

Assignment for Week 2 - KNN

Task 1 - KNN Exercise

Data Set: <https://archive.ics.uci.edu/ml/datasets/Heart+Disease>

Notice: this is the same dataset that you worked with last week. Feel free to use your cleaned up data file or you can use the one that I have posted in WorldClass

Objective: Predict heart disease in patients.

- Get to know your data, start out by data exploration. Summarized your finding.
- Divide the data into training set and test set randomly with ratio 80:20. Make prediction based on 1-nearest neighbor. What is the error rate of this approach? Report your results in a confusion matrix.
- Use different values for k , what is the optimal value of k from your experiments? Report the error rate of the optimal k value and its confusion matrix, is there any improvement (by how much) over 1-nearest neighbor?
- Is there anything else you can do to improve your model? If yes, demonstrate your approach. (Hint: there is always something that you can try, unless your accuracy score is 100%)

Task 2 - KNN Project

Data Set: <https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>

Objective: Determine whether the product is 'yes' or 'no' subscription

- Use KNN to provide an analyze of the given data set. Make sure you are answering the stated objective above.
- What is your optimal k ? What is the accuracy rate?
- Discover any insights from this analysis? (include numbers/graphs corresponding to your reasoning)

Deliverables:

Upload your notebook's .ipynb file (this assignment can be done on 1 notebook or 2. Your preference.) Please include the jupyter notebook file and a pdf printout of your notebook.

Important: Make sure your provide complete and thorough explanations for all of your analysis. You need to defend your thought process and reasoning.

Introduction

Week 2 is about KNN method. KNN is a supervised machine learning algorithm that can be used to solve classification and regression problems. This week we are learning how to part, part 1 is exercise on Heart disease data from week 1. The second part is the project where we use KNN on Bank data.

Dataset

For the exercise, we are going to use Heart disease dataset which is set of records of people with there diagnosis of heart disease. The data contains 15 variables where 'num' is the target and the other 14 are our features. those features are various from chest pain to number of years of smoking. As for the project, the dataset is Bank dataset which contains personal information such as marital status and education level. This dataset contains 17 columns where 'y' is the target. Because this data has many categorical variable, we need to convert them to numeric in order to fit them in KNN model.

exercise 1

```
# Loading required libraries
import pandas as pd
import sklearn
from sklearn.neighbors import KNeighborsRegressor
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.decomposition import PCA

# plotting
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib
sns.set()

# Loading cleaned Heart disease dataset.
heartd = pd.read_csv('heart_disease_data_clean.csv')

# Data is clean and numeric. No cleaning needed except for dropping unnecessary columns for KNN model
# which will be determined after seeing correlation.

heartd.head(10)

# Out[3]:
   age  sex  cp  trestbps  chol  cigs  yrs  fbs  famhist  restecg  thalach  exang  thal  num
0  63   1   1   145  233  500  200  0   1   2  150  0  6  0
1  67   1   4   160  286  400  400  0   1   2  108  1  3  2
2  67   1   4   120  229  200  350  0   1   2  129  1  7  1
3  37   1   3   130  250  0  0  0  0   1   0  187  0  3  0
4  41   0   2   130  204  0  0  0  0   1   2  172  0  3  0
5  56   1   2   120  236  200  200  0   1   0  178  0  3  0
6  62   0   4   140  268  0  0  0  0   1   2  160  0  3  3
7  57   0   4   120  354  0  0  0  0   1   0  163  1  3  0
8  63   1   4   130  254  0  0  0  0   2  145  1  7  2
9  53   1   4   140  203  200  250  1   1   2  157  0  7  1

# In [4]:
heartd.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 282 entries, 0 to 281
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0  age         282 non-null      int64
 1  sex         282 non-null      int64
 2  cp          282 non-null      int64
 3  trestbps    282 non-null      int64
 4  chol        282 non-null      float64
 5  cigs        282 non-null      float64
 6  yrs         282 non-null      int64
 7  fbs         282 non-null      int64
 8  famhist     282 non-null      int64
 9  restecg     282 non-null      int64
10  thalach     282 non-null      int64
11  exang       282 non-null      int64
12  thal        282 non-null      int64
13  num         282 non-null      int64
memory usage: 131.0 KB

# In [5]:
heartd.describe()

# Out[5]:
   age      sex      cp      trestbps      chol      cigs      yrs      fbs      famhist      restecg      thalach
count  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000
mean    63.000000  0.677005    3.163121  131.159035  247.705674  16.836011  15.347364  1.048936  0.620567  1.014184  149.826241
std     9.503083  0.468338  0.955405  16.739821  64.179771  18.876511  15.276814  0.356658  0.486108  0.998118  22.737594
min      29.000000  0.000000  0.000000  94.000000  126.000000  0.000000  0.000000  0.000000  0.000000  0.000000  88.000000
25%     48.000000  0.000000  3.000000  120.000000  213.000000  0.000000  0.000000  0.000000  0.000000  0.000000  133.250000
50%     55.000000  1.000000  3.000000  130.000000  244.000000  11.976385  15.000000  0.000000  1.000000  2.000000  153.500000
75%     61.000000  1.000000  4.000000  140.000000  277.000000  33.000000  15.000000  0.000000  1.000000  2.000000  165.750000
max     77.000000  1.000000  4.000000  170.000000  360.000000  75.000000  54.000000  1.000000  1.000000  2.000000  202.000000

# In [7]:
plt.figure(figsize=(16,6))
_ = sns.heatmap(heartd.corr(), vmin=-1, vmax=1, annot=True)

# Out[7]:
   age      sex      cp      trestbps      chol      cigs      yrs      fbs      famhist      restecg      thalach
age      1  -0.091  0.077  0.28  0.21  -0.13  0.071  0.1  0.026  0.13  -0.39  0.087  0.1  0.21
sex      -0.091  1  0.023  0.05  0.17  0.33  0.25  0.054  -0.071  0.033  -0.053  0.19  0.39  0.24
cp        0.077  0.023  1  -0.077  0.1  0.086  0.081  -0.072  0.11  0.061  -0.32  0.36  0.26  0.38
trestbps  0.28  0.05  -0.077  1  0.18  -0.028  0.05  0.13  0.046  0.13  -0.036  0.48  0.11  0.14
chol      0.21  -0.17  0.1  0.18  1  -0.092  -0.022  0.022  0.097  0.15  -0.031  0.11  0.024  0.11
cigs      -0.13  0.33  0.086  -0.026  -0.092  1  0.65  0.04  -0.026  0.028  0.0065  0.45  0.2  0.046
yrs       0.071  0.25  0.081  0.05  -0.022  0.65  1  0.021  0.062  0.087  0.0075  0.16  0.22  0.078
fbs       0.1  0.054  -0.072  0.13  0.022  0.04  -0.021  1  0.055  0.054  0.019  0.028  0.07  0.04
famhist   0.026  0.071  0.11  0.046  0.097  -0.028  0.062  0.042  1  0.005  -0.004  0.017  0.049
restecg   0.13  0.033  0.061  0.13  0.15  -0.028  0.017  0.004  0.055  1  -0.006  0.089  0.015  0.19
thalach   -0.39  -0.053  -0.32  -0.036  -0.031  -0.065  0.075  0.019  -0.006  -0.086  1  0.38  -0.25  -0.4
exang     0.087  0.19  0.36  0.48  0.11  0.045  0.016  0.028  0.0014  0.089  0.38  1  0.34  0.4
thal      0.1  0.39  0.26  0.11  0.024  0.02  0.22  0.07  0.017  0.015  -0.25  0.34  1  0.61
num       0.21  0.24  0.38  0.14  0.11  0.046  0.078  0.046  0.19  0.04  0.4  0.61  0.4  1

# Since maximum correlation with num is 0.51, i am going to drop very low correlated variables. In this case, cigarettes consumption and number of years of smoking are below 0.07. Also fbs is 0.04 and Family history is 0.049. Those variables have very low correlation with the target compared to other independent variables which most of them have low correlation with num. However, I cannot drop most of them and only have 1 or 2 variables.
```

```
# trimming our data set
# cigs, years, fbs and famhist
heartd.drop(['cigs','years','fbs','famhist'],axis=1,inplace=True)

# Out[8]:
KeyError                                Traceback (most recent call last)
<ipython-input-9-f5b2e1cd6623> in <module>
      2 # trimming our data set
      3 # cigs, years, fbs, and famhist
----> 3 heartd.drop(['cigs','years','fbs','famhist'],axis=1,inplace=True)

C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\frame.py in drop(self, labels, axis, index, columns, level, inplace, errors)
    1307         """
    1308         return super().drop(
    1309             labels=labels,
    1310             axis=axis,
C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\generic.py in drop(self, labels, axis, index, columns, level, inplace, errors)
    1152         if labels is not None:
    1153             obj = obj._drop_axis(labels, axis, level=level, errors=errors)
    1154         if inplace:
C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\generic.py in _drop_axis(self, labels, axis, level, errors)
    1187             new_axis = axis.drop(labels, level=level, errors=errors)
    1188         else:
    1189             new_axis = axis.drop(labels, errors=errors)
    1190             result = self.reindex(**{axis.name: new_axis})
C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\indexes\base.py in drop(self, labels, errors)
    5589         if mask.any():
    5590             # if errors is "ignore":
--> 5591             raise KeyError(f"{labels[mask]} not found in axis")
    5592             # if errors is "indexer":
    5593             return self.delete(indexer)
KeyError: ["cigs" "years" "fbs" "famhist"] not found in axis

# In [10]:
heartd.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 282 entries, 0 to 281
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0  age         282 non-null      int64
 1  sex         282 non-null      int64
 2  cp          282 non-null      int64
 3  trestbps    282 non-null      int64
 4  chol        282 non-null      int64
 5  restecg     282 non-null      int64
 6  thalach     282 non-null      int64
 7  exang       282 non-null      int64
 8  thal        282 non-null      int64
 9  num         282 non-null      int64
memory usage: 22.8 KB

# In [11]:
# Defining our features and targets
cols = heartd.columns
target_col = 'num'
feat_cols = [c for c in cols if c != target_col]

# there is nothing magical about the X and y notation here.
# however, it seems to be a fairly standard notation, so we will use is here
X = heartd[feat_cols].values
y = heartd[target_col].values

# In [12]:
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# In [13]:
# Make prediction based on 1-nearest neighbor.
# What is the error rate of this approach?
# Report your results in a confusion matrix.

# define and fit our model
model = KNeighborsRegressor(n_neighbors=1, n_jobs=-1)
model.fit(X_train, y_train)

# Out[13]:
KNeighborsRegressor(n_jobs=-1, n_neighbors=1)

# In [14]:
# gather the predictions that our model made for our test set
preds = model.predict(X_test)

# display the actuals and predictions for the test set
print('Actuals for test data set')
print(y_test)
print('Predictions for test data set')
print(preds)

Actuals for test data set:
[0 0 1 4 1 3 0 0 1 4 3 1 0 3 0 4 0 0 0 1 2 0 0 0 0 3 0 2 0 2 0 3 0 0 2
 3 0 1 4 4 1 0 0 0 0 3 0 3 1 0 3 1 0 2 0 3]
Predictions for test data set:
[1. 0. 0. 1. 0. 3. 0. 0. 3. 4. 2. 0. 2. 0. 2. 1. 0. 1. 0. 1. 3. 0. 1. 0.
 1. 0. 0. 2. 0. 1. 3. 0. 2. 0. 2. 0. 0. 1. 4. 0. 0. 0. 1. 0. 0. 0. 0.
 1. 0. 0. 0. 0. 2. 1. 3. 0. 1.]

# In [15]:
# compare the two sets
differs = y_test - preds
print('Differences between the two sets')
print(differs)

Differences between the two sets
[-2, -1, 0, 1, 2, 0, -3, 1, -2, 0, 1, 1, -2, 3, 0, 2, -1, 3,
 1, -1, -3, 1, 0, 1, 0, 2, 0, 1, -1, 1, 0, 0, 1, 0, 0, 0,
 1, -1, 0, 1, 4, 4, 1, -1, 0, 0, 0, 3, -1, 0, 3, 1, 0, 1,
 -1, 0, 2, 1]

# In [16]:
# R^2 (coefficient of determination) regression score function.
# Best possible score is 1.0 and it can be negative
# (because the model can be arbitrarily worse).
# If a constant model that always predicts the expected value of y,
# disregarding the input features, would get a R^2 score of 0.0.
from sklearn.metrics import r2_score
print(r2_score(y_test, preds))

-0.1974050046339202

# In [17]:
# Explained variance regression score function
# Best possible variance is 1.0, lower values are worse.
from sklearn.metrics import explained_variance_score
print(explained_variance_score(y_test, preds))

-0.13561940067964162

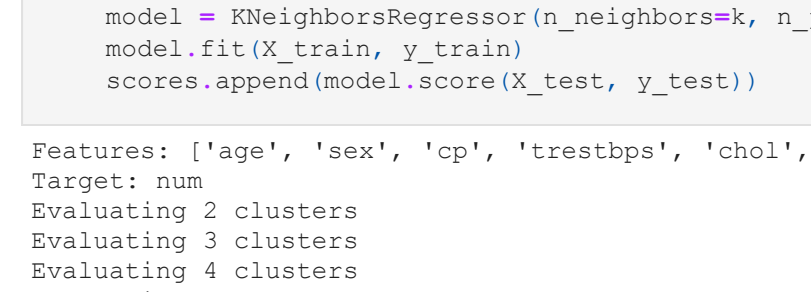
# In [18]:
# Determining the optimal number of clusters
scores = []
# remember the ending number for range is not inclusive
for k in range(2, 20):
    # output to let us know where we are
    print(f'Evaluating {k} clusters')
    # n_jobs=-1 will use all processors on your system
    model = KNeighborsRegressor(n_neighbors=k, n_jobs=-1)
    model.fit(X_train, y_train)
    scores.append(k).score(X_test, y_test)

Features: ['age', 'sex', 'cp', 'trestbps', 'chol', 'restecg', 'thalach', 'exang', 'thal']
Target: num
Evaluating 2 clusters
Evaluating 3 clusters
Evaluating 4 clusters
Evaluating 5 clusters
Evaluating 6 clusters
Evaluating 7 clusters
Evaluating 8 clusters
Evaluating 9 clusters
Evaluating 10 clusters
Evaluating 11 clusters
Evaluating 12 clusters
Evaluating 13 clusters
Evaluating 14 clusters
Evaluating 15 clusters
Evaluating 16 clusters
Evaluating 17 clusters
Evaluating 18 clusters
Evaluating 19 clusters

# In [19]:
scores

Out[19]:
(-0.1027571825764595,
 0.02276204945319768,
 0.09342401013146461,
 0.0526413456904563,
 0.11734503151083815,
 0.0974518093094527,
 0.0462233549582962,
 0.0179807273534691,
 0.07192307692307709,
 0.0363584854357052,
 0.01072327606832216,
 0.04849438719831545,
 0.0359104052043658,
 0.08182670586109032,
 0.0688305244667802,
 0.07219047128681585,
 0.08526413345690453)

# In [20]:
# display the results
plt.plot(range(2, 20), scores)
plt.scatter(range(2, 20), scores)
plt.grid()
plt.xticks(range(2, 20))
```



```
# define and fit our model
model = KNeighborsRegressor(n_neighbors=6, n_jobs=-1)
model.fit(X_train, y_train)

# gather the predictions that our model made for our test set
preds = model.predict(X_test)

# display the actuals and predictions for the test set
print('Actuals for test data set')
print(y_test)
print('Predictions for test data set')
print(preds)

Actuals for test data set:
[0 0 1 3 0 0 1 1 3 1 0 3 0 4 0 0 0 1 2 0 0 0 0 3 0 2 0 2 0 3 0 0 2
 3 0 1 4 4 1 0 0 0 0 3 0 3 1 0 3 1 0 2 0 3]
Predictions for test data set:
[1. 0. 0. 1. 0. 3. 0. 0. 3. 4. 2. 0. 2. 0. 2. 1. 0. 1. 0. 1. 3. 0. 1. 0.
 1. 0. 0. 2. 0. 1. 3. 0. 2. 0. 2. 0. 0. 1. 4. 0. 0. 0. 1. 0. 0. 0. 0.
 1. 0. 0. 0. 0. 2. 1. 3. 0. 1.]

# In [21]:
differs = y_test - preds
print('Differences between the two sets')
print(differs)

Differences between the two sets:
[1.16666667 0.33333333 0. 0.66666667 2. -1.16666667
 -1.16666667 0.33333333 1. -1.16666667 2.4. -1.26666667 1.66666667
 0.66666667 2. -1.66666667 0.66666667 1.5 0. 0.33333333
 0.33333333 1.83333333 0.33333333 0.33333333 1.33333333 0.
 0.33333333 0.66666667 0. 0.66666667 0.33333333 1.5
 0.66666667 0.66666667 2.5 1.16666667 -0.5 1.
 1.16666667 -0.66666667 1.66666667]
r2_score: 0.1174050046339202

# In [24]:
# visualize model
pc = PCA()
tr_pca = pc.fit_transform(X_train)
te_pca = pc.transform(X_test)

# In [26]:
# indexing of a matrix in Python (numpy) is [rows, cols]
plt.scatter(tr_pca[:, 0], tr_pca[:, 1], c=y_train)
plt.colorbar()

# In [27]:
plt.scatter(te_pca[:, 0], te_pca[:, 1], c=y_test)
plt.colorbar()

# In [28]:
plt.scatter(te_pca[:, 0], te_pca[:, 1], c=preds)
plt.colorbar()

# In [29]:
plt.scatter(te_pca[:, 0], te_pca[:, 1], c=differs)
plt.colorbar()

# In [30]:
# normalize data to improve model maybe
heartd.describe()

# Out[30]:
   age      sex      cp      trestbps      chol      restecg      thalach      exang      thal      num
count  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000  282.000000
mean    63.113408  0.677305    3.163121  131.159035  247.705674  1.014184  149.826241  0.326241  4.677305  0.927801
std     9.541948  0.468338  0.955405  16.739821  64.179771  0.998118  22.737594  0.469670  1.936386  1.224894
min      29.000000  0.000000  0.000000  94.000000  126.000000  0.000000  88.000000  0.000000  3.000000  0.000000
25%     48.000000  0.000000  3.000000  120.000000  213.000000  0.000000  133.250000  0.000000  3.000000  0.000000
50%     55.000000  1.000000  3.000000  130.000000  244.000000  2.000000  153.500000  0.000000  3.000000  0.000000
75%     61.000000  1.000000  4.000000  140.000000  277.000000  2.000000  165.750000  1.000000  7.000000  2.000000
max     77.000000  1.000000  4.000000  170.000000  360.000000  2.000000  202.000000  1.000000  7.000000  4.000000

# In [32]:
heartd.head(10)

# Out[32]:
   age  sex  cp  trestbps  chol  restecg  thalach  exang  thal  num
0  63   1   1   145  233  500  200  0   1   2  150  0  6  0
1  67   1   4   160  286  400  400  0   1   2  108  1  3  2
2  67   1   4   120  229  200  350  0   1   2  129  1  7  1
3  37   1   3   130  250  0  187  0  3  0
4  41   0   2   130  204  0  172  0  3  0
5  56   1   2   120  236  200  178  0  3  0
6  62   0   4   140  268  0  160  0  3  3
7  57   0   4   120  354  0  163  1  3  0
8  63   1   4   130  254  0  147  0  7  2
9  53   1   4   140  203  200  155  1  7  1

# In [33]:
# we only want to normalize our feature columns in the dataset and I don't want to resplit the dataset.
# if I do, I will have to re-split the X_train and X_test only
from sklearn import preprocessing
min_max_scaler = preprocessing.MinMaxScaler()
X_tr_norm = min_max_scaler.fit_transform(X_train)
X_te_norm = min_max_scaler.fit_transform(X_test)

# In [34]:
new_heartd_tr = pd.DataFrame(X_tr_norm, columns=feat_cols)
new_heartd_te = pd.DataFrame(X_te_norm, columns=feat_cols)
new_heartd_tr['num'] = y_train
new_heartd_te['num'] = y_test
new_heartd_tr.head(10)

# Out[34]:
   age  sex  cp  trestbps  chol  restecg  thalach  exang  thal  num
0  0.458333  0  0.000000  0.078995  0.508547  0  0.168421  0.0  0.0  0
1  0.687500  1  0.000000  0.578947  0.717949  0  0.059811  0.0  0.0  2
2  0.416667  1  0.333333  0.437589  0.594291  0  0.025632  0.0  1.0  0
3  0.812500  1  0.666667  0.315789  0.645299  0  0.057620  0.0  1.0  0
4  0.625000  1  0.333333  0.605263  0.405883  0  0.066667  1.0  0.0  0
5  0.833333  1  0.333333  0.394737  0.576923  0  0.046912  0.0  1.0  1
6  0.562500  1  0.333333  0.324105  0.470085  0  0.078947  0.0  0.0  0
7  0.479167  1  0.333333  0.324105  0.850427  0  0.073842  0.0  0.0  0
8  0.583333  1  0.666667  0.736842  0.179487  0  0.054386  0.0  0.0  0
9  0.583333  1  1.000000  0.605263  0.282051  0  0.052616  0.0  0.75  0

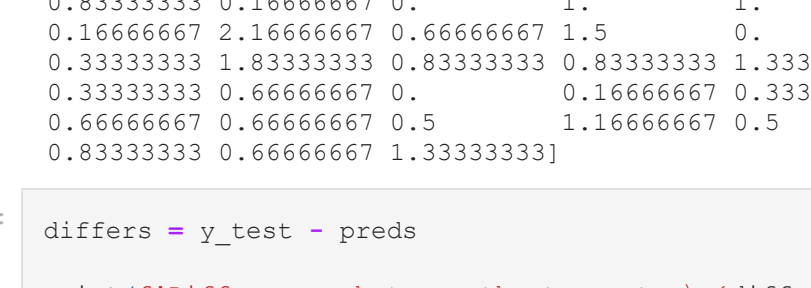
# In [35]:
new_heartd_te = pd.DataFrame(X_te_norm, columns=feat_cols)
new_heartd_tr['num'] = y_train
new_heartd_te.head(10)

# Out[35]:
   age  sex  cp  trestbps  chol  restecg  thalach  exang  thal  num
0  0.675676  1  0.666667  0.384615  0.450  0  0.536842  0.0  1.0  0
1  0.189189  0  0.666667  0.046154  0.000  0  0.082105  0.0  0.0  0
2  0.432652  1  0.666667  0.384615  0.525  0  0.821053  0.0  0.0  0
3  0.765757  1  0.000000  0.507692  0.705  1  0.831579  0.0  0.0  1
4  0.575688  0  0.333333  0.476923  0.890  1  0.600000  0.0  0.0  3
5  1.000000  0  0.333333  0.230769  0.640  1  0.273684  0.0  0.0  0
6  0.378378  0  0.666667  0.538462  0.835  1  0.494737  0.0  0.0  0
7  0.434232  1  1.000000  0.538462  0.310  1  0.631579  1.0  1.0  1
8  0.567568  1  1.000000  0.630769  0.385  0  0.105263  0.0  1.0  1
9  0.648649  1  1.000000  0.507692  0.125  1  0.315789  1.0  0.0  4

# In [36]:
scores_norm = []
print(f'Features: {feat_cols} \nTarget: {target_col}')
# remember the ending number for range is not inclusive
for k in range(1, 20):
    # output to let us know where we are
    print(f'Evaluating {k} clusters')
    # n_jobs=-1 will use all processors on your system
    model_norm = KNeighborsRegressor(n_neighbors=k, n_jobs=-1)
    model_norm.fit(X_tr_norm, y_train)
    scores_norm.append(k).score(X_te_norm, y_test)

Features: ['age', 'sex', 'cp', 'trestbps', 'chol', 'restecg', 'thalach', 'exang', 'thal']
Target: num
Evaluating 1 clusters
Evaluating 2 clusters
Evaluating 3 clusters
Evaluating 4 clusters
Evaluating 5 clusters
Evaluating 6 clusters
Evaluating 7 clusters
Evaluating 8 clusters
Evaluating 9 clusters
Evaluating 10 clusters
Evaluating 11 clusters
Evaluating 12 clusters
Evaluating 13 clusters
Evaluating 14 clusters
Evaluating 15 clusters
Evaluating 16 clusters
Evaluating 17 clusters
Evaluating 18 clusters
Evaluating 19 clusters

# In [37]:
# display the results
plt.plot(range(1, 20), scores_norm)
plt.scatter(range(1, 20), scores_norm)
plt.grid()
plt.xticks(range(1, 20))
```



```
# define and fit our model with k=12
model_norm = KNeighborsRegressor(n_neighbors=15, n_jobs=-1)
model_norm.fit(X_tr_norm, y_train)

# gather the predictions that our model made for our test set
preds_norm = model_norm.predict(X_te_norm)

# display the actuals and predictions for the test set
print('Actuals for test data set')
print(y_test)
print('Predictions for test data set')
print(preds_norm)

Actuals for test data set:
[0 0 1 3 0 0 1 1 3 1 0 3 0 4 0 0 0 1 2 0 0 0 0 3 0 2 0 2 0 3 0 0 2
 3 0 1 4 4 1 0 0 0 0 3 0 3 1 0 3 1 0 2 0 3]
Predictions for test data set:
[1.1333333
```



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45211 entries, 0 to 45210
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   age         45211 non-null    int64
 1   job         45211 non-null    int64
 2   marital     45211 non-null    int64
 3   education   45211 non-null    int64
 4   default     45211 non-null    int64
 5   balance     45211 non-null    int64
 6   housing     45211 non-null    int64
 7   loan        45211 non-null    int64
 8   contact     45211 non-null    int64
 9   day         45211 non-null    int64
10  month       45211 non-null    int64
11  duration    45211 non-null    int64
12  campaign    45211 non-null    int64
13  pdays       45211 non-null    int64
14  previous    45211 non-null    int64
15  outcome     45211 non-null    int64
16  y           45211 non-null    int64
memory usage: 5.9 MB

In [47]: bank.head(10)

Out[47]:
   age  job  marital  education  default  balance  housing  loan  contact  day  month  duration  campaign  pdays  previous  outcome  y
0    44   9     2       1       0     29       0     0     0     5     5     5     151     1     -1     0     0
2    33   2       1       1       0      2       0     0     0     5     5     5     76     1     -1     0     0
3    47   1       1       0     0    1506     0     1     0     5     5     5     92     1     -1     0     0
4    33  11     2       0     0      1       1     1     0     5     5     5    198     1     -1     0     0
5    35   4       1       3     0     231     0     1     0     5     5     5    139     1     -1     0     0
6    28   4       2       3     0     447     0     0     0     5     5     5    217     1     -1     0     0
7    42   2       0       3     1     2       0     1     0     5     5     5    380     1     -1     0     0
8    58   5       1       2     0    121     0     1     0     5     5     5     50     1     -1     0     0
9    43   9     2       1     0     593     0     1     0     5     5     5     55     1     -1     0     0

In [48]: plt.figure(figsize=(16, 6))
_ = sns.heatmap(bank.corr(), vmin=-1, vmax=1, annot=True)

age      job      marital      education      default      balance      housing      loan      contact      day      month      duration      campaign      pdays      previous      outcome      y
job      1.000000      0.0222      0.0191      0.018      0.19      0.016      0.026      0.0091      0.093      0.0046      0.0046      0.024      0.00013      0.014      0.025
marital  0.0222      1.0000      0.042      0.0007      0.0021      0.016      0.047      0.039      0.0053      0.051      0.012      0.009      0.019      0.015      0.025      0.046
education 0.0191      0.0007      0.016      1.0000      0.0073      0.002      0.002      0.009      0.008      0.011      0.014      0.0073      0.0038      0.016      0.015      0.047
default  0.016      0.0007      0.0021      0.0073      1.0000      0.006      0.006      0.007      0.015      0.004      0.015      0.011      0.017      0.013      0.018      0.042
balance  0.19      0.016      0.0021      0.0073      0.006      1.0000      0.041      0.019      0.027      0.0045      0.005      0.022      0.015      0.0034      0.017      0.035      0.053
housing  0.016      0.003      0.047      0.009      0.007      0.004      1.0000      0.041      0.011      0.011      0.022      0.012      0.011      0.023      0.011      0.04      0.068
loan     0.026      0.0091      0.0046      0.0073      0.006      0.041      0.019      1.0000      0.009      0.011      0.005      0.005      0.011      0.009      0.015      0.015      0.015
contact  0.0091      0.0053      0.0021      0.015      0.0045      0.005      0.022      0.011      1.0000      0.011      0.011      0.011      0.011      0.011      0.011      0.011      0.011
day      0.011      0.011      0.011      0.011      0.011      0.011      0.011      0.011      1.0000      0.011      0.011      0.011      0.011      0.011      0.011      0.011      0.011
month     0.014      0.0073      0.0038      0.016      0.015      0.011      0.011      0.005      0.005      1.0000      0.011      0.011      0.011      0.011      0.011      0.011      0.011
duration  0.0073      0.0038      0.0038      0.016      0.015      0.011      0.011      0.005      0.005      0.011      1.0000      0.011      0.011      0.011      0.011      0.011      0.011
campaign  0.0046      0.0046      0.0046      0.0046      0.0046      0.0046      0.0046      0.0046      0.0046      0.0046      0.0046      1.0000      0.011      0.011      0.011      0.011      0.011
pdays    0.025      0.015      0.015      0.015      0.015      0.015      0.015      0.015      0.015      0.015      0.015      0.015      1.0000      0.011      0.011      0.011      0.011
previous  0.014      0.0055      0.025      0.022      0.044      0.035      0.031      0.04      0.27      0.082      0.036      0.013      0.11      0.79      0.48      1.0000      0.011
outcome   0.025      0.04      0.046      0.047      0.022      0.053      0.14      0.068      0.15      0.028      0.019      0.39      0.073      0.1      0.093      0.22      1.0000
y         0.025      0.04      0.046      0.047      0.022      0.053      0.14      0.068      0.15      0.028      0.019      0.39      0.073      0.1      0.093      0.22      1.0000

In [49]: # trimming our data set
bank.drop(['age', 'job', 'marital', 'education', 'month', 'day', 'y'], axis=1, inplace=True)

In [50]: cols = bank.columns
target_col = 'y'
feat_cols = [c for c in cols if c != target_col]

# there is nothing magical about the X and y notation here.
# however, it seems to be a fairly standard notation, so we will use is here
X = bank[feat_cols].values
y = bank[target_col].values

In [51]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2)

In [54]: # define and fit our model
model = KNeighborsClassifier(n_neighbors=2, n_jobs=-1)
model.fit(X_train, y_train)

Out[54]: KNeighborsClassifier(n_jobs=-1, n_neighbors=2)

In [55]: # gather the predictions that our model made for our test set
preds = model.predict(X_test)

# display the actuals and predictions for the test set
print('Actuals for test data set')
print(y_test)
print('Predictions for test data set')
print(preds)

Actuals for test data set
[0 0 ... 0 0 0]
Predictions for test data set
[0 0 ... 0 0 0]

In [56]: #compare the two sets for 'subscribed'
differs = y_test - preds
print('Differences between the two sets')
print(differs)

Differences between the two sets
[0 0 ... 0 0 0]

In [57]: # R^2 (coefficient of determination) regression score function.
# Best possible score is 1.0 and it can be negative
# (because the model can be arbitrarily worse).
# A constant model that always predicts the expected value of y,
# regardless of the input features, would get a R^2 score of 0.0.
from sklearn.metrics import r2_score

print(r2_score(y_test, preds))

-0.17576494715312996

In [58]: # Explained variance regression score function
# Best possible score is 1.0, lower values are worse.
from sklearn.metrics import explained_variance_score

print(explained_variance_score(y_test, preds))

-0.12024859165847525

In [59]: scores = []
print(f'Features: {feat_cols} \nTarget: {target_col}')

# remember the ending number for range is not inclusive
for k in range(2, 20):
    # output so let us know where we are
    print(f'Evaluating {k} clusters')

    # n_jobs=-1 will use all processors on your system
    model = KNeighborsRegressor(n_neighbors=k, n_jobs=-1)
    model.fit(X_train, y_train)
    scores.append(model.score(X_test, y_test))

Features: ['default', 'balance', 'housing', 'loan', 'contact', 'duration', 'campaign', 'pdays', 'previous', 'poutcome']
Target: y
Evaluating 2 clusters
Evaluating 3 clusters
Evaluating 4 clusters
Evaluating 5 clusters
Evaluating 6 clusters
Evaluating 7 clusters
Evaluating 8 clusters
Evaluating 9 clusters
Evaluating 10 clusters
Evaluating 11 clusters
Evaluating 12 clusters
Evaluating 13 clusters
Evaluating 14 clusters
Evaluating 15 clusters
Evaluating 16 clusters
Evaluating 17 clusters
Evaluating 18 clusters
Evaluating 19 clusters

In [60]: scores

Out[60]:
[-0.12735849943022897,
 -0.022771892071443034,
 0.04026101499003165,
 0.0852381935978396,
 0.1054075656965444,
 0.1246872867049445,
 0.1385871600063255,
 0.14989702988045117,
 0.15714956212988265,
 0.1607683816366695,
 0.16654398428200312,
 0.1717795090767929,
 0.17683477569997308,
 0.1787299742066205,
 0.184674632688379,
 0.19018824895153637,
 0.18923267805278986,
 0.1943152913046287]

In [61]: # display the results
plt.plot(range(2, 20), scores)
plt.scatter(range(2, 20), scores)
plt.grid()
_ = plt.xticks(range(2, 20))

In [62]: # define and fit our model
model = KNeighborsRegressor(n_neighbors=17, n_jobs=-1)
model.fit(X_train, y_train)

# gather the predictions that our model made for our test set
preds = model.predict(X_test)

# display the actuals and predictions for the test set
print('Actuals for test data set')
print(y_test)
print('Predictions for test data set')
print(preds)

Actuals for test data set
[0 0 ... 0 0 0]
Predictions for test data set
[0. 0.05882353 -0.05882353 ... 0.05882353 -0.05882353 0. ]

In [63]: differs = y_test - preds

print(f'Differences between the two sets:\n{differs}')

print(f'r2_score: {r2_score(y_test, preds)}')

Differences between the two sets:
[0. 0.05882353 -0.05882353 ... -0.05882353 -0.05882353 0. ]
r2_score: 0.19018824895153637

In [64]: #visualization
pc = PCA()
tr_pca = pc.fit_transform(X_train)
te_pca = pc.transform(X_test)

In [65]: # indexing of a matrix in Python (numpy) is (rows, cols)
_ = plt.scatter(tr_pca[:, 0], tr_pca[:, 1], c=y_train)
_ = plt.colorbar()

In [66]: _ = plt.scatter(te_pca[:, 0], te_pca[:, 1], c=y_test)
_ = plt.colorbar()

In [67]: _ = plt.scatter(te_pca[:, 0], te_pca[:, 1], c=preds)
_ = plt.colorbar()

In [68]: _ = plt.scatter(te_pca[:, 0], te_pca[:, 1], c=differs)
_ = plt.colorbar()

In [69]: # we only want to normalize our feature columns in the dataset and I don't want to resplit the dataset.
from sklearn import preprocessing

min_max_scaler = preprocessing.MinMaxScaler()
X_tr_norm = min_max_scaler.fit_transform(X_train)
X_te_norm = min_max_scaler.transform(X_test)

In [70]: new_bank_tr = pd.DataFrame(X_tr_norm, columns=feat_cols)
new_bank_tr['y'] = y_train
new_bank_tr.head(10)

Out[71]:
   default  balance  housing  loan  contact  duration  campaign  pdays  previous  outcome  y
0    0.0  0.080511    0.0  1.0    0.014640  0.000000  0.000000  0.000000  0.000000  0.0  0
1    0.0  0.095650    0.0  1.0    0.054697  0.016129  0.000000  0.000000  0.000000  0.0  0
2    0.0  0.107694    0.0  1.0    0.026434  0.048387  0.000000  0.000000  0.000000  0.0  0
3    0.0  0.071723    0.0  1.0    0.076251  0.161290  0.000000  0.000000  0.000000  0.0  0
4    0.0  0.104561    1.0  1.0    0.053680  0.016129  0.000000  0.000000  0.000000  0.0  1
5    0.0  0.080511    0.0  1.0    0.053070  0.080645  0.000000  0.000000  0.000000  0.0  0
6    0.0  0.076417    0.0  0.0    0.043514  0.016129  0.000000  0.000000  0.000000  0.0  0
7    0.0  0.077788    1.0  1.0    0.030297  0.000000  0.000000  0.000000  0.000000  0.0  0
8    0.0  0.084787    1.0  1.0    0.069134  0.000000  0.055505  0.000000  0.000000  1.0  1
9    0.0  0.084787    1.0  1.0    0.069134  0.000000  0.055505  0.000000  0.000000  1.0  1

In [71]: new_auto_te = pd.DataFrame(X_te_norm, columns=feat_cols)
new_auto_te['y'] = y_test
new_auto_te.head(10)

Out[71]:
   default  balance  housing  loan  contact  duration  campaign  pdays  previous  outcome  y
0    0.0  0.032552    0.0  1.0    0.0  0.060320  0.000000  0.000000  0.000000  0.000000  0
1    0.0  0.062850    1.0  1.0    0.0  0.026076  0.018519  0.000000  0.000000  0.000000  0
2    0.0  0.032137    0.0  1.0    1.0  0.071002  0.000000  0.000000  0.000000  0.000000  0
3    0.0  0.044349    1.0  1.0    1.0  0.029707  0.000000  0.399763  0.017241  0.666667  0
4    0.0  0.028963    1.0  0.0    1.0  0.038014  0.018519  0.000000  0.000000  0.000000  0
5    0.0  0.038031    1.0  1.0    1.0  0.040842  0.037037  0.000000  0.000000  0.000000  0
6    0.0  0.037546    1.0  0.0    1.0  0.029846  0.092593  0.000000  0.000000  0.000000  0
7    0.0  0.028073    1.0  1.0    1.0  0.026076  0.055556  0.000000  0.000000  0.000000  0
8    0.0  0.026847    0.0  1.0    1.0  0.016023  0.000000  0.000000  0.000000  0.000000  0
9    0.0  0.023968    1.0  0.0    0.0  0.184417  0.018519  0.000000  0.000000  0.000000  0

In [74]: scores_norm = []
print(f'Features: {feat_cols} \nTarget: {target_col}')

# remember the ending number for range is not inclusive
for k in range(2, 30):
    # output so let us know where we are
    print(f'Evaluating {k} clusters')

    # n_jobs=-1 will use all processors on your system
    model_norm = KNeighborsRegressor(n_neighbors=k, n_jobs=-1)
    model_norm.fit(X_tr_norm, y_train)
    scores_norm.append(model_norm.score(X_te_norm, y_test))

Features: ['default', 'balance', 'housing', 'loan', 'contact', 'duration', 'campaign', 'pdays', 'previous', 'poutcome']
Target: y
Evaluating 2 clusters
Evaluating 3 clusters
Evaluating 4 clusters
Evaluating 5 clusters
Evaluating 6 clusters
Evaluating 7 clusters
Evaluating 8 clusters
Evaluating 9 clusters
Evaluating 10 clusters
Evaluating 11 clusters
Evaluating 12 clusters
Evaluating 13 clusters
Evaluating 14 clusters
Evaluating 15 clusters
Evaluating 16 clusters
Evaluating 17 clusters
Evaluating 18 clusters
Evaluating 19 clusters
Evaluating 20 clusters
Evaluating 21 clusters
Evaluating 22 clusters
Evaluating 23 clusters
Evaluating 24 clusters
Evaluating 25 clusters
Evaluating 26 clusters
Evaluating 27 clusters
Evaluating 28 clusters
Evaluating 29 clusters

In [76]: # display the results
plt.plot(range(2, 30), scores_norm)
plt.scatter(range(2, 30), scores_norm)
plt.grid()
_ = plt.xticks(range(2, 30))

In [77]: # define and fit our model with k=12
model_norm = KNeighborsRegressor(n_neighbors=21, n_jobs=-1)
model_norm.fit(X_tr_norm, y_train)

# gather the predictions that our model made for our test set
preds_norm = model_norm.predict(X_te_norm)

# display the actuals and predictions for the test set
print('Actuals for test data set')
print(y_test)
print('Predictions for test data set')
print(preds_norm)

Actuals for test data set
[0 0 ... 0 0 0]
Predictions for test data set
[0. 0. 0.0952381 ... 0. 0. 0. ]

In [78]: differs_norm = y_test - preds_norm

print(f'Differences between the two sets:\n{differs_norm}')

print(f'r2_score: {r2_score(y_test, preds_norm)}')

Differences between the two sets:
[0. 0. 0. -0.0952381 ... 0. 0. 0. ]
r2_score: 0.19855738338320994

In [79]: # confusion matrix
from sklearn.metrics import confusion_matrix

In [80]: confusion_matrix(y_test, preds)

ValueError                                Traceback (most recent call last)
<ipython-input-80-d8a7f0b74565> in <module>
----> 1 confusion_matrix(y_test, preds)

C:\ProgramData\Anaconda3\Lib\site-packages\sklearn\utils\validation.py in inner_f(*args, **kwargs)
    62         if extra_args <= 0:
----> 63             return f(*args, **kwargs)
    64
    65         # extra_args > 0

C:\ProgramData\Anaconda3\Lib\site-packages\sklearn\metrics\classification.py in confusion_matrix(y_true, y_pred, labels, sample_weight, normalize)
    295     """
--> 296     y_type, y_true, y_pred = check_targets(y_true, y_pred)
    297     if y_type not in ('binary', 'multiclass'):
    298         raise ValueError("%s not supported" % y_type)

C:\ProgramData\Anaconda3\Lib\site-packages\sklearn\metrics\classification.py in check_targets(y_true, y_pred)
    90         if len(y_type) > 1:
----> 91             raise ValueError("Classification metrics can't handle a mix of (%s) "
    92                               "and (%s) targets." % format(y_true, y_type))
    93
    94
ValueError: Classification metrics can't handle a mix of binary and continuous targets
```