# Task -- Diagnostics

In the first homework, we had two issues with the classifiers we built. Namely (1) the data were not shuffled, and (2) the labels were highly imbalanced. Both of these made it difficult to effectively build an accurate classifier. Here we'll try and correct for those issues using the Bankruptcy dataset.

## Problem 1

Download and parse the bankruptcy data. We'll use the `5year.arff` file. Code to read the data is available in the stub. Train a logistic regressor (e.g. `sklearn.linear model.LogisticRegression`) with regularization coefficient C = 1.0. Report the accuracy and Balanced Error Rate (BER) of your classifier (1 mark).

In this problem, we read the dataset and use `sklearn.linear.model.LogisticRegression` to predict. In our setting, the accuracy and Balanced Error Rate is 0.966 and 0.481 respectively.

```python
# ============ Environment Setup ============
# ignore the warnings
import warnings
warnings.filterwarnings("ignore")
# import some packages
from sklearn import linear_model
from sklearn.metrics import accuracy_score
import numpy as np

# Read the dataset
def myReadData(path):
    f = open(path, 'r')
    while not '@data' in f.readline():
        pass
    dataset = []
    for l in f:
        if '?' in l: # Missing entry
            continue
        l = l.split(',')
        values = [1] + [float(x) for x in l]
        values[-1] = values[-1] > 0 # Convert to bool
        dataset.append(values)

    # Data setup
```

```
25        X = [values[:-1] for values in dataset]
26        y = [values[-1] for values in dataset]
27
28        return X, y
29
30  def myEvaluation(model, X, y):
31        # get predictions from model
32        pred = model.predict(X)
33
34        TP_ = np.logical_and(pred, y)
35        FP_ = np.logical_and(pred, np.logical_not(y))
36        TN_ = np.logical_and(np.logical_not(pred), np.logical_not(y))
37        FN_ = np.logical_and(np.logical_not(pred), y)
38
39        TP = sum(TP_)
40        FP = sum(FP_)
41        TN = sum(TN_)
42        FN = sum(FN_)
43
44        # Accuracy
45        correct = pred == y
46        print("Accuracy: %.3f" % np.mean(correct))
47        # BER
48        print("Balanced Error Rate: %.3f" %  (1 - 0.5 * (TP / (TP + FN) + TN
    / (TN + FP))))
49
50  X, y = myReadData("./Homework2/data/5year.arff")
51  # define a logistic regression
52  model = linear_model.LogisticRegression(C=1.0)
53  # fit a model
54  model.fit(X, y)
55  myEvaluation(model, X, y)
56
```

```
1  Accuracy: 0.966
2  Balanced Error Rate: 0.481
```

# Problem 3

Shuffle the data, and split it into training, validation, and test splits, with a 50/25/25% ratio. Using the class weight='balanced' option, and training on the training set, report the training/validation/test accuracy and BER (1 mark).

First, we split the dataset following 50 / 25 / 25 % ratio. And then, I train my model on the training set and get the performance on training, validation and test datasets.

- For training: Acc is 0.779 and Ber is 0.242
- For validating: Acc is 0.781 and Ber is 0.203

- For test: Acc is 0.780 and Ber is 0.262

```python
# import packages
import random
random.seed(5583)

def myShuffle(X, y):
    """ Shuffle the data """
    Xy = list(zip(X,y))
    random.shuffle(Xy)
    X = [d[0] for d in Xy]
    y = [d[1] for d in Xy]
    return X, y

# Train/validation/test splits
def mySplits(X, y):
    N = len(y)

    Ntrain = int(0.5*len(y))
    Nvalid = int(0.25*len(y))
    Ntest = int(0.25*len(y))

    Xtrain = X[:Ntrain]
    Xvalid = X[Ntrain:Ntrain+Nvalid]
    Xtest = X[Ntrain+Nvalid:]

    ytrain = y[:Ntrain]
    yvalid = y[Ntrain:Ntrain+Nvalid]
    ytest = y[Ntrain+Nvalid:]

    return Xtrain, Xvalid, Xtest, ytrain, yvalid, ytest

X, y = myReadData("./Homework2/data/5year.arff")
X, y = myShuffle(X, y)
Xtrain, Xvalid, Xtest, ytrain, yvalid, ytest = mySplits(X, y)
# define a logistic regression
model = linear_model.LogisticRegression(C=1.0, class_weight='balanced')
# fit a model
model.fit(Xtrain, ytrain)
print("=============== Training ===============")
myEvaluation(model, Xtrain, ytrain)    # myEvaluation function was
defined in Problem 1
print("=============== Validating ===============")
myEvaluation(model, Xvalid, yvalid)
print("=============== Testing ===============")
myEvaluation(model, Xtest, ytest)

```

```
1  =============== Training ===============
2  Accuracy: 0.779
3  Balanced Error Rate: 0.242
4  =============== Validating ===============
5  Accuracy: 0.781
6  Balanced Error Rate: 0.203
7  =============== Testing ===============
8  Accuracy: 0.780
9  Balanced Error Rate: 0.262
```

# Problem 4

Implement a complete regularization pipeline with the balanced classifier. Consider values of C in the range $\{10^{-4}, 10^{-3}, \ldots, 10^3, 10^4\}$. Report (or plot) the train, validation, and test BER for each value of C. Based on these values, which classifier would you select (in terms of generalization performance) and why (1 mark)?
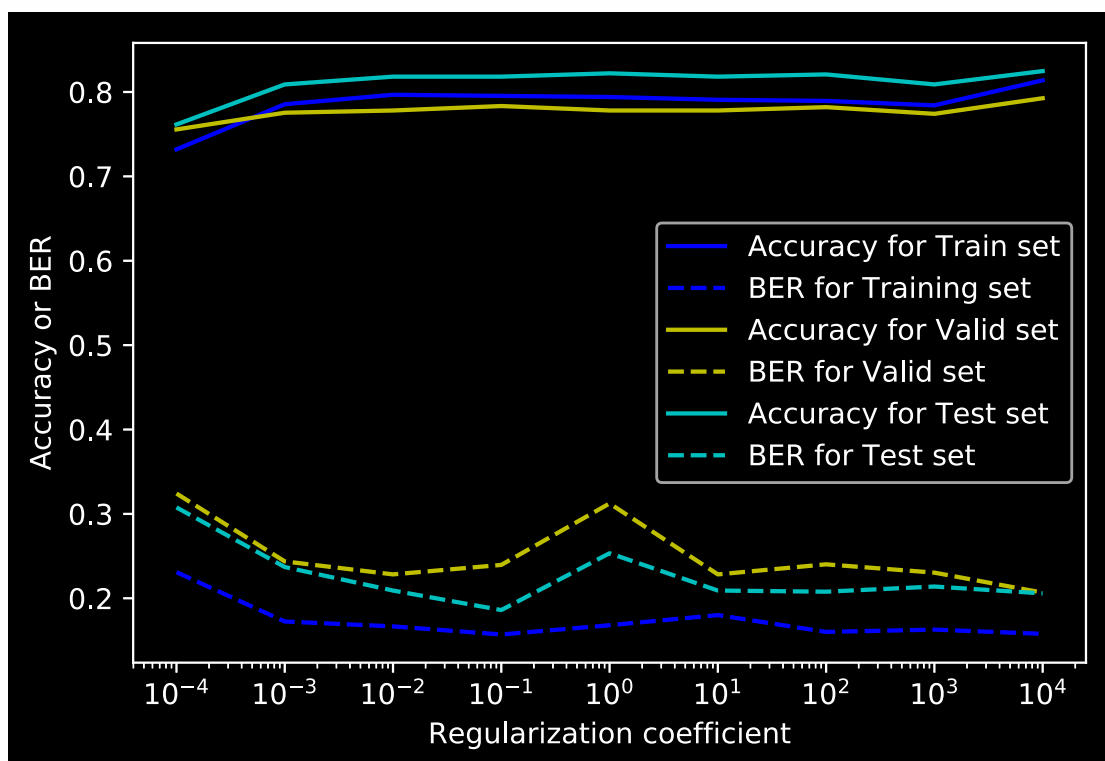
```python
1  import matplotlib.pyplot as plt
2
3  def myEvaluation(model, X, y, acc, ber):
4      # get predictions from model
5      pred = model.predict(X)
6
7      TP_ = np.logical_and(pred, y)
8      FP_ = np.logical_and(pred, np.logical_not(y))
9      TN_ = np.logical_and(np.logical_not(pred), np.logical_not(y))
10     FN_ = np.logical_and(np.logical_not(pred), y)
11
12     TP = sum(TP_)
13     FP = sum(FP_)
14     TN = sum(TN_)
15     FN = sum(FN_)
16
17     correct = pred == y
18     # Accuracy
19     acc.append(np.mean(correct))
20     # BER
21     ber.append(1 - 0.5 * (TP / (TP + FN) + TN / (TN + FP)))
22
23  def regPipeline(X, y, regs, trainAcc, trainBer, validAcc, validBer,
    testAcc, testBer):
24      X, y = myShuffle(X, y)
25      Xtrain, Xvalid, Xtest, ytrain, yvalid, ytest = mySplits(X, y)
26      for reg in regs:
27          model = linear_model.LogisticRegression(C=reg,
    class_weight='balanced')
28          # fit a model
29          model.fit(Xtrain, ytrain)
```

```
30          myEvaluation(model, Xtrain, ytrain, trainAcc, trainBer)
31          myEvaluation(model, Xvalid, yvalid, validAcc, validBer)
32          myEvaluation(model, Xtest, ytest, testAcc, testBer)
33
34  trainAcc, trainBer, validAcc, validBer, testAcc, testBer = [], [], [],
    [], [], []
35  regs = [10**N for N in [-4, -3, -2, -1, 0, 1, 2, 3, 4]]
36  X, y = myReadData("./Homework2/data/5year.arff")
37  regPipeline(X, y, regs, trainAcc, trainBer, validAcc, validBer, testAcc,
    testBer)
38
39  plt.plot(regs, trainAcc, 'b-', label='Accuracy for Train set')
40  plt.plot(regs, trainBer, 'b--', label='BER for Training set')
41  plt.plot(regs, validAcc, 'y-', label='Accuracy for Valid set')
42  plt.plot(regs, validBer, 'y--', label='BER for Valid set')
43  plt.plot(regs, testAcc, 'c-', label='Accuracy for Test set')
44  plt.plot(regs, testBer, 'c--', label='BER for Test set')
45  plt.legend()
46  plt.xlabel('Regularization coefficient')
47  plt.xscale('log')
48  plt.ylabel('Accuracy or BER')
49  plt.show()
50
51
```



Based on the plot, we will choose $10^4$ since it has the lowest BER for the validation set and meanwhile keeps a high accuracy on both the training and validation set.

# Problem 6

The sample weight option allows you to manually build a balanced (or imbalanced) classifier by assigning different weights to each datapoint (i.e., each label y in the training set). For example, we would assign equal weight to all samples by fitting:

```
1  weights = [1.0] * len(ytrain)
2  mod = linear_model.LogisticRegression(C=1, solver='lbfgs')
3  mod.fit(Xtrain, ytrain, sample_weight=weights)
```

(note that you should use the lbfgs solver option, and need not set class weight='balanced' in this case). Assigning larger weights to (e.g.) positive samples would encourage the logistic regressor to optimize for the True Positive Rate. Using the above code, compute the F$\beta$ score (on the test set) of your (unweighted) classifier, for $\beta = 1$ and $\beta = 10$. Following this, identify weight vectors that yield better performance (compared to the unweighted vector) in terms of the F1 and F10 scores (2 marks).

```
1  def FBeta(model, Xtest, ytest, beta):
2      pred = model.predict(Xtest)
3      retrieved = sum(pred)
4      relevant = sum(ytest)
5      intersection = sum([y and p for y, p in zip(ytest,pred)])
6      precision = intersection / retrieved
7      recall = intersection / relevant
8      F_Beta = (1+beta**2)*(precision*recall)/((beta**2)*precision+recall)
9      return F_Beta
10
11  print("=============== None weight =============== ")
12  weights = [1.0] * len(ytrain)
13  model = linear_model.LogisticRegression(C=1.0, solver='lbfgs')
14  model.fit(Xtrain, ytrain, sample_weight=weights)
15  print("F Score is %.3f" % FBeta(model, Xtest, ytest, 1), " for Beta 1")
16  print("F Score is %.3f" % FBeta(model, Xtest, ytest, 10), " for Beta 10")
17
18  print("=============== Balanced weight =============== ")
19  model = linear_model.LogisticRegression(C=1.0, solver='lbfgs',
       class_weight='balanced')
20  model.fit(Xtrain, ytrain)
21  print("F Score is %.3f" % FBeta(model, Xtest, ytest, 1), " for Beta 1")
22  print("F Score is %.3f" % FBeta(model, Xtest, ytest, 10), " for Beta 10")
23
24  print("=============== My weight =============== ")
25  true_ratio = len([1.0 for y in ytrain if y is True])/len(ytrain)
26  false_ratio = 1 - true_ratio
27  alphas = [2**i for i in range(-10, 8)]
28  Beta_1, Beta_10 = [], []
29  for alpha in alphas:
30      weights = [false_ratio if y is True else true_ratio*alpha for y in
       ytrain]
31      model = linear_model.LogisticRegression(C=1.0, solver='lbfgs')
32      model.fit(Xtrain, ytrain, sample_weight=weights)
```
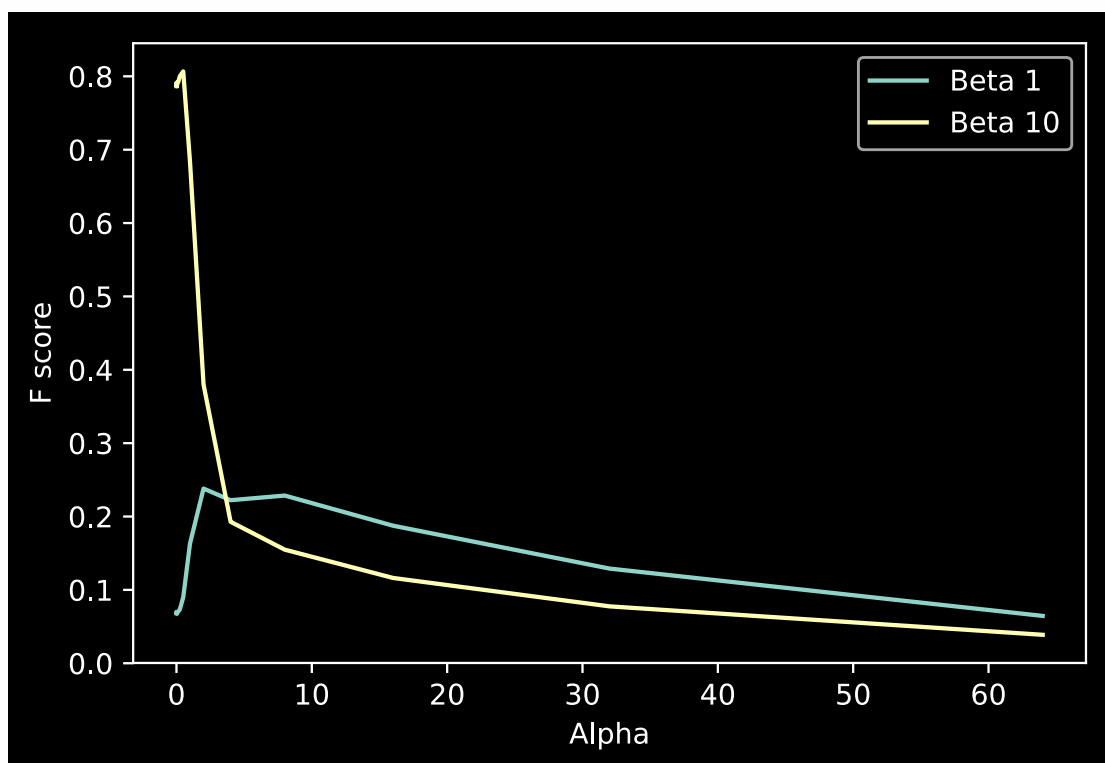
```
33        Beta_1.append(FBeta(model, Xtest, ytest, 1))
34        Beta_10.append(FBeta(model, Xtest, ytest, 10))
35
36  plt.plot(alphas, Beta_1, label='Beta 1')
37  plt.plot(alphas, Beta_10, label='Beta 10')
38  plt.legend()
39  plt.xlabel('Alpha')
40  plt.ylabel('F score')
41  plt.show()x
```

```
1  =============== None weight ===============
2  F Score is 0.242   for Beta 1
3  F Score is 0.155   for Beta 10
4  =============== Balanced weight ===============
5  F Score is 0.162   for Beta 1
6  F Score is 0.683   for Beta 10
7  =============== My weight ===============
```



We define the following weights to make the F score better.

```
1  true_ratio = len([1.0 for y in ytrain if y is True])/len(ytrain)
2  false_ratio = 1 - true_ratio
3  weights = [false_ratio if y is True else true_ratio*alpha for y in ytrain]
```

We find when $\alpha = 0.5$, we will achieve the best performance on f10 as 0.807, when $\alpha = 3$, we got the best performance on f1 as 0.259, which is better than the unweighted setting as f1 = 0.242 and f10 = 0.155.

# Tasks -- Dimensionality Reduction

Next we'll consider using PCA to build a lower-dimensional feature vector to do prediction.

## Problem 7

Following the stub code, compute the PCA basis on the training set. Report the first PCA component (i.e., pca.components [0])

The first component is:

```
1  from sklearn.decomposition import PCA
2
3  pca = PCA()
4  pca.fit(Xtrain)
5  print(pca.components_[0])
```

```
1   [-1.76636055e-18 -3.88774766e-08  3.91178524e-07  1.62093982e-06
2     9.81592130e-06  6.78648530e-04 -1.52568558e-06  2.31596680e-06
3     9.02354365e-06 -1.17215537e-06 -1.28587101e-07  3.39889517e-07
4     1.77000649e-06  6.98953943e-07  2.31596680e-06 -8.56732994e-03
5     1.67758042e-06  9.72886226e-06  2.31596680e-06  7.33390001e-07
6     9.30880022e-05  3.39868608e-07  3.08547249e-07  6.55226551e-07
7     5.00167483e-07  9.05428198e-07  1.48632063e-06 -1.76037474e-04
8     4.80818029e-05  4.94340305e-06 -2.36427176e-06  7.21586824e-07
9    -3.93427090e-04  7.03688492e-06 -2.38631665e-06  2.95736451e-07
10   -1.60518909e-06  2.00372698e-03 -7.06985096e-07  2.90046692e-07
11    4.62573681e-06 -3.24927526e-06  6.05977982e-07  1.33117229e-04
12    4.00277499e-05 -2.76291185e-06  7.12800934e-06  1.39281569e-04
13    3.77066322e-07  6.37742025e-07  7.26183476e-06 -1.14274329e-06
14   -1.00850645e-06  4.13873094e-06  4.80711918e-05  9.99960858e-01
15    2.53489811e-07  1.39209629e-06 -3.57102739e-07 -7.40231098e-07
16   -1.77887524e-04 -2.26509347e-05 -3.70642492e-04  8.73621271e-06
17   -1.27704672e-05]
```
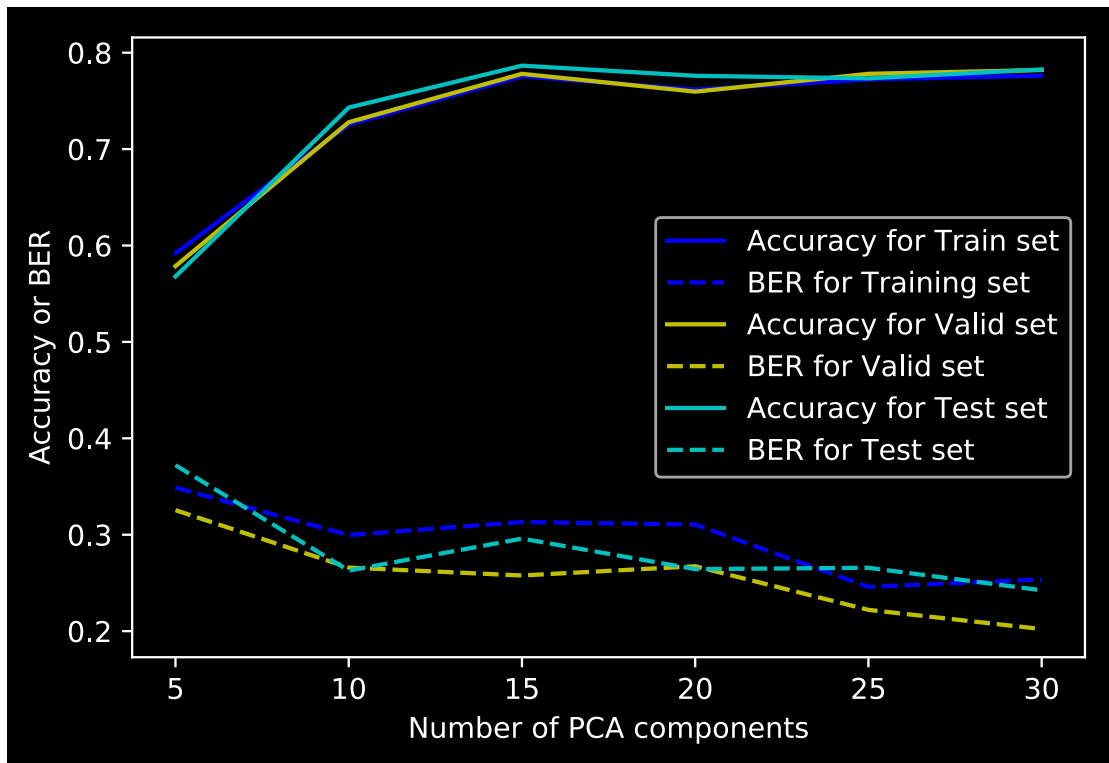
## Problem 8

Next we'll train a model using a low-dimensional feature vector. By representing the data in the above basis, i.e.:

```
1  Xpca_train = np.matmul(Xtrain, pca.components_.T)
2  Xpca_valid = np.matmul(Xvalid, pca.components_.T)
3  Xpca_test = np.matmul(Xtest, pca.components_.T)
```

compute the validation and test BER of a model that uses just the first N components (i.e., dimensions) for $N = 5, 10, \ldots, 25, 30$. Again use class weight='balanced' and C = 1.0 (2 marks).

```
1  def getLowDimension(N, Xtrain, Xvalid, Xtest):
2      Xpca_train = np.matmul(Xtrain, pca.components_[:N].T)
3      Xpca_valid = np.matmul(Xvalid, pca.components_[:N].T)
4      Xpca_test = np.matmul(Xtest, pca.components_[:N].T)
5      return Xpca_train, Xpca_valid, Xpca_test
6
7  def regPipeline(Xtrain, Xvalid, Xtest, ytrain, yvalid, ytest, trainAcc,
   trainBer, validAcc, validBer, testAcc, testBer):
8      model = linear_model.LogisticRegression(C=1.0,
   class_weight='balanced')
9      model.fit(Xtrain, ytrain)
10      myEvaluation(model, Xtrain, ytrain, trainAcc, trainBer)
11      myEvaluation(model, Xvalid, yvalid, validAcc, validBer)
12      myEvaluation(model, Xtest, ytest, testAcc, testBer)
13
14  trainAcc, trainBer, validAcc, validBer, testAcc, testBer = [], [], [],
   [], [], []
15  pca = PCA()
16  pca.fit(Xtrain)
17  Ns = list(range(5, 35, 5))
18  for N in Ns:
19      Xtrain_lowd, Xvalid_lowd, Xtest_lowd = getLowDimension(N, Xtrain,
   Xvalid, Xtest)
20      regPipeline(Xtrain_lowd, Xvalid_lowd, Xtest_lowd, ytrain, yvalid,
   ytest, trainAcc, trainBer, validAcc, validBer, testAcc, testBer)
21
22  plt.plot(Ns, trainAcc, 'b-', label='Accuracy for Train set')
23  plt.plot(Ns, trainBer, 'b--', label='BER for Training set')
24  plt.plot(Ns, validAcc, 'y-', label='Accuracy for Valid set')
25  plt.plot(Ns, validBer, 'y--', label='BER for Valid set')
26  plt.plot(Ns, testAcc, 'c-', label='Accuracy for Test set')
27  plt.plot(Ns, testBer, 'c--', label='BER for Test set')
28  plt.legend()
29  plt.xlabel('Number of PCA components')
30  plt.ylabel('Accuracy or BER')
31  plt.show()
```

Here we know $N = 30$ is the best choice for this experiment, because $N = 30$ can remain more informatipon of the orginal features (dimension is 65) and also is good for generalization.