

In [319]: stemmed_message = vocab_table.where('Word', 'alternating').column('Stem').item(0)
stemmed_message

Out[319]: 'altern'

In [320]: ok.grade("q1_1_1");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

Question 1.1.2

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

In [321]: unstemmed_run = vocab_table.where('Stem', 'run').column('Word')
unstemmed_run

Out[321]: array(['runs', 'running', 'run', 'runned', 'runnings'], dtype='<U17')

In [322]: ok.grade("q1_1_2");

Running tests

Test summary

Passed: 2
Failed: 0
[oooooooooook] 100.0% passed

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).

In [323]:

```
# 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.
len_stem = vocab_table.apply(len, 'Stem')
len_word = vocab_table.apply(len, 'Word')
newtbl = vocab_table.with_column('Added Length', len_word - len_stem)

most_shortened_table = newtbl.where('Added Length', max(newtbl.column('Added Length'))).sort('Word')
descending = True)
most_shortened = most_shortened_table.column('Word').item(0)

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

Out[323]:

Stem	Word
respons	responsibilities

In [324]: ok.grade("q1_1_3");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

1.2. Splitting the dataset

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

- First, we want to *train* movie genre classifiers.
- 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 [325]:

```
# 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

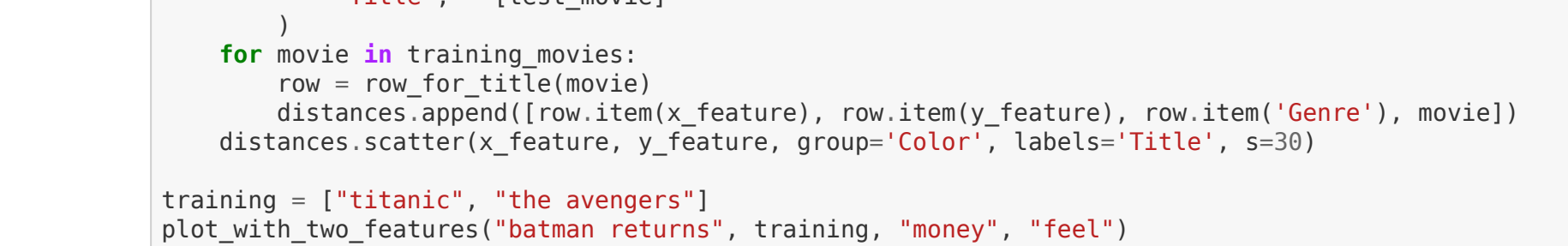
Draw a horizontal bar chart with two bars that show the proportion of action movies in each dataset. Complete the function `action_proportion` first; it should help you create the bar chart.

In [326]:

```
def action_proportion(table):
    # Return the proportion of movies in a table that have the Action genre.
    return table.group('Genre').barh('Genre')

action_proportion(movies)
```

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



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-KLUD-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.)

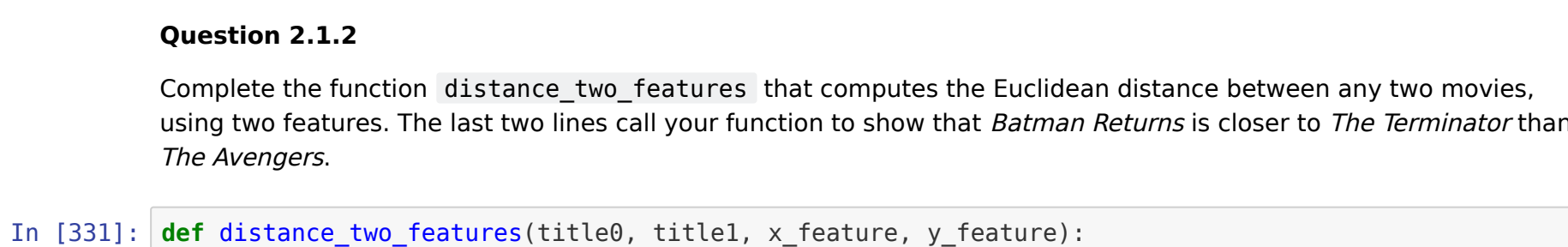
A third movie, *The Avengers* (in the training set), is 0 "money" and 0.001115 "feel".

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 [327]:

```
# 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('Color'), 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.

In [328]:

```
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 + (batman.item("feel") - avengers.item("feel"))**2)**0.5
action_distance
```

Out[328]: 0.0029437356216780243

In [329]: ok.grade("q1_2_1");

Running tests

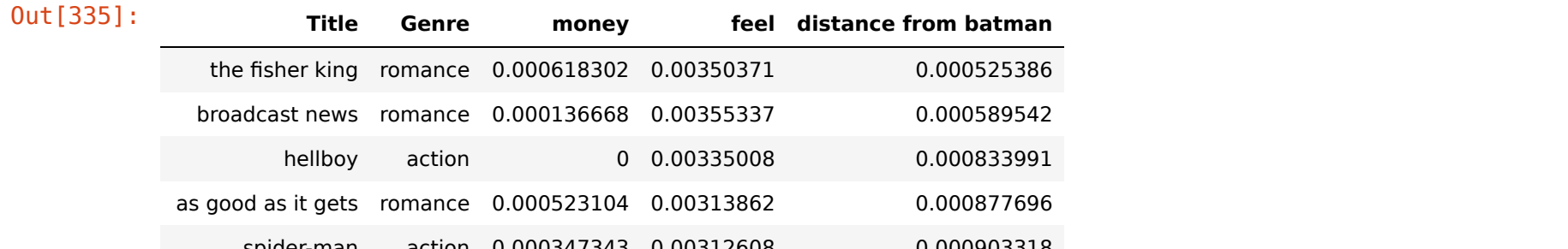
Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

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 [330]:

```
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*.

In [331]:

```
def distance_two_features(title0, title1, x_feature, y_feature):
    """Compute the distance between two movies with titles title0 and title1
    Only the features named x_feature and y_feature are used when computing the distance.
    """
    row0 = row_for_title(title0)
    row1 = row_for_title(title1)
    distance = ((row0.item(x_feature) - row1.item(x_feature))**2 + (row0.item(y_feature) - row1.item(y_feature))**2)**0.5
    return distance

for movie in make_array("the terminator", "the avengers"):
    movie_distance = distance_two_features(movie, "batman returns", "money", "feel")
    print(movie, "distance:", movie_distance)

the terminator distance: 0.0018531387547749904
the avengers distance: 0.0029437356216780243
```

In [332]: ok.grade("q2_1_2");

Running tests

Test summary

Passed: 2
Failed: 0
[oooooooooook] 100.0% passed

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!

In [333]:

```
def distance_from_batman_returns(movie):
    """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")
```

In [334]: ok.grade("q2_1_3");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

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**.

In [335]:

```
# The staff solution took multiple lines.
tbl_with_distance = train_movies.select('title', 'Genre', 'money', 'feel').with_column('distance from batman',
distance_two_features('batman returns', title))
close_movies = tbl_with_distance.take(np.arange(1,8))
close_movies
```

Title	Genre	money	feel	distance from batman
the fisher king	romance	0.000618302	0.00350371	0.000525386
broadcast news	romance	0.000136668	0.00355337	0.000589542
hellboy	action	0	0.00335008	0.000833991
as good as dead	romance	0.000523104	0.00313862	0.000877696
spider-man	action	0.000347343	0.00312608	0.000903318
harold and maude	romance	0	0.00302343	0.0011235
the wedding date	romance	0.00127227	0.00318066	0.00113631

In [336]: ok.grade("q2_1_4");

Running tests

Test summary

Passed: 2
Failed: 0
[oooooooooook] 100.0% passed

Question 2.1.5

Next, we'll classify *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.

In [337]:

```
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.select(label).group(label).sort('count', descending=True).column(label).item(0)

# 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[337]: 'romance'

In [338]: ok.grade("q2_1_5");

Running tests

Test summary

Passed: 2
Failed: 0
[oooooooooook] 100.0% passed

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 [339]:

```
_ = ok.submit()
```

Saving notebook... Saved 'project3.ipynb'.
Submit... 100% complete
Submission successful for user: epergerberkeley.edu
URL: https://okpy.org/cal/data8/fa19/project3/submissions/99v25P

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.

In [340]:

```
def distance(features_array1, features_array2):
    """Compute the Euclidean distance between two arrays of feature values.
    distance = (sum((features_array1 - features_array2)**2))**0.5
    return distance

tbl_without_firstrows = train_movies.drop('Title', 'Genre', 'Year', 'Rating', '# Votes', '# Words')
distance_first_to_second = distance(np.array(tbl_without_firstrows.row(0)), np.array(tbl_without_firstrows.row(1)))
distance_first_to_second
```

Out[340]: 0.0386948978499248

In [341]: ok.grade("q3_1_1");

Running tests

Test summary

Passed: 3
Failed: 0
[oooooooooook] 100.0% passed

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 we lead us directly to 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 3.1.1

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.

- The word is uncommon in both action and romance movies
- The word is common in action movies and uncommon in romance movies
- The word is uncommon in action movies and common in romance movies
- The word is common in both action and romance movies
- 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.

In [342]: bottom_left = 1 #for all questions below are 5 bc plot talks about average, not individual movies.

In [343]: ok.grade("q3_0_1");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

Question 3.0.2

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

In [344]: bottom_right = 3

In [345]: ok.grade("q3_0_2");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

Question 3.0.3

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

In [346]: top_right = 4

In [347]: ok.grade("q3_0_3");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

Question 3.0.4

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

In [348]: top_left = 2

In [349]: ok.grade("q3_0_4");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

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:

- It is an action movie.
- It is a romance movie.

In [350]: movie_genre_guess = 1

In [351]: ok.grade("q3_0_5");

Running tests

Test summary

Passed: 1
Failed: 0
[oooooooooook] 100.0% passed

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.

In [352]:

```
# Set my_20_features to an array of 20 features (strings that are column labels)
my_20_features = make_array('start', 'new', 'hear', 'rememb', 'show', 'same', 'walk', 'next', 'fe', 'understand', 'another', 'car', 'afraid', 'troub', 'came', 'best', 'pay', 'street', 'women', 'y')
my_20_features = make_array('captain', 'ship', 'cop', 'head', 'hour', 'power', 'job', 'weve', 'got', 'angry', 'mom', 'mother', 'hello', 'happy', 'school', 'huh', 'ladi', 'beauti', 'marri')
my_20_features = make_array('marri', 'captain', 'power', 'run', 'world', 'move', 'three', 'nice', 'boy', 'wouldnt', 'enough', 'miss', 'head', 'cours', 'done', 'turn', 'yourself', 'home', 'job', 'knew')
# Select the 20 features of interest from both the train and test datasets
train_20 = train_movies.select(my_20_features)
test_20 = test_movies.select(my_20_features)
```

In [353]: ok.grade("q3_1_1");

Running tests

Test summary

Passed: 5
Failed: 0
[oooooooooook] 100.0% passed

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.

I looked for words that were furthest from the trend line to make the distinction clear on whether they belong in an action or romance film, and I focused on words towards the top right so that the occurrence was high enough for the word to have occurred more than once, making it a better predictor.

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 [354]:

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

Movie:

Title	Genre
the mummy	action

Features:

marri	captain	power	run	world	move	three	nice	boy	wouldnt	enough	
0.000321027	0	0	0	0.000321027	0	0.000642055	0.000321027	0.000321027	0.000321027	0	0.0

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 [355]:

```
# 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 train_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 the features of the movies
    table."
    counts_matrix = np.asmatrix(train_table.columns.transposed())
    diff = np.tile(np.array(test_row), [counts_matrix.shape[0], 1]) - counts_matrix
    sums = np.sum(diff**2, axis=1)
    distances = np.sqrt(sums)
    eps = np.random.uniform(size=distances.shape)*1e-10 #noise for tie break
    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**.

In [356]:

```
# The staff solution took multiple lines of code.
fast_array = fast_distances(test_20.row(0), train_20)
genre_and_distances = Table().with_column('Genre', train_movies.column('Genre')).with_column('Distance', fast_array).sort('Distance')
genre_and_distances
```

Out[356]:

Genre	Distance
romance	0.00169342
romance	0.00172382
romance	0.00184509
action	0.00186555
romance	0.00187958
romance	0.00196601
romance	0.00196691
action	0.00203951
action	0.00205723
...	(195 rows omitted)

In [357]: ok.grade("q3_1_3");

Running tests

Test summary

Passed: 4
Failed: 0
[oooooooooook]


```
In [359]: ok.grade("q3_1_4");
```

Running tests

```
-----
Test summary
  Passed: 2
  Failed: 0
[ooooooooook] 100.0% passed
```

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(8)`).
- 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'`).

```
In [360]: def classify(test_row, train_rows, train_labels, k):
        """Return the most common class among k nearest neighbors to test_row."""
        distances = fast_distances(test_row, train_rows)
        genre_and_distances = Table().with_column('Genre', train_labels).with_column('Distance', distances).sort('Distance').take(np.arange(k))
        classified = genre_and_distances.group('Genre').sort('count', descending = True).column('Genre').item(0)
        return classified
```

```
In [361]: ok.grade("q3_2_1");
```

Running tests

```
-----
Test summary
  Passed: 22
  Failed: 0
[ooooooooook] 100.0% passed
```

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.

```
In [362]: # The staff solution first defined a row called king_kong_features.
king_kong_features = test_movies.where('Title', 'king kong').select(my_20_features).row(0)
king_kong_genre = classify(king_kong_features, train_20, train_movies.column('Genre'), 11)
king_kong_genre

Out[362]: 'action'
```

```
In [363]: ok.grade("q3_2_2");
```

Running tests

```
-----
Test summary
  Passed: 1
  Failed: 0
[ooooooooook] 100.0% passed
```

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.

```
In [364]: def classify_feature_row(row):
        return classify(row, train_20, train_movies.column('Genre'), 11)

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

Out[364]: 'romance'
```

```
In [365]: ok.grade("q3_2_3");
```

Running tests

```
-----
Test summary
  Passed: 11
  Failed: 0
[ooooooooook] 100.0% passed
```

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.

```
In [366]: test_guesses = test_20.apply(classify_feature_row)
proportion_correct = sum(test_guesses == test_movies.column('Genre'))/len(test_guesses)
proportion_correct

Out[366]: 0.7297297297297297
```

```
In [367]: ok.grade("q3_3_1");
```

Running tests

```
-----
Test summary
  Passed: 1
  Failed: 0
[ooooooooook] 100.0% passed
```

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.

```
In [368]: # 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.select('Title', 'Genre').with_column('Was correct', test_guesses == test_movies.column('Genre'))
test_movie_correctness.sort('Was correct', descending = True).show()
```

Title	Genre	Was correct
american outlaws	action	True
solaris	romance	True
logan's run	action	True
48 hrs.	action	True
independence day	action	True
my girl 2	romance	True
the bourne identity	action	True
the apartment	romance	True
rear window	romance	True
an officer and a gentleman	romance	True
blade	action	True
superman iv: the quest for peace	action	True
dune	action	True
the majestic	romance	True
crime spree	action	True
the fantastic four	action	True
mr. deeds goes to town	romance	True
blade ii	action	True
cruel intentions	romance	True
star wars	action	True
starship troopers	action	True
the hudsucker proxy	romance	True
king kong	action	True
sleepless in seattle	romance	True
the bourne supremacy	action	True
witness	romance	True
x-men	action	True
batman returns	action	False
legend	romance	False
the negotiator	action	False
the cider house rules	romance	False
top gun	romance	False
conspiracy theory	action	False
badlands	romance	False
body of evidence	romance	False
men in black	action	False
the mummy	action	False

```
In [369]: ok.grade("q3_3_2");
```

Running tests

```
-----
Test summary
  Passed: 3
  Failed: 0
[ooooooooook] 100.0% passed
```

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.

The movies above that were classified wrong were mostly crime/mystery related movies, and many of them have romantic relationships as part of their plot.

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 [370]: # 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 classify(row, train_20, train_movies.column('Genre'), 12)
test_guesses_1 = test_20.apply(another_classifier)
test_correctness = test_movies.select('Title', 'Genre').with_column('Was Correct', test_movies.column('Genre')==test_guesses_1)
test_correctness.group('Was Correct')
```

```
Out[370]: Was Correct  count
         False         9
         True        28
```

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.

The classifier improved by increasing the `k` number of nearest neighbors.

Question 4.3

Briefly describe what you tried to improve your classifier.

I tried putting in larger and smaller values of `k` until I found the one that resulted in the least number of False predictions

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

```
In [371]: _ = ok.submit()
```

Saving notebook... Saved 'project3.ipynb'.
Submit... 100% complete
Submission successful for user: eperegberkeley.edu
URL: <https://okpy.org/cal/data8/fa19/project3/submissions/k2Womx>

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](#)
- [Train/test split and cross-validation](#)
- [More information about k-nearest neighbors](#)
- [Overfitting](#)

In future data science classes, such as Data Science 100, you'll learn 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](#)
- [Tensorflow information](#)

...and many more!

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