

Task-D: Collinear features and their effect on linear models

```
In [11]: %matplotlib inline
import warnings
warnings.filterwarnings('ignore')
import pandas as pd
import numpy as np
from sklearn.datasets import load_iris
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import GridSearchCV
import seaborn as sns
import matplotlib.pyplot as plt

In [12]: data = pd.read_csv('task_d.csv')

In [29]: data.head()

Out[29]:   x     y     z   x*x   2*y   2*z+3*x   w  target
0 -0.581066  0.841837 -1.012978 -0.604025  0.841837 -0.665927 -0.536277  0
1 -0.894309 -0.207835 -1.012978 -0.883052 -0.207835 -0.917054 -0.522364  0
2 -1.207552  0.212034 -1.082312 -1.150918  0.212034 -1.166507  0.205738  0
3 -1.364174  0.002099 -0.943643 -1.280666  0.002099 -1.266540 -0.665720  0
4 -0.737687  1.051772 -1.012978 -0.744934  1.051772 -0.792746 -0.735054  0
```

```
In [38]: X = data.drop(['target'], axis=1).values
Y = data['target'].values
```

Doing perturbation test to check the presence of collinearity

Task: 1 Logistic Regression

1. Finding the Correlation between the features
 - a. check the correlation between the features
 - b. plot heat map of correlation matrix using seaborn heatmap
2. Finding the best model for the given data
 - a. Train Logistic regression on data(X,Y) that we have created in the above cell
 - b. Find the best hyper parameter alpha with hyper parameter tuning using k-fold cross validation (grid search CV or random search CV make sure you choose the alpha in log space)
 - c. Create a new Logistic regression with the best alpha (search for how to get the best hyper parameter value), name the best model as 'best_model'
3. Getting the weights with the original data
 - a. train the 'best_model' with X, Y
 - b. Check the accuracy of the model 'best_model_accuracy'
 - c. Get the weights W using best_model.coef_
4. Modifying original data
 - a. Add a noise(order of 10^{-2}) to each element of X
 - b. Train the same 'best_model' with data (X', Y)
 - c. Check the accuracy of the model 'best_model_accuracy_edited'
 - d. Get the weights W' using best_model.coef_
5. Checking deviations in metric and weights
 - a. find the difference between 'best_model_accuracy_edited' and 'best_model_accuracy'
 - b. find the absolute change between each value of W and W' $\Rightarrow |(W-W')|$
 - c. print the top 4 features which have higher % change in weights compare to the other feature

Task: 2 Linear SVM

1. Do the same steps (2, 3, 4, 5) we have done in the above task 1.

Do write the observations based on the results you get from the deviations of weights in both Logistic Regression and linear SVM

TASK 1. LOGISTIC REGRESSION.

1. Finding the Correlation between the features

- a. check the correlation between the features

b. plot heat map of correlation matrix using seaborn heatmap

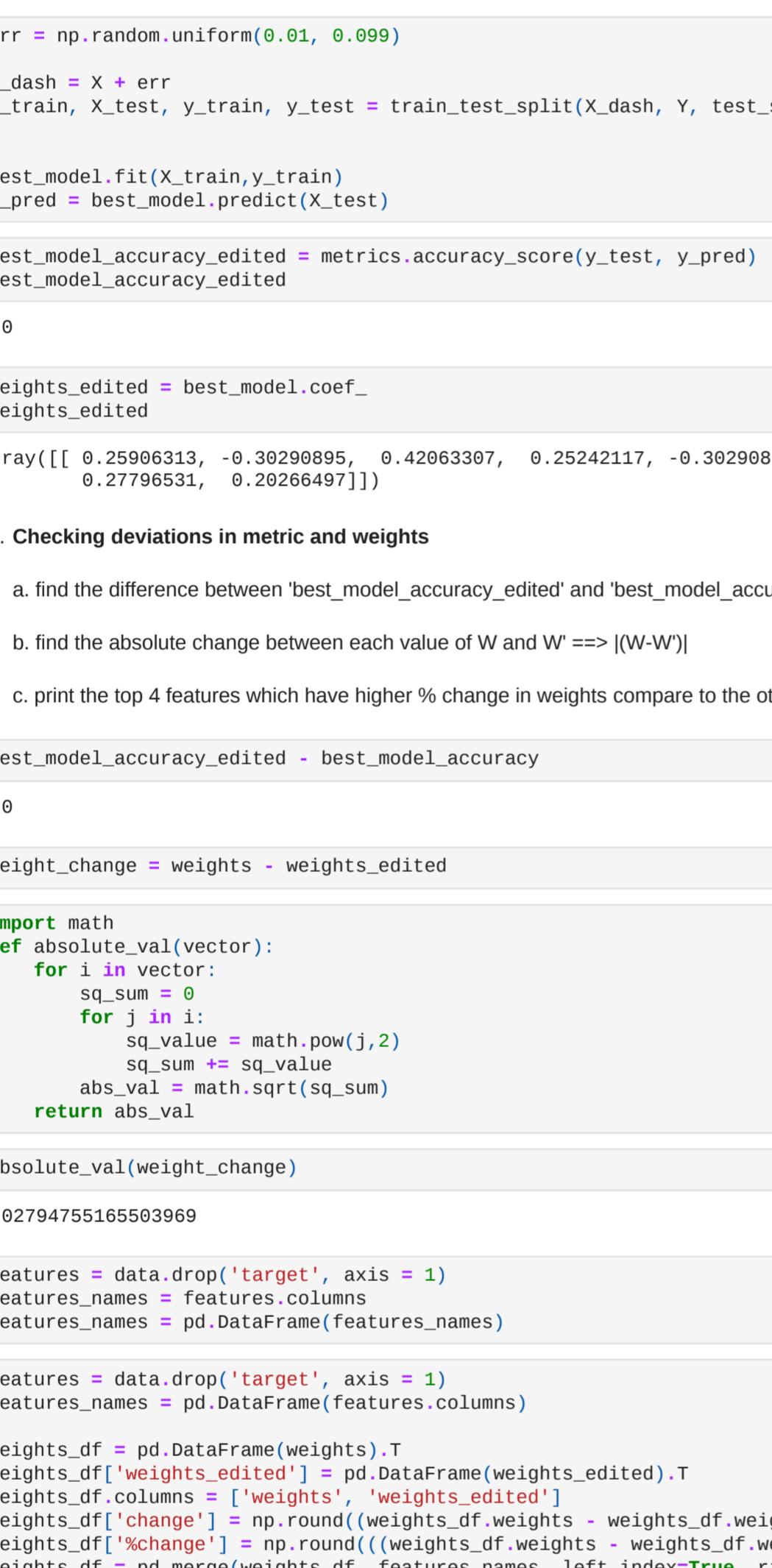
```
In [33]: X_df = data.drop(['target'], axis=1)
```

```
covariance = pd.DataFrame(X_df.cov())
```

```
plt.figure(figsize=(12,8))
sns.heatmap(covariance, annot=True, fmt="f", cmap="YlGnBu")
plt.plot()
```

```
...
spearman_cov = data.corr(method='spearman')
plt.figure(figsize=(12,8))
sns.heatmap(spearman_cov, annot=True, fmt="f", cmap="YlGnBu")
plt.plot()
```

```
Out[65]: 'nspearman_cov = data.corr(method=\''spearman'\')\nplt.figure(figsize=(12,8))\nnsns.heatmap(spearman_cov, annot=True, fmt="f", cmap="YlGnBu")\nplt.plot()\n'
```



1. Finding the best model for the given data

- a. Train Logistic regression on data(X,Y) that we have created in the above cell
- b. Find the best hyper parameter alpha with hyper parameter tuning using k-fold cross validation (grid search CV or random search CV make sure you choose the alpha in log space)

c. Create a new Logistic regression with the best alpha (search for how to get the best hyper parameter value), name the best model as 'best_model'

```
In [66]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import metrics
```

```
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size = 0.33, stratify = Y)
```

```
lr = LogisticRegression()
```

```
alphas = np.logspace(-2,3,10, endpoint=True, dtype='float')
param_grid = {'penalty':['l1','l2'],
             'C': np.logspace(-2,3,10, endpoint=True, dtype='float')}
```

```
clf = GridSearchCV(lr, param_grid=param_grid, cv=3)
```

```
clf.fit(X_train,y_train)
```

```
print("Best Parameters Are: ",clf.best_params_)
```

```
Best Parameters Are: { 'C': 0.03593813663804628, 'penalty': 'l2'}
```

```
In [67]: best_model = clf.best_estimator_
best_model
```

```
Out[67]: LogisticRegression(C=0.03593813663804628)
```

1. Getting the weights with the original data

- a. train the 'best_model' with X, Y

b. Check the accuracy of the model 'best_model_accuracy'

c. Get the weights W using best_model.coef_

```
In [68]: best_model.fit(X_train,y_train)
```

```
y_pred = best_model.predict(X_test)
```

```
best_model_accuracy = metrics.accuracy_score(y_test, y_pred)
```

```
best_model_accuracy
```

```
Out[68]: 1.0
```

```
In [69]: weights = best_model.coef_
weights
```

```
Out[69]: array([[ 0.25141997, -0.31088651,  0.40387112,  0.24372977, -0.31088651,
```

```
           0.26811928,  0.21457786]])
```

1. Modifying original data

- a. Add a noise(order of 10^{-2}) to each element of X and get the new data set X' ($X' = X + e$)

b. Train the same 'best_model' with data (X', Y)

c. Check the accuracy of the model 'best_model_accuracy_edited'

d. Get the weights W' using bestmodel.coef_

```
In [70]: err = np.random.uniform(0.01, 0.099)
```

```
X_dash = X + err
X_train, X_test, y_train, y_test = train_test_split(X_dash, Y, test_size = 0.33, stratify = Y)
```

```
best_model.fit(X_train,y_train)
```

```
y_pred = best_model.predict(X_test)
```

```
best_model_accuracy_edited = metrics.accuracy_score(y_test, y_pred)
```

```
best_model_accuracy_edited
```

```
Out[68]: 1.0
```

```
In [69]: weights_edited = best_model.coef_
weights_edited
```

```
Out[69]: array([[ 0.25090631, -0.30290985,  0.42063307,  0.25242117, -0.30290985,
```

```
           0.27798531,  0.20266497]])
```

1. Checking deviations in metric and weights

- a. find the difference between 'best_model_accuracy_edited' and 'best_model_accuracy'

b. find the absolute change between each value of W and W' $\Rightarrow |(W-W')|$

c. print the top 4 features which have higher % change in weights compare to the other feature

```
In [73]: best_model_accuracy_edited - best_model_accuracy
```

```
Out[73]: 0.0
```

```
Out[61]: weight_change = weights - weights_edited
```

```
Out[61]: 0.030630891044988456
```

```
In [62]: features = data.drop('target', axis = 1)
features_names = pd.DataFrame(features.columns)
```

```
weights_df = pd.DataFrame(weights.T)
weights_df['weights_edited'] = pd.DataFrame(weights_edited.T)
```

```
weights_df['change'] = np.round((weights_df.weights - weights_df.weights_edited), 6)
```

```
weights_df['df'] = pd.merge(weights_df, features_names, left_index=True, right_index=True)
```

```
weights_df['columns'] = ['weights', 'weights_edited', 'change', 'Rank']
```

```
weights_df
```

feature weights weights_edited change %change Rank

```
0 x 0.251420 0.250906 -0.007643 -0.039998 4.0
```

```
1 y -0.310887 -0.302909 -0.007978 2.566070 5.5
```

```
2 z 0.403871 0.420633 -0.016762 -4.150322 1.0
```

```
3 x*x 0.243730 0.252421 -0.008091 -3.565998 3.0
```

```
4 2*y -0.310887 -0.302909 -0.009846 2.566070 5.5
```

```
5 2*z+3*x 0.268119 -0.277985 -0.007179 4.321027 5.0
```

```
6 w 0.214578 0.202665 0.011913 5.551778 7.0
```

Top features with highest percentage change are:

```
In [49]: top_features = weights_df.nlargest(4, 'Rank')
```

```
Out[49]: feature weights weights_edited change %change Rank
```

```
3 x*x 0.145937 0.138740 0.001974 4.93167 7.0
```

```
0 x 0.147764 0.140835 0.006925 4.689193 6.0
```

```
5 2*z+3*x 0.156498 0.148769 0.006719 4.321027 5.0
```

```
2 z 0.199976 0.189796 0.002000 1.000136 4.0
```

```
In [50]: print("Top features with highest percentage change are: {}".format(list(top_features.feature)))
```

Top features with highest percentage change are: ['x*x', 'x', '2*z+3*x', 'z']

OBSERVATIONS:

- From Covariance Matrix, it can be observed that some of the features are highly correlated as their correlation value is high. Example (y & 2y), (z, 2z+3xx) (XX, X) etc.
- After performing perturbation test top features with highest percentage change are: [x*x, x', 2z+3xx', z]. Which indicates that these features are collinear/multilinear.

- As we can express these features in terms of other in linear fashion i.e. we can establish linear relationship between these features. E.g. X^2X is a square of X, their weights can be arbitrarily changed hence we cannot used their weights directly as a feature importance, therefore other methods such as "Forward Selection Method" can be recommended in order to decide feature importance.

Task: 2 Linear SVM

1. Do the same steps (2, 3, 4, 5) we have done in the above task 1.

```
In [51]: from sklearn.svm import SVC
```

```
clf_svm = SVC(kernel='linear')
```

```
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size = 0.33, stratify = Y)
```

```
clf_svm = SVC(kernel='linear')
```

```
alphas = np.logspace(-2,3,10, endpoint=True, dtype='float')
```

```
param_grid = {'penalty':['l1','l2'],
             'C': np.logspace(-2,3,10, endpoint=True, dtype='float')}
```

```
clf = GridSearchCV(clf_svm, param_grid=param_grid, cv=3)
```

```
clf.fit(X_train,y_train)
```

```
print("Best Parameters Are: ",clf.best_params_)
```

```
Best Parameters Are: { 'C': 0.03593813663804628, 'penalty': 'l2'}
```

```
In [52]: best_model = clf.best_estimator_
best_model
```

```
Out[52]: SVC(C=0.03593813663804628, kernel='linear')
```

```
In [54]: best_model.fit(X_train,y_train)
```

```
y_pred = best_model.predict(X_test)
```

```
best_model_accuracy = metrics.accuracy_score(y_test, y_pred)
```

```
best_model_accuracy
```

```
Out[54]: 1.0
```

```
In [69]: weights = best_model.coef_
weights
```

```
Out[69]: array([[ 0.16074191, -0.31088651,  0.40387112,  0.24372977, -0.31088651,
```

```
           0.26811928,  0.21457786]])
```

1. Modifying original data

- a. Add a noise(order of 10^{-2}) to each element of X and get the new data set X' ($X' = X + e$)

b. Train the same 'best_model' with data (X', Y)

c. Check the accuracy of the model 'best_model_accuracy_edited'

d. Get the weights W' using bestmodel.coef_

```
In [70]: err = np.random.uniform(0.01, 0.099)
```

```
X_dash = X + err
X_train, X_test, y_train, y_test = train_test_split(X_dash, Y, test_size = 0.33, stratify = Y)
```

```
best_model.fit(X_train,y_train)
```

```
y_pred = best_model.predict(X_test)
```