

today: Gaussian networks
factor analysis & PCA
• VAE

Gaussian networks

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

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

put it in exponential family

$$\leq \Sigma^{-1} \boxed{\frac{-x^T x}{2}} + \underbrace{\Sigma^{-1} \mu}_{\triangleq n} \boxed{x} - \frac{1}{2} \mu^T \Sigma^{-1} \mu$$

sufficient statistic

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

canonical parameter
 $\Lambda \triangleq \Sigma^{-1}$
precision matrix

$$\mu = \Sigma n = \Lambda^{-1} n$$

$$\text{canonical parameter } \tilde{n}\left(\begin{pmatrix} \mu \\ \Sigma \end{pmatrix}\right) = \begin{pmatrix} \eta \\ \Delta \end{pmatrix} = \begin{pmatrix} \Sigma^{-1} \mu \\ \Sigma^{-1} \end{pmatrix}$$

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

$$\Omega = \{(\eta, \Delta) : \eta \in \mathbb{R}^P, \Delta \geq 0, \Delta = \Delta^T, \Delta \in \mathbb{R}^{P \times P}\}$$

useful exercise: check $\mathbb{E}_\eta A(\eta, \Delta) = \mathbb{E}[x] = \mu$

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

UGM viewpoint

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

$\rho \in \mathcal{S}(G)$ where $E \triangleq \{i, j\} \text{ s.t. } \lambda_{ij} \neq 0\}$

zeros in precision matrix \Rightarrow cond. indep. properties

$$\text{"Gaussian network"} \quad p(x) = \prod_{i,j} \frac{1}{Z_{i,j}} \pi_{i,j}(x_i, x_j) \prod_i \pi_i(x_i)$$

quick Schur-complement digression

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

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

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

"Schur complement of Σ "
w.r.t. to Σ_{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: $|\Sigma| = |\Sigma_{11}| \cdot |\Sigma / \Sigma_{11}| = |\Sigma_{22}| \cdot (\Sigma / \Sigma_{22})$

$$p(x_1, x_2) = \frac{1}{\sqrt{(2\pi)^{\dim p_1} |\Sigma_{11}|}} \exp\left(-\frac{1}{2} (x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1)\right) \cdot p(x_1)$$

$$\frac{1}{\sqrt{(2\pi)^{\dim p_2} |\Sigma / \Sigma_{11}|}} \exp\left(-\frac{1}{2} (x_2 - \mu_2 - b(x_1))^T (\Sigma / \Sigma_{11})^{-1} (x_2 - \mu_2 - b(x_1))\right) p(x_2 | x_1)$$

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

mean parameterization
of marginal on x_1
and conditional $x_2 | x_1$, $\begin{cases} \mu_1^m = \mu_1 \\ \Sigma_1^m = \Sigma_{11} \end{cases}$ } super simple! } param.
for marginal on x_1

$\begin{cases} \mu_{2|1}^{\text{cond.}} = \mu_2 + b(x_1) \\ \Sigma_{2|1}^{\text{cond.}} = \Sigma / \Sigma_{11} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \end{cases}$ } param.
for cond. $x_2 | x_1$

in canonical param. $\Lambda_{2|1}^{\text{cond.}} = -\Lambda_{22}$ (simple)

$$\Lambda_{2|1}^{\text{cond.}} = \Lambda_2 - \Lambda_{21} x_1$$

$$\Lambda_1^m = \Lambda_1 - \Lambda_{12} \Lambda_{22}^{-1} \Lambda_{21} \quad (\text{more complicated})$$

$$\Lambda_1^m = \Lambda_{11} - \Lambda_{12} \Lambda_{22}^{-1} \Lambda_{21} = 1/\Lambda_{22}$$

for example: block $\begin{matrix} \Sigma_{i,j} \\ \hline I \end{matrix}$ | rest

$$\text{cov}(X_I | X_{\text{rest}}) = \Sigma_{I|\text{rest}} = \Lambda_{I|\text{rest}}^{-1} = \Lambda_{II}^{-1} = (\Lambda_{ii} \ \Lambda_{ij} \\ \Lambda_{ji} \ \Lambda_{jj})^{-1}$$

if $\Lambda_{ij} = 0$ $\Lambda_{II} = \begin{pmatrix} \Lambda_{ii} & 0 \\ 0 & \Lambda_{jj} \end{pmatrix}$

$$\text{if } \lambda_{ij} = 0 \quad \Lambda_{II} = \begin{pmatrix} \Lambda_{ii} & 0 \\ 0 & \Lambda_{jj} \end{pmatrix}$$

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

$$\Rightarrow \boxed{x_i \perp\!\!\!\perp x_j \mid x_{rest}}$$

(also true by Markov property of UGM)

15h28

Factor analysis:

latent variable model

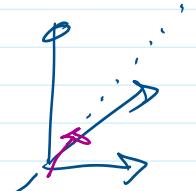
$$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array} \quad z \in \mathbb{R}^k$$

learn "latent representation" in \mathbb{R}^k

or
dimensionality reduction $k < d$

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



$$W = \begin{bmatrix} w_1 & \dots & w_k \end{bmatrix} \quad W^T W = I_k \quad (\text{by orthonormality})$$

$$WW^T \neq I_d$$

$$P_w \triangleq WW^T \quad P_w^2 = WW^T W W^T = P_w \quad P_w x = W(W^T x)$$

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

$$= \begin{pmatrix} w_1 & \dots & w_k \end{pmatrix} \begin{pmatrix} \langle w_1, x \rangle \\ \vdots \\ \langle w_k, x \rangle \end{pmatrix}$$

$$= \sum_k w_k \langle w_k, x \rangle = Wz$$

$$z = W^T x \quad \hookrightarrow \text{lower-dimensional representation}$$

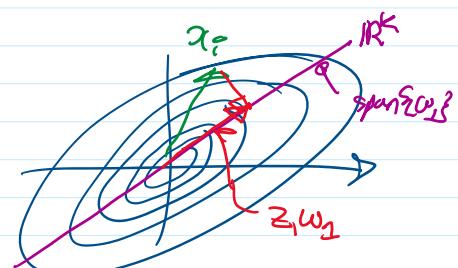
PCA $\min_{W \in \mathbb{R}^{d \times k}} \sum_{i=1}^n \|x_i - W W^T x_i\|_2^2$

$W^T W = I_k$ $\text{col}(W) \triangleq \text{principal subspace}$

$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$

$$\begin{aligned} & \|X^T - WW^T X^T\|_F^2 \\ &= \|(Id - WW^T)X^T\|_F^2 \end{aligned}$$

$$= \text{tr}(\underbrace{X(Id - WW^T)^T}_{\perp X^T X} (Id - WW^T) X^T)$$



$$\frac{1}{n} \sum_i x_i x_i^T = \frac{1}{n} \sum_i x_i x_i^T$$

empirical covariance
of x when
 $\sum_i x_i = 0$
(mean = 0)

$$= \text{tr} \left(X \underbrace{(\text{Id} - W W^T)^T}_{\text{Id} - P_W} (\text{Id} - W W^T) X^T \right)$$

$$= \text{tr} \left(X (\text{Id} - P_W) X^T \right) = \text{tr} \left(X^T X (\text{Id} - P_W) \right) = \text{tr} (\text{cst.}) - \text{tr} \left(X^T X P_W \right)$$

min rec. error \Leftrightarrow

$$\text{maximizing } \text{tr}(X^T X W W^T) = \sum_k w_k^T X^T X w_k$$

"analyzing rows of
PCA"

max sum of
empirical covariances
of new representation

④ W is not unique, only $\text{col}(W)$

$$\text{e.g., } \tilde{W} = W R \text{ where } R^T R = R R^T = I_K$$

$$\tilde{W} \tilde{W}^T = W R R^T W^T$$

Factor analysis \rightarrow simplest generative model

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

$$x = Wz + \mu + \varepsilon$$

noise

$\varepsilon \perp \! \! \! \perp z$, $\varepsilon \sim N(0, D)$

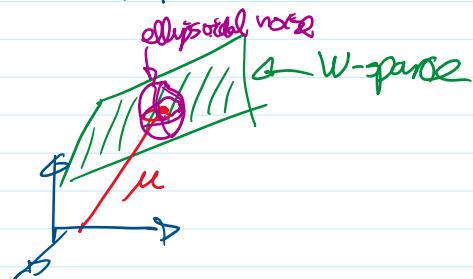
diag matrix

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

$p(x)$ is Gaussian

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

$$\max_{\|w\|=1} w^T z = \lambda_{\max} \|w\|^2$$



$$\mathbb{E}[x] = \mathbb{E}[\underbrace{\mathbb{E}[x|z]}_{Wz + \mu}] = 0 + \mu = \mu$$

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

indep

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

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

$$= WW^T + D$$

equivalent model on $\boxed{x \sim N(\mu, WW^T + D)}$

low rank covariance prior

diagonal $\rightarrow d$ degrees of freedom

estimate $W \& D \& \mu$ by MLE

\rightsquigarrow do EM (latent variable model)

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

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

$$\mathbb{E}[z|x] = W^T(WW^T + D)^{-1}(x - \underbrace{\mu_x}_{\text{mean of } z})$$

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

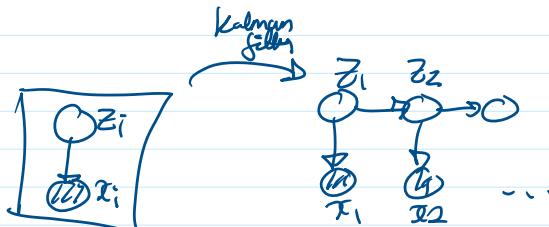
$$\lim_{\sigma \rightarrow 0} W^T(WW^T + \sigma^2 I)^{-1} = W^T \leftarrow \text{pseudo-inverse}$$

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

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

Kalman filter:

factor analysis



move to state space model: unroll intime (HMM style)

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



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

diagonal noise

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

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

"decoder"

MLE → use EM

$\Rightarrow p(z|x)$ is intractable

\Rightarrow approximate with variational approach

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

$$z|x \sim N(\underbrace{\mu_\phi(x)}_{\text{output of NN}}, \sigma_\phi^2(x))$$

output of NN "encoder"

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

$$= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\phi}(x|z)] - Kullback-Leibler divergence$$

$$\text{allow re-parameterization trick}$$

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

allows "re-parameterization trick" $\sim N(0,1)$

$$z|x \rightarrow \mu_\phi(x) + \sigma_\phi^2(x) \cdot \epsilon$$

$p(z)$
 $N(0, I_K)$

- 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) => 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