

A Gentle Introduction to Threshold-Moving for Imbalanced Classification



machinelearningmastery.com/threshold-moving-for-imbalanced-classification

Classification predictive modeling typically involves predicting a class label.

Nevertheless, many machine learning algorithms are capable of predicting a probability or scoring of class membership, and this must be interpreted before it can be mapped to a crisp class label. This is achieved by using a threshold, such as 0.5, where all values equal or greater than the threshold are mapped to one class and all other values are mapped to another class.

For those classification problems that have a severe class imbalance, the default threshold can result in poor performance. As such, a simple and straightforward approach to improving the performance of a classifier that predicts probabilities on an imbalanced classification problem is to tune the threshold used to map probabilities to class labels.

In some cases, such as when using ROC Curves and Precision-Recall Curves, the best or optimal threshold for the classifier can be calculated directly. In other cases, it is possible to use a grid search to tune the threshold and locate the optimal value.

In this tutorial, you will discover how to tune the optimal threshold when converting probabilities to crisp class labels for imbalanced classification.

After completing this tutorial, you will know:

- The default threshold for interpreting probabilities to class labels is 0.5, and tuning this hyperparameter is called threshold moving.
- How to calculate the optimal threshold for the ROC Curve and Precision-Recall Curve directly.
- How to manually search threshold values for a chosen model and model evaluation metric.

Kick-start your project with my new book [Imbalanced Classification with Python](#), including *step-by-step tutorials* and the *Python source code* files for all examples.

Let's get started.

- **Update Feb/2020:** Fixed typo in Specificity equation.
- **Update Jan/2021:** Updated links for API documentation.



A Gentle Introduction to Threshold-Moving for Imbalanced Classification
Photo by [Bruna cs](#), some rights reserved.

Tutorial Overview

This tutorial is divided into five parts; they are:

1. Converting Probabilities to Class Labels
2. Threshold-Moving for Imbalanced Classification
3. Optimal Threshold for ROC Curve
4. Optimal Threshold for Precision-Recall Curve
5. Optimal Threshold Tuning

Converting Probabilities to Class Labels

Many machine learning algorithms are capable of predicting a probability or a scoring of class membership.

This is useful generally as it provides a measure of the certainty or uncertainty of a prediction. It also provides additional granularity over just predicting the class label that can be interpreted.

Some classification tasks require a crisp class label prediction. This means that even though a probability or scoring of class membership is predicted, it must be converted into a crisp class label.

The decision for converting a predicted probability or scoring into a class label is governed by a parameter referred to as the “*decision threshold*,” “*discrimination threshold*,” or simply the “*threshold*.” The default value for the threshold is 0.5 for normalized predicted probabilities or scores in the range between 0 or 1.

For example, on a binary classification problem with class labels 0 and 1, normalized predicted probabilities and a threshold of 0.5, then values less than the threshold of 0.5 are assigned to class 0 and values greater than or equal to 0.5 are assigned to class 1.

- Prediction < 0.5 = Class 0
- Prediction >= 0.5 = Class 1

The problem is that the default threshold may not represent an optimal interpretation of the predicted probabilities.

This might be the case for a number of reasons, such as:

- The predicted probabilities are not calibrated, e.g. those predicted by an SVM or decision tree.
- The metric used to train the model is different from the metric used to evaluate a final model.
- The class distribution is severely skewed.
- The cost of one type of misclassification is more important than another type of misclassification.

Worse still, some or all of these reasons may occur at the same time, such as the use of a neural network model with uncalibrated predicted probabilities on an imbalanced classification problem.

As such, there is often the need to change the default decision threshold when interpreting the predictions of a model.

... almost all classifiers generate positive or negative predictions by applying a threshold to a score. The choice of this threshold will have an impact in the trade-offs of positive and negative errors.

— Page 53, Learning from Imbalanced Data Sets, 2018.

Want to Get Started With Imbalance Classification?

Take my free 7-day email crash course now (with sample code).

Click to sign-up and also get a free PDF Ebook version of the course.

Threshold-Moving for Imbalanced Classification

There are many techniques that may be used to address an imbalanced classification problem, such as resampling the training dataset and developing customized version of machine learning algorithms.

Nevertheless, perhaps the simplest approach to handle a severe class imbalance is to change the decision threshold. Although simple and very effective, this technique is often overlooked by practitioners and research academics alike as was noted by Foster Provost in his 2000 article titled “[Machine Learning from Imbalanced Data Sets](#).”

The bottom line is that when studying problems with imbalanced data, using the classifiers produced by standard machine learning algorithms without adjusting the output threshold may well be a critical mistake.

— [Machine Learning from Imbalanced Data Sets 101](#), 2000.

There are many reasons to choose an alternative to the default decision threshold.

For example, you may use [ROC curves](#) to analyze the predicted probabilities of a model and ROC AUC scores to compare and select a model, although you require crisp class labels from your model. How do you choose the threshold on the ROC Curve that results in the best balance between the true positive rate and the false positive rate?

Alternately, you may use precision-recall curves to analyze the predicted probabilities of a model, precision-recall AUC to compare and select models, and require crisp class labels as predictions. How do you choose the threshold on the Precision-Recall Curve that results in the best balance between precision and recall?

You may use a probability-based metric to train, evaluate, and compare models like log loss ([cross-entropy](#)) but require crisp class labels to be predicted. How do you choose the optimal threshold from predicted probabilities more generally?

Finally, you may have different costs associated with false positive and false negative misclassification, a so-called cost matrix, but wish to use and evaluate cost-insensitive models and later evaluate their predictions use a cost-sensitive measure. How do you choose a threshold that finds the best trade-off for predictions using the cost matrix?

Popular way of training a cost-sensitive classifier without a known cost matrix is to put emphasis on modifying the classification outputs when predictions are being made on new data. This is usually done by setting a threshold on the positive class, below which the negative one is being predicted. The value of this threshold is optimized using a validation set and thus the cost matrix can be learned from training data.

— Page 67, [Learning from Imbalanced Data Sets](#), 2018.

The answer to these questions is to search a range of threshold values in order to find the best threshold. In some cases, the optimal threshold can be calculated directly.

Tuning or shifting the decision threshold in order to accommodate the broader requirements of the classification problem is generally referred to as “*threshold-moving*,” “*threshold-tuning*,” or simply “*thresholding*.”

It has been stated that trying other methods, such as sampling, without trying by simply setting the threshold may be misleading. The threshold-moving method uses the original training set to train [a model] and then moves the decision threshold such that the minority class examples are easier to be predicted correctly.

— Pages 72, Imbalanced Learning: Foundations, Algorithms, and Applications, 2013.

The process involves first fitting the model on a training dataset and making predictions on a test dataset. The predictions are in the form of normalized probabilities or scores that are transformed into normalized probabilities. Different threshold values are then tried and the resulting crisp labels are evaluated using a chosen evaluation metric. The threshold that achieves the best evaluation metric is then adopted for the model when making predictions on new data in the future.

We can summarize this procedure below.

- 1. Fit Model on the Training Dataset.
- 2. Predict Probabilities on the Test Dataset.
- 3. For each threshold in Thresholds:
 - 3a. Convert probabilities to Class Labels using the threshold.
 - 3b. Evaluate Class Labels.
 - 3c. If Score is Better than Best Score.
 - 3ci. Adopt Threshold.
 - 4. Use Adopted Threshold When Making Class Predictions on New Data.

Although simple, there are a few different approaches to implementing threshold-moving depending on your circumstance. We will take a look at some of the most common examples in the following sections.

Optimal Threshold for ROC Curve

A ROC curve is a diagnostic plot that evaluates a set of probability predictions made by a model on a test dataset.

A set of different thresholds are used to interpret the true positive rate and the false positive rate of the predictions on the positive (minority) class, and the scores are plotted in a line of increasing thresholds to create a curve.

The false-positive rate is plotted on the x-axis and the true positive rate is plotted on the y-axis and the plot is referred to as the Receiver Operating Characteristic curve, or ROC curve. A diagonal line on the plot from the bottom-left to top-right indicates the “*curve*” for a no-skill classifier (predicts the majority class in all cases), and a point in the top left of the plot indicates a model with perfect skill.

The curve is useful to understand the trade-off in the true-positive rate and false-positive rate for different thresholds. The area under the ROC Curve, so-called ROC AUC, provides a single number to summarize the performance of a model in terms of its ROC Curve with a value between 0.5 (no-skill) and 1.0 (perfect skill).

The ROC Curve is a useful diagnostic tool for understanding the trade-off for different thresholds and the ROC AUC provides a useful number for comparing models based on their general capabilities.

If crisp class labels are required from a model under such an analysis, then an optimal threshold is required. This would be a threshold on the curve that is closest to the top-left of the plot.

Thankfully, there are principled ways of locating this point.

First, let's fit a model and calculate a ROC Curve.

We can use the make_classification() function to create a synthetic binary classification problem with 10,000 examples (rows), 99 percent of which belong to the majority class and 1 percent belong to the minority class.

```
1 ...
2 # generate dataset
3 X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
4 n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
```

We can then split the dataset using the train_test_split() function and use half for the training set and half for the test set.

```
1 ...
2 # split into train/test sets
3 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
4 stratify=y)
```

We can then fit a LogisticRegression model and use it to make probability predictions on the test set and keep only the probability predictions for the minority class.

```
1 ...
2 # fit a model
3 model = LogisticRegression(solver='lbfgs')
4 model.fit(trainX, trainy)
5 # predict probabilities
6 lr_probs = model.predict_proba(testX)
7 # keep probabilities for the positive outcome only
8 lr_probs = lr_probs[:, 1]
```

We can then use the roc_auc_score() function to calculate the true-positive rate and false-positive rate for the predictions using a set of thresholds that can then be used to create a ROC Curve plot.

```
1 ...
2 # calculate scores
3 lr_auc = roc_auc_score(testy, lr_probs)
```

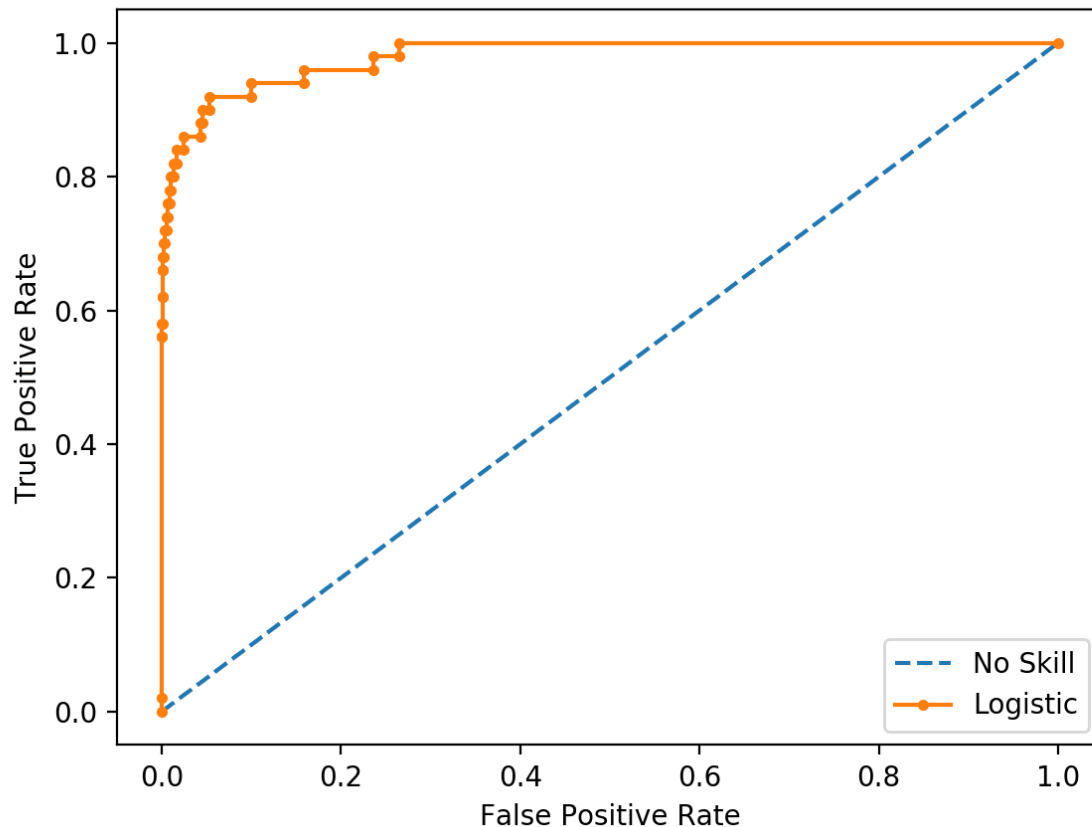
We can tie this all together, defining the dataset, fitting the model, and creating the ROC Curve plot. The complete example is listed below.

```
1  # roc curve for logistic regression model
2  from sklearn.datasets import make_classification
3  from sklearn.linear_model import LogisticRegression
4  from sklearn.model_selection import train_test_split
5  from sklearn.metrics import roc_curve
6  from matplotlib import pyplot
7  # generate dataset
8  X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
9  n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
10 # split into train/test sets
11 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
12 stratify=y)
13 # fit a model
14 model = LogisticRegression(solver='lbfgs')
15 model.fit(trainX, trainy)
16 # predict probabilities
17 yhat = model.predict_proba(testX)
18 # keep probabilities for the positive outcome only
19 yhat = yhat[:, 1]
20 # calculate roc curves
21 fpr, tpr, thresholds = roc_curve(testy, yhat)
22 # plot the roc curve for the model
23 pyplot.plot([0, 1], [0, 1], linestyle='--', label='No Skill')
24 pyplot.plot(fpr, tpr, marker='.', label='Logistic')
25 # axis labels
26 pyplot.xlabel('False Positive Rate')
27 pyplot.ylabel('True Positive Rate')
28 pyplot.legend()
29 # show the plot
    pyplot.show()
```

Running the example fits a logistic regression model on the training dataset then evaluates it using a range of thresholds on the test set, creating the ROC Curve

We can see that there are a number of points or thresholds close to the top-left of the plot.

Which is the threshold that is optimal?



ROC Curve Line Plot for Logistic Regression Model for Imbalanced Classification

There are many ways we could locate the threshold with the optimal balance between false positive and true positive rates.

Firstly, the true positive rate is called the Sensitivity. The inverse of the false-positive rate is called the Specificity.

- Sensitivity = $\text{TruePositive} / (\text{TruePositive} + \text{FalseNegative})$
- Specificity = $\text{TrueNegative} / (\text{FalsePositive} + \text{TrueNegative})$

Where:

- Sensitivity = True Positive Rate
- Specificity = $1 - \text{False Positive Rate}$

The Geometric Mean or G-Mean is a metric for imbalanced classification that, if optimized, will seek a balance between the sensitivity and the specificity.

$$\text{G-Mean} = \sqrt{\text{Sensitivity} * \text{Specificity}}$$

One approach would be to test the model with each threshold returned from the call `roc_auc_score()` and select the threshold with the largest G-Mean value.

Given that we have already calculated the Sensitivity (TPR) and the complement to the Specificity when we calculated the ROC Curve, we can calculate the G-Mean for each threshold directly.

```
1 ...
2 # calculate the g-mean for each threshold
3 gmeans = sqrt(tpr * (1-fpr))
```

Once calculated, we can locate the index for the largest G-mean score and use that index to determine which threshold value to use.

```
1 ...
2 # locate the index of the largest g-mean
3 ix = argmax(gmeans)
4 print('Best Threshold=%f, G-Mean=%.3f' % (thresholds[ix], gmeans[ix]))
```

We can also re-draw the ROC Curve and highlight this point.

The complete example is listed below.

```

1  # roc curve for logistic regression model with optimal threshold
2  from numpy import sqrt
3  from numpy import argmax
4  from sklearn.datasets import make_classification
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import train_test_split
7  from sklearn.metrics import roc_curve
8  from matplotlib import pyplot
9  # generate dataset
10 X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
11 n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
12 # split into train/test sets
13 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
14 stratify=y)
15 # fit a model
16 model = LogisticRegression(solver='lbfgs')
17 model.fit(trainX, trainy)
18 # predict probabilities
19 yhat = model.predict_proba(testX)
20 # keep probabilities for the positive outcome only
21 yhat = yhat[:, 1]
22 # calculate roc curves
23 fpr, tpr, thresholds = roc_curve(testy, yhat)
24 # calculate the g-mean for each threshold
25 gmeans = sqrt(tpr * (1-fpr))
26 # locate the index of the largest g-mean
27 ix = argmax(gmeans)
28 print('Best Threshold=%f, G-Mean=%.3f' % (thresholds[ix], gmeans[ix]))
29 # plot the roc curve for the model
30 pyplot.plot([0,1], [0,1], linestyle='--', label='No Skill')
31 pyplot.plot(fpr, tpr, marker='.', label='Logistic')
32 pyplot.scatter(fpr[ix], tpr[ix], marker='o', color='black', label='Best')
33 # axis labels
34 pyplot.xlabel('False Positive Rate')
35 pyplot.ylabel('True Positive Rate')
36 pyplot.legend()
37 # show the plot
    pyplot.show()

```

Running the example first locates the optimal threshold and reports this threshold and the G-Mean score.

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

In this case, we can see that the optimal threshold is about 0.016153.

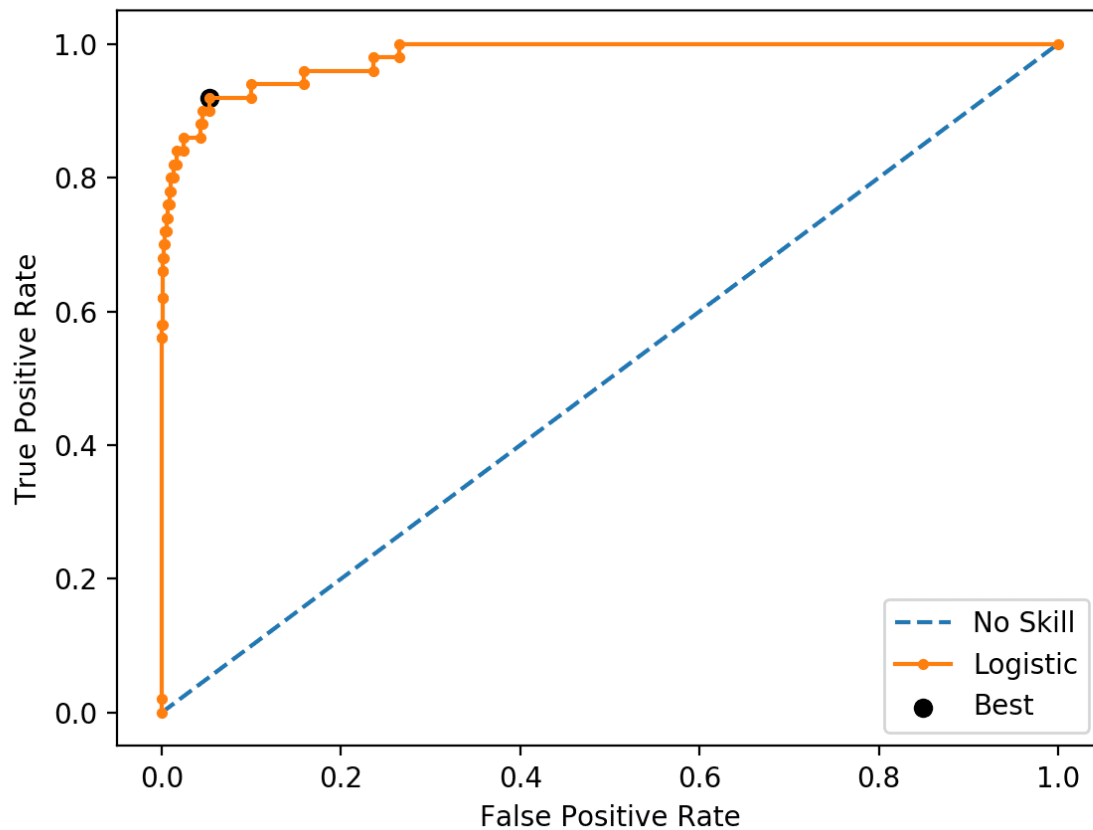
```

1  Best Threshold=0.016153, G-Mean=0.933

```

The threshold is then used to locate the true and false positive rates, then this point is drawn on the ROC Curve.

We can see that the point for the optimal threshold is a large black dot and it appears to be closest to the top-left of the plot.



ROC Curve Line Plot for Logistic Regression Model for Imbalanced Classification With the Optimal Threshold

It turns out there is a much faster way to get the same result, called the Youden's J statistic.

The statistic is calculated as:

$$J = \text{Sensitivity} + \text{Specificity} - 1$$

Given that we have Sensitivity (TPR) and the complement of the specificity (FPR), we can calculate it as:

$$J = \text{Sensitivity} + (1 - \text{FalsePositiveRate}) - 1$$

Which we can restate as:

$$J = \text{TruePositiveRate} - \text{FalsePositiveRate}$$

We can then choose the threshold with the largest J statistic value. For example:

```

1  ...
2  # calculate roc curves
3  fpr, tpr, thresholds = roc_curve(testy, yhat)
4  # get the best threshold
5  J = tpr - fpr
6  ix = argmax(J)
7  best_thresh = thresholds[ix]
8  print('Best Threshold=%f' % (best_thresh))

```

Plugging this in, the complete example is listed below.

```

1  # roc curve for logistic regression model with optimal threshold
2  from numpy import argmax
3  from sklearn.datasets import make_classification
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.model_selection import train_test_split
6  from sklearn.metrics import roc_curve
7  # generate dataset
8  X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
9  n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
10 # split into train/test sets
11 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
12 stratify=y)
13 # fit a model
14 model = LogisticRegression(solver='lbfgs')
15 model.fit(trainX, trainy)
16 # predict probabilities
17 yhat = model.predict_proba(testX)
18 # keep probabilities for the positive outcome only
19 yhat = yhat[:, 1]
20 # calculate roc curves
21 fpr, tpr, thresholds = roc_curve(testy, yhat)
22 # get the best threshold
23 J = tpr - fpr
24 ix = argmax(J)
25 best_thresh = thresholds[ix]
   print('Best Threshold=%f' % (best_thresh))

```

We can see that this simpler approach calculates the optimal statistic directly.

```

1  Best Threshold=0.016153

```

Optimal Threshold for Precision-Recall Curve

Unlike the ROC Curve, a precision-recall curve focuses on the performance of a classifier on the positive (minority class) only.

Precision is the ratio of the number of true positives divided by the sum of the true positives and false positives. It describes how good a model is at predicting the positive class. Recall is calculated as the ratio of the number of true positives divided by the sum of the true positives and the false negatives. Recall is the same as sensitivity.

A precision-recall curve is calculated by creating crisp class labels for probability predictions across a set of thresholds and calculating the precision and recall for each threshold. A line plot is created for the thresholds in ascending order with recall on the x-axis and precision on the y-axis.

A no-skill model is represented by a horizontal line with a precision that is the ratio of positive examples in the dataset (e.g. $TP / (TP + TN)$), or 0.01 on our synthetic dataset. A perfect skill classifier has full precision and recall with a dot in the top-right corner.

We can use the same model and dataset from the previous section and evaluate the probability predictions for a logistic regression model using a precision-recall curve. The `precision_recall_curve()` function can be used to calculate the curve, returning the precision and recall scores for each threshold as well as the thresholds used.

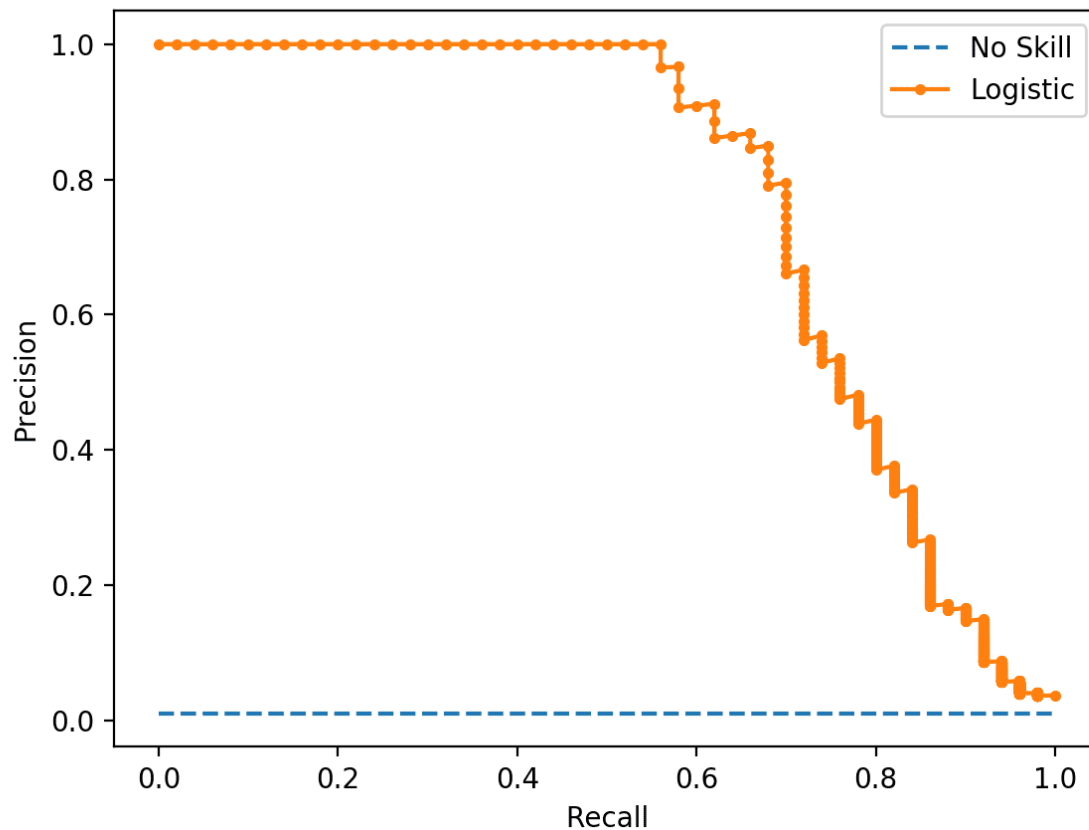
```
1 ...
2 # calculate pr-curve
3 precision, recall, thresholds = precision_recall_curve(testy, yhat)
```

Tying this together, the complete example of calculating a precision-recall curve for a logistic regression on an imbalanced classification problem is listed below.

```
1 # pr curve for logistic regression model
2 from sklearn.datasets import make_classification
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import precision_recall_curve
6 from matplotlib import pyplot
7 # generate dataset
8 X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
9 n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
10 # split into train/test sets
11 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
12 stratify=y)
13 # fit a model
14 model = LogisticRegression(solver='lbfgs')
15 model.fit(trainX, trainy)
16 # predict probabilities
17 yhat = model.predict_proba(testX)
18 # keep probabilities for the positive outcome only
19 yhat = yhat[:, 1]
20 # calculate pr-curve
21 precision, recall, thresholds = precision_recall_curve(testy, yhat)
22 # plot the roc curve for the model
23 no_skill = len(testy[testy==1]) / len(testy)
24 pyplot.plot([0, 1], [no_skill, no_skill], linestyle='--', label='No Skill')
25 pyplot.plot(recall, precision, marker='.', label='Logistic')
26 # axis labels
27 pyplot.xlabel('Recall')
28 pyplot.ylabel('Precision')
29 pyplot.legend()
30 # show the plot
31 pyplot.show()
```

Running the example calculates the precision and recall for each threshold and creates a precision-recall plot showing that the model has some skill across a range of thresholds on this dataset.

If we required crisp class labels from this model, which threshold would achieve the best result?



Precision-Recall Curve Line Plot for Logistic Regression Model for Imbalanced Classification

If we are interested in a threshold that results in the best balance of precision and recall, then this is the same as optimizing the F-measure that summarizes the harmonic mean of both measures.

$$\text{F-Measure} = (2 * \text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

As in the previous section, the naive approach to finding the optimal threshold would be to calculate the F-measure for each threshold. We can achieve the same effect by converting the precision and recall measures to F-measure directly; for example:

```
1 ...
2 # convert to f score
3 fscore = (2 * precision * recall) / (precision + recall)
4 # locate the index of the largest f score
5 ix = argmax(fscore)
6 print('Best Threshold=%f, F-Score=%.3f' % (thresholds[ix], fscore[ix]))
```

We can then plot the point on the precision-recall curve.

The complete example is listed below.

```
1  # optimal threshold for precision-recall curve with logistic regression model
2  from numpy import argmax
3  from sklearn.datasets import make_classification
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.model_selection import train_test_split
6  from sklearn.metrics import precision_recall_curve
7  from matplotlib import pyplot
8  # generate dataset
9  X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
10 n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
11 # split into train/test sets
12 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
13 stratify=y)
14 # fit a model
15 model = LogisticRegression(solver='lbfgs')
16 model.fit(trainX, trainy)
17 # predict probabilities
18 yhat = model.predict_proba(testX)
19 # keep probabilities for the positive outcome only
20 yhat = yhat[:, 1]
21 # calculate roc curves
22 precision, recall, thresholds = precision_recall_curve(testy, yhat)
23 # convert to f score
24 fscore = (2 * precision * recall) / (precision + recall)
25 # locate the index of the largest f score
26 ix = argmax(fscore)
27 print('Best Threshold=%f, F-Score=%.3f' % (thresholds[ix], fscore[ix]))
28 # plot the roc curve for the model
29 no_skill = len(testy[testy==1]) / len(testy)
30 pyplot.plot([0,1], [no_skill,no_skill], linestyle='--', label='No Skill')
31 pyplot.plot(recall, precision, marker='.', label='Logistic')
32 pyplot.scatter(recall[ix], precision[ix], marker='o', color='black', label='Best')
33 # axis labels
34 pyplot.xlabel('Recall')
35 pyplot.ylabel('Precision')
36 pyplot.legend()
37 # show the plot
38 pyplot.show()
```

Running the example first calculates the F-measure for each threshold, then locates the score and threshold with the largest value.

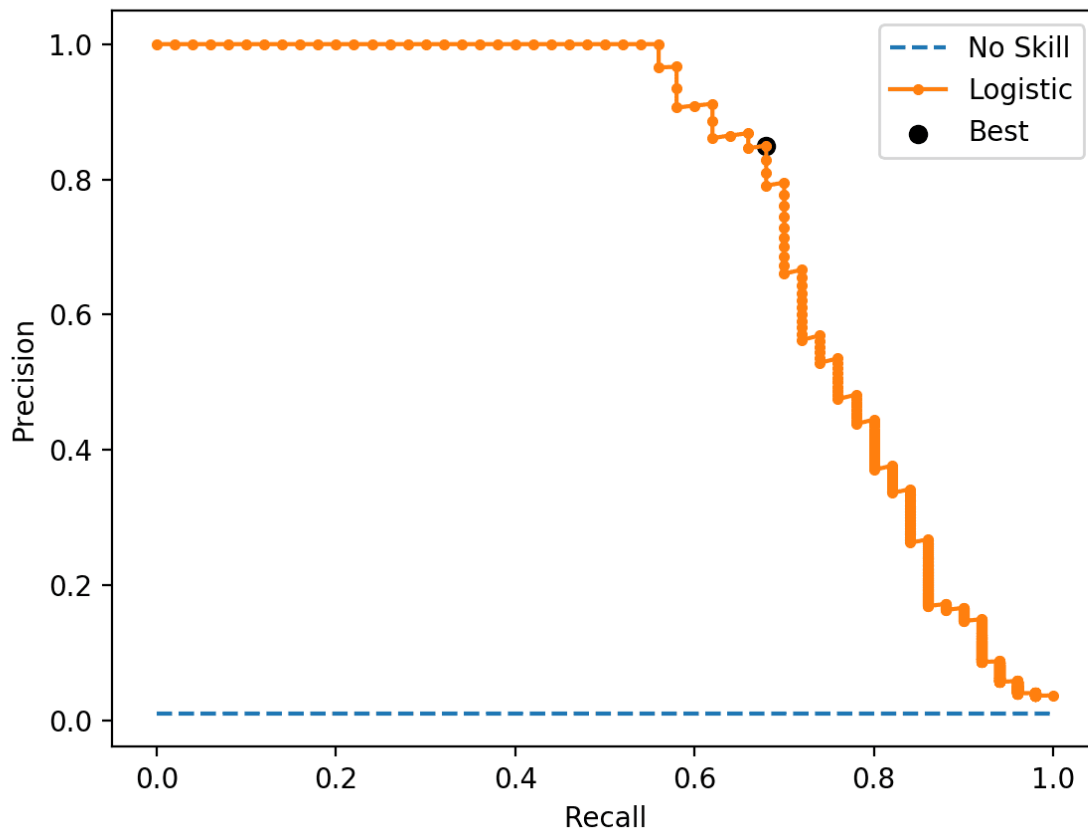
Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

In this case, we can see that the best F-measure was 0.756 achieved with a threshold of about 0.25.

1 Best Threshold=0.256036, F-Score=0.756

The precision-recall curve is plotted, and this time the threshold with the optimal F-measure is plotted with a larger black dot.

This threshold could then be used when making probability predictions in the future that must be converted from probabilities to crisp class labels.



Precision-Recall Curve Line Plot for Logistic Regression Model With Optimal Threshold

Optimal Threshold Tuning

Sometimes, we simply have a model and we wish to know the best threshold directly.

In this case, we can define a set of thresholds and then evaluate predicted probabilities under each in order to find and select the optimal threshold.

We can demonstrate this with a worked example.

First, we can fit a logistic regression model on our synthetic classification problem, then predict class labels and evaluate them using the F-Measure, which is the harmonic mean of precision and recall.

This will use the default threshold of 0.5 when interpreting the probabilities predicted by the logistic regression model.

The complete example is listed below.

```
1  # logistic regression for imbalanced classification
2  from sklearn.datasets import make_classification
3  from sklearn.linear_model import LogisticRegression
4  from sklearn.model_selection import train_test_split
5  from sklearn.metrics import f1_score
6  # generate dataset
7  X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
8  n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
9  # split into train/test sets
10 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
11 stratify=y)
12 # fit a model
13 model = LogisticRegression(solver='lbfgs')
14 model.fit(trainX, trainy)
15 # predict labels
16 yhat = model.predict(testX)
17 # evaluate the model
18 score = f1_score(testy, yhat)
    print('F-Score: %.5f' % score)
```

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

Running the example, we can see that the model achieved an F-Measure of about 0.70 on the test dataset.

```
1  F-Score: 0.70130
```

Now we can use the same model on the same dataset and instead of predicting class labels directly, we can predict probabilities.

```
1  ...
2  # predict probabilities
3  yhat = model.predict_proba(testX)
```

We only require the probabilities for the positive class.

```
1  ...
2  # keep probabilities for the positive outcome only
3  probs = yhat[:, 1]
```

Next, we can then define a set of thresholds to evaluate the probabilities. In this case, we will test all thresholds between 0.0 and 1.0 with a step size of 0.001, that is, we will test 0.0, 0.001, 0.002, 0.003, and so on to 0.999.

```

1  ...
2  # define thresholds
3  thresholds = arange(0, 1, 0.001)

```

Next, we need a way of using a single threshold to interpret the predicted probabilities.

This can be achieved by mapping all values equal to or greater than the threshold to 1 and all values less than the threshold to 0. We will define a *to_labels()* function to do this that will take the probabilities and threshold as an argument and return an array of integers in {0, 1}.

```

1  # apply threshold to positive probabilities to create labels
2  def to_labels(pos_probs, threshold):
3  return (pos_probs >= threshold).astype('int')

```

We can then call this function for each threshold and evaluate the resulting labels using the *f1_score()*.

We can do this in a single line, as follows:

```

1  ...
2  # evaluate each threshold
3  scores = [f1_score(testy, to_labels(probs, t)) for t in thresholds]

```

We now have an array of scores that evaluate each threshold in our array of thresholds.

All we need to do now is locate the array index that has the largest score (best F-Measure) and we will have the optimal threshold and its evaluation.

```

1  ...
2  # get best threshold
3  ix = argmax(scores)
4  print('Threshold=%.3f, F-Score=%.5f' % (thresholds[ix], scores[ix]))

```

Tying this all together, the complete example of tuning the threshold for the logistic regression model on the synthetic imbalanced classification dataset is listed below.

```

1  # search thresholds for imbalanced classification
2  from numpy import arange
3  from numpy import argmax
4  from sklearn.datasets import make_classification
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import train_test_split
7  from sklearn.metrics import f1_score
8
9  # apply threshold to positive probabilities to create labels
10 def to_labels(pos_probs, threshold):
11     return (pos_probs >= threshold).astype('int')
12
13 # generate dataset
14 X, y = make_classification(n_samples=10000, n_features=2, n_redundant=0,
15 n_clusters_per_class=1, weights=[0.99], flip_y=0, random_state=4)
16 # split into train/test sets
17 trainX, testX, trainy, testy = train_test_split(X, y, test_size=0.5, random_state=2,
18 stratify=y)
19 # fit a model
20 model = LogisticRegression(solver='lbfgs')
21 model.fit(trainX, trainy)
22 # predict probabilities
23 yhat = model.predict_proba(testX)
24 # keep probabilities for the positive outcome only
25 probs = yhat[:, 1]
26 # define thresholds
27 thresholds = arange(0, 1, 0.001)
28 # evaluate each threshold
29 scores = [f1_score(testy, to_labels(probs, t)) for t in thresholds]
30 # get best threshold
31 ix = argmax(scores)
    print('Threshold=%.3f, F-Score=%.5f' % (thresholds[ix], scores[ix]))

```

Running the example reports the optimal threshold as 0.251 (compared to the default of 0.5) that achieves an F-Measure of about 0.75 (compared to 0.70).

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

You can use this example as a template when tuning the threshold on your own problem, allowing you to substitute your own model, metric, and even resolution of thresholds that you want to evaluate.

```

1  Threshold=0.251, F-Score=0.75556

```