

[Follow](#)

592K Followers · Editors' Picks · Features · Deep Dives · Grow · Contribute · About

Derivative of the Softmax Function and the Categorical Cross-Entropy Loss

A simple and quick derivation

Thomas Kurbiel Apr 23 · 6 min read

...

In this short post, we are going to compute the Jacobian matrix of the softmax function. By applying an elegant computational trick, we will make the derivation super short. Using the obtained Jacobian matrix, we will then compute the gradient of the categorical cross-entropy loss.

Softmax Function

The main purpose of the softmax function is to grab a vector of arbitrary real numbers and turn it into probabilities:

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} \left\{ \begin{array}{c} 1.1 \rightarrow \\ 2.2 \rightarrow \\ 0.2 \rightarrow \\ -1.7 \rightarrow \end{array} s_i = \frac{e^{z_i}}{\sum_l e^{z_l}} \begin{array}{c} \rightarrow 0.224 \\ \rightarrow 0.672 \\ \rightarrow 0.091 \\ \rightarrow 0.013 \end{array} \right\} \mathbf{s} = \begin{pmatrix} s_1 \\ s_2 \\ s_3 \\ s_4 \end{pmatrix}$$

(Image by author)

The exponential function in the formula above ensures that the obtained values are non-negative. Due to the normalization term in the denominator the obtained values sum to 1. Furthermore, all values lie between 0 and 1. An important property of the softmax function is that it preserves the rank order of its input values:

$$-1.7 < 0.2 < 1.1 < 2.2 \rightarrow 0.013 < 0.091 < 0.224 < 0.672$$

Jacobian of the Softmax Function

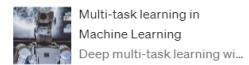
Formally, the softmax function is a so called *vector function*, which takes a vector as input and produces a vector as output:

$$\text{softmax: } \mathbb{R}^n \rightarrow \mathbb{R}^n$$

Therefore, when talking about the derivative of the softmax function, we actually talk about its Jacobian matrix, which is the matrix of all first-order partial derivatives:

$$J_{\text{softmax}} = \begin{pmatrix} \frac{\partial s_1}{\partial z_1} & \frac{\partial s_1}{\partial z_2} & \dots & \frac{\partial s_1}{\partial z_n} \\ \frac{\partial s_2}{\partial z_1} & \frac{\partial s_2}{\partial z_2} & \dots & \frac{\partial s_2}{\partial z_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial s_n}{\partial z_1} & \frac{\partial s_n}{\partial z_2} & \dots & \frac{\partial s_n}{\partial z_n} \end{pmatrix}$$

Related



Multi-task learning in Machine Learning
Deep multi-task learning wi...



Data Noise and Label Noise in Machine Learning



Contrasting Contrastive Learning Approaches



Class weights for categorical loss

where

$$s_i = \frac{e^{z_i}}{\sum_{l=1}^n e^{z_l}}, \quad \forall i = 1, \dots, n$$

Notice, how each output of the softmax function depends on all the input values (due to the denominator). For this reason the off-diagonal elements of the Jacobian aren't zero.

Since the outputs of the softmax function are strictly positive values, we can make the following derivation super short, by applying the following trick: instead of taking the partial derivative of the output, we take the partial derivative of the log of the output (also called "logarithmic derivative"):

$$\frac{\partial}{\partial z_j} \log(s_i) = \frac{1}{s_i} \cdot \frac{\partial s_i}{\partial z_j}$$

where the expression on the right-hand side follows directly from the chain rule. Next, we rearrange the upper formula and obtain:

$$\frac{\partial s_i}{\partial z_j} = s_i \cdot \frac{\partial}{\partial z_j} \log(s_i)$$

The left-hand side is exactly the partial derivative we're looking for. As we will shortly see, the right-hand side simplifies the computation of the derivative, such that we don't require the quotient rule of derivatives. We must first take the logarithm of s :

$$\log s_i = \log \left(\frac{e^{z_i}}{\sum_{l=1}^n e^{z_l}} \right) = z_i - \log \left(\sum_{l=1}^n e^{z_l} \right)$$

The partial derivative of the resulting expression is:

$$\frac{\partial}{\partial z_j} \log s_i = \frac{\partial z_i}{\partial z_j} - \frac{\partial}{\partial z_j} \log \left(\sum_{l=1}^n e^{z_l} \right)$$

Let's have a look at the first term on the right-hand side:

$$\frac{\partial z_i}{\partial z_j} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

which can be concisely written using the indicator function $1\{\cdot\}$. The indicator function takes on a value of 1 if its argument is true, and 0 otherwise.

The second term on the right-hand side can be evaluated by applying the chain rule:

$$\frac{\partial}{\partial z_j} \log s_i = 1\{i = j\} - \frac{1}{\sum_{l=1}^n e^{z_l}} \cdot \left(\frac{\partial}{\partial z_j} \sum_{l=1}^n e^{z_l} \right)$$

In the step above we used the derivative of the natural logarithm:

$$\frac{d}{dx} \log(x) = \frac{1}{x}$$

Obtaining the partial derivative of the sum is trivial:

$$\frac{\partial}{\partial z_j} \sum_{l=1}^n e^{z_l} = \frac{\partial}{\partial z_j} [e^{z_1} + e^{z_2} + \dots + e^{z_j} + \dots + e^{z_n}] = \frac{\partial}{\partial z_j} [e^{z_j}] = e^{z_j}$$

Plugging the result into the formula yields:

$$\frac{\partial}{\partial z_j} \log s_i = 1\{i=j\} - \frac{e^{z_j}}{\sum_{l=1}^n e^{z_l}} = 1\{i=j\} - s_j$$

Finally, we have to multiply the upper expression with s , as shown at the beginning of this section:

$$\frac{\partial s_i}{\partial z_j} = s_i \cdot \frac{\partial}{\partial z_j} \log(s_i) = s_i \cdot (1\{i=j\} - s_j)$$

This concludes our derivation. We've obtained a formula for all the elements of the Jacobi matrix (both diagonal and off-diagonal). For the special case of $n = 4$ we get:

$$J_{softmax} = \begin{pmatrix} s_1 \cdot (1 - s_1) & -s_1 \cdot s_2 & -s_1 \cdot s_3 & -s_1 \cdot s_4 \\ -s_2 \cdot s_1 & s_2 \cdot (1 - s_2) & -s_2 \cdot s_3 & -s_2 \cdot s_4 \\ -s_3 \cdot s_1 & -s_3 \cdot s_2 & s_3 \cdot (1 - s_3) & -s_3 \cdot s_4 \\ -s_4 \cdot s_1 & -s_4 \cdot s_2 & -s_4 \cdot s_3 & s_4 \cdot (1 - s_4) \end{pmatrix}$$

See how the diagonal elements differ from the off-diagonal elements.

Categorical Cross-Entropy Loss

Categorical cross-entropy loss is closely related to the softmax function, since it's practically only used with networks with a softmax layer at the output. Before we formally introduce the categorical cross-entropy loss (often also called softmax loss), we shortly have to clarify two terms: **multi-class classification** and **cross-entropy**.

Classification problems can be subdivided into the following two categories:

- **multi-class classification**, where each sample belongs to only one class (mutually exclusive)
- **multi-label classification**, where each sample may belong to multiple classes (or to no class)

The categorical cross-entropy loss is exclusively used in multi-class classification tasks, where each sample belongs exactly to one of the C classes. The true label assigned to each sample consists hence of a single integer value between 0 and $C-1$. The label can be represented by an one-hot encoded vector of size C , which has the value one for the correct class and zero everywhere else, see example below for $C = 4$:

$$\text{label} = 2 \rightarrow \mathbf{y} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$$

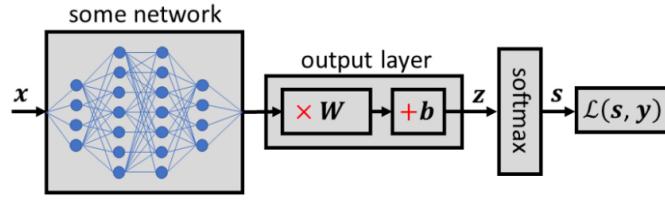
Cross-entropy takes as input two discrete probability distributions (simply vectors whose elements lie between 0,...,1 and sum to 1) and outputs a single real-valued (!) number representing the similarity of both probability distributions:

$$\mathcal{L}(\mathbf{y}, \mathbf{s}) = - \sum_i^C y_i \cdot \log(s_i)$$

where C denotes the number of different classes and the subscript i denotes i -th element of the vector. The smaller the cross-entropy, the more similar the two probability distributions are.

When cross-entropy is used as loss function in a multi-class classification task, then y is fed with the one-hot encoded label and the probabilities generated by the softmax layer are put in s . This way round we won't take the logarithm of zeros, since mathematically softmax will never really produce zero values.

By minimizing the loss during training, we basically force the predicted probabilities to gradually resemble the true one-hot encoded vectors.



In order to kick off the backpropagation process, as described in [this post](#), we have to calculate the derivative of the loss w.r.t to *weighted input* z of the output layer, see figure above:

$$\frac{\partial \mathcal{L}}{\partial z_j} = -\frac{\partial}{\partial z_j} \sum_{i=1}^C y_i \cdot \log(s_i) = -\sum_{i=1}^C y_i \cdot \frac{\partial}{\partial z_j} \log(s_i) = -\sum_{i=1}^C y_i \cdot \frac{\partial s_i}{\partial z_j}$$

Let us plug in the derivatives we have obtained in the last section:

$$\frac{\partial \mathcal{L}}{\partial z_j} = -\sum_{i=1}^C \frac{y_i}{s_i} \cdot s_i \cdot (1\{i=j\} - s_j) = -\sum_{i=1}^C y_i \cdot (1\{i=j\} - s_j)$$

We see, that j -th summand is the only one where $1\{i=j\} = 1$, so we pull it out of the sum:

$$\frac{\partial \mathcal{L}}{\partial z_j} = \left[-\sum_{i \neq j} y_i \cdot (1\{i=j\} - s_j) \right] - y_j \cdot (1\{j=j\} - s_j)$$

For all remaining summands we have $1\{i=j\} = 0$:

$$\frac{\partial \mathcal{L}}{\partial z_j} = \left[\sum_{i \neq j} y_i \cdot s_j \right] - y_j \cdot (1 - s_j) = \left[\sum_{i \neq j} y_i \cdot s_j \right] - y_j + y_j \cdot s_j$$

Now we can pull the j -th summands again inside the sum:

$$\frac{\partial \mathcal{L}}{\partial z_j} = \sum_{i=1}^C y_i \cdot s_j - y_j = s_j \cdot \sum_{i=1}^C y_i - y_j = s_j - y_j$$

In the last step we used the fact, that y is a probability distribution and sums to 1. In the concise vector notation we get:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}} = \mathbf{s} - \mathbf{y}$$

More readings

[Eli Bendersky's website](#)

[Raúl Gómez blog](#)

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

 [Get this newsletter](#)

Emails will be sent to tamalmondal495@gmail.com.
[Not you?](#)

Data Science Machine Learning Artificial Intelligence Optimization

 157  3



[Follow](#)

More from Towards Data Science

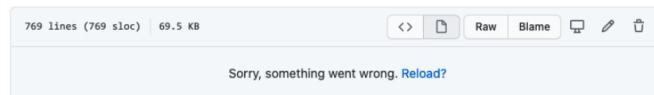
Your home for data science. A Medium publication sharing concepts, ideas and codes.

Sidney Kung · Apr 23 ★

Jupyter Notebook Not Rendering on GitHub? Here's a Simple Solution.

You can use this website instead!

Have you ever tried to preview your Jupyter Notebook file and received this pesky error message after about 20+ seconds of waiting?



About 70% of the time, we see this message. As if it's a gamble to see whether your `.ipynb` file will render or not.

Good news, this “Sorry...

[Read more · 2 min read](#)

 118 



Share your ideas with millions of readers. [Write on Medium](#)

Andrew Rothman · Apr 23 ★

GLMs Part I: A Rigorous Mathematical Formulation

Intuition for Unifying Theory of GLMs with Derivations in Canonical and Non-Canonical Forms





Photo by [ThisisEngineering.RAEng](#) on [Unsplash](#)

1: Background and Motivation

Generalized Linear Models (GLMs) play a critical role in fields including Statistics, Data Science, Machine Learning, and other computational sciences. This class of models are a generalization of ordinary linear regression for certain response variable types with error distribution models other than a normal distribution.

To refresh our memories on...

[Read more - 8 min read](#)

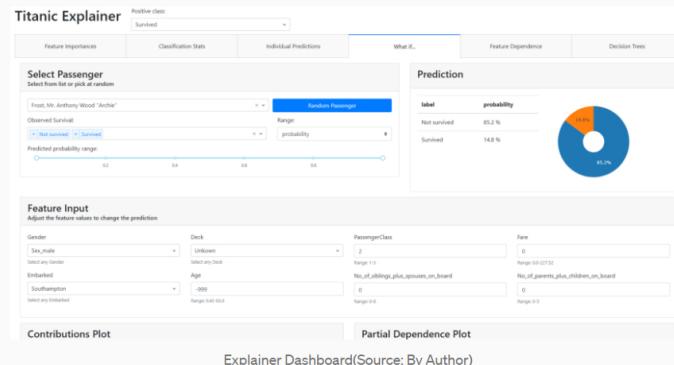
85 Q 1

↑ ↓ ...

Himanshu Sharma · Apr 22 ★

Machine Learning Model Dashboard

Creating dashboards to interpret machine learning model



Nowadays, creating a machine learning model is easy because of different python libraries that are in the market like sklearn, lazypredict, etc. These libraries are easy to use and used to create different types of models along with different types of visualizations and finding out the model performance. ...

[Read more · 4 min read](#)

361 2

...

Barr Moses · Apr 22

The Right Way to Measure ROI on Data Quality

Introducing a better metric for calculating the cost of bad data at your company



Image courtesy of Barr Moses.

To quote a friend, “Building your data stack without factoring in [data quality](#) is like buying a Ferrari but keeping it in the garage.”

In this article, guest columnist [Francisco Alberini](#), Product Manager at [Monte Carlo](#), introduces a better way measure the cost of bad data on your company.

Last...

[Read more · 6 min read](#)

120

...

Gianluca Malato · Apr 22 ★

How to calculate Heikin Ashi candles in Python for trading

Let's see how to use this powerful trading tool in Python





Example of Heikin Ashi candlestick chart. Image by the author.

Note from Towards Data Science's editors: While we allow independent authors to publish articles in accordance with our [rules and guidelines](#), we do not endorse each author's contribution. You should not rely on an author's works without seeking professional advice. See our [Reader Terms](#) for details.

Quantitative traders know that...

[Read more - 5 min read](#)

54

...

[Read more from Towards Data Science](#)



[About](#) [Write](#) [Help](#) [Legal](#)