Dimensionality Reduction (Contd.)

Piyush Rai

Introduction to Machine Learning (CS771A)

October 9, 2018

Intro to Machine Learning (CS771A)

< ロ > < 回 > < 回 > < 回 > < 回 >

• Quiz graded and scores sent



- Quiz graded and scores sent
- Homework 3 out. Due on Oct 31, 11:59pm. Please start early.



イロト イロト イモト イモト

- Quiz graded and scores sent
- Homework 3 out. Due on Oct 31, 11:59pm. Please start early.
- We will finish homework 1 and 2 grading soon

・ロト ・回ト ・モト ・モト

- Quiz graded and scores sent
- Homework 3 out. Due on Oct 31, 11:59pm. Please start early.
- \bullet We will finish homework 1 and 2 grading soon
- Start thinking about your course project (if not working on it already)

Recap: Dimensionality Reduction - The Compression View



Intro to Machine Learning (CS771A)

DxK

Dx1

3

メロト メぼト メヨト

Recap: Dimensionality Reduction - The Compression View



A probabilistic model that maps a low-dim z via a linear mapping to generate a high-dim x

$$egin{array}{rcl} oldsymbol{z}_n &\sim & \mathcal{N}(\mathbf{0}, \mathbf{I}_{\mathcal{K}}) \ oldsymbol{x}_n | oldsymbol{z}_n &\sim & \mathcal{N}(\mathbf{W} oldsymbol{z}_n, \sigma^2 oldsymbol{I}_D) \end{array}$$



A probabilistic model that maps a low-dim z via a linear mapping to generate a high-dim x



A probabilistic model that maps a low-dim z via a linear mapping to generate a high-dim x



通び

A probabilistic model that maps a low-dim z via a linear mapping to generate a high-dim x



Many improvements possible (non-Gaussian distributions, nonlinear mappings, etc)

EN DOC

・ロト ・雪ト ・ヨト・

Recap: Such Models Can Learn to Generate Real-Looking Data..

- Learn the model parameters from training data $\{x_1, \ldots, x_N\}$, e.g., using MLE
- Generate a random z from p(z) and a random new sample x conditioned on that z using p(x|z)



• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\mathbf{x}_n | \mathbf{z}_n)$, or minimize its negative



• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \boldsymbol{z}_n)$, or minimize its negative

$$\mathcal{L}(\mathbf{Z},\mathbf{W},\sigma^2) = \frac{1}{2\sigma^2} \sum_{n=1}^{N} ||\mathbf{x}_n - \mathbf{W}\mathbf{z}_n||^2 + \frac{ND}{2} \log(2\pi\sigma^2)$$



• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \boldsymbol{z}_n)$, or minimize its negative

$$\begin{split} \mathcal{L}(\mathbf{Z},\mathbf{W},\sigma^2) &= \frac{1}{2\sigma^2}\sum_{n=1}^N ||\boldsymbol{x}_n - \mathbf{W}\boldsymbol{z}_n||^2 + \frac{ND}{2}\log(2\pi\sigma^2) \\ &= \frac{1}{2\sigma^2}||\mathbf{X} - \mathbf{Z}\mathbf{W}^\top||_F^2 + \frac{ND}{2}\log(2\pi\sigma^2) \end{split}$$



• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \boldsymbol{z}_n)$, or minimize its negative

$$\begin{split} \mathcal{L}(\mathbf{Z},\mathbf{W},\sigma^2) &= \frac{1}{2\sigma^2}\sum_{n=1}^{N}||\boldsymbol{x}_n-\mathbf{W}\boldsymbol{z}_n||^2+\frac{ND}{2}\log(2\pi\sigma^2) \\ &= \frac{1}{2\sigma^2}||\mathbf{X}-\mathbf{Z}\mathbf{W}^\top||_F^2+\frac{ND}{2}\log(2\pi\sigma^2) \end{split}$$

 \bullet For known $\sigma^2,$ learning PPCA boils down to solve

$$\{\hat{\mathbf{Z}}, \hat{\mathbf{W}}\} = \arg\min_{\mathbf{Z}, \mathbf{W}} ||\mathbf{X} - \mathbf{Z}\mathbf{W}^{\top}||_{F}^{2}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ >

• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \boldsymbol{z}_n)$, or minimize its negative

$$\begin{aligned} \mathcal{L}(\mathbf{Z},\mathbf{W},\sigma^2) &= \frac{1}{2\sigma^2}\sum_{n=1}^{N}||\boldsymbol{x}_n-\mathbf{W}\boldsymbol{z}_n||^2 + \frac{ND}{2}\log(2\pi\sigma^2) \\ &= \frac{1}{2\sigma^2}||\mathbf{X}-\mathbf{Z}\mathbf{W}^\top||_F^2 + \frac{ND}{2}\log(2\pi\sigma^2) \end{aligned}$$

 \bullet For known $\sigma^2,$ learning PPCA boils down to solve

$$\{\hat{\mathbf{Z}}, \hat{\mathbf{W}}\} = \arg\min_{\mathbf{Z}, \mathbf{W}} ||\mathbf{X} - \mathbf{Z}\mathbf{W}^{\top}||_{F}^{2}$$

• Similar to doing matrix factorization of X by minimizing the reconstruction error

(日) (四) (日) (日) (日)

• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \boldsymbol{z}_n)$, or minimize its negative

$$\begin{split} \mathcal{L}(\mathbf{Z},\mathbf{W},\sigma^2) &= \frac{1}{2\sigma^2}\sum_{n=1}^{N}||\boldsymbol{x}_n-\mathbf{W}\boldsymbol{z}_n||^2+\frac{ND}{2}\log(2\pi\sigma^2) \\ &= \frac{1}{2\sigma^2}||\mathbf{X}-\mathbf{Z}\mathbf{W}^\top||_F^2+\frac{ND}{2}\log(2\pi\sigma^2) \end{split}$$

 \bullet For known $\sigma^2,$ learning PPCA boils down to solve

$$\{\hat{\mathbf{Z}}, \hat{\mathbf{W}}\} = \arg\min_{\mathbf{Z}, \mathbf{W}} ||\mathbf{X} - \mathbf{Z}\mathbf{W}^{\top}||_{F}^{2}$$

- Similar to doing matrix factorization of X by minimizing the reconstruction error
- Can solve it using ALT-OPT (Z given W, and W given Z)

• One way: Maximize (conditional) log-likelihood $\sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \boldsymbol{z}_n)$, or minimize its negative

$$\begin{split} \mathcal{L}(\mathbf{Z},\mathbf{W},\sigma^2) &= \frac{1}{2\sigma^2}\sum_{n=1}^{N}||\boldsymbol{x}_n-\mathbf{W}\boldsymbol{z}_n||^2 + \frac{ND}{2}\log(2\pi\sigma^2) \\ &= \frac{1}{2\sigma^2}||\mathbf{X}-\mathbf{Z}\mathbf{W}^\top||_F^2 + \frac{ND}{2}\log(2\pi\sigma^2) \end{split}$$

 \bullet For known $\sigma^2,$ learning PPCA boils down to solve

$$\{\hat{\mathbf{Z}}, \hat{\mathbf{W}}\} = \arg\min_{\mathbf{Z}, \mathbf{W}} ||\mathbf{X} - \mathbf{Z}\mathbf{W}^{\top}||_{F}^{2}$$

- Similar to doing matrix factorization of X by minimizing the reconstruction error
- Can solve it using ALT-OPT (**Z** given **W**, and **W** given **Z**)
- Another (better) way: will be to do a proper MLE on $\log p(x_n)$

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$



[†] Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{R}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$



[†]Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{F}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$

where \mathbf{U}_{K} is $D \times K$ matrix of top K eigvecs of \mathbf{S} ,



・ロト ・回ト ・モト ・モト

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{F}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$

where $\mathbf{U}_{\mathcal{K}}$ is $D \times \mathcal{K}$ matrix of top \mathcal{K} eigvecs of \mathbf{S} , $\mathbf{L}_{\mathcal{K}}$: $\mathcal{K} \times \mathcal{K}$ diagonal matrix of top \mathcal{K} eigvals $\lambda_1, \ldots, \lambda_{\mathcal{K}}$,

[†] Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{F}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$

where \mathbf{U}_{K} is $D \times K$ matrix of top K eigvecs of \mathbf{S} , \mathbf{L}_{K} : $K \times K$ diagonal matrix of top K eigvals $\lambda_{1}, \ldots, \lambda_{K}$, \mathbf{R} is a $K \times K$ arbitrary rotation matrix

A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

[†] Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{F}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$

where \mathbf{U}_{K} is $D \times K$ matrix of top K eigvecs of \mathbf{S} , \mathbf{L}_{K} : $K \times K$ diagonal matrix of top K eigvals $\lambda_{1}, \ldots, \lambda_{K}$, \mathbf{R} is a $K \times K$ arbitrary rotation matrix (equivalent to PCA for $\mathbf{R} = \mathbf{I}$ and $\sigma^{2} \rightarrow 0$)

[†] Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{F}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$

where \mathbf{U}_{K} is $D \times K$ matrix of top K eigvecs of \mathbf{S} , \mathbf{L}_{K} : $K \times K$ diagonal matrix of top K eigvals $\lambda_{1}, \ldots, \lambda_{K}$, \mathbf{R} is a $K \times K$ arbitrary rotation matrix (equivalent to PCA for $\mathbf{R} = \mathbf{I}$ and $\sigma^{2} \rightarrow 0$)

• Need to do eigen-decomposition of $D \times D$ data covariance matrix **S**. EXPENSIVE!!!

[†] Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

• Doing MLE for PPCA requires maximizing

$$\log p(\mathbf{X}|\Theta) = -\frac{N}{2}(D\log 2\pi + \log |\mathbf{C}| + \operatorname{trace}(\mathbf{C}^{-1}\mathbf{S}))$$

where **S** is the data covariance matrix, $\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-1}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^{\top}$ and $\mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^{2}\mathbf{I}$

 \bullet Assuming both ${\bf W}$ and σ^2 as unknowns, their MLE solution is given by

$$\mathbf{W}_{ML} = \mathbf{U}_{K} (\mathbf{L}_{K} - \sigma_{ML}^{2} \mathbf{I})^{1/2} \mathbf{F}$$
$$\sigma_{ML}^{2} = \frac{1}{D - K} \sum_{k=K+1}^{D} \lambda_{k}$$

where \mathbf{U}_{K} is $D \times K$ matrix of top K eigvecs of \mathbf{S} , \mathbf{L}_{K} : $K \times K$ diagonal matrix of top K eigvals $\lambda_{1}, \ldots, \lambda_{K}$, \mathbf{R} is a $K \times K$ arbitrary rotation matrix (equivalent to PCA for $\mathbf{R} = \mathbf{I}$ and $\sigma^{2} \rightarrow 0$)

- Need to do eigen-decomposition of $D \times D$ data covariance matrix **S**. EXPENSIVE!!!
- Also, can't do MLE like this if each x_n has some missing entries

[†] Probabilistic Principal Component Analysis (Tipping and Bishop, 1999)

Intro to Machine Learning (CS771A)

・ロト ・日下・ ・日下・ ・日下

• Using EM for MLE for PPCA has several benefits



メロト メロト メヨト メヨト

Intro to Machine Learning (CS771A)

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition



- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)

イロン イロン イヨン イヨン

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)
- EM does MLE by maximizing the expected CLL

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)
- EM does MLE by maximizing the expected CLL

$$\{\mathbf{W}, \sigma^2\} = \arg \max_{\mathbf{W}, \sigma^2} \mathbb{E}_{\rho(\mathbf{Z}|\mathbf{X}, \mathbf{W}, \sigma^2)}[\log p(\mathbf{X}, \mathbf{Z}|\mathbf{W}, \sigma^2)]$$

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)
- EM does MLE by maximizing the expected CLL

$$\{\mathbf{W}, \sigma^2\} = \arg \max_{\mathbf{W}, \sigma^2} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \mathbf{W}, \sigma^2)}[\log p(\mathbf{X}, \mathbf{Z}|\mathbf{W}, \sigma^2)]$$

• This is done by iterating between the following two steps

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)
- EM does MLE by maximizing the expected CLL

$$\{\mathbf{W}, \sigma^2\} = \arg \max_{\mathbf{W}, \sigma^2} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \mathbf{W}, \sigma^2)}[\log p(\mathbf{X}, \mathbf{Z}|\mathbf{W}, \sigma^2)]$$

- This is done by iterating between the following two steps
 - E Step: For n = 1, ..., N, infer the posterior $p(z_n | x_n)$ given current estimate of $\Theta = (\mathbf{W}, \sigma^2)$

(日)

- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)
- EM does MLE by maximizing the expected CLL

$$\{\mathbf{W}, \sigma^2\} = \arg \max_{\mathbf{W}, \sigma^2} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \mathbf{W}, \sigma^2)}[\log p(\mathbf{X}, \mathbf{Z}|\mathbf{W}, \sigma^2)]$$

- This is done by iterating between the following two steps
 - E Step: For n = 1, ..., N, infer the posterior $p(z_n | x_n)$ given current estimate of $\Theta = (\mathbf{W}, \sigma^2)$

$$p(\boldsymbol{z}_n | \boldsymbol{x}_n, \boldsymbol{\mathsf{W}}, \sigma^2) = \mathcal{N}(\boldsymbol{\mathsf{M}}^{-1} \boldsymbol{\mathsf{W}}^\top \boldsymbol{x}_n, \sigma^2 \boldsymbol{\mathsf{M}}^{-1}) \qquad (\text{where } \boldsymbol{\mathsf{M}} = \boldsymbol{\mathsf{W}}^\top \boldsymbol{\mathsf{W}} + \sigma^2 \boldsymbol{\mathsf{I}}_{\mathcal{K}})$$

(日)
- Using EM for MLE for PPCA has several benefits
 - No need to do expensive eigen-decomposition
 - Works even when x_n may have some missing entries (HW3 has a problem related to this)
- EM does MLE by maximizing the expected CLL

$$\{\mathbf{W}, \sigma^2\} = \arg \max_{\mathbf{W}, \sigma^2} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \mathbf{W}, \sigma^2)}[\log p(\mathbf{X}, \mathbf{Z}|\mathbf{W}, \sigma^2)]$$

- This is done by iterating between the following two steps
 - E Step: For n = 1, ..., N, infer the posterior $p(z_n | x_n)$ given current estimate of $\Theta = (\mathbf{W}, \sigma^2)$

$$p(\boldsymbol{z}_n | \boldsymbol{x}_n, \boldsymbol{\mathsf{W}}, \sigma^2) = \mathcal{N}(\boldsymbol{\mathsf{M}}^{-1} \boldsymbol{\mathsf{W}}^\top \boldsymbol{x}_n, \sigma^2 \boldsymbol{\mathsf{M}}^{-1}) \qquad (\text{where } \boldsymbol{\mathsf{M}} = \boldsymbol{\mathsf{W}}^\top \boldsymbol{\mathsf{W}} + \sigma^2 \boldsymbol{\mathsf{I}}_{\mathcal{K}})$$

• M Step: Maximize the expected CLL $\mathbb{E}[p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. \mathbf{W}, σ^2

イロト 不得 とうせい うけい

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$



イロト イポト イヨト イヨト

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. \mathbf{W} and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^\top\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top]\right]^{-1}$$

メロト メロト メヨト メヨト

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. W and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_{n} \mathbb{E}[\mathbf{z}_{n}]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\top}]\right]^{-1}$$

• To compute **W**, we need two posterior expectations $\mathbb{E}[\boldsymbol{z}_n]$ and $\mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top]$

・ロト ・四ト ・日ト ・日・

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. W and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}]\right]^{-1}$$

- To compute **W**, we need two posterior expectations $\mathbb{E}[\boldsymbol{z}_n]$ and $\mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top]$
- These can be easily obtained from the posterior $p(\mathbf{z}_n | \mathbf{x}_n)$ computed in E step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^{\top} \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^{\top} \mathbf{W} + \sigma^2 \mathbf{I}_K$$

A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. W and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}]\right]^{-1}$$

- To compute **W**, we need two posterior expectations $\mathbb{E}[\boldsymbol{z}_n]$ and $\mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top]$
- These can be easily obtained from the posterior $p(\mathbf{z}_n | \mathbf{x}_n)$ computed in E step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^{\top} \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^{\top} \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^{\top} \mathbf{x}_n$$

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. W and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}]\right]^{-1}$$

- To compute **W**, we need two posterior expectations $\mathbb{E}[\boldsymbol{z}_n]$ and $\mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top]$
- These can be easily obtained from the posterior $p(\mathbf{z}_n | \mathbf{x}_n)$ computed in E step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1}\mathbf{W}^{\top}\mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1}\mathbf{W}^{\top}\mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}] = \mathbb{E}[\mathbf{z}_n]\mathbb{E}[\mathbf{z}_n]^{\top} + \operatorname{cov}(\mathbf{z}_n)$$

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. W and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}]\right]^{-1}$$

- To compute **W**, we need two posterior expectations $\mathbb{E}[\boldsymbol{z}_n]$ and $\mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top]$
- These can be easily obtained from the posterior $p(\mathbf{z}_n | \mathbf{x}_n)$ computed in E step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1}\mathbf{W}^{\top}\mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^{\top}\mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1}\mathbf{W}^{\top}\mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}] = \mathbb{E}[\mathbf{z}_n]\mathbb{E}[\mathbf{z}_n]^{\top} + \operatorname{cov}(\mathbf{z}_n) = \mathbb{E}[\mathbf{z}_n]\mathbb{E}[\mathbf{z}_n]^{\top} + \sigma^2 \mathbf{M}^{-1}$$

• The expected complete data log-likelihood $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$

$$= -\sum_{n=1}^{N} \left\{ \frac{D}{2} \log \sigma^{2} + \frac{1}{2\sigma^{2}} ||\boldsymbol{x}_{n}||^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\boldsymbol{z}_{n}]^{\top} \mathbf{W}^{\top} \boldsymbol{x}_{n} + \frac{1}{2\sigma^{2}} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}] \mathbf{W}^{\top} \mathbf{W}) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[\boldsymbol{z}_{n} \boldsymbol{z}_{n}^{\top}]) \right\}$$

• Taking the derivative of $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$ w.r.t. W and setting to zero

$$\mathbf{W} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}]\right]^{-1}$$

- To compute **W**, we need two posterior expectations $\mathbb{E}[\boldsymbol{z}_n]$ and $\mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top]$
- These can be easily obtained from the posterior $p(\boldsymbol{z}_n | \boldsymbol{x}_n)$ computed in E step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \operatorname{cov}(\mathbf{z}_n) = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

• Note: The noise variance σ^2 can also be estimated (take deriv., set to zero..)

• Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n - \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$



・ロト ・四ト ・ヨト ・ヨト

- Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$
- **E** step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current **W** and σ^2 . Compute exp. for the M step



A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A

- Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$



- Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\boldsymbol{z}_{n}|\boldsymbol{x}_{n}, \boldsymbol{\mathsf{W}}) = \mathcal{N}(\boldsymbol{\mathsf{M}}^{-1}\boldsymbol{\mathsf{W}}^{\top}\boldsymbol{x}_{n}, \sigma^{2}\boldsymbol{\mathsf{M}}^{-1}) \text{ where } \boldsymbol{\mathsf{M}} = \boldsymbol{\mathsf{W}}^{\top}\boldsymbol{\mathsf{W}} + \sigma^{2}\boldsymbol{\mathsf{I}}_{K}$$
$$\mathbb{E}[\boldsymbol{z}_{n}] = \boldsymbol{\mathsf{M}}^{-1}\boldsymbol{\mathsf{W}}^{\top}\boldsymbol{x}_{n}$$
$$\mathbb{E}[\boldsymbol{z}_{n}\boldsymbol{z}_{n}^{\top}] = \operatorname{cov}(\boldsymbol{z}_{n}) + \mathbb{E}[\boldsymbol{z}_{n}]\mathbb{E}[\boldsymbol{z}_{n}]^{\top} = \mathbb{E}[\boldsymbol{z}_{n}]\mathbb{E}[\boldsymbol{z}_{n}]^{\top} + \sigma^{2}\boldsymbol{\mathsf{M}}^{-1}$$



(a)

- Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$\begin{aligned} \rho(\boldsymbol{z}_n | \boldsymbol{x}_n, \boldsymbol{\mathsf{W}}) &= \mathcal{N}(\boldsymbol{\mathsf{M}}^{-1} \boldsymbol{\mathsf{W}}^\top \boldsymbol{x}_n, \sigma^2 \boldsymbol{\mathsf{M}}^{-1}) & \text{where } \boldsymbol{\mathsf{M}} = \boldsymbol{\mathsf{W}}^\top \boldsymbol{\mathsf{W}} + \sigma^2 \boldsymbol{\mathsf{I}}_{\mathcal{K}} \\ \mathbb{E}[\boldsymbol{z}_n] &= \boldsymbol{\mathsf{M}}^{-1} \boldsymbol{\mathsf{W}}^\top \boldsymbol{x}_n \\ \mathbb{E}[\boldsymbol{z}_n \boldsymbol{z}_n^\top] &= \operatorname{cov}(\boldsymbol{z}_n) + \mathbb{E}[\boldsymbol{z}_n] \mathbb{E}[\boldsymbol{z}_n]^\top = \mathbb{E}[\boldsymbol{z}_n] \mathbb{E}[\boldsymbol{z}_n]^\top + \sigma^2 \boldsymbol{\mathsf{M}}^{-1} \end{aligned}$$

• M step: Re-estimate W and σ^2

- Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \operatorname{cov}(\mathbf{z}_n) + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

• M step: Re-estimate W and σ^2

$$\mathbf{W}_{new} = \left[\sum_{n=1}^{N} \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\top}]\right]^{-1}$$

- Specify K, initialize W and σ^2 randomly. Also center the data $(\mathbf{x}_n = \mathbf{x}_n \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \operatorname{cov}(\mathbf{z}_n) + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

• M step: Re-estimate W and σ^2

$$\mathbf{W}_{new} = \left[\sum_{n=1}^{N} \mathbf{x}_{n} \mathbb{E}[\mathbf{z}_{n}]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n}\mathbf{z}_{n}^{\top}]\right]^{-1}$$

$$\sigma_{new}^{2} = \frac{1}{ND} \sum_{n=1}^{N} \left\{ ||\mathbf{x}_{n}||^{2} - 2\mathbb{E}[\mathbf{z}_{n}]^{\top} \mathbf{W}_{new}^{\top} \mathbf{x}_{n} + \operatorname{tr}\left(\mathbb{E}[\mathbf{z}_{n}\mathbf{z}_{n}^{\top}] \mathbf{W}_{new}^{\top} \mathbf{W}_{new}\right) \right\}$$

- Specify K, initialize W and σ^2 randomly. Also center the data $(x_n = x_n \frac{1}{N} \sum_{n=1}^{N} x_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \operatorname{cov}(\mathbf{z}_n) + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

• M step: Re-estimate W and σ^2

$$\mathbf{W}_{new} = \left[\sum_{n=1}^{N} \mathbf{x}_{n} \mathbb{E}[\mathbf{z}_{n}]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\top}]\right]^{-1}$$

$$\sigma_{new}^{2} = \frac{1}{ND} \sum_{n=1}^{N} \left\{ ||\mathbf{x}_{n}||^{2} - 2\mathbb{E}[\mathbf{z}_{n}]^{\top} \mathbf{W}_{new}^{\top} \mathbf{x}_{n} + \operatorname{tr}\left(\mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\top}] \mathbf{W}_{new}^{\top} \mathbf{W}_{new}\right) \right\}$$

• Set $\mathbf{W} = \mathbf{W}_{new}$ and $\sigma^2 = \sigma_{new}^2$. If not converged (monitor $p(\mathbf{X}|\Theta)$), go back to E step

- Specify K, initialize W and σ^2 randomly. Also center the data $(x_n = x_n \frac{1}{N} \sum_{n=1}^{N} x_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \text{ where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \operatorname{cov}(\mathbf{z}_n) + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

• M step: Re-estimate W and σ^2

$$\mathbf{W}_{new} = \left[\sum_{n=1}^{N} \mathbf{x}_{n} \mathbb{E}[\mathbf{z}_{n}]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\top}]\right]^{-1}$$

$$\sigma_{new}^{2} = \frac{1}{ND} \sum_{n=1}^{N} \left\{ ||\mathbf{x}_{n}||^{2} - 2\mathbb{E}[\mathbf{z}_{n}]^{\top} \mathbf{W}_{new}^{\top} \mathbf{x}_{n} + \operatorname{tr}\left(\mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\top}] \mathbf{W}_{new}^{\top} \mathbf{W}_{new}\right) \right\}$$

• Set $\mathbf{W} = \mathbf{W}_{new}$ and $\sigma^2 = \sigma_{new}^2$. If not converged (monitor $p(\mathbf{X}|\Theta)$), go back to E step

 Note: For σ² = 0, this EM algorithm can also be used to efficiently solve standard PCA (note that this EM algorithm doesn't require any eigen-decomposition)

- Specify K, initialize W and σ^2 randomly. Also center the data $(x_n = x_n \frac{1}{N} \sum_{n=1}^{N} x_n)$
- E step: For each *n*, compute $p(\mathbf{z}_n | \mathbf{x}_n)$ using current W and σ^2 . Compute exp. for the M step

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_n$$
$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$
$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \operatorname{cov}(\mathbf{z}_n) + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

• M step: Re-estimate W and σ^2

$$\mathbf{W}_{new} = \left[\sum_{n=1}^{N} \mathbf{x}_{n} \mathbb{E}[\mathbf{z}_{n}]^{\top}\right] \left[\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n}\mathbf{z}_{n}^{\top}]\right]^{-1}$$

$$\sigma_{new}^{2} = \frac{1}{ND} \sum_{n=1}^{N} \left\{ ||\mathbf{x}_{n}||^{2} - 2\mathbb{E}[\mathbf{z}_{n}]^{\top} \mathbf{W}_{new}^{\top} \mathbf{x}_{n} + \operatorname{tr}\left(\mathbb{E}[\mathbf{z}_{n}\mathbf{z}_{n}^{\top}] \mathbf{W}_{new}^{\top} \mathbf{W}_{new}\right) \right\}$$

• Set $\mathbf{W} = \mathbf{W}_{new}$ and $\sigma^2 = \sigma_{new}^2$. If not converged (monitor $p(\mathbf{X}|\Theta)$), go back to E step

- Note: For σ² = 0, this EM algorithm can also be used to efficiently solve standard PCA (note that this EM algorithm doesn't require any eigen-decomposition)
- Missing entries of x_n can be estimated in the E step as $p(x_n^{miss}|x_n^{obs})$

• The PPCA model, for each $\boldsymbol{x}_n, n = 1, \dots, N$, can also be written as

$$\mathbf{x}_n = \mu + \mathbf{W} \mathbf{z}_n + \epsilon_n$$
 where $\epsilon_n \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_D)$



・ロト ・四ト ・ヨト ・ヨト

• The PPCA model, for each $\boldsymbol{x}_n, n = 1, \dots, N$, can also be written as

$$\boldsymbol{x}_n = \mu + \boldsymbol{W} \boldsymbol{z}_n + \epsilon_n$$
 where $\epsilon_n \sim \mathcal{N}(0, \sigma^2 \boldsymbol{I}_D)$

• The marginal distribution is

$$p(\boldsymbol{x}_n) = \mathcal{N}(\mu, \boldsymbol{W} \boldsymbol{W}^\top + \sigma^2 \boldsymbol{I}_D)$$

イロト イロト イモト イモト

• The PPCA model, for each $\boldsymbol{x}_n, n = 1, \dots, N$, can also be written as

$$\boldsymbol{x}_n = \mu + \boldsymbol{W} \boldsymbol{z}_n + \epsilon_n$$
 where $\epsilon_n \sim \mathcal{N}(0, \sigma^2 \boldsymbol{I}_D)$

• The marginal distribution is

$$p(\boldsymbol{x}_n) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\mathsf{W}}\boldsymbol{\mathsf{W}}^\top + \sigma^2 \boldsymbol{\mathsf{I}}_D)$$

• The MLE of μ is simply $\frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n$

• The PPCA model, for each $\boldsymbol{x}_n, n = 1, \dots, N$, can also be written as

$$\boldsymbol{x}_n = \mu + \boldsymbol{W} \boldsymbol{z}_n + \epsilon_n$$
 where $\epsilon_n \sim \mathcal{N}(0, \sigma^2 \boldsymbol{I}_D)$

• The marginal distribution is

$$p(\boldsymbol{x}_n) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\mathsf{W}}\boldsymbol{\mathsf{W}}^\top + \sigma^2 \boldsymbol{\mathsf{I}}_D)$$

- The MLE of μ is simply $\frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n$
- \bullet So we can simply subtract μ from each observation and assume

$$\boldsymbol{x}_n = \boldsymbol{W}\boldsymbol{z}_n + \boldsymbol{\epsilon}_n$$

... and apply PPCA without μ

• Several option to select the "best" K, e.g.,



メロト メポト メヨト メヨト

- Several option to select the "best" K, e.g.,
 - Look at AIC/BIC criteria (NLL + KD or NLL + $K \log D$) and pick the one with smallest K



メロト メロト メヨト メヨト

- Several option to select the "best" K, e.g.,
 - Look at AIC/BIC criteria (NLL + KD or NLL + $K \log D$) and pick the one with smallest K
 - Use sparsity inducing priors on W and/or z_n (set K to some large value; the unnecessary columns of **W** will "turn off" automatically as they will be shrunk to zero during inference)



W

Effect: Only few columns of W will have entries with significant magnitudes

イロト イポト イヨト イヨト



- Several option to select the "best" K, e.g.,
 - Look at AIC/BIC criteria (NLL + KD or NLL + $K \log D$) and pick the one with smallest K
 - Use sparsity inducing priors on **W** and/or z_n (set K to some large value; the unnecessary columns of **W** will "turn off" automatically as they will be shrunk to zero during inference)



Effect: Only few columns of **W** will have entries with significant magnitudes

イロト イヨト イヨト イヨト

• Compute the marginal likelihood (or its approximation) for each K and choose the best model

- Several option to select the "best" K, e.g.,
 - Look at AIC/BIC criteria (NLL + KD or NLL + $K \log D$) and pick the one with smallest K
 - Use sparsity inducing priors on **W** and/or z_n (set K to some large value; the unnecessary columns of **W** will "turn off" automatically as they will be shrunk to zero during inference)



Effect: Only few columns of **W** will have entries with significant magnitudes

A D > A B > A E > A E > \

- Compute the marginal likelihood (or its approximation) for each K and choose the best model
- Nonparametric Bayesian methods (allow K to grow with data)

DOC SE

• Compression/dimensionality reduction is a natural application (use z_n instead of x_n)



イロト イポト イヨト イヨト

- Compression/dimensionality reduction is a natural application (use z_n instead of x_n)
- Also used for learning low-dim. "good" features z_n from high-dim noisy features x_n



< ロ > < 回 > < 回 > < 回 > < 回 >

- Compression/dimensionality reduction is a natural application (use z_n instead of x_n)
- Also used for learning low-dim. "good" features z_n from high-dim noisy features x_n
 - Note that this is different from feature selection (z_n is a transformed version of x_n , not a subset)



・ロト ・回ト ・モト ・モト

- Compression/dimensionality reduction is a natural application (use z_n instead of x_n)
- Also used for learning low-dim. "good" features z_n from high-dim noisy features x_n
 - Note that this is different from feature selection $(z_n \text{ is a transformed version of } x_n, \text{ not a subset})$
- Learning the noise variance enables "image denoising": $\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \epsilon_n$; $\mathbf{W}\mathbf{z}_n$ is the "clean" part



イロト イポト イヨト イヨト

- Compression/dimensionality reduction is a natural application (use z_n instead of x_n)
- Also used for learning low-dim. "good" features z_n from high-dim noisy features x_n
 - Note that this is different from feature selection $(z_n \text{ is a transformed version of } x_n, \text{ not a subset})$
- Learning the noise variance enables "image denoising": $\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \epsilon_n$; $\mathbf{W}\mathbf{z}_n$ is the "clean" part



• Ability to fill-in missing data enables "image inpainting" (left: image with 80% missing data, middle: reconstructed, right: original)



イロト イロト イヨト イヨト

Mixture of PPCA

• May be appropriate if data also exists in clusters (suppose M > 1 clusters)



メロト メポト メヨト メヨト

Intro to Machine Learning (CS771A)

Mixture of PPCA

- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$



イロト イロト イヨト イヨト

Mixture of PPCA

- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$
- Can use *M* such PPCA models $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M$ (one per cluster) for the entire data
- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$
- Can use *M* such PPCA models $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M$ (one per cluster) for the entire data



メロト メぼト メヨト

- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$
- Can use *M* such PPCA models $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M$ (one per cluster) for the entire data



• Mixtures of PPCA can be seen as playing several roles

メロト メポト メヨト メヨト

- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$
- Can use *M* such PPCA models $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M$ (one per cluster) for the entire data



- Mixtures of PPCA can be seen as playing several roles
 - Jointly learning clustering and dimensionality reduction

- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$
- Can use *M* such PPCA models $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M$ (one per cluster) for the entire data



- Mixtures of PPCA can be seen as playing several roles
 - Jointly learning clustering and dimensionality reduction
 - Nonlinear dimensionality reduction

- May be appropriate if data also exists in clusters (suppose M > 1 clusters)
- Data in each cluster (say m) can have its own "local" PPCA model defined by $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}$
- Can use M such PPCA models $\{\mu_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M$ (one per cluster) for the entire data



- Mixtures of PPCA can be seen as plaving several roles
 - Jointly learning clustering and dimensionality reduction
 - Nonlinear dimensionality reduction
 - A flexible probability density model: Mixture of low-rank Gaussians

• For mixture of PPCA, the generative story for each observation x_n is as follows



メロト メポト メヨト メヨト

- For mixture of PPCA, the generative story for each observation x_n is as follows
 - Generate its cluster id as

 $\boldsymbol{c}_n \sim \text{multinoulli}(\pi_1, \ldots, \pi_M)$



イロト イロト イモト イモト

- For mixture of PPCA, the generative story for each observation x_n is as follows
 - Generate its cluster id as

 $\boldsymbol{c}_n \sim \text{multinoulli}(\pi_1, \ldots, \pi_M)$

• Generate latent variable $\boldsymbol{z}_n \in \mathbb{R}^K$ as

 $oldsymbol{z}_n \sim \mathcal{N}(oldsymbol{0}, oldsymbol{\mathsf{I}}_{\mathcal{K}})$



・ロト ・ 日 ト ・ モ ト ・ モ ト

- For mixture of PPCA, the generative story for each observation x_n is as follows
 - Generate its cluster id as

 $\boldsymbol{c}_n \sim \text{multinoulli}(\pi_1, \ldots, \pi_M)$

• Generate latent variable $\boldsymbol{z}_n \in \mathbb{R}^{K}$ as

 $oldsymbol{z}_n \sim \mathcal{N}(oldsymbol{0}, oldsymbol{\mathsf{I}}_{\mathcal{K}})$

• Generate obervation $\boldsymbol{x}_n \in \mathbb{R}^D$ from the \boldsymbol{c}_n^{th} PPCA/FA model

$$\mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\mu_{c_n}} + \mathbf{W_{c_n}} \mathbf{z}_n, \sigma_{c_n}^2 \mathbf{I}_D)$$

A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

- For mixture of PPCA, the generative story for each observation x_n is as follows
 - Generate its cluster id as

 $\boldsymbol{c}_n \sim \text{multinoulli}(\pi_1, \ldots, \pi_M)$

• Generate latent variable $\boldsymbol{z}_n \in \mathbb{R}^K$ as

 $oldsymbol{z}_n \sim \mathcal{N}(oldsymbol{0}, oldsymbol{\mathsf{I}}_{\mathcal{K}})$

• Generate obervation $\boldsymbol{x}_n \in \mathbb{R}^D$ from the \boldsymbol{c}_n^{th} PPCA/FA model

$$\mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\mu_{c_n}} + \mathbf{W_{c_n}} \mathbf{z}_n, \sigma^2_{\mathbf{c_n}} \mathbf{I}_D)$$

• Each PPCA model has its separate mean μ_{c_n} (not needed when M = 1 if data is centered)

(日) (四) (三) (三) (三)

- For mixture of PPCA, the generative story for each observation x_n is as follows
 - Generate its cluster id as

 $\boldsymbol{c}_n \sim \text{multinoulli}(\pi_1, \ldots, \pi_M)$

• Generate latent variable $\boldsymbol{z}_n \in \mathbb{R}^K$ as

 $oldsymbol{z}_n \sim \mathcal{N}(oldsymbol{0}, oldsymbol{\mathsf{I}}_{\mathcal{K}})$

• Generate obervation $\boldsymbol{x}_n \in \mathbb{R}^D$ from the \boldsymbol{c}_n^{th} PPCA/FA model

$$\mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\mu_{c_n}} + \mathbf{W_{c_n}} \mathbf{z}_n, \sigma_{c_n}^2 \mathbf{I}_D)$$

- Each PPCA model has its separate mean μ_{c_n} (not needed when M = 1 if data is centered)
- Exercise: What will be the marginal distribution of x_n , i.e., $p(x_n|\Theta)$?

- For mixture of PPCA, the generative story for each observation x_n is as follows
 - Generate its cluster id as

 $\boldsymbol{c}_n \sim \text{multinoulli}(\pi_1, \ldots, \pi_M)$

• Generate latent variable $\boldsymbol{z}_n \in \mathbb{R}^K$ as

 $oldsymbol{z}_n \sim \mathcal{N}(oldsymbol{0}, oldsymbol{\mathsf{I}}_{\mathcal{K}})$

• Generate obervation $\boldsymbol{x}_n \in \mathbb{R}^D$ from the \boldsymbol{c}_n^{th} PPCA/FA model

$$\mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\mu_{c_n}} + \mathbf{W_{c_n}} \mathbf{z}_n, \sigma^2_{\mathbf{c_n}} \mathbf{I}_D)$$

- Each PPCA model has its separate mean μ_{c_n} (not needed when M = 1 if data is centered)
- Exercise: What will be the marginal distribution of x_n , i.e., $p(x_n|\Theta)$?
- Exercise: Use EM in this model to learn the parameters and latent variables

(Classic) Principal Component Analysis



• A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)



メロト メロト メヨト メヨト

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as



メロト メロト メヨト メヨト

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
 - Learning projection directions that capture maximum variance in data



< ロ > < 回 > < 回 > < 回 > < 回 >

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
 - Learning projection directions that capture maximum variance in data
 - Learning projection directions that result in smallest reconstruction error

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
 - Learning projection directions that capture maximum variance in data
 - Learning projection directions that result in smallest reconstruction error
- Can also be seen as changing the basis in which the data is represented (and transforming the features such that new features become decorrelated)



- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
 - Learning projection directions that capture maximum variance in data
 - Learning projection directions that result in smallest reconstruction error
- Can also be seen as changing the basis in which the data is represented (and transforming the features such that new features become decorrelated)



• Also related to other classic methods, e.g., Factor Analysis (Spearman, 1904)



PCA as Maximizing Variance



- Consider projecting $\boldsymbol{x}_n \in \mathbb{R}^D$ on a one-dim subspace (basically, a line) defined by $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x_n along a one-dim subspace u₁ = u₁[⊤]x_n (location of the green point along the purple line representing u₁)



- Consider projecting $\boldsymbol{x}_n \in \mathbb{R}^D$ on a one-dim subspace (basically, a line) defined by $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x_n along a one-dim subspace u₁ = u₁[⊤]x_n (location of the green point along the purple line representing u₁)



• Mean of projections of all the data: $\frac{1}{N}\sum_{n=1}^{N} u_1^\top x_n = u_1^\top (\frac{1}{N}\sum_{n=1}^{N} x_n) = u_1^\top \mu$

- Consider projecting $\boldsymbol{x}_n \in \mathbb{R}^D$ on a one-dim subspace (basically, a line) defined by $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x_n along a one-dim subspace u₁ = u₁[⊤]x_n (location of the green point along the purple line representing u₁)



- Mean of projections of all the data: $\frac{1}{N}\sum_{n=1}^{N} u_1^\top x_n = u_1^\top (\frac{1}{N}\sum_{n=1}^{N} x_n) = u_1^\top \mu$
- Variance of the projected data ("spread" of the green points)

$$\frac{1}{N}\sum_{n=1}^{N}\left(\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}-\boldsymbol{u}_{1}^{\top}\boldsymbol{\mu}\right)^{2}=\frac{1}{N}\sum_{n=1}^{N}\left\{\boldsymbol{u}_{1}^{\top}(\boldsymbol{x}_{n}-\boldsymbol{\mu})\right\}^{2}=\boldsymbol{u}_{1}^{\top}\boldsymbol{\mathsf{S}}\boldsymbol{u}_{1}$$

- Consider projecting $\boldsymbol{x}_n \in \mathbb{R}^D$ on a one-dim subspace (basically, a line) defined by $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x_n along a one-dim subspace u₁ = u₁[⊤]x_n (location of the green point along the purple line representing u₁)



- Mean of projections of all the data: $\frac{1}{N}\sum_{n=1}^{N} u_1^\top x_n = u_1^\top (\frac{1}{N}\sum_{n=1}^{N} x_n) = u_1^\top \mu$
- Variance of the projected data ("spread" of the green points)

$$\frac{1}{N}\sum_{n=1}^{N}\left(\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}-\boldsymbol{u}_{1}^{\top}\boldsymbol{\mu}\right)^{2}=\frac{1}{N}\sum_{n=1}^{N}\left\{\boldsymbol{u}_{1}^{\top}(\boldsymbol{x}_{n}-\boldsymbol{\mu})\right\}^{2}=\boldsymbol{u}_{1}^{\top}\boldsymbol{\mathsf{S}}\boldsymbol{u}_{1}$$

• **S** is the $D \times D$ data covariance matrix: $\mathbf{s} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^{\top}$. If data already centered $(\boldsymbol{\mu} = 0)$ then $\mathbf{s} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^{\top} = \frac{1}{N} \mathbf{X}^{\top} \mathbf{X}$



• We want \boldsymbol{u}_1 s.t. the variance of the projected data is maximized $\arg\max_{\boldsymbol{u}_1} \, \boldsymbol{u}_1^\top \mathbf{S} \boldsymbol{u}_1$





• We want u_1 s.t. the variance of the projected data is maximized

$$\operatorname{arg\,max}_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \mathbf{S} \boldsymbol{u}_1$$

• To prevent trivial solution (max var. = infinite), assume $||\boldsymbol{u}_1|| = 1 = \boldsymbol{u}_1^\top \boldsymbol{u}_1$

・ロト ・回 ト ・ヨト ・ヨ



• We want \boldsymbol{u}_1 s.t. the variance of the projected data is maximized

$$\arg \max_{\boldsymbol{u}_1} \ \boldsymbol{u}_1^\top \mathbf{S} \boldsymbol{u}_1$$

- To prevent trivial solution (max var. = infinite), assume $||\boldsymbol{u}_1|| = 1 = \boldsymbol{u}_1^\top \boldsymbol{u}_1$
- We will find \boldsymbol{u}_1 by solving the following constrained opt. problem

$$\arg \max_{\boldsymbol{u}_1} \ \boldsymbol{u}_1^\top \mathbf{S} \boldsymbol{u}_1 + \lambda_1 (1 - \boldsymbol{u}_1^\top \boldsymbol{u}_1)$$

where λ_1 is a Lagrange multiplier

A D D A A B D A B D A B B

• The objective function: $\arg \max_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \boldsymbol{S} \boldsymbol{u}_1 + \lambda_1 (1 - \boldsymbol{u}_1^\top \boldsymbol{u}_1)$



メロト メロト メヨト メヨト

- The objective function: $\arg \max_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \mathbf{S} \boldsymbol{u}_1 + \lambda_1 (1 \boldsymbol{u}_1^\top \boldsymbol{u}_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

 $\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$



・ロト ・ 日 ト ・ モ ト ・ モ ト

- The objective function: $\arg \max_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \mathbf{S} \boldsymbol{u}_1 + \lambda_1 (1 \boldsymbol{u}_1^\top \boldsymbol{u}_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

 $\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$

• Thus u_1 is an eigenvector of **S** (with corresponding eigenvalue λ_1)

< ロ > < 回 > < 回 > < 回 > < 回 >

- The objective function: $\arg \max_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \boldsymbol{S} \boldsymbol{u}_1 + \lambda_1 (1 \boldsymbol{u}_1^\top \boldsymbol{u}_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

 $\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$

- Thus \boldsymbol{u}_1 is an eigenvector of **S** (with corresponding eigenvalue λ_1)
- But which of **S**'s (*D* possible) eigenvectors it is?

< ロ > < 回 > < 回 > < 回 > < 回 >

- The objective function: $\arg \max_{u_1} u_1^\top \mathbf{S} u_1 + \lambda_1 (1 u_1^\top u_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

 $\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$

- Thus \boldsymbol{u}_1 is an eigenvector of **S** (with corresponding eigenvalue λ_1)
- But which of **S**'s (*D* possible) eigenvectors it is?
- Note that since $\boldsymbol{u}_1^\top \boldsymbol{u}_1 = 1$, the variance of projected data is

$$\boldsymbol{u}_1^{\top} \mathbf{S} \boldsymbol{u}_1 = \lambda_1$$

- The objective function: $\arg \max_{u_1} u_1^\top \mathbf{S} u_1 + \lambda_1 (1 u_1^\top u_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

$$\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$$

- Thus \boldsymbol{u}_1 is an eigenvector of **S** (with corresponding eigenvalue λ_1)
- But which of **S**'s (*D* possible) eigenvectors it is?
- Note that since $\boldsymbol{u}_1^\top \boldsymbol{u}_1 = 1$, the variance of projected data is

$$\boldsymbol{u}_1^{\top} \mathbf{S} \boldsymbol{u}_1 = \lambda_1$$

• Var. is maximized when u_1 is the (top) eigenvector with largest eigenvalue

- The objective function: $\arg \max_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \boldsymbol{S} \boldsymbol{u}_1 + \lambda_1 (1 \boldsymbol{u}_1^\top \boldsymbol{u}_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

$$\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$$

- Thus \boldsymbol{u}_1 is an eigenvector of **S** (with corresponding eigenvalue λ_1)
- But which of **S**'s (*D* possible) eigenvectors it is?
- Note that since $\boldsymbol{u}_1^\top \boldsymbol{u}_1 = 1$, the variance of projected data is

$$\boldsymbol{u}_1^{\top} \mathbf{S} \boldsymbol{u}_1 = \lambda_1$$

- Var. is maximized when u_1 is the (top) eigenvector with largest eigenvalue
- The top eigenvector u_1 is also known as the first Principal Component (PC)

- The objective function: $\arg \max_{\boldsymbol{u}_1} \boldsymbol{u}_1^\top \boldsymbol{S} \boldsymbol{u}_1 + \lambda_1 (1 \boldsymbol{u}_1^\top \boldsymbol{u}_1)$
- Taking the derivative w.r.t. \boldsymbol{u}_1 and setting to zero gives

$$\mathbf{S}\boldsymbol{u}_1 = \lambda_1 \boldsymbol{u}_1$$

- Thus \boldsymbol{u}_1 is an eigenvector of **S** (with corresponding eigenvalue λ_1)
- But which of **S**'s (*D* possible) eigenvectors it is?
- Note that since $\boldsymbol{u}_1^\top \boldsymbol{u}_1 = 1$, the variance of projected data is

$$\boldsymbol{u}_1^{\top} \mathbf{S} \boldsymbol{u}_1 = \lambda_1$$

- Var. is maximized when u_1 is the (top) eigenvector with largest eigenvalue
- The top eigenvector u_1 is also known as the first Principal Component (PC)
- Other directions can also be found likewise (with each being orthogonal to all previous ones) using the eigendecomposition of S (this is PCA)

• Center the data (subtract the mean $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$ from each data point)


- Center the data (subtract the mean $\mu = \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{x}_n$ from each data point)
- $\bullet\,$ Compute the covariance matrix ${\boldsymbol{S}}$ using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^{ op} \mathbf{X}$$

(日) (四) (三) (三) (三)

- Center the data (subtract the mean $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$ from each data point)
- Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = rac{1}{N} \mathbf{X}^{ op} \mathbf{X}$$

 $\bullet\,$ Do an eigendecomposition of the covariance matrix ${\bf S}\,$

- Center the data (subtract the mean $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$ from each data point)
- Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = rac{1}{N} \mathbf{X}^{ op} \mathbf{X}$$

- $\bullet\,$ Do an eigendecomposition of the covariance matrix ${\bf S}$
- Take first K leading eigenvectors $\{\boldsymbol{u}_k\}_{k=1}^{K}$ with eigenvalues $\{\lambda_k\}_{k=1}^{K}$

(日)

- Center the data (subtract the mean $\mu = \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{x}_n$ from each data point)
- Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = rac{1}{N} \mathbf{X}^{ op} \mathbf{X}$$

- $\bullet\,$ Do an eigendecomposition of the covariance matrix ${\bf S}$
- Take first K leading eigenvectors $\{\boldsymbol{u}_k\}_{k=1}^{K}$ with eigenvalues $\{\lambda_k\}_{k=1}^{K}$
- The final K dim. projection/embedding of data is given by

 $\mathbf{Z} = \mathbf{X}\mathbf{U}$

where $\mathbf{U} = [\mathbf{u}_1 \ \dots \ \mathbf{u}_K]$ is $D \times K$ and embedding matrix \mathbf{Z} is $K \times N$

(日) (四) (三) (三) (三)



• We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$



- We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$
- $\mathbf{U} = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]$ is $N \times N$, each $\boldsymbol{u}_n \in \mathbb{R}^N$ a left singular vector of \mathbf{X}



- We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$
- $\mathbf{U} = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]$ is $N \times N$, each $\boldsymbol{u}_n \in \mathbb{R}^N$ a left singular vector of \mathbf{X}
 - U is orthonormal: $u_n^\top u_{n'} = 0$ for $n \neq n'$, and $u_n^\top u_n = 1 \Rightarrow UU^\top = I_N$



- We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$
- $\mathbf{U} = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]$ is $N \times N$, each $\boldsymbol{u}_n \in \mathbb{R}^N$ a left singular vector of \mathbf{X}
 - U is orthonormal: $u_n^{\top} u_{n'} = 0$ for $n \neq n'$, and $u_n^{\top} u_n = 1 \Rightarrow UU^{\top} = I_N$
- Λ is $N \times D$ with only min(N, D) diagonal entries (all positive) singular values (decreasing order)

< ロ > < 回 > < 回 > < 回 > < 回 >



- We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$
- $\mathbf{U} = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]$ is $N \times N$, each $\boldsymbol{u}_n \in \mathbb{R}^N$ a left singular vector of \mathbf{X}
 - **U** is orthonormal: $u_n^{\top} u_{n'} = 0$ for $n \neq n'$, and $u_n^{\top} u_n = 1 \Rightarrow \mathbf{U}\mathbf{U}^{\top} = \mathbf{I}_N$
- A is $N \times D$ with only min(N, D) diagonal entries (all positive) singular values (decreasing order)
- $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_D]$ is $D \times D$, each $\mathbf{v}_d \in \mathbb{R}^D$, a right singular vector of \mathbf{X}



- We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$
- $\mathbf{U} = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]$ is $N \times N$, each $\boldsymbol{u}_n \in \mathbb{R}^N$ a left singular vector of \mathbf{X}
 - U is orthonormal: $u_n^{\top} u_{n'} = 0$ for $n \neq n'$, and $u_n^{\top} u_n = 1 \Rightarrow UU^{\top} = I_N$
- A is $N \times D$ with only min(N, D) diagonal entries (all positive) singular values (decreasing order)
- $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_D]$ is $D \times D$, each $\mathbf{v}_d \in \mathbb{R}^D$, a right singular vector of \mathbf{X}
 - V is orthonormal: $\mathbf{v}_d^\top \mathbf{v}_{d'} = 0$ for $d \neq d'$, and $\mathbf{v}_d^\top \mathbf{v}_d = 1 \Rightarrow \mathbf{V}\mathbf{V}^\top = \mathbf{I}_D$



- We can represent any matrix **X** of size $N \times D$ using SVD as $\mathbf{X} = \mathbf{U} \wedge \mathbf{V}^{\top}$
- $\mathbf{U} = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]$ is $N \times N$, each $\boldsymbol{u}_n \in \mathbb{R}^N$ a left singular vector of \mathbf{X}
 - U is orthonormal: $u_n^{\top} u_{n'} = 0$ for $n \neq n'$, and $u_n^{\top} u_n = 1 \Rightarrow UU^{\top} = I_N$
- A is $N \times D$ with only min(N, D) diagonal entries (all positive) singular values (decreasing order)
- $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_D]$ is $D \times D$, each $\mathbf{v}_d \in \mathbb{R}^D$, a right singular vector of \mathbf{X}
 - V is orthonormal: $\mathbf{v}_d^\top \mathbf{v}_{d'} = 0$ for $d \neq d'$, and $\mathbf{v}_d^\top \mathbf{v}_d = 1 \Rightarrow \mathbf{V}\mathbf{V}^\top = \mathbf{I}_D$

• Note: If X is symmetric then it is known as eigenvalue decomposition (and U = V in that case)

Low-Rank Approximation via SVD

• Can also expand the SVD expression as

$$\mathbf{X} = \sum_{k=1}^{\min(N,D)} \lambda_k \boldsymbol{u}_k \boldsymbol{v}_k^\top$$



Low-Rank Approximation via SVD

• Can also expand the SVD expression as

$$\mathbf{X} = \sum_{k=1}^{\min(N,D)} \lambda_k \boldsymbol{u}_k \boldsymbol{v}_k^ op$$

• Can write a rank-K approximation of **X** (where $K \ll \min(N, D)$) as

$$\mathbf{X} \approx \mathbf{\hat{X}} = \sum_{k=1}^{K} \lambda_{k} \mathbf{u}_{k} \mathbf{v}_{k}^{\top} \neq \mathbf{U}_{K} \Lambda_{K} \mathbf{V}_{K}^{\top}$$

$$\approx \mathbf{v}_{K} \mathbf{v}_{K} \mathbf{v}_{K}^{\top}$$

メロト メロト メヨト メヨト

• A linear projection method



イロト イポト イヨト イヨト

Intro to Machine Learning (CS771A)

- A linear projection method
 - Won't work well if data can't be approximated by a linear subspace

- A linear projection method
 - Won't work well if data can't be approximated by a linear subspace
 - But PCA/PPCA can be kernelized (Kernel PCA or Gaussian Process Latent Variable Models)



- A linear projection method
 - Won't work well if data can't be approximated by a linear subspace
 - But PCA/PPCA can be kernelized (Kernel PCA or Gaussian Process Latent Variable Models)
- Variance based projection directions can sometimes be suboptimal (e.g., if we want to preserve class separation, e.g., when doing classification)



- A linear projection method
 - Won't work well if data can't be approximated by a linear subspace
 - But PCA/PPCA can be kernelized (Kernel PCA or Gaussian Process Latent Variable Models)
- Variance based projection directions can sometimes be suboptimal (e.g., if we want to preserve class separation, e.g., when doing classification)



- PCA relies on eigendecomposition of an $D \times D$ covariance matrix
 - Can be slow if done naïvely. Takes $O(D^3)$ time
 - Many faster methods exists (e.g., Power Method)

- A linear projection method
 - Won't work well if data can't be approximated by a linear subspace
 - But PCA/PPCA can be kernelized (Kernel PCA or Gaussian Process Latent Variable Models)
- Variance based projection directions can sometimes be suboptimal (e.g., if we want to preserve class separation, e.g., when doing classification)



- PCA relies on eigendecomposition of an $D \times D$ covariance matrix
 - Can be slow if done naïvely. Takes $O(D^3)$ time
 - Many faster methods exists (e.g., Power Method)
 - Note: PPCA doesn't suffer from this issue (EM can be very efficient!)

A D > A D >

- How to compute singular vectors (SVD) power method
- Nonlinear Dimensionality Reduction
- Supervised Dimensionality Reduction
- Dimensionality Reduction for Visualization

< ロ > < 回 > < 回 > < 回 > < 回 >