

Imbalanced Classes

- when there is more of one class than another in a classification task
- common in real world datasets
- Ex: credit card fraud
 - very small number of fraud transactions relative to total transactions

Dealing With Imbalanced Classes

- Stratified Sampling
- Random Undersampling
- Random Oversampling
- Oversample Synthetic Minority Items
 - SMOTE
 - ADASYN
- Other methods

Stratified Sampling

In [2]: `from sklearn.model_selection import StratifiedKFold`

```
X = np.ones(10)
y = [0, 0, 0, 0, 1, 1, 1, 1, 1, 1]

skf = StratifiedKFold(n_splits=3)
for train, test in skf.split(X, y):
    print("%s %s" % (train, test))
```

```
[2 3 6 7 8 9] [0 1 4 5]
[0 1 3 4 5 8 9] [2 6 7]
[0 1 2 4 5 6 7] [3 8 9]
```

Random Sampling

- Randomly Oversample minority class
- Randomly Undersample majority class

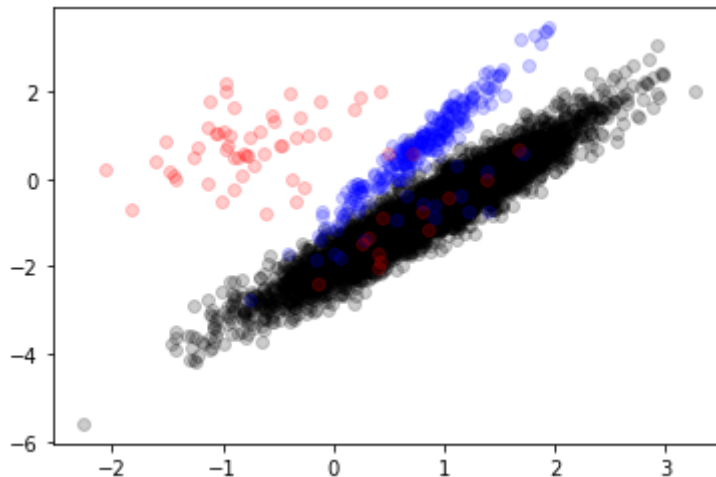
Example Dataset

```
In [3]: from sklearn.datasets import make_classification
from collections import Counter
X, y = make_classification(n_samples=5000, n_features=2, n_informative=2,
                           n_redundant=0, n_repeated=0, n_classes=3,
                           n_clusters_per_class=1,
                           weights=[0.01, 0.05, 0.94],
                           class_sep=0.8, random_state=0)

Counter(y).items()
```

```
Out[3]: dict_items([(2, 4674), (1, 262), (0, 64)])
```

```
In [4]: plt.scatter(X[y==2,0],X[y==2,1],c='k', alpha=.2);
plt.scatter(X[y==1,0],X[y==1,1],c='b', alpha=.2);
plt.scatter(X[y==0,0],X[y==0,1],c='r', alpha=.2);
```



Using imblearn

```
In [5]: # conda install -c conda-forge -n eods-f20 imbalanced-learn
```

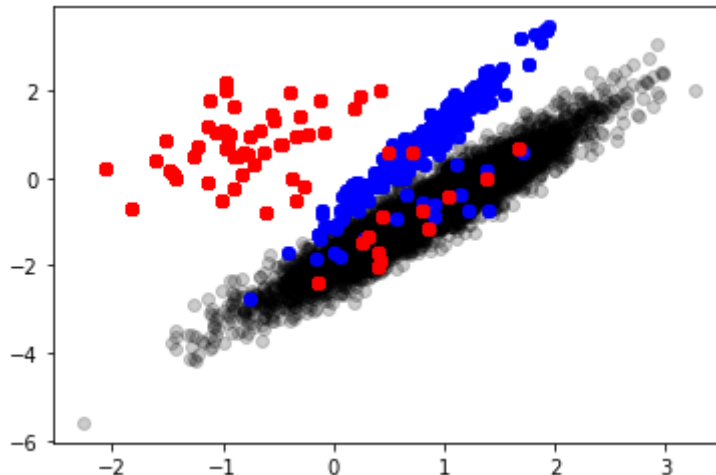
Random Oversampling

```
In [6]: from imblearn.over_sampling import RandomOverSampler
```

```
ros = RandomOverSampler(random_state=0)  
X_r, y_r = ros.fit_sample(X, y)  
Counter(y_r).items()
```

```
Out[6]: dict_items([(2, 4674), (1, 4674), (0, 4674)])
```

```
In [7]: plt.scatter(X_r[y_r==2,0],X_r[y_r==2,1],c='k', alpha=.2);  
plt.scatter(X_r[y_r==1,0],X_r[y_r==1,1],c='b', alpha=.2);  
plt.scatter(X_r[y_r==0,0],X_r[y_r==0,1],c='r', alpha=.2);
```



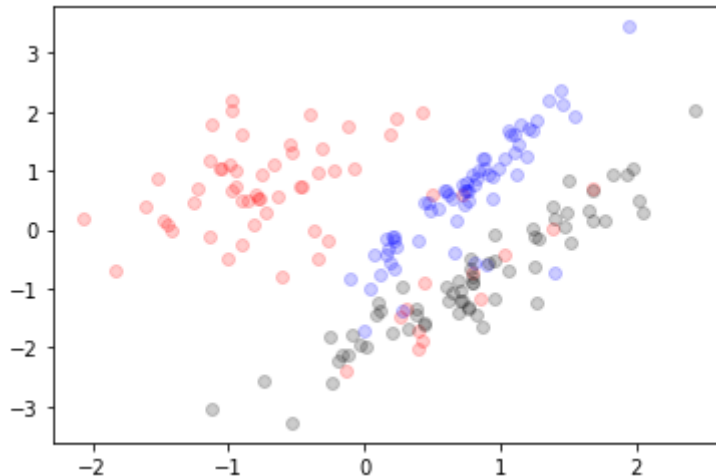
Random Undersampling

In [8]: `from imblearn.under_sampling import RandomUnderSampler`

```
rus = RandomUnderSampler(random_state=0)
X_r, y_r, = rus.fit_sample(X, y)
Counter(y_r).items()
```

Out[8]: `dict_items([(0, 64), (1, 64), (2, 64)])`

In [9]: `plt.scatter(X_r[y_r==0,0],X_r[y_r==0,1],c='r', alpha=.2);`
`plt.scatter(X_r[y_r==1,0],X_r[y_r==1,1],c='b', alpha=.2);`
`plt.scatter(X_r[y_r==2,0],X_r[y_r==2,1],c='k', alpha=.2);`



Oversample Synthetic Minority Items

- SMOTE: Synthetic Minority Oversampling
- ADASYN: Adaptive Synthetic Minority Oversampling

SMOTE: Synthetic Minority Oversampling

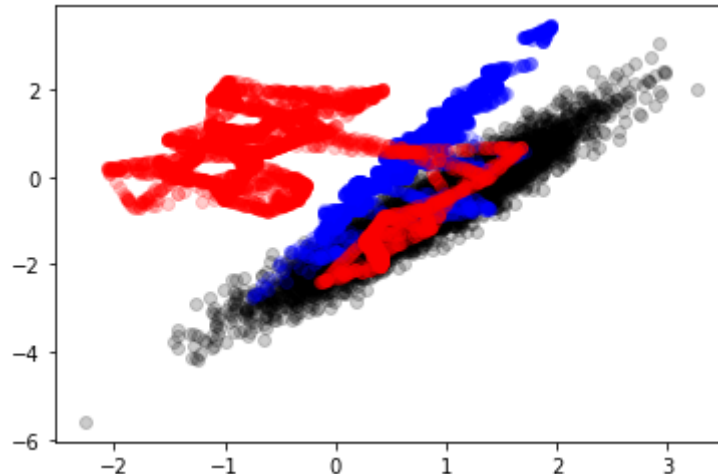
- Create new synthetic points between existing points

```
In [10]: from imblearn.over_sampling import SMOTE
```

```
X_r, y_r = SMOTE().fit_sample(X, y)  
Counter(y_r).items()
```

```
Out[10]: dict_items([(2, 4674), (1, 4674), (0, 4674)])
```

```
In [11]: plt.scatter(X_r[y_r==2,0],X_r[y_r==2,1],c='k', alpha=.2);  
plt.scatter(X_r[y_r==1,0],X_r[y_r==1,1],c='b', alpha=.2);  
plt.scatter(X_r[y_r==0,0],X_r[y_r==0,1],c='r', alpha=.2);
```



ADASYN: Adaptive Synthetic Minority Oversampling

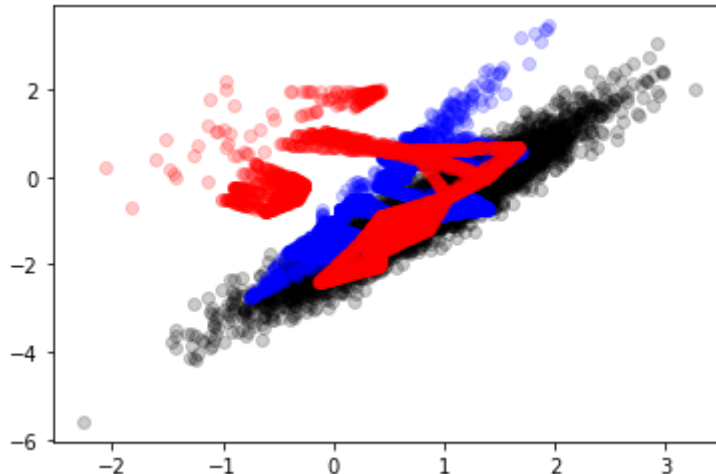
- Create new synthetic points between existing points *where classes overlap*

```
In [12]: from imblearn.over_sampling import ADASYN
```

```
X_r, y_r = ADASYN().fit_sample(X, y)  
Counter(y_r).items()
```

```
Out[12]: dict_items([(2, 4674), (1, 4662), (0, 4673)])
```

```
In [13]: plt.scatter(X_r[y_r==2,0],X_r[y_r==2,1],c='k', alpha=.2);  
plt.scatter(X_r[y_r==1,0],X_r[y_r==1,1],c='b', alpha=.2);  
plt.scatter(X_r[y_r==0,0],X_r[y_r==0,1],c='r', alpha=.2);
```



Other methods for dealing with imbalanced classes

- Adjust class weight (sklearn)
- Adjust decision threshold (sklearn)
- Treat as anomaly detection
- Buy more data

See https://imbalanced-learn.readthedocs.io/en/stable/auto_examples/over-sampling/plot_comparison_over_sampling.html (https://imbalanced-learn.readthedocs.io/en/stable/auto_examples/over-sampling/plot_comparison_over_sampling.html).