

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

Gaussian networks

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

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

$$\exp\left(-\frac{1}{2} (x^T \mu x - x \mu^T + \mu \mu^T)\right)$$

sufficient statistics

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

canonical parameters

$$\Sigma \triangleq \Sigma^{-1}$$

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

$$\underbrace{\frac{1}{2} \text{tr}(\Sigma^{-1} x x^T)}_{\Delta n} + \underbrace{\frac{1}{2} \text{tr}(\Sigma^{-1} \mu \mu^T)}_{\Delta \mu} - \frac{1}{2} \mu \mu^T$$

$$\text{canonical parameters } \hat{n} \left(\begin{pmatrix} \mu \\ \Sigma \end{pmatrix} \right) = \begin{pmatrix} n \\ \Sigma \end{pmatrix} = \begin{pmatrix} \Sigma^{-1} \mu \\ \Sigma^{-1} \end{pmatrix}$$

$$p(x; n, \Sigma) = \exp(n^T x + \underbrace{\frac{1}{2} n^T \Sigma^{-1} n}_{\text{constant}} - \underbrace{\frac{1}{2} \log(2\pi) - \frac{1}{2} \log|\Sigma|}_{\text{constant}})$$

$$\Omega = \{ (n, \Sigma) : n \in \mathbb{R}^P, \Sigma \succ 0, \Sigma = \Sigma^T, \Sigma \in \mathbb{R}^{P \times P} \} \quad A(n, \Sigma)$$

useful exercise: check $\mathbb{E}_n A(n, \Sigma) = \mathbb{E}[x] = \mu$

$$\mathbb{E}_\Sigma A(n, \Sigma) = \mathbb{E}\left[\frac{x x^T}{2}\right]$$

UGM viewpoint

$$p(x; \mu, \Lambda) = \exp\left(-\frac{1}{2} \sum_{ij} \Lambda_{ij} x_i x_j + \sum_i n_i x_i - A(n, \Lambda)\right)$$

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

zeros in precision matrix \Rightarrow cond. indep property

$$\text{"Gaussian network"} \quad p(x) = \prod_{i,j} \prod_{i,j} \psi_{ij}(x_i, x_j) \prod_i \psi_i(x_i)$$

quick Schur-complement digression

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

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

$$M \triangleq \Sigma / \Sigma_{11} \triangleq \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \quad \text{"Schur complement of } \Sigma \text{ w.r.t. to } \Sigma_{11}$$

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

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

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

$$p(x_1, x_2) = \frac{1}{\sqrt{(2\pi)^2 |\Sigma_{11}|}} \exp\left(-\frac{1}{2} (x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1)\right) \quad \left. \begin{array}{l} p(x_1) \\ p(x_2 | x_1) \end{array} \right\}$$

$$\frac{1}{\sqrt{(2\pi)^2 |\Sigma_{22}|}} \exp\left(-\frac{1}{2} (x_2 - \mu_2 - b(x_1))^T (\Sigma / \Sigma_{11})^{-1} (x_2 - \mu_2 - b(x_1))\right)$$

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

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

$$\begin{aligned} \mu_1^M &= \mu_1 \\ \Sigma_1^M &= \Sigma_{11} \end{aligned} \quad \left. \begin{array}{l} \text{super simple?} \\ \text{param.} \end{array} \right\} \text{for marginal on } x_1$$

$$\begin{aligned} \mu_{2|1} &= \mu_2 + b(x_1) \\ \Sigma_{2|1} &= \Sigma / \Sigma_{11} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \end{aligned} \quad \left. \begin{array}{l} \text{param.} \\ \text{fr cond. } x_2 | x_1 \end{array} \right\}$$

in can. param.

$$\Lambda_{2|1}^{\text{can}} = -\Lambda_{22} \quad (\text{simple})$$

$$\mu_{2|1}^{\text{can}} = \mu_2 - \Lambda_{21} x_1$$

$$\mu_1^M = \mu_1 - \Lambda_{12} \Lambda_{22}^{-1} \mu_2 \quad (\text{more complicated})$$

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

for example: block $\underbrace{\Sigma_{ij}}_{I} \}_{rest}$

$$\Lambda = \begin{pmatrix} \Lambda_{RR} & \Lambda_{RI} \\ \Lambda_{IR} & \Lambda_{II} \end{pmatrix}$$

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

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

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

w/

$$\text{get } \Sigma_{IR, \text{est}} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Rightarrow \boxed{\underline{\underline{X_i \perp\!\!\!\perp X_j}} \quad \underline{\underline{X_{\text{rest}}}}}$$

(also true by Markov property of OGM)

Factor analysis:

latent variable model



$$z \in \mathbb{R}^k$$

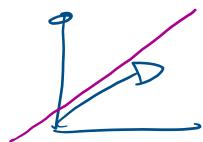
$$x \in \mathbb{R}^d$$

Learn a "latent representation" in \mathbb{R}^k
or
dimensionality reduction $k \ll d$

15h28

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



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

$$W W^T \stackrel{?}{=} I_d$$

$$P_W \stackrel{\Delta}{=} W W^T \quad P_W^2 = W W^T W W^T = P_W$$

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

$$P_W x = W(W^T x)$$

$$= \begin{pmatrix} | & \cdots & | \\ w_1 & \cdots & w_k \\ | & \cdots & | \end{pmatrix} \begin{pmatrix} \langle w_1, x \rangle \\ \langle w_2, x \rangle \\ \vdots \\ \langle w_k, x \rangle \end{pmatrix}$$

$$= \sum_k w_k \langle w_k, x \rangle = W z$$

$$z = W^T x$$

↳ lower dimensional representation

PCA

$$\min_{W \in \mathbb{R}^{d \times k}} \sum_i \|x_i - \underbrace{W W^T x_i}_{z_i}\|_2^2$$

latent auto-encoder

$$W^T W = I_k \quad \text{col}(W) \stackrel{\Delta}{=} \text{principal subspace}$$

$$X = \begin{pmatrix} | & & | \\ -x_1^T & \cdots & x_n^T \\ | & & | \end{pmatrix}$$

$$\|X^T - W W^T X^T\|_F^2$$

$$= \|(I_d - W W^T) X^T\|_F^2$$

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

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

min rec. error \Leftrightarrow

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

"analysis view"

max

long vector

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

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

(mean=0)

"analysis view
of PCA"

max
sum of empirical
covariances of
new representations

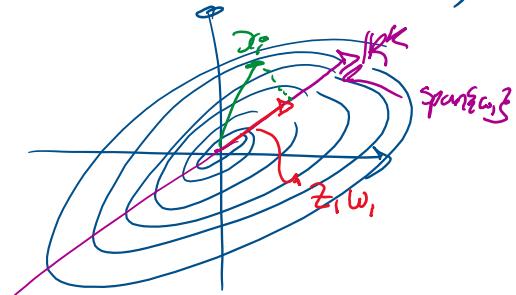
long
vector
 $\mathbf{r}(\mathbf{z}_{i,k})_{i,j}$

④ W is not unique, only col(W)

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

$$\tilde{W}\tilde{W}^T = WRR^TW^T$$

(compute sol'n of PCA \Rightarrow top eigenvectors of X^TX)



Factor analysis \rightarrow simplest generative model

$$\mathbf{z} \sim N(0, I_K)$$

$$\mathbf{x} = W\mathbf{z} + \mu + \epsilon$$

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

diagonal matrix

$$X|Z \sim N(Wz + \mu, D)$$

$p(z)$ is Gaussian

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

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

$$= \underbrace{\text{cov}(Wz, Wz)}_{=WW^T} + \underbrace{\text{cov}(\epsilon, \epsilon)}_{=D}$$

$$= WW^T + D$$

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

diagonal $\rightarrow d$ degrees of freedom

low rank covariance structure $\otimes \otimes$

estimate W, D, μ by MLE

\rightsquigarrow do EM (latent variable model)

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

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

Probabilistic PCA: special case of factor analysis

where suppose $D = \sigma^2 I$

+

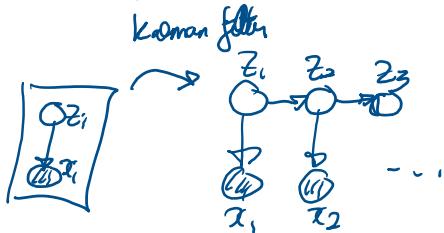
where suppose $\Sigma = \sigma^2 I$

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

this suggests that PCA is limit of PPCA as $\sigma \rightarrow 0$ $= W^T$ if $W^T W = I_K$

Kalman filter

factor analysis



move to state space model; unroll time (HMM style)

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

→ doing "sum-product" alg. in HMM $p(z_t | x_{1:t})$

get "Kalman filter" alg.

Variational auto-encoder

generalization of factor analysis



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

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

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

"decoder"

MLE → use EM

↳ $p(z|x)$ is intractable

→ 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 NN "encoder"

variational
decoder

$$\text{In EM, } \log p(x) \geq E_q [\log p(x|z)] + H(q)$$

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

allows "reparameterization trick" $\varepsilon \sim N(0, I_K)$

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

- 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