

Chapter 6

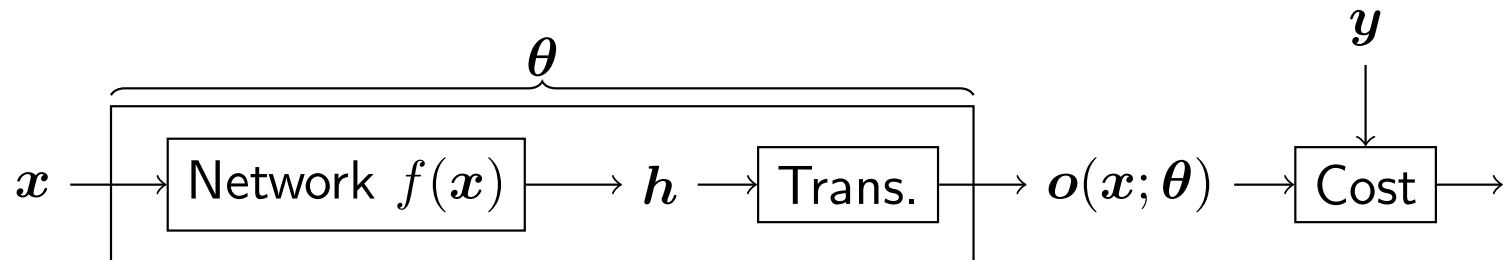
Deep Feedforward Networks

Learning XOR

- Study by yourself

Gradient-based Learning

- General setting



- x : Inputs
- $f(x)$: Feedforward network
- h : Hidden units computed by $f(x)$
- Trans.: Output layer transforming h to output $o(x; \theta)$
- $o(x; \theta)$: Output units parameterized by model parameters θ
- Cost: A function of ground-truth y and output o to be minimized w.r.t. model parameters θ

Cost Functions

- The maximum likelihood (ML) principle provides a guide for designing cost functions
- If we define a conditional distribution $p(y|\mathbf{x})$ as a distribution over y parameterized by the network outputs $\mathbf{o}(\mathbf{x}; \boldsymbol{\theta})$

$$p(y|\mathbf{x}) \triangleq p(y; \mathbf{o}(\mathbf{x}; \boldsymbol{\theta}))$$

- Then, the ML principle suggests we take the negative log-likelihood

$$-\log p(y; \mathbf{o}(\mathbf{x}; \boldsymbol{\theta}))$$

as the cost function to be minimized w.r.t. model parameters $\boldsymbol{\theta}$

Learning Conditional Statistics

- If we define

$$p(y|\mathbf{x}) \triangleq \mathcal{N}(y; o(\mathbf{x}; \boldsymbol{\theta}), \mathbf{I}),$$

- Then, minimizing the negative log-likelihood yields

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} E_{\mathbf{x}, y \sim p_{\text{data}}} \|y - o(\mathbf{x}; \boldsymbol{\theta})\|^2$$

- With sufficient capacity, the network will learn the conditional mean

$$o(\mathbf{x}; \boldsymbol{\theta}^*) \approx E[y|\mathbf{x}],$$

which predicts the mean value of y for each \mathbf{x}

- By the same token, if we define

$$p(y|\mathbf{x}) \triangleq \text{Laplace}(y; o(\mathbf{x}; \boldsymbol{\theta}), \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|y - o(\mathbf{x}; \boldsymbol{\theta})|}{\gamma}\right),$$

- Then, minimizing the negative log-likelihood yields

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} E_{\mathbf{x}, y \sim p_{\text{data}}} |y - o(\mathbf{x}; \boldsymbol{\theta})|$$

- With sufficient capacity, the network will learn the conditional median

$$o(\mathbf{x}; \boldsymbol{\theta}^*) \approx \text{Median}[y|\mathbf{x}],$$

which predicts the median value of y for each \mathbf{x}

- In other words, $o(\mathbf{x}; \boldsymbol{\theta}^*)$ satisfies

$$\int_{-\infty}^{o(\mathbf{x}; \boldsymbol{\theta}^*)} p(y|\mathbf{x}) dy = \int_{o(\mathbf{x}; \boldsymbol{\theta}^*)}^{\infty} p(y|\mathbf{x}) dy$$

Learning Gaussian Output Distributions

- The n -dimensional Gaussian distribution $\mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is given by

$$\mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

- A conditional Gaussian can be learned by treating $o(\mathbf{x}; \boldsymbol{\theta})$ as the means

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; o(\mathbf{x}; \boldsymbol{\theta}), \mathbf{I})$$

- In this model, the outputs $o(\mathbf{x}; \boldsymbol{\theta})$ often take a linear form

$$o(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}^T \mathbf{h}(\mathbf{x}) + \mathbf{b}$$

- As such, they are known as *linear output units*

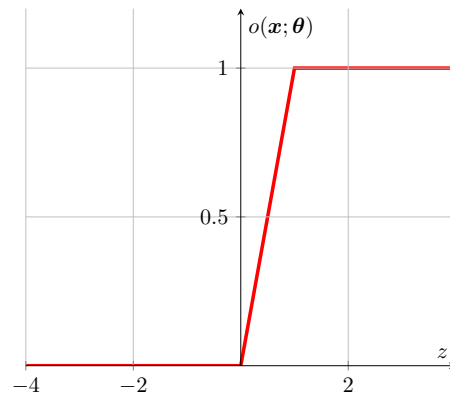
Learning Bernoulli Output Distributions

- A conditional Bernoulli can be learned by having $o(\mathbf{x}; \boldsymbol{\theta})$ as the very distribution parameter

$$p(y|\mathbf{x}) = o(\mathbf{x}; \boldsymbol{\theta})^y (1 - o(\mathbf{x}; \boldsymbol{\theta}))^{1-y}, \quad y \in \{0, 1\}$$

- The $o(\mathbf{x}; \boldsymbol{\theta})$ must be in $[0, 1]$ in order to be a valid parameter, one trivial implementation being

$$o(\mathbf{x}; \boldsymbol{\theta}) = \max \{0, \min \{1, z\}\}, \quad z = \mathbf{w}^T \mathbf{h}(\mathbf{x}) + b$$



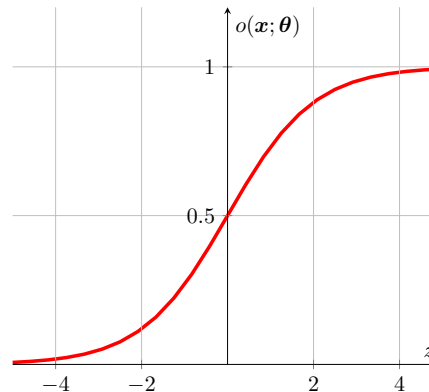
- This however may be problematic for the gradient of the log-likelihood

w.r.t. \mathbf{w}, b is zero when $z = \mathbf{w}^T h(\mathbf{x}) + b$ is outside the unit interval

$$-\nabla_{\mathbf{w}, b} \log p(y|\mathbf{x}) = -\frac{\partial \log p(y|\mathbf{x})}{\partial o(\mathbf{x}; \boldsymbol{\theta})} \underbrace{\frac{\partial o(\mathbf{x}; \boldsymbol{\theta})}{\partial z}}_{=0} \nabla_{\mathbf{w}, b} z$$

- The gradient-based learning may fail to learn \mathbf{w}, b properly
- An alternative approach is to use sigmoid units combined with the negative log-likelihood function

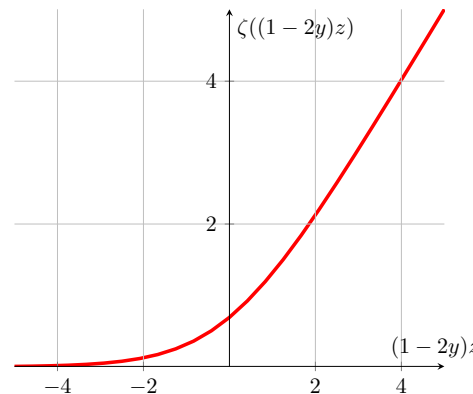
$$o(\mathbf{x}; \boldsymbol{\theta}) = \sigma(z), \text{ with } z = \mathbf{w}^T h(\mathbf{x}) + b, \sigma(z) = \frac{e^z}{1 + e^z}$$



- By definition, the negative log-likelihood can be computed as

$$\begin{aligned}
 -\log p(y|\mathbf{x}) &= -\log o(\mathbf{x}; \boldsymbol{\theta})^y (1 - o(\mathbf{x}; \boldsymbol{\theta}))^{1-y} \\
 &= -\log \sigma((2y - 1)z) \\
 &= \zeta((1 - 2y)z)
 \end{aligned}$$

where $y \in \{0, 1\}$ and $\zeta(\cdot)$ is the softplus function



- Ideally, we want $z = \mathbf{w}^T \mathbf{h}(\mathbf{x}) + b$ to be very positive (respectively, negative) when the ground-true $y = 1$ (respectively, $y = 0$)
- This suggests that samples correctly classified have very negative

$(1 - 2y)z$; in other words, these samples will not contribute to the gradient computation

$$-\nabla_{\mathbf{w},b} \log p(y|\mathbf{x}) = -\frac{\partial \log p(y|\mathbf{x})}{\partial \zeta} \underbrace{\frac{\partial \zeta}{\partial (1 - 2y)z}}_{=0} \frac{\partial (1 - 2y)z}{\partial z} \nabla_{\mathbf{w},b} z$$

- On the other hand, samples incorrectly classified have very positive $(1 - 2y)z$, which leads to strong gradients
- When sigmoid units are combined with other cost functions, e.g. mean squared error, the gradient vanishing problem may occur
- *The choice of output units is tightly coupled with that of cost functions*

Multinoulli Random Variables

- A Multinoulli random variable x has n possible values with distribution

$$P(x; \alpha) = \prod_{i=1}^n (\alpha_i)^{\mathbf{1}_{x=i}}, \quad x = 1, 2, \dots, n$$

where

$$\alpha_i \in [0, 1], \quad \forall i$$

$$\sum_{i=1}^n \alpha_i = 1$$

$$\mathbf{1}_{x=i} = \begin{cases} 1, & \text{if } x = i \\ 0, & \text{if } x \neq i \end{cases}$$

- When using a one-hot vector (or the 1-of-K coding) x , in which only one element equals to 1 with the others being 0, to represent the n

possible values, the Multinoulli distribution is given by

$$P(\mathbf{x}; \boldsymbol{\alpha}) = \prod_{i=1}^n (\alpha_i)^{x_i}, \quad x_i \in [0, 1]$$

- Observe that Bernoulli is a special case of Multinoulli with $n = 2$

Learning Multinoulli Output Distributions

- A conditional Multinoulli can be learned by having its distribution parameters be modeled by softmax output units $o(\mathbf{x}; \boldsymbol{\theta})$

$$p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n (o(\mathbf{x}; \boldsymbol{\theta})_i)^{y_i},$$

where \mathbf{y} is a one-hot vector and

$$o(\mathbf{x}; \boldsymbol{\theta})_i = \text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_{j=1}^n \exp(z_j)}$$

$$\mathbf{z} = \mathbf{W}^T \mathbf{h}(\mathbf{x}) + \mathbf{b}$$

- When the difference between the maximal element of \mathbf{z} and the others becomes large, softmax becomes a form of winner-take-all, namely, one output is nearly 1 and the others are nearly 0
- As such, it is viewed as a softened version of $\arg \max$ with a one-hot

representation

- Softmax overparameterizes the distribution: the n -th probability may be obtained by subtracting the first $n - 1$ probabilities from 1
- One may require that one element of \mathbf{z} be 0; this leads to sigmoid when $n = 2$

$$\sigma(z) = \frac{\exp(z)}{\exp(z) + \exp(0)} = \frac{1}{1 + \exp(-z)}$$

- As with sigmoid units, softmax units almost always come with the log-likelihood function

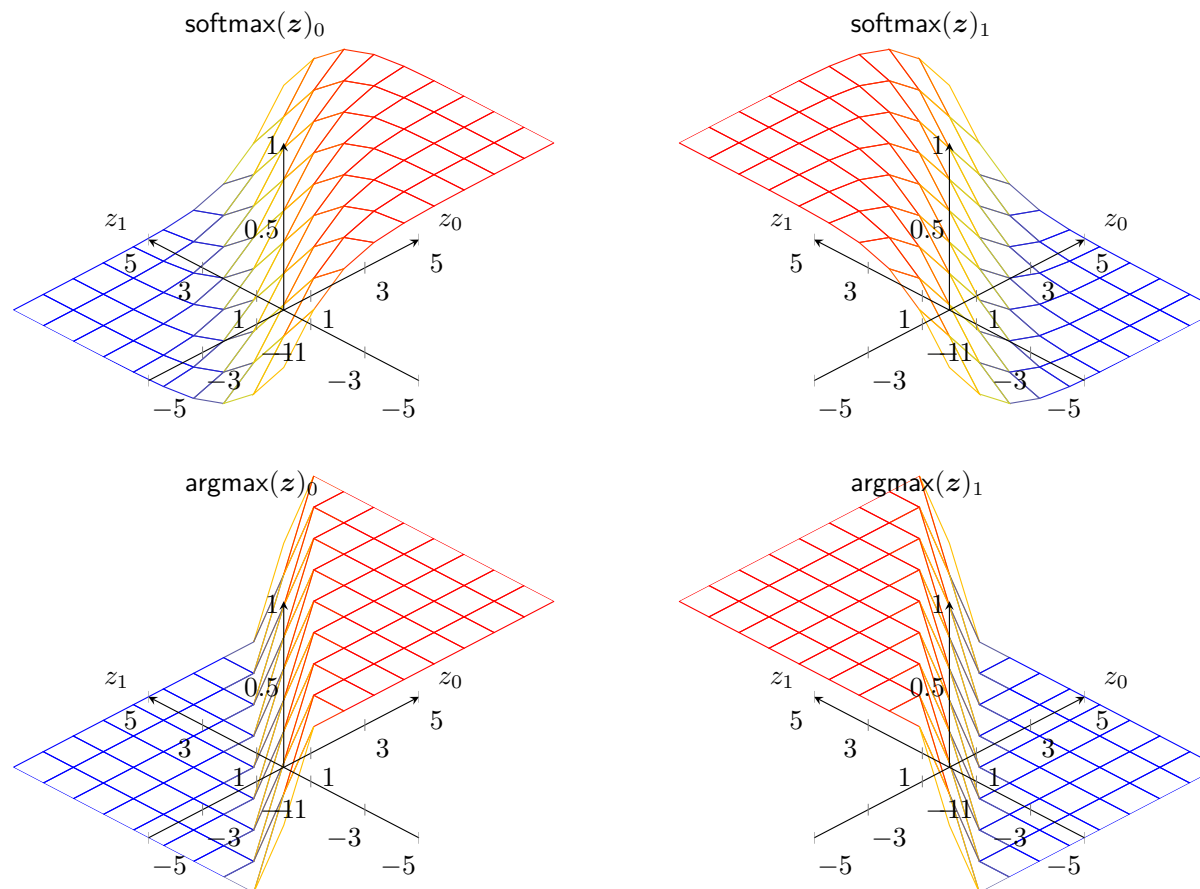
$$-\log p(\mathbf{y}|\mathbf{x}) = -\sum_{i=1}^n y_i \log(o(\mathbf{x}; \boldsymbol{\theta})_i) = -\sum_{i=1}^n y_i (z_i - \log \sum_{j=1}^n \exp(z_j))$$

- Assuming $y_i = 1$ and $y_j = 0, j \neq i$, we have

$$-\log p(\mathbf{y}|\mathbf{x}) = -(z_i - \log \sum_{j=1}^n \exp(z_j))$$

- This suggests that when the ground truth $y_i = 1$, minimizing the negative log-likelihood amounts to maximizing z_i and penalizes the most active incorrect prediction z_j if $z_j \gg z_k, k \neq j$
- On the other hand, samples correctly classified (i.e., $z_i \gg z_k, k \neq i$) will contribute little to the gradient computation; in this case, $(z_i - \log \sum_{j=1}^n \exp(z_j)) \approx 0$
- Observe that the log undoes the saturating effect of the softmax

Softmax vs. Argmax



Relative Frequencies

- With sufficient capacity, minimizing the negative log-likelihood will

$$\begin{aligned}
 -E_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p(\mathbf{y} | \mathbf{x}) &= -E_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \left[\sum_{i=1}^n y_i \log o(\mathbf{x}; \boldsymbol{\theta})_i \right] \\
 &= - \sum_{j=1}^m \sum_{i=1}^n y_i^{(j)} \log o(\mathbf{x}^{(j)}; \boldsymbol{\theta})_i,
 \end{aligned}$$

drive output units to approximate relative frequencies in training data

$$o(\mathbf{x}; \boldsymbol{\theta})_i \approx \frac{\sum_{j=1}^m \mathbf{1}_{\mathbf{x}^{(j)} = \mathbf{x}, y_i^{(j)} = 1}}{\sum_{j=1}^m \mathbf{1}_{\mathbf{x}^{(j)} = \mathbf{x}}}$$

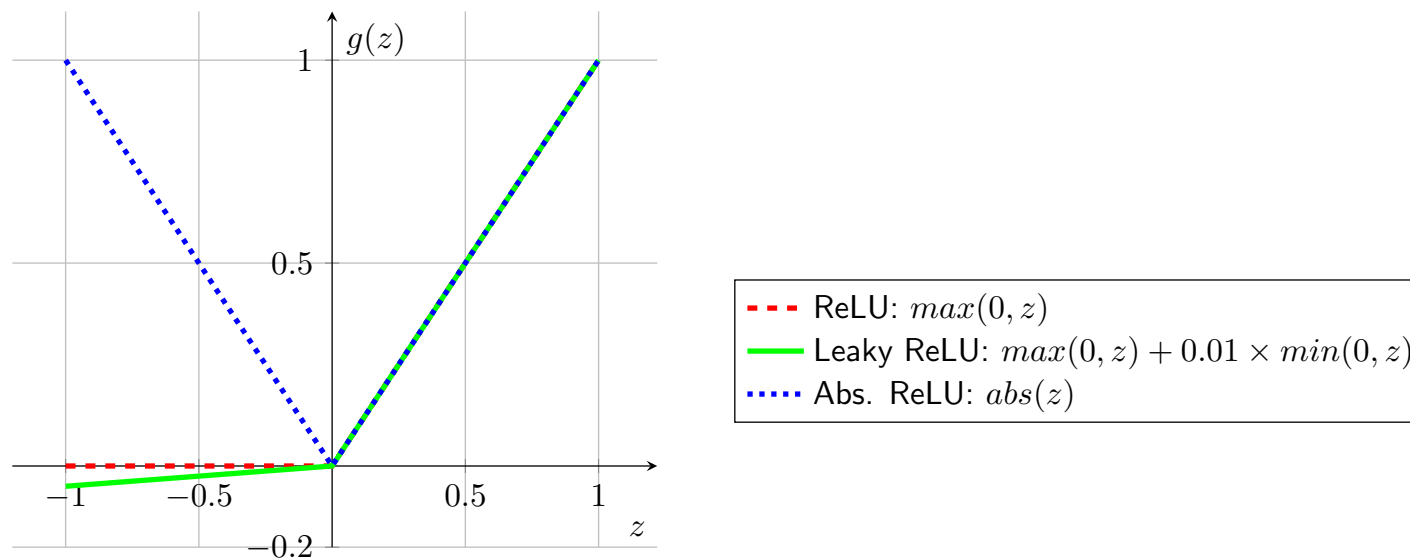
- To see this, solve the following constrained optimization problem:

$$\arg \min_{o(\mathbf{x}; \boldsymbol{\theta})_i} \left[- \sum_{j: \mathbf{x}^{(j)} = \mathbf{x}} \sum_{i=1}^n y_i^{(j)} \log o(\mathbf{x}; \boldsymbol{\theta})_i \right] \quad \text{s.t.} \quad \sum_{i=1}^n o(\mathbf{x}; \boldsymbol{\theta})_i = 1$$

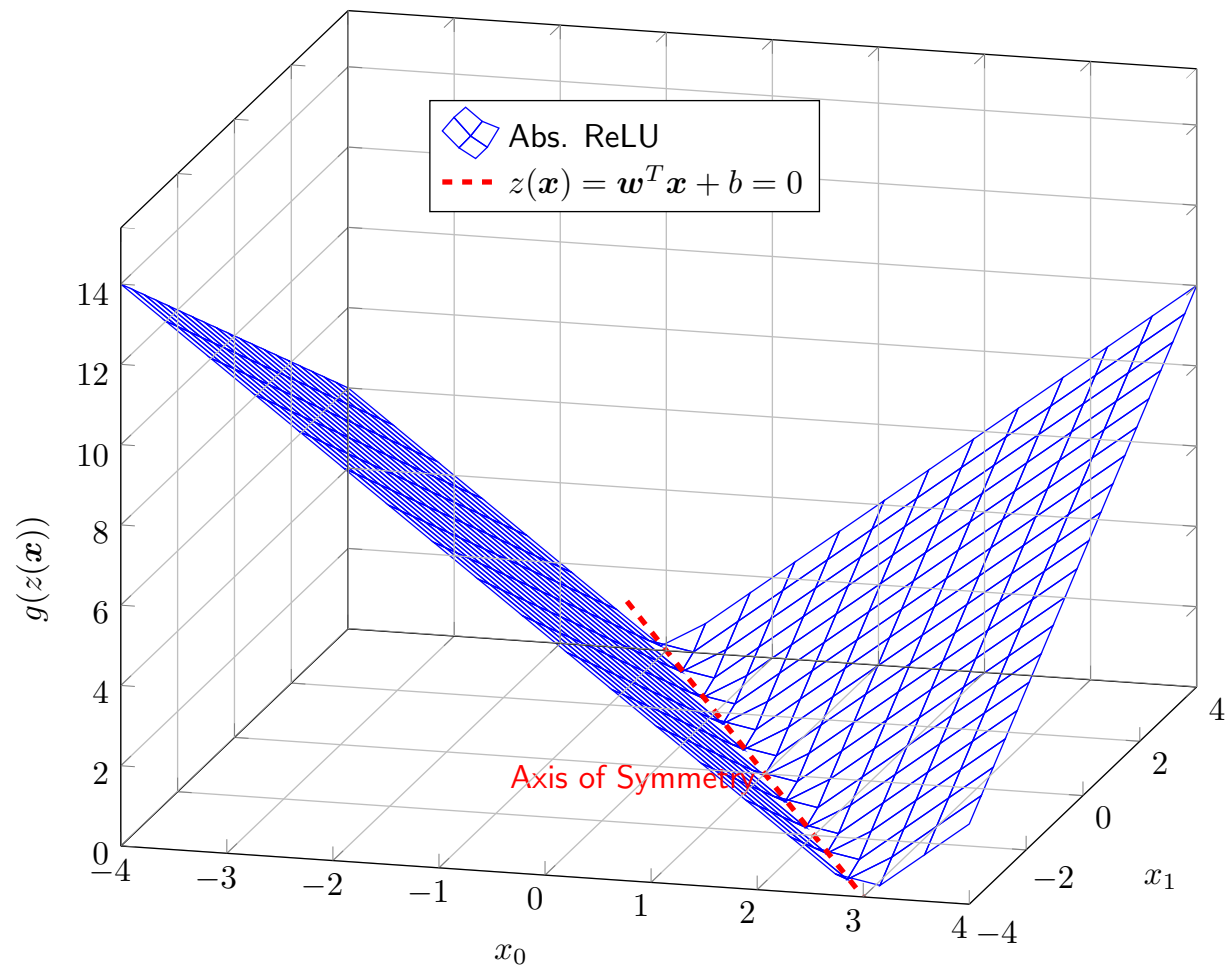
- Essentially, we are using the network to encode/approximate relative frequencies for different combinations of x, y
- How should we interpret the values of $o(x; \theta)$ if we input an x that is not included in training data?

Hidden Units and Activation Functions

- Computation of hidden units h
 1. Accepting inputs x
 2. Applying affine transformation $z = \mathbf{W}^T x + b$
 3. Applying element-wise non-linear mapping $h = g(z)$
- Typical $g(z)$, a.k.a. activation functions



Activation Functions in Input Space



Maxout Units

- Formulation

$$g(\mathbf{z})_i = \max_{j \in \mathbb{G}(i)} z_j$$

where

$$\mathbb{G}(i) = \{(i-1)k + 1, (i-1)k + 2, \dots, (i-1)k + k\}$$

- Maxout can implement many typical activation functions

- Given

$$z_1 = \mathbf{w}_1^T \mathbf{x} + b_1$$

$$z_2 = \mathbf{w}_2^T \mathbf{x} + b_2$$

- ReLU

$$g(z_1) = \max(z_1, z_2)$$

with $z_2 = 0$ by $(\mathbf{w}_2, b_2) = \mathbf{0}$

– Leaky ReLU

$$g(z_1) = \max(z_1, z_2)$$

with $z_2 = 0.01 \times z_1$ by $(\mathbf{w}_2, b_2) = 0.01 \times (\mathbf{w}_1, b_1)$

- Maxout can be seen as learning the activation function, in the sense that (\mathbf{w}_1, b_1) and (\mathbf{w}_2, b_2) can both be learned
- Maxout can approximate any piecewise linear, convex function in input space \mathbf{x} , when more pieces $\{z_i\}$ are input
- ReLU and all these variants exhibit linear behavior

Why Linearity?

- Toy problem: Training a network with only one hidden unit

$$\mathbf{x} \xrightarrow[\text{Affine}]{f(\mathbf{x})=\mathbf{w}^T \mathbf{x}+b} z \xrightarrow[\text{Activation}]{g(z)} h \xrightarrow[\text{Likelihood}]{c(h)} J$$

- Partial derivative of J w.r.t model parameter w_i is given by

$$\frac{\partial J}{\partial w_i} = \frac{\partial J}{\partial h} \times \underbrace{\frac{\partial h}{\partial z}}_{\text{Linearity}} \times \frac{\partial z}{\partial w_i} = c'(h) \times \underbrace{g'(z)}_{\text{Linearity}} \times x_i$$

- Linearity in activation (i.e. $\partial h / \partial z$ is some non-zero constant) helps ensure a well-behaved gradient w.r.t. model parameters

Universal Approximation Theorem

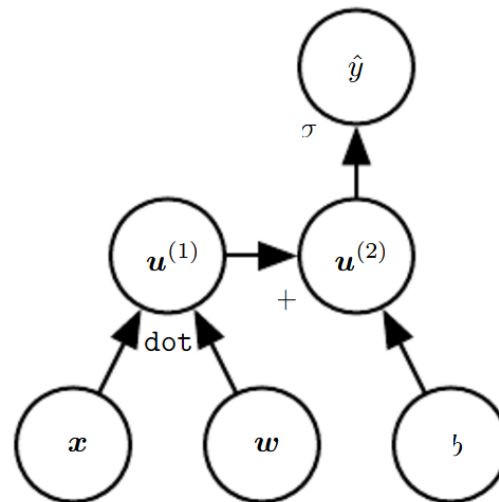
- (Hornik et al., 1989; Cybenko, 1989) A feedforward network with one linear output layer and at least one hidden layer with squashing activation function (e.g sigmoid) and enough hidden units can approximate any Borel measurable function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with any degree of accuracy
- Any continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ on a closed and bounded subset of \mathbb{R}^n is Borel measurable
- (Leshno et al., 1993) The theorem holds true for networks with rectified linear activation functions

How Many Are Enough?

- (Montufar et al., 2014) Functions representable with a deep rectifier net can require an exponential number of hidden units with a shallow network of one hidden layer (Check the paper for details)
- Empirically, deeper models can use fewer units to represent the desired function and generalize better

Computational Graphs

- To formalize computation as a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$
 - A node $v \in \mathcal{V}$ indicates a variable
 - A directed edge $e \in \mathcal{E}$ from x to y indicates that y is computed by applying an operation to x
- Example: Logistic regression prediction $\hat{y} = \mathbf{w}^T \mathbf{x} + b$



- Observe that some nodes represent model parameters w, b

Chain Rule of Calculus

- To compute the derivative dz/dx of a function $z(x)$ formed by the composition of functions $z(x) = f(g(x))$

$$x \xrightarrow{g(x)} y \xrightarrow{f(y)} z$$

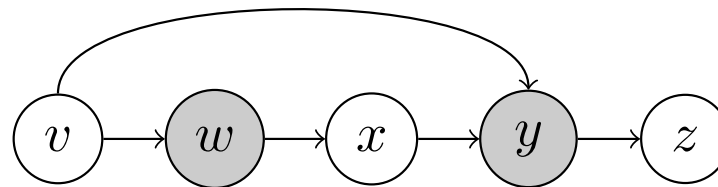
- $x \in \mathbb{R}$ is a real number
- $f, g : \mathbb{R} \rightarrow \mathbb{R}$ are real-valued functions of single variable
- The chain rule states that

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

- As an extension, we have for the following graph

$$\frac{dz}{dv} = \frac{dz}{dy} \frac{dy}{dv} + \frac{dz}{dw} \frac{dw}{dv} = \sum_{n:v \in Pa(n)} \frac{dz}{dn} \frac{dn}{dv}$$

where $Pa(n)$ is the set of nodes that are parents of n



- To verify the result requires another chain rule from calculus

$$z = f(y_1, y_2), \quad y_1 = g_1(x), \quad y_2 = g_2(x)$$

$$\frac{dz}{dx} = \frac{dz}{dy_1} \frac{dy_1}{dx} + \frac{dz}{dy_2} \frac{dy_2}{dx}$$

Vector Case

- z is a scalar, and \mathbf{x}, \mathbf{y} are vectors

$$\mathbf{x}_{m \times 1} \xrightarrow{g(\mathbf{x})} \mathbf{y}_{n \times 1} \xrightarrow{f(\mathbf{y})} z_{1 \times 1}$$

- Applying the chain rule leads to

$$\frac{\partial z}{\partial x_i} = \sum_{j=1}^n \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}, \quad i = 1, 2, \dots, m$$

- In matrix notation,

$$\begin{bmatrix} \frac{\partial z}{\partial x_1} \\ \frac{\partial z}{\partial x_2} \\ \vdots \\ \frac{\partial z}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \frac{\partial y_2}{\partial x_n} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} \frac{\partial z}{\partial y_1} \\ \frac{\partial z}{\partial y_2} \\ \vdots \\ \frac{\partial z}{\partial y_m} \end{bmatrix}$$

- This is recognized as a **Jacobian-gradient product**

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z,$$

where $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ (abbrev. $J_{\mathbf{y}, \mathbf{x}}$) is known as **Jacobian matrix** with

$$\left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{i,j} = \frac{\partial y_i}{\partial x_j}$$

As an example, when $\mathbf{y} = g(\mathbf{x}) = \mathbf{W}\mathbf{x}$ (i.e. $y_i = \sum_j w_{i,j} x_j$),

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{W}$$

Matrix Case

- z is a scalar, and \mathbf{X}, \mathbf{Y} are matrices

$$\mathbf{X}_{m \times n} \xrightarrow{g(\mathbf{X})} \mathbf{Y}_{s \times k} \xrightarrow{f(\mathbf{Y})} z_{1 \times 1}$$

- The chain rule suggests that

$$\frac{\partial z}{\partial x_{i,j}} = \sum_{s,k} \frac{\partial z}{\partial y_{s,k}} \frac{\partial y_{s,k}}{\partial x_{i,j}}, \quad \forall i, j$$

- More generally, when \mathbf{X}, \mathbf{Y} are tensors (high-dimensional arrays),

$$\nabla_{\mathbf{X}} z = \sum_j \left(\frac{\partial z}{\partial Y_j} \right) \nabla_{\mathbf{X}} Y_j,$$

where $\mathbf{Y} = g(\mathbf{X})$, $z = f(\mathbf{Y})$ and \mathbf{X} is treated as if it were a vector

• **Example 1:** Assume $\mathbf{Y} = g(\mathbf{X}) = \mathbf{W}\mathbf{X}$

- Let $\nabla_{\mathbf{Y}} z$ denote a matrix with its element (i, j) given by $\partial z / \partial Y_{i,j}$
- And $\nabla_{\mathbf{X}} z$ be a matrix with its element (i, j) given by $\partial z / \partial X_{i,j}$
- Observe that each **column** $\mathbf{Y}_{:,j}$ of \mathbf{Y} is a function of the **corresponding column** $\mathbf{X}_{:,j}$ in \mathbf{X} , i.e. $\mathbf{Y}_{:,j} = \mathbf{W}\mathbf{X}_{:,j}$
- We apply the Jacobian-gradient product of vector form to obtain

$$\nabla_{\mathbf{X}_{:,j}} z = \mathbf{W}^T \nabla_{\mathbf{Y}_{:,j}} z \Rightarrow \underbrace{\nabla_{\mathbf{X}} z = \mathbf{W}^T \nabla_{\mathbf{Y}} z}$$

• **Example 2:** Assume $\mathbf{Y} = g(\mathbf{X}) = \mathbf{X}\mathbf{W}$

- Observe that each **row** $\mathbf{Y}_{i,:}$ of \mathbf{Y} is a function of the **corresponding row** $\mathbf{X}_{i,:}$ in \mathbf{X} , i.e., $\mathbf{Y}_{i,:} = \mathbf{X}_{i,:}\mathbf{W}$, or equivalently, $\mathbf{Y}_{i,:}^T = \mathbf{W}^T \mathbf{X}_{i,:}^T$
- Applying the Jacobian-gradient product of vector form yields

$$\nabla_{\mathbf{X}_{i,:}^T} z = \mathbf{W} \nabla_{\mathbf{Y}_{i,:}^T} z \Rightarrow \underbrace{\nabla_{\mathbf{X}} z = (\nabla_{\mathbf{Y}} z) \mathbf{W}^T}$$

Back-Propagation

- Toy problem: To compute the derivative of J w.r.t. x

$$x \xrightarrow{f(x)} z \xrightarrow{g(z)} h \xrightarrow{c(h)} J$$

- From the chain rule, we have

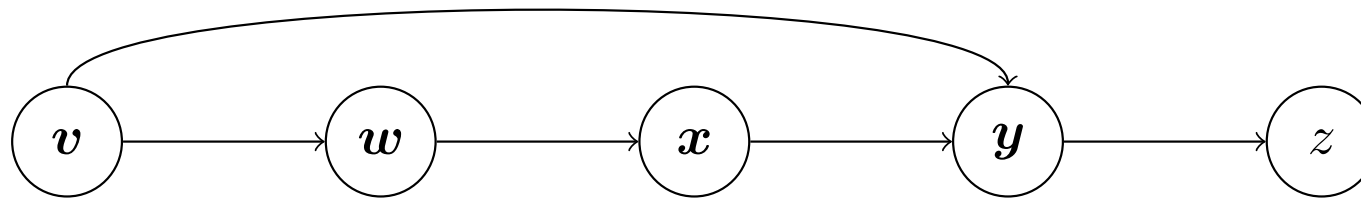
$$\begin{aligned}\frac{\partial J}{\partial x} &= \frac{\partial J}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial x} \\ &= c'(h)g'(z)f'(x) \\ &= c'(g(f(x)))g'(f(x))f'(x)\end{aligned}$$

- There are two possible implementations
 - The one based on the last equality incurs redundant subcomputation (e.g. $f(x)$)
 - The other following the penultimate equality requires z, h be pre-computed through forward propagation

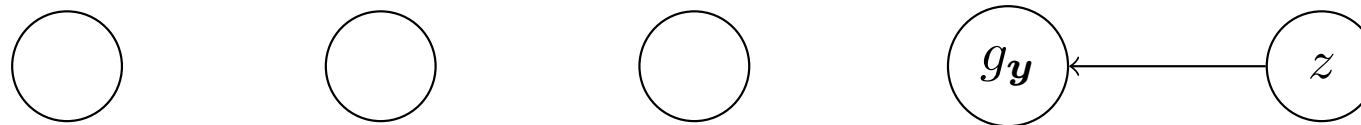
- Now, assuming z, h have been pre-computed, one way to compute the derivatives of J w.r.t. all variables x, z, h is to proceed in the order of
 1. $\partial J / \partial h = c'(h)$
 2. $\partial J / \partial z = (\partial J / \partial h) g'(z)$
 3. $\partial J / \partial x = (\partial J / \partial z) f'(x)$and, at each step, keep the result for subsequent use
- This technique is known as the **back-propagation (backprop)** method

General Back-Propagation

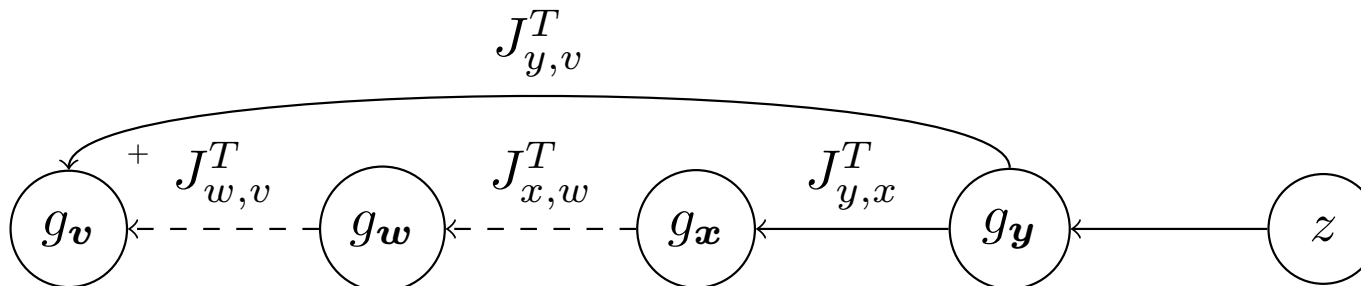
- To compute the gradient of z w.r.t. all its ancestors y, x, w, v



- Compute gradient w.r.t. every parent of z



- Travel backward, multiply current gradient by Jacobian recursively



- Sum gradients from different paths

- In symbols, we have

$$g_{\mathbf{y}} = \nabla_{\mathbf{y}} z$$

$$g_{\mathbf{x}} = \nabla_{\mathbf{x}} z = J_{\mathbf{y},\mathbf{x}}^T g_{\mathbf{y}}$$

$$g_{\mathbf{w}} = \nabla_{\mathbf{w}} z = J_{\mathbf{x},\mathbf{w}}^T g_{\mathbf{x}}$$

$$g_{\mathbf{v}} = \nabla_{\mathbf{v}} z = J_{\mathbf{w},\mathbf{v}}^T g_{\mathbf{w}} + J_{\mathbf{y},\mathbf{v}}^T g_{\mathbf{y}}$$

where

$$J_{\mathbf{y},\mathbf{x}} = \partial \mathbf{y} / \partial \mathbf{x}$$

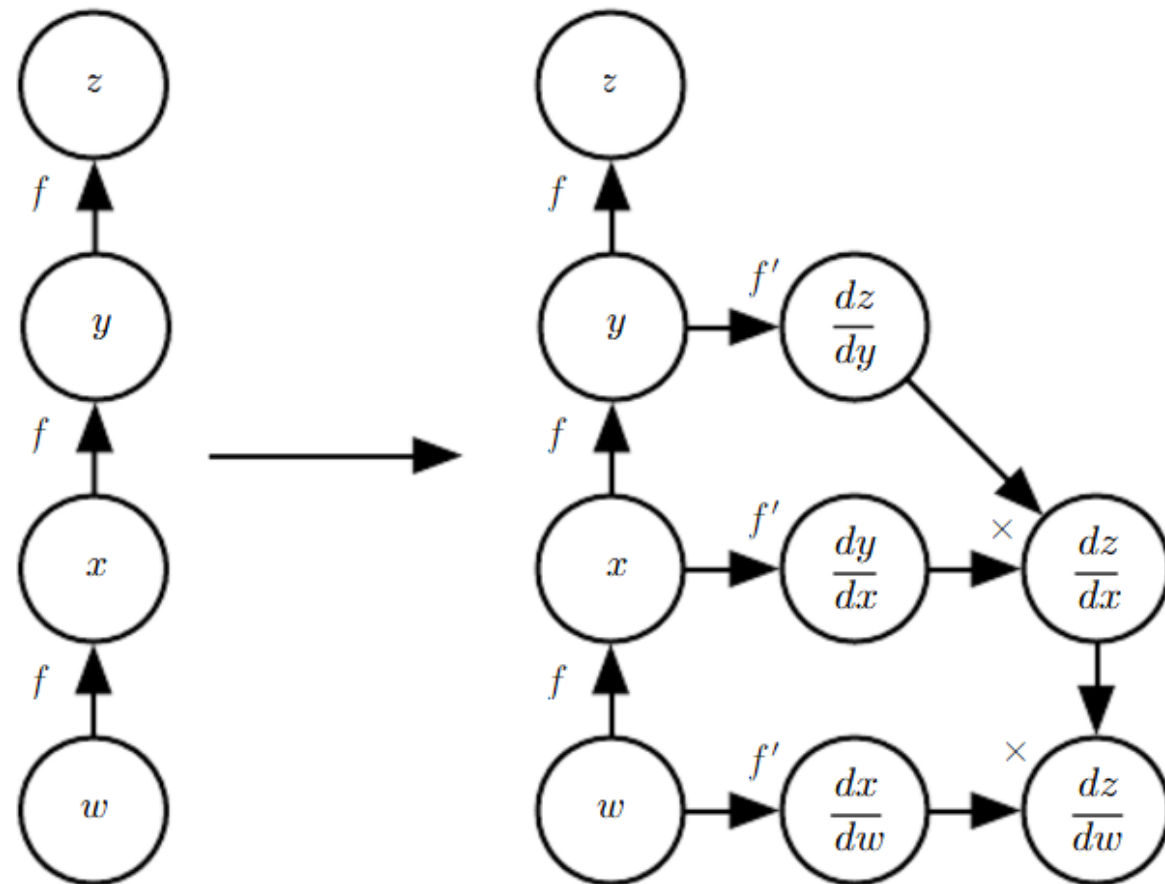
$$J_{\mathbf{x},\mathbf{w}} = \partial \mathbf{x} / \partial \mathbf{w}$$

$$J_{\mathbf{w},\mathbf{v}} = \partial \mathbf{w} / \partial \mathbf{v}$$

$$J_{\mathbf{y},\mathbf{v}} = \partial \mathbf{y} / \partial \mathbf{v}$$

- The spirit of this procedure can extend to cases where $\mathbf{y}, \mathbf{x}, \mathbf{w}, \mathbf{v}$ are matrices, tensors, vectors, scalars, or their mixing combinations

Backprop in Computational Graphs



Backprop for MLP Training

- Input \mathbf{X} in mini-batch form

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_0^T \\ \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_{m-1}^T \end{bmatrix}$$

- One layer of hidden features \mathbf{H} with ReLU

$$\mathbf{U}^{(1)} = \mathbf{X}\mathbf{W}^{(1)},$$

$$\mathbf{H} = \max\{0, \mathbf{U}^{(1)}\}$$

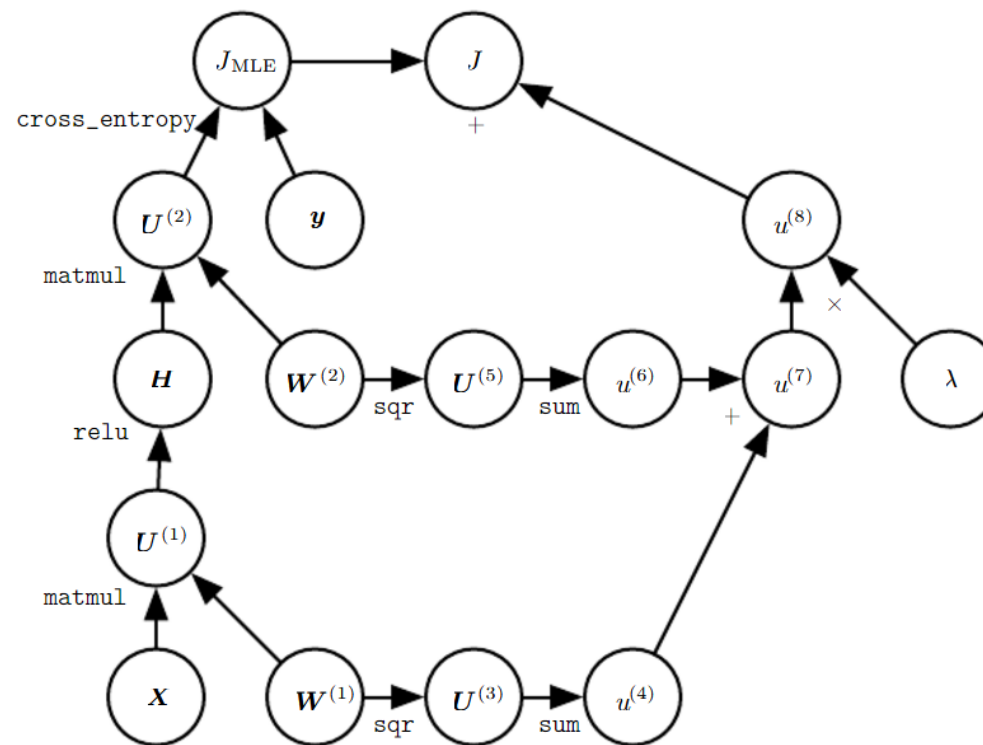
- One layer of outputs $\mathbf{U}^{(2)}$ (before normalization)

$$\mathbf{U}^{(2)} = \mathbf{H}\mathbf{W}^{(2)}$$

- Objective: To minimize the cross-entropy with weight decay

$$J = J_{MLE} + \lambda \left(\sum_{i,j} (W_{i,j}^{(1)})^2 + \sum_{i,j} (W_{i,j}^{(2)})^2 \right)$$

- Computational graph



- Backprop: To compute $\nabla_{\mathbf{W}^{(1)}} J$ and $\nabla_{\mathbf{W}^{(2)}} J$
 - Two paths from J to the weights (only one path illustrated)
 - Assume $\nabla_{\mathbf{U}^{(2)}} J = \mathbf{G}$
 - Then $\nabla_{\mathbf{W}^{(2)}} J = \mathbf{H}^T \mathbf{G}$ (cf. **Example 1** in Matrix Case)
 - Similarly, $\nabla_{\mathbf{H}} J = \mathbf{G} \mathbf{W}^{(2)T}$ (cf. **Example 2** in Matrix Case)
 - Tracing back further, we have $\nabla_{\mathbf{U}^{(1)}} J = \mathbf{G}'$ by zeroing out elements in $\nabla_{\mathbf{H}} J$ corresponding to entries of $\mathbf{U}^{(1)}$ less than zero
 - Again, $\nabla_{\mathbf{W}^{(1)}} J = \mathbf{X}^T \mathbf{G}'$ (cf. **Example 1** in Matrix Case)

Review

- Maximum likelihood as cost functions
- Output units (linear/sigmoid/softmax) as model parameters
- Cost functions and output units: gradient vanishing issues
- Hidden units and activation functions
- Universal approximation theorem: depth vs. width
- Chain rules of calculus
- Back propagation