## CSW182（2021）· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

视频·B 站 [ 扫码或点击链接 ]
https://www.bilibili.com/video/BV1Ff4y1n7ar

课件 & 代码·博客 [ 扫码或点击链接 ]
http://blog.showmeai.tech/berkeley-csw182

Berkeley

循环神经网络

可视化

梯度策略

Q-Learning

风格迁移

模仿学习

元学习

计算机视觉

机器学习基础

生成模型

卷积网络

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS231n |

### # Awesome AI Courses Notes Cheatsheets· 持续更新中

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]

# Backpropagation

Designing, Visualizing and Understanding Deep Neural Networks
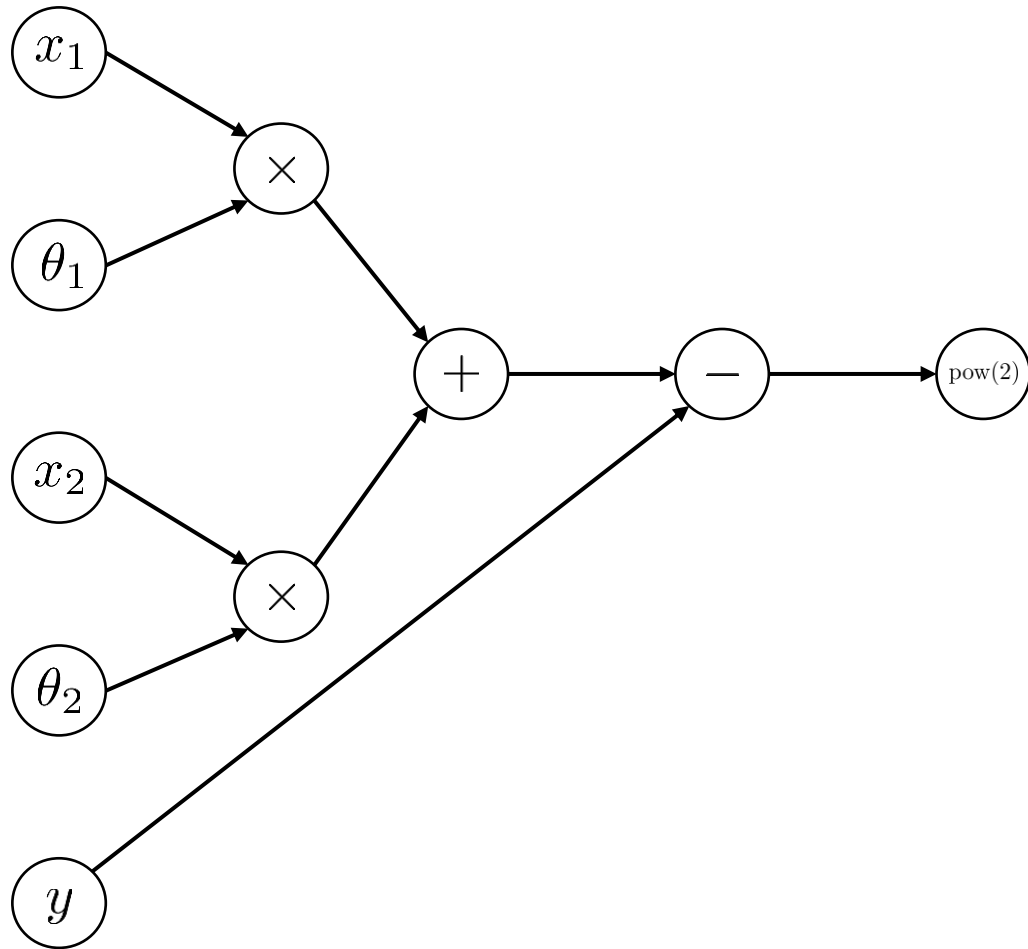
# CS W182/282A

Instructor: Sergey Levine
UC Berkeley

# Neural networks

# Drawing computation graphs



what **expression** does this compute?

equivalently, what **program** does this correspond to?
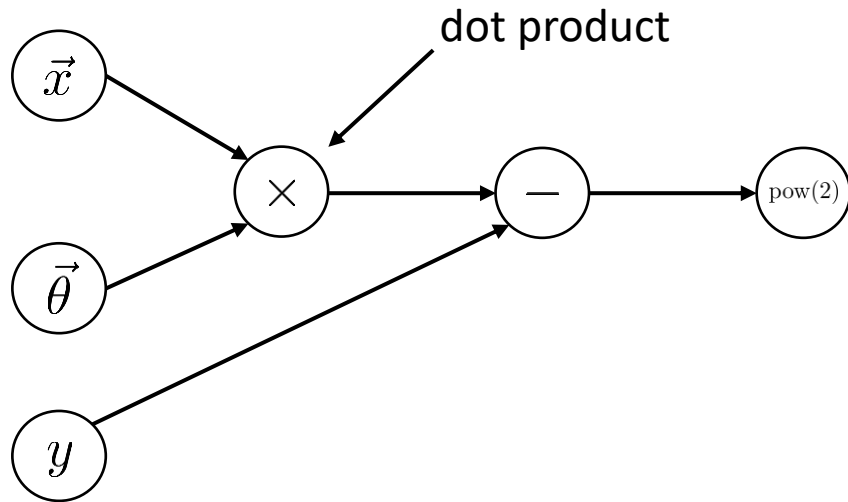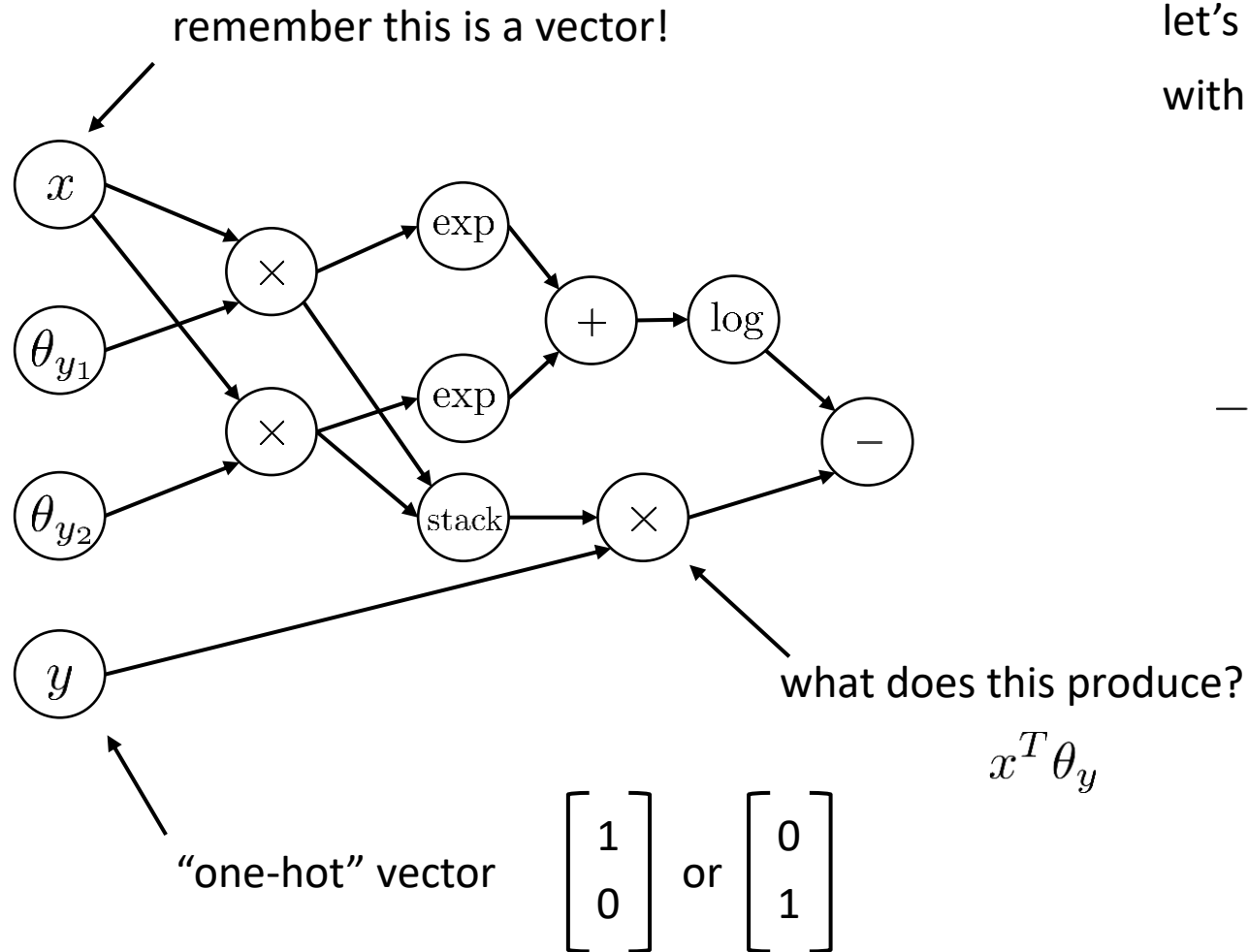
$$||(x_1\theta_1 + x_2\theta_2) - y||^2$$

this is a **MSE loss** with a **linear regression** model

**neural networks** are **computation graphs**

if we design **generic tools** for computation graphs, we can train **many kinds** of neural networks

# Drawing computation graphs

a simpler way to draw the same thing:

what **expression** does this compute?

equivalently, what **program** does this correspond to?

$$||(x_1\theta_1 + x_2\theta_2) - y||^2$$

this is a **MSE loss** with a **linear regression** model

dot product

$\vec{x}$

$\vec{\theta}$

$y$

$\times$

$-$

pow(2)

**neural networks** are **computation graphs**

if we design **generic tools** for computation graphs, we can train **many kinds** of neural networks

I'll drop the $\vec{\phantom{x}}$ decorator from now on...

# Logistic regression

remember this is a vector!

let's draw the computation graph for **logistic regression**

with the negative log-likelihood loss

$$p_\theta(y|x) = \frac{\exp(x^T\theta_y)}{\sum_{y'}\exp(x^T\theta_{y'})}$$

$$-\log p_\theta(y|x) = -x^T\theta_y + \log\sum_{y'}\exp(x^T\theta_{y'})$$

what does this produce?

$$x^T\theta_y$$

"one-hot" vector $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ or $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$

# Logistic regression

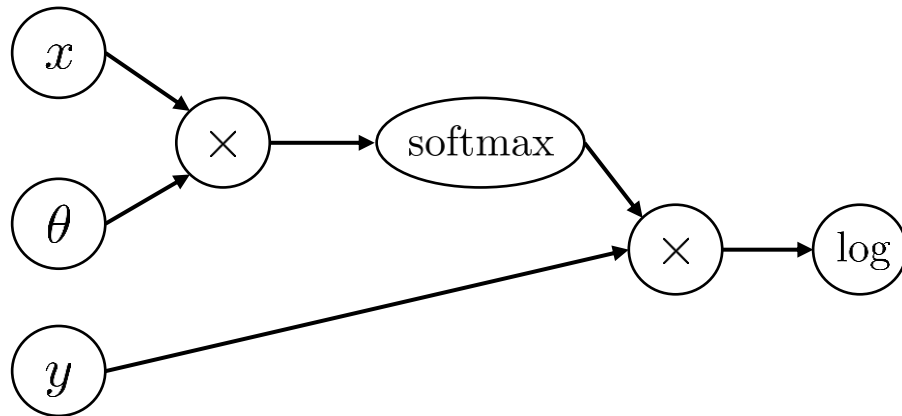$$p_\theta(y|x) = \frac{\exp(x^T\theta_y)}{\sum_{y'}\exp(x^T\theta_{y'})}$$

a simpler way to draw the same thing:

$$-\log p_\theta(y|x) = -x^T\theta_y + \log\sum_{y'}\exp(x^T\theta_{y'})$$

$$f_\theta(x) = \begin{bmatrix} x^T\theta_{y_1} \\ x^T\theta_{y_2} \\ \cdots \\ x^T\theta_{y_m} \end{bmatrix} \qquad f_\theta(x) = \theta x$$

matrix

$$\begin{bmatrix} \theta_{y_1} \\ \theta_{y_2} \\ \theta_{y_3} \end{bmatrix} \times \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} x^T\theta_{y_1} \\ x^T\theta_{y_2} \\ \cdots \\ x^T\theta_{y_m} \end{bmatrix}$$

$$p_\theta(y = i|x) = \text{softmax}(f_\theta(x))[i] = \frac{\exp(f_{\theta,i}(x))}{\sum_{j=1}^m \exp(f_{\theta,j}(x))}$$
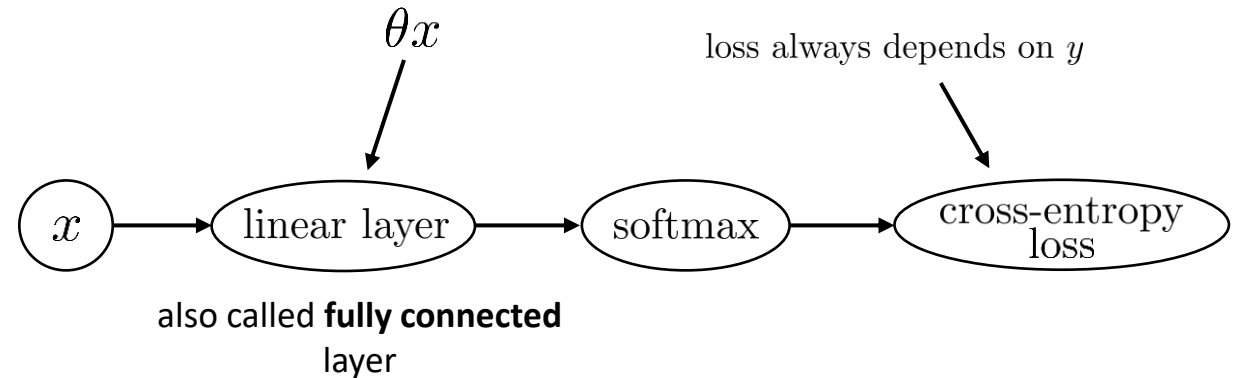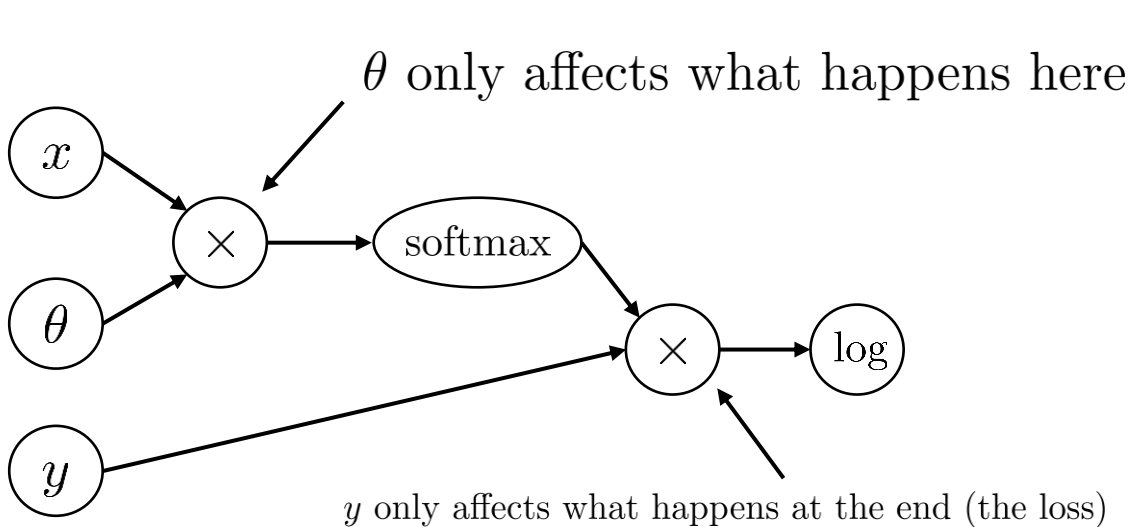
# Drawing it even *more* concisely

Notice that we have **two types** of variables:

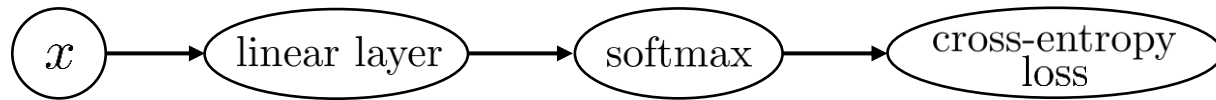data (e.g., $x, y$), which serves as input or target output

parameters (e.g., $\theta$)          the parameters *usually* affect one specific operation

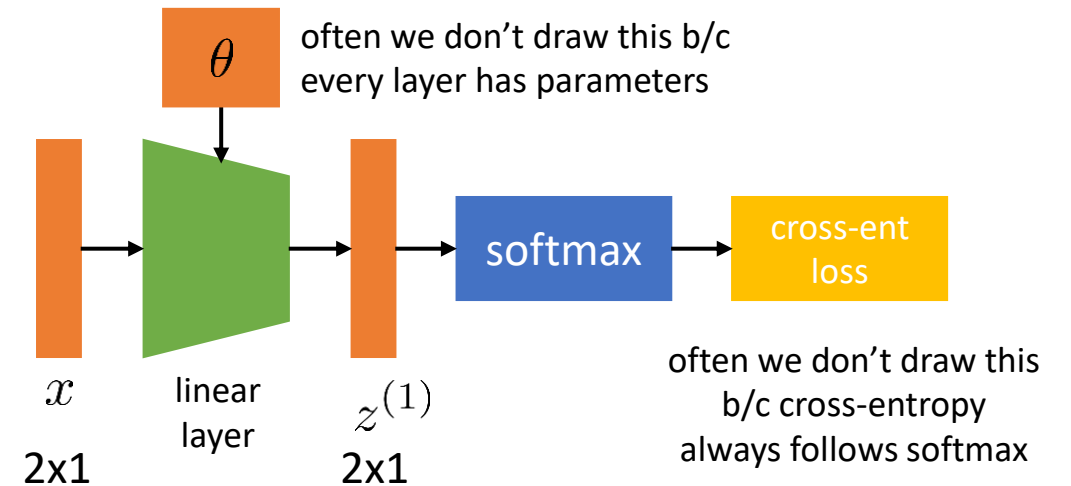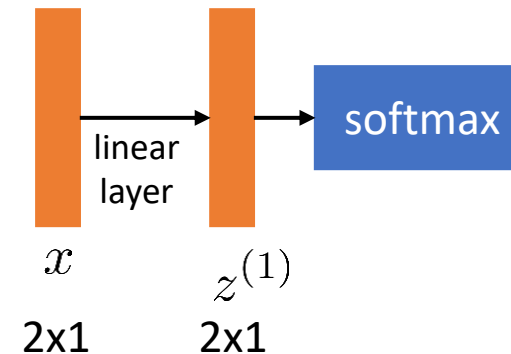(though there is often *parameter sharing*, e.g., conv nets – more on this later)

$\theta$ only affects what happens here

$x$

$\times$ → softmax

$\theta$

$y$ → $\times$ → log

*y only affects what happens at the end (the loss)*

$\theta x$

loss always depends on $y$

$x$ → linear layer → softmax → cross-entropy loss

also called **fully connected** layer

# Neural network diagrams

(simplified) computation graph diagram

neural network diagram



often we don't draw this b/c every layer has parameters

often we don't draw this b/c cross-entropy always follows softmax

simplified drawing:

# Logistic regression with features

pop quiz: what is the dimensionality of $\theta$?

$$\text{softmax}(x^T\theta)$$

$$\phi(x) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{pmatrix}$$

$$\text{softmax}(\phi(x)^T\theta)$$

$x_1$

$x_2$

$x_1{}^2$

$x_2{}^2$

$x_1x_2$

# Learning the features

which layer

$w_1^{(1)}$

which feature
= rows of weight **matrix**

**Problem:** how do we represent the learned features?

**Idea:** what if each feature is a (binary) logistic regression output?

$$\phi_1(x) = \mathrm{softmax}(x^T w_1^{(1)}) = \frac{1}{1 + \exp(-x^T w_1^{(1)})}$$

$$W^{(1)} \begin{array}{|c|} \hline w_1^{(1)} \\ \hline w_2^{(1)} \\ \hline w_3^{(1)} \\ \hline \end{array}$$

aside: I'll switch to use $w$ or $W$ instead of $\theta$ here
$\theta - all$ parameters of the model
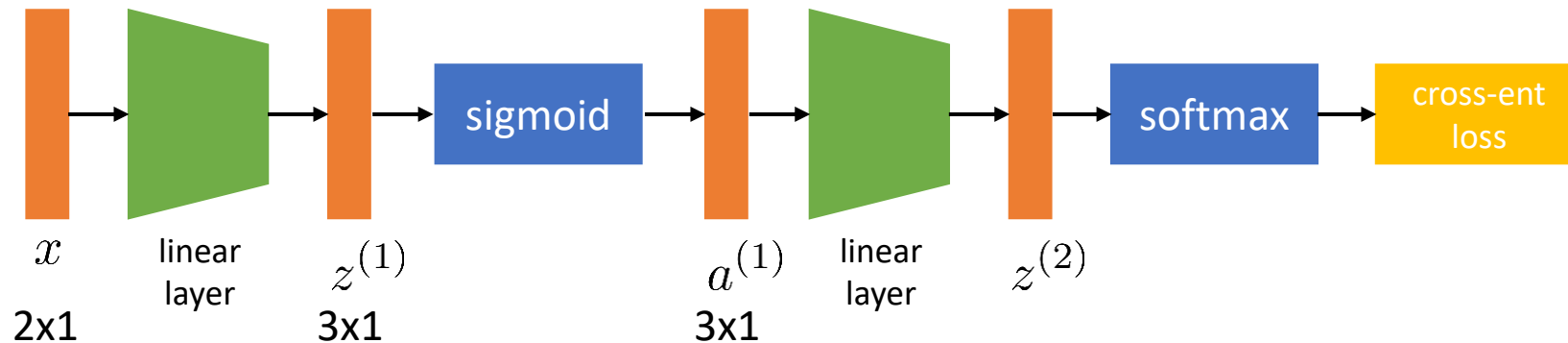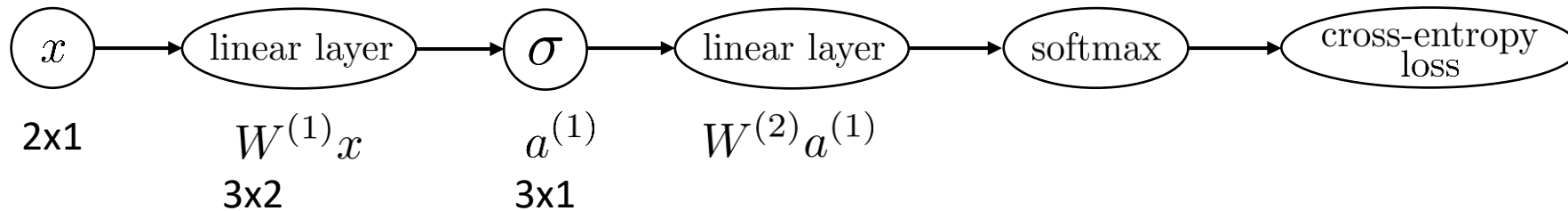$w_1^{(1)} -$ weights (a.k.a. parameters) of feature 1 at layer 1

$$\phi(x) = \begin{pmatrix} \mathrm{softmax}(x^T w_1^{(1)}) \\ \mathrm{softmax}(x^T w_2^{(1)}) \\ \mathrm{softmax}(x^T w_3^{(1)}) \end{pmatrix} = \sigma(W^{(1)} x)$$

**per-element** sigmoid
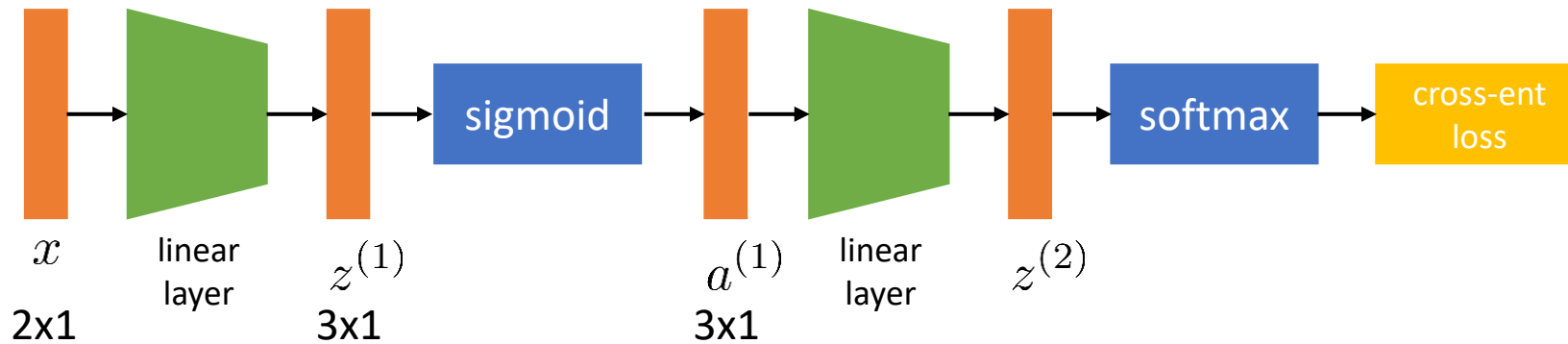
**not** the same as softmax
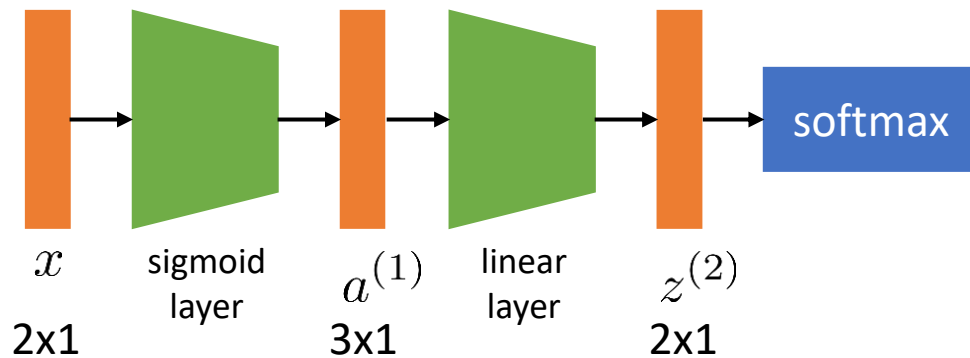
each feature is independent

# Let's draw this!

$$\phi(x) = \begin{pmatrix} \mathrm{softmax}(x^T w_1^{(1)}) \\ \mathrm{softmax}(x^T w_2^{(1)}) \\ \mathrm{softmax}(x^T w_3^{(1)}) \end{pmatrix} = \sigma(W^{(1)}x) \qquad p(y|x) = \mathrm{softmax}(\phi(x)^T \theta)$$

$x$ → linear layer → $\sigma$ → linear layer → softmax → cross-entropy loss

2x1     $W^{(1)}x$     $a^{(1)}$     $W^{(2)}a^{(1)}$

3x2     3x1

$x$    linear layer    $z^{(1)}$    sigmoid    $a^{(1)}$    linear layer    $z^{(2)}$    softmax    cross-ent loss

2x1     3x1     3x1

# Simpler drawing



$x$    linear layer    $z^{(1)}$    sigmoid    $a^{(1)}$    linear layer    $z^{(2)}$    softmax    cross-ent loss

2x1    3x1    3x1

simpler way to draw the same thing:

$x$    sigmoid layer    $a^{(1)}$    linear layer    $z^{(2)}$    softmax

2x1    3x1    2x1

even simpler:

$x$    sigmoid layer    $a^{(1)}$    linear layer    $z^{(2)}$    softmax

2x1    3x1    2x1

# Doing it multiple times



$x$ → linear layer → $\sigma$ → linear layer → $\sigma$ → linear layer → $\sigma$ → linear layer → softmax

2x1    $W^{(1)}x$    $a^{(1)}$    $W^{(2)}a^{(1)}$    $a^{(2)}$    $W^{(3)}a^{(2)}$    $a^{(3)}$    $W^{(4)}a^{(3)}$

3x2    3x1    3x3    3x1    3x3    3x1

sigmoid layer → sigmoid layer → sigmoid layer → linear layer → softmax

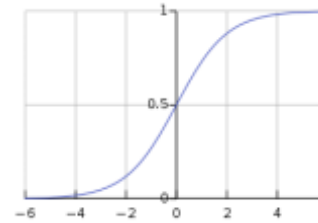$x$    $a^{(1)}$    $a^{(2)}$    $a^{(3)}$    $z^{(4)}$

2x1    3x1    3x1    3x1    2x1

# Activation functions

$$\phi_1(x) = \mathrm{softmax}(x^T w_1^{(1)}) = \frac{1}{1 + \exp(-x^T w_1^{(1)})}$$



we don't have to use a **sigmoid**!

a wide range of non-linear functions will work

these are called **activation functions**

we'll discuss specific choices later
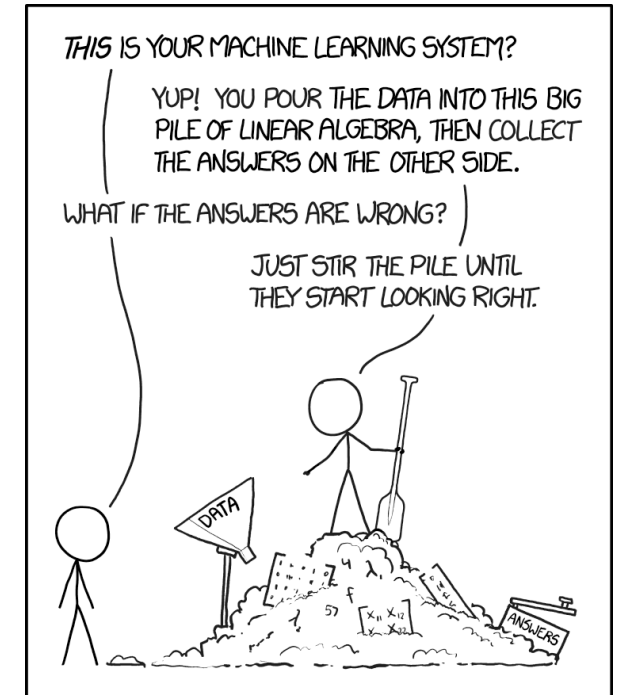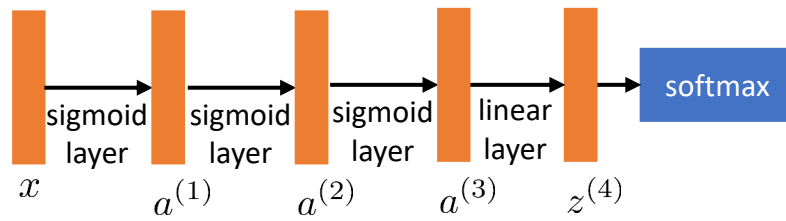
why **non-linear?**

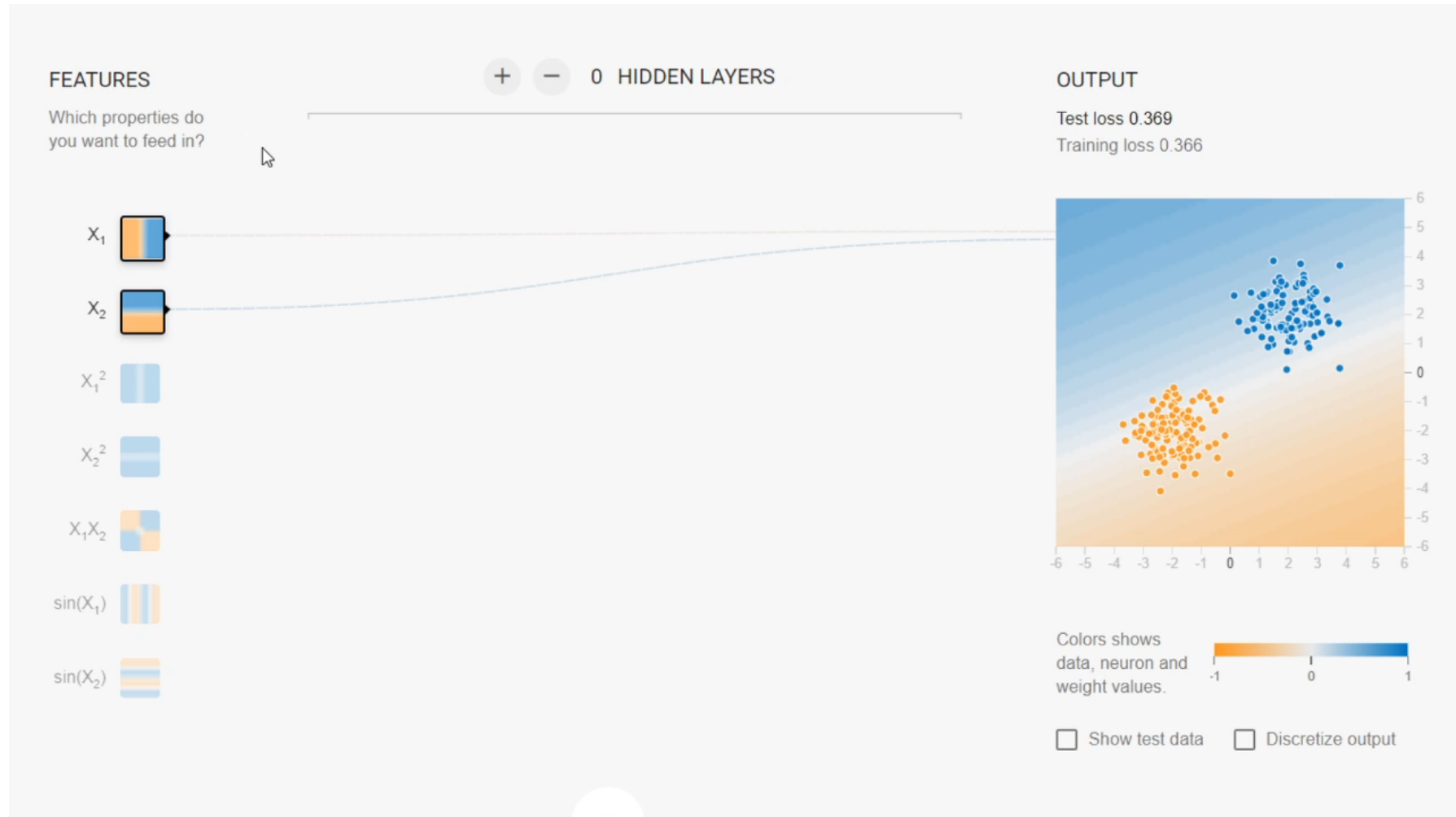$$a^{(2)} = \sigma(W^{(2)} \sigma(W^{(1)} x))$$

if $\sigma(z) = z$, then...

$$a^{(2)} = W^{(2)} W^{(1)} x = Mx$$

multiple linear layers = one linear layer

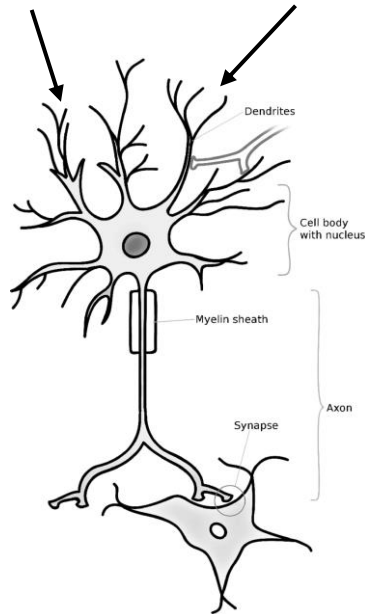enough layers = we can represent anything (so long as they're nonlinear)



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

# Demo time!



Source: https://playground.tensorflow.org/
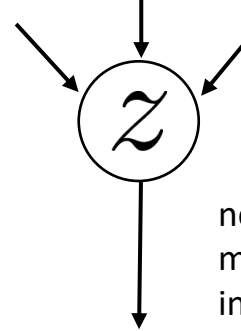
# Aside: what's so neural about it?

dendrites receive signals from other neurons

neuron "decides" whether to fire based on incoming signals

axon transmits signal to downstream neurons

artificial "neuron" sums up signals from upstream neurons (also referred to as "units")

$$z = \sum_i a_i$$

upstream activations

neuron "decides" how much to fire based on incoming signals

$$a = \sigma(z)$$

activations transmitted to downstream units
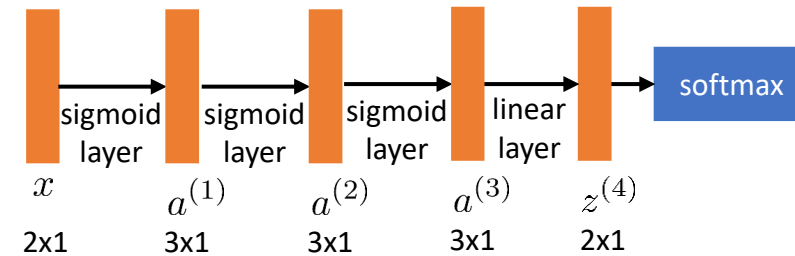
activation function

# Training neural networks

# What do we need?

1. Define your **model class**



2. Define your **loss function**

negative log-likelihood, just like before

3. Pick your **optimizer**

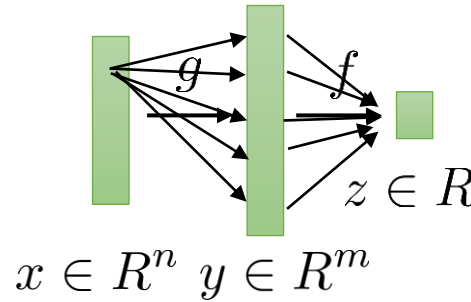stochastic gradient descent

what do we need?

$$\nabla_\theta \mathcal{L}(\theta) = \begin{pmatrix} \dfrac{d\mathcal{L}(\theta)}{d\theta_1} \\ \dfrac{d\mathcal{L}(\theta)}{d\theta_2} \\ \vdots \\ \dfrac{d\mathcal{L}(\theta)}{d\theta_n} \end{pmatrix}$$

4. Run it on a big GPU

# Aside: chain rule

Chain rule: $\quad x \xrightarrow{g} y \xrightarrow{f} z$

$$\frac{d}{dx}f(g(x)) = \frac{dz}{dx} = \frac{dy}{dx}\frac{dz}{dy}$$

Jacobian of $g$  Jacobian of $f$

$x \in R^n \ y \in R^m$

$z \in R$

## High-dimensional chain rule

$$\frac{d}{dx_i}f(g(x)) = \sum_{j=1}^{m}\frac{dy_j}{dx_i}\frac{dz}{dy_j} = \frac{dy}{dx_i}\frac{dz}{dy}$$

sum over all dimensions of $y$

row $1 \times m$   col $m \times 1$

$$\frac{d}{dx}f(g(x)) = \frac{dy}{dx}\frac{dz}{dy}$$

mat $n \times m$   col $m \times 1$

Row or column?

In this lecture:

In some textbooks:

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

$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$

$\left(\frac{dy}{dx}\right)_{ij} = \frac{dy_j}{dx_i}$

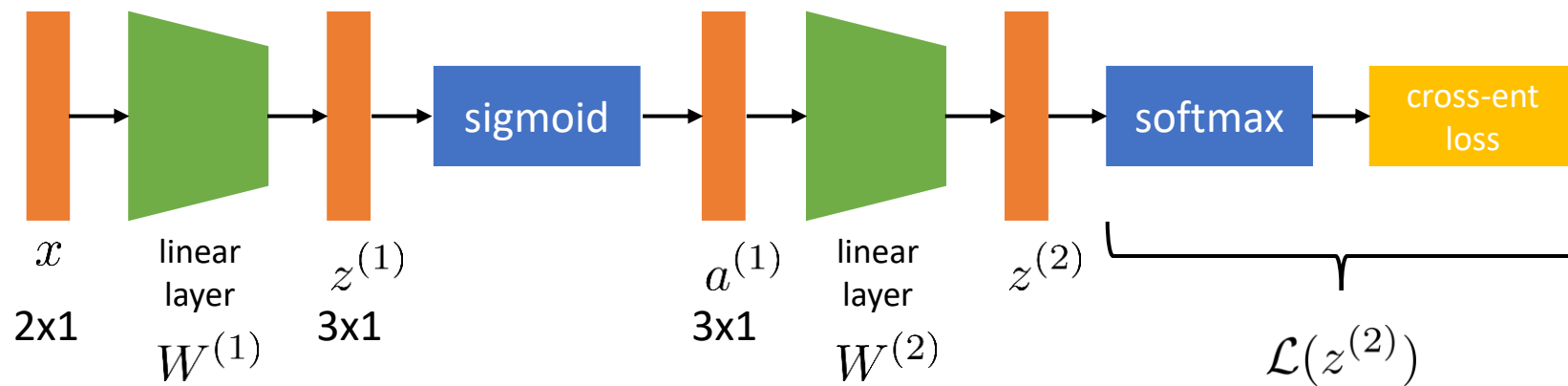$y \in R^m \qquad \frac{dz}{dy} \in R^m$

Just two different conventions!

# Chain rule for neural networks

A neural network is just a composition of functions

So we can use chain rule to compute gradients!



$x$    linear layer    $z^{(1)}$    $a^{(1)}$    linear layer    $z^{(2)}$

2x1    $W^{(1)}$    3x1    3x1    $W^{(2)}$    $\mathcal{L}(z^{(2)})$

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \frac{da^{(1)}}{dz^{(1)}} \frac{dz^{(2)}}{da^{(1)}} \frac{d\mathcal{L}}{dz^{(2)}} \qquad\qquad \frac{d\mathcal{L}}{dW^{(2)}} = \frac{dz^{(2)}}{dW^{(2)}} \frac{d\mathcal{L}}{dz^{(2)}}$$

# Does it work?

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \frac{da^{(1)}}{dz^{(1)}} \frac{dz^{(2)}}{da^{(1)}} \frac{d\mathcal{L}}{dz^{(2)}}$$

We **can** calculate each of these Jacobians!

Example:

$$z^{(2)} = W^{(2)} a^{(1)}$$

$$\frac{dz^{(2)}}{da^{(1)}} = W^{(2)T}$$

Why might this be a **bad** idea?

if each $z^{(i)}$ or $a^{(i)}$ has about $n$ dims...

each Jacobian is about $n \times n$ dimensions

matrix multiplication is $O(n^3)$

do we care?

AlexNet has layers with 4096 units...

# Doing it more efficiently

this product is cheap: $O(n^2)$

this product is expensive

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \frac{da^{(1)}}{dz^{(1)}} \frac{dz^{(2)}}{da^{(1)}} \frac{d\mathcal{L}}{dz^{(2)}}$$

$n \times n$        $n \times 1$

this is **always** true because
the loss is scalar-valued!

**Idea:** start on the right

compute $\frac{dz^{(2)}}{da^{(1)}} \frac{d\mathcal{L}}{dz^{(2)}} = \delta$ first

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \frac{da^{(1)}}{dz^{(1)}} \delta$$
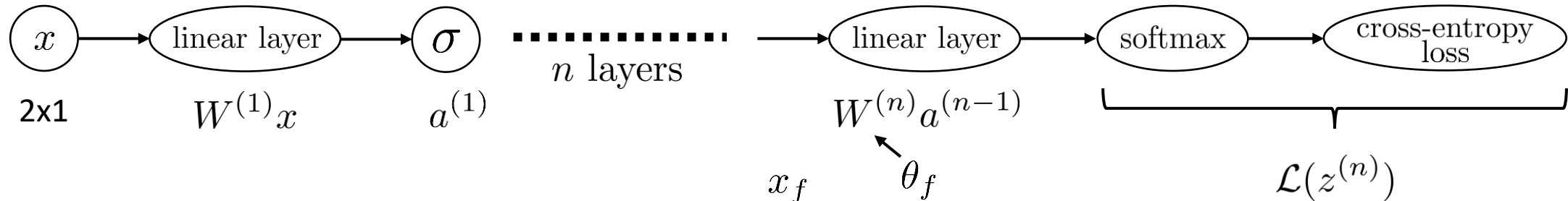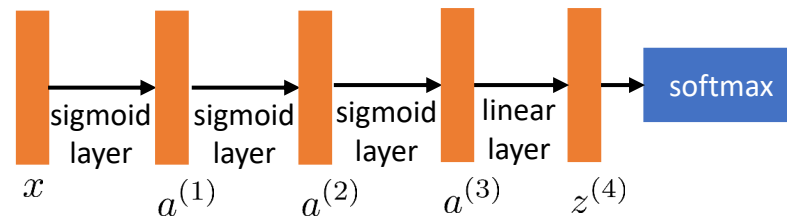
this product is cheap: $O(n^2)$

compute $\frac{da^{(1)}}{dz^{(1)}} \delta = \gamma$

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \gamma$$

this product is cheap: $O(n^2)$

# The backpropagation algorithm

"Classic" version



forward pass: calculate each $a^{(i)}$ and $z^{(i)}$ $\quad a^{(n-1)} \longrightarrow f \longrightarrow z^{(n-1)}$

backward pass:

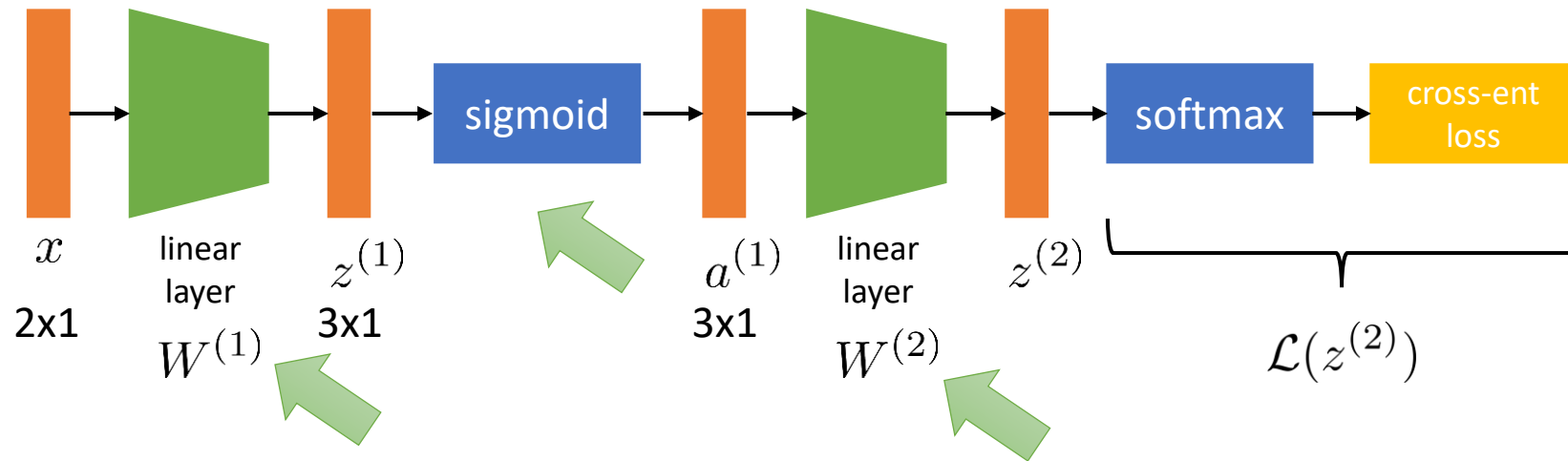initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$

for each $f$ with input $x_f$ & params $\theta_f$ from end to start:

$$\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f}\delta$$

$$\delta \leftarrow \frac{df}{dx_f}\delta$$

# Let's walk through it...



$$\frac{d\mathcal{L}}{dW^{(2)}} = \frac{dz^{(2)}}{dW^{(2)}} \underbrace{\frac{d\mathcal{L}}{dz^{(2)}}}_{\delta}$$

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \underbrace{\frac{da^{(1)}}{dz^{(1)}} \frac{dz^{(2)}}{da^{(1)}} \frac{d\mathcal{L}}{dz^{(2)}}}_{\delta}$$

forward pass: calculate each $a^{(i)}$ and $z^{(i)}$

backward pass:
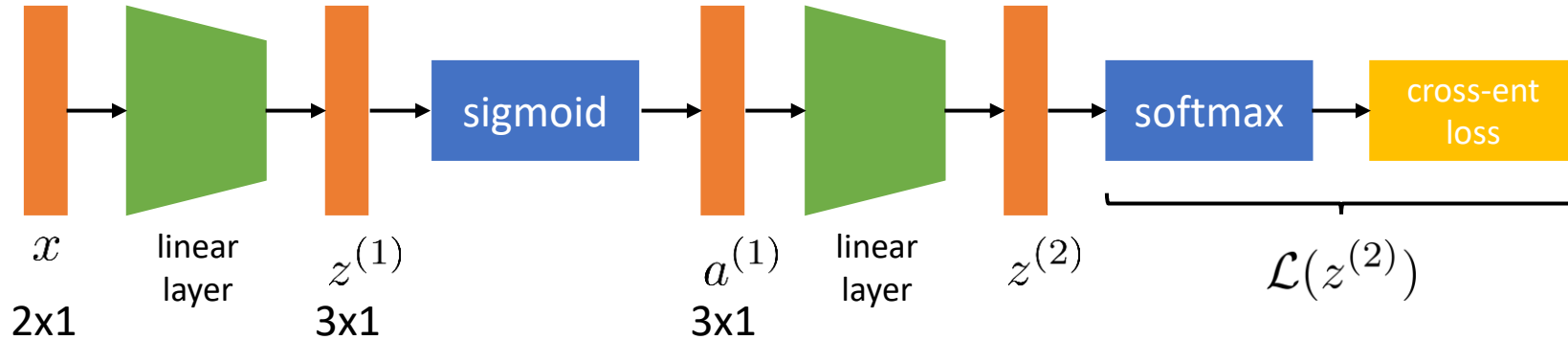
initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$

for each $f$ with input $x_f$ & params $\theta_f$ from end to start:

$$\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f} \delta$$

$$\delta \leftarrow \frac{df}{dx_f} \delta$$

# Practical implementation

# Neural network architecture details



$x$

linear layer

$z^{(1)}$

sigmoid

$a^{(1)}$

linear layer

$z^{(2)}$

softmax

cross-ent loss

$\mathcal{L}(z^{(2)})$

2x1

3x1

3x1

Some things we should figure out:
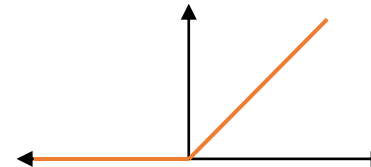
How many layers?

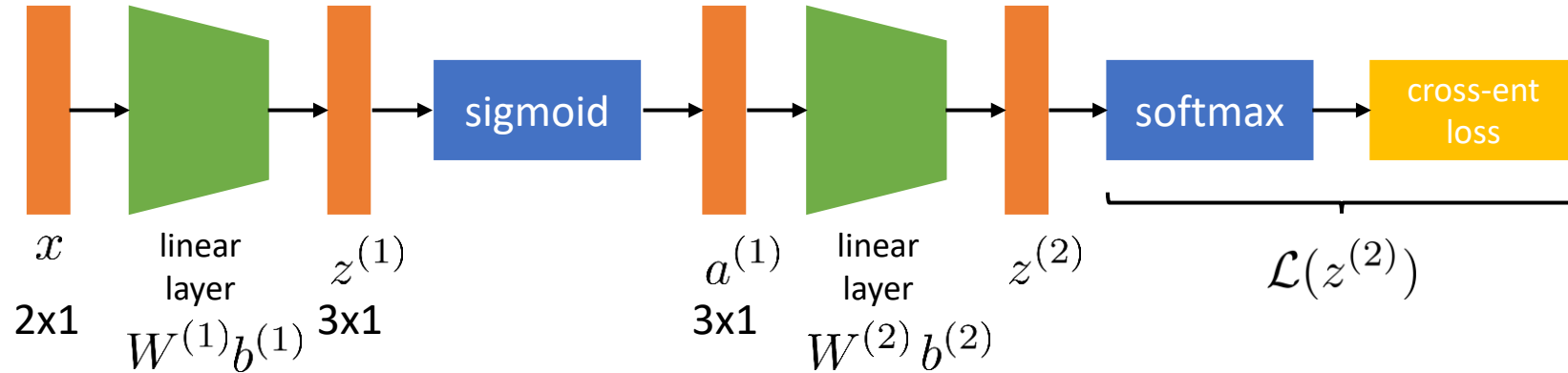How big are the layers?

What type of **activation function**?

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

$$\mathrm{ReLU}(x) = \max(0, x)$$
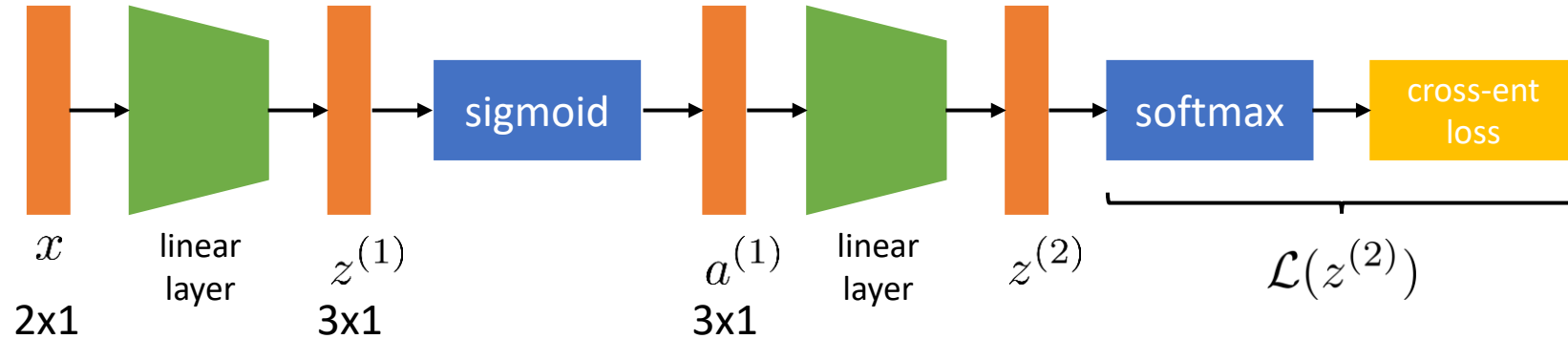
# Bias terms



Linear layer:

$$z^{(i+1)} = W^{(i)} a^{(i)}$$

problem: if $a^{(i)} = \vec{0}$, we always get 0...

Solution: add a "bias":

has nothing to do with bias/variance bias

$$z^{(i+1)} = W^{(i)} a^{(i)} + b^{(i)}$$

additional parameters in each linear layer

# What else do we need for backprop?



forward pass: calculate each $a^{(i)}$ and $z^{(i)}$

backward pass:
initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$

for each $f$ with input $x_f$ & params $\theta_f$ from end to start:
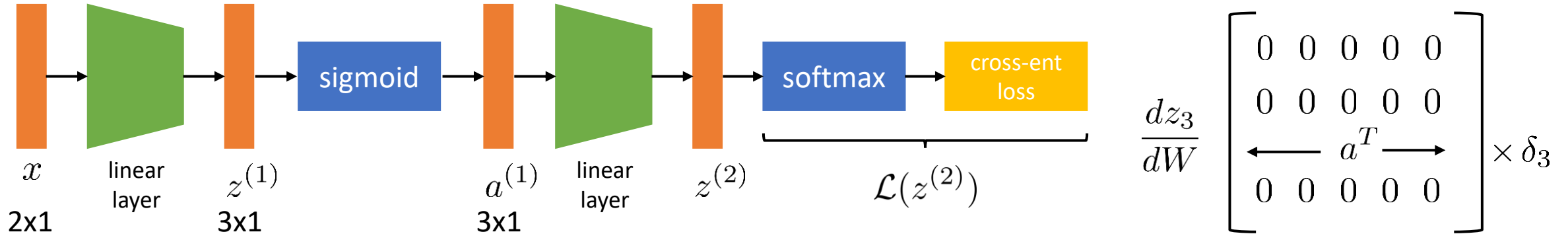
$$\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f}\delta$$

$$\delta \leftarrow \frac{df}{dx_f}\delta$$

for each function, we need to compute:

$$\frac{df}{d\theta_f}\delta \qquad \frac{df}{dx_f}\delta$$

linear layer

softmax + cross-entropy

sigmoid
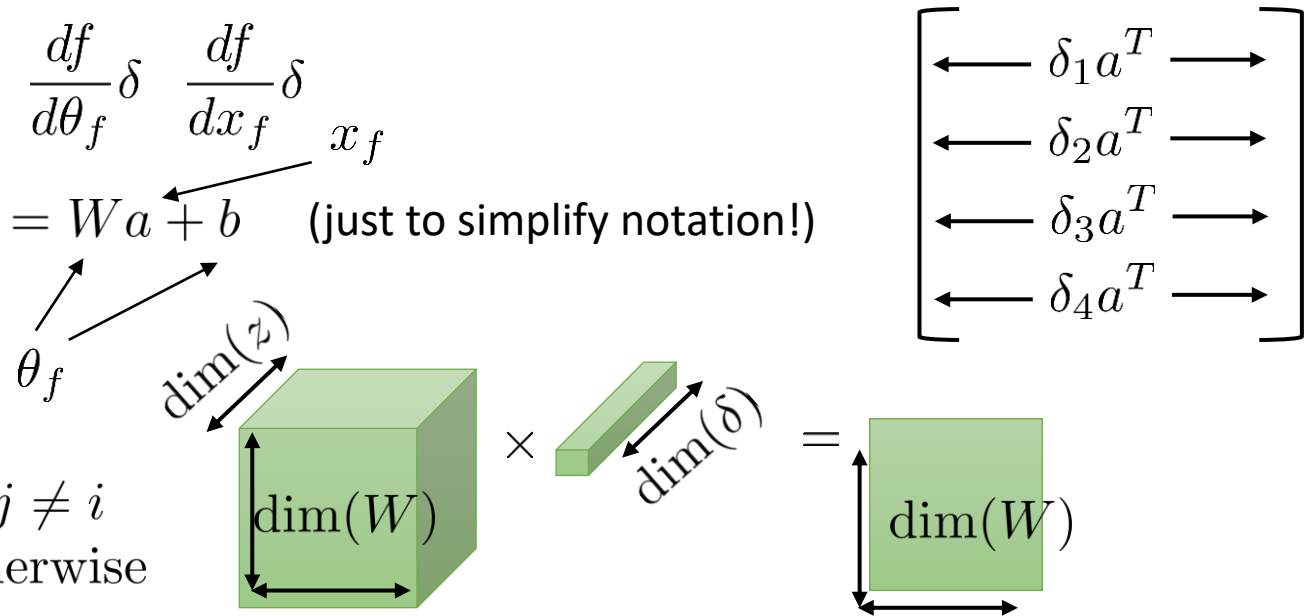
ReLU

# Backpropagation recipes: linear layer



for each function, we need to compute: $\dfrac{df}{d\theta_f}\delta \qquad \dfrac{df}{dx_f}\delta$
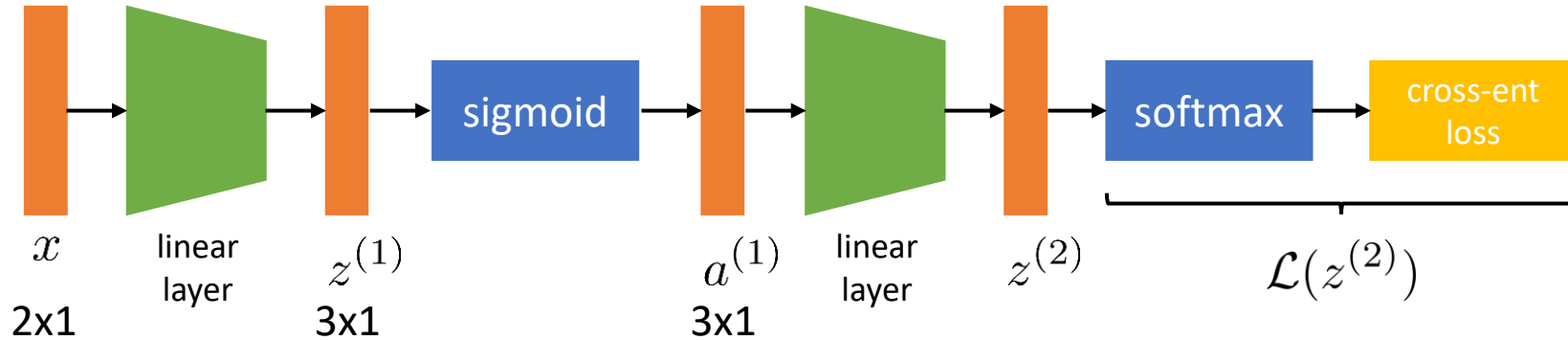
linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)} \qquad z = Wa + b$ (just to simplify notation!)

$\dfrac{dz}{dW}\delta = \sum_i \dfrac{dz_i}{dW}\delta_i = \delta a^T$

$z_i = \sum_k W_{ik}a_k + b_i \qquad \dfrac{dz_i}{dW_{jk}} = \begin{cases} 0 \text{ if } j \neq i \\ a_k \text{ otherwise} \end{cases}$

# Backpropagation recipes: linear layer



for each function, we need to compute: $\dfrac{df}{d\theta_f}\delta$ $\quad$ $\dfrac{df}{dx_f}\delta$

linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)}$ $\qquad$ $z = Wa + b$ $\qquad$ (just to simplify notation!)

$$\frac{dz}{db}\delta = \delta$$

$$z_i = \sum_k W_{ik}a_k + b_i \qquad \frac{dz_i}{db_j} = \text{Ind}(i = j) \qquad \frac{dz}{db} = \mathbf{I}$$
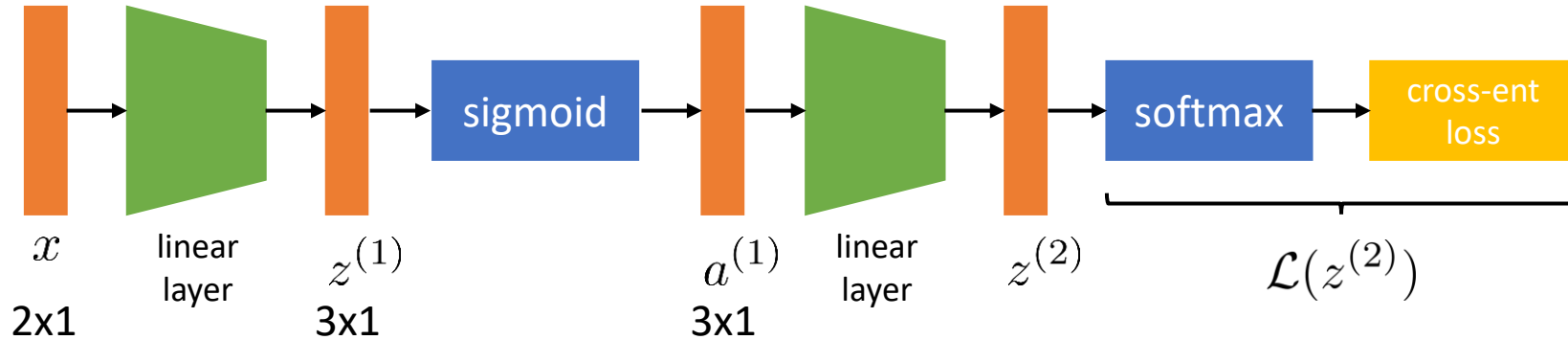
# Backpropagation recipes: linear layer



for each function, we need to compute: $\dfrac{df}{d\theta_f}\delta \quad \dfrac{df}{dx_f}\delta$

linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)} \qquad z = Wa + b$ (just to simplify notation!)

$$\frac{dz}{da}\delta = W^T\delta$$

$$z_i = \sum_k W_{ik}a_k + b_i \quad \frac{dz_i}{da_k} = W_{ik} \quad \frac{dz}{da} = W^T \qquad \left(\frac{dy}{dx}\right)_{ij} = \frac{dy_j}{dx_i}$$

# Backpropagation recipes: linear layer



$x$  
2x1

linear layer

$z^{(1)}$  
3x1

sigmoid

$a^{(1)}$  
3x1

linear layer

$z^{(2)}$

softmax

cross-ent loss

$\mathcal{L}(z^{(2)})$

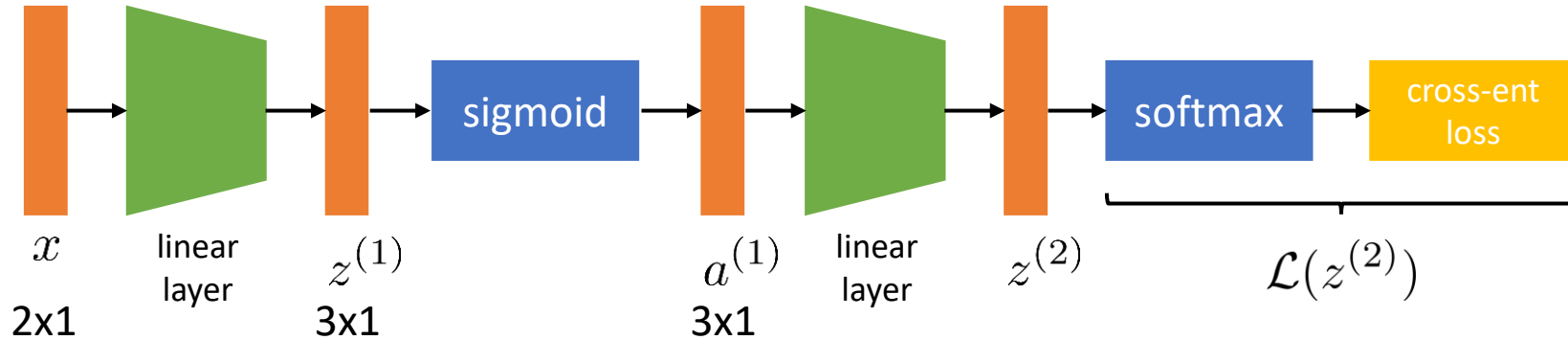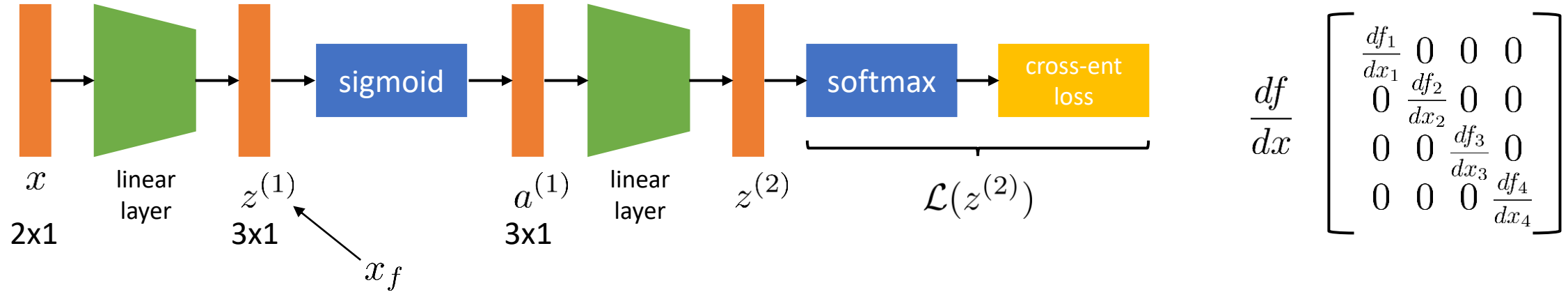for each function, we need to compute: $\dfrac{df}{d\theta_f}\delta$ $\dfrac{df}{dx_f}\delta$

linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)}$ $z = Wa + b$ (just to simplify notation!)

$$\frac{dz}{da}\delta = W^T\delta \qquad \frac{dz}{dW}\delta = \delta a^T \qquad \frac{dz}{db}\delta = \delta$$

$\underbrace{\qquad\qquad}$ $\underbrace{\qquad\qquad\qquad\qquad}$

$\dfrac{df}{dx_f}\delta$ $\dfrac{df}{d\theta_f}\delta$

# Backpropagation recipes: sigmoid



$$\frac{df}{dx} \begin{bmatrix} \frac{df_1}{dx_1} & 0 & 0 & 0 \\ 0 & \frac{df_2}{dx_2} & 0 & 0 \\ 0 & 0 & \frac{df_3}{dx_3} & 0 \\ 0 & 0 & 0 & \frac{df_4}{dx_4} \end{bmatrix}$$
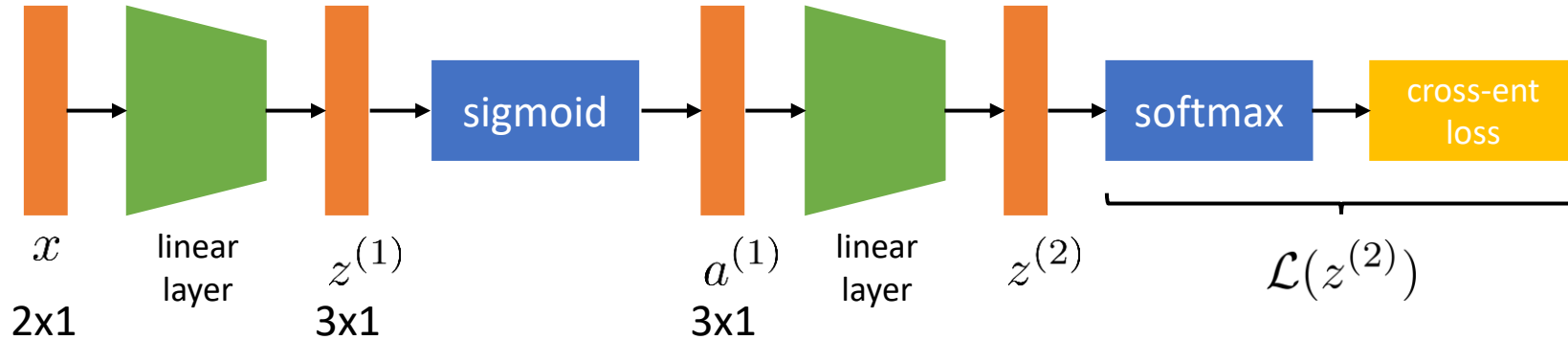
for each function, we need to compute: $\quad \dfrac{df}{d\theta_f}\delta \quad \dfrac{df}{dx_f}\delta$

$$\sigma(z_i) = \frac{1}{1+\exp(-z_i)} \qquad \frac{df_i}{dz_i} = \underbrace{\frac{\exp(-z_i)}{1+\exp(-z_i)}}\underbrace{\frac{1}{1+\exp(-z_i)}} = (1-\sigma(z_i))\sigma(z_i)$$

$$\left(\frac{df}{dz}\delta\right)_i = (1-\sigma(z_i))\sigma(z_i)\delta_i \qquad \underbrace{\frac{1+\exp(-z_i)}{1+\exp(-z_i)} - \frac{1}{1+\exp(-z_i)}}_{1-\sigma(z_i)} \quad \overset{\sigma(z_i)}{}$$
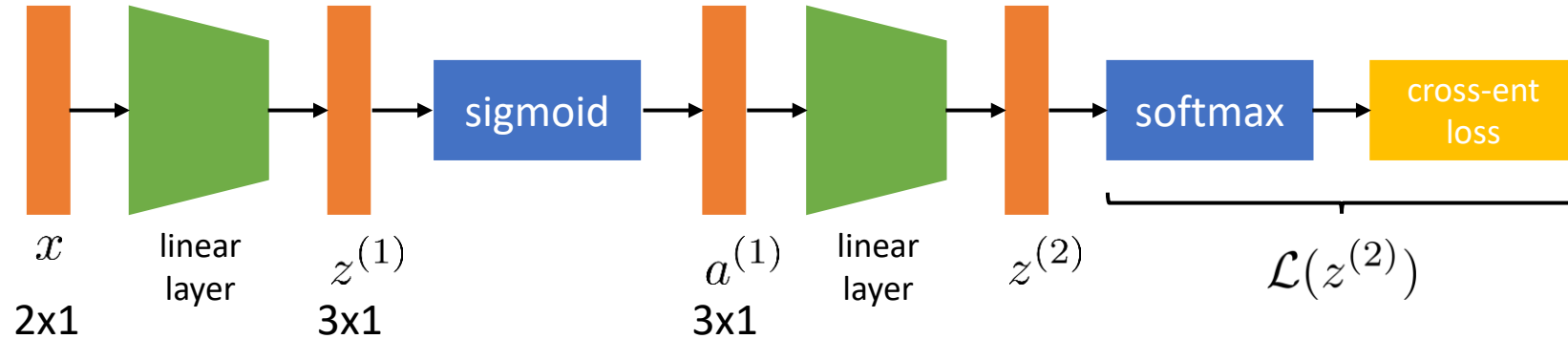
# Backpropagation recipes: ReLU



for each function, we need to compute: $\dfrac{df}{d\theta_f}\delta$ $\dfrac{df}{dx_f}\delta$

$$f_i(z_i) = \max(0, z_i) \qquad \frac{df_i}{dz_i} = \mathrm{Ind}(z_i \geq 0)$$

$$\left(\frac{df}{dz}\delta\right)_i = \mathrm{Ind}(z_i \geq 0)\delta_i$$

# Summary



forward pass: calculate each $a^{(i)}$ and $z^{(i)}$

backward pass:
initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$
for each $f$ with input $x_f$ & params $\theta_f$ from end to start:

$$\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f}\delta$$

$$\delta \leftarrow \frac{df}{dx_f}\delta$$

for each function, we need to compute:

$$\frac{df}{d\theta_f}\delta \qquad \frac{df}{dx_f}\delta$$

linear layer
softmax + cross-entropy
sigmoid
ReLU

## CSW182 (2021)· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

**视频·B 站 [ 扫码或点击链接 ]**
https://www.bilibili.com/video/BV1Ff4y1n7ar

**课件 & 代码·博客 [ 扫码或点击链接 ]**
http://blog.showmeai.tech/berkeley-csw182

Berkeley
循环神经网络
可视化
梯度策略
Q-Learning
风格迁移
模仿学习
元学习
计算机视觉
机器学习基础
生成模型
卷积网络

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS231n |

### # Awesome AI Courses Notes Cheatsheets· 持续更新中

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]