2]:	<pre>import numpy as np import pandas as pd from sklearn.model_selection import StandardScaler from sklearn.neighbors import KNeighborsRegressor import matplotlib.pyplot as plt # set default figure size plt.rcParams['figure.figsize'] = [8.0, 6.0]</pre>
3]:	<pre># set default figure size plt.rcParams['figure.figsize'] = [8.0, 6.0] # code in this cell from: # https://stackoverflow.com/questions/27934885/how-to-hide-code-from-cells-in-ipython-notebook-visualized from IPython.display import HTML HTML('''<script> code_show=true; function code_toggle() { if (code_show) { \$('div.input').hide(); } }</pre></th></tr><tr><th>3]: </th><th><pre>} else { \$('div.input').show(); } code_show = !code_show } \$(document).ready(code_toggle); </script> <form action="javascript:code_toggle()"><input click="" code.<="" display="" here="" hide="" pre="" the="" to="" type="submit" value="Click here to display/hide the code."/></form></pre>
	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 Machin learning with Scikit-Learn and TensorFlow' for information on the dataset. df = pd.read_csv("https://raw.githubusercontent.com/grbruns/cst383/master/housing.csv") df.info() <class 'pandas.core.frame.dataframe'=""></class>
	RangeIndex: 20640 entries, 0 to 20639 Data columns (total 10 columns): # Column Non-Null Count Dtype
	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. from IPython.display import Image from IPython.core.display import HTML
5]:	<pre>Image(url= " #Flickr source for "Houses going down" : https://www.flickr.com/photos/59937401@N07/5474464467</pre>
7]:	df.describe() longitude latitude housing_median_age total_rooms total_bedrooms population households median_income median_income median_income total_rooms 20640.000000 20640.
	std 2.003532 2.135952 12.585558 2181.615252 421.385070 1132.462122 382.329753 1.899822 min -124.350000 32.540000 1.000000 2.000000 1.000000 3.000000 1.000000 0.499900 25% -121.800000 33.930000 18.000000 1447.750000 296.000000 787.000000 280.000000 2.563400 50% -118.490000 34.260000 29.000000 2127.000000 435.000000 1166.000000 409.000000 3.534800 75% -118.010000 37.710000 37.000000 3148.000000 647.000000 1725.000000 605.00000 4.743250 max -114.310000 41.950000 52.000000 39320.000000 6445.000000 35682.000000 6082.000000 15.000100
	Missing Data Notice that 207 houses are missing their total_bedroom info: print(df.isnull().sum()) df[df['total_bedrooms'].isnull()] longitude
3]:	total_rooms 0 total_bedrooms 207 population 0 households 0 median_income 0 median_house_value 0 ocean_proximity 0 dtype: int64 longitude latitude housing_median_age total_rooms total_bedrooms population households median_income median_house_value 0 0
	341 -122.17 37.75 38.0 992.0 NaN 732.0 259.0 1.6196 857 538 -122.28 37.78 29.0 5154.0 NaN 3741.0 1273.0 2.5762 1734 563 -122.24 37.75 45.0 891.0 NaN 384.0 146.0 4.9489 247 696 -122.10 37.69 41.0 746.0 NaN 387.0 161.0 3.9063 1784
	20372 -118.88 34.17 15.0 4260.0 NaN 1701.0 669.0 5.1033 4107 20460 -118.75 34.29 17.0 5512.0 NaN 2734.0 814.0 6.6073 2587 20484 -118.72 34.28 17.0 3051.0 NaN 1705.0 495.0 5.7376 2186 207 rows × 10 columns Let's drop these instances for now
,	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. # for repeatability
]:	<pre># ToT Tepeatability np.random.seed(42) # select the predictor variables and target variables to be used with regression predictors = ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', #dropping categortical features, such as ocean_proximity, including spatial ones such as long/lat. target = 'median_house_value' X = df[predictors].values y = df[target].values</pre>
3]:	<pre># 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] # 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=</pre>
:]:	<pre># scaled version X_train = scaler.fit_transform(X_train_raw) X_test = scaler.transform(X_test_raw) # sanity check print(X_train.shape) print(X_train[:3])</pre> (7000, 8)
	[[1.22783551 -1.3492796
	<pre>def rmse(predicted, actual): return np.sqrt(((predicted - actual)**2).mean()) baseline = rmse(y_train.mean(), y_test).round(1) print(f"test, rmse baseline: {baseline}")</pre>
	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 RM. the RMSE computed using the test data set. When using the KNN algorithm, use algorithm='brute' to get the basic KNN algorithm.
7]:	<pre>knn = KNeighborsRegressor(algorithm="brute") knn.fit(X_train, y_train) predictions = knn.predict(X_test) default = rmse(predictions, y_test).round(1) print(f"test RMSE, default hyperparameters: {default}") test RMSE, default hyperparameters: 62448.9</pre> 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 RI is the RMSE computed using the training data. Use the 'brute' algorithm, and Euclidean distance, which is the default. You need to add get_train_testrmse() function. def get_train_test_rmse(regr, X_train, X_test, y_train, y_test): regr.fit(X_train, y_train) test = regr.predict(X test)
]:	<pre>rmse_te = rmse(test, y_test) train = regr.predict(X_train) rmse_tr = rmse(train, y_train) return rmse_tr, rmse_te n = 30 test_rmse = [] train rmse = []</pre>
	<pre>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</pre>
	<pre># sanity check print('Test RMSE when k = 3: {:0.1f}'.format(test_rmse[1])) 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 to test RMSE value obtained when k=11.</pre> <pre>def get best(ks, rmse):</pre>
	<pre>best_rmse = min(rmse) index = rmse.index(best_rmse) best_k = ks[index] return best_k, best_rmse best_k, best_rmse = get_best(ks, test_rmse) print('best k = {}, best test RMSE: {:0.1f}'.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.</pre>
]:	<pre>plt.plot(ks, test_rmse, label="Test") plt.plot(ks, train_rmse, label="Train") plt.title("Test and Training RMSE vs. K") plt.xlabel("K") plt.ylabel("RMSE") plt.legend() <matplotlib.legend.legend 0x2ab012d1340="" at=""> Test and Training RMSE vs. K</matplotlib.legend.legend></pre>
	Test and Training RMSE vs. K 70000 - 50000 - 40000 - 30000 -
	20000 - 10000
;	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. RMSE for training data seems to always be lower than error for test RMSE. Professor Memo confirmed that this is because test has not previously seen by the model. RMSE for both the training and test predictions appear to converge as K increases. The test RMSE starts high, falls until it reaches the best k value at k = 7, and slowly increases afterwards.
	Impact of noise predictors In class we heard that the KNN performance goes down if useless "noisy predictors" are present. These are predictor that don't help in make 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 Xtest.
	<pre>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, the new noisy predictor is already scaled. Don't modify X_train and Xtest however, as you will be using them again. X_te = X_test.copy() X_tr = X_train.copy() k = 10</pre>
	<pre>k = 10 test_rmse = [] train_rmse = [] for i in range(5): print(i, ' ', end='') regr = KNeighborsRegressor(n_neighbors=k, algorithm="brute", p=2) rmse_tr, rmse_te = get_train_test_rmse(regr, X_tr, X_te, y_train, y_test) train_rmse.append(rmse_tr) test_rmse.append(rmse_te) X_te = add_noise_predictor(X_te) X_tr = add_noise_predictor(X_tr)</pre>
	print ("done") 0 1 2 3 4 done 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 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. percent_increase = []
]:	<pre>for i in test_rmse: percent_increase.append(100 * (i - test_rmse[0]) / test_rmse[0]) plt.plot(range(5), percent_increase) plt.grid(visible=True) plt.title("Percent increase in test RMSE by a number of noise predictors") plt.ylabel("Percent increase in test RMSE") plt.xlabel("Number of noise predictors")</pre> Text(0.5, 0, 'Number of noise predictors')
	Percent increase in test RMSE by a number of noise predictors 17.5 15.0 10.0 Percent increase in test RMSE by a number of noise predictors 10.0
	10.0 7.5 5.0 2.5 0.0 0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5 4.0 Number of noise predictors
	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 experiment 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. $n = 30$
	<pre>raw_test_rmse = [] raw_train_rmse = [] ks = np.arange(1, n+1, 2) for k in ks: print(k, ' ', end='') regr = KNeighborsRegressor(n_neighbors=k, algorithm='brute',p=2) rmse_tr, rmse_te = get_train_test_rmse(regr, X_train_raw, X_test_raw, y_train, y_test) raw_train_rmse.append(rmse_tr) raw_test_rmse.append(rmse_te) print('done')</pre> 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. best_k, best_rmse = get_best(ks, raw_test_rmse) print('best k = {}, best test RMSE: {:0.1f}'.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.
]:	<pre>plt.plot(ks, raw_test_rmse, label="Raw Test") plt.plot(ks, raw_train_rmse, label="Raw Train") plt.title("Unscaled Test and Training RMSE vs. K") plt.xlabel("K") plt.ylabel("RMSE") plt.legend() <matplotlib.legend.legend 0x2ab7bc73ca0="" at=""> Unscaled Test and Training RMSE vs. K Raw Test</matplotlib.legend.legend></pre>
	100000 - Raw Train 80000 - 60000 -
	40000 - 20000 - 20000 - 5 10 15 20 25 30 K
	Comments Reflect on what happened and provide some short commentary, as in previous sections. The shape of the graph appears to be similar to the scaled version, but with higher values of RMSE. The best K value is also higher, at 9 instead of 7. I feel that this probably impacts performance as higher values of K means higher burden in terms of calculation to find distances per point. Higher values in the unscaled training data may also mean more memory is being used. 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, kdtree, and brute. In each case, k=10 and use Euclidean distance. algorithm = ['ball_tree', 'kd_tree', 'brute'] results = [] for algo in algorithm:
	<pre>for algo in algorithm: n = 10 test_rmse = [] train_rmse = [] ks = np.arange(1, n+1, 2) for k in ks: regr = KNeighborsRegressor(n_neighbors=k, algorithm=algo,p=2) 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(k, ' ', end='') best_k, best_rmse = get_best(ks, test_rmse)</pre>
.]:	<pre>for res in results: plt.bar(res[0], res[2]) plt.title("RMSE by KNN algorithm") Text(0.5, 1.0, 'RMSE by KNN algorithm') RMSE by KNN algorithm 60000</pre>
	50000 - 40000 - 30000 -
	20000 - 10000 -
	Comments As usual, reflect on the results and add comments.
	ball_tree kd_tree brute Comments
	Comments As usual, reflect on the results and add comments. It seems as though each algorithm is equally accurate even before rounding the RMSE values. The primary visible difference seems to be time for execution, with brute being the slowest. According to this medium article: https://outline.com/YU7nSM, the difference in speed due to the use of a binary tree data structure in the KD and Ball tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms). 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 'distant 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 algority and Euclidean distance. weights = ['uniform', 'distance'] results = ['uniform', 'distance']
	Comments As usual, reflect on the results and add comments. It seems as though each algorithm is equally accurate even before rounding the RMSE values. The primary visible difference seems to be time for execution, with brute being the slowest. According to this medium article: https://outline.com/YU7nSM, the difference in speece due to the use of a binary tree data structure in the KD and Ball tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms). 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 'distant 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 algority and Euclidean distance. weights = ['uniform', 'distance'] results = ['luniform', 'distance'] results = ['luniform', 'distance'] results = ['luniform', 'distance'] results = ['RNeighborsRegressor(n_neighbors=k, weights=weight,p=2) rmse tr, rmse tese get train_test_rmse(regr, X_train, X_test, y_train, y_test) train_rmse.append(rmse_tr) test_rmse.append(rmse_tr) test_rmse.append(rmse_tr)
]:	Comments As usual, reflect on the results and add comments. It seems as though each algorithm is equally accurate even before rounding the RMSE values. The primary visible difference seems to be time for execution, with brute being the slowest. According to this medium article: https://outline.com/YU7nSM, the difference in speed due to the use of a binary tree data structure in the KD and Ball tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms). 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 'distant 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 algority and Euclidean distance. weights = ['uniform', 'distance'] results = [] for weight in weights: n = 10 test_mse = [] train_mse = [] ks = np.arange(1, n+1, 2) for k in ks: regr = KNeighborsRegressor(n_neighbors=k, weights=weight,p=2) regr = KNeighborsRegressor(n_neighbors=k, weights=weight,p
3]:	Comments As usual, reflect on the results and add comments. It seems as though each algorithm is equally accurate even before rounding the RMSE values. The primary visible difference seems to be time for execution, with brute being the slowest. According to this medium article, https://outline.com/YU7nSM, the difference in speed due to the use of a binary tree data structure in the KD and Ball tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms). Impact of weighting It was briefly mentioned in lecture that there is a variant of KNNI 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 'distant 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 algority and Euclidean distance. weights = ['uniform', 'distance'] results = ['uniform', 'distance'] results = ['uniform', 'distance'] for keight in weights: n = 10 test rase = [] train_trase = [] train_trase = [] train_trase = [] train_trase = [] train_trase_append(trase_te) print(k, '', ende'') best_k, best_trase_append(trase_te) print(k, '', ende'') best_k, best_trase_append(trase_te) print(k', '', ende'')
	Comments As usual, reflect on the results and add comments. It seems as though each algorithm is equally accurate even before rounding the RMSE values. The primary visible difference seems to be time for execution, with brute being the slowest. According to this medium article. https://outline.com/YUTASA, the difference in speed due to the use of a binary tree data structure in the KD and Ball tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms.) 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 KNeighborsRegressori) has two possible values: uniform' and distance for 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. **Velopits** = \text{"uniform"}, 'distance']
2]:	Comments As usual, reflect on the results and add comments. Its seems as though each algorithm is equally accurate even before rounding the RMSE values. The primary visible difference seems to be time for execution, with brute being the slowest. According to this medium article: https://outline.com/PU/DSM, the difference seems to be time for execution, with brute being the slowest. According to this medium article: https://outline.com/PU/DSM, the difference seems to be time for execution, with brute being the slowest. According to this medium article: https://outline.com/PU/DSM, the difference is speed due to the use of a binary tree data structure in the KD and Ball tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms and possibly the dimensionality and structure of data (when considering performance differences between the two tree algorithms). Impact of weight time the tree is a various of KNN in which training points are given more weight when they are closer to the point for which a prediction is to be mode. The weight' parameter of KNeighborsRegressor() has two possible values: uniform and Vistan Choloring in the basic algorithm. Auan on experiment similar to the previous one. Compute the test RMSE for uniform and distance weighting. Using k = 10, the brute algorithms and Excludent distance. **weights = ("uniform", "distance") **results = ("uniform", "distance") **