

Estimation in Unnormalised Models

Z (👿)

Maximum Likelihood



- write down likelihood as function of θ
- (try to) maximise this function
- easy? (conceptually, in practice)
- hard? (non-standard models)
- variants (penalised, Bayes, ...)

$\hat{\theta}_{MLE}$

Cramér-Rao, Fisher Efficiency



- consider asymptotic (co)variance of estimator
- among unbiased estimators, can't beat MLE
- same for biased (more conditions)
- open and shut case? (asymptotically)



Quick proof of Cramér-Rao

$$\hat{a} = \nabla_{\theta} \langle a, \theta \rangle$$

$$= \nabla_{\theta} \int P(x|\theta) \langle a, T(x) \rangle dx$$

$$= \int P(x|\theta) \langle a, T(x) \rangle \nabla_{\theta} \log P(x|\theta) dx$$


$$= \left(\int P(x|\theta) \nabla_{\theta} \log P(x|\theta) T(x)^{\top} dx \right) a$$

$$\Rightarrow I = \text{Cov}(\nabla_{\theta} \log P(x|\theta), T(x)) \Rightarrow \text{apply CS}$$

Challenges of MLE

- tractability (convexity, computability, ...)
- robustness / sensitivity (misspecification)
- identifiability (parametrisation, symmetries)

When is MLE hard / weird?

- Gamma, Beta (too  to code up Γ, ψ ?)
- Mixture of Gaussians (non-identifiable, non-convex)
- Latent Variable Models (likelihood = integral)
- Non-Regular Models (uniform, constrained)

Alternatives to MLE

- method of moments (GMM, GEE, QL, ...)
- one-step estimator (\sqrt{n} + Newton)
- robustified, Z-/M-estimation
- (model-specific, more fancy stuff, ...)

Unnormalised Models

- "energy-based", "doubly-intractable"

$$P(x|\theta) = \frac{f(x, \theta)}{z(\theta)}$$

- usually $\dim x \gg 1$ (otherwise, could integrate)
- from DAGs to Factor Graphs, causes to interactions

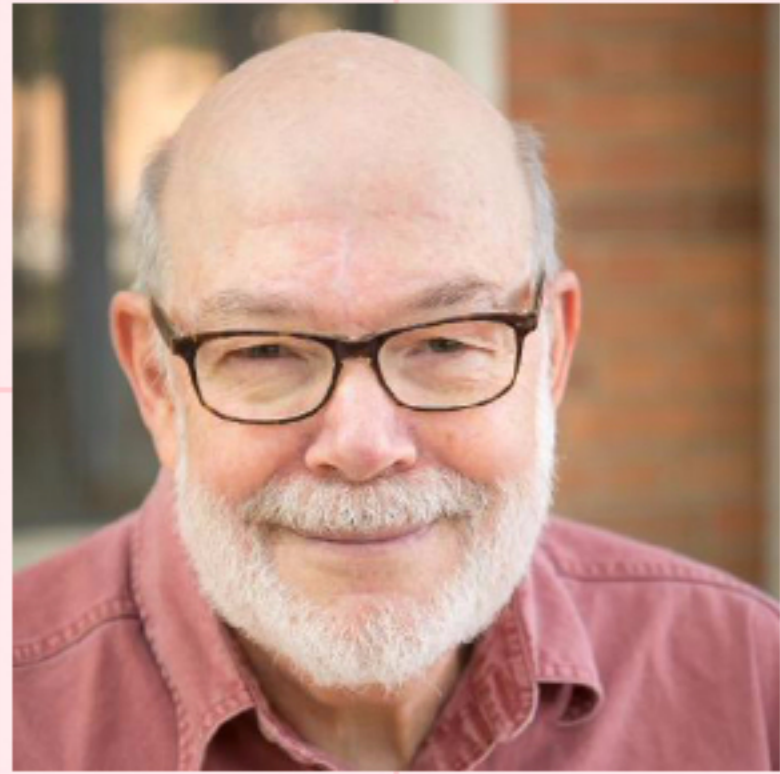
Unnormalised Models

- Ising, (Deep, Restricted) Boltzmann Machine
- (Gaussian, Hidden, Sequential) Markov Random Field
- Text Models, Image Models (Field of Experts)
- ERGMs, Stochastic Block Model, Random Networks
- (Kernel) Exponential Family

Can it get worse?

- latent variables as well! (HMRF, RBM)
- conditionally-unnormalised as well! (LATKES)
- high-dimensional! ($d \approx n$)
- nonparametric! ($\log f = \text{NN, Kernel, ...}$)
- but, already pretty hard ...

MLE Objective, Algorithms



- (MC)MC-MLE / Importance Sampling

$$\log Z(\theta) \approx \log \left(\frac{1}{M} \sum_{j=1}^M \frac{f(y_j, \theta)}{q(y_j)} \right)$$



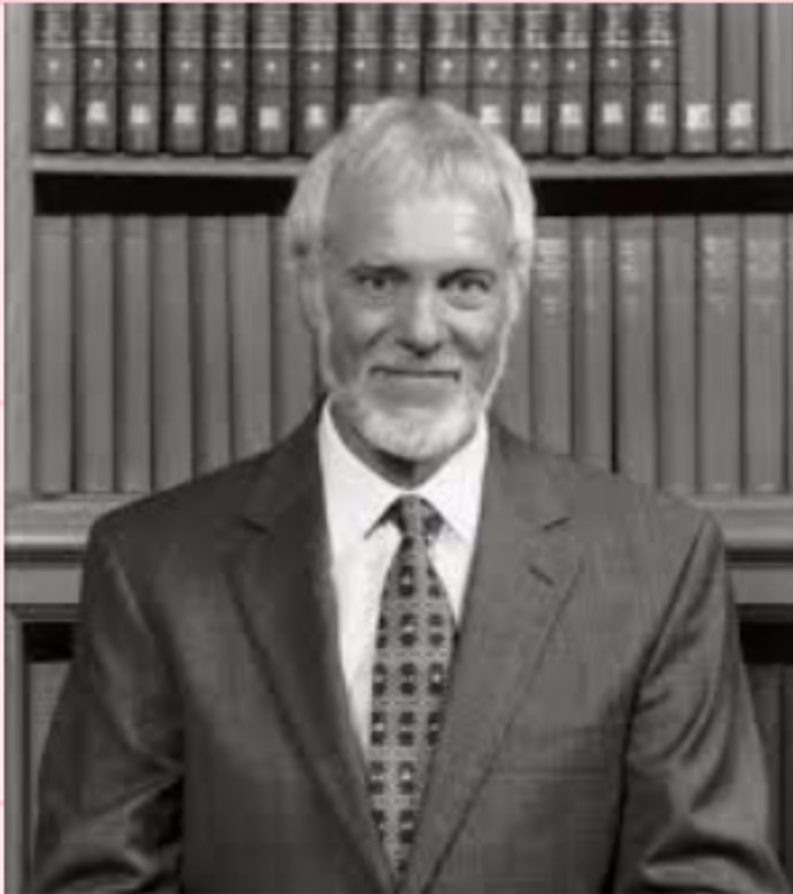
- Stochastic Approximation, Contrastive Divergence

$$\nabla_{\theta} \log Z(\theta) = \mathbb{E}_{\theta} [\nabla_{\theta} \log f(x, \theta)]$$

Non-MLE Objectives

- that damn normalising constant!
- can we make it cancel out?
- tricks: differences, ratios, components
- `_some_` part of the model is tractable?
- safety check: support of P

Graphical Methods



- $\mathbb{P} (X_A | X_B)$ tractable for some A, B
- Pseudo-Likelihood ($A = \{ i \}, B = V \setminus \{ i \}$)
- Composite Likelihood (arbitrary A, B)
- applicable for MRFs (no sparsity needed!)
- often convex

Pseudo-Likelihood Example

- Consider pairwise Markov Random Field

$$P(x|\theta) = \frac{\exp\left(\sum_{i \sim j} \theta_{ij} x_i x_j\right)}{Z(\theta)}$$

$$\Rightarrow P(x_i | x_{-i}, \theta) = \text{Ber}(x_i | \sigma((\Theta x)_{-i}))$$

- In practice: constrain / regularise θ
- Remark: Belief Propagation for Sub-Trees

Difference Methods

- $\nabla_x \log P (x | \theta) = \nabla_x \log f (x, \theta)$
 - \leadsto no $Z (\theta)!$ (c.f. OLD / MCMC)
- requires smoothness of model
- Score Matching, Stein Discrepancies
- different complexities, both often convex

Score Matching Objective



- Score Matching Objective (Q = Data)

$$\mathbb{E}_Q \left[\left| \nabla_x \log Q(x) - \nabla_x \log P(x|\theta) \right|^2 \right]$$

$$\rightsquigarrow \mathbb{E}_Q \left[2 \Delta_x \log P(x|\theta) + \left| \nabla_x \log P(x|\theta) \right|^2 \right] + c$$


$$B \left(\nabla_x \log P \rightarrow \nabla_x \log Q \right)$$

KSD Objective

- Kernel Stein Discrepancy Objective (Q = Data)

$$\begin{aligned} D_p(Q)^2 &= \sup_{h \in \mathcal{H}} \mathbb{E}_Q \left[(\mathcal{L}^p h)(x) \right]^2 \\ &= \mathbb{E}_{Q \otimes Q} \left[(\mathcal{L}_x^p \mathcal{L}_y^p K)(x, y) \right] \end{aligned}$$

Rational Methods

- $P(y | \theta) / P(x | \theta) = f(y, \theta) / f(x, \theta)$
- Ratio Matching
- Let $Q(x)$ be known, classify $Q(x)$ vs $P(x | \theta)$
 - Optimal: $\sigma(\log Q(x) - \log f(x, \theta) - \log Z(\theta))$
 - Noise-Contrastive Estimation 
- Stein Density Ratio Estimation (a bit of both)

Ratio Matching Objective

$$\begin{aligned} & \mathbb{E}_{\mathcal{Q}} \left[\left(\beta \left(\frac{\mathcal{Q}(\phi x)}{\mathcal{Q}(x)} \right) - \beta \left(\frac{P(\phi x | \theta)}{P(x | \theta)} \right) \right)^2 \right] + \text{symmetrise} \\ &= \mathbb{E}_{\mathcal{Q}} \left[\left(1 - \beta \left(\frac{P(\phi x | \theta)}{P(x | \theta)} \right) \right)^2 \right] + \text{const.} \end{aligned}$$

Noise-Contrastive Estimation



$$\begin{array}{llll}
 x_1, \dots, x_N & \stackrel{\text{iid}}{\sim} & P(x|\theta) & z=1 \\
 y_1, \dots, y_M & \stackrel{\text{iid}}{\sim} & Q(y) & z=0 \quad \nu = m/n
 \end{array}$$

$$\begin{aligned}
 P(z=1|u) &= \frac{1}{1 + \nu Q(u)/P(u|\theta)} \\
 &= \sigma\left(\log\left\{\frac{Q(u)}{f(u,\theta)}\right\} + c\right)
 \end{aligned}$$

↙ unknown

- often convex in (θ, c)

Related / Frontiers

- Denoising Autoencoders, Denoising Score Matching
- Score Estimation (SBGM, DDPM, ...)
- Learned Stein Discrepancies (beyond Kernels)
- Hybrids with other approaches

