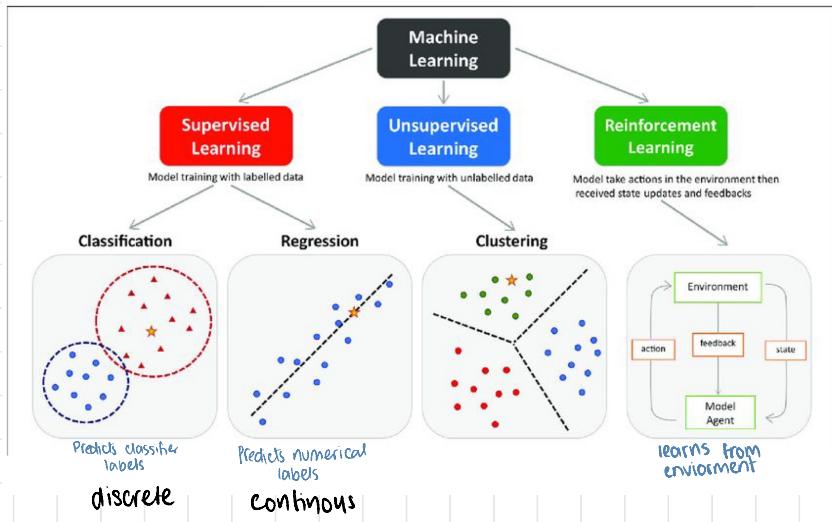
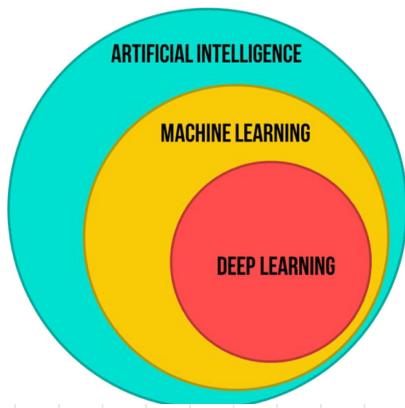


Deep Learning for Perception

MACHINE LEARNING RECAP

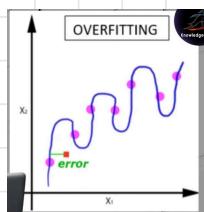


- D. Is it possible to use ReLU activation in the last layer for (i) classification and (ii) Clustering?
For unsupervised learning, yes. For classification NO

MACHINE LEARNING TERMS

OVERTFITTING

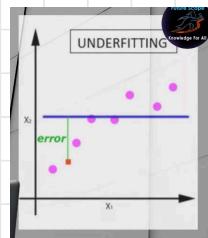
- ↳ memorises noise instead of patterns
- ↳ learns training data too well
- ↳ training accuracy high
- ↳ testing accuracy too low
- ↳ doesn't work well on unseen data



Learns more features than req.
Night train diameter and material of chair
So won't classify other chairs as a chair

UNDERFITTING

- ↳ model too simple to capture underlying trends
- ↳ training accuracy too low
- ↳ high training error



might only learn 1 feature (4 legs)
So will classify dog as a chair too as both 4 legs

- ↳ when less sample size

Training

- ↳ training weights
- ↳ changes your parameters

↳ weights
↳ biasness

Cross Validation

- ↳ avoids over and under fitting
- ↳ tests on different folds
- ↳ reduces biases by working on different splits in data
- ↳ averages out your accuracy
- ↳ gives validation accuracy

Validation

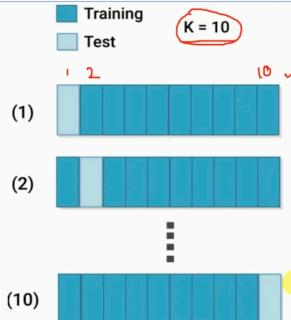
- ↳ used to optimize hyperparameters

epoch activation function

Testing

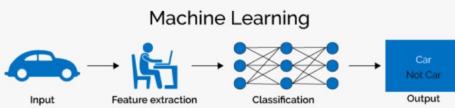
- ↳

K-Fold Cross Validation in Machine Learning



Machine Learning

- ↳ Input
- ↳ feature extraction → done manually
even more tips → so can have biasness
- ↳ Classification
- ↳ Output



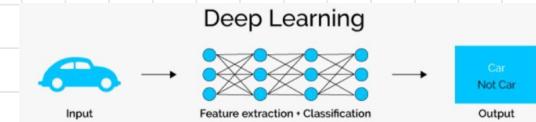
- ↳ makes the pic black and white
- ↳ then identifies edges
- ↳ then fragments → edge detector

- ↳ These are hard engineered features that are
- ↳ time consuming → manually select relevant features
- ↳ brittle → isn't generalized features will fail when data distribution changes
- ↳ not scalable in practice → when dataset grows manually selecting becomes infeasible

what your eyes can't do
Deep learning model can't do either

Deep Learning

- ↳ Input
- ↳ feature extraction + classification
- ↳ Output



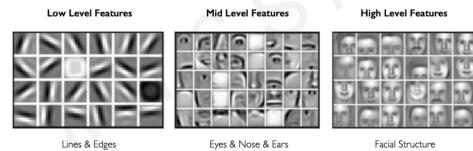
Deep Neural Networks



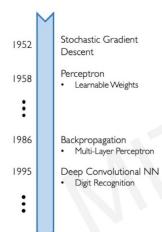
use some cv filters matrix → gradient to generate edge detector

it learns underlying features directly from data automatically

HIERARCHICAL FEATURES



Why Now?



Neural Networks date back decades, so why the resurgence?

1. Big Data

- Larger Datasets
- Easier Collection & Storage



2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



3. Software

- Improved Techniques
- New Models
- Toolboxes

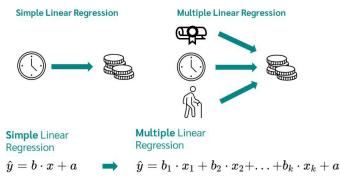


Linear Regression

↳ Output is continuous $-\infty \rightarrow +\infty$

↳ Predicts outcomes by learning

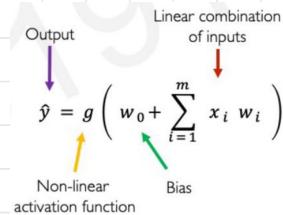
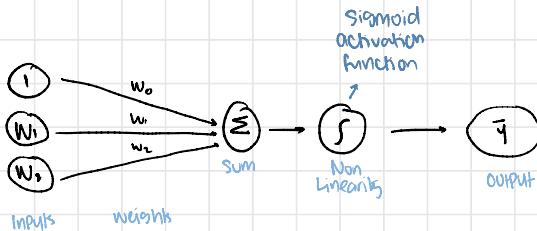
Weights and outcomes through Gradient Descent



Logistic Regression

↳ Using activation function in linear regression

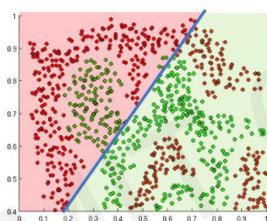
Forward Propagation



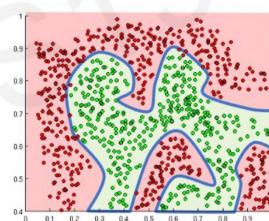
Activation Function

↳ used to introduce non-linearity

↳ chooses which nodes should be activated



Linear activation functions produce linear decisions no matter the network size

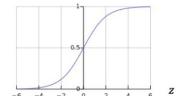


Activation Functions

$$\hat{y} = g(w_0 + X^T W)$$

- Example: sigmoid function

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



Common Activation Functions

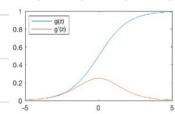
Sigmoid Function

$$\sigma = \frac{1}{1 + e^{-n}}$$

$\xrightarrow{-\infty} 0 \xrightarrow{\text{class 0}}$
 $\xrightarrow{0} 0.5 \xrightarrow{\text{class 1}}$

$$\delta = a(1-a)$$

\downarrow derivative



$$g(z) = \frac{1}{1 + e^{-z}}$$

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

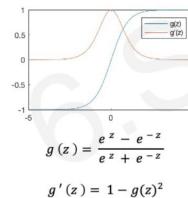
`tf.math.sigmoid(z)`

Hyperbolic Tangent

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

$$\delta = 1 - f(u)^2$$

\downarrow derivative



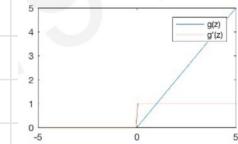
`tf.math.tanh(z)`

ReLU

↳ makes -ve value 0

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases} \rightarrow \text{derivative}$$



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

`tf.nn.relu(z)`

↳ for unsupervised

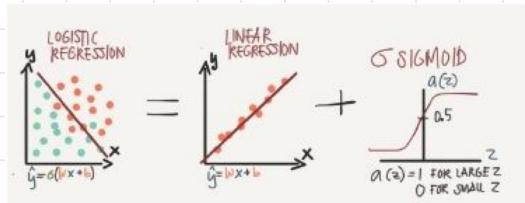
3. Is it possible to use ReLU in the last layer?

Ans: For unsupervised learning, yes. For classification NO

LOGISTIC REGRESSION

Let's outcome variable is binary

↳ basically activation function in linear regression



Steps to learn Parameters

1. Loss function / cost function

↳ measures how well models

Predictions \hat{y} match the actual y

↳ gives LOSS SCORE that model tries to minimize during training

2. Optimizer algorithm

↳ uses loss score to

↳ iteratively update weights

LOSS FUNCTION OR COST FUNCTION

Conditional maximum likelihood estimation prefers the correct class labels of the training examples to be more likely.

We choose the parameters w , b that maximize the log probability of the true y labels in the training data given the observations x .

LOSS FUNCTION

if $y=1 \rightarrow p(y|u) = \hat{y}$

if $y=0 \rightarrow p(y|u) = 1-\hat{y}$

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

LOSS FUNCTION TYPES

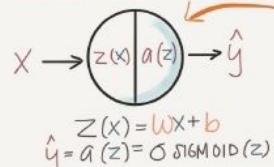
1. Cross Entropy multiple class classification

2. Binary Cross Entropy logistic regression, binary classification

3. MSE regression tasks, numerical labels

- C. Briefly discuss the training and testing phase of Logistic Regression.
training: we train the system (specifically the weights w and b) using stochastic gradient descent and the cross-entropy loss.
test: Given a test example x we compute $p(y|x)$ and return the higher probability label $y = 1$ or $y = 0$.

PUTTING IT ALL TOGETHER



1. Forward Propagation

↳ calculate \hat{y}

2. Backward Propagation

↳ Gradient Descent + update w and b

REPEAT UNTIL IT CONVERGES

MAXIMUM LIKELIHOOD ESTIMATION

$$P(y=1) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

For samples labeled as '1', we try to estimate $(w$ and $b)$ such that the product of all probability $p(x)$ is as close to 1 as possible.

And for samples labeled as '0', we try to estimate $(w$ and $b)$ such that the product of all probability is as close to 0 as possible in other words $(1 - p(x))$ should be as close to 1 as possible.

for samples labelled as 1: $\prod_{s \text{ in } y_i=1} p(x_i)$

for samples labelled as 0: $\prod_{s \text{ in } y_i=0} (1 - p(x_i))$

On combining both, $(w$ and $b)$ parameters should be such that the product of both of these products is maximum over all elements of the dataset.

$$L(\beta) = \prod_{s \text{ in } y_i=1} p(x_i) + \prod_{s \text{ in } y_i=0} (1 - p(x_i))$$

This function is the one we need to optimize and is called the likelihood function.

LOG Likelihood

↳ evaluates how well a model explains observed data

↳ goal is to maximise the log likelihood in training

Since there are only two discrete outcomes (1 or 0), this is a Bernoulli distribution. $\rightarrow 2$ diff outcomes

$$p(y|x) \text{ for one observation: } p(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$

$$\begin{aligned} \log p(y|x) &= \log [\hat{y}^y (1-\hat{y})^{1-y}] \\ &= y \log \hat{y} + (1-y) \log (1-\hat{y}) \end{aligned}$$

whatever values maximize a probability will also maximize the log of the probability.

means minimize loss function like cross entropy

$p(y|x)$ for n observations:

$$L(\beta) = \prod_{i=1}^n \hat{y}^{y_i} (1-\hat{y})^{1-y_i}$$

$$l(\beta) = \sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

where, $l(\beta)$ is log-likelihood

use ln instead of log in calculator

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

Derived from log likelihood
in classification tasks

LOSS FUNCTION TYPES

1. CROSS ENTROPY

↳ measures diff b/w 2 probability distributions

True Distribution (actual labels)
Predicted Distribution

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

Log likelihood should be maximized.

In order to turn this into loss function (something that we need to minimize), we'll just flip the sign.

negative log likelihood loss = cross-entropy loss.

$$LCE(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$$LCE(w, b) = -[y \log \sigma(w \cdot x + b) + (1-y) \log (1 - \sigma(w \cdot x + b))]$$

True label	Predicted Probability	LOSS
1	0.9	$-\log(0.9) = 0.1054$
0	0.2	$-\log(0.8) = 0.223$
1	0.8	$-\log(0.2) = 0.223$
0		$-\log(0.1) = 0.5108$

$$0.1054 + 0.223 + 0.223 + 0.5108$$

sum of losses

$$\frac{1.0624}{4} = 0.2656$$

2. Binary Cross Entropy

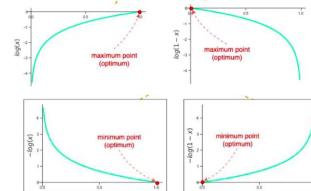
↳ used for binary classification problems

↳ 2 classes $\{0, 1\}$

Binary cross-entropy is often calculated as the average cross-entropy across all data examples.

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

$$LCE(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$$



3. Mean Squared Error (MSE)

- ↳ used in regression tasks
- ↳ measures the avg of squared diff b/w actual and predicted values

$$L = (\text{Actual} - \text{Predicted})^2$$

↳ if more than 1 take avg

WHY NOT MSE?

- ↳ struggles with classification involving probabilities

1. Treats all errors same

e.g. if class 1 is cat and

↳ model Predicts 0.01% confidence → bad prediction

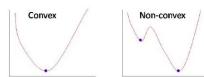
↳ model Predicts 0.99% confidence → not great

BUT MSE will treat both errors similarly even though 0.01 is a much worse prediction

2. Leads to slow and bad convergence

- ↳ MSE is a non-convex function

- ↳ Cross entropy is a convex function



Gradient descent will converge into global minimum only if the function is convex.

1. Actual label: 1, Predicted: 0

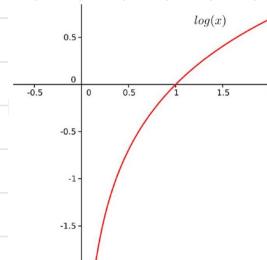
$$\text{MSE} = (1-0)^2 = 1$$

$$\text{Log Loss} = -(1 \log(0) + 0 \log(1)) = \infty \text{ loss}$$

2. Actual label: 1, Predicted: 1

$$\text{MSE} = (1-1)^2 = 0$$

$$\text{Log Loss} = -(1 \log(1) + 0 \log(0)) = 0$$

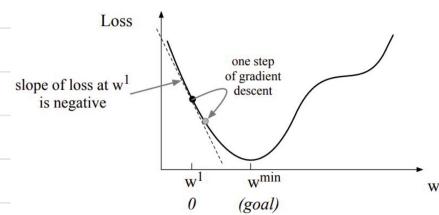
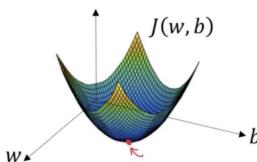


Gradient Descent Algo

- ↳ The Goal is to

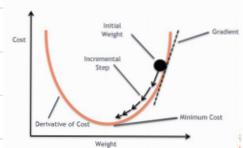
- ↳ minimises loss function by finding optimal weights $\underline{\text{for logistic regression}}$ $\theta = w, b$

Gradient Descent

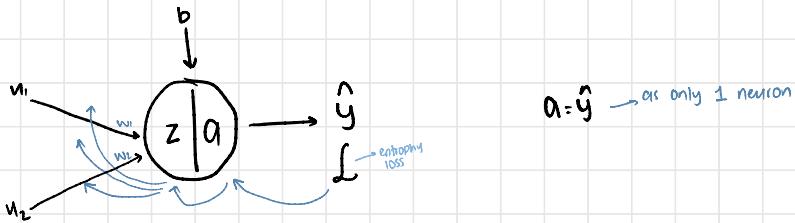


Steps

- ↳ choose starting point
- ↳ calculate gradient



Gradient-Backward Propogation



Entropy Loss

$$J = -[y \log(a) + (1-y) \log(1-a)]$$

$$\frac{\delta L}{\delta a} = -y \times \frac{1}{a} (1) - (1-y) \times \frac{1}{1-a} (-1)$$

$$\frac{\delta L}{\delta a} = \frac{-y}{a} + \frac{1-y}{1-a}$$

Activation Function

$$\hat{y} \leftarrow a = \text{sigmoid}(z)$$

$$a = \frac{1}{1+e^{-z}}$$

Forward Propogation

$$z = w_1 u_1 + w_2 u_2 + b$$

$$\frac{\delta z}{\delta w_1} = u_1$$

$$\frac{\delta a}{\delta z} = a(1-a) \rightarrow \text{derivative given in paper}$$

Updating w_1

$$\begin{aligned} \frac{\delta L}{\delta w_1} &= \frac{\delta L}{\delta a} \cdot \frac{\delta a}{\delta z} \cdot \frac{\delta z}{\delta w_1} \\ &= \left[\frac{-y}{a} + \frac{1-y}{1-a} \right] \times a(1-a) \times u_1 \\ &= \left[\frac{-y(1-a) + a(1-y)}{a(1-a)} \right] \times a(1-a) \times u_1 \\ &= \left[-y + a^2 + a - a^2 \right] \times u_1 \end{aligned}$$

Updating w_2

⇒ same steps only change w_1 to w_2

$$\frac{\delta L}{\delta w_2} = (a-y)u_2$$

$$w_2' = w_2 - \alpha \frac{\delta L}{\delta w_2}$$

$$w_1' = w_1 - \alpha \frac{\delta L}{\delta w_1}$$

$$b' = b - \alpha \frac{\delta L}{\delta b}$$

$$\text{Here, } (a-y) = \frac{\delta L}{\delta z}$$

$$\frac{\delta z}{\delta w_1} = u_1$$

$$\frac{\delta z}{\delta w_2} = u_2$$

$$\text{where, } \frac{\delta L}{\delta b} = (a-y)$$

$$\frac{\delta z}{\delta b} = 1$$

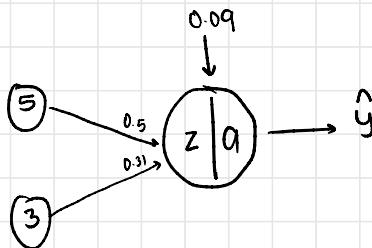
$$\frac{\delta a}{\delta z} = a(1-a)$$

Question

$$w_1: 0.5 \quad u_1: 5 \quad \alpha = 0.01$$

$$w_2: 0.31 \quad u_2: 3 \quad y = 1$$

$$b: 0.09$$



$$Z = w_1 u_1 + w_2 u_2 + b$$

$$= 0.5(5) + 0.31(3) + 0.09$$

$$= 3.52$$

$$a \cdot f(z) = \frac{1}{1+e^{-z}}$$

$$= \frac{1}{1+e^{-3.52}}$$

$$= 0.97$$

$$f = -y \log(a)$$

$$= -1 \log(0.97)$$

$$= 0.0305$$

$$\frac{\delta L}{\delta w_1} = (a-y)u_1$$

$$= -0.15$$

$$\frac{\delta L}{\delta w_2} = (a-y)u_2$$

$$= -0.09$$

$$\frac{\delta L}{\delta b} = a - y$$

$$= -0.03$$

$$w_1^+ = w_1 - \alpha \frac{\delta L}{\delta w_1}$$

$$w_1^+ = 0.5 - (0.01)(-0.15)$$

$$= 0.5015$$

$$w_2^+ = w_2 - \alpha \frac{\delta L}{\delta w_2}$$

$$w_2^+ = 0.31 - (0.01)(-0.09)$$

$$= 0.3109$$

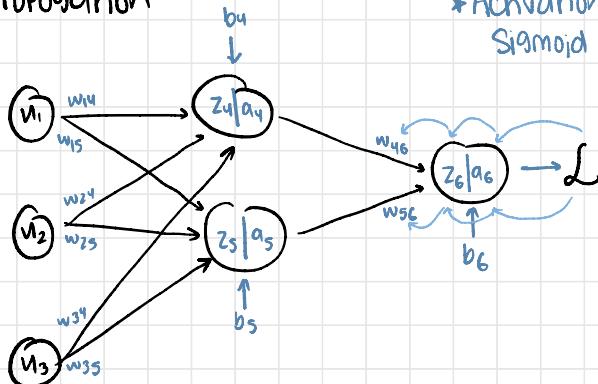
$$b^+ = b - \alpha \frac{\delta L}{\delta b}$$

$$b^+ = 0.09 - (0.01)(-0.03)$$

$$= 0.0903$$

Question → Backpropagation

* Activation function is Sigmoid in all layers



↳ For log loss and cross entropy both

$$L = -[y \log(a) + (1-y) \log(1-a)]$$

$$\frac{\delta L}{\delta a_6} = -[y \log(a_6) + (1-y) \log(1-a_6)]$$

$$= -y \times \frac{1}{a_6} (1) - (1-y) \times \frac{1}{1-a_6} (-1)$$

$$= \frac{-y}{a_6} + \frac{(1-y)}{1-a_6}$$

$$z_6 = a_4 w_{46} + a_5 w_{56} + b_6$$

$$\frac{\delta z_6}{\delta w_{46}}$$

$$\frac{\delta a_6}{\delta z_6} = a_6(1-a_6)$$

→ derivative given in paper

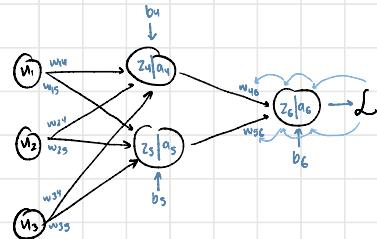
$$\begin{aligned} \frac{\delta L}{\delta w_{46}} &= \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta w_{46}} \\ &= \left[\frac{-y}{a_6} + \frac{(1-y)}{1-a_6} \right] \cdot a_6(1-a_6) \cdot a_4 \\ &= \left[\frac{-y(1-a_6) + a_6(1-y)}{a_6(1-a_6)} \right] \cdot a_6(1-a_6) \cdot a_4 \\ &= [-y + a_6y + a_6 - a_6y] a_4 \end{aligned}$$

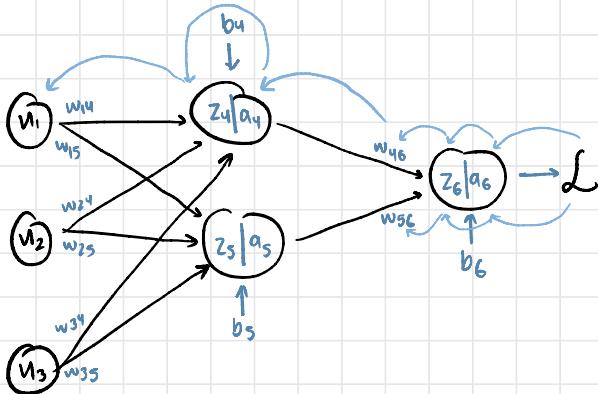
$$\frac{\delta L}{\delta w_{46}} = (a_6 - y) a_4$$

$$\frac{\delta L}{\delta w_{56}} = (a_6 - y) a_5$$

$$-(y-a_6) a_6 (1-a_6) a_4$$

$$\frac{\delta L}{\delta b_6} = a_6 - y$$





L for hidden layer

$$\frac{\delta L}{\delta w_{14}} = \frac{\delta L}{\delta a_4} \cdot \frac{\delta a_4}{\delta z_4} \cdot \frac{\delta z_4}{\delta w_{14}}$$

$\cdot (a_6 - y) w_{46} \cdot a_4(1-a_4) \cdot v_1$

$$\frac{\delta L}{\delta a_4} = \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta a_4}$$

$\cdot (a_6 - y) w_{46}$

Substitute

$$\frac{\delta L}{\delta w_{14}} \cdot (a_6 - y) w_{46} \cdot a_4(1-a_4) \cdot v_1$$

$$\frac{\delta L}{\delta w_{24}} = (a_6 - y) w_{46} \cdot a_4(1-a_4) \cdot v_2$$

$$\frac{\delta L}{\delta w_{34}} = (a_6 - y) w_{46} \cdot a_4(1-a_4) \cdot v_3$$

$$\frac{\delta L}{\delta b_4} = (a_6 - y) w_{46} \cdot a_4(1-a_4)$$

$$\frac{\delta L}{\delta w_{15}} = (a_6 - y) w_{56} \cdot a_5(1-a_5) \cdot v_1$$

$$\frac{\delta L}{\delta w_{25}} = (a_6 - y) w_{56} \cdot a_5(1-a_5) \cdot v_2$$

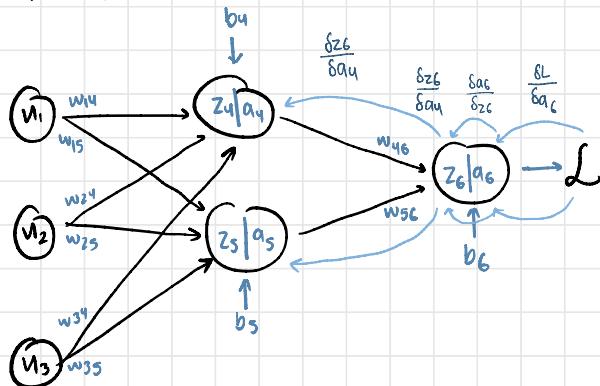
$$\frac{\delta L}{\delta w_{35}} = (a_6 - y) w_{56} \cdot a_5(1-a_5) \cdot v_3$$

$$\frac{\delta L}{\delta b_5} = (a_6 - y) w_{56} \cdot a_5(1-a_5)$$

Question → Backpropagation squared

$$L = \frac{1}{2} (y - a_6)^2$$

$$\begin{aligned}\frac{\delta L}{\delta a_6} &= \frac{1}{2} [2(-1)(y - a_6)] \\ &= -(y - a_6) \\ &= a_6 - y\end{aligned}$$



$$z_4 = w_{14}u_1 + w_{24}u_2 + w_{34}u_3 + b_4$$

α_4 : Sigmoid (z_4)

$$z_5 = w_{15}u_1 + w_{25}u_2 + w_{35}u_3 + b_5$$

α_5 : Sigmoid (z_5)

$$z_6 = w_{46}a_4 + w_{56}a_5 + b_6$$

α_6 : Sigmoid (z_6)

$$L = \frac{1}{2} (y - a_6)^2$$

$$\begin{aligned}\frac{\delta L}{\delta w_{46}} &= \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta w_{46}} \\ &= -(y - a_6) \cdot \alpha_6(1 - \alpha_6) \cdot a_4\end{aligned}$$

$$w_{46}^+ = w_{46} - \eta \frac{\delta L}{\delta w_{46}}$$

→ learning rate

$$L = \frac{1}{2} (y - a_6)^2$$

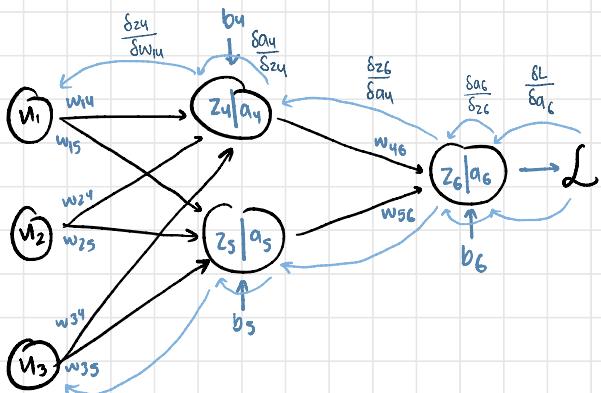
$$\begin{aligned}\frac{\delta L}{\delta w_{56}} &= \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta w_{56}} \\ &= -(y - a_6) \cdot \alpha_6(1 - \alpha_6) \cdot a_5\end{aligned}$$

$$w_{56}^+ = w_{56} - \eta \frac{\delta L}{\delta w_{56}}$$

$$L = \frac{1}{2} (y - a_6)^2$$

$$\begin{aligned}\frac{\delta L}{\delta b_6} &= \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta b_6} \\ &= -(y - a_6) \cdot \alpha_6(1 - \alpha_6)\end{aligned}$$

$$b_6^+ = b_6 - \eta \frac{\delta L}{\delta b_6}$$



↳ For Hidden Layer

$$\frac{\delta L}{\delta w_{14}} = \frac{\delta L}{\delta a_4} \cdot \frac{\delta a_4}{\delta z_4} \cdot \frac{\delta z_4}{\delta w_{14}}$$

$$\frac{\delta L}{\delta a_4} = \frac{\delta L}{\delta z_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta a_4}$$

$$= -(y - a_6) a_6 (1 - a_6) \cdot w_{46}$$

Substitute

$$\frac{\delta L}{\delta w_{14}} = -(y - a_6) a_6 (1 - a_6) \cdot w_{46} \cdot a_4 (1 - a_4) \cdot v_1$$

$$w_{14}^t = w_{14} - \eta \frac{\delta L}{\delta w_{14}}$$

$$\frac{\delta L}{\delta w_{24}} = -(y - a_6) a_6 (1 - a_6) \cdot w_{46} \cdot a_4 (1 - a_4) \cdot v_2$$

$$w_{24}^t = w_{24} - \eta \frac{\delta L}{\delta w_{24}}$$

$$\frac{\delta L}{\delta w_{34}} = -(y - a_6) a_6 (1 - a_6) \cdot w_{46} \cdot a_4 (1 - a_4) \cdot v_3$$

$$w_{34}^t = w_{34} - \eta \frac{\delta L}{\delta w_{34}}$$

$$\frac{\delta L}{\delta b_4} = -(y - a_6) a_6 (1 - a_6) \cdot w_{46} \cdot a_4 (1 - a_4)$$

$$b_4^t = b_4 - \eta \frac{\delta L}{\delta b_4}$$

$$\frac{\delta L}{\delta w_{15}} = -(y - a_6) a_6 (1 - a_6) \cdot w_{56} \cdot a_5 (1 - a_5) \cdot v_1$$

$$w_{15}^t = w_{15} - \eta \frac{\delta L}{\delta w_{15}}$$

$$\frac{\delta L}{\delta w_{25}} = -(y - a_6) a_6 (1 - a_6) \cdot w_{56} \cdot a_5 (1 - a_5) \cdot v_2$$

$$w_{25}^t = w_{25} - \eta \frac{\delta L}{\delta w_{25}}$$

$$\frac{\delta L}{\delta w_{35}} = -(y - a_6) a_6 (1 - a_6) \cdot w_{56} \cdot a_5 (1 - a_5) \cdot v_3$$

$$w_{35}^t = w_{35} - \eta \frac{\delta L}{\delta w_{35}}$$

$$\frac{\delta L}{\delta b_5} = -(y - a_6) a_6 (1 - a_6) \cdot w_{56} \cdot a_5 (1 - a_5)$$

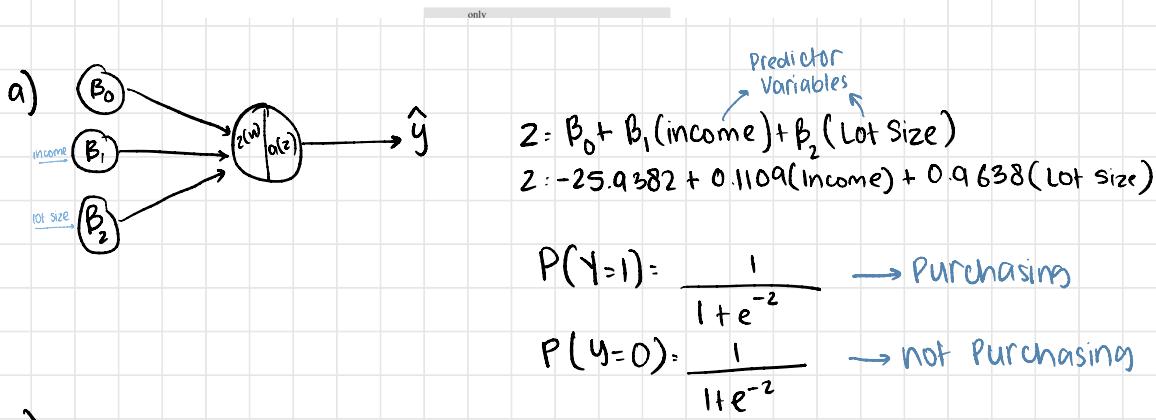
$$b_5^t = b_5 - \eta \frac{\delta L}{\delta b_5}$$

2. Consider a classification problem where we are building a logistic regression classifier. The task is to predict if a given city has a risk of a disease epidemic or not. The data is defined using two input

Practice Problem 6

- Construct a logistic regression model with two predictors for the riding mower example with $\beta_0 = -25.9382$, $\beta_1 = 0.1109$, $\beta_2 = 0.9638$, where β_1 and β_2 are for the "Income" and "Lot_Size" variables, respectively.
- Using the logistic regression model with probability cutoff = 0.75, classify the following 6 customers as "Owner" or "Nonowner": if $p \geq 0.75$ then the case as a "Owner". Present the results in a classification matrix.

Customer#	Income	Lot_Size	Ownership
1	60.0	18.4	Owner
2	64.8	21.6	Owner
3	84.0	17.6	Nonowner
4	59.4	16.0	Nonowner
5	108.0	17.6	Owner
6	75	19.6	Nonowner



Customer#	Income	Lot_Size	Ownership
1	60.0	18.4	Owner
2	64.8	21.6	Owner
3	84.0	17.6	Nonowner
4	59.4	16.0	Nonowner
5	108.0	17.6	Owner
6	75	19.6	Nonowner

Z	y	
-1.55028	0.175	NOT
2.0662	0.887	YES
0.34028	0.584	NOT
-3.92994	0.019	NOT
3.00188	0.9527	YES
1.26978	0.78	YES

DEEP LEARNING TERMS

INPUT LAYER

- ↳ first layer of NN
- ↳ has input features

Neurons/Nodes

↳ each node performs

- ↳ a weighted sum of inputs
- ↳ applies an activation function
- ↳ to produce an output

EPOCH

- ↳ 1 complete pass through the entire training set

Activation Function

- ↳ applied to the output of neuron to introduce non-linearity
- ↳ helps network learn complex rls.

OUTPUT LAYER

- ↳ final layer of NN
- ↳ produces Prediction
- ↳ no. of neurons of output layer depends on the type of task

↳ classification with multiple classes

SHALLOW NN

- ↳ a NN with only 1 or few hidden layers
- ↳ simple networks
- ↳ struggle with complex problems

PARAMETER

- ↳ internal values weights
biases
- ↳ these are learned by model during training

HIDDEN LAYER

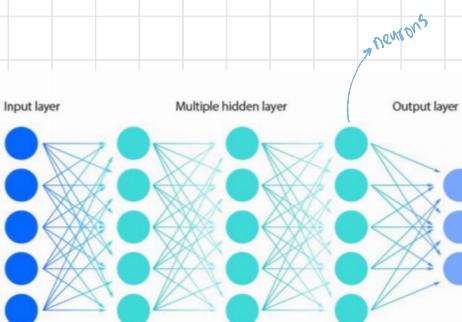
- ↳ layers b/w input and output
- ↳ they process data by applying weights, biases, activation functions to learn patterns from input data

DEEP NN

- ↳ a NN with multiple hidden layers
- ↳ good for complex problems
- ↳ are the foundation of Deep Learning

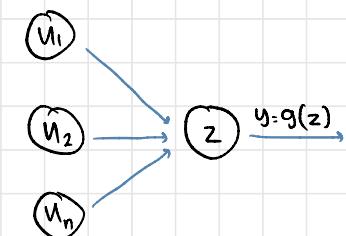
HYPERPARAMETER

- ↳ external settings learning rate
no of layers
no of neurons
- ↳ these are set before training
- ↳ controls the learning process



1. Is it true that the number of neurons in the output layer should match the number of classes i.e. where the number of classes are greater than 2.
Ans: False: Depends upon output coding. For 4 class problem, two output neurons are enough
2. What are hyper-parameters? Name any three.
Ans: Hyperparameters are configuration settings in machine learning models. Three examples are: Learning Rate, Number of Hidden Layers and Neurons, Batch Size
3. Is it possible to use ReLU in the last layer?
Ans: For unsupervised learning, yes. For classification NO
4. Write a derivative equation for log loss with L2 regularization.
Ans: $\frac{\partial \text{y}}{\partial (1-\text{y})} + (-\lambda w)$
5. What is the Difference Between Epoch, Batch, and Iteration in Deep Learning?
Ans: Epoch - Represents one iteration over the entire dataset (everything put into the training model).
Batch - Refers to when we cannot pass the entire dataset into the neural network at once, so we divide the dataset into several batches.
Iteration - If we have 10,000 images as data and a batch size of 200. Then an epoch should run 50 iterations (10,000 divided by 50)

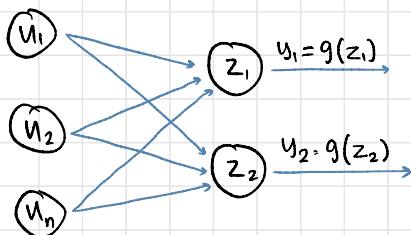
SIMPLE Perceptron



$$Z = W_0 + \sum_{j=1}^m u_j w_j$$

MULTI OUTPUT Perceptron

↳ As all inputs are densely connected to all outputs
these layers called dense layers



$$Z_i = W_{0,i} + \sum_{j=1}^m u_j w_{j,i}$$

Dense layer from scratch

```
class MyDenseLayer(tf.keras.layers.Layer):
    def __init__(self, input_dim, output_dim):
        super(MyDenseLayer, self).__init__()

        # Initialize weights and bias
        self.W = self.add_weight([input_dim, output_dim])
        self.b = self.add_weight([1, output_dim])

    def call(self, inputs):
        # Forward propagate the inputs
        z = tf.matmul(inputs, self.W) + self.b

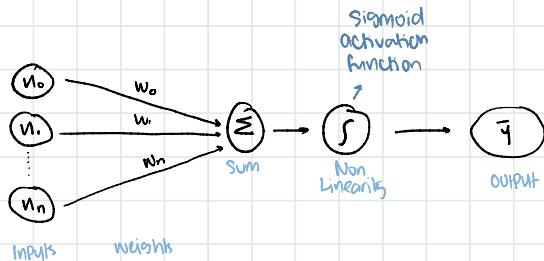
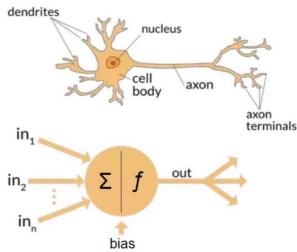
        # Feed through a non-linear activation
        output = tf.math.sigmoid(z)

        return output
```

NEURAL NETWORKS

STRUCTURE OF NEURON / Perceptron

mathematical model of a neuron



Artificial Neural Network

↳ an adaptive system

Parameters can be changed during training

↳ suit the problem

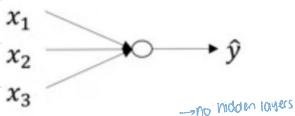
↳ a parallel system which is capable of resolving

Paradigms that linear computing can't resolve

L Layered NN

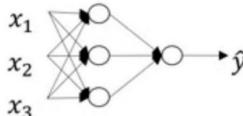
↳ no. of hidden layers + output layer

1 layered NN



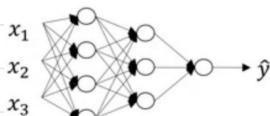
logistic regression

2 layered NN



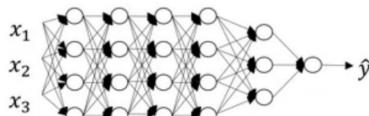
1 hidden layer

3 layered NN



2 hidden layers

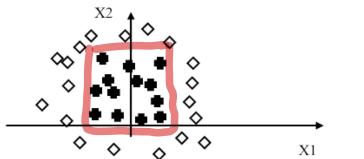
6 layered NN



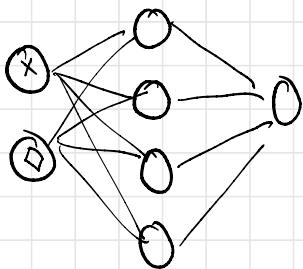
5 hidden layers

Question # 2**[CLO 1](2 marks)**

Suppose that we want to build a neural network that classifies two-dimensional data (i.e., $X = [x_1, x_2]$) into two classes: diamonds and crosses. We have a set of training data that is plotted as follows:



Draw a network that can solve this classification problem. Justify your choice of the number of nodes and the architecture.

**1. Input Layer:**

- 2 neurons (one for each feature x_1 and x_2).

2. Hidden Layer:

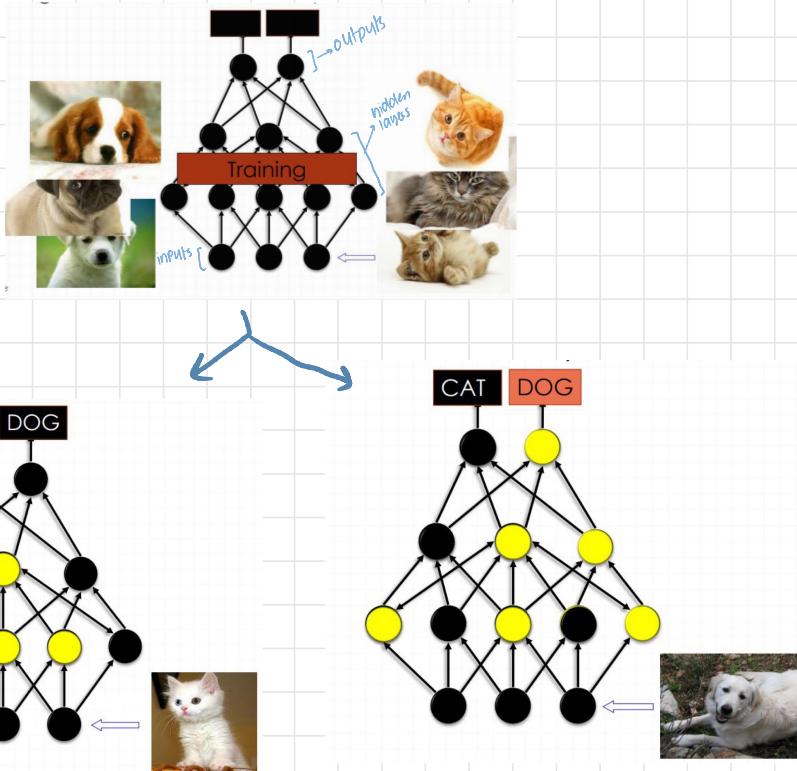
- 4 neurons (to capture the decision boundary).
- The hidden neurons form 4 linear boundaries, defining a **square-like decision region** around the crosses.
- Activation function: ReLU (Rectified Linear Unit) is commonly used.

3. Output Layer:

- 1 neuron (for binary classification).
- Activation function: Sigmoid ($\sigma(x) = \frac{1}{1+e^{-x}}$) to produce a probability output.

Deep Learning

- ↳ a set of machine learning algos based on multi-layer networks



Shallow NN

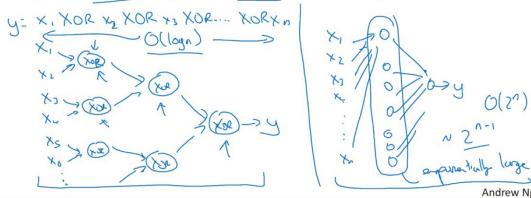
- ↳ consists of 1 or 2 layer NN

e.g. Logistic Regression

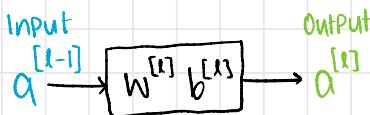
- ↳ It underfits the data

Circuit theory and deep learning

Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.



Forward Pass through layer l



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

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

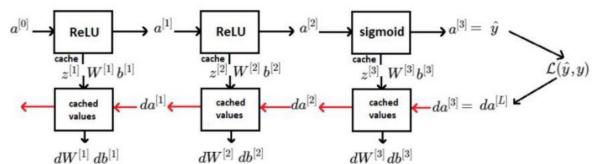
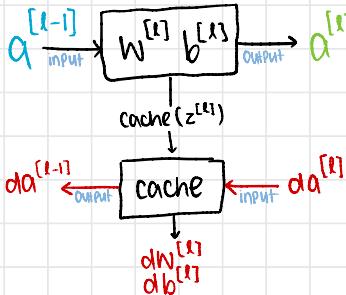
an activation function in the layer

Backward Pass



It is good to cache the value of $z^{[l]}$ for calculations in backward pass

Forward Pass and Backward Pass of layer L



$$\mathbf{W}^{[l]} := \mathbf{W}^{[l]} - \alpha \mathbf{dW}^{[l]}$$

$$b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

$$da^{[L]} = -\frac{y}{a} + \frac{1-y}{1-a}$$

Multilayer Backpropagation

- ↳ NN learning for classification/numerical prediction using backpropagation algo

↳ Inputs

- ↳ a data set "D"
 - ↳ training tuples
 - ↳ their associated target values
- ↳ the learning rate " η "
- ↳ a multilayer feed forward network

↳ Output

- ↳ a trained NN

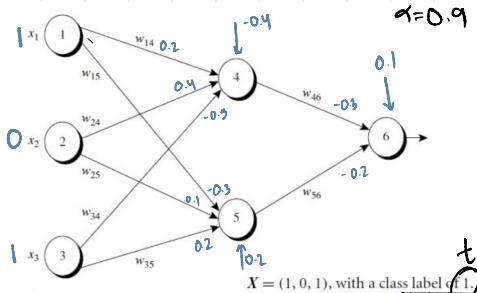
Weights Initialization

- ↳ they are initialized to small random no.
- ↳ each unit has a bias associated to it
 - ↳ biases also initialized to small random no.

Convergence condition

- ↳ Terminating condition
- ↳ Training stops when
 - ↳ all the Δw_{ij} in the prev epoch are below a specified threshold
 - ↳ The % of tuples misclassified in the prev epoch is below some threshold
 - ↳ a prespecified no. of epoch has expired

Worked Example → No loss function Given



Initial Input, Weight, and Bias Values

x_1	x_2	x_3	w_{14}	w_{15}	w_{24}	w_{25}	w_{34}	w_{35}	w_{46}	w_{56}	θ_4	θ_5	θ_6
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

↳ Feed Forward

$$z_4 = 1(0.2) - 0.5 - 0.4 = -0.7$$

$$z_5 = 1(-0.3) + 1(0.2) + 0.2 = 0.1$$

$$z_6 = 0.332(-0.3) + 0.525(-0.2) + 0.1 = -0.105$$

0

0.332

0.525

0.474

$$w^t = w + \alpha(e_j, o_i)$$

$$\frac{1}{1+e^{(-0.7)}}$$

↳ Calculation of error at each Node

↳ Output layer 6

$$\begin{aligned} \text{error}_j &= o_j (1 - o_j)(t - o_j) \\ &= 0.474(1 - 0.474)(1 - 0.474) \\ &= 0.1311 \end{aligned}$$

↳ Hidden layer 4

$$\begin{aligned} \text{error}_j &= o_j (1 - o_j)(e_j)(w_{jk}) \\ \text{error}_j &= 0.332(1 - 0.332)(0.1311)(-0.3) \\ &= -0.0087 \end{aligned}$$

↳ Hidden layer 5

$$\begin{aligned} \text{error}_j &= o_j (1 - o_j)(e_j)(w_{jk}) \\ \text{error}_j &= 0.525(1 - 0.525)(0.1311)(-0.2) \\ &= -0.0065 \end{aligned}$$

3

↳ Calculation for weight and biasness

$$w_{46}^t = -0.3 + 0.9(0.1311)(0.332) = -0.261$$

w_{46}

α

o_4

Calculations for Weight and Bias Updating

Weight or Bias	New Value
w_{46}	$-0.3 + (0.9)(0.1311)(0.332) = -0.261$
w_{56}	$-0.2 + (0.9)(0.1311)(0.525) = -0.138$
w_{14}	$0.2 + (0.9)(-0.0087)(1) = 0.192$
w_{15}	$-0.3 + (0.9)(-0.0065)(1) = -0.306$
w_{24}	$0.4 + (0.9)(-0.0087)(0) = 0.4$
w_{25}	$0.1 + (0.9)(-0.0065)(0) = 0.1$
w_{34}	$-0.5 + (0.9)(-0.0087)(1) = -0.508$
w_{35}	$0.2 + (0.9)(-0.0065)(1) = 0.194$
θ_6	$0.1 + (0.9)(0.1311) = 0.218$
θ_5	$0.2 + (0.9)(-0.0065) = 0.194$
θ_4	$-0.4 + (0.9)(-0.0087) = -0.408$

Classification of unknown datapoint (x)

- ↳ Input x into the trained NN
- ↳ Compute the net input and output of each unit → no backpropagation needed
- ↳ If Multiple output nodes → One per class
 - ↳ The node with highest output value determines the predicted class
- ↳ If One output node
 - ↳ Output $\geq 0.5 \rightarrow$ Positive class
 - ↳ Output $< 0.5 \rightarrow$ Negative class

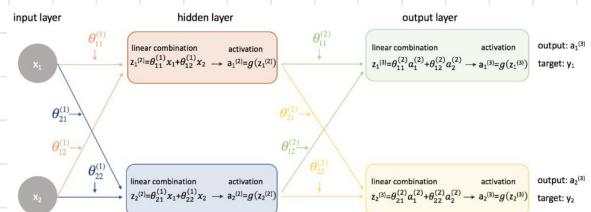
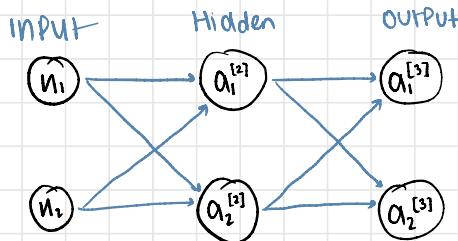
Critique of ANN

- **Parameter tuning:** Has a no of parameters that requires empirical determination of parameters like network structure.
 - **Long training times:** Suitable for applications where long training is feasible.
 - **Poor interpretability:** Difficult to understand the meaning behind learned weights and hidden units, making them less desirable for data mining initially.
- trial and error

Adv of ANN

- **Noise tolerance:** Can handle noisy data well.
- **Generalization:** Can classify patterns not seen during training.
- **Flexibility:** Useful when relationships between attributes and classes are unknown.
- **Continuous data:** Suitable for continuous-valued inputs and outputs (unlike most decision trees).
- **Real-world success:** Effective in applications like handwritten character recognition, medical diagnosis, and text pronunciation.
- **Parallel processing:** Neural networks are inherently parallel, allowing for faster computation.
- **Non-linear tasks:** Can perform tasks that linear classifiers cannot.
- **Function approximation:** Multilayer feed-forward networks can approximate any function given enough hidden units and training data.

Neural Network with Multi Output

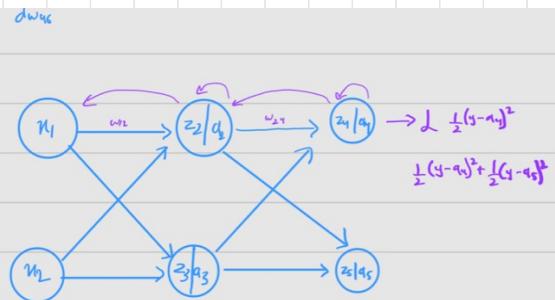


NN with Multi Output

$$\frac{dl}{dw_{16}} = \frac{dl}{da_1} \cdot \frac{day}{da_1} \cdot \frac{dza_1}{da_2} \cdot \frac{dza_2}{da_2} \cdot \frac{dza_2}{dw_{16}}$$

+

$$\frac{dl}{dw_{15}} \cdot \frac{das}{das} \cdot \frac{dza_5}{da_2} \cdot \frac{dza_2}{da_2} \cdot \frac{dza_2}{dw_{15}}$$

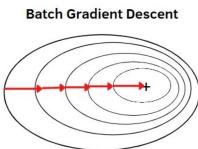


$$\frac{\delta L}{\delta w_{14}} = \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta a_4} \cdot \frac{\delta a_4}{\delta z_4} \cdot \frac{\delta z_4}{\delta w_{14}}$$

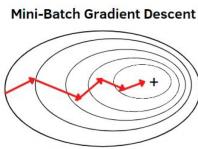
$$- (y_1 - a_6) a_6 (1 - a_6) a_4 \cdot \underbrace{\frac{a_4 (1 - a_4) (1)}{6z_4}}_{\delta z_4}$$

$$\frac{\delta L}{w_{35}} \cdot -(y_1 - a_6) a_6 (1 - a_6) w_{56} \cdot a_5 (1 - a_5) (n_3)$$

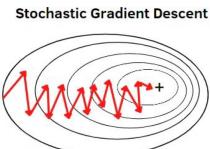
Gradient Descent Variations



Batch Gradient Descent



Mini-Batch Gradient Descent



Stochastic Gradient Descent

1. Batch Gradient Descent

- ↳ It won't update weights until every training sample has been assessed
- ↳ It is slow
- ↳ entire procedure is called cycle and training epoch

↳ aka Vanilla Gradient Descent

2. Stochastic Gradient Descent (SGD)

- ↳ doesn't wait for anyone
 - ↳ updates weights after every sample
- PROS:
- ↳ 100 samples → 100 updates
 - ↳ It is faster than batch
 - ↳ Regular updates gives an accurate idea of weights improvement
- CONS:
- ↳ can produce noisy gradients which can cause error rate to fluctuate rather than gradually go down

3. Mini-Batch Gradient Descent

- ↳ divides dataset into batches
 - ↳ updates after a batch
 - ↳ updates after every 10 samples
 - ↳ merges batch and SGD
- PROS:
- ↳ It strikes a balance b/w batch's effectiveness and SGD's durability

A standard feed-forward neural network is trained by Mini-Batch Gradient Descent. The network has 9 input neurons, 5 output neurons and a hidden layer, H. Hidden layer uses ReLU activation while output layer uses softmax activation function. The network is fully connected, but only the hidden layer uses biases. The network has a total of 300 trainable parameters, $p_i \in P$. During a particular session, it performs the forward and backward processes of on a small batch size of 40 samples out of total 200 observations.

1. How many nodes are in hidden layer H?
2. How many total internal calculations of gradients of the form $\partial L / \partial p_i$ are performed in one mini batch where p_i is any trainable parameter?
3. How many total parameter updates (of the p_i) are performed in a single mini batch?

Solution:

$$1. 9H + H + 5H = 300, so H = 20$$

$$2. 300 \times 40 = 12,000$$

$$3. 300, each p_i one time$$



SOFTMAX Activation Function for Multiclass

$$\frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$$z = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} \left\{ \begin{array}{l} 1.1 \rightarrow s_i = \frac{e^{z_i}}{\sum_l e^{z_l}} \rightarrow 0.224 \\ 2.2 \rightarrow \quad \quad \quad \quad \quad \rightarrow 0.672 \\ 0.2 \rightarrow \quad \quad \quad \quad \quad \rightarrow 0.091 \\ -1.7 \rightarrow \quad \quad \quad \quad \quad \rightarrow 0.013 \end{array} \right\} s = \begin{pmatrix} s_1 \\ s_2 \\ s_3 \\ s_4 \end{pmatrix}$$

SOFTMAX ACTIVATION FUNCTION FOR MULTICLASS - DERIVATIVE

$$J_{softmax} = \begin{pmatrix} \frac{\partial s_1}{\partial z_1} & \frac{\partial s_1}{\partial z_2} & \dots & \frac{\partial s_1}{\partial z_n} \\ \frac{\partial s_2}{\partial z_1} & \frac{\partial s_2}{\partial z_2} & \dots & \frac{\partial s_2}{\partial z_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial s_n}{\partial z_1} & \frac{\partial s_n}{\partial z_2} & \dots & \frac{\partial s_n}{\partial z_n} \end{pmatrix}$$

$$J_{softmax} = \begin{pmatrix} s_1 \cdot (1 - s_1) & -s_1 \cdot s_2 & -s_1 \cdot s_3 & -s_1 \cdot s_4 \\ -s_2 \cdot s_1 & s_2 \cdot (1 - s_2) & -s_2 \cdot s_3 & -s_2 \cdot s_4 \\ -s_3 \cdot s_1 & -s_3 \cdot s_2 & s_3 \cdot (1 - s_3) & -s_3 \cdot s_4 \\ -s_4 \cdot s_1 & -s_4 \cdot s_2 & -s_4 \cdot s_3 & s_4 \cdot (1 - s_4) \end{pmatrix}$$

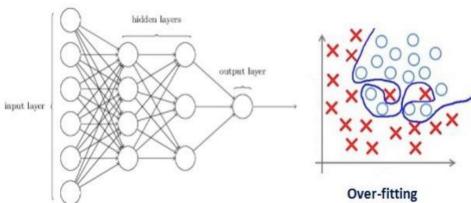
REGULARIZATION

OVERFITTING

- ↳ learns training data too well
- ↳ training accuracy high
- ↳ testing accuracy too low
- ↳ doesn't work well on unseen data

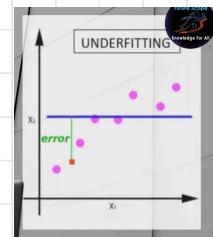


Learns more features than req.
Might learn diameter and material of chair
So wont classify other chairs as a chair



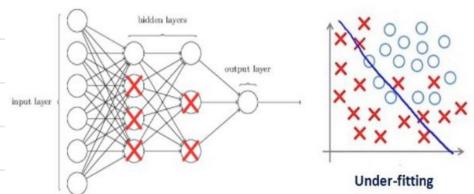
UNDERFITTING

- ↳ training accuracy too low
- ↳ high training error

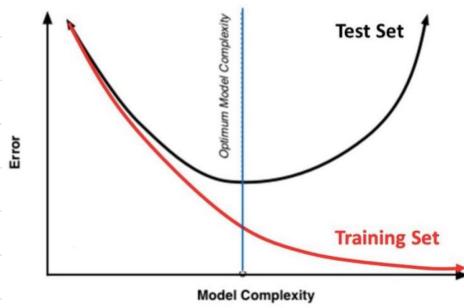


Might only learn 1 feature 4 legs
So will classify dog as a chair too

- ↳ increase the length of parameters/weights
- ↳ train longer



Training Vs. Test Set Error



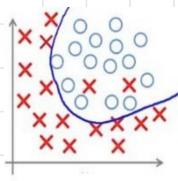
REGULARIZATION

- ↳ techniques that improve NN's generalization ability by modifying learning algo
- ↳ improves models performance
- ↳ Penalizes parameters that cause overfitting or underfitting

Assume that our regularization coefficient is so high that some of the weight matrices are nearly equal to zero

This will result in a much simpler linear network and slight underfitting of the training data.

In ML it penalizes co-efficients
In DLP it penalizes weights



Such a large value of the regularization coefficient is not that useful

We need to optimize the value of regularization coefficient in order to obtain a well-fitted model as shown in the image below

Regulation Techniques

- ↳ L1 L2 regularization
 - Penalizes parameters that cause overfitting
 - reduce model complexity
- ↳ Data Augmentation → Get more training data
- ↳ Early stopping
 - stops training before tolerance goes down
- ↳ Dropout
 - randomly disables a fraction of neurons during each forward

L1 Regularization

- ↳ adds absolute value of magnitude of the coefficient to loss function
 - ↳ if a feature doesn't serve any purpose it penalizes the weight to approach equal to 0

↳ Feature Selection

Modified loss function = Loss function + $\lambda \sum_{i=1}^n |w_i|$

$$E = \frac{1}{2} * \sum (t_k - o_k)^2 + \lambda * \sum |w_i|$$

squared error
L1 weight penalty
?

$\frac{\partial E}{\partial w_{jk}}$ gradient

$$\Delta w_{jk} = -1 * \eta * [x_j * (o_k - t_k) * o_k * (1 - o_k)] \pm \lambda$$

learning rate
signal

$$w_{jk} = w_{jk} + \Delta w_{jk}$$

L1 L2 Regularization

- ↳ combination of L1 and L2
- ↳ adds an extra hyperparameter that controls ratio of L1 and L2 regularization

L2 Regularization

- ↳ adds square magnitude of the coefficient to loss function
- ↳ shrinks parameter by spreading errors along all the terms
- ↳ reduces overfitting

Modified loss function = Loss function + $\lambda \sum_{i=1}^n w_i^2$

$$E = \frac{1}{2} * \sum (t_k - o_k)^2 + \frac{\lambda}{2} * \sum w_i^2$$

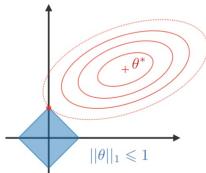
plain error
weight penalty
?

elegant math
simple math

$$\frac{\partial E}{\partial w_{jk}}$$
 gradient

$$\Delta w_{jk} = \eta * [x_j * (o_k - t_k) * o_k * (1 - o_k)] + [\lambda * w_{jk}]$$

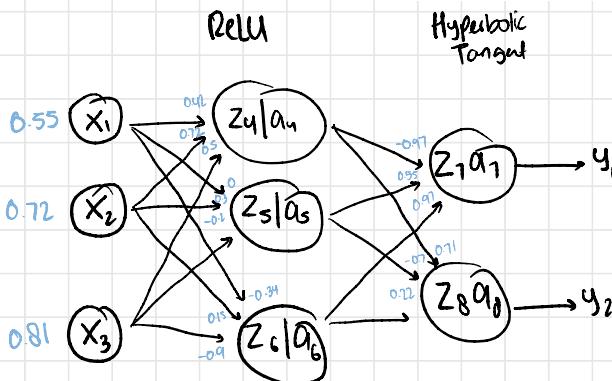
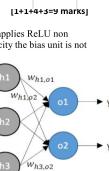
learning rate
signal

LASSO	Ridge	Elastic Net
<ul style="list-style-type: none"> - Shrinks coefficients to 0 - Good for variable selection 	Makes coefficients smaller	Tradeoff between variable selection and small coefficients
$... + \lambda \theta _1$ $\lambda \in \mathbb{R}$	$... + \lambda \theta _2^2$ $\lambda \in \mathbb{R}$	$... + \lambda [(1 - \alpha) \theta _1 + \alpha \theta _2^2]$ $\lambda \in \mathbb{R}, \alpha \in [0, 1]$

Q3) Learning Rate = 0.5

Regularization coefficient = 0.03

Q3:
 Consider the following feed forward neural network.
 Consider the given two-layer architecture, where each node in the layer L1 applies ReLU non-linearity and L2 uses the Hyperbolic Tangent activation function, for simplicity the bias unit is not considered for the given network.
 $X = \{[0.55], [0.72], [0.81]\}$
 weights of L1
 $W1 = [[0.42, 0.72, 0.51], [0.3, -0.2], [-0.34, 0.15, -0.9]]$
 weights of L2
 $W2 = [[-0.97, 0.55, 0.97], [0.71, -0.7, 0.22]]$
 Learning Rate = 0.5.
 Regularization Coefficient = 0.03



Feed Forward

$$Z_4 = (0.55)(0.42) + (0.72)(0.12) + (0.81)(0.5) = 1.1544 \quad a_4 = 1.1544$$

$$Z_5 = (0.55)(0) + (0.72)(0.3) + (0.81)(-0.2) = 0.094$$

$$Z_6 = (0.55)(-0.34) + (0.72)(0.15) + (0.81)(-0.9) = -0.808$$

$$Z_7 =$$

- Compute the output of each node in the network for forward propagation.

$$Z(h1)=1.1544 \Rightarrow a(h1)=1.1544$$

$$Z(h2)=0.116 \Rightarrow a(h2)=0.116$$

$$Z(h3)=-0.4848 \Rightarrow a(h3)=0$$

$$Z(o1)=1.0559 \Rightarrow a(o1)=0.784$$

$$Z(o2)=0.7384 \Rightarrow a(o2)=0.63819$$

- Compute the error if $y = [-0.98, 0.75]$

$$E = \frac{1}{2} ((-0.784+0.98)^2 + (0.63819-0.75)^2) = 0.02545$$

- Perform backpropagation, and compute the derivatives i.e. $dWh1O1, dWh2O2, dWx0h1$. Calculate the updated weights.

$$D h1o1 = -(-0.98+0.784)(1-0.784^2) 1.1544 = 0.0871$$

$$Wh1o1 = -0.97 - 0.5^*0.0871 = -1.01355$$

$$D h2o2 = -(0.75-0.63819)(1-0.63819^2) 0.116 = -0.00768$$

$$Wh2o2 = -0.7 - 0.5^*0.00768 = -0.69616$$

$$D X0h1 = -(-0.98+0.784)(1-0.784^2) (-0.97)(1)(0.55) + (-0.75-0.63819)(1-0.63819^2) (0.71)(1)(0.55)$$

$$D X0h1 = -0.0662$$

$$W X0h1 = 0.42 - 0.5^* -0.0662 = 0.4531$$

- Perform backpropagation with **L1 regularization** and compute the derivatives for $dWh2O1, dWh3O2$. Calculate the updated weights.

$$D h2o1 = -(-0.98+0.784)(1-0.784^2) 0.116 + 0.03 = 0.03876$$

$$Wh2o1 = 0, -0.5^*0.03876 = -0.01938$$

$$D h3o2 = 0 + 0.03 = 0.03$$

$$Wh3o2 = 0.22 - 0.5^*0.03 = 0.205$$

$$E = \frac{1}{2} (y - \hat{y})^2 \cdot \lambda$$

$$w_{12}^+ = w_{12} - \alpha \left(\frac{\delta L}{\delta w_{12}} + \lambda \right)$$

DATA AUGMENTATION

↳ synthesizing data

↳ increase training data by

↳ Resiz

↳ Horizontal/Vertical flip

↳ Rotate

↳ Add noise

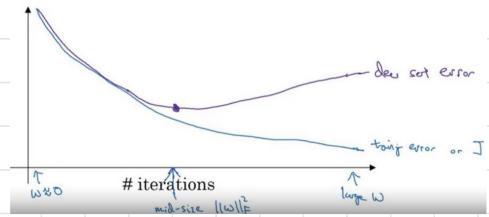
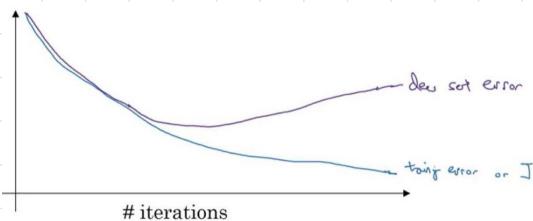
↳ Deform

↳ Modify colors



EARLY STOPPING

↳ stops training when performance on the validation set starts to degrade



DROPOUT

↳ forgets 20% - 50% of the learning

↳ randomly disables a fraction of neurons during each forward pass

↳ on backward pass only update those weights

↳ it reduces co-dependence among features

The network with dropout during a single forward pass

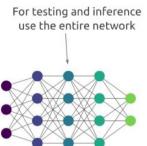


Nodes set to zero during forward passes

Dropout is the equivalent of training several independent, smaller networks on the same task. The final model is like an ensemble of smaller networks, reducing variance and providing more robust predictions.

Dropout
Network regularization

For each forward pass during training, set the output of each node to zero with probability P .

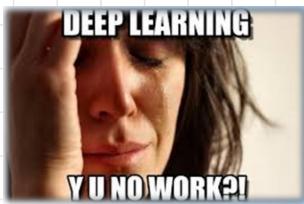


For testing and inference use the entire network

DROPOUT

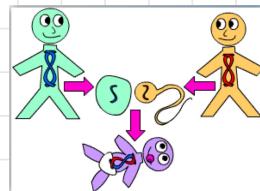
Motivation.

- ↳ NN can sometimes overfit
- ↳ Solution is to use Bayesian Framework
 - ↳ but not feasible if network too big
- ↳ we need something faster



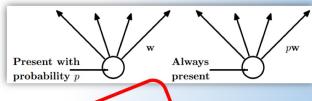
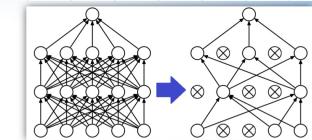
INSPIRATION

- ↳ Genes work well another random set of genes
- ↳ Similarly, dropout suggests that each unit should work with a random set of other samples



DROPOUT

- ↳ At each training iteration
 - ↳ each unit is retained with a probability P
- ↳ At test
 - ↳ the network is used as a whole
 - ↳ the weights are scaled down by a factor of P
- ↳ In practice
 - ↳ dropout trains 2^n networks (n - no. of units)

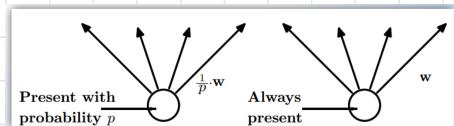


$p = 0.5$
MAKES THE
CHARM!



DROPOUT - an equivalent method

- ↳ At training
 - ↳ weights are scaled up by a factor of $\frac{1}{P}$
- ↳ At test
 - ↳ no scaling is applied
- ↳ This method is used in TensorFlow
 - ↳ `tf.nn.dropout(u, keep_prob=p)`



Why apply dropout on units?

- ↳ with all their in and outgoing arcs
and not just on the arc's themselves

Dropout is applied to entire neurons (units) rather than individual connections (arcs) to prevent co-adaptation, encourage independent feature learning, and improve generalization.

Dropping full neurons forces the network to learn redundant and distributed representations, making it more robust.

It also simplifies implementation and has been empirically shown to work better.

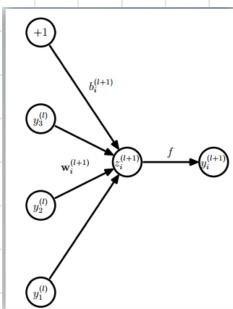


Dropout - Model Description

Feed Forward w/o Dropout

$$z_i^{(l+1)} = w_i^{(l+1)} y_i^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$



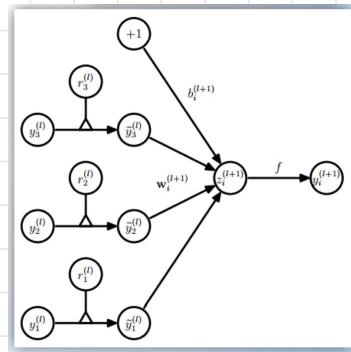
Feed Forward with Dropout

$$r_i^{(l)} \sim \text{Bernoulli}(p)$$

$$\tilde{y}_i^{(l)} = r_i^{(l)} \cdot y_i^{(l)}$$

$$z_i^{(l+1)} = w_i^{(l+1)} \tilde{y}_i^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$



Backpropagation with Dropout

- ↳ is only performed on trained network

for each training iteration

- ↳ At training

↳ a case which does not use a parameter

contributes a gradient of '0' for that parameter

- ↳ A test

↳ the weights are scaled as:

$$W_{\text{test}}^{(l)} = p \cdot W^{(l)}$$

Experimental Results

1. Street View House Numbers (SVHN)

↳ Dropout is also applied to convolution layers

↳ all hidden units are ReLUs



Method	Error %
Binary Features (WDCH) (Netzer et al., 2011)	36.7
HOG (Netzer et al., 2011)	15.0
Stacked Sparse Autoencoders (Netzer et al., 2011)	10.3
KMeans (Netzer et al., 2011)	9.4
Multi-stage Conv Net with average pooling (Sermanet et al., 2012)	9.06
Multi-stage Conv Net + L2 pooling (Sermanet et al., 2012)	5.36
Multi-stage Conv Net + L4 pooling + padding (Sermanet et al., 2012)	4.90
Conv Net + max-pooling	3.95
Conv Net + max pooling + dropout in fully connected layers	3.02
Conv Net + stochastic pooling (Zeiler and Fergus, 2013)	2.80
Conv Net + max pooling + dropout in all layers	2.55
Conv Net + maxout (Goodfellow et al., 2013)	2.47
Human Performance	2.0

2. CIFAR-10 & CIFAR-100

↳ images drawn from 10 and 100 categories respectively

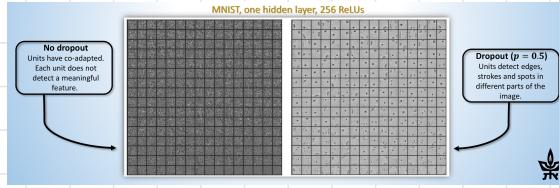


Method	CIFAR-10	CIFAR-100
Conv Net + max pooling (hand tuned)	15.60	43.48
Conv Net + stochastic pooling (Zeiler and Fergus, 2013)	15.13	42.51
Conv Net + max pooling (Snoek et al., 2012)	14.98	-
Conv Net + max pooling + dropout fully connected layers	14.32	41.26
Conv Net + max pooling + dropout in all layers	12.61	37.20
Conv Net + maxout (Goodfellow et al., 2013)	11.68	38.57

Effect Of Dropout On

1. Learned Features

- ↳ w/o dropout
 - ↳ units tend to compensate for mistakes of other units
 - ↳ leads to overfitting
 - ↳ as these co-adaption do not generalize to unseen data
- ↳ with dropout
 - ↳ reduced overfitting
 - ↳ as co-adaption are prevented by making presence of other hidden layer unreliable



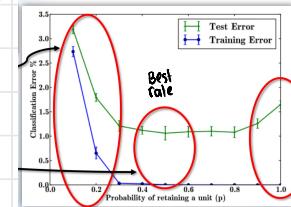
2. Dropout rate P

- ↳ An architecture of is used on MINST dataset
- ↳ The dropout rate P is changed from

small numbers to 1.0

most units are dropped ↴

↳ no dropout



↳ High rate of dropout ($P < 0.3$)

↳ Training error high → Underfitting

↳ very few units are turned on during training

↳ Best rate of dropout ($P = 0.5$)

↳ Training error low

↳ Test error is low

↳ Best

↳ NO dropout ($P = 1.0$)

↳ Training error low

↳ Test error is high → Overfitting

3. Dataset Size

- ↳ An architecture of is used on MINST dataset

↳ Extremely small data set

↳ dropout does not improve error rate

↳ makes it worse sometimes

↳ Average to large dataset

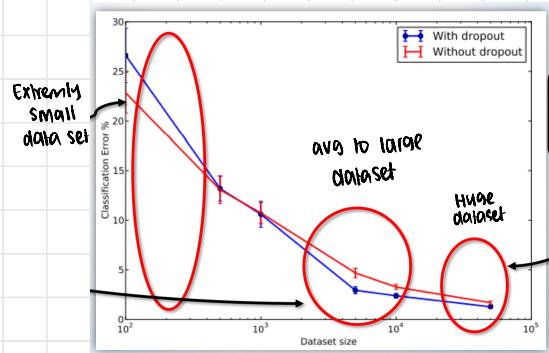
↳ dropout improves error rate

↳ Huge dataset

↳ dropout barely improves error rate

↳ dataset is big enough so overfitting

isn't an issue



Weight Decay

↳ a regularization technique

↳ limiting the growth of weight in the network

↳ a term is added to the original loss function

penalizing large weights

$$L_{\text{new}}(w) = L_{\text{old}}(w) + \frac{1}{2} \lambda \|w\|_2^2$$

↳ An architecture of is used on MNIST dataset with diff regularizations

Method	Test Classification error %
L2	1.62
L2 + L1 applied towards the end of training	1.60
L2 + KL-sparsity	1.55
Max-norm	1.35
Dropout + L2	1.25
Dropout + Max-norm	1.05

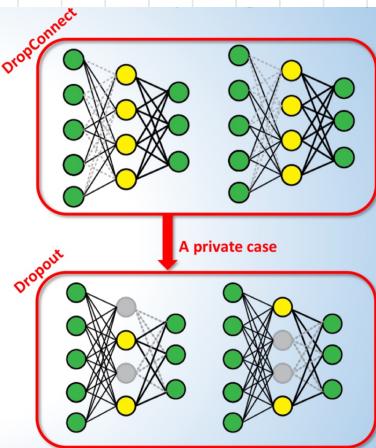
DropConnect

↳ a regularization technique

↳ It generalizes dropout by suggesting dropping out weights

↳ In Dropout P=0.5

↳ usually gives the best result



→ pushes weights towards smaller values to prevent overfitting

Dropout vs Weight Decay

↳ Dropout Adv

↳ it is scale free

↳ it doesn't penalize the use of large weights when needed

↳ is invariant to parameter scaling

↳ it is unaffected if

↳ weights in a certain layer are upscaled by a constant C

↳ and weights in another layer are downscalled by a constant C

↳ implying dropout does not have isolated minima

DropConnect

Averaging Technique:

- DropConnect asserts:

$$\begin{aligned} z^{l+1} &= (M * W)y^l \\ E_M[z^{l+1}] &= pW y^l \\ V_M[z^{l+1}] &= p(1-p)(W * W)(y^l * y^l) \end{aligned}$$

Input to the activation function

A weighted sum of Bernoulli variables. Can be approximated by a Gaussian

Statistics of the Gaussian

- At test:

- Draw samples of z^{l+1} according to the Gaussian distribution.
- Feed the samples into the activation function ($f(z^{l+1})$).
- Average.

Cons

↳ training is slower

↳ implementation is more complicated than Dropouts'

↳ In their Paper, it has been proven to work mostly when using more than 1 network

COMPOSITION OF

1. NO-drop
2. Dropout
3. DropConnect

MNIST dataset

↳ 2 layers → 800 neurons each

MNIST

neuron	model	error(%) 5 network	voting error(%)
relu	No-Drop	1.62 ± 0.037	1.40
	Dropout	1.28 ± 0.040	1.20
	DropConnect	1.20 ± 0.034	1.12
sigmoid	No-Drop	1.78 ± 0.037	1.74
	Dropout	1.38 ± 0.039	1.36
	DropConnect	1.55 ± 0.046	1.48
tanh	No-Drop	1.65 ± 0.026	1.49
	Dropout	1.58 ± 0.053	1.55
	DropConnect	1.36 ± 0.054	1.35

CIFAR-10

↳ 4 layers

↳ Dropout & DropConnect

Applied only on final layer

CIFAR-10

model	error(%)
No-Drop	23.5
Dropout	19.7
DropConnect	18.7

SVHN

↳ same architecture as CIFAR-10

↳ due to large training set size

all methods nearly achieve
nearly same performance

SVHN

model	error(%) 5 network	voting error(%)
No-Drop	2.26 ± 0.072	1.94
Dropout	2.25 ± 0.034	1.96
DropConnect	2.23 ± 0.039	1.94

↳ In DropConnect's Paper on MNIST dataset

↳ DropConnect achieves the lowest error rate

recorded so far in MNIST

↳ results are equivalent for

↳ Drop Connect and no dropout



Really?

crop	rotation scaling	model	error(%) 5 network	voting error(%)
no	no	No-Drop	0.77 ± 0.051	0.67
		Dropout	0.59 ± 0.039	0.52
		DropConnect	0.63 ± 0.035	0.57
yes	no	No-Drop	0.50 ± 0.098	0.38
		Dropout	0.39 ± 0.039	0.35
		DropConnect	0.39 ± 0.047	0.32
yes	yes	No-Drop	0.30 ± 0.035	0.21
		Dropout	0.28 ± 0.016	0.27
		DropConnect	0.28 ± 0.032	0.21

DROPOUT

PROS

- ↳ good and fast regularization method
- ↳ if dataset is avg-large it excels
- ↳ achieves better results than weight decay

CONS

- ↳ slow to train → 2-3 times slower w/o dropout

- ↳ if dataset is big enough it doesn't help much

DROP Connect

- ↳ is a generalization of dropout

CONS

- ↳ training is slower than dropout
- ↳ its superiority over dropout is unclear
- ↳ has complications in implementation

SUMMARY

- ↳ use a small dropout value → 20% - 50% of neurons

- ↳ a probability too

- ↳ low → has minimal effect

- ↳ high → under-learning by network

- ↳ use a larger network → dropout then gives better performance

- ↳ use dropout on incoming and visible units

- ↳ use of dropout at each layer → has good results

- ↳ use large learning rate with decays → increase by factor of 10-100 and large learning momentum → 0.9/0.99

- ↳ constrain the size of network weights

- ↳ large learning rate → large network weights

- ↳ imposing a constraint on the size of network weights → improves results

ishma hafeez
notes

repst
rept