

- Topics:
- Model Calibration
    - prior-predictive checks
    - posterior-predictive checks
  - Predictive Accuracy
    - log posterior predictive density (lppd)
    - cross validation
  - Approximate inference
    - ↳ what to do when we can't compute posteriors?
    - ↳ Parametric approximations
    - ↳ Monte Carlo methods

Assignments: Quiz 4 (credit for attendance)

↳ use HW 4 to help complete

Model calibration : To check if the model has fit well to the observed data

↳ we introduced the posterior-predictive distribution

$$\begin{aligned}
 p(\tilde{x} | x_{1:n}) &= \int p(\tilde{x}, \theta | x_{1:n}) d\theta \\
 &= \int p(\tilde{x} | \theta, x_{1:n}) \underbrace{p(\theta | x_{1:n})}_{\text{posterior}} d\theta \\
 &= \int p(\tilde{x} | \theta) p(\theta | x_{1:n}) d\theta
 \end{aligned}$$

new samples (pointing to  $\tilde{x}$ )  
 observed historical data (under  $x_{1:n}$ )  
 (chain rule) (under the first equality)

$\tilde{x} \perp\!\!\!\perp x_{1:n}$  given  $\theta$

$$= \mathbb{E} \left[ p(\tilde{x} | \theta) \mid X_{1:n} \right]$$

Point estimation:  $\hat{\theta} = g(X_1, \dots, X_n)$   
 observe  $x_{1:n}$

$$\hat{\theta} = g(x_1, \dots, x_n)$$

↓ plug in

$$x \sim p(x | \hat{\theta})$$

only one realization of the parameter

Bayesian estimation

$$x \sim \mathbb{E} \left[ p(x | \theta) \mid X_{1:n} \right]$$

$$\theta^{(1)}, \dots, \theta^{(s)} \text{ iid } \sim p(\theta | X_{1:n})$$

$$\approx \frac{1}{S} \sum_{s=1}^S p(x | \theta^{(s)})$$

more than 1 realization of the parameter

- ① Prior - predictive check
- ② Posterior - predictive check

used as means to verify that our Bayesian model can simulate the data

Prior-predictive: 
$$p(\tilde{x}) = \int p(\tilde{x}, \theta) d\theta$$

$$= \int p(\tilde{x} | \theta) p(\theta) d\theta$$

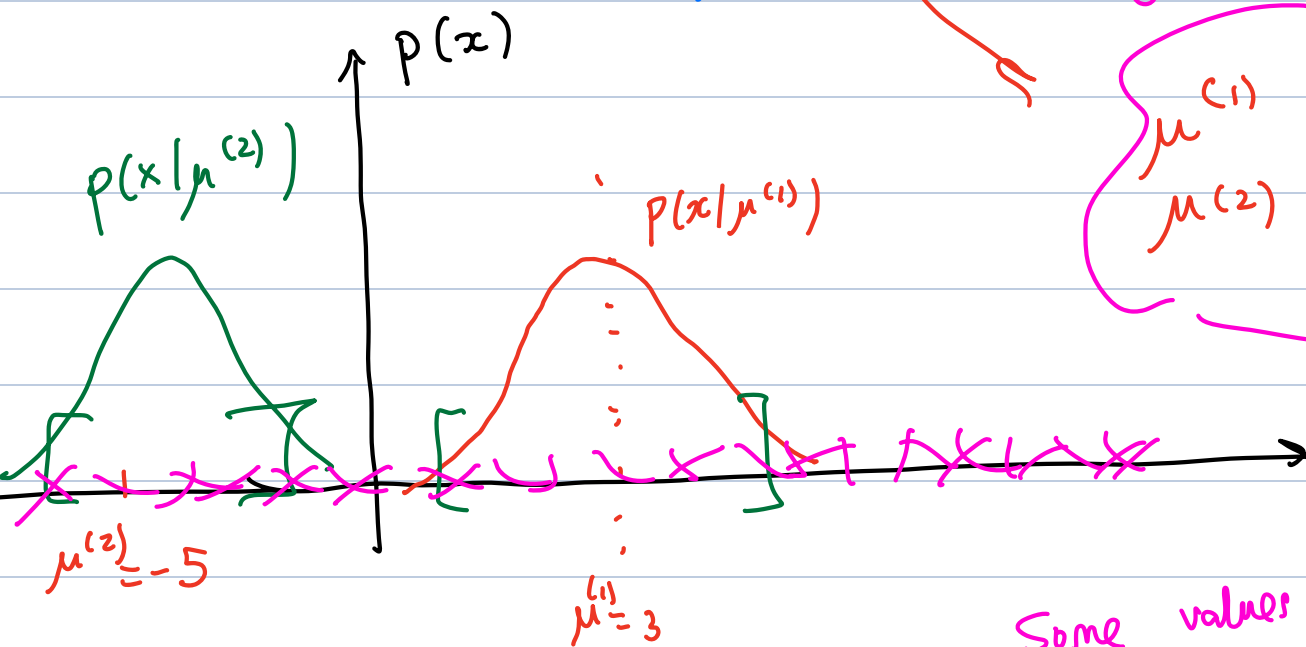
$$= \mathbb{E} [ p(\tilde{x} | \theta) ]$$

$$x \sim \mathcal{N}(\mu, 1)$$

$$\mu \sim \mathcal{N}(5, 100)$$

*sd = 10*

This expectation is taken w.r.t. the prior dist. of  $\theta$



$$\mu^{(1)} = 3$$

$$\mu^{(2)} = -5$$

Some values of  $\mu$  won't be

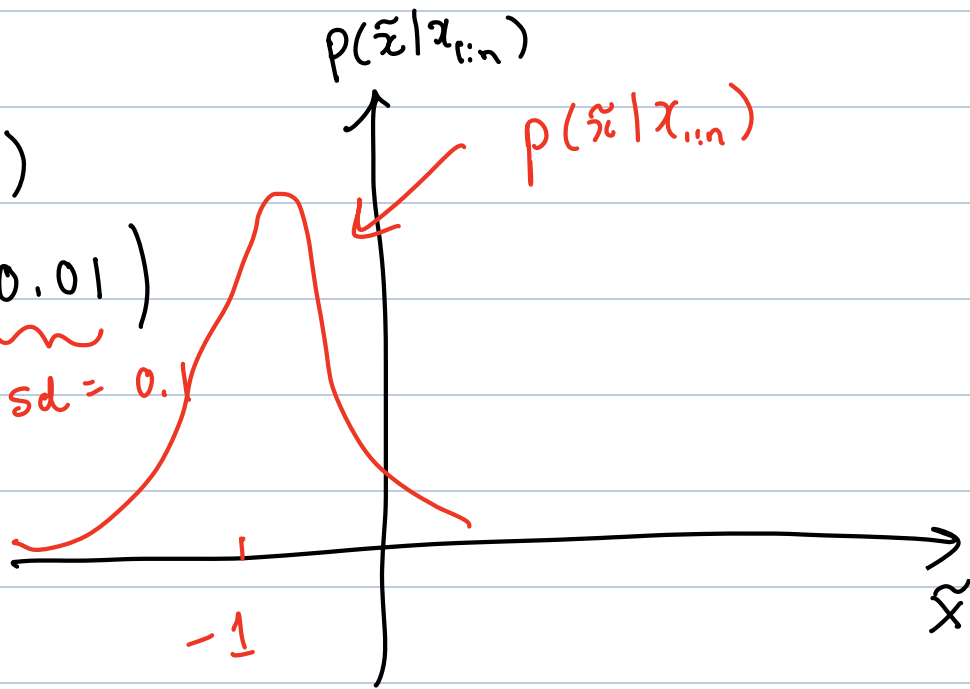
sampled with high probability

## Posterior-predictive check

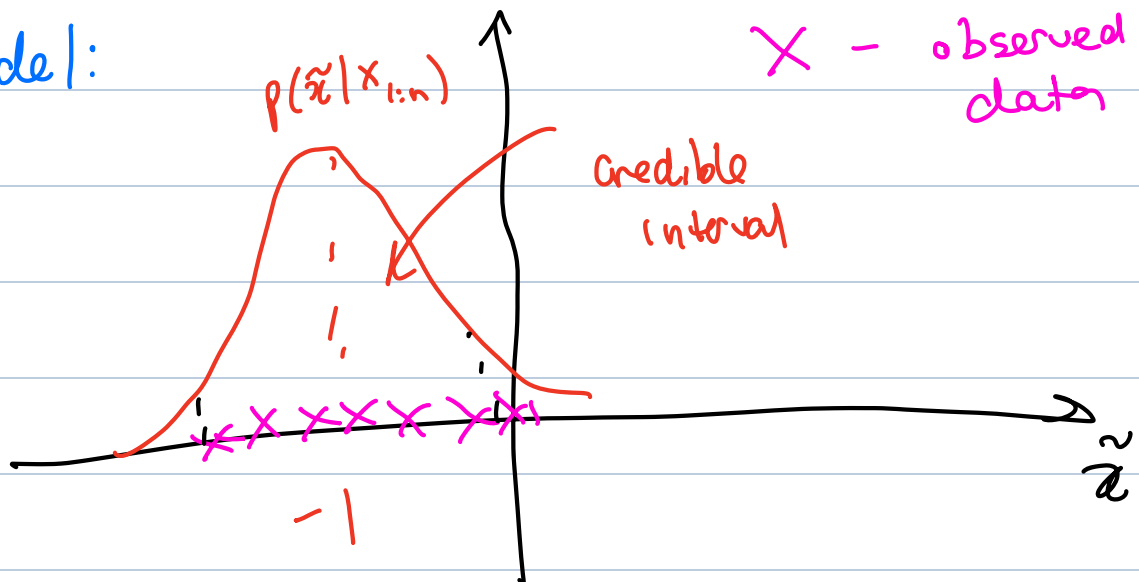
- A way for us to measure how well a model is calibrated after we have done inference

$$x \sim \mathcal{N}(\mu, 1)$$
$$\mu | x_{1:n} \sim \mathcal{N}(-1, 0.01)$$

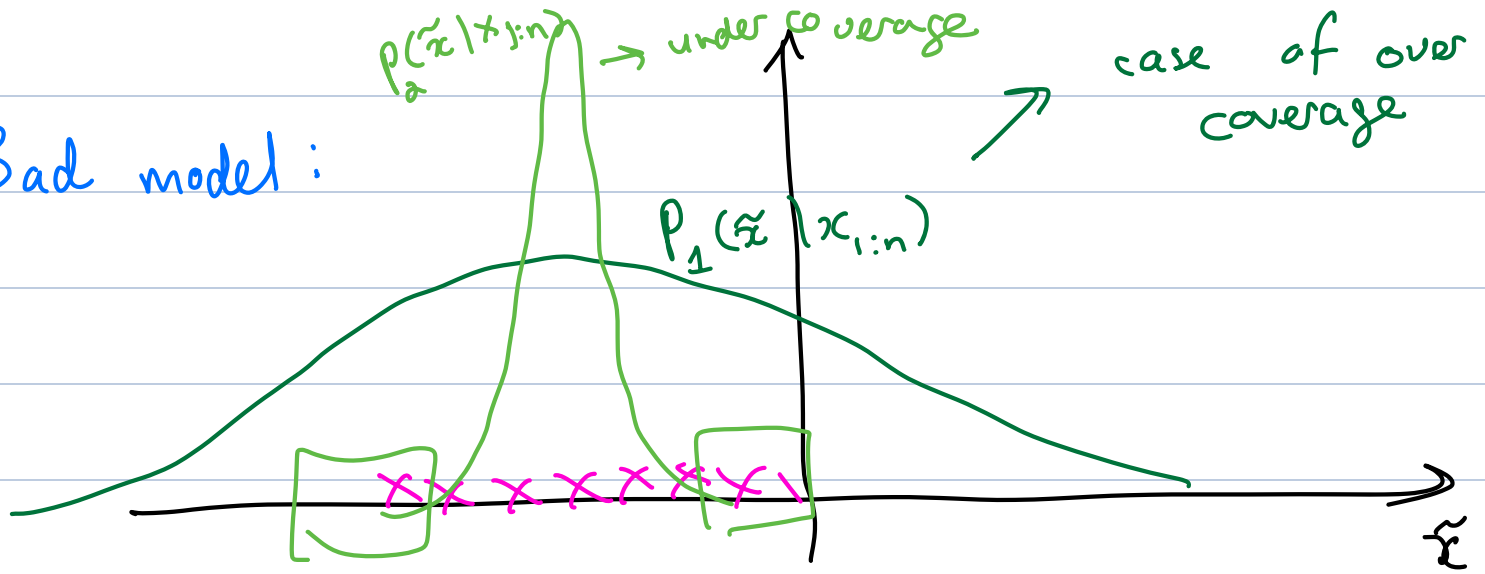
sd = 0.1



Good model:



Bad model:



Credible interval:

↳ a metric/estimator Bayesian statisticians use to assert the probability of something happening in a Bayesian model

$$p(\theta | x_{1:n})$$

Where does the 95% of the probability mass of the posterior sit?

$$I = [a - \mu, a + \mu] \quad \text{s.t.}$$

$$IP(\theta \in I | x_{1:n}) = 0.95$$

95%  
CI

→ Ways to get this:

→ Approximate posterior with a Gaussian and use the Z-scores

→ Simulate samples and empirically evaluate the credible intervals

## Predictive Accuracy

- Before, we were using the posterior-predictive distribution  $p(\tilde{x} | x_{1:n})$  to check the calibration of our model
- We may also want to see how well our model generalizes to **unseen** observations

Metric: log posterior predictive density (lppd)

↳ Observe new data  $\tilde{x}_*$

↳ Evaluate the predictive distribution

$$L_{ppd} = \log p(\tilde{x}_* | X_{1:n})$$

one  $\tilde{x}_*$

$$= \frac{1}{S} \sum_{s=1}^S \log p(\tilde{x}_* | \theta^{(s)})$$

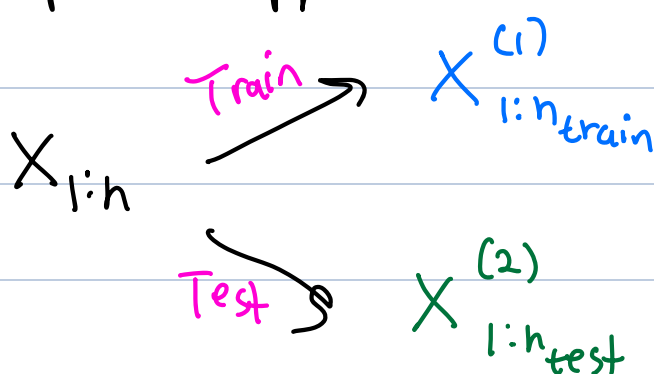
$\theta^{(s)} \sim p(\theta | X_{1:n})$

→ the higher this metric, the better the performance of the model

→ evaluate this always on out-of-sample data

Question: In practice, we don't want to assess this metric after observing new data. We want to know it after we train the model!

① Simplest approach



} union of these datasets should be  $X_{1:n}$

↳ train model on  $X_{1:n_{\text{train}}}^{(1)}$

↳ evaluate  $L_{\text{ppd}}$  on  $X_{1:n_{\text{test}}}^{(2)}$

Problems with this approach:

→ depends on the size of the split

→ if  $n_{\text{train}}$  is small, then we don't have enough data to train the model

→ if  $n_{\text{test}}$  is small, the  $L_{\text{ppd}}$  estimator will have high variance

→ depends on the way we split the data

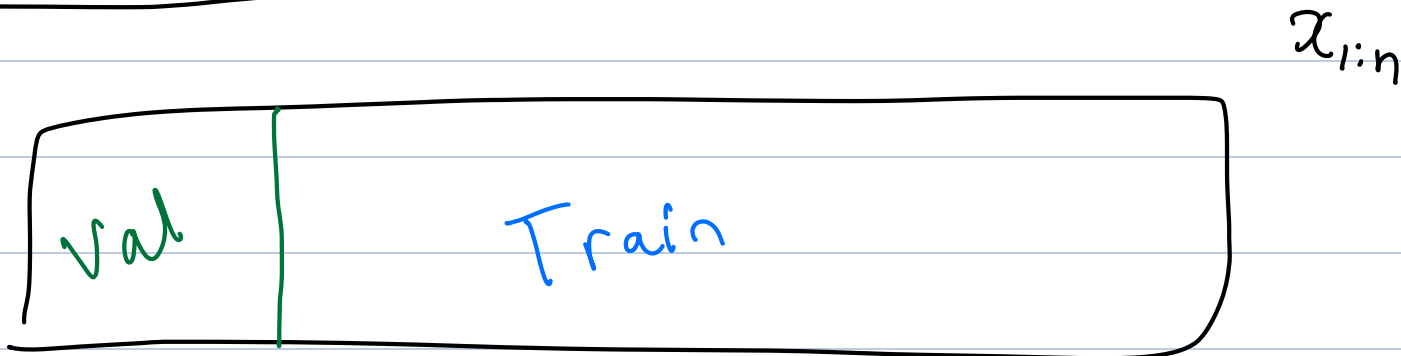
In practice, we evaluate the average performance

Average  $L_{\text{ppd}}$ :  $\mathbb{E}_{p(x)} [L_{\text{ppd}}]$

$\{x_1^{(2)}, \dots, x_{n_{\text{test}}}^{(2)}\}$  Small number of samples

Average lppd:  $\approx \frac{1}{n_{\text{test}}} \sum_{n=1}^{n_{\text{test}}} \text{lppd}(x_2^{(n)})$

Cross-validation



$p(\theta | \underbrace{x_{\text{train}}^{(i)}})$

Evaluate average lppd  $^{(i)} = \frac{1}{n_{\text{test}}^{(i)}} \sum_{n=1}^{n_{\text{test}}^{(i)}} \text{lppd}(x_{\text{val},n}^{(i)})$

↓ splitting and evaluation is repeated for K folds

Average lppd  $= \frac{1}{K} \sum_{k=1}^K \frac{1}{n_{\text{test}}^{(k)}} \sum_{n=1}^{n_{\text{test}}^{(k)}} \text{lppd}(x_{\text{val},n}^{(k)})$

↪ averaging the result over many possible train/validation splits

→ When  $n_{\text{test}}^{(k)} = 1$  for all  $k$ , this is called leave-one-out cross validation (LOOCV)

Bayesian

## Approximate<sup>^</sup> Inference

- We have seen that Bayesian methods are capable of:
  1. Incorporating domain knowledge into the inference procedure via the prior  $p(\theta)$
  2. Capture uncertainty about the parameter and about predictions using Bayesian inference

We have only dealt with models where the posterior is analytically tractable:

$$p(\theta | x_{1:n}) \propto p(x_{1:n} | \theta) p(\theta) \propto \tilde{\pi}(\theta)$$

The  $\tilde{\pi}(\theta)$  usually has had the form of a dist. we already know. For more general choices of  $q(\theta)$ , we cannot analytically determine the posterior.

↳ Restricted to only using certain priors if we want to obtain an analytical form to the posterior

Approximate Bayesian inference aims to obtain an approximate posterior (with nice properties) via the following two approaches:

### ① Parametric approximations

- choose a distribution that we know and estimate its parameters so that it is close enough to the true posterior
- Laplace approximation
- Variational inference

# Laplace Approximation

→ Assume the posterior is a Gaussian

$$p(\theta | x_{1:n}) \approx N(\theta | \hat{\mu}_\theta, \hat{\sigma}_\theta^2)$$

→ Want to choose  $\hat{\mu}_\theta$  and  $\hat{\sigma}_\theta^2$  such that

as  $n \rightarrow \infty$ , the approximation is perfect

→ Bernstein von-Mises theorem

$$\theta | X_{1:n} \xrightarrow{d} N\left(\underbrace{\hat{\theta}_{\text{MAP}}}_{\downarrow n \rightarrow \infty \hat{\theta}_{\text{MLE}}}, \underbrace{I_{\text{Bayes}}^{-1}(\hat{\theta}_{\text{MAP}})}_{\downarrow n \rightarrow \infty 0}\right)$$

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta | X_{1:n})$$

Bayesian  
information

$$I_{\text{Bayes}}(\hat{\theta}_{\text{MAP}}) = -\mathbb{E} \left[ \frac{\partial^2 \log p(\theta, X_{1:n})}{\partial \theta^2} \bigg| \theta = \hat{\theta}_{\text{MAP}} \right]$$

$$I(\hat{\theta}_{\text{MAP}}) + J(\hat{\theta}_{\text{MAP}})$$

Fisher information

information  
of the prior

→ In practice, we evaluate the MAP and then approximate the Bayesian information to form the Laplace approximation

→ Since posteriors are asymptotically normal, we expect this to be a good approximation, when  $n$  is large.

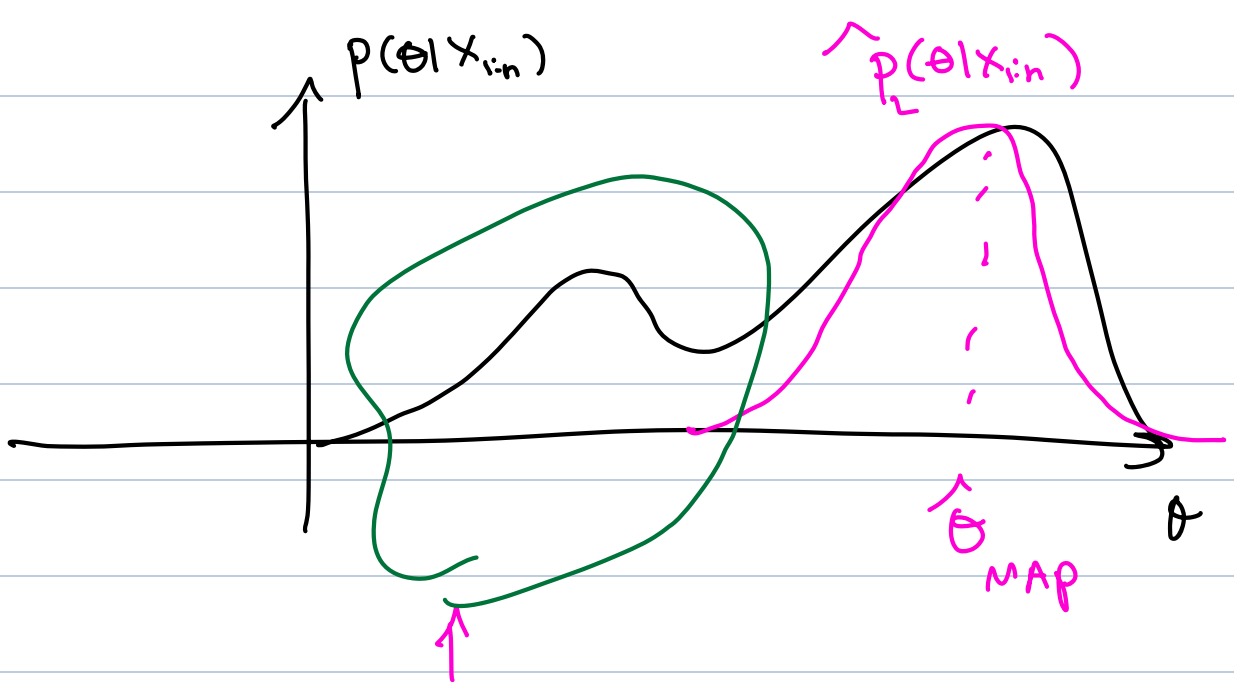
### Limitations

(1.) No convergence guarantees for derived estimators:

$$\int \theta p(\theta | X_{1:n}) d\theta \neq \mathbb{E}[\theta | X_{1:n}]$$

replaced with Laplace approximation

(2.) Approximation is bad when  $n$  is small and the likelihood is non-Gaussian



will not capture the other modes in the distribution

→ We seek a more theoretically rigorous approach to obtaining the posterior distribution  
 → Monte Carlo methods!

## ② Monte Carlo (MC) methods

- Let  $\pi(\theta) \triangleq p(\theta | X_{1:n})$  denote the distribution of interest.
- MC methods aim to approximate expectations taken w.r.t. to the target density  $\pi(\theta)$

- Let's say we are interested in computing:

$$\mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)] = \int h(\theta) \pi(\theta) d\theta$$

- Standard Monte Carlo:

- Simulate  $M$  samples from  $\pi(\theta)$ :

$$\theta^{(m)} \stackrel{\text{i.i.d.}}{\sim} \pi(\theta)$$

- Approximate the integral using a sample average:

$$\mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)] \approx \frac{1}{M} \sum_{m=1}^M h(\theta^{(m)})$$

$$\hat{I}_{MC}(h)$$

$$\hat{I}_{MC}(h) = \frac{1}{M} \sum_{m=1}^M h(\theta^{(m)})$$

random sample (above  $\theta^{(m)}$ )  
random sample (below  $\theta^{(m)}$ )

$$\xrightarrow{\text{a.s.}} \mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)]$$

if  $h(\theta)$  is integrable w.r.t.  $\pi(\theta)$ :

$$\int |h(\theta)| \pi(\theta) d\theta < \infty$$

So, if we can draw samples from  $\pi(\theta)$ , we can approximate quantity related to  $\pi(\theta)$

PROBLEM: We cannot draw samples from  $\pi(\theta)$

We utilize more clever sampling techniques to approximately generate samples from  $\pi(\theta)$

- ① Rejection sampling
- ② Importance sampling
- ③ MCMC

## Rejection Sampling:

- Setup: Target distribution:  $\pi(\theta)$   
Proposal distribution:  $q(\theta)$

$$M \geq \sup_{\theta \in \Theta} \frac{\pi(\theta)}{q(\theta)}$$

Supremum

- Algorithm:
  - ① Sample  $\theta_* \sim q(\theta)$
  - ② Sample  $u \sim \mathcal{U}(0, 1)$
  - ③ Check if  $u < \frac{\pi(\theta_*)}{M q(\theta_*)}$ :

Accept the sample  
 $\theta = \theta_*$

Otherwise:

reject the sample

Conditions:  $M \geq \sup_{\theta \in \Theta} \frac{\pi(\theta)}{q(\theta)}$

The support of  $\pi(\theta)$   
should be a subset  
of the support of  
 $q(\theta)$

$$q(\theta) > 0 \text{ whenever } \pi(\theta) > 0$$

Prove that the samples that are accepted are samples from  $\pi(\theta)$ :

$$\begin{aligned}
 F_{\theta}(\tilde{\theta}) &= \mathbb{P}(\theta \leq \tilde{\theta}) \quad \text{not a RV} \\
 &= \mathbb{P}(\theta_* \leq \tilde{\theta} \mid \text{accepted } \theta_*) \\
 &= \mathbb{P}(\theta_* \leq \tilde{\theta}, U \leq \frac{\pi(\theta_*)}{Mq(\theta_*)}) \\
 &= \mathbb{P}(U \leq \frac{\pi(\theta_*)}{Mq(\theta_*)})
 \end{aligned}$$

Our goal is to show that this is the CDF of the target distribution!

$$\begin{aligned}
 \textcircled{1} \mathbb{P}(\theta_* \leq \tilde{\theta}, U \leq \frac{\pi(\theta_*)}{Mq(\theta_*)}) \\
 = \int \mathbb{P}(\theta_* \leq \tilde{\theta}, U \leq \frac{\pi(\theta_*)}{Mq(\theta_*)} \mid \theta_* = \gamma)
 \end{aligned}$$

$$q(x) dx$$

$$= \int \mathbb{I}(x \leq \tilde{\theta}) \underbrace{\mathbb{P}\left(U \leq \frac{\pi(x)}{M q(x)}\right)} q(x) dx$$

$$= \int_{-\infty}^{\tilde{\theta}} \underbrace{\mathbb{P}\left(U \leq \frac{\pi(x)}{M q(x)}\right)} q(x) dx$$

$$= \int_{-\infty}^{\tilde{\theta}} \frac{\pi(x)}{M q(x)} q(x) dx$$

$$= \frac{1}{M} \int_{-\infty}^{\tilde{\theta}} \pi(x) dx$$

$$= \frac{1}{M} \mathbb{P}_{\pi}(\theta \leq \tilde{\theta})$$

This is the  
CDF of  
 $\pi(\theta)$

$$\textcircled{2.} \quad \mathbb{P} \left( U \leq \frac{\pi(\theta_*)}{M q(\theta_*)} \right)$$

$$= \int \mathbb{P} \left( U \leq \frac{\pi(\theta_*)}{M q(\theta_*)} \mid \theta_* = \gamma \right) q(\gamma) d\gamma$$

$$= \int \frac{\pi(\gamma)}{M q(\gamma)} q(\gamma) d\gamma$$

$$= \frac{1}{M} \underbrace{\int \pi(\gamma) d\gamma}_{=1} = \frac{1}{M}$$

Plugging  $\textcircled{1.}$  and  $\textcircled{2.}$  together, we find that

$$\underbrace{F(\theta \leq \tilde{\theta})}_{\text{distribution of our generated samples}} = F_{\tilde{\pi}}(\theta \leq \tilde{\theta})$$

distribution  
of our generated  
samples

The idea is to generate  $M$  samples from rejection sampling and then apply MC to obtain estimators of expected values (taken w.r.t.  $\pi(\theta)$ )

## Limitations

① The probability of accepting a sample:

$$\begin{aligned} P(\text{accept}) &= P\left(U \leq \frac{\pi(\theta^*)}{M q(\theta^*)}\right) \\ &= \frac{1}{M} \end{aligned}$$

→ we must choose  $M \geq \sup_{\theta \in \Theta} \frac{\pi(\theta)}{M q(\theta)}$

→ if we choose  $M$  too large, we will reject a lot of samples and rejection sampling becomes inefficient

$$\pi(\theta) \triangleq p(\theta | X_{1:n}) \propto \underbrace{p(X_{1:n} | \theta)}_{\tilde{\pi}(\theta)/Z} p(\theta)$$

→ Best  $M$  will also depend on  $q(\theta)$

# Importance Sampling

$$\mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)] = \int h(\theta) \pi(\theta) d\theta$$

Idea: Rewrite this expectation as an expectation taken w.r.t. another distribution  $q(\theta)$

$$= \int h(\theta) \pi(\theta) \frac{q(\theta)}{q(\theta)} d\theta$$

*Importance weight*

$$= \int h(\theta) \left( \frac{\pi(\theta)}{q(\theta)} \right) q(\theta) d\theta$$

*integrand*

$\rightarrow w(\theta) = \frac{\pi(\theta)}{q(\theta)}$

$$= \mathbb{E}_{\theta \sim q(\theta)} [h(\theta) w(\theta)]$$

# Algorithm (Importance Sampling)

$$\theta^{(1)}, \dots, \theta^{(M)} \stackrel{\text{i.i.d.}}{\sim} q(\theta)$$

$$\hat{I}_{\text{IS}}(h) = \frac{1}{M} \sum_{m=1}^M \underbrace{h(\theta^{(m)}) w(\theta^{(m)})}_{q(\theta^{(m)})}$$

a.s.  $\longrightarrow \mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)]$  by SLLN

Condition for convergence:

$$\int |h(\theta)| w(\theta) q(\theta) d\theta < \infty$$

Limitations: (1)  $\text{Var}[\hat{I}_{\text{IS}}(h)] \propto \text{Var}\left[\frac{h(x) \pi(x)}{q(x)}\right]$

If  $q$  is poorly chosen, this variance can be large

The optimal choice:

$$q(\theta) \propto |h(\theta)| \pi(\theta)$$

② Bayesian inference:

$$\pi(\theta) = \frac{\tilde{\pi}(\theta)}{Z} \quad \begin{array}{l} \text{we don't} \\ \text{know} \\ Z \end{array}$$

plug back into our estimator  $\tilde{w}(\theta^{(m)})$

$$\hat{I}_{IS}(h) = \frac{1}{M} \sum_{m=1}^M h(\theta^{(m)}) \frac{\tilde{\pi}(\theta^{(m)})}{q(\theta^{(m)})} Z$$

$$= \frac{1}{M Z} \sum_{m=1}^M h(\theta^{(m)}) \tilde{w}(\theta^{(m)})$$

$$Z = \int \tilde{\pi}(\theta) d\theta \approx \mathbb{E}_{\theta \sim q(\theta)} [\tilde{w}(\theta)]$$

a.s.  $\rightarrow Z$

$$= \frac{1}{M} \sum_{m=1}^M \left( \frac{1}{M} \sum_{m=1}^M \tilde{w}(\theta^{(m)}) \right) \sum_{m=1}^M h(\theta^{(m)}) \tilde{w}(\theta^{(m)})$$

a.s.  $\rightarrow \theta$       a.s.  $\rightarrow \mathbb{E}_{\theta \sim q(\theta)} [h(\theta) \tilde{w}(\theta)]$

$$X_n \xrightarrow{\text{a.s.}} \mu_x$$

$$Y_n \xrightarrow{\text{a.s.}} \mu_y$$

$$\frac{X_n}{Y_n} \xrightarrow{P} \frac{\mu_x}{\mu_y}$$

$$= \frac{1}{M} \sum_{m=1}^M h(\theta^{(m)}) \tilde{w}(\theta^{(m)})$$

a.s.

By Slutsky's theorem

$\xrightarrow{P}$

$$\mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)]$$

even if the constant  $Z$  is unknown

- Next time,  $\rightarrow$  MCMC (Next week)
- $\rightarrow$  Linear models (Next week)
- $\rightarrow$  Hypothesis test (Rest of course)

$$\hat{I}_{SNIS}(h) = \frac{1}{\left( \frac{1}{M} \sum_{m=1}^M \tilde{w}(\theta^{(m)}) \right)} \sum_{m=1}^M h(\theta^{(m)}) \tilde{w}(\theta^{(m)})$$

↓ a.s.

Z

↓ a.s.

$$\mathbb{E}_{\theta \sim q(\theta)} [h(\theta) \tilde{w}(\theta)]$$

$$\xrightarrow{P} \frac{1}{Z} \mathbb{E}_{\theta \sim q(\theta)} [h(\theta) \tilde{w}(\theta)] \rightarrow \frac{\pi(\theta)}{q(\theta)}$$

$$= \mathbb{E}_{\theta \sim q(\theta)} \left[ h(\theta) \frac{\tilde{w}(\theta)}{Z} \right]$$

$$= \int h(\theta) \frac{\pi(\theta)}{q(\theta)} q(\theta) d\theta$$

$$= \mathbb{E}_{\theta \sim \pi(\theta)} [h(\theta)]$$