

The key idea is to randomly drop units (along with their connections) from the neural network during training. This prevents units from co-adapting too much. During training, dropout samples from an exponential number of "thinned" networks. At test time, it is easy to approximate the effect of averaging the predictions of all these thinned networks by simply using a single unthinned network that has smaller weights. This greatly reduces variance and gives a major improvement over other regularization methods. We show that dropout improves the performance of neural networks on supervised learning tasks in vision, speech recognition, document classification, and computational biology, obtaining state-of-the-art results on many benchmark datasets. Keywords: neural networks, regularization, model combination, deep learning

1. Introduction

Deep neural networks contain multiple non-linear hidden layers and this makes them very expressive models that can learn very complicated relationships between their inputs and outputs. With limited training data, however, many of these complicated relationships will be the result of sampling noise so they will exist in the training set but not in real test data, even if it is drawn from the same distribution. This leads to overfitting and many methods have been developed for reducing it. These include stopping the training as soon as performance on a validation set starts to get worse, introducing weight penalties of various kinds such as L1 and L2 regularization and soft weight sharing (Nowlan and Hinton 1992). Srivastava, Hinton, Saxe, and Sengco (2014).

(a) Standard Neural Net

(b) After applying dropout. Its posterior probability given the training data. This can sometimes be approximated well for simple or small models (Xiong et al. 2011; Salakhutdinov and Mnih 2008) but we would like to approach the performance of the Bayesian gold standard using considerably less computation. We propose to do this by approximating an equally weighted geometric mean of the predictions of an exponential number of learned models that share parameters. Dropout is a technique that provides a way of approximately combining exponentially many neural network architectures. The term "dropout" refers to dropping out units (hidden and visible) in a neural network. By dropping a unit out, we mean temporarily removing it from the network along with all its incoming and outgoing connections as shown in Figure 1. 1930

Dropout

Present with probability

p
 w
 -

(a) At training time

Always present

p
 w
 -

(b) At test time

Figure 2:

Left

: A unit at training time that is present with probability

p
 and is disconnected to units
 in the next layer with weights
 w

. Applying dropout to a neural network amounts to sampling a "thinned" network from it. The thinned network consists of all the units that survived dropout (Figure 1b). A neural net with

n units can be seen as a collection of 2^n possible thinned neural networks. These networks all share weights so that the total number of parameters is $O(n^2)$ or less. For each presentation of each training case a new thinned network is sampled and trained. So training a neural network with dropout can be seen as training a collection of 2^n thinned networks with extensive weight sharing where each thinned network gets trained very rarely if at all. At test time it is not feasible to explicitly average the predictions from exponentially many thinned models. However a very simple approximate averaging method works well in practice. The idea is to use a single neural net at test time without dropout. The weights of this network are scaled-down versions of the trained weights. If a unit is retained with probability p during training the outgoing weights of that unit are multiplied by p at test time as shown in Figure 2. This ensures that for any hidden unit the expected output (under the distribution used to drop units at training time) is the same as the actual output at test time. By doing this scaling 2^n networks with shared weights can be combined into a single neural network to be used at test time. We found that training a network with dropout and using this approximate averaging method at test time leads to a lower generalization error on a wide variety of problems compared to training with other regularization methods. The idea of dropout is not limited to feed-forward neural nets. It can be more broadly applied to graphical models such as Boltzmann Machines. In this paper we introduce the dropout Restricted Boltzmann Machine model and compare it to standard Restricted Boltzmann Machines (RBM). Our experiments show that dropout RBMs are better than standard RBMs in certain respects.

3. Related Work

Dropout can be interpreted as a way of regularizing a neural network by adding noise to its hidden units. The idea of adding noise to the states of units has previously been used in the context of Denoising Autoencoders (DAEs) by Vincent et al. (2008, 2010) where noise

Dropout is added to the input units of an autoencoder and the network is trained to reconstruct the noise-free input. Our work extends this idea by showing that dropout can be applied in the hidden layers as well and that it can be interpreted as a form of model averaging. We also show that adding noise is not only useful for unsupervised feature learning but can also be extended to supervised learning problems. In fact our method can be applied to other neuron-based architectures for example Boltzmann Machines. While 5% noise typically works best for DAEs we found that our weight scaling procedure applied at test time enables us to use much higher noise levels. Dropping out 20% of the input units and 50% of the hidden units was often found to be optimal. (2013) also explored deterministic regularizers corresponding to exponential-family noise

distributions including dropout (which they refer to as "blankout noise"). However they apply noise to the inputs and only explore models with no hidden layers. Wang and Manning (2013) proposed a method for speeding up dropout by marginalizing dropout noise. Chen et al. (2012) explored marginalization in the context of denoising autoencoders. This can be seen as minimizing Roweis (2006); Deke et al. (2010) explored an alternate setting where the loss is minimized when an adversary gets to pick which units to drop. Here instead of a noisy distribution the maximum number of units that can be dropped is L . However this work also does not explore models with hidden units.

4. Model Description

This section describes the dropout neural network model. Consider a neural network with

L hidden layers. Let

$1 \leq l \leq L$ denote the index of the hidden layers of the network. Let

$\mathbf{z}^{(l)}$ denote the vector of inputs into layer l

$\mathbf{y}^{(l)}$ denote the vector of outputs from layer l

$\mathbf{y}^{(0)}$ is the input). $\mathbf{W}^{(l)}$

and $\mathbf{b}^{(l)}$ are the weights and biases at layer l

. The feed-forward operation of a standard neural network (Figure 3a) can be described as (for

$1 \leq l \leq L$)

;;::;L

l
g
and
anyhiddenunit

i
)

z

(

l

+1)

i

=

w

(

l

+1)

i

y

l

+

b

(

l

+1)

i

;

y

(

l

+1)

i

=

f

(

z

(

l

+1)

i

)

;

where

f

isanyactivationfunctionforexample

f

(

x

)=1

=

(1+exp(

x
 $\left. \right)$). With dropout the feed-forward operation becomes (Figure 3b)
 r
 $($
 1
 $)$
 j
Bernoulli(
 p
 $)$
 $;$
 e
 y
 $($
 1
 $)$
 $=$
 r
 $($
 1
 $)$

 y
 $($
 1
 $)$
 $;$
 z
 $($
 1
 $+1)$
 i
 $=$
 w
 $($
 1
 $+1)$
 i
 e
 y
 l
 $+$
 b
 $($
 1
 $+1)$
 i
 $;$
 y
 $($
 1
 $+1)$
 i

=
f
(
z
(
l
+1)
i
)
:
1933

Srivastava Hinton Krizhevsky Sutskever and Salakhutdinov

(a) Standard network

(b) Dropout network

Figure 3:

Comparison of the basic operations of a standard and dropout network. Here

denotes an element-wise product. For any layer

l
r
(
l
)
is a vector of independent
Bernoulli random variables each of which has probability
p
of being 1. This vector is
sampled and multiplied element-wise with the output of that layer

y
(
l
)
to create the
thinned outputs

e
y
(
l
)
. The thinned outputs are then used as input to the next layer. This
process is applied at each layer. This amounts to sampling a sub-network from a larger
network. For learning the derivatives of the loss function are backpropagated through the
sub-network. At test time the weights are scaled as

W
(
l
)
test
=
pW
(

1
)
as shown in Figure 2. One particular form of regularization was found to be especially useful for dropout|
constraining the norm of the incoming weight vector at each hidden unit to be upper
bounded by a constant

c
. In other words if
w
represents the vector of weights incident
on any hidden unit the neural network was optimized under the constraint
 $\|w\|_2 \leq c$

c
. This
constraint was imposed during optimization by projecting
w
onto the surface of a ball of
radius
c
whenever
w
went out of fit. This is also called max-norm regularization since it
implies that the maximum value that the norm of any weight can take is
c
. The constant
1934

Dropout

c
is a tunable hyperparameter which is determined using a validation set. Max-norm
regularization has been previously used in the context of collaborative (Srebro and
Shraibman 2005). It typically improves the performance of stochastic gradient descent
training of deep neural networks even when no dropout is used. 5.2 Unsupervised Pretraining
Neural networks can be pretrained using stacks of RBMs (Hinton and Salakhutdinov 2006)
autoencoders (Vincent et al. 2010) or Deep Boltzmann Machines (Salakhutdinov and Hin-
ton 2009). Pretraining is a new way of making use of unlabeled data. Pretraining
followed by with backpropagation has been shown to give performance
boosts over random initializations in certain cases. CIFAR-10 and CIFAR-100: Tiny natural images (F
Google Street View (Netzer et al. 2011). 1935

Srivastava Hinton Krizhevsky Sutskever and Salakhutdinov

Alternative Splicing dataset: RNA features for predicting alternative genes splicing
(Xiong et al. 2011). Data Set Domain Dimensionality Training Set Test Set
MNIST Vision 784 (28

28 grayscale) 60K 10K
SVHN Vision 3072 (32

32 color) 600K 26K

CIFAR-10/100 Vision 3072(32

32color) 60K 10K

ImageNet (ILSVRC-2012) Vision 65536(256

256color) 1.2M 150K

TIMIT Speech 2520(120-dim 21 frames) 1.1M frames 58K frames

Reuters-RCV1 Text 2000 200K 200K

Alternative Splicing Genetics 10142932733

Table 1:

Overview of the datasets used in this paper. 6.1.1 MNIST

Method

Unit

Type

Architecture

Error

%

Standard Neural Net (Simard et al. 2003) Logistic 2 layers 800 units 1.60

SVM Gaussian kernel NANA 1.40

Dropout NN Logistic 3 layers 1024 units 1.35

Dropout NN ReLU 3 layers 1024 units 1.25

Dropout NN + max-norm constraint ReLU 3 layers 1024 units 1.06

Dropout NN + max-norm constraint ReLU 3 layers 2048 units 1.04

Dropout NN + max-norm constraint ReLU 2 layers 4096 units 1.01

Dropout NN + max-norm constraint ReLU 2 layers 8192 units 0.95

Dropout NN + max-norm constraint (Goodfellow et al. 2013)

Maxout

2 layers (5

240)

units

0.94

DBN + ng (Hinton and Salakhutdinov 2006) Logistic 500-500-2000 1.18

DBM + (Salakhutdinov and Hinton 2009) Logistic 500-500-2000 0.96

DBN + dropout Logistic 500-500-2000 0.92

DBM + dropout tuning Logistic 500-500-2000

0.79

Table 2:

Comparison of model on MNIST. The MNIST dataset consists of 28

28 pixel handwritten digit images. The task is

to classify the images into 10 digit classes. Table 2 compares the performance of dropout

with other techniques. The best performing neural networks for the permutation invariant

1936

Dropout

setting that does not use dropout or unsupervised pretraining achieve an error of about

1

:

60% (Simard et al. 2003). With dropout the error reduces to 1

:

35%. Replacing logistic

units with linear units (ReLU) (Jarrett et al. 2009) further reduce the error to

1

:

25%. Adding max-norm regularization again reduces it to 1

:

06%. Increasing the size of

the network leads to better results. A neural network with 2 layers and 8192 units per layer

gets down to 0.95% error. Note that this network has more than 65 million parameters and

is being trained on a dataset of size 60000. Training a network of this size to give good

generalization error is very hard with standard regularization methods and early stopping. (2013) showed that

0.94% by replacing ReLU units with maxout units. All dropout nets use

p

= 0

:

5 for hidden

units and

p

= 0

:

8 for input units. More experimental details can be found in Appendix B.1. 6.1.2 Street View House Numbers

The Street View House Numbers (SVHN)

Dataset (Netzer et al. 2011) consists of

color images of house numbers collected by

Google Street View. Figure 5 shows some examples of images from this dataset. The

part of the dataset that we use in our experiments consists of 32

32 color images roughly

centered on a digit in a house number. The task is to identify that digit. (1989). The best architecture that we found

fully connected hidden layers. All hidden units were ReLUs. Each convolutional layer was

1937

Srivastava Hinton Krizhevsky Sutskever and Salakhutdinov

Method Error %

Binary Features (WDCH) (Netzer et al. 2011) 36.7

HOG (Netzer et al. 2011) 15.0

Stacked Sparse Autoencoders (Netzer et al. 2011) 10.3

KMeans (Netzer et al. 2011) 9.4

Multi-stage ConvNet with average pooling (Sermanet et al. 2012) 9.06

Multi-stage ConvNet + L2 pooling (Sermanet et al. 2012) 5.36

Multi-stage ConvNet + L4 pooling + padding (Sermanet et al. 2012) 4.90

ConvNet + max-pooling 3.95

ConvNet + max-pooling + dropout in fully connected layers 3.02

ConvNet + stochastic pooling (Zeiler and Fergus 2013) 2.80

ConvNet + max-pooling + dropout in all layers 2.55

ConvNet + maxout (Goodfellow et al. 2013)

2.47

Human Performance 2.0

Table 3:

Results on the Street View House Numbers dataset. Dropout was applied to all the layers of the network with the

den unit being

p

= (0

:

9
;
0
:
75
;
0
:
75
;
0
:
5
;
0
:
5
;
0
:
5)for thet layers of the network (going
from input to convolutional layers to fully connected layers). Max-norm regularization was
used for weights in both convolutional and fully connected layers. Table 3 compares the
results obtained by the methods. We find that convolutional nets outperform other
methods. The best performing convolutional nets that do not use dropout achieve an error
rate of 3

95%. Adding dropout only to the fully connected layers reduces the error to 3%.

The additional gain in performance obtained by adding dropout in the convolutional layers (3

02%to2

55%) is worth noting. One may have presumed that since the convolutional layers don't have a lot of parameters so it's not a problem and therefore dropout would not have much effect. However, dropout in the lower layers still helps because it provides noisy inputs for the higher fully connected layers which prevents them from overfitting.

6.1.3 CIFAR-10 and CIFAR-100

The CIFAR-10 and CIFAR-100 datasets consist of 32

32colorimagesdrawnfrom10

and 100 categories respectively. Figure 5b shows some examples of images from this dataset. A detailed description of the dataset input preprocessing network architectures and other experimental details is given in Appendix B.3. Table 4 shows the error rate obtained by methods on these datasets. Without any data augmentation Snoek et al. (2012) used Bayesian hyperparameter optimization to find the best architecture for CIFAR-10. Using dropout in the fully connected layers reduces that to 14.32% and adding dropout in every layer further reduces the error to 12.61%. Goodfellow et al. (2013) showed that the error is further reduced to 11.68% by replacing ReLU units with maxout units. On CIFAR-100 dropout reduces the error from 43.48% to 37.20% which is a huge improvement. 1938

Dropout

(a) Street View House Numbers (SVHN)

(b) CIFAR-10

Figure 5:

Samples from image datasets. Each row corresponds to a category. Method CIFAR-10 CIFAR-100

ConvNet+maxpooling (hand-tuned) 15.60 43.48

ConvNet+stochastic pooling (Zeiler and Fergus 2013) 15.13 42.51

ConvNet+maxpooling (Snoek et al. 2012) 14.98-

ConvNet+maxpooling+dropout fully connected layers 14.32 41.26

ConvNet+maxpooling+dropout in all layers 12.61

37.20

ConvNet+maxout (Goodfellow et al. 2013)

11.68

38.57

Table 4:

Error rates on CIFAR-10 and CIFAR-100. 6.1.4 ImageNet

ImageNet is a dataset of over 15 million labeled high-resolution images belonging to roughly

22,000 categories. Starting in 2010 as part of the Pascal Visual Object Challenge, an annual

competition called the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) has

been held. A subset of ImageNet with roughly 1,000 images in each of 1,000 categories is

used in this challenge. Since the number of categories is rather large, it is conventional to

report two error rates: top-1 and top-5, where the top-5 error rate is the fraction of test

images for which the correct label is not among the labels considered most probable by

the model. Figure 6 shows some predictions made by our model on a few test images. ILSVRC-2010 is the only

most of four experiments were performed on this dataset. Table 5 compares the performance

of methods. Convolutional nets with dropout outperform other methods by a large

margin. The architecture and implementation details are described in detail in Krizhevsky

et al. (2012). Model Top-1 Top-5

Sparse Coding (Lin et al. 2010) 47.1 28.2

SIFT+Fisher Vectors (Sanchez and Perronnin 2011) 45.7 25.7

ConvNet+dropout (Krizhevsky et al. 2012) 37.5 17.0

Table 5:

Results on the ILSVRC-2010 test set. Model

Top-1

(val)

Top-5

(val)

Top-5

(test)

SVM on Fisher Vectors of Dense SIFT and Color Statistics -- 27.3

Avg of over FV of SIFT, LBP, GIST and CSIFT -- 26.2

ConvNet+dropout (Krizhevsky et al. 2012) 40.7 18.2-

Avg of 5 ConvNets+dropout (Krizhevsky et al. 2012) 38.1 16.4 16.4

Table 6:

Results on the ILSVRC-2012 validation/test set. Method Phone Error Rate %

NN (6 layers) (Mohamed et al. 2010) 23.4

Dropout NN (6 layers) 21.8

DBN-pretrained NN (4 layers) 22.7

DBN-pretrained NN (6 layers) (Mohamed et al. 2010) 22.4

DBN-pretrained NN (8 layers) (Mohamed et al. 2010) 20.7

mcRBM-DBN-pretrained NN (5 layers) (Dahl et al. 2010) 20.5

DBN-pretrained NN (4 layers)+dropout

19.7

DBN-pretrained NN (8 layers)+dropout

Table 7:

Phone error rate on the TIMIT core test set. 6.4 Comparison with Bayesian Neural Networks Dropout can be seen as a way of doing an equally-weighted averaging of exponentially many models with shared weights. On the other hand Bayesian neural networks (Neal 1996) are the proper way of doing model averaging over the space of neural network structures and parameters. In dropout each model is weighted equally whereas in a Bayesian neural network each model is weighted taking into account the prior and how well the model fits the data which is the more correct approach. Bayesian neural networks are extremely useful for solving problems in domains where data is scarce such as medical diagnosis, genetics, drug discovery and other computational biology applications. However Bayesian neural networks are slow to train and to scale to very large network sizes. Besides it is expensive to get predictions from many large networks at test time. On the other hand dropout neural networks are much faster to train and use at test time. In this section we report experiments that compare Bayesian neural networks with dropout neural networks on a small dataset where Bayesian neural networks are known to perform well and obtain state-of-the-art results. The aim is to analyze how much dropout loses compared to Bayesian neural networks. The dataset that we use (Xiong et al. 2011) is a task to predict the occurrence of alternative splicing based on RNA features. Alternative splicing is a cause of cellular diversity in mammalian tissues. Predicting the

1941

Srivastava Hinton Krizhevsky Sutskever and Salakhutdinov

Method Code Quality (bits)

Neural Network (early stopping) (Xiong et al. 2011) 440

Regression PCA (Xiong et al. 2011) 463

SVM PCA (Xiong et al. 2011) 487

Neural Network with dropout 567

Bayesian Neural Network (Xiong et al. 2011)

623

Table 8:

Results on the Alternative Splicing Data Set. occurrence of alternative splicing in certain tissues under certain conditions is important for understanding many human diseases. Given the RNA features the task is to predict the probability of three splicing related events that biologists care about. The evaluation metric is Code Quality which is a measure of the negative KL divergence between the target and the predicted probability distributions (higher is better). Appendix B.6 includes a detailed description of the dataset and this performance metric. (2011) used Bayesian neural networks for this task. As expected Bayesian neural networks perform better than dropout. However we see that dropout improves significantly upon the performance of standard neural networks and outperforms all other methods. The challenge in this dataset is to prevent overfitting since the size of the training set is small. One way to prevent overfitting is to use standard techniques such as SVMs or logistic regression can be used. However with dropout we were able to prevent overfitting without the need to do dimensionality reduction. The dropout networks are very large (1000s of hidden units) compared to a few tens of units in the Bayesian network. This shows that dropout has a strong regularizing effect.

6.5 Comparison with Standard Regularizers

Several regularization methods have been proposed for preventing overfitting in neural networks. These include L2 weight decay (more generally Tikhonov regularization (Tikhonov 1943)), lasso (Tibshirani 1996), KL-sparsity and max-norm regularization. Dropout can be seen as another way of regularizing neural networks. In this section we compared dropout with some of these regularization methods using the MNIST dataset. The values of the hyperparameters associated with each regularization method (sparsity, dropout rate, max-norm upper bound) were obtained using a validation set. We found that dropout combined with max-norm regularization gives the lowest generalization error. 7.1 on Features

(a) Without dropout

(b) Dropout with

p
 $=0$

:

5.7.2 on Sparsity

(a) Without dropout

(b) Dropout with

p
 $=0$

:

5. (a) Keeping

n

(b) Keeping

p
 n

Figure 9:

of changing dropout rates on MNIST. Figure 9b shows the test error obtained as a function of

p

. We notice that the magnitude

of errors for small values of

p

has reduced by a lot compared to Figure 9a (for

p
 $=0$

:

1 it fell

from 2.7% to 1.7%). Values of

p

that are close to 0.6 seem to perform best for this choice

of

p
 n

but our usual default value of 0

:

5 is close to optimal. The results of these experiments are

shown in Figure 10. The network was given

dataset of size 100 500 1K 5K 10K

and 50K chosen randomly from the MNIST

training set. The same network architec-

ture (784-1024-1024-2048-10) was used for

all datasets. Dropout with

p
 $=0$

:

5 was per-

formed at all the hidden layers and

p
 $=0$

:

8

at the input layer. It can be observed that

for extremely small datasets (100 500)

dropout does not give any improvements. 1946

Dropout

We again use the MNIST dataset and do by averaging the predictions

of

k

randomly sampled neural networks. Figure 11 shows the test error rate obtained for

values of

k

. This is compared with the error obtained using the weight scaling

method (shown as a horizontal line). It can be seen that around

k

$= 50$ the Monte-Carlo

method becomes as good as the approximate method. Thereafter the Monte-Carlo method

is slightly better than the approximate method but well within one standard deviation of

it. This suggests that the weight scaling method is a fairly good approximation of the true

model average.

8.1 Model Description

Consider an RBM with visible units

v

$2f$

0

;

1

g

D

and hidden units

h

$2f$

0

;

1

g

F

. It

the following probability distribution

P

(

h

;

v

;

) =

1

Z

(

)

\exp (

v

$>$

W

h

+

a
 $>$
 h
 $+$
 b
 $>$
 v
 $)$
 $:$
 Where

$=$
 f
 $W;$
 a
 $;$
 b
 g
 represents the model parameters and
 Z
 is the partition function. Dropout RBMs are RBMs augmented with a vector of binary random variables
 r
 2
 f
 0
 $;$
 1
 g
 F
 .Each random variable
 r
 j
 takes the value 1 with probability
 p
 independent of
 others. If
 r
 j
 takes the value 1 the hidden unit
 h
 j
 is retained otherwise it is dropped from
 the model. The joint distribution by a Dropout RBM can be expressed as
 P
 $($
 r
 $;$
 h
 $;$
 v
 $;$
 $p;$
 $)=$

$$\begin{aligned}
 &P \\
 &(\quad \\
 &\quad r \\
 &\quad ; \\
 &\quad p \\
 &)\quad P \\
 &(\quad h \\
 &\quad ; \\
 &\quad v \\
 &\quad j \\
 &\quad r \\
 &\quad ; \\
 &\quad) \\
 &\quad ; \\
 &P \\
 &(\quad r \\
 &\quad ; \\
 &\quad p \\
 &)= \\
 &F \\
 &Y \\
 &j \\
 &=1 \\
 &p \\
 &r \\
 &j \\
 &(1 \\
 & \\
 &p \\
 &)\quad 1 \\
 & \\
 &r \\
 &j \\
 &\quad ; \\
 &P \\
 &(\quad h \\
 &\quad ; \\
 &\quad v \\
 &\quad j \\
 &\quad r \\
 &\quad ; \\
 &)= \\
 &1 \\
 &Z \\
 &0 \\
 &(\quad
 \end{aligned}$$

```

;
r
)
exp(
v
>
W
h
+
a
>
h
+
b
>
v
)
F
Y
j
=1
g
(
h
j
;r
j
)
;
g
(
h
j
;r
j
)=
1
(
r
j
=1)+
1
(
r
j
=0)
1
(
h
j
=0)
:
Z

```

0
 $($
 $;$
 r
 $)$ is the normalization constant. g
 $($
 h
 j
 $;$
 r
 j
 $)$ imposes the constraint that if
 r
 j
 $= 0$
 h
 j
must be 0. The distribution over
 h
conditioned on
 v
and
 r
is factorial
 P
 $($
 h
 j
 r
 $;$
 v
 $) =$
 F
 Y
 j
 $= 1$
 P
 $($
 h
 j
 j
 r
 j
 $;$
 v
 $)$
 $;$
 P
 $($
 h
 j
 $= 1$
 j
 r

$$\prod_{j=1}^J \prod_{i=1}^I p_{ij}$$

$$\prod_{j=1}^J \prod_{i=1}^I p_{ij}$$

Srivastava Hinton Krizhevsky Sutskever and Salakhutdinov
 (a) Without dropout
 (b) Dropout with

$$p_{ij}$$

5. The distribution over

conditioned on \mathbf{h}
 is the same as that of an RBM

$$P(\mathbf{v} | \mathbf{h}) = \prod_{i=1}^I \prod_{j=1}^J p_{ij}$$

$$P(\mathbf{v} | \mathbf{h}) = \prod_{i=1}^I \prod_{j=1}^J p_{ij}$$

$v_i = 1$
 j
 h
 $) = 0$
 $@$
 a
 i
 $+$
 X
 j
 W
 ij
 h
 j
 1
 A
 $:$
 Conditioned on
 r
 the distribution over
 f
 v
 $;$
 h
 g
 is same as the distribution that an RBM
 would impose except that the units for which

r
 j
 $= 0$ are dropped from
 h

. Therefore the
 Dropout RBM model can be seen as a mixture of exponentially many RBMs with shared
 weights each using a subset of
 h
 . Figure 12a shows features learned by a binary RBM with 256 hidden units. Figure 12b
 shows features learned by a dropout RBM with the same number of hidden units. Features
 1948

Dropout
 (a) Without dropout
 (b) Dropout with

p
 $= 0$
 $:$
 5. 9. Marginalizing Dropout
 Dropout can be seen as a way of adding noise to the states of hidden units in a neural
 network. In this section we explore the class of models that arise as a result of marginalizing
 this noise. These models can be seen as deterministic versions of dropout. In contrast to
 standard ("Monte-Carlo") dropout these models do not need random bits and it is possible

to get gradients for the marginalized loss functions. In this section we explore these models. Deterministic algorithms have been proposed that try to learn models that are robust to feature deletion at test time (Globerson and Roweis 2006). Marginalization in the context of denoising autoencoders has been explored previously (Chen et al. 2012). The marginalization of dropout noise in the context of linear regression was discussed in Srivastava (2013). Wang and Manning (2013) investigated input noise distributions and

Srivastava, Hinton, Krizhevsky, Sutskever, and Salakhutdinov (2013) describe how the regularizers obtained by marginalizing this noise. Wager et al. (2013) describe how dropout can be seen as an adaptive regularizer.

Let \mathbf{y} be a vector of N targets. Linear regression tries to find

weights \mathbf{w} that minimize

$$\sum_{j=1}^N \|\mathbf{X}_j \mathbf{w} - y_j\|^2$$

where \mathbf{X}_j is the input vector for the j -th example. When the input \mathbf{X}_j is dropped out such that any input dimension is retained with probability p , the input can be expressed as

$$\mathbf{X}_j = \mathbf{D}_j \mathbf{X}_j^0$$

where

$$D_{ij} = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{otherwise} \end{cases}$$

\mathbf{D} is a random matrix with elements D_{ij} drawn from a Bernoulli distribution with parameter p .

Let \mathbf{X} be the matrix of input vectors, \mathbf{y} be the vector of targets, and \mathbf{w} be the weight vector. The regularized loss function is

p
)and

denotes an element-wise product. Marginalizing the noise
the objective function becomes
minimize

w
 E
 R
Bernoulli(p)

\hat{y}
 y

(
 R

X
)
 w
 \hat{y}
2

:
This reduces to
minimize
 w
 \hat{y}
 y

pX
 w
 \hat{y}
2
+
 p
(1

p
)
 \hat{y}

w
 \hat{y}
2
;
where $= (\text{diag}($
 X
>
 X
))
1
=

which takes the value 1 with probability

p

and 0 otherwise. This idea can be generalized

by multiplying the activations with random variables drawn from other distributions. We

recently discovered that multiplying by a random variable drawn from

N

$(1$

$;$

$1)$ works just

as well or perhaps better than using Bernoulli noise. This new form of dropout amounts

to adding a Gaussian distributed random variable with zero mean and standard deviation

equal to the activation of the unit. That is each hidden activation

h

i

is perturbed to

h

i

$+$

h

i

r

where

r

$(0$

$;$

$1)$ or equivalently

h

i

r

0

where

r

0

$(1$

$;$

$1)$. We can generalize

this to

r

0

$(1$

2

) where

becomes an additional hyperparameter to tune just like

p

was in the standard (Bernoulli) dropout. The expected value of the activations remains

unchanged therefore no weight scaling is required at test time. For the Gaussian multiplicative noise if we set

2

$= (1$

p

)

$= p$

we end up multiplying

h
 i
by a random variable

r
 g
where

E
[
 r
 g
]=1 and
 Var

[
 r
 g
]=(1

p
)
= p
. Therefore both
forms of dropout can be set up so that the random variable being multiplied by has the
same mean and variance. However given these and second order moments

r
 g
has the
highest entropy and

r
 b
has the lowest. Both these extremes work well although preliminary
experimental results shown in Table 10 suggest that the high entropy case might work
slightly better. For each layer the value of
in the Gaussian model was set to be

q
1

p
 p
using the
 p
from the corresponding layer in the Bernoulli model. For input layers the choice depends on the kind of input. For
patches or speech frames) a typical value is 0

:
8. For hidden layers the choice of

p
is coupled
with the choice of number of hidden units

n
. Smaller

p
requires big

n
which slows down

the training and lead to underfitting

Dropout may not produce enough dropouts to prevent overfitting

Appendix B. Detailed Description of Experiments and Data Sets

This section describes the network architectures and training details for the experimental results reported in this paper. The code for reproducing these results can be obtained from <http://www.cs.toronto.edu/>

Dropout

The implementation is GPU-based. We used the excellent CUDA libraries `cudamat` (Mnih 2009) and `cuda-convnet` (Krizhevsky et al. 2012) to implement our networks.

Srivastava, Hinton, Krizhevsky, Sutskever, and Salakhutdinov

The architectures shown in Figure 4 include all combinations of 2, 3, and 4 layer networks with 1024 and 2048 units in each layer. Thus there are six architectures in all. For all the architectures (including the ones reported in Table 2) we used

$p = 0.5$ in all hidden layers

and

$p = 0.8$ in the input layer. A momentum of 0.9

was used and weight constraints with

$c = 2$

was used in all the layers.

B.2 SVHN

The SVHN dataset consists of approximately 600,000 training images and 26,000 test images. The training set consists of two parts: a standard labeled training set and another set of unlabeled examples that are easy. A validation set was constructed by taking examples from both the parts. Two-thirds of it were taken from the standard set (400 per class) and one-third from the extra set (200 per class) at a total of 60,000 samples. This same process is used by Sermanet et al. (2012). The inputs were RGB pixels normalized to have zero mean and unit variance. Other preprocessing techniques such as global or local contrast normalization or ZCA whitening did not give any noticeable improvements. Each convolutional layer has a

5x5 receptive field applied with a stride of 1 pixel. Each max pooling layer pools 3x3

regions at strides of 2 pixels. The convolutional layers are followed by two fully connected hidden layers having 2048 units each. All units use the linear activation function. Dropout was applied to all the layers of the network with the probability of retaining the unit being

$p = 0.5$

in the hidden layers and $p = 0.8$ in the input layer.

0
:
75
;
0
:
75
;
0
:
5
;
0
:
5
;
0
:
5)for the
layers of the network (going from input to convolutional layers to fully connected
layers). In addition the max-norm constraint with
c
=4 was used for all the weights. A
momentum of 0

:
0.95 was used in all the layers. These hyperparameters were tuned using a
validation set. Since the training set was quite large we did not combine the validation
set with the training set for training. We reported test error of the model that had
smallest validation error. B.4 TIMIT
The open source Kaldi toolkit (Povey et al. 2011) was used to preprocess the data into log-
banks. A monophone system was trained to do a forced alignment and to get labels for
speech frames. Dropout neural networks were trained on windows of 21 consecutive frames
to predict the label of the central frame. No speaker dependent operations were performed. Max-norm constraint
c
=4 was used in all the layers. A momentum of 0.95 with a
high learning rate of 0.1 was used. The learning rate was decayed as

0
(1+
t=T
)
1
. For
DBN pretraining we trained RBMs using CD-1. The variance of each input unit for the
Gaussian RBM was set to 1. For the DBN with dropout we found that in
order to get the best results it was important to use a smaller learning rate (about 0.01). B.5 Reuters
The Reuters RCV1 corpus contains more than 800,000 documents categorized into 103
classes. These classes are arranged in a tree hierarchy. We created a subset of this dataset
consisting of 402,738 articles and a vocabulary of 2000 words comprising of 50 categories
in which each document belongs to exactly one class. The data was split into equal sized
training and test sets. We tried many network architectures and found that dropout gave
improvements in accuracy over all of them. However the improvement was

not as that for the image and speech datasets. This might be explained by the fact that this dataset is quite big (more than 200,000 training examples) and is not a very serious problem. **B.6 Alternative Splicing**
 The alternative splicing dataset consists of data for 3665 cassette exons, 1014 RNA features and 4 tissue types derived from 27 mouse tissues. For each input, the target consists of 4 softmax units (one for each tissue type). Each softmax unit has 3 states (inc

exc

nc) which are of the biological importance. For each softmax unit, the aim is to predict a distribution over these 3 states that matches the observed distribution from wet lab experiments as closely as possible. The evaluation metric is Code Quality, which is a

j data points

j
 X

i
 $= 1$

X

t

2

tissue types

X

s

2f

inc exc nc

g

p

s

$i; t$

$\log($

q

s

t

$($

r

i

$)$

\cdot

p

s

$)$

$;$

where

p

s

$i; t$

is the target probability for state

s

and tissue type

t
 i input
 i
 $;$
 q
 s
 t
 $($
 r
 i
 $)$ is the
predicted probability for state
 s
in issue type
 t
for input
 r
 i
and
 p
 s
is the average of
 p
 s
 $i; t$
over
 i
and
 t
. A value of
 p
 $= 0$
:
5 was used for the hidden layer and
 p
 $= 0$
:
7 for the input layer. Max-norm
regularization with high decaying learning rates was used. Results were averaged across the
same 5 folds used by Xiong et al. (2011). Machine Learning
81(2):149–178 2010. K. Jarrett K. Kavukcuoglu M. Ranzato and Y. LeCun. What is the best multi-stage
architecture for object recognition? In
Proceedings of the International Conference on
Computer Vision (ICCV'09)
. IEEE 2009. Neural Computa-
tion
1(4):541–551 1989. Proceedings of the National Academy of Sciences
107(4):1452–1457 2010. P. Sermanet S. Chintala and Y. LeCun. Convolutional neural networks applied to
numbers digit In
International Conference on Pattern Recognition (ICPR
2012)
2012. Journal of the Royal
Statistical Society. Series B. Methodological

58(1):267{2881996. Bioinformatics
27(18):2554{25622011.