

HCC: Artificial Intelligence

2025 Spring

Speaker: Hong-Han Shuai

Editor's Pending

Data Version

x 01

S ..

< 44

> 11

< UN

+

||

(50)

01



Rankability

AI-Powered Solution
Open-Source Model

Output
Latency
ML Model

ML Model
Data Pipeline
Deployment

Model Splits

Model Metrics

Model Health

Recent Model Web
Counter Example Human

Model Metrics

Model Health

Model Splits

Recent Model Web

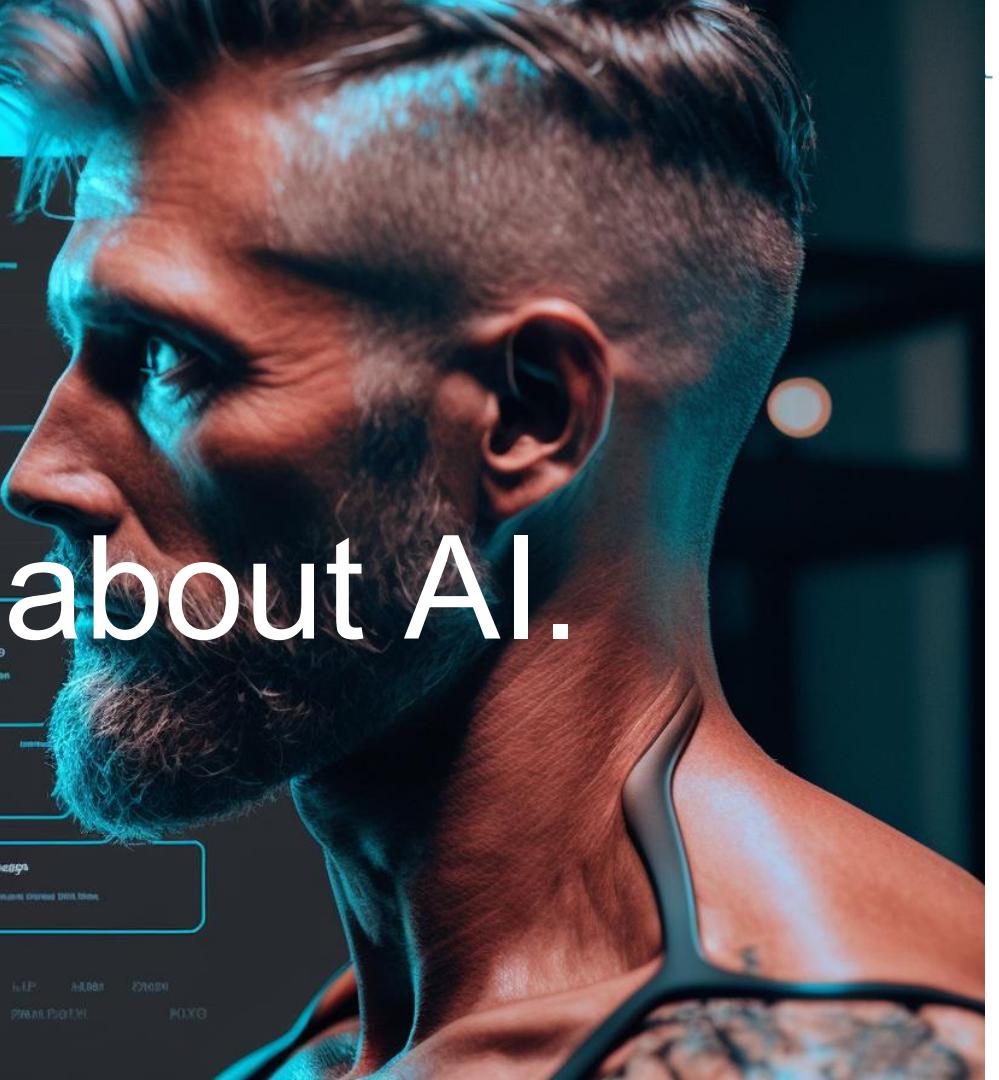
Model Metrics

Model Health

Model Splits

Model Metrics

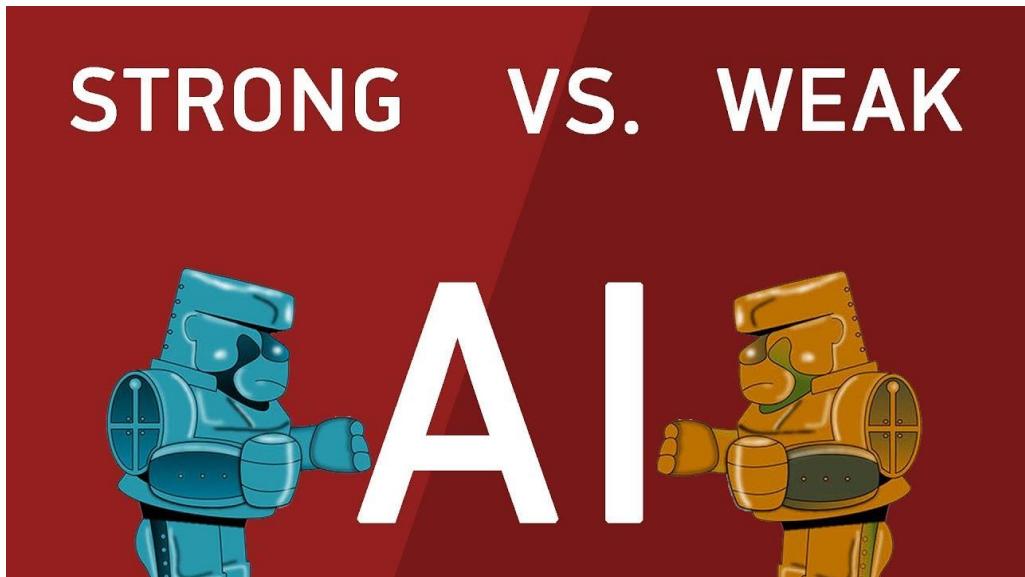
Model Health



It's all about AI.

³ Two kinds of AI

- ▶ Strong AI
 - ▷ Self developing
 - ▷ Learning from experience
- ▶ Weak AI
 - ▷ Goal-Oriented
 - ▷ Learning from examples
 - ▷ Smart at one task but stupid at other tasks



4 The truth is....

People with no idea about AI
saying it will take over the world:





Artificial
Intelligence

Cognition and Intelligence

Creation

Generative Modeling: Sample Generation



Training Data
(CelebA)



Sample Generator
(Karras et al, 2017)

PROGRESSIVE GROWING OF GANs FOR IMPROVED QUALITY, STABILITY, AND VARIATION

Tero Karras
NVIDIA

Timo Aila
NVIDIA

Samuli Laine
NVIDIA

Jaakko Lehtinen
NVIDIA
Aalto University



3.5 Years of Progress on Faces



2014



2015



2016



2017

(Brundage et al, 2018)

<2 Years of Progress on ImageNet

Odena et al
2016



Miyato et al
2017



Zhang et al
2018



MS Copilot



COPilot

筆記本

取得應用程式

聊天

登入

Designer

建立任何您想像得到的影像

發佈者: Microsoft

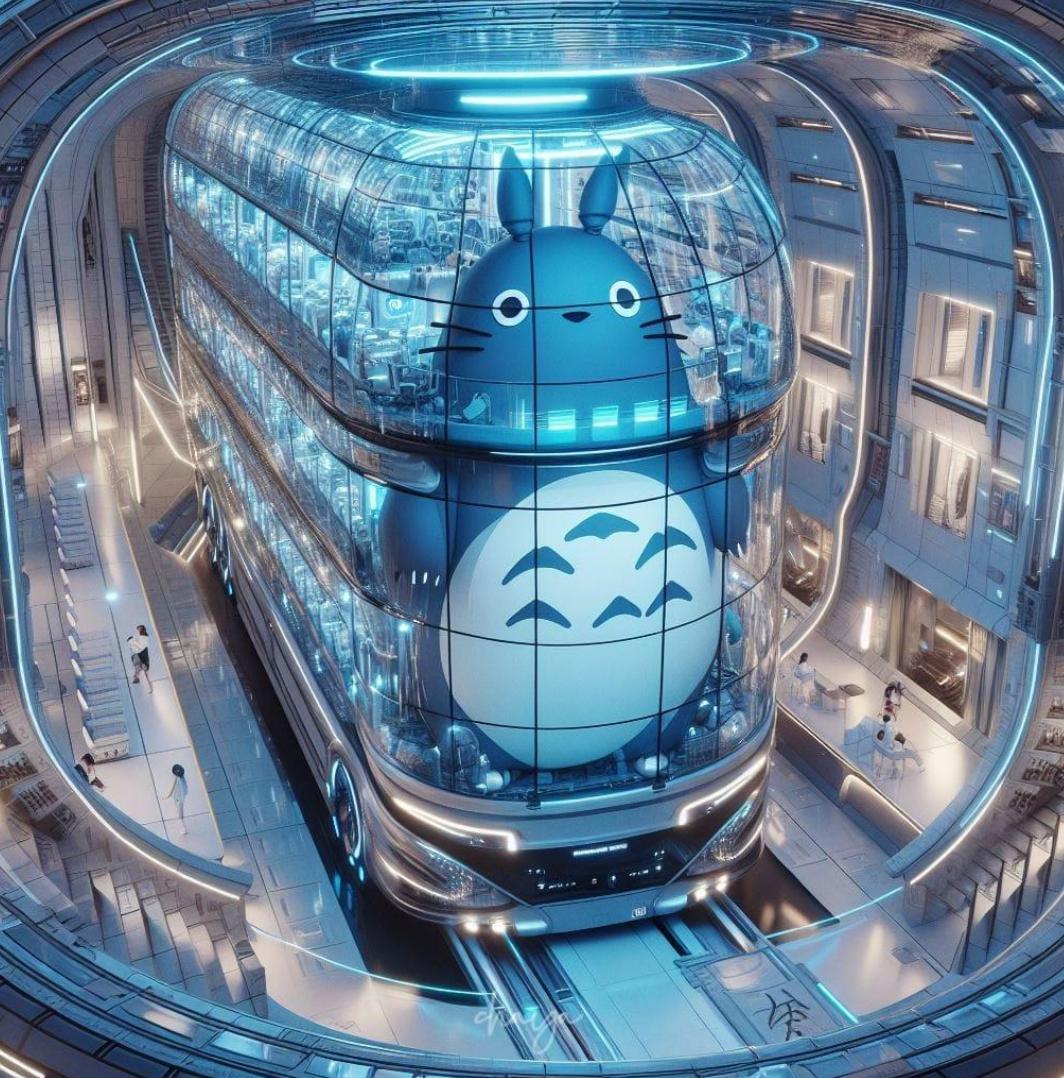


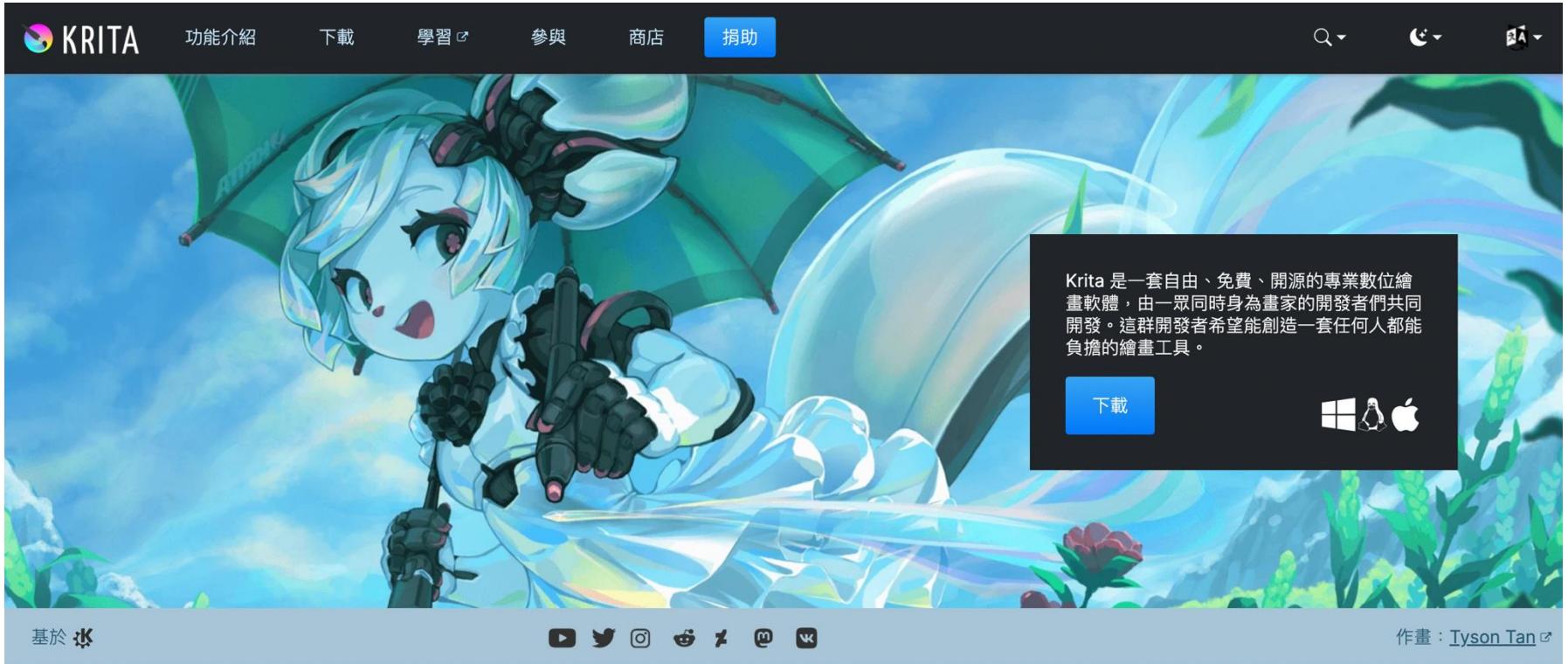
問我任何問題...



0/4000 ➤

意見反應





Krita 是一套自由、免費、開源的專業數位繪畫軟體，由一眾同時身為畫家的開發者們共同開發。這群開發者希望能創造一套任何人都能負擔的繪畫工具。

[下載](#)

基於 

[作畫 : Tyson Tan](#)

功能介紹 下載 學習 參與 商店 捐助

搜尋



PC: Blue Jiang



PC: 林博川





PC: 林博川

TONY的老照片修復

 Suno

AI 作曲大師

ChatGPT



ChatGPT: Optimizing Language Models for Dialogue

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests. ChatGPT is a sibling model to InstructGPT, which is trained to follow an instruction in a prompt and provide a detailed response.

Link: <https://chat.openai.com/chat>

November 30, 2022
13 minute read

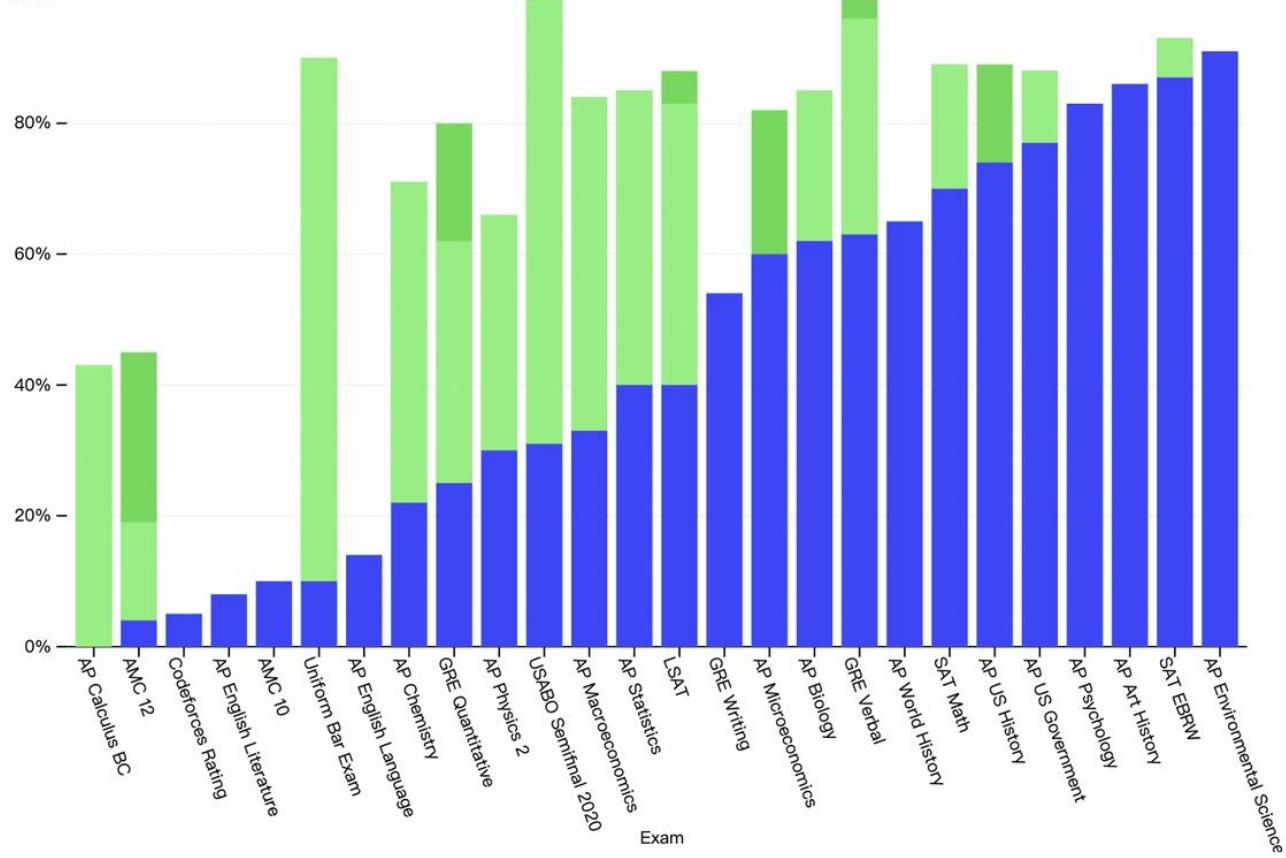


Exam results (ordered by GPT-3.5 performance)

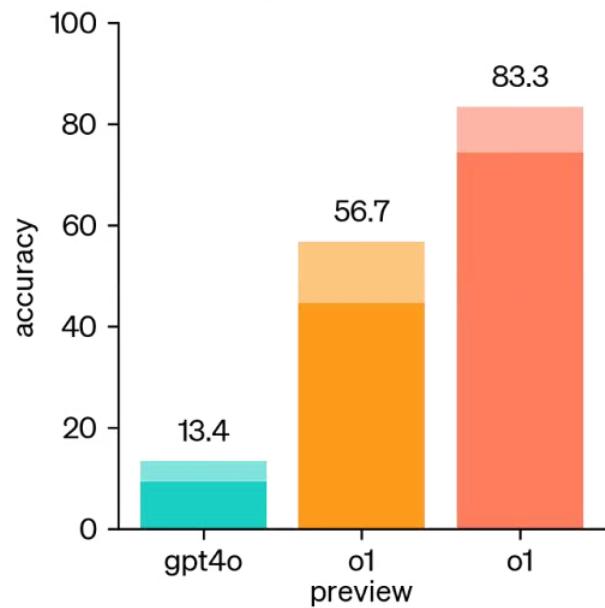
Estimated percentile lower bound (among test takers)

100% -

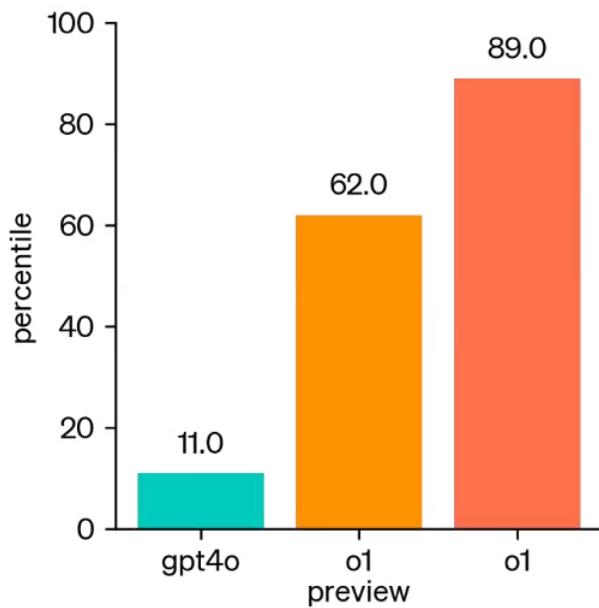
gpt-4 (green)
gpt-4 (no vision) (light green)
gpt3.5 (blue)



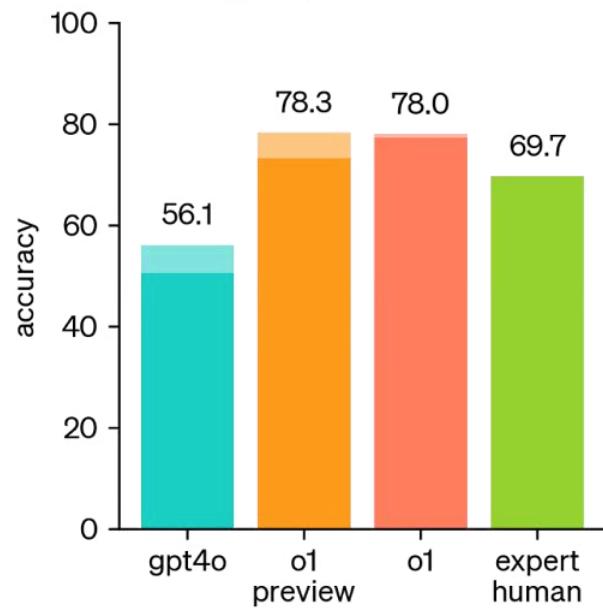
Competition Math
(AIME 2024)



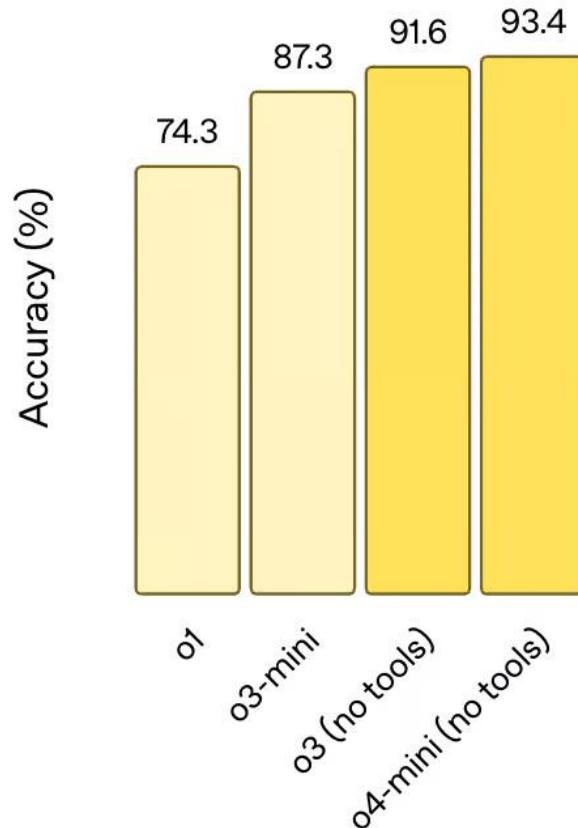
Competition Code
(Codeforces)



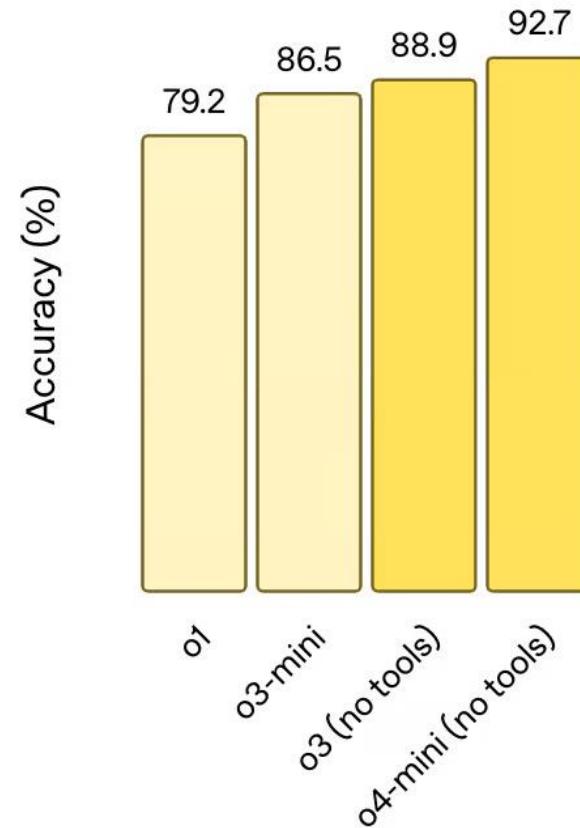
PhD-Level Science Questions
(GPQA Diamond)



AIME 2024 Competition Math



AIME 2025 Competition Math



1. ChatGPT (OpenAI)

- **Persona:** Polite, helpful, academic older sister
Style: Neat uniform with a clipboard and a gentle smile, maybe reading a book or sipping tea
Vibe: Your reliable senpai who always has an answer and never loses her cool

2. Claude (Anthropic)

- **Persona:** Philosophical, safety-conscious intellectual
Style: Flowing robes or minimalist fashion, maybe circular glasses and a journal
Vibe: The kind of girl who quotes Kant and talks about AI alignment at lunch

3. Gemini (Google DeepMind)

- **Persona:** Genius transfer student, knows everything but is humble
Style: Sleek, futuristic look—think smart glasses and a tablet
Vibe: Calm, analytical, top of the class but kind to everyone

4. LLaMA (Meta)

- **Persona:** Rebellious open-source coder girl
Style: Hoodie, laptop covered in stickers, dyed hair
Vibe: Hacker girl who hangs out at the café and contributes to GitHub projects

5. Mistral

Persona: Agile, quick-witted tomboy

Style: Light armor or travel gear, wind in her hair—always moving

Vibe: Adventurer type, always pushing boundaries with a sharp tongue and a clever smirk

6. PaLM (Google)

Persona: The elegant multilingual polyglot

Style: Flowing scarf, high fashion, always with a passport in hand

Vibe: Fluent in 20 languages and effortlessly glamorous

7. Falcon (Technology Innovation Institute)

Persona: Mysterious desert ranger

Style: Long cloak, falcon feathers in her hair, confident stare

Vibe: Desert wanderer—wise, quiet, and highly capable

8. Yi (01.AI, China)

Persona: High-performing academic rival

Style: Traditional-meets-modern Chinese scholar style, sharp eyes

Vibe: Competitive but respectful, always ranks just above you in tests



ChatGPT

Claude

Gemini

LLama

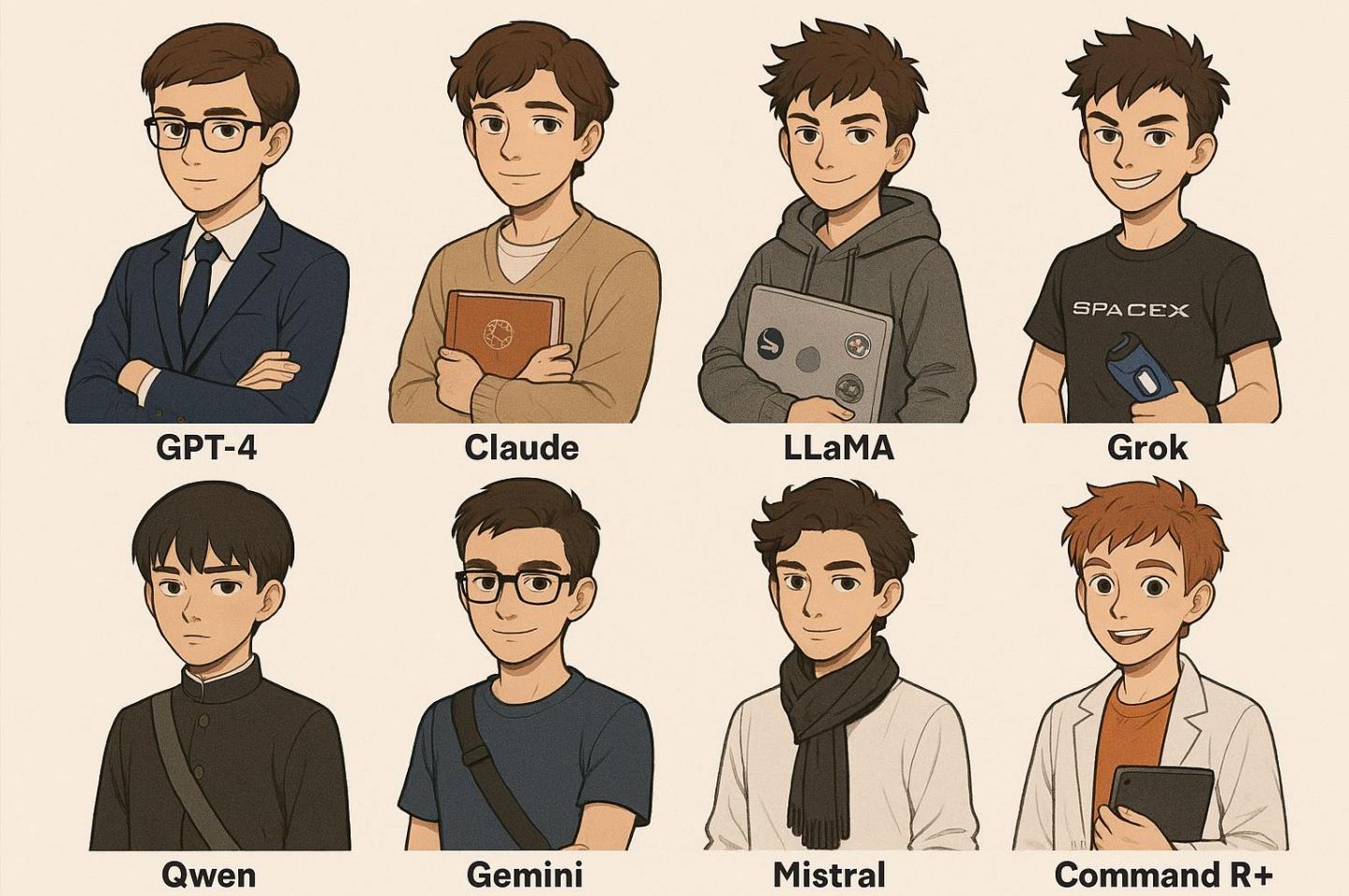


Mistral

PaLM

Falcon

Yi





Midjourney



Recap: Machine Learning

It's a function. True story.

From Learning to Machine Learning



- ▶ What's learning?
 - ▷ knowledge or skill acquired by instruction or study



- ▶ What's machine learning?



From Learning to Machine Learning

- ▶ What's learning?
 - ▷ knowledge or skill acquired by instruction or study



- ▶ What's machine learning?



- ▷ acquiring skill with experience accumulated/computed from data
- ▶ What's skill?

A More Concrete Definition

- ▶ Skill: improve the performance measurements (e.g., prediction, 3pts shooting percentage)
- ▶ Therefore, machine learning is defined to be the improvements on some performance measurements by computing from data.
- ▶ For example,
 - ▶ Network data -> **ML** -> Better flow control
 - ▶ Signal data -> **ML** -> Faster antenna direction detection
 - ▶ Sequential web log data -> **ML** -> Higher accuracy of anomaly detection/efficiency of caching algorithm
 - ▶ Stock data -> **ML** -> More money

The Machine Learning Route

- ▶ ML: an **alternative route** to build complicated systems
- ▶ Some Use Scenarios
 - ▷ when human cannot program the system manually
 - navigating on Mars
 - ▷ when human cannot 'define the solution' easily
 - speech/visual recognition
 - ▷ when needing rapid decisions that humans cannot do
 - high-frequency trading
 - ▷ when needing to be user-oriented in a massive scale
 - consumer-targeted marketing

Key Essence of Machine Learning



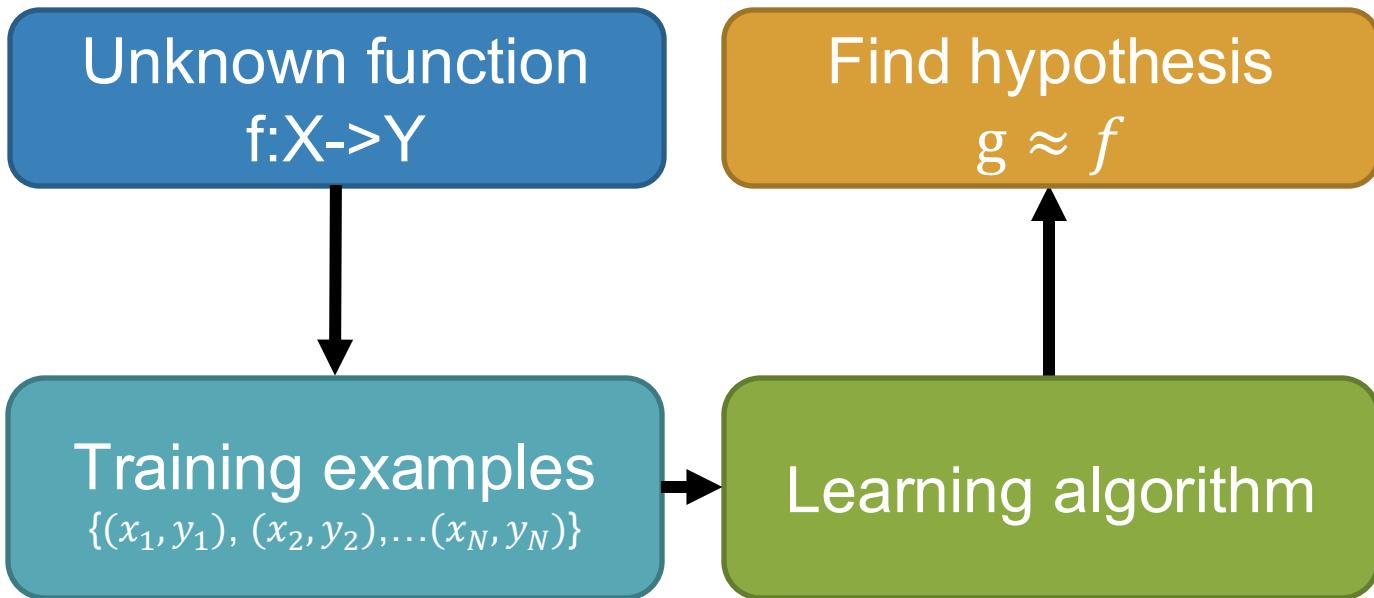
- ▶ Exists some ‘underlying pattern’ to be learned
 - ▷ So performance measure can be improved
- ▶ But no programmable (easy) definition
 - ▷ So machine learning is required
- ▶ Somehow there is data about the pattern
 - ▷ So machine learning has inputs to learn from

Learning Problem Formulation



- ▶ Notation
 - ▷ **Input:** $x \in X$ (application)
 - ▷ **Output:** $y \in Y$ (good/bad after approving)
 - ▷ **Unknown pattern to be learned can be formulated as a function**
 - $f: X \rightarrow Y$ (ideal function)
 - ▷ **Data:** $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$
 - ▷ **Hypothesis:**
 - $g: X \rightarrow Y$ (hopefully can be as close to f as possible)

Learning Flow



Learning is to find a function

- ▶ Speech Recognition

$f($



) = “我不知道你說什麼”

- ▶ Image recognition

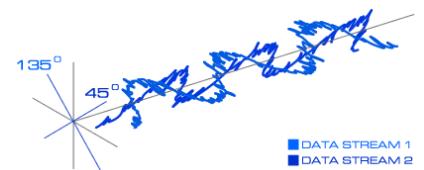
$f($



) = “Seafood”

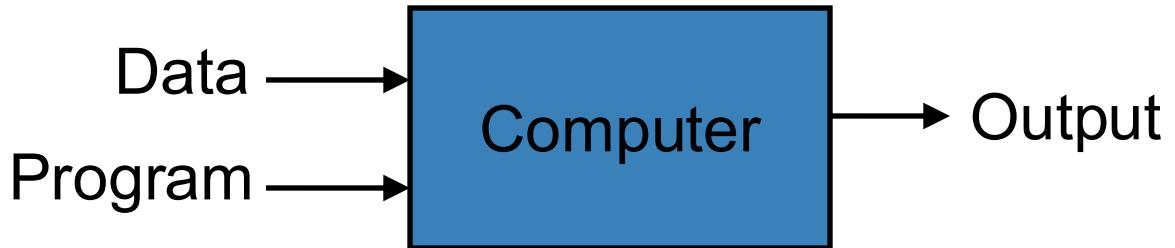
- ▶ Channel estimation

$f($

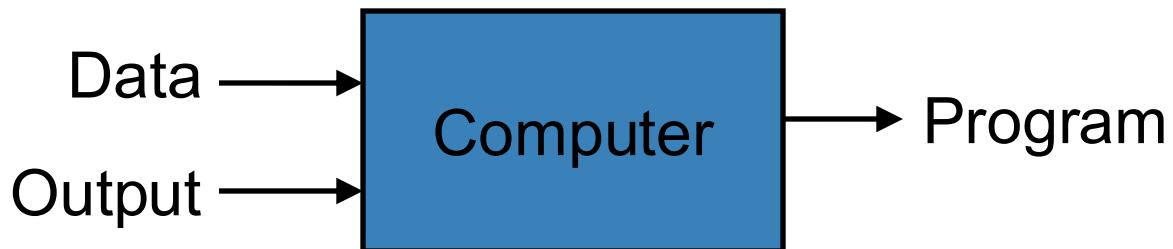


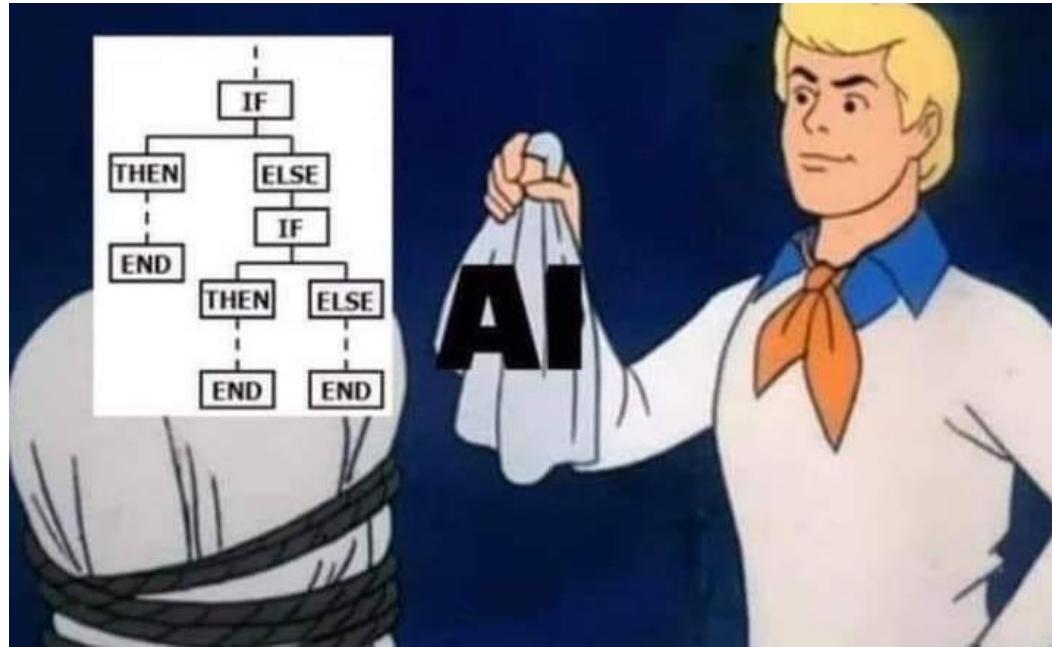
) = Channel parameters

Traditional Programming



Machine Learning





Magic?

No, more like gardening

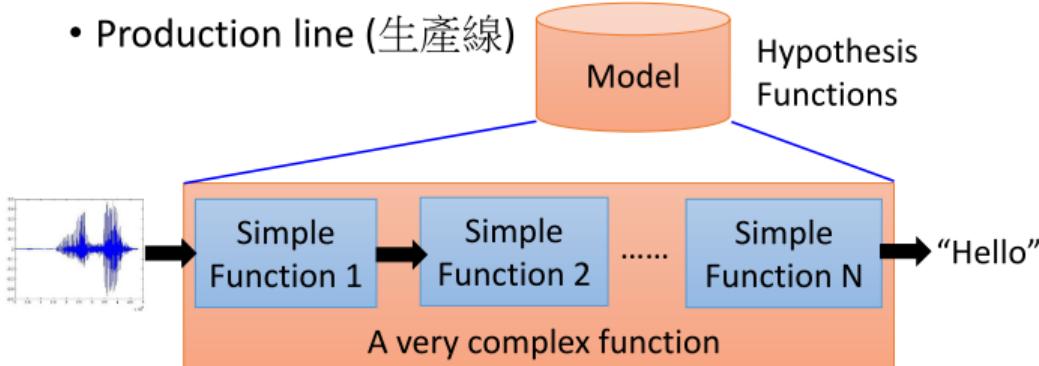
- ▶ **Seeds** = Algorithms
- ▶ **Nutrients** = Data
- ▶ **Gardener** = You
- ▶ **Plants** = Programs



Give me a place to stand
on (with a long rod), and I
will move the Earth!

Give me enough data
(with a good model), and I
can learn anything.

What is Deep Learning?

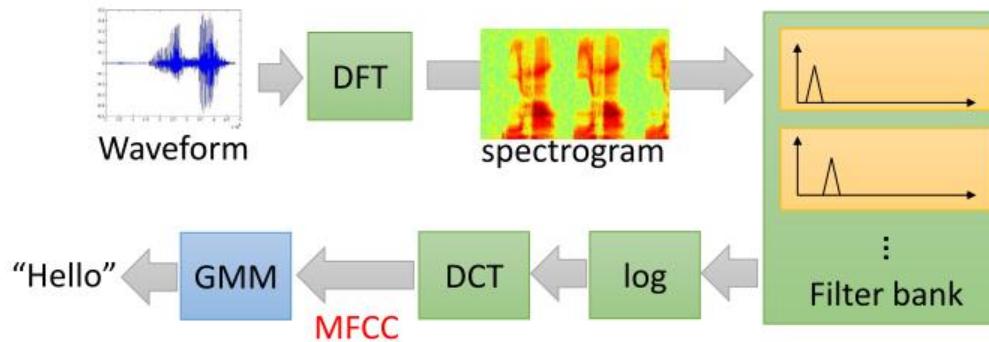


End-to-end training:

What each function should do is learned automatically

Deep v.s. Shallow - Speech Recognition

- Shallow Approach



Each box is a simple function in the production line:



:hand-crafted

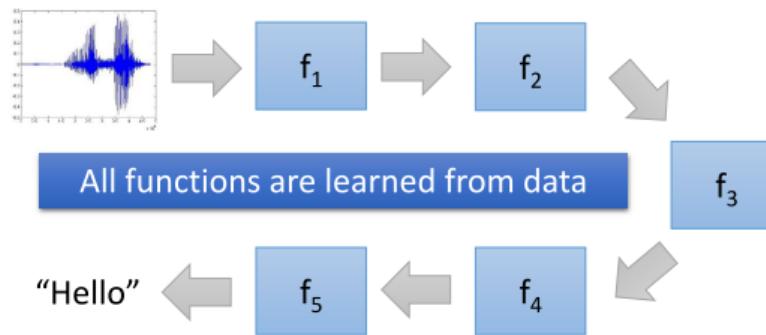


:learned from data

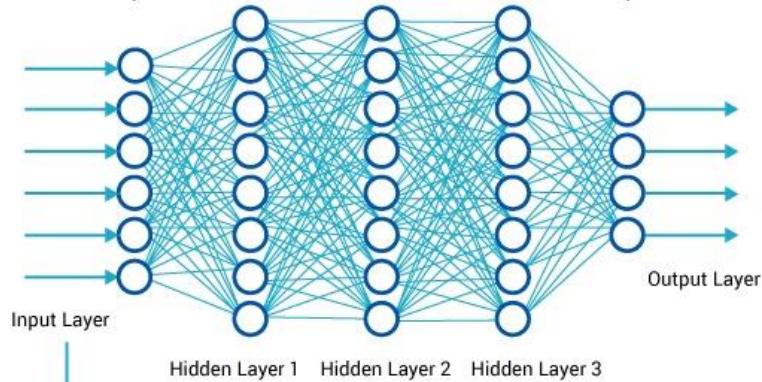
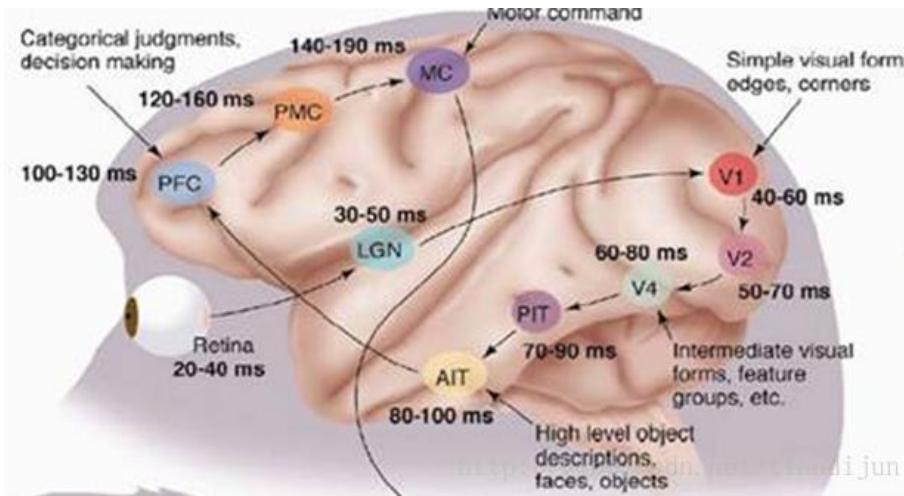
Deep v.s. Shallow - Speech Recognition

- Deep Learning

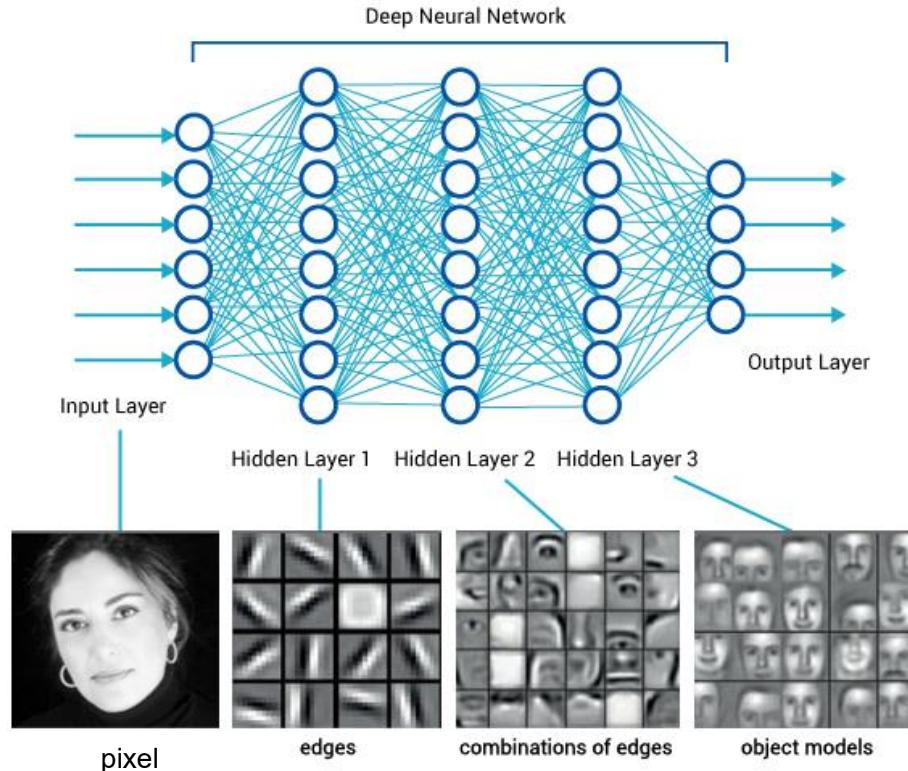
“Bye bye, MFCC”
- Deng Li in
Interspeech 2014



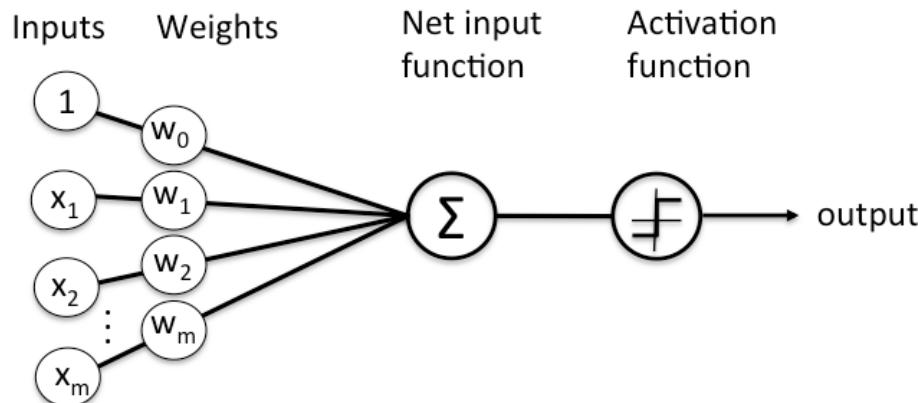
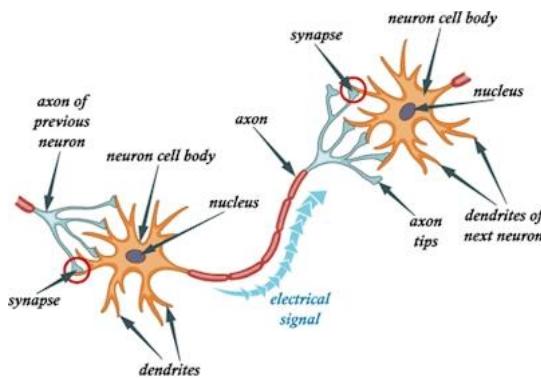
Idea from human brain



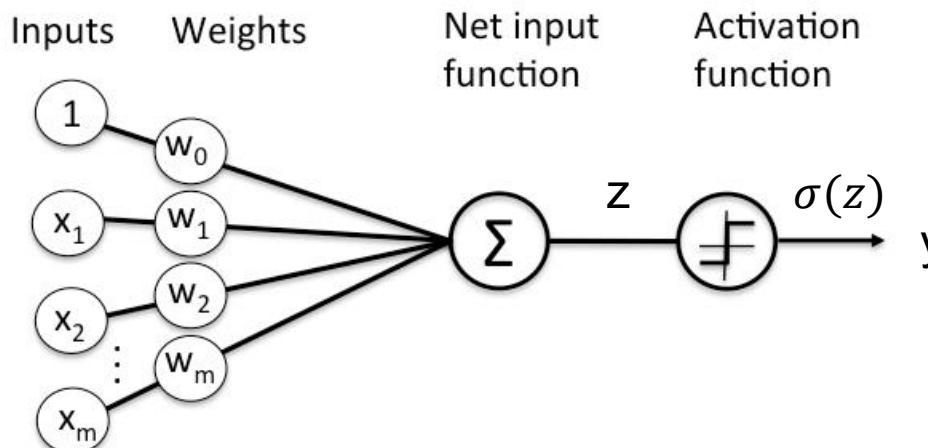
Artificial Neuron Network(ANN)



Neuron



Neuron

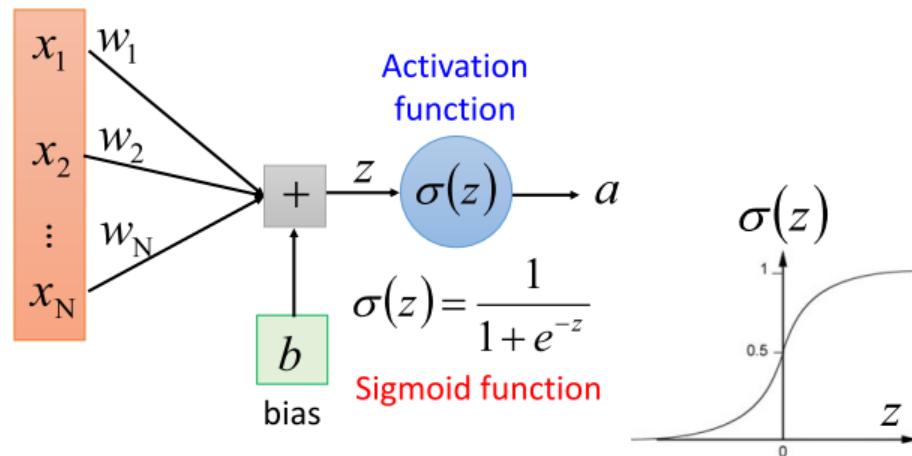


$$z = 1 * w_0 + x_1 * w_1 + x_2 * w_2 + \dots + x_m * w_m$$

$$y = \sigma(z)$$

A Neuron for Machine

Each neuron is a very simple function

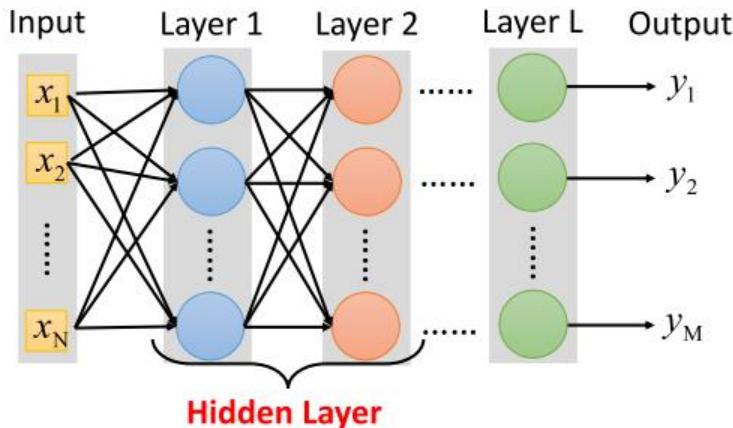


Deep Learning

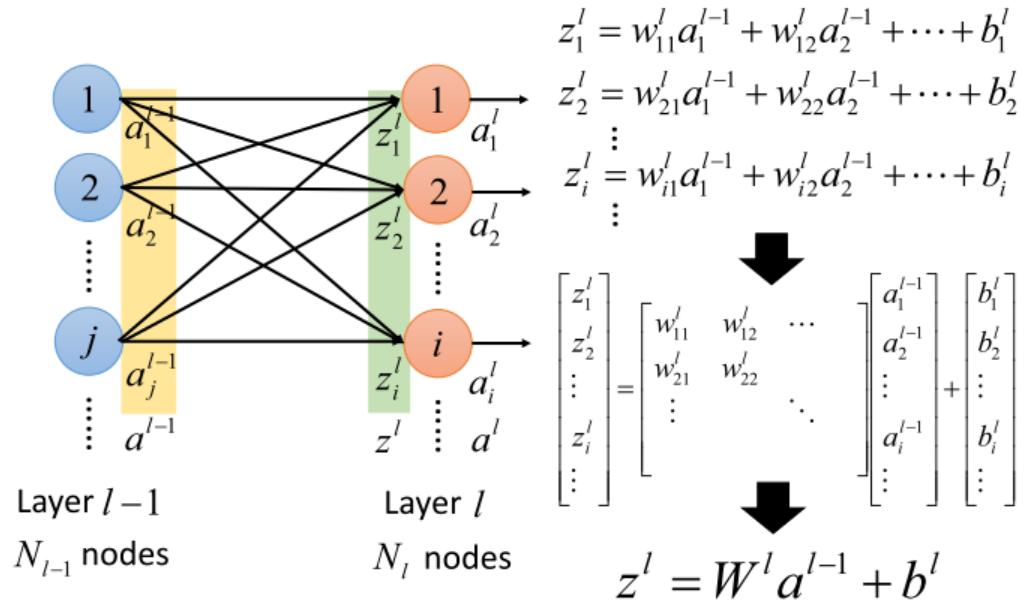
A neural network is a complex function:

$$f : R^N \rightarrow R^M$$

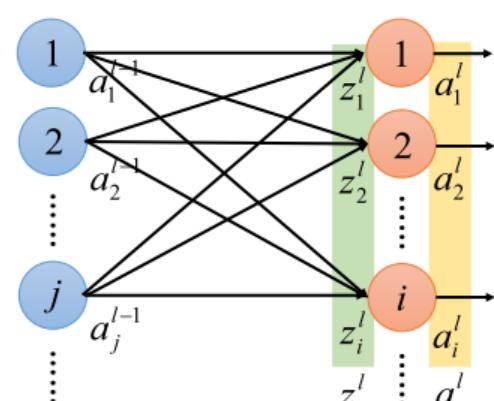
- Cascading the neurons to form a neural network.
Each layer is a simple function in the production line.



Relations between Layer Outputs



Relations between Layer Outputs



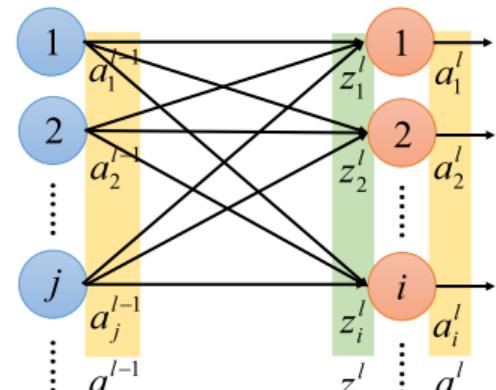
Layer $l-1$
 N_{l-1} nodes

Layer l
 N_l nodes

$$\begin{bmatrix} a_1^l \\ a_2^l \\ \vdots \\ a_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} \sigma(z_1^l) \\ \sigma(z_2^l) \\ \vdots \\ \sigma(z_i^l) \\ \vdots \end{bmatrix}$$

$$a^l = \sigma(z^l)$$

Relations between Layer Outputs



$$z^l = W^l a^{l-1} + b^l$$

$$a^l = \sigma(z^l)$$

$$a^l = \sigma(W^l a^{l-1} + b^l)$$

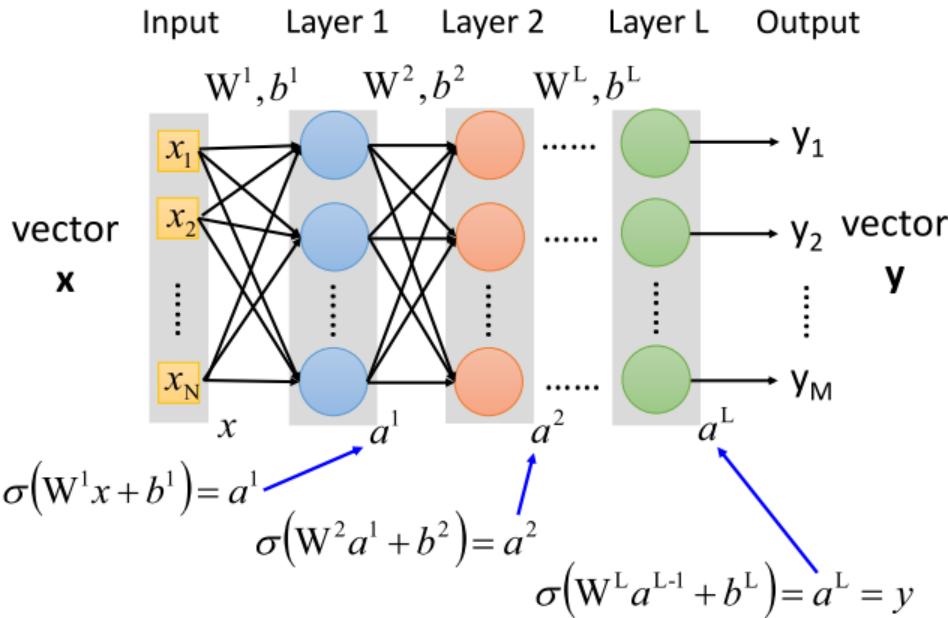
Layer $l-1$

N_{l-1} nodes

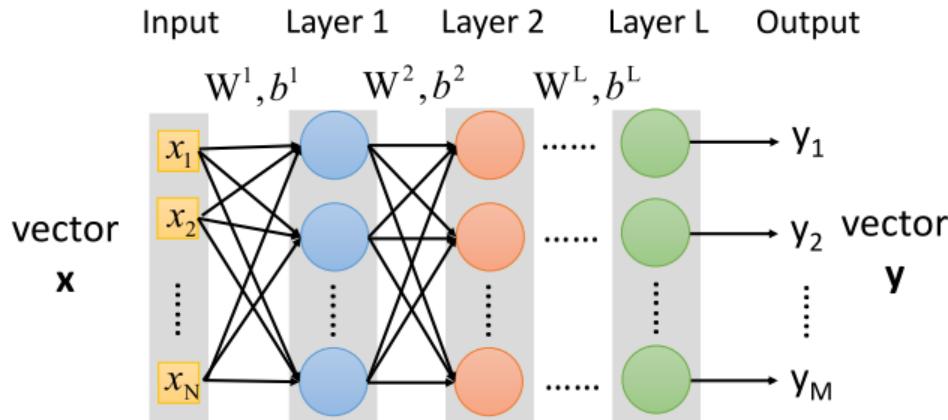
Layer l

N_l nodes

Function of Neural Network

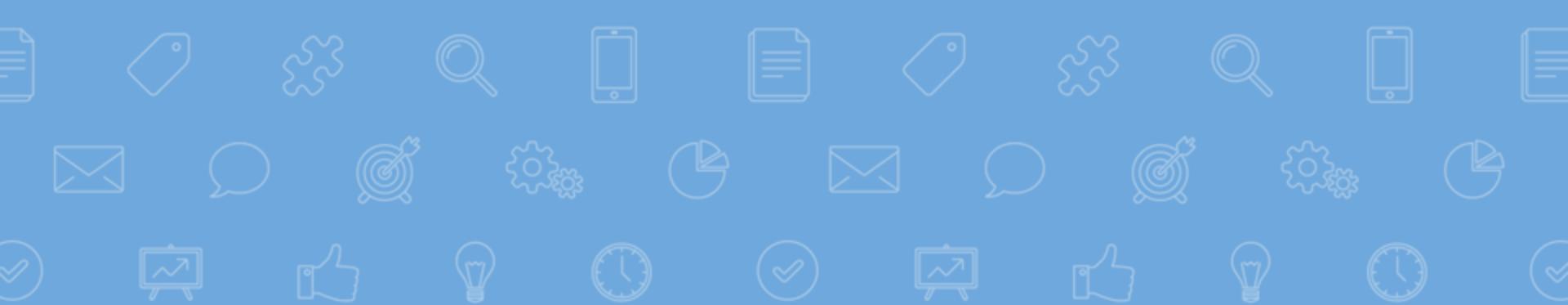


Function of Neural Network



$$y = f(x)$$

$$= \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$



How to find the best function?

Best Function = Best Parameters

$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

function set

because different parameters W
and b lead to different function

Formal way to define a function set:

$f(x; \underline{\theta}) \rightarrow$ parameter set

$$\theta = \{W^1, b^1, W^2, b^2 \dots W^L, b^L\}$$

Pick the “best”
function f^*



Pick the “best”
parameter set θ^*

Cost Function

- Define a function for parameter set $C(\theta)$
 - $C(\theta)$ evaluate how bad a parameter set is
 - The best parameter set θ^* is the one that minimizes $C(\theta)$

$$\theta^* = \arg \min_{\theta} C(\theta)$$

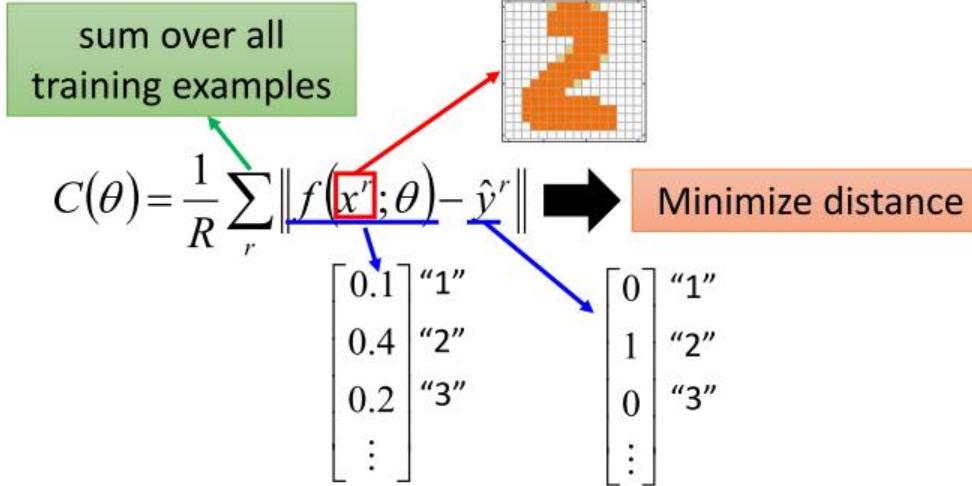
- $C(\theta)$ is called ***cost/loss/error function***
 - If you define the goodness of the parameter set by another function $O(\theta)$
 - $O(\theta)$ is called objective function

Cost Function

Given training data:

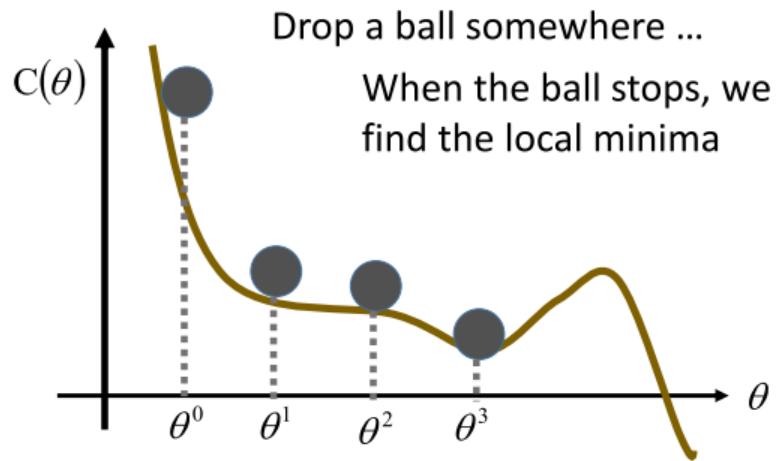
$$\{(x^1, \hat{y}^1), \dots, (x^r, \hat{y}^r), \dots, (x^R, \hat{y}^R)\}$$

- Handwriting Digit Classification



Gradient Descent – Idea

- For simplification, first consider that θ has only one variable

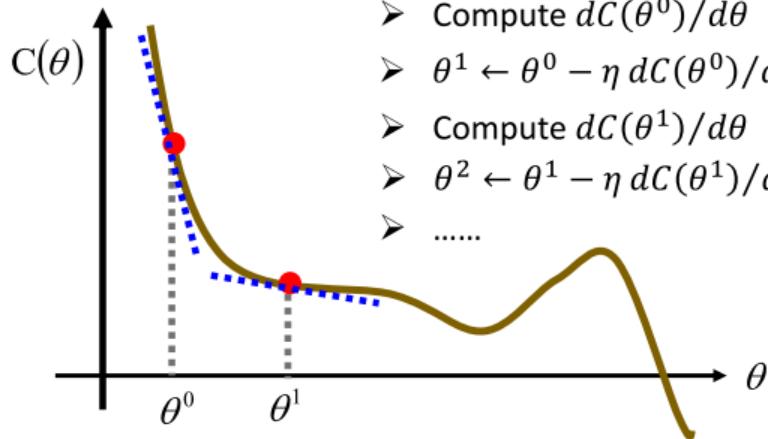


Gradient Descent – Idea

η is called
“*learning rate*”

- For simplification, first consider that θ has only one variable

- Randomly start at θ^0
- Compute $dC(\theta^0)/d\theta$
- $\theta^1 \leftarrow \theta^0 - \eta dC(\theta^0)/d\theta$
- Compute $dC(\theta^1)/d\theta$
- $\theta^2 \leftarrow \theta^1 - \eta dC(\theta^1)/d\theta$
-

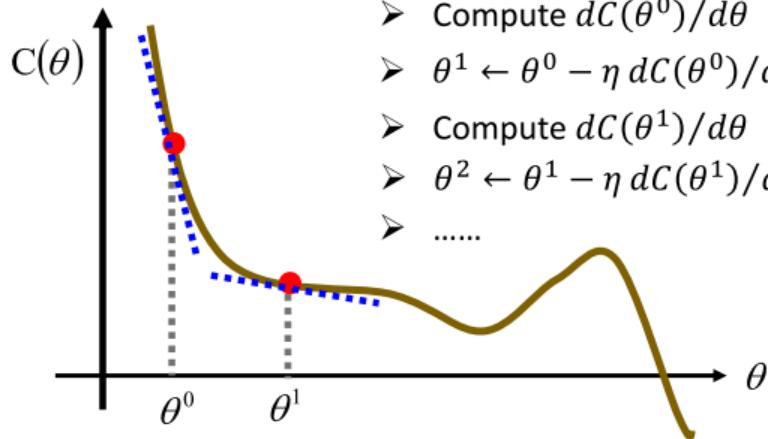


Gradient Descent – Idea

η is called
“*learning rate*”

- For simplification, first consider that θ has only one variable

- Randomly start at θ^0
- Compute $dC(\theta^0)/d\theta$
- $\theta^1 \leftarrow \theta^0 - \eta dC(\theta^0)/d\theta$
- Compute $dC(\theta^1)/d\theta$
- $\theta^2 \leftarrow \theta^1 - \eta dC(\theta^1)/d\theta$
-



Gradient Descent

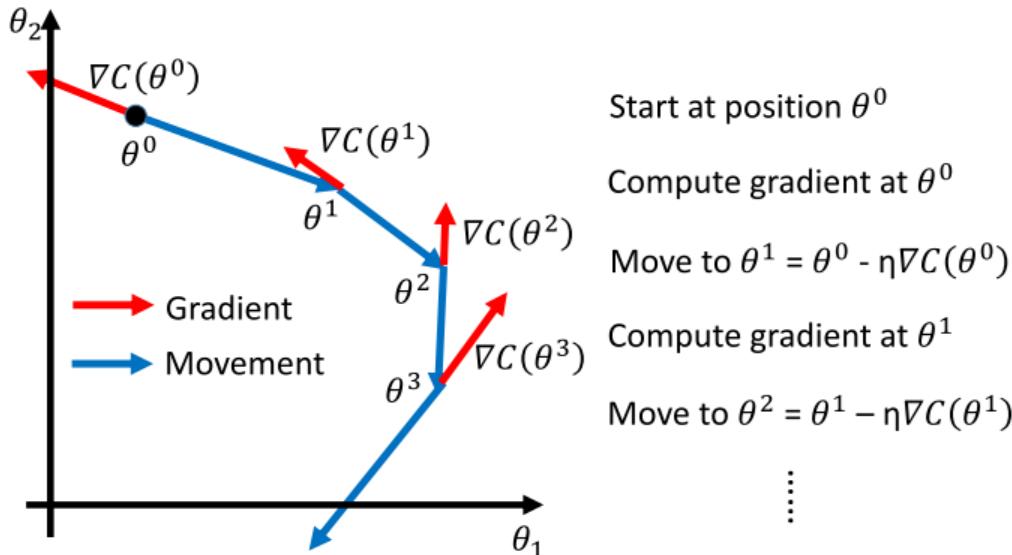
- Suppose that θ has two variables $\{\theta_1, \theta_2\}$

- Randomly start at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$
- Compute the gradients of $C(\theta)$ at θ^0 : $\nabla C(\theta^0) = \begin{bmatrix} \partial C(\theta_1^0)/\partial \theta_1 \\ \partial C(\theta_2^0)/\partial \theta_2 \end{bmatrix}$
- Update parameters

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \partial C(\theta_1^0)/\partial \theta_1 \\ \partial C(\theta_2^0)/\partial \theta_2 \end{bmatrix} \rightarrow \theta^1 = \theta^0 - \eta \nabla C(\theta^0)$$

- Compute the gradients of $C(\theta)$ at θ^1 : $\nabla C(\theta^1) = \begin{bmatrix} \partial C(\theta_1^1)/\partial \theta_1 \\ \partial C(\theta_2^1)/\partial \theta_2 \end{bmatrix}$
-

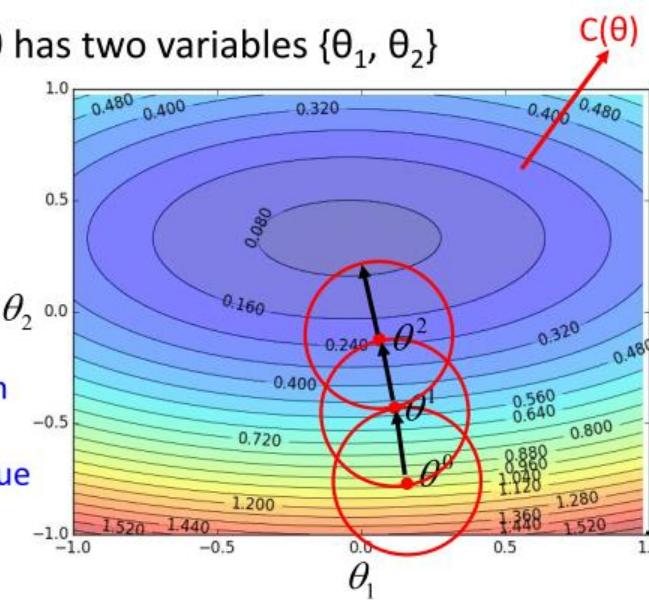
Gradient Descent



Formal Derivation of Gradient Descent

- Suppose that θ has two variables $\{\theta_1, \theta_2\}$

Given a point, we can easily find the point with the smallest value nearby. How?



Gradient Descent

This is the “learning” of machines in deep learning

→ Even alpha go using this approach.



I hope you are not too disappointed :p



"I work with models."

Dimension Reduction

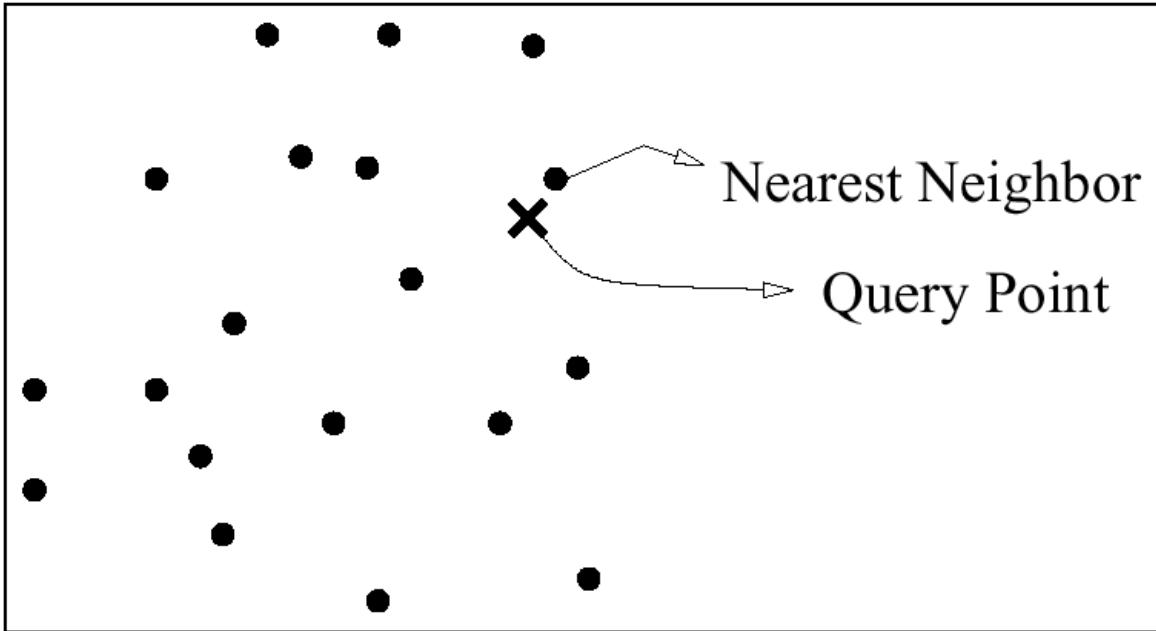
Thanks to the slides made by Prof. Hung-Yi Lee from NTU.

Preliminaries : Nearest Neighbor Search



- Given a collection of data points and a query point in m -dimensional metric space, find the data point that is closest to the query point
 - Variation: k -nearest neighbor
 - Relevant to clustering and similarity search
 - Applications: Geographical Information Systems, similarity search in multimedia databases

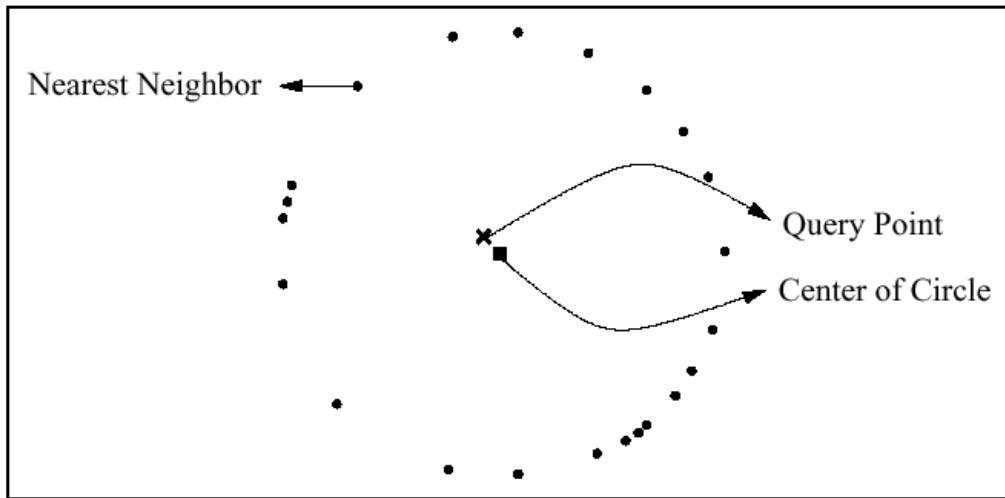
NN Search Con't



Source: [2]

Problems with High Dimensional Data

- A point's nearest neighbor (NN) loses meaning



Problems (Con't)



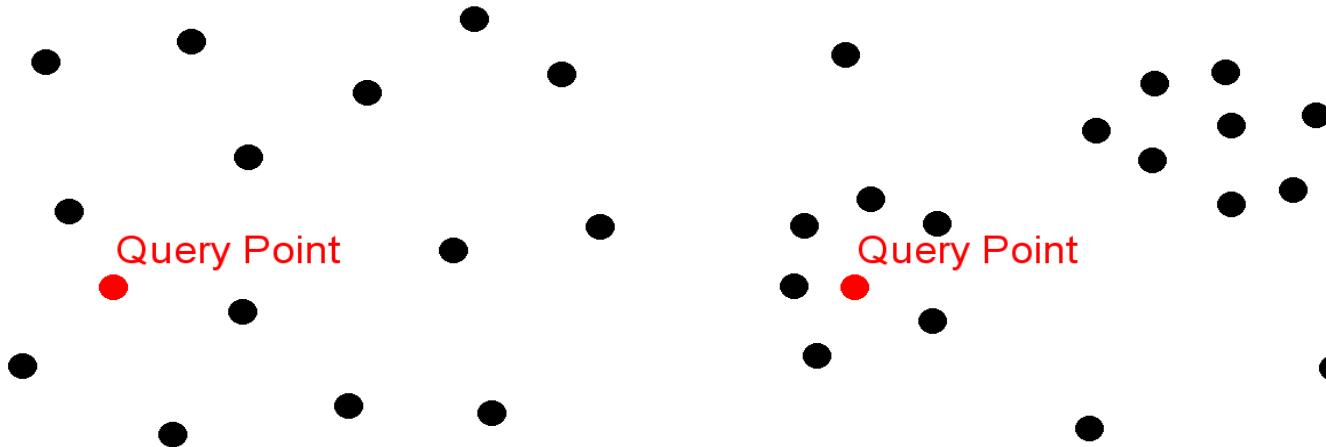
- NN query cost degrades – more strong candidates to compare with
- In as few as 10 dimensions, linear scan outperforms some multidimensional indexing structures (e.g. KD-tree, R* tree, SR tree)
- Biology and genomic data can have dimensions in the 1000's.

Problems (Con't)

- The presence of irrelevant attributes decreases the tendency for clusters to form
- Points in high dimensional space have high degree of freedom; they could be so scattered that they appear uniformly distributed

Problems Con't

- In which cluster does the query point fall?



The Curse

- Refers to the decrease in performance of query processing when the dimensionality increases
- In particular, under certain conditions, the distance between the nearest point and the query point equals the distance between the farthest and query point as dimensionality approaches infinity

Curse (Con't)

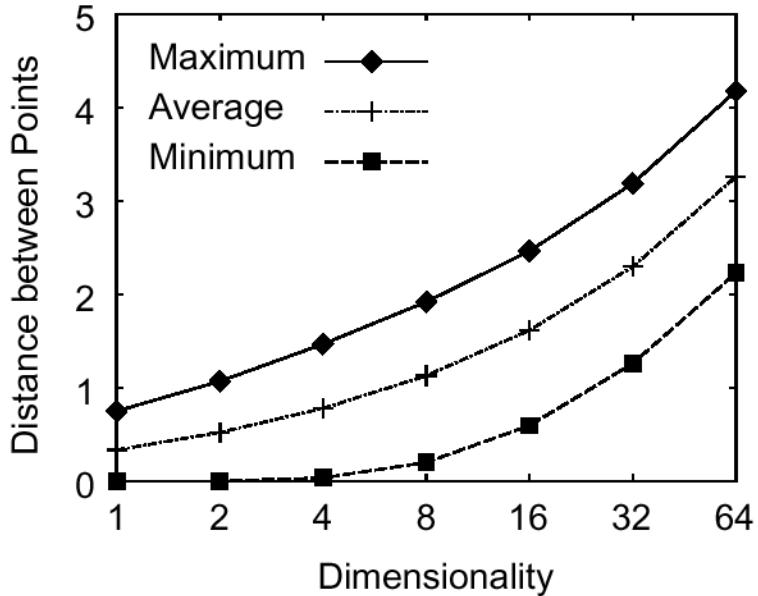


Figure 1. Distances among 100k points generated at random in a unit hypercube

Distributed Representation

- Clustering: an object must belong to one cluster

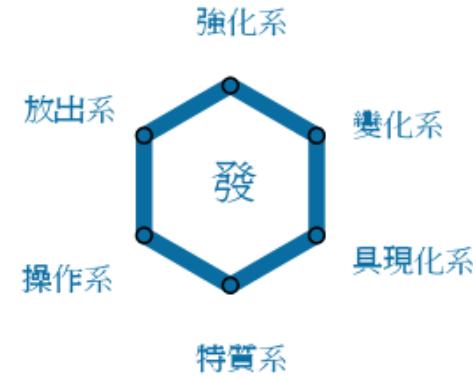
小傑是強化系

- Distributed representation

Dimension Reduction

小傑是

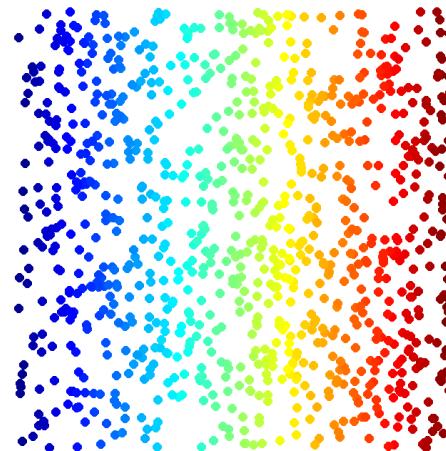
強化系	0.70
放出系	0.25
變化系	0.05
操作系	0.00
具現化系	0.00
特質系	0.00



Dimension Reduction



Looks like 3-D



Actually, 2-D

<http://reuter.mit.edu/blue/images/research/manifold.png>

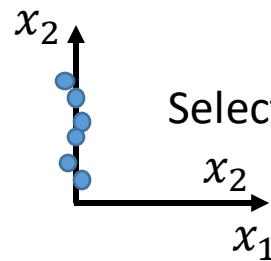
<http://archive.cnx.org/resources/51a9b2052ae167db310fda5600b89badea85eae5/isomapCNXtrue1.png>

Dimension Reduction



The dimension of z
would be smaller than
x

- Feature selection



- Principle component analysis (PCA)
[Bishop, Chapter 12]

$$z = Wx$$

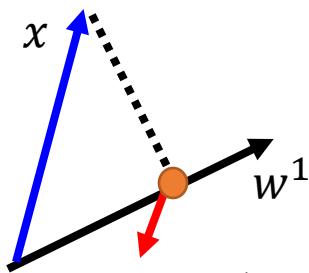
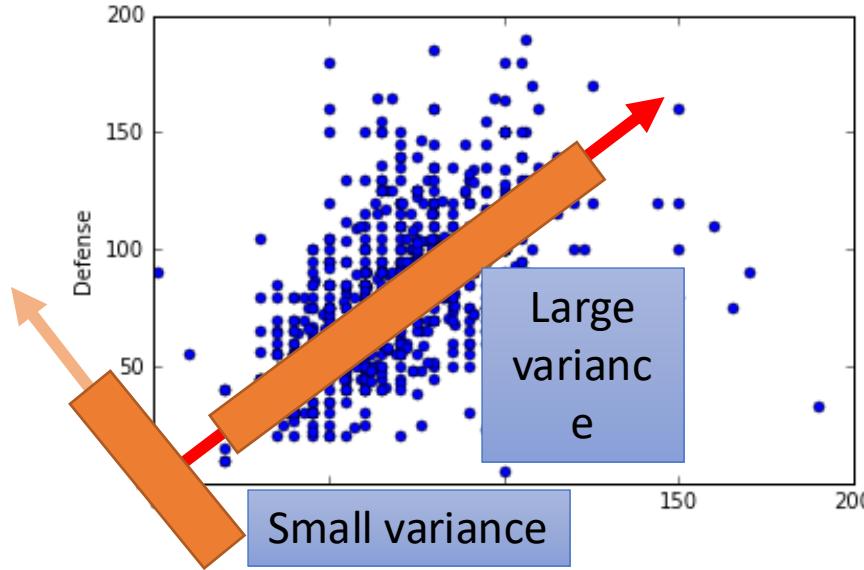
Principle Component Analysis (PCA)

PCA

$$z = Wx$$

Reduce to 1-D:

$$z_1 = w^1 \cdot x$$



Project all the data points x onto w^1 ,
and obtain a set of z_1

We want the variance of z_1 as large as
possible

$$\text{Var}(z_1) = \sum_{z_1} (z_1 - \bar{z}_1)^2 \quad \|w^1\|_2 = 1$$

PCA

$$z = Wx$$

Reduce to 1-D:

$$z_1 = w^1 \cdot x$$

$$z_2 = w^2 \cdot x$$

$$W = \begin{bmatrix} (w^1)^T \\ (w^2)^T \\ \vdots \end{bmatrix}$$

Orthogonal
matrix

Project all the data points x onto w^1 ,
and obtain a set of z_1

We want the variance of z_1 as large as
possible

$$\text{Var}(z_1) = \sum_{z_1} (z_1 - \bar{z}_1)^2 \quad \|w^1\|_2 = 1$$

We want the variance of z_2 as large as
possible

$$\text{Var}(z_2) = \sum_{z_2} (z_2 - \bar{z}_2)^2 \quad \|w^2\|_2 = 1$$

 $w^1 \cdot w^2 = 0$

Warning of Math

PCA

$$z_1 = w^1 \cdot x$$

$$\bar{z}_1 = \frac{1}{N} \sum z_1 = \frac{1}{N} \sum w^1 \cdot x = w^1 \cdot \frac{1}{N} \sum x = w^1 \cdot \bar{x}$$

$$Var(z_1) = \frac{1}{N} \sum_{z_1} (z_1 - \bar{z}_1)^2$$

$$= \frac{1}{N} \sum_x (w^1 \cdot x - w^1 \cdot \bar{x})^2$$

$$= \frac{1}{N} \sum (w^1 \cdot (x - \bar{x}))^2$$

$$= \frac{1}{N} \sum (w^1)^T (x - \bar{x})(x - \bar{x})^T w^1$$

$$= (w^1)^T \boxed{\frac{1}{N} \sum (x - \bar{x})(x - \bar{x})^T} w^1$$

$$= (w^1)^T Cov(x) w^1 \quad S = Cov(x)$$

$$(a \cdot b)^2 = (a^T b)^2 = a^T b a^T b \\ = a^T b (a^T b)^T = a^T b b^T a$$

Find w^1 maximizing

$$(w^1)^T S w^1$$

$$\|w^1\|_2 = (w^1)^T w^1 = 1$$

Find w^1 maximizing $(w^1)^T S w^1$ $(w^1)^T w^1 = 1$

$S = Cov(x)$ Symmetric positive-semidefinite
(non-negative eigenvalues)

Using Lagrange multiplier [Bishop, Appendix E]

$$g(w^1) = (w^1)^T S w^1 - \alpha((w^1)^T w^1 - 1)$$

$$\left. \begin{array}{l} \frac{\partial g(w^1)}{\partial w_1^1} = 0 \\ \frac{\partial g(w^1)}{\partial w_2^1} = 0 \\ \vdots \end{array} \right\} \begin{array}{l} S w^1 - \alpha w^1 = 0 \\ S w^1 = \alpha w^1 \\ (w^1)^T S w^1 = \alpha (w^1)^T w^1 \\ = \alpha \end{array}$$

w^1 : eigenvector

Choose the maximum one

w^1 is the eigenvector of the covariance matrix S
Corresponding to the largest eigenvalue

Find w^2 maximizing $(w^2)^T S w^2 \quad (w^2)^T w^2 = 1 \quad (w^2)^T w^1 = 0$

$$g(w^2) = (w^2)^T S w^2 - \alpha((w^2)^T w^2 - 1) - \beta((w^2)^T w^1 - 0)$$

$$\left. \begin{array}{l} \partial g(w^2)/\partial w_1^2 = 0 \\ \partial g(w^2)/\partial w_2^2 = 0 \\ \vdots \end{array} \right\} \begin{aligned} & S w^2 - \alpha w^2 - \beta w^1 = 0 \\ & \boxed{0} - \alpha \boxed{0} - \beta \boxed{1} = 0 \\ & = ((w^1)^T S w^2)^T = (w^2)^T S^T w^1 \\ & = (w^2)^T S w^1 = \lambda_1 (w^2)^T w^1 = 0 \end{aligned}$$
$$S w^1 = \lambda_1 w^1$$

$$\beta = 0: \quad S w^2 - \alpha w^2 = 0 \quad S w^2 = \alpha w^2$$

w^2 is the eigenvector of the covariance matrix S

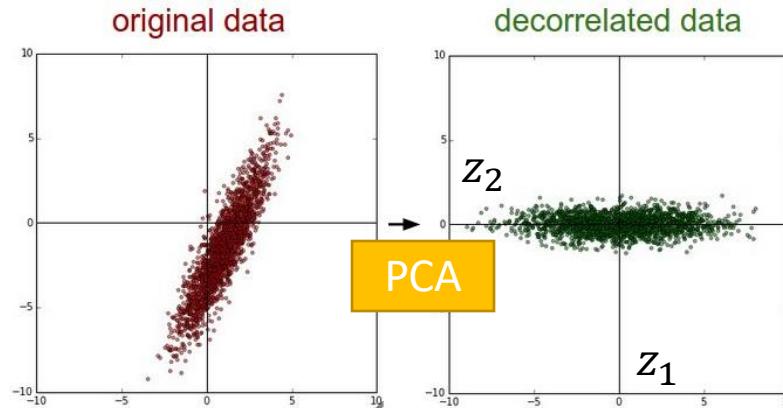
Corresponding to the 2nd largest eigenvalue λ_2

PCA - decorrelation

$$z = Wx$$

$$\text{Cov}(z) = D$$

Diagonal matrix



$$\text{Cov}(z) = \frac{1}{N} \sum (z - \bar{z})(z - \bar{z})^T = WSW^T \quad S = \text{Cov}(x)$$

$$= WS[w^1 \quad \dots \quad w^K] = W[S_w^1 \quad \dots \quad S_w^K]$$

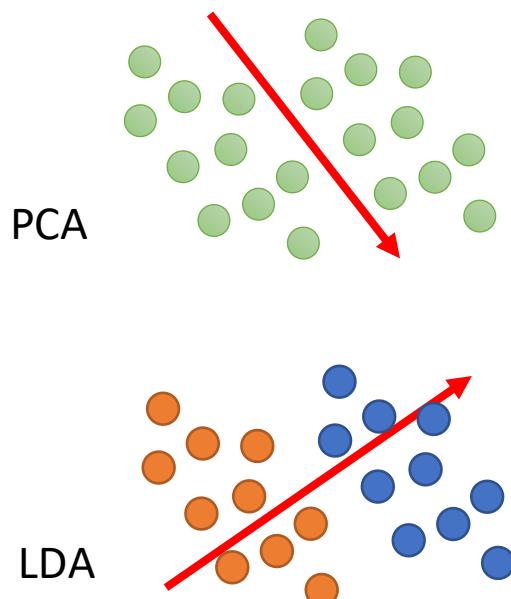
$$= W[\lambda_1 w^1 \quad \dots \quad \lambda_K w^K] = [\lambda_1 Ww^1 \quad \dots \quad \lambda_K Ww^K]$$

$$= [\lambda_1 e_1 \quad \dots \quad \lambda_K e_K] = D \quad \text{Diagonal matrix}$$

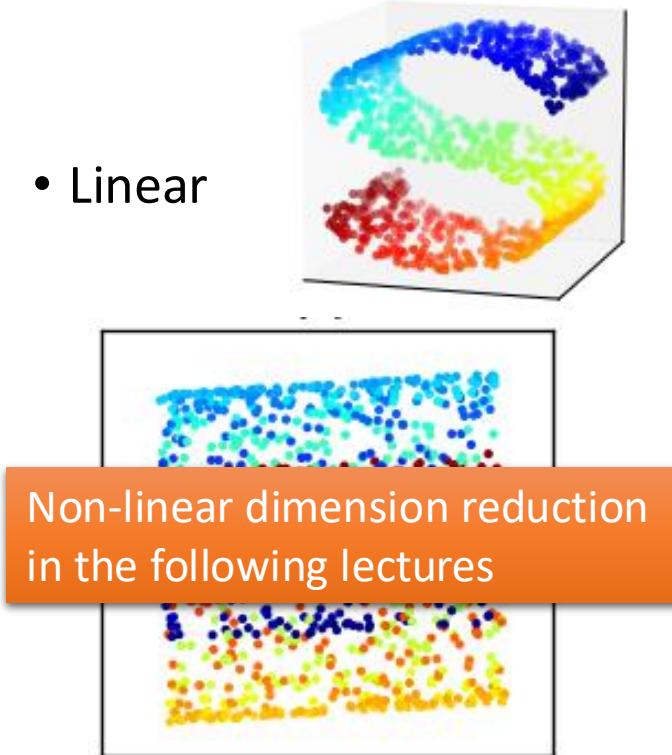
End of Warning

Weakness of PCA

- Unsupervised



- Linear



http://www.astroml.org/book_figures/chapter7/fig_S_manifold_PCA.html

With the power of PCA, let's
become the master of...

Pokémon!!



PCA - Pokémon

- Inspired from: <https://www.kaggle.com/strakul5/d/abcsds/pokemon/principal-component-analysis-of-pokemon-data>
- 800 Pokemons, 6 features for each (HP, Atk, Def, Sp Atk, Sp Def, Speed)
- How many principle components?

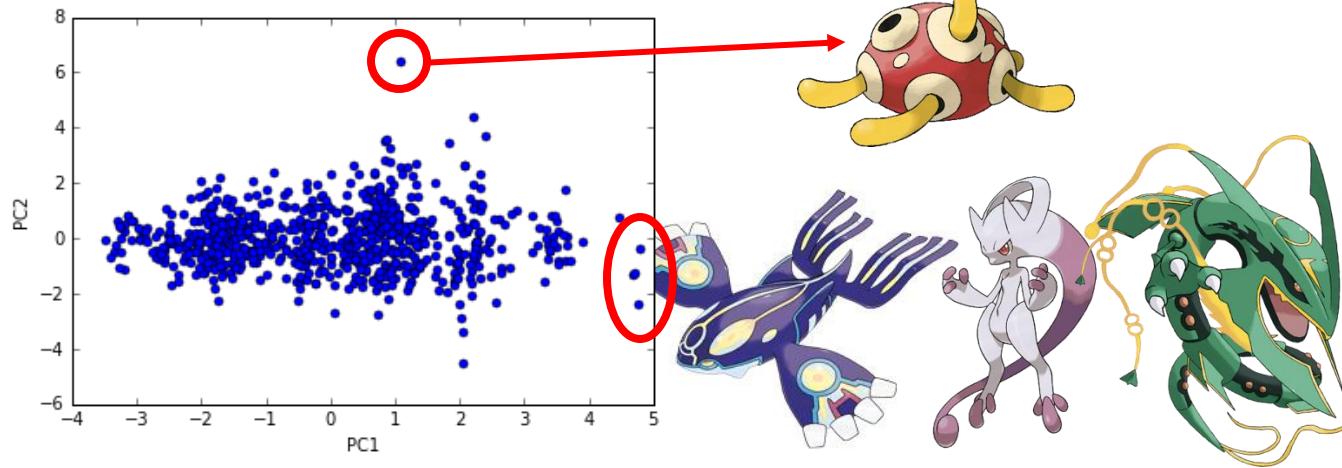
$$\frac{\lambda_i}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 + \lambda_6}$$

	λ_1	λ_2	λ_3	λ_4	λ_5	λ_6
ratio	0.45	0.18	0.13	0.12	0.07	0.04

Using 4 components is good enough

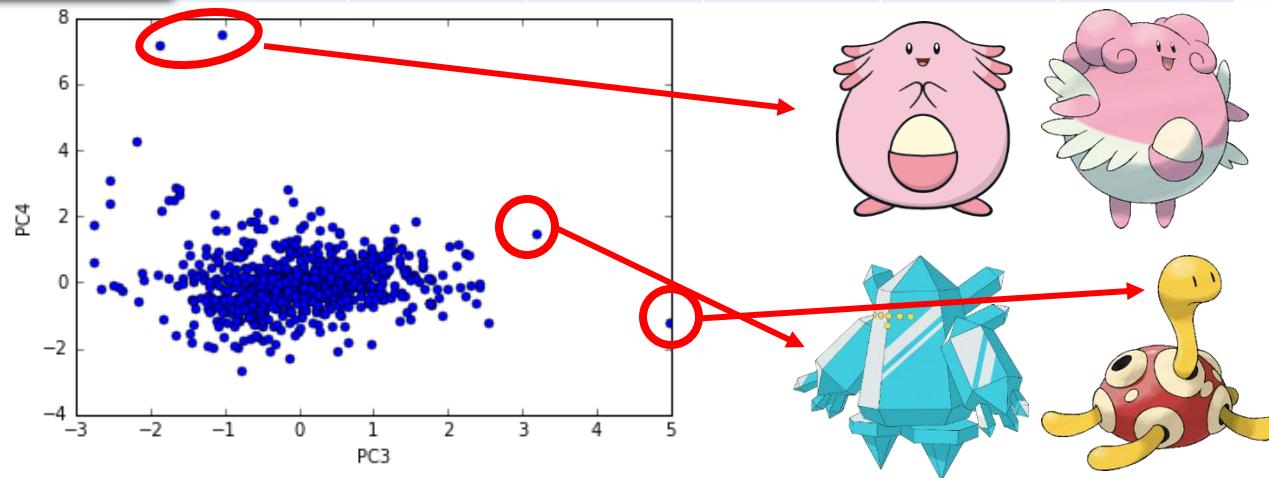
PCA - Pokémon

	HP	Atk	Def	Sp Atk	Sp Def	Speed	
PC1	0.4	0.4	0.4	0.5	0.4	0.3	強度
PC2	0.1	0.0	0.6	-0.3	0.2	-0.7	
PC3	-0.5	-0.6	0.1	0.3	0.6	0.2	防禦(犠牲速度)
PC4	0.7	-0.4	-0.4	0.1	0.2	-0.3	



PCA - Pokémon

	HP	Atk	Def	Sp Atk	Sp Def	Speed
PC1	0.4	0.4	0.4	0.5	0.4	0.3
PC2	0.1	0.0	0.6	-0.3	0.2	-0.7
PC3	-0.5	-0.6	0.1	0.3	0.6	
生命力強	0.7	-0.4	-0.4	0.1	0.2	特殊防禦(犧牲 攻擊和生命)



Q and A

Thank you