

Feature Engineering



Mass Street
University
PER EDUCATIONEM PROGRESSUS

What we will cover

- Feature Engineering
- Fine Tuning
- Multiclass
- Multilabel
- Multi-output



Feature Engineering

- Anything that changes our data.
- Very broad term but important to learn how to do all facets of it.



JSON Data

- Heavily nested data.
- Uses dictionaries to hold data.
- Can call on dictionary keys and values.

```
{
  "business_id": "PK6aSizckHFWk8i0xt5DA",
  "full_address": "400 Waterfront Dr E\nHomestead\nHomestead, PA 15120",
  "hours": {},
  "open": true,
  "categories": [
    "Burgers",
    "Fast Food",
    "Restaurants"
  ],
  "city": "Homestead",
  "review_count": 5,
  "name": "McDonald's",
  "neighborhoods": [
    "Homestead"
  ],
  "longitude": -79.910032,
  "state": "PA",
  "stars": 2,
  "latitude": 40.412086,
  "attributes": {
    "Take-out": true,
    "Wi-Fi": "free",
    "Drive-Thru": true,
    "Good For": {
      "dessert": false,
      "latenight": false,
      "lunch": false,
      "dinner": false,
      "breakfast": false,
      "brunch": false
    },
    "Caters": false,
    "Noise Level": "average",
    "Takes Reservations": false,
    "Delivery": false
  }
}
```



CSV Data

Comma instead of period interprets as one additional value

Indication of unknown value

```
5.7,2.9,4.2,?,Iris-versicolor
6.2,2.9,4.3,1.3,Iris-versicolor
5,1,2.5,3.0,1.1,Iris-versicolor
5.7,2.8,4.1,1.3,5.1,Iris-versicolor
6.3,3.3,6.0,2.5,Iris-virginica
N/A,2.7,5.1,1.9,Iris-virginica
7.1,3.0,5.9,2.1
```

Additional value for an observation

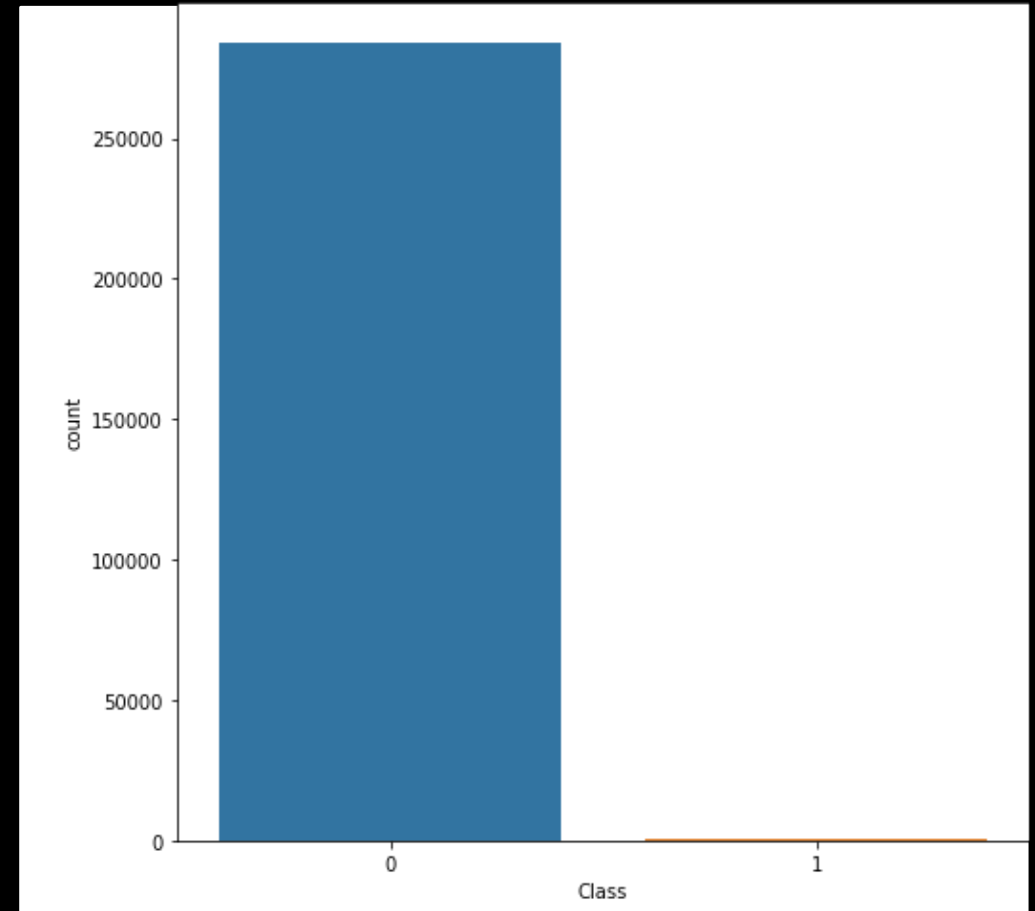
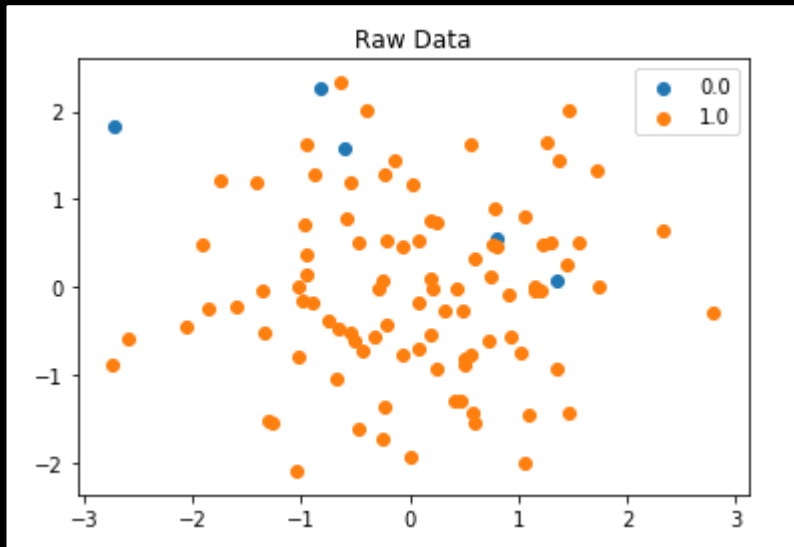
Real value missing, replaced with string

Missing class (insufficient parameters)



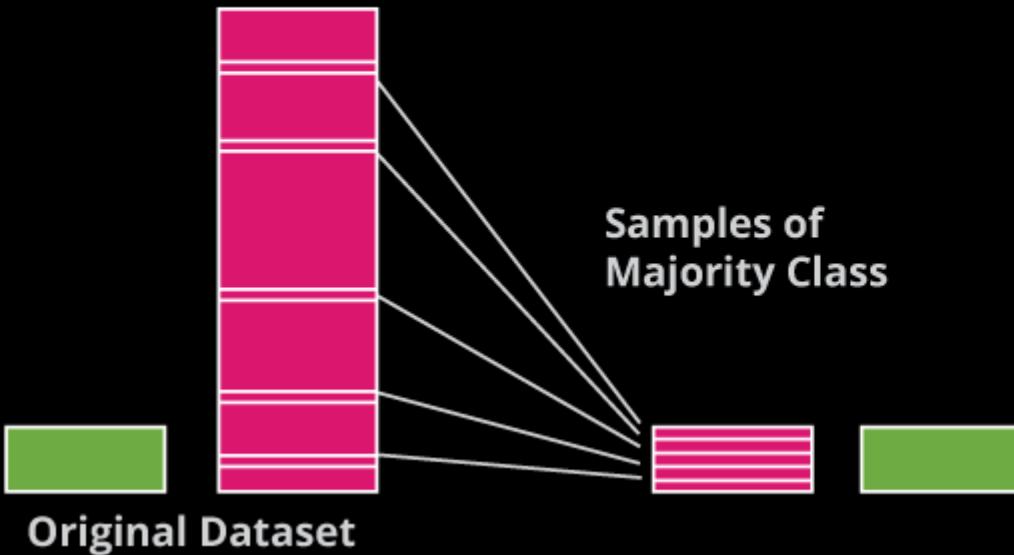
Imbalanced Data

- Data that is heavily skewed to one class or one class having much less data than the other classes.

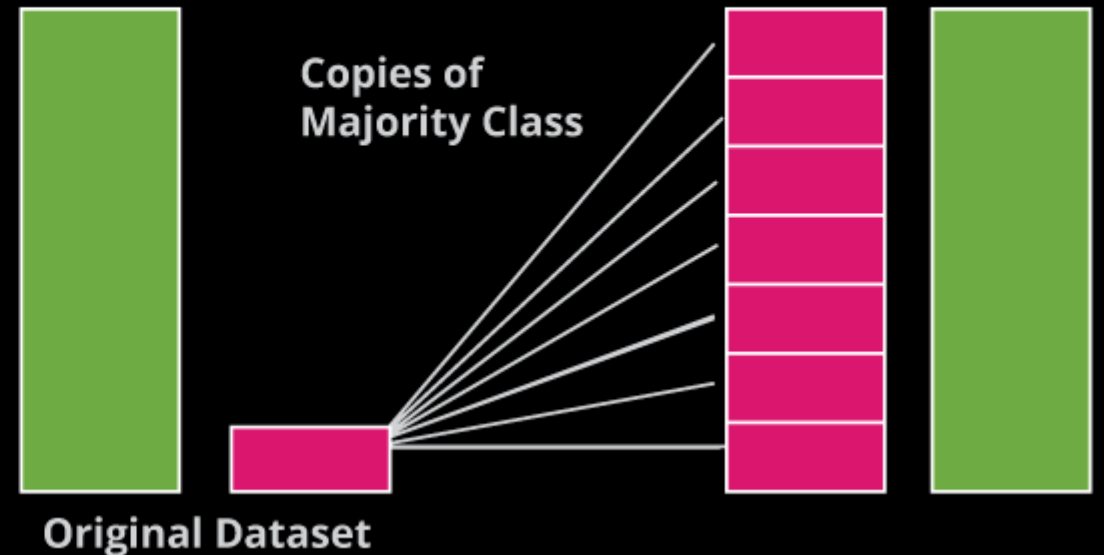


Under and Over Sampling

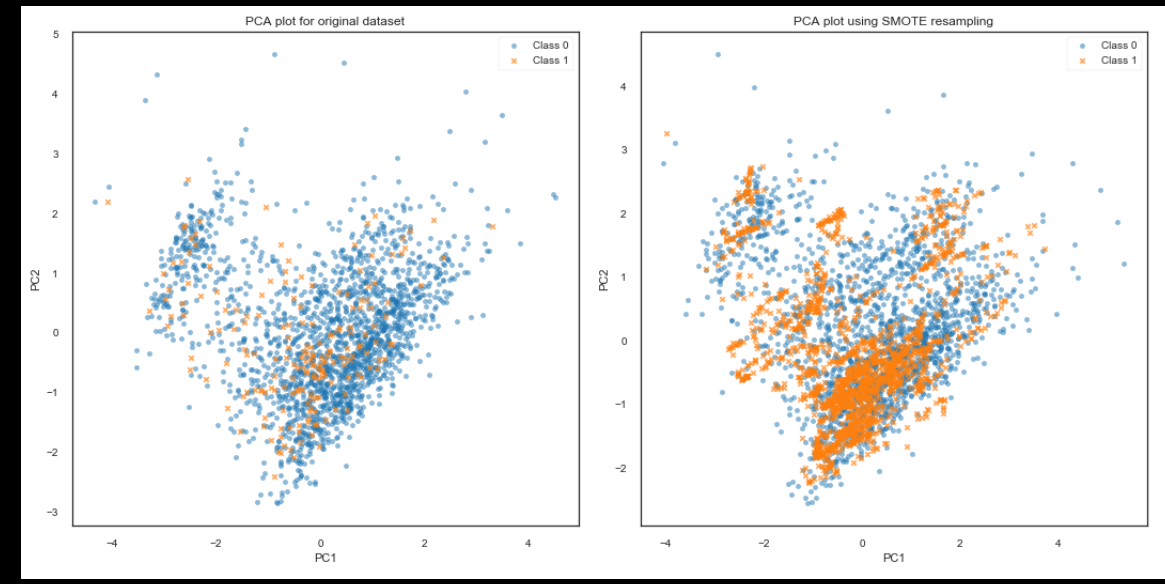
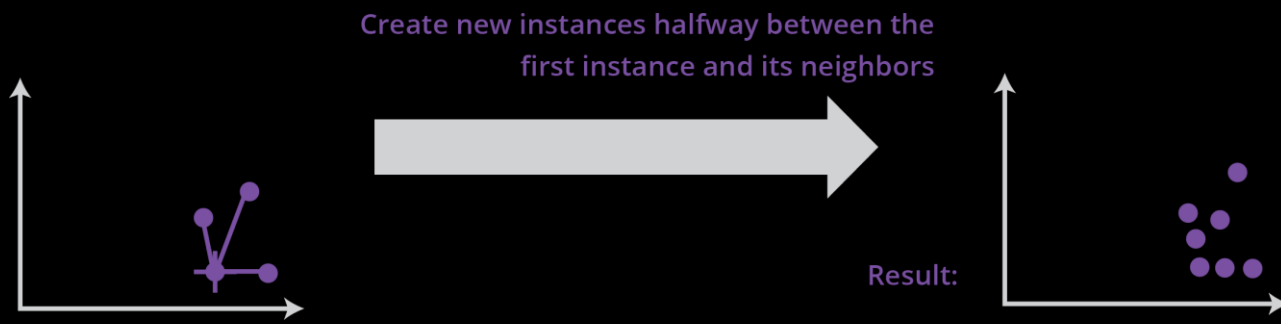
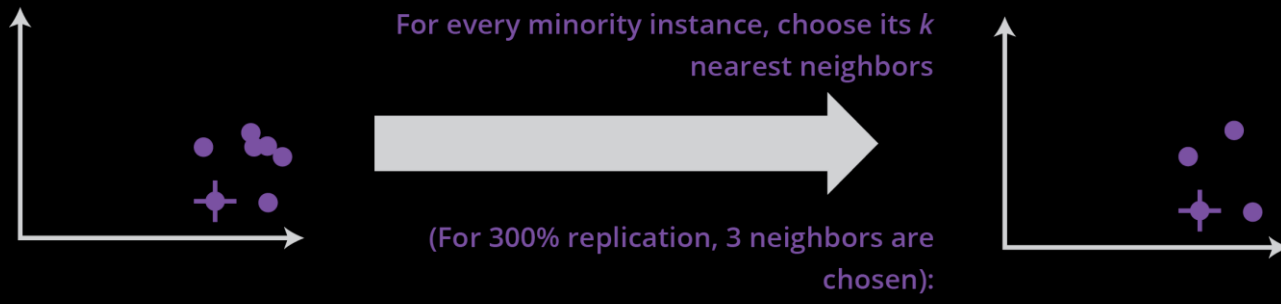
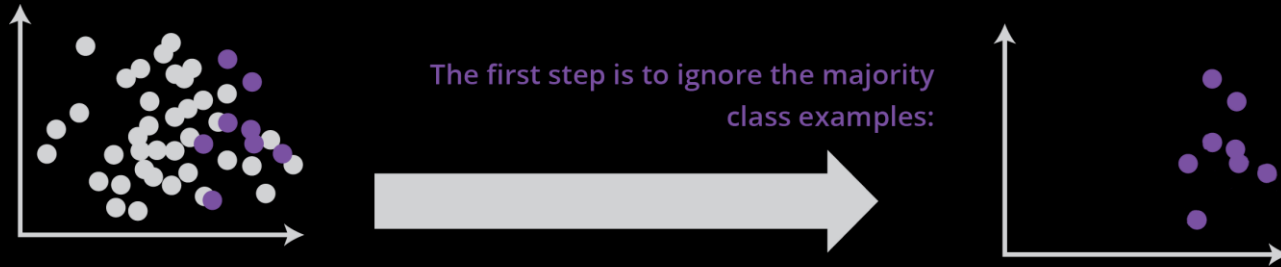
Undersampling



Oversampling



SMOTE



Class weighting (feature weighting)

- Gives importance to certain classes.
- Penalizes our models more for mistakes on classes that are important.
- Is usually built in to SciKit Learn and Keras' models.



Metrics

Metric	Formula
True Positive Rate, Recall	$\frac{TP}{TP+TN}$
False Positive Rate	$\frac{FP}{FP+TN}$
Precision	$\frac{TP}{TP+FP}$
Accuracy	$\frac{TP+TN}{TP+TN+FP+FN}$
F-Measure	$\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$



Ensemble sampling (Bootstrapping)

- Creates many smaller datasets to work with to get a good distribution of the data.
- You can either do K-Fold or Hold testing on this method.
- Works best a majority of the time.

