

Principal Component Analysis with Missing Data

葉冠麟

June 16, 2006

1 PCA review

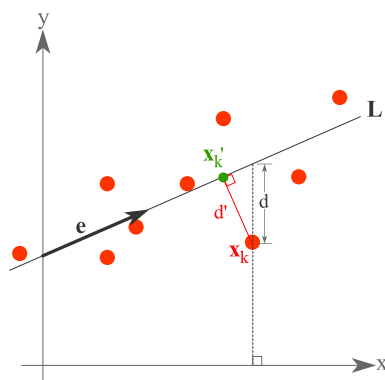


圖 1: 我們試著以一條直線 L 來表示空間中的點, 以線性迴歸的方式, 我們會累加 d 取最小值, 然而我們應該累加的是 d' , 而這正是 PCA 的處理方式。

「我們要如何用一個點來代表二維空間中散佈的點?」這是一個簡單的數學問題, 許多人可以憑著直覺正確地回答:「取平均值 (mean)。」我們再進一步地問, 如果是一條線呢? 「如何用一條線來代表二維空間中散佈的點?」用國中數學來思考這個問題, 我們會很直覺地聯想到線性規劃。線性規劃的精神在於「找一條直線 L , 使得每一個點到 L 垂直於 x 軸 (或 y 軸) 的距離 d , 累加起來最小。」我們國中的時候很滿足這樣的答案, 但顯然這樣的解是有問題的, 為什麼計算的不是垂直於直線 L 的距離 d' ? 當開始有這樣的疑慮產生, 就是引入 PCA 最佳的時機。因為 PCA 正是以 Least Squared Error 的角度來求解。

1.1 Principal Component Analysis (PCA)

(這部份的數學推導全引自 **Pattern Classification, Wiley**)

PCA的使用並不被限制在二維的空間，我們可以把問題提升到 m 維空間。我們依然在空間中灑點，這些點代表的是我們的資料，分別是 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ ，皆為 m 維的向量 (vector)。

我們先問個簡單的問題，「我們如何找到一個 m 維的向量，足以代表灑在空間中的 n 個點？」「代表」這個字詞其實用得不是很準確，比較嚴謹的問法應該是，我們如何找到 m 維的向量 \mathbf{x}_0 ，使得

$$J_0(\mathbf{x}_0) = \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{x}_k\|^2, \quad (1)$$

最小。其中 $J_0(\mathbf{x}_0)$ 是我們的 squared-error criterion function。定義 \mathbf{m} 為 n 筆資料的平均值，即

$$\mathbf{m} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k, \quad (2)$$

則

$$\begin{aligned} J_0(\mathbf{x}_0) &= \sum_{k=1}^n \|(\mathbf{x}_0 - \mathbf{m}) - (\mathbf{x}_k - \mathbf{m})\|^2 \\ &= \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{m}\|^2 - 2 \sum_{k=1}^n (\mathbf{x}_0 - \mathbf{m})^t (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{m}\|^2 - 2(\mathbf{x}_0 - \mathbf{m})^t \sum_{k=1}^n (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{m}\|^2 + \underbrace{\sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2}_{\text{independent of } \mathbf{x}_0}. \end{aligned} \quad (3)$$

其中 $\sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2$ 與 \mathbf{x}_0 無關，所以略去不看，此時若要讓 $J_0(\mathbf{x}_0)$ 最小，顯然 $\mathbf{x}_0 = \mathbf{m}$ 。這結果很符合直覺，我們的確可以用平均值來代替整組資料。

我們現在把問題變得有趣些，我們試著從 \mathbf{x}_0 拉出一條直線 \mathbf{L} ，目標依然是使 squared-error 最小。如圖 1 所示，每一點 \mathbf{x}_k 都能在直線 \mathbf{L} 上找到相對應的 \mathbf{x}'_k ，若是直線 \mathbf{L} 沿著單位向量 \mathbf{e} 走，則

$$\mathbf{x}'_k = \mathbf{m} + a_k \mathbf{e}, \quad (4)$$

其中 a_k 為縮放項。我們重新定義 squared-error criterion function ,

$$J_1(a_1, \dots, a_n, \mathbf{e}) = \sum_{k=1}^n \|(\mathbf{x}'_k - \mathbf{x}_k)\|^2 = \sum_{k=1}^n \|((\mathbf{m} + a_k \mathbf{e}) - \mathbf{x}_k)\|^2, \quad (5)$$

在 J_1 中未知的變數有 a_1, a_2, \dots, a_k , 以及 \mathbf{e} , 我們先假設 \mathbf{e} 已知, 然後看看 a_k 必須等於多少才能使得 J_1 最小。

$$\begin{aligned} J_1(a_1, \dots, a_n, \mathbf{e}) &= \sum_{k=1}^n \|((\mathbf{m} + a_k \mathbf{e}) - \mathbf{x}_k)\|^2 = \sum_{k=1}^n \|(a_k \mathbf{e} - (\mathbf{x}_k - \mathbf{m}))\|^2 \\ &= \sum_{k=1}^n a_k^2 \|\mathbf{e}\|^2 - 2 \sum_{k=1}^n a_k \mathbf{e}^t (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2. \end{aligned} \quad (6)$$

因為 $\sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2$ 與 a_k 無關, 且 \mathbf{e} 為單位向量, 即 $\|\mathbf{e}\|^2 = 1$, 我們可以另外定義

$$J'_1(a_1, \dots, a_n, \mathbf{e}) = \sum_{k=1}^n a_k^2 - 2 \sum_{k=1}^n a_k \mathbf{e}^t (\mathbf{x}_k - \mathbf{m}), \quad (7)$$

將 J'_1 分別對 a_k 做偏微分, 我們可以得到

$$a_k = \mathbf{e}^t (\mathbf{x}_k - \mathbf{m}). \quad (8)$$

這是一個很有趣的結果, 因為這代表如果你已經知道 \mathbf{e} , 將空間中任一點 \mathbf{x}_k 投射到直線 \mathbf{L} 上, 只需要將原座標為移後與 \mathbf{e}^t 做內積, 就可以得到空間轉換後的新座標 $\mathbf{x}'_k = a_k$ 。也就是說, 不論是原來存在的點, 或是之後才加入的, 都可以做相同的空間轉換。不過, 到底 \mathbf{e} 是什麼? 既然我們已經知道 a_k 該等於多少, 我們先將 J_1 改寫成較為精簡的模式

$$\begin{aligned} J_1(\mathbf{e}) &= \sum_{k=1}^n a_k^2 \|\mathbf{e}\|^2 - 2 \sum_{k=1}^n a_k^2 + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= - \sum_{k=1}^n [\mathbf{e}^t (\mathbf{x}_k - \mathbf{m})]^2 + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= - \sum_{k=1}^n \mathbf{e}^t (\mathbf{x}_k - \mathbf{m}) (\mathbf{x}_k - \mathbf{m})^t \mathbf{e} + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= - \mathbf{e}^t \mathbf{S} \mathbf{e} + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2. \end{aligned} \quad (9)$$

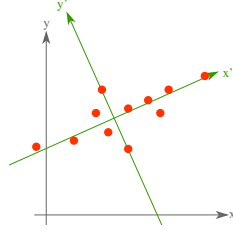


圖 2: 深灰色的座標表示原座標軸, 綠色的座標是轉換後的結果。

其中, $\mathbf{S} = \sum_{k=1}^n \mathbf{e}^t (\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^t \mathbf{e}$, 我們將之稱為 scatter matrix, 他的樣子幾乎同等於 covariance matrix, 只是少除了分母 $n - 1$ 。推導到這一步, 我們似乎只要做個微分就能求出 \mathbf{e} , 但有一點要特別注意, \mathbf{e} 是單位向量, $\|\mathbf{e}\|^2 = 1$ 的限制必須在求解時一併考慮進去。所以我們引入 Lagrange Multipliers 的方法來求 \mathbf{e} , 使得 $\mathbf{e}^t \mathbf{S} \mathbf{e}$ 最大 (為要使 J_1 最小)。令

$$u = \mathbf{e}^t \mathbf{S} \mathbf{e} - \lambda(\mathbf{e}^t \mathbf{e} - 1) \quad (10)$$

並將 u 對 \mathbf{e} 做偏微分,

$$\frac{\partial u}{\partial \mathbf{e}} = 2\mathbf{S}\mathbf{e} - 2\lambda\mathbf{e}, \quad (11)$$

令其值為 0, 則得

$$\mathbf{S}\mathbf{e} = \lambda\mathbf{e}. \quad (12)$$

這是典型的 eigen decomposition 的問題, 已超出本文討論的範圍, 我們只要知道, 原來我們所要找的 \mathbf{e} 其實就是空間中所有點形成的 covariance matrix 經過 eigen decomposition 所求得的 eigen vector。

「要如何表示空間中的點?」我們從點討論到線, 問題依然可以提升到更高的維度, 但原則就只有一個 — 對 covariance matrix 做 eigen decomposition。理論上 m 維的空間我們就可以找到 m 個 eigen vector, 我們若是對各個向量作投影, 就可以把點轉換到另一個比較符合分布情形的座標空間, 如圖 2所示。

1.2 降維

我們目前對 PCA 的認識是建立在 Least-Squared Error 上, 也就是說, 若是我們把點投射至 PCA 找到的向量上, 以這些向量為新的座標軸, 則我們可以發現投影點

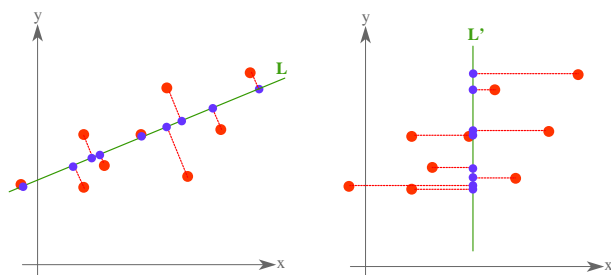


圖 3: 相較於 L' , L 會是我們比較希望看到的投影方式, 因為經過座標轉換後, 右圖的投影點有重疊, 使得我們無法單從一個維度辨別在原座標系統中不同位置的兩個點。

都會緊臨座標軸, 以致於我們原本要計較的 error, 也就是新座標軸中的座標, 值都不大。但是這樣的座標轉換對我們而言有什麼意義?

其實 PCA 名氣會這麼大, 是因為它可以拿來做資料壓縮。我們之前一直在說的「空間中的點」, 其實就只是一個向量, 他可以是一筆資料。一個簡單的想法, 一個 100 維的向量若是可以精簡成 10 維, 就可以將儲存空間壓縮成原來的 1/10。問題就在於怎麼降維才能儘可能地保留原來資料的特性。簡單的說, 維度降下來以後, 我們希望原本分開的點仍然要分開, 就如圖 3 所示, 左側圖示那樣的投影才是理想的投影。

「為什麼右側的投影結果會比較差?」圖 3 右側中的投影之所以會重疊, 直覺的想法, 是因為座標軸之間有相關, 以致於拿掉其中一個軸, 就會使得座標點無法單獨被定義。不過這在 PCA 是不成問題的, 因為 PCA 各軸之間是相互獨立的。我們底下就來證明這件事。

首先, 給定空間中散佈的 n 個 m 維的點, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, 其平均值

$$\mathbf{m} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k. \quad (13)$$

令 $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 - \mathbf{m} & \mathbf{x}_2 - \mathbf{m} & \dots & \mathbf{x}_n - \mathbf{m} \end{bmatrix}^T$ 。則 covariance matrix

$$\mathbf{S} = \frac{\mathbf{X} \mathbf{X}^T}{n-1}. \quad (14)$$

我們接著對 \mathbf{S} 做 eigen decomposition, 得到 m 個 eigen vector $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_m$, 以及相對應的 m 個 eigen value $\lambda_1, \lambda_2, \dots, \lambda_m$ 。令 $\mathbf{E} = \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_m \end{bmatrix}$, 則

$$\begin{aligned}
\mathbf{S}\mathbf{E} &= \begin{bmatrix} \mathbf{S}\mathbf{e}_1 & \mathbf{S}\mathbf{e}_2 & \dots & \mathbf{S}\mathbf{e}_m \end{bmatrix} \\
&= \begin{bmatrix} \lambda_1 \mathbf{e}_1 & \lambda_2 \mathbf{e}_2 & \dots & \lambda_m \mathbf{e}_m \end{bmatrix} \\
&= \begin{bmatrix} \lambda_1 \mathbf{e}_{11} & \lambda_2 \mathbf{e}_{21} & \dots & \lambda_m \mathbf{e}_{m1} \\ \lambda_1 \mathbf{e}_{12} & \lambda_2 \mathbf{e}_{22} & \dots & \lambda_m \mathbf{e}_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1 \mathbf{e}_{1m} & \lambda_2 \mathbf{e}_{2m} & \dots & \lambda_m \mathbf{e}_{mm} \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{e}_{11} & \mathbf{e}_{21} & \dots & \mathbf{e}_{m1} \\ \mathbf{e}_{12} & \mathbf{e}_{22} & \dots & \mathbf{e}_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{e}_{1m} & \mathbf{e}_{2m} & \dots & \mathbf{e}_{mm} \end{bmatrix} \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_m \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_m \end{bmatrix} \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_m \end{bmatrix} \\
&= \mathbf{E}\mathbf{D} \\
\Rightarrow \mathbf{S} &= \mathbf{E}\mathbf{D}\mathbf{E}^{-1} \tag{15}
\end{aligned}$$

其中

$$\mathbf{D} = \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_m \end{bmatrix} \tag{16}$$

我們現在把空間中的點投射到 eigen vector 所張出來的空間, 則每個 \mathbf{x}_k 都會得到新的座標 \mathbf{y}_k , 令 $\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1 & \mathbf{y}_2 & \dots & \mathbf{y}_n \end{bmatrix}^T$ 則

$$\mathbf{Y} = \mathbf{E}^T \mathbf{X}, \tag{17}$$

而新座標所形成的 covariance matrix

$$\begin{aligned}
 S_Y &= \frac{Y Y^T}{n-1} \\
 &= \frac{1}{n-1} E^T X X^T E \\
 &= E^T \left(\frac{X X^T}{n-1} \right) E \\
 &= E^T S E \\
 &= E^T E D E^{-1} E = D
 \end{aligned} \tag{18}$$

這顯示了一個珍貴的情報：當我們做完座標轉換，covariance matrix 會形成一個對角矩陣，也就是說，座標系統中任兩軸間的相關度是 0！

不過這就代表我們可以隨便挑掉幾個軸來降維嗎？當然不是！我們希望取自相關度高的軸，因為這表示該軸對整組資料有比較高的代表性。這時候我們回頭去看一下矩陣 D ，也就是投影過後的 covariance matrix，對角線上擺得其實就是 eigen value，也就是說，要計較座標軸的 variance，其實就只要看 eigen value！所以事情變得再簡單不過了，當我們做完 eigen decomposition，我們看要把資料降成幾維，假設是 k 維，就挑相對應的 eigen value 最大的 k 個 eigen vector，然後把資料點一一投影上去，就得到新的座標值。

1.3 實作

假設我們有 n 筆 m 維的資料，分別是 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ ，我們的目標是把筆資料降維成 k 維。底下是實作的步驟。

1. 求出 covariance matrix。

令 $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{bmatrix}^T$ ， $\mathbf{m} = \sum_{i=1}^n \mathbf{x}_i$ 。計算 $\mathbf{S} = \frac{1}{n-1}(\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^t$ 。

2. 對 \mathbf{S} 做 eigen decomposition。

得到 m 個 eigen vector $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ ，以及相對應的 eigen value $\lambda_1, \lambda_2, \dots, \lambda_n$

。

3. 選出前 k 大的 eigen value 所對應的 eigen vector , e_1, e_2, \dots, e_k 。並投影資料點。

假設每一筆資料 \mathbf{x}_i 投影後得到 $\mathbf{x}'_i = \begin{bmatrix} \mathbf{x}'_{i1} & \mathbf{x}'_{i2} & \dots & \mathbf{x}'_{ik} \end{bmatrix}$, 則 $\mathbf{x}'_{ij} = \mathbf{e}_j^t(\mathbf{x}_i - \mathbf{m})$ 。

2 PCA with missing data

前面會花這麼多篇幅介紹 PCA , 是因為若是不討論 PCA , 我們要探討的主題 PCA with missing data 就不完整。PCA with missing data 嚴格來說只多引入了一個步驟, 其餘的與 PCA 如出一轍。PCA with missing data 要多處理的一步, 就如同它名稱所揭示的, 就是資料遺失的問題。也就是說我們可能拿到如

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & * & * & \mathbf{x}_5 & * \end{bmatrix} \quad (19)$$

這樣的資料。因為資料通常不容易收集, 我們不可能因為資料不完整就全部作廢。但很明顯的, 其餘的處理不說, 光是 PCA 就會遇上麻煩, 因為無法拿這樣的資料去算 covariance matrix 。所以 PCA with missing data 多出的一步就是要還原資料。最偷懶的方式, 就是直接把遺失的資料代換成現存資料的平均值, 不過當資料遺失過多, 平均值似乎也不能代表什麼。所以我們底下要介紹一個方法能比較正確的補回資料。

2.1 Wiberg's method

跟之前一樣, 假設我們有 n 筆 m 維的資料, 分別是 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, 令 $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{bmatrix}^T \in \Re^{m \times n}$, 只是我們目前無法直接拿 \mathbf{X} 來用, 因為裡頭有許多的「洞」, 也就是遺失的資料。給個具體的例子, 我們很可能會需要處理像

$$\begin{bmatrix} \mathbf{x}_1^{(1)} & \mathbf{x}_2^{(1)} & \mathbf{x}_3^{(1)} & \mathbf{x}_4^{(1)} & \mathbf{x}_5^{(1)} & \mathbf{x}_6^{(1)} & * & * & * & * & * & * \\ \mathbf{x}_1^{(2)} & \mathbf{x}_2^{(2)} & \mathbf{x}_3^{(2)} & \mathbf{x}_4^{(2)} & * & * & \mathbf{x}_7^{(2)} & \mathbf{x}_8^{(2)} & * & * & * & * \\ \mathbf{x}_1^{(3)} & \mathbf{x}_2^{(3)} & * & * & * & \mathbf{x}_6^{(3)} & * & \mathbf{x}_8^{(3)} & \mathbf{x}_9^{(3)} & \mathbf{x}_{10}^{(3)} & * & * \\ * & * & * & * & * & * & \mathbf{x}_7^{(4)} & \mathbf{x}_8^{(4)} & \mathbf{x}_9^{(4)} & \mathbf{x}_{10}^{(4)} & \mathbf{x}_{11}^{(4)} & \mathbf{x}_{12}^{(4)} \end{bmatrix}$$

這樣坑坑洞洞的矩陣。我們的任務就是要找到一個矩陣 \mathbf{Y} , 使得它可以很把洞填滿, 而且盡量近似於原來的資料。不過什麼叫近似原來的資料? 資料都遺失了, 我們怎麼

知道他原來長什麼樣子？我們目前唯一能要求的，就是有的資料希望能夠保留下來。那沒有的資料呢？我們就必須使用類似 em 的方法來求解。

因為資料是殘缺的，我們沒有辦法得到資料的平均值，為了與之前確切的平均值作區分，我們以 $\boldsymbol{\mu}$ 來表示，代表平均值的 maximum likelihood approximation。於是，我們的問題以比較正式的方式寫下，就是要求 $\mathbf{Y} \in \Re^{m \times n}$ ，使得

$$\Phi = \|\mathbf{X} - \mathbf{e}^T \boldsymbol{\mu} - \mathbf{Y}\| \quad (20)$$

最小，其中 $\mathbf{e} = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}^T$ 。在之前的討論中，我們會先把平均值減去，為的是免去位移向量的動作，不過現在平均值未知，所以必須當作未知數一併求解。

SVD¹ 告訴我們，任何一個 $m \times n$ 的矩陣 $\mathbf{M} \in \Re^{m \times n}$ 都可以分解成

$$\mathbf{M} = \mathbf{U} \mathbf{S} \mathbf{V}^T, \quad (21)$$

其中 \mathbf{U} 和 \mathbf{V} 分別是 $m \times r$ 和 $n \times r$ 的正交矩陣， $\mathbf{S} = \text{diag}(\sigma_i)$ 是一個 $r \times r$ 的對角矩陣。當然， \mathbf{Y} 也可以作 SVD 分解，假設

$$\mathbf{Y} = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T, \quad (22)$$

我們另外令

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_m \end{bmatrix}^T = \tilde{\mathbf{U}} \tilde{\mathbf{S}}^{\frac{1}{2}}, \quad (23)$$

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_n \end{bmatrix}^T = \tilde{\mathbf{V}} \tilde{\mathbf{S}}^{\frac{1}{2}}, \quad (24)$$

$$\Rightarrow \mathbf{u}^T \mathbf{v} = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^T = \mathbf{Y}. \quad (25)$$

所以我們就不需要去解整個 \mathbf{Y} 矩陣，取而代之的，我們只要計算 $\mathbf{u} \in \Re^m$ 和 $\mathbf{v} \in \Re^n$ 兩個向量，於是原來 $m \times n$ 的大矩陣，只剩下 $m + n$ 個變數。

所以現在我們要求的有三個向量，分別是 $\mathbf{u} \in \Re^m$ ， $\mathbf{v} \in \Re^n$ 和 $\boldsymbol{\mu} \in \Re^m$ ，而我們已知的只有部分的資料。當然，遺失的資料已經超出我們掌控的範圍，我們只能根據現存的資料來求解。所以我們的問題需要改寫成

$$\min \phi = \frac{1}{2} \sum_I (\mathbf{X}_{mn} - \boldsymbol{\mu}_n - \mathbf{u}_m^T \mathbf{v}_n)^2, \quad (26)$$

$$I = \{(m, n) : \mathbf{X}_{mn} \text{ is observed}\}.$$

¹詳情可以參考 <http://www.uwlax.edu/faculty/will/svd/index.html>

與 Φ 相比, 我們現在要對每一個在 \mathbf{X} 中有值的 entry, 去累加他的誤差, 並設法使得誤差最小。附帶一提, 其中的 $\mathbf{u}_m^T \mathbf{v}_n$ 其實就是 \mathbf{Y}_{mn} 。

既然我們只使用部分的資料, 且是各個 entry 獨立的計算誤差, 我們實在沒有必要維持原來矩陣的形式, 我們可以把有值的 entry 拉成一條向量, 舉個例子,

$$\begin{bmatrix} 1 & * & 2 \\ 3 & * & * \\ * & 4 & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}.$$

顯而易見的, 若是 \mathbf{Y} 中有 p 個 entry 是有值的, 我們就可以得到一向量 $\mathbf{w} = [\Omega_1 \ \dots \ \Omega_{|\Omega|}] \in \mathbb{R}^p$, 其中 $\Omega = \{\omega = X_{ij} | X_{ij} \text{ is observed}\}$ 。當然, 經過這樣的調整, $\boldsymbol{\mu}$ 也必須改變, 定義

$$\hat{\boldsymbol{\mu}} = \begin{bmatrix} \hat{\boldsymbol{\mu}}_1 \\ \vdots \\ \hat{\boldsymbol{\mu}}_{|\Omega|} \end{bmatrix}, \hat{\boldsymbol{\mu}}_k = \boldsymbol{\mu}_i \text{ if } \Omega_k \in \mathbf{X}_i. \quad (27)$$

承接上一個例子, 若是該矩陣相對應的 $\boldsymbol{\mu}$ 為平均值, 且假設 $\boldsymbol{\mu} = [1 \ 2 \ 3]^T$, 則 $\hat{\boldsymbol{\mu}} = [1 \ 1 \ 2 \ 3 \ 3]^T$ 。基本上, 我們目前所做的努力無非是想讓原來的 ϕ 在拿掉遺失的資料後, 仍可以維持原來的等式。所以我們也必須對 \mathbf{v} 和 \mathbf{u} 做點小修正, 定義

$$\hat{\mathbf{v}} = \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_n \end{bmatrix} \in \mathbb{R}^{r \times 1}, \hat{\mathbf{u}} = \begin{bmatrix} \hat{\mathbf{u}}_1 \\ \dots \\ \hat{\mathbf{u}}_m \end{bmatrix} \in \mathbb{R}^{(r+1) \times 1}, \hat{\mathbf{u}}_i = \begin{bmatrix} \mathbf{u}_i \\ \boldsymbol{\mu}_i \end{bmatrix} \quad (28)$$

, 最後我們引入兩個矩陣 $\mathbf{B} \in \mathbb{R}^{|\Omega| \times rn}$ 和 $\mathbf{G} \in \mathbb{R}^{|\Omega| \times (r+1)m}$, 使得

$$\phi = \frac{1}{2} \mathbf{f}^T \mathbf{f}. \quad (29)$$

$$\mathbf{f} = \mathbf{w} - \hat{\boldsymbol{\mu}} - \mathbf{B}\hat{\mathbf{v}} = \mathbf{w} - \mathbf{G}\hat{\mathbf{u}} \quad (30)$$

定義

$$B = \begin{bmatrix} B_{11} & B_{12} & \dots & B_{1n} \\ B_{21} & B_{22} & \dots & B_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ B_{|\Omega|1} & B_{|\Omega|2} & \dots & B_{|\Omega|n} \end{bmatrix},$$

$$B_{ij} = \begin{cases} \mathbf{u}_m & \text{if } \mathbf{w}_i = \mathbf{X}_{mn} \text{ and } j = n, \\ 0 & \text{otherwise.} \end{cases} \in \mathbb{R}^{|\mathbf{u}|} \quad (31)$$

$$G = \begin{bmatrix} G_{11} & G_{12} & \dots & G_{1m} \\ G_{21} & G_{22} & \dots & G_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ G_{|\Omega|1} & G_{|\Omega|2} & \dots & G_{|\Omega|m} \end{bmatrix},$$

$$G_{ij} = \begin{cases} \begin{bmatrix} \mathbf{v}_n & -1 \end{bmatrix} & \text{if } \mathbf{w}_i = \mathbf{X}_{mn} \text{ and } j = m, \\ 0 & \text{otherwise.} \end{cases} \in \mathbb{R}^{|\mathbf{v}+1|} \quad (32)$$

爲了求 ϕ 最小，我們將 ϕ 分別對 $\hat{\mathbf{u}}$ 和 $\hat{\mathbf{v}}$ 做偏微分。令

$$\frac{\partial \phi}{\partial \hat{\mathbf{u}}} = \mathbf{B}^T \mathbf{B} \hat{\mathbf{v}} - \mathbf{B}^T (\mathbf{w} - \hat{\boldsymbol{\mu}}) = 0 \quad (33)$$

$$\frac{\partial \phi}{\partial \hat{\mathbf{v}}} = \mathbf{G}^T \mathbf{G} \hat{\mathbf{u}} - \mathbf{B}^T \mathbf{w} = 0 \quad (34)$$

$$\Rightarrow \hat{\mathbf{v}} = \mathbf{B}^+ (\mathbf{w} - \hat{\boldsymbol{\mu}}) \quad (35)$$

$$\hat{\mathbf{u}} = \mathbf{G}^+ \mathbf{w} \quad (36)$$

其中 \mathbf{B}^+ 和 \mathbf{G}^+ 分別是 \mathbf{B} 和 \mathbf{G} 的 pseudo-inverse 。

所以若是要實作，我們必須先猜測一組 $\hat{\mathbf{u}}, \hat{\mathbf{v}}$ 和 $\hat{\boldsymbol{\mu}}$ 並分成兩邊，一邊是 $\hat{\mathbf{u}}$ 和 $\hat{\boldsymbol{\mu}}$ ，另一邊則單放 $\hat{\mathbf{v}}$ 。一開始先固定其中一邊，假設是 $\hat{\mathbf{v}}$ ，則我們可以兜出 \mathbf{G} ，所以可以求出 $\hat{\mathbf{u}}$ 。反之，固定另一邊，則可以求出 $\hat{\mathbf{v}}$ 。反覆相同的步驟，直到收斂爲止。當我們求出了 $\hat{\mathbf{u}}$ 和 $\hat{\mathbf{v}}$ ，資料重建的工作也就大功告成了，之後就就可以進行 PCA 的正常程序，所以就不贅述了。

2.2 實作

1. 起始 $\hat{\mathbf{v}}$ 。

2. 更新 $\hat{\boldsymbol{v}} = \boldsymbol{B}^+(\boldsymbol{w} - \hat{\boldsymbol{\mu}})$
3. 更新 $\hat{\boldsymbol{u}} = \boldsymbol{G}^+\boldsymbol{w}$
4. 若是收斂就停止, 否則回到步驟 2。