

17 - maximum likelihood learning of undirected GM

Ex: what is key algorithm for PGMs?

Q: EM

- 10-708: see lots of algorithms; develop taste and understanding
 - PGMs: do inference on unobserved; then apply completely observed tools (heuristic)
 - better researchers dig out foundations
- eg. EM as coordinate 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

Mot for BNs:

- (*) most important: - $\theta_{ijk}^{ml} = \frac{n_{ijk}}{\sum_{i,j,k} n_{ijk}}$ (Q: looks like counts / empirical probability)
 (however log GM is)
- (*) due to factorisability of GM
- Q: Does this apply to UGM?

~~not for undirected GMs~~

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

$$P(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c) \quad Z = \sum_{x_1, \dots, x_n} \prod_{c \in C} \psi_c(x_c)$$

• Z - partition fn

- normalisation constant of product of unnormalised potentials

$$\prod_c \psi_c(x_c)$$

Ex: suppose ψ contains hidden parameters (e.g. strength of config in clique)
 given complete
 and you want to estimate i observations
 (potentials)

Q) Can you compute parameters within ϕ ?
easily

try to engage with
life of flight

- Q: No; you will have the following:- $\frac{1}{2} \prod_{c \in C} \phi_c(x_c, z)$

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

$$\prod_{c \in C} \sum_z \phi_c(x_c, z)$$

(it's not probability; but assume
an analogous operation)

$$- \phi_c(x_c) \\ X \quad \text{ARE}$$

- so product sum difficult for ML estimation

(*) Coupling is the answer

- Not compute the parameters wrong

→ you have written
latent variables

Ex: Nothing explicit to optimise against (if doing MLE); as latent variables
you will have unknown parameters; and hence $\frac{1}{2} \prod_{c \in C} \phi_c(x_c)$ will be
be unknown

Ex: Some graphical models structures can also be described by UGM.

(*) Log-likelihood for UGMs with tabular clique potentials.

Sufficient stats:

- UGM (V, E) : the no. times config z i.e. $X = z$ is observed in a data

set $D = \{z_1, \dots, z_N\}$ can be represented as follows:-

define $m(z) = \sum_n \delta(z, z_n)$ (total count) ①

② ③

$$m(z_c) = \sum ? m(z) \quad (\text{clique count}) \quad ④$$

- 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 over total counts (W@R)

Assume discreteness

$$\text{log-likelihood: } p(\theta | \Theta) = \sum_c \sum_{x_c} m(x_c) \log \psi_c(x_c) - N \log Z \quad (*)$$

(W@R): check you understand how log-like is specified (quick)

ex: log-like: - sum over all possible configurations of \mathbf{x}

- use delta function to choose 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. counts) obtained from data

(*) do you understand how log-like and observations / sufficient statistics for ms are related?

Q: What is θ (parameters?) (so you remember L3!)

(W@R): refresh perm. of ms!

- we are computer scientists; not mathematicians
- do not matter to theoreticians (actually dealing with messiness)
- You were close \rightarrow CPDs; but do not obey constraints of prob.
- an unnormalised table of nos $\psi_c(x_c)$ - x_c associated with a no. e.g. α

(*) derivative of LL:-

- standard calculus:

$$\frac{\partial L}{\partial \psi_c(x_c)} = \frac{m(x_c)}{\psi_c(x_c)}$$

2nd term: (W@R) - Review this (quick)

conditions on clique marginals

- get optimal ψ_c^* \rightarrow 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 $P^k(\mathbf{x}_k)$; we want the potential function (estimates of) of each clique.

- these are not the same in UCMS (✓)

(*) 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 UCMS

Ex: previous iterations relied on these concepts
(decision tree style questions)

- triangulated
- cliquepotentials defined on maximal cliques
- full tables or compact

} - skip
- see Koller
for historical
account

2 workhorse algorithms (most insightful)

• IPF (Iterative proportional fitting) - MRFS to tabular form.

→ less behind
every config

• GJS (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 \rightarrow 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) (P6) :- 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 fixed point equation (endow identity with time component)
for ψ_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 (actable from data)
our current version of estimated marginal; based
on your model. (derivable from $\psi_c^{(t)}(x_c)$)

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

Algorithmic question: i) Does it converge etc?

Properties of IPF updates:

IPF is a fixed point program over time; but also our potential functions

(*) W) (P7) : Our potentials

(*) A co-ordinate ascent algorithm; attaining an optima in a parti. direction
when other directions fixed.

- converge somewhere.

- (*) Also known as l-projection (dist from one space to another
where only one potential is allowed
to change)

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

(*) understood
via KL divergence view \rightarrow arises w/in V.I/D.L (Jordan II)

KL divergence view

- MC can be reframed as KL divergence
- coordinate ascent chanc. 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 distn into:-

x_c and $x_{\bar{c}}$ c1 - component of c

- to carry out a parti. potential clique
- (*) combine & KL with conditional chain rule $\textcircled{1} \textcircled{2} \textcircled{3} \textcircled{4}$: review ✓

IPF minimises KL divergence

- (*) changing ψ_c (clique potential) has no effect on c.d.
(and term unaffected)

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

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

① ② quick review

- IPF
- start with random guess of potential nos. $\psi_c^{(0)}(z_c)$
- multiply by a ratio, $\frac{\hat{p}(z_c)}{p^{(0)}(z_c)}$ (proportional)
- convexity only qualifies whether local or global (in this context)
- initial random no. generate 100 times and run (to deal with convergence/local/global)

- ex: 4th year PhD did not know how to run k-means
- you want to know theorem; but you have to get it working (Q)
 - (Q): implement these?

- potential functions $\psi_C(x_C)$

- some issues: - spell-checking example
 - build affinity models of charact. streams
 - e.g. consec. appearances of 3 streams
 - use this criterion to score likelihood
- $\psi_C(x_C) = \psi_C(x_1, x_2, x_3)$
- 'xyz' and 'ink'
- define potential function over a triplet of characters
- 26³ different features and put a function on it?
- 26⁵

(*) above shows infeasibility of tabular potentials (with respect to enum. of joint)
feature based clique pot. (Q): note instantiation of cliques in practice

(*) Features

- Q: find a way of appropriately compressing the granularities
- Handcrafted feature design \rightarrow role of human knowledge in A.I.
- use feature-engine to save on rep. cost?
- distinct from current, overcomplete ML models (e.g. gigantic placeholders)
- ex: against the idea of human knowledge being ignored in ML

(*) micropotentials for OR
- Have K features and weights define our 3 charact. potential

$$f_c(a_1, a_2, a_3) = \exp \left\{ \sum_{k=1}^K \theta_k f_k(a_1, a_2, a_3) \right\}$$

(*) overall potential (clique) is exp weighted sum of micropot.

(*) micropotent. distinct from tabular potentials.

(*) K parameters over K features \rightarrow more compact

Combining features

(*) sliding window / overlapping sliding window

(*) Note how we can modify standard Gibbs. rep. for exp.
- measure of exponential/GUMS

(*) not entirely clear how to apply IPF in this case due to
coupling of estimated θ_k and assigned $f_k(a_1, a_2, a_3)$

Mix of feature based UGMs:

- scaled likelihood:- $\hat{l}(\theta; D) = l(\theta; D)/N = \frac{1}{N} \sum_n \log p(x_n | \theta)$

$$= \sum_x \hat{p}(x) \log p(x | \theta)$$

$$= \sum_x \hat{p}(x) \sum_i \theta_i f_i(x) - \log Z(\theta) \quad (i) \quad (ii)$$

(*) calculate derivatives \rightarrow not fruitful.

Ex: Nonlinearities causing issues (e.g. log-linearization)

- unlike if so argument can be exposed to linear attack

- log has linear upper bound

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

$$- Bound holds \forall \mu : \mu = Z^{-1}(\theta^{(t)})$$

fixed point
it. strategy

- assume this (there is a previous
version of Z)

(*) GIS derivation ② A7: Review

- define $\Delta\theta_i^{(t)} = \theta_i - \hat{\theta}_i^{(t)}$ and introduce
- still nasty: \rightarrow every

(*) note exp of weighted sum:- $\exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\}$

Ex: we make distinction between weight θ_k and $(\Delta\theta_i^{(t)})$ and features $f_i(x)$; }
algebraically the same.

- treat f_i as weights; $\Delta\theta_i^{(t)}$ as arguments
- (*) impose prob. constraints (normally applying to weights) to
 f_i our assumed "weights".

(*) $\exp(\cdot)$ is convex \rightarrow use Jensen's

- Algebraic trick often used in ML
 - getting $\sum_i f_i(x) \exp(\Delta\theta_i^{(t)})$ \rightarrow only linearly coupled with others.

(*) use lower bound of scaled LL:- GIS

(*) use calculus:-

(*) writing update steps:-

$$\text{(*) Note: } - \frac{\sum_x \hat{p}(x) f_i(x)}{\sum_x p^{(t)}(x) f_i(x)}$$

"weighted sum of feature"
by emp. prob. $\hat{p}(x)$ "

"feature weighted
by inferred probability $p^t(x)$ "

(*) Iterative re-scaling \rightarrow connection with IPF.

(*) Summary GIS/IPF

- fixed point iterations on LL obj.
- one for tabular, feature based

(*) where does exponential come from? (0: A move from Gauss \rightarrow Gibbs?)

⑥ ⑦: review exp. family form (IS)

(*) note at MLE; expectations of sufficient statistics
model match ^{emp.} feature average

- Note eq.

(*) Begin with exp. family \rightarrow get a consequence.

- reversible - impose constraint on distri: to you can't give me n. expectat. of arbitrary ones
- maximum entropy - fixed feat. exp: $\sum_x p(x) f_i(x) = x_i$: feature must match (from data)
 - encode complex few assumpt. about model
 - entropy as amt as poss. of randomness / amt of assumptions made

$$\max_{\mathbf{p}} H(\mathbf{p}(\mathbf{x})) = - \sum_{\mathbf{x}} \mathbf{p}(\mathbf{x}) \log \mathbf{p}(\mathbf{x})$$

$$\text{s.t. } \sum_{\mathbf{x}} \mathbf{p}(\mathbf{x}) f_i(\mathbf{x}) = x_i$$

$$\sum_{\mathbf{x}} \mathbf{p}(\mathbf{x}) = 1 \quad \rightarrow \mathbf{p}(\mathbf{x}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \boldsymbol{\theta}_i f_i(\mathbf{x}) \right\}$$

(*) variational definition: - define a distib. as a solution to a constrained optimisation problem.

⑧ ⑨: review lagrangian sol.

(*) natural consequence gives: an exponential family distri.

(*) Benefit of information theoretic principles of ML \rightarrow ⑩ Explore

- more general max entropy method

- incorporate prior distib on \mathbf{x} ; reflect it. ($h(\mathbf{x})$)

- estimated distib has least addit. assumptions from priors

- use KL-divergence rather than entropy

$$\min_p \text{KL}(p(x) || h(x)) = \sum_x p(x) \log \frac{p(x)}{h(x)} = -H(p) - \sum_x p(x) \log h(x)$$

s.t. $\sum_x p(x) f_i(x) = \alpha_i$

$$\sum_x p(x) = 1$$

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

(*) constraints from data

(**) where do constraints α_i come from?

- data itself is the constraint

- (*) (**): automatic consistency?

- geometric interp; general process:-

either:-

i) Assume all exponential family distns as "model":-

$$E = \{ p(x) : p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp \left\{ \sum_i \theta_i f_i(x) \right\} \}$$

^{oR} ii) Assume all distns satisfying model constraints

$$M = \{ p(x) : \sum_x p(x) f_i(x) = \sum_x \hat{p}(x) f_i(x) \}$$

do not acknowledge
roots

- information geometry Pythagorean theorem:-

→ inspires V.I; deep gen. models

(*) Summary

(*) exp family viewed as a sol. to variational exp \rightarrow maximum entropy

(*)

Supplementary → structure learning (see supp.) / 2020/

case study: CRFs (graphy) - at CMU

- insight of explicit modelling

- latterly paper → impressive; interpretable, motivation

(*) local normalisability is a double-edged sword

(- makes computing simple)

(?) an - model of HMM

(*) what you want is global normalisability

- use scores rather than enforcing local normalisability

- use potentials

$$\exp \left\{ \sum \theta_i f_i \right\}$$

- features w.r.t. nodes

- use human knowledge for features

- Art of modelling → totally complex.