

Adaptive Perturbation-Based Gradient Estimation for Discrete Latent Variable Models

Pasquale Minervini p.minervini@ed.ac.uk www.neuralnoise.com @pminervini



Mathias Niepert Luca Franceschi



Discrete Representations in Deep Neural Models



Discrete Representations in Deep Neural Models

• Combination of *discrete optimisation* algorithms (e.g., an ILP solver, Dijkstra's shortest path algorithm) and arbitrary neural modules





Discrete Representations in Deep Neural Models

 Combination of *discrete optimisation* algorithms (e.g., an ILP solver, Dijkstra's shortest path algorithm) and arbitrary neural modules



• Combination of *complex discrete probability distributions* (e.g., discrete attention distributions over the input, graph-structured representations) with arbitrary neural





1. A neural model assigns weights to cells in a map





1. A neural model assigns weights to cells in a map







- 1. A neural model assigns weights to cells in a map
- 2. Weights are used as input to a shortest path algorithm





- 1. A neural model assigns weights to cells in a map
- 2. Weights are used as input to a shortest path algorithm
- 3. The solver returns a shortest path which is then used by the downstream neural model for e.g., deciding the next action, or computing a loss function





Learning to Explain

1. A neural model assigns weights to words in the input

in the aroma, coffee and chocolate that is quite pronounced. in the taste, coffee, dry chocolate and a hit of hoppy bitterness. a small bite and medium bodied mouthfeel, with a dry roasty coffee in the aftertaste. a nice coffee and chocolate taste, nice hop presence, really freakin good.

input text



Learning to Explain

- 1. A neural model assigns weights to words in the input
- 2. Weights are used as parameters of a discrete distributions with a k-subset constraint

in the aroma , coffee and chocolate that is quite pronounced . in the taste , coffee , dry chocolate and a hit of hoppy bitterness . a small bite and medium bodied mouthfeel , with a dry roasty coffee in the aftertaste . a nice coffee and chocolate taste , nice hop presence , really freakin good .

input text



in the aroma , coffee and chocolate that is quite pronounced . in the taste , coffee , dry chocolate and a hit of hoppy bitterness . a small bite and medium bodied mouthfeel , with a dry roasty coffee in the aftertaste . a nice coffee and chocolate taste , nice hop presence , really freakin good .

Weight θ assigned to each token



Learning to Explain

- 1. A neural model assigns weights to words in the input
- 2. Weights are used as parameters of a discrete distributions with a k-subset constraint
- 3. A subset of k input words is sampled from this distribution, and used by a classification model to produce a prediction

in the aroma , coffee and chocolate that is quite pronounced . in the taste , coffee , dry chocolate and a hit of hoppy bitterness . a small bite and medium bodied mouthfeel , with a dry roasty coffee in the aftertaste . a nice coffee and chocolate taste , nice hop presence , really freakin good .

input text



in the aroma , coffee and chocolate that is quite pronounced . in the taste , coffee , dry chocolate and a hit of hoppy bitterness . a small bite and medium bodied mouthfeel , with a dry roasty coffee in the aftertaste . a nice coffee and chocolate taste , nice hop presence , really freakin good .

Weight θ assigned to each token



Sample discretely *exactly k* tokens



















We consider the problem of back-propagating through a combinatorial solver, i.e., a (black-box) component that solves the following problem:

Examples of combinatorial solvers:

- Top-k functions
- Shortest path algorithms
- Maximum spanning three algorithms
- Scheduling algorithms
- Supply chain optimisation algorithms
- .. and many more





Backpropagating through Combinatorial Solvers





Backpropagating through Combinatorial Solvers



Problem: How do we estimate the gradient of the loss $\nabla_{\theta} \mathscr{L}(\hat{\mathbf{y}}, \mathbf{y})$?



Combinatorial Solvers and Exponential Families

Implicit MLE is based on observing that calling a **combinatorial solver** $\mathbf{z} \leftarrow \text{solve}(\theta)$ on some **input weights** $\theta \in \mathbb{R}^n$ to obtain a **discrete solution** $\mathbf{z} \in \{0,1\}^n$ is equivalent to computing the MAP state of an **exponential family distribution** $p(\mathbf{z}; \theta)$:

e.g., top-k, shortest path algorithm, maximum spanning tree algorithm..

$$z \leftarrow solve(\theta)$$



Combinatorial Solvers and Exponential Families

Implicit MLE is based on observing that calling a **combinatorial solver** $\mathbf{z} \leftarrow \text{solve}(\theta)$ on some **input weights** $\theta \in \mathbb{R}^n$ to obtain a **discrete solution** $\mathbf{z} \in \{0,1\}^n$ is equivalent to computing the MAP state of an **exponential family distribution** $p(\mathbf{z}; \theta)$:

e.g., top-k, shortest path algorithm, maximum spanning tree algorithm..

$$z \leftarrow solve(\theta)$$

Maximum a Posteriori (MAP) estimation of z wrt. $p(\cdot; \theta)$

$$\mathbf{z} \leftarrow \arg \max_{\overline{\mathbf{z}}} p(\overline{\mathbf{z}}; \theta)$$



Combinatorial Solvers and Exponential Families

Implicit MLE is based on observing that calling a **combinatorial solver** $\mathbf{z} \leftarrow \text{solve}(\theta)$ on some **input weights** $\theta \in \mathbb{R}^n$ to obtain a **discrete solution** $\mathbf{z} \in \{0,1\}^n$ is equivalent to computing the MAP state of an **exponential family distribution** $p(\mathbf{z}; \theta)$:





22

Implicit MLE in a nutshell







23

Implicit MLE in a nutshell







24

Implicit MLE in a nutshell



$$\theta = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$



Implicit MLE in a nutshell





Implicit MLE in a nutshell





Constructing the Target Distribution $q(\mathbf{z}; \theta') = p(\mathbf{z}; \theta - \lambda \nabla_{\mathbf{z}} L(f_{\mathbf{u}}(\overline{\mathbf{z}}), y)$ How do we select the perturbation size λ ?



Constructing the Target Distribution $q(\mathbf{z}; \theta') = p(\mathbf{z}; \theta - \lambda \nabla_{\mathbf{z}} L(f_{\mathbf{u}}(\overline{\mathbf{z}}), y)$ How do we select the perturbation size λ ?

If $\lambda \approx 0$, the IMLE estimate of $\nabla_{\theta} \mathscr{L}(\hat{\mathbf{y}}, \mathbf{y})$ is **0**

If $\lambda \gg 0$, the IMLE estimate of $\nabla_{\theta} \mathscr{L}(\hat{\mathbf{y}}, \mathbf{y})$ is highly biased

How can we automatically select λ during training?



Impact of the Perturbation Size λ







Impact of the Perturbation Size λ





Impact of the Perturbation Size λ





Automatically Selecting λ

Very simple idea:



Automatically Selecting λ

Very simple idea:

• Initialise $\lambda \leftarrow 0$



Automatically Selecting λ

Very simple idea:

- Initialise $\lambda \leftarrow 0$
- During training, **adaptively change** λ until the gradient estimates satisfy some **desired sparsity criteria**



Relationships with Finite Difference Methods

Single-sample IMLE gradient estimate:

$$\nabla_{\theta} L \approx \frac{1}{\lambda} \left[\mathsf{MAP}(\theta + \epsilon) - \mathsf{MAP}(\theta - \lambda \nabla_{\mathbf{z}} L + \epsilon) \right]$$

Can be seen as a **forward** (one-sided) **finite difference approximation** in the form $f'(x) \approx [f(x+h) - f(x)]/h$.



Relationships with Finite Difference Methods

Single-sample IMLE gradient estimate:

$$\nabla_{\theta} L \approx \frac{1}{\lambda} \left[\mathsf{MAP}(\theta + \epsilon) - \mathsf{MAP}(\theta - \lambda \nabla_{\mathbf{z}} L + \epsilon) \right]$$

Can be seen as a **forward** (one-sided) **finite difference approximation** in the form $f'(x) \approx [f(x+h) - f(x)]/h$.

We can obtain a better approximation using a **centred** (two-sided) **difference formula**:

$$\nabla_{\theta} L \approx \frac{1}{2\lambda} \left[\mathsf{MAP}(\theta + \lambda \nabla_{\mathbf{z}} L) + \epsilon) - \mathsf{MAP}(\theta - \lambda \nabla_{\mathbf{z}} L) + \epsilon) \right]$$



Results - Estimator Bias and Variance





Results - Manual vs. Automatic λ **Selection**

Value of $\lambda \times$ Similarity to the true gradient 1.0 0.8 IMLE (Forward, S = 1000) **Cosine Similarity** IMLE (Central, S = 1000) 0.6 AIMLE (Forward, S = 1000) AIMLE (Central S = 1000) 0.4 STE (*S* = 1000) SFE (S = 1000)0.2 0.0 3



Results - Learning to Explain

 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 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

in the aroma, coffee and chocolate that is quite pronounced . in the taste, coffee, dry chocolate and a hit of hoppy bitterness . a small bite and medium bodied mouthfeel, with a dry roasty coffee in the aftertaste . a nice coffee and chocolate taste, nice hop presence, really freakin good.



in the aroma, coffee and chocolate that is quite pronounced. in the taste, coffee, dry chocolate and a hit of hoppy bitterness. a small bite and medium bodied mouthfeel, with a dry roasty coffee in the aftertaste. a nice coffee and chocolate taste, nice hop presence, really freakin good.





Results - Learning to Explain

Learning to Explain - Aroma, K = 10





Results - Discrete Variational Auto-Encoder





Results - Neural Relational Inference





Results - Neural Relational Inference

Observed dynamics		Interaction graph		Predicted dynamics		
T=10			T=20			
Edge Distribution	ELBO	Edge Prec.	Edge Rec.	ELBO	Edge Prec.	Edge Rec.
SST (Hard) IMLE (Forward) IMLE (Central) AIMLE (Forward) AIMLE (Central)	$\begin{array}{c} -2301.47 \pm 85.86 \\ -2289.94 \pm 4.31 \\ -2341.71 \pm 41.68 \\ \textbf{-1877.90} \pm \textbf{277.53} \\ -2018.39 \pm 357.16 \end{array}$	$\begin{array}{c} 33.75 \pm 9.44 \\ 23.94 \pm 0.03 \\ 43.95 \pm 7.22 \\ 55.23 \pm 11.86 \\ 29.32 \pm 6.89 \end{array}$	$\begin{array}{c} 60.40 \pm 23.23 \\ 95.75 \pm 0.14 \\ 43.95 \pm 7.22 \\ 55.23 \pm 11.86 \\ 41.83 \pm 21.51 \end{array}$	$\begin{array}{c} -3407.89 \pm 221.53 \\ -3820.68 \pm 25.32 \\ -3447.29 \pm 550.38 \\ \textbf{-1884.83} \pm \textbf{124.62} \\ -1999.57 \pm 856.27 \end{array}$	$\begin{array}{c} 57.40 \pm 17.87 \\ 20.28 \pm 0.12 \\ 40.25 \pm 14.26 \\ 40.48 \pm 4.25 \\ 70.89 \pm 24.77 \end{array}$	$\begin{array}{c} 70.42 \pm 8.22 \\ 20.28 \pm 0.12 \\ 40.25 \pm 14.26 \\ 40.48 \pm 4.25 \\ 83.73 \pm 1.31 \end{array}$



Summary

Implicit MLE [Niepert et al., NeurIPS 2021] allows you to back-propagate through arbitrary combinatorial solvers, achieving seamless neuro-symbolic integration

• Highly dependent on the perturbation size hyper-parameter λ



Summary

Implicit MLE [Niepert et al., NeurIPS 2021] allows you to back-propagate through arbitrary combinatorial solvers, achieving seamless neuro-symbolic integration

• Highly dependent on the perturbation size hyper-parameter λ

In this work, we propose **Adaptive IMLE**, a method for adaptively selecting λ during training based on analysing the sparsity of the gradient estimates

- Off-the-shelf: just add a decorator to your combinatorial solver to use it as a neural network layer
- Easy-to-use library available at github.com/EdinburghNLP/torch-adaptive-imle