
From Softmax to Sparsemax: A Sparse Model of Attention and Multi-Label Classification

André F. T. Martins^{†‡}
Ramón F. Astudillo^{†◊}

ANDRE.MARTINS@UNBABEL.COM

RAMON@UNBABEL.COM

[†]Unbabel Lda, Rua Visconde de Santarém, 67-B, 1000-286 Lisboa, Portugal

[‡]Instituto de Telecomunicações (IT), Instituto Superior Técnico, Av. Rovisco Pais, 1, 1049-001 Lisboa, Portugal

[◊]Instituto de Engenharia de Sistemas e Computadores (INESC-ID), Rua Alves Redol, 9, 1000-029 Lisboa, Portugal

Abstract

We propose sparsemax, a new activation function similar to the traditional softmax, but able to output sparse probabilities. After deriving its properties, we show how its Jacobian can be efficiently computed, enabling its use in a network trained with backpropagation. Then, we propose a new smooth and convex loss function which is the sparsemax analogue of the logistic loss. We reveal an unexpected connection between this new loss and the Huber classification loss. We obtain promising empirical results in multi-label classification problems and in attention-based neural networks for natural language inference. For the latter, we achieve a similar performance as the traditional softmax, but with a selective, more compact, attention focus.

1. Introduction

The softmax transformation is a key component of several statistical learning models, encompassing multinomial logistic regression (McCullagh & Nelder, 1989), action selection in reinforcement learning (Sutton & Barto, 1998), and neural networks for multi-class classification (Bridle, 1990; Goodfellow et al., 2016). Recently, it has also been used to design attention mechanisms in neural networks, with important achievements in machine translation (Bahdanau et al., 2015), image caption generation (Xu et al., 2015), speech recognition (Chorowski et al., 2015), natural language understanding (Hermann et al., 2015; Sukhbaatar et al., 2015; Rocktäschel et al., 2015), and computation learning (Graves et al., 2014; Grefenstette et al., 2015).

There are a number of reasons why the softmax transfor-

mation is so appealing: it is simple to evaluate and differentiate, and it can be turned into a (concave) log-likelihood function by taking the logarithm of its output. Alternatives proposed in the literature, such as the Bradley-Terry model (Bradley & Terry, 1952; Zadrozny, 2001), the multinomial probit (Albert & Chib, 1993), the spherical softmax (Ollivier, 2013; de Brébisson & Vincent, 2015), or softmax approximations (Bouchard, 2007), while advantageous in certain scenarios, lack some of these convenient properties.

In this paper, we propose the **sparsemax transformation**. Sparsemax has the distinctive feature that it can return sparse posterior distributions, assigning zero probability to some output variables. This property makes it appealing for filtering large output spaces, predicting multiple labels, or as a component to identify which of a group of variables are relevant for a decision, making the model more interpretable. Crucially, this is done while preserving most of the attractive properties of softmax: we show that sparsemax is also simple to evaluate, it is even cheaper to differentiate, and that it can be turned into a convex loss function.

To sum up, our contributions are as follows:

- We formalize the new sparsemax transformation, derive its properties, and show how it can be efficiently computed (§2.1–2.3). We show that in the binary case sparsemax reduces to a hard sigmoid (§2.4).
- We derive the Jacobian of sparsemax, comparing it to the softmax case, and show that it can lead to faster gradient backpropagation (§2.5).
- We propose the **sparsemax loss**, a new loss function that is the sparsemax analogue of logistic regression (§3). We show that it is convex, everywhere differentiable, and can be regarded as a multi-class generalization of the Huber classification loss (Huber, 1964; Zhang, 2004).
- We apply the sparsemax loss to train multi-label linear classifiers (which predict a *set* of labels instead of a single label) on benchmark datasets (§4.1–4.2).

- Finally, we devise a neural selective attention mechanism using the sparsemax transformation, applying to a natural language inference problem (§4.3).

2. The Sparsemax Transformation

2.1. Definition

Let $\Delta^{K-1} := \{\mathbf{p} \in \mathbb{R}^K \mid \mathbf{1}^\top \mathbf{p} = 1, \mathbf{p} \geq \mathbf{0}\}$ be the $(K-1)$ -dimensional simplex. We are interested in maps from \mathbb{R}^K to Δ^{K-1} . Such maps are useful to convert a vector of real weights (e.g., label scores) to a probability distribution (e.g. posterior probabilities of labels). The classical example is the **softmax**, defined componentwise as:

$$\text{softmax}_i(\mathbf{z}) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}. \quad (1)$$

A limitation of the softmax transformation is that the resulting probability distribution always has full support, i.e., $\text{softmax}_i(\mathbf{z}) \neq 0$ for every \mathbf{z} and i . This is a disadvantage in applications where a sparse probability distribution is desired, in which case it is common to define a threshold below which small probability values are truncated to zero.

In this paper, we propose as an alternative the following transformation, which we call **sparsemax**:

$$\text{sparsemax}(\mathbf{z}) := \underset{\mathbf{p} \in \Delta^{K-1}}{\text{argmin}} \|\mathbf{p} - \mathbf{z}\|^2. \quad (2)$$

In words, sparsemax returns the Euclidean projection of the input vector \mathbf{z} onto the probability simplex. This projection is likely to hit the boundary of the simplex, in which case $\text{sparsemax}(\mathbf{z})$ becomes sparse. We will see that sparsemax retains most of the important properties of softmax, having in addition the ability of producing sparse distributions.

2.2. Closed-Form Solution

Many algorithms have been proposed to project onto the simplex (Michelot, 1986; Pardalos & Kovoor, 1990; Duchi et al., 2008). We start by recalling the result that such projections correspond to a soft-thresholding operation. Below, we denote $[K] := \{1, \dots, K\}$ and $[t]_+ := \max\{0, t\}$.

Proposition 1 *The solution of Eq. 2 is of the form:*

$$\text{sparsemax}_i(\mathbf{z}) = [z_i - \tau(\mathbf{z})]_+, \quad (3)$$

where $\tau : \mathbb{R}^K \rightarrow \mathbb{R}$ is the (unique) function that satisfies $\sum_j [z_j - \tau(\mathbf{z})]_+ = 1$ for every \mathbf{z} . Furthermore, τ can be expressed as follows. Let $z_{(1)} \geq z_{(2)} \geq \dots \geq z_{(K)}$ be the sorted coordinates of \mathbf{z} , and define $k(\mathbf{z}) := \max \left\{ k \in [K] \mid 1 + kz_{(k)} > \sum_{j \leq k} z_{(j)} \right\}$. Then,

$$\tau(\mathbf{z}) = \frac{\left(\sum_{j \leq k(\mathbf{z})} z_{(j)} \right) - 1}{k(\mathbf{z})} = \frac{\left(\sum_{j \in S(\mathbf{z})} z_j \right) - 1}{|S(\mathbf{z})|}, \quad (4)$$

Algorithm 1 Sparsemax Evaluation

Input: \mathbf{z}

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

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

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

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

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

Proof: See App. A.1 in the supplemental material. ■

In essence, Prop. 1 states that all we need for evaluating the sparsemax is to compute the threshold $\tau(\mathbf{z})$; all coordinates above this threshold (the ones in the set $S(\mathbf{z})$) will be shifted by this amount, and the others will be truncated to zero. We call τ in Eq. 4 the **threshold function**: this piecewise linear function will play an important role in the sequel. Alg. 1 illustrates a naïve $O(K \log K)$ algorithm that uses Prop. 1 for evaluating the sparsemax. More elaborate $O(K)$ algorithms exist based on partitioning and median-finding (Blum et al., 1973; Pardalos & Kovoor, 1990).¹

2.3. Basic Properties

We now highlight some properties that are common to softmax and sparsemax. Let $z_{(1)} := \max_k z_k$, and denote by $A(\mathbf{z}) := \{k \in [K] \mid z_k = z_{(1)}\}$ the set of maximal components of \mathbf{z} . We define the indicator vector $\mathbb{1}_{A(\mathbf{z})}$, whose k th component is 1 if $k \in A(\mathbf{z})$, and 0 otherwise. We further denote by $\gamma(\mathbf{z}) := z_{(1)} - \max_{k \notin A(\mathbf{z})} z_k$ the gap between the maximal components of \mathbf{z} and the second largest. We let $\mathbf{0}$ and $\mathbf{1}$ be vectors of zeros and ones, respectively.

Proposition 2 *The following properties hold for $\rho \in \{\text{softmax}, \text{sparsemax}\}$.*

1. $\rho(\mathbf{0}) = \mathbf{1}/K$ and $\lim_{\epsilon \rightarrow 0^+} \rho(\epsilon^{-1}\mathbf{z}) = \mathbb{1}_{A(\mathbf{z})}/|A(\mathbf{z})|$ (uniform distribution, and distribution peaked on the maximal components of \mathbf{z} , respectively). For sparsemax, the last equality holds for any $\epsilon \leq \gamma(\mathbf{z}) \cdot |A(\mathbf{z})|$.
2. $\rho(\mathbf{z}) = \rho(\mathbf{z} + c\mathbf{1})$, for any $c \in \mathbb{R}$ (i.e., ρ is invariant to adding the same constant to each coordinate).
3. $\rho(\mathbf{P}\mathbf{z}) = \mathbf{P}\rho(\mathbf{z})$ for any permutation matrix \mathbf{P} (i.e., ρ commutes with permutations).
4. If $z_i \leq z_j$, then $0 \leq \rho_j(\mathbf{z}) - \rho_i(\mathbf{z}) \leq \eta(z_j - z_i)$, where $\eta = \frac{1}{2}$ for softmax, and $\eta = 1$ for sparsemax.

Proof: See App. A.2 in the supplemental material. ■

¹In practice, there are *observed* linear time algorithms which are often a better choice. See Condat (2014, Table 1), who provide an empirical comparison among several of these algorithms.

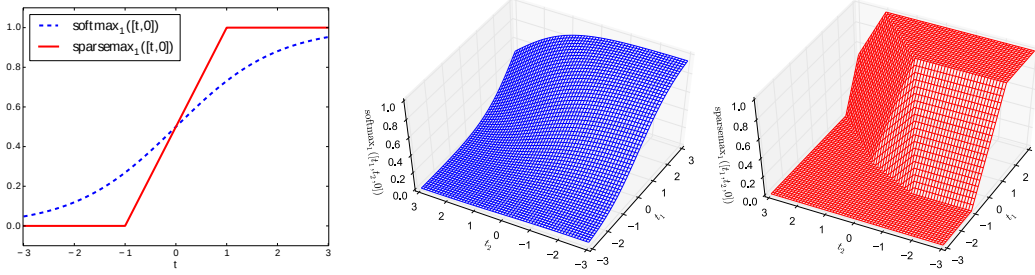


Figure 1. Comparison of softmax and sparsemax in 2D (left) and 3D (two rightmost plots).

Interpreting ϵ as a “temperature parameter,” the first clause of Prop. 2 shows that the sparsemax has the same “zero-temperature limit” behaviour as the softmax, but without the need of making the temperature arbitrarily small. We will use this fact in the experiments in §4.2 to control the degree of sparsity achieved with the sparsemax activation.

Prop. 2 is reassuring, since it shows that the sparsemax transformation, despite being defined very differently from the softmax, has a similar behaviour and preserves the same invariances. Note that some of these properties are not satisfied by other proposed replacements of the softmax: for example, the spherical softmax (Ollivier, 2013), defined as $\rho_i(\mathbf{z}) := z_i^2 / \sum_j z_j^2$, does not satisfy properties 2 and 4.

2.4. Two and Three-Dimensional Cases

For the two-class case, it is well known that the softmax activation becomes the logistic (sigmoid) function. More precisely, if $\mathbf{z} = (t, 0)$, then $\text{softmax}_1(\mathbf{z}) = \sigma(t) := (1 + \exp(-t))^{-1}$. We next show that the analogous in sparsemax is the “hard” version of the sigmoid. In fact, using Prop. 1, Eq. 4, we have that, for $\mathbf{z} = (t, 0)$,

$$\tau(\mathbf{z}) = \begin{cases} t - 1, & \text{if } t > 1 \\ (t - 1)/2, & \text{if } -1 \leq t \leq 1 \\ -1, & \text{if } t < -1, \end{cases} \quad (5)$$

and therefore

$$\text{sparsemax}_1(\mathbf{z}) = \begin{cases} 1, & \text{if } t > 1 \\ (t + 1)/2, & \text{if } -1 \leq t \leq 1 \\ 0, & \text{if } t < -1. \end{cases} \quad (6)$$

Fig. 1 provides an illustration for the two and three-dimensional cases. For the latter, we parameterize $\mathbf{z} = (t_1, t_2, 0)$ and plot $\text{softmax}_1(\mathbf{z})$ and $\text{sparsemax}_1(\mathbf{z})$ as a function of t_1 and t_2 . We can see that sparsemax is piecewise linear, but asymptotically similar to the softmax.

2.5. Jacobian of Sparsemax

The Jacobian matrix of a transformation ρ , $J_\rho(\mathbf{z}) := [\partial \rho_i(\mathbf{z}) / \partial z_j]_{i,j}$, is of key importance to train models with gradient-based optimization. We next derive the Jacobian

of the sparsemax activation, but before doing so, let us recall how the Jacobian of the softmax looks like. We have

$$\begin{aligned} \frac{\partial \text{softmax}_i(\mathbf{z})}{\partial z_j} &= (\delta_{ij} e^{z_i} \sum_k e^{z_k} - e^{z_i} e^{z_j}) / (\sum_k e^{z_k})^2 \\ &= \text{softmax}_i(\mathbf{z})(\delta_{ij} - \text{softmax}_j(\mathbf{z})), \end{aligned} \quad (7)$$

where δ_{ij} is the Kronecker delta, which evaluates to 1 if $i = j$ and 0 otherwise. Letting $\mathbf{p} = \text{softmax}(\mathbf{z})$, the full Jacobian can be written in matrix notation as

$$\mathbf{J}_{\text{softmax}}(\mathbf{z}) = \mathbf{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top, \quad (8)$$

where $\mathbf{Diag}(\mathbf{p})$ is a matrix with \mathbf{p} in the main diagonal.

Let us now turn to the sparsemax case. The sparsemax is differentiable everywhere except at splitting points \mathbf{z} where the support set $S(\mathbf{z})$ changes, *i.e.*, where $S(\mathbf{z}) \neq S(\mathbf{z} + \epsilon \mathbf{d})$ for some \mathbf{d} and infinitesimal ϵ .² From Eq. 3, we have that:

$$\frac{\partial \text{sparsemax}_i(\mathbf{z})}{\partial z_j} = \begin{cases} \delta_{ij} - \frac{\partial \tau(\mathbf{z})}{\partial z_j}, & \text{if } z_i > \tau(\mathbf{z}), \\ 0, & \text{if } z_i \leq \tau(\mathbf{z}). \end{cases} \quad (9)$$

From Eq. 4, the gradient of the threshold function τ is:

$$\frac{\partial \tau(\mathbf{z})}{\partial z_j} = \begin{cases} \frac{1}{|S(\mathbf{z})|} & \text{if } j \in S(\mathbf{z}), \\ 0, & \text{if } j \notin S(\mathbf{z}). \end{cases} \quad (10)$$

Note that $j \in S(\mathbf{z}) \Leftrightarrow z_j > \tau(\mathbf{z})$. Therefore we obtain:

$$\frac{\partial \text{sparsemax}_i(\mathbf{z})}{\partial z_j} = \begin{cases} \delta_{ij} - \frac{1}{|S(\mathbf{z})|}, & \text{if } i, j \in S(\mathbf{z}), \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

Let \mathbf{s} be an indicator vector whose i th entry is 1 if $i \in S(\mathbf{z})$, and 0 otherwise. We can write the Jacobian matrix as

$$\mathbf{J}_{\text{sparsemax}}(\mathbf{z}) = \mathbf{Diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top / |S(\mathbf{z})|. \quad (12)$$

It is instructive to compare Eqs. 8 and 12. We may regard the Jacobian of sparsemax as the Laplacian of a graph whose elements of $S(\mathbf{z})$ are fully connected. To compute

²For those points, we can take an arbitrary matrix in the set of generalized Clarke’s Jacobians (Clarke, 1983), the convex hull of all points of the form $\lim_{t \rightarrow \infty} \mathbf{J}_{\text{sparsemax}}(\mathbf{z}_t)$, where $\mathbf{z}_t \rightarrow \mathbf{z}$.

it, we only need $S(\mathbf{z})$, which can be obtained in $O(K)$ time with the same algorithm that evaluates the sparsemax.

Often, *e.g.*, in the gradient backpropagation algorithm, it is not necessary to compute the full Jacobian matrix, but only the product between the Jacobian and a given vector \mathbf{v} . In the softmax case, from Eq. 8, we have:

$$\mathbf{J}_{\text{softmax}}(\mathbf{z}) \cdot \mathbf{v} = \mathbf{p} \odot (\mathbf{v} - \bar{v}\mathbf{1}), \quad \text{with } \bar{v} := \sum_j p_j v_j, \quad (13)$$

where \odot denotes the Hadamard product; this requires a linear-time computation. For the sparsemax case, we have:

$$\mathbf{J}_{\text{sparsemax}}(\mathbf{z}) \cdot \mathbf{v} = \mathbf{s} \odot (\mathbf{v} - \hat{v}\mathbf{1}), \quad \text{with } \hat{v} := \frac{\sum_{j \in S(\mathbf{z})} v_j}{|S(\mathbf{z})|}. \quad (14)$$

Interestingly, if $\text{sparsemax}(\mathbf{z})$ has already been evaluated (*i.e.*, in the forward step), then so has $S(\mathbf{z})$, hence the nonzeros of $\mathbf{J}_{\text{sparsemax}}(\mathbf{z}) \cdot \mathbf{v}$ can be computed in only $O(|S(\mathbf{z})|)$ time, which can be sublinear. This can be an important advantage of sparsemax over softmax if K is large.

Computational efficiency. In sum, the computational needs of softmax and sparsemax are compared as follows:

- At training time, sparsemax can backpropagate gradients faster due to the sparsity (as stated above).
- At inference time, the softmax forward pass is faster, but asymptotically both are linear time. In our experiments (§4), even with the $O(K \log K)$ Algorithm 1, similar runtimes were achieved with softmax and sparsemax.

3. A Loss Function for Sparsemax

Now that we have defined and analyzed the sparsemax transformation, we will use it to design a new loss function that resembles the logistic loss, but can yield sparse posterior distributions. Later (in §4.1–4.2), we apply this loss to label proportion estimation and multi-label classification.

3.1. Logistic Loss

Let $\mathcal{D} := \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be a dataset, where each $\mathbf{x}_i \in \mathbb{R}^D$ is an input vector and $y_i \in \{1, \dots, K\}$ is a target label. We consider empirical risk minimization problems of the form

$$\begin{aligned} \text{minimize} \quad & \frac{\lambda}{2} \|\mathbf{W}\|_F^2 + \frac{1}{N} \sum_{i=1}^N L(\mathbf{W}\mathbf{x}_i + \mathbf{b}; y_i), \\ \text{w.r.t.} \quad & \mathbf{W} \in \mathbb{R}^{K \times D}, \quad \mathbf{b} \in \mathbb{R}^K, \end{aligned} \quad (15)$$

where L is a loss function, \mathbf{W} is a matrix of weights, and \mathbf{b} is a bias vector. The loss function associated with the softmax is the logistic loss (or negative log-likelihood):

$$\begin{aligned} L_{\text{softmax}}(\mathbf{z}; k) &= -\log \text{softmax}_k(\mathbf{z}) \\ &= -z_k + \log \sum_j \exp(z_j), \end{aligned} \quad (16)$$

where $\mathbf{z} = \mathbf{W}\mathbf{x}_i + \mathbf{b}$, and $k = y_i$ is the “gold” label. The gradient of this loss is, invoking Eq. 7,

$$\nabla_{\mathbf{z}} L_{\text{softmax}}(\mathbf{z}; k) = -\delta_k + \text{softmax}(\mathbf{z}), \quad (17)$$

where δ_k denotes the delta distribution on k , $[\delta_k]_j = 1$ if $j = k$, and 0 otherwise. This is a well-known result; when plugged into a gradient-based optimizer, it leads to updates that move probability mass from the distribution predicted by the current model (*i.e.*, $\text{softmax}_k(\mathbf{z})$) to the gold label (via δ_k). Can we have something similar for sparsemax?

3.2. Sparsemax Loss

A nice aspect of the log-likelihood (Eq. 16) is that adding up loss terms for several examples, assumed i.i.d, we obtain the log-probability of the full training data. Unfortunately, this idea cannot be carried out to sparsemax: now, some labels may have *exactly* probability zero, so any model that assigns zero probability to a gold label would zero out the probability of the entire training sample. This is of course highly undesirable. One possible workaround is to define

$$L_{\text{sparsemax}}^\epsilon(\mathbf{z}; k) = -\log \frac{\epsilon + \text{sparsemax}_k(\mathbf{z})}{1 + K\epsilon}, \quad (18)$$

where ϵ is a small constant, and $\frac{\epsilon + \text{sparsemax}_k(\mathbf{z})}{1 + K\epsilon}$ is a “perturbed” sparsemax. However, this loss is non-convex, unlike the one in Eq. 16.

Another possibility, which we explore here, is to construct an alternative loss function whose *gradient* resembles the one in Eq. 17. Note that the gradient is particularly important, since it is directly involved in the model updates for typical optimization algorithms. Formally, we want $L_{\text{sparsemax}}$ to be a differentiable function such that

$$\nabla_{\mathbf{z}} L_{\text{sparsemax}}(\mathbf{z}; k) = -\delta_k + \text{sparsemax}(\mathbf{z}). \quad (19)$$

We show below that this property is fulfilled by the following function, henceforth called the **sparsemax loss**:

$$L_{\text{sparsemax}}(\mathbf{z}; k) = -z_k + \frac{1}{2} \sum_{j \in S(\mathbf{z})} (z_j^2 - \tau^2(\mathbf{z})) + \frac{1}{2}, \quad (20)$$

where τ^2 is the square of the threshold function in Eq. 4. This loss, which has never been considered in the literature to the best of our knowledge, has a number of interesting properties, stated in the next proposition.

Proposition 3 *The following holds:*

1. $L_{\text{sparsemax}}$ is differentiable everywhere, and its gradient is given by the expression in Eq. 19.
2. $L_{\text{sparsemax}}$ is convex.
3. $L_{\text{sparsemax}}(\mathbf{z} + c\mathbf{1}; k) = L_{\text{sparsemax}}(\mathbf{z}; k)$, $\forall c \in \mathbb{R}$.

4. $L_{\text{sparsemax}}(\mathbf{z}; k) \geq 0$, for all \mathbf{z} and k .
5. The following statements are all equivalent: (i) $L_{\text{sparsemax}}(\mathbf{z}; k) = 0$; (ii) $\text{sparsemax}(\mathbf{z}) = \delta_k$; (iii) margin separation holds, $z_k \geq 1 + \max_{j \neq k} z_j$.

Proof: See App. A.3 in the supplemental material. ■

Note that the first four properties in Prop. 3 are also satisfied by the logistic loss, except that the gradient is given by Eq. 17. The fifth property is particularly interesting, since it is satisfied by the hinge loss of support vector machines. However, unlike the hinge loss, $L_{\text{sparsemax}}$ is everywhere differentiable, hence amenable to smooth optimization methods such as L-BFGS or accelerated gradient descent (Liu & Nocedal, 1989; Nesterov, 1983).

3.3. Relation to the Huber Loss

Coincidentally, as we next show, the sparsemax loss in the binary case reduces to the Huber classification loss, an important loss function in robust statistics (Huber, 1964).

Let us note first that, from Eq. 20, we have, if $|S(\mathbf{z})| = 1$,

$$L_{\text{sparsemax}}(\mathbf{z}; k) = -z_k + z_{(1)}, \quad (21)$$

and, if $|S(\mathbf{z})| = 2$, $L_{\text{sparsemax}}(\mathbf{z}; k) =$

$$-z_k + \frac{1 + (z_{(1)} - z_{(2)})^2}{4} + \frac{z_{(1)} + z_{(2)}}{2}, \quad (22)$$

where $z_{(1)} \geq z_{(2)} \geq \dots$ are the sorted components of \mathbf{z} . Note that the second expression, when $z_{(1)} - z_{(2)} = 1$, equals the first one, which asserts the continuity of the loss even though $|S(\mathbf{z})|$ is non-continuous on \mathbf{z} .

In the two-class case, we have $|S(\mathbf{z})| = 1$ if $z_{(1)} \geq 1 + z_{(2)}$ (unit margin separation), and $|S(\mathbf{z})| = 2$ otherwise. Assume without loss of generality that the correct label is $k = 1$, and define $t = z_1 - z_2$. From Eqs. 21–22, we have

$$L_{\text{sparsemax}}(t) = \begin{cases} 0 & \text{if } t \geq 1 \\ -t & \text{if } t \leq -1 \\ \frac{(t-1)^2}{4} & \text{if } -1 < t < 1, \end{cases} \quad (23)$$

whose graph is shown in Fig. 2. This loss is a variant of the Huber loss adapted for classification, and has been called “modified Huber loss” by Zhang (2004); Zou et al. (2006).

3.4. Generalization to Multi-Label Classification

We end this section by showing a generalization of the loss functions in Eqs. 16 and 20 to multi-label classification, *i.e.*, problems in which the target is a non-empty set of labels $\mathcal{Y} \in 2^{[K]} \setminus \{\emptyset\}$ rather than a single label.³ Such problems have attracted recent interest (Zhang & Zhou, 2014).

³Not to be confused with “multi-class classification,” which denotes problems with a target variable $y \in [K]$ where $K > 2$.

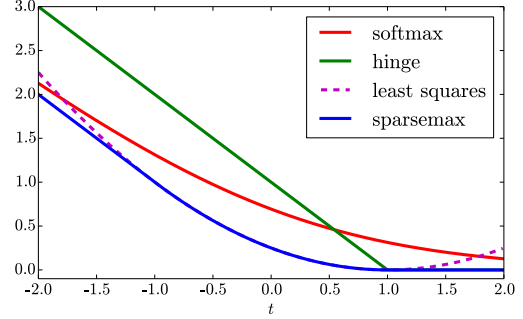


Figure 2. Comparison between the sparsemax loss and other commonly used losses for binary classification.

More generally, we consider the problem of estimating **sparse label proportions**, where the target is a probability distribution $\mathbf{q} \in \Delta^{K-1}$, such that $\mathcal{Y} = \{k \mid q_k > 0\}$. We assume a training dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{q}_i)\}_{i=1}^N$, where each $\mathbf{x}_i \in \mathbb{R}^D$ is an input vector and each $\mathbf{q}_i \in \Delta^{K-1}$ is a target distribution over outputs, assumed sparse.⁴ This subsumes single-label classification, where all \mathbf{q}_i are delta distributions concentrated on a single class. The generalization of the multinomial logistic loss to this setting is

$$\begin{aligned} L_{\text{softmax}}(\mathbf{z}; \mathbf{q}) &= \mathbf{KL}(\mathbf{q} \parallel \text{softmax}(\mathbf{z})) \\ &= -\mathbf{H}(\mathbf{q}) - \mathbf{q}^\top \mathbf{z} + \log \sum_j \exp(z_j), \end{aligned} \quad (24)$$

where $\mathbf{KL}(\cdot \parallel \cdot)$ and $\mathbf{H}(\cdot)$ denote the Kullback-Leibler divergence and the Shannon entropy, respectively. Note that, up to a constant, this loss is equivalent to standard logistic regression with soft labels. The gradient of this loss is

$$\nabla_{\mathbf{z}} L_{\text{softmax}}(\mathbf{z}; \mathbf{q}) = -\mathbf{q} + \text{softmax}(\mathbf{z}). \quad (25)$$

The corresponding generalization in the sparsemax case is:

$$L_{\text{sparsemax}}(\mathbf{z}; \mathbf{q}) = -\mathbf{q}^\top \mathbf{z} + \frac{1}{2} \sum_{j \in S(\mathbf{z})} (z_j^2 - \tau^2(\mathbf{z})) + \frac{1}{2} \|\mathbf{q}\|^2, \quad (26)$$

which satisfies the properties in Prop. 3 and has gradient

$$\nabla_{\mathbf{z}} L_{\text{sparsemax}}(\mathbf{z}; \mathbf{q}) = -\mathbf{q} + \text{sparsemax}(\mathbf{z}). \quad (27)$$

We make use of these losses in our experiments (§4.1–4.2).

4. Experiments

We next evaluate empirically the ability of sparsemax for addressing two classes of problems:

1. Label proportion estimation and multi-label classification, via the sparsemax loss in Eq. 26 (§4.1–4.2).
2. Attention-based neural networks, via the sparsemax transformation of Eq. 2 (§4.3).

⁴This scenario is also relevant for “learning with a probabilistic teacher” (Agrawala, 1970) and semi-supervised learning (Chapelle et al., 2006), as it can model label uncertainty.

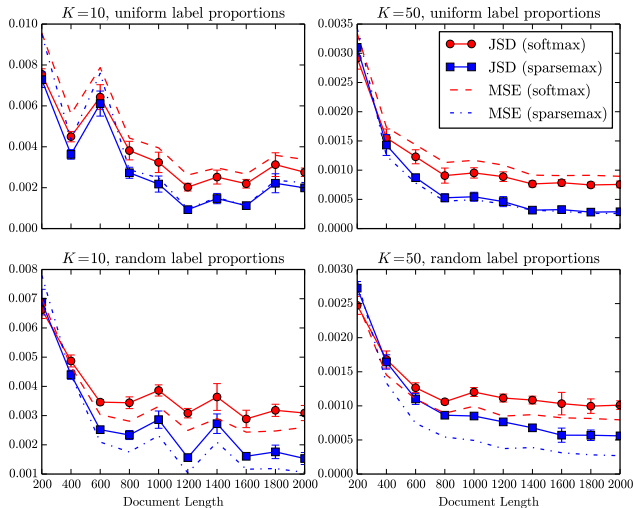


Figure 3. Simulation results for label proportion estimation with uniform (top) and random mixtures (bottom). Shown are the mean squared error and Jensen-Shannon divergence as a function of the document length, for the logistic and the sparsemax estimators.

4.1. Label Proportion Estimation

We show simulation results for sparse label proportion estimation on synthetic data. Since sparsemax can predict sparse distributions, we expect its superiority in this task. We generated datasets with 1,200 training and 1,000 test examples. Each example emulates a “multi-labeled document”: a variable-length sequence of word symbols, assigned to multiple topics (labels). We pick the number of labels $N \in \{1, \dots, K\}$ by sampling from a Poisson distribution with rejection sampling, and draw the N labels from a multinomial. Then, we pick a document length from a Poisson, and repeatedly sample its words from the mixture of the N label-specific multinomials. We experimented with two settings: **uniform mixtures** ($q_{k_n} = 1/N$ for the N active labels k_1, \dots, k_N) and **random mixtures** (whose label proportions q_{k_n} were drawn from a flat Dirichlet).⁵ We set the vocabulary size to be equal to the number of labels $K \in \{10, 50\}$, and varied the average document length between 200 and 2,000 words. We trained models by optimizing Eq. 15 with $L \in \{L_{\text{softmax}}, L_{\text{sparsemax}}\}$ (Eqs. 24 and 26). We picked the regularization constant $\lambda \in \{10^j\}_{j=-9}^0$ with 5-fold cross-validation.

Results are shown in Fig. 3. We report the mean squared error (average of $\|\mathbf{q} - \mathbf{p}\|^2$ on the test set, where \mathbf{q} and \mathbf{p} are respectively the target and predicted label posteriors) and the Jensen-Shannon divergence (average of $\mathbf{JS}(\mathbf{q}, \mathbf{p}) := \frac{1}{2}\mathbf{KL}(\mathbf{q} \parallel \frac{\mathbf{p} + \mathbf{q}}{2}) + \frac{1}{2}\mathbf{KL}(\mathbf{p} \parallel \frac{\mathbf{p} + \mathbf{q}}{2})$).⁶ We observe that the two

⁵Note that, with uniform mixtures, the problem becomes essentially multi-label classification.

⁶Note that the KL divergence is not an appropriate metric here,

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

Table 2. Micro (left) and macro-averaged (right) F₁ 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

losses perform similarly for small document lengths (where the signal is weaker), but as the average document length exceeds 400, the sparsemax loss starts outperforming the logistic loss consistently. This is because with a stronger signal the sparsemax estimator manages to identify correctly the support of the label proportions \mathbf{q} , contributing to reduce both the mean squared error and the JS divergence. This occurs both for uniform and random mixtures.

4.2. Multi-Label Classification on Benchmark Datasets

Next, we ran experiments in five benchmark multi-label classification datasets: the four small-scale datasets used by Koyejo et al. (2015),⁷ and the much larger Reuters RCV1 v2 dataset of Lewis et al. (2004).⁸ For all datasets, we removed examples without labels (*i.e.* where $\mathcal{Y} = \emptyset$). For all but the Reuters dataset, we normalized the features to have zero mean and unit variance. Statistics for these datasets are presented in Table 1.

Recent work has investigated the consistency of multi-label classifiers for various micro and macro-averaged metrics (Gao & Zhou, 2013; Koyejo et al., 2015), among which a plug-in classifier that trains independent binary logistic regressors on each label, and then tunes a probability threshold $\delta \in [0, 1]$ on validation data. At test time, those labels whose posteriors are above the threshold are predicted to be “on.” We used this procedure (called LOGISTIC) as a baseline for comparison. Our second baseline (SOFTMAX) is a multinomial logistic regressor, using the loss function

since the sparsity of \mathbf{q} and \mathbf{p} could lead to $-\infty$ values.

⁷Obtained from <http://mulan.sourceforge.net/datasets-mlc.html>.

⁸Obtained from <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multilabel.html>.

in Eq. 24, where the target distribution \mathbf{q} is set to uniform over the active labels. A similar probability threshold p_0 is used for prediction, above which a label is predicted to be “on.” We compare these two systems with the sparsemax loss function of Eq. 26. We found it beneficial to scale the label scores \mathbf{z} by a constant $t \geq 1$ at test time, before applying the sparsemax transformation, to make the resulting distribution $\mathbf{p} = \text{sparsemax}(t\mathbf{z})$ more sparse. We then predict the k th label to be “on” if $p_k \neq 0$.⁹

We optimized the three losses with L-BFGS (for a maximum of 100 epochs), tuning the hyperparameters in a held-out validation set (for the Reuters dataset) and with 5-fold cross-validation (for the other four datasets). The hyperparameters are the regularization constant $\lambda \in \{10^j\}_{j=-8}^2$, the probability thresholds $\delta \in \{.05 \times n\}_{n=1}^{10}$ for LOGISTIC and $p_0 \in \{n/K\}_{n=1}^{10}$ for SOFTMAX, and the coefficient $t \in \{0.5 \times n\}_{n=1}^{10}$ for SPARSEMAX.

Results are shown in Table 2. Overall, the performances of the three losses are all very similar, with a slight advantage of SPARSEMAX, which attained the highest results in 4 out of 10 experiments, while LOGISTIC and SOFTMAX won 3 times each. In particular, sparsemax appears better suited for problems with larger numbers of labels.

4.3. Neural Networks with Attention Mechanisms

We now assess the suitability of the sparsemax transformation to construct a “sparse” neural attention mechanism.

We ran experiments on the task of natural language inference, using the recently released SNLI 1.0 corpus (Bowman et al., 2015), a collection of 570,000 human-written English sentence pairs. Each pair consists of a premise and an hypothesis, manually labeled with one the labels ENTAILMENT, CONTRADICTION, or NEUTRAL. We used the provided training, development, and test splits.

The architecture of our system, shown in Fig. 4, is the same as the one proposed by Rocktäschel et al. (2015). We compare the performance of four systems: NOATTENTION, a (gated) RNN-based system similar to Bowman et al. (2015); LOGISTICATTENTION, an attention-based system with independent logistic activations; SOFTATTENTION, a near-reproduction of the Rocktäschel et al. (2015)’s attention-based system; and SPARSEATTENTION, which replaces the latter softmax-activated attention mechanism by a sparsemax activation.

We represent the words in the premise and in the hypothesis with 300-dimensional GloVe vectors (Pennington et al.,

⁹The ability of using this “inverse temperature” t to control the sparsity is enabled by Prop. 2, item 1. We may regard this as an hyperparameter analogous to the thresholds δ and p_0 used in the LOGISTIC and SOFTMAX baselines, respectively.

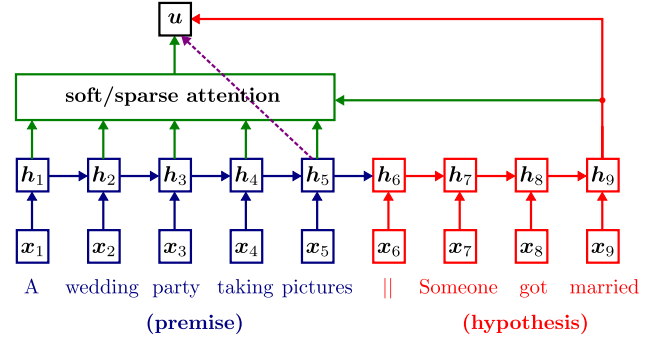


Figure 4. Network diagram for the NL inference problem. The premise and hypothesis are both fed into (gated) RNNs. The NOATTENTION system replaces the attention part (in green) by a direct connection from the last premise state to the output (dashed violet line). The LOGISTICATTENTION, SOFTATTENTION and SPARSEATTENTION systems have respectively independent logististics, a softmax, and a sparsemax-activated attention mechanism. In this example, $L = 5$ and $N = 9$.

2014), not optimized during training, which we linearly project onto a D -dimensional subspace (Astudillo et al., 2015).¹⁰ We denote by $\mathbf{x}_1, \dots, \mathbf{x}_L$ and $\mathbf{x}_{L+1}, \dots, \mathbf{x}_N$, respectively, the projected premise and hypothesis word vectors. These sequences are then fed into two recurrent networks (one for each). Instead of long short-term memories, as Rocktäschel et al. (2015), we used gated recurrent units (GRUs, Cho et al. 2014), which behave similarly but have fewer parameters. Our premise GRU generates a state sequence $\mathbf{H}_{1:L} := [\mathbf{h}_1 \dots \mathbf{h}_L] \in \mathbb{R}^{D \times L}$ as follows:

$$\mathbf{z}_t = \sigma(\mathbf{W}^{xz} \mathbf{x}_t + \mathbf{W}^{hz} \mathbf{h}_{t-1} + \mathbf{b}^z) \quad (28)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}^{xr} \mathbf{x}_t + \mathbf{W}^{hr} \mathbf{h}_{t-1} + \mathbf{b}^r) \quad (29)$$

$$\bar{\mathbf{h}}_t = \tanh(\mathbf{W}^{xh} \mathbf{x}_t + \mathbf{W}^{hh} (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}^h) \quad (30)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \mathbf{h}_{t-1} + \mathbf{z}_t \bar{\mathbf{h}}_t, \quad (31)$$

with model parameters $\mathbf{W}^{\{xz, xr, xh, hz, hr, hh\}} \in \mathbb{R}^{D \times D}$ and $\mathbf{b}^{\{z, r, h\}} \in \mathbb{R}^D$. Likewise, our hypothesis GRU (with distinct parameters) generates a state sequence $[\mathbf{h}_{L+1}, \dots, \mathbf{h}_N]$, being initialized with the last state from the premise (\mathbf{h}_L). The NOATTENTION system then computes the final state \mathbf{u} based on the last states from the premise and the hypothesis as follows:

$$\mathbf{u} = \tanh(\mathbf{W}^{pu} \mathbf{h}_L + \mathbf{W}^{hu} \mathbf{h}_N + \mathbf{b}^u) \quad (32)$$

where $\mathbf{W}^{pu}, \mathbf{W}^{hu} \in \mathbb{R}^{D \times D}$ and $\mathbf{b}^u \in \mathbb{R}^D$. Finally, it predicts a label \hat{y} from \mathbf{u} with a standard softmax layer. The SOFTATTENTION system, instead of using the last premise state \mathbf{h}_L , computes a weighted average of premise words

¹⁰We used GloVe-840B embeddings trained on Common Crawl (<http://nlp.stanford.edu/projects/glove/>).

Table 3. Accuracies for the natural language inference task. Shown are our implementations of a system without attention, and with logistic, soft, and sparse attentions.

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

Table 4. Examples of sparse attention for the natural language inference task. Nonzero attention coefficients are marked in bold. Our system classified all four examples correctly. The examples were picked from Rocktäschel et al. (2015).

A boy rides on a camel in a crowded area while talking on his cellphone. Hypothesis: <i>A boy is riding an animal.</i> [entailment]
A young girl wearing a pink coat plays with a yellow toy golf club. Hypothesis: <i>A girl is wearing a blue jacket.</i> [contradiction]
Two black dogs are frolicking around the grass together . Hypothesis: <i>Two dogs swim in the lake.</i> [contradiction]
A man wearing a yellow striped shirt laughs while seated next to another man who is wearing a light blue shirt and clasp ing his hands together. Hypothesis: <i>Two mimes sit in complete silence.</i> [contradiction]

with an attention mechanism, replacing Eq. 32 by

$$z_t = \mathbf{v}^\top \tanh(\mathbf{W}^{pm} \mathbf{h}_t + \mathbf{W}^{hm} \mathbf{h}_N + \mathbf{b}^m) \quad (33)$$

$$\mathbf{p} = \text{softmax}(\mathbf{z}), \text{ where } \mathbf{z} := (z_1, \dots, z_L) \quad (34)$$

$$\mathbf{r} = \mathbf{H}_{1:L} \mathbf{p} \quad (35)$$

$$\mathbf{u} = \tanh(\mathbf{W}^{pu} \mathbf{r} + \mathbf{W}^{hu} \mathbf{h}_N + \mathbf{b}^u), \quad (36)$$

where $\mathbf{W}^{pm}, \mathbf{W}^{hm} \in \mathbb{R}^{D \times D}$ and $\mathbf{b}^m, \mathbf{v} \in \mathbb{R}^D$. The LOGISTICATTENTION system, instead of Eq. 34, computes $\mathbf{p} = (\sigma(z_1), \dots, \sigma(z_L))$. Finally, the SPARSEATTENTION system replaces Eq. 34 by $\mathbf{p} = \text{sparsemax}(\mathbf{z})$.

We optimized all the systems with Adam (Kingma & Ba, 2014), using the default parameters $\beta_1 = 0.9, \beta_2 = 0.999$, and $\epsilon = 10^{-8}$, and setting the learning rate to 3×10^{-4} . We tuned a ℓ_2 -regularization coefficient in $\{0, 10^{-4}, 3 \times 10^{-4}, 10^{-3}\}$ and, as Rocktäschel et al. (2015), a dropout probability of 0.1 in the inputs and outputs of the network.

The results are shown in Table 3. We observe that the soft and sparse-activated attention systems perform similarly, the latter being slightly more accurate on the test set, and that both outperform the NOATTENTION and LOGIS-

TICATTENTION systems.¹¹

Table 4 shows examples of sentence pairs, highlighting the premise words selected by the SPARSEATTENTION mechanism. We can see that, for all examples, only a small number of words are selected, which are key to making the final decision.¹² Compared to a softmax-activated mechanism, which provides a dense distribution over all the words, the sparsemax activation yields a compact and more interpretable selection, which can be particularly useful in long sentences such as the one in the bottom row.

5. Conclusions

We introduced the sparsemax transformation, which has similar properties to the traditional softmax, but is able to output sparse probability distributions. We derived a closed-form expression for its Jacobian, needed for the backpropagation algorithm, and we proposed a novel “sparsemax loss” function, a sparse analogue of the logistic loss, which is smooth and convex. Empirical results in multi-label classification and in attention networks for natural language inference attest the validity of our approach.

The connection between sparse modeling and interpretability is key in signal processing (Hastie et al., 2015). Our approach is distinctive: it is not the model that is assumed sparse, but the label posteriors that the model parametrizes. Sparsity is also a desirable (and biologically plausible) property in neural networks, present in rectified units (Glorot et al., 2011) and maxout nets (Goodfellow et al., 2013).

There are several avenues for future research. The ability of sparsemax-activated attention to select only a few variables to attend makes it potentially relevant to neural architectures with random access memory (Graves et al., 2014; Grefenstette et al., 2015; Sukhbaatar et al., 2015), since it offers a compromise between soft and hard operations, maintaining differentiability. Sparsemax is also appealing for hierarchical attention: if we define a top-down product of distributions along the hierarchy, the sparse distributions produced by sparsemax will automatically prune the hierarchy, leading to computational savings. A possible disadvantage of sparsemax over softmax is that it seems less GPU-friendly, since it requires sort operations or median-finding algorithms. There is, however, recent work providing efficient implementations of these algorithms on GPUs (Alabi et al., 2012).

¹¹Rocktäschel et al. (2015) report scores slightly above ours: they reached a test accuracy of 82.3% for their implementation of SOFTATTENTION, and 83.5% with their best system, a more elaborate word-by-word attention model. Differences in the former case may be due to distinct word vectors and the use of LSTMs instead of GRUs.

¹²On the development set, we observed that only 24.6% of the premise words were selected.

Acknowledgements

We would like to thank the three anonymous reviewers, and also Tim Rocktäschel for answering various implementation questions, and Mário Figueiredo and Chris Dyer for helpful comments. This work was partially supported by Fundação para a Ciência e Tecnologia (FCT), through contracts UID/EEA/50008/2013 and UID/CEC/50021/2013, the LearnBig project (PTDC/EEI-SII/7092/2014), and the GoLocal project (grant CMUPERI/TIC/0046/2014).

References

- Agrawala, Ashok K. Learning with a Probabilistic Teacher. *IEEE Transactions on Information Theory*, 16(4):373–379, 1970.
- Alabi, Tolu, Blanchard, Jeffrey D, Gordon, Bradley, and Steinbach, Russel. Fast k-Selection Algorithms for Graphics Processing Units. *Journal of Experimental Algorithmics (JEA)*, 17:4–2, 2012.
- Albert, James H and Chib, Siddhartha. Bayesian Analysis of Binary and Polychotomous Response Data. *Journal of the American statistical Association*, 88(422):669–679, 1993.
- Astudillo, Ramon F, Amir, Silvio, Lin, Wang, Silva, Mário, and Trancoso, Isabel. Learning Word Representations from Scarce and Noisy Data with Embedding Subspaces. In *Proc. of the Association for Computational Linguistics (ACL), Beijing, China, 2015*.
- Bahdanau, Dzmitry, Cho, Kyunghyun, and Bengio, Yoshua. Neural Machine Translation by Jointly Learning to Align and Translate. In *International Conference on Learning Representations*, 2015.
- Blum, Manuel, Floyd, Robert W, Pratt, Vaughan, Rivest, Ronald L, and Tarjan, Robert E. Time Bounds for Selection. *Journal of Computer and System Sciences*, 7(4): 448–461, 1973.
- Bouchard, Guillaume. Efficient Bounds for the Softmax Function and Applications to Approximate Inference in Hybrid Models. In *NIPS Workshop for Approximate Bayesian Inference in Continuous/Hybrid Systems*, 2007.
- Bowman, Samuel R, Angeli, Gabor, Potts, Christopher, and Manning, Christopher D. A Large Annotated Corpus for Learning Natural Language Inference. In *Proc. of Empirical Methods in Natural Language Processing*, 2015.
- Bradley, Ralph Allan and Terry, Milton E. Rank Analysis of Incomplete Block Designs: The Method of Paired Comparisons. *Biometrika*, 39(3-4):324–345, 1952.
- Bridle, John S. Probabilistic Interpretation of Feedforward Classification Network Outputs, with Relationships to Statistical Pattern Recognition. In *Neurocomputing*, pp. 227–236. Springer, 1990.
- Chapelle, Olivier, Schölkopf, Bernhard, and Zien, Alexander. *Semi-Supervised Learning*. MIT Press Cambridge, 2006.
- Cho, Kyunghyun, Van Merriënboer, Bart, Gulcehre, Caglar, Bahdanau, Dzmitry, Bougares, Fethi, Schwenk, Holger, and Bengio, Yoshua. Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation. In *Proc. of Empirical Methods in Natural Language Processing*, 2014.
- Chorowski, Jan K, Bahdanau, Dzmitry, Serdyuk, Dmitriy, Cho, Kyunghyun, and Bengio, Yoshua. Attention-based Models for Speech Recognition. In *Advances in Neural Information Processing Systems*, pp. 577–585, 2015.
- Clarke, Frank H. *Optimization and Nonsmooth Analysis*. New York, Wiley, 1983.
- Condat, Laurent. Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, pp. 1–11, 2014.
- de Brébisson, Alexandre and Vincent, Pascal. An Exploration of Softmax Alternatives Belonging to the Spherical Loss Family. *arXiv preprint arXiv:1511.05042*, 2015.
- Duchi, J., Shalev-Shwartz, S., Singer, Y., and Chandra, T. Efficient Projections onto the L1-Ball for Learning in High Dimensions. In *Proc. of International Conference of Machine Learning*, 2008.
- Gao, Wei and Zhou, Zhi-Hua. On the Consistency of Multi-Label Learning. *Artificial Intelligence*, 199:22–44, 2013.
- Glorot, Xavier, Bordes, Antoine, and Bengio, Yoshua. Deep Sparse Rectifier Neural Networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 315–323, 2011.
- Goodfellow, Ian, Bengio, Yoshua, and Courville, Aaron. Deep Learning. Book in preparation for MIT Press, 2016. URL <http://goodfeli.github.io/dlbook/>.
- Goodfellow, Ian J, Warde-Farley, David, Mirza, Mehdi, Courville, Aaron, and Bengio, Yoshua. Maxout Networks. In *Proc. of International Conference on Machine Learning*, 2013.
- Graves, Alex, Wayne, Greg, and Danihelka, Ivo. Neural Turing Machines. *arXiv preprint arXiv:1410.5401*, 2014.

- Grefenstette, Edward, Hermann, Karl Moritz, Suleyman, Mustafa, and Blunsom, Phil. Learning to Transduce with Unbounded Memory. In *Advances in Neural Information Processing Systems*, pp. 1819–1827, 2015.
- Hastie, Trevor, Tibshirani, Robert, and Wainwright, Martin. *Statistical Learning with Sparsity: the Lasso and Generalizations*. CRC Press, 2015.
- Hermann, Karl Moritz, Kocisky, Tomas, Grefenstette, Edward, Espeholt, Lasse, Kay, Will, Suleyman, Mustafa, and Blunsom, Phil. Teaching Machines to Read and Comprehend. In *Advances in Neural Information Processing Systems*, pp. 1684–1692, 2015.
- Huber, Peter J. Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35(1):73–101, 1964.
- Kingma, Diederik and Ba, Jimmy. Adam: A Method for Stochastic Optimization. In *Proc. of International Conference on Learning Representations*, 2014.
- Koyejo, Sanmi, Natarajan, Nagarajan, Ravikumar, Pradeep K, and Dhillon, Inderjit S. Consistent Multilabel Classification. In *Advances in Neural Information Processing Systems*, pp. 3303–3311, 2015.
- Lewis, David D, Yang, Yiming, Rose, Tony G, and Li, Fan. RCV1: A New Benchmark Collection for Text Categorization Research. *The Journal of Machine Learning Research*, 5:361–397, 2004.
- Liu, Dong C and Nocedal, Jorge. On the Limited Memory BFGS Method for Large Scale Optimization. *Mathematical programming*, 45(1-3):503–528, 1989.
- McCullagh, Peter and Nelder, John A. *Generalized Linear Models*, volume 37. CRC press, 1989.
- Michelot, C. A Finite Algorithm for Finding the Projection of a Point onto the Canonical Simplex of \mathbb{R}^n . *Journal of Optimization Theory and Applications*, 50(1):195–200, 1986.
- Nesterov, Y. A Method of Solving a Convex Programming Problem with Convergence Rate $O(1/k^2)$. *Soviet Math. Doklady*, 27:372–376, 1983.
- Ollivier, Yann. Riemannian Metrics for Neural Networks. *arXiv preprint arXiv:1303.0818*, 2013.
- Pardalos, Panos M. and Kuvshinov, Naina. An Algorithm for a Singly Constrained Class of Quadratic Programs Subject to Upper and Lower Bounds. *Mathematical Programming*, 46(1):321–328, 1990.
- Pennington, Jeffrey, Socher, Richard, and Manning, Christopher D. Glove: Global Vectors for Word Representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12:1532–1543, 2014.
- Penot, Jean-Paul. Conciliating Generalized Derivatives. In Demyanov, Vladimir F., Pardalos, Panos M., and Batsyn, Mikhail (eds.), *Constructive Nonsmooth Analysis and Related Topics*, pp. 217–230. Springer, 2014.
- Rocktäschel, Tim, Grefenstette, Edward, Hermann, Karl Moritz, Kočiský, Tomáš, and Blunsom, Phil. Reasoning about Entailment with Neural Attention. *arXiv preprint arXiv:1509.06664*, 2015.
- Sukhbaatar, Sainbayar, Szlam, Arthur, Weston, Jason, and Fergus, Rob. End-to-End Memory Networks. In *Advances in Neural Information Processing Systems*, pp. 2431–2439, 2015.
- Sutton, Richard S and Barto, Andrew G. *Reinforcement Learning: An Introduction*, volume 1. MIT press Cambridge, 1998.
- Xu, Kelvin, Ba, Jimmy, Kiros, Ryan, Courville, Aaron, Salakhutdinov, Ruslan, Zemel, Richard, and Bengio, Yoshua. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *International Conference on Machine Learning*, 2015.
- Zadrozny, Bianca. Reducing Multiclass to Binary by Coupling Probability Estimates. In *Advances in Neural Information Processing Systems*, pp. 1041–1048, 2001.
- Zhang, Min-Ling and Zhou, Zhi-Hua. A Review on Multi-Label Learning Algorithms. *Knowledge and Data Engineering, IEEE Transactions on*, 26(8):1819–1837, 2014.
- Zhang, Tong. Statistical Behavior and Consistency of Classification Methods Based on Convex Risk Minimization. *Annals of Statistics*, pp. 56–85, 2004.
- Zou, Hui, Zhu, Ji, and Hastie, Trevor. The Margin Vector, Admissible Loss and Multi-class Margin-Based Classifiers. Technical report, Stanford University, 2006.