

(日本語は下にあります)

We apply Logistic Regression model on Hand-written digits recognition task.

In logistic regression model, an occur probability of an event is represented by a logistic function. For example, in two-class problem, logistic sigmoid function is commonly used. In multi-class problem, softmax function is known offering a well performance.

In this text, we introduce how a posteriori probability can be represented in logistic function. Next, we imply linear transformation of each category to input patterns and estimate their posteriori probabilities.

Consider a classification task with the number of class C . Using maximum a posteriori probability rule, the output category of an input pattern \mathbf{x} is the category y where $p(y = i|\mathbf{x})$ is maximum.

A posteriori probability $p(y = i|\mathbf{x})$ is,

$$p(y = i|\mathbf{x}) = \frac{p(\mathbf{x}|y = i)p(y = i)}{\sum_{j=1}^C p(\mathbf{x}|y = j)p(y = j)} \quad (1)$$

Here, let

$$a_i = \ln p(\mathbf{x}|y = i)p(y = i) \quad (2)$$

Therefore,

$$\begin{aligned} p(y = i|\mathbf{x}) &= \frac{p(\mathbf{x}|y = i)p(y = i)}{\sum_{j=1}^C p(\mathbf{x}|y = j)p(y = j)} \\ &= \frac{\exp(a_i)}{\sum_{j=1}^C \exp(a_j)} \\ &= \text{softmax}(a_i) \end{aligned} \quad (3)$$

$\text{softmax}(a_i) = \frac{\exp(a_i)}{\sum_{j=1}^C \exp(a_j)}$ is called softmax function. When $C = 2$, it is known as logistic sigmoid function. From (3), a posteriori probability can be represented by a logistic function.

In multi-class task, we use softmax function. $p(y = i|\mathbf{x}) = \text{softmax}(a_i)$. To estimate this probability, we estimate a_i by linear transformation of input pattern \mathbf{x} in to i -th category. This means, $a_i \approx \mathbf{w}_i^T \mathbf{x}$, where parameter \mathbf{w}_i is estimated by maximum likelihood method.

Hence, a posteriori probability is $p(y = i|\mathbf{x}) \approx \text{softmax}(\mathbf{w}_i^T \mathbf{x}) = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{j=1}^C \exp(\mathbf{w}_j^T \mathbf{x})} \stackrel{\text{def}}{=} \pi_i(\mathbf{x})$

Let C be the number of categories, $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be sample data set. Likelihood function is

$$p(y_1, y_2, \dots, y_N | \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) = \prod_{n=1}^N \prod_{k=1}^C p(y = k | \mathbf{x}_n)^{\mathbb{I}[y_n=k]} = \prod_{n=1}^N \prod_{k=1}^C \pi_{nk}^{\mathbb{I}[y_n=k]} \quad (4)$$

where

$$\mathbb{I}[p] = \begin{cases} 1 & \text{if } p \text{ is true} \\ 0 & \text{if } p \text{ is false} \end{cases}$$

$$\pi_{nk} = p(y = k | \mathbf{x}_n) = \frac{\exp(a_{nk})}{\sum_{j=1}^C \exp(a_{nj})}, \quad a_{nj} = \mathbf{w}_j^T \mathbf{x}_n$$

Log-likelihood can be written as

$$L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) = \ln p(y_1, y_2, \dots, y_N | \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) = \sum_{n=1}^N \sum_{k=1}^C \mathbb{I}[y_n = k] \ln \pi_{nk} = \sum_{n=1}^N \sum_{k=1}^C t_{nk} \ln \pi_{nk} \quad (5)$$

where

$$t_{nk} = \mathbb{I}[y_n = k]$$

If $\hat{\mathbf{w}}_j$ is an estimator of \mathbf{w}_j then $\left. \frac{\partial}{\partial \mathbf{w}_j} L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) \right|_{\mathbf{w}_j = \hat{\mathbf{w}}_j} = 0$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}_j} &= \sum_{n=1}^N \sum_{k=1}^C t_{nk} \frac{\partial L}{\partial \mathbf{w}_j} \ln \pi_{nk} \\ &= \sum_{n=1}^N \sum_{k=1}^C t_{nk} \frac{\partial \pi_{nk}}{\partial a_{nj}} \frac{\partial a_{nj}}{\partial \mathbf{w}_j} = \sum_{n=1}^N \sum_{k=1}^C t_{nk} (\delta_{kj} - \pi_{nj}) \mathbf{x}_n \\ &= \sum_{n=1}^N t_{nj} \mathbf{x}_n - \sum_{n=1}^N \pi_{nj} \mathbf{x}_n \sum_{k=1}^C t_{nk} = \sum_{n=1}^N (t_{nj} - \pi_{nj}) \mathbf{x}_n = 0 \end{aligned} \quad (6)$$

We cannot analytically solve (6). Here, we use gradient ascent learning method to maximize the log-likelihood function.

Learning Algorithm :

1. Initialize \mathbf{w}_j
2. Update \mathbf{w}_j by 3. while the change of \mathbf{w}_j is less than a small value ε
3. $\mathbf{w}_j := \mathbf{w}_j + \eta \frac{\partial L}{\partial \mathbf{w}_j} (= \mathbf{w}_j + \eta \sum_{n=1}^N (t_{nj} - \pi_{nj}) \mathbf{x}_n)$, η : learning rate

手書き数字認識の問題にロジスティック回帰モデルを適応してみる。

ロジスティック回帰モデルでは、事象の生起確率をロジスティック関数で表現する。2クラス問題の場合は、ロジスティックシグモイド関数でモデル化する。多クラス問題では、その拡張であるソフトマックス関数でモデル化する。そのため、まず、各クラスの事後確率がロジスティック関数（ソフトマックス関数）で表現できることを示す。その後、入力パターンに対して、各クラスへの線形変換を導入し、事後確率を計算する。

C クラス問題を考える際、最大事後確率則より、入力パターン \mathbf{x} の予測クラスは、クラス i の事後確率 $p(y = i|\mathbf{x})$ が最大となるクラスを選べば良い。

事後確率 $p(y = i|\mathbf{x})$ は、

$$p(y = i|\mathbf{x}) = \frac{p(\mathbf{x}|y = i)p(y = i)}{\sum_{j=1}^C p(\mathbf{x}|y = j)p(y = j)} \quad (1)$$

となるが、

$$a_i = \ln p(\mathbf{x}|y = i)p(y = i) \quad (2)$$

とおけば、

$$\begin{aligned} p(y = i|\mathbf{x}) &= \frac{p(\mathbf{x}|y = i)p(y = i)}{\sum_{j=1}^C p(\mathbf{x}|y = j)p(y = j)} \\ &= \frac{\exp(a_i)}{\sum_{j=1}^C \exp(a_j)} \\ &= \text{softmax}(a_i) \end{aligned} \quad (3)$$

$\text{softmax}(a_i) = \frac{\exp(a_i)}{\sum_{j=1}^C \exp(a_j)}$ はソフトマックス関数と呼ばれる。 $C = 2$ の場合、その関数はロジスティックシグモイド関数となる。(3)より、クラス分類問題の事後確率がロジスティック関数で表現できることが示される。

多クラス問題において、ソフトマックス関数で事後確率をモデル化する。すなわち、事後確率 $p(y = i|\mathbf{x}) = \text{softmax}(a_i)$ とモデル化し、この事後確率を推定するために、 a_i を、入力パターン \mathbf{x} のクラス i への線形変換で近似する。すなわち、 $a_i \approx \mathbf{w}_i^T \mathbf{x}$ とし、パラメータ \mathbf{w}_i は訓練データを利用し、最尤推定で求まる。

この近似より、事後確率 $p(y = i|\mathbf{x}) \approx \text{softmax}(\mathbf{w}_i^T \mathbf{x}) = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{j=1}^C \exp(\mathbf{w}_j^T \mathbf{x})} \stackrel{\text{def}}{=} \pi_i(\mathbf{x})$ となる。

今クラス数を C 、訓練データを $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ とする。尤度関数は、

$$p(y_1, y_2, \dots, y_N | \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) = \prod_{n=1}^N \prod_{k=1}^C p(y = k | \mathbf{x}_n)^{\mathbb{I}_{y_n=k}} = \prod_{n=1}^N \prod_{k=1}^C \pi_{nk}^{\mathbb{I}_{y_n=k}} \quad (4)$$

ただし、

$$\mathbb{I}[p] = \begin{cases} 1 & \text{if } p \text{ is true} \\ 0 & \text{if } p \text{ is false} \end{cases} \quad \text{である。}$$

$$\pi_{nk} = p(y = k | \mathbf{x}_n) = \frac{\exp(a_{nk})}{\sum_{j=1}^C \exp(a_{nj})}, \quad a_{nj} = \mathbf{w}_j^T \mathbf{x}_n$$

対数尤度関数は、

$$L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) = \ln p(y_1, y_2, \dots, y_N | \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) = \sum_{n=1}^N \sum_{k=1}^C \mathbb{I}_{y_n=k} \ln \pi_{nk} = \sum_{n=1}^N \sum_{k=1}^C t_{nk} \ln \pi_{nk} \quad (5)$$

ここで、 $t_{nk} = \mathbb{I}_{y_n=k}$ とおいた。

各パラメータ \mathbf{w}_j の最尤推定 $\hat{\mathbf{w}}_j$ で、 $\left. \frac{\partial}{\partial \mathbf{w}_j} L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C) \right|_{\mathbf{w}_j=\hat{\mathbf{w}}_j} = 0$ が成立する。

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}_j} &= \sum_{n=1}^N \sum_{k=1}^C t_{nk} \frac{\partial}{\partial \mathbf{w}_j} \ln \pi_{nk} \\ &= \sum_{n=1}^N \sum_{k=1}^C t_{nk} \frac{\partial \pi_{nk}}{\partial a_{nj}} \frac{\partial a_{nj}}{\partial \mathbf{w}_j} = \sum_{n=1}^N \sum_{k=1}^C t_{nk} (\delta_{kj} - \pi_{nj}) \mathbf{x}_n \\ &= \sum_{n=1}^N t_{nj} \mathbf{x}_n - \sum_{n=1}^N \pi_{nj} \mathbf{x}_n \sum_{k=1}^C t_{nk} = \sum_{n=1}^N (t_{nj} - \pi_{nj}) \mathbf{x}_n = 0 \end{aligned} \quad (6)$$

(6)の方程式は解析的な解法で解けないので、対数尤度関数を最大化するために、山登り法（最急勾配降下法）を適用する。

学習アルゴリズム：

4. \mathbf{w}_j の初期化として適当な値を与える。
5. \mathbf{w}_j が変化しない(変化分が非常に小さい)まで3.を繰り返す。
6. $\mathbf{w}_j := \mathbf{w}_j + \eta \frac{\partial L}{\partial \mathbf{w}_j} (= \mathbf{w}_j + \eta \sum_{n=1}^N (t_{nj} - \pi_{nj}) \mathbf{x}_n)$, η : 学習率