

# **Mathematics of Big Data, I**

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

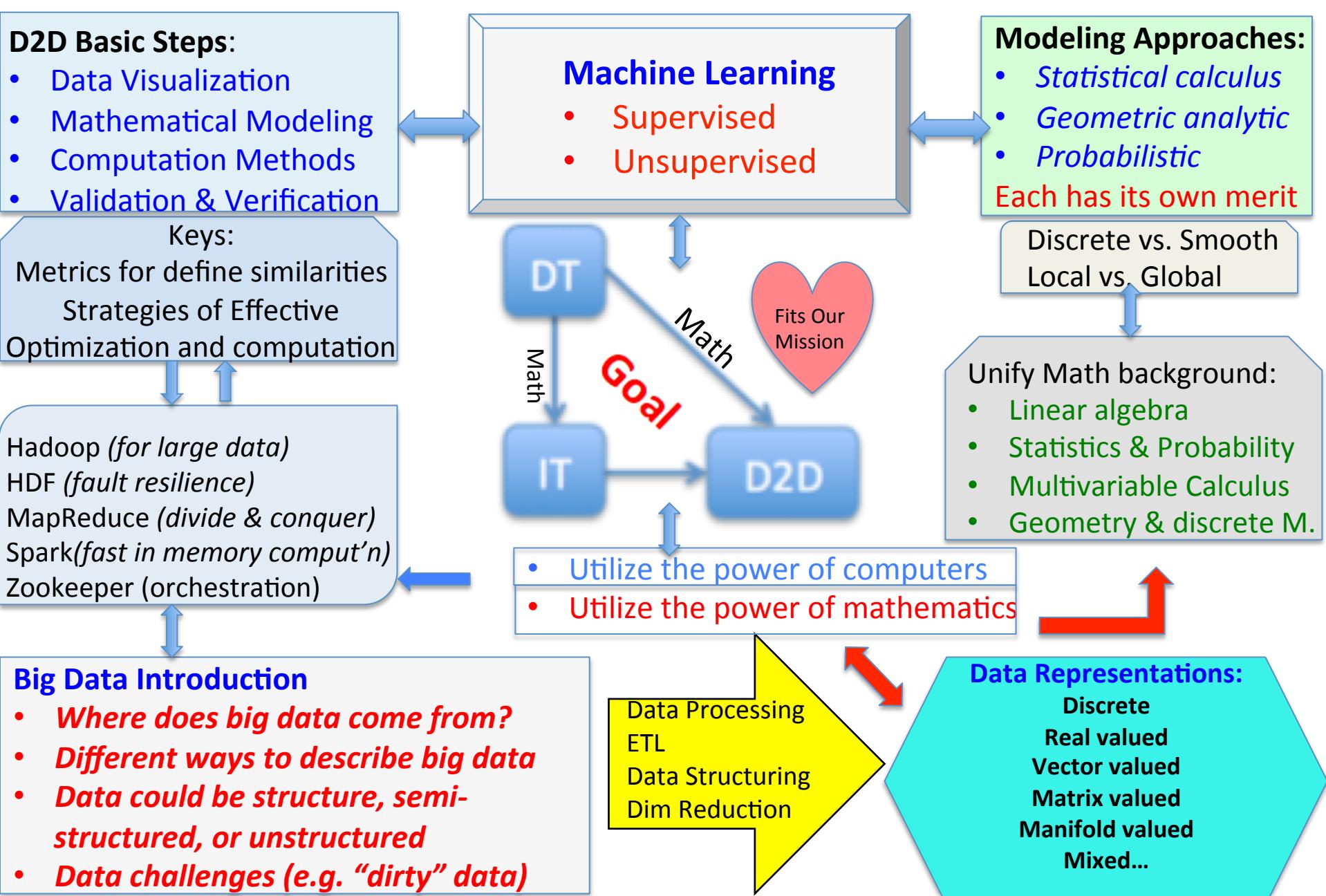
**Weiqing Gu**

Professor of Mathematics  
Director of the Mathematics Clinic

Harvey Mudd College  
Summer 2017

<https://math189su17.github.io/project.html>  
<https://github.com/math189su17/math189su17.github.io>

# 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

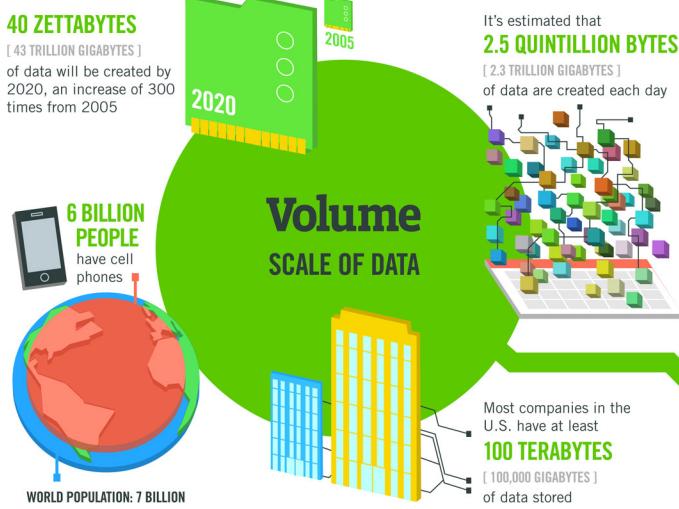
People

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)**



The New York Stock Exchange captures

**1 TB OF TRADE INFORMATION** during each trading session

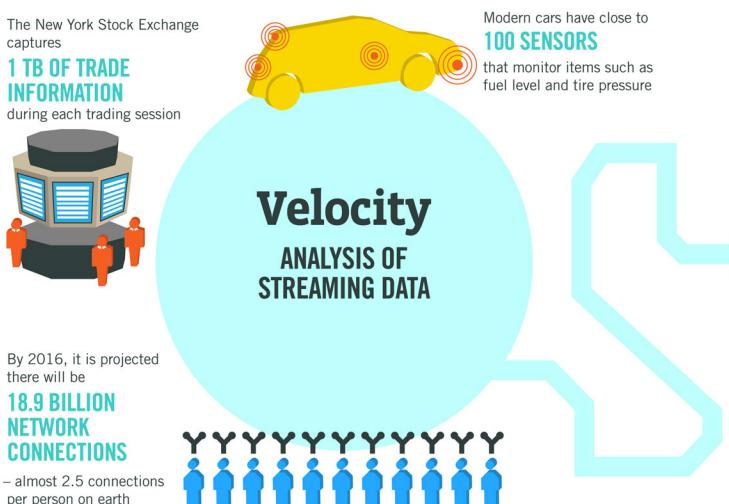


By 2016, it is projected there will be

**18.9 BILLION NETWORK CONNECTIONS**

– almost 2.5 connections per person on earth

Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS



# 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]



**30 BILLION PIECES OF CONTENT**

are shared on Facebook every month



**Variety DIFFERENT FORMS OF DATA**

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)

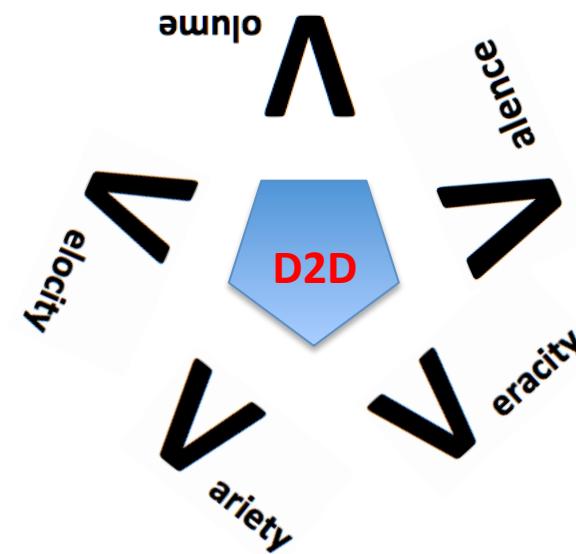
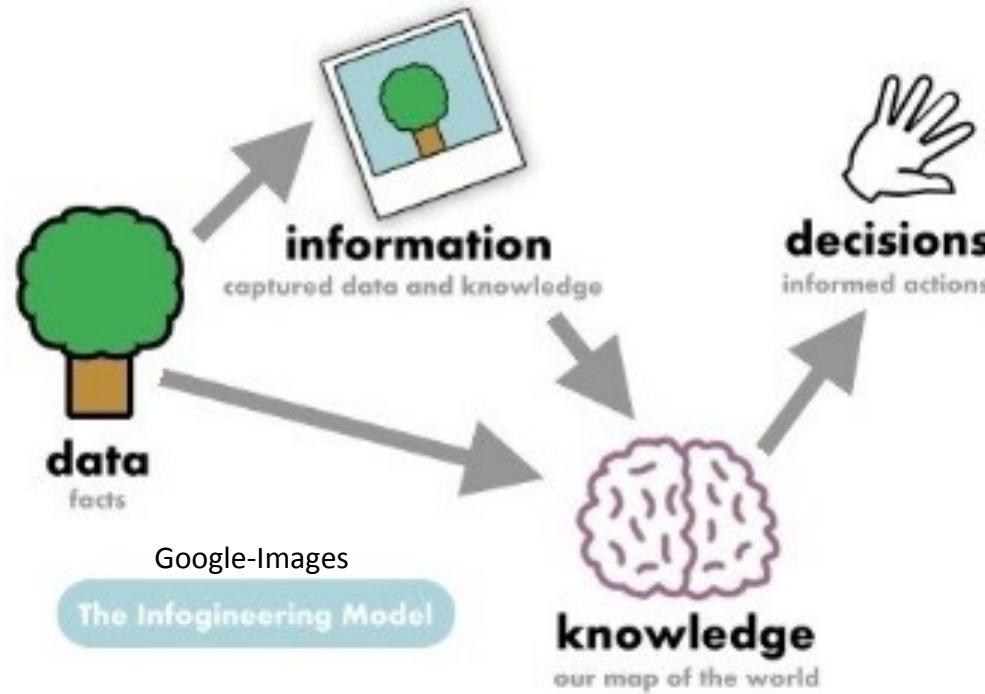
V  
olume

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

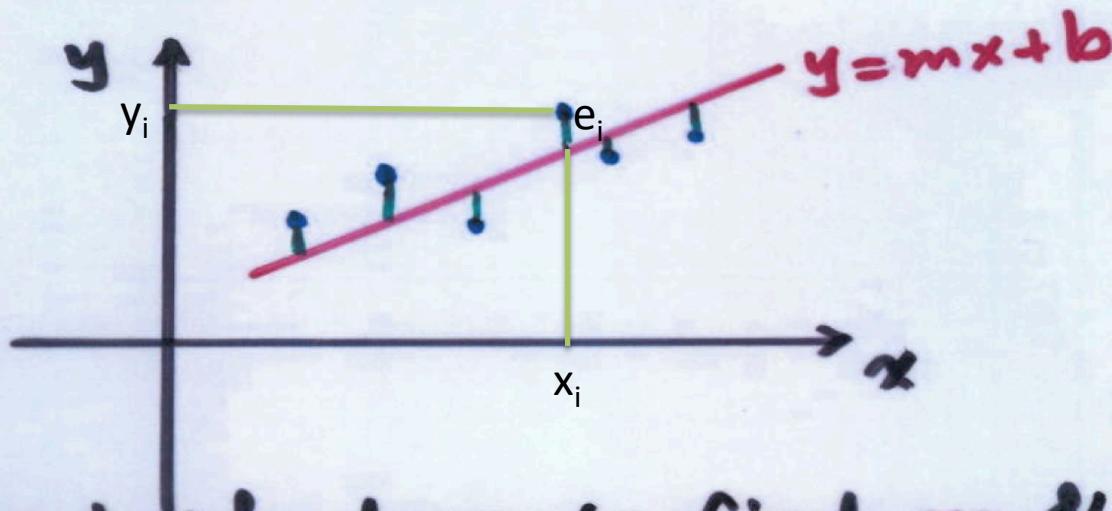
## Modeling Approaches:

- *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?  
 $\sum |e_i|$  is labeled *L<sub>1</sub> norm*

# **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  $m, b$ :

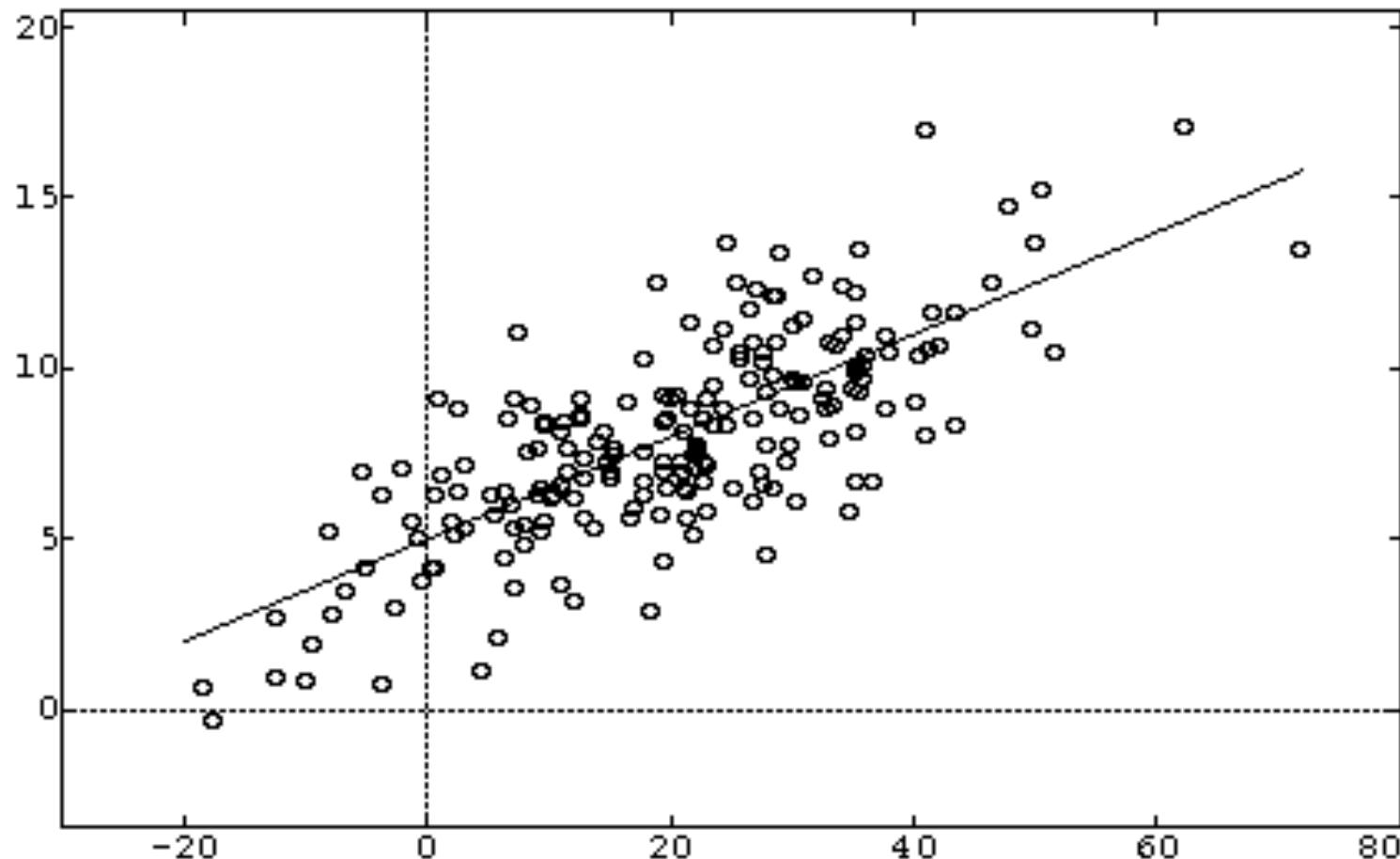
$$m = \frac{n \sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

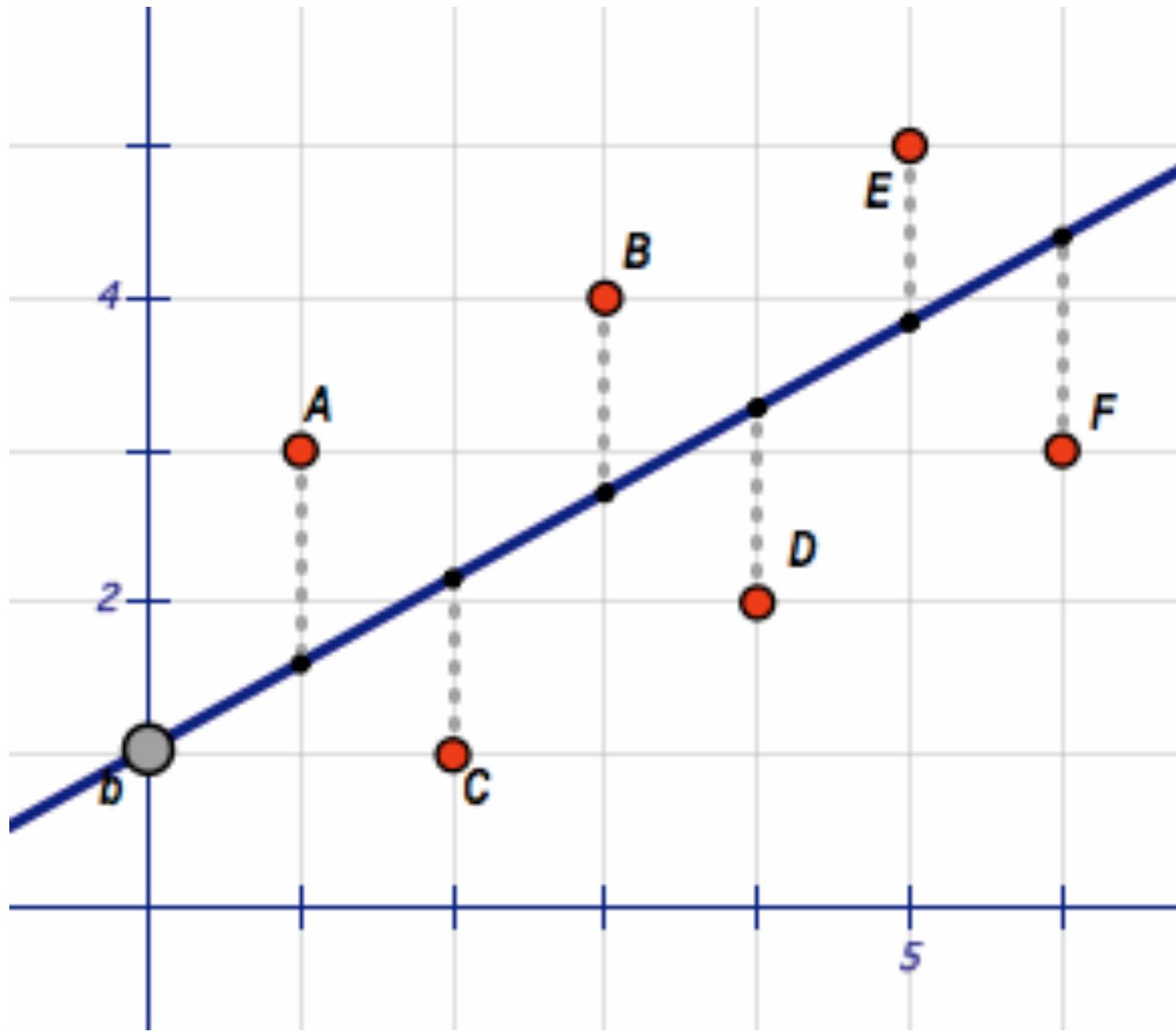
$$b = \frac{(\sum_{i=1}^n x_i^2)(\sum_{i=1}^n y_i) - (\sum_{i=1}^n x_i)(\sum_{i=1}^n x_i y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

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:

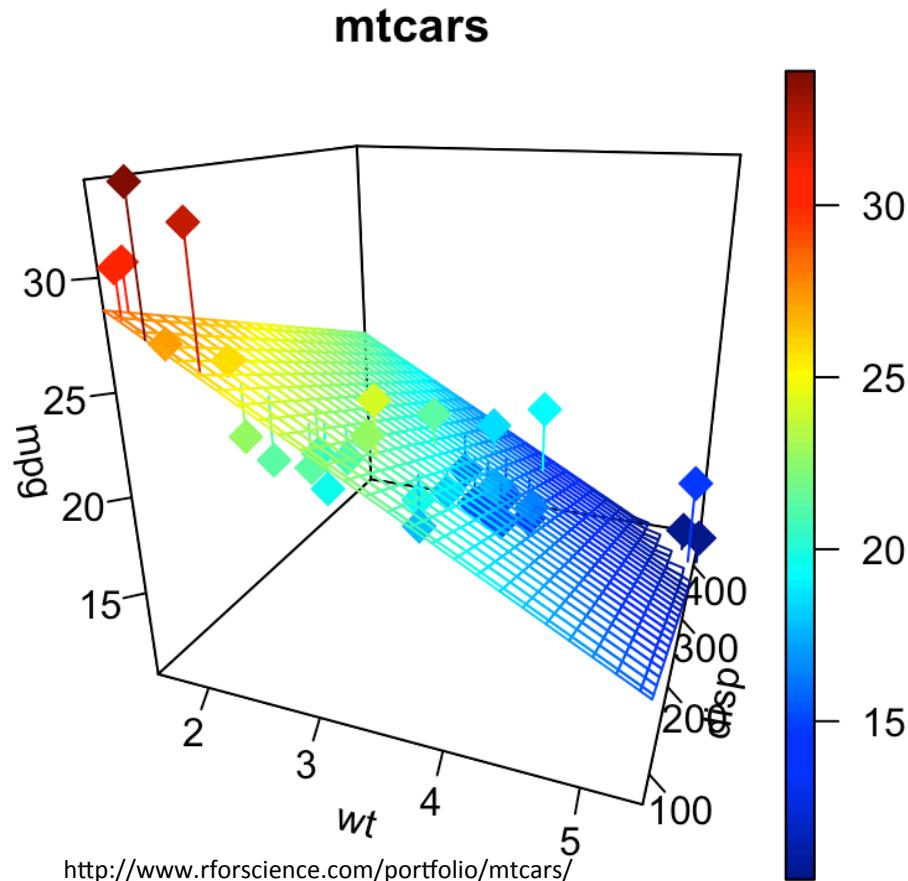
$x_i$	$y_i$	$\bar{x}_i$	$x_i y_i$
0	1	0	0
2	3	4	6
3	6	9	18
4	8	16	32

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

- 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?



# 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

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

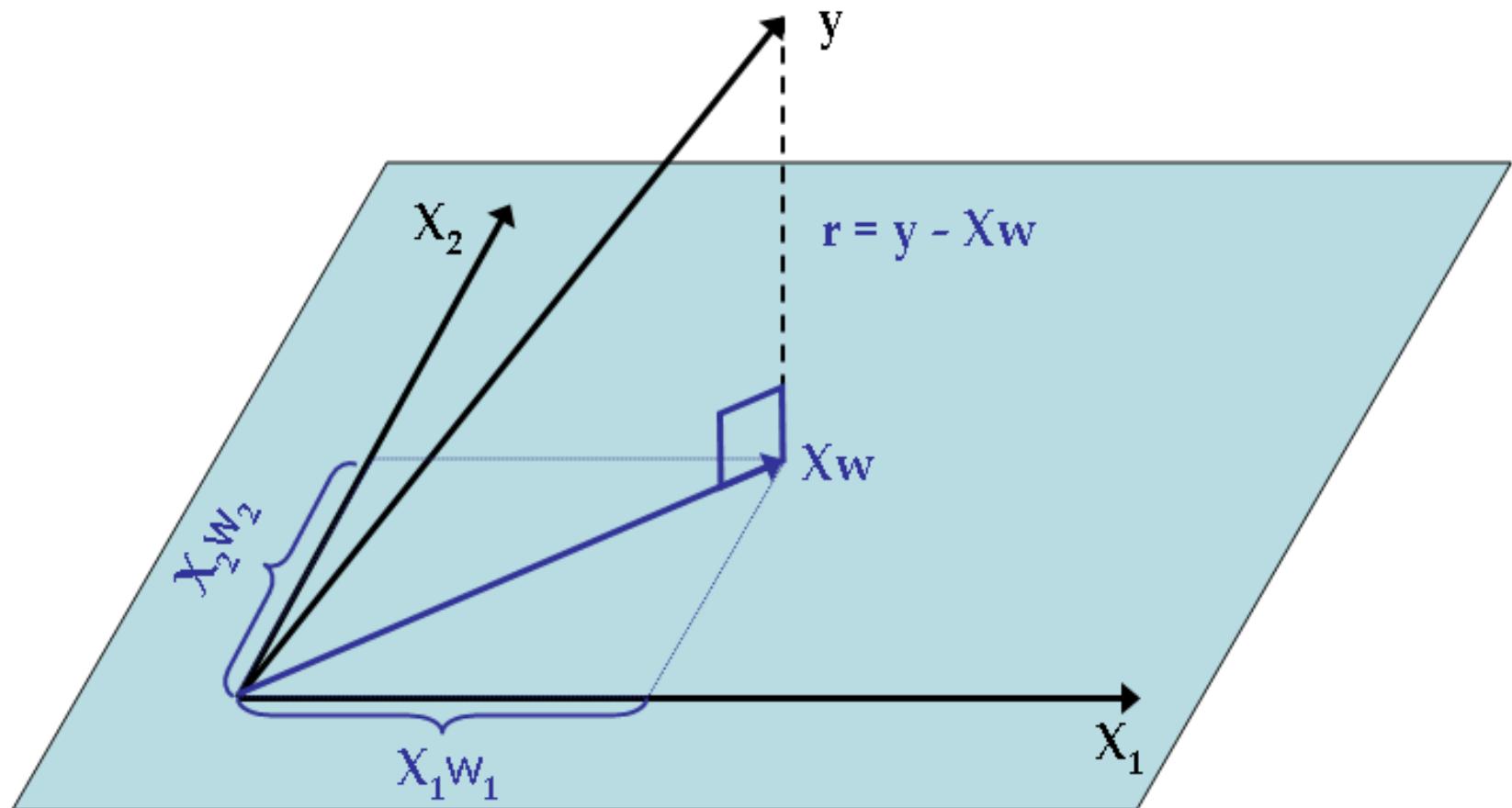
$$\rightarrow \mathbf{r} = (\mathbf{y} - \mathbf{X}\mathbf{w})$$

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!*



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

### **3. Probabilistic Approach (Maximal Likelihood)**

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

Again we get the same solution!

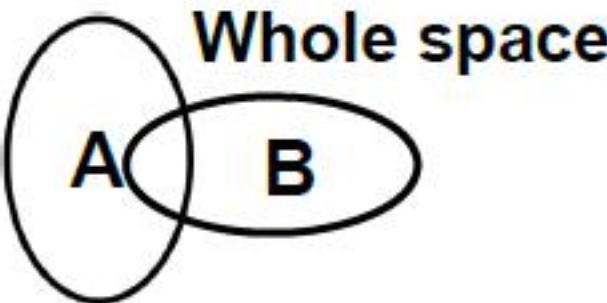
$$\theta = (X^T X)^{-1} X^T \vec{y}.$$

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

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

- Back up slides

# Visualize Bayes' Theorem



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

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

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

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

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

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

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

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

## 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!

# A Big Picture of Derivatives (By Prof. Gu)

Type of functions	Type of Derivatives	Notations	Pictures	Meanings	Remarks	Formulas
$f: \mathbb{R} \rightarrow \mathbb{R}$	Derivative of $f(x)$ at $x_0$ .	$\frac{df}{dx} _{x_0}$ or $f'(x_0)$		Slope at $(x_0, f(x_0))$ of the curve		The tangent line at $(x_0, f(x_0))$ , $y = f(x_0) + f'(x_0)(x - x_0)$
$f: \mathbb{R}^2 \rightarrow \mathbb{R}$ or $\mathbb{R}^n \rightarrow \mathbb{R}$	Partial derivatives w.r.t $x$ or $y$	$\frac{\partial f}{\partial x} _{(a,b)} = f_x$ $\frac{\partial f}{\partial y} _{(a,b)} = f_y$		$f_x =$ slope of graph in $x$ -direction $f_y =$ slope of graph in $y$ -direction	Similarly, $f_y =$ slope of graph in $y$ -direction	(If exists) The tangent plane at $(a, b, f(a, b))$ : $z = f(a, b) + f_x(a, b)(x-a) + f_y(a, b)(y-b)$
$f: \mathbb{R}^3 \rightarrow \mathbb{R}$ or $\mathbb{R}^n \rightarrow \mathbb{R}$	Higher order partials: Here: 2 <sup>nd</sup> partials	$\frac{\partial^2 f}{\partial x^2} _{(a,b)} = f_{xx}$ $\frac{\partial^2 f}{\partial y^2} _{(a,b)} = f_{yy}$		$f_{xx} =$ concavity in $x$ direction $(f_y)_x =$ rate of change of $f_y$ as $x$ increases	$f_{xx} > 0 :$ $f_{yy} < 0 :$	Laplace's eqn: $\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} + \frac{\partial^2 f}{\partial z^2} = 0$
$f: \mathbb{R}^2 \rightarrow \mathbb{R}$ or $\mathbb{R}^n \rightarrow \mathbb{R}$	Mixed partials	$\frac{\partial^2 f}{\partial x \partial y} _{(a,b)} = f_{xy}$		$f_{xy} =$ rate of change of $f_y$ as $x$ increases	In the first picture, as $x$ increases, $f_y$ increases from neg to positive $\Rightarrow f_{xy} = 0$	$f_{xy} = f_{yx}$
$f: \mathbb{R}^n \rightarrow \mathbb{R}$	Directional derivative of $\mathbb{R}$ -valued function of $n$ -variable	$D_u f(\vec{a})$		Rate change of $f$ in the direction of $u$ .	Any directional derivative is completely determined by the direction and gradient.	$D_u f(\vec{a}) = \nabla f(\vec{a}) \cdot \vec{u}$
$f: \mathbb{R}^n \rightarrow \mathbb{R}$	Gradient of real valued function of $2$ or $3$ variables	$\nabla f(\vec{a})$		$\frac{\partial f(a, b)}{\  \nabla f(a, b) \ } \hat{u}, a$ direction of steepest ascent. i.e. where $D_u f(a)$ is maximized.	If $S$ is a surface given by: $f(x, y, z) = C$ then an equation for tangent plane to $S$ at $x_0$ : $\nabla f(x_0) \cdot (\vec{x} - \vec{x}_0) = 0$	$\nabla = i \frac{\partial}{\partial x} + j \frac{\partial}{\partial y} + k \frac{\partial}{\partial z}$ $\nabla f = i \frac{\partial f}{\partial x} + j \frac{\partial f}{\partial y} + k \frac{\partial f}{\partial z}$
$\vec{F}: \mathbb{R}^n \rightarrow \mathbb{R}^n$	Divergence of a vector field	$\text{div}(\vec{F}) = \nabla \cdot \vec{F}$		Measurement of the "net mass flow" of $\vec{F}$ in or out at a pt.	or $\text{div } \vec{F} =$ rate of expansion per unit volume.	$\text{div } \vec{F} = \nabla \cdot \vec{F} = \frac{\partial F_1}{\partial x_1} + \frac{\partial F_2}{\partial x_2} + \dots + \frac{\partial F_n}{\partial x_n}$
$\vec{F}: \mathbb{R}^3 \rightarrow \mathbb{R}^3$	Curl of a vector field	$\text{curl}(\vec{F}) = \nabla \times \vec{F}$		A tiny twig or paddle at $x \in \mathbb{R}^3$ will spin around the axis in dir <sup>n</sup> of vector $\text{curl } \vec{F}$ obeying R-H rule w/ angular velocity $\pm \text{curl } \vec{F}(x)$ and angular speed $\pm  \text{curl } \vec{F}(x) $ radius/sec.	$\text{curl } (\vec{F}) = \nabla \times \vec{F} = \begin{vmatrix} i & j & k \\ \frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} & \frac{\partial}{\partial x_3} \\ F_1 & F_2 & F_3 \end{vmatrix}$	

# 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:
- $f, g: \mathbb{R} \rightarrow \mathbb{R}$ :  $(f \cdot g)' = f' \cdot g + f \cdot g'$
- $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$ :  $\nabla(f \cdot g) = \nabla f \cdot g + f \cdot \nabla g$
- $\mathbf{f}, \mathbf{g}: \mathbb{R}^n \rightarrow \mathbb{R}^n$ :  $(\mathbf{f} \cdot \mathbf{g})' = \mathbf{f}' \cdot \mathbf{g} + \mathbf{f} \cdot \mathbf{g}'$

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

