

LAPORAN TENGAH/AKHIR
MAGANG & STUDI INDEPENDEN BERSERTIFIKAT
EDUCATION: STUDENT PERFORMANCE
Di PT Digitalisasi Pemuda Indonesia

Diajukan untuk memenuhi persyaratan kelulusan

Program MSIB MBKM

oleh :

FORDINAND HALOMOAN PASARIBU / 188160061



TEKNIK INFORMATIKA
UNIVERSITAS MEDAN AREA

2021 / 2022

Lembar Pengesahan Teknik Informatika
EDUCATION: STUDENT PERFORMANCE
Di PT Digitalisasi Pemuda Indonesia

oleh :

FORDINAND HALOMOAN PASARIBU / 188160061

disetujui dan disahkan sebagai

Laporan Magang atau Studi Independen Bersertifikat Kampus Merdeka

Medan, 17 Januari 2022

Pembimbing Magang atau Studi Independen Teknik Informatika Universitas
Medan Area.



Rizki Muliono, S.Kom, M.Kom

NIP: 0109038902

Lembar Pengesahan

EDUCATION: STUDENT PERFORMANCE

Di PT Digitalisasi Pemuda Indonesia

oleh :

FORDINAND HALOMOAN PASARIBU / 188160061

disetujui dan disahkan sebagai

Laporan Magang atau Studi Independen Bersertifikat Kampus Merdeka

Bandung, 17 Januari 2022

Mentor



Ignasius Frans D.S.T.N, S.Kom

Abstraksi

Laporan ini bertujuan untuk memahami nilai ujian siswa yang berpengaruh pada Jenis Kelamin, Suku/Ras, Tingkat Pendidikan Orang Tua, Makan Siang, dan Kursus persiapan ujian. Kemudian Membuat model *Machine Learning* untuk memprediksi nilai siswa Analisis ini menggunakan variabel independen yaitu Jenis Kelamin, Suku/Ras, Tingkat Pendidikan Orang Tua, Makan Siang, dan Kursus persiapan ujian. Variabel dependennya yaitu keterangan lulus atau tidak lulus. Data yang digunakan merupakan data yang sudah di sediakan di situs Kaggle. Adapun algoritma yang diterapkan ialah *KNearest Neighbors* (KNN), *Support Vector Machine* (SVM), *Multilayer Perceptron* (MLP). Alur pemodelannya yaitu: Pengumpulan Data, *Exploratory Data Analysis* (EDA), *Data Preprocessing*, dan Evaluasi Model. Hasil dari yang sudah dilakukan oleh penulis menunjukkan bahwa penerapan algoritma SVM memiliki akurasi sebesar 98%, dan variabel seperti makan siang sangat berpengaruh terhadap nilai ujian siswa.

Kata Kunci: Nilai ujian, *Machine Learning*, KNN, SVM, MLP

Kata Pengantar

Puji dan syukur saya sampaikan kepada Tuhan Yang Maha Esa karena telah memberikan rahmat dan karunia-Nya kepada saya, sehingga berhasil menyelesaikan laporan ini tepat pada waktunya tentang Tugas Proyek Tengah/Akhir Magang dan Studi Independen Bersertifikat (MSIB) dengan judul *Education: Student Performance*. Laporan ini adalah salah satu syarat kelulusan peserta MSIB dalam program Kampus Merdeka yang dibuat oleh Kementerian Pendidikan, Kebudayaan, Riset, dan Teknologi Republik Indonesia. Saya telah banyak memperoleh bantuan dan bimbingan dari berbagai pihak, maka pada kesempatan ini saya menyampaikan ucapan terima kasih kepada :

1. Kak Aditya Soleh, Selaku Officer at Digital Skola.
2. Kak Agil Haykal, Kak Romansya Setyo, dan Pak Ganjar Alfian , Selaku Tutor Digital Skola.
3. Kak Shavira Tisna dan Kak Qoriyana Nurselvi, Selaku Representative at Digital Skola.
4. Kepada mentor kelompok Fantastic Five yaitu Kak Ignasius Frans, dan teman-teman: Denta Winardi, Rahmatul Fajri, Nur Afni, dan Zulva Amaliya.

Terima kasih atas bantuan, bimbingan, dukungan, fasilitas yang telah diberikan kepada saya. Semoga Tuhan Yang Maha Esa membalas semua kebaikan yang telah diberikan kepada saya. Saya menyadari bahwa laporan ini jauh dari kata kesempurnaan, oleh karena itu kritikan dan saran yang bersifat membangun selalu saya harapkan demi kesempurnaan laporan ini. Akhir kata saya sampaikan terima kasih kepada pembaca yang sekiranya telah meluangkan waktunya untuk membaca laporan ini seraya memajukan selangkah lagi pengetahuan tentang isi laporan ini.

Medan, 17 Januari 2022

Fordinand Halomoan Pasaribu

NIM: 188160061

Daftar Isi

Lembar Pengesahan Teknik Informatika	i
Lembar Pengesahan	ii
Abstraksi	iii
Kata Pengantar	iv
Daftar Isi.....	v
Daftar Gambar.....	vi
Bab I Pendahuluan.....	1
I.1. Latar belakang.....	1
I.2. Lingkup	2
I.3. Tujuan	3
Bab II Digital Skola.....	4
II.1. Struktur Organisasi	4
II.2. Lingkup Pekerjaan	6
II.3. Deskripsi Pekerjaan	6
II.4. Jadwal Magang dan Studi Independen Bersertifikat	10
Bab III <i>Education: Student Performance</i>	11
III.1. Deskripsi/Konteks.....	11
III.2. Proses Model <i>Machine Learning</i>	11
III.3. Rekomendasi Model <i>Machine Learning</i>	21
Bab IV Penutup.....	22
IV.1. Kesimpulan	22
IV.2. Saran	22
References	vii
Lampiran A. TOR Digital Skola	A-1
Lampiran B. Log Activity	B-1
Lampiran C. Dokumen Teknik	C-1

Daftar Gambar

Gambar 1 Tampilan dataset student performance.....	11
Gambar 2 Persentase data per kolom	12
Gambar 3 Race/Ethnicity berdasarkan Gender	13
Gambar 4 Parental level of education berdasarkan Gender	14
Gambar 5 Race/Ethnicity berdasarkan Nilai	14
Gambar 6 Penyebaran nilai berdasarkan gender	15
Gambar 7 Preprocessing	16
Gambar 8 Sebelum feature engineering	17
Gambar 9 Feature Engineering	17
Gambar 10 Label Encoder	18
Gambar 11 Feature Scaling	19
Gambar 12 Korelasi antar kolom	21

Bab I Pendahuluan

I.1. Latar belakang

Dalam rangka menyiapkan mahasiswa menghadapi perubahan sosial, budaya, dunia kerja dan kemajuan teknologi yang pesat, kompetensi mahasiswa harus disiapkan untuk lebih gayut dengan kebutuhan zaman. *Link and match* tidak saja dengan dunia industri dan dunia kerja tetapi juga dengan masa depan yang berubah dengan cepat. Perguruan Tinggi dituntut untuk dapat merancang dan melaksanakan proses pembelajaran yang inovatif agar mahasiswa dapat meraih capaian pembelajaran mencakup aspek sikap, pengetahuan, dan keterampilan secara optimal dan selalu relevan.

Kampus Merdeka diharapkan dapat menjadi jawaban atas tuntutan tersebut. Kampus Merdeka merupakan wujud pembelajaran di perguruan tinggi yang otonom dan fleksibel sehingga tercipta kultur belajar yang inovatif, tidak mengekang, dan sesuai dengan kebutuhan mahasiswa. Proses pembelajaran dalam Kampus Merdeka merupakan salah satu perwujudan pembelajaran yang berpusat pada mahasiswa (*student centered learning*) yang sangat esensial. Pembelajaran dalam Kampus Merdeka memberikan tantangan dan kesempatan untuk pengembangan inovasi, kreativitas, kapasitas, kepribadian, dan kebutuhan mahasiswa, serta mengembangkan kemandirian dalam mencari dan menemukan pengetahuan melalui kenyataan dan dinamika lapangan seperti persyaratan kemampuan, permasalahan riil, interaksi sosial, kolaborasi, manajemen diri, tuntutan kinerja, target dan pencapaiannya.

Digital Skola merupakan perusahaan yang baru didirikan oleh tiga anak bangsa yang bertujuan mendorong pertumbuhan talenta digital di Indonesia dengan menghadirkan pusat edukasi nonformal secara daring mengenai kecakapan digital yaitu *data scientist*, *data engineer* dan *digital marketing*. Digital Skola berkomitmen mewujudkan pelatihan *skill digital* yang inklusif dengan biaya terjangkau, waktu fleksibel, pengajar kredibel dan kurikulum sesuai kebutuhan industri. Dalam pelaksanaan pelatihannya, Digital Skola bekerja sama dengan tutor-tutor yang merupakan praktisi di bidang masing-masing.

Dari segi penyusunan kurikulum serta metode belajar, Digital Skola memberikan jaminan bahwa semua peserta dari latar belakang apapun akan mampu menguasai ilmu yang diajarkan. Kurikulum pun disusun dengan cermat sesuai kebutuhan industri agar lulusan pelatihan dari Digital Skola menguasai kemampuan untuk langsung dapat bekerja. Berikut adalah visi & misi Digital Skola :

1. Visi: Mengembangkan wawasan dan kemampuan angkatan kerja di Indonesia untuk menjadi generasi.
2. Misi: Mengadakan pelatihan informal keahlian digital dengan biaya terjangkau, waktu fleksibel, pengajar kredibel dan kurikulum sesuai kebutuhan industri.

I.2. Lingkup

- Nama Kegiatan.
- Latar Belakang.
- Tujuan Kegiatan.
- Jadwal dan Waktu Pelaksanaan.
- Peserta.
- Penyelenggara.

I.3. Tujuan

Data Science meliputi pembelajaran individu dan implementasi proyek dalam bentuk tim. Pada pembelajaran individu, setiap peserta akan mengikuti kelas secara *interactive live-online meeting* dimana peserta akan mendengar secara langsung serta dapat berkonsultasi dengan *expert* serta pembimbing terkait materi yang dipelajarinya. Hasil yang akan diperoleh, yaitu:

1. Mampu memahami dasar-dasar statistika yang digunakan pada domain *data science*.
2. Mampu menguasai bahasa pemrograman Python termasuk *libraries* yang digunakan.
3. Mampu mengerjakan *dataset* yang diberikan serta membuat *modelling* nya.
4. Mampu membuat model *machine learning* sesuai kebutuhan.
5. Mampu membuat visualisasi data secara efektif.
6. Memahami materi yang terkait dengan *Analytics*, yang meliputi metodologi *data science*, Numpy, *dataframe*, statistik, visualisasi data, dan *business intelligence*.

Bab II Digital Skola

II.1. Struktur Organisasi

Tim Digital Skola

1. Aditya Soleh (Chief Executive Officer).
2. Zico Alaia Akbar (Chief Operating Officer).
3. Stephanie Octavia (Chief Business Development Officer).
4. Indah Salimin (Senior Copywriter).
5. Ismaya R. P. (Graphic Designer).
6. Qoriyana Nurselvi (Product Officer).
7. Neti Rosmayanti (Customer Relations Officer).
8. Faliha Ishma (Product Officer).
9. Siti Meiranti Nabilah (Account Executive).
10. Adellia Anggun Trisnawati (Customer Relations Officer).
11. Shavira Tisna (Representative).

Tutor Digital Skola

1. Agil Haykal (Analytics).
2. Mohammad Aprialdi (Data Scientist).
3. Faldi Sulistiawan (Data Scientist).
4. Farhan Reza (Analytics).
5. Muhammad Farid (Data Analyst).
6. Romansya Setyo (Business Intelligence Analyst).
7. Calvin Nixcon (Data Analyst).

8. Addinul Masri (Data Analyst).
9. Dr. Eng. Ganjar Alfian (Assistant Professor at one of the campus in Korea).
10. Ari Sulistiyo Prabowo (Senior Data Analyst).
11. Anastasia Yoan (Digital Marketing Specialist).
12. Claudia Natasya (Senior Sales & Marketing).
13. Gabriella Ovina (Head of Digital Minshare).
14. Irene Bergosa (Internet Marketing Lead).
15. Jovanka Audria (Content Specialist).
16. Laurensia Octavia (Founder Go Digital Consultant).
17. Deasy Natalia (Chief Operations Officer).
18. Salman Al Farisiy (Head of Marketing Communications).
19. Sherly Fusin (Senior Digital & Media Planning Manager).
20. Jaka Darmawan (SEO Specialist).
21. Nana Kusuma (Senior SEM).
22. Shona Firdaus (Performance Marketing Specialist).
23. Muhammad Haekal S.A (Senior Digital Marketing).
24. Ali Firdaus Chifari (Data Engineering).
25. Dimasta Juniar (Data Engineering).
26. Rio Harapan P (Big Data Engineering).
27. Rizki Dermawan (Data Engineering).
28. Ignes Nathania Widjaja (Product Designer).
29. Fahmi Kusuma Aji (UX Designer).
30. Areta Selena Khrista (UX Designer).

31. Mohamad Arief Bagusprastyo (UI/UX Designer).
32. Hary Nugraha (Software Engineering).
33. Dedi Ananto (Full-Stack Web Developer).

II.2. Lingkup Pekerjaan

Tim Digital Skola

1. Aditya Soleh (Chief Executive Officer).
2. Zico Alaia Akbar (Chief Operating Officer).
3. Qoriyana Nurselvi (Product Officer).
4. Shavira Tisna (Representative).

Tutor Program MSIB Data Science Digital Skola

1. Agil Haykal (Analytics).
2. Farhan Reza (Analytics).
3. Romansya Setyo (Business Intelligence Analyst).
4. Dr. Eng. Ganjar Alfian (Assistant Professor at one of the campus in Korea).

II.3. Deskripsi Pekerjaan

2.3.1. Nama Kegiatan

Dataset & Modelling Project

2.3.2. Latar Belakang

Machine Learning adalah metode yang digunakan untuk membuat program yang bisa belajar dari data. Berbeda dengan program komputer biasa yang statis, program *machine learning* adalah program yang dirancang untuk mampu belajar sendiri. Cara belajar program *machine learning* mengikuti cara belajar manusia, yakni belajar dari contoh-contoh. *Machine learning* akan mempelajari pola dari contoh-contoh yang dianalisa, untuk menentukan jawaban dari pertanyaan-pertanyaan berikutnya. Memang tidak semua masalah bisa dipecahkan dengan program *machine learning*. Namun, seringkali algoritma yang sifatnya kompleks, ternyata bisa dipecahkan dengan sangat simpel oleh *machine learning*. Adapun alur kerja *Machine Learning* yaitu:

1. Perencanaan Model dan Pengumpulan data

Tentukan tujuan dari pembuatan *machine learning* dengan penguraian masalah yang ingin diselesaikan. Kemudian kumpulkan dataset yang relevan dengan permasalahan dan tujuan pemodelan. *Dataset* akan memberikan kontribusi besar terhadap akurasi model yang akan dibuat.

2. Persiapan Data

- *Training* data yaitu data yang akan digunakan mesin untuk belajar.
- *Testing* data yaitu data yang digunakan untuk prediksi.
- *Validation* data yaitu data bentuk validasi model hasil belajar.

3. Training Model

Gunakan algoritma yang sesuai dengan kebutuhan dan masalah yang ingin diselesaikan. Cari model yang terbaik dengan menggunakan *training* data. Pastikan fitur-fitur yang digunakan sebagai indikator sesuai dengan tujuan pemodelan.

4. Evaluasi Model

Gunakan metrik evaluasi untuk mengukur performa model yang dihasilkan. Kemudian, gunakan *validation* data untuk merepresentasikan performa model saat digunakan pada data yang belum dipelajari, sehingga hasilnya dapat digunakan untuk proses *tuning*.

5. Parameter Tuning

Lakukan *tuning* parameter untuk meningkatkan performa model yang dihasilkan. Contoh parameter yang dapat disesuaikan: jumlah tahapan *training*, *learning rate*, dan lain-lain. Kemudian, lakukan evaluasi kembali pada model yang sudah di *tuning*.

6. Lakukan Prediksi

Ketika model yang diharapkan sudah sesuai dengan kebutuhan maka Langkah selanjutnya adalah melakukan uji prediksi menggunakan *testing* data. Tujuannya adalah mendapatkan gambaran yang lebih baik mengenai performa model ketika digunakan untuk memprediksi data di masa depan.

2.3.3. Tujuan Kegiatan

Adapun tujuan kegiatan *dataset & modelling project* untuk :

1. Prediksi yaitu melakukan prediksi berupa nilai, probabilitas maupun data dan kemudian merekomendasikan hasilnya untuk digunakan sebagai alat bantu pengambil keputusan maupun secara langsung digunakan secara otomatis oleh sistem.
2. Deskripsi yaitu menampilkan pola data untuk dianalisa dan penemuan masalah.

2.3.4. Jadwal dan Waktu Pelaksanaan

Tanggal : 23 Agustus 2021 – 28 Januari 2022

Jadwal Kelas :

- Selasa: jam 19.15 – 21.15 WIB
- Kamis: jam 19.15 – 21.15 WIB
- Sabtu: jam 13.00 – 15.00 WIB

Tempat : Online Via Zoom.

2.3.5. Peserta

Jumlah peserta kegiatan kelas SIB B sekitar 49 orang, dengan didominasi oleh mahasiswa dari program studi Informatika, Sistem Informasi, Ilmu Komputer, Geofisika, Geologi Elektro, Matematika, Akuntansi, Statistika, dan Manajemen.

2.3.6. Penyelenggara

Penyelenggara dari kegiatan ini adalah PT Digitalisasi Pemuda Indonesia (Digital Skola).

II.4. Jadwal Magang dan Studi Independen Bersertifikat

Sesi pembelajaran online sebanyak 60.

Pembelajaran dan sesi bimbingan sebanyak 600 jam.

Tanggal : 23 Agustus 2021 – 28 Januari 2022

Jadwal Kelas	
SIB A	SIB B
Senin : 19.15 - 21.15 WIB	Selasa: 19.15 - 21.15 WIB
Rabu : 19.15 - 21.15 WIB	Kamis: 19.15 - 21.15 WIB
Jumat : 19.15 - 21.15 WIB	Sabtu: 13.00 - 15.00 WIB

Tempat : Online Via Zoom.

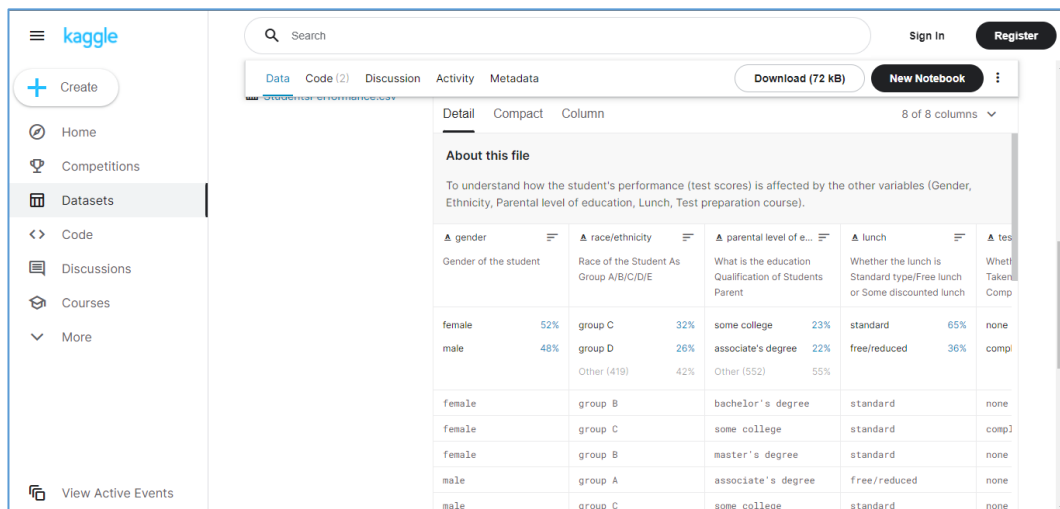
Bab III *Education: Student Performance*

III.1. Deskripsi/Konteks

Bagaimana memahami kinerja siswa(nilai ujian) yang dipengaruhi oleh variabel lain (Jenis Kelamin, Suku/Ras, Tingkat Pendidikan Orang Tua, Makan Siang, dan Kursus persiapan ujian). Adapun nilai-nilai ujian yang diprediksi dapat dipengaruhi oleh variabel lain meliputi Nilai Matematika, Nilai Membaca, dan Nilai Menulis. Kemudian Membuat beberapa model *machine learning* untuk memprediksi nilai siswa. Adapun dataset yang penulis gunakan dalam MSIB ini adalah berasal dari situs <https://www.kaggle.com/adithyabshetty100/student-performance>

III.2. Proses Model *Machine Learning*

3.2.1 Pengumpulan Data



The image shows a screenshot of the Kaggle website interface for the 'student performance' dataset. The page includes a search bar, navigation tabs (Data, Code, Discussion, Activity, Metadata), and a 'Download (72 KB)' button. The main content area displays a table with columns for gender, race/ethnicity, parental level of education, lunch type, and test scores. A summary table is visible below the main data table.

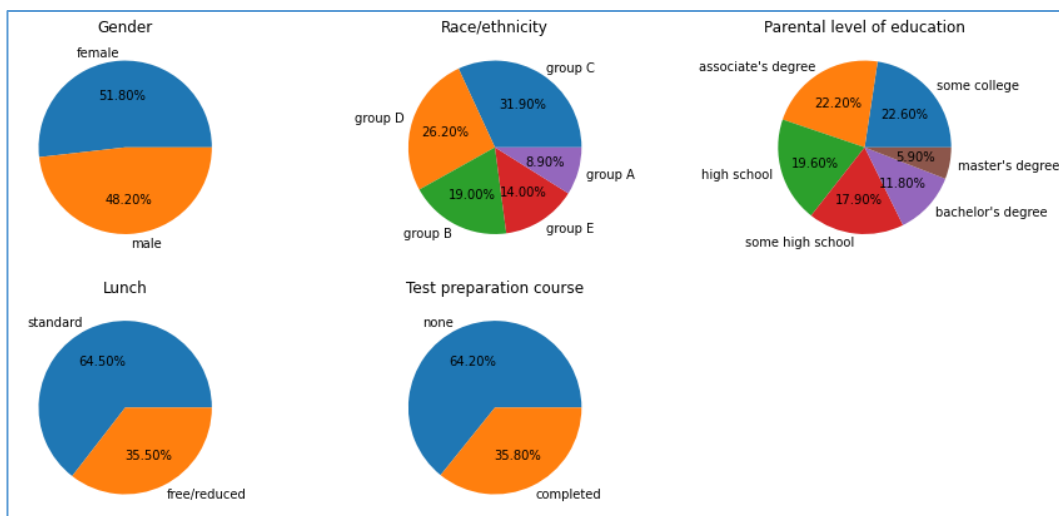
gender	race/ethnicity	parental level of e...	lunch	tes
female	group C	some college	standard	none
male	group D	associate's degree	free/reduced	compi
	Other (419)	Other (552)		
female	group B	bachelor's degree	standard	none
female	group C	some college	standard	compi
female	group B	master's degree	standard	none
male	group A	associate's degree	free/reduced	none
male	group C	some college	standard	none

Gambar 1 Tampilan *dataset student performance*

Langkah pertama dalam pembuatan *machine learning* yaitu mengumpulkan data. Semakin banyak dan semakin baik kualitas data yang kita punya, performa *machine learning* yang kita buat akan semakin baik. Ada beberapa metode dalam mengumpulkan data seperti *web scrapping*, *data mining*, atau dari situs website yang sudah di sediakan seperti *kaggle.com*, *data.go.id* dan lain sebagainya. Adapun data yang digunakan dalam *project* MSIB ialah melalui situs *kaggle*.

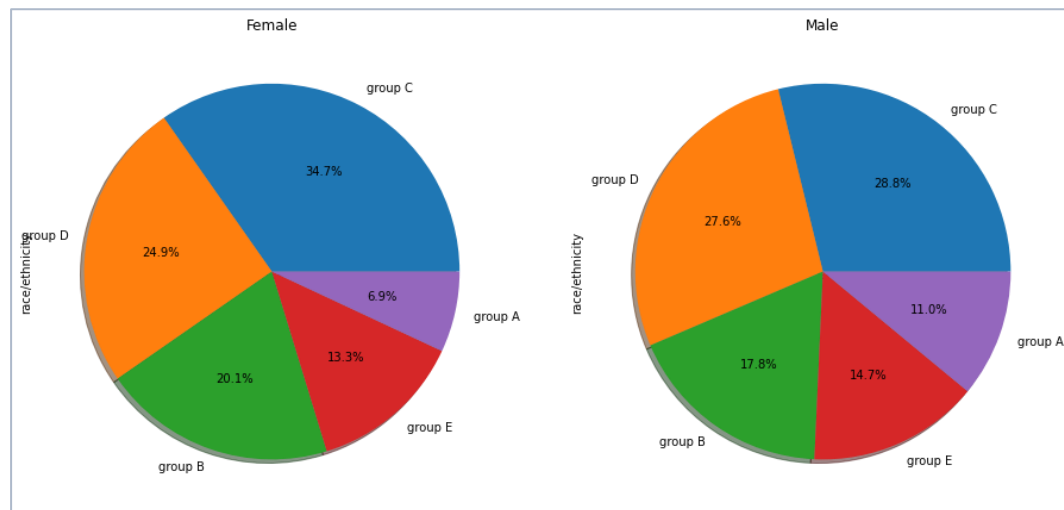
3.2.2 Exploratory Data Analysis (EDA)

EDA atau dikenal pula dengan analisis data eksploratif merupakan pendekatan analisis untuk suatu data guna membuat gambaran keseluruhan (*summary*) data sehingga mudah untuk dipahami. EDA memungkinkan analisis memahami isi data yang digunakan, mulai dari distribusi, frekuensi, korelasi dan lainnya. Dalam prakteknya, *curiosity* sangat penting dalam proses ini, pemahaman konteks data juga diperhatikan, karena akan menjawab masalah masalah dasar.



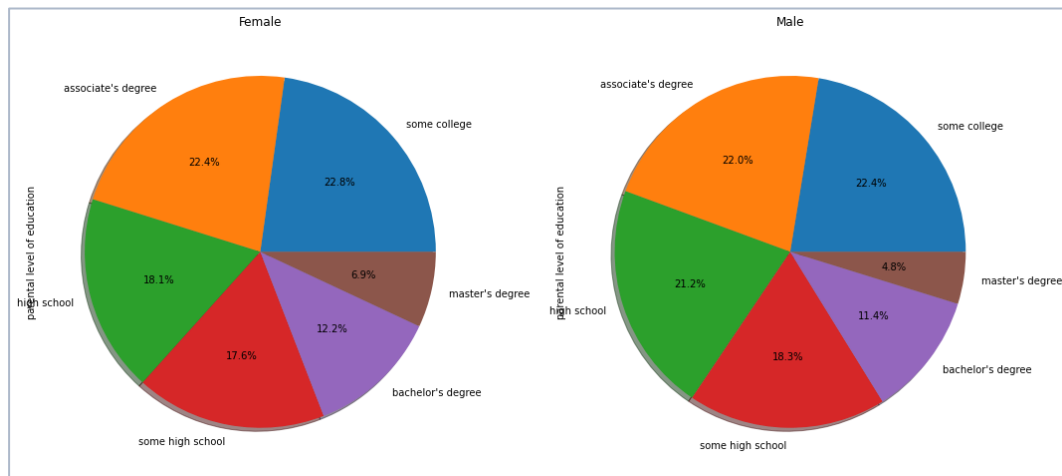
Gambar 2 Persentase data per kolom

Pada gambar 2 kita dapat melihat selisih perbandingan antar data pada tiap – tiap kolom, yaitu pada kolom *gender* didominasi oleh *female*, kolom *race/ethnicity* didominasi oleh *group C*, kolom *parental level of education* didominasi oleh *Some College* dan memiliki selisih 0,40% dengan *associate's degree*, kemudian kolom *lunch* didominasi oleh *standard*, dan terakhir kolom *test preparation course* didominasi oleh *none*(tidak mengikuti kursus).



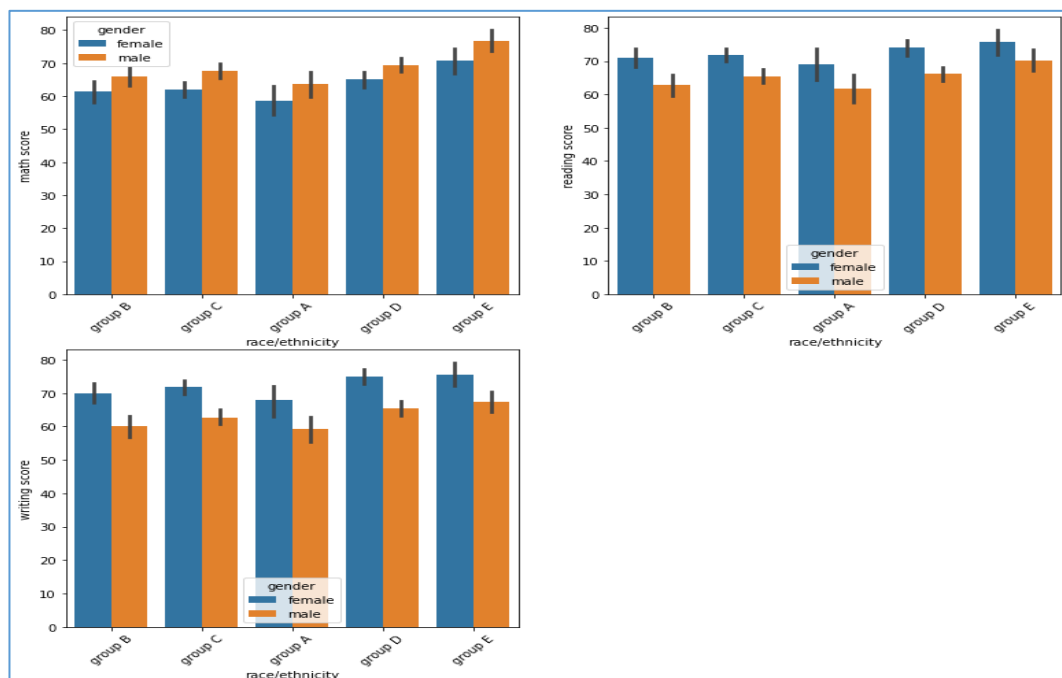
Gambar 3 *Race/Ethnicity* berdasarkan *Gender*

Pada gambar 3 bertujuan untuk membandingkan persentase *gender female* dan *male* terhadap *race/ethnicity*, dan hasil menunjukkan bahwa *group C* pada *gender female* lebih unggul sebesar 34,7%, kemudian pada *group D* terhadap *gender male* lebih unggul sebesar 27,6%, dan terakhir *group B* didominasi oleh *gender female* kembali dengan sebesar 20,1%.



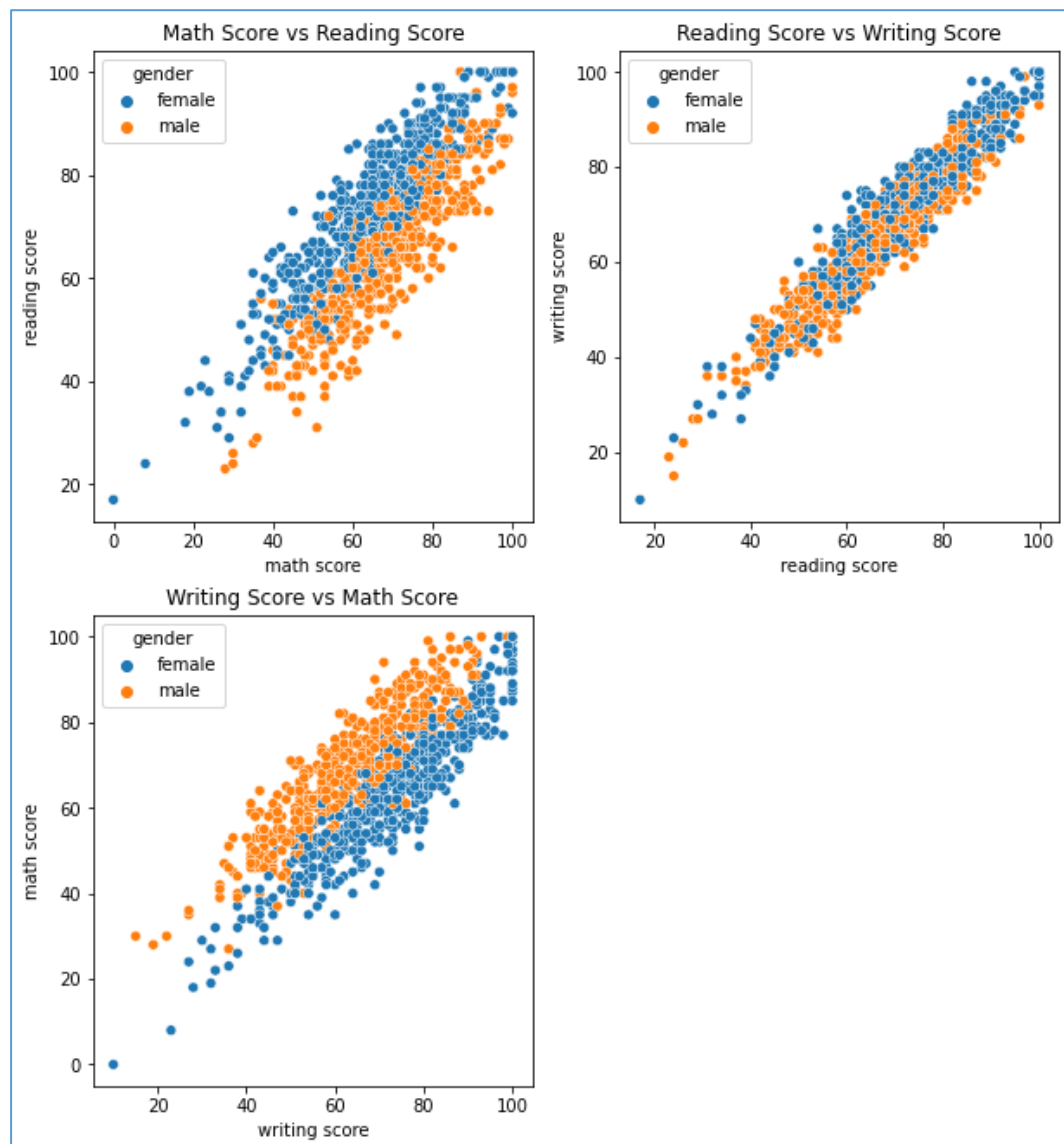
Gambar 4 Parental level of education berdasarkan Gender

Pada gambar 4 terlihat bahwa perbandingan *female* dengan *male* pada *associate's degree* memiliki selisih sebesar 0,4%, kemudian pada *Some College* juga selisih sebesar 0,4%, dan terakhir *high school* pada *male* sebesar 21,2%.



Gambar 5 Race/Ethnicity berdasarkan Nilai

Gambar 5: dapat kita lihat bahwa grup dengan *score math* paling tinggi dimiliki oleh *Group E*, kemudian dengan *score reading* tertinggi dimiliki oleh *Group D* dan *Group E*, dan terakhir pada *score writing* tertinggi dimiliki oleh *Group D* dan *Group E*.



Gambar 6 Penyebaran nilai berdasarkan *gender*

Gambar 6 yaitu variabel *male* dan *female* pada *math score* saling berdekatan(tidak beraturan), variabel *male* sedikit dominan. Kemudian variabel *male* dan *female* pada *reading score* saling menindih sehingga variabel *male* menjadi sangat dominan. Dan variabel *male* dan *female* pada *writing score* saling berdekatan(beraturan), 2 variabel seimbang.

3.2.3 Data Preprocessing

Data preprocessing adalah teknik awal *data mining* untuk mengubah *raw data*(data mentah) menjadi format dan informasi yang lebih efisien dan bermanfaat. Format pada *raw data* yang diambil dari berbagai macam sumber seringkali mengalami error, *missing value*, dan tidak konsisten. Sehingga, perlu dilakukan pembenahan format agar hasil data mining tepat dan akurat. Pada kasus *student performance* tidak memiliki data yang kosong, *error*, data terduplikat, dan lain sebagainya.

```
# Cek Apakah ada data yang missing value
df.isnull().sum()

gender                0
race/ethnicity        0
parental level of education  0
lunch                 0
test preparation course  0
math score            0
reading score         0
writing score         0
dtype: int64

[ ] df.duplicated().sum()

0
```

Gambar 7 Preprocessing

Feature Engineering adalah bagaimana kita menggunakan pengetahuan kita dalam memilih *features* atau membuat *features* baru agar model *machine learning* dapat bekerja lebih akurat dalam memecahkan masalah.

math score	reading score	writing score
72	72	74
69	90	88
90	95	93
47	57	44
76	78	75

Gambar 8 Sebelum *feature engineering*

Pada kasus proyek diatas penulis menambah kolom seperti *total score*, *average score*, dan yang terakhir adalah keterangan yang berdasarkan *average score* yang akan menentukan lulus atau tidaknya, berikut dibawah adalah gambar sesudah melakukan *feature engineering*.

math score	reading score	writing score	total_score	avg_nilai	keterangan
72	72	74	218	72.666667	lulus
69	90	88	247	82.333333	lulus
90	95	93	278	92.666667	lulus
47	57	44	148	49.333333	tidak lulus
76	78	75	229	76.333333	lulus

Gambar 9 *Feature Engineering*

Label Encoder digunakan untuk mengonversi data kategorikal, atau data teks, menjadi angka, yang dapat lebih dipahami oleh model. Berikut adalah hasil setelah di *encoder*.

		Test Preparation course	Hasil Encoder
		Completed	0
		None	1
		Lunch	Hasil Encoder
		Free/Reduced	0
		Standard	1
		Parental level of education	Hasil Encoder
		Associate's Degree	0
		Bachelor's Degree	1
		High School	2
		Master's Degree	3
		Some College	4
		Some High School	5
Gender	Hasil Encoder		
Female	0		
Male	1		
Keterangan	Hasil Encoder		
Lulus	0		
Tidak Lulus	1		
Race/Ethnicity	Hasil Encoder		
Group A	0		
Group B	1		
Group C	2		
Group D	3		
Group E	4		

Gambar 10 *Label Encoder*

Train-Test Split adalah membagi dataset menjadi *train set* dan *test set*, atau dengan kata lain, data yang digunakan untuk proses *training* dan *testing* merupakan kumpulan data yang berbeda. Dalam membagi *test size* pada umumnya sebesar 20% - 30%. *Feature Scaling* adalah suatu cara untuk membuat *numerical data* pada dataset memiliki rentang nilai (*scale*) yang sama. Tidak ada lagi satu variabel data yang mendominasi variabel data lainnya. Adapun alasan penulis memilih *standarization* dikarenakan cocok pada data yang sifatnya terdistribusi normal.

Standardisation	Normalisation
$x_{\text{stand}} = \frac{x - \text{mean}(x)}{\text{standard deviation}(x)}$	$x_{\text{norm}} = \frac{x - \min(x)}{\max(x) - \min(x)}$

Gambar 11 *Feature Scaling*

3.2.4. *Evaluation and Recommendation Model*

Machine learning model adalah algoritma *machine learning* yang sebelumnya telah dilakukan proses pelatihan/*training* dengan data latih tertentu sehingga dia siap digunakan untuk melakukan prediksi terhadap data baru. Adapun algoritma yang diterapkan ialah *KNearest Neighbors* (KNN) adalah algoritma yang berfungsi untuk melakukan klasifikasi suatu data berdasarkan data pembelajaran (*train data sets*), yang diambil dari k tetangga terdekatnya (*nearest neighbors*). Dengan k merupakan banyaknya tetangga terdekat. Lalu yang kedua adalah algoritma *Support Vector Machine* (SVM) digunakan untuk mencari *hyperplane* terbaik dengan memaksimalkan jarak antar kelas. Dan yang terakhir adalah algoritma *Multilayer Perceptron* (MLP) adalah salah satu permodelan dalam teknologi jaringan saraf tiruan (JST) dengan karakteristik memiliki nilai bobot yang lebih baik dari pada pemodelan yang lain, sehingga menghasilkan klasifikasi yang lebih akurat. Berikut adalah hasil setelah melakukan pemodelan dengan beberapa algoritma.

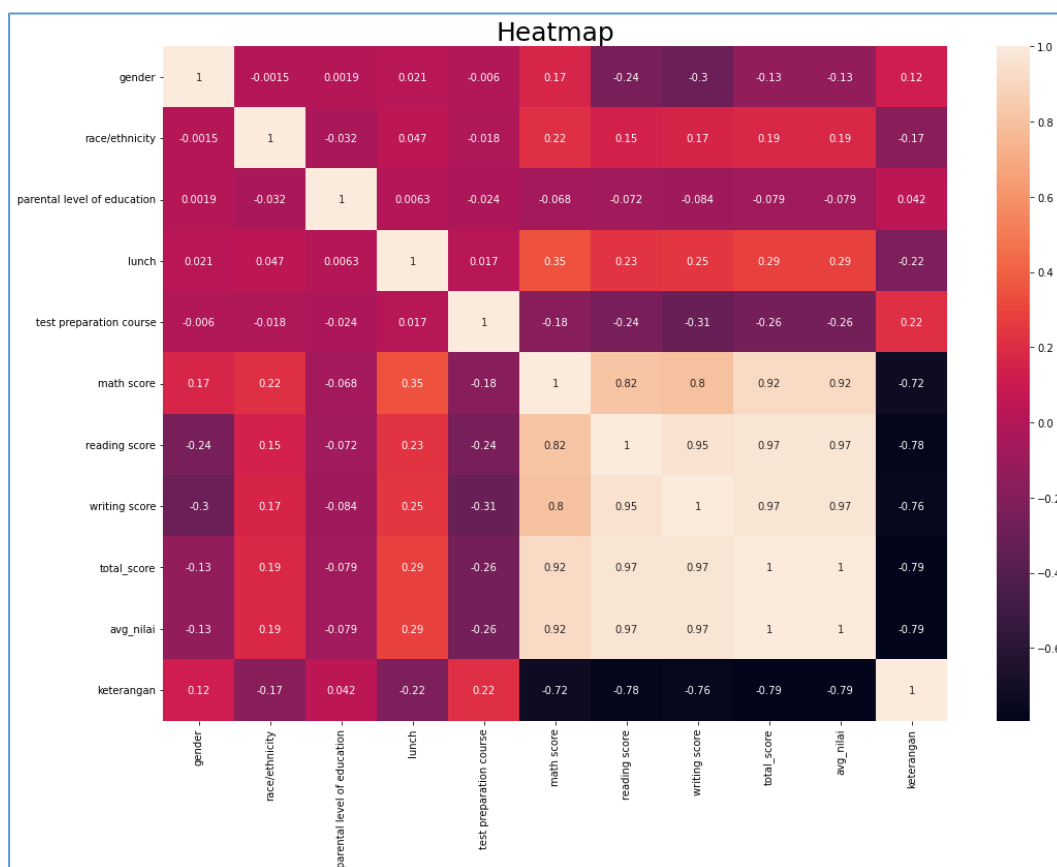
<i>Train-Test Spli</i>	KNN	SVM	MLP
<i>TRAINING</i>	<i>Accuracy : 97.38 % Precision : 97.43 % Recall : 97.66 %</i>	<i>Accuracy : 98.62 % Precision : 97.71 % Recall : 99.77 %</i>	<i>Accuracy : 99.88 % Precision : 99.77 % Recall : 100.0 %</i>
<i>TESTING</i>	<i>Accuracy : 92 % Precision : 92.66 % Recall : 92.66 %</i>	<i>Accuracy : 98.0 % Precision : 98.17 % Recall : 98.17 %</i>	<i>Accuracy : 99.5 % Precision : 100.0 % Recall : 99.08 %</i>

Kemudian untuk mengotomatisasi/memaksimalkan kinerja model maka penulis menggunakan *machine learning pipeline*, dan berikut adalah hasilnya.

KNN	SVM	MLP
<i>Accuracy : 92. % Precision : 92.66 % Recall : 92.66 %</i>	<i>Accuracy : 98.0 % Precision : 98.17 % Recall : 98.17 %</i>	<i>Accuracy : 99 % Precision : 99 % Recall : 99 %</i>

III.3. Rekomendasi Model *Machine Learning*

Berdasarkan hasil penerapan pemodelan *machine learning* yang sudah penulis lakukan maka merekomendasikan model *Support vektor machine* (SVM), karena perbandingan nilai *accuracy* untuk evaluasi model dan pemodelan *pipeline* tidak berubah atau tidak mengalami perubahan yang sangat jauh artinya model tersebut bisa diterapkan.



Gambar 12 Korelasi antar kolom

Kemudian berdasarkan *heatmap* diatas, maka korelasi *lunch* dan *math score* memiliki nilai paling tinggi dibandingkan variabel lain, dalam arti nilai matematika juga dipengaruhi oleh makan siang.

Bab IV Penutup

IV.1. Kesimpulan

Kesimpulan yang dapat diambil dari pelaksanaan MSIB di PT Digitalisasi Pemuda Indonesia atau Digital Skola antara lain:

1. Mampu beradaptasi dengan baik sehingga terjalin hubungan antar peserta dan tim Digital Skola.
2. Tim Digital Skola mengadakan 2 proyek utama yaitu prediksi klasifikasi dan prediksi *time series forecasting*.
3. Hasil yang telah dilakukan oleh penulis pada studi kasus *education student performance* bisa diterapkan bagi siswa yang ingin memprediksi apakah nilai berpengaruh pada variabel lain, misal suku/ras, jenis kelamin, pendidikan orang tua, dan lain sebagainya.

IV.2. Saran

Adapun saran dari pelaksanaan MSIB di PT Digitalisasi Pemuda Indonesia atau Digital Skola antara lain:

1. Untuk pengerjaan proyek akhir penulis berharap diadakannya proyek klasifikasi berupa text ataupun gambar.
2. Perlu adanya perbaikan pemodelan *machine learning* pada studi kasus diatas, misalnya menambah *feature* baru agar model bisa banyak belajar pada data baru. Kemudian bisa menambah *hyperparameter tuning* agar kinerja model lebih maksimal.
3. Tentunya penulis sudah menyadari jika dalam penyusunan laporan ini masih banyak ada kesalahan serta jauh dari kata sempurna. Adapun nantinya penulis akan segera melakukan perbaikan susunan laporan ini dengan menggunakan pedoman dari beberapa sumber dan kritik yang bisa membangun dari para pembaca.

References

- [1] C. Indonesia, *Kampus Merdeka ala Mas Nadiem Makarim di Era Revolusi Industri 4.0*, 2020.
- [2] C. Darujati and A. B. Gumelar, *Pemanfaatan Teknik Supervised Untuk Klasifikasi Teks Bahasa Indonesia*, pp. 1-5, 2012.
- [3] “Direktorat Jenderal Pendidikan Tinggi Kementerian Pendidikan dan Kebudayaan,” *Buku Panduan Merdeka Belajar - Kampus Merdeka*, 2020.
- [4] “Kementerian Pendidikan dan Kebudayaan,” *Merdeka Belajar: Kampus Merdeka*, 2020.
- [5] G. Kosala, A. Harjoko and S. Hartati, “Proceedings of the International Conference on Video and Image Processing,” *License Plate Detection Based on Convolutional Neural Network: Support Vector Machine (CNN-SVM)*, pp. 1-5, 2017.
- [6] Kusri and E. T. Luthfi, *Algoritma Data Mining*, 2009.
- [7] H. Lubaba, *Merdeka dalam Berpikir*, 2020.
- [8] Y. Makdori, *Kemdikbud Rilis Konsep Kampus Merdeka untuk Perguruan Tinggi*, 2020.
- [9] Noertjahyana, Agustinus and Yulia, “Studi Analisa Pelatihan Jaringan Syaraf Tiruan dengan dan tanpa Algoritma Genetika,” *Jurnal Informatika*, vol. 3, pp. 13-18, 2002.
- [10] P. A. Octaviani, Y. Wilandari and D. Ispriyanti, “Penerapan Metode Klasifikasi Support Vector Machine (SVM) Pada Data Akreditasi Sekolah Dasar (SD) Di Kabupaten Magelang,” *Jurnal Gaussian*, pp. 811-820, 2014.
- [11] K. Teknomo, “K-Nearest Neighbours Tutorial, K-Tetangga Terdekat Tutorial,” 2010.
- [12] R. A. Wattimena, “Kemerdekaan Pikiran,” 2012.

Lampiran A. TOR Digital Skola

- A.1. Quiz*
- B.1. Homework*
- C.1. Learning Progress Review*
- D.1. Personal Branding*
- E.1. Dataset & Modelling*

Capaian Pembelajaran Lulusan

Topik	Tingkat Kompetensi	Detail Pembelajaran	Durasi Total Jam (Belajar/ Bimbingan & Tugas)	Sumber Daya Pembelajaran	Cara Penilaian
<i>Analytic</i>	Memahami materi yang terkait dengan <i>Analytics</i> , yang meliputi metodologi <i>data science</i> , <i>Numpy</i> , <i>dataframe</i> , statistik, visualisasi data, dan <i>business intelligence</i> .	<ol style="list-style-type: none"> 1. <i>Introduction to data Science</i> 2. <i>Data Science Methodology</i> 3. <i>Introduction to Numpy</i> 4. <i>Introduction and basic dataframe (Pandas)</i> 5. <i>Intermediate dataframe</i> 6. <i>Advance dataframe</i> 7. <i>Business intelligence</i> 	60	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay</i>/koding) 3. Merangkum materi
<i>Data & Database</i>	Memahami segala hal tentang data dan <i>database</i> , termasuk di dalamnya <i>Structured Query Language (SQL)</i> , <i>version control</i> , <i>Application Programming Interface (API)</i> , and <i>database programming</i>	<ol style="list-style-type: none"> 1. <i>Introduction to data and database</i> 2. <i>Basic SQL</i> 3. <i>Intermediate SQL</i> 4. <i>Advance SQL</i> 5. <i>Versioning/version control</i> 6. <i>Data & Database</i> 7. <i>API</i> 	73	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay</i>/koding) 3. Merangkum materi

<i>Python Programming</i>	Menguasai bahasa pemrograman Python termasuk <i>libraries</i> yang digunakan	<ol style="list-style-type: none"> 1. <i>Introduction to python and programming</i> 2. <i>Basic programming I: conditions</i> 3. <i>Basic programming II: Iteration</i> 4. <i>Basic Programming III: Array and Other data types</i> 5. <i>Basic Programming IV: functions</i> 	52	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay</i>/koding) 3. Merangkum materi
<i>Statistics</i>	Memahami dasar- dasar statistika yang digunakan pada domain <i>data science</i>	<ol style="list-style-type: none"> 1. <i>Basic statistic</i> 2. <i>Intermediate statistic</i> 3. <i>Advanced statistic</i> 	31	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay</i>/koding) 3. Merangkum materi
<i>Data Visualization</i>	Mampu membuat visualisasi data secara efektif	<ol style="list-style-type: none"> 1. <i>Introduction to data visualization</i> 2. <i>Intermediate visualization</i> 3. <i>Data visualization exercise</i> 	42	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay</i>/koding) 3. Merangkum materi

<i>Machine Learning</i>	Mampu membuat model <i>machine learning</i> sesuai kebutuhan	<ol style="list-style-type: none"> 1. <i>Introduction to data mining</i> 2. <i>Introduction to machine learning</i> 3. <i>Data preprocessing for ML (python)</i> 4. <i>Advanced data preprocessing for ML (python)</i> 5. <i>Classification I</i> 6. <i>Classification II</i> 7. <i>Regression</i> 8. <i>Unsupervised learning</i> 9. <i>Evaluation metrics & model selection</i> 10. <i>Advanced machine learning topics</i> 	104	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay/koding</i>) 3. Merangkum materi
<i>Artificial Intelligence</i>	Memahami cara penggunaan teknologi AI dan mampu mengimplementasikannya pada bisnis	<ol style="list-style-type: none"> 1. <i>Introduction to Artificial Intelligence</i> 2. <i>Advanced data preprocessing for text</i> 3. <i>Advanced data preprocessing for image</i> 4. <i>neural network</i> 5. <i>Classification for text dataset</i> 6. <i>Classification for image dataset</i> 7. <i>Time Series forecasting</i> 8. <i>Unsupervised Learning II</i> 9. <i>Semi-supervised learning</i> 10. <i>Association rules</i> 11. <i>Outlier detection method</i> 12. <i>Recommender system</i> 	124	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Kuis 2. PR (<i>essay/koding</i>) 3. Merangkum materi

Soft Skill	<p><i>Communication & Presentation Skill</i> : Mampu berkomunikasi secara efektif, baik komunikasi individual maupun komunikasi kepada publik</p>	<p>Pemahaman mengenai:</p> <ol style="list-style-type: none"> 1. <i>Story Telling</i> 2. <i>Effective presentation</i> 	36	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	<p>Peserta diberikan tugas berupa:</p> <ol style="list-style-type: none"> 1. Kuis 2. Merangkum Materi
	<p><i>Analytical & Critical Thinking</i>: Mampu melihat masalah dari sisi user serta memberikan solusi yang tepat</p>	<p>Pemahaman mengenai:</p> <ol style="list-style-type: none"> 1. <i>Design Thinking</i> 2. <i>Design thinking steps</i> 3. <i>Design thinking tools</i> 4. <i>Assume a beginner's mindset</i> 5. <i>Ask the 5 whys</i> 6. <i>Empthy map</i> 		<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	<p>Peserta diberikan tugas berupa:</p> <ol style="list-style-type: none"> 1. Kuis 2. Merangkum Materi
	<p><i>Career Coaching & Mentoring</i>: Mampu membuat CV yang efektif dan mampu membuat kesan yang baik pada saat di <i>interview</i></p>	<p>Pemahaman mengenai:</p> <ol style="list-style-type: none"> 1. <i>How to make a powerful CV</i> 2. <i>How to make a good impression in an interview</i> 3. <i>How to express your desired benefits in an interview</i> 		<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	<p>Peserta diberikan tugas berupa:</p> <ol style="list-style-type: none"> 1. Kuis 2. Merangkum Materi

<i>Dataset & Modelling Project</i>	Mampu mengerjakan <i>Project dataset</i> yang diberikan serta membuat <i>modeling</i> -nya	Pemahaman mengenai: <ol style="list-style-type: none"> 1. <i>How to understand the dataset given</i> 2. <i>How to identify which activities should be done</i> 3. <i>using exploratory and visualization</i> 4. <i>How to do data pre-processing</i> 5. <i>How to develop model which relevant with the problem</i> 6. <i>How to evaluate the model</i> 	97	<ol style="list-style-type: none"> 1. <i>Interactive live-online learning</i> (dengan <i>Expert</i> maupun Pembimbing) 2. Silabus pembelajaran (PPT) 3. <i>Video recording</i> (pengulangan materi) 	Peserta diberikan tugas berupa: <ol style="list-style-type: none"> 1. Mengerjakan <i>dataset</i> dan membuat <i>modeling</i> hasil <i>modeling</i> yangtelah dibuat 2. Menjawab pertanyaanpenilai (juri)
--	--	--	----	--	---

Lampiran B. Log Activity

Minggu/Tgl	Kegiatan	Hasil
24 Agustus 2021	<i>Class Introduction</i>	Selesai
26 Agustus 2021	<i>Introduction to Data Science</i>	Selesai
28 Agustus 2021	<i>Data Science Methodology.</i>	Selesai
31 Agustus 2021	<i>Introduction to Data and Database</i>	Selesai
2 September 2021	<i>Basic SQL</i>	Selesai
4 September 2021	<i>Intermediate SQL</i>	Selesai
7 September 2021	<i>Advanced SQL</i>	Selesai
9 September 2021	<i>Versioning/Version Control</i>	Selesai
11 September 2021	<i>Introduction to Python and Programming.</i>	Selesai
14 September 2021	<i>Basic Programming I: Conditions,</i>	Selesai
16 September 2021	<i>Basic Programming II: Iteration</i>	Selesai

18 September 2021	<i>Basic Programming III: Array and Other Data Types.</i>	Selesai
21 September 2021	<i>Basic Programming IV: Functions</i>	Selesai
23 September 2021	<i>Database Programming</i>	Selesai
25 September 2021	<i>Introduction to Numpy</i>	Selesai
28 September 2021	<i>Introduction and Basic Dataframe (Pandas)</i>	Selesai
30 September 2021	<i>Kaggle Project Stage 1</i>	Selesai
2 Oktober 2021	<i>Analytical & Critical Thinking</i>	Selesai
5 Oktober 2021	<i>Intermediate Dataframe I</i>	Selesai
7 Oktober 2021	<i>Advanced Dataframe</i>	Selesai
9 Oktober 2021	<i>API</i>	Selesai
12 Oktober 2021	<i>Basic Statistics,</i>	Selesai
14 Oktober 2021	<i>Intermediate Statistics</i>	Selesai

16 Oktober 2021	<i>Advanced Statistics</i>	Selesai
19 Oktober 2021	<i>Introduction to Data Visualization,</i>	Selesai
21 Oktober 2021	<i>Intermediate Visualization</i>	Selesai
23 Oktober 2021	<i>Data Visualization Exercises (Advance).</i>	Selesai
26 Oktober 2021	<i>Introduction to Data Mining</i>	Selesai
28 Oktober 2021	<i>Introduction to Machine Learning</i>	Selesai
30 Oktober 2021	<i>Data Preprocessing for ML (python).</i>	Selesai
2 November 2021	<i>Advanced Data Preprocessing for ML (python)</i>	Selesai
4 November 2021	<i>Classification I</i>	Selesai
6 November 2021	<i>Classification II</i>	Selesai
9 November 2021	<i>Regression</i>	Selesai
11 November 2021	<i>Unsupervised Learning</i>	Selesai
13 November 2021	<i>Communication & Presentation Skill</i>	Selesai

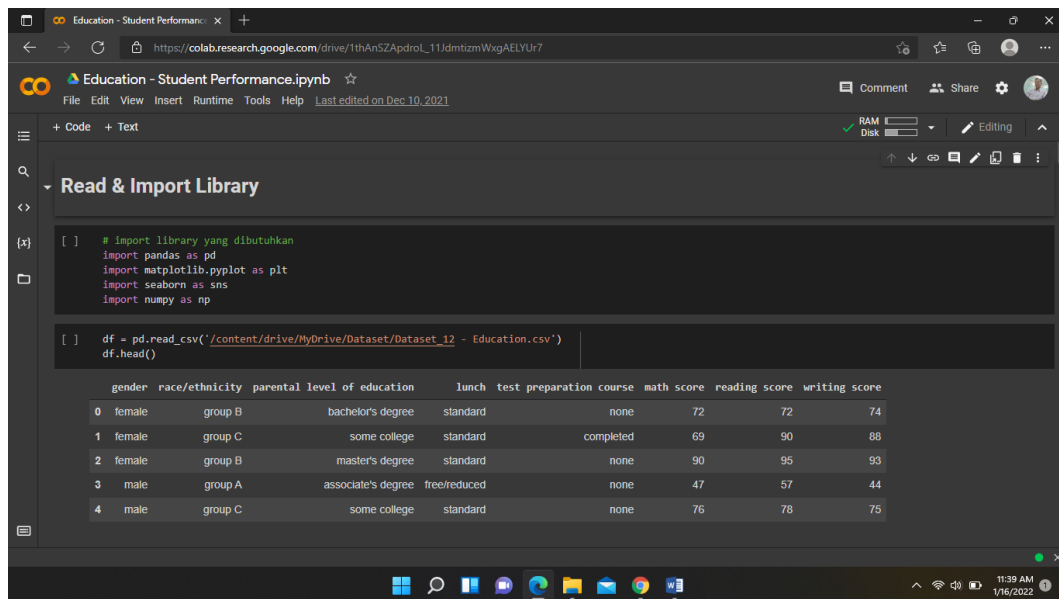
16 November 2021	<i>Evaluation Metrics and Model Selection</i>	Selesai
18 November 2021	<i>Advanced ML Topics</i>	Selesai
20 November 2021	<i>Business Intelligence</i>	Selesai
23 November 2021	<i>Mentor Experience Sharing</i>	Selesai
25 November 2021	<i>HR Practitioner Sharing</i>	Selesai
27 November 2021	<i>Kaggle Project Presentation</i>	Selesai
29 November 2021	<i>Kaggle Project Presentation</i>	Selesai
2 Desember 2021	<i>Kaggle Project Presentation</i>	Selesai
5 Desember 2021	<i>Introduction to Artificial Intelligence.</i>	Selesai
7 Desember 2021	<i>Advanced Data Preprocessing for Text,</i>	Selesai
9 Desember 2021	<i>Advanced Data Preprocessing for Image</i>	Selesai

11 Desember 2021	<i>Neural Network</i>	Selesai
14 Desember 2021	<i>Classification for Text Dataset</i>	Selesai
16 Desember 2021	<i>Classification for Image Dataset,</i>	Selesai
18 Desember 2021	<i>Time series Forecasting</i>	Selesai
21 Desember 2021	<i>Unsupervised Learning II</i>	Selesai
23 Desember 2021	<i>Semi-supervised Learning</i>	Selesai
28 Desember 2021	<i>Association Rules</i>	Selesai
30 Desember 2021	<i>Outlier Detection Method</i>	Selesai
4 Januari 2022	<i>Recommender System</i>	Selesai
20 Januari 2022	<i>Kaggle Project Presentation</i>	Selesai
22 Januari 2022	<i>Kaggle Project Presentation</i>	Selesai
25 Januari 2022	<i>Kaggle Project Presentation</i>	Selesai
26 Januari 2022	<i>Finalisasi Draft Laporan</i>	Selesai
27 Januari 2022	<i>Finalisasi Draft Laporan</i>	Selesai

Lampiran C. Dokumen Teknik

Berikut ini adalah *screenshot* hasil *google colab* dalam pengerjaan *project* program MSIB Digital Skola:

C.1. *Import libraries and Read Dataset*



The screenshot shows a Google Colab notebook titled "Education - Student Performance.ipynb". The code cell contains the following Python code:

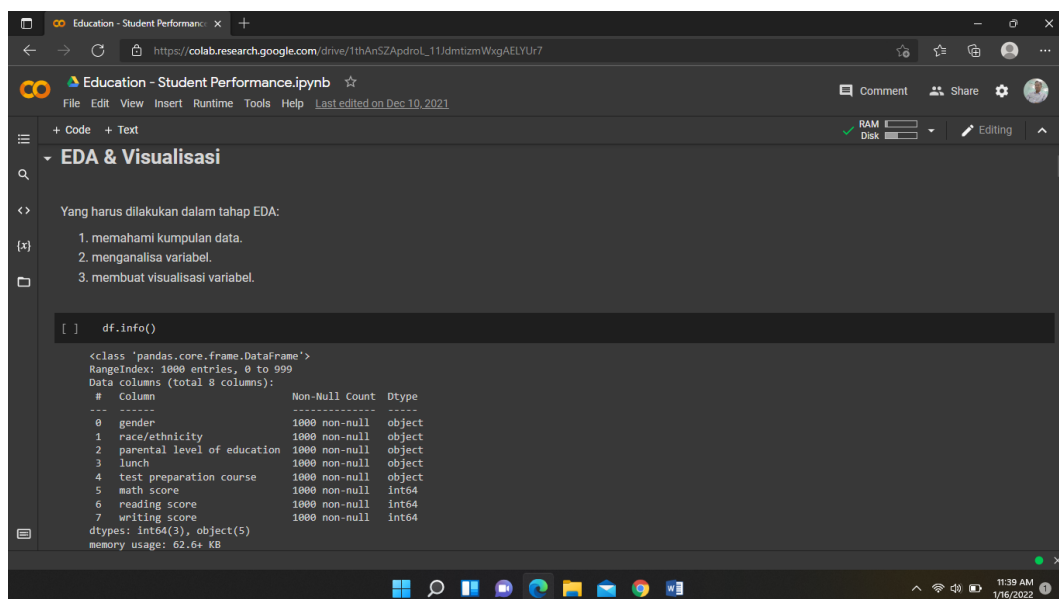
```
[ ] # Import library yang dibutuhkan
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

[ ] df = pd.read_csv('/content/drive/MyDrive/Dataset/Dataset_12 - Education.csv')
df.head()
```

The output of the code is a preview of the dataset, showing the first five rows of a DataFrame with 8 columns: gender, race/ethnicity, parental level of education, lunch, test preparation course, math score, reading score, and writing score.

	gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score
0	female	group B	bachelor's degree	standard	none	72	72	74
1	female	group C	some college	standard	completed	69	90	88
2	female	group B	master's degree	standard	none	90	95	93
3	male	group A	associate's degree	free/reduced	none	47	57	44
4	male	group C	some college	standard	none	76	78	75

I.2. *Exploratory Data Analysis & Visualization*

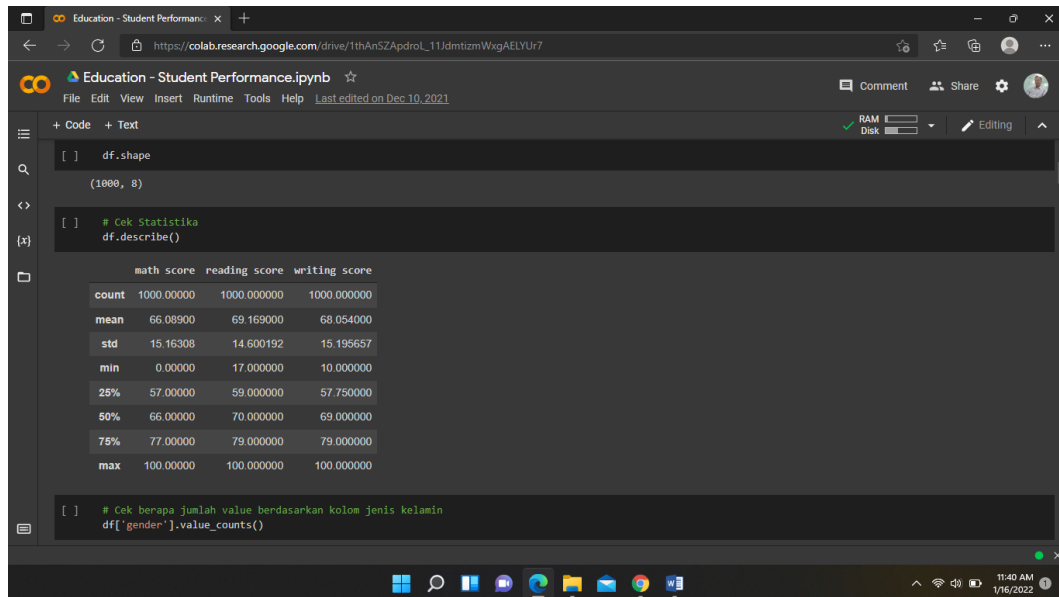


The screenshot shows a Google Colab notebook titled "Education - Student Performance.ipynb". The code cell contains the following Python code:

```
[ ] df.info()
```

The output of the code is the information of the DataFrame, showing the data types and non-null counts for each column.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 8 columns):
 #   column              Non-Null Count  Dtype
---  -
 0   gender              1000 non-null  object
 1   race/ethnicity      1000 non-null  object
 2   parental level of education  1000 non-null  object
 3   lunch               1000 non-null  object
 4   test preparation course  1000 non-null  object
 5   math score          1000 non-null  int64
 6   reading score       1000 non-null  int64
 7   writing score        1000 non-null  int64
dtypes: int64(3), object(5)
memory usage: 62.6+ KB
```



The screenshot shows a Jupyter Notebook interface with the following code and output:

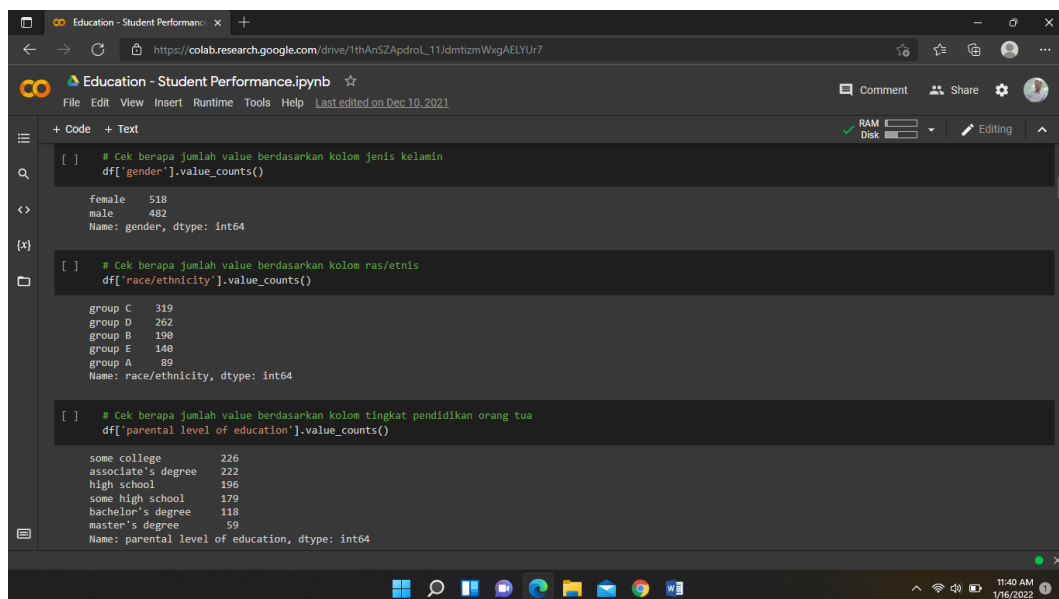
```
[ ] df.shape
```

```
(1000, 8)
```

```
[ ] # Cek Statistika
df.describe()
```

	math score	reading score	writing score
count	1000.00000	1000.000000	1000.000000
mean	66.08900	69.169000	68.054000
std	15.16308	14.600192	15.195657
min	0.00000	17.000000	10.000000
25%	57.00000	59.000000	57.750000
50%	66.00000	70.000000	69.000000
75%	77.00000	79.000000	79.000000
max	100.00000	100.000000	100.000000

```
[ ] # Cek berapa jumlah value berdasarkan kolom jenis kelamin
df['gender'].value_counts()
```



The screenshot shows a Jupyter Notebook interface with the following code and output:

```
[ ] # cek berapa jumlah value berdasarkan kolom jenis kelamin
df['gender'].value_counts()
```

```
female    518
male      482
Name: gender, dtype: int64
```

```
[ ] # Cek berapa jumlah value berdasarkan kolom ras/etnis
df['race/ethnicity'].value_counts()
```

```
group C    319
group D    262
group B    190
group E    140
group A     89
Name: race/ethnicity, dtype: int64
```

```
[ ] # Cek berapa jumlah value berdasarkan kolom tingkat pendidikan orang tua
df['parental level of education'].value_counts()
```

```
some college    226
associate's degree    222
high school    196
some high school    179
bachelor's degree    118
master's degree     59
Name: parental level of education, dtype: int64
```

```

[ ] # Cek berapa jumlah value berdasarkan kolom makan siang
df['lunch'].value_counts()

standard      645
free/reduced  355
Name: lunch, dtype: int64

[ ] # Cek berapa jumlah value berdasarkan kolom kursus persiapan ujian
df['test_preparation_course'].value_counts()

none          642
completed    358
Name: test preparation course, dtype: int64

Pie Chart

[ ] female_data = df[(df['gender'] == 'female')]
male_data = df[(df['gender'] == 'male')]

[ ] # Berdasarkan tiap tiap kolom
plt.figure(figsize=(15, 7)) # Membuat ukuran
plt.subplot(2, 3, 1) #2 baris, 3 kolom, plot 1
data = df['gender'].value_counts()
label = data.index

```

```

[ ] # Berdasarkan tiap tiap kolom
plt.figure(figsize=(15, 7)) # Membuat ukuran
plt.subplot(2, 3, 1) #2 baris, 3 kolom, plot 1
data = df['gender'].value_counts()
label = data.index

plt.pie(data, labels=label, autopct='%.2f%%')
plt.title('Gender')

plt.subplot(2, 3, 2) #2 baris, 3 kolom, plot 2
data = df['race/ethnicity'].value_counts()
label = data.index

plt.pie(data, labels=label, autopct='%.2f%%')
plt.title('Race/ethnicity')

plt.subplot(2, 3, 3) #2 baris, 3 kolom, plot 3
data = df['parental level of education'].value_counts()
label = data.index

plt.pie(data, labels=label, autopct='%.2f%%')
plt.title('Parental level of education')

plt.subplot(2, 3, 4) #2 baris, 3 kolom, plot 4
data = df['lunch'].value_counts()
label = data.index

```

```

Education - Student Performance.ipynb
https://colab.research.google.com/drive/1thAnSZApdrol_11JdmtizmWxgAELYUr7

data = df['race/ethnicity'].value_counts()
label = data.index

plt.pie(data, labels=label, autopct='%1.2f%%')
plt.title('Race/ethnicity')

plt.subplot(2, 3, 3) #2 baris, 3 kolom, plot 3
data = df['parental level of education'].value_counts()
label = data.index

plt.pie(data, labels=label, autopct='%1.2f%%')
plt.title('Parental level of education')

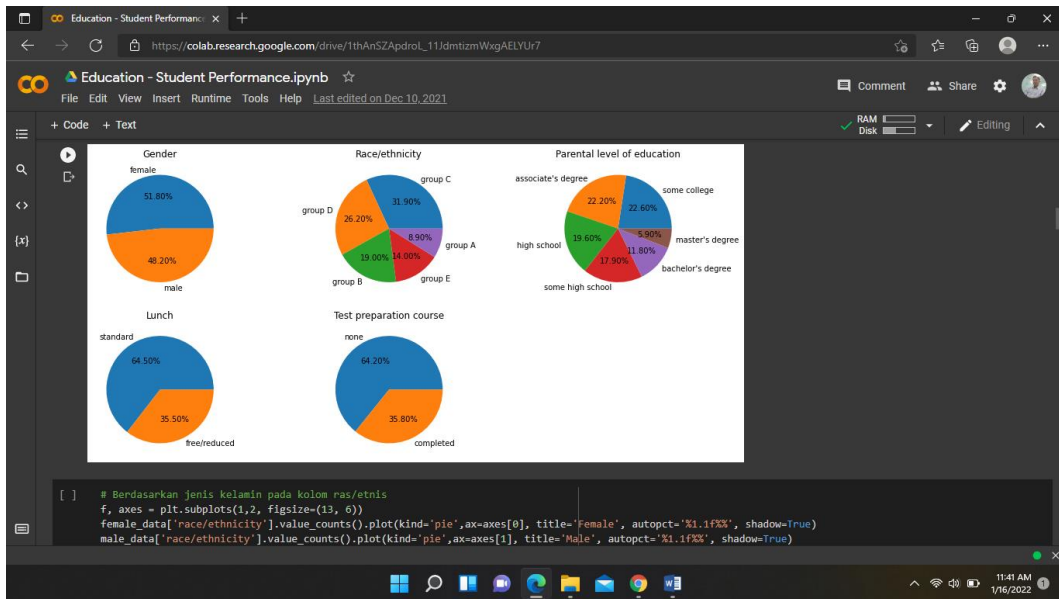
plt.subplot(2, 3, 4) #2 baris, 3 kolom, plot 4
data = df['lunch'].value_counts()
label = data.index

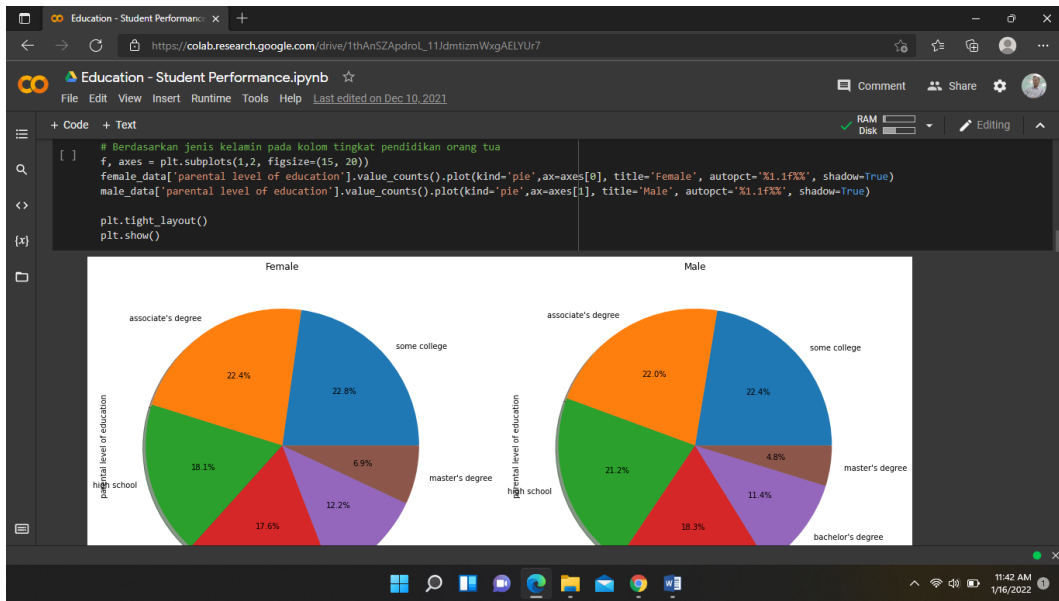
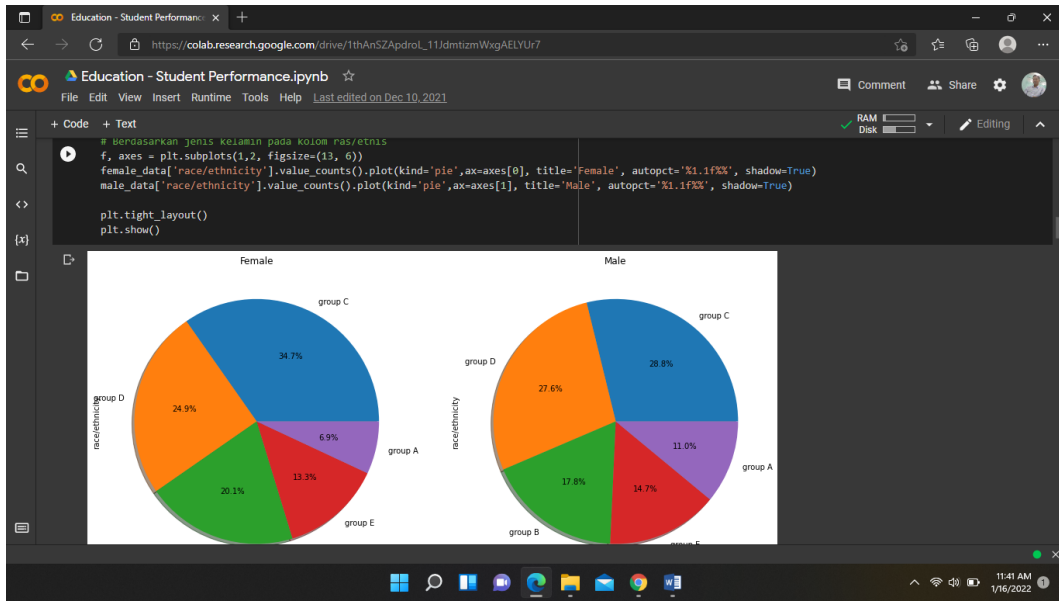
plt.pie(data, labels=label, autopct='%1.2f%%')
plt.title('Lunch')

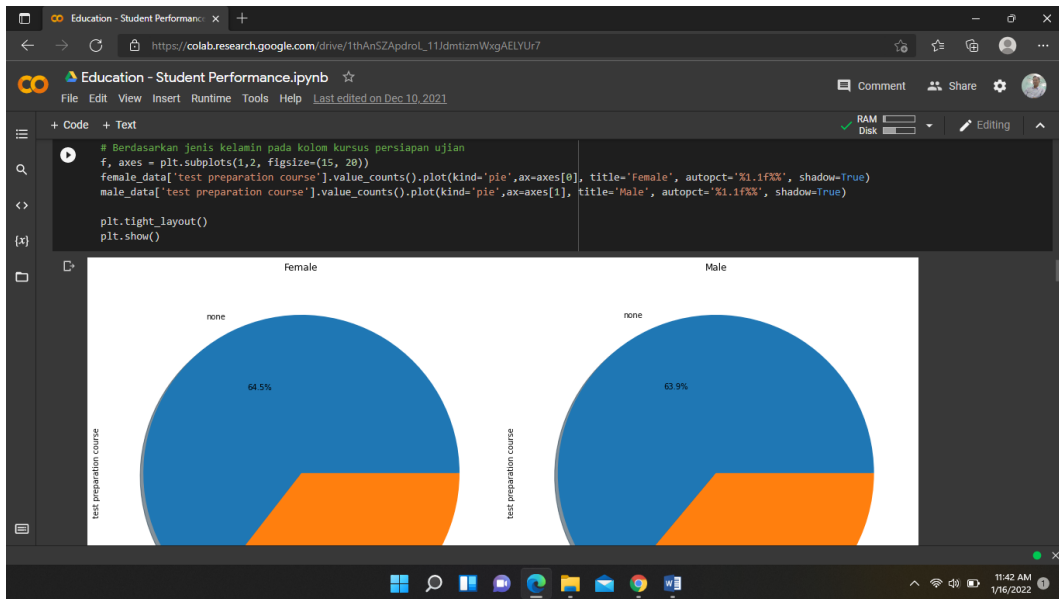
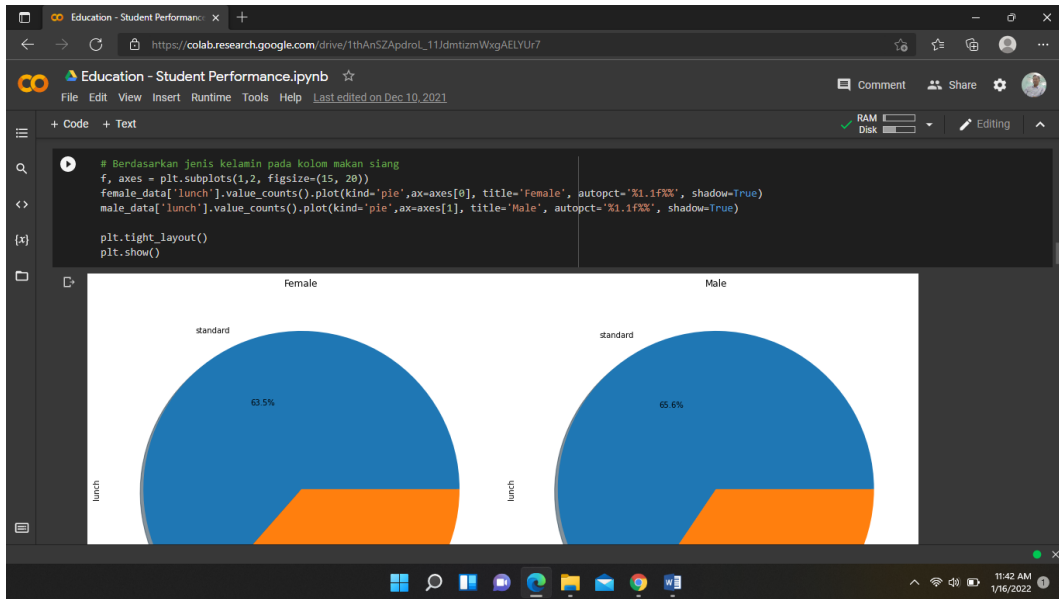
plt.subplot(2, 3, 5) #2 baris, 3 kolom, plot 5
data = df['test preparation course'].value_counts()
label = data.index

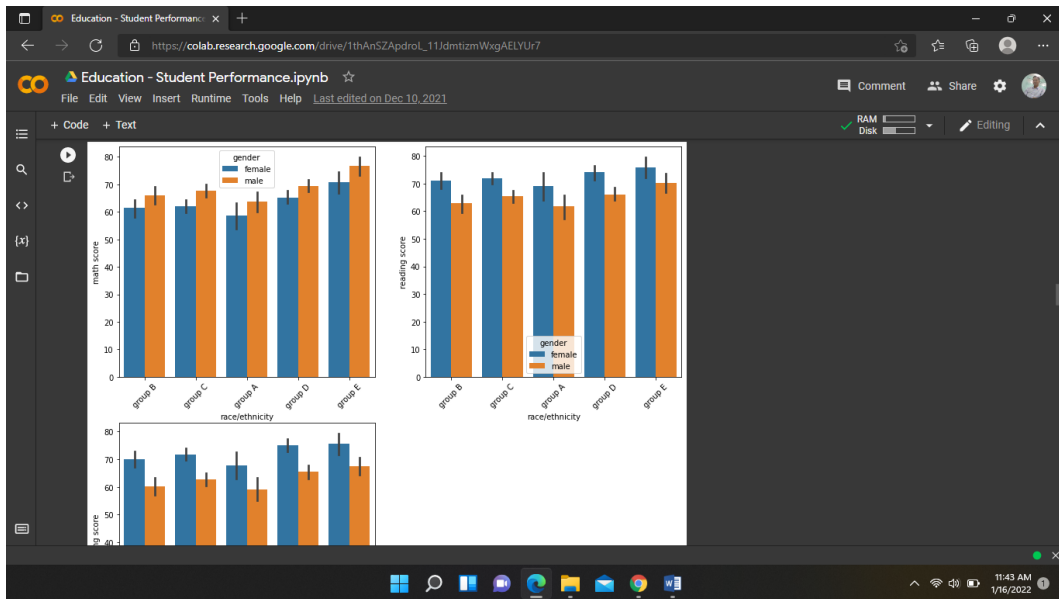
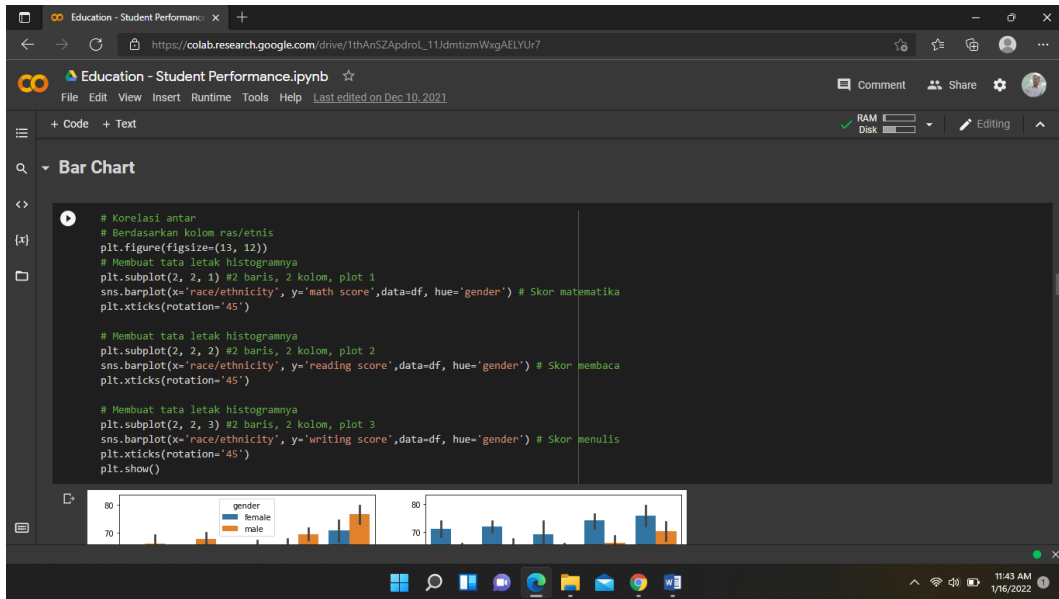
plt.pie(data, labels=label, autopct='%1.2f%%')
plt.title('Test preparation course')
plt.show()

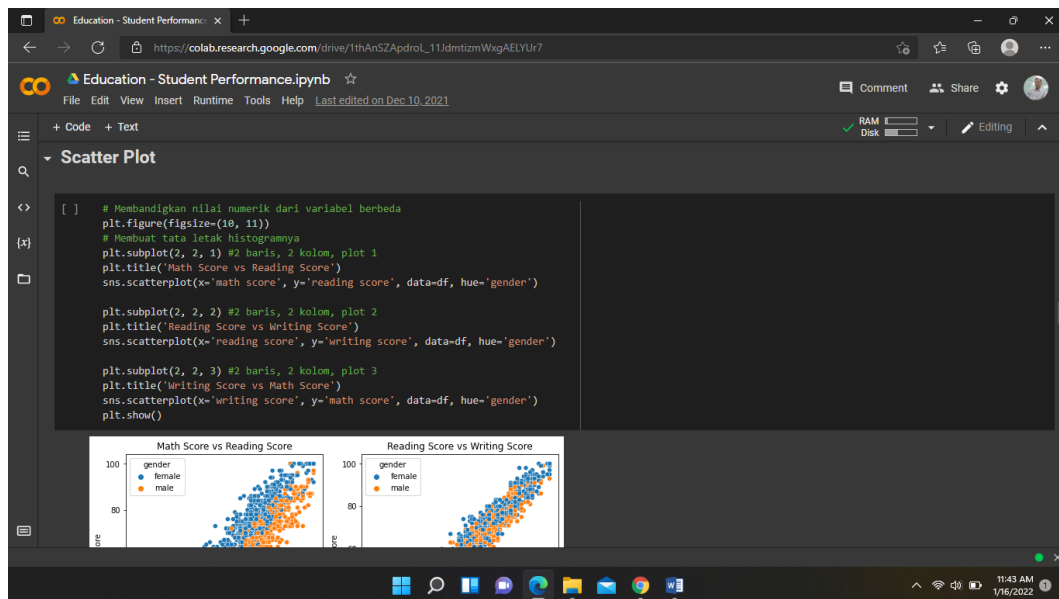
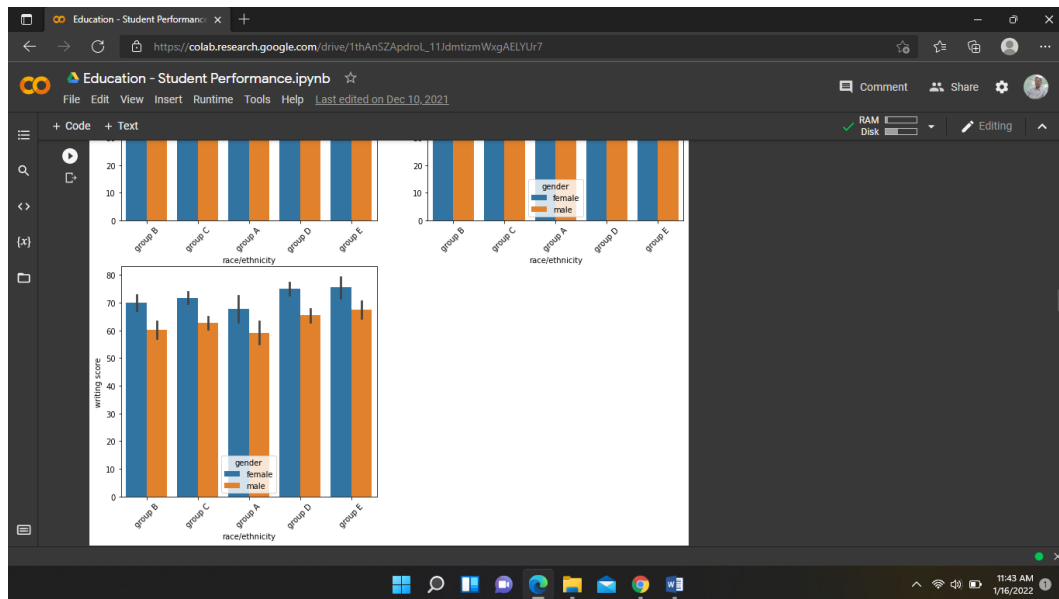
```

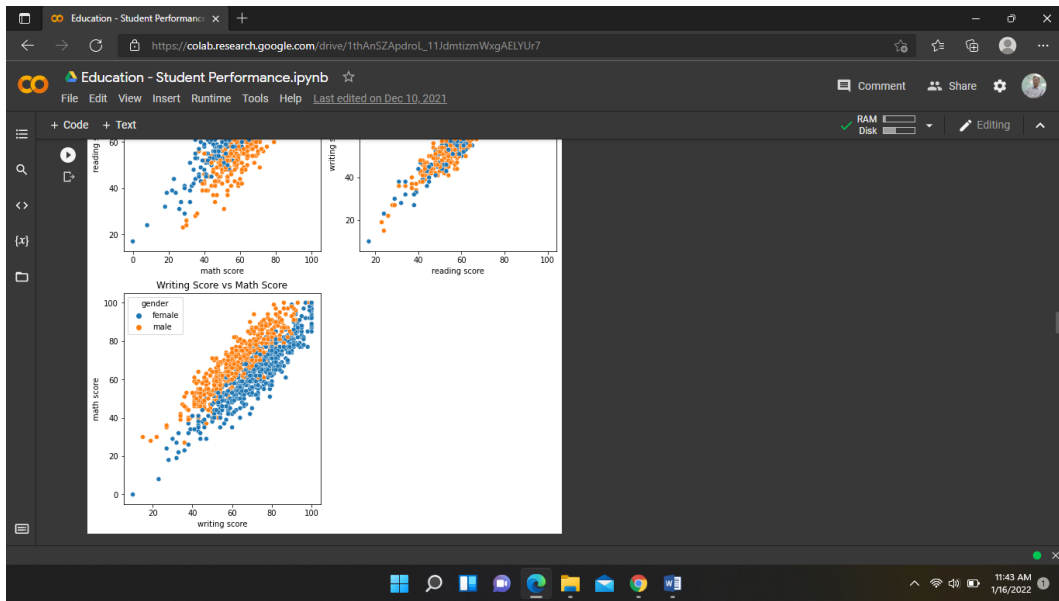
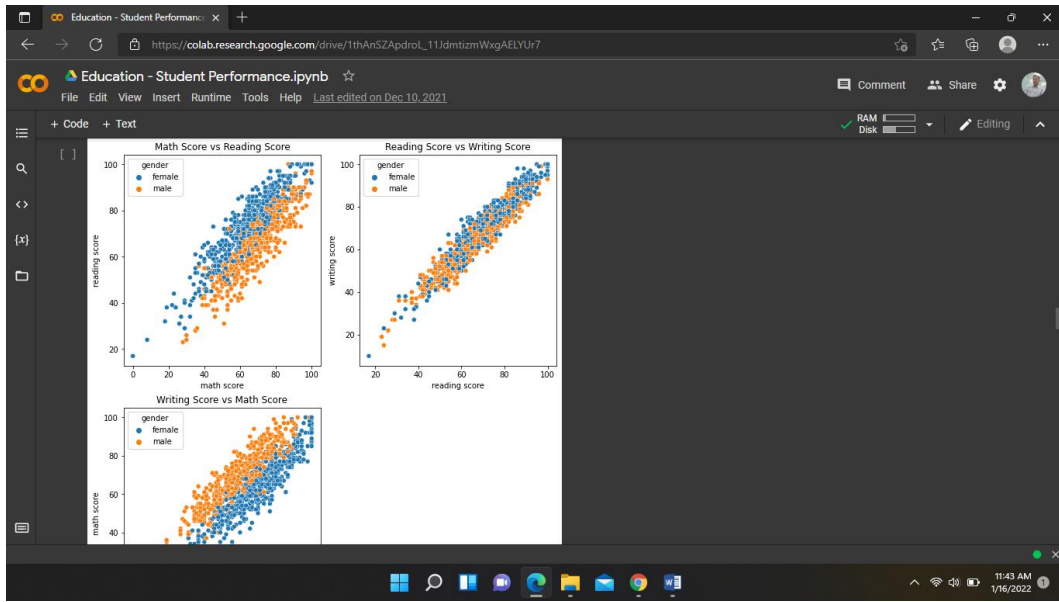




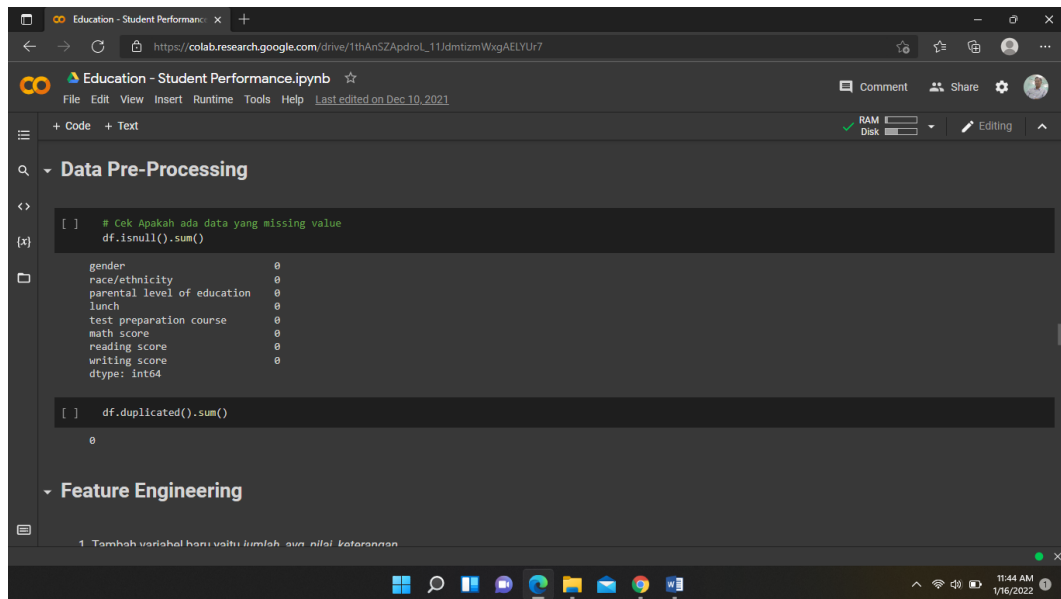








Data Preprocessing



```

[ ] # Cek Apakah ada data yang missing value
df.isnull().sum()

gender          0
race/ethnicity  0
parental level of education  0
lunch           0
test preparation course  0
math score      0
reading score   0
writing score   0
dtype: int64

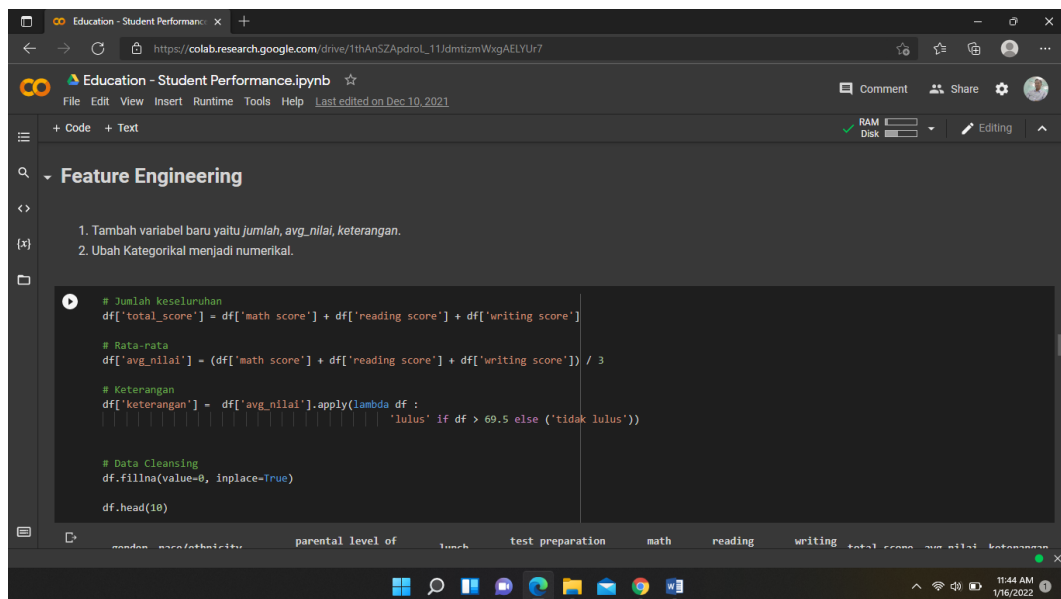
[ ] df.duplicated().sum()

0

```

1. Tambah variabel baru yaitu jumlah dan nilai ketangan

Feature Engineering



```

1. Tambah variabel baru yaitu jumlah, avg_nilai, keterangan.
2. Ubah Kategorikal menjadi numerikal.

# Jumlah keseluruhan
df['total_score'] = df['math score'] + df['reading score'] + df['writing score']

# Rata-rata
df['avg_nilai'] = (df['math score'] + df['reading score'] + df['writing score']) / 3

# Keterangan
df['keterangan'] = df['avg_nilai'].apply(lambda df :
                                         'lulus' if df > 69.5 else ('tidak lulus'))

# Data Cleansing
df.fillna(value=0, inplace=True)

df.head(10)

```

Education - Student Performance.ipynb

```

# Data Cleansing
df.fillna(value=0, inplace=True)
df.head(10)

```

	gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score	total_score	avg_nilai	keterangan
0	female	group B	bachelor's degree	standard	none	72	72	74	218	72.666667	lulus
1	female	group C	some college	standard	completed	69	90	88	247	82.333333	lulus
2	female	group B	master's degree	standard	none	90	95	93	278	92.666667	lulus
3	male	group A	associate's degree	free/reduced	none	47	57	44	148	49.333333	tidak lulus
4	male	group C	some college	standard	none	76	78	75	229	76.333333	lulus
5	female	group B	associate's degree	standard	none	71	83	78	232	77.333333	lulus
6	female	group B	some college	standard	completed	88	95	92	275	91.666667	lulus
7	male	group B	some college	free/reduced	none	40	43	39	122	40.666667	tidak lulus
8	male	group D	high school	free/reduced	completed	64	64	67	195	65.000000	tidak lulus
9	female	group B	high school	free/reduced	none	38	60	50	148	49.333333	tidak lulus

Education - Student Performance.ipynb

```

# Label Encoder
from sklearn.preprocessing import LabelEncoder

# Menyalin / copy dataframe agar dataframe awal tetap utuh
df = df.copy()

label = LabelEncoder()

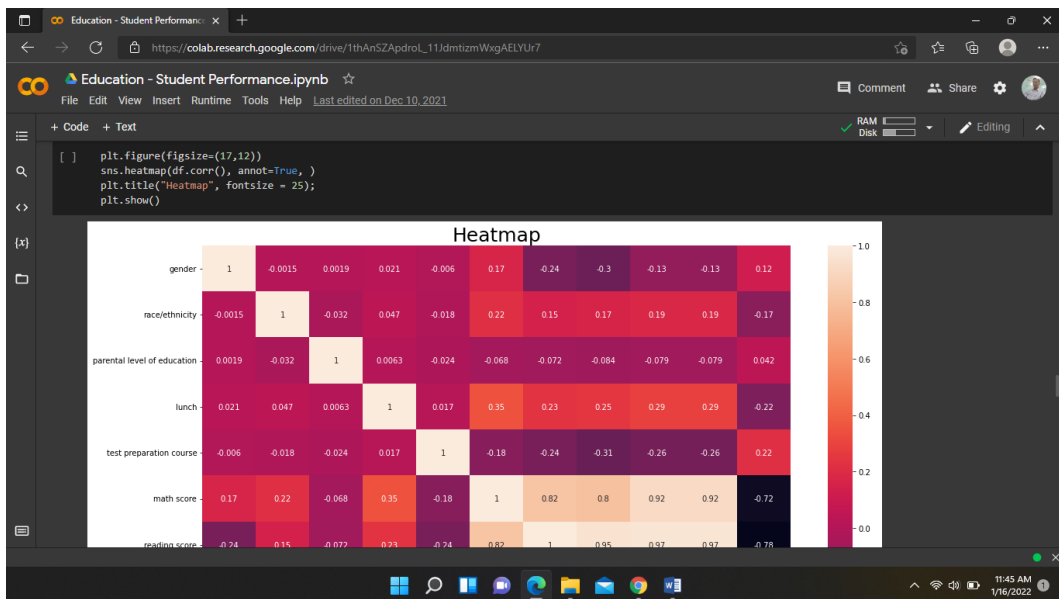
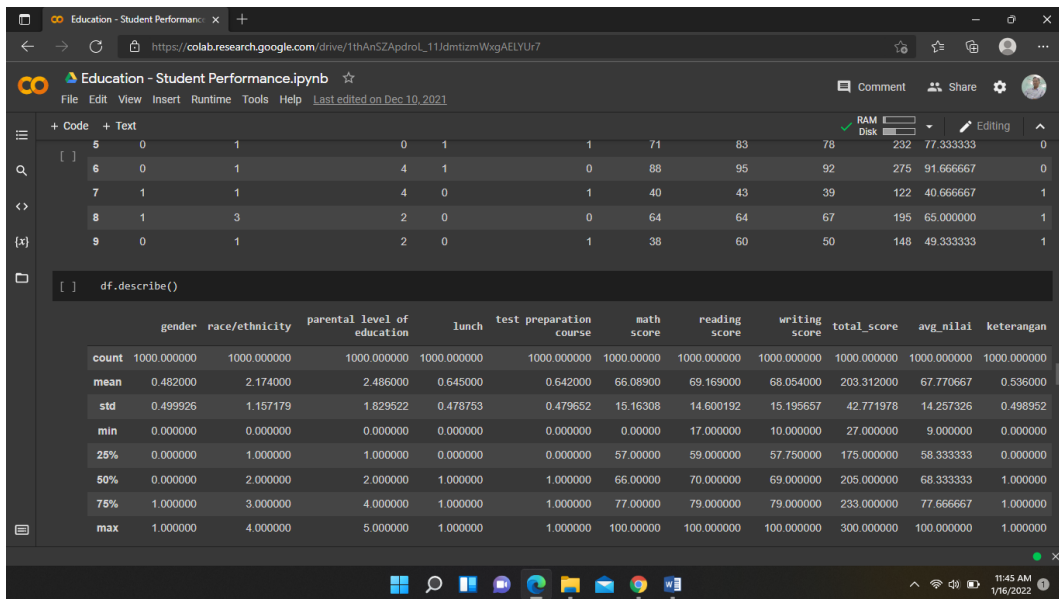
# Membuat list dari nama kolom data kategori
categorical_data = ['gender', 'race/ethnicity', 'parental level of education', 'lunch', 'test preparation course', 'keterangan']

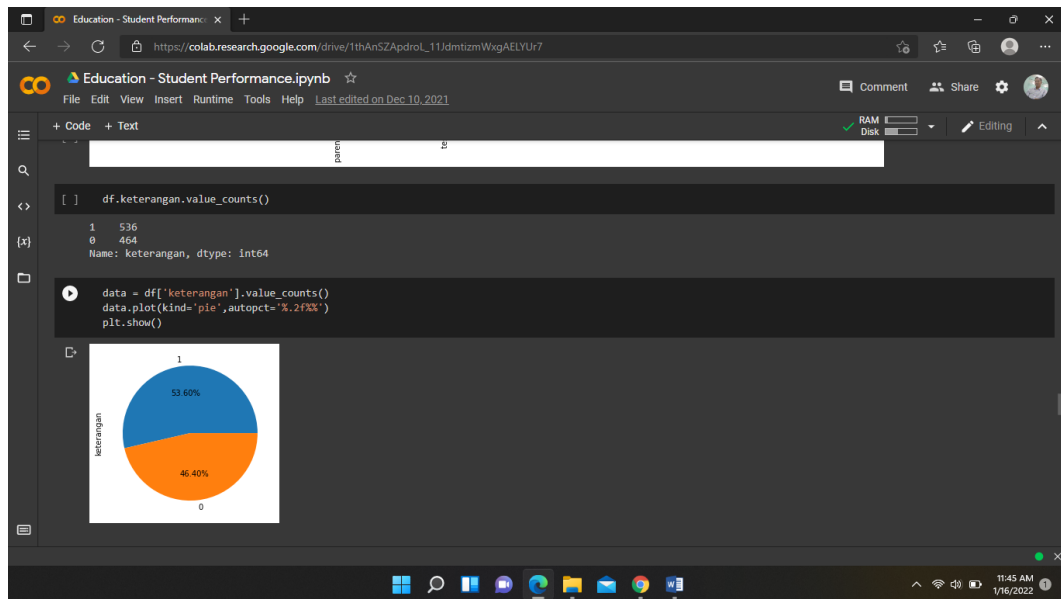
# Mengubah setiap data kategori menjadi numerik dengan encoder
for kolom in categorical_data:
    df[kolom] = label.fit_transform(df[kolom])

df.head(10)

```

	gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score	total_score	avg_nilai	keterangan
0	0	1	1	1	1	72	72	74	218	72.666667	0
1	0	2	4	1	0	69	90	88	247	82.333333	0
2	0	1	3	1	1	90	95	93	278	92.666667	0
3	1	0	0	0	1	47	57	44	148	49.333333	1
4	1	2	4	1	1	76	78	75	229	76.333333	0





Train-Test Split

The screenshot shows a Google Colab notebook titled "Education - Student Performance.ipynb". The code cell contains the following Python code:

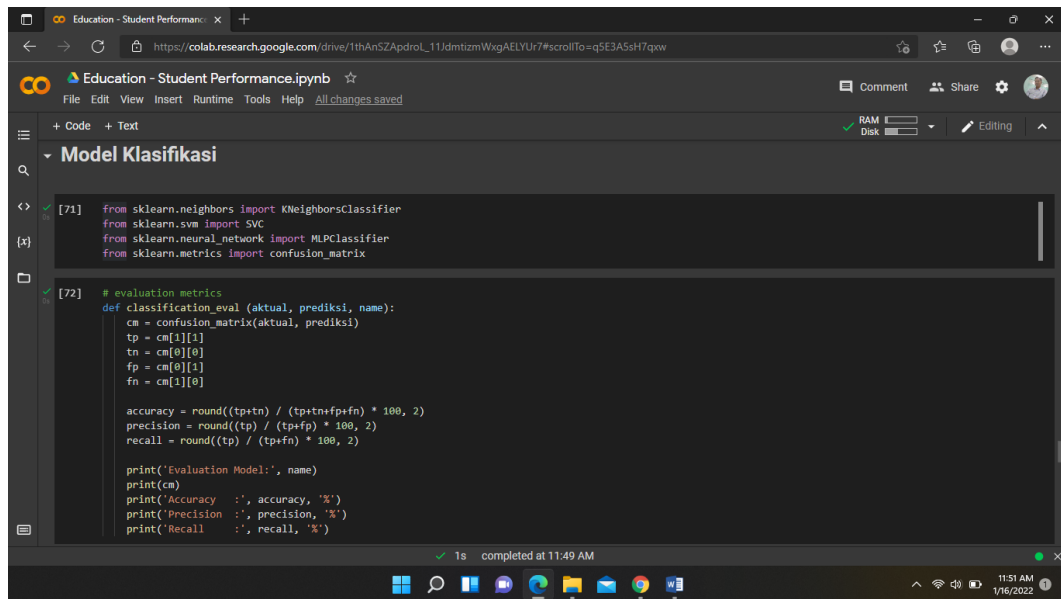
```
# define X & y
X = df.drop(['keterangan'], axis=1)
y = df['keterangan']

from sklearn.model_selection import train_test_split
# split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) # agar splitting tidak berubah, diberi random_state

# scaling
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

The code is organized into sections: "Train-Test Split" and "Model Klasifikasi".

Evaluation Model Machine Learning



The screenshot shows a Jupyter Notebook interface with a code cell titled "Model Klasifikasi". The code defines a function named `classification_eval` that takes three arguments: `aktual`, `prediksi`, and `name`. It uses `sklearn.metrics.confusion_matrix` to generate a confusion matrix. The function then calculates accuracy, precision, and recall based on the matrix elements. Finally, it prints the model name, the confusion matrix, and the calculated metrics.

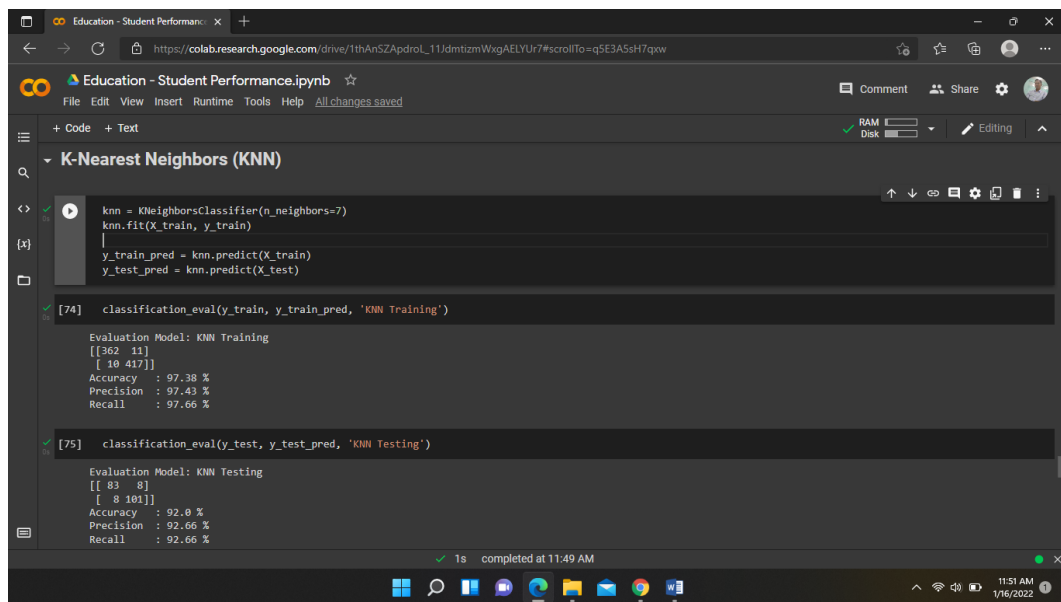
```
[71] from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.neural_network import MLPClassifier
from sklearn.metrics import confusion_matrix

[72] # evaluation metrics
def classification_eval (aktual, prediksi, name):
    cm = confusion_matrix(aktual, prediksi)
    tp = cm[1][1]
    tn = cm[0][0]
    fp = cm[0][1]
    fn = cm[1][0]

    accuracy = round((tp+tn) / (tp+tn+fp+fn) * 100, 2)
    precision = round(tp / (tp+fp) * 100, 2)
    recall = round(tp / (tp+fn) * 100, 2)

    print('Evaluation Model:', name)
    print(cm)
    print('Accuracy   :', accuracy, '%')
    print('Precision   :', precision, '%')
    print('Recall       :', recall, '%')
```

1s completed at 11:49 AM



The screenshot shows a Jupyter Notebook interface with a code cell titled "K-Nearest Neighbors (KNN)". The code creates a `KNeighborsClassifier` with `n_neighbors=7`, fits it to training data, and then uses it to predict on both training and testing data. The `classification_eval` function from the previous cell is used to evaluate the performance of the model on both datasets.

```
knn = KNeighborsClassifier(n_neighbors=7)
knn.fit(X_train, y_train)

y_train_pred = knn.predict(X_train)
y_test_pred = knn.predict(X_test)

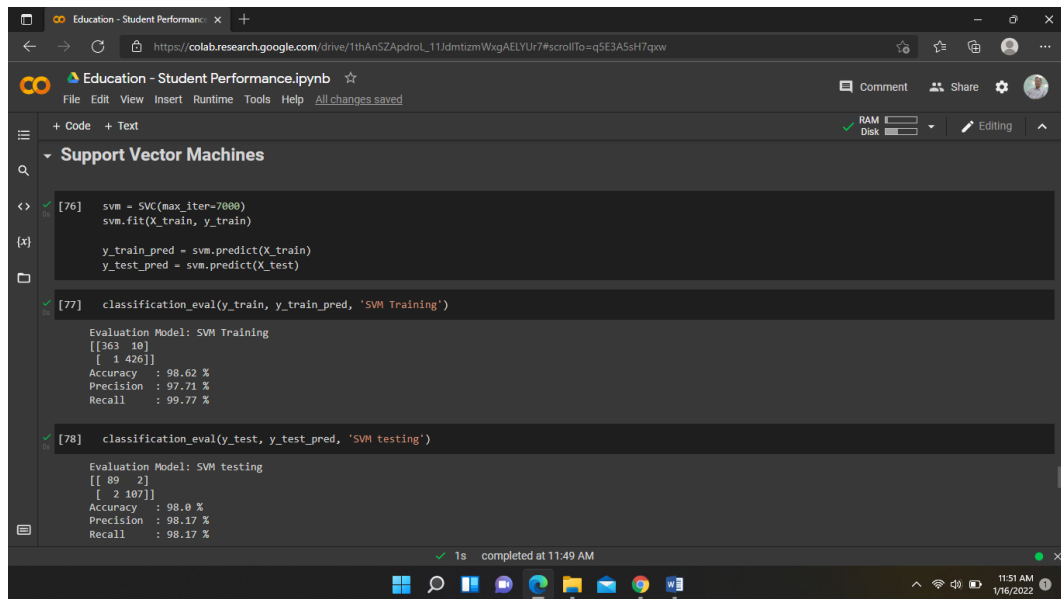
[74] classification_eval(y_train, y_train_pred, 'KNN Training')

Evaluation Model: KNN Training
[[362 11]
 [ 10 417]]
Accuracy   : 97.38 %
Precision   : 97.43 %
Recall     : 97.66 %

[75] classification_eval(y_test, y_test_pred, 'KNN Testing')

Evaluation Model: KNN Testing
[[ 83  8]
 [  8 101]]
Accuracy   : 92.0 %
Precision   : 92.66 %
Recall     : 92.66 %
```

1s completed at 11:49 AM



The screenshot shows a Google Colab notebook interface. The browser address bar displays the URL: `https://colab.research.google.com/drive/1thAnSZApdrol_11JdmtizmWxgAELYur7#scrollTo=q5E3A5sH7qww`. The notebook title is "Education - Student Performance.ipynb". The code editor shows the following code:

```
[76] svm = SVC(max_iter=7000)
      svm.fit(X_train, y_train)
      y_train_pred = svm.predict(X_train)
      y_test_pred = svm.predict(X_test)

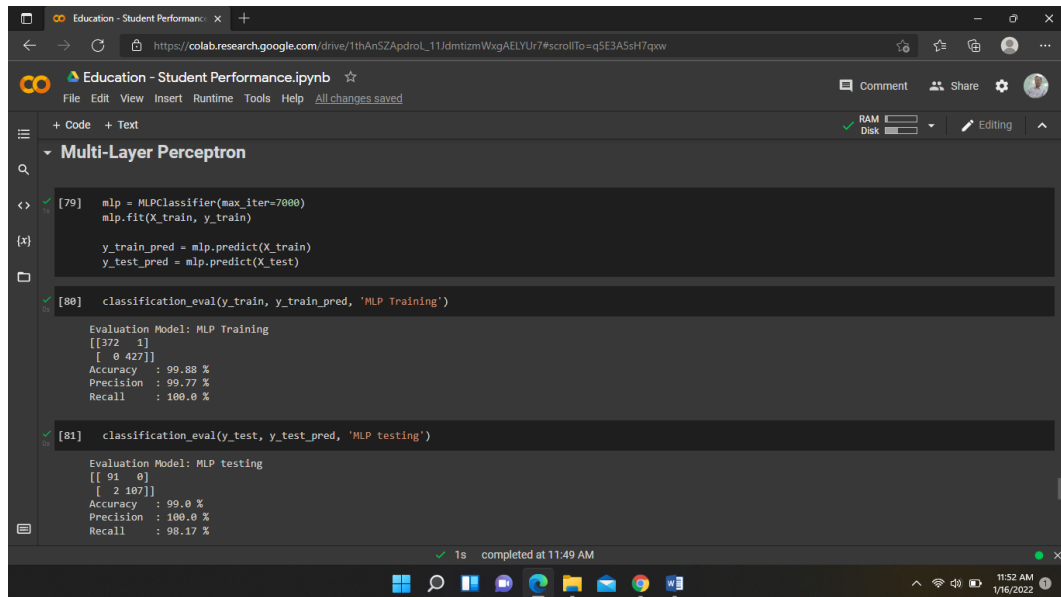
[77] classification_eval(y_train, y_train_pred, 'SVM Training')

Evaluation Model: SVM Training
[[363 10]
 [ 1 426]]
Accuracy : 98.62 %
Precision : 97.71 %
Recall : 99.77 %

[78] classification_eval(y_test, y_test_pred, 'SVM testing')

Evaluation Model: SVM testing
[[ 89  2]
 [ 2 107]]
Accuracy : 98.0 %
Precision : 98.17 %
Recall : 98.17 %
```

The status bar at the bottom indicates "1s completed at 11:49 AM". The system tray shows the time as 11:51 AM on 1/16/2022.



The screenshot shows a Google Colab notebook interface. The browser address bar displays the URL: `https://colab.research.google.com/drive/1thAnSZApdrol_11JdmtizmWxgAELYur7#scrollTo=q5E3A5sH7qww`. The notebook title is "Education - Student Performance.ipynb". The code editor shows the following code:

```
[79] mlp = MLPClassifier(max_iter=7000)
      mlp.fit(X_train, y_train)
      y_train_pred = mlp.predict(X_train)
      y_test_pred = mlp.predict(X_test)

[80] classification_eval(y_train, y_train_pred, 'MLP Training')

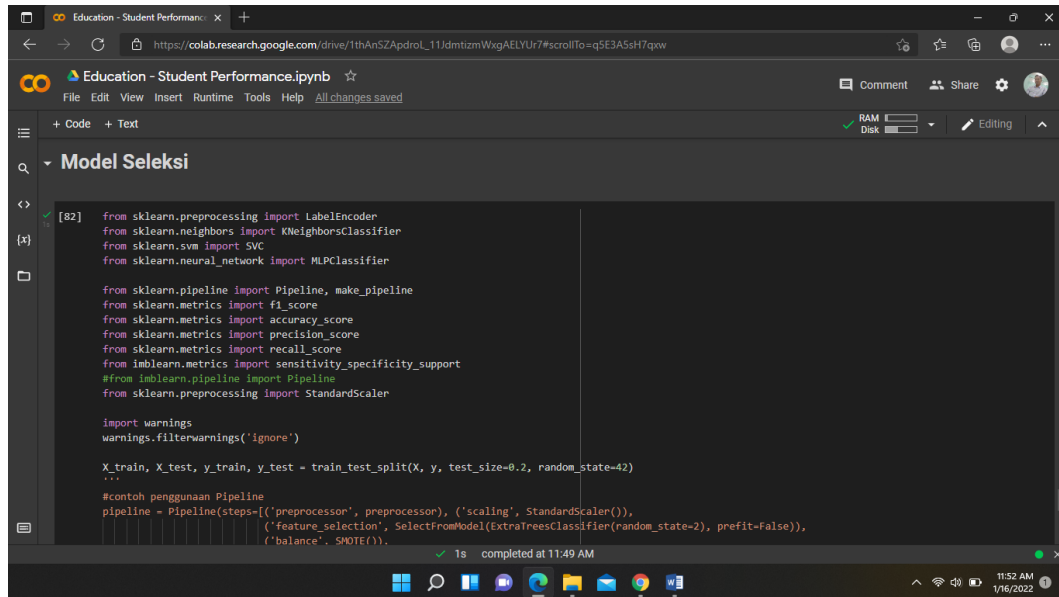
Evaluation Model: MLP Training
[[372  1]
 [ 0 427]]
Accuracy : 99.88 %
Precision : 99.77 %
Recall : 100.0 %

[81] classification_eval(y_test, y_test_pred, 'MLP testing')

Evaluation Model: MLP testing
[[ 91  0]
 [ 2 107]]
Accuracy : 99.0 %
Precision : 100.0 %
Recall : 98.17 %
```

The status bar at the bottom indicates "1s completed at 11:49 AM". The system tray shows the time as 11:52 AM on 1/16/2022.

Model Selection Machine Learning Pipeline



```

[82] from sklearn.preprocessing import LabelEncoder
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.neural_network import MLPClassifier

from sklearn.pipeline import Pipeline, make_pipeline
from sklearn.metrics import f1_score
from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score
from sklearn.metrics import recall_score
from imblearn.metrics import sensitivity_specificity_support
#from imblearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler

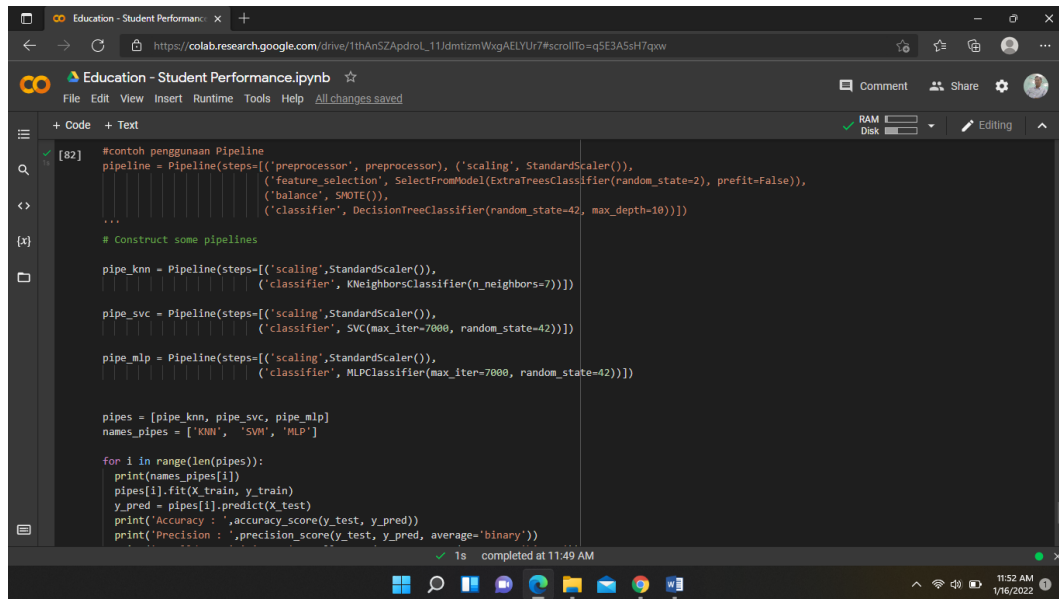
import warnings
warnings.filterwarnings('ignore')

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
...

#contoh penggunaan Pipeline
pipeline = Pipeline(steps=[('preprocessor', preprocessor), ('scaling', StandardScaler()),
                           ('feature_selection', SelectFromModel(ExtraTreesClassifier(random_state=2), prefit=False)),
                           ('balance', SMOTE())])

```

1s completed at 11:49 AM



```

[82] #contoh penggunaan Pipeline
pipeline = Pipeline(steps=[('preprocessor', preprocessor), ('scaling', StandardScaler()),
                           ('feature_selection', SelectFromModel(ExtraTreesClassifier(random_state=2), prefit=False)),
                           ('balance', SMOTE()),
                           ('classifier', DecisionTreeClassifier(random_state=42, max_depth=10))])

# Construct some pipelines
pipe_knn = Pipeline(steps=[('scaling', StandardScaler()),
                           ('classifier', KNeighborsClassifier(n_neighbors=7))])

pipe_svc = Pipeline(steps=[('scaling', StandardScaler()),
                           ('classifier', SVC(max_iter=7000, random_state=42))])

pipe_mlp = Pipeline(steps=[('scaling', StandardScaler()),
                           ('classifier', MLPClassifier(max_iter=7000, random_state=42))])

pipes = [pipe_knn, pipe_svc, pipe_mlp]
names_pipes = [ 'KNN', 'SVM', 'MLP' ]

for i in range(len(pipes)):
    print(names_pipes[i])
    pipes[i].fit(X_train, y_train)
    y_pred = pipes[i].predict(X_test)
    print('Accuracy : ', accuracy_score(y_test, y_pred))
    print('Precision : ', precision_score(y_test, y_pred, average='binary'))

```

1s completed at 11:49 AM

```

for i in range(len(pipes)):
    print(names_pipes[i])
    pipes[i].fit(X_train, y_train)
    y_pred = pipes[i].predict(X_test)
    print('Accuracy : ',accuracy_score(y_test, y_pred))
    print('Precision : ',precision_score(y_test, y_pred, average='binary'))
    print('Recall/ sensitivity : ',recall_score(y_test, y_pred, average='binary'))
    print('F1 : ',f1_score(y_test, y_pred, average='binary'))
    sens, spec, sup = sensitivity_specificity_support(y_test,y_pred, average='binary')
    print('Specificity : ',spec)
    print('-----')
    print('')

KNN
Accuracy : 0.92
Precision : 0.926605504587156
Recall/ sensitivity : 0.926605504587156
F1 : 0.926605504587156
Specificity : 0.9128879128879121
-----

SVM
Accuracy : 0.98
Precision : 0.981651376146789
Recall/ sensitivity : 0.981651376146789
F1 : 0.981651376146789
Specificity : 0.978021978021978
-----
  
```

1s completed at 11:49 AM

```

print('-----')
sens, spec, sup = sensitivity_specificity_support(y_test,y_pred, average='binary')
print('Specificity : ',spec)
print('-----')
print('')

KNN
Accuracy : 0.92
Precision : 0.926605504587156
Recall/ sensitivity : 0.926605504587156
F1 : 0.926605504587156
Specificity : 0.9128879128879121
-----

SVM
Accuracy : 0.98
Precision : 0.981651376146789
Recall/ sensitivity : 0.981651376146789
F1 : 0.981651376146789
Specificity : 0.978021978021978
-----

MLP
Accuracy : 0.99
Precision : 0.9908256880733946
Recall/ sensitivity : 0.9908256880733946
F1 : 0.9908256880733946
Specificity : 0.989010989010989
-----
  
```

1s completed at 11:49 AM