# Module 2: The Modeling Process

# Module 2 Objectives:

**At the conclusion of this week, you should be able to:**
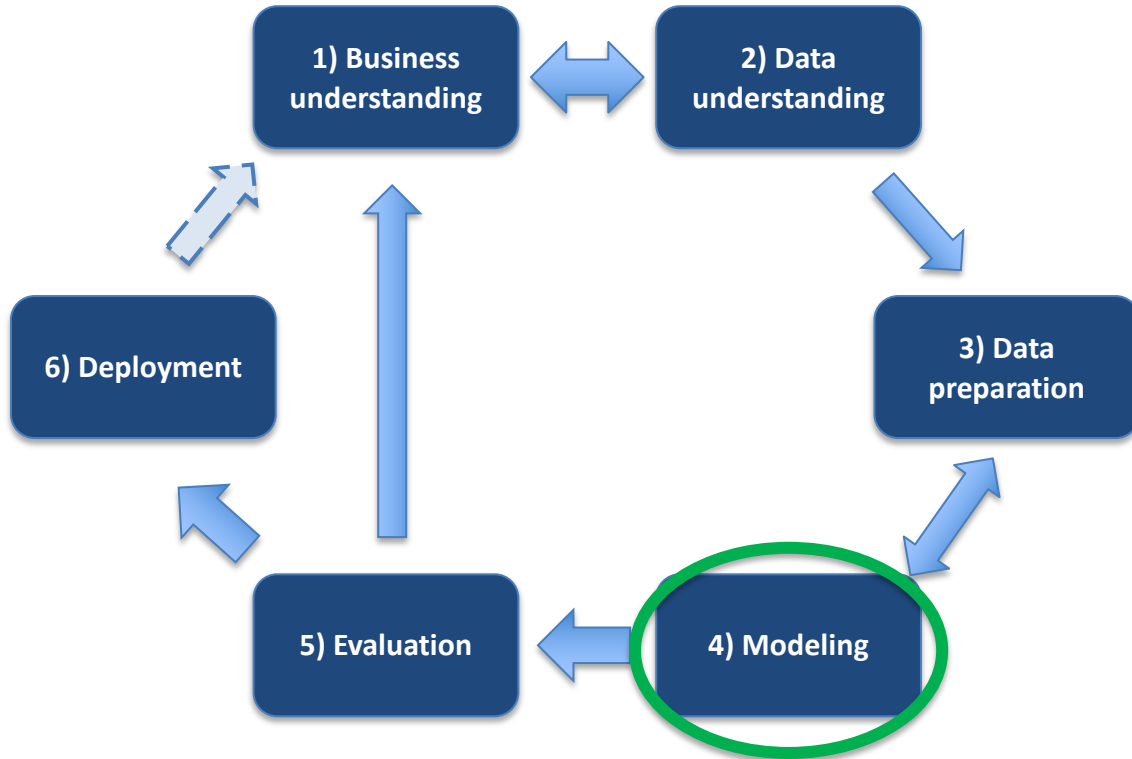
1) Describe the steps to develop a ML model

2) Explain the bias-variance tradeoff

3) Identify possible sources of data leakage and strategies to prevent it

# Building a Model

# CRISP-DM Process

# Creating a Model

Past Observations

| | Neighbor-hood | School district | Square footage | Number of bedrooms | Year built |
|---|---|---|---|---|---|
| House 1 | Weycroft | Wake | 3400 | 4 | 2010 |
| House 2 | Horton Creek | Wake | 4200 | 5 | 2008 |
| House 3 | Cary Park | Chatham | 3250 | 4 | 2012 |
| ... | ... | ... | ... | ... | ... |

Targets

| Market sale price |
|---|
| $612,000 |
| $675,000 |
| $520,000 |
| ... |

Training

**Model** $y_i = f(X_i) + \varepsilon$

New data

Prediction

# Components of a Model



**Selection of features**

...tions

**Loss (cost) function**

Targets

| | Number of bedrooms | Year built |
|---|---|---|
| | | 2010 |
| | | 2008 |
| House 3 | Cary Park | Chatham | 3250 | 4 | 2012 |
| ... | ... | ... | ... | ... | ... |

| Market sale price |
|---|
| $612,000 |
| $675,000 |
| $520,000 |
| ... |

Training

**Model** $y_i = f(X_i) + \varepsilon$

**Choice of algorithm**

**Values for hyperparameters**

# Modeling Process

# Feature Selection

# What are Features?

**Features**

| | Neighbor-hood | School district | Square footage | Number of bedrooms | Year built |
|---|---|---|---|---|---|
| House 1 | Weycroft | Wake | 3400 | 4 | 2010 |
| House 2 | Horton Creek | Wake | 4200 | 5 | 2008 |
| House 3 | Cary Park | Chatham | 3250 | 4 | 2012 |
| … | … | … | … | … | … |

# How to Define Features

Case Study: Outage Prediction

# Methods of Feature Selection

- Domain expertise

- Visualization

- Statistical correlations

- Modeling

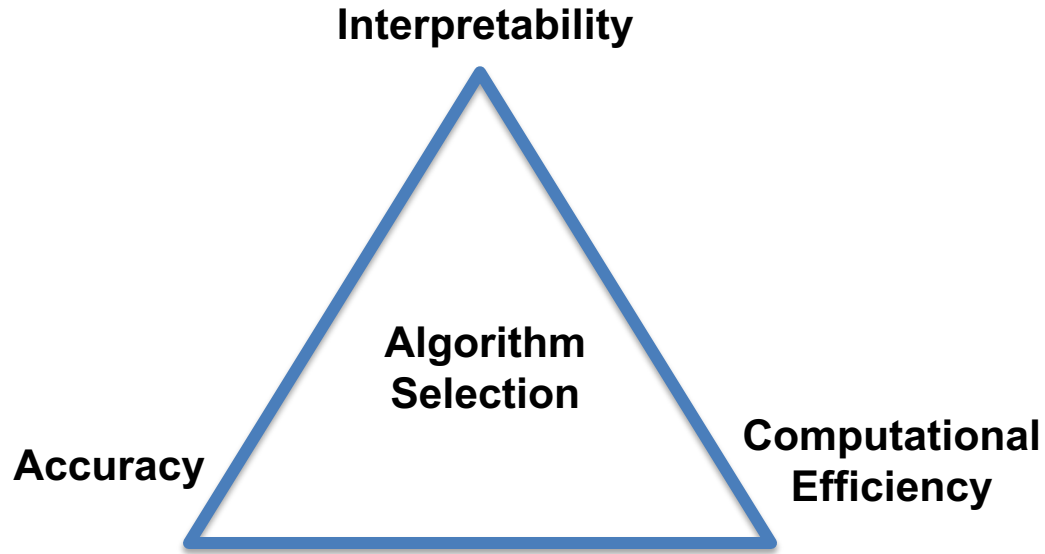Including too few features is usually much worse than including too many!

Algorithm Selection
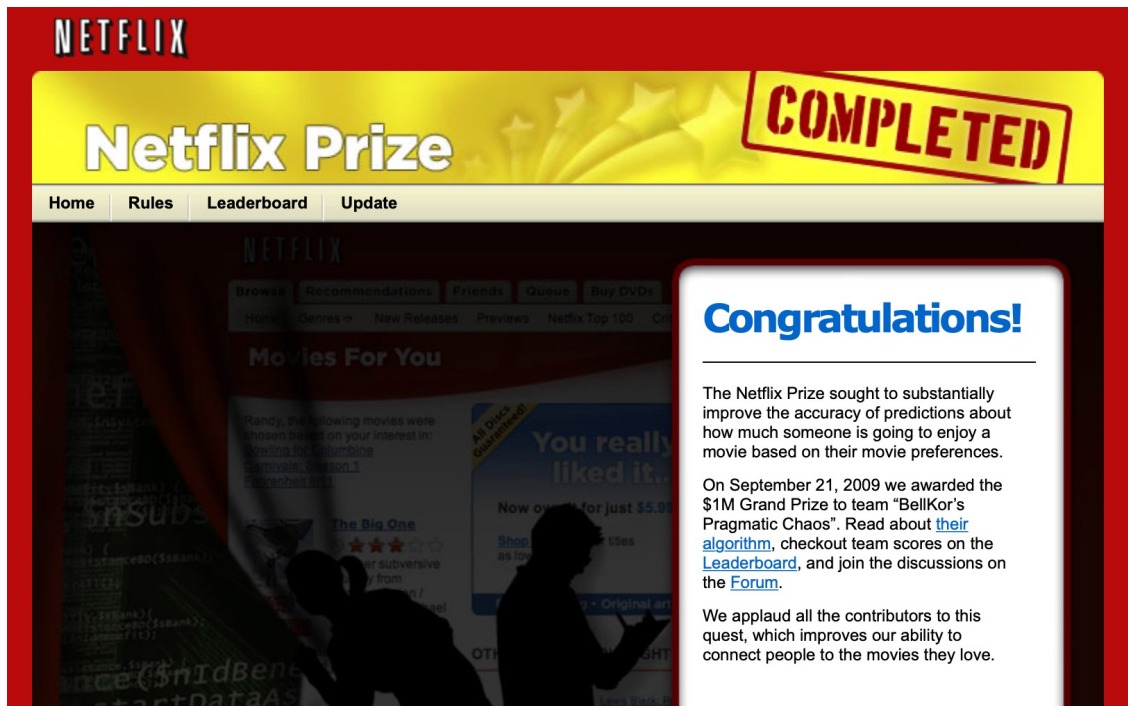
Duke
PRATT SCHOOL of ENGINEERING

# Algorithm Selection

**"No free lunch theorem"**

# Algorithm Selection

# Netflix Example

# Bias – Variance Tradeoff

Duke
PRATT SCHOOL of
ENGINEERING

# Model Complexity

# Bias and Variance

- **Bias** is error introduced by modeling a real life problem using a simpler model that is unable to fully capture the underlying patterns in data

- **Variance** refers to the sensitivity of the model to small fluctuations in the data, because it models fine patterns which may just be noise

**Low bias**    **High bias**

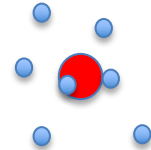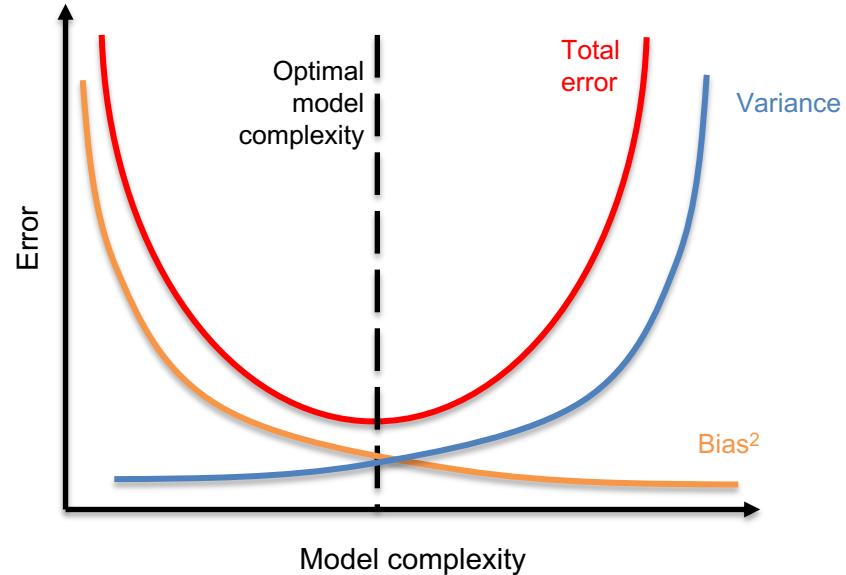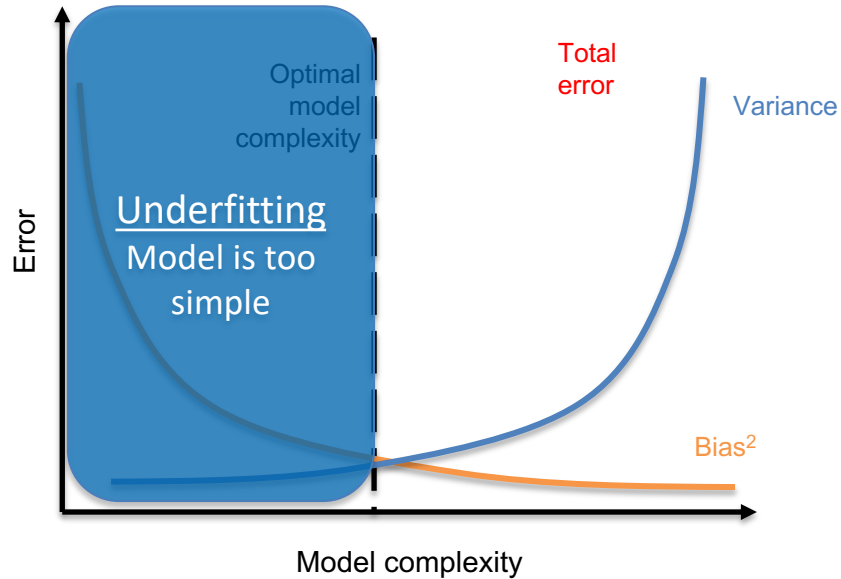**Low variance**    **High variance**

● Target    ● Predictions

# Bias – Variance Tradeoff

- Simpler models often have **higher** bias and **lower** variance

- Complex models typically have **lower** bias but **higher** variance

- Total Error $=$ Bias$^2$ + Var + $\sigma_e^2$

# Underfitting vs. Overfitting

# Underfitting vs. Overfitting

# Underfitting vs. Overfitting

**Underfitting**
Model is too simple

**Good Fit**
Model fits well, with some error

**Overfitting**
Model is too complex



*Image source: Scikit Learn documentation*

# Test & Validation Sets

Duke

PRATT SCHOOL of
ENGINEERING

# Training & Test Sets

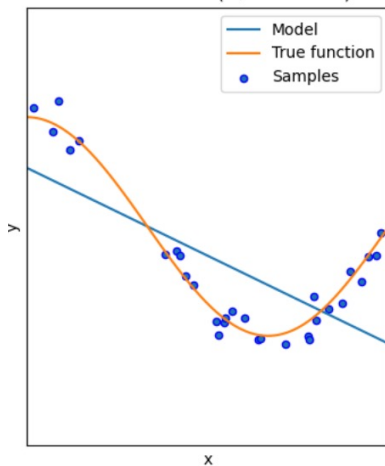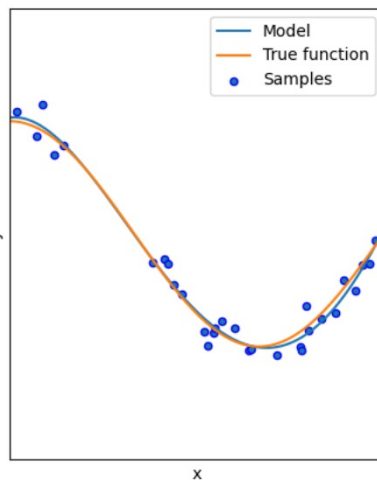- Goal of predictive modeling is to create a model that makes accurate predictions on new unseen data

- We cannot estimate performance on data we do not have, so instead we split our data into two sets

  - **Training set** - build and train the model

  - **Test set** – Evaluate model performance performance

Training  (~80-90%)

Test (~10-20%)

Use for model building & training

Use to evaluate final model

# Data Leakage

- **"Data leakage"** occurs when some of our test set data "leaks" into model building and influences the development of the model

- For example, if we use all of our data to select our features, or compare algorithms

- This **invalidates the estimated performance** of the model and causes it to be overoptimistic
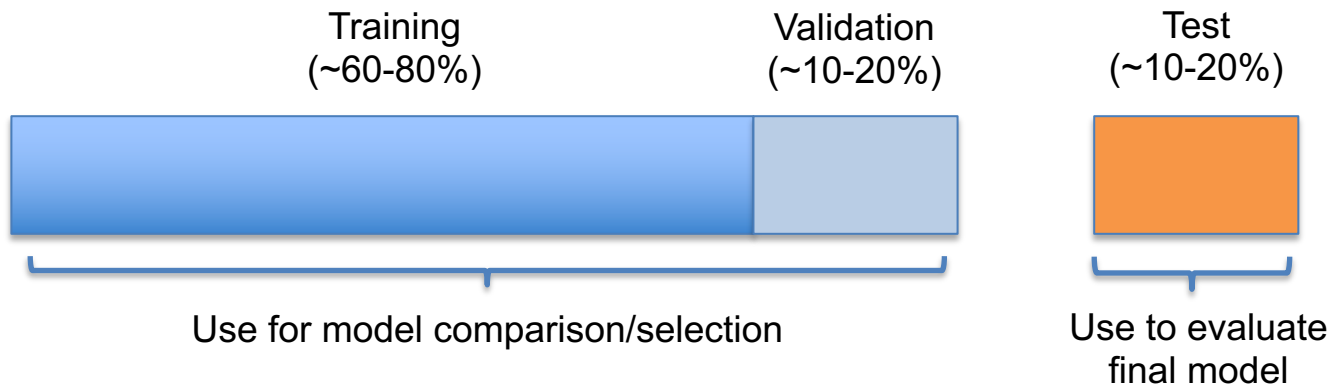
# Validation Sets

- Often we want to compare models to select the optimal model
- If we use the test set to compare model performance, it is not longer an unbiased indicator of performance
- Instead, we split our training set further into training and validation sets
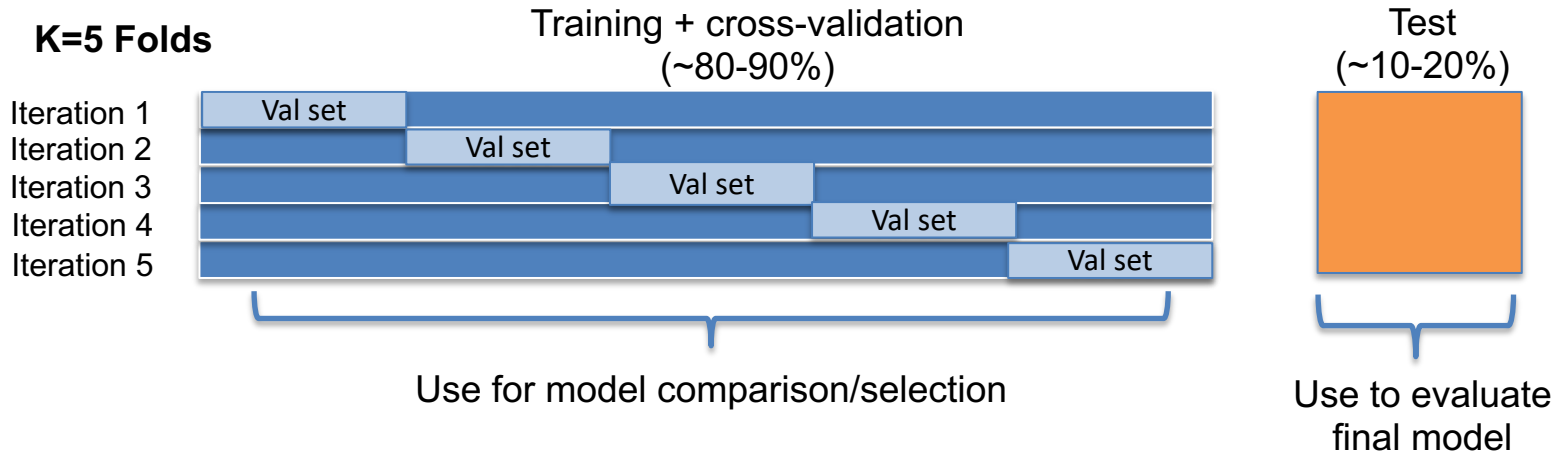- We use the validation set for model selection, and report performance on the test set



Training (~60-80%)    Validation (~10-20%)    Test (~10-20%)

Use for model comparison/selection    Use to evaluate final model

# K-Folds Cross Validation

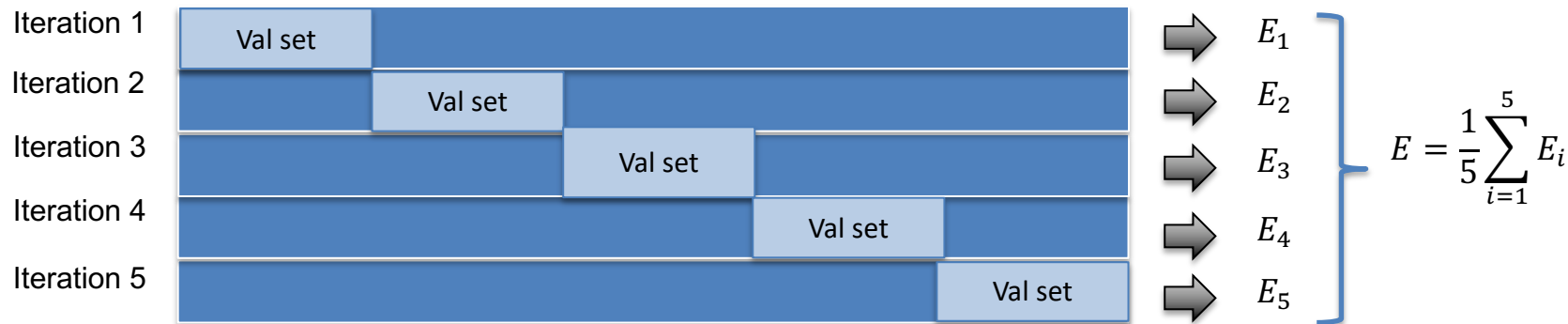Rather than using a fixed validation set, we train and run the model(s) multiple times, each time using a different subset ("fold") as the validation set



**K=5 Folds**

Training + cross-validation
(~80-90%)

Test
(~10-20%)

Iteration 1 — Val set
Iteration 2 — Val set
Iteration 3 — Val set
Iteration 4 — Val set
Iteration 5 — Val set

Use for model comparison/selection

Use to evaluate final model

# K-Folds Cross Validation

We calculate the error on the validation fold for each iteration, and then average them together to get the average error

**K=5 Folds**



$$E = \frac{1}{5} \sum_{i=1}^{5} E_i$$

# Benefits of Cross Validation

- Maximizes the data available for training the model – important for small datasets

- Provides a better evaluation of how well the model can generalize to new data – validation performance is not biased by choice of datapoints to use for validation

**Wrap-up**

# Wrap Up

- Modeling process is just one piece of the CRISP-DM process

- Model complexity comes from features, algorithm and hyperparameters

- Underfitting and overfitting are common modeling issues

- Test sets and validation sets ensure we properly select and evaluate models