

Review of Independent Component Analysis

Element of Statistical Learning Book, Chapter 14

Hanchao Zhang

November 9, 2021

1 Latent Variables an Factor Analysis

The singular decomposition

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \quad (1)$$

We can write $\mathbf{S} = \sqrt{N}\mathbf{U}$ and $\mathbf{A}^T = \frac{\mathbf{D}\mathbf{V}^T}{\sqrt{N}}$
and we have

$$\mathbf{X} = \mathbf{S}\mathbf{A}^T = \sqrt{N}\mathbf{U}\frac{\mathbf{D}\mathbf{V}^T}{\sqrt{N}} = \mathbf{U}\mathbf{D}\mathbf{V}^T \quad (2)$$

Where \mathbf{S} and \mathbf{X} have mean 0, and \mathbf{U} is an orthogonal matrix. We can interpret the SVD or the corresponding principal component analysis as an estimate of a latent variable model.

$$\begin{aligned} X_1 &= a_{11}S_1 + a_{12}S_2 + \cdots + a_{1p}S_p \\ X_2 &= a_{21}S_1 + a_{22}S_2 + \cdots + a_{2p}S_p \\ &\vdots \\ X_p &= a_{p1}S_1 + a_{p2}S_2 + \cdots + a_{pp}S_p \end{aligned} \quad (3)$$

However, for any orthogonal matrix \mathbf{R} , we can write

$$\begin{aligned} \mathbf{X} &= \mathbf{A}\mathbf{S} \\ &= \mathbf{A}\mathbf{R}^T\mathbf{R}\mathbf{S} \\ &= \mathbf{A}^*\mathbf{S}^* \end{aligned} \quad (4)$$

Hence there are many such decompositions and it is therefore impossible to identify any particular latent variable as unique underlying sources. The classical factor analysis model has the form ($q < p$)

$$\begin{aligned}
X_1 &= a_{11}S_1 + a_{12}S_2 + \cdots + a_{1q}S_q + \varepsilon_1 \\
X_2 &= a_{21}S_1 + a_{22}S_2 + \cdots + a_{2q}S_q + \varepsilon_2 \\
&\vdots \qquad \qquad \qquad \vdots \\
X_p &= a_{p1}S_1 + a_{p2}S_2 + \cdots + a_{pq}S_q + \varepsilon_p
\end{aligned} \tag{5}$$

or

$$X = \mathbf{A}S + \varepsilon \tag{6}$$

Typically the S_j and ε_j are modeled as Gaussian random variables and the model is fit by maximum likelihood.

2 Independent Component Analysis

The ICA model has the form:

$$\begin{aligned}
X_1 &= a_{11}S_1 + a_{12}S_2 + \cdots + a_{1p}S_p \\
X_2 &= a_{21}S_1 + a_{22}S_2 + \cdots + a_{2p}S_p \\
&\vdots \qquad \qquad \qquad \vdots \\
X_p &= a_{p1}S_1 + a_{p2}S_2 + \cdots + a_{pp}S_p
\end{aligned} \tag{7}$$

or

$$X = \mathbf{A}S \tag{8}$$

where the S_i are assumed to be statistically independent rather than uncorrelated. Intuitively, uncorrelated determines the second cross moment $\text{Cov}(X)$, and statistically independent determines all cross moments.

We wish to recover the matrix \mathbf{A} in $X = \mathbf{A}S$. Without loss of generality, we can assume that X has already been whitened to have $\text{Cov}(X) = \mathbf{I}$, which implies that \mathbf{A} is orthogonal. Solving the ICA problem amounts to finding an orthogonal \mathbf{A} such that the components of the vector random variables $S = \mathbf{A}^T X$ are independent and non-Gaussian.

Many of the popular approaches to ICA are based on entropy. The differential entropy H of a random variable Y with density $g(y)$ is given by

$$H(Y) = E[-\log g(Y)] = - \int g(y) \log g(y) dy \tag{9}$$

The quantity $I(Y)$ is called the *Kullback-Leibler* distance or *mutual information*

$$I(Y) = \sum_{j=1}^p H(Y_j) - H(Y) \tag{10}$$

This is the measurement of *Kullback-Leibler* distance between the density $g(y)$ of Y and its independence version $\prod_{j=1}^p g_j(y_j)$, where $g_j(y_j)$ is the marginal density of Y_j . Now, if X has covariance \mathbf{I} , and $Y = \mathbf{A}^T X$ with \mathbf{A} orthogonal, then we will have

$$I(Y) = \sum_{j=1}^p H(Y_j) - H(Y) \quad (11)$$

$$= \sum_{j=1}^p H(Y_j) - H(X) - \log |\det \mathbf{A}| \quad (12)$$

$$= \sum_{j=1}^p H(Y_j) - H(X) \quad (13)$$

since we have $Y \sim f_Y(y)$, and $Y = \mathbf{A}^T X$, where we can get $X = \mathbf{A}Y$, and the *PDF* of X is $f_X(x) = f_Y(x) \|\mathbf{A}\|$

$$H(X) = E[-\log f_X(x)] = -E \left[\log (\|\mathbf{A}\| \cdot f_Y(x)) \right] \quad (14)$$

$$= -E \left[\log \|\mathbf{A}\| + \log f_Y(x) \right] \quad (15)$$

$$= -\log \|\mathbf{A}\| - E[-\log f_Y(x)] \quad (16)$$

$$= -\log \|\mathbf{A}\| + H[Y] \quad (17)$$

$$= H[Y] \quad (18)$$

For convenience, rather than using the entropy $H(Y_j)$, Hyvarinen and Oja (2000) uses the negentropy measure $J(Y_j)$ defined by

Note

a

$$J(Y_j) = H(Z_j) - H(Y_j) \quad (19)$$

where Z_j is a Gaussian random variable with the same variance as Y_j . They proposed simple approximations to negentropy which can be computed and optimized on the data.

$$J(Y_j) \approx \left(E[G(Y_j)] - E[G(Z_j)] \right)^2 \quad (20)$$

where $G(u) = \frac{1}{a} \log \cosh(au)$ for $1 < a < 2$.

With pre-whitened data, this amounts to looking for components that are as independent as possible.

3 Direct Approach of Independent Component Analysis by a Joint Product Density

Independent component have by definition a joint product density

$$f_S(s) = \prod_{j=1}^p f_j(s_j) \quad (21)$$

And in the spirit of representing departures from Gaussianity, we represent each f_j as

$$f_j(s_j) = \phi(s_j) \exp\{g_j(s_j)\} = \frac{1}{\sqrt{2\pi}} \exp\{\frac{1}{2}s_j^2\} \cdot \exp\{g_j(s_j)\} \quad (22)$$

the log-likelihood for the observed data $X = \mathbf{A}S$ is

$$\ell(\mathbf{A}\{g_j\}_{j=1}^p; \mathbf{X}) = \sum_{i=1}^N \sum_{j=1}^p [\log \phi_j(a_j^T x_i) + g_j(a_j^T x_i)] \quad (23)$$

which we want to maximize subjected to \mathbf{A} orthogonal and g_j result in density function. So, we instead maximize a regularized version

$$\sum_{j=1}^p \left[\sum_{i=1}^N [\log \phi_j(a_j^T x_i) + g_j(a_j^T x_i)] - \underbrace{\int \phi(t) e^{g_j(t)} dt}_{\text{density}} - \lambda_j \underbrace{\int \{g_j'''(t)\}^2(t) dt}_{\text{splines penalty}} \right] \quad (24)$$

The first integral that controls for density is problematic and requires an approximation. We construct a fine grid of L values s_ℓ^* in increments Δ covering the observed value s_i and count the number of s_i in the resulting bins:

$$y_\ell^* = \frac{\#s_i \in (s_i^* - \frac{\Delta}{2}, s_i^* + \frac{\Delta}{2})}{N} \quad (25)$$

Step 2 (a), we can then approximate the penalized likelihood by

$$\sum_{\ell=1}^L \left\{ y_\ell^* [\log(\phi(s_\ell^*)) + g(s_\ell^*)] - \Delta \phi(s_\ell^*) e^{g(s_\ell^*)} \right\} - \lambda \int g'''^2(s) ds \quad (26)$$

and in practice, we set all λ_j to the same.

Step 2(b), we optimize the \mathbf{A} with respect to the penalized likelihood function. Only the first terms in the sum involve \mathbf{A} , and since \mathbf{A} is orthogonal, ϕ do not depend on \mathbf{A} . Hence, we need to maximize

$$C(\mathbf{A}) = \frac{1}{N} \sum_{j=1}^p \sum_{i=1}^N \hat{g}_j(a_j^T x_i) = \sum_{j=1}^p C_j(a_j) \quad (27)$$

1. for each j update

$$a_j \leftarrow E \left\{ X \hat{g}_j'(a_j^T X) - E[g_j''(a_j^T X)] a_j \right\} \quad (28)$$

2. Orthogonalize \mathbf{A} using the symmetric square-root transformation $(\mathbf{A}\mathbf{A}^T)^{\frac{1}{2}}\mathbf{A}$. Let UDV^T be the singular decomposition of \mathbf{A} , we will have

$$(\mathbf{A}\mathbf{A}^T)^{\frac{1}{2}}\mathbf{A} = (UDV^TV D^T U^T)^{\frac{1}{2}}UDV^T \quad (29)$$

$$= (UD^2U^T)^{\frac{1}{2}}UDV^T \quad (30)$$

$$= D^{-1}UDV^T \quad (31)$$

$$\mathbf{A} \leftarrow UV^T \quad (32)$$