

Lecture 24 - scribbles

Friday, December 1, 2017
12:56

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

Gaussian networks:

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

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

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

precision matrix

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

$$= \langle \Sigma^{-1}, x x^T \rangle \quad \text{linear form} \quad -2 \langle \Sigma^{-1} \mu, x \rangle \quad \triangleq \eta$$

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

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

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

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

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

$A(m, \Lambda)$

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

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

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

$$p \in \mathcal{J}(\mathcal{G}) \text{ where } E \hat{=} \{ \lambda_{ij} \} \text{ s.t. } \Lambda_{ij} \neq 0 \}$$

Zeros in precision matrix \Rightarrow cond. indep. properties \oplus

"Gaussian network"

quick Schur-complement digression

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

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

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

Schur-complement of Σ with respect to Σ_{11}

use this to derive

"Woodbury-Sherman-Morrison inversion formula"

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

using above,

$p(x_1, x_2) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_{11}|}} \exp(- (x_1 - \mu_1) \Sigma_{11}^{-1} (x_1 - \mu_1)) \} p(x_1)$

block of diag. pt

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

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

mean parameterization of marginal and conditionals:

$\mu_1^m = \mu_1$
 $\Sigma_1^m = \Sigma_{11}$

} marginal on x_2

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

} conditional of x_2 given x_1

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

in canonical parameterization

$\Lambda_{2|2}^{cond} = \Lambda_{22}$ (simple)

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

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

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

block $\begin{matrix} x_2 \\ \sim \\ \{i,j\} \\ \uparrow \\ \uparrow \end{matrix}$ | $\begin{matrix} x_2 \\ \sim \\ \text{rest} \end{matrix}$

$$\text{cov}(X_I | X_{\text{rest}}) = \sum_{I \in I | \text{rest}} = \Lambda_{II}^{-1}$$

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

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

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

(also true by Markov property of UGM)

Factor analysis:

↑ latent variable model $z \in \mathbb{R}^k$ learning a "latent representation"
 ↓ $x \in \mathbb{R}^d$ or dimensionality reduction

PCA for dimensionality reduction:

Synthesis view: find a k -orthormal vectors in \mathbb{R}^d w_1, \dots, w_k

s.t. the projection of x on $\text{span}\{w_1, \dots, w_k\}$ is a good approximation of x

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

what about $W W^T$?

$D \triangleq W W^T$ \hookrightarrow projection matrix on $\text{span}\{w_1, \dots, w_k\} = \text{col}(W)$

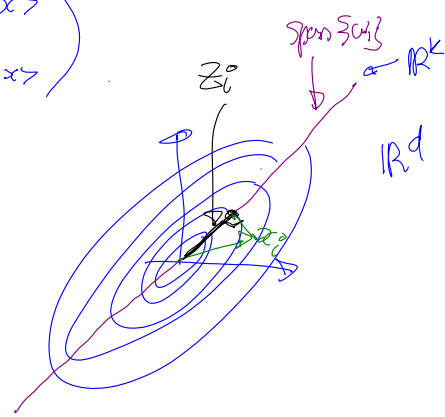
$P_W \triangleq WW^T$ \hookrightarrow projection matrix on span $\{w_1, \dots, w_k\} = \text{col}(W)$

$P_W^2 = \underbrace{WW^T}_{I_k} WW^T = P_W$

$$P_W x = WW^T x = (w_1 \dots w_k) \begin{pmatrix} \langle w_1, x \rangle \\ \vdots \\ \langle w_k, x \rangle \end{pmatrix}$$

$$= \sum_k w_k \underbrace{\langle w_k, x \rangle}_{(z)_k}$$

get lower dim. representation $z = W^T x \in \mathbb{R}^k$ $P_W x = Wz$



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

$W^T W = I_k$

$\text{col}(W) \triangleq$ "principal subspace"

"Synthesis view"

$W = WR$ where R is k-rotations $R^T R = R R^T = I_k$

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

then $\tilde{W} \tilde{W}^T = W \underbrace{R^T R}_{I_k} W^T = WW^T$

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

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

$$= \|(I - P_W) X^T\|_F^2$$

$$= \text{tr}(X(I - P_W)^T (I - P_W) X^T)$$

$$= \text{tr}(X(I - P_W) X^T)$$

$$= \text{tr}(X^T X (I - P_W))$$

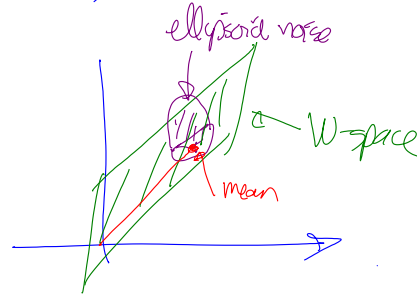
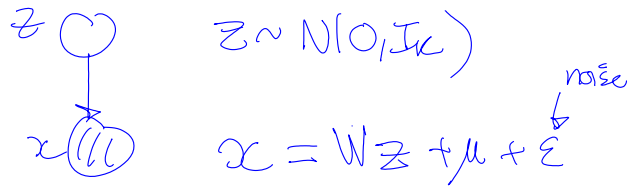
min rec. error \Leftrightarrow maximize $\text{tr}(X^T X WW^T) = \sum_{i=1}^k \lambda_i$

"analysis view of PCA"

min rec. error \Leftrightarrow maximize $\text{tr}(X^T X W W^T) = \sum_k w_k^T X^T X w_k$ view of PCA
 Maximizing variance in new representation

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

Factor analysis \rightarrow generative model (latent variable)



$\epsilon \perp z, \epsilon \sim N(0, D)$
 D \downarrow
 a diag matrix

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

$p(z)$ is Gaussian, $E[X] = E[E[X|Z]]$
 $E[Wz + \mu] = 0 + \mu = \mu$
 $\text{Cov}(X, X) = \text{Cov}(Wz + \mu + \epsilon, Wz + \mu + \epsilon)$
 $\quad \quad \quad \uparrow$
 $\quad \quad \quad \text{indep.}$
 $= \text{Cov}(Wz, Wz) + \text{Cov}(\epsilon, \epsilon)$
 $\quad \quad \quad \underbrace{\hspace{2cm}} \quad \quad \quad \underbrace{\hspace{2cm}}$
 $\quad \quad \quad W^T \text{Cov}(z) W \quad \quad \quad D$
 $= WW^T + D$

new joint mult. for $x, z \sim N(\mu, WW^T + D)$

equivalent model for x : $x \sim N(\mu, \underbrace{WW^T}_{d \times k} + \underbrace{D}_{\text{diagonal}})$
 $\Rightarrow d$ degrees of freedom

get $p(z|x)$:

$$\begin{pmatrix} x \\ z \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_x \\ \mu_z \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_{zz} \end{pmatrix} \right)$$

$$\begin{aligned} \mu_x &= \mu & \Sigma_{xx} &= WW^T + D \\ \mu_z &= 0 & \Sigma_{zz} &= I_k \end{aligned}$$

$\text{cov}(x, z) = W$ (exercise)

$$\begin{aligned} \mathbb{E}[z|x] &= \mu_z + \Sigma_{zz}^{-1} \Sigma_{zx} (x - \mu_x) \\ &= 0 + W^T (W^T W + D)^{-1} (x - \mu_x) \end{aligned}$$

probabilistic PCA: suppose $D = \sigma^2 I_d$

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

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

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

⊛ ^{for} factor analysis, estimate $\hat{W}, \hat{D}, \hat{\mu}$ using ML

→ use EM → need $p(z|x)$

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

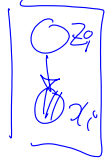
$x | \Theta \sim \text{Mult}(W\Theta, n)$)

"discrete version of PPCA"

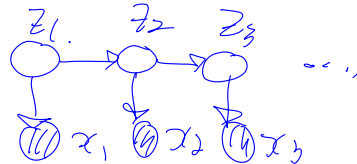
note that \hat{W} is not identifiable (can only be identified up to rotation)

Kalman filter:

factor analysis



state space model: unrolled in time:

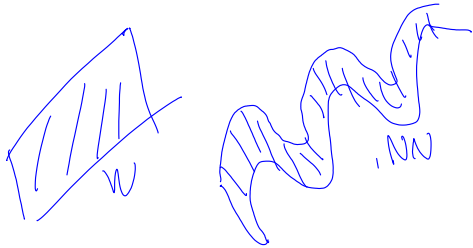


Kalman filter: $z_t | z_{t-1} \sim N(Az_{t-1}, B)$

Variational auto-encoder:

$z \sim N(0, I_K)$

generalization of factor analysis where $x|z \sim N(\mu_w(z), \sigma_w^2(z))$ "decoder"



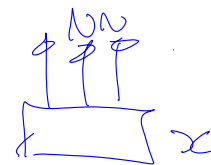
$\mu_w(z)$ ← output of NN

$p(z|x)$ is intractable \Rightarrow use variational approach

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

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

$q_\phi(z|x)$



in EM,

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

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

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