

Mathematics of Big Data, I

Lecture 1: Introduction of Big Data & Overview of Big Data Analytics

Weiqing Gu

Professor of Mathematics
Director of the Mathematics Clinic

Harvey Mudd College
Summer 2021

<https://math189bigdata.github.io/>

Textbook:

All members of the class will be required to obtain the following text:

Kevin Patrick Murphy, ***Machine Learning: a Probabilistic Perspective***. MIT Press, 2012.

Grading:

- 5% Reading Summary
- 35% Homework
- 20% Midterm (Project or Exam, TBD)
- 40% Final Project
- [Up to 5% Extra Credit]

Course Requirements and Evaluation:

- ***Reading Presentations***

All readings are compulsory, but some are more compulsory than others.

To encourage the goal of reading active research in the field, we will assign each non-Murphy reading to a group of two students who will write a summary of 1-2 pages to be turned in at the start of class. Each student will do approxiamately two summaries in total. They must be clear and demonstrate that you have read the paper with a high degree of confidence. Credit will be given on a 0-10 scale for each summary. Your summary should be done at a high level, and should focus on the main point of the readings (i.e. avoid complicated math). As long as your summary is reasonable, you will be given full credit.

Homework!

- ***Homework***

The homework is due every week at the beginning of each lecture. There will be two parts for each assignment: math and coding. The homework is split approximately evenly between mathematical analysis and extension of our course material and application of algorithms to real world data.

For coding: You are highly recommended to use Python3. For each problem, the starter code and the sample solution are implemented in Python3. All the results and graphs for the sample solutions were produced under Python 3.5.2 under macOS Sierra; different versions of Python or system environment may produce different results. You are also welcome to use Jupyter Notebooks, but the starter code is not provided in notebook format.

Numpy and **Pandas** are two important python libraries to know for coding assignment for this course. You might also want to look at **Matplotlib** for generating plots. If you never used these libraries before, make sure you check out the tutorials online before starting the first assignment.

Example of Data Set

[https://physionet.org/content/mimiciii
-demo/1.4/](https://physionet.org/content/mimiciii-demo/1.4/)

Abstract

MIMIC-III is a large, freely-available database comprising deidentified health-related data associated with over 40,000 patients who stayed in critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012 [1]. The MIMIC-III Clinical Database is available on PhysioNet (doi: [10.13026/C2XW26](https://doi.org/10.13026/C2XW26)). Though deidentified, MIMIC-III contains detailed information regarding the care of real patients, and as such requires credentialing before access. To allow researchers to ascertain whether the database is suitable for their work, we have manually curated a demo subset, which contains information for 100 patients also present in the MIMIC-III Clinical Database. Notably, the demo dataset does not include free-text notes.

- At the end of each lecture, the head grader will give you some instruction on how to start to write your code and what would be some of the expected challenges for the next coding assignment.

Note:

- 1) When doing the coding problem for each homework set, you are not allowed to use any machine learning algorithms implemented by external libraries, such as LinearRegression in sklearn. However, you may use these algorithms in your final project.**

- 2) Each homework has both pdf and tex versions. To have the tex files successfully compiled, make sure that you have downloaded both macros.tex and hmcpset.cls and put them and the hw tex file under same folder.
If you have any questions with regard to the compilation of the tex files, feel free to ask the grutors for help.**

- 3) For each coding problem, please submit your code to GitHub; please print out any graph or printing statements and submit them with the written part.**

Exams

- *Midterm*

The midterm will either be a take-home exam covering all topics seen in the first week of the course or a project where you will apply the methods learned in the first half of the course (TBD).

- *Final Project*

The final project is the largest component of the course. Each student will discover, explore, and attack a real world problem of your choosing. The detailed description and requirements for the final project can be found under the "Final Project" tab.

- *GitHub*

As we stated in the course overview, students are expected to become comfortable with Github. Hence, each student is required to create a Github account for coding assignment submission and final project submission. If you already have a Github account, that's perfect. If not, please create a personal Github account and go over the tutorials online.

Note: Please make sure to send the username of your Github account to TA for homework grading.

Classroom Policies:

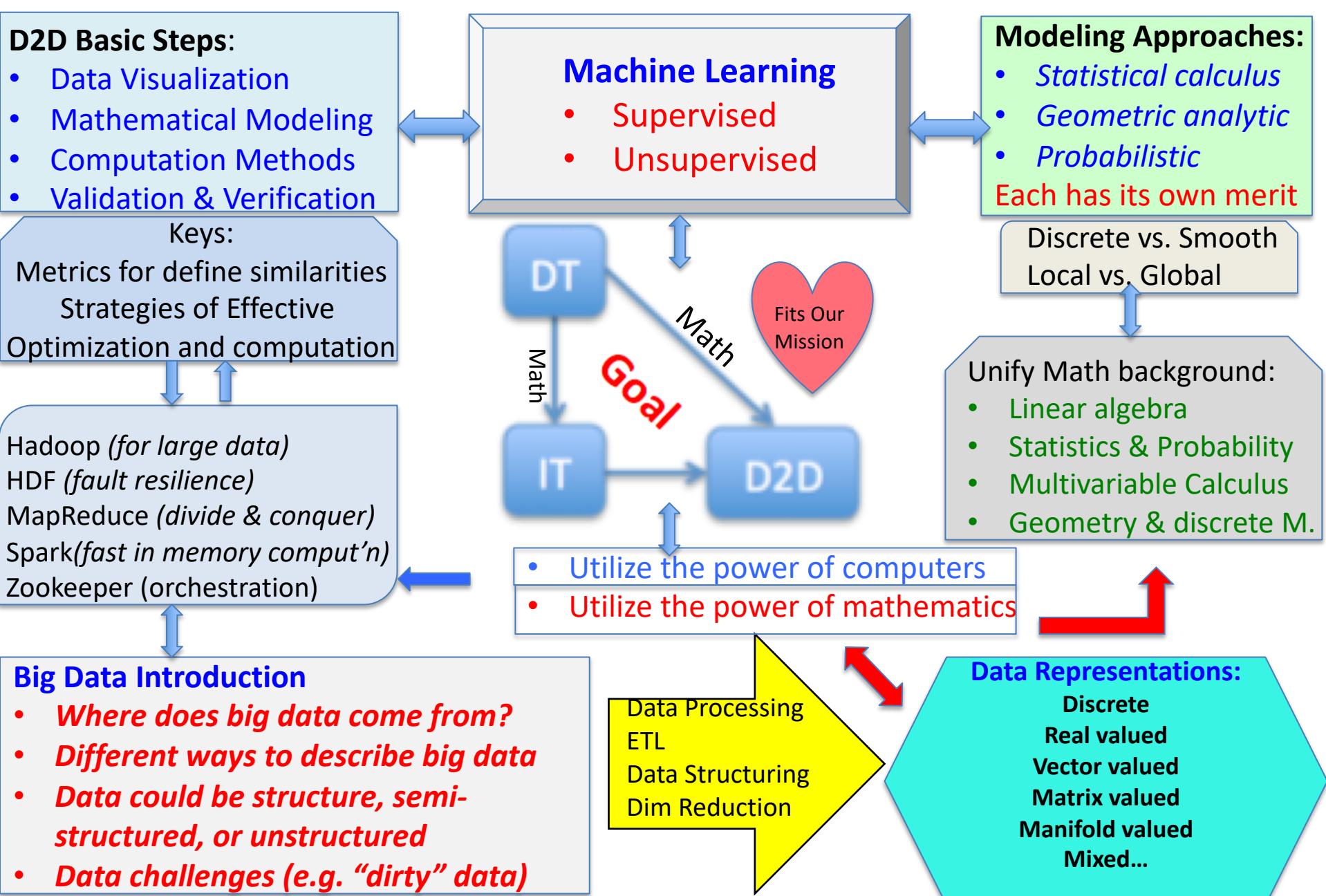
- ***Attendance***

Attendance for each lecture is mandatory and is expected of all class members. If you're going to miss a lecture, it is necessary for you to inform the instructor as soon as possible. You are also responsible for obtaining notes from another class member.

- ***Devices***

You are welcome to use your computer or tablet for note-taking (the PowerPoint slides will also be posted shortly after the lecture for your convenience).

A Big Picture of Mathematics of Big Data, I



Today's Lecture

- First: Big data introduction (answer first two questions)
 - Big Data Introduction
 - *Where does big data come from?*
 - *Different ways to describe big data*
- Second: Use linear regression as an example to give an overview of big data analytics

Modeling Approaches:

- *Statistical calculus*
- *Geometric analytic*
- *Probabilistic*

Each has its own merit

- Note:

*Mathematics of Big Data (in academic) ==
Big Data Analytics (in industry).*

First: Introduction of Big Data

- *Where does big data come from?*

Organizations

Machines

People

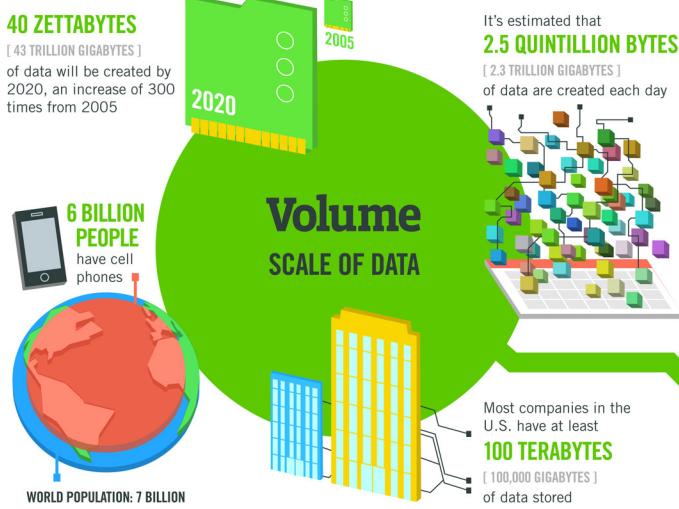
Data is not new. But the scale has been changed!
The way how people using data has been transformed!

Types of big data

1. Structured data (e.g. often Generated by organizations)
2. Semi-structured data (e.g. Generated by machine with manual records)
3. Unstructured data (often Generated by people)

• What exactly is big data?

- Does “big” here mean “big volume”?
- In fact, there are 5 “V”s to describe big data.
 - **Volume (Size)**
 - **Velocity (Speed)**
 - **Variety (Types)**
 - **Veracity (Quality)**
 - **Valence (Relationships)**



By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** – almost 2.5 connections per person on earth

The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015 **4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES [161 BILLION GIGABYTES]



By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

4 BILLION+ HOURS OF VIDEO are watched on YouTube each month



400 MILLION TWEETS are sent per day by about 200 million monthly active users

Poor data quality costs the US economy around **\$3.1 TRILLION A YEAR**



1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions

27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate

Veracity UNCERTAINTY OF DATA

IBM

Data to Decision (D2D)

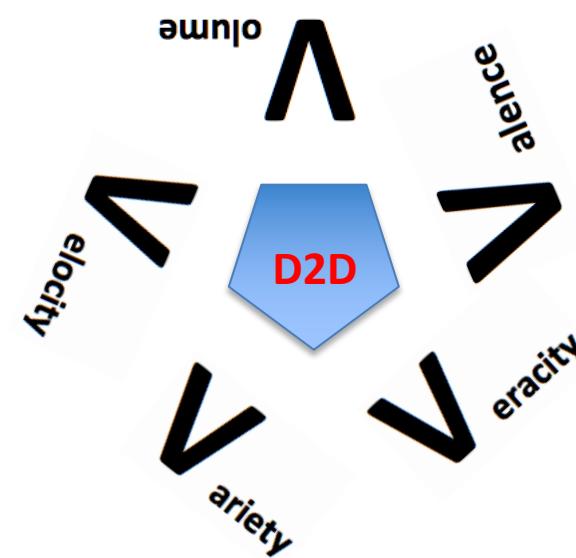
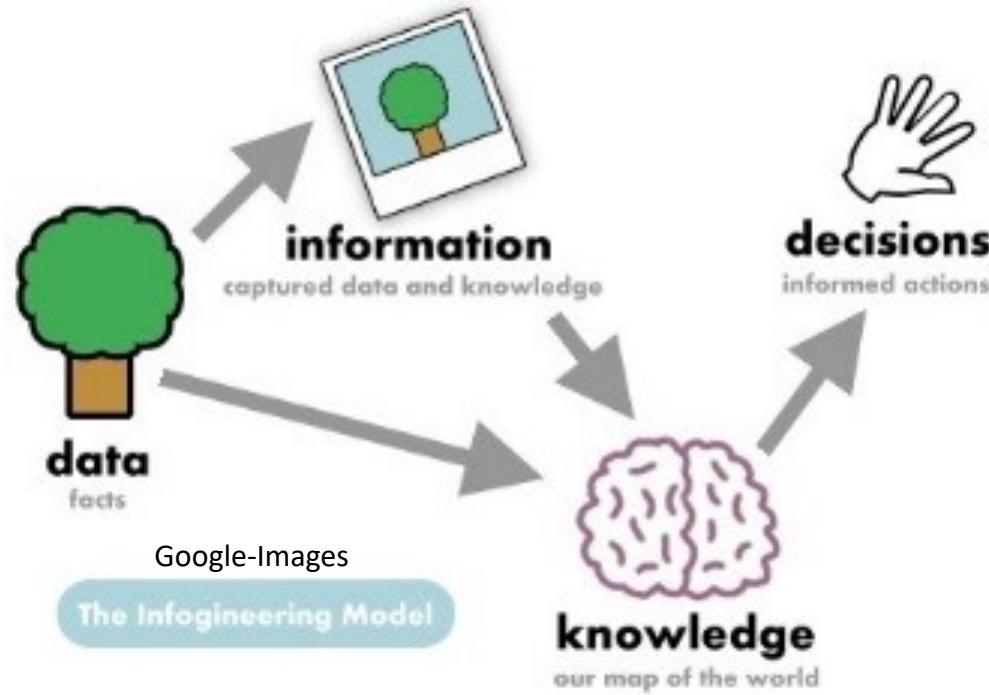
V
olume

V
elocity

V
ariety

V
eracity

V
alence



Second for today: Analytic Approaches

- Use “linear regression” as an example to give an overview of big data analytics

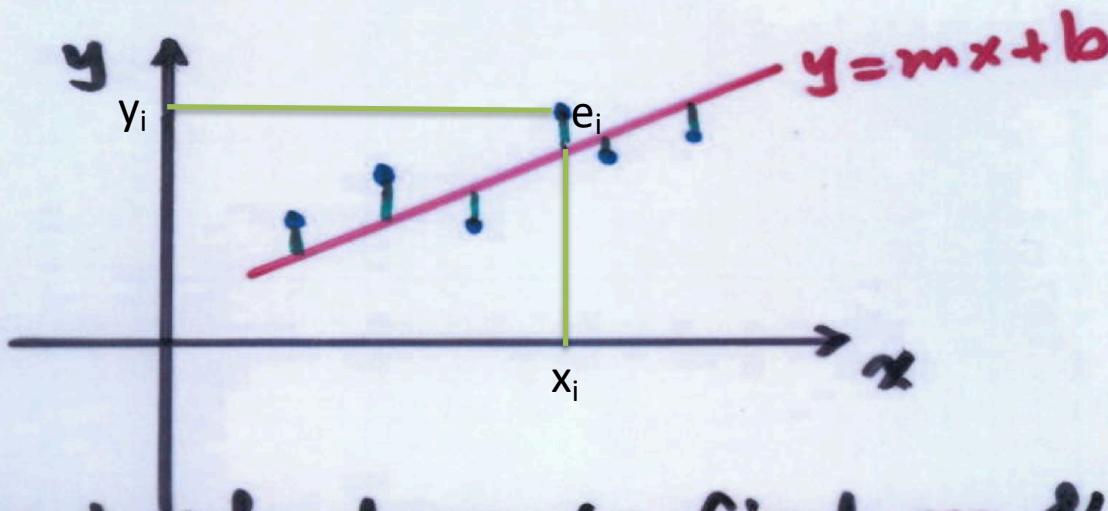
Modeling Approaches:

- *Statistical calculus*
- *Geometric analytic*
- *Probabilistic*

Each has its own merit

1. Statistical Calculus Approach (Classical Least Square Approximation)

Suppose we have data pts (x_i, y_i) and want to find the line $y = mx + b$ which best describes the data.



The problem boils down to find m & b .

The error between one point and the line is

$$e_i = y_i - (mx_i + b)$$

Our objective is minimizing the total error.

- However, the errors e_i , some could be positive and some could be negative. A simple sum of the errors would not work well.
- Can you think about an example why not working well?
- How to fix this problem?
- Instead we consider the following **objective or cost function**:
- $J(m,b) = \sum (e_i)^2 = \sum (y_i - mx_i - b)^2$
- Can we use $\sum |e_i|$ instead?

L₂ norm

$$\begin{aligned}Y &= mx + b \\&= (b, m) \cdot (1, x)\end{aligned}$$

L₁ norm

$$\begin{aligned}\text{In general, } Y &= b + m_1x_1 + m_2x_2 + \dots + m_kx_k \\&= (b, m_1, m_2, \dots, m_k) \cdot (1, x_1, x_2, \dots, x_k)\end{aligned}$$

Goal: Find m and b to minimize the cost function J

- How?
- Set all partials equal to zero!
- Work out the details with the students on the board.

Obtained solution using Cramer's rule

- Give a linear system:

$$\begin{cases} a_1x + b_1y = c_1 \\ a_2x + b_2y = c_2 \end{cases}$$

- Write it into matrix form:

$$\begin{bmatrix} a_1 & b_1 \\ a_2 & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

**Assume the coefficient matrix is invertible,
i.e. the $\det = a_1b_2 - b_1a_2$ is nonzero.** Then

$$x = \frac{\begin{vmatrix} c_1 & b_1 \\ c_2 & b_2 \end{vmatrix}}{\begin{vmatrix} a_1 & b_1 \\ a_2 & b_2 \end{vmatrix}} = \frac{c_1b_2 - b_1c_2}{a_1b_2 - b_1a_2}, \quad y = \frac{\begin{vmatrix} a_1 & c_1 \\ a_2 & c_2 \end{vmatrix}}{\begin{vmatrix} a_1 & b_1 \\ a_2 & b_2 \end{vmatrix}} = \frac{a_1c_2 - c_1a_2}{a_1b_2 - b_1a_2}.$$

Close formula for Least Square Approximation

Using Cramer's rule, we get solution for m, b :

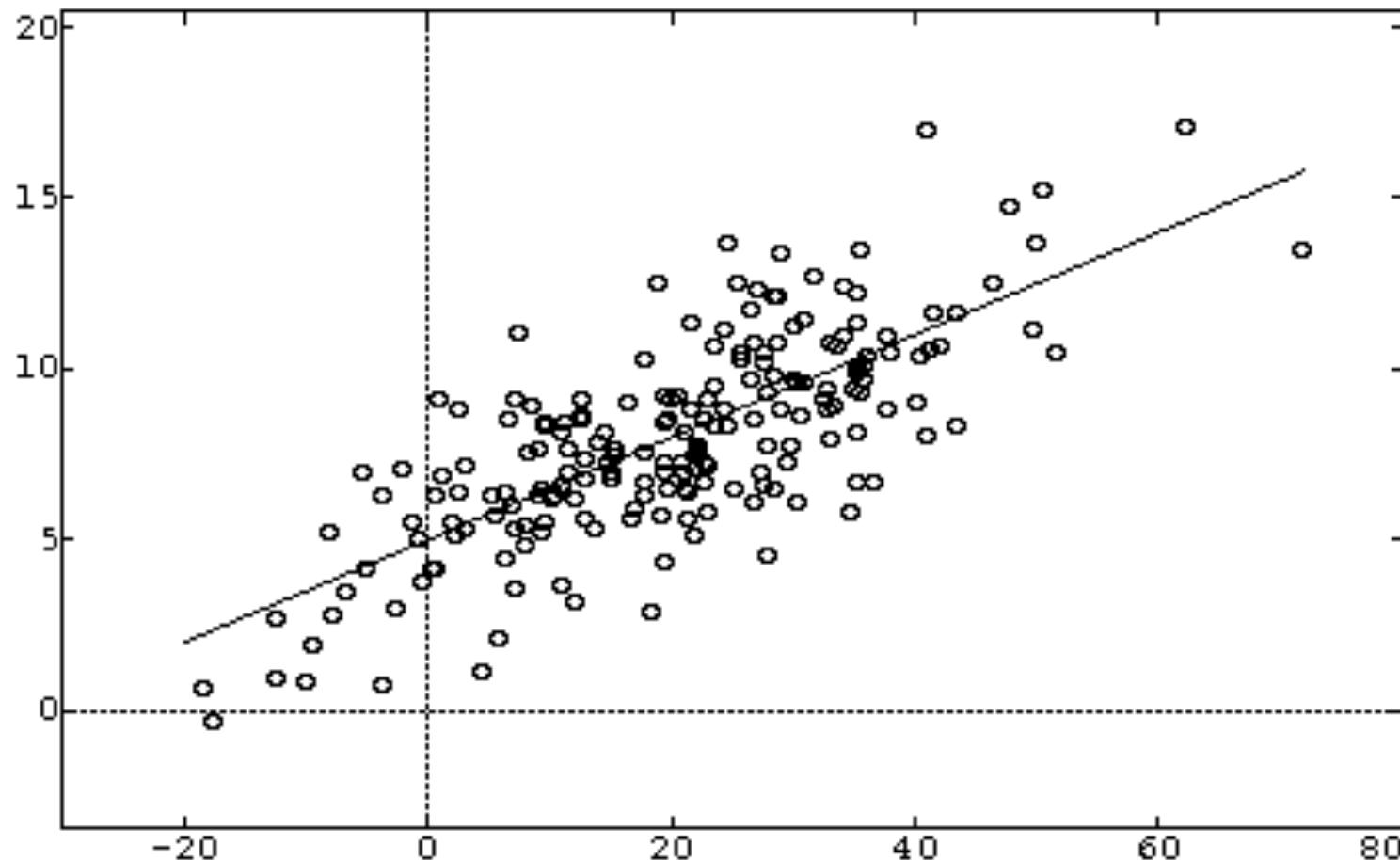
$$m = \frac{n \sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

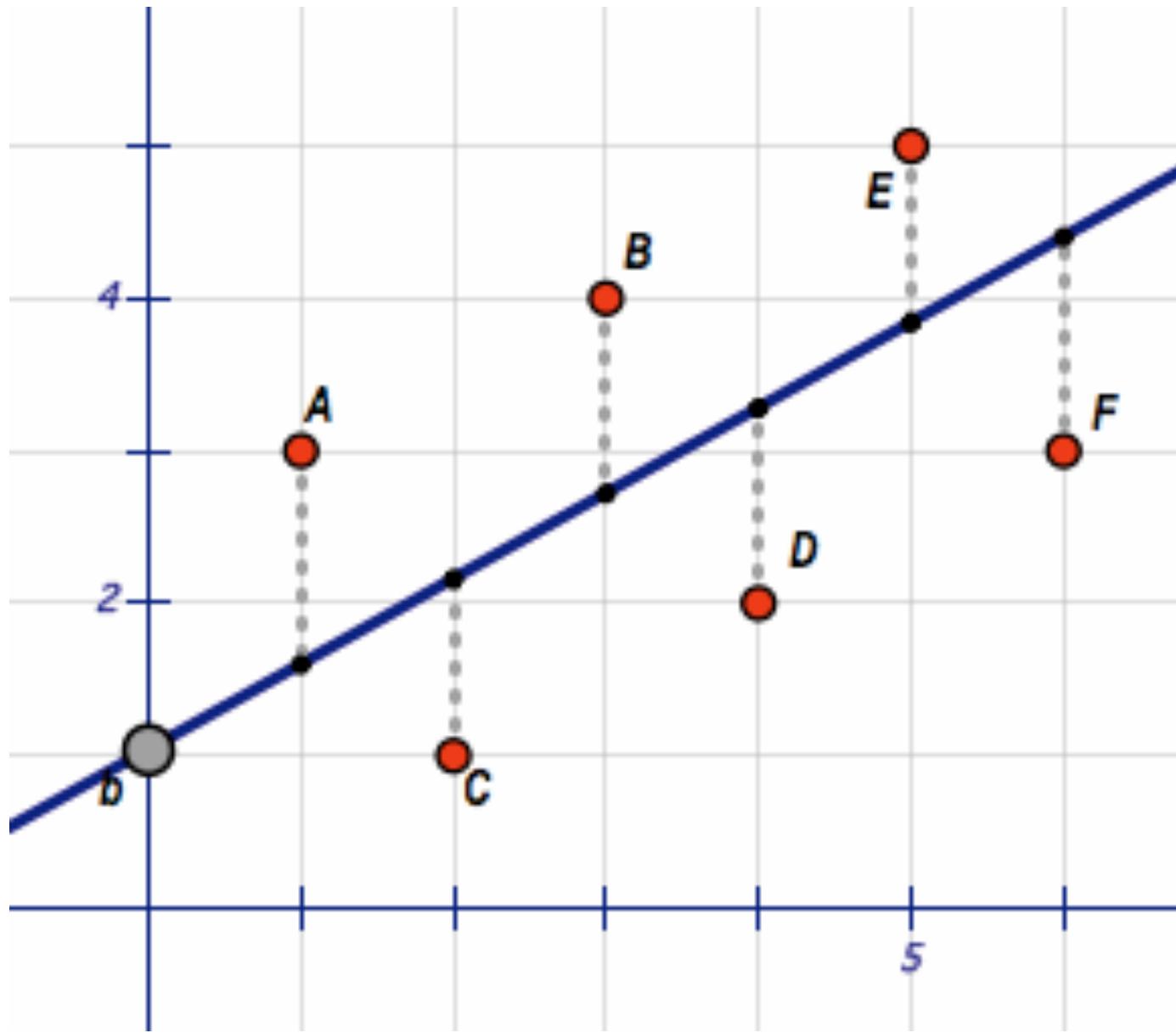
$$b = \frac{(\sum_{i=1}^n x_i^2)(\sum_{i=1}^n y_i) - (\sum_{i=1}^n x_i)(\sum_{i=1}^n x_i y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

But the formula is massy. Next we'll find a compact form of this formula.

Linear Regression

Given some data: $D = \{x_i, y_i\}$





Normal Equation for Least Square Approximation

- i.e. Representing the Least Square Solution in Matrix Form
- Work out the details with the students on the board.
- Recall the product rule:
 - $f, g: \mathbb{R} \rightarrow \mathbb{R}$: $(f \cdot g)' = f' \cdot g + f \cdot g'$
 - $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$: $\nabla(f \cdot g) = \nabla f \cdot g + f \cdot \nabla g$
 - $\mathbf{f}, \mathbf{g}: \mathbb{R}^n \rightarrow \mathbb{R}^n$: $(\mathbf{f} \cdot \mathbf{g})' = \mathbf{f}' \cdot \mathbf{g} + \mathbf{f} \cdot \mathbf{g}'$

$$\theta = (X^T X)^{-1} X^T \vec{y}.$$

Homework problem

- Given 4 points as below:

$(0, 1), (2, 3), (3, 6), (4, 8)$

- a) Find $y = mx + b$ based on Cramer's rule.

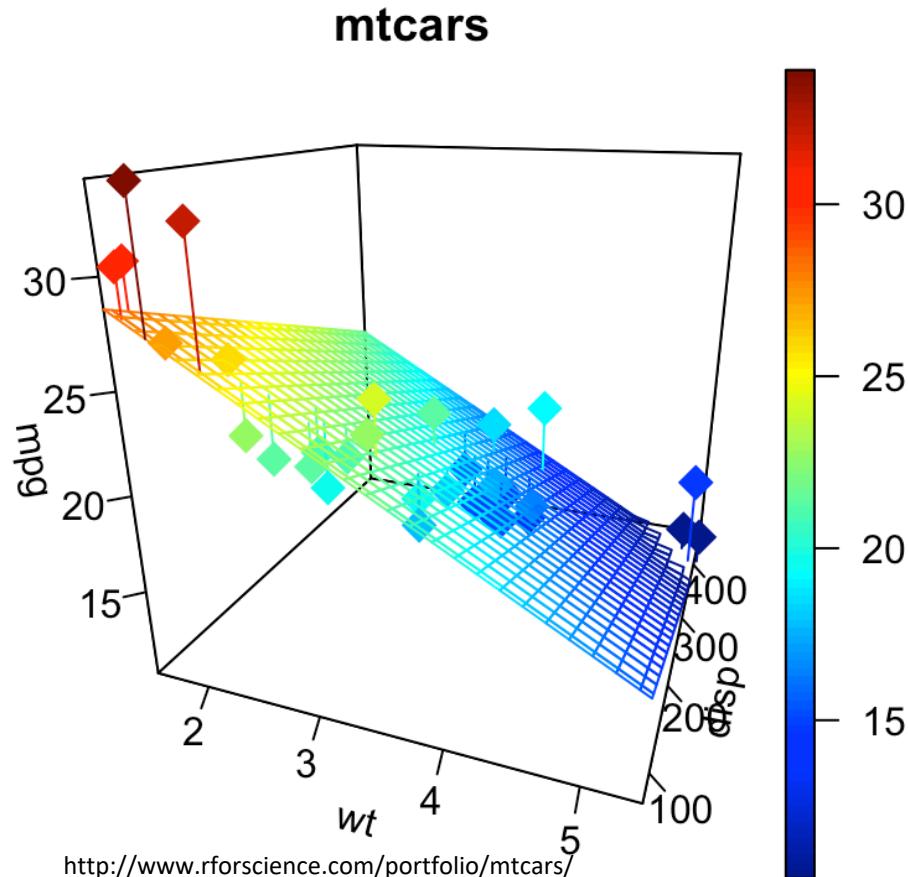
- Hint:

x_i	y_i	\bar{x}_i	$x_i y_i$
0	1	0	0
2	3	4	6
3	6	9	18
4	8	16	32

$$\sum x_i = 9 \quad \sum y_i = 18 \quad \sum x_i^2 = 29 \quad \sum x_i y_i = 56$$

- b) Use the normal formula to find the solution and compare it with that of a).
- c) Plot the data points, and draw $y = mx + b$.
- d) (All by coding) Find another 100 points near the line $y = mx + b$. Then find the least square approxim'n again & plot both the data points & the new line.

How about fit data by a plane?



Please read or review the concept in probability

Random variable

From Wikipedia, the free encyclopedia

In [probability](#) and [statistics](#), a **random variable**, **random quantity**, **aleatory variable**, or **stochastic variable** is described informally as a [variable whose values depend on outcomes of a random phenomenon](#).^[1] The formal mathematical treatment of random variables is a topic in [probability theory](#). In that context, a random variable is understood as a [measurable function](#) defined on a [probability space](#) that maps from the [sample space](#) to the [real numbers](#).^[2]

A random variable's possible values might represent the possible outcomes of a yet-to-be-performed experiment, or the possible outcomes of a past experiment whose already-existing value is uncertain (for example, because of imprecise measurements or [quantum uncertainty](#)).^[1] They may also conceptually represent either the results of an "objectively" random process (such as rolling a die) or the "subjective" randomness that results from incomplete knowledge of a quantity. The meaning of the probabilities assigned to the potential values of a random variable is not part of probability theory itself, but is instead related to philosophical arguments over the [interpretation of probability](#). The mathematics works the same regardless of the particular interpretation in use.

As a function, a random variable is required to be [measurable](#), which allows for probabilities to be assigned to sets of its potential values. It is common that the outcomes depend on some physical variables that are not predictable. For example, when tossing a fair coin, the final outcome of heads or tails depends on the uncertain physical conditions, so the outcome being observed is

Part of a series on [statistics](#)
Probability theory



[Probability](#) · [Probability axioms](#) · [Determinism](#) · [Indeterminism](#) · [Randomness](#)
[Probability space](#) · [Sample space](#) · [Event](#) · [Collectively exhaustive events](#) · [Elementary event](#) · [Mutual exclusivity](#) · [Outcome](#) · [Singleton](#) · [Experiment](#) · [Bernoulli trial](#) · [Probability distribution](#) · [Bernoulli distribution](#) · [Binomial distribution](#) · [Normal distribution](#) · [Probability measure](#) · **Random variable** ([Bernoulli process](#) · [Continuous or discrete](#) · [Expected value](#) · [Markov chain](#) · [Observed value](#) · [Random walk](#) · [Stochastic process](#)) · [Complementary event](#) · [Joint probability](#) ·

Get the same close solution by normal equation!

- Can you imagine what other cases you would get the same kind of solution?

2. Geometric Analytic Approach (Geometric Least Square)

- Work out the details with the students on the board.

Assume a linear model

$$\begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ \mathbf{r}_n \end{pmatrix} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{pmatrix} - \begin{pmatrix} \mathbf{x}_{11} & \dots & \mathbf{x}_{1m} \\ \vdots & \vdots & \vdots \\ \mathbf{x}_{n1} & \dots & \mathbf{x}_{nm} \end{pmatrix} \begin{pmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_m \end{pmatrix}$$

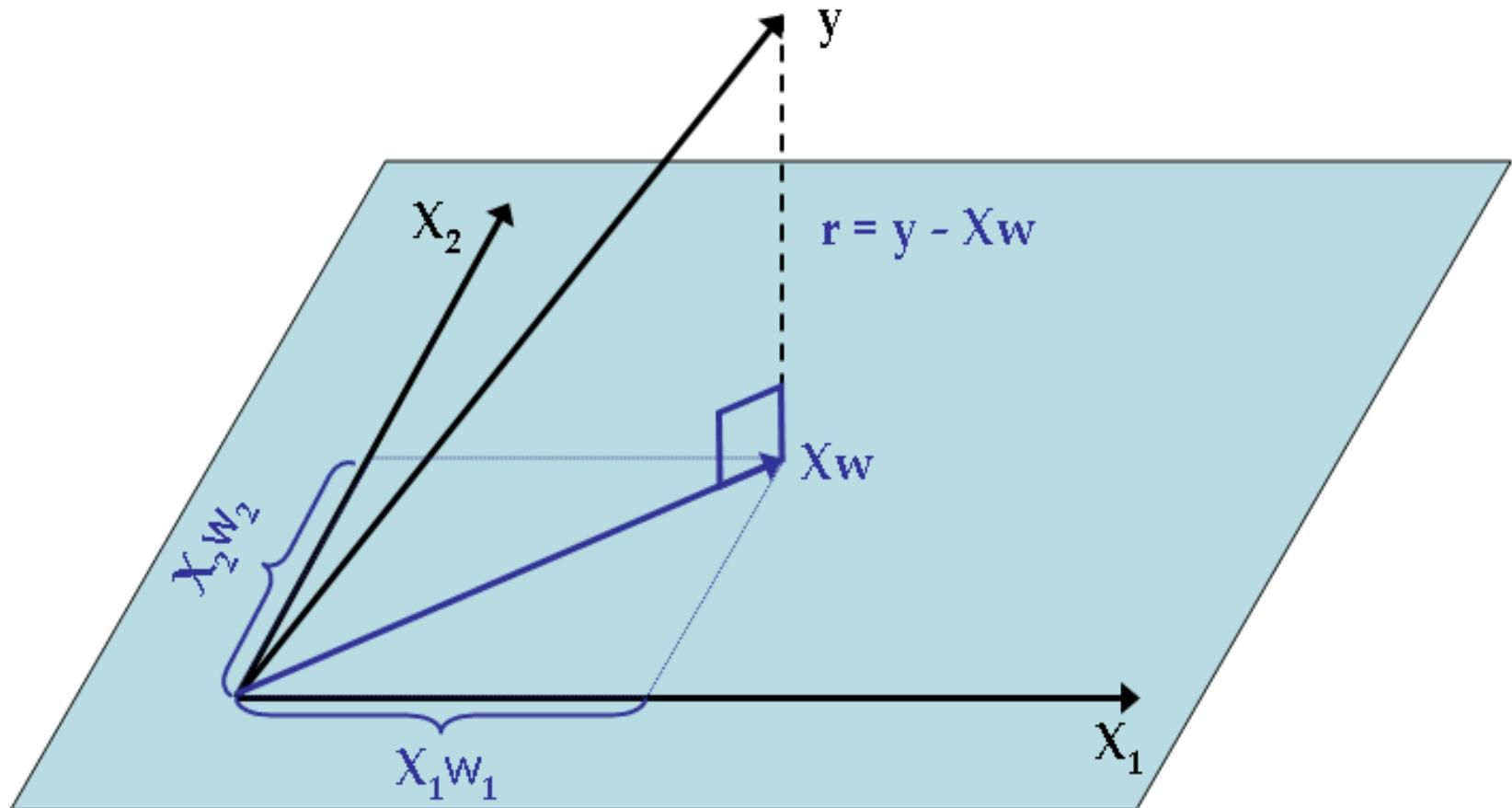
$$\rightarrow \mathbf{r} = (\mathbf{y} - \mathbf{X}\mathbf{w})$$

This is equivalent to

$$\mathbf{y}_i = \sum_j w_j x_{ij} + \mathcal{N}(0, \sigma^2) = \mathbf{x}_i \mathbf{w} + \mathcal{N}(0, \sigma^2)$$

Key in *Geometric* Least Square Approximation

Geometrically you can see the solution!



$$\mathbf{w}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Again we get the same solution!

$$\theta = (X^T X)^{-1} X^T \vec{y}.$$

Q: But what's wrong if we use Cramer's rule to solve it?

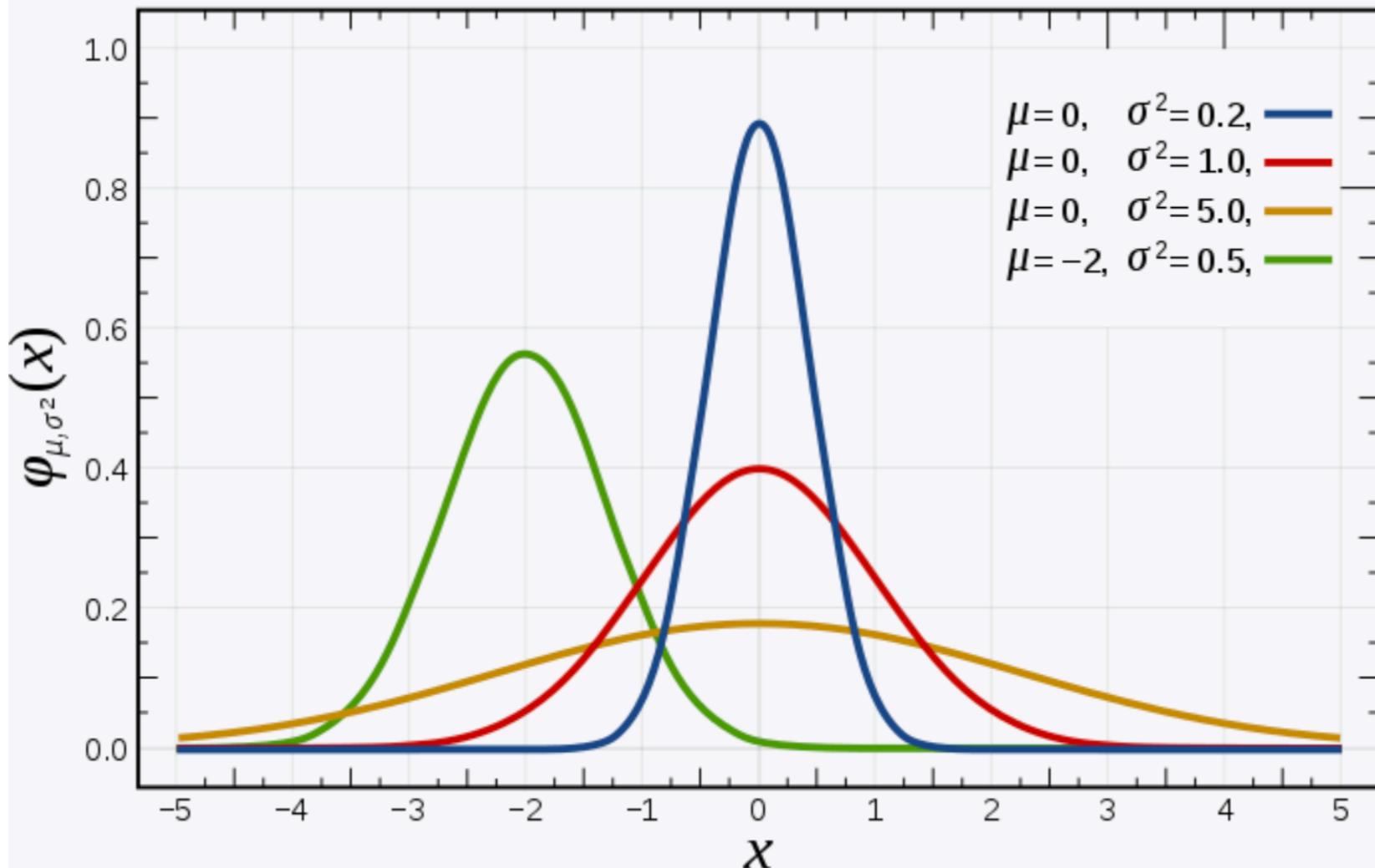
Or directly use the formula by finding the inverse $X^T X$?

3. Probabilistic Approach (Maximal Likelihood Estimation (MLE))

- Work out the details with the students on the board.

Recall Gaussian distribution

Probability density function



The red curve is the *standard normal distribution*

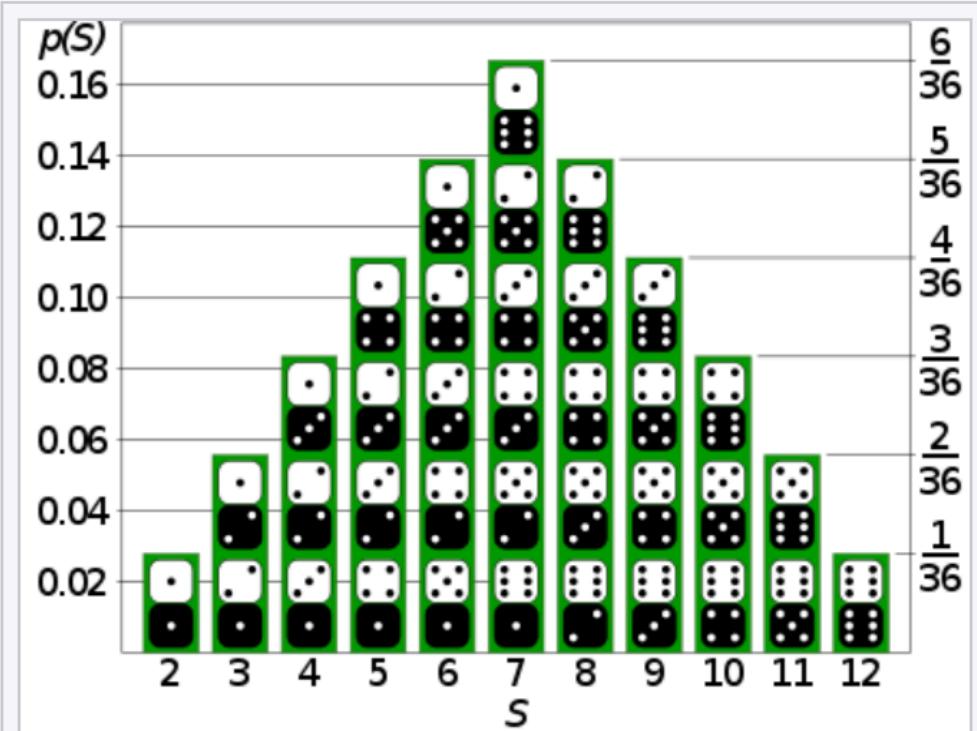
You also can add and multiply two random variables, etc. to define a function of this random variable

Dice roll [[edit](#)]

A random variable can also be used to describe the process of rolling dice and the possible outcomes. The most obvious representation for the two-dice case is to take the set of pairs of numbers n_1 and n_2 from $\{1, 2, 3, 4, 5, 6\}$ (representing the numbers on the two dice) as the sample space. The total number rolled (the sum of the numbers in each pair) is then a random variable X given by the function that maps the pair to the sum:

$$X((n_1, n_2)) = n_1 + n_2$$

Example:



If the sample space is the set of possible numbers rolled on two dice, and the random variable of interest is the sum S of the numbers on the two dice, then S is a discrete random variable whose distribution is described by the **probability mass function** plotted as the height of picture columns here.

Notation	$\mathcal{N}(\mu, \sigma^2)$
Parameters	$\mu \in \mathbb{R}$ = mean (location) $\sigma^2 > 0$ = variance (squared scale)
Support	$x \in \mathbb{R}$
PDF	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Expected Value

Let X be a random variable with a finite number of finite outcomes x_1, x_2, \dots, x_k occurring with probabilities p_1, p_2, \dots, p_k , respectively. The **expectation** of X is defined as

$$E[X] = \sum_{i=1}^k x_i p_i = x_1 p_1 + x_2 p_2 + \cdots + x_k p_k.$$

Since all probabilities p_i add up to 1 ($p_1 + p_2 + \cdots + p_k = 1$), the expected value is the **weighted average**, with p_i 's being the weights.

Special case: Average

If all outcomes x_i are **equiprobable** (that is, $p_1 = p_2 = \dots = p_k$), then the weighted average turns into the simple **average**. This is intuitive: the expected value of a random variable is the average of all values it can take; thus the expected value is what one expects to happen *on average*.

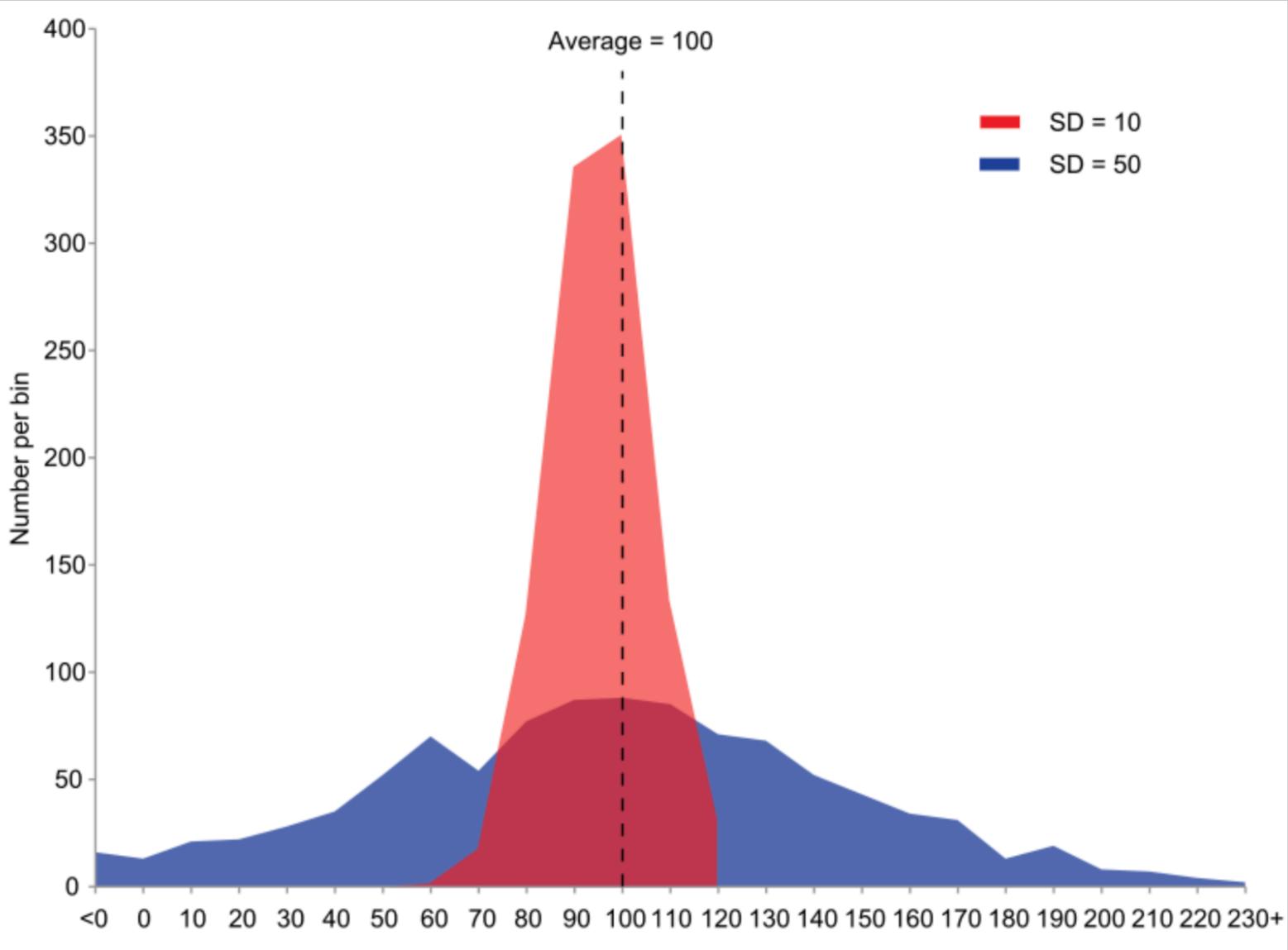
Continuous case

If X is a random variable whose **cumulative distribution function** admits a **density** $f(x)$, then the expected value is defined as the following Lebesgue integral:

$$\mathbb{E}[X] = \int_{\mathbb{R}} x f(x) dx.$$

The variance of a random variable X is the **expected value** of the squared deviation from the **mean** of X , $\mu = \text{E}[X]$:

$$\text{Var}(X) = \text{E}[(X - \mu)^2].$$



Example of samples from two populations with the same mean but different variances. The red population has mean 100 and variance 100 ($SD=10$) while the blue population has mean 100 and variance 2500 ($SD=50$). □

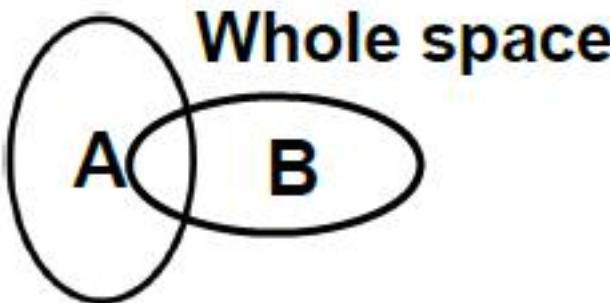
$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

Continuous case

If the random variable X represents samples generated by a continuous distribution with probability density function $f(x)$,

$$\text{Var}(X) = \sigma^2 = \int (x - \mu)^2 f(x) dx$$

Visualize Bayes' Theorem



$$P(A) = \frac{\text{Area of } A}{\text{Area of Whole space}}$$

$$P(B) = \frac{\text{Area of } B}{\text{Area of Whole space}}$$

$$P(A|B) = \frac{\text{Area of } A \cap B}{\text{Area of } B}$$

$$P(B|A) = \frac{\text{Area of } A \cap B}{\text{Area of } A}$$

$$P(A \cap B) = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}}$$

$$P(A) \times P(B|A) = \frac{\text{Area of } A}{\text{Area of Whole space}} \times \frac{\text{Area of } A \cap B}{\text{Area of } B} = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}} = P(A \cap B)$$

$$P(B) \times P(A|B) = \frac{\text{Area of } B}{\text{Area of Whole space}} \times \frac{\text{Area of } A \cap B}{\text{Area of } A} = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}} = P(A \cap B)$$

$$\Rightarrow P(B|A) = P(A|B) \times P(B) / P(A)$$

- Back up slides

Taking Partial Derivatives -for different Types of functions

30

Type 1 : $\text{IR} \rightarrow \text{IR}$ (one-to-one)
 $x \mapsto f(x)$

$$\frac{\partial f}{\partial x} = \frac{df}{dx}$$

* Type 2 : $\text{IR}^n \rightarrow \text{IR}$ (Many-to-one)
 $(x_1, x_2, \dots, x_n) \mapsto f(x_1, \dots, x_n)$

$$\left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) \stackrel{\Delta}{=} \nabla f$$

$$\nabla f(\vec{a}) = \left(\frac{\partial f}{\partial x_1}|_{\vec{a}}, \frac{\partial f}{\partial x_2}|_{\vec{a}}, \dots, \frac{\partial f}{\partial x_n}|_{\vec{a}} \right)$$

is called the gradient of f at \vec{a} .

Type 3 : $\text{IR} \rightarrow \text{IR}^m$ (one-to-many)
 $t \mapsto (f_1(t), \dots, f_m(t)) \stackrel{\Delta}{=} f(t)$

$$\begin{bmatrix} \frac{\partial f_1}{\partial t} \\ \vdots \\ \frac{\partial f_m}{\partial t} \end{bmatrix} = \begin{bmatrix} \frac{df_1}{dt} \\ \vdots \\ \frac{df_m}{dt} \end{bmatrix} \stackrel{\Delta}{=} f'(t)$$

Key Technique:
Treat each component function as many-to-one function!

* Type 4 : $\text{IR}^n \rightarrow \text{IR}^m$ (many-to-many)
 $(x_1, \dots, x_n) \mapsto (f_1(\vec{x}), \dots, f_m(\vec{x}))$

$$Df(x_1, x_2, \dots, x_n) = \underbrace{\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}}_{\text{Derivative matrix}} \leftarrow \nabla f_1(\vec{x})$$

$$\leftarrow \nabla f_2(\vec{x})$$

$$\leftarrow \nabla f_m(\vec{x})$$

You must
keep your
mind
clear
what type
of
function
you are
dealing
with!

Normal Equation for Least Square Approximation

- i.e. Representing the Least Square Solution in Matrix Form
- Work out the details with the students on the board.
- Recall the product rule:
 - $f, g: \mathbb{R} \rightarrow \mathbb{R}$: $(f \cdot g)' = f' \cdot g + f \cdot g'$
 - $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$: $\nabla(f \cdot g) = \nabla f \cdot g + f \cdot \nabla g$
 - $\mathbf{f}, \mathbf{g}: \mathbb{R}^n \rightarrow \mathbb{R}^n$: $(\mathbf{f} \cdot \mathbf{g})' = \mathbf{f}' \cdot \mathbf{g} + \mathbf{f} \cdot \mathbf{g}'$

$$\theta = (X^T X)^{-1} X^T \vec{y}.$$

Mean

- Arithmetic Mean

The *arithmetic mean* (or simply "mean") of a sample x_1, x_2, \dots, x_n , usually denoted by \bar{x} , is the sum of the sampled values divided by the number of items in the example

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right) = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

- Expected Value:

The mean of a **probability distribution** is the long-run arithmetic average value of a **random variable** having that distribution. In this context, it is also known as the **expected value**. For a **discrete probability distribution**, the mean is given by $\sum xP(x)$, where the sum is taken over all possible values of the random variable and $P(x)$ is the **probability mass function**.

Mean of a probability Distribution (Expected Value)

The mean of a **probability distribution** is the long-run arithmetic average value of a **random variable** having that distribution. In this context, it is also known as the **expected value**. For a **discrete probability distribution**, the mean is given by

$\sum xP(x)$, where the sum is taken over all possible values of the random variable and $P(x)$ is the **probability mass function**.

For a **continuous distribution**, the mean is $\int_{-\infty}^{\infty} xf(x) dx$, where $f(x)$ is the **probability density function**.

If the entries in the **column vector**

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

are **random variables**, each with finite **variance**, then the covariance matrix Σ is the matrix whose (i, j) entry is the **covariance**

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = \mathbb{E} [(X_i - \mu_i)(X_j - \mu_j)] = \mathbb{E} [X_i X_j] - \mu_i \mu_j$$

where the operator \mathbb{E} denotes the **expected (mean) value** of its argument, and

$$\mu_i = \mathbb{E}(X_i)$$

is the **expected value** of the i th entry in the vector \mathbf{X} .

Covariance Matrix

$$\Sigma = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}.$$

Note: The covariance matrix is a symmetric matrix.

In fact, a covariant matrix is also positive semi-definite.

The inverse of this matrix, Σ^{-1} , if it exists, is the inverse covariance matrix, also known as the concentration matrix or **precision** matrix.^[1]

