

Machine Learning Worksheet 10

Unsupervised Learning

Problem 1: Where would you use PCA and where ICA on real world data?

Problem 2: Consider the latent space distribution

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$

and a conditional distribution for the observed variable $\mathbf{x} \in \mathbb{R}^d$,

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Phi})$$

where $\boldsymbol{\Phi}$ is an arbitrary symmetric, positive-definite noise covariance variable. Furthermore, \mathbf{A} is a non-singular $d \times d$ matrix and $\mathbf{y} = \mathbf{A}\mathbf{x}$. Show that for the maximum likelihood solution for the parameters of the model for \mathbf{y} specific constraints on $\boldsymbol{\Phi}$ are preserved in the following two cases: (i) \mathbf{A} is a diagonal matrix and $\boldsymbol{\Phi}$ is a diagonal matrix (this corresponds to the case of Factor Analysis). (ii) \mathbf{A} is orthogonal and $\boldsymbol{\Phi} = \sigma^2 \mathbf{I}$ (this corresponds to pPCA).

The model for \mathbf{y} is a *noiseless* linear transformation. Given that the distribution of \mathbf{x} is known, we therefore know the distribution of \mathbf{y} . Because of the definitions for \mathbf{z} and $\mathbf{x}|\mathbf{z}$ we know that \mathbf{x} is a Gaussian with mean $\boldsymbol{\mu}$ and covariance $\mathbf{W}\mathbf{W}^T + \boldsymbol{\Phi}$. And thus, \mathbf{y} is also Gaussian with mean $\mathbf{A}\boldsymbol{\mu}$ and covariance $\mathbf{A}\mathbf{W}\mathbf{W}^T\mathbf{A}^T + \mathbf{A}\boldsymbol{\Phi}\mathbf{A}^T$. Now, assuming that the maximum likelihood solutions for the conditional model for \mathbf{x} are $\boldsymbol{\mu}_x$, \mathbf{W}_x and $\boldsymbol{\Phi}_x$, by simple *matching patterns* the MLE solutions for \mathbf{y} are $\mathbf{A}\boldsymbol{\mu}_x$, $\mathbf{A}\mathbf{W}_x$ and $\mathbf{A}\boldsymbol{\Phi}_x\mathbf{A}^T$. (i) If \mathbf{A} and $\boldsymbol{\Phi}$ are diagonal matrices (Factor Analysis model), the characteristics of \mathbf{x} are preserved for \mathbf{y} . Similarly (ii) if \mathbf{A} is orthogonal and $\boldsymbol{\Phi}$ a scaled identity matrix, the model characteristics are also preserved ($\mathbf{A}\boldsymbol{\Phi}_x\mathbf{A}^T = \sigma^2 \mathbf{I}$ in this case).

Problem 3: Show that in the limit $\sigma^2 \rightarrow 0$ the posterior mean for the probabilistic PCA model becomes an orthogonal projection onto the same principal subspace as in PCA.

You may use the solution for the posterior of the FA model:

$$\begin{aligned} p(\mathbf{z}_i|\mathbf{x}_i) &= \mathcal{N}(\mathbf{z}_i|\mathbf{m}_i, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma} &= (\mathbf{I} + \mathbf{W}^T\boldsymbol{\Psi}^{-1}\mathbf{W})^{-1} \\ \mathbf{m}_i &= \boldsymbol{\Sigma}(\mathbf{W}^T\boldsymbol{\Psi}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})) \end{aligned}$$

Remember, pPCA is a Factor Analysis model with $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ and \mathbf{W} orthonormal. First, we plug the special form of $\boldsymbol{\Psi}$ into the general result for the posterior mean of the latent variable \mathbf{z} :

$$\mathbf{m}_i = \boldsymbol{\Sigma}(\mathbf{W}^T\sigma^{-2}\mathbf{I}(\mathbf{x}_i - \boldsymbol{\mu}))$$

with

$$\Sigma = (\mathbf{I} + \mathbf{W}^T \sigma^{-2} \mathbf{I} \mathbf{W})^{-1} = \sigma^2 (\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W})^{-1}$$

which gives

$$\mathbf{m}_i = (\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W})^{-1} (\mathbf{W}^T (\mathbf{x}_i - \boldsymbol{\mu}))$$

With $\sigma^2 \rightarrow 0$ the maximum likelihood solution for \mathbf{W} (given in slides) converges to $\mathbf{V}_l \boldsymbol{\Lambda}_l^{1/2}$. So $(\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W})^{-1} \rightarrow \boldsymbol{\Lambda}_l^{-1}$, and thus

$$\mathbf{m}_i = \boldsymbol{\Lambda}_l^{-1/2} \mathbf{V}_l^T (\mathbf{x}_i - \boldsymbol{\mu})$$

which is a projection on the same subspace as PCA does.