

# Laporan Proyek UTS - IBDA3111

Calvin Institute of Technology

Semester ganjil 2022/2023



Oleh

Victor Chendra / 202000338 / IT & Big Data Analytics

## Students Performance in Exams Dataset

### A. Latar belakang

Dalam mata kuliah ini, yaitu *IBDA3111 Prapemrosesan & Rekayasa data*, kami diberikan sebuah proyek UTS dengan mencari dataset sendiri untuk diolah lalu dianalisa atau melakukan pemodelan.

Adapun dataset yang saya pilih untuk proyek ini adalah dataset yang berjudul Students Performance in Exams Dataset yang diambil dari situs [kaggle.com](https://www.kaggle.com/datasets/whenamancodes/students-performance-in-exams).

Berikut adalah tautan dataset:

<https://www.kaggle.com/datasets/whenamancodes/students-performance-in-exams>

*Students Performance in Exams* dataset adalah sebuah dataset yang menggambarkan tentang peforma ujian siswa. Dalam dataset berisikan tentang nilai-nilai ujian siswa di sekolah. Adapun tabel-tabel yang terdapat di dalam dataset.

gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score
string	string	string	string	string	integers	integers	integers
male	group A	high school	standard	completed	67	67	63
female	group D	some college	standard	none	59	60	50
female	group C	high school	standard	none	63	73	66
female	group C	bachelor's degree	standard	none	51	64	63
male	group D	high school	free/reduced	completed	82	83	77

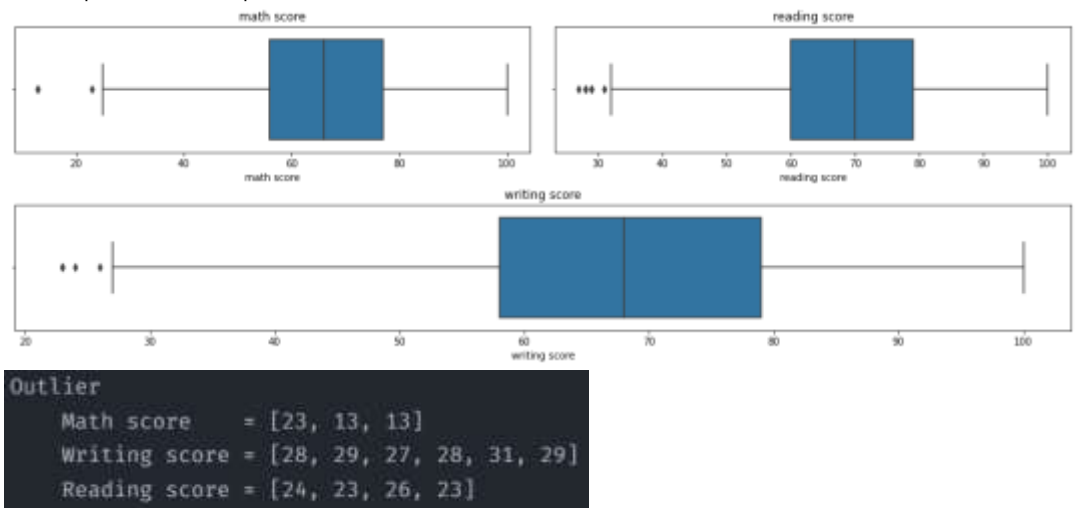
Terdapat 8 kolom, 5 string 3 integers dan 1000 baris data. Dataset ini dapat digunakan untuk analisa statistika maupun membangun pemodelan untuk menebak nilai ujian seorang siswa. Dataset ini juga memberikan gambaran, apakah pengaruh latar belakang pendidikan orang tua, *test preparation course*, dan lain-lain terhadap kinerja ujian siswa.

Adapun masalah yang ingin saya ingin selesaikan secara *general* adalah:

- How effective is the test preparation course?
- Which major factors contribute to test outcomes?
- What would be the best way to improve student scores on each test?
- What patterns and interactions in the data can you find?

## B. Data

- a. Pemrosesan data yang saya lakukan dalam proyek ini yaitu, 2 metode *data cleaning* dan 2 metode *feature selection* yang mana merupakan kriteria yang telah ditentukan dalam proyek ini. Dua metode *data cleaning* tersebut adalah:
- Delete duplicated rows*:  
Terdapat 1 baris data yang terduplikasi
  - Identify outlier and imputation the outlier*:  
Terdapat 13 data pencilan dalam dataset



Data pencilan tersebut, diimputasi ke batas bawah mereka masing-masing.

Kemudian untuk 2 metode feature selection:

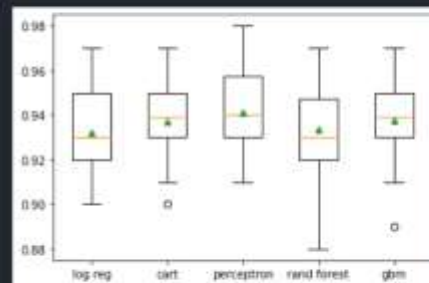
- RFE
  - Chi-square
- b. Pendekatan yang saya lakukan dalam proyek ini adalah analisis statistika dan visualisasi. Alasan saya memilih analisis statistika untuk proyek/dataset ini adalah karena menurut saya analisa statistika lebih menyenangkan dan lebih menunjukkan sebuah hasil/gambaran dataset tersebut. Walaupun sebenarnya bisa melakukan pemodelan yaitu menebak/prediksi atau klasifikasi (*pass or fail*), tetapi kurang tepat saja dilakukan pada dataset ini dibandingkan dengan analisis statistika.

## C. Solusi

- a. Langkah-langkah yang saya lakukan dalam mengerjakan proyek ini:
- Import library yang dibutuhkan
  - Load and check dataset

- iii. Data cleaning → 1. Delete duplicated rows → 2. Identify outlier & impute to it's lower bound value
- iv. Feature selection (dataset after data cleaning)
  1. (RFE) Add new column "average" → one hot encoding the input variables → change the output variables to PassStatus (if the value is  $\geq 60$ , it's mean pass or 1 else 0 fail) → then we select the best RFE estimator from 5 RFE models estimator (LogReg, Cart, Perceptron, RanFor, GBM) it's Perceptron.

```
> log reg 0.932 (0.019)
> cart 0.937 (0.017)
> perceptron 0.941 (0.018)
> rand forest 0.934 (0.019)
> gbm 0.938 (0.019)
```



!! As we can see.. "perceptron" has the highest accuracy, so we take "perceptron" model as our estimator

Selected features: [0, 1, 4, 5, 6, 7, 10, 12, 13, 14, 15, 16, 17, 18, 19]

Selected column according RFE model Perceptron

gender female	gender male	race/ethnicity group C	race/ethnicity group D	race/ethnicity group E	parental level of education associate's degree	parental level of education master's degree	parental level of education some high school
0	0	1	0	0	0	0	0
1	1	0	0	1	0	0	0

lunch free/reduced	lunch standard	test preparation course completed	test preparation course none	math PassStatus	read PassStatus	write PassStatus	average
0	1	1	0	1	1	1	1
1	0	0	1	0	0	0	0

2. Chi-square  
Chi-square output is the score of each input variables. → Sort it but keep the original index → then display the final dataset after features selection chi-square (best of 15 out of 20 features)

Selected column according chi-square score (best of 15 out of 20 features)

parental level of education some college	race/ethnicity group	parental level of education master's degree	parental level of education bachelor's degree	parental level of education associate's degree	race/ethnicity group	race/ethnicity group	test preparation course note
0	0	0	0	0	0	0	0
1	0	0	0	0	1	0	1
2	1	1	0	0	0	0	1

parental level of education some high school	test preparation course completed	lunch_standard	lunch_free/reduced	read_PassStatus	math_PassStatus	write_PassStatus	average
0	1	1	0	1	1	1	1
1	0	0	1	0	0	0	0
0	0	0	1	1	0	0	0

- v. Answering the problem define in description  
Untuk menjawab masalahnya, pertama kita melakukan visualisasi dataset untuk diamati lalu mengambil kesimpulan.

#### D. Kesimpulan

- Dataset diperoleh dari [kaggle.com](https://www.kaggle.com/datasets/whenamancodes/students-performance-in-exams)  
<https://www.kaggle.com/datasets/whenamancodes/students-performance-in-exams>
- Cara memperoleh dataset: mengunduh dari website
- Terdapat 2 metode *cleaning data*: Delete duplicated rows, Identify outlier & Imputation.
- Terdapat 2 metode feature selection: RFE, Chi-square
- Pendekatan analisis statistika
- Hasil observasi dataset:
  - Students with a higher parental's degrees have more chances to achieve a higher scores
  - Based on gender, female have higher score average than male
  - The females dominate the higher scores in reading and writing
  - Male outperform female in ONLY math
  - Along this dataset female and male students are equally distributed
  - "Lunch" is quite influential in performing student's exams score