

# Taming numerical imprecision by adapting the KL divergence to negative probabilities<sup>a</sup>

---

Simon Pfahler<sup>1</sup>, Peter Georg<sup>1</sup>, Rudolf Schill<sup>2</sup>, Maren Klever<sup>3</sup>, Lars Grasedyck<sup>3</sup>,  
Rainer Spang<sup>1</sup>, Tilo Wettig<sup>1</sup>

<sup>1</sup>University of Regensburg, <sup>2</sup>ETH Zürich, <sup>3</sup>RWTH Aachen University

August 28, 2024

Slides: [pfahler.online](http://pfahler.online)

---

<sup>a</sup>Pfahler et al. in *Stat. Comput.* **34** (2024)

# Introduction

---

## Introduction - Example

- Consider the SVD of a matrix:

$$p \sim \begin{pmatrix} 1 & 2 & 2 \\ 2 & 2 & 1 \\ 2 & 1 & 0.01 \end{pmatrix} \approx \begin{pmatrix} -0.62 & 0.68 & 0.40 \\ -0.66 & -0.18 & -0.73 \\ -0.42 & -0.72 & 0.56 \end{pmatrix} \begin{pmatrix} 4.51 & 0 & 0 \\ 0 & 1.63 & 0 \\ 0 & 0 & 0.14 \end{pmatrix} \begin{pmatrix} -0.62 & -0.66 & -0.42 \\ -0.68 & 0.18 & 0.72 \\ 0.40 & -0.73 & 0.56 \end{pmatrix}$$

- Truncating the smallest singular value, we get the approximation

$$\tilde{p} \sim \begin{pmatrix} 0.98 & 2.04 & 1.97 \\ 2.04 & 1.93 & 1.06 \\ 1.97 & 1.06 & -0.03 \end{pmatrix} \approx \begin{pmatrix} -0.62 & 0.68 \\ -0.66 & -0.18 \\ -0.42 & -0.72 \end{pmatrix} \begin{pmatrix} 4.51 & 0 \\ 0 & 1.63 \end{pmatrix} \begin{pmatrix} -0.62 & -0.66 & -0.42 \\ -0.68 & 0.18 & 0.72 \end{pmatrix}$$

- $\tilde{p}$  can not be used to calculate a KL divergence, even though it is an approximation of a probability distribution

## Introduction - Approach

$$D_{\text{KL}}(p\|q) = \sum_{i=1}^n p_i \log \frac{p_i}{q_i}$$

- we consider the case where  $p$  is our data and  $q$  is our model
- negative entries can occur in  $q$  due to various reasons:
  1. approximations
  2. assumptions of a theory
  3. rounding errors
- if a KL-like divergence is needed, we can
  - (a) give up
  - (b) reformulate
  - (c) devise a workaround
- this talk follows approach (c)

## **shifted KL divergence**

---

## sKL divergence - Idea

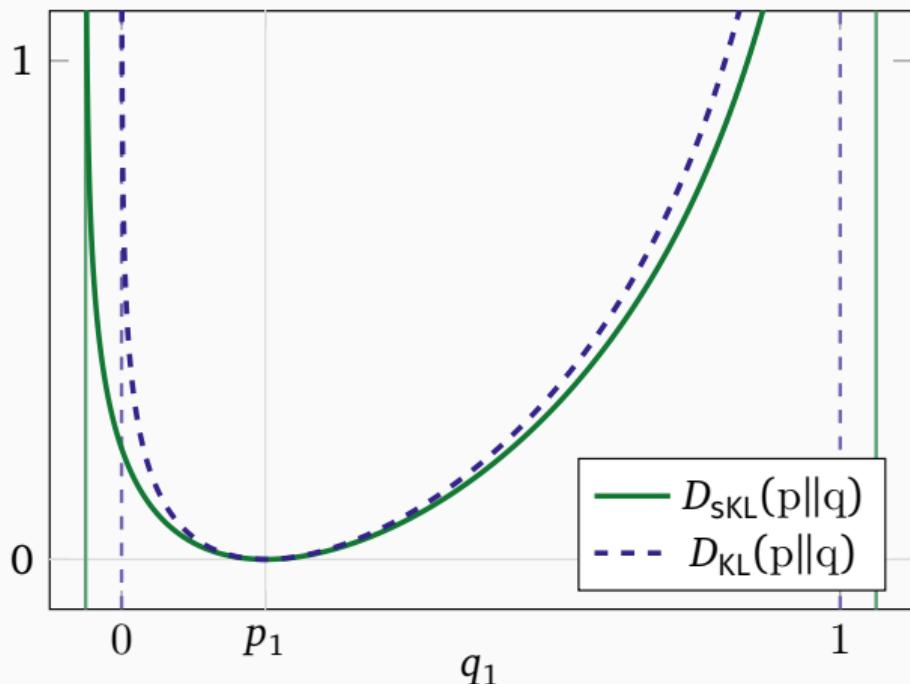
$$D_{\text{sKL}}(p\|q) = \sum_{i=1}^n (p_i + \varepsilon_i) \log \frac{p_i + \varepsilon_i}{q_i + \varepsilon_i}$$

- Negative entries cause a problem because  $\log x$  is only well-defined for  $x > 0$
- How can we adapt the KL divergence so we can use it as an objective function in an optimization task?
- Can we just shift the negative entries?  
→ Almost!
- In order to retain important properties of the KL divergence, p also has to be shifted

## sKL divergence - Properties

The sKL divergence retains many properties of the KL divergence:

- The sKL divergence is a statistical divergence, i.e.
  - (1)  $D_{\text{sKL}}(p\|q) \geq 0$
  - (2)  $D_{\text{sKL}}(p\|q) = 0 \Leftrightarrow p = q$
  - (3)  $\frac{d^2 D_{\text{sKL}}(p\|q)}{dq_i dq_j} \Big|_{q=p}$  is a positive definite matrix
- The sKL divergence is convex in the pair of its arguments



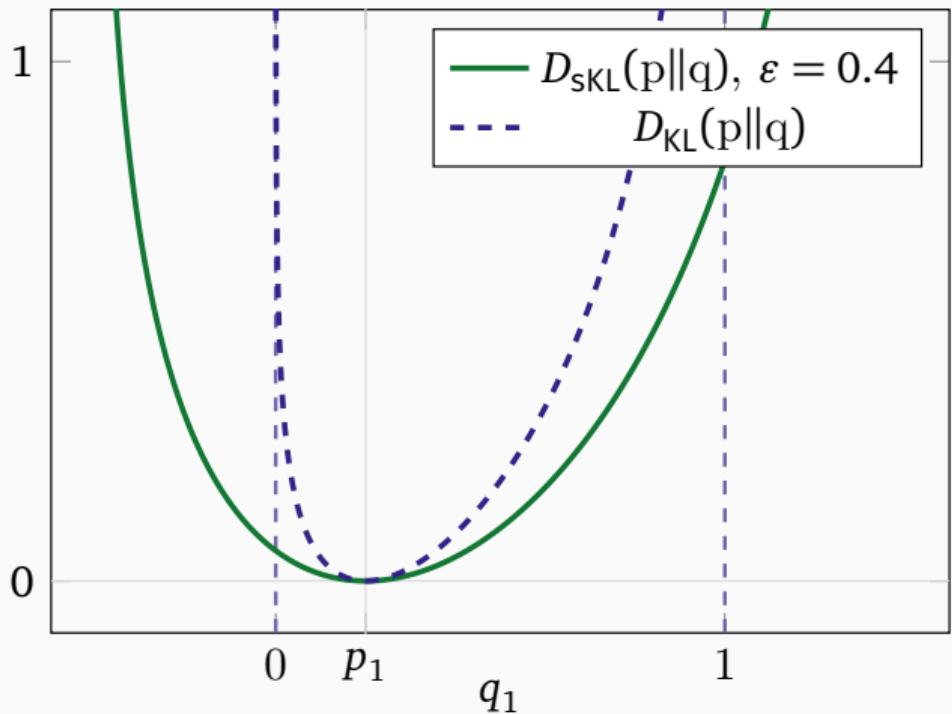
2D example with  $p = (0.2, 0.8)$ ,  $q = (q_1, 1 - q_1)$  and  $\varepsilon = (0.05, 0.05)$

## Parameter choice

---

## Parameter choice - General observations

$$D_{\text{SKL}}(p\|q) = \sum_{i=1}^n (p_i + \varepsilon_i) \log \frac{p_i + \varepsilon_i}{q_i + \varepsilon_i}$$



- Especially in high-dimensional data, often many entries of the data distribution  $p$  are zero  
→ choose  $\varepsilon_i = 0$  if  $p_i = 0$
- The gradient with respect to  $q_i$  is

$$\frac{\partial D_{\text{SKL}}(p\|q)}{\partial q_i} = -\frac{p_i + \varepsilon_i}{q_i + \varepsilon_i}$$
$$\xrightarrow{\varepsilon_i \rightarrow \infty} -1$$

→ large  $\varepsilon_i$  lead to a loss of gradient information

## Parameter choice - Static choice

$$D_{\text{SKL}}(p\|q) = \sum_{i=1}^n (p_i + \varepsilon_i) \log \frac{p_i + \varepsilon_i}{q_i + \varepsilon_i}$$

Simplest idea:  $\varepsilon_i = \begin{cases} 0, & p_i = 0, \\ \varepsilon, & \text{else} \end{cases}$

### ⊕ Pros:

- $\varepsilon$  does not depend on  $q$ , so the SKL divergence is the same function for any  $q$
- Use in higher-order optimizers is possible

### ⊖ Cons:

- A reasonable value for  $\varepsilon$  has to be determined before the SKL divergence can be used
- Negative values  $q_i < -\varepsilon$  still render the SKL divergence undefined

## Parameter choice - Dynamic choice

$$D_{\text{sKL}}(p\|q) = \sum_{i=1}^n (p_i + \varepsilon_i) \log \frac{p_i + \varepsilon_i}{q_i + \varepsilon_i}$$

Dynamic choice:  $\varepsilon_i = \begin{cases} 0, & p_i = 0 \text{ or } q_i > 0, \\ |q_i| + f(|q_i|), & \text{else} \end{cases}$

with  $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ , e.g.  $f(x) = \delta \cdot x$  with  $\delta > 0$

### ⊕ Pros:

- No prior knowledge about  $q$  is needed
- Less gradient information is lost
- Average of the sKL divergence under Gaussian noise is controlled  
→ next slide

### ⊖ Cons:

- The sKL divergence is a different function for every  $q$
- Only first-order optimizers can be used

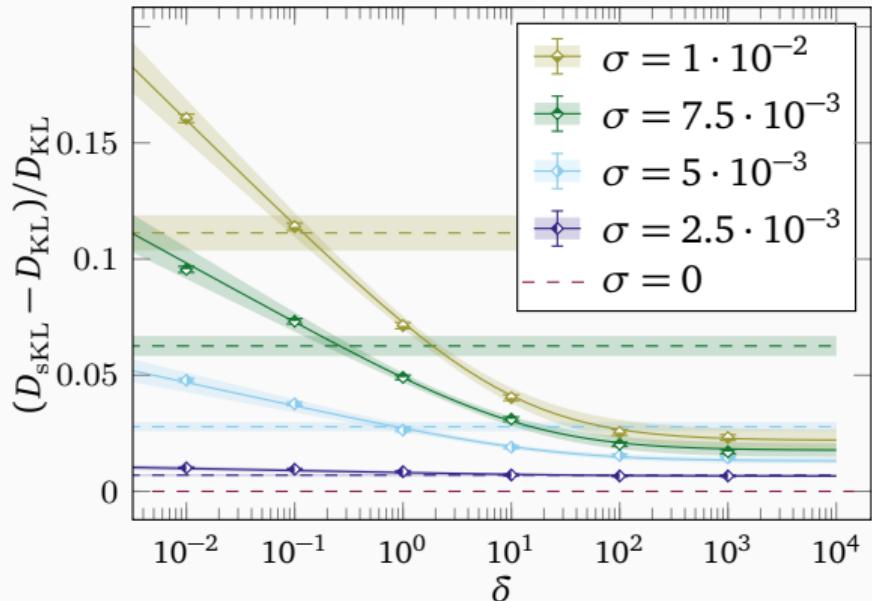
## Parameter choice - Gaussian noise

$$D_{\text{SKL}}(p\|q) = \sum_{i=1}^n (p_i + \varepsilon_i) \log \frac{p_i + \varepsilon_i}{q_i + \varepsilon_i}$$

Consider a vector  $x$  of i.i.d. Gaussian random variables with mean 0 and standard deviation  $\sigma$ .

With the dynamic choice of  $\varepsilon$ , we get

$$\begin{aligned} & \langle D_{\text{SKL}}(p\|q+x) \rangle_x \\ &= D_{\text{KL}}(p\|q) + \sigma^2 \sum_i \frac{p_i}{2q_i^2} + \mathcal{O}(\sigma^4) \end{aligned}$$



example for  $p, q \in \mathbb{R}^{10}$

numerical result

$\sigma^2$ -term

## Application

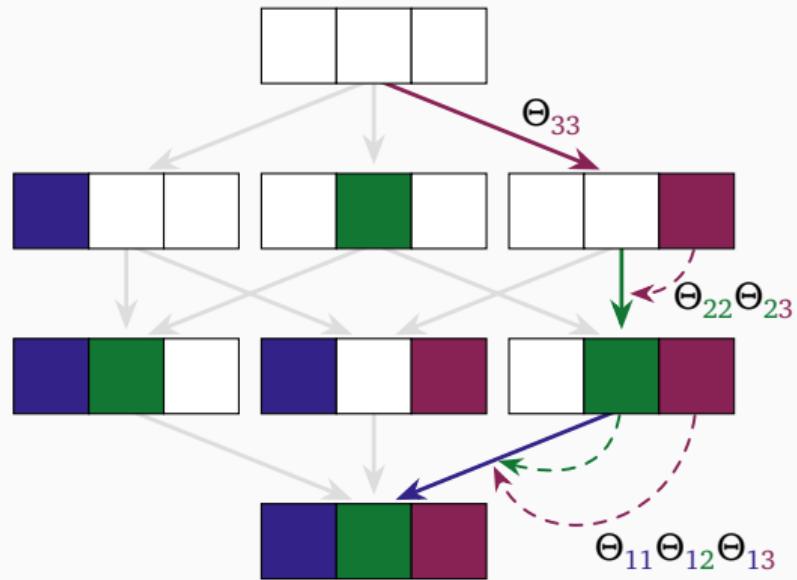
---

# Application - Mutual Hazard Networks<sup>b</sup>

- Modeling tumor progression as a Markov chain
- $d$  binary events (mutations, CNAs)  
→ 
- Tumors start healthy and accumulate events one by one with transition rate

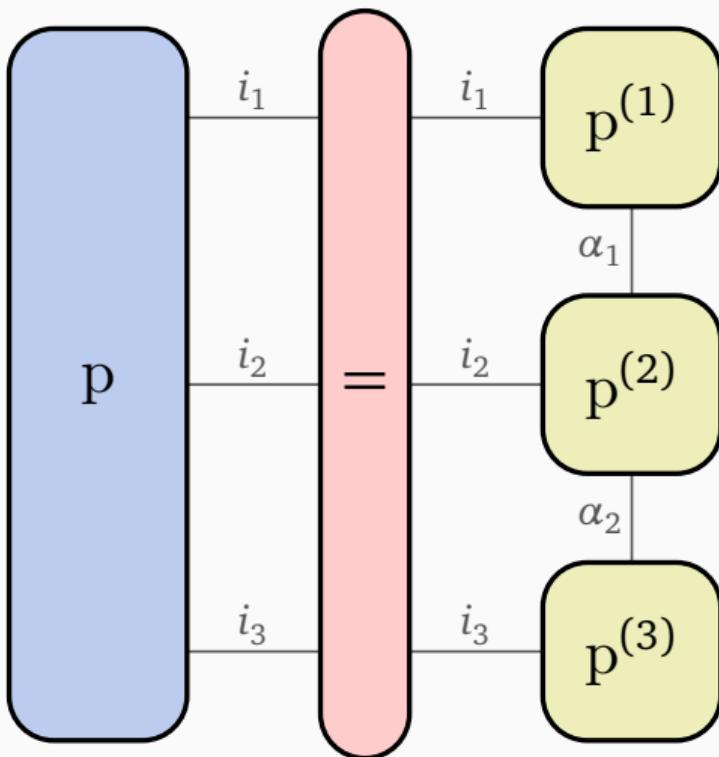
$$Q_{x^{+i}, x} = \Theta_{ii} \prod_{x_j=1} \Theta_{ij}$$

- The model is parameterized by  $\Theta \in \mathbb{R}^{d \times d}$ , which we can learn from patient data



<sup>b</sup>Schill et al. in *Bioinformatics* **36** (2020)

## Application - Mutual Hazard Networks



- The model parameters  $\Theta$  are optimized by minimizing the KL divergence of

$$p_\Theta = (I - Q_\Theta)^{-1} p(t=0)$$

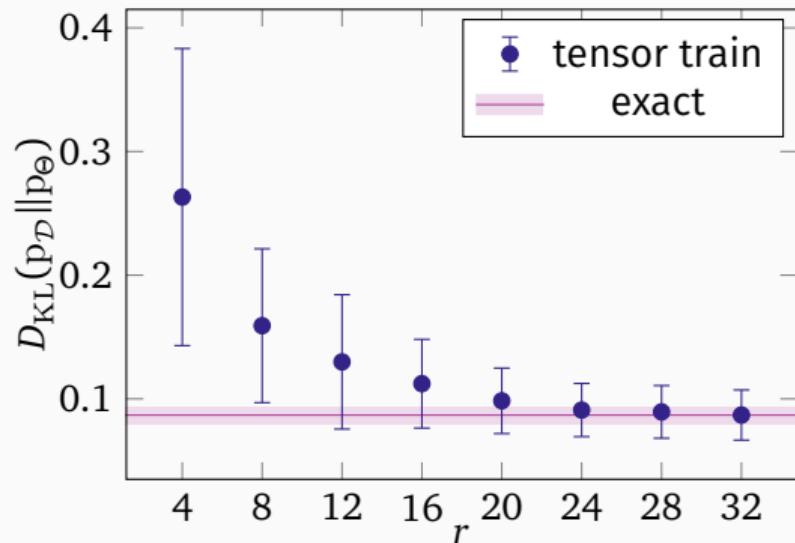
from a given patient data distribution  $p_D$

- The size of  $Q_\Theta \in \mathbb{R}^{2^d \times 2^d}$  and  $p_\Theta \in \mathbb{R}^{2^d}$  grows exponentially with  $d$   
→ Efficient format needed!
- We use the tensor train format<sup>c</sup>, which approximates the probability distribution using  $d$  3D tensors

<sup>c</sup>Oseledets in SIAM J. Sci. Comput. 33 (2011)

## Application - Results

- Tests on datasets with 20 events, so we can compare to the exact result
  - Tensor train rank  $r$  is used to control the approximation quality (tensor trains would be exact for  $r = 1024$ )
- ⇒ Results quickly converge when increasing the approximation quality



## **Summary**

---

## Summary

- The sKL divergence provides a statistical divergence measure which
  1. shares important theoretical properties with the KL divergence
  2. is compatible with approximations of probability vectors, even if negative values occur
  3. is applicable to a wide range of problems
- Optimization tasks on high-dimensional data can be performed using the sKL divergence as an objective function
- Mutual Hazard Networks with up to  $\sim 40$  events (compared to  $\sim 25$  before) can be learned using tensor trains  
→ for more events, we are working on another approach

# Team

University of Regensburg



**Tilo Wettig**

Physics



**Peter Georg**



**Marco Huber**

RWTH Aachen University



**Lars Grasedyck**



**Maren Klever**

ETH Zürich



**Rainer Spang**



**Linda Hu**



**Andreas Lösch**



**Rudolf Schill**



**Kevin Rupp**

Bioinformatics