

Artificial Intelligence

Machine Learning

Nguyễn Văn Diêu

HO CHI MINH CITY UNIVERSITY OF TRANSPORT

2025

Kiến thức - Kỹ năng - Sáng tạo - Hội nhập

Sứ mệnh - Tầm nhìn

Triết lý Giáo dục - Giá trị cốt lõi

Outline I

① ML Overview

- 1.1 Definition
- 1.2 Supervised Learning
- 1.3 Unsupervised Learning
- 1.4 Reinforcement Learning
- 1.5 ML Workflow

② Regression

- 2.1 Linear Regression Model
- 2.2 Gradient descent algorithm
- 2.3 Polynomial Regression
- 2.4 Logistic Regression

③ Classification

- 3.1 Binary Classification
- 3.2 Naïve Bayes Classifier

Outline II

- 3.3 Review Probability Distribution
- 3.4 Gaussian Naïve Bayes
- 3.5 Multinomial Naïve Bayes
- 3.6 Bernoulli Naïve Bayes

4 Decision Tree

- 4.1 Decision Tree Overview
- 4.2 Decision Tree Model
- 4.3 Regression Decision Tree algo.
- 4.4 Classification Decision Tree algo.

5 Support Vector Machine

- 5.1 Equation of a Line Summary
- 5.2 SVM Overview
- 5.3 SVM, Linear, Binary class

Outline III

- 5.4 Soft Margin
- 5.5 Kernel Trick
- 5.6 Dual Solution: Gradient Ascent

⑥ Neural Network

- 6.1 Rosenblatt's Perceptron
- 6.2 Perceptron Training Rule
- 6.3 Perceptron learning AND gate
- 6.4 XOR Problem
- 6.5 Multilayer Networks
- 6.6 Backward Propagation of Errors

⑦ Ensemble Learning

- 7.1 Wisdom of the Crowd
- 7.2 Majority Voting

7.3 Random Forest

ML Definition

Traditional Programming:

Refers to any manually created program that uses input data and runs on a computer to produce the output.



Machine Learning:

Input and output data are fed to an algorithm to create a program. This program can be used to predict future outcomes.



ML Definition

Definition of **Tom Mitchell** (Machine Learning 1997):

*"A computer program is said to **learn** from experience **E** with respect to some class of tasks **T** and performance measure **P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**."*

Checkers learning problem:

- Task **T**: Playing checkers.
- Performance measure **P**: Percent of games won against opponents.
- Training experience **E**: Playing practice games against itself.

ML Definition

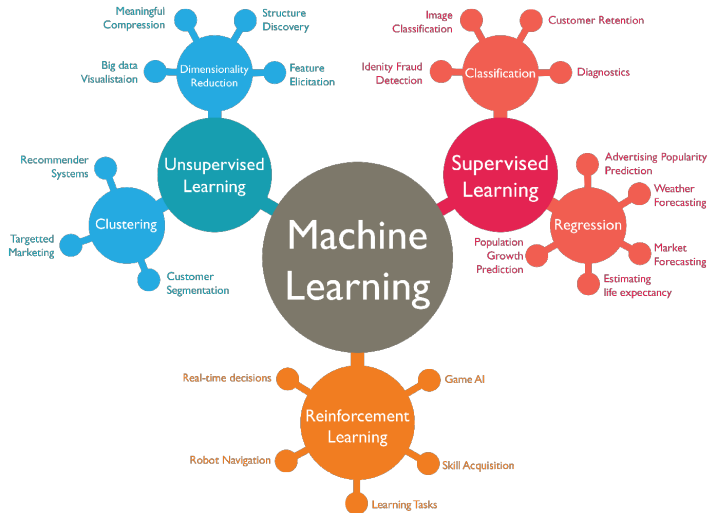
Handwriting recognition learning problem:

- Task **T**: Recognizing and classifying handwritten words within images.
- Performance measure **P**: Percent of words correctly classified.
- Training experience **E**: Database of handwritten words with given classifications.

Robot driving learning problem:

- Task **T**: Driving on public four-lane highways using vision sensors.
- Performance measure **P**: Average distance traveled before an error (as judged by human overseer).
- Training experience **E**: Sequence of images and steering commands recorded while observing a human driver. .

Type of Machine learning



Supervised Learning

Training data

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}) \dots (x^{(i)}, y^{(i)}) \dots (x^{(n)}, y^{(n)})\}$$

$$x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R} : \text{label.}$$

sample $x^{(1)} \rightarrow$	$(x_1^{(1)} \dots x_j^{(1)} \dots x_d^{(1)})$	$y^{(1)} \leftarrow \text{label}$
sample $x^{(i)} \rightarrow$	$(x_1^{(i)} \dots x_j^{(i)} \dots x_d^{(i)})$	$y^{(i)} \leftarrow \text{label}$
...
sample $x^{(n)} \rightarrow$	$(x_1^{(n)} \dots x_j^{(n)} \dots x_d^{(n)})$	$y^{(n)} \leftarrow \text{label}$

fruit	length	width	weight	label
fruit 1	165	38	172	Banana
fruit 2	218	39	230	Banana
fruit 3	76	80	145	Orange
fruit 4	145	35	150	Banana
fruit 5	90	88	160	Orange
...				
fruit n

Supervised vs. Unsupervised

fruit	length	width	weight	label
fruit 1	165	38	172	Banana
fruit 2	218	39	230	Banana
fruit 3	76	80	145	Orange
fruit 4	145	35	150	Banana
fruit 5	90	88	160	Orange
...				
fruit n

1. Unsupervised learning

Learning a model from unlabeled data.

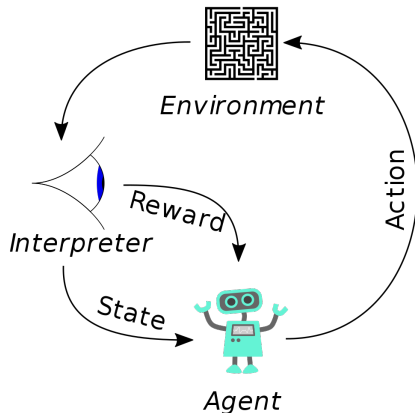
2. Supervised learning

Learning a model from labeled data.

Reinforcement Learning

3. Reinforcement learning

It's concerned with how **intelligent agents** ought to take actions in an **environment** in order to maximize the **cumulative reward**.



ML Workflow

1. Get Data

This process depends on project and data type.

2. Clean, Prepare & Manipulate Data

Real-world data often has unorganized, missing, or noisy elements. Need to clean, prepare, and manipulate the data.

After getting the data, need to convert the data sets into valid formats for ML platform.

Finally, split data into training and test data sets. The typical default is a 70/30 split between training and test sets.

3. Train Model

Train data set connects to an algorithm, this mathematical modeling to learn and develop predictions.

4. Test Model

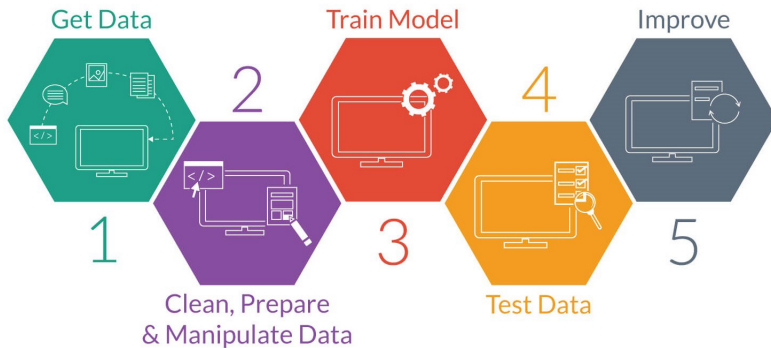
it's time to validate your trained model. Using the test data, check the model's accuracy.

5. Improve

Refine model and improve accuracy:

- Review model's results. Are there other data elements worth adding to your model to make it more accurate?
- Reconsider algorithm choice. A different algorithm may perform better.
- Adjust the parameters of chosen algorithm to improve performance.

ML Workflow



Linear Regression: history

- A very popular technique.
- Rooted in Statistics.
- Method of Least Squares used as early as 1795 by Gauss.
- Re-invented in 1805 by Legendre.
- Frequently applied in astronomy to study the large scale of the universe.
- Still a very useful tool today.

Linear Regression

Training data

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}) \dots (x^{(i)}, y^{(i)}) \dots (x^{(n)}, y^{(n)})\}$$

$$x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R} : \text{label.}$$

sample $x^{(1)} \rightarrow$	$(x_1^{(1)} \dots x_j^{(1)} \dots x_d^{(1)})$	$y^{(1)} \leftarrow \text{label}$
sample $x^{(i)} \rightarrow$	$(x_1^{(i)} \dots x_j^{(i)} \dots x_d^{(i)})$	$y^{(i)} \leftarrow \text{label}$
...
sample $x^{(n)} \rightarrow$	$(x_1^{(n)} \dots x_j^{(n)} \dots x_d^{(n)})$	$y^{(n)} \leftarrow \text{label}$

Task Learn a regression function

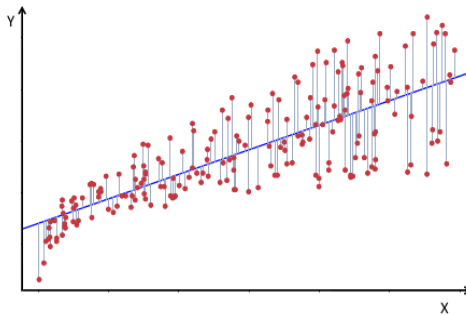
$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$f(x) = y$$

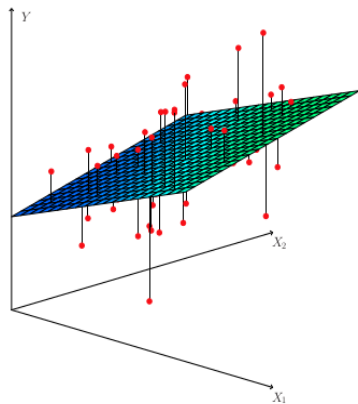
Linear Regression

A regression model is said to be linear if it is represented by a linear function.

Linear Regression



$d = 1$ **Line in \mathbb{R}^2**



$d = 2$ **Hyperplane in \mathbb{R}^3**

Linear Regression

Linear Regression Model

$$\hat{y} = f(x) = \theta_0 + \sum_{j=1}^d \theta_j x_j \quad \text{with } \theta_j \in \mathbb{R}, j \in \{1, 2, \dots, d\}$$

θ : parameters or coefficients or weights.

Learning the linear model \rightarrow learning the θ

Estimation with Least squares

Least square loss

$$loss(y^{(i)}, f(x^{(i)})) = (y^{(i)} - f(x^{(i)}))^2 = (y^{(i)} - \hat{y}^{(i)})^2$$

Minimize the loss over all samples

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

Univariate Linear Regression

Model with one feature $d = 1$

$$\hat{y} = f(x) = \theta_0 + \theta_1 x$$

Loss function

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 \quad , \quad \mathcal{L}(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \theta_1 x^{(i)})^2$$

Minimize $\mathcal{L}(\theta_0, \theta_1) \rightarrow$ **find** θ_0 , θ_1

$$\underset{\{\theta_0, \theta_1\}}{\operatorname{argmin}} \mathcal{L}(\theta_0, \theta_1) = \underset{\{\theta_0, \theta_1\}}{\operatorname{argmin}} \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \theta_1 x^{(i)})^2$$

Find θ_0 , θ_1

$$\frac{\partial \mathcal{L}}{\partial \theta_0} = 0 \quad \frac{\partial \mathcal{L}}{\partial \theta_1} = 0$$

$$\underset{\{\theta_0, \theta_1\}}{\operatorname{argmin}} \mathcal{L}(\theta_0, \theta_1)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_0} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \theta_1 x^{(i)}) \times \frac{\partial}{\partial \theta_0} (y^{(i)} - \theta_0 - \theta_1 x^{(i)})$$

$$\frac{\partial \mathcal{L}}{\partial \theta_0} = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \theta_1 x^{(i)}) \times (-1) = 0$$

$$\theta_0 = \frac{1}{n} \sum_{i=1}^n y^{(i)} - \theta_1 \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

$$\underset{\{\theta_0, \theta_1\}}{\operatorname{argmin}} \mathcal{L}(\theta_0, \theta_1)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \theta_1 x^{(i)}) \times \frac{\partial}{\partial \theta_1} (y^{(i)} - \theta_0 - \theta_1 x^{(i)})$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \theta_1 x^{(i)}) \times (-x^{(i)}) = 0$$

$$\theta_1 \sum_{i=1}^n x^{(i)2} = \sum_{i=1}^n x^{(i)} y^{(i)} - \sum_{i=1}^n \theta_0 x^{(i)}$$

Plugging θ_0 in θ_1

$$\theta_1 = \frac{\sum_{i=1}^n x^{(i)} y^{(i)} - \frac{1}{n} \sum_{i=1}^n x^{(i)} \sum_{i=1}^n y^{(i)}}{\sum_{i=1}^n x^{(i)2} - \frac{1}{n} \sum_{i=1}^n x^{(i)} \sum_{i=1}^n x^{(i)}}$$

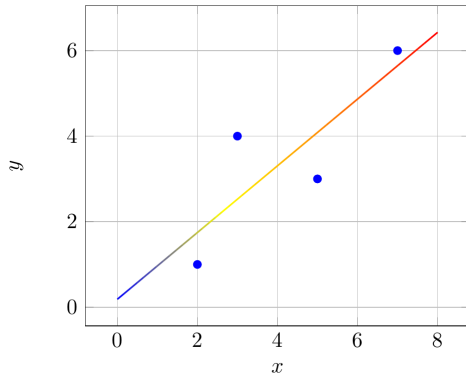
e.g.

- **Training set**

$$\mathcal{D} = \{(2, 1), (3, 4), (5, 3), (7, 6)\}$$

- **Linear Regression Model**

$$f(\theta, x) = \theta_0 + \theta_1 x$$



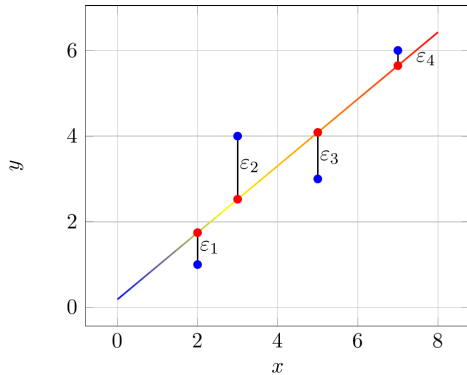
e.g.

$$\epsilon_i = (y_i, f(\theta, x_i))$$

$$\mathcal{L}_i(\theta_0, \theta_1) = \epsilon_i^2$$

$$= (y^{(i)} - f(\theta, x^{(i)}))^2$$

$$= (y^{(i)} - (\theta_0 + \theta_1 x^{(i)}))^2$$



Multivariate Linear Regression

Model with many feature $d > 1$

$$f(x) = \theta_0 + \sum_{j=1}^d \theta_j x_j$$

Minimize

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \theta_0 - \sum_{j=1}^d \theta_j x_j^{(i)})^2$$

Matrix with n sample, d feature

X Matrix $n \times (d + 1)$

y Label vector

Θ Weights vector

$$X = \begin{pmatrix} 1 & x_1^{(1)} & \cdots & x_j^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^{(i)} & \cdots & x_j^{(i)} & \cdots & x_d^{(i)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_j^{(n)} & \cdots & x_d^{(n)} \end{pmatrix}$$
$$y = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(i)} \\ \vdots \\ y^{(n)} \end{pmatrix} \quad \Theta = \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_j \\ \vdots \\ \theta_d \end{pmatrix}$$

Normal Equation

$$\mathcal{L}(\Theta) = \frac{1}{2n} \| (y - X\Theta) \|^2$$

$$\mathcal{L}(\Theta) = \frac{1}{2n} (y - X\Theta)^T (y - X\Theta)$$

$$\frac{\partial \mathcal{L}}{\partial \Theta} = -\frac{1}{n} X^T (y - X\Theta)$$

We have

$$\frac{\partial^2 \mathcal{L}}{\partial \Theta} = \frac{1}{n} X^T X : \text{positive}$$

So that is a minimum when

$$X^T (y - X\Theta) = 0$$

The unique solution

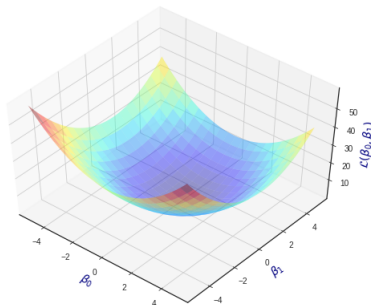
$$\Theta = (X^T X)^{-1} X^T y$$

Gradient descent algo.

$$\mathcal{D} = \{ (x^{(1)}, y^{(1)}) \dots (x^{(i)}, y^{(i)}) \dots (x^{(n)}, y^{(n)}) \}$$
$$x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}.$$

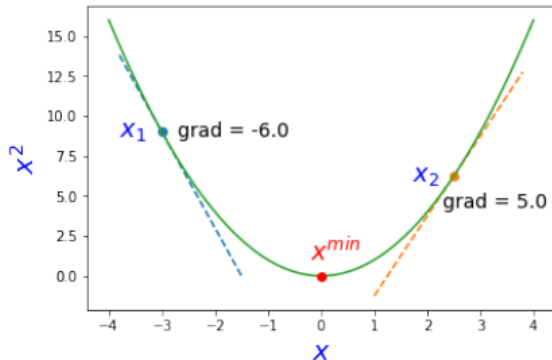
Loss Function $\mathcal{L}(\cdot)$

- n
- $x^{(i)} \in \mathbb{R}^d$
- $\mathcal{L} : k$ order



Gradient Descent

$$f(x) = x^2, \quad f'(x) = 2x$$



Taylor Series Expansion

$$1. \quad f(x + \epsilon) = f(x) + \epsilon f'(x) + \frac{\epsilon^2}{2!} f''(x) + \dots$$

$$2. \quad f(x + \epsilon) \approx f(x) + \epsilon f'(x)$$

$$3. \quad \eta > 0, \quad \epsilon = -\eta f'(x)$$

$$4. \quad f(x - \eta f'(x)) \approx f(x) - \eta f'(x)^2$$

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

Gradient descent algo.

1. x_0
2. $x \leftarrow x_0$
3. η : learning rate
4. ϵ : error convergence
5. while $|f'(x)| \geq \epsilon$ do
 $x \leftarrow x - \eta f'(x)$

Batch Gradient Descent

1. θ_0
2. $\theta \leftarrow \theta_0$
3. η : learning rate
4. ϵ : error convergence
5. while $\|\nabla \mathcal{L}(\theta)\|_2 \geq \epsilon$ do
 $\theta \leftarrow \theta - \eta \nabla \mathcal{L}(\theta)$

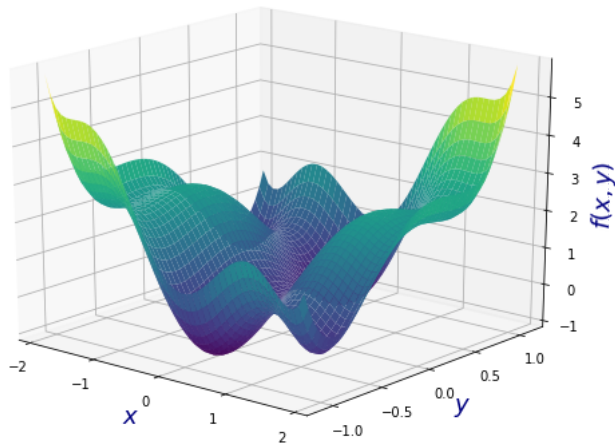
Stochastic Gradient Descent

$$\mathcal{L}(\theta) = \frac{1}{2n} \sum_{i=1}^n (\theta^T X^{(i)} - y^{(i)})^2$$

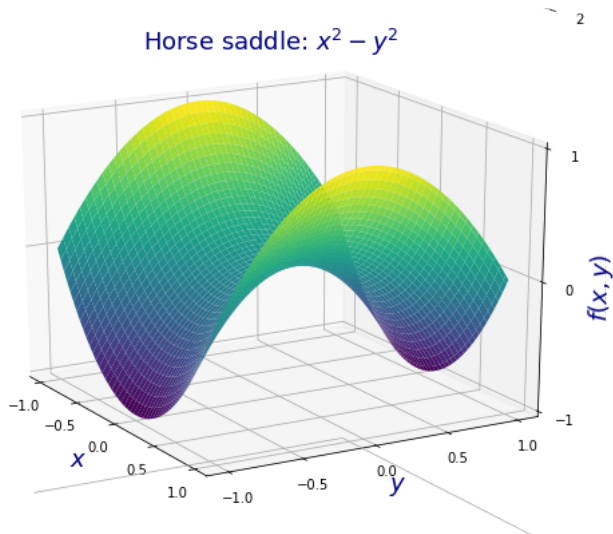
$$\mathcal{L}(\theta, X^{(k)}) = \frac{1}{2} (\theta^T X^{(k)} - y^{(k)})^2$$

Non-Convex Function

Six-hump Camel: $(4 - 2.1x^2 + \frac{x^4}{3})x^2 + xy + (4y^2 - 4)y^2$

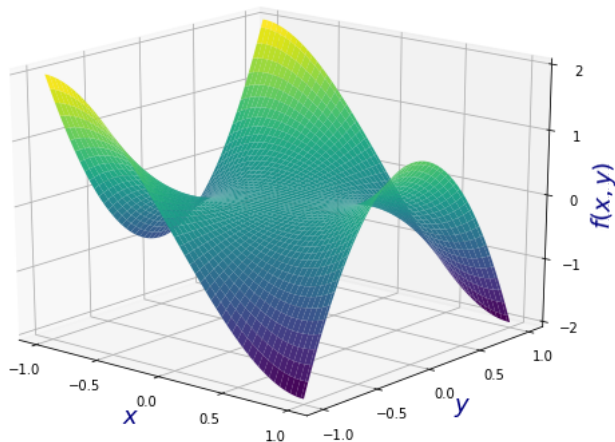


Non-Convex Function

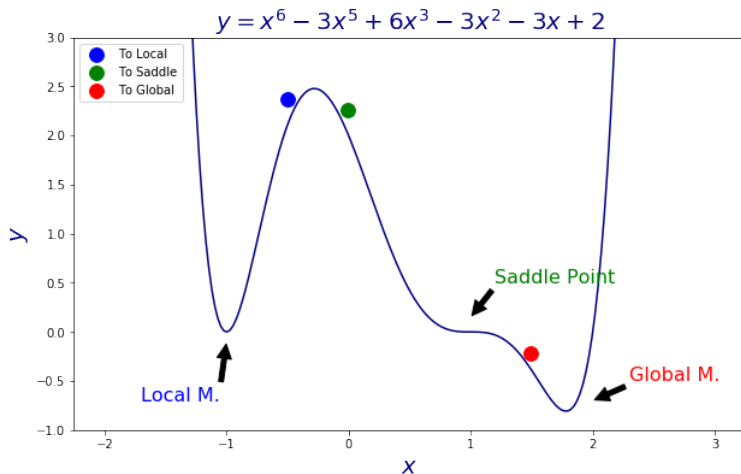


Non-Convex Function

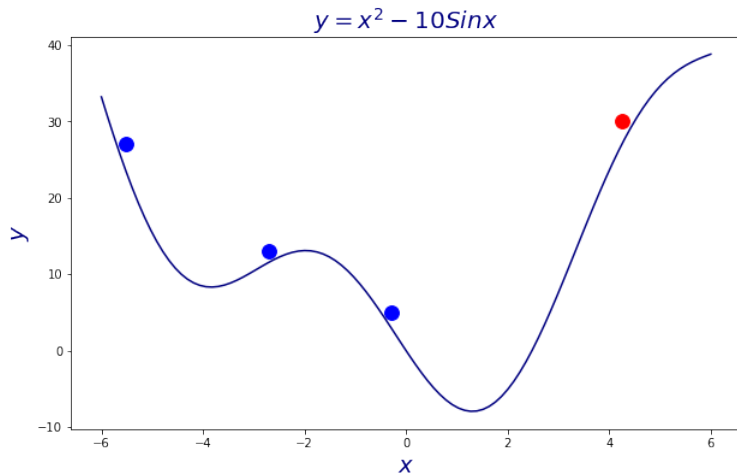
Monkey saddle: $x^3 - 3xy^2$



Non-Convex Function



Momentum



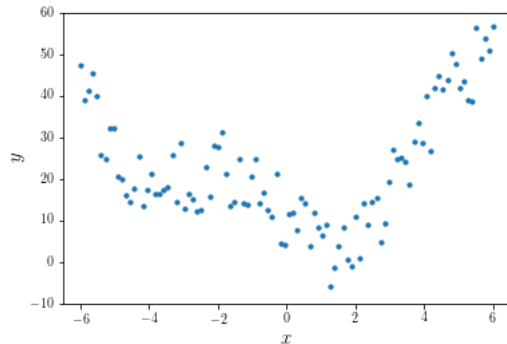
Momentum Algorithm

1. $\theta_0 ; \theta \leftarrow \theta_0$
2. $\eta ; \epsilon$
3. $v_0 ; v \leftarrow v_0$
4. γ : momentum
5. while $\|\nabla \mathcal{L}(\theta)\|_2 \geq \epsilon$ do
 - $v \leftarrow \gamma v + \eta \nabla \mathcal{L}(\theta)$
 - $\theta \leftarrow \theta - v$

Polynomial Regression

Sample with $d = 1$

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}) \dots (x^{(i)}, y^{(i)}) \dots (x^{(n)}, y^{(n)})\}$$



Polynomial Regression

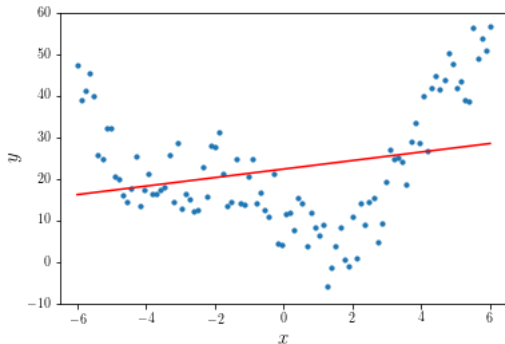
Samples

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}) \dots (x^{(i)}, y^{(i)}) \dots (x^{(n)}, y^{(n)})\}$$

$$x^{(i)} \in \mathbb{R}, y^{(i)} \in \mathbb{R}$$

Simple Linear Regression

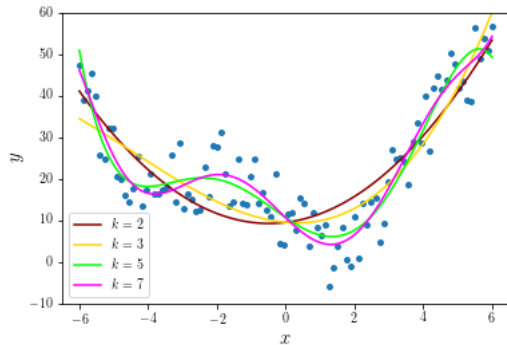
$$f(x) = \theta_0 + \theta_1 x$$



Polynomial Regression

Polynomial k -order

$$f(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_k x^k$$



Polynomial Regression

Solving the problem Polynomial Regression is a model used when the response variable is non-linear.

$$f(x) = \theta_0 + \sum_{i=1}^k \theta_i x^i$$

Convert to Multivariate Linear Regression

with $d = k$ feature.

$$x_i = x^i$$

$$f(x) = \theta_0 + \sum_{i=1}^k \theta_i x_i$$

Multivariate Polynomial Regression

So far we have

- Simple Linear Regression $f(x) = \theta_0 + \theta_1 x$
- Multivariate Linear Regression $f(x) = \theta_0 + \sum_{i=1}^d \theta_i x_i$
- Polynomial Regression $f(x) = \theta_0 + \sum_{i=1}^k \theta_i x_i^i$

Multivariate Polynomial Regression

- Second Order

$$f(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_{11} x_1^2 + \theta_{22} x_2^2 + \theta_{12} x_1 x_2$$

- General Order Big issue.

To solve this issue, it can be mapped to a higher order space of independent variables called as the feature space (eg. Kernel method ...)

Logistic Regression for Classification

Logistic model (Logit model): Probability of a class label in dataset.

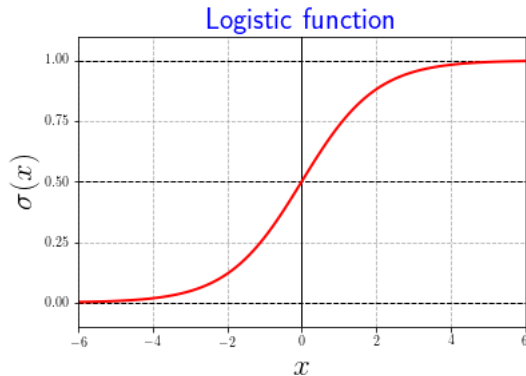
Logistic function

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

$$x \rightarrow +\infty, \sigma(x) \rightarrow 1$$

$$x \rightarrow -\infty, \sigma(x) \rightarrow 0$$

– *Continuous, has a first derivative.*



Probability return

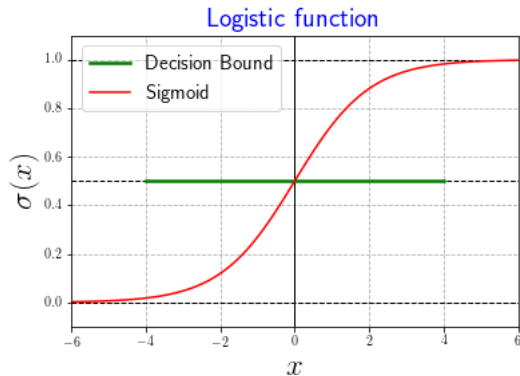
Logistic return probability score between 0 and 1

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

$$Prob \in (0, 1)$$

$$Prob \geq 0.5, \text{ class} = 1$$

$$Prob < 0.5, \text{ class} = 0$$



Logistic Regression Model

$$\mathcal{D} = \left\{ (x^{(i)}, y^{(i)}) \right\}_1^n, \quad x^{(i)} \in \mathbb{R}^d, \quad y^{(i)} \in \{0, 1\}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \in \mathbb{R}^{d+1} \quad x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^{d+1}$$

Model

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

dot product $\theta \cdot x \equiv \sum_j \theta_j x_j$

Logistic Regression Model

e.g.

$x = (x_1, x_2)$, $y = \{0, 1\}$. imagine we know θ .

$$z = \theta \cdot x = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

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

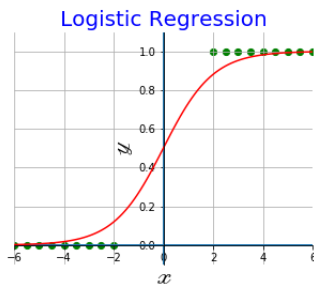
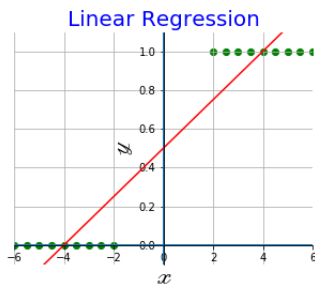
Decision boundary = .5

e.g. $Prob(class = 1)$ return .4, only have 40% chance "class 1" or this observation as "class 0".

Linear Regression Vs. Logistic Regression

- Linear Regression: *Continuous output.*
- Logistic Regression: *Constant output.*
- Linear Regression: *Using Ordinary Least Squares (OLS).*
- Logistic Regression: *Using Maximum Likelihood Estimation (MLE).*

Consider dataset which two class $\{0, 1\}$



Logistic Regression Model

Model

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

$$f(x) \in (0, 1)$$

$f(x)$: Probability.

Probability return by Model

$$Prob(y = 1 \mid x; \theta) = f(x)$$

$$Prob(y = 0 \mid x; \theta) = 1 - f(x)$$

Model base on probability

$$Prob(y \mid x; \theta) = f(x)^y (1 - f(x))^{1-y}$$

The Likelihood

One sample i

$$Prob(y^{(i)} | x^{(i)}; \theta) = f(x^{(i)})^{y^{(i)}} (1 - f(x^{(i)}))^{1-y^{(i)}}$$

Loss function $Loss(\theta)$ for all sample.

n training samples were generated independently.

Likelihood

$$Loss(\theta) = \prod_{i=1}^n Prob(y^{(i)} | x^{(i)}; \theta)$$

$$Loss(\theta) = \prod_{i=1}^n f(x^{(i)})^{y^{(i)}} (1 - f(x^{(i)}))^{1-y^{(i)}}$$

Logarithm of Likelihood

- Logarithm turns a product into a sum.
- It avoid the issue of small number(typically for probability).

$$\mathcal{L}(\theta) = \log \text{Loss}(\theta)$$

$$= \log \prod_{i=1}^n f(x^{(i)})^{y^{(i)}} (1 - f(x^{(i)}))^{1-y^{(i)}}$$

$$= \sum_{i=1}^n \log \left\{ f(x^{(i)})^{y^{(i)}} (1 - f(x^{(i)}))^{1-y^{(i)}} \right\}$$

$$= \sum_{i=1}^n \left\{ \log f(x^{(i)})^{y^{(i)}} + \log(1 - f(x^{(i)}))^{1-y^{(i)}} \right\}$$

$$= \sum_{i=1}^n \left\{ y^{(i)} \log f(x^{(i)}) + (1 - y^{(i)}) \log(1 - f(x^{(i)})) \right\}$$

Maximization into a Minimization

Maximize the Likelihood

$$\mathcal{L}(\theta) = \sum_{i=1}^n \left\{ y^{(i)} \log f(x^{(i)}) + (1 - y^{(i)}) \log(1 - f(x^{(i)})) \right\}$$

- Negative log-likelihood (NLL).
- Use gradient descent algo.

So we minimize the negative $\mathcal{L}(\theta)$ with

$$\mathcal{L}(\theta) = - \sum_{i=1}^n \left\{ y^{(i)} \log f(x^{(i)}) + (1 - y^{(i)}) \log(1 - f(x^{(i)})) \right\}$$

Minimize the negative Likelihood

$$\mathcal{L}(\theta) = \sum_{i=1}^n \left\{ -y^{(i)} \log f(x^{(i)}) - (1 - y^{(i)}) \log(1 - f(x^{(i)})) \right\}$$

$$\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta) = \sum_{i=1}^n \left\{ -y^{(i)} \frac{1}{\sigma(\theta \cdot x^{(i)})} + (1 - y^{(i)}) \frac{1}{1 - \sigma(\theta \cdot x^{(i)})} \right\} \\ \times \frac{\partial}{\partial \theta_j} \sigma(\theta \cdot x^{(i)})$$

since $\frac{\partial}{\partial \theta_j} \sigma(\theta \cdot x^{(i)}) = \sigma(\theta \cdot x^{(i)}) (1 - \sigma(\theta \cdot x^{(i)})) \frac{\partial}{\partial \theta_j} \theta \cdot x^{(i)}$, so

$$= \sum_{i=1}^n \left\{ \sigma(\theta \cdot x^{(i)}) - y^{(i)} \right\} \frac{\partial}{\partial \theta_j} \theta \cdot x^{(i)}$$

Since $\theta \cdot x^{(i)} = \theta_0 + \sum_{j=1}^d \theta_j x_j^{(i)} \Rightarrow \frac{\partial}{\partial \theta_j} \theta \cdot x^{(i)} = x_j^{(i)}$, so

$$= \sum_{i=1}^n (f(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Types of Logistic Regression

- **Binary Logistic Regression** Binary class, e.g. Spam or Not Spam, Cancer or No Cancer.
- **Multinomial Logistic Regression** Many class, e.g. predicting the type of Wine.
- **Ordinal Logistic Regression** Many ordinal class, e.g. restaurant or product rating from 1 to 5.

Binary Classification

Classification

$$\mathcal{D} = \left\{ (x^{(i)}, y^{(i)}) \right\}_1^n, \quad x^{(i)} \in \mathbb{R}^d, \quad y^{(i)} : \text{discrete.}$$

Binary Classification

$$y^{(i)} \in \{0, 1\} \text{ or } y^{(i)} \in \{-1, 1\}$$

e.g.

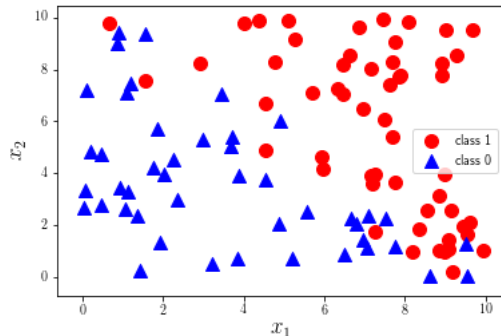
$$\left\{ (x^{(i)}, y^{(i)}) \right\}_1^n, \quad x^{(i)} \in \mathbb{R}^2, \quad y^{(i)} \in \{0, 1\}$$

x_1	x_2	y
8.625	0.058	0
3.828	0.723	0
7.150	3.899	1
6.477	8.198	1
1.922	1.331	0

e.g. Binary Classification

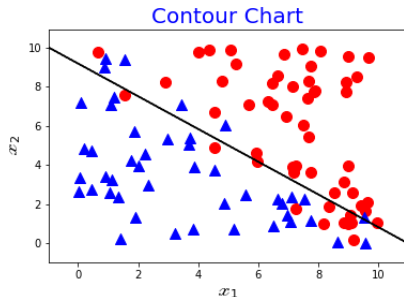
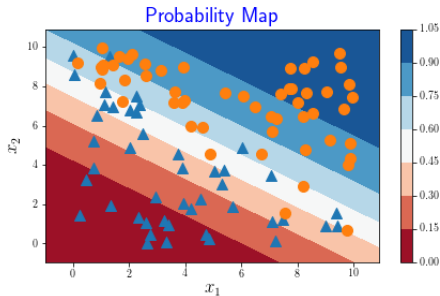
e.g. chart of dataset above

$$\left\{ (x^{(i)}, y^{(i)}) \right\}_1^n, x^{(i)} \in \mathbb{R}^2, y^{(i)} \in \{0, 1\}$$

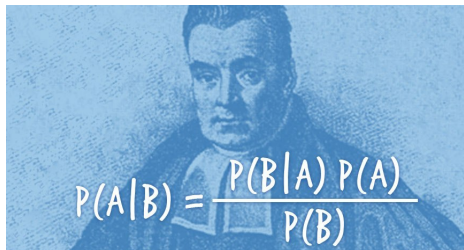


e.g. Logistic Regression

- This model predict probability of two class:
- Contour chart base on probability.



Bayes's Theorem



$p(A)$: *Prior*

$p(A|B)$: *Posterior*

$p(B|A)$: *Likelihood*

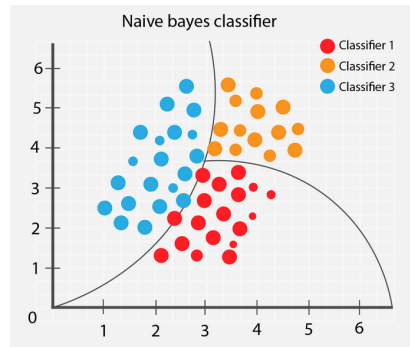
$p(B)$: *Evidence*

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

$$P(\text{class}/\text{data}) = \frac{P(\text{data}/\text{class}) \times P(\text{class})}{P(\text{data})}$$

Naïve Bayes Classifier

$$P(class/data) = \frac{P(data/class) \times P(class)}{P(data)}$$



Review: Joint Probability

Tossing two coins: Independent

- A: Means the first coin lands face up
- B: Means the second coin lands face up
 - $p(A) = p(B) = 0.5$
 - $p(A \text{ and } B) = p(A) p(B) = 0.25$
 - $p(B|A) = p(B)$

Events are not independent

- A: Mean it rains today
- B: Means it rains tomorrow
 - It rained today, it more likely rain tomorrow
 - $p(B|A) > p(B)$
 - $p(A \text{ and } B) = p(A) p(B|A)$

e.g. Cookie problem

- *Suppose there are two bowls of cookies*
 - + Bowl 1:
 - 30 vanilla
 - 10 chocolate
 - + Bowl 2:
 - 20 vanilla
 - 20 chocolate
- *Now suppose you choose*
 - + One of the bowls at random
 - + Without looking, select a cookie at random

This is a conditional probability

$$p(\text{Bowl 1} \mid \text{vanilla})$$

$$\begin{aligned} p(\text{vanilla} \mid \text{Bowl 1}) &= 3/4 \\ &\neq p(\text{Bowl 1} \mid \text{vanilla}) \end{aligned}$$

Bayes's Theorem

Any events A and B

- $p(A \text{ and } B) = p(B \text{ and } A)$
 - $p(A \text{ and } B) = p(A) p(B|A)$
 - $p(B \text{ and } A) = p(B) p(A|B)$
- $$\Rightarrow p(B) p(A|B) = p(A) p(B|A)$$

Bayes's Theorem

$$p(A|B) = \frac{p(B|A) p(A)}{p(B)}$$

Cookie problem

- + B_1 : Hypothesis of cookie came from Bowl 1
- + V : Vanilla cookie

$$p(B_1|V) = \frac{p(V|B_1) p(B_1)}{p(V)}$$

e.g. Cookie problem

$$p(B_1|V) = \frac{p(V|B_1) p(B_1)}{p(V)}$$

$p(B_1)$: Probability chose Bowl 1

$$p(B_1) = 1/2$$

$p(V|B_1)$: Probability vanilla cookie from Bowl 1

$$p(V|B_1) = 3/4$$

$p(V)$: Probability vanilla cookie from either bowl

$$p(V) = 5/8$$

$$p(B_1|V) = \frac{(3/4) (1/2)}{(5/8)} = 3/5$$

e.g. Elderly Fall and Death

- Elderly person is died: 10%
- Elderly people falling: 5%
- All elderly people die, they had fall: 7%

Probability that elderly people die when they fall?

$$P(Die|Fall) = \frac{P(Fall|Die) \times P(Die)}{P(Fall)}$$

$$P(Die) = 0.10$$

$$P(Fall) = 0.05$$

$$P(Fall|Die) = 0.07$$

$$P(Die|Fall) = \frac{0.07 \times 0.10}{0.05}$$

$$P(Die|Fall) = \mathbf{0.14}$$

- If an elderly person falls
- There is a 14% probability that they will die from the fall

e.g. Email and Spam Detection

- Email receive is spam: 2%
- Spam detector accuracy: 99%
- When an email is not spam, it will mark it as spam: 0.1%

Probability that fact spam email in spam folder?

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

$$P(B) = P(B|A) \times P(A) + P(B|not A) \times P(not A)$$

e.g. Email and Spam Detection

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

$$P(B) = P(B|A) \times P(A) + P(B|\text{not } A) \times P(\text{not } A)$$

$$P(A|B) = P(\text{Spam}|\text{Detected}) = ?$$

$$P(B|A) = P(\text{Detected}|\text{Spam}) = 0.99$$

$$P(A) = P(\text{Spam}) = 0.02$$

$$P(\text{not } A) = 1 - P(\text{Spam}) = 0.98$$

$$P(B|\text{not } A) = P(\text{Detected}|\text{not Spam}) = 0.001$$

$$P(\text{Spam}|\text{Detected}) = \frac{0.99 \times 0.02}{0.99 \times 0.02 + 0.001 \times 0.98} = 0.952$$

– Probability fact spam email in spam folder, is **95.2%**.

e.g. Liars and Lie Detectors

- Lie Detector test persons: *if positive result \Rightarrow they are lying.*
- People are tested:
 - + Telling the truth: 98%
 - + Liars: 2%
- Liar people is tested: positive result 72%
- When the machine says they are *not lying*: this is true 97%

Probability that they are indeed lying?

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

$$P(B) = P(B|A) \times P(A) + P(B|\text{not } A) \times P(\text{not } A)$$

e.g. Liars and Lie Detectors

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

$$P(B) = P(B|A) \times P(A) + P(B|\text{not } A) \times P(\text{not } A)$$

$$P(A|B) = P(\text{Lying}|\text{Positive}) = ?$$

$$P(B|A) = P(\text{Positive}|\text{Lying}) = 0.72$$

$$P(A) = P(\text{Lying}) = 0.02$$

$$P(\text{not } A) = 1 - P(\text{Lying}) = 0.98$$

$$P(\text{not } B|\text{not } A) = P(\text{not Positive}|\text{not Lying}) = 0.097$$

$$P(B|\text{not } A) = 1 - P(\text{not } B|\text{not } A) = 1 - 0.097 = 0.03$$

$$P(\text{Lying}|\text{Positive}) = \frac{0.72 \times 0.02}{0.72 \times 0.02 + 0.03 \times 0.98} = 0.328$$

— *Probability fact lying when positive test result, is **32.8%**.*

e.g. Medical test

- People have a certain genetic defect: 1% – Testing to genetic defect (true positives): 90%
- Testing have false positives: 9.6%

Probability genetic defect when get a positive test result?

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B|A) \times P(A) + P(B|\text{not } A) \times P(\text{not } A)}$$

$$P(A|B) = P(\text{GeneticDefect}|\text{Positive}) = ?$$

$$P(B|A) = P(\text{Positive}|\text{GeneticDefect}) = 0.9$$

$$P(A) = P(\text{GeneticDefect}) = 0.01$$

$$P(\text{not } A) = 1 - P(\text{GeneticDefect}) = 0.99$$

$$P(B|\text{not } A) = P(\text{Positive}|\text{not GeneticDefect}) = 0.096$$

$$P(\text{GeneticDefect}|\text{Positive}) = \frac{0.9 \times 0.01}{0.9 \times 0.01 + 0.096 \times 0.99} = 0.0865$$

- *Probability faulty gene on positive result, is **8.65%**.*

e.g. Breast Cancer test

- Women over 50 have breast cancer: 1%
- Women who have breast cancer, had positive result test: 90%
- Women will have false positives: 8%

Probability woman has cancer if she has a positive result?

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B|A) \times P(A) + P(B|not A) \times P(not A)}$$

$$P(A|B) = P(Cancer|Positive) = ?$$

$$P(B|A) = P(Positive|Cancer) = 0.9$$

$$P(A) = P(Cancer) = 0.01$$

$$P(not A) = 1 - P(Cancer) = 0.99$$

$$P(B|not A) = P(Positive|not Cancer) = 0.08$$

$$P(Cancer|Positive) = \frac{0.9 \times 0.01}{0.9 \times 0.01 + 0.08 \times 0.99} = 0.10$$

- *Probability cancer, given a positive test result, is **10%**.*

Naïve Bayes Classifier

Given dataset $D = \{ (x^{(i)}, y^{(i)}) \}_{i=1}^n$, $x^{(i)} \in \mathbb{R}^d$, $y^{(i)} \in C$

given (x, y) , $x \in \mathbb{R}^d$, find $y \in C$, with maximum $p(y|x)$

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

- $p(y)$: Prior probability of class y in dataset D
(we have C class)
- $p(y|x)$: Posterior probability of class y given **one** evidence
 $x = (x_1, x_2, \dots, x_d)$
- $p(x|y)$: Likelihood which is the probability of evidence given
class $y \in C$
- $p(x)$: Prior probability of **one** evidence in D
 $x = (x_1, x_2, \dots, x_d)$

Naïve Bayes Model

$$p(y|x_1, x_2, \dots, x_d) = \frac{p(y) p(x_1, x_2, \dots, x_d|y)}{p(x_1, x_2, \dots, x_d)}$$

(x_1, x_2, \dots, x_d) are stochastically independent, given y :

$$p(x_1, x_2, \dots, x_d|y) = p(x_1|y) p(x_2|y) \dots p(x_d|y)$$

$$p(y|x_1, x_2, \dots, x_d) = \frac{p(y) \prod_{i=1}^d p(x_i|y)}{p(x_1, x_2, \dots, x_d)}$$

$p(x_1, x_2, \dots, x_d)$ is constant given the Dataset,

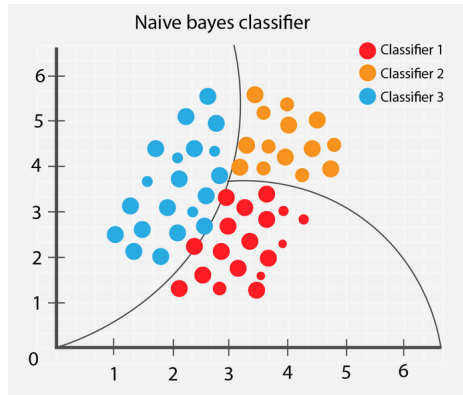
$$p(y|x_1, x_2, \dots, x_d) \propto p(y) \prod_{i=1}^d p(x_i|y)$$

Algorithm

$$\hat{y} = \underset{y \in C}{\operatorname{argmax}} p(y) \prod_{i=1}^d p(x_i|y)$$

Naïve Bayes Classifier algorithm

$$\hat{y} = \underset{y \in C}{\operatorname{argmax}} p(y) \prod_{i=1}^d p(x_i|y)$$



e.g.

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

Give a new instance x :

Outlook	Temperature	Humidity	Windy	Play
<i>sunny</i>	<i>cool</i>	<i>high</i>	<i>true</i>	?

$y = \text{yes}$

$$\begin{aligned}
 p(\text{yes}) &= \frac{9}{14} \\
 p(\text{sunny}|\text{yes}) &= \frac{2}{9} \\
 p(\text{cool}|\text{yes}) &= \frac{3}{9} \\
 p(\text{high}|\text{yes}) &= \frac{3}{9} \\
 p(\text{true}|\text{yes}) &= \frac{3}{9} \\
 p(y = \text{yes}) &= \frac{9}{14} \times \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \\
 &= 0.00529
 \end{aligned}$$

$y = \text{no}$

$$\begin{aligned}
 p(\text{no}) &= \frac{5}{14} \\
 p(\text{sunny}|\text{no}) &= \frac{3}{5} \\
 p(\text{cool}|\text{no}) &= \frac{1}{5} \\
 p(\text{high}|\text{no}) &= \frac{4}{5} \\
 p(\text{true}|\text{no}) &= \frac{1}{5} \\
 p(y = \text{no}) &= \frac{5}{14} \times \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{1}{5} \\
 &= 0.02057
 \end{aligned}$$

$p(y = \text{no}) > p(y = \text{yes})$. Predict result:

Outlook	Temperature	Humidity	Windy	Play
<i>sunny</i>	<i>cool</i>	<i>high</i>	<i>true</i>	no

Probability Distribution

- A probability distribution is a statistical function that describes all the possible values that a random variable can take within a given range.
- Probability distribution depends on factors:

Mean expected value of the distribution.

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

Variance describes the spread of observation from the mean.

$$\text{Population: } \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x - \mu)^2 \quad \text{Sample: } s^2 = \frac{1}{n-1} \sum_{i=1}^n (x - \mu)^2$$

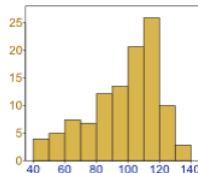
Standard deviation describes the normalized spread of observations from the mean.

$$\sigma$$

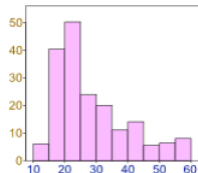
Probability Distribution

Data can be "distributed" (spread out) in different ways

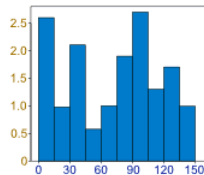
It can be spread out more on the left



Or more on the right



Or it can be all jumbled up



Normal Distribution

Many cases, data tends to be around a central value

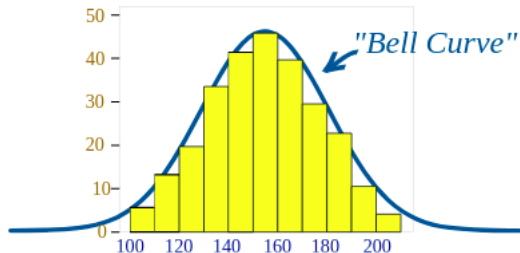
σ : Standard deviation of x

μ : Mean of x

$\pi \approx 3.14159\dots$

$e \approx 2.71828\dots$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$



Normal Distribution

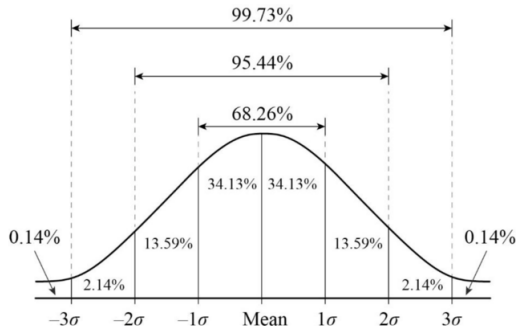
Many things closely follow a Normal Distribution:

- Heights of people
- Size of things produced by machines
- Errors in measurements
- Blood pressure
- Marks on a test

Normal Distribution

The 68 - 95 - 99.7 Rule

- Mean μ
- Standard deviation σ
- Approximately 68% observations fall within σ of the mean μ .
- Approximately 95% observations fall within 2σ of σ .
- Approximately 99.7% observations fall within 3σ of σ .



Gaussian Naïve Bayes

If features are continuous values, the likelihood of the features is assumed to be Gaussian:

$$P(x_i|y) = \frac{1}{\sigma_y \sqrt{2\pi}} e^{-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}}$$

Algorithm

$$\hat{y} = \underset{y \in C}{\operatorname{argmax}} p(y) \prod_{i=1}^d p(x_i|y)$$

Very easy!

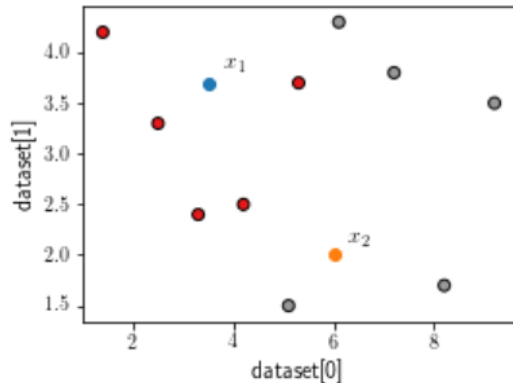
e.g. Gaussian Naïve Bayes Classifier

dataset:

(3.3 , 2.4 , 0),
(2.5 , 3.3 , 0),
(4.2 , 2.5 , 0),
(1.4 , 4.2 , 0),
(5.3 , 3.7 , 0),
(6.1 , 4.3 , 1),
(9.2 , 3.5 , 1),
(7.2 , 3.8 , 1),
(5.1 , 1.5 , 1),
(8.2 , 1.7 , 1)

$$x_1 = (3.5 , 3.7)$$

$$x_2 = (6 , 2)$$



	$y = 0$	$y = 1$
feature 1		
μ	3.34	7.16
σ	1.50	1.63
feature 2		
μ	3.22	2.96
σ	0.77	1.28

$$x_1 = (3.5, 3.7)$$

$$y = 0$$

$$p(x_1|y) = p(y) \times p(x_{10} = 3.5|y) \times p(x_{11} = 3.7|y) = 0.5 \times 0.26 \times 0.43 = 0.056$$

$$y = 1$$

$$p(x_1|y) = p(y) \times p(x_{10} = 3.5|y) \times p(x_{11} = 3.7|y) = 0.5 \times 0.02 \times 0.26 = 0.026$$

So, predict Class of x_1 is $y = 0$

Multinomial Naïve Bayes

Used in text classification

$\{(D, y)\}$, D : Document , y Class (Category of D).

$\theta_y = (\theta_{y1}, \dots, \theta_{yd})$ The distribution vector for class y
 d : number of feature = size of the vocabulary.

$\theta_{yi} = P(x_i \mid y)$ The probability of feature i appearing in a sample belonging to class y .

θ_y is estimated by a smoothed version of maximum likelihood, i.e. relative frequency counting:

$$\hat{\theta}_{yi} = \frac{N_{yi} + \alpha}{N_y + \alpha d}$$

$N_{yi} = \sum_{x \in T} x_i$: number of times feature i appears in a sample of class y in the training set T .

$N_y = \sum_{i=1}^d N_{yi}$: total count of all features for class y .

$\alpha \geq 0$ prevents zero probabilities.

Setting $\alpha = 1$: Laplace smoothing, while $\alpha < 1$ Lidstone smoothing.

e.g. Text classifier

No	Doc	Class
1	All Attended, All Exercises	A+
2	Attended, Exercises	A
3	No Attended	E
4	No Exercises	E
test	Attended, No Exercises	?

e.g. Text classifier

Feature = { All, Attended, Exercises, No }

$$\hat{y} = \underset{y \in C}{\operatorname{argmax}} p(y) \prod_{i=1}^d p(x_i|y)$$

$$P(A+) = 1/4 \quad P(A) = 1/4 \quad P(E) = 1/2$$

$$P(\text{Attended} \mid A+) = (1 + 1) / (4 + 1*4) = 2/8$$

$$P(\text{No} \mid A+) = (0 + 1) / (4 + 1*4) = 1/8$$

$$P(\text{Exercises} \mid A+) = (1 + 1) / (4 + 1*4) = 2/8$$

$$P(\text{Attended} \mid A) = (1 + 1) / (2 + 1*4) = 2/6$$

$$P(\text{No} \mid A) = (0 + 1) / (2 + 1*4) = 1/6$$

$$P(\text{Exercises} \mid A) = (1 + 1) / (2 + 1*4) = 2/6$$

$$P(\text{Attended} \mid E) = (1 + 1) / (4 + 1*4) = 2/8$$

$$P(\text{No} \mid E) = (2 + 1) / (4 + 1*4) = 3/8$$

$$P(\text{Exercises} \mid E) = (1 + 1) / (4 + 1*4) = 2/8$$

e.g. Text classifier

$$\hat{y} = \underset{y \in C}{\operatorname{argmax}} p(y) \prod_{i=1}^d p(x_i|y)$$

$$P(A+ \mid \text{test}) = 1/4 * 2/8 * 1/8 * 2/8 \approx 0.002$$

$$P(A \mid \text{test}) = 1/4 * 2/6 * 1/6 * 2/6 \approx 0.005$$

$$P(E \mid \text{test}) = 1/2 * 2/8 * 3/8 * 2/8 \approx 0.01 \quad \checkmark$$

Attended, No Exercises : E class.

Bernoulli Naïve Bayes

Multivariate Bernoulli distributions

$$D = \{ (x, y) \} , \quad x \in \mathbb{R}^d , \quad y \in C$$

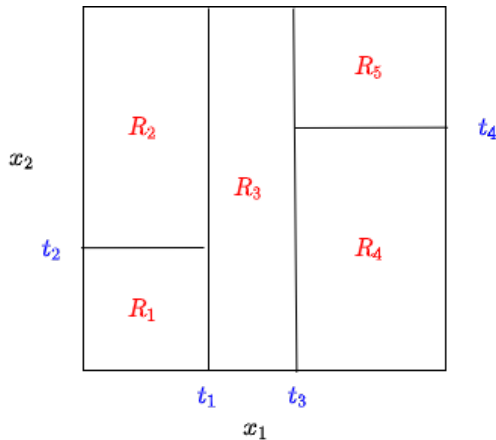
x_i binary-valued (Bernoulli, boolean) variable.

The decision rule for Bernoulli naive Bayes is based on

$$P(x_i|y) = P(i|y)x_i + (1 - P(i|y))(1 - x_i)$$

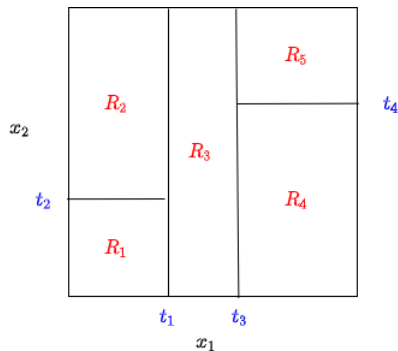
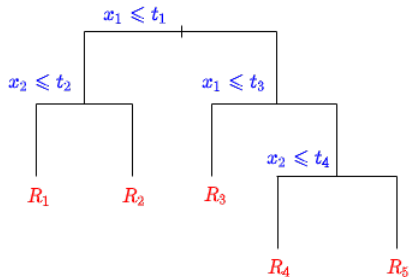
Decision Tree

e.g. split $x(x_1, x_2)$ to 5 region



Decision Tree

$$\{x(x_1, x_2), y\}$$



$$\hat{f}(x) = \sum_{m=1}^5 c_m I\{ (x_1, x_2) \in R_m \}$$

e.g. Regression Decision Tree

Dataset

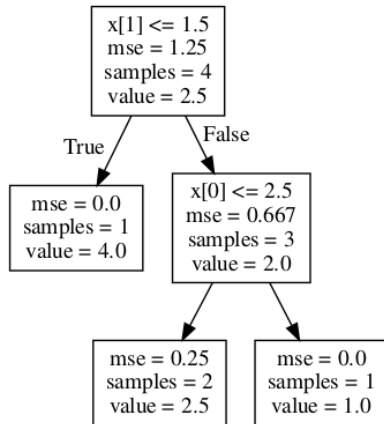
$$x = \{ (2, 4), (3, 3), (1, 2), (4, 1) \}$$

$$y = \{ 3, 1, 2, 4 \}$$

Predict

$$(2.3, 2) \Rightarrow \hat{y} = 2.5$$

$$(2.5, 1) \Rightarrow \hat{y} = 4.0$$



e.g. Classification Decision Tree

Dataset

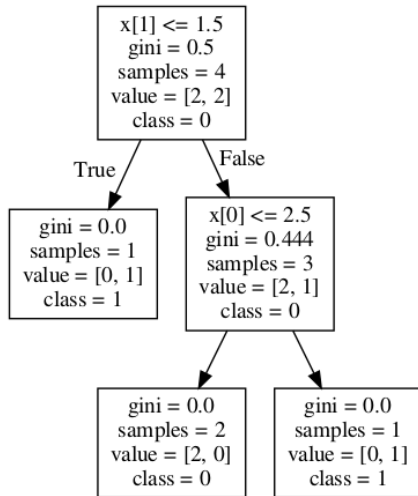
$x = \{ (2, 4), (3, 3), (1, 2), (0, 1) \}$

$y = \{ 0, 1, 0, 1 \}$

Predict

$(3.2, 1.4) \Rightarrow \hat{y} = 1$

$(2.3, 2.4) \Rightarrow \hat{y} = 0$



Decision Tree Model

- Dataset: $\{(x_i, y_i)\}_1^n$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$ (or $y_i \in \text{Category}$)
- Partition: M regions R_1, R_2, \dots, R_M
- Model response as a constant c_m in each region:
- I : Indicator function

$$f(x) = \sum_{m=1}^M c_m I\{x \in R_m\}$$

- Data at node R with n samples
- Candidate split (j, s) : feature j , threshold s

$$R(j, s) \rightarrow R_1(j, s) + R_2(j, s)$$

Left $R_1(j, s) = \{ (x, y) \mid x_j \leq s \}$

Right $R_2(j, s) = \{ (x, y) \mid x_j > s \}$

Decision Tree algorithm

Loss Function

- $H()$: **Loss Function**
- $G()$: **Information Gain**

$$G(R(j, s)) = \frac{n_1}{n} H(R_1(j, s)) + \frac{n_2}{n} H(R_2(j, s))$$

$$(j, s)^* = \underset{(j, s)}{\operatorname{argmin}} G(R(j, s))$$

Recurse

$R_1(j, s)$, $R_2(j, s)$

Until maximum allowable depth is reached

Decision Tree use for both classification and regression tasks.

Regression Decision Tree algo.

CART algorithm

Loss function base on Mean Squared Error (MSE or \mathcal{L}_2 error)

$$MSE(.) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

The best $f(x_i)$ is average of y_i on region R

$$c_k = \text{ave}(y_i \mid x_i \in R_k)$$

$$H(R_k) = MSE(R_k) = \frac{1}{n} \sum_{y \in R_k} (y - c_k)^2$$

$$G(R(j, s)) = \left\{ \frac{n_1}{n} \frac{1}{n_1} \sum_{y \in R_1} (y - c_1)^2 + \frac{n_2}{n} \frac{1}{n_2} \sum_{y \in R_2} (y - c_2)^2 \right\}$$

Regression Decision Tree algo.

CART algorithm

$$G(R(j, s)) = \frac{1}{n} \left\{ \sum_{y \in R_1} (y - c_1)^2 + \sum_{y \in R_2} (y - c_2)^2 \right\}$$

$$(j, s)^* = \underset{(j, s)}{\operatorname{argmin}} \left\{ \sum_{y \in R_1(j, s)} (y - c_1)^2 + \sum_{y \in R_2(j, s)} (y - c_2)^2 \right\}$$

Regression Decision Tree algo.

CART(R , *stop!*)

1. *list* $W = \{ \}$
2. *for all* j : *feature* x_j
 - ▶ *sort Domain* $\{x_j\}$
 - ▶ *for all* $t_k \in \text{Domain } \{x_j\}$
 - *choose* s : $s = \frac{(t_k + t_{k+1})}{2}$
 - $w(j, s) = \sum_{y \in R_1(j, s)} (y - c_1)^2 + \sum_{y \in R_2(j, s)} (y - c_2)^2$
 - *add* $w(j, s)$ *to list* W
3. $w(j, s) = \min\{W\}$
4. **CART**($R_1(j, s)$, *stop!*) , **CART**($R_2(j, s)$, *stop!*)

Regression Tree Algorithm

Problem

- How large should we grow the tree?
- Very large tree might overfit the data.
- Small tree might not show the important structure.
- Optimal tree size, we can choose from the dataset.
- Real domain x !. We have to partition it.

Classification Decision Tree algo.

Loss function base on proportion p_{mk} of class k in region R_m with N_m observations:

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k) = \frac{n_k}{N_m}$$

- **Gini**

$$H(R_m) = \sum_{k=1}^K p_{mk}(1 - p_{mk}) = 1 - \sum_{k=1}^K p_{mk}^2$$

- **Entropy**

$$H(R_m) = - \sum_{k=1}^K p_{mk} \log p_{mk}$$

- **Misclassification**

$$H(R_m) = 1 - \max(p_{mk})$$

CART algorithm

gini

$$(j, s)^* = \underset{(j, s)}{\operatorname{argmin}} \left\{ \frac{n_1}{n} H(R_1) + \frac{n_2}{n} H(R_2) \right\}$$

$$(j, s)^* = \underset{(j, s)}{\operatorname{argmin}} \left\{ \frac{n_1}{n} \left(1 - \sum_{k=1}^K p_{1k}^2 \right) + \frac{n_2}{n} \left(1 - \sum_{k=1}^K p_{2k}^2 \right) \right\}$$

$$(j, s)^* = \underset{(j, s)}{\operatorname{argmin}} \left\{ 1 - \left(\frac{n_1}{n} \sum_{k=1}^K p_{1k}^2 + \frac{n_2}{n} \sum_{k=1}^K p_{2k}^2 \right) \right\}$$

$$(j, s)^* = \underset{(j, s)}{\operatorname{argmax}} \left\{ n_1 \sum_{k=1}^K p_{1k}^2 + n_2 \sum_{k=1}^K p_{2k}^2 \right\}$$

CART algorithm (gini)

CART(R , *stop!*)

1. *list* $W = \{ \}$
2. *for all* j : *feature* x_j
 - ▶ *sort Domain* $\{x_j\}$
 - ▶ *for all* $t_k \in \text{Domain } \{x_j\}$
 - *choose* s : $s = \frac{(t_k + t_{k+1})}{2}$
 - $w(j, s) = n_1 \sum_{k=1}^K p_{1k}^2 + n_2 \sum_{k=1}^K p_{2k}^2$
 - *add* $w(j, s)$ *to list* W
3. $w(j, s) = \max\{W\}$
4. **CART**($R_1(j, s)$, *stop!*) , **CART**($R_2(j, s)$, *stop!*)

e.g. Classification Decision Tree

Dataset

$$x = \{ 3, 1, 0, 4, 2 \}$$

$$y = \{ 2, 0, 2, 1, 0 \}$$

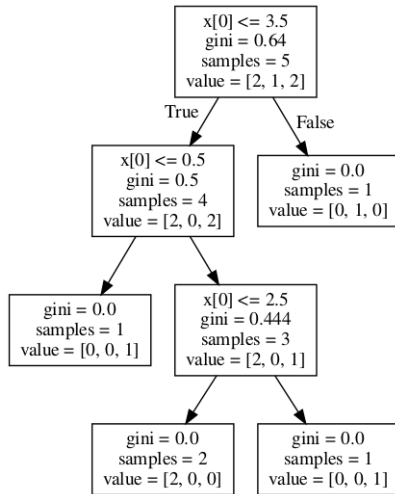
Predict

$$value = [n_0, n_1, n_2]$$

number of class

$$(x = 4.5) \Rightarrow \hat{y} = 1$$

$$(x = 1.3) \Rightarrow \hat{y} = 0$$



Equation of a Line

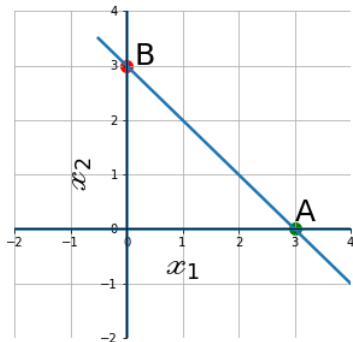
in \mathbb{R}^2 Line from 2 points.

$$\begin{aligned}\frac{x_1 - x_1^A}{x_1^A - x_1^B} &= \frac{x_2 - x_2^A}{x_2^A - x_2^B} \\ \iff \frac{x_1 - 3}{3 - 0} &= \frac{x_2 - 0}{0 - 3} \\ \iff x_2 &= -x_1 + 3 \\ \iff x_1 + x_2 - 3 &= 0\end{aligned}$$

General Line (*dot product*)

$$\bar{w} \cdot \bar{x} + b = 0$$

with $\bar{w} = (w_1, w_2)$, $\bar{x} = (x_1, x_2)$



$$A(x_1^A, x_2^A), B(x_1^B, x_2^B)$$

Normal and Direction vector

in \mathbb{R}^2

$$\overline{w} \cdot \overline{x} + b = 0$$

$\overline{w} = (w_1, w_2)$: Normal vector.

$\overline{s} = (w_2, -w_1)$: Direction vector.

$$\text{Norm of } \overline{w} \quad \|\overline{w}\| = \sqrt{w_1^2 + w_2^2}$$

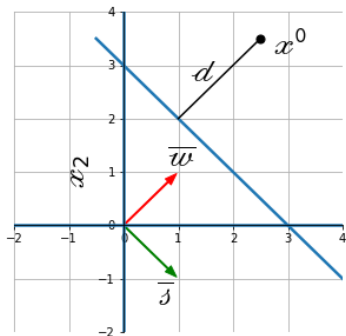
d : distance from x^0 to line.

$$d = \frac{|\overline{w} \cdot \overline{x}^0 + b|}{\|\overline{w}\|}$$

$$\Rightarrow \text{distance from } (0,0) \text{ to line} = \frac{|b|}{\|\overline{w}\|}$$

$$\overline{u} \cdot \overline{v} = u_1 v_1 + u_2 v_2 + \dots + u_n v_n \quad (\text{algebraic})$$

$$\overline{u} \cdot \overline{v} = \|\overline{u}\| \|\overline{v}\| \cos(\alpha) \quad (\text{geometric})$$



Distance between Parallel Lines

in \mathbb{R}^2

$$\ell_1 : \bar{w} \cdot \bar{x} + b_1 = 0$$

$$\ell_2 : \bar{w} \cdot \bar{x} + b_2 = 0$$

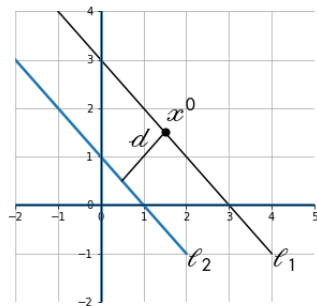
ℓ_1 parallel with ℓ_2

$$x_0 \in \ell_1 \Rightarrow \bar{w} \cdot \bar{x}^0 + b_1 = 0$$

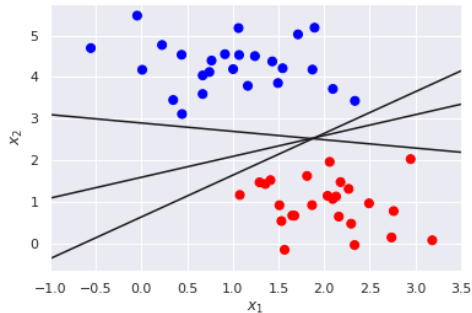
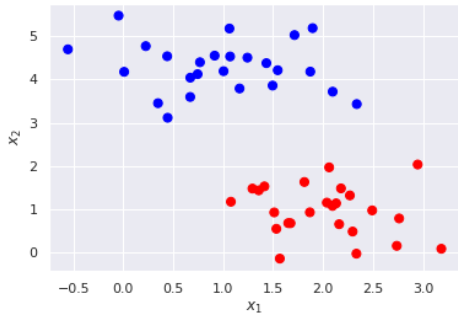
d : distance from x^0 to ℓ_2

$$d = \frac{|\bar{w} \cdot \bar{x}^0 + b_2|}{\|\bar{w}\|} = \frac{|b_2 - b_1|}{\|\bar{w}\|}$$

$$d = \frac{|b_2 - b_1|}{\|\bar{w}\|}, \quad \text{unit vector } \hat{w}, \quad \hat{w} = \frac{\bar{w}}{\|\bar{w}\|}$$

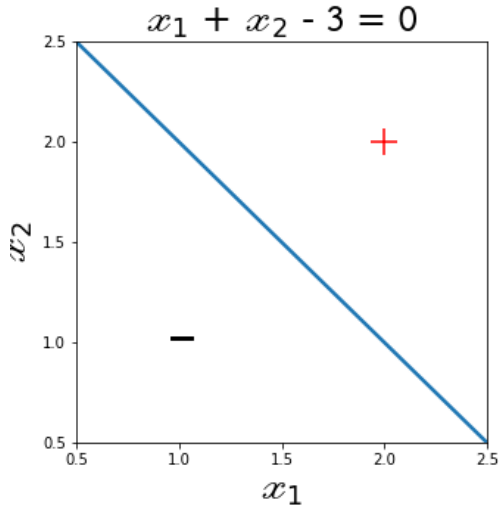
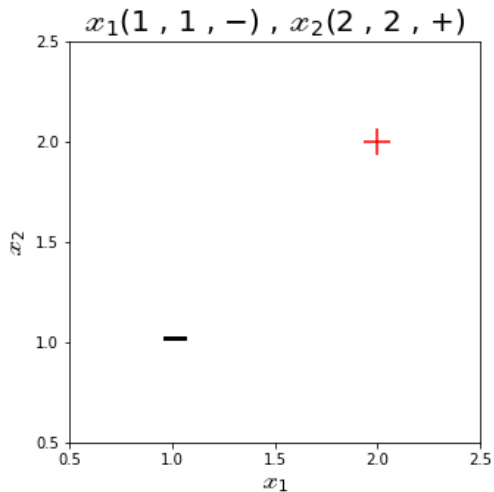


SVM Overview



- dataset from sklearn. 2 class.
- many lines can separate 2 class.

SVM Overview



SVM, Linear, Binary class

$$\mathcal{D} = \left\{ (x_i, y_i) \right\}_1^n, \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$$

Key idea: find widest separating “street”

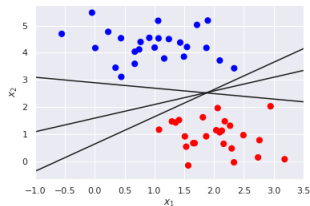
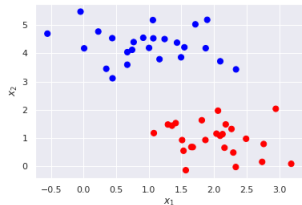
$$h^0 : w \cdot x + b = 0$$

$$h^+ : w \cdot x + b = +1$$

$$h^- : w \cdot x + b = -1$$

Model

$$f(x) = \text{sign}(w \cdot x + b)$$



SVM, Linear, binary classifiers

$$\forall (x_i, y_i)$$

$$h^+ : w \cdot x_i + b - 1 \geq 0, \quad y_i = +1$$

$$h^- : w \cdot x_i + b + 1 \leq 0, \quad y_i = -1$$

\Rightarrow

$$w \cdot x_i + b \geq +1, \quad y_i = +1$$

$$w \cdot x_i + b \leq -1, \quad y_i = -1$$

\Rightarrow

1st constrain

$$y_i(w \cdot x_i + b) \geq 1$$

Distance from h^+ to h^-

$$d = \frac{|(b-1)-(b+1)|}{\|w\|} = \frac{2}{\|w\|}$$

2nd constrain

$$\text{Maximize } \frac{2}{\|w\|} \Rightarrow \text{minimize } \frac{1}{2} \|w\|^2$$

SVM, Optimization with constraints

$$\underset{w,b}{\operatorname{argmin}} \quad \frac{1}{2} \|w\|^2$$

$$\text{s.t.} \quad y_i(w \cdot x_i + b) - 1 \geq 0 \quad , \quad \forall (x_i, y_i) \in \mathcal{D}$$

Lagrange Multipliers

$$\begin{aligned} \mathcal{L}(w, b, \alpha) &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i(w \cdot x_i + b) - 1] \\ &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i w \cdot x_i - \sum_{i=1}^n \alpha_i y_i b + \sum_{i=1}^n \alpha_i \quad (1) \end{aligned}$$

- Lagrange multiplier α_i
- Satisfies the Karush–Kuhn–Tucker (KKT)

$$\alpha_i [y_i(w \cdot x_i + b) - 1] = 0$$

$$\alpha_i \geq 0$$

- Minimize \mathcal{L} respect to w , b
- and maximize \mathcal{L} respect to α

SVM, Optimization with constraints

Minimize \mathcal{L} respect to w, b

$$\frac{\partial \mathcal{L}}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i \quad (2), \quad \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \quad (3)$$

and maximize \mathcal{L} respect to α use (1) and substitute by (2), (3):

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i w \cdot x_i - \sum_{i=1}^n \alpha_i y_i b + \sum_{i=1}^n \alpha_i \\ \mathcal{L}_{dual} &= \frac{1}{2} w \cdot w - w \cdot \left(\sum_{i=1}^n \alpha_i y_i x_i \right) - b \sum_{i=1}^n \alpha_i y_i + \sum_{i=1}^n \alpha_i \\ &= -\frac{1}{2} w \cdot w + \sum_{i=1}^n \alpha_i \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i \cdot x_j \end{aligned}$$

Dual problem

$$\begin{aligned} \underset{\alpha}{\operatorname{argmax}} \quad \mathcal{L}_{\text{dual}} &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i \cdot x_j \quad (4) \\ \text{s.t.} \quad &\sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0 \end{aligned}$$

w, b

$$\text{KKT:} \quad \alpha_i [y_i (w \cdot x_i + b) - 1] = 0$$

Two case:

1. $\alpha_i = 0$, or
2. $y_i (w \cdot x_i + b) - 1 = 0 \Rightarrow y_i (w \cdot x_i + b) = 1$

$\Rightarrow \alpha > 0$ then $y_i (w \cdot x_i + b) = 1$

and x_i must be a support vector

if a point is not a support vector, then $\alpha_i = 0$

SVM Linear Classifier

$$w = \sum_{\alpha_i > 0} \alpha_i y_i x_i$$

b_i per support vector:

$$\alpha_i (y_i (w \cdot x_i + b_i) - 1) = 0$$

$$y_i (w \cdot x_i + b_i) = 1$$

$$b_i = \frac{1}{y_i} - w \cdot x_i$$

$$b_i = y_i - w \cdot x_i$$

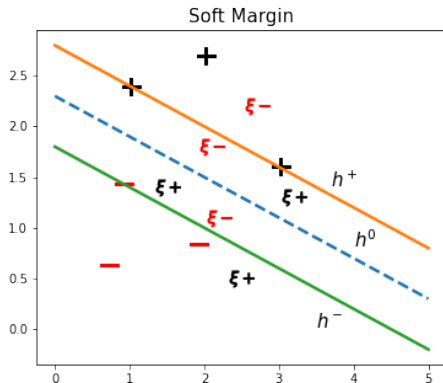
$$b = \underset{\alpha_i > 0}{avg} \{b_i\}$$

Linear Classifier

$$\hat{y} = \text{sign}(w \cdot x + b)$$

Soft Margin

- x_i away from hyperlane
- slack variables ξ_i for x_i
- $\xi_i \geq 0$
- $\xi_i = 0$, x_i at least $\frac{1}{\|w\|}$ away from hyperlane
- $0 < \xi_i < 1$ x_i within the margin
- $\xi_i \geq 1$ x_i appears on the wrong side



$$y_i(w \cdot x_i + b) \geq 1 - \xi_i$$

Soft Margin

$$\begin{aligned} \underset{w, b, \xi_i}{\operatorname{argmin}} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(w \cdot x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad \forall x_i \in \mathcal{D} \end{aligned}$$

$$C > 0$$

- large $C \sim$ small margin,
- small $C \sim$ large margin.

Soft Margin, Lagrange Multipliers

$$\begin{aligned} \underset{\alpha}{\operatorname{argmax}} \quad \mathcal{L}_{dual} &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i \cdot x_j \quad (4) \\ \text{s.t.} \quad &\sum_{i=1}^n \alpha_i y_i = 0 \\ &0 \leq \alpha_i \leq C \end{aligned}$$

w vector and b

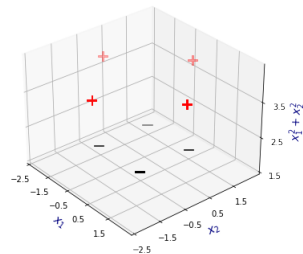
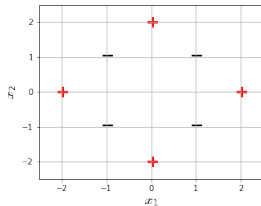
$$w = \sum_{\alpha_i > 0} \alpha_i y_i x_i$$

$$b = \operatorname{avg}_{\alpha_i > 0} \{b_i\}$$

Linear Classifier

$$\hat{y} = \operatorname{sign}(w \cdot x + b)$$

Kernel Trick



e.g. Simple Kernel Trick: $\mathcal{R}^2 \rightarrow \mathcal{R}^3$

$$\phi(x)$$

$$\phi(x) = (x_1, x_2, x_1^2 + x_2^2)$$

Kernel Trick, Lagrange Multipliers

$$\begin{aligned} \underset{\alpha}{\operatorname{argmax}} \quad \mathcal{L}_{dual} &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad &\sum_{i=1}^n \alpha_i y_i = 0 \\ &0 \leq \alpha_i \leq C \end{aligned} \quad (5)$$

w vector and b

$$w = \sum_{\alpha_i > 0} \alpha_i y_i \phi(x_i)$$

$$b = \operatorname{avg}_{\alpha_i > 0} \{b_i\}$$

Kernel SVM Classifier

$$\hat{y} = \operatorname{sign}(w \cdot \phi(x) + b)$$

Dual Solution: Gradient Ascent

(5) \Rightarrow

$$\mathcal{L}(\alpha_k) = \alpha_k - \frac{1}{2} \alpha_k^2 y_k^2 K(x_k \cdot x_k) - \alpha_k y_k \sum_{\substack{i=1 \\ i \neq k}}^n \alpha_i y_i K(x_i \cdot x_k)$$

$$\nabla \mathcal{L}(\alpha) = \left(\frac{\partial \mathcal{L}(\alpha)}{\partial \alpha_1}, \frac{\partial \mathcal{L}(\alpha)}{\partial \alpha_2}, \dots, \frac{\partial \mathcal{L}(\alpha)}{\partial \alpha_n} \right)$$

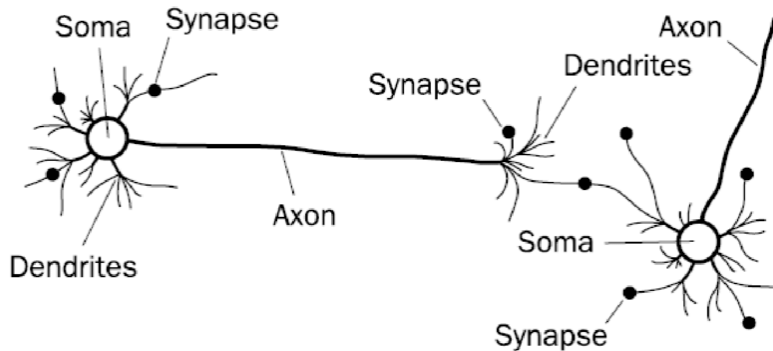
$$\frac{\partial \mathcal{L}}{\partial \alpha_k} = 1 - y_k \left(\sum_{i=1}^n \alpha_i y_i \tilde{K}(x_i \cdot x_k) \right)$$

Gradient ascent

$$\alpha_{t+1} = \alpha_t + \eta \nabla \mathcal{L}(\alpha_t)$$

This algorithm can demo by scikit-learn.

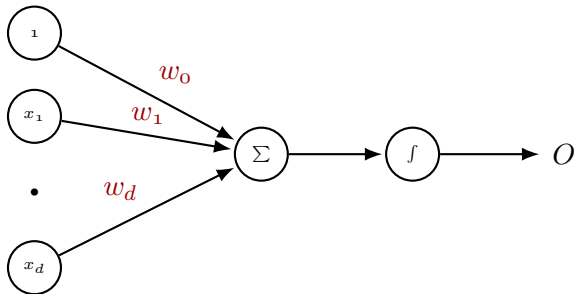
Biological Neuron



Biological Neuron vs. Artificial Neuron

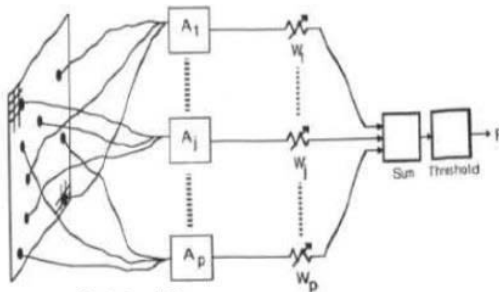
Biological Neuron	Artificial Neuron
Cell Nucleus (Soma)	Node
Dendrites	Input
Synapse	Weights or interconnections
Axon	Output

Artificial Neuron



Rosenblatt's Perceptron

Perceptron (1957)

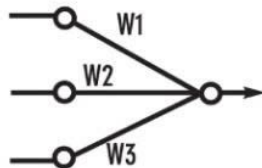


Frank Rosenblatt
(1928-1971)

Original Perceptron

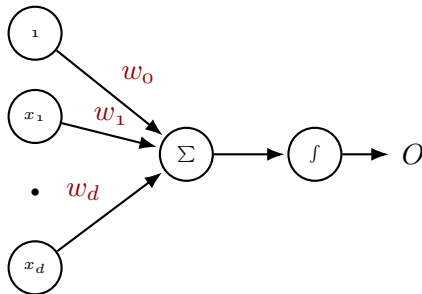
(From *Perceptrons* by M. L. Minsky and S. Papert,
1969, Cambridge, MA: MIT Press. Copyright 1969
by MIT Press.)

Simplified model:



Rosenblatt's Perceptron

$$\mathcal{D} = \left\{ (x^{(i)}, y^{(i)}) \right\}_1^n, \quad x^{(i)} \in \mathbb{R}^d, \quad y^{(i)} \in \{\mathcal{C}_1, \mathcal{C}_2\} \equiv \{0, 1\}$$



Input: $x = [+1, x_1, x_2, \dots, x_d]^T$ **Weight:** $w = [w_0, w_1, w_2, \dots, w_d]^T$

$$\hat{y}^{(i)} = f(x^{(i)}) = \begin{cases} 1 & \text{if } \sum_{j=0}^d w_j x_j^{(i)} > 0 \\ 0 & \text{otherwise} \end{cases}$$

Perceptron Training Rule

Initialization: Initial weight vector w and learning rate η , set to some small value (e.g., 0.1)

$$w = random \quad , \quad \eta \in (0, 1]$$

Loop for n epoch and *not Convergence*
for each training $(x^{(i)}, y^{(i)})$ in dataset \mathcal{D}

Activation function:

$$\text{sign}\left(\sum_{j=0}^d w_j x_j^{(i)}\right) = \begin{cases} 1 & \text{if } \sum_{j=0}^d w_j x_j^{(i)} > 0 \\ 0 & \text{otherwise} \end{cases}$$

Class response: $\hat{y}^{(i)} = \text{sign}\left(\sum_{j=0}^d w_j x_j^{(i)}\right)$

Modify weight vector: $w = w + \eta(y^{(i)} - \hat{y}^{(i)})x^{(i)}$

End

Perceptron Training Rule

$$w = w + \eta(y^{(i)} - \hat{y}^{(i)})x^{(i)}$$

- This learning procedure can be proven to convergen to a weight vector that correctly classifies all training samples (Minsky and Papert 1969).
- If the data are not linearly separable, convergence is not assured.

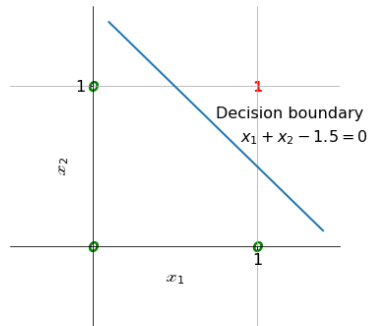
Perceptron Training Rule

AND gate



How to learn?

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



Perceptron learning AND gate

Initial Weight vector: $w = (-0.5, 0.4, 1)$, Learning rate: $\eta = 0.1$

$(1, x)$	w	$w^T x$	\hat{y}	y	$w + \eta(y - \hat{y})x$
$(1, 0, 0)$	$(-0.5, 0.4, 1)$	-0.5	0	0	<i>not update w</i>
$(1, 0, 1)$	$(-0.5, 0.4, 1)$	0.5	1	0	$(-0.5, 0.4, 1) + 0.1(0 - 1)(1, 0, 1)$
$(1, 1, 0)$	$(-0.6, 0.4, 0.9)$	-0.2	0	0	<i>not update w</i>
$(1, 1, 1)$	$(-0.6, 0.4, 0.9)$	0.7	1	1	<i>not update w</i>
$(1, 0, 0)$	$(-0.6, 0.4, 0.9)$	-0.6	0	0	<i>not update w</i>
$(1, 0, 1)$	$(-0.6, 0.4, 0.9)$	0.3	1	0	$(-0.6, 0.4, 0.9) + 0.1(0 - 1)(1, 0, 1)$
$(1, 1, 0)$	$(-0.7, 0.4, 0.8)$	-0.3	0	0	<i>not update w</i>
$(1, 1, 1)$	$(-0.7, 0.4, 0.8)$	0.5	1	1	<i>not update w</i>
$(1, 0, 0)$	$(-0.7, 0.4, 0.8)$	-0.7	0	0	<i>not update w</i>
$(1, 0, 1)$	$(-0.7, 0.4, 0.8)$	0.1	1	0	$(-0.7, 0.4, 0.8) + 0.1(0 - 1)(1, 0, 1)$
$(1, 1, 0)$	$(-0.8, 0.4, 0.7)$	-0.4	0	0	<i>not update w</i>
$(1, 1, 1)$	$(-0.8, 0.4, 0.7)$	0.3	1	1	<i>not update w</i>

Perceptron learning AND gate

$(1, x)$	w	$w^T x$	\hat{y}	y	$w + \eta(y - \hat{y})x$
$(1, 0, 0)$	$(-0.8, 0.4, 0.7)$	-0.8	0	0	<i>not update w</i>
$(1, 0, 1)$	$(-0.8, 0.4, 0.7)$	-0.1	0	0	<i>not update w</i>
$(1, 1, 0)$	$(-0.8, 0.4, 0.7)$	-0.4	0	0	<i>not update w</i>
$(1, 1, 1)$	$(-0.8, 0.4, 0.7)$	0.3	1	1	<i>not update w</i>

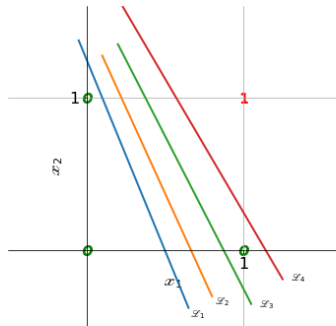
Learning result \mathcal{L}_4

$$\mathcal{L}_1 : 0.4x_1 + x_2 - 0.5 = 0$$

$$\mathcal{L}_2 : 0.4x_1 + 0.9x_2 - 0.6 = 0$$

$$\mathcal{L}_3 : 0.4x_1 + 0.8x_2 - 0.7 = 0$$

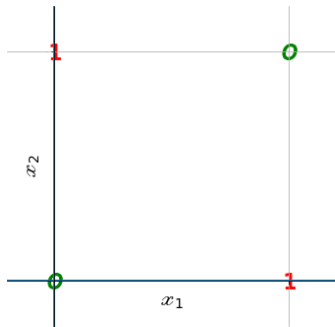
$$\mathcal{L}_4 : 0.4x_1 + 0.7x_2 - 0.8 = 0$$



Perceptron XOR problem



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



**Two class $\{0, 1\}$ cannot be separated by a single linear line.
XOR gate cannot be solved by Rosenblatt's perceptron.**

Perceptron XOR problem

Solution 1

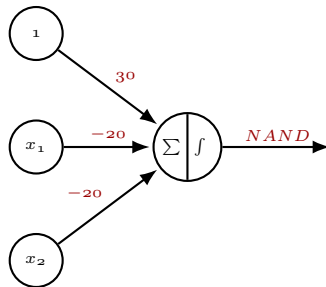
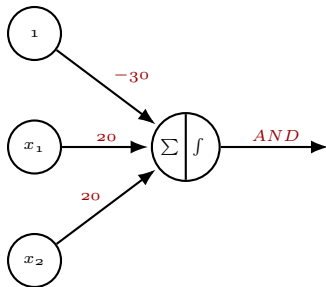
x_1	x_2	y	$\overline{x_1}x_2$	$x_1\overline{x_2}$	$\overline{x_1}x_2 + x_1\overline{x_2}$
0	0	0	0	0	0
0	1	1	1	0	1
1	0	1	0	1	1
1	1	0	0	0	0

Solution 2

x_1	x_2	y	$\overline{x_1x_2}$	$x_1 + x_2$	$\overline{x_1x_2}(x_1 + x_2)$
0	0	0	1	0	0
0	1	1	1	1	1
1	0	1	1	1	1
1	1	0	0	1	0

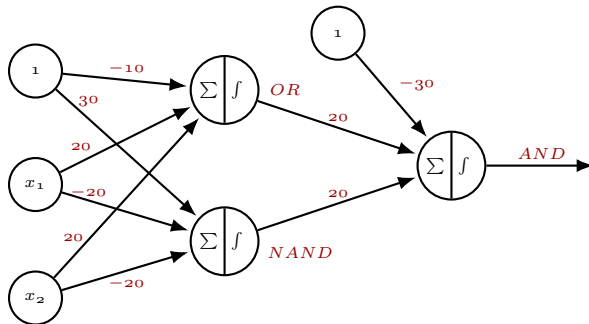
Perceptron XOR problem

Perceptron for Solution 2



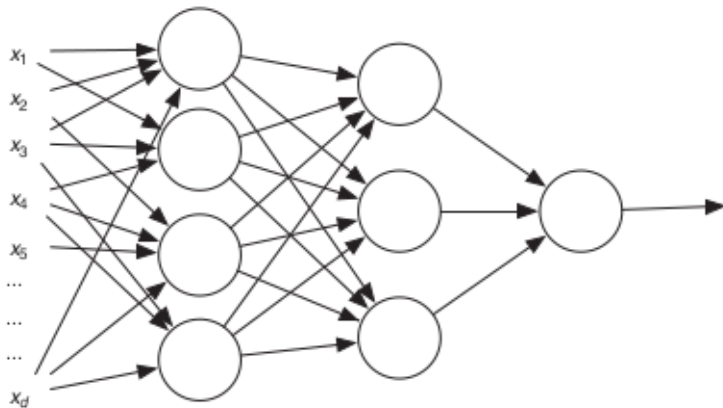
Perceptron XOR problem

Connect three gate

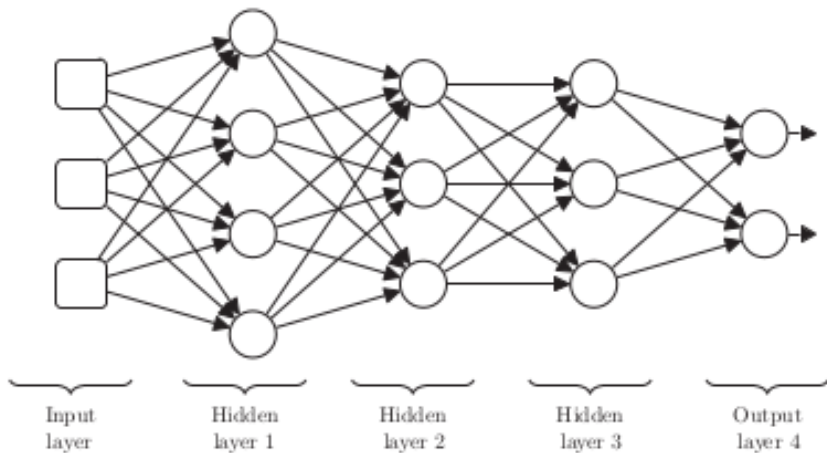


XOR as a combination of 3 basic perceptrons

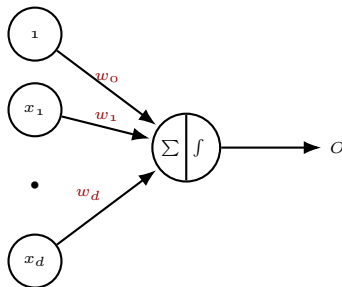
Multilayer Networks



Multilayer Networks



Activation function

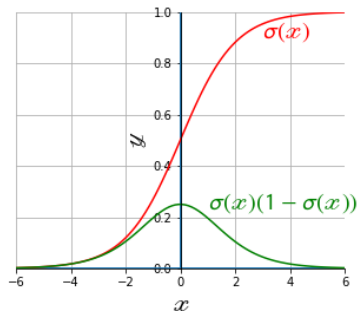


$$f(x^{(i)}) = g \left(\sum_{j=0}^d w_j x_j^{(i)} \right)$$

Activate function g

step (sign), sigmoid, tanh, relu, softmax, ...

Sigmoid activation function



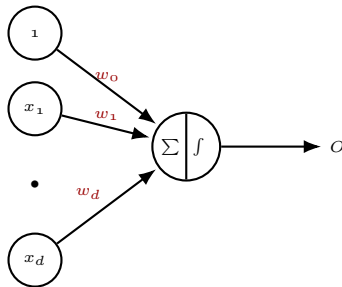
- $e = 2.71828...$ Continuous functions, differentiable, use for Gradient descent.

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

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

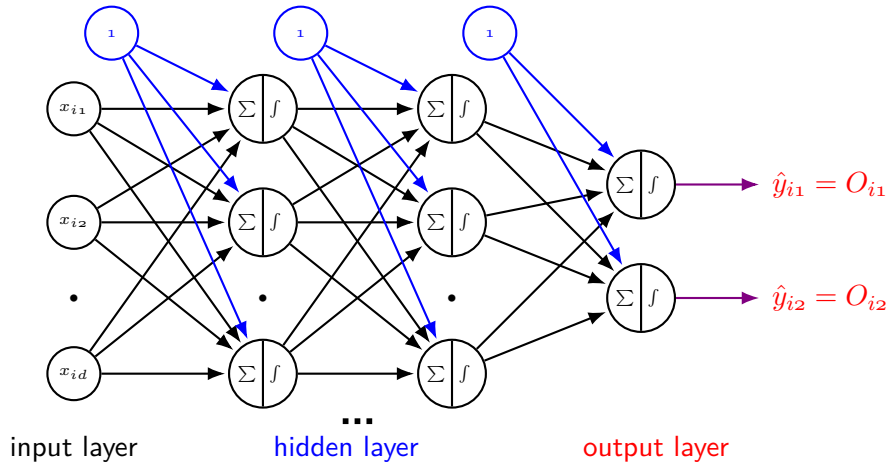
- Useful for tiny change in the weight.

Perceptron (one neuron)



$$Net = \sum_{j=0}^d w_j x_j ; \quad O = g(Net) ; \quad g: \text{Activate Function}$$

Neural Network (many neurons)



Loss function

- Network consider \mathcal{K} outputs.
- $y_k \triangleq t_k$: Target $k^{th} \in \mathcal{K}$ of example $(x, y) \in \mathcal{D}$
- $\hat{y}_k \triangleq o_k$: Output $k^{th} \in \mathcal{K}$ of networks of $(x, y) \in \mathcal{D}$

Loss function:

Sum the errors over all of the network output units

$$\begin{aligned} Loss(w) &\triangleq E(w) = \frac{1}{2} \sum_{\mathcal{D}} \sum_{\mathcal{K}} (y_k - \hat{y}_k)^2 \\ &= \frac{1}{2} \sum_{\mathcal{D}} \sum_{\mathcal{K}} (y_k - o_k)^2 \\ &= \frac{1}{2} \sum_{\mathcal{D}} \sum_{\mathcal{K}} (t_k - o_k)^2 \end{aligned}$$

Employs gradient descent algorithm to minimize $E(w)$.

Loss function for one example

- One example $e = (x, y)$.
- After forward network, we have k out put $o_k \triangleq \hat{y}_k$.
- Error all k out put o_k through w denote by $E_e(w)$:

$$E_e(w) = \frac{1}{2} \sum_k (y_k - o_k)^2$$

Employs Gradient descent

The Backprop weight update rule is based on the gradient descent method w (abandon index e)

$$w_{ij} = w_{ij} - \eta \frac{\partial E(w)}{\partial w_{ij}}$$

η : learning rate.

Backpropagation algorithm

Notations:

- x_{ij} : the input from node i to unit j .
- w_{ij} : the weight associated with x_{ij} .
- $z_j = \sum w_{ij}x_{ij}$, weighted sum of inputs for unit j .
- o_j : output computed by unit j .
- σ : the sigmoid function.
- \mathcal{K} : the set of units in the final layer of the network.
- $\text{succ}(j)$: the set of units whose immediate inputs from the output of unit j .

Chain Rule

Chain rule to write:

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial z_j} \frac{\partial z_j}{\partial w_{ij}}$$

since,

$$\frac{\partial z_j}{\partial w_{ij}} = \left(\sum w_{ij} x_{ij} \right)' = (w_{ij} x_{ij})' + \left(\sum_{k \neq i} w_{kj} x_{kj} \right)' = x_{ij} + 0$$

so,

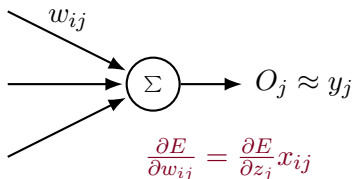
$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial z_j} x_{ij}$$

Consider two cases: the case where unit j is an output for the network, and the case where j is an internal unit.

Two cases in calculating $\frac{\partial E}{\partial z_j}$

Backpropagation algorithm

– Case 1: Neuron j is an output neuron



$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial o_j} \frac{\partial o_j}{\partial z_j}$$

$$\frac{\partial E}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_k (y_k - o_k)^2, \quad (k \in \mathcal{K})$$

$$\frac{\partial}{\partial o_j} (y_k - o_k)^2 = 0 \text{ for all output units } k \text{ except when } k = j.$$

$$\begin{aligned} \frac{\partial E}{\partial o_j} &= \frac{\partial}{\partial o_j} \frac{1}{2} (y_j - o_j)^2 \\ &= \frac{1}{2} 2 (y_j - o_j) \frac{\partial (y_j - o_j)}{\partial o_j} \\ &= -(y_j - o_j) \end{aligned}$$

Backpropagation algorithm

Since $o_j = \sigma(z_j)$

$\frac{\partial o_j}{\partial z_j} = \sigma(z_j)(1 - \sigma(z_j))$. therefore,

$$\frac{\partial o_j}{\partial z_j} = \frac{\partial \sigma(z_j)}{\partial z_j} = o_j(1 - o_j)$$

substituting, we obtain:

$$\frac{\partial E}{\partial z_j} = -(y_j - o_j) o_j (1 - o_j)$$

so

$$w_{ij} = w_{ij} + \eta (y_j - o_j) o_j (1 - o_j) x_{ij}$$

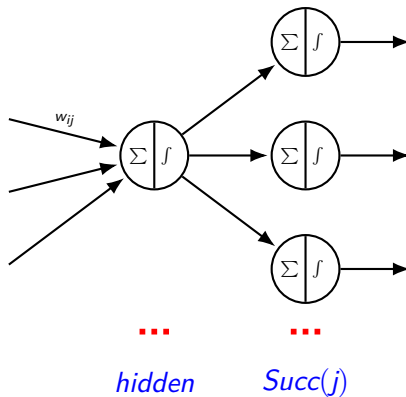
we will note:

$$\delta_j = -\frac{\partial E}{\partial z_j} = (y_j - o_j) o_j (1 - o_j)$$

$$w_{ij} = w_{ij} + \eta \delta_j x_{ij}$$

Backpropagation algorithm

- Case 2: Neuron j is a hidden neuron



$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial z_j} x_{ij}$$

Backpropagation algorithm

$$\begin{aligned}\frac{\partial E}{\partial z_j} &= \sum_{k \in \text{succ}(j)} \frac{\partial E}{\partial z_k} \frac{\partial z_k}{\partial z_j} \\&= \sum_{k \in \text{succ}(j)} -\delta_k \frac{\partial z_k}{\partial z_j} \\&= \sum_{k \in \text{succ}(j)} -\delta_k \frac{\partial z_k}{\partial o_j} \frac{\partial o_j}{\partial z_j} \\&= \sum_{k \in \text{succ}(j)} -\delta_k w_{jk} \frac{\partial o_j}{\partial z_j} \\&= \sum_{k \in \text{succ}(j)} -\delta_k w_{jk} o_j (1 - o_j)\end{aligned}$$

Backpropagation algorithm

note

$$\delta_j = -\frac{\partial E}{\partial z_j} = o_j (1 - o_j) \sum_{k \in \text{succ}(j)} \delta_k w_{jk}$$

and

$$w_{ij} = w_{ij} + \eta \delta_j x_{ij}$$

Backpropagation algorithm

Input:

(x, y) training example.

η learning rate (e.g., $\eta = 0.1$).

n_i the number of network inputs.

n_h the number of units in the hidden layer.

n_o the number of output units.

The input from unit i into unit j is denoted x_{ij} , and the weight from unit i to unit j is denoted w_{ij} .

Output:

A neural network with one input layer, one hidden layer and one output layer with n_i , n_h and n_o number of neurons respectively and all its weights.

1. Create_feedforward_network (n_i, n_h, n_o)
2. Initialize all network weights to small random numbers (e.g., between $-.05$ and $.05$).

Backpropagation algorithm

3. Repeat until convergence: *small changes in the weights*

For each training example (x, y)

3.1 **Feed forward**: Propagate example x through the network and compute the output o_j from every neuron.

3.2 **Propagate backward**: Propagate the errors backward.

Case 1 For each output neuron k , calculate its error

$$\delta_k = o_k(1 - o_k)(y_k - o_k)$$

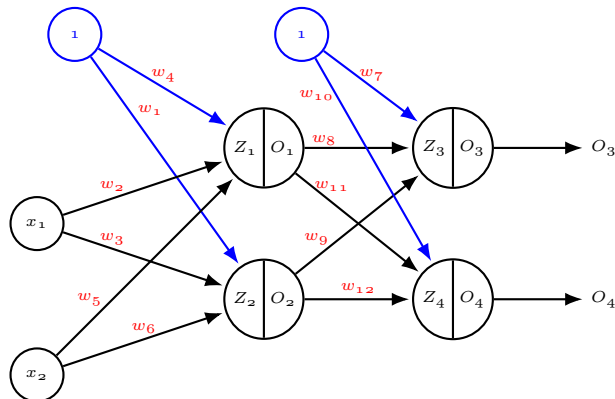
Case 2 For each hidden neuron h , calculate its error

$$\delta_h = o_h(1 - o_h) \sum_{k \in \text{succ}(h)} w_{hk} \delta_k$$

3.3. **Update each weight**:

$$w_{ij} = w_{ij} + \eta \delta_j x_{ij}$$

Example



$$(x_1 = 0.2, x_2 = 0.3) \quad (y_1 = 1, y_2 = 0)$$

$$w_1 = w_4 = 0.01, w_2 = 0.02, w_3 = 0.03, w_5 = 0.05, w_6 = 0.06$$

$$w_7 = w_{10} = 0.1, w_8 = 0.08, w_9 = 0.09, w_{11} = 0.11, w_{12} = 0.12$$

Step by step

$$\eta = 0.1$$

Feed Forward:

$$z_1 = x_1 w_2 + x_2 w_3 + w_4 = 0.023$$

$$o_1 = \sigma(z_1) = 0.506$$

$$z_2 = x_1 w_5 + x_2 w_6 + w_1 = 0.038$$

$$o_2 = \sigma(z_1) = 0.509$$

$$z_3 = o_1 w_8 + o_2 w_9 + w_7 = 0.186$$

$$o_3 = \sigma(z_3) = 0.546$$

$$z_4 = o_1 w_{11} + o_2 w_{12} + w_{10} = 0.217$$

$$o_4 = \sigma(z_4) = 0.554$$

Propagate backward:

$$\delta_3 = o_3(1 - o_3)(y_1 - o_3) = 0.112$$

$$\delta_4 = o_4(1 - o_4)(y_2 - o_4) = -0.137$$

$$\delta_1 = o_1(1 - o_1)(w_8 \delta_3 + w_{11} \delta_4) = -0.002 \quad \delta_2 = o_2(1 - o_2)(w_9 \delta_3 + w_{12} \delta_4) = -0.002$$

Step by step

Propagate backward Update weights:

$$w_7 = w_7 + \eta \delta_3 \cdot 1 = 0.112$$

$$w_8 = w_8 + \eta \delta_3 \cdot o_1 = 0.086$$

$$w_9 = w_9 + \eta \delta_3 \cdot o_2 = 0.096$$

$$w_{10} = w_{10} + \eta \delta_4 \cdot 1 = 0.085$$

$$w_{11} = w_{11} + \eta \delta_4 \cdot o_1 = 0.103$$

$$w_{12} = w_{12} + \eta \delta_4 \cdot o_2 = 0.113$$

$$w_1 = w_1 + \eta \delta_2 \cdot 1 = 0.010$$

$$w_2 = w_2 + \eta \delta_1 \cdot x_1 = 0.020$$

$$w_3 = w_3 + \eta \delta_1 \cdot x_2 = 0.030$$

$$w_4 = w_4 + \eta \delta_1 \cdot 1 = 0.010$$

$$w_5 = w_5 + \eta \delta_2 \cdot x_1 = 0.045$$

$$w_6 = w_6 + \eta \delta_2 \cdot x_2 = 0.060$$

Repeat until Convergence (small changes in the weights)

Wisdom of the Crowd

Weight-guessing contest the ox.

Sir Francis Galton (1822-1911)

book, 2004, “The Wisdom of Crowds: Why the Many Are Smarter Than the Few and How Collective Wisdom Shapes Business, Economies, Societies and Nations”.

James Michael Surowiecki.

Naturally, not all crowds are wise (for example, greedy investors of a stock market bubble)

Lior Rokach, 2010, PATTERN CLASSIFICATION USING ENSEMBLE METHODS

Wisdom of the Crowd

In order to become wise, the crowd should comply with the following criteria:

- **Diversity of opinion** Each member should have private information.
- **Independence** Members' opinions are not determined by the opinions of those around them.
- **Decentralization** Members are able to specialize and draw conclusions based on local knowledge.
- **Aggregation** Some mechanism exists for turning private judgments into a collective decision.

Lior Rokach, 2010, PATTERN CLASSIFICATION USING ENSEMBLE METHODS

Majority Voting

e.g. Ensemble of ten classifiers (using ten learning method)

Classifier	A score	B score	C score	Selected Label
1	0.2	0.7	0.1	B
2	0.1	0.1	0.8	C
3	0.2	0.3	0.5	C
4	0.1	0.8	0.1	B
5	0.2	0.6	0.2	B
6	0.6	0.3	0.1	A
7	0.25	0.65	0.1	B
8	0.2	0.7	0.1	B
9	0.2	0.2	0.8	C
10	0.4	0.3	0.3	A

Voted result

	Class A	Class B	Class C
Votes	2	5	3

Majority Voting

x : sample.

$y_k(x)$: classification of the k 'th classifier.

$I\{y = c\}$: Indicator function.

$$class(x) = \underset{c_i \in dom(y)}{argmax} \left(\sum_k I \{ y_k(x) = c_i \} \right)$$

Random Forest

