

目錄

深度學習

Part I: 基礎

Ch1 Introduction

Ch2 Linear Algebra *

Ch3 Probability & Information Theory *

Ch4 Numerical Optimization *

Part II: 理論

Ch5 Learning Theory & Regularization

Ch6 Probabilistic Models

Ch7 Non-parametric methods & SVMs *

Ch8 Cross-validation & Ensemble

Ch9 Large Scale Machine Learning

Part III: 神經網絡

Ch10 Neural Networks: Design

Ch11 Neural Networks: Optimization & Regularization

Ch12 Convolutional Neural Networks

Ch13 Recurrent Neural Networks

Ch14 Unsupervised Learning

Part IV: 強化學習與遷移學習

Ch15 Supervised Learning & Transfer Learning *

Ch16 Reinforcement Learning

Ch17 Deep Reinforcement Learning *

* 自研閱讀建議與題解

目錄:

Part I: 数学基礎

Ch1 Introduction to ML & DL

Ch2 Linear Algebra*

Ch3 Probability & Information Theory*

Ch4 Numerical Optimization*

Part II: 通論

Ch5: Learning Theory & Regularization

Ch6: Probabilistic Models

Ch7: Non-parametric methods & SVMs*

Ch8: Cross-validation & Ensembling

Ch9: Large Scale Machine Learning

Part III: 神經網路

Ch10 Neural Networks: Design

Ch11: Neural Networks: Optimization & Regularization

Ch12: Convolutional Neural Networks

Ch13: Recurrent Neural Networks

Ch14 Unsupervised Learning

Part IV: 強化學習 & 半監督學習

Ch15. Semisupervised Learning & Transfer Learning*

Ch16. Reinforcement Learning

Ch17. Deep Reinforcement Learning*

*: 自行閱讀講義 & 觀看 OCW 影片

Ch1 Introduction to ML & DL

- What is ML?

1. Data 的表示方法:

一个 dataset 可表為 $X = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$, 其中

$x^{(i)} \in \mathbb{R}^D, y^{(i)} \in \mathbb{R}^K$

$\begin{cases} x^{(i)}: \text{稱為 data} \\ y^{(i)}: \text{稱為 label} \end{cases}$

也有「無 label 之 dataset」 $X = \{x^{(i)}\}_{i=1}^N$

ex. MNIST dataset (手寫數字辨識)

$$\begin{cases} x^{(1)} = \boxed{6} \in \mathbb{R}^{32 \times 32} \\ y^{(1)} = e^{(6)} = (0, 0, 0, 0, 0, 1, 0, 0, 0, 0) \end{cases}$$

$$\vdots$$

2. Learning 的種類:

(1) Supervised: dataset 有 label \Rightarrow 切割 dataset, 一部分用來 learn, 一部分拿來 predict

(2) Unsupervised \rightarrow 只有 $\{x^{(i)}\}_{i=1}^N$, learn X 之 pattern

(3) Reinforcement \rightarrow learn from "good/bad" feedback (來自外界/人類/其他模型)

3. ML 的步驟

(1) collect data, 清資料, pre-processing, data exploration (最麻煩的一步就是(1))

\Rightarrow 得 $X = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$ or $\{x^{(i)}\}_{i=1}^N$

(2) 選擇一個由 w parametrized 之 model $\{f(\cdot; w)\}$ 與一個 cost function $C(w; X)$

⇒ $C(w; X)$ 之物理意義為「 f 可以對好的解釋 training data」

(3) Training ⇒ 找一組 w^* s.t.

$$w^* = \operatorname{argmin}_w C(w; X)$$

(4) Testing ⇒ 評估「經 training 後的 f 」 $\hat{=} f^*$ 表現多好

(5) 拿 f^* 去解決實際生活的問題

* 狹義地說，「套模仔」指的是只會 (5) 的人。
廣義來說，在 (2) 的時候若是用別人開發的 f ，也是「套模仔」。但實務表明，其實好的訓練 data 和好的訓練方法才是最重要的，因此我們更應關注 data 之品質和新的 training 方法

二. What is Deep Learning

Def. A deep learning model is a **multi-layer** (deep layer) ML model $f(\cdot; w)$ represent by

$$\hat{y} = f^{(L)}(\dots f^{(2)}(f^{(1)}(x; w^{(1)}); w^{(2)}) \dots; w^{(L)})$$

or by the block diagram

$$x \rightarrow [f^{(1)}(\cdot; w^{(1)})] \rightarrow [f^{(2)}(\cdot; w^{(2)})] \rightarrow \dots \rightarrow [f^{(L)}(\cdot; w^{(L)})] \rightarrow \hat{y}$$

1. 优点：
 < 学会自己 pre-processing data
 可学会較複雜之 pattern ex. visual objects

2. 缺点：
 需大量 data \rightarrow 成本高!
 computation cost \rightarrow **time** 成本高!

買 GPU

可以跟教授說
「我的 model 还在跑」

Ch5 Learning Theory & Regularization

DATE

- Learning Theory

Def. Learning (Mitchell, 1997)

一了程式在執行任務 T 時, 若其表現(以 P 測量)隨經驗 E 改善, 則稱此程式會從 E 中學習

$T \rightarrow$ supervised, unsupervised, RL 等項目

$P \rightarrow$ accuracy, error rate, MSE, ...

$E \rightarrow$ dataset (X, y) or (X)

給定 training set of size N , f 為 data 中可能的一種對應關係, 定義 empirical error/risk on N sample 為:

$$C_N(w) \text{ or } C_N[f] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f(x^{(i)}, w), y^{(i)})$$

令 $P(x, y)$ 為 data 實際之 pdf, 定義 generalization error/risk 為 (可由 test error 判斷!)

$$C(w) \text{ or } C[f] = \int \text{loss}(f(x; w), y) dP(x, y)$$

$\Rightarrow C$ or C_N 本質是 loss function 之期望值, 物理意義是 f 有多麼 generalize

令 \mathcal{F} 為一泛函空間(可以是我們挑的), 我們想找 $f^* = \arg \min_{f \in \mathcal{F}} C[f]$, 不過顯然我們只能藉學習方法找出

$$f_N = \arg \min_{f \in \mathcal{F}} C_N[f]$$

Learning Theory 最主要的問題之一就是: "How to characterize $C[f_N] = \int \text{loss}(f_N(x; w), y) dP(x, y)$?"
即, learn 到的 function f_N 有多麼 generalize?

Q: 若函數 f 之 $C_N[f]$ 值很小, 則 $C[f]$ 亦很小嗎?

A: No!

Q: ML 一定找得到 f^* 嗎?

A: NO!

Thm. No Free Lunch Thm

Averaged over all possible data generating distribution, every classification algorithm has the same error rate when classifying unseen points.

→ 沒有一種演算法比其他演算法更能在任意的 task 表現更好
 → 我們更在乎「對生成 data 之 pmf $P(x, y)$ 而言, 哪一個演算法生成了更好的函數 f_N 」?

至於如何討論 $C[f_N]$, 有 Bounding Method 與 Decomposition method 兩種。

二. Bounding Method

Def. Bayes error: $\min_f C[f] \triangleq C[f^*]$

⇒ $C[f^*] > 0$ if there is randomness in $P(y|x)$

ex. 考慮 $y = \sin x + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. 取 $\mathcal{F} = \mathcal{P}_{10}$ (10 次及以下多項式形成之空間)

⇒ 如何找 $f_N \in \mathcal{F}$ st. $C[f_N]$ 夠靠近 $C[f^*]$?

Def. Bounding method: 定義

$$\mathcal{E}_0 = C[f_N] - C[f^*] = \underbrace{C[f_{\#}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\#}^*]}_{\mathcal{E}_{\text{est}}}$$

其中 $f_{\#}^*$ 為 \mathcal{F} 中最好的 model.

② \mathcal{E}_{app} : approximation error, 物理意義為「model 中最好的 function 跟實際上真正的 function 之差」

③ \mathcal{E}_{est} : estimation error: 「從 data 中 learn 到的 function 和 model 中最好的 function 之差」

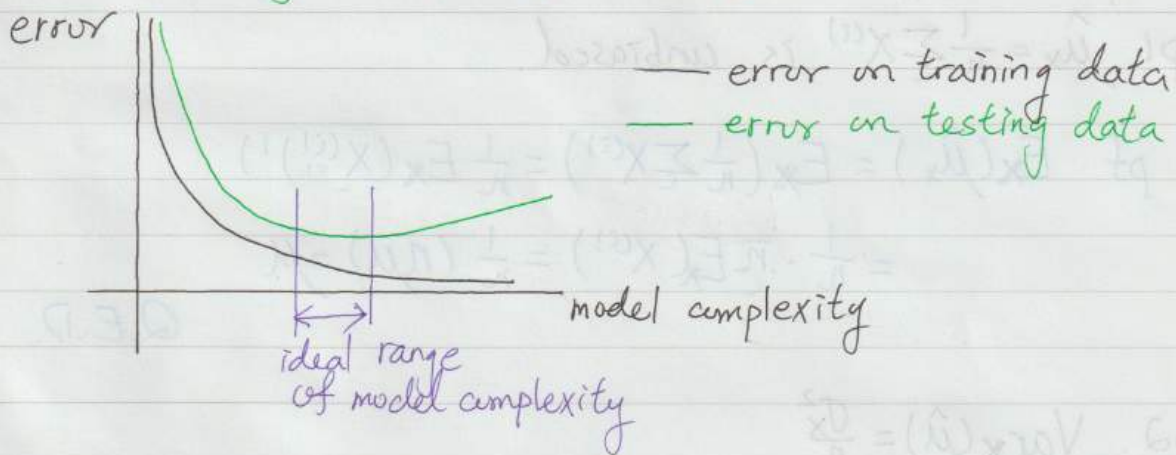
Q: How to reduce E_{app} ?

A: 顯然可以藉由挑 - 了更 **複雜的 F** 減少之

Q: How to reduce E_{est} ?

A: ① **Simpler** model (Complexity of $F \downarrow$)
 ② Larger training set.

⇒ 當 model 太簡單 ⇒ high E_{app} ⇒ underfitting
 { high training error; high testing error
 " 複雜 ⇒ high E_{est} ⇒ overfitting
 low training error; high testing error.



三. Decomposition Methods:

1. 前言: No (weak) assumption on data distribution is made ⇒ these bounds are too loose to quantify $E[f_n]$.

⇒ decompose $E[f_n]$ into multiple meaningful terms

假設: $\left\{ \begin{array}{l} \text{特定之 loss function} \\ \text{data generating distribution } P(x, y) \end{array} \right.$
 特定

(能透過數學運算完成 decomposition)

⇒ 需要統計學中 point estimation 的知識

2. Point estimation.

Def. 設 dataset $X = \{X^{(1)}, \dots, X^{(n)}\}$ 為 n 了 iid 樣本 of r.v. X . $X \sim$ 以 θ 為母數之分布, 則: a point estimator or statistic is $\hat{\theta}_n \triangleq g(X^{(1)}, \dots, X^{(n)})$ (function of the data). $\hat{\theta}_n$ is called the estimate of θ .

Def. The bias of an estimator $\text{bias}(\hat{\theta}_n) = E_X(\hat{\theta}_n) - \theta$
 The variance " $\text{Var}_X(\hat{\theta}_n) = E_X[(\hat{\theta}_n - E_X[\hat{\theta}_n])^2]$
 We say the estimator is unbiased if $\theta = E_X(\hat{\theta}_n)$

Prop 1. $\hat{\mu}_X = \frac{1}{n} \sum_{i=1}^n X^{(i)}$ is unbiased.

$$\begin{aligned} \text{pf. } E_X(\hat{\mu}_X) &= E_X\left(\frac{1}{n} \sum_{i=1}^n X^{(i)}\right) = \frac{1}{n} E_X\left(\sum_{i=1}^n X^{(i)}\right) \\ &= \frac{1}{n} \cdot \sum_{i=1}^n E_X(X^{(i)}) = \frac{1}{n} (n\mu) = \mu \end{aligned}$$

Q.E.D.

Prop 2. $\text{Var}_X(\hat{\mu}) = \frac{\sigma_X^2}{n}$

$$\begin{aligned} \text{pf. } \text{Var}(\hat{\mu}) &= E_X[(\hat{\mu} - E_X[\hat{\mu}])^2] \\ &= E[\hat{\mu}^2 - 2\hat{\mu}\mu + \mu^2] \\ &= E[\hat{\mu}^2] - \mu^2 \\ &= E\left[\left(E\left[\frac{1}{n} \sum_{i=1}^n X^{(i)}\right]\right)^2\right] - \mu^2 \\ &= E\left[\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n X^{(i)} X^{(j)}\right] - \mu^2 \\ &= \frac{1}{n^2} \left[\sum_{i=j} E[X^{(i)} X^{(j)}] + \sum_{i \neq j} E[X^{(i)} X^{(j)}] \right] - \mu^2 \\ &= \frac{1}{n^2} \left[\sum_{i=1}^n E[X^{(i)^2}] + n(n-1) E[X^{(i)}] E[X^{(j)}] \right] - \mu^2 \\ &= \frac{1}{n} E[X^2] + \frac{(n-1)}{n} \mu^2 - \mu^2 \end{aligned}$$

$$= \frac{1}{n} (E[X^2] - \mu^2)$$

$$= \frac{1}{n} \sigma_x^2$$

Q.E.D.

Prop 3. $\hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X^{(i)} - \hat{\mu}_x)^2$ is the unbiased estimator

pf. $E_X[\hat{\sigma}^2] = E\left[\frac{1}{n} \sum_{i=1}^n (X^{(i)} - \hat{\mu})^2\right]$

$$= E\left[\frac{1}{n} \left(\sum_{i=1}^n X^{(i)2} - 2 \sum_{i=1}^n X^{(i)} \hat{\mu} + \sum_{i=1}^n \hat{\mu}^2\right)\right]$$

$$= E\left[\frac{1}{n} \left(\sum_{i=1}^n X^{(i)2} - n \hat{\mu}^2\right)\right]$$

$$= E[X^2] - E[\hat{\mu}^2]$$

$$= E[(X-\mu)^2 + 2X\mu - \mu^2] - E[\hat{\mu}^2]$$

$$= (\sigma^2 + \mu^2) - (\text{Var}[\hat{\mu}] + E[\hat{\mu}]^2)$$

$$= \sigma^2 + \mu^2 - \frac{1}{n} \sigma^2 - \mu^2$$

$$= \frac{n-1}{n} \sigma^2$$

$\Rightarrow \frac{1}{n-1} \sum_{i=1}^n (X^{(i)} - \hat{\mu}_x)^2$ is an unbiased estimator of σ_x .
Q.E.D.

Prop 4. Mean square error $E[(\hat{\theta}_n - \theta)^2] = \text{Var}_X(\hat{\theta}_n) + \text{bias}(\hat{\theta}_n)^2$

pf. $E_X[(\hat{\theta}_n - \theta)^2] = E[(\hat{\theta}_n - E[\hat{\theta}_n] + E[\hat{\theta}_n] - \theta)^2]$

$$= E[(\hat{\theta}_n - E[\hat{\theta}_n])^2 + (E[\hat{\theta}_n] - \theta)^2 + 2(\hat{\theta}_n - E[\hat{\theta}_n])(E[\hat{\theta}_n] - \theta)]$$

$$= E[(\hat{\theta}_n - E[\hat{\theta}_n])^2] + E[(E[\hat{\theta}_n] - \theta)^2]$$

$$+ 2E[(\hat{\theta}_n - E[\hat{\theta}_n])(E[\hat{\theta}_n] - \theta)]$$

$$= E[(\hat{\theta}_n - E[\hat{\theta}_n])^2] + (E[\hat{\theta}_n] - \theta)^2 + 2 \cdot 0 \cdot (E[\hat{\theta}_n] - \theta)$$

$$= \text{Var}_X(\hat{\theta}_n) + \text{bias}(\hat{\theta}_n)^2$$

Q.E.D.

3. Consistency

Def. An estimator is **(weak consistent)** \Leftrightarrow

$\lim_{n \rightarrow \infty} \hat{\theta}_n$ converges in probability to θ , or $\forall \varepsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| \geq \varepsilon) = 0 \quad (\text{or } \lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| < \varepsilon) = 1)$$

Thm. Weak law of Large number

The sample mean $\hat{\mu}_x = \frac{1}{n} \sum_{i=1}^n X^{(i)}$ is a consistent estimator of μ_x . i.e.

$$\lim_{n \rightarrow \infty} P(|\hat{\mu}_{x,n} - \mu_x| < \varepsilon) = 1$$

Thm. Strong law of large number

$P\left(\lim_{n \rightarrow \infty} \hat{\mu}_{x,n} = \mu_x\right) = 1$. i.e. $\hat{\mu}_x$ is a strong consistent estimator.

4. Decomposing generalization error

現考慮 $f_N = \arg \min_{f \in \mathcal{F}} C_N[f]$, 我們可以分解 $C[f_N]$ 之期望值,

取 loss function 為 MSE, true label $y = f^*(x) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. 首先, 由

$$\begin{aligned} E_{\mathcal{X}}[C[f_N]] &= E_{\mathcal{X}}\left[\int \text{loss}(f_N(x) - y) dP(x, y)\right] \\ &= E_{\mathcal{X}, x, y}[\text{loss}(f_N(x) - y)] \\ &= E_{\mathcal{X}}\left[E_{\mathcal{X}, y}[\text{loss}(f_N(x) - y) \mid X = x]\right] \quad (*) \end{aligned}$$

代入 MSE 與 $\mu = 0$, 可得以下結果:

Thm. 在以上之假設中, $E_{x,y}[\text{loss}(f_N(x)-y)|X=x] = \sigma^2 + \text{Var}_x[f_N(x)|x] + \text{bias}[f_N(x)|x]^2$

$$\begin{aligned}
 \text{pf. } E_{x,y}[\text{loss}(f_N(x)-y)|x] &= E_{x,y}[(f_N(x)-y)^2|x] \\
 &= E_{x,y}[y^2 + f_N(x)^2 - 2f_N(x)y|x] \\
 &= E_y[y^2|x] + E_x[f_N(x)^2|x] - 2E_{x,y}[f_N(x)y|x] \\
 &= (\text{Var}_y[y|x] + E_y[y|x]^2) + (\text{Var}_x[f_N(x)|x] + E_x[f_N(x)|x]^2) \\
 &\quad - 2E_y[y|x]E_x[f_N(x)|x] \\
 &= \text{Var}_y[y|x] + \text{Var}_x[f_N(x)|x] + (E_x[f_N(x)|x] - E_y[y|x])^2 \\
 &= \text{Var}_y[y|x] + \text{Var}_x[f_N(x)|x] + E_x[f_N(x) - f^*(x)|x]^2 \\
 &= \sigma^2 + \text{Var}_x[f_N(x)|x] + \text{bias}[f_N(x)|x]^2
 \end{aligned}$$

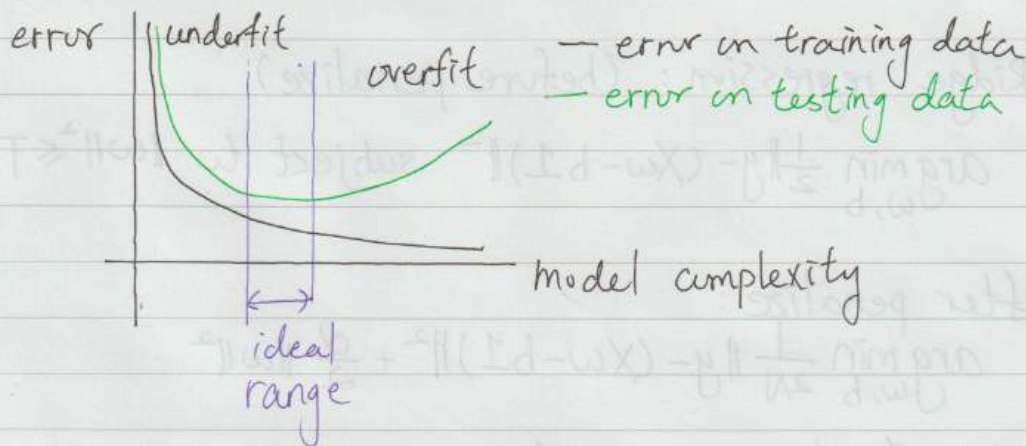
Q.E.D.

$$\begin{aligned}
 \text{Cor. } E_x[C[f_N]] &= E_x(E_{x,y}[\text{loss}(f_N(x)-y)|x]) \\
 &= E_x(\sigma^2 + \text{Var}_x[f_N(x)|x] + \text{bias}[f_N(x)|x]^2)
 \end{aligned}$$

(1) $\sigma^2 \rightarrow$ 無法 avoid !! \because 源自 data 自身 randomness

(2) trade off between variance & bias! (圖例: PPT p27)

* 結論: 到目前為止, 我們討論了無論要用哪種演算法 train model, 都要想好此 model f_N 會有哪些 trade off 以及如何挑出最「適宜」的 model



四. Regularization

目前我們都是由 generalization error

$$E[f_N] = \int \text{loss}(f_N(x; w), y) dP(x, y)$$

探討 f_N 有多 generalize. 不過實際上我們只能由最小化 empirical error

$$C_N[f] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f(x^{(i)}; w), y_i)$$

來 train 出 model f_N . 即 $f_N = \underset{f \in \mathcal{F}}{\text{argmin}} C_N[f]$. 我們稱任何改善 f_N generalizability 的手段為 regularization

on cost function: **weight decay**
during training process: **validation**

Occam's razor: among equal-performing models, the simplest one should be selected.

1. Weight decay:

我們以 f_N 之參數 w 之範數定義 f_N 之複雜程度
這是由於 $\|w\| = 0 \leftrightarrow w = 0$ 為隨機猜答案

\Rightarrow **penalize on norm of w** , 即在 cost function 後加上 $\|w\|$.

ex. Ridge regression: (before penalize)

$$\underset{w, b}{\text{argmin}} \frac{1}{2} \|y - (Xw - b\mathbf{1})\|^2 \text{ subject to } \|w\|^2 \leq T.$$

after penalize:

$$\underset{w, b}{\text{argmin}} \frac{1}{2N} \|y - (Xw - b\mathbf{1})\|^2 + \frac{\alpha}{2} \|w\|^2$$

$\alpha > 0$ 為 hyper parameter.

$\left\{ \begin{array}{l} \text{大 } \alpha: \text{ prefer simple model} \\ \text{小 } \alpha: \text{ care empirical error (低-真)} \end{array} \right.$

不去对 b regularize 之原因為 y 未必有標準化过。

Def. LASSO (least absolute shrinkage and selection operator) \rightarrow penalize $\|w\|_1$,

$$\arg \min_{w, b} \frac{1}{2N} \|y - (Xw - b1)\|^2 + \alpha \|w\|_1$$

\rightarrow sparse w ! 原因: $\|w\|_1$ 之頂尖落在坐標軸上, $\| \cdot \|^2$ 之等位面為高維球面 \rightarrow 兩者交點易在坐標軸上。

Def. Elastic Net: combines Ridge and LASSO

$$\arg \min_{w, b} \frac{1}{2N} \|y - (Xw - b1)\|^2 + \alpha (\beta \|w\|_1 + \frac{1-\beta}{2} \|w\|^2)$$

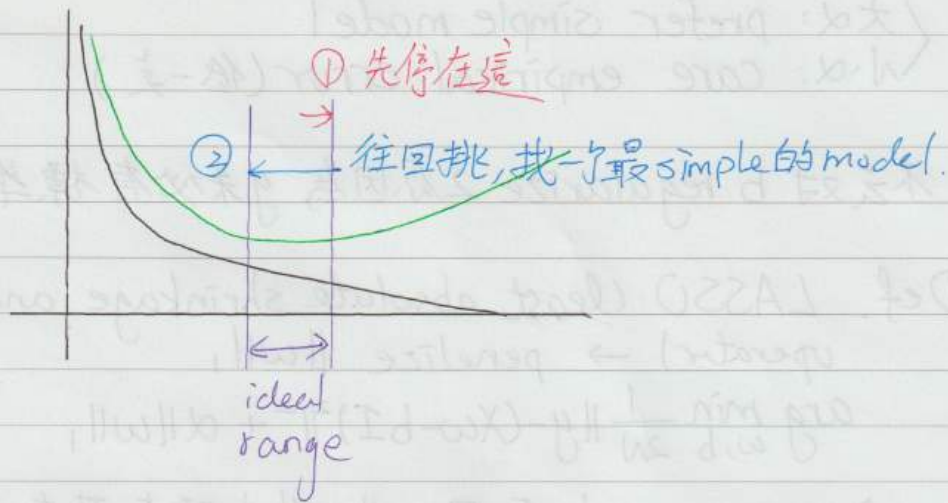
2. Validation

Def. The constants fixed in a model is called **hyperparameters**

ex. degree p in polynomial regression, α in LASSO...

物理意義 \Rightarrow 改變 hyperparameter 之值如同改變 model complexity!

\Rightarrow 根據 Occam's razor, 我們應逐漸增加 model complexity, 直至 overfit 發生前一刻停止,



此外, 「衡量是否 overfit」以及「往回挑」時, 不應以 test data 測試之, 否則只是讓 model 對 testing set overfit 而已.

⇒ 要從 training set 切出一個 validation set 來選 hyperparameter!

Ch6 Probabilistic Models

NO.

DATE

1. Probabilistic Models

Given a dataset $X = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$, Model F :
a collection of functions parametrized by Θ

Goal: new data point $x' \rightarrow$ train a function f s.t.
 $\hat{y} = f(x'; \Theta)$ is closet to the correct label y'

Def. Probabilistic models:

We write the above function $f(x'; \Theta)$ as
 $P(y=y' | X=x')$. The prediction is made by
 $\hat{y} = \underset{y}{\operatorname{argmax}} P(y=y' | X=x'; \Theta)$

In such model, we have to find Θ .

我們有兩種方法找 Θ :

<法 I>: Maximum likelihood Estimation (MLE)

設 $X \sim$ 以 Θ 為母數之分佈, 則求

$$\Theta_{ML} = \underset{\Theta}{\operatorname{argmax}} P(X|\Theta)$$

<法 II>: Maximum Posteriori Estimation

$$\Theta_{MAP} = \underset{\Theta}{\operatorname{argmax}} P(\Theta|X)$$

$$= \underset{\Theta}{\operatorname{argmax}} \frac{P(X|\Theta)P(\Theta)}{P(X)}$$

$$= \underset{\Theta}{\operatorname{argmax}} P(X|\Theta)P(\Theta)$$

得病 Θ 出
現症狀 X 之
機率?

pmf (or pdf)

ex. Θ 為病, $P(\Theta)$ 表所有病
分佈之 pmf

2. Maximum Likelihood Estimation

Assumption (ML 起手式):

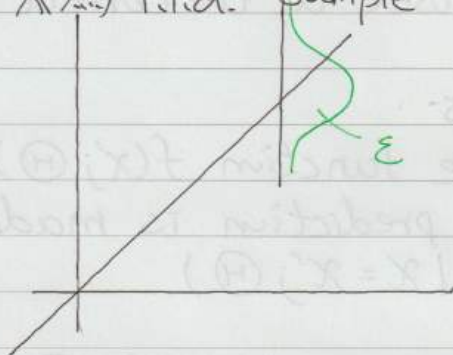
(1) \exists a deterministic f^* s.t. $y = f^*(x) + \varepsilon$

(2) $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$

(3) 做線性迴歸, 假設 $f^*(x; w^*) = w^{*T}x$

(4) X 已標準化

(5) X 為 i.i.d. sample



\Rightarrow 目標: 找出一個 w 使得 $\hat{y} = \underset{y}{\operatorname{argmax}} P(y | X=x; w) = w^T x$

方法: MLE of w^* : $w_{ML} = \underset{w}{\operatorname{argmax}} P(X|w)$

仿照統計學中的做法,

$$P(X|w) = \prod_{i=1}^N P(x^{(i)}, y^{(i)} | w)$$

$$= \prod_{i=1}^N P(y^{(i)} | x^{(i)}, w) P(x^{(i)} | w)$$

$$= \prod_{i=1}^N P(y^{(i)} | x^{(i)}, w) P(x^{(i)})$$

$$= \prod_{i=1}^N g(y^{(i)}; \mu = w^T x^{(i)}, \sigma^2) P(x^{(i)}) \quad (g \text{ 為 } \mathcal{N}(w^T x^{(i)}, \sigma^2) \text{ 的 pdf})$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2} (y^{(i)} - w^T x^{(i)})^2\right) P(x^{(i)})$$

$$\Rightarrow \operatorname{argmax}_w P(X|w) = \operatorname{argmax}_w \log P(X|w)$$

$$= \operatorname{argmax}_w \log \left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2} (y^{(i)} - w^T x^{(i)})^2\right) P(x^{(i)}) \right]$$

$$= \arg \max_w N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - w^T x^{(i)})^2 + \sum_i P(x^{(i)})$$

$$= \arg \min_w \sum_i (y^{(i)} - w^T x^{(i)})^2$$

\Rightarrow solve analytically or by SGD.

三. MLE 之应用: Logistic regression.

Assumption: $(y | x) \sim b(1, p)$, 即

$$P(y | x; p) = p^y (1-p)^{1-y}, \quad y' = \frac{y+1}{2} \in \{0, 1\}$$

(由於是要做 binary classification, 因此 y 值只有 0, 1 兩種可能)

目標: 定義 logistic function 為 $\sigma(z) = \frac{1}{1+e^{-z}}$, 且令

$z = w^T x$, 想找出一 w s.t.

$$\hat{y} = \arg \max_y P(y | x; w) = \text{sign}(w^T x)$$

$$\begin{aligned} \text{其中 } P(y | x; w) &= \sigma(z)^y (1-\sigma(z))^{1-y} \\ &= \sigma(w^T x)^y (1-\sigma(w^T x))^{1-y} \end{aligned}$$

方法: MLE of w : $w_{ML} = \arg \max_w P(X | w)$

$$\Rightarrow \arg \max_w P(X | w) = \arg \max_w \log P(X | w)$$

$$= \arg \max_w \log \prod_i P(y^{(i)} | x^{(i)}; w) P(x^{(i)} | w)$$

$$= \arg \max_w \log \prod_i \sigma(w^T x^{(i)})^{y^{(i)}} (1-\sigma(w^T x^{(i)}))^{1-y^{(i)}}$$

$$= \arg \max_w \sum_i y^{(i)} \log \sigma(w^T x^{(i)}) + (1-y^{(i)}) \log (1-\sigma(w^T x^{(i)}))$$

$$= \arg \max_w \sum_i y^{(i)} \log \frac{1}{1+e^{-w^T x^{(i)}}} + (1-y^{(i)}) \log \frac{e^{w^T x^{(i)}}}{1+e^{-w^T x^{(i)}}}$$

$$= \arg \max_w \sum_c y^{(c)} \log(1 + e^{-w^T x^{(c)}}) + (1 - y^{(c)}) (w^T x^{(c)} - \log(1 + e^{w^T x^{(c)}}))$$

$$= \arg \max_w \sum_c \underbrace{y^{(c)} w^T x^{(c)} - \log(1 + e^{w^T x^{(c)}})}_{(*)}$$

\therefore We can prove that $(*)$ is concave of w

\therefore SGD finds global optimal.

IV. Maximum A Posteriori Estimation (MAP)

$$\text{We solve } w_{\text{MAP}} = \arg \max_w P(w|X) = \arg \max_w P(X|w) \underbrace{P(w)}_{\text{pdf}}$$

\rightarrow 給某些 w 权重

Let's apply MAP on linear regression. We assume

1. $w \sim N(0, \beta^{-1} I)$

$$\Rightarrow \arg \max_w P(w|X) = \arg \max_w \log P(w|X)$$

$$= \arg \max_w \log P(X|w) P(w)$$

by MLE

$$= \arg \max_w \log P(X|w) + \log P(w)$$

$$= \arg \max_w \sum_c (y^{(c)} - w^T x^{(c)})^2 + \log \frac{1}{(2\pi)^D \det(\beta^{-1} I)} \exp\left(-\frac{1}{2} w^T (\beta^{-1} I)^T w\right)$$

$$= \arg \max_w \sum_c (y^{(c)} - w^T x^{(c)})^2 - \beta w^T w$$

$$= \arg \max_w \sum_c (y^{(c)} - w^T x^{(c)})^2 - \beta \|w\|^2$$

\Rightarrow Ridge regression $\hat{=}$ weight decay term!

2. $w \sim \text{Laplace}(0, b)$

note that Laplace distribution 的 pdf 为

$$g(w|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$$

$$\begin{aligned} \Rightarrow \operatorname{argmax}_{\omega} P(\omega | X) &= \operatorname{argmax}_{\omega} \log(X | \omega) + \log P(\omega) \\ &= \operatorname{argmax}_{\omega} \sum_c (y^{(c)} - \omega^T x^{(c)})^2 + \log \frac{1}{2b} \exp\left(-\frac{|\omega|}{b}\right) \\ &= \operatorname{argmax}_{\omega} \sum_c (y^{(c)} - \omega^T x^{(c)})^2 - \frac{1}{b} |\omega| \end{aligned}$$

\Rightarrow LASSO!

* some remarks

Thm. (Consistency)

The ML estimator \hat{H}_{ML} is consistent.

Thm. (Cramér-Raw Lower bound)

At a fixed large number N of examples, no consistent estimator of \hat{H} has a lower expected MSE than the ML estimator \hat{H}_{ML} .

Remark. data 沒有很多時 要用 MAP \Rightarrow introduce bias, reduce variance

Ch8 Cross Validation & Ensembling

— Cross Validation

Notation: X : training set

$X^{(i)}$: " 之第 i 个 fold

N : size

$N^{(i)}$: 之第 i 个 fold 之 size

C_{cv} : cross validation error

f : model

1. CV 之用途
 - Hyperparameter tuning
 - Performance reporting

Algorithm. K -fold CV

1. 将 X 分成 K 个等份之 fold $X^{(i)}$, $i=1, 2, \dots, K$
2. For $i=1, \dots, K$, 以 $X-X^{(i)}$ train 出 $f_{-N^{(i)}}$
3. 令 $C_{cv} = \frac{1}{K} \sum_{i=1}^K C[f_{-N^{(i)}}]$, 其中 $C[\cdot]$ 为 test error. 在 $X^{(i)}$ 上

Algorithm. $K \times M$ nested CV.

1. 将 X 分成 K 等份之 fold $X^{(i)}$, $i=1, 2, \dots, K$.
2. For $i=1, 2, \dots, K$,
 - (1) 以 $X^{(i)}$ 为 test fold, $X-X^{(i)}$ 为 training fold, 将 $X-X^{(i)}$ 再分为 M 个等份
 - (2) 对所有 hyperparameter 的候选组合, 以 $X-X^{(i)}$ 做 M -fold CV
 - (3) 以所找到最好的 hyperparameter 和整个 $X-X^{(i)}$ train 出 $C[f_{-N^{(i)}}]$
3. 令 $C_{cv} = \frac{1}{K} \sum_{i=1}^K C[f_{-N^{(i)}}]$ 为 test error. $C[\cdot]$ 为 $X^{(i)}$ 上之 test error.

2. 設 $X \sim$ 以 θ 為母數之分布, $\hat{\theta}_n$ 為由 n 次來自 X 之 sample 估計之 θ . 則定義 mean square error of $\hat{\theta}_n$ 為: $MSE(\hat{\theta}_n) = E_x[(\hat{\theta}_n - \theta)^2]$

Prop. $E_x[(\hat{\theta}_n - \theta)^2] = \text{Var}_x(\hat{\theta}_n) + (\text{bias}(\hat{\theta}_n))^2$

pf. $E_x[(\hat{\theta}_n - \theta)^2] = E[(\hat{\theta}_n - E[\hat{\theta}_n] + E[\hat{\theta}_n] - \theta)^2]$
 $= E[(\hat{\theta}_n - E[\hat{\theta}_n])^2 + (E[\hat{\theta}_n] - \theta)^2 + 2(\hat{\theta}_n - E[\hat{\theta}_n])(E[\hat{\theta}_n] - \theta)]$
 $= E[(\hat{\theta}_n - E[\hat{\theta}_n])^2] + (E[\hat{\theta}_n] - \theta)^2$
 $= \text{Var}_x[\hat{\theta}_n] + \text{bias}(\hat{\theta}_n)^2$

Q.E.D.

Cor. $\because C_{cv}$ is an estimator of $E_x[C[f_N]]$
 $(\because C[f_N^{(s)}] \text{ as well})$

$$\therefore MSE(C_{cv}) = E_x[C_{cv} - E_x[C[f_N]]]^2$$

$$= \text{Var}_x[C_{cv}] + \text{bias}[C_{cv}]^2$$

3. 物理意義: 以 5-Fold 與 10-Fold 比較,

$$(1) \begin{cases} \text{Var}: 10\text{-Fold} > 5\text{-Fold} \\ \text{bias}: < < \end{cases}$$

(2) MSE of an unbiased estimator is its variance

二. Fold 數之選擇

Lemma. $\text{bias}[C_{cv}] = \text{bias}[C[f_N^{(s)}]]$, $s=1, 2, \dots, K$

$$\text{Var}_x[C_{cv}] = \frac{1}{K} \text{Var}[C[f_N^{(s)}]] + \frac{2}{K^2} \sum_{i < j} \text{Cov}[C[f_N^{(i)}], C[f_N^{(j)}]]$$

$$\begin{aligned}
 \text{pf. bias}[C_{cv}] &= E_x[C_{cv}] - E_x[C[f_N]] \\
 &= E\left[\sum_{i=1}^K \frac{1}{K} C[f_{-N^{(i)}}]\right] - E[C[f_N]] \\
 &= \frac{1}{K} \sum_i E[C[f_{-N^{(i)}}]] - E[C[f_N]] \\
 &= E[C[f_{-N^{(s)}}]] - E[C[f_N]], \forall s=1, \dots, K \\
 &= \text{bias}[C[f_{-N^{(s)}}]], \forall s=1, 2, \dots, K
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}_x[C_{cv}] &= \text{Var}\left[\sum_{i=1}^K \frac{1}{K} C[f_{-N^{(i)}}]\right] \\
 &= \frac{1}{K^2} \text{Var}\left[\sum_{i=1}^K C[f_{-N^{(i)}}]\right] \\
 &= \frac{1}{K^2} \left(\sum_{i=1}^K \text{Var}[C[f_{-N^{(i)}}]] + 2 \sum_{i=1}^K \sum_{j=i}^K \text{Cov}_x[C[f_{-N^{(i)}}], C[f_{-N^{(j)}}]] \right) \\
 &= \text{所求}
 \end{aligned}$$

Q.E.D.

⇒ 可觀察到 $\therefore \text{bias}[C[f_{-N^{(i)}}]]$ 與 $\text{Var}[C[f_{-N^{(i)}}]]$ 之間必有 trade off

$\therefore \text{bias}[C_{cv}]$ 與 $\text{Var}[C_{cv}]$ 亦有 trade off

但由 lemma 2, 只要使 $f_{-N^{(i)}}$ 與 $f_{-N^{(j)}}$ uncorrelated 就能使 $\text{Var}[C_{cv}] \downarrow$.

1. K 很大時:

$\therefore f_{-N^{(s)}}$ is trained in more examples

$\therefore \text{bias}[C[f_{-N^{(s)}}]], \text{Var}[C[f_{-N^{(s)}}]] \downarrow$

$\text{Cov}[C[f_{-N^{(i)}}], C[f_{-N^{(j)}}]] \uparrow$ ($\because X-X^{(i)}, X-X^{(j)}$ are similar)

2. K 很小時 \Rightarrow 恰相反

3. Leave-one-out CV (取 $K=N$)
 \Rightarrow 用於 dataset 非常小的時候
 \therefore 當 dataset 太小時, $\text{bias}[C_{L^1}]$ & $\text{Var}[C_{L^1}]$
 dominate $\text{MSE}[C_{L^1}]$
 \therefore 必 **overfit**
 \Rightarrow 取 $K=N$, 以減少 **Variance**

三. Voting

Def. Voting is a linear combination of the predictions of base-learners for each x :

$$\tilde{y}_k = \sum_{j=1}^L w_j \hat{y}_k^{(j)}, \quad w_j > 0, \quad \sum_{j=1}^L w_j = 1$$

ex. Sum rule: $\tilde{y}_k = \frac{1}{L} \sum_{j=1}^L \hat{y}_k^{(j)}$

minimum: $\tilde{y}_k = \min_j \hat{y}_k^{(j)}$

Q: Why Voting Works?

Prop. Assume that each $\hat{y}^{(j)}$ has the expected value $E_x[\hat{y}^{(j)}|x]$ and Variance $\text{Var}_x[\hat{y}^{(j)}|x]$. If $\hat{y}^{(i)}$ and $\hat{y}^{(j)}$ are uncorrelated, $\forall i \neq j$, then the variance can be reduced, but the bias is same

pf. Set $w_j = \frac{1}{L}$, then

$$\begin{aligned} \therefore E_x[\tilde{y}|x] &= E\left(\sum_j \frac{1}{L} \hat{y}^{(j)}|x\right) \\ &= \frac{1}{L} \sum_j E[\hat{y}^{(j)}|x] \\ &= E[\hat{y}^{(1)}|x] \end{aligned}$$

\therefore The expected value does not change.

⇒ The bias doesn't change

$$\therefore \text{Var}_x(\tilde{y}|x) = \text{Var}\left(\sum_j \frac{1}{L} \hat{y}^{(j)} \mid x\right)$$

$$= \frac{1}{L^2} \text{Var}\left(\sum_j \hat{y}^{(j)} \mid x\right)$$

$$= \frac{1}{L} \text{Var}(\hat{y}^{(j)} \mid x) + \frac{2}{L^2} \sum_{l \neq j} \text{Cov}(\hat{y}^{(l)}, \hat{y}^{(j)} \mid x)$$

∴ If $\hat{y}^{(i)}$ and $\hat{y}^{(j)}$ are uncorrelated, the variance can be reduced
 " " positively correlated, " increases.

Q.E.D.

Remark. 實際上 $\hat{y}^{(j)}$ 之間不大可能 i.i.d. 比如只差在 activate function

四. Bagging (bootstrap aggregating)

目標: 使 base-learners 不同

⇒ 使用「稍微不同的」data 去 train

Algorithm. Bagging.

1. 自 X 中取 N 个 data, 取後放回, 重複 L 次, 得 L 个

training set $X^{(j)}$, $j=1, \dots, L$

2. 以 $X^{(j)}$ train L 个不同的 model

五. Boosting

1. 目標: train **complementary** base-learners

∴ 要 train uncorrelated 的 base-learner

可能要用很 unstable 的方法.

2. 基本的 Boosting 演算法
 \Rightarrow train the next learner on the mistakes of the previous learners

比如: 考虑 binary classification $d^{(j)}(x) \in \{1, -1\}$

Def. A weak learner has error prob. less than $1/2$ (better than random guess)
 "strong" arbitrarily small error prob.

Algorithm. Boosting: (binary classification)

- Training:
1. 将 X 分成三份 $X^{(1)}, X^{(2)}, X^{(3)}$
 2. 用 $X^{(1)}$ train 出 $d^{(1)}$, 再拿 $d^{(1)}$ 去 predict $X^{(2)}$
 3. 以 $X^{(2)}$ 中 $d^{(1)}$ 分错的去 train $d^{(2)}$
与其
 4. 拿 $d^{(1)}, d^{(2)}$ 去 predict $X^{(3)}$
 5. 用 $X^{(3)}$ 中 $d^{(1)}, d^{(2)}$ predict 出来不同的去 train $d^{(3)}$

- Testing:
1. 以 $d^{(1)}, d^{(2)}$ 预测一个 point
 2. 若 $d^{(1)}(x) = d^{(2)}(x)$, 则以此为 prediction; 反之, 则以 $d^{(3)}(x)$ 为 prediction.

3. AdaBoost

Notation: $\Pr^{(i,j)} = P((x^{(i)}, y^{(i)}) \text{ is drawn to train the } j^{\text{th}} \text{ base-learner } d^{(j)})$
 $E^{(j)} = \text{error rate of } d^{(j)} \text{ on its training set}$
 $\triangleq \sum_i \Pr^{(i,j)} \cdot \mathbb{1}(y^{(i)} \neq d^{(j)}(x^{(i)}))$

- (1) 概念: uses the same training set over and over again

(2). Modify the prob of drawing the instances as a function of error rate $\epsilon^{(j)}$

Algorithm.

Training 1. Initialize $Pr^{(i,j)} = \frac{1}{N}$

2. For all j :

(1) 以 $Pr^{(i,j)}$ 之机率自 X 中抽出 N 个 example

(2) 以此 train $d^{(j)}$

(3) 若 $\epsilon^{(j)} \geq 1/2$, 則不再新增 base-learner

(4) 令 $\alpha_j = \frac{1}{2} \log\left(\frac{1-\epsilon^{(j)}}{\epsilon^{(j)}}\right) > 0$

$Pr^{(i,j+1)} = Pr^{(i,j)} \cdot \exp(-\alpha_j y^{(i)} d^{(j)}(x^{(i)})), \forall i$

(5) Normalize $Pr^{(i,j+1)}, \forall i$

$Pr^{(i,j+1)} \leftarrow \left(\sum_i Pr^{(i,j+1)}\right)^{-1}$

Testing. 1. Given x , 算出 $\hat{y}^{(1)}, \dots, \hat{y}^{(j)}, \dots, \forall j$

2. Final prediction

$$\tilde{y} = \sum_j \alpha_j d^{(j)}(x)$$

(4) Adaboost Work 的原因.

Def. ^① The margin of a prediction of an example $(x^{(i)}, y^{(i)}) \in X$ as

margin $(x^{(i)}, y^{(i)}) \triangleq y^{(i)} f(x^{(i)})$ 越遠離 0 越好

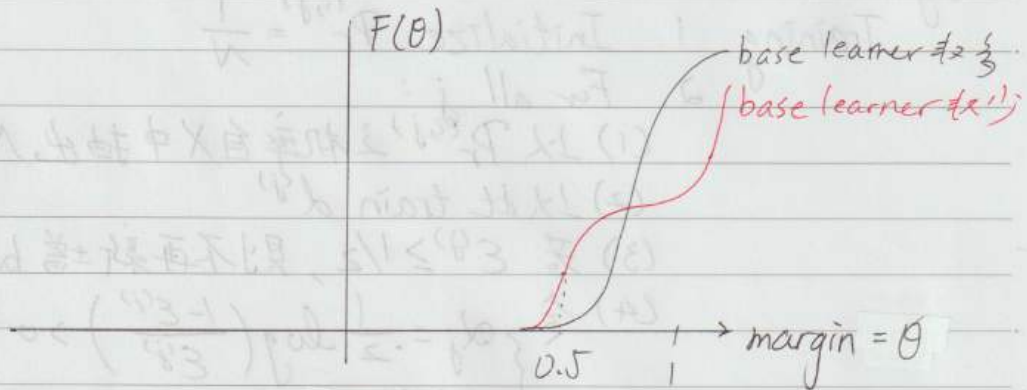
$$= \sum_{j \in \{j | y^{(i)} = d^{(j)}(x^{(i)})\}} \alpha_j - \sum_{j \in \{j | y^{(i)} \neq d^{(j)}(x^{(i)})\}} \alpha_j$$

② The confidence of $x^{(i)}$ is $f(x^{(i)})$

→ 離 decision boundary 越遠, confidence 越高

則: margin 之 cdf 可由

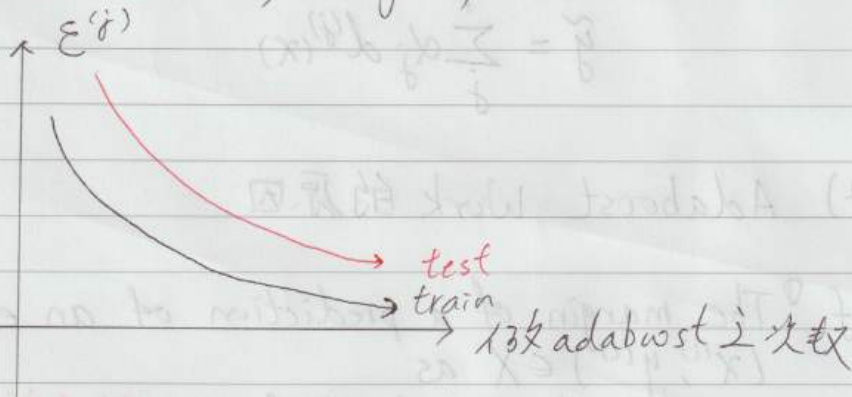
$$F(\theta) = P_X(y^{(i)} f(x^{(i)}) \leq \theta) \approx \frac{|\{(x^{(i)}, y^{(i)}) : y^{(i)} f(x^{(i)}) \leq \theta\}|}{|X|}$$



\Rightarrow base learner 少, 仍有一些 training data 之 margin < 0.5

" 多, 几乎没有 "

\Rightarrow base learner 數 \uparrow , margin \uparrow



\Rightarrow 無 overfit 之問題!

Ch9 Large Scale Machine Learning

一. Big Data

1. Big Data 的特性 — 4V

(1) Volume \rightarrow TB 甚至 PB

(2) Variety \rightarrow data 可能 structured or unstructured

(3) Velocity \rightarrow 每一單位新增很多資料

(4) Veracity: \because 資料搜集自動化

\therefore 用 implicit 的方法 (ex. 瀏覽器)

\Rightarrow 可能有搜不到的情形

\Rightarrow 每個 domain 的資料可能不是那麼健全/precise

2. Deep learning 的優勢:

(1) NN 很深時 \Rightarrow representation learning

Assumption 1. Smooth assumption: If $x \approx x^{(i)} \in X$, then $f(x; w) \sim f(x^{(i)}; w)$, w is the parameter of the model

2. 深層的 factor 可以 recursively 用淺層的 factor 來表示.

By assumption 2 \Rightarrow 深的 NN 相較淺的 NN 只需要較少的 data 就能 train 出來

\Rightarrow 消除 **curse of dimension**!

\Rightarrow 其他 model 很少見!

(2) 此 factor 為 nonlinear! \Rightarrow 未必都一定要用 deep learning
(\because By no free lunch thm)

⇒ When the function f to learn is
 { complex or have composite pattern ⇒ deep learning
 { simple or linear ⇒ other large-scale ML techniques

3. Curse of dimension. 假设 data 是散落在 \mathbb{R}^n 的 bin (小格子) 中
 ⇒ data 维数越高, 则为了要 interpolate, 所需的 example 维数指数增加!

二. Deep Learning 的特征

1. 可平行化:

(1) Data parallelism ⇒ 每个 CPU 拿一块数据同时做 SGD (以 partitioned data), 再给 parameter server 去更新 Δw

(2) model Parallelism ⇒ 每个 CPU 看到的到完整的数据, 但只负责 train 一小部份之 model

2. model 效能受 train 的时间限制

Recall from learning theory:

Notation: f_N : The function learned from N examples X

f^* : The true function

F : The model space

f_F^* : The best model we can have in F .

Def. ① The empirical risk $C_N[f_N] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f_N(x^{(i)}), y^{(i)})$

② The Expected risk: $C[f_N] = \int \text{loss}(f(x), y) dP(x, y)$
 → 无法 minimize!

③ $f_F^* \triangleq \arg \min_{f \in F} C[f]$

④ The excess error $\mathcal{E} = C[f_N] - C[f^*]$.

在 learning theory 中, 把 \mathcal{E} 拆成 approximation error 和 estimation error ($\mathcal{E}_{app}, \mathcal{E}_{est}$)

$$\mathcal{E} = \underbrace{C[f_F^*] - C[f^*]}_{\mathcal{E}_{app}} + \underbrace{C[f_N] - C[f_F^*]}_{\mathcal{E}_{est}}$$

⇒ 但受 training time 限制, 沒我們 train 出來的 model 為 \tilde{f}_N , 則 $C_N[\tilde{f}_N] \leq C_N[f_N] + \rho$.

$$\Rightarrow \mathcal{E} = \underbrace{C[f_F^*] - C[f^*]}_{\mathcal{E}_{app}} + \underbrace{C[f_N] - C[f_F^*]}_{\mathcal{E}_{est}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{opt}}$$

⇒ To reduce

- $\mathcal{E}_{app} \rightarrow$ Choose a larger model
- $\mathcal{E}_{est} \rightarrow$ increase N or choose a smaller model
- $\mathcal{E}_{opt} \rightarrow$ train 久一點 or 挑一個比較快的 optimizer.

(1) small scale ML tasks

∵ time is not an issue

∴ \mathcal{E}_{opt} 可很小

⇒ trade off between \mathcal{E}_{app} and \mathcal{E}_{est}

(2) large scale ML tasks

∴ \mathcal{E}_{opt} 要 ↓

∴ 用 SGD

∴ N 大

∴ \mathcal{E}_{est} 小

⇒ $\mathcal{E}_{app} \downarrow$, 故 prefer large model
為了使

三. NN 与 Gaussian Process 之關係: \Rightarrow omit

這部份太複雜, 有概念說很大又很寬的 NN 在 train 到很久的話會呈現 Gaussian process 就可以了.

$$[1, 1, 0] - [1, 1, 0] + [1, 1, 0] - [1, 1, 0] = 0$$

cost cost

$$[1, 1, 0] - [1, 1, 0] + [1, 1, 0] - [1, 1, 0] + [1, 1, 0] - [1, 1, 0] = 0$$

cost cost cost

not reduce

Cost \rightarrow Choose a larger model \rightarrow cost \uparrow
 Cost \rightarrow increase N or choose a smaller model \rightarrow cost \downarrow
 Cost \rightarrow train R or the optimizer \rightarrow cost \downarrow

\Rightarrow small scale ML tasks
 ; this is not an issue
 ; cost is high
 \Rightarrow trade off between cost and cost

(a) large scale ML tasks
 ; cost is high
 ; cost is high
 \Rightarrow cost \downarrow
 ; cost is high
 \Rightarrow copy the large model
 ; cost

Ch10 Neural Networks: Design

NO.

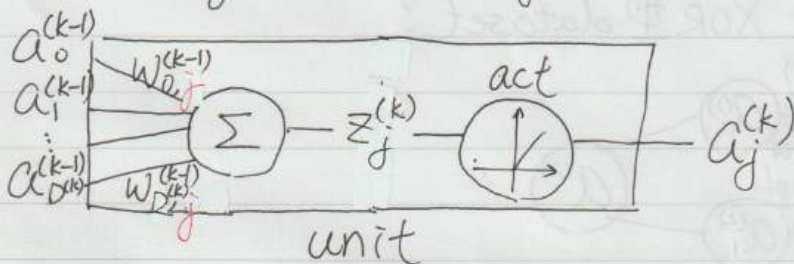
DATE

一. 什麼是神經網路

Def. Suppose parameters $\theta^{(1)}, \dots, \theta^{(L)}$ can be learn from training set X . A **feed forward neural networks** is a function composition

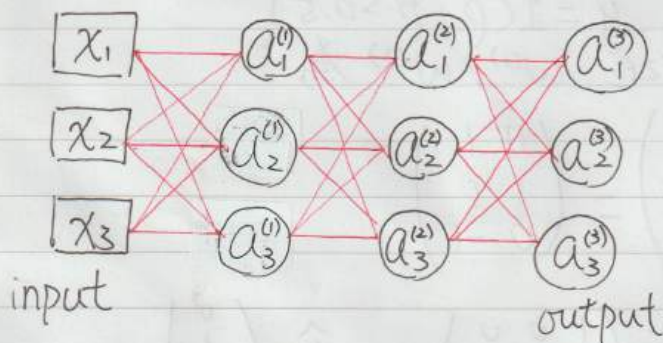
$$\hat{y} = f^{(L)}(\dots f^{(2)}(f^{(1)}(x; \theta^{(1)}); \theta^{(2)}); \dots); \theta^{(L)})$$

- ① Each function is called a **layer**
- ② Let $\theta^{(k)} = (w^{(k)}, b^{(k)})$, and $f^{(k)}$ be a **nonlinear function**. Suppose $a^{(k)} \in \mathbb{R}^{D^{(k)}}$ is the output of $f^{(k)}$, and $f^{(k)}$ **acts elementwisely**, then we denote $f^{(k)}$ by $\text{act}^{(k)}$, and $a^{(k)} = \text{act}^{(k)}(w^{(k)T} a^{(k-1)} + b^{(k)})$. Such $f^{(k)}$ is called an **activation function** ($\mathbb{R} \rightarrow \mathbb{R}$)
- ③ Each $f_j^{(k)} = \text{act}^{(k)}(w_{:j}^{(k)T} a^{(k-1)})$ is called a **unit**



常直接以 " $a_j^{(k)}$ " 表達 unit.

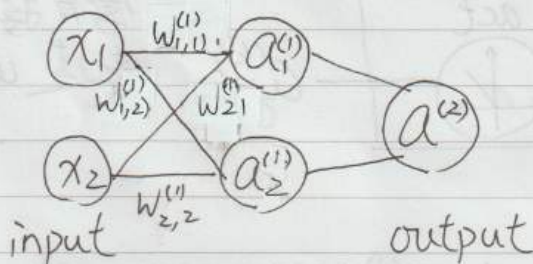
ex. A neural network $\hat{y} = f^{(3)}(f^{(2)}(f^{(1)}(x; \theta^{(1)}); \theta^{(2)}); \theta^{(3)})$. $x \in \mathbb{R}^3$, each output $a^{(k)} \in \mathbb{R}^3$, $\hat{y} = a^{(3)}$.



* 我們常忽略 "1" 這個每層都有的 unit 以節省空間; 事實上, 放 "1" 是代表把 $w_{:j}^{(k)T} a^{(k-1)} + b^{(k)}$ 的 $b^{(k)}$ 融合進 $w_{:j}^{(k)T}$ 裡.

- Remark. ① 並沒有規定每個 layer 之 output $a^{(k)}$ 的維度要一樣
- ② output layer 做 binary classification \rightarrow 可只用一個 unit; 做一般 classification 可用與 class 數相同的 unit \Rightarrow output \hat{y} 可視為一 pmf 或機率分佈
- ③ act 函數在 layer 數 > 3 時不可用 logistic function, 要用非線性函數, 否則 train 不起來
- ④ NN 中間除了 output layer 以外的 layer 稱為 hidden layer, 功用如同做了 *feature engineering* 一般.

ex. 給定 $X = \{(x_1, x_2, y) : (0, 0, 0), (0, 1, 1), (1, 0, 1), (1, 1, 0)\}$
 以及 NN: $\hat{p} = f^{(2)}(f^{(1)}(x; \theta^{(1)}); \theta^{(2)})$, 試說明此 NN 如何學習此 XOR 型 dataset?



-sol- 取 $a^{(1)} = \max(0, w^{(1)T}x)$

$$a^{(2)} = \sigma(w^{(2)T}a^{(1)}) \cong \hat{p}$$

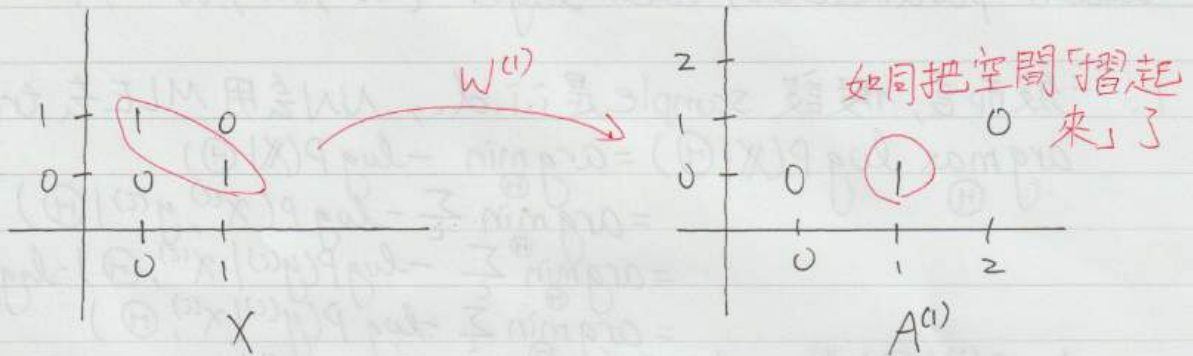
prediction: $\hat{y} = 1(\hat{p}; \hat{p} > 0.5)$

先隨機 initialize $w^{(1)}, w^{(2)}$ 為

$$W^{(1)} = \begin{pmatrix} w_{0,1}^{(1)} & w_{0,2}^{(1)} \\ w_{1,1}^{(1)} & w_{1,2}^{(1)} \\ w_{2,1}^{(1)} & w_{2,2}^{(1)} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix},$$

$$W^{(2)} = \begin{pmatrix} -1 \\ 2 \\ -4 \end{pmatrix}, \quad X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}, \quad \hat{y} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

$$\text{則 } XW^{(1)} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{pmatrix}$$



$$\text{令 } A^{(1)} = \left[1(\max(0, (XW^{(1)})_{ij})) \right]$$

$$= \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 2 & 1 \end{pmatrix}$$

$$\Rightarrow a^{(2)} = \sigma(A^{(1)}W^{(2)})$$

$$= \sigma \left[\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 2 & 1 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \\ -4 \end{pmatrix} \right] = \sigma \left[\begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix} \right]$$

$$= \begin{pmatrix} \sigma(-1) \\ \sigma(1) \\ \sigma(1) \\ \sigma(-1) \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix} \neq$$

不过這是我们「猜到」了一个好的 $W^{(1)}$, $W^{(2)}$. 一般而言, 如何訓練一个NN呢?

二. NN之訓練: Back propagation

Given examples: $X = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$
learn parameters of each layer $\{W^{(1)}, \dots, W^{(L)}\}$.

1. 一般而言, 假設 sample 是 i.i.d., NN 會用 MLE 去 train

$$\begin{aligned} \arg \max_{\Theta} \log P(X|\Theta) &= \arg \min_{\Theta} -\log P(X|\Theta) \\ &= \arg \min_{\Theta} \sum_{i=1}^N -\log P(x^{(i)}, y^{(i)}|\Theta) \\ &= \arg \min_{\Theta} \sum_{i=1}^N -\log P(y^{(i)}|x^{(i)}, \Theta) - \log P(x^{(i)}|\Theta) \\ &= \arg \min_{\Theta} \sum_{i=1}^N -\log P(y^{(i)}|x^{(i)}, \Theta) \end{aligned}$$

令 $C^{(i)}(\Theta)$ 為 $-\log P(y^{(i)}|x^{(i)}, \Theta)$, 稱為 *cost function*

ex. 設 $P(y=1|x) \sim B(1, p)$, $x \in \mathbb{R}^D$, $y \in \{0, 1\}$.

各層 output $a^{(L)} = \sigma(z^{(L)})$, 則 cost function 為:

$$\begin{aligned} C^{(i)}(\Theta) &= -\log P(y^{(i)}|x^{(i)}; \Theta) \\ &= -\log (a^{(L)} y^{(i)} (1-a^{(L)})^{1-y^{(i)}}) \\ &= -\log (\sigma(z^{(L)}) y^{(i)} (1-\sigma(z^{(L)}))^{1-y^{(i)}}) \\ &= -\log \sigma[(2y^{(i)}-1)z^{(L)}] \\ &= \xi[(1-2y^{(i)})z^{(L)}], \xi \text{ 為 softplus function} \end{aligned}$$

2. 大部份 NN 都可用 SGD 去解 $\arg \min_{\Theta} \sum_{i=1}^N C^{(i)}(\Theta)$

Algorithm: SGD

1. initialize $\Theta^{(0)}$ randomly
2. randomly partition X into *minibatches* of size M
3. $\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \eta \nabla_{\Theta} \sum_{i=1}^M C^{(i)}(\Theta^{(t)})$
4. back to 2, repeat until convergence.

優點在於收斂快, 且只要進行乘法與加法, 可以丟給 GPU, 但前提是要有了一個算 $\nabla_{\Theta} \sum_{i=1}^N C^{(i)}(\Theta^{(i)})$ 之好方法

3. 反向傳播

總之，我們要算出 $\nabla_{\Theta} \sum_{n=1}^M C^{(n)}(\Theta^{(t)}) = \sum_{n=1}^M \nabla_{\Theta} C^{(n)}(\Theta^{(t)})$

又， $\Theta = \{W^{(1)}, W^{(2)}, \dots, W^{(L)}\}$ ，方便起見，

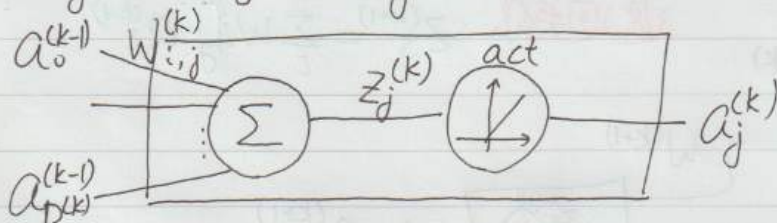
$$\therefore W^{(k)} = [W_{ij}^{(k)}]$$

$$= \begin{pmatrix} W_{0,0}^{(k)} & W_{0,1}^{(k)} & \dots & W_{0,D^{(k)}}^{(k)} \\ W_{1,0}^{(k)} & W_{1,1}^{(k)} & \dots & W_{1,D^{(k)}}^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ W_{D^{(k)},0}^{(k)} & W_{D^{(k)},1}^{(k)} & \dots & W_{D^{(k)},D^{(k)}}^{(k)} \end{pmatrix}$$

故，可將 $\nabla_{\Theta} C^{(n)}(\Theta^{(t)})$ 寫成 $\frac{\partial C^{(n)}}{\partial W_{ij}^{(k)}}$ ，可看成一個超長之

向量，只要 i, j, k 都能算就好了。由 chain rule:

$$\frac{\partial C^{(n)}}{\partial W_{ij}^{(k)}} = \frac{\partial C^{(n)}}{\partial z_j^{(k)}} \cdot \frac{\partial z_j^{(k)}}{\partial W_{ij}^{(k)}}$$



(1) Forward pass: 解 $\frac{\partial z_j^{(k)}}{\partial W_{ij}^{(k)}}$

$$\text{When } k=1, z_j^{(1)} = \sum_c W_{ij}^{(1)} x_c^{(n)}$$

$$\Rightarrow \frac{\partial z_j^{(1)}}{\partial W_{ij}^{(1)}} = x_c^{(n)}$$

$$\text{When } k>1, z_j^{(k)} = \sum_c W_{ij}^{(k)} a_c^{(k-1)}$$

$$\Rightarrow \frac{\partial z_j^{(k)}}{\partial W_{ij}^{(k)}} = a_c^{(k-1)}$$

⇒ 從最靠近 input 之層開始，逐層往後計算偏微分。

(2). Backward pass: 解 $\frac{\partial C^{(n)}}{\partial z_j^{(k)}} \triangleq \delta_j^{(k)}$

When $k=L \Rightarrow \delta^{(L)}$ 取決於最後一層如何設計!

ex. binary classification:

$$\delta_j^{(L)} = \frac{\partial C^{(n)}}{\partial z_j^{(L)}} = \frac{\partial \mathcal{L}(1 - z_j^{(n)}) z_j^{(L)}}{\partial z_j^{(L)}} = \sigma((1 - z_j^{(n)}) z_j^{(L)}) (1 - 2z_j^{(n)})$$

When $k < L \Rightarrow$ 考慮

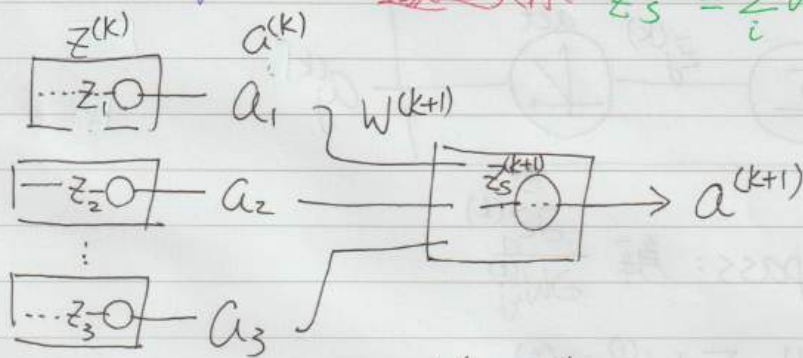
$$\delta_j^{(k)} = \frac{\partial C^{(n)}}{\partial z_j^{(k)}} = \frac{\partial C^{(n)}}{\partial a_j^{(k)}} \cdot \frac{\partial a_j^{(k)}}{\partial z_j^{(k)}}$$

$$= \frac{\partial C^{(n)}}{\partial a_j^{(k)}} \cdot \text{act}'(z_j^{(k)}) \quad \text{by 多維 chain rule}$$

$$= \left(\sum_s \frac{\partial C^{(n)}}{\partial z_s^{(k+1)}} \cdot \frac{\partial z_s^{(k+1)}}{\partial a_j^{(k)}} \right) \text{act}'(z_j^{(k)})$$

$$= \left(\sum_s \delta_s^{(k+1)} \cdot W_{j,s}^{(k+1)} \right) \text{act}'(z_j^{(k)})$$

遞迴關係 $z_s^{(k+1)} = \sum_i W_{i,s}^{(k)} a_i^{(k)}$



(deepest)

故我們可從最靠近 output 那一層開始, 逐層往回算 $\delta_j^{(k)}$

$$\delta_j^{(k)} = \left(\sum_s \delta_s^{(k+1)} W_{j,s}^{(k+1)} \right) \text{act}'(z_j^{(k)})$$

tensortflow 會直接寫好
ex. `def (x^2): return 2x`

由此一來, 組合 forward 跟 backward pass, 就能寫出 train function 了:

(3) Backprop with minibatch size $M=1$.

1. Input: $(x^{(n)}, y^{(n)})$ and $\Theta^{(t)}$.

2. Forward pass:

$$a^{(0)} \leftarrow [1, x^{(n)}]^T$$

for $k=1:L$

$$z^{(k)} \leftarrow W^{(k)T} a^{(k-1)}$$

$$a^{(k)} \leftarrow \text{act}(z^{(k)})$$

end.

* batch 設太大
對 training 有害!

3. Backward pass:

compute $\delta^{(L)}$

for $k=L-1:1$

$$\delta^{(k)} \leftarrow \left(\sum_s \dot{\delta}_s^{(k+1)} w_{j,s}^{(k+1)} \right) \text{act}'(z_j^{(k)})$$

end

4. Return $\frac{\partial C^{(n)}}{\partial W^{(k)}} = \left[\frac{\partial z_j^{(k)}}{\partial w_{ij}^{(k)}} \right]_{ij}^{(k)} = \left[a_i^{(k-1)} \delta_j^{(k)} \right]_{ij}^{(k)}$

(4) Backprop with minibatch size $M > 1$

\Rightarrow see PPT p.25.

三. Neuron Design for output layers

對於不同的任務，我們喜歡給予不同的 loss function，連帶的，最後一層 unit 之設計也要不一樣，以便於計算 $\delta^{(L)}$

1. Cross entropy.

Def. 令 $\hat{P}(y|x)$ 為 NN 之 output，則對 dataset 之分佈 $(x, y) \sim \text{Empirical}(X)$ 而言，cross entropy (or KL Divergence) 為 $-E_{(x,y) \sim \text{Emp}(X)} [\log \hat{P}(x|y)]$

⇒ 我們希望 NN 能找到

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} -E_{(x,y) \sim \text{Emp}(X)} (\log \hat{P}(y|x))$$

2. Sigmoid units \leftrightarrow $B(1, p)$.

若 $P(y=1|x) \sim B(1, p)$, $y \in \{0, 1\}$ 且 $p \in (0, 1)$, 則我們會把最後一層 (只有一個 unit) 設計成

$$\hat{p} = a^{(L)} = \sigma(z^{(L)}) = \frac{e^{z^{(L)}}}{e^{z^{(L)}} + 1}$$

$$\begin{aligned} \Rightarrow \delta^{(L)} &= \frac{\partial C^{(n)}}{\partial z^{(L)}} = \frac{\partial}{\partial z^{(L)}} -\log \hat{P}(y^{(n)} | x^{(n)}; \hat{\theta}) \\ &= (1 - 2y^{(n)}) \sigma[(1 - 2y^{(n)}) z^{(L)}] \end{aligned}$$

$$(1) \delta^{(L)} \rightarrow 0 \Leftrightarrow \begin{cases} y^{(n)} = 1, z^{(L)} \gg 0 \\ y^{(n)} = 0, z^{(L)} < 0 \end{cases}$$

(2). $C^{(n)}$ saturates only when \hat{p} is correct

(3). Sigmoid units 是專為 SGD 設計的

3. Softmax units \leftrightarrow Categorical output distribution

Def. 令 $p = (p_1, \dots, p_k)$, $\sum_i p_i = 1$, 則定義 categorical output distribution 為 $\text{Cat}(p)$, 其 pmf 為

$$\begin{cases} P(y=e_1) = p_1 \\ P(y=e_2) = p_2 \\ \vdots \end{cases}$$

$$P(y=e_k) = p_k \quad i^{\text{th}}$$

其中 $e_i = (0, \dots, 0, 1, 0, \dots, 0)$ (one-hot encoding).

設 $P(y|x) \sim \text{Cat}(p)$. 我們常用 **softmax** 放在最後一層:

$$\hat{p}_j = a_j^{(L)} = \text{softmax}(z^{(L)})_j = \frac{e^{z_j^{(L)}}}{\sum_{i=1}^K e^{z_i^{(L)}}}$$

$$\left(\text{or } = \frac{e^{z_j^{(L)}}}{\sum_{i=1}^{K-1} e^{z_i^{(L)}} + 1} \right)$$

$$\Rightarrow \delta_j^{(L)} = \frac{\partial C^{(n)}}{\partial z_j^{(L)}} = \frac{\partial}{\partial z_j^{(L)}} -\log \left(\prod_i \hat{p}_i^{1(y^{(n)}; y^{(n)}=i)} \right)$$

$$\left\{ \begin{array}{l} \textcircled{1} y^{(n)}=j \\ \textcircled{2} y^{(n)}=i \neq j \end{array} \right. = -\frac{\partial \log \hat{p}_j}{\partial z_j^{(L)}} = -\frac{1}{\hat{p}_j} (\hat{p}_j - \hat{p}_j^2) = \hat{p}_j - 1$$

$$= -\frac{\partial \log \hat{p}_i}{\partial z_j^{(L)}} = -\frac{1}{\hat{p}_i} (-\hat{p}_i \hat{p}_j) = \hat{p}_j$$

* $\delta_j^{(L)} \rightarrow 0$ if \hat{p}_j is correct.

4. Linear units \leftrightarrow Gaussian

Assume $P(y|x) \sim \mathcal{N}(\mu, \Sigma)$, 則我們用 linear units:

$$a^{(L)} = \hat{\mu} = z^{(L)}$$

$$\Rightarrow \delta^{(L)} = \frac{\partial C^{(n)}}{\partial z^{(L)}} = \frac{\partial}{\partial z^{(L)}} -\log g(y^{(n)}; \hat{\mu}, \Sigma), \text{ 設 } \Sigma = I,$$

$$\text{則 } \delta^{(L)} = \frac{\partial}{\partial z^{(L)}} \|y^{(n)} - z^{(L)}\|^2$$

此種 unit 不會 saturate, 但也無法 train 太準.

四. Neuron design for hidden layers

hidden layer 也可以用不同的 activation function 去設計; 各種設計有不同需要考量的重負。

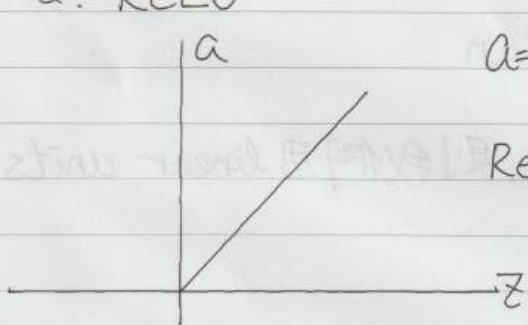
1. 梯度消失問題

(1) 背景: backward pass 中; $\delta_j^{(k)} = \sum (\delta_s^{(k+1)} \cdot w_{j,s}^{(k+1)}) \text{act}'(z_j^{(k)})$
 \Rightarrow 若 $\text{act}'(\cdot) < 1$, 則 $\delta_j^{(k)}$ 將隨 $k \uparrow$ 而越來越小
 \Rightarrow 超出系統精度後, 直接變成 0!!

(2) 即使沒變成 0 $\Rightarrow \delta_j^{(k)}$ 太小, 難以 train, 拖累時間!

(3) logistic function 有此問題, 故不應以其做為 hidden layer 的 activation function

2. ReLU



$$a = \text{ReLU}(z) = \begin{cases} z^{(k)}, & z^{(k)} \geq 0 \\ 0, & z^{(k)} < 0 \end{cases}$$

$$\text{ReLU}'(z) \equiv \begin{cases} 1, & z^{(k)} > 0 \\ 0, & z^{(k)} < 0 \\ \text{隨機給 } 0 \text{ 或 } 1, & z^{(k)} = 0 \end{cases}$$

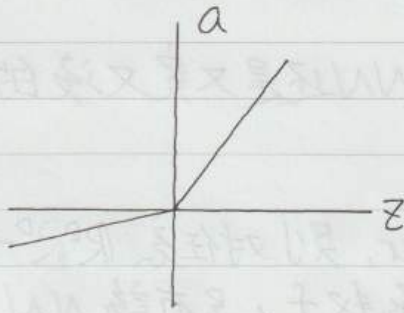
\Rightarrow 解決梯度消失問題

(1) $\text{ReLU}''(\cdot) \equiv 0$ on \mathbb{R} , 以消除 second order effect

(2) 缺點: $\forall \delta_j^{(k)} = 0$, 對應之 weight $w_{ij}^{(k)}$ 不被 update!

$$\left(\frac{\partial C^{(n)}}{\partial w_{ij}^{(k)}} \right) = \delta_j^{(k)} \frac{\partial z_j^{(k)}}{\partial w_{ij}^{(k)}}$$

(3) 改善方法: 使用 **Leakly ReLU**



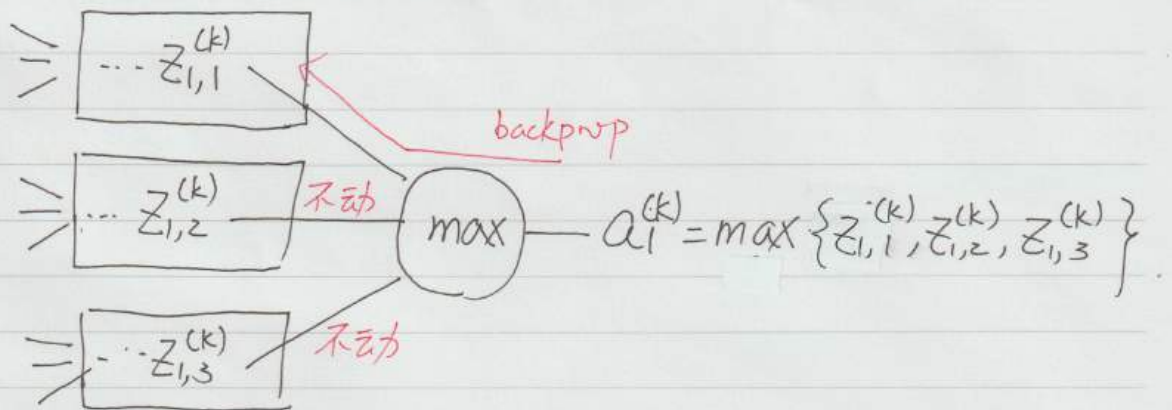
$$\text{act}(z^{(k)}) = \max(\alpha z^{(k)}, z^{(k)}),$$

for some $\alpha \in \mathbb{R}$

\Rightarrow $\begin{cases} \alpha \text{ fixed (先给定): Leakly Relu} \\ \alpha = -1: \text{absolute value rectification} \\ \alpha \text{ learned by GD: parametric relu} \end{cases}$

3. Maxout unit:

$$\text{act}(z^{(k)})_j = \max_s z_{j,s}^{(k)}$$



(1) backprop 時, 只要調最大那個就可以了

(2) 优点: 避免「catastrophic forgetting」

(3) 缺点: 需大量 data, 否則要做 regularization

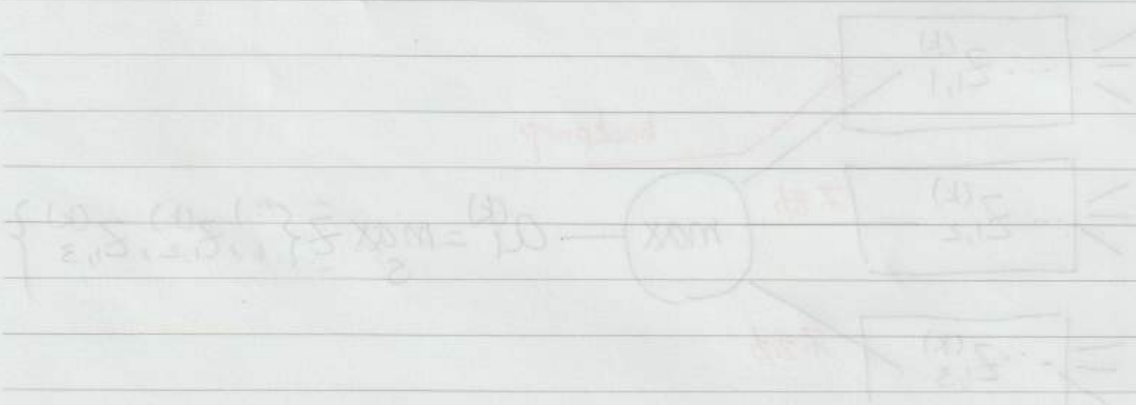
4. NN之結構選擇

Q: 我們應該 train 一個又深又窄的 NN 還是又寬又淺的 NN?

Thm. (Universal Approximation Thm)

一個 NN 若至少有一個 hidden layer, 則對任意定義在 \mathbb{R}^D 上 closed and bounded 子集的連續函數 f , 只要該 NN 足夠寬, 就能逼近 f .

問題是 data 有限, 未必做的到! 而深的 NN 可能可以把 data 分佈之 space 摺起來, 建立起很多 piecewise linear 的 region \Rightarrow better generalizability.



Ch11 Neural Networks: Optimization & Regularization

DATE

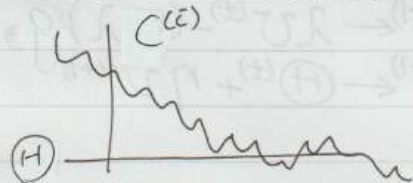
11-1 Optimization

Recall: 用 NN train - function, 我們希望找出
 $\arg\min_{\Theta} C(\Theta) = \sum_c C^{(c)}(\Theta)$

一. 常見的困難

1. SGD 常走不出 local minima 和 saddle point

(1) NN 之 cost function 呈鉅齒狀



(2) 小 NN 比較會遇到此問題, 大 NN 還好

2. 若 $C^{(c)}$ 太 ill-conditioned, 即 $\frac{\lambda_{\max}}{\lambda_{\min}}$ 太大, 則訓練緩慢

成因: SGD 會繞遠路和抖動 (PPT p6)

3. 若 $C^{(c)}$ 跟本沒有 global minimum

⇒ SGD 會永遠往同一方向衝下去

解決方案: initialization 要做好 (做得好, 一萬層的 NN 都 train 的起來)

4. training 小撇步

(1) 對 input 做 standardize

⇒ prevent dominating features
↳ improves conditioning

(2) Initialize all weights to small random values
 → 否則会有很多一樣的 weight, unit 之功能重複, 使效果不彰

(3) Early stop: prevent overfitting.

二. Momentum

1. update rule of SGD: $\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \eta \nabla_{\Theta} C(\Theta^{(t)})$

⇒ 定義 momentum 為 $v^{(t+1)} \leftarrow \lambda v^{(t)} - (1-\lambda) \nabla_{\Theta} C(\Theta^{(t)})$
 $\Theta^{(t+1)} \leftarrow \Theta^{(t)} + \eta v^{(t+1)}$

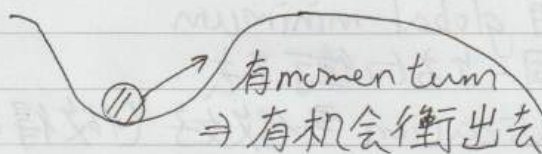
2. Nesterov momentum

$$\tilde{\Theta}^{(t+1)} \leftarrow \Theta^{(t)} + \eta v^{(t)}$$

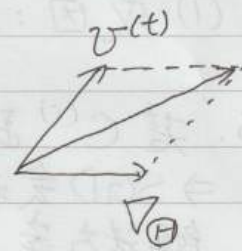
$$v^{(t+1)} \leftarrow \lambda v^{(t)} - (1-\lambda) \nabla_{\Theta} C(\tilde{\Theta}^{(t+1)})$$

$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} + \eta v^{(t+1)}$$

3. 优点: ① 避免卡在 local minimum



② 更快收斂



4. 缺点: 沒有 minima 時沒有幫助

三. AdaGrad & RMSProp & Adam

1. 前言: 卡在 saddle point \rightarrow 沒救, initialize 做好一美
 { flat valley \rightarrow 把 learning rate 弄得 adaptive 試着

ex. 在陡峭的方向把 learning rate 刪小一美
 “平坦” “大”



2. AdaGrad 的 update rule:

$$\begin{aligned} r^{(t+1)} &\leftarrow r^{(t)} + g^{(t)} \odot g^{(t)} \\ \Theta^{(t+1)} &\leftarrow \Theta^{(t)} - \frac{\eta}{\sqrt{r^{(t+1)}}} \odot g^{(t)} \end{aligned}, \quad g^{(t)} = \nabla_{\Theta} C(\Theta)$$

其中: $\frac{\eta}{\sqrt{r^{(t+1)}}} = \frac{\eta}{\sqrt{t+1}} \odot \frac{1}{\sqrt{\frac{1}{t+1} r^{(t+1)}}}$

$$= \frac{\eta}{\sqrt{t+1}} \odot \frac{1}{\sqrt{\frac{1}{t+1} \sum_{i=0}^t g^{(i)} \odot g^{(i)}}}$$

learning rate $\downarrow \Rightarrow$ SGD 不亂跳

∇ 大, 此項小
 ∇ 小, “大”

3. 缺美: 實際上效果不大好 \Rightarrow 改良: RMSProp, Adam

4. RMSProp:

$$\begin{aligned} r^{(t+1)} &\leftarrow \lambda r^{(t)} + (1-\lambda) g^{(t)} \odot g^{(t)} \\ \Theta^{(t+1)} &\leftarrow \Theta^{(t)} - \frac{\eta}{\sqrt{r^{(t+1)}}} \odot g^{(t)} \end{aligned}$$

5. Adam

$$\begin{aligned} v^{(t+1)} &\leftarrow \lambda_1 v^{(t)} + (1 - \lambda_1) g^{(t)} \\ r^{(t+1)} &\leftarrow \lambda_2 r^{(t)} + (1 - \lambda_2) g^{(t)} \odot g^{(t)} \\ \Theta^{(t+1)} &\leftarrow \Theta^{(t)} + \frac{\eta}{\sqrt{r^{(t+1)}}} \odot v^{(t+1)} \end{aligned}$$

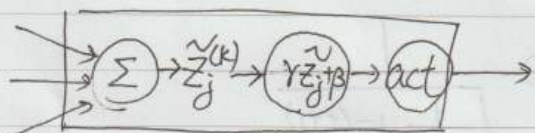
四. Batch Normalization

1. 前言: 若 $C(\Theta)$ 太过 ill-conditioned, 則 gradient-based 的方法 (如 SGD) 效率不好 (PPT p19-21)

2. 考慮把 $a^{(k)}$ standardized, 則 $g_L^{(t)} = \frac{\partial C}{\partial w^{(L)}}(\Theta^{(t)})$ 更可能讓 $C \downarrow$

\Rightarrow Batch normalization: $\tilde{a}_i^{(k)} = \frac{a_i^{(k)} - \mu^{(k)}}{\sigma^{(k)}}$, $\forall i$

3. 也可以把 $z_i^{(k)}$ 標準化, 並修改 unit 為



其中 $\tilde{z}_j^{(k)} = \frac{z_j^{(k)} - \mu^{(k)}}{\sigma^{(k)}}$, $\forall j$.

\Rightarrow 利用 backprop 去 learn γ, β 以及 $w^{(k)}$

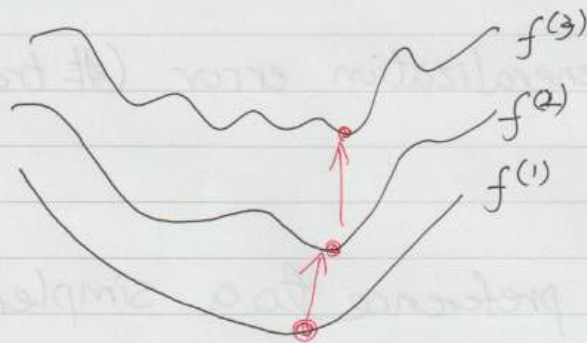
五. Continuation methods & Curriculum Learning

1. 如何找比較好的 initialization of Θ ?

(1) train 一个 NN 很多次, 每次都使用不同的 random initial point \Rightarrow 选一个最好的

(2) 先設計一个簡單的 cost function $f^{(1)}$, 再以对 $f^{(1)}$ 之最佳解 $\Theta^{(1)}$ 做为稍複雜之 cost function $f^{(2)}$ 之 initial point.

⇒ 此方法称、*continuation methods*



2. 不过 cost function 說不定不是 convex! 此外,
对 缺乏 minima 之 NN 也没啥用
⇒ DL 很少用

3. Curriculum learning (or shaping)

⇒ 先学簡單的問題

用簡單的問題先 train 出 weight

→ 当成下一个問題之 initial weight

ex. 要做 0★ ~ 5★ 評論的判斷器

→ 先用 0★ v.s. 5★ train 出 model,

→ 当成 0★, 3★, 5★ 判斷器的 initial weight

→ 当成你問題的 weight, train 出 model.

* 本章中較有用的方法:

- { standardize
- { good initialization
- { Adam
- { batch normalization

11-2 Regularization

NO.

DATE / /

一. Regularization

1. 目的: 減少 generalization error (≠ training error)

2. 方法:

(1) expressing preference to a simpler model

(2) provide different perspective on how to explain the training data

(3) Encoding prior knowledge

Q: data 很多時仍要 regularization 嗎?

A: 要! 真實世界的 data 不是用類似 model 的東西產生的!
要不然就是更複雜的 model 生出來的!

二. Weight Decay

Def. Weight Decay

When training an NN parameterized by Θ , weight decay is to add norm penalties when solving $\arg \min_{\Theta} C(\Theta) + \alpha \Omega(\Theta)$, where Ω can be such as L_1 or L_2 norm.

1. 好處: 避免 hidden layer 的 feature dominant, 即一顆 unit 有太大的 $z_j^{(k)}$.

2. Explicit norm penalties:

$\arg \min_{\Theta} C(\Theta)$ subject to $\Omega(\Theta) \leq R$

3. Projective SGD \rightarrow 解 explicit num penalties

\Rightarrow 先 update $\Theta^{(t+1)}$, 若 $\Theta^{(t+1)} \notin \text{feasible set}$, 則將其投影到 feasible set 上

\Rightarrow 伏美: 避免 **dead unit**
 (因 learning rate 过大造成的不穩定)

4. 最好一次使用 explicit constraints + reprojection + large learning rate (Hintum)

三. Data Augmentation

1. 前言: 理論上有越多 data 越好, 但假 data 很難製造!

\Rightarrow 不过, 假設今天要分辨手寫數字, 那把「5」轉 10° , 还是「5」

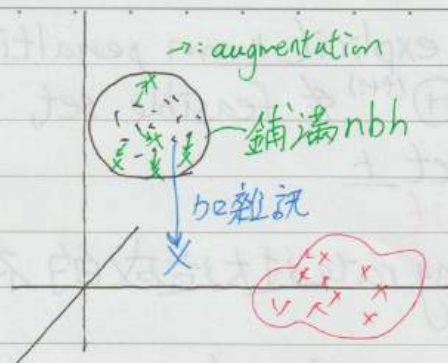
\Rightarrow 這種讓 dataset 「膨脹」的方法稱為 data augmentation

2. 在做 image classification 時, 由 scaling, translating, rotating, flipping... 生成新的 data (x, y)

** 不可 apply 會改變 $y^{(i)}$ 之 augmentation! 如: 把「6」轉成「9」

3. 加雜訊也是一種 augmentation, 但 NN 對雜訊有莫太敏感 (PPT p49)

這是由於 data augmentation 的原理是把同-class 之 data 所存在之 neighborhood 給連起來鋪滿
 (∴ dataset 相對 task 很大時可用 augmentation)



但加了雜訊後, data 跑到了「空洞」中, 而 NN 在空洞中的表現與人類很不一樣

4. Noise injection: data $(x^{(i)}, y^{(i)})$, 加入 random noise 到 $x^{(i)}$ 中, 使 cost function 對 weight 的一些小變化較不敏感
* 也可加到 hidden layer

④. Dropout

1. Ensemble method

(1) can improve generalizability by offering different explanations to x

(2) 分類:

- { Voting: 使 variance ↓
- { bagging: resample x ⇒ 讓 voter 不能串通, dependent ↓
- { boosting: 在不 overfitting 的前提下使預測信心 ↑

2. Drop out 的做法: SGD training 時, 每次 load 了 minibatch, 而每個 unit 有 α 之機率被关掉不運作
通常 input units 取 $\alpha=0.8$, hidden units 取 $\alpha=0.5$

3. 原理: dropout 等同是一種 **feature-based bagging**

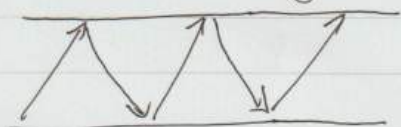
∴ ① resample input & latent features

② 等同 voter 間共享參數

4. 優點: improve generalization beyond ensemble

五. 其他方法:

1. cyclic learning rate: 使用 *learning rate schedule*



⇒ learning rate 在 - 給定形狀之範圍內 振盪.

2. manifold regularization

⇒ 假設 data 是 高維空間中的 *低維流形* K
被鑲嵌

⇒ $\forall x \in K$, 找出其 tangent space, 令其 basis 為 $\{v^{(i,j)}\}_j$ (tangent vector).

⇒ Tangent Prop 由以下方式求 classifier f :

$$Q[f] = \sum_{i,j} \nabla_x f(x^{(i)})^T v^{(i,j)}$$

但求 $v^{(i,j)}$ 可能比較困難, 要依靠 domain knowledge 或一樣用 learn 的.

3. Domain-Specific Prior knowledge

⇒ 針對不同 task 去設計 NN

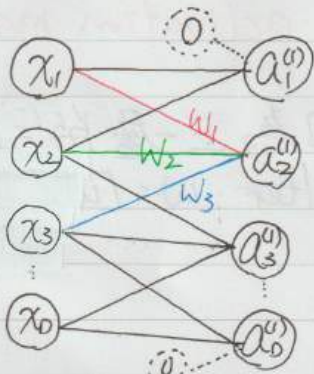
ex. Word 2 vec, CNN, ...

Ch12 Convolutional Neural Networks

DATE

一. 何謂 CNN

1. 先考慮只有一層的 NN, 接成下面的形式



- (1) 每一個 $a_i^{(1)}$ 都只接對應位置之上, 下, 的 unit
- (2) 對 layer 1 的每一個 unit, 規定所有連到其上的 weight 都要和其他人一樣

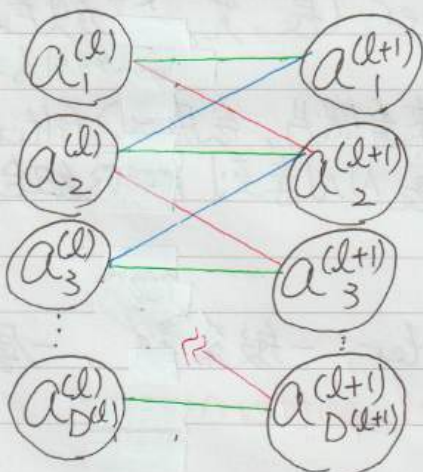
→ 實際上, \rightarrow 為一組

⇒ 每個第一層的 unit, 拿到的是同一組 weight; 對於邊緣的 unit, 假裝有一個 "0" 連到上面, 稱為 zero padding

(3) 因此, $a_i^{(1)} = \text{act}^{(1)}\left([x_{i-1} \ x_i \ x_{i+1}] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + b^{(1)}\right), \forall i=1, \dots, n. (\text{定 } x_0 = x_{n+1} = 0)$

- ⇒ 要學的 weight 只有 w_1, w_2, w_3 三個!
- ⇒ 比 fully connected layer 簡單多了!

2. 現在, 把 1 之 NN 推廣到任兩層間



- (1) 每一個 $a_i^{(l+1)}$ 都接 $k^{(l)}$ 個上-layer 之 unit, 分別是 $a_{i-\frac{k^{(l)}}{2}}^{(l)}, \dots, a_i^{(l)}, \dots, a_{i+\frac{k^{(l)}}{2}}^{(l)}$

(2) 不夠接的那些 (位在頭尾之) unit, 進行 zero padding

(3) 每個 layer $l+1$ 的 unit 所連到的 weight 皆為 $[w_1, \dots, w_{k^{(l)}}]$

⇒ 只有 $w_1, \dots, w_{k^{(l)}}$ 這 $k^{(l)}$ 個 weight 要 learn!

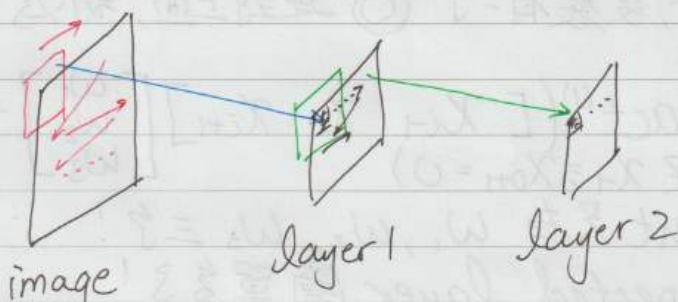
同理, $a_i^{(l+1)} = \text{act}^{(l+1)}\left([a_{i-\frac{k^{(l)}}{2}}^{(l)}, \dots, a_i^{(l)}, \dots, a_{i+\frac{k^{(l)}}{2}}^{(l)}] \begin{bmatrix} w_1 \\ \vdots \\ w_{k^{(l)}} \end{bmatrix} + b^{(l+1)}\right)$

Def. 由 2. 所定義的 layer 稱為 **convolution layer**
 $[w_1, \dots, w_{k^{(l)}}]$ 稱為 **filter** 或 **kernel**
 $act^{(l)}(\cdot)$ 稱為 **feature map** 或 **activation map**.

Remark. 會稱為 convolution 即是因為下層 y_i 是上層 $x = [x_1, \dots, x_n]$ 與 filter $w = [w_1, \dots, w_L]$ 做了 convolution

$$y_i = \sum_s x_s w_{i-s}$$

3. 2D 圖片的 convolution layer \Rightarrow 當然也做的到



$$(x * w)_{ij} = \sum_u \sum_v x_{u,v} w_{i-u, j-v}$$

反正總有辦法連起來~

4. 3D 圖片的 convolution layer (長, 寬, 三個顏色的 channel)
 \rightarrow 也可稱為 2D conv.

\Rightarrow 考量人類不會「一個一個顏色」來看照片, 會用 3D filter 去掃這 3D input, 以此讓 NN 學到 RGB 組合之 feature, 據說效果較好

5. 一層也不一定只能有一個 filter, 一般來說, 一層中有

一個 filter \Rightarrow detect 一個 **local pattern**
 多個 \Rightarrow " 多個 "

(1) 多個 filter 可以藉由增加下一層的厚度來完成

ex. 第 l 層有 $H \times W \times C$ 了 unit, 上有 C 了 $k^{(l)} \times k^{(l)} \times 3$ 的 filter

⇒ 下一層有 $W^{(l+1)} \times H^{(l+1)} \times C$ 个 unit, 如此继续
 { filter 大小為 $K^{(l+1)} \times K^{(l+1)} \times C^{(l+1)}$

(2). 优美 ⇒ 使很多个 local pattern 被叠起来
 ⇒ 更易学到新的 pattern

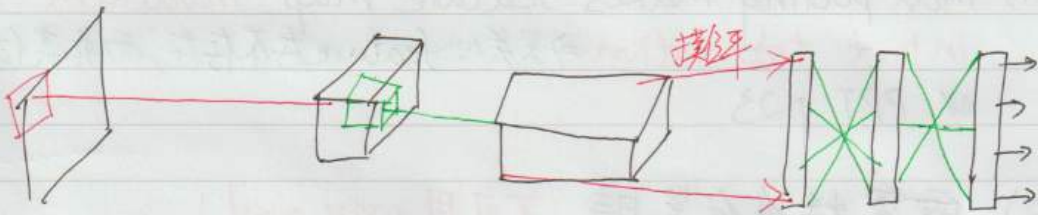
(3). 按理, 越深層的 convolution layer, 其中一个 unit 所对应原图片的范围 (称为 receptive field) 也越大

⇒ 越深層的 layer, 学到的是越 global 的 pattern.

6. 至此, 可以给出一个 CNN 的定义了:

Def. A Convolutional Neural network (CNN) is an NN having at least one convolution layer.

一个常见的 CNN 架构: 用很多 convolution layer 做完 pattern learning 后, 将其最后一层 **摊平** 成一维的 layer (**flatten layer**), 接著经过数个 fully connected layer 后 output.



二. Pooling layer

Def. A pooling layer is a layer **downsampling** the feature map.

1. 原理: 圖片的“細節”即使解析度變差了, 也不會讓人覺得人臉變得不是人臉, 樹不是樹.....

2. 方法:

(1) Max pooling: detecting edges; textures
ex. 偵測笑臉

(2) Average: detecting brightness or contrast

ex. 晴天/陰天之判別

s		k		k	
1	1	2	4		
5	6	7	8		
3	2	1	0		
1	2	3	4		

max pool →

6	8
3	4

* 一次取 $k \times k$ 個
做 pooling, k 稱為
filter 的大小

avg pool →

3.25	5.25
2	2

做完後移動 s 格, 再
pool 一次, s 稱為
stride

3. backprop 時, 在 forward pass 先記下 max 之 index (如上图所示), backward pass 時就只處理那些 index

4. 特異:

(1) max pooling makes feature map invariant to input translation ⇒ 只關心 feature 存不存在, 而非其位置
ex. PPT p23

(2) 需要精準位置時, 不可用 max pool!
ex. 自駕車

(3) maxpool 層本身沒有需要 train 的參數

* Exercise: 試求下列 CNN 每一層共有几个 unit? 共有几个要被 train 的 weight?

Input image: $256 \times 256 \times 3$

layer	種類	參數
1	Convolution	16 filter, $K=4$
2	Max pooling	$K=4$, $\text{stride}=4$
3	Convolution	32 filter, $K=4$
4	Maxpooling	$K=4$, $\text{stride}=4$
5	flatten	
6	fully connected	10个 unit
7.	output	softmax, 10个 unit

- sol- 1. \therefore 有 16 个 filter, 且 input 为 $256 \times 256 \times 3$
 \therefore 有 $256 \times 256 \times 16$ 个 unit, weight 共有 $4 \times 4 \times 3 \times 16$ 个.
2. maxpool 没有要 train 的 weight, 又 $K=4$, $\text{stride}=4$
 \Rightarrow 共有 $(256/4) \times (256/4) \times 16$ 个 unit = $64 \times 64 \times 16$ 个
3. 同 1, 共有 $64 \times 64 \times 32$ 个 unit, $4 \times 4 \times 16 \times 32$ 个 weight
4. 共 $(64/4) \times (64/4) \times 32 = 16 \times 16 \times 32$ 个 unit
5. \therefore 单纯摊平而已
 \therefore 无 weight, unit 共 $16 \times 16 \times 32 = 8192$ 个
6. \therefore fully connected
 \therefore 共 $8192 \times 10 = 81920$ 个 weight
7. \therefore 单纯 output \therefore 无 weight

layer	unit	weight
1	$256 \times 256 \times 16$	$4 \times 4 \times 3 \times 16$
2	$64 \times 64 \times 16$	0
3	$64 \times 64 \times 32$	$4 \times 4 \times 16 \times 32$
4	$16 \times 16 \times 32$	0
5	8192	0
6	10	81920
7	10	0

#

三. 重要的CNN 与其贡献

1. LeNet (LeCun, 1998)

⇒ 奠定至今 CNN 之架构基础

2. AlexNet (2012)

⇒ 首次提出 dropout, 第一个使用 ReLU 的 CNN,

3. VGG (2014)

⇒ 解决了 style transfer, 也是第一个「深度」NN

4. ResNet (2016)

⇒ 解决了梯度消失问题

Ch13 Recurrent Neural Network

NO.

DATE

一. RNN的架構

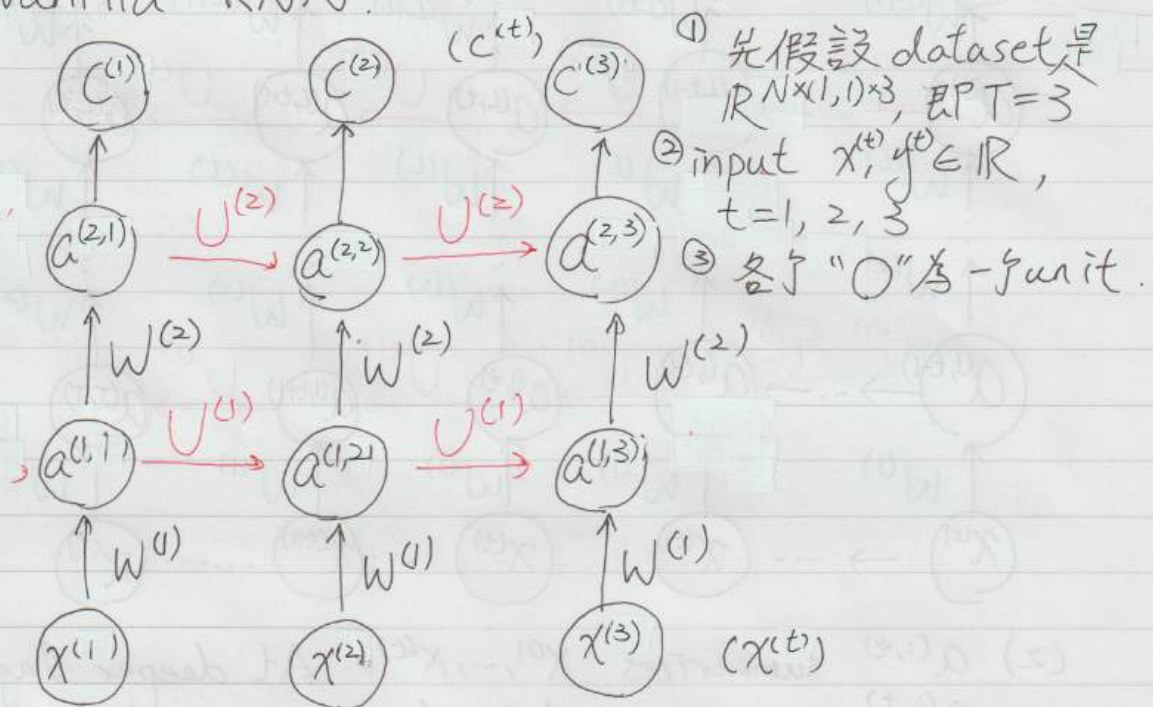
1. Sequential Data:

data point (x, y) 不一定是 i.i.d. ! 例如隨時間而變的 data: 語音, 文章, ... 此類 data set 記為 $X = \{X^{(n)}\}_n \subset \mathbb{R}^{N \times (D, K) \times T}$ ($D: x^{(n,t)}$ 之長度, $K: y^{(n,t)}$ 之長度)

(1) 每筆 data $X^{(n)} = \{(x^{(n,t)}, y^{(n,t)})\}$ 是一個 sequence, 常被簡寫為 $(x^{(t)}, y^{(t)})$

(2) T 稱為 horizon, $x^{(n)}, y^{(n)}$ 之 T 可能不同

2. Vanilla RNN.



(1) $\because y^{(t)}$ depends on $x^{(1)}, \dots, x^{(t)}$

我們希望拿到一些 data 中 time invariant 之 feature
 \Rightarrow 橫向, 直向的 weight 不隨時間而變!!

(2) 可知 $a^{(k,t)}$ ^{各 layer} = $\text{act}(z^{(k,t)})$
 $= \text{act}(U^{(k)} a^{(k,t-1)} + W^{(k)} a^{(k-1,t)})$

$k(\text{layer}) = 1, 2, 3$, $t(\text{time}) = 1, 2, 3$.

3. vanilla RNN: 考虑更一般的情形

每个 element 为 vector 对, 所形成长度为 T 的 sequence

$$\{ X = \{ X^{(n)} \}_n \in \mathbb{R}^{N \times (D, K) \times T}$$

$$\{ X^{(n)} = (\chi^{(n,t)}, y^{(n,t)})_t, \chi^{(n,t)} \in \mathbb{R}^D, y^{(n,t)} \in \mathbb{R}^T, X^{(n)} \in \mathbb{R}^{(D, K) \times T}$$

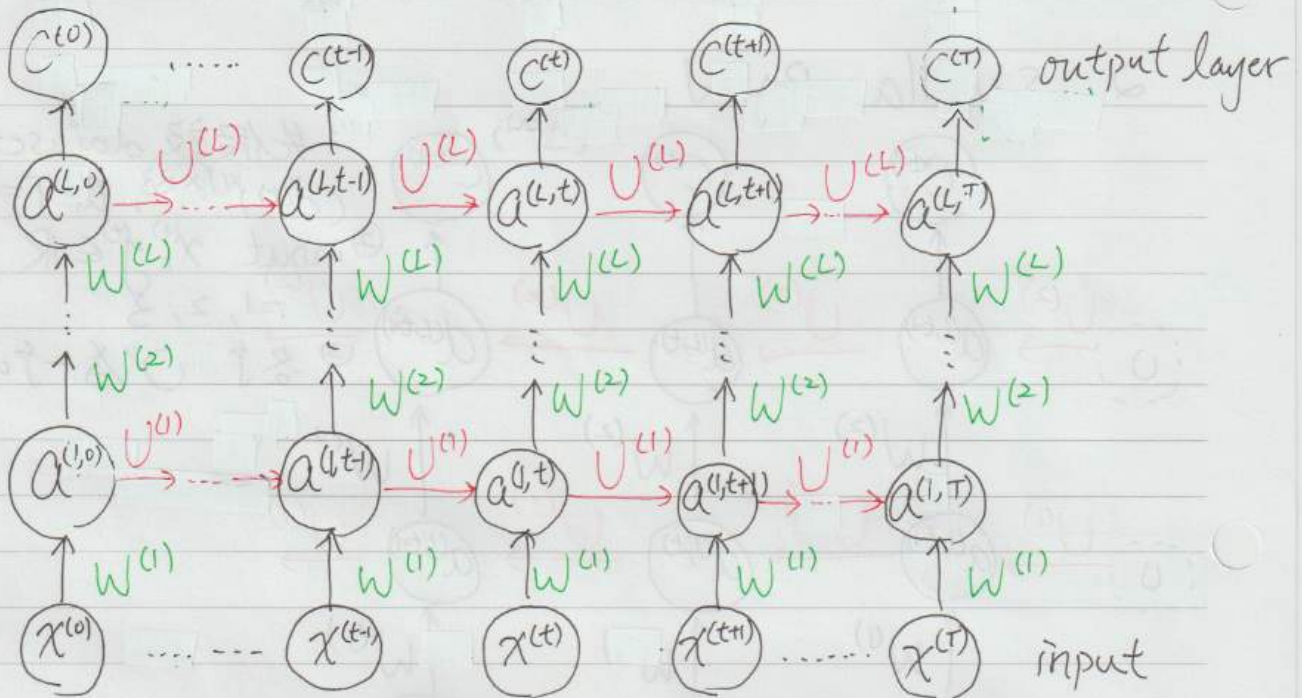
故对每一时刻 t, input $\chi^{(n,t)}$ 与 label $y^{(n,t)}$ 分别为 D 维及 K 维之 vector

↳ 简称为 $\chi^{(t)}, y^{(t)}$

⇒ model 如下:

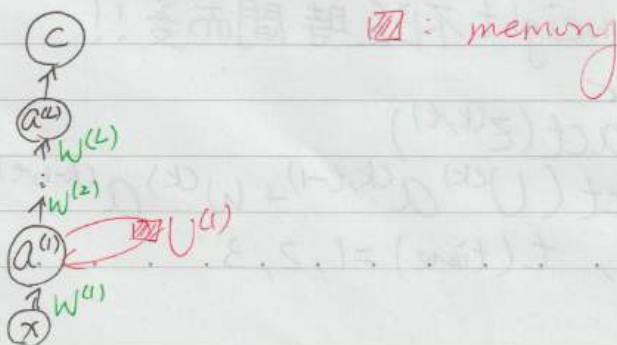
$$(1) a^{(k,t)} = \text{act}(z^{(k,t)}) = \text{act}(U^{(k)} a^{(k,t-1)} + W^{(k)} a^{(k-1,t)})$$

(bias term omitted)



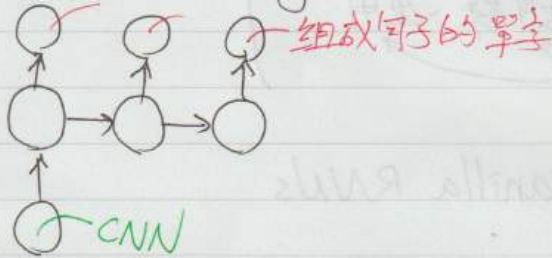
(2) $a^{(:,t)}$ summarizes $\chi^{(0)}, \dots, \chi^{(t)}$. At deeper layers, $a^{(:,t)}$ give more abstract summarizations

(3) 为方便起见, 我们可能会把上图之 RNN fold 成下图

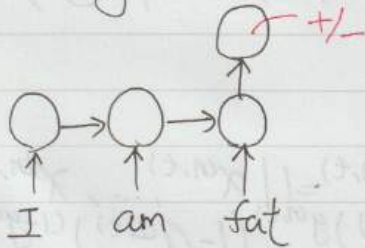


二. 不同种类的RNN

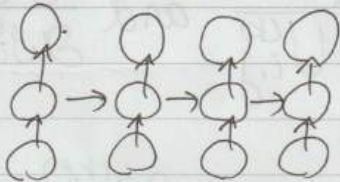
1. One to Many ex. Image Capturing



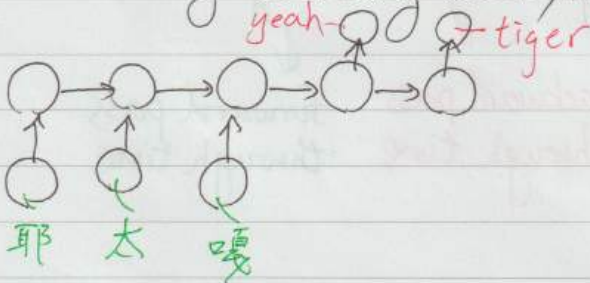
2. Many to one ex. Sentiment Analysis



3. Many 2 Many (Synced) ex. Language modeling



4. Many to Many (unsynched) ex. Machine Translation



三. RNN 之訓練

(思考 | 看到各種神經網路
→ 架構 · 變體 · 訓練 · 應用)

1. Cost function of vanilla RNNs by MLE,

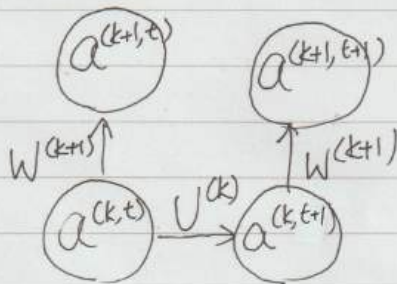
$$\begin{aligned} \operatorname{argmin}_{\Theta} C(\Theta) &= \operatorname{argmin}_{\Theta} -\log P(X|\Theta) \\ &= \operatorname{argmin}_{\Theta} -\sum_n \sum_t \log P(y^{(n,t)} | x^{(n,t)}, \dots, x^{(n,1)}, \Theta) \\ &= \operatorname{argmin}_{\Theta} -\sum_n \sum_t C^{(n,t)}(\Theta) \end{aligned}$$

ex. in binary classification, $P(y^{(n,t)} = 1 | x^{(n,t)}, \dots, x^{(n,1)}) \sim B(1, a^{(k,t)}) \Rightarrow C^{(n,t)}(\Theta) = (a^{(k,t)})^{y^{(n,t)}} (1 - a^{(k,t)})^{(1 - y^{(n,t)})}$

2. Using SGD, $\Theta^{(s+1)} \leftarrow \Theta^{(s)} - \eta \nabla_{\Theta} \sum_n \sum_t C^{(n,t)}(\Theta^{(s)})$
 \Rightarrow we need to evaluate $\frac{\partial C^{(n,t)}}{\partial U_{i,j}^{(k)}}$ and $\frac{\partial C^{(n,t)}}{\partial W_{i,j}^{(k)}}$

(1) $\frac{\partial C^{(n,t)}}{\partial W_{i,j}^{(k)}}$ \Rightarrow similar to NN

$$(2) \frac{\partial C^{(n,t)}}{\partial U_{i,j}^{(k)}} = \frac{\partial C^{(n,t)}}{\partial z_j^{(k,t)}} \cdot \frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}} \triangleq \delta_j^{(k,t)} \frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}}$$



backward pass
through time

forward pass
through time

3. Forward pass through time \Rightarrow 解 $\frac{\partial z_j^{(k,t)}}{\partial U_{ij}^{(k)}}$

$$\because z_j^{(k,t)} = \sum_i W_{ij}^{(k)} a_i^{(k,t-1)} + \sum_i U_{ij}^{(k)} a_i^{(k,t-1)}$$

$$\therefore \frac{\partial z_j^{(k,t)}}{\partial U_{ij}^{(k)}} = a_i^{(k,t-1)}$$

\Rightarrow 可從最 shallow (靠近 input) 以及時間最早的 unit 開始求解

4. Backward pass through time \Rightarrow 解 $\delta_j^{(k,t)} = \frac{\partial C^{(n,t)}}{\partial z_j^{(k,t)}}$

$$\because \delta_j^{(k,t)} = \frac{\partial C^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial C^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}}$$

$$= \frac{\partial C^{(n,t)}}{\partial a_j^{(k,t)}} \text{act}'(z_j^{(k,t)})$$

$$= \text{act}'(z_j^{(k,t)}) \left(\sum_s \frac{\partial C^{(n,t)}}{\partial z_s^{(k+1,t)}} \cdot \frac{\partial z_s^{(k+1,t)}}{\partial a_j^{(k,t)}} + \sum_s \frac{\partial C^{(n,t)}}{\partial z_s^{(k,t+1)}} \cdot \frac{\partial z_s^{(k,t+1)}}{\partial a_j^{(k,t)}} \right)$$

$$= \left(\sum_s \delta_s^{(k+1,t)} W_{j,s}^{(k+1)} + \sum_s \delta_s^{(k,t+1)} U_{j,s}^{(k)} \right) \text{act}'(z_j^{(k,t)})$$

\Rightarrow 可從最 deep (靠近 output) 及時間最晚之 unit 算回去以求 $\delta_j^{(k,t)}$

5. 組合 2~4, 我們可以得出 train RNN 的方法:

(1) 使用普通的 Backprop 求 $\frac{\partial C^{(n,t)}}{\partial W_{ij}^{(k)}}$

(2) 使用 **Backprop through time** 求 $\frac{\partial C^{(n,t)}}{\partial U_{ij}^{(k)}}$ (BPTT)

$\left\langle \begin{array}{l} \text{forward pass through time (single)} \rightarrow \text{求 } \frac{\partial z_j^{(k,t)}}{\partial U_{ij}^{(k)}} \\ \text{backward pass through time (multiple)} \rightarrow \text{求 } \delta_j^{(k,t)} \end{array} \right.$

四. 訓練RNN會遇到的問題&改善方法.

1. Exploding/Vanishing Gradient Problem

(1) 前言: 由前述 training 的討論, 可知 $a^{(k,c)}$ 與 $a^{(k,j)}$ 、 $\delta^{(k,c)}$ 與 $\delta^{(k,j)}$ 是直接與 $(U^{(k)})^{j-c}$ 有關. 更甚者, 假設 (w.l.o.g.) 忽略 activation function, 則

$$\begin{cases} a^{(k,j)} = (U^{(k)T})^{j-c} a^{(k,c)} \\ \delta^{(k,c)} = (U^{(k)})^{j-c} \delta^{(k,j)} \end{cases}$$

\Rightarrow 若 $a^{(k,c)}$, $a^{(k,j)}$ 相距太遠, 即 $j-c$ 這個時間很大, 會引發一些 optimization problem! 這了原由稱為 *long-term dependency*

(2). 若對 $U^{(k)}$ 做 eigendecomposition

$$\Rightarrow (U^{(k)})^{j-c} = Q \begin{pmatrix} \lambda_1^{j-c} & & \\ & \dots & \\ & & \lambda_{D^{(k)}}^{j-c} \end{pmatrix} Q^T$$

\Rightarrow When $j-c$ 很大, 則

$\left\{ \begin{array}{l} \lambda_s^{j-c} \text{ 趨大} \Rightarrow \text{Exploding gradient} \\ \lambda_s^{j-c} \text{ 趨小} \Rightarrow \text{vanishing gradient} \end{array} \right.$

2. Cost surface: 通常很陡, 跟懸崖一樣

\Rightarrow 解法:

(1) Nesterov momentum

(2) Gradient clipping: 若梯度超過一定值就把多出來的部份砍掉 \Rightarrow *effective in practice!*

(3) RMSProp

(3) 解決方案: Sigmoid activation function

$\left\{ \begin{array}{l} \text{bounded range in the forward pass} \\ \text{act}'(\cdot) < 1 \text{ in backward pass} \end{array} \right.$

但只能解決 $U^{(k)}$ 之 exploding 問題, vanishing

解决不了! 还会 introduce vanishing gradient of $W^{(k)}$

3. LSTM:

(1) idea: create **shortcut** in each neuron
 \Rightarrow error signals flows backward more smoothly

(2) 方法:

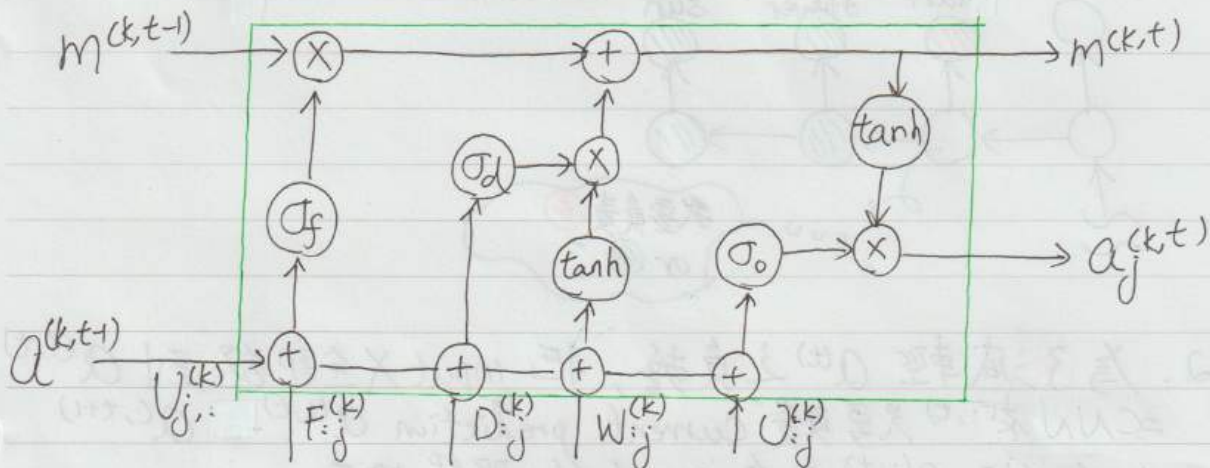
tanh activation

forget gate σ_f : 决定要不要忘掉 $m^{(k,t-1)}$

input gate σ_d : " 信号存第 j 个 activation 到 $m^{(k,t)}$

output gate σ_o : " output 第 j 个 activation

\downarrow unit of LSTM



(3) equation:

(Left as an exercise!)

(4) 优美: LSTM + gradient clipping 非常有效

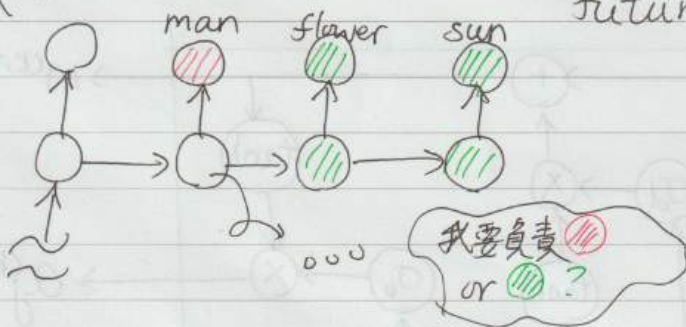
4. 平行化 train: 由於要 pass through time, RNN 的 training 是無法被平行化的。不過可藉由 output to hidden 或 teacher forcing 來替代

五. Attention — RNN 訓練的最佳助手

1. 前言: 在 RNN 中, 一個 hidden representation $a^{(l,t)}$ 需要負責 current prediction $a^{(l,t)}$ 以及 所有未來的 feature $a^{(l,t+1)}, \dots, a^{(l,T)}$

⇒ $a^{(l,t)}$ 的 trade off: 究竟它該

represent features for current prediction?
future?



2. 為了減輕 $a^{(l,t)}$ 之負擔, 把 input X 全部餵到 $a^{(l,t)}$
 ⇒ 每個 $a^{(l,t)}$ 只要負責 current prediction $a^{(l,t)}$ 與 $a^{(l,t+1)}$
 ⇒ 想了解 $a^{(l,t)}$ 在看 input 的哪些地方

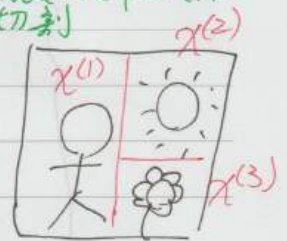
3. 方法: attention mechanism

(1) 假設 input 可被切成不同 part $\{x^{(i)}\}$

(2) 令 attention vector 為 $b^{(t)}$ 是由 $a^{(l,t-1)}$ 計算而來 (方法後述)

(3) feed $a^{(l,t)}$ with the weighted input $c^{(t)} = \sum_i b_i^{(t)} x^{(i)}$, 若我們實際畫出 $c^{(t)}$, 會發現對應之 $x^{(i)}$ 變得較明亮
 ⇒ attention! (PPT p50).

實務上會用 CNN 的 filter 之 activation part 來做切割



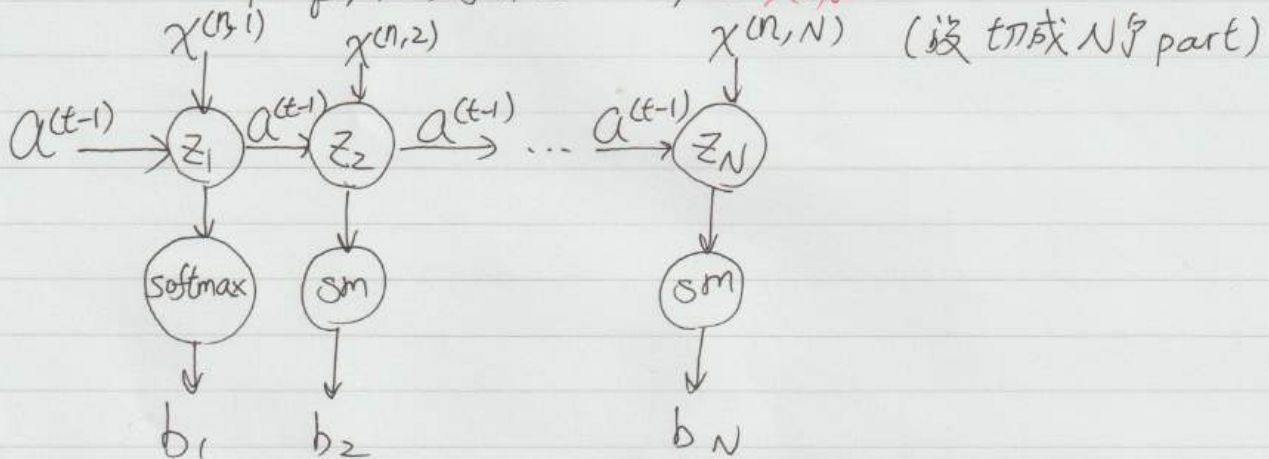
4. Computing the attention vector:

(1) 对每个 part $x^{(i)}$, 以 $a^{(L-1, t-1)}$ 来算分数 (match score). 算分的方法可以用 NN:

$$z_i = \text{act}(p^T a^{(L-1, t-1)} + q^T x^{(i)} + r)$$

* { jointly trained with the main RNN

p, q, r 对不同的 i, t 共用



(2) 接著用 softmax normalize $\{z_i\}$

$$b_i = \text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

Ch14 Unsupervised Learning

一. Unsupervised learning 的觀念

1. Data mining: Dataset $X = \{x^{(i)}\}_i$, $x^{(i)} \sim \text{i.i.d. } x$
 \Rightarrow no labels! no supervision

ex. K-means, dimension reduction (PCA), generative

2. Self-supervised learning

(1) 目標: learn 一个 model 使其可以 "fill the blank"

(2) 手法: 假設 X 是一个高维空间裡的 manifold
 \Rightarrow 把这个 manifold 学起来

ex. Mnist dataset. $\mathbb{R}^{28 \times 28} = \mathbb{R}^{784}$, 但数字只存在於 \mathbb{R}^{784}
 中一个很小的 subspace 上

\Rightarrow learn 每个 data 之 n.b.h. 的 tangent space 的基底

(3) 重要範例:

A. Word2vec \Rightarrow 将一句话中每个字的 semantic 以
 一向量 h 来表示

B. Doc2vec \Rightarrow 将文件中每一段给予一特别之 ID,
 並以整篇文章为 input

\Rightarrow 学到没有用文字叙述的 context (上下文关系)

C. Filling Image \Rightarrow 将照片一块挖掉並要 model
 補上

\because 如果把挖掉的区域当成 label y , 則 $\dim(y)$ 太大

\therefore 可用 RNN 解

\Rightarrow 一个 pixel 只 depend 附近之 label

\Rightarrow 比 RNN 是可 train 的!

Def. A manifold is a topological space that are
 linear locally

⇒ 相鄰之間可由 linear transform 近似，但太遠就不行！

Remark. 對每個 manifold 上的點而言，其 tangent space 為 tangent vectors 的 span

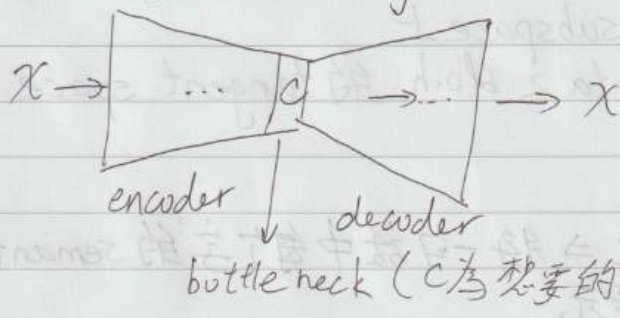
二. Autoencoders

1. 由一個 encoder 與 decoder 組成，目標為輸入 x 給 encoder 後使 decoder 能 regenerate x

(1) 可當成 label 為 x 自身之 task

(2) encoder 尾端與 decoder 前端應為 bottleneck，否則不會把「有用的東西」embed 進去 C

⇒ 會得到 trivially identical mapping



(3) Cost function:

$$\arg \min_{\Theta} -\log P(X | \Theta) = \arg \min_{\Theta} -\sum_{i=1}^N \log P(x^{(i)} | \Theta)$$

(4) 第 j 層 sigmoid output unit $a_j^L = \hat{p}_j$, $\forall x_j \sim B(1, p_j)$
L 層之第

$$\Rightarrow P(x_j^{(i)} | \Theta) = (a_j^{(L)})^{x_j^{(i)}} (1 - a_j^{(L)})^{(1 - x_j^{(i)})}$$

(5) 第 L 層之 output $a^{(L)} \triangleq \hat{\mu}$, $x \sim \mathcal{N}(\mu, \Sigma)$
↳ μ 之估計

$$\Rightarrow -\log P(x^{(i)} | \Theta) = \|x^{(i)} - a^{(L)}\|^2$$

ex. 以 MNIST train - autoencoder, 並將各數字學到的 embedding 稱為 C_0, C_1, \dots, C_9

⇒ 將 train 好的 decoder 拆下，並以數字 1 的某 data $x_1^{(i)}$ 的 embedding $C_1^{(i)}$ 為 decoder 之輸入，預期為：

$$\text{decoder}(C_1^{(i)}) = \boxed{1} = x_1^{(i)}$$

同理，以數字 2 的某 data $x_2^{(j)}$ (embedding $C_2^{(j)}$)

$$\text{decoder}(C_2^{(j)}) = \boxed{2} = x_2^{(j)}$$

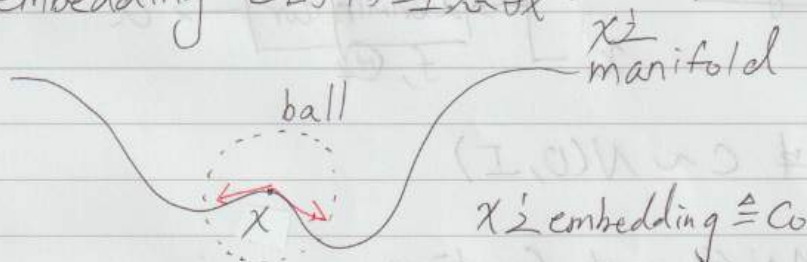
若以 $\frac{C_1^{(i)} + C_2^{(j)}}{2}$ 為 decoder 之 input，竟得到

$$\text{decoder}\left(\frac{C_1^{(i)} + C_2^{(j)}}{2}\right) = \boxed{?}$$

⇒ 介於 $\boxed{1}$ - $\boxed{2}$ 之間的圖形

⇒ 此 autoencoder 確實拿到了 manifold!

2. embedding C 的物理意義：



設 $x \in$ 某 manifold，則希望 input 在該 ball 內的話，
encoder 之 output 為 C_0

⇒ C 為 manifold 之 Δ 一科 coordinate

⇒ 可 regularize on 「local change of x 對應到的 change of C 」

⇒ 以 Jacobian $J(x) = \frac{\partial C}{\partial x}$ 來建模，regularize on

$$\Omega(C) \triangleq \sum_{i=1}^N \left\| \frac{\partial C^{(i)}}{\partial x^{(i)}} \right\|_F^2$$

⇒ 對 $J(x)$ 做 SVD，得 $J(x) = UDV^T$ ，則可以最大的 singular values in D 所對應之 V 的 row 為 manifold 在 x 之 tangent vector

⇒ 此 tangent vector 即能使 C 變化最大的方向，且延此方向稍微移動仍在 manifold 上

⇒ 將 tangent space 聯繫起來，可視為 manifold 之近似

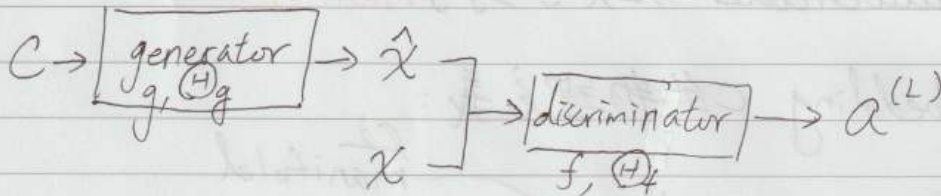
三. Generative Adversarial Networks

1. 目標: generate data points from **random codes**

2. 構造:

(1) Generator g : 生成 image

(2) Discriminator: 分類 real image 和 generator 生成之 image 之不同



其中 $c \sim N(0, I)$

3. GAN 的 cost function:

$\arg \min_{\theta_g} \max_{\theta_f} \log P(X | \theta_g, \theta_f)$ ⇒ 使 X 盡可能 $\approx X$, 令 discriminator 分類不出來

$$= \arg \min_{\theta_g} \max_{\theta_f} \sum_{i=1}^N \log f(x^{(i)}) + \sum_{j=1}^N \log(1 - f(g(c^{(j)})))$$

$$= \arg \min_{\theta_g} \max_{\theta_f} \sum_{i=1}^N \log \hat{p}^{(i)} + \sum_{j=1}^N \log(1 - \hat{p}^{(j)})$$

⇒ $\begin{cases} \hat{p}^{(i)} \text{ depends on } \theta_f \text{ only} \\ \hat{p}^{(j)} \text{ " both } \theta_f \text{ and } \theta_g \end{cases}$

4. 如何 train GAN?

Algorithm. train GAN by SGD

1. Initialize Θ_g, Θ_f

2. At each iteration:

① 对 f overfit
② " g update 是错的

(1) 将 Θ_g 固定, 並重複 K 次以下步驟:

(i) Sample N real point $\{x^{(i)}\}_i$ from X

(ii) " $c \sim \mathcal{N}(0, 1)$

(iii) update Θ_f

$$\Theta_f \leftarrow \Theta_f + \eta \nabla_{\Theta_f} \left(\sum_{i=1}^N \log f(x^{(i)}) + \sum_{j=1}^N \log(1 - f(g(c^{(j)}))) \right)$$

(2) 将 Θ_f 固定, 執行一次以下步驟:

(i) Sample N $c \sim \mathcal{N}(0, 1)$

(ii) update Θ_g

$$\Theta_g \leftarrow \Theta_g - \eta \nabla_{\Theta_g} \left(\sum_{i=1}^N \log(1 - f(g(c^{(i)}))) \right)$$

5. GAN 很難 train!

(1) 很難 converge: "The update of Θ_g and Θ_f may cancel each other's progress"

(2) mode collapsing \Rightarrow train 不起來

$\because K$ 很小時, alternate SGD 不知道 $\min_{\Theta_g} \max_{\Theta_f}$ 之差別

$\therefore g$ 可能生成「騙得過 f , 但騙不過人類」的 output

\Rightarrow solutions: $\left\{ \begin{array}{l} \text{minibatch discrimination} \\ \text{unrolled GANs} \end{array} \right.$

(3). Balance between g and f :

Too large K : f overfit data

\Rightarrow g update for "wrong" target f

$\Rightarrow \nabla_{\Theta_g}$ vanishes

Too small K : g updated for "meaningless" f

\Rightarrow solution: Wasserstein GAN

6. Information theory & GAN

Thm. Let P_{data} be the distribution of x , P_g be the distribution of $\hat{x} = g(c)$. Then the max term in the cost function of GAN actually measures the JS divergence $\triangleq C^*$

$$D_{JS}(P_{data} \parallel P_g) = \frac{1}{2} D_{KL}(P_{data} \parallel Q) + \frac{1}{2} D_{KL}(P_g \parallel Q),$$

where $Q = \frac{1}{2}(P_g + P_{data})$

lemma 1. For each x , assuming that f has infinite capacity, then we can find f maximizing $P_{data}(x) \log f(x) + P_g(x) \log(1-f(x))$ to have C^* — (1)

$$\begin{aligned} \text{pf. } C^* &= \max_{\mathcal{H}_f} \left(\frac{1}{N} \sum_{i=1}^N \log f(x^{(i)}) + \frac{1}{N} \sum_{j=1}^N \log(1-f(g(c^{(j)}))) \right) \\ &= \max_{\mathcal{H}_f} \left(\frac{1}{N} \sum_{i=1}^N \log f(x^{(i)}) + \frac{1}{N} \sum_{j=1}^N \log(1-f(\hat{x}^{(j)})) \right) \\ &\triangleq \max_{\mathcal{H}_f} \left(E_{x \sim P_{data}} [\log f(x)] + E_{x \sim P_g} [\log(1-f(x))] \right) \\ &= \max_{\mathcal{H}_f} \int_x P_{data}(x) \log f(x) dx + \int_x P_g(x) \log(1-f(x)) dx \\ &= \max_{\mathcal{H}_f} \int_x \left(P_{data}(x) \log f(x) + P_g(x) \log(1-f(x)) \right) dx \end{aligned}$$

Q.E.D.

lemma 2. The f maximize (1) is

$$f^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \in [0, 1]$$

pf. differentiate (1) w.r.t. x .

pf of Thm: By lemma 1 and 2,

$$\begin{aligned}
 C^* &= \max_{\mathbb{H}_f} \int_{\mathcal{X}} (P_{\text{data}}(x) \log f(x) + P_g(x) \log(1-f(x))) dx \\
 &= \int_{\mathcal{X}} P_{\text{data}}(x) \log \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} dx + \int_{\mathcal{X}} P_g(x) \log \left(1 - \frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx \\
 &= -2 \log 2 + \int_{\mathcal{X}} P_{\text{data}}(x) \log \frac{P_{\text{data}}(x)}{(P_{\text{data}}(x) + P_g(x))/2} dx \\
 &\quad + \int_{\mathcal{X}} P_g(x) \log \left(1 - \frac{P_g(x)}{(P_{\text{data}}(x) + P_g(x))/2}\right) dx \\
 &= -2 \log 2 + 2D_{\text{JS}}(P_{\text{data}} \parallel P_g)
 \end{aligned}$$

Q.E.D.

Cor. The cost function of GAN is

$$\arg \min_{\mathbb{H}_g} D_{\text{JS}}(P_{\text{data}} \parallel P_g)$$

⇒ 若 data 维数高, $x, g(z) \in$ low dimensional manifold,
 则 $\because P_g$ and P_{data} 只在 manifold 上定义
 $\therefore \{x \mid P_g(x) \neq 0 \text{ and } P_{\text{data}}(x) \neq 0\}$ 可能很小
 ⇒ GAN 难 train!

四. GAN 之变体

1. Wasserstein GAN (W-GAN)

⇒ 以 modify 的 cost function 輔以 training 時的 gradient clip 來找 $P_g(x)$ 和 $P_{\text{data}}(x)$ 皆 $\neq 0$ 之集合

2. Improved W-GAN

⇒ 加了 gradient penalty 在 W-GAN 的 cost function 上, 使 convergence 更快

3. CapsuleNet \Rightarrow discriminator detects the existence of patterns more than their relative positions

\Rightarrow 繞生成的 image 不要那麼「怪」

4. DC-GAN \Rightarrow 找一 class 畫出來的東西包含事先畫好的几筆

5. Conditional GAN \rightarrow text to image synthesis \rightarrow 已被 diffusion model 取代

6. 其他 GAN 可做的事:

(1) 增加 image 之 resolution

(2) image to image translation. ex. 空拍圖 \rightarrow 地圖

五. Concluding Remarks: Ch12 ~ Ch14 所介紹的 NN, 在今日 (2024) 几乎都被更先進的 model 取代了。

Past & current 之比較 \rightarrow (voice/music)

model \ domain	sound	image	text
CNN	transformer	transformer	無
RNN	少見 ^(*)	diffusion model	LLM
GAN	LLM ^(**)	diffusion model	LLM

*: 通常 sound 都用 CNN 以時頻譜做訓練

** : text to music

今日這些更 powerful 的商用 model 在不公開其架構

\Rightarrow 懂得如何操作 (prompt engineering) 更重要!

Ch16 Reinforcement Learning

NO.

DATE

一. 何謂強化學習

Def. Given an agent who sees **states** $s^{(t)}$, take actions $a^{(t)}$, and receives **reward** $R^{(t)}$, a **reinforcement learning** algorithm learns the best **policy** $\pi^*(s^{(t)}) = a^{(t)}$ that maximizes the total reward $\sum_{t} R^{(t)}$

1. Environment: 定義 S 為 environment 之集合, $s^{(t)}$ 為時刻 t agent 所在之 state.
* S 不隨時間改變!

2. Reward $R^{(t)}$ 取決於 $s^{(t+1)}, s^{(t)}, \dots$ 以及 $a^{(t)}, a^{(t-1)}, \dots$
把所有 action 所形成之集合稱為 **action space** A .

ex. 機器人走迷宮遊戲

$S = \{(1,1), (1,2), (1,3), (2,1), (2,3), \dots, (4,3)\}$

$A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$

若機器人一開始在 $(3,1)$, 接著依序經過 $(3,2), (3,3), (4,3)$, 則:

$$\begin{cases} s^{(0)} = (3,1), a^{(0)} = \uparrow, R^{(0)} = 0 \\ s^{(1)} = (3,2), a^{(1)} = \uparrow, R^{(1)} = 0 \\ s^{(2)} = (3,3), a^{(2)} = \rightarrow, R^{(2)} = 0 \\ s^{(3)} = (4,3), R^{(3)} = 1 \end{cases}$$

3				+1
2		///		-1
1			*	
	1	2	3	4



二. Markov Process

Def. A **random process** is a collection of time-indexed random variables $\{S^{(t)}\}_t$

Def. A **Markov process** $\{S^{(t)}\}$ is a random process satisfy the **Markov property**:

$$P(S^{(t+1)} | S^{(t)}, S^{(t-1)}, \dots) = P(S^{(t+1)} | S^{(t)})$$

Def. A **Markov decision process** (MDP) is defined by

- (1) the state space S and the action space A
- (2) start state $s^{(0)}$
- (3) The **transition distribution** $P(s' | s; a)$ (fixed over time t)
- (4) the deterministic reward function $R(s, a, s') \in \mathbb{R}$
- (5) A **discount factor** $\gamma \in [0, 1]$
- (6) The horizon $H \in \mathbb{N}$ (can be infinite)

ex. 机器人迷宫游戏中, 可定义 $P((3, 2) | (3, 1); \uparrow) = 0.8$ (马车可能失灵), $R((4, 1), \uparrow, (4, 2)) = -1$

(7). Given a policy $\pi(s) = a$, an MDP proceeds as

$$s^{(0)} \xrightarrow{a^{(0)}} s^{(1)} \xrightarrow{a^{(1)}} \dots \xrightarrow{a^{(H-1)}} s^{(H)}$$

(8) The accumulative reward

$$= R(s^{(0)}, a^{(0)}, s^{(1)}) + \gamma R(s^{(1)}, a^{(1)}, s^{(2)}) + \dots + \gamma^{H-1} R(s^{(H-1)}, a^{(H-1)}, s^{(H)})$$

Remark. $\because 0 < \gamma < 1 \therefore$ 越后面的 reward 越小. γ 决定该算法要前瞻 or 短视, 且 γ 为 hyperparameter

Def. Given a policy π , the expected accumulative reward collected by taking actions following π can be expressed by

$$V_{\pi} = E_{S^{(0)}, \dots, S^{(H)}} \left[\sum_{t=0}^H \gamma^t R(S^{(t)}, \pi(S^{(t)}), S^{(t+1)}) \mid \pi \right]$$

Goal: find the **optimal policy** $\pi^* = \arg \max_{\pi} V_{\pi}$

≡. Value iteration

Def. The **optimal value function** is the maximum expected accumulative reward when starting from state s and acting optimally for h steps

$$V^{*(h)}(s) = \max_{\pi} E_{S^{(0)}, \dots, S^{(h)}} \left[\sum_{t=0}^h \gamma^t R(S^{(t)}, \pi(S^{(t)}), S^{(t+1)}) \mid S^{(0)} = s; \pi \right]$$

1. Value iteration of finite horizon

Having $V^{*(h-1)}(s)$ for each s , we can solve π^* by

$$\pi^* = \arg \max_{\pi} E_{S^{(0)}, \dots, S^{(H)}} \left[\sum_{t=0}^H \gamma^t R(S^{(t)}, \pi(S^{(t)}), S^{(t+1)}) \mid \pi \right]$$

$$= \arg \max_a \sum_{S'} P(S' | S; a) \left[R(S, a, S') + \gamma V^{*(h-1)}(S') \right], \forall S.$$

於是有了 recurrence \Rightarrow 必有 sub-optimal

\Rightarrow 想到 **dynamic programming**

$$\Rightarrow \begin{cases} h=H-1, & V^{*(H-1)}(s) = \dots \\ h=H-2, & V^{*(H-2)}(s) = \dots \\ \vdots & \vdots \\ h=0, & V^{*(0)}(s) = \dots \\ h=-1, & V^{*(-1)}(s) = 0 \end{cases}$$

故可寫下下列 sudo code:

Algorithm: Value iteration (Finite Horizon)

Input: S, A, P, R, γ, H

step 1. $\forall s \in S$, initialize $V^*(s) \leftarrow 0$

step 2. For $h=0$ to $H-1$

For $s \in S$

$$V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s;a) [R(s,a,s') + \gamma V^*(s')]$$

step 3. For $s \in S$ do

$$\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s;a) [R(s,a,s') + \gamma V^*(s')]$$

step 4. Return $\pi^*(s), \forall s \in S$

ex. 机器人走迷宫游戏中, 设 $P(s'|s;a)$ 的 noise 为 0.2, 即指令为 "↑" 时, 真的往上之概率为 0.8, 但往左 or 右之概率各为 0.1; 往下之概率为 0, 依此类推. 若任一方向有障碍物, 则「留在原地」之概率为这些有障碍之方向的概率和. 给定 $\gamma=0.9$, 求 $h=2$ 时, $V^{*(2)}(3,3)=?$

-sol- 当 $h=0$, 只有 $(2,3), (3,2), (3,3), (4,3)$ 会影响

$$\Rightarrow V^{*(1)}(3,3) = \max \left\{ \begin{array}{l} \uparrow: 0.9 \times (0.8 \times 0 + 0 \times 0 + 0.1 \times 0 + 0.1 \times 1) \\ \downarrow: 0.9 \times (0 \times 0 + 0.8 \times 0 + 0.1 \times 0 + 0.1 \times 1) \\ \rightarrow: 0.9 \times (0.1 \times 0 + 0.1 \times 0 + 0 \times 0 + 0.8 \times 1) \\ \leftarrow: 0.9 \times (0.1 \times 0 + 0.1 \times 0 + 0 \times 0 + 0.8 \times 1) \end{array} \right\}$$

$$= 0.72$$

\Rightarrow 当 $h=2$, 有

$$V^{*(2)}(3,3) = \max \left\{ \begin{array}{l} \uparrow: 0.9 \times (0.8 \times 0.72 + 0 \times 0 + 0.1 \times 0 + 0.1 \times 1) \\ \downarrow: 0.9 \times (0 \times 0.72 + 0.8 \times 0 + 0.1 \times 0 + 0.1 \times 1) \\ \leftarrow: 0.9 \times (0.1 \times 0.72 \end{array} \right.$$

理論上, 6步以內就可以 cover 所有格子, 但也可以把 H 設到 100, PPT p17 是 $H=100$ 之結果

2. Bellman optimality equation:

$$\text{當 } h \rightarrow \infty, V^*(s) = \max_a \sum_{s'} P(s'|s;a) [R(s,a,s') + \gamma V^*(s')], \forall s$$

由此可知, optimal policy $\pi^*(s) = \arg \max_a V^*(s)$, $\forall s$ 在 $h \rightarrow \infty$ 時

- (1) *stationary*: optimal action 傾向不再變化
- (2) *memoryless*: 與 $S^{(0)}$ 無美

Algorithm: Value Iteration (Infinite Horizon)

Input: $S, A, P, R, \gamma, H = \infty$

step 1. $\forall s \in S, V^*(s) \leftarrow 0$

step 2. Do

For $s \in S$

$$V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s;a) [R(s,a,s') + \gamma V^*(s')]$$

until $V^*(s)$ converge

step 3. For $s \in S$

$$\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s;a) [R(s,a,s') + \gamma V^*(s')]$$

step 4. Return $\pi^*(s), \forall s \in S$.

3. The convergence of infinite horizon.

Thm. When $H \rightarrow \infty$, Value iteration converges and gives the optimal policy π^*

pf. Note that the reward can be negative

$$\Rightarrow V^*(s) - V^{*(H)}(s) = \gamma^{H+1} R(s^{(H+1)}, a^{(H+1)}, s^{(H+2)}) \\ + \gamma^{H+2} R(s^{(H+2)}, a^{(H+2)}, s^{(H+3)}) \\ + \dots$$

$$\leq \gamma^{H+1} |R(s^{(H+1)}, a^{(H+1)}, s^{(H+2)})| + \dots + \gamma^Q |R(s^{(Q)}, a^{(Q)}, s^{(Q+1)})|$$

Define $R_{\max} = \max \{ |R(s^{(i)}, a^{(i)}, s^{(i+1)})| : i \in N \}$

$$\Rightarrow V^*(s) - V^{*(H)}(s) \leq \gamma^{H+1} R_{\max} + \gamma^{H+2} R_{\max} + \dots$$

$$= \frac{\gamma^{H+1}}{1-\gamma} R_{\max}$$

$$\Rightarrow \lim_{H \rightarrow \infty} V^*(s) - V^{*(H)}(s) = 0$$

$$\Rightarrow V^*(s) = \lim_{H \rightarrow \infty} V^{*(H)}(s)$$

喜長遠目標

比較會坏 \Rightarrow 喜保險

Q.E.D.

4. γ 值的選擇: (PPT p21)

$$\left\{ \begin{array}{l} \gamma \text{ 大, noise 大} \rightarrow \text{prefer distant exit, avoiding risk.} \\ \gamma \text{ 大, noise 小} \rightarrow \text{" distant " , risking.} \\ \gamma \text{ 小, noise 大} \rightarrow \text{" close " , avoiding risk.} \\ \gamma \text{ 小, noise 小} \rightarrow \text{" close " , risking.} \end{array} \right.$$

短視
 \Rightarrow 喜較近目標

比較不會坏
 \Rightarrow 喜犯難

IV. Policy iteration

Def. For a policy π , we define its value function

$$V_{\pi}(s) = E_{s^{(0)}, \dots} \left[\sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) \mid s^{(0)} = s, \pi \right]$$

即照著 π 行事所得的 expected accumulated reward.

1. 如果要用這種方式去找 $\pi^* = \arg \max_{\pi} V_{\pi}$, 那演算法理應長得像:

Input: $S, A, P, R, \gamma, H \rightarrow \infty$

step 1. $\forall s \in S$, initialize $\pi(s)$ randomly

step 2. Do

Evaluate $V_{\pi}(s)$, $\forall s$

Improve π s.t. $\forall s \in S$, $V_{\pi}(s)$ becomes better until $\pi(s)$ converge

step 3. Return $\pi(s)$, $\forall s \in S$.

2. 考慮 Bellman expectation equation (policy iteration 版)

$$V_{\pi}(s) = \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')], \forall s \in S$$

(infinite!)

故, 要計算 $V_{\pi}(s)$ 有 2 種方式:

(1) $V_{\pi}(s)$ 是 S 之方程式, 又 $s \in S$, 故解 $|S|$ 條聯立方程式 ($O(|S|^3)$)

(2) 用 dynamic programming: (numerical friendly)
initialize $V_{\pi}(s) = 0, \forall s \in S$

$$V_{\pi}(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')], \forall s$$

3. How to improve π .

\Rightarrow We want to find $\hat{\pi}$ s.t. $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$.

Prup. We should update π by the rule

$$\hat{\pi}(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_{\pi}(s')]$$

下一步往哪更好?

pf. Consider $\forall s \in S$,

$$V_{\pi}(s) = \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')]$$

$$\leq \sum_{s'} P(s'|s; \hat{\pi}(s)) [R(s, \hat{\pi}(s), s') + \gamma V_{\pi}(s')]$$

$$\leq \sum_{s'} P(s'|s; \hat{\pi}(s)) [R(s, \hat{\pi}(s), s') + \gamma \sum_{s''} P(s''|s'; \hat{\pi}(s')) [R(s', \hat{\pi}(s'), s'') + \gamma V_{\pi}(s'')]]$$

$\leq \dots$

$$\leq E_{s', s'', \dots} [R(s, \hat{\pi}(s), s') + \gamma R(s', \hat{\pi}(s'), s'') + \dots \mid s^{(0)} = s; \hat{\pi}]$$

$$= V_{\hat{\pi}}(s)$$

Q.E.D.

4. 由以上討論, 可寫出 Policy Iteration 演算法:

Input: $S, A, P, R, \gamma, H \rightarrow \infty$

step 1. Initialize $\pi(s)$, $\forall s \in S$

step 2. Do

For $s \in S$, $V_{\pi}(s) \leftarrow 0$

Do

For $s \in S$

$$V_{\pi}(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')]$$

until $V_{\pi}(s)$ converge

For $s \in S$

$$\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_{\pi}(s')]$$

until $\pi(s)$ converge

step 3. Return $\pi(s)$.

Thm. When $H \rightarrow \infty$, policy iteration converges and gives the optimal policy π^* .

五. Reinforcement Learning 的概念

1. 剛討論的為 MDP \Rightarrow

(1) 可能無法 model environment as an MDP

$\Rightarrow P(s'|s;a)$ 可能未知!

(2) Reward $R(s,a,s')$ 可能未知!

2. 對策:

(1) Explore: perform actions randomly

\Rightarrow estimate $P(s'|s;a)$ and $R(s,a,s')$
by collecting samples

(2) Exploit: learn $P(s'|s;a)$ and $R(s,a,s')$

$\Rightarrow \pi^*$ can be computed using value/policy iteration

3. Model-based 與 Model free RL 之差異

假設要估算 $E[f(x)|y] = \sum_x P(x|y)f(x)$ given samples

$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots$, 要怎麼做呢?

(1) Model-based 的方法會直接 estimate

$$\hat{P}(x|y) = \frac{\#(x^{(i)}=x \text{ and } y^{(i)}=y)}{\#(y^{(i)}=y)}$$

$$\text{接著算 } \hat{E}[f(x)|y] = \sum_x \hat{P}(x|y)f(x)$$

\Rightarrow 先把要用到的參數準備好

(2) Model-free \Rightarrow 直接算 $\hat{E}[f(x)|y]$

$$\hat{E}[f(x)|y] = \frac{1}{\#(y^{(i)}=y)} \sum_{i \in \{i | x^{(i)}=x, y^{(i)}=y\}} f(x^{(i)})$$

\Rightarrow 由 sampling theory, 會 work!

4. Model-free RL 之挑戰

(1) Value iteration 的 Algorithm 中 $V^*(s)$ 的 update

$$\Rightarrow \textcircled{1} V^*(s) = \max_{\pi} E \left[\sum_t \gamma^t R^{(t)} \mid s^{(0)} = s; \pi \right]$$

除非真的 sample 到 true max, 否則無法 estimate!

$$\textcircled{2} \pi^*(s) \leftarrow \operatorname{argmin}_a \sum_{s'} P(s'|s;a) [R(s,a,s') + \gamma V_{\pi}(s)]$$

→ need model to solve

(2) Policy iteration:

① 和 (1) Value iteration 相比, $V_{\pi}(s)$ 能以 MC estimation 去估計

② $\pi(s)$ 仍需 model 去 solve

六. Model-based RL

1. 演算法:

Algorithm. Model-based RL using MC estimation

1. Use some **exploration policy** π' to perform one or more episodes/trails

→ Each sample records samples of $P(s'|s;a)$, $R(s,a,s')$ from start to terminal state

$$(1) s^{(0)} \xrightarrow{\pi'(s^{(0)})} s^{(1)} \xrightarrow{\pi'(s^{(1)})} \dots \xrightarrow{\pi'(s^{(H-1)})} s^{(H)}$$

$$(2) R(s^{(0)}, \pi'(s^{(0)}), s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, \pi'(s^{(H-1)}), s^{(H)})$$

2. Estimate $P(s'|s;a)$ and $R(s,a,s')$

$$(1) \hat{P}(s'|s;a) = \frac{\# \text{ times the action } a \text{ takes state } s \text{ to } s'}{\# \text{ times action } a \text{ is taken in state } s}$$

(2) $R(s,a,s') = \text{Avg. of reward values received}$
when a takes s to s'

3. Update the exploitation policy π

4. Repeat 1-3, but gradually mix π into π'

2. 困難矣:

(1) $P(s'|s,a)$ and $R(s,a,s')$ 要估的量可能很多!
 $\rightarrow s, a$, 甚至是連續的

(2) 假如 $P(s'|s,a)$ 太佳 \Rightarrow 難被 sample 到
 \Rightarrow 可能找不到最好的 policy

\Rightarrow 若 $R(s,a,s')$ depends on s , 則 $R(s,a,s')$ 可能估的很爛!

3. 實務上 model-based 很難 work, 故常直接以 model-free 之方法去算 $V^*(s)$ or $V_\pi(s)$.

t. SARSA 演算法. (model free)

1. Q function:

Def. Q function for π :

$$Q_\pi(s,a) = E_{s^0, \dots} [R(s,a,s^{(1)}) + \sum_{t=1}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)})]; s, a, \pi]$$

s.t. $V_\pi(s) = Q_\pi(s, \pi(s))$ with recurrence

$$Q_\pi(s,a) = \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma Q_\pi(s', \pi(s'))],$$

$\forall s, a$

Algorithm. Policy iteration based on Q_π

Input: MPP($S, A, P, R, \gamma, H \rightarrow \infty$)

Output: $\pi(s)$'s, $\forall s$

1. For each state s , initialize $\pi(s)$ randomly

2. Repeat until $\pi(s)$ converge

(1) Initialize $Q_\pi(s,a) = 0, \forall s, a$

(2) repeat until $Q_\pi(s,a)$ is converged:

$$Q_\pi(s,a) \leftarrow \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma Q_\pi(s', \pi(s'))]$$

(3) $\pi(s) \leftarrow \arg \max_a Q_\pi(s,a)$

(1) 此 Algorithm 之好處:

① 2-(2) 中, 可用 MC estimation 估 $Q_{\pi}(s, a)$

$$Q_{\pi}(s, a) = E\left[R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)}\right]$$

② 2-(3) 中, $\pi(s)$ 之 update 不需要 model.

(2) 缺點: 花費很多 trials, π 才 improve - 衰² \Rightarrow 太慢!

2. Temporal Difference Estimation: (TDE)

Algorithm. TDE of $Q_{\pi}(s, a)$:

① $\hat{Q}_{\pi}(s, a) \leftarrow$ random value, $\forall s, a$

② Repeat until conv., $\forall a^{(t)}$

$$\hat{Q}_{\pi}(s^{(t)}, a^{(t)}) \leftarrow \hat{Q}_{\pi}(s^{(t)}, a^{(t)}) + \eta \left[(R(s^{(t)}, a^{(t)}, s^{(t+1)}) + \gamma \hat{Q}_{\pi}(s^{(t+1)}, \pi(s^{(t+1)}))) - \hat{Q}_{\pi}(s^{(t)}, a^{(t)}) \right] \quad (*)$$

(1) (*) 可進一步改寫為

$$\hat{Q}_{\pi}(s, \pi) \leftarrow \eta (R(s, a, s') + \gamma \hat{Q}_{\pi}(s', \pi(s'))) + (1 - \eta) \hat{Q}_{\pi}(s, a)$$

$\Rightarrow \eta \in [0, 1]$ 之物理一為 **forget rate**

(2) $\eta \rightarrow 0$ 時, $\hat{Q}_{\pi}(s, a)$ degenerates to the avg.

of accumulative reward 且 $\hat{Q}_{\pi}(s, a) \rightarrow Q_{\pi}(s, a)$

3. SARSA 的步驟

Algorithm. SARSA

Input: S, A, γ

Output: $\pi^*(s), \forall s$

1. Initialize $Q_{\pi}(s, a)$ arbitrarily, $\forall s, a$

2. For each episode, do:

(1) Set s to initial state

(2) repeat until s is terminal state

① $a \leftarrow \arg \max_a Q_{\pi}(s, a')$ (先做動作)

② observe s' and reward $R(s, a, s')$

③ $Q_{\pi}(s, a) \leftarrow Q_{\pi}(s, a) + \eta [(R(s, a, s') + \gamma Q_{\pi}(s', \pi(s')))]$

- $Q_n(s^{(t)}, a^{(t)})$] (再做更新)

④ $S \leftarrow S'$

(1) 每次決定下一動作 a 時, policy 都會改進

(2) 將 V 改成 Q -function 是一大躍進, 因為不僅走通了 value & policy iteration, 且增快速度。

(3). SARSA 的收斂性:

Thm. SARSA converges and gives the optimal policy π^* almost surely if

(1) π is greedy in the limit with infinite exploration

(2) η is small enough eventually, but not decreasing too fast. Furthermore, η should satisfy $\sum \eta^{(t)} = \infty$ and $\sum \eta^{(t)^2} < \infty$.

Remark. (1) Greedy in the limit: the policy $\pi \rightarrow$ exploitation / greedy policy (in the limit)

(2) Infinite exploration: all (s, a) pairs are visited infinite times

\rightarrow SARSA 的第 2(2) ① 無法保證!

4. Exploration Strategies:

現希望 mix-in / exploration 的策略使 infinite exploration with π 可達成

(1). ϵ -Greedy Strategy

\Rightarrow At every time step:

ξ 以 ϵ 之機率隨機移動 (explore)

"(1- ϵ)" update exploitation policy 並隨之 act (exploit)

(2) Softmax Strategy

① 如果 a 曾经得到了多 accumulative reward, 则就更常使用 a

⇒ 在 SARSA 中, 可由此 distribution 来 sample a

$$P(a|s) = \frac{\exp(Q_\pi(s,a)/t)}{\sum_{a'} \exp(Q_\pi(s,a')/t)}, \forall a$$

temperature

② $t_{\text{高}}$: exploration
 $t_{\text{低}}$: exploitation

(3) Exploration Function

⇒ to explore areas with fewest samples.

⇒ 比如: 在 SARSA 中, 定义 exploration function

$$f(g, n) = g + \frac{K}{n}$$

其中, g 为 Q -value 之 estimate

n 为 the number of samples for the estimate

K 为某常数

⇒ 将 $Q_\pi(s,a)$ 的 update 改用

$$Q_\pi(s,a) \leftarrow Q_\pi(s,a) + \eta ([R(s,a,s') + \gamma f(Q_\pi(s',a'), \text{count}(s',a'))] - Q_\pi(s,a))$$

⇒ Infinite exploration; exploit once exploring enough.

1. Q-learning (model free)

Def. (Optimal Q function)

$$Q^*(s,a) = \max_{\pi} E_{s^{(0)}, \dots} [R(s,a,s') + \sum_{t=1}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s, a, \pi]$$

s.t. $V^*(s) = \max_a Q^*(s,a)$ with the recurrence

$$Q^*(s,a) = \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma \max_{a'} Q^*(s',a')], \forall s.$$

Def. Optimal value function.

$$V^*(s) = \max_{\pi} E_{S^{(0)}, \dots} \left[\sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) \mid S^{(0)} = s, \pi \right]$$

with recurrence

$$V^*(s) = \max_a \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma V^*(s')]]$$

Remark. ① OVF 之物理-1-1: 当初始 state 是 s 且 acting optimally onward 时, 其 accumulative reward 之期望值最大

② OQF: "在 taking action a 时,"

⇒ 将 value iteration 改高为 Q-value iteration

1. Q-value iteration

Algorithm. Q-value iteration.

Input: MDP $(S, A, P, R, \gamma, H \rightarrow \infty)$

Output: $\pi^*(s) \forall s$

1. Initialize $Q^*(s, a) = 0, \forall s, a$

2. Repeat until $Q^*(s, a)$ converge

$$Q^*(s, a) \leftarrow \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]]$$

3. For each s do

$$\pi^*(s) \leftarrow \arg \max_a Q^*(s, a)$$

2. Temporal Difference Estimation for exploitation policy π :

(1) $Q^*(s, a) \leftarrow$ random value, $\forall s, a$

(2) Repeat until converge, $\forall a^{(t)}$

$$Q^*(s^{(t)}, a^{(t)}) \leftarrow \hat{Q}^*(s^{(t)}, a^{(t)}) + \eta [R(s^{(t)}, a^{(t)}, s^{(t+1)}) + \gamma \max_a Q^*(s^{(t+1)}, a) - \hat{Q}^*(s^{(t)}, a^{(t)})]$$

3. Q-Learning 的步驟

Algorithm. Q-Learning

Input: S, A, r

Output: $\pi^*(s), \forall s$

1. Initialize $Q^*(s, a)$ arbitrarily

2. For each episode do

(1) Set s to initial state

(2) Repeat until s is terminal state

① Taking action a from s using some "exploration" policy π' derived from Q^* (ex. ϵ -greedy)

② Observe s' and reward $R(s, a, s')$

③ $Q^*(s, a) \leftarrow Q^*(s, a) + \eta (R(s, a, s') + \gamma \max_a Q^*(s', a') - Q^*(s, a))$

④ $s \leftarrow s'$

4. Q-learning 之收敛性

Thm. Q-learning converges and gives the **optimal policy** π^* if

(1) All states and actions are visited infinitely often (π' has explored enough)

(2) η satisfies $\sum_t \eta^{(t)} = \infty, \sum_t \eta^{(t)^2} < \infty$.

5. SARSA 与 Q-learning 之比较

Def. ① Off-policy: π updated toward a greedy policy indep. with π'

② On-policy: π updated to improve (and **depend on**) π'

ex. off-policy: Q-learning
on-policy: SARSA

1. 實務上: SARSA 較保守, 不願犯錯, 但 Q-learning 較願意冒險, '在 exploration policy 改變時可持續學習: 死的是別人(其他元)', 願意冒險

2. 空間複雜度: SARSA 與 Q-learning 皆為 $O(|S||A|)$
 \Rightarrow memory usage 可能很大, 要注意

*. Concluding remark: RL 除了在遊戲 AI 方面之應用外, 今日之 LLM 亦將其應用於訓練中, 其重要性可見一斑