

Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity

Gonçalo Correia Instituto de Telecomunicações, Lisbon

Vlad Niculae Ivl, University of Amsterdam

Wilker Aziz ILLC, University of Amsterdam

André Martins Instituto de Telecomunicações & LUMLIS & Unbabel

Latent Variable Models

Latent variable z can be

Latent Variable Models

Latent variable z can be continuous



Source: Bouges et al., 2013

Latent Variable Models

Latent variable z can be **continuous**, **discrete**



Source: valleyeyecareaz.com

Latent Variable Models

Latent variable z can be continuous, discrete, or structured



Source: Liu et al., 2015

Training Discrete or Structured Latent Variable Models

Latent variable z can be

Training Discrete or Structured Latent Variable Models

Latent variable z can be discrete



Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**



Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z



Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z



Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z

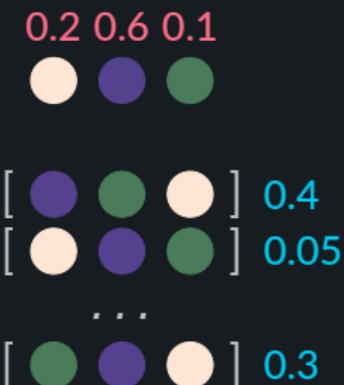


Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z

$\ell(x, z; \theta)$: downstream loss: ELBO, Log-Likelihood, (...)



Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z

$\ell(x, z; \theta)$: downstream loss: ELBO, Log-Likelihood, (...)

To train, we need to compute the following expectation:



Training Discrete or Structured Latent Variable Models

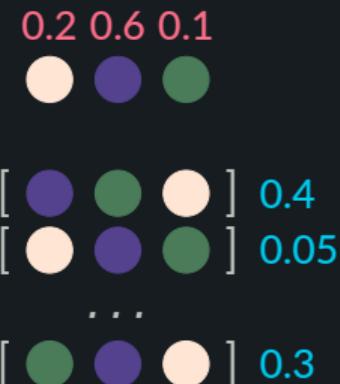
Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z

$\ell(x, z; \theta)$: downstream loss: ELBO, Log-Likelihood, (...)

To train, we need to compute the following expectation:

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta)$$

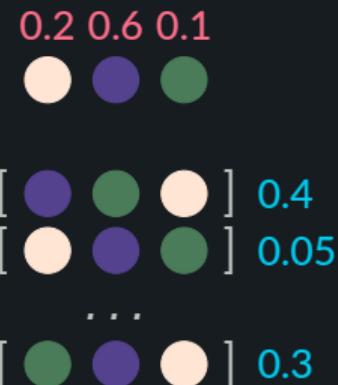


Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z

$\ell(x, z; \theta)$: downstream loss: ELBO, Log-Likelihood, (...)



To train, we need to compute the following expectation:

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta)$$

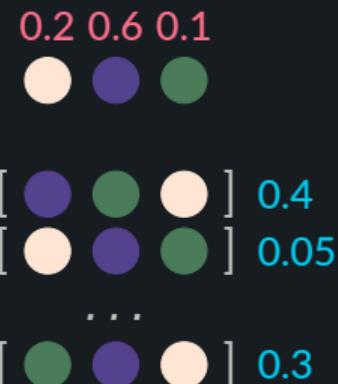
If \mathcal{Z} is **large**, this sum can get very expensive due to $\ell(x, z; \theta)$! 🍻

Training Discrete or Structured Latent Variable Models

Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: distribution over possible z

$\ell(x, z; \theta)$: downstream loss: ELBO, Log-Likelihood, (...)



To train, we need to compute the following expectation:

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta)$$

If \mathcal{Z} is **combinatorial**, this can be intractable to compute!



Current Solutions

If \mathcal{Z} is large, exact gradient computation is prohibitive

Current Solutions

If \mathcal{Z} is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE)—unbiased but high variance

Current Solutions

If \mathcal{Z} is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE)—unbiased but high variance

Another option: Gumbel-Softmax—continuous relaxation, biased estimation

Current Solutions

If \mathcal{Z} is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE)—unbiased but high variance

Another option: Gumbel-Softmax—continuous relaxation, biased estimation

New option: use sparsity! 

Current Solutions

If \mathcal{Z} is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE)—unbiased but high variance

Another option: Gumbel-Softmax—continuous relaxation, biased estimation

New option: use sparsity! 

no need for sampling → no variance

Current Solutions

If \mathcal{Z} is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE)—unbiased but high variance

Another option: Gumbel-Softmax—continuous relaxation, biased estimation

New option: use sparsity! 

no need for sampling → no variance

no relaxation into the continuous space

Taking a step back...

Does the expectation over possible z need to be expensive?

Taking a step back...

Does the expectation over possible z need to be expensive?

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\ &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \pi(z_2|x, \theta) \ell(x, z_2; \theta) + \dots \\ &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \pi(z_N|x, \theta) \ell(x, z_N; \theta)\end{aligned}$$

Taking a step back...

Does the expectation over possible z need to be expensive?

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\ &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \pi(z_2|x, \theta) \ell(x, z_2; \theta) + \dots \\ &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \pi(z_N|x, \theta) \ell(x, z_N; \theta)\end{aligned}$$

Usually we normalize π with softmax $\propto \exp(s) \Rightarrow \pi(z_i|x, \theta) > 0$

Sparse normalizers

We use **sparsemax**, **top- k sparsemax** and **SparseMAP** to allow efficient marginalization

Sparse normalizers

We use **sparsemax**, **top- k sparsemax** and **SparseMAP** to allow efficient marginalization

These functions are able to assign **probabilities of exactly zero!**

Sparse normalizers

We use **sparsemax**, **top- k sparsemax** and **SparseMAP** to allow efficient marginalization

These functions are able to assign **probabilities of exactly zero!**

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\ &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \underbrace{\pi(z_2|x, \theta)}_{=0} \ell(x, z_2; \theta) + \dots \\ &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \underbrace{\pi(z_N|x, \theta)}_{=0} \ell(x, z_N; \theta)\end{aligned}$$

Sparse normalizers

We use **sparsemax**, **top- k sparsemax** and **SparseMAP** to allow efficient marginalization

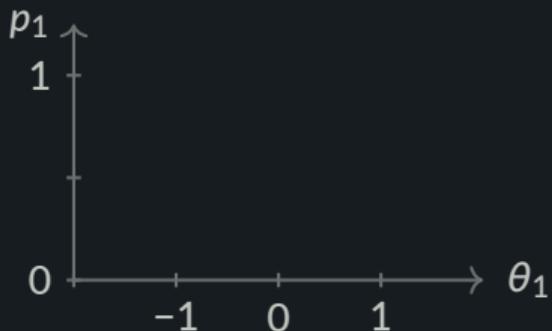
These functions are able to assign **probabilities of exactly zero!**

$$\begin{aligned}
 \mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\
 &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \underbrace{\pi(z_2|x, \theta) \ell(x, z_2; \theta)}_{=0} + \dots \\
 &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \underbrace{\pi(z_N|x, \theta) \ell(x, z_N; \theta)}_{=0}
 \end{aligned}$$

No need for computing $\ell(x, z; \theta)$ for all $z \in \mathcal{Z}$!

Discrete, unstructured case: sparsemax

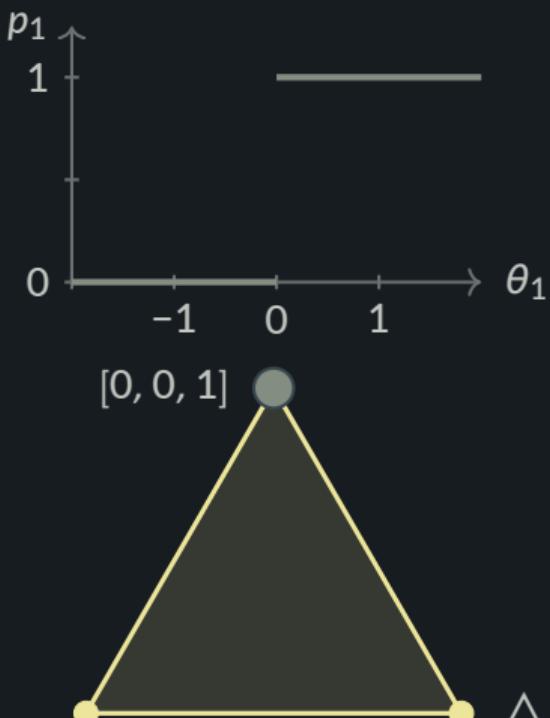
$$\boldsymbol{\pi}_\Omega(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} - \Omega(\mathbf{p})$$



Discrete, unstructured case: sparsemax

$$\boldsymbol{\pi}_\Omega(s) = \arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top s - \Omega(\boldsymbol{p})$$

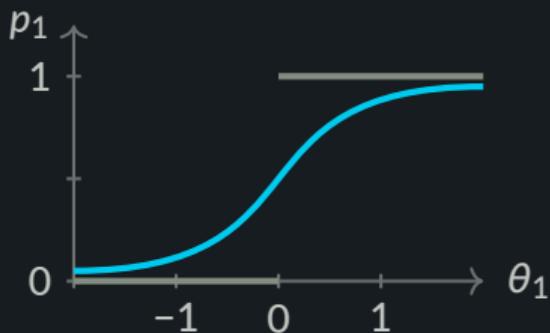
- argmax: $\Omega(\boldsymbol{p}) = 0$ (*no smoothing*)



Discrete, unstructured case: sparsemax

$$\boldsymbol{\pi}_\Omega(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} - \Omega(\mathbf{p})$$

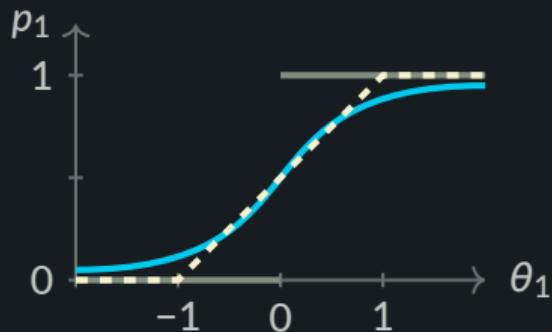
- argmax: $\Omega(\mathbf{p}) = 0$ (*no smoothing*)
- softmax: $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$



Discrete, unstructured case: sparsemax

$$\pi_{\Omega}(s) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T s - \Omega(\mathbf{p})$$

- argmax: $\Omega(\mathbf{p}) = 0$ (*no smoothing*)
- softmax: $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax: $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$



Semi-Supervised VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

Semi-Supervised VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \text{ discrete inference network} \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$


- Semi-Supervised VAE on MNIST: z is a possible digit

Semi-Supervised VAE

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z)$$

discrete inference network

sum over
the 10 digits

Gaussian VAE

- Semi-Supervised VAE on MNIST: z is a possible digit
- Train this with 10% labeled data

Semi-Supervised VAE

Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	94.75\pm.002	1
SFE+	96.53\pm.001	2
NVIL	96.01\pm.002	1
Gumbel	95.46\pm.001	1
<i>Marginalization</i>		
Dense	96.93\pm.001	10
Sparse (proposed)	96.87\pm.001	1.01\pm0.01

Semi-Supervised VAE

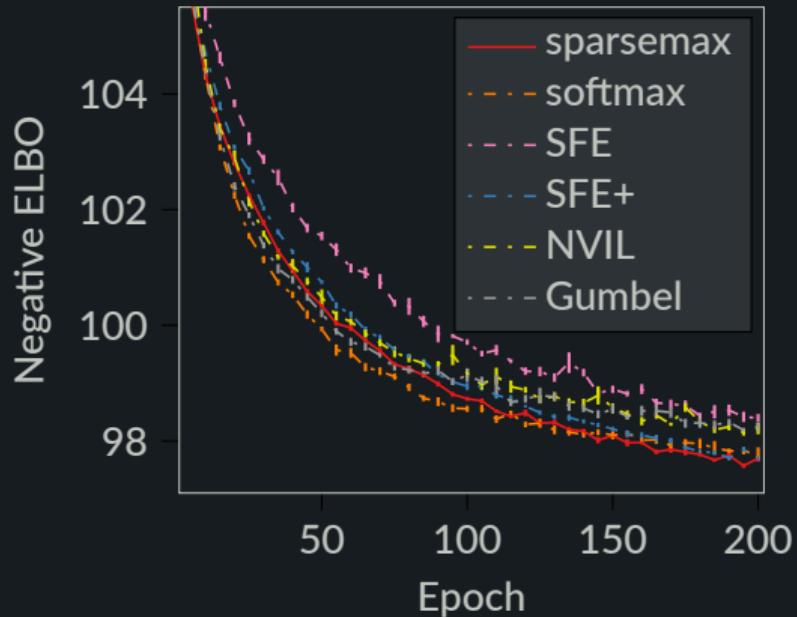
Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	94.75\pm.002	1
SFE+	96.53\pm.001	2
NVIL	96.01\pm.002	1
Gumbel	95.46\pm.001	1
<i>Marginalization</i>		
Dense	96.93\pm.001	10
Sparse (proposed)	96.87\pm.001	1.01\pm0.01

Semi-Supervised VAE

Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	94.75\pm.002	1
SFE+	96.53\pm.001	2
NVIL	96.01\pm.002	1
Gumbel	95.46\pm.001	1
<i>Marginalization</i>		
Dense	96.93\pm.001	10
Sparse (proposed)	96.87\pm.001	1.01\pm0.01

Semi-Supervised VAE

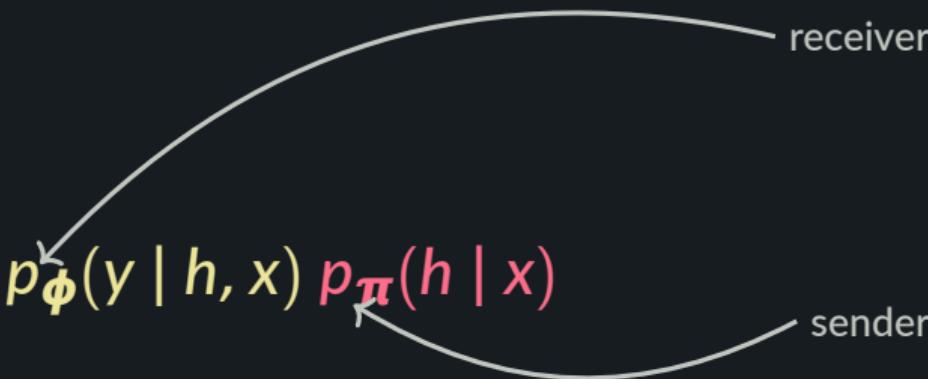
Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	$94.75 \pm .002$	1
SFE+	$96.53 \pm .001$	2
NVIL	$96.01 \pm .002$	1
Gumbel	$95.46 \pm .001$	1
<i>Marginalization</i>		
Dense	$96.93 \pm .001$	10
Sparse (proposed)	$96.87 \pm .001$	1.01 ± 0.01



Emergent communication

$$\begin{aligned} p(y \mid x) &= \sum_{h \in \mathcal{H}} p_{\phi}(y \mid h, x) \, p_{\pi}(h \mid x) \\ &= \mathbb{E}_{h \sim p_{\pi}(h|x)} p_{\phi}(y \mid h, x) \end{aligned}$$

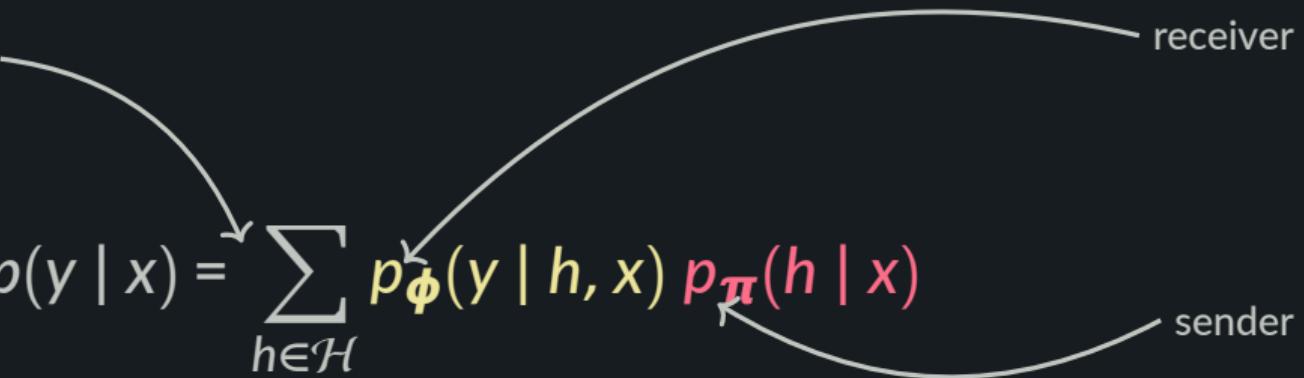
Emergent communication

$$\begin{aligned} p(y | x) &= \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x) \\ &= \mathbb{E}_{h \sim p_{\pi}(h|x)} p_{\phi}(y | h, x) \end{aligned}$$


- Emergent communication: h is a word from a big vocabulary.

Emergent communication

sum over
all possible messages

$$\begin{aligned} p(y | x) &= \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x) \\ &= \mathbb{E}_{h \sim p_{\pi}(h|x)} p_{\phi}(y | h, x) \end{aligned}$$


- Emergent communication: h is a word from a big vocabulary.

Emergent Communication

... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 ± 2.84	1
NVIL	37.04 ± 1.61	1
Gumbel	23.51 ± 16.19	1
ST Gumbel	27.42 ± 13.36	1

Marginalization

Emergent Communication

... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 ± 2.84	1
NVIL	37.04 ± 1.61	1
Gumbel	23.51 ± 16.19	1
ST Gumbel	27.42 ± 13.36	1
<i>Marginalization</i>		
Dense	93.37 ± 0.42	256

Emergent Communication

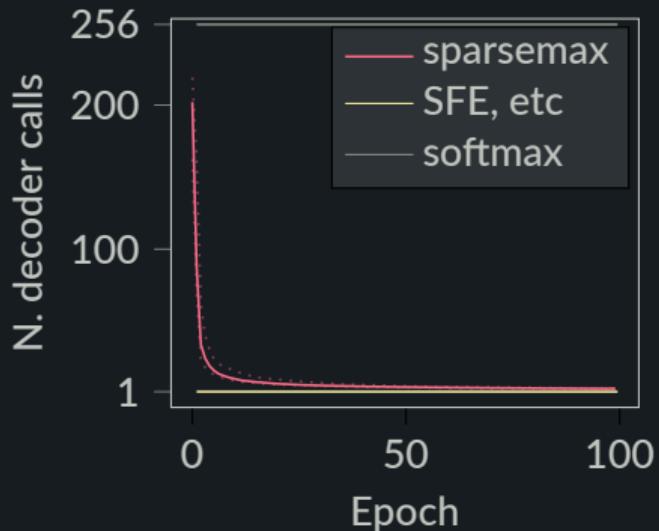
... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 ± 2.84	1
NVIL	37.04 ± 1.61	1
Gumbel	23.51 ± 16.19	1
ST Gumbel	27.42 ± 13.36	1
<i>Marginalization</i>		
Dense	93.37 ± 0.42	256
Sparse	93.35 ± 0.50	3.13 ± 0.48

Emergent Communication

... but make it harder: $|\mathcal{Z}| = 256$, $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 ± 2.84	1
NVIL	37.04 ± 1.61	1
Gumbel	23.51 ± 16.19	1
ST Gumbel	27.42 ± 13.36	1
<i>Marginalization</i>		
Dense	93.37 ± 0.42	256
Sparse	93.35 ± 0.50	3.13 ± 0.48



Limitations

- Mostly (and eventually) very sparse.
But fully dense worst case.
- For the same reason, sparsemax cannot handle structured z .

Limitations

- Mostly (and eventually) very sparse.
But fully dense worst case.
- For the same reason, sparsemax cannot handle structured z.

One solution: top-k sparsemax

$$k\text{-sparsemax}(s) = \arg \min_{p \in \Delta, \|p\|_0 \leq k} \|p - s\|_2^2$$

Limitations

- Mostly (and eventually) very sparse.
But fully dense worst case.
- For the same reason, sparsemax cannot handle structured z .

One solution: **top-k sparsemax**

$$k\text{-sparsemax}(s) = \arg \min_{\mathbf{p} \in \Delta, \|\mathbf{p}\|_0 \leq k} \|\mathbf{p} - s\|_2^2$$

- Non-convex but easy: sparsemax over the k highest scores (Kyrillidis et al., 2013).
- Top-k oracle available for some structured problems.
- Certificate: if at least one of the top- k h gets $p(h) = 0$, **k -sparsemax = sparsemax!**
thus, for latent variables: biased early on, but it goes away.

 Δ  \mathcal{M}

$$\mathcal{M} := \text{conv} \left\{ \mathbf{a}_h : h \in \mathcal{H} \right\}$$



Δ



\mathcal{M}

$$\begin{aligned}\mathcal{M} &:= \text{conv} \left\{ \mathbf{a}_h : h \in \mathcal{H} \right\} \\ &= \left\{ \mathbf{A} \mathbf{p} : \mathbf{p} \in \Delta \right\}\end{aligned}$$



Δ



\mathcal{M}

$$\begin{aligned}\mathcal{M} &:= \text{conv} \left\{ \mathbf{a}_h : h \in \mathcal{H} \right\} \\ &= \left\{ \mathbf{A} \mathbf{p} : \mathbf{p} \in \Delta \right\} \\ &= \left\{ \mathbb{E}_{H \sim p} \mathbf{a}_H : \mathbf{p} \in \Delta \right\}\end{aligned}$$



- $\text{argmax}_{p \in \Delta} p^T s$



Δ



\mathcal{M}

-

$$\text{argmax}_{\boldsymbol{p} \in \Delta} \arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s$$



$$\text{MAP}_{\boldsymbol{\mu} \in \mathcal{M}} \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$$

-



-

$$\text{argmax}_{\boldsymbol{p} \in \Delta} \arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s$$

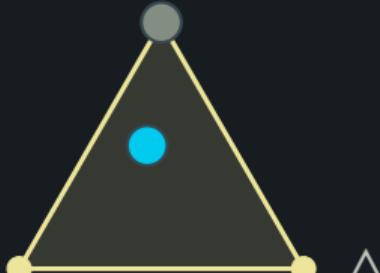


$$\text{MAP}_{\boldsymbol{\mu} \in \mathcal{M}} \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$$

-



- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s}$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s} + H(\boldsymbol{p})$



- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s}$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s} + H(\boldsymbol{p})$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **marginals** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



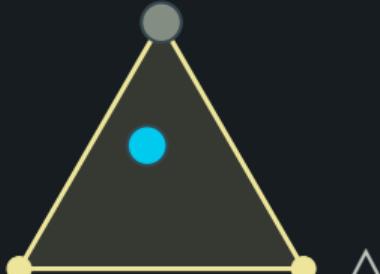
- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s}$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s} + H(\boldsymbol{p})$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **marginals** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s}$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s} + H(\boldsymbol{p})$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **marginals** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s}$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T \boldsymbol{s} + H(\boldsymbol{p})$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **marginals** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s + H(\boldsymbol{p})$
- **sparsemax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s - 1/2 \|\boldsymbol{p}\|^2$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **marginals** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



- **argmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s$
- **softmax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s + H(\boldsymbol{p})$
- **sparsemax** $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^T s - 1/2 \|\boldsymbol{p}\|^2$
- **MAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$
- **marginals** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$
- **SparseMAP** $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$



Bit-vector VAE

Key takeaways

References I

-  Bouges, Pierre, Thierry Chateau, Christophe Blanc, and Gaëlle Loosli (Dec. 2013). "Handling missing weak classifiers in boosted cascade: application to multiview and occluded face detection". In: *EURASIP Journal on Image and Video Processing* 2013, p. 55. DOI: [10.1186/1687-5281-2013-55](https://doi.org/10.1186/1687-5281-2013-55).
-  Correia, Gonçalo M., Vlad Niculae, Wilker Aziz, and André FT Martins (2020). "Efficient marginalization of discrete and structured latent variables via sparsity". In: *Proc. NeurIPS*.
-  Kyriolidis, Anastasios, Stephen Becker, Volkan Cevher, and Christoph Koch (2013). "Sparse projections onto the simplex". In: *Proc. ICML*.
-  Lazaridou, Angeliki, Alexander Peysakhovich, and Marco Baroni (2017). "Multi-agent cooperation and the emergence of (natural) language". In: *Proc. ICLR*.
-  Liu, Ziwei, Ping Luo, Xiaogang Wang, and Xiaoou Tang (Dec. 2015). "Deep Learning Face Attributes in the Wild". In: *Proceedings of International Conference on Computer Vision (ICCV)*.
-  Martins, André FT and Ramón Fernandez Astudillo (2016). "From softmax to sparsemax: A sparse model of attention and multi-label classification". In: *Proc. of ICML*.
-  Niculae, Vlad and Mathieu Blondel (2017). "A regularized framework for sparse and structured neural attention". In: *Proc. of NeurIPS*.
-  Niculae, Vlad, André FT Martins, Mathieu Blondel, and Claire Cardie (2018). "SparseMAP: Differentiable sparse structured inference". In: *Proc. of ICML*.