# Data Science Tutorial

Eliezer Kanal – *Technical Manager, CERT*

Daniel DeCapria – *Data Scientist, ETC*

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA  15213

**Software Engineering Institute** | **Carnegie Mellon University**

# About us



### Eliezer Kanal

*Technical Manager, CERT*

## Recent projects:

- ML-based Malware Classifier

- Network traffic analysis

- Cybersecurity questionnaire optimization



### Daniel DeCapria

*Data Scientist, ETC*

## Recent projects:

- Cyber risk situational dashboard

- Big Learning benchmarks

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

2

# Today's presentation – a tale of two roles

The call center manager

*Introduction to
data science capabilities*

The master carpenter

*Overview of the
data science toolkit*

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

3

# Call center manager

## First day on job… welcome!

Goal:   Reduce costs

Task:   Keep calls short!

Data:

    Average call time:        5.14 minutes (5:08)… very long!

    Number of employees:    300

    Average calls per day:    ~28,000

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**4**

# Call center manager – *Gather data*

Get the data!

- Where is it?

- What will you use to analyze it?

- How accurate it is?

- How complete is it?

- Is it too big to easily read?

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**5**

# Data cleaning = 90% of the work

2 weeks (10 days) = 9 cleaning, 1 analyzing

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**6**

# Cleaning the Data – *Structuring the Data*

**Goal**: Organize data in a table, where…

Columns = descriptor (age, weight, height)

Row = individual, complete records



How can you get data out of these documents?



Less structure

More structure

**Software Engineering Institute** | **Carnegie Mellon University**

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

7

# Cleaning the Data

Even when you think your data should be clean, it might not be…

**Please tell us how many years of experience you have had working in the following domains. Enter a whole number.**

**Machine Learning**

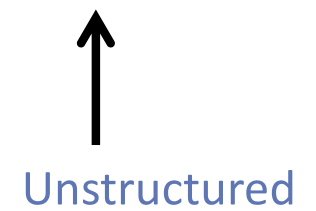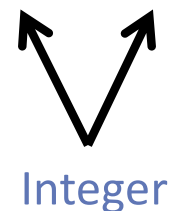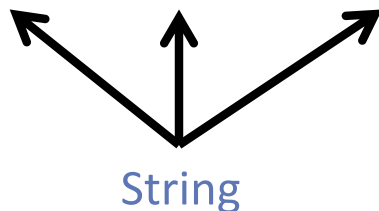| 0.5 | 2 | 1 | 0 | 1/2 | none | 0 semesters | 6 months |

**Computer Science**

| 1.5 | this semester | 3 | 2 | 1 | 0 | 6 | 5 | 4 | 8 | 11 | second |
| .5 | 6 months |

**Mathematics**

| 0.5 | .333 | 22 | 3 | some background in calculus | 2 | 1 | 0 | 6 |
| 5 | 4 | 8 | 10+ | 16 | fourth | 10 | 11 years | 7 semesters | 3.5 |

**Software Engineering Institute** | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**8**

# Cleaning the Data – *Call Center Example*

| Name | Mgr | Dir | Call Length | Phone Line | Problem solved? | Comment |
|------|-----|-----|-------------|------------|-----------------|---------|
| Beth Jones | Dan Thomas | Anne Kim | 1:30 | 1 | Y | (5) … |
| Beth Jones | Dan Thomas | Anne Kim | 1:52 | 3 | Y | … |
| (1) Jones, Beth | Dan Thomas | Anne Kim | (1) 90 | 2 | Y | … |
| Tom Keane | Mark Ryan | Tim Pike | 88 | 2 | (1) N | … |
| Tom Keane | (2) Mark Ryan | Tim Pike | 144 | 3 | No | … |
| Tom Keane | Kevin Wood | Tim Pike | 200 | (4) | Yes | … |
| Tom Keane | Kevin Wood | Tim Pike | 94511 | 2 | No | … |
| (6) Tom Keane | Kevin Wood | Tim Pike | (3) 421 | 2 | Yes | … |

String       Integer       "Nominal"       Unstructured

Software Engineering Institute | Carnegie Mellon University

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

2017 SEI Data Science in Cybersecurity Symposium
Approved for Public Release; Distribution is Unlimited

9

| | 1 CRIM | 2 ZN | 3 INDUS | 4 CHAS | 5 NOX | 6 RM | 7 AGE | 8 DIS | 9 RAD | 10 TAX | 11 PTRATIO | 12 B | 13 LSTAT | 14 MEDV |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.0063 | 18 | 2.3100 | 0 | 0.5380 | 6.5750 | 65.2000 | 4.0900 | 1 | 296 | 15.3000 | 396.9000 | 4.9800 | 24 |
| 2 | 0.0273 | 0 | 7.0700 | 0 | 0.4690 | 6.4210 | 78.9000 | 4.9671 | 2 | 242 | 17.8000 | 396.9000 | 9.1400 | 21.6000 |
| 3 | 0.0273 | 0 | 7.0700 | 0 | 0.4690 | 7.1850 | 61.1000 | 4.9671 | 2 | 242 | 17.8000 | 392.8300 | 4.0300 | 34.7000 |
| 4 | 0.0324 | 0 | 2.1800 | 0 | 0.4580 | 6.9980 | 45.8000 | 6.0622 | 3 | 222 | 18.7000 | 394.6300 | 2.9400 | 33.4000 |
| 5 | 0.0691 | 0 | 2.1800 | 0 | 0.4580 | 7.1470 | 54.2000 | 6.0622 | 3 | 222 | 18.7000 | 396.9000 | 5.3300 | 36.2000 |
| 6 | 0.0299 | 0 | 2.1800 | 0 | 0.4580 | 6.4300 | 58.7000 | 6.0622 | 3 | 222 | 18.7000 | 394.1200 | 5.2100 | 28.7000 |
| 7 | 0.0883 | 12.5000 | 7.8700 | 0 | 0.5240 | 6.0120 | 66.6000 | 5.5605 | 5 | 311 | 15.2000 | 395.6000 | 12.4300 | 22.5000 |
| 8 | 0.1446 | 12.5000 | 7.8700 | 0 | 0.5240 | 6.1720 | 96.1000 | 5.9505 | 5 | 311 | 15.2000 | 396.9000 | 19.1500 | 27.1000 |
| 9 | 0.2112 | 12.5000 | 7.8700 | 0 | 0.5240 | 5.6310 | 100 | 6.0821 | 5 | 311 | 15.2000 | 386.6300 | 29.9300 | 16.5000 |
| 10 | 0.1700 | 12.5000 | 7.8700 | 0 | 0.5240 | 6.0040 | 85.9000 | 6.5921 | 5 | 311 | 15.2000 | 386.7100 | 17.1000 | 18.9000 |
| 11 | 0.2249 | 12.5000 | 7.8700 | 0 | 0.5240 | 6.3770 | 94.3000 | 6.3467 | 5 | 311 | 15.2000 | 392.5200 | 20.4500 | 15 |
| 12 | 0.1175 | 12.5000 | 7.8700 | 0 | 0.5240 | 6.0090 | 82.9000 | 6.2267 | 5 | 311 | 15.2000 | 396.9000 | 13.2700 | 18.9000 |
| 13 | 0.0938 | 12.5000 | 7.8700 | 0 | 0.5240 | 5.8890 | 39 | 5.4509 | 5 | 311 | 15.2000 | 390.5000 | 15.7100 | 21.7000 |

# Exploratory Data Analysis (EDA)

- Mean

- Median

- Standard deviation

- Histograms!

# Distributions

- The majority of data will follow SOME distribution

  - Weight of all Americans: *Gaussian*

    

  - phone call length: *Exponential*

    

- Determining distribution is a common Data Science task

- Multidimensional outliers: Insider Threat example



Image Copyright 2001-2016 The Apache Software Foundation. See Copyright slide for more details.

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

12

# EDA – *Smart visualizations*

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

Software Engineering Institute | Carnegie Mellon University

14

TRANSPORT SYSTEM

RIVER

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

15

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

16

# Brief interruption



**Lovely visuals and all, but THIS isn't data science! Where's the fancy predicting the future and whatnot?**

**Skeptics in the audience**

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**17**

# Brief interruption

Data Science helps you
use data to get results.
*This is it.*

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**18**

# Call center manager – *call duration histogram*

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**19**

# Call Center manager – *Insights!*

Strategy update:

- Goodbye "reduce call time"

- Hello "reduce callbacks"

How to measure?

"callbacks" isn't currently captured

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**20**

# Feature Engineering

*Need more useful data?*

**Create it yourself!**



"When you put it like that, it makes complete sense."

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**21**

# Feature Engineering

- Feature Engineering: coming up with new, useful (i.e., informative) data

  o mean, sums, medians, etc.

  o $x^2$, xy, sqrt(xy), etc.

- Our case:

  o # of callbacks

  o Call during peak time?

  o Overall agent performance? (combination of factors)

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**22**

Software Engineering Institute | Carnegie Mellon University

# The role of Listening in Data Science

Data science finds hidden patterns in data

Experts know what data & patterns are important

## Talk to subject matter experts



Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

23

# Call Center manager – *Predictive analytics*

Can we predict staffing levels…

- …one day ahead?

- …one week ahead?

- …one month ahead?

Can we determine what types of calls to expect…

- …for a product we haven't had before?

- …for a market we've never seen before?

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**24**

# Example Predictive Analytics Questions

## Predicting Current Unknowns

Online:  Which ads are malicious?

Security:  Is the bank transaction fraudulent?

IC:  Which names map to the same person (entity resolution)?

## Predicting Future Events

Retail:  What will be the new trend of merchandise that a company should stock?

Security:  Where will a hacker next attack our network?

IC:  Who will become the next insider threat?

## Determining Future Actions

Sales:  How can a company increase sales revenues?

Health:  What actions can be taken to prevent the spread of flu?

IC:  How will a vulnerability patch affect our knowledge/preparedness for future attacks?

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**25**

# Call Center manager – *Predictive analytics*

Many techniques available, explored in next section

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**Software Engineering Institute** | **Carnegie Mellon University**

**26**

# Call Center manager – *Review*

Get data



Clean data



EDA, Visualization



Interpretation



Action!



Prediction

**Because we know our data,** we can ask…

- …more intelligent questions

- …action-oriented questions

- …questions that can be answered

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**28**

*This slide intentionally left blank*

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**Software Engineering Institute** | **Carnegie Mellon University**

**29**

# The master carpenter



"The right tool for the job"

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**30**

# Feature Engineering – *Part 2*

*"With the wrong wood, I can make nothing"*

The fuel of data science is data

Data preparation is critical

Data quality $\gg$ algorithm choice

That will come up…

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**31**

# Types of Machine Learning Algorithms

Classification

- Naïve Bayes

- Logistic Regression

- Decision Trees

- K-Nearest Neighbors

- Support Vector Machines

Regression

- Linear Regression

- Support Vector Machines

Clustering

- K-Means Clustering

# Types of Machine Learning Algorithms

Applications: Everywhere

- Banking

- Weather

- Sports scores

- Economics

- Environmental science

- Cybersecurity

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**33**

# Linear Regression – *Prediction*

Problem:

If I have examples of X and Y, when I learn a new X, can I predict Y?

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**34**

Software Engineering Institute | Carnegie Mellon University

# Linear Regression – *Prediction*

Solution: Find the line that is closest to every point

Said differently: Find the line that the SUM of all errors is smallest

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**35**

# Linear Regression – *Prediction*

Three dimensions,
same concept

HUNDREDS of dimensions,
same concept

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**36**

Software Engineering Institute | Carnegie Mellon University

# Linear Regression

Very widely used

- Simple to implement

- Quick to run

- Easy to interpret

- Works for many problems

- First identified in early 1800's; very well studied

When applicable:

- Works best with numeric data (usually)

- Works for predicting specific numeric outcome

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**37**

# Logistic Regression – *Classification*



Idea: Classification using a *discriminative* model

- Predict future behavior based on existing labeled data

- Draws a line to assign labels

Mainly used for binary classification: either "red" or "blue"

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**38**

# Logistic Regression – *Classification*

Look at *distribution*, what's likely based on current data

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**39**

# Logistic Regression

Three dimensions, same concept

HUNDREDS of dimensions, same concept

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**40**

Software Engineering Institute | Carnegie Mellon University

# Classification: Support Vector Machine

Idea: The optimal classifier is the one that is the farthest from both classes

# Classification: Support Vector Machine

Idea: The optimal classifier is the one that is the farthest from both classes

# Classification: Support Vector Machine

Algorithm:

- Find lines like before

- Assign a cost to misclassified data points based on distance from the classification line

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**43**

# Classification: Decision Trees

Idea: Instead of drawing a single complicated line through the data, draw many simpler lines.

Algorithm:

- Scan through all values of all features to find the one that "helps the most" to determine what data gets what label.

- Divide the data based on that value, and then repeat recursively on each part.

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**44**

# Classification: Decision Trees

Idea: Instead of drawing a single complicated line through the data, draw many simpler lines.

Algorithm:

- Scan through all values of all features to find the one that "helps the most" to determine what data gets what label.

- Divide the data based on that value, and then repeat recursively on each part.

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
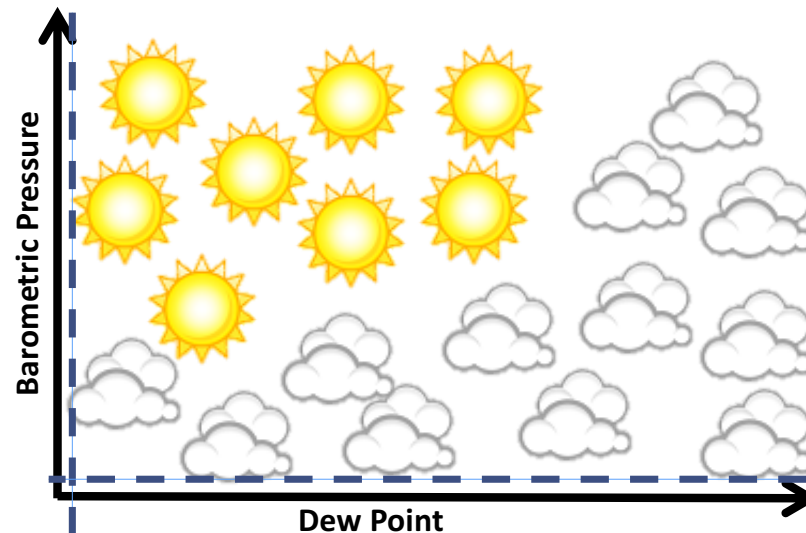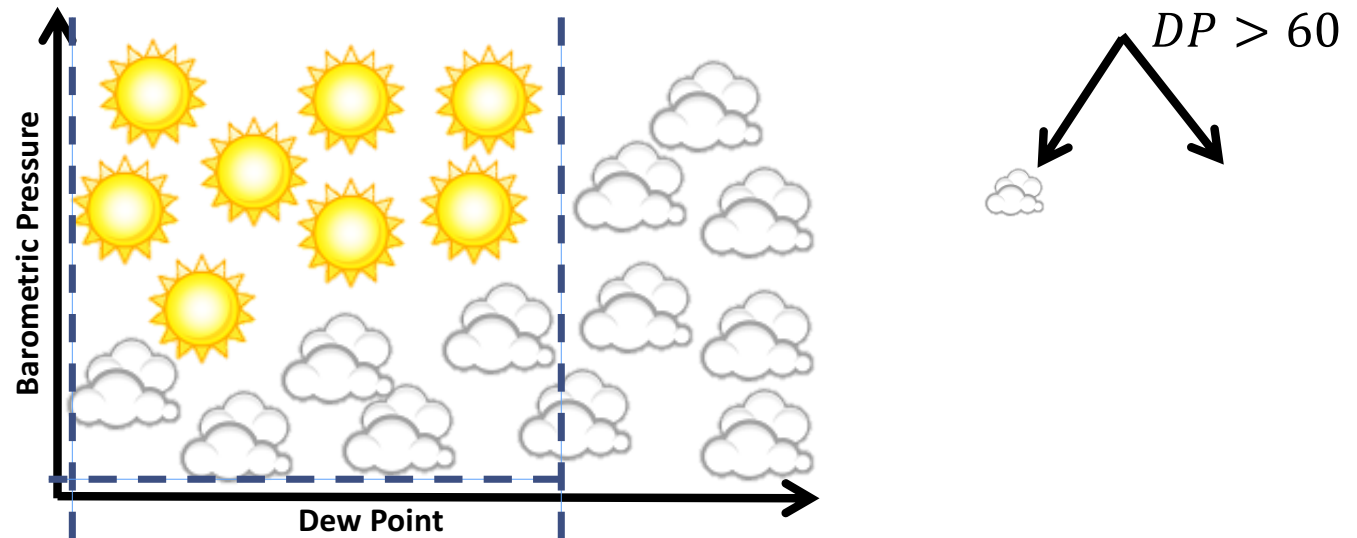Approved for Public Release; Distribution is Unlimited

**45**

# Classification: Decision Trees

Idea: Instead of drawing a single complicated line through the data, draw many simpler lines.

Algorithm:

- Scan through all values of all features to find the one that "helps the most" to determine what data gets what label ("information gain").

- Divide the data based on that value, and then repeat recursively on each part.
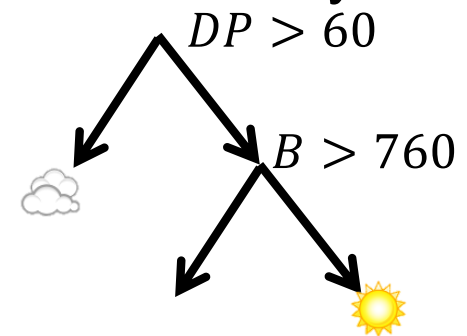
Software Engineering Institute | Carnegie Mellon University

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

2017 SEI Data Science in Cybersecurity Symposium
Approved for Public Release; Distribution is Unlimited

46

# Classification: Decision Trees

Benefits:

- Works well when small.

- Very easy to understand!

Challenges:

- Trees overfit easily

- Very sensitive to data; Random Forests

Software Engineering Institute | Carnegie Mellon University

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

2017 SEI Data Science in Cybersecurity Symposium
Approved for Public Release; Distribution is Unlimited

47

# Classification: K-Nearest Neighbors

Idea: A new point is likely to share the same label as points around it.

Algorithm:

- Pick constant k as number of neighbors to look at.

- For each new point, vote on new label using the k neighbor labels.

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**48**

# Classification: K-Nearest Neighbors

Idea: A new point is likely to share the same label as points around it.

Algorithm:

- Pick constant k as number of neighbors to look at.

- For each new point, vote on new label using the k neighbor labels.

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
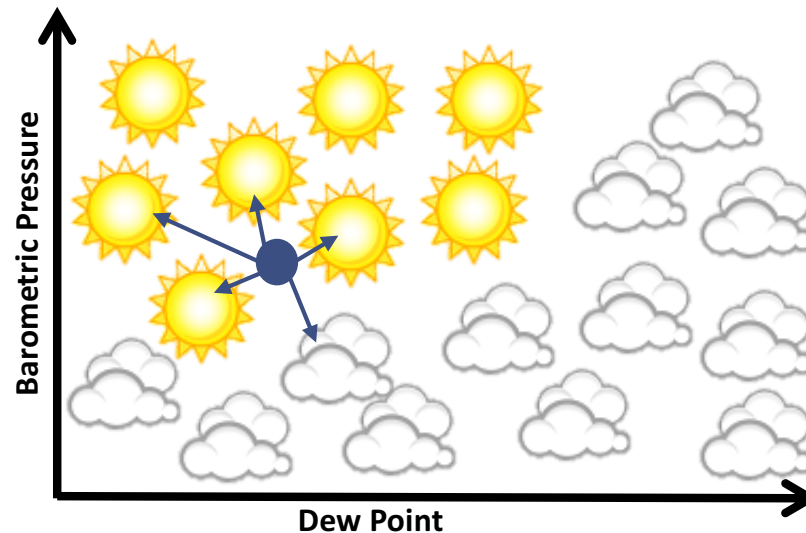Approved for Public Release; Distribution is Unlimited

**49**

# Classification: K-Nearest Neighbors

Idea: A new point is likely to share the same label as points around it.

Algorithm:

- Pick constant k as number of neighbors to look at.

- For each new point, vote on new label using the k neighbor labels.
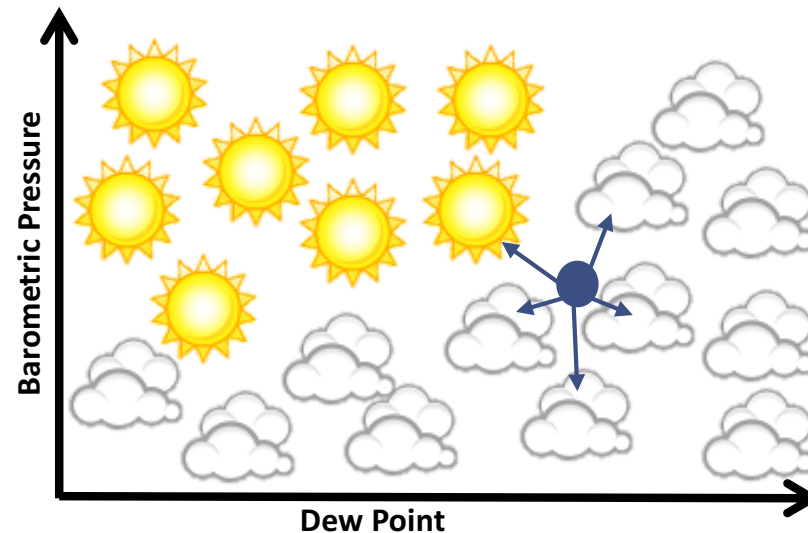
Software Engineering Institute | Carnegie Mellon University

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

2017 SEI Data Science in Cybersecurity Symposium
Approved for Public Release; Distribution is Unlimited

50

# Classification: K-Nearest Neighbors

Works well when

- there is a good distance metric and weighting function to vote on classification

Challenges:

- Not a smooth classifier; points near each other may get classified differently

- Must search all your data every time you want to classify a new point

- When k is small (1,2,3,4), essentially it is overfitting to the data points

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**51**

# Clustering

- Unsupervised learning

- Structure of un-labled data

- Organize records into groups based on some similarity measure

- Cluster is the collection of records which are similar

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**52**

# Clustering: K-means

Idea: Find the clusters by minimizing distances of cluster centers to data.

Algorithm:

- Instantiate k distinct random guesses $\mu_i$ of the cluster centers

- Each data point classifies itself as the $\mu_i$ it is closest to it

- Each $\mu_i$ finds the centroid of the points that were closest to it and jumps there

- Repeat until centroids don't move

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
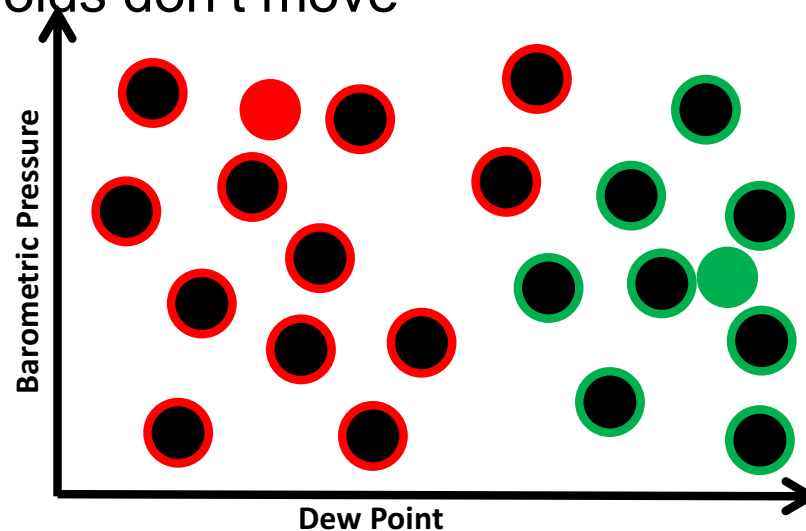Approved for Public Release; Distribution is Unlimited

**53**

# Clustering: K-means

Idea: Find the clusters by minimizing distances of cluster centers to data.

Algorithm:

- Instantiate k distinct random guesses $\mu_i$ of the cluster centers
- Each data point classifies itself as the $\mu_i$ it is closest to it
- Each $\mu_i$ finds the centroid of the points that were closest to it and jumps there
- Repeat until centroids don't move

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
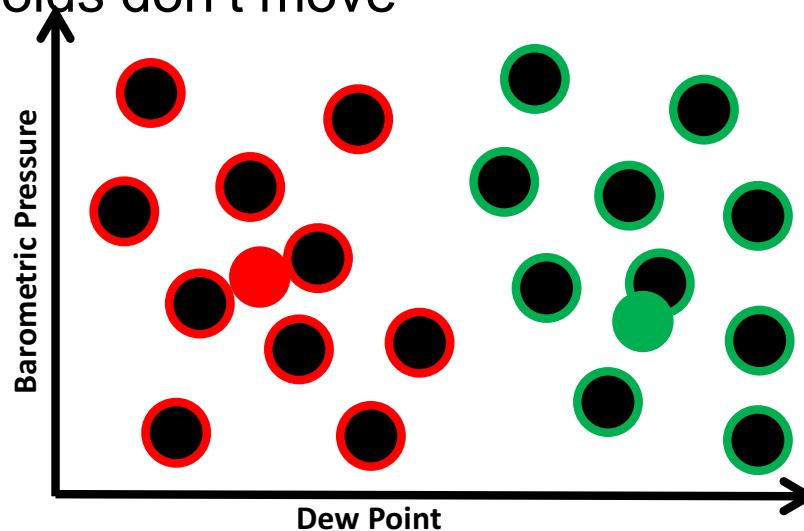Approved for Public Release; Distribution is Unlimited

**54**

# Clustering: K-means

Idea: Find the clusters by minimizing distances of cluster centers to data.

Algorithm:

- Instantiate k distinct random guesses $\mu_i$ of the cluster centers
- Each data point classifies itself as the $\mu_i$ it is closest to it
- Each $\mu_i$ finds the centroid of the points that were closest to it and jumps there
- Repeat until centroids don't move

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
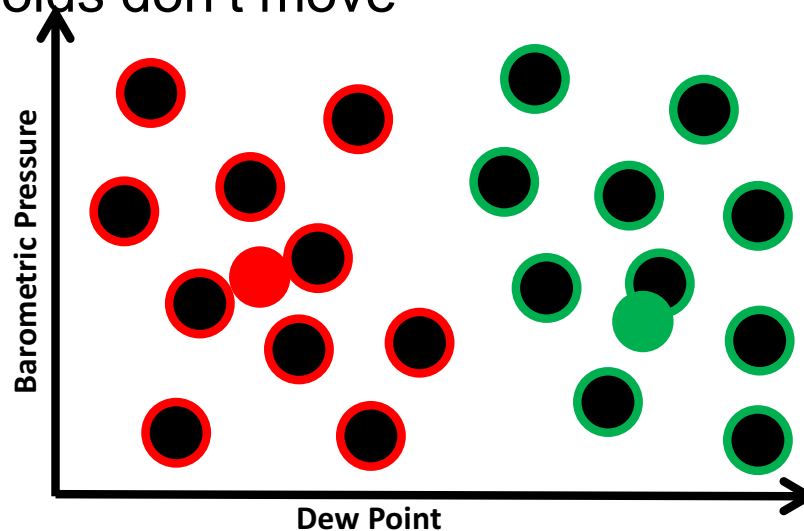Approved for Public Release; Distribution is Unlimited

55

# Clustering: K-means

Idea: Find the clusters by minimizing distances of cluster centers to data.

Algorithm:

- Instantiate k distinct random guesses $\mu_i$ of the cluster centers

- Each data point classifies itself as the $\mu_i$ it is closest to it

- Each $\mu_i$ finds the centroid of the points that were closest to it and jumps there

- Repeat until centroids don't move

Software Engineering Institute | Carnegie Mellon University

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

2017 SEI Data Science in Cybersecurity Symposium
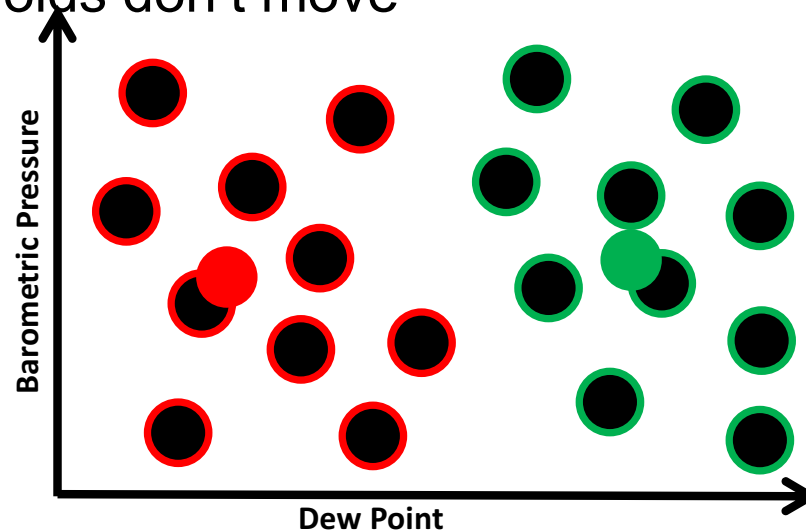Approved for Public Release; Distribution is Unlimited

56

# Clustering: K-means

Idea: Find the clusters by minimizing distances of cluster centers to data.

Algorithm:

- Instantiate k random guesses $\mu_i$ of the clusters

- Each data point classifies itself as the $\mu_i$ it is closest to it

- Each $\mu_i$ finds the centroid of the points that were closest to it and jumps there

- Repeat until centroids don't move

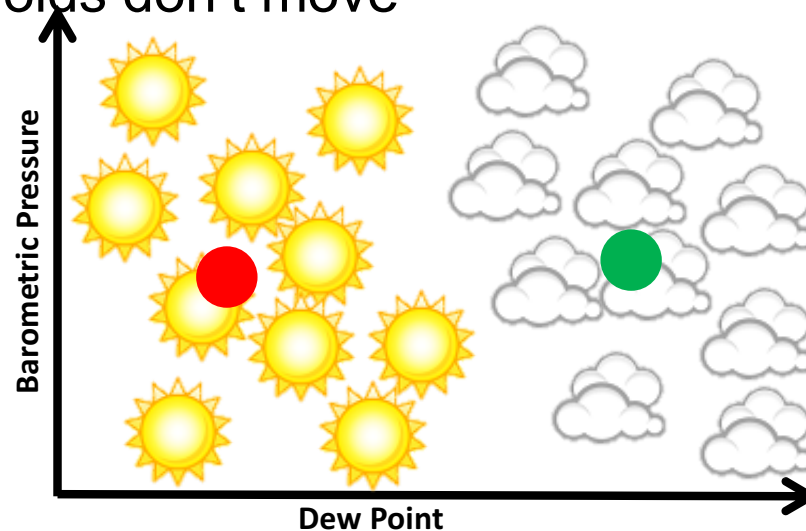**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**57**

# Clustering: K-means

Works well when

- there is a good distance metric between the points

- the number of clusters is known in advance

Challenges:

- Clusters that overlap or are not separable are difficult to cluster correctly.

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**58**

# Influencers

Goal: Detect the people who control or distribute information through a network.

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**59**

# Influencers: Degree Centrality

Idea: Influential people have a lot of people watching them.

Equation

- Degree centrality = number of directed edges to the node

    - High degree centrality people are those with large numbers of followers.

- If undirected graph, transform to bi-directional and compute

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
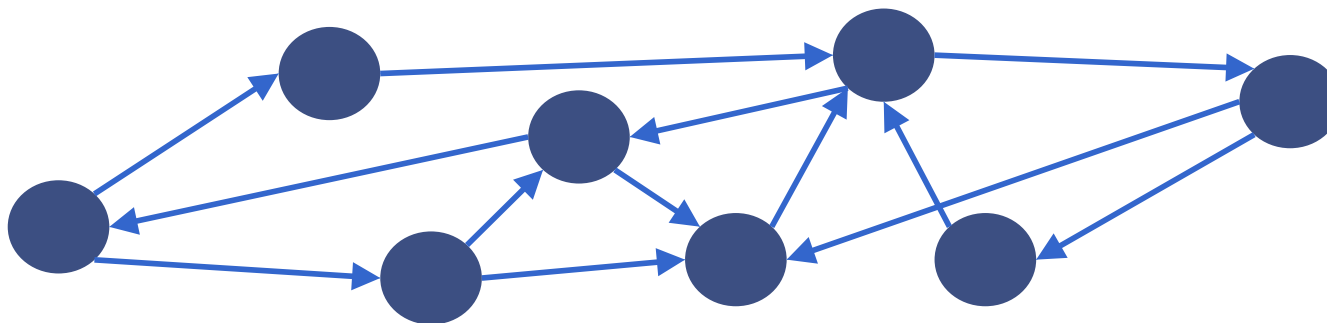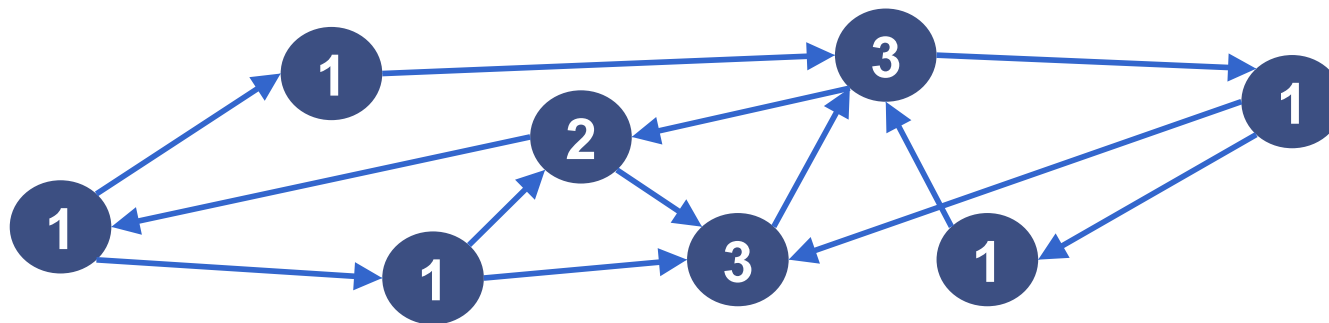Approved for Public Release; Distribution is Unlimited

60

# Influencers: Degree Centrality

Idea: Influential people have a lot of people watching them.

Equation

- Degree centrality = number of directed edges to the node

  - High degree centrality people are those with large numbers of followers.

- If undirected graph, transform to bi-directional and compute

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**61**

# Influencers: Betweenness Centrality

Idea: Influential people are "information brokers" who connect different groups of people.

Algorithm

- Find all shortest paths from all nodes to all other nodes in the graph.

- Betweenness centrality for a node = sum over all start and end nodes of the number of shortest paths in the graph that include the node

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**62**

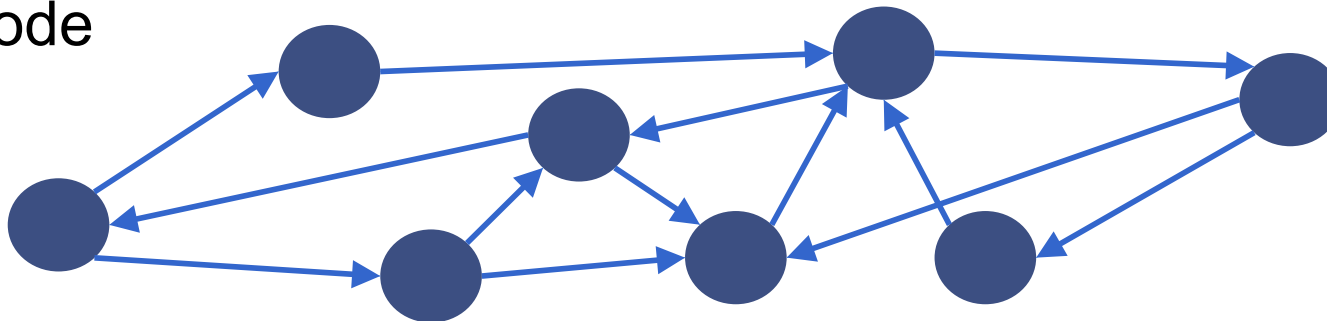Software Engineering Institute | Carnegie Mellon University

# Influencers: Betweenness Centrality

Idea: Influential people are "information brokers" who connect different groups of people.

Algorithm

- Find all shortest paths from all nodes to all other nodes in the graph.

- Betweenness centrality for a node = sum over all start and end nodes of the number of shortest paths in the graph that include the node
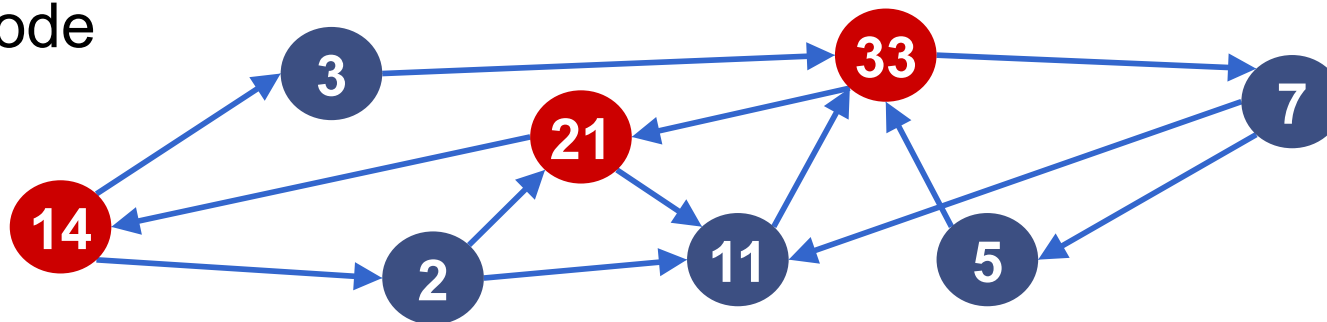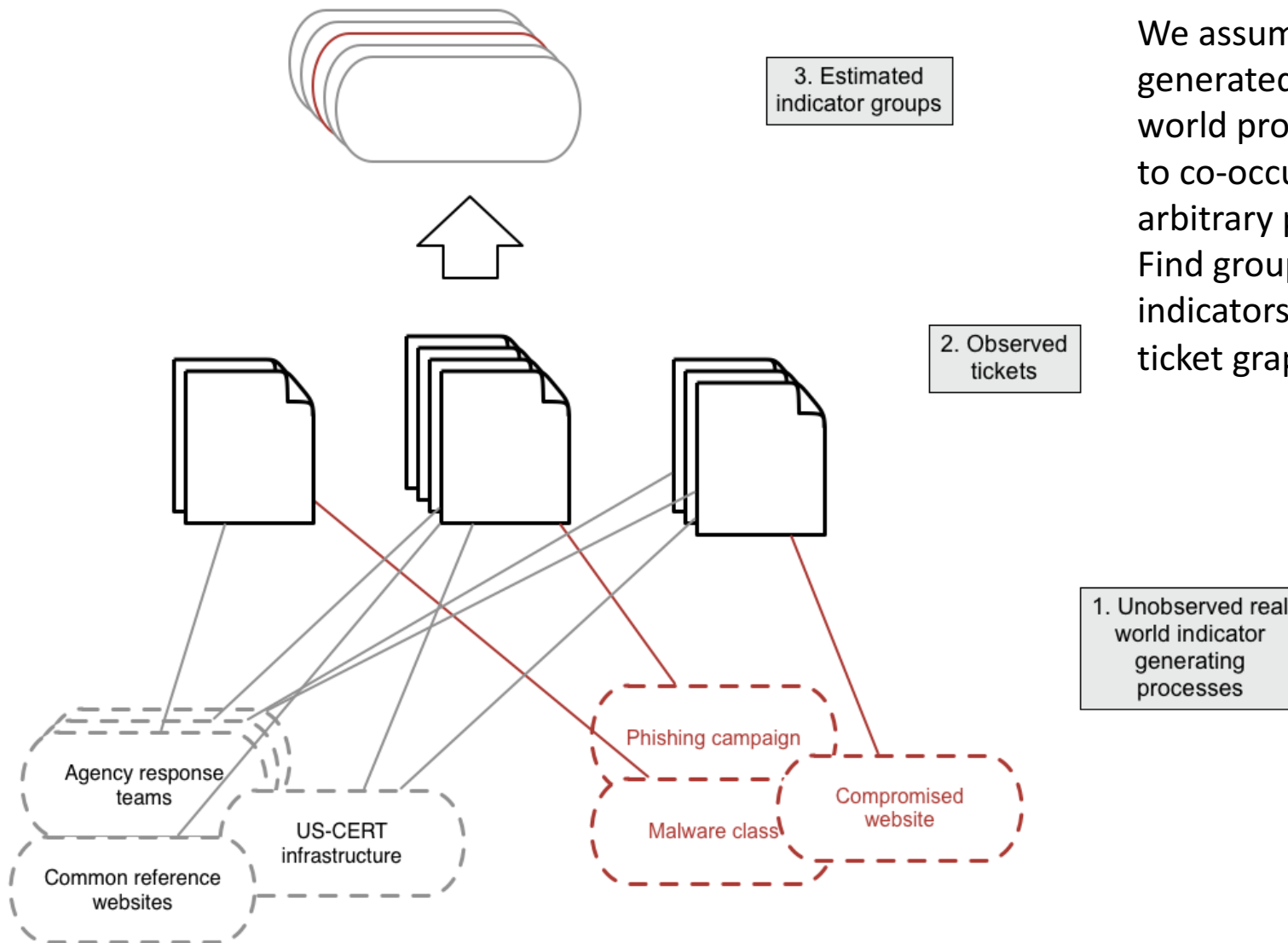
# Indicator communities



3. Estimated indicator groups

2. Observed tickets

1. Unobserved real world indicator generating processes

Agency response teams

Common reference websites

US-CERT infrastructure

Phishing campaign

Malware class

Compromised website

But what if we aren't starting with a reference indicator? We assume that indicators generated by a coherent real world process will be more likely to co-occur in tickets than arbitrary pairs of indicators. Find groups of highly similar indicators in complete indicator-ticket graph.

Software Engineering Institute | Carnegie Mellon University

Data Science Tutorial
August 10, 2017
© 2017 Carnegie Mellon University

2017 SEI Data Science in Cybersecurity Symposium
Approved for Public Release; Distribution is Unlimited

**64**

# Indicator-ticket graph



A subset of the ticket-indicator graph
*(for a small set of selected indicators)*

- Tickets are grey triangles
- Indicators are black circles
- Edges connect tickets to the indicators they contain

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**65**

# Machine Learning Is Growing

Preferred approach for many problems

- Speech recognition

- Natural language processing

- Medical diagnosis

- Robot control

- Sensor networks

- Computer vision

- Weather prediction

- Social network analysis

- AlphaGO, Watson Jeopardy!

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**66**

*This slide also intentionally left blank, just like the earlier one*

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**67**

Software Engineering Institute | Carnegie Mellon University

# What we did today



| Name | Mgr | Dir | Length | Line | Solved? | Comment |
|------|-----|-----|--------|------|---------|---------|
| Beth Jones | Dan Thomas | Anne Kim | 1:30 | 1 | Y | ⑤ ... |
| Beth Jones | Dan Thomas | Anne Kim | 1:52 | 3 | Y | ... |
| Jones, Beth | Dan Thomas | Anne Kim | 90 | 2 | Y | ... |
| Tom Keane | Mark Ryan | Tim Pike | 88 | 2 | N | ... |



Average (5:08)

Call duration (minutes)

"When you put it like that, it makes complete sense."

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**68**

# What we did today

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**69**

# Data Science helps you use data to get results.

**Software Engineering Institute** | **Carnegie Mellon University**

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**70**

**Eliezer Kanal**

Technical Manager

(412) 268-5204

ekanal@sei.cmu.edu

**Daniel DeCapria**

Data Scientist

(412) 268-2457

djdecapria@sei.cmu.edu

Software Engineering Institute | Carnegie Mellon University

**Data Science Tutorial**
August 10, 2017
© 2017 Carnegie Mellon University

**2017 SEI Data Science in Cybersecurity Symposium**
Approved for Public Release; Distribution is Unlimited

**71**