

Penanganan Imbalanced Dataset dan Seleksi Fitur dalam Prediksi Customer Churn

1. Pendahuluan

Customer churn merupakan permasalahan penting dalam industri layanan pelanggan. Untuk memprediksi churn, kita menggunakan model Logistic Regression. Tantangan utama dalam kasus ini adalah ketidakseimbangan kelas (class imbalance), yaitu jumlah pelanggan yang churn jauh lebih sedikit daripada yang tidak. Oleh karena itu, dilakukan beberapa pendekatan penanganan class imbalance dan seleksi fitur.

Untuk source code dan dataset bisa di akses di [repository github](#).

2. Dataset

Dataset yang digunakan adalah data pelanggan dari perusahaan telekomunikasi, di mana target variabel adalah churn_num (0 = tidak churn, 1 = churn). Dataset telah melalui tahap encoding dan preprocessing.

3. Penanganan Imbalanced Class

Berikut adalah empat pendekatan yang digunakan:

3.1. Method 1: Logistic Regression dengan class_weight='balanced'

```
1 logreg_balanced = LogisticRegression(**params)
2 logreg_balanced.fit(X_train, y_train)
3 y_pred_balanced = logreg_balanced.predict(X_test)
4 logit_roc_auc_balanced = roc_auc_score(y_test, y_pred_balanced)
5 print(classification_report(y_test, y_pred))
6 print(f"The area under the curve is: {logit_roc_auc_balanced}")
7
8 with mlflow.start_run():
9     mlflow.log_params(params)
10    mlflow.log_metric("area_under_curve", logit_roc_auc_balanced)
11    mlflow.set_tag("Training Info", "
Basic Logistic regression model with balanced class")
12    signature_balanced = infer_signature(X_train, logreg_balanced.
predict(X_test))
13    model_info_balanced = mlflow.sklearn.log_model(
14        sk_model=logreg_balanced,
15        artifact_path="balanced_log_reg",
16        signature=signature_balanced,
17        input_example=X_train,
18        registered_model_name="balanced_log_reg",
19    )
```

Pendekatan ini mengatur bobot kelas secara otomatis berdasarkan distribusi data

3.2. Method 2: Custom Class Weights

```
1  params = {
2      "random_state" : 13,
3      "class_weight" : {0:90, 1:10}
4  }
5  logreg_custom_ = LogisticRegression(**params)
6  logreg_custom_.fit(X_train, y_train)
7  y_pred_custom_ = logreg_custom_.predict(X_test)
8  logit_roc_auc_custom_ = roc_auc_score(y_test, y_pred_custom_)
9  print(classification_report(y_test, y_pred_custom_))
10 print(f"The area under the curve is: {logit_roc_auc_custom_}")
11
12 with mlflow.start_run():
13     mlflow.log_params(params)
14     mlflow.log_metric("area_under_curve", logit_roc_auc_custom_)
15     mlflow.set_tag("Training Info", "
16     Basic Logistic regression model with custom class weights")
17     signature_balanced = infer_signature(X_train, logreg_balanced.
18     predict(X_test))
19     model_info_balanced = mlflow.sklearn.log_model(
20         sk_model=logreg_custom_,
21         artifact_path="custom_class_weights_log_reg",
22         signature=signature_balanced,
23         input_example=X_train,
24         registered_model_name="custom_class_weights_log_reg",
25     )
```

Bobot disesuaikan secara manual untuk memperbesar kontribusi minoritas.

3.3. Method 3: Resampling Data

```
1 params = {  
2     "random_state":13  
3 }  
4 logreg_resampled = LogisticRegression(**params)  
5  
6 X=resampled_df[resampled_df[resampled_df.select_dtypes(include=np.number).  
7     columns.tolist()].columns.difference(['customer_id','churn_num'])]  
8 y=resampled_df['churn_num']  
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,  
10     random_state=0)
```

Minoritas (kelas 1) di-upsample agar jumlahnya sama dengan mayoritas.

3.4. Method 4: SMOTE (Synthetic Minority Over-sampling Technique)

```

1 from imblearn.over_sampling import SMOTE
2 sm = SMOTE(random_state=0, sampling_strategy=1.0)
3
4 X=df_enc_sel[df_enc_sel[df_enc_sel.select_dtypes(include=np.number).
5 columns.tolist()].columns.difference(['customer_id','churn_num'])]
6 y=df_enc_sel['churn_num']
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3
8 , random_state=0)
9
10 X_train_smote, y_train_smote = sm.fit_resample(X_train, y_train)
11
12 params = {
13     "random_state":13
14 }
15 logreg_smote = LogisticRegression(**params)
16 logreg_smote.fit(X_train_smote, y_train_smote)
17 y_pred_smote = logreg_smote.predict(X_test)
18 logit_roc_auc_smote = roc_auc_score(y_test, y_pred_smote)
19 print(classification_report(y_test, y_pred_smote))
20 print(f"The area under the curve is: {logit_roc_auc_smote}")
21
22 with mlflow.start_run():
23     mlflow.log_params(params)
24     mlflow.log_metric("area_under_curve", logit_roc_auc_smote)
25     mlflow.set_tag("Training Info", "
26     Logistic regression model with SMOTE to fix imbalanced data")
27     signature_smote = infer_signature(X_train_smote, logreg_balanced.
28 predict(X_test))
29     model_info_smote = mlflow.sklearn.log_model(
30         sk_model=logreg_smote,
31         artifact_path="logreg_smote",
32         signature=signature_smote,
33         input_example=X_train_smote,
34         registered_model_name="logreg_smote",
35     )

```

SMOTE membuat sampel sintetis berdasarkan tetangga terdekat dari kelas minoritas.

4. Seleksi Fitur

4.1. Variance Threshold

Fitur dengan variansi rendah dihapus, kemudian training dilakukan dengan SMOTE dan Logistic Regression.

4.2. Recursive Feature Elimination (RFE)

```
1 from sklearn.feature_selection import RFE
2 rfe = RFE(model)
3 rfe = rfe.fit(X, y)
4 selected_cols_rfe = X.columns[rfe.support_]
5 print(list(selected_cols_rfe))
6
7 X_rfe_cols=df_enc_sel[['SeniorCitizen', 'contract_code', 'dependents_code',
8   'multiple_lines_code', 'online_backup_code',
9   'online_security_code', 'paperless_billing_code', '
10  phone_service_code', 'tech_support_code']]
11 y=df_enc_sel['churn_num']
12 X_train_rfe, X_test_rfe, y_train_rfe, y_test_rfe = train_test_split(
13   X_rfe_cols, y, test_size=0.3, random_state=0)
14 sm = SMOTE(random_state=0, sampling_strategy=1.0)
15 X_train_rfe, y_train_rfe = sm.fit_resample(X_train_rfe, y_train_rfe)
16 params = {
17     "random_state":13
18 }
```

RFE memilih fitur paling penting berdasarkan performa model.

5. Evaluasi Model

Evaluasi dilakukan menggunakan klasifikasi report dan ROC AUC score. Berikut beberapa hasil:

```

['SeniorCitizen', 'contract_code', 'dependents_code', 'multiple_lines_code', 'online_backup_code',
  precision    recall  f1-score   support

     0.0         0.91     0.69     0.78     1555
     1.0         0.48     0.81     0.61      555

 accuracy         0.72     2110
 macro avg         0.70     0.75     0.69     2110
weighted avg         0.80     0.72     0.74     2110

The area under the curve is: 0.7506792966600041
Successfully registered model 'logreg_rfe_t'.
2025/04/28 19:31:41 INFO mlflow.store.model_registry.abstract_store: Waiting up to 300 seconds for
Created version '1' of model 'logreg_rfe_t'.
precision    recall  f1-score   support

     0.0         0.83     0.61     0.71     1555
     1.0         0.38     0.66     0.48      555

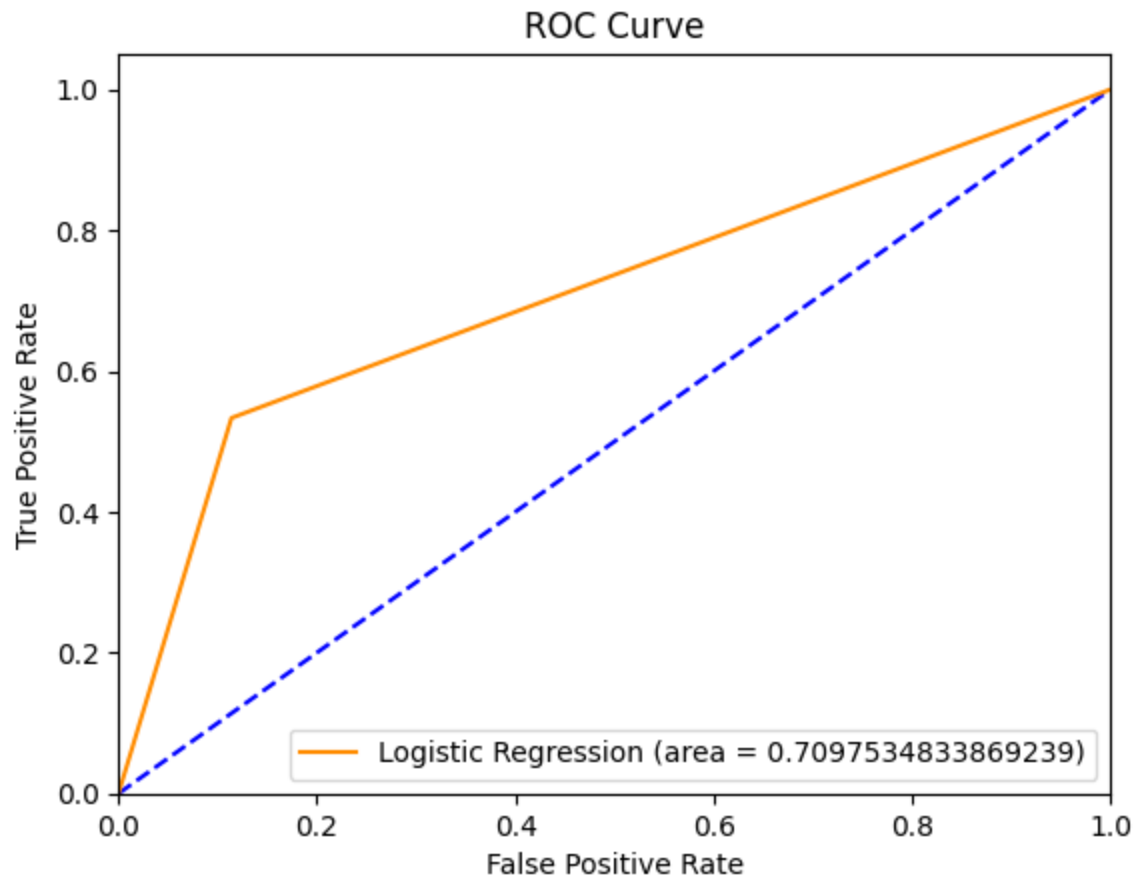
 accuracy         0.62     2110
 macro avg         0.60     0.63     0.59     2110
weighted avg         0.71     0.62     0.65     2110

The area under the curve is: 0.6337157092784103
Successfully registered model 'logreg_smote_t'.
2025/04/28 19:31:13 INFO mlflow.store.model_registry.abstract_store: Waiting up to 300 seconds for
Created version '1' of model 'logreg_smote_t'.

```

6. Visualisasi

Confusion Matrix dan ROC Curve (Contoh)



7. Logging dengan MLflow

Semua model dilog menggunakan MLflow untuk tracking parameter, metrik, dan artefak model.

```
1 with mlflow.start_run():
2     mlflow.log_params(params)
3     mlflow.log_metric("area_under_curve", logit_roc_auc_rfe_t)
4     mlflow.set_tag("Training Info", "
5     Logistic regression model with SMOTE & RFE for feature selection")
6     signature_rfe_t = infer_signature(X_train_rfe, logreg_balanced.
7     predict(X_test))
8     model_info_rfe_t = mlflow.sklearn.log_model(
9         sk_model=logreg_rfe_t,
10        artifact_path="logreg_rfe_t",
11        signature=signature_rfe_t,
12        input_example=X_train_rfe,
13        registered_model_name="logreg_rfe_t",
14    )
```

8. Kesimpulan

SMOTE dan penggunaan `class_weight='balanced'` terbukti membantu meningkatkan performa model dalam kasus data tidak seimbang. Rekomendasi terbaik adalah menggabungkan teknik balancing dengan seleksi fitur seperti RFE untuk hasil optimal.