

# Machine Learning for Design Engineers



---

---

---

---



## 5. Overview of lecture content

Lecture title	Instructor	Content
1. The magic of machine learning	Robert Shorten	Overview of course Introduction to team and assessment Some background material Reading list Background material for laboratories
2. The machine learning portfolio	Robert Shorten	What is Machine learning? Types of machine learning? Basic mathematical background Sprint on Markov chains
3. Key ideas in one lecture	Robert Shorten	The basics explained in an elementary manner.
4. Practicalities of real-life machine learning.	Robert Shorten	Some of the practical issues in dealing with data.
5. ANNs ✓	Robert Shorten	An introduction to artificial neural networks.
6. ANNs ✓	Robert Shorten	Learning algorithms for artificial neural networks.
7. CNNs and RNNs	Pietro Ferraro	An overview of convolutional neural networks
8. SVM ✓	Robert Shorten	An overview of support vector machines and Kernel methods
9. Supervised learning	Robert Shorten	An introduction to K-means clustering
10. Reinforcement learning	P. Ferraro	What are Markov Decision processes?
11. Reinforcement learning	<u>P. Ferraro</u>	What is Q-learning?
13. Spectral graph theory.	Robert Shorten	An introduction to machine learning over graphs.
14. Spectral graph theory	Robert Shorten	An introduction to machine learning over graphs.
15. Spectral graph theory	Robert Shorten	An introduction to machine learning over graphs.
15. Advanced topic	Ms. Quan Zhao	What is fairness in AI?
Algorithmic fairness		
16. Advanced topic	Martin Mevissen – IBM Research	Quantum machine learning – an introduction
Quantum AI		

Optimal transport

# Lecture ①

ML is good at:

- Regression Problems
- Pattern Recognition Problems
- Sequential Decision Making Problems
- Finding Features of Graphs

## Context of Design Engineering

- Designing Products (i.e. how virtual objects react in real-time)
- Part of a designed system
- ML Algorithms

## Types of Learning

- Supervised Learning (w/ labels)
- Non-Supervised Learning (w/o labels)
- Exploration-Based Learning (finding Optimal Policies)
- Machine Learning on Graphs

## Books:

1. The hundred-page machine learning book, A. Burkov  
*Surface Level*
2. Hands-on machine learning with SciKit-learn, Keras and TensorFlow, A. Geron
3. Decision making under uncertainty, M. Kovaldorfer
4. Google's page rank and beyond, A. Langville, C. Meyer
5. Mathematics for Machine Learning, M. Davenport et. al.

# Lecture ②

## 2-types of Machine Learning

- Supervised (Regression & Classification)
- Unsupervised

### Supervised learning

- provided labels

### Unsupervised Learning

- learns on its own
- used in clustering applications

### Semi-supervised Learning

- combination of labeled and unlabeled

### Reinforcement Learning

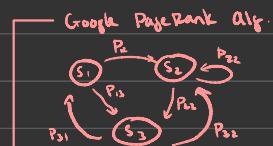
- "reward" system to optimise based on outcome

Built on 4 areas of mathematics: Optimisation, Probability, Linear Algebra, Markov Chains

#### - Lagrange Multipliers

- Probability Theory (random variable, mass & density function)
- Markov Chain (essential for graphs and reinforcement learning)

'dynamic system that evolves probabilistically with NO memories



# Lecture 3

- linear regression
- gradient descent
- building a simple predictor
- building a classifier

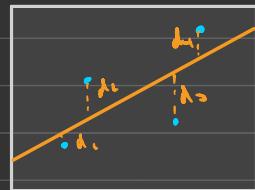
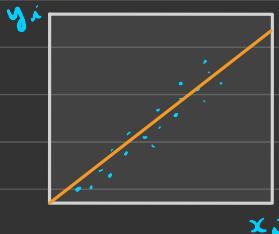
Predict  $y(x)$  as  $x(i)$  increases:

- fit a line  $g(x) = mx(x) + c$

- least squares problem (optimization)

$$J(m, c) = d_1^2 + d_2^2 + \dots + d_n^2$$

$$= \sum_{i=1}^n (mx_i + c - y_i)^2$$



$$\theta = (U^T U)^{-1} U^T Y$$

derivation:  $Y = [y^{(1)}, y^{(2)}, \dots, y^{(n)}]^T$ ,  $\theta = [c, m, 1; c, m, 1; \dots; c, m, 1]^T$

if linear relationship,  $Y = U\theta \Rightarrow U^T Y = U^T U\theta$

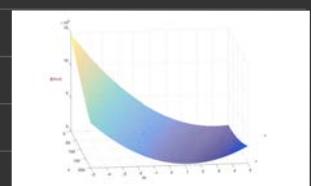
if  $U^T U$  is invertible,  $(U^T U)^{-1} U^T Y = \theta$

assumes: 1) linear relationship

2) squared distance ( $d_i^2$ )

$\Rightarrow U^T U$  is invertible (multiple unique solutions; dataset is rich enough)

can be used for nonlinear problems (i.e.  $y = mx^2, e^x$ )



$J(m, c)$  as a function of  $m, c$

## Gradient Descent

Method:

Gives initial value  $m_0, c_0 \rightarrow$  calculate

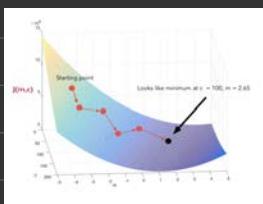
gradient  $\frac{\partial J}{\partial m}, \frac{\partial J}{\partial c}$  on the plot:

$\frac{\partial J}{\partial m}, \frac{\partial J}{\partial c} \Rightarrow$  update estimates of parameter:

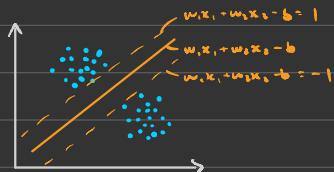
$$m \leftarrow m_0 - \alpha_m \frac{\partial J}{\partial m}; c \leftarrow c_0 - \alpha_c \frac{\partial J}{\partial c} \text{ where}$$

$\alpha_m$  and  $\alpha_c$  are small constants

repeat  $\Rightarrow$



## Classification



Minimize  $w_1^2 + w_2^2$  subject to the constraints:

$$w_1x_1 + w_2x_2 - b \geq 1 \text{ if } y(i)=1$$

$$w_1x_1 + w_2x_2 - b \leq -1 \text{ if } y(i)=-1$$

## Lecture ④

### Anatomy of Machine Learning

- every ML algorithm has the following steps:

- 1) cost function
- 2) optimisation criterion based on a loss function
- 3) algorithm to fit the optimisation

Gradient Descent ~ method of choice to implement optimisation

cost function:  $J(w, c) = \frac{1}{n} \left( \sum_{i=1}^n (y(i)x(i) + c - y(i))^2 \right)$  given n data points

$$\frac{\partial J}{\partial w} = \frac{1}{n} \sum_{i=1}^n 2x(i)(y(i) - (wx(i) + c))$$

$$\frac{\partial J}{\partial c} = \frac{1}{n} \sum_{i=1}^n 2(y(i) - (wx(i) + c))$$

Practical Comments: Consider 3 things

- getting data ready for learning ~ label the dataset ~ normalization procedure
  - ~ deal w/ missing data (data imputation) ~ learning (training, test, validation test)
- selecting features and the model ~ too few parameters  $\rightarrow$  underfitting  $\rightarrow$  high bias
  - ~ too many parameters  $\rightarrow$  overfitting  $\rightarrow$  high variance ~ bias variance dilemma
  - ~ overfitting solution (regularisation) aka adding a term that constrains the parameter. (i.e. Ridge, LASSO)
- choosing the appropriate learning algorithm ~ many factors (i.e. hardware, continuous, training speed, type of data & problem, nonlinearities, memory constraints)

# Lecture ⑤

- ANN - artificial neural networks

~ useful to solve regression and classification tasks

~ best known is Multi-layer Perceptron

- Linear Separation

$$f(w_0x_0 + w_1x_1 + w_2x_2 + w_3) = y \quad \begin{cases} y=1, & w_0x_0 + w_1x_1 + w_2x_2 + w_3 > 0 \\ y=-1, & \text{else} \end{cases} \quad \text{f labels each class 1 or -1}$$

(Issue) Single layer could not separate classes for following cases

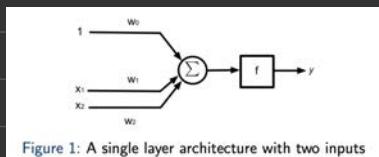


Figure 1: A single layer architecture with two inputs

$$f(w_0x_0 + w_1x_1 + w_2x_2 + w_3) = y_1,$$

$$f(w_0x_0 + w_1x_1 + w_2x_2 + w_3) = y_2$$

$$f(w_0y_1 + w_1y_2 + w_3) = y_3$$

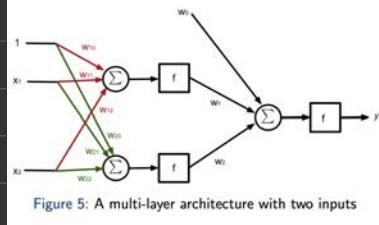
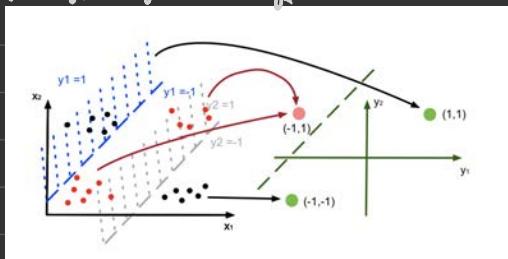


Figure 5: A multi-layer architecture with two inputs

Problems of the difficulty to apply gradient descent to learn the parameters  $\Rightarrow$  MLP

Multi-layer Perceptron

MLP - multi-layer perceptron

- When non-differentiable  $\Rightarrow f(z) = \frac{1}{1+e^{-z}}$  logarithmic or  
 $f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$  tanh function

$\Rightarrow$  differential approximation to rectified linear function

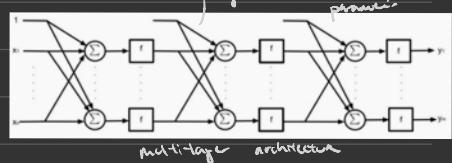
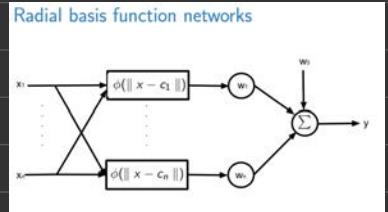
Network is a non-linear system - difficult to understand and certify.  
outputs are also nonlinear

Radial Basis function Networks (RBF) - single layer network

$$y = \sum_{i=1}^n w_i R_i(x) + w_0 \quad R_i(x) = \phi(\|x - c_i\|) \quad c_i : \text{RBF centres}$$

- estimating weights ( $w$ ) is a linear regression problem
- many optimisation methods take the form

## Lecture 6



- how to find parameter of network  $\rightarrow$  gradient descent
- 3 ways to calculate gradient:
  1. Numerical Differentiation
  2. Symbolic Manipulation
  3. Algorithmic Differentiation

Numerical  $\frac{\partial J}{\partial w_i} = \frac{J(w_i + h) - J(w_i)}{h}$  for small  $h$

easy to program not exact, how to choose  $h$ ?

Symbolic exact method by hand tedious, complicated formulae

Automatic Back propagation and partial derivatives

Backpropagation - minimize error iteratively. Chain rule and partial der.

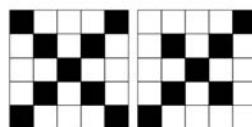
stochastic vs batch

All derivations are in Lecture 6 \*

# Lecture ⑦

## CNN - convolutional neural networks

A human would have no issues saying that these two images look very similar.



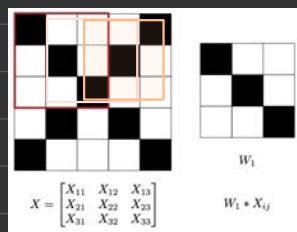
Why not LNN?

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = [1 \ 0 \ 0 \ 1]$$

\* loses spatial feature such as diagonals. applies to other forms.

## Convolutional Kernels (filters)

$$X_{ij} * W_1 = \sum_i \sum_j w_{ij} x_{ij}$$

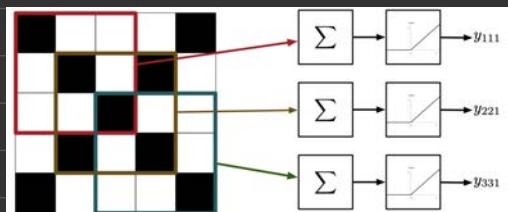


1	-1	-1	-1	1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
1	-1	-1	-1	1

$$* \begin{bmatrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{bmatrix} = \begin{bmatrix} 9 & -3 & 1 \\ -3 & 3 & -3 \\ 1 & -3 & 9 \end{bmatrix}$$

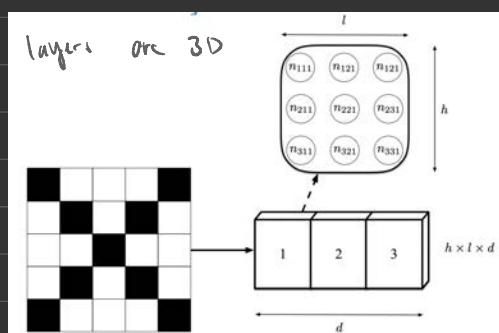
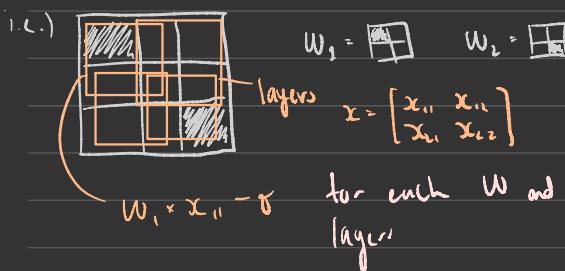
higher <sup>hit</sup> matches

## Convolutional Layers

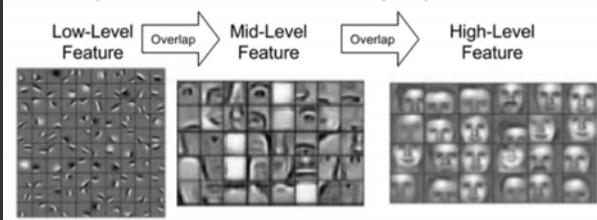


$$h(s) \begin{cases} s, & \text{if } s \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

$$y_i = \sigma \left( \sum_j w_{ij} x_j + b_i \right)$$



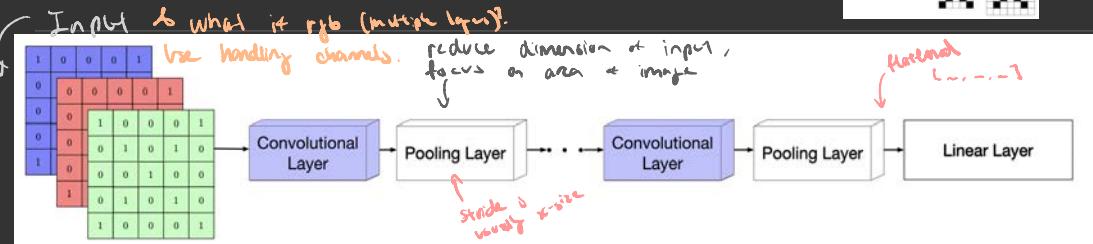
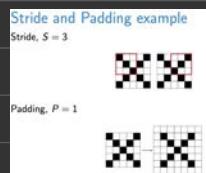
### Feature Map in Convolutional Neural Networks (CNN)



- layer is composed of  $h \times w$  neurons
- output is fed to next layer
- "feature map" set of neurons that share the same feature

- **Kernel Dimensions ( $k_{xi} \times k_{yi}$ ):** Dimensions of Matrix  $W_i$
- **Stride (s):** distance btw current and next operation
- **Padding (p):** filling zeros at 0 on edges to avoid cropping effects.

CNN architecture  $\xrightarrow[\text{number of channels}]{} [ ] \Rightarrow \text{CNN} \Rightarrow \dots$



### Pooling Layer

There are different types  
(MaxPool, AvgPool, stride almost always =  $x \times z$ )  
→ Trial and Error

9	-3	1
-3	3	-7
1	-3	9

ReLU  $\Rightarrow$

9	0	1
0	3	0
1	0	9

Pooling max  $\Rightarrow$

9	3
3	9

- eliminating parts that are not too important
- reduce dimension



add two lines of 0s since  
dimension is odd

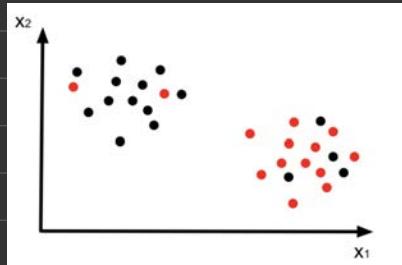
9	0	1	0
0	3	0	0
1	0	9	0
0	0	0	0

$\Rightarrow$

9	1
1	9

## Support Vector Machine

- limitation ~ constraint in every data point (too many constraints)
- ~ noisy data
- ~ non-linearities (not separable by a line)



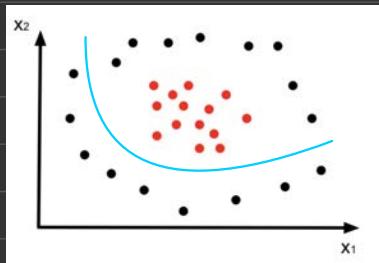
Noisy Data

how to deal w/ it

hinge function

$\max(0, 1 - y(i)(w_1x_1(i) + w_2x_2(i) - b))$

↳ non-differentiable  $\rightarrow$  subgradient



$$b^2 \leq w_1^2 x_1^2 + w_2^2 x_2^2 + w_3^2 x_3^2$$

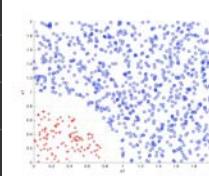


Figure 1: No linear separation

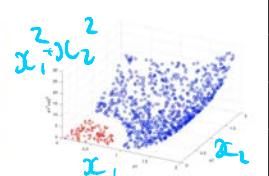


Figure 2: Linear separation

Non-Linearity

- adding another dimension, existing in 3 dimensions

Minimize over all data points

$$(w_1^2 + w_2^2) + \frac{1}{\mu} \sum_{i=1}^n \max(0, 1 - y(i)(w_1x_1(i) + w_2x_2(i) - b))$$

# Lecture 9

## Unsupervised learning and K-means Clustering

~ dealing w/ unlabelled data

~ problems regularly:

clustering - group similar patterns into clusters

anomaly detection - finding outliers

density estimation - estimating PDF for random points given dataset.

What is k-means algorithm?

### Inputs

$k$ , number of clusters

$X$ , training data  $\{x(1), x(2), x(3), \dots, x(M)\}$

- 1) random initialisation, choose  $k$  points to begin for centre ( $k < M$ )
- 2) continue updating centre for each cluster assignment.
- 3) repeat until assignment doesn't change.

$c(i)$  = index of cluster when  $x(i)$  is assigned

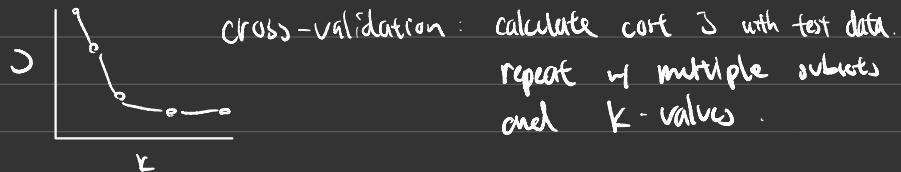
$\mu(j)$  = centre of cluster  $j$

$\mu(c(i))$  = cluster centre for  $x(i)$

$$\text{minimize } J(c(1), \dots, c(M), \mu(1), \dots, \mu(k)) = \frac{1}{M} \sum_{i=1}^M \|x(i) - \mu(c(i))\|^2$$

can converge to a global optimum instead.

elbow method ← when choosing # of clusters



# Markov Chain



# Lecture 10

- dynamic system that evolves probabilistically
- NO MEMORY



- If Probabilities do not change (i.e wrt time)  $\rightarrow$  homogeneous
- If it does  $\rightarrow$  non-homogeneous

## (1) Markov Property

$$P(X_{k+1}=S_j, X_k=S_i, \dots, X_0=S_1) = P(X_{k+1}=S_j, X_k=S_i)$$

## Kemeny Constant

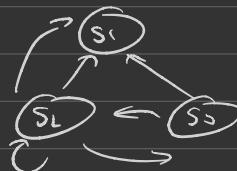
$$K_i = \sum_{j=1}^n \pi_j m_{ij} \quad K = \sum_{i=1}^n \frac{1}{1-\pi_i} \quad \text{with } j=1 \text{ bc } \lambda=1$$

*rainbow arrow*

- Can be used to find bridges by removing nodes
- Identifying different communities: calculate eigenvalues and see how close they are to 1

## Sharing Economy

- Markov Chains are useful for modelling macroscopic behaviour of systems (i.e. congestions, how it'd evolve)



state transition  
diagram

$$P = \begin{bmatrix} 0 & 0.5 & 0.5 \\ 0.1 & 0 & 0.9 \\ 0.5 & 0.5 & 0 \end{bmatrix}$$

matrix

$$\vec{x}^{(1)} = \vec{x}^{(0)} P \quad , \quad \vec{x}^{(2)} = \vec{x}^{(1)} P = \vec{x}^{(0)} P^2$$

this is the 'iteration', where  $k$  in  $\vec{x}^{(k)}$  denotes index.

Stochastic: sums to 1 bc of probability

Non-negative: special properties, diagonalizable for eigenvalues

## Finding Bridges

- means first passage time (MFPT)

- how well a state is connected to a graph

- kemeny constant

- how well connected is the graph

# Bias in Machine Learning

3 types (Fairness, Transparency, Privacy/Security)

discriminatory outcomes hard to know if change after outcome issue for data

i.e.) Bank loan (using AI to speculate)

~ is it related to gender, race, etc.

~ own or hidden, AI can guess community (i.e. living in China town)

~ rejected by AI system (can it explain why?)

~ COMPAS in US Courts (predicted likelihood)

~ 2016, bias against African American

~ reasoning:

- under representation (minority subgroups)

- historical bias (i.e. past gender stereotypes)

- model used is to minimize sum of loss (majority group takes larger share, minimize for majority)

## Fairness

- removing sensitive attribute does NOT work

- Statistical Parity (average fairness)  $P(F=f | A=0) = P(F=f | A=1)$  binary sensitive variable

improving - Independence

- Separation Same chance regardless of A

- Substitution

- Individual Fairness (treating similar individuals similarly)

- Causal Fairness

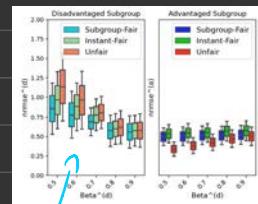
(i.e. Counterfactual: changing attributes to see if "accepted")

## Addressing Fairness

1) Pre-procedural (pre-train) - resampling

2) In-procedural (in-train) - minimizing problems

3) Post-procedural (after)



suffer from higher loss

Reinforcement Learning : optimal policy

Markov decision process: model the decision process  
- to maximize reward, etc.

$$V^*(s) = E_{\pi} [G_t \mid S_t = s]$$

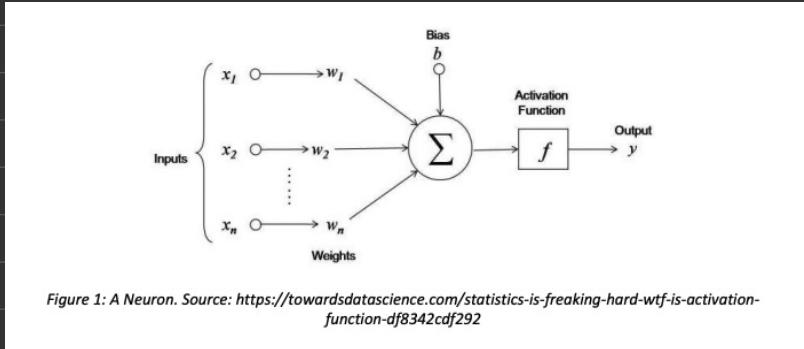
$$= R(s, \pi(s)) + \sum_{s' \in S} P(s, s', \pi(s')) r V^*(s')$$

Dynamic Programming: Bellman to find converging estimates

# LAB NOTES

# Lab 1 - Neural Network Basics

- Understand and code a simple neuron
  - Understand how a neuron learns
  - Understand its limitations



for this lab,  $n=2$  (number of dimensions)

⇒ Set up the weight at the neurons

$\Rightarrow$  if  $n=2$ ,  $g(x_1, x_2) = y$  with  $y \in \{0, 1\}$

$\Rightarrow$  Associated w/ "Truth Tables" OR AND  $\neg$

1

Teach neuron to learn: update weights  
w gradient descent and error function.

$$E = (\hat{y} - o)^2$$

$\uparrow \quad \uparrow$

Correct output      Current output

$$w_i' = w_i + n \left( \frac{dE}{dw_i} \right)$$

$$b' = b + n \left( \frac{dE}{db} \right)$$

$\Rightarrow$

learning rate

OR		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	1

AND		
$x_1$	$x_2$	$y$
0	0	0
0	1	0
1	0	0
1	1	1

$$E = (\hat{y} - o)^2$$

↑  
Correct output      ↑  
Current output

$$w'_i = w_i + \eta \left( \frac{dE}{dw_i} \right)$$

$$b' = b + \eta \left( \frac{dE}{db} \right)$$

learning rate

$$\sigma = \frac{1}{1 + e^{-(w^T x + b)}}$$

$$\frac{d\sigma}{dx} = \sigma(x) (1 - \sigma(x))$$

$$x = w^T x + b$$

$$w'_i = w_i + \eta \frac{dE}{dw_i}$$

$$\frac{dE}{dw} = \frac{dE}{d\sigma} \frac{d\sigma}{dx} \underbrace{\frac{dx}{dw}}_{= x}$$

$$\frac{dE}{d\sigma} = 2(\hat{y} - \sigma)$$

$$\frac{dE}{dw} = \frac{dE}{d\sigma} \times \frac{d\sigma}{dw}$$

$$= 2(\hat{y} - \sigma) \underbrace{\frac{\sigma(1 - \sigma)}{dw ds}}_{ds} \frac{x}{x}$$

$$\left\{ \begin{array}{l} \frac{dx}{dw} = x \\ \frac{d\sigma}{dx} = \sigma(x)(1 - \sigma(x)) \\ \therefore \frac{d\sigma}{dw} = (\sigma(x)(1 - \sigma(x)))x \end{array} \right.$$

$$W \leftarrow W + \eta \left[ \frac{dE}{dw_i} \right] \rightarrow \underbrace{\text{diff\_W.T}}_{\text{X.T. dot (dW)}} \rightarrow \text{X.T. dot (dW)}$$

$$\therefore X^T \cdot dW = \sum_{i=1}^n X^T \cdot dW_i = X_1 \cdot dW_1 + X_2 \cdot dW_2 + \dots + X_n \cdot dW_n$$

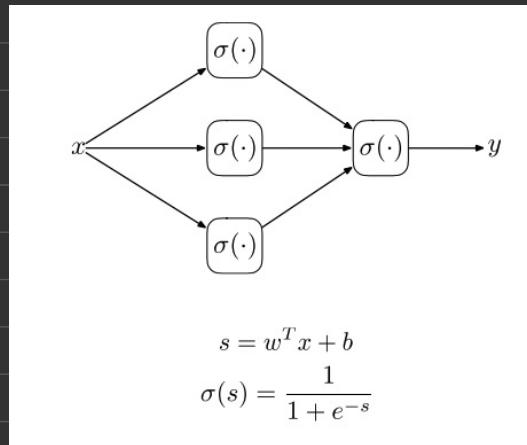
$$b'_i = b_i + \eta \frac{dE}{db}$$

$$\frac{dE}{db} = \frac{dE}{d\sigma} \frac{d\sigma}{dx} \frac{dx}{db} \quad \frac{dx}{db} = 1$$

$$= \frac{dE}{d\sigma} \frac{d\sigma}{dx} = 2(\hat{y} - \sigma) \underbrace{\frac{\sigma(1 - \sigma)}{dB ds}}_{ds}$$

## Lab 2 - Feedforward Network & Backpropagation

- Understand Backpropagation
- Write a neural network with  $\geq 1$  hidden layers
- Solve the XOR
- Understand how to build general classifiers



Use backpropagation due to hidden layer

$$E = (z - y)^2 \quad w_i' = w_i + \eta \frac{dE}{dw_i} \quad b' = b + \eta \frac{dE}{db}$$

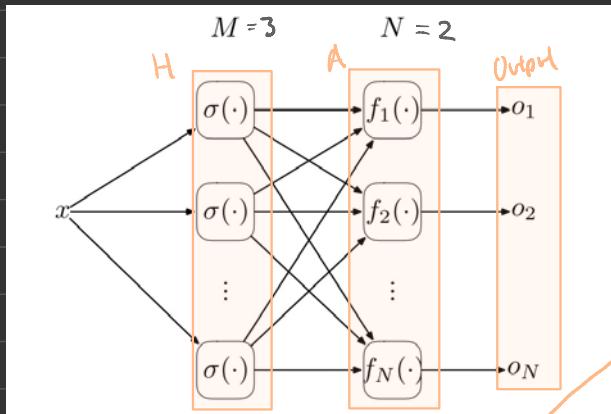
$$y = \sigma(z) \quad \sigma(s) = \frac{1}{1 + e^{-s}} \quad z = w^T x + b \quad \frac{d\sigma}{ds} = \sigma(1 - \sigma)$$

hidden layer  $db = b \cdot \eta \frac{dE}{db}$   $\frac{dE}{db} = \frac{dE}{ds} \frac{d\sigma}{ds} \frac{ds}{db} = (z - y) \cdot \sigma(1 - \sigma) \cdot (1)$

final layer  $\frac{dE}{dw} = \frac{dE}{ds} \frac{d\sigma}{ds} \frac{ds}{dw} = (z - y) \cdot \sigma(1 - \sigma) \cdot (x)$

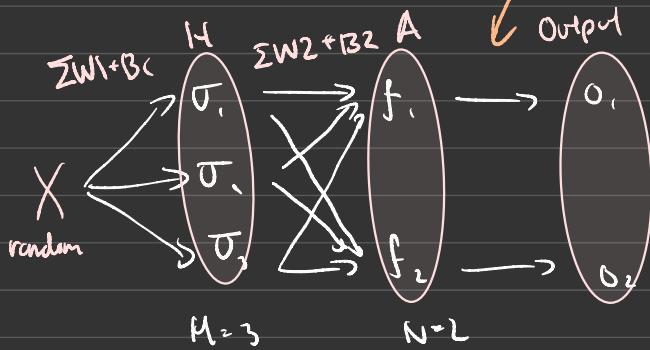
$$db = b \cdot \eta \frac{dE}{db} \cdot$$

Unlike Lab 1, sigm function is not used to update.



$$f_i = (\sum \sigma_i w_i) + b_2$$

$$o_i = f_i(s) = \frac{e^{s_i}}{\sum_{j=1}^n e^{s_j}}$$



$$A = \Sigma W H \rightarrow B$$

$$s = w^T x + b$$

Cross-Entropy Loss

$$\mathcal{H} = -\sum z_j \log(o_j)$$

new error function

$$\frac{dE}{db} = \frac{dE}{do} \frac{do}{ds} \frac{ds}{db} =$$

$$\frac{dE}{db} = \frac{dE}{do} \frac{do}{ds} \frac{ds}{db} =$$

# Lab 3 - PyTorch and Deep Networks

- Understand PyTorch to write a Neural Network
  - Write a Neural Network of Multiple Hidden Layers.

Pytorch "Autograd" uses tensor to automatically derive and keep track of operations.

dataset : MNIST

teach network to properly  
classify digits

000000000000000000000000  
1111111111111111111111  
2222222222222222222222  
3333333333333333333333  
4444444444444444444444  
5555555555555555555555  
6666666666666666666666  
7777777777777777777777  
888888888888888888888888  
999999999999999999999999

- Input is  $28 \times 28$  matrix (2-dimensional tensor)
    - ↳ corresponds to image pixels.
  - Output is integer b/w 0-9

Room E  
15:45 ~ 55  
March 10<sup>th</sup>

# Oral Prep

2 Questions

- > Class
- > Lab stuff