

Introduction of Machine Learning

No.

Date

→ 幫機器找到好的 model to loss function

人工智慧 → 目標：讓機器同人一樣聰明

機器學習 → 手段：寫程式讓機器具學習能力

深度學習 → 機器學習的一種

如何賦予機器智能？

hand-coded rule: 用程式事先設定 → 無法超越人類

Machine learning: Looking for a function that fits the input datas.

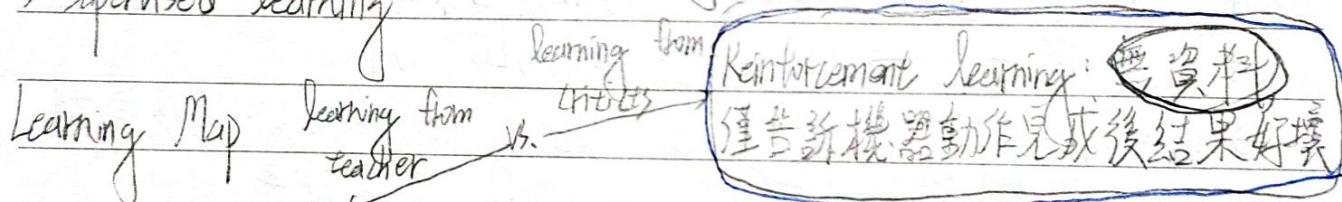
Steps:

1. prepare a "Model" (a set of functions)

2. Give "Training Data" → the AI judge the functions in the set according to the Training Data

3. test whether the AI can correctly work

→ supervised learning



Supervised Learning

Regression → 輸出是數值

Semi-supervised Learning:

Data → part labeled/part unlabeled

Classification | Binary: 二元分類 → Yes or No

Multi-class: 將資料分為多類

Transfer Learning:

use old models for similar problems

Linear Function: 線性函數

Non-linear Function → more complex

ex: Deep learning, SVM, decision tree, K-NN

↳ Structure learning → 輸出有結構性的資料

Unsupervised learning:

machine self-learning
only outputs, no inputs

* Training Data: Input/Output pairs → for target function
→ Output = Label

My scenario ~~my~~ tasks my method
charculture

Regression: 輸出數值 (Scalar)

Ex: 寶可夢進化 CP 值預測

→ input: 寶可夢 (X) output: 進化後 CP 值 (y)

Step 1: Model = a set of function

$$(1) y = b + w \cdot X_{CP} \rightarrow \text{Linear model: } y = \text{bias} + \sum w_i x_i \text{ 的各種數值}$$

bias: 截距 weight: 權重

Step 2: Goodness of function \rightarrow Loss function (L)

→ input: a function; output: how bad it is

$$\rightarrow L(f) = L(w, b) = \sum_{n=1}^N (y_n - (b + w \cdot X_{CP}^n))^2 \rightarrow \text{估測誤差}$$

→ 這次 loss function 僅考慮 y 與 \hat{y} 的距離 \rightarrow 圖示以 w, b 為兩軸 $L(w, b)$ 第三軸

Step 3: Best Function \rightarrow Gradient Descent

→ 由 w 判斷 \downarrow 好壞: y 軸 $L(w)$; x 軸 w 正: $w \downarrow$

→ random 選 $-w_0$ \rightarrow compute $\frac{dL}{dw}|_{w=w_0}$ 負: $w \uparrow$

→ $w_{new} = w_{old} \pm \eta \frac{dL}{dw}|_{w=w_{old}} \rightarrow$ (eta): learning rate = 加減幅度常數

→ 多次循環 \rightarrow 到 local minimum \rightarrow 在 regression 就是 global

→ 考慮 w 及 b 為第二個因素, 同理 \rightarrow 偏微分 (分別微分)

$$L(w, b) = \sum_{n=1}^{10} (y_n - (b + w \cdot X_{CP}^n))^2$$

$$\frac{\partial L}{\partial w} = \sum_{n=1}^{10} 2(y_n - b - w \cdot X_{CP}^n)(-X_{CP}^n); \frac{\partial L}{\partial b} = \sum_{n=1}^{10} 2(y_n - b - w \cdot X_{CP}^n)(-1)$$

優化:

1. more complex Model:

error | training | testing | \leftarrow 越複雜 \rightarrow training error \downarrow

一次	31.9	35.1	testing error 不一定
----	------	------	-------------------

二次	15.4	16.4	\leftarrow overfitting \rightarrow model 太複雜 \rightarrow error 低
----	------	------	--

三次	15.3	16.1	\leftarrow underfitting \rightarrow model 太簡單 \rightarrow error 高
----	------	------	---

四次	14.9	28.8	
----	------	------	--

五次	12.8	232.5	
----	------	-------	--

如何在現實中 testing data Error 更小？

- N-fold validation：將 training data 分為 N 份，每次選取一份作為 validation set (模擬 testing data)，其他作為 training set
- 將所有 model 跑過並輸出 Error 再取 average error
- 取 an error 最小的嘗試 public testing data
- 這種結果較能應在其他 data 上的 error

* Gradient Descent tips:

tip 1: tuning your learning rate

[太小：下降太慢]

[太大：無法走到最低點]

1. 畫兩張次數 update 下降趨勢圖： No. of parameter updates

tip 2: Adaptive Learning Rates → 一系列：Ada...

Maybe → Reduce learning rate by some factor every epochs.

∴ Learning rate cannot be one-size-fit-all

∴ Give different parameters different learning rates

Ex:

$$\text{Adagrad: } w^{t+1} = w^t - \eta^t g^t \rightarrow w^{t+1} = w^t - \frac{\eta^t}{\sum_{j=0}^t (g_j)^2} g^t$$

$$\Rightarrow \sigma^2 = \sqrt{\left(\sum_{j=0}^t (g_j)^2\right)} = \sqrt{\left[\sum_{j=0}^t (g_j)^2\right] \frac{1}{t+1}} \quad \eta^t = \frac{\eta}{t+1} \Rightarrow \frac{\eta^t}{\sigma^2} = \frac{\eta}{\sqrt{t+1}}$$

but $w^{t+1} = w^t - \frac{\eta^t}{\sum_{j=0}^t (g_j)^2} g^t$ gradient ↑ step ↑ why?

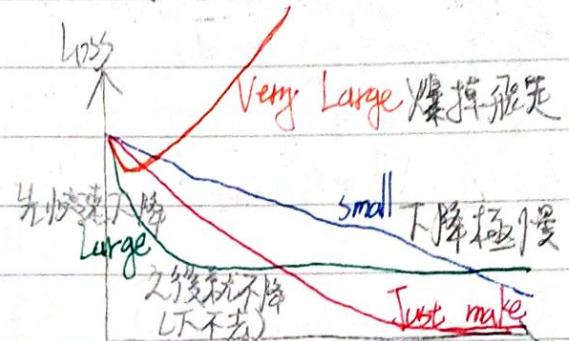
Intuitive Reason: how surprise? 讓參數反差

∴ Gradient descent 認為微分下與最低點距離↑ → 只在二次式微分成立

$$\therefore \text{best step} = \frac{1 - \text{一次微分}}{\text{二次微分}} \rightarrow \text{Adagrad: } w^{t+1} \leftarrow w^t - \frac{\eta(g^t) \text{ 一次微分}}{\left(\sum_{j=0}^t (g_j)^2\right) \text{ 二次微分}}$$

→ 為何不直接計算二次微分？

數據量↑ 可能耗時太多



搜集更多資料

2. hidden factor → 僅

→ 使用 δ (delta): species 的 true v false \Rightarrow 不同物種決定

→ 考慮 weight - height - Hp ... \Rightarrow redesign Model

加入 $y = y' + Hp$ 一次式 + weight = 二次式 + height = 二次式

\Rightarrow 結果: overfitting!

若不考慮 bias? \Rightarrow bias 只影響 function 的平滑

→ Regularization (Step 2) 不影響平滑程度 常數

Loss function 改為: $L = \frac{1}{n} (y_n - b + \sum w_i x_i)^2 + \lambda (\sum w_i^2)$

\Rightarrow 期待更平滑 function \Rightarrow 結果波動更小，較不受 noise 影響

\Rightarrow 入 α training error \Rightarrow 更考慮 function 本質 & 雜訊

Where does the error come from [due to bias: 偏心率

解答函式

[due to variance: 精準度

$f = \hat{f}$ (寶可夢) \Rightarrow only Niantic knows

but 可由 training data find f^* \Rightarrow 最佳解 靠近於 f

low variance

high

\Rightarrow Simpler model is less influenced by the sampled data \rightarrow lower variance

\rightarrow 極端例子: $f = 0$ 常數 \Rightarrow 完全不受 data 影響

training data loss function 大

\rightarrow large bias \Rightarrow 重新找 function

testing data loss function 大

\rightarrow large variance \Rightarrow More data

$\checkmark \Rightarrow$ regularization \rightarrow 強制曲線平滑 但 可能傷害 bias

1. 利用 weigh 寫定 regularization 程度

- Loss = sum over all training example
- ~~Stochastic Gradient Descent~~ * Gradient Descent: $\theta^{t+1} = \theta^t - \eta \nabla L(\theta^t)$
- Loss for only one example: $L^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$
- 只隨機考慮一個取樣點
- 在每個 example 都更新一次參數 → 20 times faster
(天下武功，唯快不破 XD)

- tip 3: Feature Scaling: 希望讓不同輸入元素分佈差不多
- 讓各個 Feature 的 Loss surface 更接近正圓
- 數據標準化 $\Rightarrow \text{variation} = 0$

Gradient Descent Theory

* Warning of math:

if 取 \downarrow 隨機初始值 \Rightarrow can easily find the point with the smallest value nearby.
but how?

→ Taylor Series 泰勒級數 (泰勒展開式):

Let $h(x)$ be any function infinitely differentiable around $x = x_0$

$$h(x) = \sum_{k=0}^{\infty} \frac{h^{(k)}(x_0)}{k!} (x-x_0)^k \Rightarrow \text{when } x \rightarrow x_0 \Rightarrow h(x) \approx h(x_0) + h'(x_0)(x-x_0)$$

高次項都 $\rightarrow 0$

→ Multivariable Taylor Series: if $x, y \rightarrow x_0, y_0$

$$h(x, y) \approx h(x_0, y_0) + \frac{\partial h(x_0, y_0)}{\partial x} (x-x_0) + \frac{\partial h(x_0, y_0)}{\partial y} (y-y_0)$$

若範圍夠小，可用泰勒展開式簡化 loss function

$$L(\theta) \approx L(a, b) + \frac{\partial L(a, b)}{\partial \theta_1} (\theta_1 - a) + \frac{\partial L(a, b)}{\partial \theta_2} (\theta_2 - b) = S + U(\theta_1 - a) + V(\theta_2 - b)$$

求 $U\Delta\theta_1 + V\Delta\theta_2$ 之 $\min \sqrt{\Delta\theta_1^2 + \Delta\theta_2^2} \leq d^2$

$$\begin{bmatrix} \Delta\theta_1 \\ \Delta\theta_2 \end{bmatrix} = -h \begin{bmatrix} a \\ b \end{bmatrix} \Rightarrow \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - h \begin{bmatrix} U \\ V \end{bmatrix}$$

若沒有就不成立

More Limitation of Gradient Descent

1. Stuck at local minima

2. Stuck at saddle point $\Rightarrow \partial L / \partial w \geq 0$

實作時可能遇幾次

3. Stuck at a point that's very slow at the plateau $\Rightarrow \partial L / \partial w \geq 0$ 停止點

Classification:

Ex: 寶可夢分類

input: 一隻寶可夢 \rightarrow 轉為可輸入數值 vector

output: 寶可夢種類 ($\times 18$)

方法 1: 用 regression 裝煤? if binary 極端數值 (離二分界太遠者) 會造成分界並非最佳 \rightarrow 會嘗試調整分界平復子極端數值
if mutipul \rightarrow 極難做到 \wedge 並非理想方法

\rightarrow Ideal Alternatives

Function Model:

Binary: a Model contains if...else

Loss function: $L(f) = \sum \delta(f(x^n) \neq y^n)$ \rightarrow 判斷是否正確

Find the best function: ex: perception, SVM

\rightarrow Two classes \rightarrow Binary: Class 1, Class 2

Given a x : $P(L_1|x) = \frac{P(X|L_1)P(L_1)}{P(X|L_1)P(L_1) + P(X|L_2)P(L_2)}$ \rightarrow 看 X 比較可能來自哪個 class

= Generative Model: $P(a) = P(a|L_1)P(L_1) + P(a|L_2)P(L_2)$

\rightarrow Feature: 考慮 Df 跟 $sd(Df)$ 畫出分布圖

\rightarrow 假設點都是 Gaussian Distribution (常態分布) 取樣

\rightarrow 計算 μ (= mean 平均) 及 Σ (分布狀況 = 聚集 \leftrightarrow 分散)

\rightarrow 即可算新點出現機率 \hookrightarrow covariance 跟異數

To find best function: Maximum Likelihood

$$\hat{L}(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) \cdots f_{\mu, \Sigma}(x^m) \rightarrow \text{水系}$$

$$f_{\mu, \Sigma}(x) = \left(\frac{1}{(2\pi)^{\frac{m}{2}}} \right) \left(\frac{1}{|\Sigma|^{\frac{1}{2}}} \right) \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$

$$\rightarrow \text{Class 1: water } \mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 814 & 321 \\ 321 & 929 \end{bmatrix}$$

$$\text{Class 2: normal } \mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 841 & 422 \\ 422 & 685 \end{bmatrix}$$

\rightarrow Now we can do classification: if $P(L_1|x) > 0.5 \Rightarrow x$ belongs to class 1

results \Rightarrow testing data: 47% accuracy

All? total, hp, att, sp att, de, sp de, speed (7 features)

μ^1, μ^2 : 7-dim vector; Σ^1, Σ^2 : 7×7 matrices \Rightarrow 54% accuracy

於是嘗試 \Rightarrow 讓 2 class 有相同 Σ : μ^1, μ^2 相同, $\Sigma = \frac{74}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$
boundary 由曲線 \rightarrow linear \Rightarrow accuracy: 54% \rightarrow 73%

總結 Three steps:

1. Function Set (Model)

2. Goodness of a function: The mean μ and covariance Σ that maximize the likelihood (the probability of generating data)

3. Find the best function

\rightarrow Probability Distribution (你喜歡的就好 XD)

\Rightarrow binary features: Bernoulli distribution

\Rightarrow all the dimensions are independent: Naive Bayes Classifier

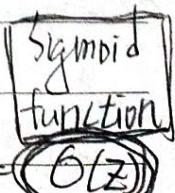
$$\rightarrow \text{Posterior Probability: } P(L_1|x) = \frac{P(x|L_1)P(L_1)}{P(x|L_1)P(L_1) + P(x|L_2)P(L_2)}$$

分子分母同除分子

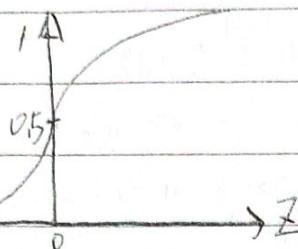
$$z = \ln \frac{P(x|L_1)P(L_1)}{P(x|L_2)P(L_2)}$$

\rightarrow Sigmoid function

$$G(z) = \frac{1}{1 + \exp(-z)}$$



* Sigmoid function $\sigma(z)$: if $x \rightarrow \infty, y \rightarrow 1$; $x \rightarrow -\infty, y \rightarrow 0$



$$z = \ln \frac{P(x|L_1)P(L_1)}{P(x|L_2)P(L_2)} = \ln \frac{P(x|L_1)}{P(x|L_2)} + \ln \frac{P(L_1)}{P(L_2)}$$

$$= \left(\frac{1}{N_1} \frac{1}{|\Sigma^1|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\} \right) + \frac{N_1}{N_1 + N_2}$$

$$= \frac{1}{N_2} \frac{1}{|\Sigma^2|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\} - \frac{N_2}{N_1 + N_2}$$

$$= \ln \left[\frac{1}{|\Sigma^1|^{\frac{1}{2}}} - \frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)] \right] + \frac{N_1}{N_2}$$

$$\Rightarrow (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) = x^T (\Sigma^1)^{-1} x - x^T (\Sigma^1)^{-1} \mu^1 - (\mu^1)^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$= x^T (\Sigma^1)^{-1} x - 2(\mu^1)^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} \mu^1$$

$$\Rightarrow (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) = x^T (\Sigma^2)^{-1} x - 2(\mu^2)^T (\Sigma^2)^{-1} x + (\mu^2)^T (\Sigma^2)^{-1} \mu^2$$

$$\therefore z = \ln \frac{1}{|\Sigma^1|^{\frac{1}{2}}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 + \frac{1}{2} x^T (\Sigma^2)^{-1} x$$

$$- (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

if $\Sigma_1 = \Sigma_2 = \Sigma$

$$z = \underbrace{[(\mu^1 - \mu^2)^T \Sigma^{-1}] x}_{w} - \underbrace{\frac{1}{2} \Sigma^{-1} (\mu^1 \mu^1 - \mu^2 \mu^2)}_{b} + \ln \frac{N_1}{N_2} \Rightarrow P(L_1|x) = \sigma(w \cdot x + b)$$

∴ In generative model, we estimate $N_1, N_2, \mu^1, \mu^2, \Sigma \Rightarrow$ we have w and b
So how about directly finding w and b?

Logistic Regression vs. Linear Regression

Step 1:

$$f_{w,b}(x) = \sigma \left(\sum_i w_i x_i + b \right)$$

Output: 0 ~ 1

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

Output: any value

Step 2: Training data $\rightarrow (X^n, y^n)$

y^n : 1 for class 1, 0 for class 2

$$L(f) = \sum_n L(f(x^n), y^n)$$

$$\therefore w^*, b^* = \arg \max_{w,b} L(w,b) = w^*, b^* = \arg \min_{w,b} -\ln L(w,b)$$

y^n : a real number

$$L(f) = \frac{1}{2} \sum_n L(f(x^n) - y^n)^2$$

* Cross entropy: $L(f(x^n), \hat{y}^n) = -[\hat{y}^n \ln f_w b(x^n) + (1 - \hat{y}^n) \ln (1 - f_w b(x^n))]$

Q: Why not also use square error as linear regression?

$$f_w b(z) = G(z) = \frac{1}{1 + e^{-z}}$$

Step 3: 其實跟 linear regression 的 gradient descend 一樣

$$\frac{\partial L(w, b)}{\partial w_i} = \sum_n - [\hat{y}^n \frac{\partial \ln f_w b(x^n)}{\partial z} + (1 - \hat{y}^n) \frac{\partial \ln (1 - f_w b(x^n))}{\partial z}] x_i^n$$

$$\frac{\partial \ln f_w b(x)}{\partial w_i} = \frac{\partial \ln f_w b(x)}{\partial z} \frac{\partial z}{\partial w_i} = \frac{\partial \ln G(z)}{\partial z} = \frac{1}{G(z)} \frac{\partial G(z)}{\partial z}$$

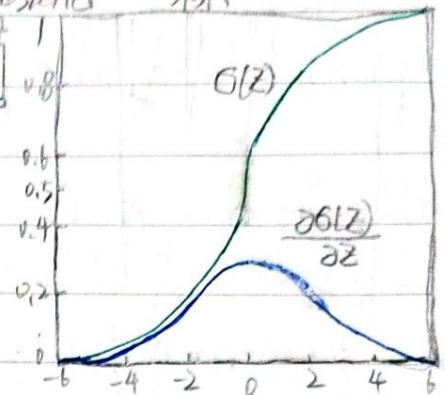
$$\Rightarrow \frac{\partial \ln (1 - f_w b(x))}{\partial w_i} = \frac{\partial \ln (1 - f_w b(x))}{\partial z} \frac{\partial z}{\partial w_i} = x_i^n \cdot \frac{1}{1 - G(z)} G(z)(1 - G(z))$$

$$\frac{\partial \ln (1 - G(z))}{\partial z} = \frac{-1}{1 - G(z)} \frac{\partial G(z)}{\partial z} = \frac{-1}{1 - G(z)} G(z)(1 - G(z))$$

$$\Rightarrow \sum_n - [\hat{y}^n (1 - f_w b(x^n)) x_i^n + (1 - \hat{y}^n) f_w b(x^n) x_i^n] = \sum_n - [\hat{y}^n - \hat{y}^n f_w b(x^n) - f_w b(x^n) + \hat{y}^n f_w b(x^n)] x_i^n$$

$\hat{y}^n - f_w b(x^n)$ Larger difference, larger update

∴ Step 3 → logistic 跟 linear 一樣，重點在 Step 2



Logistic Regression + Square Error

$$\text{Step 1: } f_w b(x) = G(\sum w_i x_i + b)$$

Step 2: Training data: (X^n, \hat{y}^n) , $\hat{y}^n = 1$ for class 1, 0 for class 2

$$L(f) = \frac{1}{2} \sum_n (f_w b(x^n) - \hat{y}^n)^2$$

$$\text{Step 3: } \frac{\partial (f_w b(x) - \hat{y})^2}{\partial w_i} = 2(f_w b(x) - \hat{y}) \frac{\partial f_w b(x)}{\partial z} \frac{\partial z}{\partial w_i}$$

$$\hat{y} = 0 \quad \Rightarrow 2(f_w b(x) - \hat{y}) f_w b(x) (1 - f_w b(x)) x_i$$

if $f_w b(x) = 1$ (far from target) $\hat{y} = 0$ (close to target) $\Rightarrow \frac{\partial L}{\partial w_i} = 0$

∴ 可能剛好在離目標很遠的地方卡住 → Cross Entropy 無此問題

Generative vs. Discriminative \Rightarrow 壓 Data 大量，需直接用 Data 判斷

\rightarrow Discriminative 精準度更高 but 若 Data 不夠，有範本還未被列舉出

v 训练 Data 中有 noise v 训练 Data 有多種形式

\rightarrow Generative 可能精準度更高 (二階補功能 XD)

• Generative \Rightarrow 機率性的; Discriminative \Rightarrow 完全根據 Data 讀練

* Benefit of generative model

- with the assumption of probability distribution, less training data is needed.
- with the assumption of probability distribution, more robust to the noise.
- priors and class-dependent probabilities can be estimated from different sources.

Multi-class Classification

$$z_i \rightarrow e^{z_i} \rightarrow \sum_{j=1}^k e^{z_j} \rightarrow y_i = e^{z_i} / \sum_{j=1}^k e^{z_j} \rightarrow \text{softmax} \xrightarrow{\text{Cross Entropy}} \hat{g}(\text{target})$$

$$\Rightarrow \text{if } X \in \text{class1: } \hat{g} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}; \text{ class2: } \hat{g} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \dots \text{以此推類}$$

Limitation of Logistic Regression

若有 4 筆 data \Rightarrow [Class 1: (0,0), (1,0) \Rightarrow 無法分類]

用 Feature Transformation? (矩陣轉換) \rightarrow Not always easy to find a good one (不聰明 (X))

✓ Cascading logistic regression models \Rightarrow 結合多个 "

\Rightarrow 機器自己 transform (V)

1. 將每個 logistic regression model 為 neuron

2. 多个 neuron 組在一起 = neural network

3. 就是 deep learning (X)

Deep learning

— Ups and downs \Rightarrow History

1958 Perceptron (linear model)

1969 Perception has limitation

1980s Multi-layer perception

\Rightarrow Do not have significant difference from DNN today

1986 Backpropagation

\Rightarrow Usually more than 3 hidden layers is not helpful

1989 1 hidden layer's "good enough", why deep?

2006 RBM initialization

2009 GPU 加速

2011 Start to be popular in speech recognition

2012 win ILSVRC image competition

Three Steps for Deep Learning (actually \neq Three Steps for Machine Learning)

Step 1: Neural Network \Rightarrow you need to decide the network structure to let a good function in your function set

Neuron = Logistic Regression \Rightarrow Neuron 前後 content of neuron network

\Rightarrow Different connection leads to different network

\rightarrow Network parameter θ = all weights and bias in the "neurons"

Ex: fully connected feedforward network

\Rightarrow a function, input vector, output vector

not fully connected

\Rightarrow Given network structure, a defined function set

special structure

= input layer + "hidden layers" + output layers \therefore Deep = Many hidden layers

home (year) \rightarrow AlexNet (2012) \rightarrow VGG (2014) \rightarrow GoogleNet (2014) \rightarrow Residual Net

layers - error rate 8 - 16.4% 19 - 7.3% 22 - 6.7% 2015) 152 - 3.5%

\Rightarrow 建算其實 = 一連串 matrix operation n 在 hidden layers 最後通过 softmax

Feature extractor $\xrightarrow{\text{replacing}}$ Feature engineering n output layer = Multi-class classifier

Step 2: Goodness of a function \Rightarrow Gradient descent

Step 3: Find the best function

(given enough hidden neurons)

but, any continuous function f can be realized by a network with one layer.

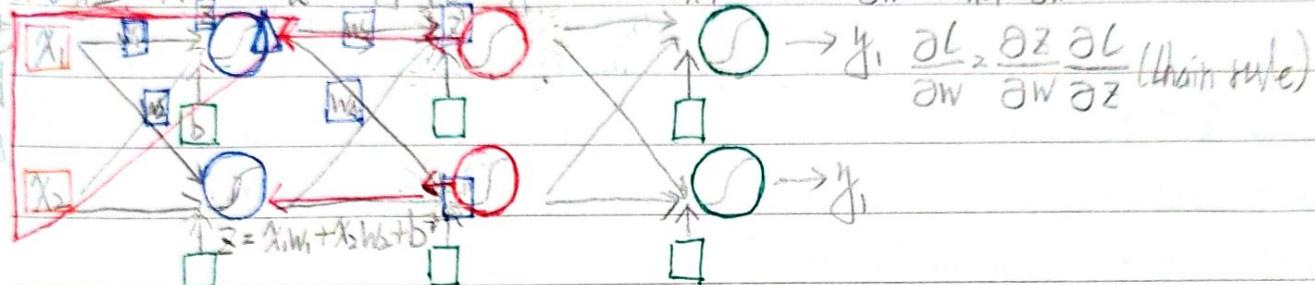
→ 提問 why "Deep" not "Flat" neural network?

Backpropagation: to compute the gradients efficiently (Same as Gradient Descent)

⇒ Chain Rule: $x: y = g(x), z = h(y) \rightarrow \Delta x \rightarrow \Delta y \rightarrow \Delta z \Rightarrow \frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$

$x = g(s), y = h(x), z = k(x, y): \Delta s \rightarrow \Delta y \rightarrow \Delta z \Rightarrow \frac{\partial z}{\partial s} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial s}$

流程: $x^n \rightarrow \text{NN} \rightarrow y^{n+1} \rightarrow \dots \rightarrow L(n) = \sum_{i=1}^k l_i(n) \Rightarrow \frac{\partial L(n)}{\partial w} = \frac{\partial L(n)}{\partial w}$



⇒ Forward pass: Compute $\frac{\partial z}{\partial w}$ for all parameters

Backward pass: Compute $\frac{\partial L}{\partial z}$ for all activation function inputs z

Forward pass: $\frac{\partial z}{\partial w_1} = x_1 \wedge \frac{\partial z}{\partial w_2} = x_2 \Rightarrow$ The value of the input connected by the weight
= 所有 weight 之偏微分, 即為其前面連接那一值

要計算每個 weight 之偏微分, 只需 inputs 然後計算每個 neuron 之 output

Backward pass: $\frac{\partial L}{\partial z} = \frac{\partial a}{\partial z} \frac{\partial L}{\partial a} = \frac{\partial z'}{\partial a} \frac{\partial L}{\partial z'} + \frac{\partial z''}{\partial a} \frac{\partial L}{\partial z''} \dots$ (Chain rule)

$$\therefore \frac{\partial L}{\partial z} = \delta'(z) [w_3 \frac{\partial L}{\partial z'} + w_4 \frac{\partial L}{\partial z''}]$$

⇒ 可想為 neuron 之 input = $\frac{\partial L}{\partial z'}$; $\frac{\partial L}{\partial z''}$, 分別乘上 w_3 ; w_4 後相加, 再乘上 $\delta'(z)$, 即等於 $\frac{\partial L}{\partial z}$

⇒ $\delta'(z)$ is a constant because z is already determined in the forward pass.

Label 1: Output layer ⇒ output = network 的 計算 $\frac{\partial L}{\partial z} = \frac{\partial l_1}{\partial z} \frac{\partial L}{\partial l_1}$ (Chain rule)

⇒ 只需知道 activation function ⇒ y_i 對 L 影響 → loss function 定義方式決定

$\frac{\partial L}{\partial z} \Rightarrow$ 同理

Case 2: Not output layer: Compute $\frac{\partial L}{\partial z}$ recursively → until output layer

= 由 output layer 反向計算至欲求層

⇒ Compute $\frac{\partial L}{\partial z}$ for all activation function inputs z * from the output layer