Improved Mean and Variance Approximations for Belief Net Responses via Network Doubling

Peter Hooper¹, Yasin Abbasi-Yadkori², Russ Greiner², and Bret Hoehn² ¹Department of Mathematical & Statistical Sciences, ²Department of Computing Science, University of Alberta

Introduction

Query: what is probability of A = a, given (B = b, Y = y)?

 θ_{a^1} θ_{a^2} (

 $\begin{bmatrix} \theta_{b^1|a^1} & \theta_{b^2|a^1} \\ \theta_{b^1|a^2} & \theta_{b^2|a^2} \end{bmatrix}$ $\begin{bmatrix} (\beta_{a^1}, \sigma_{a^1}^2) \\ (\beta_{a^2}, \sigma_{a^2}^2) \end{bmatrix} \begin{pmatrix} Y \\ Y \end{pmatrix}$

- A, B binary, Y Gaussian, naïve Bayes, $B \perp Y \mid A$
- $\theta_a = P(A = a), \ \theta_{b|a} = P(B = b \mid A = a), \ Y \mid A = a \sim N(\beta_a, \sigma_a^2)$
- Θ = vector containing all parameters $\theta_a, \theta_{b|a}, \beta_a, \sigma_a$

$$q(\boldsymbol{\Theta}) = P(A = a \mid B = b, Y = y, \boldsymbol{\Theta}) = \frac{\theta_a \theta_{b|a} p(y \mid \beta_a, \sigma_a)}{\sum_{a_1} \theta_{a_1} \theta_{b|a_1} p(y \mid \beta_{a_1}, \sigma_{a_1})} \bullet$$

Bayesian paradigm: Θ random, want posterior distribution of $q(\Theta)$

- Approx mean: replace $\theta_a, \theta_{b|a}$ by means; density by Student's t
- Approx variance: delta method [2]
- uses approx local linearity of function q
- Approx distribution of $q(\Theta)$: Beta

Novel approach to approximations via network doubling:



• For replicates $(A_1, B_1, Y_1) \perp (A_2, B_2, Y_2) \mid \Theta$

$$q(\mathbf{\Theta})^2 = P(A_1 = A_2 = a \mid B_1 = B_2 = b, Y_1 = Y_2 = y, \mathbf{\Theta})$$

= a query in doubled network

- Compute approx mean of $q(\mathbf{\Theta})^2$
- Variance = $E\{q(\boldsymbol{\Theta})^2\} [E\{q(\boldsymbol{\Theta})\}]^2$
- Doubling method uses minor changes in conditioning events
- Doubling works for discrete, continuous, and hybrid networks
- Our paper emphasizes implementation for discrete networks:
- asymptotic analysis and adjustments to improve accuracy
- numerical results and computational issues.

Assumptions and Notation

- Discrete network
- Θ = vector of all CPtable parameters
- $\theta_{b|a} = P\{B = b \mid A = a, \Theta\}$, where A is vector of parents of B
- Conjugate prior: global and local independence, Dirichlet
- Complete data \mathcal{D}
- $q(\boldsymbol{\Theta}) = E\{w(\boldsymbol{H}) \mid \boldsymbol{E} = \boldsymbol{e}, \boldsymbol{\Theta}\}$
- -H, E are vectors of hypothesis and evidence variables
- w is an indicator function, e.g., $I_h(H)$

Important Formula

Plug-in approximation of posterior mean $E\{q(\boldsymbol{\Theta}) \mid \mathcal{D}\}$ [1]

 $E\{q(\boldsymbol{\Theta}) \mid \mathcal{D}, \boldsymbol{e}\} = E\{w(\boldsymbol{H}) \mid \boldsymbol{E} = \boldsymbol{e}, \mathcal{D}\} = q(E\{\boldsymbol{\Theta} \mid \mathcal{D}\}) \quad (1)$

• (\mathcal{D}, e) = augmented data, \mathcal{D} plus partial observation E• Global independence needed, not local independence or Dirichlet

Network Doubling

Apply (1) to doubled network:

• Replicates $(B_1, A_1, H_1, E_1) \perp (B_2, A_2, H_2, E_2) | \Theta$ • Doubled CPtable entries: $\theta^*_{b_1b_2|a_1a_2} = \theta_{b_1|a_1}\theta_{b_2|a_2}$ $H^* = (H_1, H_2), E^* = (E_1, E_2), e^* = (e, e)$ • $w^*(\boldsymbol{H}^*) = w(\boldsymbol{H}_1)w(\boldsymbol{H}_2)$ • Result:

$$q(\boldsymbol{\Theta})^2 = E\{w^*(\boldsymbol{H}^*) \mid \boldsymbol{E}^* = \boldsymbol{e}^*, \boldsymbol{\Theta}\}$$
$$q(\boldsymbol{\Theta}) = E\{w(H_1) \mid \boldsymbol{E}^* = \boldsymbol{e}^*, \boldsymbol{\Theta}\}$$

• Approx $\operatorname{Var}\{q(\boldsymbol{\Theta}) \mid \mathcal{D}\}$ by

$$\operatorname{Var}\{q(\boldsymbol{\Theta}) \,|\, \mathcal{D}, \boldsymbol{e}^*\} = E\{q(\boldsymbol{\Theta})^2 \,|\, \mathcal{D}, \boldsymbol{e}^*\} - [E\{q(\boldsymbol{\Theta}) \,|\, \mathcal{D}, \boldsymbol{e}^*\}]^2 \quad (2$$

• Doubled network retains global independence, so (1) valid:

$$E\{q^*(\boldsymbol{\Theta}^*) \,|\, \mathcal{D}, \boldsymbol{e}^*\} = q^*(E\{\boldsymbol{\Theta}^* \,|\, \mathcal{D}\}) \tag{3}$$

• (\mathcal{D}, e^*) = augmented data, \mathcal{D} plus partial observations (E_1, E_2) • Compute $E\{\Theta^* | \mathcal{D}\}$ from Dirichlet means and covariances • Compute (2) using (3) with $q^*(\Theta^*) = q(\Theta)^2$ and $q^*(\Theta^*) = q(\Theta)$

Adjustments

Adjusted mean and variance approximations:

Means	Variances
$\hat{q}_1 = E\{q(\boldsymbol{\Theta}) \mathcal{D}, \boldsymbol{e}\}$	$\hat{v}_1 = \text{delta method [2]}$
$\hat{q}_2 = E\{q(\boldsymbol{\Theta}) \mathcal{D}, \boldsymbol{e}, \boldsymbol{e}\}$	$\hat{v}_2 = \operatorname{Var}\{q(\boldsymbol{\Theta}) \mid \mathcal{D}, \boldsymbol{e}, \boldsymbol{e}\} \text{ in expression (2)}$
$\hat{q}_3 = \hat{q}_1 - (\hat{q}_2 - \hat{q}_1)$	$\hat{v}_3 = $ expression (5)
$\hat{q}_4 = \hat{q}_1 - \hat{\sigma}_{qr}/\mu_r$	$\hat{v}_4 = $ expression (6)

• Simple adjustments (\hat{q}_3, \hat{v}_3) , more complex (\hat{q}_4, \hat{v}_4) • $Q = q(\boldsymbol{\Theta})$ and $R = P(\boldsymbol{E} = \boldsymbol{e} \mid \boldsymbol{\Theta})$

• Formulae involve moments of posterior distribution of (Q, R)• \hat{q}_4 and \hat{v}_4 use Cov(Q, R) approximation:

$$\hat{\sigma}_{qr} = \frac{(\hat{q}_2 - \hat{q}_1)\mu_r(\mu_r^2 + \sigma_{rr})\{\mu_r(1 - \mu_r) + \sigma_{rr}\}}{\mu_r^3(1 - \mu_r) + \mu_r(1 - 2\mu_r)\sigma_{rr} - \sigma_{rr}^2}$$
(4)

• Adjusted variance approximations obtained by iteratively solving

$$\hat{v}_3 := \frac{\hat{v}_2 + 2(\hat{q}_2 - \hat{q}_1)^2}{1 + 4(\hat{q}_2 - \hat{q}_1)(1 - 2\hat{q}_3)/\{\hat{q}_3(1 - \hat{q}_3) + \hat{v}_3\}} \tag{4}$$

$$\hat{v}_4 := \frac{(\mu_r^2 + \sigma_{rr})\{v_2 + (q_2 - q_4)^2\} - 2\sigma_{qr}^2}{\mu_r^2 + \sigma_{rr} + 4\mu_r \hat{\sigma}_{qr}(1 - 2\hat{q}_4) / \{\hat{q}_4(1 - \hat{q}_4) + \hat{v}_4\}}$$
(6)

Asymptotic Analysis

As effective sample size $m \to \infty$ with $E\{\Theta \mid D\}$ fixed:

- query mean μ_q remains constant
- query variance $\sigma_{qq} = O(m^{-1})$ • errors $\hat{q}_j - \mu_q = O(m^{-1})$ for j = 1, 2 and $O(m^{-3/2})$ for j = 3, 4
- relative errors $(\hat{v}_j \sigma_{qq})/\sigma_{qq} = O(m^{-1})$ for j = 1, 2, 3, 4

- Complexity of Variable Elimination (VE) Algorithm is $O(d^r)$ -d = upper bound on variable domain size -r = upper bound on size of a factor generated by VE
- uating a query, so has corresponding computational complexity
- Doubling method uses technique for variance that is similar to eval-• Doubled CPtables are larger (squared number of rows and columns), so computational complexity of VE is $O(d^{2r})$
- Delta method retains $O(d^r)$ complexity, so faster in large networks
- Sometimes exploit network or query structure for polynomial time

- simulate 10^6 replicates of Θ from posterior
- evaluate each $q(\Theta)$
- calculate sample mean and variance

- BDe posteriors, hyperparameters determined by m and $E\{\Theta \mid D\}$
- Examples from three small networks, each with one vector $E\{\boldsymbol{\Theta} \mid \mathcal{D}\}$ and $m \in \{20, 50, 100, 200, 500\}$:
- Two naïve Bayes networks (NB-2 and NB-4 with 2 and 4 binary features plus the binary root variable H); E = all children of H, e varies over all combinations (2² for NB-2, 2⁴ for NB-4)
- Diamond network $\langle \rangle$, 4 binary variables, 108 distinct queries with one hypothesis variable

Results consistent with asymptotics.

• Query mean:



Computational Issues

• Belief net inference is NP-complete

Numerical Results

- \hat{q}_i and \hat{v}_i compared with empirical estimates of μ_q and σ_{qq} :
- Examples differ with respect to network structure, posterior, and query

$-\hat{q}_3 \approx \hat{q}_4$, both more accurate than \hat{q}_1

• Query Variance:

- $-\hat{v}_3 \approx \hat{v}_4$, both more accurate than \hat{v}_1
- NB networks, \hat{v}_2 accuracy similar to \hat{v}_1
- Diamond network, \hat{v}_2 accuracy similar to \hat{v}_3 and \hat{v}_4
- Variation in empirical variances noticeable for m = 500



Continuous Variables

- parents, coefficients and variance depending on discrete parents
- Conjugate prior is normal-(inverse chi-square)
- t densities for normal (bivariate t for doubled network)
- Work in progress on computational issues

References

- 309-347.
- for belief net inference. Artificial Intelligence 172: 483–513.

• Doubling applies to continuous and hybrid networks (parents of discrete variables must be discrete) but implementation more complex • Continuous variables Gaussian, mean related linearly to continuous

• Predictive distribution substitutes means for CPtable parameters and

[1] G. Cooper and E. Herskovits (1992). A Bayesian method for the induction of probabilistic networks from data. *Machine Learning* **9**:

[2] T. Van Allen, A. Singh, R. Greiner, and P. Hooper (2008). Quantifying the uncertainty of a belief net response: Bayesian error-bars

UAI June 2009