



Who Is The Killer?

Piraeus Vice Pattern Recognition Project

Student 1: Βασίλαινα Μαρία (ID: ΠΙ23015)

Student 2: Γρηγοράσκου Τερέζα (ID: ΠΙ22037)

Student 3: Λιβέρης Φώτης (ID: ΠΙ23104)

Department of Informatics
University of Piraeus

February 20, 2026

Abstract

In this project, we addressed the killer identification problem as a multiclass classification task. We first estimated class-specific means and covariances using Maximum Likelihood Estimation (MLE), and used them to build a Gaussian Bayes classifier. We then implemented a linear classifier, a non-linear SVM (RBF kernel), a two-hidden-layer MLP neural network, and an unsupervised PCA + k-means approach with majority-vote label mapping. All preprocessing (standardization and one-hot encoding) was fitted on the TRAIN split and applied to VAL and TEST to avoid data leakage. PCA was used both for visualization and as a preprocessing step for k-means. On the validation set, non-linear models (SVM and MLP) achieved the best performance, outperforming the Gaussian Bayes and linear classifiers. The PCA + k-means approach achieved lower accuracy, indicating that natural clusters in feature space do not perfectly align with killer identities. Overall, the features that proved to be of the most importance were those related with the victim and the crime scene. Last but not least, modeling nonlinear decision boundaries significantly improved classification performance.

Collectively, in order to classify data and reach our final prediction we followed the following steps.

- MLE (Maximum Likelihood Estimation)
- Gaussian Bayes classifier
- Linear classifier
- SVM (Support Vector Machine)
- MLP (Multi-Layer Perceptron)
- PCA (Principal Component Analysis)
- k-means

Contents

1	Introduction	4
2	Data description	6
2.1	Feature overview	6
3	Results	7
Q1:	Exploratory analysis	7
	Explanation	9
Q2:	Gaussian MLE per killer	10
	Explanation	14
Q3:	Multiclass Gaussian Bayes classifier	15
	Explanation	16
Q4:	Linear classifier	17
	Explanation	18
Q5:	Support Vector Machine	19
	Explanation	21
Q6:	Multi-Layer Perceptron	22
	Explanation	23
Q7:	Principal Component Analysis	24
	Explanation	25
Q8:	k -means in PCA space	26
	Explanation	27
4	Discussion and conclusions	29
5	LLM prompts and responses	32

1 Introduction

In the current project we were asked to identify the most likely killer for each crime incident in the “Piraeus Vice” dataset. This is formulated as a multiclass classification problem, where each of the incidents must be related to one of the 8 possible killers based on the available features. So, the task is to build models that can accurately predict the killer ID.

Each row in the dataset corresponds to a single crime incident and contains the following features:

- The hour that the crime was committed
- The latitude and the longitude, together forming the coordinates of the crime scene
- The victim’s age
- The prevailing temperature at the time the crime was committed
- The prevailing humidity at the time the crime took place
- The distance between the crime scene and the nearest police precinct
- The population density of the area
- The weapon type (knife, handgun 9mm, revolver .38, shotgun, blunt object, unknown)
- The scene type (street, residence, business, other)
- The prevailing weather conditions (clear, rain, snow, fog, unknown)
- The victim gender (male, female)
- The incident ID
- The partition type (TRAIN, VAL, TEST)
- The killer ID (1 ... 8)

Comment: Further analysis will follow in section 2.

The main steps of analysis we are instructed to use are:

- Q1: Exploratory Data Analysis
- Q2: Parameter Estimation (MLE)
- Q3: Gaussian Bayes Classifier
- Q4: Linear Classifier
- Q5: Support Vector Machine (SVM)
- Q6: Multilayer Perceptron (MLP)
- Q7: Principal Component Analysis (PCA)
- Q8: PCA + k-means

The following picture shows a high-level overview of the killer identification process.

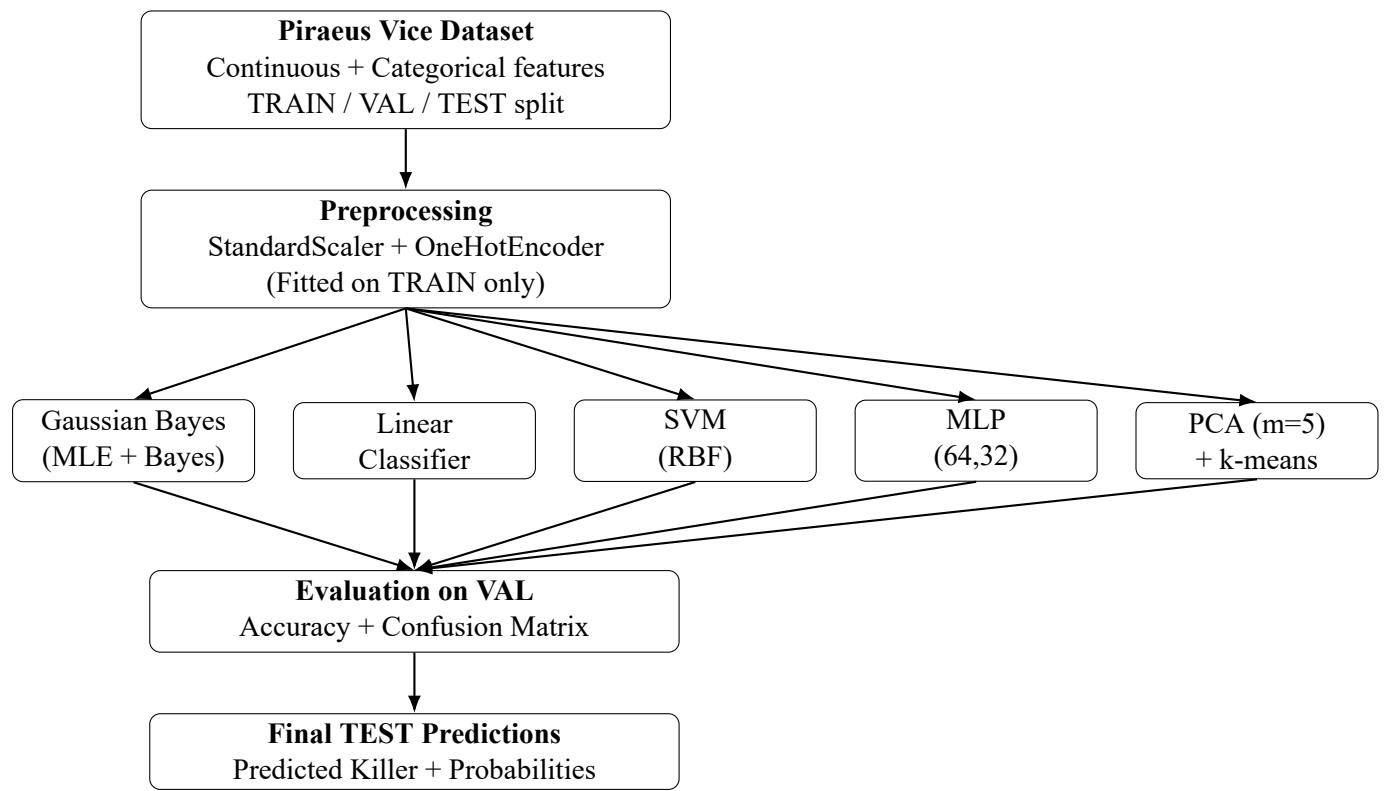


Figure 1: High-level overview of the killer identification pipeline.

2 Data description

The csv file we were given consists of a total of 4.800 rows, each one representing one crime incident. As we have already mentioned, it consists of 15 columns, one for each feature. An overview of these features and their corresponding types follows in section 2.1. The dataset contains a predefined split column that is either TRAIN, VAL or TEST. This means that the data is already partitioned for modeling and evaluation. As far as we can tell the dataset is not standardised or normalised.

2.1 Feature overview

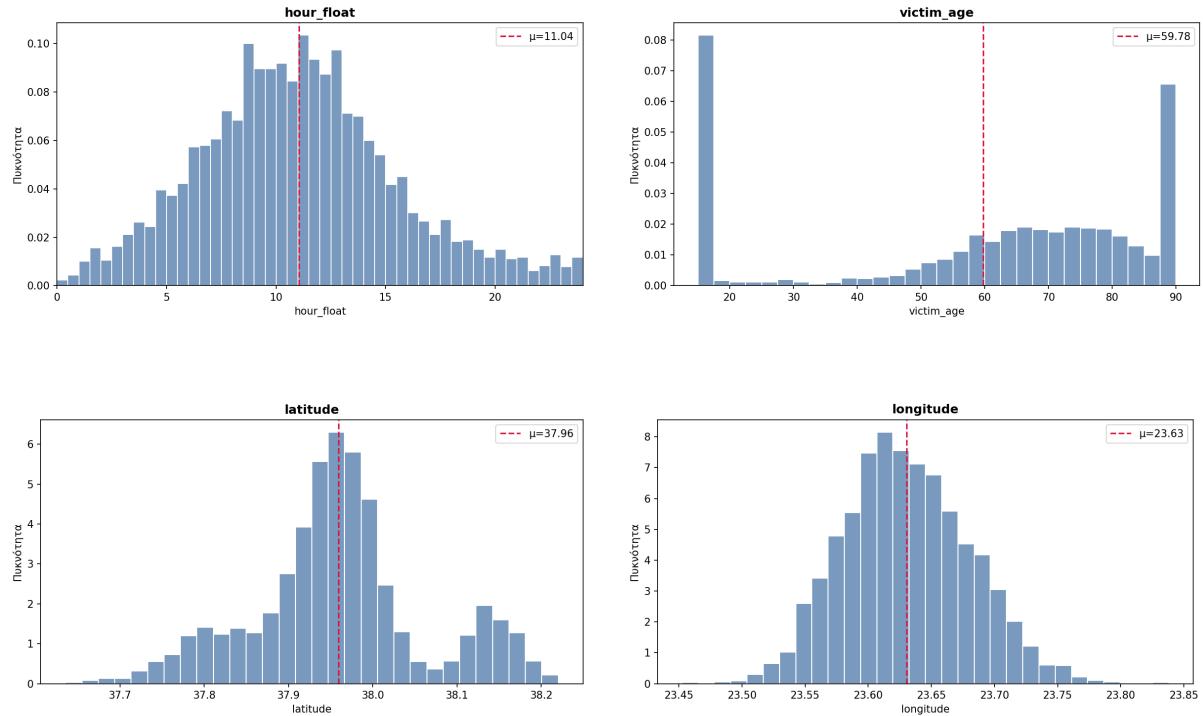
Table 1: Feature summary.

Feature name	Type
incident_id	identifier
hour_float	continuous
latitude	continuous
longitude	continuous
victim_age	continuous
temp_c	continuous
humidity	continuous
dist_precinct_km	continuous
pop_density	continuous
weapon_code	categorical
scene_type	categorical
weather	categorical
vic_gender	categorical (binary)
split	categorical
killer_id	identifier

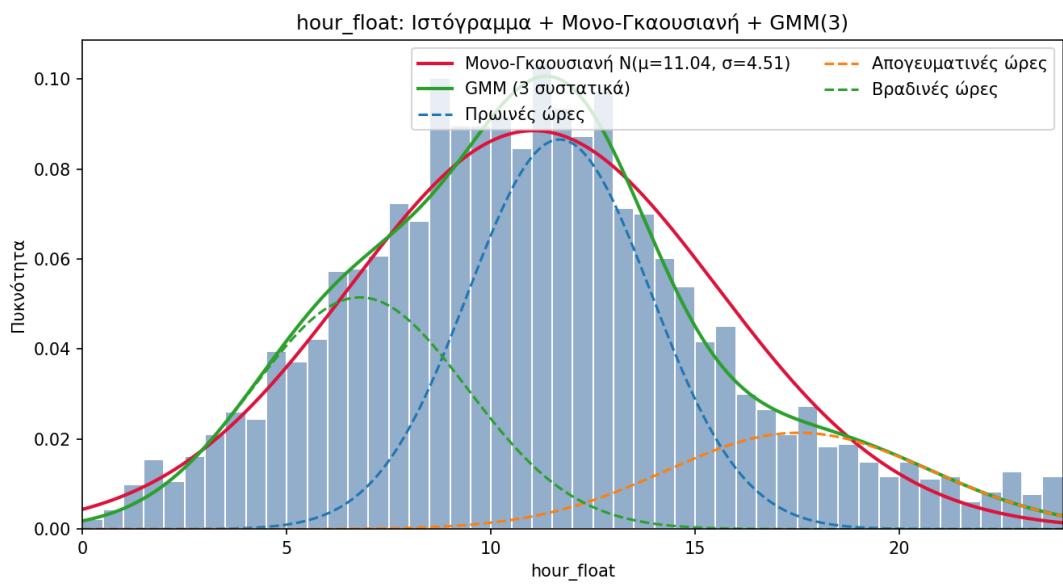
3 Results

Q1: Exploratory analysis

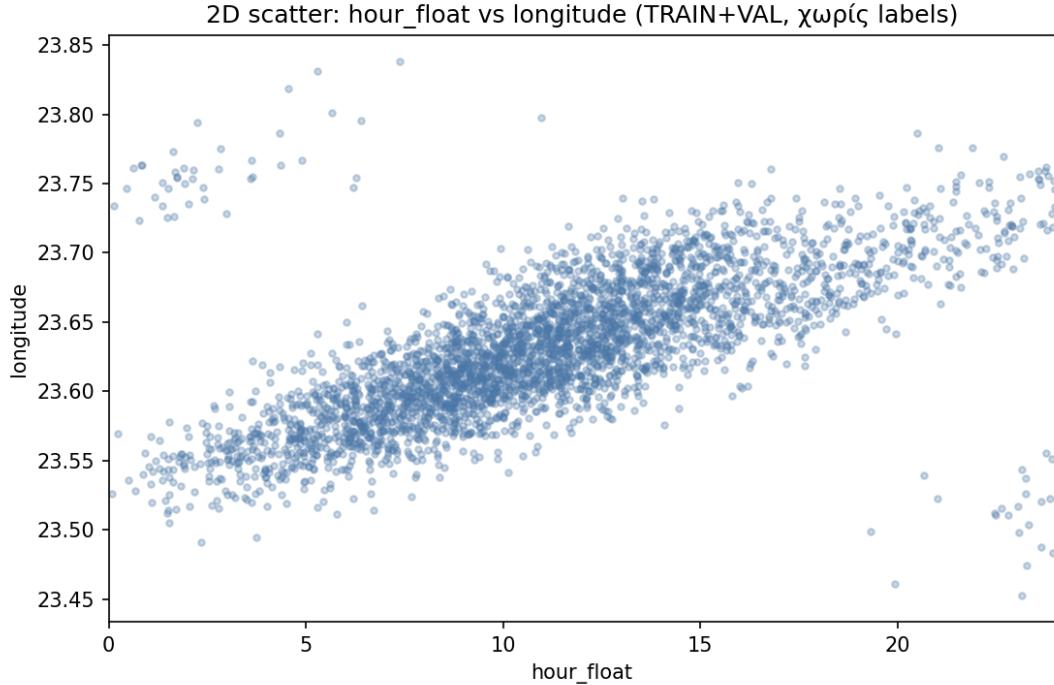
The images that follow are the histograms for the key continuous variables (hour, age, latitude, longitude) that we produced.



Below is the comparison of the single Gaussian fit and a mixture of Gaussians for `hour_float`.



Also, we provide you with the additional two-dimensional plot involving `hour_float` and `longitude`. using TRAIN+VAL and without using any labels, as requested.



Our observations around the 4 key continuous variables are the following:

- Firstly, about the variable `hour_float` the incidents are at their peak during daytime, which is somewhat unexpected.
- As far as the variable `victim_age` is concerned it can easily be observed that the selected victims were either of a very young age or of a very old. With that said, we can safely come to the conclusion that they were targeted because they belonged to the most vulnerable age groups.
- The histogram of the variable `latitude` presents some fluctuations that are at their core closer to its average value.
- Lastly, concerning the variable `longitude`, by observing both the corresponding histogram as well as the 2D scatter, it is clear that the incident density is also thicker closer to its average value.

Regarding the single Gaussian, we believe that it is inadequate. We can easily see that it is not symmetric. But it is not solely that. It also ignores some of the peaks. As a result it fails to represent correctly the given data. By splitting it to the 3-component Gaussian Mixture we achieve better interpretation of the times of the day as well as better representation of the data as it focuses on distinct times of the day.

Explanation

We have decided to use the library GaussianMixture that proved to be very useful. Also, we created the following function in order to be able to create the histograms. As we can see it requires 5 parameters, 3 of which have default values in case they don't receive any. The parameters are:

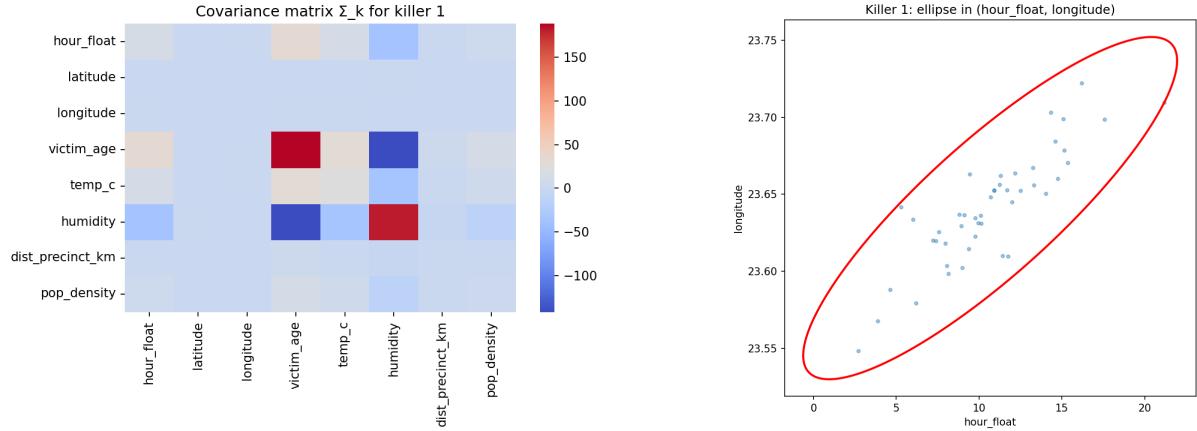
- series: a pandas series containing numerical data,
- title: it describes the column we are focusing on,
- bins: the number or range of lines that the histogram will consist of,
- kde: it controls whether the result will be a histogram or a smooth curve that represents the distribution of the data,
- xlim: it limits the range of the x-axis.

```
1 def plot_hist_with_stats(series, title, bins=30, kde=False, xlim=None):
2     s = series.dropna().values
3     fig, ax = plt.subplots(figsize=(8, 4.5))
4     sns.histplot(s, bins=bins, stat="density", kde=kde, edgecolor="white",
5                   color="#4C78A8", ax=ax)
6     if xlim:
7         ax.set_xlim(xlim)
8         ax.set_title(title, fontweight="bold")
9         ax.set_xlabel(title)
10        ax.set_ylabel("Πυκνότητα")
11        mu, sigma = np.mean(s), np.std(s, ddof=1)
12        ax.axvline(mu, color="crimson", linestyle="--", label=f" $\mu={mu:.2f}$ ")
13        ax.legend()
14    plt.tight_layout()
15    return fig, ax, mu, sigma
```

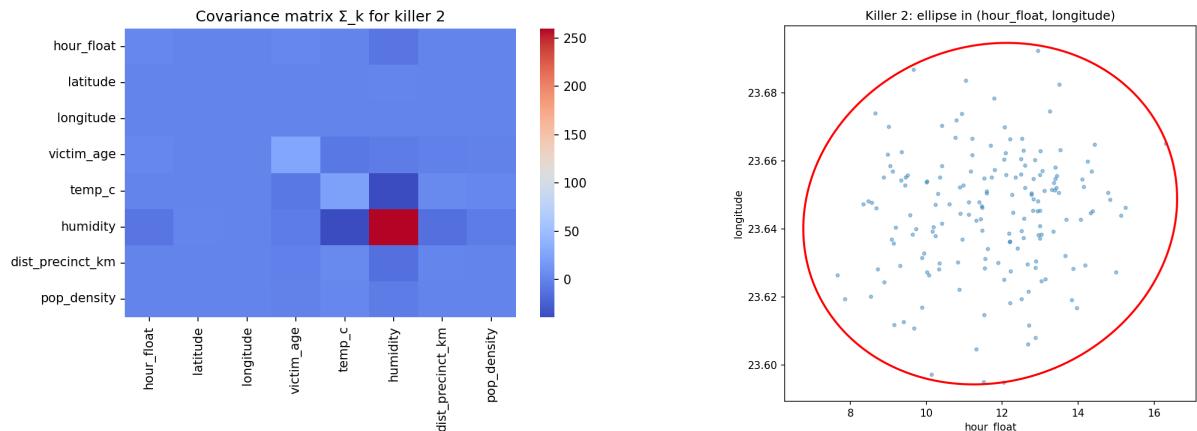
Q2: Gaussian MLE per killer

Below we present you one covariance heatmap and one 2D projection with ellipses per killer, as requested.

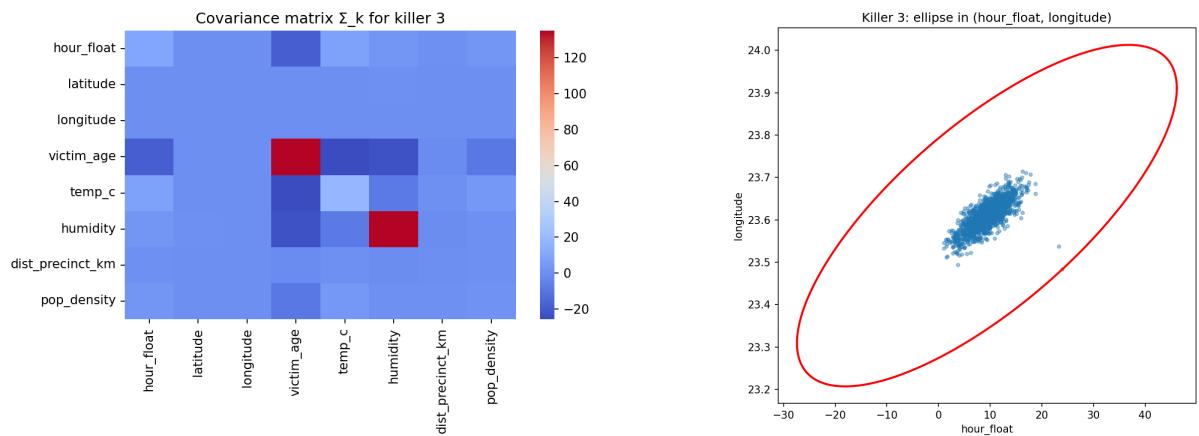
- Killer 1:



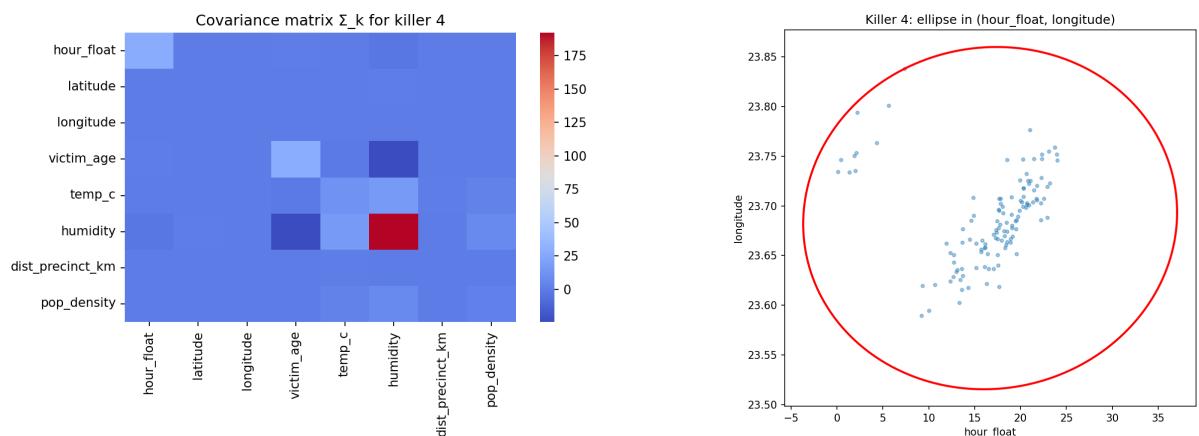
- Killer 2:



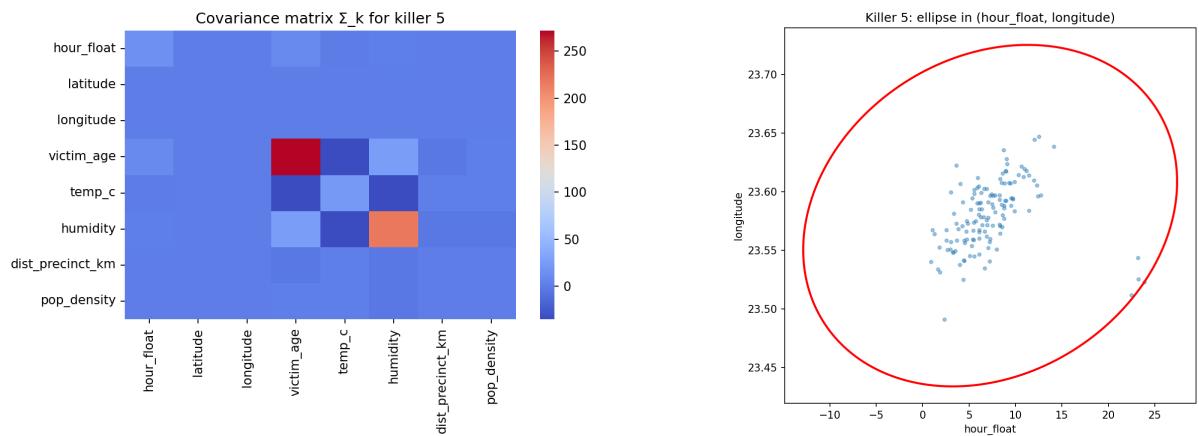
- Killer 3:



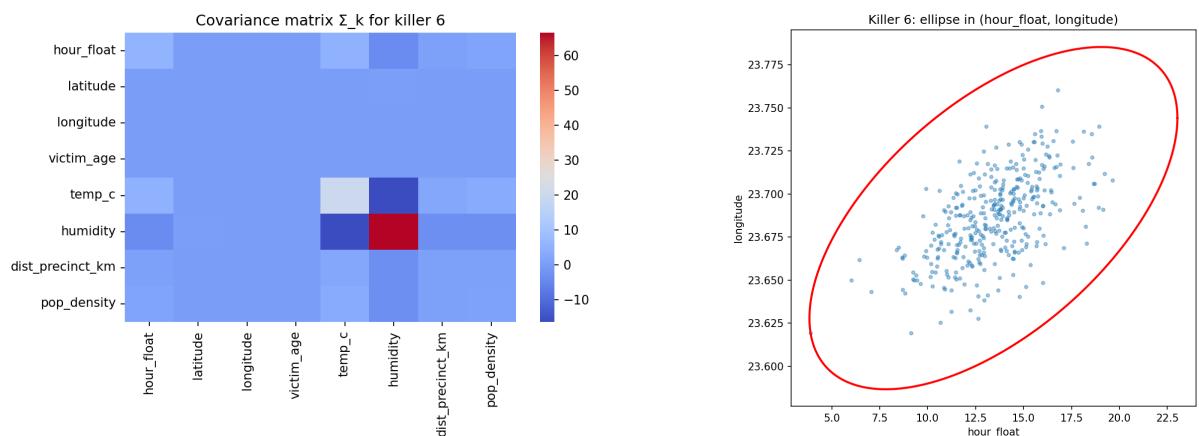
- Killer 4:



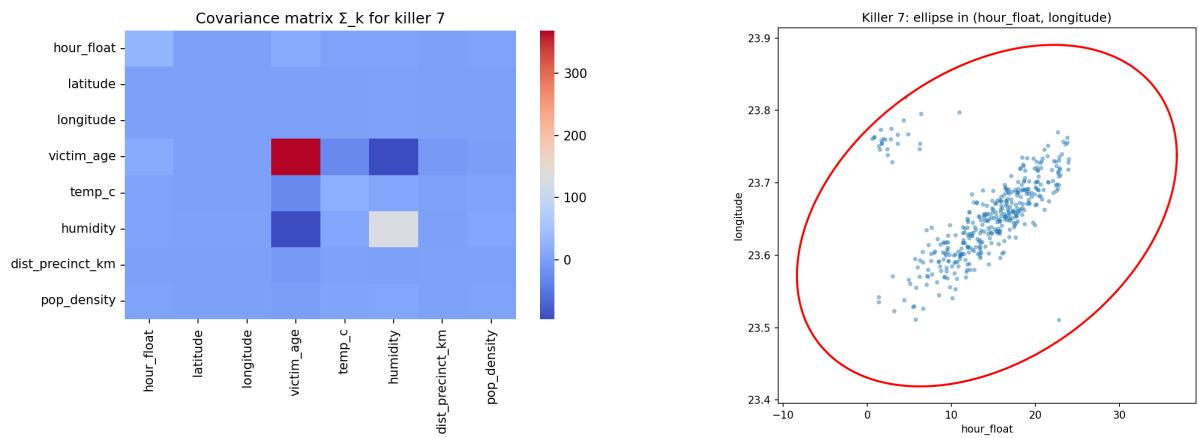
- Killer 5:



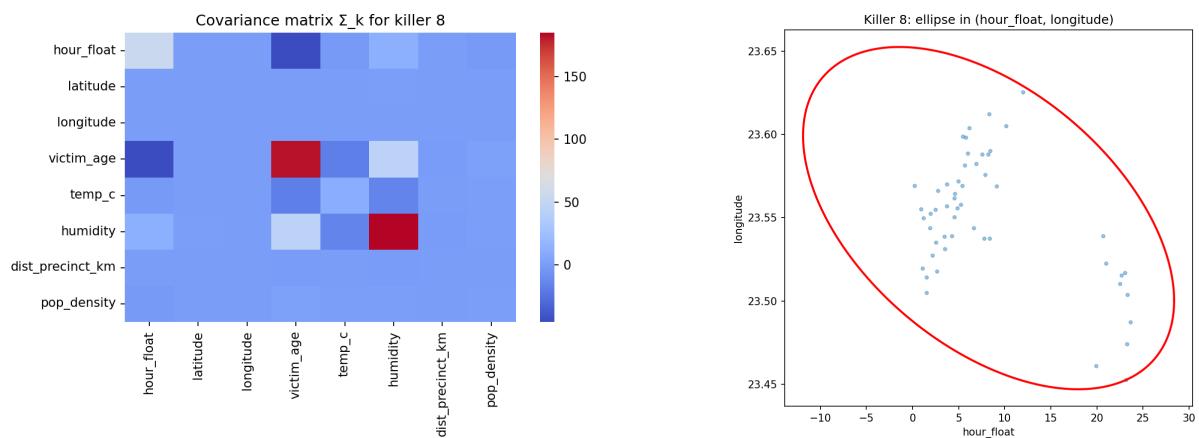
- Killer 6:



- Killer 7:



- Killer 8:



Explanation

To estimate the mean and covariance of the continuous features for each killer on TRAIN we assumed that, for each killer k , the continuous feature vector follows a multivariate Gaussian distribution. To estimate the parameters we used Maximum Likelihood Estimation (MLE), restricting ourselves only to the TRAIN split, where the killer labels are known. For each killer k , we collected all TRAIN incidents committed by that killer. This way, each killer's profile starts to assemble. Lastly, we visualise our results for better interpretation.

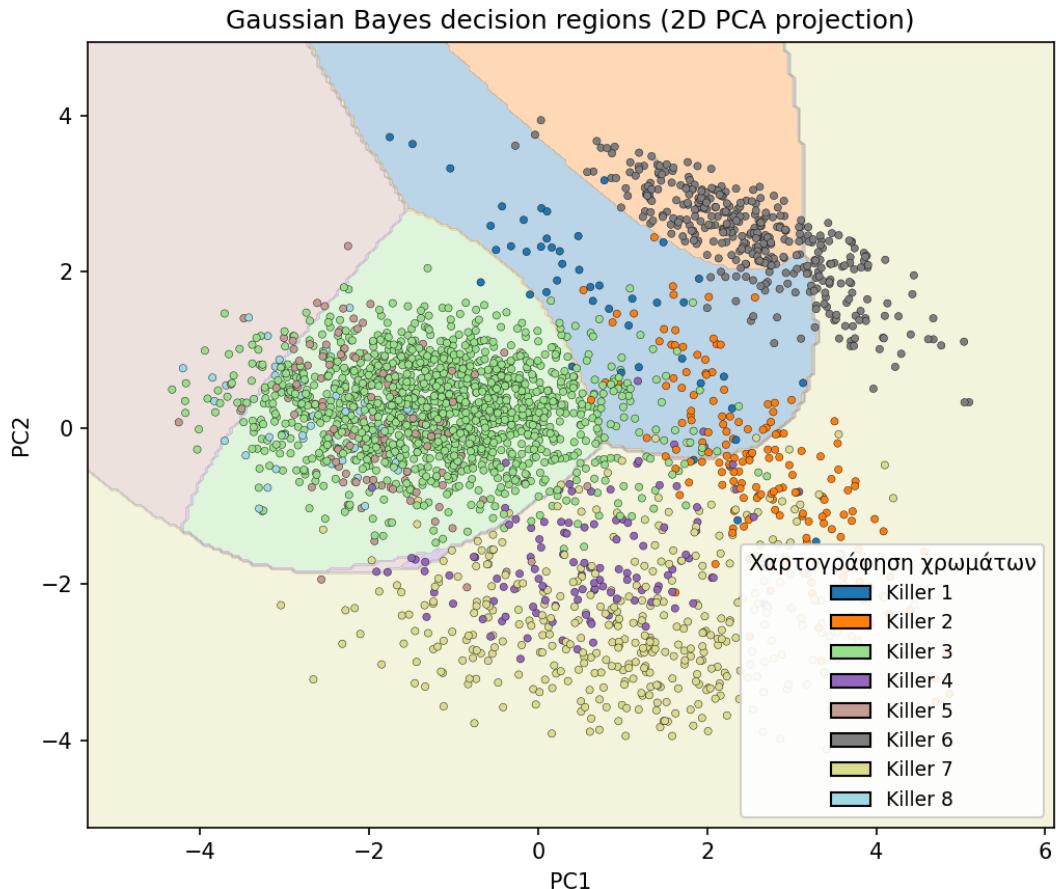
Q3: Multiclass Gaussian Bayes classifier

We evaluated the TRAIN accuracy and the VAL accuracy and concluded that they are equal to 89% and 90% respectively. Here is the requested confusion matrix on VAL.

Table 2: Confusion matrix for the Gaussian Bayes classifier on VAL.

True / Pred	1	2	3	4	5	6	7	8
1	13	1	2	0	0	1	0	0
2	1	59	1	0	0	1	0	0
3	3	0	474	0	5	0	7	2
4	0	0	2	38	1	0	4	0
5	0	0	32	1	11	0	1	3
6	0	0	0	0	0	133	0	0
7	1	2	0	3	0	0	137	0
8	0	0	15	0	2	0	0	2

We can see that the matrix produced is mostly diagonal and the killers with the highest presence are the killers 3, 6 and 7, which will be explained due to the accuracy of Q8 in the end of the assignment. Below follows the two-dimensional projection, visualising the decision regions induced by the Bayes classifier, with TRAIN points coloured by true killer. Each curved boundary separates the feature regions of each killer.



Explanation

We implemented a multiclass Gaussian Bayes classifier using the parameters estimated in Q2. For each killer k, we modelled the continuous feature vector as a multivariate Gaussian distribution. For each incident x we calculated the long-posterior probability. Lastly, we evaluated the classifier on the VAL split by calculating the overall accuracy and of course prosuced the corresponding confusion matrix, as shown in thee code below.

```
1 def predict_proba(x):
2     # returns vector of posterior probabilities for one x
3     log_scores = []
4     for k in killers:
5         lg = np.log(pi[k]) + log_gaussian(
6             x,
7             mu_dict[k],
8             invcov[k],
9             logdet[k]
10        )
11        log_scores.append(lg)
12    log_scores = np.array(log_scores)
13    # softmax
14    exp_scores = np.exp(log_scores - np.max(log_scores))
15    return exp_scores / exp_scores.sum()
16
17 def predict_class(x):
18     pp = predict_proba(x)
19     return killers[np.argmax(pp)]
20
21 # Predict on TRAIN
22 ytrain_pred = np.array([predict_class(x) for x in Xtrain])
23 train_acc = accuracy_score(ytrain, ytrain_pred)
24
25 # Predict on VAL
26 yval_pred = np.array([predict_class(x) for x in Xval])
27 val_acc = accuracy_score(yval, yval_pred)
28
29 print("TRAIN accuracy:", train_acc)
30 print("VAL accuracy:", val_acc)
31 print("VAL confusion matrix:")
32 print(confusion_matrix(yval, yval_pred))
```

Q4: Linear classifier

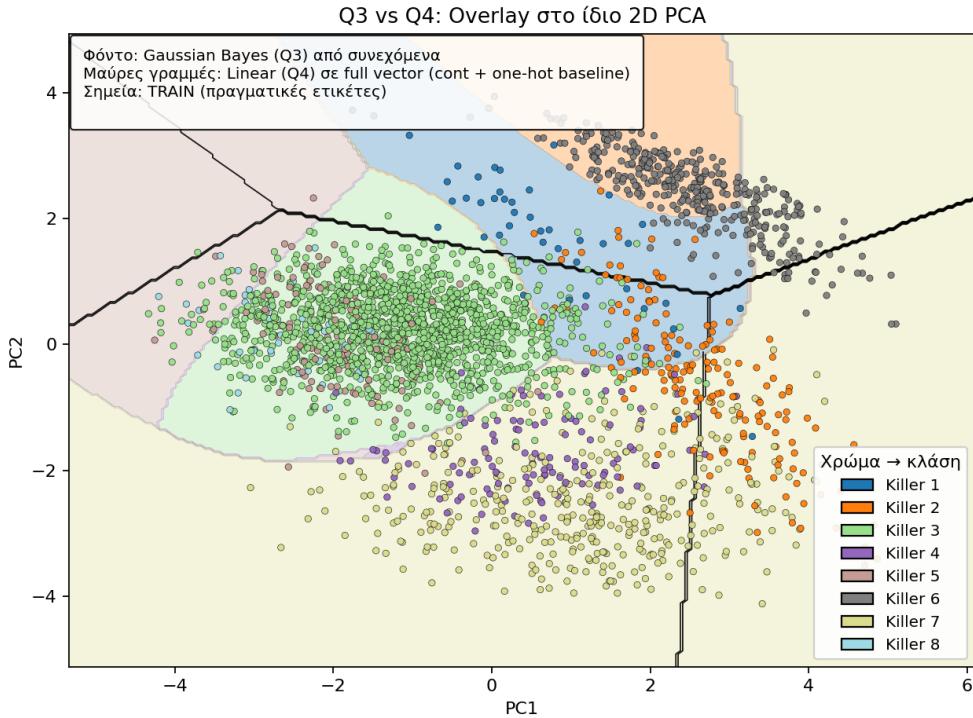
This time, the TRAIN accuracy and the VAL accuracy are equal to 77% and 78% respectively. We need to highlight that these numbers are lower than Q3's. That could mean that the Gaussian assumption is more accurate. Here is the confusion matrix on VAL.

Table 3: Confusion matrix for the Gaussian Bayes classifier on VAL.

True / Pred	1	2	3	4	5	6	7	8
1	0	1	9	0	2	4	1	0
2	1	38	17	0	0	3	3	0
3	0	4	473	0	4	1	9	0
4	1	0	4	0	0	0	40	0
5	0	1	28	1	14	0	4	0
6	0	3	12	0	0	117	1	0
7	3	4	26	0	1	1	108	0
8	0	0	12	0	4	0	3	0

Comparing it to Q3's matrix it obviously does not resemble a diagonal matrix which may mean that we have strayed from our end goal. While some numbers have slightly changed(e.g. 474 -> 473), most of them have completely changed. Also, the last column consists entirely of zeros. However, the killers with the highest presence remain the same.

Also, in the 2D PCA projection used in Q3, we overlayed the approximate linear decision boundaries and the result is the following.



When compared to Q3's PCA projection one can easily see that the regions are separated by lines and clearly determine which one corresponds to each killer. Some may correspond to more than one killers. This happens because the linear model is not able to fully adapt to the shape of the data the way that the Gaussian Bayes does.

Explanation

We treated the killer identification as a discriminative multiclass classification problem using the full vector consisting of the 8 continuous features (hour_float, latitude, longitude, victim_age, temp_c, humidity, dist_precinct_km and pop_density) and 4 more categorical features (weapon_code, scene_type, weather, vic_gender). We implemented the linear classifier given. We trained the model on the TRAIN split using MSE (Mean Squared Error loss). Lastly, we evaluate the accuracy and the confusion matrix and produce the 2D PCA overlaying the one from Q3.

```
1 Xtr_t = torch.tensor(Xfull_train, dtype=torch.float32).to(device) # (N,
2   ↵ d_full)
3 Xva_t = torch.tensor(Xfull_val,      dtype=torch.float32).to(device)
4 ytr_1h = torch.eye(S,
5   ↵ dtype=torch.float32)[torch.tensor(ytrain_idx)].to(device)
6
7 d_full = Xfull_train.shape[1]
8 linear = nn.Linear(d_full, S).to(device)
9 criterion = nn.MSELoss()
10 opt = optim.Adam(linear.parameters(), lr=1e-3)
11
12 linear.train()
13 for epoch in range(1500):
14     opt.zero_grad()
15     out = linear(Xtr_t)           # logits (S)
16     loss = criterion(out, ytr_1h) # SSE με one-hot
17     loss.backward()
18     opt.step()
19     if epoch % 300 == 0:
20         print(f"[Q4] epoch {epoch}, loss={loss.item():.4f}")
21
22 linear.eval()
23 with torch.no_grad():
24     tr_logits = linear(Xtr_t).cpu().numpy()
25     va_logits = linear(Xva_t).cpu().numpy()
26     tr_pred_idx = tr_logits.argmax(1)
27     va_pred_idx = va_logits.argmax(1)
28     print("Q4 TRAIN acc:", accuracy_score(ytrain_idx, tr_pred_idx))
29     print("Q4 VAL acc:",   accuracy_score(yval_idx,   va_pred_idx))
30     print("Q4 VAL confusion:\n", confusion_matrix(yval_idx, va_pred_idx))
```

Note to teacher: As you may have noticed if you have run our code, on the terminal also appear the Q3's results. Everything on terminal is exclusively for our convenience in order to make sure all things run smoothly and are executed correctly. It should not concern you as we have provided you with Q4's results in the previous section.

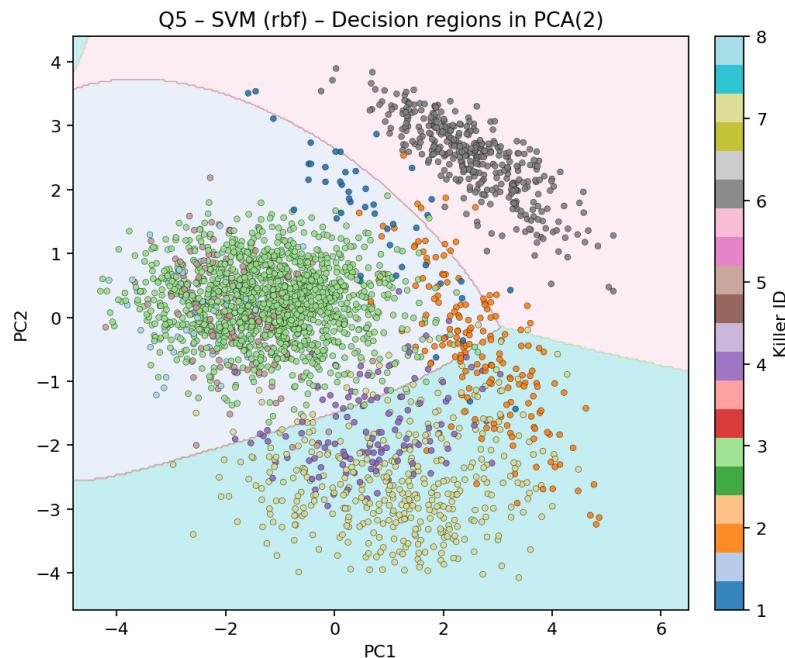
Q5: Support Vector Machine

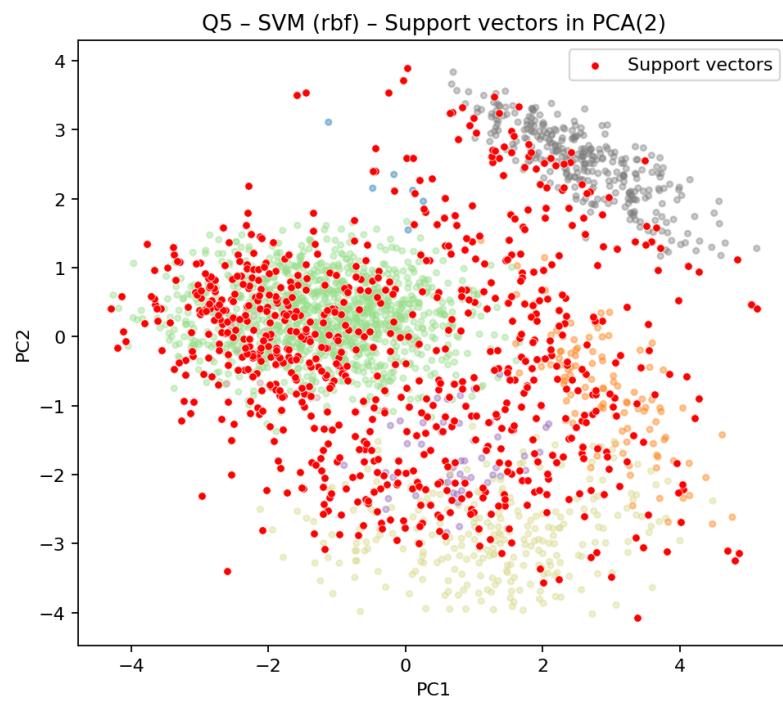
The VAL accuracy for the RBF method is 94%. The confusion matrix is the following. Once again it slightly resembles a diagonal matrix which means we are getting closer.

Table 4: Confusion matrix for the Gaussian Bayes classifier on VAL.

True / Pred	1	2	3	4	5	6	7	8
1	14	0	3	0	0	0	0	0
2	0	62	0	0	0	0	0	0
3	0	0	484	0	1	0	2	4
4	0	0	1	35	1	0	8	0
5	0	0	14	0	30	0	1	3
6	0	0	0	0	0	133	0	0
7	0	0	0	1	1	0	141	0
8	0	0	6	0	3	0	0	10

In the images that follow we can see the decision regions as well as the support vectors.





Highlighted in red are distributed the incidents that determine the characteristics of the killers' profiles.

Explanation

Our SVM setup is formed by choosing either RBF or Polynomial kernel. We prefered RBF as it allows non-linear decision boundaries. The code is modified for either one.

```
1 def small_search_space(kernel: str):
2     if kernel == "rbf":
3         Cs = [0.3, 1, 3]
4         gammas = ["scale", 0.1, 0.03]
5         return [(c, g, None, None) for c in Cs for g in gammas]
6     else:
7         Cs = [0.3, 1, 3]
8         gammas = ["scale", 0.1]
9         degrees = [2]
10        coef0s = [0, 1]
11        return [(c, g, d, c0) for c in Cs for g in gammas for d in degrees
12             for c0 in coef0s]
```

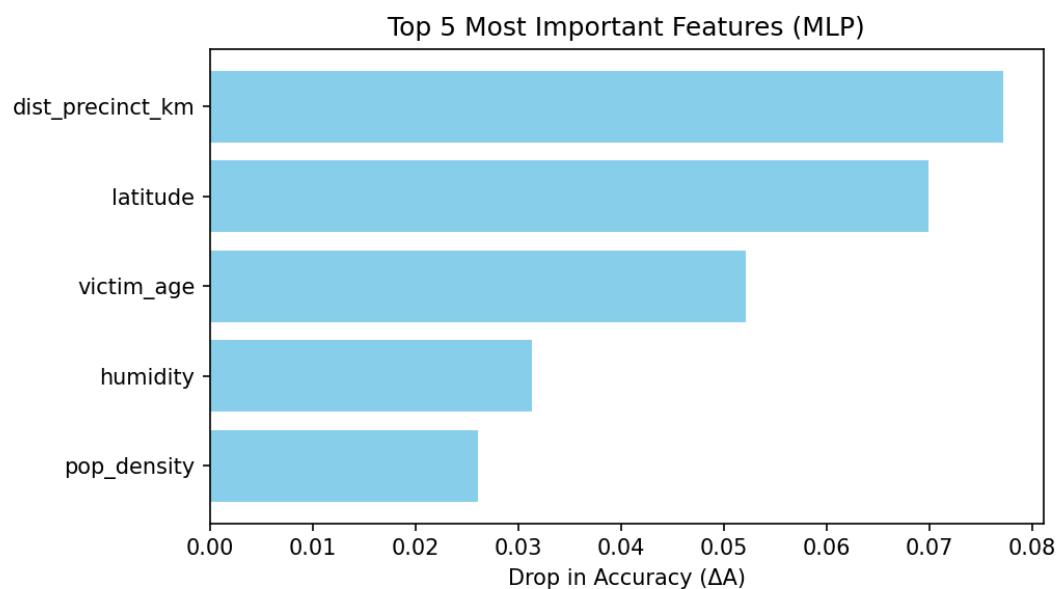
We continued with tuning the hyperparameters using a small manual grid search on VAL. Each candidate model was trained on TRAIN and evaluated on VAL. We chose the one that seemed more suitable to us based on the highest validation accuracy. Lastly, since the SVM is a binary classifier due to inheritance, we used One vs Rest strategy. This way, there is one binary SVN trained per killer, which means we have 8 classifiers, each one able to separate class k from the rest and as far as prediction is concerned, it is the class with the highest score that is selected. All preprocessing was applied inside a Pipeline thus ensuring there will be no data leakage.

Q6: Multi-Layer Perceptron

The VAL accuracy is equal to 94% (same as the SVM's).

The features of most importance seem to be (in ascending order):

1. dist_precinct_km $\Delta A = 0.07$
2. latitude $\Delta A = 0.06$
3. victim_age $\Delta A = 0.05$
4. humidity $\Delta A = 0.03$
5. pop_density $\Delta A = 0.02$



Explanation

We implemented a Multilayer Perceptron (MLP) as classifier using the following code:

```
1 mlp = MLPClassifier(
2     hidden_layer_sizes=(64, 32),      # 2 hidden layers
3     activation="relu",
4     solver="adam",
5     max_iter=200,
6     random_state=0,
7     early_stopping=True
8 )
```

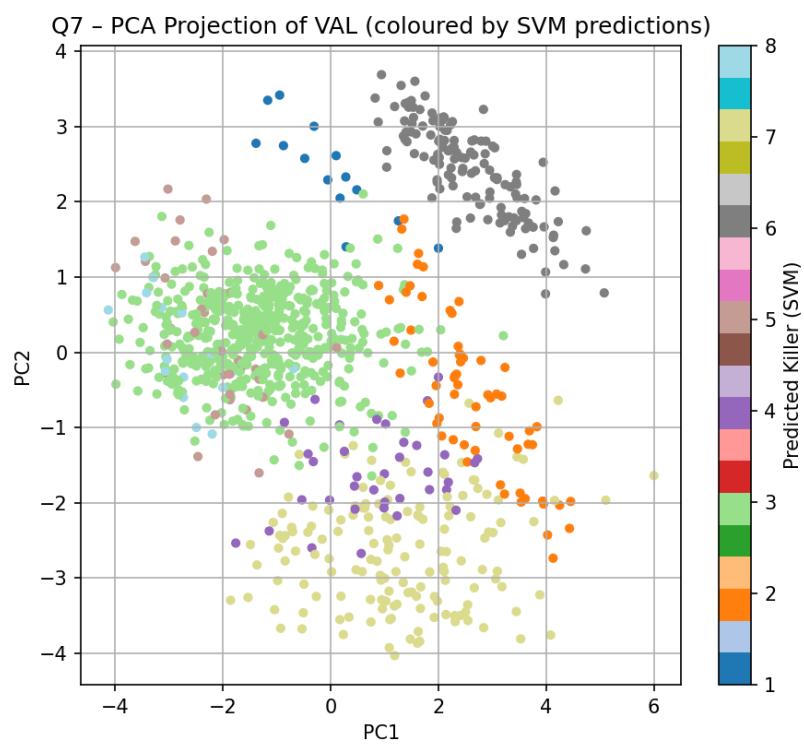
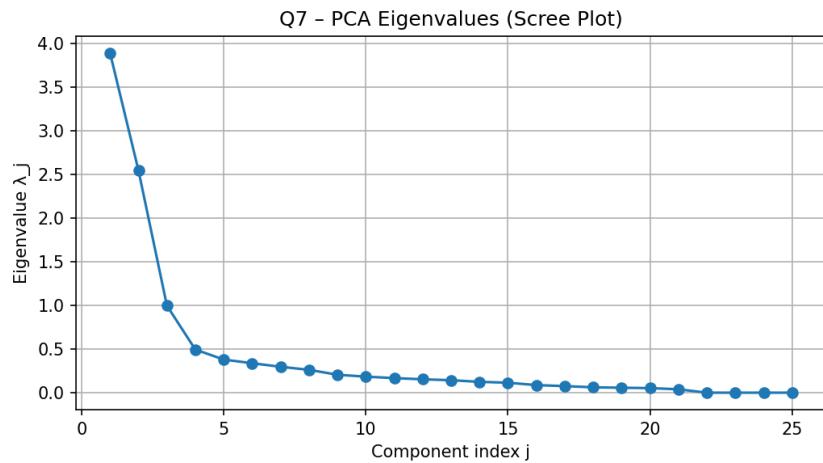
We continued by once again training the model on TRAIN, calculated VAL accuracy etc. Lastly, feature importance was estimated using permutation importance.

```
1 Z_val = model.named_steps["pre"].transform(X_val)
2
3 importances = []
4 rng = np.random.default_rng(0)
5
6 for j in range(Z_val.shape[1]):
7     Zp = Z_val.copy()
8     Zp[:, j] = rng.permutation(Zp[:, j])  # permute one feature
9     pred = model.named_steps["mlp"].predict(Zp)
10    Aj = accuracy_score(y_val, pred)
11    importances.append((all_feat[j], A_base - Aj))
12
13 # Sort by importance
14 importances.sort(key=lambda x: x[1], reverse=True)
15 top5 = importances[:5]
16
17 print("\nTop 5 Most Important Features:")
18 for name, imp in top5:
19     print(f"{name:35s} ΔA = {imp}")
```

The code presented computes the permutation feature importance by measuring how much the model's accuracy decreases when each feature is randomly permuted. First, the validation data are transformed using the preprocessing step of the pipeline so that the MLP receives inputs in the same format as during training. Then, for each feature, a copy of the transformed validation set is created and the values of that feature are randomly shuffled, breaking its relationship with the target variable. The modified dataset is passed to the trained MLP to obtain new predictions, and the resulting accuracy is compared with the original validation accuracy. The drop in accuracy quantifies the importance of that feature: a larger decrease indicates that the feature is more critical for the model's performance. Finally, all features are sorted by their impact, and the five most important ones are printed.

Q7: Principal Component Analysis

Below we provide you with a plot of eigenvalues versus component index and a PC1–PC2 scatter plot coloured by predicted killer labels, as requested.



Explanation

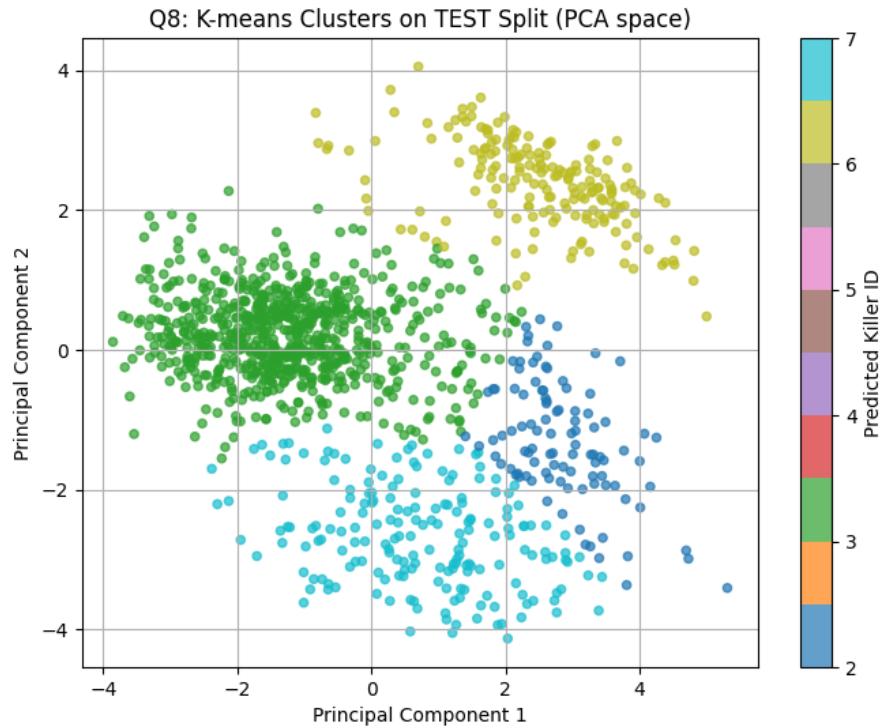
Before applying PCA, we started off by preprocessing the features this way once again ensuring no data leakage.

```
1 pre = ColumnTransformer([
2     ("scaler", StandardScaler(), CONT),
3     ("onehot", OneHotEncoder(handle_unknown="ignore", sparse_output=False),
4         CAT)
])
```

Subsequently, PCA was fitted on TRAIN only and then applied to both TRAIN and VAL. Firstly, we let PCA run without any restrictions on the number of components. However, for visualization purposes only we chose to set the parameter n_components to 2 to achieve projection in R^2 .

Q8: k -means in PCA space

Below, we provide you with the plot with our final predicted killer labels C_i .



As you can see we came to the conclusion that there are 4 predicted killers in k -means clustering. These are killers 2, 3, 6 and 7. We would like to remind you that in Q3 we had predicted killers 3, 6 and 7. This means that even though the VAL accuracy of k -means is equal to 81%, it is still very close to the Gaussian Bayes and SVM models which have better accuracies.

Explanation

To start with, the data was once again preprocessed, fitted on TRAIN and applied on VAL and TEST. The PCA was applied on TRAIN only. The number of components was set to 2 (based on Q7's results). So, we ran k-means with k equal to S, where equals to the nnumber of killers(in our case 8). K-means is unsupervised, so the clusters do not correspond directly to the killer IDs. In order to convert clusters into class labels, we used majority voting on the TRAIN set. For each cluster, we identified all TRAIN samples assigned to that cluster, examined their true killer labels, and assigned each cluster the label of the most frequent killer in that cluster. For evaluation on VAL, samples were assigned to clusters using the fitted k-means model, mapped to killer IDs using the majority-vote mapping, and accuracy was of course calculated. Finally, for the TEST set, we produced a probability distribution over the clusters by calculating the Euclidean distance from each sample to all centroids. Distances were converted to scores using score = -distance, and a softmax approach was applied to obtain normalized probabilities.

```
1 test_clusters = kmeans.predict(Z_test_m)
2 test_pred = np.array([cluster_to_killer[c] for c in test_clusters])
3
4 distances = kmeans.transform(Z_test_m)
5
6 # 2. Μετατροπή αποστάσεων σε πιθανότητες (Softmax approach)
7 # Χρησιμοποιούμε exp(-dist) ώστε οι μικρές αποστάσεις να δίνουν μεγάλες
    ↪ πιθανότητες
8 inv_distances = np.exp(-distances)
9 # Κανονικοποίηση ώστε το άθροισμα κάθε γραμμής να είναι 1 [cite: 91]
10 probs_clusters = inv_distances / np.sum(inv_distances, axis=1,
    ↪ keepdims=True)
11
12 # 3. Δημιουργία DataFrame
13 submission = pd.DataFrame({"incident_id": test["incident_id"],
    ↪ "predicted_killer": test_pred})
14
15 # 4. Υπολογισμός πιθανοτήτων για κάθε killer_id (1 έως S)
16 for k in range(1, S + 1):
    # Βρίσκουμε ποια clusters (q) αντιστοιχούν στον δολοφόνο k [cite: 218]
17     matching_clusters = [q for q, killer in cluster_to_killer.items() if
        ↪ killer == k]
18
19     if matching_clusters:
        # Αθροίζουμε τις πιθανότητες των clusters που ανήκουν στον δολοφόνο
        ↪ k
        submission[f"p_killer_{k}"] = probs_clusters[:, matching_clusters].sum(axis=1)
20     else:
        # Αν ο δολοφόνος k δεν κέρδισε κανένα cluster, του δίνουμε μια πολύ
        ↪ μικρή
        # βασική πιθανότητα (π.χ. από το πλησιέστερο cluster) για να μην
        ↪ είναι 0
        submission[f"p_killer_{k}"] = 1e-5
21
22     # Επανεξισορρόπηση (Normalization) για να αθροίζουν ακριβώς στο 1
        ↪ [cite: 91]
23 prob_cols = [f"p_killer_{k}" for k in range(1, S + 1)]
24 submission[prob_cols] =
    ↪ submission[prob_cols].div(submission[prob_cols].sum(axis=1), axis=0)
```

This code for the TEST set predicts clusters and labels similarly, then derives class probabilities by first obtaining distances of each example to all cluster centroids (`kmeans.transform`). Distances are converted to soft probabilities with a softmax over negative distances ($\exp(-\text{distance})$), so nearer centroids receive higher probability mass. These cluster-level probabilities are aggregated to class-level probabilities by summing over all clusters mapped to the same killer ID. The classes with no winning cluster get a tiny fallback probability to avoid exact zeros($1e-5$). Finally, the per-row probabilities are renormalized to sum to 1, and a `submission.csv` is created containing the `incident_id`, the predicted label, and the calibrated probabilities `p_killer_1` ... `p_killer_S`.

4 Discussion and conclusions

Let's summarise the VAL accuracies of these methods up until now:

Table 5: Validation accuracy for different models.

Model	VAL accuracy
Gaussian Bayes	0.905
Linear classifier	0.782
SVM (RBF kernel)	0.948
MLP	0.944
PCA + k-means	0.813

So, in ascending order of performance we have:

1. SVM (Q5)
2. MLP (Q6)
3. Gaussian Bayes (Q3)
4. Linear classifier (Q4)
5. PCA + k-means (Q8)

We have concluded that the more helpful features in order to model behavior were geographic features (latitude and longitude), victim-related information (their age and gender), information about the place where the crime was committed and the prevailing weather conditions (weather, temperature and scene type) and of course the weapon type (not necessarily in that order). With the use of all presented diagrams, confusion matrixes and the final file submision.csv, we highlighted the 4 killers (2, 3, 6 and 7) as the killers who are most likely to have committed the most crimes, especially in TEST. Potentially, more complex generative models could be included and explored. Additionally, there are other factors to be analyzed and taken into consideration, such as crime frequency patterns, killer's motive or mental health disorders. Overall, the current study highlights the importance of combining proper preprocessing, nonlinear modeling and feature analysis when trying to solve real-world pattern recognition problems.

Code organisation

We implemented the code in a modular way, creating one script per question. Each script is responsible for producing the results and diagrams of its corresponding question. The diagrams from each question are stored in a separate folder. The dataset is also included in the main folder. Lastly, the prediction file `submission.csv` is also included in the main folder.

```
pattern_recognition
├── Q1_Diagrams
├── Q2_Diagrams
├── Q3_Diagrams
├── Q4_Diagrams
├── Q5_Diagrams
├── Q6_Diagrams
├── Q7_Diagrams
└── Q8_Diagrams
├── crimes.csv
├── q1.py
├── q2.py
├── q3.py
├── q4.py
├── q5.py
├── q6.py
├── q7.py
└── q8.py
└── submission.csv
```

Libraries and data load

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import norm
from sklearn.mixture import GaussianMixture

# Φόρτωση & φιλτράρισμα
csv_path = "crimes.csv"
df = pd.read_csv(csv_path)

# Κρατάμε μόνο TRAIN+VAL
df = df[df["split"].isin(["TRAIN", "VAL"])].copy()

# Τυπική επιβεβαίωση ότι έχουμε όλες τις απαιτούμενες στήλες
required_cols = ["hour_float", "victim_age", "latitude",
                  "longitude"]
missing = [c for c in required_cols if c not in df.columns]
assert not missing, f"Δείπουν στήλες: {missing}"
```

In this specific script we load the required libraries and the data from the `crimes.csv` file. We split the data into TRAIN and VAL as it is required for Q1, a similar method is used for all the other questions. Lastly, we confirm that all the required columns are present in the dataset. This is a common practice to avoid errors later on.

Briefly describe your Python code structure (scripts, notebooks, helper modules) and how to reproduce:

- the main results and figures in this report;
- the prediction file `submission.csv`.

5 LLM prompts and responses

The contribution of an AI, specifically Copilot, in this project was purely auxiliary. In general, we consulted Copilot in order to provide us with formulas, types, libraries and assistance in diagrams' creation. In particular, we list all prompts and the corresponding responses we used during the preparation of this project.

Comment: As far as Q3 is concerned, we asked Copilot how we can create a 2D PCA projection, as we were facing some problems. Also, we asked how we can include a legend box that would contain the killers and their corresponding colors. Unfortunately, we are unable to find the corresponding prompt and response.

In Q4 we asked Copilot how we can make a 2D PCA projection that overlays the one created in Q3.

Prompt: Πως μπορούν να δημιουργηθούν PCA και πως μπορώ να τα κάνω overlay?

Response:

Καλή ερώτηση — πάμε καθαρά, περιεκτικά και πρακτικά, όπως θα το έγραφες και σε εργασία. 📈

Θα σου εξηγήσω:

- πώς δημιουργούνται τα PCA (Θεωρία + βήματα)
- πώς τα υπολογίζεις στην πράξη (Python)
- πώς κάνεις overlay (π.χ. διαφορετικά σύνολα / predictions στο ίδιο PCA plot)

1 Πώς δημιουργούνται τα PCA (λογική – βήματα)

To PCA (Principal Component Analysis) είναι γραμμικός μετασχηματισμός που:

- μειώνει τη διάσταση των δεδομένων
- διατηρεί όσο το δυνατόν περισσότερη διασπορά (variance)

Βήματα PCA

- Σύλλογη δεδομένων
 - Πίνακας $X \in \mathbb{R}^{N \times d}$
- Κανονικοποίηση (απαραίτητη)

2. Κανονικοποίηση (απαραίτητη)

$$X_{std} = \frac{X - \mu}{\sigma}$$

3. Υπολογισμός συνδιακύμανσης

$$\Sigma = \frac{1}{N} X_{std}^T X_{std}$$

4. Ιδιοτιμές & ιδιοδιανύσματα

$$\Sigma v_i = \lambda_i v_i$$

5. Ταξινόμηση ιδιοδιανυσμάτων

- κατά φθινουσα ιδιοτιμή

6. Projection

$$Z = X_{std} W$$

όπου $W = [v_1, v_2, \dots, v_m]$

2 Πώς υλοποιείς PCA στην πράξη (σωστός τρόπος για εργασία)

♦ PCA μόνο στο TRAIN (ΠΟΛΥ ΣΗΜΑΝΤΙΚΟ)

```
python
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
```

⌚ Αντιγραφή κώδικα

python

```
X_train = train[features].values
X_val = val[features].values
```

⌚ Αντιγραφή κώδικα

♦ Standardization

```
python
scaler = StandardScaler()
X_train_std = scaler.fit_transform(X_train)
X_val_std = scaler.transform(X_val)
```

⌚ Αντιγραφή κώδικα

! ΝΟΤΕ fit στο VAL ή TEST

♦ PCA fit στο TRAIN

```
python
pca = PCA(n_components=2)
z_train = pca.fit_transform(X_train_std)
z_val = pca.transform(X_val_std)
```

⌚ Αντιγραφή κώδικα

♦ Τι κρατάς για ανάλυση

```
python
pca.components_ # κύριοι οξενες
pca.explained_variance_ratio_
pca.explained_variance_
```

⌚ Αντιγραφή κώδικα

3 Πώς κάνεις PCA overlay (το σημαντικότερο)

🔥 Τι σημαίνει overlay

• ίδιο PCA space

• διαφορετικά σύνολα ή labels

• κοινό σύστημα αξόνων

● Overlay TRAIN vs VAL

```
python
plt.figure(figsize=(6,6))

plt.scatter(z_train[:,0], z_train[:,1],
            c='blue', alpha=0.4, label='TRAIN')
plt.scatter(z_val[:,0], z_val[:,1],
            c='red', alpha=0.4, label='VAL')

plt.xlabel('PCA 1')
plt.ylabel('PCA 2')
plt.title('PCA overlay: TRAIN vs VAL')
plt.show()
```

● Overlay true labels vs predictions

```
python
plt.figure(figsize=(6,6))

plt.scatter(z_val[:,0], z_val[:,1],
            c=y_val, cmap='tab10',
            alpha=0.6)

plt.title('PCA (colored by true killer)')
plt.xlabel('PCA 1')
plt.ylabel('PCA 2')
plt.show()

python
plt.figure(figsize=(6,6))

plt.scatter(z_val[:,0], z_val[:,1],
            c=y_pred, cmap='tab10',
            alpha=0.6)
```

⌚ Αντιγραφή κώδικα

```

plt.title("PCA (colored by predicted killer")
plt.xlabel("PC1")
plt.ylabel("PC2")
plt.show()

👉 Autό είναι overlay σύγκρισης (ιδιο PCA space)

```

● Overlay πολλών μοντέλων

```

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

plt.scatter(z_val[:,0], z_val[:,1],
            c=y_pred_svm, marker='o', label='SVM', alpha=0.5)

plt.scatter(z_val[:,0], z_val[:,1],
            c=y_pred_bayes, marker='x', label='Bayes', alpha=0.5)

plt.legend()
plt.title("Overlay predictions on PCA")
plt.show()

```

● Πώς το περιγράφεις σε εργασία

«Το PCA υπολογίστηκε αποκλειστικά στο σύνολο εκπαίδευσης μετά από κανονικοποίηση. Το ίδιο PCA transformation εφαρμόστηκε στα σύνολα validation/test, επιτρέποντας άμεσο οπτικό συγκριτικό έλεγχο (overlay) των δεδομένων στο ίδιο υποχώριο.»

⚠ Κλασικά λάθη (να τα αποφύγεις)

- ✗ PCA fit σε όλο το dataset
- ✗ PCA χωρίς standardization
- ✗ Διαφορετικό PCA για κάθε split
- ✗ σύγκριση διαφορετικών PCA axes

✳ Σύνδεση με τα Q της εργασίας

- Q3-Q6: PCA overlay για decision regions
- Q7: Scree plot + επιλογή m
- Q8: PCA → k-means → overlay clusters

Prompt: Τί είναι το rbf kernel και πως μπορώ να το χρησιμοποιήσω σε python?
Response:

Καλή ερώτηση — το RBF kernel είναι από τα πιο βασικά και εξεταστικά θέματα σε SVM. Πάμε ξεκάθαρα, δομήμενα και με παράδειγμα Python :

❶ Τί είναι το RBF kernel (θεωρία)

RBF (Radial Basis Function) kernel = Gaussian kernel

Χρησιμοποιείται στα Support Vector Machines για να:

- χειρίζονται μη γραμμικά διαχωρισμά δεδομένα
- μεταφέρουν τα δεδομένα σε υψηλότερη διάσταση όπου ο διαχωρισμός γίνεται γραμμικός

❷ Μαθηματικός ορισμός

$$K(x, x') = \exp(-\gamma \|x - x'\|^2)$$

όπου:

- x, x' : δύο σημεία δεδομένων
- $\| \cdot \|$: Ευκλείδεια απόσταση
- $\gamma > 0$: παράμετρος πλάτους Gaussian

❸ Διαισθητική ερμηνεία

- Κοντινά σημεία \rightarrow kernel ≈ 1
- Μακρινά σημεία \rightarrow kernel ≈ 0
- Δημιουργεί τοπικές περιοχές επιρροής
- Επιτρέπει καμπύλες (μη γραμμικές) αποφάσεις

❹ Πότε χρησιμοποιούμε RBF kernel

- ✓ Όταν τα δεδομένα δεν διαχωρίζονται γραμμικά
- ✓ Όταν δεν γνωρίζουμε τη συστήμα γραμμικότητα
- ✓ Όταν τα features είναι συνεχή
- ✓ Default επιλογή σε SVM

❺ Παράμετροι RBF SVM (πολύ σημαντικό)

- C – Regularization
- Μεγάλο C \rightarrow αυστηρός διαχωρισμός (overfitting)
- Μικρό C \rightarrow πιο «μαλακά» δρώσι (underfitting)

❻ γ (gamma)

- Μεγάλο γ \rightarrow πολύ στενές Gaussian \rightarrow overfitting
- Μικρό γ \rightarrow πολύ φερδιές Gaussian \rightarrow underfitting

❼ Πώς χρησιμοποιώ RBF kernel σε Python (scikit-learn)

• Βασική χρήση

```

python
from sklearn.svm import SVC

svm_rbf = SVC(kernel='rbf', C=1.0, gamma='scale')
svm_rbf.fit(X_train, y_train)

y_pred = svm_rbf.predict(X_val)

```

❽ Τι σημαίνει `gamma='scale'`

❽ Τι σημαίνει `gamma='scale'`

$$\gamma = \frac{1}{d \cdot \text{Var}(X)}$$

✓ ασφαλής default επιλογή

❾ Σωστή ροή για εργασία / εξετάσεις

• Standardization (ΑΠΑΡΑΙΤΗΤΟ)

```

python
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
X_train_std = scaler.fit_transform(X_train)
X_val_std = scaler.transform(X_val)

```

◆ Εκπαίδευση SVM

```
python
svm = SVC(kernel='rbf', C=10, gamma=0.1)
svm.fit(X_train_std, y_train)
```

Αντιγραφή κώδικα

◆ Grid Search για C και γ (όπως ζητάει εργασία)

```
python
from sklearn.model_selection import GridSearchCV

param_grid = {
    'C': [0.1, 1, 10, 100],
    'gamma': [0.01, 0.1, 1]
}

svm = SVC(kernel='rbf')

grid = GridSearchCV(svm, param_grid, cv=5)
grid.fit(X_train_std, y_train)

best_model = grid.best_estimator_
```

Αντιγραφή κώδικα

7 Οπτικοποίηση RBF decision regions σε PCA (overlay)

```
python
Z_train = pca.fit_transform(X_train_std)

svm.fit(Z_train, y_train)

# μετά grid για decision boundary
```

Το RBF δημιουργεί καμπύλα decision boundaries στο PCA χώρο

Prompt: Στο q6, τι είναι το baseline accuracy και πως κάνουμε shuffle?

Response:

Πολύ καλή ερώτηση — στο Q6 αυτό είναι κλασικό σημείο που μπερδεύει. Πάμε ξεκάθαρα, εξεταστικά και πρακτικά, ακριβώς όπως ζητάει η εκφύνηση 📈

1 Τι είναι το baseline accuracy (Q6)

👉 Ορισμός

Baseline accuracy = η αρχική απόδοση του εκπαιδευμένου μοντέλου, χωρίς καμία τροποποίηση στα δεδομένα.

Στο Q6:

- είναι η accuracy του MLP στο validation set
- με όλα τα features όπως είναι

• Γιατί τη χρειαζόμαστε

Χρησιμοποιείται ως:

- σημείο αναφοράς
- για να μετρήσουμε πόσο σημαντικό είναι κάθε feature

Αν χαλάσαι ένα feature και η accuracy πέσει πολύ → το feature είναι σημαντικό

♦ Παράδειγμα

text

Baseline accuracy = 0.78

Αντιγραφή κώδικα

2 Τι είναι το Permutation Feature Importance

Στο Q6 δεν ζητά:

- weights
- gradients
- SHAP

Ζητά Permutation Importance, δηλαδή:

«Χαλάμε» ένα feature ανακατεύοντας τις τιμές της και βλέπουμε πόσο πέφτει η accuracy.

4 Τι σημαίνει shuffle feature

- Shuffle = τυχαιά αναδιάταξη των τιμών
- κρατάει ίδια κατανομή
- καταστρέφεις τη σχέση feature ↔ label
- Παρόδειγμα

```
text
python
import numpy as np

drops = []

for j in range(X_val.shape[1]):
    X_shuffled = X_val.copy()

    # shuffle MNO to feature j
    np.random.shuffle(X_shuffled[:, j])

    acc = accuracy_score(y_val, model.predict(X_shuffled))
    drops.append(baseline_acc - acc)
```

Πριν: [22, 35, 41, 29]
μετά: [41, 22, 29, 35]

5 Πώς γίνεται το shuffle (βήμα-βήμα)

- Βήμα 1: Υπολόγιση baseline

```
python
baseline_acc = accuracy_score(y_val, model.predict(X_val))
```

6 Βήμα 2: Ήταν κάθε feature

```
python
import numpy as np

drops = []

for j in range(X_val.shape[1]):
    X_shuffled = X_val.copy()

    # shuffle MNO to feature j
    np.random.shuffle(X_shuffled[:, j])

    acc = accuracy_score(y_val, model.predict(X_shuffled))
    drops.append(baseline_acc - acc)
```

Αντιγραφή κώδικα

- Τι υπολογίζεις

importance_j = baseline acc_↓ - accuracy μετά το shuffle

5 Πώς το παρουσιάζεις στο Q6

- Πίνακας

Feature	Accuracy μετά shuffle	Drop
hour_float	0.61	0.17
longitude	0.73	0.05

6 Bar plot (ζητείται)

```
python
plt.barh(feature_name, drops)
plt.xlabel("Accuracy drop")
plt.title("Permutation Feature Importance")
plt.show()
```

7 Τι ΔΕΝ είναι baseline accuracy (συχνά λάθος)

- ✗ accuracy random classifier
- ✗ accuracy majority class
- ✗ accuracy στο TRAIN
- ✗ accuracy μετά από shuffle
- ✓ μόνο: validation accuracy του κανονικού μοντέλου

Το Q6 ζητά:

- MLP
- σύγκριση με άλλα μοντέλα
- feature importance με permutation
- bar plot top-5 features

Όλα βασίζονται στο baseline accuracy

Prompt: Αντός ο κώδικας απαντά στο ερώτημα;

```
1 acc = accuracy_score(y_val, pred_val)
2 print(f"VAL Accuracy: {acc}")
3
4 # =====
5 # Permutation Feature Importance
6 # =====
7
8 # Baseline
9 accuracy_A_base = acc
10
11 # Get feature names after one-hot encoding
12 ohe = model.named_steps["pre"].named_transformers_["onehot"]
13 ohe_names = ohe.get_feature_names_out(CAT)
14 all_feat = CONT + list(ohe_names)
15
16 # Transform VAL once
17 Z_val = model.named_steps["pre"].transform(X_val)
18
19 importances = [] rng = np.random.default_rng(0)
20
21 for j in range(Z_val.shape[1]):
22     Zp = Z_val.copy()
23     Zp[:, j] = rng.permutation(Zp[:, j]) # permute one feature
24     pred = model.named_steps["mlp"].predict(Zp)
25     Aj = accuracy_score(y_val, pred)
26     importances.append((all_feat[j], A_base - Aj))
27
28 # Sort by importance
29 importances.sort(key=lambda x: x[1], reverse=True)
```

```

30 top5 = importances[:5]
31
32 print("\nTop 5 Most Important Features:")
33 for name, imp in top5:
34     print(f"{name:35s} ΔA = {imp}")
35
36 # =====
37 # Plot top 5
38 # =====
39
40 names = [x[0] for x in top5]
41 values = [x[1] for x in top5]
42
43 plt.figure(figsize=(7,4))
44 plt.barh(names, values, color="skyblue")
45 plt.gca().invert_yaxis()
46 plt.title("Top 5 Most Important Features (MLP)")
47 plt.xlabel("Drop in Accuracy (ΔA)")
48 plt.tight_layout()
49 plt.savefig("q6_top5_features.png", dpi=150)

```

Response:

Τι ζητάει το Q6 και αν το καλύπτεις

- ♦ Q6 ζητά:
 1. Baseline accuracy στο VAL
 2. Permutation (shuffle) κάθε feature
 3. Υπολογισμό πτώσης accuracy
 4. Ranking features
 5. Bar plot για top-5
- Ο κώδικας σου τα κάνει ΌΛΑ.

Ανάλυση του κώδικα (γιατί είναι σωστός)

Baseline accuracy

```
python
acc = accuracy_score(y_val, pred_val)
A_base = acc
```

Σωστό:

- Accuracy στο VAL
- Χωρίς τρωτοποίηση δεδομένων
- Αυτό είναι ακριβώς το baseline accuracy

Feature names μετά το one-hot

```
python
ohe = model.named_steps["pre"].named_transformers_["onehot"]
ohe_names = ohe.get_feature_names_out(CAT)
all_feat = CONT + list(ohe_names)
```

Πολύ σωστό και προσεγγένειο:

- Κρατάς σωστά ονόματα features
- Περιλαμβάνει continuous + one-hot categorical
- Απαιτείται για σωστό report

Transform VAL μία φορά (σωστό)

```
python
Z_val = model.named_steps["pre"].transform(X_val)
```

Πολύ σωστό:

- Δεν ξανακάνει fit
- Όδιο preprocessing με training
- Αποφεύγεις data leakage

Permutation (shuffle) feature-by-feature

```
python
for j in range(Z_val.shape[1]):
    Zp = Z_val.copy()
    Zp[:, j] = np.random.permutation(Zp[:, j])
```

Αυτό είναι ακριβώς το permutation importance:

- Shuffle μόνο ένα feature

Όλα τα άλλα μένουν σταθερά

- Καταστρέφεται η σχέση feature ↔ label

Accuracy μετά το shuffle

```
python
pred = model.named_steps["mlp"].predict(Zp)
Aj = accuracy_score(y_val, pred)
importances.append((all_feat[j], A_base - Aj))
```

Μαθηματικά σωστό:

$$\Delta A_j = A_{\text{baseline}} - A_{\text{shuffled}}$$

Ranking & Top-5

```
python
importances.sort(key=lambda x: x[1], reverse=True)
top5 = importances[:5]
```

Bar plot (όπως ζητείται)

```
python
plt.barh(names, values)
plt.title("Top 5 Most Important Features (MLP)")
```

✓ Καλύπτει πλήρως το ζητούμενο visualization

Τελική επιμηγορία

- NAI: ο κώδικάς σου απαντά πλήρως και σωστά στο Q6
- ✓ Ορθός οριζόντιος baseline accuracy
- ✓ Σωστό permutation feature importance
- ✓ Σωστό ranking & plot
- ✓ Report-ready

Μικρές παρατηρήσεις (προαιρετικές, όχι υποχρεωτικές)

◆ (1) Reproducibility (το έχεις ήδη 🎉)

```
python
rng = np.random.default_rng(0)
```

✓ Πολύ καλό — μπορείς να το αναφέρεις στο report

◆ (2) Θεωρητική διευκρίνιση (μόνο στο κείμενο)

Μπορείς να γράψεις:

«Για τα one-hot encoded χαρακτηριστικά, το *permutation importance* αναφέρεται σε μεμονωμένες διαδικασίες κατηγορίες.»

(δεν χρειάζεται αλλαγή κώδικα)