

- today:
- exp. family
  - estimation DGM/UGM
  - sampling

moment matching different than MLE in exp. family

gamma dist. has  $T(x) = \begin{bmatrix} \log x \\ x \end{bmatrix}$   
 so moment matching with  $\tilde{T}(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$  will yield different estimate than MLE

example 2: 1d Gaussian

$X \sim N(\mu, \sigma^2)$   $X = \mathbb{R}$   $\Theta = (\mu, \sigma^2)$  "moment parametrization"

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

$$= \exp\left(-\frac{x^2}{2} \underbrace{\left[\frac{1}{\sigma^2}\right]}_{\eta_2} + x \underbrace{\left[\frac{\mu}{\sigma^2}\right]}_{\eta_1} - \underbrace{\left[\frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)\right]}_{A(\eta)}\right)$$

$$T(x) = \begin{bmatrix} x \\ -\frac{x^2}{2} \end{bmatrix} \quad \eta(\theta) = \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}$$

$\eta_2 = \frac{1}{\sigma^2} = \text{precision} > 0$

$\eta_1 = \eta_2 \cdot \mu$

$h(\eta) = 1$  (but some people use  $h(\eta) = \frac{1}{\sqrt{2\pi}}$  for Gaussian)

$\Omega = \{(\eta_1, \eta_2) : \eta_2 > 0, \eta_1 \in \mathbb{R}\}$

[we'll see later: multivariate Gaussian  $T(x) = \begin{bmatrix} x \\ -\frac{xx^T}{2} \end{bmatrix}$   $\begin{cases} \Omega = \Sigma^{-1} \\ \eta = \Omega \mu = \Sigma^{-1} \mu \end{cases}$ ]

example 3: discrete UGM

let  $p \in \mathcal{P}(G)$ ,  $G$  is undirected with  $\psi_c(x_c) > 0 \forall c, x_c$

$$p(x) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(x_c) = \exp\left(\sum_c \log \psi_c(x_c) - \log Z\right)$$

$$= \exp\left(\sum_{c \in \mathcal{C}} \sum_{y_c \in X_c} \underbrace{\mathbb{1}\{x_c = y_c\}}_{T_{c,y_c}(x)} \underbrace{\log \psi_c(y_c)}_{\eta_{c,y_c}} - \log Z\right)$$

$$T(x) = \begin{pmatrix} \vdots \\ \mathbb{1}\{x_c = y_c\} \\ \vdots \end{pmatrix} \leftarrow \begin{matrix} y_c \in X_c \\ c \in \mathcal{C} \end{matrix} \quad X_c = \sum_i (y_i)_{i \in c} : y_i \in X_i$$

$\eta(\theta) = \left( \eta_{c,y_c} \right)_{c \in \mathcal{C}, y_c \in X_c}$  [not a minimal representation]

$$\eta(\theta) = \left( \begin{array}{c} \vdots \\ \log \mu_c(y_c) \\ \vdots \end{array} \right)_{\substack{y_c \in X_c \\ c \in \mathcal{C}}} \quad [\text{not a minimal representation}]$$

notes: a) Mult(x) is a special case where <sup>have</sup> Complete graph (1 big clique)

b) feature perspective: instead of using all possible indicators  $\mathbb{1}_{\{y_c = x_c\}}$  you could use a subset for a task

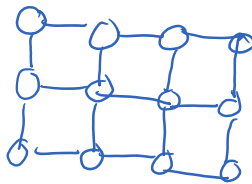
for example: suppose  $x$  is a sentence  
 $x_i$  is "word"

feature on  $x_i \& x_{i+1}$  e.g.  $\mathbb{1}_{\{x_i \text{ is a verb, } x_{i+1} \text{ is "run"}\}}$

→ much smaller set of parameters

c) binary Ising model

$$x_i \in \{0, 1\} \quad |C| \leq 2$$



suppose use nodes & pairs (edges) as cliques

⇒ dimension of  $T(x)$  is  $2|V| + 4|E|$  → "overparameterized" exp. family

$$\sum_{y_c} T_{c,y_c}(x) = 1 \quad \text{for any } c$$

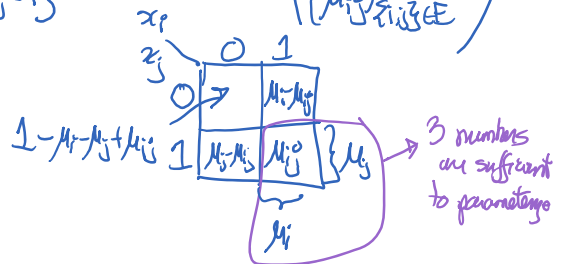
→ not a min. exp. family

\* a minimal representation

$$\text{is } T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{i,j \in E} \end{pmatrix}_{\substack{\mu_i \\ \mu_{ij}}} \rightarrow \text{dim: } |V| + |E|$$

$$\mathbb{1}_{\{x_i=1, x_j=1\}}$$

$$E T(x) = \begin{pmatrix} (\mu_i)_{i \in V} \\ (\mu_{ij})_{i,j \in E} \end{pmatrix}$$



properties of  $A$ :

$$\bullet \nabla_m A(m) = \mathbb{E}_{p(x;m)} [T(x)] \triangleq \mu(m) \quad \text{"moment vector"} \quad (\text{for } m \in \text{int}(\Omega))$$

$$\bullet \left( \frac{\partial^2 A(m)}{\partial \eta_i \partial \eta_j} \right) \text{ is } \mathbb{E}_{p(x;m)} [ [T(x) - \mu(m)] [T(x) - \mu(m)]^T ] = \text{cov}(T(x))$$

(proof as exercise)

"cumulant generating fun."

16h02

Estimation of parameters DGM/UGM

DGM:  
(fully observed)

parametric family  $\mathcal{P}_{\Theta} = \{ p_{\theta}(x) = \prod_i p(x_i | x_{\setminus i}, \theta_i) \}$   
 $\theta = (\theta_1, \dots, \theta_{|M|})$   
 $\Theta = \Theta_1 \times \Theta_2 \times \dots \times \Theta_{|M|}$

independent parameters  
ie. no tying parameters

$\Rightarrow$  MLE decouples in  $|M|$  independent MLE problems

$$p(\text{data} | \theta) = \prod_{i=1}^n p(x^{(i)} | \theta) = \prod_{i=1}^n \prod_{j=1}^{|M|} p(x_j^{(i)} | x_{\setminus j}^{(i)}, \theta_j)$$

$$\log [ ] = \sum_{j=1}^{|M|} \underbrace{\left( \sum_{i=1}^n \log p(x_j^{(i)} | x_{\setminus j}^{(i)}, \theta_j) \right)}_{\ell(\theta_j)}$$

example: for discrete R.V.  
(multinomial conditions)

$\Rightarrow \theta_j^{\text{MLE}} = \text{proportion of observations}$

$$\hat{p}(x_j = k | x_{\setminus j} = \text{stuff}) = \frac{\#(x_j = k, x_{\setminus j} = \text{stuff})}{\#(x_{\setminus j} = \text{stuff})}$$

(fully observed DGM is relatively easy)

$\otimes$  if have latent variables (ie. unobserved variables)  
 $\Rightarrow$  use EM. (like HMM)

UGM:

example for expo family

$$p(x | m) = \exp\left(\sum_C m_C T_C(x_C) - A(m)\right)$$

$\rightarrow$  unlike in a DGM,  $\log p(x | m)$  does not separate as  $\sum_C f_C(m_C)$

gradient ascent on log-likelihood

$$\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)} | m) = \sum_C m_C^T \left( \frac{1}{n} \sum_{i=1}^n T_C(x_C^{(i)}) \right) - \frac{1}{n} A(m)$$

$$\nabla_{m_C} [ ] = \hat{\mu}_C - \mu_C(m) \xrightarrow{\text{MC}} \mathbb{E}_{p(x; m)} [T_C(x_C)]$$

to compute this, need inference

e.g. Ising model  $T_{ij}(x_i, x_j) = x_i \cdot x_j$  ( $x_i \in \{-1, 1\}$ )

$$\mathbb{E}[T_{ij}] = \mu_{ij} = p(x_i = 1, x_j = 1 | m)$$

here need approx. inference  $\xrightarrow{\text{samplng}}$  eg. Gibbs sampling

here need approx. inference  
 sampling  
 variational method  
 eg [Gibbs sampling] [mean field]

## Approximate inference - Sampling

example: NP hard to do exact inference in Ising model  $\rightarrow$  need approximation

why sampling?  $X = (X_1, \dots, X_p)$

a) simulation:  $X^{(i)} \sim P$

b) approximate marginal  $p(x_i)$

$\rightarrow$  special case of expectations

consider:  $f: \mathbb{R}^p \rightarrow \mathbb{R}$

we want to approximate  $\mu = \mathbb{E}_P[f(x)]$

special case: if  $f(x) \triangleq \mathbb{1}\{x_A = x_A\}$   $\mathbb{E}_P[f(x)] = p(x_A = x_A)$

Monte Carlo integration / estimation  $\rightarrow$  appears in physics, applied math., ML, statistics, ...

to approximate  $\mu = \mathbb{E}_P[f(x)]$

MC estimation alg:  $\cdot$  n samples  $X^{(i)} \stackrel{iid}{\sim} P$   
 $\cdot$  estimate  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(x^{(i)}) = \mathbb{E}_{\hat{P}_n}[f(x)]$

properties: 1) unbiased estimator  $\mathbb{E}[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(x^{(i)})] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$   
 expectation over  $(x^{(i)})_{i=1}^n$   
 emp. dist.  $\hat{P}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x^{(i)} = x\}$   
 this is still true if  $X^{(i)}$  are dependent

2) expected error (ls-error)  $\mathbb{E}[\|\hat{\mu} - \mu\|_2^2] = \mathbb{E}[\langle \frac{1}{n} \sum_{i=1}^n f(x^{(i)}) - \mu, \frac{1}{n} \sum_{j=1}^n f(x^{(j)}) - \mu \rangle]$   
 $\text{tr}(\text{cov}(\hat{\mu}, \hat{\mu})) = \mathbb{E}[\frac{1}{n^2} \sum_{i,j} \langle f(x^{(i)}) - \mu, f(x^{(j)}) - \mu \rangle]$

by independence  $\Rightarrow$  off-diagonal terms are zero

ie.  $\langle \mathbb{E}[f(x^{(i)})] - \mu, \mathbb{E}[f(x^{(j)})] - \mu \rangle = \langle \mu - \mu, \mu - \mu \rangle = 0$

$= \frac{1}{n^2} \sum_{i=j} \mathbb{E}[\langle f(x^{(i)}) - \mu, f(x^{(i)}) - \mu \rangle]$   
 $= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[\|f(x^{(i)}) - \mu\|_2^2] \triangleq \sigma^2$   
 $= \text{tr}(\text{cov}(f(x), f(x)))$

$= \frac{\sigma^2}{n}$

$$\boxed{\mathbb{E}[\| \hat{\mu} - \mu \|^2] = \frac{\sigma^2}{n}}$$

$$= \frac{d\sigma^2}{n}$$

$$= \frac{1}{n} \text{tr}(\text{cov}(f(x), f(x)))$$

note: not dimension in rate

(apart  $\sigma^2$   
which could depend  
implicitly)

e.g.  $f(x) = X$   
 $X_j \sim N(0, \sigma^2)$

$$\mathbb{E}[\|f(x) - \mu\|^2] = d\sigma^2$$