

Bias/Variance and Cross-Validation

Review: Linear Regression

Overfitting and Underfitting

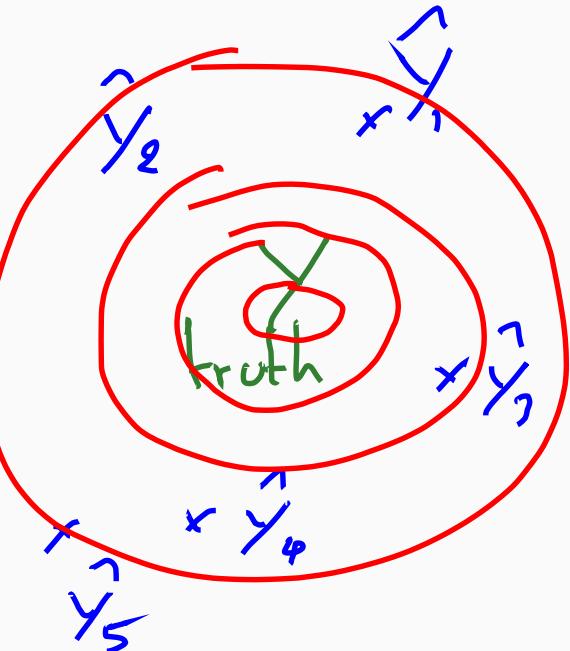
The Bias/Variance Tradeoff

Cross-Validation

K-fold Cross-Validation

Subset Selection of Predictors

Bias and Variance

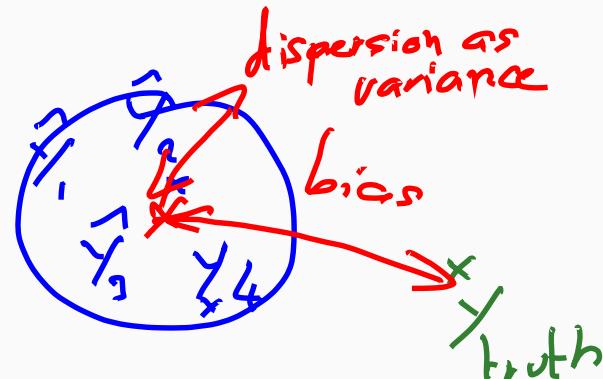


average all $\hat{y}_i = \hat{Y}$ (truth)

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \hat{Y}$$

UNBIASED

$$E[\hat{Y}] = \hat{Y}$$



BIASED

Bias and Variance

measurement bias

ML algo eg linear regression
bias: assuming linear relationship between X and Y

+ data scientists,
preference of ML algo
feature engineering expertise

inductive bias

restrictive bias
set of models the ML algo will consider

e.g., stepwise regression

preference bias
guide the learning process to prefer some models over other models

One Goal of Data Science: Make Future Predictions

One goal is to make accurate *predictions* on future (unseen) data.

1. Define a business goal.

e.g. make Tesla cars the most dependable vehicles on the market

refine the question

2. Collect training data.

e.g. Tesla cars' event logs + historical record of parts replaced

3. Train a model.

e.g. **features**: event statistics, **target**: time until failure

4. Deploy the model.

e.g. monitor cars' events in real time, send mechanics to replace parts that will soon fail

Review: Linear Regression

We assume the world is built on linear relationships. Under that assumption, we can model the relationship between *features* and a *target* like this:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

target
outcome
dependent variable
regressand ...
estimated $\rightarrow \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p + \hat{\epsilon}$

feature
predictor
independent input

coefficients
model's parameters

sampling error
(noise irreducible...)

Review: Linear Regression

We assume the world is built on linear relationships. Under that assumption, we can model the relationship between *features* and a *target* like this:

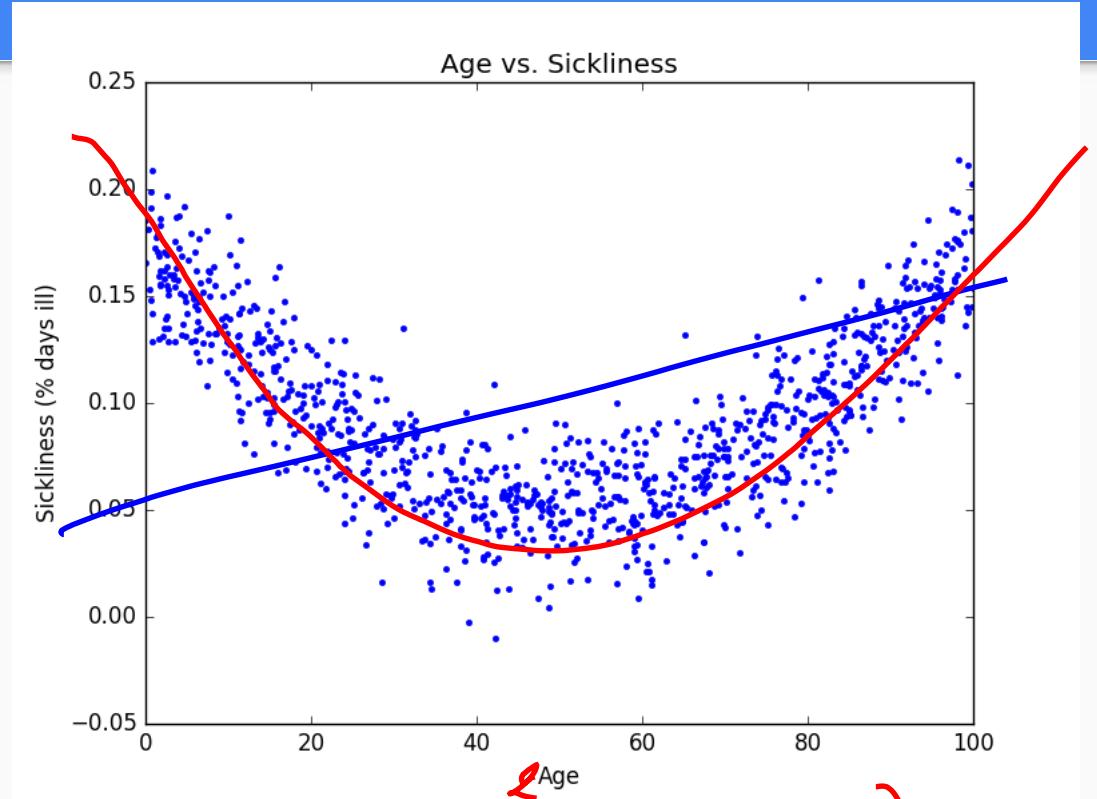
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

The diagram illustrates the components of a linear regression equation. The target variable Y is highlighted in yellow. The model parameters $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ are highlighted in blue. The features X_1, X_2, \dots, X_p are highlighted in green. The sampling error ϵ is highlighted in red.

Review: Linear Regression

We can make linear regression non-linear by inserting extra “interaction” features or higher-order features.

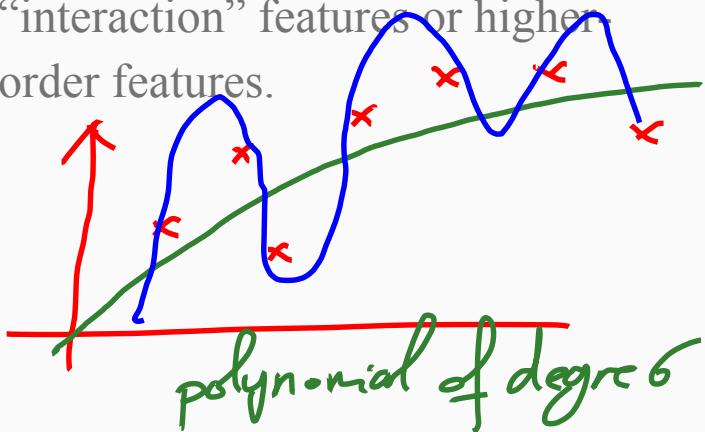
eg.
 $\text{sickliness} = \beta_0 + \beta_1 \text{age}$
feature engineering age²
 $\text{sickliness} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2$



$$R^2 = \frac{\text{variance explained}}{\text{variance total}} = \frac{\sum y_i^2}{\sum (y_i - \bar{y})^2} (= \rho_{\text{correlation}})$$

Review: Linear Regression

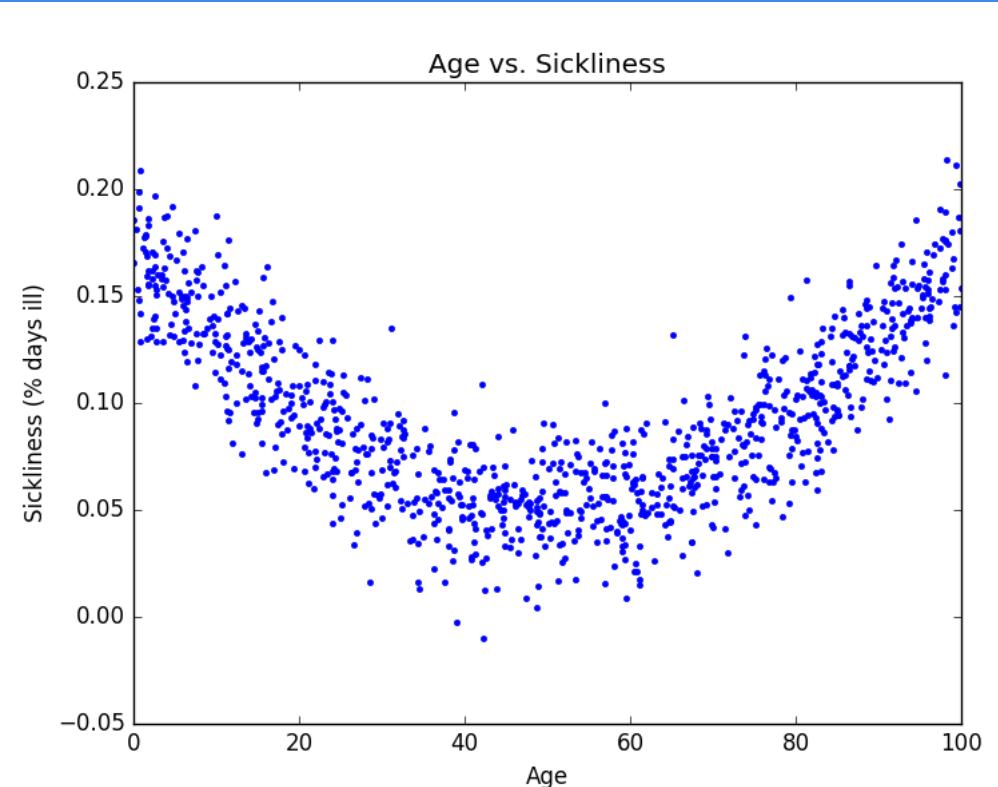
We can make linear regression non-linear by inserting extra “interaction” features or higher order features.



Examples:

$$\text{Sickliness} = \beta_0 + \beta_1 * \text{age}$$

$$\text{Sickliness} = \beta_0 + \beta_1 * \text{age} + \beta_2 * \text{age}^2$$



$$Y = \hat{Y} + \varepsilon$$

$$\text{Var}(Y) = \text{Var}(\hat{Y}) + \text{Var}(\varepsilon) + 2\text{Cov}(\hat{Y}, \varepsilon)$$

$$\sigma_Y^2 = \sigma_{\hat{Y}}^2 + \sigma_{\varepsilon}^2$$

random noise

$$R^2 = \frac{\sigma_{\hat{Y}}^2}{\sigma_Y^2} = 1 - \frac{\sigma_{\varepsilon}^2}{\sigma_Y^2}$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2$$

$$= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2$$

significance of
 this matter
 not just
 β_1, β_2

$$Y = \underbrace{\beta_0}_{\text{intercept}} + \underbrace{\beta_1}_{\text{slope}} X_1$$

Is R² all that matters?

We *could* just keep inserting interaction features until R² = 1.

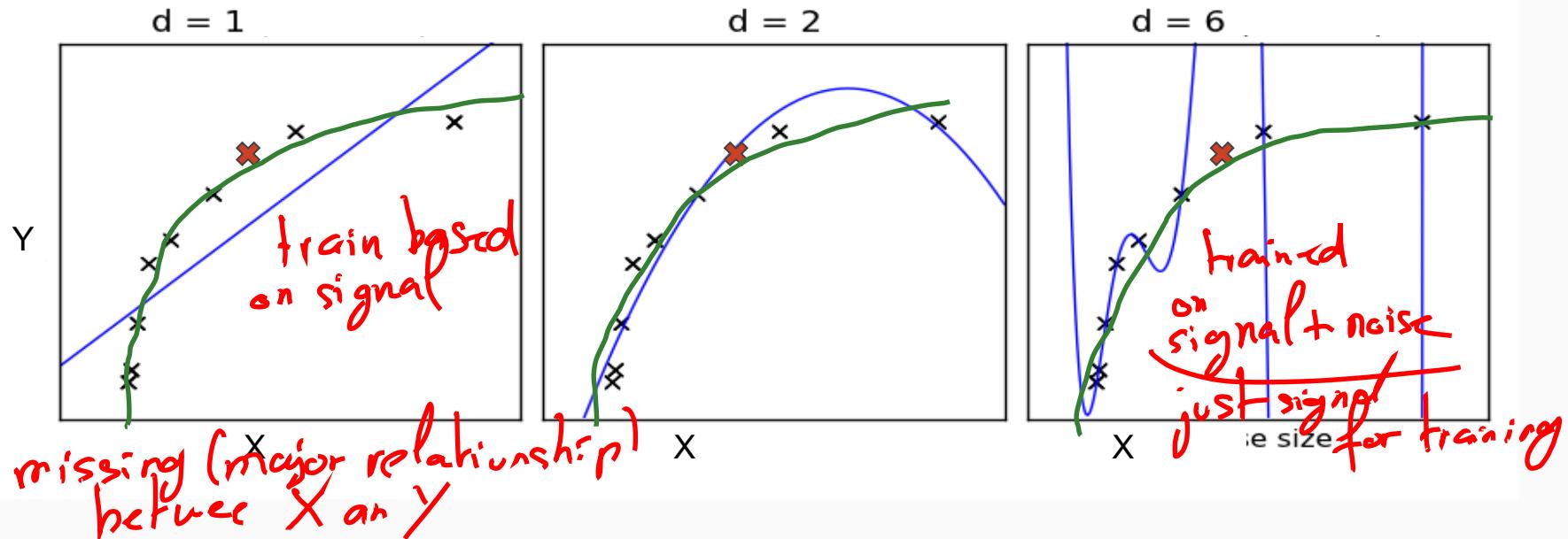
Boom. I solved data science. Here's my idea:

```
def train_super_awesome_perfect_model(X, y):  
    while True:  
        model = LinearRegression()  
        model.fit(X, y)  
        if calculate_r2(model, X, y) >= 0.999:  
            return model  
        else:  
            X = insert_random_interaction_feature(X)
```

Why is this a
bad idea?

oh yeah!

Oh the woes of overfitting...



Underfitting and Overfitting

Underfitting: The model doesn't fully capture the relationship between predictors and the target. The model has *not* learned the data's signal. The model is not flexible enough.

→ What should we do if our model underfits the data?

Overfitting: The model has tried to capture the sampling error. The model has learned the data's signal *and* the noise. I.e., the model attributes to signal that which is truly noise. The model is too flexible.

→ What should we do if our model overfits the data?

The Bias/Variance Tradeoff

We assume the true predictor/target relationship is given by an unknown function plus some sampling error:

$$Y = f(X) + \epsilon$$

*y = f + ε
f and ε are independent*

We estimate the true (unknown) function by fitting a model over the training set

$$\hat{Y} = \hat{f}(X)$$

Let's evaluate this model using a test observation (x_0, y_0) drawn from the population.

What is the model's expected squared prediction error on this test observation?

$$y_0 \rightarrow y \quad \text{MSE} \quad f_{\text{truth}} \quad E[(y_0 - \hat{f}(x_0))^2] = \dots$$

y₀ → y → f_{truth} → f(x₀) is fixed → f(x₀) is fixed → constant

The Bias/Variance Tradeoff

$$E[X^2] \longrightarrow \text{Var}[X] = E[X]$$

$$\begin{aligned}\text{Var}[X] &= E[(X - E(X))^2] \\ &= E[X^2 - 2E[X]X + E[X]^2] \\ &= E[X^2] - \underbrace{2E[X]E[X]}_{\text{const}} + E[X]^2\end{aligned}$$

$$\begin{aligned}E[ax+by] &= aE[x]+bE[y] \\ E[ax] &= aE[x]\end{aligned}$$

$$= E[X^2] - E[X]^2$$

$$E[X^2] = \text{Var}[X] + E[X]^2 \quad (\text{a})$$

The Bias/Variance Tradeoff

$$E[y] = f + \epsilon$$
$$E[y] = E[f + \epsilon] = E[f] + E[\epsilon]$$

o noise

$$E[f(x_o)] = f(x_o)$$

(b)

$$\boxed{E[y] = f}$$

$$\text{Var}[y] = E[(y - E[y])^2] = E[\epsilon^2]$$

$$\boxed{\text{Var}[y] = \text{Var}[\epsilon]}$$

(c)

$$\text{Var}[\epsilon] + E[\epsilon_o]^2$$

(a)

The Bias/Variance Tradeoff

$$\begin{aligned} E[\hat{y}] &= E[(\hat{f} + \varepsilon)] \\ &= E[\hat{f}] + E[\varepsilon] \\ &= \underbrace{E[\hat{f}]}_{\text{est}} + \underbrace{E[\varepsilon]}_{\cancel{\text{est}}} \quad \cancel{E[\varepsilon]\varepsilon = 0} \quad E[XY] = E[X]E[Y] \\ &= \hat{f} E[\hat{f}] \quad \text{iff } X \text{ and } Y \text{ are independent} \end{aligned}$$

$$E[\hat{y}] = \hat{f} E[\hat{f}] (d)$$

The Bias/Variance Tradeoff

$$MSE = E[(y - \hat{f})^2]$$

$$= E[y^2 - 2\hat{f}y + \hat{f}^2]$$

(a)

$$\text{Var}[y] + E[y]^2$$

Variance of
irreducible error

$$- 2E[\hat{f}y] + E[\hat{f}]^2$$

$$- 2E[\hat{f}]$$

(d)

$$+ E[\hat{f}]$$

(a)

$$\text{Var}[\hat{f}] + E[\hat{f}]^2$$

Variance of
model prediction

$$MSE = \text{Var}[\hat{f}] + (f - 2E[\hat{f}])^2 + E[\hat{f}]^2 + \text{Var}[\epsilon]$$

$$(f - E[\hat{f}])^2 \rightarrow E[(f - \hat{f})]^2 \rightarrow \text{bias}^2$$

The Bias/Variance Tradeoff

Our model's expected squared prediction error will depend on (1) the variability of y_0 and (2) the variability of the training set used to train our model. We can break this into three pieces:

$$E[(y_o - \hat{f}(x_0))^2] = \dots = \text{Var}(\hat{f}(x_0)) + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\epsilon)$$

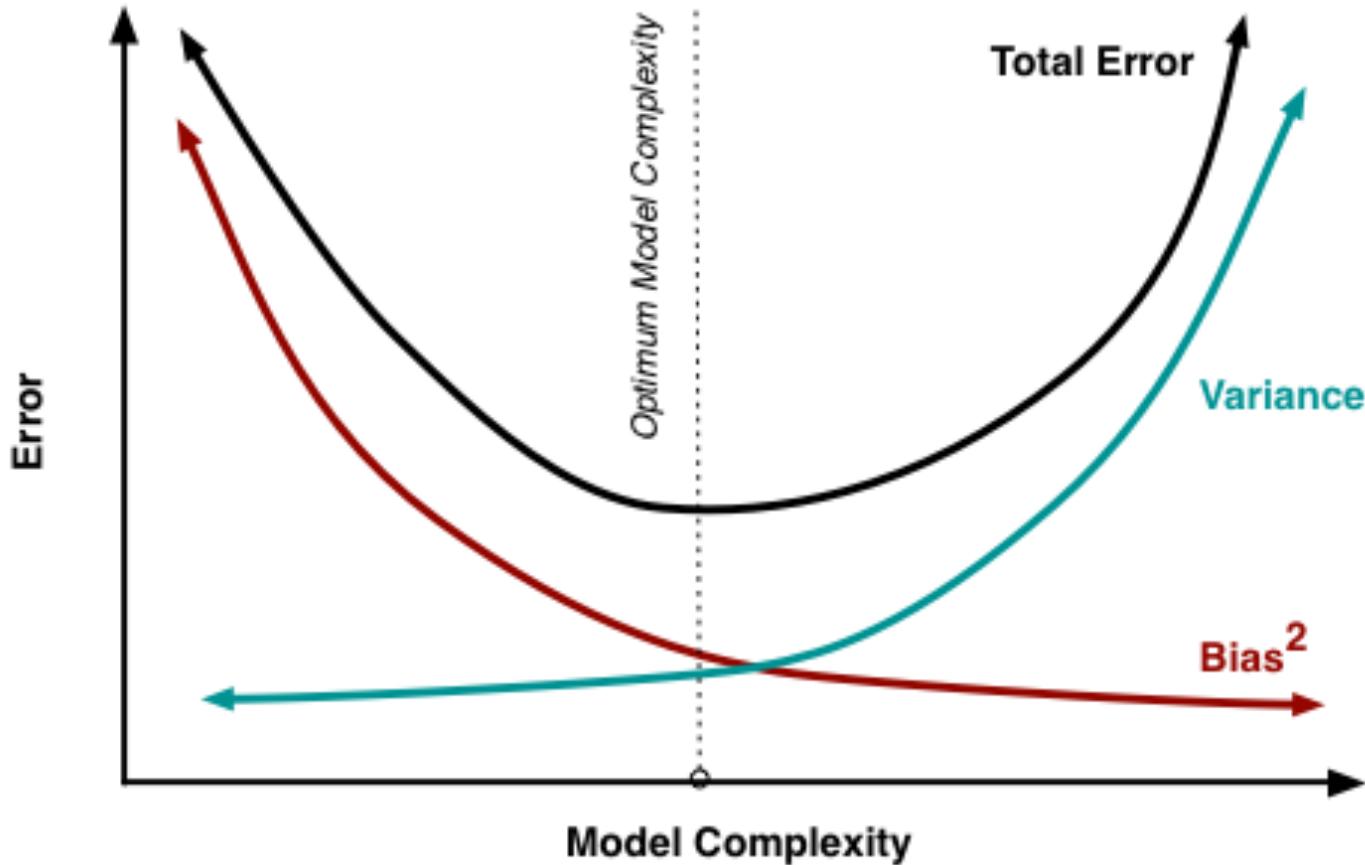
The variance of our model's prediction of x_0 over all possible training sets

The difference between the true prediction and our model's average prediction over all possible training sets

$$\blacktriangleright \text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0)$$

The variance of the irreducible error.

The Bias/Variance Tradeoff



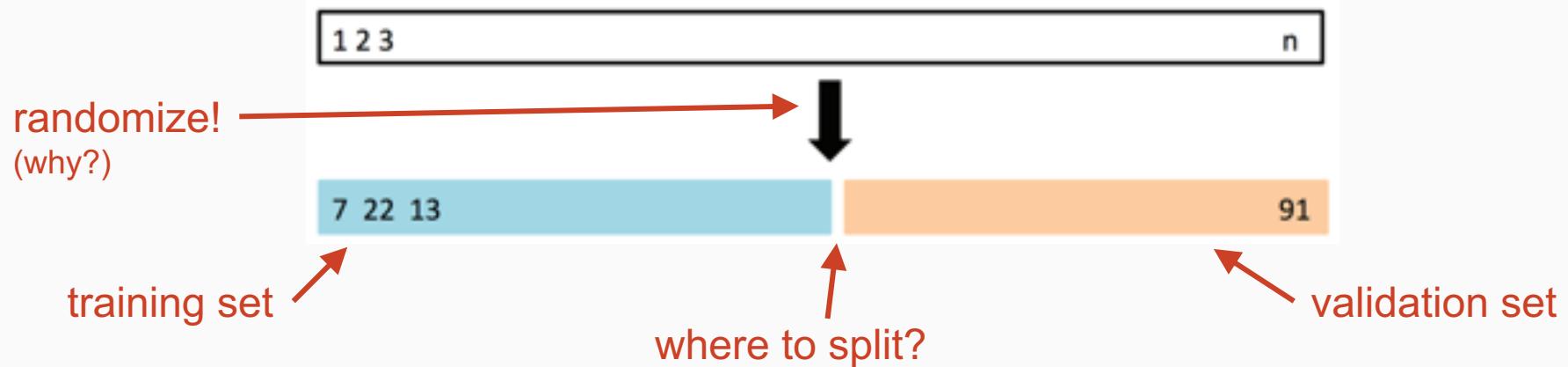
How is the bias/variance tradeoff related to underfitting and overfitting?

How can we find the best tradeoff point?
I.e. The optimum model complexity

Cross-Validation

Main idea: Don't use all your data for training.

Instead: Split your data into a “training set” and a “validation set”.



Cross-Validation

1. Split your data into training/validation sets.

70/30 or 90/10 splits are commonly used

2. Use the training set to train several models of varying complexity.

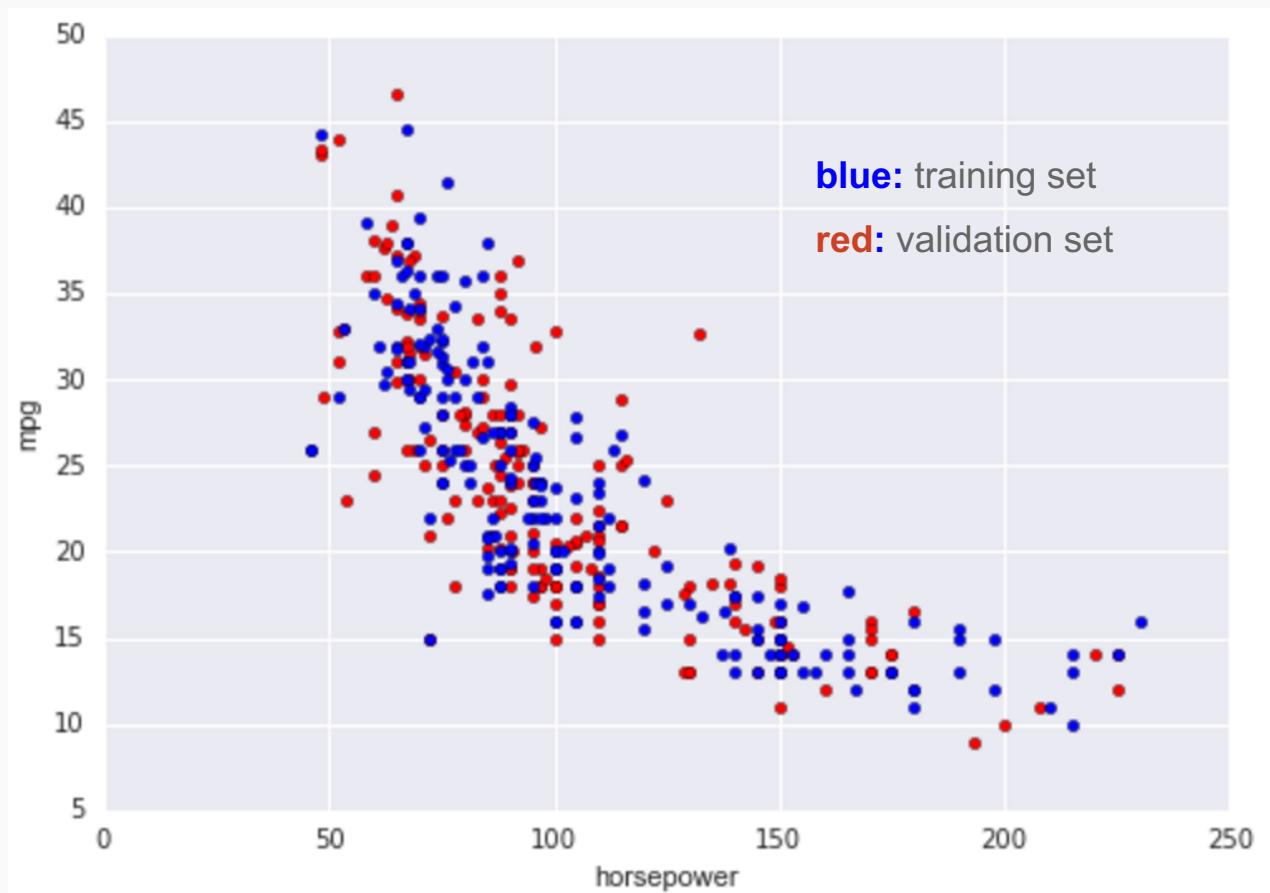
e.g. linear regression (w/ and w/out interaction features), neural nets, decision trees, etc.
(we'll talk about hyperparameter tuning, grid search, and feature engineering later)

3. Evaluate each model using the validation set.

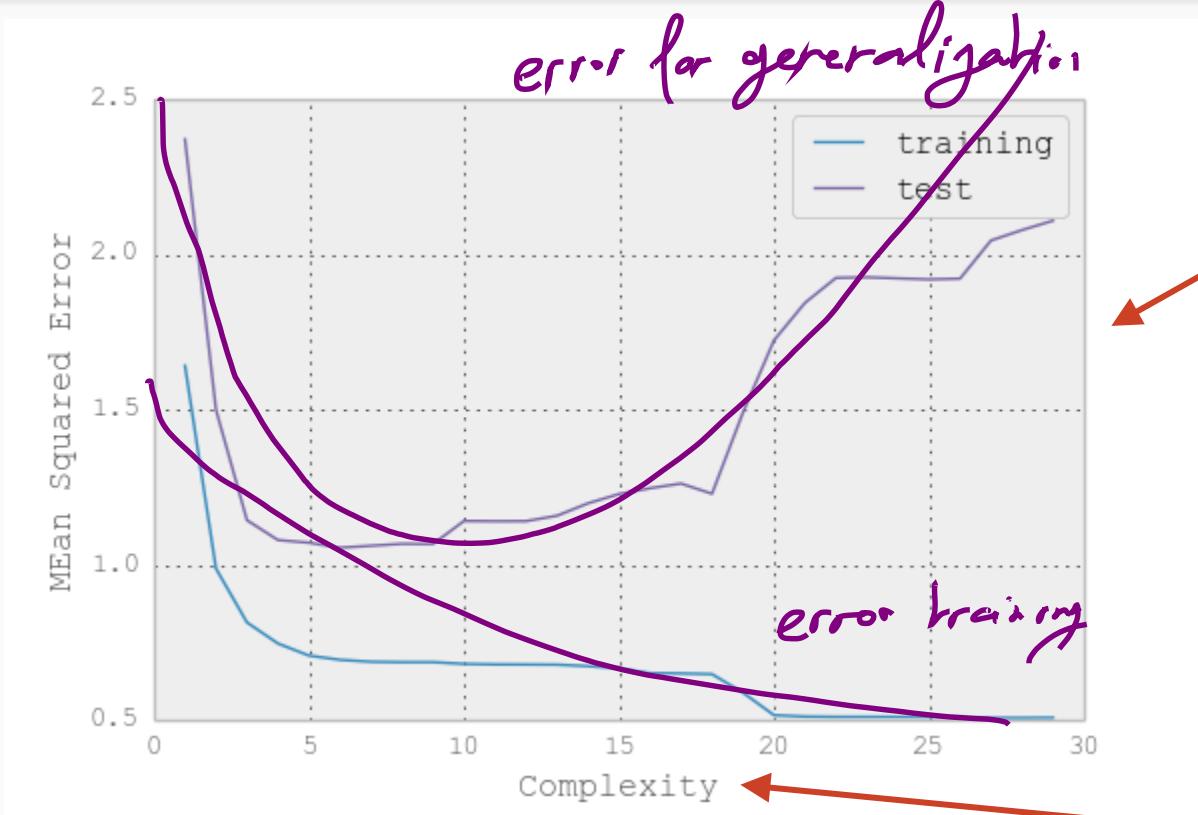
calculate R^2 , MSE, accuracy, or whatever you think is best

4. Keep the model that performs best over the validation set.

Let's predict MPG from horsepower



Cross-Validation Example



You will see
this shape all
the time!

You will wrestle
with the
bias/variance
tradeoff
constantly...

E.g. linear regression w/
varying degree of
polynomial

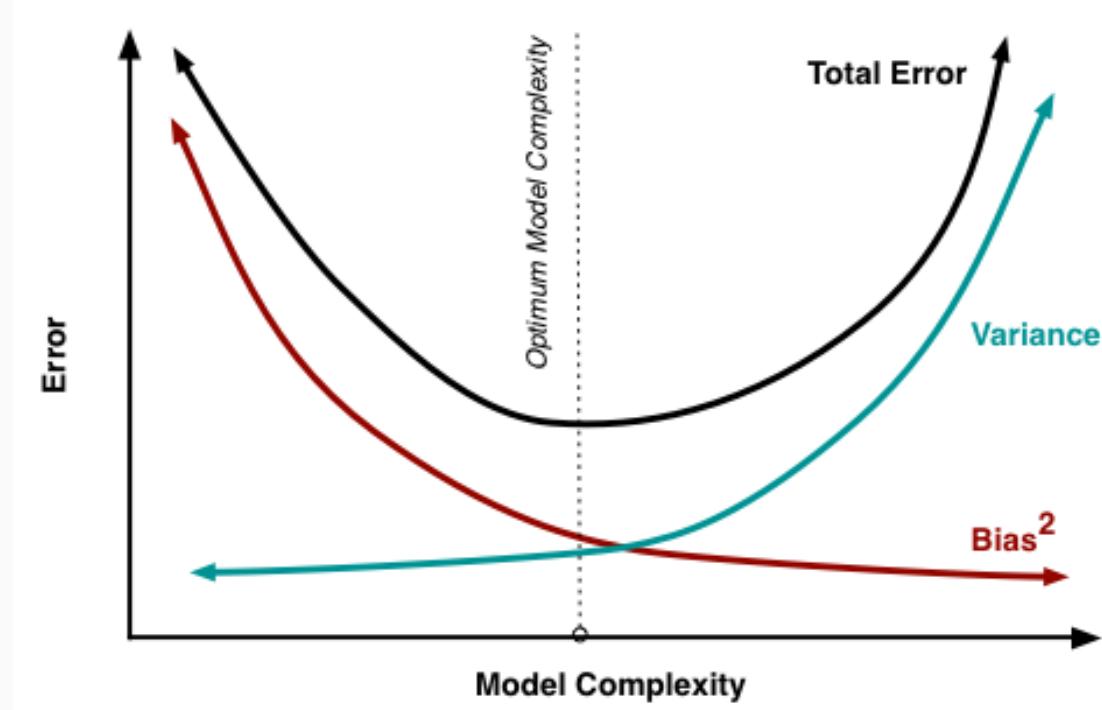
Recall our goal: Making accurate future predictions

Fitting the training set perfectly is easy.

How?

Fitting future (unseen) data is *not* easy.

Cross validation helps us choose a model that performs well on unseen data.

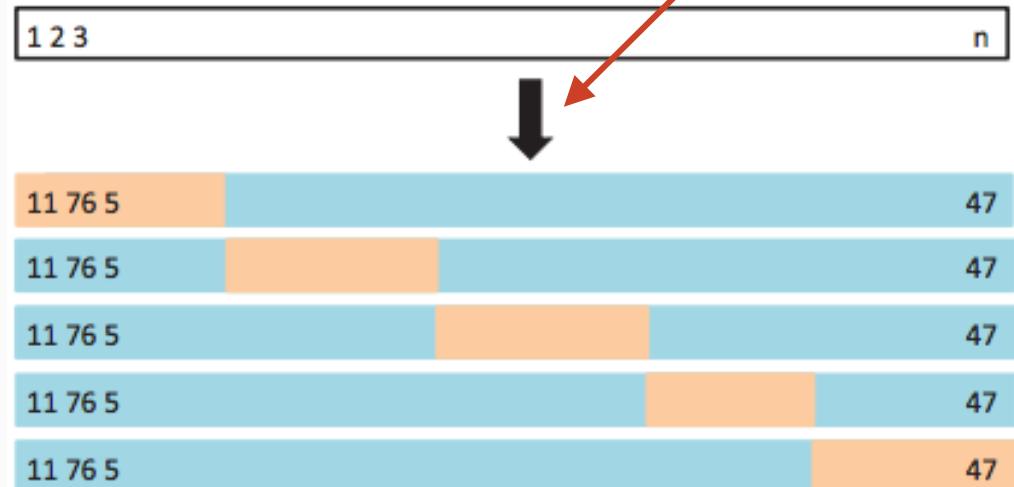


k-Fold Cross-Validation

1. Split the dataset into k “folds”.

Commonly, $k=5$ or $k=10$

2. Train using $(k-1)$ folds.
Validate using the one left-out fold. Record a validation metric such as RSS or accuracy.



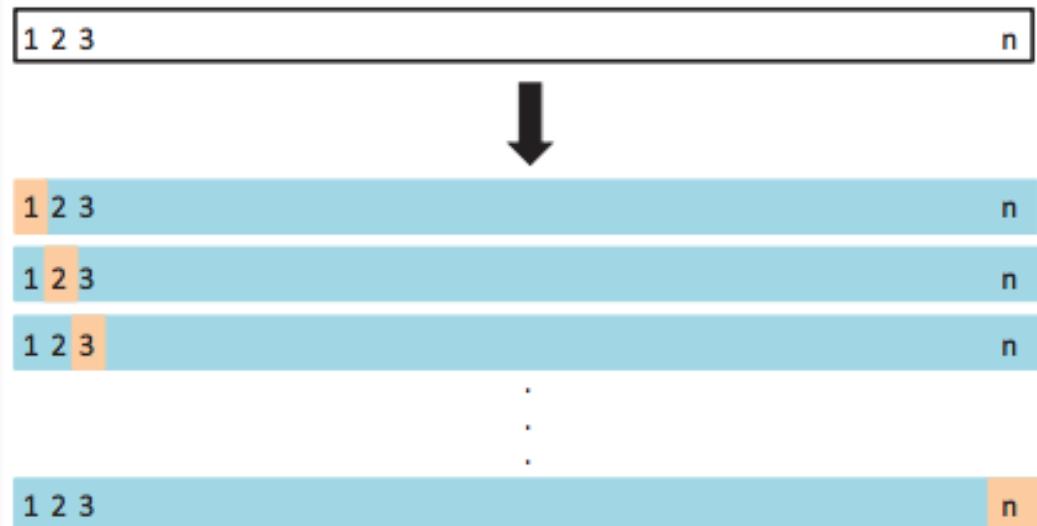
3. Train k models, leaving out a different fold for each one.
4. Average the validation results.

Leave-one-out Cross-Validation

Assume we have n training examples.

A special case of k-fold CV is when $k=n$. This is called *leave-one-out cross-validation*.

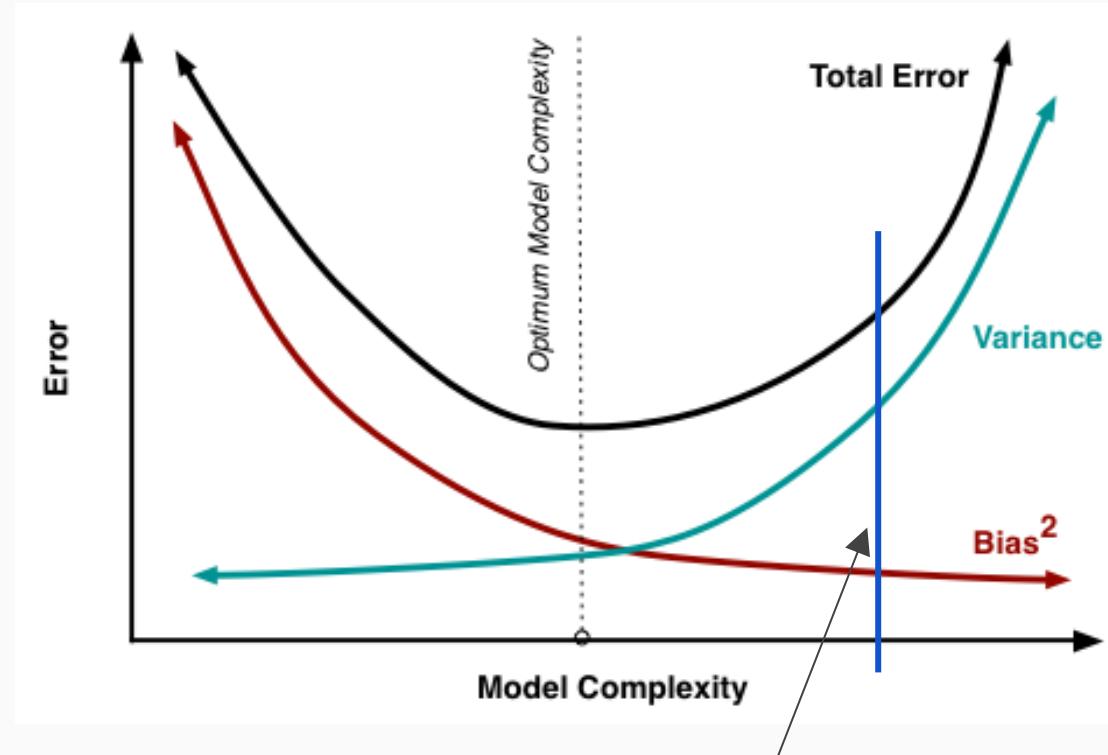
Useful (only) if you have a tiny dataset where you can't afford a large validation set.



Overfitting in high dimensions is easy, even with simple models.

If our data has high dimensionality (many many predictors), then it becomes easy to overfit the data.

Even linear regression might be too complex of a model for high dimensional data (and the smaller the dataset, the worse this problem is).



Linear regression in high dimensions

“HELP, my model is overfitting!”

You have a few options.

1. **Get more data:** not always possible/practical
2. **Subset Selection:** keep only a subset of your predictors (i.e, dimensions)
3. **Regularization:** restrict your model's parameter space
4. **Dimensionality Reduction:** project the data into a lower dimensional space

Subset Selection



Best subset: Try every model. Every possible combination of p predictors

Computationally intensive. 2^p possible subsets of p predictors

High chance of finding a “good” model by random chance.

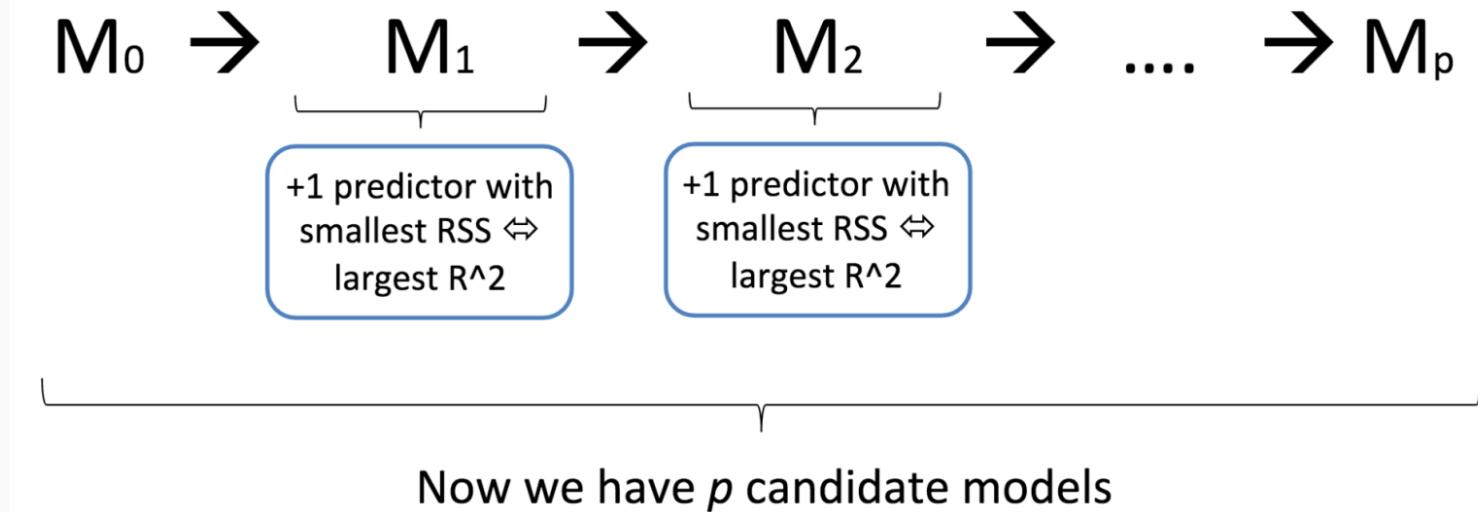
... A sort-of monkeys-Shakespeare situation ...

Stepwise: Iteratively pick predictors to be in/out of the final model.

Forward, backward, forward-backward strategies

Handwritten notes:
P features
2 choices for each feature include or not
 $2 \times 2 \times \dots \times 2$
first second last
feature
gp

Forward Stepwise Selection



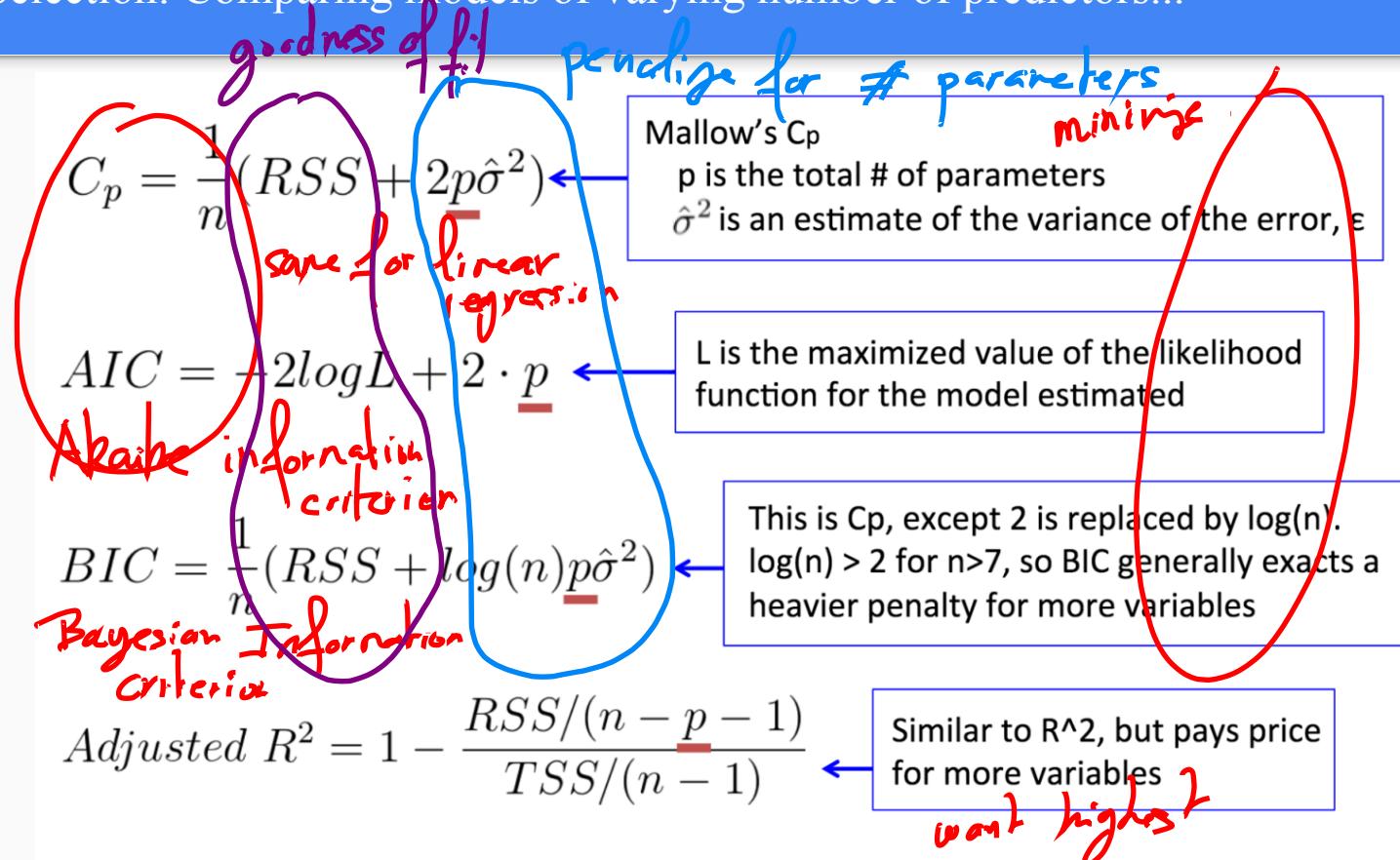
Are RSS and R^2 good ways to decide amongst the resulting p candidates?

Answer: Don't use RSS or R^2 for this part.

Use Mallow's C_p , or AIC, or BIC, or Adjusted R^2 .

... or just use cross-validation with any error measurement.

Subset Selection: Comparing models of varying number of predictors...



Side Note: Can show AIC and Mallow's Cp are equivalent for linear case

Subset Selection: Comparing models of varying number of predictors...

```

=====
OLS Regression Results
=====
Dep. Variable: y R-squared: 0.933
Model: OLS Adj. R-squared: 0.928
Method: Least Squares F-statistic: 211.8
Date: Mon, 03 Nov 2014 Prob (F-statistic): 6.30e-27
Time: 14:45:06 Log-Likelihood: -34.438
No. Observations: 50 AIC: 76.88
Df Residuals: 46 BIC: 84.52
Df Model: 3
Covariance Type: nonrobust
=====

      coef  std err      t      P>|t|   [95.0% Conf. Int.]
-----
x1      0.4687  0.026  17.751  0.000      0.416  0.522
x2      0.4836  0.104   4.659  0.000      0.275  0.693
x3     -0.0174  0.002  -7.507  0.000     -0.022 -0.013
const    5.2058  0.171  30.405  0.000      4.861  5.550
=====
Omnibus: 0.655 Durbin-Watson: 2.896
Prob(Omnibus): 0.721 Jarque-Bera (JB): 0.360
Skew: 0.207 Prob(JB): 0.835
Kurtosis: 3.026 Cond. No. 221.
=====
```