

# Dimensionality Reduction

(Part 2)

18-698 / 42-632

Neural Signal Processing  
Prof. Byron Yu

## C) Probabilistic PCA (PPCA)

- So far, we have formulated PCA as a linear projection of data into lower dimensional space.
- Here, we will show that PCA can be expressed as the maximum likelihood solution of a probabilistic latent variable model.

In other words, we will construct a latent variable model for which, when we maximize  $p(X|\theta)$  with respect to  $\theta$ , the resulting  $\theta$  will be the PC directions.

### C.1) Advantages of PPCA over conventional PCA

- PPCA assigns probabilities to data, so we can select the dimensionality  $M$  of low-d space and compare to other models using cross-validated likelihoods.
- PPCA has an explicit noise model, so it is able to more effectively denoise data than PCA.

- If data dimensionality  $D$  is large, diagonalization in conventional PCA is costly  $O(D^3)$ . If we only need top eigenvectors, we can compute them more efficiently using PPCA.

- Because PPCA is a probabilistic model, it can deal with missing data, PCA cannot.

$$x_1 = \begin{bmatrix} \otimes \\ \otimes \end{bmatrix} \quad x_2 = \begin{bmatrix} \otimes \\ \otimes \end{bmatrix} \quad x_3 = \begin{bmatrix} \otimes \\ \otimes \end{bmatrix} \dots$$

$\otimes$  denotes missing data

- PPCA represents a constrained form of Gaussian distribution

$$\begin{bmatrix} & 0 \\ & & \\ 0 & & \end{bmatrix}$$

Diagonal covariance  
is often too constrained

$$\begin{bmatrix} \Sigma \\ \end{bmatrix}$$

Full covariance  
is often too flexible  
(easy to overfit)

PPCA provides nice compromise between diagonal and full covariance.

- Because PPCA is a probabilistic model, we can easily propose extensions, such as mixtures of PPCA models (analogous to mixtures of Gaussians).
- PPCA is a generative model, so we can generate samples from the distribution.

Note: Many of these advantages of PPCA over conventional PCA are typical advantages of probabilistic over non-probabilistic models.

## C.2) Generative model for PPCA

$\underline{x} \in \mathbb{R}^D$  is high-dimensional observed data

$\underline{z} \in \mathbb{R}^M$  is low-dimensional latent variable

$$\begin{cases} P(\underline{z}) = N(\underline{z} \mid \underline{0}, \mathbf{I}) & \text{"state model"} \\ P(\underline{x} \mid \underline{z}) = N(\underline{x} \mid \underbrace{\mathbf{W}\underline{z} + \underline{\mu}}_{\text{"observation noise"}}, \underbrace{\sigma^2 \mathbf{I}}_{\text{"observation noise"}} \end{cases} \quad (7)$$

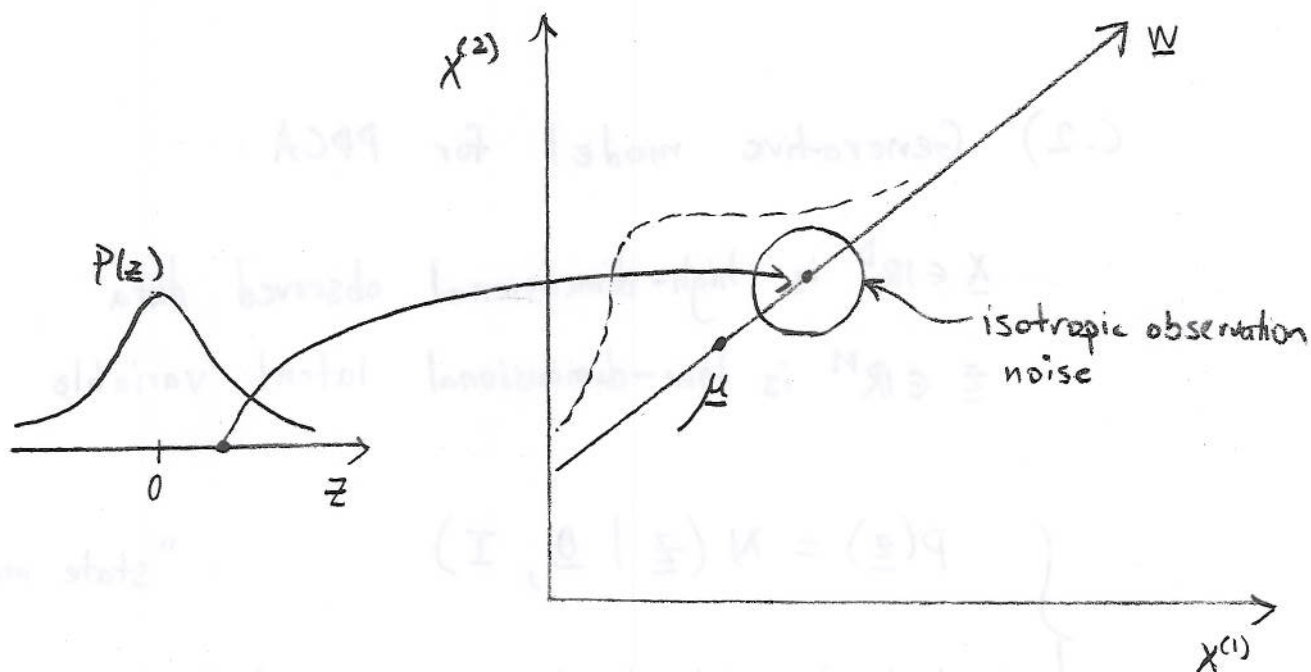
$\nwarrow$   
DxM matrix

### Relationship to PCA:

- If we fit this model to data  $x_1, \dots, x_N$  (i.e., we fit the model parameters  $\theta = \{W, \mu, \sigma^2\}$ ), the columns of  $W$  will span the principal component space (i.e., the space spanned by the columns of  $U_M$  in PCA)
- In the limit  $\sigma^2 \rightarrow 0$ , PPCA low-d projections approach PCA low-d projections.

### Illustration of PPCA generative model:

Let  $D=2$  and  $M=1$ .



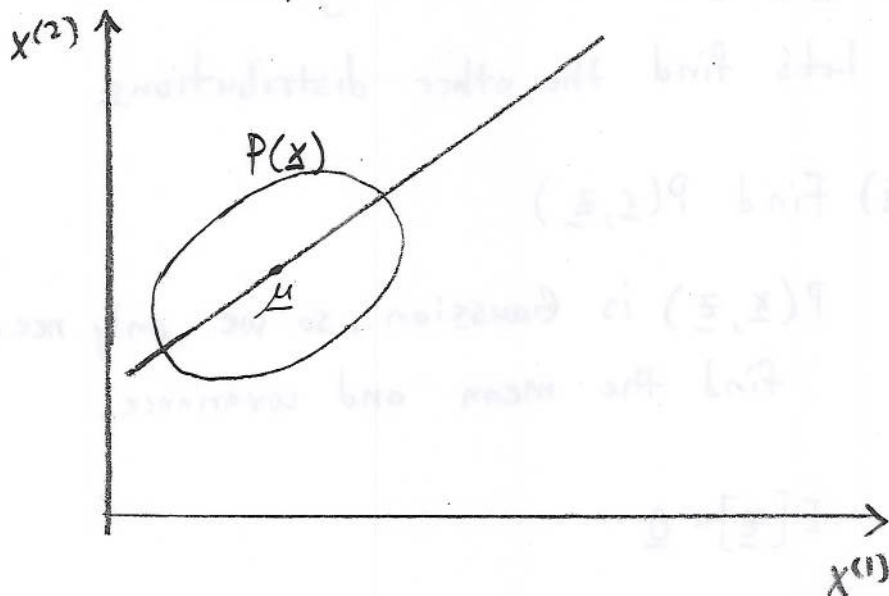


To generate from PPCA model:

- 1) Draw latent variable  $\underline{z} \sim N(\underline{0}, \underline{I})$ .
- 2) Project  $\underline{z}$  into high-dimensional data space
- 3) Add observation noise  $\underline{\epsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$ ,  
where  $\underline{\epsilon} \in \mathbb{R}^D$ .

$$\underline{x} = \underline{W}\underline{z} + \underline{\mu} + \underline{\epsilon}. \quad (8)$$

If we generate many simulated data points  $\underline{x}$ ,  
the data points will have a distribution  $P(\underline{x})$



Notes:

More  
details  
later

- $\underline{W}$  for PPCA is closely related to  $\underline{U}_M$  for PCA.
- $\underline{z}$  for PPCA is closely related to  $\underline{z}$  for PCA.
- $\underline{\mu}$  for PPCA is identical to  $\underline{\mu}$  for PCA.

C.3) PPCA is a linear-Gaussian model

- The variables  $\underline{x}$  and  $\underline{z}$  are linearly related (see (8))
- All marginal, conditional, and joint distributions are Gaussian.

$P(\underline{z}), P(\underline{x})$ : Marginal distributions

$P(\underline{x}|\underline{z}), P(\underline{z}|\underline{x})$ : Conditional distributions

$P(\underline{x}, \underline{z})$ : Joint distribution

$P(\underline{z})$  and  $P(\underline{x}|\underline{z})$  are given (see (7))

Let's find the other distributions.

i) Find  $P(\underline{x}, \underline{z})$

$P(\underline{x}, \underline{z})$  is Gaussian, so we only need to find the mean and covariance.

$$E[\underline{z}] = \underline{0}$$

From (8),

$$\begin{aligned} E[\underline{x}] &= E[W\underline{z} + \underline{\mu} + \underline{\epsilon}] \\ &= W E[\underline{z}] + \underline{\mu} + E[\underline{\epsilon}] \\ &= \underline{\mu} \end{aligned}$$

$$\text{Cov}(\underline{z}) = \mathbf{I}$$

$$\begin{aligned} \text{Cov}(\underline{x}) &= E[\underline{x}\underline{x}^T] - E[\underline{x}]E[\underline{x}]^T \\ &= E[(\underline{W}\underline{z} + \underline{\mu} + \underline{\epsilon})(\underline{W}\underline{z} + \underline{\mu} + \underline{\epsilon})^T] - \underline{\mu}\underline{\mu}^T \quad \text{from (8)} \\ &= E\left[\cancel{\underline{W}\underline{z}\underline{z}^T\cancel{W^T}} + \cancel{\underline{\mu}\underline{z}^T\cancel{W^T}} + \cancel{\underline{\epsilon}\underline{z}^T\cancel{W^T}} + \cancel{W\cancel{\underline{z}}\underline{\mu}^T} + \cancel{\underline{\mu}\underline{\mu}^T} + \cancel{\underline{\epsilon}\underline{\mu}^T} + \cancel{W\cancel{\underline{z}}\underline{\epsilon}^T} + \cancel{\underline{\mu}\cancel{\underline{\epsilon}}^T} + \underline{\epsilon}\underline{\epsilon}^T\right] - \underline{\mu}\underline{\mu}^T \\ &= \underline{W}E[\underline{z}\underline{z}^T]\underline{W}^T + \cancel{\underline{\mu}\underline{\mu}^T} + E[\underline{\epsilon}\underline{\epsilon}^T] - \cancel{\underline{\mu}\underline{\mu}^T} \\ &= \underline{W}\underline{W}^T + \sigma^2\mathbf{I} \end{aligned}$$

$$\begin{aligned} \text{Cov}(\underline{x}, \underline{z}) &= E[\underline{x}\underline{z}^T] - E[\underline{x}]E[\underline{z}]^T \\ &= E[(\underline{W}\underline{z} + \underline{\mu} + \underline{\epsilon})\underline{z}^T] \\ &= \underline{W}E[\underline{z}\underline{z}^T] + \underline{\mu}E[\underline{z}^T] + E[\underline{\epsilon}\underline{z}^T] \\ &= \underline{W} \end{aligned}$$

$$\begin{bmatrix} \underline{z} \\ \underline{x} \end{bmatrix} \sim N\left(\begin{bmatrix} \underline{0} \\ \underline{\mu} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & \underline{W}^T \\ \underline{W} & \underline{W}\underline{W}^T + \sigma^2\mathbf{I} \end{bmatrix}\right) \quad (9)$$

ii) Find  $P(\underline{x})$

From (9),

$$\underline{x} \sim N(\underline{\mu}, WW^T + \sigma^2 I) \quad (10)$$

iii) Find  $P(\underline{z} | \underline{x})$

In general, we would use Bayes rule.

An easier way is to apply the results of conditioning for Gaussian random variables (see PRML Section 2.3.1)

$$\left[ \begin{array}{l} \text{If } \underline{x} = \begin{bmatrix} x_a \\ x_b \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \right), \\ \\ P(x_a | x_b) \text{ is Gaussian with the following} \\ \text{mean and covariance.} \\ E[x_a | x_b] = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b) \\ \text{cov}(x_a | x_b) = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \end{array} \right]$$

$$E[\underline{z} | \underline{x}] = \underline{0} + W^T C^{-1} (\underline{x} - \underline{\mu})$$

$$\text{cov}(\underline{z} | \underline{x}) = I - W^T C^{-1} W,$$

$$\text{where } C = WW^T + \sigma^2 I$$



Thus,

$$\underline{z} | \underline{x} \sim N(W^T C^{-1}(\underline{x} - \underline{\mu}), \underbrace{I - W^T C^{-1} W}_{\text{note that covariance does not depend on } \underline{x}}) \quad (11)$$

note that covariance  
does not depend on  $\underline{x}$

#### 6.4) EM algorithm for PPCA

Goal: Maximize  $\log p(\{\underline{x}\} | \theta)$  w.r.t.  $\theta$ ,

where  $\theta = \{W, \underline{\mu}, \sigma^2\}$ .

EM will find the sample mean exactly for  $\underline{\mu}$  and attempt to find  $W$  and  $\sigma^2$  such that

$$\text{Sample Covariance} \rightarrow S \approx WW^T + \sigma^2 I.$$

This should make sense intuitively from (10).

EM is trying to match the mean and covariance of the data.

For simplicity here, we will fix  $\underline{\mu}$  to be the sample mean, which is the maximum likelihood solution.

E-step:

$P(\underline{z}_n | \underline{x}_n)$  is shown in (11).

M-step:

$$\begin{aligned}\log P(X, Z) &= \sum_{n=1}^N \log P(x_n, z_n) \\&= \sum_{n=1}^N \left( \log P(x_n | z_n) + \log P(z_n) \right) \\&= \sum_{n=1}^N \left( -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log(\sigma^2)^D - \frac{1}{2} (x_n - Wz_n - \mu)^T (\sigma^2 I)^{-1} (x_n - Wz_n - \mu) \right. \\&\quad \left. - \frac{M}{2} \log(2\pi) - \frac{1}{2} \log 1 - \frac{1}{2} z_n^T z_n \right)\end{aligned}$$

$$\begin{aligned}Q &= E_z [\log P(X, Z)] \\&= \sum_{n=1}^N \left\{ -\frac{D}{2} \log(2\pi \sigma^2) - \frac{1}{2\sigma^2} E[(x_n - \mu) - Wz_n]^T [(x_n - \mu) - Wz_n] \right. \\&\quad \left. - \frac{M}{2} \log(2\pi) - \frac{1}{2} E[z_n^T z_n] \right\} \\&= \sum_{n=1}^N \left\{ -\frac{D}{2} \log(2\pi \sigma^2) - \frac{1}{2\sigma^2} \left( (x_n - \mu)^T (x_n - \mu) - E[z_n^T W^T (x_n - \mu)] \right. \right. \\&\quad \left. \left. - E[(x_n - \mu)^T W z_n] + E[z_n^T W^T W z_n] \right) \right. \\&\quad \left. - \frac{M}{2} \log(2\pi) - \frac{1}{2} E[z_n^T z_n] \right\}\end{aligned}$$

$$= \sum_{n=1}^N \left\{ -\frac{D}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \left( (x_n - \mu)^T (x_n - \mu) - E[\underline{z}_n]^T W^T (x_n - \mu) \right. \right. \\ \left. \left. - (x_n - \mu)^T W E[\underline{z}_n] + \text{Tr}(W^T W E[\underline{z}_n \underline{z}_n^T]) \right) \right. \\ \left. - \frac{M}{2} \log(2\pi) - \frac{1}{2} E[\underline{z}_n^T \underline{z}_n] \right\}$$

$$\frac{\partial Q}{\partial W} = \sum_{n=1}^N \left\{ -\frac{1}{2\sigma^2} \left( -(x_n - \mu) E[\underline{z}_n]^T - (x_n - \mu) E[\underline{z}_n]^T \right. \right. \\ \left. \left. + 2W (E[\underline{z}_n \underline{z}_n^T]) \right) \right\} = 0$$

$$W \left( \sum_{n=1}^N E[\underline{z}_n \underline{z}_n^T] \right) = \sum_{n=1}^N (x_n - \mu) E[\underline{z}_n]^T$$

$$W_{\text{new}} = \left( \sum_{n=1}^N (x_n - \mu) E[\underline{z}_n]^T \right) \left( \sum_{n=1}^N E[\underline{z}_n \underline{z}_n^T] \right)^{-1} \quad (12)$$

$$\frac{\partial Q}{\partial \sigma^2} = \sum_{n=1}^N \left\{ -\frac{D}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{2} \cdot \frac{1}{(\sigma^2)^2} ( \cdot ) \right\} = 0$$

$$-ND + \frac{1}{\sigma^2} \sum_{n=1}^N ( \cdot ) = 0$$

$$\sigma^2 = \frac{1}{ND} \sum_{n=1}^N ( \cdot )$$

We can simplify the update for  $\sigma^2$  by plugging in  $W_{\text{new}}$ .

$$\sigma^2 = \frac{1}{ND} \left\{ \text{Tr} \left( \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T \right) - \text{Tr} \left( \sum_{n=1}^N (x_n - \mu) E[z_n]^T \cdot W^T \right) \right. \\ \left. - \text{Tr} \left( W \sum_{n=1}^N E[z_n] (x_n - \mu)^T \right) + \text{Tr} \left( W \sum_{n=1}^N E[x_n x_n^T] W^T \right) \right\}$$

$$(13) \quad \sigma_{\text{new}}^2 = \frac{1}{ND} \text{Tr} \left( \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T - W_{\text{new}} \sum_{n=1}^N E[z_n] (x_n - \mu)^T \right)$$

[ See physical analogy of EM for PPCA, PRML Figure 12.12. ]

The expressions for  $W_{\text{new}}$  and  $\sigma_{\text{new}}^2$  should make sense intuitively.

Consider the linear regression problem:

$$x_n = w z_n + \mu$$

Minimizing mean squared error with respect to  $w$ ,

$$w^* = \frac{\sum_{n=1}^N (x_n - \mu) z_n}{\sum_{n=1}^N z_n^2} \quad (\text{compare to (12)})$$

Using this  $w^*$ , minimum mean squared error is

$$\frac{1}{N} \left( \sum_{n=1}^N (x_n - \mu)^2 - w^* \sum_{n=1}^N z_n (x_n - \mu) \right) \quad (\text{compare to (13)})$$



## C.5) Relating PPCA to PCA

- PC directions

The columns of  $W$  for PPCA span the same space as that spanned by the columns of  $U_M$ .

The difference between  $W$  and  $U_M$  is that the columns of  $U_M$  are orthonormal and ordered based on amount of variance explained. In general, the columns of  $W$  are neither orthonormal nor ordered.

To obtain  $U_M$  from  $W$ , apply the singular value decomposition (SVD) to  $W$ . The SVD is a generalization of diagonalization to non-square matrices.

$$W = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_M \\ | & | & & | \end{bmatrix} \begin{bmatrix} d_1 & & & 0 \\ & d_2 & & \\ & & \ddots & \\ 0 & & & d_M \end{bmatrix} \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \dots & v_M \\ | & | & & | \end{bmatrix}^T$$

$\tilde{U} \qquad \qquad \tilde{D} \qquad \qquad \tilde{V}^T$

$(D \times M)$

$(M \times M)$

$(M \times M)$

columns of  $\tilde{U}$   
orthonormal

diagonal  
matrix

columns of  $\tilde{V}$   
orthonormal

$$d_1 \geq d_2 \geq \dots \geq d_M \geq 0$$

"singular values"

Now,  $\hat{U}$  for PPCA is identical to  $U_M$  for PCA.

- Low-dimensional projections

In PPCA, the low-d. projection corresponding to  $W$  is  $E[\underline{z}_n | \underline{x}_n] = W^T C^{-1} (\underline{x}_n - \underline{\mu})$  from (11).

The same point has high-d coordinates:

$$\begin{aligned} & W E[\underline{z}_n | \underline{x}_n] + \underline{\mu} \quad (\text{from (8)}) \\ &= \underbrace{\hat{U} \hat{D} \hat{V}^T E[\underline{z}_n | \underline{x}_n]}_{\text{call this } \hat{\underline{z}}_n} + \underline{\mu} \end{aligned}$$

There are several important reasons why  $\hat{\underline{z}}_n$  is easier to interpret than  $E[\underline{z}_n | \underline{x}_n]$ . All of the reasons stem from the fact that the columns of  $\hat{U}$  are orthonormal and ordered, while those of  $W$  are not.

i)  $\hat{\underline{z}}_n$  has the same units as  $\underline{x}_n$  (as in PCA)

ii) The dimensions of  $\hat{\underline{z}}_n$  are ordered (as in PCA), whereas  $E[\underline{z}_n | \underline{x}_n]$  is subject to arbitrary rotations and exchange-of-dimensions in latent space.

iii)  $\hat{\underline{z}}_n$  can be easily compared to PCA low-d projection.

- How does the low-dimensional projection for PPCA ( $\tilde{\mathbf{z}}_n$ ) relate to that for PCA ( $\mathbf{U}_M^T (\mathbf{x}_n - \boldsymbol{\mu})$ )?

$$\tilde{\mathbf{z}}_n = \begin{bmatrix} \frac{\lambda_1 - \sigma^2}{\lambda_1} & & 0 \\ & \ddots & \\ 0 & & \frac{\lambda_M - \sigma^2}{\lambda_M} \end{bmatrix} \mathbf{U}_M^T (\mathbf{x}_n - \boldsymbol{\mu}) \quad (14)$$

(We won't show this here)

In other words, PPCA projection is the PCA projection shrunk towards the origin, since

$$0 \leq \frac{\lambda_i - \sigma^2}{\lambda_i} \leq 1 \quad i=1, \dots, M.$$

As  $\sigma^2 \rightarrow 0$ , PPCA becomes identical to PCA.

- Intuition for PPCA vs. PCA

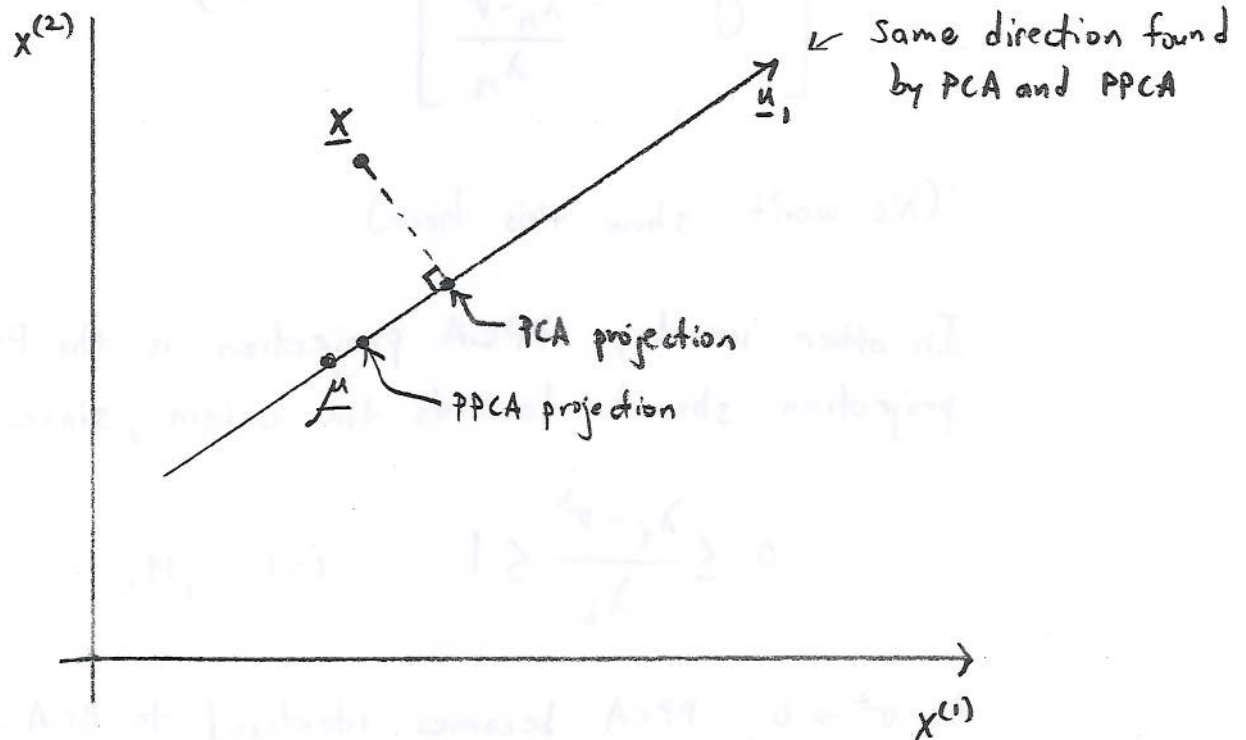
Each model tries to explain deviations of  $\mathbf{x}$  from  $\boldsymbol{\mu}$ .

PPCA can explain this using a combination of the latent variable  $\mathbf{z}$  and the observation noise  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$  (see (8)). But how much of each?

As  $\sigma^2$  increases (more observation noise), the proportion of the deviation attributed to observation noise increases and  $\mathbf{z}$  gets shrunk towards the mean.

PCA is the limit  $\sigma^2 \rightarrow 0$ , where there is no observation noise and all deviations of  $\underline{x}$  from  $\underline{\mu}$  must be explained using the latent variable  $\underline{z}$ .

In other words, PPCA is more effective at denoising the data than PCA.



### C.6) Returning to advantages of PPCA over PCA

- Dimensionality  $M$  of latent space for PPCA can be selected using cross-validated likelihoods, where  $P(\underline{x})$  given in (10).
- PPCA defines a constrained Gaussian in (10), that's a useful compromise between a Gaussian



with diagonal covariance (overconstrained in some settings) and a Gaussian with full covariance (underconstrained in some settings):

$$\text{cov}(\underline{x}) = WW^T + \sigma^2 I$$

## D) Factor Analysis (FA)

D.1) Motivation: PPCA assumes isotropic observation noise. In some settings, we would like each dimension of  $\underline{x}$  to have a different level of observation noise.

The only difference between FA and PPCA is that instead of  $\sigma^2 I$  in (7), the observation noise covariance is a diagonal matrix  $\Psi$ .

For FA,

$$\underline{x} \sim N(\underline{\mu}, WW^T + \Psi),$$

which is similar to (10).

We can define an EM algorithm that is nearly identical to that for PPCA. After replacing all

instances of  $\sigma^2 \mathbf{I}$  with  $\Psi$ , (13) becomes

$$\Psi_{\text{new}} = \frac{1}{N} \text{diag} \left\{ \sum_{n=1}^N (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T - W_{\text{new}} \sum_{n=1}^N E[\mathbf{z}_n](\mathbf{x}_n - \mu)^T \right\}$$

where  $\text{diag}\{\cdot\}$  zeroes all off-diagonal elements.

## D.2) Comparing FA and PPCA

Because of the non-isotropic observation noise, FA will identify different directions in the data space compared to PCA/PPCA.

As with PPCA, we will want to orthonormalize the columns of  $W$  for interpretability.

- PCA/PPCA is invariant to rotations in the data space, whereas FA is not.

The reason is the FA observation noise must be axis-aligned.

- FA is invariant to a component-wise rescaling of the data, whereas PCA/PPCA is not.

The reason is that PCA/PPCA uses the same  $\sigma^2$  for each component of  $\mathbf{x}$ .

Example of component-wise rescaling:

$$\underline{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} 3x^{(1)} \\ \frac{1}{2}x^{(2)} \end{bmatrix}$$

## Appendix

Matrix derivatives:

$$\frac{d}{dX} (\underline{a}^T X \underline{b}) = \underline{a} \underline{b}^T$$

$$\frac{d}{dX} (\underline{a}^T X^T \underline{b}) = \underline{b} \underline{a}^T$$

$$\frac{d}{dX} \text{Tr}(X^T X A) = X(A + A^T)$$

Matrix inversion lemma:

Inverting  $C = WW^T + \sigma^2 I$  directly in (11) is a costly  $O(D^3)$  operation. Instead, can apply the matrix inversion lemma:

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

now, matrix to be inverted is  $M \times M$ ,  
which is  $O(M^3)$  operation

Can use same trick for FA.