Aidan Gomez

In collaboration with: Ivan Zhang, Kevin Swersky, Yarin Gal, Geoff Hinton







Neural Networks

Briefly:

A neural network is a composition of functions (layers) similar to f(x) below.

$$f_j(\mathbf{x}) = \phi(\sum_i W_{j,i} \cdot x_i)$$



Targeted Dropout Sparsification

Immense Redundancy

- Neural networks are typically drastically overparameterized.
- Yet, this seems necessary to facilitate training.
- Ideally, we'd like to be able to remove as much of this redundancy as possible. How?





Identifying Important Subnetworks

- We want to remove portions of our networks; But, which portion?
- We need some measure for determining importance
- Judging weight importance may appear to be a complicated matter:
 - many subsets of our network to choose from (combinatorial explosion)
 - measuring subset importance can be extremely expensive

Idea:

Look at a Taylor expansion of the change in loss after deleting a subset "d" of the network

$$\left| \mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d}) \right| \approx \left| -\nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3) \right|$$

- Network parameter vector
- **d** Deletion vector ($\mathbf{d}_i = \Theta_i$ for deletion, 0s elsewhere)
- H Hessian of loss

LeCun et al. (1990)

Idea:

Look at a Taylor expansion of the change in loss after deleting a subset "d" of the network

$$\left| \mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d}) \right| \approx \left| -\nabla_{\Theta} \mathscr{L}^{\top} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\top} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3) \right|$$
vanishes



- **d** Deletion vector ($\mathbf{d}_i = \Theta_i$ for deletion, 0s elsewhere)
- H Hessian of loss

LeCun et al. (1990)

Idea:

Look at a Taylor expansion of the change in loss after deleting a subset "d" of the network

$$\left| \mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d}) \right| \approx \left| -\nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3) \right|$$



- Network parameter vector
- **d** Deletion vector ($\mathbf{d}_i = \Theta_i$ for deletion, 0s elsewhere)
- Hessian of loss

LeCun et al. (1990)

$$\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d}) \Big| \approx \Big| - \nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3)$$

$$\mathbf{H} = \begin{bmatrix} \frac{\partial \mathscr{L}}{\partial \Theta_1 \Theta_1} & \frac{\partial \mathscr{L}}{\partial \Theta_1 \Theta_2} & \cdots \\ \frac{\partial \mathscr{L}}{\partial \Theta_2 \Theta_1} & \frac{\partial \mathscr{L}}{\partial \Theta_2 \Theta_2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \quad \mathbf{d} = \begin{bmatrix} 0 & 0 & \Theta_3 & 0 & \Theta_5 & \cdots \end{bmatrix}$$

$$\left|\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d})\right| \approx \left| -\nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3) \right|$$

$$\left|\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d})\right| \approx \left| -\nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3) \right|$$

$$\begin{array}{c} \Theta_{2} & \Theta_{1} \\ \\ \Theta_{2} & \left[\frac{\partial \mathscr{L}}{\partial \Theta_{2} \Theta_{2}} & \frac{\partial \mathscr{L}}{\partial \Theta_{2} \Theta_{1}} & \cdots \right] \\ \\ \Theta_{1} & \left[\frac{\partial \mathscr{L}}{\partial \Theta_{1} \Theta_{2}} & \frac{\partial \mathscr{L}}{\partial \Theta_{1} \Theta_{1}} & \cdots \right] \\ \\ \vdots & \vdots & \ddots \end{array} \right]$$

 $\left| \mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d}) \right| \approx \left| -\nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3) \right|$

 $\mathbf{\Omega}$

 $\mathbf{\Omega}$

	01	Θ_2	03	04
Θ_1	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_4}$
Θ_2	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_4}$
Θ_3	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_4}$
Θ_4	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_4}$
				:

Ω

 $\mathbf{\Omega}$

$$|\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d})| \approx |-\nabla_{\Theta}\mathscr{L}^{\top} \cdot \mathbf{d} + 0.5 \,\mathbf{d}^{\top}\mathbf{H}\mathbf{d} + O(||\mathbf{d}||^3)$$

$\mathbf{d}^{\mathsf{T}}\mathbf{H}\mathbf{d}$

0 0	Θ ₃	Θ_4	
-----	----------------	------------	--

$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_4}$	0
$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_4}$	0
$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_4}$	Θ ₃
$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_1}$	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_2}$	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_3}$	$\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_4}$	Θ_4

$$|\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d})| \approx |-\nabla_{\Theta}\mathscr{L}^{\top} \cdot \mathbf{d} + 0.5 \,\mathbf{d}^{\top}\mathbf{H}\mathbf{d} + O(||\mathbf{d}||^3)$$

$\mathbf{b}^{\top}(\mathbf{d}^{\top}\odot\mathbf{H}\odot\mathbf{d})\mathbf{b}$

0 0 1	1
-------	---

•	$0\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_1} 0$	$0\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_2} 0$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_3} 0$	$\Theta_4 rac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_4} 0$	0
	$0\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_1} 0$	$0\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_2} 0$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_3} 0$	$\Theta_4 \frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_4} 0$	0
	$0\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_1} \Theta_3$	$0\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_2} \Theta_3$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_3} \Theta_3$	$\Theta_4 \frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_4} \Theta_3$	1
	$0\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_1} \Theta_4$	$0\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_2} \Theta_4$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_3} \Theta_4$	$\Theta_4 rac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_4} \Theta_4$	1

$$|\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d})| \approx |-\nabla_{\Theta}\mathscr{L}^{\top} \cdot \mathbf{d} + 0.5 \,\mathbf{d}^{\top}\mathbf{H}\mathbf{d} + O(||\mathbf{d}||^3)$$

$\mathbf{b}^{\top}(\mathbf{d}^{\top}\odot\mathbf{H}\odot\mathbf{d})\mathbf{b}$

0 0	1 1
-----	-----

•	$0\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_1} 0$	$0\frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_2} 0$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_3} 0$	$\Theta_4 rac{\partial \mathscr{L}}{\partial \Theta_1 \partial \Theta_4} 0$	0
	$0\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_1} 0$	$0\frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_2} 0$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_3} 0$	$\Theta_4 \frac{\partial \mathscr{L}}{\partial \Theta_2 \partial \Theta_4} 0$	0
	$0\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_1} \Theta_3$	$0\frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_2} \Theta_3$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_3} \Theta_3$	$\Theta_4 \frac{\partial \mathscr{L}}{\partial \Theta_3 \partial \Theta_4} \Theta_3$	1
	$0\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_1} \Theta_4$	$0\frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_2} \Theta_4$	$\Theta_3 \frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_3} \Theta_4$	$\Theta_4 \frac{\partial \mathscr{L}}{\partial \Theta_4 \partial \Theta_4} \Theta_4$	1

$$\mathscr{L}(\Theta) - \mathscr{L}(\Theta - \mathbf{d}) \Big| \approx \Big| - \nabla_{\Theta} \mathscr{L}^{\mathsf{T}} \cdot \mathbf{d} + 0.5 \, \mathbf{d}^{\mathsf{T}} \mathbf{H} \mathbf{d} + O(\|\mathbf{d}\|^3)$$

$\mathbf{b}^{\mathsf{T}}(\Theta^{\mathsf{T}}\odot\mathbf{H}\odot\Theta)\mathbf{b}$

0	0	1	1	
		1	1	



Visualizing Importance

- We can test our assumptions by inspecting the coadaptation matrix
- Position (*i*, *j*) of the matrix represents the dependence of parameter *i* on parameter *j*.
- Pruning a set of weights
 {*l*, *m*, *n*, ...} will change
 the loss proportional to the
 sum of the *l*, *m*, and *n*th
 columns



 $\left[\boldsymbol{\Theta}^{\top} \odot \mathbf{H} \odot \boldsymbol{\Theta}\right]_{i,j} = \boldsymbol{\Theta}_i \odot \mathbf{H}_{i,j} \odot \boldsymbol{\Theta}_j$

Visualizing Importance

- First 25% of weights are those with the largest magnitude.
- Remaining 75% are those we intend to prune.
- We hope that very little importance is placed on the 75% deemed unimportant.
- Ideally, all importance would be concentrated in the top-left 16th



 $\left[\boldsymbol{\Theta}^{\top} \odot \mathbf{H} \odot \boldsymbol{\Theta}\right]_{i,j} = \boldsymbol{\Theta}_i \odot \mathbf{H}_{i,j} \odot \boldsymbol{\Theta}_j$

Unregularized Model

- First 25% of weights are those with the largest magnitude.
- Remaining 75% are those we intend to prune.
- We hope that very little importance is placed on the 75% deemed unimportant.
- Ideally, all importance would be concentrated in the top-left 16th



$$\left[\Theta^{\top} \odot \mathbf{H} \odot \Theta\right]_{i,j} = \Theta_i \odot \mathbf{H}_{i,j} \odot \Theta_j$$

- It's clear that some simplifications are necessary:
 - Assume group/weight-wise independence
 - Assume a first order approx. is "good enough"
 - Assume a zeroth order approx. is "good enough"

- It's clear that some simplifications are necessary:
 - Assume group/weight-wise independence
 - Assume a first order approx. is "good enough"
 - Assume a zeroth order approx. is "good enough"

Saves computing the full Hessian, much less the **many** Hessianvector products originally required.

Now we need only compute a block-diagonal or diagonal approximation.

This technique is known as Optimal Brain Damage (LeCun et al., 1990)

- It's clear that some simplifications are necessary:
 - Assume group/weight-wise independence
 - Assume a first order approx. is "good enough"
 - Assume a zeroth order approx. is "good enough"

Despite the claim on a previous slide that $\nabla_{\Theta} \mathscr{L} \to 0$, Pavlo Molchanov and colleagues suggest that:

- yes, ∇_Θℒ · d^T is uninformative in expectation, but...
- the variance of this quantity correlates with the local stability of the loss, and
- we have access to this variance since $\mathbb{E}[|\nabla_{\Theta} \mathscr{L} \cdot \mathbf{d}^{\top}|] \propto \sigma$

Molchanov et al. (2017)

- It's clear that some simplifications are necessary:
 - Assume group/weight-wise independence
 - Assume a first order approx. is "good enough"
 - Assume a zeroth order approx. is "good enough"

Assume that one can judge the importance of a weight by its magnitude.

In practice this works quite well.

But far from perfect...

probably: Rumelhart (1988)







Idea: Inspired by Hinton's description of dropout reducing coadaptation between units, use dropout to reduce dependance of the identified important subnetwork on its complement in the network.



Apply dropout to <u>only</u> the smallest k%

Idea: Inspired by Hinton's description of dropout reducing coadaptation between units, use dropout to reduce dependance of the identified important subnetwork on its complement in the network.

After training delete those weights





```
def targeting fn(inputs, k):
  shape = tf.shape(inputs)
  size = tf.to int32(tf.reduce_prod(shape[:-1]))
  inputs = tf.reshape(inputs, [size, shape[-1]])
 transpose = tf.transpose(inputs)
 thres = tf.contrib.framework.sort(tf.abs(transpose), axis=1)[:, k]
 mask = tf.to float(tf.abs(inputs) <= thres[None, :])</pre>
 return tf.reshape(mask, shape)
def targeted dropout(inputs, targ rate, keep prob, is training):
  dim = tf.to float(tf.shape(inputs)[-1])
 k = tf.round(dim * targ rate)
 if not is training and do prune:
    drop rate = 1. - keep prob
    k = tf.round(dim * targ rate * drop rate)
 mask = targeting fn(inputs, k)
 mask = tf.cast(mask, inputs.dtype)
 if is training:
    return inputs * (1 - mask) + tf.nn.dropout(inputs, keep prob) * mask
  else:
    return inputs * (1 - mask)
```

Unregularized Model

- First 25% of weights are those with the largest magnitude.
- Remaining 75% are those we intend to prune.
- We hope that very little importance is placed on the 75% deemed unimportant.
- Ideally, all importance would be concentrated in the top-left 16th



$$\left[\Theta^{\top} \odot \mathbf{H} \odot \Theta\right]_{i,j} = \Theta_i \odot \mathbf{H}_{i,j} \odot \Theta_j$$

- Targeted Dropout achieves precisely what we had hoped for:
 - decoupling the "important" sub-network from the "unimportant" one.



 $\left[\Theta^{\top} \odot \mathbf{H} \odot \Theta\right]_{i,i} = \Theta_i \odot \mathbf{H}_{i,j} \odot \Theta_j$

- Targeted Dropout achieves precisely what we had hoped for:
 - decoupling the "important" sub-network from the "unimportant" one.



$$\left[\boldsymbol{\Theta}^{\top} \odot \mathbf{H} \odot \boldsymbol{\Theta}\right]_{i,j} = \boldsymbol{\Theta}_i \odot \mathbf{H}_{i,j} \odot \boldsymbol{\Theta}_j$$

Weight-level Results

		none	dropout $\alpha = 0.25$	targeted $\alpha = 0.66, \gamma = 0.75$	targeted $\alpha = 0.5, \gamma = 0.5$
	0 %	94.30	94.11	93.87	94.38
ge	10%	94.16	94.05	93.85	94.31
nta	20%	94.19	93.97	93.84	94.27
cel	30%	94.21	93.90	93.89	94.27
Der	40%	93.93	93.59	93.81	94.27
le I	50%	93.31	92.00	93.84	94.33
un	60%	91.50	88.23	93.89	93.92
Id	70%	82.89	58.07	93.84	92.19
	80%	38.35	10.66	92.31	74.31
	90%	12.76	9.92	46.57	16.67

 We baseline against three recent techniques: L1 regularisation, variational dropout, and smallify (to come). Results on CIFAR-10 using a ResNet-32.

Unit-level Results

		none	dropout $\alpha = 0.25$	targeted $\alpha = 0.66, \gamma = 0.75$	targeted $\alpha = 0.5, \gamma = 0.5$
	0 %	94.29	92.93	90.55	93.27
ge	10%	90.38	84.98	90.83	92.87
nta	20%	74.38	21.27	89.88	91.95
ce]	30%	36.49	10.32	87.35	89.84
per	40%	12.64	12.52	85.39	86.41
le]	50%	10.19	10.13	80.84	67.01
un	60%	11.32	9.95	71.97	11.70
Id	70%	10.14	9.98	55.98	9.98
	80%	10.01	9.91	10.02	9.94
	90%	10.12	9.95	10.07	9.95

 We baseline against three recent techniques: L1 regularisation, variational dropout, and smallify (to come). Results on CIFAR-10 using a ResNet-32.

Ramping Results

		targeted $\alpha = 0.66, \gamma = 0.75$		$\underset{alpha=0.99,\gamma=0.99}{\text{ramp targ}}$		$\underset{alpha=0.99,\gamma=0.99}{\text{ramp targ}}$
	0 %	93.87	90%	88.70	98.5%	88.74
ĝe	10%	93.85	91%	88.74	98.6%	88.82
nta	20%	93.84	92%	88.75	98.7%	88.79
ce]	30%	93.89	93%	88.75	98.8%	88.74
Der	40%	93.81	94%	88.80	98.9%	88.70
le I	50%	93.84	95%	88.73	99.0%	88.71
un	60%	93.89	96%	88.74	99.1%	87.63
Id	70%	93.84	97%	88.67	99.2%	67.25
	80%	92.31	98%	88.70	99.3%	51.86
	90%	46.57	99%	88.75	99.4%	11.59

• Idea: slowly ramp up the targeting proportion from 0% to some large percent (99% in our experiments)

Bitrot + Exprot Quantization



Single-precision floating-point (float32)



Truncated half-precision floating-point (bfloat16)



Single-precision floating-point (float32)



Truncated half-precision floating-point (bfloat16)



8-bit fixed-point (Q3.4)



Bitrot

Single-precision floating-point (float32)



Idea: Apply dropout to the bits of the mantissa in order to gradually anneal the network from a float32 to a bfloat16.

Bitrot

Single-precision floating-point (float32)



Apply Nested Dropout

Idea: Apply dropout to the bits of the mantissa in order to gradually anneal the network from a float32 to a bfloat16.

Nested Dropout



Nested Dropout



Apply Nested Dropout

Nested Dropout



Apply Nested Dropout

Borrowing from *ramping targeted dropout*, instead of sampling indices, we simply anneal upwards from zero over the course of training.

Single-precision floating-point (float32)



Truncated half-precision floating-point (bfloat16)



8-bit fixed-point (Q3.4)



Exprot

8-bit fixed-point (Q3.4)



Our Q3.4 can't represent any number with magnitude larger than $8 - 2^{-4}$ or smaller than 2^{-4} (but greater than zero)

Idea: We can gradually constrain the exponents of our numbers to be as close to 0 as possible simply by minimizing $|\log_2(W)|$ and stochastically perturbing the exponent towards zero.

Bitrot

```
def bitrot(inputs, targ bits, keep prob, is training):
  shifted = tf.bitwise.right shift(inputs, targ bits)
 rotten = tf.bitwise.left shift(shifted, targ bits)
 mask = tf.random uniform(tf.shape(inputs)) <= keep prob</pre>
 if is training:
    return inputs * mask + rotten * (1 - mask)
  else:
    return rotten
def exprot(inputs, exp shift, keep prob, is training):
  log2 = tf.log(inputs) / tf.log(2.)
  shift = -tf.sign(log2) * tf.minimum(exp shift, log2)
 rotten = inputs * 2**(shift)
 mask = tf.random uniform(tf.shape(inputs)) <= keep prob</pre>
 if is_training:
    return inputs * mask + rotten * (1 - mask), tf.abs(log2)
  else:
    return rotten
```

Bitrot + Exprot Results

	Acc.
Naive (float32)	93.48
Bitrot (4 mant.)	92.19
Bitrot (3 mant.)	87.37
Bitrot + Exprot (Q4.4)	93.36
Bitrot + Exprot (Q3.4)	~20.00

More to come!

More precisely, much more in next couple months. Still, go give it a try yourself!

Conclusion

- Early pruning
- Better sparsity support

Keep and eye out for the full version of the Targeted Dropout paper in the coming weeks!

Check out for.ai for opportunities to get involved in ML projects!