# Approximating Class Posteriors with a Model Function whose Parameters Are Estimated using Maximum Likelihood Estimation Method

Given data set $D = \{(x_1, \ell_1), \ldots, (x_N, \ell_N)\}$ and model $h_i(x; \theta)$ for $P(\ell(x) = i | x)$, use MLE to optimize (train) $\theta$.

$$\hat{\theta}_{ML} = \underset{\theta}{\arg\max}\ p(D|\theta) = \underset{\theta}{\arg\max}\ \ln p(D|\theta)$$

$$= \underset{\theta}{\arg\max}\ \ln\left\{ p\left((x_1, \ell_1), \ldots, (x_N, \ell_N) | \theta\right) \right\}$$

Assume iid samples $\hookleftarrow$

$$= \underset{\theta}{\arg\max}\ \ln \prod_{n=1}^{N} p(x_n, \ell_n | \theta)$$

$$= \underset{\theta}{\arg\max}\ \sum_{n=1}^{N} \ln p(x_n, \ell_n | \theta)$$

Bayes Rule $\hookleftarrow$

$$= \underset{\theta}{\arg\max}\ \sum_{n=1}^{N} \left( \ln p(\ell_n | x_n \theta) + \ln p(x_n | \theta) \right)$$

$x_n \perp\!\!\!\perp \theta \hookleftarrow$

$N$ constant

$$= \underset{\theta}{\arg\max}\ \sum_{n=1}^{N} \ln p(\ell_n | x_n \theta) + \cancel{\sum_{n=1}^{N} \ln p(x_n)}$$

Now introduce the model ~~[scribbled]~~ $p(\ell_n|x_n, \theta)$
$$= h_{\ell_n}(x_n, \theta)$$

$$\hat{\theta}_{ML} = \underset{\theta}{\arg\max} \sum_{n=1}^{N} \ln h_{\ell_n}(x_n, \theta)$$

Our model $h: \mathbb{R}^d \to \mathbb{R}^c$ where $x \in \mathbb{R}^d$ and $\ell \in \{1, -, c\}$ is a multi-input multi-output function



$$h_1(x, \theta) \approx P(\ell = 1 | x)$$
$$\vdots$$
$$h_c(x, \theta) \approx P(\ell = c | x)$$

$$\hat{\theta}_{ML} = \underset{\theta}{\arg\min} -\frac{1}{N} \sum_{n=1}^{N} \ln h_{\ell_n}(x_n, \theta)$$
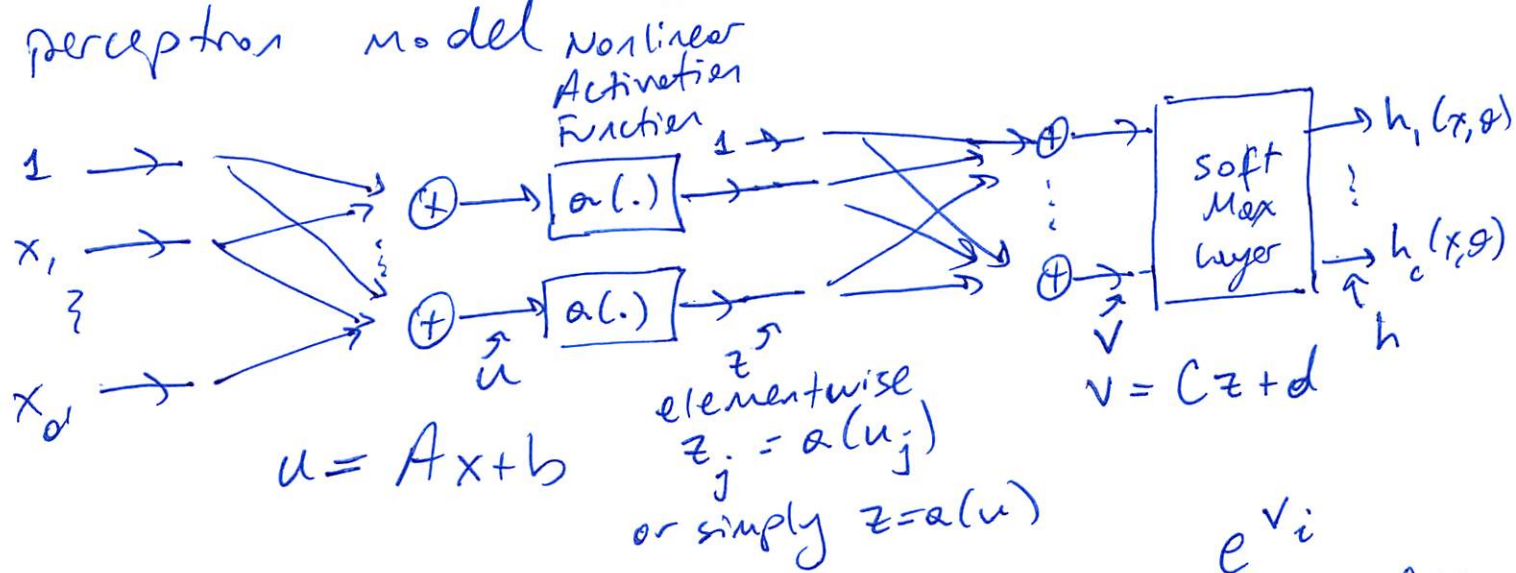
If we define $y_{(\ell)_n} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \} \ 0\text{'s} \\ \leftarrow 1 \text{ at } \ell_n^{th} \text{ entry} \\ \} \ 0\text{'s} \end{array}$

$$\hat{\theta}_{ML} = \underset{\theta}{\arg\min} -\frac{1}{N} \sum_{n=1}^{N} \left[ \sum_{\ell=1}^{c} y_{(\ell_n)_\ell} \ln h_\ell(x_n, \theta) \right]$$

Note that for the inner summation only one term is non-zero. This expression warrants the phrase minimum-cross-entropy training, which is commonly used. It allows us to also use soft labels in $y$ (e.g. if labels are uncertain).

# Simple MLP Model

Consider the following simple multilayer perceptron model

Nonlinear Activation Function



$$u = Ax+b$$

elementwise $z_j = a(u_j)$ or simply $z = a(u)$

$$v = Cz+d$$

The softmax layer produces $h_i = \dfrac{e^{v_i}}{\sum\limits_{j=1}^{c} e^{v_j}}$ for $i \in \{1, \ldots, c\}$

Let's denote it as $h = m(v)$.

The overall model function is

$$h(x, \theta) = m\Big( d + C\, a\, ( b + Ax)\Big)$$

1st layer

2nd layer

soft max layer

$\theta = $ vectorize $\{ b, A, d, C\}$   all parameters.

# Cross-Entropy Loss

for sample $(x_n, l_n)$
cross-entropy loss is this

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} \left[ - \sum_{l=1}^{c} y_{(l_n)_l} \ln h_l(x_n, \theta) \right]$$

For $y_l \neq 0$

$$= \underset{\theta}{\operatorname{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} \left[ \sum_{l=1}^{c} y_{(l_n)_l} \ln y_{(l_n)_l} - \sum_{l=1}^{c} y_{(l_n)_l} \ln h_l(x_n, \theta) \right]$$

$$= \underset{\theta}{\operatorname{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} \left[ \sum_{l=1}^{c} y_{(l_n)_l} \ln \frac{y_{(l_n)_l}}{h_l(x_n, \theta)} \right]$$

$$= \underset{\theta}{\operatorname{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} D_{KL} \left( y_{(l_n)} \| h(x_n, \theta) \right)$$

KL - divergence between
desired $y$ and model output $h$

So if $y_{(l_n)}$ is a desired posterior distribution
for $x_n$, it turns out, under the assumptions
we made, the model trained with MLE is
minimizing average KLD between desired and
model-output posterior values.

Approximating a function that maps $x$ to $y$ using Max Likelihood parameter estimation under the additive Gaussian noise model.

Given dataset $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$

where $x \in \mathbb{R}^d$, $y \in \mathbb{R}^m$; and model

$$y \approx h(x; \theta) + v \quad \text{with} \quad v \sim \mathcal{N}(0, \Sigma_v),$$

find $\hat{\theta}_{ML}$ that maximizes the likelihood of $D$.

$$\hat{\theta}_{ML} = \arg\max_{\theta} \; p(D|\theta) = \arg\max_{\theta} \; \ln p(D|\theta)$$

$$\overset{\text{iid}}{\underset{\text{samples}}{\downarrow}} \quad = \arg\max_{\theta} \; \sum_{n=1}^{N} \ln p(x_n, y_n | \theta)$$

$$\overset{\text{Bayes}}{\underset{\text{Rule}}{\downarrow}} \quad = \arg\max_{\theta} \; \sum_{n=1}^{N} \left[ \ln p(y_n | x_n, \theta) + \ln p(x_n | \theta) \right]$$

$$x_n \perp\!\!\!\perp \theta \downarrow \quad = \arg\max_{\theta} \; \sum_{n=1}^{N} \ln p(y_n | x_n, \theta) + \overset{\text{constant}}{\underset{n=1}{\sum^{N} \ln p(x_n)}}$$

$$\underset{\text{Model}}{\overset{\text{substitute}}{\downarrow}} \quad = \arg\max_{\theta} \; \sum_{n=1}^{N} \ln \left[ (2\pi)^{-m/2} |\Sigma_v|^{-1/2} e^{-\frac{1}{2}(y_n - h(x_n, \theta))^T } \right.$$
$$\left. -\Sigma_v^{-1}(y_n - h(x_n, \theta)) \right.$$

Note that our model $y \approx h(x, \theta) + v$ with
$v \sim \mathcal{N}(0, \Sigma_v)$ implies $y_n | x_n, \theta \sim \mathcal{N}(h(x_n, \theta), \Sigma_v)$
and we used this in the previous step.

Now simplify:

$$\hat{\theta}_{ML} = \arg\max_\theta \; \underbrace{N \ln (2\pi)^{-M/2}}_{constant} + \underbrace{N \ln |\Sigma_v|^{-1/2}}_{constant}$$

$$- \frac{1}{2} \sum_{n=1}^{N} (y_n - h(x_n, \theta))^T \Sigma_v^{-1} (y_n - h(x_n, \theta))$$

$\times \frac{-2}{N} \Downarrow$

$$= \arg\min_\theta \; \frac{1}{N} \sum_{n=1}^{N} (y_n - h(x_n, \theta))^T \Sigma_v^{-1} (y_n - h(x_n, \theta))$$

Now let's assume $\Sigma_v = \sigma_v^2 I$.

$$\hat{\theta}_{ML} = \arg\min_\theta \; \frac{1}{N\sigma_v^2} \sum_{n=1}^{N} (y_n - h(x_n, \theta))^T (y_n - h(x_n, \theta))$$

$$= \arg\min_\theta \; \frac{1}{N} \sum_{n=1}^{N} (y_n - h(x_n, \theta))^T (y_n - h(x_n, \theta))$$

$$= \arg\min_\theta \; \frac{1}{N} \sum_{n=1}^{N} \underbrace{\| y_n - h(x_n, \theta) \|_2^2}_{}$$

Average (Mean) Squared Error

## Simple MLP Model

$$h(x, \theta) = d + C a(b + Ax)$$

1st layer with nonlinear activation $a(\cdot)$

Linear 2nd layer

The number of perceptrons in the 1st layer can be adjusted using model selection procedures (e.g. cross-validation). For the activation function soft versions of ReLu can be used (e.g. ELu).

If the values of $y$ are known to be bounded, we can use a nonlinear activation function for the second layer, in order to make sure $h(\cdot)$ produces values in the appropriate range.

**Fact.** If $P_{XY}(x, y)$ is the joint pdf of $X$ and $Y$, then $\hat{Y}_{MSE}(x) = E_{XY}\{Y \mid X=x\}$, the conditional expectation of $Y$ given $X=x$ is the minimum-MSE estimator of $Y$ from $X$.

So our MLE-model parameters yield a function that approximates this conditional expectation.