

Solution

1 Learning of a multivariate Gaussian density model

The likelihood to be maximized is given as the joint probability density function for sample $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ of n independent identically distributed normal random variables

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \prod_{i=1}^n f(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)\right). \quad (1)$$

Maximizing this likelihood is equivalent to maximizing the following log-likelihood function

$$\begin{aligned} \mathcal{L}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) &= \log \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)\right) \\ &= -\frac{nd}{2} \log 2\pi - \frac{n}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k). \end{aligned} \quad (2)$$

1.1 ML estimation for the expectation (mean vector)

Taking the derivative with respect to $\boldsymbol{\mu}_k$ is straightforward and is given by

$$\begin{aligned} \frac{\partial \mathcal{L}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} &= \frac{\partial -\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} = -\frac{1}{2} \frac{\sum_{i=1}^n \partial (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} \\ &= -\frac{1}{2} \sum_{i=1}^n \frac{\partial (\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{\partial \boldsymbol{\mu}_k} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) + (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \frac{\partial \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} \end{aligned} \quad (3)$$

$$= -\frac{1}{2} \sum_{i=1}^n (\boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k))^T \frac{\partial (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} + (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \frac{\partial \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} \quad (4)$$

$$= -\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} \frac{\partial (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} + (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \frac{\partial \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} \quad (5)$$

$$= -\frac{1}{2} \sum_{i=1}^n -(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} - (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} \quad (6)$$

$$= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} \quad (7)$$

Here we used standard results as follows

(3) $\frac{\partial YZ}{\partial X} = \frac{\partial Y}{\partial X} Z + Y \frac{\partial Z}{\partial X}$

(4) a scalar is equal to its transpose, that is for the scalar $\mathbf{a}^T \mathbf{b}$ we have $\mathbf{a}^T \mathbf{b} = (\mathbf{a}^T \mathbf{b})^T$

(5) $(AB)^T = B^T A$

(6) $(\boldsymbol{\Sigma}_k^{-1})^T = \boldsymbol{\Sigma}_k^{-1}$, $\boldsymbol{\Sigma}_k^{-1}$ being symmetric

Setting to 0, multiplying both sides by Σ_k and we get the ML estimate for μ_k , that is

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

which is the sample mean. This is an unbiased estimator for μ_k since

$$\mathbb{E}[\hat{\mu}_k] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{X}_i] = \frac{1}{n} \sum_{i=1}^n \mu_k = \mu_k$$

1.2 ML estimation for the covariance matrix

The derivative with respect to Σ_k^{-1} is given by

$$\begin{aligned} \frac{\partial \mathcal{L}(\mu_k, \Sigma_k)}{\partial \Sigma_k^{-1}} &= \frac{\partial \left(-\frac{n}{2} \log |\Sigma_k| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right)}{\partial \Sigma_k^{-1}} \\ &= \frac{n}{2} \frac{\partial \log |\Sigma_k^{-1}|}{\partial \Sigma_k^{-1}} - \frac{1}{2} \sum_{i=1}^n \frac{\partial (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k)}{\partial \Sigma_k^{-1}} \end{aligned} \quad (8)$$

$$= \frac{n}{2} \frac{1}{\Sigma_k^{-1}} - \frac{1}{2} \sum_{i=1}^n \frac{\partial \text{trace} \left((\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right)}{\partial \Sigma_k^{-1}} \quad (9)$$

$$= \frac{n}{2} \frac{1}{\Sigma_k^{-1}} - \frac{1}{2} \sum_{i=1}^n \frac{\partial \text{trace} \left((\mathbf{x}_i - \mu_k) (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} \right)}{\partial \Sigma_k^{-1}} \quad (10)$$

$$= \frac{n}{2} \Sigma_k - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu_k) (\mathbf{x}_i - \mu_k)^T \quad (11)$$

Here we used standard results as follows

(8) $\frac{\partial \log |A|}{\partial A} = \frac{1}{A}$

(9) $\mathbf{x}^T A \mathbf{x} = \text{trace}(\mathbf{x}^T A \mathbf{x})$

(10) $\text{trace}(\mathbf{x}^T A \mathbf{x}) = \text{trace}(\mathbf{x} \mathbf{x}^T A)$

(11) $\frac{\partial \text{trace}(BA)}{\partial A} = B^T$

Finally, setting to zero and using the ML estimate $\hat{\mu}_k$ for μ_k yields the maximum likelihood estimate

$$\hat{\Sigma}_k = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_k) (\mathbf{x}_i - \hat{\mu}_k)^T \quad (12)$$

which is simply the sample covariance matrix. We can see that the ML estimator of Σ_k

$$\hat{\Sigma}_k = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \hat{\mu}_k) (\mathbf{X}_i - \hat{\mu}_k)^T$$

is biased. Indeed, we can show similarly as we have seen for the variance (last year), that $\mathbb{E}[\hat{\Sigma}_k] = \frac{n-1}{n} \Sigma_k$. An unbiased estimator is therefore

$$\hat{\Sigma}_k = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \hat{\mu}_k) (\mathbf{X}_i - \hat{\mu}_k)^T$$

2 Learning of a binary logistic regression model

The conditional likelihood of \mathbf{w} for the binary logistic regression model given the labeled training set of independent observations $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ is given by

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}) = \prod_{i=1}^n p(y_i | \mathbf{x}_i; \mathbf{w}) \quad (13)$$

By using the fact that

$$\begin{aligned} p(y_i = 1 | \mathbf{x}_i; \mathbf{w}) &= \pi(\mathbf{x}_i; \mathbf{w}) \\ p(y_i = 0 | \mathbf{x}_i; \mathbf{w}) &= 1 - \pi(\mathbf{x}_i; \mathbf{w}) \end{aligned}$$

we can write

$$p(y_i | \mathbf{x}_i; \mathbf{w}) = \pi(\mathbf{x}_i; \mathbf{w})^{y_i} (1 - \pi(\mathbf{x}_i; \mathbf{w}))^{1-y_i}$$

whith $y_i \in \{0, 1\}$.

The likelihood (13) can therefore be rewritten as

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}) = \prod_{i=1}^n \pi(\mathbf{x}_i; \mathbf{w})^{y_i} (1 - \pi(\mathbf{x}_i; \mathbf{w}))^{1-y_i} \quad (14)$$

and the log-likelihood is then given by

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= \log p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}) \\ &= \sum_{i=1}^n \log \left[\pi(\mathbf{x}_i; \mathbf{w})^{y_i} (1 - \pi(\mathbf{x}_i; \mathbf{w}))^{1-y_i} \right] \\ &= \sum_{i=1}^n y_i \log \pi(\mathbf{x}_i; \mathbf{w}) + (1 - y_i) \log (1 - \pi(\mathbf{x}_i; \mathbf{w})) \end{aligned} \quad (15)$$

This log-likelihood is convex but can not be maximized in a closed form. The Newton-Raphson (NR) algorithm is generally used to maximize it. A single NR update is given by :

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \left[\frac{\partial^2 \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} \right]^{-1} \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \quad (16)$$

where the Hessian and the gradient of $\mathcal{L}(\mathbf{w})$ (which are respectively the second and first derivative of $\mathcal{L}(\mathbf{w})$) are evaluated at $\mathbf{w} = \mathbf{w}^{(t)}$, l being the iteration number.

The gradient is calculated as follows.

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} &= \sum_{i=1}^n y_i \frac{\partial \log \pi(\mathbf{x}_i; \mathbf{w})}{\partial \mathbf{w}} + (1 - y_i) \frac{\partial \log (1 - \pi(\mathbf{x}_i; \mathbf{w}))}{\partial \mathbf{w}} \\ &= \sum_{i=1}^n y_i \frac{\partial \log \pi(\mathbf{x}_i; \mathbf{w})}{\partial \mathbf{w}} + (1 - y_i) \frac{\partial \log (1 - \pi(\mathbf{x}_i; \mathbf{w}))}{\partial \mathbf{w}} \\ &= \sum_{i=1}^n y_i \frac{\frac{\partial \pi(\mathbf{x}_i; \mathbf{w})}{\partial \mathbf{w}}}{\pi(\mathbf{x}_i; \mathbf{w})} + (1 - y_i) \frac{\frac{\partial (1 - \pi(\mathbf{x}_i; \mathbf{w}))}{\partial \mathbf{w}}}{(1 - \pi(\mathbf{x}_i; \mathbf{w}))} \end{aligned} \quad (17)$$

The partial derivative of the logistic function $\pi(\mathbf{x}_i; \mathbf{w})$ with respect to \mathbf{w} is given by

$$\begin{aligned}
\frac{\partial \pi(\mathbf{x}_i; \mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial \frac{\exp(\mathbf{w}^T \mathbf{x}_i)}{1 + \exp(\mathbf{w}^T \mathbf{x}_i)}}{\partial \mathbf{w}} \\
&= \frac{\frac{\partial \exp(\mathbf{w}^T \mathbf{x}_i)}{\partial \mathbf{w}} (1 + \exp(\mathbf{w}^T \mathbf{x}_i)) - \exp(\mathbf{w}^T \mathbf{x}_i) \frac{\partial (1 + \exp(\mathbf{w}^T \mathbf{x}_i))}{\partial \mathbf{w}}}{(1 + \exp(\mathbf{w}^T \mathbf{x}_i))^2} \\
&= \frac{\mathbf{x}_i \exp(\mathbf{w}^T \mathbf{x}_i) (1 + \exp(\mathbf{w}^T \mathbf{x}_i)) - \exp(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i \exp(\mathbf{w}^T \mathbf{x}_i)}{(1 + \exp(\mathbf{w}^T \mathbf{x}_i))^2} \\
&= \frac{\exp(\mathbf{w}^T \mathbf{x}_i)}{1 + \exp(\mathbf{w}^T \mathbf{x}_i)} \mathbf{x}_i - \frac{\exp(\mathbf{w}^T \mathbf{x}_i) \exp(\mathbf{w}^T \mathbf{x}_i)}{(1 + \exp(\mathbf{w}^T \mathbf{x}_i))^2} \mathbf{x}_i \\
&= \pi(\mathbf{x}_i; \mathbf{w}) \mathbf{x}_i - \pi(\mathbf{x}_i; \mathbf{w}) \pi(\mathbf{x}_i; \mathbf{w}) \mathbf{x}_i \\
&= \pi_k(\mathbf{x}_i; \mathbf{w}) (1 - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i
\end{aligned} \tag{18}$$

The gradient (17) is therefore given by

$$\begin{aligned}
\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} &= \sum_{i=1}^n y_i \frac{\pi_k(\mathbf{x}_i; \mathbf{w}) (1 - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i}{\pi(\mathbf{x}_i; \mathbf{w})} - (1 - y_i) \frac{\pi_k(\mathbf{x}_i; \mathbf{w}) (1 - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i}{(1 - \pi(\mathbf{x}_i; \mathbf{w}))} \\
&= \sum_{i=1}^n y_i (1 - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i - (1 - y_i) \pi_k(\mathbf{x}_i; \mathbf{w}) \mathbf{x}_i \\
&= \sum_{i=1}^n (y_i (1 - \pi(\mathbf{x}_i; \mathbf{w})) - (1 - y_i) \pi_k(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i \\
&= \sum_{i=1}^n (y_i - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i
\end{aligned} \tag{19}$$

which can be formulated in a matrix form as

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{X}^T (\mathbf{y} - \mathbf{p})$$

where \mathbf{X} is the $n \times (d + 1)$ matrix whose rows are the input vectors \mathbf{x}_i , \mathbf{y} is the $n \times 1$ column vector whose elements are the indicator variables y_i

$$\mathbf{y} = (y_1, \dots, y_n)^T$$

and \mathbf{p} is the $n \times 1$ column vector of logistic probabilities corresponding to the i th input

$$\mathbf{p} = (\pi(\mathbf{x}_1; \mathbf{w}), \dots, \pi(\mathbf{x}_n; \mathbf{w}))^T.$$

The hessian (matrix of second partial derivatives) is then calculated as

$$\begin{aligned}
\frac{\partial^2 \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} &= \frac{\partial \sum_{i=1}^n (y_i - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i}{\partial \mathbf{w}^T} \\
&= - \sum_{i=1}^n \mathbf{x}_i \frac{\partial \pi(\mathbf{x}_i; \mathbf{w})}{\partial \mathbf{w}^T} \\
&= - \sum_{i=1}^n \mathbf{x}_i \pi_k(\mathbf{x}_i; \mathbf{w}) (1 - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i^T \\
&= - \sum_{i=1}^n \pi_k(\mathbf{x}_i; \mathbf{w}) (1 - \pi(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i \mathbf{x}_i^T
\end{aligned} \tag{20}$$

which can be formulated in a matrix form as

$$\frac{\partial^2 \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = -\mathbf{X}^T \mathbf{W} \mathbf{X}$$

where \mathbf{W} is the $n \times n$ diagonal matrix whose diagonal elements are $\pi(\mathbf{x}_i; \mathbf{w})(1 - \pi(\mathbf{x}_i; \mathbf{w}))$ for $i = 1, \dots, n$.

The NR algorithm (16) in this case can therefore be reformulated from the Equations (20) and (21) as

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}^{(t)}) \\ &= (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \left[\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X} \mathbf{w}^{(t)} + \mathbf{X}^T (\mathbf{y} - \mathbf{p}^{(t)}) \right] \\ &= (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \left[\mathbf{W}^{(t)} \mathbf{X} \mathbf{w}^{(t)} + (\mathbf{y} - \mathbf{p}^{(t)}) \right] \\ &= (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{(t)} \mathbf{y}^* \end{aligned}$$

where $\mathbf{y}^* = \mathbf{X} \mathbf{w}^{(t)} + (\mathbf{W}^{(t)})^{-1} (\mathbf{y} - \mathbf{p}^{(t)})$ which yields the Iteratively Reweighted Least Squares (IRLS) algorithm.