# CS 6316 Machine Learning

## Review of Linear Algebra and Probability

Yangfeng Ji

Department of Computer Science
University of Virginia

UNIVERSITY *of* VIRGINIA | ENGINEERING
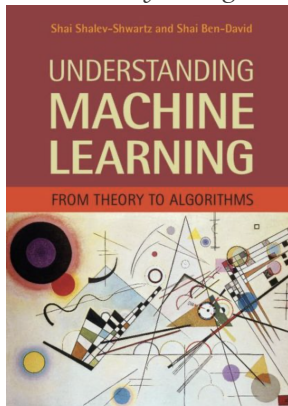
# Overview

# Course Information

# Instructors

- Yangfeng Ji
  - Office hour: Wednesday 11 AM - 12 PM
  - Office: Rice 510
- Hanjie Chen (TA)
  - Office hour: Tuesday and Thursday 1 PM – 2 PM
  - Office: Rice 442
- Kai Lin (TA)
  - Office hour: TBD

Understand the basic concepts and models from the computational perspective

To

- provide a wide coverage of basic topics in machine learning
  - Example: PAC learning, linear predictors, SVM, boosting, $k$NN, decision trees, neural networks, etc
- discuss a few fundamental concepts in each topic
  - Example: learnability, generalization, overfitting/underfitting, VC dimension, max margins methods, etc.

Shalev-Shwartz and Ben-David. *Understanding Machine Learning: From Theory to Algorithms.* 2014[1]



---

[1] https://www.cse.huji.ac.il/~shais/UnderstandingMachineLearning/index.html

# Outline

This course will cover the basic materials on the following topics

1. Learning theory
2. Linear classification and regression
3. Model selection and validation
4. Boosting and support vector machines
5. Neural networks
6. Clustering and dimensionality reduction

The following topics will not be the emphasis of this course

- ▶ Statistical modeling
  - ▶ Statistical Learning and Graphical Models by Farzad Hassanzadeh
- ▶ Deep learning
  - ▶ Deep Learning for Visual Recognition by Vicente Ordonez-Roman

# Reference Courses

For fans of machine learning:

- ▶ Shalev-Shwartz. Understanding Machine Learning. 2014
- ▶ Mohri. Foundations of Machine Learning. Fall 2018

# Reference Books

For fans of machine learning:

- ▶ Hastie, Tibshirani, and Friedman. The Elements of Statistical Learning (2nd Edition). 2009
- ▶ Murphy. Machine Learning: A Probabilistic Perspective. 2012
- ▶ Bishop. Pattern Recognition and Machine Learning. 2006
- ▶ Mohri, Rostamizadeh, and Talwalkar. Foundations of Machine Learning. 2nd Edition. 2018

- Homeworks (75%)
  - Five homeworks, each of them worth 15%
- Final project (22%)
  - Project proposal: 5%
  - Midterm report: 5%
  - Final project presentation: 6%
  - Final project report: 6%
- Class attendance (3%): we will take attendance at three randomly-selected lectures. Each is worth 1%

# Grading Policy

The final grade is threshold-based instead of percentage-based

| Point range | Letter grade |
|---|---|
| [99 100] | A+ |
| [94 99) | A |
| [90 94) | A- |
| [88 90) | B+ |
| [83 88) | B |
| [80 83) | B- |
| [74 80) | C+ |
| [67 74) | C |
| [60 67) | C- |

# Late Penalty

- Homework submission will be accepted up to 72 hours late, with 20% deduction per 24 hours on the points as a penalty
- It is usually better if students just turn in what they have in time
- Submission will not be accepted if more than 72 hours late
- Do not submit the wrong homework — late penalty will be applied if resubmit after deadline

Plagiarism, examples are

- ▶ in a homework submission, copying answers from others directly (even, with some minor changes)
- ▶ in a report, copying texts from a published paper (even, with some minor changes)
- ▶ in a code, using someone else's functions/implementations without acknowledging the contribution

- Course webpage

  `http://yangfengji.net/uva-ml-course/`

  which contains all the information you need about this course.

- Piazza

  `https://piazza.com/virginia/spring2020/cs6316/home`

# Basic Linear Algebra

Consider the following system of equations

$$4x_1 - 5x_2 = -13$$
$$-2x_1 + 3x_2 = 9 \tag{1}$$

In matrix notation, it can be written as a more compact from

$$\mathbf{A}x = b \tag{2}$$

with

$$\mathbf{A} = \begin{bmatrix} 4 & -5 \\ -2 & 3 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad b = \begin{bmatrix} -13 \\ 9 \end{bmatrix} \tag{3}$$

# Basic Notations
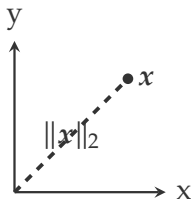
$$\mathbf{A} = \begin{bmatrix} 4 & -5 \\ -2 & 3 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad b = \begin{bmatrix} -13 \\ 9 \end{bmatrix}$$

- $\mathbf{A} \in \mathbb{R}^{m \times n}$: a matrix with $m$ rows and $n$ columns
  - The element on the $i$-th row and the $j$-th column is denoted as $a_{i,j}$
- $x \in \mathbb{R}^n$: a vector with $n$ entries. By convention, an $n$-dimensional vector is often thought of as matrix with $n$ rows and 1 column, known as a column vector.
  - The $i$-th element is denoted as $x_i$

# Vector Norms

- A norm of a vector $\|x\|$ is informally a measure of the "length" of the vector.
- Formally, a norm is any function $f : \mathbb{R}^n \to \mathbb{R}$ that satisfies four properties
    1. $f(x) \geq 0$ for any $x \in \mathbb{R}^n$
    2. $f(x) = 0$ if and only if $x = 0$
    3. $f(ax) = |a| \cdot f(x)$ for any $x \in \mathbb{R}^n$
    4. $f(x + y) \leq f(x) + f(y)$, for any $x, y \in \mathbb{R}^n$

# $\ell_2$ Norm

The $\ell_2$ norm of a vector $x \in \mathbb{R}^n$ is defined as

$$\|x\|_2 = \sqrt{\sum_{i=1}^{n} x_i^2} \tag{4}$$
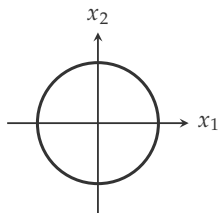


*Exercise*: prove $\ell_2$ norm satisfies all four properties

The $\ell_1$ norm of a vector $x \in \mathbb{R}^n$ is defined as

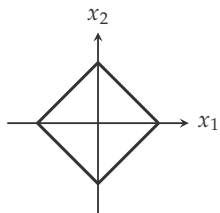$$\|x\|_1 = \sum_{i=1}^{n} |x_i| \tag{5}$$

For a two-dimensional vector $x = (x_1, x_2) \in \mathbb{R}^2$, which of the following plot is $\|x\|_1 = 1$?



(a)                    (b)                    (c)
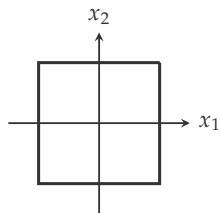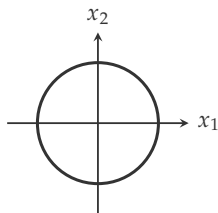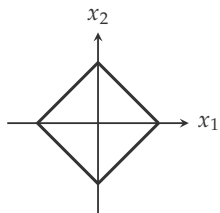
For a two-dimensional vector $x = (x_1, x_2) \in \mathbb{R}^2$, which of the following plot is $\|x\|_1 = 1$? Answer: (b)



|     (d)     |     (e)     |     (f)     |

# Dot Product

The dot product of $x, y \in \mathbb{R}^n$ is defined as

$$\langle x, y \rangle = x^\mathsf{T} y = \sum_{i=1}^{n} x_i y_i \tag{6}$$

where $x^\mathsf{T}$ is the transpose of $x$.

- $\|x\|_2^2 = \langle x, x \rangle$
- If $x = (0, 0, \ldots, \underbrace{1}_{x_i}, \ldots, 0)$, then $\langle x, y \rangle = y_i$

- If $x$ is an unit vector ($\|x\|_2 = 1$), then $\langle x, y \rangle$ is the projection of $y$ on the direction of $x$

# Cauchy-Schwarz Inequality

For all $x, y \in \mathbb{R}^n$

$$|\langle x, y \rangle| \leq \|x\|_2 \|y\|_2 \tag{7}$$

with equality if and only if $x = \alpha y$ with $\alpha \in \mathbb{R}$

**Proof**:

Let $\tilde{x} = \frac{x}{\|x\|_2}$ and $\tilde{y} = \frac{y}{\|y\|_2}$, then $\tilde{x}$ and $\tilde{y}$ are both unit vectors.

Based on the geometric interpretation on the previous slide, we have

$$\langle \tilde{x}, \tilde{y} \rangle \leq 1 \tag{8}$$

if and only if $\tilde{x} = \tilde{y}$.

The Forbenius norm of a matrix $\mathbf{A} = [a_{i,j}] \in \mathbb{R}^{m \times n}$ denoted by $\| \cdot \|_F$ is defined as

$$\|\mathbf{A}\|_F = \Big( \sum_i \sum_j a_{i,j}^2 \Big)^{1/2} \tag{9}$$

▶ The Frobenius norm can be interpreted as the $\ell_2$ norm of a vector when treating $\mathbf{A}$ as a vector of size $mn$.

## Two Special Matrices

▶ The identity matrix, denoted as $\mathbf{I} \in \mathbb{R}^{n \times n}]$, is a square matrix with ones on the diagonal and zeros everywhere else.

$$\mathbf{I} = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \tag{10}$$

▶ A diagonal matrix, denoted as $\mathbf{D} = \mathrm{diag}(d_1, d_2, \ldots, d_n)$, is a matrix where all non-diagonal elements are 0.

$$\mathbf{D} = \begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_n \end{bmatrix} \tag{11}$$

The *inverse* of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is denoted as $\mathbf{A}^{-1}$, which is the unique matrix such that

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I} = \mathbf{A}\mathbf{A}^{-1} \tag{12}$$

- ▶ Non-square matrices do not have inverses (by definition)
- ▶ Not all square matrices are invertible
- ▶ The solution of the linear equations in Eq. (1) is $x = \mathbf{A}^{-1}b$

# Orthogonal Matrices

▶ Two vectors $x, y \in \mathbb{R}^n$ are orthogonal if $\langle x, y \rangle = 0$



▶ A square matrix $\mathbf{U} \in \mathbb{R}^{n \times n}$ is orthogonal, if all its columns are orthogonal to each other *and* normalized (orthonormal)

$$\langle u_i, u_j \rangle = 0, \|u_i\| = 1, \|u_j\| = 1 \tag{13}$$

for $i, j \in [n]$ and $i \neq j$

▶ Furthermore, $\mathbf{U}^\mathsf{T}\mathbf{U} = \mathbf{I} = \mathbf{U}\mathbf{U}^\mathsf{T}$, which further implies $\mathbf{U}^{-1} = \mathbf{U}^\mathsf{T}$

A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is defined as

$$\mathbf{A}^\mathsf{T} = \mathbf{A} \tag{14}$$

or, in other words,

$$a_{i,j} = a_{j,i} \quad \forall i, j \in [n] \tag{15}$$

Comments

▶ The identity matrix $\mathbf{I}$ is symmetric
▶ A diagonal matrix is symmetric

Every symmetric matrix $\mathbf{A}$ can be decomposed as

$$\mathbf{A} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\mathsf{T} \tag{16}$$

with

- $\boldsymbol{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$ as a diagonal matrix (Slide 25)

- $\mathbf{Q}$ is an orthogonal matrix (Slide 27)

- *Exercise*: if $\mathbf{A}$ is invertible, show $\mathbf{A}^{-1} = \mathbf{U}\boldsymbol{\Lambda}^{-1}\mathbf{U}^\mathsf{T}$ with $\Lambda^{-1} = \mathrm{diag}(\frac{1}{\lambda_1}, \ldots, \frac{1}{\lambda_n})$

# Symmetric Positive Semidefinite Matrices

A symmetric matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ is positive semidefinite if and only if

$$x^{\mathsf{T}} \mathbf{P} x \geq 0 \qquad (17)$$

for all $x \in \mathbb{R}^n$.

# Symmetric Positive Semidefinite Matrices

A symmetric matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ is positive semidefinite if and only if

$$x^\top \mathbf{P} x \geq 0 \tag{17}$$

for all $x \in \mathbb{R}^n$.

Eigen decomposition (Slide 29) of $\mathbf{P}$ as

$$\mathbf{P} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top \tag{18}$$

with $\Lambda = \mathrm{diag}(\lambda_1, \ldots, \lambda_n)$ and

$$\lambda_i \geq 0 \tag{19}$$

# Symmetric Positive Definite Matrices

A symmetric matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ is positive definite if and only if

$$x^\mathsf{T} \mathbf{P} x > 0 \qquad (20)$$

for all $x \in \mathbb{R}^n$.

- Eigen values of $\mathbf{P}$, $\Lambda = \mathrm{diag}(\lambda_1, \ldots, \lambda_n)$ with

$$\lambda_i > 0 \qquad (21)$$

- *Exercise*: if one of the eigen values $\lambda_i < 0$, show that you can also find a vector $x$ such that $x^\mathsf{T} \mathbf{P} x < 0$

The identity matrix $\mathbf{I}$ is

- a diagonal matrix?
- a symmetric matrix?
- an orthogonal matrix?
- a positive (semi-)definite matrix?

Further reference [Kolter and Do, 2015]

The identity matrix $\mathbf{I}$ is

- a diagonal matrix? ✓
- a symmetric matrix? ✓
- an orthogonal matrix? ✓
- a positive (semi-)definite matrix? ✓

Further reference [Kolter and Do, 2015]

# Probability Theory

The probability of landing heads is 0.52

**Frequentist** Probability represents the *long-run frequency* of an event

- ▶ If we flip the coin many times, we expect it to land heads about 52% times

**Frequentist** Probability represents the *long-run frequency* of an event

▶ If we flip the coin many times, we expect it to land heads about 52% times

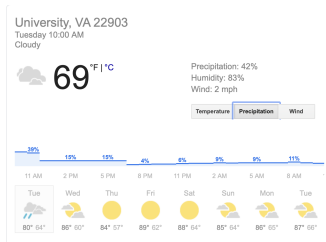**Bayesian** Probability quantifies our *(un)certainty* about an event

▶ We believe the coin is 52% of chance to land head on the next toss

Example scenarios of Bayesian interpretation of probability:

# Binary Random Variables

- **Event** $X$. Such as
  - *the coin will lead head on the next toss*
  - *it will rain tomorrow*
- Sample space of $X \in \{\text{false}, \text{true}\}$ or for simplicity $\{0, 1\}$

# Binary Random Variables

- **Event** $X$. Such as
  - *the coin will lead head on the next toss*
  - *it will rain tomorrow*
- Sample space of $X \in \{\text{false}, \text{true}\}$ or for simplicity $\{0, 1\}$
- Probability $P(X = x)$ or $P(x)$
- Let $X$ be the event that *the coin will lead head on the next toss*, then the probability from the previous example is

$$P(X = 1) = 0.52 \tag{22}$$

Given the binary random variable $X$ and its sample space as $\{0, 1\}$

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

with a single parameter $\theta$ as

$$\theta = P(X = 1)$$



Jacob Bernoulli

- Let $X$ be the number of heads
- Sample space of $X \in \{0, 1, 2\}$

- Let $X$ be the number of heads
- Sample space of $X \in \{0, 1, 2\}$
- Assume we use the same coin, the probability distribution of $X$
  - $P(X = 0) = (1 - \theta)^2$

# Tossing a Coin Twice?

- Let $X$ be the number of heads
- Sample space of $X \in \{0, 1, 2\}$
- Assume we use the same coin, the probability distribution of $X$
  - $P(X = 0) = (1 - \theta)^2$
  - $P(X = 2) = \theta^2$

- Let $X$ be the number of heads
- Sample space of $X \in \{0, 1, 2\}$
- Assume we use the same coin, the probability distribution of $X$
  - $P(X = 0) = (1 - \theta)^2$
  - $P(X = 2) = \theta^2$
  - $P(X = 1) = \theta(1 - \theta) + (1 - \theta)\theta = 2\theta(1 - \theta)$

Consider a general case, in which we toss the coin $n$ times, then the random variable $Y$ can be formulated as a binomial distribution

$$P(Y = k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \qquad (23)$$

where

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

is the binomial coefficient and

$$n! = n \cdot (n-1) \cdot (n-2) \cdots 1$$

How to define the corresponding random variable?

- ▶ $X \in \{1, 2, 3, 4, 5, 6\}$
- ▶ $X \in \{100000, 010000, 001000, 000100, 000010, 000001\}$

$$P(\mathbf{X} = \mathbf{x}) = \prod_{k=1}^{6} (\theta_k)^{x_k} \tag{24}$$

where

- $\mathbf{x} = (x_1, x_2, \ldots, x_6)$
- $x_k \in \{0, 1\}$, and
- $\{\theta_k\}_{k=1}^{6}$ are the parameters of this distribution, which is also the probability of side $k$ showing up.

## Multinomial Distribution

Repeat the previous event $n$ times, the corresponding probability distribution is modeled as

$$P(X = x) = \binom{n}{x_1 \cdots x_K} \prod_{k=1}^{K} \theta_k^{x_k} \tag{25}$$

where $x = (x_1, \ldots, x_K)$ and each $x_k \in \{0, 1, 2, \ldots, n\}$ indicates the number of times that side $k$ showing up.

$$\binom{n}{x_1 \cdots x_K} = \frac{n!}{x_1! \cdots x_K!}$$

The sum of $\{x_k\}$ follows the constraint:

$$\sum_{k=1}^{K} x_k = n$$

# Gaussian Distribution

A random variable $X \in \mathbb{R}$ is said to follow a normal (or Gaussian) distribution $\mathcal{N}(\mu, \sigma^2)$ if its probability density function is given by

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \tag{26}$$

- $\mu$: mean
- $\sigma^2$: variance
- Probability of $X \in [a, b]$: $P(a \leq X \leq b) = \int_a^b f(x)dx$

# Gaussian Distribution (II)

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \tag{27}$$

There examples of Gaussian distributions



- ▶ Blue: $\mathcal{N}(0,1)$ (standard normal distribution)
- ▶ Red: $\mathcal{N}(0,2)$
- ▶ Green: $\mathcal{N}(1,1)$

# Probability of Two Random Variables

Modeling two random variables together with a **joint** distribution

$$P(X, Y) \tag{28}$$

Related concepts

- Independence
- Conditional probability and chain rule
- Bayes rule

## Independence

**Definition** Two random variable $X$ and $Y$ are independent with each other, if we can represent the joint probability as the product of their marginal distributions for *any* values of $X$ and $Y$, or mathematically,

$$P(X, Y) = P(X) \cdot P(Y) \tag{29}$$

Marginal distributions

$$
\begin{aligned}
P(X) &= \sum_Y P(X, Y) \tag{30} \\
P(Y) &= \sum_X P(X, Y) \tag{31}
\end{aligned}
$$

## Independence

**Definition** Two random variable $X$ and $Y$ are independent with each other, if we can represent the joint probability as the product of their marginal distributions for *any* values of $X$ and $Y$, or mathematically,

$$P(X, Y) = P(X) \cdot P(Y) \tag{29}$$

Marginal distributions

$$P(X) = \sum_Y P(X, Y) \tag{30}$$

$$P(Y) = \sum_X P(X, Y) \tag{31}$$

► $X$: whether it is cloudy

► $Y$: whether it will rain

| $P(X \cap Y)$ | $X = 0$ | $X = 1$ |
|---------------|---------|---------|
| $Y = 0$ | 0.35 | 0.15 |
| $Y = 1$ | 0.05 | 0.45 |

## Conditional Probability

Conditional probability of $Y$ given $X$

$$P(Y \mid X) = \frac{P(X, Y)}{P(X)} \tag{32}$$

Example: document classification

- ▶ $X$: a document
- ▶ $Y$: the label of this document

A special case: if $X$ and $Y$ are independent

$$P(Y \mid X) = P(Y) \tag{33}$$

Intuitively, it means *Knowing X does not provide any new information about Y*

| $P(X, Y)$ | $X = 0$ | $X = 1$ |
|-----------|---------|---------|
| $Y = 0$   | 0.35    | 0.15    |
| $Y = 1$   | 0.05    | 0.45    |

- $X$: whether it is cloudy
- $Y$: whether it will rain

- $P(Y \mid X = 1)$:
  - $P(Y = 0 \mid X = 1) = 0.25,$
  - $P(Y = 1 \mid X = 1) = 0.75$

- $X$: whether it is cloudy
- $Y$: whether it will rain

| $P(X, Y)$ | $X = 0$ | $X = 1$ |
|-----------|---------|---------|
| $Y = 0$   | 0.35    | 0.15    |
| $Y = 1$   | 0.05    | 0.45    |

- $P(Y \mid X = 1)$:
    - $P(Y = 0 \mid X = 1) = 0.25$,
    - $P(Y = 1 \mid X = 1) = 0.75$
- $P(Y)$: $P(Y = 0) = P(Y = 1) = 0.5$

The probability density function of a multivariate Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is defined as

$$f(x) = \frac{1}{(2\pi)^{n/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(x - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1}(x - \boldsymbol{\mu})\right) \quad (34)$$

where

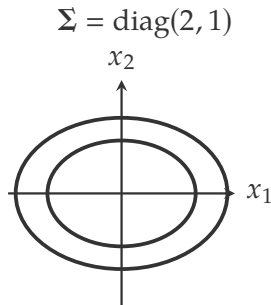- $\boldsymbol{\mu}$ is the $n$-dimensional mean vector and
- $\boldsymbol{\Sigma}$ is the $n \times n$ covariance matrix.

## Covariance Matrix $\Sigma$

Assume $\mu = 0$, the probability density function is

$$f(x) \propto \exp\left(-\frac{1}{2}x^{\mathsf{T}}\Sigma^{-1}x\right) \tag{35}$$

In general, $\Sigma$ is required to be a symmetric positive definite matrix



$\Sigma = I$

$\Sigma = \mathrm{diag}(2, 1)$

# Sampling from Gaussians



(a)  (b)

(a) $: \boldsymbol{\Sigma} = \mathbf{I}$

(b) $: \boldsymbol{\Sigma} = \text{diag}(2, 1)$

*Exercise*: Sample from an arbitrary Gaussian distribution

# Sum Rule

Given two random variables $X$ and $Y$ describing the same experiment, without any additional assumption we have

$$P(X \cup Y) = P(X) + P(Y) - P(X \cap Y) \tag{36}$$

▶ If $X \cap Y = \emptyset$, then

$$P(X \cap Y) = 0 \quad \text{and} \quad P(X \cup Y) = P(X) + P(Y) \tag{37}$$

▶ *Exercise*: Prove the following inequality by generalizing the sum rule in

$$P(\cup_{i=1}^{n} X_i) \le \sum_{i=1}^{n} P(X_i) \tag{38}$$

This inequality is called the union bound.

# Chain Rule

Any joint probability of two random variable can be decomposed as

$$P(X, Y) = P(X) \cdot P(Y \mid X) = P(Y) \cdot P(X \mid Y) \qquad (39)$$

No independence assumption is needed

## Chain Rule

Any joint probability of two random variable can be decomposed as

$$P(X, Y) = P(X) \cdot P(Y \mid X) = P(Y) \cdot P(X \mid Y) \tag{39}$$

No independence assumption is needed

The chain rule can be easily generalized

$$
\begin{aligned}
P(X_1, X_2, \cdots, X_k) &= P(X_1)P(X_2, \cdots, X_k \mid X_1) \\
&= P(X_1)P(X_2 \mid X_1)P(X_3, \cdots, X_k \mid X_2, X_1) \\
&= P(X_1)P(X_2 \mid X_1)P(X_3 \mid X_2, X_1) \cdots \\
&\quad P(X_k \mid X_1, \cdots, X_{k-1}) \tag{40}
\end{aligned}
$$

## Inverse Probability

Given

- ▶ $P(Y)$: prior probability, and
- ▶ $P(X \mid Y)$: conditional probability of $X$ given $Y$,

we can compute the probability $P(Y \mid X)$ using Bayes' rule as

$$P(Y \mid X) = \frac{P(Y)P(X \mid Y)}{P(X)} \qquad (41)$$

where

$$P(X) = \sum_Y P(Y)P(X \mid Y) \qquad (42)$$

Two random variables, alarm $A$ and burglar $B$

▶ $P(A = 1 \mid B = 1) = 0.99$: burglar happens, alarm rings

▶ $P(A = 1 \mid B = 0) = 0.001$: burglar does not happen, alarm rings

▶ $P(B = 1) = 0.01$: burglar rate

Question: if the alarm rang, what is the probability of a burglar happened?

$$P(B = 1 \mid A = 1) \tag{43}$$

## Example: The burglar alarm (II)

- $P(A = 1 \mid B = 1) = 0.99$: burglar happens $\Rightarrow$ alarm rings
- $P(A = 1 \mid B = 0) = 0.001$: burglar does not happen $\Rightarrow$ alarm rings
- $P(B = 1) = 0.01$: burglar rate

Question: if the alarm rang, what is the probability of a burglar happened?

$P(B = 1 \mid A = 1)$

$$= \frac{P(B = 1)P(A = 1 \mid B = 1)}{P(A = 1 \mid B = 1)P(B = 1) + P(A = 1 \mid B = 0)P(B = 0)}$$

$$= \frac{0.01 \times 0.99}{(0.01 \times 0.99) + (0.001 \times (1 - 0.01))}$$

$$\approx 0.91$$

- ▶ $P(A = 1 \mid B = 1) = 0.99$: burglar happens $\Rightarrow$ alarm rings
- ▶ $P(A = 1 \mid B = 0) = 0.001$: burglar does not happen $\Rightarrow$ alarm rings
- ▶ $P(B = 1) = 0.01$: burglar rate

Question: if the alarm rang, what is the probability of a burglar happened?

$$P(B = 1 \mid A = 1)$$
$$= \frac{P(B = 1)P(A = 1 \mid B = 1)}{P(A = 1 \mid B = 1)P(B = 1) + P(A = 1 \mid B = 0)P(B = 0)}$$
$$= \frac{0.01 \times 0.99}{(0.01 \times 0.99) + (0.001 \times (1 - 0.01))}$$
$$\approx 0.91$$

Further Question: What if $P(A = 1 \mid B = 0) = 0.01$?

## Expectation

The expectation or expected value of a function $h(x)$ with respect to a probability distribution $P(X)$ is defined as

$$E[h(x)] = \sum_x P(x)h(x) \tag{44}$$

# Expectation

The expectation or expected value of a function $h(x)$ with respect to a probability distribution $P(X)$ is defined as

$$E[h(x)] = \sum_x P(x)h(x) \tag{44}$$

## The number of ice creams [Eisenstein, 2018]

- ▶ If it is sunny, Lucia will eat four ice creams
- ▶ If it is rainy, she will eat only one ice cream
- ▶ There is a 90% chance it will be rainy

The expected number of ice creams she will eat is

$$(1 - 0.9) \times 4 + 0.9 \times 1 = 1.3 \tag{45}$$

# Mean

▶ Let $h(x) = x$, the expectation is the mean value of the random variable $X$ (discrete random variable)

$$E[X] = \sum_x x P(x) \tag{46}$$

or, (continuous random variable)

$$E[X] = \int_x x f(x) \tag{47}$$

▶ A Bernoulli distribution $P(X)$ with the parameter $\theta$, $P(X = x) = \theta^x (1 - \theta)^{(1-x)}$

$$E[X] = 1 \cdot \theta + 0 \cdot (1 - \theta) = \theta \tag{48}$$

The variance of a random variable gives a measure of how much the values of this random variable vary

$$
\begin{aligned}
\text{Var}[X] &= E\left[(X - E[X])^2\right] \\
&= E\left[X^2 - 2XE[X] + E[X]^2\right] \\
&= E[X^2] - 2E[X]E[X] + E[X]^2 \\
&= E[X^2] - E[X]^2 \tag{49}
\end{aligned}
$$

A Bernoulli distribution $P(X)$ with the parameter $\theta$,
$P(X = x) = \theta^x (1 - \theta)^{(1-x)}$

$$\text{Var}[X] = E\left[X^2\right] - E\left[X\right]^2 = p - p^2 \tag{50}$$

*Exercise*: Compute the mean and variance of a categorical distribution

# Statistical Estimation

Statistics is, in a certain sense, the inverse of probability theory.

- ▶ Observed: values of random variables
- ▶ Unknown: the model
- ▶ Task: infer the model from the observed data

For a probability $P(X; \theta)$ with $\theta$ as the unknown parameter, likelihood-based estimation with observations $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ requires two steps

1. Define a likelihood function with observations
2. Optimize the likelihood function to estimate $\theta$

## Likelihood Function

The likelihood function of $\theta$ is defined as

$$L(\theta) = \prod_{i=1}^{n} P(x^{(i)}; \theta) \tag{51}$$

Alternatively, we often use the log-likelihood function to avoid the numerical issues

$$
\begin{aligned}
\ell(\theta) &= \log L(\theta) \\
&= \sum_{i=1}^{n} \log P(x^{(i)}; \theta)
\end{aligned}
\tag{52}
$$

Maximum Likelihood Estimation: a method of estimating the parameter by maximizing the (log-)likelihood function

$$\hat{\theta} = \underset{\theta}{\mathrm{argmax}}\ \ell(\theta) \tag{53}$$

Usually, this can be done with the following equation

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^{n} \frac{\partial \log P(x^{(i)}; \theta)}{\partial \theta} = 0 \tag{54}$$

## Example: Bernoulli Distribution

Consider a Bernoulli distribution $P(X; \theta)$ with the parameter $\theta = P(X = 1; \theta)$ unknown

$$P(X = x; \theta) = \theta^x (1 - \theta)^{(1-x)} \tag{55}$$

Consider a Bernoulli distribution $P(X; \theta)$ with the parameter $\theta = P(X = 1; \theta)$ unknown

$$P(X = x; \theta) = \theta^x (1 - \theta)^{(1-x)} \tag{55}$$

With $n$ observations $\{x^{(1)}, x^{(2)}, \ldots, x^{(n)}\}$, the likelihood function is

$$
\begin{aligned}
\ell(\theta) &= \sum_{i=1}^{n} \log P(x^{(i)}; \theta) \\
&= \sum_{i=1}^{n} \{x^{(i)} \log \theta + (1 - x^{(i)}) \log(1 - \theta)\} \tag{56}
\end{aligned}
$$

The derivative with respect to $\theta$

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^{n} \{\frac{x^{(i)}}{\theta} - \frac{1 - x^{(i)}}{1 - \theta}\} \tag{57}$$

The derivative with respect to $\theta$

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^{n} \{ \frac{x^{(i)}}{\theta} - \frac{1 - x^{(i)}}{1 - \theta} \} \tag{57}$$

Let $\frac{\partial \ell(\theta)}{\partial \theta} = 0$, we have

$$\theta = \frac{\sum_{i=1}^{n} x^{(i)}}{n} \tag{58}$$

Assume the $n = 7$ observations are

$$\{0, 1, 1, 0, 0, 1, 0\},$$

then

$$\theta = \frac{3}{7} \tag{59}$$

Further Reference [Murphy, 2012, Chap 5 & 6]

## Example: Bernoulli Distribution (III)

Assume the $n = 7$ observations are

$$\{0, 1, 1, 0, 0, 1, 0\},$$

then

$$\theta = \frac{3}{7} \tag{59}$$

**Likelihood Principle**: With $x$ observed, all relevant information of inferring $\theta$ is contained in the likelihood function.

Further Reference [Murphy, 2012, Chap 5 & 6]

# Reference

Eisenstein, J. (2018).
*Natural Language Processing*.
MIT Press.

Kolter, Z. and Do, C. (2015).
Linear algebra review and reference.

Murphy, K. P. (2012).
*Machine learning: a probabilistic perspective*.
MIT press.