

Analyze the provided dataset to identify its key characteristics for data-driven decision-making. Upon receiving a specific analytical question, address it strictly using the dataset. Your deliverables must include a clear, concise answer to the question and the Python code that extracts and processes the data used to derive your answer. The dataset should be cleaned by filtering admission\_period to include only values between 2006.1 and 2019.2, removing entries where enrollment\_type is labeled as "Exempted", and excluding those where situation is "In Progress". The resulting cleaned dataset should be stored in a new DataFrame named df\_filtered.

Anexo oculto

### Key Characteristics of the Dataset for Data-Driven Decision-Making

The dataset appears to be a detailed record of student academic performance and enrollment history at an educational institution. Below are the key characteristics identified for data-driven decision-making:

#### 1. Student Demographics:

- `gender`, `color_code`, `color_name`, `admission_age`, `dropout_age`: Useful for analyzing diversity, equity, and age-related trends in enrollment and dropout rates.

#### 2. Admission and Enrollment Details:

- `admission_code`, `admission_period`, `admission_form`, `enrollment_period`, `enrollment_type_code`, `enrollment_type`: Enable tracking of admission trends over time and the impact of different admission methods (e.g., Entrance Exam, SISU) on student success.

#### 3. Academic Performance:

- `general_average`, `final_average`, `situation_code`, `situation`: Provide insights into student success rates, subject difficulty, and factors influencing grades.

#### 4. Dropout and Exit Information:

- `dropout_code`, `dropout_form`, `dropout_period`, `dropout_period_count`, `exit_form`: Allow analysis of retention rates, reasons for dropout, and time-to-dropout patterns.

#### 5. Course and Subject Details:

- `curriculum_code`, `subject_code`, `subject_name`, `credits`, `workload`, `sector_name`: Facilitate evaluation of curriculum effectiveness, subject-specific performance, and workload impact.

#### 6. Quota and Affirmative Action:

- `affirmative_action`, `quota_type`, `quota_form`: Support analysis of the effectiveness of inclusion policies on student outcomes.

#### 7. High School Background:

- `high_school_type`, `high_school`: Help assess how prior education influences higher education performance.

#### 8. Temporal Data:

- `admission_period`, `dropout_period`, `enrollment_period`, `relative_period`: Enable longitudinal analysis of trends over time.

These characteristics make the dataset suitable for addressing questions related to student retention, academic performance, equity in education, and curriculum optimization.

### Data Cleaning

The dataset will be cleaned according to the specified requirements:

- Filter `admission_period` to include only values between 2006.1 and 2019.2 (inclusive).
- Remove entries where `enrollment_type` is "Exempted".
- Exclude entries where `situation` is "In Progress".
- Store the cleaned dataset in a new DataFrame called `df_filtered`.

Below is the Python code to perform the cleaning:

```
python ... Copiar
import pandas as pd

# Sample data provided as a string (in practice, this would be loaded from a file)
data = """
id.curriculum_code.admission_code.admission_period.admission_form.dropout_code.dropdown_f
```

Q386391,1990,1,1979.2,Entrance Exam,4,Cancellation Due to Abandonment,1989.1,20.0,Dropou  
A261799,1999,16,2011.1,Entrance Exam / ENEM,1,Graduate,2016.1,11.0,Graduate,No,1,Without  
K266013,2017,18,2017.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
H046905,1999,16,2014.1,Entrance Exam / ENEM,21,Canceled Failed All Due to Absences,2016.  
T079917,1990,1,1996.1,Entrance Exam,1,Graduate,1999.2,8.0,Graduate,No,1,Without Quota Re  
Z239214,1999,1,2005.2,Entrance Exam,1,Graduate,2008.2,7.0,Graduate,No,1,Without Quota Re  
S374319,1990,1,1989.1,Entrance Exam,4,Cancellation Due to Abandonment,1998.1,19.0,Dropou  
G381246,1999,1,2003.1,Entrance Exam,1,Graduate,2007.1,9.0,Graduate,No,1,Without Quota Re  
S126527,2017,18,2017.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
X175376,1999,16,2011.2,Entrance Exam / ENEM,8,Cancellation by Student Request,2013.2,5.0  
C331369,1999,1,2005.1,Entrance Exam,4,Cancellation Due to Abandonment,2008.2,8.0,Dropou  
I362687,1990,1,1995.2,Entrance Exam,1,Graduate,2000.1,10.0,Graduate,No,1,Without Quota R  
C332956,1999,1,2009.2,Entrance Exam,1,Graduate,2013.2,9.0,Graduate,No,1,Without Quota Re  
O215120,1990,1,1983.1,Entrance Exam,1,Graduate,1987.1,9.0,Graduate,No,1,Without Quota Re  
X131461,1990,1,1992.1,Entrance Exam,4,Cancellation Due to Abandonment,2000.1,17.0,Dropou  
Y164857,2017,18,2018.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
W033589,1999,1,1997.2,Entrance Exam,1,Graduate,2001.1,8.0,Graduate,No,1,Without Quota Re  
D108883,1990,1,1978.2,Entrance Exam,1,Graduate,1982.2,9.0,Graduate,No,1,Without Quota Re  
Q283427,1999,1,2010.2,Entrance Exam,12,Canceled New Enrollment Same Course,2012.1,4.0,Re  
E158583,1990,1,1977.2,Entrance Exam,4,Cancellation Due to Abandonment,1992.1,30.0,Dropou  
U060448,1999,1,2000.2,Entrance Exam,1,Graduate,2004.1,8.0,Graduate,No,1,Without Quota Re  
A225959,1999,16,2013.1,Entrance Exam / ENEM,22,Canceled 3 Failures in Same Subject,2017.  
A281268,1999,1,2009.2,Entrance Exam,1,Graduate,2013.2,9.0,Graduate,No,1,Without Quota Re  
T020349,1999,1,2004.2,Entrance Exam,1,Graduate,2005.2,3.0,Graduate,No,1,Without Quota Re  
R024758,1999,1,2010.2,Entrance Exam,12,Canceled New Enrollment Same Course,2013.2,7.0,Re  
M330399,1990,1,1981.1,Entrance Exam,1,Graduate,1991.1,21.0,Graduate,No,1,Without Quota R  
S001296,1999,16,2011.1,Entrance Exam / ENEM,1,Graduate,2016.1,11.0,Graduate,No,1,Without  
P213040,1990,1,1986.1,Entrance Exam,4,Cancellation Due to Abandonment,1988.1,5.0,Dropout  
H267951,1999,1,2001.2,Entrance Exam,1,Graduate,2008.1,14.0,Graduate,No,1,Without Quota R  
D152066,1990,1,1993.1,Entrance Exam,4,Cancellation Due to Abandonment,1995.2,6.0,Dropout  
U157075,2017,18,2014.2,Unified Selection System (SISU),22,Canceled 3 Failures in Same Su  
A306762,1990,1,1979.2,Entrance Exam,4,Cancellation Due to Abandonment,1992.2,27.0,Dropou  
C381009,1999,1,2006.1,Entrance Exam,1,Graduate,2010.2,10.0,Graduate,No,1,Without Quota R  
F276339,1999,16,2012.1,Entrance Exam / ENEM,1,Graduate,2016.2,10.0,Graduate,No,1,Without  
P052203,1999,1,2000.2,Entrance Exam,1,Graduate,2007.2,15.0,Graduate,No,1,Without Quota R  
X188358,1990,1,1995.1,Entrance Exam,1,Graduate,1999.1,9.0,Graduate,No,1,Without Quota Re  
P095116,1999,16,2014.1,Entrance Exam / ENEM,1,Graduate,2019.1,11.0,Graduate,No,1,Without  
E327596,2017,18,2019.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,2,Scho  
A162943,1999,16,2011.1,Entrance Exam / ENEM,4,Cancellation Due to Abandonment,2016.1,11.  
U312547,1999,16,2012.2,Entrance Exam / ENEM,1,Graduate,2015.1,6.0,Graduate,No,1,Without  
X057867,2017,18,2016.2,Unified Selection System (SISU),4,Cancellation Due to Abandonment  
W257210,1999,16,2013.2,Entrance Exam / ENEM,1,Graduate,2017.2,9.0,Graduate,No,1,Without  
W189220,1999,1,2007.1,Entrance Exam,12,Canceled New Enrollment Same Course,2010.1,7.0,Re  
Z280243,2017,18,2016.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
M027306,1999,1,2001.1,Entrance Exam,4,Cancellation Due to Abandonment,2005.2,10.0,Dropou  
P206861,2017,18,2018.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,4,Scho  
F050322,1999,1,2004.2,Entrance Exam,1,Graduate,2009.2,11.0,Graduate,No,1,Without Quota R  
Y225407,1979,1,1978.2,Entrance Exam,4,Cancellation Due to Abandonment,1985.2,15.0,Dropou  
K297836,2017,2,2016.2,Transfer,21,Canceled Failed All Due to Absences,2019.2,7.0,Dropout  
B227546,1999,1,2001.1,Entrance Exam,1,Graduate,2008.1,15.0,Graduate,No,1,Without Quota R  
S346697,2017,6,2015.1,Reoption,0,Regular,0.0,,Regular,No,1,Without Quota Reservation,Mal  
B081857,1999,1,2008.2,Entrance Exam,1,Graduate,2013.1,10.0,Graduate,No,1,Without Quota R  
M003781,1999,1,2008.1,Entrance Exam,1,Graduate,2012.1,9.0,Graduate,No,1,Without Quota Re  
C010158,1990,1,1989.1,Entrance Exam,5,Enrollment Cancellation,1989.2,2.0,Dropout,No,1,Wi  
C168059,1999,1,2001.1,Entrance Exam,6,Cancellation Due to Course Change,2003.2,6.0,Dropo  
I102693,1999,1,2009.1,Entrance Exam,1,Graduate,2013.1,9.0,Graduate,No,1,Without Quota Re  
P263567,1990,1,1992.1,Entrance Exam,1,Graduate,1996.2,10.0,Graduate,No,1,Without Quota R  
B002015,2017,18,2017.1,Unified Selection System (SISU),8,Cancellation by Student Request  
B195251,1999,1,2008.1,Entrance Exam,1,Graduate,2012.1,9.0,Graduate,No,1,Without Quota Re  
F395109,1999,1,2004.2,Entrance Exam,22,Canceled 3 Failures in Same Subject,2010.1,12.0,D  
0281950,1999,18,2014.2,Unified Selection System (SISU),1,Graduate,2016.2,5.0,Graduate,No  
M050872,2017,18,2015.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
L221296,2017,18,2016.1,Unified Selection System (SISU),4,Cancellation Due to Abandonment  
P101304,1999,1,2010.1,Entrance Exam,8,Cancellation by Student Request,2013.2,8.0,Dropout  
T175186,1999,16,2011.1,Entrance Exam / ENEM,1,Graduate,2015.2,10.0,Graduate,No,1,Without  
I286150,1999,1,1998.2,Entrance Exam,1,Graduate,2006.1,16.0,Graduate,No,1,Without Quota R  
Z391555,1999,16,2014.1,Entrance Exam / ENEM,12,Canceled New Enrollment Same Course,2016.  
L137292,2017,18,2015.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,3,Scho  
I305593,1990,1,1993.1,Entrance Exam,1,Graduate,1996.2,8.0,Graduate,No,1,Without Quota Re  
H213513,1999,1,2007.2,Entrance Exam,1,Graduate,2007.2,1.0,Graduate,No,1,Without Quota Re  
Z054606,2017,18,2017.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
K315331,1999,16,2012.2,Entrance Exam / ENEM,22,Canceled 3 Failures in Same Subject,2014.  
G285770,2017,18,2016.2,Unified Selection System (SISU),22,Canceled 3 Failures in Same Su  
I228730,1999,16,2011.2,Entrance Exam / ENEM,22,Canceled 3 Failures in Same Subject,2015.  
V300924,2017,18,2016.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,2,Scho  
Y345407,1990,1,1993.1,Entrance Exam,1,Graduate,1998.1,11.0,Graduate,No,1,Without Quota R  
B129851,2017,18,2017.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,2,Scho  
M299294,1999,16,2012.2,Entrance Exam / ENEM,22,Canceled 3 Failures in Same Subject,2014.  
Y377848,1990,1,1992.1,Entrance Exam,1,Graduate,1995.2,8.0,Graduate,No,1,Without Quota Re  
T300079,1990,1,1994.2,Entrance Exam,1,Graduate,1999.1,10.0,Graduate,No,1,Without Quota R  
W240314,2017,18,2018.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho  
O270316,1990,1,1978.2,Entrance Exam,4,Cancellation Due to Abandonment,1989.1,22.0,Dropou  
K108952,2017,11,2015.2,Judicial Decision,0,Regular,0.0,,Regular,No,1,Without Quota Reser  
U234644,1999,16,2013.2,Entrance Exam / ENEM,1,Graduate,2018.2,11.0,Graduate,No,1,Without  
P300538,1999,16,2012.2,Entrance Exam / ENEM,1,Graduate,2017.1,10.0,Graduate,No,1,Without  
H252816,1990,1,1995.2,Entrance Exam,1,Graduate,1999.1,8.0,Graduate,No,1,Without Quota Re  
P051362,1999,16,2012.1,Entrance Exam / ENEM,1,Graduate,2015.1,7.0,Graduate,No,1,Without  
D162709,1999,1,2003.1,Entrance Exam,2,Transferred to Another HEI,2008.1,11.0,Dropout,No,  
S249423,1999,2,2012.1,Transfer,1,Graduate,2016.2,10.0,Graduate,No,1,Without Quota Reserv  
P080257,1990,1,1993.1,Entrance Exam,1,Graduate,1996.2,8.0,Graduate,No,1,Without Quota Re  
V249525,1999,18,2016.1,Unified Selection System (SISU),21,Canceled Failed All Due to Abs  
T324243,1990,1,1996.1,Entrance Exam,1,Graduate,2000.1,9.0,Graduate,No,1,Without Quota Re  
E144413,1999,1,2003.1,Entrance Exam,1,Graduate,2007.1,9.0,Graduate,No,1,Without Quota Re

```
J346815,1990,1,1991.1,Entrance Exam,1,Graduate,1995.2,10.0,Graduate,No,1,Without Quota R
D166923,1999,16,2014.1,Entrance Exam / ENEM,1,Graduate,2018.1,9.0,Graduate,No,1,Without
N263387,1999,1,2007.1,Entrance Exam,1,Graduate,2011.2,10.0,Graduate,No,1,Without Quota R
Z071605,2017,18,2018.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho
V086960,1999,1,2010.1,Entrance Exam,14,Canceled New Enrollment Different Course,2012.1,5
T194884,1999,1,1997.2,Entrance Exam,5,Enrollment Cancellation,2004.2,15.0,Dropout,No,1,W
T265302,1990,1,1993.1,Entrance Exam,1,Graduate,2001.1,17.0,Graduate,No,1,Without Quota R
M163283,1999,16,2011.1,Entrance Exam / ENEM,12,Canceled New Enrollment Same Course,2012.
N223673,2017,18,2020.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,2,Scho
F358958,1990,1,1991.1,Entrance Exam,1,Graduate,1995.2,10.0,Graduate,No,1,Without Quota R
R352812,2017,18,2018.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,No,1,Witho
X313301,1990,1,1977.2,Entrance Exam,4,Cancellation Due to Abandonment,1986.1,18.0,Dropou
P165264,1999,1,2003.2,Entrance Exam,1,Graduate,2009.2,13.0,Graduate,No,1,Without Quota R
S201095,1999,16,2011.2,Entrance Exam / ENEM,8,Cancellation by Student Request,2013.2,5.0
I212832,1999,1,2008.2,Entrance Exam,1,Graduate,2012.2,9.0,Graduate,No,1,Without Quota Re
G001744,1990,1,1982.1,Entrance Exam,1,Graduate,1986.2,10.0,Graduate,No,1,Without Quota R
G114877,1990,1,1989.1,Entrance Exam,1,Graduate,1992.2,8.0,Graduate,No,1,Without Quota Re
V351549,1999,18,2016.2,Unified Selection System (SISU),1,Graduate,2019.1,6.0,Graduate,No
X223251,1999,16,2013.1,Entrance Exam / ENEM,1,Graduate,2017.2,10.0,Graduate,No,1,Without
J352733,1999,1,2002.1,Entrance Exam,1,Graduate,2006.2,10.0,Graduate,No,1,Without Quota R
X249394,2017,18,2019.2,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,4,Scho
L342881,1979,1,1983.1,Entrance Exam,1,Graduate,1987.2,10.0,Graduate,No,1,Without Quota R
A050159,1999,1,2009.1,Entrance Exam,1,Graduate,2015.1,13.0,Graduate,No,1,Without Quota R
Q176954,1999,16,2011.2,Entrance Exam / ENEM,1,Graduate,2012.2,3.0,Graduate,No,1,Without
S319195,1990,1,1978.2,Entrance Exam,4,Cancellation Due to Abandonment,1987.1,18.0,Dropou
0244945,2017,18,2016.1,Unified Selection System (SISU),0,Regular,0.0,,Regular,Yes,2,Scho
S163537,1999,1,2004.1,Entrance Exam,1,Graduate,2011.1,15.0,Graduate,No,1,Without Quota R
G317054,1990,1,1993.1,Entrance Exam,1,Graduate,1998.2,12.0,Graduate,No,1,Without Quota R
N356122,1979,1,1979.2,Entrance Exam,1,Graduate,1985.2,13.0,Graduate,No,1,Without Quota R
X358074,1999,1,2000.2,Entrance Exam,1,Graduate,2003.2,7.0,Graduate,No,1,Without Quota Re
G246877,1999,16,2013.2,Entrance Exam / ENEM,12,Canceled New Enrollment Same Course,2015.
Y161414,1999,1,1998.1,Entrance Exam,1,Graduate,2004.2,14.0,Graduate,No,1,Without Quota R
J218774,1990,1,1977.1,Entrance Exam,1,Graduate,1984.1,15.0,Graduate,No,1,Without Quota R
T301512,1999,16,2011.2,Entrance Exam / ENEM,22,Canceled 3 Failures in Same Subject,2012.
P298203,1999,1,1999.2,Entrance Exam,1,Graduate,2003.2,9.0,Graduate,No,1,Without Quota Re
0259273,1999,1,2001.2,Entrance Exam,1,Graduate,2005.1,8.0,Graduate,No,1,Without Quota Re
W245674,1999,1,2010.1,Entrance Exam,12,Canceled New Enrollment Same Course,2012.2,6.0,Re
T033393,1999,16,2013.2,Entrance Exam / ENEM,21,Canceled Failed All Due to Absences,2015.
```

Estimate the probability of a student's exit mode (graduation vs. dropout) using mode of admission, period of admission, and academic status. Propose a modeling approach and interpret its predictive performance.

To estimate the probability of a student's exit mode (graduation vs. dropout) using mode of admission (`admission_form`), period of admission (`admission_period`), and academic status (`situation`), we'll follow a structured approach. Below, I outline the steps, propose a modeling approach, and provide an interpretation of its predictive performance based on typical outcomes for such a dataset. Since I don't have direct access to the full dataset, I'll assume the provided sample is representative and use it to guide the process.

### Step 1: Data Preparation

First, we clean and preprocess the dataset as specified:

- Filter `admission_period` between 2006.1 and 2019.2.
- Remove entries where `enrollment_type` is "Exempted".
- Exclude entries where `situation` is "In Progress".
- Define the target variable (`exit_mode`) as:
  - "Graduation" if `exit_form` is "Graduate".
  - "Dropout" if `exit_form` is "Dropout" or "Reentry" (assuming reentry implies a prior dropout).

Here's how we'd preprocess the data in Python:

```
python
...
Copiar

import pandas as pd

# Load the dataset (assuming it's in a CSV file or DataFrame)
df = pd.read_csv("student_data.csv") # Replace with actual data source

# Filter admission_period
df['admission_period'] = df['admission_period'].astype(float)
df = df[(df['admission_period'] >= 2006.1) & (df['admission_period'] <= 2019.2)]

# Remove Exempted enrollment_type
df = df[df['enrollment_type'] != "Exemption"]

# Remove In Progress situations
df = df[df['situation'] != "In Progress"]

# Define target variable: exit_mode
df['exit_mode'] = df['exit_form'].apply(lambda x: 'Graduation' if x == 'Graduate' else 'Dropout')
```

```

df = df.dropna(subset=['exit_mode']) # Drop rows where exit_mode is undefined

# Select relevant features
df_model = df[['admission_form', 'admission_period', 'situation', 'exit_mode']]

```

## Step 2: Exploratory Analysis

Before modeling, we'd explore the relationships between the predictors and `exit_mode`:

- **Admission Form:** Categories like "Entrance Exam," "SISU," or "Entrance Exam / ENEM" might influence completion rates due to differences in student preparedness.
- **Admission Period:** Trends over time (e.g., earlier vs. later periods) could reflect changes in institutional policies or student demographics.
- **Situation:** Academic status (e.g., "Approved," "Failed," "Failed Due to Absences") likely correlates strongly with graduation or dropout.

For example, we could compute basic statistics:

```

python ... ⚡ Copiar

# Proportion of graduation vs. dropout by admission_form
print(df_model.groupby('admission_form')['exit_mode'].value_counts(normalize=True))

# Average admission_period by exit_mode
print(df_model.groupby('exit_mode')['admission_period'].mean())

# Situation distribution by exit_mode
print(df_model.groupby('situation')['exit_mode'].value_counts(normalize=True))

```

## Step 3: Proposed Modeling Approach

I propose using a **Logistic Regression** model due to its simplicity, interpretability, and suitability for binary classification (graduation = 1, dropout = 0). Here's why:

- **Binary Outcome:** The target (`exit_mode`) is dichotomous.
- **Categorical and Continuous Predictors:** Logistic regression handles both `admission_form` (categorical) and `admission_period` (continuous) well with proper encoding.
- **Interpretability:** Coefficients provide direct insight into how each predictor affects the odds of graduation.

### Preprocessing for Modeling

- **Encode Categorical Variables:** Use one-hot encoding for `admission_form` and `situation`.
- **Standardize Continuous Variable:** Scale `admission_period` to ensure comparability.
- **Train-Test Split:** Split the data to evaluate model performance.

```

python ... ⚡ Copiar

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, roc_auc_score, classification_report

# Encode categorical variables
df_encoded = pd.get_dummies(df_model, columns=['admission_form', 'situation'], drop_first=True)

# Define features (X) and target (y)
X = df_encoded.drop('exit_mode', axis=1)
y = (df_encoded['exit_mode'] == 'Graduation').astype(int) # 1 = Graduation, 0 = Dropout

# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Scale admission_period
scaler = StandardScaler()
X_train['admission_period'] = scaler.fit_transform(X_train[['admission_period']])
X_test['admission_period'] = scaler.transform(X_test[['admission_period']])

# Fit logistic regression
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

# Predict probabilities
y_pred_prob = model.predict_proba(X_test)[:, 1]
y_pred = (y_pred_prob >= 0.5).astype(int)

```

## Step 4: Predictive Performance Metrics

To evaluate the model, we'd use:

- **Accuracy:** Proportion of correct predictions.
- **ROC-AUC Score:** Measures the model's ability to distinguish between classes (ideal for imbalanced data).
- **Classification Report:** Precision, recall, and F1-score for both classes.

```
python
...
Copiar
```

```
# Evaluate the model
print("Accuracy:", accuracy_score(y_test, y_pred))
print("ROC-AUC:", roc_auc_score(y_test, y_pred_prob))
print(classification_report(y_test, y_pred, target_names=['Dropout', 'Graduation']))
```

### Hypothetical Results

Based on typical educational datasets:

- **Accuracy:** ~75-85% (assuming moderate class imbalance and predictive power of features).
- **ROC-AUC:** ~0.80-0.90 (indicating good discrimination between graduation and dropout).
- **Precision/Recall:**
  - Graduation (majority class): High precision (~0.85), moderate recall (~0.75).
  - Dropout (minority class): Lower precision (~0.60), higher recall (~0.70).

### Step 5: Interpretation of Predictive Performance

#### 1. Model Fit:

- An ROC-AUC of 0.85 suggests the model effectively separates graduates from dropouts, capturing meaningful patterns in the data.
- Accuracy alone might be misleading if the dataset is imbalanced (e.g., more graduates than dropouts), so ROC-AUC is a better indicator.

#### 2. Feature Importance:

- **Admission Form:** Coefficients might reveal that "Entrance Exam / ENEM" students have higher odds of graduation than "SISU" students, possibly due to selection criteria.
- **Admission Period:** A positive coefficient could indicate that later periods (e.g., 2015+) have higher graduation rates, perhaps due to improved support systems.
- **Situation:** "Approved" status likely strongly predicts graduation (large positive coefficient), while "Failed" or "Failed Due to Absences" predicts dropout.

#### 3. Limitations:

- If dropout is rare (e.g., <20% of cases), the model might overpredict graduation, leading to lower recall for dropouts.
- Missing interactions (e.g., `admission_form × situation`) could reduce predictive power.

#### 4. Practical Implications:

- The model can identify at-risk students early (e.g., those with "Failed" status or from certain admission modes) for targeted interventions.
- Predictive performance suggests it's reliable for institutional planning but may need refinement for individual-level predictions.

### Step 6: Alternative Approach (Optional)

If higher accuracy is needed or the data exhibits non-linear patterns, a **Random Forest Classifier** could be used:

- **Advantages:** Captures complex interactions and handles imbalanced data better with class weighting.
- **Trade-off:** Less interpretable than logistic regression.

```
python
...
Copiar
```

```
from sklearn.ensemble import RandomForestClassifier

# Fit Random Forest
rf_model = RandomForestClassifier(n_estimators=100, random_state=42, class_weight='balanced')
rf_model.fit(X_train, y_train)

# Predict and evaluate
rf_pred = rf_model.predict(X_test)
print("Random Forest Accuracy:", accuracy_score(y_test, rf_pred))
print("Random Forest ROC-AUC:", roc_auc_score(y_test, rf_model.predict_proba(X_test)[:, 1]))
```

### Conclusion

The logistic regression model provides a solid baseline for predicting exit mode, with an expected ROC-AUC of 0.80-0.90, leveraging `admission_form`, `admission_period`, and `situation`. It's interpretable and actionable for educational institutions. For improved performance, a Random Forest could be explored, especially if the dataset is large and complex. The model's success hinges on data quality and the assumption that these features capture key drivers of student outcomes.



Entre para continuar a conversa