# validation of SBI (contd.)

- posterior checks: - calibration of true parameters $Y^*$ relative to predicted $p(Y|X^{obs})$
  - compare $p(Y|X)$ to $p^*(Y|X)$
  - check diversity of $p(Y|X)$

- posterior predictive checks:

  - sample from posterior $\{\hat{Y}_k \sim p(Y|X^{obs})\}_{k=1}^M$

  - use predicted parameters is simulation inputs $\{\hat{X}_k = \phi(\hat{Y}_k, \eta)\}_{k=1}^M$

  - compare $\hat{X}_k$ to true $X^{obs} \Rightarrow$ should be similar

    - traditional: $\|\hat{X}_k - X^{obs}\|_2^2$

    $\Big[$ also used to solve inverse problem: $\hat{Y} = \arg\min_Y \|\phi(Y) - X^{obs}\|_2^2$

    disadvantages: - non-linear optimization finds local optimum
      - disregards uncertainty / ambiguity of solution
      - may be ill-posed $\Rightarrow$ add a regularizer to make
        well-posed $\Rightarrow$ possible bias $\Big]$

not good when $\hat{x}_u$ are (correctly!) very diverse

$\Rightarrow$ check calibration of $x^{oss}$ relative to $\{\hat{x}_u\}$



$x_{oss}$ $\qquad$ [ $x_j \equiv$ time points of a dynamic system ]

$\hat{x}_u$

$[ Y \equiv$ parameters of dynamic system,

95% confidence region of $\hat{x}_u$ $\qquad$ e.g virus: infection rate

should contain $x^{oss}$ most of the time $\qquad$ duration of disease $\}$ later ]

if posterior has lower variance than prior ($=$ data $x^{oss}$ gave us information)

$\qquad$ posterior predictive scenarios or more accurate than prior predictive ones

$\qquad$ $\Rightarrow$ better predictions of future outcomes

does not work so easily when behavior is non-stationary, i.e. $Y$ is time-dependent!

gets worse if outcomes $X$ influence $Y$ (feedback)

- **self-consistency error**  train two networks - $p(X|Y)$  (simulation surrogate)

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ - $p(Y|X)$  (posterior)

if networks are correct, they fullfill Bayes rule

$\quad$ fix  $Y_i \sim p^s(Y)$  and  $X_i = \phi(Y_i, \eta)$

for true distributions, we have  $\quad p^s(Y) \cdot p^s(X|Y) = p^s(X) \, p^s(Y|X)$

same should hold for our models  $\quad p^s(Y) \, p(X|Y) = p(X) \, p(Y|X)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ with  $p(X) = \int p(X|Y) \, p^s(Y) \, dY$

rearrange Bayes and take logarithms

$$p(X) = \frac{p^s(Y) \, p(X|Y)}{p(Y|X)} \implies$$

for X fixed.  $\quad \underbrace{\log p(X) = \log p^s(Y) + \log p(X|Y) - \log p(Y|X)}_{} = const$

$\quad\quad\quad\quad\quad\quad\quad\quad$ must be constant for __all__  $Y \sim p(Y|X)$

if models are consistent, for fixed $X$ we have

$$\text{Var}_{y \sim p(Y|X)} \left[ \log p^S(Y) + \log p(X|Y) - \log p(Y|X) \right] = 0$$

$\Rightarrow$ measure variance as a quality score

$\Rightarrow$ even better: use Variance as an additional loss    [Schmitt et al. 2023]

$\Rightarrow$ can also be used to determine $\log p(X)$ <u>without</u> solving integral

$$\log p(x) = E_{y \sim p(Y|X)} \left[ \log p^S(Y) + \log p(X|Y) - \log p(Y|X) \right]$$

$\Rightarrow$ it is beneficial to train surrogate $p(X|Y)$ and posterior $p(Y|X)$ <u>jointly</u>

JANA method         [Radev et al. 2023]

# • sensitivity analysis

- how much do the solutions change, when
  - training set
  - simulation      } change
  - learned model
  - observed data

- example:
  - if we have a lot of prior knowledge $\hat{=}$ highly confident in $p^s(Y)$
    $\Rightarrow$ new data should change our oppinion only a little
    $$kL\left[ p(Y|x) \parallel p^s(Y) \right] \approx 0 \qquad \text{if not} \Rightarrow \text{prior was worse}$$
    $\quad\quad\quad$ $\lfloor$ or MMD $\qquad\qquad\qquad\qquad\qquad$ than we thought

  - if we do not have much prior knowledge, $p^s(Y)$ is "uninformative"
    $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ [should not impose prejudice]
    $\qquad\qquad\qquad\qquad\qquad\qquad$ e.g. Jeffrey's prior

    $\Rightarrow$ the posterior should depend only very little on prior, most
    $\qquad$ information should come from data

    $\Rightarrow$ sequential Bayesian updating  data arrive in batches $X^{(1)}, X^{(2)}, \ldots, X^{(t)},$

$$p(Y \mid X^{(1)}) \propto p(X^{(1)} \mid Y)\, p(Y) \qquad \text{iteration } 1$$

$$\vdots$$

$$p(Y \mid X^{(1)}, \ldots, X^{(t)}) \propto p(X^{(t)} \mid Y) \cdot \underset{\uparrow}{p(Y \mid X^{(1)}, \ldots, X^{(t-1)})} \qquad \text{iteration } t$$

$$KL\left[ p(Y \mid X^{(1)} \ldots X^{(t)}) \| p(Y \mid X^{(1)}, \ldots X^{(t-1)}) \right] \to 0 \quad \text{as } t \to \infty \qquad \begin{array}{l} \text{use posterior at } t-1 \text{ as} \\ \text{prior for } t \end{array}$$

- especially elegant with conjugate priors ($p(Y)$, $p(Y \mid X^{(1)}, \ldots$) are from same distribution family)

- so far, no cheap algorithm to do this with neural networks
  (iterative alg. is not much cheaper than training from scratch with $X^{(1)}, \ldots, X^{(t)}$)

---

in amortized SBI: hyper-parameter aware training

- apply parametric variation to training process <u>and</u> tell networks about
  current parameters $\qquad X \quad$ parameter values of training variation

$$Y \leftarrow \boxed{cNF} \leftarrow z$$

<u>example</u> power-scaling of the prior $\qquad p^s(Y) \longrightarrow \dfrac{p^s(Y)^\alpha}{\zeta(\alpha)}$ $\qquad$ with $\zeta(\alpha=1)=1$

for many standard distributions
analytic formula exist

$\alpha > 1 \Rightarrow$ prior gets sharper (more informative)
$\alpha < 1 \Rightarrow$ prior gets less informative

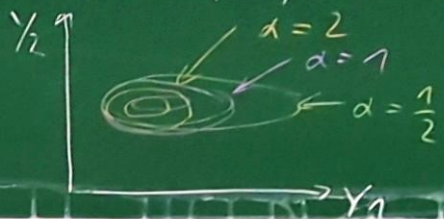e.g $N(0, \sigma^2)^\alpha = N(0, \dfrac{\sigma^2}{\alpha})$

[ $\alpha$ is inverse temperature in Gibb's distribution $\exp\left(-\dfrac{\text{energy}}{kT}\right)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \alpha = \dfrac{1}{T}$

§ during training $\qquad \alpha = p^s(\alpha)$, $Y \sim \dfrac{p^s(Y)^\alpha}{\zeta(\alpha)}$ $\qquad X = \emptyset(Y, \eta)$

$\alpha \in [\frac{1}{2}, 2]$ $\qquad$ additional prior for $\alpha$ (sometimes difficult to choose)

$\qquad\qquad\qquad p^s(\alpha) = \text{uniform}(\frac{1}{2}, 1, 2)$

network is <u>told</u> current value of $\alpha$ $\Rightarrow$ learns posterior $p(Y | X, \alpha)$

∘ during inference $\qquad$ create $p(Y | X^{obs}, \alpha)$ for different values of $\alpha$ and compare
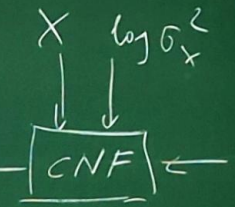
example from epidemiology $\qquad \frac{1}{2}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\alpha = 2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\alpha = 1$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \alpha = \frac{1}{2}$

$Y_1$ depends strongly on prior

$Y_2$ is robust (outcome dominated by the data)

$\qquad Y_1$

example 2 : Noise Net $\{$ kan et al. 2023 $\}$

$\quad$ X may have variable levels of noise

$\qquad \sigma_X^2 \sim p^s(\sigma_X^2) \quad \Rightarrow \quad X = \phi(Y) + N(0, \sigma_X^2) \qquad$ additive noise $y \longleftarrow \boxed{CNF} \longleftarrow$

$\qquad X \quad \log \sigma_X^2$

- train with different $\sigma_X^2$ $\quad$ • inference: estimate actual noise level of $X^{obs}$

error of $\mathbb{E}[Y]$ $\qquad \longleftarrow$ standard training with fixed low noise $\sigma^2$

$\qquad \longleftarrow$ Noise Net



$\sigma_{obs}^2 \qquad \sigma^2$

use same principle for all types of perturbations during training

$\quad$ - prior scaling $\quad$ - likelihood variations $\quad$ - data augmentation/perturbation

etc. $\Rightarrow$ hyper-parameter aware network amortizes over permutations

$\Rightarrow$ sensitivity analysis at inference time cheap (does not require additional training) -

makes many repetitions (as required for trustworthy sensitivity analysis) tractable

$\qquad [$ Elsemüller et al. 2023 $]$

# References:

Maximum mean discrepancy (MMD): [Gretton et al. 2012] A Kernel Two-Sample Test, https://jmlr.csail.mit.edu/papers/v13/gretton12a.html
https://www.onurtunali.com/ml/2019/03/08/maximum-mean-discrepancy-in-machine-learning.html

Frechet Inception Distance (FID score): https://en.wikipedia.org/wiki/Fr%C3%A9chet_inception_distance

Uni- and bivariate posterior plots of model vs. MCMC: [Radev et al. 2020] BayesFlow: Learning complex stochastic models with invertible neural networks, https://arxiv.org/abs/2003.06281, (figure 5)

Calibration via uniformity of a histogram: [Radev et al. 2020] BayesFlow (figures 6 and 7)

Calibration via empirical CDF and random walk: [Säilynoja et al. 2022] Graphical Test for Discrete Uniformity and its Applications in Goodness of Fit Evaluation and Multiple Sample Comparison, https://arxiv.org/abs/2103.10522, see also [Elsemüller et al. 2023]

Proper scoring rules: [Gneiting & Rafterty 2007] Strictly proper scoring rules, prediction, and estimation, https://doi.org/10.1198/016214506000001437

Density and coverage metrics: [Naeem et al. 2020] Reliable Fidelity and Diversity Metrics for Generative Models, https://arxiv.org/abs/2002.09797

Vendi score: [Friedman & Dieng 2022] The Vendi Score: A Diversity Evaluation Metric for Machine Learning, https://arxiv.org/abs/2210.02410

Sensitivity analysis: [Elsemüller et al. 2023] Sensitivity-Aware Amortized Bayesian Inference, https://arxiv.org/abs/2310.11122

Sequential Bayesian updating, conjugate priors: [Kessler et al. 2023] On Sequential Bayesian Inference for Continual Learning, https://arxiv.org/abs/2301.01828

NoiseNet: [Kang et al. 2023] Noise-Net: determining physical properties of $H_{II}$ regions reflecting observational uncertainties, https://doi.org/10.1093/mnras/stad072

Self-consistency via Bayes rule: [Schmitt et al. 2023] Leveraging Self-Consistency for Data-Efficient Amortized Bayesian Inference, https://arxiv.org/abs/2310.04395