

Evaluating a classification model

- What is the purpose of **model evaluation**, and what are some common evaluation procedures?
- What is the usage of **classification accuracy**, and what are its limitations?
- How does a **confusion matrix** describe the performance of a classifier?
- What **metrics** can be computed from a confusion matrix?
- How can you adjust classifier performance by **changing the classification threshold**?
- What is the purpose of an **ROC curve**?
- How does **Area Under the Curve (AUC)** differ from classification accuracy?

Review of model evaluation

- Need a way to choose between models: different model types, tuning parameters, and features
- Use a **model evaluation procedure** to estimate how well a model will generalize to out-of-sample data
- Requires a **model evaluation metric** to quantify the model performance

Model evaluation procedures

1. Training and testing on the same data

- Rewards overly complex models that "overfit" the training data and won't necessarily generalize

2. Train/test split

- Split the dataset into two pieces, so that the model can be trained and tested on different data
- Better estimate of out-of-sample performance, but still a "high variance" estimate
- Useful due to its speed, simplicity, and flexibility

3. K-fold cross-validation

- Systematically create "K" train/test splits and average the results together
- Even better estimate of out-of-sample performance
- Runs "K" times slower than train/test split

Model evaluation metrics

- **Regression problems:** Mean Absolute Error, Mean Squared Error, Root Mean Squared Error
- **Classification problems:** Classification accuracy

Classification accuracy

Pima Indians Diabetes dataset (<https://www.kaggle.com/uciml/pima-indians-diabetes-database>)

```
In [1]: # read the data into a pandas DataFrame
import pandas as pd
path = 'data/pima-indians-diabetes.data'
col_names = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi', 'pedigree', 'age', 'label']
pima = pd.read_csv(path, header=None, names=col_names)
```

```
In [2]: # print the first 5 rows of data
pima.head()
```

Out[2]:

	pregnant	glucose	bp	skin	insulin	bmi	pedigree	age	label
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

Question: Can we predict the diabetes status of a patient given their health measurements?

```
In [3]: # define X and y
feature_cols = ['pregnant', 'insulin', 'bmi', 'age']
X = pima[feature_cols]
y = pima.label
```

```
In [4]: # split X and y into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
```

```
In [5]: # train a logistic regression model on the training set
from sklearn.linear_model import LogisticRegression
logreg = LogisticRegression(solver='liblinear')
logreg.fit(X_train, y_train)
```

Out[5]: 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='l2', random_state=None, solver='liblinear', tol=0.0001, verbose=0, warm_start=False)

```
In [7]: # make class predictions for the testing set
y_pred_class = logreg.predict(X_test)
```

Classification accuracy: percentage of correct predictions

```
In [8]: # calculate accuracy
from sklearn import metrics
print(metrics.accuracy_score(y_test, y_pred_class))

0.6927083333333334
```

Null accuracy: accuracy that could be achieved by always predicting the most frequent class

```
In [9]: # examine the class distribution of the testing set (using a Pandas Series method)
y_test.value_counts()
```

```
Out[9]: 0    130
        1     62
        Name: label, dtype: int64
```

```
In [10]: # calculate the percentage of ones
y_test.mean()
```

```
Out[10]: 0.3229166666666667
```

```
In [11]: # calculate the percentage of zeros
1 - y_test.mean()
```

```
Out[11]: 0.6770833333333333
```

```
In [12]: # calculate null accuracy (for binary classification problems coded as 0/1)
max(y_test.mean(), 1 - y_test.mean())
```

```
Out[12]: 0.6770833333333333
```

```
In [13]: # calculate null accuracy (for multi-class classification problems)
y_test.value_counts().head(1) / len(y_test)
```

```
Out[13]: 0    0.677083
        Name: label, dtype: float64
```

Comparing the **true** and **predicted** response values

```
In [14]: # print the first 25 true and predicted responses
print('True:', y_test.values[0:25])
print('Pred:', y_pred_class[0:25])
```

```
True: [1 0 0 1 0 0 1 1 0 0 1 1 0 0 0 0 1 0 0 0 1 1 0 0 0]
Pred: [0 0 0 0 0 0 0 1 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0]
```

Conclusion:

- Classification accuracy is the **easiest classification metric to understand**
- But, it does not tell you the **underlying distribution** of response values
- And, it does not tell you what **"types" of errors** your classifier is making

Confusion matrix

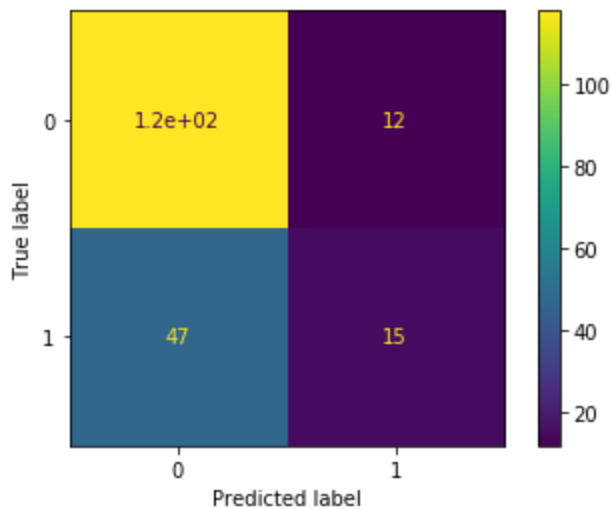
Table that describes the performance of a classification model

```
In [15]: # IMPORTANT: first argument is true values, second argument is predicted values
print(metrics.confusion_matrix(y_test, y_pred_class))
```

```
[[118  12]
 [ 47  15]]
```

```
In [21]: from sklearn.metrics import ConfusionMatrixDisplay
cm=metrics.confusion_matrix(y_test, y_pred_class)
disp = ConfusionMatrixDisplay(cm,display_labels=logreg.classes_)
disp.plot()
```

```
Out[21]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x256c5e4d048>
```



- Every observation in the testing set is represented in **exactly one box**
- It's a 2x2 matrix because there are **2 response classes**


Basic terminology

- **True Positives (TP):** we *correctly* predicted that they *do* have diabetes
- **True Negatives (TN):** we *correctly* predicted that they *don't* have diabetes
- **False Positives (FP):** we *incorrectly* predicted that they *do* have diabetes (a "Type I error")
- **False Negatives (FN):** we *incorrectly* predicted that they *don't* have diabetes (a "Type II error")

```
In [16]: # print the first 25 true and predicted responses
print('True:', y_test.values[0:25])
print('Pred:', y_pred_class[0:25])
```

```
True: [1 0 0 1 0 0 1 1 0 0 1 1 0 0 0 0 1 0 0 0 1 1 0 0 0]
Pred: [0 0 0 0 0 0 0 1 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0]
```

```
In [17]: # save confusion matrix and slice into four pieces
confusion = metrics.confusion_matrix(y_test, y_pred_class)
TP = confusion[1, 1]
TN = confusion[0, 0]
FP = confusion[0, 1]
FN = confusion[1, 0]
```

 Large confusion matrix

Metrics computed from a confusion matrix

Classification Accuracy: Overall, how often is the classifier correct?

```
In [18]: print((TP + TN) / (TP + TN + FP + FN))
print(metrics.accuracy_score(y_test, y_pred_class))

0.6927083333333334
0.6927083333333334
```

Classification Error: Overall, how often is the classifier incorrect?

- Also known as "Misclassification Rate"

```
In [19]: print((FP + FN) / (TP + TN + FP + FN))
print(1 - metrics.accuracy_score(y_test, y_pred_class))

0.3072916666666667
0.3072916666666663
```

Sensitivity: When the actual value is positive, how often is the prediction correct?

- How "sensitive" is the classifier to detecting positive instances?
- Also known as "True Positive Rate" or "Recall"

```
In [20]: print(TP / (TP + FN))
print(metrics.recall_score(y_test, y_pred_class))

0.24193548387096775
0.24193548387096775
```

Specificity: When the actual value is negative, how often is the prediction correct?

- How "specific" (or "selective") is the classifier in predicting positive instances?

```
In [21]: print(TN / (TN + FP))
```

```
0.9076923076923077
```

False Positive Rate: When the actual value is negative, how often is the prediction incorrect?

```
In [22]: print(FP / (TN + FP))
```

```
0.09230769230769231
```

Precision: When a positive value is predicted, how often is the prediction correct?

- How "precise" is the classifier when predicting positive instances?

```
In [23]: print(TP / (TP + FP))  
print(metrics.precision_score(y_test, y_pred_class))
```

```
0.5555555555555556
```

```
0.5555555555555556
```

Many other metrics can be computed: F1 score, Matthews correlation coefficient, etc.

Conclusion:

- Confusion matrix gives you a **more complete picture** of how your classifier is performing
- Also allows you to compute various **classification metrics**, and these metrics can guide your model selection

Which metrics should you focus on?

- Choice of metric depends on your **business objective**
- **Spam filter** (positive class is "spam"): Optimize for **precision or specificity** because false negatives (spam goes to the inbox) are more acceptable than false positives (non-spam is caught by the spam filter)
- **Fraudulent transaction detector** (positive class is "fraud"): Optimize for **sensitivity** because false positives (normal transactions that are flagged as possible fraud) are more acceptable than false negatives (fraudulent transactions that are not detected)

Adjusting the classification threshold

```
In [24]: # print the first 10 predicted responses  
logreg.predict(X_test)[0:10]
```

```
Out[24]: array([0, 0, 0, 0, 0, 0, 0, 1, 0, 1])
```

```
In [25]: # print the first 10 predicted probabilities of class membership
logreg.predict_proba(X_test)[0:10, :]
```

```
Out[25]: array([[0.63247571, 0.36752429],
 [0.71643656, 0.28356344],
 [0.71104114, 0.28895886],
 [0.5858938 , 0.4141062 ],
 [0.84103973, 0.15896027],
 [0.82934844, 0.17065156],
 [0.50110974, 0.49889026],
 [0.48658459, 0.51341541],
 [0.72321388, 0.27678612],
 [0.32810562, 0.67189438]])
```

```
In [26]: # print the first 10 predicted probabilities for class 1
logreg.predict_proba(X_test)[0:10, 1]
```

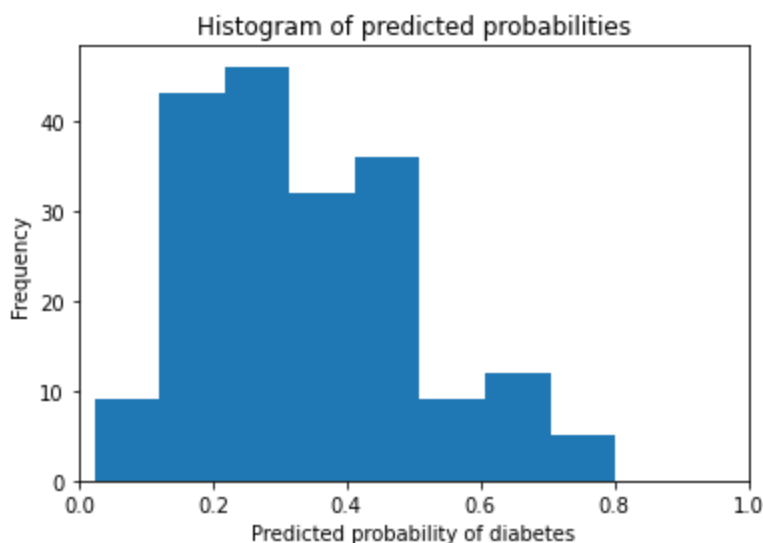
```
Out[26]: array([0.36752429, 0.28356344, 0.28895886, 0.4141062 , 0.15896027,
 0.17065156, 0.49889026, 0.51341541, 0.27678612, 0.67189438])
```

```
In [27]: # store the predicted probabilities for class 1
y_pred_prob = logreg.predict_proba(X_test)[: , 1]
```

```
In [28]: # allow plots to appear in the notebook
%matplotlib inline
import matplotlib.pyplot as plt
```

```
In [29]: # histogram of predicted probabilities
plt.hist(y_pred_prob, bins=8)
plt.xlim(0, 1)
plt.title('Histogram of predicted probabilities')
plt.xlabel('Predicted probability of diabetes')
plt.ylabel('Frequency')
```

```
Out[29]: Text(0, 0.5, 'Frequency')
```



Decrease the threshold for predicting diabetes in order to **increase the sensitivity** of the classifier

```
In [30]: # predict diabetes if the predicted probability is greater than 0.3
from sklearn.preprocessing import binarize
y_pred_class = binarize([y_pred_prob], threshold=0.3)[0]
```

```
In [31]: # print the first 10 predicted probabilities
y_pred_prob[0:10]
```

```
Out[31]: array([0.36752429, 0.28356344, 0.28895886, 0.4141062 , 0.15896027,
               0.17065156, 0.49889026, 0.51341541, 0.27678612, 0.67189438])
```

```
In [32]: # print the first 10 predicted classes with the lower threshold
y_pred_class[0:10]
```

```
Out[32]: array([1., 0., 0., 1., 0., 0., 1., 1., 0., 1.])
```

```
In [33]: # previous confusion matrix (default threshold of 0.5)
print(confusion)
```

```
[[118  12]
 [ 47  15]]
```

```
In [34]: # new confusion matrix (threshold of 0.3)
print(metrics.confusion_matrix(y_test, y_pred_class))
```

```
[[80 50]
 [16 46]]
```

```
In [35]: # sensitivity has increased (used to be 0.24)
print(46 / (46 + 16))
```

```
0.7419354838709677
```

```
In [36]: # specificity has decreased (used to be 0.91)
print(80 / (80 + 50))
```

```
0.6153846153846154
```

Conclusion:

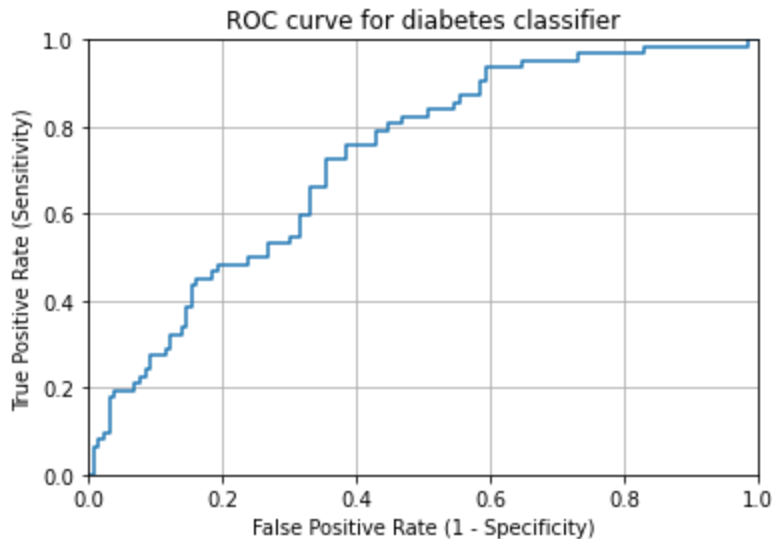
- **Threshold of 0.5** is used by default (for binary problems) to convert predicted probabilities into class predictions
- Threshold can be **adjusted** to increase sensitivity or specificity
- Sensitivity and specificity have an **inverse relationship**

ROC Curves and Area Under the Curve (AUC)

Question: Wouldn't it be nice if we could see how sensitivity and specificity are affected by various thresholds, without actually changing the threshold?

Answer: Plot the ROC curve!


```
In [37]: # IMPORTANT: first argument is true values, second argument is predicted probabilities
fpr, tpr, thresholds = metrics.roc_curve(y_test, y_pred_prob)
plt.plot(fpr, tpr)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.title('ROC curve for diabetes classifier')
plt.xlabel('False Positive Rate (1 - Specificity)')
plt.ylabel('True Positive Rate (Sensitivity)')
plt.grid(True)
```



- ROC curve can help you to **choose a threshold** that balances sensitivity and specificity in a way that makes sense for your particular context
- You can't actually **see the thresholds** used to generate the curve on the ROC curve itself

```
In [38]: # define a function that accepts a threshold and prints sensitivity and specificity
def evaluate_threshold(threshold):
    print('Sensitivity:', tpr[thresholds > threshold][-1])
    print('Specificity:', 1 - fpr[thresholds > threshold][-1])
```

```
In [39]: evaluate_threshold(0.5)

Sensitivity: 0.24193548387096775
Specificity: 0.9076923076923077
```

```
In [40]: evaluate_threshold(0.3)

Sensitivity: 0.7258064516129032
Specificity: 0.6153846153846154
```

AUC is the **percentage** of the ROC plot that is **underneath the curve**:

```
In [41]: # IMPORTANT: first argument is true values, second argument is predicted probabilities
print(metrics.roc_auc_score(y_test, y_pred_prob))

0.7245657568238213
```

- AUC is useful as a **single number summary** of classifier performance.
- If you randomly chose one positive and one negative observation, AUC represents the likelihood that your classifier will assign a **higher predicted probability** to the positive observation.
- AUC is useful even when there is **high class imbalance** (unlike classification accuracy).

```
In [42]: # calculate cross-validated AUC
from sklearn.model_selection import cross_val_score
cross_val_score(logreg, X, y, cv=10, scoring='roc_auc').mean()
```

```
Out[42]: 0.7378233618233618
```

Confusion matrix advantages:

- Allows you to calculate a **variety of metrics**
- Useful for **multi-class problems** (more than two response classes)

ROC/AUC advantages:

- Does not require you to **set a classification threshold**
- Still useful when there is **high class imbalance**

Building a Machine Learning workflow

- Why should you use a Pipeline?
- How do you encode categorical features with OneHotEncoder?
- How do you apply OneHotEncoder to selected columns with ColumnTransformer?
- How do you build and cross-validate a Pipeline?
- How do you make predictions on new data using a Pipeline?
- Why should you use scikit-learn (rather than pandas) for preprocessing?

Step 1: Load the dataset

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('http://bit.ly/kaggletrain')
```

```
In [3]: df.shape
```

```
Out[3]: (891, 12)
```

Step 2: Select features

```
In [4]: df.columns
```

```
Out[4]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',  
              'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],  
              dtype='object')
```

```
In [5]: df.isna().sum()
```

```
Out[5]: PassengerId    0  
Survived             0  
Pclass              0  
Name                0  
Sex                 0  
Age                177  
SibSp               0  
Parch              0  
Ticket             0  
Fare               0  
Cabin             687  
Embarked           2  
dtype: int64
```

```
In [6]: df = df.loc[df.Embarked.notna(), ['Survived', 'Pclass', 'Sex', 'Embarked']]
```

```
In [7]: df.shape
```

```
Out[7]: (889, 4)
```

```
In [8]: df.isna().sum()
```

```
Out[8]: Survived    0  
Pclass      0  
Sex         0  
Embarked    0  
dtype: int64
```

```
In [9]: df.head()
```

```
Out[9]:
```

	Survived	Pclass	Sex	Embarked
0	0	3	male	S
1	1	1	female	C
2	1	3	female	S
3	1	1	female	S
4	0	3	male	S

Step 3: Cross-validate a model with one feature

```
In [10]: X = df.loc[:, ['Pclass']]  
y = df.Survived
```

```
In [11]: X.shape
```

```
Out[11]: (889, 1)
```

```
In [12]: y.shape
```

```
Out[12]: (889,)
```

```
In [13]: from sklearn.linear_model import LogisticRegression
```

```
In [14]: logreg = LogisticRegression()
```

```
In [15]: from sklearn.model_selection import cross_val_score
```

```
In [16]: cross_val_score(logreg, X, y, cv=5, scoring='accuracy').mean()
```

```
Out[16]: 0.6783406335301212
```

```
In [17]: y.value_counts(normalize=True)
```

```
Out[17]: 0    0.617548  
1    0.382452  
Name: Survived, dtype: float64
```

Step 4: Encode categorical features

```
In [18]: df.head()
```

```
Out[18]:
```

	Survived	Pclass	Sex	Embarked
0	0	3	male	S
1	1	1	female	C
2	1	3	female	S
3	1	1	female	S
4	0	3	male	S

```
In [19]: # dummy encoding of categorical features  
from sklearn.preprocessing import OneHotEncoder  
ohe = OneHotEncoder(sparse=False)
```

```
In [20]: ohe.fit_transform(df[['Sex']])
```

```
Out[20]: array([[0., 1.],
                [1., 0.],
                [1., 0.],
                ...,
                [1., 0.],
                [0., 1.],
                [0., 1.]])
```

```
In [21]: ohe.categories_
```

```
Out[21]: [array(['female', 'male'], dtype=object)]
```

```
In [22]: ohe.fit_transform(df[['Embarked']])
```

```
Out[22]: array([[0., 0., 1.],
                [1., 0., 0.],
                [0., 0., 1.],
                ...,
                [0., 0., 1.],
                [1., 0., 0.],
                [0., 1., 0.]])
```

```
In [23]: ohe.categories_
```

```
Out[23]: [array(['C', 'Q', 'S'], dtype=object)]
```

Step 5: Cross-validate a Pipeline with all features

```
In [24]: X = df.drop('Survived', axis='columns')
```

```
In [25]: X.head()
```

```
Out[25]:
```

	Pclass	Sex	Embarked
0	3	male	S
1	1	female	C
2	3	female	S
3	1	female	S
4	3	male	S

```
In [26]: # use when different features need different preprocessing
from sklearn.compose import make_column_transformer
```

```
In [27]: column_trans = make_column_transformer(
    (OneHotEncoder(), ['Sex', 'Embarked']),
    remainder='passthrough')
```

```
In [28]: column_trans.fit_transform(X)
```

```
Out[28]: array([[0., 1., 0., 0., 1., 3.],
                [1., 0., 1., 0., 0., 1.],
                [1., 0., 0., 0., 1., 3.],
                ...,
                [1., 0., 0., 0., 1., 3.],
                [0., 1., 1., 0., 0., 1.],
                [0., 1., 0., 1., 0., 3.]])
```

```
In [29]: # chain sequential steps together
from sklearn.pipeline import make_pipeline
```

```
In [30]: pipe = make_pipeline(column_trans, logreg)
```

```
In [31]: # cross-validate the entire process
# thus, preprocessing occurs within each fold of cross-validation
cross_val_score(pipe, X, y, cv=5, scoring='accuracy').mean()
```

```
Out[31]: 0.7727924839713071
```

Step 6: Make predictions on "new" data

```
In [32]: # added empty cell so that the cell numbering matches the video
```

```
In [33]: X_new = X.sample(5, random_state=99)
X_new
```

```
Out[33]:
```

	Pclass	Sex	Embarked
599	1	male	C
512	1	male	S
273	1	male	C
215	1	female	C
790	3	male	Q

```
In [34]: pipe.fit(X, y)
```

```
Out[34]: Pipeline(steps=[('columntransformer',
                           ColumnTransformer(remainder='passthrough',
                                                transformers=[('onehotencoder',
                                                                OneHotEncoder(),
                                                                ['Sex', 'Embarked'])])),
                           ('logisticregression', LogisticRegression())])
```

```
In [35]: pipe.predict(X_new)
```

```
Out[35]: array([1, 0, 1, 1, 0])
```

Recap

```
In [36]: import pandas as pd
from sklearn.compose import make_column_transformer
from sklearn.preprocessing import OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.model_selection import cross_val_score
```

```
In [37]: df = pd.read_csv('http://bit.ly/kaggletrain')
df = df.loc[df.Embarked.notna(), ['Survived', 'Pclass', 'Sex', 'Embarked']]
X = df.drop('Survived', axis='columns')
y = df.Survived
```

```
In [38]: column_trans = make_column_transformer(
    (OneHotEncoder(), ['Sex', 'Embarked']),
    remainder='passthrough')
logreg = LogisticRegression(solver='lbfgs')
```

```
In [39]: pipe = make_pipeline(column_trans, logreg)
```

```
In [40]: cross_val_score(pipe, X, y, cv=5, scoring='accuracy').mean()
```

```
Out[40]: 0.7727924839713071
```

```
In [41]: X_new = X.sample(5, random_state=99)
```

```
In [42]: pipe.fit(X, y)
pipe.predict(X_new)
```

```
Out[42]: array([1, 0, 1, 1, 0])
```