

Lab Sessions Day 3

Exercise 1

Linear regression

#Make two vector X and y

```
X=np.array([1,2,4,3,5])
```

```
y=np.array([1,3,3,2,5])
```

#With simple linear regression we want to model our data as follows:

$y = B_0 + B_1 * x$

#We can start off by estimating the value for B1 as:

$B_1 = \frac{\sum((X_i - \text{mean}(X)) * (y_i - \text{mean}(y)))}{\sum((X_i - \text{mean}(X))^2)}$

```

# Create a scatter plot of X vs y
plt.scatter(X, y)

# Calculate the mean of X and y
mean_X = np.mean(X)
mean_y = np.mean(y)

# Calculate the numerator and denominator for B1
numerator = 0
denominator = 0

for i in range(len(X)):
    numerator += (X[i] - mean_X) * (y[i] - mean_y)
    denominator += (X[i] - mean_X)**2

# Calculate B1
B1 = numerator / denominator

# Print B1
print(B1)
```

#We can calculate B0 using B1 and some statistics from our dataset, as follows:

$B_0 = \text{mean}(y) - B_1 * \text{mean}(X)$

```

# Calculate B0
B0 = mean_y - B1 * mean_X

# Print B0
print(B0)
```

#Making Predictions (y_{hat} is a predicted y)

$y_{\text{hat}} = B_0 + B_1 * X$

```

# Making predictions
y_hat = B0 + B1 * X

# Print y_hat
print(y_hat)
```

#Evaluation

RMSE = sqrt(sum((y_hat_i - yi)^2)/n)



Exercise 2

Logistic regression

```
from sklearn.linear_model import LogisticRegression
```

```
from sklearn.datasets import load_breast_cancer
```

```
from sklearn.model_selection import train_test_split
```

```
cancer=load_breast_cancer()
```

```
X_train,X_test,y_train,y_test=train_test_split(cancer.data,cancer.target, stratify=cancer.target, random_state=42)
```

```
#####default C=1#####
```

```
lgr=LogisticRegression().fit(X_train,y_train)
```

```
print("training set score: %f" % lgr.score(X_train, y_train))
```

```
print("\n"test set score: %f" % lgr.score(X_test, y_test))
```

```
#####increase C to 100#####
```

```
lgr100=LogisticRegression(C=100).fit(X_train,y_train)
```

```
print("\n"training set score of lgr100: %f" % lgr100.score(X_train, y_train))
```

```
print("\n"test set score of lgr100: %f" % lgr100.score(X_test, y_test))
```

Change C value and compare the performance metric

```
#####decrease C to 0.01#####
```

```
lgr001=LogisticRegression(C=0.01).fit(X_train,y_train)
```

```
print("\n"training set score of lgr001: %f" % lgr001.score(X_train, y_train))
```

```
print("\n"test set score of lgr001: %f" % lgr001.score(X_test, y_test))
```

```
import matplotlib.pyplot as plt
```

```

plt.plot(lgr.coef_.T,'o',label='C=1')
plt.plot(lgr100.coef_.T,'+',label='C=100')
plt.plot(lgr001.coef_.T,'-',label='C=0.01')
plt.xticks(range(cancer.data.shape[1]),cancer.feature_names,rotation=90)
plt.ylim(-5,5)
plt.legend()
plt.show()

```

####If we desire a more interpretable model, using L1 regularization might help
 ####As LogisticRegression applies an L2 regularization by default, the result
 ####looks similar to Ridge in Figure ridge_coefficients. Stronger regularization
 ####pushes coefficients more and more towards zero, though coefficients never
 ####become exactly zero.

```

import numpy as np
import math
n=np.arange(-2,3)
print(n)
r=pow(float(10),n)
print(r)
for C in r:
    lr_l1=LogisticRegression(C=C,penalty="l1").fit(X_train,y_train)
    print("\n""Training Accuracy of L1 LogRegress with C=%f: %f"%(C,lr_l1.score(X_train,y_train)))
    print("\n""Test Accuracy of L1 LogRegress with C=%f: %f"%(C,lr_l1.score(X_test,y_test)))
    plt.plot(lr_l1.coef_.T,'o',label="C=%f"%C)
plt.xticks(range(cancer.data.shape[1]),cancer.feature_names,rotation=90)
plt.ylim(-5,5)
plt.legend(loc='best')
plt.show()

```