# AI4FutureWorkForce Data Processing

Exploration of predicting course completion using sklearn and Tensorflow

## Dependencies

These are all the libraries we need to run this notebook

```python
1  import pandas as pd
2  import plotly.offline as plt
3  import _pickle as pickle
4  import os
5  import numpy as np
6
7  import tensorflow as tf
8
9  import tensorflow.keras as keras
10
11 import plotly.graph_objs as go
12 import itertools
13
14 from sklearn.naive_bayes import GaussianNB
15 from sklearn import svm
16 from sklearn.metrics import accuracy_score
17 from sklearn import linear_model
18 from sklearn.ensemble import RandomForestClassifier
19 from sklearn.model_selection import GridSearchCV
20 from sklearn import preprocessing
21 from sklearn.model_selection import train_test_split
22 from collections import defaultdict
23 from sklearn.metrics import classification_report
24 from numpy.random import seed
25 seed(42)
26 from tensorflow import set_random_seed
27 set_random_seed(42)
28 from tensorflow.python.client import device_lib
29
30 plt.init_notebook_mode(connected=True)
```

```python
1  # Verify that the TF device we're using is what we expect
2  device_lib.list_local_devices()
```

```
[name: "/device:CPU:0"
 device_type: "CPU"
 memory_limit: 268435456
 locality {
 }
 incarnation: 5457231463178488201, name: "/device:XLA_CPU:0"
 device_type: "XLA_CPU"
 memory_limit: 17179869184
 locality {
 }
 incarnation: 7095777417938979866
 physical_device_desc: "device: XLA_CPU device"]
```

## Load Data

```
1 # Whether or not to use raw or processed dfs
2 PREPROCESS = 0
3 # Whether to fill fee columns or just encode as NaN
4 FILL = 0
5 SHUFFLE = True
6
7 NUM_CLASSES = 2
```

```
1  # List all columns of the tables in data
2  column_list = ['Age (Birthday Masked)','Income','Education',
3                 'MAX(Learner Test Score)','Primary Interest In Course',
4                 'Hours Coded','How Many Hours A Week Can You Commit To Class',
5                 'Promise Zone Indicator','Hacker Rank Score', 'Location; number']
6
7  # List desired columns for train/test/validation
8  desired_columns = ['Age (Birthday Masked)','Income','Education',
9                 'MAX(Learner Test Score)','Primary Interest In Course',
10                'Hours Coded','How Many Hours A Week Can You Commit To Class',
11                'Promise Zone Indicator','Hacker Rank Score', 'Location; number']
12
13 # Load raw data and run pre-processing
14 if PREPROCESS == 1:
15     df_data = pickle.load( open( "df_data.p", "rb" ) )
16     df_labels = pickle.load( open( "df_labels.p", "rb" ) )
17
18     df_validata = pickle.load( open( "df_validata.p", "rb" ) )
19     df_valilabels = pickle.load( open( "df_valilabels.p", "rb" ) )
20
21     df_valilabels = df_valilabels.add(1)
22     df_labels = df_labels.add(1)
23
24     X_train, X_test, X_val, y_train, y_test, y_val = data_prepare.prepare
25     (column_list, desired_columns, df_data, df_validata, df_labels,
26      df_valilabels, SHUFFLE, FILL)
27
28     print(X_train.shape, X_test.shape, X_val.shape)
29
30 # Otherwise load pre-processed data
31 else:
32     X_train = pickle.load( open( "X_train.p", "rb" ) )
33     X_test = pickle.load( open( "X_test.p", "rb" ) )
34     X_val = pickle.load( open( "X_val.p", "rb" ) )
35
36     y_train = pickle.load( open( "y_train.p", "rb" ) )
37     y_test = pickle.load( open( "y_test.p", "rb" ) )
38     y_val = pickle.load( open( "y_val.p", "rb" ) )
```

Checking that the dataframe looks like it should:

```
1 X_train.head()
```

| | Age (Birthday Masked) | Income | MAX(Learner Test Score) | Hours Coded | How Many Hours A Week Can You Commit To Class | Promise Zone Indicator | Location; number |
|---|---|---|---|---|---|---|---|
| 1341 | 0.353798 | 1.832671 | 0.988618 | -0.055229 | -2.027239 | -0.405866 | 0.571825 |
| 1649 | -1.263828 | -0.299421 | 0.018359 | 0.959409 | 0.381102 | -0.405866 | 0.571825 |
| 998 | -0.940303 | -1.152258 | 1.312038 | 0.959409 | 0.381102 | -0.405866 | 0.571825 |

```
1 y_train
```

```
array([1, 1, 1, ..., 1, 0, 1])
```

# Classification Functions

```python
1  class EpochTrack(keras.callbacks.Callback):
2      def on_epoch_end(self, epoch, logs):
3          if epoch % 100 == 0:
4              print('')
5          print('.', end='')
6
7
8  def eval_model(model, test_data, test_labels, VERBOSE, NUM_CLASSES):
9      test_labels_sparse = keras.utils.to_categorical(test_labels, NUM_CLASSES)
10     [loss, mae, acc] = model.evaluate(test_data, test_labels_sparse,
11                                       verbose=VERBOSE)
12
13     print("Mean Abs Error:\t{:7.2f}".format(mae * 1000))
14     print("Loss:\t\t", loss)
15     print("Accuracy:\t", acc)
16
17     preds = model.predict(test_data)
18     y_preds = preds.argmax(axis=-1)
19     print(classification_report(test_labels, y_preds))
20
21     return acc, preds
22
23
24 def plot_history(history):
25     # Define each data series
26     trace1 = go.Scatter(x=history.epoch, y=np.array
27                         (history.history['mean_absolute_error']),
28                         name='Training Loss')
29     trace2 = go.Scatter(x=history.epoch, y=np.array
30                         (history.history['val_mean_absolute_error']),
31                         name='Val Loss')
32
33     # Add each series
34     data = [trace1, trace2]
35
36     # Define graph layout
37     layout = go.Layout(
38         title='Training History',
39         xaxis=dict(title='Epoch'),
40         yaxis=dict(title='Mean Abs Error'))
41
42     fig = go.Figure(data=data, layout=layout)
43     plt.iplot(fig)
44
45 def show_incorrect(y_classes, data, act):
46     incorrects = np.nonzero(y_classes.reshape((-1,)) != act)
47     test_X = data.iloc[incorrects]
48
49     test_y = pd.DataFrame(y_classes[tuple(incorrects)])
50     test_y.columns = ['Prediction']
51
52     act_y = pd.DataFrame(act[tuple(incorrects)])
53     act_y.columns = ['Actual']
54
55     test_X.index = range(len(test_X))
56
57     frames = [test_X, test_y, act_y]
58     df_pred = pd.concat(frames, axis=1)
59
60     return df_pred
61
62
63 def skl_evaluate(model, desired_columns, X, y):
```

```
64      score = model.score(X, y.astype(int))
65
66      preds = model.predict(X)
67      print(classification_report(y, preds))
68
69      print("Accuracy on dataset:\t %f\n" % score)
70
71      importances = model.feature_importances_
72
73      indices = np.argsort(importances)[::-1]
74
75      # Print the feature ranking
76      print("Feature ranking:")
77      for f in range(X.shape[1]):
78          print("%d. %s \t(%f)" % (f + 1, desired_columns[indices[f]], importances[indi
79
80      return score
```

## Classify

```
1 nb = GaussianNB()
2 nb.fit(X_train, y_train.astype(int))
```

⮕  GaussianNB(priors=None, var_smoothing=1e-09)

```
1 X_test.iloc[0, :]
```

⮕  Age (Birthday Masked)                        -1.048144
    Income                                        0.553416
    MAX(Learner Test Score)                       0.341779
    Hours Coded                                  -0.055229
    How Many Hours A Week Can You Commit To Class  1.585273
    Promise Zone Indicator                       -0.405866
    Location; number                              0.571825
    Name: 522, dtype: float64

```
1 nb_score = nb.score(X_test, y_test.astype(int))
2 print("Naive Bayes accuracy on test:\t %f" % nb_score)
3
4 nb_score_val = nb.score(X_val, y_val.astype(int))
```

⮕  Naive Bayes accuracy on test:     0.568765

## SVM

```
1 clf = svm.SVC(kernel='linear', C=1, random_state=42)
2 clf.fit(X_train, y_train.astype(int))
```

⮕  SVC(C=1, cache_size=200, class_weight=None, coef0=0.0,
        decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
        kernel='linear', max_iter=-1, probability=False, random_state=42,
        shrinking=True, tol=0.001, verbose=False)

```
1 clf_score = clf.score(X_test, y_test.astype(int))
2 print("SVM accuracy on test:\t %f" % clf_score)
3
```

```
4  clf_score_val = clf.score(X_val, y_val.astype(int))
```

➥   SVM accuracy on test:    0.596737

## Linear Regression

```
1  linr = linear_model.LinearRegression(n_jobs = -1)
2  linr.fit(X_train, y_train.astype(int))
```

➥   LinearRegression(copy_X=True, fit_intercept=True, n_jobs=-1, normalize=False)

```
1  linr_score = linr.score(X_test, y_test.astype(int))
2  print("Linear Regression coefficient of determination, 1 is perfect prediction:\t %f"
3
4  linr_score_val = linr.score(X_val, y_val.astype(int))
```

➥   Linear Regression coefficient of determination, 1 is perfect prediction:        0.

## Random Forest

```
1  rf = RandomForestClassifier(criterion='gini', max_depth=5,
2                              min_samples_leaf=5, min_samples_split=2,
3                              n_estimators = 220, oob_score=True,
4                              max_features=0.5, n_jobs = -1, random_state=42)
5
6  rf.fit(X_train, y_train.astype(int))
```

➥   RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                           max_depth=5, max_features=0.5, max_leaf_nodes=None,
                           min_impurity_decrease=0.0, min_impurity_split=None,
                           min_samples_leaf=5, min_samples_split=2,
                           min_weight_fraction_leaf=0.0, n_estimators=220,
                           n_jobs=-1, oob_score=True, random_state=42, verbose=0,
                           warm_start=False)

```
1  print("Random Forest - Test set")
2  rf_score = skl_evaluate(rf, desired_columns, X_test, y_test)
3
4  print("Random Forest - Validation set")
5  rf_score_val = skl_evaluate(rf, desired_columns, X_val, y_val)
```

➥

```
Random Forest - Test set
              precision    recall  f1-score   support

           0       0.65      0.79      0.72       241
           1       0.64      0.46      0.54       188

    accuracy                           0.65       429
   macro avg       0.64      0.63      0.63       429
weighted avg       0.65      0.65      0.64       429
```

```
Accuracy on dataset:     0.648019
```

```
Feature ranking:
1. Education    (0.275820)
2. Age (Birthday Masked)        (0.238993)
3. MAX(Learner Test Score)      (0.222282)
4. Income       (0.123961)
5. Primary Interest In Course   (0.076364)
6. How Many Hours A Week Can You Commit To Class        (0.034888)
7. Hours Coded  (0.027692)
Random Forest - Validation set
              precision    recall  f1-score   support

           0       0.78      0.85      0.82       308
           1       0.51      0.39      0.44       121
```

## ▾ Random Forest Gridsearch Optimisation

```
weighted avg       0.70      0.72      0.71       429
```

```python
1  rf=RandomForestClassifier(random_state=42)
2
3  param_grid = {
4      'n_estimators': [100, 120, 140, 160, 180, 200, 220, 240, 260],
5      'min_samples_leaf': [3, 5, 7],
6      'min_samples_split': [2, 3, 4, 5, 6],
7      'max_depth' : [5, 10, 15, 20, 25, 30, 35, 40, 45],
8      'criterion' :['gini', 'entropy']
9  }
10
11
12 CV_rf = GridSearchCV(estimator=rf, param_grid=param_grid, cv= 5, n_jobs = -1,
13                 verbose = 2)
14 CV_rf.fit(X_train, y_train.astype(int))
15 CV_rf.best_params_
```

This is sample output from this step:

```
Fitting 5 folds for each of 2430 candidates, totalling 12150 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done   37 tasks      | elapsed:    8.8s
[Parallel(n_jobs=-1)]: Done  158 tasks      | elapsed:   33.4s
[Parallel(n_jobs=-1)]: Done  361 tasks      | elapsed:  1.3min
[Parallel(n_jobs=-1)]: Done  644 tasks      | elapsed:  2.3min
[Parallel(n_jobs=-1)]: Done 1009 tasks      | elapsed:  3.7min
[Parallel(n_jobs=-1)]: Done 1454 tasks      | elapsed:  5.5min
[Parallel(n_jobs=-1)]: Done 1981 tasks      | elapsed:  7.7min
[Parallel(n_jobs=-1)]: Done 2588 tasks      | elapsed: 10.2min
[Parallel(n_jobs=-1)]: Done 3277 tasks      | elapsed: 13.0min
[Parallel(n_jobs=-1)]: Done 4046 tasks      | elapsed: 16.2min
```

```
[Parallel(n_jobs=-1)]: Done 4897 tasks       | elapsed: 19.7min
[Parallel(n_jobs=-1)]: Done 5828 tasks       | elapsed: 23.5min
[Parallel(n_jobs=-1)]: Done 6841 tasks       | elapsed: 27.4min
[Parallel(n_jobs=-1)]: Done 7934 tasks       | elapsed: 32.2min
[Parallel(n_jobs=-1)]: Done 9109 tasks       | elapsed: 37.4min
[Parallel(n_jobs=-1)]: Done 10364 tasks       | elapsed: 43.0min
[Parallel(n_jobs=-1)]: Done 11701 tasks       | elapsed: 49.0min
[Parallel(n_jobs=-1)]: Done 12150 out of 12150 | elapsed: 50.9min finished
```

The optimum parameters are: {'criterion': 'gini', 'max_depth': 5, 'min_samples_leaf': 5, 'min_samples_split': 2, 'n_estimators': 220}

```python
1 CV_rf_score = CV_rf.score(X_test, y_test.astype(int))
2 print("Optimised Random Forest accuracy on test:\t %f" % CV_rf_score)
```

## Random Forest - Find Optimal Parameters

```python
1  def skl_run_combos(combo_columns, model, df_data, df_labels, SHUFFLE, FILL):
2
3      cool_list = list()
4
5      X_train_combo = X_train
6      X_test_combo = X_test
7
8      y_train_combo = y_train
9      y_test_combo = y_test
10
11     for L in range(0, len(combo_columns) + 1):
12         combinations = list(itertools.combinations(combo_columns, L))[1:]
13         for combo in combinations:
14             combo = list(combo)
15
16             model.fit(X_train_combo[combo], y_train)
17             rf_score = model.score(X_test_combo[combo], y_test)
18             output = str(rf_score) + " " + str(combo)
19             print(output)
20             cool_list.append(output)
21
22     cool_list.sort()
23
24     with open('outputs.txt', 'w') as f:
25         for item in cool_list:
26             f.write("%s\n" % item)
27
28 def normalise_df(train_df, test_df, AXIS, val_df=1):
29     mu = train_df.mean(axis=AXIS)
30     sd = train_df.std(axis=AXIS)
31
32     train_df = (train_df - mu) / sd
33     test_df = (test_df - mu) / sd
34     val_df = (val_df - mu) / sd
35
36     return train_df, test_df, val_df
37
```

```python
1 df_data = pickle.load( open( "df_data.p", "rb" ) )
2 df_labels = pickle.load( open( "df_labels.p", "rb" ) )
3
4 combo_columns = ['Age (Birthday Masked)','Income','Education',
```

```
 5                     'MAX(Learner Test Score)','Primary Interest In Course',
 6                     'Hours Coded','How Many Hours A Week Can You Commit To Class',
 7                     'Promise Zone Indicator','Hacker Rank Score',
 8                     'Location; number']
 9
10 model = RandomForestClassifier(criterion='gini', max_depth=5,
11                                min_samples_leaf=5, min_samples_split=2,
12                                n_estimators = 220, oob_score=True,
13                                max_features=0.5, n_jobs = -1, random_state=42)
14
15 print('Accuracy on test, [columns used]')
16 skl_run_combos(combo_columns, model, df_data, df_labels,
17                               SHUFFLE, FILL)
```

This is a sample output from this step:

```
0.648018648 Income', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator', 'Location; number']
0.648018648 Age (Birthday Masked)', 'Income', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator', 'Location; number']
0.648018648 'Age (Birthday Masked)', 'Income', 'Education', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator', 'Hacker Rank Score', 'Location; number']
0.645687646 'Income', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator', 'Hacker Rank Score', 'Location; number']
0.645687646 'Age (Birthday Masked)', 'Income', 'Education', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator', 'Hacker Rank Score', 'Location; number']
0.643356643 'Income', 'Hours Coded', 'Hacker Rank Score', 'Location; number']
0.643356643 'Age (Birthday Masked)', 'Income', 'Education', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class']
0.643356643 'Age (Birthday Masked)', 'Education', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Hacker Rank Score', 'Location; number']
0.641025641 'Income', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Hacker Rank Score', 'Location; number']
0.641025641 'Age (Birthday Masked)', 'Income', 'Education', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Hacker Rank Score']
0.641025641 'Age (Birthday Masked)', 'Education', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Hacker Rank Score']
0.641025641 'Income', 'Education', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator']
0.641025641 'Income', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Hacker Rank Score', 'Location; number']
0.641025641 'Age (Birthday Masked)', 'Income', 'Education', 'MAX(Learner Test Score)', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Location; number']
0.641025641 'Age (Birthday Masked)', 'Education', 'Hours Coded', 'How Many Hours A Week Can You Commit To Class', 'Promise Zone Indicator', 'Hacker Rank Score', 'Location; number']
```

## Multi Layer Perceptron

```python
 1 def build_mlp(input, NUM_CLASSES):
 2     model = keras.Sequential([
 3         keras.layers.Dense(8, activation='relu', input_shape=(input.shape[1],)),
 4         keras.layers.Dense(5, activation='relu'),
 5 #         keras.layers.Dense(3, activation='relu'),
 6         keras.layers.Dense(NUM_CLASSES, activation='softmax')
 7     ])
 8
 9     model.compile(loss='categorical_crossentropy', optimizer='adam',
10                   metrics=['mae', 'acc'])
11
12     model.summary()
13
14     return model
15
16
17 def train_mlp(dataframe, labels, model, checkpoint_path, EPOCHS):
18     # Limit the training when there are multiple epochs with little change loss
19     # The patience parameter is the amount of epochs to check for improvement
20     early_stop = keras.callbacks.EarlyStopping(monitor='val_loss',
21                                                patience=4000)
22
23     # Create checkpoint callback
24     cp_callback = keras.callbacks.ModelCheckpoint(checkpoint_path,
25                                                   save_weights_only=True,
26                                                   verbose=0)
27
28     # Track the training statistics
29     history = model.fit(dataframe, labels, epochs=EPOCHS,
30                         validation_split=0.2, verbose=0,
31                         callbacks=[early_stop, EpochTrack(), cp_callback])
32
33     print("\nEpochs: {}".format(len(history.epoch)))
34     plot_history(history)
35
36     return model
```

```python
 1 y_train_sparse = keras.utils.to_categorical(y_train, NUM_CLASSES)
```

```python
 2
 3 TRAIN = 0
 4
 5 if TRAIN == 1:
 6     checkpoint_path = "training_1/cp.ckpt"
 7     checkpoint_dir = os.path.dirname(checkpoint_path)
 8
 9     # Train MLP and save to checkpoints
10     %time mlp = train_mlp(X_train, y_train_sparse, build_mlp(X_train, NUM_CLASSES), cl
11
12 else:
13     checkpoint_path = "training_1/cp.ckpt"
14
15     checkpoint_dir = os.path.dirname(checkpoint_path)
16
17     # Load trained weights from the checkpoint path
18     mlp = build_mlp(X_train, NUM_CLASSES)
19     mlp.load_weights(checkpoint_path)
20
21 # Test
22 print("\nMLP - Test set")
23 acc_test, preds_test = eval_model(mlp, X_test, y_test, 0, NUM_CLASSES)
24
25 print("\nMLP - Validation set")
26 acc_val, preds_val = eval_model(mlp, X_val, y_val, 0, NUM_CLASSES)
```

⯈

```
WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/tensorflow/python/op
Instructions for updating:
Colocations handled automatically by placer.
```
_____
Layer (type)                    Output Shape                 Param #

## MLP every combination of parameters

dense_1 (Dense)                    (None, 5)                    45

```python
1  col_names = ['Age (Birthday Masked)','Income','Education',
2               'MAX(Learner Test Score)','Primary Interest In Course',
3               'Hours Coded','How Many Hours A Week Can You Commit To Class'
4               ,'Promise Zone Indicator','Hacker Rank Score','Location; number']
5
6  cool_list = list()
7
8  y_train_sparse = keras.utils.to_categorical(y_train, NUM_CLASSES)
9  y_val_sparse = keras.utils.to_categorical(y_val, NUM_CLASSES)
10
11 # with open('results.txt', 'w') as f:
12 #    for L in range(0, len(col_names)+1):
13 #        combinations = list(itertools.combinations(col_names, L))[4:]
14 #        for combo in combinations:
15 #            combo = list(combo)
16
17 #            mlp = train_mlp(X_train[combo], y_train_sparse,
18 #                      build_mlp(X_train[combo], NUM_CLASSES),
19 #                      checkpoint_path, EPOCHS=1000)
20 #            acc_val = eval_model(mlp, X_val[combo], y_val_sparse, 0, NUM_CLASSES)
21
22 #            output = str(acc_val) + " " + str(combo) + "\n"
23 #            print(output)
24 #            cool_list.append(output)
25 #            f.write(str(cool_list))
```

MLP - Validation set

## Test Results Summary

precision    recall  f1-score    support

```python
1  print("Naive Bayes accuracy on validation:\t\t %f" % nb_score_val)
2  print('')
3  print("SVM accuracy on validation:\t\t\t %f" % clf_score_val)
4  print('')
5  print("Linear Regression accuracy on validation:\t %f" % linr_score_val)
6  print('')
7  print("Random Forest accuracy on validation:\t\t %f" % rf_score_val)
8  print('')
9  print("MLP accuracy on validation:\t\t\t %f" % acc_val)
```

```
Naive Bayes accuracy on validation:              0.641026

SVM accuracy on validation:                      0.727273

Linear Regression accuracy on validation:        0.051558

Random Forest accuracy on validation:            0.722611

MLP accuracy on validation:                      0.668998
```

## Insight

Let's sanity check the results by comparing model predictions to the real test results.

```
1 predictions = preds_test
2
3 y_classes = predictions.argmax(axis=-1)
```

When comparing predictions to true labels we'd expect to see mostly true values for high scoring classifiers

```
1 y_classes = y_classes.flatten()
2
3 # Match is 'True' when values are the same
4 match = [y_classes[x] == y_test[x] for x in range(len(y_classes))]
5
6 # Combine and show
7 df_compare = pd.DataFrame({'Predictions': y_classes, 'Actual': y_test,
8                            'Match': match})
9
10 print('Show performance on first 5 rows')
11 df_compare.head(5)
```

⌷→ Show performance on first 5 rows

|   | Predictions | Actual | Match |
|---|---|---|---|
| **0** | 1 | 1 | True |
| **1** | 0 | 0 | True |
| **2** | 1 | 1 | True |
| **3** | 0 | 1 | False |
| **4** | 1 | 1 | True |

We can also obtain the set of rows that were incorrectly predicted

```
1 df_pred = show_incorrect(y_classes, X_test, y_test)
2
3 df_pred.sort_values('Prediction', inplace=True)
4 print('Show first 5 rows that had wrong predictions')
5 print('Values are encoded as per earlier steps')
6 df_pred.head(5)
```

⌷→ Show first 5 rows that had wrong predictions
   Values are encoded as per earlier steps

|   | Age (Birthday Masked) | Income | MAX(Learner Test Score) | Hours Coded | How Many Hours A Week Can You Commit To Class | Promise Zone Indicator | Location; number |
|---|---|---|---|---|---|---|---|
| **0** | 1.647899 | 1.832671 | -0.628481 | -0.055229 | 0.381102 | -0.405866 | 0.571825 |
| **125** | -0.185411 | -1.578677 | 0.988618 | -0.055229 | -0.823068 | -0.405866 | 0.571825 |
| **124** | -0.293252 | -0.299421 | 0.341779 | 1.974047 | 0.381102 | -0.405866 | 0.571825 |

```
1
```