Our final classifier per a larger dataset may better understand ch	attempt to build a classificerformed consistently on to be necessary to general naracteristics of incorrectly	unseen test data, achievi lize predictions to other in y predicted students woul	ng a cross-validation nstitutions or countrie ld still be beneficial.	training score of 0	0.71, with a similal model is close to	ar test score. Altho o supporting dropo	ugh the model's accur	acy is moderate, it	performs consistently. Go the data was collected,	Siven that the data we though further rese	as collected from a sir arch to improve perfori	ngle inst
difficulties is valuable. The goal of this analogous methods Data The data set is create. Repository and can be	titutions worldwide face the not only for supporting in a large state of the by Mónica Vieira Maribe found here. The data of the student. There are	individual students in ach ademic dropout and failur tins, Jorge Machado, Luís contains demographic, er	ieving their goals but re in higher education s Baptista and Valent	t also for institution n by applying mach	es aiming to implo nine learning tecl	ement strategies the	at support and guide sat-risk students early i	tudents who may their academic jo	be at risk of academic faurney, enabling institution	nilure or dropout. ons to implement target	geted support strategie	es. achine Lo
Analysis The k-nearest neight rate, GDP, Previous hyperparameter K was Results &	of the student. There are store (k-nn) algorithm was qualification, Mother qual as chosen using 5-fold crussion	s used to build a classifica ification Mother occupation oss validation. All numer	on, Father qualification	on, Father occupati dardized just prior	on columns were to model fitting.	e used to fit the mo	del. Data was split wit gorical features as the	h 80% being partit	ioned into the training se	et and 20% being pa	artitioned into the test s	
Args: file_path (st Returns: bool: True if """ # Check if fi if file_path. return Tr # Try to read try: pd.read_c return Tr except Except return Fa	d the file using panels of the file using pane	e. Eile, False otherwise	.	not a CSV file								
<pre>if check_csv(file print(f"{file else: print(f"{file/data/raw/data.c # Import data df = pd.read_csv(# Remove extra '\ df.rename(columns) # Remove ' from columns</pre>	e_path): e_path} is a CSV file e_path} is not a CSV esv is a CSV file. ('/data/raw/data.cs \t' from the column r s = {"Daytime/evening column name to prever columns.str.replace("	<pre>file.") sv', delimiter = ';') name g attendance\t" : "Da nt issues with Altain</pre>	aytime/evening att	tendance"}, inp	lace = True)\							
"Applicat	zion mode": pa.Columr 2, 5, 7, 10, 15, 16, 39, 42, 43, 44, 51, zion order": pa.Colum 1, 2, 3, 4, 5, 6, 9]) pa.Column (int, pa.	nullable=True), n(int, pa.Check.isin() 17, 18, 26, 53, 57])), nn(int, pa.Check.isin() 0, 9085, 9119, 9130, 9773, 9853, 9991]), pa.Column(int, pa.Check.isin() 10, 20, 20, 20, 20, 20, 20, 20, 20, 20, 2	(m) 9147, 9238, nullable=True), Check.isin(c.isin(39, 40, 42, 43])),								
"Nacional [1, 2] 100, "Mother of [1, 2] 34, "Father of [1, 2] 26, 43, "Mother of [0, 1] 132, 192, "Father of	lity": pa.Column(int, 2, 6, 11, 13, 14, 17, 101, 103, 105, 108, qualification": pa.Co 2, 3, 4, 5, 6, 9, 10, 35, 36, 37, 38, 39, qualification": pa.Co 2, 3, 4, 5, 6, 9, 10, 27,29, 30, 31, 33, 344]), nullable=True) occupation": pa.Colum 1, 2, 3, 4, 5, 6, 7, 134, 141, 143, 144, 193, 194]), nullable occupation": pa.Colum 1, 2, 3, 4, 5, 6, 7, 134, 141, 143, 144, 193, 194]), nullable occupation": pa.Colum 1, 2, 3, 4, 5, 6, 7, 134, 141, 143, 144, 193, 194]), nullable occupation": pa.Colum	21, 22, 24, 25, 26, 109]), nullable=Trublumn(int, pa.Check.id), 11, 12, 14, 18, 19, 40, 41, 42, 43, 44]) plumn(int, pa.Check.id), 11, 12, 13, 14, 18, 34, 35, 36, 37, 38, 36, 37, 38, 36, 37, 38, 36, 37, 38, 37, 38, 37, 38, 38, 38, 38, 38, 38, 38, 38, 38, 38	ne), isin(22, 26, 27, 29, n, nullable=True) isin(19, 20, 22, 25, 39, 40, 41, 42, n(22, 123, 125, 131, 173, 175, 191, n(,								
112, 144, 181, "Admission "Displace "Education "Debtor": "Tuition "Gender": "Scholars	1, 2, 3, 4, 5, 6, 7, 114, 121, 122, 123, 151, 152, 153, 154, 182, 183, 192, 193, on grade": pa.Column(int, pa.c	124, 131, 132, 134, 161, 163, 171, 172, 194, 195]), nullable (float, pa.Check.betwonullable=True), pa.Check.isin([0, 1]) pa.Column(int, pa.Check.isin([0, 1]), ra.Column(int, pa.Check.isin([0, 1]), ra.Column(int, pa.Check.isin([0, 1]), ra.mn(int, pa.Check.isin([0, 1])), ra.mn([0, 1])), ra.mn([0, 1]), ra.mn([0, 1])), ra.mn([0, 1])), ra.mn([0, 1])), ra.mn([0, 1]), ra.mn([0, 1])), ra.mn([0, 1])), ra.mn([0, 1])), ra.mn([0, 1])), ra.mn([0, 1])), ra.mn	135, 141, 143, 174, 175, le=True), ween(0, 200), 1, nullable=True) neck.isin([0, 1]), letrue), nullable=True), ck.isin([0, 1]), nullable=True), in([0, 1]),									
"Curricul "Curricul "Curricul "Curricul float "Curricul int, "Curricul	lar units 1st sem (er lar units 1st sem (gr t, pa.Check.between(Clar units 1st sem (wi nullable=True), lar units 2nd sem (cr nullable=True),	redited) ": pa.Column (nrolled) ": pa.Column (valuations) ": pa.Column (pproved) ": pa.Column (cade) ": pa.Column (d), 20), nullable=True ithout evaluations) ":	<pre>(int, nullable=True), (int, nullable=True), umn(int, nullable=True) (int, nullable=True),</pre>),								
"Curricul int, "Curricul int, "Curricul "Curricul float "Curricul int, "Unemploy "Inflatio "GDP": pa	lar units 2nd sem (er nullable=True), lar units 2nd sem (ev nullable=True), lar units 2nd sem (aplar units 2nd sem (grat, pa.Check.between(Clar units 2nd sem (winullable=True), yment rate": pa.Column (fa.Column(float, nullable=Column(float, nullable=Column(str, pa.Copout', 'Enrolled', 'encolumn', 'encolum	proved) ": pa.Column (umn((int, nullable=Tro e), : pa.Column(rue),	ue),								
pa.Check	ication Application Cour mode order	Daytime/evening attendance questions at the state of the	Previous qualifica ualification (gra	tion Nacionality (nde)	M other qualification qua	Father uni	ricular Curricular ts 2nd units 2nd sem sem edited) (enrolled)	nits 2nd sem (evaluations)	sem sem (grade) oproved) 0 0.000000	Curricular units 2nd sem Uner (without evaluations)	rate rate 10.8 1.4	GDP 1.74 D
1 1 2 1 3 1 4 2 4419 1 4420 1 4421 1	15 1 92 1 5 90 17 2 97 39 1 80 1 6 97 1 2 97 1 1 95 1 1 91	70 1 73 1 14 0 73 1 73 1 00 1	1 12 1 10 1 12 1 12 1 15	30.0 1 22.0 1 22.0 1 22.0 1 20.0 1 25.0 1 20.0 105 34.0 1 30.0 1	1 37 38 37 1 1 37 37	3 37 37 38 1 1 37	0 6 0 6 0 6 0 6 0 6 0 6 0 6	6 0 10 6 8 6 9	6 13.666667 0 0.000000 5 12.400000 6 13.000000 5 12.666667 2 11.000000 1 13.500000	0 0 0 0 0 0	10.8	0.79 Gra 1.74 D 3.12 Gra 0.79 Gra 4.06 Gra 2.02 D 0.79 D
#Check Target/res #Calculate observed_frequence print("Observed F # Calculate total total_students = # Define the expe	nns sponse variable follorved frequencies in the cies = df['Target'].vertequencies:\n", observed frequencies:\n", observed frequencies for ected frequencies for ected frequencies for extension of the context of the	73 1 Dows expected distribute 'Target' column value_counts() erved_frequencies) (observations) The a uniform distribute of the content of the counts of t	1 15	52.0 22	38	37 37	0 5 0 6	6	5 12.000000 6 13.000000	0		1.70 Gra
<pre>print("\nExpected from scipy.stats # Perform the Chi chi2_stat, p_valu print(f"\nChi2 St print(f"P-Value: Observed Frequenci Target Graduate 2209 Dropout 1421 Enrolled 794</pre>	<pre>i-Square goodness-of- ue = chisquare(observ tat: {chi2_stat}") {p_value}")</pre>	rm Distribution):", e	expected_frequenc:	ies)	encies)							
Chi2 Stat: 681.803 P-Value: 8.8773772 # Checking for ar chi2_results = {} features = ['Admi # Perform chi-squ for feature in feature in feature cont	Les (Uniform Distribut 33453887884 277253749e-149 momalous correlations } ission grade', 'Age a	s between target vari at enrollment', 'GDP' eature	iable and a subse			57]						
<pre># Perform chi chi2_results[# Check if the te for feature in ch np.testing.as # Checking for ar pearson_results = feature_pairs = 1 # Perform Pearsor for f1, f2 in feature</pre>	i-square test [feature] = chi2_cont est statistic is appr ni2_results: ssert_almost_equal(ch nomalous correlations = {} list(combinations(feature) ature_pairs:	roximately equal to pail no pa	_table).statistic	lues								
<pre># Check if the te for pair in pears np.testing.as</pre>	<pre>lts[f1+'_'+f2] = pear est statistic is appr son_results: ssert_almost_equal(pe ain_test_split(df, tr</pre>	roximately equal to pearson_results[pair],	ore-calculated va. pearson_validat									
Application mode	r attendance cation cation (grade) tion tion	17										
Father occupation Admission grade Displaced Educational spect Debtor Tuition fees up to Gender Scholarship holde Age at enrollment International Curricular units	ial needs to date er t 1st sem (credited)	44 585 2 2 2 2 2 2 46 2 21										
Curricular units Unemployment rate	1st sem (enrolled) 1st sem (evaluation: 1st sem (approved) 1st sem (grade) 1st sem (without evaluation: 2nd sem (credited) 2nd sem (enrolled) 2nd sem (evaluation: 2nd sem (approved) 2nd sem (grade) 2nd sem (without evaluation:	22 s) 34 22 693 aluations) 11 19 21 s) 29 20 685 aluations) 10										
<pre>Inflation rate GDP Target dtype: int64</pre> # Group feature to	types based on featur ures = ["Application "Marital stat "Previous qua "Father quali "Educational	9 10 3 re description from s	Nacionality", "Genoration", "Daytime/evenoration", qualification", poccupation", "Dispotor", "Tuition fe	ening attendanc "Mother occup placed",	ation",							
<pre># Plot numerical alt.Chart(train).</pre>	"Curricular uni "Curricular uni "Curricular uni "Curricular uni "Curricular uni "Unemployment r	its 1st sem (credited its 1st sem (evaluati its 1st sem (grade)", its 2nd sem (credited its 2nd sem (evaluati its 2nd sem (grade)", rate", "Inflation rat	d)", "Curricular to tons)", "Curricular unit "Curricular unit d)", "Curricular tons)", "Curricular unit "Curricular unit "Curricular unit	units 1st sem (ar units 1st se ts 1st sem (wit units 2nd sem (ar units 2nd se	enrolled)", m (approved)" hout evaluation enrolled)", m (approved)"	ons)",						
y = alt.Y('Ta				•••								
Dropout – Enrolled – Graduate – 16 24 32	35 155 175 cation (grade) (binned)	Dropout - Enrolled - Graduate - 0 2 4 6 8	5 155 175 grade (binned)									
Dropout – Enrolled – Graduate – 0 2 4 6 8 Curricular units 1s Dropout – Enrolled – Graduate –	10 12 14 16 18 20 22 st sem (enrolled) (binned)	Dropout – Enrolled – Graduate – 0 4 8 12 16 2 Curricular units 1st s Dropout – Enrolled – Graduate – Graduate –	sem (credited) (binned) 0 24 28 32 36 40 44 em (evaluations) (binned))								
Dropout - Enrolled - Graduate - 0.0 2.0 4.0	10 12 14 16 18 20 t sem (approved) (binned) 6.0 8.0 10.012.0 n (without evaluations) (binned)	Dropout – Enrolled – O 2 4 6 8 Curricular units 2nd Dropout – Enrolled – En	8 10 12 14 16 18 st sem (grade) (binned) 10 12 14 16 18 I sem (credited) (binned)									
Dropout – Enrolled – Graduate – 0 2 4 6 8	10 12 14 16 18 20 22 ad sem (enrolled) (binned) 8 10 12 14 16 1820 d sem (approved) (binned)	Dropout - Enrolled - Graduate - 0 2 4 6 8	16 20 24 28 32 sem (evaluations) (binned 10 12 14 16 18 ad sem (grade) (binned)	0								
Enrolled – Graduate – 0.0 2.0 4.0 Curricular units 2nd sem Dropout – Enrolled – Graduate – -0.8 0.0 0.8	6.0 8.0 10.012.0 (binned) 1.6 2.4 3.2 (binned)	Enrolled - Graduate - 7.5 9.5 11 Unemploym Dropout - Enrolled - Graduate4.5 -2.5 -	1.5 13.5 15.5 ent rate (binned)									
<pre># Plot categorica alt.Chart(train). x = alt.X(alt y = 'count()' column = alt.</pre>	al and boolean feature.mark_bar().encode(c.repeat()).type('qua',	antitative'), tle = None),										
).resolve_scale(y = 'independ').repeat(categorical_f columns = 1) Dropout	features,	Enrolled \$6.000 - 0.00	Graduate	••••								
Application of Dropout	order Applic	Count Count	2 4 6 Application order Graduate									
500 5,000 Course Dropout	Count of Records	000 10,000 ourse Enrolled 50 100 0	5,000 10,000 Course Graduate									
Nacionali Dropout	Sp 400 - 400 - 200	cionality Enrolled 0.5 1.0 Gender	Nacionality Graduate									
Dropout Sp 1,000 500 Marital state Dropout	Connt of Records turns Mari	Formula of the collection of t	Graduate 2 4 6 Marital status Graduate									
Application r			20 40 60 plication mode Graduate									
Dropout		0.5 1.0 0.5 ening attendance Dayti	0 0.5 1.0 ime/evening attendance									
Previous qualif	40 fication Previous	20 40 Count of	20 40 revious qualification Graduate									
Mother qualifi		0 40 0 0	20 40 her qualification Graduate									
Mother occup Dropout	pation Mother of	occupation Motorolled	100 200 ther occupation Graduate									
Pather qualified Dropout Dropo	Cation Father question 50 100 0 1	rolled Tolled	20 40 ner qualification Graduate 100 200 her occupation									
Propout Space of the propout of the	pation Father of the patient of the	Enrolled Sport 1,000 0.5 1.0 0.5 splaced	Graduate 0 0.5 1.0 Displaced									
Dropout 1,000 1,000 0.0 0.5 Educational spec	1.0 Sial needs Educations	Inrolled Spool 1,000 O.5 1.0 al special needs Enrolled	Graduate 0 0.5 1.0 cational special needs Graduate									
500 0.5 Debtor	Count of Records	Oos 1.00 Oo										
Dropout	1.0 0.0 Tuition for	o.5 1.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	0 0.5 1.0 uition fees up to date Graduate									
Scholarship h	nolder Schola	0.5 1.0 0.0	0 0.5 1.0 Scholarship holder Graduate									
	of Categorical Variable properties of model ["Course", "Nacionali" "Inflation rate", "G" "Mother occupation",		International mployment rate", ification", "Mothe		n",							
<pre># Features to dro drop_features = [</pre>	captures 17 different cours y and Gender are rem ough Unemployment r evious qualification g. "5 - Higher Education - drop(columns=['Target'] pp(columns=['Target']	oved to avoid racial and rate, Inflation rate on, Mother qualification Doctorate" is ranked high	gender bias e and GDP are con ation, Mother oc	ntinous values, they	her qualific	ation and Fath	er occupation hav	e cluster patterns	out it's unclear what the	pattern represents s	since the ranking of edu	ucation I
# Features to dropping: Course only of Nactionality In Figure 2, Property are arbitrary. E.g. X_train = train.org	op(columns=['Target'] arget'] sor ake_column_transformeler(), numeric_featur	er(ces),										
Figure 2. Distribution # Features to dro drop_features = [Course only of the Nactionality of the Nactio		,	rdScaler(),									
Figure 2. Distribution # Features to dro drop_features = [Course only of the image of the ima	<pre>ke_pipeline(, assifier(n_neighbors=</pre>	Standa: ['Prev: '(grad 'Admi:	ssion grade',	sem '								
Figure 2. Distribution # Features to dro drop_features = [Course only of the image of the ima	<pre>xe_pipeline(, assifier(n_neighbors= X_train, y_train) ('columntransformer',</pre>	Standa: ['Prev: '(grad' 'Admi: 'Age d' 'Curr: '(cred' 'Curr: '(enred' 'Curr: '(eval: 'Curr: '(app: 'Curr: '(grad' 'Curr:	de)', ssion grade', at enrollment', icular units 1st dited)', icular units 1st olled)', icular units 1st luations)', icular units 1st roved)', icular units 1st roved)', icular units 2st	sem ' sem '								
Figure 2. Distribution # Features to dro drop_features = [Course only of the image of the ima	<pre>xe_pipeline(, assifier(n_neighbors= X_train, y_train) ('columntransformer',</pre>	Standa: ['Prev. '(grad' 'Admi: 'Age a' 'Curr: '(cred' 'Curr: '(enr) '(eva. 'Curr: '(app: 'Curr: '(grad' 'Unemj' 'Infla' 'Prev: 'Mothe	de)', ssion grade', at enrollment', icular units 1st dited)', icular units 1st olled)', icular units 1st luations)', icular units 1st roved)', icular units 1st de icular units 2nd de)', icular units 2nd hout evaluations) ployment rate', ation rate', 'GDP , 'drop', se', 'Nacionality er', ployment rate', ation rate', 'GDP ious qualificatio er qualification'	sem ' sem ' sem ' sem ' sem ' ' sem ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' '								
# Use Randomizeds param_distribution # Use Randomizeds param_distribution # Use Randomizeds param_distribution # Build the pipel my_pipeline = mak preprocessor, KNeighborsCla) my_pipeline.fit() Pipeline(steps=[<pre>ce_pipeline(, assifier(n_neighbors= K_train, y_train) ('columntransformer', ColumnTransformer(t: ### ColumnTransformer to the column t</pre>	Standa: ['Prev: '(grade: 'Admi: 'Age a' 'Curr: '(crede: 'Curr: '(enred: 'Curr: '(eval: 'Curr: '(app: 'Curr: '(grade: 'Curr: '(de)', ssion grade', at enrollment', icular units 1st dited)', icular units 1st olled)', icular units 1st luations)', icular units 1st roved)', icular units 2st de icular units 2nd de)', icular units 2nd hout evaluations) ployment rate', ation rate', 'GDP , 'drop', se', 'Nacionality er', ployment rate', ation rate', 'GDP ious qualification' er occupation', er qualification'	sem ' sem ' sem ' sem ' sem ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' '								
# Use Randomizeds param_distribution # Was Randomizeds **In Figure 1, althorate train.com y_train = trai	<pre>ke_pipeline(</pre>	Standa: ['Prev. '(grad. 'Admi: 'Age a' 'Curr. '(cred. 'Curr. '(enr. '(eva. 'Curr. '(app. 'Curr. '(grad. 'Curr. '(grad. 'Curr. '(grad. 'Curr. '(with. 'Unemp. 'Infla. ('drop'. ['Cour. 'Gend. 'Unemp. 'Infla. 'Prev. 'Mothe. 'Mothe. 'Fathe. 'Fathe. 'Fathe. 'er', KNeighborsClass. 'erparameters cs': randint(1, 30)	de)', ssion grade', at enrollment', icular units 1st dited)', icular units 1st olled)', icular units 1st luations)', icular units 1st roved)', icular units 2st de icular units 2nd de)', icular units 2nd hout evaluations) ployment rate', ation rate', 'GDP , 'drop', se', 'Nacionality er', ployment rate', ation rate', 'GDP ious qualification' er occupation', er qualification'	sem ' sem ' sem ' sem ' sem ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' '								

Predicting academic success of undergraduate students