

Draft Handout 5: Artificial neural networks

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the Artificial neural network as a model and procedure in classical and Bayesian framework. Motivation, set-up, description, computation, implementation, tricks. We focus on the Feedforward network.

Reading list & references:

- (1) Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 20 Neural Networks
- (2) Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.
 - Ch. 5 Neural Networks
- (3) Bishop, C. M. (1995). Neural networks for pattern recognition. Oxford university press.
 - Ch. 4 The multi-layer perceptron
- (4) LeCun, Y., Bottou, L., Orr, G. B., & Müller, K. R. (2002). Efficient backprop. In Neural networks: Tricks of the trade (pp. 9-50). Berlin, Heidelberg: Springer Berlin Heidelberg.

1. INTRO AND MOTIVATION ¹

Note 1. Artificial Neural Networks (NN) are statistical models which have mostly been developed from the algorithmic perspective of machine learning. They were originally created as an attempt to model the act of thinking by modeling neurons in a brain. In ML, NN are used as global approximators.

Note 2. The original biological motivation for feed-forward NN stems from McCulloch & Pitts (1943) who published a seminal model of a NN as a binary thresholding device in discrete time, i.e.

$$n_j(t) = 1 \left(\sum_{\forall i \rightarrow j} w_{j,i} n_i(t-1) > \theta_i \right)$$

where the sum is over neuron i connected to neuron j ; $n_j(t)$ is the output of neuron i at time t and $0 < w_{j,i} < 1$ are attenuation weights. Thus the effect is to threshold a weighted sum of the inputs at value θ_i . Perhaps, such a mathematical model involving compositions of interconnected non-linear functions could be able to mimic human's learning mechanism and be implemented in a computing environment (with faster computational abilities) with purpose to discover patterns, make predictions, cluster, classify, etc...

¹In this Section, formulas not needed to be memorized.

Remark 3. Mathematically, NN are rooted in the classical theorem by Kolmogorov stating (informally) that every continuous function $h(\cdot)$ on $[0, 1]^d$ can be written as

$$(1.1) \quad h(x) = \sum_{i=1}^{2d+1} F_i \left(\sum_{j=1}^d G_{i,j}(x_j) \right)$$

where $\{G_{i,j}\}$ and $\{F_i\}$ are continuous functions whose form depends on f . Perhaps, one may speculate that functions $\{G_{i,j}\}$ and $\{F_i\}$ can be approximated by sigmoids or threshold functions of the form $\sigma(w^\top x)$ allowing the number of the tunable coefficients w to be high enough such that they can represent any function -hence the property of NN as global approximators.

Example 4. Consider a regression problem with predictive rule $h : \mathbb{R}^d \rightarrow \mathbb{R}^q$, suitable for cases where the examples (data) consist of input $x \in \mathbb{R}^d$, and output targets $y \in \mathbb{R}^q$. A 2 layer artificial neural network is

$$h_k(x) = \sigma_{(2)} \left(w_{(2),k,0} + \sum_{\forall j} w_{(2),k,j} \sigma_{(1)} \left(w_{(1),j,0} + \sum_{\forall i} w_{(1),j,i} x_i \right) \right)$$

for $k = 1, \dots, q$. One may choose with $\sigma_2(\alpha) = \alpha$, $\sigma_1(\alpha) = 1/(1 + \exp(-\alpha))$. The only left here is to learn unknown parameters $\{w_{(\cdot),\cdot,\cdot}\}$.

2. FEEDFORWARD NEURAL NETWORK (MATHEMATICAL SET-UP)

Note 5. An (Artificial) Neural Network (NN) is an interconnection of several sigmoid or threshold functions associated arranged in layers, such that the outputs of one layer form the input of the next layer. Formally the structure of these interconnections can be depicted as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ whose nodes \mathcal{V} correspond to vertices/neurons and edges \mathcal{E} correspond to links between them.

Note 6. A Feed-Forward Neural Network (FFNN) or else multi-layer perceptron is a special case of NN whose vertices can be numbered so that all connections go from a neuron (vertex) to one with a higher number. Hence, neurons (vertices) have one-way connections to other neurons such that the output of a lower numbered neuron feeds the input of the higher numbered neuron (in a forward manner). The neurons can be arranged in layers so that connections go from one layer to a later layer. FFNN can be depicted by a directed acyclic graph², $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. (See Figure 2.1).

²(its vertices can be numbered so that all connections go from one vertex to another with a higher number)

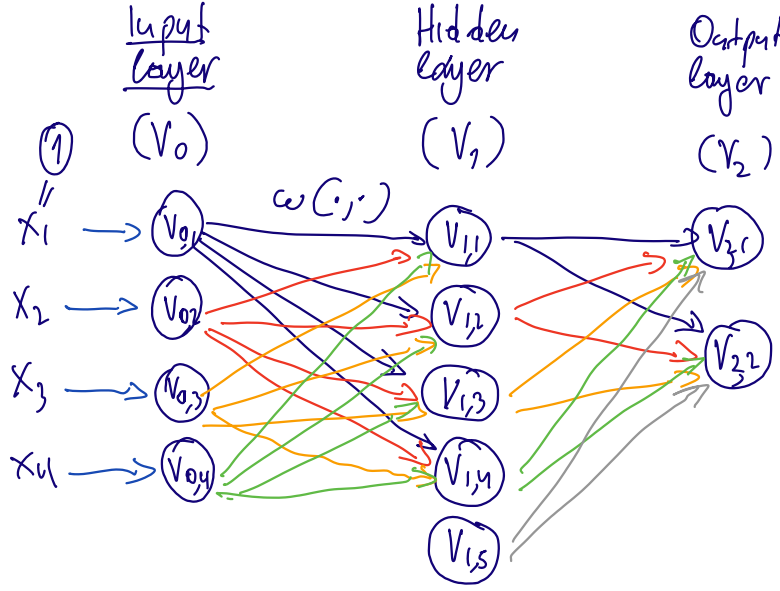


FIGURE 2.1. Feed forward neural network (1 hidden layer)

Note 7. We assume that the network is organized in layers. The set of nodes is decomposed into a union of (nonempty) disjoint subsets

$$V = \cup_{t=0}^T V_t$$

such that every edge in \mathcal{E} connects a node from V_t to a node from V_{t+1} , for $t = 1, \dots, T$.

Note 8. The first layer V_0 is called **input layer**. If x has d dimensions, then the first layer V_0 contains d nodes.

Note 9. The last layer V_T is called **output layer**. If y has q dimensions, then the last layer V_T contains q nodes.

Note 10. The intermediate layers $\{V_1, \dots, V_{T-1}\}$ are called **hidden layers**.

Note 11. In a neural network, the nodes of the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ correspond to neurons.

Notation 12. The **i -th neuron of the t -th layer** is denoted as $v_{t,i}$.

Note 13. The output of neuron i in the input layer V_0 is simply x_i that is $o_{0,i}(x) = x_i$ for $i = \{1, \dots, d\}$.

Note 14. Each edge in the graph $(v_{t,j}, v_{t+1,i})$ links the output of some neuron $v_{t,j}$ to the input of another neuron $v_{t+1,i}$; i.e. $(v_{t,j}, v_{t+1,i}) \in \mathcal{E}$.

Note 15. We define a weight function $w : \mathcal{E} \rightarrow \mathbb{R}$ over the edges \mathcal{E} .

Note 16. Activation of neuron i at hidden layer 1 is the weighted sum of the outputs $o_{0,i}(x) = x_i$ of the neurons in V_0 which are connected to $v_{1,i}$ where weighting is according to function w , that is

$$(2.1) \quad \alpha_{1,i}(x) = \sum_{\forall j:(v_{0,j}, v_{1,i}) \in \mathcal{E}} w((v_{0,j}, v_{1,i})) x_i$$

Note 17. Each single neuron $v_{t,i}$ is modeled as a simple scalar function, $\sigma_t(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$, called **activation function** at layer t .

Note 18. We denote by $o_{t,i}(x) := \sigma_{t-1}(\alpha_{t-1,i}(x))$ **the output of neuron** $v_{t,i}$ when the network is fed with the input x .

Remark 19. Activation of neuron i at layer t is the weighted sum of the outputs $o_{t-1,j}(x)$ of the neurons in V_{t-1} which are connected to $v_{t,i}$ where weighting is according to function w , that is

$$(2.2) \quad \alpha_{t,i}(x) = \sum_{\forall j:(v_{t-1,j}, v_{t,i}) \in \mathcal{E}} w((v_{t-1,j}, v_{t,i})) o_{t-1,j}(x)$$

Note 20. The input of a neuron is obtained by taking a weighted sum of the outputs of all the neurons connected to it, where the weighting is according to w .

Note 21. The input to $v_{t+1,i}$ is activation $\alpha_{t+1,i}(x)$ namely a weighted sum of the outputs $o_{t,j}(x)$ of the neurons in V_t which are connected to $v_{t+1,i}$, where weighting is according to w . The output of $v_{t+1,i}$ is the application of the activation function $\sigma_{t+1}(\cdot)$ on its input $\alpha_{t+1,i}(x)$. –That's why it is called **Feed forward Neural Network**.

Summary 22. The feed-forward NN formula in a layer by layer manner is performed (defined) according to the following recursion.

At $t = 0$,: for $i = 1, \dots, |V_0|$

$$o_{0,i}(x) := x_i$$

At $t = 0, \dots, T-1$,: for $i = 1, \dots, |V_{t+1}|$

$$\begin{aligned} \alpha_{t+1,i}(x) &= \sum_{\forall j:(v_{t,j}, v_{t+1,i}) \in \mathcal{E}} w((v_{t,j}, v_{t+1,i})) o_{t,j}(x) \\ o_{t+1,i}(x) &= \sigma_{t+1}(\alpha_{t+1,i}(x)) \end{aligned}$$

Note 23. Depth of the NN is the number of the layers T .

Note 24. Size of the network is the number $|V|$.

Note 25. Width of the NN is the number $\max_{\forall t}(|V_t|)$.

Note 26. The architecture of the neural network is defined by the triplet $(\mathcal{V}, \mathcal{E}, \sigma_t)$.

Note 27. The neural network can be fully specified by the quadruplet $(\mathcal{V}, \mathcal{E}, \sigma_t, w)$.

Example 28. Figure 2.1 denotes a NN with depth 2, size 11, width 5. The neuro with no incoming edges has $o_{1,5} = \sigma(0)$.

Notation 29. To easy the notation, we denote the weights as $w_{(t+1),i,j} := w((v_{t,j}, v_{t+1,i}))$. Using this notation, $w_{(t+1),i,j} = 0$ is equivalent in (2.2) to $(v_{t,j}, v_{t+1,i}) \notin \mathcal{E}$ and means that the link $v_{t,j} \rightarrow v_{t+1,i}$ is not in the network.

Note 30. Often a **constant neuron** $v_{t,0}$ (at each layer t and $i = 0$) which outputs 1; i.e. $o_{0,0}(x) = 1$ and $o_{t,0}(x) = 1$. The corresponding weight $w_{(t),k,0}$ is called **bias**. This resembles to the constant term in the linear regression.

Example 31. (Cont. Example 4) The 2 layer neural network

$$(2.3) \quad h_k(x) = \sigma_{(2)} \left(w_{(2),k,0} + \sum_{\forall j} w_{(2),k,j} \sigma_{(1)} \left(w_{(1),j,0} + \sum_{\forall i} w_{(1),j,i} x_i \right) \right)$$

can be written according to the recursion in Summary 22 as

- Input layer

$$o_{(0),i}(x) = \begin{cases} 1 & i = 0 \\ x_i & i = 1, \dots, d \end{cases}$$

- Hidden layer

$$\begin{aligned} \alpha_{(1),j}(x) &= w_{(1),j,0} + \sum_{\forall i} w_{(1),j,i} x_i \\ o_{(1),j}(x) &= \sigma_{(1)}(\alpha_{(1),j}(x)) \\ &= \sigma_{(1)} \left(w_{(1),j,0} + \sum_{\forall i} w_{(1),j,i} x_i \right) \end{aligned}$$

- Output layer

$$\begin{aligned} \alpha_{(2),k}(x) &= w_{(2),k,0} + \sum_{\forall j} w_{(2),k,j} o_{(2),k}(x) \\ o_{(2),k}(x) &= \sigma_{(2)}(\alpha_{(2),k}(x)) \\ &= \sigma_{(2)} \left(w_{(2),k,0} + \sum_{\forall j} w_{(2),k,j} o_{(2),k}(x) \right) \end{aligned}$$

If $h_k(x)$ returns values in \mathbb{R} , we can choose $\sigma_{(2)}$ as the identity function i.e., $\sigma_{(2)}(\alpha) = \alpha$. Note that (2.3) is also presented more compact

$$h_k(x) = \sigma_{(2)} \left(\sum_{\forall j} w_{(2),k,j} \sigma_{(1)} \left(\sum_{\forall i} w_{(1),j,i} x_i \right) \right)$$

by considering the constant neuron associated with first x which is set equal to 1, e.g $x_1 = 1$, similar to the linear regression models.

Note 32. Activation functions σ_t are non-increasing functions often sigmoids or threshold functions. Their choice is problem-dependent. some examples

- Identity function: $\sigma(\alpha) = \alpha$ cannot be used in hidden layers
- Threshold sigmoid: $\sigma(\alpha) = 1 (\alpha > 0)$
- Logistic sigmoid: $\sigma(\alpha) = \frac{1}{1+\exp(-\alpha)}$
- Rectified linear unit: $\text{RELU}(\alpha) = \max(\alpha, 0)$

Example 33. Examples on the choice of the activation function at the output layer T :

- In the univariate regression problem with prediction rule $h(\cdot) \in \mathbb{R}$, and examples $z_i = (x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, we can choose $\sigma_T(\alpha) = \alpha$ to get $h(\alpha) = \sigma_T(\alpha) = \alpha$.
- In the binary logistic regression problem with prediction rule $h(\cdot) \in [0, 1]$ and examples $z_i = (x_i, y_i) \in \mathbb{R}^d \times \{0, 1\}$, we can choose $\sigma_T(\alpha) = \frac{1}{1+\exp(-\alpha)}$ to get $h(\alpha) = \sigma_T(\alpha) = \frac{1}{1+\exp(-\alpha)} = \frac{\exp(\alpha)}{1+\exp(\alpha)}$.

3. LEARNING NEURAL NETWORKS

Note 34. Assume we are interested in a prediction rule $h_{\mathcal{V}, \mathcal{E}, \sigma, w} : \mathbb{R}^{|V_0|} \rightarrow \mathbb{R}^{|V_T|}$ which is modeled as a feed-forward Neural Network with $(\mathcal{V}, \mathcal{E}, \sigma, w)$; that is

$$h_{\mathcal{V}, \mathcal{E}, \sigma, w}(x) = o_T(x)$$

where $o_T = (o_{T,1}, \dots, o_{T,|V_T|})^\top$ is according to the summary 22.

Note 35. Learning the architecture $(\mathcal{V}, \mathcal{E}, \sigma)$ of a neural network is a model selection task (similar to the variable selection in linear regression).

Note 36. We assume that the architecture $(\mathcal{V}, \mathcal{E}, \sigma)$ of the neural network $(\mathcal{V}, \mathcal{E}, \sigma, w)$ is fixed (given), and that there is interest in learning the weight function $w : \mathcal{E} \rightarrow \mathbb{R}$ or equivalently in vector form the vector of weights $\{w_{(t+1),i,j}\}$ where $w_{(t+1),i,j} := w((v_{t,j}, v_{t+1,i}))$.

Note 37. The class of hypotheses is

$$\mathcal{H}_{\mathcal{V}, \mathcal{E}, \sigma} = \{h_{\mathcal{V}, \mathcal{E}, \sigma, w} : \text{for all } w : \mathcal{E} \rightarrow \mathbb{R}\}$$

for given $(\mathcal{V}, \mathcal{E}, \sigma)$.

Notation 38. To simplify notation we will use h_w instead of $h_{\mathcal{V}, \mathcal{E}, \sigma, w}$.

Note 39. Assume that there is available a training set of examples (data-set) $\mathcal{S} = \{z_i = (x_i, y_i) ; i = 1, \dots, n\}$ with $x_i \in \mathcal{X} = \mathbb{R}^{|V_0|}$ and $y_i \in \mathcal{Y}$.

Note 40. To learn the unknown w , we need to specify a loss function $\ell(w, z)$ at some value of weight vector $w \in \mathbb{R}^{|\mathcal{E}|}$ and at some example $z = (x, y)$.

Example 41. For instance, loss function can be

- (1) $\ell(w, z) = \frac{1}{2} \|h_w(x) - y\|_2^2$ based on a norm
- (2) $\ell(w, z) = -\log(f(y|w))$ based on a pdf/pmf of the sampling distribution $f(y|w)$ of the target y given the weight vector w ; eg in a regression problem $y \sim N(h_w(x), \sigma^2)$.

Definition 42. Error function is a performance measure that can be defined as

$$(3.1) \quad \text{EF}(w|\{z_i^*\}) = \sum_{i=1}^n \ell(w, z_i^*)$$

where $\mathcal{S}^* = \{z_i^* = (x_i^*, y_i^*); i = 1, \dots, n^*\}$ is a set of examples which does not necessarily need to be the training set \mathcal{S} .

Example 43. Following we present some paradigms of NN implemented in applications seen before and with respect to quantities Notes 32, 34, 40, 42

Case 1. (Regression problem) Assume we wish to predict the mapping $x \xrightarrow{h(\cdot)} y$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathbb{R}$. Assume a training data-set $\{z_i = (x_i, y_i)\}$.

- The output activation function is the identity function $\sigma_T(a) = a$. This is because $h_w(x) = o_T(x) = \sigma_T(\alpha_T(x))$ by Summary 22 and Note 34. Since h_w returns in \mathbb{R} and the output activation $\alpha_T(x)$ returns in \mathbb{R} , then mapping $\sigma_T(a) = a$ suffices.
- A suitable loss can be

$$\ell(w, z = (x, y)) = \frac{1}{2} (h_w(x) - y)^2$$

Alternatively, if I consider a statistical model as

$$y_i|x, w \sim N(\mu_i, \beta^{-1}), \text{ where } \mu_i = h_w(x_i)$$

for some fixed $\beta > 0$, the loss can be

$$\begin{aligned} \ell(w, z = (x, y)) &= -\log(N(y|h_w(x), \beta^{-1})) \\ &= -\left(-\frac{1}{2} \log\left(\frac{1}{2\pi}\right) - \frac{1}{2} \log(\beta^{-1}) - \frac{1}{2} \frac{(h_w(x) - y)^2}{\beta^{-1}}\right) \\ &= \frac{1}{2} \beta (h_w(x) - y)^2 + \text{const} \dots \end{aligned}$$

- The Error function is

$$\text{EF}(w|z) = \sum_{i=1}^n \ell(w, z_i = (x_i, y_i)) = \frac{1}{2} \beta \sum_{i=1}^n (h_w(x_i) - y_i)^2$$

Case 2. (Multi-output regression problem) Assume we wish to predict the mapping $x \xrightarrow{h(\cdot)} y$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}^q$. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathbb{R}^q$. Assume a training data-set $\{z_i = (x_i, y_i)\}$.

- The output activation function is the identity function $\sigma_T(a) = a$. This is because $h_{w,k}(x) = o_{T,k}(x) = \sigma_T(\alpha_{T,k}(x))$ for $k = 1, \dots, q$ by Summary 22 and Note 34. Since h_w returns in \mathbb{R}^q and the output activation is $\alpha_T(x)$ in \mathbb{R} , then mapping $\sigma_T(a) = a$ suffices.

- A suitable loss can be

$$\ell(w, z = (x, y)) = \frac{1}{2} \|h_w(x) - y\|_2^2 = \frac{1}{2} \sum_{k=1}^q (h_{w,k}(x) - y_k)^2$$

Alternatively, if I consider a statistical model as

$$y_i|x_i, w \sim N(\mu, \Sigma), \text{ where } \mu_i = h_w(x_i)$$

for some fixed $\Sigma > 0$, the loss can be

$$\begin{aligned} \ell(w, z = (x, y)) &= -\log(N(y|h_w(x), \Sigma)) \\ &= -\left(-\frac{q}{2} \log\left(\frac{1}{2\pi}\right) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (h_w(x) - y)^\top \Sigma^{-1} (h_w(x) - y)\right) \\ &= \frac{1}{2} (h_w(x) - y)^\top \Sigma^{-1} (h_w(x) - y) + \text{const} \dots \end{aligned}$$

- The Error function is

$$\begin{aligned} \text{EF}(w|z) &= \sum_{i=1}^n \ell(w, z_i = (x_i, y_i)) \\ &= \frac{1}{2} \sum_{i=1}^n (h_{w,k}(x_i) - y_{k,i})^\top \Sigma^{-1} (h_{w,k}(x_i) - y_{k,i}) \end{aligned}$$

where $y_{k,i}$ is the k -th dimension of the i -th example in the dataset.

Case 3. (Binary classification problem) Assume we wish to classify objects with features $x \in \mathbb{R}^d$ in 2 categories. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow (0, 1)$ as a classification probability i.e, $h_w(x) = \Pr(x \text{ belongs to class 1})$. Assume a training data-set $\{z_i = (x_i, y_i)\}$ with $y_i \in \{0, \dots, q\}$ labeling the class.

- A suitable output activation function can be the logistic sigmoid

$$\sigma_T(a) = \frac{1}{1 + \exp(-a)}$$

This is because $h_w(x) = o_T(x) = \sigma_T(\alpha_T(x))$ by Summary 22 and Note 34. Since h_w returns in $(0, 1)$ and the output activation is $\alpha_T(x)$ in \mathbb{R} , then the aforesaid mapping suffices.

- If I consider a statistical model as

$$y_i|x_i, w \sim \text{Bernoulli}(p_i), \text{ where } p_i = h_w(x_i)$$

with mass function

$$f(y|x, w) = h_w(x)^y (1 - h_w(x))^{1-y}$$

the loss can be

$$\begin{aligned} \ell(w, z = (x, y)) &= -\log(\text{Bernoulli}(y_i|h_w(x))) \\ &= -(y \log(h_w(x)) + (1 - y) (1 - \log(h_w(x)))) \end{aligned}$$

- The Error function given a set of examples $\{z_i^* = (x_i^*, y_i^*)\}$ is

$$\begin{aligned} \text{EF}(w|z^*) &= \sum_{i=1}^n \ell(w, z_i^* = (x_i^*, y_i^*)) \\ &= - \sum_{i=1}^n y_i^* \log(h_w(x_i^*)) - (1 - y_i^*) (1 - \log(h_w(x_i^*))) \end{aligned}$$

Case 4. (Multi-class classification problem) Assume we wish to classify objects with features $x \in \mathbb{R}^d$ in q categories. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathcal{P}$, with $\mathcal{P} = \left\{ \varpi \in (0, 1)^q : \sum_{j=1}^q \varpi_j = 1 \right\}$ and $h_w = (h_{w,1}, \dots, h_{w,q})^\top$, as a classification probability i.e, $h_{w,k}(x) = \Pr(x \text{ belongs to class } k)$. Assume a training data-set $\{z_i = (x_i, y_i)\}$ with $y_i \in \{0, 1\}$ labeling the class. Then it is

- A suitable output activation function can be the softmax function

$$(3.2) \quad \sigma_T(a_k) = \frac{\exp(a_k)}{\sum_{k'=1}^q \exp(a_{k'})}, \text{ for } k = 1, \dots, q$$

This is because $h_{w,k}(x) = o_{T,k}(x) = \sigma_{T,k}(\alpha_{T,k}(x))$ by Summary 22 and Note 34. Since h_w is essentially a probability vector and the output activation is $\alpha_{T,k}(x)$ in \mathbb{R} the aforesaid mapping suffices. Unfortunately, 3.2 is invariant to additive transformations $a \leftarrow a + c$ i.e, $\sigma_T(a_k) = \sigma_T(a_k + \text{const})$.

- Another suitable output activation function can be

$$(3.3) \quad \tilde{\sigma}_T(a_k) = \frac{\exp(a_k)}{1 + \sum_{k'=1}^{q-1} \exp(a_{k'})}, \text{ for } k = 1, \dots, q-1$$

This is because $h_{w,k}(x) = o_{T,k}(x) = \tilde{\sigma}_{T,k}(\alpha_{T,k}(x))$ by Summary 22 and Note 34. Since h_w is essentially a probability vector and the output activation is $\alpha_{T,k}(x)$ in \mathbb{R} the aforesaid mapping suffices as $h_{w,k}(x) = o_{T,k}(x)$ for $k = 1, \dots, q-1$ and $h_{w,q}(x) = 1 - \sum_{k=1}^{q-1} h_{w,k}(x)$. The advantage of (3.3) compared to (3.2) is that the former uses one less output neurons (less unknown parameters to learn). Also (3.3) is not invariant to additive transformations $a \leftarrow a + c$ unlike 3.2 i.e, $\tilde{\sigma}_T(a_k) \neq \tilde{\sigma}_T(a_k + \text{const})$.

- If I consider a statistical model as

$$y_i|x, w \sim \text{Multinomial}(p_i), \text{ where } p_i = h_w(x_i)$$

with mass function

$$f(y_i|x, w) = \prod_{k=1}^q h_{w,k}(x_i)^{y_{i,k}}$$

The loss can be

$$\begin{aligned} \ell(w, z = (x, y)) &= -\log(\text{Multinomial}(y|h_w(x))) \\ &= - \sum_{k=1}^q y_k \log(h_{w,k}(x)) \end{aligned}$$

- The Error function given a set of examples $\{z_i^* = (x_i^*, y_i^*)\}$ is

$$\begin{aligned} \text{EF}(w|z^*) &= \sum_{i=1}^n \ell(w, z_i^* = (x_i^*, y_i^*)) \\ &= - \sum_{i=1}^n \sum_{k=1}^q y_{k,i} \log(h_{w,k}(x_i)) \end{aligned}$$

where $y_{k,i}$ is the k -th dimension of the i -th example in the dataset.

4. CLASSICAL LEARNING OF NEURAL NETWORK

Note 44. Our purpose is to find optimal $h_w \in \mathcal{H}_{\mathcal{V}, \mathcal{E}, \sigma}$ under loss $\ell(\cdot, \cdot)$ and against training data-set $\mathcal{S} = \{z_i = (x_i, y_i); i = 1, \dots, n\}$.

4.1. Standard. Essentially, this is an optimization problem, where the objective is to minimize either the risk function $R_g(w)$ or the empirical risk function $\hat{R}_S(w)$.

Problem 45. Compute optimal $w^* \in \mathbb{R}^{|\mathcal{E}|}$ by minimizing the risk function $R_g(w)$

$$(4.1) \quad w^* = \arg \min_{w \in \mathcal{H}} (R_g(w)) = \arg \min_{w \in \mathcal{H}} (\mathbb{E}_{z \sim g} (\ell(w, z)))$$

Problem 46. Compute optimal $w^* \in \mathbb{R}^{|\mathcal{E}|}$ by minimizing the empirical risk function $\hat{R}_S(w)$

$$(4.2) \quad w^* = \arg \min_{w \in \mathcal{H}} (\hat{R}_S(w)) = \arg \min_{w \in \mathcal{H}} \left(\frac{1}{n} \sum_{i=1}^n \ell(w, z_i) \right)$$

4.2. Shrinkage.

Note 47. Often neural network models are over-parameterized, in the sense that the dimensionality of $w \in \mathbb{R}^{|\mathcal{E}|}$ is too large, with negative consequences in its predictability. This can be addressed by learning the architecture NN, precisely by learning which of the edges of the FFNN significantly contribute to the NN model and should be kept, and which do not contribute and may be removed (or be inactive).

Note 48. To address Note 47, one can resort to shrinkage methods e.g., LASSO, Ridge, which indirectly allow edge/weight selection/elimination by shrinking the values of the weights $\{w_{(t),i,j}\}$ towards zero, and setting some of them as $w_{(t),i,j} = 0$ if their absolute value is small enough. This is based on the observation in (2.2) i.e., $w_{(t),i,j} = 0$ is equivalent to $(v_{t,j}, v_{t+1,i}) \notin \mathcal{E}$ which implies that the link $v_{t,j} \rightarrow v_{t+1,i}$ is not active (essentially not in the neural network).

Problem 49. Find $h_w \in \mathcal{H}_{\mathcal{V}, \mathcal{E}, \sigma}$ (essentially compute $w \in \mathbb{R}^{|\mathcal{E}|}$) under loss $\ell(\cdot, \cdot)$ and shrinkage term $J(w; \lambda)$, and against training data-set $\mathcal{S} = \{z_i = (x_i, y_i); i = 1, \dots, n\}$. Compute $w^* \in \mathbb{R}^{|\mathcal{E}|}$, according to the minimization:

$$(4.3) \quad w^* = \arg \min_{w \in \mathcal{H}} (R_g(w) + J(w; \lambda))$$

$$(4.4) \quad = \arg \min_{w \in \mathcal{H}} (\mathbb{E}_{z \sim g} (\ell(w, z) + J(w; \lambda)))$$

and set $w_{t,i,j}^* = 0$ if $w_{t,i,j}^*$ is less than a threshold user specific value $\xi > 0$ i.e. $|w_{t,i,j}^*| < \xi$.

Problem 50. Popular shrinkage terms $J(w; \lambda)$ are

- Ridge: $J(w; \lambda) = \lambda \|w\|_1$ Term 1
- LASSO: $J(w; \lambda) = \lambda \|w\|_2^2$ Term 1
- Elastic net: $J(w; \lambda = (\lambda_1, \lambda_2)) = \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2$

5. STOCHASTIC GRADIENT IMPLEMENTATION FOR TRAINING NN

Note 51. Training a neural network model is usually a high-dimensional problem (essentially w has high dimensionality to provide a better approximation) and a big-data problem (essentially we need a large number of training examples to learn a large number of weights). For this reason, Stochastic Gradient Descent (and its variations) is a suitable stochastic learning algorithm.

Note 52. To address Problem 45, the recursion of the SGD with batch size m is

$$w^{(t+1)} = w^{(t)} - \eta_t \frac{1}{m} \sum_{j=1}^m \partial_w \ell(w^{(t)}, z_j^{(t)})$$

Note 53. To address the Problem 49, the recursion of the SGD with batch size m is

$$w^{(t+1)} = w^{(t)} - \eta_t \left[\frac{1}{m} \sum_{j=1}^m \partial_w \ell(w^{(t)}, z_j^{(t)}) + \partial_w J(w; \lambda) \right]$$

for some positive λ which is user specified, or chosen via cross validation.

Note 54. The learning problem associated to the Neural network model is (almost always) non-convex due to the non-convex loss with respect to the w 's. Upon implementing SGD in the learning problem of Neural Network, the theoretical results in Section 4 (Handout 2) and Section 4 (Handout 3) will not be effective due to the violation of the assumptions.

Note 55. Practical guidelines for the use of SGD in the learning problem of NN:

- (1) Utilize a learning rate η_t that changes over the iterations. The choice of the sequence η_t is more significant. In practice, it is tuned by a trial and error manner: given a validation data-set \mathcal{S}^* you may perform cross validation based on Error Function (3.1).
- (2) Re-run the SGD procedure several times and by using different settings (learning rate η_t , batch size m) and different seed $w^{(0)}$ (randomly choosen) each time. Possibly, by luck, at some trial, we will initialize the SGD process with a random seed $w^{(0)}$ producing a trace leading to a good local minimum w^* .
- (3) The output w_{SGD}^* returned by SGD is the best discovered w tested by using a performance measure (Error Function) using a validation set $\mathcal{S}^* = \{z_i^* = (x_i^*, y_i^*); i = 1, \dots, n^*\}$; Eg.

$$w_{\text{SGD}}^* = \arg \min_{w^{(t)}} \left(\text{EF} \left(w^{(t)} | \{z_i^*\} \right) \right).$$

5.1. Error backpropagation (to compute $\nabla_w \ell(w, z)$).

Note 56. The error backpropagation procedure is an efficient algorithm for the computation of the gradient $\nabla_w \ell(w, z)$ of the loss function $\ell(w, z)$ at some value of w and some example $z = (x, y)$ as required for the implementation of stochastic gradient based algorithm for the learning problems under consideration.

Assumption 57. *Error backpropagation assumes that $\ell(w, z)$ is differentiable at w for each value of z ; $\nabla_w \ell(w, z)$ exists.*

Notation 58. Let $V_t = \{v_{t,1}, \dots, v_{t,k_t}\}$ be the t -th layer of a NN and $k_t = |V_t|$ the number of neuron in layer $t = 1, \dots, T$.

Notation 59. Let $o_t = (o_{t,1}, \dots, o_{t,k_t})^\top$ be the vector of the outputs of the t -th layer of a NN.

Notation 60. Let $\alpha_t = (\alpha_{t,1}, \dots, \alpha_{t,k_t})^\top$ be the vector of the activations of the t -th layer of a NN.

Notation 61. We denote as $\ell_t(\cdot)$ the loss function $\ell(w, z)$ as a function of the output o_t at t -th layer.

- E.g it is $\ell_T(\xi) = \frac{1}{2}(\xi - y)^2$ if $\ell(w, z) = \frac{1}{2}(h_w(x) - y)^2$ since $h_w(x) = o_T(x)$.

Algorithm 62. (*Error backpropagation*)

Requires: The NN $(\mathcal{V}, \mathcal{E}, \sigma, w)$ with the values of the a weight vector $w \in \mathbb{R}^{|\mathcal{E}|}$, and example value $z = (x, y)$

Returns: $\nabla_w \ell(w, z) = \left(\frac{\partial}{\partial w_{t,i,j}} \ell(w, z); \forall t, i, j \right)$

Initialize:

$$\text{Set } w_{t+1,j,i} = \begin{cases} w(v_{t,j}, v_{t+1,i}) & \text{if } (v_{t,j}, v_{t+1,i}) \in \mathcal{E} \\ 0 & \text{if } (v_{t,j}, v_{t+1,i}) \notin \mathcal{E} \end{cases}$$

Forward pass:

(1) For $i = 1, \dots, d$

Set:

$$o_{0,i} = x_i$$

(2) For $t = 1, \dots, T$

For $i = 1, \dots, k_t$

Compute:

$$\alpha_{t,i} = \sum_{j=1}^{k_{t-1}} w_{t,i,j} o_{t-1,j}$$

Compute:

$$o_{t,i} = \sigma_t(\alpha_{t,i})$$

Backward pass:

(1) Set^a:

$$\delta_T = \frac{d\ell_T}{do_T}(o_T)$$

(2) For $t = T - 1, \dots, 1$

For $i = 1, \dots, k_t$

Compute:

$$\delta_{t,i} = \sum_{j=1}^{k_{t+1}} w_{t+1,j,i} \delta_{t+1,j} \sigma'(\alpha_{t+1,j})$$

Output:

for each edge $(v_{t-1,j}, v_{t,i}) \in \mathcal{E}$

Set,

$$\frac{\partial}{\partial w_{t,i,j}} \ell(w, z) = \delta_{t,i} \sigma'(\alpha_{t,i}) o_{t-1,j}$$

for $i = 1, \dots, k_t$ and $j = 1, \dots, k_{t-1}$.

^aEg if $\ell(w, z) = \frac{1}{2} \sum_{k=1}^q (h_k(x) - y_k)^2$ then $\delta_T = \frac{d\ell_T}{do_T}(o_T) = o_T - y$, or otherwise $\frac{\partial \ell_T}{\partial o_{T,k}}(o_T) = o_{T,k} - y_k$ for all $k = 1, \dots, q$

Remark 63. Several sources in the literature use the notation

$$\tilde{\delta}_t := \frac{d\ell_t}{d\alpha_t} = \frac{d\ell_t}{do_t} \frac{do_t}{d\alpha_t}$$

instead of

$$\delta_t = \frac{d\ell_t}{do_t}$$

used in Algorithm 62. Then the backward pass and output steps in Algorithm 62 can be equivalently restated as

Backward pass:

(1) For $t = T$,
 Set:

$$\tilde{\delta}_T = \frac{d\ell_T}{d\alpha_T} (o_T)$$

(2) For $t = T - 1, \dots, 1$
 For $i = 1, \dots, k_t$
 Compute:

$$\tilde{\delta}_{t,i} = \sigma_t(\alpha_{t,k}) \sum_{j=1}^{k_{t+1}} w_{t+1,j,i} \tilde{\delta}_{t+1,j}$$

Output:

- for each edge $(v_{t-1,j}, v_{t,i}) \in \mathcal{E}$
 Set

$$\frac{\partial}{\partial w_{t,i,j}} \ell(w, z) = \tilde{\delta}_{t,i} o_{t-1,j}$$

Example 64. Consider the multi-output regression problem, assume a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathbb{R}^q$ with $h_w = (h_{w,1}, \dots, h_{w,q})$ and

$$h_k(x) = \sigma_2 \left(\sum_{j=1}^c w_{2,k,j} \sigma_1 \left(\sum_{i=1}^d w_{1,j,i} x_i \right) \right)$$

with activation functions $\sigma_2(a) = a$, and $\sigma_1(a) = \tanh \left(\frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} \right)$, and loss $\ell(w, z) = \frac{1}{2} \sum_{k=1}^q (h_k(x) - y_k)^2$ for example $z = (x, y)$. Perform the error back propagation steps to compute the elements of $\nabla_w \ell(w, z)$ at some w and z .

Solution.

Forward pass:

Set: $o_{0,i} = x_i$ for $i = 1, \dots, d$

Compute:

at $t = 1$: for $j = 1, \dots, c$

comp: $\alpha_{1,j} = \sum_{i=1}^d w_{1,i,j} x_i$

comp: $o_{1,j} = \tanh(\alpha_{1,j})$

at $t = 2$: for $k = 1, \dots, q$

comp: $\alpha_{2,k} = \sum_{j=1}^c w_{2,k,j} o_{1,j}$

comp: $o_{2,k} = \alpha_{2,k}$

get: $h_k = o_{2,k}$

Note that $\frac{d}{d\xi}\sigma_1(\xi) = 1 - (\sigma_1(\xi))^2$ and that $\frac{d}{d\xi}\sigma_2(\xi) = 1$.

Backward pass:

at $t = 2$: for $k = 1, \dots, q$

comp: $\tilde{\delta}_{2,k} = \frac{\partial}{\partial \alpha_{2,k}} \ell_T = \sum_{j=1}^q \frac{\partial \ell_T}{\partial o_{2,j}} \frac{\partial o_{2,j}}{\partial \alpha_{2,k}} = \frac{\partial \ell_T}{\partial o_{2,k}} \frac{\partial o_{2,k}}{\partial \alpha_{2,k}} = h_k - x_k$

at $t = 1$: for $j = 1, \dots, c$

comp:

$$\begin{aligned} \tilde{\delta}_{1,j} &= \frac{d}{d\xi} \sigma_1(\xi) \Big|_{\xi=\alpha_{1,j}} \sum_{k=1}^q w_{2,k,j} \tilde{\delta}_{2,k} \\ &= \left(1 - (o_{1,j})^2\right) \sum_{k=1}^q w_{2,k,j} \tilde{\delta}_{2,k} \end{aligned}$$

Output:

$$\frac{\partial \ell(w, z)}{\partial w_{1,j,i}} = \tilde{\delta}_{1,j} x_i \text{ and } \frac{\partial \ell(w, z)}{\partial w_{2,k,j}} = \tilde{\delta}_{2,k} o_{2,j}$$

5.2. Design of Error back-propagation (Algorithm 62).

Notation 65. Let $\ell_t : \mathbb{R}^{k_t} \rightarrow \mathbb{R}$ be the loss function of the sub-network defined by layers $\{V_1, \dots, V_T\}$ as a function of the outputs o_t of the neurons in V_t .

Notation 66. Let W_t be the $k_{t-1} \times k_t$ matrix of the weights of the t -th layer of a NN such as $[W_t]_{i,j} := w_{t,i,j} = w(v_{t,j}, v_{t+1,i})$.

Notation 67. We introduce notation, such that if $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ and $a \in \mathbb{R}^d$, then $\sigma(a) \in \mathbb{R}^d$ is a d -dimensional vector such that $\sigma(a) = (\sigma(a_1), \dots, \sigma(a_d))$.

Notation 68. The vector of the outputs $o_t \in \mathbb{R}^{k_t}$ of the neurons of V_t can be written as $o_t = \sigma_t(\alpha_t)$; aka $o_{t,j} = \sigma_t(\alpha_{t,j})$, for $j = 1, \dots, k_t$.

Notation 69. The vector of activations, aka vector of the inputs, of the neurons of V_t can be written as $\alpha_t = W_t o_{t-1}$.

Note 70. Hence,

$$(5.1) \quad \ell_t(o_t) = \ell_t(\sigma_t(\alpha_t)) = \ell_t(\sigma_t(W_t o_{t-1}))$$

Notation 71. To facilitate differential calculus, let O_{t-1} be the $k_t \times (k_{t-1} k_t)$ matrix

$$O_{t-1} = \begin{bmatrix} o_{t-1}^\top & 0 & \dots & 0 \\ 0 & \ddots & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & o_{t-1}^\top \end{bmatrix}$$

and let $w_t \in \mathbb{R}^{k_{t-1}k_t}$ is the vector $w_t = (W_{t,1}, \dots, W_{t,k_{t-1}})^\top$ by concatenating the rows of matrix W_t . Hence

$$(5.2) \quad W_t o_{t-1} = O_{t-1} w_{t-1}$$

Note 72. Hence, from (5.1) and (5.2), it is

$$\ell_t(o_t) = \ell_t \left(\underbrace{\sigma_t \left(\overbrace{O_{t-1} w_t}^{=\alpha_t} \right)}_{=o_t} \right)$$

Note 73. By chain rule it is

$$(5.3) \quad \begin{aligned} \frac{d}{dw_t} \ell_t(o_t) &= \frac{d\ell_t}{do_t} \frac{do_t}{d\alpha_t} \frac{d\alpha_t}{dw_t} \\ &= \frac{d\ell_t}{do_t} \text{diag} \left(\left. \frac{d}{d\xi} \sigma_t(\xi) \right|_{\xi=\alpha_t} \right) O_{t-1} \end{aligned}$$

Notation 74. Let as define

$$\delta_t = \frac{d}{do_t} \ell_t(o_t)$$

and name them as ‘errors’.

Note 75. Then (5.3), becomes

$$\begin{aligned} \frac{d\ell_t}{dw_t} &= \delta_t \text{diag} \left(\left. \frac{d}{d\xi} \sigma_t(\xi) \right|_{\xi=\alpha_t} \right) O_{t-1} \\ &= \left(\delta_{t,1} \sigma'(\alpha_{t,1}) o_{t-1,1}^\top, \dots, \delta_{t,k_t} \sigma'(\alpha_{t,k_t}) o_{t-1,k_t}^\top \right) \end{aligned}$$

Note 76. I will find a way to compute $\{\delta_t\}$ recursively from V_T to V_0 . For $t = T$, it is

$$\delta_T = \frac{d}{do_T} \ell_T(o_T) = \frac{d}{do_T} \ell(o_T, z)$$

Note that

$$\ell_t(o_t) = \ell_{t+1} \left(\underbrace{\sigma_{t+1} \left(\overbrace{W_{t+1} o_t}^{=\alpha_{t+1}} \right)}_{=o_{t+1}} \right)$$

Then for $t = T-1, \dots, 1$, it is

$$\begin{aligned} \delta_t &= \frac{d\ell_t}{do_t} = \frac{d\ell_{t+1}}{do_t} \\ &= \frac{d\ell_{t+1}}{do_{t+1}} \frac{do_{t+1}}{d\alpha_{t+1}} \frac{d\alpha_{t+1}}{do_t} \\ &= \delta_{t+1} \text{diag} \left(\left. \frac{d}{d\xi} \sigma_{t+1}(\xi) \right|_{\xi=\alpha_{t+1}} \right) W_{t+1} \end{aligned}$$

Note 77. Hence, for $t = 1, \dots, T$ it is

$$\frac{d\ell(w, z)}{dw_t} = \delta_t \text{diag} \left(\left. \frac{d}{d\xi} \sigma_t(\xi) \right|_{\xi=\alpha_t} \right) O_{t-1}$$

or element wise, for $t = 1, \dots, T$, $j = 1, \dots, k_{t-1}$, $i = 1, \dots, k_t$ it is

$$\frac{\partial \ell(w, z)}{\partial w_{t,i,j}} = \delta_{t,i} \sigma'_t(\alpha_{t,i}) O_{t-1,j}$$

Remark 78. Error back-propagation idea can also be used to efficiently compute the gradient $\nabla_x \ell(w, z = (x, y))$ of the loss function $\ell(w, z = (x, y))$ with respect to the inputs x . The idea is the same, use chain rule to find a recursive procedure.

5.3. Computation of the Hessian for preconditioning.

Note 79. Recall (Handout 3, Algorithm 43), that SGD may be improved by using a suitable preconditioner $P_t > 0$ as

$$w^{(t+1)} = w^{(t)} - \eta_t P_t \nabla_w \ell(w^{(t)}, z^{(t)})$$

such a preconditioner can be the $P_t := [H_t + \epsilon I_d]^{-1}$ where H_t is the Hessian of $\ell(w^{(t)}, z^{(t)})$ and $\epsilon > 0$.

Note 80. The computation of the Hessian can be done by error propagation

$$[H_t]_{i,j} = \left. \frac{\partial^2}{\partial w_i \partial w_j} f(w) \right|_{w=w^{(t)}}$$

We will not go further to exact computations.

Note 81. The exact computation of the Hessian H_t of the loss $\ell(\cdot, z)$ in NN setting is computationally expensive and hence approximations are often used (with hope to work well). One can implement the general purpose AdaGrad (Section 7.1, Handout 3). In what follows we present other alternatives tailored to the NN model.

(...)

APPENDIX A. ABOUT GRAPHS

Definition 82. A directed graph is an ordered pair $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ comprising

- a set of nodes \mathcal{V} (or vertices), where a nodes are abstract objects, and
- a set of edges $\mathcal{E} = \{(v, u) | v \in \mathcal{V}, u \in \mathcal{V}, v \neq u\}$ (or directed edges, directed links, arrows) which are ordered pairs of vertices (that is, an edge is associated with two distinct vertices).

Definition 83. An edge-weighted graph or a network is a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ equipped with a weight function $w : \mathcal{E} \rightarrow \mathbb{R}$ that assigns a number (the weight) to each edge $e \in \mathcal{E}$.

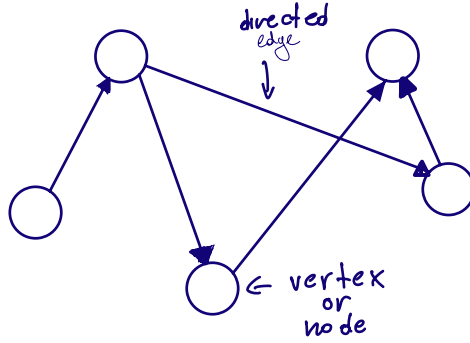


FIGURE A.1. A directed graph

APPENDIX B. ABOUT PARTIAL DERIVATIVES

Note 84. Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

Definition 85. The partial derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at the point $a = (a_1, \dots, a_n) \in U \subseteq \mathbb{R}^n$ with respect to the i -th variable is denoted as

$$\frac{\partial f}{\partial x_i}(a) \text{ or } \frac{\partial f}{\partial x_i}(x) \Big|_{x=a}$$

and defined as

$$\begin{aligned} \frac{\partial f}{\partial x_i}(x) \Big|_{x=a} &= \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_{i-1}, a_i + h, a_{i+1}, \dots, a_n) - f(a_1, \dots, a_{i-1}, a_i, a_{i+1}, \dots, a_n)}{h} \\ &= \lim_{h \rightarrow 0} \frac{f(a + he_i) - f(a)}{h} \end{aligned}$$

where e_i is a $0 - 1$ vector with only one ace in the i -th location.

Remark 86. Essentially, Definition 85 says that the partial derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at the point $a = (a_1, \dots, a_n) \in U \subseteq \mathbb{R}^n$ with respect to the i -th variable is

$$\frac{\partial f}{\partial x_i}(a) = \frac{dg}{dh}(h) \Big|_{h=0}$$

the derivative of function $g(h) := f(a_1, \dots, a_{i-1}, a_i + h, a_{i+1}, \dots, a_n)$ at value 0.

APPENDIX C. ABOUT JACOBIAN

Note 87. Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

Note 88. The Jacobian of f at $w \in \mathbb{R}^n$, is denoted as $J_w(f)$ and it is an $m \times n$ matrix with (i, j) elements such as

$$[J_w(f)]_{i,j} = \frac{d}{dw_j} f_i(w)$$

for $i = 1, \dots, m$ and $j = 1, \dots, n$.

Note 89. To align with standard notation in statistics, we will use the notation

$$\frac{df}{dw}, \text{ or } \frac{d}{dw} f(w)$$

for $J_w(f)$.

Note 90. Some properties

- Let functions $f(w) = Aw$, for matrix $A \in \mathbb{R}^{m,n}$ and vector $w \in \mathbb{R}^n$.
- Let functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^k \rightarrow \mathbb{R}^n$. The Jacobian of the composition function $f \circ g : \mathbb{R}^k \rightarrow \mathbb{R}^m$ with $f(w) = f(g(w))$ at w is

$$\left. \frac{d}{d\xi} f \circ g(\xi) \right|_{\xi=w} = \left. \frac{df}{dg}(g) \right|_{g=g(w)} \left. \frac{dg}{d\xi}(\xi) \right|_{\xi=w}$$

or more compactly

$$\frac{d}{dw} f \circ g = \frac{df}{dg} \frac{dg}{dw}$$

Example 91. Let $g : \mathbb{R}^k \rightarrow \mathbb{R}^n$ with $g(w) = Aw$ where matrix $A \in \mathbb{R}^{n,k}$. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$. Compute $\frac{d}{dw} \sigma \circ g(w)$.

Hint: Adopt notation $\sigma(a) = (\sigma(a_1), \dots, \sigma(a_n))$ when $a \in \mathbb{R}^d$.

Solution. It is

$$\begin{aligned} \frac{d}{dw} \sigma \circ g(w) &= \left. \frac{d}{dg} \sigma(g) \right|_{g=g(w)} \frac{d}{dw} g(w) \\ &= \left. \frac{d}{dg} \sigma(g) \right|_{g=Aw} \frac{d}{dw} g(w) \\ &= \text{diag}(\sigma(Aw)) A \end{aligned}$$

This is because

$$\left[\frac{d}{dw} g(w) \right]_{i,j} = \left[\frac{d}{dw} Aw \right]_{i,j} = \frac{\partial}{\partial w_j} \sum_{k=1}^k A_{i,k} w_k = A_{i,j}$$

and because

$$\left[\frac{d}{dg} \sigma(g) \right]_{i,j} = \left[\frac{d(\sigma(g_1), \dots, \sigma(g_n))}{d(g_1, \dots, g_n)} \right]_{i,j} = \frac{\partial}{\partial g_j} \sigma(g_i) = \begin{cases} \sigma(g_i), & i = j \\ 0, & i \neq j \end{cases}$$