Nama Dosen : Teguh Iman Hermanto, M.Kom

Mata Kuliah : Machine Learning 1

Pembahasan : Exploratory Data Analysis (EDA)

Pokok Pemb : - Mengenal Library ML

- Mengenal Statistik Deskriptif

- Mengenal EDA Data Numerik

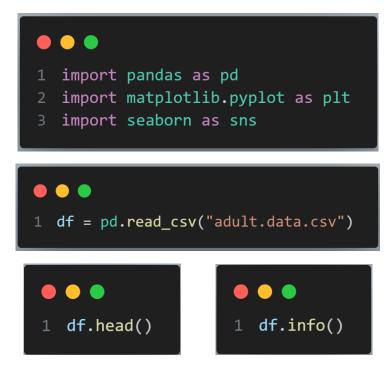
- Mengenal EDA Data Kategori

- Mengenal EDA data Multivariabel

- Mengenal Univariate Analysis

- Mengenal Bivariate Analysis

1. Mengenal Library Machine Learning



2. Mengenal Statistik Deskriptif

```
# Menghitung statistik deskriptif untuk kolom numerik
numerical_stats = df.describe()
print("Statistik Deskriptif untuk Kolom Numerik:\n", numerical_stats)
```

```
# Menghitung frekuensi untuk kolom kategorikal
categorical_stats = df.select_dtypes(include=['category']).apply(lambda x: x.value_counts(normalize=True) * 100)
print("\nPersentase Frekuensi untuk Kolom Kategorikal:\n", categorical_stats)
```

```
# Menampilkan jumlah missing values pada tiap kolom
missing_values = df.isnull().sum()
print("\nJumlah Missing Values pada tiap kolom:\n", missing_values)
```

3. Mengenal Bentuk EDA dengan Data Numerik

```
# Histogram
plt.figure(figsize=(8, 6))
sns.histplot(df['age'], kde=True)
plt.title('Histogram of Age')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()
```

```
# Box Plot
plt.figure(figsize=(8, 6))
sns.boxplot(x='sex', y='hours-per-week', data=df)
plt.title('Box Plot of Hours per Week by Sex')
plt.xlabel('Sex')
plt.ylabel('Hours per Week')
plt.show()
```

```
# Scatter Plot
plt.figure(figsize=(8, 6))
sns.scatterplot(x='age', y='capital-gain', data=df)
plt.title('Scatter Plot of Age vs Capital Gain')
plt.xlabel('Age')
plt.ylabel('Capital Gain')
plt.show()
```

4. Mengenal Bentuk EDA dengan Data Kategori

```
# Line Plot (Example: Average capital gain by age)
avg_capital_gain_by_age = df.groupby('age')['capital-gain'].mean()
plt.figure(figsize=(8, 6))

plt.plot(avg_capital_gain_by_age.index, avg_capital_gain_by_age.values)
plt.title('Line Plot of Average Capital Gain by Age')

plt.xlabel('Age')
plt.ylabel('Average Capital Gain')
plt.show()
```

```
# Bar Chart
# Contoh: Melihat jumlah orang berdasarkan tingkat pendidikan
deducation_counts = df['education'].value_counts()
plt.figure(figsize=(10, 5))
plt.bar(education_counts.index, education_counts.values)
plt.xlabel('Tingkat Pendidikan')
plt.ylabel('Jumlah Orang')
plt.title('Jumlah Orang Berdasarkan Tingkat Pendidikan')
plt.xticks(rotation=45, ha='right')
plt.show()
```

```
# Pie Chart
# Pie Chart
# Contoh: Melihat persentase jenis kelamin
sex_counts = df['sex'].value_counts()

# plt.figure(figsize=(6, 6))
plt.pie(sex_counts.values, labels=sex_counts.index, autopct='%1.1f%%', startangle=90)
plt.title('Persentase Jenis Kelamin')
plt.show()
```

```
# Donut Chart
# Donut Chart
# Contoh: Melihat persentase ras
race_counts = df['race'].value_counts()
plt.figure(figsize=(6, 6))
plt.pie(race_counts.values, labels=race_counts.index, autopct='%1.1f%%', startangle=90, wedgeprops=dict(width=0.4))
plt.title('Persentase Ras')
plt.show()
```

5. Mengenal Bentuk EDA dengan Data Multivariabel

```
# Create a heatmap
# Select numerical features for the heatmap
numerical_features = ['age', 'fnlwgt', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']
heatmap_data = df[numerical_features].corr()

plt.figure(figsize=(10, 8))
sns.heatmap(heatmap_data, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Heatmap')
plt.show()
```

```
# Create a pair plot
2 sns.pairplot(df[['age', 'education-num', 'hours-per-week', 'capital-gain', 'capital-loss']])
3 plt.show()
```

```
# Create a bubble chart
# Create a bubble chart
bubble_data = df[['age', 'hours-per-week', 'capital-gain']]

plt.figure(figsize=(10, 8))
plt.scatter(x='age', y='hours-per-week', s='capital-gain', data=bubble_data, alpha=0.5)
plt.xlabel('Age')
plt.ylabel('Hours per Week')
plt.title('Bubble Chart: Age, Hours per Week, and Capital Gain')
plt.show()
```

6. Mengenal Univarate Analysis

```
# Univariate Analysis for Age
plt.figure(figsize=(8, 6))
sns.histplot(df['age'], kde=True)
plt.title('Distribution of Age')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()
```

```
# Univariate Analysis for Education (Categorical)
plt.figure(figsize=(10, 6))
sns.countplot(x='education', data=df)
plt.title('Count of Individuals by Education Level')
plt.xlabel('Education Level')
plt.ylabel('Count')
plt.xticks(rotation=45, ha='right')
plt.show()
```

```
# Univariate Analysis for Salary (Categorical)
plt.figure(figsize=(6, 6))
sns.countplot(x='salary', data=df)
plt.title('Count of Individuals by Salary')
plt.xlabel('Salary')
plt.ylabel('Count')
plt.show()
```

7. Mengenal Bivariate Analysis

```
1 # 1. Age vs. Salary
2 # Deskripsi: Melihat hubungan antara umur dan tingkat pendapatan.
3 # Kita mengharapkan bahwa orang yang lebih tua mungkin memiliki pendapatan yang lebih tinggi.
4 plt.figure(figsize=(8, 6))
5 sns.boxplot(x='salary', y='age', data=df)
6 plt.title('Hubungan antara Umur dan Tingkat Pendapatan')
7 plt.show()
```

```
# 2. Education vs. Salary
2 # Deskripsi: Melihat hubungan antara tingkat pendidikan dan pendapatan.
3 # Kita mengharapkan bahwa orang dengan pendidikan lebih tinggi cenderung memiliki pendapatan lebih tinggi.
4 plt.figure(figsize=(12, 6))
5 sns.countplot(x='education', hue='salary', data=df)
6 plt.title('Hubungan antara Tingkat Pendidikan dan Tingkat Pendapatan')
7 plt.xticks(rotation=45, ha='right')
8 plt.show()
```

```
# 3. Hours-per-week vs. Salary
# 3. Hours-per-week vs. Salary
# Deskripsi: Melihat hubungan antara jumlah jam kerja per minggu dan pendapatan.
# Kita mengharapkan bahwa orang yang bekerja lebih banyak jam cenderung memiliki pendapatan lebih tinggi.
plt.figure(figsize=(8, 6))
# sns.boxplot(x='salary', y='hours-per-week', data=df)
# plt.title('Hubungan antara Jam Kerja per Minggu dan Tingkat Pendapatan')
# plt.show()
```

```
1 # 4. Sex vs. Salary
2 # Deskripsi: Melihat perbedaan pendapatan antara pria dan wanita.
3 plt.figure(figsize=(8, 6))
4 sns.countplot(x='sex', hue='salary', data=df)
5 plt.title('Perbedaan Tingkat Pendapatan Berdasarkan Jenis Kelamin')
6 plt.show()
```

```
# 5. Occupation vs. Salary
# Deskripsi: Melihat hubungan antara jenis pekerjaan dan pendapatan.
# Kita mengharapkan bahwa beberapa jenis pekerjaan memiliki pendapatan yang lebih tinggi daripada yang lain.
# plt.figure(figsize=(12, 6))
# sns.countplot(x='occupation', hue='salary', data=df)
# plt.title('Hubungan antara Jenis Pekerjaan dan Tingkat Pendapatan')
# plt.xticks(rotation=45, ha='right')
# plt.show()
```

```
# 6. Correlation Matrix

# Deskripsi: Melihat korelasi antara variabel-variabel numerik.

# Kita dapat melihat variabel mana yang memiliki korelasi kuat dengan variabel lainnya.

numeric_cols = df.select_dtypes(include=['number']).columns

correlation_matrix = df[numeric_cols].corr()

plt.figure(figsize=(10, 8))

sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")

plt.title('Correlation Matrix')

plt.show()
```

Latihan

- 1. Tampilkan Plot Education berdasarkan Salarynya
- 2. Tampilkan jumlah pendidikan Bachelors dengan gaji <=50k
- 3. Tampilkan data occupation dengan gaji diatas 50k
- 4. Tampilkan Plot occupation dengan gaji diatas 50k
- 5. Tampilkan rata-rata hours per week untuk occupation sales?
- 6. Tampilkan rata-rata hours per week untuk occupation Prof-Speciality?