

An Introduction to Artificial Neural Network

Ramaseshan Ramachandran

Standard Algorithms
Human/Machine Learning
Classification
Linear Models for Classification
Decision Boundary
Artificial Neural Network

Limitations of Perceptron
Loss Function
Cost Function
Gradient Descent

- ① Logistic Regression
- ② References

An algorithm is a sequence of instructions to solve a problem

- ▶ The steps to solve problems are well defined
- ▶ Steps are coded in some ordered sequence to transform the input from one form to another
- ▶ Rules are unambiguous
- ▶ Sufficient Knowledge is available to fully solve the problem

- ▶ There are problems whose solutions cannot be formulated using standard rule-based algorithms
- ▶ Problems that require subtle inputs cannot be solved using standard algorithmic approach - face recognition, speech recognition, hand-written character recognition, etc
- ▶ Finding Examples and using experience gained in similar situations are useful
- ▶ Examples provide certain underlying patterns
- ▶ Patterns give the ability to predict some outcome or help in constructing an approximate model
- ▶ **Learning** is the key to the ambiguous world

- ▶ Quantitative expression of the real world observations
- ▶ Mathematical Models translate the real world observations into mathematical expressions
- ▶ A machine learning model is a mathematical construct trained to recognize patterns in data
- ▶ It makes predictions or classifications based on those patterns
- ▶ Helps in understanding and interpreting information for decision making/classification
- ▶ Data driven decisions under uncertainty
- ▶ Model Parameters store vast amount of information captured from the data
 - ▶ Lexical, linguistic and semantic knowledge in the case of NLP
- ▶ Access the *knowledge* (learned from the data) of the Model by using the context

DISCRIMINATIVE MODEL

- ▶ A class of models that focusses on learning the decision boundary between different classes or outcomes
- ▶ It models the conditional probability distribution of the target variable given the input variables, $p(y|\mathbf{x})$.
- ▶ In the case of classification, this is $p(C_k|\mathbf{x})$
- ▶ OR, Learn a function, *discriminant function*, that maps inputs \mathbf{x} directly into decisions

GENERATIVE MODEL

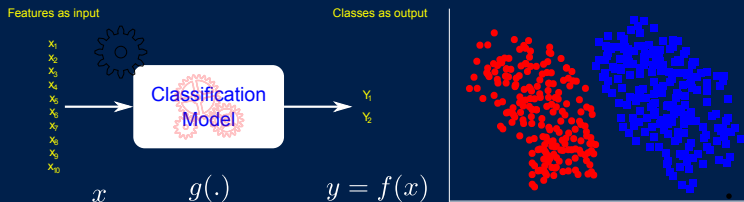
- ▶ If a function models input and output into probability distributions, then it is in general known as a generative model
- ▶ Given a training data, generate new samples similar to the training samples from the same distribution
 - ▶ **Language model** - generating the next word given the context It is possible to generate synthetic data points (or new sentences) using these distributions

Classification is the task of assigning predefined dis-joint categories to objects

- ▶ Detect Spam emails
- ▶ Find the set of mobile phones $< \text{Rs.}10000$ and received 5* reviews
- ▶ Identify the category of the incoming document as sports, politics, entertainment or business
- ▶ Determine whether a movie review is a positive or negative review

DEFINITION OF CLASSIFICATION

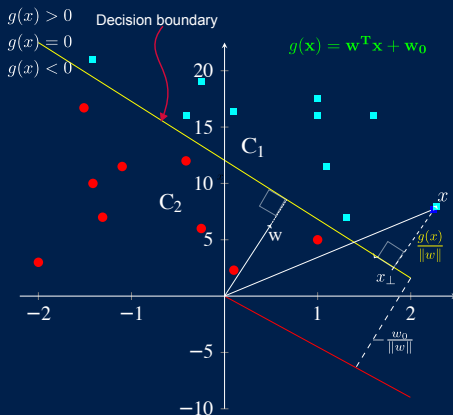
- ▶ The input is a collection of records
- ▶ Each record is represented by a tuple (\mathbf{x}, \mathbf{y})
- ▶ $\mathbf{x} = x_1, x_2, \dots, x_n$ and $\mathbf{y} = y_1, y_2, \dots, y_n$ are the input features and the classes respectively
- ▶ $\mathbf{x} \in \mathbf{R}^2$ is a vector - the set of observed variables
- ▶ (\mathbf{x}, \mathbf{y}) are related by an unknown function. The goal is to estimate the unknown function $g(\cdot)$, also known as a classifier function, such that $g(\mathbf{x}) = \mathbf{f}(\mathbf{x}), \forall \mathbf{x}$



WHAT DOES THE CLASSIFIER FUNCTION DO?

Assuming that we have a linearly separable X , the linear classifier function $g(\cdot)$ implements decision rule

- ▶ Fitting a straight line to a given data set requires two parameters (w_0 and w)
- ▶ The decision rule divides the data space into two sub-spaces - separating two classes using a boundary
- ▶ The distance of the boundary from the origin = $\frac{w_0}{\|w\|}$
- ▶ Distance of any point from the boundary = $d = \frac{g(x)}{\|w\|}$

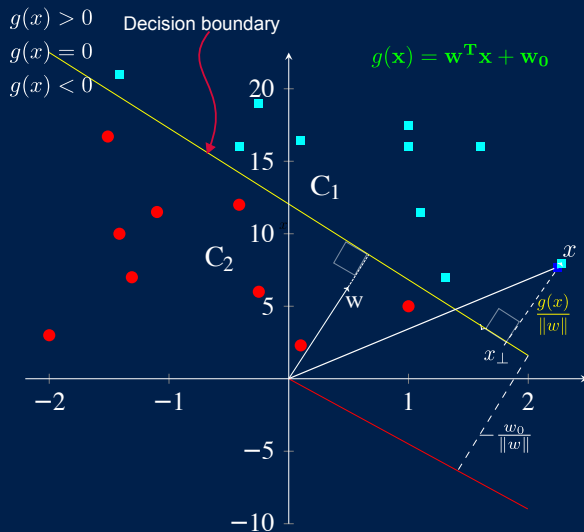


LINEAR MODELS FOR CLASSIFICATION[1]

The goal of classification is to take a vector X and assign it to one of the N discrete classes \mathbb{C}_n , where $n = 1, 2, 3, \dots, N$.

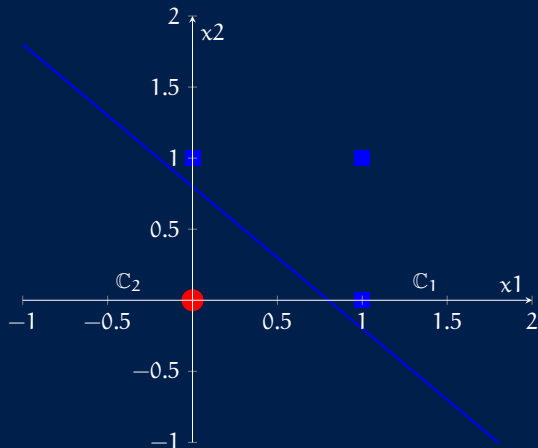
- ▶ The classes are disjoint and an input is assigned to only one class
- ▶ The input space is divided into *decision regions*
- ▶ The boundaries are called as *decision boundaries* or *decision surfaces*
- ▶ In general, if the input space is N dimensional, then $g(x)$ would define an $N - 1$ hyperplane

GEOMETRY OF THE LINEAR DISCRIMINANT FUNCTION



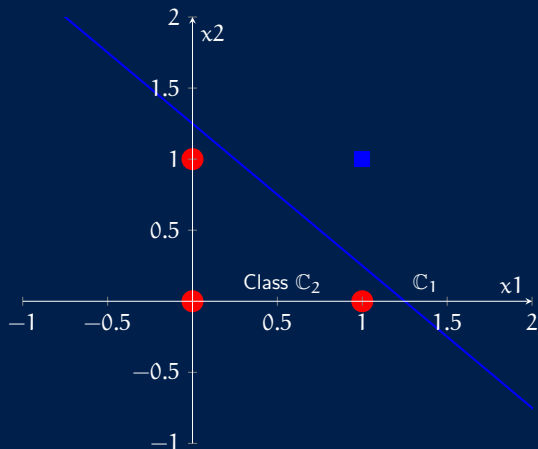
DECISION BOUNDARY FOR OR GATE

The decision regions are separated by a hyperplane and it is defined by $g(x) = 0$. This separates linearly separable classes \mathbb{C}_1 and \mathbb{C}_2



DECISION BOUNDARY FOR AND GATE

The decision regions are separated by a hyperplane and it is defined by $g(x) = 0$. This separates linearly separable classes \mathbb{C}_1 and \mathbb{C}_2

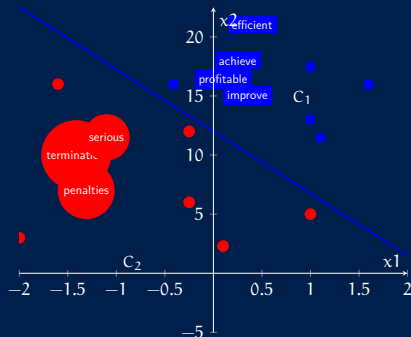


DECISION BOUNDARY FOR SENTIMENTS

Let us consider some positive and negative sentiment terms which are contained in two classes C_P and C_N

$C_P = [\text{achieve efficient improve profitable}] = +1$

$C_N = [\text{termination penalties misconduct serious}] = -1$



How do we build/develop models?

Let $x \in \mathbb{R}^n$ denote the real valued random data points as input

Let $e \in \mathbb{R}^n$ be the real valued random output with a joint distribution of $p(x, y)$

Let us find a $f(x)$ to predict y given x

Let us define a loss function $L(y, f(x))$ for penalizing errors during prediction

Most common loss function is the squared error loss. The expected error prediction is $EPE(f) = E(y - f(x))^2$

A point-wise expected error prediction is a conditional expectation and is written as $f(x) = E(y | x)$ - this is the regression function

Assumption of least squares - $f(x)$ is approximated by a global linear function

LOGISTIC REGRESSION

Supervised Learning

Known relationship (x, y) : x is the data and y is the label

Goal: Learn a the relationship and identify a function to map $x \rightarrow y$ **Model:** An optimized learning algorithm to learn the relationship

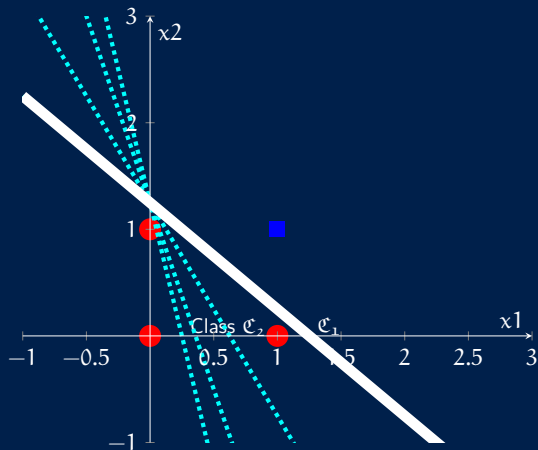
Examples: Classification, regression, senetence generation, object detection, semantic segmentation, image captioning, etc.

Unsupervised Learning

Known relationship: Only data x . Input-output relationship is not known

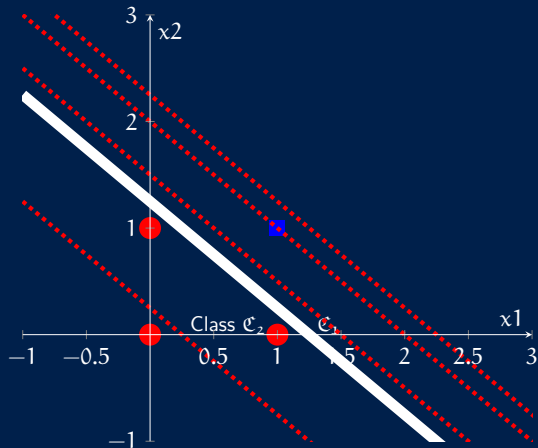
Examples:Clustering, PCA, dimensionality reduction

DECISION BOUNDARY- VARIATION OF w_j

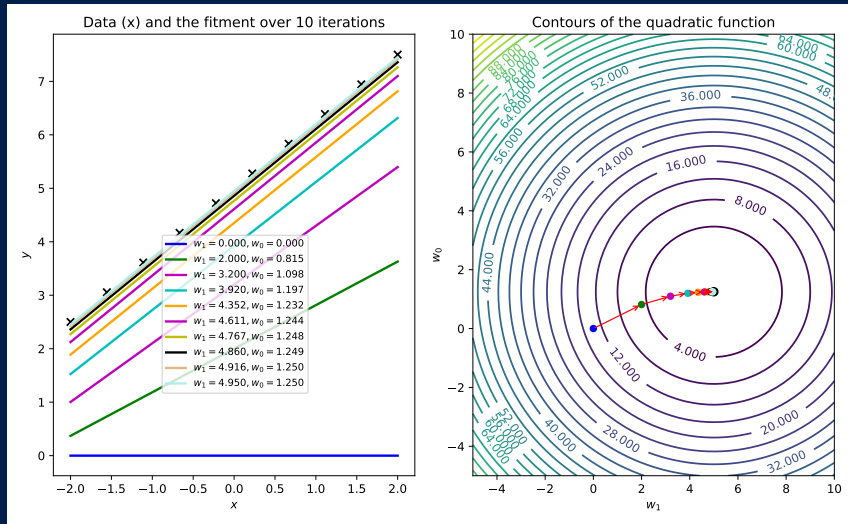


DECISION BOUNDARY - VARIATION OF BIAS

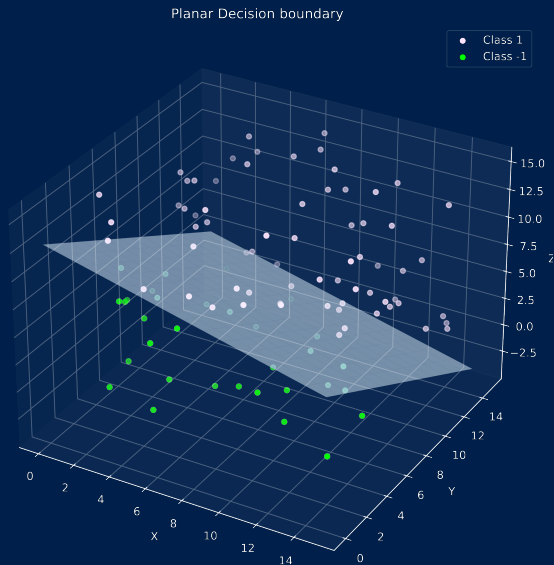
The contribution of bias to the creation of the decision boundary



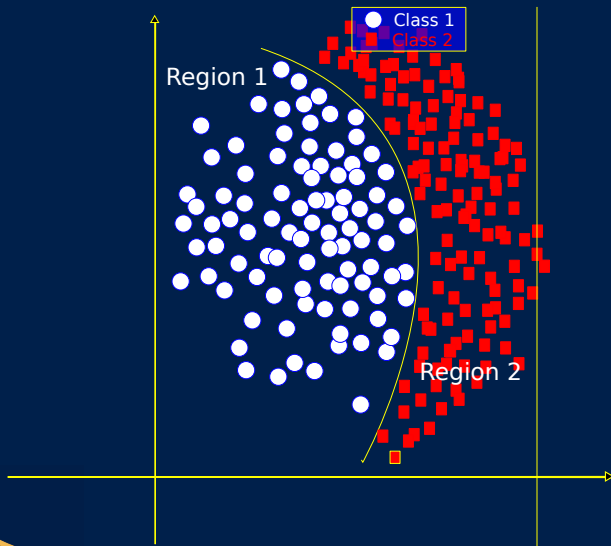
DECISION BOUNDARY AND GRADIENT DESCENT



DECISION SURFACE FOR A 3-D VECTOR



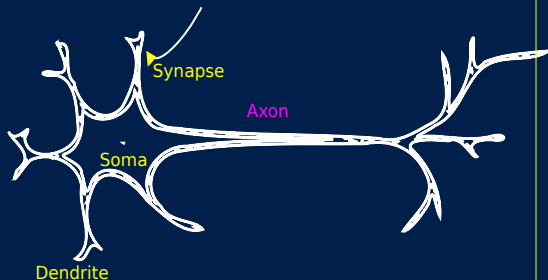
LINEARLY SEPARABLE?



Is this separable?

A	B	C
0	0	0
0	1	1
1	0	1
1	1	0

NEURAL NETWORK



- ▶ Each individual neuron can form thousands of links with other neurons.
- ▶ A typical brain has well over 100 trillion synapses

- ▶ Functionally related neurons connect to each other to form neural networks
- ▶ The electro-chemical connections between neurons are not static
- ▶ The more signals sent between two neurons, the stronger the connection grows and with each new experience and each remembered event, the brain slightly re-wires its physical structure.
- ▶ Our brains form a million new connections for every second of our lives

LAWS OF ASSOCIATION

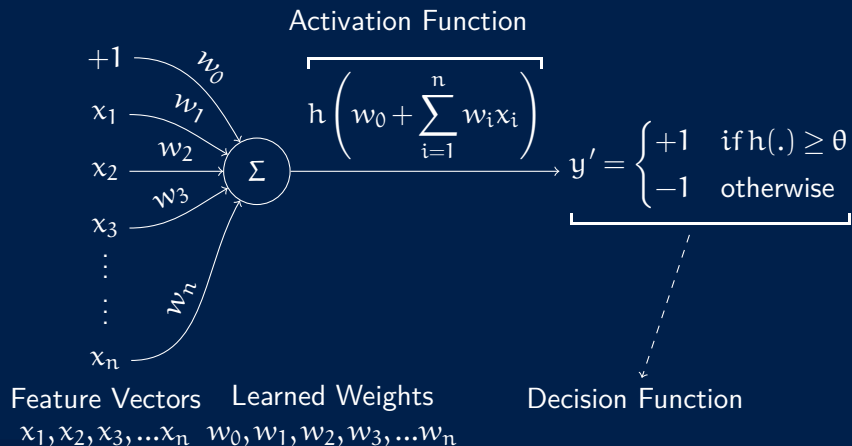
Aristotle's attempts on fundamental laws of learning and memory

The law of similarity	If two things are similar, the thought of one will tend to trigger the thought of the other - word2vec If you recollect one birthday, you may find yourself thinking about others as well
The law of contrast	Seeing or recalling something may also trigger the recollection of something completely opposite
The law of contiguity	Things or events that occur close to each other in space or time tend to get linked together in the mind If you found a snake in the corner of the street, every time you cross the corner, you tend to look for one. Events are conditioned based on the time and space
The law of frequency	The more often two things or events are linked, the more powerful will be that association - think of next word prediction - strength of the association decides who is

PERCEPTRON

Neuron	Perceptron
Biological	A mathematical model of a biological neuron
Dendrites receive electrical signals	Perceptron receives mathematical values as input
Electro-chemical signals between Dendrites and axons	The weighted sum represents the total strength of the signal
The electro-chemical signals are not static	Weights change during the training process

PERCEPTRON



PERCEPTRON LEARNING

- ▶ Perceptron learns the weights
- ▶ They are adjusted until the output is consistent with the target output in the training examples
- ▶ $w^{(k+1)} \propto (y - \hat{y})$
- ▶ The weights are updated as below
$$w_j^{(k+1)} = w_j^{(k)} - \eta(y_i - \hat{y}^{(k)})x_{ij}$$
where $w^{(k)}$ is the weight parameter associated with the i^{th} input at k^{th}

iteration

η is the learning parameter and x_{ij} is the j^{th} attribute of the i^{th} training sample

- ▶ If $(y - \hat{y}) \cong 0$, no prediction error
- ▶ During the training the weights contributing most to the error require adjustments

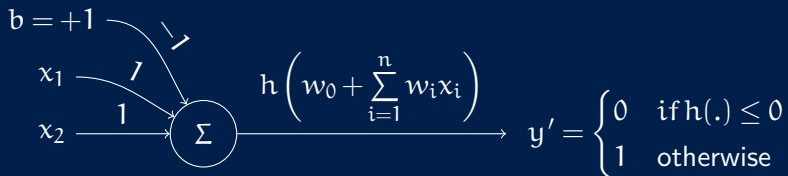
ALGORITHM FOR PERCEPTRON LEARNING

- 1: Total number of input vectors = k
- 2: Total number of features = n
- 3: Learning parameter $\eta = 0.01$, where $0 < \eta < 1$
- 4: epoch¹ count $t = 1, j = 1$
- 5: Initialize weights w_i with random numbers
- 6: Initialize the input layer with \vec{x}_j
- 7: Calculate the output using $\sum w_i x_i + w_0$
- 8: Calculate the error $(y - \hat{y})$.
- 9: Update the weights $w_j(t+1) = w_j - \eta(y - \hat{y})x_j$
- 10: Repeat steps 7 and 9 until: the error is less than θ or a predetermined number of epochs have been completed.

To provide a stable weight update for this step, $w_j(t+1) = w_j - \eta(y - \hat{y})x_j$, we require a small η . This results in slow learning. Bigger η would be good for fast learning. What are the problems? . What is the compromise?

¹An epoch is one complete presentation of the data set to be learned to a learning machine.

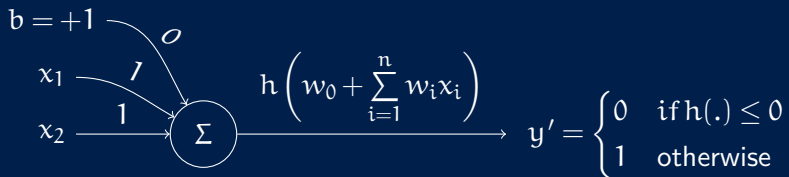
LOGICAL AND



Input x_1	Input x_2	$x_1.w_1 + x_2.w_2 + b.w_b$	output- y
0	0	$0.1+0.1-1$	0
0	1	$0.1+1.1-1$	0
1	0	$1.1+0.1-1$	0
1	1	$1.1+1.1-1$	1

Here, the perceptron is already trained and the learned weights are shown in the diagram

LOGICAL OR



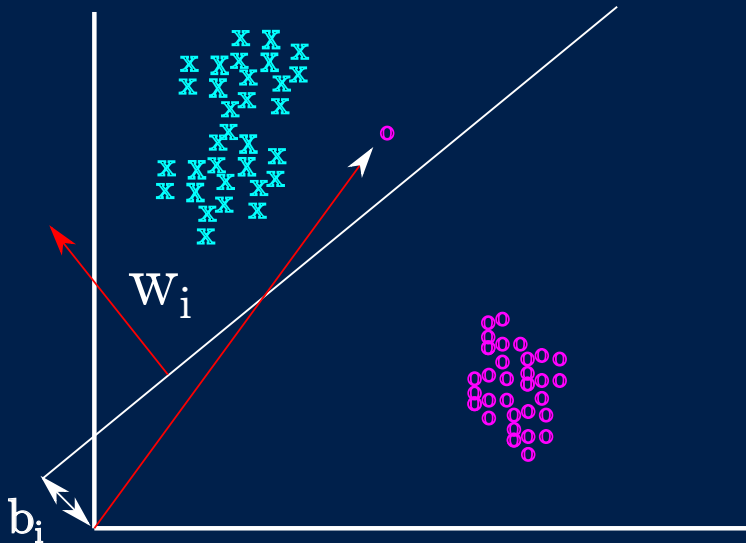
Input x_1	Input x_2	$x_1.w_1 + x_2.w_2 + b.w_b$	output- y

SENTIMENT ANALYSIS - USING PERCEPTRON

- ▶ Ability to classify reviews as positive or negative
- ▶ Positive and negative words for training
- ▶ Glove word embedding as features - input
 - ▶ 50 element word embedding²
 - ▶ Training Data generated using the intersection of the sentiment word list and word embedding from Glove

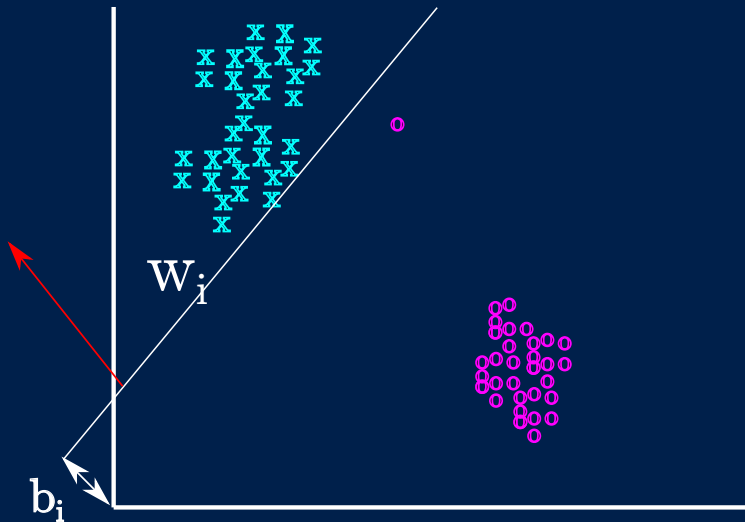
²data from <https://nlp.stanford.edu/projects/glove/>

PERCEPTRON LEARNING

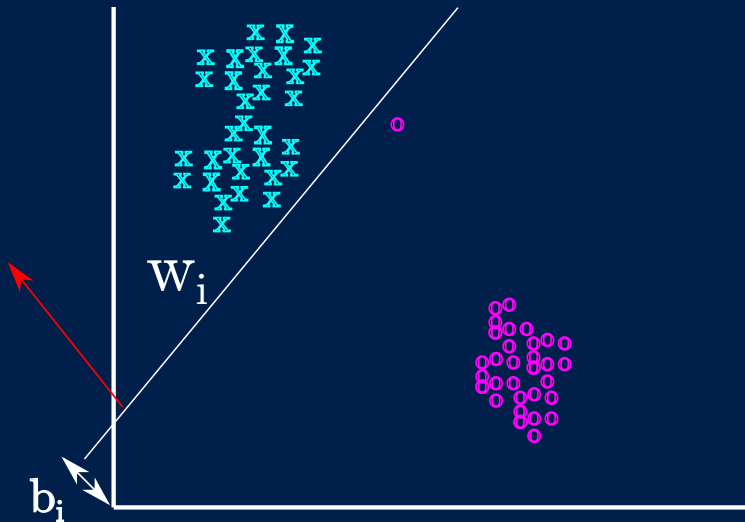


Figure

PERCEPTRON LEARNING



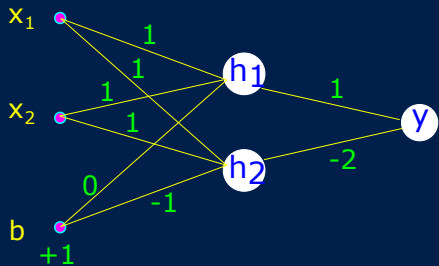
PERCEPTRON LEARNING



PERCEPTRON LIMITATIONS

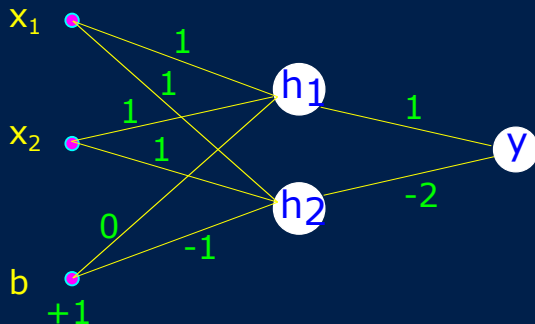
- ▶ It is based on the linear combination of fixed basis functions
- ▶ Updates the model only based on misclassification
- ▶ Documents that are linearly separable are classified

LOGICAL XOR



Input x_1	Input x_2	output- y
0	0	0
0	1	1
1	0	1
1	1	0

LOGICAL XOR



Input x_1	Input x_2	b	h_1	h_2	output- y
0	0	1	$f(0) = 0$	$f(-1) = 0$	0
0	1	1	$f(1) = 1$	$f(0) = 0$	1
1	0	1	$f(1) = 1$	$(f0) = 0$	1
1	1	1	$f(2) = 2$	$(f1) = 1$	0

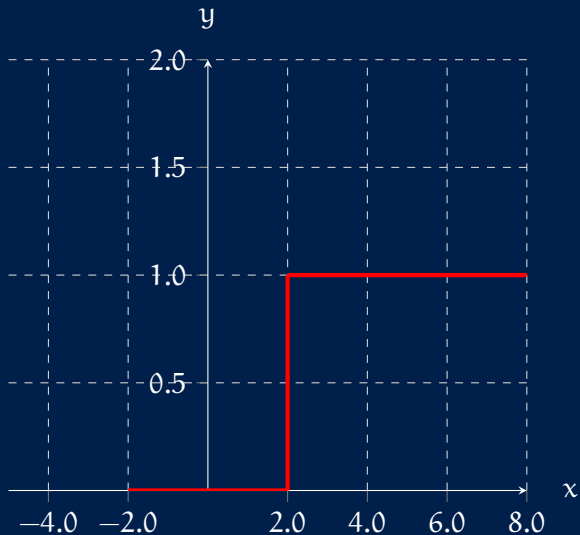
INTUITION

- ▶ Input space is transformed into hidden space
- ▶ Hidden layer represents the input layer
- ▶ Learns automatically the input representation and patterns
- ▶ $(0,1)$ and $(1,0)$ are merged into one in the h-space
- ▶ Patterns yielding similar results are merged into one
- ▶ Dimensionality reduction
- ▶ Learns to extract meaningful/latent features or representations from the input data
- ▶ Transforms the input information and create a representation in the new space
- ▶ Uses activation functions (e.g., sigmoid, tanh, or ReLU) to introduce non-linearity into the network, allowing it to capture complex relationships between input and output
- ▶ Are hidden layer neurons joining piecewise linear representations to create non-linear boundaries?

ACTIVATION FUNCTIONS

- ▶ Hard threshold
- ▶ Sigmoid
- ▶ Tanh
- ▶ ReLu - Rectified Linear Unit
- ▶ Leaky ReLu
- ▶ Softmax

HARD THRESHOLD



SIGMOID 1/3

Let us consider two classes c_1 and c_2 . The posterior probability for c_1 using Bayes theorem is

$$\begin{aligned} p(c_1|w) &= \frac{p(w|c_1)p(c_1)}{p(w|c_1)p(c_1) + p(w|c_2)p(c_2)} \\ &= \frac{1}{1 + \frac{p(w|c_2)p(c_2)}{p(w|c_1)p(c_1)}} \end{aligned}$$

$$\text{where } x = \log \left(\frac{p(w|c_1)p(c_1)}{p(w|c_2)p(c_2)} \right)$$

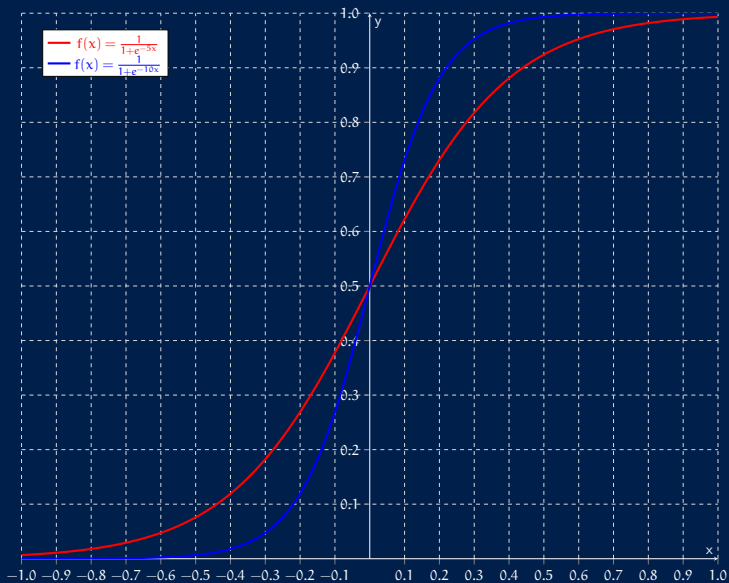
or

$$\sigma(x) = \frac{1}{1 + \exp(-x)} \Rightarrow \text{Sigmoid}$$

$$\sigma(-x) = 1 - \sigma(x)$$

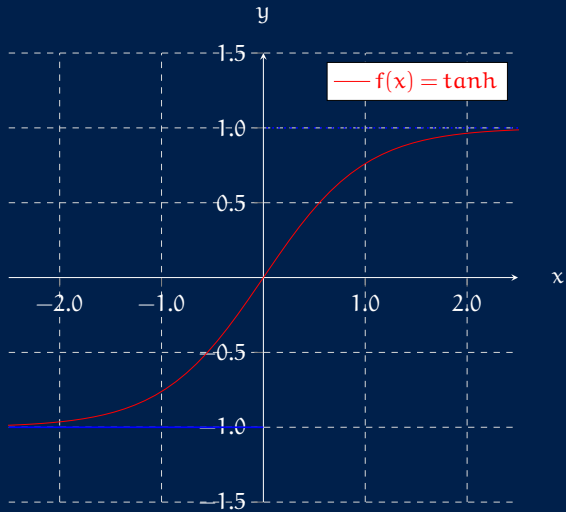
- ▶ The sigmoid is a non-linear function
- ▶ Better than hard threshold function as it squashes the net output into the range $[0, 1]$
- ▶ The values closer to the tails become 0 or 1
- ▶ In some cases, the values quickly saturate at 0 or 1
- ▶ At the bottom tail, most values become zero during the training and hence the most important aspect of learning of neural network is inhibited
- ▶ Sigmoid outputs are not zero-centered - $[0, 1]$
- ▶ It is undesirable to have all the values squashed near the tails, where the gradient is 0

SIGMOID 3/3



TANH

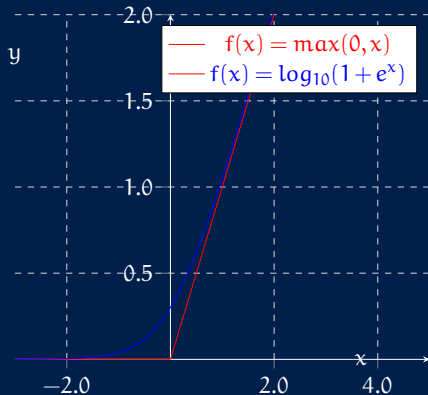
This is a zero based non-linear function.



RELU - RECTIFIED LINEAR UNIT

- ▶ There is a continuous gradient for the neurons to be in active state
- ▶ Produces a non-zero gradient for values closer to zero
- ▶ Leaky ReLU
$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ 0.01x, & \text{otherwise} \end{cases}$$
- ▶ Produces efficient propagation of the gradient
- ▶ Computationally efficient
- ▶ Scale invariant
- ▶ Unbounded and not zero centered

Learning rate (usually, very small) has to be fine tuned to minimize the death of neurons



MULTICLASS DECISION FUNCTION

- ▶ All linear classifier are used for binary classifications
- ▶ In NLP problems, we need to identify more than two classes
 - ▶ Document classification
 - ▶ Sentiment Analysis - positive, negative, neutral and non-sentiment word
- ▶ We need a decision function that predicts more than two classes by providing appropriate values
- ▶ An extension of the case function would be hard to manage

- ▶ Need a function that takes as input a vector of of size with N real numbers, and normalizes it into a K classes.
- ▶ Need a function that normalizes the net output and classes well separated (ideal condition)
- ▶ Need a function that fits the classes using probability and distributes the probability density

$$\begin{aligned} p(c_j|x) &= \frac{p(x|c_j)p(c_j)}{\sum_k (p(x_k|c_k)p(c_k))} \\ &= \frac{\exp(x_k)}{\sum_k \exp(x_k)} \Rightarrow \text{Softmax} \end{aligned}$$

ACTIVATION FUNCTION - PYTHON CODE

```
import numpy as np

def sigmoid(X,W,b):
    return 1.0/(1.0+ np.exp(-(np.dot(W.T,X)+b)))

def tanh(X,W,b):
    z = np.exp(-(np.dot(W.T,X)+b))
    return (np.exp(z) - np.exp(-z))/(np.exp(z) + np.exp(-z))

def relu(X,W,b):
    return np.maximum(np.dot(W.T,X) + b,0)

def softmax(X,W,b):
    return np.exp(np.dot(W,X)+b)/np.sum(z_exp)
```

```
if __name__ == '__main__':  
    W=np.array([0.1, 0.2, 0.6])  
    X=np.array([0.2, 0.1, 0.3])  
    b=1.5  
  
    print(sigmoid(X,W,b))      #0.848128836343  
    print(tanh(X,W,b))         #0.177176480951  
    print(relu(X,W,b))         #1.72  
  
    W = np.array([[1,2,3],[2,3,8],[1,5,7]])  
  
    print(softmax(X,W,b))
```

OBJECTIVE FUNCTION

- ▶ Provides information related to how well the model is in sync with the original data or gold standard
- ▶ Refers to a function used for optimization in a machine learning model, whether through minimization or maximization.
- ▶ In the context of machine learning, it usually refers to minimizing errors
- ▶ In reinforcement learning, it could refer to maximizing the reward

LOSS FUNCTION

$$\hat{y} = \sum_{i=1}^m w_i^T x_i + w_0 \quad (1)$$

where \hat{y} is the predicted value

w_0 is the bias

x is the input vector w is the weight vector

if y is the target, then the loss function is defined as a squared function

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

The main idea is to reduce the residual $(y - \hat{y})$. When the value of L becomes negligible, we have predicted the vector to belong to a known class.

The loss function computes the error for a single training example

Cross Entropy Loss - the distance between output distribution of the model and the original distribution

COST FUNCTION

Let us assume that we have a set of vectors for training. In the case of the sentiment analysis, these are the vectors obtained using any of the word embedding methods, representing the sentiment words. We could also use one-hot vectors representing the same

$$\mathbf{X} = [x_1, x_2, x_3, \dots, x_N] \quad (3)$$

Combining the model parameters $w_0, w_1, w_2, w_3, \dots, w_n$ with the loss function, we get a **Cost function**, averaged over all the input training samples

$$J(\theta) = \frac{1}{2} \sum_i L(y_i, \hat{y}_i)^2 \quad (4)$$

- ▶ The loss is a function of prediction and target values of a single training example
- ▶ The cost is a function of model parameters and bias - average of the loss over all training samples

GRADIENT DESCENT

- ▶ GD iteratively used to adjust the weights and as a result to minimize the cost function
- ▶ Initialize the weights to random values
- ▶ Iteratively adjust the weights in the direction of the steepest descent or in the direction that most decreases the cost function. To update the weights in the steepest descent, a learning parameter η is used

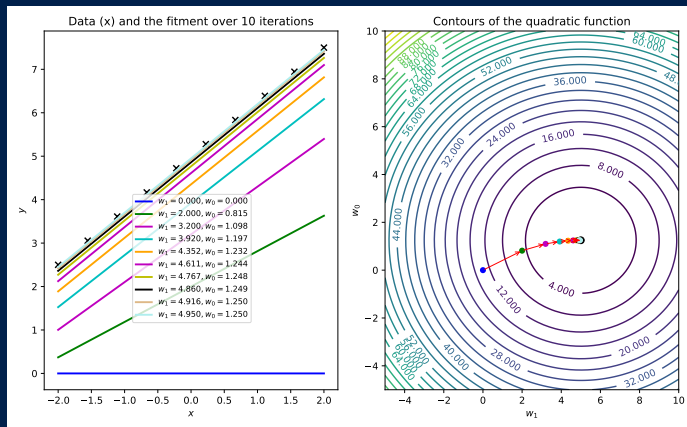
$$w_j \leftarrow w_j - \eta \frac{\partial J(\theta)}{\partial w_j} \quad (5)$$

$$= w_j - \eta \sum_{i=1}^N x_j^i (y - \hat{y}) \quad (6)$$

where η is the learning parameter and usually takes the value between 0.01 and 0.001

GD ADVANTAGES

- ▶ Iterative
- ▶ Computationally efficient
- ▶ Generic and could be used to solve even non-linear equations
- ▶ Suitable for large models
- ▶ It works!
- ▶ It is very slow when it reaches close to the local minima



SEQUENTIAL NATURE OF DATA

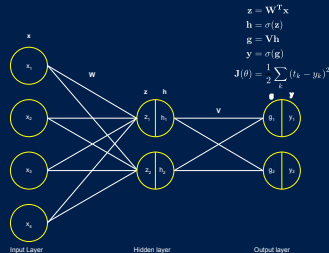
- ▶ Speech
- ▶ Documents
- ▶ Videos
- ▶ Weather forecast
- ▶ Financial - Stock market

PROPERTIES OF ANN

- ▶ Massively parallel distributed structure
- ▶ Ability to learn
- ▶ Ability to learn from training samples
- ▶ Ability to find latent patterns in the data
- ▶ Generalize and associate data

BACKPROPAGATION MODEL

The goal of backpropagation is to change the weights so that the *estimated target* \approx target, thereby minimizing the error for each neuron and the network as a whole.

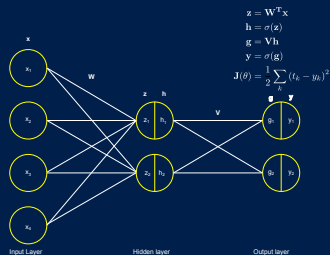


The goal is to minimize

$$J(\theta) = 1/2 \left((t_1 - y_1)^2 + (t_2 - y_2)^2 \right) \quad (7)$$

- * We want to adjust weights coming in and going out of hidden layer so that $t - y$ is minimized
- * $\Delta W \propto -\frac{\partial J(\theta)}{\partial W}$

BACK PROPAGATION MODEL - FORWARD PASS



$$z_1 = \mathbf{w}^T \mathbf{x} + b_1 \quad (8)$$

$$z_2 = \mathbf{w}^T \mathbf{x} + b_1 \quad (9)$$

$$z_1 = x_1 \cdot w_{11} + x_2 \cdot w_{21} + \dots + b_1 \quad (10)$$

$$z_2 = x_2 \cdot w_{12} + x_2 \cdot w_{22} + \dots + b_1 \quad (11)$$

$$h_1 = \sigma(z_1) = \frac{1}{1 + e^{-z_1}} \quad (12)$$

$$h_2 = \sigma(z_2) = \frac{1}{1 + e^{-z_2}} \quad (13)$$

$$g_1 = h_1 * v_{11} + h_2 \cdot v_{21} \quad (14)$$

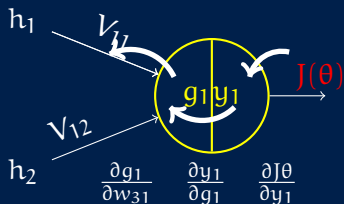
$$g_2 = h_2 * v_{12} + h_2 \cdot v_{22} \quad (15)$$

$$y_1 = \sigma(g_1) = \frac{1}{1 + e^{-g_1}} \quad (16)$$

$$y_2 = \sigma(g_2) = \frac{1}{1 + e^{-g_2}} \quad (17)$$

BACKWARD PASS—ADJUST HIDDEN-OUTPUT LAYER WEIGHTS

$$J(\theta) = 1/2 \left((t_1 - y_1)^2 + (t_2 - y_2)^2 \right)$$



$$\frac{\partial J(\theta)}{\partial v_{11}} = -\alpha \left(\frac{\partial J(\theta)}{\partial y_1} \frac{\partial y_1}{\partial g_1} \frac{\partial g_1}{\partial v_{11}} \right) \quad (18)$$

$$\frac{\partial J(\theta)}{\partial y_1} = -(t_1 - y_1) \quad (19)$$

$$\frac{\partial y_1}{\partial g_1} = y_1(1 - y_1) \quad (20)$$

$$\frac{\partial g_1}{\partial v_{11}} = h_1 \quad (21)$$

$$\frac{\partial J(\theta)}{\partial v_{11}} = -\alpha \left(\frac{\partial J(\theta)}{\partial y_1} \frac{\partial y_1}{\partial g_1} \frac{\partial g_1}{\partial v_{11}} \right) \quad (22)$$

$$= \alpha(t_1 - y_1)y_1(1 - y_1)h_1 \quad (23)$$

Now, the weights can be updated by $v_{11}^{t+1} = v_{11}^t - \eta * \frac{\partial J(\theta)}{\partial v_{11}^t}$. In the same fashion, compute the error to be propagated back to the other weights.

BACKWARD PASS – ADJUST INPUT-HIDDEN LAYER WEIGHTS

$$\frac{\partial J(\theta_1)}{\partial w_{11}} = -\alpha \frac{\partial J(\theta_1)}{\partial y_1} \frac{\partial y_1}{\partial g_1} \frac{\partial g_1}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial w_{11}} \quad (24)$$

$$\frac{\partial J(\theta_2)}{\partial w_{11}} = -\alpha \frac{\partial J(\theta_2)}{\partial y_2} \frac{\partial y_2}{\partial g_2} \frac{\partial g_2}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial w_{11}} \quad (25)$$

$$\frac{\partial J(\theta)}{\partial w_{11}} = \frac{\partial J(\theta_1)}{\partial w_{11}} + \frac{\partial J(\theta_2)}{\partial w_{11}} \quad (26)$$

$$\frac{\partial J(\theta_1)}{\partial y_1} = -(t_1 - y_1) \quad (27)$$

$$\frac{\partial y_1}{\partial g_1} = y_1(1 - y_1) \quad (28)$$

$$\frac{\partial g_1}{\partial h_1} = v_{11} \quad \frac{\partial z_1}{\partial w_{11}} = x_1 \quad (29)$$

$$\frac{\partial h_1}{\partial z_1} = \sigma(z_1)(1 - \sigma(z_1)) \quad (30)$$

$$\frac{\partial J(\theta_2)}{\partial y_2} = \dots \quad (31)$$

Now, input-hidden layer weights can be updated using $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta * \frac{\partial \mathbf{J}(\theta)}{\partial \mathbf{w}}$

Once trained,

- ▶ The hidden layer of a trained model is a lookup table
- ▶ Hidden weights act as an associative memory and capture the relational similarities

INTRODUCTION TO LOGISTIC REGRESSION

- ▶ **Logistic regression** is a classification algorithm used to predict the probability of a binary outcome.
- ▶ We want to model the conditional probability $p(y = 1|x)$ as a function of x for the binary output variable y and making $p(.)$ be a linear function of x won't work
- ▶ Consider using $\log p(x)$ as a linear function - unbounded in one direction Consider $\frac{p}{1-p}$ is even better. This is known as logistic transformation of **logit**
- ▶ The prediction is based on the logistic function (sigmoid):

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

where x is the input vector, w is the parameter vector.

- ▶ Logistic regression outputs a probability and classifies based on a threshold (usually 0.5).

- ▶ **Word2Vec** is a neural network model that learns word embeddings.
- ▶ Words are mapped to high-dimensional vectors where semantically similar words have similar vector representations.
- ▶ The word embeddings are pre-trained and can be used as features in downstream tasks.
- ▶ Embeddings transform sparse word representations (one-hot vectors) into dense, lower-dimensional continuous vectors.

Embedding: $\mathbf{w} \in \mathbb{R}^d$

REFERENCES

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006. ISBN: 0387310738.