

微積分 + AI Deep Learning + PyTorch

註記:

括號符號 $(scalar)$, $[tensor]$, $||norm||$

如：向量可使用 $v_I = [v_i]_{i \in I}$

如：矩陣可使用 $a_{IJ} = [a_{ij}]_{(i,j) \in I \times J}$ 表示

單一大寫字母代表集合，如： $I, J, K, L \dots$ 代表指標集，如太多維可使用 I^1, I^2, \dots

對於取用同個指標集 I 子集的元素，可使用 i, i', i'' ，如太多個可使用 $i^{(2)}, i^{(3)} \dots$

紅色代表為資料流，運算結果流

藍色代表模型參數 (*parameters*)，訓練時會受到最佳化演算法改變

綠色代表模型超參數 (*hyper-parameters*)，不會受訓練影響

紫色括號 $[P] \in \{0, 1\}$ ，代表 Iverson Bracket

torch.nn

1. Linear Layers

$$Linear(\mathbf{x}_J; \mathbf{w}_{IJ}, \mathbf{b}_I) := \mathbf{w}_{IJ} \cdot \mathbf{x}_J + \mathbf{b}_I = \left[\sum_{j \in J} \mathbf{w}_{ij} \mathbf{x}_j + \mathbf{b}_i \right]_{i \in I}$$

Gradient :

$$\nabla_{\mathbf{w}_{IJ}} Linear = [\mathbf{x}_j]_{(i,j) \in I \times J}$$
$$\nabla_{\mathbf{b}_I} Linear = \mathbf{1}_I$$

2. Nonlinear Activations

$$Softmax(\mathbf{x}_I) = \left[\frac{e^{\mathbf{x}_i}}{\|e^{\mathbf{x}_I}\|_1} \right]_{i \in I} \quad \text{其中分母 } \|e^{\mathbf{x}_I}\|_1 := \sum_{i' \in I} e^{\mathbf{x}_{i'}}$$

$$Tanh(\mathbf{x}_I) := \left[\frac{e^{\mathbf{x}_i} - e^{-\mathbf{x}_i}}{e^{\mathbf{x}_i} + e^{-\mathbf{x}_i}} \right]_{i \in I}$$

Gradient :

$$\nabla_{\mathbf{x}_I} \text{Tanh}(\mathbf{x}_I) = [1 - \text{Tanh}^2(\mathbf{x}_i)]_{i \in I}$$

$$\text{Sigmoid}(\mathbf{x}_I) := \left[\frac{1}{1 + e^{-\mathbf{x}_i}} \right]_{i \in I}$$

Gradient :

$$\nabla_{\mathbf{x}_I} \text{Sigmoid}(\mathbf{x}_I) = \text{Sigmoid}(\mathbf{x}_I) \odot (1_I - \text{Sigmoid}(\mathbf{x}_I))$$

$$\text{LogSoftmax}(\mathbf{x}_I) := \ln \circ \text{Softmax}(\mathbf{x}_I) = \left[\ln \left(\frac{e^{\mathbf{x}_i}}{\|e^{\mathbf{x}_I}\|_1} \right) \right]_{i \in I}$$

$$\text{Softplus}(\mathbf{x}_I; \beta) = \left[\frac{1}{\beta} \ln(1 + e^{\beta \mathbf{x}_i}) \right]_{i \in I} \quad \text{註：}\beta \text{ 每個維度都是同個參數}$$

$$\text{Softsign}(\mathbf{x}_I) = \left[\frac{\mathbf{x}_i}{1 + |\mathbf{x}_i|} \right]_{i \in I}$$

$$\text{Threshold}(\mathbf{x}_I, \alpha) = \left[\mathbf{x}_i \cdot [\mathbf{x}_i > \alpha] + \alpha \cdot [\mathbf{x}_i \leq \alpha] \right]_{i \in I}$$

$$\text{ReLU}(\mathbf{x}_I) = \left[\max(0, \mathbf{x}_i) \right]_{i \in I}$$

$$\text{PReLU}(\mathbf{x}_I) = \left[\max(0, \mathbf{x}_i) + a \cdot \min(0, \mathbf{x}_i) \right]_{i \in I}$$

註：激發函數通常寫成 $\sigma(\mathbf{x}_I)$

3. Dropout Layers

$$\text{Dropout}(\mathbf{x}_I; p) = \left[0 \cdot [\#_i \leq p] + \mathbf{x}_i \cdot [\#_i > p] \right]_{i \in I} \quad \text{其中 } \#_I \sim \text{uniform}\left((0, 1)^{|I|}\right) \text{ 為隨機向量}$$

4. Sparse Layers

$$\text{Embedding}(\mathbf{z}_J; \mathbf{I}, \mathbf{D}) := \left[\mathbf{w}_{\mathbf{z}_j d} \right]_{J \times D} \quad \text{其中 } \mathbf{z}_J \in I^{|J|}, \text{ 參數矩陣為 } \mathbf{w}_{ID} := [w_{id}]_{(i,d) \in I \times D}$$

註：在 NLP(Natural Language Processing)領域

I 代表詞種類(words)集合

$|D|$ 為 Word Embedding Dimension

不同的詞，可用 $1, 2, 3, 4 \dots |I|$ 編號，即 $I \xleftrightarrow{1-1} \{1, 2, 3, \dots |I|\}$

$word_i$ 的詞向量($word2vec$)即為 $[\omega_{id}]_{d \in D}$

一句有 $|J|$ 個詞的句子 $= z_J = [z_j] = [\text{第 } j \text{ 個詞}]_{j=1,2,\dots,|J|}$

$Embedding(\text{句子}) = [\text{詞實向量}] = \text{實矩陣}$

核心量化概念：詞 \rightarrow 向量，句子 \rightarrow 矩陣(有順序概念)， j 代表位置， d 代表詞特徵，詞特徵是演算法學來的!!

5.Distance Functions

$$\text{CosineSimilarity}(u_I, v_I) = \frac{u_I \cdot v_I}{\max(\|u_I\|_2 \|v_I\|_2, \epsilon)} = \frac{\sum_{i \in I} u_i v_i}{\max\left(\sqrt{\sum_{i \in I} u_i^2 \sum_{i \in I} v_i^2}, \epsilon\right)}$$

$$\text{PairwiseDistance}(u_I; p) := \|u_I\|_p = \left(\sum_{i \in I} |x_i|^p\right)^{\frac{1}{p}}$$

y_I^{pred} 代表經由數學模型計算後的預測向量

y_I^{target} 代表原始資料的目標向量(正確答案)

$|\mathcal{B}|$ 代表 batchsize

$y_{ib}^{pred}, y_{ib}^{target}$ 代表樣本 b 的 y 值

註: 模型架構好以後 Pytorch 支援 Batch Input !!

$$\text{Model}([x_{Ib}]_{b \in \mathcal{B}}) := \left[\text{Model}(x_{Ib}) \right]_{b \in \mathcal{B}}$$

6. Loss Functions (Sum Overall Batch Samples)

$$\text{L1Loss}\left(\left[(y_I^{pred}, y_I^{target})\right]_{b \in \mathcal{B}}\right) = \sum_{b \in \mathcal{B}} \left(\frac{1}{|I|} \sum_{i \in I} |y_{ib}^{pred} - y_{ib}^{target}| \right)$$

$$\text{MSELoss}\left(\left[(y_I^{pred}, y_I^{target})\right]_{b \in \mathcal{B}}\right) = \sum_{b \in \mathcal{B}} \left(\frac{1}{|I|} \sum_{i \in I} (y_{ib}^{pred} - y_{ib}^{target})^2 \right)$$

$$\text{CrossEntropyLoss}\left(\left[(y_I^{pred}, y_I^{target})\right]_{b \in \mathcal{B}}\right) = H\left(y_I^{target}, \text{Softmax}(y_I^{pred})\right) = - \sum_{i \in I} y_i^{target} \ln \left(\frac{y_i^{pred}}{\sum_{i' \in I} y_{i'}^{pred}} \right)$$

註： y_I^{target} is one hot encoding, API use integer input

$$\text{CRF}(s_{IY}; \omega_{YY}) = - \left(\sum_{i=1}^{|I|} s_{iy_i} + \sum_{i=1}^{|I|-1} \omega_{y_i, y_{i+1}} \right)$$

註： $|I|$ 為句子長度， $|Y|$ 為 Label 種類集， s_{i, y_i} 又稱為 emission score， y_I 為 target label vector

註 2：Bi-LSTM輸出為 x_{ID} 向量， $|D|$ 為 hidden dimension，需要再作線性轉換 $S_{IY}, \text{Linear}(x_{ID}) = s_{IY}$

7.Recurrent Layers (Share Weight Matrix)

h_D 為接口 (hidden dimension), h_D^0 可fixed或可加入一起學習

$$RNNCell(\mathbf{x}_I, h_D; \omega_{DI}, \omega_{DD}, b_D) = \tanh(\omega_{DI} \mathbf{x}_I + \omega_{DD} h_D + b_D)$$

$$LSTMCell(\mathbf{x}_I, c_D, h_D; \overbrace{\omega_{DI}^{x \rightarrow i}, \omega_{DD}^{h \rightarrow i}, \omega_{DI}^{x \rightarrow f}, \omega_{DD}^{h \rightarrow f}, \omega_{DI}^{x \rightarrow g}, \omega_{DD}^{h \rightarrow g}, \omega_{DI}^{x \rightarrow o}, \omega_{DD}^{h \rightarrow o}}^{8 \text{ weight matrixs}}, \underbrace{b_D^i, b_D^f, b_D^g, b_D^o}_{4 \text{ bias vectors}})$$

結構細節:

$$\begin{aligned} f_D &:= \sigma(\omega_{DI}^{x \rightarrow f} \mathbf{x}_I + \omega_{DD}^{h \rightarrow f} h_D + b_D^f) \\ i_D &:= \sigma(\omega_{DI}^{x \rightarrow i} \mathbf{x}_I + \omega_{DD}^{h \rightarrow i} h_D + b_D^i) \\ o_D &:= \sigma(\omega_{DI}^{x \rightarrow o} \mathbf{x}_I + \omega_{DD}^{h \rightarrow o} h_D + b_D^o) \\ g_D &:= \tanh(\omega_{DI}^{x \rightarrow g} \mathbf{x}_I + \omega_{DD}^{h \rightarrow g} h_D + b_D^g) \\ LSTMCell(\mathbf{x}_I, c_D, h_D; \dots) &= o_D \odot \tanh(f_D \odot c_D + i_D \odot g_D) \end{aligned}$$

$$GRUCell(\mathbf{x}_I, h_D; \omega_{DI}^{x \rightarrow r}, \omega_{DD}^{h \rightarrow r}, \omega_{DD}^{x \rightarrow n}, \omega_{DI}^{x \rightarrow z}, \omega_{DD}^{h \rightarrow z}, \omega_{DD}^{h \rightarrow n}, b_D^r, b_D^n, b_D^z)$$

結構細節:

$$\begin{aligned} r_D &= \sigma(\omega_{DI}^{x \rightarrow r} \mathbf{x}_I + \omega_{DD}^{h \rightarrow r} h_D + b_D^r) \\ z_D &= \sigma(\omega_{DI}^{x \rightarrow z} \mathbf{x}_I + \omega_{DD}^{h \rightarrow z} h_D + b_D^z) \\ n_D &= \tanh(\omega_{DD}^{x \rightarrow n} \mathbf{x}_I + r_D \odot (\omega_{DD}^{h \rightarrow n} h_D) + b_D^n) \\ GRUCell(\mathbf{x}_I, h_D; \dots) &= (1_D - z_D) \odot h_D + z_D \odot n_D \end{aligned}$$

