









DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

JOURNALISM

**EVENTS** 

SUBMIT

#### **Towards Data** Science

Sharing concepts, ideas, and codes.

Following

# Dealing with Imbalanced Data A guide to effectively handling imbalanced datasets in Python



Tara Boyle Follow Feb 4 · 5 min read







Upgrade





DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

JOURNALISM

492

Fraud

**EVENTS** 

SUBMIT

#### Science

Sharing concepts, ideas, and codes.

Following



Photo by Ales Nesetril on Unsplash

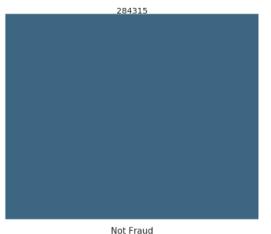
Imbalanced classes are a common problem in machine learning classification where there are a disproportionate ratio of observations in each class. Class imbalance can be found in many different areas including medical diagnosis, spam filtering, and fraud detection. In this guide, we'll look at five possible ways to handle an imbalanced class problem. Important Note: This guide will focus solely on addressing imbalanced classes and will not addressing other important machine learning steps including, but not limited to, feature selection or hyperparameter tuning.

#### Data

We will use the Credit Card Fraud Detection Dataset available on Kaggle. The dataset is high imbalanced, with only 0.17% of transactions being classified as fraudulent. The full notebook can be found here.

Our objective will be to correctly classify the minority class of fraudulent transactions.

### Distribution of Transactions



The Problem with Imbalanced Classes

Most machine learning algorithms work best when the number of samples in each class are about equal. This is because most algorithms are designed to maximize accuracy and reduce error. The Problem with  ${\bf Accuracy}$ 

VISUALIZATION





**EVENTS** 

JOURNALISM





SUBMIT

Towards Data

Science

DATA SCIENCE

Sharing concepts, ideas, and codes.

Following





```
# setting up testing and training sets
     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=27)
     \# DummyClassifier to predict only target 0
     dummy = DummyClassifier(strategy='most_frequent').fit(X_train, y_train)
10
     dummy pred = dummy.predict(X test)
     # checking unique labels
     print('Unique predicted labels: ', (np.unique(dummy_pred)))
14
     # checking accuracy
     print('Test score: ', accuracy_score(y_test, dummy_pred))
17
     Unique predicted labels: [0]
     Test score: 0.9981461194910255
dummy_classifier.py hosted with ♥ by GitHub
                                                                                             view raw
```

PROGRAMMING

We got an accuracy score of 99.8% — And without even training a model! Let's compare this to

logistic regression, an actual trained classifier.

MACHINE LEARNING

```
# Modeling the data as is
     # Train model
     lr = LogisticRegression(solver='liblinear').fit(X train, y train)
     # Predict on training set
     lr_pred = lr.predict(X_test)
     # Checking accuracy
     accuracy_score(y_test, lr_pred)
         0.9992
10
     # Checking unique values
     predictions = pd.DataFrame(lr_pred)
     predictions[0].value_counts()
              71108
         1
                 94
Ir_baseline.py hosted with ♥ by GitHub
                                                                                              view raw
```









DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

AI JOURNALISM

EVENTS

SUBMIT

# Towards Data Science

Sharing concepts, ideas, and codes.

Following



327



can be very misleading. Metrics that can provide better insight include:

- Confusion Matrix: a table showing correct predictions and types of incorrect predictions.
- **Precision:** the number of true positives divided by all positive predictions. Precision is also called Positive Predictive Value. It is a measure of a classifier's exactness. Low precision indicates a high number of false positives.
- Recall: the number of true positives divided by the number of positive values in the test data.
   Recall is also called Sensitivity or the True Positive Rate. It is a measure of a classifier's completeness. Low recall indicates a high number of false negatives.
- **F1: Score:** the weighted average of precision and recall.

Let's see what happens when we apply these F1 and recall scores to our logistic regression from

```
above.

1  # f1 score

2  f1_score(y_test, lr_pred)

3   0.7522

4

5  # recall score

6  recall_score(y_test, lr_pred)

7   0.6439

change_metric.py hosted with ♥ by GitHub
```

These scores don't look quite so impressive. Let's see what other methods we might try to improve our new metrics.

#### 2. Change the algorithm

While in every machine learning problem, it's a good rule of thumb to try a variety of algorithms, it can be especially beneficial with imbalanced datasets. Decision trees frequently perform well on imbalanced data. They work by learning a hierarchy of if/else questions and this can force both









Deta Sale

DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

AI

JOURNALISM

**EVENTS** 

SUBMIT

#### Towards Data Science

Sharing concepts, ideas, and codes.

Following



327



change\_algorithm.py hosted with ♥ by GitHub

view raw

While our accuracy score is slightly lower, both F1 and recall have increased as compared to logistic regression! It appears that for this specific problem, random forest may be a better choice of model.

#### 3. Resampling Techniques — Oversample minority class

Our next method begins our resampling techniques.

Oversampling can be defined as adding more copies of the minority class. Oversampling can be a good choice when you don't have a ton of data to work with.

We will use the resampling module from Scikit-Learn to randomly replicate samples from the minority class.

#### **Important Note**

Always split into test and train sets BEFORE trying oversampling techniques! Oversampling before splitting the data can allow the exact same observations to be present in both the test and train sets.

This can allow our model to simply memorize specific data points and cause overfitting and poor

```
generalization to the test data.
      from sklearn.utils import resample
      # Separate input features and target
      y = df.Class
      X = df.drop('Class', axis=1)
      # setting up testing and training sets
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=27)
       # concatenate our training data back together
      X = pd.concat([X_train, y_train], axis=1)
      # separate minority and majority classes
      not_fraud = X[X.Class==0]
      fraud = X[X.Class==1]
      # upsample minority
 18
      fraud_upsampled = resample(fraud,
                                 replace=True, # sample with replacement
 20
                                 n_samples=len(not_fraud), # match number in majority class
                                 random_state=27) # reproducible results
      # combine majority and upsampled minority
```









Deste Sei

DATA SCIENCE

MACHINE LEARNING

410440

PROGRAMMING

VISUALIZATION

AI JOURNALISM

EVENTS

SUBMIT

# Towards Data Science

Sharing concepts, ideas, and codes.

Following



327



upsampled.py hosted with ♥ by GitHub view raw

After resampling we have an equal ratio of data points for each class! Let's try our logistic regression

```
again with the balanced training data.
      # trying logistic regression again with the balanced dataset
      y_train = upsampled.Class
      X_train = upsampled.drop('Class', axis=1)
       upsampled = LogisticRegression(solver='liblinear').fit(X_train, y_train)
      upsampled pred = upsampled.predict(X test)
      # Checking accuracy
      accuracy score(y test, upsampled pred)
          0.9807
      # f1 score
      f1 score(y test, upsampled pred)
          0.1437
      recall_score(y_test, upsampled_pred)
 18
          0.8712
 upsampled_Ir.py hosted with ♥ by GitHub
                                                                                               view raw
```

Our recall score increased, but F1 is much lower than with either our baseline logistic regression or random forest from above. Let's see if undersampling might perform better here.

#### 4. Resampling techniques — Undersample majority class

Undersampling can be defined as removing some observations of the majority class. Undersampling can be a good choice when you have a ton of data -think millions of rows. But a drawback is that we are removing information that may be valuable. This could lead to underfitting and poor generalization to the test set.

We will again use the resampling module from Scikit-Learn to randomly remove samples from the majority class.

```
# still using our separated classes fraud and not_fraud from above

# downsample majority

not_fraud_downsampled = resample(not_fraud,

replace = False, # sample without replacement
```

VISUALIZATION





**EVENTS** 

JOURNALISM





SUBMIT

Towards Data Science DATA SCIENCE

Sharing concepts, ideas, and codes.

Following



\_

```
12 # checking counts
13 downsampled.Class.value_counts()
14 1 360
15 0 360

downsampled.py hosted with ♥ by GitHub

view raw
```

Again, we have an equal ratio of fraud to not fraud data points, but in this case a much smaller

PROGRAMMING

quantity of data to train the model on. Let's again apply our logistic regression.

MACHINE LEARNING

```
# trying logistic regression again with the undersampled dataset
     y_train = downsampled.Class
     X_train = downsampled.drop('Class', axis=1)
     undersampled = LogisticRegression(solver='liblinear').fit(X_train, y_train)
     undersampled_pred = undersampled.predict(X_test)
     # Checking accuracy
     accuracy_score(y_test, undersampled_pred)
         0.9758
     # f1 score
     f1_score(y_test, undersampled_pred)
16
         0.1171
17
     recall_score(y_test, undersampled_pred)
         0.8636
downsampled_Ir.py hosted with ♥ by GitHub
```

Undersampling underperformed oversampling in this case. Let's try one more method for handling imbalanced data.

#### 5. Generate synthetic samples

A technique similar to upsampling is to create synthetic samples. Here we will use imblearn's SMOTE or Synthetic Minority Oversampling Technique. SMOTE uses a nearest neighbors algorithm to generate new and synthetic data we can use for training our model. Again, it's important to generate the new samples only in the training set to ensure our model generalizes well to unseen data.

```
1 from imblearn.over sampling import SMOTE
```

VISUALIZATION





JOURNALISM



**EVENTS** 



SUBMIT

Towards Data Science

Sharing concepts, ideas, and codes.

Following



327



```
# setting up testing and training sets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=27)

sm = SMOTE(random_state=27, ratio=1.0)

X_train, y_train = sm.fit_sample(X_train, y_train)

smote.py hosted with \( \subseteq \) by GitHub
view raw
```

After generating our synthetic data points, let's see how our logistic regression performs.

PROGRAMMING

```
smote = LogisticRegression(solver='liblinear').fit(X_train, y_train)

smote_pred = smote.predict(X_test)

function

checking accuracy

accuracy_score(y_test, smote_pred)

o.9858

function

function

function

function

checking accuracy

accuracy_score(y_test, smote_pred)

o.9858

function

function

function

checking accuracy

accuracy_score(y_test, smote_pred)

o.9858

function

function

checking accuracy

accuracy

function

checking accuracy

function

checking accuracy

accuracy

function

checking accuracy

accuracy

function

checking accuracy

accuracy

function

checking accuracy

function

checking accuracy

accuracy

function

checking accuracy

function

function

checking accuracy

function

checking accuracy

function

functi
```

Our F1 score is increased and recall is similar to the upsampled model above and for our data here outperforms undersampling.

#### Conclusion

DATA SCIENCE

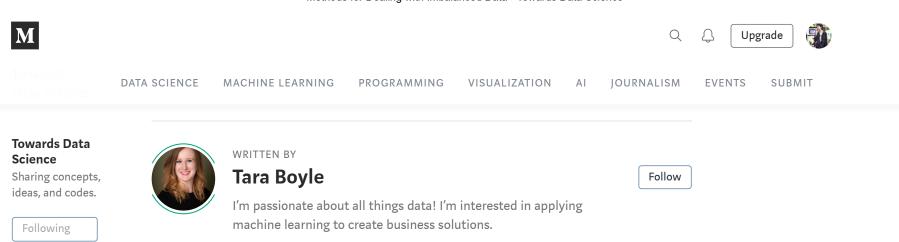
MACHINE LEARNING

We explored 5 different methods for dealing with imbalanced datasets:

- 1. Change the performance metric
- 2. Change the algorithm
- 3. Oversample minority class
- 4. Undersample majority class
- 5. Generate synthetic samples

It appears for this particular dataset random forest and SMOTE are among the best of the options we tried here.

These are just some of the many possible methods to try when dealing with imbalanced datasets, and not an exhaustive list. Some others methods to consider are collecting more data or choosing different resampling ratios — you don't have to have exactly a 1:1 ratio! You should always try several approaches and then decide which is best for your problem.









## **Towards Data Science**

Sharing concepts, ideas, and codes.

See responses (3)

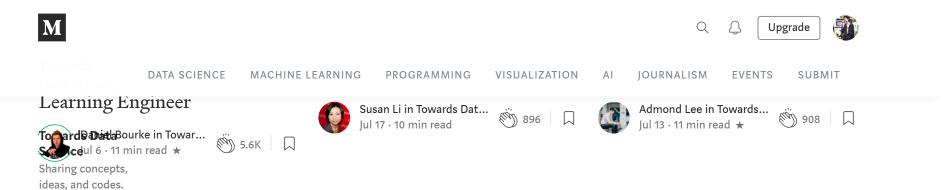
## **More From Medium**

More from Towards Data Science

More from Towards Data Science

More from Towards Data Science

Following



Following

(間) 327