

Bayes Classifier with Gaussians, DINO-Features and PCA

- K Classes • For each image, we can encode it into a $\underline{x} \in \mathbb{R}^d$ feature representation of dimensionality d.
- We have some train data $\{\underline{x}_i, y_i\}_{i=1}^N$ for N samples and each one has label $y_i \in \{0, 1, \dots, (K-1)\}$
- So, our data matrix $\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_N \end{bmatrix}_{N \times d}$
- And, our ground truth "vector": $\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$
- Since classes are mostly defined by names e.g. "airplane", "car", and our $y_i \in \mathbb{R}$, we need a mapping, simply:
 $\underline{\text{label_map}} = \{\text{"airplane"}: 0, \text{"audf": 1}, \dots\}$
- For our matrix operations, we will use the one-hot-encoding for the labels \underline{y} , e.g. if $K=4$ and $\underline{y} = [2, 0, 3, 2]$ our one hot encoding matrix is:
- In general: $\underline{Y} \in \mathbb{R}^{N \times K}$ $\underline{Y} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}_{4 \times 4}$ Class 1 Class 2 Class 3 Class 4
Sample 1 Sample 2 Sample 3 Sample 4
- In our code, in one-hot-encode, given \underline{y} and K (num classes) we create a zeros matrix first: $\underline{Y} = [1 \ 0 \ 3 \ 2]$
 - 1) $\underline{Y} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$
 - 2) $\underline{Y} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$
 - 2) $\underline{Y} = \begin{bmatrix} [0, 1, 2, 3], [1, 0, 3, 2] \end{bmatrix} = I \rightarrow \underline{Y} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$

"rows"
 $\text{np.arange}(N)$

• For Bayes Classifier, we are interested in computing the Posterior:

$$P(x=k|x=x_i) = \frac{P(x=x_i|y=k) \cdot P(y=k)}{P(x=x_i)}$$

Given Sample data
x is prob. it is class K
→ Posterior

Prior
our class K prob. of observing x_i

- At the end of the day, for a data x_i , in order to assign it to a class K , we just compute:

$$k_0 = \arg \max_{k \in C} P(x=k|x=x_i)$$

$$= \arg \max_{k \in C} \frac{P(x=x_i|y=k) P(y=k)}{P(x_i)}$$

Since this does not depend on k we can drop it.

$$k_0 = \arg \max_{k \in C} P(x=x_i|y=k) P(y=k)$$

- For the prior of class k_0 , we can estimate it empirically just by counting of all the dataset samples N , how many belong to class K .

$$P(y=k) = \frac{N_k}{N}$$

- In our code, we can easily do this by summing over the rows of our one-hot-encoded y and dividing by N .

$$\begin{aligned} \text{Prior} &= \frac{1}{4} [1 \ 1 \ 1 \ 1] = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \xrightarrow{\text{Sum}} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \\ &= [0.25 \ 0.25 \ 0.25 \ 0.25] \underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}}_{Y} \end{aligned}$$

- The Conditional prob. $P(X=x_i | Y=k)$: thus we can set to a Gaussian for which we need to determine its Mean \underline{N}_k and Covariance $\underline{\Sigma}_k^1$. Given our dataset with labels we can estimate these are:

$$\boxed{\begin{aligned}\underline{N}_k &= \frac{1}{N_k} \sum_{\substack{i=1 \\ Y_i=k \\ N_k}}^N x_i \\ \underline{\Sigma}_k^1 &= \frac{1}{N_k-1} \sum_{\substack{i=1 \\ Y_i=k \\ N_k}}^N (x_i - \underline{N}_k) (x_i - \underline{N}_k)^T\end{aligned}}$$

In Matrix Form:

Mean: $\underline{N}_k = \frac{1}{N_k} \begin{bmatrix} \text{---} & x_1 & \text{---} \\ \text{---} & x_2 & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & x_N & \text{---} \end{bmatrix} = \frac{1}{N_k} \mathbf{y} \quad [\forall i, y[i] = k] \Rightarrow \text{mean}[0]$

* Mean over all rows. * Mean only over x_i which $y[i] = k$.

$\mathbf{y} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

($k=1$) Only sample 0 and 3.

Covariance:

$$\underline{\Sigma}_k^1 = \frac{1}{N_k-1} \underline{\tilde{X}}_k^T \underline{\tilde{X}}_k \quad \text{where}$$

Contains only points belonging to class k . $\underline{\tilde{X}}_k \in \mathbb{R}^{N_k \times d}$ is around \underline{N}_k centered points / data

$$\rightarrow \underline{\tilde{X}}_k = \underline{X}_k - \underline{N}_k$$

* For numerical stability since later in Gaussian Formula we need to invert $\underline{\Sigma}_k^1$ to avoid singular matrix

We add: $E \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$ to $\underline{\Sigma}_k^1$: $\hat{\underline{\Sigma}}_k^1 = \underline{\Sigma}_k^1 + E \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

- The Multi-Variate Gaussian formula reads:

$$f_K = \frac{1}{(2\pi)^{d/2} |\Sigma_K|^{1/2}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_K)^T \Sigma_K^{-1} (\underline{x} - \underline{\mu}_K)\right)$$

→ This is $p(x|y=k)$

- We remember we wanna: $\arg \max_{K \in C} p(x|y=k) p(y=k)$

→ Since $\log(x) > 0$ we can also:

$$\arg \max_{r \in C} \log p(x|y=k) p(y=k)$$

Log-odds Cond.
(Gaussian for us)

$$= \log p(x|y=k) + \log(p(y=k)) \xrightarrow{\text{Log of Prior.}}$$

$$= \log \frac{1}{(2\pi)^{d/2} |\Sigma_K|^{1/2}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_K)^T \Sigma_K^{-1} (\underline{x} - \underline{\mu}_K)\right) + \log(p(y=k))$$

$$= \log\left(\frac{1}{(2\pi)^{d/2} |\Sigma_K|^{1/2}}\right) - \frac{1}{2} (\underline{x} - \underline{\mu}_K)^T \Sigma_K^{-1} (\underline{x} - \underline{\mu}_K) + \log(p(y=k))$$

$$= \underbrace{\log\left(\frac{1}{(2\pi)^{d/2}}\right)}_{\text{Const. Term.}} \log\left(|\Sigma_K|^{-1/2}\right) - \frac{1}{2} (\underline{x} - \underline{\mu}_K)^T \Sigma_K^{-1} (\underline{x} - \underline{\mu}_K) + \log(p(y=k))$$

Same for all K,

so we can drop it.

$$f_K = -\frac{1}{2} \log(|\Sigma_K|) - \frac{1}{2} (\underline{x} - \underline{\mu}_K)^T \Sigma_K^{-1} (\underline{x} - \underline{\mu}_K) + \log(p(y=k))$$

$\arg \max_{r \in C}$

↑
determinant
of Cov. Matrix

↑
Inverse
of Cov. Matrix.

IA

- All in all, after we fit Gaussians to predict class for a new data point \underline{x}_i , we compute for every class K f_K with A formula and pick the one with max. value.

$$\underline{f} = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \end{bmatrix} \stackrel{\text{e.g.}}{=} \begin{bmatrix} 2 \\ 4 \\ 7 \\ 1 \end{bmatrix}$$

$$K_{\text{pred}} = \arg \max (\underline{f}, i) = 3$$

Extracting Features for Images:

- Assume we wanna classify images belonging to C classes e.g. "airplane", "car"...
- In older days, we might construct for N images our feature matrix as: $\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_N \end{bmatrix}$ where \underline{x}_i is the flattened RGB pixel data.
 - ↳ Issues: 1) For 1024×1024 RGB Images (3 channels) this results in $\underline{x}_i \in \mathbb{R}^{3,14 \text{ millions}}$ crazy large dimensions.
 - ↳ Cov. Matrix will be $3.14 \times 3.14 \text{ Mill.}$ Shape ----
- Nowadays, we can use powerful Vision Foundation Models such as DINO to extract Features from our images



→ DINO → $\underline{x}_i \in \mathbb{R}^{315}$ Feature
Feature Encoder

- Still, to further compress the feature dimension without losing too much information, we can apply Principal Component Analysis **PCA** which maximizes the variance after projecting the data to a lower dimensional space, so we could go from e.g. $\underline{X} \in \mathbb{R}^{315} \rightarrow \underline{X} \in \mathbb{R}^{50}$. This saves us computations and numerical stability issues when computing $\underline{\Sigma}_{\text{in}}^{-1}$ and $\det(\underline{\Sigma}_{\text{in}})$.