

Factor Analysis and Kalman Filtering (11/2/04)

Lecturer: Michael I. Jordan

Scribes: Byung-Gon Chun and Sunghoon Kim

1 Factor Analysis

Factor analysis is used for dimensionality reduction. It finds a subspace where the data lie. Figure 1 shows the geometry of the factor analysis model. In the figure, μ is the mean or the centroid of manifold, Λ is the coordinate system, and Ψ is the noise variance (represented by the sphere in the figure). The latent Gaussian variable $X_n \sim \mathcal{N}(0, I)$, and the observed variable $Y_n \sim \mathcal{N}(\mu + \Lambda x_n, \Psi)$ where Ψ is a diagonal covariance matrix. When λ_1 and λ_2 are basis vectors, $\Lambda = (\lambda_1 \ \lambda_2)$. $\hat{\mu} + \hat{\Lambda} \hat{x}_n$ is the expectation of Y_n .

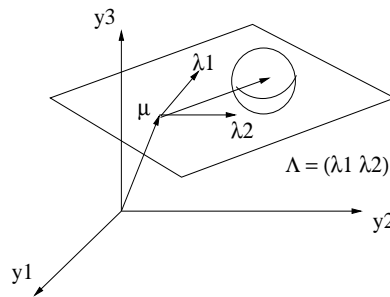


Figure 1: The geometry of the factor analysis model.

Why do we need dimensionality reduction? One application of factor analysis is compression. In the k-means algorithm, we can represent data with centroid coordinates and residuals. Likewise, with factor analysis, we can compress data using the coordinates of the subspace and residuals. Also, it can be used to predict regression. In high dimensions, there can be a large amount of noise. With dimensionality reduction, we can perform prediction with a small number of states, in effect reducing the impact of the noise.

2 Principal Component Analysis

Principal component analysis (PCA) finds the direction of maximal variance, and is a more simplistic model than factor analysis. We have data $\mathcal{D} = \{y_1, \dots, y_N\}$. First, we center the data, that is, subtract the sample mean from data. Let the modified data be \tilde{y}_n . $\tilde{y}_n = y_n - \frac{1}{N} \sum_{n=1}^N y_n$. Then, we find θ to maximize $J = \sum_{n=1}^N (\theta^T \tilde{y}_n)^2$ such that $\|\theta\| = 1$. This is a simple optimization problem. Let's manipulate J to achieve θ to maximize J .

$$J = \sum_{n=1}^N (\theta^T \tilde{y}_n)^2 \tag{1}$$

$$= \sum_{n=1}^N (\theta^T \tilde{y}_n \tilde{y}_n^T \theta) \quad (2)$$

$$= \theta^T \left(\sum_{n=1}^N \tilde{y}_n \tilde{y}_n^T \right) \theta \quad (3)$$

$$= \theta^T Y \theta \quad (4)$$

$\hat{\theta}$ that maximizes J meets $Y\hat{\theta} = \lambda\hat{\theta}$. As you can see, when λ is the largest eigenvalue, $\hat{\theta}$ is the first eigenvector.

We compare factor analysis (FA) with principal component analysis (PCA). FA is insensitive to scale factors. Since FA captures the underlying coordinate systems, FA does not change the direction of projections when Y^i has different variance. The uncertainty is given by covariance matrix in FA. When Y_n is far off from the subspace, the covariance increases and the sphere in Figure 1 grows. On the contrary, PCA is sensitive to scale factors, since the eigenvectors are affected by scale factors. Whether to use FA or PCA depends on goals. If the goal is compression, the scale factor does not matter much, and PCA can be effective.

3 Maximum Likelihood Estimation of the Factor Analysis Model

In this section, we derive the maximum likelihood estimation of the factor analysis model. $\hat{\mu}_{ML}$ is the sample mean and we subtract $\hat{\mu}_{ML}$ from data. The graphical model is presented in Figure 2. In the model, unconditional mean and unconditional variance are as follows.

$$\text{unconditional mean} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (5)$$

$$\text{unconditional covariance} = \begin{pmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{pmatrix} \quad (6)$$

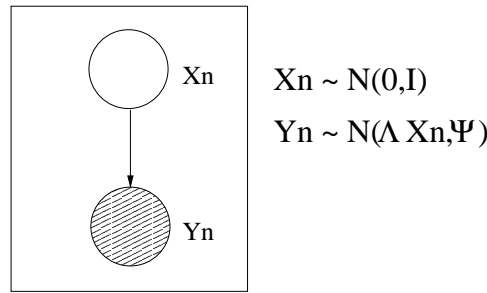


Figure 2: The factor analysis model as a graphical model.

Given complete data, the complete log likelihood is a product of Gaussian distributions. The complete log likelihood is:

$$\ell_c(\Lambda, \Psi) = -\frac{N}{2} \log |\Psi| - \frac{N}{2} \text{tr}(S\Psi^{-1}) \quad (7)$$

$$\text{where } S = \frac{1}{N} \sum_{n=1}^N (Y_n - \Lambda X_n)(Y_n - \Lambda X_n)^T \quad (8)$$

Now, we take the conditional expectation of the complete log likelihood, conditioning on the observed data and the current parameter vector. $\langle \cdot \rangle$ denotes the conditional expectation. The expected complete log

likelihood is:

$$\langle \ell_c(\Lambda, \Psi) \rangle = -\frac{N}{2} \log |\Psi| - \frac{N}{2} \text{tr}(\langle S \rangle \Psi^{-1}) \quad (9)$$

$$\langle S \rangle = \frac{1}{N} \sum_{n=1}^N \{Y_n Y_n^T - Y_n \langle X_n^T \rangle \Lambda^T - \Lambda \langle X_n \rangle Y_n^T + \Lambda \langle X_n X_n^T \rangle \Lambda^T\} \quad (10)$$

In the E-step of the algorithm, sufficient statistics we need are $\langle X_n \rangle$ and $\langle X_n X_n^T \rangle$.

$$\langle X_n \rangle = E[X_n | Y_n] = \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} Y_n \quad (11)$$

$$\langle X_n X_n^T \rangle = \text{Var}(X_n | Y_n) + E[X_n | Y_n] E[X_n | Y_n]^T = I - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda + E[X_n | Y_n] E[X_n | Y_n]^T \quad (12)$$

In the M-step of the algorithm, we update Λ and Ψ using the updated sufficient statistics.

$$\Psi^{(t+1)} = \left(\sum_{n=1}^N Y_n \langle X_n^T \rangle \right) \left(\sum_{n=1}^N \langle X_n X_n^T \rangle \right)^{-1} \quad (13)$$

Recall $\hat{\theta} = (X^T X)^{-1} X^T Y$.

$$\Psi^{(t+1)} = \frac{1}{N} \text{diag} \left\{ \sum_{n=1}^N (Y_n - \Lambda^{(t+1)} \langle X_n \rangle) Y_n^T \right\} \quad (14)$$

Here, $\Lambda^{(t+1)} \langle X_n \rangle$ is the best guess and $Y_n - \Lambda^{(t+1)} \langle X_n \rangle$ is the residual. In the above, we derive the E-step and the M-step for estimation.

4 State Space Models

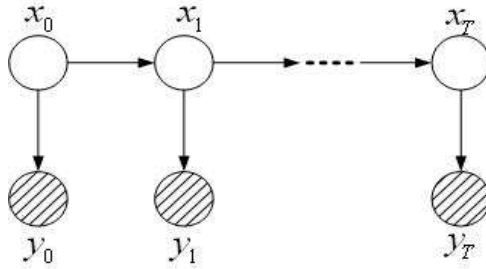


Figure 3: The SSM as a graphical model

The state nodes in the factor analysis model are continuous, vector-valued nodes endowed with a Gaussian probability distribution. To develop a dynamical generalization of the factor analysis model we must represent the transition between the nodes at successive moments in time. Perhaps the simplest choice that we can make is to allow the mean of the state at time $t + 1$ to be a linear function of the state at time t .

$$x_{t+1} = Ax_t + Gw_t$$

where w_t is a "noise" term — a Gaussian random variable that is independent of w_s for $s < t$, and thus independent of x_t . We assume that w_t has zero mean and covariance matrix Q and endow the initial state,

probability distribution $p(x_t | y_0, \dots, y_t)$ must be Gaussian.

We write $\hat{x}_{t|t}$ to denote the mean of x_t conditioned on the partial sequence y_0, \dots, y_t . The covariance matrix of x_t conditioned on y_0, \dots, y_t is denoted $P_{t|t}$; thus:

$$\begin{aligned}\hat{x}_{t|t} &\triangleq E[x_t | y_0, \dots, y_t] \\ P_{t|t} &\triangleq E[(x_t - \hat{x}_{t|t})(x_t - \hat{x}_{t|t})^T | y_0, \dots, y_t]\end{aligned}$$

We assume that we have already calculated $p(x_t | y_0, \dots, y_t)$ that is, we have calculated $\hat{x}_{t|t}$ and $p(x_t | y_0, \dots, y_t)$. We wish to carry this distribution forward into the fragment on the right, where we condition on y_0, \dots, y_{t+1} . We decompose the transformation into two steps:

$$\begin{aligned}\text{time update:} & \quad p(x_t | y_0, \dots, y_t) \rightarrow p(x_{t+1} | y_0, \dots, y_t) \\ \text{measurement update:} & \quad p(x_{t+1} | y_0, \dots, y_t) \rightarrow p(x_{t+1} | y_0, \dots, y_{t+1})\end{aligned}$$

Thus, in the time update step, we simply propagate the distribution forward one step in time, calculating the new mean and covariance based on the old mean and covariance, but conditioning on no new measurements (i.e., no new output). In the measurement update step, we incorporate the new measurement y_{t+1} and update the probability distribution for x_{t+1} . The overall result is a transformation from $\hat{x}_{t|t}$ and $P_{t|t}$ to $\hat{x}_{t+1|t+1}$ and $P_{t+1|t+1}$.

Given $\hat{x}_{t|t}$, compute $\hat{x}_{t+1|t}$

$$\begin{aligned}\hat{x}_{t+1|t} &= E[x_{t+1} | y_0, \dots, y_t] \\ &= E[Ax_t + Gw_t | y_0, \dots, y_t] \\ &= AE[x_t | y_0, \dots, y_t] \\ &= A\hat{x}_{t|t}\end{aligned}$$