## L7 - Maximum likelihood learning of undirected GM

EX: what is key algorithm for POGMs?

Ⓐ: EM

- 10-708: see lots of algorithms; develop taste and understanding
- POGMS: DO inference on unobserved; then apply completely observed tools
   (heuristic)
- better researchers dig out foundations

eg. EM as co-ordinate ascent algorithm (characterising it this
                             way can place it in a class)

EX: see rationale behind algorithm

___

EX: use graphical models to pull together local structures

___

MLE for BNS:

(*) most important:-   $\theta_{ijk}^{ML} = \dfrac{n_{ijk}}{\sum_{i,j',k} n_{ij'k}}$      (Ⓐ: looks like counts/
                                           empirical probability )

(however
   lage GM is)

(*) due to factorisibility of DGM

Ⓠ: Does this apply to UGM?

___

MLE for undirected GMS

(*) UGM: Hammesley-Clifford means we can define a UGM in terms
         of a Gibbs distribution and partition function

$$p(x_1,...,x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c) \qquad Z = \sum_{x_1,...,x_n} \prod_{c \in C} \psi_c(x_n)$$

· $Z$ - partition fn
  - normalisation constant of product of unnormalised potentials

     $\prod \psi_c(x_c)$                                pcc

EX: suppose $\psi$ contains hidden parameters (e.g. strength of config
        (potentials)                                in clique)       given complete
                                     and you want to estimate ; observations

@ Can you compute parameters within $\psi$?

easily

- ⓥ: No; you will have the following:- $\frac{1}{Z}\prod_{c \in C} \psi_c(x_c, z)$

- so parameters will appear inside clique potential functions
- You will need to marginalise along the lines of

$$\prod_{c \in C} \sum_{z} \psi_c(x_c, z)$$

($\psi_c$ not probability; but assume an analogous operation)

$X \quad \begin{matrix} -\psi_c(x_c) \\ ARE \end{matrix}$

- so product-sum difficult for ML estimation

- Not complete the parameters wrong → you have written latent variables

(4) coupling is the answer

EX: Nothing explicit to optimise against (if doing MLE); as you will have unknown parameters; and hence $\frac{1}{Z}\prod_{c \in C}\psi_c(x_c)$ will be. be unknown

EX: Some graphical models structures can also be described by U.G.M.

(*) log-likelihood for UGMS with tabular clique potentials.

sufficient stats:

- UGM $(V, E)$: the no. times config $x$ i.e. $\underline{X} = \underline{x}$ is observed in a data set $D = \{x_1, ..., x_N\}$ can be represented as follows:-

Define

$\cdot m(x) = \sum_n \partial(x, x_n)$ (total count) ⑥π

$\cdot m(x_c) = \sum_? m(x)$ (clique count) ⑥π

Ⓦ ④

- Total counts - no. of time a configuration appears in dataset
- clique counts - no. of times a particular configuration within a clique appears in the dataset

(*) Clique counts obtained by marginalising our total counts Ⓦ Ⓐ2

: Assume discreteness

Log-likelihood:- $p(D|\theta) = \sum_c \sum_{\underline{x_c}} m(x_c) \log \psi_c(x_c) - N \log Z$     (*)

Ⓦ Ⓐ3 : check you understand how log-like is specified (quick)

ex: Log-like :- Sum over all possible <u>configurations</u> of $\underline{x}$
       - use delta function to clamp those values of $x | x_n$ (?)
       which are consistent with your observations of the data,
       count 1 everytime you see it.

   - connects log-likelihood with ms (i.e. <u>counts</u>) obtained from <u>data</u>

(*) Do you understand how log-like and observations / sufficient statistics
for UGMs are related?
       ①

Q: What is $\theta$ (parameters?) (so you remember L3!)     Ⓦ Ⓐ4 : Refresh
                                                        param. of
                                                        UGMs !

   - we are compute scientists, not mathematicians
   - Do not matter to theoreticians (actually dealing with messiness)
   - You were close → CPDS ; but do not obey constraints of prob.
   - An <u>unnormalised table of nos</u>    $\psi_c(x_c)$ - $x_c$ associated with
                                                      a no. e.g $\alpha$

---

(*) Derivative of LL :-
- standard calculus :

1st term     $\dfrac{\partial L_1}{\partial \psi_c(x_c)} = \dfrac{m(x_c)}{\psi_c(x_c)}$

2nd term : Ⓦ Ⓐ5 - Review this (quick)

conditions on clique marginals

- get optimal $\psi_c^*$ → find it vanishes

(*) At ML setting of parameters, for each clique, model marginals equal to observed marginals (empirical counts)

EX: Only get marginal probability of a clique $p^*_{mle}(x_c)$; we want the potential function (estimates of) of each clique.
  - These are not the same in UGMS (✓)

(*) Only provides condition that must be satisfied when we have ML param; does not specify how to get ML param.

EX: serious work into doing this

## MLE for UGMS

EX: Previous iterations relied on these concepts (decision tree style questions)

  - triangulated
  - clique potentials defined on maximal cliques
  - full tables or compact

⎫ - skip
⎬ - see Koller for historical account

2 workhorse algorithms (most insightful)    ⎰ nos. behind eigenconfig

• IPF (Iterative proportional fitting)  -  MRFs tabular pot.

- GIS (Generalised iterative scaling)  -  MRFs with features potential

· EX: key is how the differences in problem scopes yield to differences in algorithmic approach
  algebraic tricks → make problem easier (significantly)

## IPF

· identity from LL optimisation → anti-climactic
- How to recover from this?

(*): From LL:

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

(W)(PB): - Review this algo.

$$p^*_{mle}(x_c) = \frac{m(x_c)}{N} = \hat{p}(x_c)$$

(*) Derive:-

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

---

· Turn identity into a fixed point equation (endow identity with time component)
  for $\psi_c$

$$\psi_c^{(t+1)} = \psi_c^{(t)}(x_c)\frac{\hat{p}(x_c)}{p^{(t)}(x_c)}$$

(*) update fn: $\frac{\hat{p}(x_c)}{p^{(t)}(x_c)}$ — proportion of empirical marginal (countable from data) over current version of estimated marginal; based on your model. (derivable from $\psi_c^{(t)}(x_c)$)

(*) In UGM; even with observed data; have to do inference

Additional question:   i) Does it converge etc.?

---

Properties of IPF updates:

IPF is a fixed point program over time; but also over potential functions
(*) N?: our potentials
(*) A co-ordinate ascent algorithm; attaining an optima in a parti· direction
     when other directions fixed.

- convergence somewhere.
- (*) Also known as I-projection (distn from one space to another
                          where only one potential is allowed
                              to change)

- our space of possible distri families
   (via attained via max-entropy)

(*) understood
via KL divergence view → comes up in V.I/D.L (Jordan 11)

_____

KL divergence view

- ML can be reframed as KL divergence
- coordinate ascent charac. of IPF through KL divergence (via info theory)

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

- Partition arguments of distri into:-

$X_c$ and $X_{c1}$        c1 - complement of c

- To carve out a parti. potential clique

(*) combine $\ominus$ KL with conditional chain rule ⊖ Ⓐ⑤ : review ✓

_____

IPF minimises KL divergence

(*) changing $\psi_c$ (clique potential) has no effect on c.d.
(2nd term unaffected)

(*)  $KL(\hat{p}(x) \| p(x|\theta)) = KL(\hat{p}(x_c) \| p(x_c|\theta)) + \sum_c \hat{p}(x_c) KL(\hat{p}(x_{c1}|x_c) \| p(x_{c1}|x_c))$

- i.e. setting $p(x_c) = \hat{p}(x_c)$

Ⓗ Ⓐ⑥ : quick review

_____

· IPF
- start with random guess of potential nos. $\psi_c^{(0)}(x_c)$
- multiply by a ratio. $\frac{\hat{p}(x_c)}{p^{(t)}(x_c)}$        (proportional)

- convexity only qualifies whether local or global (in this context)
- initialise random no. generator 100 times and run (to deal with conveg/local/global)

-