

11-695: AI Engineering

Other Techniques II

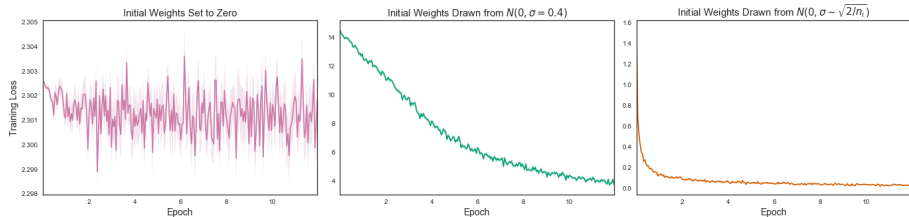
LTI/SCS

Spring 2020

① Initialization

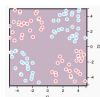
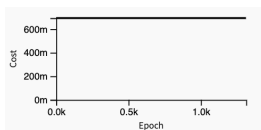
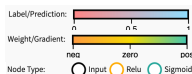
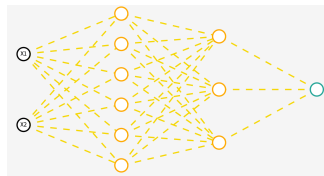
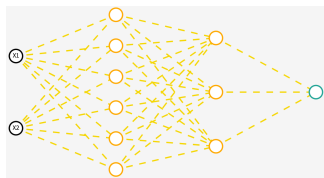
② Weights regularization

③ Dropout



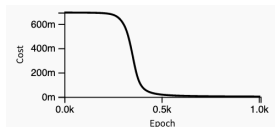
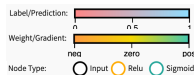
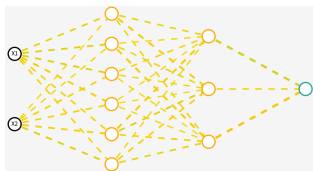
- Training NNs is difficult!¹
- Weight init has an important role in training NNs

¹<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

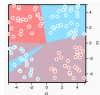
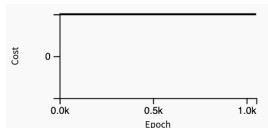
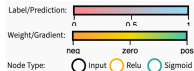
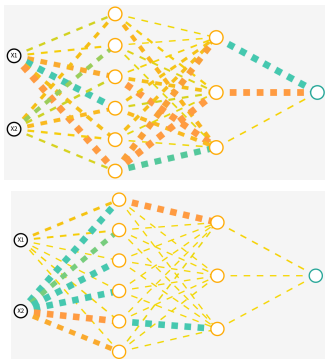


- A big NO, would create a symmetry between hidden layers²
- Bias can be initialized zeros, because weights take care of avoiding symmetry

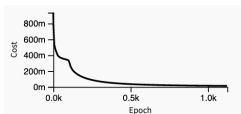
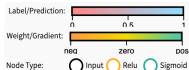
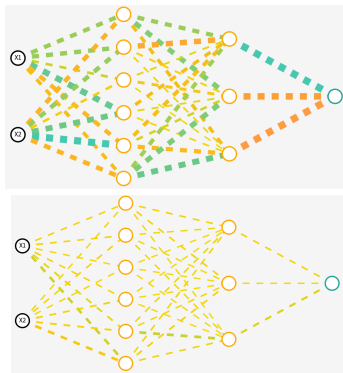
²<https://arxiv.org/pdf/1206.5533.pdf>



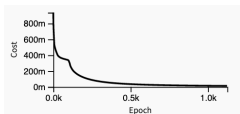
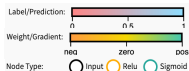
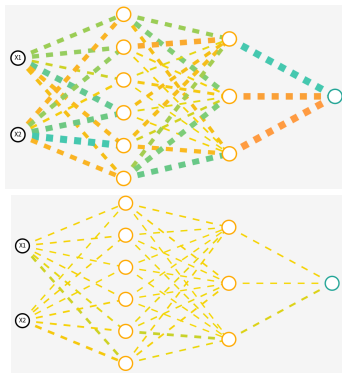
- Gradient vanishing problem



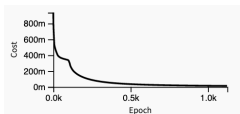
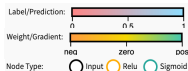
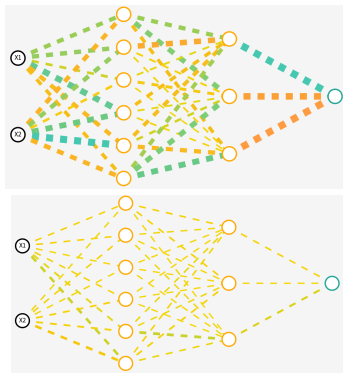
- Gradient exploding problem



- Would have no such problems



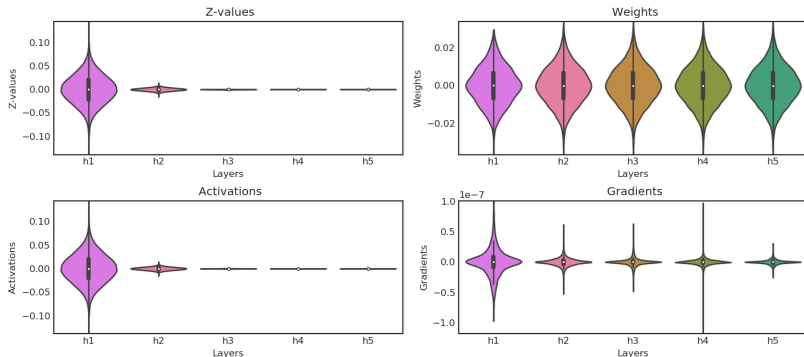
- Would have no such problems
- But what it means by being “proper”



- Would have no such problems
- But what it means by being “proper”
 - Mean should center around zero
 - Variance should not shift between layers

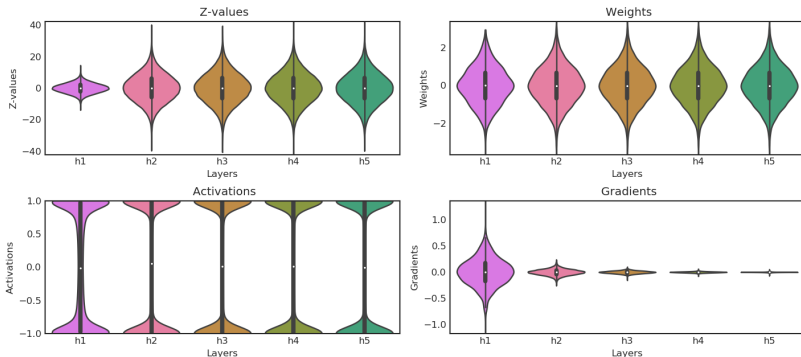
Image Credit: <https://www.deeplearning.ai/ai-notes/initialization/>

Activation: tanh - Initializer: Normal $\sigma = 0.01$ - Epoch 0



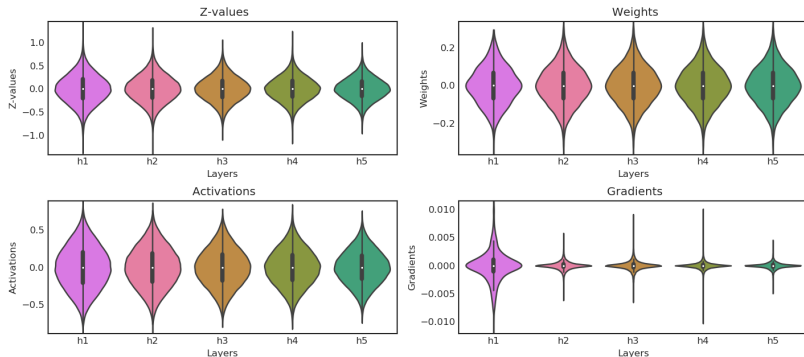
- Normally from $\mathbf{N}(\mathbf{0}, \sigma^2)$
- Variance *matters*

Activation: tanh - Initializer: Normal $\sigma = 1.00$ - Epoch 0



- Normally from $\mathbf{N}(\mathbf{0}, \sigma^2)$
- Variance *matters*

Activation: tanh - Initializer: Normal $\sigma = 0.10$ - Epoch 0



- Normally from $\mathcal{N}(\mathbf{0}, \sigma^2)$
- Variance *matters*

- Forward

$$\sigma(W) = \sqrt{\frac{1}{\text{fan_in}}} \quad (\text{normal}) \quad \sigma(W) = \sqrt{\frac{3}{\text{fan_in}}} \quad (\text{uniform})$$

- Backward

$$\sigma(W) = \sqrt{\frac{1}{\text{fan_out}}} \quad (\text{normal}) \quad \sigma(W) = \sqrt{\frac{3}{\text{fan_out}}} \quad (\text{uniform})$$

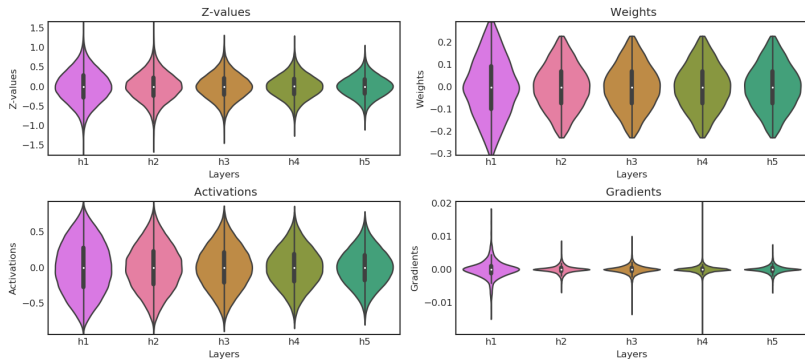
- Combined:

$$\sigma(W) = \sqrt{\frac{2}{\text{fan_in} + \text{fan_out}}} \quad (\text{normal})$$

$$\sigma(W) = \sqrt{\frac{6}{\text{fan_in} + \text{fan_out}}} \quad (\text{uniform})$$

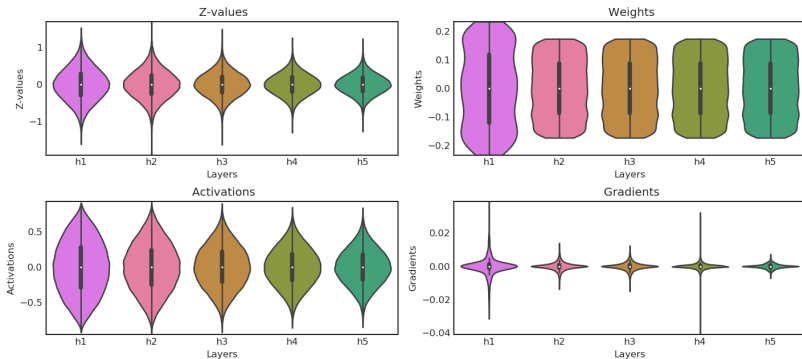
³<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

Activation: tanh - Initializer: Glorot Normal - Epoch 0



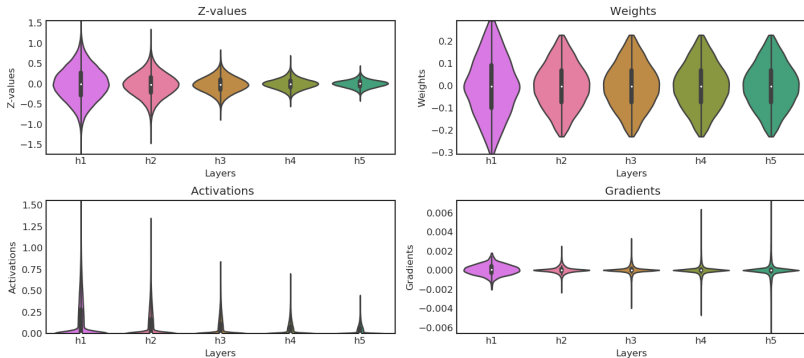
- Normal distribution

Activation: tanh - Initializer: Glorot Uniform - Epoch 0



- Uniform distribution

Activation: relu - Initializer: Glorot Normal - Epoch 0



- Does it work for ReLU?

- Forward

$$\sigma(W) = \sqrt{\frac{2}{\text{fan_in}}} \quad (\text{normal}) \quad \sigma(W) = \sqrt{\frac{6}{\text{fan_in}}} \quad (\text{uniform})$$

- Backward

$$\sigma(W) = \sqrt{\frac{2}{\text{fan_out}}} \quad (\text{normal}) \quad \sigma(W) = \sqrt{\frac{6}{\text{fan_out}}} \quad (\text{uniform})$$

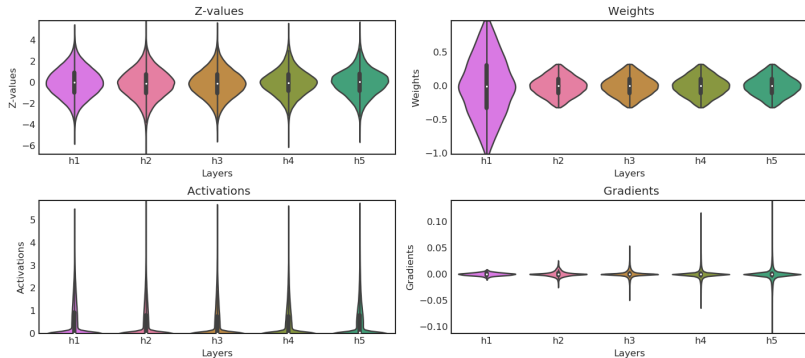
- Combined:

$$\sigma(W) = \sqrt{\frac{4}{\text{fan_in} + \text{fan_out}}} \quad (\text{normal})$$

$$\sigma(W) = \sqrt{\frac{12}{\text{fan_in} + \text{fan_out}}} \quad (\text{uniform})$$

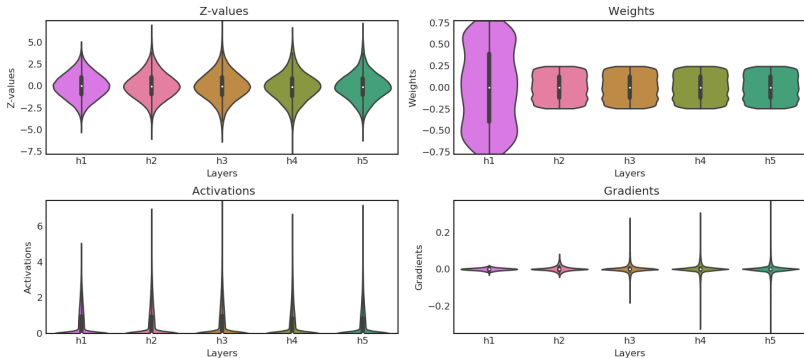
⁴<https://arxiv.org/pdf/1502.01852.pdf>

Activation: relu - Initializer: He Normal - Epoch 0



- Normal Distribution

Activation: relu - Initializer: He Uniform - Epoch 0



- Uniform Distribution

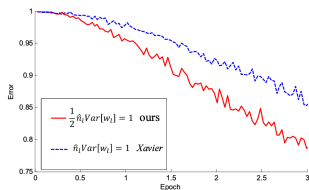


Figure 2. The convergence of a **22-layer** large model (B in Table 3). The x-axis is the number of training epochs. The y-axis is the top-1 error of 3,000 random val samples, evaluated on the center crop. We use ReLU as the activation for both cases. Both our initialization (red) and “Xavier” (blue) [7] lead to convergence, but ours starts reducing error earlier.

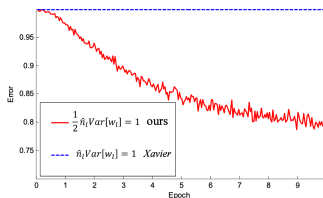


Figure 3. The convergence of a **30-layer** small model (see the main text). We use ReLU as the activation for both cases. Our initialization (red) is able to make it converge. But “Xavier” (blue) [7] completely stalls - we also verify that its gradients are all diminishing. It does not converge even given more epochs.

- Better than Glorot inits in their settings

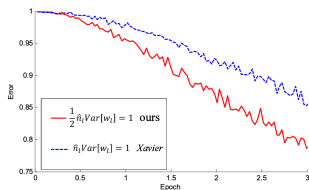


Figure 2. The convergence of a **22-layer** large model (B in Table 3). The x-axis is the number of training epochs. The y-axis is the top-1 error of 3,000 random val samples, evaluated on the center crop. We use ReLU as the activation for both cases. Both our initialization (red) and “Xavier” (blue) [7] lead to convergence, but ours starts reducing error earlier.

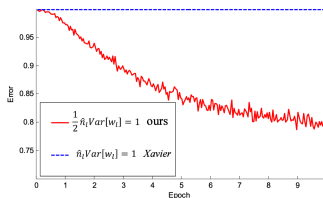
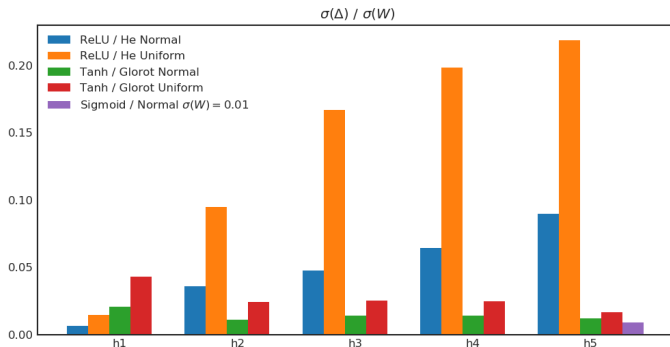


Figure 3. The convergence of a **30-layer** small model (see the main text). We use ReLU as the activation for both cases. Our initialization (red) is able to make it converge. But “Xavier” (blue) [7] completely stalls - we also verify that its gradients are all diminishing. It does not converge even given more epochs.

- Better than Glorot inits in their settings
- Many in the community agree, too!



- It varies
- Glorot works well with sigmoid or tanh
- He works well with ReLU and is usually used in vision

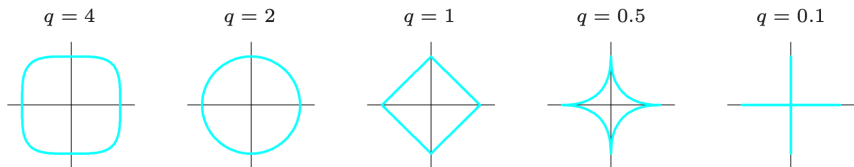
Image Credit: Daniel Godoy

tf.keras.layers.Conv2D example

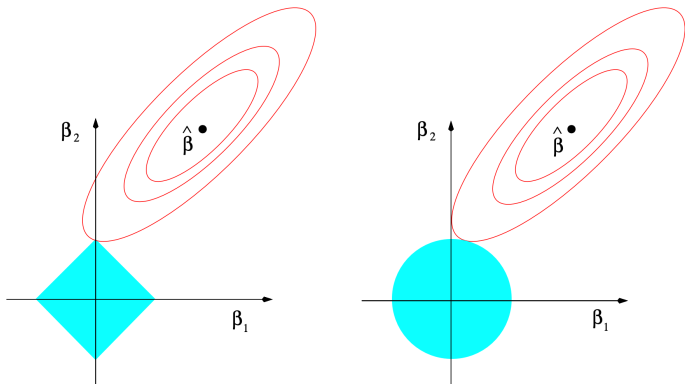
```
1  __init__(
2      filters,
3      kernel_size,
4      strides=(1, 1),
5      padding='valid',
6      data_format=None,
7      dilation_rate=(1, 1),
8      activation=None,
9      use_bias=True,
10     kernel_initializer='glorot_uniform',
11     bias_initializer='zeros',
12     kernel_regularizer=None,
13     bias_regularizer=None,
14     activity_regularizer=None,
15     kernel_constraint=None,
16     bias_constraint=None,
17     **kwargs
18 )
```

- API: `tf.keras.initializers`

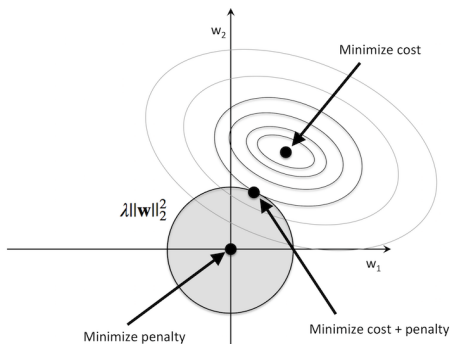
- 1 Initialization
- 2 Weights regularization
- 3 Dropout



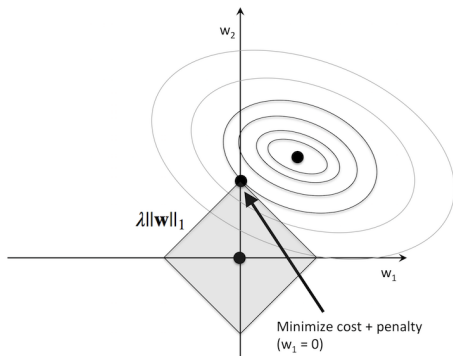
- Is common in practice which helps generalization
- Inject weights constraint and optimize the loss within that constraint
- Incorporate domain prior knowledge (recall: MAP)
- Can penalize based on L_q -norm of the weights: $\|\theta\|_q$



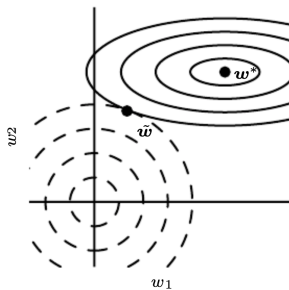
- L1 and L2 are simple, classical but usually works well in practice



- L2: Loss change: $\frac{\lambda}{2} \|\theta\|_2^2$, Gradients change: $\eta \lambda \theta$
- Analytically sound, and so very popular in practice



- L1: Loss change: $\frac{\lambda}{2} \|\theta\|_1$, (sub)Gradients change: $\eta \lambda \text{sign}(\theta)$
- Makes weights sparse, *unlike* L2



- Weight decaying: $\theta^{(t+1)} \leftarrow \alpha\theta^{(t)} - \eta\nabla_{\theta}l(\theta^{(t)})$ with small α
- L2 regularization: $\tilde{l}(\theta^{(t)}) = l(\theta^{(t)}) + \frac{\lambda}{2}\|\theta^{(t)}\|_2^2$
- They are the same for SGD but not for other adaptive methods
- AdamW⁵ implements weight decaying properly for Adam

⁵<https://arxiv.org/pdf/1711.05101.pdf>

tf.keras.regularizers example

```
1 # init object
2 my_l1 = tf.keras.regularizers.l1(l=0.01)
3 my_l2 = tf.keras.regularizers.l2(l=0.02)
4 my_l1l2 = tf.keras.regularizers.l1_l2(l1=0.01, l2=0.01) # L1 + L2 penalties
5
6 dense = tf.keras.layers.Dense(kernel_initializer='ones',
7                               kernel_regularizer=my_l1,
8                               bias_regularizer=my_l2, ... )
9
10 conv = tf.keras.layers.Conv2D(kernel_regularizer=my_l2,
11                               activity_regularizer=my_l1l2, ...)
```

- Normally, pass a `tf.keras.regularizers` object into a NN layer

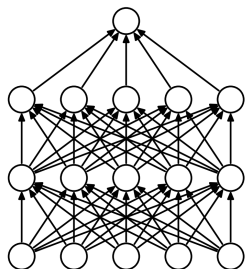
custom_regularizer

```
1 @tf.keras.utils.register_keras_serializable(package='Custom', name='l2')
2 class L2Regularizer(tf.keras.regularizers.Regularizer):
3     def __init__(self, l2=0.): # pylint: disable=redefined-outer-name
4         self.l2 = l2
5
6     def __call__(self, x):
7         return self.l2 * tf.math.reduce_sum(tf.math.square(x))
8
9     def get_config(self):
10        return {'l2': float(self.l2)}
11
12 dense = tf.keras.layers.Dense(kernel_initializer='ones',
13                               kernel_regularizer=L2Regularizer(l2=0.5),
14                               ... )
```

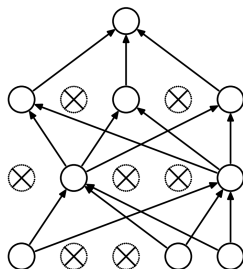
- Need to overwrite `__call__` function properly

⁶https://www.tensorflow.org/api_docs/python/tf/keras/regularizers/Regularizer

- 1 Initialization
- 2 Weights regularization
- 3 Dropout



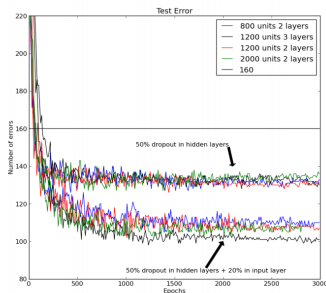
(a) Standard Neural Net



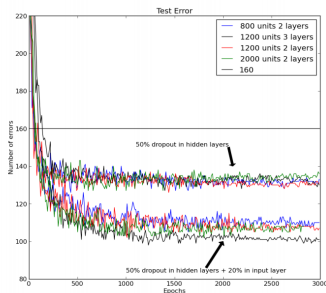
(b) After applying dropout.

- One of the most efficient regularization techniques
- Idea: Randomly drop neurons during training to reduce overfitting

⁷<http://www.cs.toronto.edu/~hinton/absps/dropout.pdf>

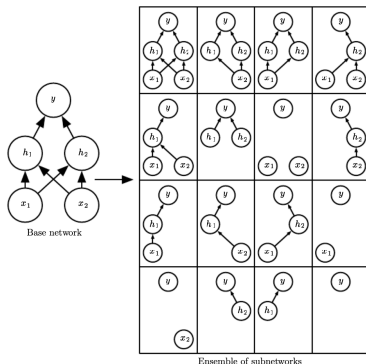


- Dropping probability: $\Lambda \sim Ber(p)$ (note: some use keeping prob)
- Training: each activation is changed as follow (both directions):
 - Origin: $h_i = \phi(W_i^T h_{i-1} + b_i)$
 - With dropout: $\Lambda_i \sim Ber(p)$, $h_i = \phi(W_i^T \Lambda_i h_{i-1} + b_i)$



- Dropping probability: $\Lambda \sim Ber(p)$ (note: some use keeping prob)
- Training: each activation is changed as follow (both directions):
 - Origin: $h_i = \phi(W_i^T h_{i-1} + b_i)$
 - With dropout: $\Lambda_i \sim Ber(p)$, $h_i = \phi(W_i^T \Lambda_i h_{i-1} + b_i)$
- Test: dropout is turned off, and *scaled* (reduced) by a factor of p

Image Credit: Geoffrey Hinton *et al.*



- Has ensembling effect
- Difference:
 - Exponentially subnetworks, only a small fraction gets trained
 - All subnetworks share different subsets of params

- Noise distribution (*e.g.* Bernoulli⁸ as above): $\lambda_i \sim p(\lambda)$
- Noise diagonal matrix: $\Lambda_i = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{|h_{i-1}|})$
- With dropout: $h_i = \phi\left(W_i^T \Lambda_i h_{i-1} + b_i\right)$
- Loglikelihood is changed to: Expectation of Loglikelihood:

$$\mathcal{L}(\mathbf{X}, \mathbf{y}, \{W_i\}_{i=1}^n) = \mathbb{E}_{p(\lambda)} \left[\log p\left(\mathbf{y} \mid \mathbf{X}, \{W_i\}_{i=1}^n, \{\Lambda_i\}_{i=1}^{n-1}\right) \right]$$

- The *only* assumption: Gaussian prior for all weights

$$W_i \sim \mathbf{N}(\mathbf{0}, \sigma^2)$$

⁸Others such as Gaussian, Half-Cauchy, ... also work well: <https://arxiv.org/pdf/1806.05975.pdf>

⁹<https://arxiv.org/pdf/1810.04045.pdf>

- Independent scalar random variable: $\lambda \sim p(\lambda)$
- A Gaussian variable with zero mean: $W \sim \mathbf{N}(\mathbf{0}, \sigma^2)$
- What is the distribution of λW ?

¹⁰<https://www.jstor.org/stable/2984774>

- Independent scalar random variable: $\lambda \sim p(\lambda)$
- A Gaussian variable with zero mean: $W \sim \mathbf{N}(\mathbf{0}, \sigma^2)$
- What is the distribution of λW ?

$$\theta \stackrel{d}{=} \lambda W \sim \mathbf{N}(\mathbf{0}, \lambda^2 \sigma^2)$$

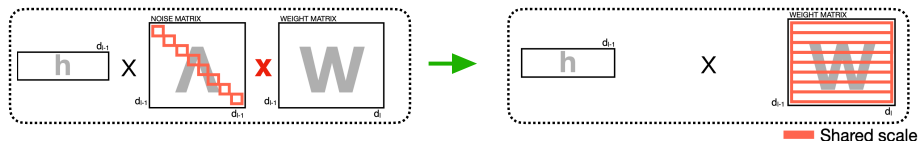
¹⁰<https://www.jstor.org/stable/2984774>

- Independent scalar random variable: $\lambda \sim p(\lambda)$
- A Gaussian variable with zero mean: $W \sim \mathbf{N}(\mathbf{0}, \sigma^2)$
- What is the distribution of λW ?

$$\theta \stackrel{d}{=} \lambda W \sim \mathbf{N}(\mathbf{0}, \lambda^2 \sigma^2)$$

- Super-Gaussian distributions can be represented as GSMs such as horseshoe, student-t, Laplace.

¹⁰<https://www.jstor.org/stable/2984774>



- Review the output with dropout:

$$\begin{aligned}
 h_i &= \phi \left(\underbrace{W_i^T \Lambda_i}_{\text{scale mixtures}} h_{i-1} + b_i \right) \\
 &= \phi \left(\widehat{W}_i^T h_{i-1} + b_i \right) \quad (\text{original type})
 \end{aligned}$$

with $\widehat{W}_{i,j} \sim \mathbf{N}(\mathbf{0}, \lambda_i^2 \sigma^2)$, and $\lambda_i \sim p(\lambda_i)$

- This interpretation concurs with Automatic Relevance Determination (ARD) by Mackay¹¹ and Neal¹²

¹¹ <https://bayes.wustl.edu/MacKay/pred.pdf>

¹² <http://www.cs.toronto.edu/pub/radford/thesis.pdf>