

Course number: 80240743

Deep Learning

Xiaolin Hu (胡晓林)

Dept. of Computer Science and Technology

Tsinghua University

Topic 2: Basic Knowledge

Xiaolin Hu

Department of Computer Science and Technology
Tsinghua University

Outline

- Math basics
- Machine learning basics
- Neural networks in the early stage

Math basics

LINEAR ALGEBRA

Math objects

- Scalar

- A single number, often denoted by a lower case letter without boldface, e.g., a, b, x

- Vector

- An array of numbers, often denoted by a lowercase letter with boldface, e.g., $\mathbf{a}, \mathbf{b}, \mathbf{x}$

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}$$

- Matrix


- A 2D array of numbers, often denoted by an uppercase letter with boldface, e.g., $\mathbf{A}, \mathbf{B}, \mathbf{X}$

$$\mathbf{A} = \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix}$$

- Tensor

- An n -D array of numbers, often denoted like this:
e.g., $\mathbf{A}, \mathbf{B}, \mathbf{X}$

$$\mathbf{A} = \left(\begin{pmatrix} A_{1,1,1} & A_{1,2,1} \\ A_{2,1,1} & A_{2,2,1} \end{pmatrix}, \begin{pmatrix} A_{1,1,2} & A_{1,2,2} \\ A_{2,1,2} & A_{2,2,2} \end{pmatrix} \right)$$



But I sometimes
may not follow
these conventions

Simple operations

- Matrix transpose: \mathbf{A}^\top

$$\mathbf{A} = \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix} \quad \mathbf{A}^\top = \begin{pmatrix} A_{1,1} & A_{2,1} \\ A_{1,2} & A_{2,2} \end{pmatrix}$$

- A vector can be viewed as a special matrix

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \quad \mathbf{a}^\top = (a_1, a_2, a_3)$$

- Matrix product: if $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$, then $\mathbf{C} = \mathbf{AB}$ with shape $m \times p$ and $C_{i,j} = \sum_k A_{i,k} B_{k,j}$
- Elementwise product (Hadamard product): $\mathbf{C} = \mathbf{A} \odot \mathbf{B}$ where the 3 matrices are of the same shape and $C_{i,j} = A_{i,j} B_{i,j}$

Simple operations

- Properties

- $A(B + C) = AB + AC$
- $A(BC) = (AB)C$
- $AB \neq BA$ but $x^\top y = y^\top x$
- $(AB)^\top = B^\top A^\top$

- Identity matrix: $I \in \mathbb{R}^{n \times n}$

- For $x \in \mathbb{R}^n$, we have $Ix = x$
- For $A \in \mathbb{R}^{n \times m}$, we have $IA = A$

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

- Matrix inverse of a square matrix A is a matrix denoted by A^{-1} such that $A^{-1}A = I$

Linear equations

- A system of linear equations

$$\mathbf{Ax} = \mathbf{b}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{x} \in \mathbb{R}^n$, or equivalently

$$A_{1,1}x_1 + A_{1,2}x_2 + \cdots + A_{1,n}x_n = b_1$$

$$A_{2,1}x_1 + A_{2,2}x_2 + \cdots + A_{2,n}x_n = b_2$$

...

$$A_{m,1}x_1 + A_{m,2}x_2 + \cdots + A_{m,n}x_n = b_m$$

- If $m = n$ and \mathbf{A}^{-1} exists, then there is exactly one solution \mathbf{x}^* to the above equations
- If $m < n$, there are infinite number of solutions **underdetermined**
- If $m > n$, the solution may not exist! **overdetermined**

Norms

- Norm is used to measure the “size” of a vector or matrix
- The L_p norm of a vector is defined as

$$||\mathbf{x}||_p = \left(\sum_i |x_i|^p \right)^{1/p}$$

for $p \in \mathbb{R}, p > 1$

- The L_2 norm is known as the **Euclidean norm**, and $||\mathbf{x}||_2 = \sqrt{\mathbf{x}^\top \mathbf{x}}$
- The L_1 norm is simplified to $||\mathbf{x}||_1 = \sum_i |x_i|$
- The L_∞ norm is simplified to $||\mathbf{x}||_\infty = \max_i |x_i|$

Norms

- A norm is any function f that satisfies
 - $f(\mathbf{x}) \geq 0$
 - $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = 0$
 - $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (the triangle inequality)
 - $f(\alpha \mathbf{x}) = |\alpha|f(\mathbf{x}), \forall \alpha \in \mathbb{R}$
- Frobenius norm of a matrix

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$$

Math basics

PROBABILITY THEORY

Random variable

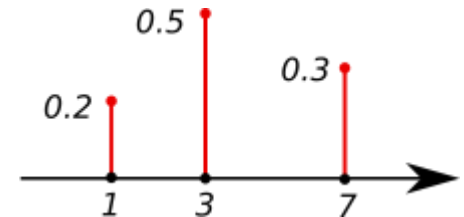
- A **random variable** is a variable that can take on different values randomly
 - Denote the random variable by x and its two possible values by x_1 and x_2
 - For vectors, we write the random variable as \mathbf{x} and one of its values as \mathbf{x}
 - **Discrete** versus **continuous**



Probability distribution

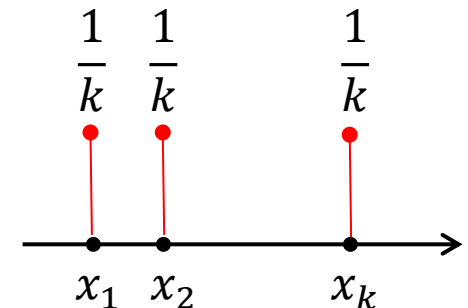
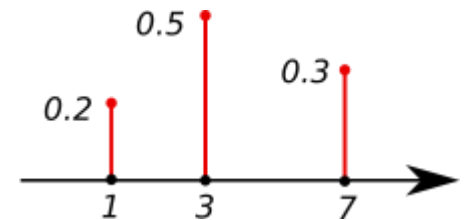
A **probability distribution** is a distribution of how likely a random variable or a set of random variables is to take on each of its possible states

- A probability distribution over **discrete** variables may be described using a **probability mass function** (PMF)
 - The prob that $x = x$ is denoted as $P(x)$ or $P(x = x)$
 - $x \sim P(x)$ specify which distribution x follows
- Joint probability
 - $P(x = x, y = y)$ or $P(x, y)$ denotes the prob that $x = x$ and $y = y$ simultaneously



Probability mass function

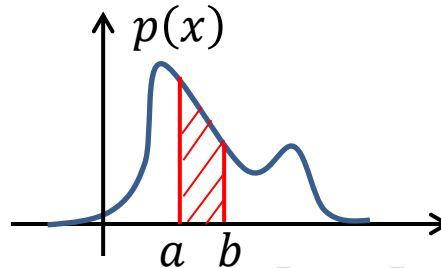
- To be a PMF of a random variable x , a function P must satisfy:
 - The domain of P must be the set of all possible states of x
 - $\forall x \in \mathcal{X}, 0 \leq P(x) \leq 1$
 - $\sum_{x \in \mathcal{X}} P(x) = 1$
- Uniform distribution
 - Consider a single discrete random variable x with k different states
 - $P(x = x_i) = \frac{1}{k}, \forall i$



Probability density function

- A probability distribution over **continuous** variables may be described using a **probability density function** (PDF)
- To be a PDF, a function p must satisfy the following properties

- The domain of p must be the set of all possible states of x
- $\forall x \in \mathbf{x}, p(x) \geq 0$
- $\int p(x)dx = 1$



Can $p(x) > 1$?

- The prob that x lies in the interval $[a, b]$ is given by

$$\int_a^b p(x)dx$$

- Note $p(x)$ does not give the prob of a specific state directly

Marginal probability

Suppose we know the prob distribution over a set of variables. The prob distribution over just a subset of them is known as the **marginal prob distribution**

- Let $P(x, y)$ denote the prob distribution of discrete random variables x and y , then

$$P(x = x) = \sum_y P(x = x, y = y)$$

- Let $p(x, y)$ denote the PDF of continuous random variables x and y , then

$$p(x) = \int P(x, y) dy$$

Conditional probability

The **conditional probability** is the probability of some event, given that some other event has happened

- The conditional prob that $y = y$ given $x = x$ is denoted by $P(y = y|x = x)$, which can be calculated as

$$P(y = y|x = x) = P(y = y, x = x)/P(x = x)$$

- The chain rule

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} | x^{(1)}, \dots, x^{(i-1)})$$

- **Exercise:** With this definition, is the following correct?

$$P(a, b, c) = P(a|b, c)P(b|c)P(c)$$

Independence and conditional independence

Two random variables x and y are **independent** if their joint probability distribution can be expressed as a product of two factors, one involving x and one involving y :

$$\forall x \in \mathcal{X}, y \in \mathcal{Y}, \\ p(x = x, y = y) = p(x = x)p(y = y)$$

Denoted by $x \perp y$

Two random variables x and y are **conditionally independent** given a random variable z if the conditional probability distribution over x and y factorizes in this way:

$$\forall x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}, \\ p(x = x, y = y | z = z) = p(x = x | z = z)p(y = y | z = z)$$

Denoted by $x \perp y | z$

Expectation

The **expectation**, or **expected value**, of some function $f(x)$ w.r.t. a prob distribution $P(x)$ is the average value that f takes on when x is drawn from P

- For discrete variables

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x)$$

- For continuous variables

$$\mathbb{E}_{x \sim P}[f(x)] = \int p(x)f(x)dx$$

- If the identity of the distribution is clear, we may write $\mathbb{E}_x[f(x)]$
- Expectation is **linear**: if α and β do not depend on x , then

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)]$$

Variance and covariance

- The **variance** of a function $f(x)$

$$\text{Var}(f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

- It measures how much the value of the function f varies when x is sampled from its distribution

- The **covariance** of two functions $f(x)$ and $g(x)$

$$\text{Cov}(f(x), g(y)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$

- It measures how much two values are linearly related to each other, as well as the scale of these variables
- High absolute values of the covariance mean that the values change very much and are both far from their respective means

- The **covariance matrix** of a random vector $\mathbf{x} \in \mathbb{R}^n$ is an $n \times n$ matrix

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j)$$

f and g can be identity functions



Common prob distributions

- **Bernoulli distribution**: over a single binary random variable

$$P(x = 1) = \phi, P(x = 0) = 1 - \phi$$

$$P(x = x) = \phi^x (1 - \phi)^{1-x}$$

$$\mathbb{E}_x[x] = \phi, \text{Var}(x) = \phi(1 - \phi)$$

- **Multinoulli or categorical distribution**: over a single discrete variable with k different states where k is finite

$$P(x = i | \mathbf{p}) = p_i$$

where $\mathbf{p} \in [0,1]^k$ and $\sum_{i=1}^k p_i = 1$

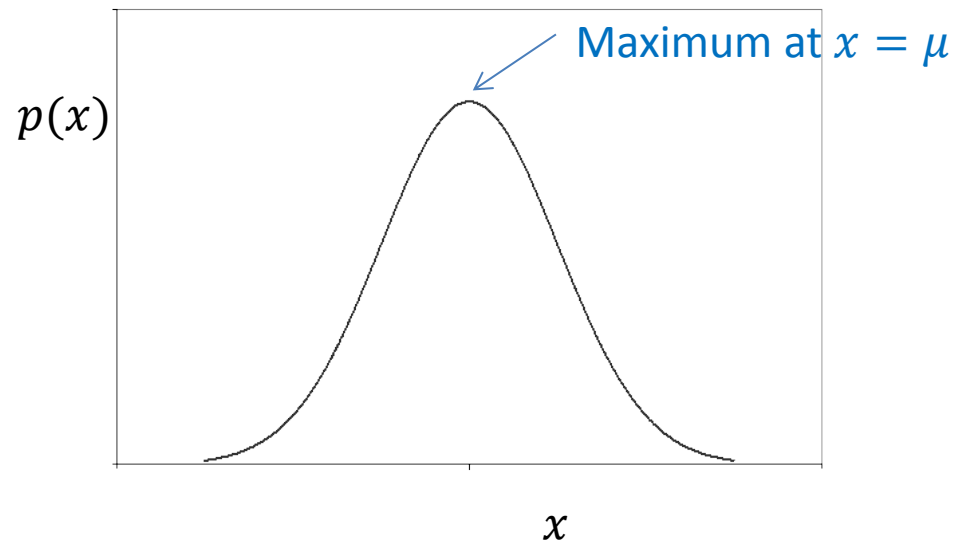
Common prob distributions

- **Gaussian distribution** or **normal distribution**: over a continuous variable

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- Mean: $\mathbb{E}[x] = \mu$
- Variance: $\text{Var}[x] = \sigma^2$
- Standard deviation: σ

The *central limit theorem* shows that the sum of many independent variables is approximately normally distributed



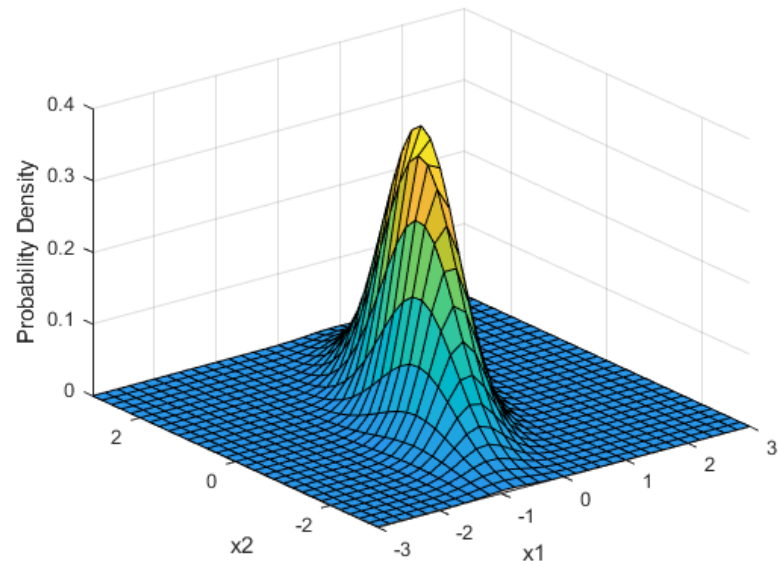
Common prob distributions

- **Multivariate normal distribution:** over a continuous vector

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) = \sqrt{\frac{1}{(2\pi)^2 \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

- Mean: $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$
- Covariance: $\text{Cov}[\mathbf{x}] = \boldsymbol{\Sigma}$

Isotropic Gaussian distribution: $\boldsymbol{\Sigma} = \alpha \mathbf{I}$
where α is a scalar



Bayes' rule



Thomas Bayes (1702~1761)

Prior probability of x

Posterior probability of x

Probability of y given x ; likelihood

Prior probability of y

$$P(x|y) = \frac{P(x)P(y|x)}{P(y)}$$

Math basics

OPTIMIZATION

Gradient-based optimization

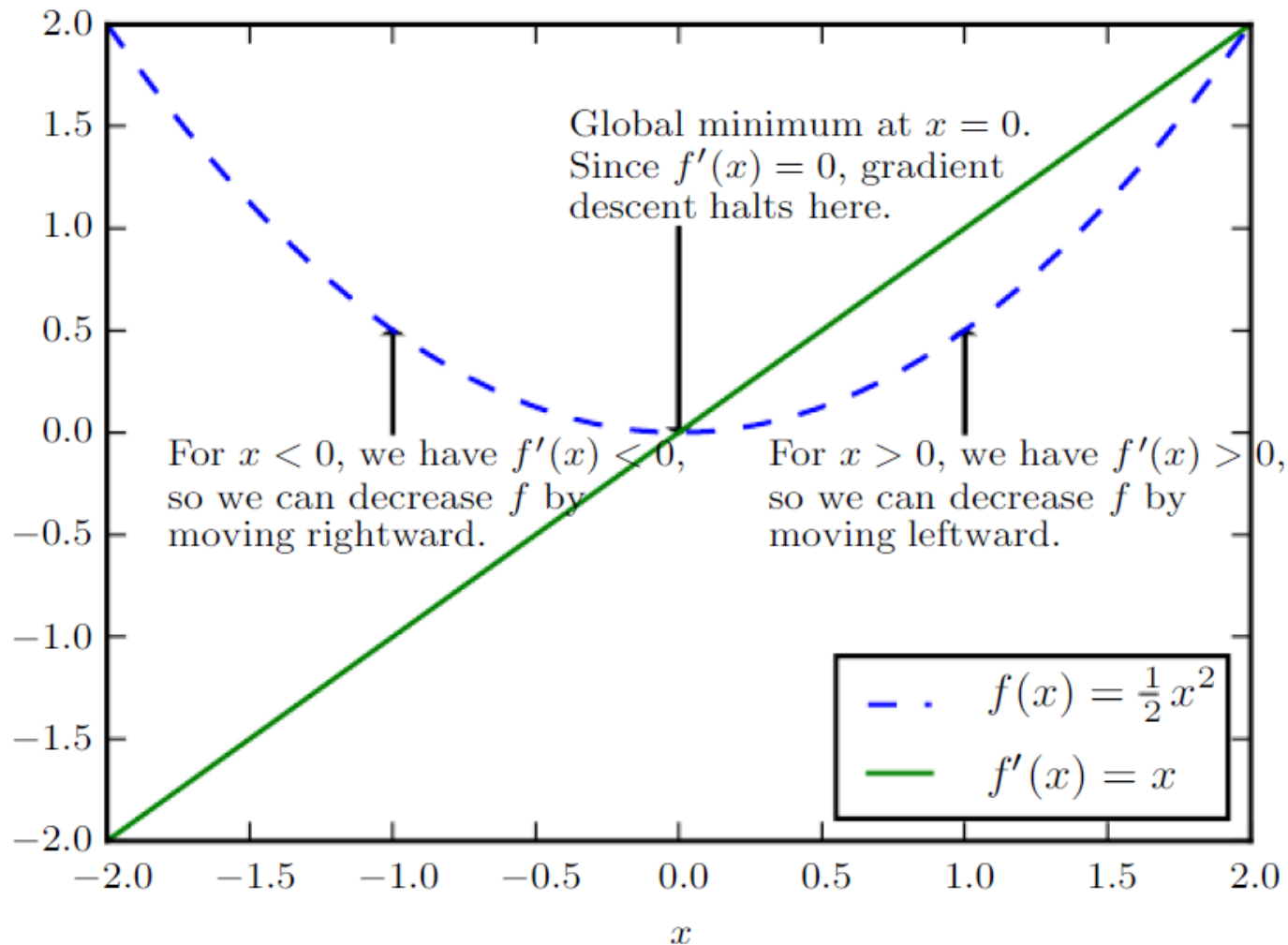
- The function we want to minimize or maximize is called **objective function**
- When we are minimizing it, we may also call it the **cost function**, **loss function**, or **error function**
- The **derivative** of a function $y = f(x)$, denoted by $f'(x)$ or $\frac{dy}{dx}$, gives the slope, or **gradient**, of f at the point x
- Gradient descent

$$x' = x - \eta f'(x)$$

where $\eta > 0$ is the **learning rate**

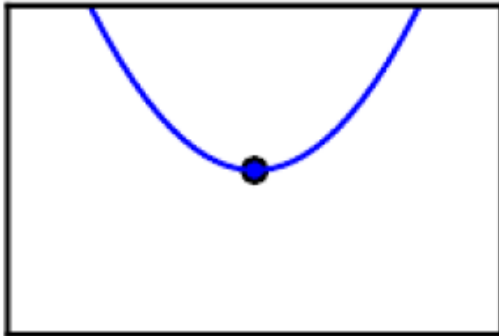
Gradient descent

$$x' = x - \eta f'(x)$$

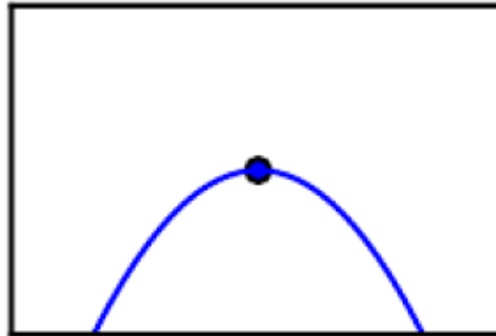


Critical points

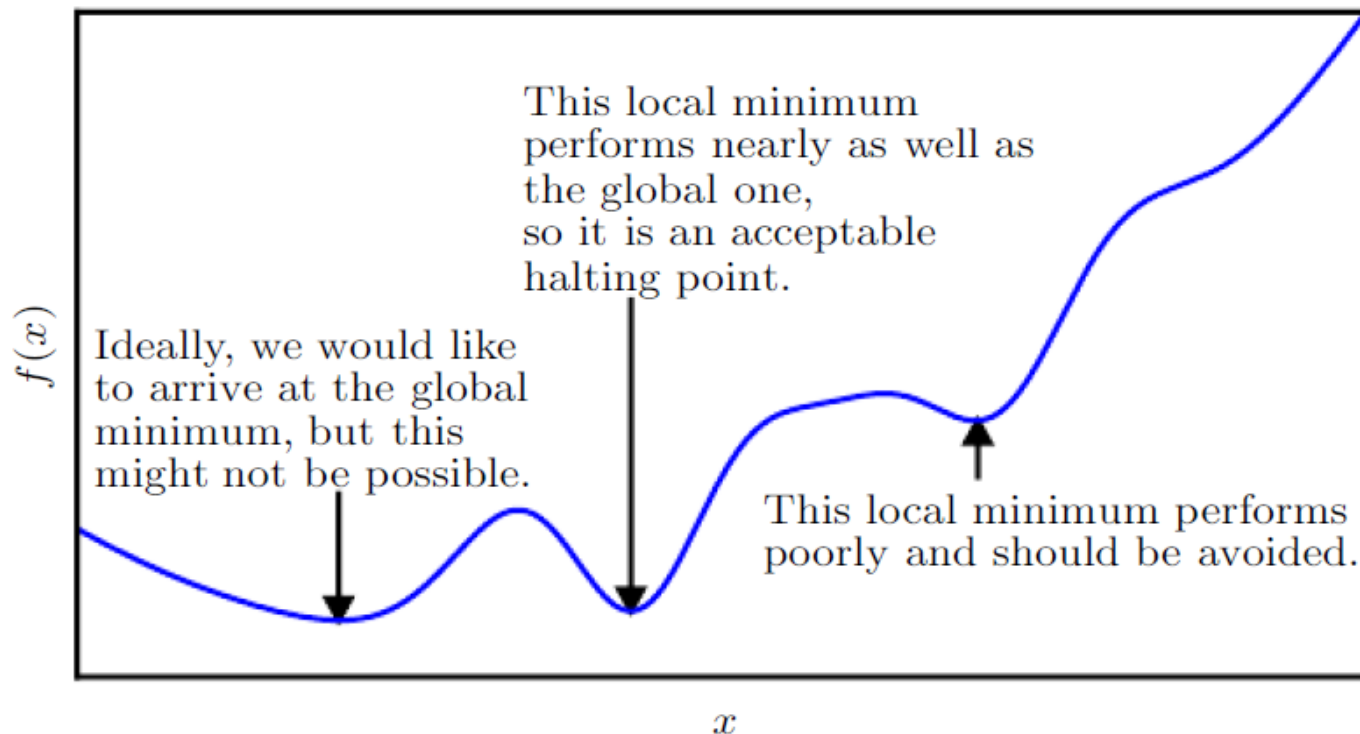
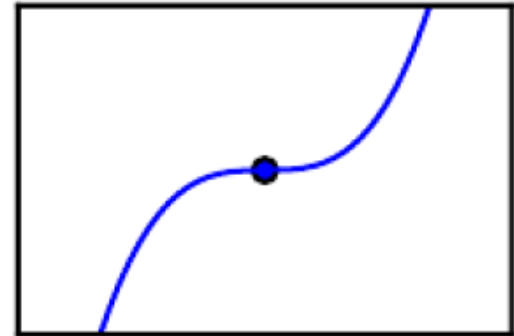
Minimum



Maximum



Saddle point



Gradient decent for multivariate functions

- For a function of a single variable $y = f(x)$, the gradient decent method is

$$x' = x - \eta f'(x)$$

- For a function $y = f(\mathbf{x})$, the **partial derivative** is denoted by $\partial f / \partial x_i$
- The gradient decent method becomes

$$\mathbf{x}' = \mathbf{x} - \eta \nabla_{\mathbf{x}} f(\mathbf{x})$$

where $\eta > 0$ and $\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \dots \\ \partial f / \partial x_n \end{pmatrix}$

Rules in calculus

- **Chain rule:** the derivative of the composition function $f(g(x))$ is

$$[f(g(x))]' = f'(g(x))g'(x)$$

or in Leibniz's notation

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

- **Product rule:** the derivative of product of two functions

$$(f \cdot g)' = f' \cdot g + f \cdot g'$$

or in Leibniz's notation

$$\frac{d}{dx}(u \cdot v) = \frac{du}{dx} \cdot v + u \cdot \frac{dv}{dx}$$

- **Quotient rule**

$$\frac{d}{dx}\left(\frac{f(x)}{g(x)}\right) = \frac{f'g - fg'}{g^2}$$

Derivative of two-step composition

Suppose we have:

- Independent input variables x_1, x_2, \dots, x_n
- Dependent **intermediate variables**, u_1, u_2, \dots, u_m , each of which is a function of x_1, x_2, \dots, x_n
- Dependent output variables w_1, w_2, \dots, w_p , each of which is a function of u_1, u_2, \dots, u_m

Then for any $i \in \{1, 2, \dots, p\}$ and $j \in \{1, 2, \dots, n\}$ we have

$$\frac{\partial w_i}{\partial x_j} = \sum_{k=1}^m \frac{\partial w_i}{\partial u_k} \frac{\partial u_k}{\partial x_j}$$

Sum over the intermediate variables

From wikipedia

Summary so far

- Linear algebra
 - Math objects: Scalars, vectors, matrices, tensors
 - Simple operations: matrix transpose, inverse, product
 - Norms: L_p norm
- Probability theory
 - Random variables: discrete, continuous
 - Prob distribution: PMF and PDF
 - Marginal probability
- Conditional probability
- Independence and conditional independence
- Expectation, variance and covariance
- Common prob distributions
- Bayes' rule
- Optimization
 - Gradient descent
 - Critical points
 - Rules in calculus

Outline

- Math basics
- Machine learning basics
- Neural networks in the early stage

Learning algorithms

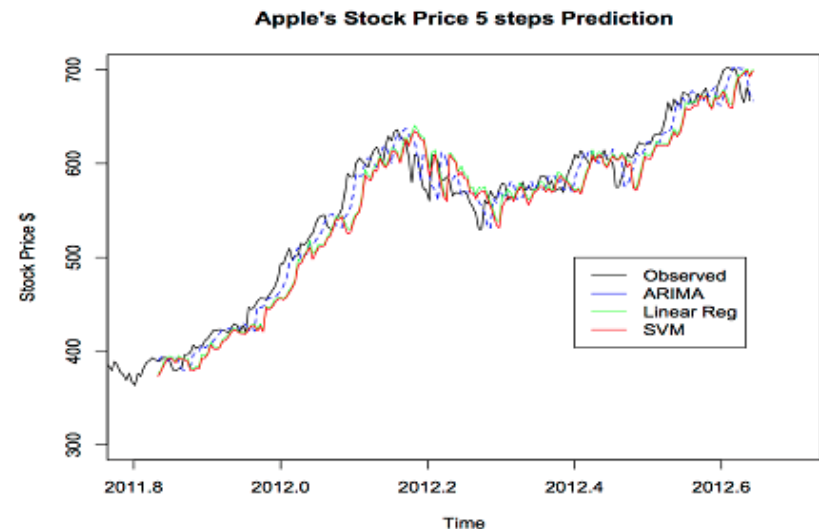
“A computer program is said to **learn** from experience E w.r.t. some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .” ---Tom Mitchell, 1997

- Machine learning (ML) tasks are usually described in terms of how the ML system should process an **example**
- An example is a collection of **features** that have been quantitatively measured from some object or event
 - Features of a bucket: color, diameter, height, material, etc
 - Features of an animal: size, shape, number of legs, , etc



The tasks T

- Classification
 - Suppose there are k categories. Find a function $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$



- Regression
 - Find a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, and m is often 1

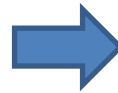
Regression results might be converted to classification results

The tasks T

- Synthesis and sampling
dataset



Synthesized using GAN

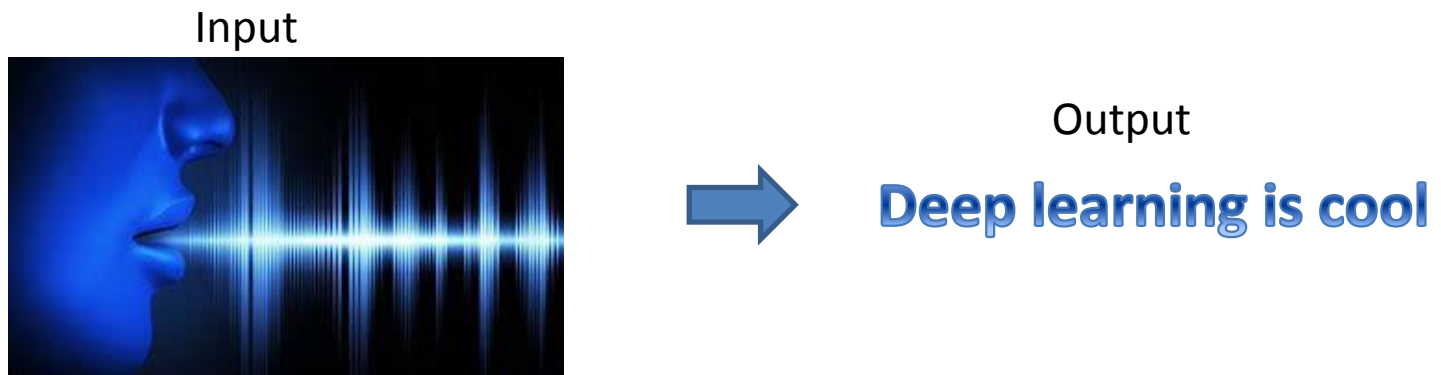


- Denoising



The tasks T

- Transcription



- Machine translation



The tasks, T

- Structured output
- Anomaly detection
- Synthesis and sampling
- Imputation of missing values
- Density estimation
- Etc.

The performance measure, P

- A performance measure is required to quantitatively evaluate the performance of a ML system
- Usually this measure P is **specific to the task T** being carried out by the system
 - Classification and transcription: accuracy or error rate
 - Regression and denoising: distance between the ground-truth and prediction
 - Synthesis, machine translation: difficult and sometimes need human evaluation
- What we are more interested in is the performance measure on a **test set** of data that is **separated** from the data used for training the system

The experience, E

- ML algorithms including deep learning algorithms can be broadly categorized as **unsupervised** or **supervised** by what kind of experience they are allowed to have during the learning process
- The algorithms experience a **dataset**, which is a collection of many **examples** or **data points** denoted by \mathbf{x}
 - We can view examples as samples of a random variable \mathbf{x}
- Unsupervised learning algorithms
 - experience a dataset containing many features, then learn useful properties or structures of this dataset
- Supervised learning algorithms
 - experience a dataset containing many features, but each example is also associated with a **label** or **target** denoted by \mathbf{y}

learn $p(\mathbf{x})$

learn $p(\mathbf{y}|\mathbf{x})$

Example: linear regression

- **Task T** : to predict y from x by outputting $\hat{y} = \mathbf{w}^\top \mathbf{x}$ x_i : feature
- **Performance P** : mean squared error of the model on the test with m test samples $\{(\mathbf{x}_i, y_i)\}^{\text{test}}$ w_i : weight

$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i (\hat{y}_i - y)^{\text{test}}$$

- **Experience E** : minimize the MSE on the training set of q samples $\{(\mathbf{x}_i, y_i)\}^{\text{train}}$

$$\text{MSE}_{\text{train}} = \frac{1}{q} \sum_i (\hat{y}_i - y)^{\text{train}}$$

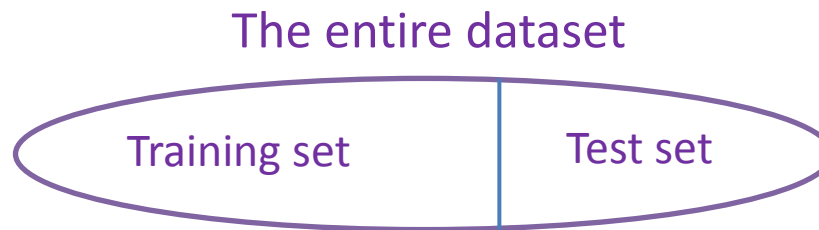
- Denote $\{(\mathbf{x}_i, y_i)\}^{\text{train}}$ collectively by $(\mathbf{X}^{\text{train}}, \mathbf{y}^{\text{train}})$, then

$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{w}} \frac{1}{q} \left\| \hat{\mathbf{y}}^{\text{train}} - \mathbf{y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \mathbf{w} = \left(\mathbf{X}^{\text{train}^\top} \mathbf{X}^{\text{train}} \right)^{-1} \mathbf{X}^{\text{train}^\top} \mathbf{y}^{\text{train}}$$

Capacity, overfitting and underfitting

- A ML algorithm must perform well on **new, previously unseen** inputs—not just on which it was trained
 - This ability is called generalization



- Smaller training error → higher model capacity
 - If the training error is too large, the model is **underfitting** the training set
- Smaller test error or generalization error → higher generalization ability
 - If the training error is very small but the test error is very large, the model is **overfitting** the training set

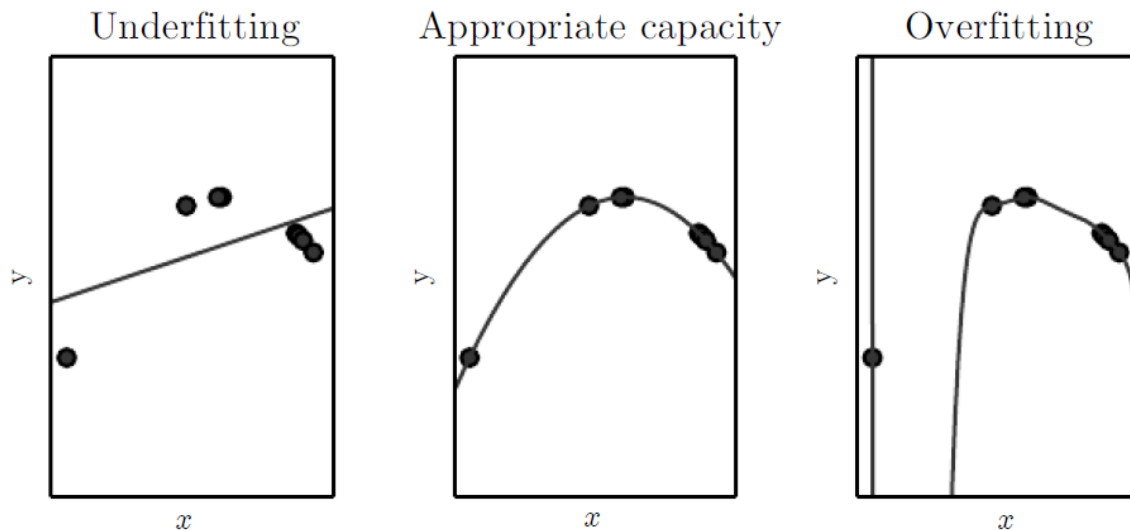
Example: polynomial regression

- Consider a regression problem in which the input x and output y are both scalars. Find a function $f: \mathbb{R} \rightarrow \mathbb{R}$ to fit the data

- $f(x) = b + wx$
- $f(x) = b + w_1x + w_2x^2$
- $f(x) = b + \sum_{i=1}^9 w_i x^i$

MSE training:

$$\min_w \frac{1}{N} \sum_{n=1}^N \|f(x^{(n)}) - y^{(n)}\|_2^2$$



General principles

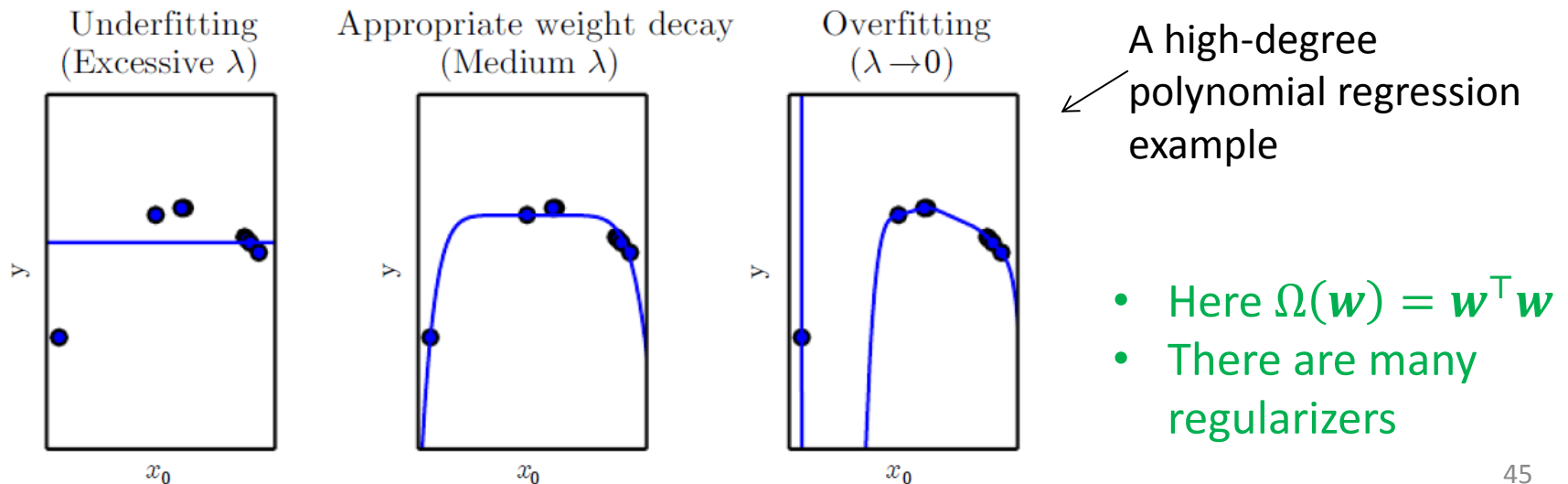
- Increase the model capacity
 - Make the training error small
- Increase the generalization ability
 - Make the gap between training error and test error small

Regularization

- To carry out a specific task, we often build a set of **preferences** into the learning algorithm, which is embodied by a **regularizer** Ω
- E.g., for polynomial regression, the total cost function becomes

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w} \leftarrow \text{Weight decay}$$

where $\lambda > 0$ is a constant.



Regularization

Regularization is any modification we make to a learning algorithm that is intended to reduce its **generalization error** but not its training error

Hyperparameters

- Many machine learning algorithms have two sets of parameters:
 - **Hyperparameters**: control the algorithm's behavior and are not adapted by the algorithm itself. They often determine the **capacity** of the model
 - **Learnable parameters** ("learnable" is often omitted): can be learned from data
- The polynomial regression algorithm $J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$
 - Hyperparameters: λ
 - Learnable parameters: \mathbf{w}
- A neural network
 - Hyperparameters: the number of layers, the number of neurons per layer, etc.
 - Learnable parameters: weights and bias in each layer

Question

- How to choose the hyperparameters considering that we cannot see the test set?

Maximum likelihood estimation

Problem definition

- Given a set of N examples $\mathbb{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$ drawn independently from the true but unknown data-generating distribution $p_{\text{data}}(\mathbf{x})$
- Find a prob distribution $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$ to approximate $p_{\text{data}}(\mathbf{x})$
- The task is to find optimal $\boldsymbol{\theta}$

- The maximum likelihood estimator for $\boldsymbol{\theta}$ is defined as

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^N p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- We usually use

$$\begin{aligned} \boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \end{aligned}$$

Log-likelihood

- where \hat{p}_{data} is the empirical distribution

Conditional log-likelihood

- Estimate a conditional probability $P(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$ in order to predict \mathbf{y} given \mathbf{x}
 - E.g. For classification, \mathbf{y} is a (discrete) random variable representing label of an input \mathbf{x}
- If \mathbf{X} represents all inputs and \mathbf{Y} all observed targets, then the **conditional maximum likelihood estimator** is

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P_{\text{model}}(\mathbf{Y}|\mathbf{X}; \boldsymbol{\theta})$$

- If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log P_{\text{model}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Stochastic gradient decent



- Minimizing the cost function over the entire training set is computationally expensive
- We often decompose the training set into smaller subsets or **minibatches** and optimize the cost function defined over individual minibatches $(\mathbf{X}^{(i)}, \mathbf{y}^{(i)})$ and take the average

$$J(\boldsymbol{\theta}) = \frac{1}{N'} \sum_{i=1}^{N'} L(\mathbf{X}^{(i)}, \mathbf{y}^{(i)}, \boldsymbol{\theta})$$

$$\mathbf{g} = \frac{1}{N'} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{N'} L(\mathbf{X}^{(i)}, \mathbf{y}^{(i)}, \boldsymbol{\theta})$$

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \mathbf{g}$$

- A total of N' minibatches
- The batchsize ranges from 1 to a few hundreds

Summary so far

- Math basics
 - Linear algebra
 - Probability theory
 - Optimization
- Machine learning basics
 - Learning algorithms, task T , performance P and experience E
 - Capacity versus generalization
 - Maximum likelihood estimation
 - SGD

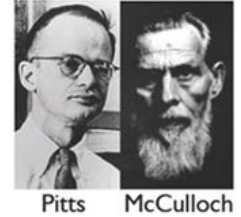
Outline

- Math basics
- Machine learning basics
- Neural networks in the early stage

Neural networks in the early stage

MCCULLOCH AND PITTS NEURON

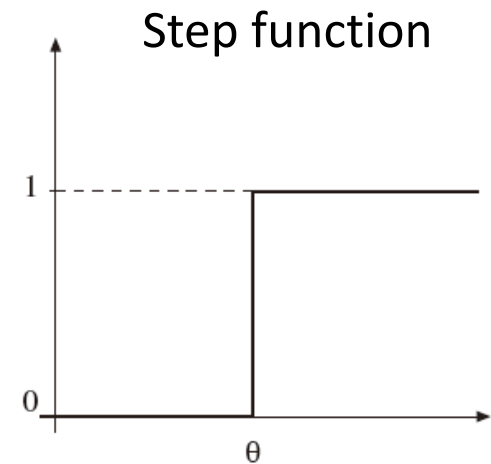
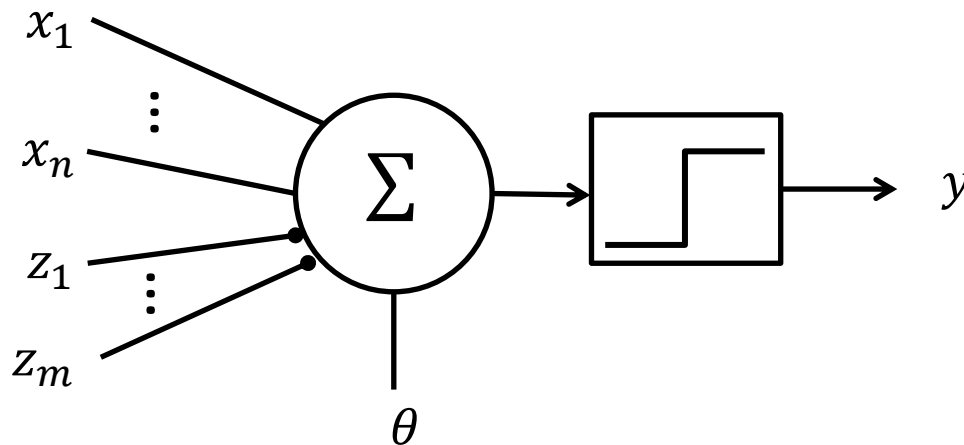
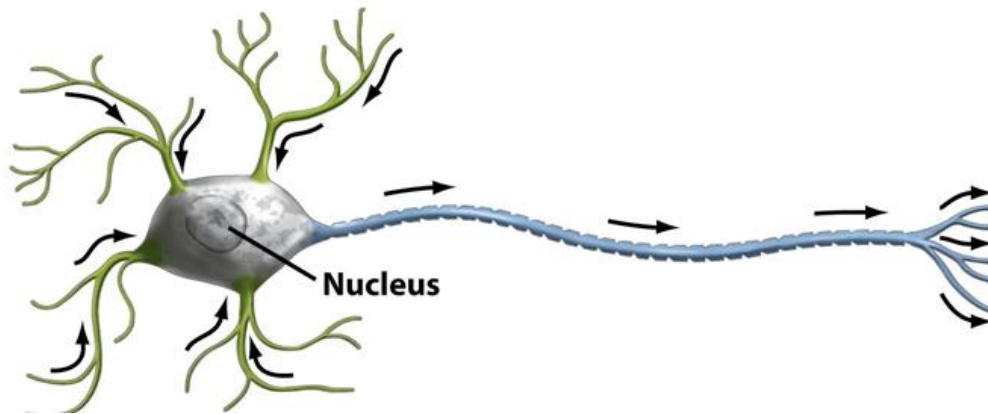
Threshold Logic Unit (TLU)



- The first artificial neuron proposed by Warren McCulloch and Walter Pitts in 1943

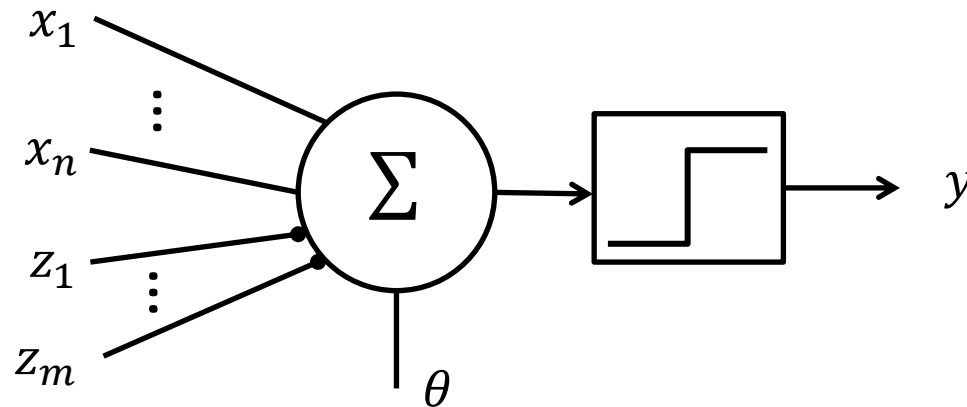
A Logical Calculus of Ideas Immanent in Nervous Activity,
Bulletin of Mathematical Biophysics, 1943
- The model was specifically targeted as a computational model of the “nerve net” in the brain
- Pitts believed that “the brain computing information digital neuron by digital neuron using the exacting implement of mathematical logic”, but the experiments on frog’s eyes opposed this belief
- Pitts burned his unpublished PhD dissertation

Threshold Logic Unit (TLU)



- Excitatory input x_i
- Inhibitory input z_i
- Binary output y_i
- Threshold θ

McCulloch–Pitts unit (M-P unit)



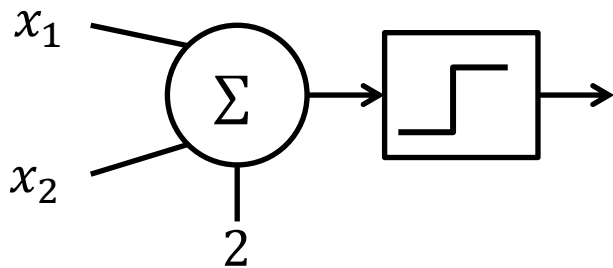
- If **at least** one of z_1, z_2, \dots, z_m is 1, the unit is inhibited and $y = 0$
- Otherwise the total excitation $T = \sum_{i=1}^n x_i$ is computed and compared with the threshold θ of the unit (if $n = 0$ then $x = 0$)
 - If $T \geq \theta$ the unit fires a 1
 - If $T < \theta$ the result is 0.
- The MP unit can be inactivated by a single inhibitory signal, as is the case with some real neurons

Synthesis of Boolean functions

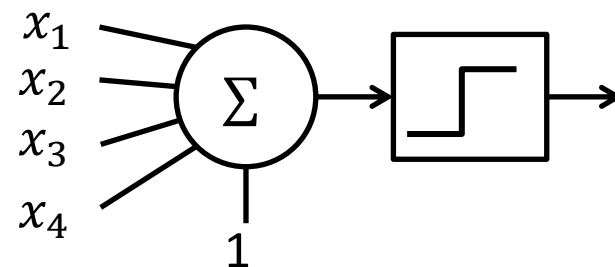
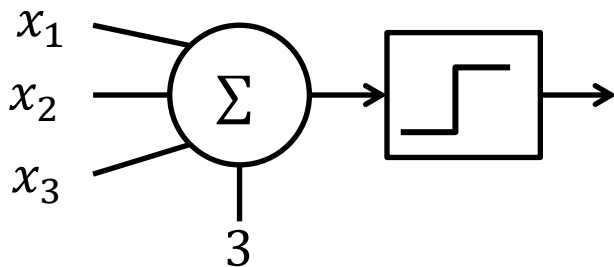
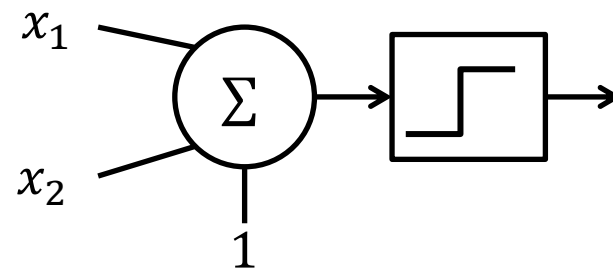
Boolean function: $\{0, 1\}^n \rightarrow \{0, 1\}$

- Conjunction, disjunction, negation

AND



OR



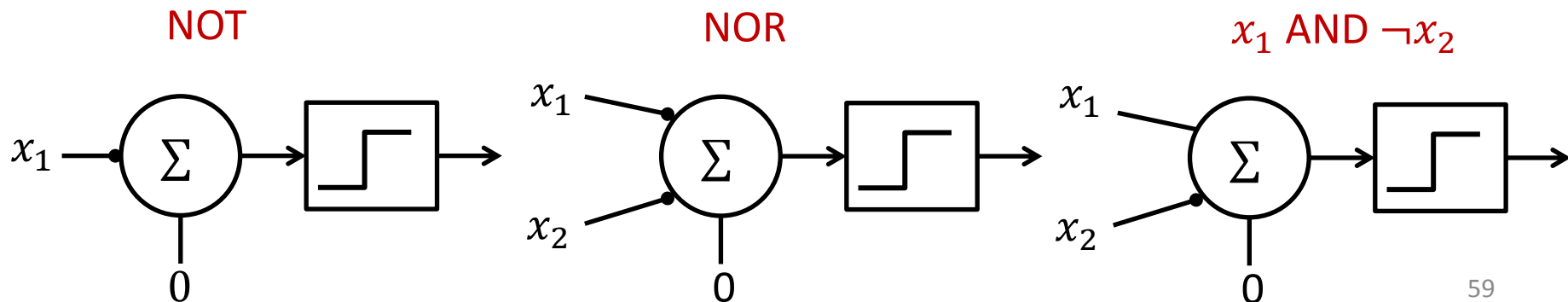
Can you implement negation using the M-P units?

Monotonic logical function

A **monotonic logical function** f of n arguments is one whose value at two given n -dimensional points $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$ is such that $f(\mathbf{x}) \geq f(\mathbf{y})$ whenever the number of ones in the input \mathbf{y} is a subset of the ones in the input \mathbf{x} .

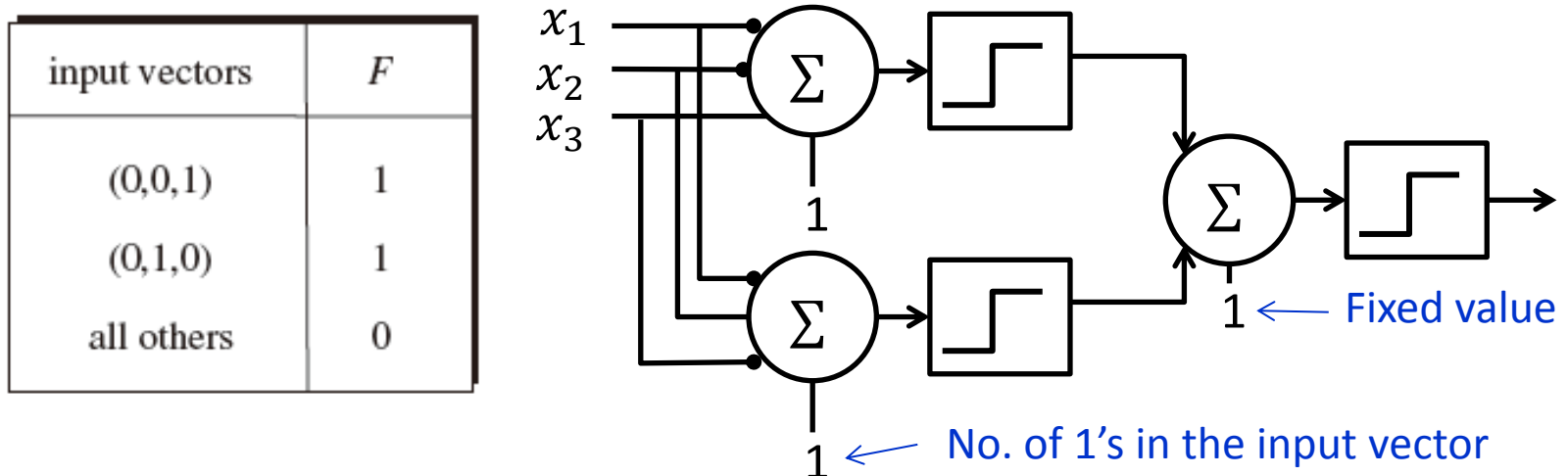
Proposition 1. *Uninhibited threshold logic elements of the McCulloch–Pitts type can only implement monotonic logical functions.*

Proof. See (Rojas: Neural Networks, 1996)



Constructive synthesis

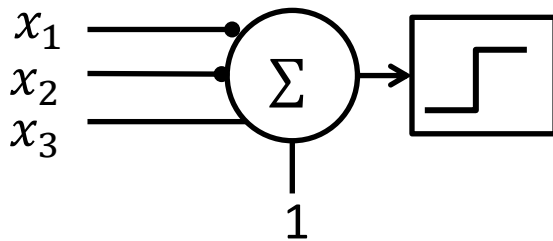
- Every logical function of n variables can be written in tabular form. Suppose there are K rows that have results 1
 - ① Use a M-P unit to represent n values which lead to the result of 1
 - ② Use a disjunction unit to connect the K M-P unit
- For example, let $n = 3$ and the table is as follows



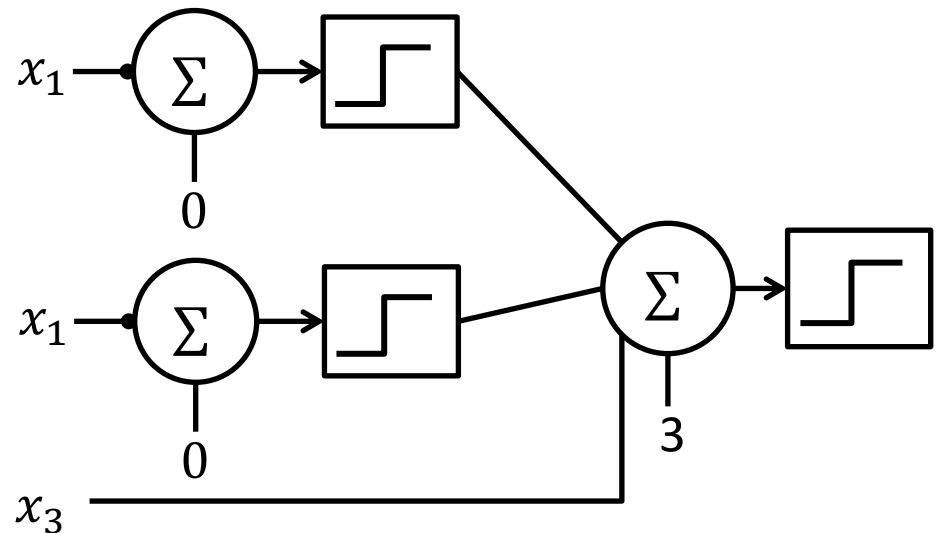
Proposition 2. Any logical function $F : \{0, 1\}^n \rightarrow \{0, 1\}$ can be computed with a M-P network of two layers.

Constructive synthesis

What are **basic components** to construct all logical functions?



Inhibitory connections in the decoders can be replaced with a negation gate



Proposition 3. *All logical functions can be implemented with a network composed of units which exclusively compute the AND, OR, and NOT functions*

Neural networks in the early stage

PERCEPTRON

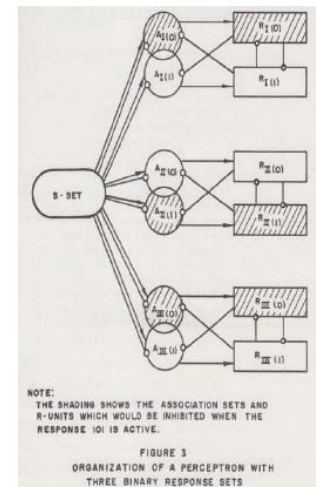
Frank Rosenblatt and Perceptron

- The quest
 - How is information about the physical world sensed by the biological system?
 - In what form is information stored and retrieved?
 - How does remembered information influence recognition and behavior?
- In 1957, Rosenblatt proposed the Perceptron algorithm

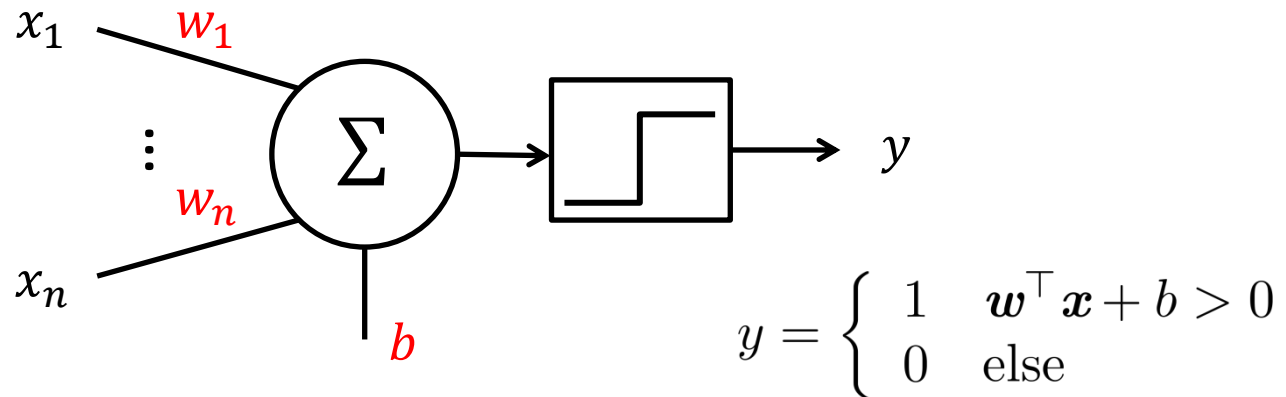
Rosenblatt (1957) The Perceptron--a perceiving and recognizing automaton. Report 85-460-1, Cornell Aeronautical Laboratory.
- Variants, applications and theoretical results were developed in 1960s
- A machine was built



1928 -1971

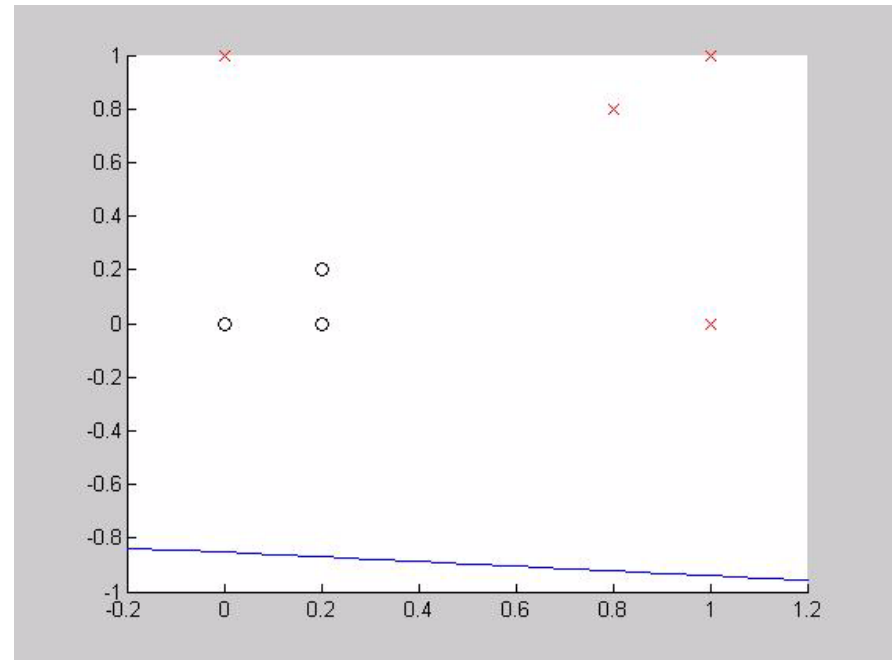
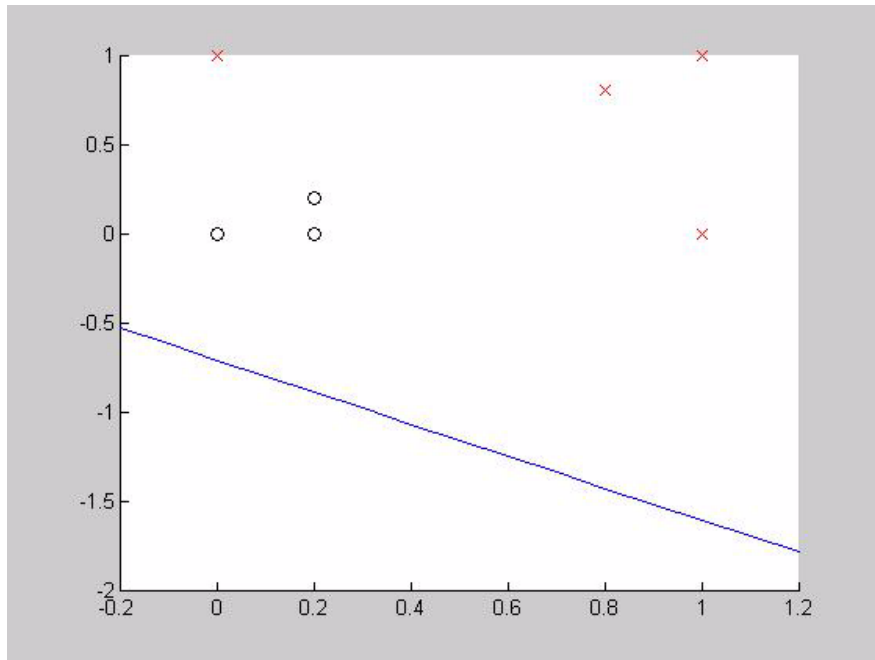


Perceptron



- Add weights to the input connections of the M-P unit
 - Propose a **supervised learning** algorithm: For each data points $\mathbf{x}^{(j)} \in R^m$ and the corresponding labels $t^{(j)}$
 - Calculate the actual output $y^{(j)}$
 - Update the weights: $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta(t^{(j)} - y^{(j)})\mathbf{x}^{(j)}$;
 $b^{\text{new}} = b^{\text{old}} + \eta(t^{(j)} - y^{(j)})$
- where $\eta > 0$ is the learning rate

Example

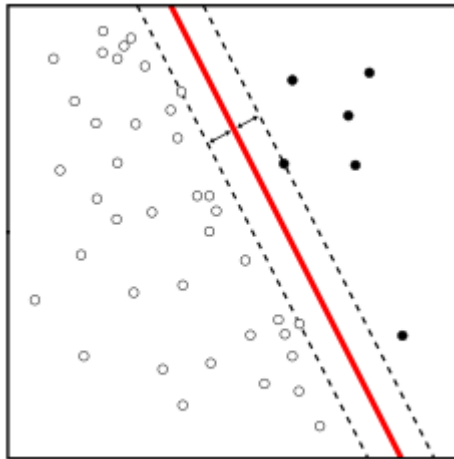


From two different sets of initial weights

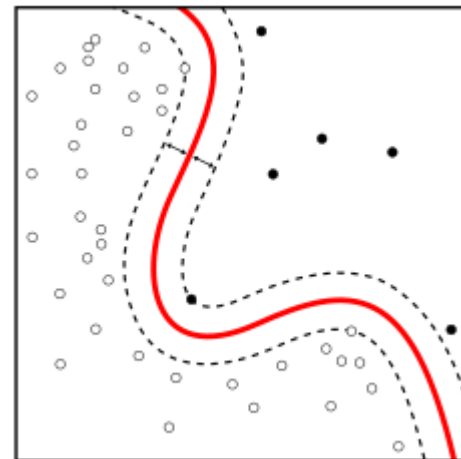
Convergence

Proposition 4: *If the training set is **linearly separable**, then the perceptron is guaranteed to converge. Furthermore, there is an upper bound on the number of times the perceptron will adjust its weights during the training.*

Proof. See (Novikoff, 1962)



linearly separable

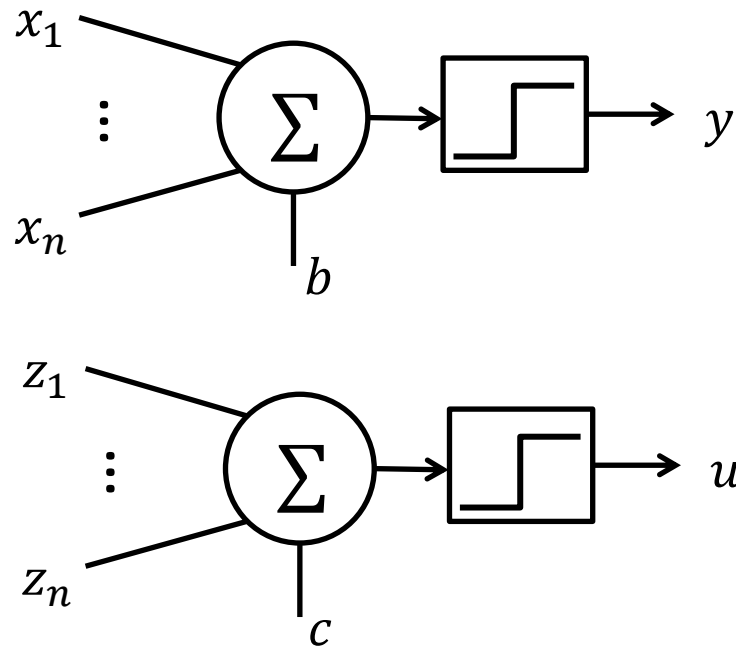


linearly non-separable

Disadvantage: the decision boundary is often close to the training samples, therefore sensitive to noise

Multiple Perceptrons in one layer

- When multiple Perceptrons are combined, each output neuron operates **independently** of all the others; thus, learning each output can be considered in isolation



Neural networks in the early stage

ADALINE

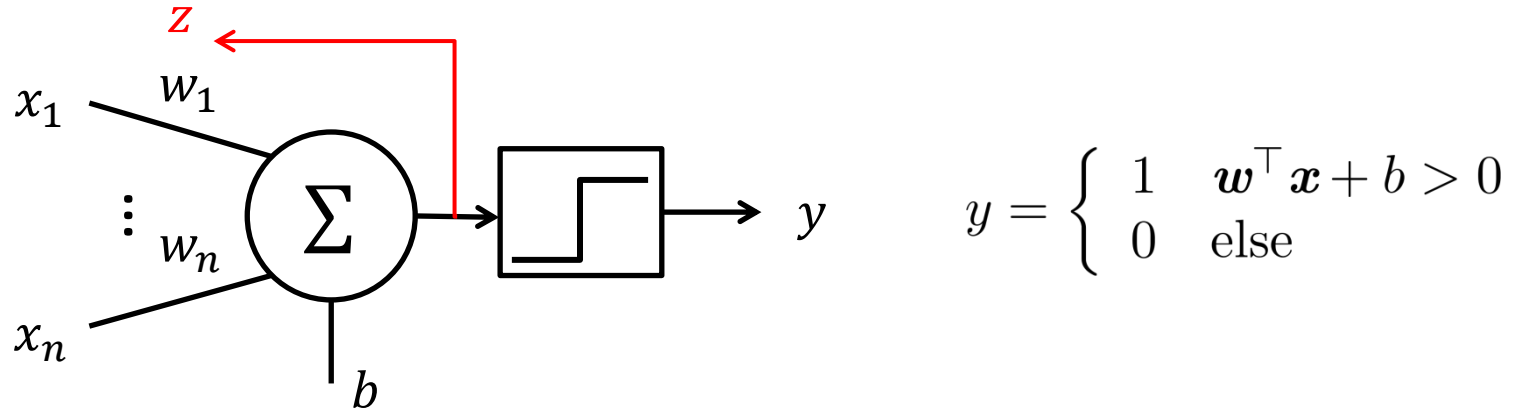
Bernard Widrow and ADALINE

- Widrow and his doctoral student Ted Hoff invented the **least mean squares filter (LMS)** adaptive algorithm in 1960
- The LMS algorithm led to the **ADALINE** and **MADALINE** and to the **backpropagation** technique
- ADALINE (*Adaptive Linear Neuron* or *Adaptive Linear Element*): A single-layer model
- LMS algorithm minimizes the mean squared error (MSE), and is a **stochastic gradient descent (SGD)** method
 - It was proposed for signal processing and achieved great success in that field, but not so successful in training multi-layer neural networks
- In early 1960 Widrow turned to study signal processing, and returned to neural networks in 1980s



Born in 1929

ADALINE

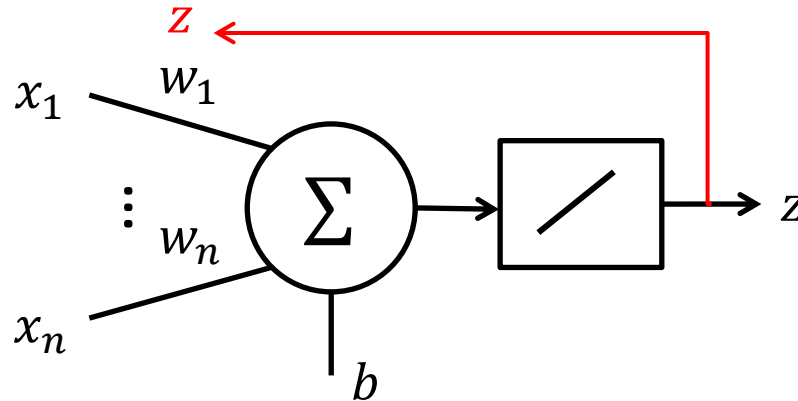


- Same architecture as Perceptron; different training algorithm
 - $z = \mathbf{w}^\top \mathbf{x} + b$ instead of y is used to adjust the weights and bias
- Minimize MSE $E = \sum_j (t^{(j)} - z^{(j)})^2$. The learning algorithm:
$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta (t^{(j)} - z^{(j)}) \mathbf{x}^{(j)}$$
$$b^{\text{new}} = b^{\text{old}} + \eta (t^{(j)} - z^{(j)})$$

where $\eta > 0$ is the learning rate

- Different names: LMS rule, Delta rule, Widrow-Hoff rule, actually SGD

Another view

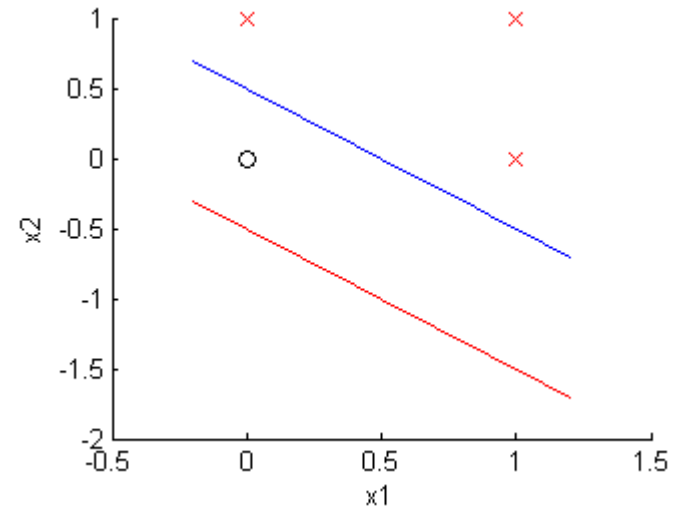
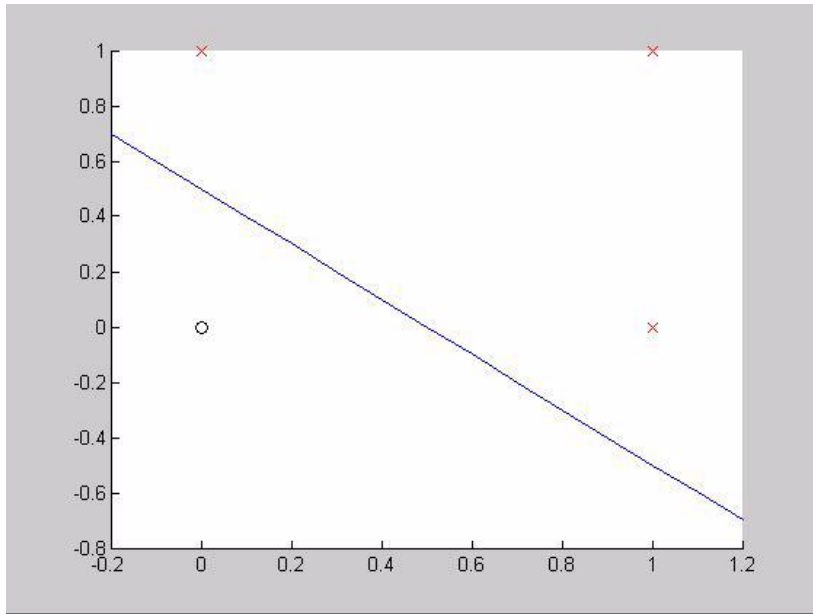


$$z = \mathbf{w}^\top \mathbf{x} + b$$

$$y = \begin{cases} 1 & z > 0 \\ 0 & \text{else} \end{cases}$$

- There is a linear activation function for the variable z
 - This is where the name *Adaptive **Linear** Neuron* comes
- The step function is only used for output y and the output is not involved in the learning process

Example



Blue: $\mathbf{w} = [1 \ 1], b = -0.5$

Red (optimal): $\mathbf{w} = [0.5 \ 0.5], b = 0.25$

Loss for blue: 0.25

Loss for Red: 0.0625

What's the problem?

Doesn't the algorithm converge to the minimum of the MSE?

MADALINE

- MADALINE: Many ADALINE
- A 3-layer (input, hidden, output), fully connected, feed-forward architecture for classification that uses ADALINE units in its hidden and output layers
- Three different training algorithms for MADALINE networks were proposed
 - The first dates back to 1962 and cannot adapt the weights of the hidden-output connection
 - The second training algorithm improved on Rule I and was described in 1988
 - The third "Rule" applied to a modified network with [sigmoid](#) activations instead of step function; it was later found to be equivalent to backpropagation

[See Wikipedia and references therein](#)

Summary

- Math basics
 - Linear algebra
 - Probability theory
 - Optimization
- Machine learning basics
 - Learning algorithms, task T , performance P and experience E
 - Capacity versus generalization
 - Maximum likelihood estimation
 - SGD
- Neural networks in the early stage
 - MP unit
 - Logic unit
 - Perceptron
 - Binary classification, convergence
 - Adaline
 - Widrow-Hoff rule

References

- Chapters 2-5 in [Deep Learning](#) by Goodfellow, Bengio and Courville, 2016, MIT Press
- MP unit material on Web Learning