

svm-gd-21mis1152

February 28, 2024

21MIS1152 Rajeev Sekar

SVM with gradient descent

```
[25]: import numpy as np
import pandas as pd
import statsmodels.api as sm
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split as tts
from sklearn.metrics import accuracy_score, recall_score, precision_score
from sklearn.utils import shuffle
```

Functions for feature selection

```
[26]: def remove_correlated_features(X):
    corr_threshold = 0.9
    corr = X.corr()
    drop_columns = np.full(corr.shape[0], False, dtype=bool)
    for i in range(corr.shape[0]):
        for j in range(i + 1, corr.shape[0]):
            if corr.iloc[i, j] >= corr_threshold:
                drop_columns[j] = True
    columns_dropped = X.columns[drop_columns]
    X.drop(columns_dropped, axis=1, inplace=True)
    return columns_dropped

def remove_less_significant_features(X, Y):
    sl = 0.05
    regression_ols = None
    columns_dropped = np.array([])
    for itr in range(0, len(X.columns)):
        regression_ols = sm.OLS(Y, X).fit()
        max_col = regression_ols.pvalues.idxmax()
        max_val = regression_ols.pvalues.max()
        if max_val > sl:
            X.drop(max_col, axis='columns', inplace=True)
            columns_dropped = np.append(columns_dropped, [max_col])
        else:
```

```

        break
    regression_ols.summary()
    return columns_dropped

```

MODEL TRAINING using SVM and SGD(stochastic gradient descent)

```

[27]: def compute_cost(W, X, Y):
    # calculate hinge loss
    N = X.shape[0]
    distances = 1 - Y * (np.dot(X, W))
    distances[distances < 0] = 0 # equivalent to max(0, distance)
    hinge_loss = regularization_strength * (np.sum(distances) / N)

    # calculate cost
    cost = 1 / 2 * np.dot(W, W) + hinge_loss
    return cost

def calculate_cost_gradient(W, X_batch, Y_batch):
    # if only one example is passed (eg. in case of SGD)
    if type(Y_batch) == np.float64:
        Y_batch = np.array([Y_batch])
        X_batch = np.array([X_batch]) # gives multidimensional array

    distance = 1 - (Y_batch * np.dot(X_batch, W))
    dw = np.zeros(len(W))

    for ind, d in enumerate(distance):
        if max(0, d) == 0:
            di = W
        else:
            di = W - (regularization_strength * Y_batch[ind] * X_batch[ind])
        dw += di

    dw = dw/len(Y_batch) # average
    return dw

def sgd(features, outputs):
    max_epochs = 5000
    weights = np.zeros(features.shape[1])
    nth = 0
    prev_cost = float("inf")
    cost_threshold = 0.01 # in percent
    # stochastic gradient descent
    for epoch in range(1, max_epochs):
        # shuffle to prevent repeating update cycles
        X, Y = shuffle(features, outputs)
        for ind, x in enumerate(X):

```

```

        ascent = calculate_cost_gradient(weights, x, Y[ind])
        weights = weights - (learning_rate * ascent)

# convergence check on 2nth epoch
    if epoch == 2 ** nth or epoch == max_epochs - 1:
        cost = compute_cost(weights, features, outputs)
        print("Epoch is: {} and Cost is: {}".format(epoch, cost))
        # stoppage criterion
        if abs(prev_cost - cost) < cost_threshold * prev_cost:
            return weights
        prev_cost = cost
        nth += 1
    return weights

```

[]: Initialising the model and iterating till we converge to the minimal cost
 ↪function (Checked for every 2^{nt} epoch)

```

[28]: def init():
    print("reading dataset...")
    # read data in pandas (pd) data frame
    data = pd.read_csv('data.csv')

    # drop last column (extra column added by pd)
    # and unnecessary first column (id)
    data.drop(data.columns[[-1, 0]], axis=1, inplace=True)

    print("applying feature engineering...")
    # convert categorical labels to numbers
    diag_map = {'M': 1.0, 'B': -1.0}
    data['diagnosis'] = data['diagnosis'].map(diag_map)

    # put features & outputs in different data frames
    Y = data.loc[:, 'diagnosis']
    X = data.iloc[:, 1:]

    # filter features
    remove_correlated_features(X)
    remove_less_significant_features(X, Y)

    # normalize data for better convergence and to prevent overflow
    X_normalized = MinMaxScaler().fit_transform(X.values)
    X = pd.DataFrame(X_normalized)

    # insert 1 in every row for intercept b
    X.insert(loc=len(X.columns), column='intercept', value=1)

    # split data into train and test set

```

```

print("splitting dataset into train and test sets...")
X_train, X_test, y_train, y_test = tts(X, Y, test_size=0.25)

# train the model
print("training started...")
W = sgd(X_train.to_numpy(), y_train.to_numpy())
print("training finished.")
print("weights are: {}".format(W))

# testing the model
print("testing the model...")
y_train_predicted = np.array([])
for i in range(X_train.shape[0]):
    yp = np.sign(np.dot(X_train.to_numpy()[i], W))
    y_train_predicted = np.append(y_train_predicted, yp)

y_test_predicted = np.array([])
for i in range(X_test.shape[0]):
    yp = np.sign(np.dot(X_test.to_numpy()[i], W))
    y_test_predicted = np.append(y_test_predicted, yp)

print("accuracy on test dataset: {}".format(accuracy_score(y_test,
↪y_test_predicted)))
print("recall on test dataset: {}".format(recall_score(y_test,
↪y_test_predicted)))
print("precision on test dataset: {}".format(recall_score(y_test,
↪y_test_predicted)))

# set hyper-parameters and call init
regularization_strength = 10000
learning_rate = 0.000001
init()

```

```

reading dataset...
applying feature engineering...
splitting dataset into train and test sets...
training started...
Epoch is: 1 and Cost is: 7109.194396019079
Epoch is: 2 and Cost is: 6460.123104376328
Epoch is: 4 and Cost is: 5445.841624060475
Epoch is: 8 and Cost is: 3863.6426428321656
Epoch is: 16 and Cost is: 2649.379971909295
Epoch is: 32 and Cost is: 1902.5057357266098
Epoch is: 64 and Cost is: 1474.7279371324562
Epoch is: 128 and Cost is: 1270.699283468124
Epoch is: 256 and Cost is: 1078.158074135148

```

```
Epoch is: 512 and Cost is: 962.9101023032541
Epoch is: 1024 and Cost is: 908.4769450273418
Epoch is: 2048 and Cost is: 894.1284662538296
Epoch is: 4096 and Cost is: 893.9173968433818
training finished.
weights are: [ 3.57181726 12.1286535 -2.929618 -10.29191262 10.08305002
-1.52580759 -7.0736469 2.19727649 -2.40074262 3.0096523
5.65326568 5.78504576 -4.59164315]
testing the model...
accuracy on test dataset: 0.9440559440559441
recall on test dataset: 0.8771929824561403
precision on test dataset: 0.8771929824561403
```

[]: