

knn_regression

October 20, 2019

```
[1]: import pandas as pd
import numpy as np
import scipy
import matplotlib.pyplot as plt
%matplotlib inline

music = pd.DataFrame()
music['duration'] = [184, 134, 243, 186, 122, 197, 294, 382, 102, 264,
                    205, 110, 307, 110, 397, 153, 190, 192, 210, 403,
                    164, 198, 204, 253, 234, 190, 182, 401, 376, 102]
music['loudness'] = [18, 34, 43, 36, 22, 9, 29, 22, 10, 24,
                    20, 10, 17, 51, 7, 13, 19, 12, 21, 22,
                    16, 18, 4, 23, 34, 19, 14, 11, 37, 42]
music['bpm'] = [105, 90, 78, 75, 120, 110, 80, 100, 105, 60,
               70, 105, 95, 70, 90, 105, 70, 75, 102, 100,
               100, 95, 90, 80, 90, 80, 100, 105, 70, 65]
```

1 KNN Regression

So far we've introduced KNN as a classifier, meaning it assigns observations to categories or assigns probabilities to the various categories. However, KNN is also a reasonable algorithm for regression. It's a simple extension of what we've learned before and just as easy to implement.

1.1 Everything's the Same

Switching KNN to a regression is a simple process. In our previous models, each of the k observations voted for a *category*. As a regression they vote instead for a *value*. Then instead of taking the most popular response, the algorithm averages all of the votes. If you have weights you perform a weighted average.

It's really that simple.

Let's go over a quick example just to confirm your understanding.

Let's stick with the world of music. Instead of trying to classify songs as rock or jazz, let's take the same data with an additional column: beats per minute, or BPM. Can we train our model to predict BPM?

First let's try to predict just in terms of loudness, as this will be easier to represent graphically.

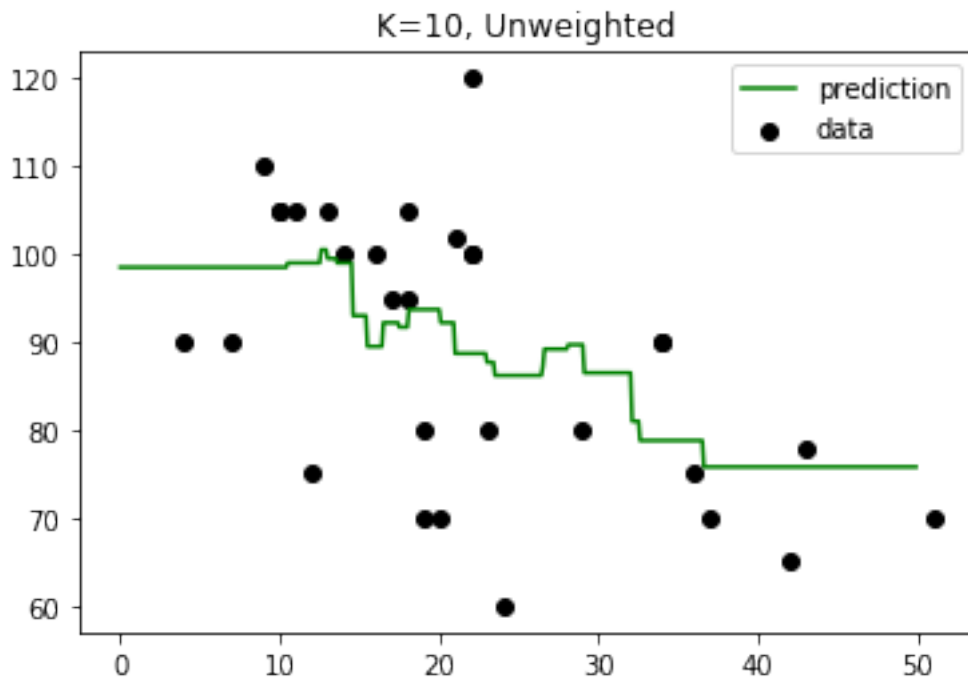
```
[2]: from sklearn import neighbors

# Build our model.
knn = neighbors.KNeighborsRegressor(n_neighbors=10)
X = pd.DataFrame(music.loudness)
Y = music.bpm
knn.fit(X, Y)

# Set up our prediction line.
T = np.arange(0, 50, 0.1)[: , np.newaxis]

# Trailing underscores are a common convention for a prediction.
Y_ = knn.predict(T)

plt.scatter(X, Y, c='k', label='data')
plt.plot(T, Y_, c='g', label='prediction')
plt.legend()
plt.title('K=10, Unweighted')
plt.show()
```



```
[3]: # Run the same model, this time with weights.
knn_w = neighbors.KNeighborsRegressor(n_neighbors=10, weights='distance')
X = pd.DataFrame(music.loudness)
```

```

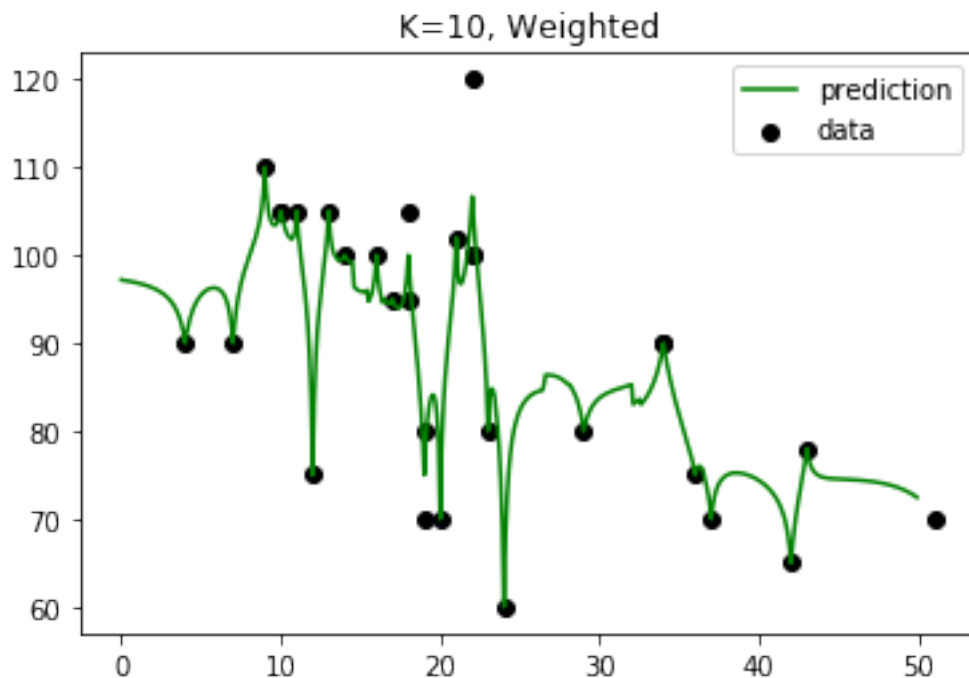
Y = music.bpm
knn_w.fit(X, Y)

# Set up our prediction line.
T = np.arange(0, 50, 0.1)[: , np.newaxis]

Y_ = knn_w.predict(T)

plt.scatter(X, Y, c='k', label='data')
plt.plot(T, Y_, c='g', label='prediction')
plt.legend()
plt.title('K=10, Weighted')
plt.show()

```



Notice how it seems like the weighted model grossly overfits to points. It is interesting that it oscillates around the datapoints. This is because the decay in weight happens so quickly.

1.2 Validating KNN

Now validating KNN, whether a regression or a classifier, is pretty much exactly the same as evaluating other classifiers or regression. Cross validation is still tremendously valuable. You can do holdouts. You even still get an R^2 value for the regression.

Why don't we validate that overfitting of the previous model with some k-fold cross validation? The test statistic given by this model is R^2 , which measures the same as in linear regression.

```
[4]: from sklearn.model_selection import cross_val_score
score = cross_val_score(knn, X, Y, cv=5)
print("Unweighted Accuracy: %0.2f (+/- %0.2f)" % (score.mean(), score.std() *
↳2))
score_w = cross_val_score(knn_w, X, Y, cv=5)
print("Weighted Accuracy: %0.2f (+/- %0.2f)" % (score_w.mean(), score_w.std() *
↳2))
```

Unweighted Accuracy: -0.18 (+/- 0.66)

Weighted Accuracy: 0.11 (+/- 0.94)

First let me state that these two models are fantastically awful. There doesn't seem to be much of a relationship. It's all very poor. However the increased variance in the weighted model is interesting.

Why don't you add the other feature and mess around with k and weighting to see if you can do any better than we've done so far?

```
[5]: ## Your model here.

# Run the same model, this time with weights.
knn_w = neighbors.KNeighborsRegressor(n_neighbors=10, weights='uniform')
X = music[['loudness', 'duration']]
Y = music.bpm
knn_w.fit(X, Y)

# Set up our prediction line.
T = music[['loudness', 'duration']]
T['loudness'] = [np.random.
↳randint(min(music['loudness']),max(music['loudness']+1)) for i in
↳T['loudness']]
T['duration'] = [np.random.
↳randint(min(music['duration']),max(music['duration']+1)) for i in
↳T['duration']]

Y_ = knn_w.predict(T)

# score = cross_val_score(knn, X, Y, cv=5)
print("Unweighted Accuracy: %0.2f (+/- %0.2f)" % (score.mean(), score.std() *
↳2))
score_w = cross_val_score(knn_w, X, Y, cv=5)
print("Weighted Accuracy: %0.2f (+/- %0.2f)" % (score_w.mean(), score_w.std() *
↳2))
```