# Softmax

## Suriya Chaudary

## 18 October 2024

My opinion on how and why softmax materialize.
Consider the transformation

$$\mathbf{p} = \underset{\mathbf{w} \in \mathbf{S}}{\arg\min} \langle \mathbf{w}, \mathbf{x} \rangle \tag{1}$$

where $\mathbf{x} \in \mathbf{R^d}, \mathbf{w} \in \mathbf{S} \subset \mathbf{R^d}$ such that $\mathbf{S} = \left\{ \mathbf{w} \mid \sum_i^d \mathbf{w}_i = 1 \right\}$, $\mathbf{p}_i \in [0,1]^d$ and $\sum_{i=1}^d \mathbf{p}_i = 1$.
$\hat{\mathbf{p}}$ with minimum entropy regularizer

$$\hat{\mathbf{p}} = \underset{\mathbf{w} \in \mathbf{S}}{\arg\min} \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{i=1}^d \mathbf{w}_i \log\left(\mathbf{w}_i\right) \tag{2}$$

Since $\mathbf{w} \in \mathbf{S}$, add a Lagrange multiplier $\lambda\left(\langle \mathbf{w}, \mathbf{1} \rangle - 1\right)$ to the objective function.

$$\hat{\mathbf{p}} = \underset{\mathbf{w} \in \mathbf{S}}{\arg\min} \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{i=1}^d \mathbf{w}_i \log\left(\mathbf{w}_i\right) + \lambda\left(\langle \mathbf{w}, \mathbf{1} \rangle - 1\right) \tag{3}$$

$$= \underset{\mathbf{w} \in \mathbf{S}}{\arg\min} \sum_{i=1}^d \mathbf{w}_i \mathbf{x}_i + \sum_{i=1}^d \mathbf{w}_i \log\left(\mathbf{w}_i\right) + \lambda\left(\sum_{i=1}^d \mathbf{w}_i - 1\right) \tag{4}$$

$$\tag{5}$$

Differentiate the objective function with respect to $\mathbf{w}_i$ and equate to 0

$$\mathbf{x}_i + 1 + \log\left(\mathbf{w}_i\right) + \lambda = 0 \tag{6}$$

$$\mathbf{w}_i^\star = \exp\left(-\mathbf{x}_i\right) \exp\left(-1 - \lambda\right) \tag{7}$$

$$= \frac{\exp\left(-\mathbf{x}_i\right)}{\exp\left(1 + \lambda\right)} \tag{8}$$

Set $\lambda$ such that $\sum_i^d \mathbf{w}_i^\star = 1$

$$\mathbf{w}_i^\star = \frac{\exp\left(-\mathbf{x}_i\right)}{\sum_{i=1}^d \exp\left(-\mathbf{x}_i\right)} \tag{9}$$

$$\hat{\mathbf{p}} = \mathbf{w}^\star \tag{10}$$

## Reference

Luca Trevison. The "Follow-the-Regularized-Leader" algorithm. Topics in computer science and optimization (Fall 2019).