# Data Science and Machine Learning for Non-Data Scientists
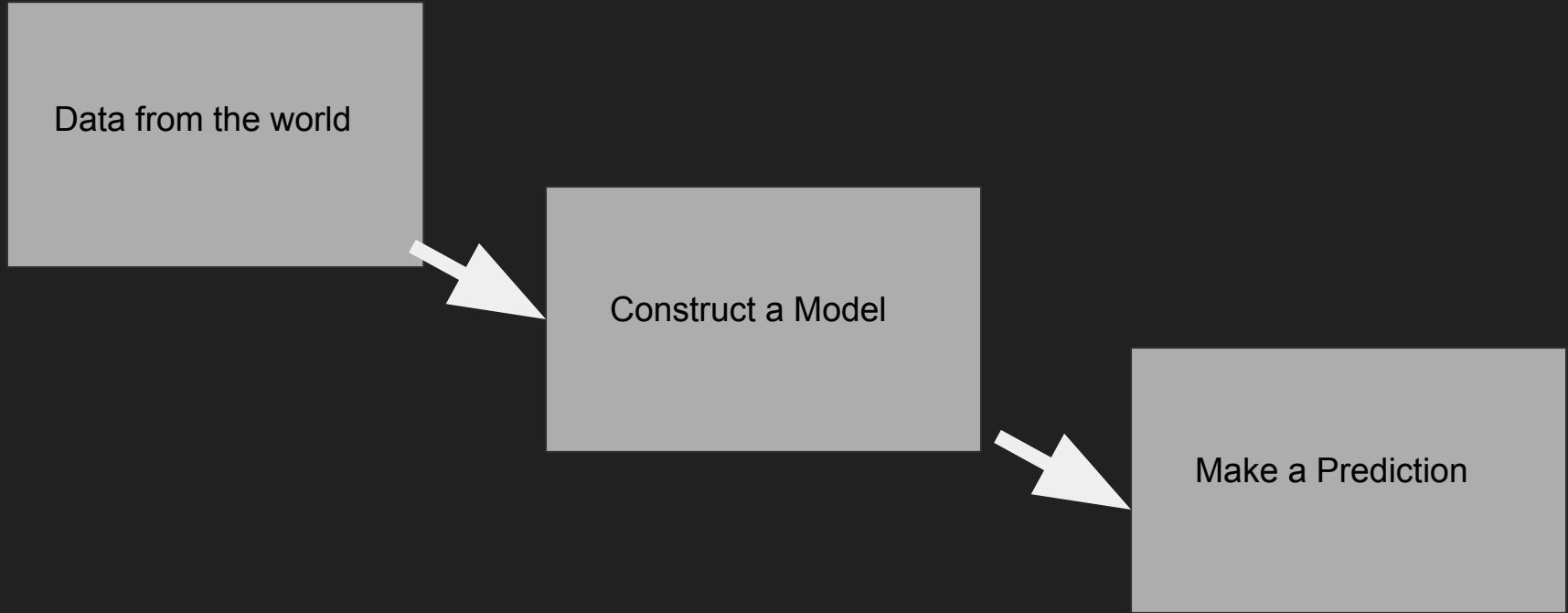
Instructor: Camille Avestruz, PhD
EFI/KICP Postdoctoral Fellow @ University of Chicago
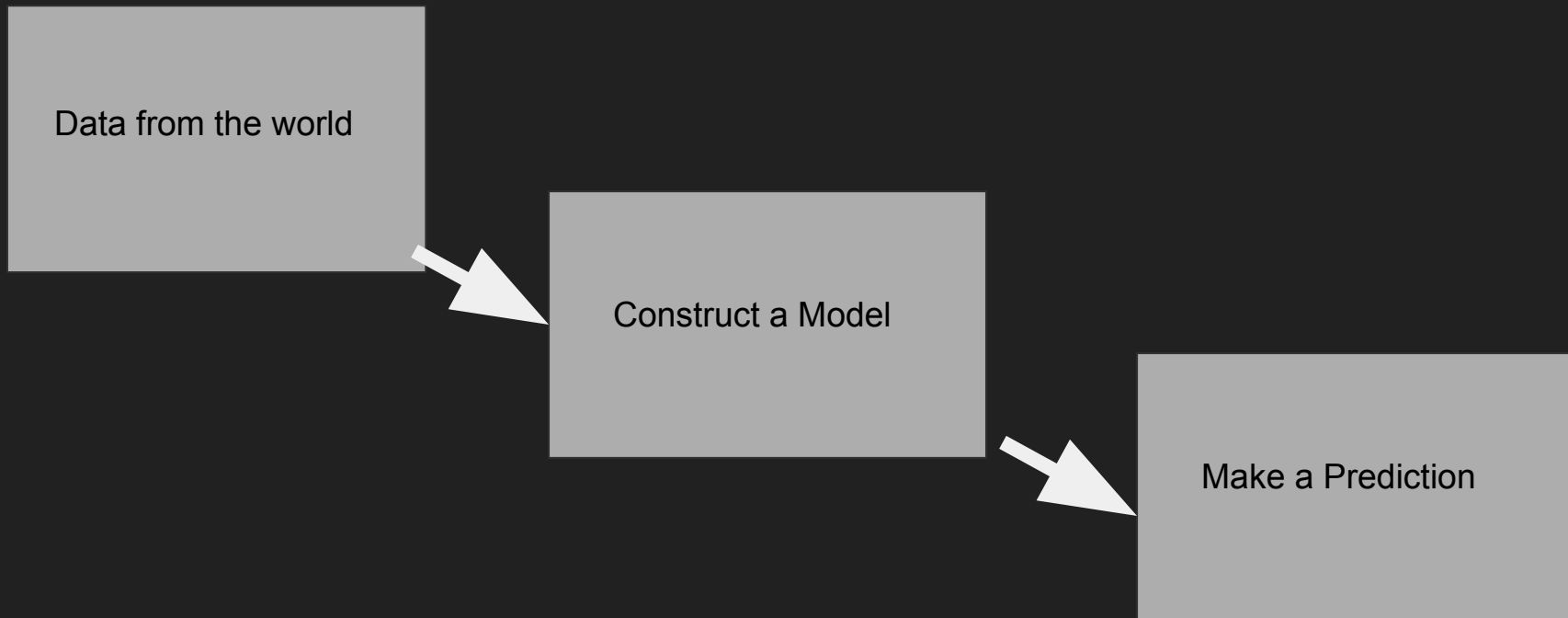co-Instructor: Daniela Huppenkothen, PhD
Assoc Director DIRAC Institute @ U. Washington

# What is data science?

Data from the world

Construct a Model

Make a Prediction

# Who does data science?

Data from the world

Construct a Model

Make a Prediction

# Who does data science?  Everyone!!!



Construct a Model

# Who does data science?
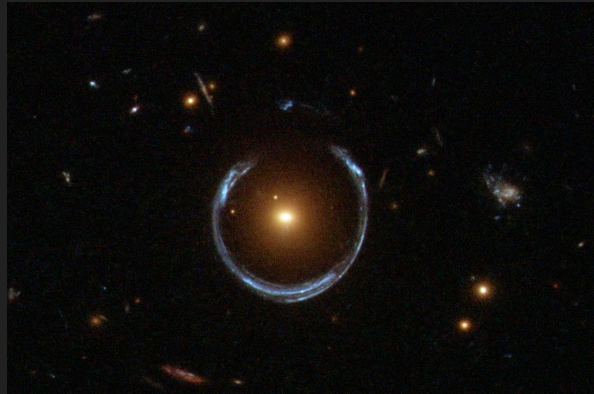
Camille Avestruz

KICP Fellow @ UChicago

Areas of research:
Astronomy, Cosmology

Data sets: Simulations "volumes of our universe", galaxies, clusters of galaxies, ...
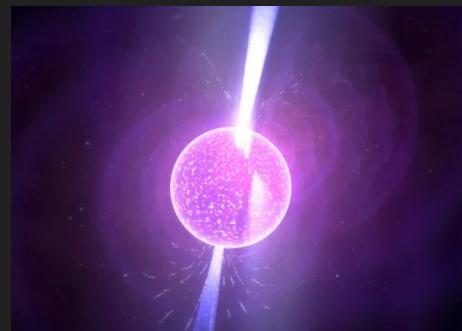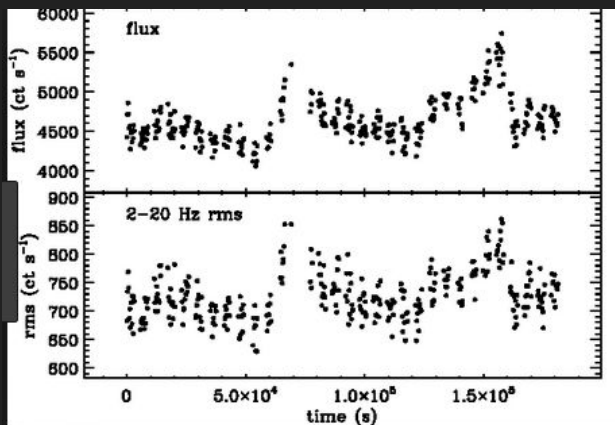
# Who does data science?

Daniela Huppenkothen

Assoc Director @ DIRAC
Institute, U. Washington

Areas of work:
Astrostatistics

Data sets: Time series of
black holes, neutron stars

# Let's do data science: (1) Look at Data

**Zootopia**
**Moonlight**
Deadpool
**Rogue One**

Zootopia
Moonlight
**Deadpool**
**Rogue One**

**Zootopia**
Moonlight
Deadpool
Rogue One

**Zootopia**
**Moonlight**
Deadpool
Rogue One

Zootopia
**Moonlight**
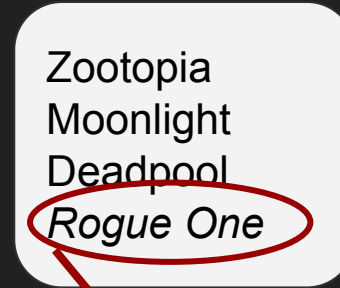Deadpool
Rogue One

Zootopia
**Moonlight**
Deadpool
**Rogue One**

**Zootopia**
**Moonlight**
**Deadpool**
**Rogue One**

**Zootopia**
Moonlight
**Deadpool**
Rogue One

# Let's do data science: (2) Build a Model

Break into groups of 4 - discuss for 5 minutes: Does the person like Rogue One?

# Let's do data science: (2) Build a Model

**Zootopia**
**Moonlight**
Deadpool
**Rogue One**

Zootopia
Moonlight
**Deadpool**
**Rogue One**

**Zootopia**
Moonlight
Deadpool
Rogue One

**Zootopia**
**Moonlight**
Deadpool
Rogue One

Zootopia
**Moonlight**
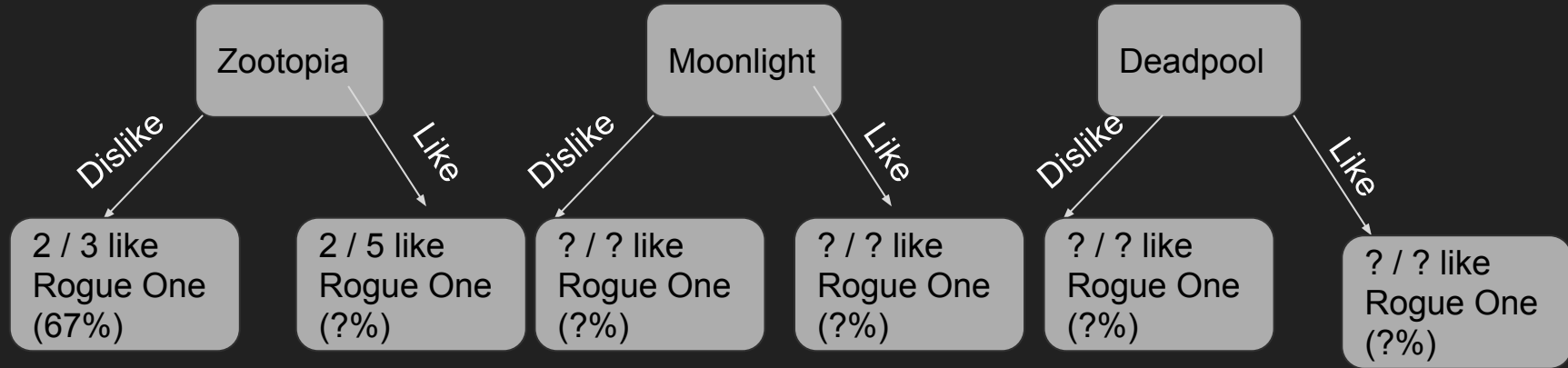Deadpool
Rogue One

Zootopia
**Moonlight**
Deadpool
**Rogue One**

**Zootopia**
**Moonlight**
**Deadpool**
**Rogue One**

**Zootopia**
Moonlight
**Deadpool**
Rogue One

**Zootopia**
**Moonlight**
Deadpool
*Rogue One* **?**

Zootopia
Moonlight
Deadpool
*Rogue One* **?**

# Let's do data science: (3) Make a Prediction

What ideas did you come up with?

# Let's do data science: (2b) Construct a "tree" model

(Fill this out in your groups)

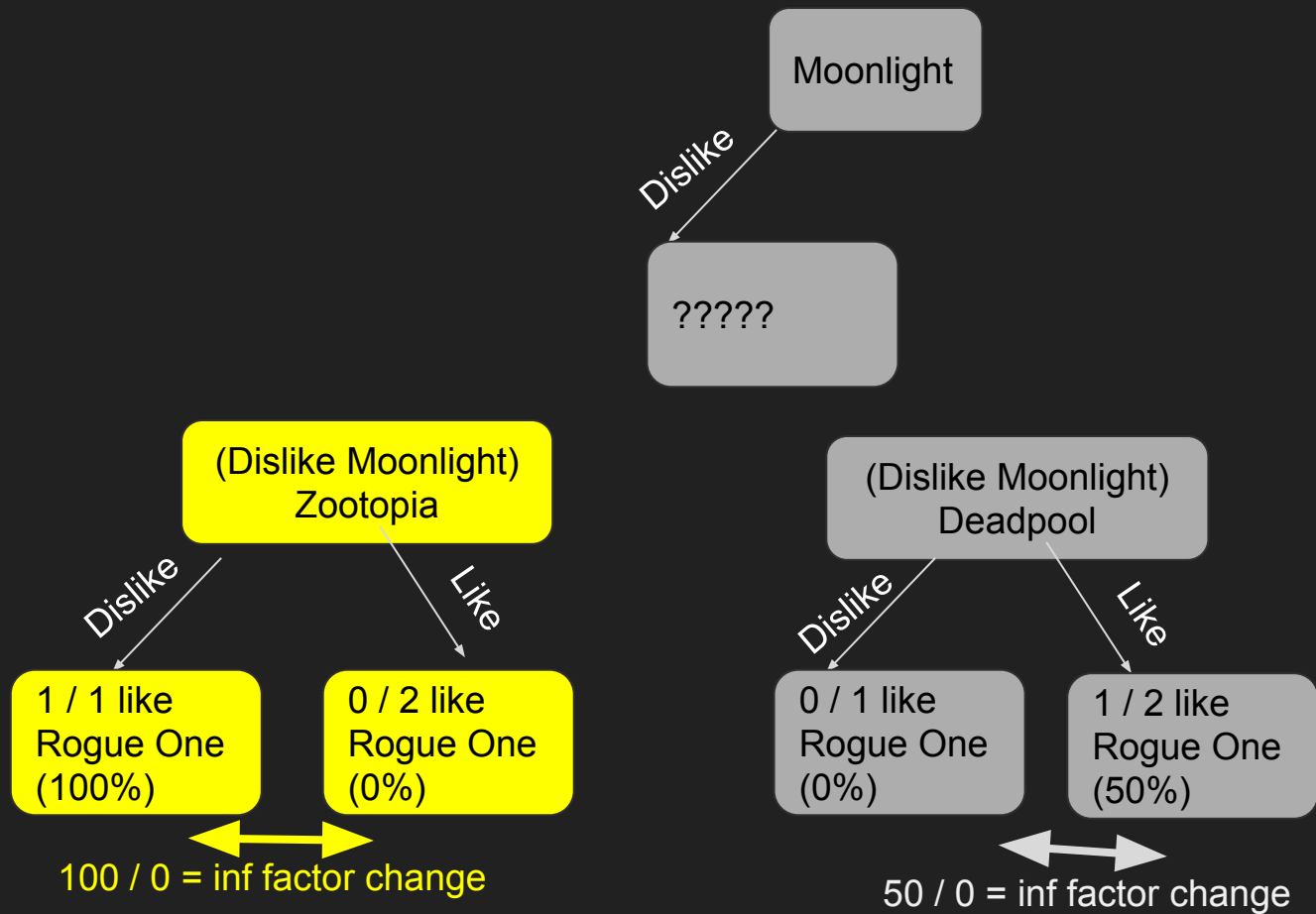# Let's do data science: (2b) Construct a "tree" model
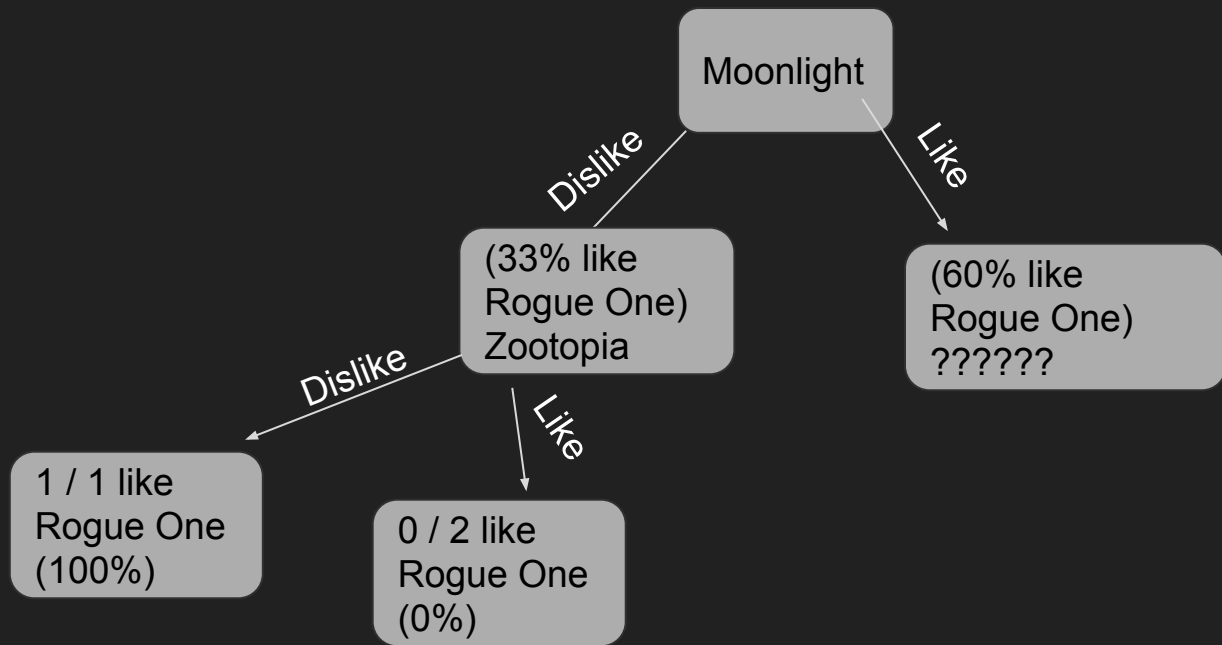
(Fill this out in your groups)

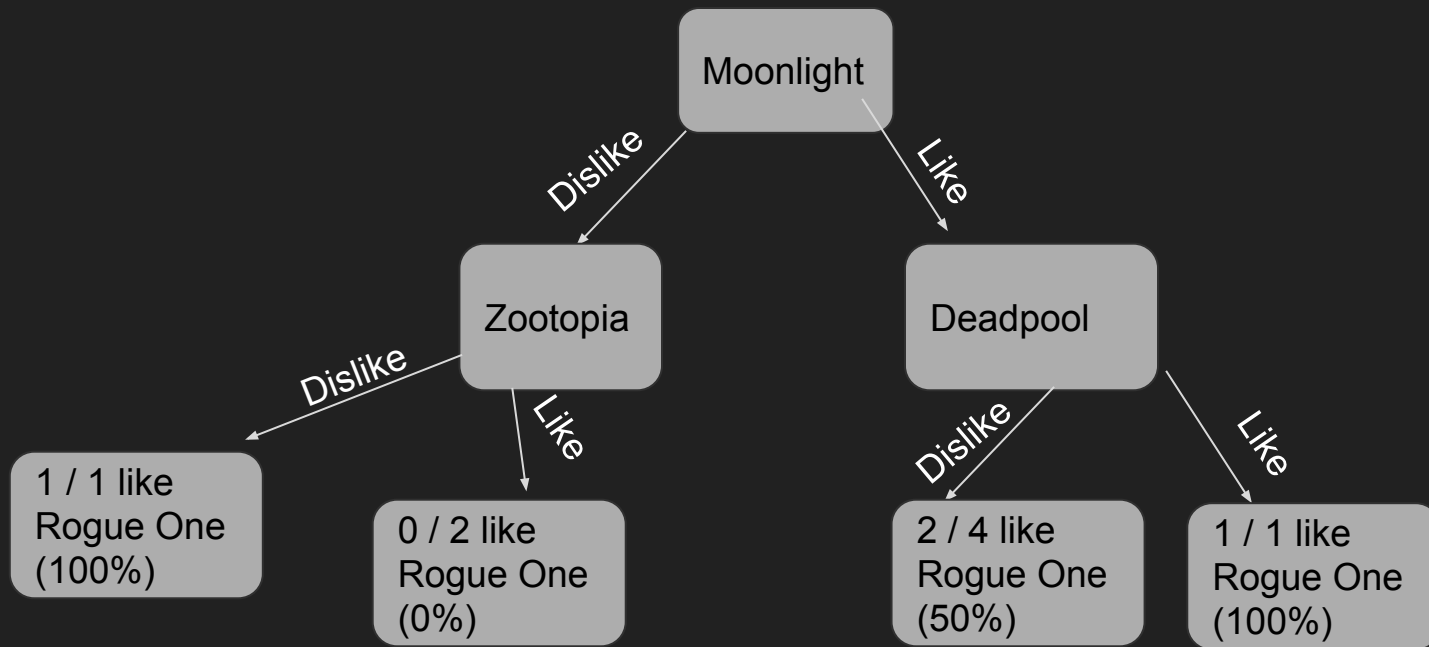# Let's do data science: (2b) A visualized decision tree

Moonlight

Dislike

?????

# Let's do data science: (2b) A visualized decision tree

Moonlight

Dislike

?????

(Dislike Moonlight)
Zootopia

(Dislike Moonlight)
Deadpool

Dislike | Like

1 / 1 like
Rogue One
(100%)

0 / 2 like
Rogue One
(0%)

Dislike | Like

0 / 1 like
Rogue One
(0%)

1 / 2 like
Rogue One
(50%)

100 / 0 = inf factor change

50 / 0 = inf factor change

# Let's do data science: (2b) A visualized decision tree

# Let's do data science: (2b) A visualized decision tree

# Let's do data science: A coding version

Data from the World
I/O: input/output
(Excel .xls file to
Pandas dataframe)

Construct a Model
scikit-learn

Make a Prediction
Visualize Results

# Let's do data science: Use built-in tools

```
In [1]:  import pandas as pd

         from sklearn.tree import DecisionTreeClassifier

         from sklearn.metrics import roc_auc_score, roc_curve

         from sklearn.tree import export_graphviz
         import graphviz

In [2]:  %pylab inline

         Populating the interactive namespace from numpy and matplotlib
```

I/O of data

Model tools (sci-kit learn)

Visualization tools

# Data science we did: Data Input

**Zootopia**
**Moonlight**
Deadpool
**Rogue One**

Zootopia
Moonlight
**Deadpool**
**Rogue One**

**Zootopia**
Moonlight
Deadpool
Rogue One

**Zootopia**
**Moonlight**
Deadpool
Rogue One

Zootopia
**Moonlight**
Deadpool
Rogue One

Zootopia
**Moonlight**
Deadpool
**Rogue One**

**Zootopia**
**Moonlight**
**Deadpool**
**Rogue One**

**Zootopia**
Moonlight
**Deadpool**
Rogue One

# Let's do data science: A coding version (I/O)

Read in the data from Excel

# Let's do data science: Build a Model

Build a decision tree (using tools from python sci-kit learn library)

```
In [4]: model = DecisionTreeClassifier()

In [5]: model.fit(X,y)

Out[5]: DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                max_features=None, max_leaf_nodes=None,
                min_impurity_split=1e-07, min_samples_leaf=1,
                min_samples_split=2, min_weight_fraction_leaf=0.0,
                presort=False, random_state=None, splitter='best')
```

# Let's do data science: Visualize your Predictive Power

```
In [7]: accuracy = float((model.predict(X) == y).sum()) / y.shape[0]
        print accuracy

        0.75

In [8]: data.assign(**{'Rogue One Prediction': model.predict(X)})

Out[8]:
```

|   | Zootopia | Moonlight | Deadpool | Rogue One | Rogue One Prediction |
|---|----------|-----------|----------|-----------|----------------------|
| 0 | 1 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 | 1 |
| 2 | 1 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 0 | 0 | 0 |
| 4 | 0 | 1 | 0 | 0 | 0 |
| 5 | 0 | 1 | 0 | 1 | 0 |
| 6 | 1 | 1 | 1 | 1 | 1 |
| 7 | 1 | 0 | 1 | 0 | 0 |

# Let's do data science: Visualize your Predictive Power

```
In [9]:  auc = roc_auc_score(y, model.predict_proba(X)[:,1])
         fpr, tpr, _ = roc_curve(y, model.predict_proba(X)[:,1])
         plot([0] + list(fpr), [0] + list(tpr), color='r')
         plot([0,1], [0,1], linestyle='--', color='k')
         text(0.25, 0.1, 'Area Under Curve = {}'.format(auc), fontsize='xx-large')
         xlabel('False Positive Rate', fontsize='xx-large')
         ylabel('True Positive Rate', fontsize='xx-large')

Out[9]:  <matplotlib.text.Text at 0x117893650>
```

# Let's do data science: Visualize your Model

```
In [11]: dot_data = export_graphviz(model, out_file=None,
                                     feature_names=data.columns[:-1],
                                     class_names=['Dislikes {}'.format(data.columns[-1]),
                                                  'Likes {}'.format(data.columns[-1])],
                                     filled=True, rounded=True)
         graphviz.Source(dot_data)
```

# Data science we did: (2) The "Machine's" Model

# Data science we did: (3) The "Machine's" Prediction

What did the machine predict?

# Data science we did: (2) The "Machine's" Model

# Data science we did: (3) The "machine's" Prediction

What did the machine predict?

**Zootopia**
**Moonlight**
Deadpool
*Rogue One*

**Dislike**

Zootopia
Moonlight
Deadpool
*Rogue One*

**?**

# Data science we did: (2) The "Machine's" Model

# Data science we did: (3) The "Machine's" Prediction

What did the machine predict?

# Do **your** coding version of data science

(1)   Identify Skills you Wish to Craft

(2)   Join or Develop a Community

(3)   Find Resources

(4)   Practice, practice, practice

# Do your coding version of data science

(1) Identify Skills you Wish to Craft

(2) Join or Develop a Community

(3) Find Resources

(4) Practice, practice, practice

# Do your coding version of data science

(1) Identify Skills you Wish to Craft
(2) Join or Develop a Community
(3) Find Resources
(4) Practice, practice, practice

## The Learning Pyramid*

**Average Retention Rates**

**Passive Teaching Methods**
- 5% Lecture
- 10% Reading
- 20% Audio-Visual
- 30% Demonstration

**Participatory Teaching Methods**
- 50% Group Discussion
- 75% Practice
- 90% Teaching Others

*Adapted from National Training Laboratories. Bethel, Maine

# Do your coding version of data science

(1)  Identify Skills you Wish to Craft
(2)  Join or Develop a Community
(3)  Find Resources
(4)  Practice, practice, practice

Study buddies, online communities, lab co-workers, fellow Amsterdam summer students!



The Learning Pyramid*

Average Retention Rates

| | | |
|---|---|---|
| Passive Teaching Methods | 5% | Lecture |
| | 10% | Reading |
| | 20% | Audio-Visual |
| | 30% | Demonstration |
| Participatory Teaching Methods | 50% | Group Discussion |
| | 75% | Practice |
| | 90% | Teaching Others |

*Adapted from National Training Laboratories. Bethel, Maine

# Do your coding version of data science

(1) Identify Skills you Wish to Craft
(2) Join or Develop a Community
(3) Find Resources
(4) Practice, practice, practice

# Online Lesson Material: Software and Data Carpentry

# Online Lesson Material

Plot data directly from a `Pandas dataframe`.

- We can also plot Pandas dataframes.
- This implicitly uses `matplotlib.pyplot`.

```
import pandas

data = pandas.read_csv('data/gapminder_gdp_oceania.csv', index_col='country')
data.loc['Australia'].plot()
plt.xticks(rotation=90)
```



<
### Plotting and Programming in Python

## Pandas DataFrames

### ❓ Overview

| | |
|---|---|
| **Teaching:** 15 min<br>**Exercises:** 15 min | **Questions**<br>• How can I do statistical analysis of tabular data?<br>**Objectives**<br>• Select individual values from a Pandas dataframe.<br>• Select entire rows or entire columns from a dataframe.<br>• Select a subset of both rows and columns from a dataframe in a single operation.<br>• Select a subset of a dataframe by a single Boolean criterion. |

### Note about Pandas DataFrames/Series

A DataFrame is a collection of Series; The DataFrame is the way Pandas represents a table, and Series is the data-structure Pandas use to represent a column.

Pandas is built on top of the Numpy library, which in practice means that most of the methods defined for Numpy Arrays apply to Pandas Series/DataFrames.

What makes Pandas so attractive is the powerful interface to access individual records of the table, proper handling of missing values, and relational-databases operations between DataFrames.

### Selecting values

To access a value at the position `[i,j]` of a DataFrame, we have two options, depending on what is the meaning of `i` in use. Remember that a DataFrame provides a *index* as a way to identify the rows of the table; a row, then, has a *position* inside the table as well as a *label*, which uniquely identifies its *entry* in the DataFrame.
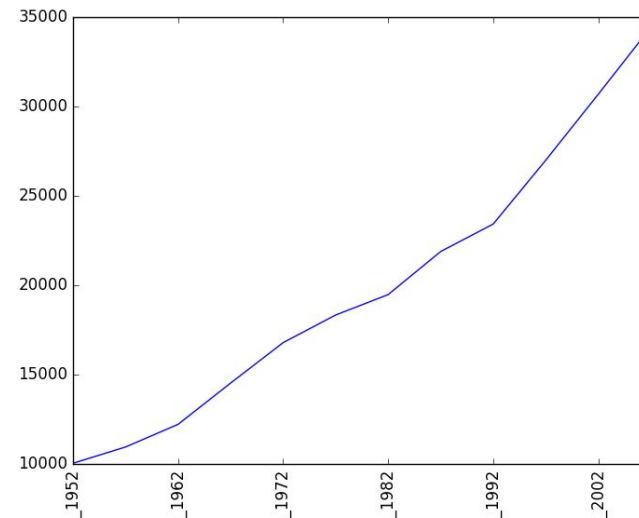
Use `DataFrame.iloc[..., ...]` to select values by their (entry) position

# Online Lesson Material: Host/Attend a Workshop

# Do your coding version of data science

(1) Identify Skills you Wish to Craft
(2) Join or Develop a Community
(3) Find Resources
(4) Practice, practice, practice

# Do your coding version of data science

(1)  Identify Skills you Wish to Craft
(2)  Join or Develop a Community
(3)  Find Resources
(4)  Practice, practice, practice



**LINK TO PRACTICE HERE**