# Leveraging ML for Analysis

# Outline

- Classifying Input

- Features, feature extraction

- Training

- Evaluation

# Types of ML

- Machine Learning (ML) is a computational approach to classifying or labeling types of input

- Two broad approaches
  - Supervised
    - The learning is based on a training set of data that has been labeled in advance (often my hand)
  - Unsupervised
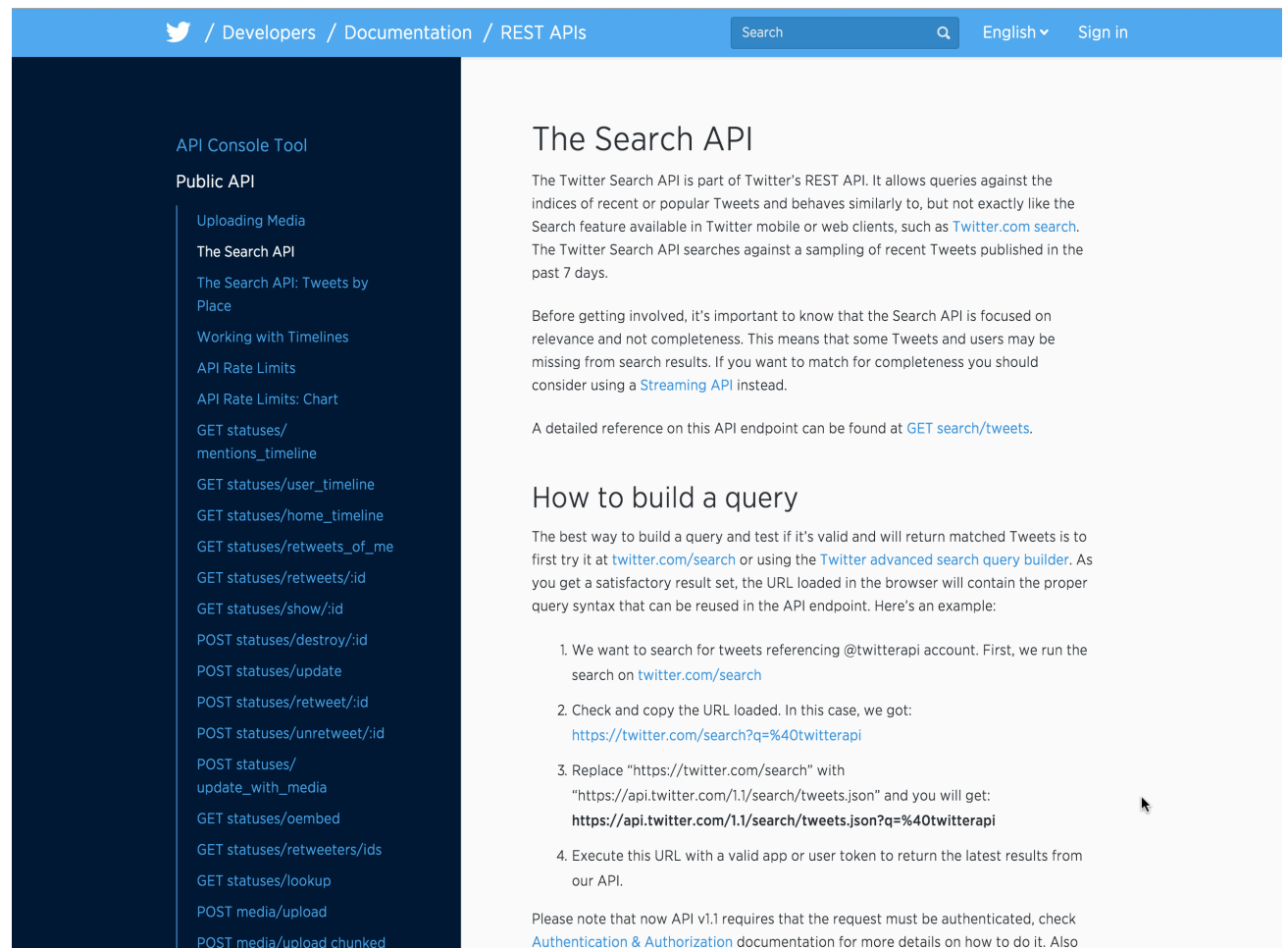    - Learning is inferred from unlabeled data

# Types of Classification/Labeling

- Binary classification
  - Answers the question does this label/classification apply?
    - Yes or No
    - Assume dichotomous labels (classes)

- Multiple classification
  - Answers the question does this input belong to one of several different categories?

# Binary Classifications

- Simple sentiment analysis
  - Is this tweet "happy" or "sad"?

- Generalize to any binary valence
  - Positive to Negative
  - Bright to Dark
  - Introverted to extroverted

- How might this fail?

# Sentiment in Twitter – a Query Operator

# Sentiment in Twitter
# A Query Operator

☐ REST API
   ☐ Search

# Sentiment in Twitter
## A Query Operator

☐ Scroll



**Resources**

Libraries

Sample code

Playbooks

Case studies

Join the community

Events

Developer terms

| | |
|---|---|
| | mentioning Twitter account "NASA". |
| politics filter:safe | containing "politics" with Tweets marked as potentially sensitive removed. |
| puppy filter:media | containing "puppy" and an image or video. |
| puppy filter:native_video | containing "puppy" and an uploaded video, Amplify video, Periscope, or Vine. |
| puppy filter:periscope | containing "puppy" and a Periscope video URL. |
| puppy filter:vine | containing "puppy" and a Vine. |
| puppy filter:images | containing "puppy" and links identified as photos, including third parties such as Instagram. |
| puppy filter:twimg | containing "puppy" and a pic.twitter.com link representing one or more photos. |
| hilarious filter:links | containing "hilarious" and linking to URL. |
| puppy url:amazon | containing "puppy" and a URL with the word "amazon" anywhere within it. |
| superhero since:2015-12-21 | containing "superhero" and sent since date "2015-12-21" (year-month-day). |
| puppy until:2015-12-21 | containing "puppy" and sent before the date "2015-12-21". |
| movie -scary :) | containing "movie", but not "scary", and with a positive attitude. |
| flight :( | containing "flight" and with a negative attitude. |
| traffic ? | containing "traffic" and asking a question. |

Please, make sure to URL encode these queries before making the request. There are several online tools to help you to do that, or you can search at twitter.com/search and copy the encoded URL from the browser's address bar. The table below shows some example mappings from search queries to URL encoded queries:

| Search query | URL encoded query |
|---|---|
| #haiku #poetry | %23haiku+%23poetry |
| "happy hour" :) | %22happy%20hour%22%20%3A%29 |

Note that the space character can be represented by "%20" or "+" sign.

**Additional parameters**

# Twitter Query Operators

| Operator | Finds tweets… |
| --- | --- |
| watching now | containing both "watching" and "now". This is the default operator. |
| "happy hour" | containing the exact phrase "happy hour". |
| love OR hate | containing either "love" or "hate" (or both). |
| beer -root | containing "beer" but not "root". |
| #haiku | containing the hashtag "haiku". |
| from:interior | sent from Twitter account "interior". |
| to:NASA | tweets authored in reply to Twitter account "NASA". |
| @NASA | mentioning Twitter account "NASA". |
| politics filter:safe | containing "politics" with Tweets marked as potentially sensitive removed. |

# Twitter Query Operators

| Operator | Finds tweets… |
|---|---|
| puppy filter:media | containing "puppy" and an image or video. |
| puppy filter:images | containing "puppy" and an image. |
| hilarious filter:links | containing "hilarious" and linking to URL. |
| superhero since:2015-12-21 | containing "superhero" and sent since date "2015-12-21" (year-month-day). |
| puppy until:2015-12-21 | containing "puppy" and sent before the date "2015-12-21". |
| movie -scary :) | containing "movie", but not "scary", and with a positive attitude. |
| flight :( | containing "flight" and with a negative attitude. |
| traffic ? | containing "traffic" and asking a question. |

# Demo

- Try out Twitter Sentiment operators

- How could we try this?

# Other Classification Problems

- Suppose you wanted to classify data using other categories?

- How would you build a classifier?
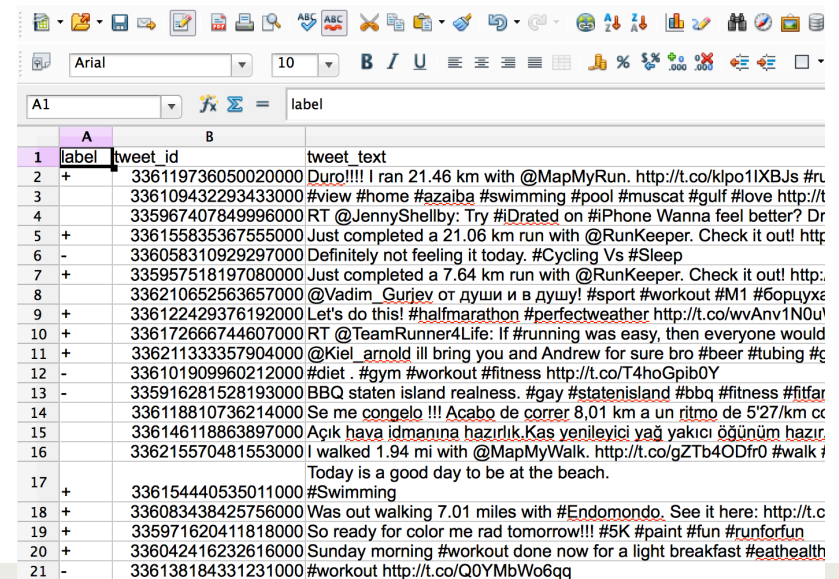
# Process for Creating a Classifier

- Collect Data
- Create a sub-sample
- Pick one (or several) classification algorithms to try
- Select key features
- Score the sub-sample, positive/negative examples
- Train Classifier
- Validate Classifier
- Apply Classifier

# Process for Creating a Classifier

- ◻ Collect Data
- ◻ Create a sub-sample
- ◻ Pick one (or several) classification algorithms to try
- ◻ Select key features
- ◻ Score the sub-sample, positive/negative examples
- ◻ Train Classifier
- ◻ Validate Classifier
- ◻ Apply Classifier

# Samples to Explore

- In hcde user module, ml directory
  - Classification.py – a basic object
  - ClassifyTweet.py – a subclass of Classification

- Sample code
  - explore_feature_selection.py
  - explore_classification.py

# Labeled CSV Tweet data

◻ fitness_label_data1.csv

  ◻ Dump – based on simple_sample.py (using the file output option)

  ◻ Labeled – positive and negative labeling

  ◻ Must have

    ◻ 'label'

    ◻ 'tweet_text'

# Labeled CSV Tweet data

- Two samples for the fitness data
    - fitness_label_data1.csv
    - fitness_label_data2.csv

# Process for Creating a Classifier

- Collect Data
- Create a sub-sample
- Pick a Classifier
- Select key features
- Score the sub-sample, positive/negative examples
- Train Classifier
- Validate Classifier
- Apply Classifier

# Feature Selection

- What are the 'features' of tweets?


- How could you decide which features are important?

# Demo Feature Selection

# Demo Classification

# Interpreting Top Features

```
Most Informative Features
  #Swimming = True          negati : positi =       4.7 : 1.0
  #gym = True               negati : positi =       4.7 : 1.0
  #fitness = True           negati : positi =       3.9 : 1.0
  #RunKeeper = True         positi : negati =       3.4 : 1.0
  completed = True          positi : negati =       3.2 : 1.0
  today. = True             negati : positi =       2.8 : 1.0
  #Workout = True           negati : positi =       2.8 : 1.0
  bring = True              negati : positi =       2.8 : 1.0
```

# Reminder

- Week 9
  - Project "Studio" class session
    - Ray and I will wander from group to group
    - Location TBD