



# Announcements

Homework 1 due tomorrow night!

Remember to do all 3 parts

- [A1: Concept]
- [A1: Programming]
- [A1: Upload]

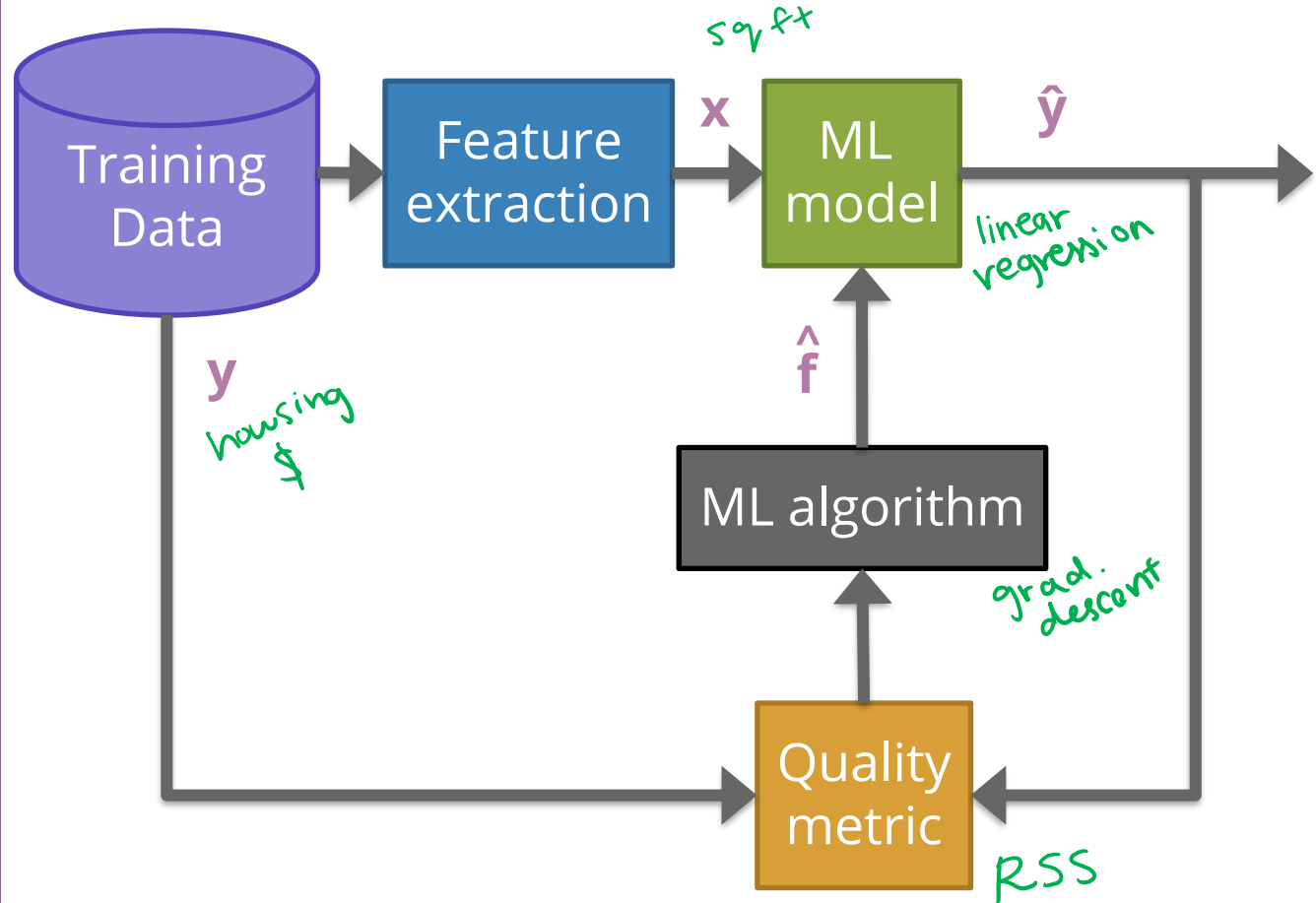
No office hours on Friday, but there is section on Thursday.

Happy 4<sup>th</sup> of July!

HW2 goes out on Wednesday and is due next Tuesday (like regular)

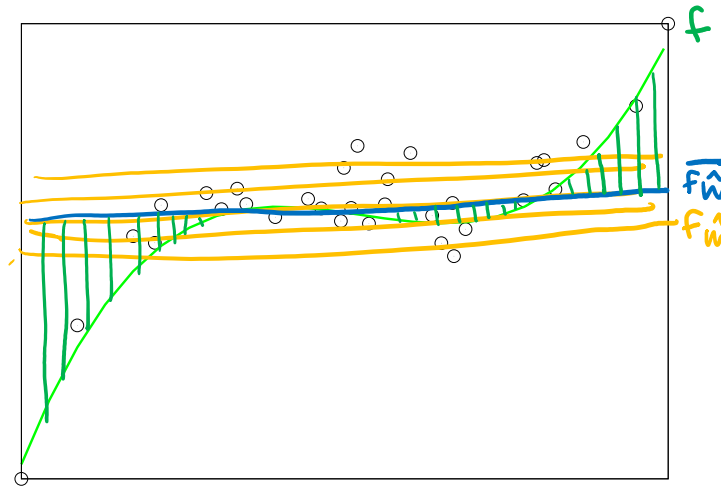
HW assignments are weighted equally.

# ML Pipeline



# Bias

A model that is too simple fails to fit the signal. In some sense, this signifies a fundamental limitation of the model we are using to fail to fit the signal. We call this type of error **bias**.

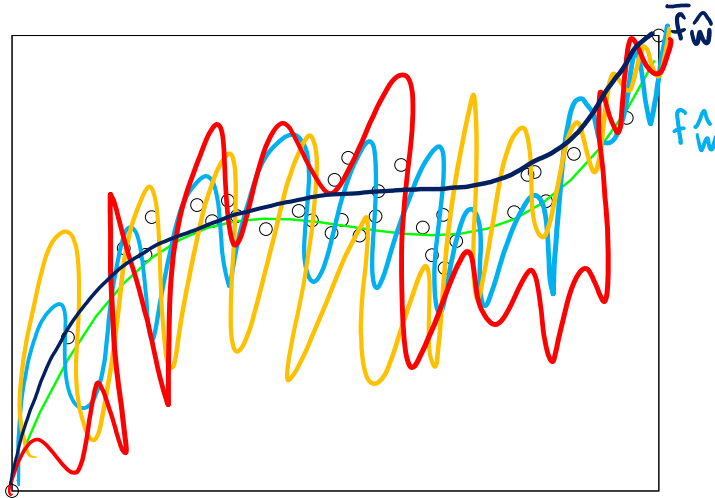


$$\text{Bias} = E \left[ \underbrace{f(x)}_{\substack{\text{true} \\ \text{function}}} - \underbrace{\overline{f_{\hat{w}}(x)}}_{\substack{\text{average} \\ \text{predictor}}} \right]$$

Low complexity (simple) models tend to have high bias.\*

# Variance

A model that is too complicated for the task overly fits to the noise. The flexibility of the complicated model makes it capable of memorizing answers rather than learning general patterns. This contributes to the error as **variance**.



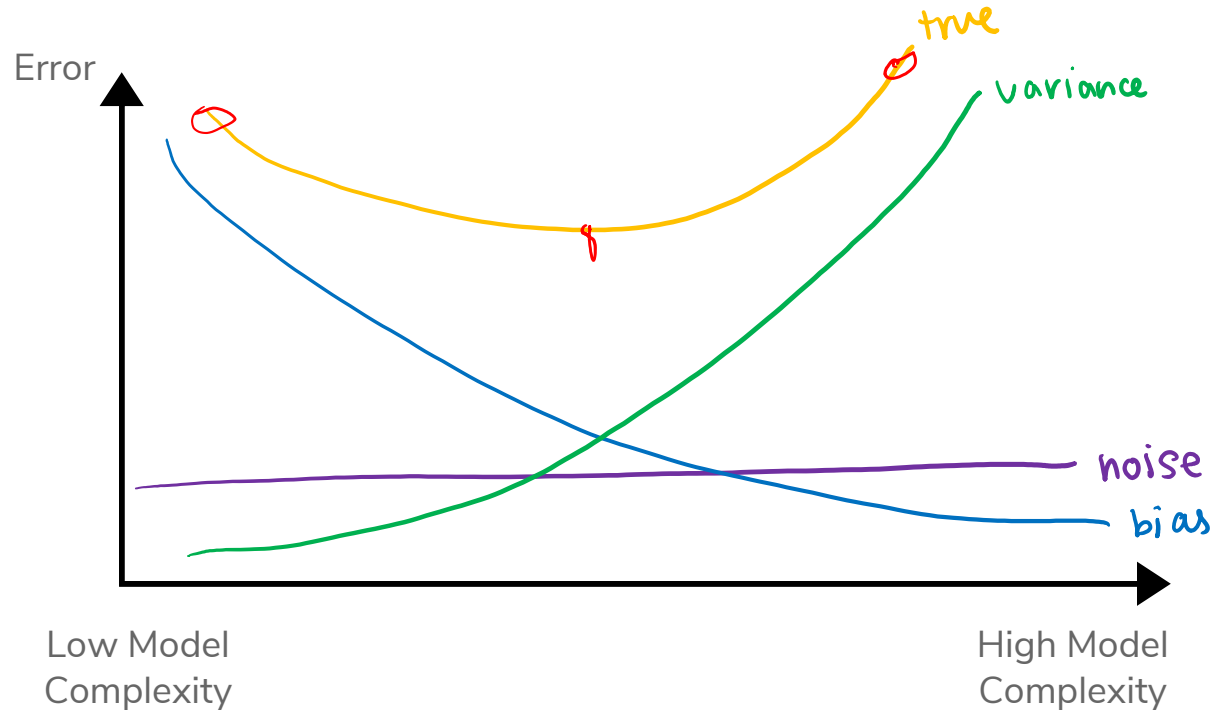
$$\text{variance} = E \left[ \underbrace{\bar{f}_{\hat{w}}(x)}_{\text{avg predictor}} - \underbrace{f_{\hat{w}}(x)}_{\text{specific predictor}} \right]^2$$

High complexity models tend to have high variance.\*

# Bias-Variance Tradeoff

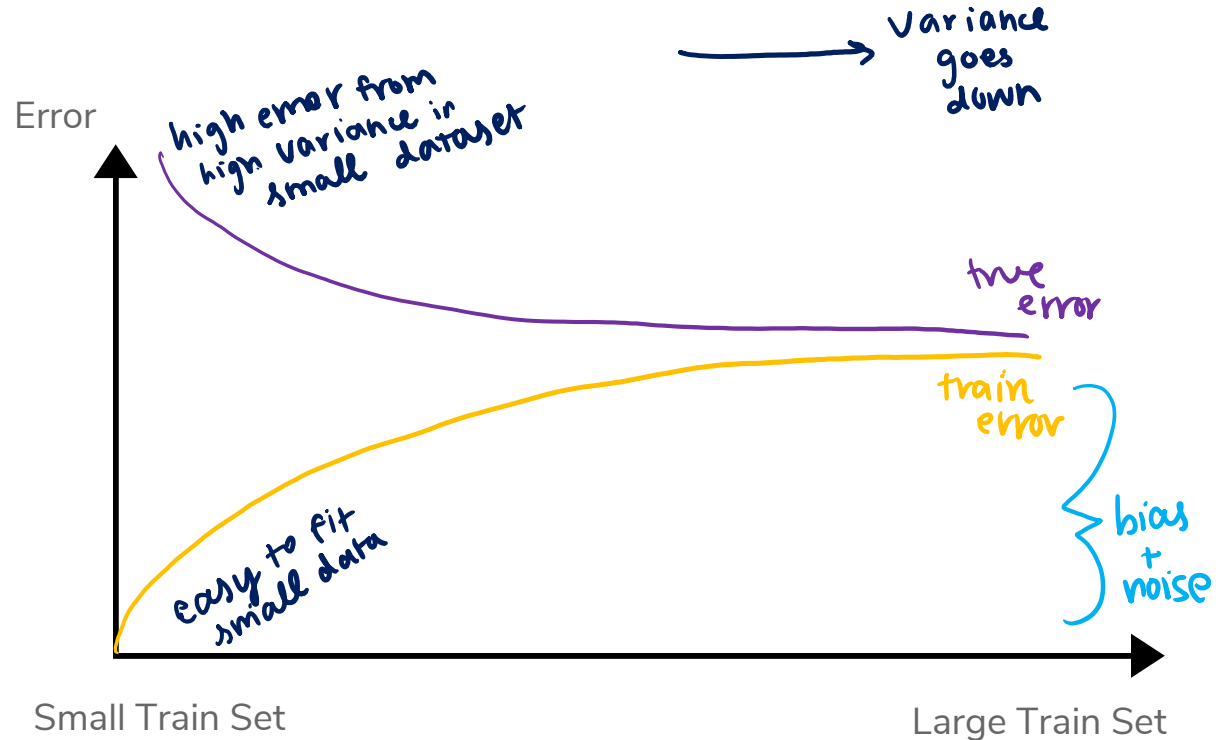
Visually, this looks like the following!

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Noise}$$



# Dataset Size

So far our entire discussion of error assumes a fixed amount of data. What happens to our error as we get more data?

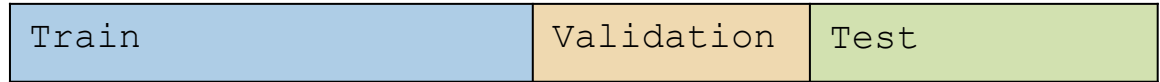


# Validation

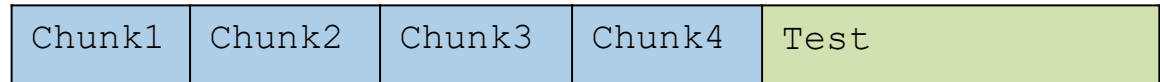
So far we have divided our dataset into train and test



Validation set: we can't use Test to choose our model complexity, so instead, break up Train into ANOTHER dataset



Cross Validation: Break up our training set into chunks

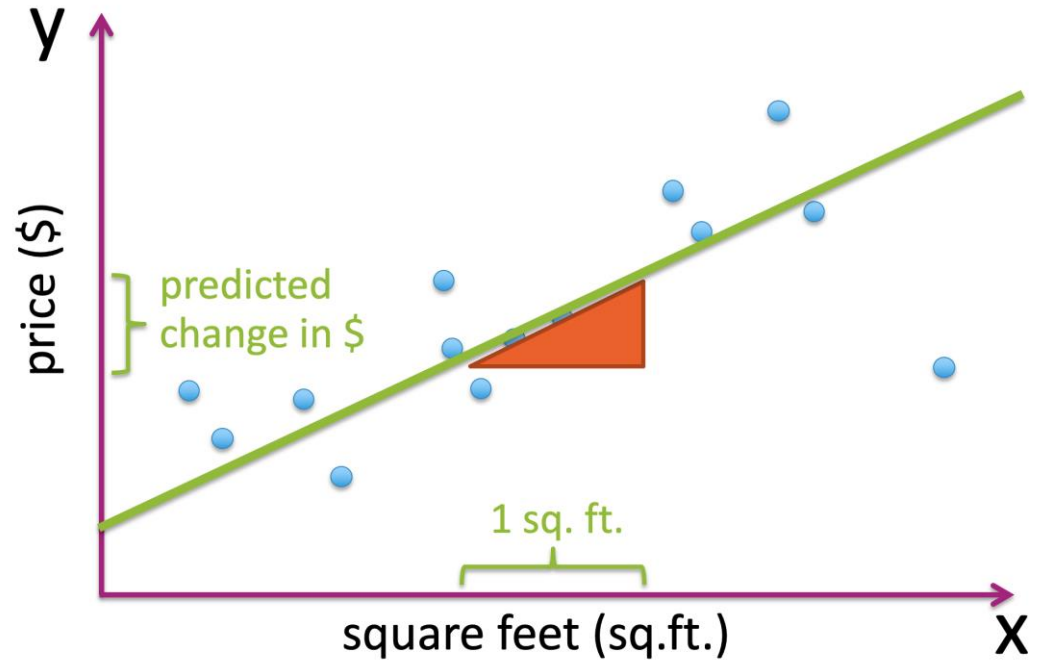




# Interpreting Coefficients

Interpreting Coefficients – Simple Linear Regression

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x$$

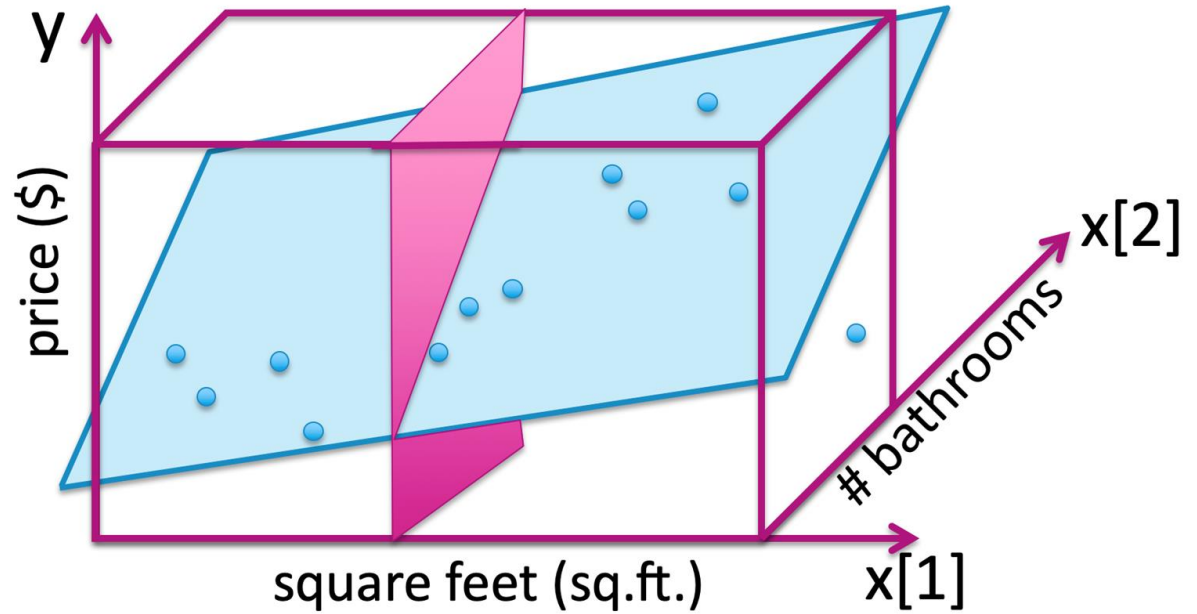


# Interpreting Coefficients

Interpreting Coefficients – Multiple Linear Regression

$$\hat{y} = \hat{w}_0 + \hat{w}_1x[1] + \hat{w}_2x[2]$$

Fix



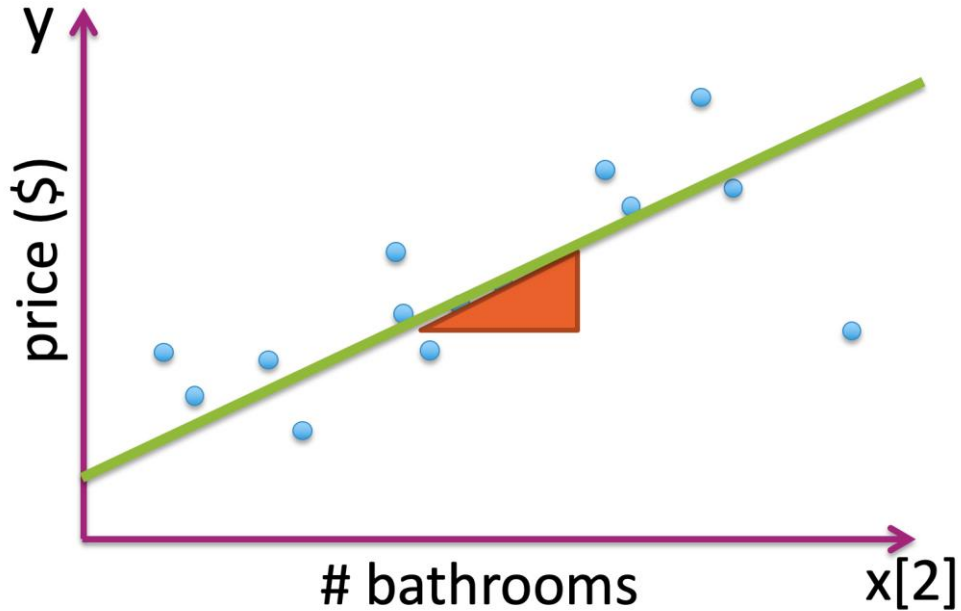
# Interpreting Coefficients

Interpreting Coefficients – Multiple Linear Regression

$$\hat{y} = \hat{w}_0 + \hat{w}_1x[1] + \hat{w}_2x[2]$$

Fix

Holding  $x[1]$  fixed!



# Interpreting Coefficients

This also extends for multiple regression with many features!

$$\hat{y} = \hat{w}_0 + \sum_{j=1}^D \hat{w}_j h_j(x)$$

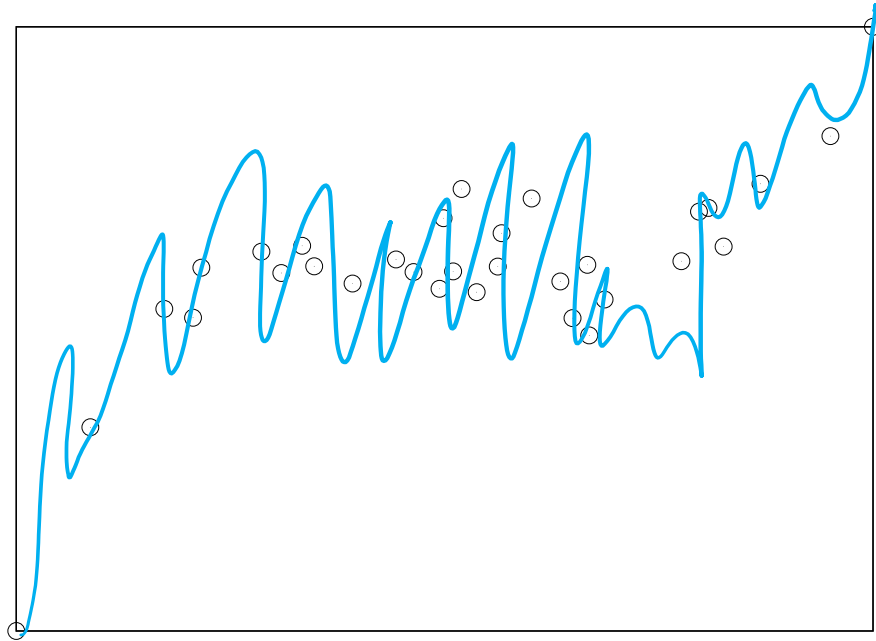
Interpret  $\hat{w}_j$  as the change in  $y$  per unit change in  $h_j(x)$  if all other features are held constant.

This is generally not possible for polynomial regression or if other features use same data input!

- Can't “fix” other features if they are derived from same input.

$$\hat{w} = [\hat{w}_0, \hat{w}_1, \dots, \hat{w}_D]$$

# Overfitting

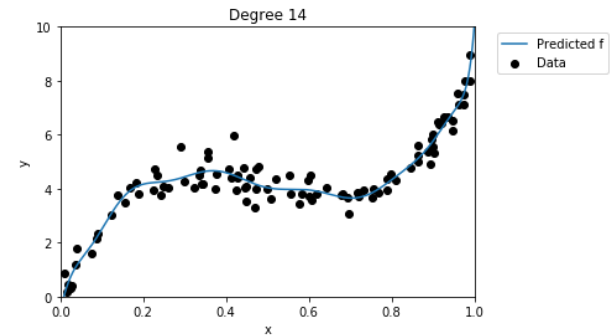
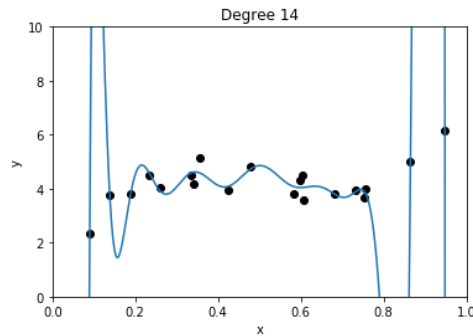


Often, overfitting is associated with very large estimated parameters  $\hat{w}$ !

# Number of Features

Overfitting is not limited to polynomial regression of large degree. It can also happen if you use a large number of features!

Why? Overfitting depends on how much data you have and if there is enough to get a representative sample for the complexity of the model.



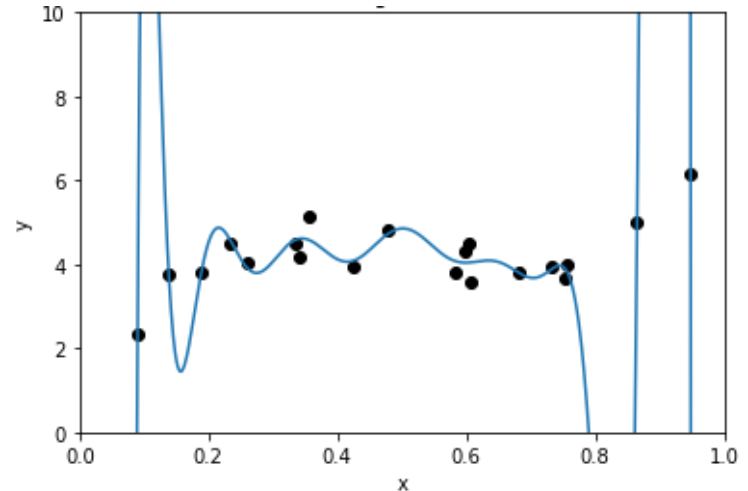
# Number of Features

How do the number of features affect overfitting?

## 1 feature

Data must include representative example of all  $(x, y)$  pairs to avoid overfitting

**HARD**



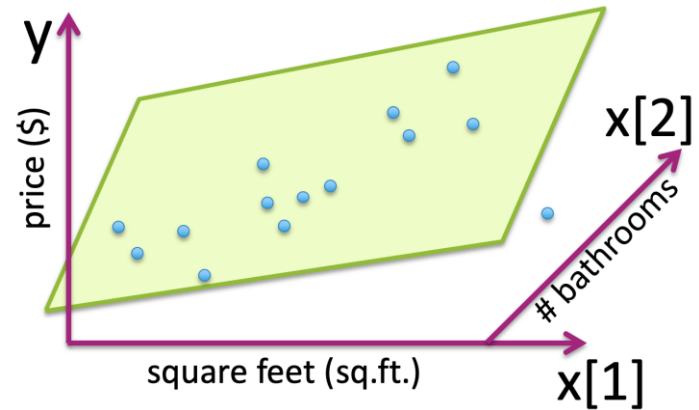
# Number of Features

How do the number of features affect overfitting?

## D features

Data must include representative example of all  $((x[1], x[2], \dots, x[D]), y)$  combos to avoid overfitting!

**MUCH HARDER!!**



Introduction to the **Curse of Dimensionality**.  
We will come back to this later in the quarter!



# Prevent Overfitting

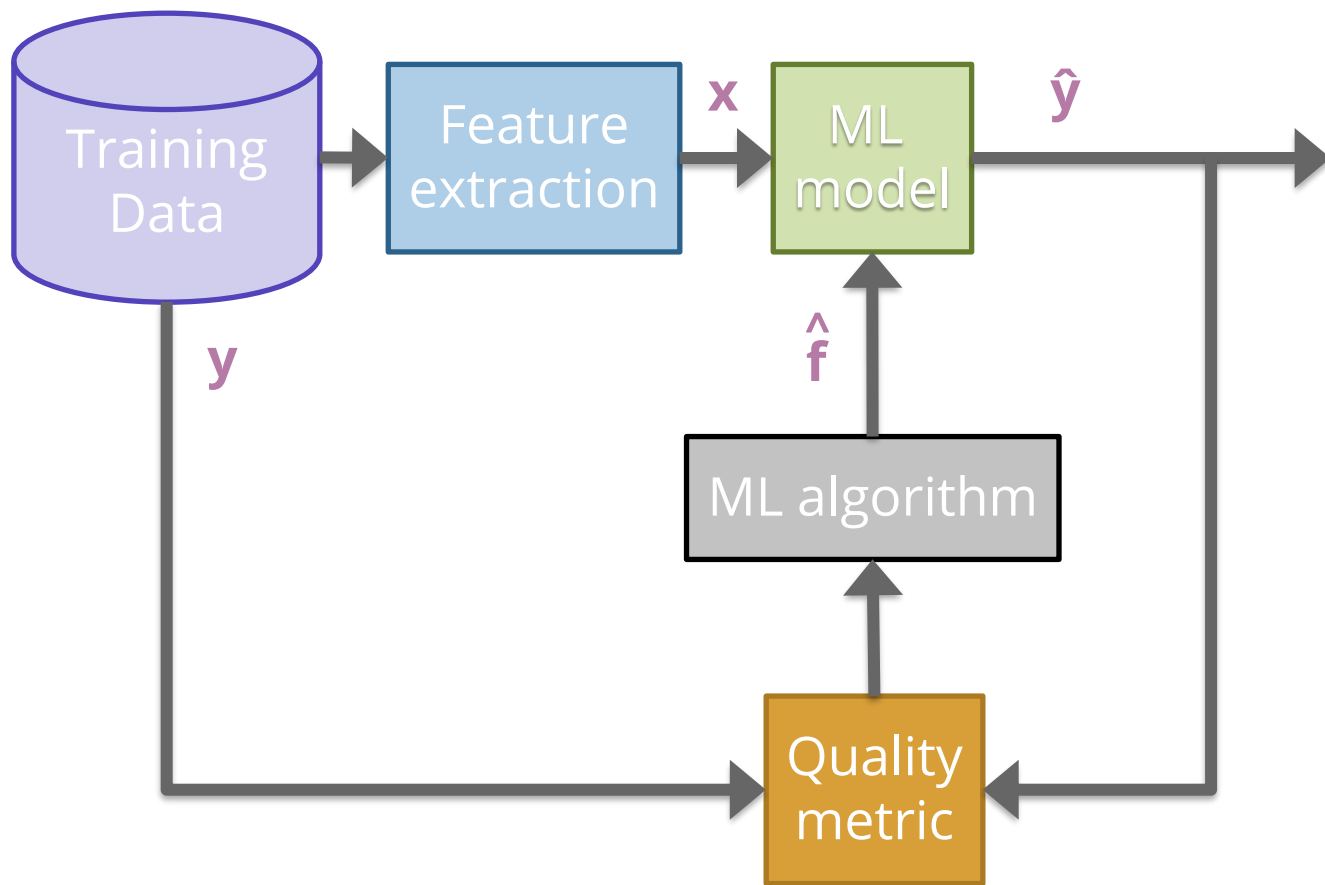
Last time, we saw we could use cross validation / validation set to pick which model complexity to use

- In the case of polynomial regression, we just chose degree  $p$
- For deciding which or how many features to use, there are a lot of choices!
  - For  $d$  inputs, there are  $2^d$  subsets of those features!

What if we use a model that wasn't prone to overfitting?

- Big Idea: Have the model self-regulate to prevent overfitting by making sure its coefficients don't get "too large"

This idea is called **regularization**.



# Regularization

Before, we used the quality metric that minimized loss

$$\hat{w} = \min_w L(w) \quad \text{RSS}$$

Change quality metric to balance loss with measure of overfitting

- $L(w)$  is the measure of fit
- $R(w)$  measures the magnitude of coefficients

$$\hat{w} = \min_w L(w) + \lambda R(w)$$

How do we actually measure the magnitude of coefficients?

# Magnitude

$$W = [w_0, w_1, \dots, w_D]$$

$R(w)$  = measure of overfitting

Come up with some number that summarizes the magnitude of the coefficients in  $w$ .

Sum?

$$R(w) = \sum_{j=0}^D w_j \quad \begin{array}{l} w_0 = 10,000 \\ w_1 = -10,003 \end{array} \rightarrow -3$$

Sum of absolute values?

$$R(w) = \sum_{j=0}^D |w_j| \triangleq \|w\|_1 \quad \text{L1 Norm}$$

Sum of squares?

$$R(w) = \sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2 \quad \text{L2 Norm}$$

# Ridge Regression

Change quality metric to minimize

$$\hat{w} = \min_w \text{RSS}(W) + \lambda \|w\|_2^2$$

$\lambda$  is a tuning parameter (hyperparameter) that changes how much the model cares about the regularization term.

What if  $\lambda = 0$ ?

$$\min_w \text{RSS}(w) \rightarrow \hat{w}_{LS} \text{ least squares}$$

What if  $\lambda = \infty$ ?

$$\cancel{w \neq \vec{0} \rightarrow \infty} \\ w = \vec{0} \quad \text{RSS}(\vec{w})$$

$\lambda$  in between?

$$0 \leq \|\hat{w}\|_2^2 \leq \|\hat{w}_{LS}\|_2^2$$

# Poll Everywhere

Think 

1.5 Minutes

How does  $\lambda$  affect the bias and variance of the model? For each underlined section, select “Low” or “High” appropriately.

When  $\lambda = 0$

*Complex*

The model has (Low / High) Bias and (Low / High) Variance.

When  $\lambda = \infty$

*Simple*

The model has (Low / High) Bias and (Low / High) Variance.

[pollev.com/cse416](https://pollev.com/cse416)



Brain Break



# Demo: Ridge Regression

See Jupyter Notebook for interactive visualization.

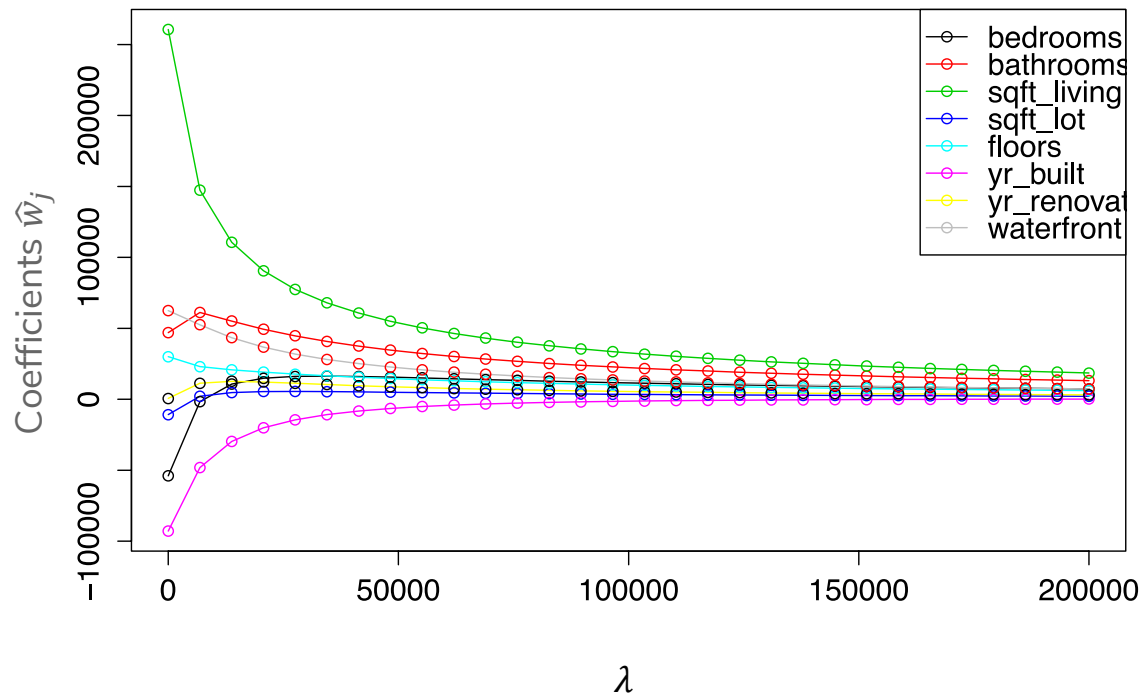
Shows relationship between

- Regression line
- Residual Sum of Squares Quality Metric
  - Also called Ordinary Least Squares
- Ridge Regression Quality Metric
- Coefficient Paths





# Coefficient Paths



Think 

2 min

[pollev.com/cse416](https://pollev.com/cse416)

### How should we choose the best value of $\lambda$ ?

- ~~▪ Pick the  $\lambda$  that has the smallest  $RSS(\hat{w})$  on the **training set**~~
- ~~▪ Pick the  $\lambda$  that has the smallest  $RSS(\hat{w})$  on the **test set**~~
- ✓ ▪ Pick the  $\lambda$  that has the smallest  $RSS(\hat{w})$  on the **validation set**
- ✗ ▪ Pick the  $\lambda$  that has the smallest  $RSS(\hat{w}) + \lambda \|\hat{w}\|_2^2$  on the **training set**
- ~~▪ Pick the  $\lambda$  that has the smallest  $RSS(\hat{w}) + \lambda \|\hat{w}\|_2^2$  on the **test set**~~
- ✗ ▪ Pick the  $\lambda$  that has the smallest  $RSS(\hat{w}) + \lambda \|\hat{w}\|_2^2$  on the **validation set**
  - ~~▪ Pick the  $\lambda$  that results in the smallest coefficients~~
  - ~~▪ Pick the  $\lambda$  that results in the largest coefficients~~
  - ~~▪ None of the above~~

# Regularization

At this point, I've hopefully convinced you that regularizing coefficient magnitudes is a good thing to avoid overfitting!

You:



We might have gotten a bit carried away, it doesn't ALWAYS make sense...

$$w = [w_0, w_{rest}] \quad w_{rest} = [w_1, \dots, w_p]$$

## The Intercept

For most of the features, looking for large coefficients makes sense to spot overfitting. The one it does not make sense for is the **intercept**.

We shouldn't penalize the model for having a higher intercept since that just means the  $y$  value units might be really high

- My demo before does this wrong and penalizes  $w_0$  as well

Two ways of dealing with this

- Change the measure of overfitting to not include the intercept

$$\min_{w_0, w_{rest}} \text{RSS}(w_0, w_{rest}) + \lambda \|w_{rest}\|_2^2$$

- Center the  $y$  values so they have mean 0
  - This means forcing  $w_0$  to be small isn't a problem

# Scaling Features

The other problem we looked over is the “scale” of the coefficients.

Remember, the coefficient for a feature increase per unit change in that feature (holding all others fixed in multiple regression)

Consider our housing example with (*sq. ft.*, *price*) of houses

- Say we learned a coefficient  $\hat{w}_1$  for that feature
- What happens if we change the unit of  $x$  to square **miles**?  
Would  $\hat{w}_1$  need to change?
  - It would need to get bigger since the prices are the same but its inputs are smaller

This means we accidentally penalize features for having large coefficients due to having small value inputs!

# Scaling Features

Fix this by **normalizing** the features so all are on the same scale!

$$\tilde{h}_j(x_i) = \frac{h_j(x_i) - \mu_j(x_1, \dots, x_N)}{\sigma_j(x_1, \dots, x_N)}$$

Where

The mean of feature  $j$ :

$$\mu_j(x_1, \dots, x_N) = \frac{1}{N} \sum_{i=1}^N h_j(x_i)$$

The standard deviation of feature  $j$ :

$$\sigma_j(x_1, \dots, x_N) = \sqrt{\frac{1}{N} \sum_{i=1}^N (h_j(x_i) - \mu_j(x_1, \dots, x_N))^2}$$

**Important:** Must scale the test data and all future data using the means and standard deviations **of the training set!**

- Otherwise the units of the model and the units of the data are not comparable!

# Recap

**Theme:** Use regularization to prevent overfitting

**Ideas:**

- How to interpret coefficients
- How overfitting is affected by number of data points
- Overfitting affecting coefficients
- Use regularization to prevent overfitting
- How L2 penalty affects learned coefficients
- Visualizing what regression is doing
- Practicalities: Dealing with intercepts and feature scaling