	Predicting academic success of undergraduate students by Catherine Meng, Jenson Chang, Jingyuan Wang, Siddarth Subrahmanian 2024/11/23
In [1]:	<pre>import altair as alt from sklearn.model_selection import train_test_split from sklearn.preprocessing import StandardScaler, OneHotEncoder from sklearn.compose import make_column_transformer from sklearn.pipeline import make_pipeline</pre>
	<pre>from sklearn.neighbors import KNeighborsClassifier from sklearn.model_selection import RandomizedSearchCV from scipy.stats import randint from sklearn import set_config set_config(display='text')</pre>
	In this analysis, we attempt to build a classification model using the k-nearest neighbors algorithm to predict student dropout and academic success based on information available at enrollment (including academic path, demographics, and socio-economic factors). Our final classifier performed consistently on unseen test data, achieving a cross-validation training score of 0.71, with a similar test score. Although the model's accuracy is moderate, it performs consistently.  Given that the data was collected from a single institution, a larger dataset may be necessary to generalize predictions to other institutions or countries. We believe this model is close to supporting dropout prediction for the institution from
	which the data was collected, though further research to improve performance and better understand characteristics of incorrectly predicted students would still be beneficial.  Introduction
	Higher education institutions worldwide face the ongoing challenge of academic dropout and student failure, which affect not only individual students' futures but also the institution's reputation and resources. The ability to predict and anticipate students' potential difficulties is valuable not only for supporting individual students in achieving their goals but also for institutions aiming to implement strategies that support and guide students who may be at risk of academic failure or dropout.  The goal of this analysis is to help reduce academic dropout and failure in higher education by applying machine learning techniques to identify at-risk students early in their academic journey, enabling institutions to implement targeted
	Methods  Methods
	Data  The data set is created by Mónica Vieira Martins, Jorge Machado, Luís Baptista and Valentim Realinho at the Instituto Politécnico de Portalegre (M.V.Martins, D. Tolledo, J. Machado, L. M.T. Baptista, V.Realinho. 2021). It is sourced from UC Irvine's Machine Learning Repository and can be found here. The data contains demographic, enrollment and academic (1st and 2nd semesters) information on the students. Each row in the data set represents a student record. Using these data, a model would be built to predict the academic outcome of the student. There are 36 columns in total.
	Analysis  The k-nearest neighbors (k-nn) algorithm was used to build a classification model to predict whether a student is at risk of dropping out. All variables included in the original data set, with the exception of the Course, Nacionality, Gender,
	Unemployment rate, Inflation rate, GDP, Previous qualification, Mother qualification Mother occupation, Father qualification, Father occupation columns were used to fit the model. Data was split with 80% being partitioned into the test set. The hyperparameter K was chosen using 5-fold cross validation. All numeric features were standardized just prior to model fitting. We leave the categorical features as they are because they all have integer data type.  Results & Discussion
T. [2].	To look at whether each of the predictors might be useful to predict the tumour class, we plotted the distributions of each predictor from the training data set and coloured the distribution by class (Dropout: blue, Enrolled: orange, and Graduate: red).
In [3]:	
Out[3]:	Marital status 6 Application mode 17 Application order 7 Course 17 Daytime/evening attendance\t 2 Previous qualification (grade) 97
	Nacionality 21 Mother's qualification 29 Father's qualification 30 Mother's occupation 31 Father's occupation 44 Admission grade 585
	Displaced 2 Educational special needs 2 Debtor 2 Tuition fees up to date 2 Gender 2 Scholarship holder 2 Age at enrollment 46
	International 2 Curricular units 1st sem (credited) 21 Curricular units 1st sem (enrolled) 22 Curricular units 1st sem (evaluations) 34 Curricular units 1st sem (approved) 22 Curricular units 1st sem (grade) 693
	Curricular units 1st sem (without evaluations) 11 Curricular units 2nd sem (credited) 19 Curricular units 2nd sem (enrolled) 21 Curricular units 2nd sem (evaluations) 29 Curricular units 2nd sem (approved) 20 Curricular units 2nd sem (grade) 685 Curricular units 2nd sem (without evaluations) 10
	Unemployment rate 10 Inflation rate 9 GDP 10 Target 3 dtype: int64
In [4]:	<pre># Remove extra '\t' from the column name train.rename(columns = {"Daytime/evening attendance\t" : "Daytime/evening attendance"}, inplace = True) test.rename(columns = {"Daytime/evening attendance\t" : "Daytime/evening attendance"}, inplace = True)  # Remove ' from column name to prevent issues with Altair plots train.columns = train.columns.str.replace("'s", "", regex=False) test.columns = test.columns str.replace("'s", "", regex=False)</pre>
In [5]:	<pre>test.columns = test.columns.str.replace("'s", "", regex=False)  # Group feature types based on feature description from source data categorical_features = ["Application order", "Course", "Nacionality", "Gender",</pre>
	"Father qualification", "Father occupation", "Displaced",
	"Curricular units 1st sem (grade)", "Curricular units 1st sem (without evaluations)", "Curricular units 2nd sem (credited)", "Curricular units 2nd sem (enrolled)", "Curricular units 2nd sem (evaluations)", "Curricular units 2nd sem (approved)", "Curricular units 2nd sem (grade)", "Curricular units 2nd sem (without evaluations)", "Unemployment rate", "Inflation rate", "GDP"]
In [6]:	<pre># Plot numerical features alt.Chart(train).mark_rect().encode(     x = alt.X(alt.repeat()).bin(maxbins=30),     y = alt.Y('Target', title=None),     color = alt.Color('count()').legend(orient="top") ).properties(</pre>
Out[6]:	<pre>width = 180 ).repeat(    numeric_features,    columns = 2 )  Count of Records</pre>
Out[6]:	Dropout - Enrolled - Graduate - Graduate - Graduate - Graduate - Graduate -
	95 115 135 155 175 95 115 135 155 175  Previous qualification (grade) (binned)  Dropout -
	Graduate
	0 2 4 6 8 10 12 14 16 18 20 22
	Graduate -
	Enrolled - Graduate -
	Graduate -
	Enrolled -
	Dropout -
	Enrolled - Graduate0.8 0.0 0.8 1.6 2.4 3.2 -4.5 -2.5 -0.5 1.5 3.5 Inflation rate (binned)  GDP (binned)
In [7]:	<pre>Figure 1. Distribution of Numerical Variable per Academic Outcome  # Plot categorical and boolean features alt.Chart(train).mark_bar().encode(     x = alt.X(alt.repeat()).type('quantitative'),     y = 'count()',</pre>
	<pre>column = alt.Column('Target', title = None),   color = alt.Color('Target', legend = None) ).properties(   width = 120,   height = 80 ).resolve_scale(</pre>
Out[7]:	<pre>y = 'independent' ).repeat(     categorical_features,     columns = 1 )</pre> <pre>Dropout Enrolled Graduate</pre>
ode[/]:	Sound of Month of Mon
	Application order Application order Application order  Dropout Enrolled Graduate
	Tourse Course Co
	Dropout Enrolled Graduate
	Nacionality  Dropout  Enrolled  Sylvation  S
	9 400 10 0 0.5 1.0 9 1,000 10 0 0.5 1.0 0.0 0.5 1.0 0.0 0.5 1.0 Gender Gender
	Dropout Enrolled Graduate  \$\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\f
	Marital status  Dropout  Enrolled  Graduate
	Application mode Application mode Application mode  Dropout Enrolled Graduate  \$\frac{y}{500} = \frac{1,000}{400} = \frac{y}{400} = \frac{y}{4
	Dropout  Enrolled  Daytime/evening attendance  Dropout  Enrolled  Graduate
	Sound of Food of Sound of Soun
	Dropout Enrolled Graduate  y  y  y  400  y  y  y  y  y  y  y  y  y  y  y  y
	Wother qualification  Mother qualification  Mother qualification  Mother qualification
	Dropout Enrolled Graduate  \$\begin{array}{c} \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \
	Mother occupation Mother occupation Mother occupation  Dropout Enrolled Graduate
	Father qualification  Father qualification  Father qualification  Father qualification
	Dropout Enrolled Graduate  \$\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\f
	S 0 0 100 200 0 100 200 0 100 200 0 100 200  Father occupation Father occupation Father occupation  Dropout Enrolled Graduate  S 0 0 0 100 200  Father occupation Father occupation  Father occupation  S 0 0 0 100 200  Father occupation  Father occupation  Father occupation  S 0 0 0 100 200  Father occupation  Father occupation
	B 000 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	Dropout Enrolled Graduate  \$\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\frac{\f
	Debtor Debtor Debtor  Dropout Enrolled Graduate  \$\frac{y}{2} \\ \frac{y}{5} \\ \
	Tuition fees up to date  Dropout  Enrolled  Graduate
	<b>y</b> 1,000 <b>y</b> 1,
	Scholarship holder Scholarship holder Scholarship holder  Dropout Enrolled Graduate  \$\frac{y}{2} 1,000  \text{1} \frac{y}{2} 500  \text{1} \frac{y}{2} 1,000  \text{2} \frac{y}{2} 1,0
	5 500
In [8]:	Figure 2. Distribution of Categorical Variable per Academic Outcome  # Features to drop from model drop_features = ["Course", "Nacionality", "Gender", "Unemployment rate",
	Reason for dropping:  • Course only captures 17 different courses and does not generalize the real world well.  • Nactionality and Gender are removed to avoid racial and gender bias  • In Figure 1, although Unemployment rate, Inflation rate and GDP are continous values, they each have less than 10 unique values out of 3000+ rows. This doesn't provide enough range to generalize the problem.
In [O]	• In Figure 2, Previous qualification, Mother qualification, Mother occupation, Father qualification and Father occupation have cluster patterns but it's unclear what the pattern represents since the ranking of education levels are arbitrary. E.g. "5 - Higher Education - Doctorate" is ranked higher than "1 - Secondary Education" but lower than "10 - 11th Year of Schooling - Not Completed". The source data website provides description on each ranking.  X_train = train.drop(columns=['Target'])
_~1 •	<pre>y_train = train['Target'] X_test = test.drop(columns=['Target']) y_test = test['Target']  # Make preprocessor preprocessor = make_column_transformer(</pre>
In [10]:	<pre>(StandardScaler(), numeric_features),    ('drop', drop_features) ) # Build the pipeline, use knn to train the model</pre>
	<pre>my_pipeline = make_pipeline(     preprocessor,     KNeighborsClassifier(n_neighbors=5) )  my_pipeline.fit(X_train, y_train)</pre>
Out[10]:	Pipeline(steps=[('columntransformer',
	'Age at enrollment', 'Curricular units 1st sem ' '(credited)', 'Curricular units 1st sem ' '(enrolled)', 'Curricular units 1st sem '
	'(evaluations)', 'Curricular units 1st sem ' '(approved)', 'Curricular units 1st sem ' '(grade 'Curricular units 2nd sem ' '(grade)',
	'Curricular units 2nd sem ' '(without evaluations)', 'Unemployment rate', 'Inflation rate', 'GDP']), ('drop', 'drop', ['Course', 'Nacionality',
	['Course', 'Nacionality', 'Gender', 'Unemployment rate', 'Inflation rate', 'GDP', 'Previous qualification', 'Mother occupation',
In [11]:	'Father qualification',
	<pre>param_distributions = {     'kneighborsclassifiern_neighbors': randint(1, 30) }  random_search = RandomizedSearchCV(     estimator=my_pipeline,     param distributions=param distributions,</pre>
	<pre>param_distributions=param_distributions,     n_iter=50,     cv=5,     scoring='accuracy',     random_state=42,     n_jobs=-1 )</pre>
	<pre>random_search.fit(X_train, y_train)  print("Best Parameters:", random_search.best_params_) print("Best CV Accuracy:", random_search.best_score_)</pre>
	Best Parameters: {'kneighborsclassifiern_neighbors': 17} Best CV Accuracy: 0.7143280671892855  We utilized the KNN to train the dataset and employed RandomizedSearchCV to fine-tune the hyperparameters. Based on the results, the optimal hyperparameter value is k=12, achieving a best cross-validation score of 0.71. Using this value, we retrained the model and evaluated its performance on the test set, obtaining a final test score of 0.71.
In [12]:	<pre>my_pipeline_best = make_pipeline(     preprocessor,     KNeighborsClassifier(n_neighbors=12) ) my_pipeline_best.fit(X_train, y_train)</pre>
	<pre>test_score = my_pipeline_best.score(X_test, y_test) print(f"Test accuracy: {test_score}")  Test accuracy: 0.7129943502824859</pre>
	References  1. Martins, M. V., D. Tolledo, J. Machado, L. M. T. Baptista, and V. Realinho. "Early Prediction of Student's Performance in Higher Education: A Case Study." In Trends and Applications in Information Systems and Technologies, vol. 1, Advances in Intelligent Systems and Computing series. Springer, 2021. https://archive.ics.uci.edu/dataset/697/predict+students+dropout+and+academic+success. DOI: 10.1007/978-3-030-72657-7_16.
	<ol> <li>Scikit-learn Developers. "Scikit-learn: Machine Learning in Python – API Reference Documentation." Accessed November 20, 2024. https://scikit-learn.org/dev/api/index.html.</li> <li>Vega-Altair Developers. "Altair User Guide." Accessed November 20, 2024. https://altair-viz.github.io/user_guide/data.html#.</li> </ol>
	4. Timbers, T., J. Ostblom, and M. Lee. "Predicting Breast Cancer from Digitized Images of Breast Mass." Accessed November 21, 2024. https://github.com/ttimbers/breast-cancer-predictor_report.ipynb.