

continue U2(*) variational principle (V.P.)

ex: Highly technical \rightarrow review material with assigned readings.

- contents cannot be assimilated immediately in particular.

recap: Algebraic problems / probabilistic inference reformulated as variational inference

- define solution to partition function (quantity of interest) in variational manner.

- normally computed (part. function) via summation/integration (exponentially complex)

$$(*) A(\theta) = \sup_{\mu \in M} \{ \theta^T \mu - A^*(\mu) \}$$

- inner product of dual pair
- conjugate dual

- corresponds to canonical form in exp family M-marginal polytope
- μ - " - mean parameters

- μ - satisfies constraints as expectation of sufficient statistics $\mu = E$
negative

- A^* entropy of original distribution

- *) Solution only exists to this problem when mean params, satisfy constraints; that they lie within a domain called marginal polytope

- *) Marginal polytope - convex combination of sufficient statistics

- *) Marginal polytope - convex hull

- *) In discrete case \rightarrow convex hull
cont. case \rightarrow - " -

- *) This is key; use two-node Ising Model as example. and review.

- *) Example - 2 node Ising model

- *) Note Minkowski-Weyl rep of marginal polytope

(A1)

- (*) optimum \rightarrow stationary points
 - (*) more complex problems; we cannot use above variational principle to derive results of a qualitatively exact nature (that is the point of introducing V.P.)
 - (*) V.P. is a vehicle for approximations
-
- (*) variational principle
- (*) understand MF approxim. in context of V.P.
 - (*) note: from readings, introducing V.P. in this way shows / reformulates the intractability of summation/inference
 - it also opens possibilities on dealing with intractability of summation
→ dealing with intractability of marginal polytope and entropy.
 - Approximation schemes can be viewed as approx. the marginal polytope and entropy terms.

- (*) MF approx:-
- i) Redefine domain of mean parameters to be in a subspace of the true marginal polytope M
 - ii) redefine A^* (negative entropy) in a way that is less intractable

(*) formally : MF approx uses :-

- i) non-convex inner bound $M_{MF} \subset M$
- ii) exact form of entropy $A^* = H$

(*) Bethe approx and LBP:-

- i) polyhedral inner bound $M_{Bethe} \supset M$
- ii) non-convex Bethe approx. $\rightarrow A^* = H \approx M_{Bethe}$

(?) ex: states that entropy is approx with Bethe free energy
(surely he means Bethe entropy?)

(4) Mean-Field APPROX.

- tractable subgraphs

- MF approx relies on using a tractable subgraph approx. to original graph.

- for a graph with canonical params $\Omega := \{\theta \in \mathbb{R}^d \mid A(\theta) < +\infty\}$
or each edge

- set Ω is defined in such way that the range of log-normaliser is bounded.

- exp family with sufficient statistics ϕ defined on graph G :-

$$\mathcal{M}(G; \phi) := \{\mu \in \mathbb{R}^d \mid \exists p \text{ s.t. } E_p[\phi(X)] = \mu\}$$

- define a subgraph T_0 with no edges (still a graph - it has nodes)

- corresponding canonical parameter set:-

$$\Omega(T_0) := \{\theta \in \Omega \mid \theta_{(s,t)} = 0 \forall (s,t) \in E\}$$

$$\theta_{ij} = 0$$

- what are corresponding mean params $\mu_{i,j}$ (expected suff statistic
defined on pair (i,j))

$$\theta_{ij} = 0$$

- As $\mu_{i,j} = p(x_i, x_j)$; $\theta_{ij} = 0 \Rightarrow \mu_{i,j} = p(x_i)p(x_j)$ ($\theta_{ij} = 0$ (lack of edges)
induces independence)

- refine another subgraph T (properties vary)

$$\Omega(T) := \{\theta \in \Omega \mid \theta_{(s,t)} = 0 \forall (s,t) \notin E(T)\}$$

(i) You can always define a subgraph approximation - an approximation
to a graph by removing edges

(ii) can also give geometric approx to marginal polytope

(iii) above is high-level treatment.

(*) Mean Field methods

- In more detail:-

- for a given tractable subgraph F , a subset of canonical params is :-

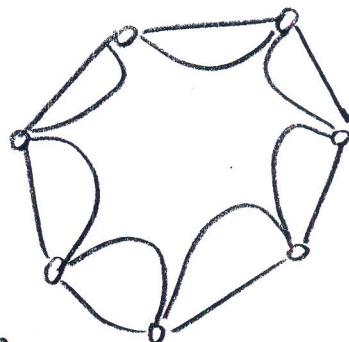
$$M(F; \phi) := \{ \tau \in \mathbb{R}^d \mid \tau = E_\theta[\phi(x)] \text{ for some } \theta \in \Omega(F) \}$$

- $\Omega(F)$ defines a new space for canonical params. corresp. to Θ constrained by subgraph.

- innerapprox to the full mag. polytope):-

$$M(F; \phi)^* \subseteq M(G; \phi)^*$$

- Recall marginal polytope is convex



- innerapproximation of marginal polytope lies within the marginal polytope, because it uses a subset (of canonical parameters associated with the tractable subgraph F)

- for MF approx, we define the (τ) parameters to lie within the innerapprox to marginal polytope:-

- relaxed problem:

$$\max_{\tau \in M_F(G)} \{ \langle \tau, \theta \rangle - A_F^*(\tau) \}$$

* Compute dual of the partition function i.e. entropy of a distribution that corresponds to a tractable subgraph with no edges

* Subgraph with no edge:-

$p(x_1, \dots, x_N)$ approximated by $Q_F(x_1, \dots)$ where Q represents joint over graph with no edges

* For graph with no edges:-

$$Q_F(x_1, \dots, x_N) = \prod_{n=1}^N q(x_n)$$

$$(*) \text{ To compute } A_{DEF}^* = H_Q - \sum H(q(x_i))$$

(A2) - unclear
- well
- tidy

(*) complex models e.g. topic models ; H_P is not tractable

(*) Stepping back;

we have a relaxed problem with an analytic exp for $A_F^*(\tau)$
and easier exp for approx. to marginal polytope $M_F(G)$

(*) Naive MF approx for Ising

- (B) - review example

- solve: take derivative w.r.t τ to get MF update equations
MF approx sol. to topic models etc is a similar method. flavour generates iterative solution.

(*) Geometry of MF approx

MF opt. is always non-convex

full space is convex \rightarrow subsets / space within cannot be convex
(marginal polytope)

(*) strict subset

(*) will an approximation that is equivalent to the solution result?

(*) No. of issues:-

i) true solutions lie within $M(G) - M_F(G)$ i.e. outside of naive approx.

ii) Non-convexity \rightarrow not guaranteed global optima

(*) Benefit: Simple algorithm that offers solution to a complex prob.
man iterative fashion

(*) sum-product / BP.

- key equations:- M.P.P., marginals
minimises Bethe free energy

* Tree-GMS

- * Bethe approx yields exact solution for tree-structured G.M.S.
- * In this context the marginal polytope \rightarrow tree polytope
- (A4) review equations
- * local consistency + global consistency constraints (for trees)
on marginal probabilities / mean parameters
- * local consistency $\not\Rightarrow$ global consistency generally

(* Decomposition of entropy for trees)

- (B) review the variational form on trees

(* Lagrangian deriv.)

- reparametrisation of constrained minim. of Bethe free energy for tree-structured G.M.S
yields exact inference; and B.P. update equations.

B.P. on arbitrary graphs

- 2 main difficulties of variational formulation:-
- In earlier MF approx, he relied on inner approx. of marginal polytope
- Variational form:

$$A(\theta) = \sup_{\mu \in M} \{ \theta^T \mu - A^*(\mu) \}$$

- Here, use a tree-based outer bound on the marginal polytope M :-

$$L(G) = \left\{ \tau \geq 0 \mid \sum_{x_s} \tau_s(x_s) = 1, \sum_{x_t} \tau_{st}(x_s, x_t) = \tau_s(x_s) \right\}$$

- only impose local consistency on vectors $\tau_s(x_s)$ and $\tau_{st}(x_s, x_t)$
which are pseudom marginals

Q: Is the space $\mathbb{H}(G)$ larger or smaller than the marginal polytope M ?

(*) Key nuance here:-

The tree structured GM and local consistency conditions inspire constraints; but we are not working in that setting, this is an arbitrary graph.

(*) There are global constraints not present here; hence $\mathbb{H}(G)$ is larger (?)

Q: don't understand

(*) For arbitrary distri, $A^*(\mu)$ is also complicated (graphs)

(*) Approximate $A^*(\mu)$ with Bethe free energy (possibly entropy)

$$-A^*(z) \approx H_{\text{Bethe}}(z) = \sum_{s \in V} H_s(z_s) - \sum_{(s,t) \in E} I_{st}(z_{st})$$

(*) H_{Bethe} only requires pairwise and singleton sufficient statistics

(*) Bethe variational Problem (marg. polyt.) (entropy)

Combining $\mathbb{H}(G)$ and $A^*(z)$ i.e. approx of M and $A^*(\mu)$:-

Bethe V.P.: $\max_{z \in \mathbb{H}(G)} \left\{ \langle \theta, z \rangle + \sum_{s \in V} H_s(z_s) - \sum_{(s,t) \in E} I_{st}(z_{st}) \right\}$

(*) MF approx uses $M \approx M_F(G)$ and $A^*(\mu) = H_q$
(more approx of marginal polytope) (true entropy)

(*) Ex: why do we use true entropy H_q for $A^*(\mu)$; whereas here we use Bethe approx to entropy?

Q
we
As we approx. P with an MF approx. as product of marginal q distri (fully factorised) ✓ correct

In loopy BP, the approximation being used was NOT explicit in the same manner as MF approx.

(*) For loopy B.P.; we only requested queries
 $E[\phi(x_i)]$ and $E[\phi(x_i, x_j)]$; with no regard for
singleton pairwise exp. what else lies in full distri. $p(x)$

(*) Treat the above, corresponding to near params. of singleton and pairwise
and define a space based on local consistency constraints; then define
approximation to entropy only using these 2 statistics.

(*) Bypasses explicit representation of approx. of P, \tilde{P} ; but the full distri
has more edges and structure in entropy; near case of this;
but set aside.

(A6) - review rest of results here.

(*) Geometry of B.P. (cont'd approx)

- for any graph; $M(G) \subseteq E(G)$; with equality $M(G) = E(G)$ if and
only if the graph is a tree

- concerning entropy; $H_{\text{Bethe}} \hat{=} H$

If graph is a tree $H_{\text{Bethe}} = H$

(A7) Tidy entropy functionals

(*) Yields insight on behavior of B.P.

- may or may not converge

(*) Review BP convergence in context of outer approx to marginal
polytope geometrically

(*) Inexactness of Bethe entropy approx.

(Q) Review

(*) Remark & Summary

Kikuchi variants rarely used nowadays

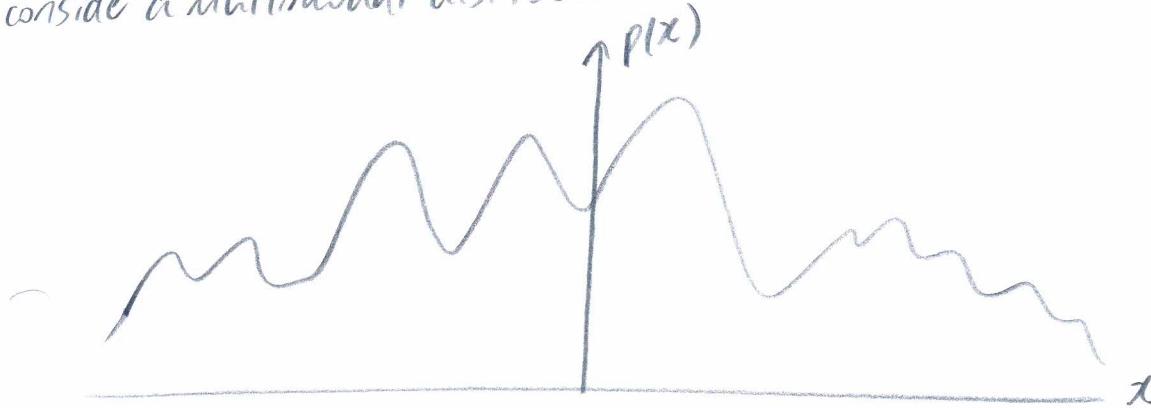
U3 - Approx Inference: Monte Carlo and sequential Monte Carlo

(*) How to represent a joint model.

- $E_p[f(x)] = \int f(x) p(x) dx$

- sample based representation:

- consider a multimodal distribution (and density): -



(?) How to integrate?

- It may not be feasible (analytically)? to integrate this.

- So: Draw samples from distri

- computationally: will use a uniform PRNG to draw a sample $u(0,1)$ and transform to get Gaussian.

(*) draw N samples $x^{(i)} \sim P$:-

$i=1, \dots, N$

compute $E_p[f(x)] = \frac{\sum_{i=1}^N f(x^{(i)})}{N}$

} Monte Carlo methods

(*) original distri P expressed as a set of samples

(*) Monte Carlo methods

(*) replace integration with summation (sample-based averages)

(*) Asymptotically exact \rightarrow more samples drawn ($N \rightarrow \infty$) arbitrarily close to expectation we wish to approx.

Challenges:

- 1) How to draw samples from distri?
- How to do this efficiently?
from an arbitrary
- 2) How to make better use of samples?
- asymptotic - how many samples can we draw ...

(*) Naive sampling

$$P(x_1, \dots, x_m) = P(x_B) P(x_A | x_B) \dots$$

$$\underline{x}^{(k)} = \begin{bmatrix} x_A \\ \vdots \\ x_m \end{bmatrix} \sim P$$

- How to draw samples?
- traverse graph
- (*) CX walks through realisation of sample in BN
- Sample r.v.s. from B.N.

- (*) Rows of table correspond to a 'sample' from B.N.
- (*) Computation of queries: check samples where condition is met.
- (*) Rare events \rightarrow conditioning on rare events?
- (*) Shows the limitations of sampling

(*) Monte Carlo methods (rejection sampling)

- attempt to get around above deficiencies
- $\pi(x) = \frac{\pi'(x)}{z}$ $\pi(x)$ - difficult to sample from
 $\pi'(x)$ - easy to evaluate
 (e.g. unnormalised clique pot. product)
- sample from simpler distri $Q(x)$

(*) $x^* \sim Q(x)$ (draw samples from Q)
 - accept x^* with probability $\frac{\pi'(x^*)}{RQ(x^*)}$ R -arbitrary constant

(*) Ratio of normalised part of given (evalatable distn) and a constant multiplied by proposal distribution.

(*) with these two steps:-

$$x^* \sim \pi(x)$$

Ex: How to prove this?

$$\begin{aligned} P(x^*) &= Q(x^*) \frac{\pi'(x^*)}{RQ(x^*)} \\ &= \frac{R\pi'(x^*)}{R \int \pi'(x^*) dx^*} = \pi(x) \end{aligned}$$

(*) A solution to sampling from a distn $\pi(x)$ which is difficult to normalise (using a proposal distn)

Pitfall: (i) geometric intuition / explanation of too many rejections ?

(*) Rejection sampling

- conceptual experiment to illustrate rejection sampling

- small differences between proposal distn and distn of interest

- caused by scaling constant can cause a large difference in rejection rate

$$R = \left(\frac{\sigma_q}{\sigma_p}\right)^d \approx \frac{1}{30,000}$$

(*) reactual vol. of 2 distns are vastly different in high-dim space \Rightarrow reject samples most of time even though distns are close

(*) solution: don't use one $\neq Q$ -distn (proposal distn) to cover the target distn, but use a piece-wise envelope

(*) Adaptive rejection sampling is an example of this

- so rejection sampling doesn't work well in a high-dim space

(*) unnormalised importance sampling

- draw, as before, samples $x^{(i)} \sim Q(x)$ from a proposal distri. N times

- instead of accepting/rejecting as earlier, compute weight based on new sample

$$w^{(i)} = \frac{P(x^{(i)})}{Q(x^{(i)})}$$

- store $\{x^{(i)}, w^{(i)}\}_{i=1}^N$ as representation of true distri. samples

(*) why is P suddenly available/tractable?

- assume in simplest case that P can be eval (up to norm. constant)

(ii) assess proof, picture

- proposal distri introduced as a dummy in proof

(*) unnormalised as weights directly comp. from likelihood ratio (do not sum to one).

ex. what if I can only evaluate true distri up to a constant? ②

e.g. no exact value of norm. constant
(define to be arbitrary α)

(*) normalised importance

sampling

- only can evaluate $p'(x) = \alpha P(x)$ (e.g. MRF)

- get α and normalisation constant:-

(*) take expectation of the ratio

$$r(x) = \frac{p'(x)}{Q(x)} \Rightarrow \langle r(x) \rangle_Q = \int \frac{p'(x)}{Q(x)} Q(x) dx = \int p'(x) dx = \alpha$$

- expectation of ratio is normalisation constant

$$\alpha = \frac{\langle r^{(i)} \rangle}{\sum_{i=1}^N \langle r^{(i)} \rangle}$$

- now draw $x^{(i)} \sim Q(x)$ $i=1, \dots, N$ (N samples)
 - compute $r^{(i)}$ for each i^{th} sample.
 - compute α via $\alpha = \sum_{i=1}^N r^{(i)}$
- (A3) - confused here
- review derivations
- Proof:
- (*) Compute expectations of function using samples drawn as above
 - (*) Instructor presentation lost me. (*) don't need to know norm const or pdf function of target
 - (*) Normalised vs unnormalised importance sampling
 - (A4) - review at home (instructor does not emphasise).

- (*) Efficiency of likelihood weighting
- So far, we have used likelihood weighting \rightarrow may not get 'the' answer in peculiar/path scenarios
 - consider practical implementation of MCMC. When do we stop sampling?
typically stop when stability/convergence observed as empir. strategy for controlling iteration/algo duration.
 - in a technically equivalent way; can view this as terminating when variance \downarrow of current solution.
 - makes sense. (i.e. convergent behaviour through variance)
 - (*) So variance in a sense can be viewed as measure of closeness to true answer.
(of current algo output)
 - (*) In importance sampling; we can have a pathological case:-
 - (A5) Review; watching + record distribution
 - (*) Example of poor proposal distri (note we do not know its poor); does not envelop true distri.
 - (*) Rerunning weights $w^{(i)}$ will correct this.
 - Solution: use heavy tailed $Q \rightarrow$ (but many samples wasted)
- weighted resampling from $\{(x^{(i)}, w^{(i)})\}_{i=1}^N$

- we resample from $\{w^{(i)}\}_{i=1}^{l_{2N}}$ N' times where $N' \gg N$.
- (*) 'amplify' distn using resa weighted resampling ⑯ Review weighted resampling

(*) weighted resampling

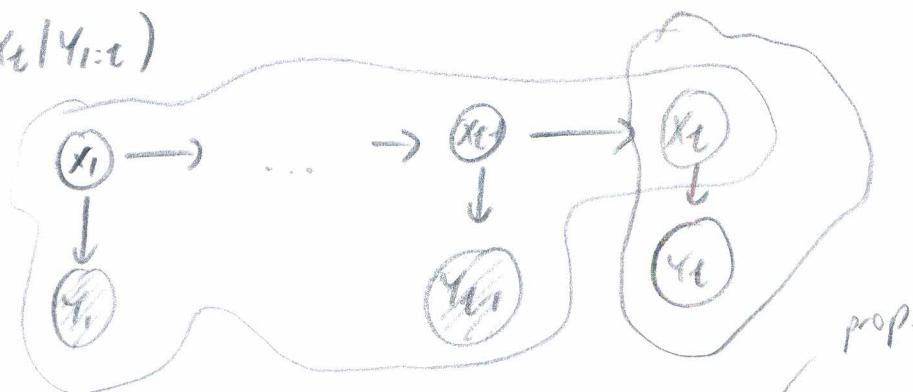
(*) Particle filtering

- use resampling idea to elegant, efficient inference.
- make a fast sampling based algorithm for SSMs.
- KF demanding to implement in high-dm space; or non-Gaussian (Gaussian approx)

(*) sketch of particle filtering

- use KF to establish recursive procedure.

$$p(x_t | y_{1:t})$$



- break down $p(x_t | y_{1:t}) \rightarrow 2$ parts.

$$p(x_t | y_{1:t}) = p(x_t | y_t, y_{1:t-1}) = \frac{p(x_t | y_{1:t-1}) p(y_t | x_t)}{\int p(x_t | y_{1:t-1}) p(y_t | x_t) dx_t}$$

weights'

- inspect $p(x_t | y_{1:t-1})$ and use this as a proposal distn. from which we sample

- represent $p(x_t | y_{1:t})$ using abot decompr:

$$\left\{ x_t^{(i)} \sim p(x_t | y_{1:t-1}), w_t^{(i)} = \frac{p(y_t | x_t^{(i)})}{\sum_{i=1}^N p(y_t | x_t^{(i)})} \right\}$$

If hence $p(x_t | y_{1:t})$ is expressed as weighted resamples from proposal distn.
 (*) goal:- use weighted sample representation of $p(x_t | y_{1:t})$ to
 accurately infer $p(x_{t+1} | y_{1:t+1})$

(*) sequential weighted resampler

time update: keep evidence; propagate forward with no new measurement

$$p(x_{t+1}|y_{1:t}) = \underbrace{\int p(x_{t+1}|x_t)}_{(i)} \underbrace{p(y_t|x_t)}_{(ii)}$$

→ given by weighted res. rep.

• i) Given via model

(*) replace integration with sampling :-

$$p(x_{t+1}|y_{1:t}) = \sum_{i=1}^N w_t^{(i)} p(x_{t+1}|x_t^{(i)})$$

- weighted sum
of transition prob.
and all. on samples of
previous states.

(*) Similar to mixture model (sampling form)

(*) uses new evidence:-

measurement update :- $p(x_{t+1}|y_{1:t+1}) =$

$$\frac{p(x_{t+1}|y_{1:t}) p(y_{t+1}|x_{t+1})}{\int p(x_{t+1}|y_{1:t}) p(y_{t+1}|x_{t+1}) dx_{t+1}}$$

(*) note (*) is computed from time update

- replace by weighted resampling version:-

⇒ $p(x_{t+1}|y_{1:t+1})$ has representation:-

$$\left\{ x_{t+1}^{(i)} \sim p(x_{t+1}|y_{1:t}), w_{t+1}^{(i)} = \frac{p(y_{t+1}|x_{t+1}^{(i)})}{\sum_{i=1}^N p(y_{t+1}|x_{t+1}^{(i)})} \right\}$$

renormalized

(*) particle filtering graph (illustrating resampling)

- size of balls correspond to weights

② - review

(*) Particle filtering

gives online results for predictive distri

(*) PF for SSMs (switching)

- allows drawing from more complex distri e.g. switching SSM
- multiple hidden distri's.

(*) Rao-Blackwellised sampling

(*) review slides (not emphasized)

(*) Statistical properties of unnormalised/normalised importance sampling.

- unnormalised importance sampling is unbiased

$$\text{i.e. } \mathbb{E}_Q[f(x)w(x)] = \mathbb{E}_P[f(x)]$$

- normalised importance sampling is biased for finite samples e.g. $M=1$

$$\mathbb{E}_Q\left[\frac{f(x')r(x')}{\sum r(x')}\right] \neq \mathbb{E}_P[f(x)]$$