

Data-X Spring 2019: Homework 06

Name :

Anish Saha

SID :

26071616

Course (IEOR 135/290) :

Machine Learning

In this homework, you will do some exercises with prediction. We will cover these algorithms in class, but this is for you to have some hands on with these in scikit-learn. You can refer - <https://github.com/ikhlaqsidhu/data-x/blob/master/05a-tools-prediction-titanic/titanic.ipynb> (<https://github.com/ikhlaqsidhu/data-x/blob/master/05a-tools-prediction-titanic/titanic.ipynb>).

Display all your outputs.

```
In [1]: import numpy as np
import pandas as pd
```

```
In [2]: # machine learning libraries
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.linear_model import Perceptron
from sklearn.tree import DecisionTreeClassifier
```

1. Read `diabetesdata.csv` file into a pandas dataframe. About the data:

1. **TimesPregnant**: Number of times pregnant
2. **glucoseLevel**: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
3. **BP**: Diastolic blood pressure (mm Hg)
4. **insulin**: 2-Hour serum insulin (mu U/ml)
5. **BMI**: Body mass index (weight in kg/(height in m)²)
6. **pedigree**: Diabetes pedigree function
7. **Age**: Age (years)
8. **IsDiabetic**: 0 if not diabetic or 1 if diabetic)

```
In [3]: #Read data & print the head
df = pd.read_csv("diabetesdata.csv")
df.head()
```

Out[3]:

	TimesPregnant	glucoseLevel	BP	insulin	BMI	Pedigree	Age	IsDiabetic
0	6	148.0	72	0	33.6	0.627	50.0	1
1	1	NaN	66	0	26.6	0.351	31.0	0
2	8	183.0	64	0	23.3	0.672	NaN	1
3	1	NaN	66	94	28.1	0.167	21.0	0
4	0	137.0	40	168	43.1	2.288	33.0	1

2. Calculate the percentage of Null values in each column and display it.

```
In [4]: df.isna().sum() / len(df)
```

```
Out[4]: TimesPregnant    0.000000
glucoseLevel    0.04271
BP    0.000000
insulin    0.000000
BMI    0.000000
Pedigree    0.000000
Age    0.042969
IsDiabetic    0.000000
dtype: float64
```

3. Split data into train_df and test_df with 15% as test.

```
In [5]: np.random.seed(999)

# idx = np.random.rand(len(df)) < 0.85
# train_df, test_df = df[idx], df[~idx]
train_df = df.sample(frac=0.85, random_state=999)
test_df = df.drop(train_df.index)
len(df), len(train_df), len(test_df)
```

Out[5]: (768, 653, 115)

4. Display the means of the features in train and test sets. Replace the null values in train_df and test_df with the mean of EACH feature column separately for train and test. Display head of the dataframes.

```
In [6]: print(train_df.mean())
        print(test_df.mean())

        train_df.fillna(train_df.mean(), inplace=True)
        test_df.fillna(test_df.mean(), inplace=True)

        print("\nTrain")
        print(train_df.head())
        print("\nTest")
        print(test_df.head())
```

```

TimesPregnant      3.834609
glucoseLevel       120.444623
BP                 68.758040
insulin            77.165391
BMI                31.970904
Pedigree           0.474689
Age                33.494382
IsDiabetic         0.349158
dtype: float64
TimesPregnant      3.904348
glucoseLevel       124.225225
BP                 71.078261
insulin            94.756522
BMI                32.115652
Pedigree           0.455904
Age                32.571429
IsDiabetic         0.347826
dtype: float64

```

Train

	TimesPregnant	glucoseLevel	BP	insulin	BMI	Pedigree	Age \
441	2	83.0	66	50	32.2	0.497	22.0
57	0	100.0	88	110	46.8	0.962	31.0
68	1	95.0	66	38	19.6	0.334	25.0
95	6	144.0	72	228	33.9	0.255	40.0
411	1	112.0	72	176	34.4	0.528	25.0

	IsDiabetic
441	0
57	0
68	0
95	0
411	0

Test

	TimesPregnant	glucoseLevel	BP	insulin	BMI	Pedigree	Age
2	8	183.000000	64	0	23.3	0.672	32.571429
11	10	168.000000	74	0	38.0	0.537	34.000000
12	10	139.000000	80	0	27.1	1.441	57.000000
14	5	166.000000	72	175	25.8	0.587	51.000000
16	0	124.225225	84	230	45.8	0.551	31.000000

	IsDiabetic
2	1
11	1
12	0
14	1
16	1

5. Split train_df & test_df into X_train, Y_train and X_test, Y_test. Y_train and Y_test should only have the column we are trying to predict, IsDiabetic.

```
In [7]: X_train, X_test = train_df.drop("IsDiabetic", axis=1), test_df.drop("IsDiabetic", axis=1)
y_train, y_test = train_df["IsDiabetic"], test_df["IsDiabetic"]
print("X_train")
print(X_train.head())
print("\nX_test")
print(X_test.head())
print("\ny_train")
print(y_train.head())
print("\ny_test")
print(y_test.head())
```

X_train

	TimesPregnant	glucoseLevel	BP	insulin	BMI	Pedigree	Age
441	2	83.0	66	50	32.2	0.497	22.0
57	0	100.0	88	110	46.8	0.962	31.0
68	1	95.0	66	38	19.6	0.334	25.0
95	6	144.0	72	228	33.9	0.255	40.0
411	1	112.0	72	176	34.4	0.528	25.0

X_test

	TimesPregnant	glucoseLevel	BP	insulin	BMI	Pedigree	Age
2	8	183.000000	64	0	23.3	0.672	32.571429
11	10	168.000000	74	0	38.0	0.537	34.000000
12	10	139.000000	80	0	27.1	1.441	57.000000
14	5	166.000000	72	175	25.8	0.587	51.000000
16	0	124.225225	84	230	45.8	0.551	31.000000

y_train

441	0
57	0
68	0
95	0
411	0

Name: IsDiabetic, dtype: int64

y_test

2	1
11	1
12	0
14	1
16	1

Name: IsDiabetic, dtype: int64

6. Use this dataset to train perceptron, logistic regression and random forest models using 15% test split. Report training and test accuracies. Try different hyperparameter values for these models and see if you can improve your accuracies.

```
In [8]: from sklearn.metrics import accuracy_score, confusion_matrix
        from sklearn.model_selection import GridSearchCV

        # 6a. Logistic Regression
        mod1 = LogisticRegression()
        mod1.fit(X_train, y_train)

        print("Logistic Regression Model Performance:")
        y_pred_train = mod1.predict(X_train)
        train_rmse = accuracy_score(y_train, y_pred_train)
        print("Training Accuracy: " + str(train_rmse))
        y_pred_test = mod1.predict(X_test)
        test_rmse = accuracy_score(y_test, y_pred_test)
        print("Test Accuracy: " + str(test_rmse))

        print("\n")

        penalty = ['l1', 'l2']
        C = np.logspace(0, 5, 10)
        hyperparameters = dict(C=C, penalty=penalty)
        clf = GridSearchCV(mod1, hyperparameters, cv=5, verbose=0)
        mod2 = clf.fit(X_train, y_train) # optimized hyperparameters

        print("Optimized Logistic Regression Model Performance:")
        y_pred_train = mod2.predict(X_train)
        train_rmse = accuracy_score(y_train, y_pred_train)
        print("Training Accuracy: " + str(train_rmse))
        y_pred_test = mod2.predict(X_test)
        test_rmse = accuracy_score(y_test, y_pred_test)
        print("Test Accuracy: " + str(test_rmse))
```

Logistic Regression Model Performance:

Training Accuracy: 0.7733537519142419

Test Accuracy: 0.7565217391304347

Optimized Logistic Regression Model Performance:

Training Accuracy: 0.77947932618683

Test Accuracy: 0.782608695652174

```

In [9]: from sklearn.neural_network import MLPClassifier

# 6b. Perceptron
mod3 = Perceptron()
mod3.fit(X_train, y_train)

print("Perceptron Model Performance:")
y_pred_train = mod3.predict(X_train)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod3.predict(X_test)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))

hyperparameters = { 'alpha': [0.0001, 0.05], 'fit_intercept': [True, False],
                    'max_iter': [100, 1000], 'penalty': penalty }
clf = GridSearchCV(mod3, hyperparameters, cv=5, verbose=0)
mod4 = clf.fit(X_train, y_train) # optimized hyperparameters

print("Optimized Perceptron Model Performance:")
y_pred_train = mod4.predict(X_train)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod4.predict(X_test)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))

print("\n")

# 6b. Multi-Layer Perceptron
mod3 = MLPClassifier()
mod3.fit(X_train, y_train)

print("Multi-Layer Perceptron Model Performance:")
y_pred_train = mod3.predict(X_train)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod3.predict(X_test)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))

print("\n")

hyperparameters = { 'hidden_layer_sizes': [(50,50,50), (50,100,50), (100,
)],
                    'alpha': [0.0001, 0.05], 'activation': ['tanh', 'relu'],
                    'learning_rate': ['constant', 'adaptive'] }
clf = GridSearchCV(mod3, hyperparameters, cv=5, verbose=0)
mod4 = clf.fit(X_train, y_train) # optimized hyperparameters

print("Optimized Multi-Layer Perceptron Model Performance:")
y_pred_train = mod4.predict(X_train)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod4.predict(X_test)

```

```
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))
```

Perceptron Model Performance:

Training Accuracy: 0.6447166921898928

Test Accuracy: 0.6434782608695652

/Library/Frameworks/Python.framework/Versions/3.6/lib/python3.6/site-packages/sklearn/linear_model/stochastic_gradient.py:128: FutureWarning: max_iter and tol parameters have been added in <class 'sklearn.linear_model.perceptron.Perceptron'> in 0.19. If both are left unset, they default to max_iter=5 and tol=None. If tol is not None, max_iter defaults to max_iter=1000. From 0.21, default max_iter will be 1000, and default tol will be 1e-3.

"and default tol will be 1e-3." % type(self), FutureWarning)

Optimized Perceptron Model Performance:

Training Accuracy: 0.6477794793261868

Test Accuracy: 0.6521739130434783

Multi-Layer Perceptron Model Performance:

Training Accuracy: 0.6707503828483921

Test Accuracy: 0.6869565217391305

Optimized Multi-Layer Perceptron Model Performance:

Training Accuracy: 0.7304747320061256

Test Accuracy: 0.6869565217391305


```

In [10]: # 6c. Random Forest
mod5 = RandomForestClassifier()
mod5.fit(X_train, y_train)

print("Random Forest Model Performance:")
y_pred_train = mod5.predict(X_train)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod5.predict(X_test)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))

print("\n")

max_depth = [int(x) for x in np.linspace(10, 110, num = 5)]
max_depth.append(None)
hyperparameters = { 'max_depth': max_depth, 'min_samples_split': [2, 5,
10],
                    'max_features': ['auto', 'sqrt'], 'bootstrap': [True
, False] }
clf = GridSearchCV(mod5, hyperparameters, cv=5, verbose=0)
mod6 = clf.fit(X_train, y_train) # optimized hyperparameters

print("Optimized Random Forest Model Performance:")
y_pred_train = mod6.predict(X_train)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod6.predict(X_test)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))

```

Random Forest Model Performance:
Training Accuracy: 0.9862174578866769
Test Accuracy: 0.7652173913043478

Optimized Random Forest Model Performance:
Training Accuracy: 0.9831546707503829
Test Accuracy: 0.7913043478260869

7. For your logistic regression model -

a . Compute the log probability of classes in `IsDiabetic` for the first 10 samples of your train set and display it. Also display the predicted class for those samples from your logistic regression model trained before.

```
In [11]: print("Log Probabilities for first 10 training samples")
print(mod2.predict_log_proba(X_train)[:10])
print("\n")
print("Predicted Class for first 10 training samples")
print(mod2.predict(X_train)[:10])
```

Log Probabilities for first 10 training samples

```
[[-0.0795963 -2.57032189]
 [-0.45018631 -1.01475664]
 [-0.03334802 -3.41738442]
 [-0.76194222 -0.62878179]
 [-0.20794316 -1.67266106]
 [-1.13052689 -0.38988181]
 [-0.04600159 -3.10199196]
 [-0.12242993 -2.16080692]
 [-0.10207376 -2.33266235]
 [-0.03616457 -3.33770305]]
```

Predicted Class for first 10 training samples

```
[0 0 0 1 0 1 0 0 0 0]
```

b . Now compute the log probability of classes in `IsDiabetic` for the first 10 samples of your test set and display it. Also display the predicted class for those samples from your logistic regression model trained before. (using the model trained on the training set)

```
In [12]: print("Log Probabilities for first 10 training samples")
print(mod2.predict_log_proba(X_test)[:10])
print("\n")
print("Predicted Class for first 10 training samples")
print(mod2.predict(X_test)[:10])
```

Log Probabilities for first 10 training samples

```
[[-1.63359437 -0.21719453]
 [-2.18604515 -0.11918929]
 [-1.74734353 -0.19144648]
 [-1.11840579 -0.39571336]
 [-0.56430783 -0.84107586]
 [-2.67552943 -0.07135676]
 [-0.3046589 -1.33702754]
 [-2.66660708 -0.07201987]
 [-0.43038092 -1.05056912]
 [-0.06876545 -2.71123949]]
```

Predicted Class for first 10 training samples

```
[1 1 1 1 0 1 0 1 0 0]
```

c . What can you interpret from the log probabilities and the predicted classes?

In the outputs above, the first column represents the log probability that the sample is of class 0 [`IsDiabetic = 0`], while the second column represents the log probability that the sample is of class 1 [`IsDiabetic = 1`]. The probability that a sample is of a certain class is computed using the formula:

$$P(\text{sample}_j \text{ is not diabetic}) = e^{a[j][0]} \mid P(\text{sample}_j \text{ is diabetic}) = e^{a[j][1]}$$

for the j^{th} sample, and a is the array displayed above

The predicted class corresponds to the column with the higher log probability (and consequently, higher probability) value – or in other words, whichever log probability value is closer to 0 since all log probability values are negative. This can be confirmed by observing the outputs above.

8. Is mean imputation is the best type of imputation (as we did in 4.) to use? Why or why not? What are some other ways to impute the data?

Mean imputation is not the best type of imputation to use. This is because it often does not preserve relationships between variables (imputed values have zero correlation with other variables), presents biased metrics of standard error and variance, and can result in a biased sample mean. The only advantage is that it preserves the sample size. Some other ways to impute the data include hot-deck imputation, cold-deck imputation, regression imputation, and multiple imputation (ex: MICE, using chained equations, for when data is randomly missing).

Extra Credit (2 pts) - MANDATORY for students enrolled in IEOR 290

9. Implement the K-Nearest Neighbours (https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm) (https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm) algorithm for $k=1$ from scratch in python (do not use KNN from existing libraries). KNN uses Euclidean distance to find nearest neighbors. Split your dataset into test and train as before. Also fill in the null values with mean of features as done earlier. Use this algorithm to predict values for 'IsDiabetic' for your test set. Display your accuracy.

```
In [13]: # K-Nearest Neighbors Classifier for k=1
class KNN_Classifier():
    def fit(self, X_train, y_train):
        self.X_train = X_train
        self.y_train = y_train

    def euclidean_dist(self, x1, x2):
        distance = 0
        for i in range(len(x1)):
            distance = distance + (x1[i] - x2[i])**2
        return distance

    def k_nearest(self, row, k=1):
        best_dist = self.euclidean_dist(row, self.X_train[0])
        best_idx = 0
        for i in range(k, len(self.X_train)):
            dist = self.euclidean_dist(row, self.X_train[i])
            if dist < best_dist:
                best_dist = dist
                best_idx = i
        return self.y_train[best_idx]

    def predict(self, X_test, k=1):
        result = []
        for row in X_test:
            label = self.k_nearest(row, k)
            result.append(label)
        return result

mod7 = KNN_Classifier()
mod7.fit(X_train.values, y_train.values)

print("MANUAL IMPLEMENTATION | K-Nearest Neighbors Model Performance:")
y_pred_train = mod7.predict(X_train.values)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod7.predict(X_test.values)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))
```

```
MANUAL IMPLEMENTATION | K-Nearest Neighbors Model Performance:
Training Accuracy: 1.0
Test Accuracy: 0.6956521739130435
```

```
In [14]: # Checking Implementation against SciKitLearn Library Implementation
from sklearn.neighbors import KNeighborsClassifier

mod8 = KNeighborsClassifier(n_neighbors=1)
mod8.fit(X_train, y_train)

print("SKLEARN IMPLEMENTATION | K-Nearest Neighbors Model Performance:")
y_pred_train = mod8.predict(X_train.values)
train_rmse = accuracy_score(y_train, y_pred_train)
print("Training Accuracy: " + str(train_rmse))
y_pred_test = mod8.predict(X_test.values)
test_rmse = accuracy_score(y_test, y_pred_test)
print("Test Accuracy: " + str(test_rmse))
```

```
SKLEARN IMPLEMENTATION | K-Nearest Neighbors Model Performance:
Training Accuracy: 1.0
Test Accuracy: 0.6956521739130435
```