# Project 3: Movie Classification

Welcome to the third project of Data 8! You will build a classifier that guesses whether a movie is romance or action, using only the number of times words appear in the movies's screenplay. By the end of the project, you should know how to:

1. Build a k-nearest-neighbors classifier.
2. Test a classifier on data.

## Logistics

**Deadline.** This project is due at 11:59pm on Friday 12/06. You can earn an early submission bonus point by submitting your completed project by 11:59 on Thursday 12/05. It's **much** better to be early than late, so start working now.

**Checkpoint.** For full credit, you must also **complete Part 2 of the project (out of 4) and submit it by 11:59pm on Friday 11/22**. You will not have lab time to work on these questions, we recommend that you start early on each part to stay on track.

**Partners.** You may work with one other partner; this partner **must** be enrolled in the same lab section as you are. Only one of you is required to submit the project. On okpy.org (http://okpy.org), the person who submits should also designate their partner so that both of you receive credit.

**Rules.** Don't share your code with anybody but your partner. You are welcome to discuss questions with other students, but don't share the answers. The experience of solving the problems in this project will prepare you for exams (and life). If someone asks you for the answer, resist! Instead, you can demonstrate how you would solve a similar problem.

**Support.** You are not alone! Come to office hours, post on Piazza, and talk to your classmates. If you want to ask about the details of your solution to a problem, make a private Piazza post and the staff will respond. If you're ever feeling overwhelmed or don't know how to make progress, email your TA or tutor for help. You can find contact information for the staff on the course website (http://data8.org/fa18/staff.html).

**Tests.** Passing the tests for a question **does not** mean that you answered the question correctly. Tests usually only check that your table has the correct column labels. However, more tests will be applied to verify the correctness of your submission in order to assign your final score, so be careful and check your work!

**Advice.** Develop your answers incrementally. To perform a complicated table manipulation, break it up into steps, perform each step on a different line, give a new name to each result, and check that each intermediate result is what you expect. You can add any additional names or functions you want to the provided cells. Also, please be sure to not re-assign variables throughout the notebook! For example, if you use max_temperature in your answer to one question, do not reassign it later on.

To get started, load `datascience` , `numpy` , `plots` , and `ok` .

```
In [ ]:  # Run this cell to set up the notebook, but please don't change it.
         import numpy as np
         import math
         from datascience import *

         # These lines set up the plotting functionality and formatting.
         import matplotlib
         matplotlib.use('Agg', warn=False)
         %matplotlib inline
         import matplotlib.pyplot as plots
         plots.style.use('fivethirtyeight')
         import warnings
         warnings.simplefilter(action="ignore", category=FutureWarning)

         # These lines load the tests.
         from client.api.notebook import Notebook
         ok = Notebook('project3.ok')
         _ = ok.auth(inline=True)
```

# 1. The Dataset

In this project, we are exploring movie screenplays. We'll be trying to predict each movie's genre from the text of its screenplay. In particular, we have compiled a list of 5,000 words that occur in conversations between movie characters. For each movie, our dataset tells us the frequency with which each of these words occurs in certain conversations in its screenplay. All words have been converted to lowercase.

Run the cell below to read the `movies` table. **It may take up to a minute to load.**

```
In [127]:  movies = Table.read_table('movies.csv')
           movies.where("Title", "the matrix").select(0, 1, 2, 3, 4, 5, 10, 30, 5005)
```

Out[127]:

| Title | Genre | Year | Rating | # Votes | # Words | it | not | fling |
|-------|-------|------|--------|---------|---------|-----|-----|-------|
| the matrix | action | 1999 | 8.7 | 389480 | 3792 | 0.030327 | 0.00870253 | 0 |

The above cell prints a few columns of the row for the action movie *The Matrix*. The movie contains 3792 words. The word "it" appears 115 times, as it makes up $\frac{115}{3792} \approx 0.030327$ of the words in the movie. The word "not" appears 33 times, as it makes up $\frac{33}{3792} \approx 0.00870253$ of the words. The word "fling" doesn't appear at all.

This numerical representation of a body of text, one that describes only the frequencies of individual words, is called a bag-of-words representation. A lot of information is discarded in this representation: the order of the words, the context of each word, who said what, the cast of characters and actors, etc. However, a bag-of-words representation is often used for machine learning applications as a reasonable starting point, because a great deal of information is also retained and expressed in a convenient and compact format. In this project, we will investigate whether this representation is sufficient to build an accurate genre classifier.

All movie titles are unique. The `row_for_title` function provides fast access to the one row for each title.

```
In [ ]: title_index = movies.index_by('Title')
        def row_for_title(title):
            """Return the row for a title, similar to the following expression (but
        faster)

            movies.where('Title', title).row(0)
            """
            return title_index.get(title)[0]

        row_for_title('the terminator')
```

For example, the fastest way to find the frequency of "hey" in the movie *The Terminator* is to access the `'hey'` item from its row. Check the original table to see if this worked for you!

```
In [129]: row_for_title('the terminator').item('hey')
```

```
Out[129]: 0.000540833
```

## Question 1.1

Set `expected_row_sum` to the number that you **expect** will result from summing all proportions in each row, excluding the first six columns.

```
BEGIN QUESTION
name: q1_1
```

```
In [130]: # Set row_sum to a number that's the (approximate) sum of each row of word
          #  proportions.
          expected_row_sum = 1 # SOLUTION
```

This dataset was extracted from a dataset from Cornell University (http://www.cs.cornell.edu/~cristian/Cornell_Movie-Dialogs_Corpus.html). After transforming the dataset (e.g., converting the words to lowercase, removing the naughty words, and converting the counts to frequencies), we created this new dataset containing the frequency of 5000 common words in each movie.

```
In [133]: print('Words with frequencies:', movies.drop(np.arange(6)).num_columns)
          print('Movies with genres:', movies.num_rows)
```

```
Words with frequencies: 5000
Movies with genres: 242
```

# 1.1. Word Stemming

The columns other than "Title", "Genre", "Year", "Rating", "# Votes" and "# Words" in the `movies` table are all words that appear in some of the movies in our dataset. These words have been *stemmed*, or abbreviated heuristically, in an attempt to make different [inflected (https://en.wikipedia.org/wiki/Inflection)](https://en.wikipedia.org/wiki/Inflection) forms of the same base word into the same string. For example, the column "manag" is the sum of proportions of the words "manage", "manager", "managed", and "managerial" (and perhaps others) in each movie. This is a common technique used in machine learning and natural language processing.

Stemming makes it a little tricky to search for the words you want to use, so we have provided another table that will let you see examples of unstemmed versions of each stemmed word. Run the code below to load it.

```
In [134]:  # Just run this cell.
           vocab_mapping = Table.read_table('stem.csv')
           stemmed = np.take(movies.labels, np.arange(3, len(movies.labels)))
           vocab_table = Table().with_column('Stem', stemmed).join('Stem', vocab_mappi
           ng)
           vocab_table.take(np.arange(1100, 1110))
```

Out[134]:

| Stem | Word |
|---|---|
| blame | blamed |
| blame | blame |
| blank | blanks |
| blank | blank |
| blank | blankness |
| blanket | blanket |
| blanket | blankets |
| blast | blasting |
| blast | blast |
| blast | blasted |

**Question 1.1.1**

Assign `stemmed_message` to the stemmed version of the word "alternating".

```
    BEGIN QUESTION
    name: q1_1_1
```

```
In [135]: stemmed_message = vocab_table.where('Word', are.equal_to('alternating')).co
          lumn('Stem').item(0) # SOLUTION
          stemmed_message
```

Out[135]: 'altern'

### Question 1.1.2

Assign `unstemmed_run` to an array of words in `vocab_table` that have "run" as its stemmed form.

```
    BEGIN QUESTION
    name: q1_1_2
```

```
In [138]: unstemmed_run = vocab_table.where('Stem', 'run').column('Word') # SOLUTION
          unstemmed_run
```

Out[138]: array(['runs', 'running', 'run', 'runned', 'runnings'], dtype='<U17')

### Question 1.1.3

Which word in `vocab_table` was shortened the most by this stemming process? Assign `most_shortened` to the word. If there are multiple words, use the word whose first letter is latest in the alphabet (so if your options are albatross or batman, you should pick batman).

It's an example of how heuristic stemming can collapse two unrelated words into the same stem (which is bad, but happens a lot in practice anyway).

```
    BEGIN QUESTION
    name: q1_1_3
```

```
In [142]:  # BEGIN SOLUTION NO PROMPT
           shortened = vocab_table.apply(len, 'Word') - vocab_table.apply(len, 'Stem')
           # SOLUTION
           most_shortened = vocab_table.with_column('Short', shortened)\
                                       .where('Short', are.equal_to(max(shortened)))\
                                       .sort('Word', descending = True)\
                                       .column(1)\
                                       .item(0)
                                        # SOLUTION
           vocab_table.where('Word', most_shortened)
           # END SOLUTION
           """ # BEGIN PROMPT
           # In our solution, we found it useful to first make an array
           # containing the number of characters that was
           # chopped off of each word in vocab_table, but you don't have
           # to do that.

           most_shortened = ...

           # This will display your answer and its shortened form.
           vocab_table.where('Word', most_shortened)
           """; # END PROMPT
```

## 1.2. Splitting the dataset

We're going to use our `movies` dataset for two purposes.

1. First, we want to *train* movie genre classifiers.
2. Second, we want to *test* the performance of our classifiers.

Hence, we need two different datasets: *training* and *test*.

The purpose of a classifier is to classify unseen data that is similar to the training data. Therefore, we must ensure that there are no movies that appear in both sets. We do so by splitting the dataset randomly. The dataset has already been permuted randomly, so it's easy to split. We just take the top for training and the rest for test.

Run the code below (without changing it) to separate the datasets into two tables.

```
In [145]:  # Here we have defined the proportion of our data
           # that we want to designate for training as 17/20ths
           # of our total dataset.  3/20ths of the data is
           # reserved for testing.

           training_proportion = 17/20

           num_movies = movies.num_rows
           num_train = int(num_movies * training_proportion)
           num_test = num_movies - num_train

           train_movies = movies.take(np.arange(num_train))
           test_movies = movies.take(np.arange(num_train, num_movies))

           print("Training: ",   train_movies.num_rows, ";",
                 "Test: ",       test_movies.num_rows)
```

```
Training:  205 ; Test:  37
```

## Question 1.2.1

Draw a horizontal bar chart with two bars that show the proportion of Action movies in each dataset. Complete the function `action_proportion` first; it should help you create the bar chart.
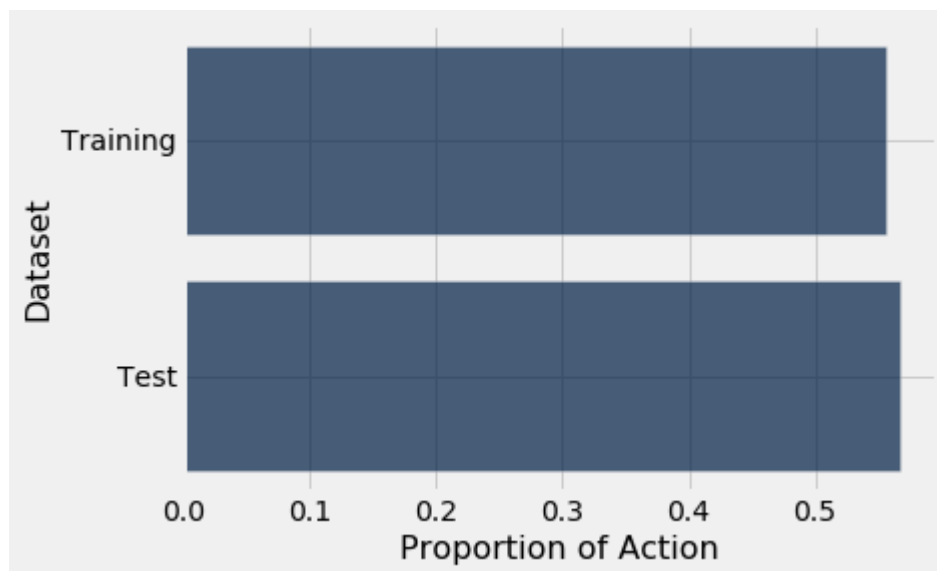
```
BEGIN QUESTION
name: q1_2_1
manual: true
image: true
```

In [146]:
```python
# BEGIN SOLUTION NO PROMPT
def action_proportion(table):
    """Return the proportion of movies in a table that have the Action genre."""
    return table.where('Genre', are.equal_to('action')).num_rows / table.num_rows

Table().with_columns(
        'Dataset', make_array('Training', 'Test'),
        'Proportion of Action', make_array(action_proportion(train_movies),
action_proportion(test_movies)))\
        .barh('Dataset')
# END SOLUTION
""" # BEGIN PROMPT
def action_proportion(table):
    # Return the proportion of movies in a table that have the Action genre.
    return ...

# The staff solution took multiple lines.  Start by creating a table.
# If you get stuck, think about what sort of table you need for barh to work
...."""; # END PROMPT
```

# 2. K-Nearest Neighbors - A Guided Example

K-Nearest Neighbors (k-NN) is a classification algorithm. Given some numerical *attributes* (also called *features*) of an unseen example, it decides whether that example belongs to one or the other of two categories based on its similarity to previously seen examples. Predicting the category of an example is called *labeling*, and the predicted category is also called a *label*.

An attribute (feature) we have about each movie is *the proportion of times a particular word appears in the movies*, and the labels are two movie genres: romance and action. The algorithm requires many previously seen examples for which both the attributes and labels are known: that's the `train_movies` table.

To build understanding, we're going to visualize the algorithm instead of just describing it.

## 2.1. Classifying a movie

In k-NN, we classify a movie by finding the `k` movies in the *training set* that are most similar according to the features we choose. We call those movies with similar features the *nearest neighbors*. The k-NN algorithm assigns the movie to the most common category among its `k` nearest neighbors.

Let's limit ourselves to just 2 features for now, so we can plot each movie. The features we will use are the proportions of the words "money" and "feel" in the movie. Taking the movie "Batman Returns" (in the test set), 0.000502 of its words are "money" and 0.004016 are "feel". This movie appears in the test set, so let's imagine that we don't yet know its genre.

First, we need to make our notion of similarity more precise. We will say that the *distance* between two movies is the straight-line distance between them when we plot their features in a scatter diagram. This distance is called the Euclidean ("yoo-KLID-ee-un") distance, whose formula is $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$.

For example, in the movie *Titanic* (in the training set), 0.0009768 of all the words in the movie are "money" and 0.0017094 are "feel". Its distance from *Batman Returns* on this 2-word feature set is $\sqrt{(0.000502 - 0.0009768)^2 + (0.004016 - 0.0017094)^2} \approx 0.00235496$. (If we included more or different features, the distance could be different.)

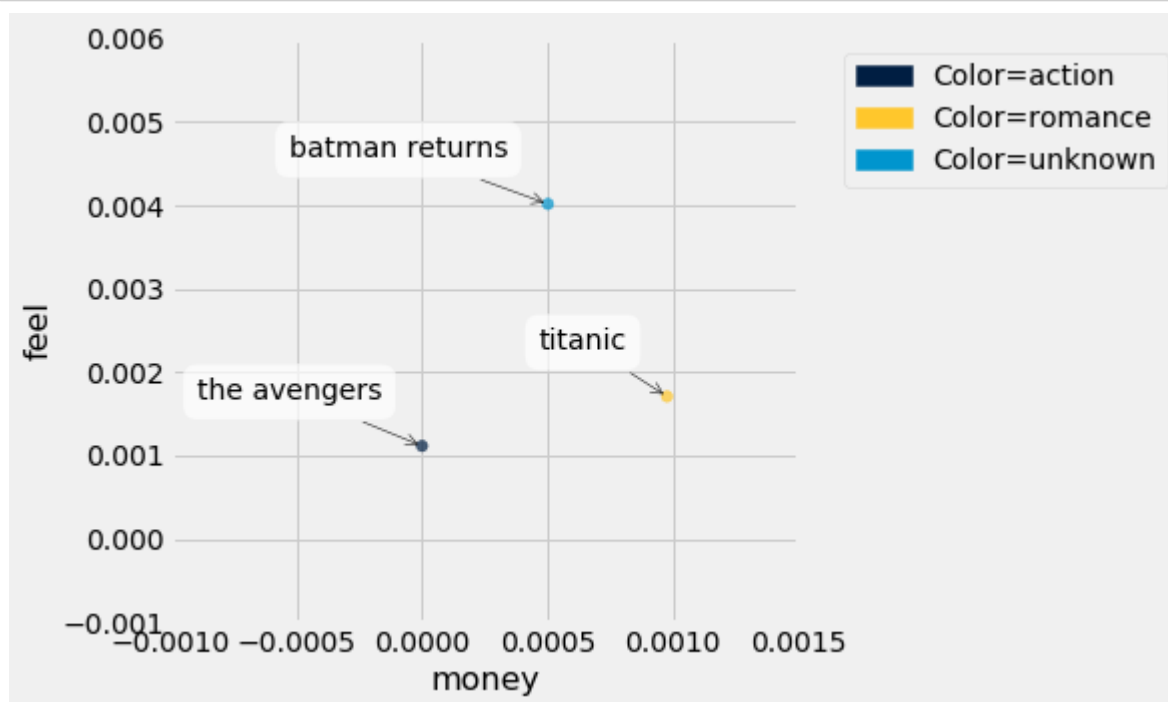A third movie, *The Avengers* (in the training set), is 0 "money" and 0.001115 "feel".

The function below creates a plot to display the "money" and "feel" features of a test movie and some training movies. As you can see in the result, *Batman Returns* is more similar to *Titanic* than to *The Avengers* based on these features. However, we know that *Batman Returns* and *The Avengers* are both action movies, so intuitively we'd expect them to be more similar. Unfortunately, that isn't always the case. We'll discuss this more later.

In [147]:
```python
# Just run this cell.
def plot_with_two_features(test_movie, training_movies, x_feature, y_feature):
    """Plot a test movie and training movies using two features."""
    test_row = row_for_title(test_movie)
    distances = Table().with_columns(
            x_feature, [test_row.item(x_feature)],
            y_feature, [test_row.item(y_feature)],
            'Color',   ['unknown'],
            'Title',   [test_movie]
        )
    for movie in training_movies:
        row = row_for_title(movie)
        distances.append([row.item(x_feature), row.item(y_feature), row.item('Genre'), movie])
    distances.scatter(x_feature, y_feature, group='Color', labels='Title', s=30)

training = ["titanic", "the avengers"]
plot_with_two_features("batman returns", training, "money", "feel")
plots.axis([-0.001, 0.0015, -0.001, 0.006]);
```

## Question 2.1.1

Compute the distance between the two action movies, *Batman Returns* and *The Avengers*, using the `money` and `feel` features only. Assign it the name `action_distance`.

**Note:** If you have a row, you can use `item` to get a value from a column by its name. For example, if `r` is a row, then `r.item("Genre")` is the value in column `"Genre"` in row `r`.

*Hint*: Remember the function `row_for_title`, redefined for you below.

```
BEGIN QUESTION
name: q1_2_1
```

```python
In [148]: title_index = movies.index_by('Title')
          def row_for_title(title):
              """Return the row for a title, similar to the following expression (but
          faster)

              movies.where('Title', title).row(0)
              """
              return title_index.get(title)[0]

          batman = row_for_title("batman returns")
          avengers = row_for_title("the avengers")
          action_distance = ((batman.item("money") - avengers.item("money"))**2 + (ba
          tman.item("feel") - avengers.item("feel"))**2)**0.5 # SOLUTION
          action_distance
```
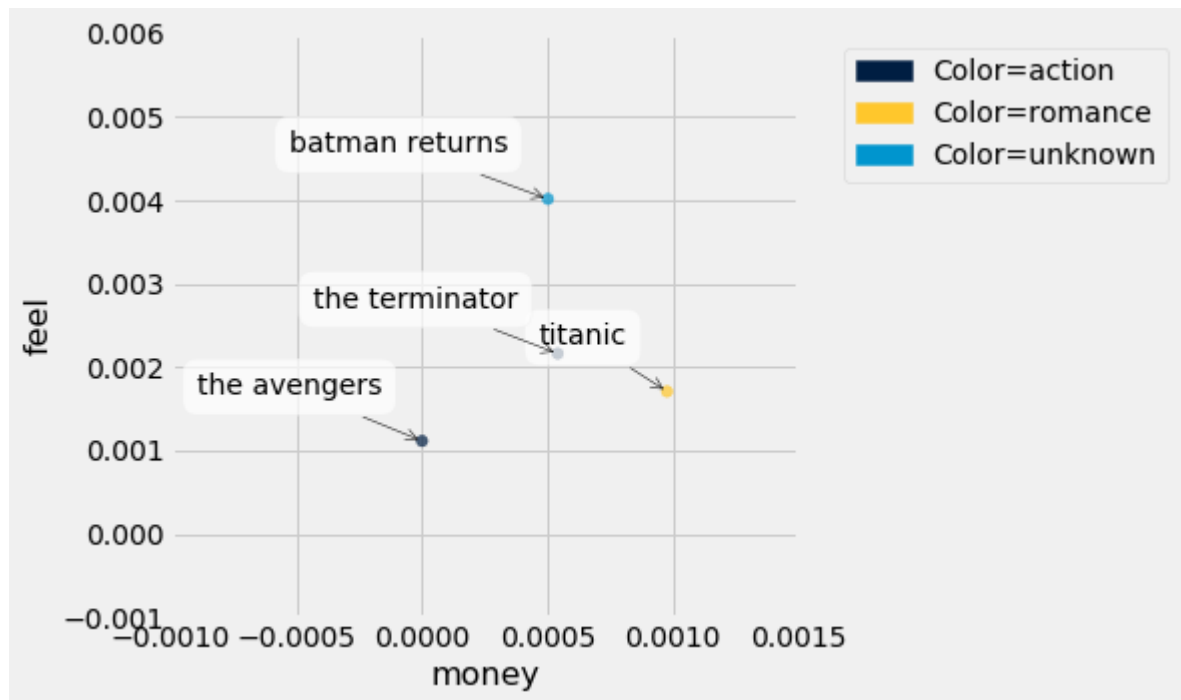
```
Out[148]: 0.0029437356216700243
```

Below, we've added a third training movie, *The Terminator*. Before, the point closest to *Batman Returns* was *Titanic*, a romance movie. However, now the closest point is *The Terminator*, an action movie.

```
In [151]:   training = ["the avengers", "titanic", "the terminator"]
            plot_with_two_features("batman returns", training, "money", "feel")
            plots.axis([-0.001, 0.0015, -0.001, 0.006]);
```



## Question 2.1.2

Complete the function `distance_two_features` that computes the Euclidean distance between any two movies, using two features. The last two lines call your function to show that *Batman Returns* is closer to *The Terminator* than *The Avengers*.

```
BEGIN QUESTION
name: q2_1_2
```

```
In [152]: def distance_two_features(title0, title1, x_feature, y_feature):
              """Compute the distance between two movies with titles title0 and title
          1

              Only the features named x_feature and y_feature are used when computing
          the distance.
              """
              row0 = row_for_title(title0) # SOLUTION
              row1 = row_for_title(title1) # SOLUTION
              return ((row0.item(x_feature) - row1.item(x_feature))**2 + (row0.item(y
          _feature) - row1.item(y_feature))**2)**0.5 # SOLUTION

          for movie in make_array("the terminator", "the avengers"):
              movie_distance = distance_two_features(movie, "batman returns", "money"
          , "feel")
              print(movie, 'distance:\t', movie_distance)
```

```
the terminator distance:        0.0018531387547749904
the avengers distance:   0.0029437356216700243
```

## Question 2.1.3

Define the function `distance_from_batman_returns` so that it works as described in its documentation.

**Note:** Your solution should not use arithmetic operations directly. Instead, it should make use of existing functionality above!

```
    BEGIN QUESTION
    name: q2_1_3
```

```
In [155]: def distance_from_batman_returns(title):
              """The distance between the given movie and "batman returns", based on
           the features "money" and "feel".

              This function takes a single argument:
                title: A string, the name of a movie.
              """

              return distance_two_features(title, "batman returns", "money", "feel")
          # SOLUTION
```

## Question 2.1.4

Using the features `"money"` and `"feel"`, what are the names and genres of the 7 movies in the **training set** closest to *Batman Returns*? To answer this question, make a table named `close_movies` containing those 7 movies with columns `"Title"`, `"Genre"`, `"money"`, and `"feel"`, as well as a column called `"distance from batman"` that contains the distance from *Batman Returns*. The table should be **sorted in ascending order by `distance from batman`**.

```
BEGIN QUESTION
name: q2_1_4
```

```
In [158]: # BEGIN SOLUTION NO PROMPT
          money_feel = train_movies.select("Title", "Genre", "money", "feel")
          close_movies = money_feel.with_column('distance from batman', money_feel.ap
          ply(distance_from_batman_returns, 'Title'))\
                                   .sort('distance from batman')\
                                   .take(np.arange(7))
          # END SOLUTION

          """ # BEGIN PROMPT
          # The staff solution took multiple lines.
          ...
          close_movies = ...
          close_movies
          """; # END PROMPT
```

## Question 2.1.5

Next, we'll clasify *Batman Returns* based on the genres of the closest movies.

To do so, define the function `most_common` so that it works as described in its documentation below.

```
BEGIN QUESTION
name: q2_1_5
```

```
In [162]: def most_common(label, table):
              """The most common element in a column of a table.

              This function takes two arguments:
                label: The label of a column, a string.
                table: A table.

              It returns the most common value in that column of that table.
              In case of a tie, it returns any one of the most common values
              """
              return table.group(label).sort('count', descending=True).column(label).
          item(0) # SOLUTION

          # Calling most_common on your table of 7 nearest neighbors classifies
          # "batman returns" as a romance movie, 5 votes to 2.
          most_common('Genre', close_movies)
```

```
Out[162]: 'romance'
```

Congratulations are in order -- you've classified your first movie! However, we can see that the classifier doesn't work too well since it categorized *Batman Returns* as a romance movie (unless you count the bromance between Alfred and Batman). Let's see if we can do better!

## Checkpoint (Due 11/22)

**Congratulations, you have reached the first checkpoint! Run the submit cell below to generate the checkpoint submission.**

To get full credit for this checkpoint, you must pass all the public autograder tests above this cell.

```
In [ ]: _ = ok.submit()
```

# 3. Features

Now, we're going to extend our classifier to consider more than two features at a time.

Euclidean distance still makes sense with more than two features. For $n$ different features, we compute the difference between corresponding feature values for two movies, square each of the $n$ differences, sum up the resulting numbers, and take the square root of the sum.

### Question 3.1

Write a function called `distance` to compute the Euclidean distance between two **arrays** of **numerical** features (e.g. arrays of the proportions of times that different words appear). The function should be able to calculate the Euclidean distance between two arrays of arbitrary (but equal) length.

Next, use the function you just defined to compute the distance between the first and second movie in the training set *using all of the features*. (Remember that the first six columns of your tables are not features.)

**Note:** To convert rows to arrays, use `np.array`. For example, if `t` was a table, `np.array(t.row(0))` converts row 0 of `t` into an array.

```
BEGIN QUESTION
name: q3_1
```

```
In [165]:   def distance(features_array1, features_array2):
                """The Euclidean distance between two arrays of feature values."""
                return sum((features_array1 - features_array2)**2) ** 0.5 # SOLUTION

            distance_first_to_second = distance(np.array(train_movies.drop(np.arange(6
            )).row(0)), np.array(train_movies.drop(np.arange(6)).row(1))) # SOLUTION
            distance_first_to_second
```

```
Out[165]:   0.03869489784999248
```

# 3.1. Creating your own feature set

Unfortunately, using all of the features has some downsides. One clear downside is *computational* -- computing Euclidean distances just takes a long time when we have lots of features. You might have noticed that in the last question!

So we're going to select just 20. We'd like to choose features that are very *discriminative*. That is, features which lead us to correctly classify as much of the test set as possible. This process of choosing features that will make a classifier work well is sometimes called *feature selection*, or, more broadly, *feature engineering*.

### Question 3.0.1

In this question, we will help you get started on selecting more effective features for distinguishing romance from action movies. The plot below (generated for you) shows the average number of times each word occurs in a romance movie on the horizontal axis and the average number of times it occurs in an action movie on the vertical axis.

 alt text

The following questions ask you to interpret the plot above. For each question, select one of the following choices and assign its number to the provided name.

1. The word is uncommon in both action and romance movies
2. The word is common in action movies and uncommon in romance movies
3. The word is uncommon in action movies and common in romance movies
4. The word is common in both action and romance movies
5. It is not possible to say from the plot

What properties does a word in the bottom left corner of the plot have? Your answer should be a single integer from 1 to 5, corresponding to the correct statement from the choices above.

```
BEGIN QUESTION
name: q3_0_1
```

In [173]: `bottom_left = 1 # SOLUTION`

## Question 3.0.2

What properties does a word in the bottom right corner have?

```
BEGIN QUESTION
name: q3_0_2
```

In [176]: `bottom_right = 3 # SOLUTION`

## Question 3.0.3

What properties does a word in the top right corner have?

```
BEGIN QUESTION
name: q3_0_3
```

In [179]: `top_right = 4 # SOLUTION`

### Question 3.0.4

What properties does a word in the top left corner have?

```
BEGIN QUESTION
name: q3_0_4
```

In [182]: `top_left = 2 # SOLUTION`

### Question 3.0.5

If we see a movie with a lot of words that are common for action movies but uncommon for romance movies, what would be a reasonable guess about the genre of the movie? Assign `movie_genre` to the number corresponding to your answer:

```
1. It is an action movie.
2. It is a romance movie.
```

```
BEGIN QUESTION
name: q3_0_5
```

In [185]: `movie_genre_guess = 1 # SOLUTION`

### Question 3.1.1

Using the plot above, choose 20 common words that you think might let you distinguish between romance and action movies. Make sure to choose words that are frequent enough that every movie contains at least one of them. Don't just choose the 20 most frequent, though... you can do much better.

You might want to come back to this question later to improve your list, once you've seen how to evaluate your classifier.

```
BEGIN QUESTION
name: q3_1_1
```

In [188]:
```python
# Set my_20_features to an array of 20 features (strings that are column la
bels)

my_20_features = make_array('i', 'the', 'to', 'a', 'it', 'and', 'that', 'o
f', 'your', 'what', 'in', 'me', 'is', 'do', 'thi', 'dont', 'he', 'for', 'i
m', 'know') # SOLUTION

# Select the 20 features of interest from both the train and test sets
train_20 = train_movies.select(my_20_features)
test_20 = test_movies.select(my_20_features)
```

This test makes sure that you have chosen words such that at least one appears in each movie. If you can't find words that satisfy this test just through intuition, try writing code to print out the titles of movies that do not contain any words from your list, then look at the words they do contain.

**Question 3.1.2**

In two sentences or less, describe how you selected your features.

```
BEGIN QUESTION
name: q3_1_2
manual: True
```

**SOLUTION:** The staff features don't work very well. It's a good idea to pick words that are common in action movies but uncommon in romance movies and vice-versa. These are points that are far away from the diagonal in the above plot.

Next, let's classify the first movie from our test set using these features. You can examine the movie by running the cells below. Do you think it will be classified correctly?

```
In [194]: print("Movie:")
          test_movies.take(0).select('Title', 'Genre').show()
          print("Features:")
          test_20.take(0).show()
```

Movie:

| Title | Genre |
|---|---|
| the mummy | action |

Features:

| i | the | to | a | it | and | that | of | you |
|---|---|---|---|---|---|---|---|---|
| 0.035313 | 0.0478331 | 0.0324238 | 0.0202247 | 0.0195827 | 0.0199037 | 0.0202247 | 0.0211878 | 0.01123 |

As before, we want to look for the movies in the training set that are most like our test movie. We will calculate the Euclidean distances from the test movie (using the 20 selected features) to all movies in the training set. You could do this with a `for` loop, but to make it computationally faster, we have provided a function, `fast_distances`, to do this for you. Read its documentation to make sure you understand what it does. (You don't need to understand the code in its body unless you want to.)

```
In [195]:  # Just run this cell to define fast_distances.

           def fast_distances(test_row, train_table):
               """Return an array of the distances between test_row and each row in tr
           ain_rows.

               Takes 2 arguments:
                 test_row: A row of a table containing features of one
                   test movie (e.g., test_20.row(0)).
                 train_table: A table of features (for example, the whole
                   table train_20)."""
               assert train_table.num_columns < 50, "Make sure you're not using all th
           e features of the movies table."
               counts_matrix = np.asmatrix(train_table.columns).transpose()
               diff = np.tile(np.array(test_row), [counts_matrix.shape[0], 1]) - count
           s_matrix
               np.random.seed(0) # For tie breaking purposes
               distances = np.squeeze(np.asarray(np.sqrt(np.square(diff).sum(1))))
               eps = np.random.uniform(size=distances.shape)*1e-10 #Noise for tie brea
           k
               distances = distances + eps
               return distances
```

## Question 3.1.3

Use the `fast_distances` function provided above to compute the distance from the first movie in the test set to all the movies in the training set, **using your set of 20 features**. Make a new table called `genre_and_distances` with one row for each movie in the training set and two columns:

- The `"Genre"` of the training movie
- The `"Distance"` from the first movie in the test set

Ensure that `genre_and_distances` is **sorted in increasing order by distance to the first test movie**.

```
BEGIN QUESTION
name: q3_1_3
```

```
In [196]:  # BEGIN SOLUTION NO PROMPT
           distances = fast_distances(test_20.row(0), train_20)

           genre_and_distances = Table().with_columns(
                                       'Genre',    train_movies.column('Genre'),
                                       'Distance', distances)\
                                   .sort('Distance')
           genre_and_distances.show(4)
           # END SOLUTION
           """ # BEGIN PROMPT
           # The staff solution took multiple lines of code.
           genre_and_distances = ...
           genre_and_distances
           """; # END PROMPT
```

| Genre | Distance |
|---|---|
| romance | 0.0144543 |
| action | 0.0144543 |
| action | 0.015057 |
| action | 0.0154917 |

... (201 rows omitted)

## Question 3.1.4

Now compute the 5-nearest neighbors classification of the first movie in the test set. That is, decide on its genre by finding the most common genre among its 5 nearest neighbors in the training set, according to the distances you've calculated. Then check whether your classifier chose the right genre. (Depending on the features you chose, your classifier might not get this movie right, and that's okay.)

```
BEGIN QUESTION
name: q3_1_4
```

```
In [201]:  # Set my_assigned_genre to the most common genre among these.
           my_assigned_genre = most_common("Genre", genre_and_distances.take(np.arange
           (5))) # SOLUTION

           # Set my_assigned_genre_was_correct to True if my_assigned_genre
           # matches the actual genre of the first movie in the test set.
           my_assigned_genre_was_correct = my_assigned_genre == test_movies.column("Ge
           nre").item(0) # SOLUTION

           print("The assigned genre, {}, was{}correct.".format(my_assigned_genre, " "
           if my_assigned_genre_was_correct else " not "))
```

The assigned genre, action, was correct.

# 3.2. A classifier function

Now we can write a single function that encapsulates the whole process of classification.

### Question 3.2.1

Write a function called `classify`. It should take the following four arguments:

- A row of features for a movie to classify (e.g., `test_20.row(0)` ).
- A table with a column for each feature (e.g., `train_20` ).
- An array of classes (e.g. the labels "romance" or "action") that has as many items as the previous table has rows, and in the same order.
- `k`, the number of neighbors to use in classification.

It should return the class a `k` -nearest neighbor classifier picks for the given row of features (the string `'romance'` or the string `'action'` ).

```
    BEGIN QUESTION
    name: q3_2_1
```

```python
In [204]: def classify(test_row, train_rows, train_labels, k):
              """Return the most common class among k nearest neigbors to test_ro
          w."""
              distances = fast_distances(test_row, train_rows)
              genre_and_distances = Table().with_columns('Genre', train_labels, 'Dist
          ance', distances).sort('Distance') # SOLUTION
              return most_common("Genre", genre_and_distances.take(np.arange(k))) # S
          OLUTION
```

### Question 3.2.2

Assign `king_kong_genre` to the genre predicted by your classifier for the movie "king kong" in the test set, using **11 neighbors** and using your 20 features.

```
    BEGIN QUESTION
    name: q3_2_2
```

```python
In [227]: # The staff solution first defined a row called king_kong_features.
          king_kong_features = movies.where("Title", "king kong").select(my_20_featur
          es).row(0) # SOLUTION
          king_kong_genre = classify(king_kong_features, train_20, train_movies.colum
          n("Genre"), 11) # SOLUTION
          king_kong_genre
```

```
Out[227]: 'action'
```

Finally, when we evaluate our classifier, it will be useful to have a classification function that is specialized to use a fixed training set and a fixed value of `k` .

**Question 3.2.3**

Create a classification function that takes as its argument a row containing your 20 features and classifies that row using the 11-nearest neighbors algorithm with `train_20` as its training set.

```
BEGIN QUESTION
name: q3_2_3
```

```
In [229]: def classify_feature_row(row):
              return classify(row, train_20, train_movies.column('Genre'), 11) # SOLU
          TION

              # When you're done, this should produce 'Romance' or 'Action'.
              classify_feature_row(test_20.row(0))
```

```
Out[229]: 'action'
```

# 3.3. Evaluating your classifier

Now that it's easy to use the classifier, let's see how accurate it is on the whole test set.

**Question 3.3.1.** Use `classify_feature_row` and `apply` to classify every movie in the test set. Assign these guesses as an array to `test_guesses` . **Then**, compute the proportion of correct classifications.

```
BEGIN QUESTION
name: q3_3_1
```

```
In [241]: test_guesses = test_20.apply(classify_feature_row) # SOLUTION
          proportion_correct = np.count_nonzero(test_guesses == test_movies.column('G
          enre')) / test_movies.num_rows # SOLUTION
          proportion_correct
```

```
Out[241]: 0.7297297297297297
```

**Question 3.3.2.** An important part of evaluating your classifiers is figuring out where they make mistakes. Assign the name `test_movie_correctness` to a table with three columns, `'Title'`, `'Genre'`, and `'Was correct'`. The last column should contain `True` or `False` depending on whether or not the movie was classified correctly.

```
BEGIN QUESTION
name: q3_3_2
```

In [244]:
```python
# Feel free to use multiple lines of code
# but make sure to assign test_movie_correctness to the proper table!
test_movie_correctness = test_movies.with_column('Was correct', test_guesses == test_movies.column('Genre')).select('Title', 'Genre', 'Was correct') # SOLUTION
test_movie_correctness.sort('Was correct', descending = True).show(5)
```

| Title | Genre | Was correct |
|---|---|---|
| batman returns | action | True |
| american outlaws | action | True |
| solaris | romance | True |
| independence day | action | True |
| my girl 2 | romance | True |

... (32 rows omitted)

**Question 3.3.3.** Do you see a pattern in the types of movies your classifier misclassifies? In two sentences or less, describe any patterns you see in the results or any other interesting findings from the table above. If you need some help, try looking up the movies that your classifier got wrong on Wikipedia.

```
BEGIN QUESTION
name: q3_3_3
manual: true
```

**SOLUTION:** The classifier tends to misclassify movies that have both action and romance in them.

At this point, you've gone through one cycle of classifier design. Let's summarize the steps:

1. From available data, select test and training sets.
2. Choose an algorithm you're going to use for classification.
3. Identify some features.
4. Define a classifier function using your features and the training set.
5. Evaluate its performance (the proportion of correct classifications) on the test set.

# 4. Explorations

Now that you know how to evaluate a classifier, it's time to build a better one.

### Question 4.1

Develop a classifier with better test-set accuracy than `classify_feature_row`. Your new function should have the same arguments as `classify_feature_row` and return a classification. Name it `another_classifier`. Then, check your accuracy using code from earlier.

You can use more or different features, or you can try different values of `k`. (Of course, you still have to use `train_movies` as your training set!)

**Make sure you don't reassign any previously used variables here**, such as `proportion_correct` from the previous question.

```
In [248]:  # To start you off, here's a list of possibly-useful features
           # Feel free to add or change this array to improve your classifier
           new_features = make_array("come", "do", "have", "heart", "make", "never",
           "now", "wanna", "with", "yo")

           train_new = train_movies.select(new_features)
           test_new = test_movies.select(new_features)

           def another_classifier(row):
               return ...
```

### Question 4.2

Do you see a pattern in the mistakes your new classifier makes? What about in the improvement from your first classifier to the second one? Describe in two sentences or less.

**Hint:** You may not be able to see a pattern.

```
BEGIN QUESTION
name: q4_2
manual: true
```

**SOLUTION:** Any reasonable student solution is fine.

**Question 4.3**

Briefly describe what you tried to improve your classifier.

```
BEGIN QUESTION
name: q4_3
manual: true
```

**SOLUTION:** You were meant to put in at least some minimal effort.

Congratulations: you're done with the required portion of the project! Time to submit.

```
In [ ]:  _ = ok.submit()
```

# 5. Other Classification Methods (OPTIONAL)

**Note**: Everything below is **OPTIONAL**. Please only work on this part after you have finished and submitted the project. If you create new cells below, do NOT reassign variables defined in previous parts of the project.

Now that you've finished your k-NN classifier, you might be wondering what else you could do to improve your accuracy on the test set. Classification is one of many machine learning tasks, and there are plenty of other classification algorithms! If you feel so inclined, we encourage you to try any methods you feel might help improve your classifier.

We've compiled a list of blog posts with some more information about classification and machine learning. Create as many cells as you'd like below--you can use them to import new modules or implement new algorithms.

Blog posts:

- Classification algorithms/methods (https://medium.com/@sifium/machine-learning-types-of-classification-9497bd4f2e14)
- Train/test split and cross-validation (https://towardsdatascience.com/train-test-split-and-cross-validation-in-python-80b61beca4b6)
- More information about k-nearest neighbors (https://medium.com/@adi.bronshtein/a-quick-introduction-to-k-nearest-neighbors-algorithm-62214cea29c7)
- Overfitting (https://elitedatascience.com/overfitting-in-machine-learning)

In future data science classes, such as Data Science 100, you'll learn about some about some of the algorithms in the blog posts above, including logistic regression. You'll also learn more about overfitting, cross-validation, and approaches to different kinds of machine learning problems.

There's a lot to think about, so we encourage you to find more information on your own!

Modules to think about using:

- Scikit-learn tutorial (http://scikit-learn.org/stable/tutorial/basic/tutorial.html)
- TensorFlow information (https://www.tensorflow.org/tutorials/)

...and many more!

```
In [ ]:  ...
```