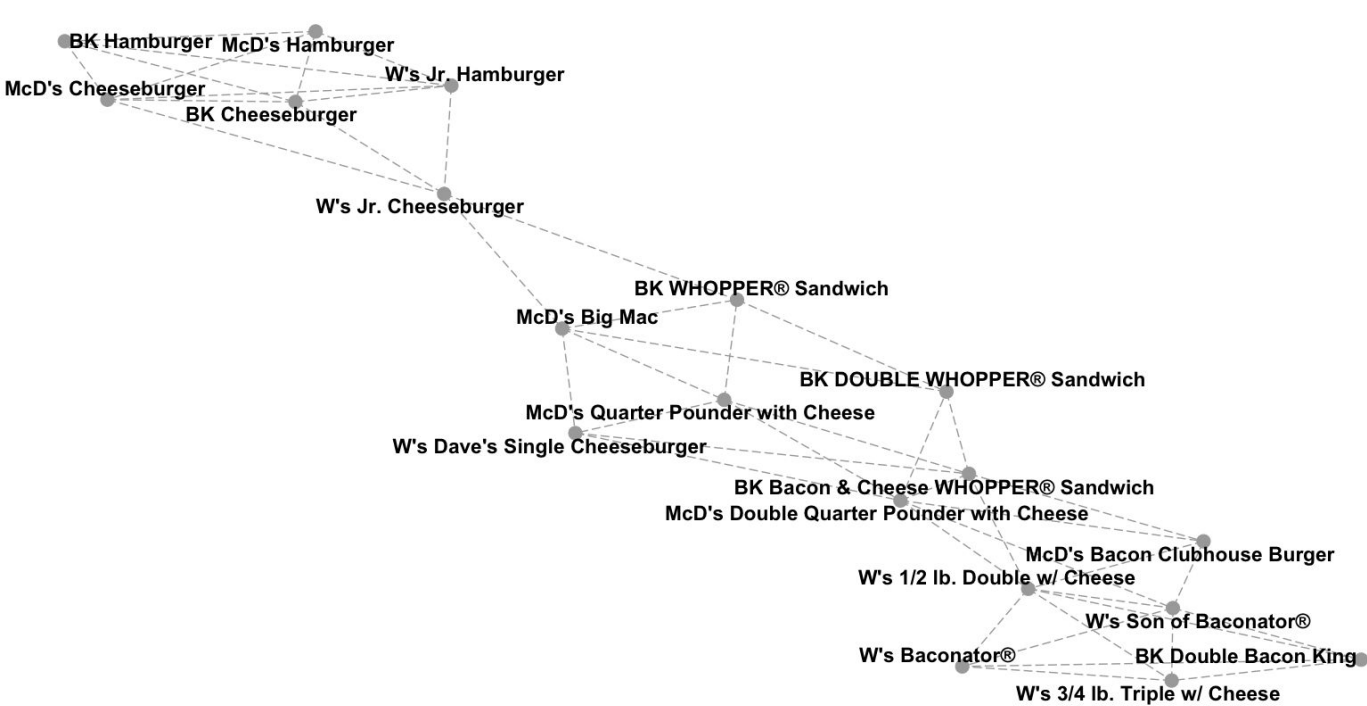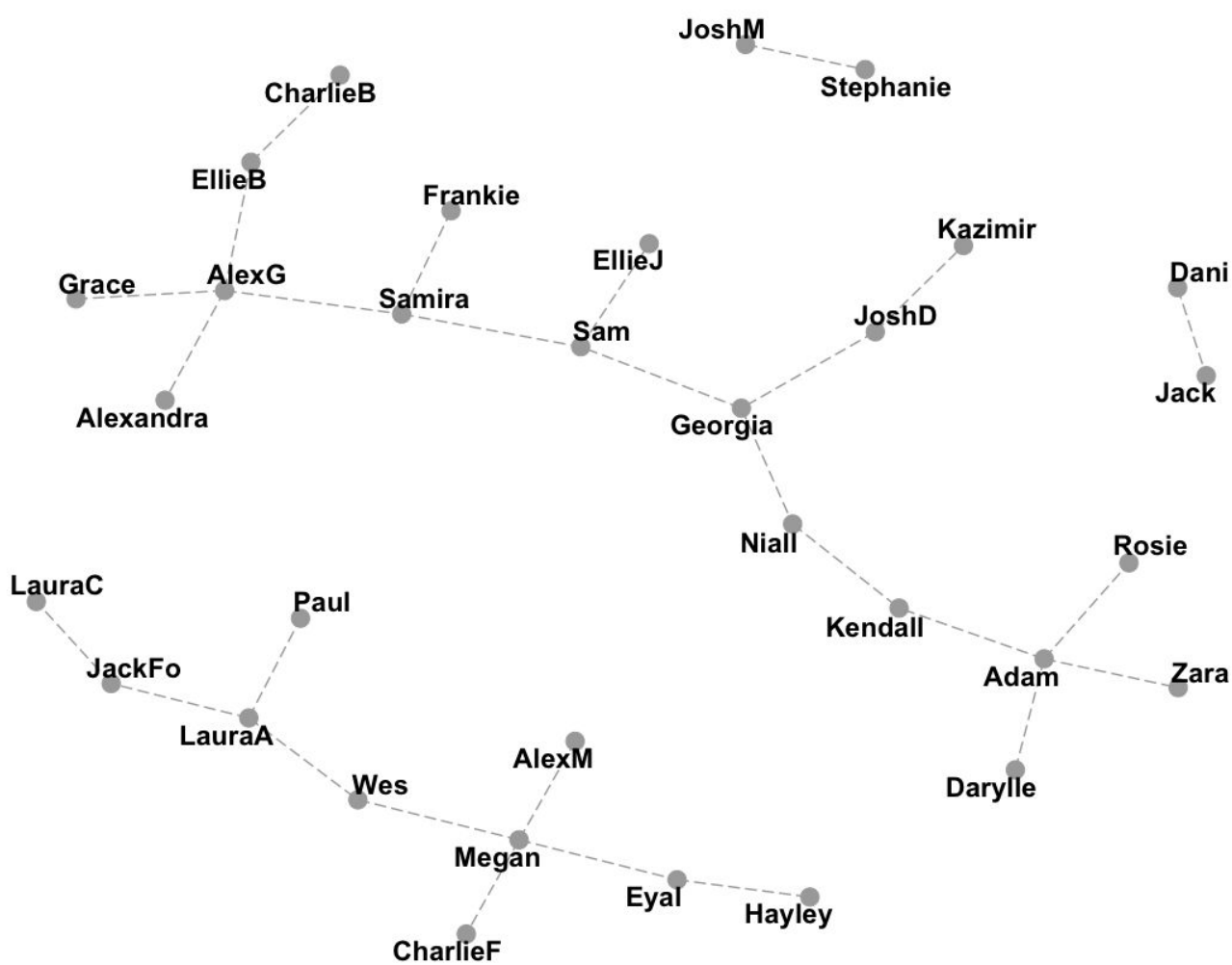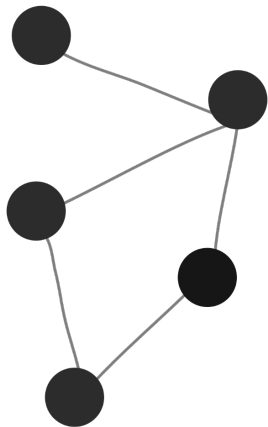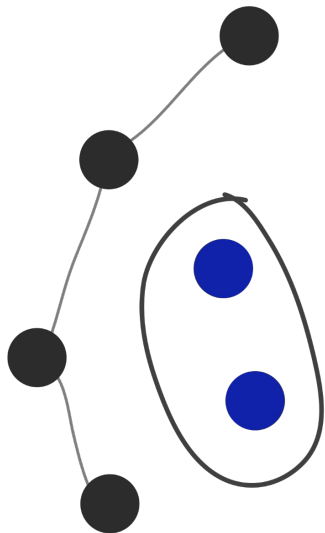# Topics In Data Science

Chelsea Parlett-Pelleriti

# Node Based Resilience Clustering

BK Hamburger
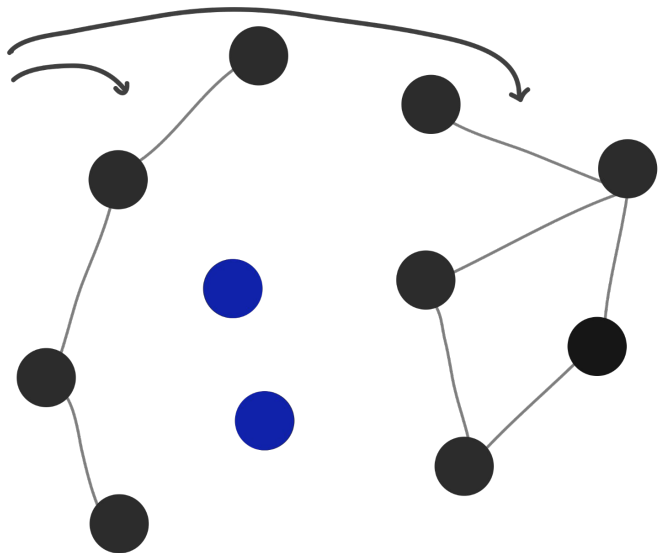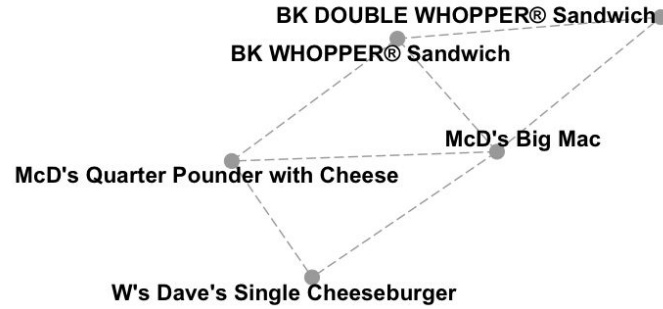
McD's Hamburger

W's Jr. Hamburger

McD's Cheeseburger

BK Cheeseburger

W's Jr. Cheeseburger

BK WHOPPER® Sandwich

McD's Big Mac

BK DOUBLE WHOPPER® Sandwich

McD's Quarter Pounder with Cheese

W's Dave's Single Cheeseburger

BK Bacon & Cheese WHOPPER® Sandwich

McD's Double Quarter Pounder with Cheese

McD's Bacon Clubhouse Burger

W's 1/2 lb. Double w/ Cheese

W's Son of Baconator®

W's Baconator®

BK Double Bacon King

W's 3/4 lb. Triple w/ Cheese
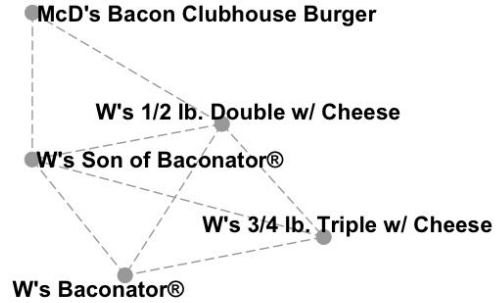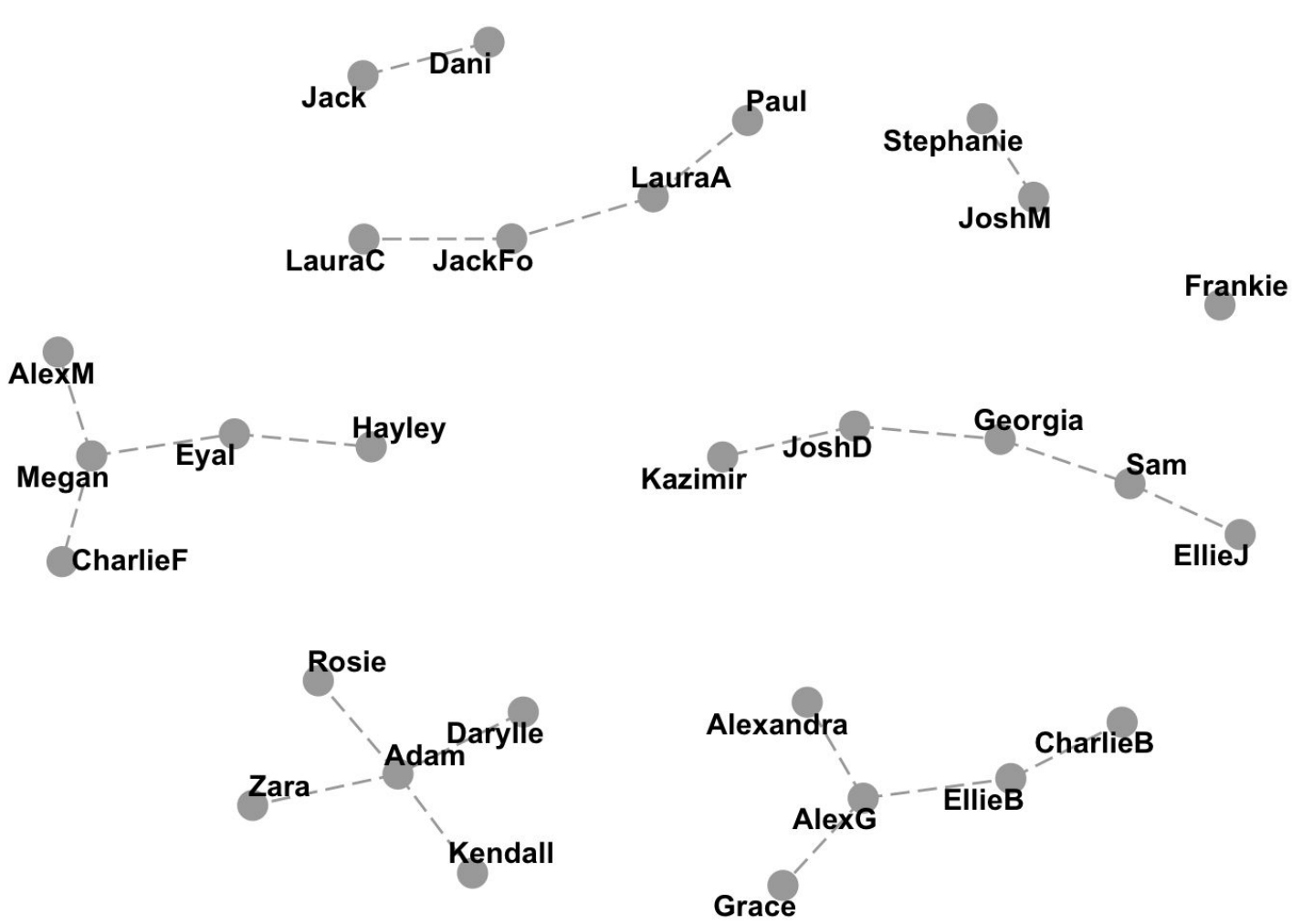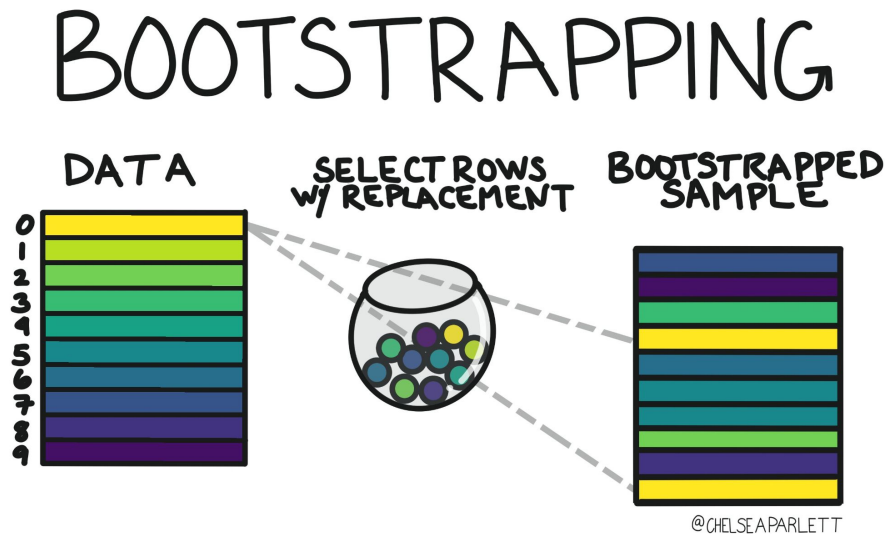
S

V-S

# Ensemble Methods
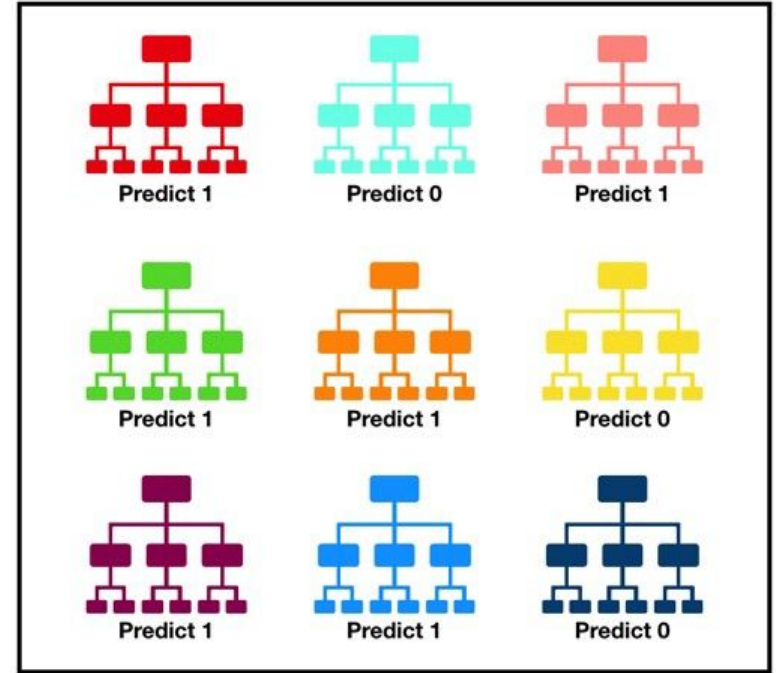
ONE vs. MANY

# Random Forest

# Random Forest

- **Bagging (Bootstrap Aggregating)**: Instead of using all of our training data to train each model in our sample we use **bootstrapping** to choose the samples (rows) we will include.
  - **Bootstrapping** is when you randomly sample data points *with replacement*, meaning that a data point can be included in your bootstrapped sample *more* than once, OR not at all.



BOOTSTRAPPING

DATA    SELECT ROWS w/ REPLACEMENT    BOOTSTRAPPED SAMPLE

0 1 2 3 4 5 6 7 8 9

@CHELSEAPARLETT

# Random Forest

- **Random Feature Selection**: Instead of using all the available *features/predictors* in our dataset for every model, for each model *we randomly choose a different subset of features to use* when training.
  - This helps our ensemble generalize, because it doesn't become overly reliant on one feature (since that feature might not appear in every model).

# Bayesian Statistics

# Bayesian Statistics

$$\underbrace{P(\theta \mid D)}_{posterior} \propto \underbrace{P(\theta)}_{prior} \underbrace{P(D \mid \theta)}_{likelihood}$$

Data + Expertise =

Inference

# Bayesian Statistics

$$\frac{\text{Data}}{\text{Inference}} =$$

# Bayesian Statistics

MINIMIZE:

$$\sum \left( x_i - \hat{x}_i \right)^2 + \lambda \sum \beta_j^2$$

how off we were

true value

model's guess

how HARSHLY we penalize

how big the coefs are

MINIMIZE:

$$\sum \left( x_i - \hat{x}_i \right)^2 + \lambda \sum |\beta_j|$$

how off we were

true value

model's guess

how HARSHLY we penalize

how big the coefs are

Bayesian Statistics

**What we believed before** + **Evidence from the data** = **New beliefs**