# Advanced Machine Learning

Likhit Nayak

# Dropout



(a) Standard Neural Net

(b) After applying dropout.

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Why should we use dropout?

To make an ensemble of neural networks, they should either have

1. Different architectures
    - Finding optimal hyperparameters for each architecture is a daunting task
    - Training each large network requires a lot of computation
2. Be trained on different data
    - There may not be enough data available to train different networks on different subsets of the data

**Dropout is a technique that addresses both these issues**

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Implementing Dropout

In a neural network with $L$ hidden layers and no dropout, the feed-forward operation can be described as:

$$
\begin{aligned}
z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)} \mathbf{y}^l + b_i^{(l+1)}, \\
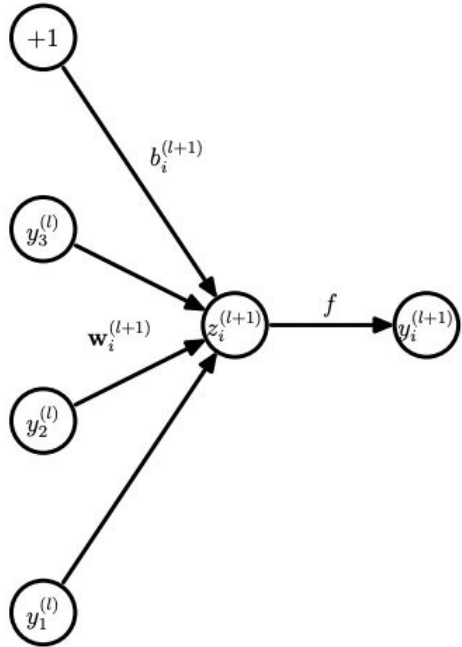y_i^{(l+1)} &= f(z_i^{(l+1)}),
\end{aligned}
$$

where $z$ is the input into a layer, $y$ is the output from a layer, $w$ and $b$ are the weights and

biases of a layer, and $f()$ is the activation function

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.
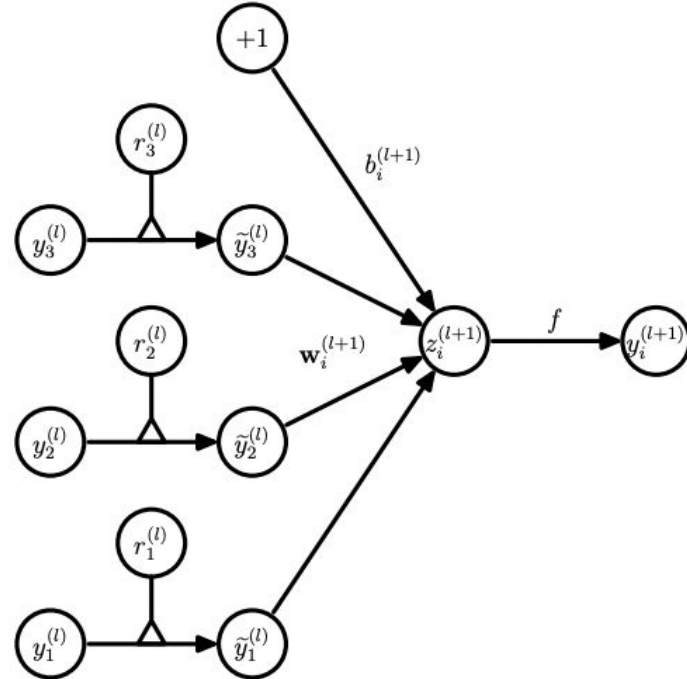
# Implementing Dropout

With dropout, the feed-forward operation becomes:

$$
\begin{aligned}
r_j^{(l)} &\sim \text{Bernoulli}(p), \\
\widetilde{\mathbf{y}}^{(l)} &= \mathbf{r}^{(l)} * \mathbf{y}^{(l)}, \\
z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)}\widetilde{\mathbf{y}}^l + b_i^{(l+1)}, \\
y_i^{(l+1)} &= f(z_i^{(l+1)}).
\end{aligned}
$$

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Implementing Dropout



(a) Standard network  (b) Dropout network

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Implementing Dropout - Inference time

At inference time, it is not feasible to explicitly average the predictions all the different models like we do with bagging. So, we use scaling:



**Present with probability $p$** ... $\mathbf{w}$

(a) At training time

**Always present** ... $p\mathbf{w}$

(b) At test time

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Implementing Dropout - Results



Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Implementing Dropout - Results

| Method | Test Classification error % |
| --- | --- |
| L2 | 1.62 |
| L2 + L1 applied towards the end of training | 1.60 |
| L2 + KL-sparsity | 1.55 |
| Max-norm | 1.35 |
| Dropout + L2 | 1.25 |
| Dropout + Max-norm | **1.05** |

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Advantages of Dropout

1.  Computationally cheap
●  During training, dropout requires only $O(n)$ computations per example per update, to generate $n$ random numbers and multiply them
●  During testing, the cost of dividing the weights (scaling) is a single operation per example
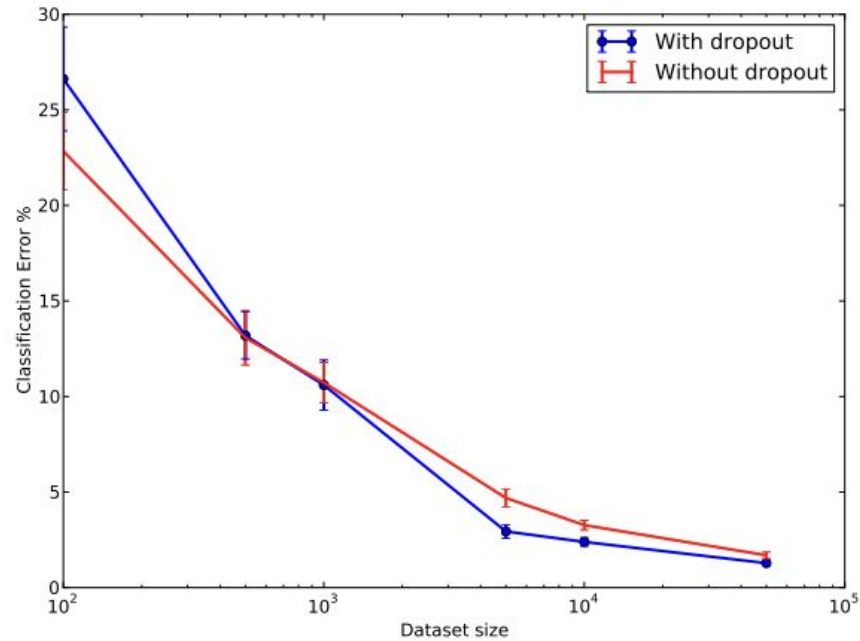
# Advantages of Dropout

2. It works well with nearly any model

| 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 | **19.7** |

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Limitations of Dropout

1.  Works better with larger training datasets



Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *The journal of machine learning research* 15.1 (2014): 1929-1958.

# Limitations of Dropout

2. The cost function (or loss function) isn't well-defined