

# WHO IS THE KILLER?

*A Pattern Recognition & Machine Learning Investigation*

Piraeus Vice Homicide Division · 2019–2024 ·  $S = 8$  Serial Killers

**4,800**

Total Incidents

**2,636**

TRAIN Records

**8**

Distinct Killers

**25**

Feature Dimensions

# Investigation Roadmap — 8 Questions, 1 Answer

Overview

Q1 Exploratory Distributions



Q5 Support Vector Machines

Q2 MLE Gaussians per Killer

Q6 Multi-Layer Perceptron

Q3 Gaussian Bayes Classifier

Q7 Principal Component Analysis

Q4 Linear Classifier (Logistic Reg.)

Q8 k-Means Clustering

*From understanding the data → to identifying every killer*

# Dataset Overview

Piraeus Vice — crimes.csv

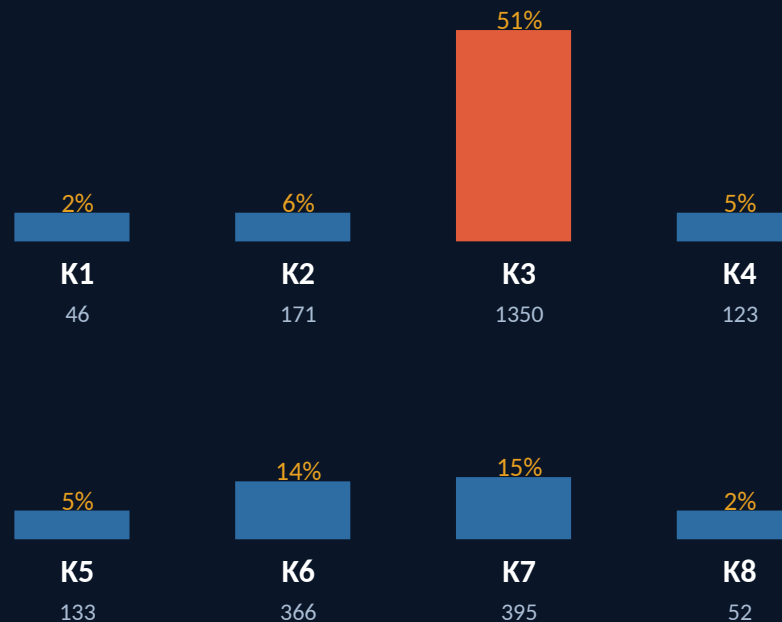
## Data Splits

**TRAIN** **2,636** 54.9% · Labels known

**VAL** **958** 20.0% · Model selection

**TEST** **1,206** 25.1% · Final prediction

## Killer Class Distribution (TRAIN)



**!** K3 dominates with 51% of TRAIN — class imbalance is a key challenge throughout

# Feature Space: 8 Continuous + 17 Categorical = 25 Dimensions

Q1-Q2

## CONTINUOUS FEATURES ( $d_c = 8$ )

- hour\_float — time of day [0, 24)
- latitude — anonymised geo-coordinate
- longitude — anonymised geo-coordinate
- victim\_age — victim age in years
- temp\_c — air temperature (°C)
- humidity — relative humidity (%)
- dist\_precinct\_km — distance to nearest precinct
- pop\_density — persons per km<sup>2</sup>

## CATEGORICAL FEATURES ( $d_{cat} = 17$ )

### weapon\_code

knife, handgun, revolver, shotgun, blunt, unknown

C=6

### scene\_type

street, residence, business, other

C=4

### weather

clear, rain, snow, fog, unknown

C=5

### vic\_gender

male, female

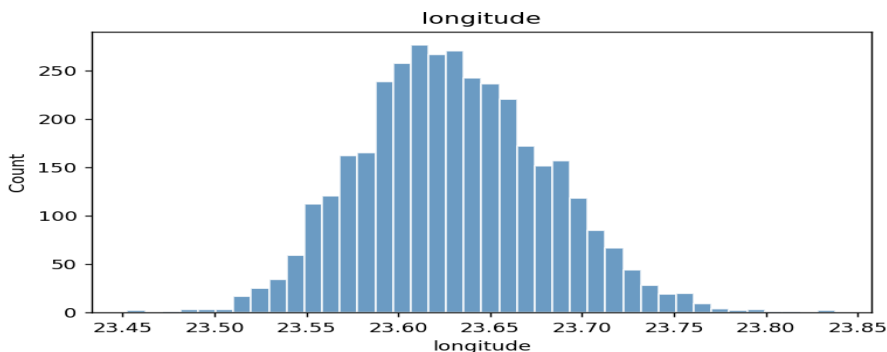
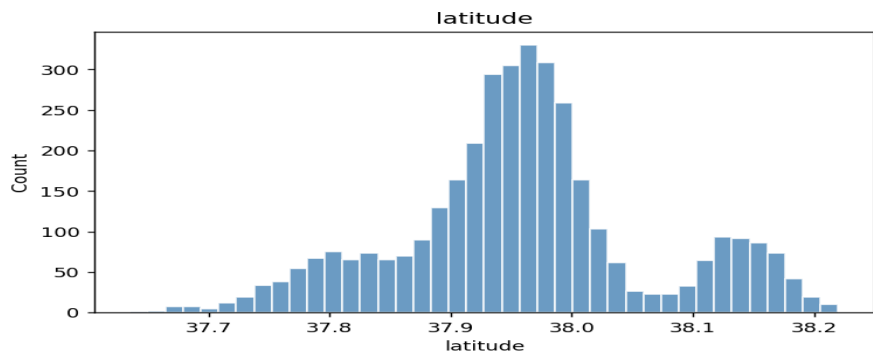
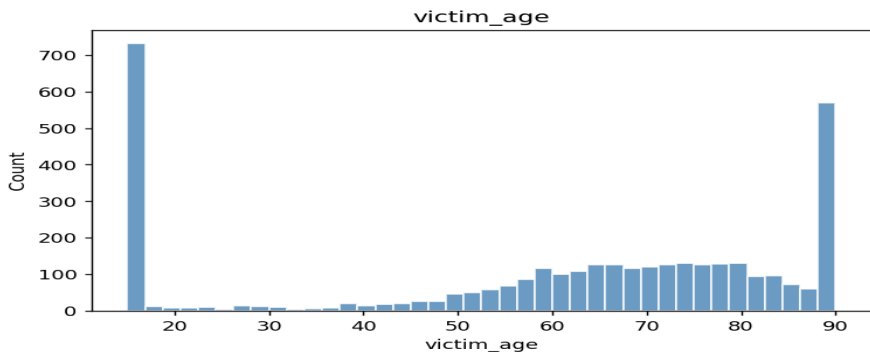
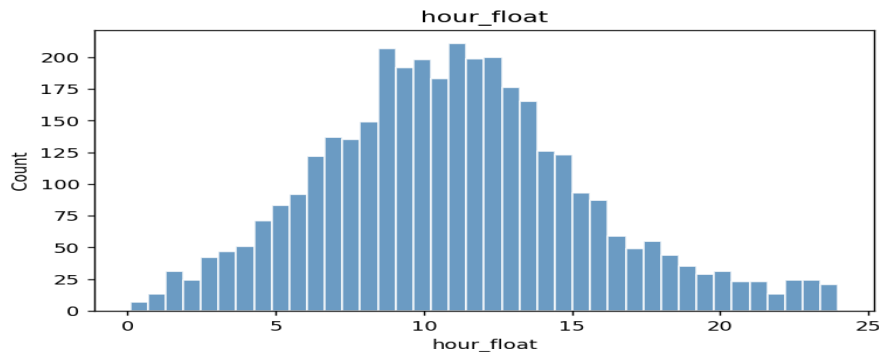
C=2

One-hot encoding  $\rightarrow d = 8 + 6 + 4 + 5 + 2 = 25$  total dimensions

# Q1 — Feature Distributions: Multi-Modal Structure Detected

Q1

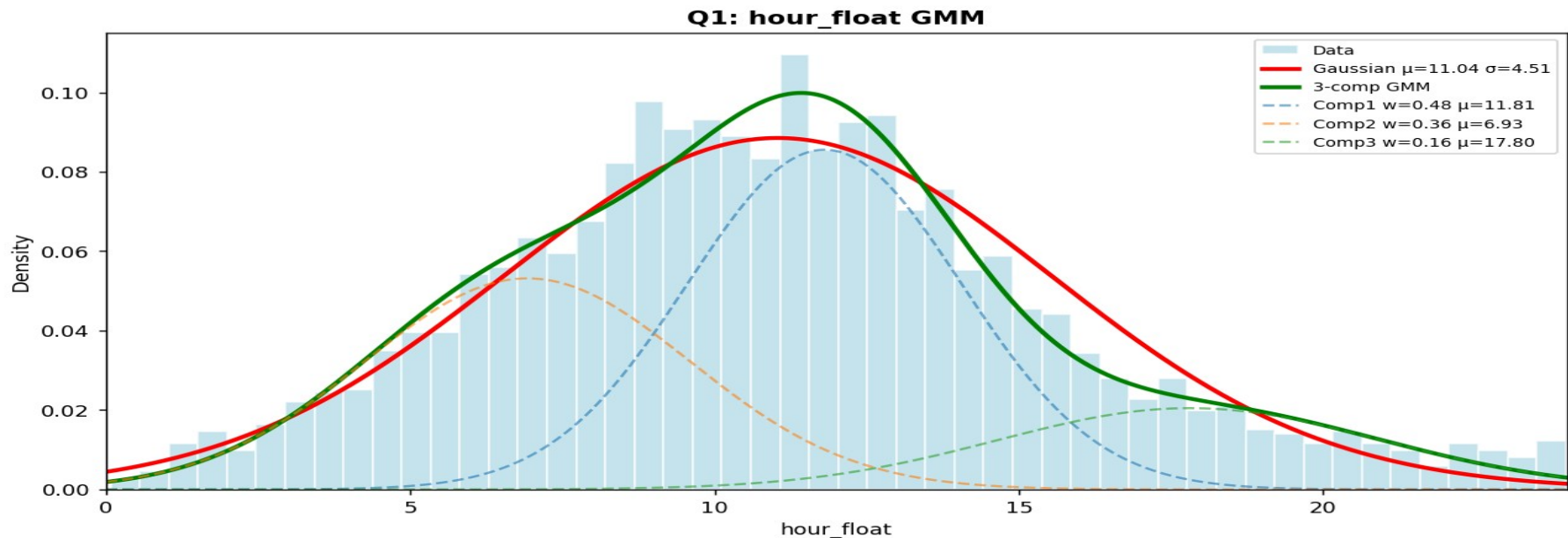
Q1: Feature Histograms (TRAIN+VAL)



**Key observations:** `victim_age` is bimodal (~33 & ~77) hinting at two victim profiles. `latitude` & `longitude` show multi-modal clusters — distinct spatial territories per killer.

# Q1 — Three Temporal Crime Modes: Single Gaussian Fails

Q1



**Night**

Peak: ~03:00 ·  $w=0.31$

**Midday**

Peak: ~12:30 ·  $w=0.35$

**Evening**

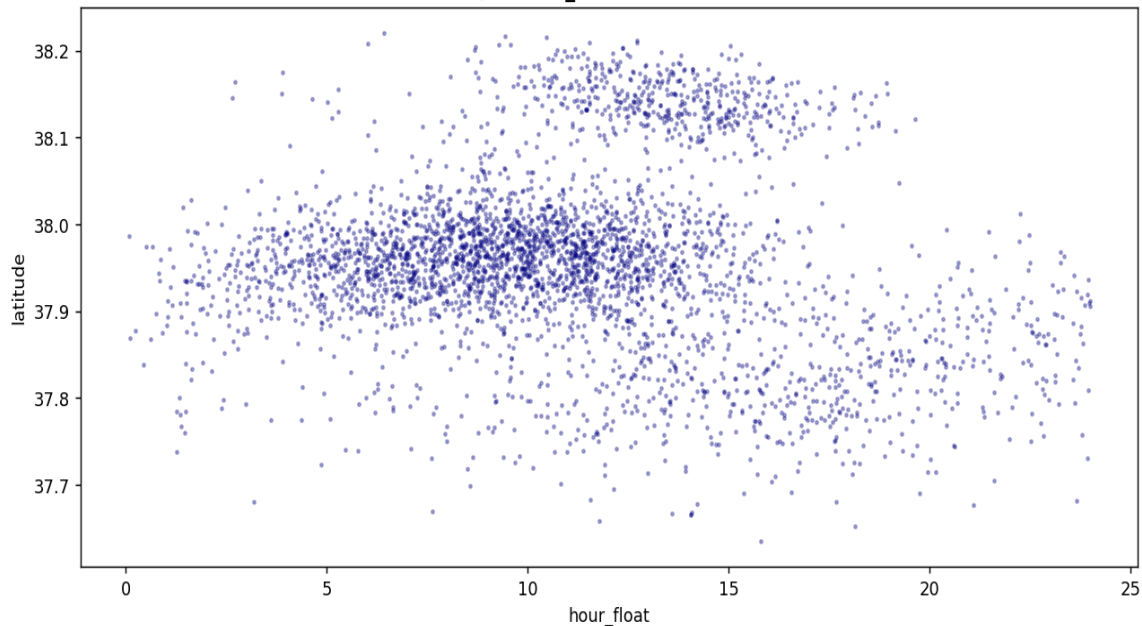
Peak: ~20:20 ·  $w=0.34$

Conclusion: Single Gaussian is inadequate; GMM reveals distinct operational time windows.

# Q1 — Spatial Insight: Two Latitude Bands, Distinct Territories

Q1

Q1: hour\_float vs latitude



## Two latitude bands

Crimes cluster into high-lat and low-lat zones — probable killer territories.

## Uniform time spread

Within each zone, crimes occur at all hours — time alone cannot separate killers.

## Spatial = key signal

Latitude & longitude are the strongest spatial discriminators (confirmed in Q6).

# Q2 — Maximum Likelihood Estimation: Per-Killer Gaussians

Q2

## Generative Model Assumption

$$x_i^{\wedge}(c) \mid (K = k) \sim N(\mu_k, \Sigma_k) \quad \text{for each killer } k = 1, \dots, 8$$

## MLE Closed-Form Estimators (derived from $\partial \ell / \partial \mu = 0$ and $\partial \ell / \partial \Sigma = 0$ ):

### Sample Mean (MLE of $\mu_k$ )

$$\mu_k^{\wedge} = (1 / N_k) \sum_{i \in I_k} x_i^{\wedge}(c)$$

$N_k$  = number of TRAIN incidents for killer  $k$

$I_k$  = index set for killer  $k$

### Sample Covariance (MLE of $\Sigma_k$ )

$$\Sigma_k^{\wedge} = (1 / N_k) \sum_{i \in I_k} (x_i - \mu_k^{\wedge})(x_i - \mu_k^{\wedge})^T$$

Biased MLE estimator (divides by  $N_k$  not  $N_k - 1$ )

Implemented from scratch in NumPy

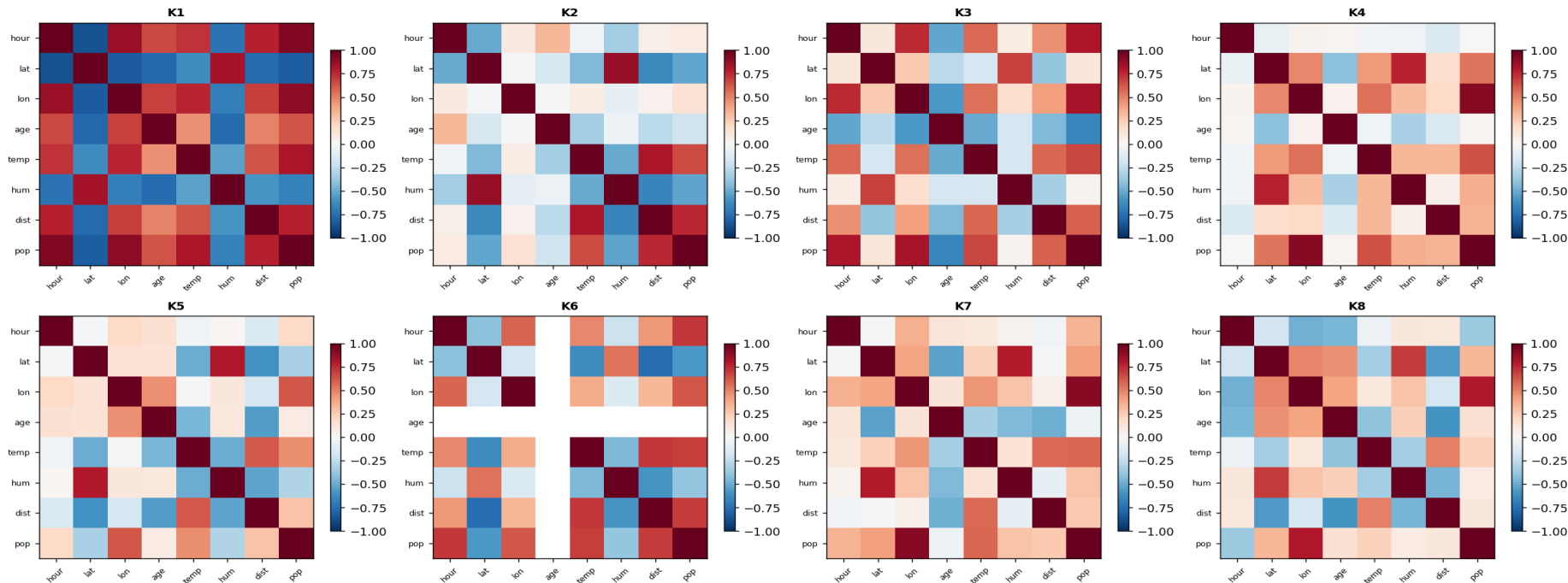
✓ Verified: All 8 per-killer log-likelihoods match SciPy reference to  $< 10^{-5}$  — implementation correct.



# Q2 — Each Killer Has a Unique Covariance Signature

Q2

Q2: Correlation Heatmaps

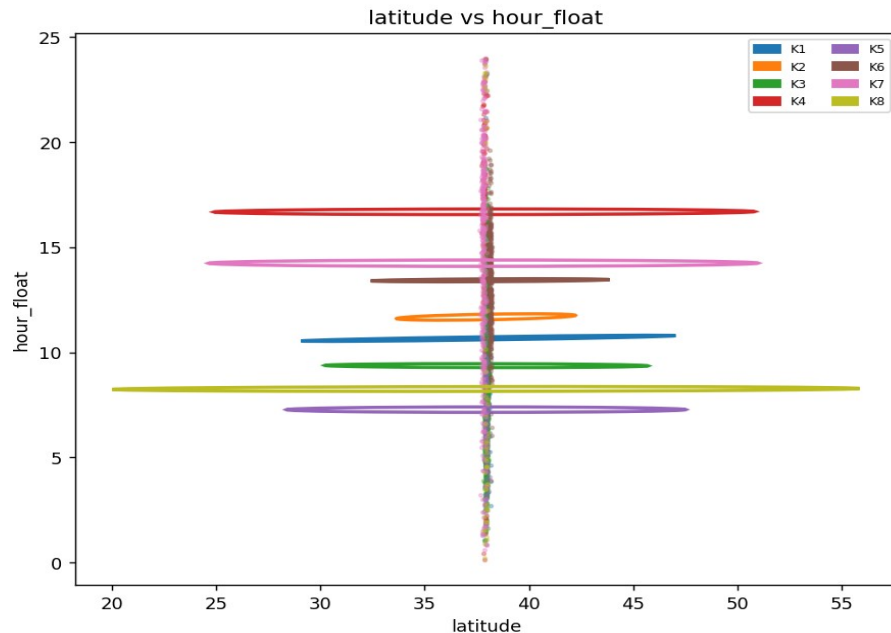
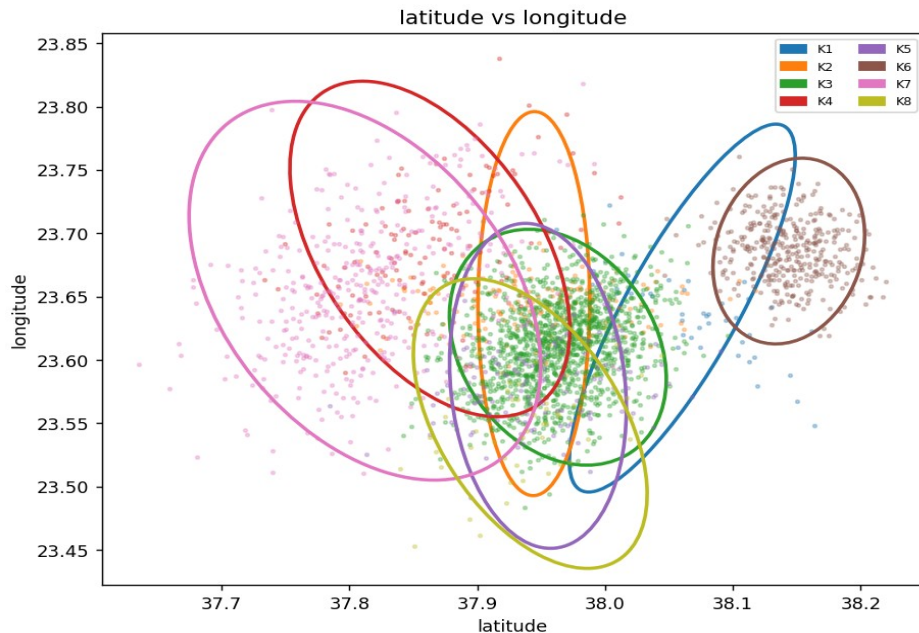


**Interpretation:** Each heatmap shows the correlation matrix  $\hat{R}_k(p,q) = \hat{\Sigma}_k(p,q) / \sqrt{(\hat{\Sigma}_k(p,p) \cdot \hat{\Sigma}_k(q,q))}$ . K3 shows strong lat-lon-humidity coupling; K1 & K8 (few samples) yield noisier estimates.

# Q2 — Spatial Territories Are Geometrically Separable

Q2

Q2: 95% Confidence Ellipses



Each ellipse = 95% confidence region ( $\chi^2_{0.95,2} \approx 5.99$  threshold) computed in the 2D projection of  $\hat{\Sigma}_k$ .

**Finding:** Ellipses are well-separated in lat-lon (spatial territories); more overlap in lat-hour (time is less discriminative).

# Q3 — Multiclass Gaussian Bayes Classifier

Q3

Using MLE estimates from Q2 + empirical priors, Bayes' theorem gives:

$$\pi_{i|k}^{\wedge} = \pi_k \cdot N(x_i^{\wedge}(c) \mid \mu_k^{\wedge}, \Sigma_k^{\wedge}) / \sum_j \pi_j \cdot N(x_i^{\wedge}(c) \mid \mu_j^{\wedge}, \Sigma_j^{\wedge})$$

Priors:  $\hat{\pi}_k = N_k / \sum_j N_j$  | Log-space computation (logsumexp) prevents underflow

**TRAIN**

**77.5%**

Sanity check

**VAL**

**75.2%**

Generalisation

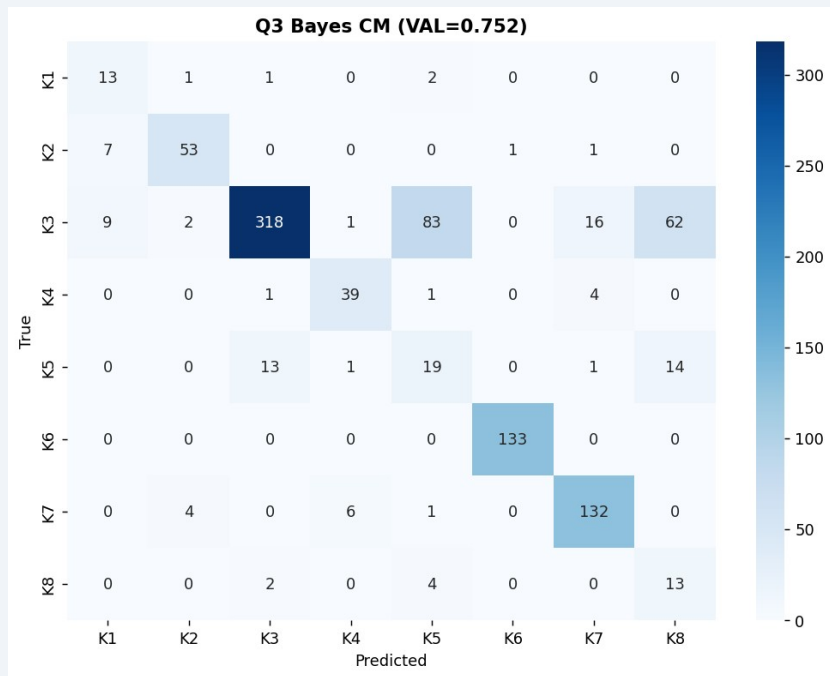
**Limitation**

Uses only 8 continuous features.  
Categorical features are ignored.

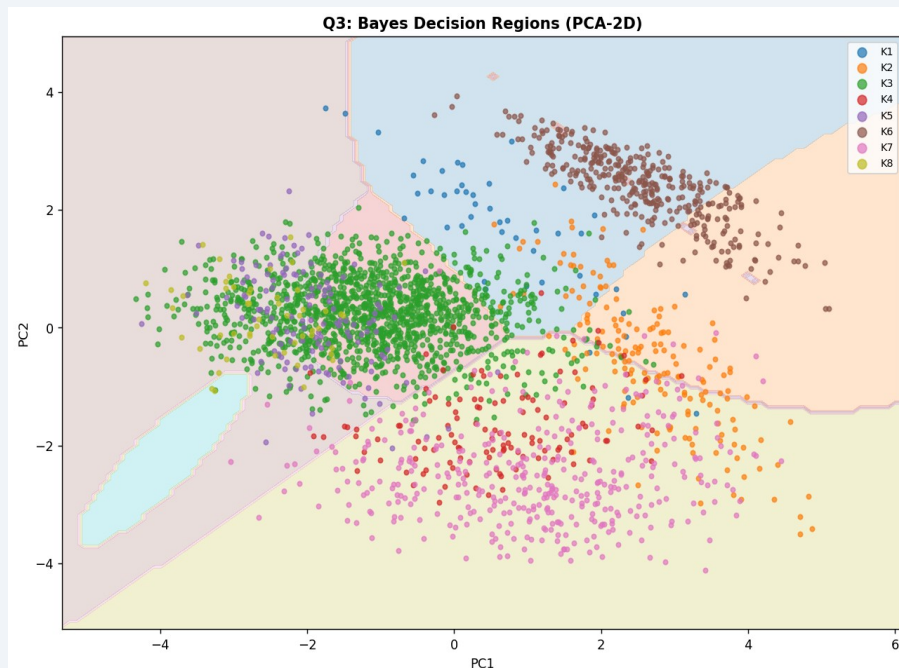
*K3 dominates predictions due to its 51% prior. Minority killers (K1, K8) suffer most.*

# Q3 — Bayes: Confusion Matrix & Decision Region (PCA-2D)

Q3



Confusion Matrix (VAL)



Decision Regions (PC1 vs PC2)

Curved ellipsoidal boundaries emerge naturally from Gaussian assumptions; K3 region (blue) dominates the central space.

# Q4 — Linear Classifier: Categorical Features

## Unlock +19pp Jump

Multinomial Logistic Regression on FULL feature vector  $x_i \in \mathbb{R}^{25}$  (continuous + one-hot):

$$\hat{c}_i = \operatorname{argmax}_k (W \cdot x_i + b)_k \quad W \in \mathbb{R}^{8 \times 25}, \quad b \in \mathbb{R}^8$$

$\ell_2$  regularisation  $C=1.0$  | Solver: L-BFGS | Features standardised (zero mean, unit variance)

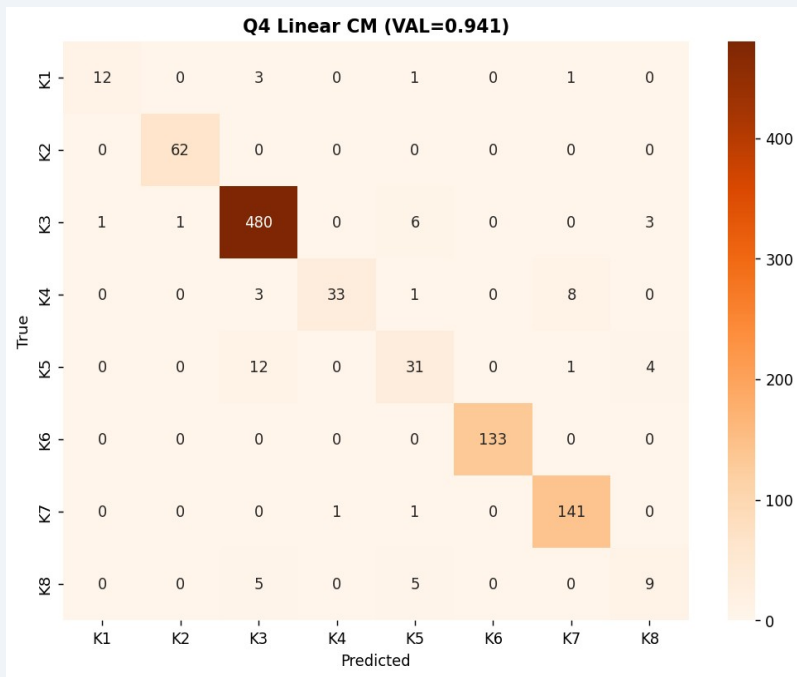
Model	Features	TRAIN	VAL	vs Bayes
Gaussian Bayes	8 continuous	77.5%	75.2%	—
Linear (LR)	25 (full)	95.9%	94.1%	+18.9pp

*Key insight: weapon\_code and scene\_type are highly killer-specific — linear model leverages these directly.*

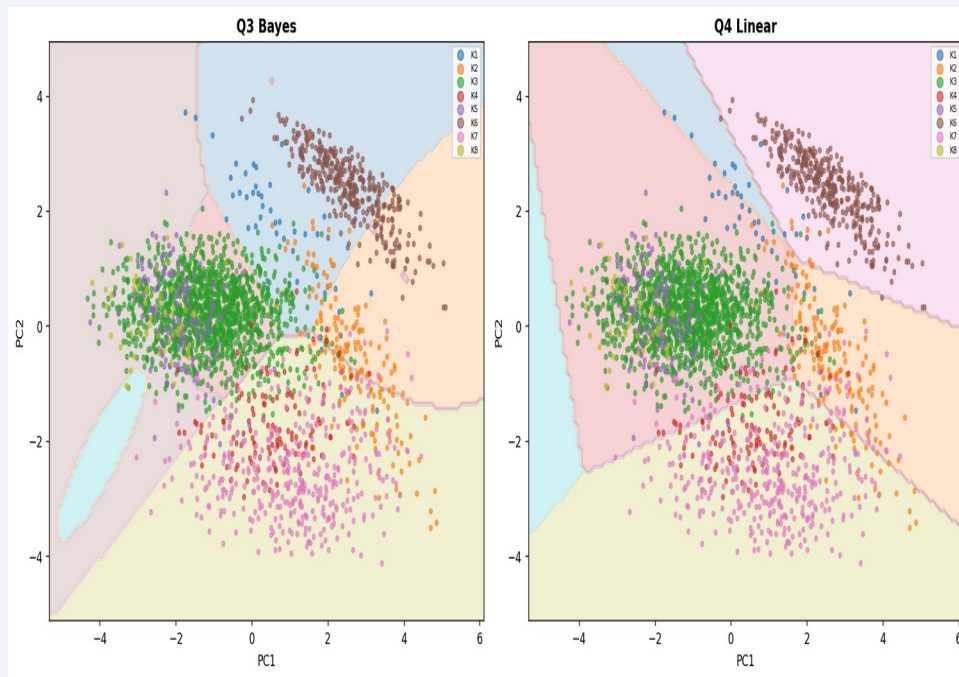
The +18.9pp gain over Bayes comes entirely from adding the 17 categorical one-hot dimensions.

# Q4 — Linear Classifier: VAL=94.1%, Bayes vs. Linear Regions

Q4



Confusion Matrix (VAL)



Decision Regions: Bayes (curved) vs. Linear (polygonal)

K1 and K8 (minority classes) now correctly classified in most cases — categorical features are decisive.

# Q5 — Support Vector Machines with RBF Kernel

Q5

One-vs-Rest multiclass SVM with RBF kernel  $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$

**C**

**10.0**

Soft-margin penalty  
(tuned on VAL)

**$\gamma$**

**scale**

$\gamma = 1/(d \cdot \text{Var}(X))$   
auto heuristic

**Kernel**

**RBF**

Non-linear  
Gaussian kernel

**Strategy**

**0vR**

8 binary SVMs  
for 8 killers

**TRAIN**

**98.8%**

Near-perfect fit  
(mild overfit)

**VAL**

**93.4%**

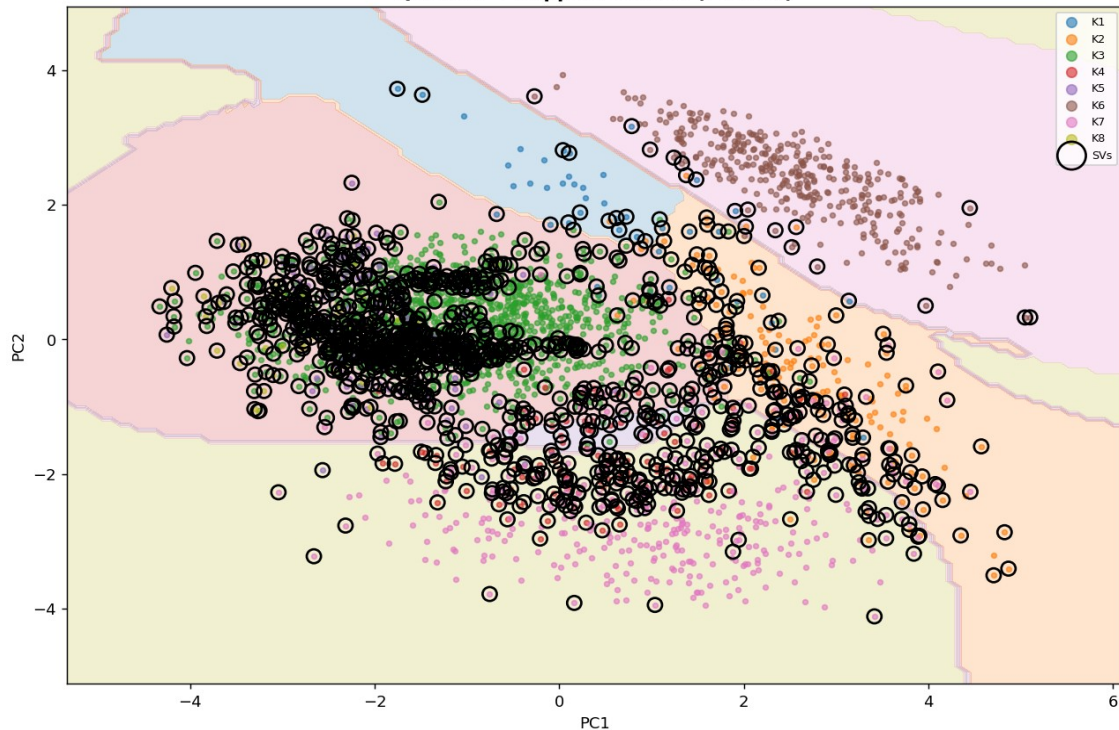
Comparable  
to Q4 Linear

*SVM achieves  $\approx$  same VAL as MLP but fits more tightly on TRAIN — the RBF kernel may be overfitting slightly.*

# Q5 — SVM Decision Regions & Support Vectors in PCA-2D

Q5

Q5: SVM + Support Vectors (PCA-2D)



## Non-linear boundaries

The RBF kernel creates curved, locally adaptive decision surfaces — impossible with a linear model.

## Support vectors (o)

Hollow circles mark the SVs: the critical incidents lying closest to the decision boundary.

## Boundary density

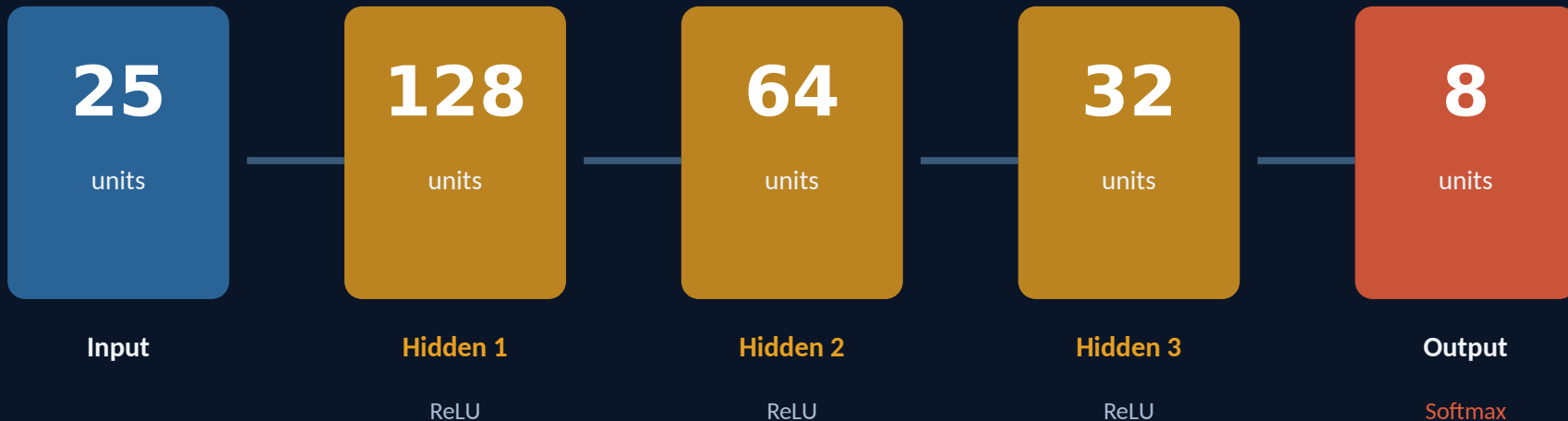
SVs concentrate most densely at K3/K6 and K3/K7 borders — the hardest classification zones.

Both the confusion matrix and VAL=93.4% confirm SVM is competitive with — but doesn't surpass — Logistic Regression.



# Q6 — Multi-Layer Perceptron: Architecture & Training

Q6



- Loss: Categorical Cross-Entropy
- Optimizer: Adam ( $\text{lr} = 10^{-3}$ )
- Weight Decay:  $\alpha = 10^{-3}$
- Early Stopping on 10% internal VAL

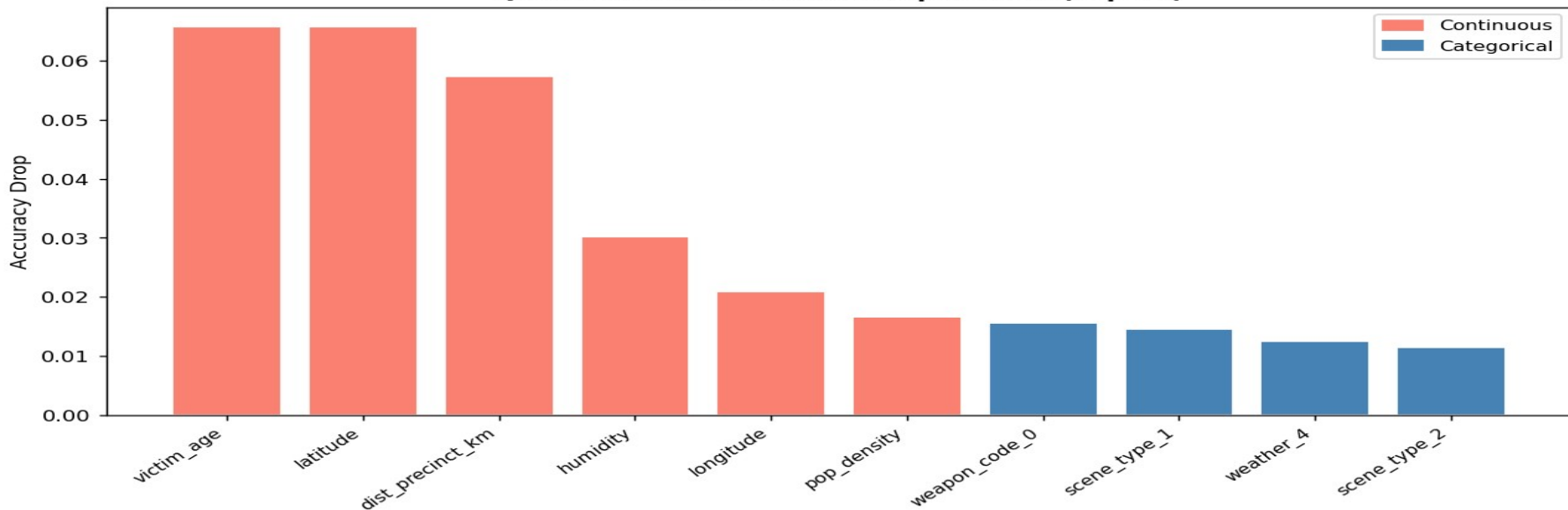
**TRAIN**  
**95.6%**

**VAL**  
**93.4%**

# Q6 — Permutation Feature Importance: What Identifies a Killer?

Q6

Q6: Permutation Feature Importance (Top 10)



## #1 victim\_age

Most powerful single feature — each killer targets a specific age profile

## #2 latitude

Geographic territory — strong killer-specific spatial clustering

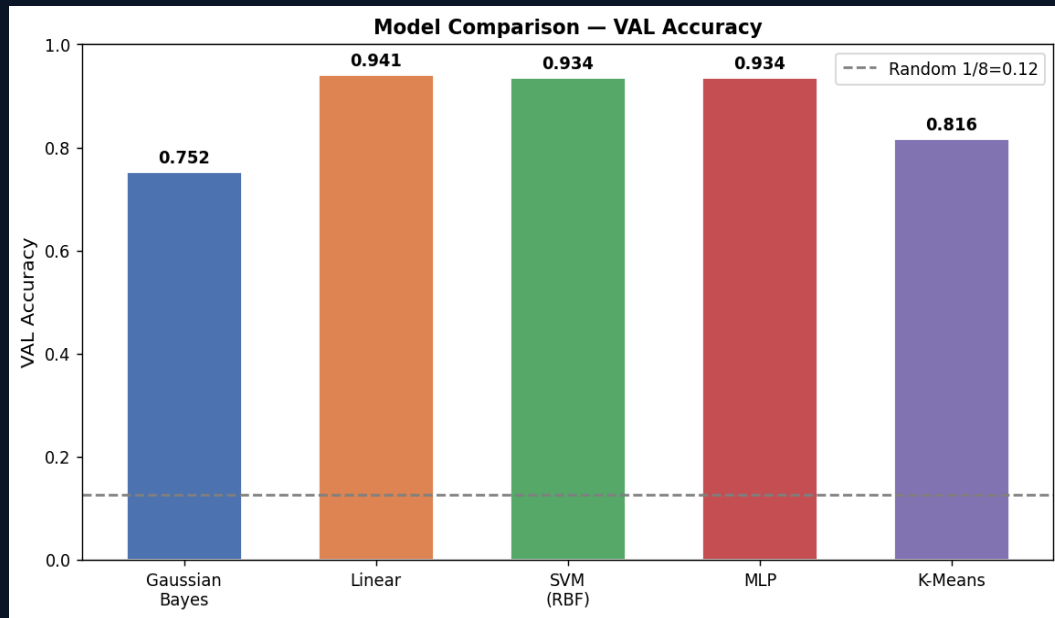
## #3 dist\_precinct\_km

Operational risk tolerance — some killers avoid precincts, others don't

*All top-5 features are continuous. Categorical features matter collectively but less individually.*

# Supervised Model Comparison — VAL Accuracy

Summary



## Gaussian Bayes

75.2%

Continuous only.  
Gaussian assumption  
holds partially.

## Linear (LR)

94.1%

Best model.  
Categorical features  
decisive.

## SVM RBF

93.4%

Competitive.  
Mild overfit  
on TRAIN.

## MLP 3-layer

93.4%

Same as SVM.  
Early stopping  
prevents overfit.

Linear decision boundaries suffice — the feature space, once encoded, is approximately linearly separable.

# Q7 — Principal Component Analysis: Finding the Latent Structure

Q7

Goal: project the 25-dimensional feature space onto a lower-dimensional 'modus operandi' space that preserves maximum variance and reveals killer cluster structure.

PCA: eigen-decompose  $\hat{\Sigma} = (1/N) \tilde{X}^T \tilde{X} = V \Lambda V^T \rightarrow z_i = V_m^T \tilde{x}_i \in \mathbb{R}^m$

Fitted on TRAIN only. Applied identically to VAL and TEST (no data leakage).

25

Input dimensions  
(standardised)

14

Components for  
90% variance

75.6%

Variance in  
first 10 PCs

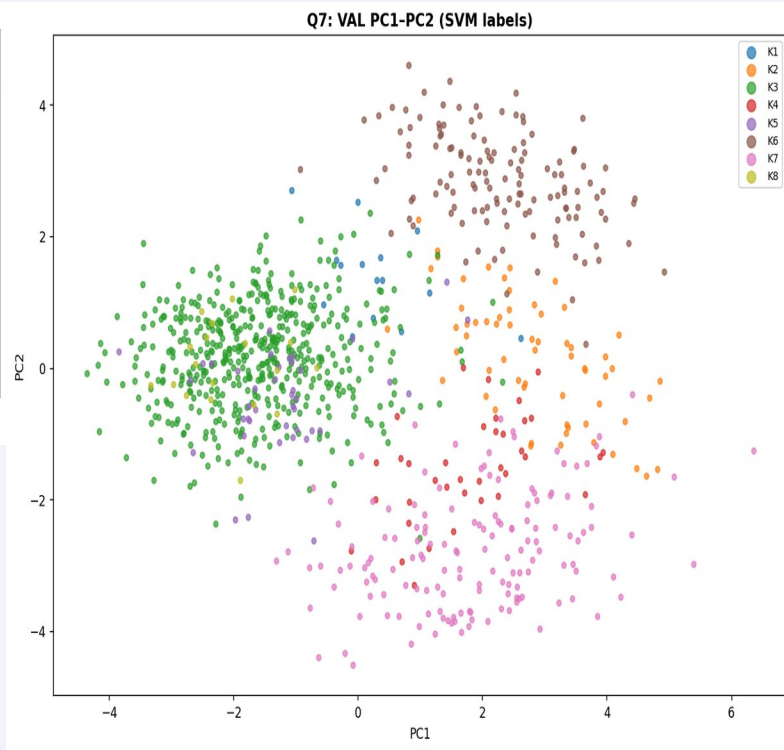
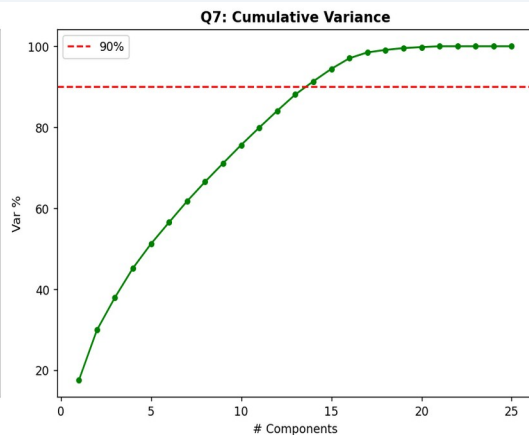
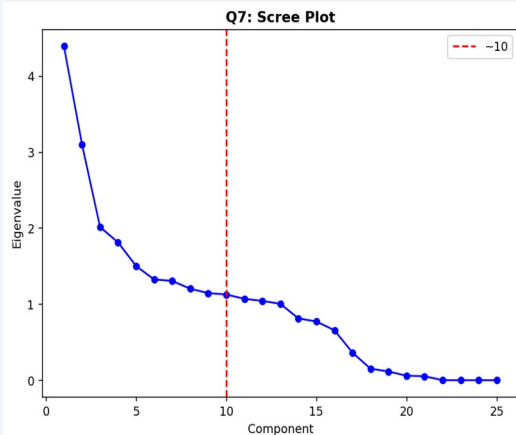
m=14

Selected latent  
dimension for Q8

*The gradual elbow in the scree plot at ~10 components suggests no sharp dimensionality boundary — variance is spread across features.*

# Q7 — Scree Plot & VAL Killer Clusters in PCA-2D Space

Q7



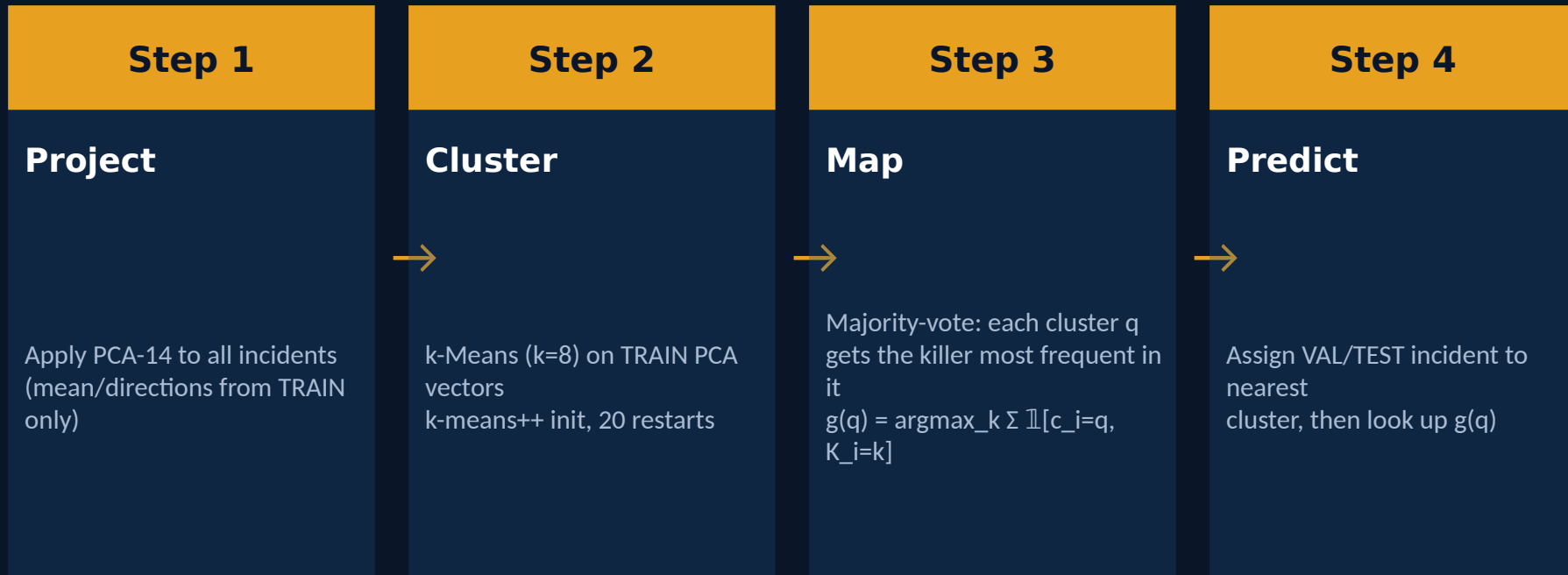
Scree plot (left): elbow at ~10, 90% threshold at m=14 components selected.

- ~8 distinct clusters visible in just 2 PCs
- K3 occupies the largest central region
- Minority killers (K1, K8) form tight clusters

Scatter = VAL | Color = SVM predicted killer

# Q8 — k-Means Clustering: Unsupervised Killer Attribution

Q8

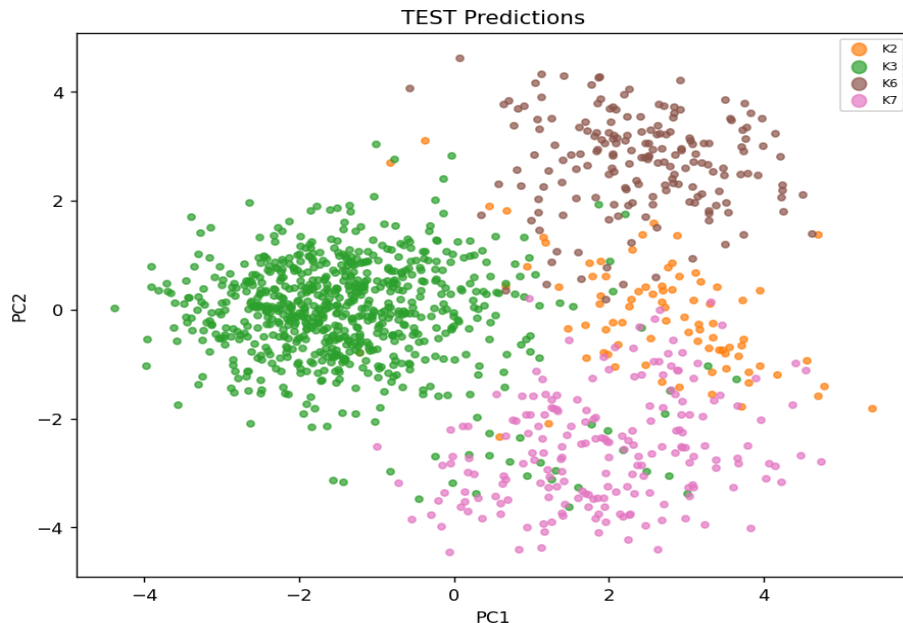
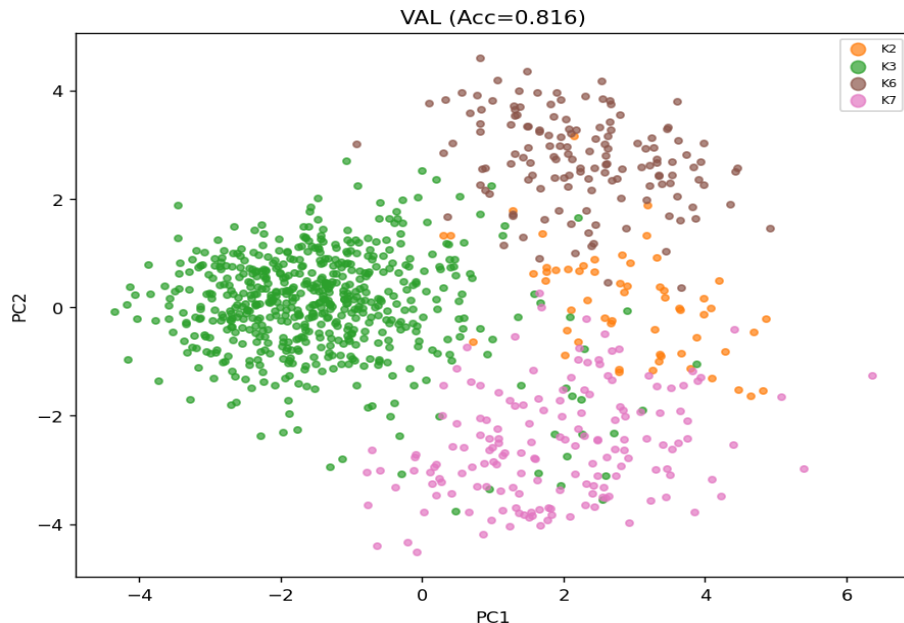


*No killer labels are used at test time — only at the mapping step (TRAIN only). Purely unsupervised inference.*

# Q8 — K-Means: 81.6% VAL Accuracy Without Test Time Labels!

Q8

Q8: K-Means Clustering



81.6  
%

VAL accuracy — unsupervised k-means beats the fully-supervised Bayes (75.2%)

Limitation: Multiple clusters mapped to K3 (the dominant class) — minority killers K1, K4, K5, K8 share a cluster or get absorbed.

# Key Forensic Insights — What the Data Tells Us

## Killers have fixed spatial territories

Latitude and longitude are the #2 most important features. Each killer operates in a geographically coherent zone visible even in unsupervised clustering.

## Victim age is the strongest single predictor

The top permutation-importance feature — each killer consistently targets victims in a specific age range (young adults vs. elderly), likely reflecting opportunistic targeting patterns.

## Categorical features collectively decisive

The +18.9pp jump from Bayes → Linear Regression demonstrates that weapon, scene type, and weather conditions together define a killer's modus operandi.

## Feature space is approximately linearly separable

The best model is the simplest (Logistic Regression). Non-linear models (SVM, MLP) offer no improvement, suggesting class boundaries are linear in the 25D feature space.

*The 81.6% unsupervised k-means result confirms: killer identity is geometrically encoded in the feature space.*



# Conclusions & Deliverables

---

Best Model

**Logistic Regression**

VAL = 94.1%

Best Unsupervised

**k-Means (PCA-14)**

VAL = 81.6%

Top Feature

**victim\_age**

Perm. imp. #1

Submission

**submission.csv**

4,800 predictions

---

The Piraeus Vice dataset is geometrically structured: killer identity is linearly separable in the 25-dimensional feature space.

Deliverables: solution\_Q1\_Q8.py · piraeus\_vice\_report.pdf · submission.csv · this slide deck