

Forward-Feeding Neural Networks (FFNNs): Canonical Transformation Functions, Their Derivatives & Their Corresponding Cost Functions

STUFF TO KNOW BY HEART - EVEN WHEN DRUNK!

1. Three most common FFNN transformation functions: **Linear**, **Logistic** and **Softmax**
2. An FFNN's **top layer** almost always involves one of these three functions
3. **Logistic** function models probabilities of 2 binary states (OFF & ON); **Softmax** models relative probabilities of more than 2 states
4. Corresponding canonical cost functions: **Squared Error**, **Cross Entropy (Binary)** and **Cross Entropy (Multi-Class)**

This note discusses the three most common/standard (“canonical”) transformation functions used in layers of FFNNs. These are **continuous real-valued** functions that have **known, mathematically-convenient partial derivatives**. The top layer of an FFNN almost always involves one of these three transformation functions, and the FFNN’s hypothesized output is evaluated by one of three corresponding canonical cost functions.

LINEAR FUNCTION

The **Linear** function transforms a “fan-in” activation matrix \mathbf{A}_{in} (dimensions: $\langle \# \text{ of cases} \rangle \times \langle \# \text{ of “in” nodes} \rangle$) to “fan-out” activation matrix \mathbf{A}_{out} (dimensions: $\langle \# \text{ of cases} \rangle \times \langle \# \text{ of “out” nodes} \rangle$) via weights \mathbf{W} (dimensions: $\langle \# \text{ of “in” nodes} \rangle \times \langle \# \text{ of “out” nodes} \rangle$) by linearly multiplying \mathbf{A}_{in} and \mathbf{W} . This is the same as the Linear Regression model.

The Linear function and its backward partial derivative functions are defined below:

$$\begin{aligned}\mathbf{A}_{\text{out}} &= f(\mathbf{A}_{\text{in}}, \mathbf{W}) = \mathbf{Z}, \text{ where } \mathbf{Z} = \mathbf{A}_{\text{in}} \mathbf{W} \\ \frac{\partial v}{\partial \mathbf{A}_{\text{in}}} &= b_A \left(\frac{\partial v}{\partial \mathbf{A}_{\text{out}}}, \mathbf{A}_{\text{in}}, \mathbf{W}, \mathbf{A}_{\text{out}} \right) = \mathbf{W} \frac{\partial v}{\partial \mathbf{A}_{\text{out}}} \\ \frac{\partial v}{\partial \mathbf{W}} &= b_W \left(\frac{\partial v}{\partial \mathbf{A}_{\text{out}}}, \mathbf{A}_{\text{in}}, \mathbf{W}, \mathbf{A}_{\text{out}} \right) = \frac{\partial v}{\partial \mathbf{A}_{\text{out}}} \mathbf{A}_{\text{in}}\end{aligned}$$

The cost function corresponding to a Linear FFNN top layer is the **Squared Error** function, which measures the Euclidean geometric distance between Hypothesized Output \mathbf{H} from the “right-answer” Target Output \mathbf{Y} . This function has a convenient partial derivative with respect to \mathbf{H} .

$$\begin{aligned}c(\mathbf{H}, \mathbf{Y}) &= \frac{1}{2} \frac{|\mathbf{H} - \mathbf{Y}|^2}{\text{No. of Cases}} \\ \frac{\partial c}{\partial \mathbf{H}} &= \frac{(\mathbf{H} - \mathbf{Y})^T}{\text{No. of Cases}}\end{aligned}$$

LOGISTIC FUNCTION

The **Logistic** function transforms a “fan-in” activation matrix \mathbf{A}_{in} (dimensions: $\langle \# \text{ of cases} \rangle \times \langle \# \text{ of “in” nodes} \rangle$) to “fan-out” activation matrix \mathbf{A}_{out} (dimensions: $\langle \# \text{ of cases} \rangle \times \langle \# \text{ of “out” nodes} \rangle$) via weights \mathbf{W} (dimensions: $\langle \# \text{ of “in” nodes} \rangle \times \langle \# \text{ of “out” nodes} \rangle$) by first linearly multiplying \mathbf{A}_{in} and \mathbf{W} , and then “squashing” the resulting values into the (0,1) unit interval. This is the same as the Logistic Regression model: the output represents probabilities, with values near 0 representing “likely OFF” and values near 1 representing “likely ON”.

The Logistic function and its backward partial derivative functions are defined below:

$$\begin{aligned}\mathbf{A}_{\text{out}} &= f(\mathbf{A}_{\text{in}}, \mathbf{W}) = \mathbf{1} ./ (\mathbf{1} + \exp(-\mathbf{Z})), \text{ where } \mathbf{Z} = \mathbf{A}_{\text{in}} \mathbf{W} \\ \frac{\partial v}{\partial \mathbf{A}_{\text{in}}} &= b_A \left(\frac{\partial v}{\partial \mathbf{A}_{\text{out}}}, \mathbf{A}_{\text{in}}, \mathbf{W}, \mathbf{A}_{\text{out}} \right) = \mathbf{W} \mathbf{B} \\ \frac{\partial v}{\partial \mathbf{W}} &= b_W \left(\frac{\partial v}{\partial \mathbf{A}_{\text{out}}}, \mathbf{A}_{\text{in}}, \mathbf{W}, \mathbf{A}_{\text{out}} \right) = \mathbf{B} \mathbf{A}_{\text{in}} \\ \text{where: } \mathbf{B} &= \frac{\partial v}{\partial \mathbf{A}_{\text{out}}} .* (\mathbf{A}_{\text{out}})^T .* (\mathbf{1} - \mathbf{A}_{\text{out}})^T\end{aligned}$$

The Target Output \mathbf{Y} for a Logistic FFNN top layer is a matrix of binary values 0 (“OFF”) and 1 (“ON”), and the cost function corresponding to a Logistic FFNN top layer is the **Cross Entropy (Binary)** function, which is defined below together with its partial derivatives:

$$\begin{aligned}c(\mathbf{H}, \mathbf{Y}) &= - \frac{\text{sumAllDims}(\mathbf{Y} .* \ln(\mathbf{H}) + (\mathbf{1} - \mathbf{Y}) .* \ln(\mathbf{1} - \mathbf{H}))}{\text{No. of Cases}} \\ \frac{\partial c}{\partial \mathbf{H}} &= - \frac{(\mathbf{Y} ./ \mathbf{H} - (\mathbf{1} - \mathbf{Y}) ./ (\mathbf{1} - \mathbf{H}))^T}{\text{No. of Cases}}\end{aligned}$$

SOFTMAX FUNCTION

The **Softmax** function generalizes the Logistic function. While the Logistic function models probabilities of 2 binary states OFF and ON, the Softmax function models **relative probabilities of more than 2 states**, with such probabilities summing to 1.

The Softmax function and its partial derivatives are defined below:

$$\begin{aligned}\mathbf{A}_{\text{out}} &= f(\mathbf{A}_{\text{in}}, \mathbf{W}) = \exp(\mathbf{Z}) ./ \text{rowwiseSum}(\exp(\mathbf{Z})), \text{ where } \mathbf{Z} = \mathbf{A}_{\text{in}} \mathbf{W} \\ \frac{\partial v}{\partial \mathbf{A}_{\text{in}}} &= b_A \left(\frac{\partial v}{\partial \mathbf{A}_{\text{out}}}, \mathbf{A}_{\text{in}}, \mathbf{W}, \mathbf{A}_{\text{out}} \right) = \mathbf{W} \mathbf{B} \\ \frac{\partial v}{\partial \mathbf{W}} &= b_W \left(\frac{\partial v}{\partial \mathbf{A}_{\text{out}}}, \mathbf{A}_{\text{in}}, \mathbf{W}, \mathbf{A}_{\text{out}} \right) = \mathbf{B} \mathbf{A}_{\text{in}} \\ \text{where } \mathbf{B} &\text{ is a complicated, cumbersome but computable function}\end{aligned}$$

The Target Output \mathbf{Y} (dimensions: <# of cases> x <# of output states>) for a Logistic FFNN top layer is a matrix of binary values 0 (“OFF”) and 1 (“ON”), with each row having only one single 1 value representing the “ON” state for the case. The cost function corresponding to a Softmax FFNN top layer is the **Cross Entropy (Multi-Class)** function, which is defined below:

$$\begin{aligned}c(\mathbf{H}, \mathbf{Y}) &= - \frac{\text{sumAllDims}(\mathbf{Y} .* \ln(\mathbf{H}))}{\text{No. of Cases}} \\ \frac{\partial c}{\partial \mathbf{H}} &= - \frac{(\mathbf{Y} ./ \mathbf{H})^T}{\text{No. of Cases}}\end{aligned}$$

IMPLEMENTATION NOTE

In actual implementation, the partial derivative $\frac{\partial c}{\partial \mathbf{H}}$ is often numerically unstable to compute for a Logistic or Softmax FFNN top layer. One convenient and stable work-around applicable for all three canonical Linear, Logistic and Softmax functions is to compute $\frac{\partial c}{\partial \mathbf{Z}}$, where $\mathbf{Z} = \mathbf{A}_{\text{in}} \mathbf{W}$ for the FFNN’s top layer. Interestingly, all three canonical functions have the very same and extremely convenient form for $\frac{\partial c}{\partial \mathbf{Z}}$:

$$\frac{\partial c}{\partial \mathbf{Z}} = \frac{(\mathbf{H} - \mathbf{Y})^T}{\text{No. of Cases}}$$