

# Self-supervised learning for Bayesian experimental design

Michael U. Gutmann

`michael.gutmann@ed.ac.uk`

School of Informatics, University of Edinburgh

Corcoran Memorial Lecture

29 January 2024

# Contents

## Research objective

- Two main goals: inference and experimental design

- Tasks are computationally intractable for simulator models

## Self-supervised learning to deal with intractability

- Link to logistic regression and Jensen-Shannon divergence

- Technical challenge: the density-chasm problem

## Application to Bayesian experimental design

- Via self-supervised learning of density ratios

- Exploiting bounds to increase computational efficiency

# DALL·E's visual summary of the talk



# Contents

## Research objective

- Two main goals: inference and experimental design

- Tasks are computationally intractable for simulator models

## Self-supervised learning to deal with intractability

- Link to logistic regression and Jensen-Shannon divergence

- Technical challenge: the density-chasm problem

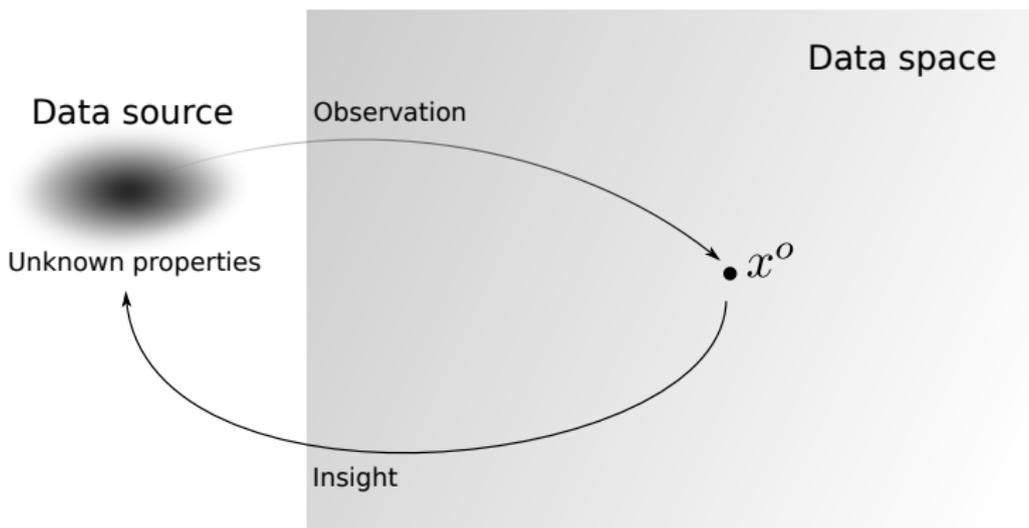
## Application to Bayesian experimental design

- Via self-supervised learning of density ratios

- Exploiting bounds to increase computational efficiency

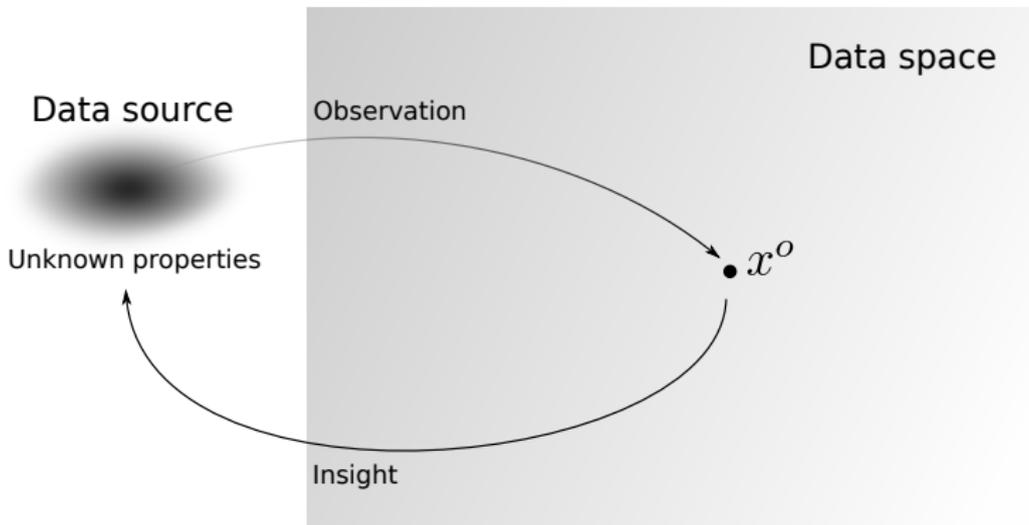
# Overall goal

- ▶ Goal: Understanding properties of some data source
- ▶ Enables predictions, decision making under uncertainty, ...



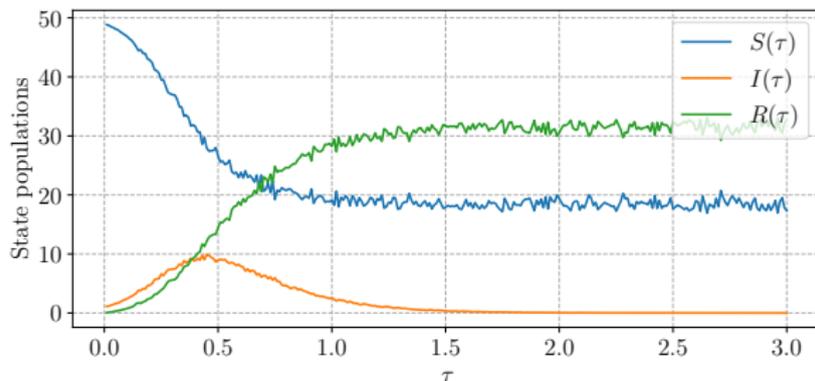
# Two fundamental tasks

- ▶ **Inference task** : Given  $\mathbf{x}_o$ , what can we robustly say about the properties of the source?
- ▶ **Experimental design task** : How to obtain a  $\mathbf{x}_o$  that is maximally useful for learning about the properties?



## Example: stochastic SIR model

- ▶ Stochastic model describing the population of susceptibles  $S(\tau)$ , infected  $I(\tau)$  and recovered  $R(\tau)$  as a function of time.
- ▶ Parameters  $\theta$ : rate of infection  $\beta$  and the rate of recovery  $\gamma$ .
- ▶ **Inference task** : determine plausible values of  $\beta$  and  $\gamma$  given some measurements of the population sizes.
- ▶ **Exp design task** : find the optimal times at which to perform the measurements to most accurately estimate  $\beta$  and  $\gamma$ .



(Figure by Steven Kleinagesse)

# Bayesian inference and design with tractable models

- ▶ Assume model is expressed as a family of pdfs  $\{p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})\}$  indexed by parameter  $\boldsymbol{\theta}$  and design variable  $\mathbf{d}$ .

# Bayesian inference and design with tractable models

- ▶ Assume model is expressed as a family of pdfs  $\{p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})\}$  indexed by parameter  $\boldsymbol{\theta}$  and design variable  $\mathbf{d}$ .
- ▶ Bayesian inference of  $\boldsymbol{\theta}$  for data  $\mathbf{x}_o$  obtained with design  $\mathbf{d}_o$ :

$$p(\boldsymbol{\theta}|\mathbf{x}, \mathbf{d}) = \frac{p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})}{p(\mathbf{x}|\mathbf{d})}p(\boldsymbol{\theta}|\mathbf{d}) \quad (1)$$

with  $\mathbf{x}$  fixed to  $\mathbf{x}_o$  and  $\mathbf{d}$  to  $\mathbf{d}_o$ .

# Bayesian inference and design with tractable models

- ▶ Assume model is expressed as a family of pdfs  $\{p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})\}$  indexed by parameter  $\boldsymbol{\theta}$  and design variable  $\mathbf{d}$ .
- ▶ Bayesian inference of  $\boldsymbol{\theta}$  for data  $\mathbf{x}_o$  obtained with design  $\mathbf{d}_o$ :

$$p(\boldsymbol{\theta}|\mathbf{x}, \mathbf{d}) = \frac{p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})}{p(\mathbf{x}|\mathbf{d})}p(\boldsymbol{\theta}|\mathbf{d}) \quad (1)$$

with  $\mathbf{x}$  fixed to  $\mathbf{x}_o$  and  $\mathbf{d}$  to  $\mathbf{d}_o$ .

- ▶ Experimental design by maximising mutual information (MI) between data  $\mathbf{x}$  and parameters  $\boldsymbol{\theta}$ :

$$\hat{\mathbf{d}} = \underset{\mathbf{d}}{\operatorname{argmax}} \operatorname{MI}_{\mathbf{d}}(\mathbf{x}, \boldsymbol{\theta}) \quad (2)$$

$$\operatorname{MI}_{\mathbf{d}}(\mathbf{x}, \boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}|\mathbf{d})} \operatorname{KL}(p(\boldsymbol{\theta}|\mathbf{x}, \mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})) \quad (3)$$

$$= \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{d})} \log \left[ \frac{p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})}{p(\mathbf{x}|\mathbf{d})} \right] \quad (4)$$

# Simulator models

- ▶ Not all models are specified as family of pdfs.
- ▶ We consider here the important class of simulator models: models that are specified via a parameterised stochastic mechanism for generating data



DALL·E's view on simulator models

# Simulator models

- ▶ Technically, a simulator model is a measurable function  $g$

$$\mathbf{x} = g(\boldsymbol{\theta}, \mathbf{d}, \boldsymbol{\omega}), \quad \boldsymbol{\omega} \sim p(\boldsymbol{\omega}) \quad (5)$$

Maps params  $\boldsymbol{\theta}$ , design variables  $\mathbf{d}$ , and “noise”  $\boldsymbol{\omega}$  to data  $\mathbf{x}$

# Simulator models

- ▶ Technically, a simulator model is a measurable function  $g$

$$\mathbf{x} = g(\boldsymbol{\theta}, \mathbf{d}, \boldsymbol{\omega}), \quad \boldsymbol{\omega} \sim p(\boldsymbol{\omega}) \quad (5)$$

Maps params  $\boldsymbol{\theta}$ , design variables  $\mathbf{d}$ , and “noise”  $\boldsymbol{\omega}$  to data  $\mathbf{x}$

- ▶ Function  $g$  is not known in closed form but implemented as a (complex) computer programme.

# Simulator models

- ▶ Technically, a simulator model is a measurable function  $g$

$$\mathbf{x} = g(\boldsymbol{\theta}, \mathbf{d}, \omega), \quad \omega \sim p(\omega) \quad (5)$$

Maps params  $\boldsymbol{\theta}$ , design variables  $\mathbf{d}$ , and “noise”  $\omega$  to data  $\mathbf{x}$

- ▶ Function  $g$  is not known in closed form but implemented as a (complex) computer programme.
- ▶ No closed form expression for  $p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})$  available

# Simulator models

- ▶ Technically, a simulator model is a measurable function  $g$

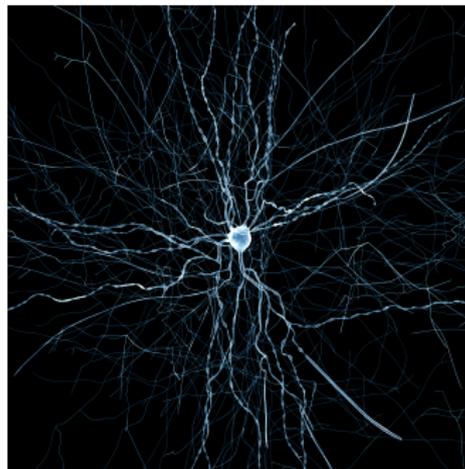
$$\mathbf{x} = g(\boldsymbol{\theta}, \mathbf{d}, \omega), \quad \omega \sim p(\omega) \quad (5)$$

Maps params  $\boldsymbol{\theta}$ , design variables  $\mathbf{d}$ , and “noise”  $\omega$  to data  $\mathbf{x}$

- ▶ Function  $g$  is not known in closed form but implemented as a (complex) computer programme.
- ▶ No closed form expression for  $p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})$  available
- ▶ Sampling data  $\mathbf{x}|\boldsymbol{\theta}, \mathbf{d} \sim p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})$  is possible

# Simulator models are widely used in the natural sciences

- ▶ Evolutionary biology:  
to model evolution
- ▶ Biochemistry:  
to model gene expression
- ▶ Neuroscience:  
to model neural processing
- ▶ Cognitive sciences:  
to model human decision  
making
- ▶ Epidemiology:  
to model the spread of an  
infectious disease
- ▶ ...



Simulated neural activity in rat somatosensory cortex  
(Figure from <https://bbp.epfl.ch/nmc-portal>)

# Research objective

- ▶ Simulator models have great modelling power.

# Research objective

- ▶ Simulator models have great modelling power.
- ▶ However, we pay the price when attempting to perform inference and experimental design:  
evaluating  $p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})$  is computationally intractable

# Research objective

- ▶ Simulator models have great modelling power.
- ▶ However, we pay the price when attempting to perform inference and experimental design:  
evaluating  $p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})$  is computationally intractable
- ▶ Paradoxical situation: we have great models from the natural sciences but cannot fully use them because we lack suitable tools to perform inference and experimental design with them.

# Research objective

- ▶ Simulator models have great modelling power.
- ▶ However, we pay the price when attempting to perform inference and experimental design:  
evaluating  $p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{d})$  is computationally intractable
- ▶ Paradoxical situation: we have great models from the natural sciences but cannot fully use them because we lack suitable tools to perform inference and experimental design with them.
- ▶ Research objective:  
Develop efficient tools for Bayesian inference and experimental design with simulator models.

# Contents

## Research objective

Two main goals: inference and experimental design

Tasks are computationally intractable for simulator models

## Self-supervised learning to deal with intractability

Link to logistic regression and Jensen-Shannon divergence

Technical challenge: the density-chasm problem

## Application to Bayesian experimental design

Via self-supervised learning of density ratios

Exploiting bounds to increase computational efficiency

## Basic idea

- ▶ Self-supervised learning is a paradigm in machine learning where labels are generated by the learning algorithm itself without manual labelling.

## Basic idea

- ▶ Self-supervised learning is a paradigm in machine learning where labels are generated by the learning algorithm itself without manual labelling.
- ▶ Will focus on “contrastive self-supervised learning”.

# Basic idea

- ▶ Self-supervised learning is a paradigm in machine learning where labels are generated by the learning algorithm itself without manual labelling.
- ▶ Will focus on “contrastive self-supervised learning”.
- ▶ The basic idea is to learn the difference between the data of interest and some reference data.

# Basic idea

- ▶ Self-supervised learning is a paradigm in machine learning where labels are generated by the learning algorithm itself without manual labelling.
- ▶ Will focus on “contrastive self-supervised learning”.
- ▶ The basic idea is to learn the difference between the data of interest and some reference data.
- ▶ Properties of the reference are typically known or not of interest; by learning the difference we focus the (computational) resources on learning what matters.

# Basic idea

- ▶ Self-supervised learning is a paradigm in machine learning where labels are generated by the learning algorithm itself without manual labelling.
- ▶ Will focus on “contrastive self-supervised learning”.
- ▶ The basic idea is to learn the difference between the data of interest and some reference data.
- ▶ Properties of the reference are typically known or not of interest; by learning the difference we focus the (computational) resources on learning what matters.
- ▶ As straightforward as

$$\underbrace{b}_{\text{reference}} + \underbrace{a - b}_{\text{difference}} \Rightarrow \underbrace{a}_{\text{interest}} \quad (6)$$

# Basic idea

- ▶ As straightforward as

$$\underbrace{b}_{\text{reference}} + \underbrace{a - b}_{\text{difference}} \Rightarrow \underbrace{a}_{\text{interest}} \quad (7)$$

# Basic idea

- ▶ As straightforward as

$$\underbrace{b}_{\text{reference}} + \underbrace{a - b}_{\text{difference}} \Rightarrow \underbrace{a}_{\text{interest}} \quad (7)$$

- ▶ Link to (log) density ratio estimation

$$\underbrace{\log p_b}_{\text{reference}} + \underbrace{\log p_a - \log p_b}_{\text{difference}} \Rightarrow \underbrace{\log p_a}_{\text{interest}} \quad (8)$$

# Basic idea

- ▶ As straightforward as

$$\underbrace{b}_{\text{reference}} + \underbrace{a - b}_{\text{difference}} \Rightarrow \underbrace{a}_{\text{interest}} \quad (7)$$

- ▶ Link to (log) density ratio estimation

$$\underbrace{\log p_b}_{\text{reference}} + \underbrace{\log p_a - \log p_b}_{\text{difference}} \Rightarrow \underbrace{\log p_a}_{\text{interest}} \quad (8)$$

- ▶ Link to Bayes' rule

$$\underbrace{\log p(\boldsymbol{\theta})}_{\text{reference}} + \underbrace{\log p(\mathbf{x}|\boldsymbol{\theta}) - \log p(\mathbf{x})}_{\text{difference}} \Rightarrow \underbrace{\log p(\boldsymbol{\theta}|\mathbf{x})}_{\text{interest}} \quad (9)$$

# Logistic loss

- ▶ Link to classification: learning differences between data sets can be seen as a classification problem.

## Logistic loss

- ▶ Link to classification: learning differences between data sets can be seen as a classification problem.
- ▶ Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be the data of interest,  $\mathbf{x}_i \sim p$  (iid), and  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  be the reference data,  $\mathbf{y}_i \sim q$  (iid).

## Logistic loss

- ▶ Link to classification: learning differences between data sets can be seen as a classification problem.
- ▶ Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be the data of interest,  $\mathbf{x}_i \sim p$  (iid), and  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  be the reference data,  $\mathbf{y}_i \sim q$  (iid).
- ▶ Label the data:  $(\mathbf{x}_i, 1)$ ,  $(\mathbf{y}_i, 0)$  and minimise the (rescaled) logistic loss  $J(h)$

$$J(h) = \frac{1}{n} \sum_{i=1}^n \log [1 + \nu \exp(-h(\mathbf{x}_i))] + \frac{\nu}{m} \sum_{i=1}^m \log \left[ 1 + \frac{1}{\nu} \exp(h(\mathbf{y}_i)) \right] \quad (10)$$

where  $\nu = n/m$  and  $h$  is a nonlinearity (e.g. neural network) that we learn.

## Logistic loss

- ▶ Link to classification: learning differences between data sets can be seen as a classification problem.
- ▶ Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be the data of interest,  $\mathbf{x}_i \sim p$  (iid), and  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  be the reference data,  $\mathbf{y}_i \sim q$  (iid).
- ▶ Label the data:  $(\mathbf{x}_i, 1)$ ,  $(\mathbf{y}_i, 0)$  and minimise the (rescaled) logistic loss  $J(h)$

$$J(h) = \frac{1}{n} \sum_{i=1}^n \log [1 + \nu \exp(-h(\mathbf{x}_i))] + \frac{\nu}{m} \sum_{i=1}^m \log \left[ 1 + \frac{1}{\nu} \exp(h(\mathbf{y}_i)) \right] \quad (10)$$

where  $\nu = n/m$  and  $h$  is a nonlinearity (e.g. neural network) that we learn.

- ▶ For large sample sizes  $n$  and  $m$  (and fixed ratio  $\nu$ ), the optimal  $h$  is

$$h^* = \log p - \log q \quad (11)$$

# Logistic loss

Two key points:

1. The optimisation is done without any constraints (e.g. normalisation). The optimal  $h$  is automatically the ratio between two *densities*

$$h^* = \log p - \log q \quad (12)$$

2. We only need samples from  $p$  and  $q$ ; we do not need their densities or a model of them (but we do need an appropriate model for the ratio)

## Logistic loss

- ▶ For large sample sizes  $n$  and  $m$ ,  $J(h) \rightarrow \bar{J}(h)$  and the corresponding minimal loss is

$$\begin{aligned}\bar{J}(h^*) &= \mathbb{E}_{\mathbf{x} \sim p} \log \left[ 1 + \nu \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] + \nu \mathbb{E}_{\mathbf{y} \sim q} \log \left[ 1 + \frac{p(\mathbf{y})}{\nu q(\mathbf{y})} \right] \\ &= -\mathbb{E}_{\mathbf{x} \sim p} \log \left[ \frac{p(\mathbf{x})}{p(\mathbf{x}) + \nu q(\mathbf{x})} \right] - \nu \mathbb{E}_{\mathbf{y} \sim q} \log \left[ \frac{\nu q(\mathbf{y})}{p(\mathbf{y}) + \nu q(\mathbf{y})} \right]\end{aligned}\tag{13}$$

## Logistic loss

- ▶ For large sample sizes  $n$  and  $m$ ,  $J(h) \rightarrow \bar{J}(h)$  and the corresponding minimal loss is

$$\begin{aligned}\bar{J}(h^*) &= \mathbb{E}_{\mathbf{x} \sim p} \log \left[ 1 + \nu \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] + \nu \mathbb{E}_{\mathbf{y} \sim q} \log \left[ 1 + \frac{p(\mathbf{y})}{\nu q(\mathbf{y})} \right] \\ &= -\mathbb{E}_{\mathbf{x} \sim p} \log \left[ \frac{p(\mathbf{x})}{p(\mathbf{x}) + \nu q(\mathbf{x})} \right] - \nu \mathbb{E}_{\mathbf{y} \sim q} \log \left[ \frac{\nu q(\mathbf{y})}{p(\mathbf{y}) + \nu q(\mathbf{y})} \right]\end{aligned}\tag{13}$$

- ▶ For  $\nu = 1$  and introducing  $m = (p + q)/2$

$$\bar{J}(h^*) = -\text{KL}(p||m) - \text{KL}(q||m) + 2 \log 2 \tag{14}$$

$$= -2\text{JSD}(p, q) + 2 \log 2 \tag{15}$$

## Logistic loss

- ▶ For large sample sizes  $n$  and  $m$ ,  $J(h) \rightarrow \bar{J}(h)$  and the corresponding minimal loss is

$$\begin{aligned}\bar{J}(h^*) &= \mathbb{E}_{\mathbf{x} \sim p} \log \left[ 1 + \nu \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] + \nu \mathbb{E}_{\mathbf{y} \sim q} \log \left[ 1 + \frac{p(\mathbf{y})}{\nu q(\mathbf{y})} \right] \\ &= -\mathbb{E}_{\mathbf{x} \sim p} \log \left[ \frac{p(\mathbf{x})}{p(\mathbf{x}) + \nu q(\mathbf{x})} \right] - \nu \mathbb{E}_{\mathbf{y} \sim q} \log \left[ \frac{\nu q(\mathbf{y})}{p(\mathbf{y}) + \nu q(\mathbf{y})} \right]\end{aligned}\tag{13}$$

- ▶ For  $\nu = 1$  and introducing  $m = (p + q)/2$

$$\bar{J}(h^*) = -\text{KL}(p||m) - \text{KL}(q||m) + 2 \log 2 \tag{14}$$

$$= -2\text{JSD}(p, q) + 2 \log 2 \tag{15}$$

- ▶ Since we are minimising the loss  $\bar{J}(h)$ , we have

$$\bar{J}(h) \geq -2\text{JSD}(p, q) + 2 \log 2 \tag{16}$$

- ▶ Rearranging, we obtain

$$\text{JSD}(p, q) \geq \log 2 - \frac{1}{2} \bar{J}(h) \quad (17)$$

$$\text{JSD}(p, q) = \log 2 - \frac{1}{2} \bar{J}(h^*) \quad (18)$$

# Logistic loss

- ▶ Rearranging, we obtain

$$\text{JSD}(p, q) \geq \log 2 - \frac{1}{2} \bar{J}(h) \quad (17)$$

$$\text{JSD}(p, q) = \log 2 - \frac{1}{2} \bar{J}(h^*) \quad (18)$$

- ▶ Contrastive learning via classification with the logistic loss corresponds to estimating the Jensen-Shannon divergence (JSD) between  $p$  and  $q$ .

# Logistic loss

- ▶ Rearranging, we obtain

$$\text{JSD}(p, q) \geq \log 2 - \frac{1}{2} \bar{J}(h) \quad (17)$$

$$\text{JSD}(p, q) = \log 2 - \frac{1}{2} \bar{J}(h^*) \quad (18)$$

- ▶ Contrastive learning via classification with the logistic loss corresponds to estimating the Jensen-Shannon divergence (JSD) between  $p$  and  $q$ .
- ▶ For a review paper on statistical applications of contrastive learning, see Gutmann, Kleinegese, and Rhodes, *Behaviormetrika*, 2022.

## Other loss functions

- ▶ In the following, I will focus on the logistic loss as done in our early work on contrastive learning for the estimation of unnormalised models, “Noise-contrastive estimation (NCE)” (Gutmann and Hyvärinen, AISTATS 2010).

# Other loss functions

- ▶ In the following, I will focus on the logistic loss as done in our early work on contrastive learning for the estimation of unnormalised models, “Noise-contrastive estimation (NCE)” (Gutmann and Hyvärinen, AISTATS 2010).
- ▶ But other loss functions can be used:
  - ▶ multinomial logistic loss (Srivastava, et al, TMLR 2023)
  - ▶ Bregman divergences (Gutmann and Hirayama, UAI 2011)
  - ▶ f-divergences (e.g. Rhodes and Gutmann, AISTATS, 2019)
  - ▶ ...

# Constructing reference data

Choice depends on the specific application of contrastive learning.

- ▶ Fit a preliminary model and keep it fixed (as often done in NCE)
- ▶ Iterative approach: fitted model becomes reference in the next iteration (as also done in our original work on NCE)
- ▶ Use other segments for time series data  
(Hyvärinen and Morioka, NeurIPS 2016)
- ▶ For Bayesian inference, use prior predictive distribution  
(Thomas et al, 2016; Thomas et al, Bayesian Analysis, 2020)
- ▶ Generate it conditionally on observed data  
(Ceylan and Gutmann, ICML 2018)
- ▶ Iterative adaptive approach with generative models: results into GANs (Goodfellow et al, NeurIPS 2014)
- ▶ Iterative adaptive approach with flexible density model such as flows (“Flow-contrastive estimation”, Gao et al, NeurIPS 2019)
- ▶ . . .

# Constructing reference data

Is there an optimal reference (“noise”) distribution?

## Constructing reference data

Is there an optimal reference (“noise”) distribution?

For parameter estimation, see the paper *The Optimal Noise in Noise-Contrastive Learning Is Not What You Think* by Omar Chehab, Alex Gramfort, Aapo Hyvarinen, at UAI, 2022.

# The density-chasm problem

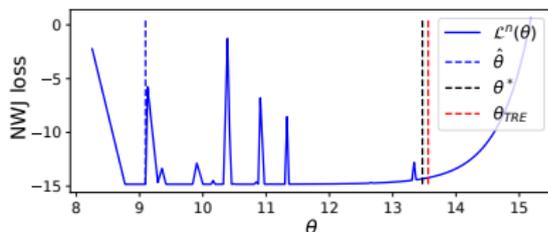
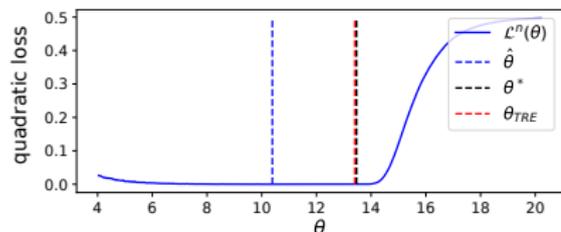
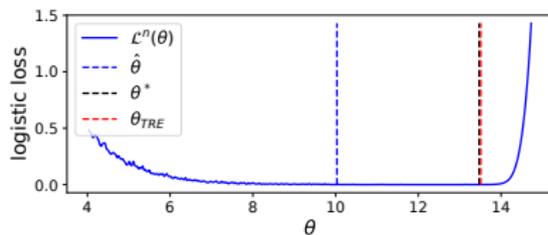
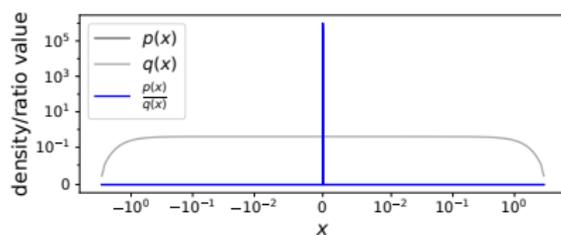
- ▶ Problem: Single ratio methods are sample inefficient if the two distributions are very different (“density chasm”)

# The density-chasm problem

- ▶ Problem: Single ratio methods are sample inefficient if the two distributions are very different (“density chasm”)
- ▶ Consider ratio between two zero-mean Gaussians. 10'000 samples from each distribution. Ratio parametrised by  $\theta \in \mathbb{R}$ .

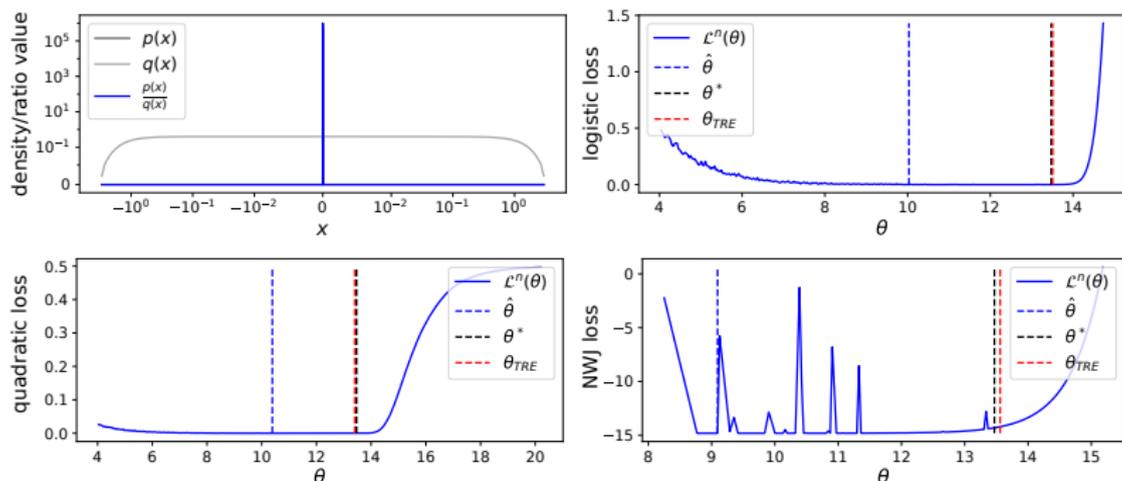
# The density-chasm problem

- ▶ Problem: Single ratio methods are sample inefficient if the two distributions are very different (“density chasm”)
- ▶ Consider ratio between two zero-mean Gaussians. 10'000 samples from each distribution. Ratio parametrised by  $\theta \in \mathbb{R}$ .



# The density-chasm problem

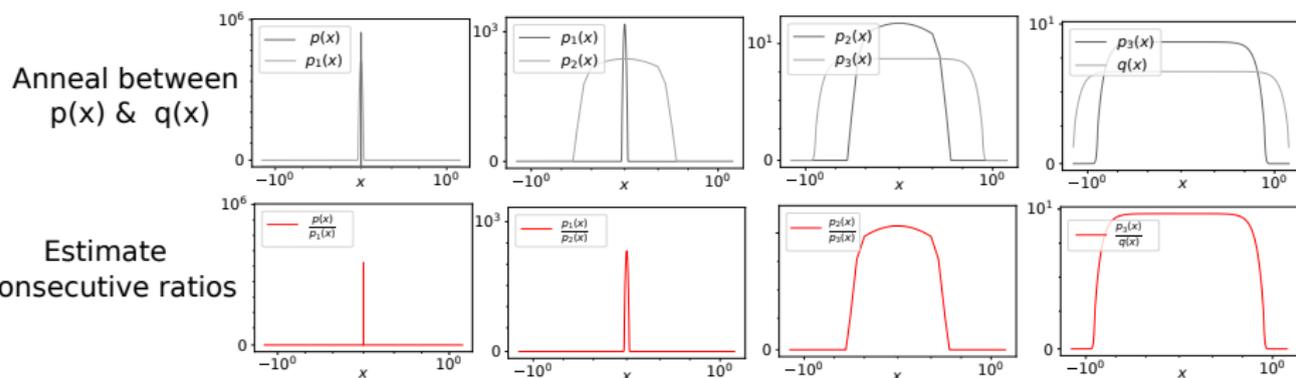
- ▶ Problem: Single ratio methods are sample inefficient if the two distributions are very different (“density chasm”)
- ▶ Consider ratio between two zero-mean Gaussians. 10'000 samples from each distribution. Ratio parametrised by  $\theta \in \mathbb{R}$ .
- ▶ Solution in red bridges the “gap” using telescopic ratio estimation (TRE) (Rhodes, Xu, and Gutmann, NeurIPS 2020)



# Telescoping density-ratio estimation (Rhodes, Xu, and Gutmann, NeurIPS 2020)

A single density-ratio fails to “bridge” the density-chasm.

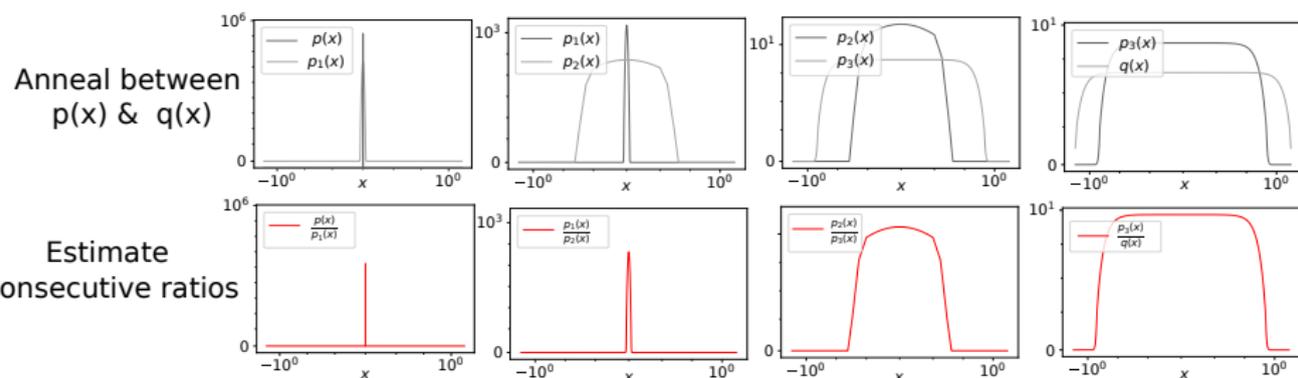
Let us thus use multiple bridges.



# Telescoping density-ratio estimation (Rhodes, Xu, and Gutmann, NeurIPS 2020)

A single density-ratio fails to “bridge” the density-chasm.

Let us thus use multiple bridges.



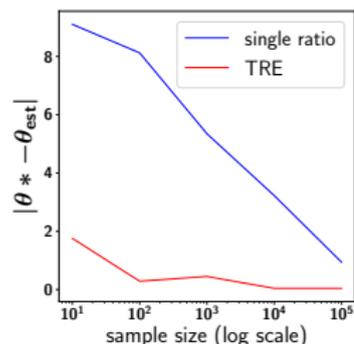
(relabel  $p \equiv p_0$  and  $q \equiv p_4$ ) and compute *telescoping* product

$$\frac{p_0(\mathbf{x})}{p_4(\mathbf{x})} = \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})} \frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} \frac{p_2(\mathbf{x})}{p_3(\mathbf{x})} \frac{p_3(\mathbf{x})}{p_4(\mathbf{x})}. \quad (19)$$

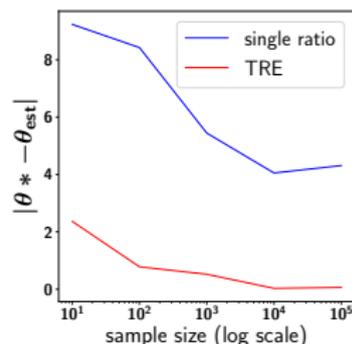
# Telescoping density-ratio estimation (Rhodes, Xu, and Gutmann, NeurIPS 2020)

Sample efficiency curves for the 1d peaked ratio experiment.

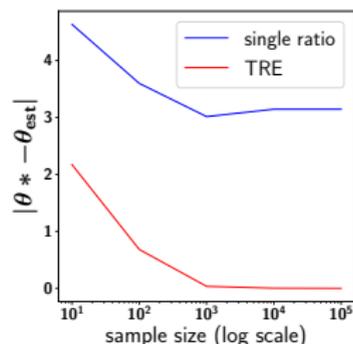
More results in the paper!



(a) Logistic loss



(b) NWJ loss

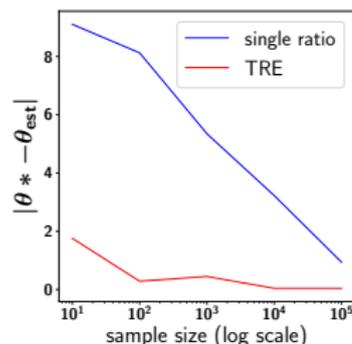


(c) Quadratic loss

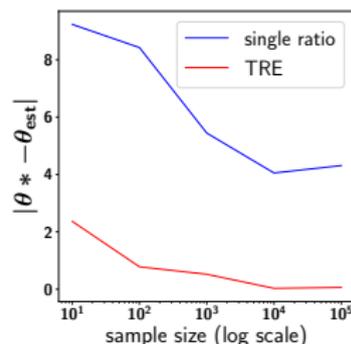
# Telescoping density-ratio estimation (Rhodes, Xu, and Gutmann, NeurIPS 2020)

Sample efficiency curves for the 1d peaked ratio experiment.

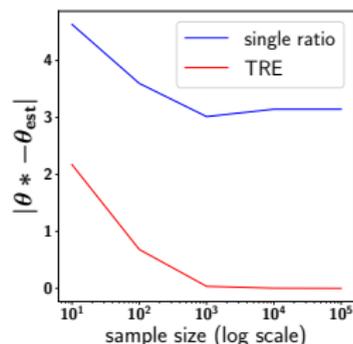
More results in the paper!



(a) Logistic loss



(b) NWJ loss



(c) Quadratic loss

For further improvements, see “Estimating the Density Ratio between Distributions with High Discrepancy using Multinomial Logistic Regression”, Srivastava et al, TMLR 2023.

# Contents

## Research objective

Two main goals: inference and experimental design

Tasks are computationally intractable for simulator models

## Self-supervised learning to deal with intractability

Link to logistic regression and Jensen-Shannon divergence

Technical challenge: the density-chasm problem

## Application to Bayesian experimental design

Via self-supervised learning of density ratios

Exploiting bounds to increase computational efficiency

# Example: stochastic SIR model

(Kleinegesse and Gutmann, AISTATS 2019; ICML 2020; arXiv:2105.04379)

(Kleinegesse, Drovandi and Gutmann, Bayesian Analysis 2020)

- ▶ Example: Stochastic SIR model with noisy observations  
Latent process: Susceptibles  $\rightarrow$  Infected  $I(t) \rightarrow$  Recovered  
Observation model:  $y(t)|\theta \sim \text{Poisson}(y; \phi I(t))$

# Example: stochastic SIR model

(Kleinegesse and Gutmann, AISTATS 2019; ICML 2020; arXiv:2105.04379)

(Kleinegesse, Drovandi and Gutmann, Bayesian Analysis 2020)

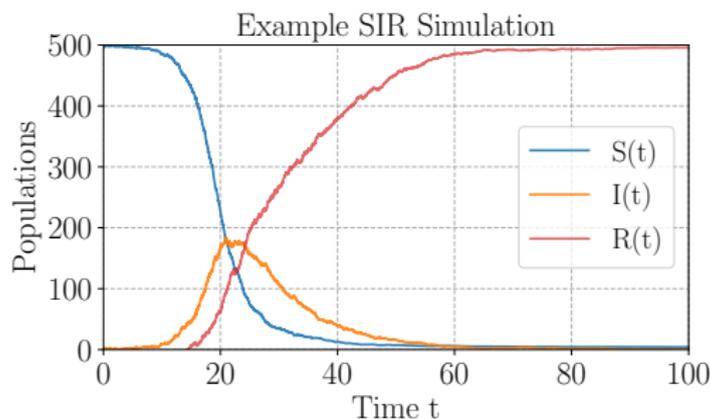
- ▶ Example: Stochastic SIR model with noisy observations  
Latent process: Susceptibles  $\rightarrow$  Infected  $I(t) \rightarrow$  Recovered  
Observation model:  $y(t)|\theta \sim \text{Poisson}(y; \phi I(t))$
- ▶ Parameters  $\theta = (\beta, \gamma)$  (infection rate and recovery rate)

# Example: stochastic SIR model

(Kleinegesse and Gutmann, AISTATS 2019; ICML 2020; arXiv:2105.04379)

(Kleinegesse, Drovandi and Gutmann, Bayesian Analysis 2020)

- ▶ Example: Stochastic SIR model with noisy observations  
Latent process: Susceptibles  $\rightarrow$  Infected  $I(t) \rightarrow$  Recovered  
Observation model:  $y(t)|\theta \sim \text{Poisson}(y; \phi I(t))$
- ▶ Parameters  $\theta = (\beta, \gamma)$  (infection rate and recovery rate)

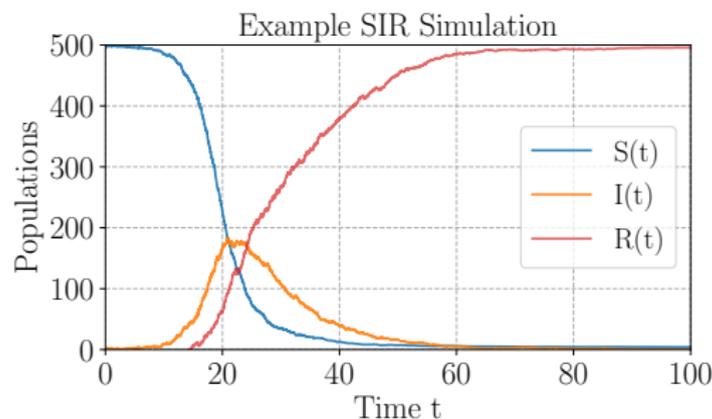


# Example: stochastic SIR model

(Kleinegesse and Gutmann, AISTATS 2019; ICML 2020; arXiv:2105.04379)

(Kleinegesse, Drovandi and Gutmann, Bayesian Analysis 2020)

- ▶ Example: Stochastic SIR model with noisy observations  
Latent process: Susceptibles  $\rightarrow$  Infected  $I(t) \rightarrow$  Recovered  
Observation model:  $y(t)|\theta \sim \text{Poisson}(y; \phi I(t))$
- ▶ Parameters  $\theta = (\beta, \gamma)$  (infection rate and recovery rate)
- ▶ Task: find the optimal times at which to take measurements to most accurately estimate  $\theta$ .



# Experimental design for simulator models

- ▶ Experimental design by maximising mutual information (MI)

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{d})} \log \left[ \frac{p(\mathbf{x} | \boldsymbol{\theta}, \mathbf{d})}{p(\mathbf{x} | \mathbf{d})} \right] \quad (20)$$

- ▶ Use contrastive self-supervised learning to estimate

$$h_{\mathbf{d}}(\mathbf{x}, \boldsymbol{\theta}) = \log p(\mathbf{x} | \boldsymbol{\theta}, \mathbf{d}) - \log p(\mathbf{x} | \mathbf{d}), \quad (21)$$

and maximise sample average of  $h_{\mathbf{d}}(\mathbf{x}, \boldsymbol{\theta})$  with respect to  $\mathbf{d}$

- ▶ Static setting: Kleingesse and Gutmann, AISTATS 2019
- ▶ Sequential setting where we update our belief about  $\boldsymbol{\theta}$  as we sequentially acquire the data: Kleingesse, Drovandi and Gutmann, Bayesian Analysis 2020

# Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \mathbb{E}_{p(\mathbf{x}, \theta | \mathbf{d})} \log \left[ \frac{p(\mathbf{x} | \theta, \mathbf{d})}{p(\mathbf{x} | \mathbf{d})} \right]$$

- ▶ Learning the ratio  $h_{\mathbf{d}}(\mathbf{x}, \theta)$  and approximating the MI is computationally costly.

# Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \mathbb{E}_{p(\mathbf{x}, \theta | \mathbf{d})} \log \left[ \frac{p(\mathbf{x} | \theta, \mathbf{d})}{p(\mathbf{x} | \mathbf{d})} \right]$$

- ▶ Learning the ratio  $h_{\mathbf{d}}(\mathbf{x}, \theta)$  and approximating the MI is computationally costly.
- ▶ But we do not need to estimate the MI accurately everywhere! Only around it's maximum.

## Experimental design for simulator models

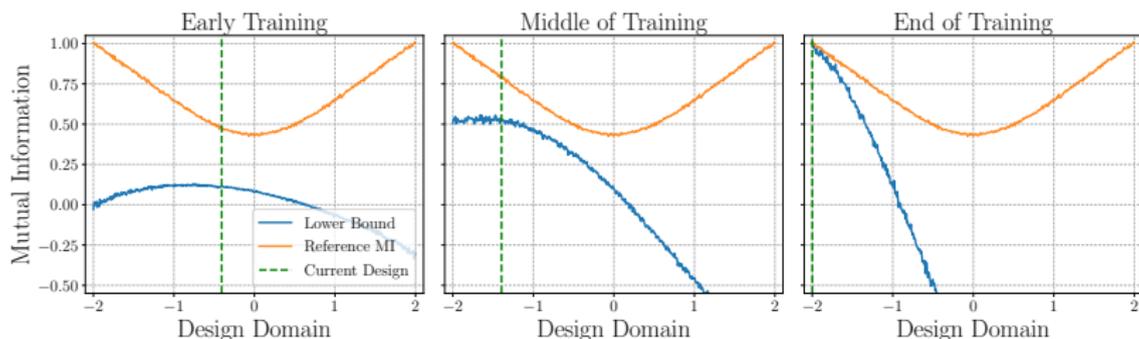
$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \mathbb{E}_{p(\mathbf{x}, \theta | \mathbf{d})} \log \left[ \frac{p(\mathbf{x} | \theta, \mathbf{d})}{p(\mathbf{x} | \mathbf{d})} \right]$$

- ▶ Learning the ratio  $h_{\mathbf{d}}(\mathbf{x}, \theta)$  and approximating the MI is computationally costly.
- ▶ But we do not need to estimate the MI accurately everywhere! Only around it's maximum.
- ▶ Let us use lower bounds on the MI (or proxy) where we concurrently tighten the bound and maximise the (proxy) MI!

# Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \mathbb{E}_{p(\mathbf{x}, \theta | \mathbf{d})} \log \left[ \frac{p(\mathbf{x} | \theta, \mathbf{d})}{p(\mathbf{x} | \mathbf{d})} \right]$$

- ▶ Learning the ratio  $h_{\mathbf{d}}(\mathbf{x}, \theta)$  and approximating the MI is computationally costly.
- ▶ But we do not need to estimate the MI accurately everywhere! Only around it's maximum.
- ▶ Let us use lower bounds on the MI (or proxy) where we concurrently tighten the bound and maximise the (proxy) MI!



(Kleinegesse and Gutmann, ICML 2020; arXiv:2105.04379)

## Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \operatorname{KL} (p(\boldsymbol{\theta}, \mathbf{x}|\mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})p(\mathbf{x}|\mathbf{d}))$$

- ▶ We can (again!) leverage logistic regression.

## Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \operatorname{KL} (p(\boldsymbol{\theta}, \mathbf{x}|\mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})p(\mathbf{x}|\mathbf{d}))$$

- ▶ We can (again!) leverage logistic regression.
- ▶ Logistic regression results in replacing the KL divergence with the JSD when measuring the MI.

$$\operatorname{JSD}(p, q; \mathbf{d}) \geq \log 2 - \frac{1}{2} \bar{J}(h; \mathbf{d}) \quad (22)$$

where  $h$  is the regression function and  $\bar{J}$  the logistic loss.

## Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \operatorname{KL} (p(\boldsymbol{\theta}, \mathbf{x}|\mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})p(\mathbf{x}|\mathbf{d}))$$

- ▶ We can (again!) leverage logistic regression.
- ▶ Logistic regression results in replacing the KL divergence with the JSD when measuring the MI.

$$\operatorname{JSD}(p, q; \mathbf{d}) \geq \log 2 - \frac{1}{2} \bar{J}(h; \mathbf{d}) \quad (22)$$

where  $h$  is the regression function and  $\bar{J}$  the logistic loss.

- ▶ Perform experimental design by maximising the negative logistic loss jointly with respect to  $h$  and  $\mathbf{d}$ .

## Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \operatorname{KL} (p(\boldsymbol{\theta}, \mathbf{x}|\mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})p(\mathbf{x}|\mathbf{d}))$$

- ▶ We can (again!) leverage logistic regression.
- ▶ Logistic regression results in replacing the KL divergence with the JSD when measuring the MI.

$$\operatorname{JSD}(p, q; \mathbf{d}) \geq \log 2 - \frac{1}{2} \bar{J}(h; \mathbf{d}) \quad (22)$$

where  $h$  is the regression function and  $\bar{J}$  the logistic loss.

- ▶ Perform experimental design by maximising the negative logistic loss jointly with respect to  $h$  and  $\mathbf{d}$ .
- ▶ Learned  $h$  provides an estimate of the posterior.

## Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \operatorname{KL} (p(\boldsymbol{\theta}, \mathbf{x}|\mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})p(\mathbf{x}|\mathbf{d}))$$

- ▶ We can (again!) leverage logistic regression.
- ▶ Logistic regression results in replacing the KL divergence with the JSD when measuring the MI.

$$\operatorname{JSD}(p, q; \mathbf{d}) \geq \log 2 - \frac{1}{2} \bar{J}(h; \mathbf{d}) \quad (22)$$

where  $h$  is the regression function and  $\bar{J}$  the logistic loss.

- ▶ Perform experimental design by maximising the negative logistic loss jointly with respect to  $h$  and  $\mathbf{d}$ .
- ▶ Learned  $h$  provides an estimate of the posterior.
- ▶ For more details and other loss functions:  
Kleinegesse and Gutmann, ICML 2020; arXiv:2105.04379

## Experimental design for simulator models

$$\hat{\mathbf{d}} = \operatorname{argmax}_{\mathbf{d}} \operatorname{KL} (p(\boldsymbol{\theta}, \mathbf{x}|\mathbf{d}) || p(\boldsymbol{\theta}|\mathbf{d})p(\mathbf{x}|\mathbf{d}))$$

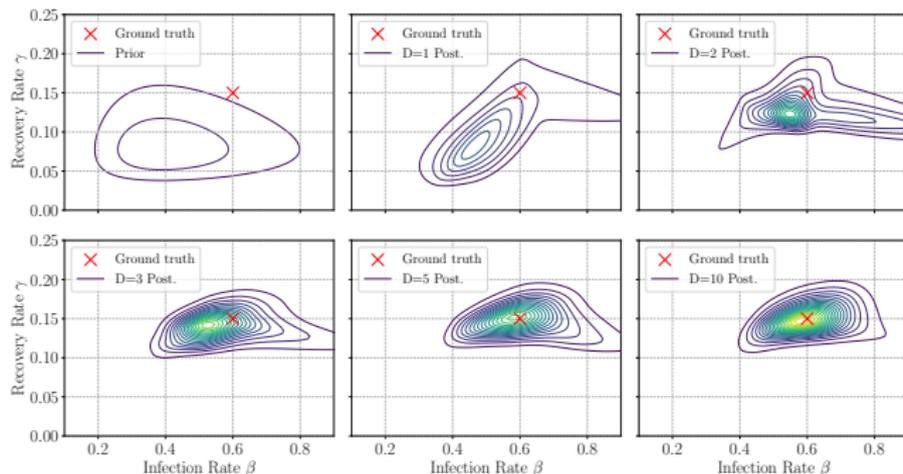
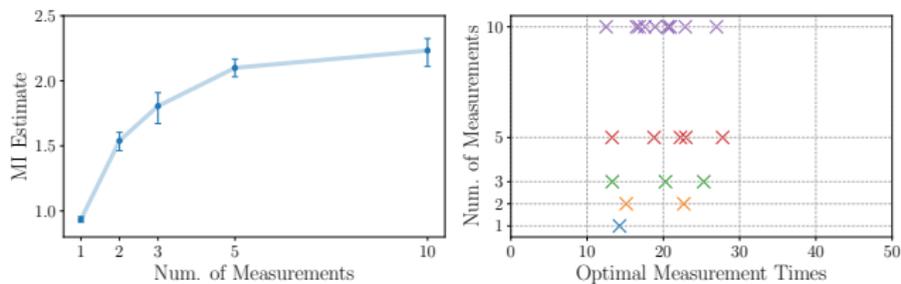
- ▶ We can (again!) leverage logistic regression.
- ▶ Logistic regression results in replacing the KL divergence with the JSD when measuring the MI.

$$\operatorname{JSD}(p, q; \mathbf{d}) \geq \log 2 - \frac{1}{2} \bar{J}(h; \mathbf{d}) \quad (22)$$

where  $h$  is the regression function and  $\bar{J}$  the logistic loss.

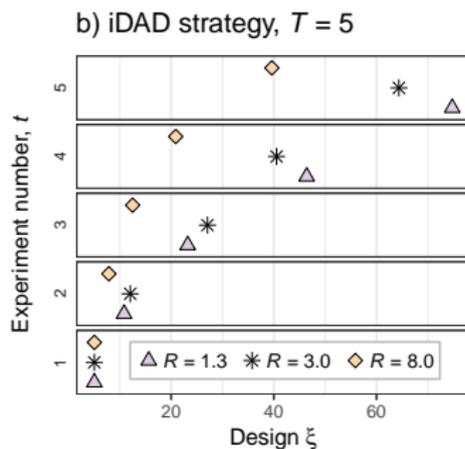
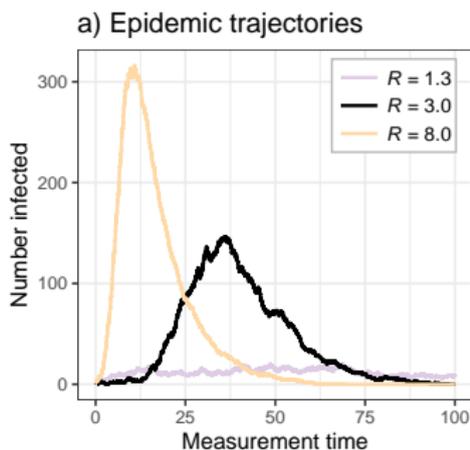
- ▶ Perform experimental design by maximising the negative logistic loss jointly with respect to  $h$  and  $\mathbf{d}$ .
- ▶ Learned  $h$  provides an estimate of the posterior.
- ▶ For more details and other loss functions:  
Kleinegesse and Gutmann, ICML 2020; arXiv:2105.04379
- ▶ For sequential setting: Ivanova et al, NeurIPS, 2021

# SIR example: static case (Kleinegge and Gutmann, ICML 2020)

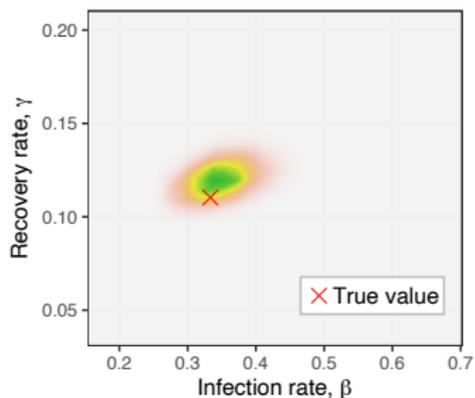


# SIR example: sequential case

(Ivanova et al, NeurIPS, 2021)



c) Estimated posterior,  $R = 3.0$



# Conclusions

## Research objective

- Two main goals: inference and experimental design

- Tasks are computationally intractable for simulator models

## Self-supervised learning to deal with intractability

- Link to logistic regression and Jensen-Shannon divergence

- Technical challenge: the density-chasm problem

## Application to Bayesian experimental design

- Via self-supervised learning of density ratios

- Exploiting bounds to increase computational efficiency

Thank you for your attention!

