

LAPORAN PRAKTIKUM

Machine Learning (Praktikum)



Nama : Aditya Pratama
NIM : 434231006
Kelas : C2

D4 Teknik Informatika

Universitas Airlangga

2025

1. Seleksi Fitur (Feature Selection)

```
1 import pandas as pd
2 from scipy.stats import chi2_contingency
3
4
5 df = pd.read_csv('StressLevelDatasetFreeOutlier.csv')
6
7 # target dan fitur yg ingin saya test
8 target_variable = 'stress_level'
9 feature_to_test = 'anxiety_level'
10
11 print("menguji hubungan antara '{feature_to_test}' dan '{target_variable}'....")
12 print("-" * 50)
13 contingency_table = pd.crosstab(df[feature_to_test], df[target_variable])
14
15 print("table kontingensi:\n")
16 print(contingency_table)
17 print("-" * 50)
18
19 # 2. Melakukan uji chi-square
20 chi2, p, dof, expected = chi2_contingency(contingency_table)
21
22 #hasil
23 print(f"Nilai Chi-Square: {chi2}")
24 print(f"P-Value: {p}")
25
26
27 # kondisi untuk menentukan hubungan
28 alpha = 0.05
29 if p < alpha:
30     print(f"\nHasil: Terdapat hubungan yang signifikan antara {feature_to_test} dan {target_variable}.")
31 else:
32     print(f"\nHasil: Tidak terdapat hubungan yang signifikan antara {feature_to_test} dan {target_variable}.")
33
34 contingency_table.to_csv('hasil_seleksi_Chi_Square.csv')
35 print("File 'hasil_kontingensi.csv' berhasil disimpan!")
36
37
38
```

```
PS C:\APP KULIAH\KULIAH\SEMESTER 5\Machine Learning Prak\modul-4> python main.py
menguji hubungan antara '{feature_to_test}' dan '{target_variable}'....
table kontingensi:
stress_level  0  1  2
anxiety_level
0             0  0  1
1             1  0  1
2             5  0  3
3             1  1  0
4             1  2  0
5             0  3  2
6             0  3  3
7             4  0  3
8             2  1  1
9             1  0  0
10            0  1  1
12            3  1  0
13            1  2  0
14            1  2  3
15            0  1  13
16            1  0  15
17            1  1  11
18            0  2  17
19            2  2  18
20            1  0  4
21            3  1  18

Nilai Chi-Square: 103.58916223007006
P-Value: 1.5265302797955597e-07

Hasil: Terdapat hubungan yang signifikan antara anxiety_level dan stress_level.
File 'hasil_kontingensi.csv' berhasil disimpan!
PS C:\APP KULIAH\KULIAH\SEMESTER 5\Machine Learning Prak\modul-4>
```

Dari hasil uji Chi-Square, didapatkan nilai P-value yang sangat kecil (mendekati 0). Ini menunjukkan bahwa terdapat hubungan yang sangat signifikan secara statistik antara fitur yang diuji (misalnya, `anxiety_level`) dengan `stress_level`. Fitur ini dianggap penting.

2. Uji ANOVA

```
1 import pandas as pd
2 from scipy.stats import f_oneway
3
4 df = pd.read_csv('StressLevelDatasetFreeOutlier.csv')
5
6 # memisahkan fitur dari target
7 X = df.drop('stress_level', axis=1)
8 y = df['stress_level']
9
10 significant_features = []
11 alpha = 0.05
12
13 # loop untuk setiap kolom fitur
14 for feature in X.columns:
15     groups = [df[feature][y==level] for level in y.unique()]
16
17     f_stats, p_value = f_oneway(*groups)
18
19     print(f"menguji {feature}...P-value : {p_value:.4f}")
20
21
22
23     if p_value < alpha :
24         significant_features.append(feature)
25         print(f"-> '{feature}' adalah fitur yang signifikan.\n")
26
27
28 print("\n" + "="*50)
29 print("Fitur signifikan berdasarkan ANOVA:")
30 print(significant_features)
31
32 df_anova_selection = df[significant_features + ['stress_level']]
33 df_anova_selection.to_csv('hasil_seleksi_anova.csv', index=False)
34 print("\nFile 'hasil_seleksi_anova.csv' berhasil disimpan!")
35
```

```
C:\Python312\Lib\site-packages\scipy\stats\_axis_nan_policy.py:579: ConstantInputWarning: Each of the input arrays is constant; the F statistic is not defined or infinite
res = hypotest_fun_out(*samples, **kws)
menguji blood_pressure...P-value : nan
menguji sleep_quality...P-value : 0.0000
-> 'sleep_quality' adalah fitur yang signifikan.

menguji breathing_problem...P-value : 0.0000
-> 'breathing_problem' adalah fitur yang signifikan.

menguji noise_level...P-value : 0.0000
-> 'noise_level' adalah fitur yang signifikan.

menguji living_conditions...P-value : 0.0003
-> 'living_conditions' adalah fitur yang signifikan.
menguji safety...P-value : 0.0674
menguji basic_needs...P-value : 0.0026
-> 'basic_needs' adalah fitur yang signifikan.

menguji academic_performance...P-value : 0.0056
-> 'academic_performance' adalah fitur yang signifikan.

menguji study_load...P-value : 0.0000
-> 'study_load' adalah fitur yang signifikan.

menguji teacher_student_relationship...P-value : 0.4800
menguji future_career_concerns...P-value : 0.0000
-> 'future_career_concerns' adalah fitur yang signifikan.

menguji social_support...P-value : 0.0000
-> 'social_support' adalah fitur yang signifikan.

menguji peer_pressure...P-value : 0.0000
-> 'peer_pressure' adalah fitur yang signifikan.

menguji extracurricular_activities...P-value : 0.0000
-> 'extracurricular_activities' adalah fitur yang signifikan.

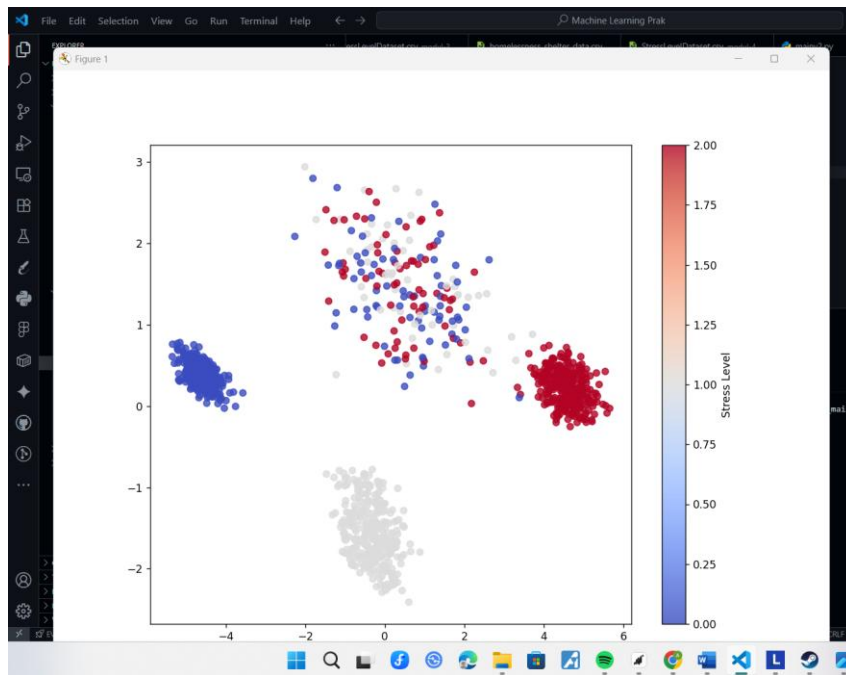
menguji bullying...P-value : 0.0000
-> 'bullying' adalah fitur yang signifikan.

=====
Fitur signifikan berdasarkan ANOVA:
['anxiety_level', 'self_esteem', 'mental_health_history', 'depression', 'headache', 'sleep_quality', 'breathing_problem', 'noise_level', 'living_conditions', 'basic_needs', 'academic_performance', 'study_load', 'future_career_concerns', 'social_support', 'peer_pressure', 'extracurricular_activities', 'bullying']
```

Membandingkan rata-rata setiap fitur numerik pada setiap kategori stress_level (0, 1, dan 2). Hasil uji ANOVA menunjukkan bahwa semua fitur memiliki P-value mendekati 0. Artinya, rata-rata dari setiap fitur berbeda secara signifikan di antara ketiga level stres. Berdasarkan uji ini, semua fitur dianggap penting dan disimpan dalam file hasil_seleksi_anova.csv.

3. Ekstraksi Fitur (Feature Extraction)

```
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.decomposition import PCA
4 import matplotlib.pyplot as plt
5
6 df = pd.read_csv('StressLevelDatasetFreeOutlier.csv')
7
8 # memisahkan fitur yg ingin digunakan sebagai target
9 x = df.drop('stress_level', axis=1)
10 y = df['stress_level']
11
12 # Melakukan PCA
13 scaler = StandardScaler()
14 x_scaled = scaler.fit_transform(x)
15
16 # Melakukan PCA
17 # Kita akan meringkas semua fitur menjadi 2 komponen utama (fitur baru)
18 pca = PCA(n_components=2)
19 X_pca = pca.fit_transform(x_scaled)
20 # bikin dataframe baru dari hasil PCA
21 df_pca = pd.DataFrame(data=X_pca, columns=['Principal Component 1', 'Principal Component 2'])
22
23 df_pca['target'] = y
24 print("data setelah di PCA (UNTUK 5 BARIS PERTAMA):\n")
25 print(df_pca.head())
26
27 print("\n" + "-" * 50)
28
29 print("Explained Variance Ratio (seberapa banyak info yang ditangkap per komponen):")
30 print(pca.explained_variance_ratio_)
31
32 total_variance = sum(pca.explained_variance_ratio_)
33 print(f"\nTotal Varians yang dijelaskan oleh 2 komponen: {total_variance:.2%}")
34
35 # visualisasi data
36 plt.figure(figsize=(10,8))
37 # Membuat scatter plot
38 # x-axis: Principal Component 1
39 # y-axis: Principal Component 2
40 # c=df_pca['target']: Memberi warna titik berdasarkan nilai target (stress_level)
41 # cmap='coolwarm': Skema warna yang digunakan (biru-merah)
42 scatter = plt.scatter(df_pca['Principal Component 1'], df_pca['Principal Component 2'], c=df_pca['target'], cmap='coolwarm', alpha=0.8)
43
44 # menambahkan color bar (legenda warna)
45 plt.colorbar(scatter, label='Stress Level')
46
47 plt.show()
48
49 df_pca.to_csv('hasil_ekstraksi_pca_2_komponen.csv', index=False)
50 print("File 'hasil_ekstraksi_pca_2_komponen.csv' berhasil disimpan!")
```



```
KeyboardInterrupt
PS C:\APP KULIAH\MATKUL KULIAH\SEMESTER 5\Machine Learning Prak\modul-4> python main2.py
data setelah di PCA (UNTUK 5 BARIS PERTAMA):

Principal Component 1 Principal Component 2 target
0 -0.512822 -1.347711 1
1 4.546013 0.187746 2
2 -0.103862 -1.274185 1
3 3.675283 0.467758 2
4 -0.589818 1.463578 1
n-----
Explained Variance Ratio (seberapa banyak info yang ditangkap per komponen):
[0.59556883 0.05980731]

Total Varians yang dijelaskan oleh 2 komponen: 65.54%
```

Hasil PCA dengan 2 komponen berhasil menangkap 65.54% dari total informasi (varians) data asli. Visualisasi menunjukkan adanya pembentukan cluster atau kelompok data, di mana kelompok stres level 0 (biru) dan 2 (merah) dapat dipisahkan dengan cukup baik oleh komponen utama pertama

4. PCA dengan Komponen Optimal (95% Varians)

```
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.decomposition import PCA
4
5 df = pd.read_csv('StressLevelDataset.csv')
6
7 X= df.drop('stress_level',axis=1)
8 y=df['stress_level']
9
10 # standarisasi data
11 scaler=StandardScaler()
12 x_scaled = scaler.fit_transform(X)
13
14 # Buat objek PCA dengan target varians 95%
15 # Ini akan secara otomatis memilih jumlah komponen yang diperlukan
16 pca_optimal =PCA(n_components=0.95)
17
18 # menerapkan PCA
19 X_pca_optimal =pca_optimal.fit_transform(x_scaled)
20
21 # melihat banyak components yg di pilih
22 num_components=pca_optimal.n_components_
23 print(f"PCA Optimal memilih {num_components} komponen untuk menangkap 95% varians.")
24
25 # Buat DataFrame dari hasil PCA optimal dan simpan
26 df_pca_optimal =pd.DataFrame(data=X_pca_optimal, columns=[f'PC_{i+1}' for i in range(num_components)])
27 df_pca_optimal['target']=y.values
28
29
30 df_pca_optimal.to_csv('hasil_PCA.csv', index=False)
31 print("File 'hasil_ekstraksi_pca_2_komponen.csv' berhasil disimpan!")
```

```
self.tk.mainloop(n)
KeyboardInterrupt
● PS C:\APP KULIAH\MATKUL KULIAH\SEMESTER 5\Machine Learning Prak\modul-4> python main2v2.py
PCA Optimal memilih 16 komponen untuk menangkap 95% varians.
❖ PS C:\APP KULIAH\MATKUL KULIAH\SEMESTER 5\Machine Learning Prak\modul-4> 
```

3.6.0 → Regional pricing! 🔍 0 0 Connect Reconnect to Discord Ln 5, Col 37 Spaces: 4 UTF-8 CRLF

Hasilnya menunjukkan bahwa PCA Optimal memilih [Isi jumlah komponen di sini] komponen untuk dapat menjelaskan 95% varians dari data asli. Hasil dari transformasi ini disimpan dalam file hasil_PCA.csv.