

# LAPORAN TUGAS 2

Implementasi Oversampling dan Undersampling Methods pada Imbalance Datasets

## Kelompok 6

Muhammad Adib Arinanda – 5115100111

Wahyu Ivan Satyagraha – 5115100155

Alfindio Muhammad Abdallah – 5115100164

### Library yang digunakan :

1. Pandas
2. Numpy
3. Scikit-learn
4. Matplotlib

### Datasets yang digunakan :

PubChem Bioassay Data Dataset

### STEP :

1. Load dataset.
2. Memisahkan atribut kelas dengan non kelas.
3. Melakukan balancing data dengan metode oversampling (SMOTE) dan undersampling (Near Miss).
4. Menghitung akurasi.
5. Melakukan Feature Reduction dengan PCA.

## HASIL BALANCING DATA DAN AKURASI :

Under-sampling using Near Miss

=====

Original dataset shape Counter({'Inactive': 47538, 'Active': 293})

Accuracy value knn dengan data Imbalance: 0.99 (+/- 0.00)

Resampled dataset shape Counter({'Active': 293, 'Inactive': 293})

=====

Accuracy value knn untuk Near Miss: 0.81 (+/- 0.12)

Over-sampling using SMOTE

=====

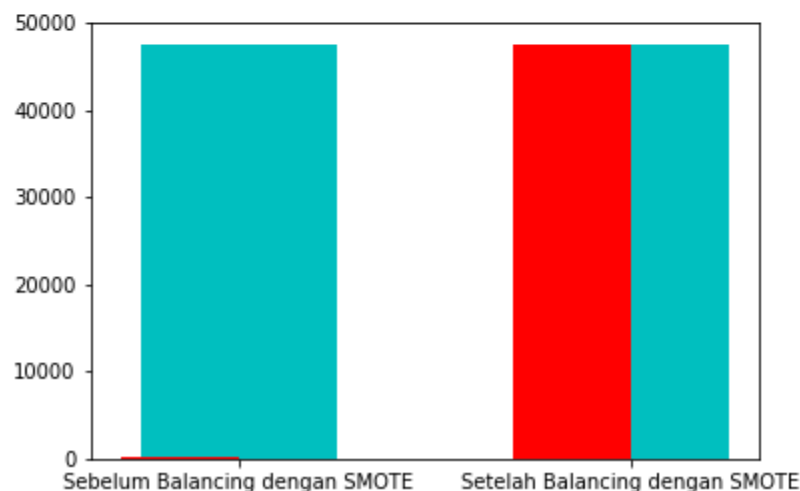
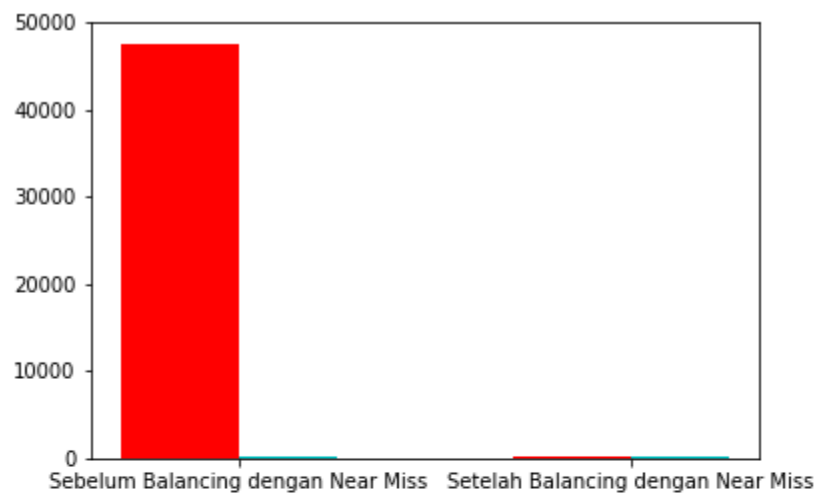
Original dataset shape Counter({'Inactive': 47538, 'Active': 293})

Resampled dataset shape Counter({'Active': 47538, 'Inactive': 47538})

=====

Accuracy value knn untuk SMOTE: 0.91 (+/- 0.02)

## PLOT :



## CODE :

```
7
8 import matplotlib.pyplot as plt
9 from collections import Counter
10 from imblearn.under_sampling import NearMiss
11 from imblearn.over_sampling import SMOTE
12 import pandas as pd
13 from sklearn.decomposition import PCA
14 from sklearn.model_selection import cross_val_score
15 from sklearn import neighbors
16
17 data = pd.read_csv('dataset.csv')
18
19 knn = neighbors.KNeighborsClassifier(n_neighbors=5, algorithm='auto')
20
21 X = data.drop('class', axis=1)
22 y = data['class']
23
24 besar = data.groupby(['class']).size()
25 besar = list(besar)
26
27 koor_x = ['Setelah Balancing dengan Near Miss', 'Sebelum Balancing dengan Near Miss']
28 koor_y = besar
29
30 besar_smote = data.groupby(['class']).size()
31 besar_smote = list(besar_smote)
32
33 koor_x_smote = ['Sebelum Balancing dengan SMOTE', 'Setelah Balancing dengan SMOTE']
34 koor_y_smote = besar_smote
35
36 pca = PCA(n_components=2)
37 X_vis = pca.fit_transform(X)
38
39 print('Under-sampling using Near Miss')
40 print('=====')
41 print('Original dataset shape {}'.format(Counter(y)))
42 #Original dataset shape Counter({1: 900, 0: 100})
43
44 nm = NearMiss(random_state=42)
45 X_under_resampled, y_under_resampled = nm.fit_sample(X, y)
46 X_und_res_vis = pca.transform(X_under_resampled)
47
48 y_res = list(y_under_resampled)
49
50 valp = y_res.count('Inactive')
51 valn = y_res.count('Active')
52
53 new_y = []
54 new_y.append(valp)
55 new_y.append(valn)
56
57 print('Resampled dataset shape {}'.format(Counter(y_under_resampled)))
58 print('=====')
```

```

61 xnya = pd.DataFrame(X_under_resampled)
62 ynya = pd.DataFrame(y_under_resampled)
63 ynya.columns = ['class']
64
65 scores3 = cross_val_score(knn,xnya, ynya['class'], cv=10)
66 print("Accuracy value knn untuk Near Miss: %0.2f (+/- %0.2f)" % (scores3.mean(), scores3.std() * 2))
67
68 print('Over-sampling using SMOTE')
69 print('=====')
70 print('Original dataset shape {}'.format(Counter(y)))
71
72 sm = SMOTE(random_state=42)
73 X_over_resampled, y_over_resampled = sm.fit_sample(X, y)
74 X_ov_res_vis = pca.transform(X_over_resampled)
75
76 y_res_smote = list(y_over_resampled)
77
78 valp_smote = y_res_smote.count('Inactive')
79 valn_smote = y_res_smote.count('Active')
80
81 new_y_smote = []
82 new_y_smote.append(valp_smote)
83 new_y_smote.append(valn_smote)
84
85 print('Resampled dataset shape {}'.format(Counter(y_over_resampled)))
86 print('=====')
87
88 xnya1 = pd.DataFrame(X_over_resampled)
89 ynya1 = pd.DataFrame(y_over_resampled)
90 ynya1.columns = ['class']
91
92 scores3 = cross_val_score(knn,xnya1, ynya1['class'], cv=10)
93 print("Accuracy value knn untuk SMOTE: %0.2f (+/- %0.2f)" % (scores3.mean(), scores3.std() * 2))
94
95 plt.figure(1)
96 plt.bar(koor_x, new_y, label = 'after Near Miss', color='c', width= 0.5, align = 'center')
97 plt.bar(koor_x, koor_y, label = 'before Near Miss', color='r', width= -0.3, align = 'edge')
98
99 plt.figure(2)
100 plt.bar(koor_x_smote, new_y_smote, label = 'after SMOTE', color='c', width= 0.5, align = 'center')
101 plt.bar(koor_x_smote, koor_y_smote, label = 'before SMOTE', color='r', width= -0.3, align = 'edge')
102

```

## KESIMPULAN

Dalam menentukan sebuah akurasi data mining, metode balancing data sangat diperlukan untuk keakuratan hasil yang di peroleh, pada data set ini, oversampling memiliki akurasi yang lebih tinggi namun dengan compile time yang lebih lama disbanding undersampling karena data yang di olah menjadi lebih banyak. Pada data imbalance akurasi yang didapat lebih tinggi namun dalam keadaan imbalance, hasil dari akurasi tersebut masih harus dipertanyakan.