

BOED

Bayesian optimal experimental design

parameter θ design ξ prior $p(\theta)$

joint $p(y, \theta | \xi) = p(y | \theta, \xi) p(\theta)$ θ not depend on ξ $\theta \rightarrow y$ (outcome)

explicit/implicit = $p(y | \xi) p(\theta | \xi, y)$ \rightarrow posterior

$\mathbb{E}_{p(\theta)} [p(y | \xi, \theta)]$

double intractable: both posterior and its entropy are hard to compute

Info Gain $IG(\xi, y) = \mathcal{H}[p(\theta)] - \mathcal{H}[p(\theta | \xi, y)]$

Expected IG $I(\xi) = \mathbb{E}_{p(y | \xi)} [IG(\xi, y)] = MI(\theta; y | \xi)$

(EIG) $= \mathbb{E}_{p(y | \xi), p(\theta | \xi, y)} [\log p(\theta | \xi, y) - \log p(\theta)] = -\mathbb{E}_{p(y | \xi)} [\mathcal{H}[p(\theta | \xi, y)]] + \mathcal{H}[p(\theta)]$

$= \mathbb{E}_{p(\theta)} [p(y | \theta, \xi) (\log p(y | \theta, \xi) - \log p(y | \xi))] = -\mathbb{E}_{p(\theta)} [\mathcal{H}[p(y | \theta, \xi)]] + \mathcal{H}[p(y | \xi)]$

optimal design $\xi^* = \operatorname{argmax}_{\xi} I(\xi)$

Explicit likelihood

estimate $\hat{p}(y | \xi) \approx \frac{1}{N} \sum_{n=1}^N p(y | \theta_n, \xi)$ $\theta_n \stackrel{iid}{\sim} p(\theta)$

by IS $\mathbb{E}_{p(\theta | \xi, y)} [\log p(y | \theta, \xi)] \approx \frac{1}{N} \sum_{n=1}^N \frac{p(y | \theta_n, \xi)}{\frac{1}{N} \sum_{n=1}^N p(y | \theta_n, \xi)} \log p(y | \theta_n, \xi)$ $y \sim p(y | \xi)$

estimate $I(\xi) \approx \frac{1}{N} \sum_{n=1}^N \log \frac{p(y_n | \theta_n, \xi)}{\hat{p}(y_n | \xi)}$ $(\theta_n, y_n) \sim p(\theta) p(y | \theta, \xi)$

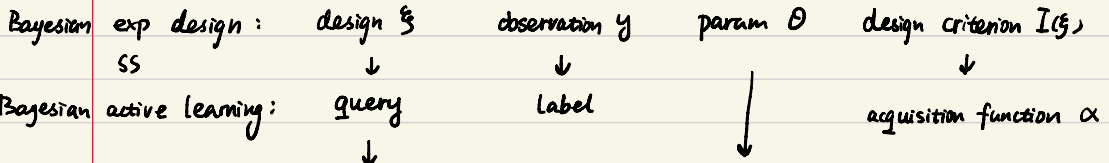
by NMC $\approx \frac{1}{N} \sum_{n=1}^N \log \frac{p(y_n | \theta_n, \xi)}{\frac{1}{M} \sum_{m=1}^M p(y_n | \theta'_m, \xi)}$ $\theta'_m \sim p(\theta)$

Implicit likelihood

ABC $\theta \sim p(\theta | y, \xi) \Leftrightarrow \tilde{y} \sim p(\cdot | \theta, \xi)$, if $\|y - \tilde{y}\| < \epsilon$ then accept θ

LFIRE logistic regression $r(\xi, \theta, y) \rightarrow p(y | \theta, \xi) / p(y | \xi)$

$\Rightarrow I(\xi) = \mathbb{E}_{p(\theta), p(y | \theta, \xi)} [\log r(\xi, \theta, y)] \approx \frac{1}{N} \sum_{n=1}^N \log f(\xi, \theta_n, y_n)$



Pooled-based active learning: unlabeled data classifier NN param
 Alg: For $t=1, 2, \dots$ do: $p(y|\theta, \xi)$: Dropout / deep ensemble

$$\xi_t = \operatorname{argmax}_{\xi \in \Xi} \alpha(\xi; D_{t-1}) \text{ or } \alpha(\xi; \theta_{t-1})$$

$y_t = \text{Human Labeller}(\xi_t)$

$$D_t = D_{t-1} \cup \{(\xi_t, y_t)\}$$

calculate $\theta_t \sim p(\theta|D_t)$ or $\theta_t = \operatorname{argmax} p(\theta|D_t)$

very similar to
 black box - opt!
 (goal is different)

greedy
 acquisition
 ||

$$I(\xi; D_t) = \mathbb{E}_{p(y|\xi, D_t)} \mathbb{E}_{p(\theta|S, y, D_t)} [\log p(\theta|\xi, y, D_t)] - \mathbb{E}_{p(\theta|D_t)} [\log p(\theta|D_t)]$$

$$\xi_t = \operatorname{argmax}_{\xi} I(\xi; D_t)$$

Bayesian active learning by disagreement (BALD score) \equiv EIG

$$\alpha_{\text{BALD}}(\xi; D_t) = \mathbb{E}_{p(\theta|D_t)} [\mathcal{H}[p(y|\xi, D_t)] - \mathcal{H}[p(y|\xi, \theta, D_t)]]$$

$$= \mathbb{E}_{p(\theta|D_t)} p(y|\xi, \theta, D_t) [\log p(y|\xi, \theta, D_t) - \log p(y|\xi, D_t)] = I(\xi; D_t)$$

Batch BALD

Stepwise
 Uncertainty Reduction

$$\alpha_{\text{SUR}}(\xi; D_t) = - \mathbb{E}_{p(y|\xi, D_t)} [\mathcal{H}[p(\theta|D_t \cup \{(\xi, y)\})]]$$

$\curvearrowright p(\theta|\xi, y, D_t)$

$$= I(\xi; D_t) + \mathbb{E}_{p(\theta|D_t)} [\log p(\theta|D_t)] = I(\xi; D_t) - \mathcal{H}[p(\theta|D_t)]$$

constant w.r.t. ξ

Entropy Search

$$\alpha_{\text{ES}}(\xi; D_t) = \mathbb{E}_{p(y|\xi, D_t)} [KL[p(\theta|D_t \cup \{(\xi, y)\}) \| b(\theta)]]$$

$$= \mathbb{E}_{p(y|\xi, D_t)} p(\theta|\xi, y, D_t) [\log p(\theta|\xi, y, D_t) - \log p(\theta) + \log p(\theta) - \log b(\theta)]$$

$$= I(\xi; D_t) + KL[p(\theta) \| b(\theta)]$$

constant w.r.t. ξ

Other acquisition functions

entropy acquisition $\alpha(\xi; D_t) = \mathcal{H}[P(y|\xi, D_t)]$ $\mathbb{E}_{P(\theta|D_t)}[P(y|\theta, \xi)]$

mean-STD acquisition $\alpha(\xi; D_t) = \frac{1}{|y|} \sum_y \sqrt{\text{Var}_{P(\theta|D_t)}[P(y|\xi, D_t)]}$

Prob of improvement $\alpha_{PI}(\xi; D_t) = \mathbb{P}(f(\xi) < \tau | D_t)$ $\tau = \max\{y_1, \dots, y_n\}$

Expected improvement $\alpha_{EI}(\xi; D_t) = \mathbb{E}[(f(\xi) - \tau)_+ | D_t]$

UCB $\alpha_{UCB-p}(\xi; D_t) = \overset{\rightarrow p\text{-quantile}}{q_p(f(\xi) | D_t)} = \mu(\xi | D_t) + \overset{\rightarrow p\text{-quantile of } N(0,1)}{\beta_p \cdot \sigma(\xi | D_t)}$

Thompson Sampling $\alpha_{TS}(\xi; D_t) = f_t(\xi)$ $f_t \sim \overset{\rightarrow GP}{p(f | D_t)}$

history

$$h_{t-1} = \{(s_i, y_i)\}_{i=1:t-1}$$

$$p(y_t | s) = \mathbb{E}_{p(\theta | h_{t-1})} [p(y_t | \theta, s)]$$

$$\begin{aligned} I_{h_{t-1}}(s) &= \mathbb{E}_{p(\theta | h_{t-1})} p(y_t | \theta, s, h_{t-1}) [\log p(y_t | \theta, s, h_{t-1}) - \log p(y_t | s, h_{t-1})] \\ &= \dots [\log p(\theta | h_{t-1}, s, y_t) - \log p(\theta | h_{t-1})] \end{aligned}$$

batch/
static design

determine all $s_{1:T}$ before the first iteration \approx one-step w/ larger design space
treating whole sequence as one experiment

variational
BOED

BA bound:

$$\hat{\mu}_{\text{posterior}}(\xi) = \mathbb{E}_{p(y, \theta | \xi)} [\log q_p(\theta | y, \xi) - \log p(\theta)] \approx \frac{1}{N} \sum_{n=1}^N \log \frac{q_p(\theta_n | y, \xi)}{p(\theta_n)} \leq I(\xi)$$

$$\hat{\mu}_{\text{marginal}}(\xi) = \mathbb{E}_{p(y, \theta | \xi)} [\log p(y | \theta, \xi) - \log q_m(y | \xi)] \approx \frac{1}{N} \sum_{n=1}^N \log \frac{p(y_n | \theta_n, \xi)}{q_m(y_n | \xi)} \geq I(\xi)$$

$$\hat{\mu}_{\text{VMC}}(\xi) = \mathbb{E}_{p(y, \theta | \xi)} \underbrace{q_w(\theta_{1:L} | s, y)}_{\text{or } q_w(\theta_{1:L} | y)} \left[\log \frac{p(y | \theta, \xi)}{\frac{1}{L} \sum_{l=1}^L \frac{p(y, \theta_l | s)}{q_w(\theta_l | y, \xi)}} \right] \stackrel{L \rightarrow \infty}{\text{IWAE bound, consistent}} \geq I(\xi)$$

$$\leq \log p(y | s) = \log \frac{1}{L} \sum_{l=1}^L \mathbb{E}_{q_w(\theta_{1:L} | s, y)} \left[\frac{p(y, \theta_l | s)}{q_w(\theta_l | y, s)} \right] \geq \mathbb{E}_{q_w(\theta_{1:L} | s, y)} \left[\log \frac{1}{L} \sum_{l=1}^L \frac{p(y, \theta_l | s)}{q_w(\theta_l | y, s)} \right]$$

$$I_{\text{ACE}}(\xi, L) = \mathbb{E}_{p(y, \theta | \xi)} q_w(\theta_{1:L} | s, y) \left[\log \frac{p(y | \theta, \xi)}{\frac{1}{L+1} \sum_{l=0}^L \frac{p(y, \theta_l | s)}{q_w(\theta_l | y, s)}} \right] \leq I(\xi)$$

by $q_w(\theta | y, \xi) \rightarrow p(\theta)$

$$I_{\text{DCE}}(\xi, L) = \mathbb{E}_{p(y, \theta | \xi)} p(\theta_{1:L}) \left[\log \frac{p(y | \theta, \xi)}{\frac{1}{L+1} \sum_{l=0}^L p(y | \theta_l, \xi)} \right] \leq I(\xi) \quad \text{InfoNCE}$$

(implicit likelihood)

$$\hat{\mu}_{m+l}(\xi) = \mathbb{E}_{p(y, \theta | \xi)} \left[\log \frac{q_l(y | \theta, \xi)}{q_m(y | \xi)} \right] \approx \frac{1}{N} \sum_{n=1}^N \log \frac{q_l(y_n | \theta_n, \xi)}{q_m(y_n | \xi)} \quad \text{not a bound of } I(\xi)$$

+ sequential

$$p(\theta) \rightarrow p(\theta | h_{t-1}) = p(\theta) \prod_{i=1}^{t-1} p(y_i | \theta, \xi_i) / \prod_{i=1}^{t-1} p(y_i | \xi_i)$$

$$\Rightarrow \hat{\mu}_{\text{post}}(\xi_t) = \mathbb{E}_{p(\theta | h_{t-1})} p(y_t | \theta, \xi_t) \left[\log q_p(\theta | y_t, \xi_t) - \log p(\theta) \prod_{i=1}^{t-1} p(y_i | \theta, \xi_i) \right] + \log p(y_{1:t-1} | s_{1:t-1})$$

constant for ξ_t

(DAD)

Deep Adaptive
Design

design function: $\xi_t = \pi_\phi^{\text{NN}}(h_{t-1})$

MI chain rule

$$\begin{aligned} \text{MI}(\theta; h_T) &= \mathcal{H}[p(\theta)] - \mathcal{H}[p(\theta|h_T)] = \sum_t \mathbb{E}[\mathcal{H}[p(\theta|h_{t-1})] - \mathcal{H}[p(\theta|h_t)]] \\ &= \sum_t \text{MI}(\theta; (y_t, s_t) | h_{t-1}) \end{aligned}$$

explicit likelihood

$$\begin{aligned} \mathcal{I}_T(\pi_\phi) &= \mathbb{E}_{p(\theta)} p(h_T | \theta, \pi_\phi) \left[\sum_{t=1}^T \mathcal{I}_{h_{t-1}}(\xi_t) \right] & p(h_T | \theta, \pi_\phi) &= \prod_t p(y_t | \theta, \xi_t) \\ &= \mathbb{E}_{p(\theta)} p(h_T | \theta, \pi_\phi) \left[\log \frac{p(h_T | \theta, \pi_\phi)}{p(h_T | \pi_\phi)} \right] & p(h_T | \pi_\phi) &= \mathbb{E}_{p(\theta)} [p(h_T | \theta, \pi_\phi)] \\ &\geq \mathbb{E}_{p(\theta_0)} p(h_T | \theta_0, \pi_\phi) p(\theta_{1:L}) \left[\log \frac{p(h_T | \theta_0, \pi_\phi)}{\prod_{t=1}^T \int_{\theta_0}^{\xi_t} p(h_T | \theta_t, \pi_\phi)} \right] =: \mathcal{L}_T^{\text{PCE}}(\pi_\phi, L) \end{aligned}$$

$$\begin{aligned} \partial_\phi \mathcal{I} &= \mathbb{E}_{p(\theta_0)} p(h_T | \theta_0, \pi_\phi) p(\theta_{1:L}) \left[\log \frac{p(h_T | \theta_0, \pi_\phi)}{\prod_{t=1}^T \int_{\theta_0}^{\xi_t} p(h_T | \theta_t, \pi_\phi)} \cdot \nabla_\phi \log p(h_T | \theta_0, \pi_\phi) \right. \\ &\quad \left. + \nabla_\phi \left(\log \frac{p(h_T | \theta_0, \pi_\phi)}{\prod_{t=1}^T \int_{\theta_0}^{\xi_t} p(h_T | \theta_t, \pi_\phi)} \right) \right] \\ &\quad \mathbb{E}[\nabla_\phi \log p(h_T | \theta_0, \pi_\phi)] = 0 \end{aligned}$$