

## KNN regression experiments

In class we learned about how KNN regression works, and tips for using KNN. For example, we learned that data should be scaled when using KNN, and that extra, useless predictors should not be used with KNN. Are these tips really correct?

In this notebook we run a bunch of tests to see how KNN is affected by the choice of  $k$ , scaling of the predictors, presence of useless predictors, and other things.

One experiment we do not run, and which would be interesting, is to see how KNN performance changes as a function of the size of the training set.

## INSTRUCTIONS

Enter code wherever you see # YOUR CODE HERE in code cells, or YOU TEXT HERE in markup cells.

```
In [85]: import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsRegressor
import matplotlib.pyplot as plt
```

```
In [86]: # set default figure size
plt.rcParams['figure.figsize'] = [8.0, 6.0]
```

```
In [87]: # code in this cell from:
# https://stackoverflow.com/questions/27934885/how-to-hide-code-from-cells-in-ipython-notebook-visualized-with-nbviewer
from IPython.display import HTML

HTML('''<script>
code_show=true;
function code_toggle() {
  if (code_show){
    $('div.input').hide();
  } else {
    $('div.input').show();
  }
  code_show = !code_show
}
$( document ).ready(code_toggle);
</script>
<form action="javascript:code_toggle()"><input type="submit" value="Click here to display/hide the code."></form>''')
```

Out[87]:

## Read the data and take a first look at it

The housing dataset is good for testing KNN because it has many numeric features. See Aurélien Géron's book titled 'Hands-On Machine learning with Scikit-Learn and TensorFlow' for information on the dataset.

```
In [88]: df = pd.read_csv("https://raw.githubusercontent.com/grbruns/cst383/master/housing.csv")
```

```
In [89]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   longitude              20640 non-null  float64
1   latitude               20640 non-null  float64
2   housing_median_age     20640 non-null  float64
3   total_rooms            20640 non-null  float64
4   total_bedrooms        20433 non-null  float64
5   population             20640 non-null  float64
6   households             20640 non-null  float64
7   median_income          20640 non-null  float64
8   median_house_value     20640 non-null  float64
9   ocean_proximity        20640 non-null  object
dtypes: float64(9), object(1)
memory usage: 1.6+ MB
```

Note that numeric features have different ranges. For example, the mean value of 'total\_rooms' is over 2,500, while the mean value of 'median\_income' is about 4. 'median\_house\_value' has a much greater mean value, over \$200,000, but we will be using it as the target variable.

```
In [90]: from IPython.display import Image
from IPython.core.display import HTML
Image(url= "data:image/jpeg;base64,/9j/4AAQSkZJRgABAQAAQABAAQ/2wCEAAoHCBIVFRgVFRYYGBgYGBgZGBgYGBgaGBoaGBgzGRgYGBgcIS41
#Flickr source for "Houses going down" : https://www.flickr.com/photos/59937401@N07/5474464467
```



```
In [91]: df.describe()
```

Out[91]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value
count	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.000000	20640.000000
mean	-119.569704	35.631861	28.639486	2635.763081	537.870553	1425.476744	499.539680	3.870671	206855.816909
std	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.899822	115395.615874
min	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.499900	14999.000000
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.563400	119600.000000
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.534800	179700.000000
75%	-118.010000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.743250	264725.000000
max	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.000100	500001.000000

### Missing Data

Notice that 207 houses are missing their *total\_bedroom* info:

```
In [92]: print(df.isnull().sum())
df[df['total_bedrooms'].isnull()]

longitude      0
latitude       0
housing_median_age  0
total_rooms    0
total_bedrooms 207
population     0
households     0
median_income  0
median_house_value  0
ocean_proximity  0
dtype: int64
```

Out[92]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
290	-122.16	37.77	47.0	1256.0	NaN	570.0	218.0	4.3750	161900.0	NEAR BAY
341	-122.17	37.75	38.0	992.0	NaN	732.0	259.0	1.6196	85100.0	NEAR BAY
538	-122.28	37.78	29.0	5154.0	NaN	3741.0	1273.0	2.5762	173400.0	NEAR BAY
563	-122.24	37.75	45.0	891.0	NaN	384.0	146.0	4.9489	247100.0	NEAR BAY
696	-122.10	37.69	41.0	746.0	NaN	387.0	161.0	3.9063	178400.0	NEAR BAY
...	...	...	...	...	...	...	...	...	...	...
20267	-119.19	34.20	18.0	3620.0	NaN	3171.0	779.0	3.3409	220500.0	NEAR OCEAN
20268	-119.18	34.19	19.0	2393.0	NaN	1938.0	762.0	1.6953	167400.0	NEAR OCEAN
20372	-118.88	34.17	15.0	4260.0	NaN	1701.0	669.0	5.1033	410700.0	<1H OCEAN
20460	-118.75	34.29	17.0	5512.0	NaN	2734.0	814.0	6.6073	258100.0	<1H OCEAN
20484	-118.72	34.28	17.0	3051.0	NaN	1705.0	495.0	5.7376	218600.0	<1H OCEAN

207 rows × 10 columns

Let's drop these instances for now

```
In [93]: df = df.dropna()
```

## Prepare data for machine learning

We will use KNN regression to predict the price of a house from its features, such as size, age and location.

We use a subset of the data set for our training and test data. Note that we keep an unscaled version of the data for one of the experiments we will run.

```
In [94]: # for repeatability
np.random.seed(42)
```

```
In [95]: # select the predictor variables and target variables to be used with regression
predictors = ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households',
#dropping categorical features, such as ocean_proximity, including spatial ones such as long/lat.
target = 'median_house_value'
X = df[predictors].values
y = df[target].values
```

```
In [96]: # KNN can be slow, so get a random sample of the full data set
indexes = np.random.choice(y.size, size=10000)
X_mini = X[indexes]
y_mini = y[indexes]
```

```
In [97]: # Split the data into training and test sets, and scale
scaler = StandardScaler()

# unscaled version (note that scaling is only used on predictor variables)
X_train_raw, X_test_raw, y_train, y_test = train_test_split(X_mini, y_mini, test_size=0.30, random_state=42)

# scaled version
X_train = scaler.fit_transform(X_train_raw)
X_test = scaler.transform(X_test_raw)
```

```
In [98]: # sanity check
print(X_train.shape)
print(X_train[:3])

(7000, 8)
[[ 1.22783551 -1.3492796   0.34639424 -0.16627017  0.11697691 -0.15874461
   0.18687025 -0.74984935]
 [ 0.62095726 -0.82169566  0.58720859 -0.11584049 -0.22077651 -0.0770853
  -0.14171346  1.12877289]
 [-1.16983102  0.7563873  -0.45632025 -0.32112946  0.02736886 -0.37395092
  -0.04890738 -0.10303138]]
```

## Baseline performance

For regression problems, our baseline is the "blind" prediction that is just the average value of the target variable. The blind prediction must be calculated using the training data. Calculate and print the test set root mean squared error (test RMSE) using this blind prediction. I have provided a function you can use for RMSE.

```
In [99]: def rmse(predicted, actual):
        return np.sqrt(((predicted - actual)**2).mean())
```

```
In [100]: rmse(y_test.mean(), y_test)
```

```
Out[100]: 112886.96075401208
```

## Performance with default hyperparameters

Using the training set, train a KNN regression model using the ScikitLearn KNeighborsRegressor, and report on the test RMSE. The test RMSE is the RMSE computed using the test data set.

When using the KNN algorithm, use `algorithm='brute'` to get the basic KNN algorithm.

```
In [101]: knn = KNeighborsRegressor(n_neighbors=5, algorithm="brute")
knn.fit(X_train, y_train)

predictions = knn.predict(X_test)
round(rmse(predictions, y_test),1)
```

```
Out[101]: 62448.9
```

## Impact of K

In class we discussed the relationship of the hyperparameter  $k$  to overfitting.

I provided code to test KNN on  $k=1$ ,  $k=3$ ,  $k=5$ , ...,  $k=29$ . For each value of  $k$ , compute the training RMSE and test RMSE. The training RMSE is the RMSE computed using the training data. Use the 'brute' algorithm, and Euclidean distance, which is the default. You need to add the `get_train_test_rmse()` function.

In [ ]:

```
In [102]: def get_train_test_rmse(regr, X_train, X_test, y_train, y_test):
    regr.fit(X_train, y_train)
    predictionsTest = regr.predict(X_test)
    predictionsTrain = regr.predict(X_train)
    return round(rmse(predictionsTrain, y_train),1), round(rmse(predictionsTest, y_test),1)
```

```
In [103]: n = 30
test_rmse = []
train_rmse = []
ks = np.arange(1, n+1, 2)
for k in ks:
    print(k, ' ', end='')
    regr = KNeighborsRegressor(n_neighbors=k, algorithm='brute')
    rmse_tr, rmse_te = get_train_test_rmse(regr, X_train, X_test, y_train, y_test)
    train_rmse.append(rmse_tr)
    test_rmse.append(rmse_te)
print('done')

1  3  5  7  9  11  13  15  17  19  21  23  25  27  29  done
```

```
In [104]: # sanity check
print('Test RMSE when k = 3: {:.1f}'.format(np.array(test_rmse)[ks==3][0]))
```

Test RMSE when k = 3: 64167.1

Using the training and test RMSE values you got for each value of  $k$ , find the  $k$  associated with the lowest test RMSE value. Print this  $k$  value and the associated lowest test RMSE value. In other words, if you found that  $k=11$  gave the lowest test RMSE, then print the value 11 and the test RMSE value obtained when  $k=11$ .

```
In [105]: def get_best(ks, rmse):
    best = np.sort(rmse)[0]
    i = ks[np.where(rmse==best)]
    return i, best

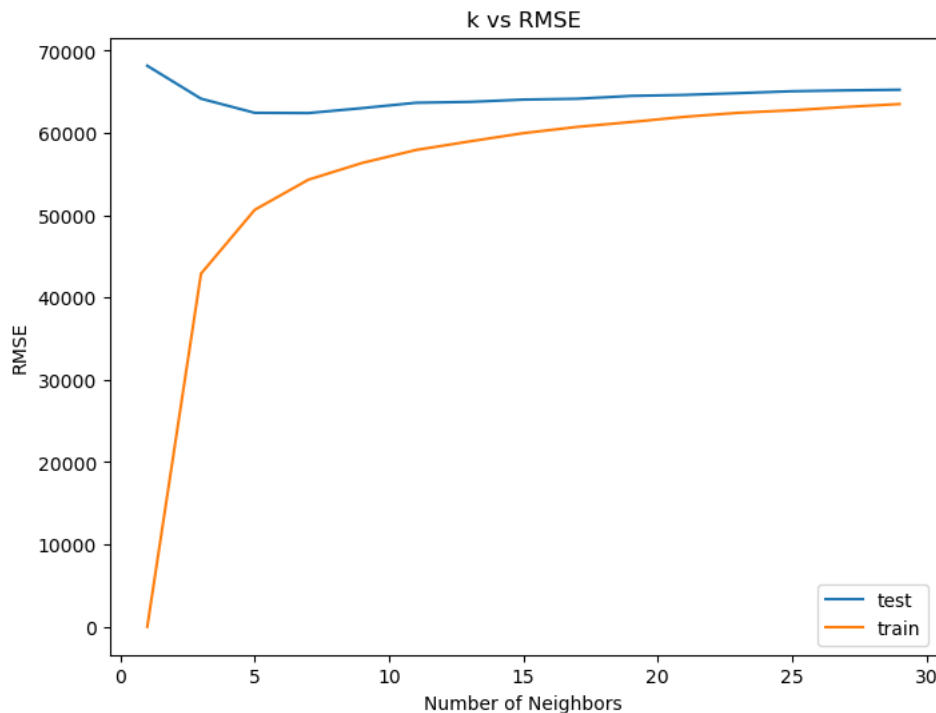
best_k, best_rmse = get_best(ks, test_rmse)
print('best k = {}, best test RMSE: {:.1f}'.format(best_k, best_rmse))
```

best k = [7], best test RMSE: 62421.5

Plot the test and training RMSE as a function of  $k$ , for all the  $k$  values you tried.

```
In [106]: plt.plot(ks, test_rmse, label='test')
plt.plot(ks, train_rmse, label='train')
plt.title('k vs RMSE')
plt.legend()
plt.xlabel('Number of Neighbors')
plt.ylabel('RMSE')
```

Out[106]: Text(0, 0.5, 'RMSE')



## Comments

*In the markup cell below, write about what you learned from your plot. I would expect two or three sentences, but what's most important is that you write something thoughtful.*

I think it's interesting that the RMSE start on opposite ends (high and low) and then start to converge as number of neighbors increase. The behavior of the test data is what we expect. At a low number of neighbors we are underfitting and outliers are influencing our predictions. The RMSE drops as we increase k and we can see when we start to overfit when RMSE starts to increase again. The behavior of the training data did not seem as intuitive to me. After thinking about it, I think it's because when we underfit we basically have exact points of the  $y_{train}$  data and as we increase k the areas start to blend thus increasing RMSE.

## Impact of noise predictors

*In class we heard that the KNN performance goes down if useless "noisy predictors" are present. These are predictors that don't help in making predictions. In this section, run KNN regression by adding one noise predictor to the data, then 2 noise predictors, then three, and then four. For each, compute the training and test RMSE. In every case, use  $k=10$  as the k value and use the default Euclidean distance as the distance function.*

*The `add_noise_predictor()` method makes it easy to add a predictor variable of random values to  $X_{train}$  or  $X_{test}$ .*

```
In [107]: def add_noise_predictor(X):
          """ add a column of random values to 2D array X """
          noise = np.random.normal(size=(X.shape[0], 1))
          return np.hstack((X, noise))
```

*Hint: In each iteration of your loop, add a noisy predictor to both  $X_{train}$  and  $X_{test}$ . You don't need to worry about rescaling the data, as the new noisy predictor is already scaled. Don't modify  $X_{train}$  and  $X_{test}$  however, as you will be using them again.*

```

In [108]: # YOUR CODE HERE
#copy xtrain and xtest
Xtr = X_train
Xts = X_test
#knn object
knn = KNeighborsRegressor(n_neighbors=10)
#result lists
noises = [0]
rmse = []
#baseline
knn.fit(Xtr, y_train)
rmse.append(round(rmse(knn.predict(Xts), y_test),1))
for i in range(8):
    #update noises nnumber list
    noises.append(i+1)
    #add noise columns
    Xtr = add_noise_predictor(Xtr)
    Xts = add_noise_predictor(Xts)
    #Train it
    knn.fit(Xtr, y_train)
    #append rmse
    rmse.append(round(rmse(knn.predict(Xts), y_test),1))

```

Plot the percent increase in test RMSE as a function of the number of noise predictors. The x axis will range from 0 to 4. The y axis will show a percent increase in test RMSE.

To compute percent increase in RMSE for  $n$  noise predictors, compute  $100 * (rmse - base\_rmse) / base\_rmse$ , where  $base\_rmse$  is the test RMSE with no noise predictors, and  $rmse$  is the test RMSE when  $n$  noise predictors have been added.

```

In [109]: # np array
y = np.array(rmse)

```

```

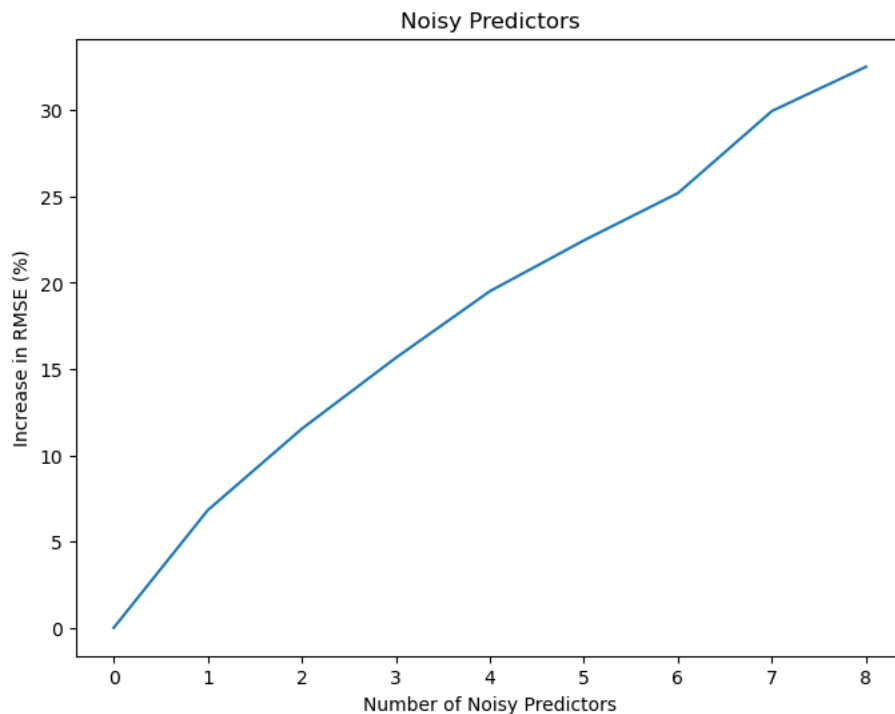
In [110]: #plt
plt.plot(noises, 100*(y - y[0])/y[0])
plt.title('Noisy Predictors')
plt.ylabel('Increase in RMSE (%)')
plt.xlabel('Number of Noisy Predictors')

```

```

Out[110]: Text(0.5, 0, 'Number of Noisy Predictors')

```



## Comments

Look at the results you obtained and add some thoughtful commentary.

As expected adding noisy predictors increases the RMSE. I thought it looked like the RMSE would converge if we extended x which I thought was suprising so I increased the iterations of noisy predictors and it seems like RMSE is increasing linearly which is more inline with what I would have predicted.

## Impact of scaling

In class we learned that we should scaled the training data before using KNN. How important is scaling with KNN? Repeat the experiments you ran before (like in the impact of distance metric section), but this time use unscaled data.

Run KNN as before but use the unscaled version of the data. You will vary  $k$  as before. Use `algorithm='brute'` and Euclidean distance.

```
In [111]: n = 30
test_rmse = []
train_rmse = []
ks = np.arange(1, n+1, 2)
for k in ks:
    print(k, ' ', end='')
    regr = KNeighborsRegressor(n_neighbors=k, algorithm='brute')
    rmse_tr, rmse_te = get_train_test_rmse(regr, X_train_raw, X_test_raw, y_train, y_test)
    train_rmse.append(rmse_tr)
    test_rmse.append(rmse_te)
print('done')
```

1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 done

Print the best  $k$  and the test RMSE associated with the best  $k$ .

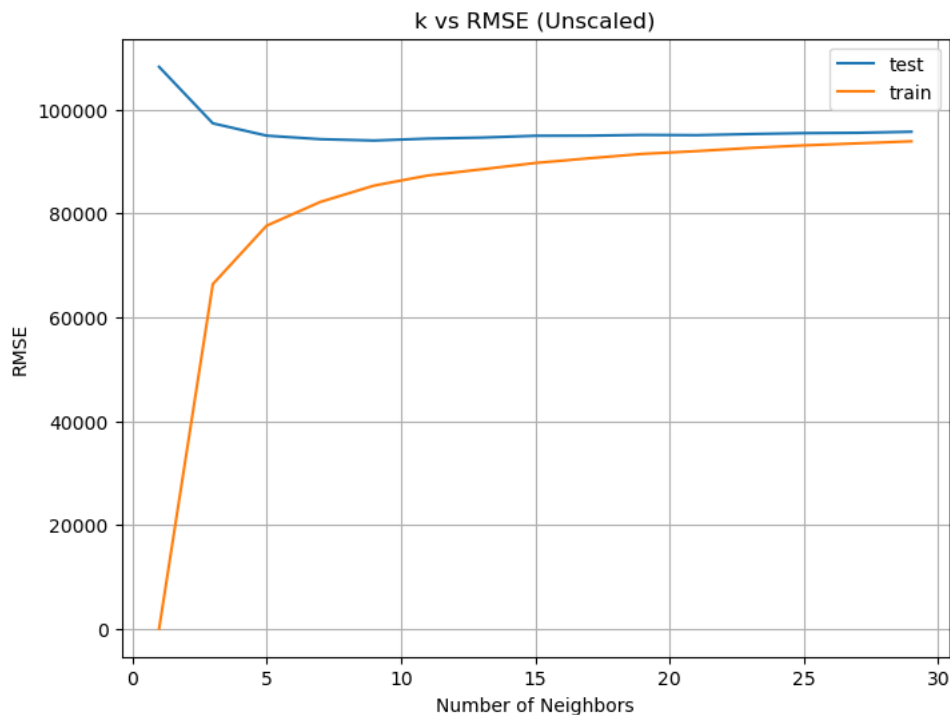
```
In [112]: best_k, best_rmse = get_best(ks, test_rmse)
print('best k = {}, best test RMSE: {:.1f}'.format(best_k, best_rmse))
```

best k = [9], best test RMSE: 94057.4

Plot training and test RMSE as a function of  $k$ . Your plot title should note the use of unscaled data.

```
In [113]: plt.plot(ks, test_rmse, label='test')
plt.plot(ks, train_rmse, label='train')
plt.title('k vs RMSE (Unscaled)')
plt.legend()
plt.grid()
plt.xlabel('Number of Neighbors')
plt.ylabel('RMSE')
```

Out[113]: Text(0, 0.5, 'RMSE')



## Comments

Reflect on what happened and provide some short commentary, as in previous sections.

The plots look similar to the unscaled data because the y axis gets scaled in the plot, but the RMSE is higher and the rate of change increases when the data is unscaled. When the training data is not scaled prediction error increases because resolution is reduced.

## Impact of algorithm

We didn't discuss in class that there are variants of the KNN algorithm. The main purpose of the variants is to be faster and to reduce that amount of training data that needs to be stored.

Run experiments where you test each of the three KNN algorithms supported by Scikit-Learn: *ball\_tree*, *kd\_tree*, and *brute*. In each case, use  $k=10$  and use Euclidean distance.

```
In [114]: def doTheThing(algo, X_train, X_test, y_train, y_test):
    n = 10
    test_rmse = []
    train_rmse = []
    ks = np.arange(1, n+1, 2)
    for k in ks:
        print(k, ' ', end='')
        regr = KNeighborsRegressor(n_neighbors=k, algorithm=algo)
        rmse_tr, rmse_te = get_train_test_rmse(regr, X_train, X_test, y_train, y_test)
        train_rmse.append(rmse_tr)
        test_rmse.append(rmse_te)
    print(algo + ' done')
    return np.array(test_rmse) , ks

#run tests
results = {}
results["ball tree"], ks = doTheThing('ball_tree', X_train, X_test, y_train, y_test)
results["kd tree"], ks = doTheThing('kd_tree', X_train, X_test, y_train, y_test)
results["brute"], ks = doTheThing('brute', X_train, X_test, y_train, y_test)

#print results and make rmselist
besto={}
bestk={}
for test in results:
    best_k, best_rmse = get_best(ks,results[test])
    print('best k = {}, best {} test RMSE: {:.1f}'.format(best_k, test, best_rmse))
    besto[test] = best_rmse
    bestk[test] = best_k

1 3 5 7 9 ball_tree done
1 3 5 7 9 kd_tree done
1 3 5 7 9 brute done
best k = [7], best ball tree test RMSE: 62421.5
best k = [7], best kd tree test RMSE: 62421.5
best k = [7], best brute test RMSE: 62421.5
```

Print the name of the best algorithm, and the test RMSE achieved with the best algorithm.

```
In [115]: lowest = np.min(np.array(list(besto.values())))
bestest = []
for test in besto:
    if besto[test] == lowest:
        bestest.append(test)
print('BEST ALGORITHM(s) with RMSE of {:.1f}:'.format(lowest))
for best in bestest:
    print('{} at k={}'.format(best, bestk[best]))

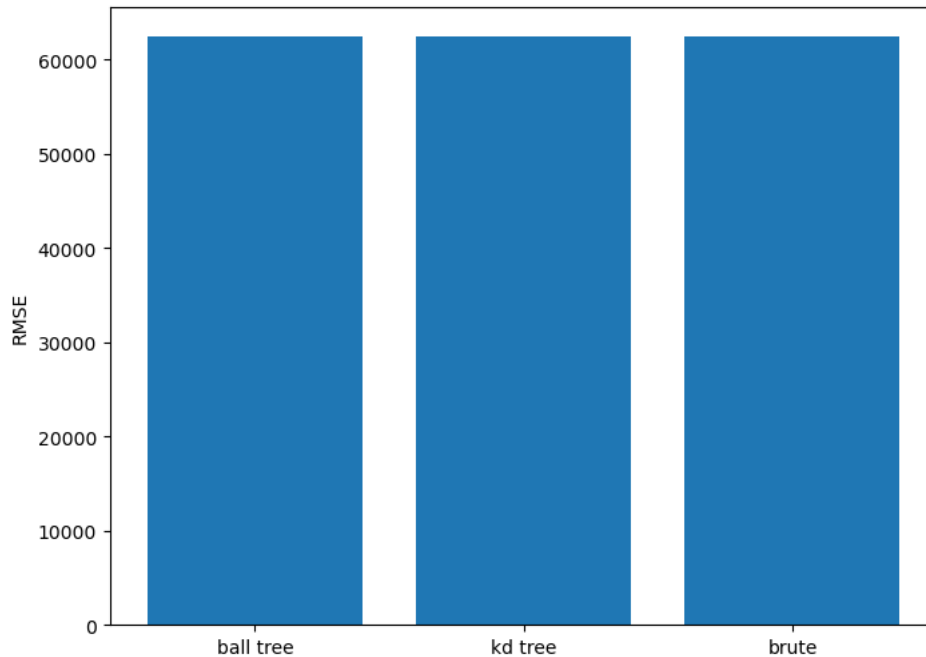
BEST ALGORITHM(s) with RMSE of 62421.5:
ball tree at k=[7]
kd tree at k=[7]
brute at k=[7]
```

Plot the test RMSE for each of the three algorithms as a bar plot.



```
In [116]: plt.bar(besto.keys(), besto.values())
plt.ylabel('RMSE')
```

```
Out[116]: Text(0, 0.5, 'RMSE')
```



## Comments

As usual, reflect on the results and add comments.

I don't think I learned anything from this section. All of the algorithms had the same results which, if I didn't have an example to reference, I would have made me think I did something wrong.

## Impact of weighting

It was briefly mentioned in lecture that there is a variant of KNN in which training points are given more weight when they are closer to the point for which a prediction is to be made. The 'weight' parameter of `KNeighborsRegressor()` has two possible values: 'uniform' and 'distance'. Uniform is the basic algorithm.

Run an experiment similar to the previous one. Compute the test RMSE for uniform and distance weighting. Using  $k = 10$ , the brute algorithm, and Euclidean distance.

```
In [117]: def weighting(w, X_train, X_test, y_train, y_test):
n = 10
test_rmse = []
train_rmse = []
ks = np.arange(1, n+1, 2)
for k in ks:
    print(k, ' ', end='')
    regr = KNeighborsRegressor(n_neighbors=k, algorithm='brute', weights=w)
    rmse_tr, rmse_te = get_train_test_rmse(regr, X_train, X_test, y_train, y_test)
    train_rmse.append(rmse_tr)
    test_rmse.append(rmse_te)
    print('{} done'.format(w))
return np.array(test_rmse) , ks

tests=['uniform', 'distance']
results={}
for test in tests:
    results[test], ks = weighting(test, X_train, X_test, y_train, y_test)
```

```
1 3 5 7 9 uniform done
1 3 5 7 9 distance done
```

Print the weighting the gave the lowest test RMSE, and the test RMSE it achieved.

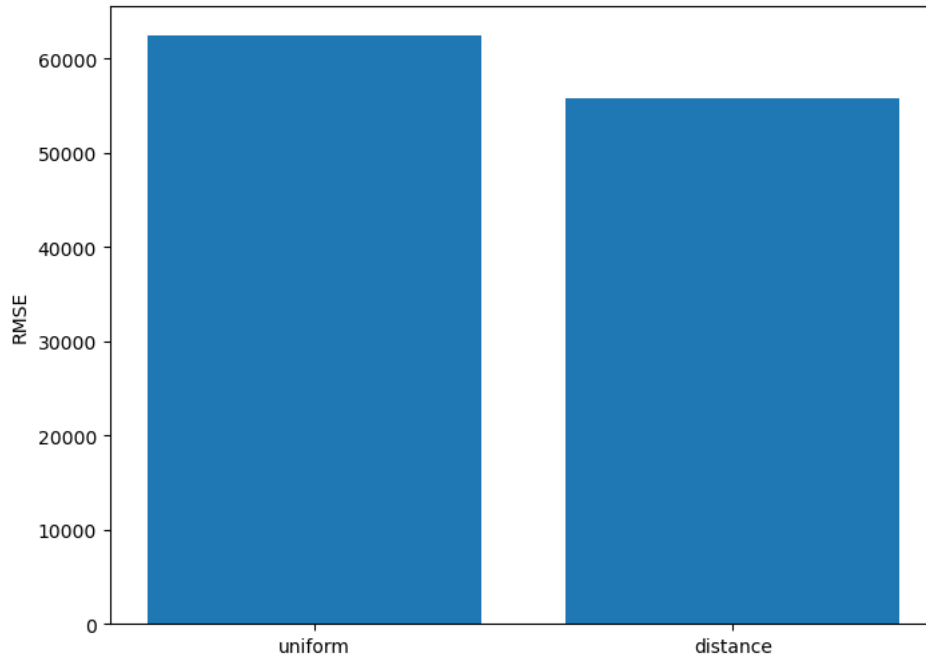
```
In [118]: best={}
          for test in results:
              best_k, best[test] = get_best(ks,results[test])
              print('best k = {}, best {} test RMSE: {:.1f}'.format(best_k,test, best[test]))

best k = [7], best uniform test RMSE: 62421.5
best k = [7], best distance test RMSE: 55800.7
```

Create a bar plot showing the test RMSE for the uniform and distance weighting options.

```
In [119]: plt.bar(best.keys(),best.values())
          plt.ylabel('RMSE')
```

```
Out[119]: Text(0, 0.5, 'RMSE')
```



## Comments

*As usual, reflect and comment.*

Based on the results from this experiment, applying weights to datapoints based on their distance from the test point improves accuracy on this data set. I would think this assumption applies to most knn datasets but I'm not sure.

## Conclusions

*Please provide at least a few sentences of commentary on the main things you've learned from the experiments you've run.*

I found these experiments to be interesting. I learned more about how to prep data for, train, and use a knn regression model. The main thing I gained from these exercise is a better understanding on how knn objects from the sklearn library work and how to test and quantify the models effectiveness.