

# classical approaches for inverse inference

- conjugate priors: choose  $p^s(Y)$  and  $p^s(X|Y)$  such that  $p^s(Y|X)$  can be analytically calculated and is in the same distribution family as  $p^s(Y)$  ( $\Rightarrow$  incremental Bayesian updating)

↑ advantage: efficient and mathematically elegant

↓ disadvantage: very unrealistic  $\Rightarrow$  big simulation gap

- likelihood-based inference:  $p^s(Y)$  and  $p^s(X|Y)$  are known (and not just  $X = \phi(Y, \eta)$ ) but  $p^s(Y|X)$  is intractable

$\Rightarrow$  create a sample  $\{Y_k \sim p^s(Y|X=x^{obs})\}_{k=1}^K$  using Markov chain Monte Carlo (MCMC) or a variant thereof (HMC)

[important relaxation:  $p^s(X|Y)$  can be unnormalized  $\Leftrightarrow \int p^s(X|Y) dx$  intractable

e.g. Gibbs distribution:  $p^s(X|Y) \propto \frac{1}{Z} \exp(-H(X, Y)/T)$

$T$ : temperature  
 $Z$ : unknown partition function  
 $H(X, Y)$ : energy of state  $X$  with parameters  $Y$

# basic MCMC algorithm

① -  $x^{obs}$  given,  $y^{(0)}$  arbitrary initial guess

- pick a "proposal distribution" for transition probs  $q(y' | y^{(t-1)}, x^{obs})$   
e.g. Gaussian (something easy to work with)

② for  $t=1, \dots, T$

(a)  $y' \sim q(y' | y^{(t-1)}, x^{obs})$  sample proposal

(b) calculate acceptance weight

$$\alpha = \frac{p^s(x^{obs} | y') p^s(y')}{p^s(x^{obs} | y^{(t-1)}) p^s(y^{(t-1)})} = \frac{p^s(y' | x^{obs})}{p^s(y^{(t-1)} | x^{obs})}$$

(evidence  $p^s(x^{obs})$  and partition fct  $Z$  cancel out!)

intractable posteriors  
↓

(c) sample  $u \sim \text{uniform}(0, 1)$  acceptance threshold

(d) if  $u \leq \alpha$ : accept  $y^{(t)} = y'$   
else: reject  $y^{(t)} = y^{(t-1)}$

Thm. for  $T \rightarrow \infty$ ,  $\{y^{(t)}\}_{t=1}^T \sim p^s(y | x^{obs})$

## classical solutions to SBI (could)

• analytic: conjugate priors

• likelihood-based inference:  $p^s(Y)$  and  $p^s(X|Y)$  are known, but posterior  $p^s(Y|X)$  is intractable

⇒ MCMC and its variants: create sample  $\{Y^{(t)} \sim p^s(Y|X^{OLS})\}_{t=1}^T$   
via a clever Markov chain  $Y^{(t)} \sim p(Y|Y^{(t-1)}, X^{OLS})$

- advantages: • very general, lot of mathematical theory (performance guarantees)  
• implemented in "probabilistic programming languages", e.g. STAN and libraries (pyro, pyMC, Edward?, Turing.jl)

- disadvantages: • not amortized, algorithm runs from scratch for each new  $X^{OSS}$   
• expensive - often,  $T$  must be very large for complete "mixing", i.e. until all modes of  $p^s(Y|X^{OSS})$  have been covered  
- often, a long "burn-in" phase is needed to forget a bad initial guess  $Y^{(0)} \Rightarrow$  throw away  $Y^{(0)} \dots Y^{(T_0)}$  hyperparameter or computed  
- samples  $Y^{(t)}$  and  $Y^{(t-1)}$  are close to each other

(even  $y^{(t)} = y^{(t-1)}$  when  $y'$  is rejected)

may bias derived statistics of the chain

$\Rightarrow$  skip  $k$  samples in the chain  $y^{(t)} \rightsquigarrow y^{(t+k)} \rightsquigarrow y^{(t+2k)}$

- often difficult to define proposal distribution  $p(y' | y^{(t-1)})$  that has low rejection rate

$\Rightarrow$  MCMC only applicable when  $\dim(Y)$  is not large and  $p^s(X|Y)$  not too slow  
making MCMC faster is hot research

• likelihood-free inference.  $p^s(X|Y)$  unknown, only implicitly defined as  $\phi_{\#} \tilde{p}(Y, \eta)$   
only samples  $X = \phi(Y, \eta)$  with  $Y, \eta \sim \tilde{p}(Y, \eta)$

Approximate Bayesian Computation (ABC)  $\hat{=}$  brute force

- requires a distance  $\text{dist}(X, X^{\text{obs}})$  for outcomes, usually hand-crafted so far

alg. (0)  $t=0$ , sample  $\{ \}$

(we later become  $\{ Y^{(t)} \sim p^s(Y | X^{\text{obs}}) \}_{t=1}^T$ )

(1) repeat until  $t=T$

(a) sample  $Y \sim p^s(Y)$  (or  $Y, \eta \sim \tilde{p}(Y, \eta)$ ), simulate  $X \sim \phi(Y, \eta)$

(b) if  $\text{dist}(X, X^{\text{obs}}) \leq \epsilon$ : sample.add( $Y$ ),  $t=t+1$ ; else reject

- advantage • works when MCMC doesn't  
• available in libraries (pyABC)

- disadvantage • very slow: if  $\epsilon$  (hyperparameter) is small,  
 $\{Y^{(t)}\}_{t=1}^T$  is a good representation of  $p^S(Y|X^{obs})$   
but rejection rate is very high and vice versa  
no amortisation, wasted simulation budget

• hard to define  $\text{dist}(X, X^{obs})$  if  $X$  is complicated or  $\text{dim}(X)$  is high  
⇒ instead design hand-crafted summary statistics  $h(X)$  and compare  
⇒ only feasible when  $\text{dim}(Y)$  is not too high  $\text{dist}(h(X), h(X^{obs}))$