

Generalizations of Gumbel Softmax Trick

Based on “Gradient Estimation with Stochastic Softmax Tricks” by Paulus M. B. et al.

Motivation

Low-level Motivation

- x is a car, the latent variable z is
 - Real vector



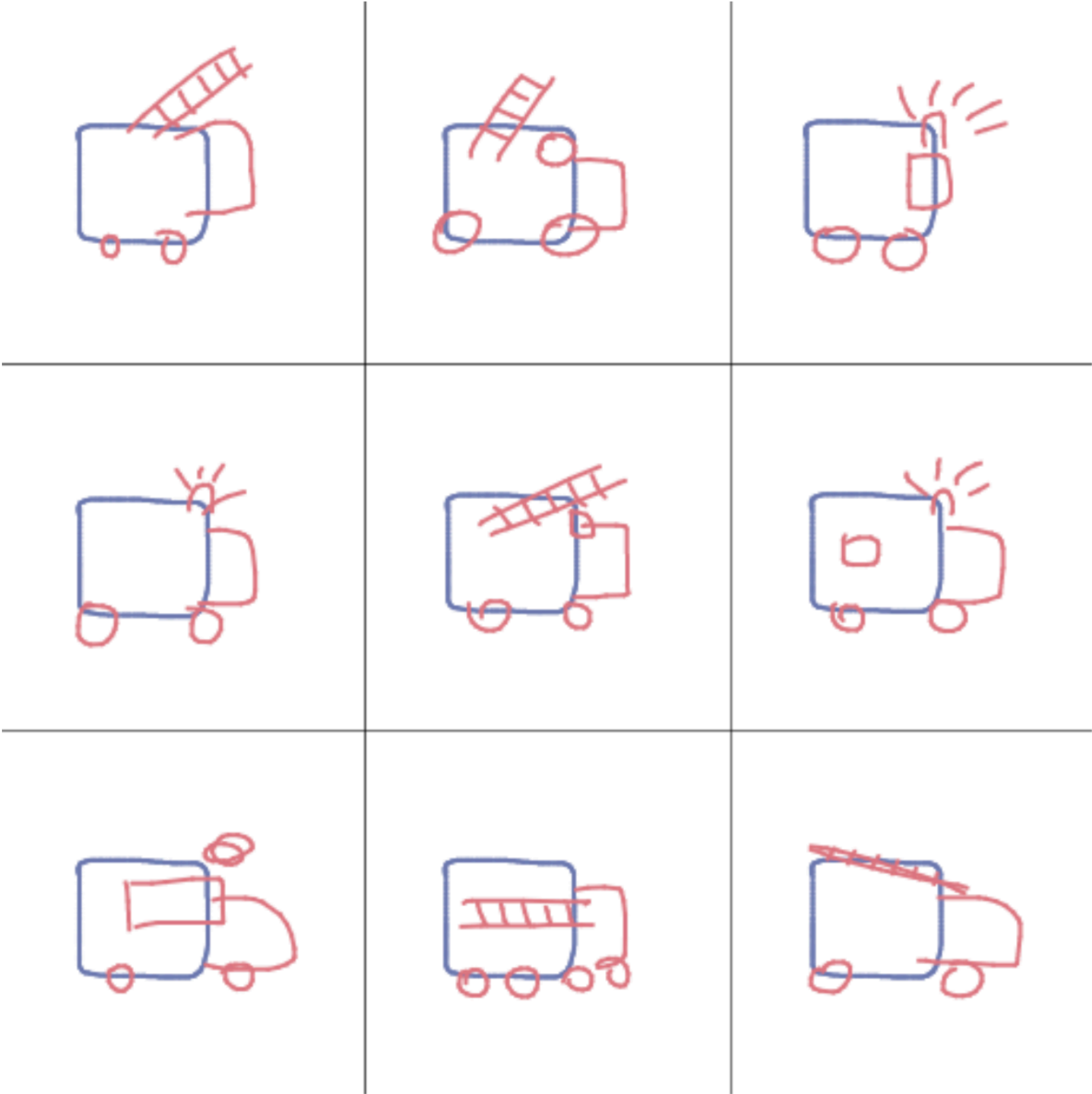
- Independent categorical

Ladder	Wheels	Siren	Visible Wheels
NO			
		YES	

- ... car scheme?

Ladder	Wheels	Siren	Visible Wheels
YES		YES	

A blue curved arrow points from the 'Wheels' column to the 'Visible Wheels' column. A red curved arrow points from the 'YES' entry in the 'Ladder' column to the right.

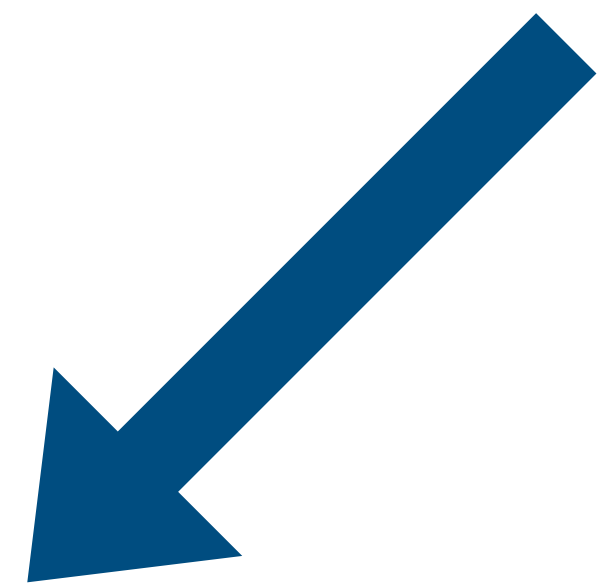


High Level Motivation

- What do we want?
 - Structured latent variables
 - Learn computational layers with discrete outputs
- Why do we want it?
 - Interpretability
 - Computational efficiency
 - Better generalization

Gumbel Softmax Trick Recap

$$\frac{d}{d\theta} \mathbb{E}_{q(z|\theta)} f(z) ?$$



REINFORCE

$$f(z) \frac{d}{d\theta} \log q(z | \theta)$$



Reparametrization trick

$$\frac{d}{d\theta} f(g_{\theta}(\epsilon))$$

Gumbel Argmax Trick

- $z \sim \text{Cat}(\text{soft max } \theta), \theta \in \mathbb{R}^d$
- Gumbel trick defines $g_\theta(\varepsilon)$ for z
 - Let $\varepsilon_i = -\log(-\log u_i)$ for $u \sim U[0,1]^d$
 - Let $g_\theta(\varepsilon) = \underset{i=1,\dots,d}{\operatorname{argmax}}(\theta + \varepsilon)$
 - Then $g_\theta(\varepsilon) \stackrel{d}{=} z$

$$\text{soft max}(\theta) = (0.1, 0.5, 0.4)$$

$$u = (0.4, 0.3, 0.7)$$

$$\begin{aligned} g_\theta(\varepsilon) &= \operatorname{arg max}(-2.2, -0.9, +0.1) \\ &= (0, 0, 1) \end{aligned}$$

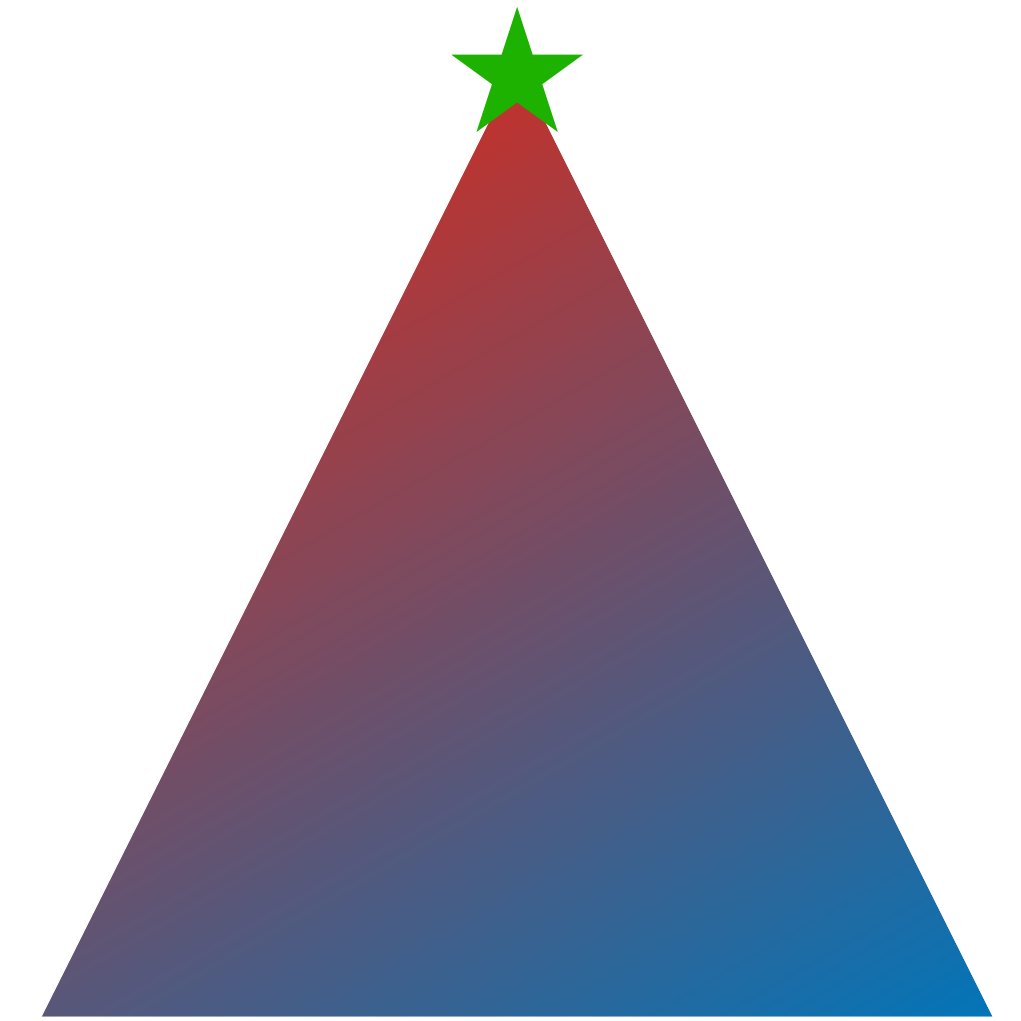
Softmax

- But $g_{\theta}(\varepsilon) = \arg \max(\theta + \varepsilon)$ is a piece-wise constant function of θ
- So $\mathbb{E}_{q(\varepsilon)} \frac{d}{d\theta} f(g_{\theta}(\varepsilon)) = 0 \neq \frac{d}{d\theta} \mathbb{E}_{q(z|\theta)} f(z)$
- **Idea:** replace $\arg \max(\cdot)$ with $\text{soft max}(\cdot)$
 - $\text{soft max}\left(\frac{\theta + \varepsilon}{T}\right) \xrightarrow{T \rightarrow 0} \arg \max(\theta + \varepsilon)$

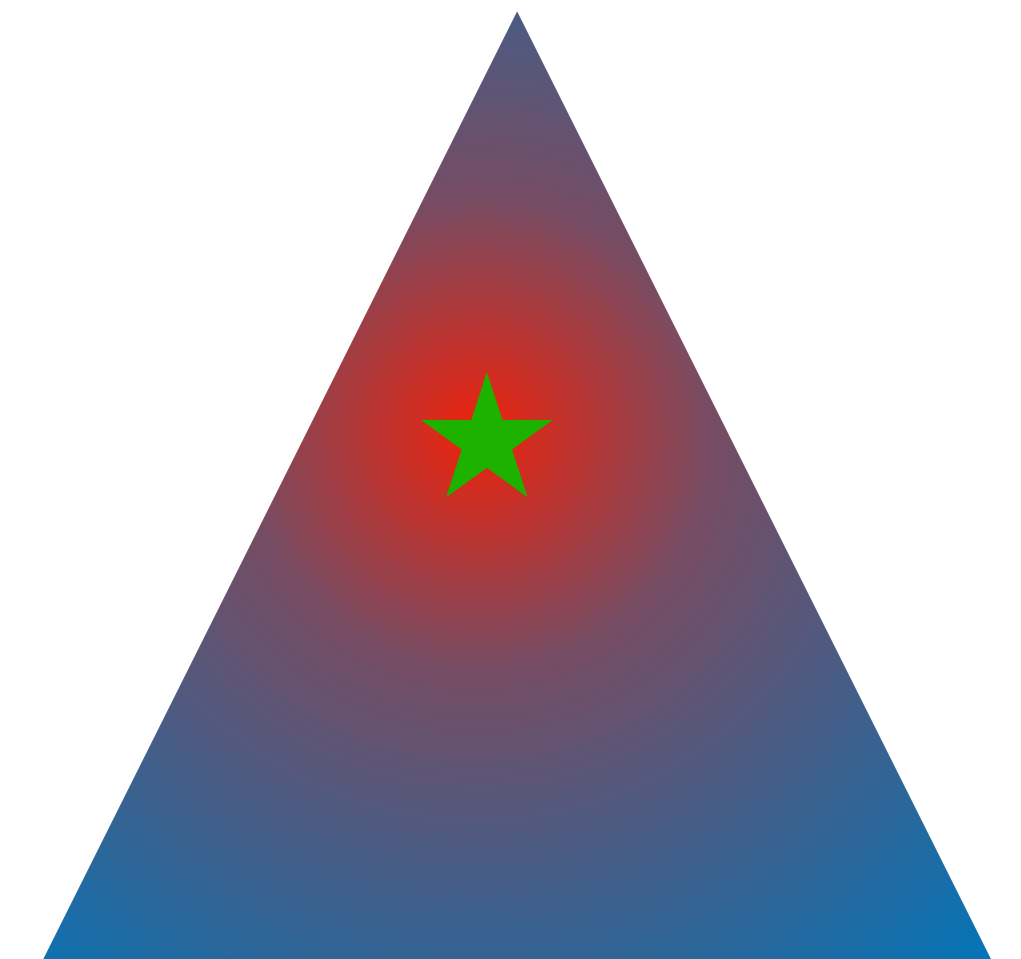
A Different View on Argmax and Softmax

- Rewrite $\operatorname{argmax}_{i=1,\dots,d} w = \operatorname{argmax}_{z \in \Delta^d} w^T z$
 - Δ^d is a convex hull of one-hots for z
- Then $\operatorname{soft\,max}(\frac{w}{T}) = \operatorname{argmax}_{z \in \Delta^d} (w^T z + \textcolor{red}{TH}(z))$
- $H(z) = - \sum_i z_i \log z_i$

arg max :



soft max :



A Different View on Gumbel Softmax Trick

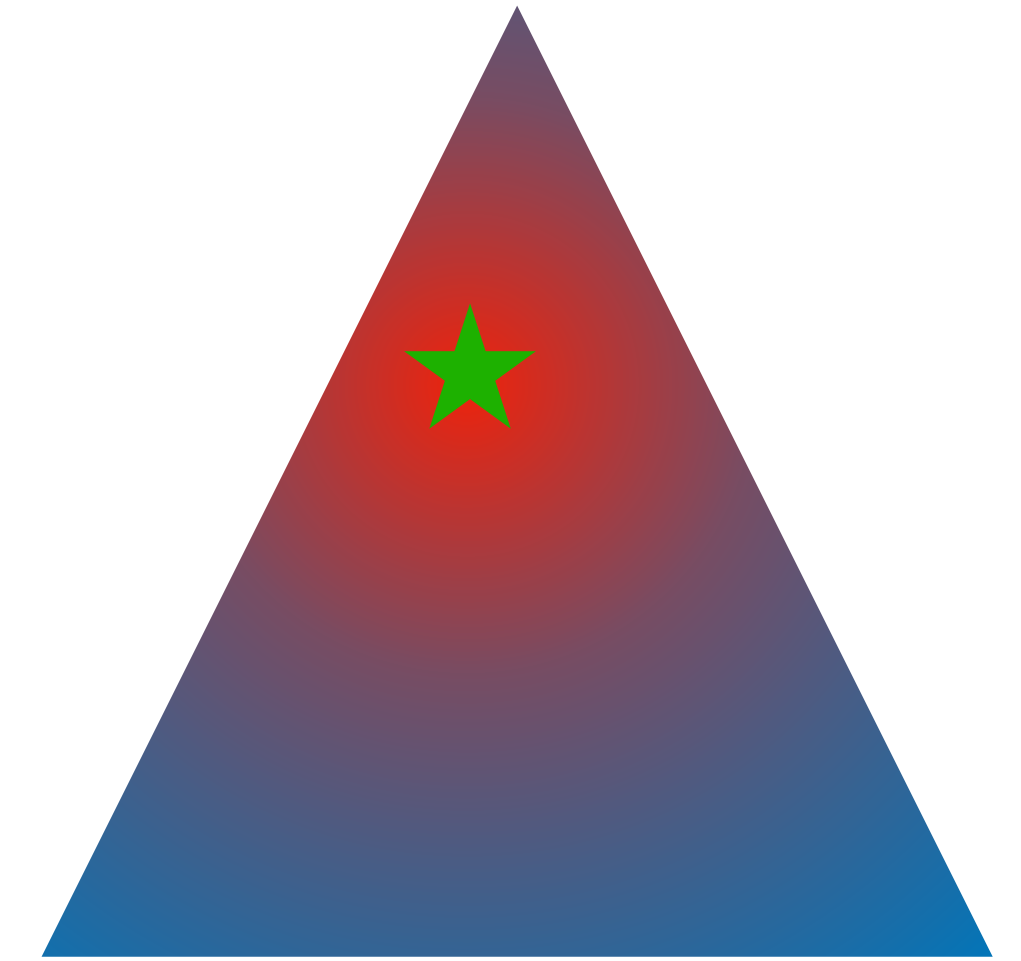
- Relaxed categorical sample is $z = \text{softmax}(\frac{\theta + \varepsilon}{T})$
- Equivalently $z = \underset{z \in \Delta^d}{\operatorname{argmax}}((\theta + \varepsilon)^T z + TH(z))$
 1. Perturb θ with ε
 2. Find arg max

Beyond Softmax

- We can consider any strongly convex function $H(\cdot)$

$$\text{sparse max}(w) = \operatorname{argmax}_{z \in \Delta^d} (w^T z - \frac{\|z\|^2}{2})$$

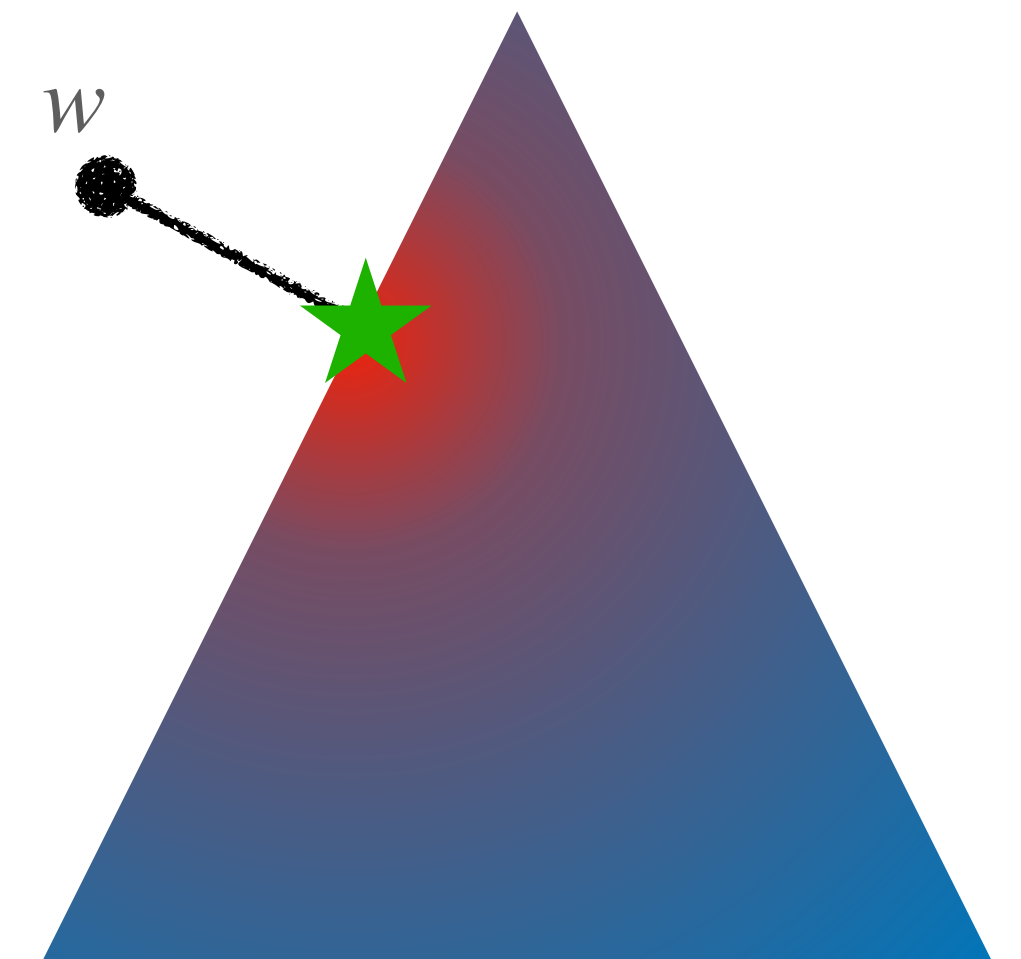
soft max :



- Finds sparse vectors:

$$w^T z - \frac{\|z\|^2}{2} = \|w - z\|^2 + \text{const}$$

soft max :



The Limitation of GST

- Time is $O(d)$
 - Perturb each θ_i and find max
- For combinatorial z the support is $d \gg 1$
 - Gumbel softmax trick is too slow

Stochastic Softmax Tricks

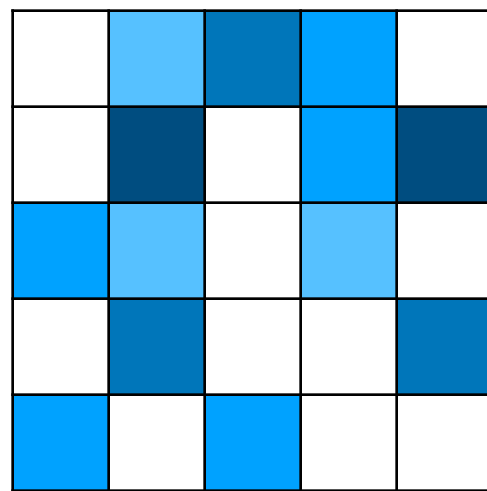
Stochastic Max Trick

	Gumbel Argmax Trick	Stochastic argmax Trick
Support	$Z = \{e_1, \dots, e_d\} \subset \mathbb{R}^d$	$Z = \{z_1, \dots, z_m\} \subset \mathbb{R}^d$
Perturbation	$w = \theta_i - \log(-\log(u_i)), u \sim U[0,1]^d$	$w = r_\theta(\varepsilon)$
Forward pass	$z = \operatorname{argmax}_{z' \in \operatorname{conv} Z} w^T z'$	$z = \operatorname{argmax}_{z' \in \operatorname{conv} Z} w^T z'$

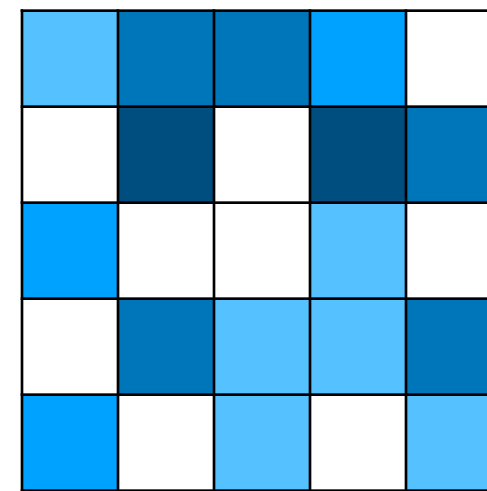
Example

- $Z = \{z_1, \dots, z_m\} \subset R^{n \times n}$ is a set of permutation matrices on n elements
- $m = n!$

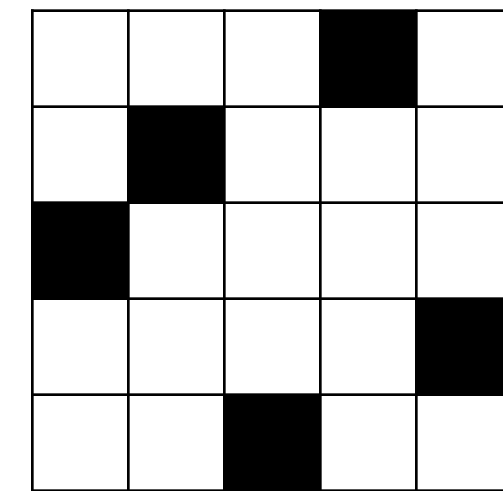
θ



$w = \theta + \varepsilon$



$\operatorname{argmax}_{z \in \operatorname{conv}(Z)} w^T z$



Regularization & Relaxation

Stochastic *Argmax* \rightarrow Stochastic *Softmax*

- Add a strongly convex regularizer $f : \mathbb{R}^d \rightarrow \{\mathbb{R}, \inf\}$

$$z = \operatorname{argmax}_{z' \in \operatorname{conv}(Z)} w^T z' \quad \rightarrow \quad z_T = \operatorname{argmax}_{z' \in \operatorname{conv}(Z)} (w^T z' - T f(z'))$$

Prop 1. If z is a.s. unique, then $\lim_{T \rightarrow 0} z_T = z$

Prop 2. z_T exists, is unique and differentiable in w

Backpropagation

- Unroll the optimizer
- Write custom backward pass
- Use finite difference approximation

$$\frac{d\mathcal{L}(z_t)}{dw} \approx \frac{z_t(w + \epsilon \partial \mathcal{L}(z_t) / \partial z_t) - z_t(w - \epsilon \partial \mathcal{L}(z_t) / \partial z_t)}{2\epsilon}$$

- $z_t(w + \epsilon \partial \mathcal{L}(z_t) / \partial z_t) = z_t + \epsilon \frac{\partial \mathcal{L}}{\partial z_t} \frac{dz_t}{dw} + o(\epsilon)$

Framework Requirements

- Inference
 - Reparametrized r.v. $w = r_{\theta}(\varepsilon)$
 - Solver for $\operatorname{argmax}_{z \in \operatorname{conv} Z} w^T z$
- Training
 - Strongly convex regularizer $f(z)$
 - Solver for $\operatorname{argmax}_{z \in \operatorname{conv} Z} (w^T z - tf(z))$

How to choose regularizer f ?

- Need to solve $\operatorname{argmax}_{z \in \operatorname{conv} Z} (w^T z - t f(z))$
- Euclidian projection $f(z) = \frac{\|z\|^2}{2}$
- Neg. entropy $f(z) = \sum z_i \log z_i$ or $f(z) = \sum_i (z_i \log z_i - (1 - z_i) \log(1 - z_i))$
- Exponential Families (next slides)

Exponential Family Reminder

- Consider $q(z | \theta) \propto I(z \in Z) \cdot \exp(z^T \theta)$
- Typical tasks:
 - MAP Inference: $\operatorname{argmax}_{z \in Z} \theta^T z$
 - Find log-partition: $A(\theta) = \log \sum_{z \in Z} \exp(\theta^T z)$
 - Nice trick: $\frac{\partial A(\theta)}{\partial \theta} = \mathbb{E}_{q(z|\theta)} z$
- Ex.: Categorical
 - $Z = \{e_1, \dots, e_d\}$
 - $A(\theta) = \log \sum \exp \theta_i$
 - $\frac{\partial A(\theta)}{\partial \theta} = \text{soft max}(\theta)$

Exponential Family Regularizer

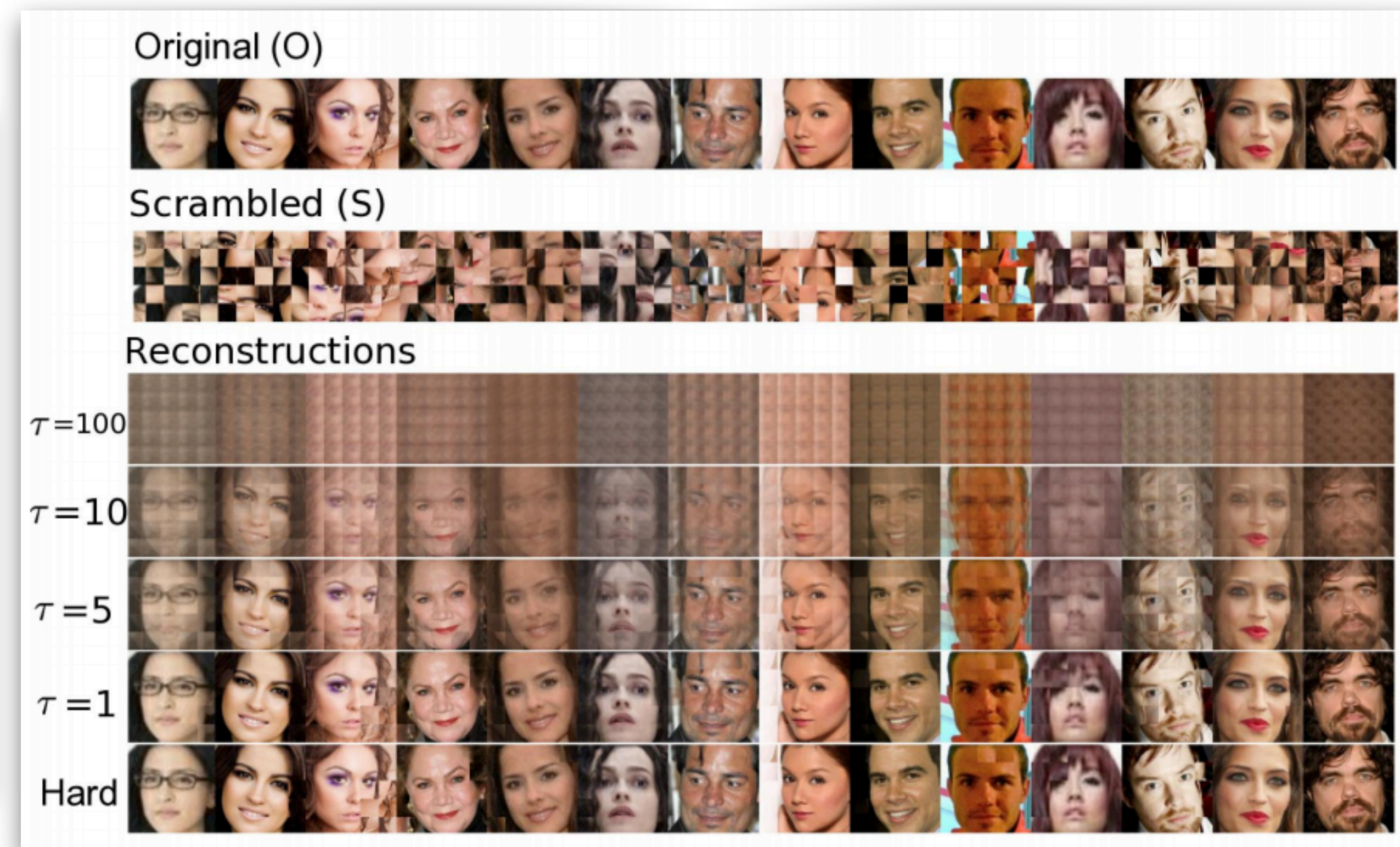
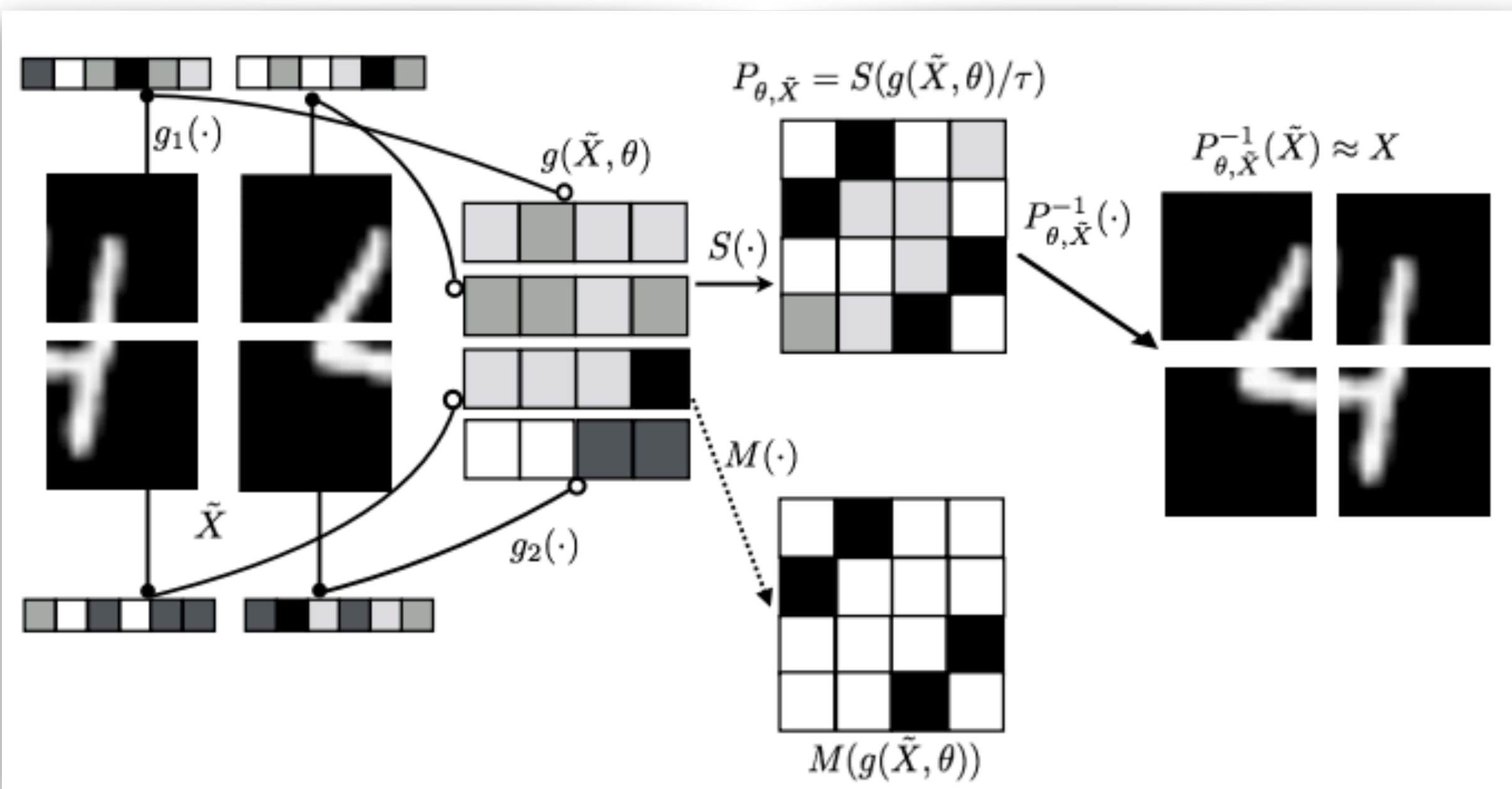
- Consider $A(\theta) = \log \sum_{z \in Z} \exp(\theta^T z)$
- Take $f(z) = A^*(z) = \sup_{\theta \in \Omega} (\theta^T z - A(\theta))$
- Softmax Trick:
 - $\operatorname{argmax}_{z \in \operatorname{conv}(Z)} (w^T z - f(z)) = \nabla_w A(w) = \mathbb{E}_{q(z|w)} z$
- Ex.: Gumbel Softmax
 - $A^*(z) = \sum_i z_i \log z_i$
 - $w = \theta - \log(-\log u)$
 - $\frac{\partial A}{\partial \theta} = \operatorname{soft max}(w)$

Applications

Gumbel Sinkhorn (1/4)

- Take permutation matrices as Z
 - Then $\text{conv}(Z)$ consists of doubly-stochastic
- **Hungarian algorithm** solves $\arg \max w^T z$
- Entropy $f(z) = \sum_{i,j} z_{i,j} \log z_{i,j}$
- **Sinkhorn algorithm** finds $\arg \max (w^T z - T \cdot f(z))$

Finding Latent Permutations: Jigsaw



Finding Latent Permutations: Connectomes

- Synthetic data $Y_t = PWP^T Y_{t-1} + \nu_t$
- Recover W

Prop. known neurons	40.%		30.%		20.%		10.%	
Difficulty	Easy	Hard	Easy	Hard	Easy	Hard	Easy	Hard
MCMC	.85	.82	.51	.44	.29	.27	.16	.12
(Linderman et al., 2017)	.97	.95	.90	.85	.77	.59	.39	.21
Gumbel-Sinkhorn	.97	.96	.92	.84	.76	.59	.44	.26
Gumbel-Sinkhorn, no regularization	.96	.93	.89	.78	.71	.52	.4	.23

k -subset Selection (2 / 4)

- $Z = \{z \in \{0,1\}^d \mid \sum z_i = k\}$



- Sort to solve $\arg \max w^T z$

- $Z = \{z \in \{0,1\}^{2d-1} \mid \sum_{i=1}^n z_i = k, z_i = z_{i-d}z_{i-d+1} \text{ for } d < i < 2d - 1\}$

- Dynamic programming for $\arg \max$

- Exponential family relaxation



BeerAdvocate Interpretability

Pours a slight tangerine orange and straw yellow. The head is nice and bubbly but fades very quickly with a little lacing. Smells like Wheat and European hops, a little yeast in there too. There is some fruit in there too, but you have to take a good whiff to get it. The taste is of wheat, a bit of malt, and a little fruit flavour in there too. Almost feels like drinking Champagne, medium mouthful otherwise. Easy to drink, but not something I'd be trying every night.

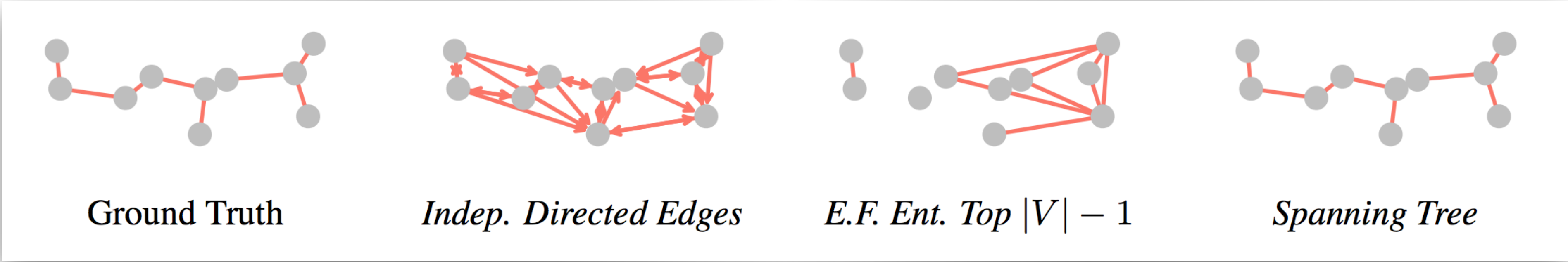
Appearance: 3.5 Aroma: 4.0 Palate: 4.5 Taste: 4.0 Overall: 4.0

Model	Relaxation	$k = 5$		$k = 10$		$k = 15$	
		MSE	Subs. Prec.	MSE	Subs. Prec.	MSE	Subs. Prec.
Simple	<i>L2X</i> [17]	3.6 ± 0.1	28.3 ± 1.7	3.0 ± 0.1	25.5 ± 1.2	2.6 ± 0.1	25.5 ± 0.4
	<i>SoftSub</i> [84]	3.6 ± 0.1	27.2 ± 0.7	3.0 ± 0.1	26.1 ± 1.1	2.6 ± 0.1	25.1 ± 1.0
	<i>Euclid. Top k</i>	3.5 ± 0.1	25.8 ± 0.8	2.8 ± 0.1	32.9 ± 1.2	2.5 ± 0.1	29.0 ± 0.3
	<i>Cat. Ent. Top k</i>	3.5 ± 0.1	26.4 ± 2.0	2.9 ± 0.1	32.1 ± 0.4	2.6 ± 0.1	28.7 ± 0.5
	<i>Bin. Ent. Top k</i>	3.5 ± 0.1	29.2 ± 2.0	2.7 ± 0.1	33.6 ± 0.6	2.6 ± 0.1	28.8 ± 0.4
	<i>E.F. Ent. Top k</i>	3.5 ± 0.1	28.8 ± 1.7	2.7 ± 0.1	32.8 ± 0.5	2.5 ± 0.1	29.2 ± 0.8
	<i>Corr. Top k</i>	2.9 ± 0.1	63.1 ± 5.3	2.5 ± 0.1	53.1 ± 0.9	2.4 ± 0.1	45.5 ± 2.7

Latent Spanning Trees (3 / 4)

- Take adjacency matrices of undirected trees as $Z \subset \mathbb{R}^{d \times d}$
- Consider $q(z \mid \theta) \propto I(z \in Z) \cdot \exp(\theta^T z)$
- **Kruskal algorithm** for $\arg \max w^T z$
- **Kirchhoff's theorem** computes $A(\theta) = \log \sum_{z \in Z} \exp(\theta^T z)$
- Use $\frac{\partial A}{\partial \theta}$ as the relaxed matrix

Dynamic Reconstruction

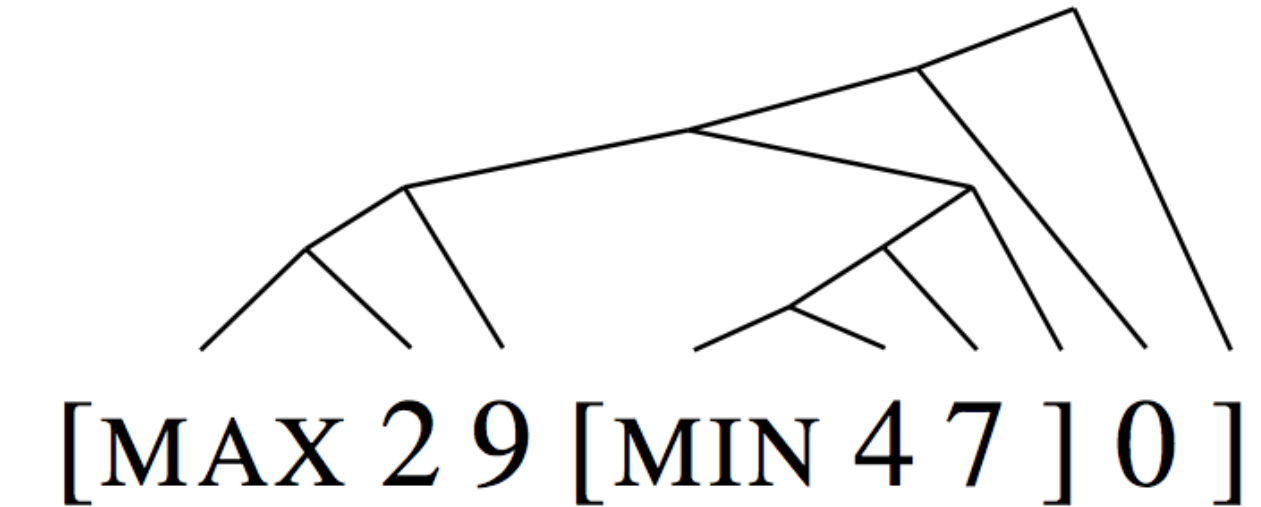


Edge Distribution	$T = 10$			$T = 20$		
	ELBO	Edge Prec.	Edge Rec.	ELBO	Edge Prec.	Edge Rec.
<i>Indep. Directed Edges</i> [38]	-1370 ± 20	48 ± 2	93 ± 1	-1340 ± 160	97 ± 3	99 ± 1
<i>E.F. Ent. Top $V - 1$</i>	-2100 ± 20	41 ± 1	41 ± 1	-1700 ± 320	98 ± 6	98 ± 6
<i>Spanning Tree</i>	-1080 ± 110	91 ± 3	91 ± 3	-1280 ± 10	99 ± 1	99 ± 1

Latent Arborescence (4 / 4)

- Take adjacency matrices of directed rooted trees as $Z \subset \mathbb{R}^{d \times d}$
- Consider $q(z \mid \theta) \propto I(z \in Z) \exp(\theta^T z)$
- **Edmonds algorithm** for $\arg \max w^T z$
- **Kirchhoff's theorem** computes $A(\theta) = \log \sum_{z \in Z} \exp(\theta^T z)$
- Use $\frac{\partial A}{\partial \theta}$ as the relaxed matrix

Unsupervised Parsing



- Corro & Titov relax Einser algorithm by replacing arg max with soft max
- Simplified ListOps dataset

Model	Edge Distribution	Task Acc.	Edge Precision	Edge Recall
LSTM	—	92.1 ± 0.2	—	—
GNN on latent graph	<i>Indep. Undirected Edges</i>	89.4 ± 0.6	20.1 ± 2.1	45.4 ± 6.5
	<i>Spanning Tree</i>	91.2 ± 1.8	33.1 ± 2.9	47.9 ± 5.2
GNN on latent digraph	<i>Indep. Directed Edges</i>	90.1 ± 0.5	13.0 ± 2.0	56.4 ± 6.7
	<i>Arborescence</i>			
	- Neg. Exp.	71.5 ± 1.4	23.2 ± 10.2	20.0 ± 6.0
	- Gaussian	95.0 ± 2.2	65.3 ± 3.7	60.8 ± 7.3
	- Gumbel	95.0 ± 3.0	75.5 ± 7.0	71.9 ± 12.4
	Ground Truth Edges	98.1 ± 0.1	100	100

Conclusions

- Learning structured latent variables is an active research direction
- *Stochastic softmax trick* generalize *gumbel softmax trick*
- Latent structure works as an *inductive bias*
- There is more to be done