

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

Gaussian networks

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

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

sufficient statistics $T(x) = \begin{pmatrix} x \\ \frac{1}{2} xx^T \end{pmatrix}$

canonical parameter \rightarrow

$$\underbrace{\langle \Sigma^{-1} \mu, x \rangle}_{\eta} + \underbrace{\langle \Sigma^{-1}, \frac{1}{2} xx^T \rangle}_{\Lambda}$$

$\Lambda \preceq \Sigma^{-1}$
precision matrix

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

canonical parameter $\tilde{\eta}(\Theta) = \begin{pmatrix} \eta \\ \Lambda \end{pmatrix} = \begin{pmatrix} \Sigma^{-1} \mu \\ \Sigma^{-1} \end{pmatrix}$

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

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

useful exercise: $\nabla_{\eta} A(\eta, \Lambda) = \mathbb{E}[x] = \mu = \Lambda^{-1} \eta$

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

UGM viewpoint:

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

$$p \in \mathcal{G}(\mathcal{G}) \text{ where } \mathcal{G} \triangleq \{ \sum_{i,j} \Lambda_{ij} x_i x_j \text{ s.t. } \Lambda_{ij} \neq 0 \}$$

Zeros in precision matrix \Rightarrow cond. indep. properties (from UGM perspective)

"Gaussian network"

$$p(x) = \prod_{i,j \in E} \Psi_{ij}(x_i, x_j)$$

quick Schur-complement digression:

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

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

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

* use this derive "Woodbury-sherman-Morrison inversion formula"

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

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

$$p(x_1, x_2) = \frac{1}{\sqrt{(2\pi)^{d_1} |\Sigma_{11}|}} \exp(-\frac{1}{2}(x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1)) \cdot \left. \right\} p(x_1)$$

$$\frac{1}{\sqrt{(2\pi)^{d_2} |\Sigma / \Sigma_{11}|}} \exp(-\frac{1}{2}(x_2 - \mu_2 - b(x_1))^T (\Sigma / \Sigma_{11})^{-1} (x_2 - \mu_2 - b(x_1))) \left. \right\} p(x_2 | x_1)$$

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

mean parameterization of marginal & conditionals

$$\left. \begin{aligned} \mu_1^M &= \mu_1 \\ \Sigma_{11}^M &= \Sigma_{11} \end{aligned} \right\} \text{super simple!}$$

} param. for marginal on x_1

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

$$\Sigma_{2|1}^{cond} = \Sigma / \Sigma_{11} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$$

} for conditional $x_2 | x_1$

in canonical param.

$$\eta_{2|1}^{cond} = \eta_{22} \quad \text{(simple)}$$

$$\eta_{21}^{cond} = \eta_{21} - \eta_{22}^{-1} \eta_{12} \eta_{11}^{-1} \eta_{11}$$

$$\eta_{11}^m = \eta_{11} - \eta_{12} \eta_{22}^{-1} \eta_{21}$$

(more complicated)

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

x_2, x_1

block $\sum_{i,j} \dots | \text{rest}$

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

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

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

15h20

(also true by Markov property of UGM)

Factor analysis

latent variable model

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

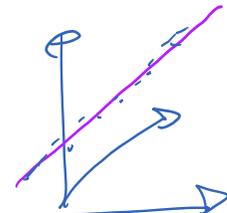
$$z \in \mathbb{R}^d$$

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

PCA for dimensionality reduction

synthesis view: find k orthonormal vectors in \mathbb{R}^d w_1, \dots, w_k

s.t. project x on span $\{w_1, \dots, w_k\}$ is a good approx of x



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

$$P_W \triangleq W W^T$$

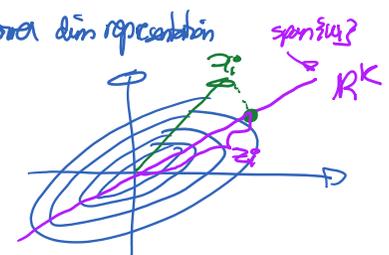
$$P_W^2 = W W^T W W^T = P_W$$

orthogonal projection matrix on span $\{w_1, \dots, w_k\} = \mathcal{C}(W)$

$$P_W x = W W^T x = \begin{pmatrix} w_1 & \dots & w_k \end{pmatrix} \begin{pmatrix} \langle w_1, x \rangle \\ \vdots \\ \langle w_k, x \rangle \end{pmatrix} = \sum_k w_k \langle w_k, x \rangle = W z$$

get lower dim representation

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



$$\text{PCA} \quad \min_{z_1, \dots, z_k} \sum_i \|x_i - W W^T z_i\|^2$$

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

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

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

$n \times d$

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

empirical covariance of x when $\sum x_i = 0$ mean = 0

$$\begin{aligned} & \|X^T - W W^T X^T\|_F^2 \\ &= \|(I - P_W) X^T\|_F^2 \\ &= \text{tr}(X (I - P_W)^T (I - P_W) X^T) \\ &= \text{tr}(X (I - P_W) X^T) \\ &= \text{tr}(X^T X (I - P_W)) \end{aligned}$$

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 in new representation

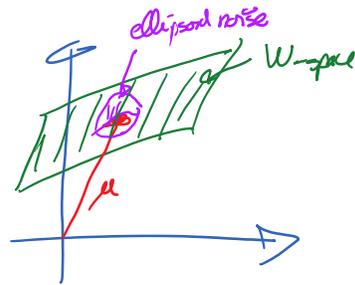
(computation of PCA \rightarrow top k e-vectors of $X^T X$)

factor analysis \rightarrow simplest generative model

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

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

Noise $\epsilon \perp z, \epsilon \sim N(0, D)$
 D \uparrow
diag. diagonal matrix



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

$p(x)$ is Gaussian; $E[x] = E[E[x|z]]$
 $E[Wz + \mu] = 0 + \mu = \mu$

$$\begin{aligned} \text{cov}(X, X) &= \text{cov}(Wz + \mu + \epsilon, Wz + \mu + \epsilon) \\ &= \text{cov}(Wz, Wz) + \text{cov}(\epsilon, \epsilon) \end{aligned}$$

\uparrow indep.

$$\overbrace{W \text{cov}(z) W^T}^{\text{Irk}} \Downarrow \\ = WW^T + D$$

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

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

estimate W, μ, D by MLE
 \rightarrow do EM (latent variable model)

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

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

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

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

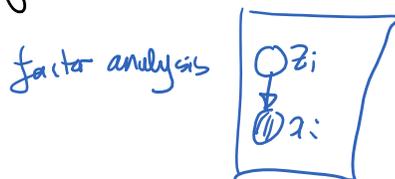
show PCA is limit as $\sigma \rightarrow 0$ of PPCA

(side note: LDA model for text is basically $G \sim \text{Dir}(\alpha)$)

$$x_i G \sim \text{Mult}(WG, n)$$

"discrete version of PPCA"

Kalman filter:



state space model: unroll in time (HMM style)



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

\rightarrow doing "sumproduct" alg. in HMM get "Kalman filtering 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 ↑ "decoder"

MLE → use EM
 ↳ $p(z|x)$ is intractable ⇒ use variational approach

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

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

in EM: $\log p(x) \geq \mathbb{E}_q[\log p(z, x)] + H(q)$
 $= \mathbb{E}_{q_\phi(z|x)}[\log p_w(z|x)] - \text{KL}(q_\phi(z|x) \parallel p(z))$

allows "reparameterization trick"

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

$\epsilon \sim N(0, 1)$

- VAE innovations:
 - share parameters ϕ 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