

# Homework 2: classification

Data source: <http://archive.ics.uci.edu/ml/datasets/Polish+companies+bankruptcy+data> **Description:** The goal of this HW is to be familiar with the basic classifiers PML Ch 3. For this HW, we continue to use Polish companies bankruptcy data Data Set from UCI Machine Learning Repository. Download the dataset and put the 4th year file (4year.arff) in your YOUR\_GITHUB\_ID/PHBS\_MLF\_2019/HW2/ I did a basic process of the data (loading to dataframe, creating bankruptcy column, changing column names, filling-in na values, training-vs-test split, standardizatio, etc). See my github.

## Preparation

### Load, read and clean

```
1 from scipy.io import arff
2 import pandas as pd
3 import numpy as np
4
5 data = arff.loadarff('./data/4year.arff')
6 df = pd.DataFrame(data[0])
7 df['bankruptcy'] = (df['class']==b'1')
8 del df['class']
9 df.columns = ['x{0:02d}'.format(k) for k in range(1,65)] + ['bankruptcy']
10 df.describe()
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	X01	X02	X03	X04	X05	X06	X07	X08
count	9791.000000	9791.000000	9791.000000	9749.000000	9.771000e+03	9791.000000	9791.000000	9773.000000
mean	0.043019	0.596404	0.130959	8.136600	6.465164e+01	-0.059273	0.059446	19.884016
std	0.359321	4.587122	4.559074	290.647281	1.475939e+04	6.812754	0.533344	698.697015
min	-12.458000	0.000000	-445.910000	-0.045319	-3.794600e+05	-486.820000	-12.458000	-1.848200
25%	0.001321	0.263145	0.020377	1.047000	-5.121700e+01	-0.000578	0.003004	0.428300
50%	0.041364	0.467740	0.199290	1.591800	-5.557600e-02	0.000000	0.048820	1.088700
75%	0.111130	0.689255	0.410670	2.880400	5.573200e+01	0.065322	0.126940	2.691000
max	20.482000	446.910000	22.769000	27146.000000	1.034100e+06	322.200000	38.618000	53209.000000

8 rows × 64 columns

```
1 sum(df.bankruptcy == True)
```

```
1 515
```

```
1 from sklearn.impute import SimpleImputer
2
3 imp_mean = SimpleImputer(missing_values=np.nan, strategy='mean')
4 x_imp = imp_mean.fit_transform(df.values)
5
```

A dll load error occurred here. Solution recorded in [my blog](#)

```

1 from sklearn.model_selection import train_test_split
2
3 X, y = X_imp[:, :-1], X_imp[:, -1]
4 X_train, X_test, y_train, y_test = \
5     train_test_split(X, y,
6                     test_size=0.3,
7                     random_state=0,
8                     stratify=y)

```

```

1 from sklearn.preprocessing import StandardScaler
2
3 stdsc = StandardScaler()
4 X_train_std = stdsc.fit_transform(X_train)
5 X_test_std = stdsc.transform(X_test)

```

## 1. Find the 2 most important features

Select the 2 most important features using LogisticRegression with L1 penalty. (Adjust C until you see 2 features)

```

1 from sklearn.linear_model import LogisticRegression
2
3 C = [1, .1, .01, 0.001]
4 cdf = pd.DataFrame()
5
6 for c in C:
7     lr = LogisticRegression(penalty='l1', C=c, solver='liblinear', random_state=0)
8     lr.fit(X_train_std, y_train)
9     print(f'[C={c}] with {lr.coef_[lr.coef_!=0].shape[0]} features: \n {lr.coef_[lr.coef_!=0]} \n') # Python >= 3.7
10    if lr.coef_[lr.coef_!=0].shape[0] == 2:
11        cdf = pd.DataFrame(lr.coef_.T, df.columns[:-1], columns=['coef'])

```

```

1 [C=1] with 41 features:
2 [-0.21124721 -0.32721186 -0.027786 -1.13272997 -0.06357798  2.33195848
3  -1.15522622 -0.00657024 -0.04684187 -1.08683927  0.1425248  0.01096755
4  -0.01922072 -0.01383184 -0.01411706 -0.10736095  0.00238513  0.24577125
5  -0.85510327  0.85209928 -0.30366778 -0.30474956 -0.0296142 -0.01749839
6  -0.04401046 -0.03613927  0.12211322 -0.0138838 -0.45429792  1.3157471
7  -0.4035633 -0.46000256  0.04269182 -0.13000407 -0.0474417 -0.11325886
8  -0.12953285 -0.56454917  0.00689488 -0.4300807 -0.08949632]
9
10 [C=0.1] with 14 features:
11 [-0.07260191 -0.00523502 -0.05866954 -0.1005855 -0.62441034 -0.04143727
12  -0.07701467 -0.00230242  0.01202513  0.01090586  0.28007452 -0.10726656
13  -0.04314757 -0.00208574]
14
15 [C=0.01] with 2 features:
16 [-0.00174613 -0.05571114]
17
18 [C=0.001] with 0 features:
19 []

```

```

1 lr = LogisticRegression(penalty='l1', C=0.01, solver='liblinear', random_state=0) # complete
2 lr.fit(X_train_std, y_train)

```

```

1 LogisticRegression(C=0.01, class_weight=None, dual=False, fit_intercept=True,
2                    intercept_scaling=1, l1_ratio=None, max_iter=100,
3                    multi_class='auto', n_jobs=None, penalty='l1',
4                    random_state=0, solver='liblinear', tol=0.0001, verbose=0,
5                    warm_start=False)

```

```

1 cdf = cdf[cdf.coef != 0]
2 cdf

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

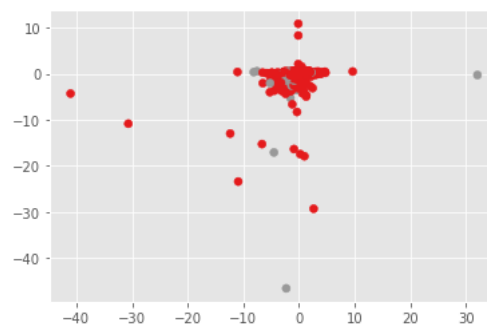
	coef
X01	-0.001746
X38	-0.055711

## redefine X\_train\_std and X\_test\_std

```
1 X_train_std = X_train_std[:, 1r.coef_[0]!=0]
2 X_test_std = X_test_std[:, 1r.coef_[0]!=0]
```

```
1 from matplotlib.colors import ListedColormap
2 import matplotlib.pyplot as plt
3 plt.style.use('ggplot')
4 plt.scatter(x=X_train_std[:,0], y=X_train_std[:,1], c=y_train, cmap='Set1')
```

```
1 <matplotlib.collections.PathCollection at 0x21ee24bfcc8>
```

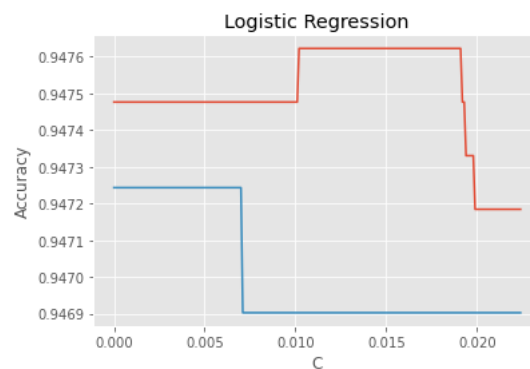


## 2. Apply LR / SVM / Decision Tree below

Using the 2 selected features, apply LR / SVM / decision tree. **Try your own hyperparameters (C, gamma, tree depth, etc)** to maximize the prediction accuracy. (Just try several values. You don't need to show your answer is the maximum.)

### LR

```
1 CLr = np.arange(0.0000000000000001, 0.0225, 0.0001)
2 acrcLr = [] # accuracy
3 for c in CLr:
4     lr = LogisticRegression(C=c,penalty='l1',solver='liblinear')
5     lr.fit(X_train_std, y_train)
6     acrcLr.append([lr.score(X_train_std, y_train), lr.score(X_test_std, y_test), c])
7 acrcLr = np.array(acrcLr)
8 plt.plot(acrcLr[:,2], acrcLr[:,0])
9 plt.plot(acrcLr[:,2], acrcLr[:,1])
10 plt.xlabel('C')
11 plt.ylabel('Accuracy')
12 plt.title('Logistic Regression')
13 plt.show()
```



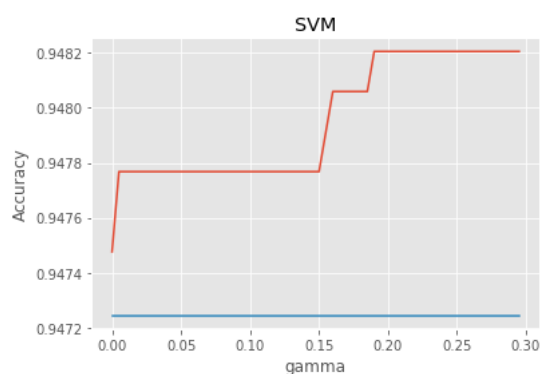
Choose `C=.01`

```
1 c = .01
2 lr = LogisticRegression(C=c,penalty='l1',solver='liblinear')
3 lr.fit(X_train_std, y_train)
4 print(f'Accuracy when [c={c}] \nTrain {lr.score(X_train_std, y_train)}\nTest {lr.score(X_test_std, y_test)}')
```

```
1 Accuracy when [c=0.01]
2 Train 0.9474759264662971
3 Test 0.9469026548672567
```

## SVM

```
1 from sklearn.svm import SVC
2 G = np.arange(0.00001, 0.3, 0.005)
3 acrcSvm = []
4 for g in G:
5     svm = SVC(kernel='rbf', gamma=g, C=1.0, random_state=0)
6     svm.fit(X_train_std, y_train)
7     acrcSvm.append([svm.score(X_train_std, y_train), svm.score(X_test_std, y_test), g])
8 acrcSvm = np.array(acrcSvm)
9 plt.plot(acrcSvm[:,2], acrcSvm[:,0])
10 plt.plot(acrcSvm[:,2], acrcSvm[:,1])
11 plt.xlabel('gamma')
12 plt.ylabel('Accuracy')
13 plt.title('SVM')
14 plt.show()
```



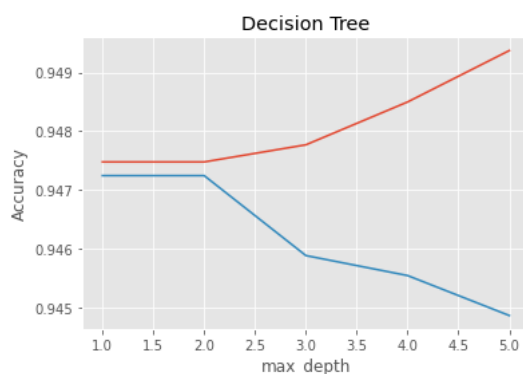
Choose `gamma = 0.2`

```
1 g = 0.2
2 svm = SVC(kernel='rbf', gamma=g, C=1.0, random_state=0)
3 svm.fit(X_train_std, y_train)
4 print(f'Accuracy when [gamma={g}] \nTrain {svm.score(X_train_std, y_train)}\nTest {svm.score(X_test_std, y_test)}')
```

```
1 Accuracy when [gamma=0.2]
2 Train 0.9482054274875985
3 Test 0.9472430224642614
```

## Decision Tree

```
1 from sklearn.tree import DecisionTreeClassifier
2 depthTree = range(1, 6)
3 acrcTree = []
4 for depth in depthTree:
5     tree = DecisionTreeClassifier(criterion='gini', max_depth=depth, random_state=0)
6     tree.fit(X_train_std, y_train)
7     acrcTree.append([tree.score(X_train_std, y_train), tree.score(X_test_std, y_test), depth])
8 acrcTree = np.array(acrcTree)
9 plt.plot(acrcTree[:,2], acrcTree[:,0])
10 plt.plot(acrcTree[:,2], acrcTree[:,1])
11 plt.xlabel('max_depth')
12 plt.ylabel('Accuracy')
13 plt.title('Decision Tree')
14 plt.show()
```



Choose `max_depth=2`:

```
1 depth = 2
2 tree = DecisionTreeClassifier(criterion='gini', max_depth=depth, random_state=0)
3 tree.fit(X_train_std, y_train)
4 print(f'Accuracy when [max_depth={depth}] \nTrain {tree.score(X_train_std, y_train)}\nTest {tree.score(X_test_std, y_test)}')
```

```
1 Accuracy when [max_depth=2]
2 Train 0.9474759264662971
3 Test 0.9472430224642614
```

### 3. Visualize the classification

Visualize your classifiers using the `plot_decision_regions` function from PML Ch. 3

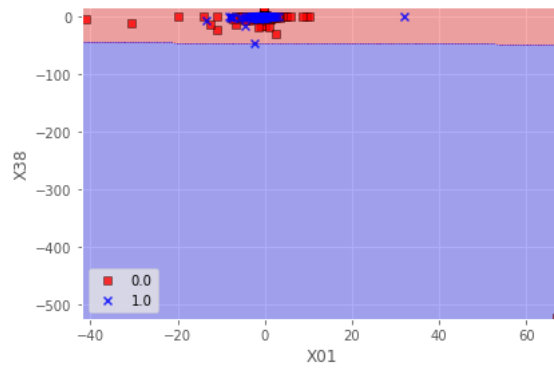
```
1 def plot_decision_regions(X, y, classifier, test_idx=None, resolution=0.02):
2
3     # setup marker generator and color map
4     markers = ('s', 'x', 'o', '^', 'v')
5     colors = ('red', 'blue', 'lightgreen', 'gray', 'cyan')
6     cmap = ListedColormap(colors[:len(np.unique(y))])
7
8     # plot the decision surface
9     x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
10    x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
11    xx1, xx2 = np.meshgrid(np.arange(x1_min, x1_max, resolution),
12                           np.arange(x2_min, x2_max, resolution))
13    Z = classifier.predict(np.array([xx1.ravel(), xx2.ravel()]).T)
14    Z = Z.reshape(xx1.shape)
15    plt.contourf(xx1, xx2, Z, alpha=0.3, cmap=cmap)
16    plt.xlim(xx1.min(), xx1.max())
17    plt.ylim(xx2.min(), xx2.max())
18
19    for idx, c1 in enumerate(np.unique(y)):
20        plt.scatter(x=X[y == c1, 0],
21                  y=X[y == c1, 1],
22                  alpha=0.8,
23                  c=colors[idx],
24                  marker=markers[idx],
25                  label=c1,
26                  edgecolor='black')
27
28    # highlight test samples
29    if test_idx:
30        # plot all samples
31        X_test, y_test = X[test_idx, :], y[test_idx]
32
33        plt.scatter(X_test[:, 0],
34                  X_test[:, 1],
35                  c='',
36                  edgecolor='black',
37                  alpha=1.0,
38                  linewidth=1,
39                  marker='o',
40                  s=100,
41                  label='test set')
```

```
1 X_combined_std = np.vstack((X_train_std, X_test_std))
2 y_combined = np.hstack((y_train, y_test))
```

## LR

`test_idx` removed on purpose

```
1 plot_decision_regions(X=X_combined_std, y=y_combined,
2                       classifier=lr)
3 plt.xlabel(cdf.index[0])
4 plt.ylabel(cdf.index[1])
5 plt.legend(loc='lower left')
6
7 plt.tight_layout()
8 #plt.savefig('images/03_01.png', dpi=300)
9 plt.show()
```

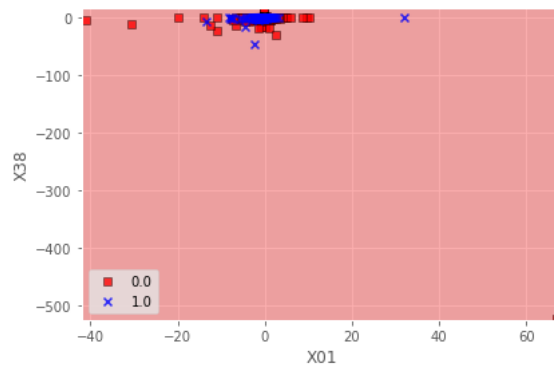


## Decision Tree

```

1 plot_decision_regions(X=X_combined_std, y=y_combined,
2                       classifier=tree)
3 plt.xlabel(cdf.index[0])
4 plt.ylabel(cdf.index[1])
5 plt.legend(loc='lower left')
6
7 plt.tight_layout()
8 #plt.savefig('images/03_01.png', dpi=300)
9 plt.show()

```



## SVM (samples)

```

1 # Visualization of all features in a SVM model is too slow
2 # Because the complexity is very high (source:https://scikit-learn.org/stable/modules/svm.html#complexity)
3 # So use random samples(n=3000) instead
4
5 samples = np.random.randint(0, len(X_combined_std), size=3000)
6 plot_decision_regions(X=X_combined_std[samples], y=y_combined[samples],
7                       classifier=svm)
8 plt.xlabel(cdf.index[0] + '[samples]')
9 plt.ylabel(cdf.index[1] + '[samples]')
10 plt.legend(loc='lower left')
11
12 plt.tight_layout()
13 #plt.savefig('images/03_01.png', dpi=300)
14 plt.show()

```

