

# Estimating Unnormalised Models by Score Matching

Michael Gutmann

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

Spring Semester 2020

# 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  $\theta$  of a parametric statistical model for a random vector  $\mathbf{x} \in \mathbb{R}^d$ .
- ▶ Given: iid data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  that are assumed to be observations of  $\mathbf{x}$  that has pdf  $p_*$
- ▶ Further notation:  $p(\xi; \theta)$  is the model pdf;  $\xi \in \mathbb{R}^d$  is a dummy variable (not a random variable).

- ▶ Assumptions:

- ▶ Model  $p(\xi; \theta)$  is known only up the partition function

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

- ▶ Evaluation of  $\tilde{p}(\xi; \theta)$  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

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

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

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

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

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

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

# The score function (in the context of score matching)

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

$$\psi(\xi; \theta) = \begin{pmatrix} \frac{\partial \log p(\xi; \theta)}{\partial \xi_1} \\ \vdots \\ \frac{\partial \log p(\xi; \theta)}{\partial \xi_d} \end{pmatrix} = \nabla_{\xi} \log p(\xi; \theta)$$

While defined in terms of  $p(\xi; \theta)$ , we also have

$$\psi(\xi; \theta) = \nabla_{\xi} \log \tilde{p}(\xi; \theta)$$

- ▶ Similarly, define the data score function as

$$\psi_*(\xi) = \nabla_{\xi} \log p_*(\xi)$$

# Definition of the SM objective function

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

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

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

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

# 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  $\psi_*$ . How to compute it?
- ▶ In fact, no need to compute it because the score matching objective function  $J_{\text{sm}}$  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

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

# Proof (general idea)

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

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

- ▶ First term does not depend on  $\boldsymbol{\psi}_*$ . The  $\psi_j$  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)

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

Use integration by parts

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

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

# Proof (not examinable)

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

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

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

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

$\psi_j$  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

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

Requirements:

- ▶ technical (from the proof):  $[p_*(\boldsymbol{\xi})\psi_j(\boldsymbol{\xi}; \boldsymbol{\theta})]_{a_j}^{b_j} = 0$ , where the  $a_j$  and  $b_j$  specify the boundaries of the data pdf  $p_*$  along dimension  $j$
- ▶ smoothness: second derivatives of  $\log \tilde{p}(\boldsymbol{\xi}; \boldsymbol{\theta})$  with respect to the  $\xi_j$  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

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

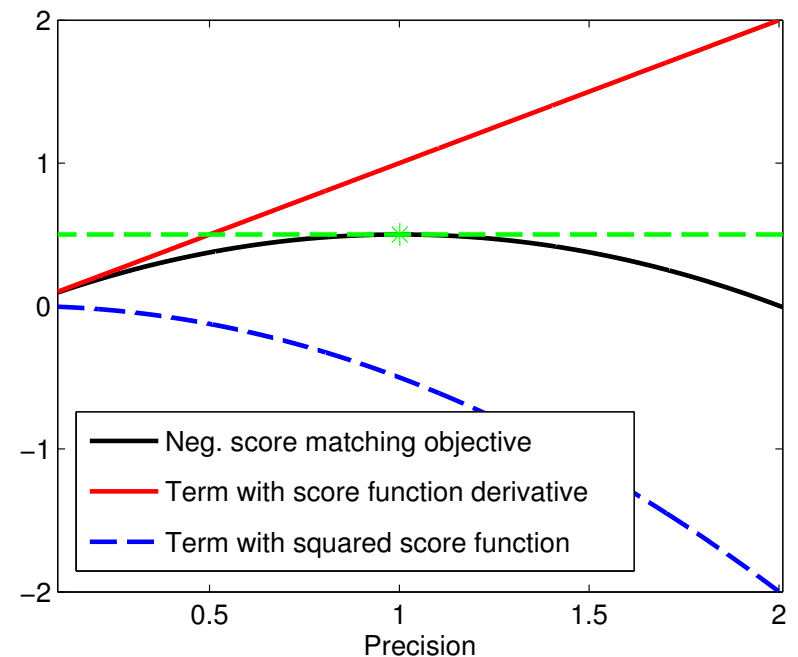
$$\psi(\xi; \theta) = \partial_{\xi} \log \tilde{p}(\xi; \theta) = -\theta\xi, \quad \partial_{\xi}^2 \psi(\xi; \theta) = -\theta.$$

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