

In order to cast PCA in a probabilistic framework, we can make a slight modification to the factor analysis model, making the conditional variances σ_i^2 equal to each other. In that case the covariance of \mathbf{x} is just $\mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}$, where σ^2 is now a scalar. This yields the conditional distribution

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}) \quad (13.5)$$

or equivalently

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \sigma\mathbf{z} \quad (13.6)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ is Gaussian noise. [Tipping and Bishop \(1999\)](#) then show an iterative EM algorithm for estimating the parameters \mathbf{W} and σ^2 .

This **probabilistic PCA** model takes advantage of the observation that most variations in the data can be captured by the latent variables \mathbf{h} , up to some small residual **reconstruction error** σ^2 . As shown by [Tipping and Bishop \(1999\)](#), probabilistic PCA becomes PCA as $\sigma \rightarrow 0$. In that case, the conditional expected value of \mathbf{h} given \mathbf{x} becomes an orthogonal projection of $\mathbf{x} - \mathbf{b}$ onto the space spanned by the d columns of \mathbf{W} , like in PCA.

As $\sigma \rightarrow 0$, the density model defined by probabilistic PCA becomes very sharp around these d dimensions spanned by the columns of \mathbf{W} . This can make the model assign very low likelihood to the data if the data does not actually cluster near a hyperplane.

13.2 Independent Component Analysis (ICA)

Independent component analysis (ICA) is among the oldest representation learning algorithms ([Herault and Ans, 1984](#); [Jutten and Herault, 1991](#); [Comon, 1994](#); [Hyvärinen, 1999](#); [Hyvärinen *et al.*, 2001a](#); [Hinton *et al.*, 2001](#); [Teh *et al.*, 2003](#)). It is an approach to modeling linear factors that seeks to separate an observed signal into many underlying signals that are scaled and added together to form the observed data. These signals are intended to be fully independent, rather than merely decorrelated from each other.¹

Many different specific methodologies are referred to as ICA. The variant that is most similar to the other generative models we have described here is a variant ([Pham *et al.*, 1992](#)) that trains a fully parametric generative model. The prior distribution over the underlying factors, $p(\mathbf{h})$, must be fixed ahead of time by the user. The model then deterministically generates $\mathbf{x} = \mathbf{W}\mathbf{h}$. We can perform a

¹See section 3.8 for a discussion of the difference between uncorrelated variables and independent variables.