

How Can We Be So Dense? The Benefits of Using Highly Sparse Representations

Subutai Ahmad¹ Luiz Scheinkman¹
Numenta, Redwood City, California, USA

Abstract

Most artificial networks today rely on dense representations, whereas biological networks rely on sparse representations. In this paper we show how sparse representations can be more robust to noise and interference, as long as the underlying dimensionality is sufficiently high. A key intuition that we develop is that the ratio of the operable volume around a sparse vector divided by the volume of the representational space decreases exponentially with dimensionality. We then analyze computationally efficient sparse networks containing both sparse weights and activations. Simulations on MNIST and the Google Speech Command Dataset show that such networks demonstrate significantly improved robustness and stability compared to dense networks, while maintaining competitive accuracy. We discuss the potential benefits of sparsity on accuracy, noise robustness, hyperparameter tuning, learning speed, computational efficiency, and power requirements.

1. Introduction

The literature on sparse representations in neural networks dates back many decades, with neuroscience as one of the primary motivations. In 1988 Kanerva proposed the use of sparse distributed memories (Kanerva, 1988) to model the highly sparse representations seen in the brain. In 1997, (Olshausen & Field, 1997) showed that incorporating sparse priors and sparse cost functions in encoders can lead to receptive field representations that are remarkably close to what is observed in the primate visual cortex. More recently (Lee et al., 2008; Chen et al., 2018) showed hierarchical sparse representations that qualitatively lead to natural looking hierarchical feature detectors. (Lee et al., 2009; Nair & Hinton, 2009; Srivastava et al., 2013; Rawlinson et al.,

2018) showed that introducing sparsity terms can sometimes lead to improved test set accuracies.

Despite the above literature the majority of neural networks today rely on dense representations. One exception is the pervasive use of dropout (Srivastava et al., 2014) as a regularizer. Dropout randomly “kills” a percentage of the units (in practice usually 50%) on every training input presentation. Variational dropout techniques tune the dropout rates individually per weight (Molchanov et al., 2017). Dropout introduces random sparse representations during learning, and has been shown to be an effective regularizer in many contexts.

In this paper we discuss certain inherent benefits of high dimensional sparse representations. We focus on robustness and sensitivity to interference. These are central issues with today’s neural network systems where even small (Szegedy et al., 2013) and large (Rosenfeld et al., 2018) perturbations can cause dramatic changes to a network’s output. We offer two main contributions. First, we analyze high dimensional sparse representations, and show that such representations are naturally more robust to noise and interference from random inputs. When matching sparse patterns, corrupted versions of a pattern are “close” to the original whereas random patterns are exponentially hard to match.

Our second contribution is an efficient construction of sparse deep networks that is designed to exploit the above properties. We implement networks where the weights for each unit in a layer randomly sample from a sparse subset of the source layer below. In addition the output of each layer is constrained such that only the k most active units are allowed to be non-zero, where k is much smaller than the number of units in that layer. In these networks, the number of non-zero products for each layer is approximately (sparsity of layer i) \times (sparse weights of layer $i + 1$). This formulation results in simple differentiable sparse layers that can be dropped into both standard linear and convolutional layers.

We demonstrate significantly improved robustness to noise for MNIST and the Google Speech Commands dataset, while maintaining competitive accuracy in the standard zero

¹. Correspondence to: Subutai Ahmad, Luiz Scheinkman
<[sahmad, lscheinkman]@numenta.com>.

noise scenario. We discuss the number of weights used by sparse networks in these datasets, and the impact of additional pruning. Our work extends the existing literature on sparse networks and pruning (see Section 5 for a comparison with some prior work). At the end of the paper we discuss some possible areas for future work.

2. High Dimensional Sparse Representations

In this section we develop some basic properties of sparse representations as they relate to noise robustness and interference. In a typical neural network an input vector is matched against a stored weight vector using a dot product. This is then followed by a threshold-like non-linearity such as $\tanh(\cdot)$ or $\text{ReLU}(\cdot)$.

Ideally we would like the outputs of each layer to be invariant to noise or corrupted inputs. When comparing two sparse vectors via a dot product, the results are unaffected by the zero components of either vector. A key quantity we consider is the ratio of the matching volume around a prototype vector divided by the volume of the whole space. The larger the match volume around a vector, the more robust it is to noise. The smaller the ratio, the less likely it is that random inputs can affect the match.

2.1. Matching Sparse Binary Vectors

We quantify the above ratio using binary vectors (following our previous work in (Ahmad & Hawkins, 2016)). In this section we show that the ratio decreases exponentially with increased dimensionality, while maintaining a large match volume. Let \mathbf{x} be a binary vector of length n , and let $|\mathbf{x}|$ denote the number of non-zero entries. The dot product $\mathbf{x}_i \cdot \mathbf{x}_j$ counts the overlap, or number of shared bits, between two such vectors. We would like to understand the probability of two vectors having significant overlap, i.e. overlap greater than some threshold θ .

We define the overlap set, $\Omega^n(\mathbf{x}_i, b, k)$, as the set of all vectors of size k that have exactly b bits of overlap with \mathbf{x}_i . The number of such vectors can be calculated as:

$$|\Omega^n(\mathbf{x}_i, b, k)| = \binom{|\mathbf{x}_i|}{b} \binom{n - |\mathbf{x}_i|}{k - b} \quad (1)$$

The left half of the above product counts all the ways we can select exactly b bits out of active bits in $|\mathbf{x}_i|$. The right half counts the number of ways we can select the remaining $k - b$ bits from the components of \mathbf{x}_i that are zero. The product of these two quantities represents the number of all vectors with exactly b bits of overlap with $|\mathbf{x}_i|$. We can now count the number of vectors that match \mathbf{x}_i , i.e. where $\mathbf{x}_i \cdot \mathbf{x}_j \geq \theta$ as:

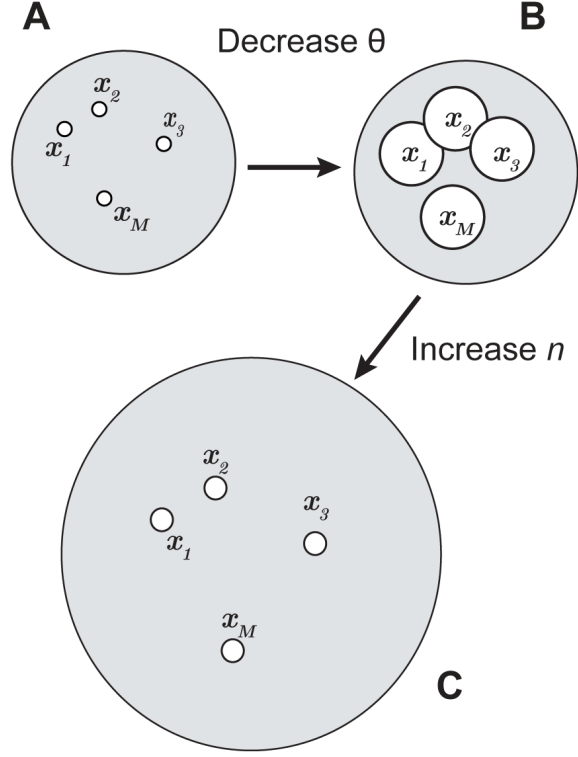


Figure 1. An illustration of the conceptual effect of decreasing the match threshold θ and increasing n , the dimensionality. The large grey circles denote the universe of possible patterns. The smaller circles each represent the set of matches around one vector. When θ is high (A), very few random vectors can match these vectors (small white circles). As you decrease θ , the set of potential matches increases (larger white circles in B). If you then increase n , the universe of possible patterns increases, and the relative sizes of the white circles shrink rapidly.

$$\sum_{b=\theta}^{|\mathbf{x}_i|} |\Omega^n(\mathbf{x}_i, b, |\mathbf{x}_j|)| \quad (2)$$

If we select vectors from a uniform random distribution, the probability of significant overlap can be calculated as:

$$P(\mathbf{x}_i \cdot \mathbf{x}_j \geq \theta) = \frac{\sum_{b=\theta}^{|\mathbf{x}_i|} |\Omega^n(\mathbf{x}_i, b, |\mathbf{x}_j|)|}{\binom{n}{|\mathbf{x}_j|}} \quad (3)$$

where $\binom{n}{|\mathbf{x}_j|}$ is the set of all possible comparison vectors.

2.2. Impact of Dimensionality and Sparsity

Two key factors in Eq. 3 are the number of non-zero components, $|\mathbf{x}_i|$, and the dimensionality, n . Figure 1 provides

an intuitive description of their impact. Assume we have M prototype vectors, and we want to match noisy versions of these vectors. Around each prototype there is a set of matching vectors. If the threshold is very high, the set of matching vectors is small (illustrated by the small circles in Figure 1A) and there will be quite a bit of space between these sets. As you decrease θ matching is less strict and you can match noisier versions of each prototype. The cost is that the chance of matching the other vectors also increases because there is less free space in between (Figure 1B). It turns out that for sparse vectors, this cost is offset as you increase n . That is, as n increases, the denominator in Eq. 3 (and the corresponding "free" space) increases much faster than the numerator. For a fixed sparsity level, you can maintain highly tolerant matches without the cost of additional false positives simply by increasing the dimensionality.

Fig 2 illustrates this trend for some example sparsities. In this figure we simulated matching with random vectors and plotted match rates with random vectors as a function of the number of active bits and the underlying dimensionality. In the simulation we repeatedly generated a random prototype vector with $|x_i| = 24$ bits on and then attempted to match against random test vectors with a bits on. We matched using a threshold θ of 12 which meant that even vectors that were up to 50% different from x_i would match. We varied a and the dimensionality of the vectors, n .

The chart shows that for sparse binary vectors, match rates with random vectors drop rapidly as the underlying dimensionality increases. The horizontal line indicates the probability of matching x_i against dense vectors, with $a = n/2$. The probability of dense matches stays relatively high and unaffected by dimensionality, indicating that both sparseness and high dimensionality are key to robust matches. In (Ahmad & Hawkins, 2016) we develop additional properties, including the probability of false negatives.

2.3. Matching Sparse Scalar Vectors

Deep networks operate on scalar vectors, and in this section we consider how the above ideas apply to sparse scalar representations. Binary and scalar vectors are similar in that the components containing zero do not affect the dot product, and thus the combinatorics in Eq. 3 are still applicable. Eq. 1 represents the set of scalar vectors where the number of non-zero multiplies in the dot product is exactly b , and Eq. 3 represents the probability that the number of non-zero multiplies is $\geq \theta$. However, an additional factor is the distribution of scalar values. If components in one vector are extremely large relative to θ , the likelihood of a significant match will be high even with a single shared non-zero component.

We wanted to see if the exponential drop in random matches for binary vectors, demonstrated by Figure 2, can be ob-

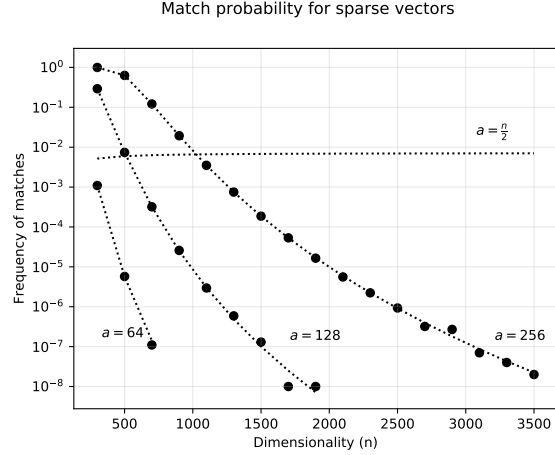


Figure 2. The probability of matches to random binary vectors (with a active bits) as a function of dimensionality, for various levels of sparsity. The probability decreases exponentially with n . Black circles denote the observed frequency of a match (based on a large number of trials). The dotted lines denote the theoretically predicted probabilities using Eq. 3.

tained using scalar vectors, and if so, the conditions under which they hold. Let x_w and x_i represent two sparse vectors such that $\|x_w\|_0$ and $\|x_i\|_0$ counts the number of non-zero entries in each. Let each non-zero component be independent and sampled from the distributions $P_{\theta_w}(x_w)$ and $P_{\theta_i}(x_i)$. The probability of a significant match is then:

$$P(x_w \cdot x_i \geq \theta) = \frac{\sum_{b=\theta}^{\|x_w\|_0} p_b |\Omega^n(x_w, b, \|x_i\|_0)|}{\binom{n}{\|x_i\|_0}} \quad (4)$$

where p_b is the probability that the dot product is $\geq \theta$ given that the overlap is exactly b components:

$$p_b = P(x_w \cdot x_i \geq \theta \mid \|x_w \cdot x_i\|_0 = b) \quad (5)$$

There does not appear to be a closed form way to compute p_b for normal or uniform distributions so we resort to simulations that mimic our network structure.

As before, we generated a large number of random vectors x_w and x_i , and plotted the frequency of random matches. With $\|x_w\|_0 = k$, we focus on simulations where the non-zero entries in x_w are uniform in $[-1/k, 1/k]$, and the non-zero entries in x_i are uniform in $S * [0, 2/k]$. We focus on this formulation because of the relationship to common network structures and weight initialization. x_w is a putative weight vector and x_i is an input vector to this layer from the

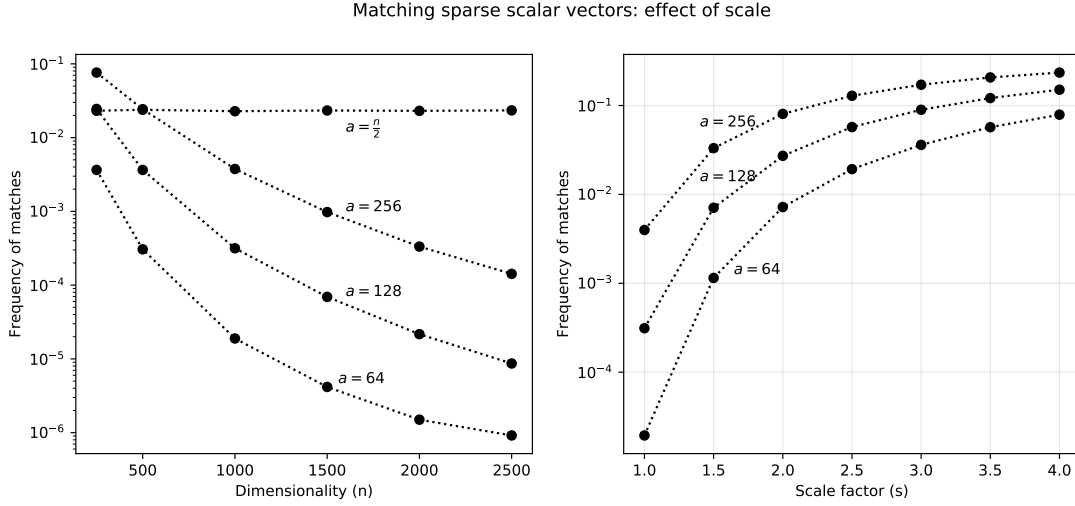


Figure 3. **Left:** The probability of matches to random scalar vectors (with a non-zero components) as a function of dimensionality, for various levels of sparsity. The probability of false matches decreases exponentially with n . Note that the probability for a dense vector, $a = \frac{n}{2}$ stays relatively high, and does not decrease with dimensionality. **Right:** The impact of scale on vector matches with a fixed $n = 1000$. The larger the scaling discrepancy, the higher the probability of a false match.

previous layer (we assume unit activations are positive, the result of a ReLU-like non-linearity). S controls the scale of \mathbf{x}_i relative to \mathbf{x}_w .

Figure 3 (left) shows the behavior with $k = 32$ and $S = 1$. We varied the activity of the input vectors $\|\mathbf{x}_i\|_0 = a$ and the dimensionality of the vectors, n . We set $\theta = E[\mathbf{x}_w \cdot \mathbf{x}_w]/2.0$. The chart demonstrates that under these conditions we can achieve robust behavior similar to that of binary vectors. Figure 3 (right) plots the effect of S on the match probabilities with a fixed $n = 1000$. As this chart shows, the error increases significantly as S increases. Taken together, these results show that the fundamental robustness properties of binary sparse vectors can also hold for sparse scalar vectors, as long as the overall scaling of vectors are in a similar range.

2.4. Non-uniform Distribution of Vectors

Eq. 3 assumes the ideal case where vectors are chosen with a uniform random distribution. With a non-uniform distribution the error rates will be higher. The more non-uniform the distribution the worse the error rates. For example, if you mostly end up observing 10 inputs, your error rates will be bounded at around 10%. Thus, to optimize error rates, it is important to be as close to a uniform distribution as possible.

3. Sparse Network Description

Here we discuss a particular sparse network implementation that is designed to exploit Eq. 3. This implementation is an extension of our previous work on the HTM Spatial Pooler, a binary sparse coding algorithm that models sparse code generation in the neocortex (Hawkins et al., 2011; Cui et al., 2017). Specifically, we formulate a version of the Spatial Pooler that is designed to be a drop-in layer for neural networks trained with back-propagation. Our work is also closely related to previous literature on k-winner take all networks (Majani et al., 1989) and fixed sparsity networks (Makhzani & Frey, 2015).

Consider a network with L hidden layers. Let \mathbf{y}^l denote the vector of outputs from layer l , respectively, with \mathbf{y}^0 as the input vector. \mathbf{W}^l and \mathbf{u}^l are the weights and biases for each layer. In a standard neural network the weights \mathbf{W}^l are typically dense and initialized using a uniform random distribution. The feed forward outputs are then calculated as follows:

$$\hat{\mathbf{y}}^l = \mathbf{W}^l \cdot \mathbf{y}^{l-1} + \mathbf{u}^l$$

$$\mathbf{y}^l = f(\hat{\mathbf{y}}^l)$$

where f is any activation function, such as $\tanh(\cdot)$ or $\text{ReLU}(\cdot)$. (Figure 4 left.)

To implement our sparse networks, we make two modifications to this basic formulation (Figure 4 right.). First, we initialize the weights using a sparse random distribution, such

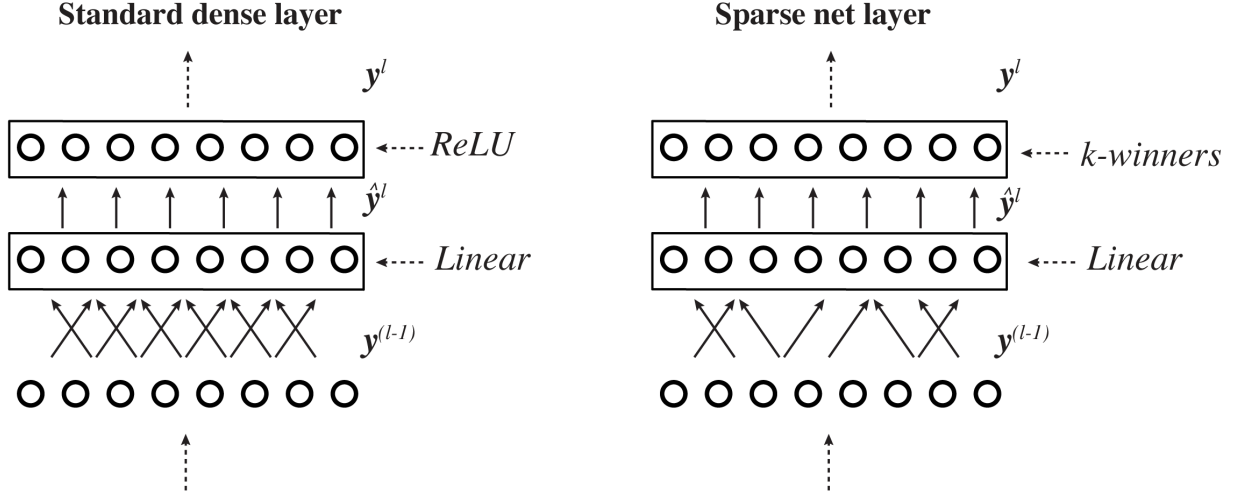


Figure 4. This figure illustrates the differences between a generic dense network layer (**left**) and a sparse network layer (**right**). In the sparse layer, the linear layer subsamples from its input layer (implemented via sparse weights, depicted with fewer arrows). In addition, the ReLU layer is replaced by a k-winners layer.

that only a fraction of the weights contain non-zero values. Non-zero weights are initialized using standard Kaiming initialization (He et al., 2015b). The rest of the connections are treated as non-existent, i.e. the corresponding weights are zero throughout the life of the network. Second, only the top- k active units within each layer are maintained in y^l , and the rest set to zero. This k -winners step is non-linear and can be thought of as a substitute for the ReLU function. Instead of a threshold of 0, the threshold here is adaptive and corresponds to the k 'th largest activation (Makhzani & Frey, 2013).

The layer can be trained using standard gradient descent. Similar to ReLU, the gradient of the layer is calculated as 1 above the threshold and 0 elsewhere. During inference we increase k by 50%, which led to slightly better accuracies. In all our simulations the last layer of each network is a standard linear output layer with log-softmax activation function.

3.1. Boosting

One practical issue with the above formulation is that it is possible for a small number of units to initially dominate and then, through learning, become active for a large percentage of patterns (this was also noted in (Makhzani & Frey, 2015; Cui et al., 2017)). Having a small number of active units negatively impacts the available representational volume. It is desirable for every unit to be equally active in order to maximize the robustness of the representation in Eq. 3.

To address this we employ a boosting term (Hawkins et al., 2011; Cui et al., 2017) which favors units that have not been active recently. We compute a running average of each unit's duty cycle (i.e. how frequently it has been one of the top k units):

$$d_i^l(t) = (1 - \alpha)d_i^l(t-1) + \alpha \cdot [i \in \text{topIndices}^l] \quad (6)$$

A boost coefficient b_i^l is then calculated for each unit based on the target duty cycle and the current average duty cycle:

$$b_i^l(t) = e^{\beta(\hat{a}^l - d_i^l(t))} \quad (7)$$

The target duty cycle \hat{a}^l is a constant reflecting the percentage of units that are expected to be active, i.e. $\hat{a}^l = \frac{k}{|y^l|}$. The boost factor, β , is a positive parameter that controls the strength of boosting. $\beta = 0$ implies no boosting ($b_i^l = 1$), and higher numbers lead to larger boost coefficients. In (Hawkins et al., 2011; Cui et al., 2017) we showed that Eq. 7 encourages each unit to have equal activation frequency and effectively maximizes the entropy of the layer.

The boost coefficients are used during the k-winners step to select which units remain active for this input. Through boosting, units which have not been active recently have a disproportionately higher impact and are more likely to win, whereas overly active units are de-emphasized. To determine the output of the layer, the non-boosted activity

Algorithm 1 k -winners layer

-
- 1: $\hat{\mathbf{y}}^l = \mathbf{w}^l \cdot \mathbf{y}^{(l-1)} + \mathbf{u}^l$
 - 2: $b_i^l(t) = e^{\beta(\hat{a}^l - d_i^l(t))}$
 - 3: $\text{topIndices}^l = \text{topk}(\mathbf{b}^l \odot \hat{\mathbf{y}}^l)$
 - 4: $\mathbf{y}^l = 0$
 - 5: $\mathbf{y}^l[\text{topIndices}^l] = \hat{\mathbf{y}}^l$
 - 6: $d_i^l(t) = (1 - \alpha)d_i^l(t-1) + \alpha \cdot [y_i^l(t) \in \text{topIndices}^l]$
-

of each winning unit is kept and the remaining units are set to zero. The duty cycle is then updated. The complete pseudo-code description for the k -winners layer is described in Algorithm 1. In our simulations we used $\beta = 1.0$ or 1.5 for all sparse simulations.

3.2. Sparse Convolutional Layers

We can apply the above algorithm to convolutional networks (CNNs) (LeCun et al., 1989). A canonical CNN layer uses a linear convolutional layer containing a number of filters, followed by a max-pooling (downsampling) layer, followed by ReLU. In order to implement sparse CNN layers, the k -winners layer is applied to the output of the max-pooling layer instead of ReLU (just as in our non-convolutional layers). However, since each filter in a CNN shares weights across the image, duty cycles are accumulated per filter. In our simulations dense and sparse CNN nets both have a hidden layer (which is dense or sparse, respectively) after the last convolutional layer, followed by a linear plus softmax layer. We used 5×5 filters throughout with a stride of 1. In our tests, the weight sparsity of CNN layers did not impact the results. We suspect this is due to the small size of each kernel and did not use sparse weights for the CNN filters in our experiments.

4. Results

4.1. MNIST

We first trained our networks on MNIST (LeCun et al., 1998). We trained both dense and sparse implementations. Each network consisted of one or two convolutional layers, followed by a hidden layer, followed by a linear + softmax output layer. Sparse nets consisted of sparse convolutional layers followed by a sparse hidden layer.

Networks were trained using standard stochastic gradient descent to minimize cross entropy loss. We used starting learning rates in the range $0.01 - 0.04$, and the learning rate was decreased by a factor between 0.5 and 0.9 after each epoch. We also tried batch normalization (Ioffe & Szegedy, 2015) and found it did not help for MNIST (it did help significantly for Google Speech Commands results - see below). For sparse networks, we used a small mini-batch size (around 4), for the first epoch only, in order

NETWORK	TEST SCORE	NOISE SCORE
DENSE CNN-1	99.14 ± 0.03	$74,569 \pm 3,200$
DENSE CNN-2	99.31 ± 0.06	$97,040 \pm 2,853$
SPARSE CNN-1	98.41 ± 0.08	$100,306 \pm 1,735$
SPARSE CNN-2	99.09 ± 0.05	$103,764 \pm 1,125$
DENSE CNN-2 SP3	99.13 ± 0.07	$100,318 \pm 2,762$
SPARSE CNN-2 D3	98.89 ± 0.13	$102,328 \pm 1,720$
SPARSE CNN-2 W1	98.2 ± 0.19	$100,322 \pm 2,082$
SPARSE CNN-2 DSW	98.92 ± 0.09	$70,566 \pm 2,857$

Table 1. MNIST results for dense and sparse architectures. We show classification accuracies and total noise scores (the total number of correct classification for all noise levels). Results are averaged over 10 random seeds, \pm one standard deviation. CNN-1 and CNN-2 indicate one or two convolutional layers, respectively.

to let duty cycle calculations update frequently and settle. Hyperparameters such as the learning rate and network size were chosen using a validation set consisting of 10,000 randomly chosen training samples. We then report final results on the test set using networks trained on the full training set.

Results Without Noise: State of the art accuracies on MNIST using convolutional neural networks (without distortions or other training augmentation) are in the range $98.3 - 99\%$ respectively¹. Table 1 (left column) lists the classification accuracies for the networks in our experiments. Our accuracies are in the same range, for both sparse and dense networks. Table 3 lists the key parameters for each of the listed networks (see also the next section for a more in-depth discussion).

Results With Noise: In order to test noise robustness we generated MNIST images with varying levels of additive noise. For each test image we randomly set $\eta\%$ of the pixels to a constant value near white (the constant value was two standard deviations over the mean pixel intensity). Figure 5 (A) shows sample images for different noise levels. We generated 11 different noise levels with η ranging between 0 and 0.5 in increments of 0.05. We also computed an overall **noise score** which counted the total number of correct classifications across all noise levels.

The right column of Table 1 shows the noise scores for each of the architectures. Networks in the top section of the table (Dense CNN-1 and Dense CNN-2) are composed of standard dense convolutional and hidden layers. Networks in the middle section (Sparse CNN-1 and Sparse CNN-2) are composed of sparse convolutional and sparse hidden layers. Networks in the last section contain a mixture of dense and sparse layers. Overall the architectures with

¹Source: <http://yann.lecun.com/exdb/mnist>

sparse layers performed significantly better on the noise score than the fully dense networks. Sparse CNN-2, the two layer completely sparse network, had the best noise score. The two fully dense networks performed substantially worse than the others on noise, even though their test accuracies were comparable. Figure 5 plots the accuracy of fully dense and sparse networks at different noise levels. Note that raw test score was not a predictor of noise robustness, suggesting that focusing on pure test set accuracy alone is not sufficient for gauging performance under adverse conditions.

Ablation studies: In order to judge the relative contributions of sparse layers we ran experiments where we replaced various sparse components with their dense counterparts, i.e. dense CNNs with sparse hidden layers, and vice versa. Dense CNN-2 SP3 contained two dense CNN layers followed by the sparse third layer from Sparse CNN-2. Sparse CNN-2 D3 contained the same CNN layers as Sparse CNN-2 followed by the dense third layer from Dense CNN-2. Sparse CNN-2 W1 was identical to Sparse CNN-2 except that the weight sparsity was 1 (i.e. fully dense weights). Sparse CNN-2 DSW contained a third layer with dense outputs, but with a weight sparsity of 0.3%.

The results of these networks are shown in the bottom third of Table 1. From a noise robustness perspective, most of the variants (except for Sparse CNN-2 DSW) performed well, better than the best pure dense network. This supports the idea that sparsity in many forms may be helpful with robustness. It is interesting to note that the standard deviation of the noise score in these variants was also higher than that of the pure sparse networks. Overall the results with mixed networks were encouraging, and suggest a clear benefit to introducing sparsity at any level.

Impact of Dropout: The above results did not use dropout (Srivastava et al., 2014), which is generally thought to improve robustness. We found that dropout did occasionally improve the robustness of dense networks, but any improvements were modest and the dropout percentage had to be tuned carefully. For sparse nets dropout consistently reduced accuracies. Even with the optimal dropout percentage, the noise scores of dense networks were significantly lower than sparse nets.

4.2. Google Speech Commands Dataset

In order to test sparse nets on a different domain, we applied them to the Google Speech Commands dataset (GSC). This audio dataset was made publicly available in 2017 (Warden, 2017) and consists of 65,000 one-second long utterances of 30 keywords spoken by thousands of individuals. The dataset contains predefined training, validation, and test sets.

Reference convolutional nets using ten of the keyword categories (plus artificial "silence" and "unknown" categories

NETWORK	TEST SCORE	NOISE SCORE
DENSE CNN-2 (DR=0.0)	96.37 \pm 0.37	8,730 \pm 471
DENSE CNN-2 (DR=0.5)	95.69 \pm 0.48	7,681 \pm 368
SPARSE CNN-2	96.65 \pm 0.21	11,233 \pm 1013
SUPER-SPARSE CNN-2	96.57 \pm 0.16	10,752 \pm 942

Table 2. Classification on Google Speech Commands for a number of architectures. We show test and noise scores, averaged over 10 random seeds, \pm one standard deviation. Dr corresponds to different dropout levels.

created during training augmentation) achieve accuracies in the range 91 – 92% (Sainath & Parada, 2015; Tang & Lin, 2017). In (Tang & Lin, 2017) they demonstrated improved accuracies in the range of 95 – 96% using residual networks (ResNets (He et al., 2015b;a)).

A Kaggle competition using GSC (also limited to 10 categories) took place between November 2017 and early 2018². For our simulations we use the preprocessing code provided by one of the top-10 contestants (Tuguldur, 2018) who achieved around 97 – 97.5% accuracies using variants of ResNet and VGG (Simonyan & Zisserman, 2014) architectures. Following this implementation, audio samples in our simulations are converted to 32-band Mel spectrograms before being fed to the network. During training we augment the data by randomly adjusting the amplitude, speed, and pitch of each training sample, and by randomly shifting and stretching samples in the frequency domain. No data augmentation is performed on the validation or test sets.

We trained dense and sparse convolutional networks, with hyperparameters chosen based on the validation set. We were able to achieve reasonable accuracies using two convolutional layers, followed by a hidden layer and then a linear + softmax output layer. Our sparse networks had sparse convolutional layers as well as a sparse hidden layer. Unlike MNIST we found that batch normalization (Ioffe & Szegedy, 2015) accelerated learning significantly, and we used it for every layer.

Using the above setup we were able to achieve test set accuracies in the range of 96.5 – 97.2% classifying the ten categories corresponding to the digits "zero" through "nine". Table 2 (left column) shows mean accuracy on the test set. Both dense and sparse networks had about the same accuracy. Dropout had a negative effect on the accuracy. Table 3 lists the key parameters in each network.

Results With Noise: As with MNIST, we again created noisy versions of the test set. For each test audio sample *A* we generated a random white noise sample and blended

²<https://www.kaggle.com/c/tensorflow-speech-recognition-challenge>

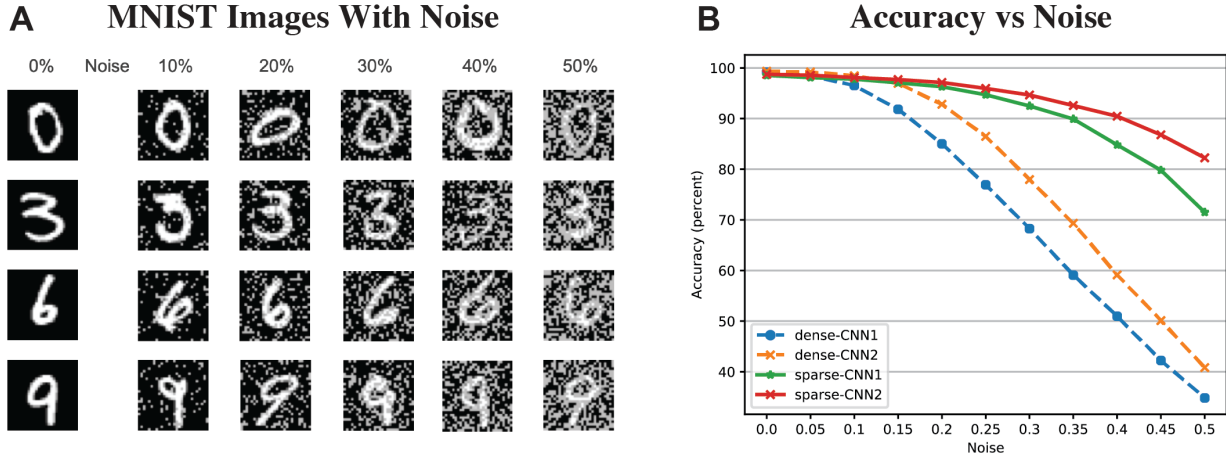


Figure 5. A. Example MNIST images with varying levels of noise. B. Classification accuracy as a function of noise level.

them together:

$$A^* = (1 - \eta)A + \eta \text{whiteNoise}$$

We generated 11 different noise levels, with η ranging from 0 to 0.5 in increments of 0.05. Our overall noise score counted the total number of classifications across all noise levels.

As can be seen in Table 2 sparse networks performed significantly better than the best dense network. We included a "Super-Sparse CNN-2" with a significantly sparser hidden layer. The hidden layer for this network had 10% weight sparsity, and a lower output sparsity (Table 3). This network had slightly lower noise score, but its score was still significantly higher than that of the dense networks. Overall these results demonstrate that the robustness of sparse networks seen with MNIST can scale to other domains.

4.3. Computational Considerations

In standard networks, the size of each weight matrix is $|\mathbf{W}^l| = |\mathbf{y}^{l-1}| |\mathbf{y}^l|$ and the order of complexity of the feed-forward operation can be approximated by the number of multiplications, $|\mathbf{y}^{l-1}|^2 |\mathbf{y}^l|$. The computational efficiency of sparse systems is closely related to the fraction of non-zeros. In our sparse hidden layers, both activations and weight values are sparse and the number of non-zero product terms in the forward computation is proportional to $k^{l-1} w^l |\mathbf{y}^{l-1}| |\mathbf{y}^l|$, where $0 < w^l \leq 1$ is the fraction of non-zero weights. In our convolutional layers, only activations values are sparse and the number of non-zero product terms in the forward computation is proportional to $k^{l-1} * K^l * K^l * |\mathbf{y}^l|$, where K^l is the kernel width of

each filter.

As an example, the number of non-zero multiplies between the first two convolutional layers in the GSC Sparse CNN-2 network is $12,544 * 1600 * 6400 = 1.23 \times 10^{10}$, about 10.5X smaller than the corresponding dense network. The number of non-zero multiplies between the second convolutional layer and the hidden layer in the same network is $200 * 640,000 * 1000 = 1.28 \times 10^{11}$, about 20X smaller than the dense network. For Super Sparse CNN-2, that ratio is 35X as compared to the dense version.

As can be seen, the number of non-zeros products is significantly smaller in the sparse net implementations. Unfortunately we found that current versions of deep learning frameworks, including PyTorch and Tensorflow do not have adequate support for sparse matrices to exploit these properties, and our implementations ran at the same speed as the corresponding dense networks. We suspect this is due to the fact that highly sparse networks are not sufficiently popular in practice. We hope that studies such as this one will encourage highly optimized sparse implementations. (Note that such optimizations may be non-trivial as the set of k -winners changes on every step.) When this becomes feasible our numbers suggest there is a strong possibility for large performance gains and/or improvements in power usage. It is also worth noting that this reduction in computational complexity does not come at a cost. Rather, our experiments showed that sparse representations can lead to improved accuracies under noisy conditions.

5. Discussion

In this paper we illustrated benefits of sparse representations. We developed intuitions and theory for the structure of vec-

NETWORK	L1 F	L1 SPARSITY	L2 F	L2 SPARSITY	L3 N	L3 SPARSITY	WT SPARSITY
MNIST							
DENSE CNN-1	30	100%			1000	100%	100%
DENSE CNN-2	30	100%	30	100%	1000	100%	100%
SPARSE CNN-1	30	9.3%			150	33.3%	30%
SPARSE CNN-2	32	8.7%	64	29.3 %	700	14.3%	30%
DENSE CNN-2 SP3	30	100%	30	100%	700	14.3%	30%
SPARSE CNN-2 D3	32	8.7%	64	29.3 %	1000	100%	100%
SPARSE CNN-2 W1	32	8.7%	64	29.3 %	700	14.3%	100%
SPARSE CNN-2 DSW	32	8.7%	64	29.3 %	1000	100%	30%
GSC							
DENSE CNN-2	64	100%	64	100%	1000	100%	100%
SPARSE CNN-2	64	9.5%	64	12.5%	1000	10%	40%
SUPER SPARSE CNN-2	64	9.5%	64	12.5%	1500	6.7%	10%

Table 3. Key parameters for each network. L1F and L2F denote the number of filters at the corresponding CNN layer. L1,2,3 sparsity indicates k/n , the percentage of outputs that were enforced to be non-zero. 100% indicates a special case where we defaulted to traditional ReLU activations. Wt sparsity indicates the percentage of weights that were non-zero. All parameters are available in the source code.

tor matching in the context of binary sparse representations. We then constructed efficient neural network formulations of sparse networks that place internal representations in the sweet spot suggested by the theory. In particular we aim to match sparse activations with sparse weights in relatively high dimensional settings. A boosting rule was used to increase the overall entropy of the internal layers in order to maximize the utilization of the representational space. We showed that this formulation increases the overall robustness of the system to noisy inputs using MNIST and the Google Speech Command Dataset. Both dense and sparse networks showed high accuracies, but the sparse nets were significantly more robust. These results suggest that it is important to look beyond pure test set performance as test accuracy by itself is not a reliable indicator of overall robustness.

Our work extends the existing literature on sparsity and pruning. A very recent theoretical paper showed that simple linear sparse networks may be more robust to adversarial attacks (Guo et al., 2018). A number of papers have shown that it is possible to effectively introduce sparsity through pruning and retraining (Han et al., 2015; Frankle & Carbin, 2018; Lee et al., 2018). The mechanisms introduced here can be seen as complementary to those techniques. Our network enforces sparse weights from the beginning by construction, and sparse weights are learned as part of the training process. In addition, we reduce the overall computational complexity by enforcing sparse activations, which in turn significantly reduces the number of overall non-zero products. This should produce significant power savings for optimized hardware implementations.

We demonstrated increased robustness in our networks whereas the papers on pruning typically do not explicitly

test robustness. It is possible that such networks are also more robust, though this remains to be tested. Pruning techniques in general are quite orthogonal to ours, and it may be feasible to combine them with the mechanisms discussed here.

In our work we did not attempt to introduce sparsity into the convolutional filters themselves. (Li et al., 2016) have shown it is sometimes possible to remove entire filters from large CNNs suggesting that sparsifying filter weights may also be possible, particularly in networks with larger filters. Introducing sparse convolutions within the context of the techniques in this paper is an area of future exploration. The techniques described here are straightforward to implement and can be extended to other architectures including RNNs. This is yet another promising area for future research.

5.1. Software

All code and experiments are available at <https://github.com/numenta/htmpapers> as open source.

Acknowledgements

We thank Jeff Hawkins, Ali Rahimi, and John Berkowitz for helpful discussions and comments.

References

- Ahmad, S., & Hawkins, J. (2016). How do neurons operate on sparse distributed representations? A mathematical theory of sparsity, neurons and active dendrites. *arXiv*, (pp. arXiv:1601.00720 [q-bio.NC]).
URL <https://arxiv.org/abs/1601.00720>

- Chen, Y., Paiton, D., & Olshausen, B. (2018). The Sparse Manifold Transform. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.) *Advances in Neural Information Processing Systems 31*, (pp. 10533–10544). Curran Associates, Inc.
- Cui, Y., Ahmad, S., & Hawkins, J. (2017). The HTM Spatial Pooler – a neocortical algorithm for online sparse distributed coding. *Frontiers in Computational Neuroscience, 11*, 111.
URL <https://www.frontiersin.org/articles/10.3389/fncom.2017.00111/abstract>
- Frankle, J., & Carbin, M. (2018). The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks.
URL <http://arxiv.org/abs/1803.03635>
- Guo, Y., Zhang, C., Zhang, C., & Chen, Y. (2018). Sparse DNNs with Improved Adversarial Robustness. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.) *Advances in Neural Information Processing Systems 31*, (pp. 240–249). Curran Associates, Inc.
- Han, S., Pool, J., Tran, J., & Dally, W. (2015). Learning both Weights and Connections for Efficient Neural Network. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, & R. Garnett (Eds.) *Advances in Neural Information Processing Systems 28*, (pp. 1135–1143). Curran Associates, Inc.
- Hawkins, J., Ahmad, S., & Dubinsky, D. (2011). Cortical Learning Algorithm and Hierarchical Temporal Memory.
URL <http://numenta.org/resources/HTM{ }CorticalLearningAlgorithms.pdf>
- He, K., Zhang, X., Ren, S., & Sun, J. (2015a). Deep Residual Learning for Image Recognition.
URL <http://arxiv.org/abs/1512.03385>
- He, K., Zhang, X., Ren, S., & Sun, J. (2015b). Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.
URL <http://arxiv.org/abs/1502.01852>
- Ioffe, S., & Szegedy, C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.
URL <http://arxiv.org/abs/1502.03167>
- Kanerva, P. (1988). *Sparse Distributed Memory*. Cambridge, MA: The MIT Press.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., & Jackel, L. D. (1989). Backpropagation Applied to Handwritten Zip Code Recognition.
- LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*.
- Lee, H., Ekanadham, C., & Ng, A. Y. (2008). Sparse deep belief net model for visual area V2. *Advances In Neural Information Processing Systems*.
- Lee, H., Grosse, R., Ranganath, R., & Ng, A. Y. (2009). Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. *Proceedings of the 26th Annual International Conference on Machine Learning - ICML '09*, (pp. 1–8).
- Lee, N., Ajanthan, T., & Torr, P. H. S. (2018). SNIP: Single-shot Network Pruning based on Connection Sensitivity.
URL <http://arxiv.org/abs/1810.02340>
- Li, H., Kadav, A., Durdanovic, I., Samet, H., & Graf, H. P. (2016). Pruning Filters for Efficient ConvNets.
URL <http://arxiv.org/abs/1608.08710>
- Majani, E., Erlanson, R., & Abu-Mostafa, Y. S. (1989). On the k-winners-take-all network. In *Advances in neural information processing systems*, (pp. 634–642).
- Makhzani, A., & Frey, B. (2013). k-Sparse Autoencoders.
URL <http://arxiv.org/abs/1312.5663>
- Makhzani, A., & Frey, B. (2015). Winner-take-all autoencoders. *Advances in Neural Information Processing*.
URL <http://papers.nips.cc/paper/5783-winner-take-all-autoencoders>
- Molchanov, D., Ashukha, A., & Vetrov, D. (2017). Variational Dropout Sparsifies Deep Neural Networks.
URL <http://arxiv.org/abs/1701.05369>
- Nair, V., & Hinton, G. E. (2009). 3D Object Recognition with Deep Belief Nets. In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, & A. Culotta (Eds.) *Advances in Neural Information Processing Systems 22*, (pp. 1339–1347). Curran Associates, Inc.
- Olshausen, B. A., & Field, D. J. (1997). Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research, 37*, 3311–3325.
- Rawlinson, D., Ahmed, A., & Kowadlo, G. (2018). Sparse Unsupervised Capsules Generalize Better.
URL <http://arxiv.org/abs/1804.06094>
- Rosenfeld, A., Zemel, R., & Tsotsos, J. K. (2018). The Elephant in the Room.
URL <http://arxiv.org/abs/1808.03305>
- Sainath, T. N., & Parada, C. (2015). Convolutional neural networks for small-footprint keyword spotting. In *Sixteenth Annual Conference of the International Speech Communication Association*.

- Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition.
URL <http://arxiv.org/abs/1409.1556>
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15, 1929–1958.
URL <http://jmlr.org/papers/v15/srivastava14a.html>
- Srivastava, R. K., Masci, J., Kazerounian, S., Gomez, F., & Schmidhuber, J. (2013). Compete to Compute. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, & K. Q. Weinberger (Eds.) *Advances in Neural Information Processing Systems* 26, (pp. 2310–2318). Curran Associates, Inc.
URL <http://papers.nips.cc/paper/5059-compete-to-compute.pdf>
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., & Fergus, R. (2013). Intriguing properties of neural networks.
URL <http://arxiv.org/abs/1312.6199>
- Tang, R., & Lin, J. (2017). Deep Residual Learning for Small-Footprint Keyword Spotting.
URL <https://arxiv.org/abs/1710.10361>
- Tuguldur, E.-O. (2018). pytorch-speech-commands.
URL <https://github.com/tugstugi/pytorch-speech-commands>
- Warden, P. (2017). Speech Commands: A public dataset for single-word speech recognition. *Dataset available from* http://download.tensorflow.org/data/speech_commands_v0.01.tar.gz.