# 1 Partial least squares model

The partial least squares model is presented as a latent variable model in Kevin Murphy's 2012 book Machine Learning: A Probabilistic Perspective.

$$p(z) = \mathcal{N}(z|0, I) \tag{1}$$
$$p(v|z, W, \mu, \sigma) = \mathcal{N}(v|Wz + \mu, \sigma^2 I) \tag{2}$$

where

$$W = \begin{pmatrix} W_y & 0 \\ W_x & B_x \end{pmatrix} \tag{3}$$
$$z = (z^s; z^x) \tag{4}$$
$$v = (y; x) \tag{5}$$
$$\mu = (\mu_y; \mu_x). \tag{6}$$

Marginalising $z$ gives

$$p(v|W, \mu, \sigma) = \int \mathcal{N}(v|Wz + \mu, \sigma^2 I) \mathcal{N}(z|0, I) \, dz \tag{7}$$
$$= \mathcal{N}(v|\mu, WW^T + \sigma^2 I) \tag{8}$$

Conditioning on $x$ gives

$$p(y|x) = \mathcal{N}(y|m_{y|x}, S_{y|x}) \tag{9}$$

where

$$C = (B_x B_x^T + W_x W_x^T + \sigma^2 I)^{-1} \tag{10}$$
$$m_{y|x} = \mu_y + W_y W_x^T C(x - \mu_x) \tag{11}$$
$$S_{y|x} = \sigma^2 I + W_y W_y^T - W_y W_x^T C W_x W_y^T \tag{12}$$

Suppose we now obtain $N$ independent observations from the model

$$v_n = (y_n; x_n), \quad 1 \le n \le N. \tag{13}$$

We wish to estimate $W, \mu$ and $\sigma$. We can do this by maximising the likelihood of the data $v = x, y$ (8) using stochastic gradient descent.

One way to validate that our implementation works is to sample data from the model and compare the estimated parameters to those that we used to sample with. Suppose we sample $v_n$ from our model parameterised by $\sigma^*, \mu_x^*, \mu_y^*, W_y^*, W_x^*, B_x^*$. We can estimate $\hat{\sigma}, \hat{\mu}_x, \hat{\mu}_y, \hat{W}_y, \hat{W}_x, \hat{B}_x$ that maximise

$$\sum_n \log p(v|\sigma, \mu, W_y, W_x, B_x) \tag{14}$$

We would love to compare the estimated parameters to the underlying parameters but unfortunately (14) is invariant to orthonormal transformations of the $W_y, W_x, B_x$. That is

$$\sum_n \log p(v_n|\sigma, \mu, W_y U_s, W_x U_s, B_x U_x) \quad = \quad \sum_n \log p(v_n|\sigma, \mu, W_y, W_x, B_x) \quad (15)$$

for any orthonormal square matrices (of suitable dimension) $U_s, U_x$. We can see this as $W$ only appears as $WW^T$ in the likelihood and $WW^T = WUU^TW^T$ for any orthonormal $U$.