

MACHINE INTELLIGENCE

UNIT - 2

Supervised Learning with
KNN, ANN, SVM

feedback/corrections: vibha@pesu.pes.edu

- VIBHA MASTI



Lazy vs Eager Learner

(i) Lazy Learner

- Data processed only when question asked
- Eg: piggy bank

(ii) Eager Learner

- Each transaction processed in real time
- Queries are instant
- Model keeps updating

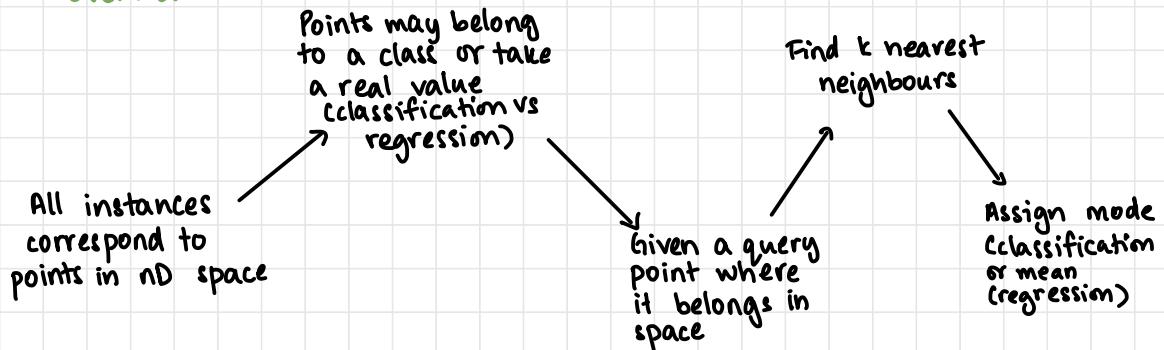
Instance BASED LEARNING

- Stores training examples and creates model when query made with new instance

K-Nearest Neighbours

- **Inductive bias:** class of instance x is most similar to class of instances closest to x
- Instance-based learning (lazy learning)

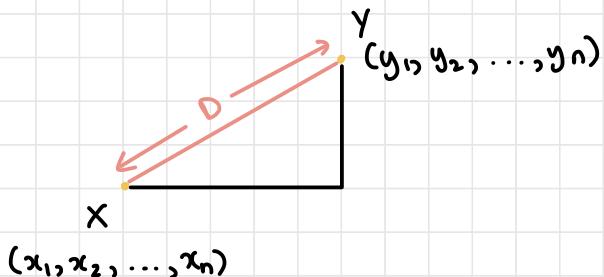
Overview



DISTANCE MEASURES

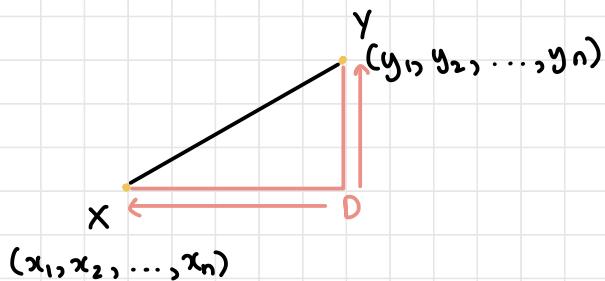
1. Euclidean distance

$$D(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^2 \right)^{1/2}$$



2. Manhattan distance

$$D(x, y) = \sum_{i=1}^n |x_i - y_i|$$



3. Minkowski distance

$$D(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

if $p=1$: Manhattan
if $p=2$: Euclidean

KNN for Regression & Classification

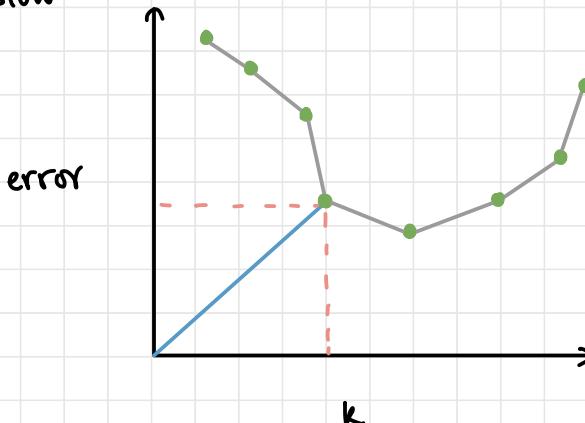
- Regression: real-valued target
- Classification: discrete-valued target

Choosing a Value of K

- K cannot be equal to the number of instances
- Usually, $K = \text{no. of classes} + 1$
 - in case of tie, decrement K by 1
- $K = \sqrt{n}$ also performs well ($n = \text{no. of instances in the dataset}$)

Elbow Method

- For every test instance, a curve of K vs error can be plotted and the min value of K with the least error is selected
- Elbow curve
- Very slow



Q: Find weight of ID 11. $k=3$. Use Euclidean Distance

ID	Height	Age	Weight	Diff in Height	Diff in Age	$\sqrt{\text{sum sq}}$
1	5	45	77	0.5	7	7.02
2	5.11	26	47	0.39	12	12.01
3	5.6	30	55	0.1	8	8.00
4	5.9	34	{ 59 }	0.4	4	4.02 ↙
5	4.8	40	{ 72 }	0.7	2	2.12 ↙
6	5.8	36	{ 60 }	0.3	2	2.02 ↙
7	5.3	19	40	0.2	19	19.00
8	5.8	28	60	0.3	10	10.00
9	5.5	23	45	0	15	15
10	5.6	32	58	0.1	6	6.00
11	5.5	38	?	0	0	0.00

- Regression (Continuous) problem: find average of k nearest values

$$\text{Weight} = 63.67$$

Error - Classification

- error = 1 if actual != predicted
- error = 0 otherwise

$$\text{Error}(E) = \frac{1}{n} \sum_{i=1}^n e_i$$

Error - Regression

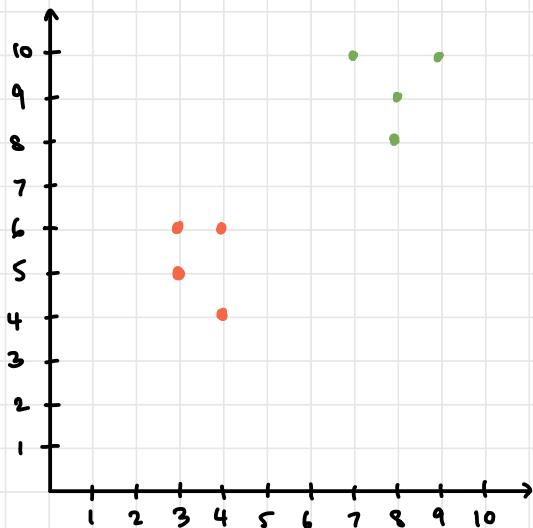
- error = $(y_{\text{pred}} - y_{\text{actual}})^2$
- MSE

Perform KNN and find error - classification

- 80:20 production and test set
- Upload prod data into memory
- Decide distance measure (Euclidean)
- Decide k value ($k=4$)
- For each data point in test, find distance wrt to all neighbours and choose 4 nearest neighbours
- Assign mode of neighbours to data point
- Calculate error and accuracy
- Produce model with achieved accuracy for a given k

TRAIN

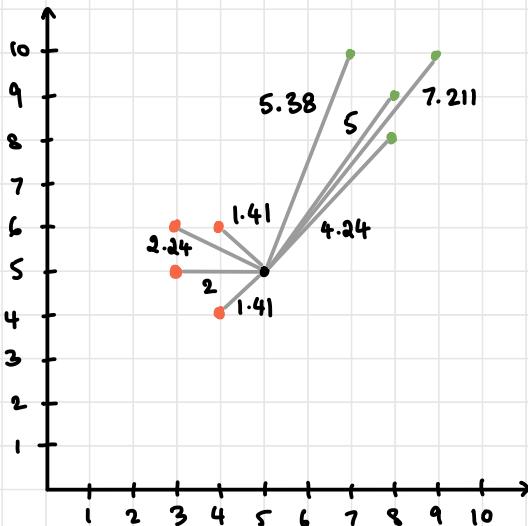
X	Y	Class
3	5	red
3	6	red
4	6	red
4	4	red
7	10	green
8	9	green
8	8	green
9	10	green



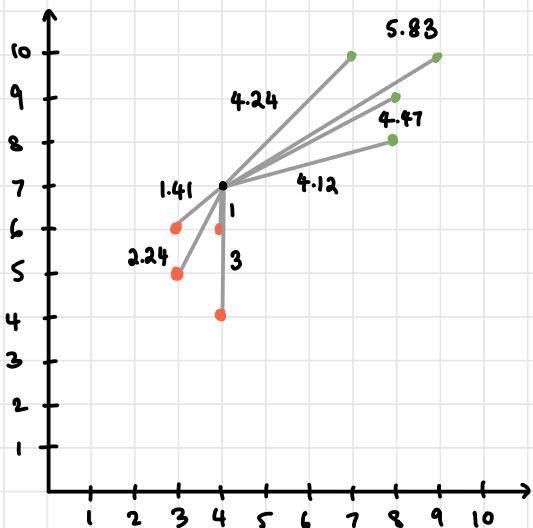
TEST

X	Y	Class
5	5	red
4	7	green

Using Euclidean Distance



$k=3$



$\text{pred} = \text{red}$
 $\text{actual} = \text{red}$
 $\text{error} = 0$

$\text{pred} = \text{red}$
 $\text{actual} = \text{green}$
 $\text{error} = 1$

Error - classification

$\text{error} = 0$ if $\text{actual} = \text{predicted}$
 $\text{error} = 1$ if $\text{actual} \neq \text{predicted}$

$$E = \frac{1}{n} \sum_{i=1}^n e_i$$

MSE - regression

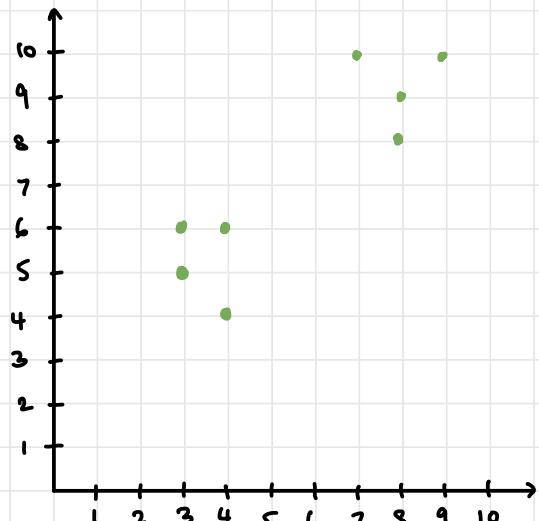
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

↑ observed
↑ predict

Perform KNN and find error - Regression

TRAIN

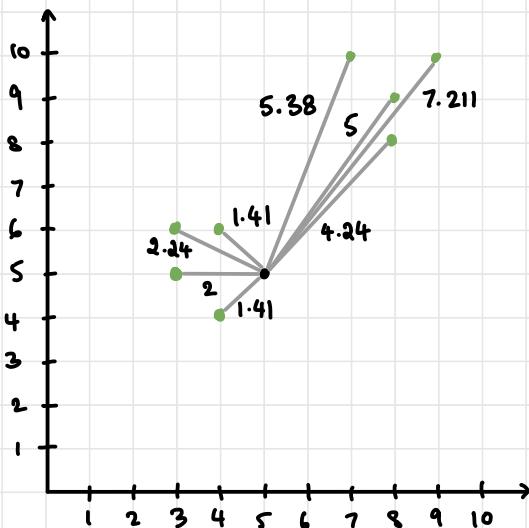
x	y	Class
3	5	3.5
3	6	2.6
4	6	2.8
4	4	3.1
7	10	8.1
8	9	8.5
8	8	9.3
9	10	9.7



TEST

x	y	class
5	5	3.7
4	7	8.4

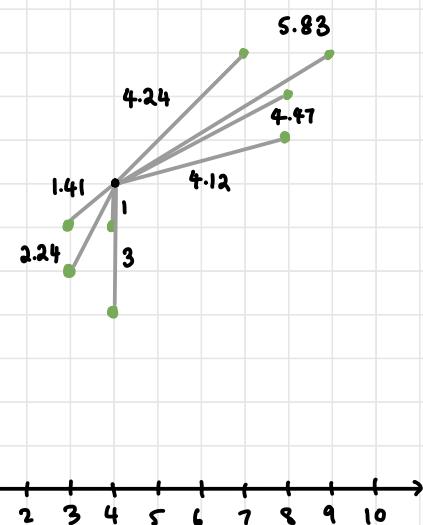
Using Euclidean Distance



$$\text{pred} = 3.13$$

$$\text{actual} = 3.7$$

K=3



$$\text{pred} = 2.967$$

$$\text{actual} = 8.4$$

$$\text{error}^2 = 0.3211$$

$$\text{error}^2 = 29.5211$$

$$\text{MSE} = 14.92$$

Distance Weighted kNN

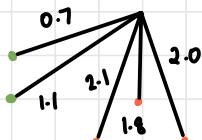
- Weigh the contribution of each of the k neighbours acc. to their distance to query point x_q .

$$\hat{f}(x_q) = \operatorname{argmax}_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i))$$

$$w_i = \frac{1}{d(x_q, x_i)}$$

distance measure used

Eg: $k=5$



$$\text{weights} = \frac{1}{\text{dists}}$$

$$\begin{aligned}\text{Support for red} &= \frac{1}{2.1} + \frac{1}{1.8} + \frac{1}{2.0} \\ &= 1.532\end{aligned}$$

$$\begin{aligned}\text{Support for green} &= \frac{1}{0.7} + \frac{1}{1.1} \\ &= 2.338\end{aligned}$$

Bias-Variance Trade-off

— K is very low

- chance of overfitting (high variance)

— K is very high

- chance for underfitting (high bias)
- label based on majority in dataset

— K neither too high or too low

- good trade-off

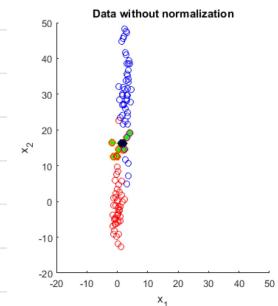
Pros of KNN

- No training required
- Only two parameters: K and distance measure

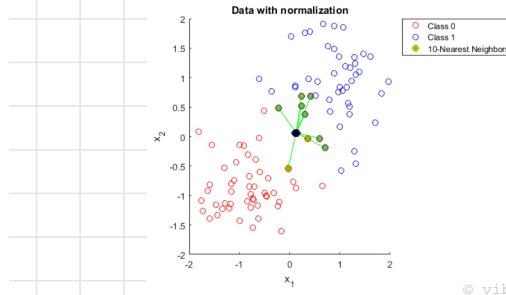
Cons of KNN

- Prediction time is computationally expensive
- Lazy algorithm — all data points to be stored
- Hard to scale
- Curse of dimensionality

Scaling



Class 0
Class 1
10-Nearest Neighbors



- Normalisation helps in uniformity
- Equal importance of each attribute
 - eg: salary & age - scales very different

ID	Age	Income(rupees)
1	25	80,000
2	30	100,000
3	40	90,000
4	30	50,000
5	40	110,000

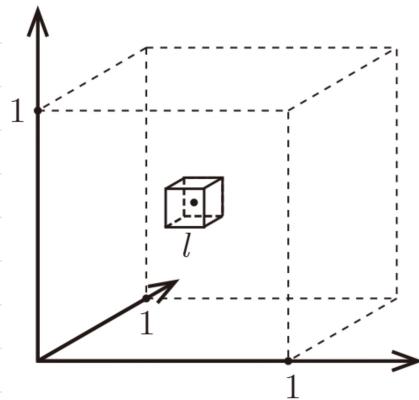
Pre normalizing

ID	Age	Income(rupees)
1	-1.192	-0.260
2	-0.447	0.608
3	1.043	0.173
4	-0.447	-1.563
5	1.043	1.042

Post normalizing

Curse of Dimensionality

- When no. of dimensions is too high, points drawn from probability distribution not close together
- In d-D: assume uniform random distribution in unit cube
- https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote02_kNN.html



- Try to find K nearest neighbours in d dimensions
- Assume all k neighbours lie in a hypercube of length ℓ (ℓ^d)

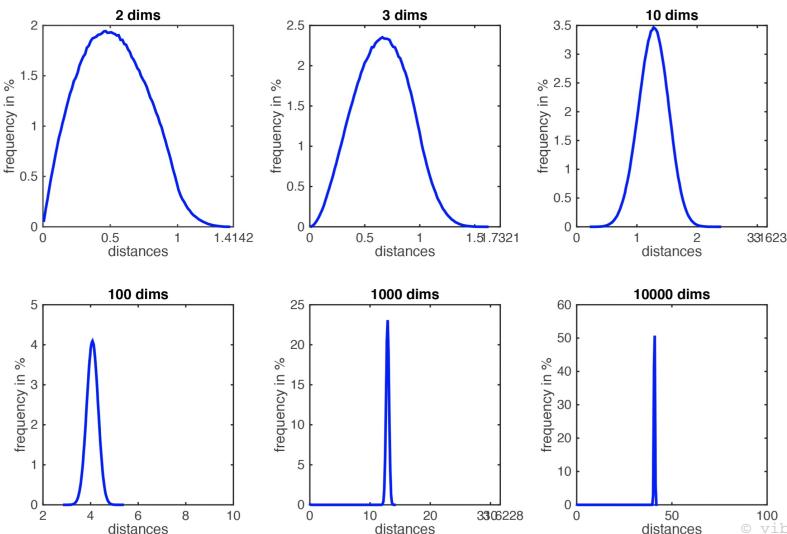
- The fraction of points in the cube ($n = \text{total points}$)

$$l^d \approx \frac{k}{n} \Rightarrow l \approx \left(\frac{k}{n}\right)^{1/d}$$

- If $n = 1000$, how big is l ?

d	l
2	0.1
10	0.63
100	0.955
1000	0.9954

- As $d \gg 0$, almost the entire space is required to find the k nearest neighbours
- Unknown point will not necessarily be 'similar' to its k nearest neighbours
- Histograms: distribution of pairwise distances between randomly distributed points as d grows



Solutions

1. Weighted attributes for calculating distances
 - eg: predicting house prices — distance gives less weightage to colour and more to locality
 - eg: predicting maths marks — distance gives less weightage to english marks
2. Cross-validation method: iteratively leaving out one of the attributes and testing
 - Best set of attributes
 - Long computational time
3. Principal Component Analysis for Dimensionality Reduction
4. Rule of thumb: 5 data instances per attribute for learning

Q: Is this a good way to calculate weights in a weighted KNN?

$$F(x, y) = \sum_{i=1}^n w_i f_i$$
$$w_i = \frac{h_i^{-p}}{\sum_{j=1}^n h_j^{-p}}$$

f_i : dataset values

h : distance between query point and a data point, say Euclidean distance.

Yes, as points closer are given heavier weights for the values of p

B:	x_1	x_2	Class	$k=3$
	7	7	bad	Euclidean distance
	7	4	bad	
	3	4	good	
O.P.:	1 3	4 7	good	No normalisation

(a) Vanilla

(b) Weighted

$$(a) |x_1 - x_{1,op}|^2 \quad |x_2 - x_{2,op}|^2 \quad \text{total}^2$$

16	0	16
16	9	25
0	9	9
4	9	13

mode = good (using vanilla kNN)

(b) Weighted

$ x_1 - x_{1,op} ^2$	$ x_2 - x_{2,op} ^2$	Euclidean	weight	class
16	0	4	0.25	bad
16	9	5	0.2	bad
0	9	3	0.33	good
4	9	$\sqrt{13}$	0.174	good

support for bad = 0.25

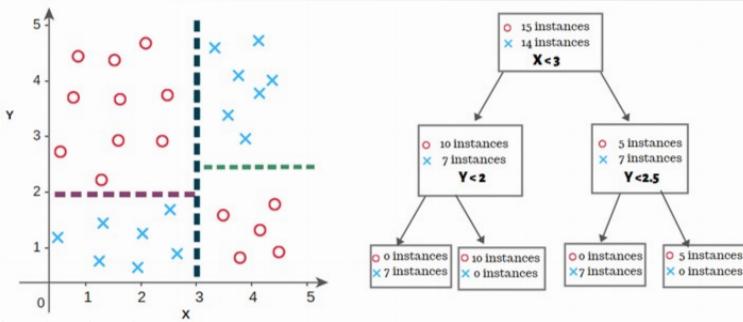
support for good = 0.504

Decision Boundaries

- Decision boundary: region of problem space where output label of classifier is ambiguous
- Separates underlying vector space into multiple sets, one for each class

Decision Boundary for Decision Trees

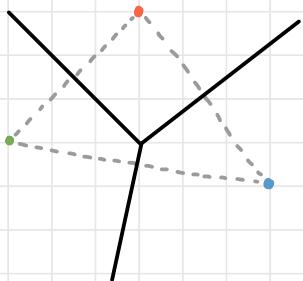
- Input space divided into axis-parallel rectangles
- Each rectangle: one class
- <https://towardsdatascience.com/decision-tree-overview-with-no-maths-66b256281e2b>



- Can represent arbitrary decision boundary

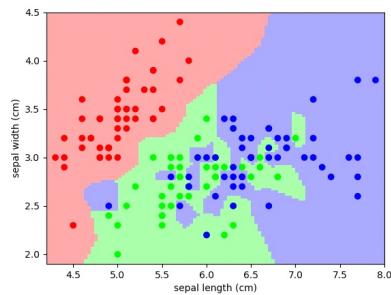
Decision Boundary for 1-KNN

- Voronoi tessellation
- Boundary drawn with points at the same distance from two training examples

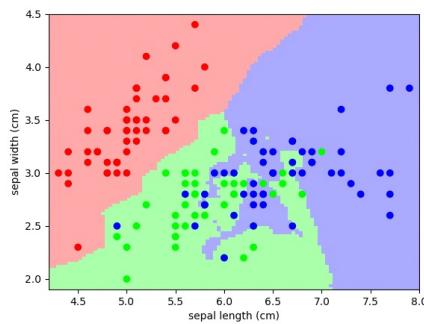


- ML Playground tool
- The decision boundary is a convex hull (fails with non-linear data)

$K=1$



$K=3$



source: scipy-lectures

- Piecewise non-linear decision boundary

Disadvantages of KNNs

- Arduous (strenuous) to find right K value
- Sensitive to noise
- Performance low for large datasets (calculating distances is linear)
- Performance low for high dimensionality (distances in each dim)
- High memory to store all training examples

Artificial Neural Network

Biological Neuron

Dendrites
cell nucleus (soma)
Axon
Synapse

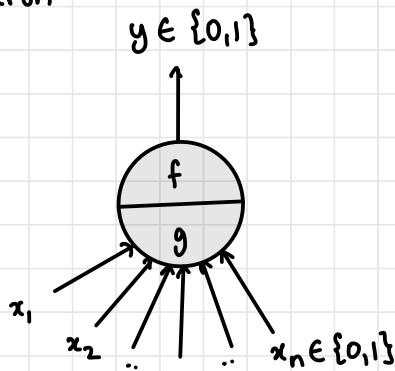
Artificial Neuron

Input
Node
Output
Interconnection

- Input layer, hidden layers, output layer

McCulloch Pitts Neuron - Logic gates

- Simple neuron

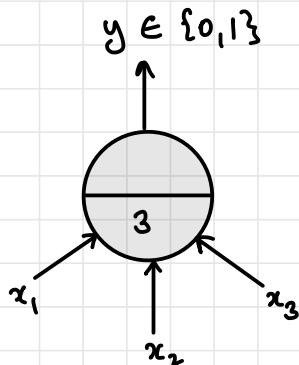


- g aggregates the input function and f takes a decision based on the aggregation
- Inputs can be excitatory or inhibitory
- If any x_i is inhibitory, $y=0$. Else

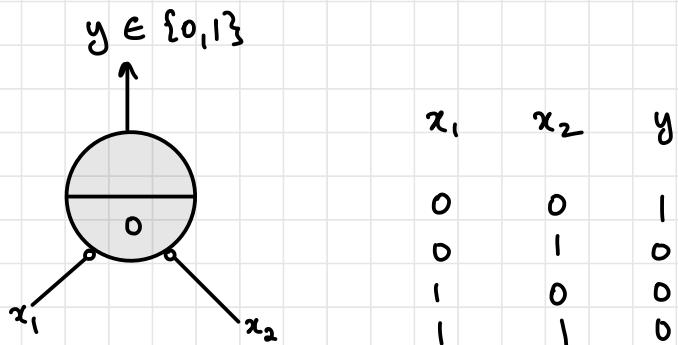
$$g(x_1, x_2, \dots, x_n) = g(x) = \sum_{i=1}^n x_i$$

$$\begin{aligned} y = f(g(x)) &= 1 && \text{if } g(x) \geq \theta \\ &= 0 && \text{if } g(x) < \theta \end{aligned}$$

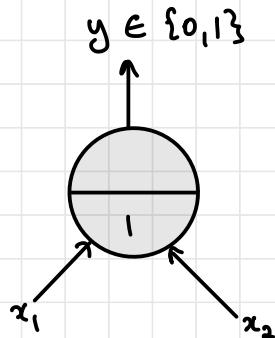
- AND function



- NOR function ($\circ \rightarrow$ inhibition)



- OR function



Geometric Representation of MP Neuron

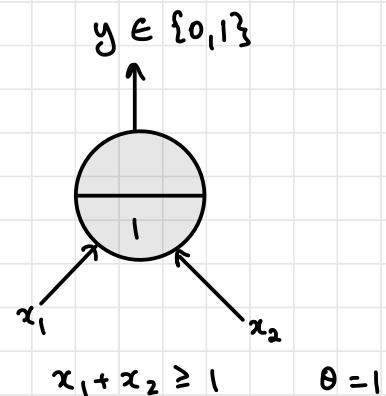
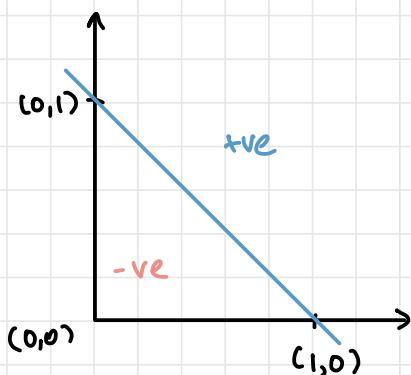
- Single MP neuron splits input points into two parts

Points lying on or above the line $\sum_{i=1}^n x_i - \theta = 0$

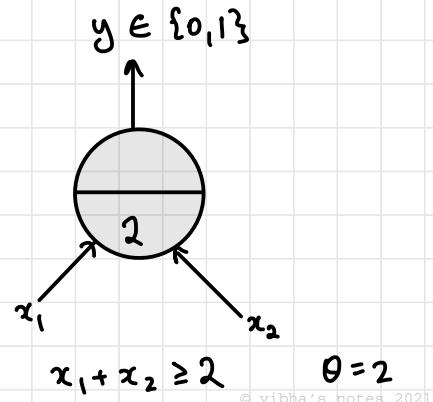
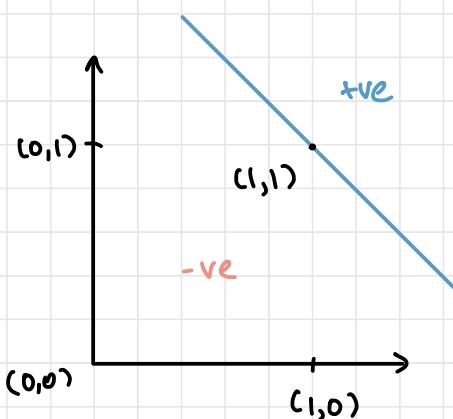
- 0: below line } must be linearly
1: above line } separable

- Cannot show XOR

(a) OR

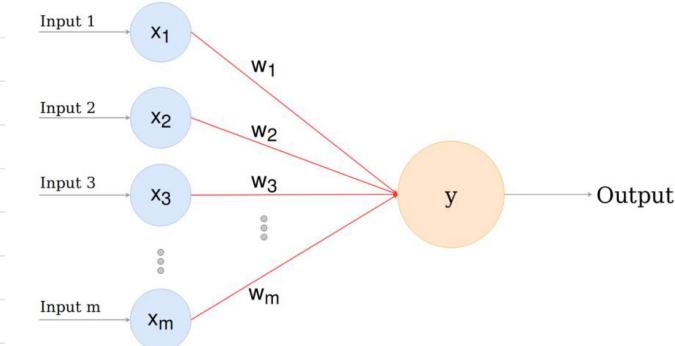


(b) AND



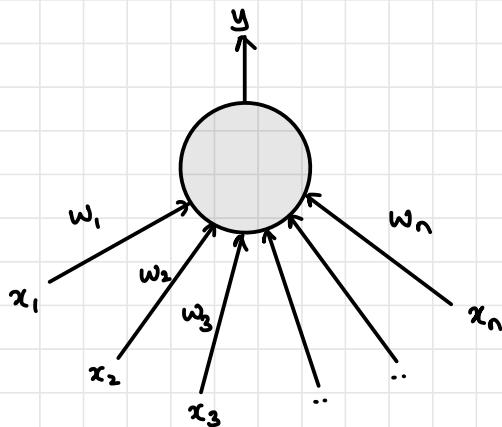
Perceptron

- More general - proposed by Frank Rosenblatt



The Perceptron Model

- Every input has an associated weight
- Mechanism for learning the weights and biases
- Inputs are numerical, not boolean



$$y = 1 \text{ if } \sum_{i=1}^n w_i x_i - \theta \geq 0 \quad \text{bias is threshold}$$

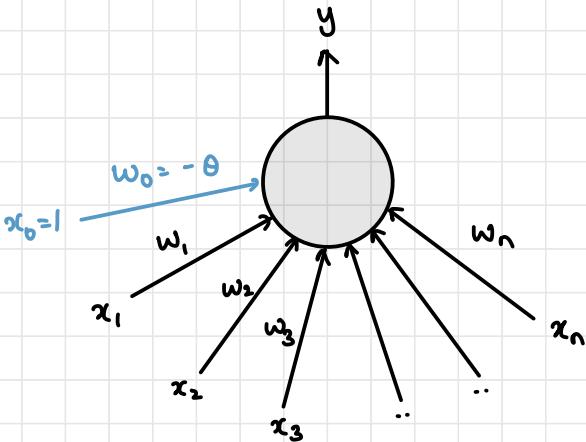
$$= 0 \text{ if } \sum_{i=1}^n w_i x_i - \theta < 0$$

A Mathematical Convention

$$y = 1 \text{ if } \sum_{i=0}^n w_i x_i \geq 0$$

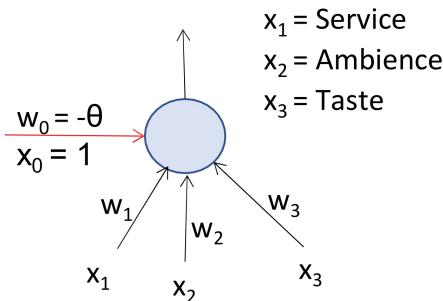
$$= 0 \text{ if } \sum_{i=0}^n w_i x_i < 0$$

Where $w_0 = -\theta$ and $x_0 = 1$
(bias)



An Intuitive Example

- Setting a threshold: eating at a restaurant



- Threshold θ maybe low or high depending on the customer's standards

Perceptron Function

$$f(x) = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

Comparing MP Neurons and Perceptrons

Differences

- Weights and threshold can be learnt for perceptron
- Inputs for perceptron maybe real-valued

Similarities

- Separate input space into 2 classes
- All inputs that give 1 lie on one side and those that give 0 lie on the other side
- Output is binary

Logical OR

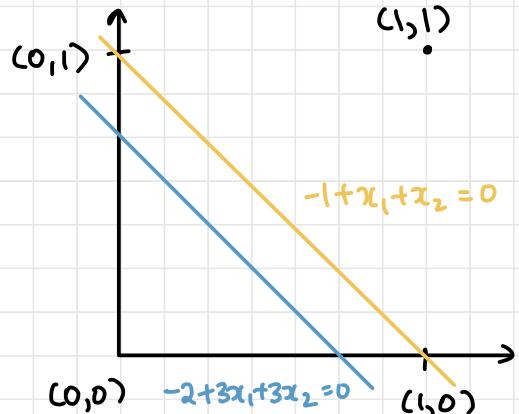
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

$$\begin{aligned} w_0 &= \text{bias} \\ w_1 & \\ w_2 & \end{aligned} \quad \left. \begin{array}{l} \text{unknown} \\ \text{unknown} \end{array} \right\}$$

Negative:

$$w_0 x_0 + w_1 x_1 + w_2 x_2 < 0$$

$$w_0 < 0 \rightarrow ①$$



Positive

$$w_0 x_0 + w_1 x_1 + w_2 x_2 \geq 0$$

$$w_0 + w_1 \geq 0 \rightarrow ②$$

$$w_0 + w_2 \geq 0 \rightarrow ③$$

$$w_0 + w_1 + w_2 \geq 0 \rightarrow ④$$

Possible solution: $w_0 = -1$, $w_1 = 1$, $w_2 = 1$

Perceptron equation:

$$\hat{y} = w_0x_0 + w_1x_1 + w_2x_2$$

$$\hat{y} = -1 + x_1 + x_2 : y = \begin{cases} 1, & \hat{y} \geq 0 \\ 0, & \hat{y} < 0 \end{cases}$$

Another possibility

$$w_0 = -2, w_1 = 3, w_2 = 3$$

$$-2 + 3x_1 + 3x_2 = 0$$

Logical AND

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

w_0 = bias
 w_1, w_2 } unknown

Negative:

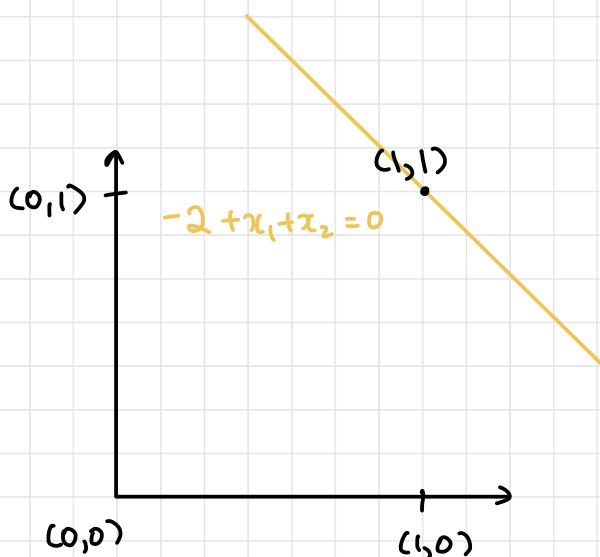
$$w_0x_0 + w_1x_1 + w_2x_2 < 0$$

$$\begin{aligned} w_0 < 0 &\rightarrow ① \\ w_0 + w_1 < 0 &\rightarrow ② \\ w_0 + w_2 < 0 &\rightarrow ③ \end{aligned}$$

Positive

$$w_0x_0 + w_1x_1 + w_2x_2 \geq 0$$

$$w_0 + w_1 + w_2 \geq 0 \rightarrow ④$$



Possible solution: $w_0 = -2, w_1 = 1, w_2 = 1$

Perceptron Learning Algorithm

Algorithm: Perceptron Learning Algorithm

```
P ← inputs with label 1;  
N ← inputs with label 0;  
Initialize w randomly;  
while !convergence do  
    Pick random x ∈ P ∪ N ;  
    if x ∈ P and w.x < 0 then  
        | w = w + x ;  
    end  
    if x ∈ N and w.x ≥ 0 then  
        | w = w - x ;  
    end  
end  
//the algorithm converges when all the  
inputs are classified correctly
```

source: (lecture 2)

<http://www.cse.iitm.ac.in/~miteshk/CS6910.html>

Check YouTube: CS7015 IITM
Deep Learning

Lecture 2.1 to 2.7



https://www.youtube.com/watch?v=4TC5s_xNKSS&list=PLZEpw4xZOspJEn5s4wUj6uA1mw2_dk0ap

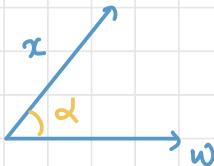
Working of Algorithm

- Initialise weights w with random vector
- Iterate over all positive & negative training examples (P ∪ N)
- If an input x belongs to P (positive training class), the value of the dot product $w \cdot x$ should be ≥ 0
 - So that perceptron predicts positive class
- If an input $x \in N$, the value of the dot product $w \cdot x$ should be < 0
 - So that perceptron predicts negative class
- Case 1: $x \in P$ and $w \cdot x < 0$: increase the values in w so next iteration of $w \cdot x$ might be ≥ 0
- Case 2: $x \in N$ and $w \cdot x \geq 0$: decrease the values in w so next iteration of $w \cdot x$ might be < 0

Why it Works

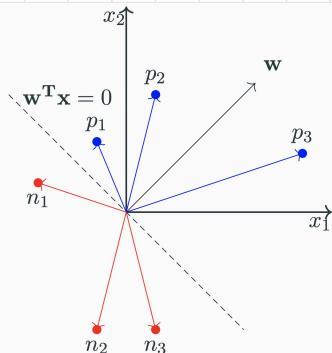
- If $x \in P$ and $w \cdot x \geq 0$, then the angle between w & x should be $\leq 90^\circ$

$$w \cdot x = |w| |x| \cos \alpha$$



$$\alpha \leq 90^\circ \Rightarrow \cos \alpha \geq 0$$

- Similarly, if $x \in N$, $w \cdot x < 0$ and $\alpha > 90^\circ$



source: IITM CS6910

Intuition for why the update works

(α_{new}) when $w_{new} = w + x$

$$\begin{aligned} \cos(\alpha_{new}) &\propto w_{new}^T x \\ &\propto (w + x)^T x \\ &\propto w^T x + x^T x \\ &\propto \cos \alpha + x^T x \end{aligned}$$

$$\cos(\alpha_{new}) > \cos \alpha$$

α is decreasing

(α_{new}) when $w_{new} = w - x$

$$\begin{aligned} \cos(\alpha_{new}) &\propto w_{new}^T x \\ &\propto (w - x)^T x \\ &\propto w^T x - x^T x \\ &\propto \cos \alpha - x^T x \end{aligned}$$

$$\cos(\alpha_{new}) < \cos \alpha$$

α is increasing

Logical XOR

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

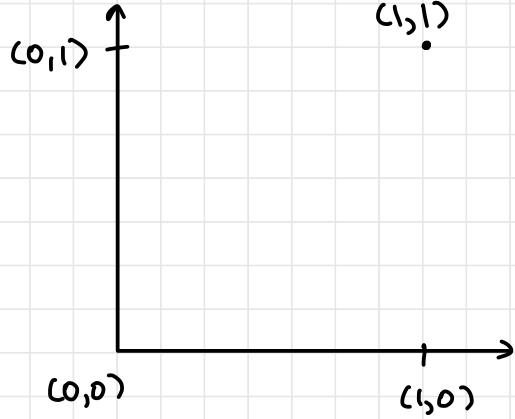
$$\begin{aligned} w_0 &= \text{bias} \\ w_1, w_2 & \quad \} \text{ unknown} \end{aligned}$$

Negative:

$$w_0 x_0 + w_1 x_1 + w_2 x_2 < 0$$

$$w_0 < 0 \rightarrow ①$$

$$w_0 + w_1 + w_2 < 0 \rightarrow ②$$



Positive

$$w_0 x_0 + w_1 x_1 + w_2 x_2 \geq 0$$

$$w_0 + w_1 \geq 0 \rightarrow ③$$

$$w_0 + w_2 \geq 0 \rightarrow ④$$

No solution! Contradiction

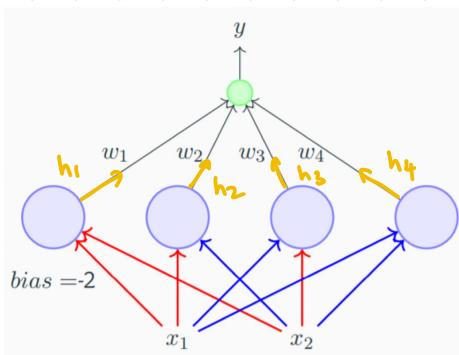
Single Perception

- Fails when data linearly inseparable
- Linear decision surface
- Problems: Prof Preet's slides (L2.2)

Multi-Layer Perceptron Network - ANN

- If data is not linearly separable, we can use a network of neurons
- Assume input True = +1, False = -1 (not 1 & 0)

hidden {
layer}



Red edge: $w = -1$
Blue edge: $w = +1$

- Can implement any Boolean function

- https://www.youtube.com/watch?v=CJr4Dst0uZE&list=PLZEpw4xZOspJEn5s4wUj6uA1mw2_dk0ap&index=17

very good video ↗

- Perceptron 1 fires when $x_1 = -1 \& x_2 = -1$
- Perceptron 2 fires when $x_1 = -1 \& x_2 = 1$
- Perceptron 3 fires when $x_1 = 1 \& x_2 = -1$
- Perceptron 4 fires when $x_1 = 1 \& x_2 = 1$

XOR Using Multi-layered Perceptron

x_1	x_2	y	h_1	h_2	h_3	h_4	$\hat{y} = \sum_{i=0}^4 w_i h_i$
-1	-1	0	1	0	0	0	$w_0 + w_1$
-1	1	1	0	1	0	0	$w_0 + w_2$
1	-1	1	0	0	1	0	$w_0 + w_3$
1	1	0	0	0	0	1	$w_0 + w_4$

Positive

$$w_0 + w_2 \geq 0$$

$$w_0 + w_3 \geq 0$$

Negative

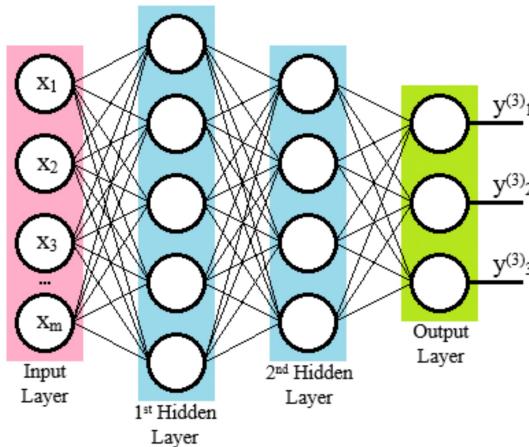
$$w_0 + w_1 < 0$$

$$w_0 + w_4 < 0$$

$$w_0 = -1, w_1 = -1, w_2 = 1, w_3 = 1, w_4 = -1$$

Multilayered ANN

- Deep ANN: more than 1 hidden layer



i) Input Layers

- attributes for one sample
- fed into network

2) Hidden Layers

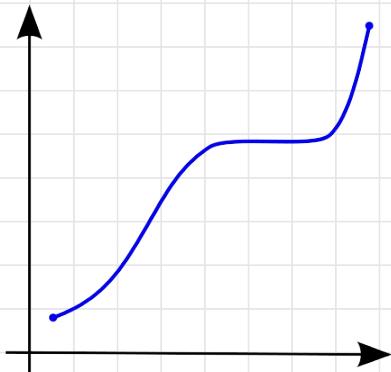
- responsible for deriving complex relationships between input and output

3) Output Layer

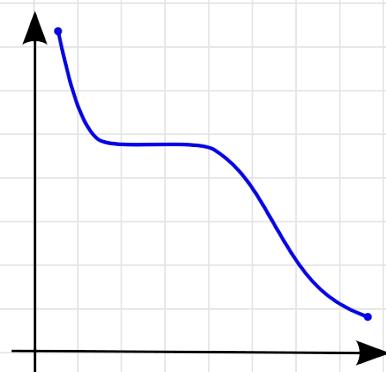
- emits output

Monotonic Functions

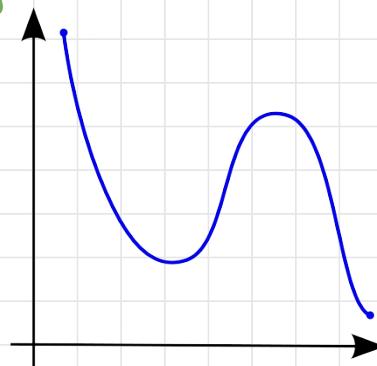
source: wikipedia



monotonically
non-decreasing



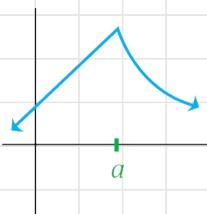
monotonically
non-increasing



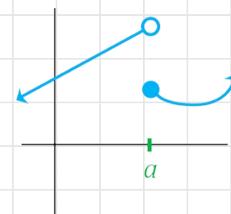
non-monotonic

Differentiability

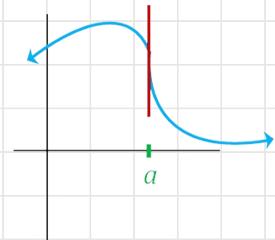
source: calc workshop



Cusp / Corner



Discontinuous



Vertical Tangent

ACTIVATION FUNCTIONS

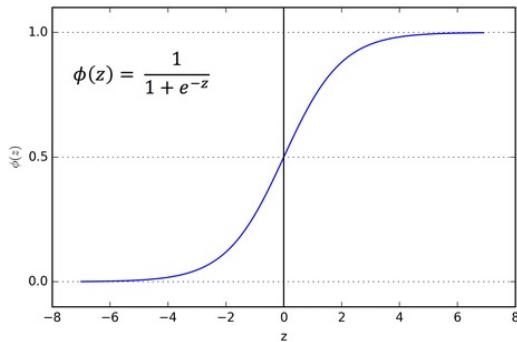
- Non-linear transformation applied to input before propagating it to the next layer of neurons

$$Y = \text{Activation} \left(\sum_{i=0}^n w_i x_i + b \right)$$

- If there are no activation functions, behaviour of NN is purely linear (linear regression)
- <https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

I. Sigmoid

$$f(x) = \frac{1}{1+e^{-x}}$$



source: towards data science

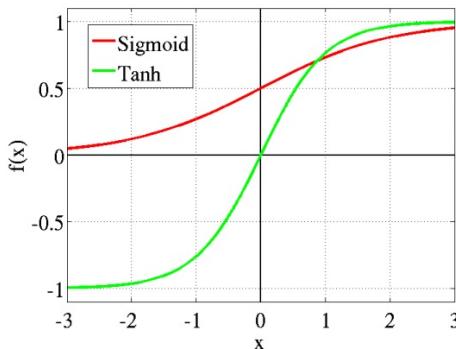
- Range . $(0, 1)$
- Differentiable, monotonic (not its derivative)
- Can cause NN to get stuck during training (problem of vanishing gradient)

2. tanh

$$f(x) = \tanh(x)$$

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- Range: $(-1, 1)$
- Better than sigmoid as negative maps to negative
- Differentiable, monotonic (not its derivative)



source: towards data science

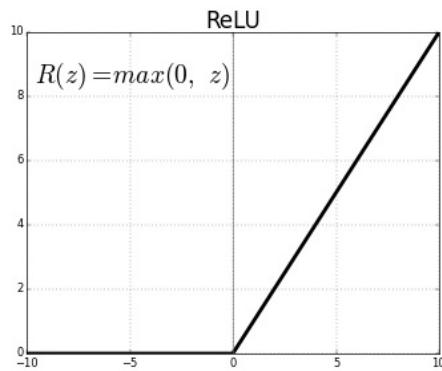
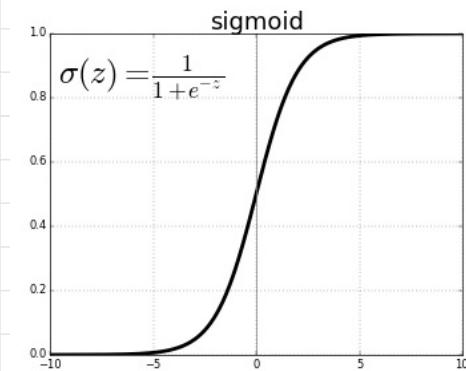
3. ReLU (Rectified Linear Unit)

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

$$f(x) = \begin{cases} x & , x > 0 \\ 0 & , x \leq 0 \end{cases}$$

$$f'(x) = \begin{cases} 1 & , x > 0 \\ 0 & , x \leq 0 \end{cases}$$

- Defined to be differentiable even though it is not at $x=0$

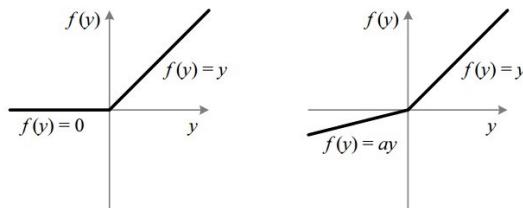


source: towards data science

- Range: $[0, \infty)$
- Does not map negative values properly — dying ReLU problem

4. Leaky ReLU

- leak increases range and tries to solve dying ReLU problem



source: towards data science

$$f(x) = \begin{cases} x & , x > 0 \\ ax & , x \leq 0 \end{cases}$$

- a is usually 0.01

$$f'(x) = \begin{cases} 1 & , x > 0 \\ \alpha & , x \leq 0 \end{cases}$$

- Defined to be differentiable even though it is not at $x=0$
- Cheatsheet (at $x=0$, arbitrary decisions made)

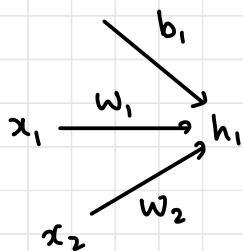
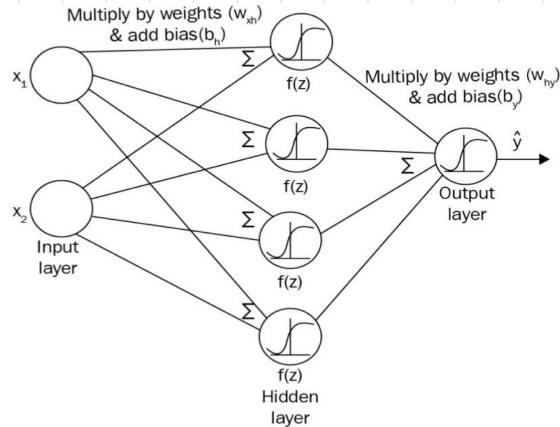
Name	Plot	Function, $f(x)$	Derivative of f , $f'(x)$	Range
Identity		x	1	$(-\infty, \infty)$
Binary step		$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	$\begin{cases} 0 & \text{if } x \neq 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$	{0, 1}
Logistic, sigmoid, or soft step		$\sigma(x) = \frac{1}{1 + e^{-x}}$ [1]	$f(x)(1 - f(x))$	$(0, 1)$
Hyperbolic tangent (tanh)		$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$1 - f(x)^2$	$(-1, 1)$
Rectified linear unit (ReLU) ^[9]		$\begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max\{0, x\} = x\mathbf{1}_{x>0}$	$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$	$[0, \infty)$
Gaussian Error Linear Unit (GELU) ^[4]		$\frac{1}{2}x \left(1 + \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) \right)$ $= x\Phi(x)$	$\Phi(x) + x\phi(x)$	$(-0.17\dots, \infty)$
Softplus ^[10]		$\ln(1 + e^x)$	$\frac{1}{1 + e^{-x}}$	$(0, \infty)$
Exponential linear unit (ELU) ^[11]		$\begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter α	$\begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ 1 & \text{if } x = 0 \text{ and } \alpha = 1 \end{cases}$	$(-\alpha, \infty)$
Scaled exponential linear unit (SELU) ^[12]		$\lambda \begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ with parameters $\lambda = 1.0507$ and $\alpha = 1.67326$	$\lambda \begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	$(-\lambda\alpha, \infty)$
Leaky rectified linear unit (Leaky ReLU) ^[13]		$\begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$	$\begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	$(-\infty, \infty)$

source: Wikipedia

https://en.wikipedia.org/wiki/Activation_function

Forward Propagation

- Total no. of layers in ANN is no. of hidden layers + output layers

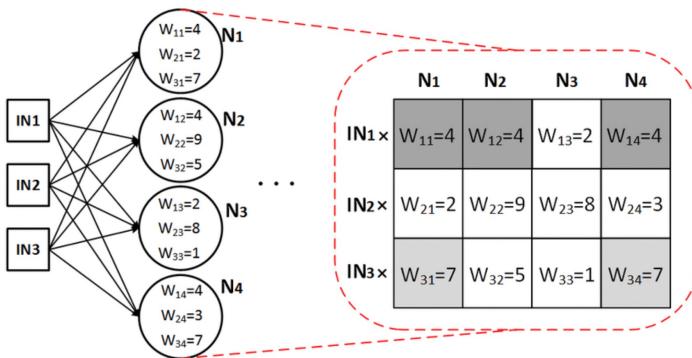


$$h_1 = w_1x_1 + w_2x_2 + b_1$$

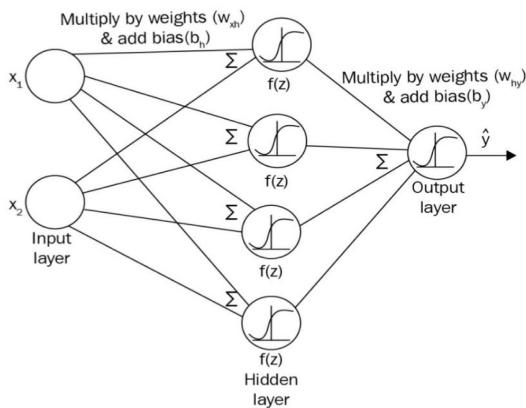
output = activation(h_1)

- Weight matrix : weights from previous layer to current layer
- Dimensions = no. of units in prev layer \times no. of units in current layer

Calculation of weight Matrix



Forward Propagation

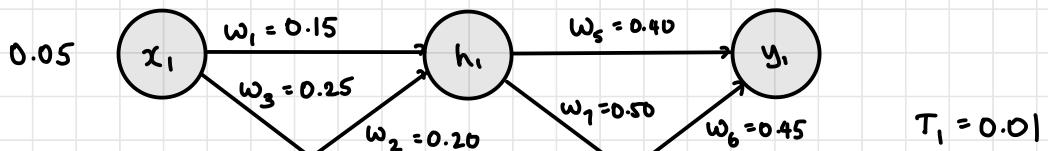


- w_{xh} : weight between input and hidden layer
- b_h : bias of hidden layer
- $Z_i = x \cdot w_{xh} + b_h$: netsum
- $a_i = \sigma(z_i)$: activation function : output of hidden layer

- w_{hy} : weight between hidden layer and output layer
- b_y : bias of output layer
- $Z_2 = a_i \cdot w_{hy} + b_y$: netsum
- $\hat{y} = \sigma(z_2)$: sigmoid activation : final output
- Loss: error function

$$L = \text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Example



$$\begin{aligned} T_1 &= 0.01 \\ T_2 &= 0.99 \end{aligned}$$

$$0.10$$

$$b_1 = 0.35$$

$$b_2 = 0.60$$

$$h_1 = (x_1 w_1 + x_2 w_2 + b_1) = (0.05 \times 0.15 + 0.10 \times 0.20 + 0.35) = 0.3775$$

$$h_2 = (x_1 w_3 + x_2 w_4 + b_2) = (0.05 \times 0.25 + 0.10 \times 0.30 + 0.60) = 0.3925$$

$$O_{h_1} = \sigma(h_1) = \frac{1}{1+e^{-0.3775}} = 0.5933$$

$$O_{h_2} = \sigma(h_2) = \frac{1}{1+e^{-0.3925}} = 0.5969$$

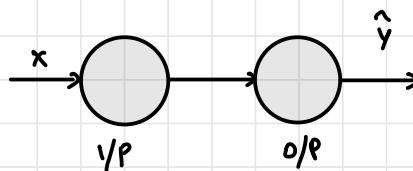
$$\hat{y}_1 = a(0_{h_1}w_5 + 0_{h_2}w_6 + b_2) = a(0.5933 \times 0.40 + 0.5969 \times 0.45 + 0.60) \\ = a(1.1059) \\ = 0.7514$$

$$\hat{y}_2 = a(0_{h_1}w_7 + 0_{h_2}w_8 + b_2) = a(0.5933 \times 0.5 + 0.5969 \times 0.55 + 0.60) \\ = a(1.2249) \\ = 0.7729$$

$$\text{Loss} = \frac{1}{2} \sum (y - \hat{y})^2 \quad [\text{MSE}] \\ = \frac{1}{2} \left((0.7514 - 0.01)^2 + (0.7729 - 0.99)^2 \right) \\ = 0.2984$$

Backpropagation

- Iteratively modify weights to minimize loss function between actual and predicted values



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

$$\hat{y} = a(z) = \frac{1}{1+e^{-z}} = (1+e^{-z})^{-1}$$

$$z = wx + b$$

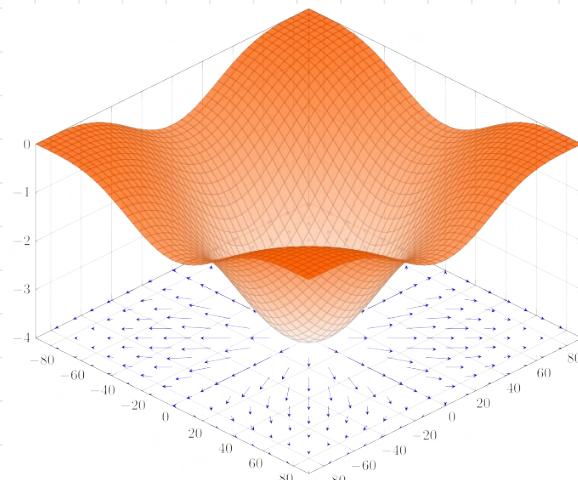
Review of Vector Calculus

nabla; $\nabla = \frac{\partial}{\partial x} \hat{i} + \frac{\partial}{\partial y} \hat{j} + \frac{\partial}{\partial z} \hat{k}$
del operator

- If ϕ is a scalar function $\phi(x, y, z)$, then $\nabla\phi$ is the gradient of ϕ

$$\nabla\phi = \frac{\partial\phi}{\partial x} \hat{i} + \frac{\partial\phi}{\partial y} \hat{j} + \frac{\partial\phi}{\partial z} \hat{k}$$

- If $\phi = f(x, y, z)$, $\nabla\phi$ is a vector that points in the direction of the steepest slope
- Gradient of $f(x, y) = -(\cos^2 x + \cos^2 y)^2$ depicted as a projected vector field on the bottom plane



source: wikipedia

$$\nabla f = \frac{\partial f}{\partial x} \hat{i} + \frac{\partial f}{\partial y} \hat{j}$$

GRADIENT DESCENT

- Minimising the cost function by taking steps in the opposite direction of the gradient

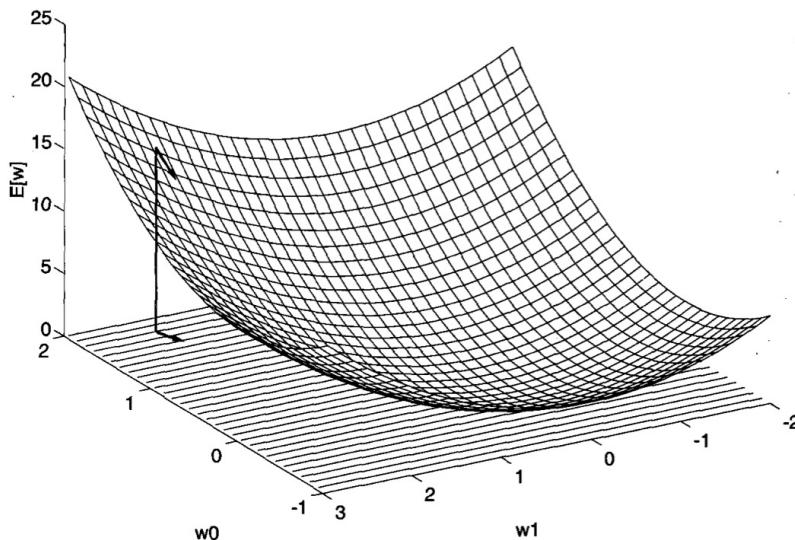


FIGURE 4.4

Error of different hypotheses. For a linear unit with two weights, the hypothesis space H is the w_0, w_1 plane. The vertical axis indicates the error of the corresponding weight vector hypothesis, relative to a fixed set of training examples. The arrow shows the negated gradient at one particular point, indicating the direction in the w_0, w_1 plane producing steepest descent along the error surface.

T2

- At each step, weight vector is altered in the direction that produces steepest descent along the error surface
- Gradient of error function wrt weight vector points in the direction of steepest ascent
- $-\nabla E$ in direction of steepest descent

Gradient Descent

- The weights are updated according to the rule

$$\vec{\omega} = \vec{\omega} + \Delta \vec{\omega}$$

- Where $\Delta \vec{\omega}$ is given by

$$\Delta \vec{\omega} = -\eta \nabla E(\vec{\omega})$$

- η : positive constant called the learning rate
- Component-Wise

$$w_i = -\eta \frac{\partial E}{\partial w_i}$$

MSE / LSE Loss Function

- For a single output unit

$$E(\vec{\omega}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- D — set of training examples
- t_d — target (expected) value
- o_d — output (predicted) value
- $\frac{1}{2}$ — to simplify calculations

- Multiple output units

$$E(\vec{\omega}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$$

Cross-Entropy Loss Function

$$E = \sum_{i=1}^n -t_i \ln(o_i)$$

- n = no. of classifications
- $o_i = a(z)$
- $z = wx + b$

Binary Cross-Entropy

$$L = -[Y \ln(\hat{y}) + (1-Y) \ln(1-\hat{y})] \quad \left. \right\} \text{both notations used}$$

$$\hat{y} = a(z)$$

$$z = wx + b$$

Derivative of L wrt w

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w}$$

$$\frac{\partial L}{\partial w} = \left(\frac{-Y}{\hat{y}} + \frac{(1-Y)}{1-\hat{y}} \right) \frac{\partial \hat{y}}{\partial w}$$

$$= \frac{\hat{y} - Y}{\hat{y}(1-\hat{y})} \underbrace{\frac{\partial \hat{y}}{\partial w}}_{\frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w}} = a(z)(1-a(z))x \quad \text{sigmoid}$$

$$\frac{\partial L}{\partial w} = \frac{\hat{y} - Y}{\hat{y}(1-\hat{y})} a(z)(1-a(z))x$$

Backpropagation Algorithm

- Employs gradient descent to minimise MSE

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$$

Terminology

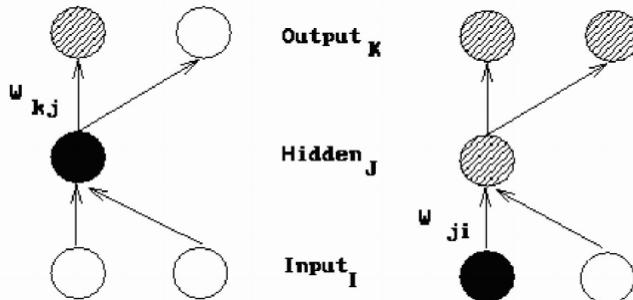
- x_{ji} = the i th input to unit j
- w_{ji} = the weight associated with the i th input to unit j
- $\text{net}_j = \sum_i w_{ji}x_{ji}$ (the weighted sum of inputs for unit j)
- o_j = the output computed by unit j
- t_j = the target output for unit j
- σ = the sigmoid function
- outputs = the set of units in the final layer of the network
- $\text{Downstream}(j)$ = the set of units whose immediate inputs include the output of unit j

T2

Derivation of SGD

- For sigmoid activation functions

<https://www.cs.swarthmore.edu/~meeden/cs81/s10/BackPropDeriv.pdf>



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

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

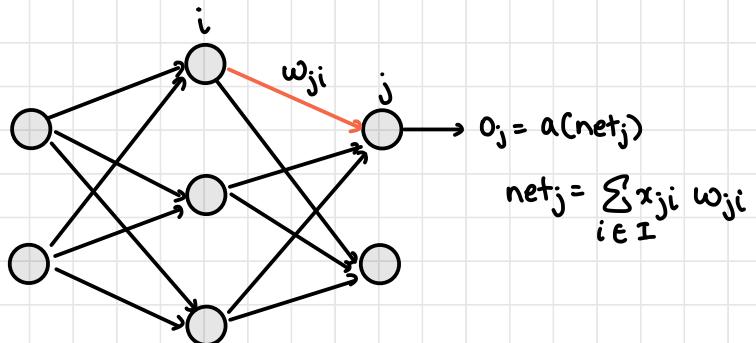
$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

Case 1: Training Rule for Output Unit Weights

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} \rightarrow ①$$

(a) (b) (c)



(a) $\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$

$$= \frac{1}{2} \sum_{k \in \text{outputs}} \frac{\partial}{\partial o_j} (t_k - o_k)^2$$

only when $j=k$, $\frac{\partial}{\partial o_j} (t_k - o_k)^2 \neq 0$

$$\frac{\partial E_d}{\partial o_j} = \frac{1}{2} \frac{\partial}{\partial o_j} (t_j - o_j)^2$$

$$\frac{\partial E_d}{\partial o_j} = -(t_j - o_j) \longrightarrow \textcircled{2}$$

(b) $\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial \sigma(\text{net}_j)}{\partial \text{net}_j}$

Derivative of $\sigma(x)$

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

$$\ln \sigma = -\ln(1 + e^{-x})$$

$$\frac{1}{\sigma} \frac{d\sigma}{dx} = \frac{1}{1 + e^{-x}} \cdot e^{-x} = \frac{e^{-x} + 1}{1 + e^{-x}} \cdot \frac{-1}{1 + e^{-x}}$$

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

$$\frac{\partial o_j}{\partial \text{net}_j} = o_j(1 - o_j) \longrightarrow \textcircled{3}$$

Combining $\textcircled{2}$ & $\textcircled{3}$

$$\frac{\partial E_d}{\partial \text{net}_j} = -\underbrace{(t_j - o_j)o_j(1 - o_j)}_{\delta_j} \longrightarrow \textcircled{4}$$

(c) $\frac{\partial \text{net}_j}{\partial w_{ji}} = \frac{\partial \sum_{i \in I} w_{ji} x_{ji}}{\partial w_{ji}} = x_{ji} \rightarrow (5)$

Combining (4) & (5) in (1)

$$\frac{\partial E_d}{\partial w_{ji}} = -(t_j - o_j) o_j (1-o_j) x_{ji}$$

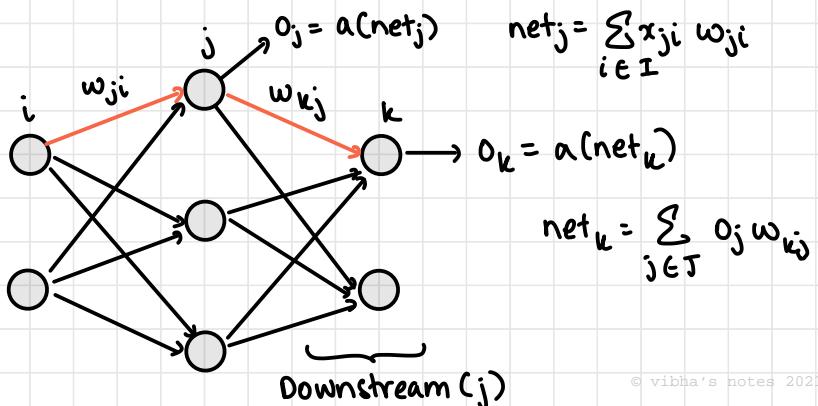
Weight Updation

$$\Delta w_{ji} = \eta (t_j - o_j) o_j (1-o_j) x_{ji} = \eta \delta_j x_{ji}$$

Case 2: Training Rule for Hidden Unit weights

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \sum_{k \in \text{Downstream}(j)} \underbrace{\frac{\partial E_d}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial \text{net}_j}}_{(a)} \underbrace{\frac{\partial \text{net}_j}{\partial w_{ji}}}_{(b)} \rightarrow (6)$$



(a)
$$\underbrace{\frac{\partial E_d}{\partial \text{net}_k}}_{-\delta_k} \frac{\partial \text{net}_k}{\partial \text{net}_j} = -\delta_k \frac{\partial \text{net}_k}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \longrightarrow ⑦$$

$$= -\delta_k w_{kj} o_j (1-o_j)$$

(b)
$$\frac{\partial \text{net}_j}{\partial w_{ji}} = x_{ji} \longrightarrow ⑧$$

Combining ⑦ & ⑧ in ⑥

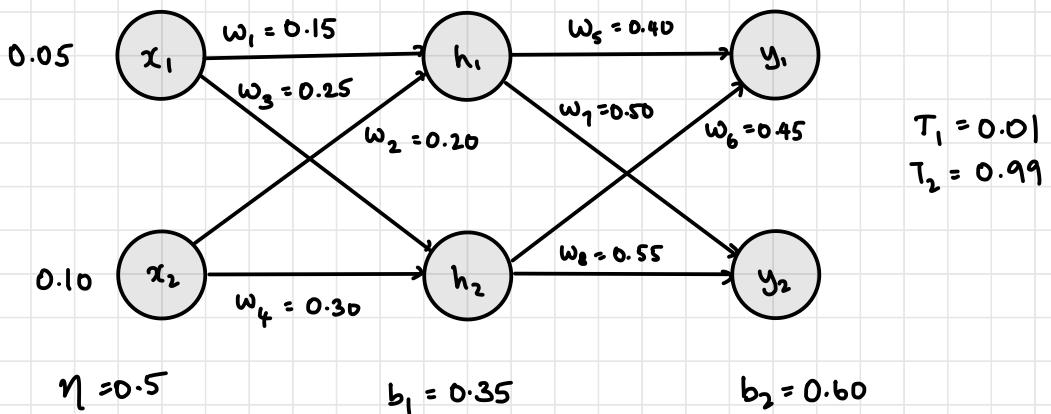
$$\frac{\partial E_d}{\partial w_{ji}} = \left(\sum_{k \in \text{Downstream}(j)} -\delta_k w_{kj} o_j (1-o_j) \right) x_{ji}$$

Weights Updation

$$\Delta w_{ji} = \eta x_{ji} \left(\sum_{k \in \text{Downstream}(j)} \delta_k w_{kj} o_j (1-o_j) \right)$$

where $\delta_k = (t_k - o_k) o_k (1-o_k)$

Q: Backpropagate once



$$\hat{y}_1 = 0.7514 \quad O_{h_1} = a(h_1) = \frac{1}{1+e^{-0.3775}} = 0.5933$$

$$\hat{y}_2 = 0.7729$$

$$\text{Loss} = 0.2984$$

$$O_{h_2} = a(h_2) = \frac{1}{1+e^{-0.3925}} = 0.5969$$

$$w_5 = w_5 + \Delta w_5$$

$$= w_5 + \eta (0.01 - 0.7514)(0.7514)(1 - 0.7514) 0.5933$$

$$= 0.40 - \eta (0.08216)$$

$$w_5 = 0.3589$$

$$w_6 = w_6 + \Delta w_6$$

$$\begin{aligned} &= w_6 + \eta (0.01 - 0.7514)(0.7514)(1 - 0.7514)(0.5969) \\ &= 0.45 + \eta (-0.0827) \end{aligned}$$

$$w_6 = 0.4087$$

$$w_7 = w_7 + \Delta w_7$$

$$\begin{aligned} &= w_7 + \eta (0.99 - 0.7729)(0.7729)(1 - 0.7729)(0.5933) \\ &= 0.50 + 0.1130 \\ &= 0.5113 \end{aligned}$$

$$w_8 = w_8 + \Delta w_8$$

$$\begin{aligned} &= w_8 + \eta (0.99 - 0.7729)(0.7729)(1 - 0.7729)(0.5969) \\ &= 0.55 + 0.1137 \\ &= 0.5614 \end{aligned}$$

$$w_i = w_i + \Delta w_i$$

$$\Delta w_i = \eta \left[\sum_{k \in \text{Down}(j)} (t_k - o_k) o_k (1 - o_k) w_{kj} \right] x_{ji} o_j (1 - o_j)$$

$$\begin{aligned} w_i &= 0.15 + 0.5 \left[(0.01 - 0.7514)(0.7514)(1 - 0.7514)(0.40) + \right. \\ &\quad \left. (0.99 - 0.7729)(0.7729)(1 - 0.7729)(0.50) \right] 0.05 \\ &\quad (0.5933)(1 - 0.5933) \end{aligned}$$

$$w_1 = 0.15 + 0.5 \times 0.01206 \times -0.0363$$

$$w_1 = 0.1498$$

$$w_2 = 0.20 + 0.5 \left[\frac{(0.01 - 0.7514)(0.7514)(1 - 0.7514)(0.40)}{(0.99 - 0.7729)(0.7729)(1 - 0.7729)(0.50)} \right] 0.10 \\ (0.5933)(1 - 0.5933)$$

$$w_2 = 0.20 + 0.01206 \times -0.0363$$

$$w_2 = 0.1996$$

$$-0.0414$$

$$w_3 = 0.25 + 0.5 \left[\frac{(0.01 - 0.7514)(0.7514)(1 - 0.7514)(0.45)}{(0.99 - 0.7729)(0.7729)(1 - 0.7729)(0.55)} \right] 0.05 \\ (0.5969)(1 - 0.5969)$$

$$w_3 = 0.2498$$

$$w_4 = 0.30 + 0.5 \left[\frac{(0.01 - 0.7514)(0.7514)(1 - 0.7514)(0.45)}{(0.99 - 0.7729)(0.7729)(1 - 0.7729)(0.55)} \right] 0.10 \\ (0.5969)(1 - 0.5969)$$

$$w_4 = 0.2995$$

Overfitting

- Fitting of noise
- Solution: validation set

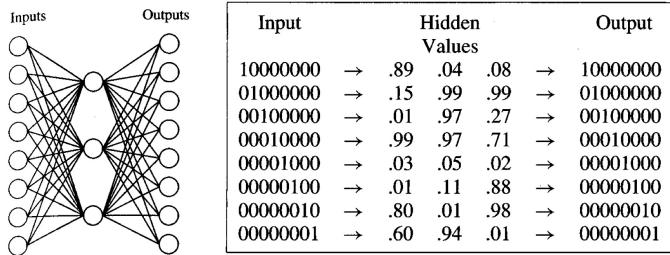


FIGURE 4.7

Learned Hidden Layer Representation. This $8 \times 3 \times 8$ network was trained to learn the identity function, using the eight training examples shown. After 5000 training epochs, the three hidden unit values encode the eight distinct inputs using the encoding shown on the right. Notice if the encoded values are rounded to zero or one, the result is the standard binary encoding for eight distinct values.

- Look up: Universal Approximation Theorem

Gradient Descent Algorithm

- Repeat until convergence (gradient = 0 at min)

$$w_j := w_j - \frac{\alpha \times \partial E(w_0, w_1)}{\partial w_j} \quad \text{for } j=0,1$$

- Update weights

α = learning rate

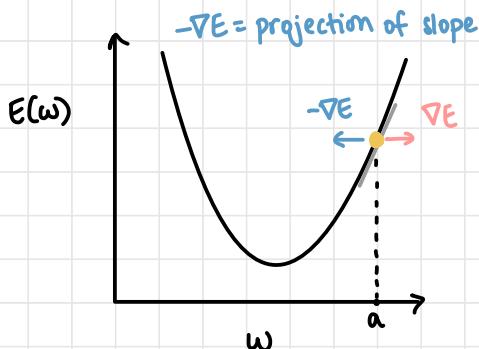
$$temp_0 = w_0 - \left(\frac{\alpha \times \partial J(w_0, w_1)}{\partial w_0} \right)$$

$$temp_1 = w_1 - \left(\frac{\alpha \times \partial J(w_0, w_1)}{\partial w_1} \right)$$

$$w_0 = temp_0$$

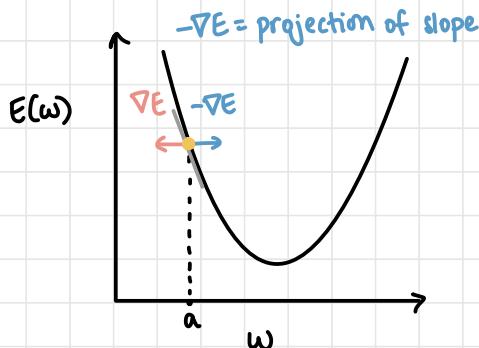
$$w_1 = temp_1$$

Visualization in 2D



$$E(w) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

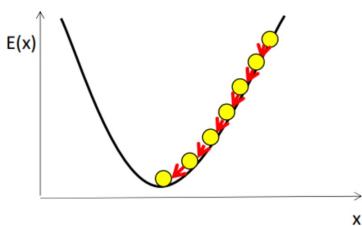
$$w = w - \alpha (\text{positive number})$$



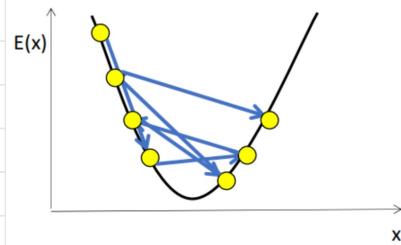
$$E(w) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$w = w - \alpha (\text{negative number})$$

Effect of Learning Rate α



too small

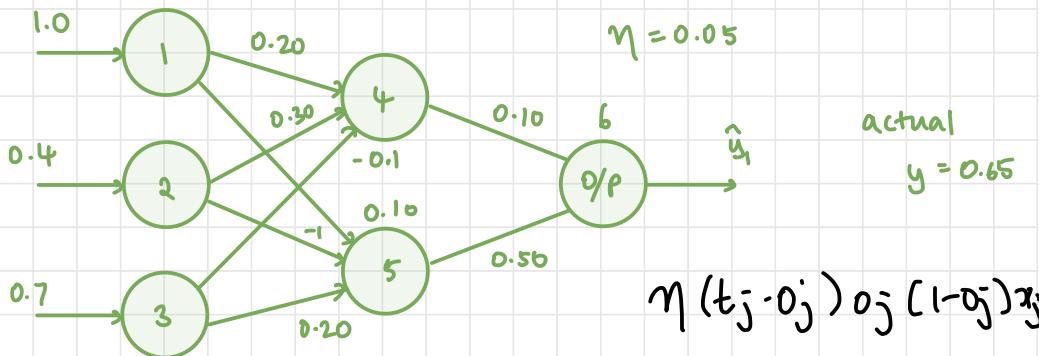


too large

Drawback of GD

- Can get stuck in local minimum

Q: Perform one pass of FP and one of BP (assume σ)



$$x = \begin{bmatrix} 1.0 \\ 0.4 \\ 0.7 \end{bmatrix} \quad w = \begin{bmatrix} 0.2 & 0.10 \\ 0.3 & -1 \\ -0.1 & 0.20 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{aligned} \text{Net} &= w^T x + b = \begin{bmatrix} 0.2 & 0.3 & -0.1 \\ 0.10 & -1 & 0.20 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.4 \\ 0.7 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0.25 \\ -0.16 \end{bmatrix} \end{aligned}$$

$$z = \sigma(x) = \begin{bmatrix} 0.5622 \\ 0.4601 \end{bmatrix}$$

$$\text{Net}_2 = w_2^T z + b_2 = [0.10 \quad 0.50] \begin{bmatrix} 0.5622 \\ 0.4601 \end{bmatrix} = 0.2863$$

$$\hat{y} = 0.5711 = o_j$$

$$y = 0.65$$

$$\text{Loss} = \frac{1}{2} (0.65 - 0.5711)^2 = 3.113 \times 10^{-3}$$

Back propagation

$$w_{64} = w_{64} + \Delta w$$

$$= 0.10 + 0.05 ((t_j - o_j) o_j (1 - o_j) x_{ji})$$

$$= 0.10 + 0.05 (0.65 - 0.5711) 0.5711 (0.5622)$$

$$= 0.10 + 0.05 \times 0.0253$$

$$= 0.101$$

$$w_{65} = w_{65} + \Delta w$$

$$= 0.50 + 0.05 (0.65 - 0.5711) 0.5711 (0.4601)$$

$$= 0.5010$$

$$w_{41} = 0.20 + \Delta w$$

$$= 0.20 + 0.05 \left[\sum_{k \in DS} [(t_k - o_k) o_k (1 - o_k) w_{kj}] x_{ji} (1 - o_j) o_j \right]$$

$$= 0.20 + 0.05 \left[(0.65 - 0.5711) 0.5711 (1 - 0.5711) (0.10) \right]$$

$$(1.0) (1 - 0.5622) 0.5622$$

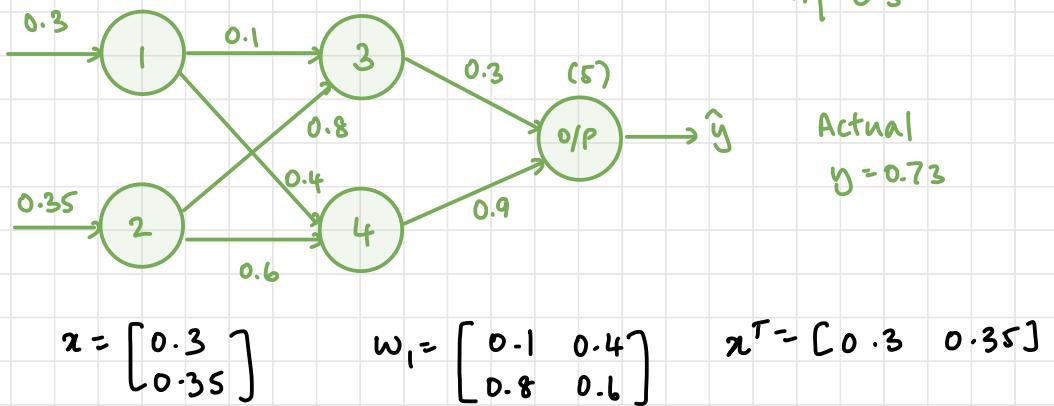
$$= 0.20 + 2.3784 \times 10^{-5}$$

$$= 0.2000$$

$$\begin{aligned}
w_{42} &= w_{42} + \Delta w \\
&= 0.30 + 0.05 \left[\sum_{k \in DS} [(t_k - o_k) o_k (1 - o_k) w_{kj}] x_{ji} (1 - o_j) o_j \right] \\
&= 0.30 + 0.05 \left[(0.65 - 0.5711) 0.5711 (1 - 0.5711) 0.1 \right] \\
&\quad (0.40) (1 - 0.5622) 0.5622 \\
&= 0.30 + 9.51 \times 10^{-6} \\
&= 0.30 \\
\\
w_{43} &= w_{43} + \Delta w \\
&= -0.1 + 0.05 \left(\sum_{k \in DS} [(t_k - o_k) o_k (1 - o_k) w_{kj}] x_{ji} (1 - o_j) o_j \right) \\
&= -0.1 + 0.05 \left[(0.65 - 0.5711) 0.5711 (1 - 0.5711) 0.1 \right] \\
&\quad (0.70)(0.5622)(1 - 0.5622) \\
&= -0.09998
\end{aligned}$$

and so on

Q: Perform one pass of FP and one of BP (assume σ) for w_{53} and w_{31}



$$\text{net}_2 = \begin{bmatrix} 0.31 \\ 0.33 \end{bmatrix} \quad a_2 = \begin{bmatrix} 0.5769 \\ 0.5818 \end{bmatrix} \quad o = 0.8261$$

Backpropagation

$$\begin{aligned} \Delta w_{53} &= \eta (t_j - o_j) o_j (1 - o_j) x_{ji} \\ &= 0.5 (0.73 - 0.8261) 0.8261 (1 - 0.8261) 0.5769 \\ &= -3.9822 \times 10^{-3} \end{aligned}$$

$$w_{53} = 0.2960$$

$$\begin{aligned} \Delta w_{31} &= \eta o_j (1 - o_j) x_{ji} \sum_{k \in \text{down}(j)} (t_k - o_k) o_k (1 - o_k) w_{kj} \\ &= 0.5 (0.5769) (1 - 0.5769) 0.3 \left[(0.73 - 0.8261) 0.8261 (1 - 0.8261) (0.3) \right] \\ &= -1.5164 \times 10^{-4} \end{aligned}$$

$$w_{31} = 0.0998$$

Mathematical Foundations for SVMs

- Build models to map input to output data
- ML models \rightarrow optimisation problem
 - minimisation of $f(x)$
 - maximisation of $f(x) \Rightarrow$ minimisation of $-f(x)$

<https://web.stanford.edu/group/sisl/k12/optimization/#lindex.md>

I. Unconstrained Optimisation

- Find x that minimises $f(x)$; global minima
- Stationary / critical point: where $f'(x) = 0$
 - local minimum
 - local maximum
 - saddle point
- Second derivative
 - $f''(x) < 0$: local maximum
 - $f''(x) > 0$: local minimum
 - $f''(x) = 0$: could be saddle point

HESSIAN MATRIX

- Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ with input vector $x \in \mathbb{R}^n$ and output scalar $f(x) \in \mathbb{R}$
- If all partial derivatives of f exist and are continuous, then $H_{n \times n}$ is defined as

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix},$$

source: wikipedia

- or $(\mathbf{H}_f)_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j}$.

- For 2 variables

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

- x is a vector \Rightarrow must find gradient
- If H_{ij} is positive definite: local min
- If H_{ij} is negative definite: local max
- If H_{ij} is indefinite: saddle point

Q: Find the Hessian of $f(x,y) = 8x - 2x^2y^2$. Also find gradient.

$$\nabla f = (8 - 4xy^2)\hat{i} - 4x^2y\hat{j} = [8 - 4xy^2 \quad -4x^2y]$$

$$H = \begin{bmatrix} -4y^2 & -8yx \\ -8yx & -4x^2 \end{bmatrix}$$

2. Constrained Optimisation

- Constraints place limits on the domain of the variables
- Eg: marathon runner

$$L(x, \lambda, \alpha) = f(x) + \sum_i \lambda_i g^i(x) + \sum_j \alpha_j h^{(j)}(x)$$

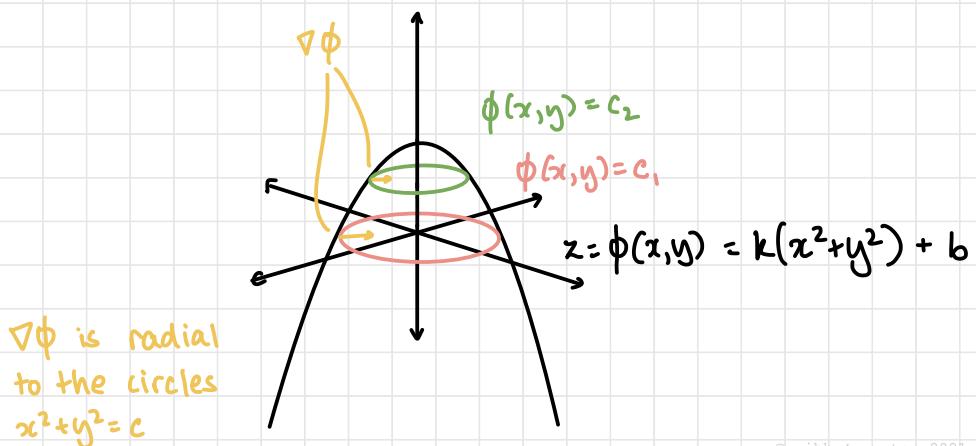
- Constrained minimum

$$\min_{x \in S} f(x) \quad \min_x \max_{\lambda, \alpha \geq 0} L(x, \lambda, \alpha)$$

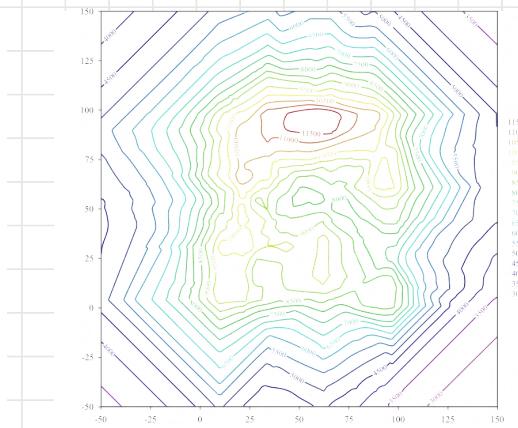
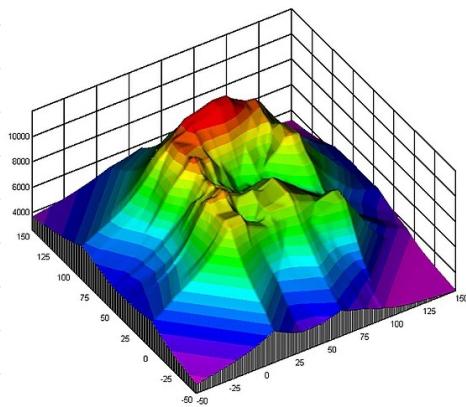
- Constraint surface places limit on value of function to be optimised

Review of Vector Calculus (part II)

- If $\phi(x, y, z)$ is a function, $\nabla \phi$ = the direction of max ascent (a vector)



- On a level surface $\phi(x,y,z) = c$, the gradient at any point is perpendicular to the surface
- This is because the value of $\nabla\phi$ along the surface is 0 (since $\phi(x,y,z)$ is constant)
- Level surfaces in higher dimensions are represented as contours



source: wikipedia

Minimization Problem

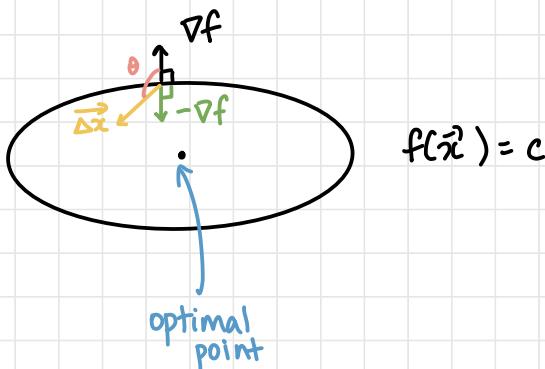
- To minimize value of a function $f(\vec{x})$ subject to the constraints $g(\vec{x}) \leq 0$ (boundary—level surface) called the feasible region
- Start with \vec{x} lying on the constraint boundary $g(\vec{x}) = 0$

Direction Requirements

1. Descent Direction

- If a vector decreases the value of the objective function $f(\vec{x})$, it is said to move in the descent direction

- Test if vector $\vec{\Delta x}$ is moving in descent direction
 - ↪ Find $\nabla f(\vec{x})$
 - ↪ Find $\nabla f(\vec{x}') \cdot \vec{\Delta x}$
- If the dot product is < 0 , the vector is moving in descent direction
- Reason: if $\vec{a} \cdot \vec{b}$ is negative, θ is obtuse. This means angle between ∇f and $\vec{\Delta x}$ is obtuse and $\vec{\Delta x}$ reduces the value of $f(\vec{x}')$



Q: $f(x) = 2x_1 + 3x_1 x_2^2$

Tell if the following vectors are in descent direction starting from $(4, -2)$

$$\nabla(x) = [2 + 3x_2^2 \quad 6x_1 x_2] \quad \text{at } (4, -2), \quad [14 \quad -48]$$

(a) $[2 \quad 1]$

$$[2 \quad 1] \cdot [14 \quad -48] = 28 - 48 = -20 \quad : \text{descent}$$

(b) $[2 \quad -1]$

$$[2 \quad -1] \cdot [14 \quad -48] = 28 + 48 = 76 \quad : \text{not descent}$$

(c) $[-2 \quad -1]$

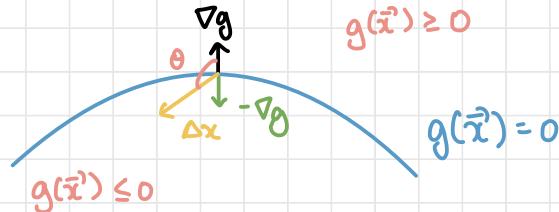
$$[-2 \quad -1] \cdot [14 \quad -48] = -28 + 48 = 20 \quad : \text{not descent}$$

(d) $[-2 \quad 1]$

$$[-2 \quad 1] \cdot [14 \quad -48] = -28 - 48 = -76 \quad : \text{descent}$$

2. Feasible direction

- Vector that will not violate the constraints is said to be moving in a feasible direction
- If the constraint is $g(\vec{x}) \leq 0$



- To test if a test vector $\vec{\Delta x}$ is moving in a feasible direction
 - ↪ Find $\nabla g(\vec{x})$
 - ↪ Find $\nabla g(\vec{x}') \cdot \vec{\Delta x}$
- If the dot product is ≤ 0 , the vector is moving in feasible direction

Q: Given the constraint $x_1 + x_2^2 - 4 \leq 0$

(a) Choose any point on the boundary of this constraint

$$(0, 2)$$

(b) Test the vectors $(3, 4)$, $(-5, 4)$, $(-3, -2)$ for feasibility

$$\nabla g = [1 \quad 2x_2] \text{ at } (0, 2) = [1 \quad 4]$$

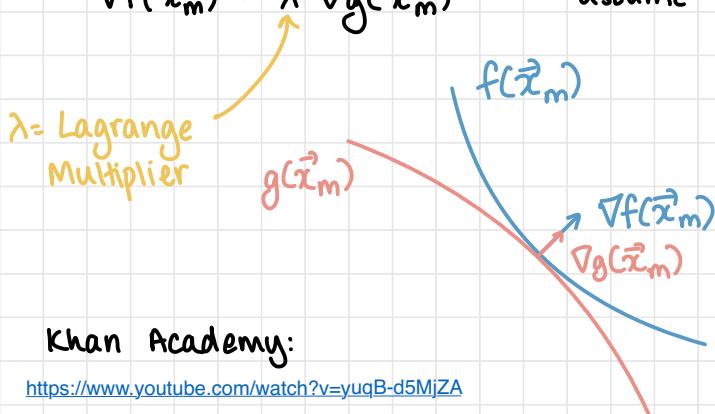
1. $(3, 4) \cdot (1, 4) = 3+16=19$ — infeasible
2. $(-5, 4) \cdot (1, 4) = -5+16=11$ — infeasible
3. $(-3, -2) \cdot (1, 4) = -3-8 = -11$ — feasible

Descent & Feasible Cones

- Overlap between set of all descent vectors and set of all feasible vectors is called set of productive vectors
- Overlap gets smaller as minimum point is approached
- Once there is no overlap, optimal point reached

Non-Linear Optimisation

- First step is to start on the boundary
- Let the optimal point be represented by vector \vec{x}_m (\vec{x}^*)
- $\nabla f(\vec{x}_m) = -\lambda \nabla g(\vec{x}_m)$ — assume $g(\vec{x})$ is an equality constraint



Lagrange Multiplier

- Define single Lagrangian function $L(X, \lambda) = f(X) + \lambda g(X)$
↳ also X^*
- To find optimal point / stationary point X_m of L wrt both X and λ , we need to solve the constrained optimization problem
- If X is d-dimensional,
 $\nabla_X L = 0 \longrightarrow d$ equations
- $\frac{\partial L}{\partial \lambda} = 0 \longrightarrow 1$ equation
- Totally: $d+1$ equations, $d+1$ unknowns

Q: Maximize $f(X) = 1 - x_1^2 - x_2^2$ subject to $g(X) = x_1 + x_2 - 1 = 0$

$$L(X, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1) = 0$$

$$\nabla_X L = \begin{pmatrix} -2x_1 + \lambda \\ -2x_2 + \lambda \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\frac{\partial L}{\partial \lambda} = x_1 + x_2 - 1 = 0$$

Solving (using calculator)

$$x_1 = \frac{1}{2} \quad x_2 = \frac{1}{2} \quad \lambda = 1$$

Optimal point X_m or $X^* = (\frac{1}{2}, \frac{1}{2})$

Multiple Constraints

- One multiplier per constraint

Karush-Kuhn-Tucker (KKT) Conditions

- Inequality: same d+1 equations plus additional constraints

$$(1) g(x) \geq 0$$

$$(2) \lambda \geq 0$$

$$(3) \lambda g(x) = 0$$

Support Vector Machines

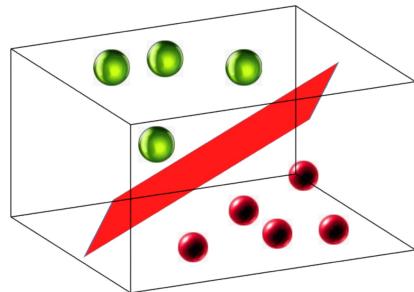
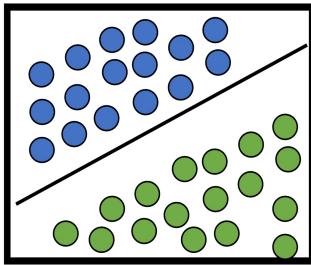
- KNN: classified based on neighbours
- SVM: divide space based on training examples and find where query point lies

Properties

1. Non linearity can be handled
2. Multi classification
3. Classification & regression
4. Computationally expensive

Linearly separable Data

- Hyperplane can separate data

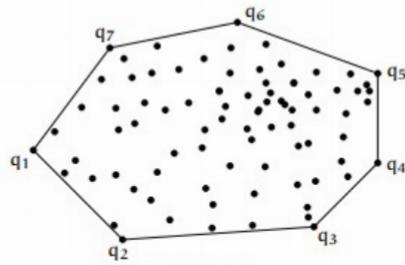


- Convex hull problem: CLRS

The convex hull of a set of points is defined as the smallest convex polygon that encloses all of the points in the set. Convex means that the polygon has no corner that is bent inwards.



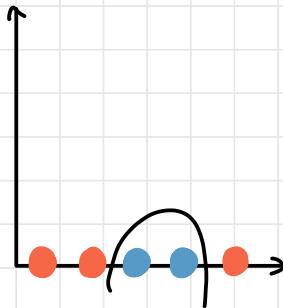
(a) Input.



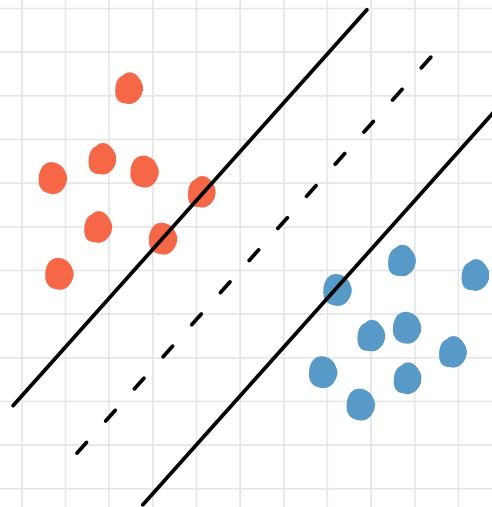
(b) Output.

<https://medium.com/@harshitsikchi/convex-hulls-explained-baab662c4e94>

Transforming Data to Linear separability



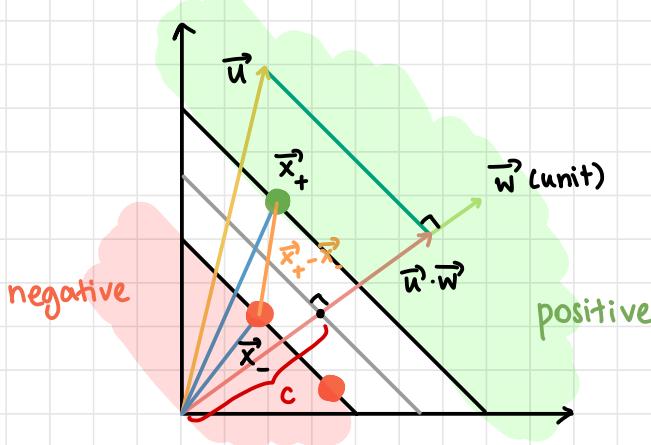
- Kernel trick : transform $x \rightarrow f(x) = x^q$
- Vapnik's idea: find widest street/channel/gutter to separate pos & neg classes



- Hyper plane with widest margin
- Watch MIT SVM: https://www.youtube.com/watch?v=_PwhiWxHK8o

Binary Classification

- Assume decision boundary & gutter found
- Let sample to be classified be represented as \vec{v} (unknown) and a unit vector \perp to the gutter be \vec{w}
- Rule for classification:
 - ↪ if projection of \vec{v} onto \vec{w} falls on positive side of street, \vec{v} is positive
 - ↪ if projection of \vec{v} onto \vec{w} falls on negative side of street, \vec{v} is negative



- In other words,

$$\begin{aligned}\vec{w} \cdot \vec{u} &\geq c \quad \text{pos} \\ \vec{w} \cdot \vec{u} &\leq c \quad \text{neg}\end{aligned}$$

$$\begin{aligned}\vec{w} \cdot \vec{u} - c &\geq 0 \quad \rightarrow \textcircled{1} \\ \vec{w} \cdot \vec{u} - c &\leq 0 \quad \rightarrow \textcircled{2}\end{aligned}$$

Let $\vec{w} = \vec{x}_+$ for +ve sample and $\vec{w} = \vec{x}_-$ for negative sample

At the boundary, let the two constants be c_1 & c_2

$$\vec{w} \cdot \vec{x}_+ - c_1 = 0 \rightarrow ③ \quad \vec{w} \cdot \vec{x}_- - c_2 = 0 \rightarrow ④$$

Let $A = -\frac{(c_2 + c_1)}{2}$ and $B = -\frac{(c_2 - c_1)}{2}$

Then $-c_1 = A - B$ and $-c_2 = A + B$

Replacing in ③ & ④

$$\vec{w} \cdot \vec{x}_+ + A - B = 0 \rightarrow ⑤ \quad \vec{w} \cdot \vec{x}_- + A + B = 0 \rightarrow ⑥$$

Dividing ⑤ & ⑥ by B

$$\frac{\vec{w}}{B} \cdot \vec{x}_+ + \frac{A}{B} - 1 = 0 \rightarrow ⑦ \quad \frac{\vec{w}}{B} \cdot \vec{x}_- + \frac{A}{B} + 1 = 0 \rightarrow ⑧$$

Redefining $\frac{\vec{w}}{B} = \vec{w}$ and $\frac{A}{B} = b$ in ⑦ & ⑧

$$\vec{w} \cdot \vec{x}_+ + b = 1 \longrightarrow ⑨$$

$$\vec{w} \cdot \vec{x}_- + b = -1 \longrightarrow ⑩$$

Converting to Inequality

$$\vec{w} \cdot \vec{x}_+ + b \geq 1 \longrightarrow ⑪$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1 \longrightarrow ⑫$$

The predicted values of y

$$y_i = 1 : +ve$$

$$y_i = -1 : -ve$$

Multiplying ⑪ & ⑫ with y_i

$$y_i (\vec{w} \cdot \vec{x}_+ + b) - 1 \geq 0 \rightarrow \textcircled{13}$$

$$y_i (\vec{w} \cdot \vec{x}_- + b) - 1 \geq 0 \rightarrow \textcircled{14}$$

Generalizing for all \vec{x} 's

$$y_i (\vec{w} \cdot \vec{x} + b) - 1 \geq 0 \rightarrow \textcircled{15}$$

Solving for \vec{x}_+ & $-\vec{x}_-$ from ⑬ & ⑭

$$\vec{x}_+ \geq (\vec{w})^{-1} (1-b) \rightarrow \textcircled{15}$$

$$-\vec{x}_- \geq (\vec{w})^{-1} (1+b) \rightarrow \textcircled{16}$$

The width of the gutter is the projection of $\vec{x}_+ - \vec{x}_-$ onto \vec{w} , if \vec{x}_+ and \vec{x}_- lie on the boundaries

$$\vec{x}_+ - \vec{x}_- = \textcircled{15} + \textcircled{16} = (\vec{w})^{-1} (2) \rightarrow \textcircled{17}$$

Width from ⑯

$$\text{width} = \frac{2(\vec{w})^{-1} \vec{w}}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$$

The goal is to maximize width, or minimize $\|\vec{w}\|$

Primal Problem

$$\text{maximize } f(w, b) = \frac{1}{2} \|w\|^2$$

subject to constraints

$$y_i(w^T x_i + b) - 1 \geq 0 \quad (\text{eq 15})$$

for $i = 1 \text{ to } m$ where $m = \text{training instances}$

Optimization

$$\text{maximize } \frac{1}{2} \|w\|^2 = \text{minimize } \frac{\|w\|^2}{2} = \text{minimize } \frac{\|w\|^2}{2}$$

Using Lagrangian Multiplier Method

Multiple constraints

$$L(x, \lambda) = f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_k g_k(x)$$

$$\text{Objective } f(w) = \frac{1}{2} \|w\|^2 \text{ to minimize}$$

$$\text{Subject to } m \text{ constraints } g(w, b) = y_i(w^T x_i + b) - 1$$

Lagrangian Function

$$L = f(w) - \sum_{i=1}^m \alpha_i g_i(w, b) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1] \rightarrow ①$$

$i = 1 \text{ to } m$

Lagrangian

Derivative wrt w

why vector? check pg 48 of L4-SVM
in Prof Preet's slides

$$\frac{\partial L}{\partial w} = \vec{w} - \sum \alpha_i y_i \vec{x}_i = 0$$

$$\vec{w} = \sum_{i=1}^m \alpha_i y_i \vec{x}_i \quad \rightarrow (2)$$

Derivative wrt b

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^m \alpha_i y_i = 0$$

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad \rightarrow (3)$$

Replacing ② & ③ in ①

$$L = \frac{1}{2} \left(\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j \right) - \left(\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j \right) - b \underbrace{\sum_{i=1}^m \alpha_i y_i}_{0} + \sum_{i=1}^m \alpha_i$$

$$L(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \left(\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j \right)$$

SVM-Dual Optimization Problem

watch: <https://www.youtube.com/watch?v=6-ntMlaJpm0>

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \left(\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \right)$$

ST constraints ($\forall i = 1 \text{ to } m$)

$$c_1: \alpha_i \geq 0$$

$$c_2: \sum_{i=1}^m \alpha_i y_i = 0$$

$$c_3: \sum_{i=1}^m \alpha_i y_i x_i x + b = 1$$

$$c_4: \sum_{i=1}^m \alpha_i y_i x_i x + b = -1$$

- Vectors that lie on boundary: support vectors ($\alpha > 0$)

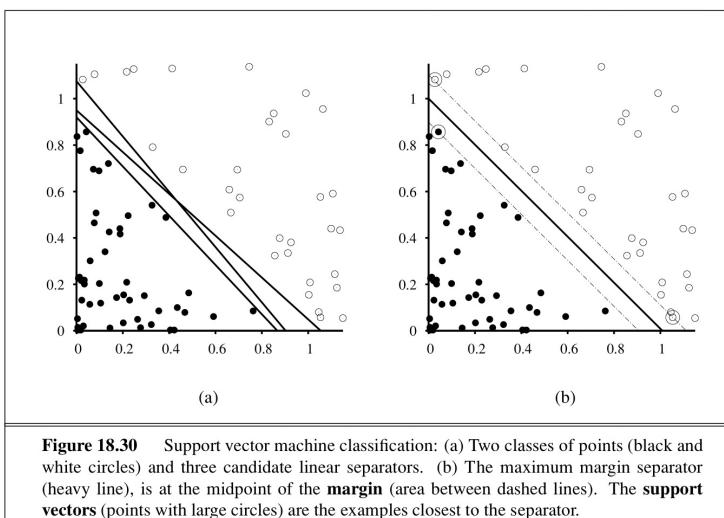


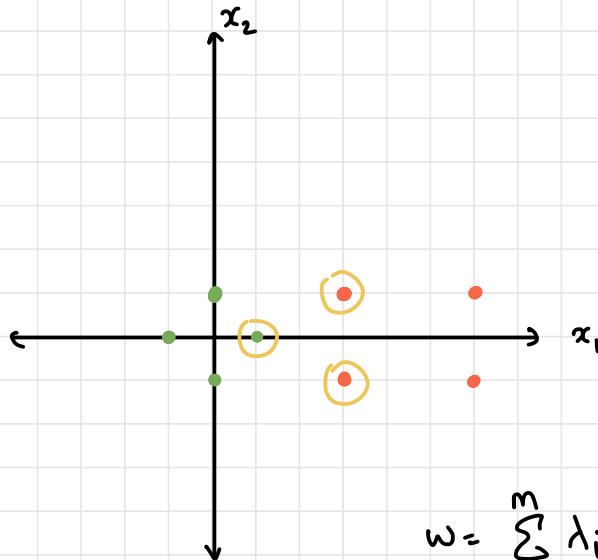
Figure 18.30 Support vector machine classification: (a) Two classes of points (black and white circles) and three candidate linear separators. (b) The maximum margin separator (heavy line), is at the midpoint of the margin (area between dashed lines). The support vectors (points with large circles) are the examples closest to the separator.

$$Q: \begin{array}{c} (3,1), (3,-1), (6,1), (6,-1) \\ +ve \end{array} \quad \left| \quad \begin{array}{c} (1,0), (0,1), (0,-1), (-1,0) \\ -ve \end{array} \right.$$

$$S_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ -ve}$$

$$S_2 = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \text{ +ve}$$

$$S_3 = \begin{bmatrix} 3 \\ -1 \end{bmatrix} \text{ +ve}$$



$$w = \sum_{i=1}^m \lambda_i y_i x_i$$

3 Lagrangian multipliers

$$c_1: \alpha_i \geq 0$$

$$c_2: \sum_{i=1}^3 \lambda_i y_i = 0 = -\lambda_1 + \lambda_2 + \lambda_3 = 0 \longrightarrow ①$$

$$c_3: \sum_{i=1}^3 \lambda_i y_i \vec{x}_i \cdot \vec{x}_1 + b = -1$$

$$c_4: \sum_{i=1}^3 \lambda_i y_i \vec{x}_i \cdot \vec{x}_2 + b = +1$$

$$c_5: \sum_{i=1}^3 \lambda_i y_i \vec{x}_i \cdot \vec{x}_3 + b = +1$$

$$c_3: \lambda_1(-1)[1 \ 0] \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \lambda_2(1)[3 \ -1] \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \lambda_3(1)[3 \ -1] \begin{bmatrix} 1 \\ 0 \end{bmatrix} + b = -1$$

$$-\lambda_1 + 3\lambda_2 + 3\lambda_3 + b = -1 \longrightarrow \textcircled{2}$$

$$c_4: \lambda_1(-1)[1 \ 0] \begin{bmatrix} 3 \\ 1 \end{bmatrix} + \lambda_2(1)[3 \ -1] \begin{bmatrix} 3 \\ 1 \end{bmatrix} + \lambda_3(1)[3 \ -1] \begin{bmatrix} 3 \\ 1 \end{bmatrix} + b = +1$$

$$-3\lambda_1 + 10\lambda_2 + 8\lambda_3 + b = 1 \longrightarrow \textcircled{3}$$

$$c_5: \lambda_1(-1)[1 \ 0] \begin{bmatrix} 3 \\ -1 \end{bmatrix} + \lambda_2(1)[3 \ -1] \begin{bmatrix} 3 \\ -1 \end{bmatrix} + \lambda_3(1)[3 \ -1] \begin{bmatrix} 3 \\ -1 \end{bmatrix} + b = +1$$

$$-3\lambda_1 + 8\lambda_2 + 10\lambda_3 + b = 1 \longrightarrow \textcircled{4}$$

$$\lambda_1 = 1/2 \quad \lambda_2 = 1/4 \quad \lambda_3 = 1/4 \quad b = -2$$

$$w = \sum_{i=1}^3 \lambda_i x_i y_i = \frac{1}{2} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} (-1) + \frac{1}{4} \begin{bmatrix} 3 \\ 1 \end{bmatrix} (1) + \frac{1}{4} \begin{bmatrix} 3 \\ -1 \end{bmatrix} (1)$$

$$= \begin{bmatrix} -1/2 \\ 0 \end{bmatrix} + \begin{bmatrix} 3/4 \\ 1/4 \end{bmatrix} + \begin{bmatrix} 3/4 \\ -1/4 \end{bmatrix}$$

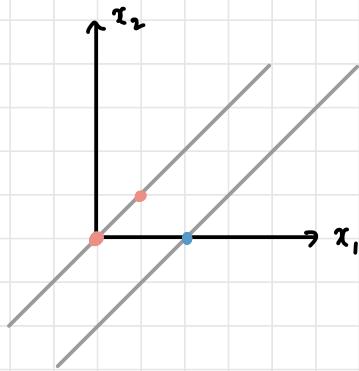
$$= \begin{bmatrix} 1 \\ 0 \end{bmatrix} \longrightarrow w x + b = 0 \text{ is line}$$

$$(1)x_1 + 0(x_2) - 2 = 0 \Rightarrow x_1 = 2$$

Q: The samples are $(0,0)$, $(1,1)$, $(2,0)$ (SVs)
 -ve -ve +ve

Test sample: $(6,0)$

How many Lagrangian Multipliers? 3 (1 for each support vect)
 Find weight & bias.



$$c_1 = \sum \lambda_i y_i = 0$$

$$-\lambda_1 - \lambda_2 + \lambda_3 = 0 \longrightarrow \textcircled{1}$$

$$\begin{aligned} c_2 = w_0 x + b &= +1 && \text{pos} \\ w_0 x + b &= -1 && \text{neg} \end{aligned}$$

$$w = \sum_{i=1}^m y_i x_i \lambda_i - x_j + b = -1$$

$$-\lambda_1 \begin{pmatrix} 0 \\ 0 \end{pmatrix} [0 \ 0] - \lambda_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} [0 \ 0] + \lambda_3 \begin{pmatrix} 2 \\ 0 \end{pmatrix} [0 \ 0] + b = -1$$

$$b = -1 \longrightarrow \textcircled{2}$$

$$-\lambda_1 \begin{pmatrix} 0 \\ 0 \end{pmatrix} [1 \ 1] - \lambda_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} [1 \ 1] + \lambda_3 \begin{pmatrix} 2 \\ 0 \end{pmatrix} [1 \ 1] + b = -1$$

$$-2\lambda_2 + 2\lambda_3 + b = -1 \quad \rightarrow \textcircled{3}$$

$$w = \sum_{i=1}^m y_i x_i \lambda_i x_j + b = +1$$

$$-\lambda_1 \begin{pmatrix} 0 \\ 0 \end{pmatrix} [2 \ 0] - \lambda_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} [2 \ 0] + \lambda_3 \begin{pmatrix} 2 \\ 0 \end{pmatrix} [2 \ 0] + b = 1$$

$$-2\lambda_2 + 4\lambda_3 + b = 1 \quad \rightarrow \textcircled{4}$$

Solving ①, ②, ③ & ④

$$\lambda_3 = \lambda_1 + \lambda_2 \quad b = -1$$

③:

$$-2\lambda_2 + 2\lambda_1 + 2\lambda_2 - 1 = -1$$

$$2\lambda_1 = 0 \Rightarrow \lambda_1 = 0$$

④.

$$-2\lambda_2 + 4\lambda_2 - 1 = 1$$

$$2\lambda_2 = 2$$

$$\lambda_2 = 1$$

$$\lambda_1 = 0 \quad \lambda_2 = 1 \quad \lambda_3 = 1 \quad b = -1$$

$$\begin{aligned} w &= \sum x_i y_i \lambda_i = \binom{0}{0}(-1)(0) + \binom{1}{1}(-1)(1) + \binom{2}{0}(1)(1) \\ &= \binom{-1}{-1} + \binom{2}{0} = \binom{1}{-1} \end{aligned}$$

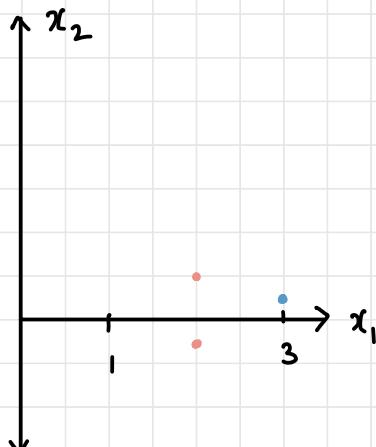
$$w = [1 \quad -1] \longrightarrow \perp \text{ to line}$$

$$x_1 - x_2 - 1 = 0 \Rightarrow x_2 = x_1 - 1$$

If $wx_i + b > 1$, +ve sample, not SV

Q: Samples: $(3, \frac{1}{4})$, $(2, -\frac{1}{4})$, $(2, \frac{1}{2})$. Find out w & b .

-ve +ve +ve



$$c_1 \quad \sum \lambda_i y_i = 0$$

$$\begin{pmatrix} 3 \\ 1/4 \end{pmatrix} - \begin{pmatrix} 2 \\ -1/4 \end{pmatrix} + \begin{pmatrix} 2 \\ 1/2 \end{pmatrix}$$

$$-\lambda_1 + \lambda_2 + \lambda_3 = 0 \quad \rightarrow ①$$

$$c_2 \quad w x + b = +1 \quad (+ve) \\ w x + b = -1 \quad (-ve)$$

$$w = \sum_{i=1}^m \lambda_i x_i y_i$$

x_1 :

$$\left(\lambda_1 \begin{pmatrix} 3 \\ 1/4 \end{pmatrix} (-1) + \lambda_2 \begin{pmatrix} 2 \\ -1/4 \end{pmatrix} (1) + \lambda_3 \begin{pmatrix} 2 \\ 1/2 \end{pmatrix} (1) \right) [3 \ 1/4] + b = -1$$

$$-\frac{145}{16} \lambda_1 + \frac{95}{16} \lambda_2 + \frac{49}{8} \lambda_3 + b = -1 \quad \rightarrow ②$$

x_2 :

$$\left(\lambda_1 \begin{pmatrix} 3 \\ 1/4 \end{pmatrix} (-1) + \lambda_2 \begin{pmatrix} 2 \\ -1/4 \end{pmatrix} (1) + \lambda_3 \begin{pmatrix} 2 \\ 1/2 \end{pmatrix} (1) \right) [2 \ -1/4] + b = +1$$

$$-\frac{95}{16} \lambda_1 + \frac{65}{16} \lambda_2 + \frac{31}{8} \lambda_3 + b = 1 \quad \rightarrow ③$$

x_3 :

$$\left(\lambda_1 \begin{pmatrix} 3 \\ 1/4 \end{pmatrix} (-1) + \lambda_2 \begin{pmatrix} 2 \\ -1/4 \end{pmatrix} (1) + \lambda_3 \begin{pmatrix} 2 \\ 1/2 \end{pmatrix} (1) \right) [2 \ 1/2] + b = +1$$

$$-\frac{49}{8} \lambda_1 + \frac{31}{8} \lambda_2 + \frac{17}{4} \lambda_3 + b = 1 \quad \rightarrow ④$$

$$\lambda_1 = 2 \quad \lambda_2 = \frac{2}{3} \quad \lambda_3 = \frac{4}{3} \quad b = 5$$

$$\omega = \sum_{i=1}^m \lambda_i x_i y_i = 2 \begin{pmatrix} 3 \\ -1/4 \end{pmatrix} (-1) + \frac{2}{3} \begin{pmatrix} 2 \\ -1/4 \end{pmatrix} (1) + \frac{4}{3} \begin{pmatrix} 2 \\ 1/2 \end{pmatrix}$$

$$\omega = \begin{pmatrix} -2 \\ 0 \end{pmatrix} = \perp \text{ to line}$$

$$b = 5$$

$$-2x_1 + 5 = 0$$

$$x_1 = \frac{5}{2}$$

NON-LINEARLY SEPARABLE DATA

- Transform data to be linearly separable (kernel)
 - can do so by mapping points to higher dim space
- Let $x = (x_1, x_2, x_3)^T$, $y = (y_1, y_2, y_3)^T$
- Assume we need to map to 9D space with function ϕ

$$\phi(x) = (x_1^2, x_1x_2, x_1x_3, x_2x_1, x_2^2, x_2x_3, x_3x_1, x_3x_2, x_3^2)^T$$

$$\phi(y) = (y_1^2, y_1y_2, y_1y_3, y_2y_1, y_2^2, y_2y_3, y_3y_1, y_3y_2, y_3^2)^T$$

- Final result required is dot product of 2 support vectors

$$\cdot \phi(x)^T \phi(y) = \sum_{i=1}^3 \sum_{j=1}^3 x_i x_j y_i y_j$$

- Computational complexity: $O(n^2)$ — here $n=3$
- Can use a **kernel trick** to avoid unnecessary computations
- Using kernel function $k(x,y) = (x^T y)^2$

$$k(x,y) = (x_1 y_1 + x_2 y_2 + x_3 y_3)^2$$

$$k(x,y) = \sum_{i=1}^3 \sum_{j=1}^3 x_i x_j y_i y_j$$

- Rewrite dual problem as

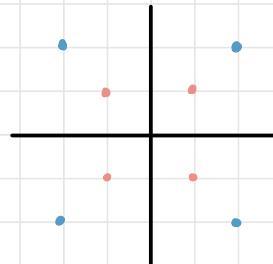
$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

Kernel Functions

- Polynomial : $k(x_1, x_2) = (x_1 \cdot x_2 + c)^d$
- Gaussian / RBF

$$Q: \left(\begin{matrix} 2 \\ 2 \\ 2 \end{matrix} \right), \left(\begin{matrix} 2 \\ -2 \\ 2 \end{matrix} \right), \left(\begin{matrix} -2 \\ 2 \\ 2 \end{matrix} \right), \left(\begin{matrix} -2 \\ -2 \\ 2 \end{matrix} \right), \left(\begin{matrix} 1 \\ 1 \\ 1 \end{matrix} \right), \left(\begin{matrix} 1 \\ -1 \\ 1 \end{matrix} \right), \left(\begin{matrix} -1 \\ 1 \\ 1 \end{matrix} \right), \left(\begin{matrix} -1 \\ -1 \\ 1 \end{matrix} \right)$$

$\underbrace{\hspace{10em}}$ +ve $\underbrace{\hspace{10em}}$ -ve



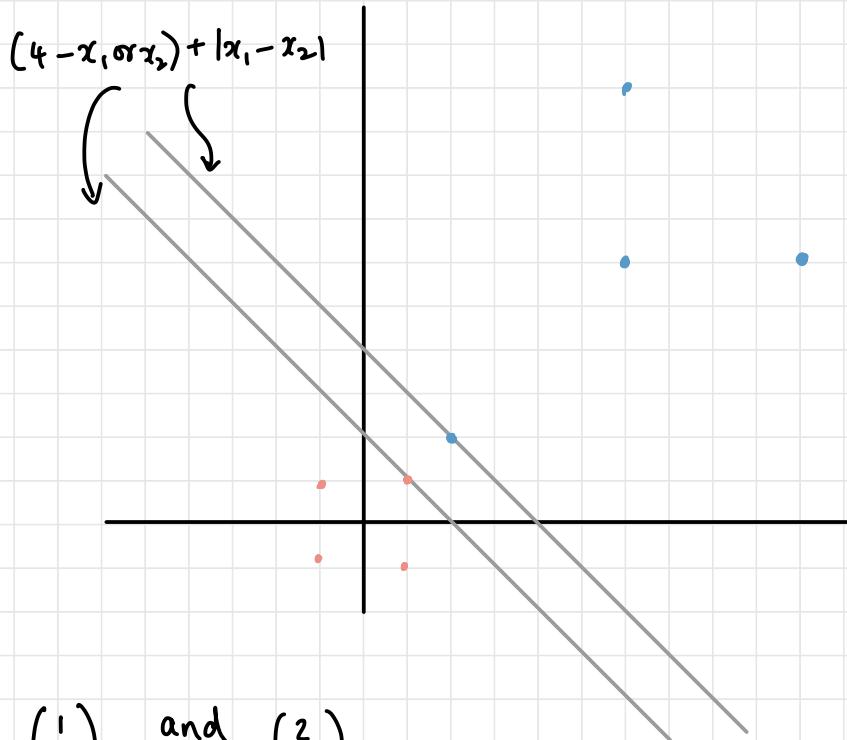
(kernel)

$$\phi \left(\begin{matrix} x_1 \\ x_2 \end{matrix} \right) = \begin{cases} \left(\begin{matrix} (4-x_2) + |x_1 - x_2| \\ (4-x_1) + |x_1 - x_2| \end{matrix} \right) & \text{if } \sqrt{x_1^2 + x_2^2} \geq 2 \\ \left(\begin{matrix} x_1 \\ x_2 \end{matrix} \right) & \text{otherwise} \end{cases}$$

New points:

$$\left(\begin{matrix} 2 \\ 2 \\ 2 \end{matrix} \right), \left(\begin{matrix} 10 \\ 6 \\ 6 \end{matrix} \right), \left(\begin{matrix} 6 \\ 10 \\ 6 \end{matrix} \right), \left(\begin{matrix} 6 \\ 6 \\ 6 \end{matrix} \right), \left(\begin{matrix} 1 \\ 1 \\ 1 \end{matrix} \right), \left(\begin{matrix} 1 \\ -1 \\ 1 \end{matrix} \right), \left(\begin{matrix} -1 \\ 1 \\ 1 \end{matrix} \right), \left(\begin{matrix} -1 \\ -1 \\ 1 \end{matrix} \right)$$

$\underbrace{\hspace{10em}}$ +ve $\underbrace{\hspace{10em}}$ -ve



SVs: $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ and $\begin{pmatrix} 2 \\ 2 \end{pmatrix}$
 -ve +ve

$$c_1 \cdot \sum_{i=1}^2 \lambda_i y_i = 0$$

$$c_2: \sum_{i=1}^2 \lambda_i y_i \vec{x}_i \cdot \vec{x}_1 + b = -1$$

$$c_3 \sum_{i=1}^2 \lambda_i y_i \vec{x}_i \cdot x_2 + b = +1$$

$$c_1: -\lambda_1 + \lambda_2 = 0 \longrightarrow ①$$

$$c_2: (\lambda_1(-1)\begin{pmatrix} 1 \\ 1 \end{pmatrix} + \lambda_2(1)\begin{pmatrix} 2 \\ 2 \end{pmatrix}) \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} + b = -1$$

$$-2\lambda_1 + 4\lambda_2 + b = -1 \rightarrow ②$$

$$c_3: (\lambda_1(-1)\begin{pmatrix} 1 \\ 1 \end{pmatrix} + \lambda_2(1)\begin{pmatrix} 2 \\ 2 \end{pmatrix}) \cdot \begin{pmatrix} 2 \\ 2 \end{pmatrix} + b = 1$$

$$-4\lambda_1 + 8\lambda_2 + b = 1 \rightarrow ③$$

Q:

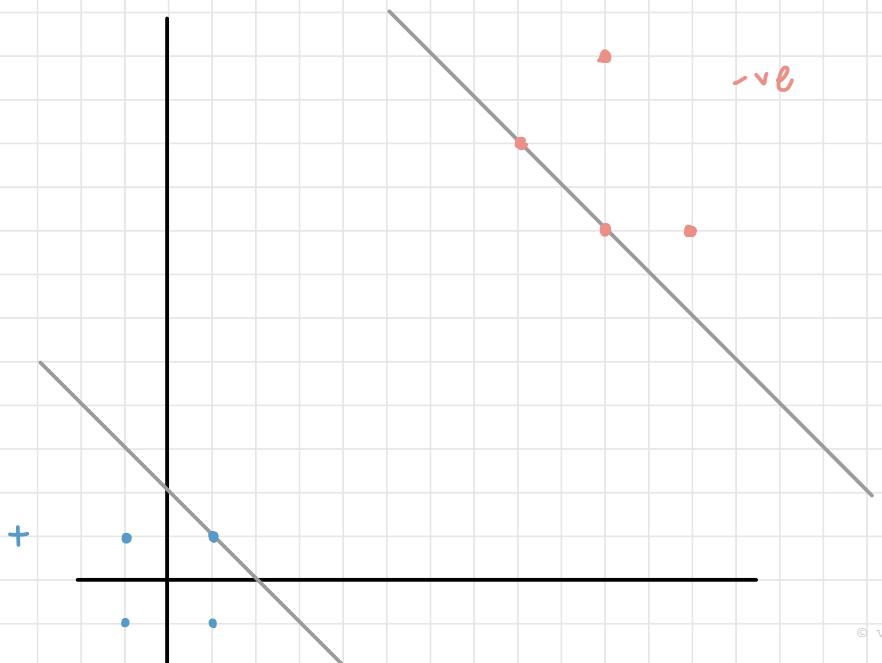
$$\phi(x_1, x_2) = \begin{cases} \left(\frac{6-x_1 + (x_1 - x_2)^2}{6-x_2 + (x_1 - x_2)^2}, \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) & \text{if } \sqrt{x_1^2 + x_2^2} \geq 2 \\ \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} & \text{otherwise} \end{cases}$$

$(1, 1), (-1, 1), (-1, -1), (1, -1)$, $\underbrace{(2, 0), (0, 2), (-2, 0), (0, -2)}$

+ve *-ve*

Transformed points

$$\left(\begin{matrix} 1 \\ 1 \end{matrix}\right), \left(\begin{matrix} -1 \\ 1 \end{matrix}\right), \left(\begin{matrix} -1 \\ -1 \end{matrix}\right), \left(\begin{matrix} 1 \\ -1 \end{matrix}\right), \left(\begin{matrix} 8 \\ 10 \end{matrix}\right), \left(\begin{matrix} 10 \\ 8 \end{matrix}\right), \left(\begin{matrix} 12 \\ 10 \end{matrix}\right), \left(\begin{matrix} 10 \\ 12 \end{matrix}\right)$$

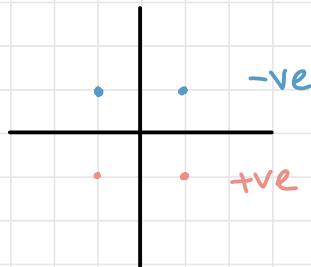


Q: Construct an SVM that computes XOR function. Use values +1 and -1 for both input & output. Map the input (x_1, x_2) into a space consisting of (x_1, x_2, x_1x_2) .

Draw the four input points in this space and the max margin separator.

What is the margin? Infer if the datapoints are separable or not.

			New plane		
x_1	x_2	y	x_1	x_2	y
-1	-1	-1	-1	1	-1
-1	1	1	-1	-1	1
1	-1	1	1	-1	1
1	1	-1	1	1	-1



$$\sum_{i=1}^m \lambda_i w_i = 0$$