



COMP [56]630– Machine Learning

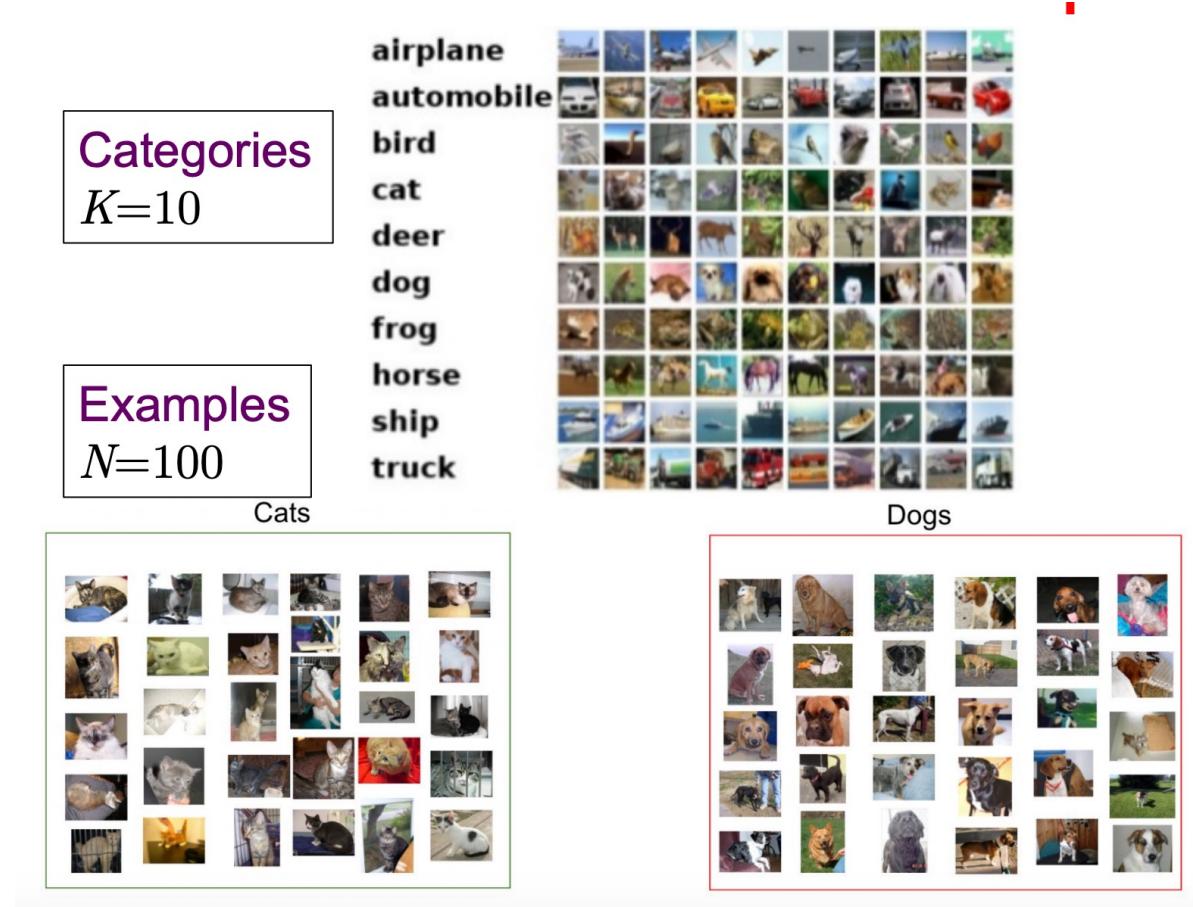
Lecture 8 – Multiclass LR, Evaluation Metrics



Multi-class problems

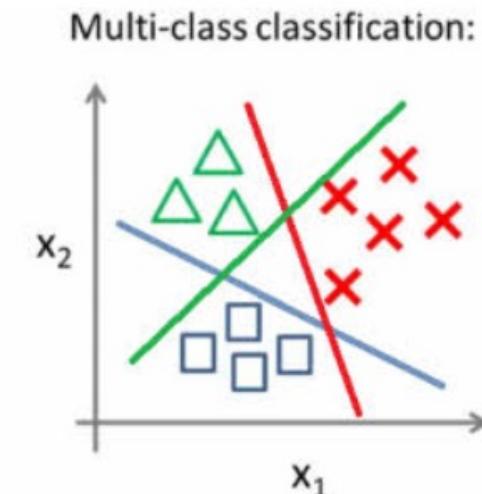
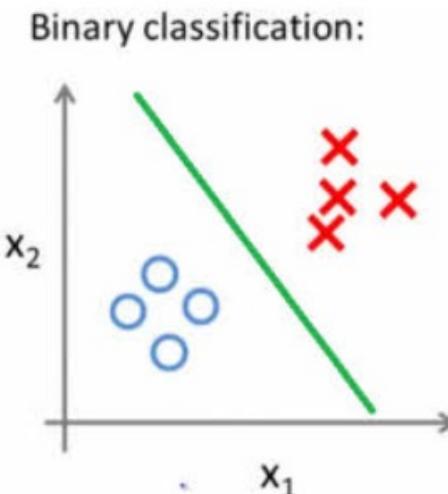


Multi-class classification



Binary vs Multi-class Classification

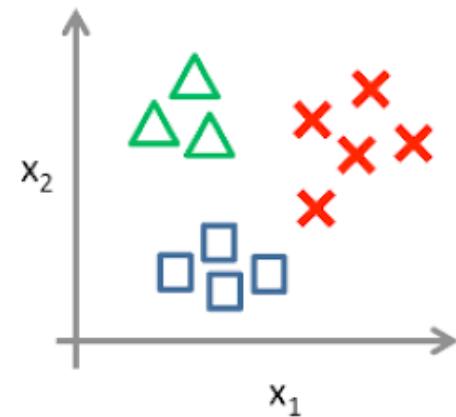
- More than two classes (binary) in our problem setting.
- Logistic regression returns a confidence score that scores the probability of whether an example is positive class or negative class.
- Binary classification requires only one classifier model.



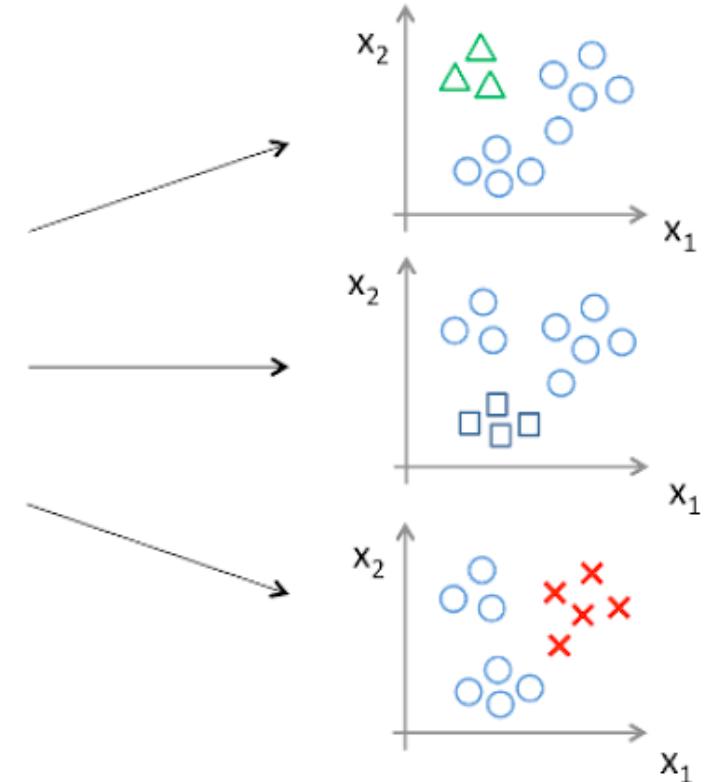
One-vs-all

- Create the N-binary classifier models, one for each class in the data.
- The number of class labels present in the dataset and the number of generated binary classifiers must be the same.
- Also called one-vs-rest classification

One-vs-all (one-vs-rest):

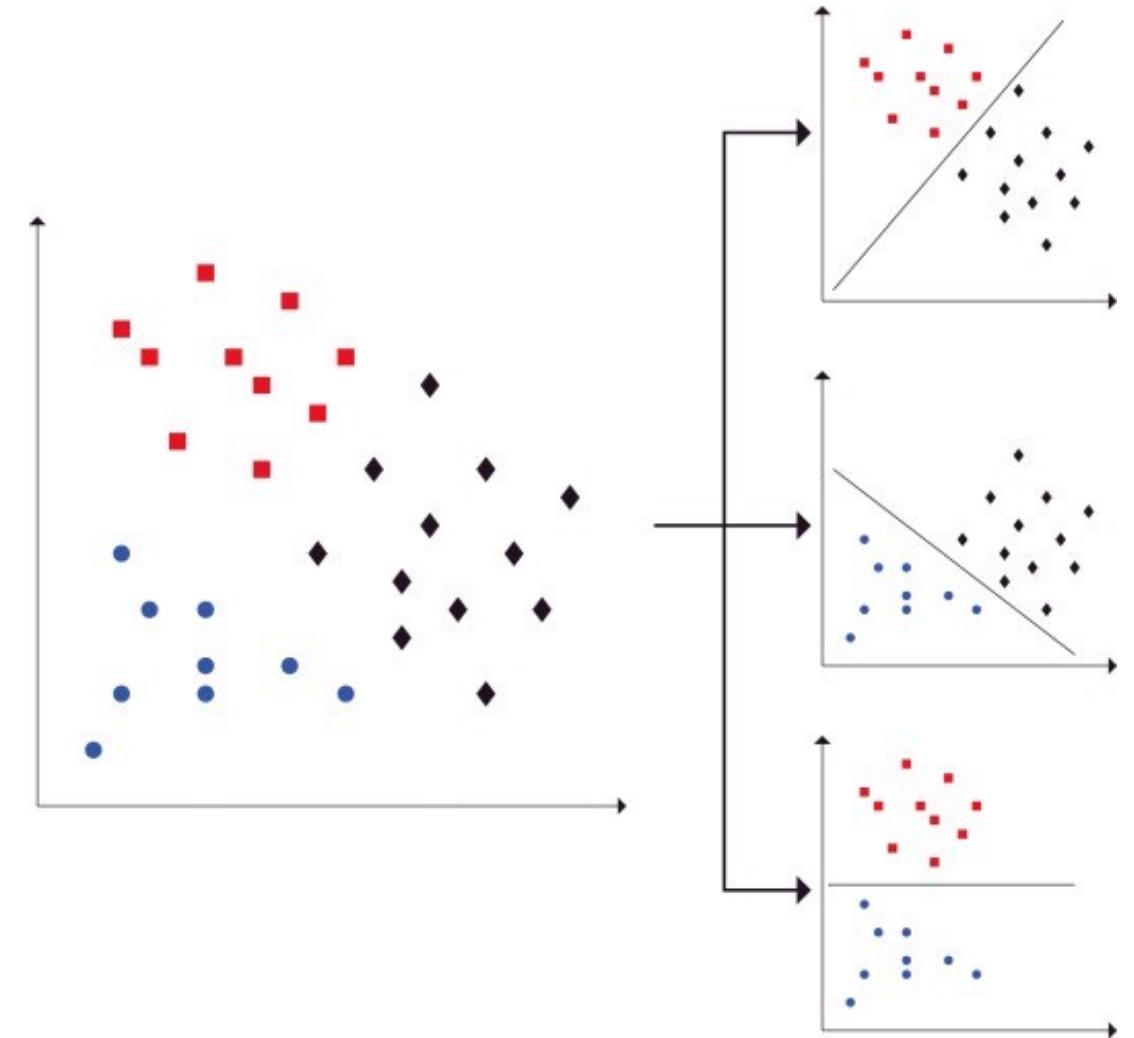


Class 1: **Green**
 Class 2: **Blue**
 Class 3: **Red**



One-vs-one

- Create the $N * (N-1)/2$ binary classifier models, one for each combination of classes in the data.
- Let us say there are three classes, A, B and C.
 - Create classifier for A-vs-B, A-vs-C, B-vs-C.
 - Each binary classifier predicts one class label.
 - Take the class with majority votes as final prediction





In the two-class case $p(C_1|\phi) = y(\phi) = \sigma(w^T\phi + b)$

where $\phi = [\phi_1, \dots, \phi_M]^T$, $w = [w_1, \dots, w_M]^T$ and $a = w^T\phi + b$ is the “activation”

For K classes, we work with soft-max function instead of logistic sigmoid (*Softmax regression*)

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where $a_k = w_k^T \phi + b_k$, $k = 1, \dots, K$
 $w_k = [w_{k1}, \dots, w_{kM}]^T$ and $a = \{a_1, \dots, a_K\}$

- We learn a set of K weight vectors $\{w_1, \dots, w_K\}$ and biases b

Arranging weight vectors as a matrix W

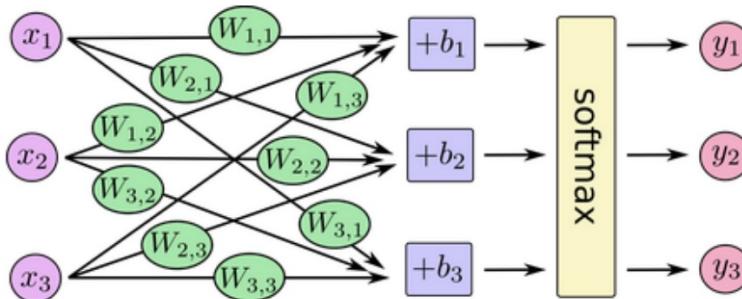
$$\mathbf{a} = W^T \phi + \mathbf{b}$$

$$W = \begin{bmatrix} \mathbf{w}_1 \\ \cdot \\ \cdot \\ \mathbf{w}_K \end{bmatrix} = \begin{bmatrix} w_{11} & w_{1M} \\ \cdot & \cdot \\ \cdot & \cdot \\ w_{K1} & w_{KM} \end{bmatrix}$$

$$\mathbf{y} = \text{softmax}(\mathbf{a})$$
$$y_i = \frac{\exp(a_i)}{\sum_{j=1}^3 \exp(a_j)}$$



Multi-class Logistic Regression



$$\mathbf{a} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\mathbf{W} = \begin{bmatrix} \mathbf{W}_1 \\ \mathbf{W}_2 \\ \mathbf{W}_3 \end{bmatrix} = \begin{bmatrix} \mathbf{W}_{1,1} & \mathbf{W}_{1,2} & \mathbf{W}_{1,3} \\ \mathbf{W}_{2,1} & \mathbf{W}_{2,2} & \mathbf{W}_{2,3} \\ \mathbf{W}_{3,1} & \mathbf{W}_{3,2} & \mathbf{W}_{3,3} \end{bmatrix}$$

$$\mathbf{y} = \text{softmax}(\mathbf{a})$$

$$y_i = \frac{\exp(a_i)}{\sum_{j=1}^3 \exp(a_j)}$$

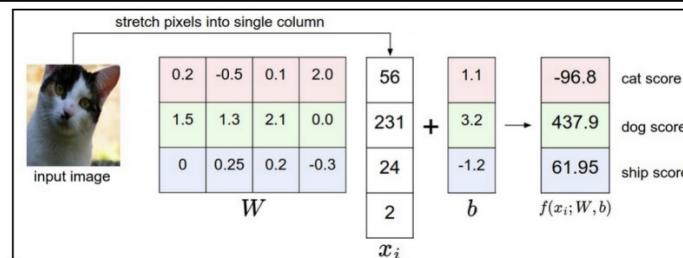
Network Computes

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \text{softmax} \begin{bmatrix} W_{1,1}x_1 + W_{1,2}x_2 + W_{1,3}x_3 + b_1 \\ W_{2,1}x_1 + W_{2,2}x_2 + W_{2,3}x_3 + b_2 \\ W_{3,1}x_1 + W_{3,2}x_2 + W_{3,3}x_3 + b_3 \end{bmatrix}$$

In matrix multiplication notation

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \text{softmax} \left(\begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \right)$$

An example





We use maximum likelihood to determine the parameters $\{w_k\}$, $k=1,..K$

The exp within softmax works very well when training using log-likelihood

$$\text{softmax}(\mathbf{a})_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

- Log-likelihood can undo the exp of softmax

$$\log \text{softmax}(\mathbf{a})_i = a_i - \log \sum_j \exp(a_j)$$

- Input a_i always has a direct contribution to cost
 - Because this term cannot saturate, learning can proceed even if second term becomes very small
- First term encourages a_i to be pushed up
- Second term encourages all a to be pushed down



Derivatives

The multiclass logistic regression model is

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

For maximum likelihood we will need the derivatives of y_k wrt all of the activations a_j

These are given by

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$

– where I_{kj} are the elements of the identity matrix



One-hot targets

▪

Classes $C_1,..C_K$ represented by 1-of- K scheme

– One-hot vector:

- class C_k is a K -dim vector or $[t_1,..,t_K]^T$, $t_i \in \{0,1\}$
 - With $K=6$, class C_3 is $(0,0,1,0,0,0)^T$ with $t_1=t_2=t_4=t_5=t_6=0$, & $t_3=1$

– The class probabilities obey

$$\sum_{k=1}^K p(C_k) = \sum_{k=1}^K t_k = 1$$

– If $p(t_k=1) = \mu_k$ then

$$p(C_k) = \prod_{k=1}^K \mu_k^{t_k} \text{ where } \boldsymbol{\mu} = (\mu_1,..,\mu_K)^T$$

e.g., probability of C_3 is

$$p([0,0,1,0,0,0]) = \mu_3$$

Why use one-hot representation?
If we used numerical categories 1,2,3...we would impute ordinality.
We can now use simpler Bernoulli instead of multinoulli



Target Matrix

Classes have values $1, \dots, K$

Each represented as a K -dimensional binary vector

We have N labeled samples

- So instead of target vector t we have a target matrix T

$$T = \begin{bmatrix} t_{11} & \cdot & \cdot & t_{1K} \\ \cdot & & & \\ \cdot & & & \\ t_{N1} & & & t_{NK} \end{bmatrix}$$

Classes →

Samples ↓

Note that t_{nk} corresponds to sample n and class k



Objective Function & Gradient

Likelihood of observations is

$$p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

$$T = \begin{bmatrix} t_{11} & \dots & t_{1K} \\ \vdots & & \vdots \\ t_{N1} & & t_{NK} \end{bmatrix}$$

- Where, for feature vector ϕ_n $y_{nk} = y_k(\phi_n) = \frac{\exp(\mathbf{w}_k^T \phi_n)}{\sum_j \exp(\mathbf{w}_j^T \phi_n)}$

Objective Function: negative log-likelihood

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

– Known as cross-entropy error for multi-class

Gradient of error function wrt parameter \mathbf{w}_j

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

using $\sum_k t_{nk} = 1$

Error x Feature Vector

$$y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \mathbf{w}_k^T \phi$$

$\frac{\partial y_k}{\partial a_j} = y_k (I_{kj} - y_j)$ where I_{kj} are elements of the identity matrix



Performance of Classifiers



Confusion Matrix

- Used to **describe the performance of a classification model** (or "classifier") on a set of test data for which the true values are known.



Confusion Matrix

- Used to **describe the performance of a classification model** (or "classifier") on a set of test data for which the true values are known.

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN



Confusion Matrix

- Used to **describe the performance of a classification model** (or "classifier") on a set of test data for which the true values are known.
- Models Four major characteristics
- **true positives (TP)**
- **true negatives (TN)**
- **false positives (FP)**
- **false negatives (FN)**

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN



Perf. Measures from Confusion Matrix

- **Accuracy:** Overall, how often is the classifier correct?
 - $(TP+TN)/\text{total}$
- **Misclassification Rate:** Overall, how often is it wrong?
 - $(FP+FN)/\text{total}$
- **Recall:** When it's actually yes, how often does it predict yes?
 - also known as "Sensitivity"

$$\textit{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- **Precision:** When it predicts yes, how often is

$$\textit{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$



Perceptron Algorithm



The Perceptron

- Consider modifying the logistic regression method to “force” it to output values that are either 0 or 1 or exactly. To do so, it seems natural to change the definition of g to be the threshold function:

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

- If we then let $h(x) = g(\theta^T x)$ as before but using this modified definition of g , and if we use the update rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$



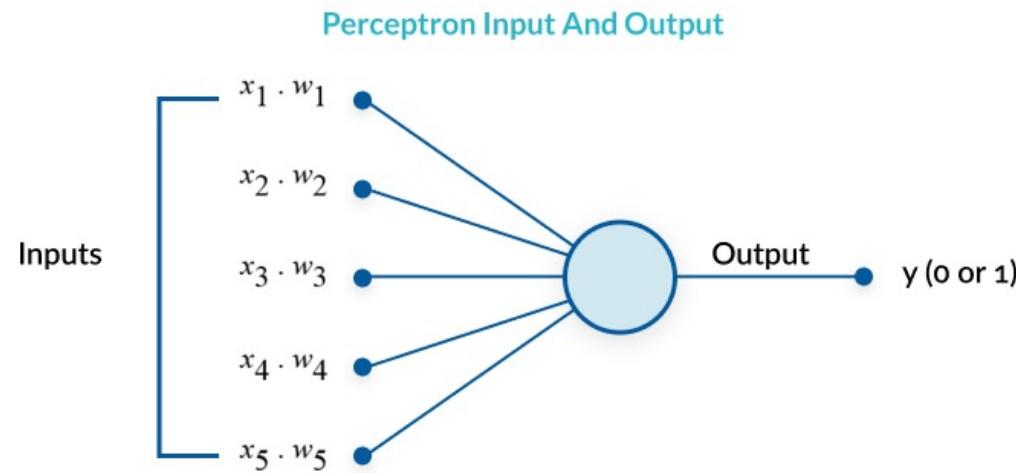
Perceptron

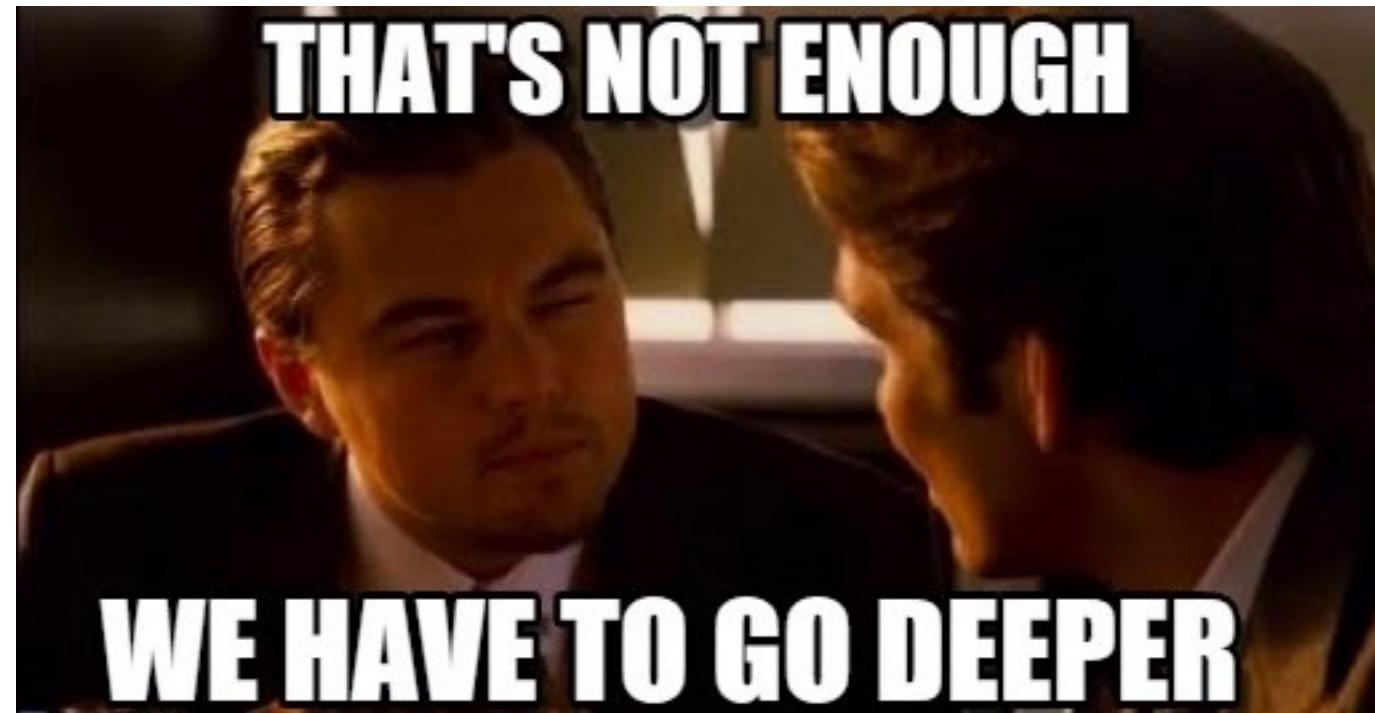
- In the 1960s, this “perceptron” was argued to be a rough model for how individual neurons in the brain work.
- Note however that even though the perceptron may be cosmetically similar to the other algorithms we talked about, it is actually a very different type of algorithm than logistic regression.
 - The hypothesis in logistic regression provides a measure of uncertainty in the occurrence of a binary outcome based on a linear model.
 - The output from a step function can of course not be interpreted as any kind of probability.
 - Since a step function is not differentiable, it is not possible to train a perceptron using the same algorithms that are used for logistic regression.



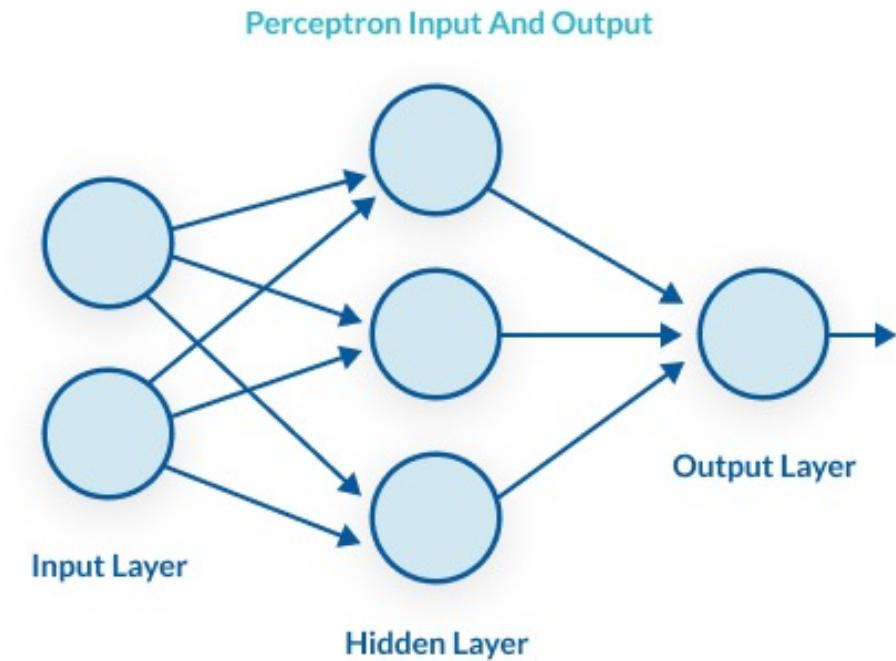
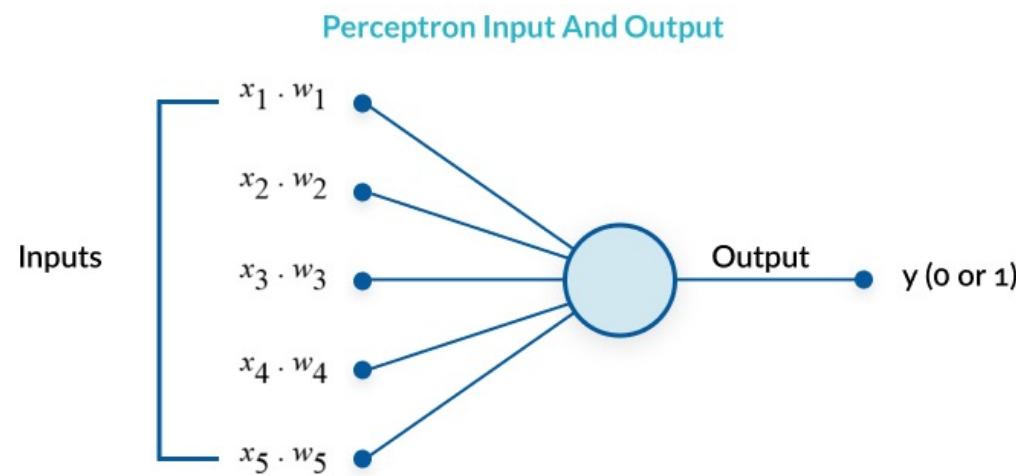
Multi-layer perceptron

Perceptron vs Multilayer Perceptron (MLP)

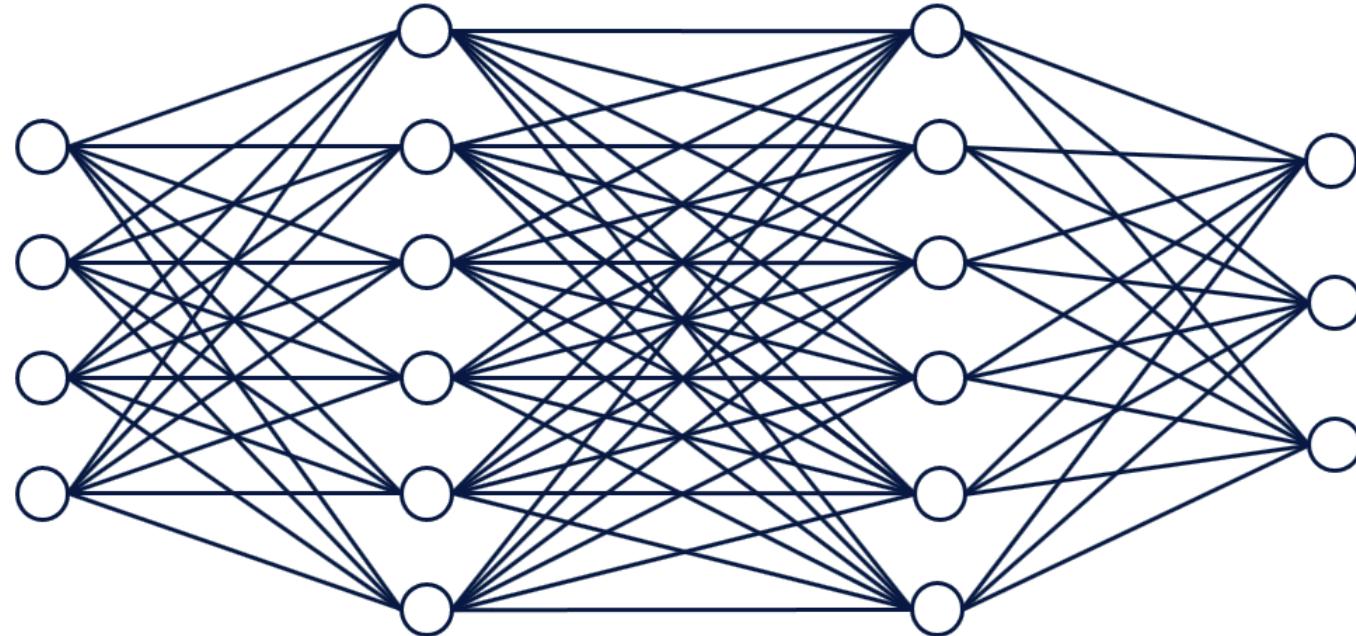




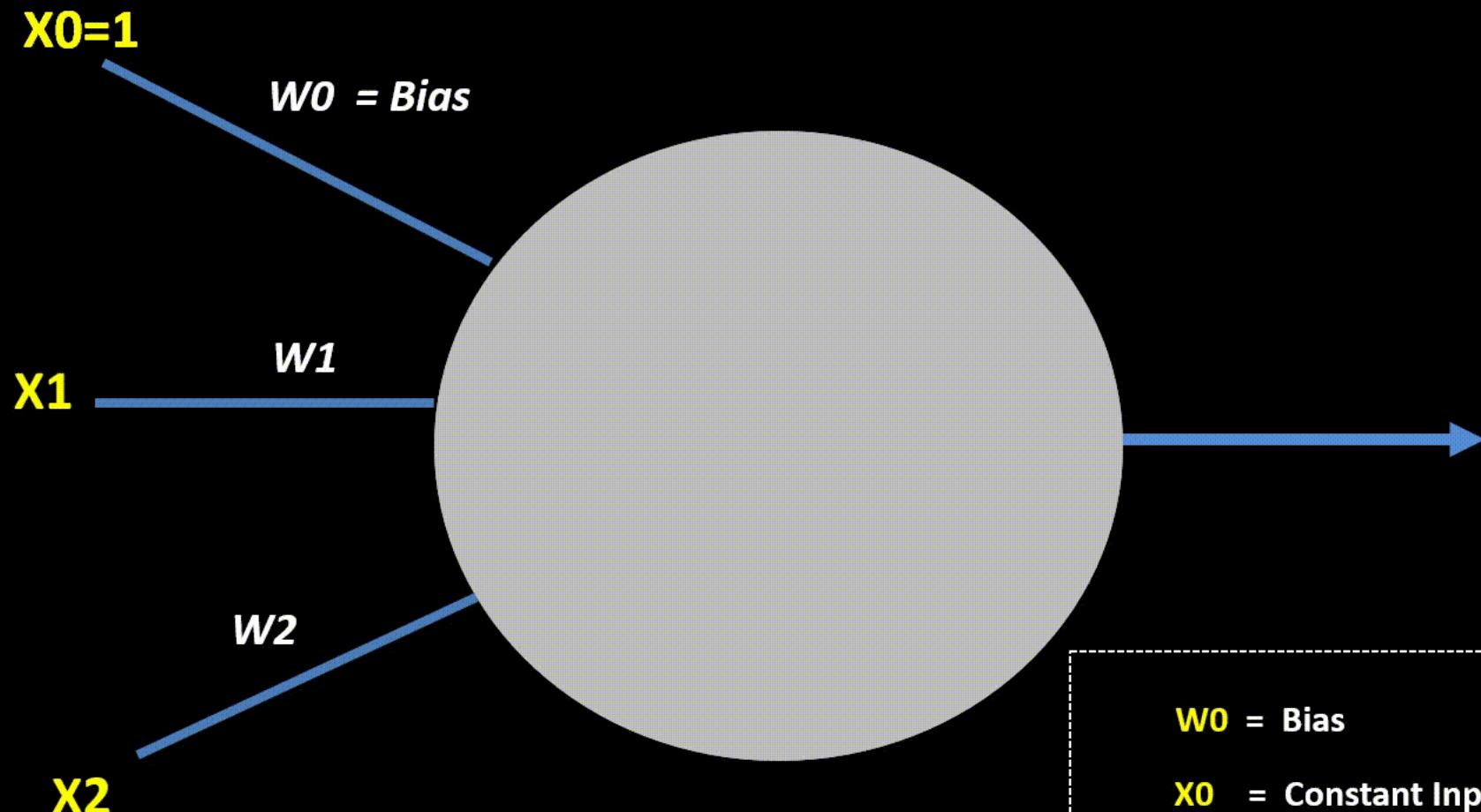
Perceptron vs Multilayer Perceptron (MLP)



Why do we need more layers?

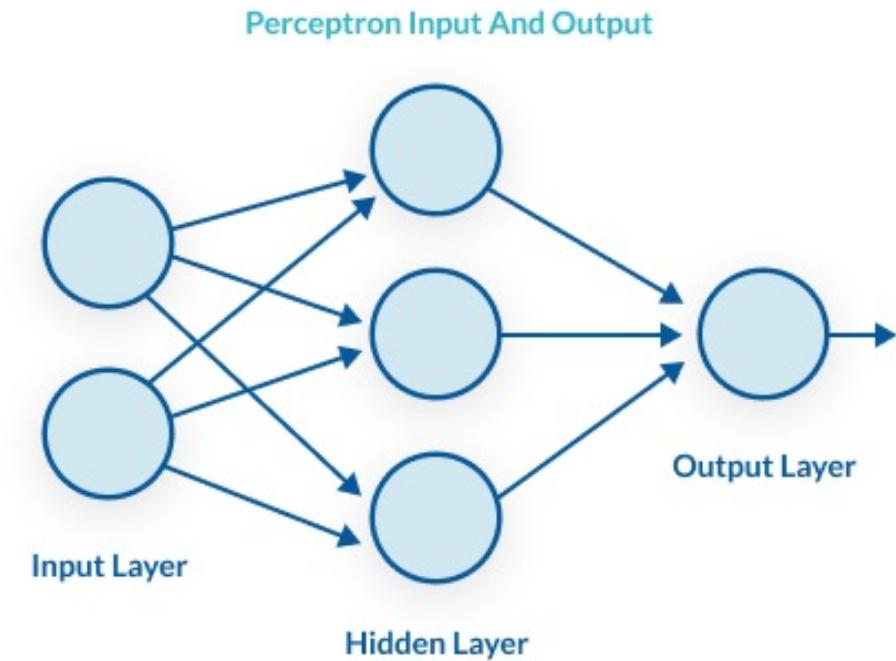
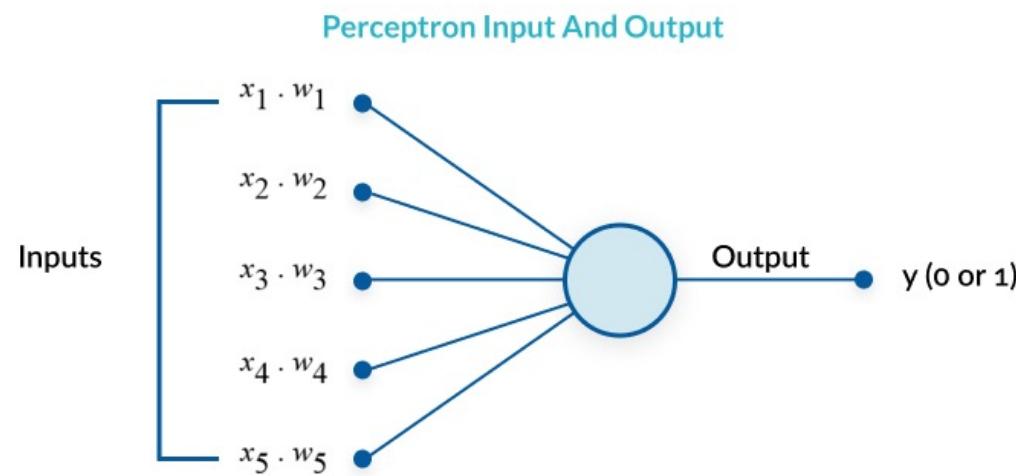


Artificial Neuron



W0 = Bias
**X0 = Constant Input 1
for Bias**
X1,X2 = Attribute Inputs
**W1,W2 = Weights of Inputs
X1,X2**

Perceptron vs Multilayer Perceptron (MLP)





Feed forward Neural Networks

A neural network can also be represented similar to linear models but basis functions are generalized

$$y(\mathbf{x}, \mathbf{w}) = f\left(\sum_{j=1}^M w_j \phi_j(\mathbf{x})\right)$$

activation function

For regression: identity function

For classification: a nonlinear function

Coefficients w_j adjusted during training

There can be several activation functions

Basis functions

$\phi_j(\mathbf{x})$ a nonlinear function of a linear combination of D inputs

its parameters are adjusted during training



Activation Functions

- Each activation a_j is transformed using differentiable nonlinear activation functions $z_j = h(a_j)$
- The z_j correspond to outputs of basis functions $\varphi_j(x)$
 - or first layer of network or hidden units
- Nonlinear functions h
- Three examples of activation functions:

1. Logistic sigmoid

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

2. Hyperbolic tangent

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

3. Rectified Linear unit

$$f(x) = \max(0, x)$$



Activation Function

- Determined by the nature of the data and the assumed distribution of the target variables
- For standard regression problems the activation function is the identity function so that $y_k = a_k$
- For multiple binary classification problems, each output unit activation is transformed using a logistic sigmoid function so that $y_k = \sigma(a_k)$
- For multiclass problems, a softmax activation function of the form:

$$\frac{\exp(a_k)}{\sum_j \exp(a_j)}$$



Visualization of activation functions

Identity	Sigmoid	TanH	ArcTan
ReLU	Leaky ReLU	Randomized ReLU	Parameteric ReLU
Binary	Exponential Linear Unit	Soft Sign	Inverse Square Root Unit (ISRU)
Inverse Square Root Linear	Square Non-Linearity	Bipolar ReLU	Soft Plus



Forward Pass

- Output of a layer l is given by

$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}$$

- Where W^l is the weights of the layer, b^l is the bias and a^l is the activation

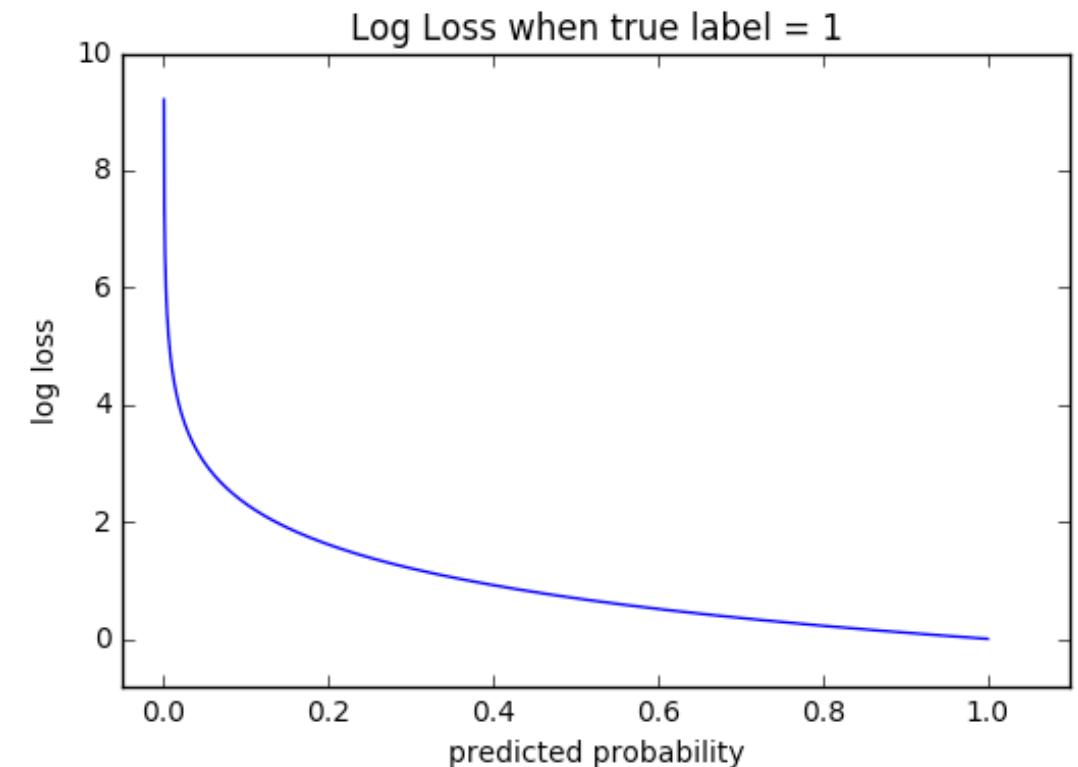
$$a^{[l]} = g^{[l]}(z^{[l]})$$

Loss Function: Log loss

- Measures the performance of a classification model whose output is a probability value between 0 and 1.

$$-(y \log(p) + (1 - y) \log(1 - p))$$

- Cross-entropy loss increases as the predicted probability diverges from the actual label.





How to train this network?

- Gradient Descent + Backprop!
- What is Backprop?
- Backprop or Backpropagation is a way to train multilayer neural networks using gradients.



Backpropagation

- Before training the neural network, select an initial value for parameters.
 - Common practice: randomly initialize the parameters to small values (e.g., normally distributed around zero; $N(0; 0:1)$)
- Next step: update the parameters.
- After a single forward pass through the neural network, the output will be a predicted value
- We can then compute the loss L , in our case the log loss

$$\mathcal{L}(\hat{y}, y) = - \left[(1 - y) \log(1 - \hat{y}) + y \log \hat{y} \right]$$



Defining Variables

Let X be the input features $[n \times f]$
 y be the output(target) labels $[n \times 1]$

Define the weights [parameters] of the neural network

First layer $\rightarrow W_1$, bias $\rightarrow b_1$

Second layer $\rightarrow W_2$, bias $\rightarrow b_2$

Third layer $\rightarrow W_3$, bias $\rightarrow b_3$



The forward pass

The output of layer 1 is

$$z_1 = w_1 \cdot x + b_1$$

$a_1 = g(z_1)$ where "g" is the activation function
The output of layer 2 is

$$z_2 = w_2 \cdot a_1 + b_2$$

$$a_2 = g(z_2)$$

The output of layer 3 is

$$z_3 = w_3 \cdot a_2 + b_3$$

$$a_3 = g(z_3)$$

The o/p of the neural network is

$$\hat{y} = a_3$$



Defining the loss

Task : Binary classification

⇒ loss is log loss or binary cross entropy

$$L = -[(1-y) \cdot \log(1-\hat{y}) + y \cdot \log(\hat{y})]$$

Learning with gradient descent



To learn the parameters

- ① Provide random values for weights w_1, w_2, w_3 & biases b_1, b_2, b_3
- ② Update the weights using gradient descent by

$$w_i := w_i - \alpha \cdot \frac{\partial L}{\partial w_i}$$

$$b_i := b_i - \alpha \cdot \frac{\partial L}{\partial b_i}$$

where i refers to the i th layer.

- ③ repeat until convergence



Finding dL/dW_3

To find $\frac{\partial L}{\partial w_3}$ to update the third layer.

$$\frac{\partial L}{\partial w_3} = \frac{\partial}{\partial w_3} \left[-[(1-y) \cdot \log(1-\hat{y}) + y \cdot \log(\hat{y})] \right]$$

$$= - \left[(1-y) \cdot \frac{\partial}{\partial w_3} (\log(1-\hat{y})) + y \cdot \frac{\partial}{\partial w_3} (\log(\hat{y})) \right]$$



Finding dL/dW_3

Remember: $\frac{d}{dz}(\log(z)) = \frac{1}{z}$



Finding dL/dW_3

$$\therefore \frac{\partial L}{\partial w_3} = - \left[\frac{(1-y)}{(1-\hat{y})} \cdot \frac{\partial}{\partial w_3} (1-\hat{y}) + \frac{y}{\hat{y}} \cdot \frac{\partial}{\partial w_3} (\hat{y}) \right]$$

$$= - \left[\frac{-(1-y)}{(1-\hat{y})} \cdot \frac{\partial \hat{y}}{\partial w_3} + \frac{y}{\hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_3} \right]$$

$$= \frac{(1-y)}{(1-\hat{y})} \cdot \frac{\partial \hat{y}}{\partial w_3} - \frac{y}{\hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_3}$$

$$= \left[\frac{(1-y)}{(1-\hat{y})} - \frac{y}{\hat{y}} \right] \cdot \frac{\partial \hat{y}}{\partial w_3}$$



Finding dL/dW_3

$$\Rightarrow \frac{\partial L}{\partial w_3} = \frac{\hat{y}(1-\hat{y}) - y(1-\hat{y})}{\hat{y}(1-\hat{y})} \cdot \frac{\partial \hat{y}}{\partial w_3}$$

$$= \frac{\hat{y} - y\hat{y} - y + y\hat{y}}{\hat{y}(1-\hat{y})} \cdot \frac{\partial \hat{y}}{\partial w_3}$$

$$= \frac{\hat{y} - y}{\hat{y}(1-\hat{y})} \cdot \frac{\partial \hat{y}}{\partial w_3}$$



Finding dL/dW_3

we know that $\hat{y} = a_3 = g(z_3)$



Finding dL/dW_3

if $g(z)$ is a sigmoid function,

$$g'(z) = g(z) \cdot (1-g(z))$$

$$\therefore \frac{\partial \hat{y}}{\partial w_3} = \frac{\partial a_3}{\partial w_3} = \frac{\partial \cdot g(z_3)}{\partial w_3}$$

$$= g(z_3) \cdot (1-g(z_3)) \cdot \frac{\partial z_3}{\partial w_3}$$

$$= a_3 (1-a_3) \cdot \frac{\partial}{\partial w_3} [w_3 a_2 + b_3]$$

$$= a_3 (1-a_3) \cdot a_2$$

$$= \hat{y} (1-\hat{y}) \cdot a_2$$



Finding dL/dW_3

$$\therefore \frac{\partial L}{\partial w_3} = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})} \cdot [\hat{y} \cdot (1-\hat{y}) \cdot a_2]$$

$$\therefore \boxed{\frac{\partial L}{\partial w_3} = (\hat{y} - y) \cdot a_2}$$

$$\boxed{\frac{\partial L}{\partial w_3} = (a_3 - y) \cdot a_2^T}$$

when providing a vectorial
solution.



Finding dL/dW_3

$$\therefore \frac{\partial L}{\partial w_3} = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})} \cdot [\hat{y} \cdot (1-\hat{y}) \cdot a_2]$$

$$\therefore \boxed{\frac{\partial L}{\partial w_3} = (\hat{y} - y) \cdot a_2}$$

$$\boxed{\frac{\partial L}{\partial w_3} = (a_3 - y) \cdot a_2^T}$$

when providing a vectorized
solution.

Similarly,

$$\boxed{\frac{\partial L}{\partial b_3} = a_3 - y}$$



Finding dL/dW_2

- We need to use the chain rule from calculus.
- Why?
- There is no direct relationship between the weights W_2 and the loss L .
- You cannot differentiate a variable by another if there is no direct relationship
 - Else it would be 0
 - Not true if the variable/function is composite



Finding dL/dW_2

We know that Loss is dependent on $\hat{y} = a_3$

$$\Rightarrow \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial ?} \cdot \frac{?}{\partial w_2}$$

We know that $a_3 = g(z_3)$

$\Rightarrow a_3$ is dependent on z_3

$$\Rightarrow \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial ?} \cdot \frac{?}{\partial w_2}$$



Finding dL/dW_2

We know that $Z_3 = W_3 a_2 + b_3$

$\Rightarrow Z_3$ is dependent on a_2

$$\Rightarrow \frac{\delta L}{\delta w_2} = \frac{\delta L}{\delta a_3} \cdot \frac{\delta a_3}{\delta Z_3} \cdot \frac{\delta Z_3}{\delta a_2} \cdot \underbrace{\frac{\delta a_2}{?}}_{?} \cdot \frac{?}{\delta w_2}$$



Finding dL/dW_2

But $a_2 = g(z_2)$

$$\Rightarrow \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \cdot ?$$

But $z_2 = w_2 \cdot a_1 + b_2$

$$\therefore \boxed{\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}}$$



Finding dL/dW_2

If we rewrite $\frac{\partial L}{\partial w_3}$ using chain rule,

$$\frac{\partial L}{\partial w_3} = \underbrace{\frac{\partial L}{\partial a_3}}_{\cdot a_3 - y} \cdot \underbrace{\frac{\partial a_3}{\partial z_3}}_{a_2^T} \cdot \underbrace{\frac{\partial z_3}{\partial w_3}}_{a_2^T}$$

Since
 $z_3 = w_3 a_2 + b_3$



Finding dL/dW_2

$$\Rightarrow \frac{\partial L}{\partial w_2} = (a_3 - y) \cdot \underbrace{\frac{\partial z_3}{\partial a_2}}_{\frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}}$$

$$\text{if } a_2 = g(z_2)$$

$$\text{then } \frac{\partial a_2}{\partial z_2} = g'(z_2)$$

$$\text{If } z_3 = w_3 \cdot a_2 + b_3$$

$$\text{Then } \frac{\partial z_3}{\partial a_2} = w_3$$

$$\text{If } z_2 = w_2 \cdot a_1 + b_2$$

$$\text{then } \frac{\partial z_2}{\partial w_2} = a_1$$



Finding dL/dW_2

$$\Rightarrow \frac{\partial L}{\partial w_2} = (a_3 - y) \cdot w_3 \cdot g'(z_2) \cdot q_1$$

Re-arranging for vectorization,

$$\frac{\partial L}{\partial w_2} = w_3^T \cdot g'(z_2) (a_3 - y) \cdot q_1^T$$

$$\frac{\partial L}{\partial b_2} = w_3^T \cdot g'(z_2) \cdot (a_3 - y)$$



Finding dL/dW_1

Use the chain rule!

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\boxed{\frac{\partial L}{\partial w_1} = (a_3 - y) \cdot w_3 \cdot g'(z_2) \cdot w_2 \cdot g'(z_1) \cdot x}$$

$$\boxed{III^{(y)} \quad \frac{\partial L}{\partial b_1} = (a_3 - y) \cdot w_3 \cdot g'(z_2) \cdot w_2 \cdot g'(z_1)}$$