

Linear model: $y_i = x_i^T w + \epsilon_i$

Least squares solution:

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

What about an offset
(a.k.a intercept)?



The regression problem in matrix notation

Linear model: $y_i = x_i^T w + \epsilon_i$

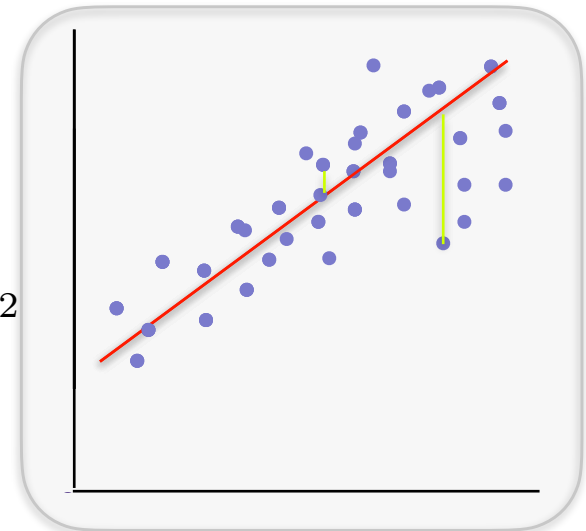
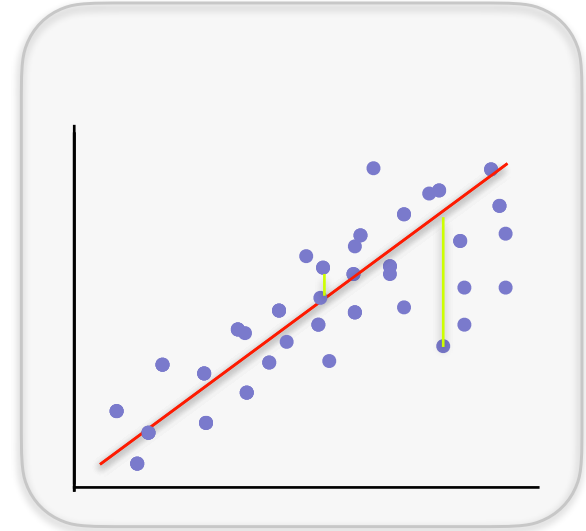
Least squares solution:

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Affine model: $y_i = x_i^T w + b + \epsilon_i$

Least squares solution:

$$\begin{aligned}\hat{w}_{LS}, \hat{b}_{LS} &= \arg \min_{w,b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2 \\ &= \arg \min_{w,b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2\end{aligned}$$



Dealing with an offset

$$\begin{aligned}\widehat{w}_{\text{LS}}, \widehat{b}_{\text{LS}} &= \arg \min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2 \\ &= \arg \min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \underbrace{(\mathbf{y} - (\mathbf{X}w + \mathbf{1}b))^T (\mathbf{y} - (\mathbf{X}w + \mathbf{1}b))}_{\mathcal{L}(w,b)}\end{aligned}$$

Set gradient w.r.t. w and b to zero to find the minima:

A reminder on vector calculus

$$f(\gamma) = (\Omega\gamma + \beta)^T (\Omega\gamma + \beta) \implies \nabla_{\gamma} f(\gamma) = 2\Omega^T(\Omega\gamma + \beta)$$

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

If $\mathbf{X}^T \mathbf{1} = 0$, if the features have zero mean,

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{b}_{LS} = (\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T \mathbf{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

If $\mathbf{X}^T \mathbf{1} = 0$,

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

If $\mathbf{X}^T \mathbf{1} = 0$,

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

$$\mu = \frac{1}{n} \mathbf{X}^T \mathbf{1}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}\mu^T$$

$$\hat{w}_{LS} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i - \mu^T \hat{w}_{LS}$$

Process for linear regression with intercept

Collect data: $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Decide on a **model**: $y_i = x_i^T w + b + \epsilon_i$

Choose a loss function - least squares

Pick the function which minimizes loss on data

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2$$

Use function to make prediction on new examples x_{new}

$$\hat{y}_{\text{new}} = x_{\text{new}}^T \hat{w}_{LS} + \hat{b}_{LS}$$

Another way of dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

reparametrize the problem as $\overline{\mathbf{X}} = [\mathbf{X}, \mathbf{1}]$ and $\overline{w} = \begin{bmatrix} w \\ b \end{bmatrix}$

$$\overline{\mathbf{X}} \overline{w} =$$

Why do we use least squares (i.e. ℓ_2 -loss)?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

$$\implies y_i \sim$$

$$\implies P(y_i; x_i, w, \sigma) =$$

Why do we use least squares (i.e. ℓ_2 -loss)?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Consider $y_i = x_i^T w + \epsilon_i$ where $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$\implies y_i \sim$

$\implies P(y_i; x_i, w, \sigma) =$

Why do we use least squares (i.e. ℓ_2 -loss)?

Maximum Likelihood Estimator:

$$\begin{aligned}\hat{w}_{\text{MLE}} &= \arg \max_w \log P(\{y_i\}_{i=1}^n; \{x_i\}_{i=1}^n, w, \sigma) \\ &= \arg \max_w -n \log(\sigma \sqrt{2\pi}) + \sum_{i=1}^n -\frac{(y_i - x_i^T w)^2}{2\sigma^2}\end{aligned}$$

Why do we use least squares (i.e. ℓ_2 -loss)?

Maximum Likelihood Estimator:

$$\begin{aligned}\hat{w}_{\text{MLE}} &= \arg \max_w \log P(\{y_i\}_{i=1}^n; \{x_i\}_{i=1}^n, w, \sigma) \\ &= \arg \max_w -n \log(\sigma \sqrt{2\pi}) + \sum_{i=1}^n -\frac{(y_i - x_i^T w)^2}{2\sigma^2} \\ &= \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2\end{aligned}$$

Recall: $\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$

$$\hat{w}_{LS} = \hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Recap of linear regression

Data $\{(x_i, y_i)\}_{i=1}^n$

Minimize the loss (Empirical Risk Minimization)

Choose a loss

e.g., ℓ_2 -loss: $(y_i - x_i^T w)^2$

Solve $\widehat{w}_{\text{LS}} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$

Maximize the likelihood (MLE)

Choose a Hypothesis class

e.g., $y_i = x_i^T w + \epsilon_i$, $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Maximize the likelihood,

$\widehat{w}_{\text{MLE}} = \arg \max_w \left\{ -n \log(\sigma \sqrt{2\pi}) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2} \right\}$

Analysis of **Error** under additive Gaussian noise

Let's suppose $y_i = x_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, then this can be written as

$$\mathbf{y} = \mathbf{X}w^* + \epsilon$$

$$\begin{aligned}\widehat{w}_{\text{MLE}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \epsilon) \\ &= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Maximum Likelihood Estimator is unbiased:

Analysis of **Error** under additive Gaussian noise

Let's suppose $y_i = x_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, then this can be written as
 $\mathbf{y} = \mathbf{X}w^* + \epsilon$

$$\begin{aligned}\widehat{\mathbf{w}}_{\text{MLE}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \epsilon) \\ &= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Covariance is:

Analysis of **Error** under additive Gaussian noise

Let's suppose $y_i = x_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, then this can be written as $\mathbf{y} = \mathbf{X}w^* + \epsilon$, and the MLE is

$$\hat{w}_{\text{MLE}} = w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

This random estimate has the following distribution:

$$\mathbb{E}[\hat{w}_{\text{MLE}}] = w^*, \text{Cov}(\hat{w}_{\text{MLE}}) = \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])(\hat{w} - \mathbb{E}[\hat{w}])^T] = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$$

$$\hat{w}_{\text{MLE}} \sim \mathcal{N}(w^*, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$$

Interpretation: consider an example with $\mathbf{x} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ -1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{bmatrix}$

The covariance of the MLE, $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$, captures how each sample gives information about the unknown w^* , but each sample gives information about for different (linear combination of) coordinates and of different quality/strength

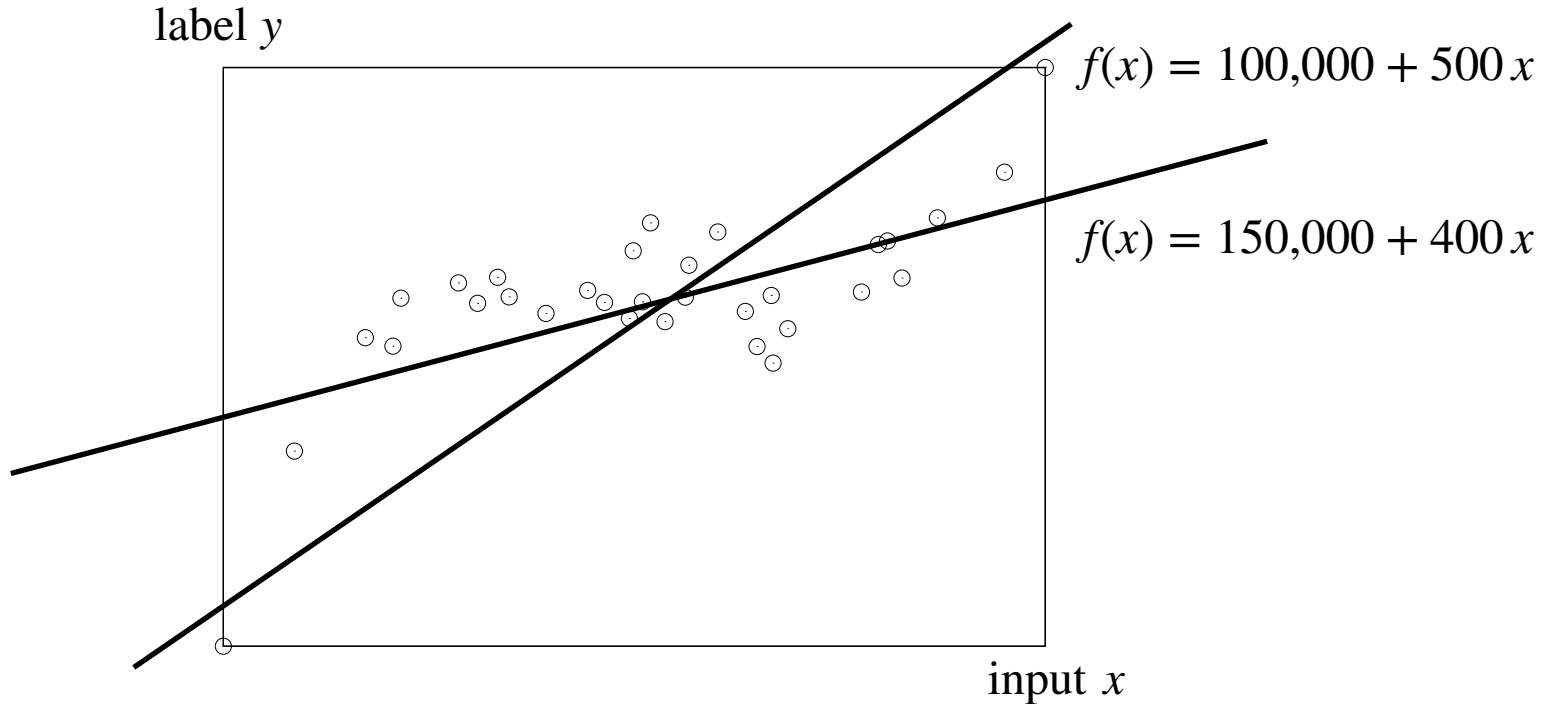
Questions?

Polynomial regression

- How to fit more complex data?



Recap: Linear Regression



- In general high-dimensions, we fit a linear model with intercept $y_i \simeq w^T x_i + b$, or equivalently $y_i = w^T x_i + b + \epsilon_i$ with model parameters $(w \in \mathbb{R}^d, b \in \mathbb{R})$ that minimizes ℓ_2 -loss

$$\mathcal{L}(w, b) = \sum_{i=1}^n \underbrace{(y_i - (w^T x_i + b))^2}_{\text{error } \epsilon_i}$$

Recap: Linear Regression

- The least squares solution, i.e. the minimizer of the ℓ_2 -loss can be written in a **closed form** as a function of data \mathbf{X} and \mathbf{y} as

As we derived in class:

$$\mu = \frac{1}{n} \mathbf{X}^T \mathbf{1}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}\mu^T$$

$$\hat{w}_{\text{LS}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

$$\hat{b}_{\text{LS}} = \frac{1}{n} \sum_{i=1}^n y_i - \mu^T \hat{w}_{\text{LS}}$$

or equivalently using straightforward linear algebra by setting the gradient to zero:

$$\begin{bmatrix} \hat{w}_{\text{LS}} \\ \hat{b}_{\text{LS}} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{X}^T \\ \mathbf{1}^T \end{bmatrix} [\mathbf{X} \ \mathbf{1}] \right)^{-1} \begin{bmatrix} \mathbf{X}^T \\ \mathbf{1}^T \end{bmatrix} \mathbf{y}$$

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

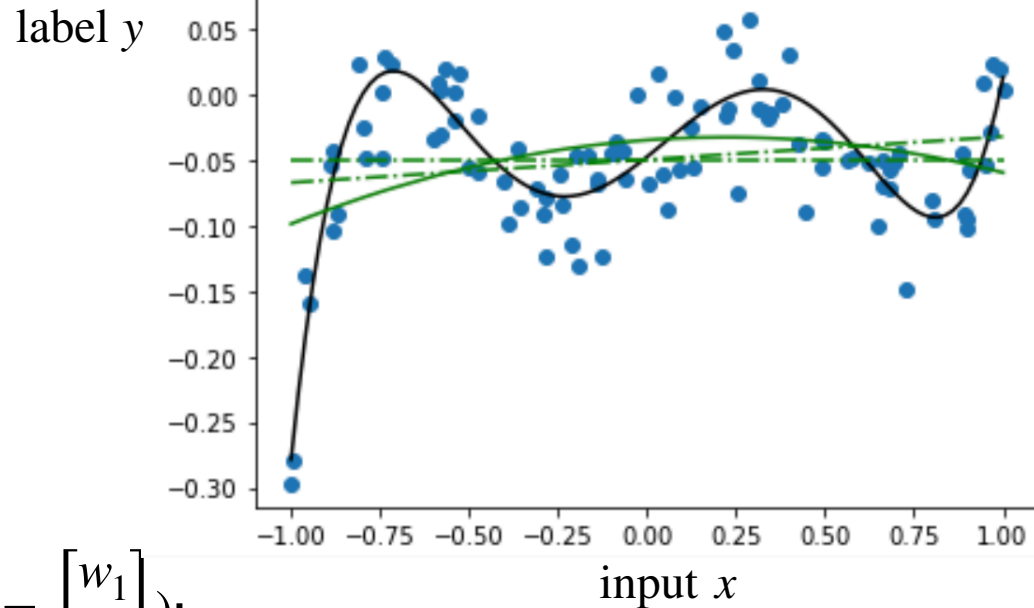
- $y_i = b + w_1 x_i + \epsilon_i$
- $\mathbf{y} = \mathbf{1}b + \mathbf{X}w_1 + \epsilon$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $y_i = b + w_1 x_i + w_2 x_i^2 + \epsilon_i$
- Define $h : \mathbb{R} \rightarrow \mathbb{R}^2$ such that $x \mapsto h(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$
- $y_i = b + h(x_i)^T w + \epsilon_i$

- Treat $h(x)$ as new input features. Let $\mathbf{H} = \begin{bmatrix} h(x_1)^T \\ \vdots \\ h(x_n)^T \end{bmatrix}$. Replace x_i by $\begin{bmatrix} x_i \\ x_i^2 \end{bmatrix}$

- $\mathbf{y} = \mathbf{1}b + \mathbf{H}w + \epsilon$



Degree- p polynomial regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $y_i = b + w_1 x_i + \epsilon_i$
- $\mathbf{y} = \mathbf{1}b + \mathbf{X}w_1 + \epsilon$

- **Degree- p model with parameter $(b, w \in \mathbb{R}^p)$:**

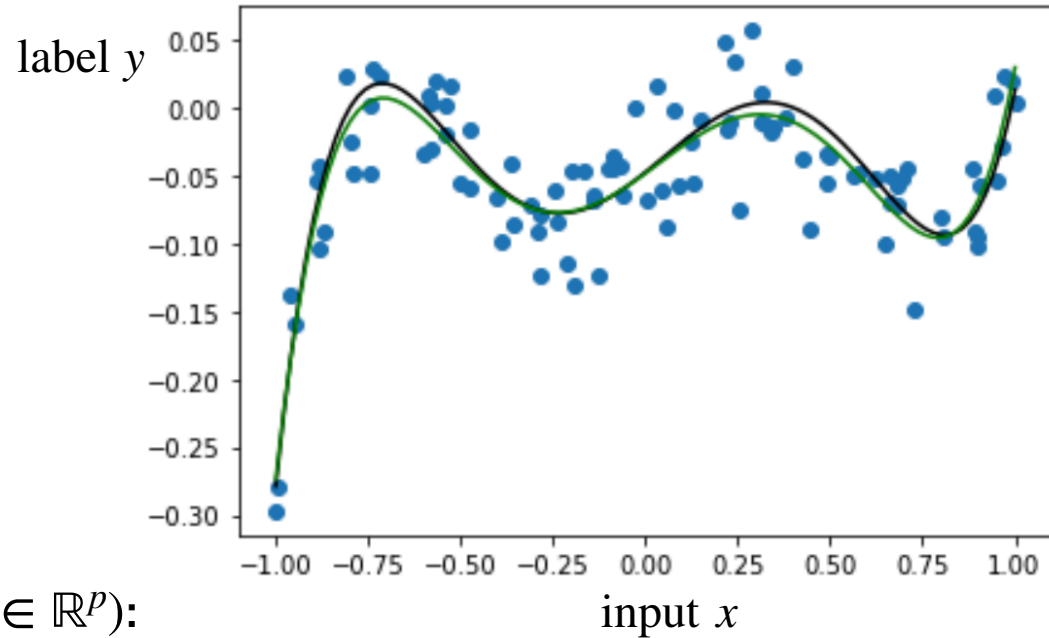
- $y_i = b + w_1 x_i + \dots + w_p x_i^p + \epsilon_i$

- Define $h : \mathbb{R} \rightarrow \mathbb{R}^p$ such that $x \mapsto h(x) = \begin{bmatrix} x \\ \vdots \\ x^p \end{bmatrix}$

- $y_i = b + h(x_i)^T w + \epsilon_i$

- Treat $h(x)$ as new input features and let $\mathbf{H} = \begin{bmatrix} h(x_1)^T \\ \vdots \\ h(x_n)^T \end{bmatrix}$

- $\mathbf{y} = \mathbf{1}b + \mathbf{H}w + \epsilon$



Degree- p polynomial regression in d -dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ & x_2^T & & \\ & \vdots & & \\ & x_n^T & & \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Degree- p model with parameter** ($b, w \in \mathbb{R}^{dp}$):

- $y_i = b + x_i^T w_1 + \cdots + (x_i^p)^T w_p + \epsilon_i$, where $x_i^p = \begin{bmatrix} x_{i1}^p \\ \vdots \\ x_{id}^p \end{bmatrix}$

- Define $h : \mathbb{R}^d \rightarrow \mathbb{R}^{dp}$ such that $x \mapsto h(x) = \begin{bmatrix} x \\ \vdots \\ x^p \end{bmatrix} \in \mathbb{R}^{dp}$

- $y_i = b + h(x_i)^T w + \epsilon_i$

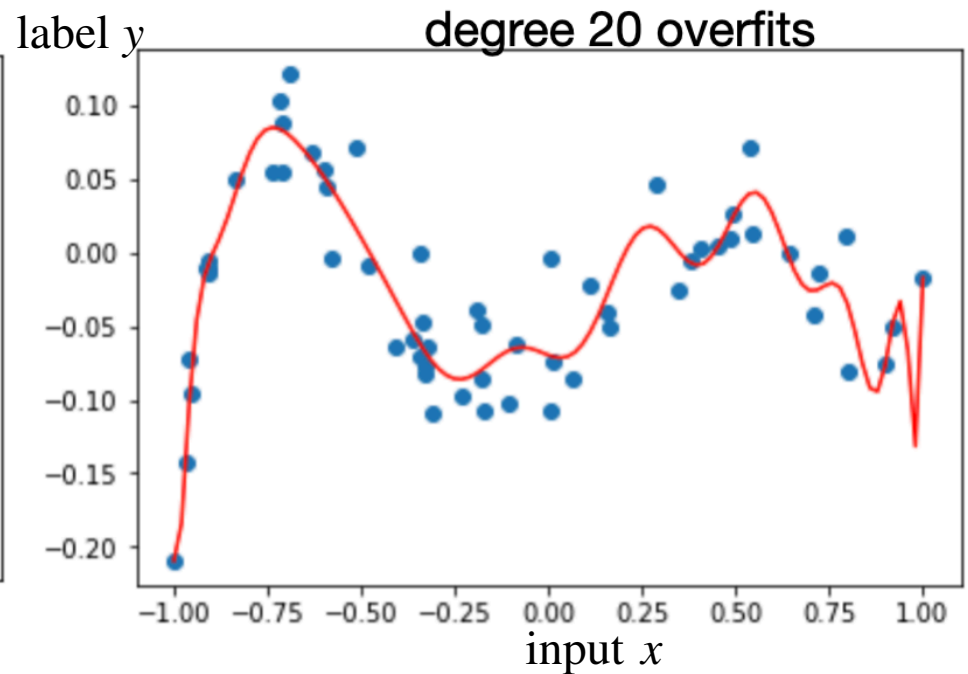
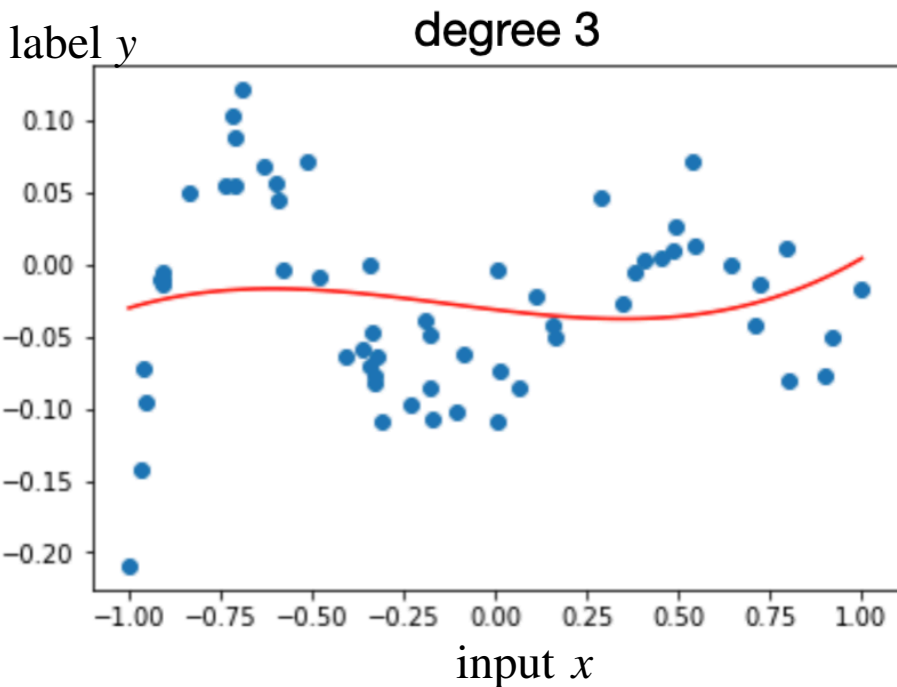
- Treat $h(x)$ as new input features and let $\mathbf{H} = \begin{bmatrix} h(x_1)^T \\ \vdots \\ h(x_n)^T \end{bmatrix}$

- $\mathbf{y} = \mathbf{1}b + \mathbf{H}w + \epsilon$

- In general, any feature $h(x)$ can be used, e.g., $\sin(ax + b)$, $e^{-b(x-a)^2}$, $\log x$, etc.

Which p should we choose?

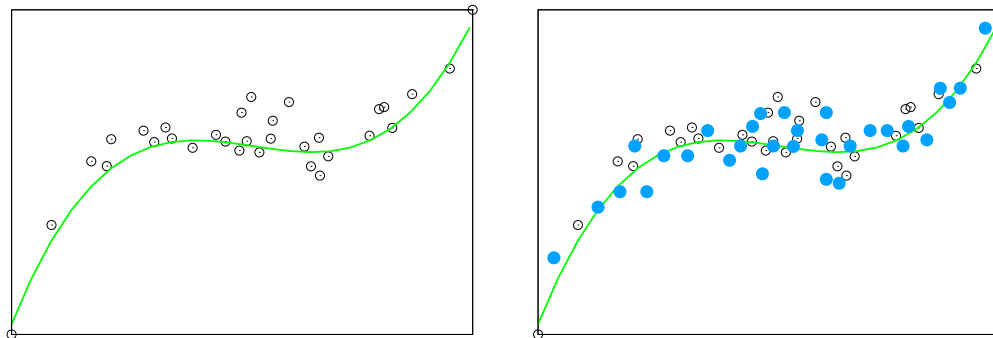
- First instance of class of models with different representation power = model complexity



- How do we determine which is better model?

Generalization

- we say a predictor **generalizes** if it performs as well on unseen data as on training data
- formal mathematical definition involves probabilistic assumptions (coming later in this week)
- 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 performs well on in-sample data but does not perform well on out-of-sample data



- **train** a cubic predictor on 32 (**in-sample**) white circles: Mean Squared Error (MSE) 174
- **predict** label y for 30 (**out-of-sample**) blue circles: MSE 192
- conclude this predictor/model generalizes, as in-sample MSE \simeq out-of-sample MSE

Split the data into **training** and **testing**

- a way to mimic how the predictor performs on unseen data
- 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)

- **training set** used to train the model

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

- **test set** used to evaluate the model

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

- this assumes that test set is similar to unseen data
- **test set should never be used in training**

We say a model w or predictor **overfits** if $\mathcal{L}_{\text{train}}(w) \ll \mathcal{L}_{\text{test}}(w)$

small training error

large training error

small test error

**generalizes well
performs well**

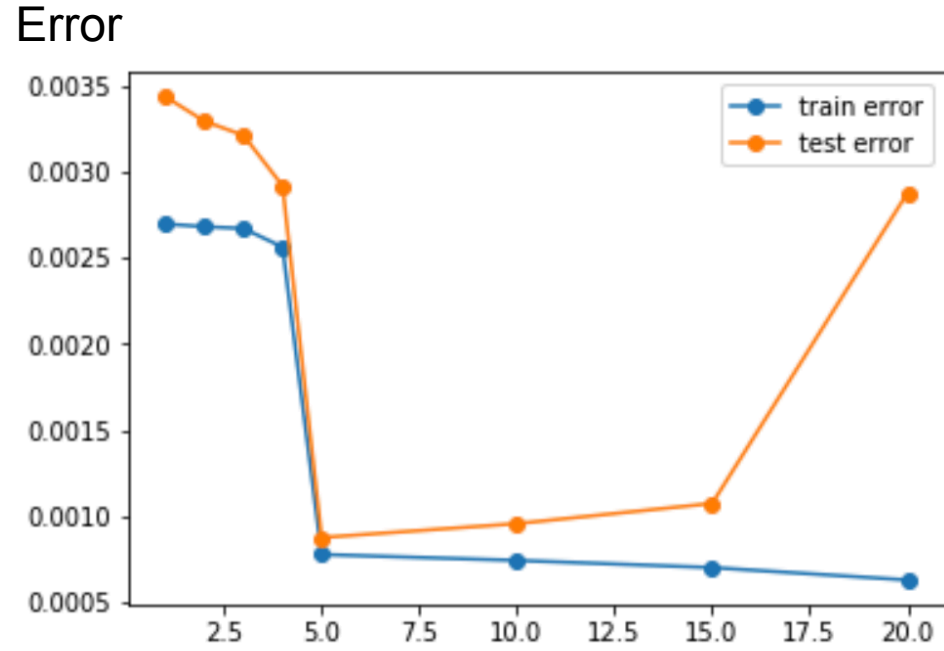
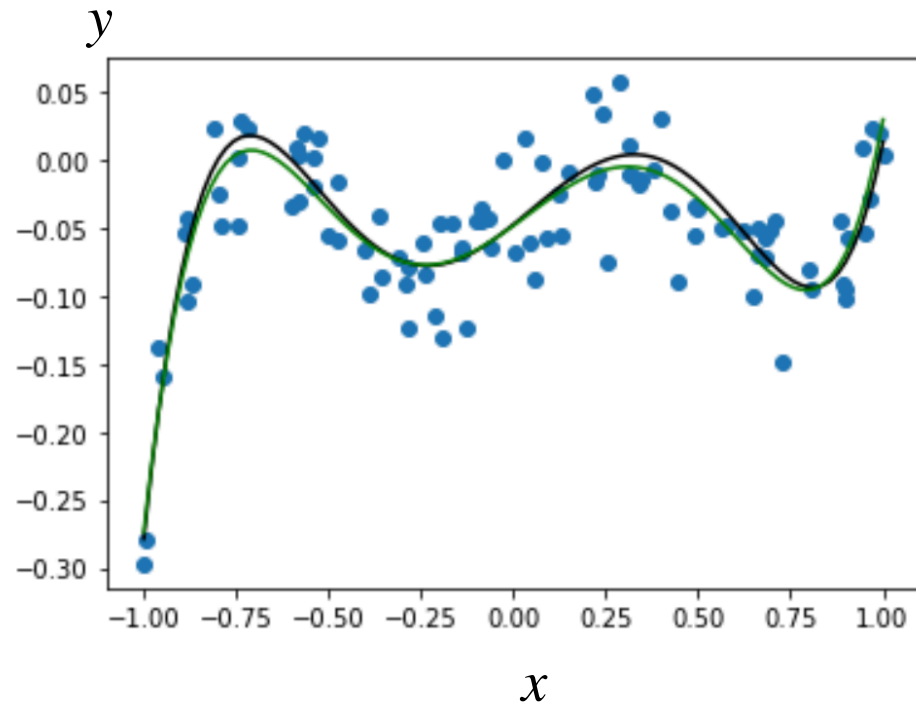
possible, but unlikely

large test error

**fails to generalize
Overfitting**

**generalizes well
performs poorly**

How do we choose which model to use?

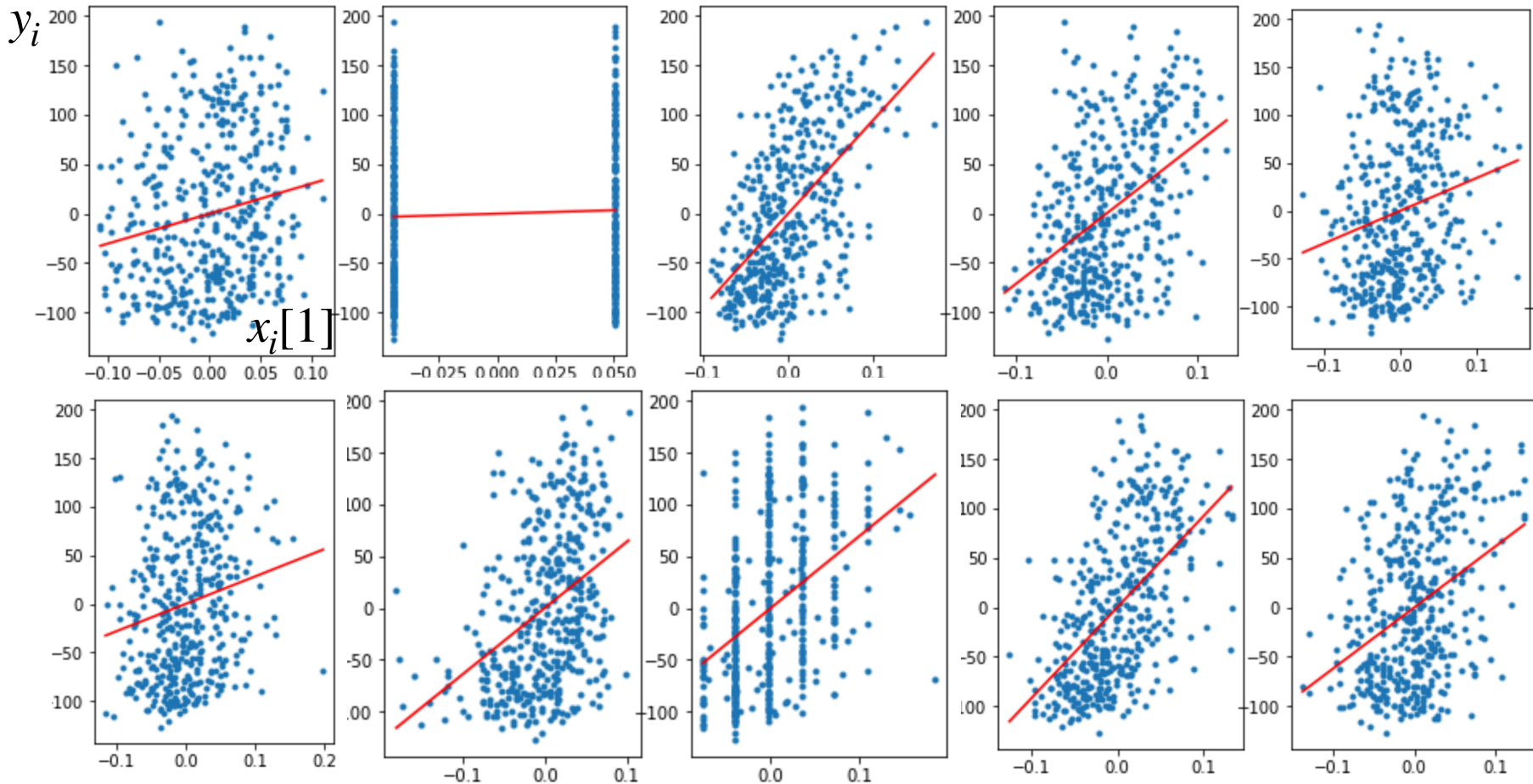



Degree- p polynomial model

1. first use 60 data points to train and 60 data points to test and train several models to get the above graph on the right
2. then choose degree $p = 5$, since it achieves **minimum test error**
3. now re-train on all 120 data points with degree 5 polynomial model

Another example: Diabetes

- Example: Diabetes
 - 10 explanatory variables
 - from 442 patients
 - we use half for train and half for validation





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

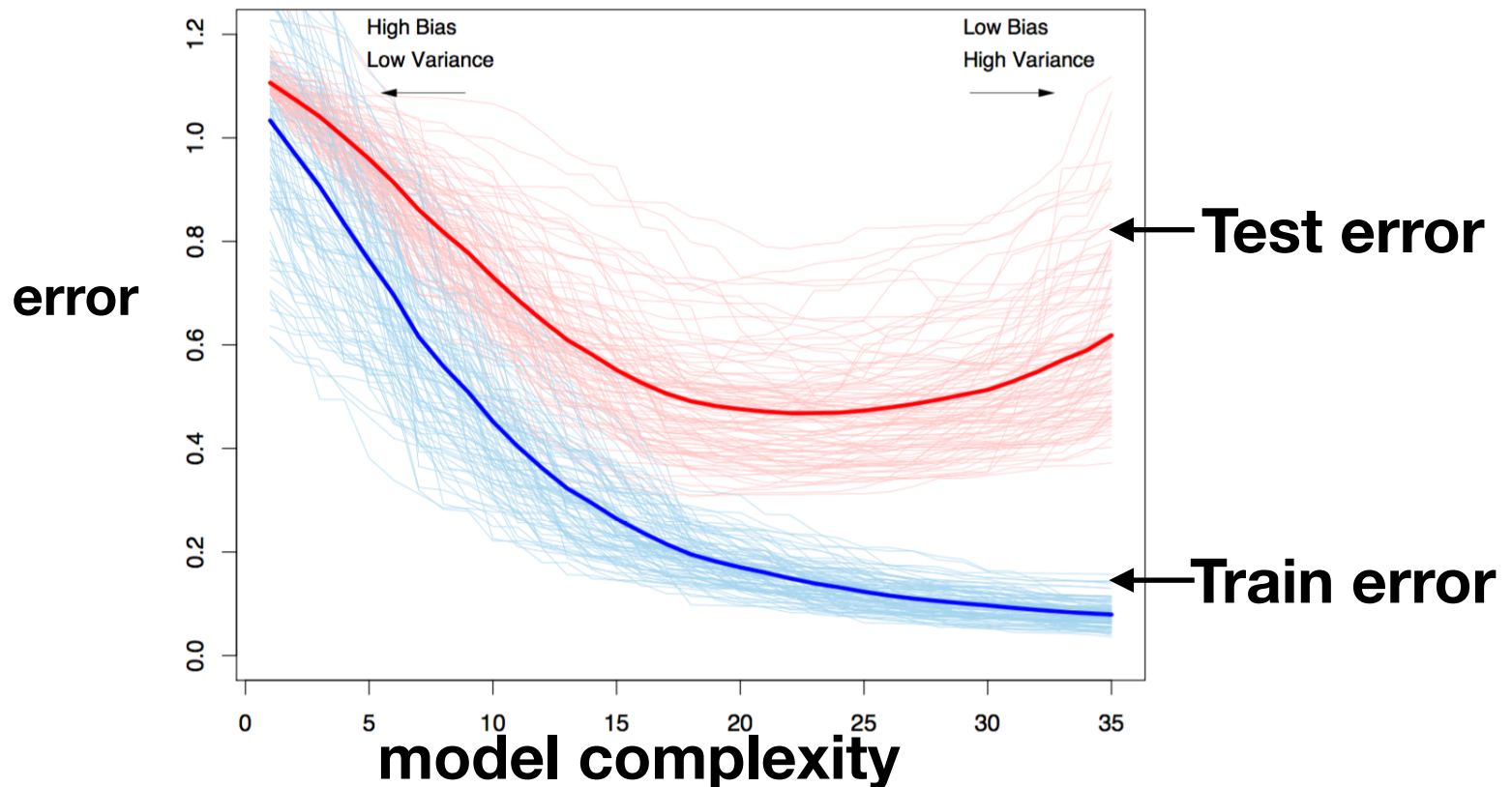
- **test MSE is the primary criteria for model selection**
- Using only 2 features (S5 and BMI), one can get very close to the prediction performance of using all features
- Combining S3 and S4 does not give any performance gain

What does the bias-variance theory tell us?

- **Train error** (random variable, randomness from \mathcal{D})
 - Use $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \sim P_{X,Y}$ to find \widehat{w}
 - Train error: $\mathcal{L}_{\text{train}}(\widehat{w}_{\text{LS}}) = \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - \widehat{w}^T x_i)^2$
- recall the **test error** is an unbiased estimator of the **true error**
- **True error** (random variable, randomness from \mathcal{D})
 - True error: $\mathcal{L}_{\text{true}}(\widehat{w}) = \mathbb{E}_{(x,y) \sim P_{X,Y}} [(y - \widehat{w}^T x)^2]$
- **Test error** (random variable, randomness from \mathcal{D} and \mathcal{T})
 - Use $\mathcal{T} = \{(x_i, y_i)\}_{i=1}^m \sim P_{X,Y}$
 - Test error: $\mathcal{L}_{\text{test}}(\widehat{w}) = \frac{1}{|\mathcal{T}|} \sum_{(x_i, y_i) \in \mathcal{T}} (y_i - \widehat{w}^T x_i)^2$
- theory explains **true error**, and hence expected behavior of the (random) **test error**

What does the bias-variance theory tell us?

- Train error is optimistically biased (i.e. smaller) because the trained model is minimizing the train error
- Test error is unbiased estimate of the true error, if test data is never used in training a model or selecting the model complexity
- Each line is an i.i.d. instance of \mathcal{D} and \mathcal{T}



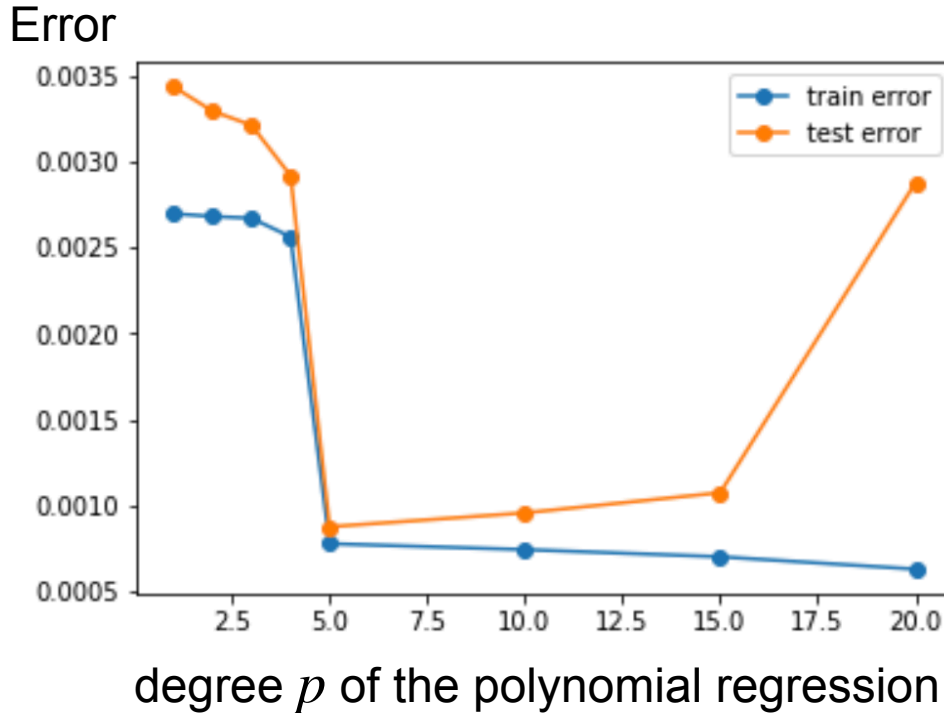
Questions?

Lecture 5: Bias-Variance Tradeoff

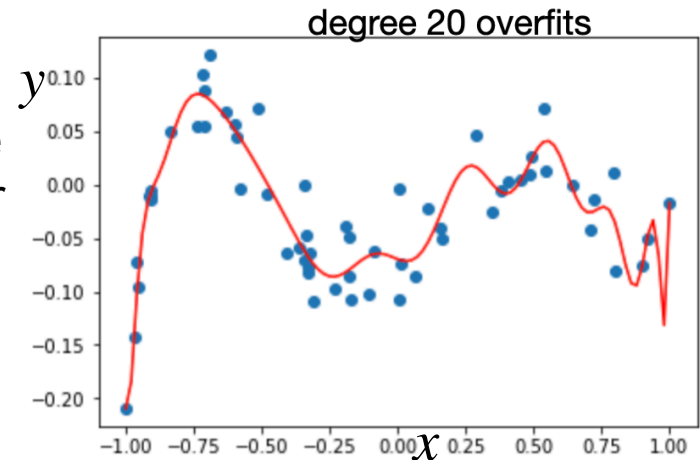
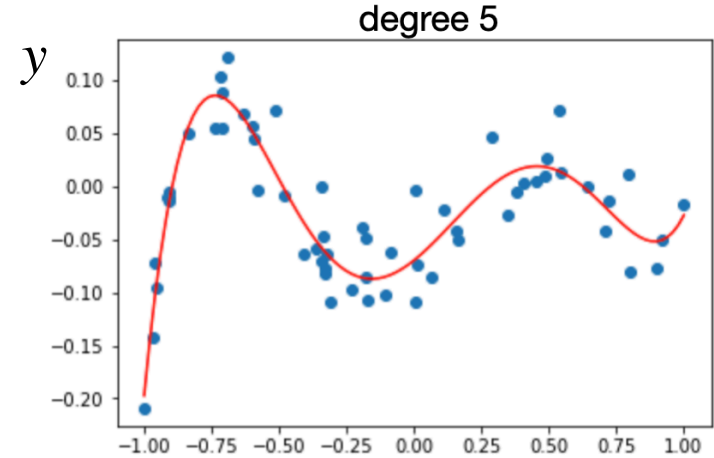
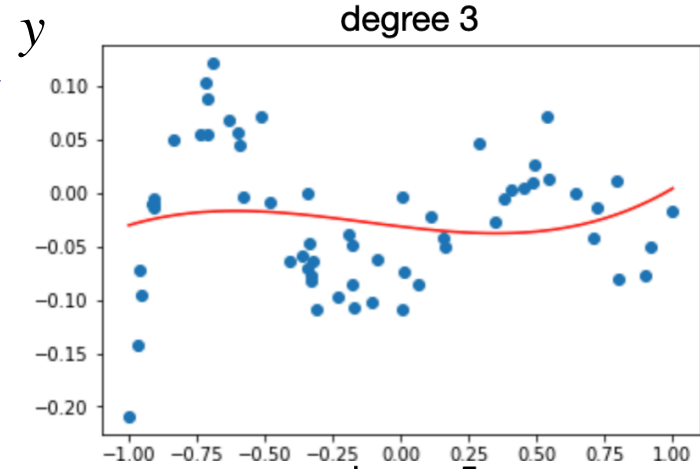
- explaining test error using theoretical analysis



Train/test error vs. complexity



- **Model complexity** e.g., degree p of the polynomial model, number of features used in diabetes example
 - Related to the dimension of the model parameter
- **Train error** monotonically decreases with model complexity
- **Test error** has a U shape



Statistical learning

Typical notation:

X denotes a random variable

x denotes a deterministic instance

- Suppose data is generated from a statistical model $(X, Y) \sim P_{X,Y}$
 - and assume we know $P_{X,Y}$ (just for now to explain statistical learning)
- Then **learning** is to find a predictor $\eta : \mathbb{R}^d \rightarrow \mathbb{R}$ that minimizes
 - the expected error $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2]$
 - think of this random (X, Y) as a new sample you will encounter when you deployed your learned model, and we care about its average performance
- Since, we do not assume anything about the function $\eta(x)$, it can take any value for each $X = x$, hence the optimization can be broken into sum (or more precisely integral) of multiple objective functions, each involving a specific value $X = x$

$$\bullet \quad \mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2] = \mathbb{E}_{X \sim P_X}[\mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 | X = x]]$$

$$= \int \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 | X = x] P_X(x) dx$$

Or for discrete X ,

$$= \sum_x P_X(x) \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 | X = x]$$

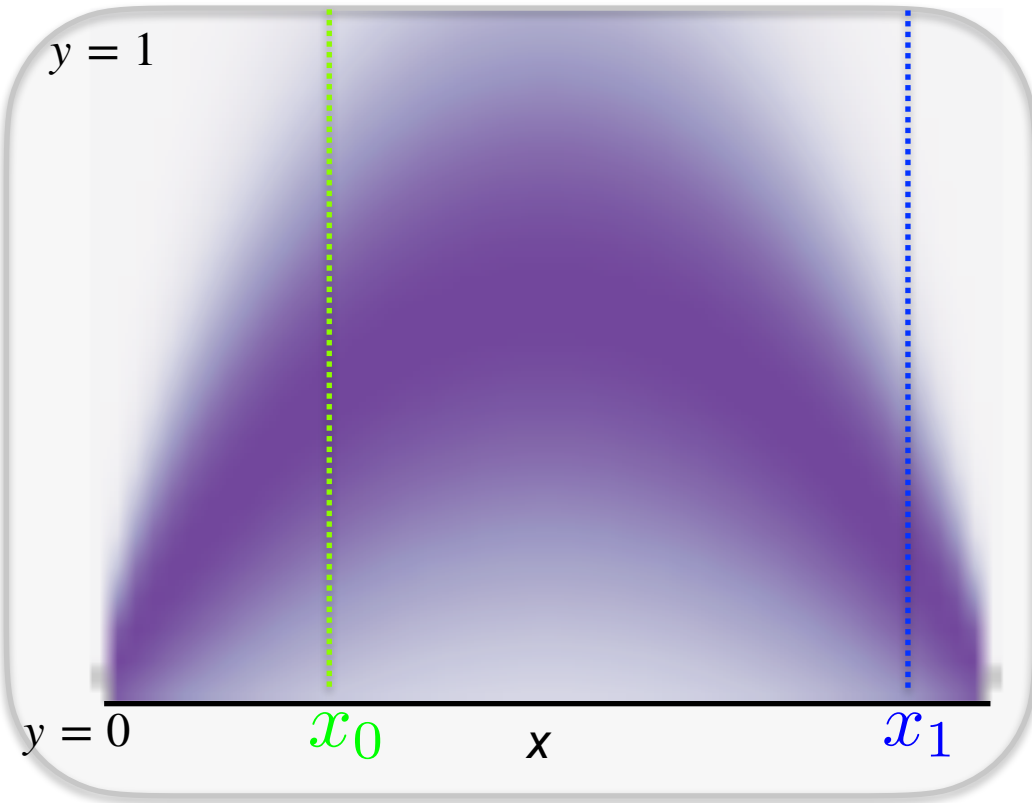
Where we used the chain rule: $\mathbb{E}_{X,Y}[f(X, Y)] = \mathbb{E}_X[\mathbb{E}_{Y|X}[f(x, Y) | X = x]]$

Statistical learning

- We can solve the optimization for each $X = x$ separately
 - $\eta(x) = \arg \min_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 | X = x]$
 - The optimal solution is $\eta(x) = \mathbb{E}_{Y \sim P_{Y|X}}[Y | X = x]$,
which is the best prediction in ℓ_2 -loss/Mean Squared Error
 - Claim: $\mathbb{E}_{Y \sim P_{Y|X}}[Y | X = x] = \arg \min_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 | X = x]$
 - Proof:
-
- Note that this optimal statistical estimator $\eta(x) = \mathbb{E}[Y | X = x]$ cannot be implemented as we do not know $P_{X,Y}$ in practice
 - This is only for the purpose of conceptual understanding

Statistical Learning

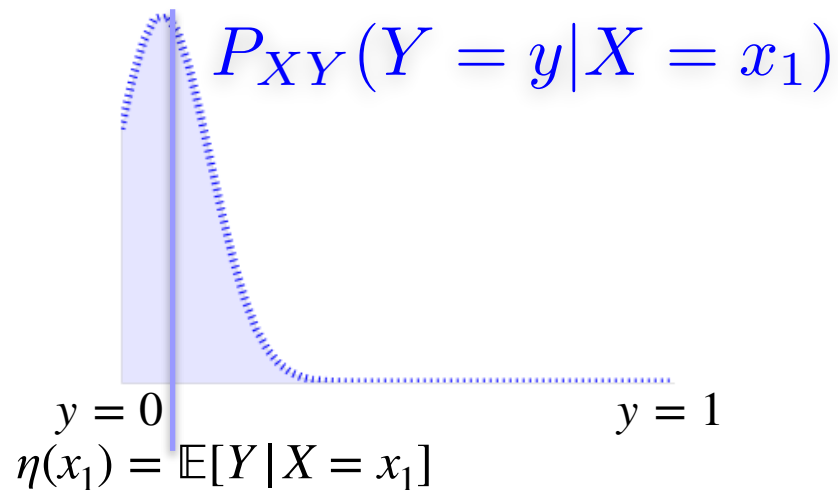
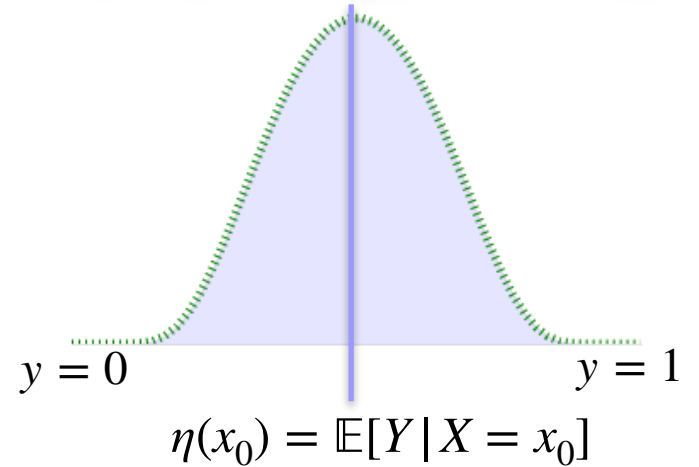
$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

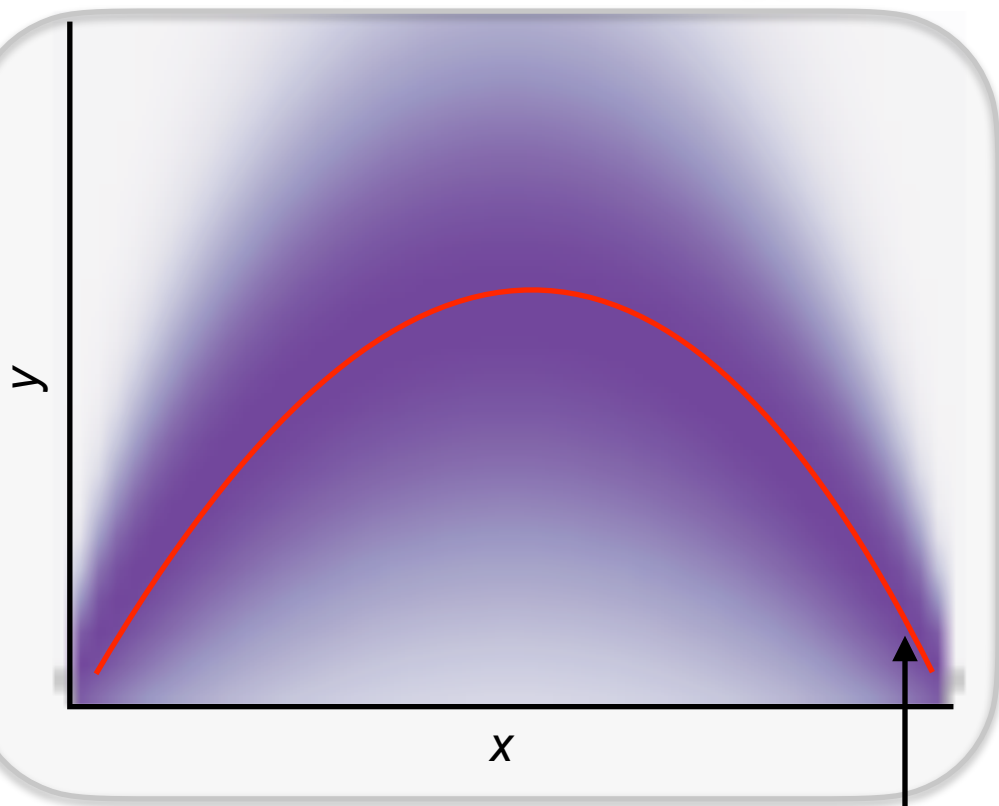
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

$$P_{XY}(Y = y|X = x_0)$$



Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

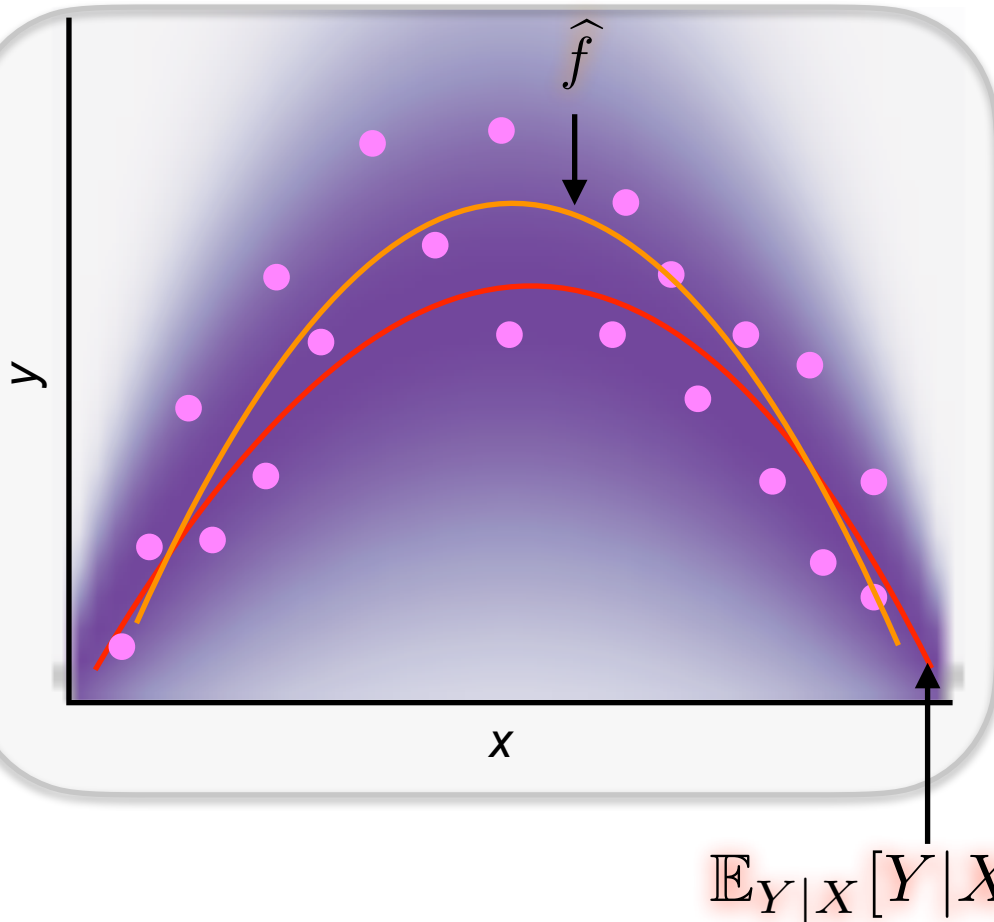
But we do not know $P_{X,Y}$

We only have samples.

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

So we need to restrict our predictor to a function class (e.g., linear, degree- p polynomial) to avoid overfitting:

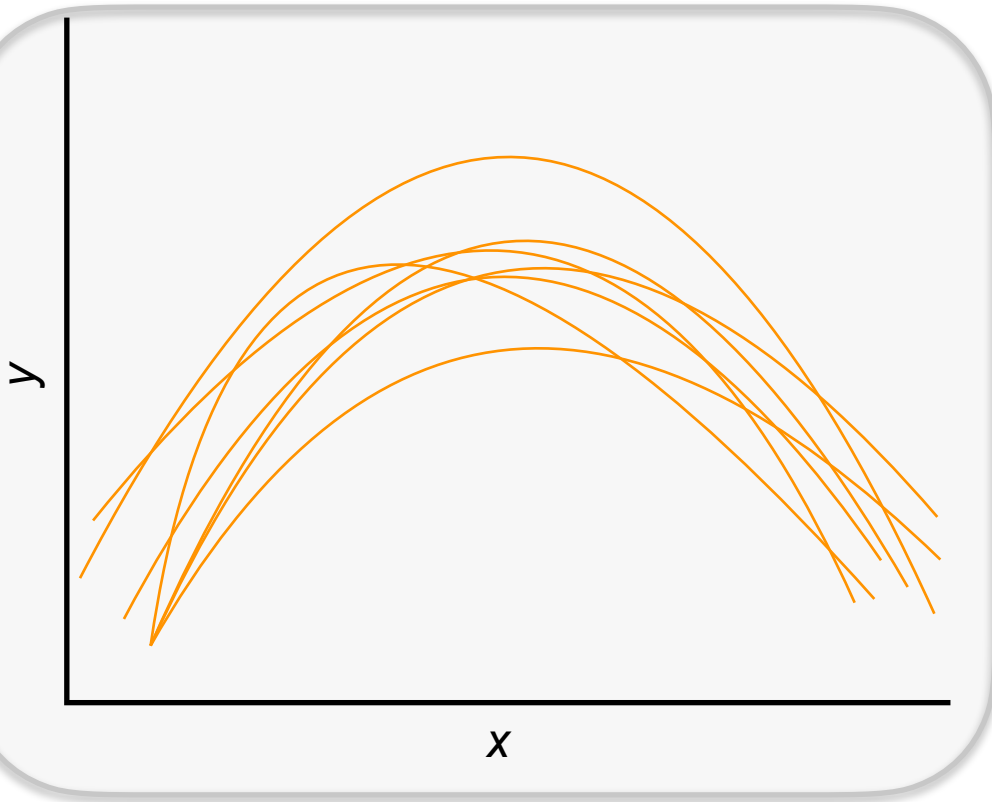
$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

We care about how our predictor performs on future unseen data

$$\text{True Error of } \hat{f}: \mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$$

Future prediction error $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$ is random
because \hat{f} is random (whose randomness comes from training data \mathcal{D})

$$P_{XY}(X = x, Y = y)$$



Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Bias-variance tradeoff

Notation:

I use predictor/model/estimate,
interchangeably

Ideal predictor

$$\eta(x) = \mathbb{E}_{Y|X}[Y | X = x]$$

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- We are interested in the **True Error** of a (random) learned predictor:

$$\mathbb{E}_{X,Y}[(Y - \hat{f}_{\mathcal{D}}(X))^2]$$

- But the analysis can be done for each $X = x$ separately, so we analyze the **conditional true error**:

$$\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathcal{D}}(x))^2 | X = x]$$

- And we care about the **average conditional true error**, averaged over training data:

$$\mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathcal{D}}(x))^2 | X = x] \right]$$

written compactly as $= \mathbb{E}[(Y - \hat{f}_{\mathcal{D}}(x))^2]$

Bias-variance tradeoff

Ideal predictor

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- **Average conditional true error:**

$$\mathbb{E}_{\mathcal{D}, Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2]$$

Bias-variance tradeoff

Ideal predictor

$$\eta(x) = \mathbb{E}_{Y|X} [Y | X = x]$$

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- **Average conditional true error:**

$$\begin{aligned} \mathbb{E}_{\mathcal{D}, Y|x} [(Y - \hat{f}_{\mathcal{D}}(x))^2] &= \mathbb{E}_{\mathcal{D}, Y|x} [(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2] \\ &= \mathbb{E}_{\mathcal{D}, Y|x} \left[(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) + (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2 \right] \\ &= \mathbb{E}_{Y|x} [(Y - \eta(x))^2] + \underbrace{2\mathbb{E}_{\mathcal{D}, Y|x} [(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x))]}_{=0} + \mathbb{E}_{\mathcal{D}} [(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] \end{aligned}$$

(this follows from independence of \mathcal{D} and (X, Y) and

$$\mathbb{E}_{Y|x} [Y - \eta(x)] = \mathbb{E}[Y | X = x] - \eta(x) = 0$$

$$= \underbrace{\mathbb{E}_{Y|x} [(Y - \eta(x))^2]}_{\text{Irreducible error}} + \underbrace{\mathbb{E}_{\mathcal{D}} [(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{Average learning error}}$$

Irreducible error

- (a) Caused by stochastic label noise in $P_{Y|X=x}$
- (b) cannot be reduced

Average learning error

- Caused by
- (a) either using too “simple” of a model or
- (b) not enough data to learn the model accurately

Bias-variance tradeoff

Ideal predictor

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

• **Average learning error:**

$$\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}}\left[\left(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x) \right)^2 \right]$$

Bias-variance tradeoff

Ideal predictor

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

• **Average learning error:**

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] &= \mathbb{E}_{\mathcal{D}} \left[\left(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2 + 2(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x)) \right. \\ &\quad \left. + (\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2 \right] \end{aligned}$$

$$= \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} + \underbrace{\mathbb{E}_{\mathcal{D}} \left[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2 \right]}_{\text{variance}}$$

biased squared

variance

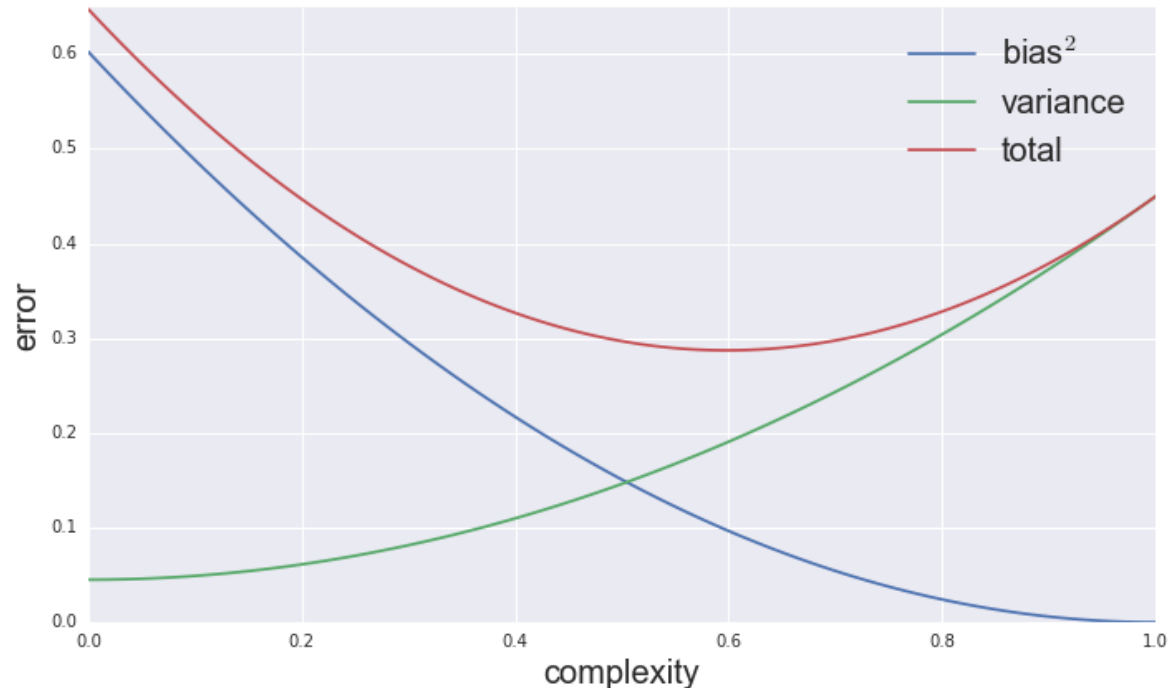
Bias-variance tradeoff

- Average conditional true error:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] &= \underbrace{\mathbb{E}_{Y|x}[(Y - \eta(x))^2]}_{\text{irreducible error}} \\ &+ \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} + \underbrace{\mathbb{E}_{\mathcal{D}}\left[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2\right]}_{\text{variance}} \end{aligned}$$

Bias squared:
measures how the predictor is mismatched with the best predictor in expectation

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

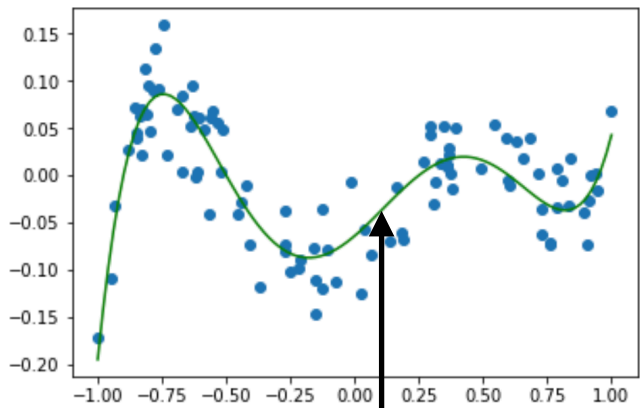


Questions?

Lecture 6: Bias-Variance Tradeoff (continued)

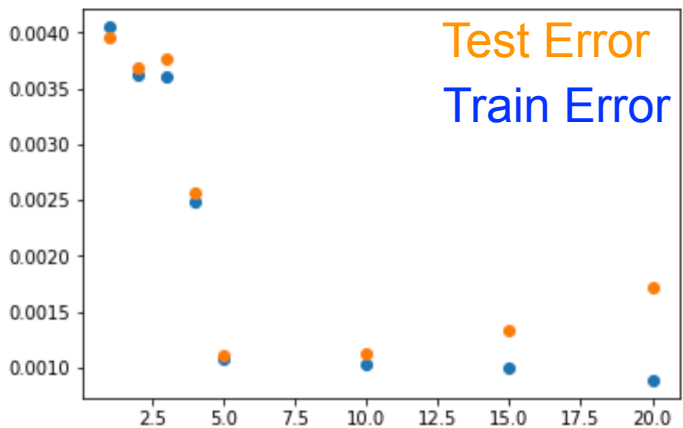


Test error vs. model complexity



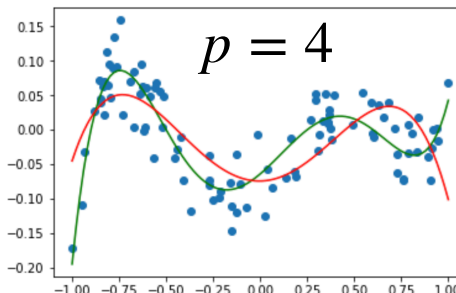
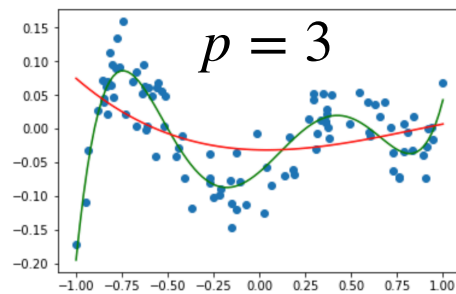
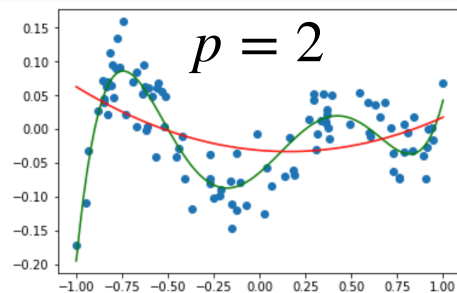
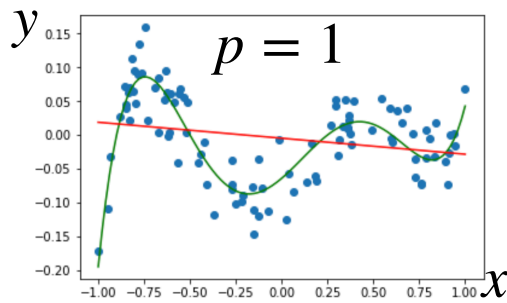
Optimal predictor $\eta(x)$ is degree-5 polynomial

Error

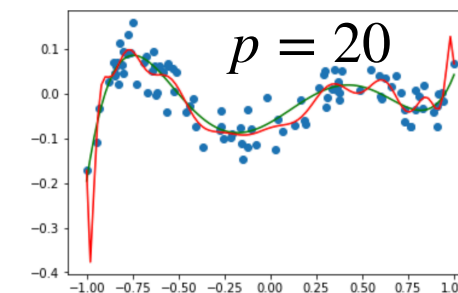
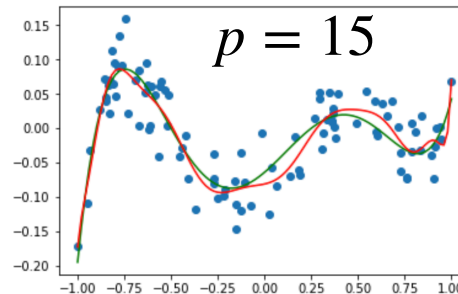
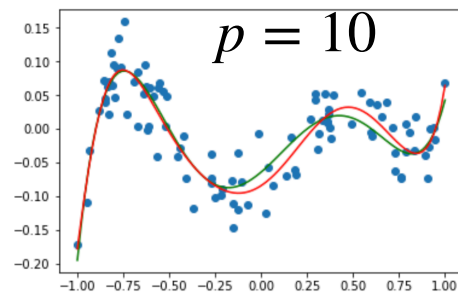
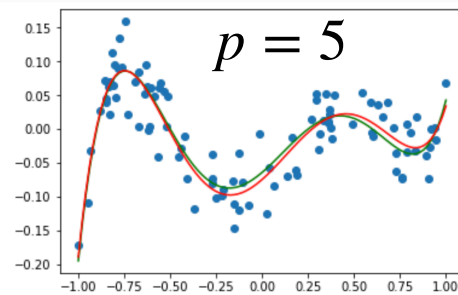


degree p of the polynomial regression

Simple model:
Model complexity is below
the complexity of $\eta(x)$

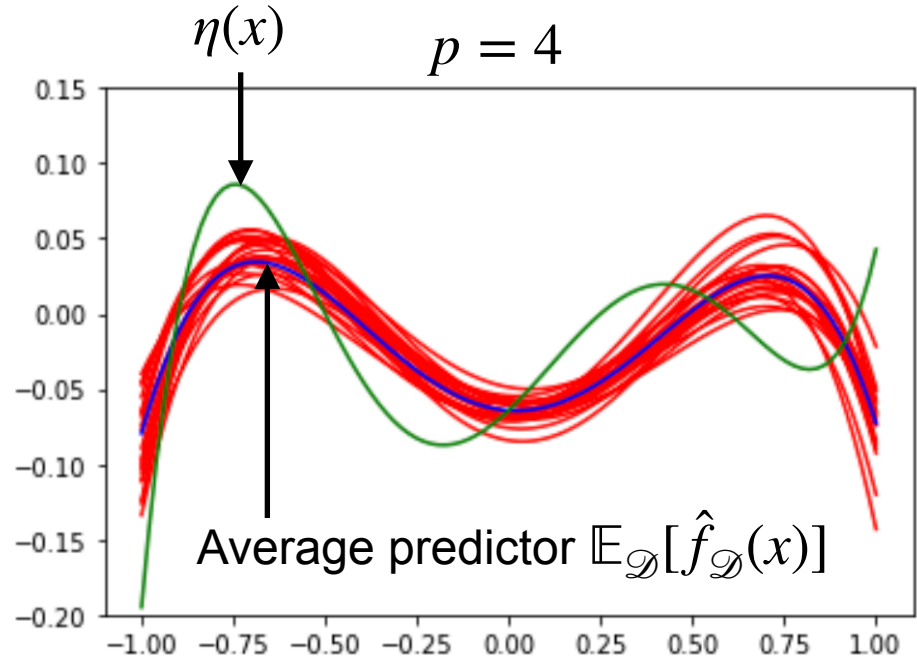
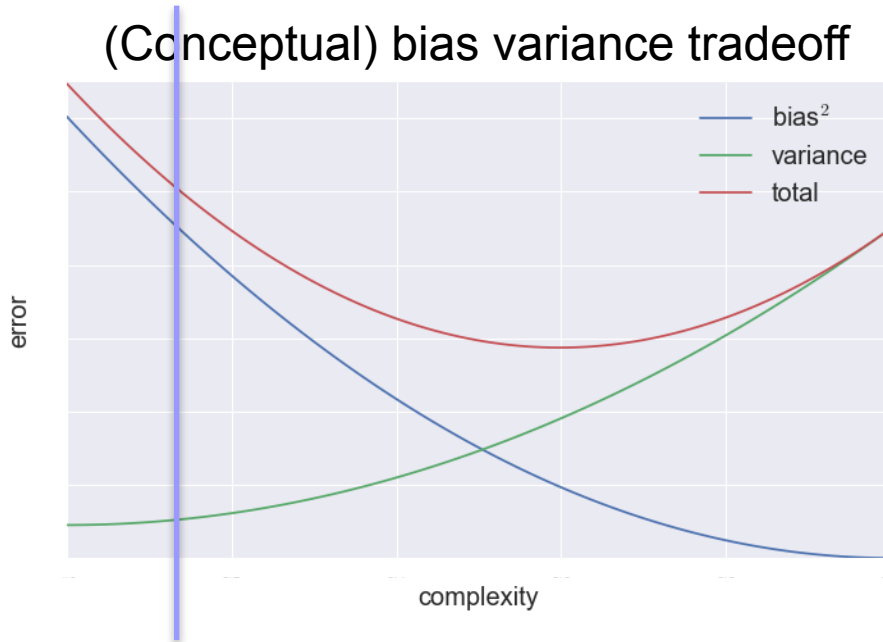


Complex model:



Recap: Bias-variance tradeoff with simple model

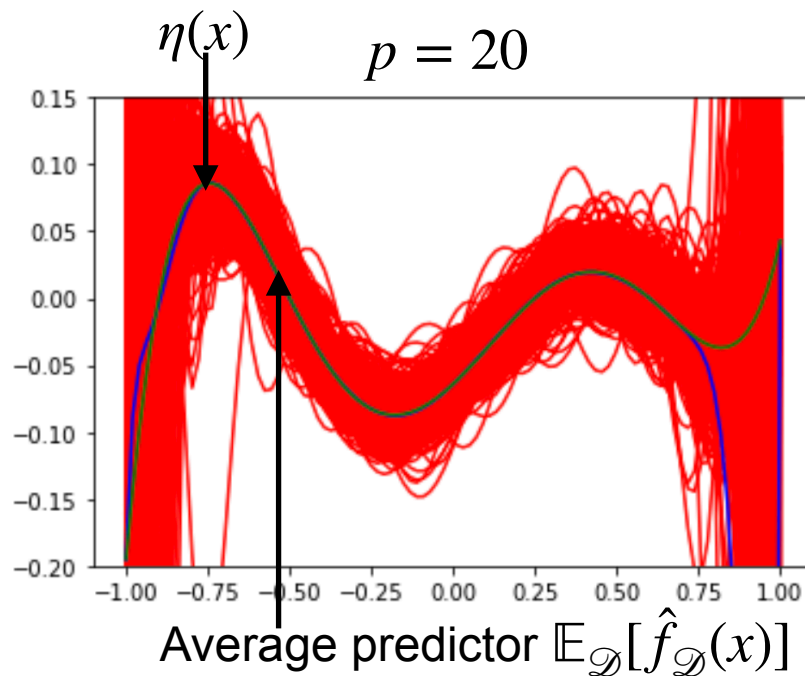
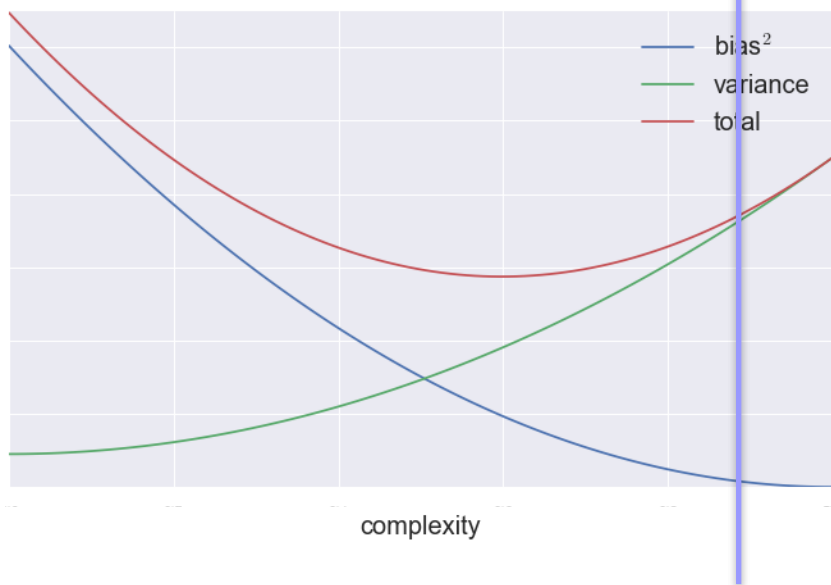
(Conceptual) bias variance tradeoff



- When model **complexity is low** (lower than the optimal predictor $\eta(x)$)
 - Bias² of our predictor, $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2$, is large
 - Variance of our predictor, $\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$, is small
 - If we have more samples, then
 - Bias
 - Variance
 - Because Variance is already small, overall test error

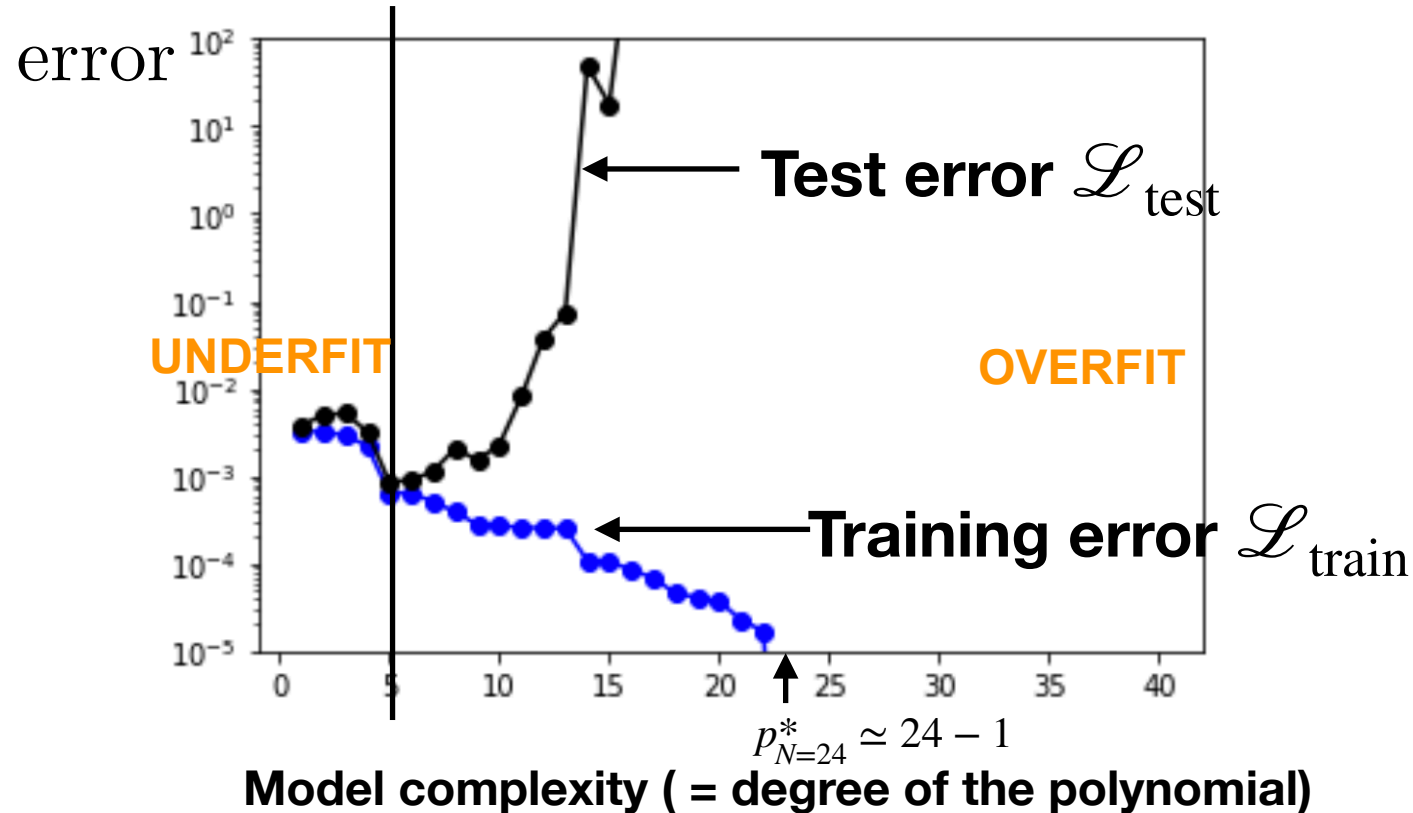
Recap: Bias-variance tradeoff with simple model

(Conceptual) bias variance tradeoff



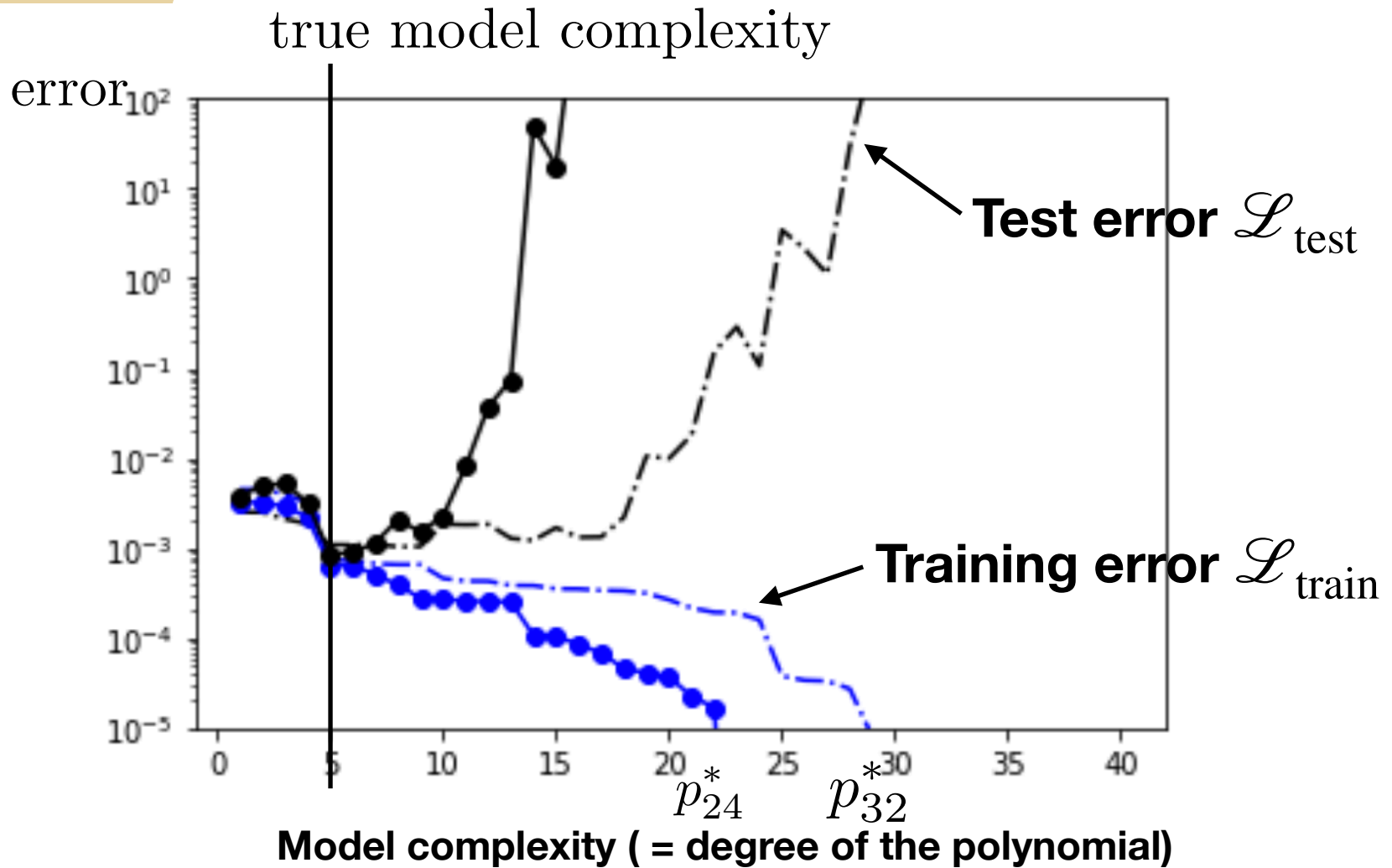
- When model complexity is high (higher than the optimal predictor $\eta(x)$)
 - Bias of our predictor, $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2$, is small
 - Variance of our predictor, $\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$, is large
 - If we have more samples, then
 - Bias
 - Variance
 - Because Variance is dominating, overall test error

- let us first fix sample size $N=30$, collect one dataset of size N i.i.d. from a distribution, 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



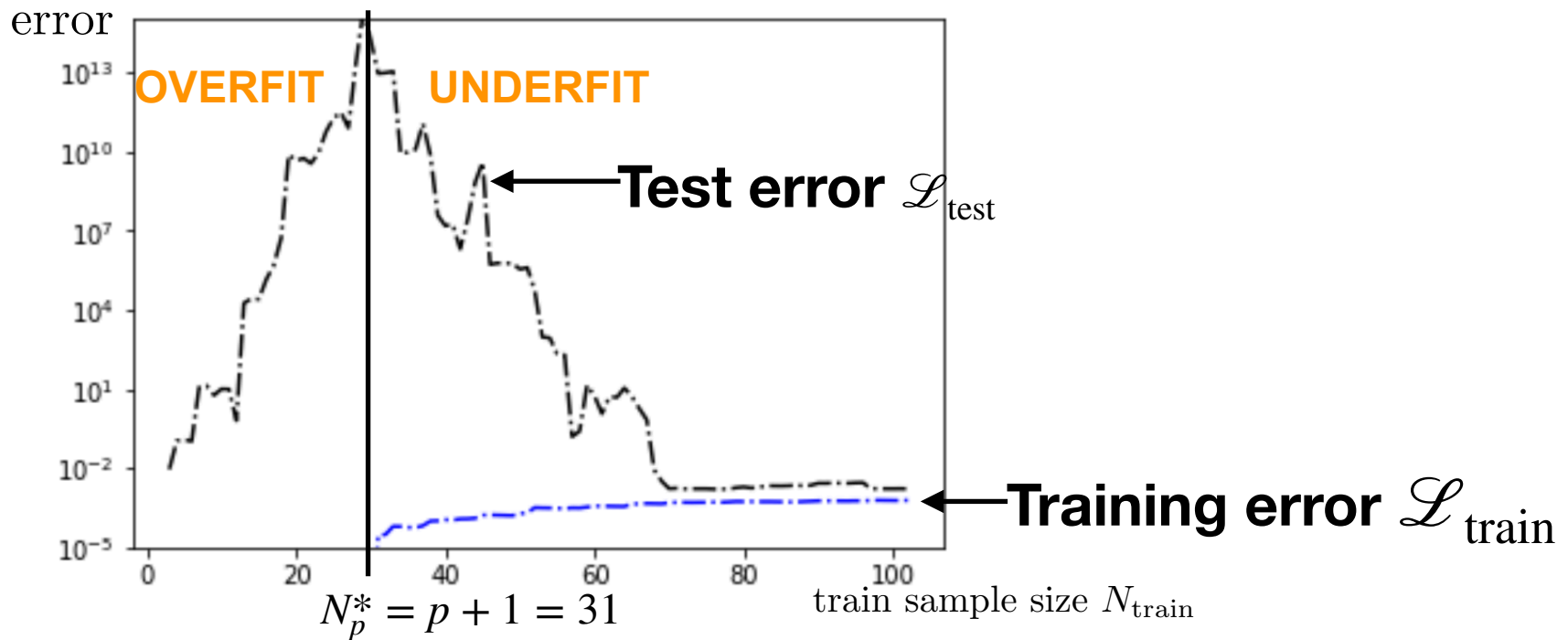
- Given sample size N there is a threshold, p_N^* , where training error is zero
- Training error is **always** monotonically non-increasing
- Test error has a trend of going down and then up, but fluctuates

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



- The threshold, p_N^* , moves right
- Training error tends to increase, because more points need to fit
- Test error tends to decrease, because Variance decreases

- 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 plot the computed MSEs for all values of train sample size N_{train} that we are interested in



- There is a threshold, N_p^* , below which training error is zero (extreme overfit)
- Below this threshold, test error is meaningless, as we are overfitting and there are multiple predictors with zero training error some of which have very large test error
- Test error tends to decrease
- Training error tends to increase

Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \epsilon$$

$$\widehat{w}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} =$$
$$=$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y | X = x] =$$

$$\hat{f}_{\mathcal{D}}(x) = x^T \widehat{w}_{\text{MLE}} =$$

Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \epsilon$$

$$\begin{aligned}\widehat{w}_{\text{MLE}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \epsilon) \\ &= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y | X = x] = x^T w^*$$

$$\widehat{f}_{\mathcal{D}}(x) = x^T \widehat{w}_{\text{MLE}} = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

- Irreducible error: $\mathbb{E}_{X,Y}[(Y - \eta(x))^2 | X = x] =$
- Bias squared: $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\widehat{f}_{\mathcal{D}}(x)])^2 =$
(is independent of the sample size!)

Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\widehat{w}_{\text{MLE}} = w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = x^T w^*$$

$$\hat{f}_{\mathcal{D}}(x) = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

- Variance: $\mathbb{E}_{\mathcal{D}} \left[\left(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2 \right] =$

Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\widehat{w}_{\text{MLE}} = w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = x^T w^*$$

$$\widehat{f}_{\mathcal{D}}(x) = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

- Variance: $\mathbb{E}_{\mathcal{D}} \left[\left(\widehat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\widehat{f}_{\mathcal{D}}(x)] \right)^2 \right] = \mathbb{E}_{\mathcal{D}} [x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x]$
 $= \sigma^2 \mathbb{E}_{\mathcal{D}} [x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x]$
 $= \sigma^2 x^T \mathbb{E}_{\mathcal{D}} [(\mathbf{X}^T \mathbf{X})^{-1}] x$
- To analyze this, let's assume that $X_i \sim \mathcal{N}(0, \mathbf{I})$ and number of samples, n , is large enough such that $\mathbf{X}^T \mathbf{X} = n\mathbf{I}$ with high probability and $\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1}] \simeq \frac{1}{n} \mathbf{I}$, then
 - Variance is $\frac{\sigma^2 x^T x}{n}$, and decreases with increasing sample size n

Questions?
