

Motivation

- Top-k and k-hot vectors are a fundamental modeling tool.
- Top-k is non-differentiable, hindering gradient-based optimization.
- Existing solutions are not as simple and scalable as those for sampling one-hot vectors.

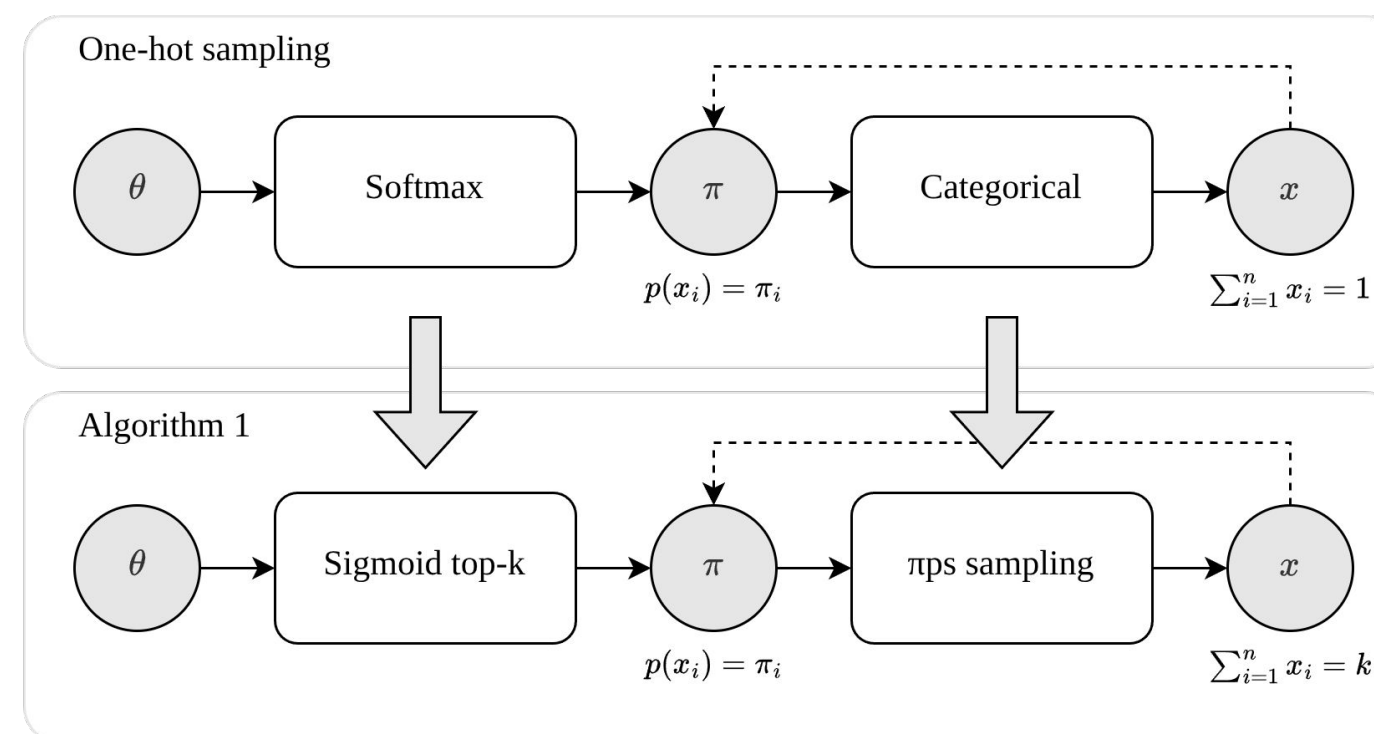
Overview

Algorithm 1 Sample k-hot

Require: $\theta \in \mathbb{R}^n$, $k \in \mathbb{N}$, $1 < k < n$
 1: $\pi \leftarrow \sigma_k(\theta)$
 2: Sample $x \in \{0, 1\}_k^n$ such that $p(x_i) = \pi_i$
 3: **return** x

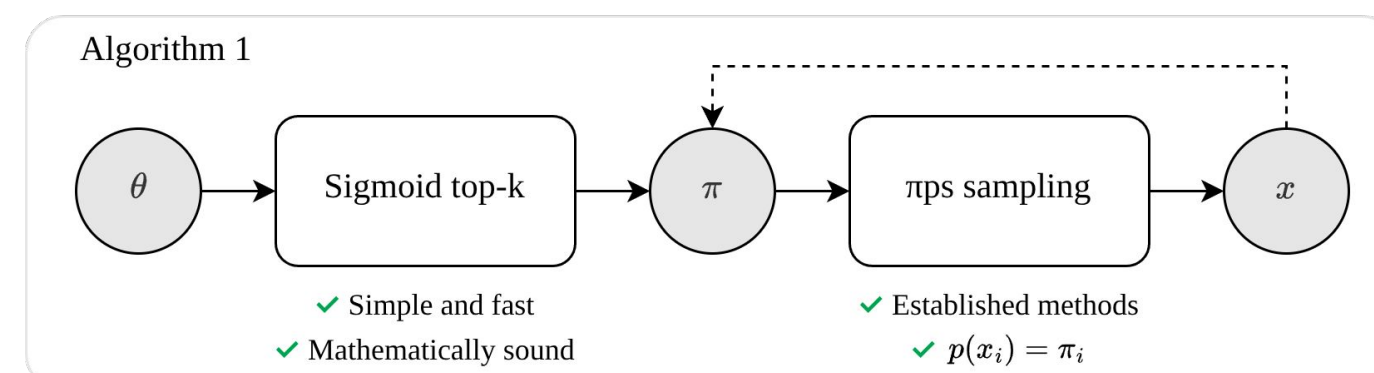
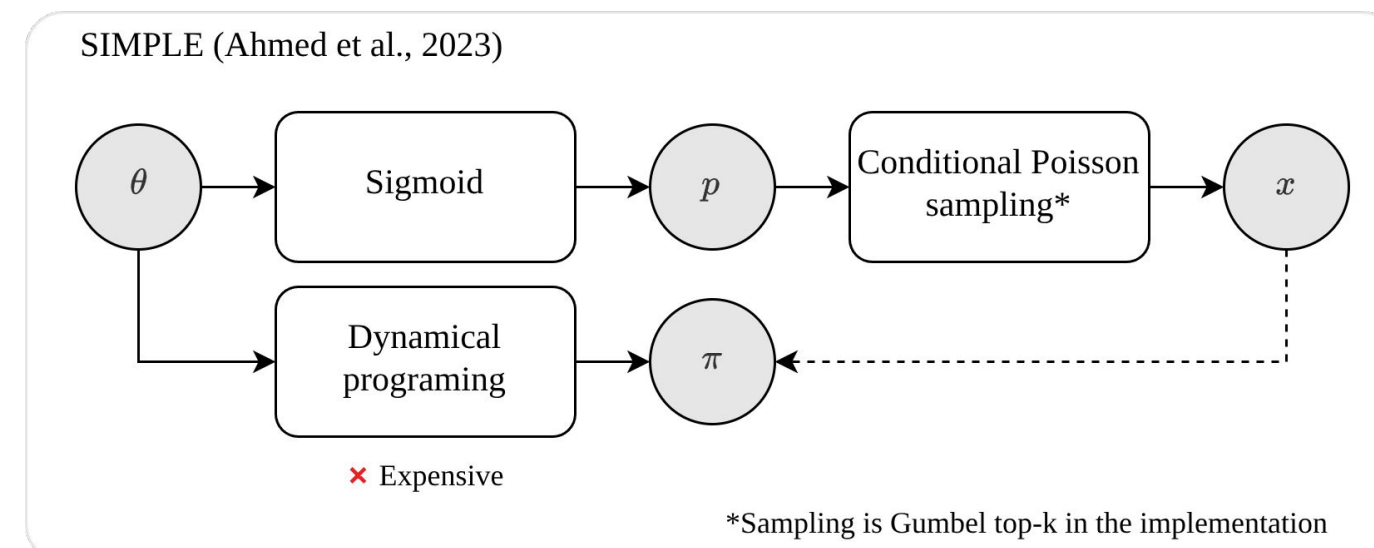
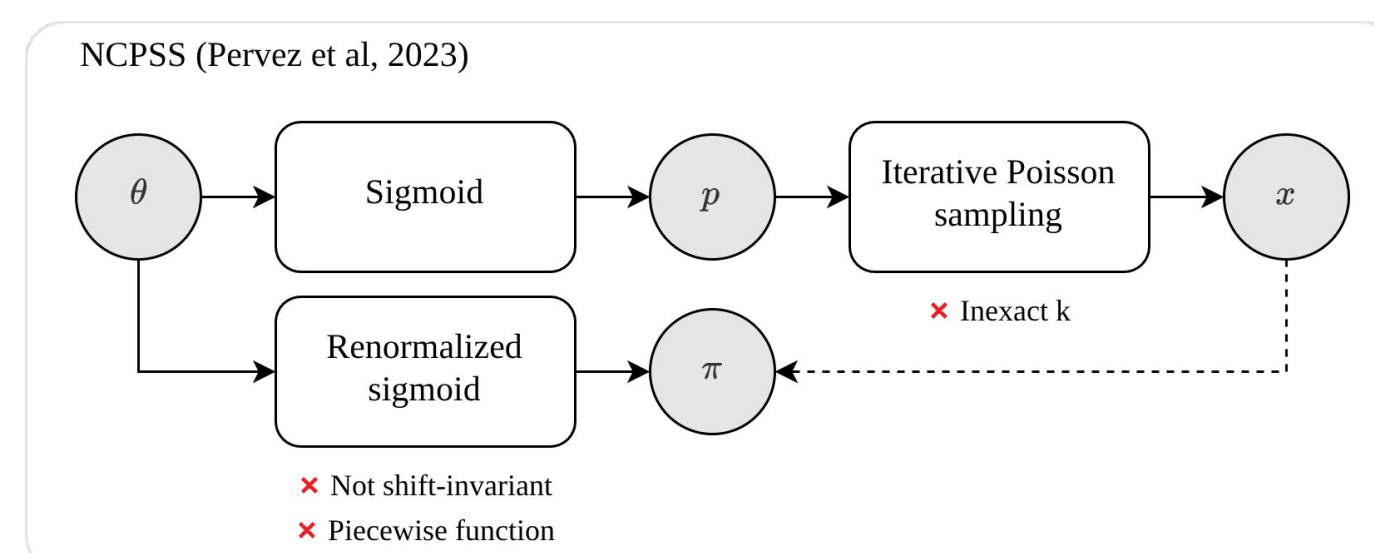
Algorithm 2 Sample relaxed k-hot

Require: $\theta \in \mathbb{R}^n$, $k \in \mathbb{N}$, $1 < k < n$
 1: Sample $g_i \sim \text{Gumbel}(0, 1)$ for $i = 1, \dots, n$
 2: $x \leftarrow \sigma_k(\theta + g)$
 3: **return** x



A simplified and improved approach

- Fast, correct, and uses a simple straight-through gradient estimate.



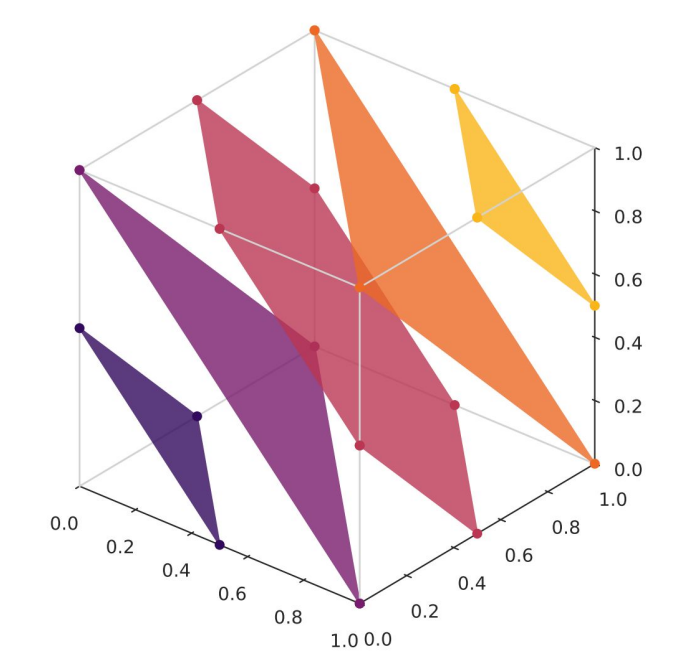
Relaxed top-k

The k-capped simplex

$$\Delta_k^{n-1} := \{\pi \in [0, 1]^n \mid \sum_{i=1}^n \pi_i = k\}$$

- Probabilities that sum to k.
- Codomain of sigmoid top-k.
- Parameter space of π ps sampling.

Generalized simplex



Sigmoid top-k

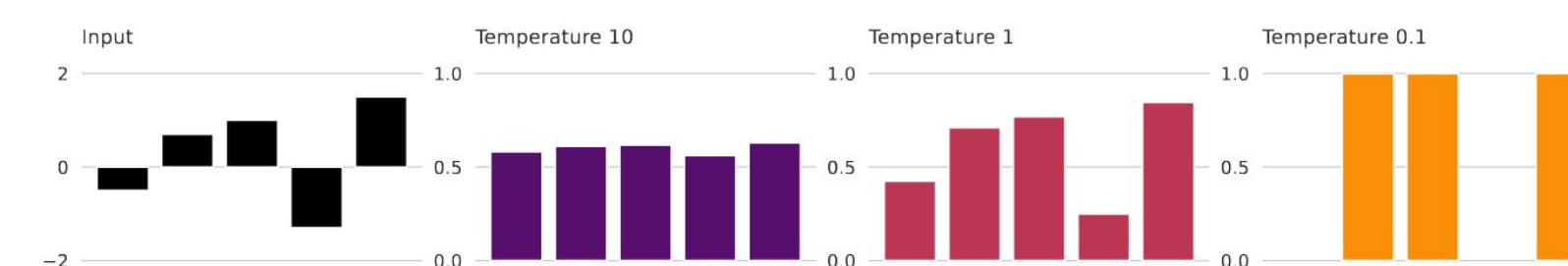
Generalized softmax

$$\sigma_k(x) := \sigma(x + c\mathbf{1}), \text{ where } c \in \mathbb{R} \text{ solves } \sum_{i=1}^n \sigma(x_i + c) = k$$

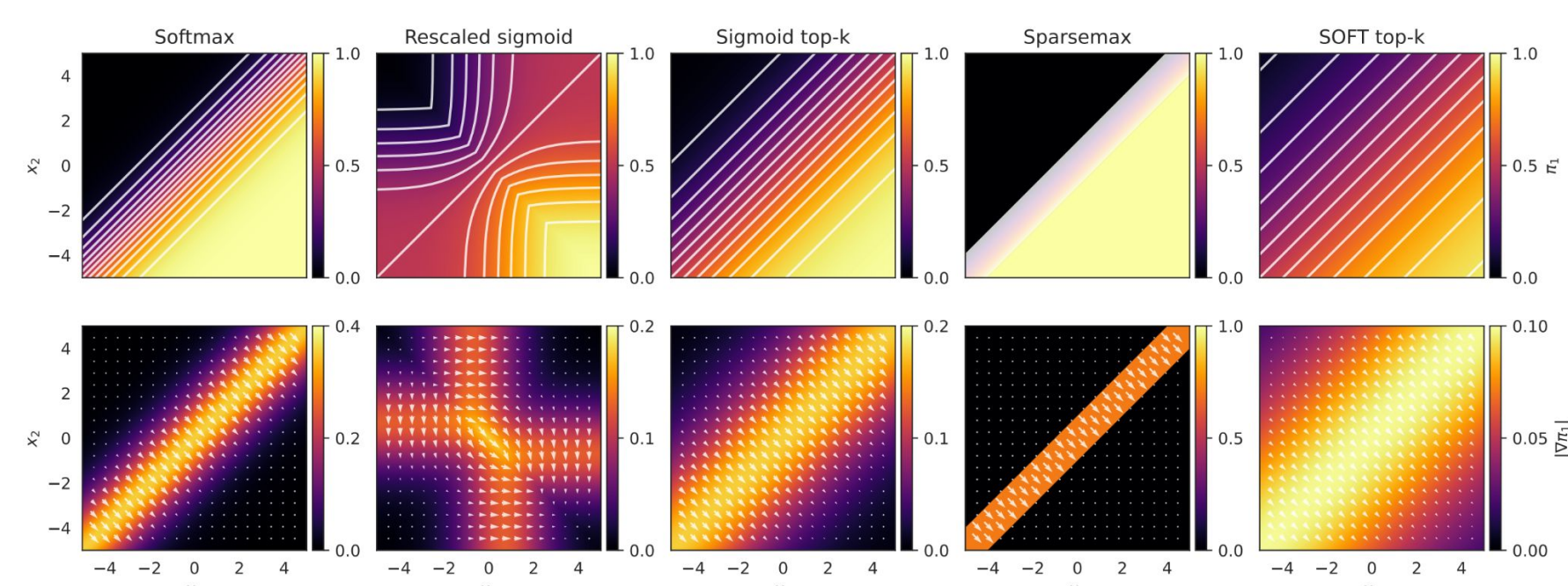
- Top-k relaxation.
- Shares many properties with softmax.
- Forward: **scalar root-finding**.
- Backward: **implicit differentiation**.

$$\sigma_k(x) = \arg \max_{\pi \in \Delta_k^{n-1}} x^\top \pi + \sum_{i=1}^n H(\pi_i)$$

- Sigmoid top-k solves an entropy-regularized optimization problem.
- It is fully differentiable, with respect to both x and k .
- It can be tempered like softmax.
- It is inverted by the logit function up to an additive constant.



Sigmoid top-k can be tempered like softmax.



Value and gradients of top-k relaxations for $n = 2$ and $k = 1$.

Sampling

π ps sampling

Generalized categorical

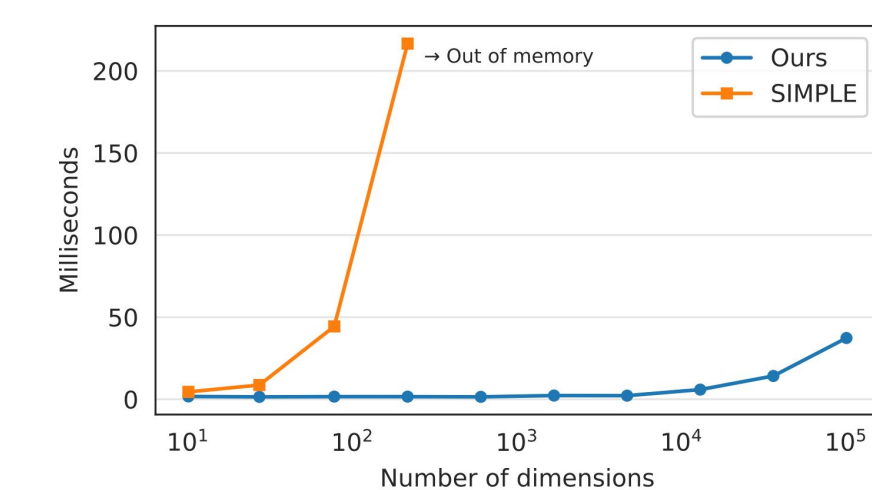
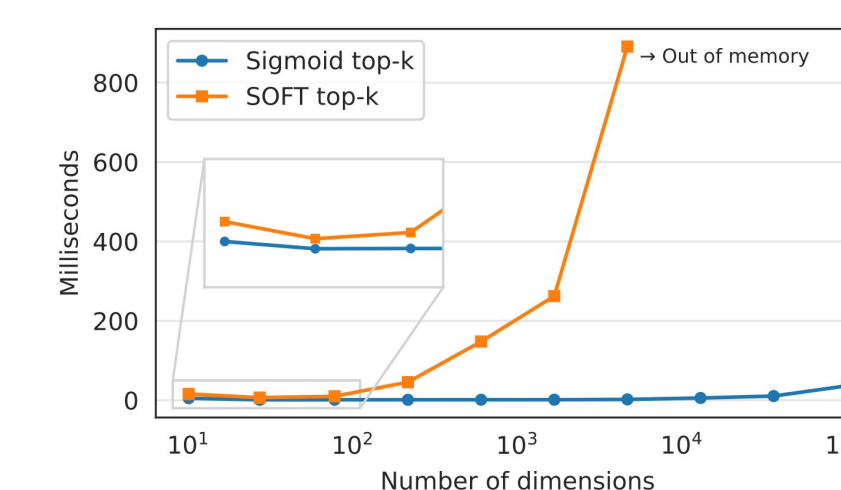
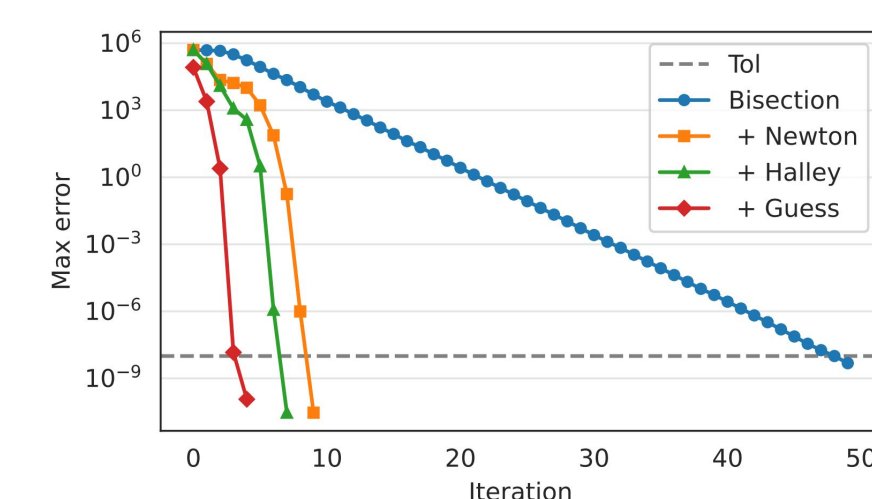
- Methods from the field of sampling design.
- Sampling proportional to size without replacement.
- Parameterized by desired **inclusion probabilities**.
- Many possible distributions (designs) and algorithms (procedures).

$$p(x_i) = \pi_i$$

- Exact inclusion probabilities.
- We prove that straight-through estimation is a first-order approximation of the true gradient in expectation.

Scalability

- Scalar root-finding is scalable.
- Guaranteed convergence.
- Scalable to high dimensions.
- A benchmark in 10^6 dim is shown to the right.
- Comparisons against two methods from prior work below.



Applications

Top-k appears in many models. Our methods can be used to make it **differentiable**, **stochastic**, or **both**.

- Routing
 - Mixture of experts (MoE)
- Sparse coding
 - Dictionary learning (k-SVD)
 - Sparse system identification (SINDy)
 - Sparse autoencoders (SAE)
- Search
 - k-nearest neighbors
 - Beam search
- Sparse networks
- Feature selection