

THE SOFTMAX FUNCTION

Given a real N -dimensional column vector \mathbf{x} , a Softmax function is defined as

$$S(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\sum_{i=1}^N e^{x_i}},$$

where the notation $e^{\mathbf{x}}$ denotes the vector $[e^{x_1} e^{x_2} \dots e^{x_N}]^T$. More compactly

$$S(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\mathbf{1}^T e^{\mathbf{x}}},$$

where $\mathbf{1}$ is a vector with all ones.

B.1 Property 1 : SHIFT INVARIANCE

$$S(\mathbf{x} + c\mathbf{1}) = S(\mathbf{x}), \quad \text{for any constant } c. \quad (34)$$

Proof: immediate because $\frac{e^{x_i+c}}{\sum_j e^{x_j+c}} = \frac{e^{x_i}}{\sum_j e^{x_j}}$.

B.2 Property 2 : THE INVERSE

Given $\mathbf{p} = S(\mathbf{x})$, the inverse function is

$$\mathbf{x} = \log \mathbf{p} + c\mathbf{1}, \quad (35)$$

where c is an undetermined constant.

This property means that the softmax function is invertible, except for a constant that cannot be determined from \mathbf{p} . Constraints on \mathbf{x} may help to obtain the constant. Note that this is also a consequence of Property 1, i.e. any constant added to the input is definitely lost.

B.3 Property 3 : SUPERPOSITION AND SCALING

$$S(\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_N \mathbf{x}_N + c\mathbf{1}) \propto S(\mathbf{x}_1)^{\alpha_1} \odot S(\mathbf{x}_2)^{\alpha_2} \odot \dots \odot S(\mathbf{x}_N)^{\alpha_N}, \quad (36)$$

where $\alpha_1, \alpha_2, \dots, \alpha_N$ are any real numbers.

Proof: From Property B.1 the constant term is irrelevant. Observing that $\frac{e^{\alpha x_i}}{\sum_{j=1}^N e^{\alpha x_j}} \propto (e^{x_i})^\alpha$ and that $\frac{e^{x_i+y_i}}{\sum_j e^{x_j+y_j}} \propto e^{x_i} e^{y_i}$, we have the result. *Very useful and very need to be noted*. Note that product in the right hand side of (36) is a distribution, but it must be normalized. More specifically, Property 2 implies that the sum becomes

$$S(\mathbf{x}_1 + \mathbf{x}_2) \propto S(\mathbf{x}_1) \odot S(\mathbf{x}_2). \quad (37)$$

Similarly a difference becomes

$$S(\mathbf{x}_1 - \mathbf{x}_2) \propto S(\mathbf{x}_1) ./ S(\mathbf{x}_2), \quad (38)$$

where $./$ denotes element-by-element division. Clearly the values that appear at the denominator must be non null. Scaling the input gives

$$(see \text{ Property 5}) \quad S(\alpha \mathbf{x}) \propto S(\mathbf{x})^\alpha, \quad (39)$$

which represent a sharper or smoother distribution with respect to $S(\mathbf{x})$ if $\alpha > 1$ or $0 < \alpha < 1$ respectively. Changing the sign to the input gives

$$S(-\mathbf{x}) \propto 1./S(\mathbf{x}), \quad (40)$$

where obviously the numbers at the denominator must be non null.

properties on the derivatives.....

PROPERTY 4 : THE JACOBIAN

The softmax is a vector function and every output depends on all inputs, because of the denominator. Therefore we must consider all the derivatives $\frac{\partial y_i}{\partial z_r} \quad \forall i, r = 1, \dots, M$, i.e. the Jacobian

$$\underline{J} = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} & \dots & \frac{\partial y_1}{\partial z_M} \\ \frac{\partial y_2}{\partial z_1} & \frac{\partial y_2}{\partial z_2} & \dots & \frac{\partial y_2}{\partial z_M} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_M}{\partial z_1} & \frac{\partial y_M}{\partial z_2} & \dots & \frac{\partial y_M}{\partial z_M} \end{bmatrix}$$

The Jacobian for the softmax function \rightarrow

$$\underline{J} = \begin{bmatrix} +y_1(1-y_1) & -y_1y_2 & \dots & -y_1y_M \\ -y_2y_1 & +y_2(1-y_2) & \dots & -y_2y_M \\ \vdots & \vdots & \ddots & \vdots \\ -y_My_1 & -y_My_2 & \dots & +y_M(1-y_M) \end{bmatrix} = \text{diag}(y) - yy^T$$

or for a generic element

$$\frac{\partial y_i}{\partial z_e} = \delta_{ie} y_i - y_i y_e$$

$(\delta_{il} \text{ is the Kronecker delta } = 0 \text{ if } i \neq l \text{ and } = 1 \text{ for } i = l)$

ASR.3

PROOF:

$$\frac{\partial y_i}{\partial z_e} = \frac{\partial}{\partial z_e} \left(\frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \right)$$

For $e = i$

$$\frac{\partial y_i}{\partial z_i} = \frac{\partial}{\partial z_i} \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} = \frac{e^{z_i} (\sum_{j=1}^n e^{z_j}) - e^{z_i} e^{z_i}}{(\sum_{j=1}^n e^{z_j})^2} = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \cdot \left(\frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \right)^2$$
$$= y_i (1 - y_i)$$

For $e \neq i$

$$\frac{\partial y_i}{\partial z_e} = \frac{\partial}{\partial z_e} \left(\frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \right) = \frac{-e^{z_i} e^{z_e}}{(\sum_{j=1}^n e^{z_j})^2} = -\frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \cdot \frac{e^{z_e}}{\sum_{j=1}^n e^{z_j}}$$
$$= -y_i y_e$$

~~MAX TO MAX BETWEEN~~

ASF.4

PROPERTY 5 Sharpness of the softmax.

From the scaling property, we have $\boxed{\text{Pf}}$

$$S(\alpha \underline{x}) \propto S(\underline{x})^{\alpha}$$

which means that we can control the sharpness of the output by changing α . More specifically for $\alpha \rightarrow \infty$ $y = S(\alpha \underline{x})$ tends to a delta on the max value.

Proof. without loss of generality, assume that x_1, x_2, \dots, x_N are ranked in descending order, so that x_1 is the maximum. The i th output can be written as

$$\frac{e^{\alpha x_i}}{\sum_{j=1}^N e^{\alpha x_j}} = \frac{e^{\alpha x_i}}{e^{\alpha x_1} (1 + e^{\alpha(x_2-x_1)} + \dots + e^{\alpha(x_N-x_1)})} = \\ = \frac{e^{\alpha(x_i-x_1)}}{1 + e^{\alpha(x_2-x_1)} + \dots + e^{\alpha(x_N-x_1)}} \quad (\star)$$

For $i \neq 1$ all the exponents $\alpha(x_i - x_1)$ go to $-\infty$

because $x_i - x_1 < 0$ and $y_i \rightarrow 0$ $\forall i \neq 1$

For $i = 1$ instead we have $\frac{y_1 = 1}{1 + e^{\alpha(x_2-x_1)} + \dots + e^{\alpha(x_N-x_1)}} \rightarrow 1$ $\boxed{\text{QED}}$

Property 5(b)

For $\alpha \rightarrow 0^+$ $S(\alpha \underline{x}) \rightarrow$ Uniform distribution.

Proof

For $\alpha \rightarrow 0^+$ all the exponents $\alpha(x_i - x_1)$ tend to zero

and

$$\frac{e^{\alpha x_i}}{\sum_{j=1}^N e^{\alpha x_j}} \rightarrow \frac{1}{N} \quad \boxed{\text{QED}}$$

A corollary of 5(a) is that if we have multiple maxima, i.e.

ASF.5

$$x_1 = x_2 = \dots = x_n > x_{n+1} > \dots > x_N$$

$$y_1, y_2, \dots, y_n \rightarrow \frac{1}{\epsilon} \quad y_{n+1}, y_N \rightarrow 0$$

Proof: Immediate from (*)

Property (5c) (dual to 5(a))

where $\alpha \rightarrow -\infty$ $S(x) \rightarrow$ Self concentrated on the maximum.

Proof: Assume that the inputs are ranked.

$$x_1 > x_2 > \dots > x_N$$

x_N is the max. Now the i th output can be written as

$$\begin{aligned} y_i &= \frac{e^{\alpha x_i}}{e^{\alpha x_N} (e^{\alpha(x_1-x_N)} + e^{\alpha(x_2-x_N)} + \dots + e^{\alpha(x_{N-1}-x_N)} + 1)} \\ &= \frac{e^{\alpha(x_i-x_N)}}{e^{\alpha(x_1-x_N)} + e^{\alpha(x_2-x_N)} + \dots + e^{\alpha(x_{N-1}-x_N)} + 1} \end{aligned}$$

All the exponentials go to zero for $i \neq N$ because $x_i - x_N > 0$ and α negative.

Therefore for $i = N \rightarrow y_i \rightarrow 1$

Observation:

The sign of α can be used to make the softmax a softmax.

ASF.6

A NUMERICALLY STABLE SOFTMAX

Since

$$S(z+c) = S(z)$$

and the outputs sum to one, the system is "overparameterized." This is because the inputs are untrained and may be arbitrarily biased to have small or large values.

A numerically more stable version (Rectified) for the softmax function is often used by fixing one of the inputs to a given value, for example $z_i = 0$. The input fixed to zero may also be different at every sample taking

$$S(z - \max z_i)$$

- In this way all exponents in

$$e^{z_j - \max z_i}$$

$$\sum_{i=1}^N e^{z_i - \max z_i}$$

are negative and therefore the exponentials are in $(0, 1]$. This prevents overflow (still prone to underflow).

- At least one of the exponents is 0, and therefore at least one of the exponentials is one. This means that at least one value does not underflow. Furthermore the denominator will be ≥ 1 preventing division by zero. We have at least one non-zero numerator and a all-zeros output caused by numerical issues.

The rectified version of the softmax rarely has an effect on learning.

ASF.7

The softmax for 2 outputs.

Note that the logistic function can be used equivalent to a two-class classifier with the softmax at the output.

$$X \rightarrow \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad y_1 = p(a_1|x) = \frac{e^{z_1}}{e^{z_1} + e^{z_2}} = \frac{1}{1 + e^{z_2 - z_1}}$$

$$y_2 = p(a_2|x) = \frac{e^{z_2}}{e^{z_1} + e^{z_2}} = \frac{1}{1 + e^{z_1 - z_2}}$$

$$1 - y_1 = p(a_0|x)$$