

Expected information gain in marginals

Consider the **marginal posterior** and **marginal prior**

$$p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d}) = \int_H \underbrace{p(\boldsymbol{\theta}, \boldsymbol{\eta}|\mathbf{y}, \mathbf{d})}_{\text{joint posterior}} d\boldsymbol{\eta}, \quad p(\boldsymbol{\theta}) = \int_H \underbrace{p(\boldsymbol{\theta}, \boldsymbol{\eta})}_{\text{joint prior}} d\boldsymbol{\eta}$$

Information gain in $\boldsymbol{\theta}$ from a single observation \mathbf{y} :

$$u_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{d}) = D_{\text{KL}} \left[\overbrace{p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d})}^{\text{marginal posterior}} \parallel \overbrace{p(\boldsymbol{\theta})}^{\text{marginal prior}} \right]$$

Expected information gain – experimental outcome is unknown beforehand:

$$U_{\boldsymbol{\theta}}(\mathbf{d}) = \mathbb{E}_{\mathbf{y}}[u_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{d})] = \int_{\mathbf{Y}} u_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{d}) p(\mathbf{y}|\mathbf{d}) d\mathbf{y}$$

Expectation is taken over the prior predictive of the data

Expected information gain in marginals

Consider the **marginal posterior** and **marginal prior**

$$p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d}) = \int_H \underbrace{p(\boldsymbol{\theta}, \boldsymbol{\eta}|\mathbf{y}, \mathbf{d})}_{\text{joint posterior}} d\boldsymbol{\eta}, \quad p(\boldsymbol{\theta}) = \int_H \underbrace{p(\boldsymbol{\theta}, \boldsymbol{\eta})}_{\text{joint prior}} d\boldsymbol{\eta}$$

Information gain in $\boldsymbol{\theta}$ from a single observation \mathbf{y} :

$$u_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{d}) = D_{\text{KL}} \left[\overbrace{p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d})}^{\text{marginal posterior}} \parallel \overbrace{p(\boldsymbol{\theta})}^{\text{marginal prior}} \right]$$

Expected information gain – experimental outcome is unknown beforehand:

$$U_{\boldsymbol{\theta}}(\mathbf{d}) = \mathbb{E}_{\mathbf{y}}[u_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{d})] = \int_{\mathbf{Y}} u_{\boldsymbol{\theta}}(\mathbf{y}, \mathbf{d}) p(\mathbf{y}|\mathbf{d}) d\mathbf{y}$$

Expectation is taken over the prior predictive of the data

Estimating the expected information gain

Challenge: No closed-form expression for expected information gain

- ▶ Nonlinear computational models
- ▶ Nontrivial priors and noise models
- ▶ Avoid large approximations (mean-field, Gaussian, etc.)

Solution: use a Monte Carlo estimate

$$\begin{aligned}U_{\theta}(d) &= \mathbb{E}_y \left[D_{\text{KL}} \left[p(\theta|y, d) \parallel p(\theta) \right] \right] \\&= \int_{\mathcal{Y}} \int_{\Theta} \log \left[\frac{p(\theta|y, d)}{p(\theta)} \right] p(\theta|y, d) d\theta p(y|d) dy \\&= \int_{\mathcal{Y}} \int_{\Theta} \left\{ \underbrace{\ln [p(y|\theta, d)]}_{\text{conditional likelihood}} - \ln [p(y|d)] \right\} p(y|\theta, d) p(\theta) d\theta dy\end{aligned}$$

Estimating the expected information gain

Challenge: No closed-form expression for expected information gain

- ▶ Nonlinear computational models
- ▶ Nontrivial priors and noise models
- ▶ Avoid large approximations (mean-field, Gaussian, etc.)

Solution: use a Monte Carlo estimate

$$\begin{aligned}U_{\theta}(d) &= \mathbb{E}_{\mathbf{y}} \left[D_{\text{KL}} [p(\theta|\mathbf{y}, d) \parallel p(\theta)] \right] \\&= \int_{\mathbf{Y}} \int_{\Theta} \log \left[\frac{p(\theta|\mathbf{y}, d)}{p(\theta)} \right] p(\theta|\mathbf{y}, d) d\theta p(\mathbf{y}|d) d\mathbf{y} \\&= \int_{\mathbf{Y}} \int_{\Theta} \left\{ \underbrace{\ln [p(\mathbf{y}|\theta, d)]}_{\text{conditional likelihood}} - \underbrace{\ln [p(\mathbf{y}|d)]}_{\text{marginal likelihood}} \right\} p(\mathbf{y}|\theta, d) p(\theta) d\theta d\mathbf{y}\end{aligned}$$

Estimating the expected information gain

Challenge: No closed-form expression for expected information gain

- ▶ Nonlinear computational models
- ▶ Nontrivial priors and noise models
- ▶ Avoid large approximations (mean-field, Gaussian, etc.)

Solution: use a Monte Carlo estimate

$$\begin{aligned}U_{\theta}(\mathbf{d}) &= \mathbb{E}_{\mathbf{y}} \left[D_{\text{KL}} [p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d}) \parallel p(\boldsymbol{\theta})] \right] \\&= \int_{\mathbf{Y}} \int_{\boldsymbol{\theta}} \log \left[\frac{p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d})}{p(\boldsymbol{\theta})} \right] p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{d}) d\boldsymbol{\theta} p(\mathbf{y}|\mathbf{d}) d\mathbf{y} \\&= \int_{\mathbf{Y}} \int_{\boldsymbol{\theta}} \left\{ \underbrace{\ln [p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{d})]}_{\text{conditional likelihood}} - \underbrace{\ln [p(\mathbf{y}|\mathbf{d})]}_{\text{marginal likelihood}} \right\} p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{d}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{y}\end{aligned}$$

Estimating the expected information gain

Challenge: No closed-form expression for expected information gain

- ▶ Nonlinear computational models
- ▶ Nontrivial priors and noise models
- ▶ Avoid large approximations (mean-field, Gaussian, etc.)

Solution: use a Monte Carlo estimate

$$\begin{aligned}U_{\theta}(\mathbf{d}) &= \mathbb{E}_{\mathbf{y}} \left[D_{\text{KL}} \left[p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{d}) \parallel p(\boldsymbol{\theta}) \right] \right] \\&= \int_{\mathbf{Y}} \int_{\boldsymbol{\theta}} \log \left[\frac{p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{d})}{p(\boldsymbol{\theta})} \right] p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{d}) d\boldsymbol{\theta} p(\mathbf{y} | \mathbf{d}) d\mathbf{y} \\&\approx \frac{1}{N} \sum_{i=1}^N \left\{ \underbrace{\ln [p(\mathbf{y}^{(i)} | \boldsymbol{\theta}^{(i)}, \mathbf{d})]}_{\text{conditional likelihood}} - \underbrace{\ln [p(\mathbf{y}^{(i)} | \mathbf{d})]}_{\text{marginal likelihood}} \right\}, \quad \begin{aligned} \boldsymbol{\theta}^{(i)} &\sim p(\boldsymbol{\theta}) \\ \mathbf{y}^{(i)} &\sim p(\mathbf{y} | \boldsymbol{\theta}^{(i)}, \mathbf{d}) \end{aligned}\end{aligned}$$

Estimating the expected information gain (importance sampling)

$$U_{\theta}(\mathbf{d}) \approx \frac{1}{N} \sum_{i=1}^N \left\{ \ln \underbrace{[p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)}, \mathbf{d})]}_{\text{conditional likelihood}} - \ln \underbrace{[p(\mathbf{y}^{(i)}|\mathbf{d})]}_{\text{marginal likelihood}} \right\}, \quad \begin{aligned} \boldsymbol{\theta}^{(i)} &\sim p(\boldsymbol{\theta}) \\ \mathbf{y}^{(i)} &\sim p(\mathbf{y}|\boldsymbol{\theta}^{(i)}, \mathbf{d}) \end{aligned}$$

Importance sampling estimates for the conditional and marginal likelihood:

$$\underbrace{p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)}, \mathbf{d})}_{\text{conditional likelihood}} = \int_H p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)}, \boldsymbol{\eta}, \mathbf{d}) p(\boldsymbol{\eta}|\boldsymbol{\theta}^{(i)}, \mathbf{d}) d\boldsymbol{\eta}$$

$$\underbrace{p(\mathbf{y}^{(i)}|\mathbf{d})}_{\text{marginal likelihood}} = \int_{\boldsymbol{\theta}} \int_H p(\mathbf{y}^{(i)}|\boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{d}) p(\boldsymbol{\theta}, \boldsymbol{\eta}|\mathbf{d}) d\boldsymbol{\eta} d\boldsymbol{\theta}$$

Biasing distributions

- ▶ $\boldsymbol{\eta}^{(i,k)} \sim q_{\text{cond}}^{(i)}$
- ▶ $\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)} \sim q_{\text{marg}}^{(i)}$

Biasing weights

- ▶ $w_{\text{cond}}^{(i,k)} = p(\boldsymbol{\eta}^{(i,k)}|\boldsymbol{\theta}^{(i)}) / q_{\text{cond}}^{(i)}(\boldsymbol{\eta}^{(i,k)})$
- ▶ $w_{\text{marg}}^{(i,j)} = p(\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)}) / q_{\text{marg}}^{(i)}(\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)})$

Estimating the expected information gain (importance sampling)

$$U_{\theta}(\mathbf{d}) \approx \frac{1}{N} \sum_{i=1}^N \left\{ \ln \underbrace{[p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)}, \mathbf{d})]}_{\text{conditional likelihood}} - \ln \underbrace{[p(\mathbf{y}^{(i)}|\mathbf{d})]}_{\text{marginal likelihood}} \right\}, \quad \begin{aligned} \boldsymbol{\theta}^{(i)} &\sim p(\boldsymbol{\theta}) \\ \mathbf{y}^{(i)} &\sim p(\mathbf{y}|\boldsymbol{\theta}^{(i)}, \mathbf{d}) \end{aligned}$$

Importance sampling estimates for the conditional and marginal likelihood:

$$\underbrace{p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)}, \mathbf{d})}_{\text{conditional likelihood}} \approx \frac{1}{M_2} \sum_{k=1}^{M_2} p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)}, \boldsymbol{\eta}^{(i,k)}, \mathbf{d}) w_{\text{cond}}^{(i,k)}$$

$$\underbrace{p(\mathbf{y}^{(i)}|\mathbf{d})}_{\text{marginal likelihood}} \approx \frac{1}{M_1} \sum_{j=1}^{M_1} p(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)}, \mathbf{d}) w_{\text{marg}}^{(i,j)}$$

Biasing distributions

- ▶ $\boldsymbol{\eta}^{(i,k)} \sim q_{\text{cond}}^{(i)}$
- ▶ $\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)} \sim q_{\text{marg}}^{(i)}$

Biasing weights

- ▶ $w_{\text{cond}}^{(i,k)} = p(\boldsymbol{\eta}^{(i,k)}|\boldsymbol{\theta}^{(i)})/q_{\text{cond}}^{(i)}(\boldsymbol{\eta}^{(i,k)})$
- ▶ $w_{\text{marg}}^{(i,j)} = p(\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)})/q_{\text{marg}}^{(i)}(\boldsymbol{\theta}^{(i,j)}, \boldsymbol{\eta}^{(i,j)})$

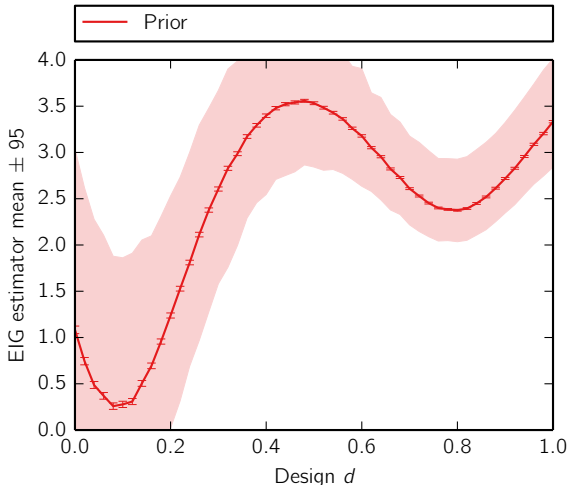
The computational problem in a nutshell:

- ▶ We must consider many realizations of the data $\mathbf{y}^{(i)}$
- ▶ For each realization, we have a posterior distribution. Compute:
 - ▶ Complete posterior normalizing constant (*marginal likelihood*)
 - ▶ Partial posterior normalizing constant (*conditional likelihood*)
- ▶ Hence, *nested* Monte Carlo sampling
- ▶ Naïve approach: use prior as a biasing distribution [Ryan 2003]
 - ▶ Can also use surrogates [Huan & M 2013] to allow very large sample sizes, but surrogate leaves an asymptotic bias

The computational problem in a nutshell:

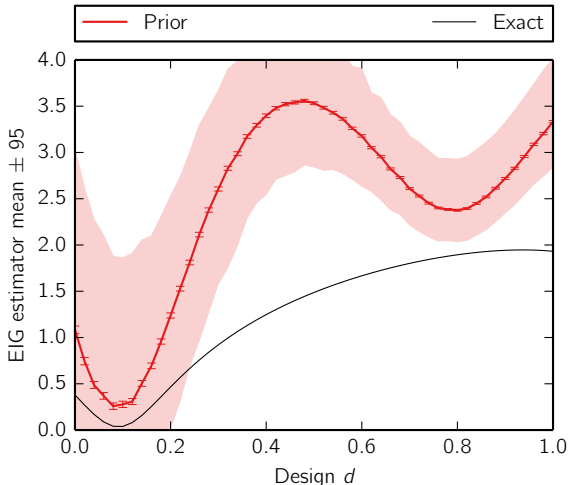
- ▶ We must consider many realizations of the data $\mathbf{y}^{(i)}$
- ▶ For each realization, we have a posterior distribution. Compute:
 - ▶ Complete posterior normalizing constant (*marginal likelihood*)
 - ▶ Partial posterior normalizing constant (*conditional likelihood*)
- ▶ Hence, *nested* Monte Carlo sampling
- ▶ Naïve approach: use prior as a biasing distribution [Ryan 2003]
 - ▶ Can also use surrogates [Huan & M 2013] to allow very large sample sizes, but surrogate leaves an asymptotic bias
- ▶ **How bad could prior sampling be?**

Sampling distribution of the naïve estimator



- ▶ $N = M_1 = M_2 = 316 \rightarrow 10^5$ samples at each design
- ▶ Bias is particularly misleading for the case of **focused** design...

Sampling distribution of the naïve estimator



- ▶ $N = M_1 = M_2 = 316 \rightarrow 10^5$ samples at each design
- ▶ Bias is particularly misleading for the case of **focused** design...

Estimator bias and variance

- ▶ Use Δ -method to analyze bias and variance of EIG estimator

$$\mathbb{E}[\hat{U}(\mathbf{d})] \approx U(\mathbf{d}) + \overbrace{\frac{A(\mathbf{d})}{M_1} - \frac{B(\mathbf{d})}{M_2}}^{\text{bias}}$$

$$\mathbb{V}[\hat{U}(\mathbf{d})] \approx \frac{C(\mathbf{d})}{N} + \frac{D(\mathbf{d})}{NM_1} + \frac{E(\mathbf{d})}{NM_2}$$

- ▶ $A(\mathbf{d}), B(\mathbf{d}), D(\mathbf{d}), E(\mathbf{d}) \sim \mathbb{E}_{\mathbf{y}}[\text{variance of inner estimators}]$
- ▶ Need to find biasing distributions $q_{\text{marg}}^{(i)}$ and $q_{\text{cond}}^{(i)}$ to reduce variance

Estimator bias and variance

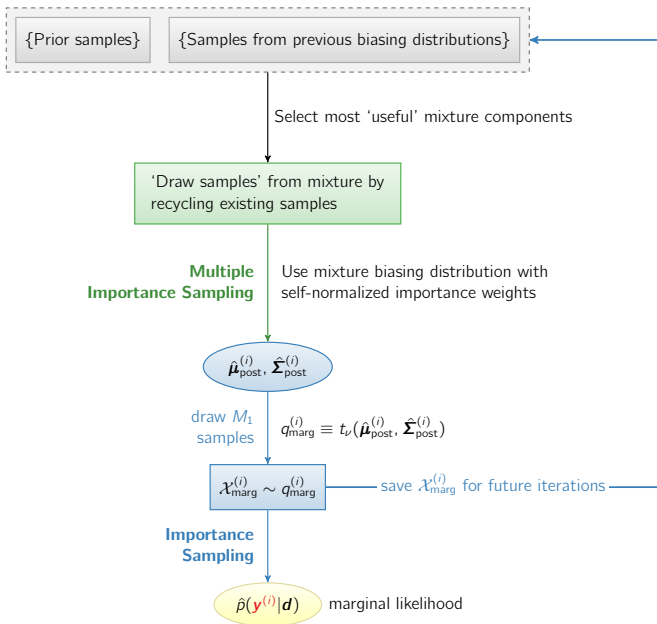
- ▶ Use Δ -method to analyze bias and variance of EIG estimator

$$\mathbb{E}[\hat{U}(\mathbf{d})] \approx U(\mathbf{d}) + \overbrace{\frac{A(\mathbf{d})}{M_1} - \frac{B(\mathbf{d})}{M_2}}^{\text{bias}}$$

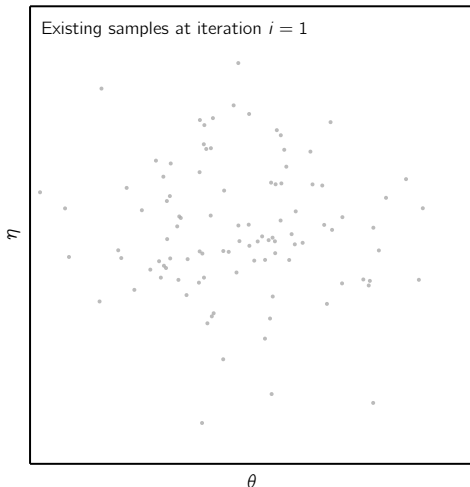
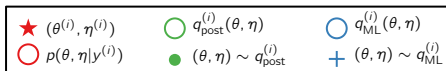
$$\mathbb{V}[\hat{U}(\mathbf{d})] \approx \frac{C(\mathbf{d})}{N} + \frac{D(\mathbf{d})}{NM_1} + \frac{E(\mathbf{d})}{NM_2}$$

- ▶ $A(\mathbf{d}), B(\mathbf{d}), D(\mathbf{d}), E(\mathbf{d}) \sim \mathbb{E}_{\mathbf{y}}[\text{variance of inner estimators}]$
- ▶ Need to find biasing distributions $q_{\text{marg}}^{(i)}$ and $q_{\text{cond}}^{(i)}$ to reduce variance
- ▶ **Idea:**
 - ▶ An individual normalizing constant is difficult to estimate. Collectively, the problem is easier!
 - ▶ Recycle existing samples/model evaluations to obtain increasingly better biasing distributions, as the outer loop proceeds
 - ▶ Preserve consistency of the estimator

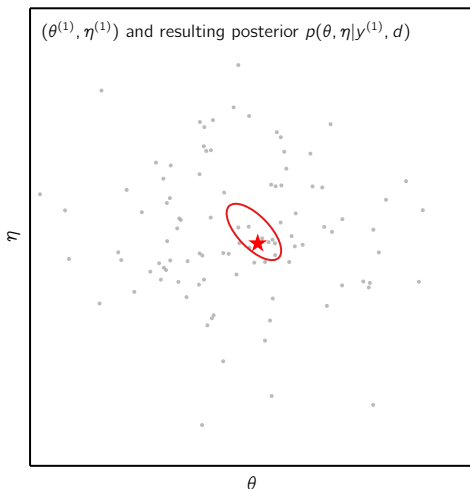
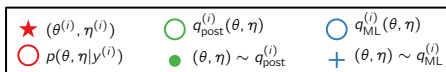
Layered multiple importance sampling



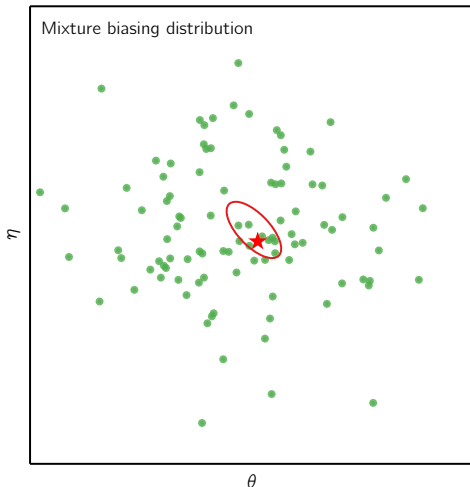
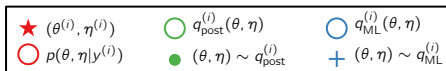
Layered multiple importance sampling



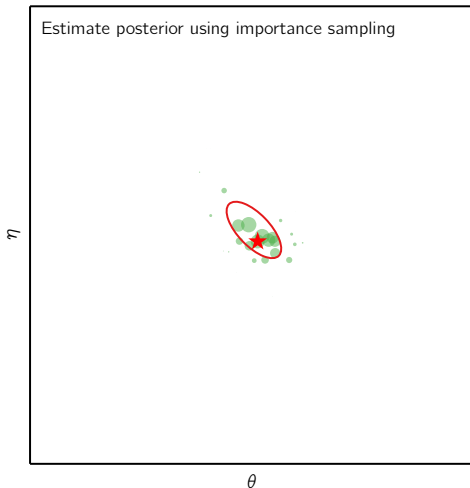
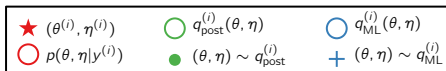
Layered multiple importance sampling



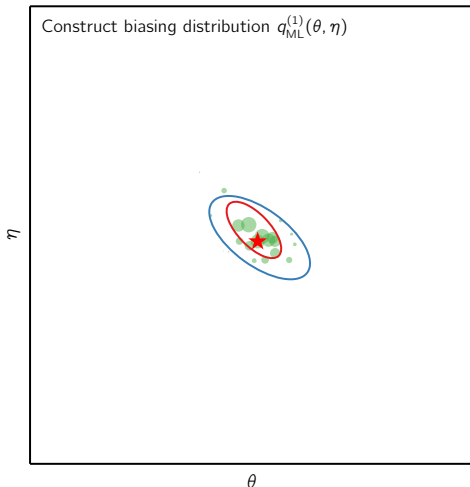
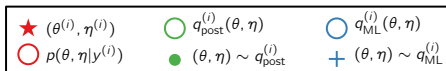
Layered multiple importance sampling



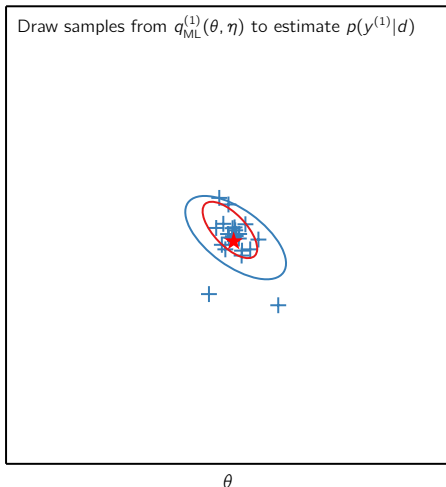
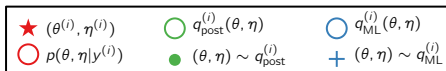
Layered multiple importance sampling



Layered multiple importance sampling



Layered multiple importance sampling



Layered multiple importance sampling

