

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

Gaussian networks

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

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

$(xx^T - \mu x^T - x \mu^T + \mu \mu^T)$

put in exponential family

sufficient statistics

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

canonical parameter

$$\begin{aligned} &\langle \underbrace{\Sigma^{-1}}_{\Lambda \triangleq \Sigma^{-1}}, -\frac{xx^T}{2} \rangle + \langle \underbrace{\Sigma^{-1} \mu}_{\eta}, x \rangle - \frac{1}{2} \mu^T \Sigma^{-1} \mu \\ &\Lambda \triangleq \Sigma^{-1} \quad \eta \quad \mu = \Sigma \eta = \Lambda^{-1} \eta \end{aligned}$$

precision matrix

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

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

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

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

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

UGM viewpoint:

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

$p \in \mathcal{L}(G)$ where $E \triangleq \{z_{i,j}\}$ c.f. $\Lambda_{ij} \neq 0$

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 =

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

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

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

"Schur complement of Σ "
w.r. to Σ_{11}

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

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

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

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

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

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

mean parameterization
of marginals &
conditionals

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

} param. of
marginal on x_1

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

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

} for conditional
 $x_2 | x_1$

in canonical
param.

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

$$\eta_{2|1}^{\text{cond.}} = \eta_2 - \Lambda_{21} x_1$$

$$\eta_1^m = \eta_1 - \Lambda_{12} \Lambda_{22}^{-1} \eta_2$$

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

(more complicated)

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

$$\text{cov}(X_i | X_{\text{rest}}) = \Sigma_{i | \text{rest}} = \Lambda_{ii}^{-1}$$

$$= (\begin{matrix} \Lambda_{ii} & \Lambda_{ij} \\ \Lambda_{ji} & \Lambda_{jj} \end{matrix})^{-1}$$

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

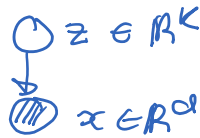
$$\Rightarrow X_i \perp\!\!\!\perp X_j | X_{\text{rest}}$$

(also true by Markov property of UGM)

15h38

Factor analysis:

latent variable model

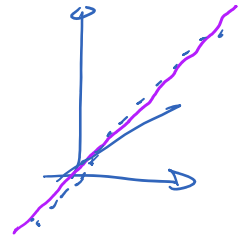


learn "latent representation"

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



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

$W W^T \neq I_d$

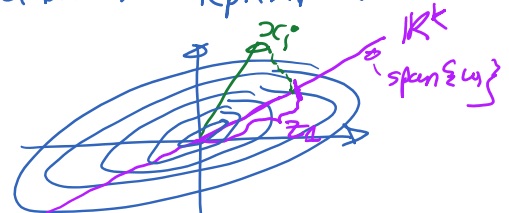
$$P_W \triangleq W W^T \quad P_W^2 = W W^T W W^T = P_W$$

\hookrightarrow orthogonal projection matrix on $\text{span}\{w_1, \dots, w_k\} = \text{col}(W)$

$$P_W x = W W^T x = \begin{pmatrix} | & \dots & | \\ w_1 & \dots & w_k \\ | & \dots & | \end{pmatrix} \begin{pmatrix} \langle w_1, x \rangle \\ \vdots \\ \langle w_k, x \rangle \end{pmatrix} = \sum_k w_k \underbrace{\langle w_k, x \rangle}_{(z)_k} = W z$$

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

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



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

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

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

e.g. $W \sim W U$ where U is orthogonal

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \|x_i - Ww_i\|_F^2 \\ & = \frac{1}{n} \sum_{i=1}^n \|x_i^T - Ww_i^T\|_F^2 \\ & = \frac{1}{n} \sum_{i=1}^n \|x_i^T - (I - P_W)x_i^T\|_F^2 \\ & = \text{tr} \left(X (I - P_W)^T (I - P_W) X^T \right) \\ & = \text{tr} \left(X (I - P_W) X^T \right) = \text{tr} (X^T X (I - P_W)) \end{aligned}$$

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

empirical covariance of x when $\mathbb{E}x_i = 0$ (mean=0)

min rec. error \Leftrightarrow maximize $\text{tr}(X^T X W W^T) = \sum_{k=1}^k w_k^T X^T X w_k$
 "analysis view of PCA" max sum of empirical variances of new representation

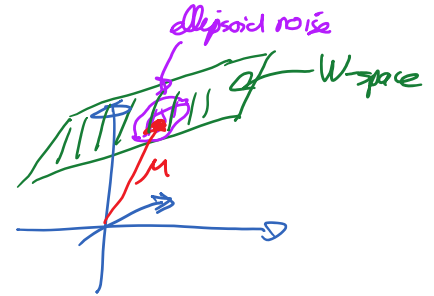
(computation of PCA \rightarrow top k -vectors of $X^T X$)

Factor analysis \rightarrow simplest generative model

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

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

$\epsilon \perp z, z \sim N(0, D)$
 D is $d \times d$ diagonal matrix



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

$p(x)$ is Gaussian; $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|Z]] = 0 + \mu = \mu$

$$\begin{aligned} \text{cov}(X, X) &= \text{cov}(Wz + \mu + \epsilon, Wz + \mu + \epsilon) \\ &= \text{cov}(Wz, Wz) + \text{cov}(\epsilon, \epsilon) \\ &= W \text{cov}(z, z) W^T + D \\ &= WW^T + D \end{aligned}$$

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

WW^T is low rank covariance piece
 D is diagonal \rightarrow d degrees of freedom

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

get $p(z|x) \rightarrow$ Gaussian with mean
 $E[z|x] = W^T(WW^T + D)^{-1}(x - \mu_0)$

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

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

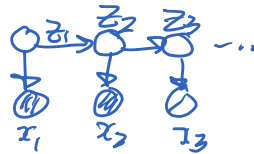
$$= W^T \quad \text{if } W^T W = I_K$$

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

(side: LDA model for text is basically $\Theta \sim \text{Dir}(\alpha)$
 $x | \Theta \sim \text{Mult}(W\Theta, n)$)
 "discrete version of PPCA"

Kalman filter

factor analysis



markov state space model: unroll in time (HMM style)

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

\rightarrow doing "sum-product" alg in HMM get "Kalman filter alg."

Variational auto-encoder

generalization of factor analysis



$z \sim N(0, I_K)$ (diagonal noise)
 $x | z \sim N(\mu_W(z), \sigma_W^2(z))$ ("encoder")
 where $\mu_W(z)$ is output of NN

MLE \rightarrow use EM
 $\hookrightarrow p(z|x)$ is intractable \Rightarrow use variational approach

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

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

in EM: $\log p(x) \approx E_q[\log p(x|z)] + H(q)$

$= E_q[\dots] - k \log \int p(x|z) dz$

$$\begin{aligned}
 \ln \mathcal{L} &= \mathbb{E}_{q(z|x)} [\log p(x|z)] - \mathbb{E}_{q(z|x)} [\log p(z)] \\
 &= \mathbb{E}_{q(z|x)} [\log p_w(x|z)] - \text{KL}(q_\phi(z|x) \parallel p(z))
 \end{aligned}$$

variance

allows "reparameterization trick"

$$z|x \rightarrow \mu_\phi(x) + \sigma_\phi(x) \cdot \xi \quad \xi \sim \mathcal{N}(0,1)$$

○ 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 $\mathcal{N}(0,1)$ 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