

# Big Data Assignment #4

## Multinomial Logistic Regression: Student Grade



# Multinomial Logistic Regression

---

Regresi logistik merupakan salah satu metode yang dapat digunakan untuk mencari hubungan variabel respon yang bersifat dichotomous (berskala nominal atau ordinal dengan dua kategori) atau polychotomous (mempunyai skala nominal atau ordinal dengan lebih dari dua kategori) dengan satu atau lebih variabel prediktor dan variabel respon bersifat kontinyu atau kategorik. Regresi logistik multinomial atau disebut juga model logit politomus adalah model regresi yang digunakan untuk menyelesaikan kasus regresi dengan variabel dependen berupa data kualitatif berbentuk multinomial (lebih dari dua kategori) dengan satu atau lebih variabel independen.





# Student Grade Dataset

## Info

Dataset berisi sekumpulan data terkait kebiasaan dari siswa sekolah menengah atas dengan mata pelajaran Matematika dan hasil nilai akhir mata pelajaran dimaksud.

## Sumber Dataset :

<https://www.kaggle.com/uciml/student-alcohol-consumption>



# Business Understanding

Ruang Lingkup Bisnis : Education

Suatu sekolah menengah atas akan melakukan perhitungan atau prediksi terkait kelulusan nilai yang akan didapatkan oleh siswanya. Dengan memperhatikan variabel-variabel yang dapat mempengaruhi nilai kelulusan, terutama pengaruh penggunaan alkohol di kehidupan sehari-hari. Output dari perhitungan ini akan menghasilkan formula untuk dapat memprediksi kelulusan dari siswa pada sekolah dimaksud.



# Data Understanding

Field Name	Description	Data Type
school	<i>Nama sekolah</i>	object
sex	<i>Jenis kelamin</i>	object
age	<i>Usia</i>	int64
address	<i>Alamat siswa</i>	object
famsize	<i>Jumlah saudara</i>	object
Pstatus	<i>Hubungan kedua orang tua</i>	object
Medu	<i>Pendidikan ibu</i>	int64
Fedu	<i>Pendidikan ayah</i>	int64
Mjob	<i>Pekerjaan ibu</i>	object
Fjob	<i>Pekerjaan ayah</i>	object



# Data Understanding

Field Name	Description	Data Type
reason	<i>Alasan memilih sekolah</i>	object
guardian	<i>Wali siswa</i>	object
traveltime	<i>Waktu tempuh dari rumah ke sekolah</i>	int64
studytime	<i>Waktu belajar mingguan</i>	int64
failures	<i>Jumlah mengalami kegagalan</i>	int64
schoolsup	<i>Support pendidikan tambahan</i>	object
famsup	<i>Support keluarga</i>	object
paid	<i>Biaya ekstra</i>	object
activities	<i>Mengikuti ekstrakurikuler</i>	object
nursery	<i>Mengikuti sekolah keperawatan</i>	object



# Data Understanding

Field Name	Description	Data Type
higher	<i>Keinginan untuk mengambil pendidikan tinggi</i>	object
internet	<i>Akses internet dirumah</i>	object
romantic	<i>Memiliki hubungan romantis</i>	int64
famrel	<i>Kualitas hubungan keluarga</i>	int64
freetime	<i>Waktu luang setelah sekolah</i>	int64
goout	<i>Menghabiskan waktu bersama teman</i>	int64
Dalc	<i>Konsumsi alkohol harian</i>	int64
Walc	<i>Konsumsi alkohol mingguan</i>	int64
health	<i>Kondisi kesehatan</i>	int64
absences	<i>Jumlah absen</i>	int64

# Data Understanding

Field Name	Description	Data Type
G1	Nilai grade 1	int64
G2	Nilai grade 2	int64
G3	Nilai grade 3 - Kelulusan	int64





# Data Preparation

## ◆ Melakukan pengecekan data

Data dilakukan pengecekan terhadap informasi seperti jumlah data (baris) dan kolom, tipe data dari setiap kolom, data yang memuat nilai null.

## ◆ Menghapus kolom dan mengecek deskripsi data

Penghapusan kolom dilakukan terhadap kolom dengan tipe data object. Dari total 32 kolom variabel independen, menjadi hanya 8 kolom variabel independen. Pengecekan deskripsi data dilakukan dengan menggunakan fungsi describe(). Hasil yang ditampilkan jumlah data, nilai min, max, standar deviasi, median, mean, kuartil bawah, kuartil tengah, dan kuartil atas.

## ◆ Melakukan normalisasi

Normalisasi dilakukan dengan menggunakan fungsi StandardScaler(), dimana nilai *mean* dijadikan 0 dan *variance* 1. Normalisasi dilakukan terhadap seluruh data independen variabel.

# Modeling (Training)

---

## ✓ Mengecek data setiap kategori

Pengecekan terhadap koefisien dan intercept yang didapatkan dengan menggunakan menggunakan fungsi MNLogit().

## ✓ Menghitung Nilai Intercept

Menghitung nilai intercept dengan menggunakan fungsi `LogisticRegression()`. Pada proses ini menghasilkan nilai intercept setiap kategori yang akan digunakan dalam perhitungan.

## ✓ Menghitung Nilai Koefisien

Nilai koefisien didapatkan dengan fungsi `coef_` pada pemodelan `LogisticRegression()` dengan parameter `multi_class='multinomial'`. Fungsi ini menghasilkan nilai koefisien pada setiap variabel independen pada setiap kategori.





# Evaluation (Testing and Accuracy)

---

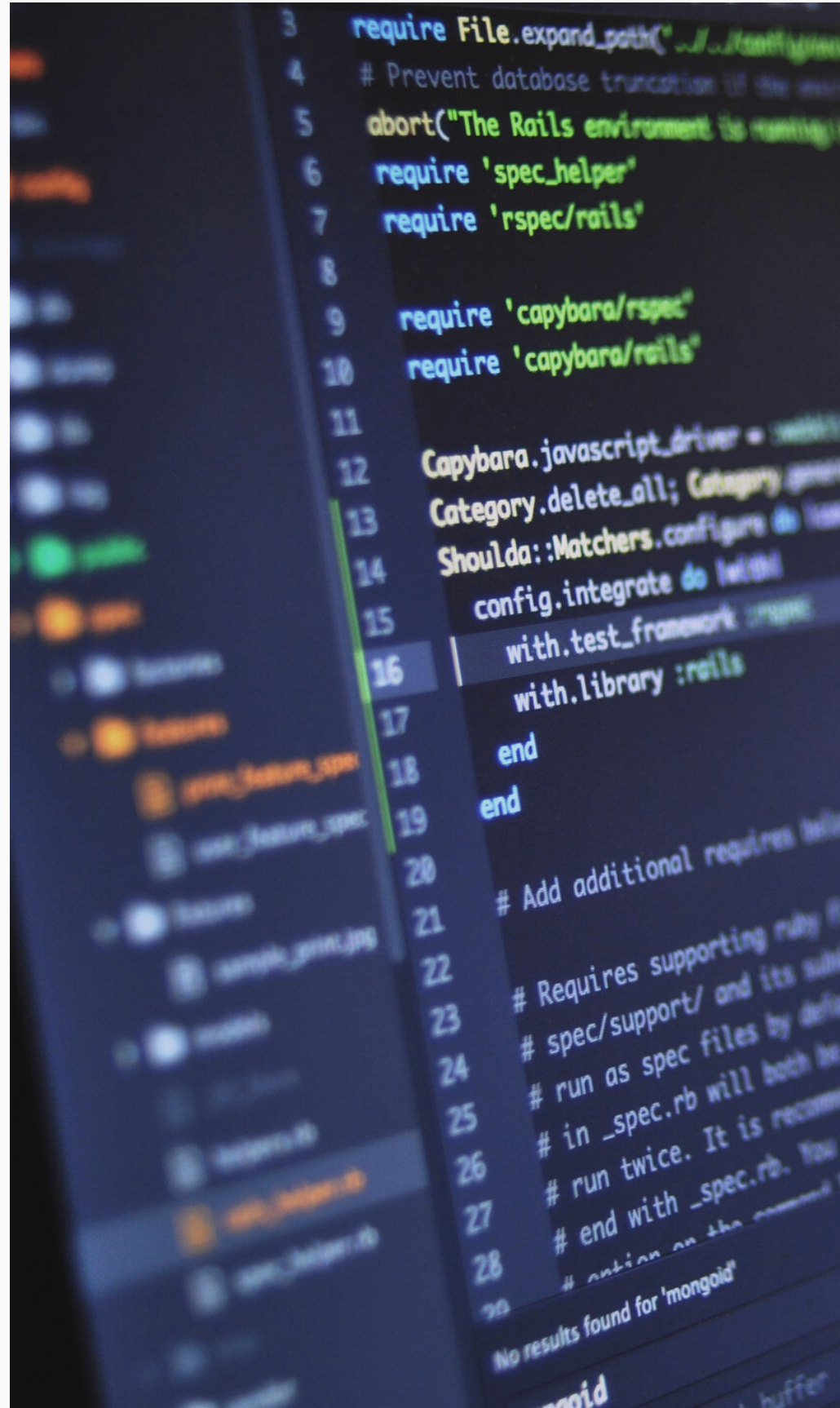
## Pemodelan Manual

Pada tahap ini, pemodelan dilakukan secara manual dengan menggunakan rumus multinomial logistic regression, dimana dalam perhitungan ini membutuhkan nilai intercept, nilai koefisien variabel independen, dan nilai setiap variabel independen.

## Penghitungan Akurasi

Pengecekan akurasi dengan menggunakan fungsi `classification_report` dan menghasilkan nilai akurasi 0.81. Nilai ini menunjukkan bahwa akurasi pemodelan yang telah dibuat sebesar 81%.





# Deployment

Tahap ini dilakukan menggunakan bahasa pemrograman Python dengan IDE Kaggle.

Code secara keseluruhan terdapat pada link Github :

<https://github.com/dyanaagustina/Learn-BigData/tree/Big-Data/Multinomial%20LR>

# Thank You

---

