

1 Semi-probabilistic model

Probabilistic framework: Let the data be generated by a joint distribution of two random variables: a feature vector $\mathbf{x} \sim X$ and a class label $y \sim Y$ from a given parametric family (θ is the parameter), with each observation being independent:

$$\mathbf{x}, y \sim f_{X,Y}(\mathbf{x}, y | \theta) \sim (X \times Y)^\ell, \quad \mathbf{x} \in \mathbb{R}^k, y = 1..N \quad (1)$$

Given a training sample, we need to estimate the parameter values and build a predictive model:

$$\hat{y}(\mathbf{x}) = a_\theta(\mathbf{x}) \quad (2)$$

Model: We will use the formalism of **semi-probabilistic models**, *i.e.*, we will consider X values as fixed and Y values as variables in the model. In other words, as a predictive model, we will build a probability distribution $f_Y(y | \mathbf{x}^*, \theta)$ for the random variable Y , parameterized by the data $\mathbf{x}^* \in X^\ell$ and the distribution shape parameters θ

Likelihood function: Let's write the joint distribution for all observed points as a product of independent distributions:

$$\begin{aligned} L(\theta) &= \prod_{(\mathbf{x}^*, y^*) \in (X, Y)^\ell} \underbrace{f(\mathbf{x} = \mathbf{x}^*, y = y^* | \theta)}_{\text{likelihood}} \\ &= \prod_{(\mathbf{x}^*, y^*) \in (X, Y)^\ell} \underbrace{\frac{f_Y(y = y^* | \mathbf{x} = \mathbf{x}^*, \theta)}{f_X(\mathbf{x} = \mathbf{x}^* | \theta)}}_{\text{predictive model } a_\theta(\mathbf{x})} \end{aligned} \quad (3)$$

$$\begin{aligned} f(\mathbf{x}, y | \theta) &= \frac{f(\mathbf{x}, y, \theta)}{f(\theta)} = \frac{f(y | \mathbf{x}, \theta) / f(\mathbf{x}, \theta)}{f(\theta)} \\ &= \frac{f(y | \mathbf{x}, \theta)}{f(\theta) \cdot f(\mathbf{x}, \theta)} = \frac{f(y | \mathbf{x}, \theta)}{f(\mathbf{x} | \theta)} \end{aligned}$$

Applying the logarithm to the likelihood function, we obtain the log-likelihood function:

$$\ell(\theta) = \sum_{(\mathbf{x}^*, y^*) \in (X, Y)^\ell} \ln \mathbb{P}(y = y^* | \mathbf{x} = \mathbf{x}^*, \theta) \rightarrow \max_{\theta} \quad (4)$$

Parameter estimation: Thus, we can find the parameter values θ and use them in the conditional distribution for prediction:

$$a_{\hat{\theta}}(\mathbf{x}') = \underbrace{f_Y(y | \mathbf{x} = \mathbf{x}', \theta = \hat{\theta})}_{\text{distribution of } Y \text{ must be chosen}} \quad (5)$$

Cross-entropy loss: From summing over specific points $y^* \in Y^\ell$, we can transition to summing over all possible values $y' = 1..N$, assuming zero probability for point y^* belonging to another class $y' \neq y^*$ and defining $0 \cdot \ln 0 \equiv 0$ by definition.

$$\begin{aligned} \ell &= \sum_{\mathbf{x} \in X^\ell} \sum_{y' \in \text{supp } Y} \mathbb{I}[y' = y^*] \cdot \ln \mathbb{P}[y = y' | \mathbf{x} = \mathbf{x}^*, \theta] \\ &\sim \sum_{\mathbf{x} \in X^\ell} \sum_{y' \in \text{supp } Y} \underbrace{\mathbb{P}[y = y' | \mathbf{x} = \mathbf{x}^*, \theta]}_{\text{smooth approximation of } \mathbb{I}} \cdot \ln \mathbb{P}[y = y' | \mathbf{x} = \mathbf{x}^*, \theta] \rightarrow \min_{\theta} \end{aligned} \quad (6)$$

This is cross-entropy loss, which can be used if the model predicts **probabilities of belonging to each class**.

The prior distribution $f_X(\mathbf{x} | \theta)$ is canceled out from the product above.

✴ The maximum likelihood method considers the prior distribution of x unknown and

unimportant (unlike in MAP), focusing solely on the conditional distribution $\mathbb{P}[y | \mathbf{x}^*, \theta]$ which serves as a model for building the algorithm $\hat{y} = a(\mathbf{x})$.

✴ Generally, when the chosen prior distribution is parameter-independent ($f_{\{X\}}(\mathbf{x} | \theta) \equiv f_{\{X\}}(\mathbf{x})$),

it naturally cancels out since it remains constant for fixed \mathbf{x}^* and does not depend on θ .