

Feature scaling

Feature scaling is a very important component of machine learning algorithms that are dependent on distance measures (euclidean, [chebyshev](#), Manhattan distance etc). Such algorithms include Linear and Logistic Regression, SVN and Neural Networks. Other algorithms like Naive Bayes or Decision Trees can work without previously scaling the features.

The reason is that unless normalization occurs, then the whole distance measure will be dominated by the greatest feature. Suppose for instance that you have a feature vector of two features, namely age and salary. Because age will be in the values raging from 15-80 and salary from 5k-10M the whole distance will be totally dominated by the salary.

To avoid this we perform feature scaling.

```
In [31]: # The standard scaler will return numpy arrays but this will not have the columns or any index
# which we need for our model! We will therefore use some intermediate variables.
scaler = StandardScaler()

train_features2 = pd.DataFrame(scaler.fit_transform(train_features))
# Set column names
train_features2.columns = train_features.columns
# Set indices
train_features2.index = train_features.index.values

test_features2 = pd.DataFrame(scaler.transform(test_features))
test_features2.columns = test_features.columns
test_features2.index = test_features.index.values

In [32]: train_features = train_features2
test_features = test_features2
```

5. Training The Model

Generally speaking, regularization is a method used to damping the overfitting of a model. We introduce a small amount of bias into the training set so that we might perform better on unseen data. We will set the model to be a lasso regularization. We have too many features that came up as a result of the one-hot encoding of the features feature.

Lasso regularization will ensure that any features that are redundant are discarded.

```
In [33]: # A random state ensures that our results are replicable
classifier = LogisticRegression(random_state=3, penalty='l1', solver='liblinear')

# classifier = LogisticRegression(random_state=3) Achieves same results. In the docs of sklearn it is m
entioned
# that regularization happens by default
classifier.fit(train_features, train_labels)

Out[33]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, l1_ratio=None, max_iter=100,
                             multi_class='auto', n_jobs=None, penalty='l1',
                             random_state=3, solver='liblinear', tol=0.0001, verbose=0,
                             warm_start=False)
```

6. Make predictions on the test data

```
In [34]: predictions = classifier.predict(test_features)

Out[34]: array([0, 0, 0, ..., 0, 1, 0], dtype=int64)
```

7. Evaluating the Model

```
In [35]: print("The accuracy of our model is {}% ".format(100* round(accuracy_score(test_labels, predictions), 2)
))
print("Precision {}% ".format(100*round(precision_score(test_labels, predictions), 2))
print("Recall {}% ".format(100* round(recall_score(test_labels, predictions), 2))
print("F1 score {}% ".format(100*round(f1_score(test_labels, predictions), 2))

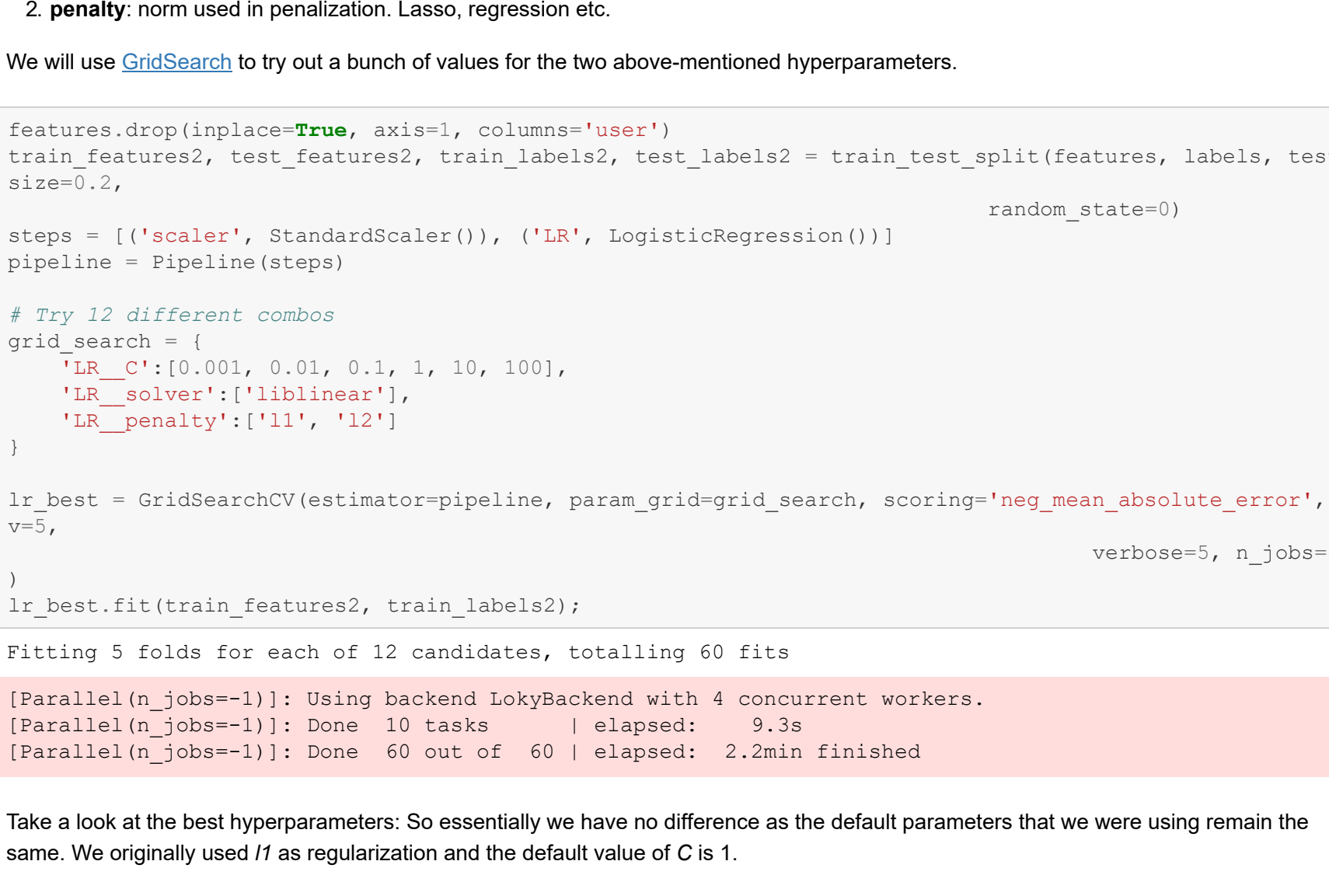
The accuracy of our model is 89.0%
Precision 97.0%
Recall 83.0%
F1 score 89.0%
```

A precision and accuracy of 89% is pretty good. The other measures also support this.

Confusion Matrix

The confusion matrix is a great tool to evaluate our model

```
In [36]: confusionMat = confusion_matrix(test_labels, predictions)
plt.figure(figsize=(9,7))
sns.heatmap(confusionMat, square=True, annot=True, fmt='d', cbar=True)
plt.xlabel("Predicted")
plt.ylabel("Ground Truth")
plt.show()
```



	0	1
0	3649	134
1	863	4095

We can indeed tell that our model is doing a good job!

The biggest problem is with those 134 users which were predicted as having enrolled but in reality they didn't (False Positives).

8. Adjusting the Model

Essentially there are three things that we can do

1. Acquiring more data: As things stand we cannot acquire more
2. Playing with Hyperparameters:
3. Choosing a different model altogether: here we only use Logistic Regression

K-Fold cross validation

[K-fold](#) cross validation partitions our data in k-folds (for instance 10) and it iterates in such a way that each fold is used for testing once with the rest of the folds being testing data.

```
In [37]: accuracies = cross_val_score(estimator=classifier, X=train_features, y=train_labels, cv=10)
print("Accuracy with a 10-fold cross validation: ", round(100*np.mean(accuracies), 2))

Accuracy with a 10-fold cross validation: 88.82

It looks like the performance is the same!
```

Hyperparameter tuning

Also using [Pipeline](#). The benefits of [pipelins](#) are that it enforces a desired order of application steps and we do not have to perform data scaling on train and test sets separately as it automatically makes sure of that during cross validation.

We will tune two hyperparameters ([here](#) the docs link):

1. **C**: The strength of regularization. A smaller C means as stronger regularization.
2. **penalty**: norm used in penalization. Lasso, regression etc.

We will use [GridSearch](#) to try out a bunch of values for the two above-mentioned hyperparameters.

```
In [38]: features.drop(inplace=True, axis=1, columns='user')
train_features2, test_features2, train_labels2, test_labels2 = train_test_split(features, labels, test_size=0.2,
                                                                                random_state=0)

pipeline = Pipeline(steps

# Try 12 different combos
grid_search = GridSearchCV(estimator=pipeline, param_grid=grid_search, scoring='neg_mean_absolute_error', c
v=5,
                           verbose=5, n_jobs=-1)

lr_best = GridSearchCV(estimator=pipeline, param_grid=grid_search, scoring='neg_mean_absolute_error', c
v=5,
                           verbose=5, n_jobs=-1)

lr_best.fit(train_features2, train_labels2):

Fitting 5 folds for each of 12 candidates, totalling 60 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks | elapsed: 9.3s
[Parallel(n_jobs=-1)]: Done 60 out of 60 | elapsed: 2.2min finished
```

Take a look at the best hyperparameters. So essentially we have no difference as the default parameters that we were using remain the same. We originally used *l1* as regularization and the default value of *C* is 1.

Hyperparameter tuning did not improve the model's performance. The results are the same with both *StandardScaler* as well as *MinMaxScaler*

```
In [39]: lr_best.best_params_

Out[39]: {'lr_c': 1, 'lr_penalty': 'l1', 'lr_solver': 'liblinear'}
```

Principal Component Analysis

[PCA](#) is a popular way to reduce the dimensionability of our data. With a feature vector of 72 features there might be redundant features. Let's use PCA to find out more.

When to use PCA?

1. You want to reduce the number of features in the feature vector but are not sure which ones to take out
2. You want to ensure that features are independent from each-other
3. You do not mind making the new features less interpretable

However, if the answer to the 3rd question is no, do not use it.

```
In [40]: print("Our feature vector consists of {} features".format(features.shape[1]))

Our feature vector consists of 72 features

In [41]: scale = StandardScaler()
scaled_features = scale.fit_transform(features)

pca = PCA(n_components=10)
pca.fit(scaled_features)

new_features = pca.transform(scaled_features)
print("Our new feature vector consists of {} orthogonal hyperplanes".format(new_features.shape[1]))

print("Total variance explained by the new features {}% ".format(round(100*sum(pca.explained_variance_ratio_), 2)))

Our new feature vector consists of 10 orthogonal hyperplanes
Total variance explained by the new features 47.47%
```

Hyperparameter Tuning with PCA

We now repeat hyperparameter tuning using PCA

```
In [42]: pipeline = [('scaler', StandardScaler()), ('PCA', PCA()), ('LR', LogisticRegression())]
pipeline = Pipeline(pipeline)

grid_search = GridSearchCV(estimator=pipeline, param_grid=grid_search, scoring='neg_mean_absolute_error', c
v=5,
                           verbose=5, n_jobs=-1)

lr_best = GridSearchCV(estimator=pipeline, param_grid=grid_search, scoring='neg_mean_absolute_error', c
v=5,
                           verbose=5, n_jobs=-1)

lr_best.fit(train_features2, train_labels2):

Fitting 3 folds for each of 160 candidates, totalling 480 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks | elapsed: 5.0s
[Parallel(n_jobs=-1)]: Done 64 tasks | elapsed: 29.5s
[Parallel(n_jobs=-1)]: Done 154 tasks | elapsed: 1.1min
[Parallel(n_jobs=-1)]: Done 280 tasks | elapsed: 2.3min
[Parallel(n_jobs=-1)]: Done 442 tasks | elapsed: 4.4min
[Parallel(n_jobs=-1)]: Done 480 out of 480 | elapsed: 6.9min finished

In [43]: lr_best.best_params_

Out[43]: {'lr_c': 0.1,
          'lr_penalty': 'l1',
          'lr_solver': 'liblinear',
          'pca_n_components': 65}

In [44]: predictions3 = lr_best.best_estimator_.predict(test_features2)
print(round(100*accuracy_score(test_labels, predictions3), 2))

88.68
```

Although PCA can grant us with a lower number of independent features (orthogonal to each-other), in this case it does not seem to improve performance.

9. Interpret the model and report results visually and numerically

```
In [45]: # Combine the results into a final data frame
final_results = pd.concat([test_uid, test_labels], axis=1).dropna()
final_results['predicted'] = predictions
# Reset the indices
final_results = final_results.reset_index(drop=True)
final_results
```

```
Out[45]:
```

	user	enrolled	predicted
0	251410	0	0
1	132536	1	0
2	38654	0	0
3	122925	0	0
4	122887	1	1
...
8738	230964	1	1
8739	128051	0	0
8740	242381	0	0
8741	247726	1	1
8742	200463	0	0

8743 rows × 3 columns

Conclusion

Our model will label every new user as highly likely or unlikely to subscribe (and all this within the first 48h). This will ensure that the company only narrows the marketing efforts on those users who appear to be unlikely to subscribe.

In this small project we saw and end-to-end-machine learning example! We saw the data analysis pipeline and finally applied a logistic regression model which achieved a precision of **89%**.