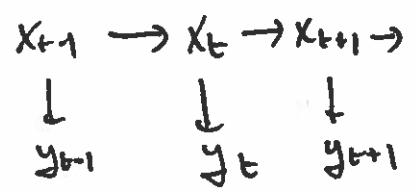


SUMMARY

- EP INFERENCE FOR GENERAL BAYESIAN NETWORKS ↖ High dimensional integrals of approx. posterior replaced by lower dimension integral
- APPROXIMATE MOMENT COMPUTATIONS USING IMPORTANCE SAMPLING ↖ REPLACE
- LEARN TO APPROXIMATE MOMENT COMPUTATIONS USING A LINEAR MODEL ↖ generalized ↖ speed up inference
- USE UNCERTAINTY IN MOMENTS TO DECIDE WHETHER TO REVERT TO IMPORTANCE SAMPLING

EP

- consider chain BNs to avoid bookkeeping



$$\begin{aligned}
 p(x_{1:t} | y_{1:t}) &= \prod_t p(x_t | x_{t-1}) p(y_t | x_t) \\
 &= \prod_t \underbrace{f_t(x_t, x_{t-1})}_{\text{approximate affect of } f_t \text{ on true posterior}}
 \end{aligned}$$

$$q(x_{1:t}) = \prod_t \alpha_t(x_t) \beta_t(x_{t-1})$$

(NB: NOT REALLY MEANFIELD)
eg. for LSSMs recovers exact (correlated) Kalman filter equations

ITERATIVE UPDATE RECIPE

1. FORM CAVITY

$$\frac{q(x_{1:t})}{\alpha_t(x_t) \beta_t(x_{t-1})}$$

2. INCLUDE NEW POTENTIAL

$$\frac{q(x_{1:t})}{\alpha_t(x_t) \beta_t(x_{t-1})} \frac{f_t(x_t, x_{t-1})}{p(x_t | x_{t-1}) p(y_t | x_t)} =$$

3. MOMENT MATCH

(match mean & variance of $q(x_{1:t})$ for function of a bn undet'd α_t & β_t)

$$\frac{q(x_{1:t})}{\alpha_t(x_t) \beta_t(x_{t-1})} p(x_t | x_{t-1}) p(y_t | x_t) \stackrel{\text{mom}}{=} \frac{q(x_{1:t})}{\alpha_t(x_t) \beta_t(x_{t-1})} \alpha_t^{\text{new}}(x_t) \beta_t^{\text{new}}(x_t)$$

As $q(x_{1:t}) = \prod_t q(x_t) \Rightarrow$ mean of $x_{t,t-1}$ already matched
+ variance

\Rightarrow moment matching reduces to

$\frac{q(x_t, x_{t-1})}{\alpha_t(x_t) \beta_t(x_{t-1})} \overset{\text{mom}}{=} \frac{\alpha_t(x_t) \beta_t(x_{t-1}) \alpha_{t-1}(x_{t-1}) \beta_{t-1}(x_t)}{\alpha_t(x_t) \beta_t(x_{t-1})} = \frac{q(x_t, x_{t-1})}{\alpha_t(x_t) \beta_t(x_{t-1})} \alpha_t^{\text{new}}(x_t) \beta_t^{\text{new}}(x_{t-1})$

$\Rightarrow \beta_{t+1}(x_t) \alpha_{t+1}(x_{t-1}) p(x_t | x_{t-1}) p(y_t | x_t) \overset{\text{mom}}{=} \beta_{t+1}(x_t) \alpha_{t-1}(x_{t-1}) \times \alpha_t^{\text{new}}(x_t) \beta_t^{\text{new}}(x_{t-1})$

Now consider update for $\alpha_t(x_t)$ ie $\langle x_t \rangle \langle x_t^2 \rangle$ moments

$\int \beta_{t+1}(x_t) \alpha_{t-1}(x_{t-1}) p(x_t | x_{t-1}) p(y_t | x_t) dx_{t-1} \overset{\text{mom}}{=} \int \beta_{t+1}(x_t) \alpha_{t-1}(x_{t-1}) \alpha_t^{\text{new}}(x_t) \beta_t^{\text{new}}(x_{t-1}) dx_{t-1}$

{ as $\int \alpha_{t-1}(x_{t-1}) \beta_t(x_{t-1}) dx_{t-1} = \text{const.}$ & does not affect 1st & 2nd moments }

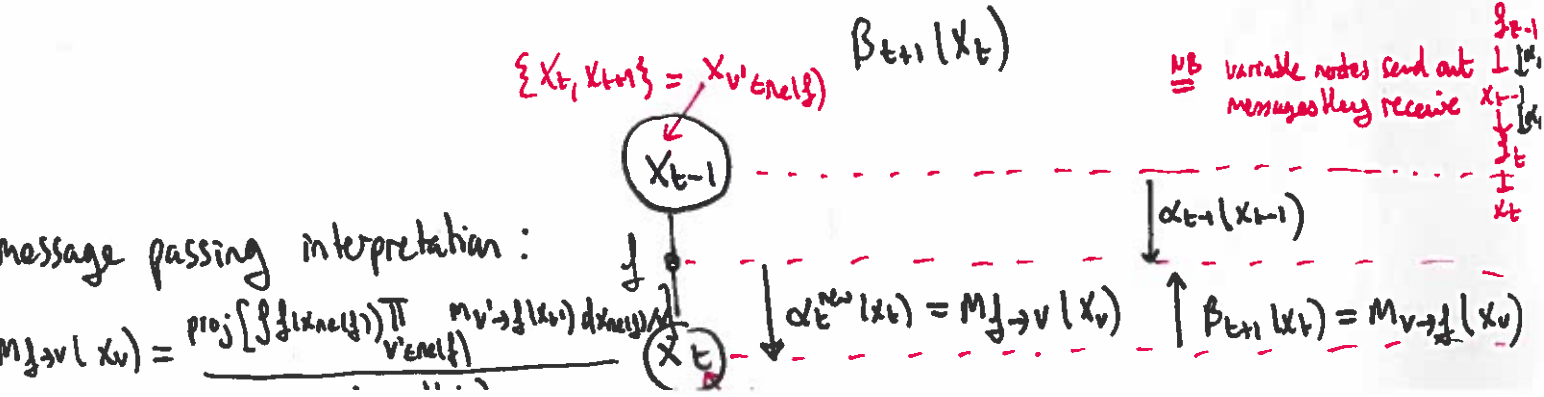
$\int \beta_{t+1}(x_t) \alpha_{t-1}(x_{t-1}) p(x_t | x_{t-1}) p(y_t | x_t) dx_{t-1} \overset{\text{mom}}{=} \beta_{t+1}(x_t) \alpha_t^{\text{new}}(x_t)$

project argument onto a Gaussian w/ same 1st & 2nd moments
 \Rightarrow ONLY REQUIRES 2D INTEGRALS

$\Rightarrow \alpha_t^{\text{new}}(x_t) =$

$\text{proj} \left[\int \beta_{t+1}(x_t) \alpha_{t-1}(x_{t-1}) p(x_t | x_{t-1}) p(y_t | x_t) dx_{t-1} \right]$

message passing interpretation:



IMPORTANCE SAMPLING FOR EP (Barthelmé & Chopin, 2011)

③

For approximations, $q(x)$, that are in the exponential family w/ moments $u(x)$
key challenge is to compute:

$$\phi = \int f(x_{n+1}) \prod_{v' \in n+1} M_{v'} \rightarrow \int f(x_v) u(x_v) dx_{n+1}$$

for Gaussian $u(x_v) = x_v^2$ & x_v

IDEA: USE IMPORTANCE SAMPLING

$$q(x_t, x_{t-1}) = p(x_t | x_{t-1}) r(x_{t-1})$$

$$\begin{aligned} \phi &= \int dx_{t:t-1} x_t^k l(x_t, x_{t-1}) = \int dx_{t:t-1} x_t^k \frac{l(x_t, x_{t-1})}{q(x_t, x_{t-1})} q(x_t, x_{t-1}) \\ &= \int dx_{t:t-1} x_t^k w(x_t, x_{t-1}) q(x_t, x_{t-1}) \approx \frac{1}{N} \sum_n (x_t^{(n)})^k w(x_t^{(n)}, x_{t-1}^{(n)}) \end{aligned}$$

where $\{x_{t-1}^{(n)}, x_t^{(n)}\} \sim q(x_t, x_{t-1})$

$$w(x_t, x_{t-1}) = \frac{l(x_t, x_{t-1})}{q(x_t, x_{t-1})} = \frac{\beta_{t+1}(x_t) \alpha_{t-1}(x_{t-1}) p(y_t | x_t) p(x_{t+1} | x_t)}{r(x_{t-1}) p(x_{t+1} | x_{t-1})}$$

paper does not consider terms like this could incorporate into q .

- Asymptotically unbiased \Rightarrow consistent

- involves importance sampling each time we need to send a message

LEARNING TO PASS MESSAGES

(4)

Hess et al 2013, Eslami 2014, Jitkrittum et al 2015
 (neural networks) no uncertainty initial tracking step \Rightarrow not online
 (random forests) uncalibrated heuristic uncertainty not online
 (this paper) calibrated uncertainty online + active

Idea: learn a mapping from $\{\alpha_{t-1}(x_{t-1}), \beta_{t+1}(x_t)\}$ to moments μ_t \therefore
 $\{\alpha_t^{(new)}(x_t), \beta_t^{(new)}(x_t)\}$

Problem

inputs: eg. Gaussian Distributions \uparrow non-standard outputs: moments

Idea: design simple linear in the parameters model w/ clever basis functions by appealing to GPs.

let $r^{(i)}$ denote incoming messages \uparrow Bayesian $r_i = \{\alpha_{t-1}^{(i)}(x_{t-1}), \beta_{t+1}^{(i)}(x_t)\}$ different examples in training set

$K_2(r^{(i)}, r^{(j)}) = e^{-\frac{1}{2\sigma^2} \int_{-\infty}^{\infty} \dots}$ similarity between lots of incoming messages

to handle distributions use:

$$K_2(\alpha^{(i)}(x), \alpha^{(j)}(x)) = \exp\left(-\frac{1}{2\sigma^2} \int (\alpha^{(i)}(x) - \alpha^{(j)}(x))^2 dx\right)$$

\uparrow
integral over products of Gaussians \Rightarrow tractable

to be slightly more general (connect to the authors' previous work) can blur messages first by eg a Gaussian

\downarrow Gaussian

$$\mu_r^{(i)}(x) = \int K_2(x, x') \alpha(x') dx'$$

inputs \vec{x}_t, \vec{x}_{t-1} → outputs $\frac{1}{\Delta}$

we use sufficient statistics to represent these distributions

we mean embeddings to represent these

Gaussian.

Step 1: blur $\int_{\vec{x}_{t-1}} \mathcal{K}(\vec{x}_t, \vec{x}_{t-1}) d\vec{x}_{t-1} = \int_{\vec{x}_{t-1}} \mathcal{K}(\vec{x}_t, \vec{x}_{t-1}) \mu(\vec{x}_{t-1}) d\vec{x}_{t-1}$

$\mu(\vec{x}_{t-1}) = \int \mathcal{K}(\vec{x}_t, \vec{x}_{t-1}) \mu(\vec{x}_t) d\vec{x}_t$

Step 2: define kernel between two probability distributions by average square difference between the distributions

$K_2(\vec{x}, \vec{s}) = e^{-\frac{\|\mu_{\vec{x}} - \mu_{\vec{s}}\|_H^2}{2\delta^2}} = e^{-\frac{1}{2\delta^2} \int_{-\infty}^{\infty} (\mu_{\vec{x}}(x) - \mu_{\vec{s}}(x))^2 dx}$

will be difference between two Gaussians whose means & variances depend on the natural parameters of $[\alpha_{t-1}^{(1)}, \beta_{t-1}^{(1)}] \parallel [\alpha_{t-1}^{(2)}, \beta_{t-1}^{(2)}]$

Step 3: do GP regression using the Gaussian

BUT Step 3 has $\mathcal{O}(N^3)$ cost \Rightarrow actually don't do GP regression at all, but approximate using sparse GP

parametric model where

Rahimi+Recht: $k(x-y) = \int K(\omega) e^{j\omega^T(x-y)} d\omega$ [Bochner's theorem]
 $= \int p^*(\omega) e^{j\omega^T(x-y)} d\omega \approx \sum_{k=1}^K e^{j\omega_k^T(x-y)} = \psi(\omega) \eta(\omega) \sim p^*(\omega)$

NOW WE USE THIS RESULT TWICE:

Now consider kernel choice: $\int \mu_{\vec{x}}(\vec{x}) \mu_{\vec{x}'}(\vec{x}') d\vec{x} = \int K_1(\vec{x}, \vec{x}') K_1(\vec{x}, \vec{x}') \mathcal{B}(\vec{x}'') d\vec{x}' d\vec{x}'' = \mathbb{E}_{\vec{x}''} \mathbb{E}_{\vec{x}'} [G_1(\vec{x}' - \vec{x}'')]$

$\int K_1(\vec{x}, \vec{x}') K_1(\vec{x}, \vec{x}'') d\vec{x} = G_1(\vec{x}' - \vec{x}'')$

approx $k(x^i - x^j)$ as $\sum_{k=1}^K \psi_k(x) \psi_k(y)$ & then can analytically integrate over $A(x)$ & $B(x)$

$$\Rightarrow \int p_A(x) p_B(x) dx = \sum_{k=1}^K \langle \psi_k(x) \rangle_{A(x)} \langle \psi_k(y) \rangle_{B(y)} = \hat{\beta}_A^T \hat{\beta}_B$$

$$- \|\hat{\beta}_K - \hat{\beta}_S\|_{\mathcal{H}}^2 \quad (\text{now just a vector norm})$$

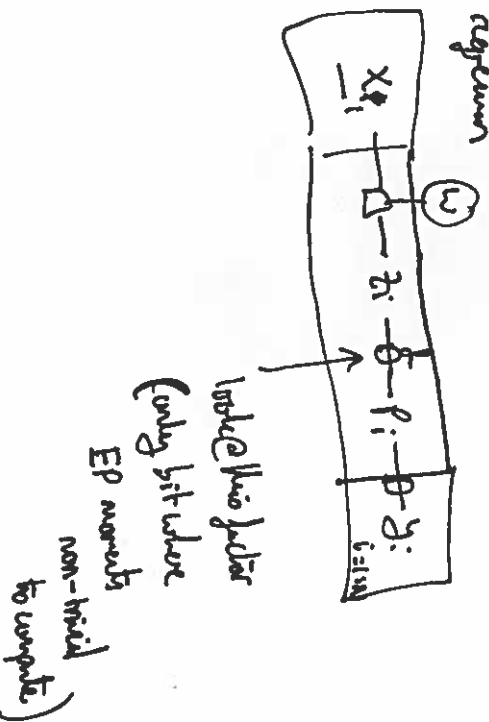
$$\therefore k(\Gamma_S) \cong e$$

Now apply RFR to this (SE) kernel!

From HERE STANDARD ON LINE LEARNING & PREDICTION ON A LINEAR - IN - THE - PARAMETERS REGRESSION MODEL.
[NB. SMART BASIS FUNCTIONS]

RESULTS

Binary logistic regression



$$p(w) = N(0, I)$$

$$p(z|w, x) = \delta(z; -w^T x_i) \quad \text{linear so can be handled exactly in EP}$$

$$p(r_i | z_i) = \delta\left(r - \frac{1}{1 + e^{-z}}\right)$$

$$m_{r_i} \rightarrow \frac{1}{2} = \text{beta } |r_i; \alpha, \beta$$

$$m_{z_i} \rightarrow \frac{1}{2} = N(z_i; \mu, \sigma^2)$$

$$p(y_i | r_i) = p_i^{y_i} (1 - p_i)^{1 - y_i} \quad \text{Bernoulli so can be handled exactly in EP}$$