

Approximate Inference

PAI Tutorial

Johannes Kirschner

October 25, 2019

What is statistical inference?

Statistical inference is the process of using data analysis to deduce properties of an underlying probability distribution. (Wikipedia)

- ▷ Use a model to relate data to a quantity of interest.

Examples:

- ▷ Maximum likelihood $x_{\text{ML}} = \arg \max_x P(Y = y|X = x)$
- ▷ **Bayesian inference** $P(X|Y) \propto P(Y|X)P(Y)$.
- ▷ Computing conditional/marginals in a Bayes Net

Bayesian Inference

The **Bayesian model** specifies:

- ▷ $P(Y|X)$, likelihood of Y given hypothesis/latent factors X .
- ▷ $P(X)$ prior probability of X (“belief”).

→ Joint distribution $P(X, Y) = P(X)P(Y|X)$.

Bayesian inference asks to compute, e.g.

- ▷ *posterior distribution* $P(X|Y = y) = \frac{P(Y=y|X)P(X)}{P(Y=y)}$
- ▷ *maximum a posteriori* estimate $x_{\text{MAP}} = \arg \max_x P(X = x|Y = y)$
- ▷ *posterior sample* $x \sim P(X|Y = y)$

But: Evidence $P(Y = y) = \sum_x P(y, x)$ is often difficult to compute.

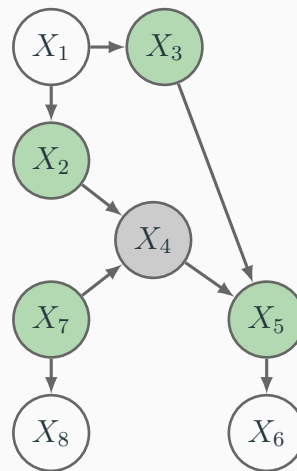
Bayesian Networks

Bayes Net:

- ▷ Defines a joint distribution $P(X_1, \dots, X_m)$
- ▷ Inference: e.g. $P(X_1 | X_4 = T)$

Queries are difficult to compute in general:

- ▷ Need to marginalize.



Approximate Inference

Exact inference can be difficult:

- ▷ Bayesian inference requires to compute the marginal $P(Y)$.
- ▷ Computing conditional/marginals in Bayes Nets not always efficient.
- ▷ Factor distribution $P(X) = \frac{1}{Z} \prod_{i=1}^m \Phi_i(X_{A_i})$.

Goal: Approximate $P(X) = \frac{1}{Z} \Phi(X)$, where $\Phi(X) \geq 0$ and $Z = \sum_x \Phi(x)$.

Two popular approaches: Variational Inference and MCMC

Goal: Approximate $P(X) = \frac{1}{Z}\Phi(X)$, where $\Phi(X) \geq 0$, $Z = \sum_x \Phi(x)$.

Variational Inference

$$Q = \arg \min_{Q \in \mathcal{Q}} \text{KL}(Q \| P)$$

- ▷ Choose a class of distributions \mathcal{Q}
- ▷ Solve optimization problem.
- ▷ Jordan et al. (1999)

Markov Chain Monte Carlo

- ▷ *Idea:* Set up a Markov chain
 - ▷ Define transition kernel $T(X, X')$
 - ▷ With stationary distribution $P(X)$
 - ▷ Metropolis et al. (1953)
- *Sample $P(X)$ from the Markov chain.*

Variational Inference:

$$Q = \arg \min_{Q \in \mathcal{Q}} \text{KL}(Q \| P)$$

Variational Inference: Evidence Lower Bound

Goal: Approximate $P(X) = \frac{1}{Z}\Phi(X)$, where $\Phi(X) \geq 0$, $Z = \sum_x \Phi(x)$.

$$\begin{aligned}\text{KL}(Q\|P) &= \sum_x Q(x) \log \frac{Q(x)}{P(x)} \\ &= \sum_x Q(x) \log Q(x) - \sum_x Q(x) \log \Phi(x) + \log(Z)\end{aligned}$$

$$\arg \min_{Q \in \mathcal{Q}} \text{KL}(Q\|P) = \arg \max_{Q \in \mathcal{Q}} \underbrace{\sum_x Q(x) \log \Phi(x) - \sum_x Q(x) \log Q(x)}_{:= \text{ELBO}(Q)}$$

Evidence Lower Bound for Bayesian Inference

Goal: Approximate $P(X|Y=y) \propto \Phi(X) = P(Y=y|X)P(X)$

$$\text{ELBO}(Q) = \sum_x Q(x) \log(P(Y|x)P(x)) - \sum_x Q(x) \log Q(x)$$

$$= \sum_x Q(x) \log(P(Y|x)) - \text{KL}(Q(X)||P(X)) = \sum_x Q(x) \log \frac{P(Y, x)}{P(x)} - \sum_x Q(x) \log \frac{Q(x)}{P(x)}$$

→ Trade-off between likelihood and prior.

$$\begin{aligned} \text{ELBO}(Q) &= \log P(Y) - \text{KL}(Q(X)||P(X|Y)) \\ &\leq \log P(Y) \end{aligned}$$

$$= \sum_x Q(x) \log(P(x|Y)P(Y)) - \sum_x Q(x) \log Q(x)$$

$$= \sum_x Q(x) \log P(Y) - \sum_x Q(x) \log \frac{Q(x)}{P(x|Y)}$$

$$= \log P(Y) - \text{KL}(Q(X)||P(X|Y))$$

→ Maximize lower-bound on evidence.

Mean Field Variational Inference

Need to choose class of distributions Q .

Mean field approximation: $Q(x) = \prod_{i=1}^m Q_i(x_i)$.

$$\text{ELBO}(Q) = \sum_x \prod_{i=1}^m Q_i(x_i) \log \Phi(x) - \sum_{i=1}^m \sum_{x_i} Q_i(x_i) \log Q_i(x_i)$$

Objective for Q_j :

$$\text{ELBO}_j(Q) = \sum_{x_j} Q_j(x_j) \sum_{x_{j-}} \prod_{i \neq j} Q_i(x_i) \log \Phi(x) - \sum_{x_j} Q_j(x_j) \log Q_j(x_j) + \text{const}$$

log exp ($\bar{z} \pi \dots$)

Exact update for Q_j :

$$Q_j^* = \arg \max_{Q_j} \text{ELBO}_j(Q) \propto \exp \left(\sum_{x_{j-}} \prod_{i \neq j} Q_i(x_i) \log \Phi(x) \right)$$

Mean Field Coordinate Ascent

CAVI (Coordinate Ascent Variational Inference):

- 1: Choose ordering of X_1, \dots, X_m
- 2: Initialize Q_1, \dots, Q_m
- 3: **repeat**
- 4: **for** $j = 1:m$ **do**
- 5: $Q_j \leftarrow \arg \max_{Q_j} \text{ELBO}_j(Q) \propto \exp \left(\sum_{x_{j-}} \prod_{i \neq j} Q_i(x_i) \log \Phi(x) \right)$
- 6: **until** converged

$$\text{ELBO}_j(Q) = \sum_{x_j} Q_j(x_j) \sum_{x_{j-}} \prod_{i \neq j} Q_i(x_i) \log \Phi(x) - \sum_x Q_j(x_j) \log Q_j(x_j)$$

Mean Field Variational Inference: Bayes Nets

Markov Blanket:

$$\text{mb}(X_j) = \text{parents}(X_j) \cup \text{children}(X_j) \cup \text{co-parents}(X_j)$$

$$\triangleright \text{co-parents}(X_j) = \text{parents}(\text{children}(X_j)) \setminus X_j$$

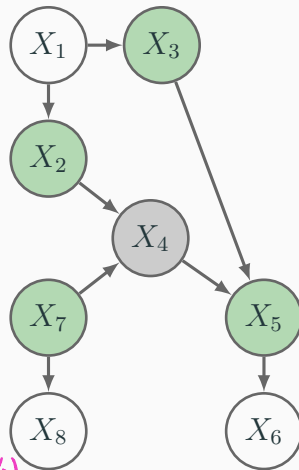
$$\triangleright \forall i \neq j, X_j \perp X_i | X_{\text{mb}(X_j)}$$

CAVI update only depends on Markov Blanket:

$$Q_j \propto \exp \left(\sum_{x_{j-}} \prod_{i \neq j} Q_i(x_i) \log \Phi(x) \right)$$

$$\propto \exp \left(\sum_{x_{\text{mb}(X_j)}} \prod_{i \in \text{mb}(X_j)} Q_i(x_i) \log P(X_j | X_{\text{mb}(X_j)}) \right)$$

只需更新 M.B. 里面的元素, $\text{mb}(X_j)$



Markov Chain Monte Carlo

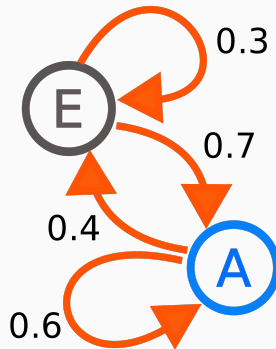
Markov Chains

Markov Chain:

- ▷ Sequence of random variables $X^{(1)}, X^{(2)}, X^{(3)}, \dots$ on \mathcal{S}
- ▷ States space $\mathcal{S} = \{S_1, \dots, S_k\}$
- ▷ Initial distribution: $P(X^{(1)} = S)$
- ▷ Transition kernel: $T(S, S') = P(X^{(t)} = S' | X^{(t-1)} = S)$
- ▷ Can be written as $k \times k$ matrix T .

Stationary distribution: $\pi(S) = \lim_{t \rightarrow \infty} P(X^{(t)} = S)$

- ▷ Satisfies $\pi(S)T(S, S') = \pi(S')$
- ▷ Matrix notation: $\pi^\top T = \pi^\top$ (left eigenvector)
- ▷ Independent of initial distribution (mild conditions)



$$T = \begin{pmatrix} 0.3 & 0.7 \\ 0.4 & 0.6 \end{pmatrix}$$

MCMC: General Idea

Goal: Approximate $P(X) = \frac{1}{Z}\Phi(X)$, where $\Phi(X) \geq 0$, $Z = \sum_x \Phi(x)$.

Markov Chain Monte Carlo (MCMC)

Idea: Construct Markov chain with stationary distribution $P(X)$.

- ▷ Simulating the Markov chain to draw samples $\sim P(X)$

Detailed balance condition:

$$\pi(S)T(S, S') = \pi(S')T(S', S)$$

- ▷ Implies that $\pi(S)$ is the stationary distribution.
- ▷ Define $T(X, X')$ that satisfies detailed balance for Φ :

$$\Phi(X)T(X, X') = \Phi(X')T(X', X)$$

MCMC: Gibbs Sampling

Gibbs Sampling:

- 1: Initialize x_1, \dots, x_m
 - 2: **repeat**
 - 3: **for** $j = 1:m$ **do**
 - 4: $x_j \leftarrow X_j \sim \frac{1}{Z_j} \Phi(X_1 = x_1, \dots, X_j, \dots, X_m = x_m)$
 - 5: **until** converged
- ▷ Draw samples conditioned on all other variables.
 - ▷ We'll show detailed balance for $j \sim \text{Uniform}(\{1, \dots, m\})$

Gibbs Sampling: Detailed Balance

Gibbs sampling generates Markov chain $X^{(1)}, X^{(2)}, X^{(3)}, \dots$

$$\triangleright X^{(t)} = (X_1^{(t)}, \dots, X_m^{(t)})$$

$$T(X, X') = \begin{cases} 0 & \text{if } X_i \neq X'_i \text{ and } X_j \neq X'_j \text{ for } i \neq j \\ \propto P(X_1, \dots, X_m) & X = X' \\ \propto \frac{1}{m} P(X_{1:j-1}, X'_j, X_{j+1:m}) & \exists j, X_j \neq X'_j \text{ and } X_i = X'_i \forall i \neq j \end{cases}$$

Detailed Balance:

- \triangleright Case 1) & 2) are easy
- \triangleright Case 3) $j, X_j \neq X'_j$ and $X_i = X'_i \forall i \neq j$

$$\underline{P(X_{1:m})} \cdot \frac{1}{m} \underline{P(X_{1:j-1}, X'_j, X_{j+1:m})} = \underline{P(X'_{1:m})} \cdot \frac{1}{m} \underline{P(X'_{1:j-1}, X_j, X'_{j+1:m})}$$

Gibbs Sampling on Bayes Nets

Markov Blanket:

$$\text{mb}(X_j) = \text{parents}(X_j) \cup \text{children}(X_j) \cup \text{co-parents}(X_j)$$

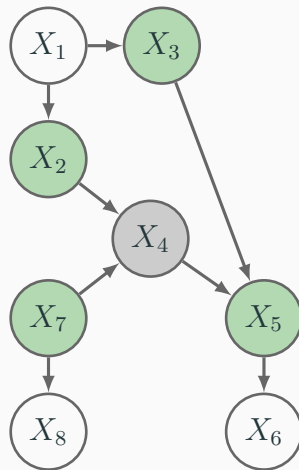
$$\triangleright \text{co-parents}(X_j) = \text{parents}(\text{children}(X_j)) \setminus X_j$$

$$\triangleright \forall i \neq j, X_j \perp X_i | X_{\text{mb}(X_j)}$$

Gibbs sample only depends on Markov Blanket:

$$X_j \sim P(X_{1:j-1} = x_{1:j-1}, X_j, X_{j+1:m} = x_{j+1:m})$$

$$\sim P(X_j, X_{\text{mb}(X_j)} = x_{\text{mb}(X_j)})$$



MCMC and Variational Inference

MCMC:

渐进.

- ▷ Asymptotically correct
- ▷ Computationally expensive, if it takes long to converge
- ▷ Difficult to detect convergence (burn-in)
- ▷ Does not generate iid samples

Variational Inference:

- ▷ Faster
- ▷ Easy to monitor optimization progress
- ▷ Class of distributions \mathcal{Q} often misspecified, i.e. $P \notin \mathcal{Q}$.
- ▷ Optimization can be difficult

Further Reading

- ▷ Bayesian Inference, MCMC and Variational Inference:

<https://towardsdatascience.com/25a8aa9bce29>

- ▷ Technical Review on Variational Inference:

<https://arxiv.org/abs/1601.00670>

- ▷ Forward & Reverse KL:

<https://blog.evjang.com/2016/08/variational-bayes.html>