

From Softmax to Sparsemax

A Sparse Model of Attention and Multi-Label Classification

André F. T. Martins and Ramón F. Astudillo

ICML 2016

Softmax

- ❖ Converting a representation vector into a posterior probabilities of labels $\mathbb{R}^K \rightarrow \Delta^{K-1}$, defined as:

$$\text{softmax}_i(z) = \frac{\exp(z_i)}{\sum_{k=1}^K \exp(z_k)}$$

- ❖ Limitation: full support, $\text{softmax}(z) > 0, \forall z$
- ❖ Requires: sparse probability distribution by assign exactly zero probability \rightarrow interpretability !

Sparsemax

- ❖ A sparse alternative: Euclidean projection of z onto the probability simplex

$$\text{sparsemax}(z) := \operatorname{argmin}_{p \in \Delta^{K-1}} \|p - z\|^2$$

- ❖ Close-Form Solution

$$\text{sparsemax}_i(z) = \max\{0, z_i - \tau\}$$

τ is a normalizing threshold function such that $\sum_j \max\{0, z_j - \tau\} = 1$

- ❖ How to compute τ ?

Sparsemax Evaluation

Algorithm 1 Sparsemax Evaluation

Input: \mathbf{z}

Sort \mathbf{z} as $z_{(1)} \geq \dots \geq z_{(K)}$

Sort the coordinates of \mathbf{z}

Find $k(\mathbf{z}) := \max \left\{ k \in [K] \mid 1 + kz_{(k)} > \sum_{j \leq k} z_{(j)} \right\}$

Define $k(\mathbf{z})$ the sparsity bound: index

Define $\tau(\mathbf{z}) = \frac{(\sum_{j \leq k(\mathbf{z})} z_{(j)}) - 1}{k(\mathbf{z})}$

Sparse logits $\{z_{(j)} \mid j \leq k(\mathbf{z})\}$ after

Output: \mathbf{p} s.t. $p_i = [z_i - \tau(\mathbf{z})]_+$.

threshold; sum of sparse logits $\sum_{j \leq k(\mathbf{z})} z_{(j)}$

$\tau(\mathbf{z})$ is the threshold function

$S(\mathbf{z}) := \left\{ j \in [K] \mid \text{sparsemax}_j(\mathbf{z}) > 0 \right\}$ the support of sparsemax

Two and Three-Dimensional Cases

- ❖ For two dimensions: $z = (t, 0)$, softmax becomes logistic (sigmoid) function as:

$$\text{softmax}_1(z) = (1 + \exp(-t))^{-1}$$

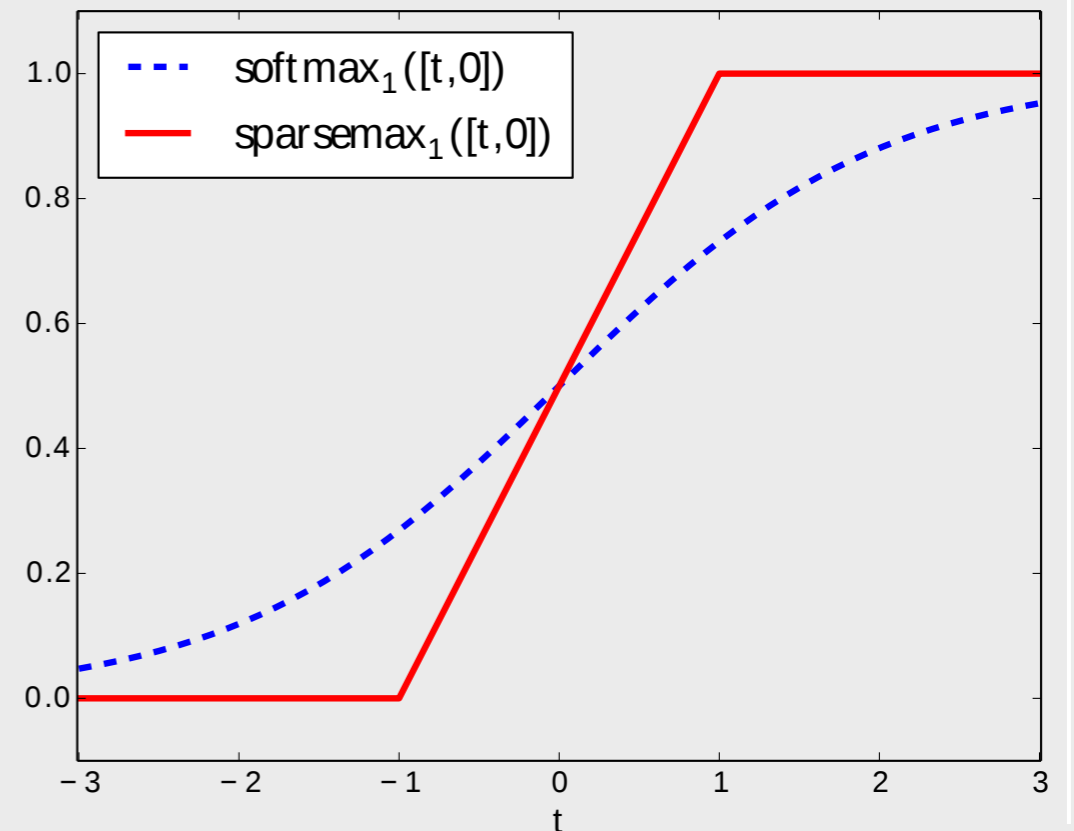
- ❖ 2D sparsemax is the “hard” version of the sigmoid

For $z = (t, 0)$

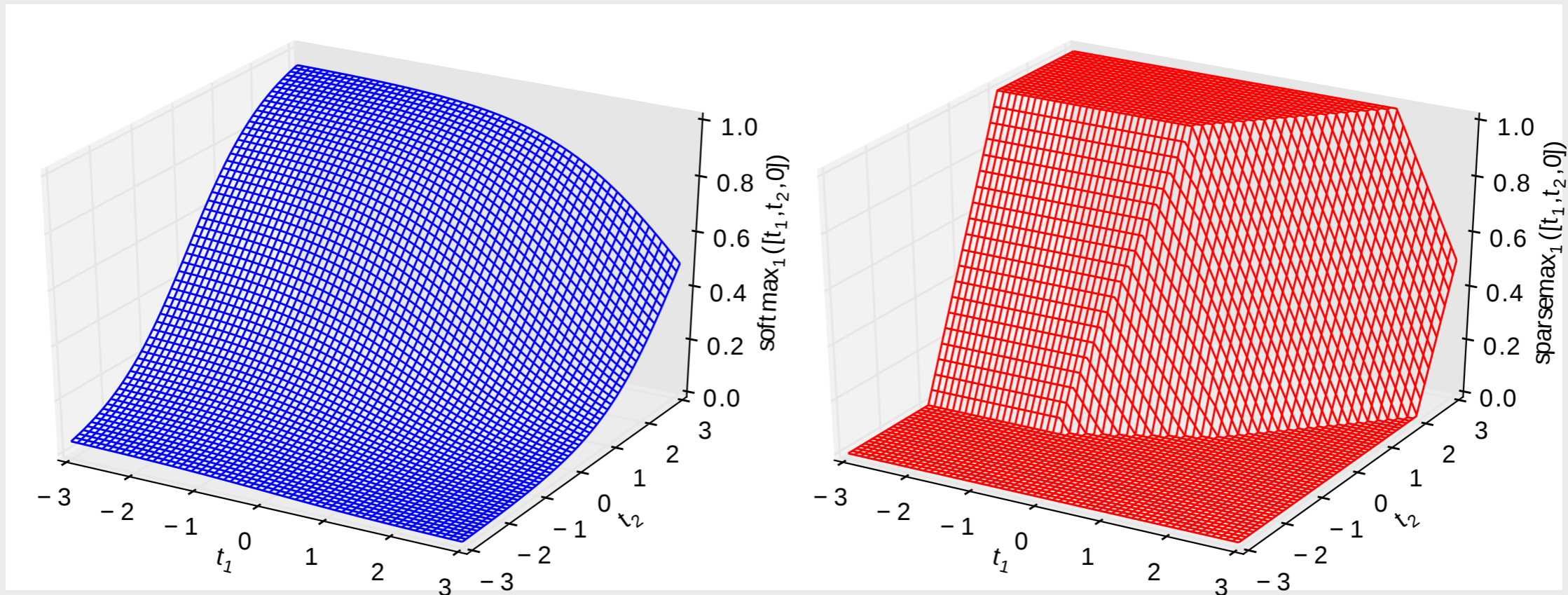
$$\tau(z) = \begin{cases} t - 1, & \text{if } t > 1 \\ (t - 1)/2, & \text{if } -1 \leq t \leq 1 \\ -1, & \text{if } t < -1 \end{cases}$$

$$\text{sparsemax}_1(z) = \begin{cases} 1, & \text{if } t > 1 \\ (t + 1)/2, & \text{if } -1 \leq t \leq 1 \\ 0, & \text{if } t < -1 \end{cases}$$

piece-wise linear



Two and Three-Dimensional Cases



5D case

```
Logits
tensor([[ 0.9549,  0.4015,  1.3101,  0.5750, -1.9022],
        [-1.7090, -0.5747, -0.1654,  0.1718,  0.1057]], device='cuda:0')

Softmax probabilities
tensor([[0.2672, 0.1536, 0.3811, 0.1827, 0.0153],
        [0.0465, 0.1447, 0.2179, 0.3052, 0.2857]], device='cuda:0')

Sparsemax probabilities
tensor([[0.3224, 0.0000, 0.6776, 0.0000, 0.0000],
        [0.0000, 0.0000, 0.1305, 0.4678, 0.4017]], device='cuda:0')
```

Gradient-based Optimization

❖ Jacobian of Softmax

$$\frac{\partial \text{softmax}_i(\mathbf{z})}{\partial z_j} = \frac{\delta_{ij} e^{z_i} \sum_k e^{z_k} - e^{z_i} e^{z_j}}{\left(\sum_k e^{z_k}\right)^2} = \text{softmax}_i(\mathbf{z}) \left(\delta_{ij} - \text{softmax}_j(\mathbf{z}) \right)$$

For matrix notation with $\mathbf{p} = \text{softmax}(\mathbf{z})$,

$$\mathbf{J}_{\text{softmax}}(\mathbf{z}) = \text{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$$

Gradient-based Optimization

❖ Jacobian of Sparsemax

$$1. \quad \frac{\partial \text{sparsemax}_i(\mathbf{z})}{\partial z_j} = \begin{cases} \delta_{ij} - \frac{\partial \tau(\mathbf{z})}{\partial z_j}, & \text{if } z_i > \tau(\mathbf{z}) \\ 0, & \text{if } z_i \leq \tau(\mathbf{z}) \end{cases}$$

$$2. \quad \frac{\partial \tau(\mathbf{z})}{\partial z_j} = \begin{cases} \frac{1}{|S(\mathbf{z})|} & \text{if } j \in S(\mathbf{z}) \\ 0, & \text{if } j \notin S(\mathbf{z}) \end{cases}$$

$$3. \quad \frac{\partial \text{sparsemax}_i(\mathbf{z})}{\partial z_j} = \begin{cases} \delta_{ij} - \frac{1}{|S(\mathbf{z})|}, & \text{if } i, j \in S(\mathbf{z}) \\ 0, & \text{otherwise} \end{cases}$$

$$4. \quad \mathbf{J}_{\text{sparsemax}}(\mathbf{z}) = \mathbf{D} \text{diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top / |S(\mathbf{z})| \quad \text{where } \mathbf{s} = \text{sparsemax}(\mathbf{z})$$

$$\text{sparsemax}_i(\mathbf{z}) = \max \{0, z_i - \tau\}$$

$$\tau(\mathbf{z}) = \frac{\left(\sum_{j \leq k(\mathbf{z})} z_{(j)} \right) - 1}{k(\mathbf{z})}$$

$$S(\mathbf{z}) := \left\{ j \in [K] \mid \text{sparsemax}_j(\mathbf{z}) > 0 \right\}$$

Loss Function: Logistic Loss

❖ Consider regularized empirical risk minimization problems

$$\text{minimize } \frac{\lambda}{2} \|\mathbf{W}\|_F^2 + \frac{1}{N} \sum_{i=1}^N L(\mathbf{W}\mathbf{x}_i + \mathbf{b}; y_i)$$

❖ Logistic loss

$$L_{\text{softmax}}(z; k) = -\log \text{softmax}_k(z) = -z_k + \log \sum_j \exp(z_j)$$

❖ Gradient

$$\nabla_z L_{\text{softmax}}(\mathbf{z}; k) = -\boldsymbol{\delta}_k + \text{softmax}(\mathbf{z})$$

$\boldsymbol{\delta}_k$ is the delta distribution on k .

Loss Function: Sparsemax Loss

- ❖ Reversing engineering the sparsemax loss

$$\nabla_{\mathbf{z}} L_{\text{sparsemax}}(\mathbf{z}; k) = -\delta_k + \text{sparsemax}(\mathbf{z})$$

- ❖ Sparsemax loss:

$$L_{\text{sparsemax}}(\mathbf{z}; k) = -z_k + \frac{1}{2} \sum_{j \in S(\mathbf{z})} \left(z_j^2 - \tau^2(\mathbf{z}) \right) + \frac{1}{2}$$

Generalization to Multi-Label Classification

❖ The multinomial logistic loss

$$L_{\text{softmax}}(\mathbf{z}; \mathbf{q}) = \mathbf{KL}(\mathbf{q} \parallel \text{softmax}(\mathbf{z})) = -\mathbf{H}(\mathbf{q}) - \mathbf{q}^\top \mathbf{z} + \log \sum_j \exp(z_j)$$

$$\text{Gradient: } \nabla_{\mathbf{z}} L_{\text{softmax}}(\mathbf{z}; \mathbf{q}) = -\mathbf{q} + \text{softmax}(\mathbf{z})$$

❖ The corresponding generalization in the sparsemax case

$$L_{\text{sparsemax}}(\mathbf{z}; \mathbf{q}) = -\mathbf{q}^\top \mathbf{z} + \frac{1}{2} \sum_{j \in S(\mathbf{z})} (z_j^2 - \tau^2(\mathbf{z})) + \frac{1}{2} \|\mathbf{q}\|^2$$

$$\text{Gradient: } \nabla_{\mathbf{z}} L_{\text{sparsemax}}(\mathbf{z}; \mathbf{q}) = -\mathbf{q} + \text{sparsemax}(\mathbf{z})$$

Experiments: Multi-Label Benchmarking

Table 1. Statistics for the 5 multi-label classification datasets.

DATASET	DESCR.	#LABELS	#TRAIN	#TEST
SCENE	IMAGES	6	1211	1196
EMOTIONS	MUSIC	6	393	202
BIRDS	AUDIO	19	323	322
CAL500	MUSIC	174	400	100
REUTERS	TEXT	103	23,149	781,265

❖ Logistic: independent
binary logistic
regressors on each label

Table 2. Micro (left) and macro-averaged (right) F_1 scores for the logistic, softmax, and sparsemax losses on benchmark datasets.

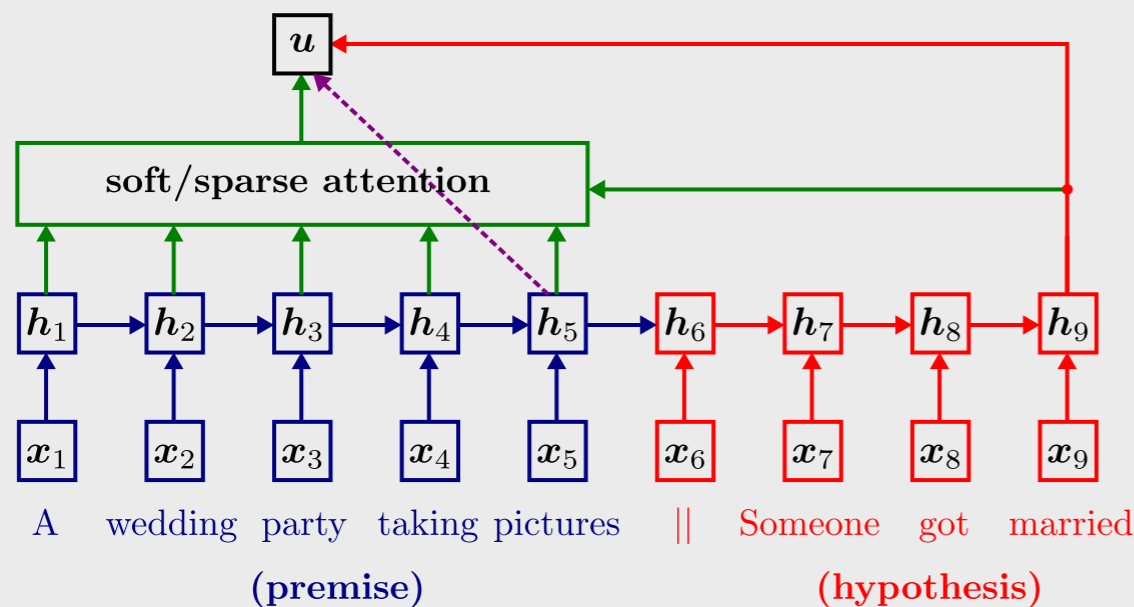
DATASET	LOGISTIC	SOFTMAX	SPARSEMAX
SCENE	70.96 / 72.95	74.01 / 75.03	73.45 / 74.57
EMOTIONS	66.75 / 68.56	67.34 / 67.51	66.38 / 66.07
BIRDS	45.78 / 33.77	48.67 / 37.06	49.44 / 39.13
CAL500	48.88 / 24.49	47.46 / 23.51	48.47 / 26.20
REUTERS	81.19 / 60.02	79.47 / 56.30	80.00 / 61.27

❖ Softmax: a multinomial
logistic regressor

❖ a slight advantage of
Sparsemax

Experiments: Neural Networks with Attention Mechanisms

❖ Sparse Attention



	DEV ACC.	TEST ACC.
NOATTENTION	81.84	80.99
LOGISTICATTENTION	82.11	80.84
SOFTATTENTION	82.86	82.08
SPARSEATTENTION	82.52	82.20

NoAttention:

$$u = \tanh(\mathbf{W}^{pu}h_L + \mathbf{W}^{hu}h_N + b^u)$$

SoftAttention:

$$z_t = v^\top \tanh(\mathbf{W}^{pm}h_t + \mathbf{W}^{hm}h_N + b^m)$$

$$p = \text{softmax}(z), \text{ where } z := (z_1, \dots, z_L)$$

A boy **rides on** a **camel** in a crowded area while talking on his cellphone.

Hypothesis: *A boy is riding an animal.* [entailment]

A young girl wearing a **pink coat** plays with a **yellow** toy golf club.

Hypothesis: *A girl is wearing a blue jacket.* [contradiction]

Two black dogs are **frolicking** around the **grass together**.

Hypothesis: *Two dogs swim in the lake.* [contradiction]

A man wearing a yellow striped shirt **laughs** while **seated next** to another **man** who is wearing a light blue shirt and **clasping** his hands together.

Hypothesis: *Two mimes sit in complete silence.* [contradiction]
