

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

mean field continuation

$$\min_{q \in Q_{MF}} KL(q \parallel p)$$

↓

$$\{q : q(x) = \prod_i q_i(x_i)\}$$

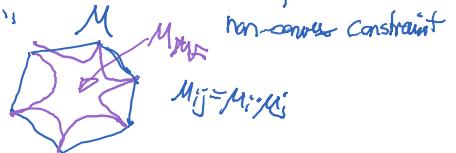
$KL(\cdot \parallel p)$ is a convex of q (as vector)

but Q_{MF} is a non-convex constraint set

⇒ can get stuck in
local minima.

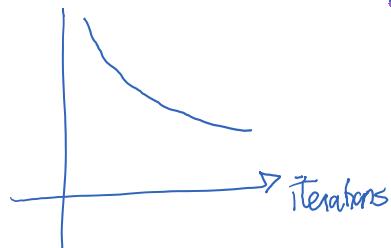
e.g., Ising model
 $M_{ij} = M_i \cdot M_j$

[see lecture 27 in 2017, for "manginal polytope"]



but can monitor progress

$$KL(q^{(t)} \parallel p) + \text{cost}$$



pros & cons of variational methods

vs. Sampling

⊕ optimization based
⇒ often faster to run
↳ easier to debug

⊖ noisy ⇒ harder to debug
mixing problems for chains

⊖ biased estimate

$$\mathbb{E}_{q(z)}[f(z)] \neq \mathbb{E}_p[f(z)]$$

⊕ unbiased estimate

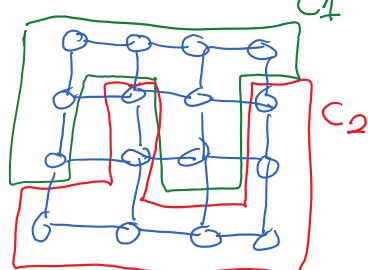
$$\mathbb{E}\left[\mathbb{E}_{q(z)}[f(z)]\right] = \mathbb{E}_p[f(z)]$$

with respect to
random sample

structured mean field:

idea $q(z) = \prod_{j=1}^k q_j(z_{C_j})$ where C_1, \dots, C_K is a partition of V

and q_j 's are tractable distributions
(for example free VGMM)



15h43

Gaussian networks

↔ ... ↔ and ↗ ... $d \times d$ ↘ ...

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

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

put in exponential family

sufficient statistics

$$T(x) = \begin{pmatrix} x \\ -\frac{x^T}{2} \end{pmatrix} \quad \begin{matrix} \text{canonical} \\ \text{parameter} \end{matrix}$$

$$\begin{aligned} & \left\langle \Sigma^{-1}, -\frac{x^T}{2} \right\rangle + \left\langle \Sigma^{-1}\mu, x \right\rangle - \frac{1}{2} \mu^T \Sigma^{-1} \mu \\ & \stackrel{\oplus n}{=} \mu^T \Sigma^{-1} \mu \\ & \mu = \Sigma n = \Sigma^{-1} n \end{aligned}$$

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

$$\mathcal{Q} = \{ (n, \Lambda) : n \in \mathbb{R}^d, \Lambda \geq 0, \Lambda = \Lambda^T, \Lambda \in \mathbb{R}^{d \times d} \}$$

$$\text{useful exercise: } \mathbb{E}[n] = \mu$$

$$\mathbb{E}[\Lambda] = \mathbb{E}\left[-\frac{x^T}{2}\right]$$

UGM viewpoint:

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

$$p \in \mathfrak{f}(G) \text{ where } E \triangleq \{ \xi_{i,j} \} \text{ s.t. } \Lambda_{i,j} \neq 0 \}$$

$$\text{"Gaussian network"} \quad p(x) = \prod_{\substack{i, j \in E \\ \xi_{i,j} \neq 0}} \psi_{i,j}(x_i, x_j) \prod_{i \in V} \psi_i(x_i)$$

quick Schur-complement digression

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

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

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

"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 : $|\Sigma| \approx |\Sigma_{11}| \cdot |\Sigma / \Sigma_{11}| = |\Sigma_{22}| |\Sigma / \Sigma_{22}|$

$$p(x_1, x_2) = \frac{1}{\sqrt{(2\pi)^{d_1} |\Sigma_{11}|}} \exp\left(-\frac{1}{2}(x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1)\right) \cdot \underbrace{\frac{1}{\sqrt{(2\pi)^{d_2} |\Sigma / \Sigma_{11}|}} \exp\left(-\frac{1}{2}(x_2 - \mu_2 - b(x_1))^T \left(\begin{matrix} \Sigma \\ \Sigma_{11} \end{matrix}\right)^{-1} (x_2 - \mu_2 - b(x_1))\right)}_{p(x_2 | x_1)}$$

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

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

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

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

in canonical param.

$$\Lambda_{211} = \Lambda_{22} \quad (\text{simple})$$

$$\Lambda_{211}^{\text{cond.}} = \Lambda_2 - \Lambda_{21} \Lambda_1$$

$$\Lambda_1^m = \Lambda_1 - \Lambda_{12} \Lambda_{22}^{-1} \Lambda_2$$

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

(more complicated)

for example : block $\begin{smallmatrix} I \\ \Sigma_{11} \end{smallmatrix}$ / rest

$$\Lambda = \begin{pmatrix} \Lambda_{RR} & \Lambda_{RZ} \\ \Lambda_{ZR} & \Lambda_{ZZ} \end{pmatrix}$$

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

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

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

(also true by Markov of UGM)

Factor analysis

latent variable model $\mathbf{y} \in \mathbb{R}^k$ learn "latent representation"

latent variable model

$$\begin{array}{c} \textcircled{z} \in \mathbb{R}^K \\ \downarrow \\ \textcircled{x} \in \mathbb{R}^d \end{array}$$

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. projecting x on $\text{span}\{w_1, \dots, w_K\}$
 is a good approx. of x



$$W = [w_1 \ w_2 \ \dots \ w_K] \quad W^T W = I_K \quad (\text{by orthogonality})$$

$$WW^T \neq I_d$$

$$P_W \triangleq WW^T \quad P_W^2 = W \underbrace{W^T}_{\perp K} WW^T = P_W$$

\hookrightarrow orthogonal projection
on $\text{span}\{w_1, \dots, w_K\}$

$$P_W x = WW^T x$$

$$= \left(\begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_K \end{array} \right) \left(\begin{array}{c} \langle w_1, x \rangle \\ \vdots \\ \langle w_K, x \rangle \end{array} \right)$$

$$= \sum_{i=1}^K w_i \underbrace{\langle w_i, x \rangle}_{(z)_i} = Wz$$

$$z = W^T x$$

PCA

$$\min_{\substack{W \in \mathbb{R}^{d \times K} \\ W^T W = I_K}} \sum_i \|x_i - WW^T x_i\|_2^2$$

col(W) \triangleq principal subspace

lower dimensional representation

$$X_{n \times d} = \begin{pmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{pmatrix}$$

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

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

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

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

Min rec. error \Leftrightarrow maximize $\text{tr}(X^T X W W^T) = \sum_K w_i^2 \sum_j x_i^j x_i^j$

"analysis view
of PCA"

max sum of empirical
variances of
new representation

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

Factor analysis \rightarrow simplest generative model

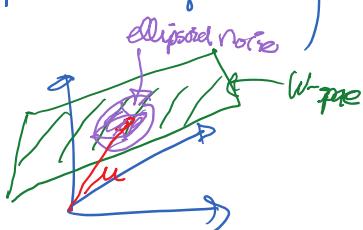
$$z \sim N(0, I_K)$$

noise

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

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

diagonal matrix



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