

Do not distribute course material

You may not and may not allow others to reproduce or distribute lecture notes and course materials publicly whether or not a fee is charged.



Topic 3 Model Selection

PROF. LINDA SELLIE

Thanks to:

- ❑ Some of the material is from Prof. Sundeep Rangan
 - This includes some slides and the motivating examples
- ❑ Some slides (the slides with the green background) are from Yaser Abu-Mostafa

Finding Parameters via Optimization

A general ML recipe

General ML problem

- ❑ Get data
- ❑ Pick a **model** with **parameters**
- ❑ Pick a **loss function**
 - Measures goodness of fit model to data
 - Function of the parameters
- ❑ Find parameters that **minimizes** loss

Multiple linear regression

- 1) Finding a way to have a more complex hypothesis class
- 2) If we have more than one hypothesis class to choose from - how do we select which one to use?

Loss function:
$$RSS(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

Select $\mathbf{w} = [w_0, w_1, w_2, \dots, w_d]^T$ to minimize $RSS(\mathbf{w})$

- In learning, our goal is to find a hypothesis that minimizes $E_{\text{out}}[g(\mathbf{x})]$ (not just $E_{\text{in}}[g(\mathbf{x})]$).
- In this lecture, we observe that choosing the model with the smallest training error doesn't work.
- Next we explore the different types of errors we make.
- We have to find a way to compare models.

Generalization

Training Error

$$E_{\text{in}}(w_0, w_1) = \underbrace{\frac{1}{N} \sum_{i=1}^N}_{\text{Average error on the } N \text{ training examples}} \underbrace{\text{error}(y^{(i)}, g(\mathbf{x}^{(i)}))}_{\substack{\text{Cost (loss) for} \\ \text{prediction not being} \\ \text{the same as true label}}} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \underbrace{(w_0 + w_1 \mathbf{x}^{(i)})}_{\text{Prediction on input } \mathbf{x}^{(i)}})^2$$

MSE over the training data is called the “in sample” error.

Generalization Error

$$E_{\text{out}}(w_0, w_1) = E_{\mathbf{x}, y} [\text{error}(y, g(\mathbf{x}))]$$

Assumption is training data is from the same distribution as the hypothesis will be used

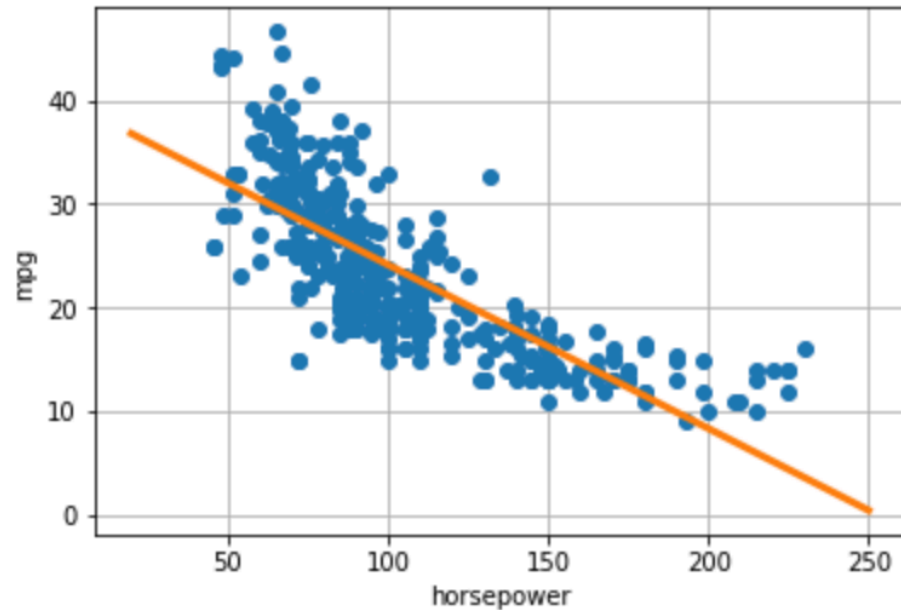
Expected error when the model is used on new examples. It is called the “out of sample error”. We cannot compute this value.

Expectation taken over all possible input/labels and the probability that input/label is seen

Outline

-
- ❑ Motivating example: What polynomial degree should a model have?
Yea!
Uh oh....
How to create a more complex hypothesis
 - ❑ Polynomial transformation
 - ❑ Underfitting and overfitting
 - ❑ Understanding error: Bias and variance
Understanding what went wrong
Understanding where the error comes from, and how to estimate $E_{\text{out}}[g(\mathbf{x})]$
 - ❑ Learning curves
 - ❑ validation and model selection
 - ❑ Model selection (with limited resources)
Our strategy
If we have many different hypothesis classes to choose from - how can we choose wisely?
And how can we estimate $E_{\text{out}}[g(\mathbf{x})]$?
 - ❑ K-fold cross validation
 - ❑ Regularization

Estimating Automobile MPG



- Found best line/hyperplane

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

- Shape appear to be nonlinear...

- To reduce E_{in} (RSS) we need something non-linear...

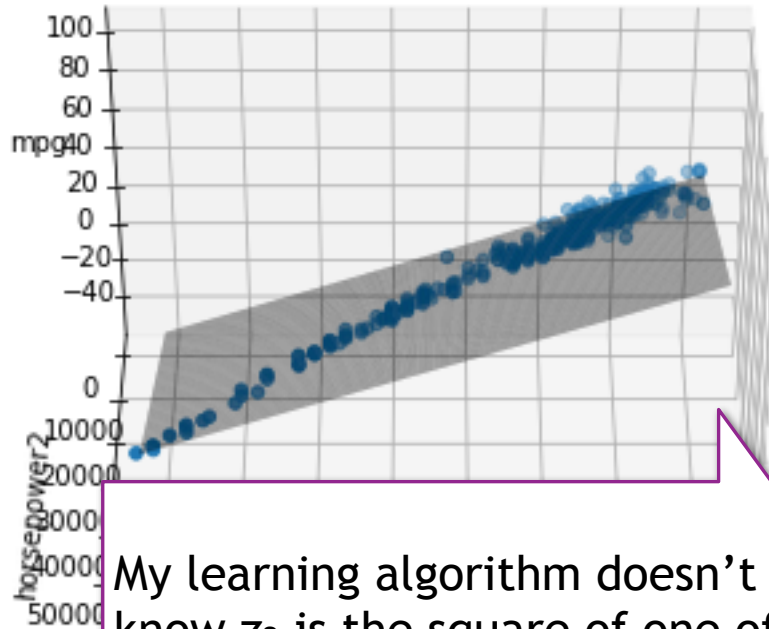
How can we get a non-linear hypothesis *easily*?

Outline

- ❑ Motivating example: What polynomial degree should a model use?
 - ➔ ❑ Polynomial transformation
 - ❑ Underfitting and overfitting
 - ❑ Understanding error: Bias and variance and noise
 - ❑ Learning curves
 - ❑ validation
 - ❑ validation and model selection
 - ❑ Model selection (with limited data)
 - ❑ Regularization
- How to create a more complex hypothesis
- Understanding where the error comes from and how to estimate $E_{\text{out}}[g(\mathbf{x})]$
- If we have many different hypothesis classes to choose from - how can we choose wisely? And what is the error of the hypothesis we chose?
- $g(\mathbf{x})$ is our hypothesis

A better hypothesis:

The R^2 value is 0.69 which is better than our previous R^2 value 0.53



My learning algorithm doesn't know z_2 is the square of one of my original features (horsepower^2). The learning algorithm only sees feature z_2

$$z_1 = \text{horsepower}$$
$$z_2 = \text{horsepower}^2$$

Trained my linear model on these features:
 z_1 and z_2
(aka horsepower and horsepower²)

$$\mathbf{x} \rightarrow \mathbf{z} = \Phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \phi_1(\mathbf{x}) \\ \phi_2(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ z_2 \end{bmatrix}$$

Learn in \mathbf{z} space with $\tilde{\mathbf{w}} = [\tilde{w}_0, \tilde{w}_1, \tilde{w}_2]$

Predict in \mathbf{z} space $\hat{y} = g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})) = \tilde{g}(\mathbf{z}) = \tilde{\mathbf{w}}^T \mathbf{z}$

$$\hat{y} = \underbrace{56.9}_{\tilde{w}_0} \cdot \underbrace{1}_{z_0} + \underbrace{(-0.466)}_{\tilde{w}_1} \cdot \underbrace{x}_{z_1} + \underbrace{0.00123}_{\tilde{w}_2} \cdot \underbrace{x^2}_{z_2}$$

$$= 56.9 \cdot 1 + (-0.466) \cdot \phi_1(\mathbf{x}) + 0.00123 \cdot \phi_2(\mathbf{x})$$


What is the feature vector in z -space of a car whose horsepower is 170 ?

$$[170]$$

$$[1, 170]^T$$

$$[1, 170, 170^2]^T$$

None of the above



$$X = \begin{bmatrix} 1 & 130.0 \\ 1 & 165.0 \\ 1 & 150.0 \\ 1 & 150.0 \\ 1 & 140.0 \\ 1 & 198.0 \\ 1 & 220.0 \\ 1 & 215.0 \\ \dots & \dots \end{bmatrix} \quad y = \begin{bmatrix} 18.0 \\ 15.0 \\ 18.0 \\ 16.0 \\ 17.0 \\ 15.0 \\ 14.0 \\ 14.0 \\ \dots \end{bmatrix}$$

Horsepower

$$X = \begin{bmatrix} 1 & 130.0 & 16900.0 \\ 1 & 165.0 & 27225.0 \\ 1 & 150.0 & 22500.0 \\ 1 & 150.0 & 22500.0 \\ 1 & 140.0 & 19600.0 \\ 1 & 198.0 & 39204.0 \\ 1 & 220.0 & 48400.0 \\ 1 & 215.0 & 46225.0 \\ \dots & \dots & \dots \end{bmatrix} \quad y = \begin{bmatrix} 18.0 \\ 15.0 \\ 18.0 \\ 16.0 \\ 17.0 \\ 15.0 \\ 14.0 \\ 14.0 \\ \dots \end{bmatrix}$$

Horsepower Horsepower²

Using our closed form solution we calculate $\tilde{\mathbf{w}} = \begin{bmatrix} 56.9 \\ -0.466 \\ 0.00123 \end{bmatrix}$

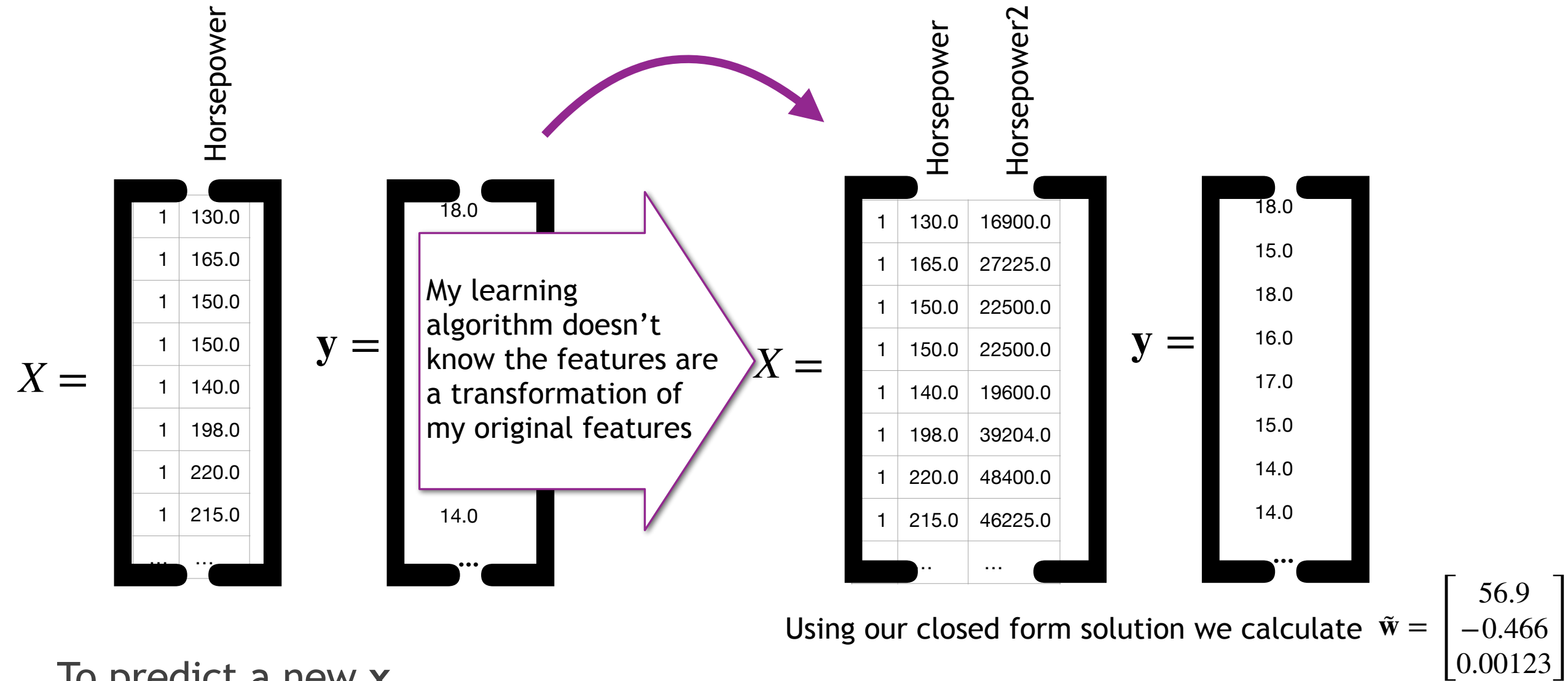
To predict a new \mathbf{x}

- transform \mathbf{x} to $\Phi(\mathbf{x})=\mathbf{z}$
- predict with $\tilde{\mathbf{w}}$ in \mathbf{z} -space

$$\hat{y} = \tilde{g}(\mathbf{z}) = \tilde{\mathbf{w}}^T \mathbf{z} = \tilde{\mathbf{w}}^T \Phi(\mathbf{x})$$

Estimated value of a car with horsepower = 170?

$$\tilde{\mathbf{w}}^T \Phi(\mathbf{x}) = \tilde{\mathbf{w}}^T \begin{bmatrix} 1 \\ 170 \\ 28900 \end{bmatrix} = [56.9 \quad -0.466 \quad 0.00123] \begin{bmatrix} 1 \\ 170 \\ 28900 \end{bmatrix}$$



To predict a new \mathbf{x}

- transform \mathbf{x} to $\Phi(\mathbf{x})=\mathbf{z}$
- predict with $\tilde{\mathbf{w}}$ in \mathbf{z} -space

$$\hat{y} = \tilde{g}(\mathbf{z}) = \tilde{\mathbf{w}}^T \mathbf{z} = \tilde{\mathbf{w}}^T \Phi(\mathbf{x})$$

Estimated value of a car with horsepower = 170?

$$\tilde{\mathbf{w}}^T \Phi(\mathbf{x}) = \tilde{\mathbf{w}}^T \begin{bmatrix} 1 \\ 170 \\ 28900 \end{bmatrix} = [56.9 \quad -0.466 \quad 0.00123] \begin{bmatrix} 1 \\ 170 \\ 28900 \end{bmatrix}$$

The General Polynomial Transform Φ_k

Polynomial basis function
polynomial features

Example: The degree-k polynomial transform over two features $\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$

k is a hyperparameter (i.e. not one of the decision variables being optimized when fitting the data)

$$\mathbf{z} = \Phi_1(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ z_2 \end{bmatrix} \quad \Phi_2(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1x_2 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \end{bmatrix} \quad \Phi_3(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1x_2 \\ x_2^2 \\ x_1^3 \\ x_1^2x_2 \\ x_1x_2^2 \\ x_2^3 \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \\ z_6 \\ z_7 \\ z_8 \\ z_9 \end{bmatrix} \quad \Phi_4(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1x_2 \\ x_2^2 \\ x_1^3 \\ x_1^2x_2 \\ x_1x_2^2 \\ x_2^3 \\ x_1^4 \\ x_1^3x_2 \\ x_1^2x_2^2 \\ x_1x_2^3 \\ x_2^4 \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \\ z_6 \\ z_7 \\ z_8 \\ z_9 \\ z_{10} \\ z_{11} \\ z_{12} \\ z_{13} \\ z_{14} \end{bmatrix} \quad \text{And so on :}$$

Square of x_2 →

Dimensionality of the features space increases rapidly

No weights in this space

General Feature Transform

$\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^{\tilde{d}}$ is also called a **feature map**

\mathcal{X} – space is \mathbb{R}^d

$$\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \\ \dots \\ x_d^{(i)} \end{bmatrix}$$

\mathcal{Z} – space is $\mathbb{R}^{\tilde{d}}$

$$\Phi(\mathbf{x}^{(i)}) = \mathbf{z}^{(i)} = \begin{bmatrix} 1 \\ \phi_1(\mathbf{x}^{(i)}) \\ \phi_2(\mathbf{x}^{(i)}) \\ \vdots \\ \phi_{\tilde{d}}(\mathbf{x}^{(i)}) \end{bmatrix} = \begin{bmatrix} 1 \\ z_1^{(i)} \\ z_2^{(i)} \\ \vdots \\ z_{\tilde{d}}^{(i)} \end{bmatrix}$$

Any function of the original features could be used

Training data : $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})$

$(\mathbf{z}^{(1)}, y^{(1)}), (\mathbf{z}^{(2)}, y^{(2)}), \dots, (\mathbf{z}^{(N)}, y^{(N)})$

No weights in original space

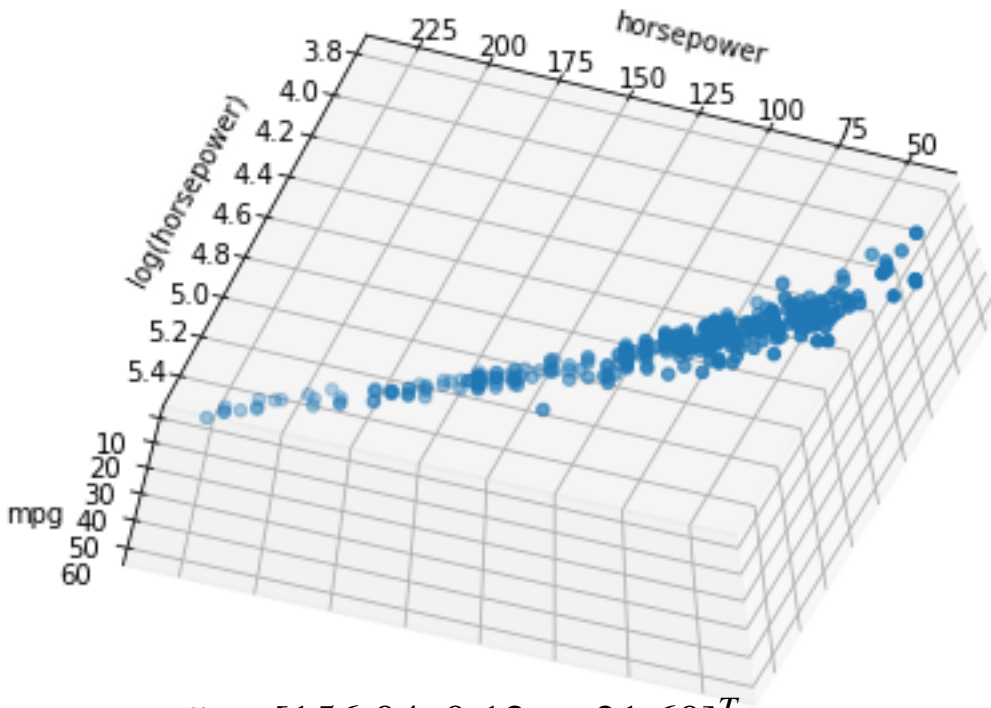
$$\hat{y} = \tilde{g}(\mathbf{z}) = \tilde{\mathbf{w}}^T \mathbf{z} = \tilde{\mathbf{w}}^T \Phi(\mathbf{x})$$

$$\tilde{\mathbf{w}} = \begin{bmatrix} \tilde{w}_0 \\ \tilde{w}_1 \\ \vdots \\ \tilde{w}_{\tilde{d}} \end{bmatrix}$$

We form a linear combination of the ϕ_j thus they are called **basis functions**

Many nonlinear features may work

$$\mathbf{x}^{(i)} \rightarrow \mathbf{z}^{(i)} = \Phi(\mathbf{x}^{(i)}) = \begin{bmatrix} 1 \\ \phi_1(\mathbf{x}^{(i)}) \\ \phi_2(\mathbf{x}^{(i)}) \end{bmatrix} = \begin{bmatrix} 1 \\ z_1^{(i)} = x_1^{(i)} \\ z_2^{(i)} = \log(x_1^{(i)}) \end{bmatrix}$$



$$\tilde{\mathbf{w}} = [156.04, 0.12, -31.60]^T$$

The R^2 value is 0.68

Polynomial Regression

- Models the relationship between the response and features as an d^{th} order polynomial

$$y = w_0 + w_1x + w_2x^2 + w_3x^3 + \dots + w_dx^d + \epsilon$$

Example is using only one feature (monomial)

- Observation: the higher the order of the polynomial, the more shapes you can fit!

- costs:

- computational complexity grows as the number of parameters grows
- chance that you *model the noise* and not the underlying parameters increases - overfitting - lose generalization

- Warning! It is always possible to perfectly fit N points with a model of order $(N-1)$. It is unlikely that such a model will provide knowledge of the unknown function or be able to predict as well on unseen data as a lower order polynomial

How can we choose which (if any) transformation to use?

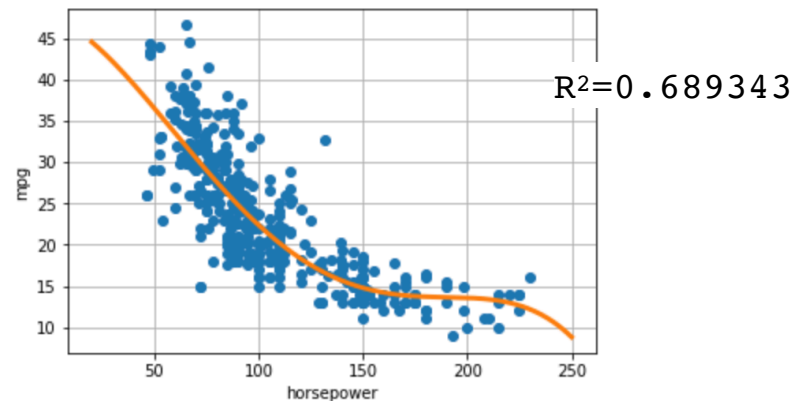
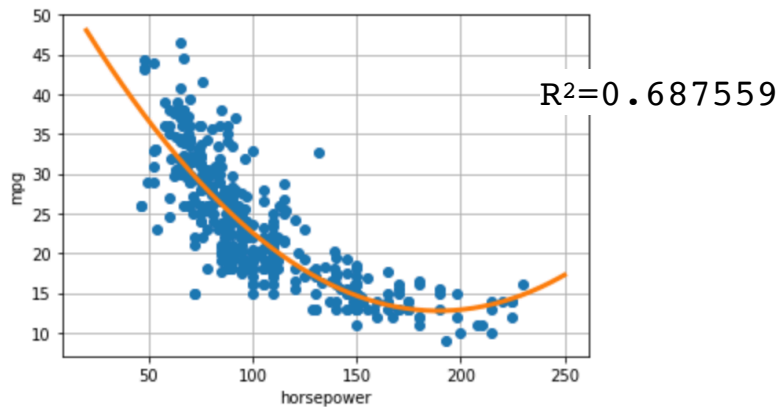
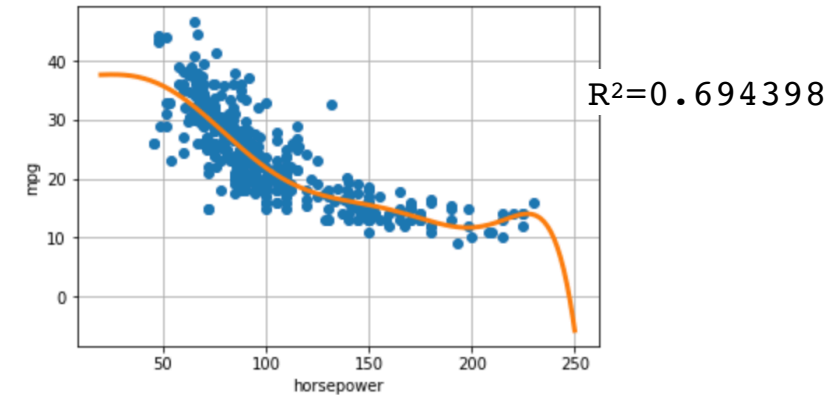
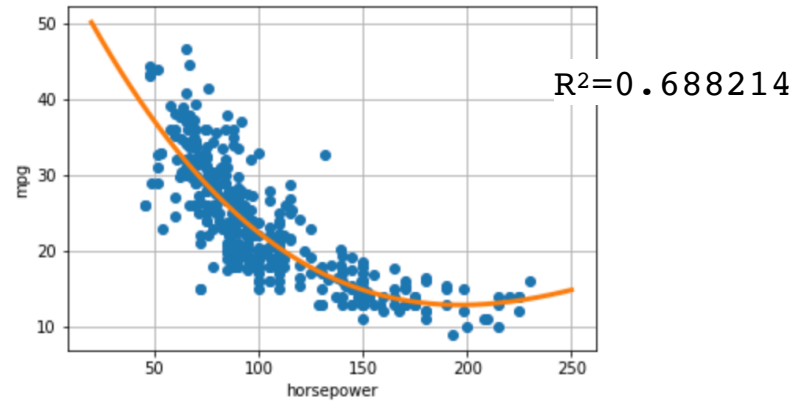
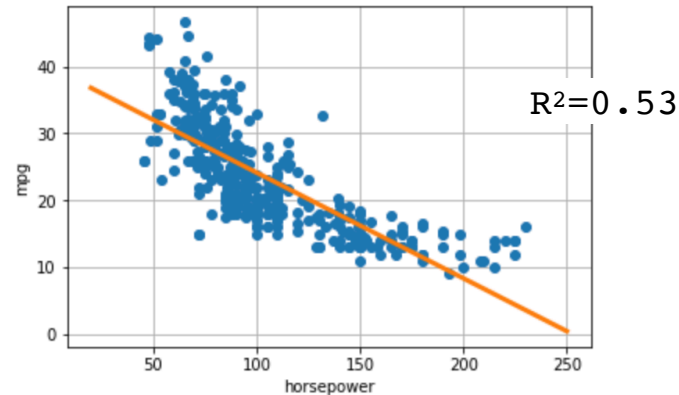
Outline

- ❑ Motivating example: What polynomial degree should a model have?
 - ❑ Polynomial transformation
 - ❑ Underfitting and overfitting
 - ❑ Understanding error: Bias and variance
 - ❑ Learning curves
 - ❑ validation and model selection
 - ❑ Model selection (with limited data)
 - ❑ K-fold cross validation
 - ❑ Regularization
-
- Yea!
- Uh oh....
- How to create a more complex hypothesis
- Understanding what went wrong
- Understanding where the error comes from, and how to estimate $E_{\text{out}}[g(\mathbf{x})]$
- Our strategy
- If we have many different hypothesis classes to choose from - how can we choose wisely? And how can we estimate $E_{\text{out}}[g(\mathbf{x})]$?

Automobile MPG

□ As we increase the degree of the polynomial, do we improve the fit of the model to the data?

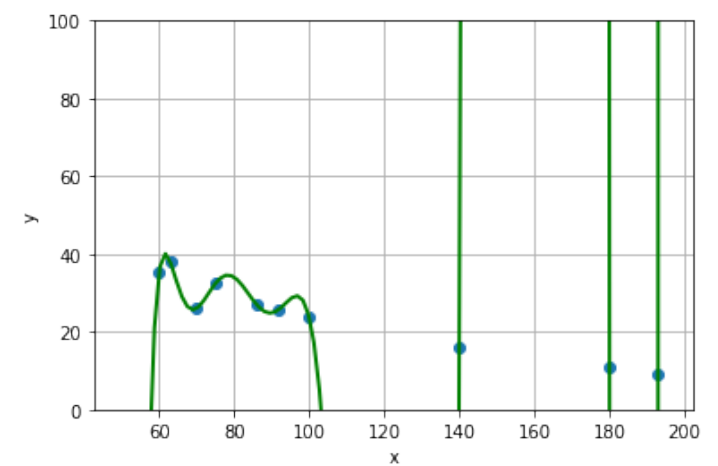
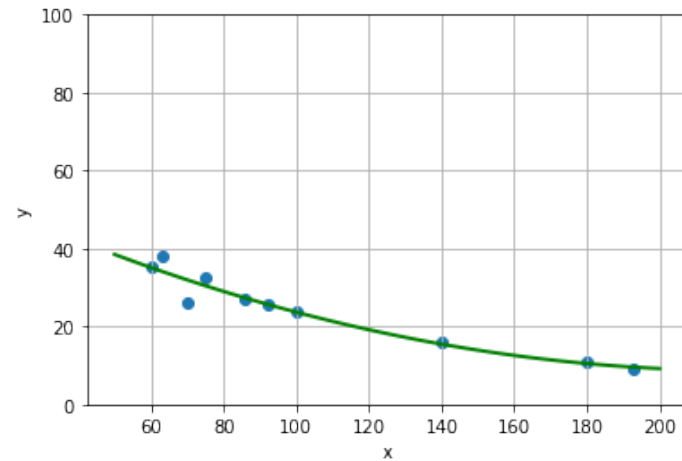
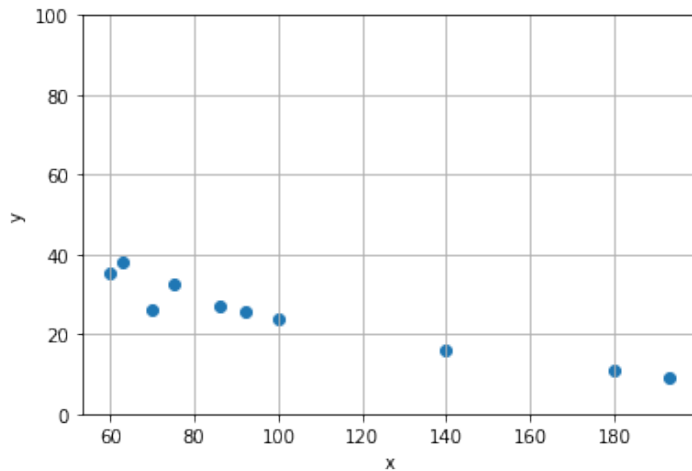
$$y = w_0 + w_1x + w_2x^2 + w_3x^3 + \dots + w_dx^d + \epsilon$$



We can keep improving our wrt our training data - but does that mean we would do better on examples not in the training data?

Example

mean = 23.4710191083



$$E_{\text{in}} \approx 0; \quad E_{\text{out}} \gg 0$$

Overfitting: Complex hypothesis that fits the training data too well.
It predicts well on patterns found in training data that won't be found in the the future data

What is the goal of machine learning?

Find a hypothesis, $g(\mathbf{x})$ that predicts the correct value of y for new values of \mathbf{x} (i.e. not in our training set).

Let $E_{out}(g)$ be the expected error our hypothesis will incur in the future. For linear regression we incurred a loss of $(g(\mathbf{x}) - y)^2$. Thus for linear regression the expected error over the input space is $E_{out}(g) = E \left[(g(\mathbf{x}) - y)^2 \right]$

We cannot compute $E_{out}(g)$.

We have been focused on minimizing $E_{in}(g)$, the cost of our training set.

For linear regression
$$E_{in}(g) = \frac{1}{N} \sum_{i=1}^N (g(\mathbf{x}^{(i)}) - y^{(i)})^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2.$$

What is the goal of machine learning?

Find a hypothesis, $g(\mathbf{x})$ that predicts the correct value of y for new values of \mathbf{x} (i.e. not in our training set).

Let $E_{out}(g)$ be the expected error over the input space is $E_{out}(g) = \mathbb{E}[(g(\mathbf{x}) - y)^2]$. The expectation taken over all possible input/labels and the probability that input/label is seen

Also called **risk** of a hypothesis (model)

We cannot compute $E_{out}(g)$.

We have been focused on minimizing $E_{in}(g)$, the cost of our training set.

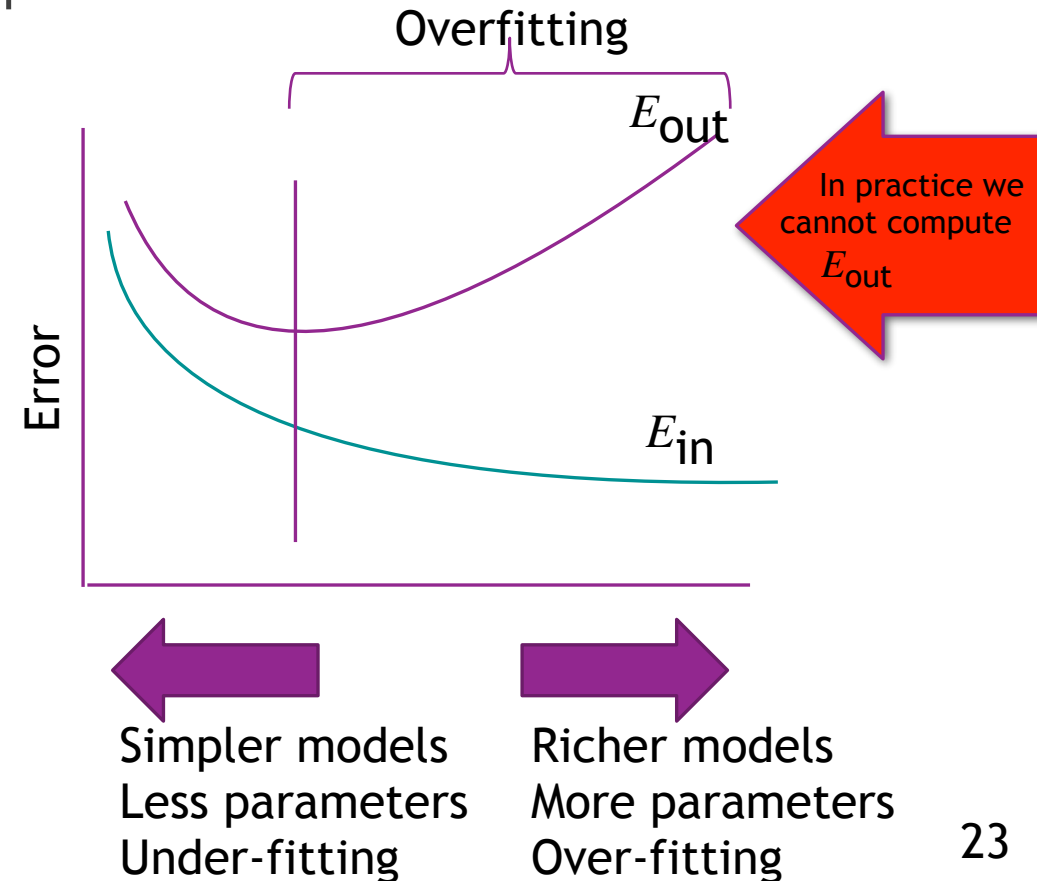
For linear regression
$$E_{in}(g) = \frac{1}{N} \sum_{i=1}^N (g(\mathbf{x}^{(i)}) - y^{(i)})^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2.$$

What can go wrong with choosing the hypothesis which has the smallest lost/cost

1. Limited Hypothesis class (model class). No function in our hypothesis class can model the data well - **biased solution**
2. Limited Data. We might model the noise and not the true pattern. Small changes to the data causes the hypothesis (model) to change - **high variance solution**

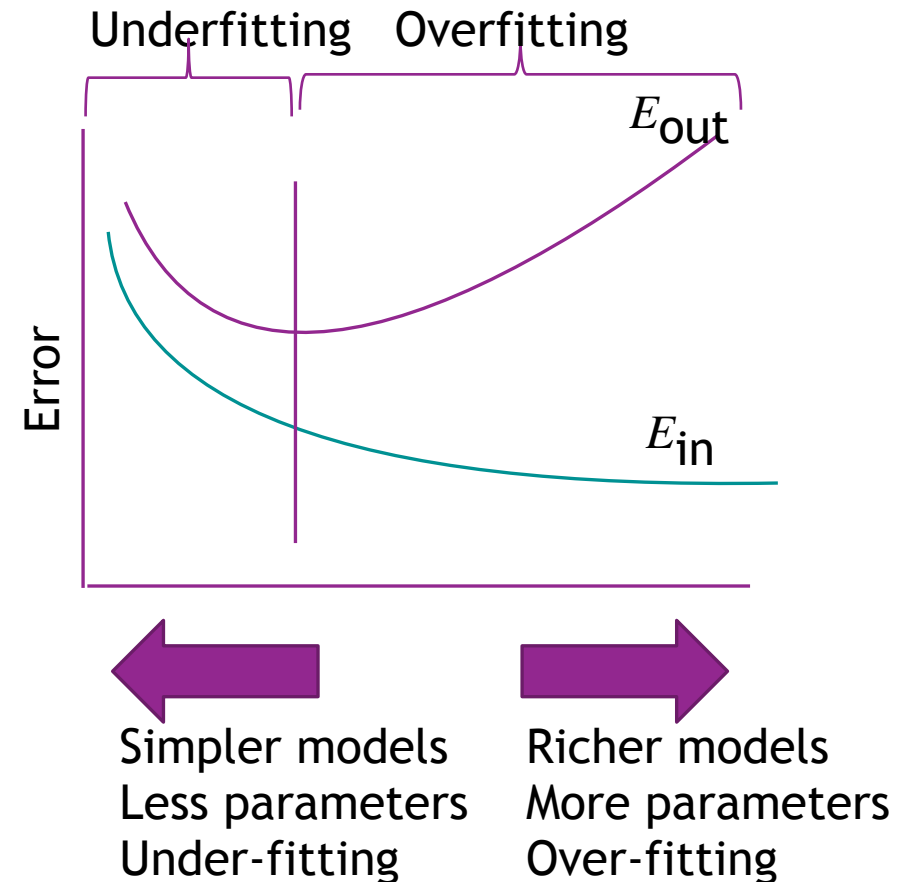
Overfitting

- If we allow a very flexible model by using complicated features or model, our model can be *too* complicated. The regression model becomes tailored to fit the noise of the training data and does not generalize well (i.e. predict well on data not in the training set, and does not accurately describe the relationship between the parameters and the outcome)
- A too complicated model will not generalize well.
- **overfitting**: The model performs worse on unseen data than a different model from the same class despite performing better on the training data
- Example: Using a degree $d=N-1$ polynomial transformation
- Training RSS (or MSE), R^2 is not a good indicator of test RSS (or MSE), R^2

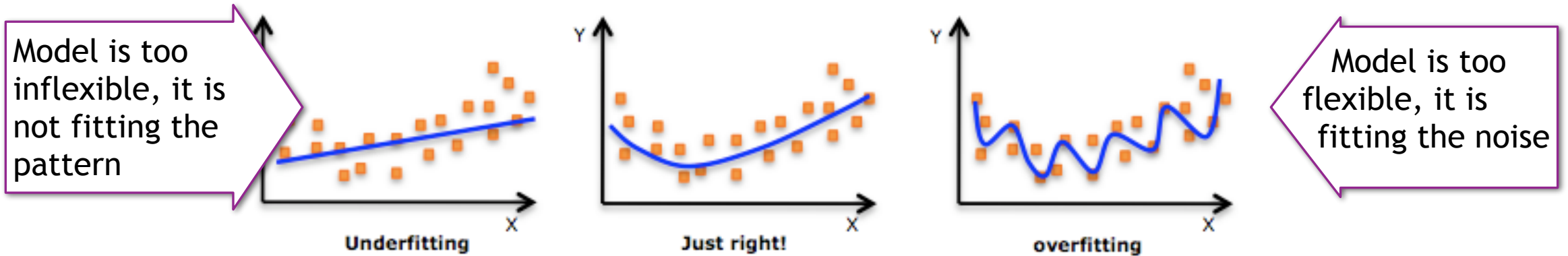


Underfitting

- ❑ The model learned does not do well on the training data and does not do well on unseen examples
- ❑ A too simple model is called *underfitting*
- ❑ Example: predicting mean of target



How Can You Tell from Data?



- ❑ Is there a way to tell what is the correct model order to use?
- ❑ Must use the data. Do not have access to the true d ?
- ❑ What happens if we guess:
 - d too big?
 - d too small?

Question

For the examples below, might we encounter a problem with the model we chose?

□ Examples:

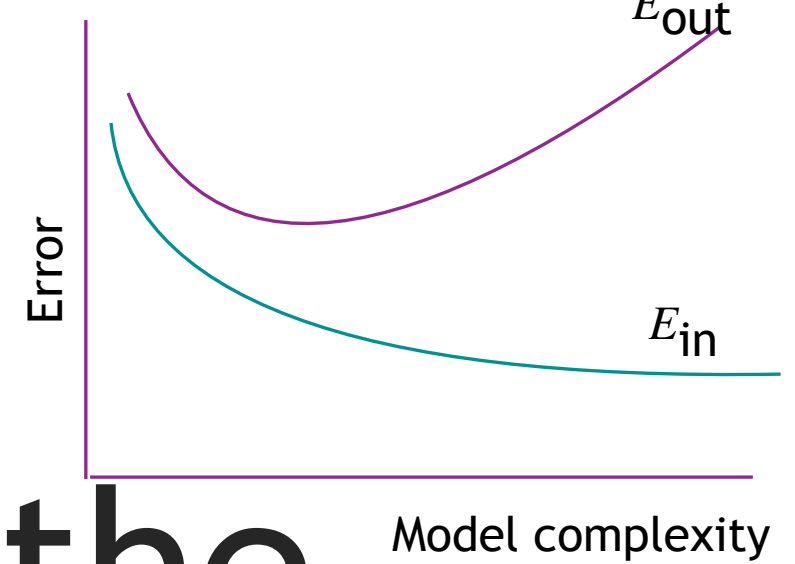
- True function $f(x) = 2 + 3x$ model class $w_0 + w_1x + w_2x^2$
- True function $f(x) = 2 + 3x + 4x^2$ model class $w_0 + w_1x$

Outline

- ❑ Motivating example: What polynomial degree should a model have?
 - ❑ Polynomial transformation
 - ❑ Underfitting and overfitting
 - ❑ Understanding error: Bias and variance
 - ❑ Learning curves
 - ❑ validation and model selection
 - ❑ Model selection (with limited data)
 - ❑ K-fold cross validation
 - ❑ Regularization
-
- The diagram illustrates the flow of the course outline. A purple arrow points to the 'Understanding error: Bias and variance' item. A red arrow labeled 'Understanding what went wrong' points from this item to the 'Learning curves' item. A red arrow labeled 'Our strategy' points from the 'Model selection (with limited data)' item to the 'K-fold cross validation' item. A red arrow labeled 'Yea! Uh oh....' points from the 'Learning curves' item to the 'validation and model selection' item. A purple bracket groups the 'validation and model selection' and 'Model selection (with limited data)' items, with a text box asking 'How to create a more complex hypothesis'. Another purple bracket groups the 'K-fold cross validation' and 'Regularization' items, with a text box asking 'Understanding where the error comes from, and how to estimate $E_{out}[g(\mathbf{x})]$ '. A third purple bracket groups the 'Model selection (with limited data)' and 'K-fold cross validation' items, with a text box asking 'If we have many different hypothesis classes to choose from - how can we choose wisely? And how can we estimate $E_{out}[g(\mathbf{x})]$ '.
- Yea! Uh oh....
- How to create a more complex hypothesis
- Understanding what went wrong
- Understanding where the error comes from, and how to estimate $E_{out}[g(\mathbf{x})]$
- Our strategy
- If we have many different hypothesis classes to choose from - how can we choose wisely? And how can we estimate $E_{out}[g(\mathbf{x})]$

Understanding the errors

THEORETICAL METRIC



How do we evaluate our model? Or choose among models (e.g. the which polynomial transformation should we choose?)

Open discussion

- We can evaluate how well it works by looking at its errors
- We would like the error to be zero on all future data. However:
 - The unseen variables means the true model has non-zero error (i.e. the world is a messy place)
 - Our hypothesis probably doesn't contain the underlying true model
 - We don't get enough data to perfectly estimate our model. We only get a finite sample of the data. The more data we receive, the more our sample is representative of underlying data and our estimates should converge

Noise

Bias

Variance



Understanding Error

$$E_{\text{out}}(g^{(D)}(\mathbf{x})) = E_{\mathbf{x},y}[(y - g^{(D)}(\mathbf{x}))^2]$$

Bias-Variance-Noise Decomposition

In predictions there are three sources of error

1. noise - irreducible error
2. bias - error of average hypothesis (estimated from N examples) from the true function $f(\mathbf{x}) + \epsilon$
3. variance - how much would the prediction for an example change if the hypothesis was fit on a different set of N points

Understanding Error

$$E_{\text{out}}(g^{(D)}(\mathbf{x})) = E_{\mathbf{x},y}[(y - g^{(D)}(\mathbf{x}))^2]$$

Bias-Variance-Noise Decomposition

Our definitions will be for the squared loss function
You can think of how to substitute other loss functions

This cannot be computed in practice
because we do not have access to the target function or the probability distribution

In predictions there are three sources of error

1. noise - irreducible error
2. bias - error of average hypothesis (estimated from N examples) from the true function $f(\mathbf{x}) + \epsilon$
3. variance - how much would the prediction for an example change if the hypothesis was fit on a different set of N points

Outline

❑ Motivating example: What polynomial degree should a

❑ Polynomial transformation

❑ Underfitting and overfitting

❑ Understanding error: Bias and variance and noise

• Bias

• Variance

• Bias and variance and noise

❑ Learning curves

❑ validation

❑ Model selection

❑ Cross validation

❑ Regularization

Yea!

Uh oh....

How to create a more complex hypothesis

Understanding where the error comes from and how to estimate $E_{\text{out}}[g(\mathbf{x})]$

Understanding what went wrong

Our strategy

If we have many different hypothesis classes to choose from - how can we choose wisely? And how to estimate $E_{\text{out}}[g(\mathbf{x})]$?

Average Hypothesis

- **Given:** N training examples $D = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
- **Learn:** If I had a different set of N training examples, I would get a different hypothesis (models) $g^{(D)}(\mathbf{x})$
- **Expected prediction (averaged over hypothesis):** $\bar{g}(\mathbf{x}) = E_D[g^{(D)}(\mathbf{x})]$

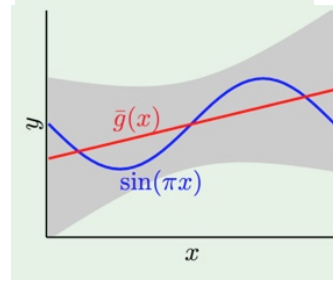
Mean prediction of the algorithm for \mathbf{x}

Intuitive approximation:

$$\bar{g}(\mathbf{x}) \approx \frac{1}{k} \sum_{i=1}^k g_i^{(D_i)}(\mathbf{x}) \quad D_1, D_2, \dots, D_k$$

Bias

Bias of the hypothesis class (not an individual hypothesis from the class)



- $\text{bias}(\mathbf{x}) = (f(\mathbf{x}) - \bar{g}(\mathbf{x}))^2$

Conceptually: squared difference from “average prediction” for \mathbf{x} , and expected label $f(\mathbf{x})$

- $\text{bias} = E_{\mathbf{x}} [(\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2]$

Bias of the hypothesis class

less flexible model then more bias

Occasionally this is called bias²

$$\approx \frac{1}{N} \sum_{i=1}^N (\bar{g}(\mathbf{x}^{(i)}) - f(\mathbf{x}^{(i)}))^2$$

When using this model class, measures how well you expect the “average prediction” to represent the true solution
We expect the bias to decrease with a more complex model

Outline

❑ Motivating example: What polynomial degree should a

❑ Polynomial transformation

❑ Underfitting and overfitting

❑ Understanding error: Bias and variance and noise

- Bias

- Variance

- Bias and variance and noise

❑ Learning curves

❑ validation

❑ Model selection

❑ Cross validation

❑ Regularization

Yea!

Uh oh....

How to create a more complex hypothesis

Understanding where the error comes from and how to estimate $E_{\text{out}}[g(\mathbf{x})]$

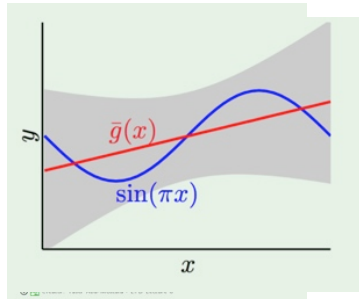
Understanding what went wrong

Our strategy

If we have many different hypothesis classes to choose from - how can we choose wisely? And how to estimate $E_{\text{out}}[g(\mathbf{x})]$?

Variance

Variance of a hypothesis class (model class)



- Variance: difference between the expected prediction and the prediction from a particular dataset

$$\bullet \text{ var}(\mathbf{x}) = E_D \{ (g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \} \approx \frac{1}{L} \sum_{\ell=1}^L (\bar{g}(\mathbf{x}) - g_{\ell}^{(D_{\ell})}(\mathbf{x}))^2$$

Conceptually: variance of a prediction for \mathbf{x} from the mean prediction

$$\text{var} = E_{\mathbf{x}} \left[E_D \left[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right] \right] \approx \frac{1}{N} \sum_{i=1}^N \frac{1}{L} \sum_{\ell=1}^L (\bar{g}(\mathbf{x}^{(i)}) - g_{\ell}^{(D_{\ell})}(\mathbf{x}^{(i)}))^2$$

less flexible model then less variance

Measures how sensitive a hypothesis class (model class) is to a specific dataset
Variance typically decreases with simpler models

Outline

❑ Motivating example: What polynomial degree should a

❑ Polynomial transformation

❑ Underfitting and overfitting

❑ Understanding error: Bias and variance and noise

- Bias

- Variance

- Bias and variance and noise

❑ Learning curves

❑ validation

❑ Model selection

❑ Cross validation

❑ Regularization

Yea!

Uh oh....

How to create a more complex hypothesis

Understanding where the error comes from and how to estimate $E_{\text{out}}[g(\mathbf{x})]$

Understanding what went wrong

Our strategy

If we have many different hypothesis classes to choose from - how can we choose wisely? And how to estimate $E_{\text{out}}[g(\mathbf{x})]$?

Where did the prediction error in our hypothesis come from?

□ Regression example: $y = f(\mathbf{x}) + \epsilon$

Deterministic

Noise $\sim N(0, \sigma)$
We are assuming the noise
has mean 0 and variance σ^2

This means $E_{\mathbf{x},y}[f(\mathbf{x}) - y] = 0$ and
 $E_{\mathbf{x},y}[(f(\mathbf{x}) - y)^2] = E_{\mathbf{x}}(\epsilon^2) = \sigma^2$

□ Goal is to understand why our *expected* hypothesis (model) does not have zero error

$$E_D[E_{\text{out}}(g^{(D)})] = E_D[E_{\mathbf{x},y}[(g^{(D)}(\mathbf{x}) - y)^2]] \neq 0$$

$E_{\text{out}}(g^{(D)})$

$E_{\mathbf{x},y}[(g^{(D)}(\mathbf{x}) - y)^2]$
expected error for they
hypothesis $g^{(D)}(\mathbf{x})$

The expected error of the
hypothesis fit using the data
set **D** on any future example

Understanding Error

Bias-Variance

Decomposition (noise free)

$$\text{bias}(\mathbf{x}) = (f(\mathbf{x}) - \bar{g}(\mathbf{x}))^2$$

$$\text{var}(\mathbf{x}) = E_D[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]$$

$$\bar{g}(\mathbf{x}) = E_D[g^{(D)}(\mathbf{x})]$$

For any constant c ,
 $E[ac] = cE[a]$
 $E[a+c] = E[a] + c$

The linearity of expectation:
 $E[a + b] = E[a] + E[b]$

$$E_D[E_{\mathbf{x}}[(g^{(D)}(\mathbf{x}) - f(\mathbf{x}))^2]]$$

$$E_{\mathbf{x}}[E_D[(g^{(D)}(\mathbf{x}) - f(\mathbf{x}))^2]] = E_{\mathbf{x}}[E_D[(\underbrace{g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x})}_A + \underbrace{\bar{g}(\mathbf{x}) - f(\mathbf{x})}_B)^2]]$$

$$= E_{\mathbf{x}}[E_D[(\underbrace{g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x})}_A)^2 + 2(\underbrace{g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x})}_A)(\underbrace{\bar{g}(\mathbf{x}) - f(\mathbf{x})}_B) + (\underbrace{\bar{g}(\mathbf{x}) - f(\mathbf{x})}_B)^2]]$$

$$= E_{\mathbf{x}}[E_D[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2] + 2E_D[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))(\bar{g}(\mathbf{x}) - f(\mathbf{x}))] + E_D[(\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2]]$$

$$= E_{\mathbf{x}}[\underbrace{E_D[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]}_{\text{variance}(\mathbf{x})} + \underbrace{2E_D[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))]}_0(\bar{g}(\mathbf{x}) - f(\mathbf{x})) + \underbrace{E_D[(\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2]}_{\text{bias}(\mathbf{x})}]$$

$$= E_{\mathbf{x}}[\text{bias}] + E_{\mathbf{x}}[\text{variance}]$$

$$= \text{bias} + \text{variance}$$

Notice that

$$E_D[(g^{(D)}(\mathbf{x}) - \bar{g}(\mathbf{x}))]$$

$$= E_D[g^{(D)}(\mathbf{x})] - \bar{g}(\mathbf{x})$$

Understanding Error

Bias-Variance-Noise Decomposition

The expected error of the hypothesis fit on a **randomly** chosen set of N training examples

$$E_{\text{out}}(g) = \text{bias} + \text{variance} + \text{noise}$$

Noise in the training set contributes to variance

Noise in the test set contributes to irreducible error

Based on averages over what is expected for a training set D

- can we lower variance without increasing too much the bias?
- can we lower bias without increasing too much the variance?

