

交差エントロピー J_{MLE} の導出方法をはっきりさせておきたい。まず多くの場合、多クラス分類において、 \mathbf{y} は $\left[y^{(1)}, \dots, y^{(m')}\right]^\top$ のような列ベクトルである。各成分 $y^{(i)}$ はクラス $k = 0, \dots, K$ のうち、データ $\mathbf{x}^{(i)}$ が属するクラスを表しているものとする。例えば $\mathbf{x}^{(i)}$ が手書き数字の「6」の画像のデータであるならば、 $K = 9$, $y^{(i)} = 6$ である。 $U^{(2)}$ の正体についてよく考え直したところ、 \mathbf{y} の推定ではなく、次のような行列になるだろうという結論に至った。

$$U^{(2)} = \begin{bmatrix} u_0^{(1)} & \dots & u_K^{(1)} \\ \vdots & & \vdots \\ u_0^{(m')} & \dots & u_K^{(m')} \end{bmatrix} = \begin{bmatrix} \mathbf{u}^{(1)\top} \\ \vdots \\ \mathbf{u}^{(m')\top} \end{bmatrix}, \quad \mathbf{u}^{(i)} = \begin{bmatrix} u_0^{(i)} \\ \vdots \\ u_K^{(i)} \end{bmatrix}, \quad i = 1, \dots, m'. \quad (1)$$

このベクトル $\mathbf{u}^{(i)}$ が、教科書の式 6.28 の $\mathbf{z} = \mathbf{W}^\top \mathbf{h} + \mathbf{b}$ に対応している。

なぜ $U^{(2)}$ が式 (1) のような行列になるかを計算によって説明する。簡単のため手書き数字認識の例を用いる。

$$X = \begin{bmatrix} x_1^{(1)} & \dots & x_{784}^{(1)} \\ \vdots & & \vdots \\ x_1^{(m')} & \dots & x_{784}^{(m')} \end{bmatrix} = \begin{bmatrix} \mathbf{x}^{(1)\top} \\ \vdots \\ \mathbf{x}^{(m')\top} \end{bmatrix}.$$

隠れ層の中にノードが n 個あると仮定すれば、 $W^{(1)} = [\mathbf{w}_1^{(1)}, \dots, \mathbf{w}_n^{(1)}]$ のように、パラメータ $W^{(1)}$ は n 個の列ベクトル $\mathbf{w}_l^{(1)}$, $l = 1, \dots, n$ を並べて作ることができる。よって $U^{(1)}$ を具体的に計算すると

$$U^{(1)} = XW^{(1)} = \begin{bmatrix} \mathbf{x}^{(1)\top} \\ \vdots \\ \mathbf{x}^{(m')\top} \end{bmatrix} [\mathbf{w}_1^{(1)}, \dots, \mathbf{w}_n^{(1)}] = \begin{bmatrix} \mathbf{x}^{(1)\top} \mathbf{w}_1^{(1)} & \dots & \mathbf{x}^{(1)\top} \mathbf{w}_n^{(1)} \\ \vdots & & \vdots \\ \mathbf{x}^{(m')\top} \mathbf{w}_1^{(1)} & \dots & \mathbf{x}^{(m')\top} \mathbf{w}_n^{(1)} \end{bmatrix}$$

となる。すなわち $U^{(1)}$ は $m' \times n$ の行列で、その成分は

$$U_{ij}^{(1)} = \mathbf{x}^{(i)\top} \mathbf{w}_j^{(1)}, \quad i = 1, \dots, m', \quad j = 1, \dots, n$$

である。次に各成分に対する活性化関数の値を並べた行列 H を計算するが、このサイズは $U^{(1)}$ と同じで $m' \times n$ である。 H を m' 個の行ベクトル $\mathbf{h}^{(i)\top} = [\varphi(\mathbf{x}^{(i)\top} \mathbf{w}_1^{(1)}), \dots, \varphi(\mathbf{x}^{(i)\top} \mathbf{w}_n^{(1)})]$, $i = 1, \dots, m'$ を縦に並べたものと解釈する。またパラメータ $W^{(2)}$ は、クラスについての列ベクトル $\mathbf{w}_k^{(2)}$, $k = 0, \dots, K$ を並べたものとするれば、行列積 $U^{(2)} = HW^{(2)}$ は列ベクトルである必要はない。

$$H = \begin{bmatrix} \mathbf{h}^{(1)\top} \\ \vdots \\ \mathbf{h}^{(m')\top} \end{bmatrix}, \quad W^{(2)} = [\mathbf{w}_0^{(2)}, \dots, \mathbf{w}_K^{(2)}]$$

であるから、

$$U^{(2)} = HW^{(2)} = \begin{bmatrix} \mathbf{h}^{(1)\top} \\ \vdots \\ \mathbf{h}^{(m')\top} \end{bmatrix} [\mathbf{w}_0^{(2)}, \dots, \mathbf{w}_K^{(2)}] = \begin{bmatrix} \mathbf{h}^{(1)\top} \mathbf{w}_0^{(2)} & \dots & \mathbf{h}^{(1)\top} \mathbf{w}_K^{(2)} \\ \vdots & & \vdots \\ \mathbf{h}^{(m')\top} \mathbf{w}_0^{(2)} & \dots & \mathbf{h}^{(m')\top} \mathbf{w}_K^{(2)} \end{bmatrix}$$

となる。すなわち $U^{(2)}$ は $m' \times (K + 1)$ の行列で、その成分は

$$U_{ik}^{(2)} = \mathbf{h}^{(i)\top} \mathbf{w}_k^{(2)}, \quad i = 1, \dots, m', \quad k = 0, \dots, K$$

である。 $U^{(2)}$ の成分を横方向に見ると、ミニバッチの各データに関する情報が並んでいて、縦方向に見ると各クラスのスコアが並んでいるような構造をしている。それを表して書き直したのが式 (1) である。

さて式 (1) の行列を用いて、交差エントロピーがどのように作られるのかを思い出したい。6.2.2.3 節「マルチヌーイ出力分布のためのソフトマックスユニット」の復習である。

i 番目のデータがクラス k に属する確率は、ソフトマックス関数を用いて次のように与えられるのだった。

$$p_k^{(i)} = \text{softmax}(\mathbf{u}^{(i)})_k = \frac{\exp(u_k^{(i)})}{\sum_{k=0}^K \exp(u_k^{(i)})}. \quad (2)$$

これは規格化されているので、あえてチルダ記号はつけていない。記号を少し整理しておく。

- $u_k^{(i)}$: i 番目のデータに対するクラス k の「スコア」
- $\mathbf{u}^{(i)}$: i 番目のデータに対する各クラスのスコアを格納したベクトル
- $p_k^{(i)}$: 各クラスのスコアから推計される、 i 番目のデータがクラス k に属する確率

ちなみに式 (2) の対数をとることで、行列 $U^{(2)}$ の各成分 $u_k^{(i)}$ が規格化されていない対数尤度になっていることを示すことができる。

$$u_k^{(i)} = \log p_k^{(i)} + \log \left(\sum_{k=0}^K \exp(u_k^{(i)}) \right) = \log \tilde{p}_k^{(i)}. \quad (3)$$

訓練データのラベル $y^{(i)}$ が k であるとする、その確率は

$$\mathbf{1}_{k=y^{(i)}} = \begin{cases} 1, & k = y^{(i)}, \\ 0, & k \neq y^{(i)} \end{cases}$$

と表せる。これと式 (2) の推計された確率を用いれば、交差エントロピーのコスト関数は次のように作ることができる。

$$J_{\text{MLE}} = -\frac{1}{m'} \sum_{i=1}^{m'} \sum_{k=0}^K \mathbf{1}_{k=y^{(i)}} \log p_k^{(i)}. \quad (4)$$

なお式 (3) の規格化されていない対数尤度 $u_k^{(i)} = \log \tilde{p}_k^{(i)}$ をそのまま交差エントロピーに用いると、式 (4) によるものとは異なる結果が得られると考えられる。なぜならば式 (3) の定数に思える部分 $\sum_{k=0}^K \exp(u_k^{(i)})$ は実際には訓練データ i とパラメータ $W^{(1)}$, $W^{(2)}$ の関数であり、 $u_k^{(i)} = \log \tilde{p}_k^{(i)}$ を用いた交差エントロピーは式 (4) のそれに $W^{(1)}$, $W^{(2)}$ の関数を付け足したものになるからである。

参考文献

- [1] Aurélien Géron, 長尾高弘訳, 「scikit-learn と TensorFlow による実践機械学習」, 株式会社オライリージャパン, 2018 年.