

Applied Artificial Intelligence

Session 13: Feedforward Multilayer Neural Networks II

Fall 2018

NC State University

Lecturer: Dr. Behnam Kia

Course Website: <https://appliedai.wordpress.ncsu.edu/>

What will we learn in this session?

- Where those cost functions come from?
- Why type of problems can I solve with neural networks?
- New homework

The Problem Definition

- Imagine we have a set of observed data:

$$x_1, x_2, \dots, x_m$$

- Assume we know which parametric probability distribution function ($p(x|w)$) has produced the data, but we don't know the parameters ($w?$).
- Question: How to Estimate the Parameters?

The Problem Definition

- Imagine we have a set of observed data:

$$x_1, x_2, \dots, x_m$$

- Assume we know which parametric probability distribution function ($p(x|w)$) has produced the data, but we don't know the parameters ($w?$).
- Question: How to Estimate the Parameters

Common format: $p_w(x)$ is usually written as $p(x|w)$ as well to indicate the dependence of probability values on the parameter w .

The Problem Definition

- Imagine we have a set of observed data:

$$x_1, x_2, \dots, x_m$$

- Assume we know which parametric probability distribution function ($p(x|w)$) has produced the data, but we don't know the parameters ($w?$).
- **Question: How to Estimate the Parameters**

Maximum Likelihood Estimation (MLE)

- MLE is an approach to estimate parameters.
- MLE: Choose parameters in a way that the observed data have the biggest probability.

$$w = \underset{w}{\operatorname{argmax}} P(\text{Data}|w)$$

- MLE tries to makes sense of observed data by making them more likely to be produced by the parametric model.

General Case

MLE says: $w_{MLE} = \underset{w}{\operatorname{argmax}} P(\text{Data}|w)$

$$\begin{aligned} P(\text{Data}|w) &= P(x_1, x_2, \dots, x_m | w) \\ &= P(x_1 | w) P(x_2 | w) \dots P(x_m | w) \quad (\text{i.i.d assumption}) \\ &= \prod_{i=1}^m P(x_i | w) \end{aligned}$$

Example: Flipping a Coin

- Assume the observed data are the outcomes of flipping a coin:
- Data: H, T, T, H, H, H, T, H, H, ..., T
- Question: What is parameter p , ($p=P(H)$)?

Example: Flipping a Coin

- Assume the observed data are the outcomes of flipping a coin:
- Data: H, T, T, H, H, H, T, H, H, ..., T
- Question: What is parameter p ?
- MLE says: $p = \operatorname{argmax}_p P(\text{Data}|p)$

Log Likelihood for Numerical Stability

$$\begin{aligned} P(\text{Data}|w) &= P(x_1, x_2, \dots, x_m|w) \\ &= P(x_1|w)P(x_2|w) \dots P(x_m|w) \quad (\text{i.i.d assumption}) \\ &= \prod_{i=1}^m P(x_i|w) \end{aligned}$$

$$\begin{aligned} \log(P(\text{Data}|w)) &= \\ &= \sum_{i=1}^m \log(P(x_i|w)) \end{aligned}$$

$$w_{MLE} = \underset{w}{\operatorname{argmax}} \log(P(\text{Data}|w))$$

MLE Example: Flipping a Coin

Calculus-based Solution

$$\begin{aligned}P(\text{Data}|p) &= P(H, T, T, \dots, T|p) \\ &= P(H|p)P(T|p)P(T|p) \dots P(T|p) \quad (\text{i.i.d. assumption}) \\ &= p^k(1-p)^{n-k}\end{aligned}$$

$$p_{MLE} = \underset{p}{\operatorname{argmax}} (P(\text{Data}|p))$$

MLE Example: Flipping a Coin

Calculus-based Solution

$$\begin{aligned}P(\text{Data}|p) &= P(H, T, T, \dots, T|p) \\ &= P(H|p)P(T|p)P(T|p) \dots P(T|p) \quad (\text{i.i.d. assumption}) \\ &= p^k(1-p)^{n-k}\end{aligned}$$

$$p_{MLE} = \underset{p}{\operatorname{argmax}} (P(\text{Data}|p))$$

- $\frac{\partial}{\partial p} P(\text{Data}|p) = 0$

MLE Example: Flipping a Coin

Calculus-based Solution

$$\begin{aligned}P(\text{Data}|p) &= P(H, T, T, \dots, T|p) \\&= P(H|p)P(T|p)P(T|p) \dots P(T|p) \quad (\text{i.i.d assumption}) \\&= p^k(1-p)^{n-k}\end{aligned}$$

$$p_{MLE} = \underset{p}{\operatorname{argmax}} (P(\text{Data}|p))$$

$$\frac{\partial}{\partial p} p^k(1-p)^{n-k} = kp^{k-1}(1-p)^{n-k} - (n-k)p^k(1-p)^{n-k-1} = 0$$

$$k(1-p) - (n-k)p = 0$$

$$k - np = 0$$

$$p = k/n$$

Supervised Machine Learning

- Imagine we have a set of training data:
 $(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)$
- And our neural network is estimating $p(y|x, w)$, but of course we don't know the optimal parameters ($w?$).
- Question: How to Estimate the Parameters?

Supervised Machine Learning

- Imagine we have a set of training data:
 $(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)$
- And our neural network is estimating $p(y|x, w)$, but of course we don't know the optimal parameters (w ?).
- Question: How to Estimate the Parameters?

$$w_{MLE} = \underset{w}{\operatorname{argmax}} \log(P(y|x, w))$$

Supervised Machine Learning

- Imagine we have a set of training data:

$$(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)$$

$$\begin{aligned} w_{MLE} &= \operatorname{argmax}_w \log P(Y|X, w) \\ &= \operatorname{argmax}_w \sum_{i=1}^m \log P(y^i | x^i, w) \end{aligned}$$

Supervised Machine Learning

Logistic Regression

- Imagine we have a set of training data:

$$(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)$$

$$\begin{aligned} w_{MLE} &= \operatorname{argmax}_w \log P(Y|X, w) \\ &= \operatorname{argmax}_w \sum_{i=1}^m \log P(y^i | x^i, w) \end{aligned}$$

Reminder: in logistic regression:

$$\begin{aligned} \hat{y}^i &\text{ is } P(\text{output} = 1 | x^i, w) \\ 1 - \hat{y}^i &\text{ is } P(\text{output} = 0 | x^i, w) \end{aligned}$$

Supervised Machine Learning

Logistic Regression

- Imagine we have a set of training data:

$$(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)$$

$$\begin{aligned} w_{MLE} &= \operatorname{argmax}_w \log P(Y|X, w) \\ &= \operatorname{argmax}_w \sum_{i=1}^m \log P(y^i | x^i, w) \end{aligned}$$

Reminder: in logistic regression:

$$\begin{aligned} \hat{y}^i &\text{ is } P(\text{output} = 1 | x^i, w) \\ 1 - \hat{y}^i &\text{ is } P(\text{output} = 0 | x^i, w) \end{aligned}$$

$$w_{MLE} = \operatorname{argmax}_w \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)]$$

Supervised Machine Learning

Logistic Regression

- Imagine we have a set of training data:

$$(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)$$

$$\begin{aligned} w_{MLE} &= \operatorname{argmax}_w \log P(Y|X, w) \\ &= \operatorname{argmax}_w \sum_{i=1}^m \log P(y^i | x^i, w) \end{aligned}$$

Reminder: in logistic regression:

$$\hat{y}^i \text{ is } P(\text{output} = 1 | x^i, w)$$

$$1 - \hat{y}^i \text{ is } P(\text{output} = 0 | x^i, w)$$

$$w_{MLE} = \operatorname{argmax}_w \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)]$$

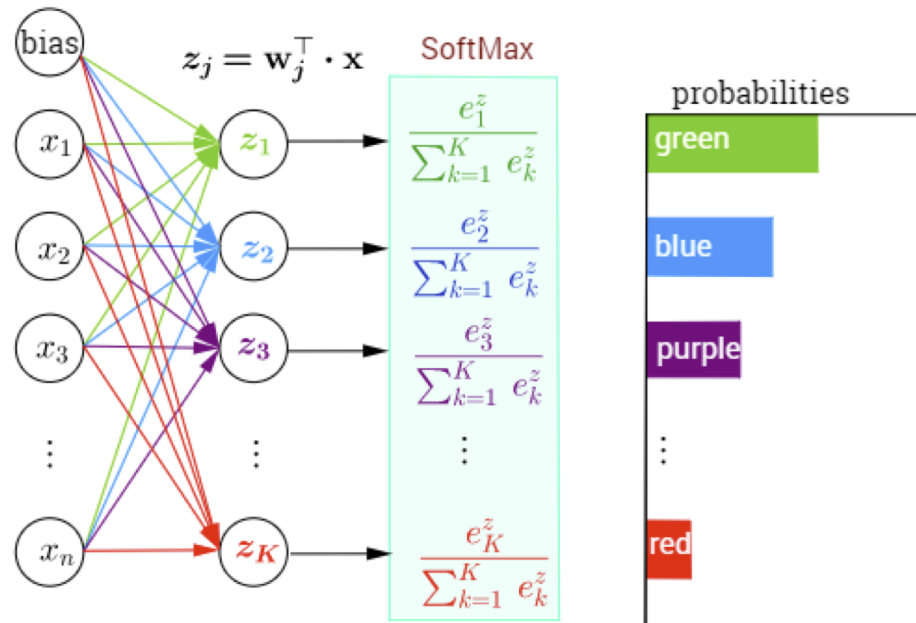
$$w_{optimal} = \operatorname{argmin}_w J(W)$$

$$= \operatorname{argmin}_w - \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)]_{19}$$

Classification with more than two categories

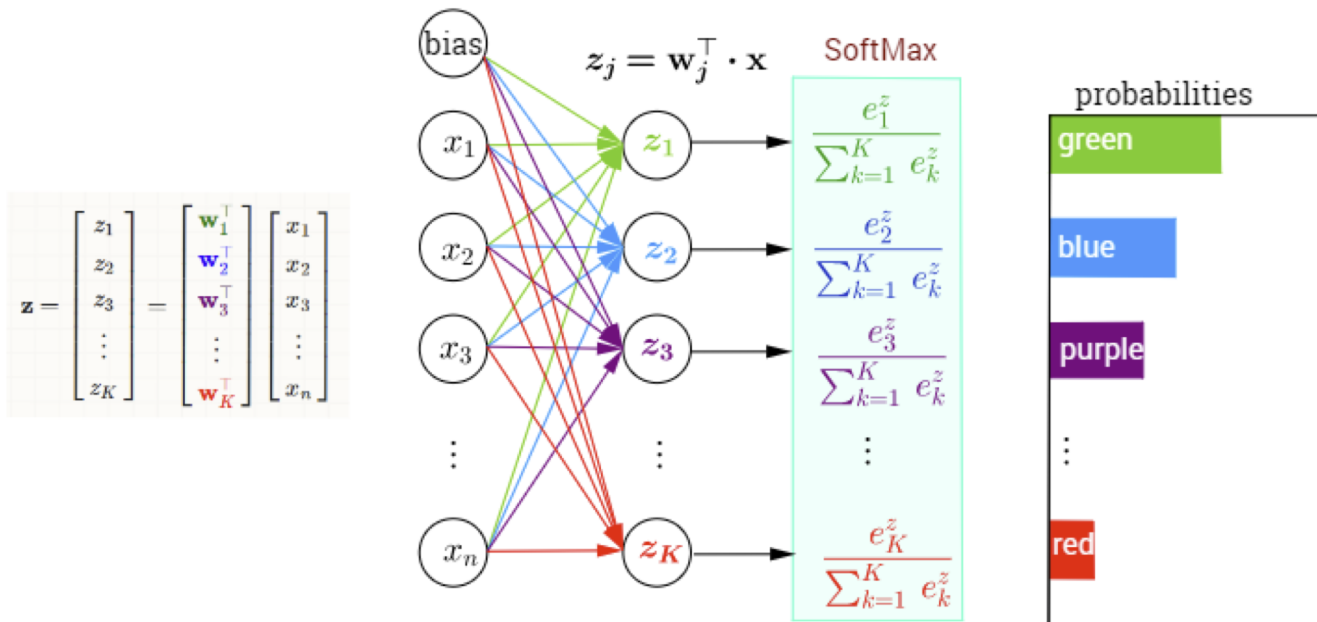
Multi-Class Classification with NN and SoftMax Function

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_K \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \mathbf{w}_3^\top \\ \vdots \\ \mathbf{w}_K^\top \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$



Classification with more than two categories

Multi-Class Classification with NN and SoftMax Function



$$w_{MLE} = \underset{w}{\operatorname{argmax}} \log P(Y|X, w)$$

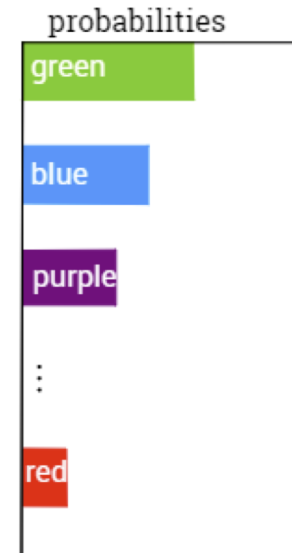
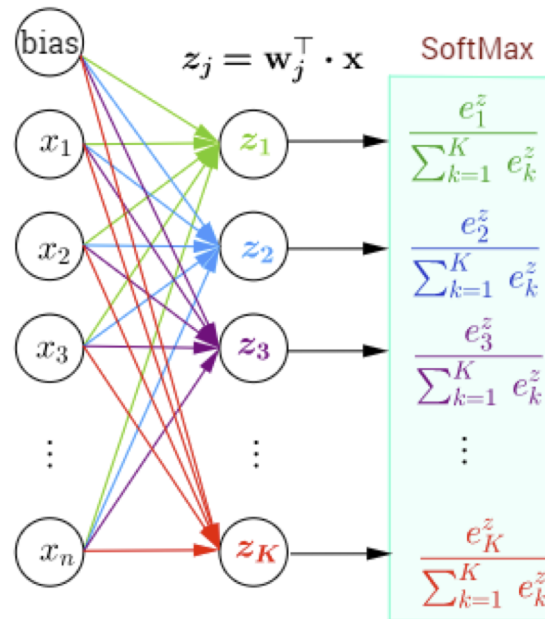
$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^m \log P(y^i|x^i, w)$$

$$J(W) = -\sum_{i=1}^m \log P(y^i|x^i, w)$$

Classification with more than two categories

Multi-Class Classification with NN and SoftMax Function

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_K \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \mathbf{w}_3^\top \\ \vdots \\ \mathbf{w}_K^\top \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$



(X₁, Green)
(X₂, Purple)

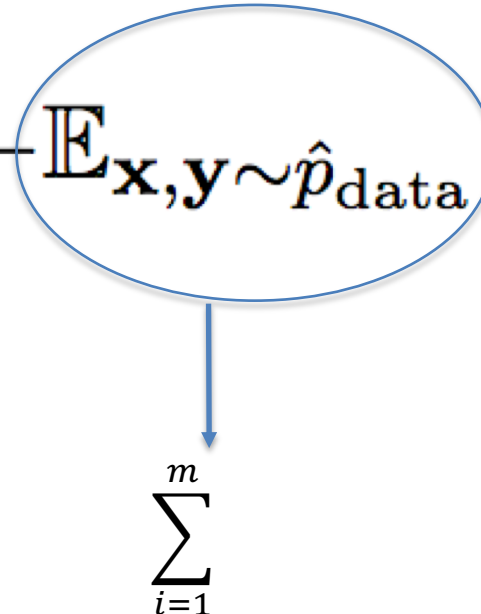
$$\begin{aligned} w_{MLE} &= \underset{w}{\operatorname{argmax}} \log P(Y|X, w) \\ &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^m \log P(y^i|x^i, w) \end{aligned}$$

$$J(W) = -\sum_{i=1}^m \log P(y^i|x^i, w)$$

Cross entropy between model distribution and the training data

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} \mid \mathbf{x}).$$

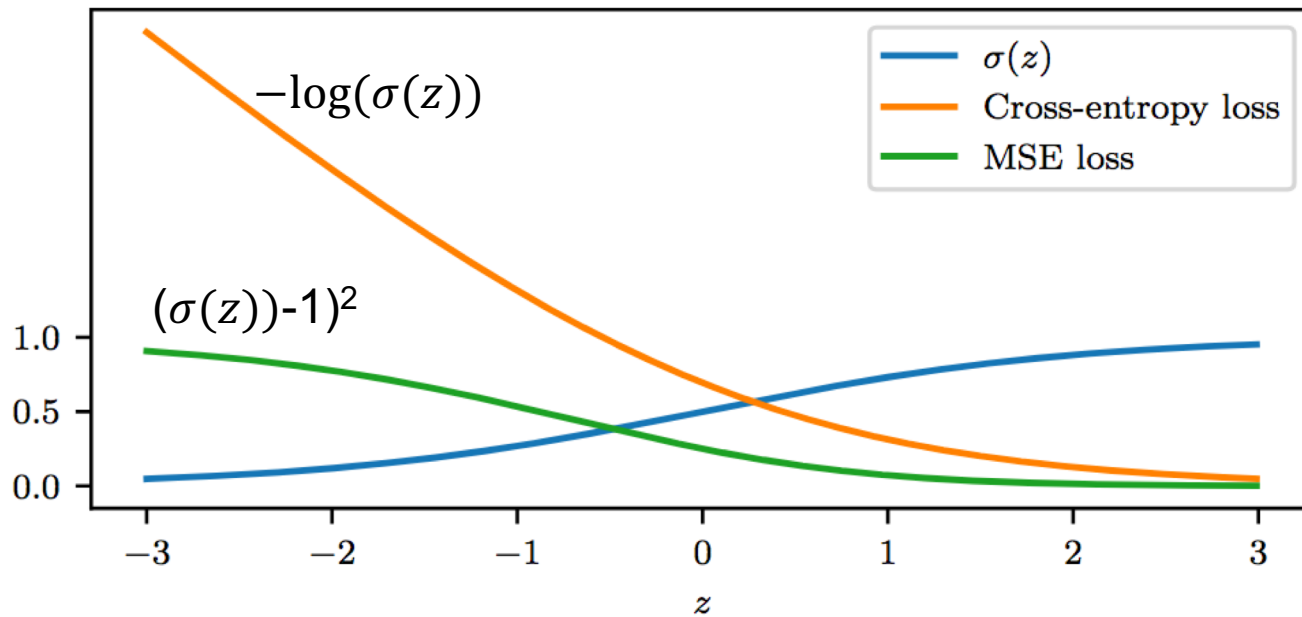
Cross entropy between model distribution and the training data

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} \mid \mathbf{x}).$$


The diagram illustrates the expansion of the expectation operator in the cross entropy formula. A blue oval encircles the expectation operator $\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}}$ in the equation above. A blue arrow points downwards from the bottom of this oval to a summation symbol $\sum_{i=1}^m$, indicating that the expectation is taken over m data points.

Output Type	Output Distribution	Output Layer	Cost Function
Binary	Bernoulli	Sigmoid	Binary cross-entropy
Discrete	Multinoulli	Softmax	Discrete cross-entropy
Continuous	Gaussian	Linear	Gaussian cross-entropy (MSE)
Continuous	Mixture of Gaussian	Mixture Density	Cross-entropy
Continuous	Arbitrary	See part III: GAN, VAE, FVBN	Various

Sigmoid output with target of 1



Neural Networks: What they can do?

Brain is a universal learning machine

Visual Projections Routed to the Auditory Pathway in Ferrets: Receptive Fields of Visual Neurons in Primary Auditory Cortex

Anna W. Roe,^a Sarah L. Pallas,^b Young H. Kwon, and Mriganka Sur

Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139

How does cortex that normally processes inputs from one sensory modality respond when provided with input from a different modality? We have addressed such a question with an experimental preparation in which retinal input is routed to the auditory pathway in ferrets. Following neonatal surgical manipulations, a specific population of retinal ganglion cells is induced to innervate the auditory thalamus and provides visual input to cells in auditory cortex (Sur et al., 1988). We have now examined in detail the visual response prop-

ulate nucleus, or MGN). Retinal afferents subsequently provide visual input to cells in primary auditory cortex (A1) of the “rewired” ferret (Sur et al., 1988) and establish a topographic visual map there (Roe et al., 1990a). We have now examined both qualitatively and quantitatively the physiological response properties of single visual units in A1 of rewired ferrets and compared them to the properties of cells in primary visual cortex (V1) of normal ferrets.

A preliminary report of these data has been published pre-

BrainPort® Vision Pro

The new BrainPort® Vision Pro is a 2nd generation oral electronic vision aid that provides electro-tactile stimulation to aid profoundly blind patients in orientation, mobility, and object recognition as an adjunctive device to other assistive methods such as the white cane or a guide dog.

BrainPort Vision Pro translates digital information from a wearable video camera into gentle electrical stimulation patterns on the surface of the tongue. Users feel moving bubble-like patterns on their tongue which they learn to interpret as the shape, size, location and motion of objects in their environment. Some users have described it as being able to “see with your tongue”.



Neural Networks

- One learning algorithm hypothesis
- Robotics, Speech Processing, Image Processing, NLP, etc, are not separate fields, and common learning algorithms can apply to all.