<div align="center">

**Lampiran**

</div>

- **User Manual Guide**
  - Deskripsi umum:

    Aplikasi ini digunakan untuk memprediksi jenis anomali dalam lalu lintas jaringan komputer dengan memanfaatkan model *Support Vector Machine* (SVM). Pendekatan klasifikasi *multiclass* dilakukan dengan skema *One-vs-One*, sedangkan proses pelatihan model menggunakan metode optimasi *Primal Estimated sub-GrAdient SOlver for SVM* (Pegasos), yang menerapkan teknik *Stochastic Gradient Descent* (SGD) untuk menyelesaikan fungsi objektif SVM secara efisien. Kernel yang digunakan adalah *Radial Basis Function* (RBF), yang memungkinkan pemisahan kelas pada data non-linear.
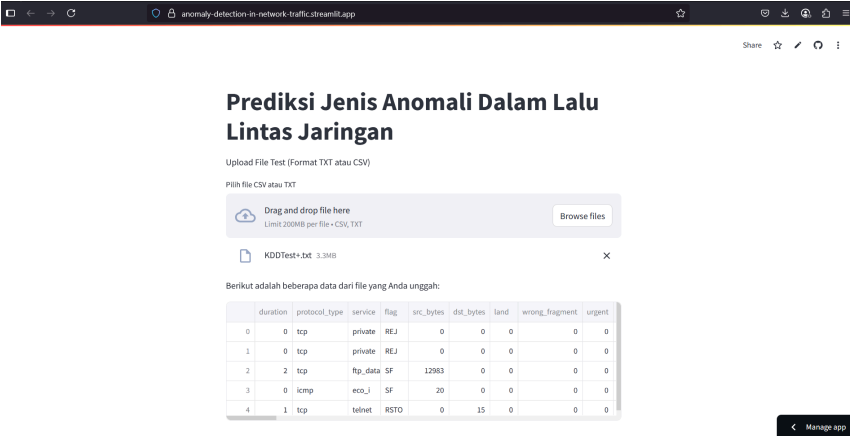  - Akses Aplikasi: https://anomaly-detection-in-network-traffic.streamlit.app
  - Fitur Utama:
    1. Upload file txt/csv
    2. Tabel hasil prediksi
    3. Visualisasi distribusi label
  - Cara Menggunakan:
    1. Download dataset KDDTest+.txt di kaggle untuk prediksi:
       https://www.kaggle.com/datasets/hassan06/nslkdd
    2. Buka url streamlit: https://anomaly-detection-in-network-traffic.streamlit.app
    3. Masukkan file KDDTest+.txt yang telah di download ke dalam streamlit
    4. Tunggu hingga hasil prediksi dan visualisasi distribusi label muncul
    5. Hasil muncul dalam bentuk tabel dan grafik
  - Contoh Tampilan:



    -

- **Listing Code**
  1. File preprosBaru.ipynb – Pemrosesan Data dan Klasifikasi Anomali
     a. Buat *header* kolom manual

```python
import pandas as pd

columns = [
    "duration", "protocol_type", "service", "flag", "src_bytes", "dst_bytes", "land",
    "wrong_fragment", "urgent", "hot", "num_failed_logins", "logged_in", "num_compromised",
    "root_shell", "su_attempted", "num_root", "num_file_creations", "num_shells",
    "num_access_files", "num_outbound_cmds", "is_host_login", "is_guest_login", "count",
    "srv_count", "serror_rate", "srv_serror_rate", "rerror_rate", "srv_rerror_rate",
    "same_srv_rate", "diff_srv_rate", "srv_diff_host_rate", "dst_host_count",
    "dst_host_srv_count", "dst_host_same_srv_rate", "dst_host_diff_srv_rate",
    "dst_host_same_src_port_rate", "dst_host_srv_diff_host_rate", "dst_host_serror_rate",
    "dst_host_srv_serror_rate", "dst_host_rerror_rate", "dst_host_srv_rerror_rate",
    "label", "difficulty"
]
```

     b. *Load* data KDDTrain+ dengan *header* kolom manual

```python
train_path = "KDDTrain+.txt"
df = pd.read_csv(train_path, names=columns)

df
```

     c. EDA
        i. Cek jumlah baris dan kolom

```
1 print(df.shape)
```

ii. Cek ringkasan informasi struktur DataFrame

```
1 df.info()
```

iii. Cek statistika deskriptif untuk kolom numerik

```
1 df.describe()
```

iv. Cek distribusi class

```
1 df['label'].value_counts(normalize=True)
```

v. Plot bar chart

```python
from matplotlib.ticker import FuncFormatter

label_counts = df['label'].value_counts()
def format_k(x, pos):
    if x >= 1000:
        return f'{x/1000:.1f}k'
    return int(x)

plt.figure(figsize=(10, 6))
bars = plt.bar(label_counts.index, label_counts.values, color='#1e77b4')

plt.title('Distribusi Jumlah Label (Bar Chart)', pad=15)
plt.xlabel('Label')
plt.ylabel('Jumlah (ribuan)')
plt.ylim(0, max(label_counts.values) * 1.1)
plt.xticks(rotation=45)

plt.gca().yaxis.set_major_formatter(FuncFormatter(format_k))

plt.tight_layout()
plt.show()
```

d. Mengelompokkan semua label menjadi 5 label

```python
1  # define attack type groups
2  dos_attacks = [
3      "back",
4      "land",
5      "neptune",
6      "pod",
7      "smurf",
8      "teardrop",
9      "apache2",
10     "udpstorm",
11     "processtable",
12     "mailbomb"
13 ]
14 probe_attacks = [
15     "satan",
16     "ipsweep",
17     "nmap",
18     "portsweep",
19     "mscan",
20     "saint"
21 ]
22 r2l_attacks = [
23     "guess_passwd",
24     "ftp_write",
25     "imap",
26     "phf",
27     "multihop",
28     "warezmaster",
29     "warezclient",
30     "spy",
31     "xlock",
32     "xsnoop",
33     "snmpguess",
34     "snmpgetattack",
35     "httptunnel",
36     "sendmail",
37     "named"
38 ]
39 u2r_attacks = [
40     "rootkit",
41     "buffer_overflow",
42     "loadmodule",
43     "perl",
44     "sqlattack",
45     "xterm",
46     "ps"
47 ]
48
49 # function to categorize attack types
50 def map_attack_type(attack):
51     if attack in dos_attacks:
52         return "DoS"
53     elif attack in probe_attacks:
54         return "Probe"
55     elif attack in r2l_attacks:
56         return "R2L"
57     elif attack in u2r_attacks:
58         return "U2R"
59     elif attack == "normal":
60         return "Normal"
61     else:
62         return "Other"
63
64 # Apply mapping function
65 df["attack_category"] = df["label"].apply(map_attack_type)
66
67 print("Attack Category Distribution (Train)")
68 print(df["attack_category"].value_counts())
69 print(f"Sum of all total attack_type from attack_category: {df['attack_category'].value_counts().sum()}")
```

e. EDA
    i. Cek korelasi matriks

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

df_numeric = df.select_dtypes(include=['number'])

correlation_matrix = df_numeric.corr()

plt.figure(figsize=(12, 10))
sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap='coolwarm', linewidths=0.5)
plt.title("Correlation Matrix (Numerical Features Only)")
plt.show()
```

ii. Encoding kolom kategorikal, lalu menghitung dan mengurutkan korelasi di kolom attack_category

```
import pandas as pd
from sklearn.preprocessing import LabelEncoder

df_encoded = df.copy()

for col in df_encoded.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    df_encoded[col] = le.fit_transform(df_encoded[col])

correlation_matrix = df_encoded.corr()

correlation_with_label = correlation_matrix['attack_category'].drop('attack_category')
correlation_sorted =
  correlation_with_label.reindex(correlation_with_label.abs().sort_values(ascending=False).index)

print(correlation_sorted)
```

iii. Cek kolom dengan nilai yang berisi 0 atau 0.00

```
zero_columns = df.columns[(df == 0).all()]

print("Kolom yang isinya 0 atau 0.00 semua:")
print(list(zero_columns))
```

iv. Cek Nan Values

```
df.isna().sum()
```

v. Cek jumlah kemunculan dari setiap kategori serangan di kolom attack_category

```
df['attack_category'].value_counts()
```

f. Encode

```
Encode

1 from sklearn.preprocessing import LabelEncoder
2
3 for col in df.select_dtypes(include=['object']).columns:
4     le = LabelEncoder()
5     df[col] = le.fit_transform(df[col])
```

g. Pisahkan fitur dan target

```
Pisahkan Fitur dan Target

1 X = df.drop(["label","attack_category", "num_outbound_cmds"], axis=1)
2 y = df["attack_category"]
```

h. Split data

```
Train Test Split

1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

i. Standarisasi

```
Standarisasi

1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_test_scaled = scaler.transform(X_test)
```

j. Handling Imbalance Data

```
Handling Imbalance Data

1 from imblearn.under_sampling import RandomUnderSampler
2 rus = RandomUnderSampler(random_state=42,sampling_strategy='auto')
3 X_train_under, y_train_under = rus.fit_resample(X_train_scaled, y_train)
```

k. Model SVM

```python
import numpy as np
from collections import Counter

class PegasosKernelSVM:
    def __init__(self, kernel='rbf', lambda_=0.01, gamma=0.1, n_iters=100):
        self.kernel = kernel
        self.lambda_ = lambda_
        self.gamma = gamma
        self.n_iters = n_iters
        self.alpha = None
        self.X = None
        self.y = None
        self.b = 0

    def _kernel_function(self, X1, X2):
        if self.kernel == 'linear':
            return np.dot(X1, X2.T)
        elif self.kernel == 'rbf':
            X1_sq = np.sum(X1 ** 2, axis=1).reshape(-1, 1)
            X2_sq = np.sum(X2 ** 2, axis=1).reshape(1, -1)
            dists = X1_sq - 2 * np.dot(X1, X2.T) + X2_sq
            return np.exp(-self.gamma * dists)
        else:
            raise ValueError("Unsupported kernel")

    def fit(self, X, y):
        n_samples = X.shape[0]
        y = np.where(y == 0, -1, 1)

        self.X = X
        self.y = y
        self.alpha = np.zeros(n_samples)
        self.b = 0

        for t in range(1, self.n_iters + 1):
            i = np.random.randint(0, n_samples)
            x_i = X[i].reshape(1, -1)
            y_i = y[i]
            K_i = self._kernel_function(self.X, x_i).flatten()
            margin = y_i * (np.sum(self.alpha * self.y * K_i) + self.b)

            eta = 1 / (self.lambda_ * t)

            if margin < 1:
                self.alpha[i] += eta
                self.b += eta * y_i

    def project(self, X):
        K = self._kernel_function(X, self.X)
        return np.dot(K, self.alpha * self.y) + self.b

    def predict(self, X):
        return np.sign(self.project(X))


class OneVsOneSVM:
    def __init__(self, kernel='rbf', lambda_=0.01, gamma=0.1, n_iters=100, max_samples_per_class=1000):
        self.kernel = kernel
        self.lambda_ = lambda_
        self.gamma = gamma
        self.n_iters = n_iters
        self.max_samples_per_class = max_samples_per_class
        self.models = {}

    def fit(self, X, y):
        self.models = {}
        self.classes_ = np.unique(y)

        for i in range(len(self.classes_)):
            for j in range(i + 1, len(self.classes_)):
                class_i = self.classes_[i]
                class_j = self.classes_[j]

                idx_i = np.where(y == class_i)[0][:self.max_samples_per_class]
                idx_j = np.where(y == class_j)[0][:self.max_samples_per_class]

                idx = np.concatenate([idx_i, idx_j])
                X_pair = X[idx]
                y_pair = y[idx]
                y_pair = np.where(y_pair == class_i, 0, 1)

                model = PegasosKernelSVM(kernel=self.kernel, lambda_=self.lambda_, gamma=self.gamma,
    n_iters=self.n_iters)
                model.fit(X_pair, y_pair)
                self.models[(class_i, class_j)] = model

    def predict(self, X):
        X = np.array(X)
        predictions = []

        for x in X:
            x = np.array(x).astype(float)
            votes = []
            for (class_i, class_j), model in self.models.items():
                pred = model.predict(x.reshape(1, -1))[0]
                winner = class_i if pred == -1 else class_j
                votes.append(winner)
            final_vote = Counter(votes).most_common(1)[0][0]
            predictions.append(final_vote)

        return np.array(predictions)
```

l. Implementasi model SVM

```
1 model = OneVsOneSVM(kernel='rbf', lambda_=0.01, gamma=0.01, n_iters=100)
2 model.fit(X_train_under, y_train_under.values)
3
4 y_pred = model.predict(X_test_scaled)
5 accuracy1 = np.mean(y_pred == y_test.values)
6 print(f"Akurasi Model 1 (Normal vs Bermasalah): {accuracy1:.4f}")
```

m. Classification report

```
1 # Daftar label kelas
2 unique_labels = np.unique(np.concatenate((y_test, y_pred)))
3 label_names = [str(label) for label in unique_labels]
4 label_indices = {label: i for i, label in enumerate(unique_labels)}
5
6 # Konversi label ke indeks
7 y_true_idx = np.array([label_indices[label] for label in y_test])
8 y_pred_idx = np.array([label_indices[label] for label in y_pred])
9
10 # Inisialisasi confusion matrix
11 n_classes = len(unique_labels)
12 conf_matrix = np.zeros((n_classes, n_classes), dtype=int)
13
14 for true, pred in zip(y_true_idx, y_pred_idx):
15     conf_matrix[true][pred] += 1
16
17 # Print classification report
18 print("\nManual Classification Report:")
19 print(f"{'Class':<10} {'Precision':>10} {'Recall':>10} {'F1-Score':>10} {'Support':>10}")
20
21 for i in range(n_classes):
22     TP = conf_matrix[i, i]
23     FP = conf_matrix[:, i].sum() - TP
24     FN = conf_matrix[i, :].sum() - TP
25     support = conf_matrix[i, :].sum()
26
27     precision = TP / (TP + FP) if (TP + FP) > 0 else 0.0
28     recall = TP / (TP + FN) if (TP + FN) > 0 else 0.0
29     f1 = 2 * precision * recall / (precision + recall) if (precision + recall) > 0 else 0.0
30
31     print(f"{label_names[i]:<10} {precision:10.2f} {recall:10.2f} {f1:10.2f} {support:10}")
32
33 # Optional: Print overall accuracy
34 accuracy = np.mean(y_pred == y_test)
35 print(f"\nOverall Accuracy: {accuracy:.4f}")
```

n. Confusion matriks

```
Confussion Matriks

 1 import numpy as np
 2 import seaborn as sns
 3 import matplotlib.pyplot as plt
 4
 5 unique_labels = np.unique(np.concatenate((y_test, y_pred)))
 6 label_names = [str(label) for label in unique_labels]
 7 label_indices = {label: i for i, label in enumerate(unique_labels)}
 8
 9 y_true_idx = np.array([label_indices[label] for label in y_test])
10 y_pred_idx = np.array([label_indices[label] for label in y_pred])
11
12 n_classes = len(unique_labels)
13 conf_matrix = np.zeros((n_classes, n_classes), dtype=int)
14
15 for true, pred in zip(y_true_idx, y_pred_idx):
16     conf_matrix[true][pred] += 1
17
18 plt.figure(figsize=(8, 6))
19 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
20             xticklabels=label_names, yticklabels=label_names)
21 plt.xlabel('Predicted Label')
22 plt.ylabel('True Label')
23 plt.title('Confusion Matrix')
24 plt.tight_layout()
25 plt.show()
```

2.  File app.py – Deploy via Streamlit

```python
1  import streamlit as st
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import joblib
6  from modelSV import *
7
8  def load_model_and_scaler(model_path, scalerFeature_path, transformFeature_path, transformTarget_path):
9      model = joblib.load(model_path)
10     scaler = joblib.load(scalerFeature_path)
11     le_F = joblib.load(transformFeature_path)
12     le_T = joblib.load(transformTarget_path)
13     return model, scaler, le_F, le_T
14
15 st.title("Prediksi Jenis Anomali Dalam Lalu Lintas Jaringan")
16 st.write("Upload File Test (Format TXT atau CSV)")
17
18 uploaded_file = st.file_uploader("Pilih file CSV atau TXT", type=["csv", "txt"])
19
20 columns = [
21     "duration", "protocol_type", "service", "flag", "src_bytes", "dst_bytes", "land",
22     "wrong_fragment", "urgent", "hot", "num_failed_logins", "logged_in", "num_compromised",
23     "root_shell", "su_attempted", "num_root", "num_file_creations", "num_shells",
24     "num_access_files", "num_outbound_cmds", "is_host_login", "is_guest_login", "count",
25     "srv_count", "serror_rate", "srv_serror_rate", "rerror_rate", "srv_rerror_rate",
26     "same_srv_rate", "diff_srv_rate", "srv_diff_host_rate", "dst_host_count",
27     "dst_host_srv_count", "dst_host_same_srv_rate", "dst_host_diff_srv_rate",
28     "dst_host_same_src_port_rate", "dst_host_srv_diff_host_rate", "dst_host_serror_rate",
29     "dst_host_srv_serror_rate", "dst_host_rerror_rate", "dst_host_srv_rerror_rate",
30     "label", "difficulty"
31 ]
32
33
34 if uploaded_file is not None:
35     if uploaded_file.name.endswith('.txt'):
36         delimiter = ","
37     elif uploaded_file.name.endswith('.csv'):
38         delimiter = ","
39
40
41     df = pd.read_csv(uploaded_file, delimiter=delimiter,names=columns)
42
43     st.write("Berikut adalah beberapa data dari file yang Anda unggah:")
44     st.write(df.head())
45
46     model, scaler, le_F, le_T = load_model_and_scaler('one_vs_one_svm_model.joblib', 'scaling.joblib',
   'feature_encoders.joblib', 'target_encoder.joblib')
47
48     target_column = 'attack_category'
49
50     categorical_columns = df.select_dtypes(include=['object']).columns.tolist()
51
52     feature_encoders = {}
53
54     for col in categorical_columns:
55         df[col] = df[col].apply(lambda x: le_F[col].transform([x])[0] if x in le_F[col].classes_ else -1)  #
   -1 atau label default
56
57     X_pred = df.drop(["label", "num_outbound_cmds"], axis=1)
58
59     X_pred_scaled = scaler.transform(X_pred)
60     y_pred = model.predict(X_pred_scaled)
61
62     y_pred_original = le_T.inverse_transform(y_pred)
63
64     y_pred_df = pd.DataFrame(y_pred_original, columns=['predicted_attack_category'])
65
66     st.write('Hasil Prediksi')
67     st.write(y_pred_df)
68
69
70     if len(df) > 10:
71         st.write('Distribusi Hasil Prediksi')
72         y_pred_counts = y_pred_df.value_counts()
73         plt.figure(figsize=(8, 6))
74         y_pred_counts.plot(kind='bar')
75         plt.title('Distribusi Hasil Prediksi')
76         plt.xlabel('Kelas')
77         plt.ylabel('Frekuensi')
78         plt.xticks(rotation=0)
79         st.pyplot(plt)
80     else:
81         pass
```