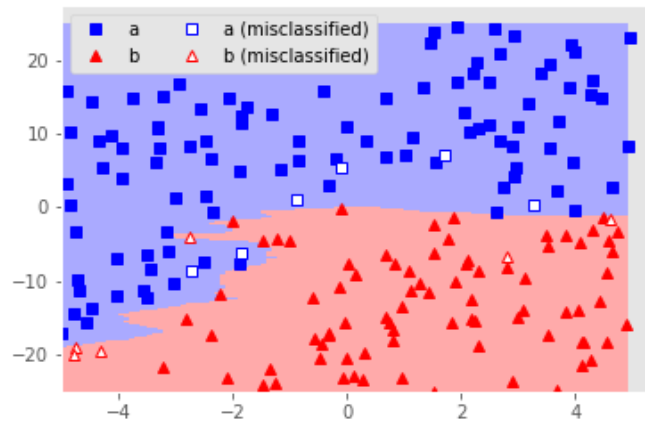


Introduction to Data Science

14 - Classification-kNN



This Lecture

- Classification concepts
- k -Nearest Neighbor Classification

Setup

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import matplotlib.pyplot as plt
        4 import sklearn as sk
        5
        6 from scipy.stats import norm
        7 from pandas import Series, DataFrame
        8 from sklearn.model_selection import train_test_split
        9 from matplotlib.colors import ListedColormap
        10
        11 plt.style.use('ggplot')
        12
        13 %matplotlib inline
```

Recall: Machine (or Statistical) Learning

- Learning structure and relationships in data
- Supervised learning
 - Goal is predicting outputs based on inputs
 - Learn from labeled exemplar data
 - E.g., regression, classification
- Unsupervised learning
 - Goal is reveal hidden structure in the data
 - Inputs but no labeled outputs
 - E.g., clustering

Regression

...we've seen. We'll return to regression in a later lecture to talk about other techniques.

Classification

In classification, we try to predict a **class label**, assigning points to groups.

Our classification training data has the same kinds of inputs as in regression, but the *target* variable is discrete, not continuous.

(Inputs could be categorical too, however. We'll see that eventually.)

E.g., the training data might include features about a borrower's financial history, credit score, owning/renting history, etc. The target might be an indicator on whether or not the borrower ultimately defaulted on their loan.

The learning task is to find a *classifier* function which will let us predict future defaults given the same information about a potential borrower.

(Classification is sometimes regression in disguise. But with discretization at the end.)

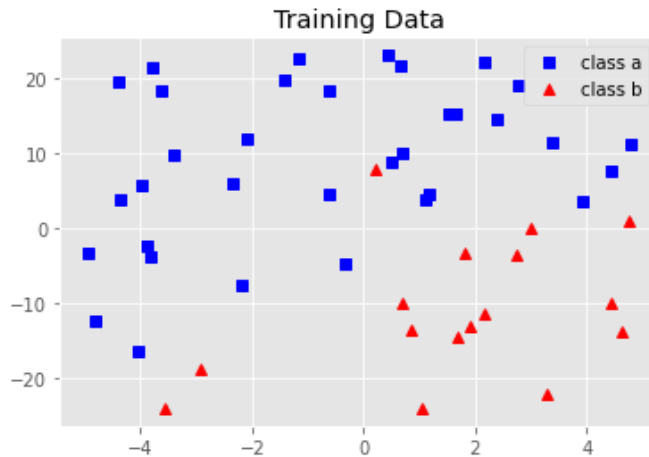
Example Problem

As with regression, we'll start with an artificially generated problem. Inputs will be points in a 2-D space, and targets will be classes simply labeled as 'a' and 'b'.

In [4]:

```
1  # ensure repeatability of this notebook
2  # (comment out for new results each run)
3  np.random.seed(0)
4
5  def f(X):
6      #  $3 + 0.5X - X^2 + 0.15X^3$ 
7      return 3 + 0.5 * X - X**2 + 0.15 * X**3
8
9  # convenience function for generating samples
10 def sample(n, fn, limits, sigma):
11
12     width = limits[1] - limits[0]
13     height = limits[3] - limits[2]
14     x = np.random.random(n) * width + limits[0]
15     y = np.random.random(n) * height + limits[2]
16
17     s = y > fn(x)
18     p = norm.cdf(np.abs(y - fn(x)), scale = sigma) # assigns p with normally distr.
19     r = np.random.random(n) # r is assigned n random variables from (0.0, 1.0).
20
21     def assign(sign, prob, rnum):
22         if sign:
23             if rnum > prob:
24                 return 'b'
25             else:
26                 return 'a'
27         else:
28             if rnum > prob:
29                 return 'a'
30             else:
31                 return 'b'
32
33     c = [assign(s[i], p[i], r[i]) for i in range(n)]
34     return DataFrame({'x' : x, 'y' : y, 'class' : c})
35
```

```
In [5]: 1 data = sample(100, f, [-5, 5, -25, 25], 5)
2 train, test = train_test_split(data, test_size = 0.5)
3 traina = train[train['class']=='a']
4 trainb = train[train['class']=='b']
5 plt.plot(traina['x'], traina['y'], 'bs', label='class a')
6 plt.plot(trainb['x'], trainb['y'], 'r^', label='class b')
7 plt.legend()
8 plt.title('Training Data')
9 plt.show()
```



***k*-Nearest Neighbor Classification**

Suppose we get a new point we haven't seen before, class unknown.

What class do we assign it to?

One of the simplest techniques is *k-nearest neighbor* classification.

Simply look at the *k* nearest points and take a majority vote on the class!

Let's try *k*-nearest neighbor on our example, using scikit-learn:

```
In [6]: 1 from sklearn import neighbors
2
3 k = 5
4
5 model = neighbors.KNeighborsClassifier(k)
6 model.fit(train[['x', 'y']], train['class'])
```

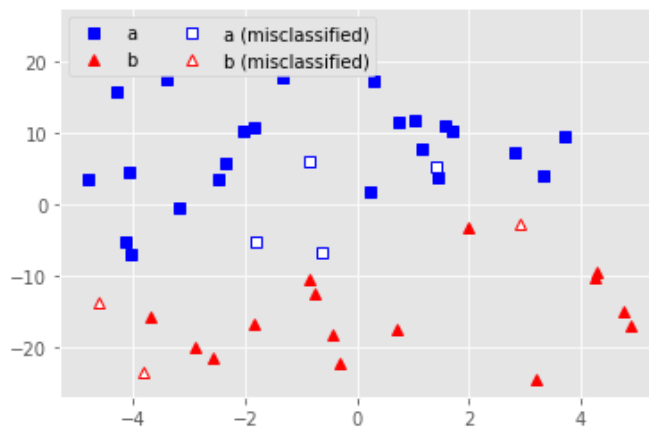
```
Out[6]: KNeighborsClassifier()
```

Plotting the results is a bit tricky. Here are a few approaches.

First, we can apply the trained model to the test data, and plot the resulting points. Correctly classified points will be shown as solid, misclassified points as bordered. Colors and shapes will match the prediction (not the ground truth).

```
In [7]: 1 def plot_predicted_1(model, test):
2         predicted = model.predict(test[['x','y']])
3         correct = test[test['class'] == predicted]
4         correct_a = correct[correct['class'] == 'a']
5         correct_b = correct[correct['class'] == 'b']
6         incorrect = test[test['class'] != predicted]
7         incorrect_a = incorrect[incorrect['class'] == 'a']
8         incorrect_b = incorrect[incorrect['class'] == 'b']
9
10        plt.plot(correct_a['x'], correct_a['y'], 'bs', label='a')
11        plt.plot(correct_b['x'], correct_b['y'], 'r^', label='b')
12        plt.plot(incorrect_a['x'], incorrect_a['y'], 'bs', markerfacecolor='w', label='a')
13        plt.plot(incorrect_b['x'], incorrect_b['y'], 'r^', markerfacecolor='w', label='b')
14        plt.legend(loc='upper left', ncol=2, framealpha=1)
```

```
In [8]: 1 plot_predicted_1(model, test)
```



One problem with this plot is that it doesn't show us the original training data, so it can be difficult to see where the misclassifications came from.

Another plot, harder to generate, shows us the "decision boundary" - the dividing line between classes.

```

In [9]: 1 def plot_predicted_2(model, test):
2         #setup colormap using hex values for blue and red
3         cmap = ListedColormap(['#8888FF', '#FF8888'])
4
5         #set up for the axis ranges we want x values and y values
6         xmin, xmax, ymin, ymax = -5, 5, -25, 25
7
8         #set up the grid for the contour ... create a rectangular grid out of an array
9         #To create a rectangular grid, we need every combination of the x and y points
10        grid_size = 0.1
11        xx, yy = np.meshgrid(np.arange(xmin, xmax, grid_size),
12                             np.arange(ymin, ymax, grid_size))
13
14        #now find y_hat using the prediction function, ravel returns a flattened 1D array
15        pp = model.predict(np.c_[xx.ravel(), yy.ravel()])
16        # pp = model.predict(np.stack((xx.flatten(), yy.flatten()), axis=1))
17
18        #filling an array for the plot function, this comprehension goes through the elements
19        #filled with predicted a or b class labels; the first part is replacing all the
20        zz = np.array([{'a':0, 'b':1}[ab] for ab in pp])
21
22        #the reshape puts the zz array in the same shape as the xx array, again so we can
23        zz = zz.reshape(xx.shape)
24
25        #create a new figure, and plot it using the x, y - zz is the class labels a=0, b=1
26        plt.figure()
27        plt.pcolormesh(xx, yy, zz, cmap = cmap)
28
29        #this calls our other function as an overlay with the actual data points - compare
30        plot_predicted_1(model, test)
31

```

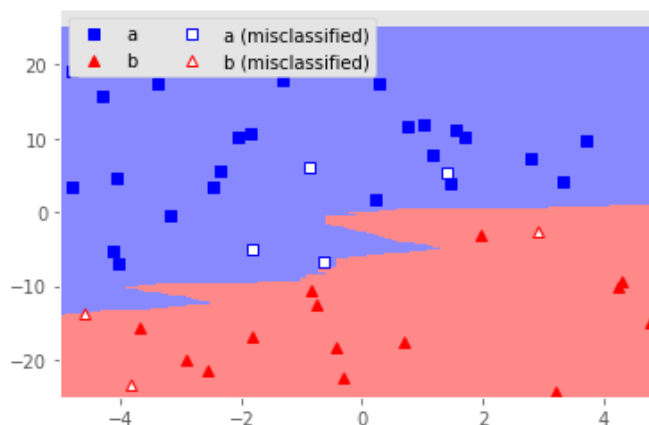
```

In [10]: 1 plot_predicted_2(model, test)
2         plt.show()

```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```

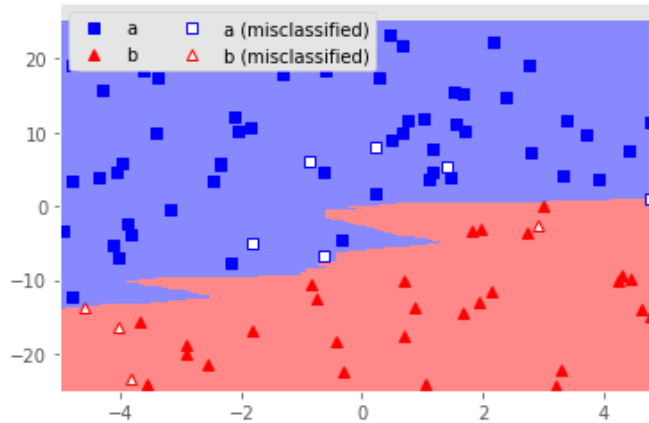


Using the functions we wrote above, we can also plot all of the data, or just the training data:

```
In [11]: 1 plot_predicted_2(model, data)
        2 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



Model Complexity with k -Nearest Neighbors

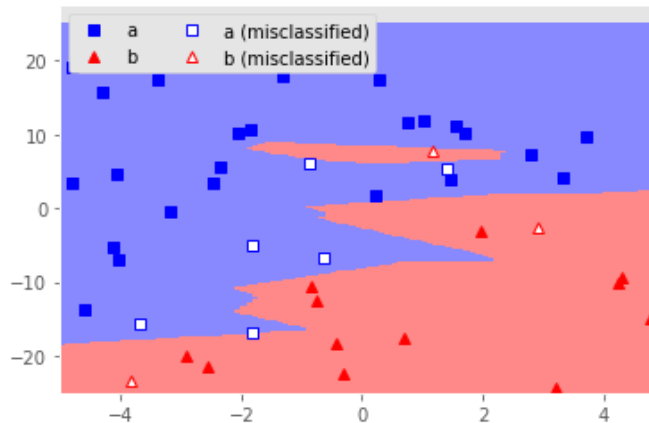
Above we used $k = 5$ to train our model. As we saw, this gave a fairly smooth decision boundary.

What happens if we decrease or increase k ?

```
In [12]: 1 k = 1
2 model = neighbors.KNeighborsClassifier(k)
3 model.fit(train[['x','y']], train['class'])
4 plot_predicted_2(model, test)
5 # ax = plt.gca()
6 # ax.axis([-25, 25, -25, 25])
7 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

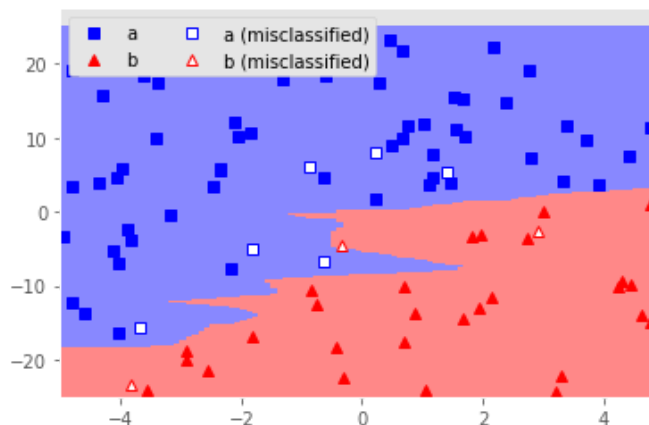
```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



```
In [13]: 1 k = 3
2 model = neighbors.KNeighborsClassifier(k)
3 model.fit(train[['x','y']], train['class'])
4 plot_predicted_2(model, data)
5 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

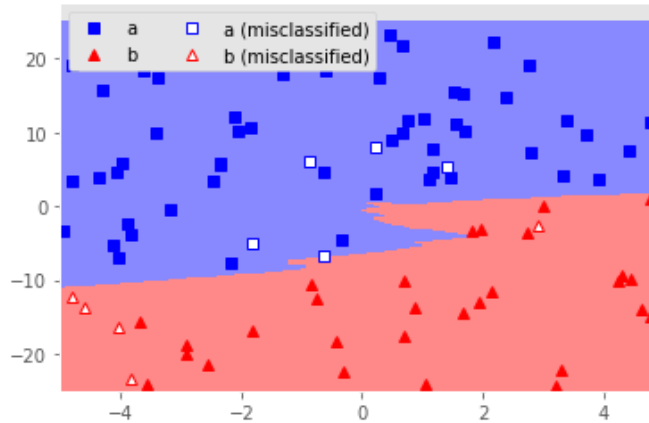
```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```




```
In [14]: 1 k = 7
2 model = neighbors.KNeighborsClassifier(k)
3 model.fit(train[['x','y']], train['class'])
4 plot_predicted_2(model, data)
5 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

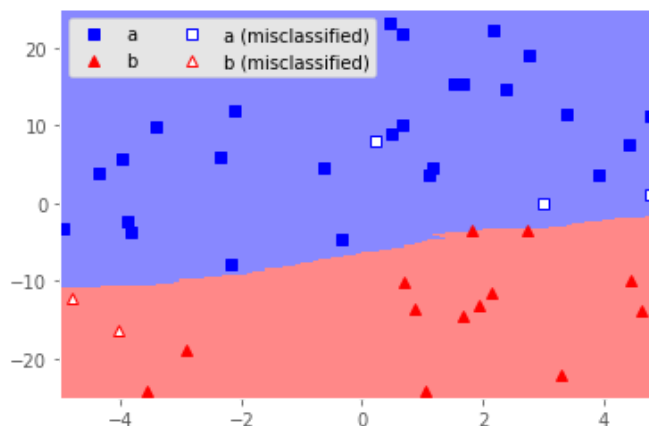
```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



```
In [15]: 1 k = 9
2 model = neighbors.KNeighborsClassifier(k)
3 model.fit(train[['x','y']], train['class'])
4 plot_predicted_2(model, train)
5 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

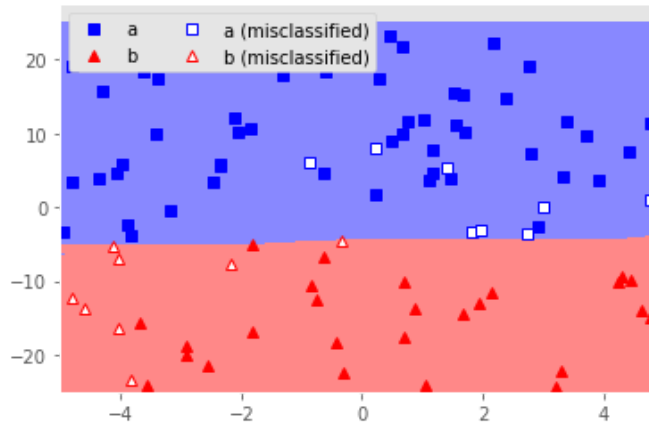
```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



```
In [16]: 1 k = 15
2 model = neighbors.KNeighborsClassifier(k)
3 model.fit(train[['x','y']], train['class'])
4 plot_predicted_2(model, data)
5 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

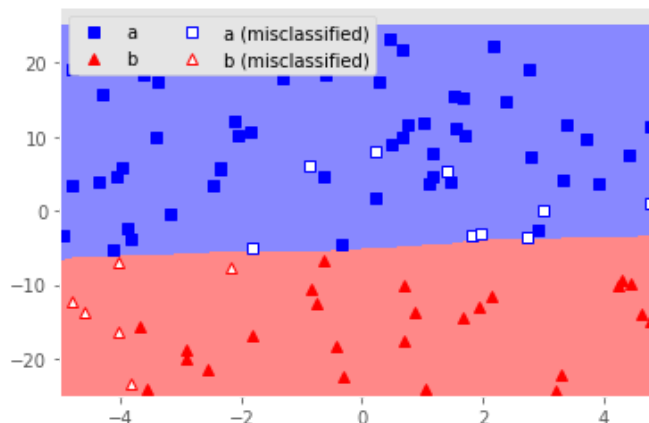
```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



```
In [17]: 1 k = 21
2 model = neighbors.KNeighborsClassifier(k)
3 model.fit(train[['x','y']], train['class'])
4 plot_predicted_2(model, data)
5 plt.show()
```

<ipython-input-9-39e5ef94d95c>:27: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



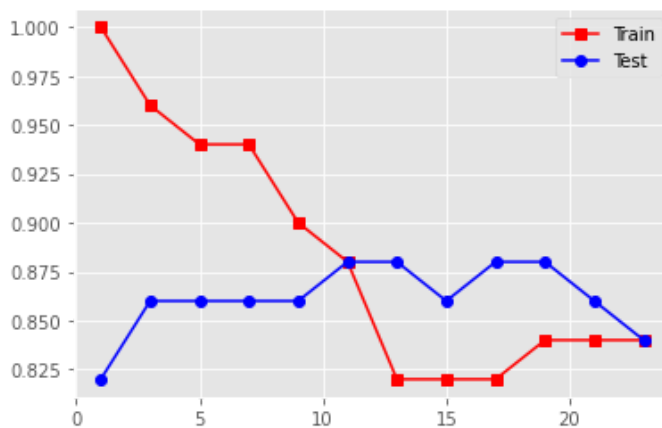
Overfitting/Underfitting

Just like with regression, classification is sensitive to the complexity of your model.

In k -nearest neighbor classification, you overfit by using too few neighbors, and underfit by using too many neighbors.

The scikit-learn nearest neighbor model can generate a score based on the accuracy of our prediction. Let's see our training and test curves:

```
In [18]: 1 kvals = range(1,25,2)
2 train_score = []
3 test_score = []
4
5 for k in kvals:
6     model = neighbors.KNeighborsClassifier(k)
7     model.fit(train[['x','y']], train['class'])
8     train_score.append(model.score(train[['x','y']], train['class']))
9     test_score.append(model.score(test[['x','y']], test['class']))
10
11 plt.plot(kvals, train_score, 'r-s', label='Train')
12 plt.plot(kvals, test_score, 'b-o', label='Test')
13 plt.legend()
14 plt.show()
```



```
In [19]: 1 conf_matrix = sk.metrics.confusion_matrix(test['class'], model.predict(test[['x','y']]))
2 print(conf_matrix)
3 model.score(test[['x','y']], test['class'])
4
```

```
[[28  3]
 [ 5 14]]
```

Out[19]: 0.84

Confusion about the confusion matrix

It's not immediately clear which dimension of the confusion matrix is the true label (class) dimension, and which is the predicted label (class).

Let's figure it out...

The sklearn documentation [says \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html): "By definition a confusion matrix, C , is such that C_{ij} is equal to the number of observations known to be in group i and predicted to be in group j ."

So, this suggests that our `conf_matrix[1,0]` value is the number of 'b' (1) samples that were misclassified as 'a' (0) samples. Let's see if that's correct...

```
In [20]: 1 true_class = test['class']
2 pred_class = model.predict(test[['x','y']])
3
4 # Check what the object types are, so we'll now what code we next...
5 print(type(true_class))
6 print(type(pred_class))
```

```
<class 'pandas.core.series.Series'>
<class 'numpy.ndarray'>
```

```
In [21]: 1 # Because true_class is a Series, not a DataFrame, we do this...
2 bool_true_b = true_class=='b'
3 # instead of something like this...
4 # bool_true_b = true_class['class']=='b'
5
6 bool_true_b_pred_a = pred_class[bool_true_b]=='a'
7
8 print(bool_true_b_pred_a)
9
10 print('')
11 print('sklearn answer: {:d}'.format(conf_matrix[1,0]))
12 print('our answer: {:d}'.format(np.sum(bool_true_b_pred_a)))
```

```
[False  True False  True False False False  True False False  True False
 False  True False False False False False]
```

```
sklearn answer: 5
our answer: 5
```

Next Time

More classification goodness!