# BALANCING SIMULATION-BASED INFERENCE FOR CONSERVATIVE POSTERIORS

**Arnaud Delaunoy**[*1], **Benjamin Kurt Miller**[*2], **Patrick Forré**[2], **Christoph Weniger**[2], **Gilles Louppe**[1]

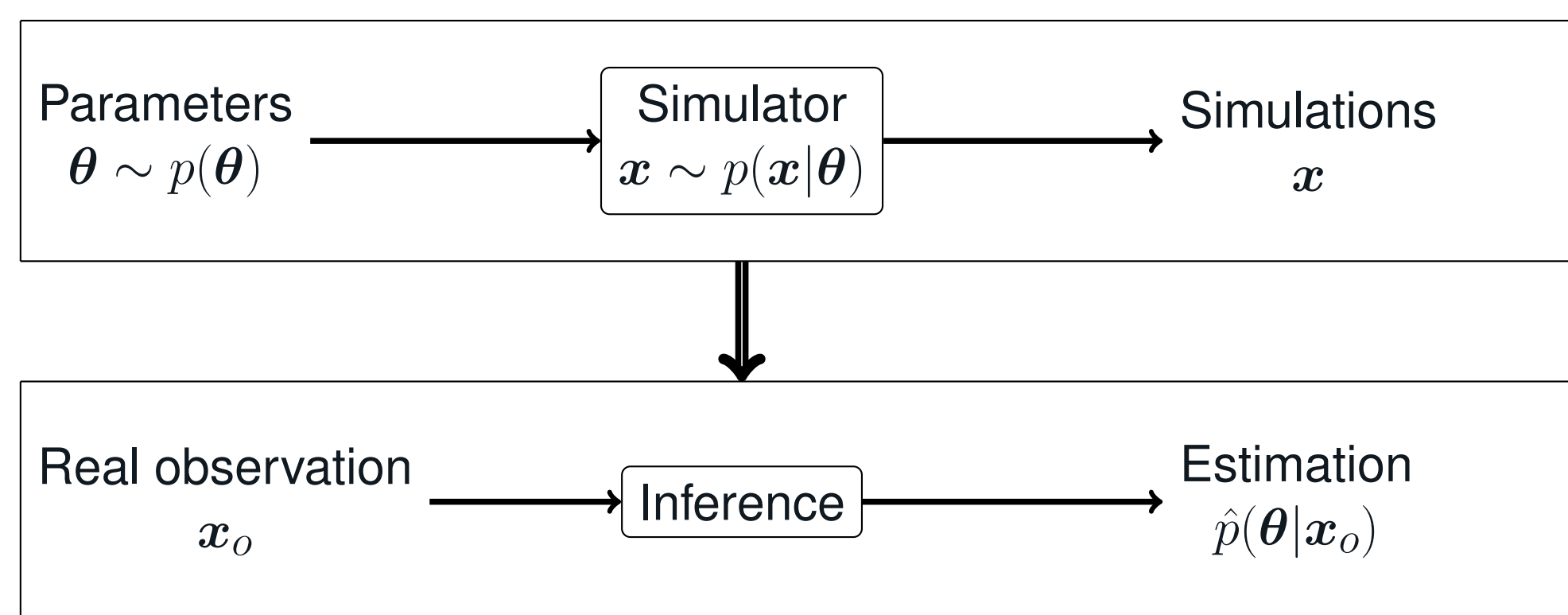[1]University of Liège, [2]University of Amsterdam
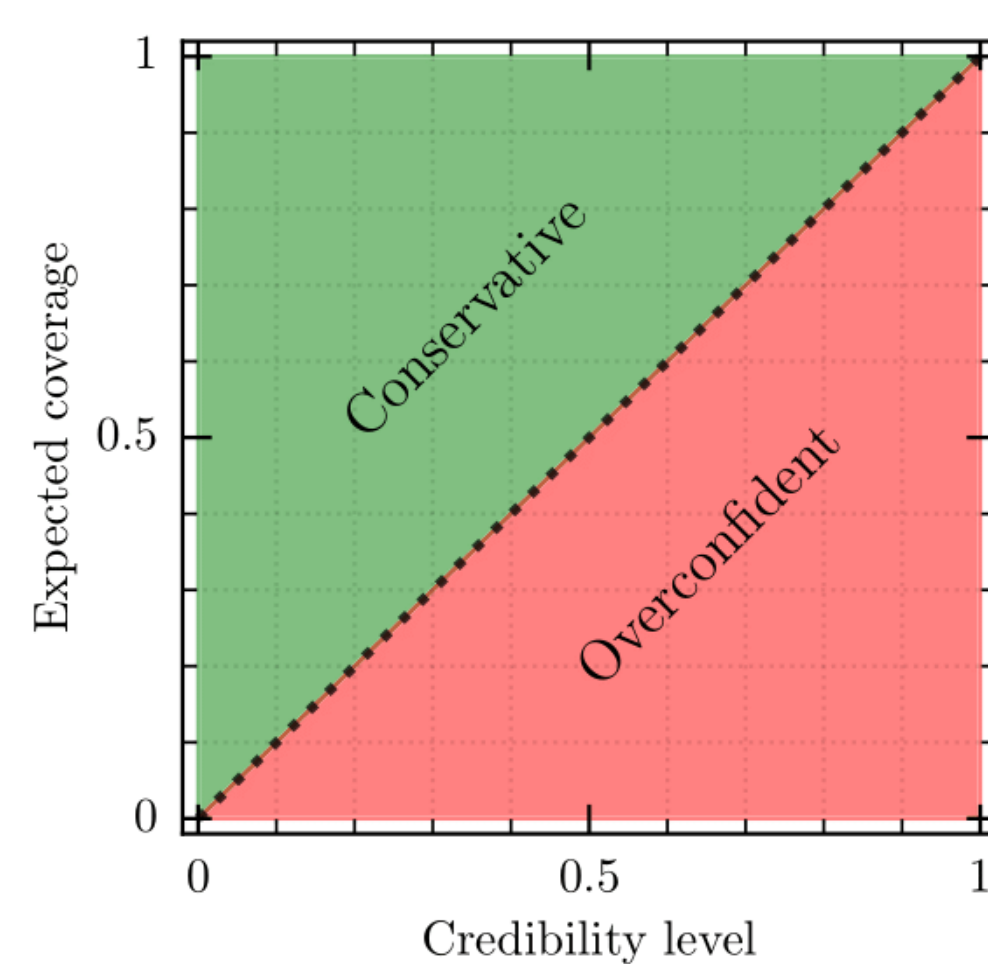
## Simulation-based inference



## Conservativeness in SBI

*Expected coverage probability* of the posterior surrogate $\hat{p}(\theta \mid x)$:

$$1 - \hat{\alpha}[\hat{p}; \alpha] := \mathbb{E}_{p(\theta, x)}\left[\mathbb{1}(\theta \in \Theta_{\hat{p}(\theta \mid x)}(1 - \alpha))\right].$$

- When $\exists \alpha' : 1 - \hat{\alpha}[\hat{p}; \alpha'] < 1 - \alpha'$, we say that $\hat{p}(\theta \mid x)$ is *overconfident*.

- Overconfidence is problematic because the surrogate tends to exclude parameter values that are actually plausible at the considered credibility level. On the other hand, extremely *underconfident* surrogates are not informative. Although there is a tradeoff, scientific applications take a cautious approach by favoring underconfidence.

- **We encourage** *conservative surrogates at credibility level* $\alpha'$**, which have** $1 - \hat{\alpha}[\hat{p}; \alpha'] \geq 1 - \alpha'$.



## Balanced Neural Ratio Estimation

Neural ratio estimation (NRE) trains a classifier $\varpi(y = 1 \mid \theta, x)$ to discriminate between jointly drawn samples, $p(\theta, x)$, and marginal samples, $p(x)p(\theta)$, i.e. sampling:

$$\pi(\theta, x \mid y) := \begin{cases} p(\theta)p(x) & y = 0 \\ p(\theta, x) & y = 1 \end{cases}$$

with marginals $\pi(y = 0) := \pi(y = 1) := \frac{1}{2}$.

Data generation for training:

$$y \sim \pi(y) := \text{Ber}(y; \frac{1}{2}) \begin{cases} y = 0 \to (\theta, x) \sim p(\theta)p(x) \\ y = 1 \to (\theta, x) \sim p(\theta, x) \end{cases}$$

Inference on $\theta \sim p(\theta)$ and $x_o$:

$$(\theta, x_o) \to \boxed{\begin{array}{c}\text{Classifier} \\ \varpi(y = 1 \mid \theta, x_o)\end{array}} \to \hat{p}(\theta \mid x_o) = \frac{\varpi(y=1 \mid \theta, x_o)}{1 - \varpi(y=1 \mid \theta, x_o)} p(\theta)$$

**Balanced Neural Ratio Estimation (BNRE) regularizes the classifier to be more conservative** by minimizing the balancing criterion (using a Lagrange multiplier) which is expressed as

$$B[\varpi] := B(w) := \left(\mathbb{E}_{p(\theta)p(x)}[\varpi(y=1 \mid \theta, x)] + \mathbb{E}_{p(\theta, x)}[\varpi(y=1 \mid \theta, x)] - 1\right)^2,$$

where $w$ are the classifier weights. This is added to *the main NRE objective, the binary cross entropy.*

## Contribution 1: a new view on the balancing criterion

The $\chi^2$ divergence is defined as

$$\chi^2(\pi(y) \| \varpi(y)) := \int \left(\frac{\varpi(y)}{\pi(y)} - 1\right)^2 \pi(y) \, dy.$$

> We identify that $B[\varpi] = \chi^2(\pi(y) \| \varpi(y))$.

- **Enforcing the balancing criterion regularizes the marginal classifier towards the target distribution for y:** $\pi(y)$**.** This new objective aims to be a building block to construct a better understanding of the balancing criterion. It revels a principle for balancing multi-class classifiers.

- **Why the $\chi^2$ divergence?** The Kullback-Leibler divergence $\text{KL}(\pi(y) \| \varpi(y))$ would be information-theoretically motivated, but it is challenging to optimize due to the $\log$ in the integrand.

## Contribution 2: extending balancing beyond NRE

Define a classifier in terms of the variational (unnormalized) posterior approximant $\hat{q}_w(\theta \mid x)$. We approximate $r(\theta, x) := \frac{p(\theta, x)}{p(\theta)p(x)} = \frac{p(\theta \mid x)}{p(\theta)}$ with $\frac{\hat{q}_w(\theta \mid x)}{p(\theta)}$ which yields the classifier

$$\varpi(y = 1 \mid \theta, x; \hat{q}_w) := \frac{\hat{q}_w(\theta \mid x)/p(\theta)}{1 + \hat{q}_w(\theta \mid x)/p(\theta)}.$$

The balancing criterion can be expressed

$$B(w) := \left(\int (\pi(\theta, x \mid y = 0) + \pi(\theta, x \mid y = 1))\varpi(y = 1 \mid \theta, x; \hat{q}_w) d\theta \, dx - 1\right)^2$$

$$= \left(\int (p(\theta)p(x) + p(\theta, x))\frac{\hat{q}_w(\theta \mid x)/p(\theta)}{1 + \hat{q}_w(\theta \mid x)/p(\theta)} d\theta \, dx - 1\right)^2$$

- **We propose BNPE** which regularizes NPE's maximum likelihood-based objective with the balance criterion to train a normalized density estimator $q_w(\theta \mid x)$. We have $\hat{q}_w(\theta \mid x) := q_w(\theta \mid x)$.

- **We propose BNRE-C** which regularizes NRE-C's multi-class, classifier-based objective with the (binary) balance criterion to train a ratio estimator. How do we define the binary classifier since NRE-C normally only defines a multi-class classifier? We do it in terms of the (unnormalized) density estimator $\hat{q}_w(\theta \mid x) := \exp \circ h_w(\theta, x)p(\theta)$ where $h_w$ is a neural network. Using the above definition, the corresponding (binary) classifier is $\varpi(y = 1 \mid \theta, x; \hat{q}_w) := \frac{\exp \circ h_w(\theta, x)p(\theta)/p(\theta)}{1 + \exp \circ h_w(\theta, x)p(\theta)/p(\theta)} = \sigma \circ h_w(\theta, x)$. This regularizing classifier is the same binary classifier as in BNRE.

## Results



> BNPE and BNRE-C produce more conservative posteriors than NPE and NRE-C, respectively.

## What if the posterior surrogate is imbalanced?

> We observed that BNPE can be harder to balance than classifier-based algorithms. We show that this can be mitigated with a proper initialization scheme.

- **Hypothesis**: In order to be balanced, BNPE needs to learn the prior as the equivalent classifier is a function of both the approximate posterior and the prior.

- **Solution**: Initialize the normalizing flow in a balanced state (close to the prior).

How?

- Use the prior as base distribution or add a transformation that maps the base distribution to the prior at the end of the normalizing flow.

- Initialize all the transformations to an identity function.

**Solving NPE "leakage"** Several NPE papers point out that the variational posterior can "leak" mass outside the prior, i.e. put estimated posterior density in a region with zero prior density. *The bijection from the support of the variational posterior to the prior support solves leakage.* It is generally applicable when such a bijection can be constructed (holds for topologically isomorphic supports).

## Going further

- The $\chi^2$ divergence, equivalent to the balancing criterion, naturally leads to the following regularization term for $K$ classes classification:

$$\chi^2(\tilde{\pi}(y) \| \tilde{\varpi}(y)) = \frac{1}{K+1} \sum_{i=0}^{K} \left(\int \tilde{\varpi}(y = i \mid \Theta, x)\left(\sum_{j=0}^{K} \tilde{\pi}(\Theta, x \mid y = j)\right) d\Theta \, dx - 1\right)^2.$$

The effect of this regularizer remains to be studied!

- We extend balancing to algorithms that provide an approximate posterior density however some methods do not fall into this framework (score-based methods, GANs, ...). Future work could reformulate our regularizer to apply to these works. It would require defining a purely sample-based (don't evaluate the estimated posterior density) version of the balancing criterion.

## Take-home messages

- The balancing criterion can be expressed as the $\chi^2$ divergence between the marginal classifier and target marginal distribution over classes. This provides a new perspective on balancing and serves as a building block for further development.

- The balancing criterion can be extended to algorithms that provide an approximate posterior density. This broadens the applicability of balancing, enabling more conservative algorithms.

- Empirically, balancing makes posteriors more conservative. Although, it may require a hyper-parameter search to find the Lagrange multiplier that yields a conservative posterior estimate.