

(画像は最後にある)

2 次元ベクトルにおける 4 クラス識別器の設計

1. はじめに

機械学習は本質的に標本 \mathbf{x} からラベル y への変換 $f: y = f(\mathbf{x}, \mathbf{w})$ (ここで \mathbf{w} はモデルパラメータ) を設定し、最適化する作業と思われる。変換 $\mathbf{y} = f(\mathbf{x}, \mathbf{w})$ が最も単純な線形関数を使用する場合は、線形識別関数が得られる。

2. 約束

訓練データセット: $\{\mathbf{x}_i | \mathbf{x}_i \in \mathbb{R}^D, 1 \leq i \leq N, N = 1000, D = 2\}$

標本 \mathbf{x}_i に対応するクラスのラベル: $y_i (y_i \in \mathbb{Z}, 1 \leq y_i \leq K, K = 4)$

訓練データセットの標本総数は N 、標本の次元は 2 、クラス総数は 4 である。

3. 線形識別関数

標本 \mathbf{x}_i とパラメータ \mathbf{w} の線形識別関数の形式は

$$s = f(\mathbf{x}_i, \mathbf{w}, b) = x_{i1}w_1 + x_{i2}w_2 + \cdots + x_{iD}w_D + b$$

ここで s はスコアで、 $\mathbf{w} = (w_1, w_2, \dots, w_D)^T \in \mathbb{R}^D$ 、 b は定数項である。 s はスカラーなので 1 クラスのスコアしか表せない。 K 個のクラス全てのスコアを得るために K 個の線形関数が必要であり

$$s_j = f(\mathbf{x}_i, \mathbf{w}, b) = w_{j1}x_{i1} + w_{j2}x_{i2} + \cdots + w_{jD}x_{iD} + b_j$$

j 番目の方程式はクラス j のスコアを表す。行列でまとめて書き換えると

$$\begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_K \end{bmatrix} = \mathbf{s}_i = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x}_i + b_1 \\ \mathbf{w}_2^T \mathbf{x}_i + b_2 \\ \vdots \\ \mathbf{w}_K^T \mathbf{x}_i + b_K \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T, b_1 \\ \mathbf{w}_2^T, b_2 \\ \vdots \\ \mathbf{w}_K^T, b_K \end{bmatrix} \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \quad (1 \leq i \leq N).$$

さらに

$$\begin{bmatrix} s_{11} & s_{21} & \cdots & s_{N1} \\ s_{12} & s_{22} & \cdots & s_{N2} \\ \vdots & \vdots & \ddots & \vdots \\ s_{K1} & s_{K2} & \cdots & s_{KN} \end{bmatrix} = [\mathbf{s}_1 \mathbf{s}_2 \cdots \mathbf{s}_N] = \begin{bmatrix} \mathbf{w}_1^T, b_1 \\ \mathbf{w}_2^T, b_2 \\ \vdots \\ \mathbf{w}_K^T, b_K \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}$$

すなわち

$$\mathbf{S} = \mathbf{W}\mathbf{X}.$$

と書ける。ここで \mathbf{S} はスコア行列。

最適化アルゴリズムはこの重み係数行列 \mathbf{W} に対して行われる。

4. 損失関数 (loss function)

損失関数は \mathbf{S} の質を評価する。 \mathbf{S} と希望されるラベルの差異が小さいほど損失関数の値が小さくなる。ベクトル \mathbf{x} に対して

$$\text{softmax}(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\mathbf{1}^T e^{\mathbf{x}}}$$

とする。関数値の成分は全て $(0,1)$ 上の実数なので「確率」の目安として捉えて良いだろう。ただし、 $\mathbf{1}$ は全ての成分が 1 の列ベクトルとする。

K 次元の列ベクトル $\hat{\mathbf{y}}_i$ を y_i の one-hot ベクトル表現として

$$L_i = -\hat{\mathbf{y}}_i^T \log \text{softmax}(\mathbf{s}_i) \quad (1 \leq i \leq N)$$

が標本 \mathbf{x}_i における損失である。全体の損失は L_i の平均値である。

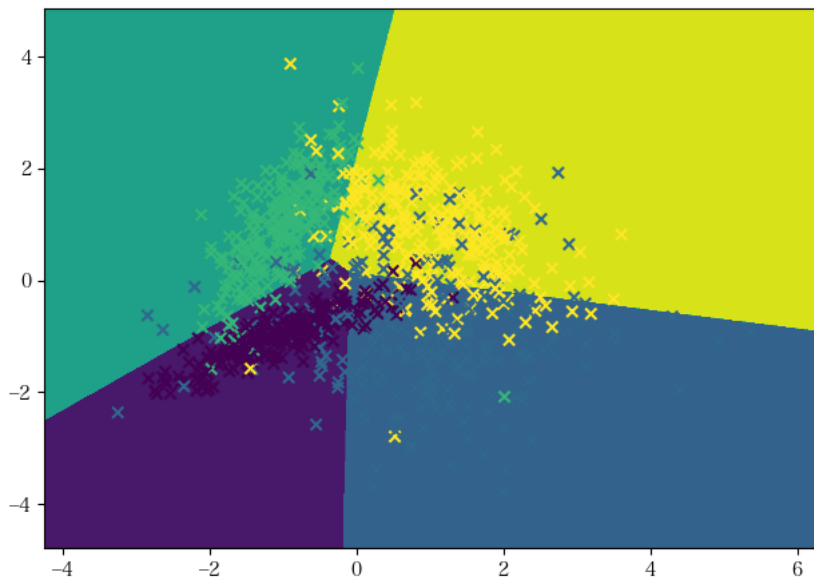
5. 重み係数の最適化

最急降下法 (gradient descent method) を使う。まずは L の勾配を求める。

$$\begin{aligned}
 L_i &= -\hat{\mathbf{y}}_i^T \log \text{softmax}(\mathbf{s}_i) & \text{よって} \\
 &= -\hat{\mathbf{y}}_i^T (\log e^{\mathbf{s}_i} - \mathbf{1} \log(\mathbf{1}^T e^{\mathbf{s}_i})) & \frac{\partial L_i}{\partial W} = (\text{softmax}(\mathbf{s}_i) - \hat{\mathbf{y}}_i) \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix}^T \\
 &= -\hat{\mathbf{y}}_i^T \log(e^{\mathbf{s}_i}) + \log(\mathbf{1}^T e^{\mathbf{s}_i}) \quad (\because \hat{\mathbf{y}}_i^T \mathbf{1} = 1) \\
 &= -\hat{\mathbf{y}}_i^T \mathbf{s}_i + \log(\mathbf{1}^T e^{\mathbf{s}_i}), & \text{grad } L = \frac{1}{N} \sum_{i=1}^N \frac{\partial L_i}{\partial W} \\
 dL_i &= -\hat{\mathbf{y}}_i^T d\mathbf{s}_i + \frac{\mathbf{1}^T (e^{\mathbf{s}_i \odot d\mathbf{s}_i})}{\mathbf{1}^T e^{\mathbf{s}_i}} \\
 &= -\hat{\mathbf{y}}_i^T d\mathbf{s}_i + \frac{(e^{\mathbf{s}_i})^T d\mathbf{s}_i}{\mathbf{1}^T e^{\mathbf{s}_i}} & = \frac{1}{N} [\text{softmax}(\mathbf{s}_1) - \hat{\mathbf{y}}_1 \quad \dots \quad \text{softmax}(\mathbf{s}_N) - \hat{\mathbf{y}}_N] \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \\ 1 \end{bmatrix} \\
 &= \text{tr}(-\hat{\mathbf{y}}_i^T d\mathbf{s}_i + \text{softmax}(\mathbf{s}_i)^T d\mathbf{s}_i) \\
 &= \text{tr}((\text{softmax}(\mathbf{s}_i)^T - \hat{\mathbf{y}}_i^T) d\mathbf{s}_i) & = \frac{1}{N} (\text{softmax}(S) - \hat{Y}) X^T. \\
 &= \text{tr}((\text{softmax}(\mathbf{s}_i) - \hat{\mathbf{y}}_i)^T dW \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix}) \\
 &= \text{tr}(\begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} (\text{softmax}(\mathbf{s}_i) - \hat{\mathbf{y}}_i)^T dW)
 \end{aligned}$$

が得られる。これにより訓練すると、線形識別可能なデータに対してほぼ 100% の識別率が達成できる。実際の訓練では、線形識別不可なデータでも 81.3% が出た。

線形判別関数による4種類分類器の分類結果 (1116201017賈書瑞)



参考文献

- [1] <https://ja.wikipedia.org/wiki/最急降下法>
- [2] <https://www.ituring.com.cn/book/tupubarticle/25626>
- [3] <https://zhuanlan.zhihu.com/p/24709748> (行列に対する微分)
- [4] <https://numpy.org/doc/stable/>
- [5] <https://blog.csdn.net/mqq9931/article/details/83829849>