# Neural Networks (Part II)

## Training:

### Setup: Given $\mathcal{D} := \{(x_1, y_1), \ldots, (x_m, y_m)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^s$

### Likelihood:

$$y_i \overset{i.i.d}{\sim} \mathcal{N}\left(y_i \mid \underbrace{f_\theta(x_i)}, \sigma^2\right) \iff y_i = f_\theta(x_i) + \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2$$

Forward Pass
$$\begin{cases} f_\theta(x_i) = H^{(L)} W^{(0)} + b^{(0)} & \text{(regression output layer)} \\ H^{(L)} = \sigma\left(H^{(L-1)} W^{(L)} + b^{(L)}\right) & \text{Model parameters to be traine} \\ \vdots & \Theta := \{\underline{\vartheta}, \sigma^2\} \\ H^{(1)} = \sigma\left(X W^{(1)} + b^{(1)}\right) & \Theta := \left\{ (W^{(0)}, W^{(L)}, \ldots, W^{(1)}), \right. \\ & \quad (b^{(0)}, b^{(L)}, \ldots, b^{(1)}), \\ & \quad \left. \underline{\sigma^2} \right\} \end{cases}$$

### Goal: Identify a set of "optimal" parameters $\Theta^*$, such that:

$$\Theta^* = \arg\max_\Theta p(y \mid x, \Theta) = \arg\min_\Theta -\log p(y \mid x, \Theta)$$

$$\underline{L(\Theta)} := -\log p(y \mid x, \Theta) = \frac{1}{2\underline{\sigma}^2} \sum_{i=1}^n \left[ f_\theta(x_i) - y_i \right]^2 + \frac{n}{2} \log(2\pi$$

↳ not <u>convex</u>

because $f_\theta$ is highly
non-linear

⊛ <u>Remark</u>: Typically we are only interested in estimating the neural network parameters $\vartheta$. For that, it suffices to minimize the sum of square errors:

(regression): $L(\vartheta) := \frac{1}{2} \sum_{i=1}^n \left[ f_\vartheta(x_i) - y_i \right]^2$  (sum of square errors)

(binary classification): $L(\vartheta) := -\sum_{i=1}^n y_i \log f_\vartheta(x_i) + (1 - y_i) \log (1 - f_\vartheta$

$$\left(\begin{array}{l}\text{multi-class}\\ \quad\text{classification}\end{array}\right) \quad L(\theta) := -\sum_{i=1}^{n}\sum_{c=1}^{d} y_{ic}\log f_{\theta}(x_i)$$

We will "train" our model

via gradient descent : $\quad \theta_{n+1} = \theta_n - \eta \nabla_\theta L(\theta_n)$

To do so, we need to compute $\nabla_\theta L(\theta_n)$ ! ... we will do so

via <u>back-propagation</u> !

## <u>Overfitting and Regularization</u>

1.) $\mathbb{L}_2$ parameter regularization / weight decay :

<u>Idea</u> : Modify the loss function:

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ regularization weight $(10^{-2} - 10^{-5})$

$$L(\theta) = \underbrace{\frac{1}{2}\sum_{i=1}^{n}\left[y_i - f_\theta(x_i)\right]^2}_{\text{MLE}} + \frac{\lambda}{2}\theta^T\theta \quad , \quad \|\theta\|_2^2 := \theta^T\theta$$

$\qquad\qquad\qquad$ corresponds to assuming

$\qquad\qquad\qquad\qquad\qquad \theta \sim \mathcal{N}(0, I)$

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\text{MAP}}$ drive the weights and bias

$\qquad\qquad\qquad\qquad\qquad$ closer to the origin.

2.) $\mathbb{L}_1$ parameter regularization :

$\qquad\qquad\qquad\qquad\qquad\qquad$ need to be tuned $(10^{-2} - 10^{-6})$ $\qquad$ prior

$\qquad\qquad\qquad\qquad\qquad\qquad\quad$ corresponds to assuming a Laplace

$$L(\theta) = \frac{1}{2}\sum_{i=1}^{n}\left[y_i - f_\theta(x_i)\right]^2 + \alpha\|\theta\|_1 \quad , \quad \|\theta\|_1 := \sum_i |\theta_i|$$

The $\mathbb{L}_1$ penalty promotes sparsity in $\theta$, i.e. discards features that

$\qquad$ are not needed.

⊕ In practice, we often use a combination of $\mathbb{L}_1/\mathbb{L}_2$ regularization.

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad H^{(\ell-1)} \quad H^{(\ell)}$

3.) Dropout regularization :

$\qquad$ Recall : $\qquad H^{(\ell)} = \sigma\left(H^{(\ell-1)} W^{(\ell)} + b^{(\ell)}\right)$

$\qquad$ With dropout, this becomes :

$$r_j^{(\ell)} \sim \text{Bernoulli}(p)$$

$$Z^{(\ell)} = r^{(\ell)} \odot H^{(\ell-1)} \quad \left(\text{element-wise multiplication will}\right.$$
$$\text{discard some features in } H^{(\ell-}$$

$$H^{(\ell)} = f\left(Z^{(\ell)} W^{(\ell)} + b^{(\ell)}\right)$$