

Chapter 1

Machine Learning - By Dorival Pedroso

Note: This chapter does *not* use the summation convention on repeated indices.

1.1 Linear Regression

Given $m = nSamples$ data points and $n = nFeatures$ features, the matrix \mathbf{X} organises the data along rows such that

$$\mathbf{X} = \begin{bmatrix} 1 & X_{00} & X_{01} & \cdots & X_{0n} \\ 1 & X_{10} & X_{11} & \cdots & X_{1n} \\ 1 & X_{20} & X_{21} & \cdots & X_{2n} \\ 1 & X_{30} & X_{31} & \cdots & X_{3n} \\ 1 & X_{40} & X_{41} & \cdots & X_{4n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{\mu 0} & X_{\mu 1} & \cdots & X_{\mu n} \end{bmatrix}_{(nSamples \times nFeatures + 1)} \quad (1.1)$$

where the columns from the second column correspond to each feature and $\mu = m - 1$. For example, X_{ij} is the value of data point i and feature j .

The vector of parameters is expressed as

$$\boldsymbol{\theta} = \begin{Bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{Bmatrix} \quad (1.2)$$

Thus, a linear regression applied to data point i results in

$$\ell_i(\boldsymbol{\theta}) = \sum_{j=0}^n X_{ij} \theta_j \quad \text{or} \quad \boldsymbol{\ell}(\boldsymbol{\theta}) = \mathbf{X} \boldsymbol{\theta} \quad (1.3)$$

and

$$\frac{\partial \ell_i}{\partial \theta_j} = X_{ij} \quad \text{or} \quad \frac{d\ell}{d\boldsymbol{\theta}} = \mathbf{X} \quad (1.4)$$

An error vector is defined by

$$\mathbf{e}(\boldsymbol{\theta}) = \boldsymbol{\ell}(\boldsymbol{\theta}) - \mathbf{y} \quad (1.5)$$

and the cost function by

$$C(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=0}^{\mu} (\ell_i - y_i)^2 = \frac{1}{2m} \mathbf{e}^T \mathbf{e} \quad (1.6)$$

thus

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{m} \sum_{i=0}^{\mu} (\ell_i - y_i) \frac{\partial \ell_i}{\partial \theta_j} = \frac{1}{m} \sum_{i=0}^{\mu} e_i X_{ij} \quad (1.7)$$

or

$$\frac{dC}{d\boldsymbol{\theta}} = \frac{1}{m} \mathbf{X}^T \mathbf{e} = \frac{1}{m} \mathbf{X}^T [\boldsymbol{\ell}(\boldsymbol{\theta}) - \mathbf{y}] \quad (1.8)$$

The minimum cost corresponds to

$$\frac{dC}{d\boldsymbol{\theta}} = 0 \quad \text{or} \quad \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{y} = 0 \quad (1.9)$$

Therefore,

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (1.10)$$

The gradient-descent update is

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \frac{dC}{d\boldsymbol{\theta}} \quad (1.11)$$

1.1.1 External bias parameter plus regularization

The previous θ_0 is now called b and the $\boldsymbol{\theta}$ contains one less component with the other indices being decreased by one. \mathbf{X} also contains one less column. The data matrix now is

$$\mathbf{X} = \begin{bmatrix} X_{00} & X_{01} & \cdots & X_{0\nu} \\ X_{10} & X_{11} & \cdots & X_{1\nu} \\ X_{20} & X_{21} & \cdots & X_{2\nu} \\ X_{30} & X_{31} & \cdots & X_{3\nu} \\ X_{40} & X_{41} & \cdots & X_{4\nu} \\ \vdots & \vdots & \ddots & \vdots \\ X_{\mu 0} & X_{\mu 1} & \cdots & X_{\mu \nu} \end{bmatrix} \quad (1.12)$$

$(nSamples \times nFeatures)$

where $\nu = n - 1$ (number of features minus one) with $\mu = m - 1$ still being the number of samples minus one.

Let's define the vector “one” of length equal to the number of samples as

$$o_i = 1 \quad \text{thus} \quad \mathbf{o} = [1 \ 1 \ 1 \ \dots \ 1]^T \quad (1.13)$$

Thus, the linear model can be written as

$$\ell_i(\boldsymbol{\theta}, b) = b o_i + \sum_{j=0}^{\nu} X_{ij} \theta_j \quad (1.14)$$

In vector notation

$$\ell(\boldsymbol{\theta}, b) = b \mathbf{o} + \mathbf{X} \boldsymbol{\theta}$$

We also define an error vector with all sample data as

$$e_i(\boldsymbol{\theta}, b) = \ell_i(\boldsymbol{\theta}, b) - y_i \quad (1.15)$$

or

$$\mathbf{e} = b \mathbf{o} + \mathbf{X} \boldsymbol{\theta} - \mathbf{y}$$

thus

$$\frac{\partial e_i}{\partial \theta_j} = \frac{\partial \ell_i}{\partial \theta_j} = X_{ij} \quad (1.16)$$

and

$$\frac{\partial e_i}{\partial b} = \frac{\partial \ell_i}{\partial b} = o_i \quad (1.17)$$

The regularized cost function is given by (with $l_i = \ell_i(\boldsymbol{\theta}, b)$)

$$C(\boldsymbol{\theta}, b) = \frac{1}{2m} \sum_{i=0}^{\mu} (\ell_i - y_i)^2 + \frac{\lambda}{2m} \sum_{i=0}^{\mu} \theta_i^2 \quad (1.18)$$

or, in vector notation,

$$C(\boldsymbol{\theta}, b) = \frac{1}{2m} \left(\mathbf{e}^T \mathbf{e} + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta} \right)$$

The gradient of C is calculated by two parts. The first part is

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{m} \sum_{i=0}^{\mu} (\ell_i - y_i) \frac{\partial \ell_i}{\partial \theta_j} + \frac{\lambda}{m} \sum_{i=0}^{\mu} \theta_i \delta_{ij} \quad (1.19)$$

or

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{m} \sum_{i=0}^{\mu} e_i X_{ij} + \frac{\lambda}{m} \theta_j \quad (1.20)$$

similarly

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{m} \sum_{i=0}^{\mu} X_{ji}^T e_i + \frac{\lambda}{m} \theta_j \quad (1.21)$$

or, in vector notation,

$$\frac{\partial C}{\partial \boldsymbol{\theta}} = \frac{1}{m} \mathbf{X}^T \mathbf{e} + \frac{\lambda}{m} \boldsymbol{\theta}$$

The second part is

$$\frac{\partial C}{\partial b} = \frac{1}{m} \sum_{i=0}^{\mu} (\ell_i - y_i) \frac{\partial \ell_i}{\partial b} = \frac{1}{m} \sum_{i=0}^{\mu} e_i o_i \quad (1.22)$$

or,

$$\frac{\partial C}{\partial b} = \frac{1}{m} \mathbf{o}^T \mathbf{e}$$

Let's define the vector

$$s_j = \sum_{i=0}^{\mu} o_i X_{ij} \equiv \text{sum}(\text{cols}(\mathbf{X})) \quad (1.23)$$

and the scalar

$$t = \sum_{i=0}^{\mu} o_i y_i \equiv \text{sum}(\text{cols}(\mathbf{y})) \quad (1.24)$$

In vector notation

$$\mathbf{s} = \mathbf{X}^T \mathbf{o} = \begin{Bmatrix} \sum_i^{\mu} X_{i0} \\ \sum_i^{\mu} X_{i1} \\ \dots \\ \sum_i^{\mu} X_{i\nu} \end{Bmatrix} \quad (1.25)$$

and

$$t = \mathbf{o}^T \mathbf{y} \quad (1.26)$$

Note that

$$\mathbf{s}^T = (\mathbf{X}^T \mathbf{o})^T = \mathbf{o}^T \mathbf{X} = [\sum_i^{\mu} X_{i0} \quad \sum_i^{\mu} X_{i1} \quad \dots \quad \sum_i^{\mu} X_{i\nu}] \quad (1.27)$$

Let's further define the vector

$$\mathbf{a} = \mathbf{X}^T \mathbf{y} \quad (1.28)$$

and the following matrices

$$\mathbf{A} = \mathbf{X}^T \mathbf{X} \quad (1.29)$$

and

$$B_{ij} = \frac{1}{m} s_i s_j \quad \text{hence} \quad \mathbf{B} = \frac{1}{m} \mathbf{s} \mathbf{s}^T \quad (1.30)$$

Note that these three quantities \mathbf{a} , \mathbf{A} and \mathbf{B} can be directly computed from the input data.

By expanding the gradient expressions, we get

$$\begin{aligned} \frac{\partial C}{\partial \boldsymbol{\theta}} &= \frac{1}{m} \left(b \mathbf{X}^T \mathbf{o} + \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{y} + \lambda \boldsymbol{\theta} \right) \\ &= \frac{1}{m} (b \mathbf{s} + \mathbf{A} \boldsymbol{\theta} - \mathbf{a} + \lambda \boldsymbol{\theta}) \end{aligned} \quad (1.31)$$

and

$$\begin{aligned} \frac{\partial C}{\partial b} &= \frac{1}{m} (b \mathbf{o}^T \mathbf{o} + \mathbf{o}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{o}^T \mathbf{y}) \\ &= \frac{1}{m} (b m + \mathbf{s}^T \boldsymbol{\theta} - t) \end{aligned} \quad (1.32)$$

The minimum is found by zeroing both partial derivatives and solving the following system for $\boldsymbol{\theta}$ and b

$$\begin{aligned} b \mathbf{s} + \mathbf{A} \boldsymbol{\theta} - \mathbf{a} + \lambda \boldsymbol{\theta} &= 0 \\ b m + \mathbf{s}^T \boldsymbol{\theta} - t &= 0 \end{aligned} \quad (1.33)$$

From the second equation we have

$$b = \frac{t}{m} - \frac{1}{m} \mathbf{s}^T \boldsymbol{\theta}$$

By substituting this result into the first equation, we obtain

$$\begin{aligned} \left(\frac{t}{m} - \frac{1}{m} \mathbf{s}^T \boldsymbol{\theta} \right) \mathbf{s} + (\mathbf{A} + \lambda \mathbf{I}) \boldsymbol{\theta} &= \mathbf{a} \\ -\frac{1}{m} \mathbf{s} (\mathbf{s}^T \boldsymbol{\theta}) + (\mathbf{A} + \lambda \mathbf{I}) \boldsymbol{\theta} &= \mathbf{a} - \frac{t}{m} \mathbf{s} \\ -\frac{1}{m} (\mathbf{s} \mathbf{s}^T) \boldsymbol{\theta} + (\mathbf{A} + \lambda \mathbf{I}) \boldsymbol{\theta} &= \mathbf{a} - \frac{t}{m} \mathbf{s} \end{aligned} \quad (1.34)$$

Therefore, the linear system to be solved is

$$(\mathbf{A} - \mathbf{B} + \lambda \mathbf{I}) \boldsymbol{\theta} = \mathbf{r}$$

where

$$\mathbf{r} = \mathbf{a} - \frac{t}{m} \mathbf{s}$$

After the computation of $\boldsymbol{\theta}$, b can be easily found.

1.2 Logistic Regression

The Logistic function is given by

$$g(z) = \frac{1}{1 + e^{-z}} \quad (1.35)$$

and the Logistic regression applied to each data point i is

$$h_i(\boldsymbol{\theta}) = g(\ell_i(\boldsymbol{\theta})) \quad \text{or} \quad h_i(\boldsymbol{\theta}) = \frac{1}{1 + e^{-\sum_j X_{ij} \theta_j}} \quad (1.36)$$

where the summation is indicated in Eq. (1.3).

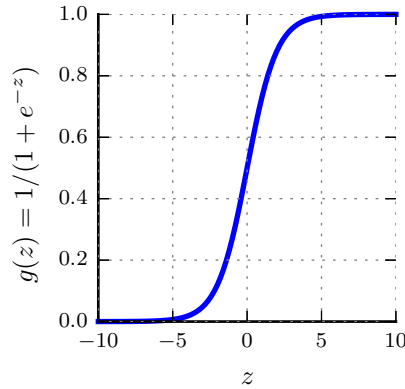


Fig. 1.1: Logistic function

Let's define p_i as

$$p_i(\boldsymbol{\theta}) = 1 + e^{-\ell_i(\boldsymbol{\theta})} \quad \text{hence} \quad h_i = (p_i)^{-1} \quad (1.37)$$

Thus (considering Eq. 1.4)

$$\frac{\partial p_i}{\partial \theta_j} = -e^{-\ell_i} \frac{\partial \ell_i}{\partial \theta_j} = -e^{-\ell_i} X_{ij} \quad (1.38)$$

Let's define q_i as

$$q_i(\boldsymbol{\theta}) = \log [p_i(\boldsymbol{\theta})] \quad (1.39)$$

Thus

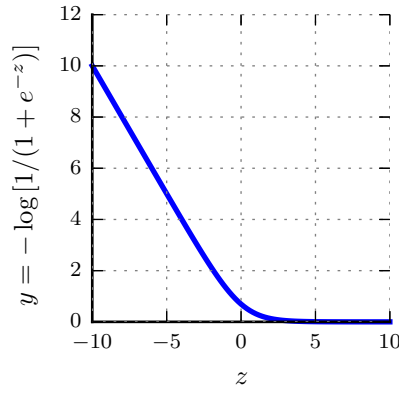
$$\frac{\partial q_i}{\partial \theta_j} = \frac{1}{p_i} \frac{\partial p_i}{\partial \theta_j} = \frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} X_{ij} \quad (1.40)$$

Note that

$$\log h_i = \log \left(\frac{1}{p_i} \right) = -\log p_i = -q_i \quad (1.41)$$

Note also that

$$\log (1 - h_i) = \log \left(\frac{p_i}{p_i} - \frac{1}{p_i} \right) = \underbrace{\log (p_i - 1)}_{-\ell_i} - \log p_i = -\ell_i - q_i \quad (1.42)$$



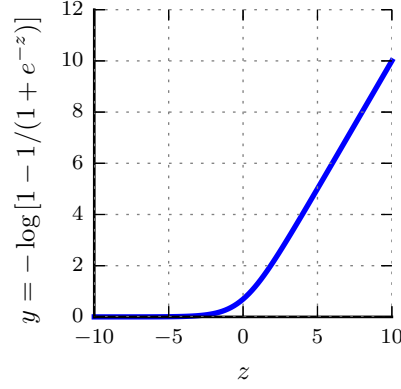
The cost function is defined as

$$C(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=0}^{\mu} c_i(\boldsymbol{\theta}) \quad (1.43)$$

where

$$\begin{aligned} c_i(\boldsymbol{\theta}) &= -y_i \log [h_i(\boldsymbol{\theta})] - (1 - y_i) \log [1 - h_i(\boldsymbol{\theta})] \\ &= y_i q_i(\boldsymbol{\theta}) + (1 - y_i) [\ell_i(\boldsymbol{\theta}) + q_i(\boldsymbol{\theta})] \\ &= y_i q_i + \ell_i + q_i - y_i \ell_i - y_i q_i \end{aligned} \quad (1.44)$$

or



$$c_i(\boldsymbol{\theta}) = q_i(\boldsymbol{\theta}) + (1 - y_i) \ell_i(\boldsymbol{\theta}) \quad (1.45)$$

The cost function can hence be written as

$$C(\boldsymbol{\theta}) = \underbrace{\frac{1}{m} \sum_{i=0}^{\mu} q_i(\boldsymbol{\theta})}_{s_q} + \sum_{i=0}^{\mu} \underbrace{\frac{1 - y_i}{m}}_{\bar{y}_i} \ell_i(\boldsymbol{\theta}) \quad (1.46)$$

or

$$C(\boldsymbol{\theta}) = s_q + \bar{\mathbf{y}}^T \boldsymbol{\ell} \quad (1.47)$$

The derivative of c_i is

$$\begin{aligned} \frac{\partial c_i}{\partial \theta_j} &= \frac{\partial q_i}{\partial \theta_j} + (1 - y_i) \frac{\partial \ell_i}{\partial \theta_j} \\ &= \left(\frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} + 1 - y_i \right) X_{ij} \\ &= \left(\frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} + \frac{1 + e^{-\ell_i}}{1 + e^{-\ell_i}} - y_i \right) X_{ij} \\ &= \left(\frac{1}{1 + e^{-\ell_i}} - y_i \right) X_{ij} \end{aligned} \quad (1.48)$$

or

$$\frac{\partial c_i}{\partial \theta_j} = (h_i(\boldsymbol{\theta}) - y_i) X_{ij} \quad (1.49)$$

The derivative of the cost function is

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{m} \sum_{i=0}^{\mu} (h_i(\boldsymbol{\theta}) - y_i) X_{ij} \quad (1.50)$$

Therefore

$$\frac{dC}{d\boldsymbol{\theta}} = \frac{1}{m} (\mathbf{h} - \mathbf{y}) \mathbf{X} \quad (1.51)$$

or

$$\frac{dC}{d\boldsymbol{\theta}} = \frac{1}{m} \mathbf{X}^T (\mathbf{h} - \mathbf{y}) \quad (1.52)$$

1.2.1 External bias parameter plus regularization

The logistic regression model with $\boldsymbol{\theta}$ and b not in the same vector is

$$h_i(\boldsymbol{\theta}, b) = g(\ell_i(\boldsymbol{\theta}, b)) \quad (1.53)$$

where ℓ_i is given by Eq. (1.14).

The auxiliary vectors are

$$p_i(\boldsymbol{\theta}, b) = 1 + e^{-\ell_i(\boldsymbol{\theta}, b)} \quad (1.54)$$

and

$$q_i(\boldsymbol{\theta}, b) = \log [p_i(\boldsymbol{\theta}, b)] \quad (1.55)$$

Thus, the cost is still given by the same expression (Eq. 1.43) with

$$c_i(\boldsymbol{\theta}, b) = q_i(\boldsymbol{\theta}, b) + (1 - y_i) \ell_i(\boldsymbol{\theta}, b) \quad (1.56)$$

Therefore, now adding regularization, we obtain

$$C(\boldsymbol{\theta}, b) = s_q(\boldsymbol{\theta}, b) + \bar{\mathbf{y}}^T \boldsymbol{\ell}(\boldsymbol{\theta}, b) + \frac{\lambda}{2m} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

with

$$s_q(\boldsymbol{\theta}, b) = \frac{1}{m} \text{sum}(\mathbf{q}) \quad \text{and} \quad \bar{y}_i = \frac{1 - y_i}{m} \quad (1.57)$$

The derivative of p_i with respect to θ_j is (no change)

$$\frac{\partial p_i}{\partial \theta_j} = -e^{-\ell_i} \frac{\partial \ell_i}{\partial \theta_j} = -e^{-\ell_i} X_{ij} \quad (1.58)$$

and the derivative of q_i with respect to θ_j is (no change)

$$\frac{\partial q_i}{\partial \theta_j} = \frac{1}{p_i} \frac{\partial p_i}{\partial \theta_j} = \frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} X_{ij} = -\frac{1}{1 + e^{\ell_i}} X_{ij} \quad (1.59)$$

The derivative of c_i with respect to θ_j is (no change)

$$\frac{\partial c_i}{\partial \theta_j} = (h_i(\boldsymbol{\theta}, b) - y_i) X_{ij} \quad (1.60)$$

Therefore, the partial derivative of the cost function with respect to $\boldsymbol{\theta}$ is (no change, except for the regularization term)

$$\frac{\partial C}{\partial \boldsymbol{\theta}} = \frac{1}{m} \mathbf{X}^T (\mathbf{h}(\boldsymbol{\theta}, b) - \mathbf{y}) + \frac{\lambda}{m} \boldsymbol{\theta}$$

The derivative of p_i with respect to b is

$$\frac{\partial p_i}{\partial b} = -e^{-\ell_i} \frac{\partial \ell_i}{\partial b} = -e^{-\ell_i} o_i \quad (1.61)$$

and the derivative of q_i with respect to b is

$$\frac{\partial q_i}{\partial b} = \frac{1}{p_i} \frac{\partial p_i}{\partial b} = \frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} o_i = -\frac{1}{1 + e^{\ell_i}} o_i \quad (1.62)$$

The derivative of c_i with respect to b is

$$\begin{aligned} \frac{\partial c_i}{\partial b} &= \frac{\partial q_i}{\partial b} + (1 - y_i) \frac{\partial \ell_i}{\partial b} \\ &= \left(\frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} + 1 - y_i \right) o_i \\ &= \left(\frac{-e^{-\ell_i}}{1 + e^{-\ell_i}} + \frac{1 + e^{-\ell_i}}{1 + e^{-\ell_i}} - y_i \right) o_i \\ &= \left(\frac{1}{1 + e^{-\ell_i}} - y_i \right) o_i \end{aligned} \quad (1.63)$$

hence

$$\frac{\partial c_i}{\partial b} = (h_i(\boldsymbol{\theta}, b) - y_i) o_i \quad (1.64)$$

Therefore, the partial derivative of the cost function with respect to b is

$$\frac{\partial C}{\partial b} = \frac{1}{m} \mathbf{o}^T (\mathbf{h}(\boldsymbol{\theta}, b) - \mathbf{y})$$

1.2.2 Hessian matrix

The first derivative of the cost function with respect to θ_i can be written as

$$\frac{\partial C}{\partial \theta_i} = \frac{1}{m} \sum_{k=0}^{\mu} X_{ik}^T e_k + \frac{\lambda}{m} \theta_i \quad (1.65)$$

where

$$e_i(\boldsymbol{\theta}, b) = h_i(\boldsymbol{\theta}, b) - y_i \quad (1.66)$$

Thus, the Hessian matrix with respect to $\boldsymbol{\theta}$ is

$$H_{ij} = \frac{\partial^2 C}{\partial \theta_i \partial \theta_j} = \frac{1}{m} \sum_{k=0}^{\mu} X_{ik}^T \frac{\partial e_k}{\partial \theta_j} + \frac{\lambda}{m} \delta_{ij} \quad (1.67)$$

The following derivative is required

$$\frac{\partial e_k}{\partial \theta_j} = \frac{\partial h_k}{\partial \theta_j} = \frac{\partial (p_k)^{-1}}{\partial \theta_j} = -\frac{1}{p_k^2} \frac{\partial p_k}{\partial \theta_j} = \frac{e^{-\ell_k}}{p_k^2} X_{kj} \quad (1.68)$$

We can show that

$$d_i \equiv \frac{e^{-\ell_i}}{p_i^2} = g(\ell_i) [1 - g(\ell_i)]$$

thus

$$\frac{\partial e_k}{\partial \theta_j} = d_k X_{kj} \quad (1.69)$$

Therefore, we obtain

$$H_{ij} = \frac{\partial^2 C}{\partial \theta_i \partial \theta_j} = \frac{1}{m} \sum_{k=0}^{\mu} X_{ik}^T d_k X_{kj} + \frac{\lambda}{m} \delta_{ij} \quad (1.70)$$

or

$$\mathbf{H} = \frac{1}{m} \mathbf{X}^T \mathbf{D} + \frac{\lambda}{m} \mathbf{I}$$

where the following matrix is defined

$$D_{ij} = d_i X_{ij}$$

We require also the cross derivative vector

$$v_i = \frac{\partial^2 C}{\partial \theta_i \partial b} = \frac{1}{m} \sum_{k=0}^{\mu} X_{ik}^T \frac{\partial e_k}{\partial b} \quad (1.71)$$

The following derivative is required

$$\frac{\partial e_k}{\partial b} = \frac{\partial h_k}{\partial b} = \frac{\partial (p_k)^{-1}}{\partial b} = -\frac{1}{p_k^2} \frac{\partial p_k}{\partial b} = \frac{e^{-\ell_k}}{p_k^2} = d_k \quad (1.72)$$

or

$$v_i = \frac{\partial^2 C}{\partial \theta_i \partial b} = \frac{1}{m} \sum_{k=0}^{\mu} X_{ik}^T d_k \quad (1.73)$$

In vector notation

$$\mathbf{v} = \frac{1}{m} \mathbf{X}^T \mathbf{d}$$

The last Hessian term is

$$\frac{\partial^2 C}{\partial b^2} = \frac{1}{m} \sum_{i=0}^{\mu} \frac{\partial h_i}{\partial b} o_i = \frac{1}{m} \sum_{i=0}^{\mu} d_i o_i \quad (1.74)$$

or

$$\frac{\partial^2 C}{\partial b^2} = \frac{1}{m} \mathbf{o}^T \mathbf{d}$$