# CS 4774 Machine Learning

## Introduction

Yangfeng Ji

Information and Language Processing Lab
Department of Computer Science
University of Virginia

# Course Instructor and Teaching Assistants

- Instructor:
  - Yangfeng Ji
  - Office: Rice 510
- Teaching assistants
  - Oishee Hoque
  - Aparna Kishore
  - Nusrat Mozumder
  - Sanchit Sinha
  - Dane Williamson

Office hours will be released on the course webpage.

- Programming and Algorithm
  - CS 2150 or CS 3100 with a grade of C- or better

# Prerequisites

- Programming and Algorithm
  - CS 2150 or CS 3100 with a grade of C- or better
- Linear Algebra
  - Math 3350 or APMA 3080 or equivalent

# Prerequisites

- Programming and Algorithm
  - CS 2150 or CS 3100 with a grade of C- or better
- Linear Algebra
  - Math 3350 or APMA 3080 or equivalent
- Probability and Statistics
  - APMA 3100, APMA 3110, MATH 3100, or equivalent

The survey results (by Jan. 18, 12 PM)

Attempts: 83 out of 83

What is the purpose of taking this course?

Please select the answers that are aligned with your expectation.

| | | | |
|---|---|---|---|
| My research needs machine learning | 6 respondents | 7 % | |
| To learn the basic idea of machine learning | 65 respondents | 78 % | |
| To have a machine learning course on my transcript | 26 respondents | 31 % | |
| To learn machine learning tools, e.g., PyTorch, Sklearn | 61 respondents | 73 % | |
| To learn how to use machine learning solving problems | 64 respondents | 77 % | |
| Machine learning is a hot topic | 52 respondents | 63 % | |
| No Answer | 1 respondent | 1 % | |

This course will cover the basic materials on the following topics

1. Introduction to learning theory
2. Linear classification and regression
3. Model selection and validation
4. Boosting
5. Optimization methods
6. Neural networks (e.g., CNN, RNN, Auto-encoders, Transformers)

The following topics will <span style="color:magenta">not</span> be the emphasis of this course

- Statistical modeling
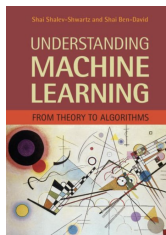  - e.g., parameter estimation, Bayesian statistics, graphical models

The following topics will not be the emphasis of this course

- ▶ Statistical modeling
  - ▶ e.g., parameter estimation, Bayesian statistics, graphical models
- ▶ Machine learning engineering
  - ▶ e.g., how to implement a classifier from end to end
  - ▶ although, we will provide some demo code for illustration purposes
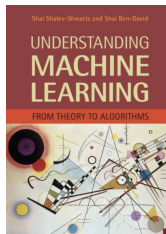
# Outline (II)

The following topics will not be the emphasis of this course

- Statistical modeling
  - e.g., parameter estimation, Bayesian statistics, graphical models
- Machine learning engineering
  - e.g., how to implement a classifier from end to end
  - although, we will provide some demo code for illustration purposes
- Advanced topics in machine learning
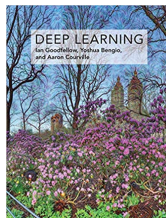  - e.g., reinforcement learning, active learning, semi-supervised learning, online learning

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

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



▶ Goodfellow, Bengio, and Courville. *Deep Learning*. 2016

For students looking for additional reading materials

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

# Homework and Grading Policy

- Homeworks (70%)
  - Five homework assignments, 14 points each

# Homework and Grading Policy

- Homeworks (70%)
  - Five homework assignments, 14 points each

- In-class Quiz (10%)
  - For the instructor to get a better understanding of students' feedback on the lectures
  - 1 point each

# Homework and Grading Policy

- Homeworks (70%)
  - Five homework assignments, 14 points each

- In-class Quiz (10%)
  - For the instructor to get a better understanding of students' feedback on the lectures
  - 1 point each

- Final project (20%)
  - There are some pre-defined problems with datasets provided
  - Students will team up (3 – 4 students per group) to solve one problem

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

- Submission will not be accepted if more than 72 hours late

- Make sure not submit wrong files
  - it is students responsbility to make sure they submit the right and complete files for each homework

- It is usually better if students just turn in what they have in time

# Violation of the Honor Code

Plagiarism, examples are

- ▶ in a homework submission, copying answers from others directly or some minor changes

- ▶ in a report, copying texts from a published paper (including, some minor changes)

- ▶ in a code, using someone else's functions/implementations without acknowledging the contribution

# Webpages

- Course webpage
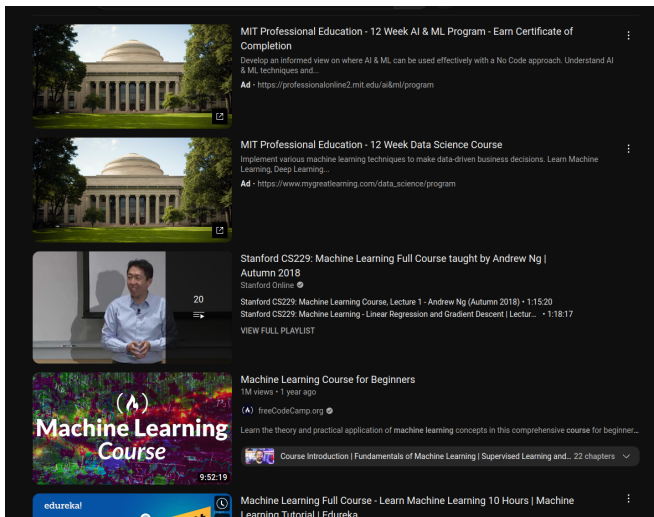
  `http://yangfengji.net/uva-ml-undergrad/`

  which contains all the information you need about this course.

- Canvas
  - For homework releasing and grading
  - For announcement, online QA, discussion, etc.

# Why Taking this Course?

# Machine Learning Courses



14

Building logistic regression classifier

# Another Example

Building a GPT model with Transformer

# What is Missing?

# What is this?



**sklearn.linear_model.LogisticRegression**

*class* sklearn.linear_model.**LogisticRegression**(*penalty='l2', *, dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0, warm_start=False, n_jobs=None, l1_ratio=None*)                                                                    [source]

Logistic Regression (aka logit, MaxEnt) classifier.

In the multiclass case, the training algorithm uses the one-vs-rest (OvR) scheme if the 'multi_class' option is set to 'ovr', and uses the cross-entropy loss if the 'multi_class' option is set to 'multinomial'. (Currently the 'multinomial' option is supported only by the 'lbfgs', 'sag', 'saga' and 'newton-cg' solvers.)

This class implements regularized logistic regression using the 'liblinear' library, 'newton-cg', 'sag', 'saga' and 'lbfgs' solvers. **Note that regularization is applied by default**. It can handle both dense and sparse input. Use C-ordered arrays or CSR matrices containing 64-bit floats for optimal performance; any other input format will be converted (and copied).

The 'newton-cg', 'sag', and 'lbfgs' solvers support only L2 regularization with primal formulation, or no regularization. The 'liblinear' solver supports both L1 and L2 regularization, with a dual formulation only for the L2 penalty. The Elastic-Net regularization is only supported by the 'saga' solver.

▶ What's the definition of this classifier?
▶ What if it does not work?
▶ What are its limitations?

# What is this?



**sklearn.linear_model.LogisticRegression**

*class* sklearn.linear_model.**LogisticRegression**(*penalty='l2', *, dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0, warm_start=False, n_jobs=None, l1_ratio=None*)                                                                              [source]

Logistic Regression (aka logit, MaxEnt) classifier.

In the multiclass case, the training algorithm uses the one-vs-rest (OvR) scheme if the 'multi_class' option is set to 'ovr', and uses the cross-entropy loss if the 'multi_class' option is set to 'multinomial'. (Currently the 'multinomial' option is supported only by the 'lbfgs', 'sag', 'saga' and 'newton-cg' solvers.)

This class implements regularized logistic regression using the 'liblinear' library, 'newton-cg', 'sag', 'saga' and 'lbfgs' solvers. **Note that regularization is applied by default**. It can handle both dense and sparse input. Use C-ordered arrays or CSR matrices containing 64-bit floats for optimal performance; any other input format will be converted (and copied).

The 'newton-cg', 'sag', and 'lbfgs' solvers support only L2 regularization with primal formulation, or no regularization. The 'liblinear' solver supports both L1 and L2 regularization, with a dual formulation only for the L2 penalty. The Elastic-Net regularization is only supported by the 'saga' solver.

▶ What's the definition of this classifier?
▶ What if it does not work?
▶ What are its limitations?

*In fact, if you explain how these parameters work and their effects, you can skip at least one third of the class lectures.*

18

To understanding machine learning, we need some mathematical and statistical knowledge

Attempts: 84 out of 84

The course materials also contain a large amount of mathematical/statistical stuff.

Please select the following answer that gives the closest description of how comfortable when you read mathematics.

| | | |
|---|---|---|
| No way for me to read any mathematics | 1 respondent | 1 % |
| Not a fan of mathematics, but I can handle it | 31 respondents | 37 % |
| I will try to avoid it, unless I have to | 9 respondents | 11 % |
| **I feel comfortable of reading mathematics** | **43 respondents** | **51 %** |

Now, let's have some fun!

Warning: you will see lots of mathematical notations.

# Basic Linear Algebra

# Linear Equations

Consider the following system of equations

$$4x_1 - 5x_2 = -13$$
$$-2x_1 + 3x_2 = 9$$

(1)

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

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

(2)

with

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

(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

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



✎ *Question for Homework*: 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)

## 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 $x$

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 $\tilde{x} = \tilde{y}$.

Given a matrix $A$ and a vector $x$, their multiplication is equivalent to performing a linear transformation on $x$

$$Ax \tag{9}$$

For example, consider the following matrix

$$A = \begin{bmatrix} 0.5 & 0 \\ 0 & 2 \end{bmatrix} \tag{10}$$

and three vectors

- $x_1^\mathsf{T} = [1, 2]$
- $x_2^\mathsf{T} = [2, 4]$
- $x_3^\mathsf{T} = [3, 6]$

Given a matrix $A$ and a vector $x$, their multiplication is equivalent to performing a linear transformation on $x$

$$Ax \tag{9}$$

For example, consider the following matrix

$$A = \begin{bmatrix} 0.5 & 0 \\ 0 & 2 \end{bmatrix} \tag{10}$$

and three vectors

- $x_1^\mathsf{T} = [1, 2]$
- $x_2^\mathsf{T} = [2, 4]$
- $x_3^\mathsf{T} = [3, 6]$

➡ *This is also what the function torch.nn.Linear means*

## 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.

$$\mathbf{I} = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \tag{11}$$

# 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.

$$\mathbf{I} = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \tag{11}$$

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

$$\mathbf{D} = \begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_n \end{bmatrix} \tag{12}$$

# Inverse

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

$$A^{-1}A = I = AA^{-1} \qquad (13)$$

▶ 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 $x = \mathbf{A}^{-1}b$

## Inverse (II)

- In matrix-vector multiplication, an inverse matrix $A^{-1}$ will reverse the linear transformation performed by $A$

$$A^{-1}Ax = x \qquad (14)$$

▶ In matrix-vector multiplication, an inverse matrix $A^{-1}$ will reverse the linear transformation performed by $A$

$$A^{-1}Ax = x \tag{14}$$

▶ When a matrix $A$ is not invertible, it means its linear transformation is not reversible

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \tag{15}$$

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

# 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)

$$\langle u_i, u_j \rangle = 0, \|u_i\| = 1, \|u_j\| = 1 \tag{16}$$

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

# 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)

$$\langle u_i, u_j \rangle = 0, \|u_i\| = 1, \|u_j\| = 1 \tag{16}$$

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 Special Case

Consider a matrix $A$ as

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \tag{17}$$

For any $x^\mathsf{T} = [x_1, x_2]$, we have

$$A \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix} \tag{18}$$

Consider a matrix $A$ as

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \tag{17}$$

For any $x^\mathsf{T} = [x_1, x_2]$, we have

$$A \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix} \tag{18}$$

The is a reflection operation. Operations like this are popularly used in computer graphics.

# Symmetric Matrices

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

$$\mathbf{A}^\top = \mathbf{A} \tag{19}$$

or, in other words,

$$a_{i,j} = a_{j,i} \quad \forall i, j \in [n] \tag{20}$$

Comments

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

$$x^\top A y \tag{21}$$

gives each dimension a different weight (importance) when computing the similarity between $x$ and $y$

# Eigen Decomposition

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

$$\mathbf{A} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^{\mathsf{T}} \tag{22}$$

with

▶ $\boldsymbol{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$ as a diagonal matrix

▶ $\mathbf{Q}$ is an orthogonal matrix

## Eigen Decomposition

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

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\mathsf{T} \tag{22}$$

with

- $\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$ as a diagonal matrix

- $\mathbf{Q}$ is an orthogonal matrix

- Consider the similarity measurement

$$x^\mathsf{T}\mathbf{A}y = x^\mathsf{T}\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\mathsf{T}y = (\mathbf{U}^\mathsf{T}x)^\mathsf{T}\mathbf{\Lambda}\mathbf{U}^\mathsf{T}y \tag{23}$$

# Eigen Decomposition

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

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\mathsf{T} \tag{22}$$

with

- $\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$ as a diagonal matrix
- $\mathbf{Q}$ is an orthogonal matrix
- Consider the similarity measurement

$$x^\mathsf{T}\mathbf{A}y = x^\mathsf{T}\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\mathsf{T}y = (\mathbf{U}^\mathsf{T}x)^\mathsf{T}\mathbf{\Lambda}\mathbf{U}^\mathsf{T}y \tag{23}$$

✎ *Question for Homework*: if a symmetric matrix $\mathbf{A}$ is invertible, show $\mathbf{A}^{-1} = \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{U}^\mathsf{T}$ with $\Lambda^{-1} = \text{diag}(\frac{1}{\lambda_1}, \ldots, \frac{1}{\lambda_n})$

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

$$x^{\top}\mathbf{P}x \geq 0 \qquad (24)$$

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

$$x^\top \mathbf{P} x \geq 0 \tag{24}$$

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

Eigen decomposition of $\mathbf{P}$ as

$$\mathbf{P} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top \tag{25}$$

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

$$\lambda_i \geq 0 \tag{26}$$

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

$$x^\mathsf{T} \mathbf{P} x > 0 \qquad (27)$$

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

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

$$\lambda_i > 0 \qquad (28)$$

# Symmetric Positive Definite Matrices

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

$$x^{\mathsf{T}} \mathbf{P} x > 0 \tag{27}$$

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

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

$$\lambda_i > 0 \tag{28}$$

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

## Review

The identity matrix $\mathbf{I}$ is

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

The identity matrix $\mathbf{I}$ is

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

# Review of 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:

- **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 $\{0, 1\}$

- **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

$$P(X = 1) = 0.52 \tag{29}$$

# Bernoulli Distribution

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

# Tossing a Coin Twice?

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

# 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 $X$
  - $P(X = 0) = (1 - \theta)^2$

# 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

$$P(Y = k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \tag{30}$$

where

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

is the binomial coefficient and

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

# Tossing a Dice



How to define the corresponding random variable?

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

# Categorical Distribution

$$P(\boldsymbol{X} = \boldsymbol{x}) = \prod_{k=1}^{6}(\theta_k)^{x_k} \tag{31}$$

where

- $\boldsymbol{x} = (x_1, x_2, \ldots, x_6)$
- $x_k \in \{0, 1\}$, 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

$$P(X = x) = \binom{n}{x_1 \cdots x_K} \prod_{k=1}^{K} \theta_k^{x_k} \qquad (32)$$

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:

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

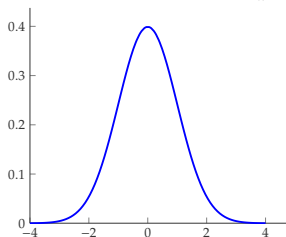# Gaussian Distribution

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

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \tag{33}$$

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

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \tag{34}$$

There examples of Gaussian distributions



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

# Gaussian Distribution (II)

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \tag{34}$$

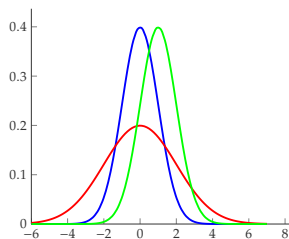There examples of Gaussian distributions



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

✎ *Question for Homework*: describe a random event with the probabilistic language

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

$$P(X, Y) \tag{35}$$

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,

$$P(X, Y) = P(X) \cdot P(Y) \tag{36}$$

Marginal distributions

$$P(X) = \sum_Y P(X, Y) \tag{37}$$

$$P(Y) = \sum_X P(X, Y) \tag{38}$$

## 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,

$$P(X, Y) = P(X) \cdot P(Y) \tag{36}$$

Marginal distributions

$$P(X) = \sum_Y P(X, Y) \tag{37}$$

$$P(Y) = \sum_X P(X, Y) \tag{38}$$

- $X$: 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 $X$

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

Example: document classification

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

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

$$P(Y \mid X) = P(Y) \tag{40}$$

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

# Conditional Probability

- $X$: 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    |

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

# Conditional Probability

▶ $X$: 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    |

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

# Conditional Probability

- $X$: 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    |

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

✎ *Question for Homework*: compute conditional probability from a given probabilistic table

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

$$f(x) = \frac{1}{(2\pi)^{n/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(x-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(x-\boldsymbol{\mu})\right) \tag{41}$$
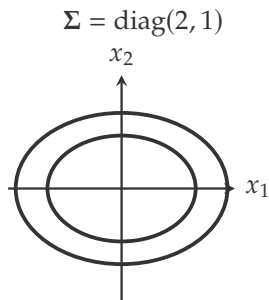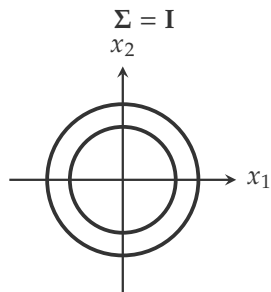
where

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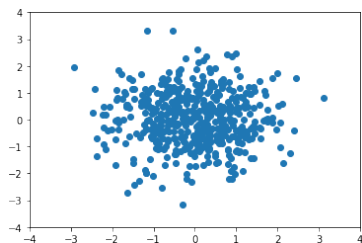- $\boldsymbol{\Sigma}$ is the $n \times n$ covariance matrix.

# Covariance Matrix $\Sigma$

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

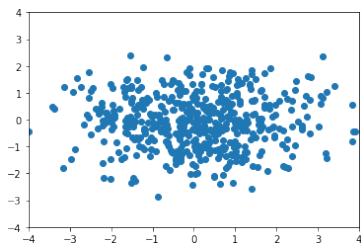$$f(x) \propto \exp\left(-\frac{1}{2}x^{\mathsf{T}}\Sigma^{-1}x\right) \tag{42}$$

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



$\Sigma = \mathbf{I}$

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

# Sampling from Gaussians



(a)  (b)

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

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

✎ *Question for Homework*: sample from a Gaussian distribution with a pre-defined mean and variance

# Thank You!