

Don't over fit II competition

This is my EE551 individual project. It is a playground prediction competition on Kaggle.

Exploratory Data Analysis(EDA)

Datacollection

In [2]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 train = pd.read_csv('train.csv')
5 test = pd.read_csv('test.csv')
```

In [3]:

```
1 train.shape
```

Out[3]:

(250, 302)

In [4]:

```
1 train.head()
```

Out[4]:

	id	target	0	1	2	3	4	5	6	7	...	290	291	
0	0	1.0	-0.098	2.165	0.681	-0.614	1.309	-0.455	-0.236	0.276	...	0.867	1.347	0.
1	1	0.0	1.081	-0.973	-0.383	0.326	-0.428	0.317	1.172	0.352	...	-0.165	-1.695	-1.
2	2	1.0	-0.523	-0.089	-0.348	0.148	-0.022	0.404	-0.023	-0.172	...	0.013	0.263	-1.
3	3	1.0	0.067	-0.021	0.392	-1.637	-0.446	-0.725	-1.035	0.834	...	-0.404	0.640	-0.
4	4	1.0	2.347	-0.831	0.511	-0.021	1.225	1.594	0.585	1.509	...	0.898	0.134	2.

5 rows × 302 columns

In [5]:

```
1 sample_submission = pd.read_csv('sample_submission.csv')
2 sample_submission.head()
```

Out[5]:

	id	target
0	250	0
1	251	0
2	252	0
3	253	0
4	254	0

In [6]:

```
1 test.head()
```

Out[6]:

	id	0	1	2	3	4	5	6	7	8	...	290	291
0	250	0.500	-1.033	-1.595	0.309	-0.714	0.502	0.535	-0.129	-0.687	...	-0.088	-2.628
1	251	0.776	0.914	-0.494	1.347	-0.867	0.480	0.578	-0.313	0.203	...	-0.683	-0.066
2	252	1.750	0.509	-0.057	0.835	-0.476	1.428	-0.701	-2.009	-1.378	...	-0.094	0.351
3	253	-0.556	-1.855	-0.682	0.578	1.592	0.512	-1.419	0.722	0.511	...	-0.336	-0.787
4	254	0.754	-0.245	1.173	-1.623	0.009	0.370	0.781	-1.763	-1.432	...	2.184	-1.090

5 rows × 301 columns

In [7]:

```
1 train.tail()
```

Out[7]:

	id	target	0	1	2	3	4	5	6	7	...	290	291
245	245	0.0	-1.199	0.466	-0.908	2.771	1.631	0.931	0.182	-0.652	...	0.724	0.177
246	246	0.0	0.237	0.233	-0.380	-1.748	0.839	-0.721	-0.114	0.005	...	0.857	0.147
247	247	0.0	1.411	-1.465	0.119	0.583	1.634	-0.207	1.173	1.622	...	-0.499	-0.455
248	248	1.0	0.620	1.040	0.184	-0.570	-0.087	-0.748	-1.559	-0.553	...	0.557	-1.494
249	249	0.0	0.489	0.403	0.139	-2.046	1.345	0.122	1.255	0.647	...	-0.025	1.305

5 rows × 302 columns

In [8]:

```
1 train.columns
```

Out[8]:

```
Index(['id', 'target', '0', '1', '2', '3', '4', '5', '6', '7',  
      ...  
      '290', '291', '292', '293', '294', '295', '296', '297', '298',  
      '299'],  
      dtype='object', length=302)
```

In [9]:

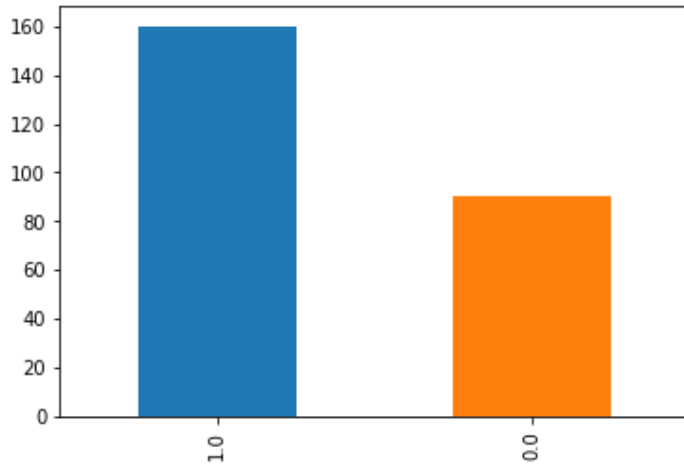
```
1 print(train.info())
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 250 entries, 0 to 249  
Columns: 302 entries, id to 299  
dtypes: float64(301), int64(1)  
memory usage: 589.9 KB  
None
```

Visualization

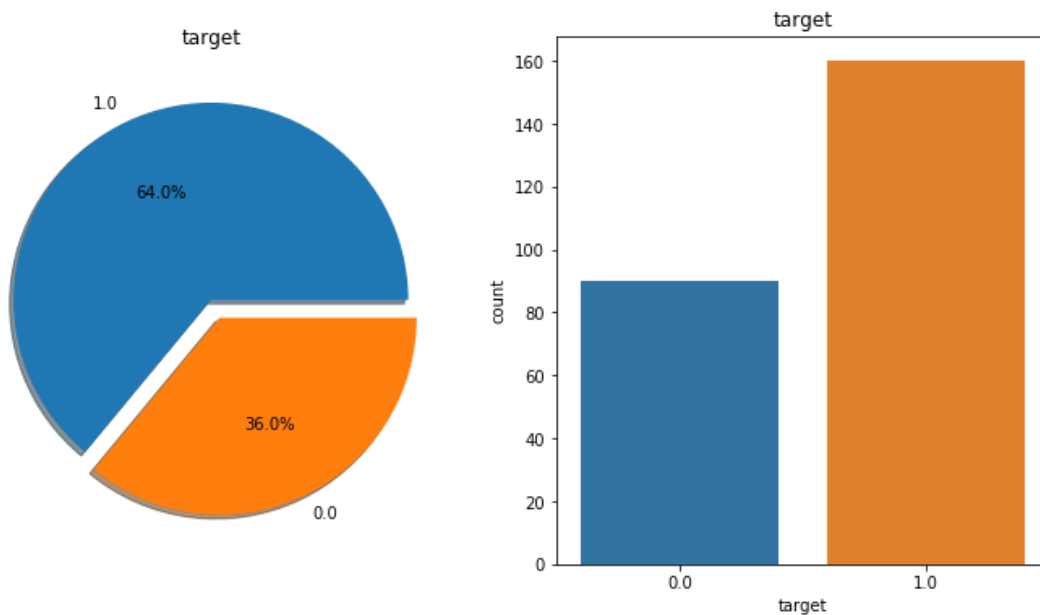
In [10]:

```
1 train['target'].value_counts().plot.bar();
```



In [11]:

```
1 import seaborn as sns
2 f,ax = plt.subplots(1,2,figsize=(12,6))
3 train['target'].value_counts().plot.pie(explode=[0,0.1],autopct='%1.1f%%',ax=ax[0])
4 ax[0].set_title('target')
5 ax[0].set_ylabel('')
6 sns.countplot('target', data = train, ax = ax[1])
7 ax[1].set_title('target')
8 plt.show()
```

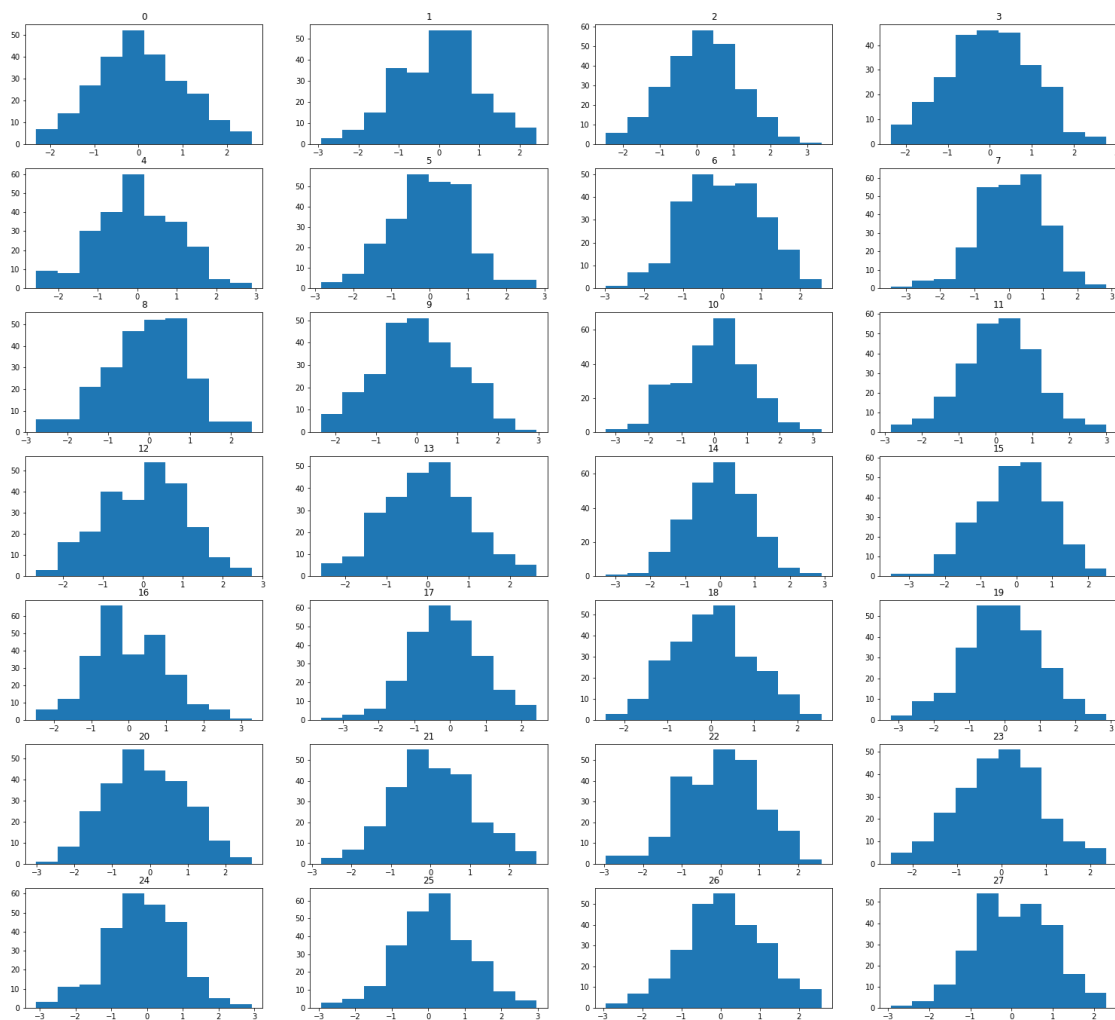


In [12]:

```

1 plt.figure(figsize = (26,24))
2 for i, col in enumerate(list(train.columns)[2:30]):
3     plt.subplot(7, 4, i+1)
4     plt.hist(train[col])
5     plt.title(col)

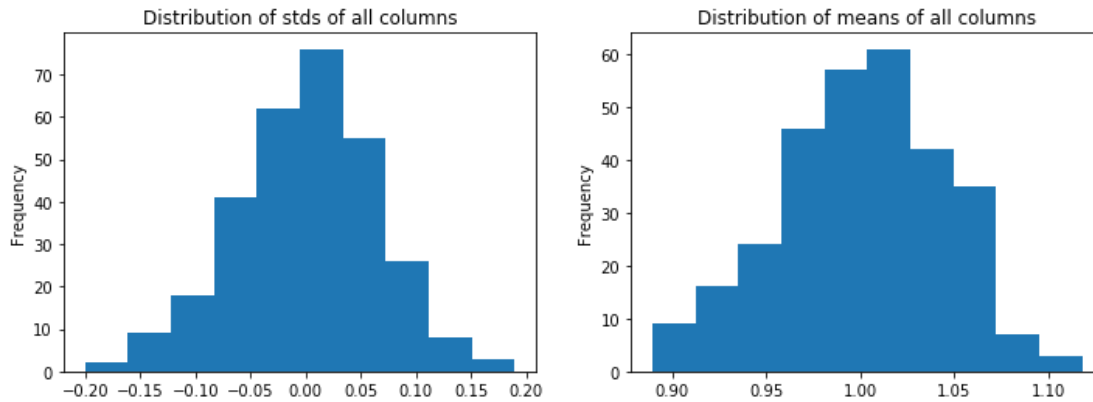
```



Values in columns are more or less similar.

In [13]:

```
1 plt.figure(figsize = (12,4))
2 plt.subplot(1,2,1)
3 train[train.columns[2:]].mean().plot('hist')
4 plt.title('Distribution of stds of all columns')
5 plt.subplot(1,2,2)
6 train[train.columns[2:]].std().plot('hist')
7 plt.title('Distribution of means of all columns')
8 plt.show()
```



Columns have mean of 0 +/- 0.15 and std of 1 +/- 0.1.

In [14]:

```
1 corr = train.corr()['target'].sort_values(ascending = False)
```

In [15]:

```
1 corr.head(10)
```

Out[15]:

```
target    1.000000
33        0.373608
65        0.293846
24        0.173096
183       0.164146
199       0.159442
201       0.142238
30        0.132705
289       0.127213
114       0.124792
Name: target, dtype: float64
```

In [16]:

```
1 corr.tail(10)
```

Out[16]:

```
16    -0.144267
194    -0.150384
id     -0.151498
189    -0.155956
80     -0.162558
73     -0.167557
295    -0.170501
91     -0.192536
117    -0.197496
217    -0.207215
Name: target, dtype: float64
```

Logistic regression

In [17]:

```
1 from sklearn.model_selection import train_test_split, learning_curve, Stratified
2 from sklearn.preprocessing import StandardScaler
3 X_train = train.drop(['id', 'target'], axis = 1)
4 y_train = train['target']
5 X_test = test.drop(['id'], axis = 1)
```

Find the best parameters for function 'LogisticRegression'.

In [18]:

```

1  from sklearn.linear_model import LogisticRegression
2  log = LogisticRegression(penalty = 'l1', random_state = 42)
3  params = {'solver': ['liblinear', 'saga'],
4            'C': [0.001, 0.1, 1, 10, 50],
5            'tol': [0.00001, 0.0001, 0.001, 0.005],
6            'class_weight': ['balanced', None]}
7  log_gs = GridSearchCV(log, params, cv = StratifiedKFold(n_splits = 5), verbose =
8
9  log_gs.fit(X_train, y_train)
10
11 log_best = log_gs.best_estimator_
12
13 print(log_best)
14 print(log_gs.best_score_)

```

Fitting 5 folds for each of 80 candidates, totalling 400 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n_jobs=-1)]: Done 34 tasks | elapsed: 2.5s

LogisticRegression(C=1, class_weight=None, dual=False, fit_intercept=True,

intercept_scaling=1, max_iter=100, multi_class='warn',
n_jobs=None, penalty='l1', random_state=42, solver='liblinear',
tol=1e-05, verbose=0, warm_start=False)

0.8177083333333334

[Parallel(n_jobs=-1)]: Done 400 out of 400 | elapsed: 5.6s finished

Define a function to plot learning curve.

In [19]:

```

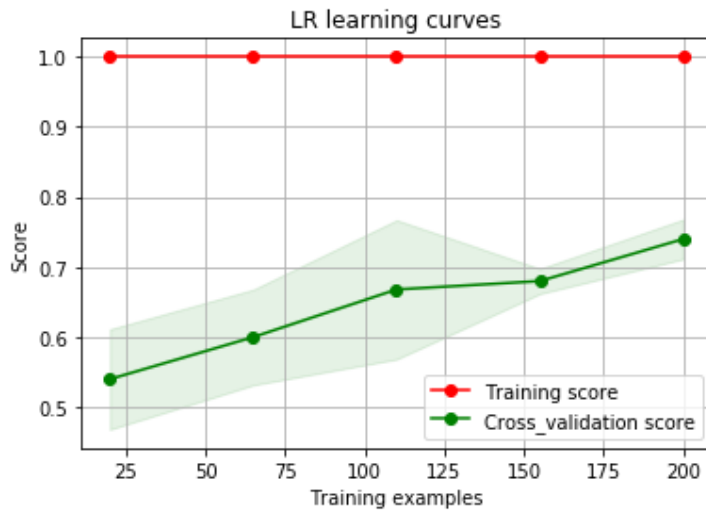
1  def plot_learning_curve(estimator, title, X, y, ylim = None, cv = None, n_jobs =
2  plt.figure()
3  plt.title(title)
4  if ylim is not None:
5      plt.ylim(*ylim)
6  plt.xlabel("Training examples")
7  plt.ylabel("Score")
8  train_sizes, train_scores, test_scores = learning_curve(estimator, X_train,
9  train_scores_mean = np.mean(train_scores, axis = 1)
10 train_scores_std = np.std(train_scores, axis = 1)
11 test_scores_mean = np.mean(test_scores, axis = 1)
12 test_scores_std = np.std(test_scores, axis = 1)
13 plt.grid()
14 plt.fill_between(train_sizes, train_scores_mean-train_scores_std, train_scores_mean+train_scores_std, color='r')
15 plt.fill_between(train_sizes, test_scores_mean-test_scores_std, test_scores_mean+test_scores_std, color='g')
16 plt.plot(train_sizes, train_scores_mean, 'o-', color = 'r', label = "Training score")
17 plt.plot(train_sizes, test_scores_mean, 'o-', color = 'g', label = "Cross-validation score")
18 plt.legend(loc = 'best')
19 return plt
20

```

Plot the learning curve of log_best.

In [20]:

```
1 learningCurve = plot_learning_curve(log_best, "LR learning curves",X_train, y_train)
```



Define a function to draw roc curve.

In [21]:

```

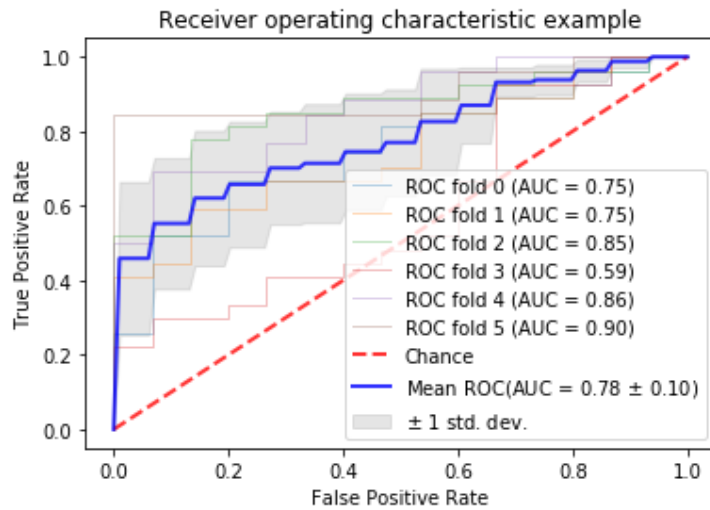
1  from sklearn.metrics import confusion_matrix, classification_report, roc_curve,
2  from scipy import interp
3  def plot_roc(clf, X = X_train, y = y_train, n = 6):
4      tprs = []
5      aucs = []
6      mean_fpr = np.linspace(0,1,100)
7      i = 0
8      classifier = clf
9      cv = StratifiedKFold(n_splits = n)
10     for train, test in cv.split(X,y):
11         probas_ = classifier.fit(X.iloc[train], y.iloc[train]).predict_proba(X.
12         fpr, tpr, thresholds = roc_curve(y[test], probas_[ :, 1])
13         tprs.append(interp(mean_fpr, fpr, tpr))
14         tprs[-1][0] = 0.0
15         roc_auc = auc(fpr,tpr)
16         aucs.append(roc_auc)
17         plt.plot(fpr, tpr, lw = 1, alpha = 0.3, label = 'ROC fold %d (AUC = %0.2f)' % (i+1, roc_auc))
18         i += 1
19     plt.plot([0, 1], [0, 1], linestyle = '--', lw =2, color = 'r', label = 'Char
20     mean_tpr = np.mean(tprs, axis = 0)
21     mean_tpr[-1] = 1.0
22     mean_auc = auc(mean_fpr, mean_tpr)
23     std_auc = np.std(aucs)
24     plt.plot(mean_fpr, mean_tpr, color = 'b', label = r'Mean ROC(AUC = %0.2f $\pm$ %0.2f)' % (mean_auc, std_auc))
25     std_tpr = np.std(tprs, axis = 0)
26     tprs_upper = np.minimum(mean_tpr + std_tpr,1)
27     tprs_lower = np.maximum(mean_tpr - std_tpr,0)
28     plt.fill_between(mean_fpr, tprs_lower, tprs_upper, color = 'grey', alpha = 0.1)
29     plt.xlim([-0.05, 1.05])
30     plt.ylim([-0.05, 1.05])
31     plt.xlabel('False Positive Rate')
32     plt.ylabel('True Positive Rate')
33     plt.title('Receiver operating characteristic example')
34     plt.legend(loc = 'lower right')
35     plt.show()
36
37

```

Plot the roc curve of log_best.

In [22]:

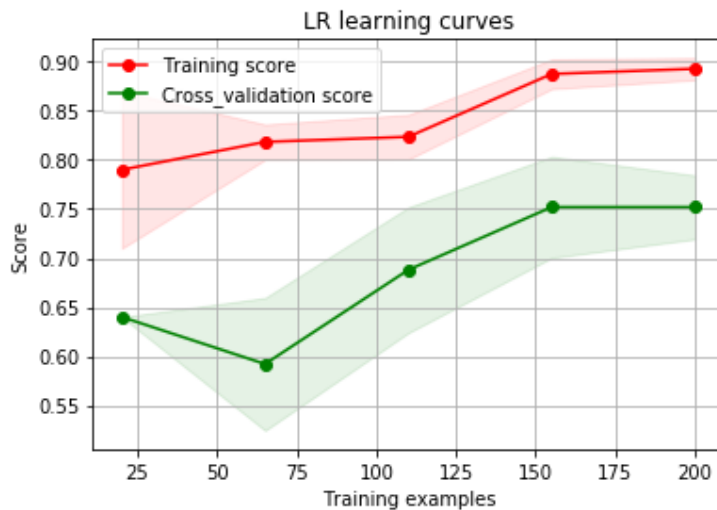
```
1 roc = plot_roc(log_best)
```



cv_score is far away from training score. It is overfitting. C is responsible for level of regularization and the smaller it is, the bigger the level of regularization it is. First try C = 0.1

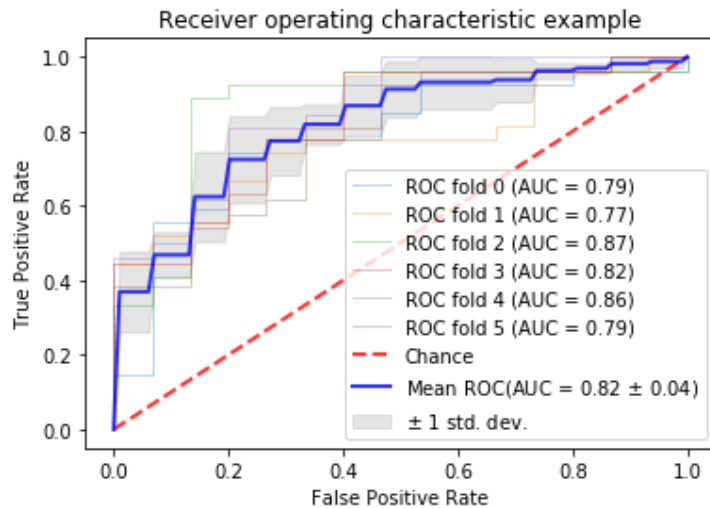
In [23]:

```
1 log_p0 = LogisticRegression(class_weight = 'balanced', penalty = 'l1', C = 0.1,
2 learningCurve0 = plot_learning_curve(log_p0, "LR learning curves", X_train, y_train)
3
```



In [24]:

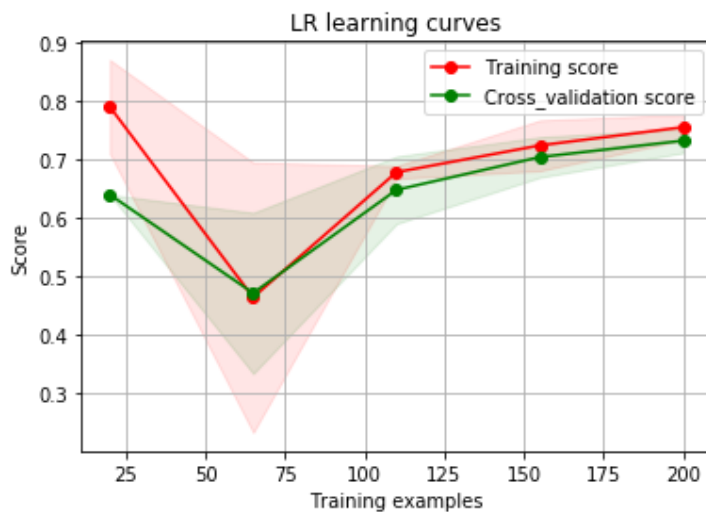
```
1 roc0 = plot_roc(log_p0)
```



Try C = 0.05.

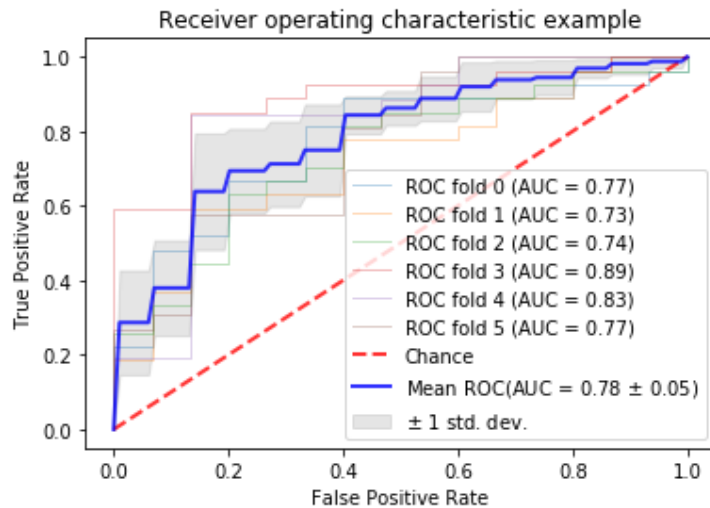
In [25]:

```
1 log_p1 = LogisticRegression(class_weight = 'balanced', penalty = 'l1', C = 0.05,
2 learningCurve1 = plot_learning_curve(log_p1, "LR learning curves", X_train, y_train,
3
```



In [26]:

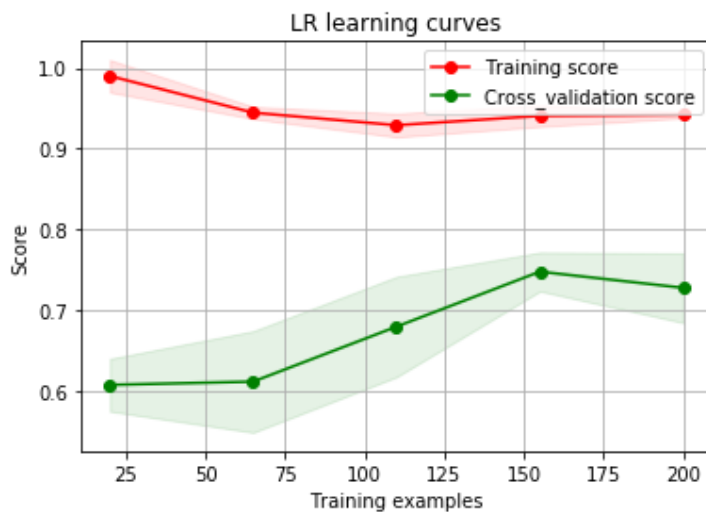
```
1 roc1 = plot_roc(log_p1)
```



Try C = 0.15.

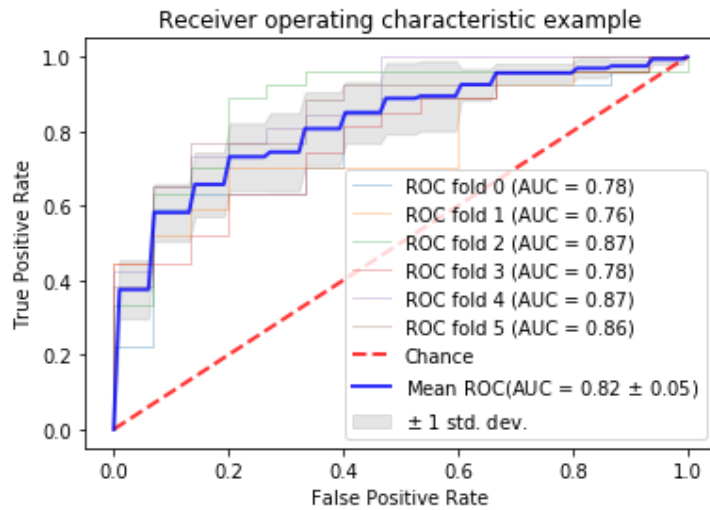
In [27]:

```
1 log_p2 = LogisticRegression(class_weight = 'balanced', penalty = 'l1', C = 0.15,  
2 learningCurve2 = plot_learning_curve(log_p2, "LR learning curves", X_train, y_train,  
3
```



In [28]:

```
1 roc2 = plot_roc(log_p2)
```



It seems like that when $C = 0.1$, the model performs best.

Output the first submission file.

In [29]:

```
1 log_p0.fit(X_train, y_train)
2 log_pred0 = log_p0.predict_proba(X_test)[:,-1]
3 submission0 = pd.DataFrame({'id':test['id'],
4                             'target':log_pred0})
5 submission0.to_csv('submission0.csv', index = False)
```

Feature Selection

Use eli5 to do the feature selection.

In [30]:

```
1 import eli5
2 eli5.show_weights(log_p0,top = 50)
```

Out[30]:

y=1.0 top features

Weight?	Feature
+0.713	x33
+0.491	x65
+0.370	<BIAS>
+0.229	x199
+0.070	x101
+0.032	x226
+0.028	x24
+0.026	x176
+0.015	x30
+0.014	x17
+0.013	x201
+0.006	x183
-0.001	x209
-0.008	x156
-0.020	x239
-0.022	x180
-0.024	x252
-0.030	x4
-0.030	x237
-0.034	x288
-0.037	x276
-0.039	x127
-0.050	x90
-0.055	x227
-0.055	x165
-0.065	x134
-0.069	x82
-0.076	x298
-0.084	x43
-0.096	x16
-0.096	x133
-0.108	x80
-0.109	x108
-0.109	x194
-0.115	x189
-0.117	x258
-0.165	x295
-0.196	x73
-0.196	x117
-0.281	x91
-0.304	x217

In [31]:

```
1 (log_p0.coef_ != 0).sum()
```

Out[31]:

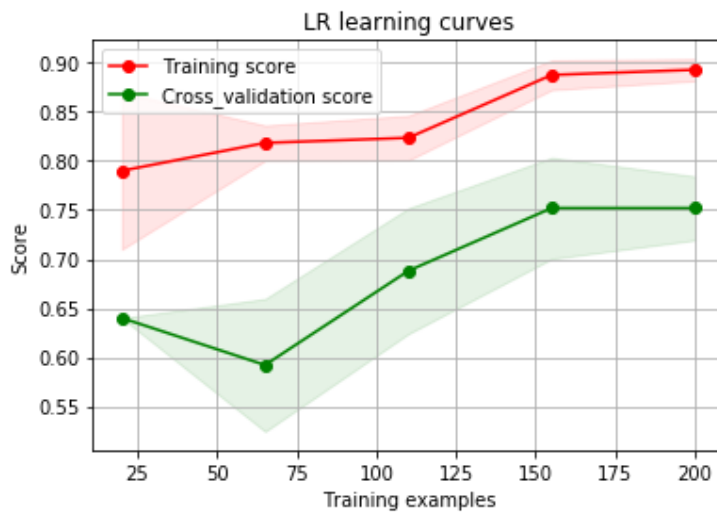
40

In [32]:

```
1 top_features = [i[1:] for i in eli5.formatters.as_dataframe.explain_weights_df(
2 X_train_new = train[top_features]
```

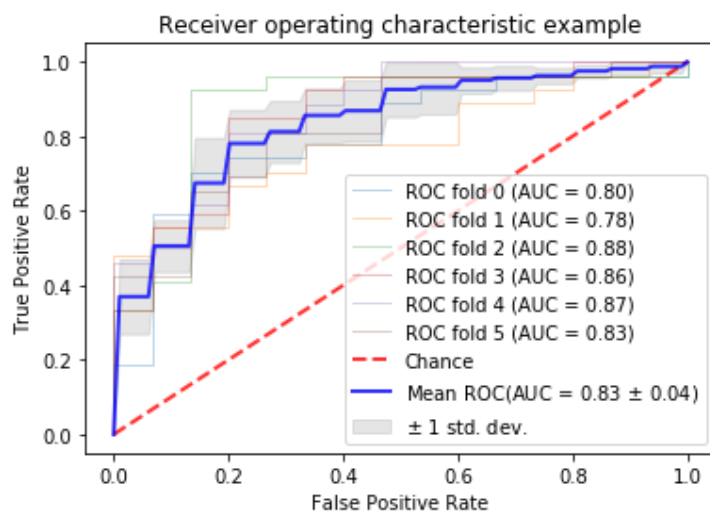
In [33]:

```
1 learningCurve3 = plot_learning_curve(log_p0, "LR learning curves", X_train_new,
```



In [34]:

```
1 roc3 = plot_roc(log_p0,X_train_new)
```



In [35]:

```
1 log_p0.fit(X_train_new, y_train)
2 X_test_new = test[top_features]
3 log_pred3 = log_p0.predict_proba(X_test_new)[: ,1]
4 submission1 = pd.DataFrame({'id':test['id'],
5                             'target':log_pred3})
6 submission1.to_csv('submission1.csv', index = False)
```


In [36]:

```
1 X_test_new.head()
```

Out[36]:

	33	65	199	101	226	24	176	30	17	201	...	80	108
0	1.988	-1.010	-0.298	1.464	0.540	0.183	-1.283	-1.003	0.764	-0.488	...	1.198	-0.639
1	0.543	-0.781	0.961	-0.981	0.476	0.475	-0.670	1.077	-2.107	-0.426	...	-0.421	0.649
2	-1.191	-0.529	0.329	0.266	-0.751	0.427	-1.331	-0.036	-1.039	-0.847	...	-1.030	1.119
3	0.542	0.754	-0.336	0.321	-1.386	0.173	-1.506	-0.374	0.857	-1.032	...	0.560	-1.875
4	0.635	-1.210	-2.235	-0.174	-1.592	1.130	2.212	0.794	0.006	1.718	...	-3.389	0.392

5 rows × 40 columns

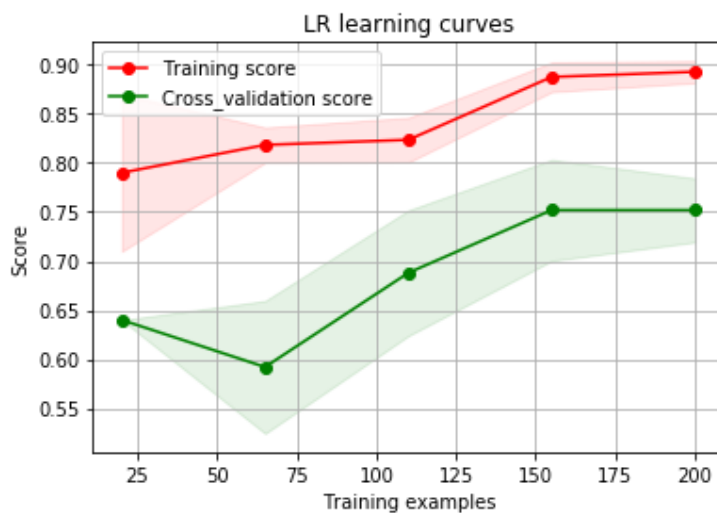
Add new statistics.

In [37]:

```
1 train['mean'] = train.mean(1)
2 train['std'] = train.std(1)
3 test['mean'] = test.mean(1)
4 test['std'] = test.std(1)
5 X_train_add = train[top_features + ['mean']]
6 X_test_add = test[top_features + ['mean']]
```

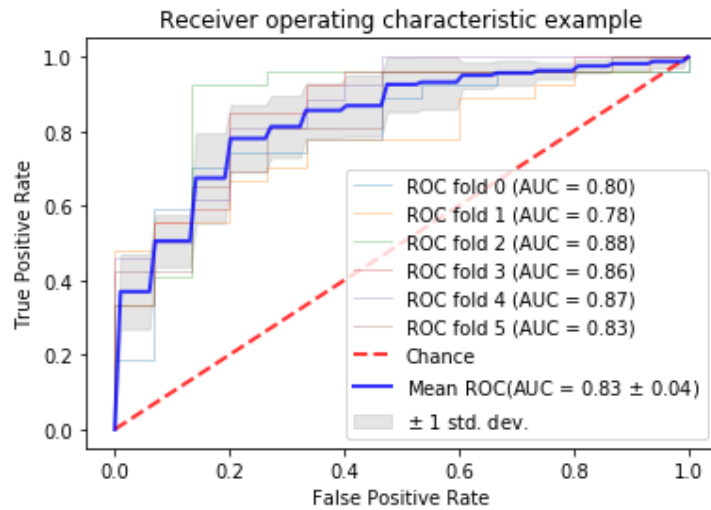
In [38]:

```
1 learningCurve4 = plot_learning_curve(log_p0, "LR learning curves", X_train_add,
```



In [39]:

```
1 roc4 = plot_roc(log_p0,X_train_add)
```



In [40]:

```
1 log_p0.fit(X_train_add, y_train)
2 log_pred4 = log_p0.predict_proba(X_test_add)[:,-1]
3 submission2 = pd.DataFrame({'id':test['id'],
4                             'target':log_pred4})
5 submission2.to_csv('submission2.csv', index = False)
```

Decison Tree

In [41]:

```

from sklearn.tree import DecisionTreeClassifier
X_train = train.drop(['id', 'target'], axis = 1)
y_train = train['target']
X_test = test.drop(['id'], axis = 1)
tree = DecisionTreeClassifier()
params = {'criterion':['gini','entropy'],
          'max_depth':[1,3,5,7,10],
          'class_weight': ['balanced', None]}
trees = GridSearchCV(tree, params, cv = StratifiedKFold(n_splits = 5), verbose = 1)
trees.fit(X_train, y_train)
tree_best = trees.best_estimator_
print(tree_best)
print(trees.best_score_)

```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,

max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=

None,

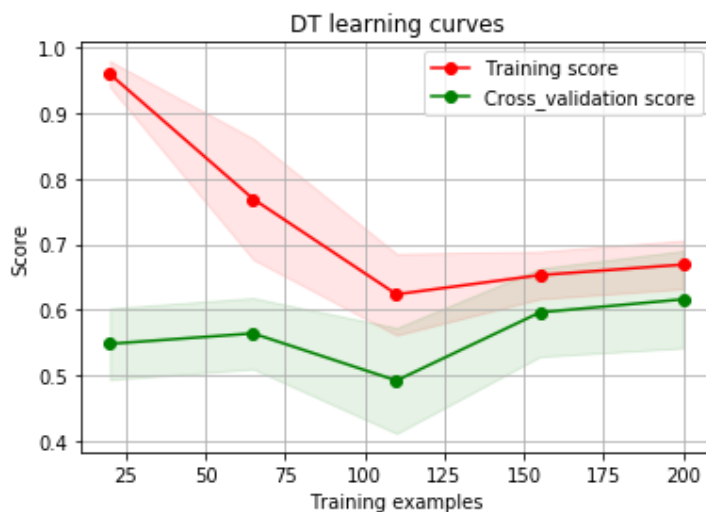
splitter='best')

0.634375

[Parallel(n_jobs=-1)]: Done 100 out of 100 | elapsed: 0.6s finished

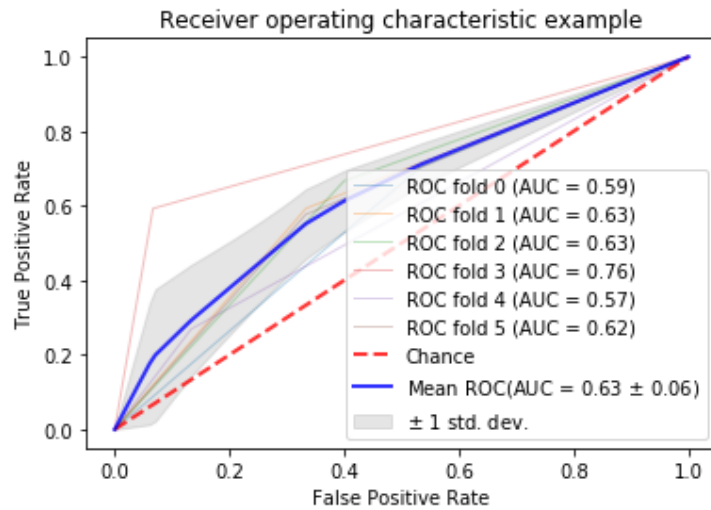
In [42]:

```
1 learningCurve5 = plot_learning_curve(tree_best, "DT learning curves", X_train, y_train)
```



In [43]:

```
1 roc5 = plot_roc(tree_best,X_train_add)
```



Decision tree is not suitable for this dataset.

Lasso Regression

In [44]:

```

1  from sklearn.linear_model import Lasso
2  X_train = train.drop(['id', 'target'], axis = 1)
3  y_train = train['target']
4  X_test = test.drop(['id'], axis = 1)
5  las = Lasso(alpha=0.031, tol=0.01, random_state=42, selection='random')
6
7  params = {
8      'alpha' : [0.022, 0.021, 0.02, 0.019, 0.023, 0.024, 0.025, 0.026, 0
9      'tol'    : [0.0013, 0.0014, 0.001, 0.0015, 0.0011, 0.0012, 0.0016, 0
10     }
11  las_ss = GridSearchCV(las, params, cv = StratifiedKFold(n_splits = 5), verbose =
12
13  las_ss.fit(X_train, y_train)
14
15  las_best = las_ss.best_estimator_
16
17  print(las_ss)
18  print(las_ss.best_score_)

```

Fitting 5 folds for each of 88 candidates, totalling 440 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

GridSearchCV(cv=StratifiedKFold(n_splits=5, random_state=None, shuffle=False),

```

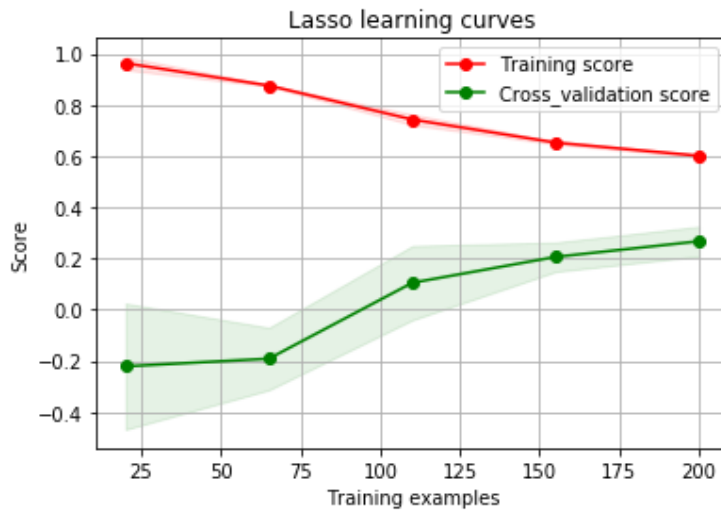
    error_score='raise-deprecating',
    estimator=Lasso(alpha=0.031, copy_X=True, fit_intercept=True, max_iter=1000,
    normalize=False, positive=False, precompute=False, random_state=42,
    selection='random', tol=0.01, warm_start=False),
    fit_params=None, iid='warn', n_jobs=-1,
    param_grid={'alpha': [0.022, 0.021, 0.02, 0.019, 0.023, 0.024,
0.025, 0.026, 0.027, 0.029, 0.031], 'tol': [0.0013, 0.0014, 0.001, 0.0015, 0.0011, 0.0012, 0.0016, 0.0017]},
    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
    scoring='roc_auc', verbose=1)
0.8350694444444444

```

[Parallel(n_jobs=-1)]: Done 440 out of 440 | elapsed: 0.6s finished

In [45]:

```
1 learningCurve6 = plot_learning_curve(las_best, "Lasso learning curves", X_train,
```

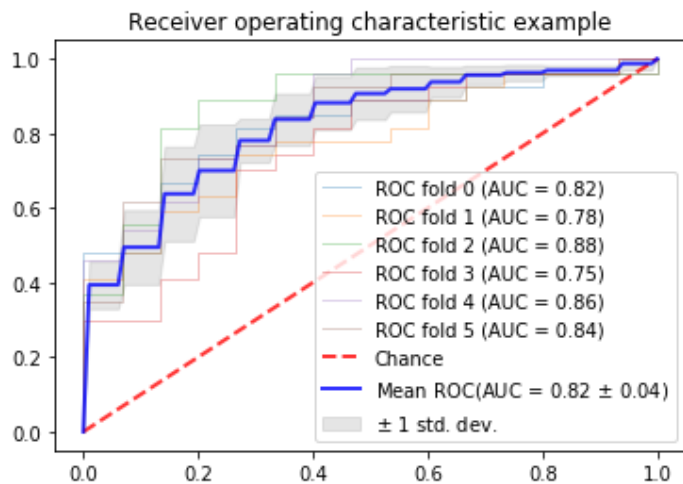


In [46]:

```
1 def plot_roc0(clf, X = X_train, y = y_train, n = 6):
2     tprs = []
3     aucs = []
4     mean_fpr = np.linspace(0, 1, 100)
5     i = 0
6     classifier = clf
7     cv = StratifiedKFold(n_splits = n)
8     for train, test in cv.split(X,y):
9         probas_ = classifier.fit(X.iloc[train], y.iloc[train]).predict(X.iloc[te
10         fpr, tpr, thresholds = roc_curve(y[test], probas_)
11         tprs.append(interp(mean_fpr, fpr, tpr))
12         tprs[-1][0] = 0.0
13         roc_auc = auc(fpr,tpr)
14         aucs.append(roc_auc)
15         plt.plot(fpr, tpr, lw = 1, alpha = 0.3, label = 'ROC fold %d (AUC = %0.2
16         i += 1
17     plt.plot([0, 1], [0, 1], linestyle = '--', lw =2, color = 'r', label = 'Chan
18     mean_tpr = np.mean(tprs, axis = 0)
19     mean_tpr[-1] = 1.0
20     mean_auc = auc(mean_fpr, mean_tpr)
21     std_auc = np.std(aucs)
22     plt.plot(mean_fpr, mean_tpr, color = 'b', label = r'Mean ROC(AUC = %0.2f $ \p
23     std_tpr = np.std(tprs, axis = 0)
24     tprs_upper = np.minimum(mean_tpr + std_tpr,1)
25     tprs_lower = np.maximum(mean_tpr - std_tpr,0)
26     plt.fill_between(mean_fpr, tprs_lower, tprs_upper, color = 'grey', alpha = 0.
27     plt.title('Receiver operating characteristic example')
28     plt.legend(loc = 'lower right')
29     plt.show()
```

In [47]:

```
1 roc6 = plot_roc0(las_best,X_train)
```



In [48]:

```
1 las_best.fit(X_train_add, y_train)
2 las_best_pred = las_best.predict(X_test_add)
3 submission3 = pd.DataFrame({'id':test['id'],
4                             'target':las_best_pred})
5 submission3.to_csv('submission3.csv', index = False)
```

In []:

```
1
```