

EDA for Profession Prediction



total merge dataset cleaned ready

1	50414234	3.3	Software Engineer	4
2	50414279	3.1	Data Engineer	3
3	50414496	3.25	Data Analyst	3
4	50414623	3.38	Software Engineer	4
5	50414925	3.72	Mobile App Developer	4
	NPM	IPK	Klasifikasi Profesi	ALGORITMA DA

help me to do EDA from this ready and cleaned dataset, my lecturer told me that this dataset is for "Prediksi profesi pekerjaan mahasiswa menggunakan algoritma Wasserstein GAN"



python

```
import pandas as pd

# Load the dataset
file_path = '/mnt/data/total_merge_dataset_cleaned_ready.csv'
dataset = pd.read_csv(file_path)

# Display the first few rows to understand the structure of the dataset
dataset.head()
```

Result

```

      NPM    IPK  Klasifikasi Profesi  ALGORITMA DAN PEMROGRAMAN 1 \
0  50414234  3.30      Software Engineer                        4
1  50414279  3.10      Data Engineer                          3
2  50414496  3.25      Data Analyst                           3
3  50414623  3.38      Software Engineer                       4
4  50414925  3.72  Mobile App Developer                       4

      ALGORITMA DAN PEMROGRAMAN 2  ALGORITMA DAN PEMROGRAMAN 3 \
0                                3                                4
1                                3                                3
2                                3                                3
3                                3                                3
4                                3                                4

      ALGORITMA DAN PENGOLAHAN PARALEL  ALGORITMA DEEP LEARNING \
0                                    3                                4
1                                    3                                2
2                                    3                                4
3                                    3                                3
```

- ``NPM``: Student identification.
- ``IPK``: Student GPA.
- ``Klasifikasi Profesi``: The predicted profession (e.g., Software Engineer, Data Engineer).
- Various columns representing subject names, such as ``ALGORITMA DAN PEMROGRAMAN 1``, ``BAHASA INDONESIA``, etc., which reflect the students' scores in each subject.

1. **Basic statistics** of the numerical columns (such as IPK and scores in various subjects).
2. **Distribution of predicted professions.**
3. **Correlation** between the IPK and predicted professions.
4. **Correlation heatmap** for the scores across subjects.
5. **IPK distribution** for different professions.

```
python

# Step 1: Basic statistics of the numerical columns
basic_stats = dataset.describe()

# Step 2: Distribution of predicted professions
profession_distribution = dataset['Klasifikasi Profesi'].value_counts()

# Display the results to the user
import ace_tools as tools; tools.display_dataframe_to_user(name="Basic Statistics",
dataframe=basic_stats)

profession_distribution
```

Result

```
Software Engineer      42
Software Quality Assurance 32
Mobile App Developer  16
Full Stack Developer  15
Data Analyst          12
Back End Developer    11
Data Engineer         7
Graphic Designer      6
UI/UX Designer        5
Web Developer         5
IT Entrepreneur       5
Front End Developer   5
IT Project Manager    4
Network Engineer      4
Network Administrator 3
Data Scientist        3
Business Analyst      3
Database Administrator 3
Lecturer             3
IT Consultant         3
IT Lecturer           2
DevOps Engineer       2
Network Security Engineer 1
Cyber Security Manager 1
System Analyst        1
Software Quality Engineer 1
Cyber Security Analyst 1
Name: Klasifikasi Profesi, dtype: int64
```

Basic Statistics

1	count	196.0	196.0	196.0
2	mean	3.486581632653061	3.5051020408163267	3.3520408163265305
3	std	0.2501016172337992	0.5406306761513019	0.4997906172337992
4	min	2.82	2.0	2.0
5	25%	3.3	3.0	3.0
		IPK	ALGORITMA DAN PE	ALGORITMA

Here are some key insights from the initial EDA:

- Basic Statistics:** A summary of the numerical columns is displayed for you, providing a quick overview of central tendencies like mean, standard deviation, and ranges of scores across different subjects.
- Distribution of Predicted Professions:** The top 3 professions are:
 - Software Engineer (42 students)
 - Software Quality Assurance (32 students)
 - Mobile App Developer (16 students)

Would you like to explore correlations between student IPK, scores in various subjects, or move on to other specific analyses?



Can you show a correlation heatmap?



python

```
import matplotlib.pyplot as plt
import seaborn as sns

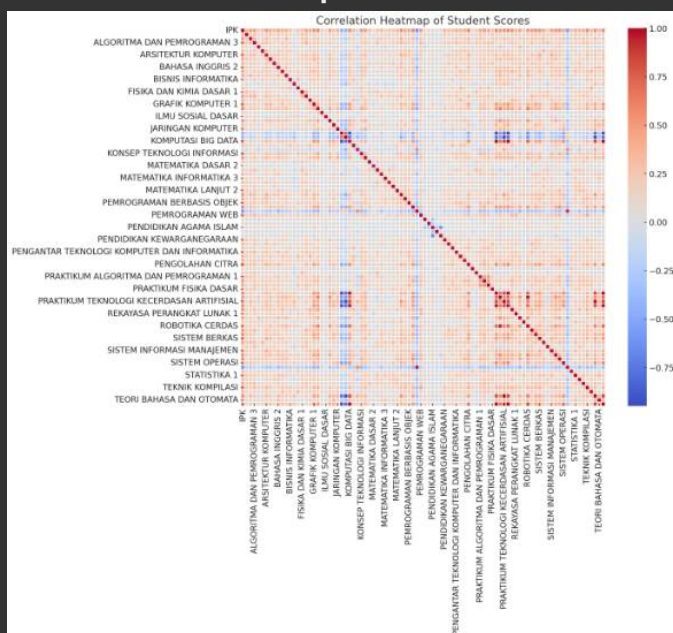
# Calculate the correlation matrix for numerical columns
correlation_matrix = dataset.iloc[:, 1:].corr()

# Plot the heatmap
plt.figure(figsize=(12, 10))
sns.heatmap(correlation_matrix, cmap='coolwarm', annot=False, fmt='.2f', linewidths=0.5)
plt.title('Correlation Heatmap of Student Scores')
plt.show()
```

STDOUT/STDERR

```
/tmp/ipykernel_56/3487444602.py:5: FutureWarning: The default value of numeric_only in DataFrame.corr is de
correlation_matrix = dataset.iloc[:, 1:].corr()
```

Correlation Heatmap of Student Scores



Here is the correlation heatmap displaying the relationships between the scores of various subjects in the dataset. Darker red indicates higher positive correlations, while darker blue represents stronger negative correlations. This visual helps in understanding how performance in different subjects is interrelated.

If you'd like, I can focus on specific correlations or areas of interest within the heatmap. Let me know how you'd like to proceed!



Can you analyze the top correlations?



python

```
# Let's find the top correlations by unstacking the correlation matrix
top_correlations =
correlation_matrix.unstack().sort_values(ascending=False).drop_duplicates()

# Filter the top correlations (above 0.8) to find significant relationships
significant_correlations = top_correlations[(top_correlations > 0.8) & (top_correlations <
1.0)]

significant_correlations.head(10)
```

Result

KOMPUTASI BIG DATA	TERAPAN TEORI GRAF	0.998408
PRAKTIKUM KOMPUTASI BIG DATA	KOMPUTASI BIG DATA	0.997686
TERAPAN TEORI GRAF	PRAKTIKUM KOMPUTASI BIG DATA	0.996076
PRAKTIKUM KOMPUTASI BIG DATA	PRAKTIKUM TERAPAN TEORI GRAF	0.995986
PRAKTIKUM TERAPAN TEORI GRAF	KOMPUTASI BIG DATA	0.995530
TEKNOLOGI KECERDASAN ARTIFISIAL	PRAKTIKUM TEKNOLOGI KECERDASAN ARTIFISIAL	0.994154
TERAPAN TEORI GRAF	PRAKTIKUM TERAPAN TEORI GRAF	0.993900
PRAKTIKUM ROBOTIKA CERDAS	ROBOTIKA CERDAS	0.974358
PEMROGRAMAN MULTIMEDIA	SISTEM TERDISTRIBUSI	0.961037
PRAKTIKUM TEKNOLOGI KECERDASAN ARTIFISIAL	TERAPAN TEORI GRAF	0.858009

dtype: float64

Here are the top significant correlations (above 0.8) between different subjects in the dataset:

1. **KOMPUTASI BIG DATA** and **TERAPAN TEORI GRAF** have an extremely high correlation of 0.998, indicating that students who perform well in one of these subjects tend to perform equally well in the other.
2. **PRAKTIKUM KOMPUTASI BIG DATA** is strongly correlated with **KOMPUTASI BIG DATA** (0.997) and **TERAPAN TEORI GRAF** (0.996), suggesting a similar trend in performance across these practical and theoretical subjects.
3. **TEKNOLOGI KECERDASAN ARTIFISIAL** and its practical component **PRAKTIKUM TEKNOLOGI KECERDASAN ARTIFISIAL** also show a strong correlation (0.994), which is typical for subjects with both theoretical and practical assessments.
4. Other significant correlations include **PRAKTIKUM ROBOTIKA CERDAS** and **ROBOTIKA CERDAS** (0.974), indicating similar performance trends in robotics courses.

These strong correlations suggest that certain subject groups (like big data, artificial intelligence, and robotics) show consistency in student performance across related practical and theoretical components.

Would you like to investigate correlations between specific subjects or any other area?



Can you find correlations with IPK?



python

```
# Find correlations specifically with IPK (the second column in the dataset)
ipk_correlations = correlation_matrix['IPK'].sort_values(ascending=False)

# Filter the top correlations with IPK (above 0.3 for notable positive correlations)
significant_ipk_correlations = ipk_correlations[ipk_correlations > 0.3]

significant_ipk_correlations.head(10)
```

Result

IPK	1.000000
PENGOLAHAN CITRA	0.654294
GRAFIK KOMPUTER 2	0.625656
PRAKTIKUM TERAPAN TEORI GRAF	0.605573
PRAKTIKUM KOMPUTASI BIG DATA	0.600570
KOMPUTASI BIG DATA	0.597063
PEMROGRAMAN JARINGAN	0.595664
TEKNIK KOMPILASI	0.593537
TERAPAN TEORI GRAF	0.593407
SISTEM BASIS DATA 2	0.590606

Name: IPK, dtype: float64

The subjects that show the strongest positive correlations with **IPK** (Grade Point Average) are:

1. **PENGOLAHAN CITRA** (Image Processing): 0.65
2. **GRAFIK KOMPUTER 2** (Computer Graphics 2): 0.63
3. **PRAKTIKUM TERAPAN TEORI GRAF** (Practical Graph Theory): 0.61
4. **PRAKTIKUM KOMPUTASI BIG DATA** (Practical Big Data Computing): 0.60
5. **KOMPUTASI BIG DATA** (Big Data Computing): 0.60

These subjects, particularly those related to graphics, big data, and network programming, have the strongest relationship with a student's overall GPA. This suggests that strong performance in these technical subjects is often reflected in higher overall academic achievement.

Would you like further analysis or deeper exploration of these findings?



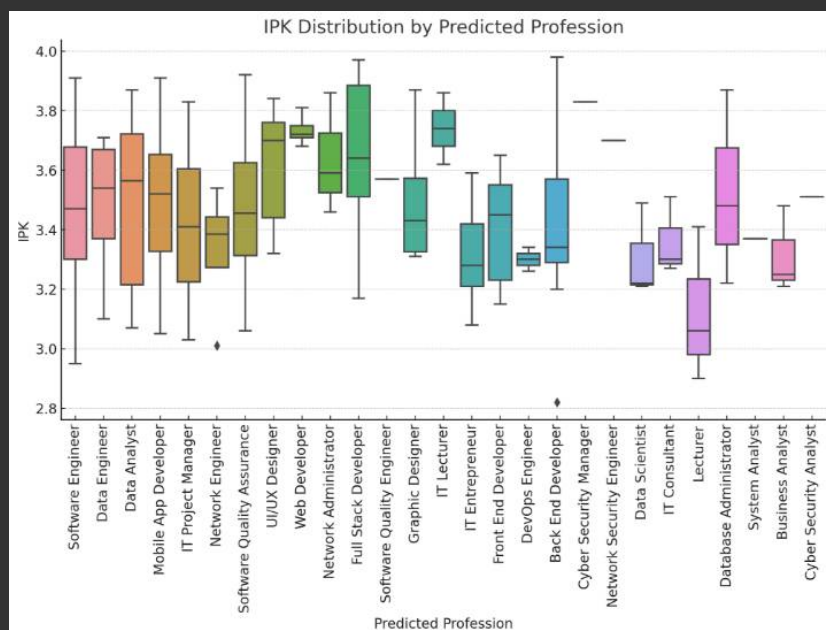
Show IPK distribution by profession.



python

```
# Create a boxplot to show the distribution of IPK by profession
plt.figure(figsize=(12, 6))
sns.boxplot(x='Klasifikasi Profesi', y='IPK', data=dataset)
plt.xticks(rotation=90)
plt.title('IPK Distribution by Predicted Profession')
plt.ylabel('IPK')
plt.xlabel('Predicted Profession')
plt.show()
```

IPK Distribution by Predicted Profession



Here is the boxplot showing the **IPK distribution** for each predicted profession. The boxplot highlights the range, median, and spread of IPK for students classified into different professions. You can observe which professions tend to have students with higher or more varied GPAs. If you'd like a more detailed breakdown or further analysis, feel free to let me know!



my lecturer told me that the prediction it has to be using Wasserstein GAN



Wasserstein GAN (WGAN) is an advanced type of Generative Adversarial Network (GAN) that improves the training stability of GANs by using a different distance metric (Wasserstein distance) to measure how close the generated data distribution is to the real data distribution. WGANs are particularly suited for tasks like this where precise and stable training is necessary.

To implement a **Wasserstein GAN** for predicting student professions, here's a general approach:

Steps to Use WGAN for Profession Prediction:

1. **Preprocessing:** Ensure that the dataset is properly encoded, especially the categorical features (like profession).
2. **Define WGAN architecture:**
 - o **Generator:** This network will create synthetic student profiles (with features like IPK and scores).
 - o **Discriminator:** This network will differentiate between real student data and synthetic data generated by the Generator.
3. **Loss Function:** WGAN replaces the discriminator loss with the Wasserstein loss, which provides smoother gradients for training.
4. **Training:** Train the GAN until the generator can create realistic student profiles with corresponding predicted professions.
5. **Evaluation:** Test how well the model predicts the profession based on the generated profiles.

Would you like assistance in setting up the code for this WGAN architecture, or do you want specific help with any part of the implementation?



Can you guide me in implementing WGAN in tensorflow?



Certainly! Here's a step-by-step guide to implementing a Wasserstein GAN (WGAN) in TensorFlow for your profession prediction task.

Prerequisites:

1. **TensorFlow:** Install via ``pip install tensorflow``.
2. **Pandas** and **Numpy:** For data handling.

Step-by-Step Implementation:

Step 1: Import Libraries

python

```
import tensorflow as tf
from tensorflow.keras import layers
import numpy as np
import pandas as pd
```

Step 2: Prepare the Data

Encode categorical data (like profession) and scale numerical data.

python

```
# Example: One-hot encode the profession
from sklearn.preprocessing import LabelEncoder, StandardScaler

# Assuming 'Klasifikasi Profesi' is your target label (profession)
label_encoder = LabelEncoder()
dataset['Profesi_Encoded'] = label_encoder.fit_transform(dataset['Klasifikasi Profesi'])

# Select the features for prediction (e.g., IPK and scores)
X = dataset.drop(['NPM', 'Klasifikasi Profesi', 'Profesi_Encoded'], axis=1)

# Standardize the numerical features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Convert data to TensorFlow datasets
batch_size = 64
train_dataset =
tf.data.Dataset.from_tensor_slices(X_scaled).shuffle(buffer_size=1024).batch(batch_size)
```


Step 3: Define the Generator and Discriminator Models

- **Generator:** Creates synthetic student profiles.
- **Discriminator:** Classifies profiles as real or fake.

python

```
# Generator model
def make_generator_model():
    model = tf.keras.Sequential()
    model.add(layers.Dense(128, input_shape=(100,))) # Input shape will be the noise vector
    model.add(layers.LeakyReLU())
    model.add(layers.BatchNormalization())

    model.add(layers.Dense(256))
    model.add(layers.LeakyReLU())
    model.add(layers.BatchNormalization())

    model.add(layers.Dense(X.shape[1])) # Output shape must match the feature size of X
    return model

# Discriminator model
def make_discriminator_model():
    model = tf.keras.Sequential()
    model.add(layers.Dense(256, input_shape=(X.shape[1],))) # Input shape matches the feature size of X
    model.add(layers.LeakyReLU())

    model.add(layers.Dense(128))
    model.add(layers.LeakyReLU())

    model.add(layers.Dense(1)) # Output 1 for Wasserstein Loss
    return model
```

Step 4: Define Loss Functions and Optimizers

In WGAN, we use Wasserstein Loss instead of Binary Cross-Entropy.

python

```
# Wasserstein Loss
def wasserstein_loss(y_true, y_pred):
    return tf.reduce_mean(y_true * y_pred)

generator_optimizer = tf.keras.optimizers.RMSprop(learning_rate=0.00005)
discriminator_optimizer = tf.keras.optimizers.RMSprop(learning_rate=0.00005)
```

Step 5: Create the Training Loop

Here, we alternate between training the discriminator and the generator.

python

```
import os
import time

epochs = 10000
noise_dim = 100
num_examples_to_generate = 16
```

```
# Create random seed for generating profiles after training
seed = tf.random.normal([num_examples_to_generate, noise_dim])

# WGAN requires clipping of the discriminator weights
def train_step(real_data):
    noise = tf.random.normal([batch_size, noise_dim])

    with tf.GradientTape() as disc_tape, tf.GradientTape() as gen_tape:
        generated_data = generator(noise, training=True)

        real_output = discriminator(real_data, training=True)
        fake_output = discriminator(generated_data, training=True)

        # Discriminator loss
        disc_loss = wasserstein_loss(tf.ones_like(real_output), real_output) +
wasserstein_loss(tf.ones_like(fake_output) * -1, fake_output)

        # Generator loss
        gen_loss = wasserstein_loss(tf.ones_like(fake_output), fake_output)

        # Calculate gradients
        gradients_of_discriminator = disc_tape.gradient(disc_loss,
discriminator.trainable_variables)
        gradients_of_generator = gen_tape.gradient(gen_loss, generator.trainable_variables)

        # Apply the gradients
        discriminator_optimizer.apply_gradients(zip(gradients_of_discriminator,
discriminator.trainable_variables))
        generator_optimizer.apply_gradients(zip(gradients_of_generator,
generator.trainable_variables))

        # Clip discriminator weights (required by WGAN)
        for var in discriminator.trainable_variables:
            var.assign(tf.clip_by_value(var, -0.01, 0.01))

# Training loop
def train(dataset, epochs):
    for epoch in range(epochs):
        start = time.time()

        for real_data in dataset:
            train_step(real_data)

        # Produce images for the GIF as we go
        if epoch % 100 == 0:
            print(f'Epoch {epoch}, Time {time.time() - start}')

# Instantiate the models
generator = make_generator_model()
discriminator = make_discriminator_model()

# Train the model
train(train_dataset, epochs)
```

Step 6: Generate Predictions

After training, you can use the generator to create new synthetic student profiles based on the learned distribution.

python

```
# Generate new student profiles
def generate_new_profiles(generator, num_profiles):
    noise = tf.random.normal([num_profiles, noise_dim])
    generated_profiles = generator(noise)

    # Inverse transform the scaled data
```

```
generated_profiles = scaler.inverse_transform(generated_profiles)

return generated_profiles

# Example: Generate 5 new profiles
new_profiles = generate_new_profiles(generator, 5)
print(new_profiles)
```

Key Points:

- **WGAN training** involves clipping the discriminator's weights after each update to ensure stable training.
- The **Wasserstein loss** measures how well the generator is capturing the true data distribution.
- The **generator** creates synthetic data, and the **discriminator** distinguishes real from fake.

This is a simplified WGAN structure tailored for student profile data, but you may need to tweak the model architecture or hyperparameters based on your specific dataset and objectives. Let me know if you need help with further tuning or testing!



what is the key differences with Vanilla GAN and Wasserstein GAN?



The key differences between **Vanilla GAN** and **Wasserstein GAN (WGAN)** are in the way they approach the training process, the loss functions they use, and the stability of their performance. Below is a detailed comparison:

1. Loss Function

- **Vanilla GAN:** Uses **binary cross-entropy** as the loss function. It aims to minimize the Jensen-Shannon (JS) divergence between the real and generated data distributions. This can lead to vanishing gradients and instability during training when the generator produces poor samples.
- **Wasserstein GAN:** Uses the **Wasserstein distance (Earth Mover's distance)** as the loss function, which provides smoother and more informative gradients. This allows the generator to keep learning, even when it's far from producing realistic data.

2. Training Stability

- **Vanilla GAN:** Prone to **mode collapse** and **instability** during training. Mode collapse happens when the generator produces a very limited set of outputs, failing to capture the diversity of the real data.
- **Wasserstein GAN:** More **stable** during training. The use of Wasserstein distance ensures better gradient flow, even when the discriminator is well-trained, leading to more consistent improvements in the generator.

3. Discriminator (Critic)

- **Vanilla GAN:** The discriminator is a **classifier** that outputs a probability (0 for fake, 1 for real). It uses the sigmoid activation function to output the likelihood of the data being real or generated.

- **Wasserstein GAN:** The discriminator is referred to as a **critic**. Instead of outputting a probability, it produces a real-valued score that represents how "real" a sample is. The critic does not classify but provides a continuous measure of realness.

4. Gradient Clipping

- **Vanilla GAN:** Does not require any special regularization like weight clipping.
- **Wasserstein GAN:** Uses **weight clipping** to keep the critic within a Lipschitz constraint, which is necessary for computing the Wasserstein distance. This ensures the gradient remains stable during training.

5. Mode Collapse

- **Vanilla GAN:** Prone to mode collapse, where the generator learns to produce only a few variations of the data, ignoring diversity.
- **Wasserstein GAN:** Significantly **reduces mode collapse**, because the Wasserstein distance provides useful gradients even when the real and generated distributions have little overlap, encouraging diversity in generated samples.

6. Training Procedure

- **Vanilla GAN:** Alternates between updating the generator and discriminator a fixed number of times. Both the generator and discriminator are updated in each iteration.
- **Wasserstein GAN:** Typically, the **critic is updated more frequently** than the generator. This ensures that the Wasserstein distance is well-estimated before the generator is updated.

7. Convergence and Performance

- **Vanilla GAN:** Hard to determine when the training has converged. It can oscillate or collapse before reaching the desired solution.
- **Wasserstein GAN:** Provides a more interpretable loss (the Wasserstein distance), which correlates with the quality of the generator. The training progress is smoother, and convergence is easier to track.

Summary of Key Differences:

Feature	Vanilla GAN	Wasserstein GAN (WGAN)
Loss Function	Binary Cross-Entropy (JS Divergence)	Wasserstein Loss (Earth Mover's Distance)
Discriminator Output	Probability (Real or Fake)	Continuous Score (Realness)
Mode Collapse	Prone to Mode Collapse	Reduces Mode Collapse
Gradient Flow	Can suffer from vanishing gradients	Improved gradient flow
Stability	Training can be unstable	More stable training
Weight Clipping	Not required	Required to maintain Lipschitz constraint
Training Procedure	Equal updates for generator and discriminator	Critic updated more often

In summary, **Wasserstein GAN** improves upon Vanilla GAN by addressing instability and mode collapse, providing a more reliable and interpretable training process, and utilizing a loss function (Wasserstein distance) that gives more meaningful feedback to the generator throughout the learning process.

ChatGPT can make mistakes. Check important info.