

KDG Code Understanding

After the network is trained,

for c in C : → set of all classes.

$$\mu_p = []$$

$$\sigma_p = []$$

$$X = X[c]$$

polytopes \leftarrow get polytopes.

get polytopes function

- get the weights & biases of the trained network.

preactivation $\rightarrow z_i = w_i a_{i-1} + b_i$

- For each layer except for the last layer \Rightarrow

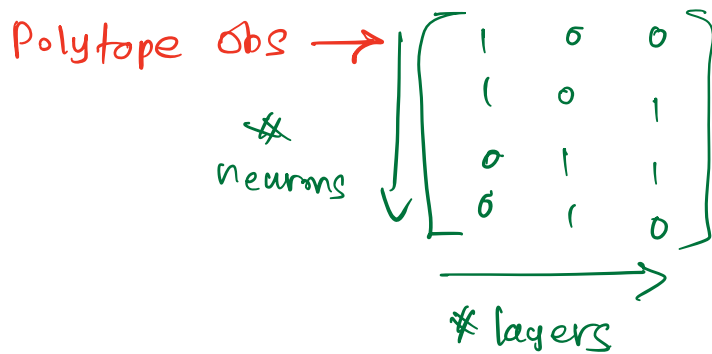
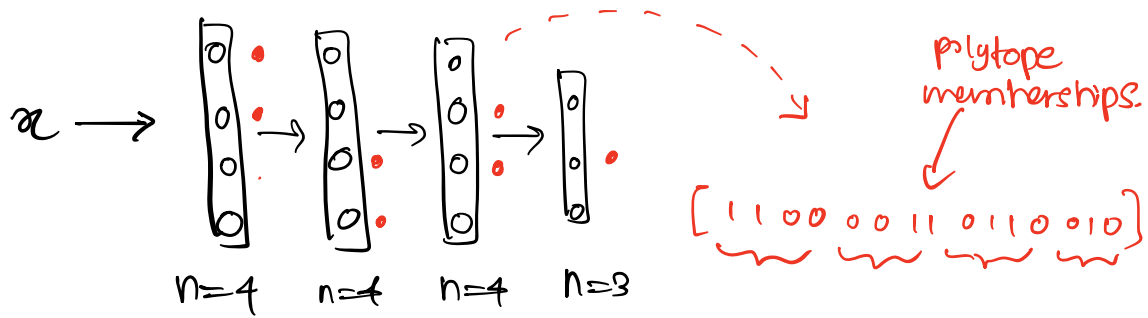
$$\mathbb{I}[z_i > 0] \quad (\text{similar to ReLU}) \quad a_i = \text{ReLU}(z_i)$$

- For the last layer.

thresholding a probability given by sigmoid/softmax layer. $\rightarrow \mathbb{I}[z_i \geq 0.5]$

$$\begin{bmatrix} 1 \\ 2 \\ 0.5 \\ 1.2 \end{bmatrix} \xrightarrow{\mathbb{I}} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

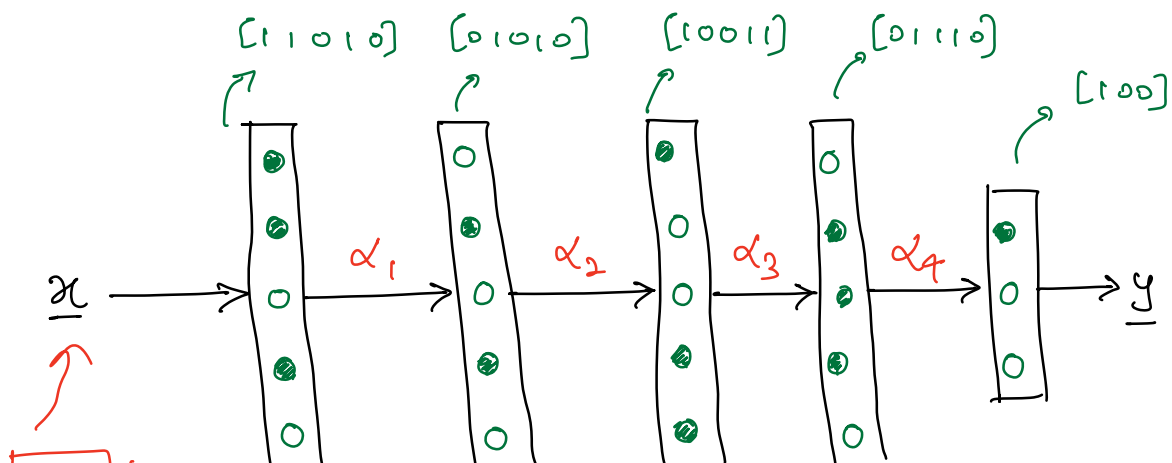
$$\begin{bmatrix} 0.8 \\ 0.9 \\ 0.1 \\ 0.2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

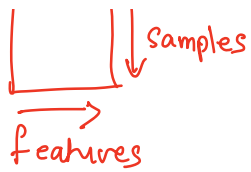


polytope memberships \rightarrow

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 2^0 \\ 2^1 \\ 2^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 5 \\ 6 \\ 2 \end{bmatrix}$$

"Get_Polytopes" Function \Rightarrow





$$a_i = \sigma(W_i a_{i-1} + b_i)$$

$$a_i^T = \sigma(a_{i-1}^T W_i^T + b_i^T)$$

$$\alpha_i = \sigma(\alpha_{i-1} W_i + \beta_i)$$

$$\text{polytope memberships} = \{ [11010], [01010], \dots, [100] \}$$

$$\text{polytope observations} = \{ 11010 | 01010 | \dots | 100 \}$$

$$\text{polytope memberships} \} = [11010 \ 01010 \ \dots \ 100] \begin{bmatrix} 2^0 \\ 2^1 \\ 2^2 \\ \vdots \\ 2^{n-1} \end{bmatrix}$$

$n \text{ digits}$

↑ returned by the get polytopes function for each data sample.

* If two data samples have similar neural network activations, they belong to the same polytope.

Gaussian Mixture Modelling

* A probabilistic model that assumes all data points are generated from a finite number of gaussian distributions with unknown params.

$$p(\phi) = \sum_{i=1}^K \phi_i \mathcal{N}(\mu_i, \Sigma_i)$$

"fit" function \Rightarrow

- * takes in the dataset $(X, y) \in \mathbb{R}^{n \times d} \times y^K$
where $K = \{1, \dots, K\}$
- * Fits a neural network to (X, y)
- * For each class label k ,
 - $P_k \leftarrow \text{_get_polytopes}(X_k)$
 - \swarrow samples with label k
 - \searrow an array with length (X_k)
contains the polytope ID of each data sample of X_k
 - for each polytope ID \Rightarrow
 - obtain the indices of all the other data samples of X_k that has the same polytope ID
 - (Hence obtain the data samples of X_k that belongs to the same polytope)
 - \nwarrow neglect the polytopes with just 1 sample.
 - Using the data samples of class k

that belongs to the polytope, fit a gaussian mixture model & obtain its means/covariance matrices (params).

"_compute_pdf" function \Rightarrow

* given a class label and a polytope ID, compute the PDF of the respective polytope group using the gaussian mixture parameters of that polytope.

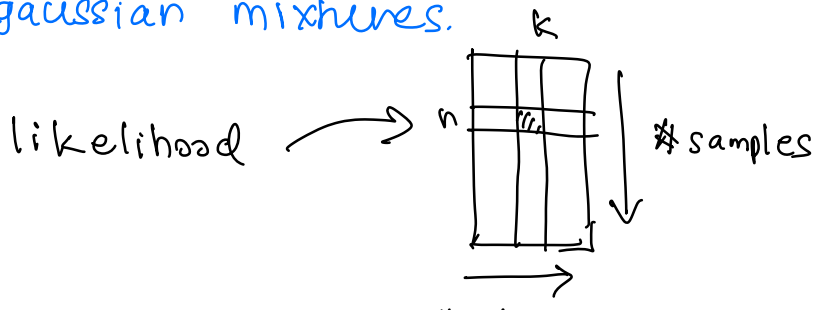
multivariate gaussian distribution

* likelihood = the PDF evaluated at a given data sample point x

\nearrow
according to m_{code}

"predict_proba" function \Rightarrow

* given the test data, obtain the posterior distribution using the stored polytopic gaussian mixtures.



* classes.

$$\left. \begin{array}{l} \text{likelihood of } k\text{th class} \\ \text{at the } n\text{th sample} \end{array} \right\} = \sum_{i=1}^{|P_k|} G(x; \mu_i, \Sigma_i)$$

$$\left. \begin{array}{l} \text{probability of } k\text{th class} \\ \text{at the } n\text{th sample} \\ P(k|x) \end{array} \right\} = \frac{\sum_{i=1}^{|P_k|} G(x; \mu_i, \Sigma_i)}{\sum_{k=1}^K \sum_{i=1}^{|P_k|} G(x; \mu_i, \Sigma_i)} //$$

∴ No need for the test data to pass through the trained neural network.

"prediction" function \Rightarrow

$$\hat{k} = \underset{k}{\operatorname{argmax}} P(k|x) \rightarrow \text{predicted class.}$$