

Chapter 12

Latent linear models

12.1 Factor analysis

One problem with mixture models is that they only use a single latent variable to generate the observations. In particular, each observation can only come from one of K prototypes. One can think of a mixture model as using K hidden binary variables, representing a one-hot encoding of the cluster identity. But because these variables are mutually exclusive, the model is still limited in its representational power.

An alternative is to use a vector of real-valued latent variables, $z_i \in \mathbb{R}^L$. The simplest prior to use is a Gaussian (we will consider other choices later):

$$p(z_i) = \mathcal{N}(z_i | \mu_0, \Sigma_0) \quad (12.1)$$

If the observations are also continuous, so $x_i \in \mathbb{R}^D$, we may use a Gaussian for the likelihood. Just as in linear regression, we will assume the mean is a linear function of the (hidden) inputs, thus yielding

$$p(x_i | z_i, \theta) = \mathcal{N}(x_i | \mathbf{W}z_i + \mu, \Psi) \quad (12.2)$$

where \mathbf{W} is a $D \times L$ matrix, known as the **factor loading matrix**, and Ψ is a $D \times D$ covariance matrix. We take Ψ to be diagonal, since the whole point of the model is to force z_i to explain the correlation, rather than baking it in to the observations covariance. This overall model is called **factor analysis** or **FA**. The special case in which $\Psi = \sigma^2 \mathbf{I}$ is called **probabilistic principal components analysis** or **PPCA**. The reason for this name will become apparent later.

12.1.1 FA is a low rank parameterization of an MVN

FA can be thought of as a way of specifying a joint density model on x using a small number of parameters. To see this, note that from Equation 4.39, the induced marginal distribution $p(x_i | \theta)$ is a Gaussian:

$$\begin{aligned} p(x_i | \theta) &= \int \mathcal{N}(x_i | \mathbf{W}z_i + \mu, \Psi) \mathcal{N}(z_i | \mu_0, \Sigma_0) dz_i \\ &= \mathcal{N}(x_i | \mathbf{W}\mu_0 + \mu, \Psi + \mathbf{W}\Sigma_0\mathbf{W}) \end{aligned} \quad (12.3)$$

From this, we see that we can set $\mu_0 = 0$ without loss of generality, since we can always absorb $\mathbf{W}\mu_0$ into μ . Similarly, we can set $\Sigma_0 = \mathbf{I}$ without loss of generality, because we can always emulate a correlated prior by using defining a new weight matrix, $\tilde{\mathbf{W}} = \mathbf{W}\Sigma_0^{-\frac{1}{2}}$. So we can rewrite Equation 12.6 and 12.2 as:

$$p(z_i) = \mathcal{N}(z_i | \mathbf{0}, \mathbf{I}) \quad (12.4)$$

$$p(x_i | z_i, \theta) = \mathcal{N}(x_i | \tilde{\mathbf{W}}z_i + \mu, \Psi) \quad (12.5)$$

We thus see that FA approximates the covariance matrix of the visible vector using a low-rank decomposition:

$$\mathbf{C} \triangleq \text{cov}[x] = \mathbf{W}\mathbf{W}^T + \Psi \quad (12.6)$$

This only uses $O(LD)$ parameters, which allows a flexible compromise between a full covariance Gaussian, with $O(D^2)$ parameters, and a diagonal covariance, with $O(D)$ parameters. Note that if we did not restrict Ψ to be diagonal, we could trivially set Ψ to a full covariance matrix; then we could set $\mathbf{W} = 0$, in which case the latent factors would not be required.

12.1.2 Inference of the latent factors

$$p(z_i | x_i, \theta) = \mathcal{N}(z_i | \mu_i, \Sigma_i) \quad (12.7)$$

$$\Sigma_i \triangleq (\Sigma_0^{-1} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \quad (12.8)$$

$$= (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \quad (12.9)$$

$$\mu_i \triangleq \Sigma_i [\mathbf{W}^T \Psi^{-1} (x_i - \mu) + \Sigma_0^{-1} \mu_0] \quad (12.10)$$

$$= \Sigma_i \mathbf{W}^T \Psi^{-1} (x_i - \mu) \quad (12.11)$$

Note that in the FA model, Σ_i is actually independent of i , so we can denote it by Σ . Computing this matrix takes $O(L^3 + L^2D)$ time, and computing each $\mu_i = \mathbb{E}[z_i | x_i, \theta]$ takes $O(L^2 + LD)$ time. The μ_i are sometimes called the **latent scores**, or **latent factors**.

12.1.3 Unidentifiability

Just like with mixture models, FA is also unidentifiable. To see this, suppose \mathbf{R} is an arbitrary orthogonal rotation matrix, satisfying $\mathbf{R}\mathbf{R}^T = \mathbf{I}$. Let us define $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$, then the likelihood function of this modified matrix is the same as for the unmodified matrix, since $\mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T + \mathbf{\Psi} = \mathbf{W}\mathbf{W}^T + \mathbf{\Psi}$. Geometrically, multiplying \mathbf{W} by an orthogonal matrix is like rotating \mathbf{z} before generating \mathbf{x} .

To ensure a unique solution, we need to remove $L(L-1)/2$ degrees of freedom, since that is the number of orthonormal matrices of size $L \times L$.²² In total, the FA model has $D + DL - L(L-1)/2$ free parameters (excluding the mean), where the first term arises from $\mathbf{\Psi}$. Obviously we require this to be less than or equal to $D(D+1)/2$, which is the number of parameters in an unconstrained (but symmetric) covariance matrix. This gives us an upper bound on L , as follows:

$$L_{\max} = \lfloor D + 0.5(1 - \sqrt{1 + 8D}) \rfloor \quad (12.12)$$

For example, $D = 6$ implies $L \leq 3$. But we usually never choose this upper bound, since it would result in overfitting (see discussion in Section 12.3 on how to choose L).

Unfortunately, even if we set $L < L_{\max}$, we still cannot uniquely identify the parameters, since the rotational ambiguity still exists. Non-identifiability does not affect the predictive performance of the model. However, it does affect the loading matrix, and hence the interpretation of the latent factors. Since factor analysis is often used to uncover structure in the data, this problem needs to be addressed. Here are some commonly used solutions:

- **Forcing \mathbf{W} to be orthonormal** Perhaps the cleanest solution to the identifiability problem is to force \mathbf{W} to be orthonormal, and to order the columns by decreasing variance of the corresponding latent factors. This is the approach adopted by PCA, which we will discuss in Section 12.2. The result is not necessarily more interpretable, but at least it is unique.
- **Forcing \mathbf{W} to be lower triangular** One way to achieve identifiability, which is popular in the Bayesian community (e.g., (Lopes and West 2004)), is to ensure that the first visible feature is only generated by the first latent factor, the second visible feature is only generated by the first two latent factors, and so on. For example, if $L = 3$ and $D = 4$, the correspond factor loading matrix is given by

$$\mathbf{W} = \begin{pmatrix} w_{11} & 0 & 0 \\ w_{21} & w_{22} & 0 \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{pmatrix}$$

We also require that $w_{jj} > 0$ for $j = 1 : L$. The total number of parameters in this constrained matrix is $D + DL - L(L-1)/2$, which is equal to the number of uniquely identifiable parameters. The disadvantage of this method is that the first L visible variables, known as the **founder variables**, affect the interpretation of the latent factors, and so must be chosen carefully.

- **Sparsity promoting priors on the weights** Instead of pre-specifying which entries in \mathbf{W} are zero, we can encourage the entries to be zero, using ℓ_1 regularization (Zou et al. 2006), ARD (Bishop 1999; Archambeau and Bach 2008), or spike-and-slab priors (Ratnayake et al. 2009). This is called sparse factor analysis. This does not necessarily ensure a unique MAP estimate, but it does encourage interpretable solutions. See Section 13.8 TODO.
- **Choosing an informative rotation matrix** There are a variety of heuristic methods that try to find rotation matrices \mathbf{R} which can be used to modify \mathbf{W} (and hence the latent factors) so as to try to increase the interpretability, typically by encouraging them to be (approximately) sparse. One popular method is known as **varimax** (Kaiser 1958).
- **Use of non-Gaussian priors for the latent factors** In Section 12.6, we will discuss how replacing $p(\mathbf{z}_i)$ with a non-Gaussian distribution can enable us to sometimes uniquely identify \mathbf{W} as well as the latent factors. This technique is known as ICA.

12.1.4 Mixtures of factor analysers

The FA model assumes that the data lives on a low dimensional linear manifold. In reality, most data is better modeled by some form of low dimensional *curved* manifold. We can approximate a curved manifold by a piecewise linear manifold. This suggests the following model: let the k 'th linear subspace of dimensionality L_k be represented by \mathbf{W}_k , for $k = 1 : K$. Suppose we have a latent indicator $q_i \in \{1, \dots, K\}$ specifying which subspace we should use to generate the data. We then sample \mathbf{z}_i from a Gaussian prior and pass it through the \mathbf{W}_k matrix (where $k = q_i$), and add noise. More precisely, the model is as follows:

$$p(q_i | \boldsymbol{\theta}) = \text{Cat}(q_i | \boldsymbol{\pi}) \quad (12.13)$$

$$p(\mathbf{z}_i | \boldsymbol{\theta}) = \mathcal{N}(\mathbf{z}_i | \mathbf{0}, \mathbf{I}) \quad (12.14)$$

$$p(\mathbf{x}_i | q_i = k, \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i | \mathbf{W}_k \mathbf{z}_i + \boldsymbol{\mu}_k, \mathbf{\Psi}) \quad (12.15)$$

²² To see this, note that there are $L-1$ free parameters in \mathbf{R} in the first column (since the column vector must be normalized to unit length), there are $L-2$ free parameters in the second column (which must be orthogonal to the first), and so on.

This is called a **mixture of factor analysers**(MFA) (Hinton et al. 1997).

Another way to think about this model is as a low-rank version of a mixture of Gaussians. In particular, this model needs $O(KLD)$ parameters instead of the $O(KD^2)$ parameters needed for a mixture of full covariance Gaussians. This can reduce overfitting. In fact, MFA is a good generic density model for high-dimensional real-valued data.

12.1.5 EM for factor analysis models

Below we state the results without proof. The derivation can be found in (Ghahramani and Hinton 1996a). To obtain the results for a single factor analyser, just set $r_{ic} = 1$ and $c = 1$ in the equations below. In Section 12.2.4 we will see a further simplification of these equations that arises when fitting a PPCA model, where the results will turn out to have a particularly simple and elegant interpretation.

In the E-step, we compute the posterior responsibility of cluster k for data point i using

$$r_{ik} \triangleq p(q_i = k | \mathbf{x}_i, \boldsymbol{\theta}) \propto \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{W}_k \mathbf{W}_k^T \boldsymbol{\Psi}_k) \quad (12.16)$$

The conditional posterior for \mathbf{z}_i is given by

$$p(\mathbf{z}_i | \mathbf{x}_i, q_i = k, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik}) \quad (12.17)$$

$$\boldsymbol{\Sigma}_{ik} \triangleq (\mathbf{I} + \mathbf{W}_k^T \boldsymbol{\Psi}_k^{-1} \mathbf{W}_k)_k^{-1} \quad (12.18)$$

$$\boldsymbol{\mu}_{ik} \triangleq \boldsymbol{\Sigma}_{ik} \mathbf{W}_k^T \boldsymbol{\Psi}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (12.19)$$

In the M step, it is easiest to estimate $\boldsymbol{\mu}_k$ and \mathbf{W}_k at the same time, by defining $\tilde{\mathbf{W}}_k = (\mathbf{W}_k, \boldsymbol{\mu}_k)$, $\tilde{\mathbf{z}} = (\mathbf{z}, 1)$, also, define

$$\tilde{\mathbf{W}}_k = (\mathbf{W}_k, \boldsymbol{\mu}_k) \quad (12.20)$$

$$\tilde{\mathbf{z}} = (\mathbf{z}, 1) \quad (12.21)$$

$$\mathbf{b}_{ik} \triangleq \mathbb{E}[\tilde{\mathbf{z}} | \mathbf{x}_i, q_i = k] = \mathbb{E}[(\boldsymbol{\mu}_{ik}; 1)] \quad (12.22)$$

$$\mathbf{C}_{ik} \triangleq \mathbb{E}[\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T | \mathbf{x}_i, q_i = k] \quad (12.23)$$

$$= \begin{pmatrix} \mathbb{E}[\mathbf{z} \mathbf{z}^T | \mathbf{x}_i, q_i = k] & \mathbb{E}[\mathbf{z} | \mathbf{x}_i, q_i = k] \\ \mathbb{E}[\mathbf{z} | \mathbf{x}_i, q_i = k]^T & 1 \end{pmatrix} \quad (12.24)$$

Then the M step is as follows:

$$\hat{\pi}_k = \frac{1}{N} \sum_{i=1}^N r_{ik} \quad (12.25)$$

$$\hat{\mathbf{W}}_k = \left(\sum_{i=1}^N r_{ik} \mathbf{x}_i \mathbf{b}_{ik}^T \right) \left(\sum_{i=1}^N r_{ik} \mathbf{x}_i \mathbf{C}_{ik}^T \right)^{-1} \quad (12.26)$$

$$\hat{\boldsymbol{\Psi}} = \frac{1}{N} \text{diag} \left[\sum_{i=1}^N r_{ik} (\mathbf{x}_i - \hat{\mathbf{W}}_{ik} \mathbf{b}_{ik}) \mathbf{x}_i^T \right] \quad (12.27)$$

Note that these updates are for vanilla EM. A much faster version of this algorithm, based on ECM, is described in (Zhao and Yu 2008).

12.1.6 Fitting FA models with missing data

In many applications, such as collaborative filtering, we have missing data. One virtue of the EM approach to fitting an FA/PPCA model is that it is easy to extend to this case. However, overfitting can be a problem if there is a lot of missing data. Consequently it is important to perform MAP estimation or to use Bayesian inference. See e.g., (Ilin and Raiko 2010) for details.

12.2 Principal components analysis (PCA)

Consider the FA model where we constrain $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$, and \mathbf{W} to be orthonormal. It can be shown (Tipping and Bishop 1999) that, as $\sigma^2 \rightarrow 0$, this model reduces to classical (nonprobabilistic) **principal components analysis**(PCA), also known as the Karhunen Loeve transform. The version where $\sigma^2 > 0$ is known as **probabilistic PCA**(PPCA) (Tipping and Bishop 1999), or sensible PCA(Roweis 1997).

12.2.1 Classical PCA

12.2.1.1 Statement of the theorem

The synthesis view of classical PCA is summarized in the following theorem.

Theorem 12.1. *Suppose we want to find an orthogonal set of L linear basis vectors $\mathbf{w}_j \in \mathbb{R}^D$, and the corresponding scores $\mathbf{z}_i \in \mathbb{R}^L$, such that we minimize the average **reconstruction error***

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 \quad (12.28)$$

where $\hat{\mathbf{x}}_i = \mathbf{W} \mathbf{z}_i$, subject to the constraint that \mathbf{W} is orthonormal. Equivalently, we can write this objective as follows

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{N} \|\mathbf{X} - \mathbf{W} \mathbf{Z}^T\|^2 \quad (12.29)$$

where \mathbf{Z} is an $N \times L$ matrix with the \mathbf{z}_i in its rows, and $\|\mathbf{A}\|_F$ is the **Frobenius norm** of matrix \mathbf{A} , defined by

$$\|A\|_F \triangleq \sqrt{\sum_{i=1}^M \sum_{j=1}^N a_{ij}^2} = \sqrt{\text{tr}(A^T A)} \quad (12.30)$$

The optimal solution is obtained by setting $\hat{W} = V_L$, where V_L contains the L eigenvectors with largest eigenvalues of the empirical covariance matrix, $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$. (We assume the x_i have zero mean, for notational simplicity.) Furthermore, the optimal low-dimensional encoding of the data is given by $\hat{z}_i = W^T x_i$, which is an orthogonal projection of the data onto the column space spanned by the eigenvectors.

An example of this is shown in Figure 12.1(a) for $D = 2$ and $L = 1$. The diagonal line is the vector w_1 ; this is called the first principal component or principal direction. The data points $x_i \in \mathbb{R}^2$ are orthogonally projected onto this line to get $z_i \in \mathbb{R}$. This is the best 1-dimensional approximation to the data. (We will discuss Figure 12.1(b) later.)

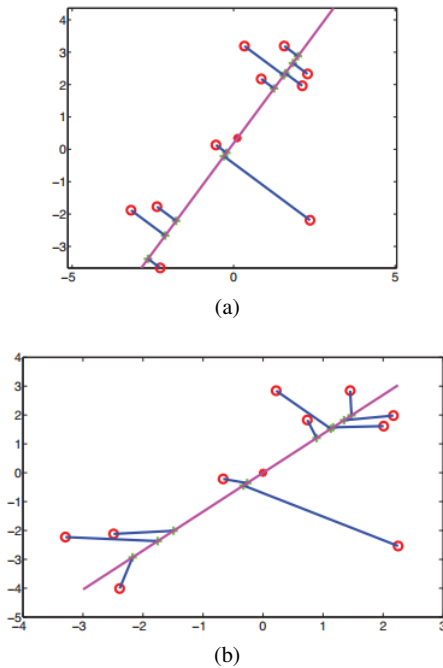


Fig. 12.1: An illustration of PCA and PPCA where $D = 2$ and $L = 1$. Circles are the original data points, crosses are the reconstructions. The red star is the data mean. (a) PCA. The points are orthogonally projected onto the line. (b) PPCA. The projection is no longer orthogonal: the reconstructions are shrunk towards the data mean (red star).

The principal directions are the ones along which the data shows maximal variance. This means that PCA can be misled by directions in which the variance is high merely because of the measurement scale. It is therefore

standard practice to standardize the data first, or equivalently, to work with correlation matrices instead of covariance matrices.

12.2.1.2 Proof *

See Section 12.2.2 of MLAPP.

12.2.2 Singular value decomposition (SVD)

We have defined the solution to PCA in terms of eigenvectors of the covariance matrix. However, there is another way to obtain the solution, based on the **singular value decomposition**, or **SVD**. This basically generalizes the notion of eigenvectors from square matrices to any kind of matrix.

Theorem 12.2. (SVD). Any matrix can be decomposed as follows

$$\underbrace{X}_{N \times D} = \underbrace{U}_{N \times N} \underbrace{\Sigma}_{N \times D} \underbrace{V^T}_{D \times D} \quad (12.31)$$

where U is an $N \times N$ matrix whose columns are orthonormal (so $U^T U = I$), V is $D \times D$ matrix whose rows and columns are orthonormal (so $V^T V = V V^T = I_D$), and Σ is a $N \times D$ matrix containing the $r = \min(N, D)$ singular values $\sigma_i \geq 0$ on the main diagonal, with 0s filling the rest of the matrix.

This shows how to decompose the matrix X into the product of three matrices: V describes an orthonormal basis in the domain, and U describes an orthonormal basis in the co-domain, and Σ describes how much the vectors in V are stretched to give the vectors in U .

Since there are at most D singular values (assuming $N > D$), the last ND columns of U are irrelevant, since they will be multiplied by 0. The **economy sized SVD**, or **thin SVD**, avoids computing these unnecessary elements. Let us denote this decomposition by $\hat{U} \hat{\Sigma} \hat{V}^T$. If $N > D$, we have

$$\underbrace{X}_{N \times D} = \underbrace{\hat{U}}_{N \times D} \underbrace{\hat{\Sigma}}_{D \times D} \underbrace{\hat{V}^T}_{D \times D} \quad (12.32)$$

as in Figure 12.2(a). If $N < D$, we have

$$\underbrace{X}_{N \times D} = \underbrace{\hat{U}}_{N \times N} \underbrace{\hat{\Sigma}}_{N \times N} \underbrace{\hat{V}^T}_{N \times D} \quad (12.33)$$

Computing the economy-sized SVD takes $O(ND \min(N, D))$ time (Golub and van Loan 1996, p254).

The connection between eigenvectors and singular vectors is the following:

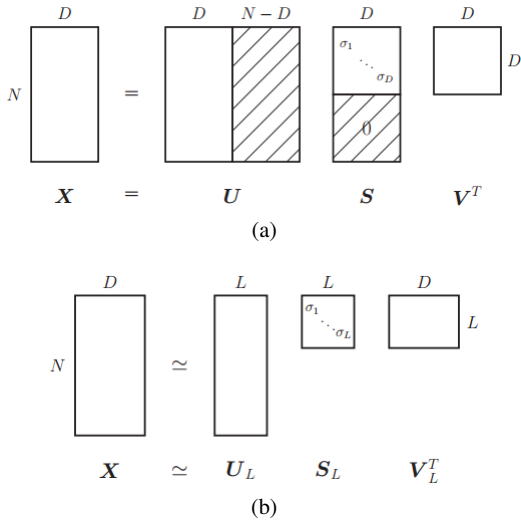


Fig. 12.2: (a) SVD decomposition of non-square matrices $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The shaded parts of $\mathbf{\Sigma}$, and all the off-diagonal terms, are zero. The shaded entries in \mathbf{U} and $\mathbf{\Sigma}$ are not computed in the economy-sized version, since they are not needed. (b) Truncated SVD approximation of rank L .

$$\mathbf{U} = \text{evec}(\mathbf{X}\mathbf{X}^T) \quad (12.34)$$

$$\mathbf{V} = \text{evec}(\mathbf{X}^T\mathbf{X}) \quad (12.35)$$

$$\mathbf{\Sigma}^2 = \text{eval}(\mathbf{X}\mathbf{X}^T) = \text{eval}(\mathbf{X}^T\mathbf{X}) \quad (12.36)$$

For the proof please read Section 12.2.3 of MLAPP.

Since the eigenvectors are unaffected by linear scaling of a matrix, we see that the right singular vectors of \mathbf{X} are equal to the eigenvectors of the empirical covariance $\hat{\mathbf{\Sigma}}$. Furthermore, the eigenvalues of $\hat{\mathbf{\Sigma}}$ are a scaled version of the squared singular values.

However, the connection between PCA and SVD goes deeper. From Equation 12.31, we can represent a rank r matrix as follows:

$$\mathbf{X} = \sigma_1 \begin{pmatrix} | \\ \mathbf{u}_1 \\ | \end{pmatrix} \begin{pmatrix} - & \mathbf{v}_1 & - \end{pmatrix} + \cdots + \sigma_r \begin{pmatrix} | \\ \mathbf{u}_r \\ | \end{pmatrix} \begin{pmatrix} - & \mathbf{v}_r^T & - \end{pmatrix}$$

If the singular values die off quickly, we can produce a rank L approximation to the matrix as follows:

$$\begin{aligned} \mathbf{X} &\approx \sigma_1 \begin{pmatrix} | \\ \mathbf{u}_1 \\ | \end{pmatrix} \begin{pmatrix} - & \mathbf{v}_1 & - \end{pmatrix} + \cdots + \sigma_L \begin{pmatrix} | \\ \mathbf{u}_L \\ | \end{pmatrix} \begin{pmatrix} - & \mathbf{v}_L^T & - \end{pmatrix} \\ &= \mathbf{U}_{:,1:L} \mathbf{\Sigma}_{1:L,1:L} \mathbf{V}_{:,1:L}^T \end{aligned} \quad (12.37)$$

This is called a **truncated SVD** (see Figure 12.2(b)).

One can show that the error in this approximation is given by

$$\|\mathbf{X} - \mathbf{X}_L\|_F \approx \sigma_L \quad (12.38)$$

Furthermore, one can show that the SVD offers the best rank L approximation to a matrix (best in the sense of minimizing the above Frobenius norm).

Let us connect this back to PCA. Let $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be a truncated SVD of \mathbf{X} . We know that $\hat{\mathbf{W}} = \mathbf{V}$, and that $\hat{\mathbf{Z}} = \mathbf{X}\hat{\mathbf{W}}$, so

$$\hat{\mathbf{Z}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V} = \mathbf{U}\mathbf{\Sigma} \quad (12.39)$$

Furthermore, the optimal reconstruction is given by $\hat{\mathbf{X}} = \mathbf{Z}\hat{\mathbf{W}}$, so we find

$$\hat{\mathbf{X}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (12.40)$$

This is precisely the same as a truncated SVD approximation! This is another illustration of the fact that PCA is the best low rank approximation to the data.

12.2.3 Probabilistic PCA

Theorem 12.3. ((*Tipping and Bishop 1999*)). Consider a factor analysis model in which $\mathbf{\Psi} = \sigma^2\mathbf{I}$ and \mathbf{W} is orthogonal. The observed data log likelihood is given by

$$\begin{aligned} \log p(\mathbf{X}|\mathbf{W}, \sigma^2\mathbf{I}) &= -\frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{i=1}^N \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i \\ &= -\frac{N}{2} \ln |\mathbf{C}| + \text{tr}(\mathbf{C}^{-1} \mathbf{\Sigma}) \end{aligned} \quad (12.41)$$

where $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ and $\mathbf{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{N} \mathbf{X}\mathbf{X}^T$. (We are assuming centred data, for notational simplicity.) The maxima of the log-likelihood are given by

$$\hat{\mathbf{W}} = \mathbf{V}(\mathbf{\Lambda} - \sigma^2\mathbf{I})^{\frac{1}{2}} \mathbf{R} \quad (12.42)$$

where \mathbf{R} is an arbitrary $L \times L$ orthogonal matrix, \mathbf{V} is the $D \times L$ matrix whose columns are the first L eigenvectors of $\mathbf{\Sigma}$, and $\mathbf{\Lambda}$ is the corresponding diagonal matrix of eigenvalues. Without loss of generality, we can set $\mathbf{R} = \mathbf{I}$. Furthermore, the MLE of the noise variance is given by

$$\hat{\sigma}^2 = \frac{1}{D-L} \sum_{j=L+1}^D \lambda_j \quad (12.43)$$

which is the average variance associated with the discarded dimensions.

Thus, as $\sigma^2 \rightarrow 0$, we have $\hat{\mathbf{W}} \rightarrow \mathbf{V}$, as in classical PCA. What about $\hat{\mathbf{Z}}$? It is easy to see that the posterior over the latent factors is given by

$$p(z_i|x_i, \hat{\theta}) = \mathcal{N}(z_i|\hat{\mathbf{F}}^{-1}\hat{\mathbf{W}}^T \mathbf{x}_i, \sigma^2 \hat{\mathbf{F}}^{-1}) \quad (12.44)$$

$$\hat{\mathbf{F}} \triangleq \hat{\mathbf{W}}^T \hat{\mathbf{W}} + \sigma^2 \mathbf{I} \quad (12.45)$$

(Do not confuse $\mathbf{F} = \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I}$ with $\mathbf{C} = \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}$.) Hence, as $\sigma^2 \rightarrow 0$, we find $\hat{\mathbf{W}} \rightarrow \mathbf{V}$, $\hat{\mathbf{F}} \rightarrow \mathbf{I}$ and $z_i \rightarrow \mathbf{V}^T \mathbf{x}_i$. Thus the posterior mean is obtained by an orthogonal projection of the data onto the column space of \mathbf{V} , as in classical PCA.

Note, however, that if $\sigma^2 \rightarrow 0$, the posterior mean is not an orthogonal projection, since it is shrunk somewhat towards the prior mean, as illustrated in Figure 12.1(b). This sounds like an undesirable property, but it means that the reconstructions will be closer to the overall data mean, $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$.

12.2.4 EM algorithm for PCA

Although the usual way to fit a PCA model uses eigenvector methods, or the SVD, we can also use EM, which will turn out to have some advantages that we discuss below. EM for PCA relies on the probabilistic formulation of PCA. However the algorithm continues to work in the zero noise limit, $\sigma^2 = 0$, as shown by (Roweis 1997).

Let $\tilde{\mathbf{Z}}$ be a $L \times N$ matrix storing the posterior means (low-dimensional representations) along its columns. Similarly, let $\tilde{\mathbf{X}} = \mathbf{X}^T$ store the original data along its columns. From Equation 12.44, when $\sigma^2 = 0$, we have

$$\tilde{\mathbf{Z}} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \tilde{\mathbf{X}} \quad (12.46)$$

This constitutes the E step. Notice that this is just an orthogonal projection of the data.

From Equation 12.26, the M step is given by

$$\hat{\mathbf{W}} = \left(\sum_{i=1}^N \mathbf{x}_i \mathbb{E}[z_i]^T \right) \left(\sum_{i=1}^N \mathbb{E}[z_i] \mathbb{E}[z_i]^T \right)^{-1} \quad (12.47)$$

where we exploited the fact that $\boldsymbol{\Sigma} = \text{cov}[z_i|x_i, \theta] = \mathbf{0}$ when $\sigma^2 = 0$.

(Tipping and Bishop 1999) showed that the only stable fixed point of the EM algorithm is the globally optimal solution. That is, the EM algorithm converges to a solution where \mathbf{W} spans the same linear subspace as that defined by the first L eigenvectors. However, if we want \mathbf{W} to be orthogonal, and to contain the eigenvectors in descending order of eigenvalue, we have to orthogonalize the resulting matrix (which can be done quite cheaply). Alternatively, we can modify EM to give the principal basis directly (Ahn and Oh 2003).

This algorithm has a simple physical analogy in the case $D = 2$ and $L = 1$ (Roweis 1997). Consider some

points in \mathbb{R}^2 attached by springs to a rigid rod, whose orientation is defined by a vector \mathbf{w} . Let z_i be the location where the i 'th spring attaches to the rod. See Figure 12.11 of MLAPP for an illustration.

Apart from this pleasing intuitive interpretation, EM for PCA has the following advantages over eigenvector methods:

- EM can be faster. In particular, assuming $N, D \gg L$, the dominant cost of EM is the projection operation in the E step, so the overall time is $O(TLND)$, where T is the number of iterations. This is much faster than the $O(\min(ND^2, DN^2))$ time required by straightforward eigenvector methods, although more sophisticated eigenvector methods, such as the Lanczos algorithm, have running times comparable to EM.
- EM can be implemented in an online fashion, i.e., we can update our estimate of \mathbf{W} as the data streams in.
- EM can handle missing data in a simple way (see Section 12.1.6).
- EM can be extended to handle mixtures of PPCA/FA models.
- EM can be modified to variational EM or to variational Bayes EM to fit more complex models.

12.3 Choosing the number of latent dimensions

In Section 11.5, we discussed how to choose the number of components K in a mixture model. In this section, we discuss how to choose the number of latent dimensions L in a FA/PCA model.

12.3.1 Model selection for FA/PPCA

TODO

12.3.2 Model selection for PCA

TODO

12.4 PCA for categorical data

In this section, we consider extending the factor analysis model to the case where the observed data is categorical rather than real-valued. That is, the data has the form

$y_{ij} \in \{1, \dots, C\}$, where $j = 1 : R$ is the number of observed response variables. We assume each y_{ij} is generated from a latent variable $z_i \in \mathbb{R}^L$, with a Gaussian prior, which is passed through the softmax function as follows:

$$p(z_i) = \mathcal{N}(z_i | \mathbf{0}, \mathbf{I}) \quad (12.48)$$

$$p(y_i | z_i, \theta) = \prod_{j=1}^R \text{Cat}(y_{ir} | \mathcal{S}(\mathbf{W}_r^T z_i + w_{0r})) \quad (12.49)$$

where $\mathbf{W}_r \in \mathbb{R}^L$ is the factor loading matrix for response j , and $\mathbf{W}_{0r} \in \mathbb{R}^M$ is the offset term for response r , and $\theta = (\mathbf{W}_r, \mathbf{W}_{0r})_{r=1}^R$. (We need an explicit offset term, since clamping one element of z_i to 1 can cause problems when computing the posterior covariance.) As in factor analysis, we have defined the prior mean to be $\mu_0 = \mathbf{0}$ and the prior covariance $\mathbf{V}_0 = \mathbf{I}$, since we can capture non-zero mean by changing w_{0j} and non-identity covariance by changing \mathbf{W}_r . We will call this categorical PCA. See Chapter 27 TODO for a discussion of related models.

In (Khan et al. 2010), we show that this model outperforms finite mixture models on the task of imputing missing entries in design matrices consisting of real and categorical data. This is useful for analysing social science survey data, which often has missing data and variables of mixed type.

12.5 PCA for paired and multi-view data

12.5.1 Supervised PCA (latent factor regression)

12.5.2 Discriminative supervised PCA

12.5.3 Canonical correlation analysis

12.6 Independent Component Analysis (ICA)

Let $\mathbf{x}_t \in \mathbb{R}^D$ be the observed signal at the sensors at time t , and $\mathbf{z}_t \in \mathbb{R}^L$ be the vector of source signals. We assume that

$$\mathbf{x}_t = \mathbf{W} \mathbf{z}_t + \epsilon_t \quad (12.50)$$

where \mathbf{W} is an $D \times L$ matrix, and $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \Psi)$. In this section, we treat each time point as an independent observation, i.e., we do not model temporal correlation (so we could replace the t index with i , but we stick with t to be consistent with much of the ICA literature). The goal is to infer the source signals, $p(\mathbf{z}_t | \mathbf{x}_t, \theta)$. In this context, \mathbf{W} is called the **mixing matrix**. If $L = D$ (number of sources =

number of sensors), it will be a square matrix. Often we will assume the noise level, $|\Psi|$, is zero, for simplicity.

So far, the model is identical to factor analysis. However, we will use a different prior for $p(\mathbf{z}_t)$. In PCA, we assume each source is independent, and has a Gaussian distribution. We will now relax this Gaussian assumption and let the source distributions be any *non-Gaussian* distribution

$$p(\mathbf{z}_t) = \prod_{j=1}^L p_j(z_{tj}) \quad (12.51)$$

Without loss of generality, we can constrain the variance of the source distributions to be $\mathbf{1}$, because any other variance can be modelled by scaling the rows of \mathbf{W} appropriately. The resulting model is known as **independent component analysis** or **ICA**.

The reason the Gaussian distribution is disallowed as a source prior in ICA is that it does not permit unique recovery of the sources. This is because the PCA likelihood is invariant to any orthogonal transformation of the sources \mathbf{z}_t and mixing matrix \mathbf{W} . PCA can recover the best linear subspace in which the signals lie, but cannot uniquely recover the signals themselves.

ICA requires that \mathbf{W} is square and hence invertible. In the non-square case (e.g., where we have more sources than sensors), we cannot uniquely recover the true signal, but we can compute the posterior $p(\mathbf{z}_t | \mathbf{x}_t, \hat{\mathbf{W}})$, which represents our beliefs about the source. In both cases, we need to estimate \mathbf{W} as well as the source distributions p_j . We discuss how to do this below.

12.6.1 Maximum likelihood estimation

In this section, we discuss ways to estimate square mixing matrices \mathbf{W} for the noise-free ICA model. As usual, we will assume that the observations have been centered; hence we can also assume \mathbf{z} is zero-mean. In addition, we assume the observations have been whitened, which can be done with PCA.

If the data is centered and whitened, we have $\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \mathbf{I}$. But in the noise free case, we also have

$$\text{cov}[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^T] = \mathbf{W} \mathbb{E}[\mathbf{z}\mathbf{z}^T] \mathbf{W}^T \quad (12.52)$$

Hence we see that \mathbf{W} must be orthogonal. This reduces the number of parameters we have to estimate from D^2 to $D(D-1)/2$. It will also simplify the math and the algorithms.

Let $\mathbf{V} = \mathbf{W}^{-1}$; these are often called the recognition weights, as opposed to \mathbf{W} , which are the generative weights.

Since $\mathbf{x} = \mathbf{W} \mathbf{z}$, we have, from Equation 2.46,

$$\begin{aligned} p_x(\mathbf{W}z_t) &= p_z(z_t)|\det(\mathbf{W}^{-1})| \\ &= p_z(\mathbf{V}\mathbf{x}_t)|\det(\mathbf{V})| \end{aligned} \quad (12.53)$$

Hence we can write the log-likelihood, assuming T iid samples, as follows:

$$\frac{1}{T} \log p(\mathcal{D}|\mathbf{V}) = \log |\det(\mathbf{V})| + \frac{1}{T} \sum_{j=1}^L \sum_{t=1}^T \log p_j(\mathbf{v}_j^T \mathbf{x}_t)$$

where \mathbf{v}_j is the j 'th row of \mathbf{V} . Since we are constraining \mathbf{V} to be orthogonal, the first term is a constant, so we can drop it. We can also replace the average over the data with an expectation operator to get the following objective

$$\text{NLL}(\mathbf{V}) = \sum_{j=1}^L \mathbb{E}[G_j(z_j)] \quad (12.54)$$

where $z_j = \mathbf{v}_j^T \mathbf{x}$ and $G_j(z) \triangleq -\log p_j(z)$. We want to minimize this subject to the constraint that the rows of \mathbf{V} are orthogonal. We also want them to be unit norm, since this ensures that the variance of the factors is unity (since, with whitened data, $\mathbb{E}[\mathbf{v}_j^T \mathbf{x}] = \|\mathbf{v}_j\|^2$, which is necessary to fix the scale of the weights. In other words, \mathbf{V} should be an orthonormal matrix.

It is straightforward to derive a gradient descent algorithm to fit this model; however, it is rather slow. One can also derive a faster algorithm that follows the natural gradient; see e.g., (MacKay 2003, ch 34) for details. A popular alternative is to use an approximate Newton method, which we discuss in Section 12.6.2. Another approach is to use EM, which we discuss in Section 12.6.3.

12.6.2 The FastICA algorithm

12.6.3 Using EM

12.6.4 Other estimation principles *