

LT review - Parameter est. VGMS

- (*) check calculations/derivations
- (*) use Jordan 9.3-9.5 + 20 to assist \rightarrow skim
- (*) skim Lafferty or tutorial

Sufficient statisticsfrom Jordan (2003)
Ch 9.3

- total count and clique counts
- \underline{x}_V - random vector associated with graph
- \underline{x}_C $C \subseteq V$ - random vectors — " subsets of nodes
- IID replicates of obs.
- $\underline{x}_{V,n}$ - n th replicate of subset C
- $\mathcal{D} = \{\underline{x}_{V,1}, \underline{x}_{V,2}, \dots, \underline{x}_{V,N}\}$ - completely observed
- Parametrise VGMM via clique potentials $\psi_C(\underline{x}_C)$ for $C \in \mathcal{C}$
- \mathcal{C} - set of cliques.
- via Hammersley-Clifford:

Joint prob: $p(\underline{x}_V | \underline{\theta}) = \frac{1}{Z} \prod_C \psi_C(\underline{x}_C)$

$\underline{\theta} = \{\psi_C(\underline{x}_C), C \in \mathcal{C}\}$ (param. ^{collec.})

 Z = normalisation

$$Z = \sum_{\underline{x}_V} \prod_C \psi_C(\underline{x}_C)$$

⊗: Z - obtained via summing (or integrating) over all configurations \underline{x}_V

The total counts

^{obs} - no. of times that configuration \underline{x}_V is observed in a dataset \mathcal{D}

$$m(\underline{x}_V) := \sum_{n=1}^N \delta(\underline{x}_V, \underline{x}_{V,n})$$

Marginal counts (clique counts)
for clique C

$$m(\underline{x}_C) := \sum_{\underline{x}_{V \setminus C}} m(\underline{x}_V)$$

$$N = \sum_{x_v} m(x_v)$$

- total no. of observations

(*) log-likelihood for MMS (tabular clique pot.) (A3)

(*) express log-likelihood in terms of counts (sufficient statistics for discrete models)

- introduce dummy var x_v (?)

$$p(x_v, n | \theta) = \prod_{x_v} p(x_v | \theta)^{\delta(x_v, x_v, n)}$$

- OK - indicator
tick that switches
on/off
depending on
distrib.

$$\delta(x_v, x_v, n) = \mathbb{1}\{x_v = x_v, n\}$$

(*) dummy variable x_v ranges across configurations of nodes rather than across datapoints.

- standard for multinomials (remember Bishop?)

Probability of observed data :- $p(D | \theta) = \prod_{n=1}^N p(x_v, n | \theta)$

$$= \prod_{n=1}^N \prod_{x_v} p(x_v | \theta)^{\delta(x_v, x_v, n)}$$

log-likelihood in terms of marginal counts

$$\ell(\theta; D) = \log p(D | \theta)$$

$$= \sum_{n=1}^N \sum_{x_v} \delta(x_v, x_v, n) \log p(x_v | \theta)$$

$$= \sum_{x_v} \sum_n \delta(x_v, x_v, n) \log p(x_v | \theta)$$

factor out
terms without
sum. avg.

$$= \sum_{x_v} m(x_v) \log p(x_v | \theta)$$

subs. $p(x_v | \theta)$

$$= \sum_{x_v} m(x_v) \log \left(\frac{1}{Z} \prod_c \psi_c(x_c) \right)$$

$$= \sum_{x_v} m(x_v) \sum_c \log \psi_c(x_c) - \sum_{x_v} m(x_v) \log Z$$

$$\textcircled{?} \downarrow = \sum_c \sum_{x_c} n(x_c) \log \psi_c(x_c) - N \log Z \quad \textcircled{?} \quad \textcircled{Q} \quad \textcircled{0/51}$$

(9.43)

- name MLE

- see earlier 9.2

- Assume 'true' and investigate

$n(x_c)$ - neg. counts \rightarrow suff. statistics

$N \log Z$ - appears (not in OGM)

MLE of ψ_c :

(*) $N \log Z$ - coupled, nonlinear set of eq. with implicit perm app.

$$\ell(\theta; D) = \sum_c \sum_{x_c} n(x_c) \log \psi_c(x_c) - N \log Z \quad (9.43)$$

(*) Partial deriv wrt $\psi_c(x_c)$ with clique C , config x_c held fixed.

$$(i) \frac{\partial}{\partial \psi_c(x_c)} n(x_c) \log \psi_c(x_c) = \frac{n(x_c)}{\psi_c(x_c)}$$

Note that D and \tilde{x} are just diff dummy indexing variables

$$\begin{aligned} (ii) \frac{\partial}{\partial \psi_c(x_c)} \log Z &= \frac{1}{Z} \frac{\partial}{\partial \psi_c(x_c)} \left\{ \sum_{\tilde{x}_D} \prod \psi_D(\tilde{x}_D) \right\} \\ &= \frac{1}{Z} \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) \frac{\partial}{\partial \psi_c(x_c)} \left(\prod_D \psi_D(\tilde{x}_D) \right) \\ &= \frac{1}{Z} \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) \prod_{D \neq C} \psi_D(\tilde{x}_D) \\ &= \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) \frac{1}{\psi_c(\tilde{x}_c)} \frac{1}{Z} \prod_D \psi_D(\tilde{x}_D) \\ &= \frac{1}{\psi_c(\tilde{x}_c)} \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) p(\tilde{x}) \\ &= \frac{p_c(x_c)}{\psi_c(x_c)} \end{aligned}$$

• note all $\psi_D(\tilde{x}_D)$ where $D \neq C$ are assumed constant wrt $\psi_c(x_c)$

$\textcircled{Q52}$
- not entirely clear crystal

yielding $\frac{\partial \ell}{\partial \psi_c(x_c)} = \frac{n(x_c)}{\psi_c(x_c)} - N \frac{p_c(x_c)}{\psi_c(x_c)}$

$\log \psi_c(x_c) > 0$

- set $\frac{\partial \ell}{\partial \psi_c(x_c)} = 0 \Rightarrow m(x_c) - N p(x_c)$

② define empirical distri $\tilde{p}(x) = \frac{m(x)}{N}$, $\tilde{p}(x_c) = \frac{m(x_c)}{N}$ is a marg. inde emp. (?)

$\Rightarrow \frac{m(x_c)}{N} = p(x_c) \Rightarrow \hat{p}_{m_c}(x_c) = p(x_c)$

(*) for each clique $c \in C$, the model marginals $p(x_c)$ must be equal to empirical marginals $\hat{p}_{m_c}(x_c)$

Jordan 2003 9.3.3. \rightarrow check MLE by inspection for decomposable graphs

(*) Iterative prop. fitting (IPF)

- IPF is not only a fixed point algo, but coordinate ascent algo

(*) use: $\frac{\tilde{p}(x_c)}{\psi_c(x_c)} = \frac{p(x_c)}{\psi_c(x_c)}$

$\Rightarrow \frac{\psi_c(x_c)}{\tilde{p}(x_c)} = \frac{\psi_c(x_c)}{p(x_c)}$

$\Rightarrow \psi_c(x_c) = \psi_c(x_c) \frac{\tilde{p}(x_c)}{p(x_c)}$

$\Rightarrow \psi_c^{(t+1)}(x_c) = \psi_c^{(t)}(x_c) \frac{\tilde{p}(x_c)}{p^{(t)}(x_c)} \quad (9.61)$

(IPF algorithm)

- Hold $\psi_c(x_c)$ fixed in numerator, denom of RHS.
- why? As $\psi_c(x_c)$ appears implicitly through $p(x_c)$
- introduce the maxim; - param at iter t $\psi_c^{(t)}(x_c)$
- joint prob based on estimates $p^{(t)}(x)$ at iter t

(*) cycle through all cliques $c \in C$, applying (9.61); one cycle - one iter.

(*) Properties of IPF

- summarised from Jordan (2003) ch 9.3.

1) marginal $p^{(t+1)}(x_c)$ is equal to empirical marginal $\tilde{p}(x_c)$

2) normalisation constant Z - constant across updates

which yields: $p^{(t+1)}(x_v) = p^{(t)}(x_{v \setminus c} | x_c) \hat{p}(x_c)$

- each IPF iteration retains previous cond. prob $p^{(t)}(x_{v \setminus c} | x_c)$ and replaces prev. marg prob $p^{(t)}(x_c)$ with new marginal $\hat{p}(x_c)$.

(*) IPF as w-ordinate ascent

- additional derivations in Jordan show this sets gradient of log likelihood to 0; and is coordinate ascent (in log-likelihood)

(*) KL divergence view of IPF (AS)

- A few key equations:

KL-divergence with marginals:

- Decompose $p(x_A, x_B) = p(x_B | x_A) p(x_A)$ and $q(x_A, x_B) = q(x_B | x_A) q(x_A)$

- KL div: $D(p(x_A, x_B) || q(x_A, x_B)) = D(p(x_A) || q(x_A)) + \sum_{x_A} p(x_A) D(p(x_B | x_A) || q(x_B | x_A))$

(*) (X): maximising likelihood equivalent to minimising following KL-div: -

$$D(\hat{p}(x) || p(x | \theta)) = \sum_x \hat{p}(x) \log \frac{\hat{p}(x)}{p(x | \theta)} \quad \hat{p}(x) - \text{empir. distn.}$$

(*) equivalence \rightarrow -ve log likelihood and KL div differ by $\sum_x \hat{p}(x) \log \hat{p}(x)$ which is independent of θ .

(*) coordinate descent in $D(\hat{p}(x) || p(x | \theta))$ (coordinates - parameters of clique potentials)

- Pick a clique C

- Adjust clique potential $\psi_C(x_c)$ to min KL-div. (X)

$$\Rightarrow D(\hat{p}(x) || p(x | \theta)) = \underbrace{D(\hat{p}(x_c) || p(x_c | \theta))}_{(I)} + \sum_{x_c} \hat{p}(x_c) \underbrace{D(\tilde{p}(x_{v \setminus c} | x_c) || p(x_{v \setminus c} | x_c, \theta))}_{(II)}$$

(*) A little dense; lecture also helps with exposition.

- changes in clique pot $\psi_C(x_c)$ does not affect $p(x_{v \setminus c} | x_c, \theta)$

\Rightarrow (I) unaffected by changes to $\psi_C(x_c)$ and minimising KL div $D(\hat{p}(x) || p(x | \theta))$ amounts to minimising (I)

(*) Minimising (1) achieved by setting $p(x_i|\theta) = \hat{p}(x_i)$ i.e. marginal to empirical marginal

(*) This what IPF achieves \rightarrow ⁽⁶¹⁾ coordinate ascent in log-like
coordinate descent in KL-div

(*) Note IPF takes form of scaling algorithm in which potentials are multiplied by a ratio of marginals.

(*) $\max l \Leftrightarrow \min KL(\hat{p}(x) \| p(x|\theta)) = \sum_x \hat{p}(x) \log \frac{\hat{p}(x)}{p(x|\theta)}$ (X)

(*) Verifying identity (A):- (KL div. decomp.) (46)

$$D(p(x_A, x_B) \| q(x_A, x_B)) = \sum_{x_A, x_B} p(x_A) p(x_B|x_A) \log \frac{p(x_A) p(x_B|x_A)}{q(x_A) q(x_B|x_A)}$$

$$= \sum_{x_A, x_B} p(x_A) p(x_B|x_A) \log \frac{p(x_A)}{q(x_A)} + \sum_{x_A, x_B} p(x_A) p(x_B|x_A) \log \frac{p(x_B|x_A)}{q(x_B|x_A)}$$

$$= \sum_{x_A, x_B} p(x_A, x_B) \log \frac{p(x_A)}{q(x_A)} + \sum_{x_A} p(x_A) \sum_{x_B} p(x_B|x_A) \log \frac{p(x_B|x_A)}{q(x_B|x_A)}$$

$$= \sum_{x_A} \log \frac{p(x_A)}{q(x_A)} \sum_{x_B} p(x_A, x_B) + \sum_{x_A} p(x_A) \sum_{x_B} p(x_B|x_A) \log \frac{p(x_B|x_A)}{q(x_B|x_A)}$$

$$= \sum_{x_A} p(x_A) \log \frac{p(x_A)}{q(x_A)} + \sum_{x_A} p(x_A) D(p(x_B|x_A) \| q(x_B|x_A))$$

$$= D(p(x_A) \| q(x_A)) + \sum_{x_A} p(x_A) D(p(x_B|x_A) \| q(x_B|x_A))$$

(*) Features and micropotentials

(*) lecture notes (5/2019) are a little clearer on feature design.

(61): Rather than define a table of values to encompass the mapping of $\psi(x_i) = \psi(x_1, x_2, x_3)$ over e.g. 26^3 possibilities

⑥ we instead use domain knowledge (e.g. a vocabulary/dictionary) to reduce representational granularity

⑦ Achieved by assigning a score to common three letter stems; or scores to 'significant' occurrences e.g. $\text{fing} = 10$, $\text{fted} = 9$ etc. (X)

⑧ remainder \rightarrow assign an arbitrarily low score to reflect low likelihood of occurrence.

⑨ since in which 'feature' is a function which is 'vacuous' over all joint settings except a few part. ones.

⑩ or use indicators?

mathematical details - define K features $f_k(c_1, c_2, c_3)$ e.g. bingting

- Assign each of K features a weight θ_k

- to avoid simult.

motiv.

A micropotential is then achieved by exponentiating $e^{\theta_k f_k(c_1, c_2, c_3)}$ estim. / setting θ_k AND f_k

⑪ A clique potential is formed by multiplying together micropotentials

- yielding: $\psi_c(x_c) = \psi_c(x_1, x_2, x_3) = e^{\theta_1 f_1} \cdot e^{\theta_2 f_2} \cdot \dots \cdot e^{\theta_K f_K}$

$$= \exp \left\{ \sum_{k=1}^K \theta_k f_k \right\}$$

- Potential is still over 26^3 settings; but only have K parameters if K features are used

⑫ we then can estimate the weights θ_k to deal with context. info

⑬ θ_k - 'strength' of feature and whether it increases or decreases clique probability

(*) Combining features

- Jordan (2003) \rightarrow suggests f_k is chosen to be indicator

- marginal probability over a clique: $p(c_1, c_2, c_3) \propto \exp \left\{ \sum_{k=1}^K \theta_k f_k \right\}$

- Addit. complexities \rightarrow overlapping features, must make; any subset of clique variables, overlapping wholeness do not alter anatomy

① unique potential as weighted sum of exponen. features.

$$\psi_c(x_c) = \exp \left\{ \sum_{i \in \mathcal{I}_c} \theta_i f_i(x_{c,i}) \right\}$$

(*) features based model

- usually, $p(x_c | \theta) = \frac{1}{z(\theta)} \prod_c \psi_c(x_c) = \frac{1}{z(\theta)} \exp \left\{ \sum_c \sum_{i \in \mathcal{I}_c} \theta_i f_i(x_{c,i}) \right\}$

- But drop explicit assoc. of features and cliques (??) and use

$$p(x | \theta) = \frac{1}{z(\theta)} \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

where $f_i(\cdot)$ - features
 θ_i - param
 $z(\theta)$ - norm. factor

② exponential family model with features as sufficient statistics

$$z(\theta) = \sum_x \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

③: estimate θ_i from data \mathcal{D} .
 (param)

rule of feature based MCMC: (see Jordan 2003) (problem setup)

- is ④

scaled log-likelihood: $\tilde{l}(\theta; \mathcal{D}) = \frac{l(\theta; \mathcal{D})}{N} = \frac{1}{N} \sum_{n=1}^N \log p(x_n | \theta)$

②②

④ $\tilde{l}(\theta; \mathcal{D}) = \sum_x \hat{p}(x) \log p(x | \theta)$

$$= \sum_x \hat{p}(x) \log \left[\frac{1}{z(\theta)} \left\{ \exp \sum_i \theta_i f_i(x) \right\} \right]$$

$$= \sum_x \hat{p}(x) \left\{ \sum_i \theta_i f_i(x) \right\} - \log z(\theta)$$

- Jordan (2003) states we use convexity of $\log(\cdot)$ to bound $\log z(\theta)$ (?)

- We then get:-

$$\log z(\theta) \leq \mu z(\theta) - \log \mu - 1$$

(O/S 3): $\log(\cdot)$ is not convex, it's concave

(O/S 4): what conseq for presentation?

(*) for now, assume this is okay \rightarrow add to
overspill and investigate at the end.

(*) Bound holds for all μ and $\mu = z^{-1}(\theta^{(t)})$

$$\Rightarrow \hat{\ell}(\theta|\mathcal{D}) \geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{z(\theta)}{z(\theta^{(t)})} - \log z(\theta^{(t)}) + 1$$

(O/S 5)

with equality at $z(\theta^{(t)})$

(*) further manipulation of scaled log-like.

define $\Delta\theta_i^{(t)} := \theta_i - \theta_i^{(t)}$, then (subs. $z(\theta)$)

$$\hat{\ell}(\theta|\mathcal{D}) \geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{z(\theta^{(t)})} \sum_x \exp\left\{\sum_i \theta_i f_i(x)\right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{z(\theta^{(t)})} \sum_x \exp\left\{\sum_i (\Delta\theta_i^{(t)} + \theta_i^{(t)}) f_i(x)\right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{z(\theta^{(t)})} \sum_x \exp\left\{\sum_i \theta_i^{(t)} f_i(x)\right\} \exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \underbrace{\sum_x \frac{1}{z(\theta^{(t)})} \exp\left\{\sum_i \theta_i^{(t)} f_i(x)\right\}}_{= p(x|\theta^{(t)})} \exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_i \theta_i \sum_x \tilde{p}(x) f_i(x) - \underbrace{\sum_x p(x|\theta^{(t)}) \exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\}}_{\text{coupled } f_i(x) \text{ and } \Delta\theta_i^{(t)}} - \log z(\theta^{(t)}) + 1$$

assume $f_i(x) \geq 0$ and $\sum_i f_i(x) = 1$

if $\exp(\cdot)$ is convex, invoke Jensen's inequality.

(O/S 5a)

$$\exp\left(\sum_i \pi_i x_i\right) \leq \sum_i \pi_i \exp(x_i) \quad \text{for } \sum_i \pi_i = 1$$

(*) Note; f_i play the role of π_i as they are positive and sum to 1.

(*) Cross-ref with lecture 7 notes

↳ f_i are being treated as 'weights' and $\Delta\theta_i^{(t)}$ as arguments

yielding the following lower bound on scaled log-likelihood: - (with params denoted)

$$\hat{\ell}(\theta; D) \geq \sum_i \theta_i \sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x}) - \sum_{\underline{x}} p(\underline{x} | \theta^{(t)}) \sum_i f_i(\underline{x}) \exp(\Delta\theta_i^{(t)}) - \log z(\theta^{(t)}) + 1$$

$$= \Lambda(\theta)$$

• We then take derivatives with respect to lower bound $\Lambda(\theta)$: -

$$\frac{\partial \Lambda}{\partial \theta_i} = \sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x}) - \exp(\Delta\theta_i^{(t)}) \sum_{\underline{x}} p(\underline{x} | \theta^{(t)}) f_i(\underline{x}) = 0$$

$$\Rightarrow \exp(\Delta\theta_i^{(t)}) = \frac{\sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x})}{\sum_{\underline{x}} p(\underline{x} | \theta^{(t)}) f_i(\underline{x})}$$

• defining (or because?) $p^t(\underline{x})$ is an unnormalised version of $p(\underline{x} | \theta^{(t)})$

$$\Rightarrow \exp(\Delta\theta_i^{(t)}) = \frac{\sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x})}{\sum_{\underline{x}} p^t(\underline{x}) f_i(\underline{x})} z(\theta^{(t)})$$

②: next bit in Jordan & Xing not clear (lecture notes + slides) from re-con; attempted to clarify

$$(*) \theta_i^{(t+1)} = \theta_i^{(t)} + \Delta\theta_i^{(t)} \Rightarrow p^{(t+1)}(\underline{x}) = p^{(t)}(\underline{x}) \left(\prod_i e^{\Delta\theta_i^{(t)} f_i(\underline{x})} \right) \quad (?) \quad \textcircled{0/57}$$

Jordan (2003):

(*) To update parameters from $\theta^{(t)}$ to $\theta^{(t+1)}$, multiply $p(\underline{x} | \theta^{(t)})$ by $e^{\Delta\theta_i^{(t)} f_i(\underline{x})} \forall i$

(*) Note: $\frac{p^{(t)}(\underline{x})}{z(\theta^{(t)})} = p(\underline{x} | \theta^{(t)})$ and \rightarrow