

# Lecture 2: Machine Learning and Deep Learning Basics



Arsalan Mosenia

Chinmay Hegde

# Linear Regression

- Now consider input:  $x \in R^d$  and output  $y \in R$  the goal is to learn

$$y = f(x) = w_d x_d + \dots + w_1 x_1 + w_0$$

# Linear Regression

Given training dataset  $X \in R^{N \times d}$  and  $Y \in R^N$

Training sample  $\rightarrow$

$$\begin{pmatrix} 1 & x_{11} & x_{12} \cdots & x_{1d} \\ \boxed{1} & \boxed{x_{21}} & \boxed{x_{22} \cdots} & \boxed{x_{2d}} \\ 1 & x_{N1} & x_{N2} \cdots & x_{Nd} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \vdots \\ \widehat{y}_N \end{pmatrix}$$

$$\widehat{Y} = XW$$

Note: for simplicity we will assume that  $X$  includes a column of 1s

# Linear Regression

$$MSE = \sum (y - \hat{y})^2 = (Y - \hat{Y})^T \cdot (Y - \hat{Y}) = (Y - XW)^T \cdot (Y - XW)$$

**Objective:**

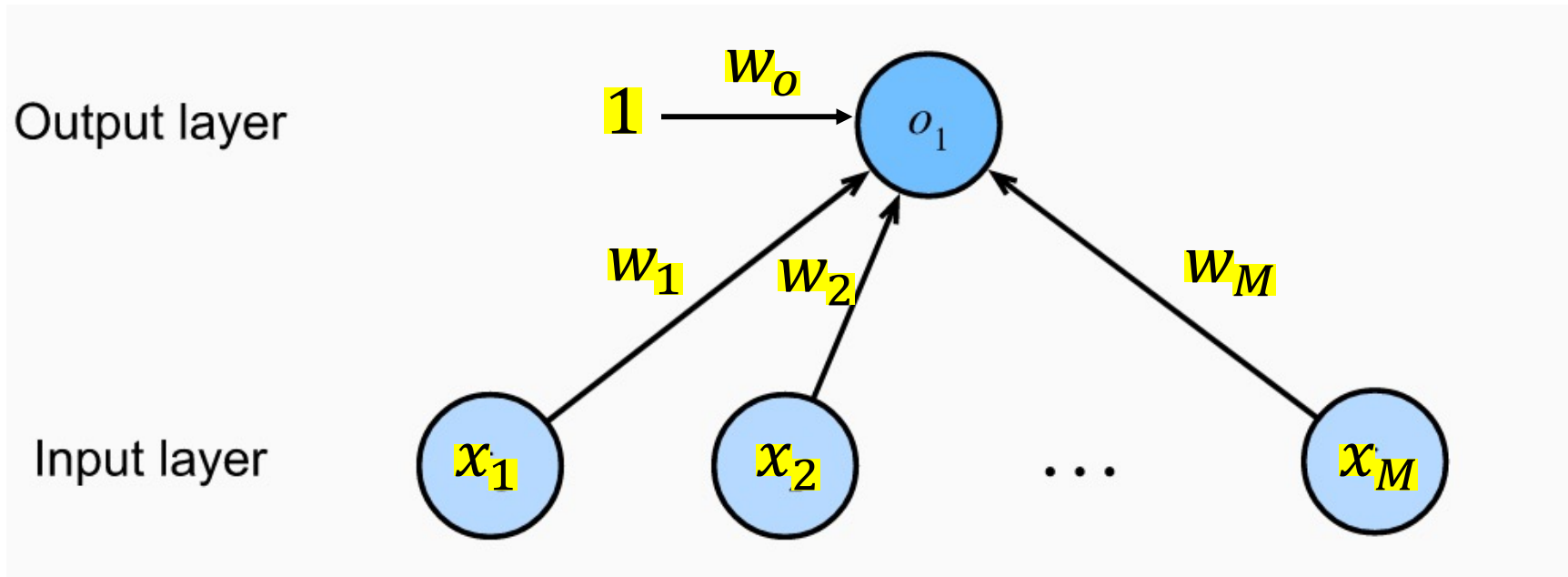
$$\min_W (Y - XW)^T \cdot (Y - XW)$$

Set derivate with respect to  $W$  to zero!

**Solution:**

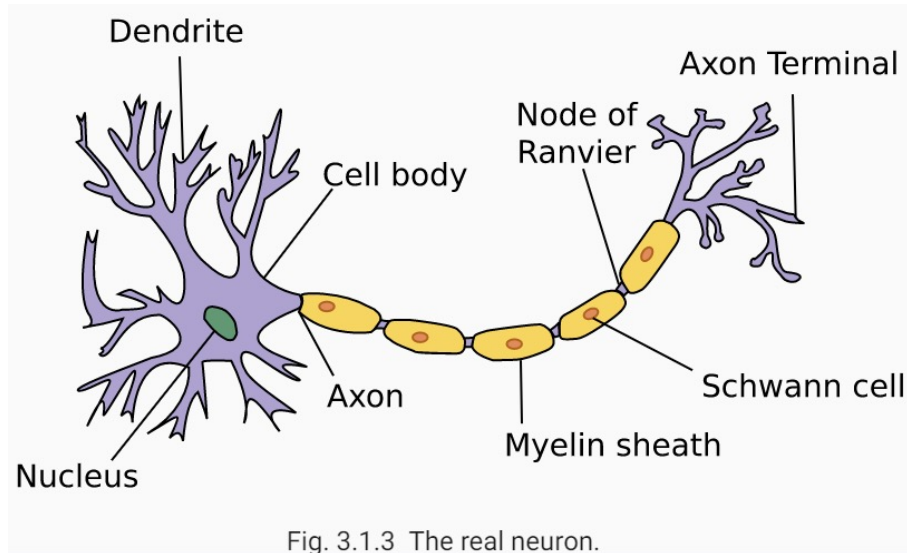
$$W^* = (X^T X)^{-1} X^T Y$$

# “Connectionist” Perspective



- Each input is connected to the output node via an edge
- Output node sums inputs multiplied with their respective weights

# A Biological Perspective



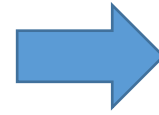
- **Dendrites**: accumulate inputs from neighboring neurons (inputs)
- **Nucleus**: performs a weighted sum of inputs, based on strengths of “synapses” + some non-linear activation
- **Axon**: output

$$y = f(x) = w_d x_d + \dots + w_1 x_1 + w_0$$

# Activation Functions

$$o = \sigma(y)$$

Non-linear “activation” function  
(example: “threshold”)



For regression, the activation function is just the identity fn.

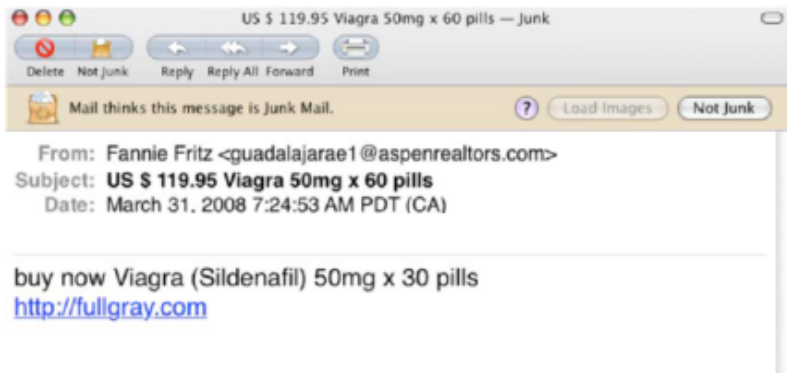
$$y = f(x) = w_d x_d + \dots + w_1 x_1 + w_0$$

- Non-linear “activation” function (example: “threshold”)

# Classification

## Task (T):

- Emails  $x \in$  all possible emails and  $y \in \{spam, non\_spam\}$  find



SPAM

$$f: x \rightarrow y$$

## Training Data:

A "training dataset" emails marked as "spam" or "non-spam"

“Supervised Learning (Classification)”

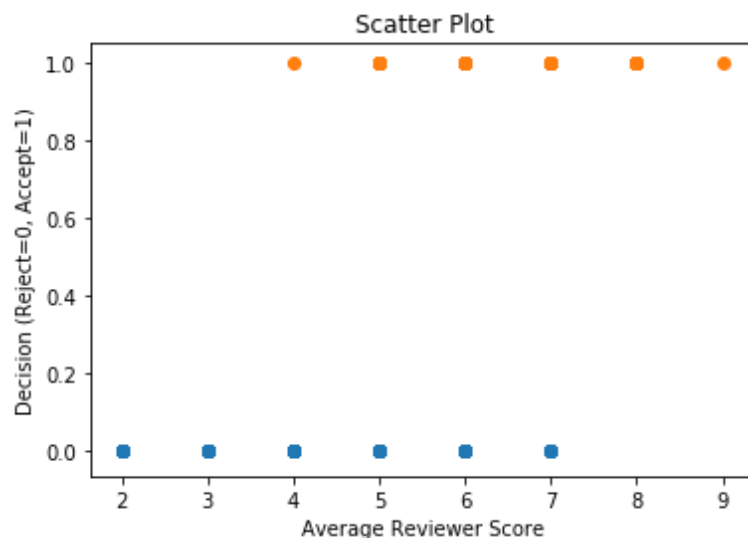


# Binary Classification

- Simplest example where  $x \in R$  and  $y \in \{0,1\}$  “Categorical variable”
- Dataset of ICLR’18 review scores vs. accept/reject decisions

TL;DR	_bibtex	abstract	authorids	authors	conf_1	conf_2	conf_3	decision	review	review_1	review_2	review_3	title
None	@article{\nsharma2018hyperedge2vec,\ntitle={H...	Data structured in form of overlapping or non-...	[sharm170@umn.edu, srjoty@ntu.edu.sg, himanshu...	[Ankit Sharma, Shafiq Joty, Himanshu Kharkwal,...	3.0	3.0	4.0	Reject	5.000000	5.0	5.0	5.0	Hyperedge2vec: Distributed Representations for...
Query-based black-box attacks on deep neural n...	@article{\nnitin2018exploring,\ntitle={Explori...	Existing black-box attacks on deep neural netw...	[abhagoji@princeton.edu, _w@eecs.berkeley.edu,...	[Arjun Nitin Bhagoji, Warren He, Bo Li, Dawn S...	4.0	3.0	4.0	Reject	6.000000	6.0	6.0	7.0	Exploring the Space of Black-box Attacks on De...
A theory and algorithmic framework for predict...	@article{\nd.2018learning,\ntitle={Learning We...	Predictive models that generalize well under d...	[fredrikj@mit.edu, kallus@cornell.edu, urish22...	[Fredrik D. Johansson, Nathan Kallus, Uri Shal...	3.0	3.0	4.0	Reject	6.666667	5.0	8.0	7.0	Learning Weighted Representations for Generali...
We prove that DNN is a recursively approximate...	@article{\nzheng2018understanding,\ntitle={Und...	Deep learning achieves remarkable generalizati...	[zhenggh@mail.ustc.edu.cn, jtsang@bjtu.edu.cn,...	[Guanhua Zheng, Jitao Sang, Changsheng Xu]	3.0	3.0	2.0	Reject	3.666667	2.0	3.0	6.0	Understanding Deep Learning Generalization by

# Binary Classification



Can you fit a linear model to this data?

TL;DR	_bibtex	abstract	authorids	authors	conf_1	conf_2	conf_3	decision	review	review_1	review_2	review_3	title
None	@article{\nsharma2018hyperedge2vec,\nntitle={H...	Data structured in form of overlapping or non-...	[sharm170@umn.edu, srjoty@ntu.edu.sg, himanshu...	[Ankit Sharma, Shafiq Joty, Himanshu Kharkwal,...	3.0	3.0	4.0	Reject	5.000000	5.0	5.0	5.0	Hyperedge2vec: Distributed Representations for...
Query-based black-box attacks on deep neural n...	@article{\nnitin2018exploring,\nntitle={Explori...	Existing black-box attacks on deep neural netw...	[abhagoji@princeton.edu, _w@eecs.berkeley.edu,...	[Arjun Nitin Bhagoji, Warren He, Bo Li, Dawn S...	4.0	3.0	4.0	Reject	6.000000	6.0	6.0	7.0	Exploring the Space of Black-box Attacks on De...
A theory and algorithmic framework for predict...	@article{\nd.2018learning,\nntitle={Learning We...	Predictive models that generalize well under d...	[fredrikj@mit.edu, kallus@cornell.edu, urish22...	[Fredrik D. Johansson, Nathan Kallus, Uri Shal...	3.0	3.0	4.0	Reject	6.666667	5.0	8.0	7.0	Learning Weighted Representations for Generali...
We prove that DNN is a recursively approximate...	@article{\nzheng2018understanding,\nntitle={Und...	Deep learning achieves remarkable generalizati...	[zhenggh@mail.ustc.edu.cn, jtsang@bjtu.edu.cn,...	[Guanhua Zheng, Jitao Sang, Changsheng Xu]	3.0	3.0	2.0	Reject	3.666667	2.0	3.0	6.0	Understanding Deep Learning Generalization by

# Logistic Regression

- First, let  $p = \Pr\{y = 1|x\}$      $\Pr\{\text{Decision}=\text{Accept}| \text{Score}\}$
- Consider the following function:  $g = \log(\frac{p}{1-p})$ 
  - $g$  is called the “logits” function
  - Logistic Regression: **fit logits function using a linear model!**

$$g = \log(\frac{p}{1-p}) = w_1x + w_0$$

# Logistic Regression

$$g = \log\left(\frac{p}{1-p}\right) = w_1x + w_0$$

- $\Pr\{\text{Decision}=\text{Accept}|\text{Score}\}$
- What is  $\Pr\{\text{Decision}=\text{Reject}|\text{Score}\}$

$$p = \frac{1}{1 + e^{-(w_1x + w_0)}}$$

$$1 - p = \frac{e^{-(w_1x + w_0)}}{1 + e^{-(w_1x + w_0)}}$$

Decision Boundary?

# How to Find Parameters? $W_1$ , $W_0$

- What is **Maximum Likelihood Estimation** (MLE)?
  - Estimating the parameters of an assumed probability distribution
  - Given some observed data.

This is achieved by maximizing a likelihood function so that, under the assumed statistical model, the observed data is most probable.

# Maximum Likelihood Estimation (MLE)

- Goal: Determine parameters for which the observed data have the highest joint probability.
- $L_n(w) = L_n(w; y) = f_n(y; w)$  where  $Y = (y_1, y_2, y_3, \dots, y_n)$
- Find  $w \rightarrow \operatorname{argmax} L_n(w; Y)$

For Independent and identically distributed random (iid) variables:

$$f_n(\mathbf{y}; \theta) = \prod_{k=1}^n f_k^{\text{univar}}(y_k; \theta) .$$

# Maximum Likelihood Estimation (MLE)

- What is the likelihood that the dataset came from our model?

$$p = \frac{1}{1 + e^{-(w_1 x + w_0)}}$$

$$1 - p = \frac{e^{-(w_1 x + w_0)}}{1 + e^{-(w_1 x + w_0)}}$$

#	X	Y
1	$x_1 = 3$	$y_1 = 0$
2	$x_2 = 8$	$y_2 = 1$
..		
..		
N	$x_N = 6$	$y_N = 1$

$$Likelihood = \frac{e^{-(3w_1 + w_0)}}{1 + e^{-(3w_1 + w_0)}} * \frac{1}{1 + e^{-(8w_1 + w_0)}} * \dots * \frac{1}{1 + e^{-(6w_1 + w_0)}}$$

# Maximum Likelihood Estimation (MLE)

- $\hat{Y}_i = P(X_i) = P(Y_i = 1 | X_i)$
- Likelihood =  $\prod_{Y_i=1} P(X_i) * \prod_{Y_i=0} (1 - P(X_i))$   
 $= \prod_{\text{for all samples}} P(X_i)^{y_i} * (1 - P(X_i))^{1-y_i}$

MLE  $\equiv$  Maximizing Log of L

$$\text{Log of Likelihood} = \sum_{i=1}^n [y_i * \text{Log}(P(x_i)) + (1 - y_i) * \text{Log}(1 - P(x_i))]$$



# Maximum Likelihood Estimation (MLE)

$$\text{Log of Likelihood} = \sum_{i=1}^n [y_i * \text{Log}(\hat{Y}_i) + (1 - y_i) * \text{Log}(1 - \hat{Y}_i)]$$

MLE  $\equiv$  Maximizing (Log of L)  $\equiv$  Minimizing ( $-\text{Log of L}$ )

$$-\text{Log}(\text{Likelihood}) = -\sum_{i=1}^n [y_i * \text{Log}(\hat{Y}_i) + (1 - y_i) * \text{Log}(1 - \hat{Y}_i)]$$

Cost function?  $-\text{Log}(\text{Likelihood})$

Loss function?  $-[y_i * \text{Log}(\hat{Y}_i) + (1 - y_i) * \text{Log}(1 - \hat{Y}_i)]$

Binary Cross-entropy Loss = Log-Loss = Logistic Regression Loss

# Maximum Likelihood Estimation (MLE)

- What is the likelihood that the dataset came from our model?

#	X	Y
1	$x_1 = 3$	$y_1 = 0$
2	$x_2 = 8$	$y_2 = 1$
..		
..		
N	$x_N = 6$	$y_N = 1$

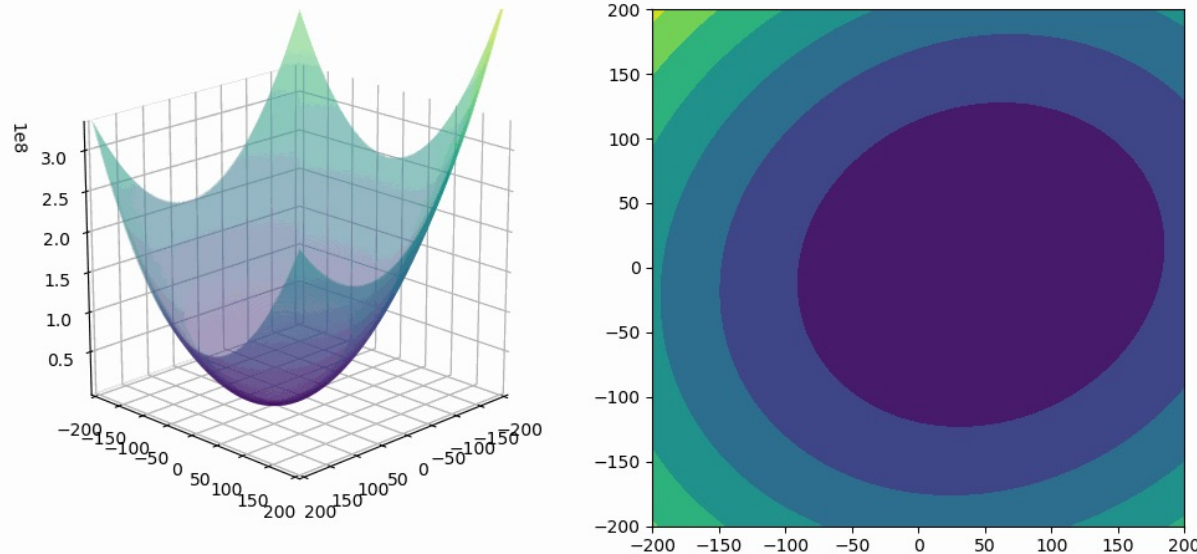
$$\text{Likelihood} = \frac{e^{-(3w_1+w_0)}}{1 + e^{-(3w_1+w_0)}} * \frac{1}{1 + e^{-(8w_1+w_0)}} * \dots * \frac{1}{1 + e^{-(6w_1+w_0)}}$$

Maximize Likelihood  
= Maximize Log-Likelihood  
= Minimize (-Log-Likelihood)

↑  
“Cost”

# How to Minimize the Cost Function?

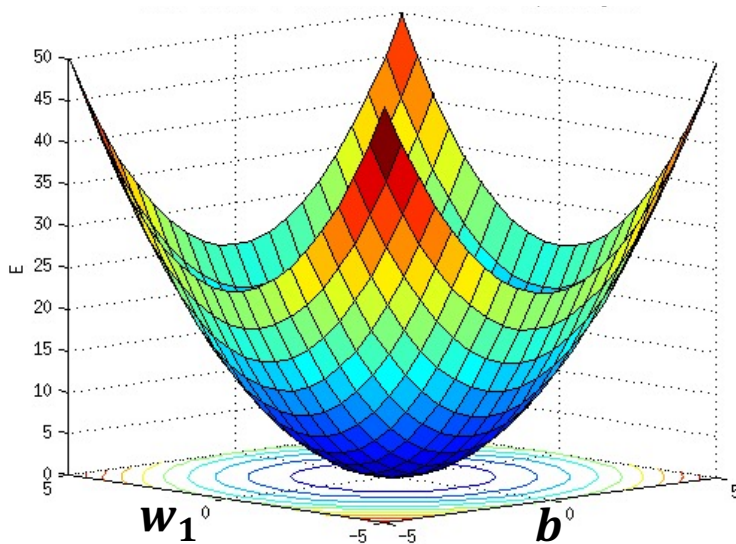
- Sometimes, we cannot derive a formula to minimize loss. (For which model and loss function were we able to do this? \_\_\_\_\_)
- Using Gradient Descent
- Repeated steps in the opposite direction of the gradient at the current point, because this is the direction of steepest descent.



# Gradient Descent for Minimizing Cost

$$J(\mathbf{w}, b) = \sum_{i=1}^N l_i(\mathbf{w}, b)$$

Gradient Update:



Step size

Partial derivative

$$(\mathbf{w}, b) \leftarrow (\mathbf{w}, b) - \frac{\eta}{N} \sum_{i=1}^N \partial_{(\mathbf{w}, b)} l_i(\mathbf{w}, b)$$

# Computing the Partial Derivative

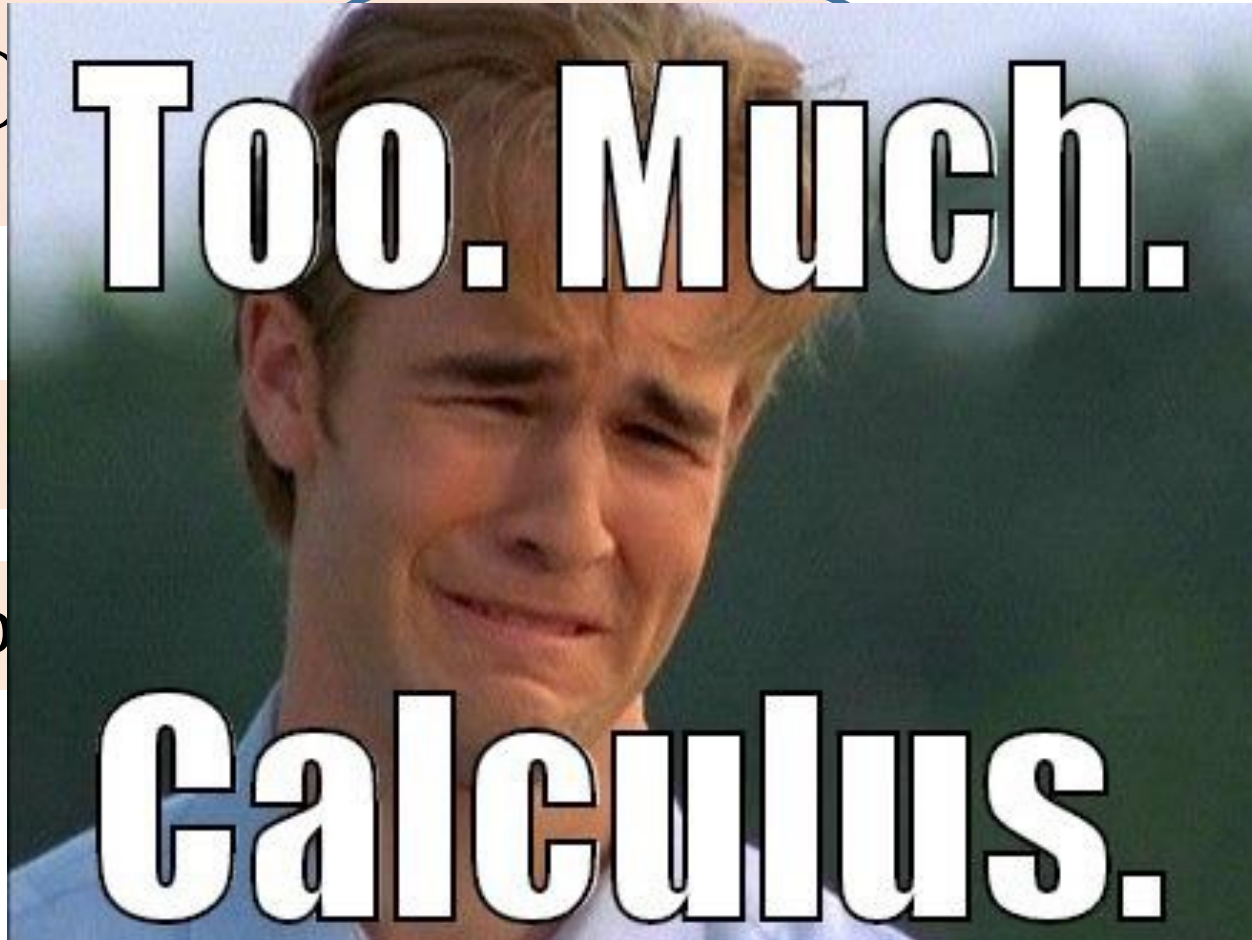
$N$

$$(\mathbf{w}, b) \leftarrow (\mathbf{w}, b)$$

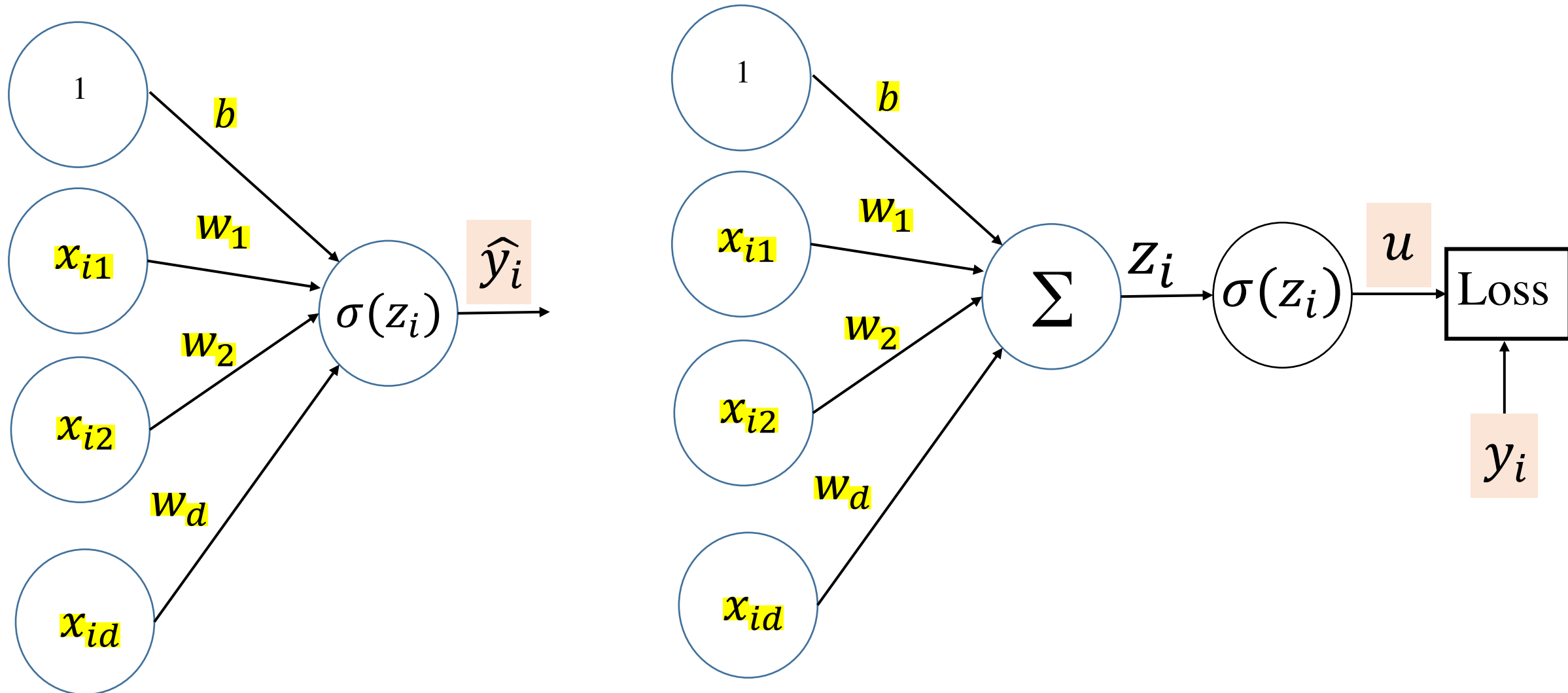
$$\partial_{(\mathbf{w}, b)} l_i(\mathbf{w}, b)$$

$$= \partial_{(\mathbf{w}, b)} [-y_i \log$$

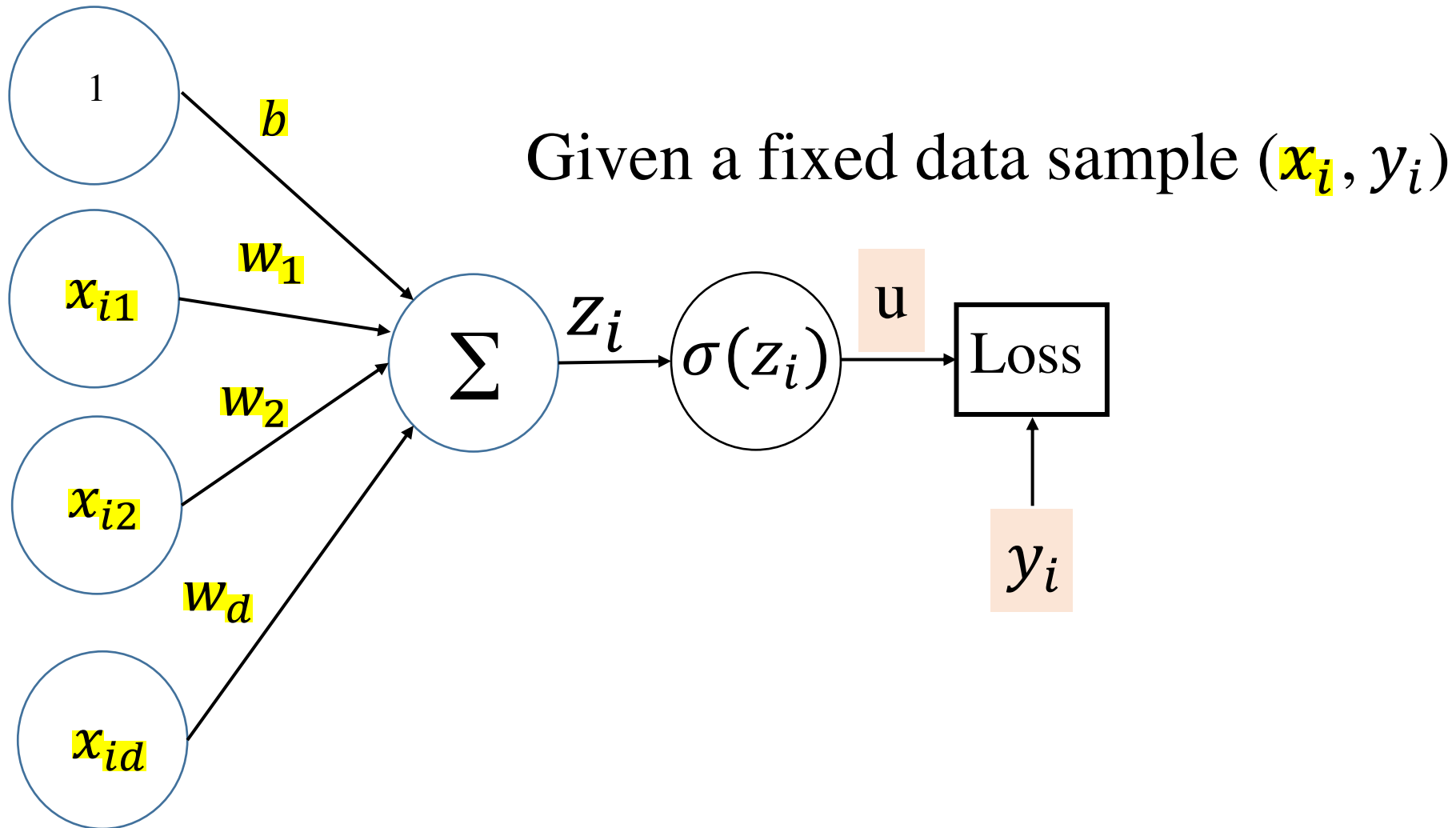
$$\sigma(\mathbf{w}^T \mathbf{x}_i + b)]$$



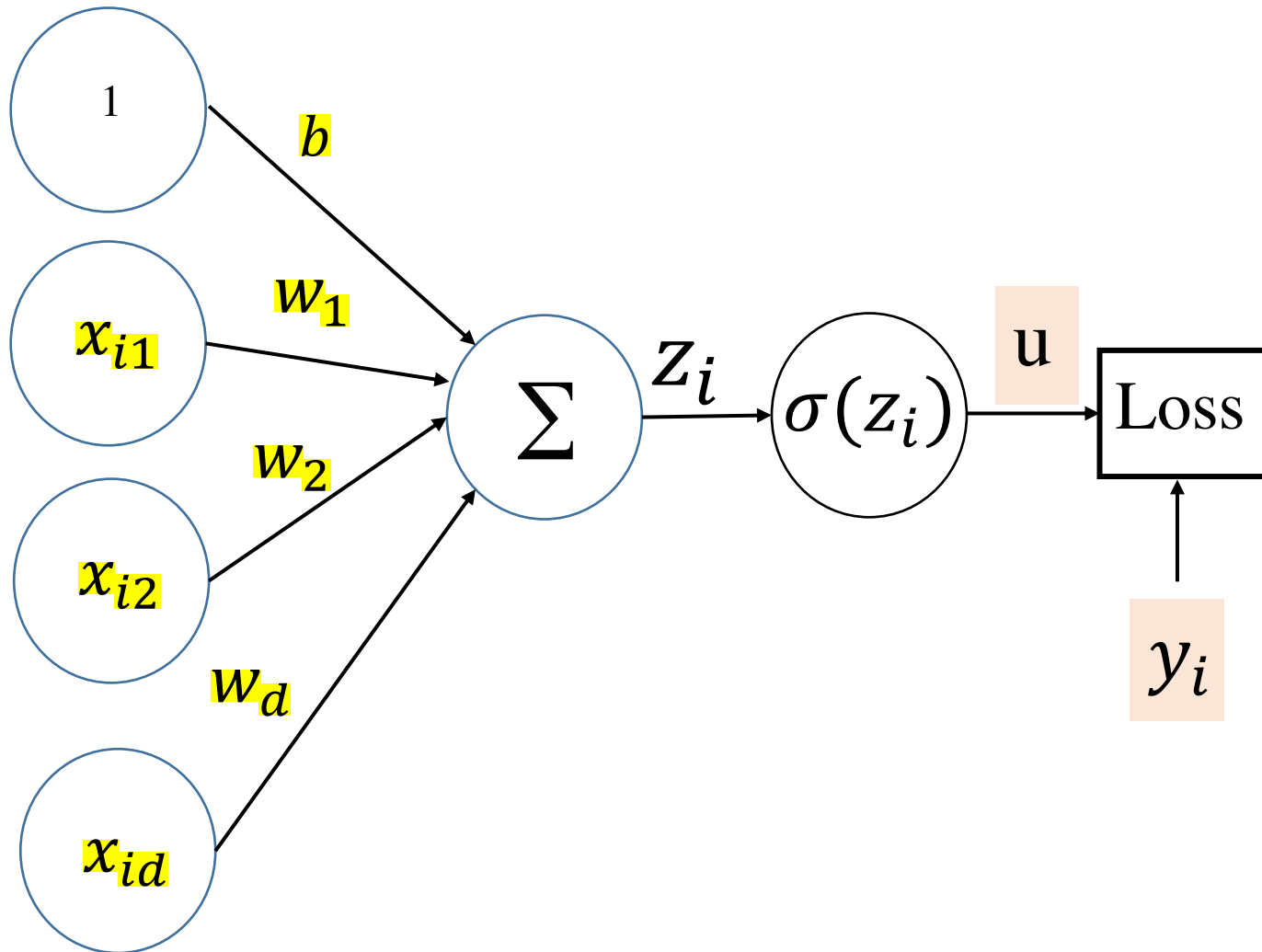
# Back Propagation: Computational Graph



# Back Propagation: Forward ()



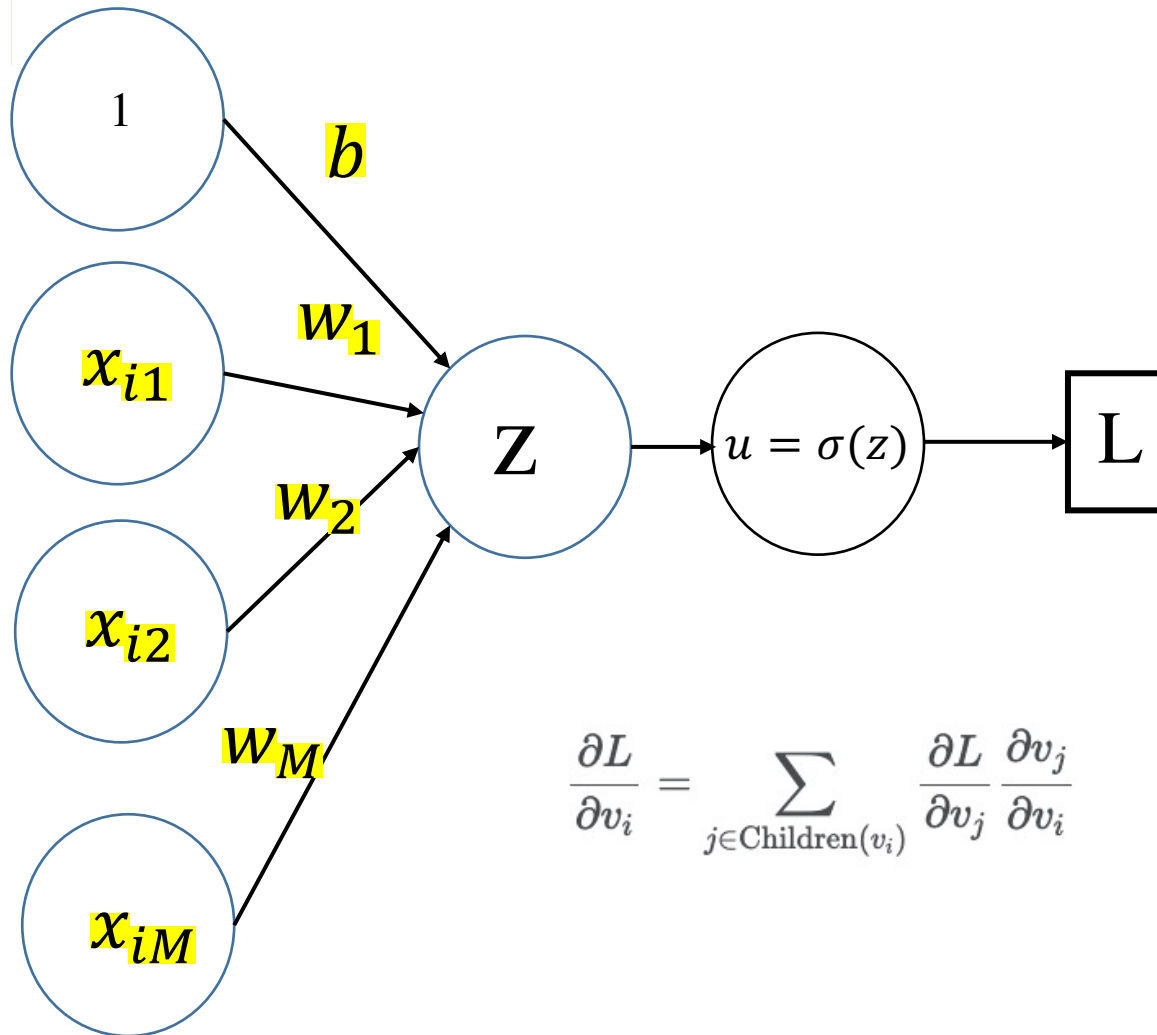
# Back Propagation: Backward()



$$\frac{\partial L}{\partial v_i} = \sum_{j \in \text{Children}(v_i)} \frac{\partial L}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$



# Back Propagation: Example



$$z = xw + b$$

$$u = \sigma(z)$$

$$\text{Loss} = -y \log(u) - (1 - y) \log(1 - u)$$

- $(\partial l) / \partial l = 1$
- $(\partial l) / \partial u = [(\partial l) / \partial l] * [(\partial l) / \partial u] = -y * \left(\frac{1}{u}\right) + (1 - y) * \frac{1}{1 - u}$
- $(\partial l) / \partial z = [(\partial l) / \partial u] * [(\partial u) / \partial z]$

$$\text{Note } \sigma'(z) = \sigma(z) * (1 - \sigma(z))$$

$$(\partial l) / \partial z = \left[ -y * \left(\frac{1}{u}\right) + (1 - y) * \frac{1}{1 - u} \right] * u(1 - u)$$

- $(\partial l) / \partial w = (\partial l) / \partial z * (\partial z) / \partial w$   

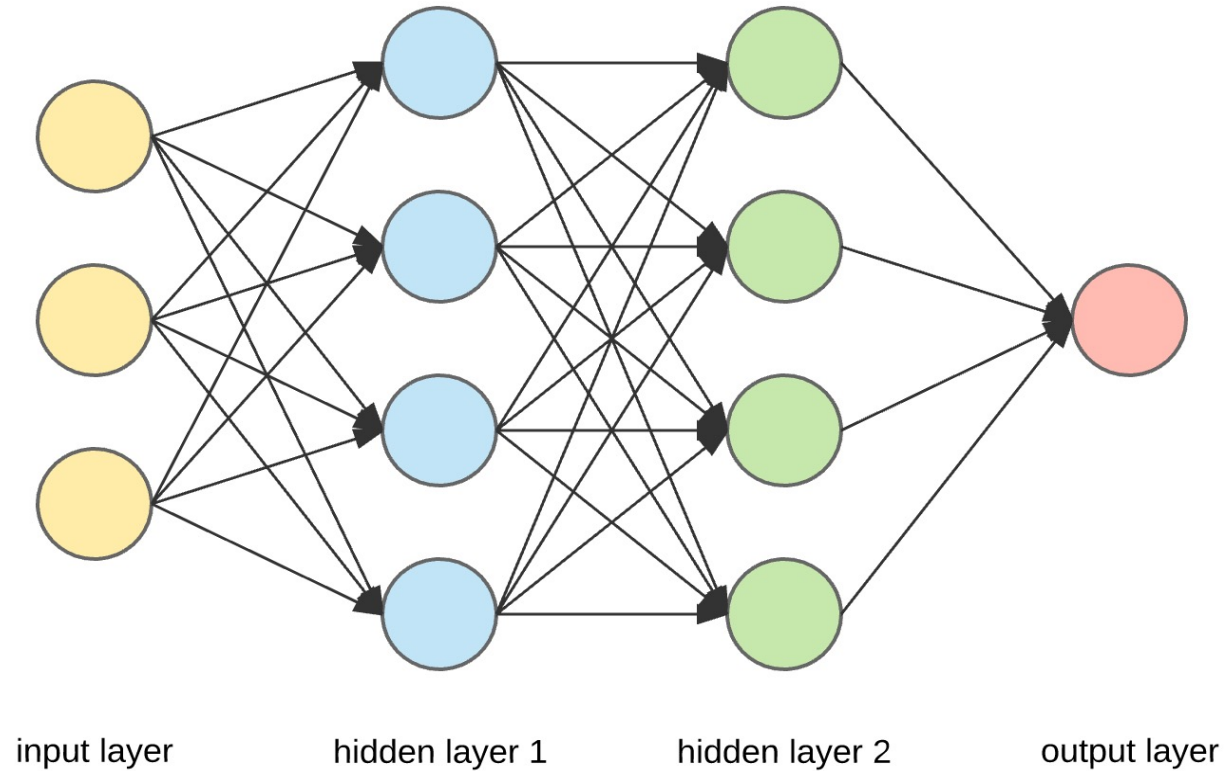
$$= \left[ -y * \left(\frac{1}{u}\right) + (1 - y) * \frac{1}{1 - u} \right] * u(1 - u) * (\partial z) / \partial w$$
  

$$= \left[ -y * \left(\frac{1}{u}\right) + (1 - y) * \frac{1}{1 - u} \right] * u(1 - u) * x$$

- $(\partial l) / \partial b = (\partial l) / \partial z * (\partial z) / \partial b$   

$$= \left[ -y * \left(\frac{1}{u}\right) + (1 - y) * \frac{1}{1 - u} \right] * u(1 - u) * 1$$

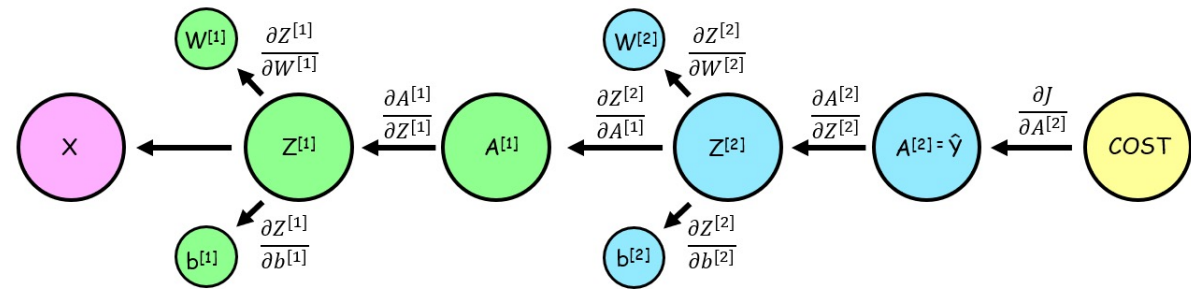
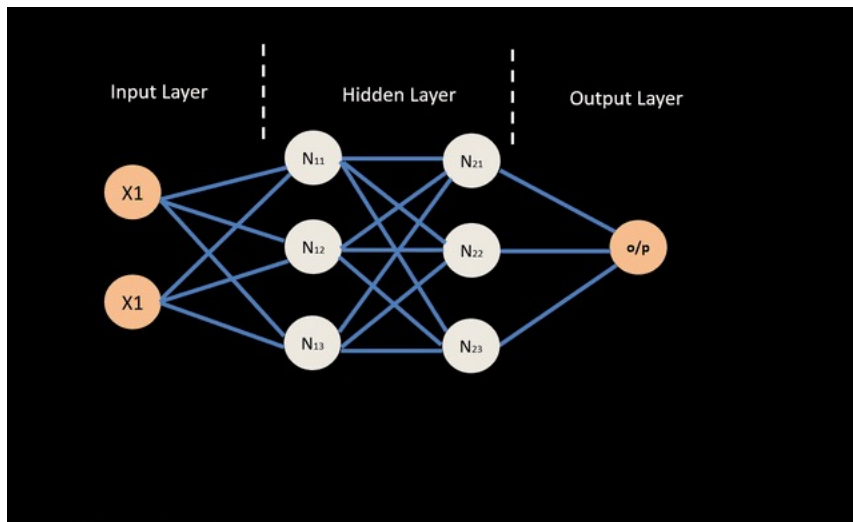
# Neural Networks



- 2-hidden layer network or 3-layer network
- **Params:** Weights/Biases
- Non-linearity? Activation functions
  - Relu, Leaky Relu, Tanh, Sigmoid

# Revisiting Back Propagation

- Every partial derivative depends cleanly on partial derivatives of child nodes
- The computation is very structured, and everything can be calculated by traversing once through the graph.
  - Gradients can be computed efficiently.



# Gradient Descent with Back Propagation

- GD is a generic way to find the local minimum
- Back propagation provides an efficient way to compute derivatives
- GD takes average of the weights suggested by examples for each step







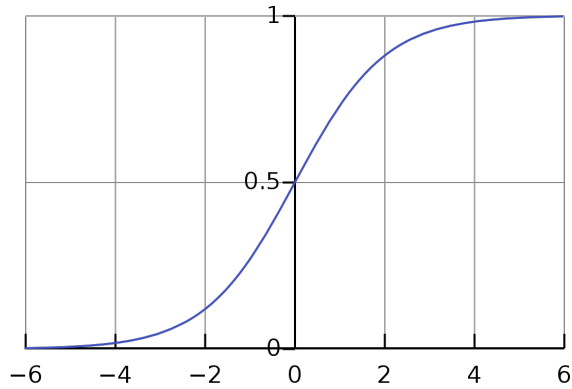
							Average over all training data ...
$w_0$	-0.08	+0.02	-0.02	+0.11	-0.05	-0.14	...
$w_1$	-0.11	+0.11	+0.07	+0.02	+0.09	+0.05	...
$w_2$	-0.07	-0.04	-0.01	+0.02	+0.13	-0.15	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$
$w_{13,001}$	+0.13	+0.08	-0.06	-0.09	-0.02	+0.04	...

Diagram illustrating the process of Gradient Descent with Back Propagation. The table shows the weights ( $w_0, w_1, w_2, \dots, w_{13,001}$ ) for different digits (2, 5, 0, 4, 1, 9) and the average weight over all training data. The weights are updated iteratively, with the final average weight for  $w_0$  being -0.08.

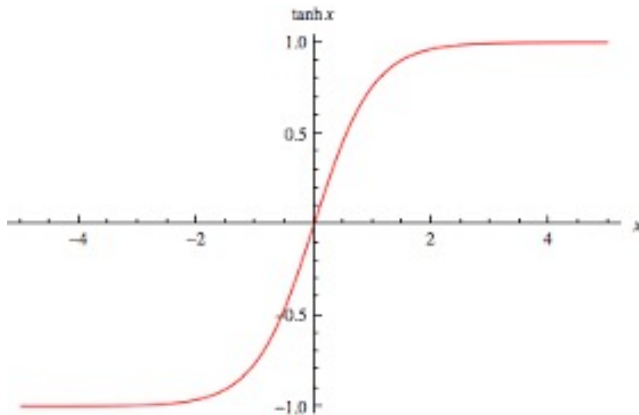
# Revisiting Activation Functions

Non-linearity? Activation functions

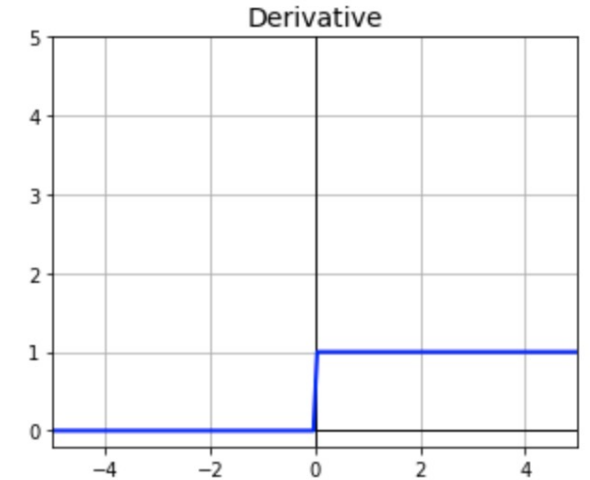
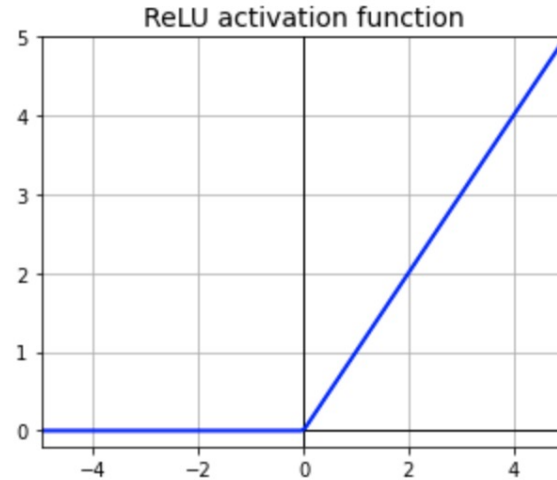
- Sigmoid



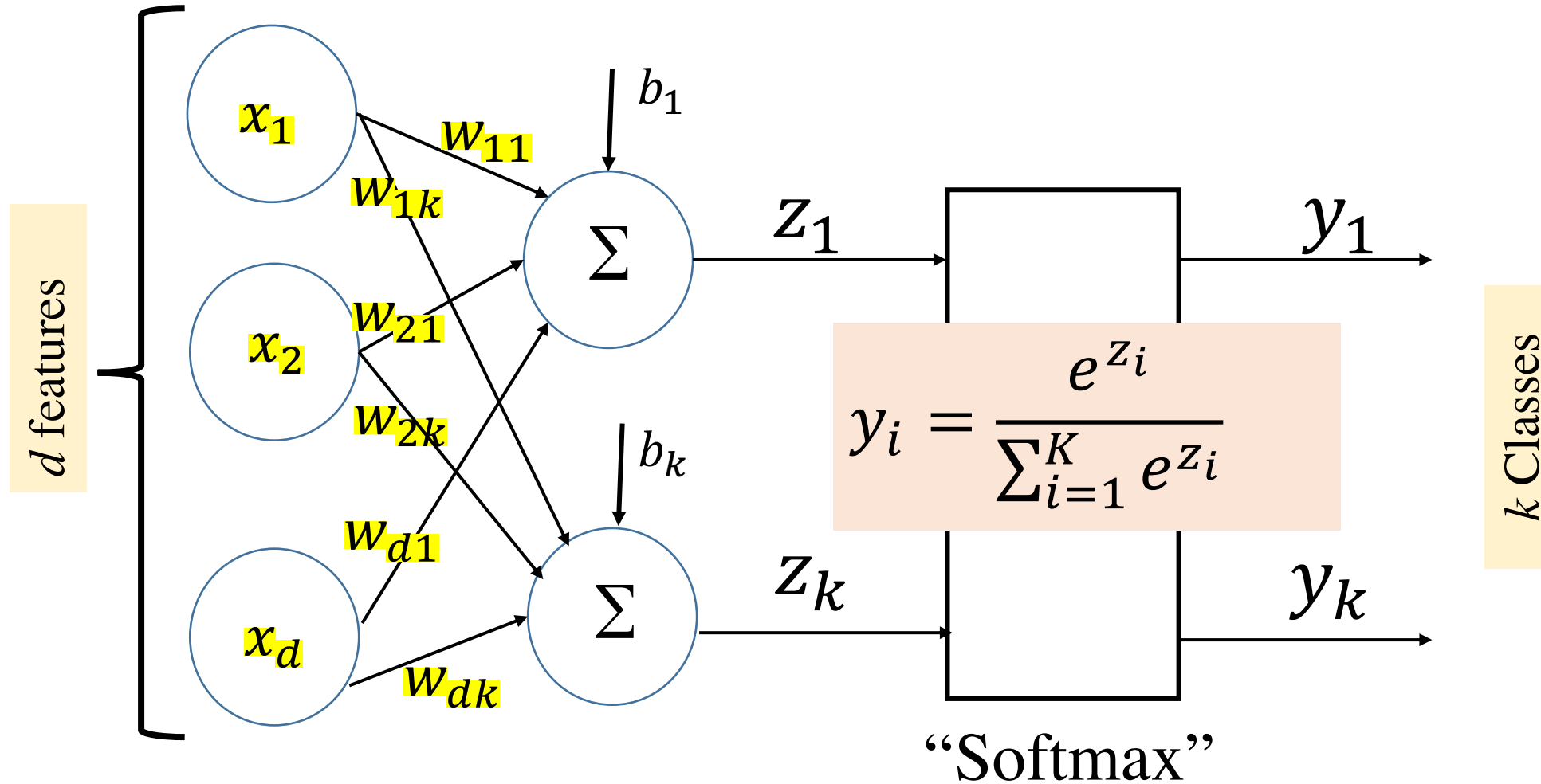
- Tanh



- Rectified Linear Unit (Relu)

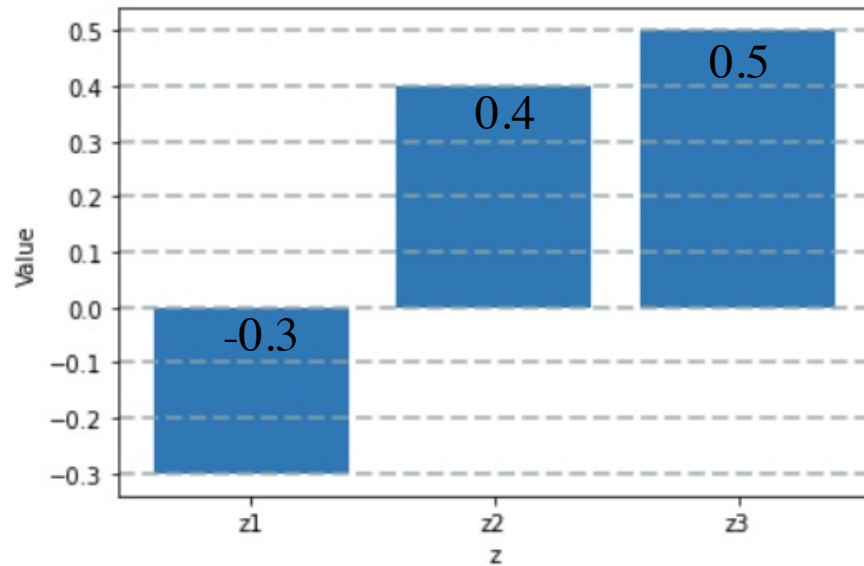


# Multi-Class Classification

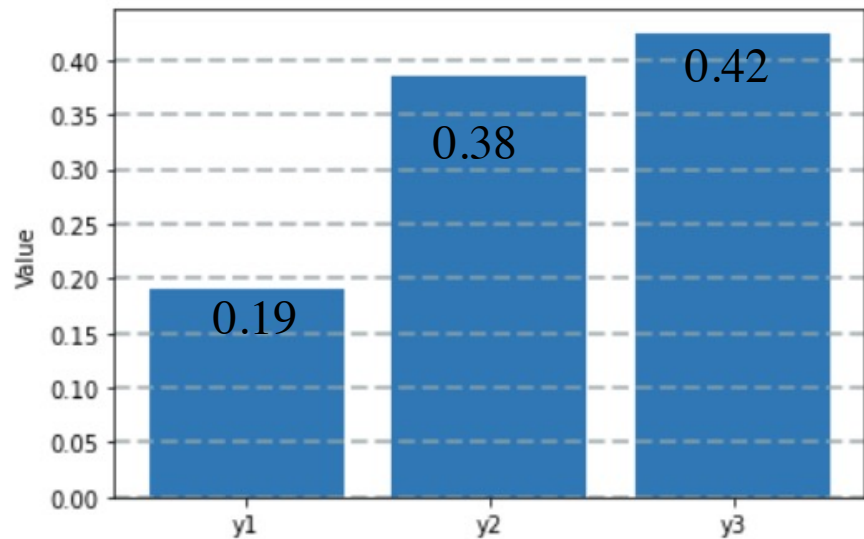


# Understanding Softmax

$$y_i = \frac{e^{z_i}}{\sum_{i=1}^K e^{z_i}}$$



The z vector has real-valued entries, including both positive and negative values

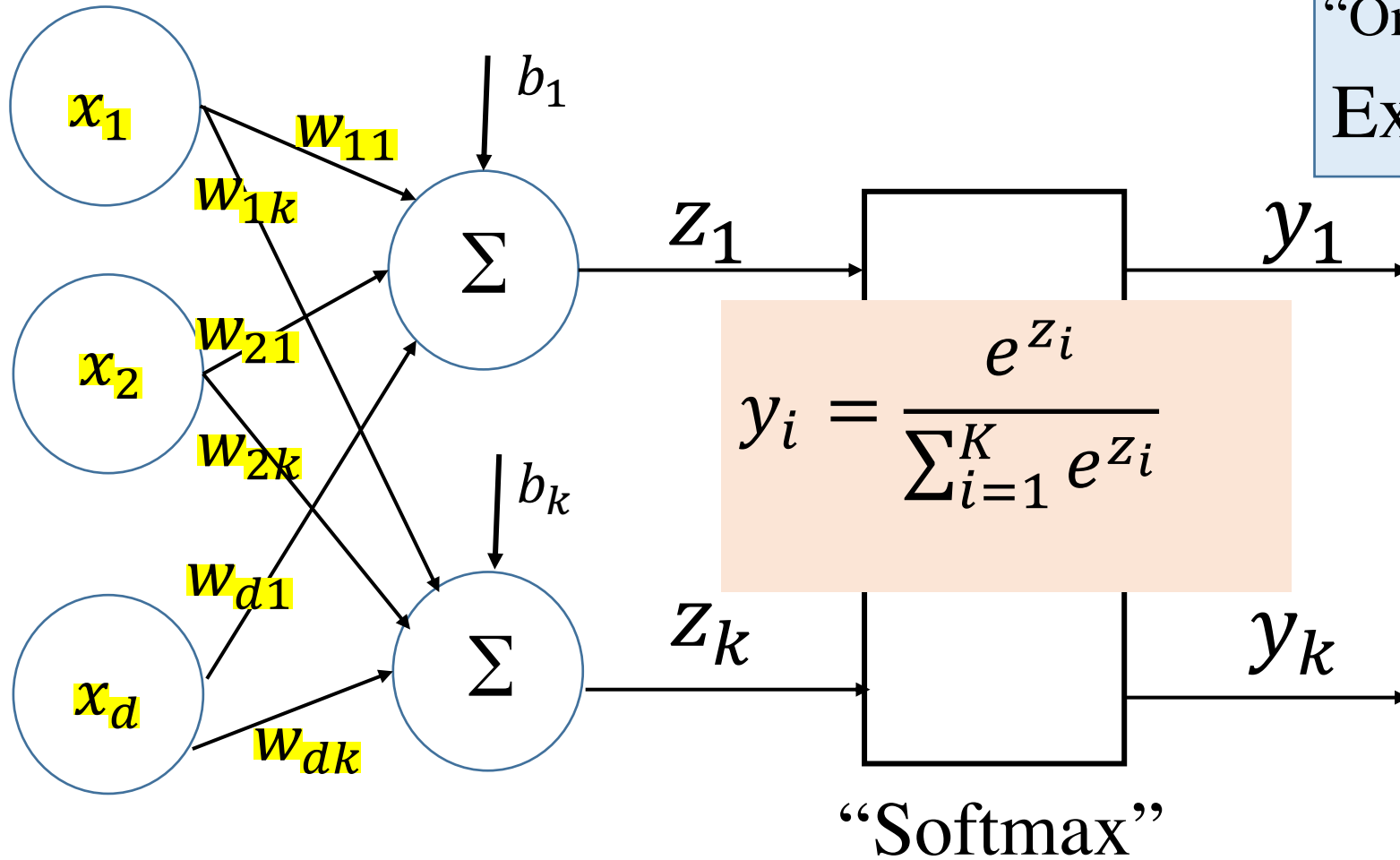


The y vector has:

- Positive real-valued entries
- Values sum to 1
- Preserves ordering of z

Can be viewed as a probability distribution over y!

# Cross-entropy Loss



Labels:  $[y_1 \dots y_k]$   
"One-hot Coded"  
Example:  $[0 \ 0 \ .. \ 1 \dots 0]$

**Cross-Entropy Loss**

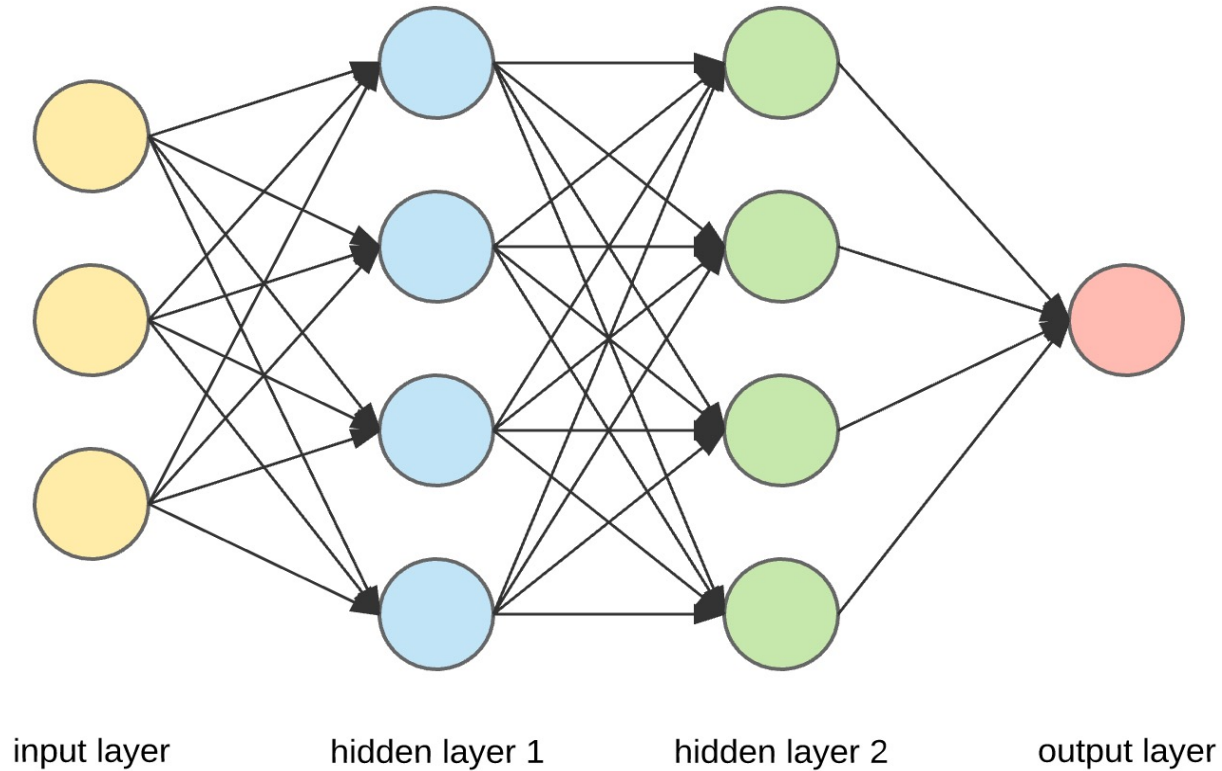
$$\sum_{i=1}^K -y_i \log(1 - \hat{y}_i)$$



# Summary

- Linear Regression
- Logistic Regression
  - MLE
  - Cross Entropy Loss
- Gradient Descent
- Backpropagation
- Neural Networks
  - Activation functions
- Multi-class Neural Networks
  - Softmax
  - One-hot encoding

# Some Questions



- What happens if we initialize all params to 1?
- What happens if we initialize all params to 0?
- Why do we need non-linearity?
- Why does Relu make training faster?