

Estimating Unnormalised Models by Score Matching

Michael U. Gutmann

Probabilistic Modelling and Reasoning (INFR11134)
School of Informatics, The University of Edinburgh

Spring Semester 2022

Program

1. Basics of score matching
2. Practical objective function for score matching

Program

1. Basics of score matching

- Basic ideas of score matching
- Objective function that captures the basic ideas but cannot be computed

2. Practical objective function for score matching

Problem formulation

- ▶ We want to estimate the parameters θ of a parametric statistical model for a random vector $\mathbf{x} \in \mathbb{R}^d$.
- ▶ Given: data $\mathbf{x}_1, \dots, \mathbf{x}_n$, iid, following p_*
- ▶ Model pdf: $p(\mathbf{x}; \theta)$
- ▶ Assumptions:
 - ▶ Model $p(\mathbf{x}; \theta)$ is known only up to the partition function

$$p(\mathbf{x}; \theta) = \frac{\tilde{p}(\mathbf{x}; \theta)}{Z(\theta)} \quad Z(\theta) = \int_{\mathbf{x}} \tilde{p}(\mathbf{x}; \theta) d\mathbf{x}$$

- ▶ Evaluation of $\tilde{p}(\mathbf{x}; \theta)$ is tractable.
 - ▶ Partition function $Z(\theta)$ cannot be computed analytically in closed form and numerical approximation is expensive.
- ▶ Goal: Estimate the model without approximating the partition function $Z(\theta)$.

Basic ideas of score matching

- ▶ Maximum likelihood estimation can be understood to find parameter values $\hat{\theta}$ so that

$$p(\mathbf{x}; \hat{\theta}) \approx p_*(\mathbf{x}) \quad \text{or} \quad \log p(\mathbf{x}; \hat{\theta}) \approx \log p_*(\mathbf{x})$$

(as measured by Kullback-Leibler divergence, see e.g. Barber 8.7)

- ▶ Instead of estimating the parameters θ by matching (log) densities, score matching identifies parameter values $\hat{\theta}$ for which the derivatives (slopes) of the log densities match

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}; \hat{\theta}) \approx \nabla_{\mathbf{x}} \log p_*(\mathbf{x})$$

- ▶ $\nabla_{\mathbf{x}} \log p(\mathbf{x}; \theta)$ does not depend on the partition function:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}; \theta) = \nabla_{\mathbf{x}} [\log \tilde{p}(\mathbf{x}; \theta) - \log Z(\theta)] = \nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}; \theta)$$

The score function (in the context of score matching)

- ▶ Define the model score function $\mathbb{R}^d \rightarrow \mathbb{R}^d$ as

$$\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial \log p(\mathbf{x}; \boldsymbol{\theta})}{\partial x_1} \\ \vdots \\ \frac{\partial \log p(\mathbf{x}; \boldsymbol{\theta})}{\partial x_d} \end{pmatrix} = \nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta})$$

While defined in terms of $p(\mathbf{x}; \boldsymbol{\theta})$, we also have

$$\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = \nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$$

- ▶ Similarly, define the data score function as

$$\boldsymbol{\psi}_*(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_*(\mathbf{x})$$

Definition of the SM objective function

- ▶ Estimate θ by minimising a distance between model score function $\psi(\mathbf{x}; \theta)$ and score function of observed data $\psi_*(\mathbf{x})$

$$\begin{aligned} J_{\text{sm}}(\theta) &= \frac{1}{2} \int_{\mathbf{x} \in \mathbb{R}^d} p_*(\mathbf{x}) \|\psi(\mathbf{x}; \theta) - \psi_*(\mathbf{x})\|^2 d\mathbf{x} \\ &= \frac{1}{2} \mathbb{E}_* \|\psi(\mathbf{x}; \theta) - \psi_*(\mathbf{x})\|^2 \end{aligned}$$

where \mathbb{E}_* denotes the expectation \mathbb{E}_{p_*} with respect to $p_*(\mathbf{x})$

- ▶ Since $\psi(\mathbf{x}; \theta) = \nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}; \theta)$ does not depend on $Z(\theta)$ there is no need to compute the partition function.
- ▶ Knowing the unnormalised model $\tilde{p}(\mathbf{x}; \theta)$ is enough.
- ▶ Expectation \mathbb{E}_* with respect to p_* can be approximated as sample average over the observed data, but what about ψ_* ?

Program

1. Basics of score matching

- Basic ideas of score matching
- Objective function that captures the basic ideas but cannot be computed

2. Practical objective function for score matching

Program

1. Basics of score matching
2. Practical objective function for score matching
 - Integration by parts to obtain a computable objective function
 - Simple example

Reformulation of the SM objective function

- ▶ In the objective function we have the score function of the data distribution ψ_* . How to compute it?
- ▶ In fact, no need to compute it because the score matching objective function J_{sm} can be expressed as

$$J_{\text{sm}}(\boldsymbol{\theta}) = \mathbb{E}_* \sum_{j=1}^d \left[\partial_j \psi_j(\mathbf{x}; \boldsymbol{\theta}) + \frac{1}{2} \psi_j^2(\mathbf{x}; \boldsymbol{\theta}) \right] + \text{const.}$$

where the constant does not depend on $\boldsymbol{\theta}$, and

$$\psi_j(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})}{\partial x_j} \quad \partial_j \psi_j(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial^2 \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})}{\partial x_j^2}$$

Proof (general idea)

- ▶ Use Euclidean distance and expand the objective function J_{sm}

$$\begin{aligned} J_{\text{sm}}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_* \|\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) - \boldsymbol{\psi}_*(\mathbf{x})\|^2 \\ &= \frac{1}{2} \mathbb{E}_* \|\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})\|^2 - \mathbb{E}_* \left[\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})^\top \boldsymbol{\psi}_*(\mathbf{x}) \right] + \frac{1}{2} \mathbb{E}_* \|\boldsymbol{\psi}_*(\mathbf{x})\|^2 \\ &= \frac{1}{2} \mathbb{E}_* \|\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})\|^2 - \sum_{j=1}^d \mathbb{E}_* [\psi_j(\mathbf{x}; \boldsymbol{\theta}) \psi_{*,j}(\mathbf{x})] + \text{const} \end{aligned}$$

- ▶ First term does not depend on $\boldsymbol{\psi}_*$. The ψ_j and $\psi_{*,j}$ are the j -th elements of the vectors $\boldsymbol{\psi}$ and $\boldsymbol{\psi}_*$, respectively. Constant does not depend on $\boldsymbol{\theta}$.
- ▶ The trick is to use integration by parts for the second term to get an objective function which does not involve $\boldsymbol{\psi}_*$.

Proof (not examinable)

$$\begin{aligned}\mathbb{E}_* [\psi_j(\mathbf{x}; \boldsymbol{\theta}) \psi_{*,j}(\mathbf{x})] &= \int_{\mathbf{x}} p_*(\mathbf{x}) \psi_{*,j}(\mathbf{x}) \psi_j(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} \\ &= \int_{\mathbf{x}} p_*(\mathbf{x}) \frac{\partial \log p_*(\mathbf{x})}{\partial x_j} \psi_j(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} \\ &= \prod_{k \neq j} \int_{x_k} \left(\int_{x_j} p_*(\mathbf{x}) \frac{\partial \log p_*(\mathbf{x})}{\partial x_j} \psi_j(\mathbf{x}; \boldsymbol{\theta}) dx_j \right) dx_k \\ &= \prod_{k \neq j} \int_{x_k} \left(\int_{x_j} \frac{\partial p_*(\mathbf{x})}{\partial x_j} \psi_j(\mathbf{x}; \boldsymbol{\theta}) dx_j \right) dx_k\end{aligned}$$

Use integration by parts

$$\begin{aligned}\int_{x_j} \frac{\partial p_*(\mathbf{x})}{\partial x_j} \psi_j(\mathbf{x}; \boldsymbol{\theta}) dx_j &= [p_*(\mathbf{x}) \psi_j(\mathbf{x}; \boldsymbol{\theta})]_{a_j}^{b_j} - \int_{x_j} p_*(\mathbf{x}) \frac{\partial \psi_j(\mathbf{x}; \boldsymbol{\theta})}{\partial x_j} dx_j \\ &= - \int_{x_j} p_*(\mathbf{x}) \frac{\partial \psi_j(\mathbf{x}; \boldsymbol{\theta})}{\partial x_j} dx_j,\end{aligned}$$

where the a_j and b_j specify the boundaries of the data pdf p_* along dimension j and where we assume that $[p_*(\mathbf{x}) \psi_j(\mathbf{x}; \boldsymbol{\theta})]_{a_j}^{b_j} = 0$.

Proof (not examinable)

If $[p_*(\mathbf{x})\psi_j(\mathbf{x}; \boldsymbol{\theta})]_{a_j}^{b_j} = 0$:

$$\begin{aligned}\mathbb{E}_* [\psi_j(\mathbf{x}; \boldsymbol{\theta})\psi_{*,j}(\mathbf{x})] &= - \prod_{k \neq j} \int_{x_k} \left(\int_{x_j} p_*(\mathbf{x}) \frac{\partial \psi_j(\mathbf{x}; \boldsymbol{\theta})}{\partial x_j} dx_j \right) dx_k \\ &= - \int_{\mathbf{x}} p_*(\mathbf{x}) \frac{\partial \psi_j(\mathbf{x}; \boldsymbol{\theta})}{\partial x_j} d\mathbf{x} \\ &= -\mathbb{E}_* [\partial_j \psi_j(\mathbf{x}; \boldsymbol{\theta})]\end{aligned}$$

so that

$$\begin{aligned}J_{\text{sm}}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_* \|\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})\|^2 - \sum_{j=1}^d -\mathbb{E}_* [\partial_j \psi_j(\mathbf{x}; \boldsymbol{\theta})] + \text{const} \\ &= \mathbb{E}_* \sum_{j=1}^d \left[\partial_j \psi_j(\mathbf{x}; \boldsymbol{\theta}) + \frac{1}{2} \psi_j^2(\mathbf{x}; \boldsymbol{\theta}) \right] + \text{const}\end{aligned}$$

Replacing the expectation / integration over the data density p_* by a sample average over the observed data gives a computable objective function for score matching.

Final method of score matching

- ▶ Given iid data $\mathbf{x}_1, \dots, \mathbf{x}_n$, the score matching estimate is

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$
$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \left[\partial_j \psi_j(\mathbf{x}_i; \boldsymbol{\theta}) + \frac{1}{2} \psi_j(\mathbf{x}_i; \boldsymbol{\theta})^2 \right]$$

ψ_j is the partial derivative of the log unnormalised model $\log \tilde{p}$ with respect to the j -th coordinate (slope) and $\partial_j \psi_j$ its second partial derivative (curvature).

- ▶ Parameter estimation with intractable partition functions without approximating the partition function.

Requirements

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \left[\partial_j \psi_j(\mathbf{x}_i; \boldsymbol{\theta}) + \frac{1}{2} \psi_j(\mathbf{x}_i; \boldsymbol{\theta})^2 \right]$$

Requirements:

- ▶ technical (from the proof): $[p_*(\mathbf{x})\psi_j(\mathbf{x}; \boldsymbol{\theta})]_{a_j}^{b_j} = 0$, where the a_j and b_j specify the boundaries of the data pdf p_* along dimension j
- ▶ smoothness: second derivatives of $\log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$ with respect to the x_j need to exist, and should be smooth with respect to $\boldsymbol{\theta}$ so that $J(\boldsymbol{\theta})$ can be optimised with gradient-based methods.

Simple example

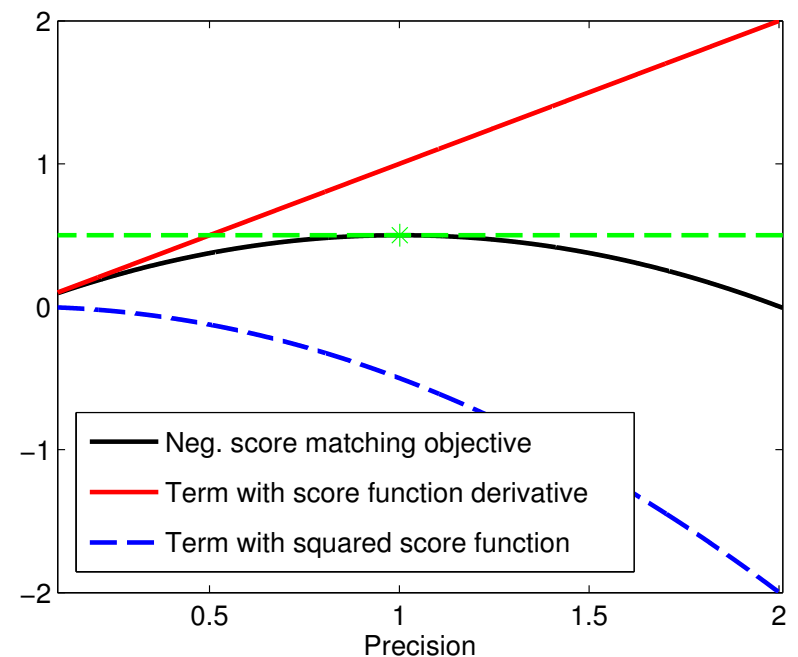
- ▶ $\tilde{p}(x; \theta) = \exp(-\theta x^2/2)$, parameter $\theta > 0$ is the precision.
- ▶ The slope and curvature of the log unnormalised model are

$$\psi(x; \theta) = \partial_x \log \tilde{p}(x; \theta) = -\theta x, \quad \partial_x \psi(x; \theta) = -\theta.$$

- ▶ If p_* is Gaussian, $\lim_{x \rightarrow \pm\infty} p_*(x)\psi(x; \theta) = 0$ for all θ .
- ▶ Score matching objective

$$J(\theta) = -\theta + \frac{1}{2}\theta^2 \frac{1}{n} \sum_{i=1}^n x_i^2$$
$$\Rightarrow \hat{\theta} = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right)^{-1}$$

- ▶ For Gaussians, same as the MLE.



Extensions

- ▶ Score matching as presented here only works for $\mathbf{x} \in \mathbb{R}^d$
- ▶ There are extensions for discrete and non-negative random variables (not examinable)
<https://www.cs.helsinki.fi/u/ahyvarin/papers/CSDA07.pdf>
- ▶ Can be shown to be part of a general framework to estimate unnormalised models (not examinable)
<https://michaelgutmann.github.io/assets/papers/Gutmann2011b.pdf>
- ▶ Overall message: in some situations, other learning criteria than maximum likelihood are preferable.

Program recap

1. Basics of score matching

- Basic ideas of score matching
- Objective function that captures the basic ideas but cannot be computed

2. Practical objective function for score matching

- Integration by parts to obtain a computable objective function
- Simple example