

- today:
- Gaussian networks
  - factor analysis & PCA
  - VAE

Gaussian networks

$$X \sim N(\mu, \Sigma) \quad \mu \in \mathbb{R}^p \quad \Sigma \in \mathbb{R}^{p \times p}, \Sigma > 0$$

$$p(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

put it in exponential family

sufficient statistics

$$T(x) = \begin{pmatrix} x \\ -\frac{xx^T}{2} \end{pmatrix}$$

canonical parameter

$$\begin{aligned} & \text{tr}\left(\Sigma^{-1} (x-\mu)(x-\mu)^T\right) \\ & \downarrow \\ & (xx^T - \mu x^T - x \mu^T + \mu \mu^T) \\ & \left\langle \Sigma^{-1}, -\frac{xx^T}{2} \right\rangle + \left\langle \Sigma^{-1} \mu, x \right\rangle - \frac{1}{2} \mu^T \Sigma^{-1} \mu \\ & \downarrow \\ & \eta \quad \mu = \Sigma \eta = \Lambda^{-1} \eta \\ & \downarrow \\ & \Lambda \triangleq \Sigma^{-1} \\ & \text{precision matrix} \end{aligned}$$

canonic parameter  $\tilde{\eta} \left( \begin{pmatrix} \mu \\ \Sigma \end{pmatrix} \right) = \begin{pmatrix} \eta \\ \Lambda \end{pmatrix} = \begin{pmatrix} \Sigma^{-1} \mu \\ \Sigma^{-1} \end{pmatrix}$

$$p(x; \eta, \Lambda) = \exp\left(\eta^T x + \left\langle \Lambda, -\frac{xx^T}{2} \right\rangle - \underbrace{\left[ \frac{1}{2} \eta^T \Lambda^{-1} \eta + \frac{p}{2} \log(2\pi) - \frac{1}{2} \log |\Lambda| \right]}_{A(\eta, \Lambda)}\right)$$

$$\Lambda = \left\{ (\eta, \Lambda) : \eta \in \mathbb{R}^p, \Lambda > 0, \Lambda = \Lambda^T, \Lambda \in \mathbb{R}^{p \times p} \right\}$$

useful exercise:  $\nabla_{\eta} A(\eta, \Lambda) = \mathbb{E}[x] = \mu$

$\nabla_{\Lambda} A(\eta, \Lambda) = \mathbb{E}\left[-\frac{xx^T}{2}\right]$

UGM viewpoint:

$$p(x; \eta, \Lambda) = \exp\left(-\frac{1}{2} \sum_{i,j} \Lambda_{ij} x_i x_j + \sum_i \eta_i x_i - A(\eta, \Lambda)\right)$$

$p \in \mathcal{G}(G)$  where  $E \triangleq \{ \xi_{ij} \}$  st.  $\Lambda_{ij} \neq 0$

zeros in precision matrix  $\Rightarrow$  cond. indep. properties  $\otimes$

"Gaussian network"  $p(x) = \prod_{i: \xi_{ij} \in E} \psi_{ij}(x_i, x_j) \prod_i \psi_i(x_i)$

quick Schur-complement digression:

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{11}^{-1} + \Sigma_{11}^{-1} \Sigma_{12} M^{-1} \Sigma_{21} \Sigma_{11}^{-1} & -\Sigma_{11}^{-1} \Sigma_{12} M^{-1} \\ -\Sigma_{11}^{-1} \Sigma_{21} \Sigma_{11}^{-1} & M^{-1} \end{pmatrix}$$

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{11}^{-1} + \Sigma_{11}^{-1} \Sigma_{12} M^{-1} \Sigma_{21} \Sigma_{11}^{-1} & -\Sigma_{11}^{-1} \Sigma_{12} M^{-1} \\ -M^{-1} \Sigma_{21} \Sigma_{11}^{-1} & M^{-1} \end{pmatrix}$$

$$M \triangleq \Sigma / \Sigma_{11} \triangleq \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$$

"Schur complement of  $\Sigma$ "  
w. r. to  $\Sigma_{11}$

$$\Sigma / \Sigma_{22} \triangleq \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

\* use this to derive the "Woodbury-Sherman-Morrison inversion formula"

property:  $|\Sigma| = |\Sigma_{11}| \cdot |\Sigma / \Sigma_{11}| = |\Sigma_{22}| \cdot |\Sigma / \Sigma_{22}|$

$$p(\underbrace{x_1}_{\text{dim } p_1}, \underbrace{x_2}_{\text{dim } p_2}) = \frac{1}{\sqrt{(2\pi)^{p_1} |\Sigma_{11}|}} \exp\left(-\frac{1}{2} (x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1)\right) \cdot \left. \frac{1}{\sqrt{(2\pi)^{p_2} |\Sigma / \Sigma_{11}|}} \exp\left(-\frac{1}{2} (x_2 - \mu_2 - b(x_1))^T (\Sigma / \Sigma_{11})^{-1} (x_2 - \mu_2 - b(x_1))\right) \right\} p(x_1, x_2)$$

where  $b(x_1) \triangleq \Sigma_{21} \Sigma_{11}^{-1} (x_1 - \mu_1)$

mean parametrization  
of marginal on  $x_1$   
and conditional  $x_2 | x_1$

$$\mu_1^M = \mu_1$$

$$\Sigma_{11}^M = \Sigma_{11}$$

super simple! } param. of marginal on  $x_1$

$$\mu_{2|1}^{\text{cond.}} = \mu_2 + b(x_1)$$

$$\Sigma_{21}^{\text{cond.}} = \Sigma / \Sigma_{11} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \quad \left. \vphantom{\Sigma_{21}^{\text{cond.}}} \right\} \text{param. for cond. } x_2 | x_1$$

in canonical param.

$$\Lambda_{21}^{\text{cond.}} = \Lambda_{22} \quad (\text{simple})$$

$$m_{2|1}^{\text{cond.}} = m_2 - \Lambda_{21} x_1$$

$$m_1^M = m_1 - \Lambda_{12} \Lambda_{22}^{-1} m_2 \quad (\text{more complicated})$$

$$\Lambda_{11}^M = \Lambda_{11} - \Lambda_{12} \Lambda_{22}^{-1} \Lambda_{21} = \Lambda_{11} / \Lambda_{22}$$

for example: block  $\underbrace{\{i, j\}}_I$  | rest

$$\Lambda = \begin{pmatrix} \Lambda_{II} & \Lambda_{I\bar{I}} \\ \Lambda_{\bar{I}I} & \Lambda_{\bar{I}\bar{I}} \end{pmatrix}$$

$$\text{Cov}(X_I | X_{\text{rest}}) = \Sigma_{I|\text{rest}} = \Lambda_{I|\text{rest}}^{-1} = \Lambda_{II}^{-1} = \begin{pmatrix} \Lambda_{ii} & \Lambda_{ij} \\ \Lambda_{ji} & \Lambda_{jj} \end{pmatrix}^{-1}$$

if  $\Lambda_{ij} = 0$  get  $\Sigma_{II|\text{rest}} = \begin{pmatrix} \Lambda_{ii}^{-1} & 0 \\ 0 & \Lambda_{jj}^{-1} \end{pmatrix}$

$$\Rightarrow X_i \perp\!\!\!\perp X_j | X_{\text{rest}}$$

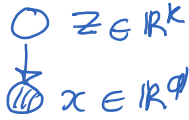
(also true by Markov property of UGM)

(also true by Markov property of UGM)

17h34

Factor analysis:

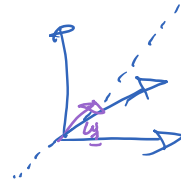
latent variable model



learn "latent representation"  
 or  
 dimensionality reduction  $k \ll d$

PCA for dimensionality reduction

synthetic view: find  $k$  orthonormal vectors in  $\mathbb{R}^d$   $w_1, \dots, w_k$   
 s.t. projection of  $x$  on  $\text{span}\{w_1, \dots, w_k\}$   
 is a good approx. of  $x$



$W = [w_1 \dots w_k]$   $W^T W = I_k$  (by orthonormality)  
 $d \times k$   $1 \times 1$   
 $W W^T \neq Id$

$P_w \triangleq W W^T$   $P_w^2 = W W^T W W^T = P_w$

$\hookrightarrow$  orthogonal projection on  $\text{span}\{w_1, \dots, w_k\}$

$P_w x = W W^T x$

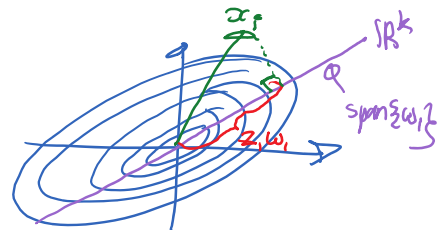
$= \begin{pmatrix} w_1^T \\ \vdots \\ w_k^T \end{pmatrix} \begin{pmatrix} \langle w_1, x \rangle \\ \vdots \\ \langle w_k, x \rangle \end{pmatrix}$   
 $= \sum_k w_k \langle w_k, x \rangle = W z$   
 $(z)_k$

PCA  $\min_{W \in \mathbb{R}^{d \times k}, W^T W = I_k} \sum_i \|x_i - W W^T x_i\|_2^2$   
 $\text{col}(W) \triangleq$  principal subspace

$z = W^T x$   
 $\hookrightarrow$  lower dimensional representation

$X = \begin{pmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{pmatrix}$   
 $n \times d$

$\|X^T - W W^T X^T\|_F^2$   
 $= \|(Id - P_w) X^T\|_F^2$   
 $= \text{tr}(X (Id - P_w)^T (Id - P_w) X^T)$   
 $= \text{tr}(X (Id - P_w) X^T) = \text{tr}(X^T X (Id - P_w))$



$W$  is not unique, only  $\text{col}(W)$   
 e.g.  $\tilde{W} = W R$  where  $R^T R = I_k$   
 then  $\tilde{W} \tilde{W}^T = W R R^T W^T = W W^T$   
 $I_k = W^T W$

$\frac{1}{n} X^T X = \frac{1}{n} \sum_i I_i x_i^T x_i$   
 $\uparrow$   
 empirical covariance of  $x$  when  $\sum x_i = 0$  (mean=0)

min rec. error  $\Leftrightarrow$  maximize  $\text{tr}(X^T X W W^T) = \sum_k w_k^T X^T X w_k$   
 "analysis view of PCA" max sum of empirical variances of new representation

(computation of PCA  $\rightarrow$  top  $k$  e-vectors of  $X^T X$ )

factor analysis  $\rightarrow$  simplest generative model

$z \sim N(0, I_k)$ , noise

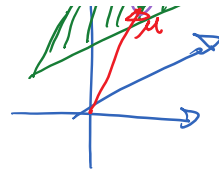


$$z \sim N(0, I_k)$$

$$x = Wz + \mu + \epsilon$$

$$\epsilon \perp z, \epsilon \sim N(0, D)$$

data  
diagonal matrix



$$x|z \sim N(Wz + \mu, D)$$

$p(x)$  is Gaussian

$$E[x] = E[E[x|z]] = 0 + \mu = \mu$$

$$\text{cov}(x, x) = \text{cov}(Wz + \mu + \epsilon, Wz + \mu + \epsilon)$$

indep.

$$= \text{cov}(Wz, Wz) + \text{cov}(\epsilon, \epsilon)$$

$$= W \underbrace{\text{cov}(z, z)}_{I_k} W^T + D$$

$$= WW^T + D$$

equivalent model on  $X \sim N(\mu, WW^T + D)$

low rank covariance piece      diagonal  $\rightarrow$  d degrees of freedom

estimate  $W, D, \mu$  by MLE

$\leadsto$  do EM (latent variable model)

get  $p(z|x) \rightarrow$  Gaussian with mean

$$E[z|x] = W^T(WW^T + D)^{-1}(x - \mu_x)$$

probabilistic PCA: special case of factor analysis where suppose  $D = \sigma^2 I$

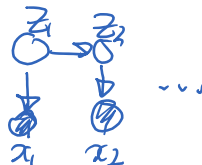
$$\lim_{\sigma \rightarrow 0} W^T(WW^T + \sigma^2 I)^{-1} = W^T \leftarrow \text{pseudoinverse}$$

$$= W^T \quad \text{if } W^T W = I_k$$

this suggests PCA is limit of PPCA as  $\sigma \rightarrow 0$

### Kalman filter

factor analysis



move to state space model: unroll in time (HMM style)

Kalman filter:  $z_t | z_{t-1} \sim N(Az_{t-1}, B)$

$\rightarrow$  define "unobserved" data in HMM

$(z_t | x)$

Kalman filter:  $z_t | z_{t-1} \sim N(Az_{t-1}, B)$

→ doing "sum-product" alg. in HMM

$p(z_t | x_{1:t})$

get "Kalman filter alg."

variational auto-encoder:

generalization of factor analysis

$W$



$z \sim N(0, I_K)$

diagonal noise

$x|z \sim N(\mu_w(z), \sigma_w^2(z))$

where  $\mu_w(z) \leftarrow$  output of NN

"decoder"

MLE → use EM

↳  $p(z|x)$  is intractable  $\Rightarrow$  approximate with variational approach

approximate  $p(z|x)$  with  $q_\phi(z|x)$

$z|x \sim N(\mu_\phi(x), \sigma_\phi^2(x))$

output of a NN

"encoder"

in EM  $\log p(x) \approx \mathbb{E}_q[\log p(x,z)] + H(q)$

$= \mathbb{E}_{q_\phi(z|x)}[\log p_w(x|z)] - KL(q_\phi(z|x) || p(z))$

$N(0, I_K)$

allows "reparameterization trick"

$z|x \rightarrow \mu_\phi(x) + \sigma_\phi^2(x) \cdot \epsilon$   
 $\epsilon \sim N(0, 1)$

o VAE innovations:

- share parameters  $\phi$  among data points for their variational approximation  $q_\phi(z|x)$
- re-parameterization trick to only have parameters appear in simple deterministic transformation, stochasticity is all left in  $N(0,1)$  noise variables (no parameters)  $\Rightarrow$  allow simple backpropagation of gradient through expectations
- for more details, see: [Slides on VAE](#) by Aaron Courville - deep learning class Winter 2017

Other skipped parts, for more details:

- see [2016 lecture 17 scribbles](#) for more info on Schur complement & block decomposition of inverse
- see [2016 lecture 18 scribbles](#) for more info on SVD, and also CCA