

Efficient Statistical Inference for Intractable Models

Michael Gutmann

<http://homepages.inf.ed.ac.uk/mgutmann>

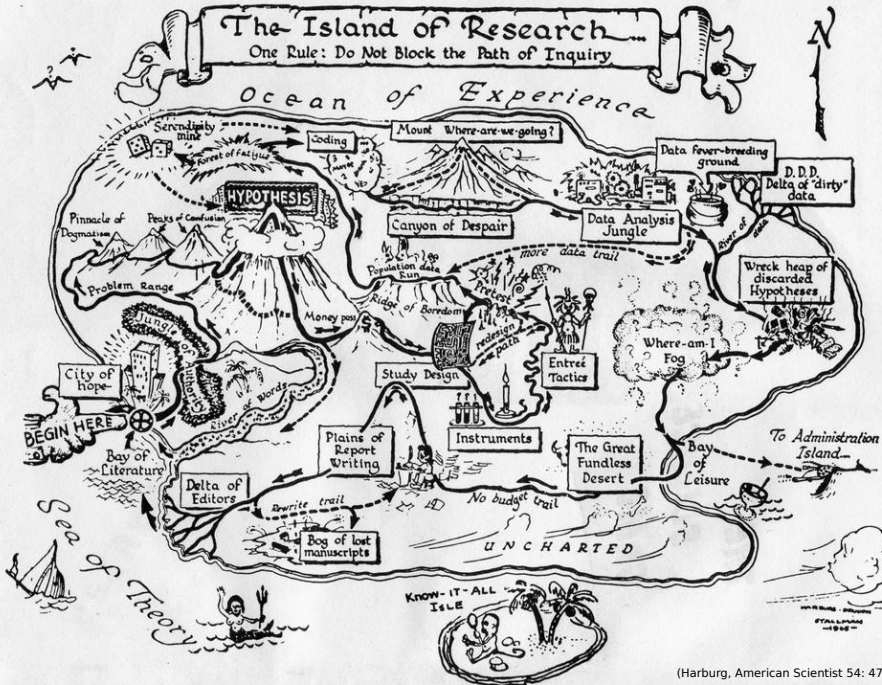
Institute for Adaptive and Neural Computation
School of Informatics, University of Edinburgh

31st August 2017

The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience



The Island of Research

One Rule: Do Not Block the Path of Inquiry

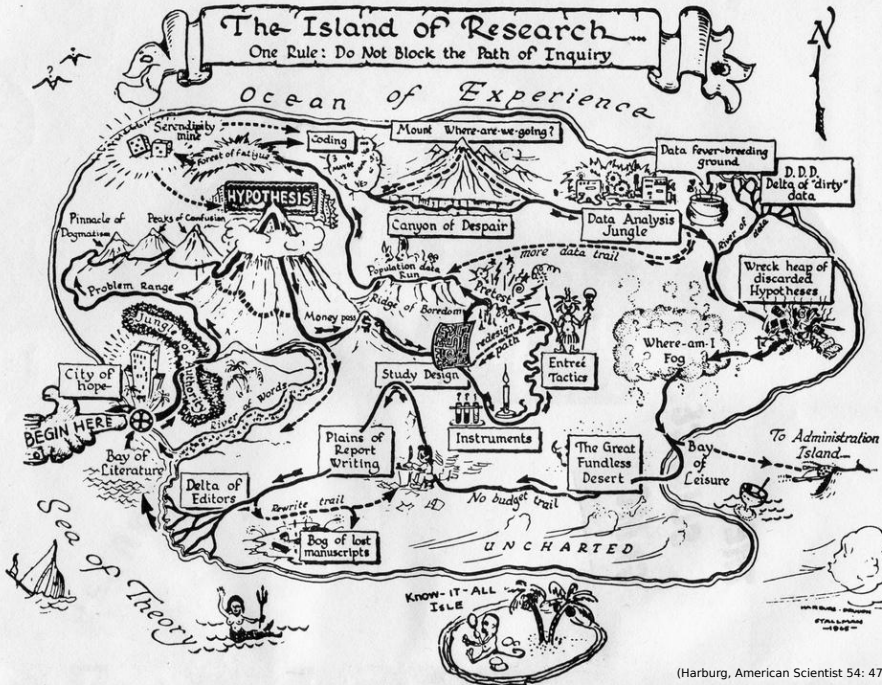
Ocean of Experience



The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience



The Island of Research

One Rule: Do Not Block the Path of Inquiry

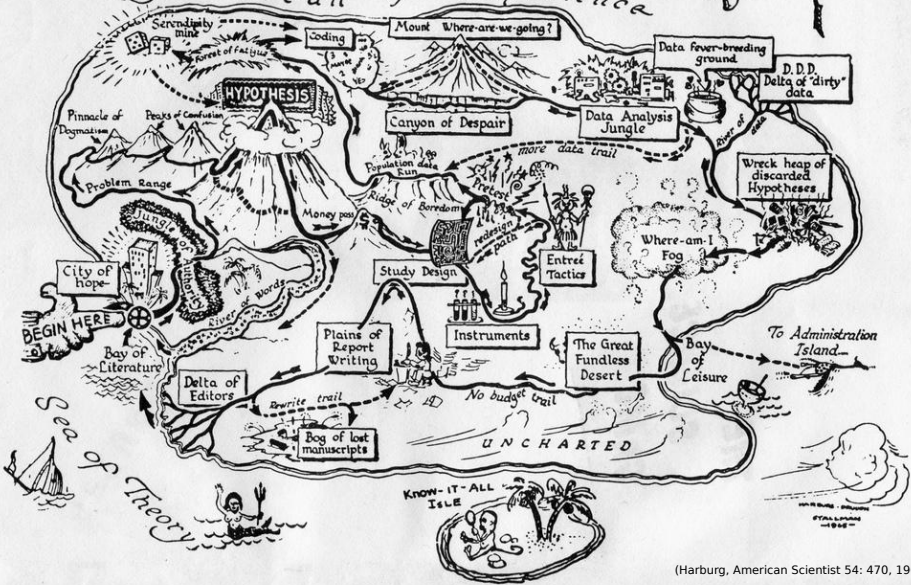
of Experience

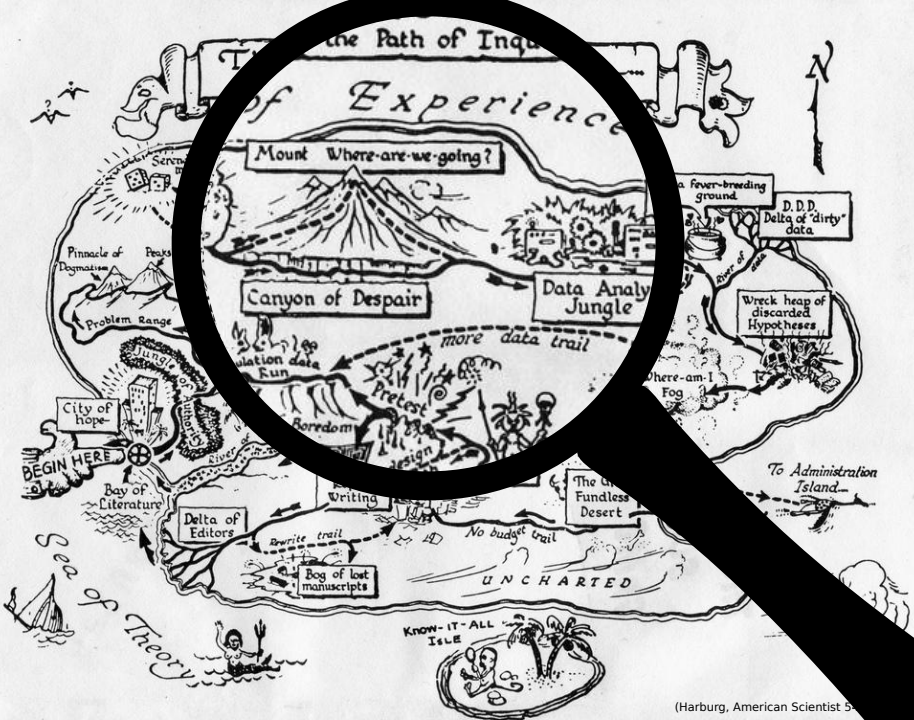


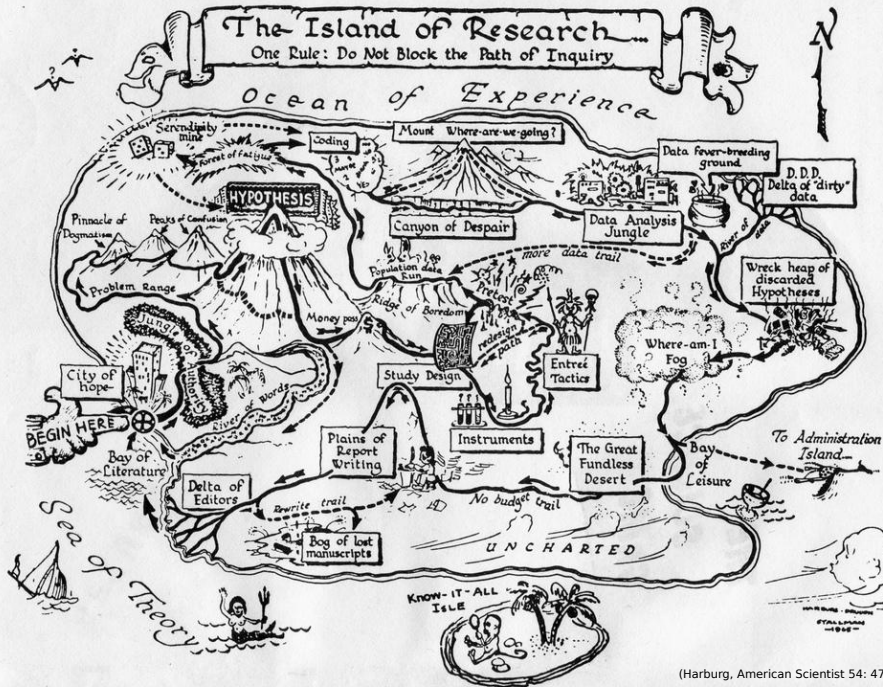
The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience







The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience



HYPOTHESIS

Mount Where-are-we-going?

Data fever-breeding ground

D.D.D. Delta of 'dirty' data

Canyon of Despair

Data Analysis Jungle

Wreck heap of discarded Hypotheses

Where-am-I Fog

Entree Tactics

The Great Fundless Desert

Bay of Leisure

To Administration Island

Delta of Editors

Plains of Report Writing

Instruments

No budget trail

Bog of lost manuscripts

UNCHARTED

KNOW-IT-ALL ISLE

STAY ALIVE - BRIGGS
1958

Sea of Theory

BEGIN HERE

Pinnacle of Dogmatism

Peaks of Confusion

Problem Range

Jungle of words

Bay of Literature

City of hope

River of words

Coding

Serenity mint

Forest of Fatigue

Population data Run

Ridge of Borrdam

Study Design

Protest

redesign path

more data trail

rewrite trail

3 dice

Money pass

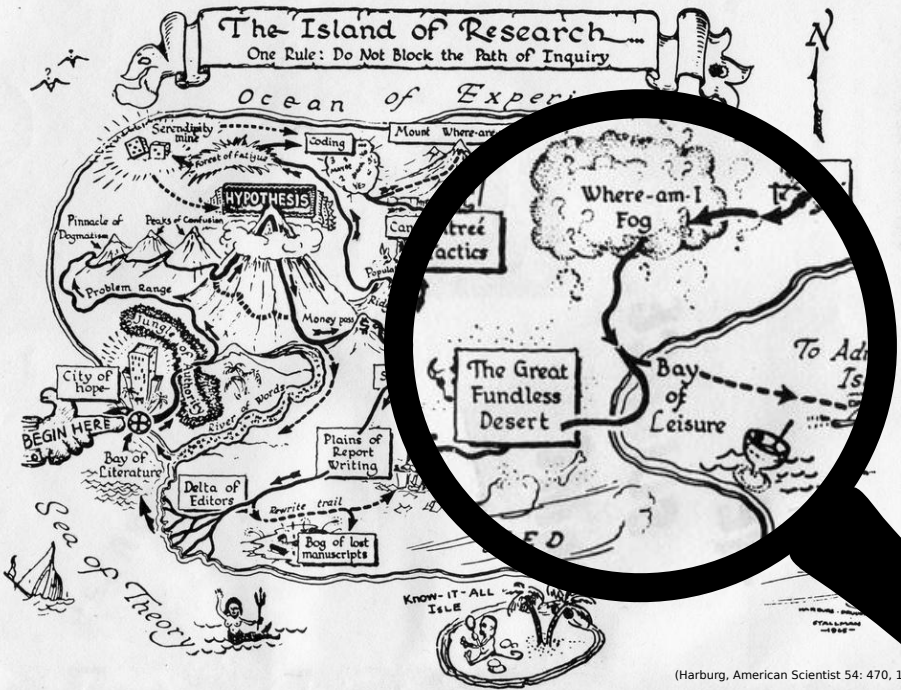
UNCLIMBABLE MOUNTAIN

UNCLIMBABLE MOUNTAIN

The Island of Research

One Rule: Do Not Block the Path of Inquiry

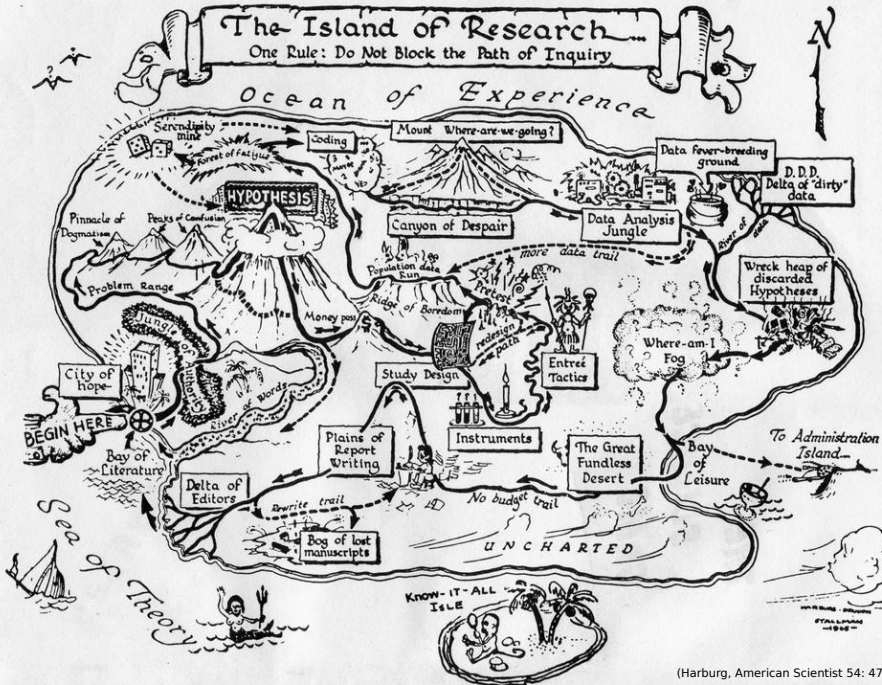
Ocean of Experience



The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience



The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience



The Island of Research

One Rule: Do Not Block the Path of Inquiry

Ocean of Experience



Progress in data science

- ▶ In the 60's, data science was very difficult.
- ▶ Today it's easier.

We have

- ▶ databases to store and access large amounts of data
- ▶ clusters to parallelise the computing
- ▶ the framework of statistical modelling and inference to provide the basic principles for analysing data.
- ▶ Challenge to further progress:
 - ▶ The basic principles do not take computational cost into account.
 - ▶ For complex data and models, exact inference is computationally impossible.
 - ▶ Good approximate solutions are needed.

We can use machine learning to perform highly efficient approximate inference for intractable models.

Introduction to statistical inference

- Likelihood function

- Case of exact inference

Models where exact inference is intractable

- Unnormalised models

- Generative models

Inference for unnormalised models

- Solution via logistic regression

- Application in unsupervised deep learning

Inference for generative models

- General overview

- Solution via logistic regression

Introduction to statistical inference

- Likelihood function

- Case of exact inference

Models where exact inference is intractable

- Unnormalised models

- Generative models

Inference for unnormalised models

- Solution via logistic regression

- Application in unsupervised deep learning

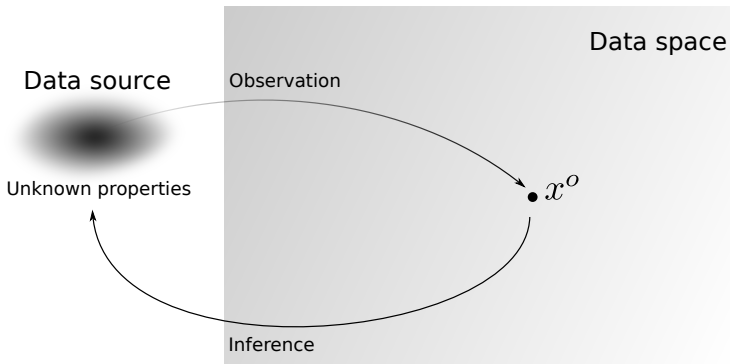
Inference for generative models

- General overview

- Solution via logistic regression

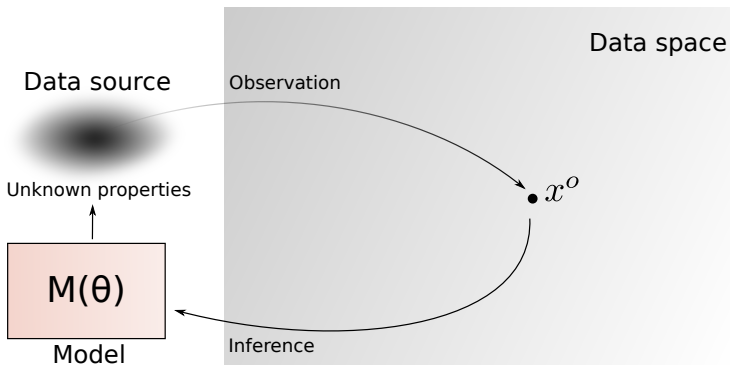
Goal of statistical inference

- ▶ Goal: Given data x^o , learn about properties of its source
- ▶ Enables decision making, predictions, ...



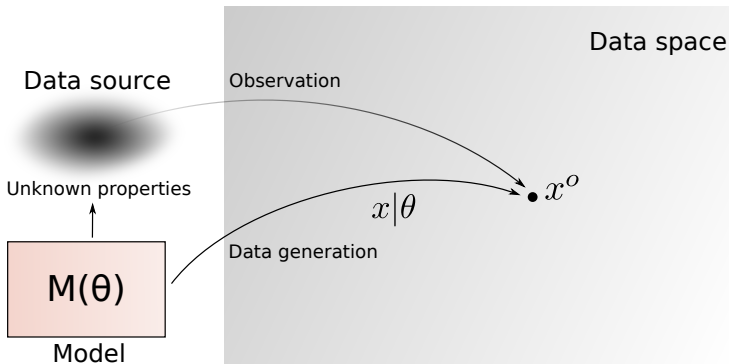
General approach

- ▶ Set up a model with potential properties θ (hypotheses)
- ▶ See which θ are in line with the observed data x^o



The likelihood function $L(\theta)$

- ▶ Measures agreement between θ and the observed data \mathbf{x}^o
- ▶ Probability to generate data like \mathbf{x}^o if hypothesis θ holds



Performing statistical inference

- ▶ If $L(\theta)$ is known, inference is straightforward
- ▶ Maximum likelihood estimation

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta) \quad (1)$$

- ▶ Bayesian inference

$$p(\theta | \mathbf{x}^o) \propto p(\theta) \times L(\theta) \quad (2)$$

posterior \propto prior \times likelihood

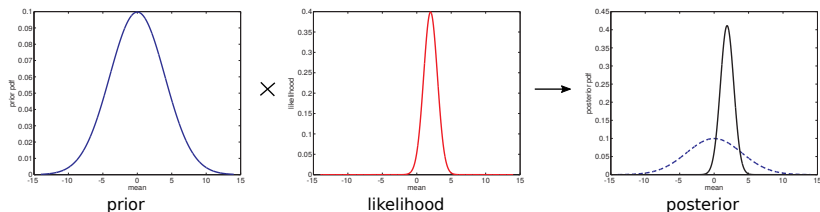
Allows us to learn from data by updating probabilities

Model specification

- ▶ Textbook: model \equiv family of probability density functions
- ▶ Probability density functions (pdfs) $p(\mathbf{x}|\theta)$ satisfy

$$\underbrace{p(\mathbf{x}|\theta) \geq 0}_{\text{non-negativity}} \quad \underbrace{\int p(\mathbf{x}|\theta) d\mathbf{x} = 1}_{\text{normalisation}} \quad (3)$$

- ▶ Likelihood function $L(\theta) \propto p(\mathbf{x}^o|\theta)$
- ▶ Closed form solutions are possible



Introduction to statistical inference

- Likelihood function

- Case of exact inference

Models where exact inference is intractable

- Unnormalised models

- Generative models

Inference for unnormalised models

- Solution via logistic regression

- Application in unsupervised deep learning

Inference for generative models

- General overview

- Solution via logistic regression

Intractable models I worked on

- ▶ Not all models are specified as family of pdfs $p(\mathbf{x}|\boldsymbol{\theta})$.
- ▶ I worked on
 1. Unnormalised models
 2. Generative models with unobserved variables
- ▶ The models are rather different, common point:

Multiple integrals needed to be computed to represent the models in terms of pdfs $p(\mathbf{x}|\boldsymbol{\theta})$.
- ▶ Solving the integrals exactly is computationally impossible. (curse of dimensionality)
 - ⇒ No model pdfs $p(\mathbf{x}|\boldsymbol{\theta})$
 - ⇒ No likelihood function $L(\boldsymbol{\theta}) \propto p(\mathbf{x}^o|\boldsymbol{\theta})$
 - ⇒ No exact inference

Unnormalised models

- ▶ Used for modelling
 - ▶ images (Markov random fields)
 - ▶ text (neural probabilistic language models)
 - ▶ social networks (exponential random graphs)
 - ▶ ...
- ▶ Specified via a non-negative function $\phi(\mathbf{x}|\boldsymbol{\theta}) \propto p(\mathbf{x}|\boldsymbol{\theta})$,

$$\int \dots \int \phi(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} = Z(\boldsymbol{\theta}) \neq 1 \quad p(\mathbf{x}|\boldsymbol{\theta}) = \frac{\phi(\mathbf{x}|\boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \quad (4)$$

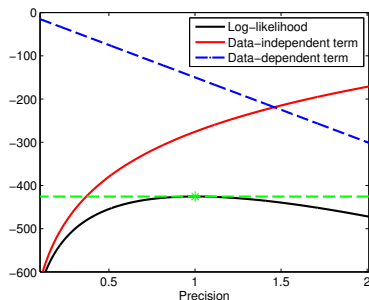
- ▶ Advantage: Specifying unnormalised models is often easier than specifying normalised models
- ▶ Disadvantage: Integral defining $Z(\boldsymbol{\theta})$, called the partition function, can generally not be computed.
 \Rightarrow Likelihood function is intractable.

Intractable partition function implies intractable likelihood

- ▶ Consider $p(x; \theta) = \frac{\phi(x; \theta)}{Z(\theta)} = \frac{\exp\left(-\theta \frac{x^2}{2}\right)}{\sqrt{2\pi/\theta}}$
- ▶ Log-likelihood function for precision $\theta \geq 0$

$$\ell(\theta) = -n \log \sqrt{\frac{2\pi}{\theta}} - \theta \sum_{i=1}^n \frac{x_i^2}{2} \quad (5)$$

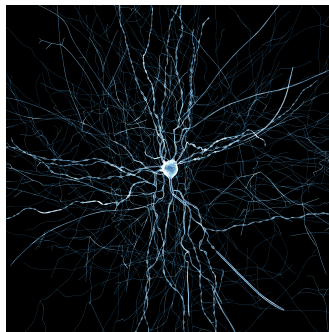
- ▶ Data-dependent (blue) and independent part (red) balance each other.
- ▶ If $Z(\theta)$ is intractable, $\ell(\theta)$ is intractable.



Generative models

- ▶ Models which specify a mechanism for generating data \mathbf{x}^o
 - ▶ e.g. stochastic dynamical systems
 - ▶ computer models / simulators of some complex biological process
 - ▶ aka: simulator-based models, implicit models, probabilistic programs

- ▶ Widely used
 - ▶ Evolutionary biology:
Simulating evolution
 - ▶ Neuroscience:
Simulating neural circuits
 - ▶ Health science:
Simulating the spread of an infectious disease



Generative models

- ▶ Advantage: detailed and realistic modelling
- ▶ Disadvantage: **likelihood function is generally intractable** due to unobserved variables.
- ▶ To compute $p(\mathbf{x}|\theta)$ one has to take into account all possible states of the unobserved variables (marginalisation)

$$p(\mathbf{x}|\theta) = \int \cdots \int p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} \quad (6)$$

- ▶ This is generally computationally impossible.

Introduction to statistical inference

- Likelihood function

- Case of exact inference

Models where exact inference is intractable

- Unnormalised models

- Generative models

Inference for unnormalised models

- Solution via logistic regression

- Application in unsupervised deep learning

Inference for generative models

- General overview

- Solution via logistic regression

Problem statement

- ▶ Task: Estimate the parameters θ of a parametric model $p(\cdot|\theta)$ of a d dimensional random vector \mathbf{x}
- ▶ Given: Data $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ (iid)
- ▶ Given: Unnormalised model $\phi(\cdot|\theta)$

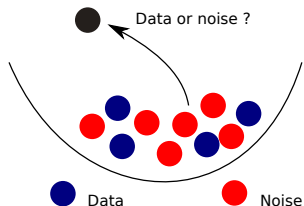
$$\int_{\xi} \phi(\xi|\theta) d\xi = Z(\theta) \neq 1 \quad p(\mathbf{x}|\theta) = \frac{\phi(\mathbf{x}|\theta)}{Z(\theta)} \quad (7)$$

Normalising partition function $Z(\theta)$ not known / computable.

Basic idea

- ▶ Formulate the estimation problem as a classification problem: observed data vs. auxiliary “noise” (with known properties)
- ▶ Successful classification \equiv learn the differences between the data and the noise
- ▶ differences + known noise properties \Rightarrow properties of the data

- ▶ Unsupervised learning by supervised learning
- ▶ We used (nonlinear) logistic regression for classification

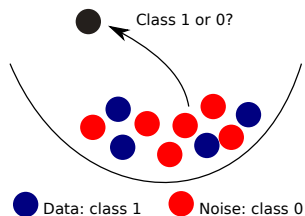


Logistic regression (1/2)

- ▶ Let $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_m)$ be a sample from a random variable \mathbf{y} with known (auxiliary) distribution p_{noise} .
- ▶ Introduce labels and form regression function:

$$P(C = 1 | \mathbf{u}; \boldsymbol{\theta}) = \frac{1}{1 + G(\mathbf{u}; \boldsymbol{\theta})} \quad G(\mathbf{u}; \boldsymbol{\theta}) \geq 0 \quad (8)$$

- ▶ Determine the parameters $\boldsymbol{\theta}$ such that $P(C = 1 | \mathbf{u}; \boldsymbol{\theta})$ is
 - ▶ large for most \mathbf{x}_i
 - ▶ small for most \mathbf{y}_i .



Logistic regression (2/2)

- ▶ Maximise (rescaled) conditional log-likelihood using the labelled data $\{(\mathbf{x}_1, 1), \dots, (\mathbf{x}_n, 1), (\mathbf{y}_1, 0), \dots, (\mathbf{y}_m, 0)\}$,

$$J_n^{\text{NCE}}(\boldsymbol{\theta}) = \frac{1}{n} \left(\sum_{i=1}^n \log P(C = 1 | \mathbf{x}_i; \boldsymbol{\theta}) + \sum_{i=1}^m \log [P(C = 0 | \mathbf{y}_i; \boldsymbol{\theta})] \right)$$

- ▶ For large sample sizes n and m , $\hat{\boldsymbol{\theta}}$ satisfying

$$G(\mathbf{u}; \hat{\boldsymbol{\theta}}) = \frac{m}{n} \frac{p_{\text{noise}}(\mathbf{u})}{p_{\text{data}}(\mathbf{u})} \quad (9)$$

is maximising $J_n^{\text{NCE}}(\boldsymbol{\theta})$. **Without any normalisation constraints.**

proof

Noise-contrastive estimation

(Gutmann and Hyvärinen, 2010; 2012)

(Gutmann and Hirayama, 2011)

- ▶ Assume unnormalised model $\phi(\cdot|\boldsymbol{\theta})$ is parametrised such that its scale can vary freely.

$$\boldsymbol{\theta} \rightarrow (\boldsymbol{\theta}; c) \qquad \phi(\mathbf{u}|\boldsymbol{\theta}) \rightarrow \exp(c)\phi(\mathbf{u}|\boldsymbol{\theta}) \qquad (10)$$

- ▶ Noise-contrastive estimation:

1. Choose p_{noise}
2. Generate auxiliary data \mathbf{Y}
3. Estimate $\boldsymbol{\theta}$ via logistic regression with

$$G(\mathbf{u}; \boldsymbol{\theta}) = \frac{m}{n} \frac{p_{\text{noise}}(\mathbf{u})}{\phi(\mathbf{u}|\boldsymbol{\theta})}. \qquad (11)$$

Noise-contrastive estimation

(Gutmann and Hyvärinen, 2010; 2012)

(Gutmann and Hirayama, 2011)

- ▶ Assume unnormalised model $\phi(\cdot|\theta)$ is parametrised such that its scale can vary freely.

$$\theta \rightarrow (\theta; c) \quad \phi(\mathbf{u}|\theta) \rightarrow \exp(c)\phi(\mathbf{u}|\theta) \quad (10)$$

- ▶ Noise-contrastive estimation:

1. Choose p_{noise}
2. Generate auxiliary data \mathbf{Y}
3. Estimate θ via logistic regression with

$$G(\mathbf{u}; \theta) = \frac{m}{n} \frac{p_{\text{noise}}(\mathbf{u})}{\phi(\mathbf{u}|\theta)}. \quad (11)$$

- ▶ $G(\mathbf{u}; \theta) \rightarrow \frac{m}{n} \frac{p_{\text{noise}}(\mathbf{u})}{p_{\text{data}}(\mathbf{u})} \Rightarrow \phi(\mathbf{u}|\theta) \rightarrow p_{\text{data}}(\mathbf{u})$

Example

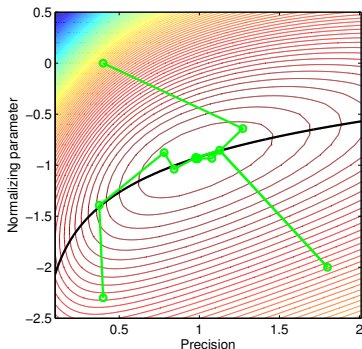
- ▶ Unnormalised Gaussian:

$$\phi(x; \theta) = \exp(\theta_2) \exp\left(-\theta_1 \frac{x^2}{2}\right), \quad \theta_1 > 0, \theta_2 \in \mathbb{R}, \quad (12)$$

- ▶ Parameters: θ_1 (precision), $\theta_2 \equiv c$ (scaling parameter)

Contour plot of $J_n^{\text{NCE}}(\theta)$:

- ▶ Gaussian noise with $\nu = m/n = 10$
- ▶ True precision $\theta_1^* = 1$
- ▶ Black: normalised models
- ▶ Green: optimisation paths



(Gutmann and Hyvärinen, 2012)

- ▶ Assume $p_{\text{data}} = p(\cdot | \theta^*)$
- ▶ Consistency: As n increases,

$$\hat{\theta}_n = \operatorname{argmax}_{\theta} J_n^{\text{NCE}}(\theta), \quad (13)$$

converges in probability to θ^* .

- ▶ Efficiency: As $\nu = m/n$ increases, for any valid choice of p_{noise} , noise-contrastive estimation tends to “perform as well” as MLE (it is asymptotically Fisher efficient).

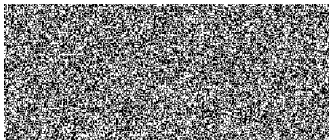
Application examples

- ▶ Models of text: e.g. Mnih and Teh, 2012, *A fast and simple algorithm for training neural probabilistic language models*
- ▶ Models of images: e.g. Gutmann and Hyvärinen, 2013, *A three-layer model of natural image statistics*
- ▶ Machine translation: e.g. Zoph et al, 2016, *Simple, fast noise-contrastive estimation for large RNN vocabularies*
- ▶ Product recommendation: e.g. Tschitschek et al, 2016, *Learning probabilistic submodular diversity models via noise contrastive estimation*

Unsupervised deep learning on natural images

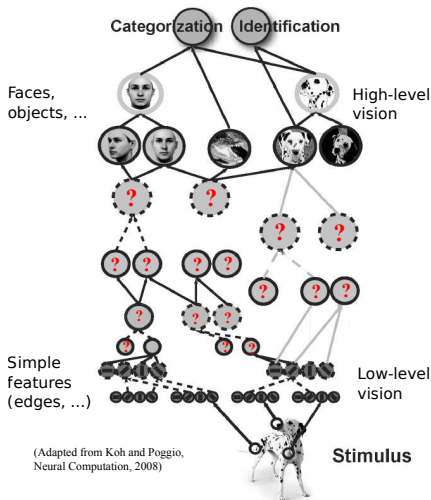


- ▶ Natural images \equiv images which we see in our environment
- ▶ Understanding their properties is important
 - ▶ for modern image processing
 - ▶ for understanding biological visual systems



Unsupervised deep learning on natural images

- ▶ Rapid object recognition by feedforward processing
- ▶ Computations in middle layers poorly understood
- ▶ Our approach: learn the computations from data
- ▶ Idea: the units indicate how probable an input image is. (up to normalisation)

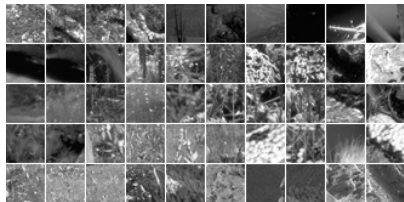


(Gutmann and Hyvärinen, 2013)

Image data

Consider two kinds of image data:

1. Image patches of size 32 by 32, extracted from larger images (left).
2. “Tiny images” dataset, converted to grey scale: complete scenes downsampled to 32 by 32 images (right)
(Torralba et al, TPAMI 2008)



Multi-layer model

- ▶ Let \mathbf{I} be a vectorised image. Processing layers:

$$\mathbf{x} = \text{gain control}(\mathbf{I})$$

$$y_i^{(1)} = \max(\mathbf{w}_i^{(1)} \cdot \mathbf{x}, 0), \quad i = 1 \dots 600$$

$$y_i^{(2)} = \log(\mathbf{w}_i^{(2)} \cdot (\mathbf{y}^{(1)})^2 + 1), \quad i = 1 \dots 100$$

$$\mathbf{z}^{(2)} = \text{gain control}(\mathbf{y}^{(2)})$$

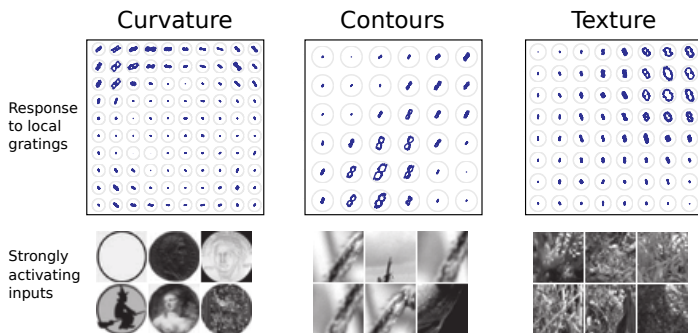
$$y_i^{(3)} = \max(\mathbf{w}_i^{(3)} \cdot \mathbf{z}^{(2)}, 0), \quad i = 1 \dots 50$$

Gain control: centring, normalising the norm after whitening, possibly dimension reduction

- ▶ The outputs $y_i^{(3)}$ define how probable an input image is.
(up to normalisation \Rightarrow unnormalised model)
- ▶ The weights are the parameters to be learned ($> 2 \cdot 10^5$ parameters)
- ▶ Only constraint: $w_{ki}^{(2)} \geq 0$.

Learned features

- ▶ 1st layer: \approx local Fourier transform (Gabor filters)
- ▶ 2nd layer: local max-pooling
- ▶ 3rd layer: emergence of units sensitive to curvature, longer contours, and texture
- ▶ Close link to neural processing in the visual cortex



(Gutmann and Hyvärinen, 2013)

Introduction to statistical inference

- Likelihood function

- Case of exact inference

Models where exact inference is intractable

- Unnormalised models

- Generative models

Inference for unnormalised models

- Solution via logistic regression

- Application in unsupervised deep learning

Inference for generative models

- General overview

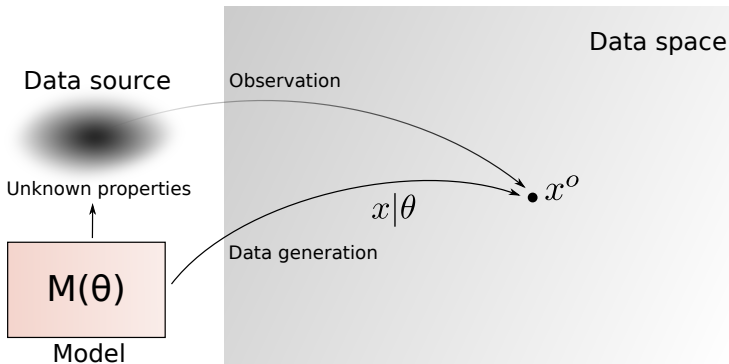
- Solution via logistic regression

Perform Bayesian inference for models where

1. the likelihood function is too costly to compute
2. sampling – simulating data – from the model is possible

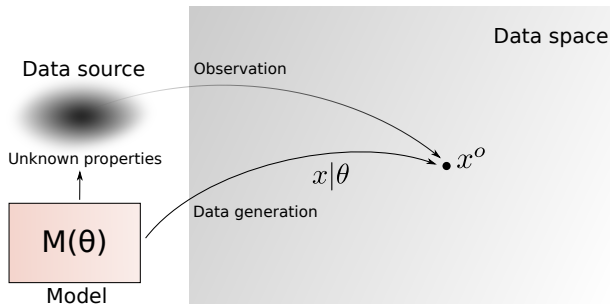
The likelihood function $L(\theta)$

- ▶ Probability that the model generates data like \mathbf{x}^o when using parameter value θ
- ▶ Generally well defined but intractable for simulator-based models



Three foundational issues

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^o$?
2. How should we compute the probability of the event $\mathbf{x}_\theta \equiv \mathbf{x}^o$?
3. For which values of θ should we compute it?



Likelihood: Probability that the model generates data like \mathbf{x}^o for parameter value θ

Approximate Bayesian computation

Recent review: Lintusaari et al (2017) “Fundamentals and recent developments in approximate Bayesian computation”, *Systematic Biology*

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ Check whether $\|T(\mathbf{x}_\theta) - T(\mathbf{x}^\circ)\| \leq \epsilon$
2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ By counting
3. For which values of θ should we compute it?
 - ⇒ Sample from the prior (or other proposal distributions)

Approximate Bayesian computation

Recent review: Lintusaari et al (2017) “Fundamentals and recent developments in approximate Bayesian computation”, Systematic Biology

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
⇒ Check whether $\|T(\mathbf{x}_\theta) - T(\mathbf{x}^\circ)\| \leq \epsilon$
2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
⇒ By counting
3. For which values of θ should we compute it?
⇒ Sample from the prior (or other proposal distributions)

Difficulties:

- ▶ Choice of summary statistics $T()$ and threshold ϵ
- ▶ Typically high computational cost

Synthetic likelihood

(Simon Wood, Nature, 2010)

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ Compute summary statistics $\mathbf{t}_\theta = T(\mathbf{x}_\theta)$
 - ⇒ Model their distribution as a Gaussian
 - ⇒ Compute likelihood function with $T(\mathbf{x}^\circ)$ as observed data
3. For which values of θ should we compute it?
 - ⇒ Use obtained “synthetic” likelihood function as part of a Monte Carlo method

Synthetic likelihood

(Simon Wood, Nature, 2010)

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^o$?
2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^o$?
 - ⇒ Compute summary statistics $\mathbf{t}_\theta = T(\mathbf{x}_\theta)$
 - ⇒ Model their distribution as a Gaussian
 - ⇒ Compute likelihood function with $T(\mathbf{x}^o)$ as observed data
3. For which values of θ should we compute it?
 - ⇒ Use obtained “synthetic” likelihood function as part of a Monte Carlo method

Difficulties:

- ▶ Choice of summary statistics $T()$
- ▶ Gaussianity assumption may not hold
- ▶ Typically high computational cost

Overview of some of my work

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ Use classification (Gutmann et al, 2014, 2017)
 2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 3. For which values of θ should we compute it?
 - ⇒ Use Bayesian optimisation (Gutmann and Corander, 2013-2016)
Compared to standard approaches: speed-up by a factor of 1000 more
-
1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ Use density ratio estimation / logistic regression
(Dutta et al, 2016, arXiv:1611.10242)

Overview of some of my work

1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ Use classification (Gutmann et al, 2014, 2017)
 2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 3. For which values of θ should we compute it?
 - ⇒ Use Bayesian optimisation (Gutmann and Corander, 2013-2016)
Compared to standard approaches: speed-up by a factor of 1000 more
-
1. How should we assess whether $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 2. How should we compute the proba of the event $\mathbf{x}_\theta \equiv \mathbf{x}^\circ$?
 - ⇒ Use density ratio estimation / logistic regression
(Dutta et al, 2016, arXiv:1611.10242)

(Dutta et al, 2016, arXiv:1611.10242)

- ▶ Frame posterior estimation as ratio estimation problem

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

- ▶ Estimating $r(\mathbf{x}, \boldsymbol{\theta})$ is the difficult part since $p(\mathbf{x}|\boldsymbol{\theta})$ unknown.
- ▶ Estimate $\hat{r}(\mathbf{x}, \boldsymbol{\theta})$ yields estimate of the likelihood function and posterior

$$\hat{L}(\boldsymbol{\theta}) \propto \hat{r}(\mathbf{x}^o, \boldsymbol{\theta}), \quad \hat{p}(\boldsymbol{\theta}|\mathbf{x}^o) = p(\boldsymbol{\theta})\hat{r}(\mathbf{x}^o, \boldsymbol{\theta}). \quad (14)$$

(Dutta et al, 2016, arXiv:1611.10242)

- ▶ Frame posterior estimation as ratio estimation problem

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

- ▶ Estimating $r(\mathbf{x}, \boldsymbol{\theta})$ is the difficult part since $p(\mathbf{x}|\boldsymbol{\theta})$ unknown.
- ▶ Estimate $\hat{r}(\mathbf{x}, \boldsymbol{\theta})$ yields estimate of the likelihood function and posterior

$$\hat{L}(\boldsymbol{\theta}) \propto \hat{r}(\mathbf{x}^\circ, \boldsymbol{\theta}), \quad \hat{p}(\boldsymbol{\theta}|\mathbf{x}^\circ) = p(\boldsymbol{\theta})\hat{r}(\mathbf{x}^\circ, \boldsymbol{\theta}). \quad (14)$$

- ▶ Often more practical to estimate log-ratio $h(\mathbf{x}, \boldsymbol{\theta}) = \log r(\mathbf{x}, \boldsymbol{\theta})$

$$\hat{L}(\boldsymbol{\theta}) \propto \exp(\hat{h}(\mathbf{x}^\circ, \boldsymbol{\theta})), \quad \hat{p}(\boldsymbol{\theta}|\mathbf{x}^\circ) = p(\boldsymbol{\theta}) \exp(\hat{h}(\mathbf{x}^\circ, \boldsymbol{\theta})) \quad (15)$$

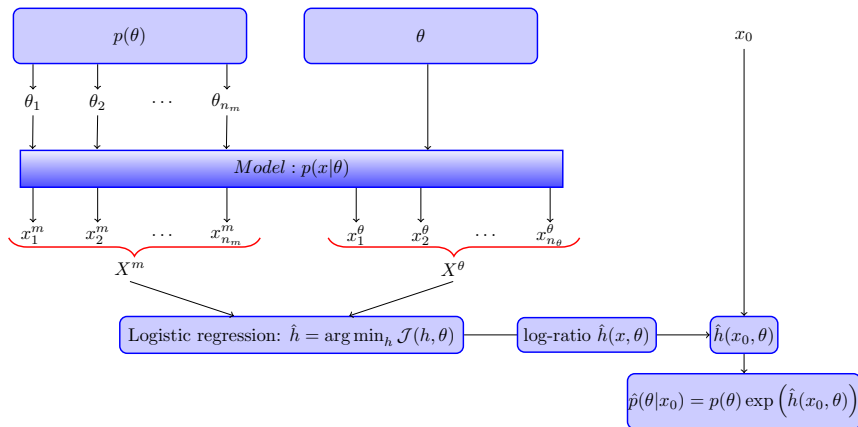
Estimating the posterior

- ▶ From theory of noise-contrastive estimation: ratio $r(\mathbf{x}, \boldsymbol{\theta})$, or log-ratio $h(\mathbf{x}, \boldsymbol{\theta})$ can be estimated by logistic regression
- ▶ Formulate classification problem with
 - ▶ one class: data sampled from $p(\mathbf{x}|\boldsymbol{\theta})$
 - ▶ other class: data sampled from marginal $p(\mathbf{x})$
- ▶ Logistic regression gives (point-wise in $\boldsymbol{\theta}$)

$$\hat{h}(\mathbf{x}, \boldsymbol{\theta}) \rightarrow \log \frac{p(\mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{x})} = \log r(\mathbf{x}, \boldsymbol{\theta}) \quad (16)$$

- ▶ We operate on synthetic data only; can generate as much data as we wish

Estimating the posterior



(Dutta et al, 2016, arXiv:1611.10242)

Auxiliary model

- ▶ We need to specify a model for h .
- ▶ For simplicity: linear model

$$h(\mathbf{x}) = \sum_{i=1}^b \beta_i \psi_i(\mathbf{x}) = \beta^\top \psi(\mathbf{x}) \quad (17)$$

where $\psi_i(\mathbf{x})$ are summary statistics

- ▶ More complex models possible
- ▶ Simple linear model leads to a generalisation of synthetic likelihood (Dutta et al, 2016, arXiv:1611.10242)
- ▶ L_1 penalty on β for weighing and selecting summary statistics

Application to ARCH model

- ▶ Model:

$$x^{(t)} = \theta_1 x^{(t-1)} + e^{(t)} \quad (18)$$

$$e^{(t)} = \xi^{(t)} \sqrt{0.2 + \theta_2 (e^{(t-1)})^2} \quad (19)$$

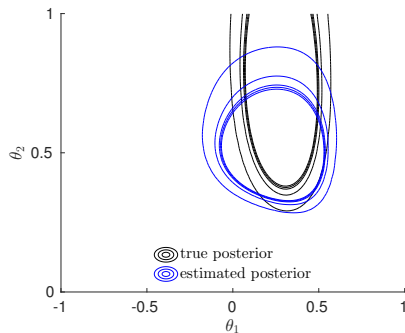
$\xi^{(t)}$ and $e^{(0)}$ independent standard normal r.v., $x^{(0)} = 0$

- ▶ 100 time points
- ▶ Parameters: $\theta_1 \in (-1, 1)$, $\theta_2 \in (0, 1)$
- ▶ Uniform prior on θ_1, θ_2

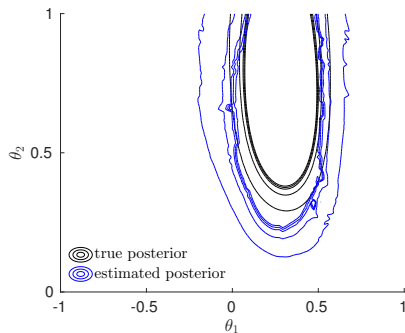
Application to ARCH model

- ▶ Summary statistics $\psi_i(\mathbf{x})$:
 - ▶ auto-correlations with lag one to five
 - ▶ all (unique) pairwise combinations of them
 - ▶ a constant
- ▶ To check robustness: 50% irrelevant summary statistics (drawn from standard normal)
- ▶ Comparison with synthetic likelihood with equivalent set of summary statistics (relevant sum. stats. only)

Example posterior

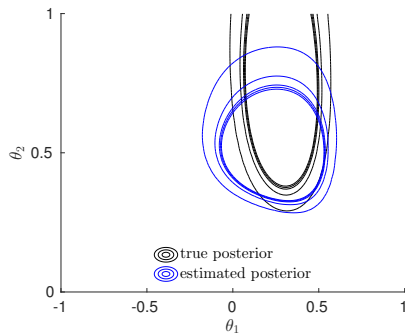


(a) synthetic likelihood

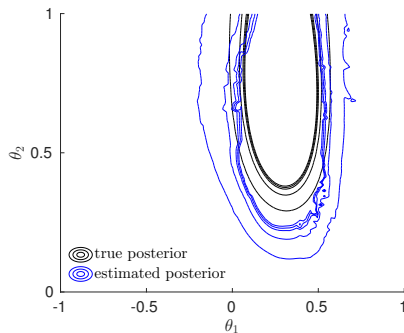


(b) proposed method

Example posterior



(c) synthetic likelihood

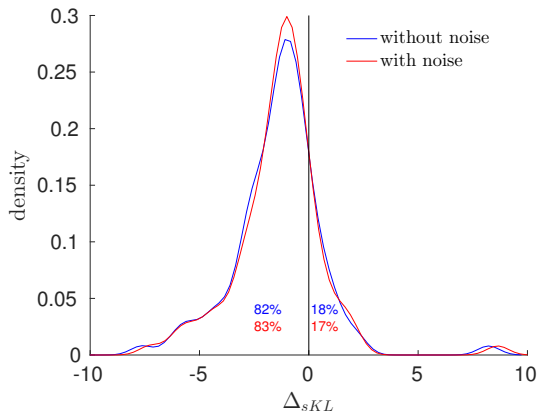


(d) proposed method subject to noise

Systematic analysis

- ▶ Symmetrised Kullback-Leibler divergence between estimated and true posterior
- ▶ Point-wise comparison with synthetic likelihood (100 data sets)

$\Delta_{sKL} = \text{SKL for proposed method} - \text{SKL for synthetic likelihood}$



Key results

For details, see arXiv:1611.10242v1

- ▶ Frame the problem of Bayesian inference with intractable generative models as ratio estimation problem
- ▶ Use logistic regression to solve the problem
- ▶ Approach includes synthetic likelihood as special case
- ▶ For **same summary statistics**, typically **more accurate inferences** than the synthetic likelihood
- ▶ Robustness to irrelevant summary statistics thanks to regularisation
- ▶ Enables selection of relevant summary statistics
- ▶ No threshold to choose (unlike in ABC)

Conclusions

- ▶ Statistical modelling and inference are part of the foundations of data science.
 - ▶ They are not concerned with computational cost.
 - ▶ Exact inference is impossible for complex models.
- ▶ Unnormalised models
 - ▶ Noise-contrastive estimation
 - ▶ Formulated the inference problem as a classification problem
- ▶ Generative models
 - ▶ General overview
 - ▶ Formulated the inference problem as a classification problem

Conclusions

- ▶ Statistical modelling and inference are part of the foundations of data science.
 - ▶ They are not concerned with computational cost.
 - ▶ Exact inference is impossible for complex models.
- ▶ Unnormalised models
 - ▶ Noise-contrastive estimation
 - ▶ Formulated the inference problem as a classification problem
- ▶ Generative models
 - ▶ General overview
 - ▶ Formulated the inference problem as a classification problem

*By re-framing inference problems,
we can use machine learning to perform highly efficient
approximate inference for intractable models.*

Maximiser of the NCE objective function

Maximiser of the NCE objective function

Proof of Equation (9)

For large sample sizes n and m , $\hat{\theta}$ satisfying

$$G(\mathbf{u}; \hat{\theta}) = \frac{m p_{\text{noise}}(\mathbf{u})}{n p_{\text{data}}(\mathbf{u})}$$

is maximising $J_n^{\text{NCE}}(\theta)$,

$$J_n^{\text{NCE}}(\theta) = \frac{1}{n} \left(\sum_{i=1}^n \log P(C = 1 | \mathbf{x}_i; \theta) + \sum_{i=1}^m \log [P(C = 0 | \mathbf{y}_i; \theta)] \right)$$

without any normalisation constraints.

Proof of Equation (9)

$$\begin{aligned} J_n^{\text{NCE}}(\boldsymbol{\theta}) &= \frac{1}{n} \left(\sum_{i=1}^n \log P(C = 1 | \mathbf{x}_i; \boldsymbol{\theta}) + \sum_{i=1}^m \log [P(C = 0 | \mathbf{y}_i; \boldsymbol{\theta})] \right) \\ &= \frac{1}{n} \sum_{t=1}^n \log P(C = 1 | \mathbf{x}_i; \boldsymbol{\theta}) + \frac{m}{n} \frac{1}{m} \sum_{t=1}^m \log [P(C = 0 | \mathbf{y}_i; \boldsymbol{\theta})] \end{aligned}$$

Fix the ratio $m/n = \nu$ and let $n \rightarrow \infty$ and $m \rightarrow \infty$. By law of large numbers, J_n^{NCE} converges to J^{NCE} ,

$$J^{\text{NCE}}(\boldsymbol{\theta}) = E_{\mathbf{x}} (\log P(C = 1 | \mathbf{x}; \boldsymbol{\theta})) + \nu E_{\mathbf{y}} (\log P(C = 0 | \mathbf{y}; \boldsymbol{\theta})) \quad (20)$$

With $P(C = 1 | \mathbf{x}; \boldsymbol{\theta}) = \frac{1}{1+G(\mathbf{x}; \boldsymbol{\theta})}$ and $P(C = 0 | \mathbf{y}; \boldsymbol{\theta}) = \frac{G(\mathbf{y}; \boldsymbol{\theta})}{1+G(\mathbf{y}; \boldsymbol{\theta})}$ we have

$$\begin{aligned} J^{\text{NCE}}(\boldsymbol{\theta}) &= - E_{\mathbf{x}} \log(1 + G(\mathbf{x}; \boldsymbol{\theta})) + \nu E_{\mathbf{y}} \log G(\mathbf{y}; \boldsymbol{\theta}) - \\ &\quad \nu E_{\mathbf{y}} \log(1 + G(\mathbf{y}; \boldsymbol{\theta})) \end{aligned} \quad (21)$$

Consider the objective $J^{\text{NCE}}(\theta)$ as a function of $H = \log G$ rather than θ ,

$$\begin{aligned} \mathcal{J}^{\text{NCE}}(H) &= -E_{\mathbf{x}} \log(1 + \exp H(\mathbf{x})) + \nu E_{\mathbf{y}} H(\mathbf{y}) - \nu E_{\mathbf{y}} \log(1 + \exp H(\mathbf{y})) \\ &= - \int p_{\text{data}}(\xi) \log(1 + \exp H(\xi)) d\xi + \nu \int p_{\text{noise}}(\xi) H(\xi) d\xi \\ &\quad - \nu \int p_{\text{noise}}(\xi) \log(1 + \exp H(\xi)) d\xi \\ &= - \int (p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)) \log(1 + \exp H(\xi)) d\xi + \\ &\quad \nu \int p_{\text{noise}}(\xi) H(\xi) d\xi \end{aligned}$$

We now expand $\mathcal{J}^{\text{NCE}}(H + \epsilon q)$ around H for an arbitrary function q and a small scalar ϵ .

With

$$\begin{aligned}\log(1 + \exp [H(\xi) + \epsilon q(\xi)]) &= \log(1 + \exp H(\xi)) + \frac{\epsilon q(\xi)}{1 + \exp(-H(\xi))} \\ &+ \frac{\epsilon^2}{2} \frac{q(\xi)}{1 + \exp(-H(\xi))} \frac{q(\xi)}{1 + \exp(H(\xi))} \\ &+ O(\epsilon^3)\end{aligned}$$

we have

$$\begin{aligned}\mathcal{J}^{\text{NCE}}(H + \epsilon q) &= - \int (p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)) \log(1 + \exp H(\xi)) d\xi \\ &- \epsilon \int \frac{p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)}{1 + \exp(-H(\xi))} q(\xi) d\xi \\ &- \frac{\epsilon^2}{2} \int \frac{p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)}{1 + \exp(-H(\xi))} \frac{q(\xi)^2}{1 + \exp(H(\xi))} d\xi \\ &+ \nu \int p_{\text{noise}}(\xi) H(\xi) d\xi + \epsilon \nu \int p_{\text{noise}}(\xi) q(\xi) d\xi + O(\epsilon^3)\end{aligned}$$

Collecting terms gives:

$$\begin{aligned} \mathcal{J}^{\text{NCE}}(H + \epsilon q) &= \mathcal{J}^{\text{NCE}}(H) + \\ &\epsilon \int \left(\nu p_{\text{noise}}(\xi) - \frac{p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)}{1 + \exp(-H(\xi))} \right) q(\xi) d\xi \\ &\quad - \frac{\epsilon^2}{2} \int \frac{p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)}{1 + \exp(-H(\xi))} \frac{q(\xi)^2}{1 + \exp(H(\xi))} d\xi + O(\epsilon^3) \end{aligned}$$

The second-order term is negative for all (non-trivial) q and H .

The first-order term is zero for all q if and only if

$$\nu p_{\text{noise}}(\xi) = \frac{p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)}{1 + \exp(-H^*(\xi))}$$

$$\nu p_{\text{noise}}(\xi) + \nu p_{\text{noise}}(\xi) \exp(-H^*(\xi)) = p_{\text{data}}(\xi) + \nu p_{\text{noise}}(\xi)$$

$$\exp(-H^*(\xi)) = \frac{p_{\text{data}}(\xi)}{\nu p_{\text{noise}}(\xi)}$$

which shows that $\hat{\theta}$ such that $G(\xi; \hat{\theta}) = \exp(H^*(\xi)) = \nu \frac{p_{\text{noise}}}{p_{\text{data}}}$ is maximising $\mathcal{J}^{\text{NCE}}(\theta)$.

[back](#)