# Picking Intro

School of Information Studies
Syracuse University

# Patterns in Games and Prices

- Obtain

- Scrub

- Explore

- Model

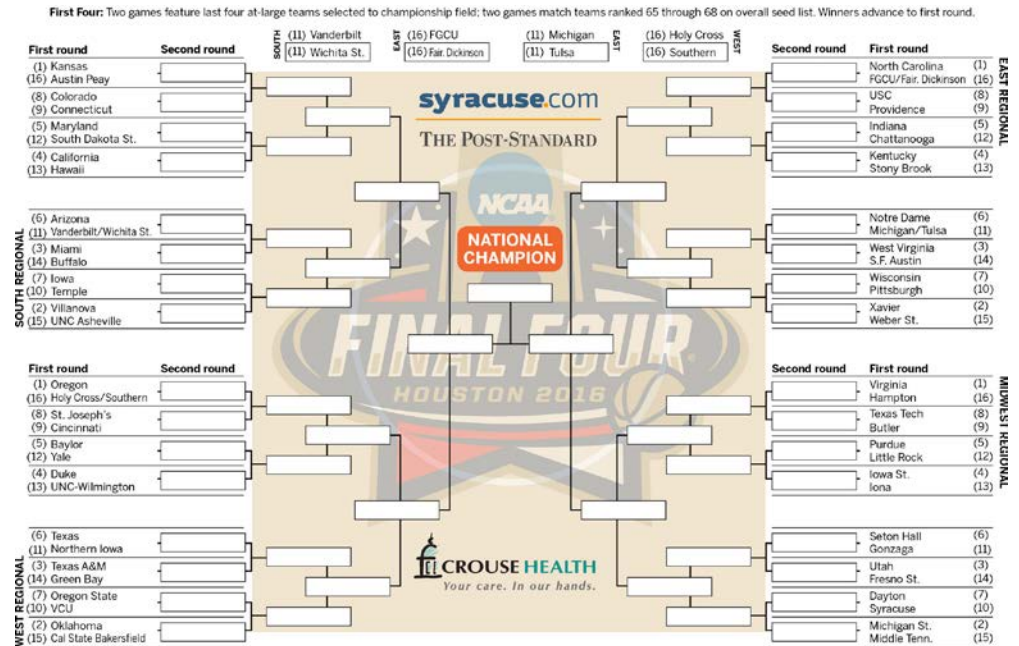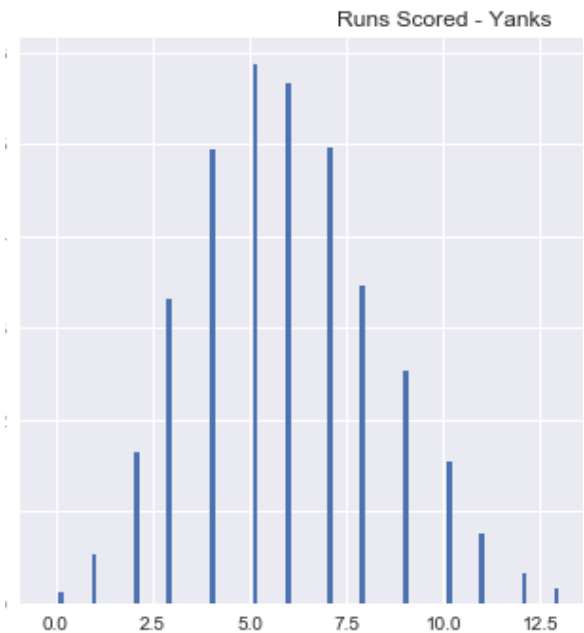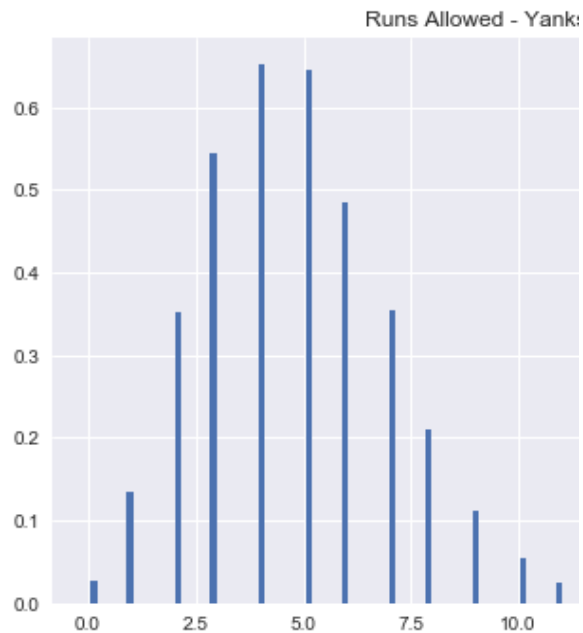- iNterpret

School of Information Studies
Syracuse University

# Our Challenge This Week?

School of Information Studies
Syracuse University

# Using Distributions to Pick a Winner

School of Information Studies
Syracuse University

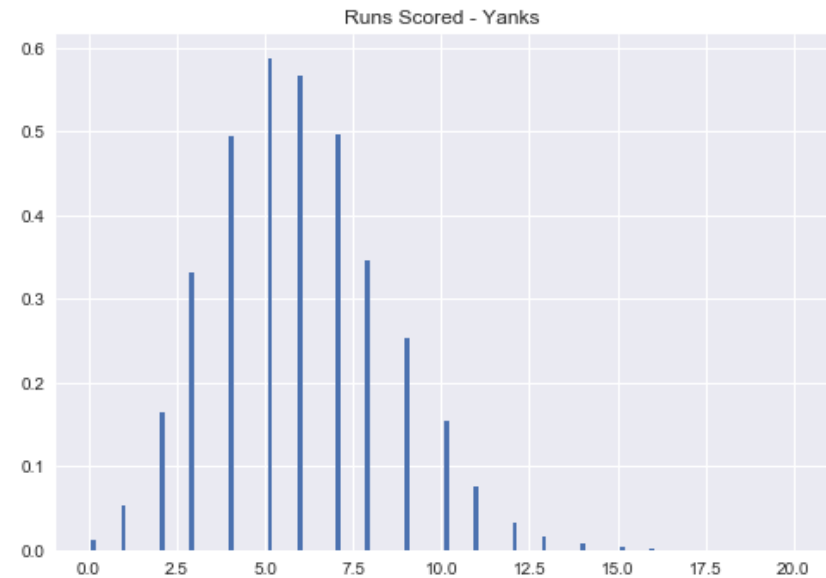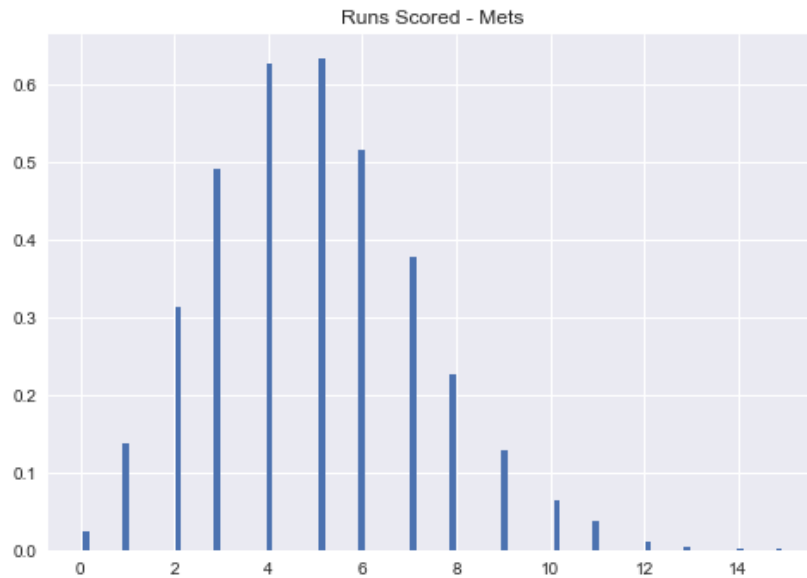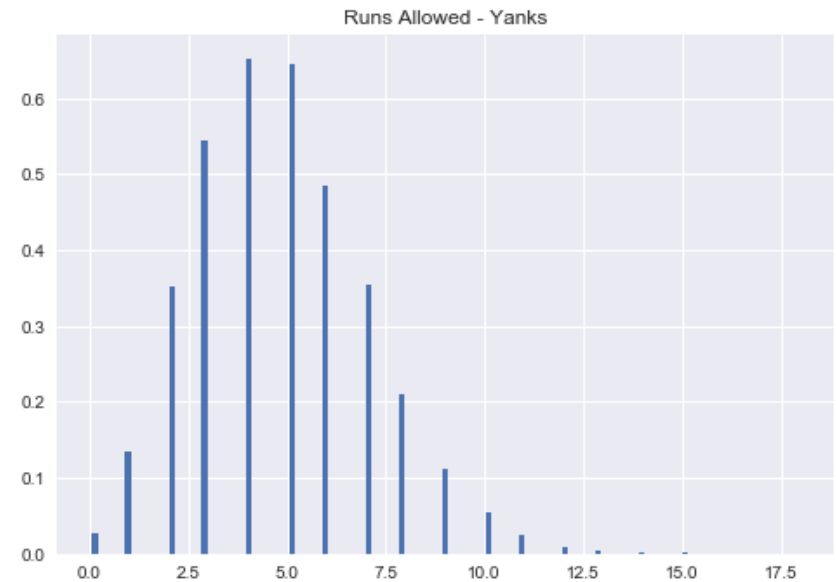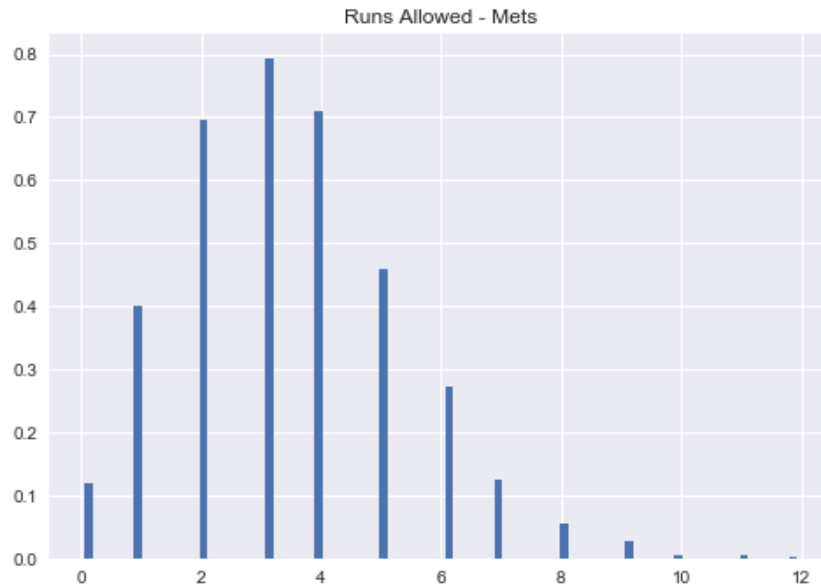# But How?

School of Information Studies
Syracuse University

Data Review

School of Information Studies
Syracuse University

# Runs Scored



Runs Scored - Mets

Runs Scored - Yanks

School of Information Studies
Syracuse University

# Runs Allowed

School of Information Studies
Syracuse University

# Scoring Distributions



Runs Allowed - Mets



Runs Allowed - Yanks

School of Information Studies
Syracuse University

# Picking Winners

School of Information Studies
Syracuse University

# Poisson Distribution

School of Information Studies
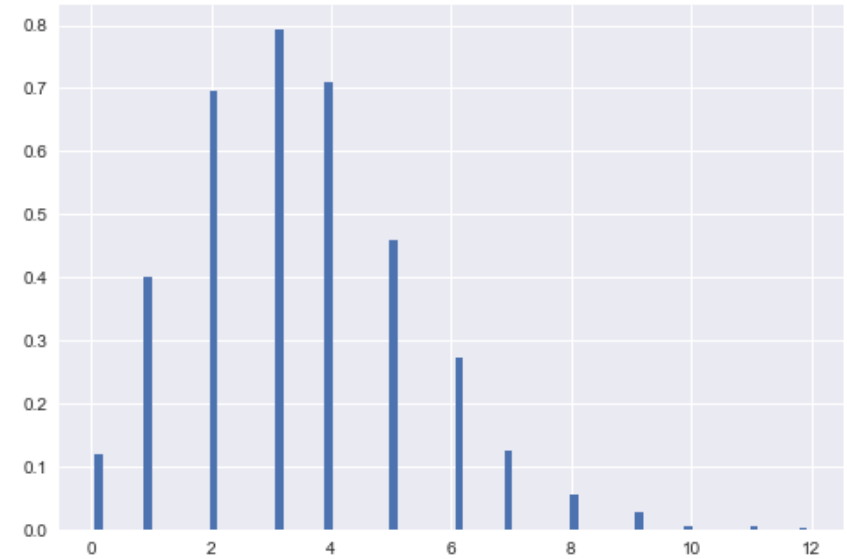Syracuse University

# Simulating Runs Scored

School of Information Studies
Syracuse University

# Sports Simulation



Simulation Results

| Runs Scored by Away Team | | Runs Scored by Home Team |
|---|---|---|
| 6 | **11** → 7 | 7 |
| 4 | | 7 |
| 12 | | 10 |
| 7 | **4** **8** | 3 |
| 8 | | 2 |
| 1 | | 7 |
| 4 | **1** **6** | 11 |
| 9 | | 8 |
| 3 | | 4 |
| 1 | • • • | 4 |
| 4 | | 1 |
| 4 | | 3 |
| 8 | | 9 |
| 11 | | 4 |
| 7 | **2** **1** | 1 |
| 2 | | 8 |
| 3 | | 7 |

Source: Adapted from Miller (2005).

School of Information Studies
Syracuse University

# Poisson Distribution

- Good approximation for count responses

$$P(Y = y) = \frac{e^{-\mu}\mu^y}{y!}$$



Runs Scored - Yanks

- Occurrence of events during certain time interval
- Arrival rate problems

School of Information Studies
Syracuse University

# Poisson Distribution (cont.)

MetAwayScore =
np.random.poisson(4.97, 10000)

MetAwayDefend =
np.random.poisson(3.45, 10000)

YankHomeScore =
np.random.poisson(5.97, 10000)

YankHomeDefend =
np.random.poisson(4.84, 10000)

plt.hist(MetAwayScore, bins='auto', rwidth = .5, normed=True)

plt.title("Runs Scored – Mets")
plt.show()


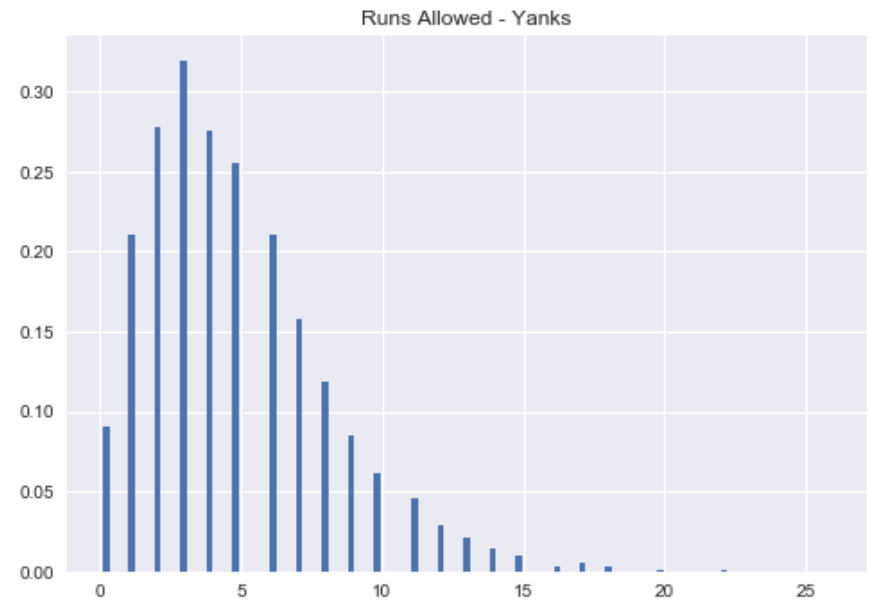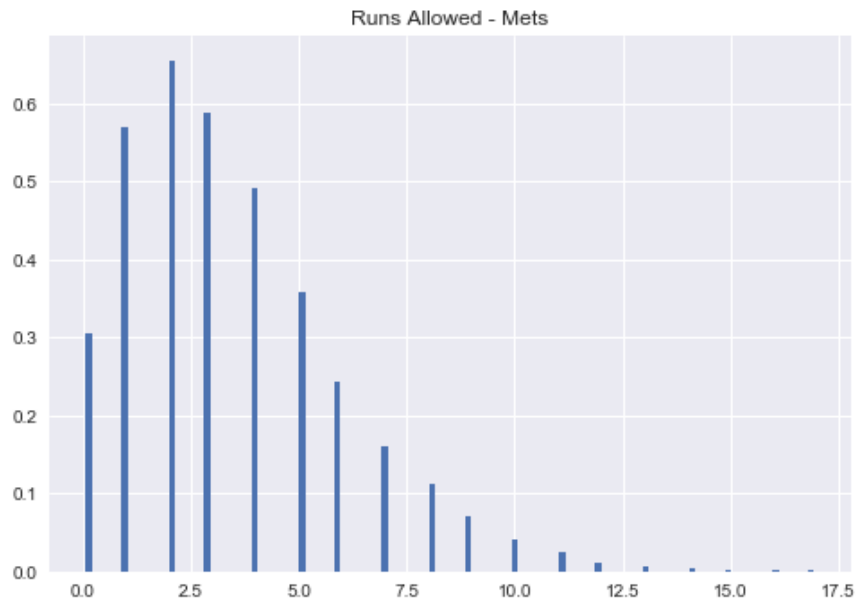
Runs Scored - Mets

School of Information Studies
Syracuse University

# Negative Binomial Distribution

School of Information Studies
Syracuse University

# Simulating Runs Allowed

School of Information Studies
Syracuse University

# Sports Simulation

## Simulation Results

| Runs Scored by Away Team | Runs Allowed by Home Team | | | Runs Scored by Home Team | Runs Allowed by Away Team |

7 → 5 → 6

2 → 6 → 4

3 → 1 → 2

4 ← 6 ← 5

5 ← 7 ← 6

2 ← 0 ← 1

Source: Adapted from Miller (2005).

School of Information Studies
Syracuse University

# Negative Binomial

- Alternative approximation for count responses

$$P(Z = z)$$
$$= \binom{z - 1}{k - 1} p^k (1 - p)^{z-k}$$

- Generalization of Poisson distribution
- Rare event problems

Runs Allowed - Yanks
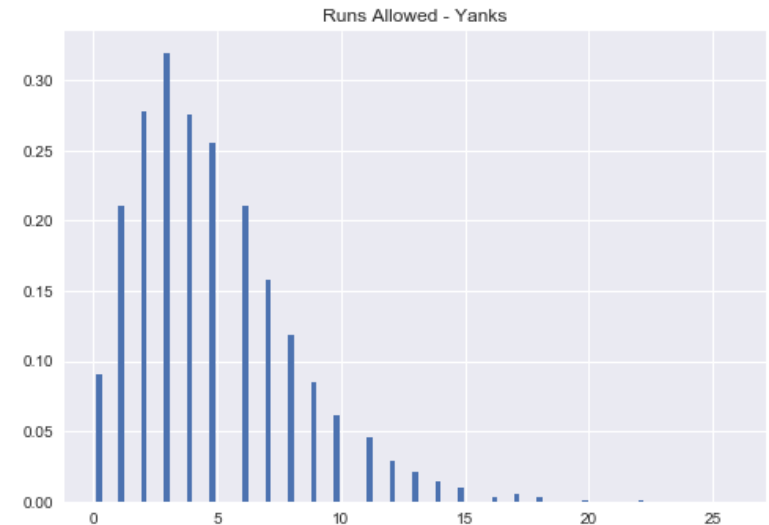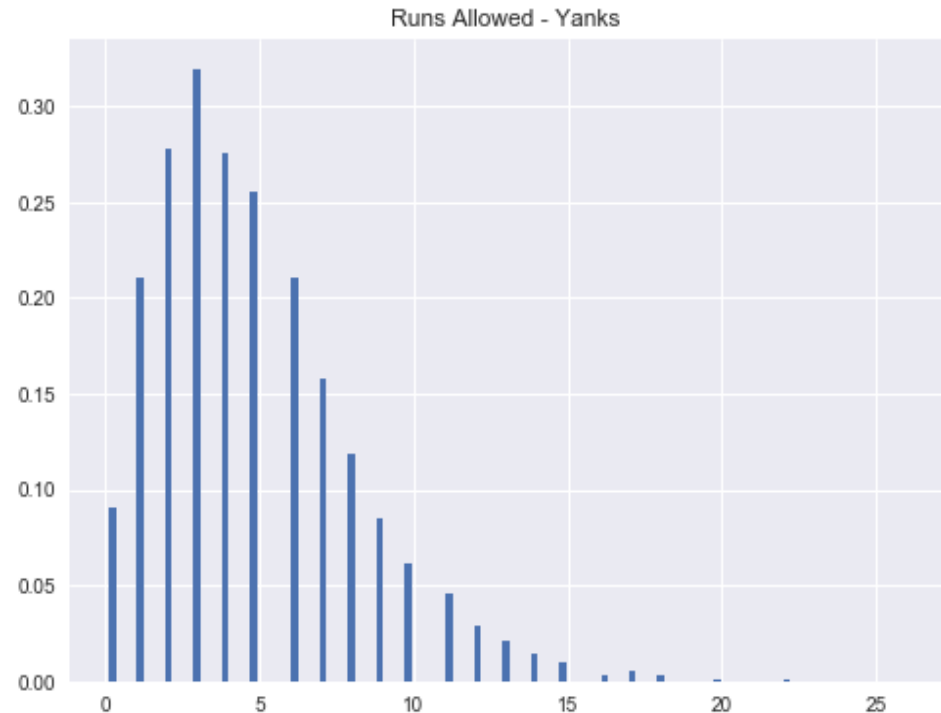
# Negative Binomial

- Alternative approximation for count responses

$$P(Z = z) = \binom{z - 1}{k - 1} p^k (1 - p)^{z-k}$$

- Generalization of Poisson distribution

- Rare event problems


Runs Allowed - Yanks

School of Information Studies
Syracuse University

# Negative Binomial (cont.)
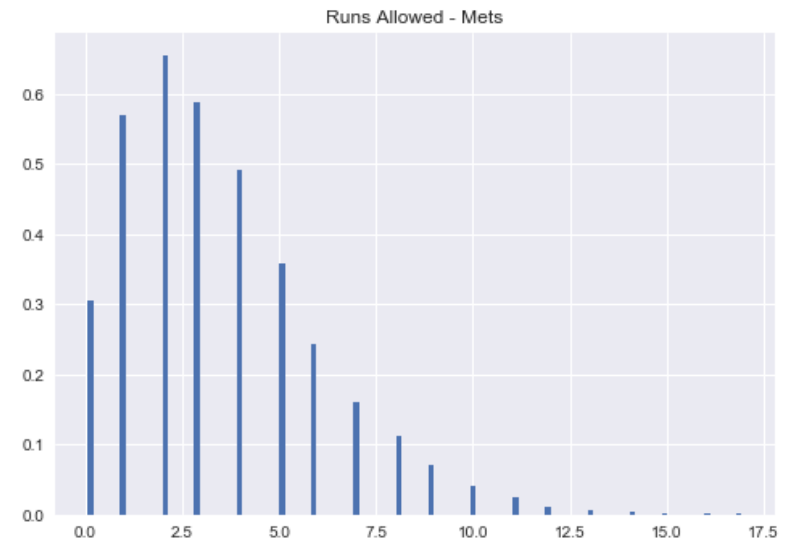
MetAwayScore =
np.random.negative_binomial(4, mas, 10000)

MetAwayDefend =
np.random.negative_binomial(4, mad, 10000)

YankHomeScore =
np.random.negative_binomial(4, yhs, 10000)

YankHomeDefend =
np.random.negative_binomial(4, yhd, 10000)

plt.hist(MetAwayScore, bins='auto', rwidth = .5,
normed=True)
plt.title("Runs Scored – Mets")
plt.show()



Runs Allowed - Mets

School of Information Studies
Syracuse University

# Applications

School of Information Studies
Syracuse University

# Additional Applications

School of Information Studies
Syracuse University
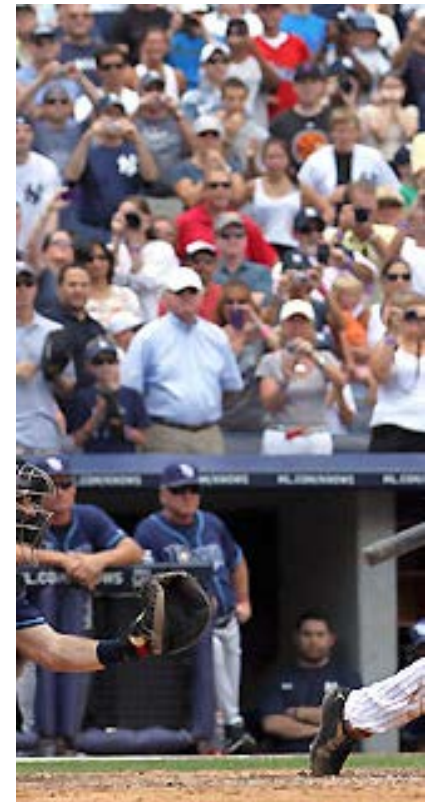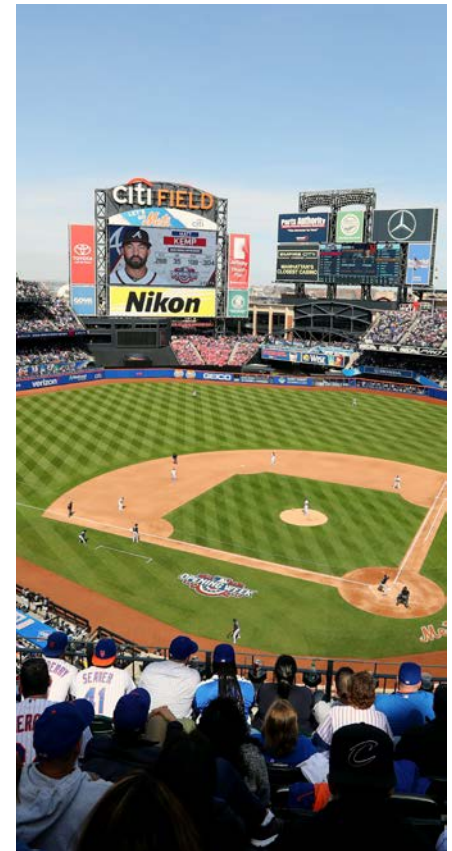
# Baseball Prospectus

- Predicting performance before the season
  - Nate Silver
  - PECOTA
- Variations
  - Military war games
  - Film releases
  - Associate performance

School of Information Studies
Syracuse University
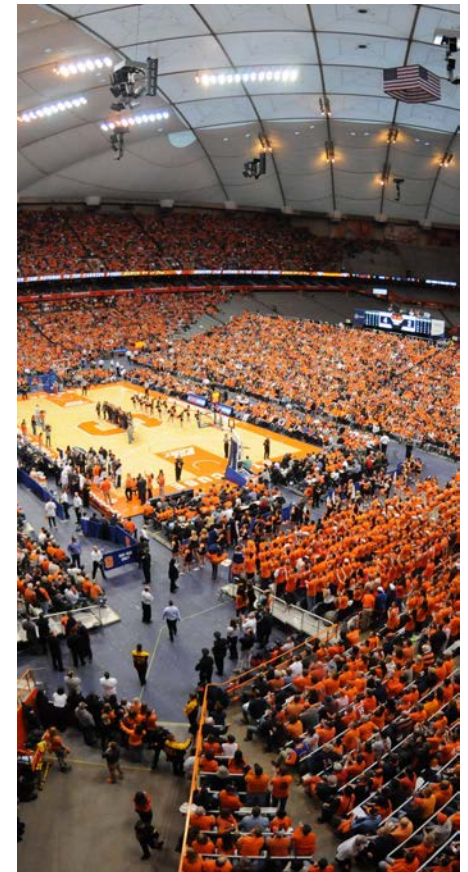
# *Moneyball* Problem

- Calculating the value of a player
  - Billy Beane
  - Individual summary stats applied to team performance
- Variations
  - Associate performance
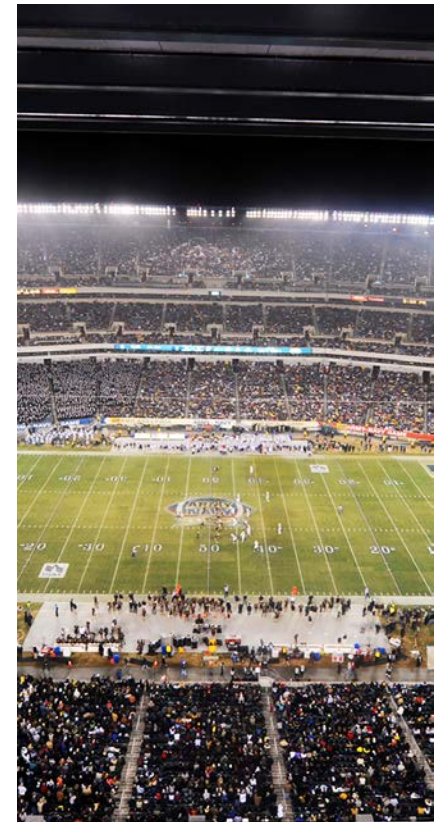  - Client conversion
  - Customer lifetime value

School of Information Studies
Syracuse University

# Coaching Problems

- Utilizing player performance
  - Microanalysis of game
  - Sabermetrics
- Variations
  - Football
  - Basketball

School of Information Studies
Syracuse University

# Bowl Championship Series

- Predicting team performance against unknown opponent
  - Strength of schedule
  - Ensemble approaches
- Variations
  - March Madness
  - Product deployment
  - Recommendation engine

School of Information Studies
Syracuse University

# Billy Waters Problem

- Predicting the winning team in the next game?
  - Human expertise
  - Simulation
- Variations
  - March Madness
  - Film release
  - Product deployment

School of Information Studies
Syracuse University

# Picking II Intro

School of Information Studies
Syracuse University

# Patterns in Games and Prices

- Obtain

- Scrub

- Explore

- Model

- iNterpret

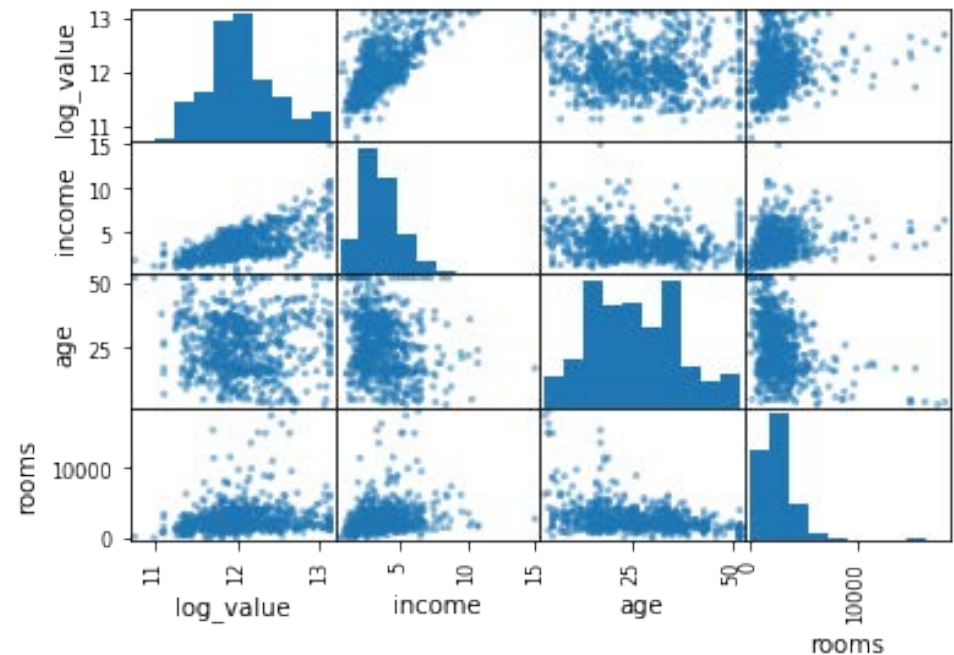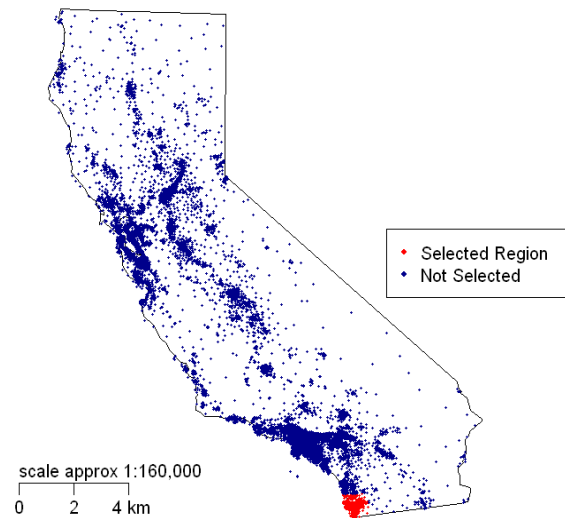School of Information Studies
Syracuse University

# Our Challenge Now?

School of Information Studies
Syracuse University

# Modeling Housing Prices

School of Information Studies
Syracuse University
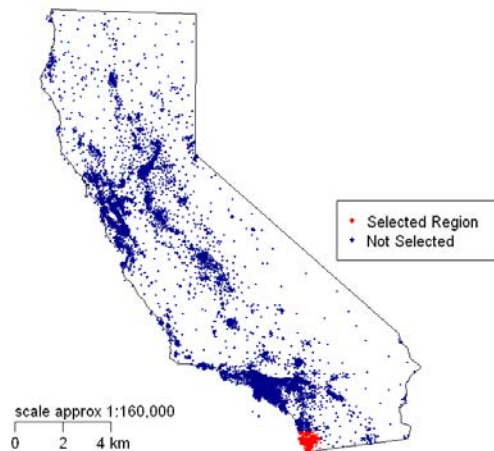
# But How?

School of Information Studies
Syracuse University

Data Review II

School of Information Studies
Syracuse University

# Housing Data

School of Information Studies
Syracuse University

# Feature Correlation

School of Information Studies
Syracuse University

# Feature Correlation (cont.)

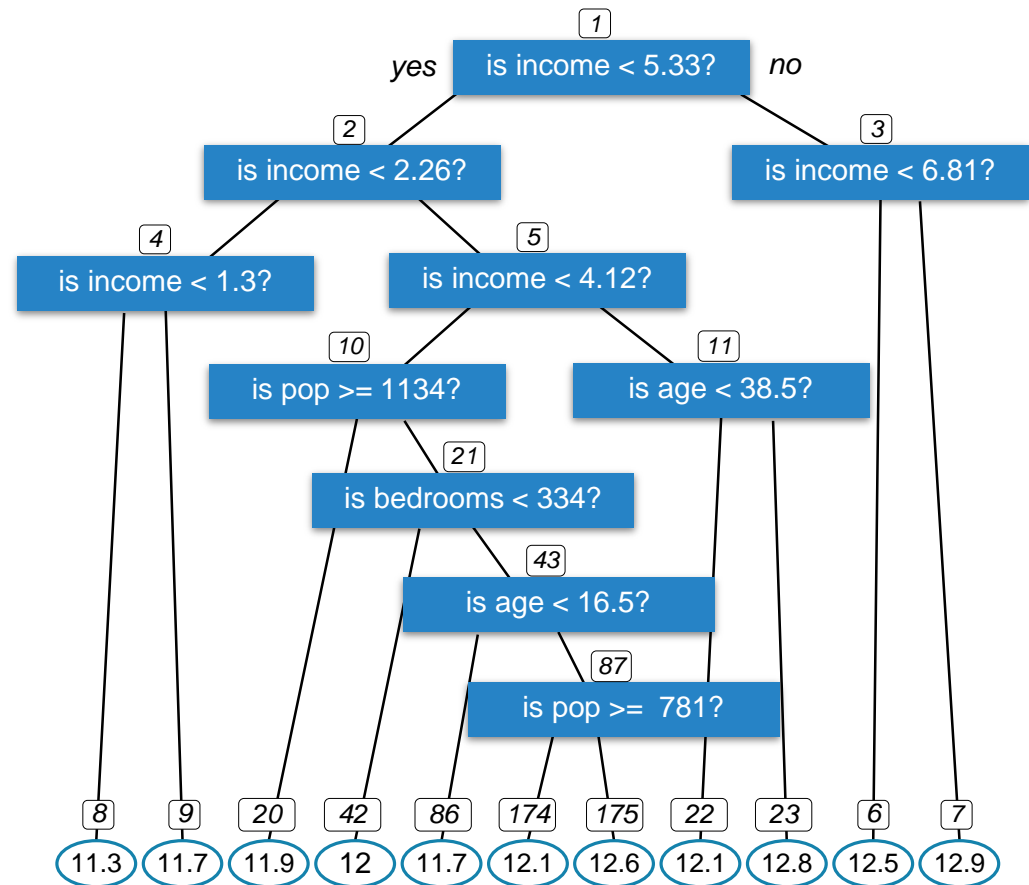School of Information Studies
Syracuse University

# Picking Values

# Trees Forests

School of Information Studies
Syracuse University

# Picking a Tree in the Forest
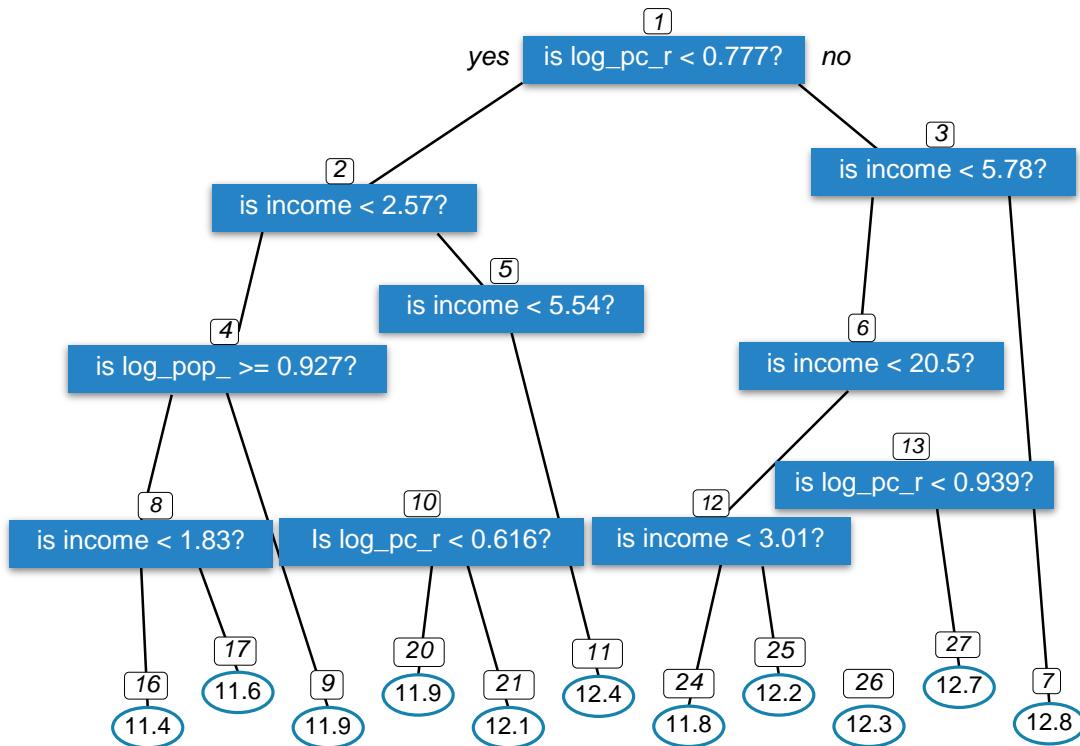
School of Information Studies
Syracuse University
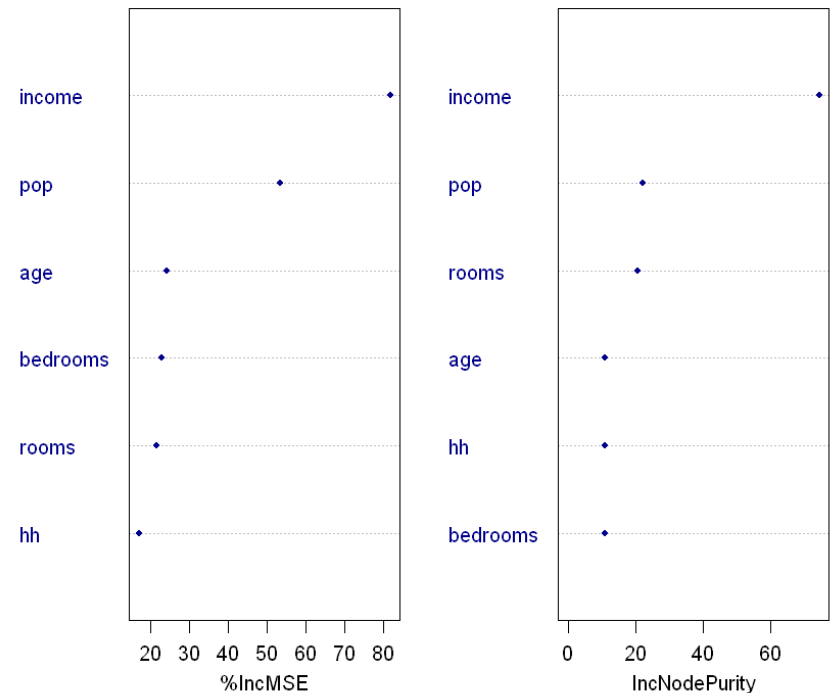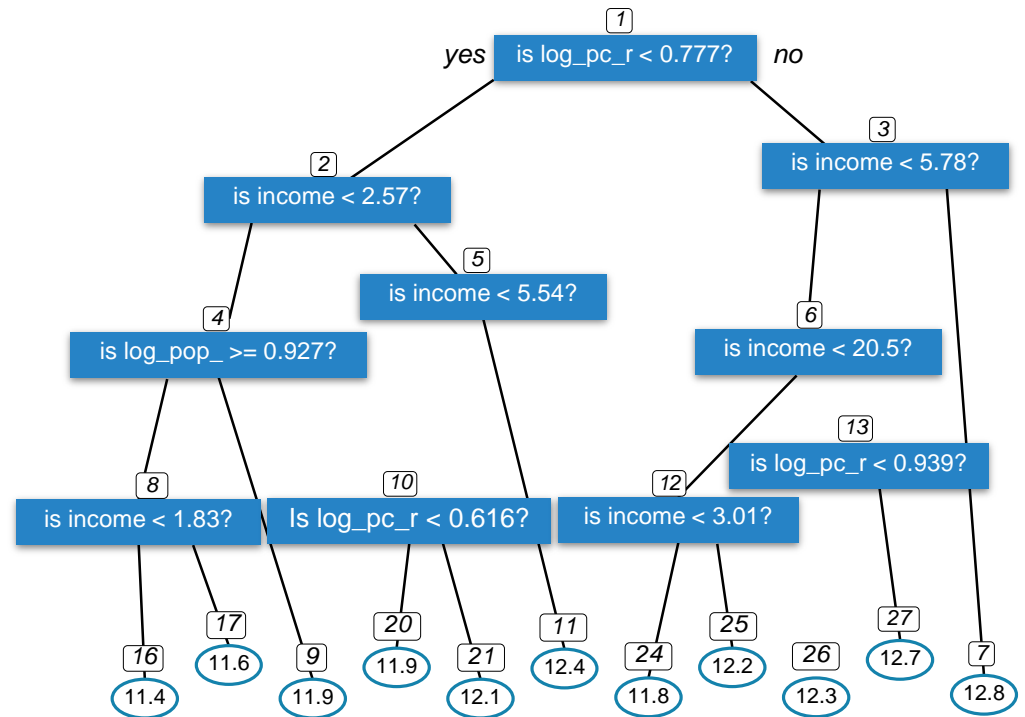
# Decision Trees

- Key advantage is interpretability

- Partition the space into simple regions to achieve best fit

- Pruning methods control the size of the tree

**School of Information Studies**
Syracuse University

# Random Forests

- Ensemble method using multiple decision trees

- Recursive partitioning on the training set

- Effective with large number of explanatory variables

# Random Forests (cont.)

- Provides interpretability through use of one tree from set

- Significant difference in performance between train and test indicates overfitting

- Individual explanatory variables can still be inferred

School of Information Studies
Syracuse University