

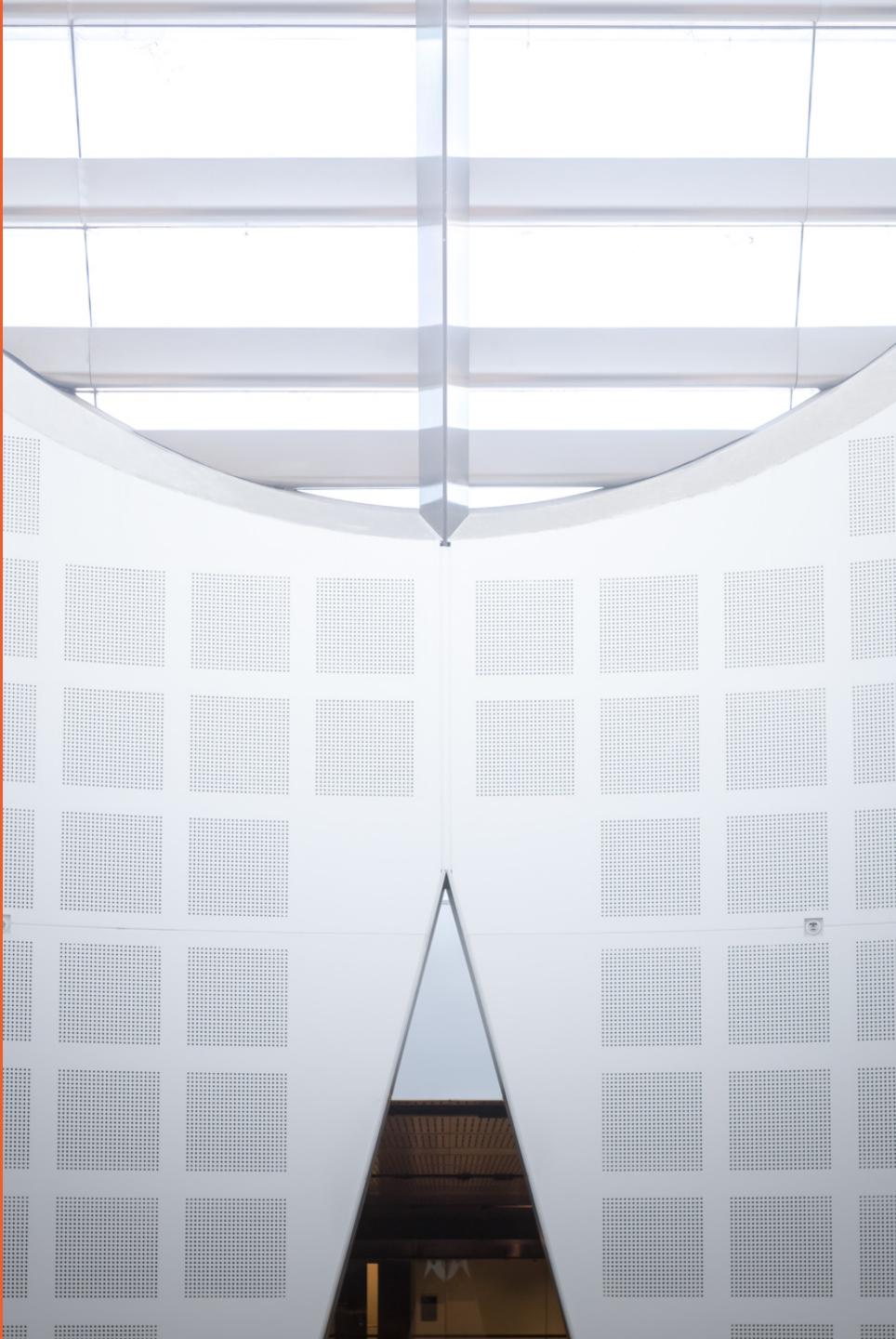
# Regularizations for Deep Models

Dr Chang Xu

School of Computer Science



THE UNIVERSITY OF  
SYDNEY



# Online Teaching

❖ No face-to-face teaching (lecture and lab).

❖ Lectures on Zoom

Join URL: <https://uni-sydney.zoom.us/j/200769779>

Lectures will be recorded and put on Canvas.

❖ Tutorials on Zoom

Find your lab at:

<https://canvas.sydney.edu.au/courses/21514/pages/unit-information>

# Before the quiz

- ❖ No lecture and lab on 27 Mar (Week 5).
- ❖ **All students** attend the ProctorU-based **online quiz**.

## What to Expect

- ✓ Be in a private, well-lit room
- ✓ Clear your workspace
- ✓ Have your photo I.D. ready
- ✓ Have a reflective surface handy
- ✗ Close all third-party programs
- ✗ No duplicate computer monitors
- ✗ No mobile devices
- ✓ No non-religious head coverings

<https://www.proctoru.com/proctoru-live-resource-center>



When you are **creating your account** and scheduling, if you need further help with scheduling, please visit the Live chat representatives, accessible from your ProctorU account homepage on the bottom right,

# Before the quiz

<https://canvas.sydney.edu.au/courses/23380/pages/preparing-for-your-online-test-one-hour-before-your-test>

[2020\\_UNIV\\_STUDENTS\\_ONLINE\\_TEST\\_SUPPORT\\_SITE](#) > Pages > Preparing for your online test (one hour before your test)

e  
les  
nments

## Preparing for your online test (one hour before your test)

[About your online test](#) | [Preparing for your online test](#) | [Get Help](#)

One week before the test

One day before the test

One hour before the test

### One hour before the test

Please read through the steps below.

**1. Clear work area and under desk.**

Except for allowable items as specified for your test (eg calculator). This will be checked by the proctor before your test begins.

**2. Have your identification document ready.**

You can use your passport, Sydney student card, or Australian government-issued identification.

**3. Remember to have your mirror ready.**

To perform a check of the monitor/laptop when instructed by your proctor. If you do not have a mirror you can use your mobile phone, but we recommend the former.

**4. Remind housemates or others not to interrupt you.**

And to avoid using the internet while you complete your exam!

**5. Go to the bathroom beforehand.**

If you need to go to the bathroom during the test you can, but the action of leaving the room will be recorded and reviewed.

**6. Restart your computer.**

To minimise technical issues during test.

**7. Close other software programs and other browser tab.**

And turn off pop-up notifications.

## In the quiz

- ❖ Name on ProctorU: **COMP 5329 Midterm**
- ❖ *If you need help*, tutors will be in their Zoom rooms.  
<https://canvas.sydney.edu.au/courses/21514/pages/unit-information>
- ❖ Effective Time:  
Fri 03/27/20 3:00 PM - Fri 03/27/20 5:00 PM (Sydney Time)  
**Please schedule your testing time in this window on ProctorU.**
- ❖ Durations: 60 minutes.  
Please visit the Live chat representatives, accessible from your ProctorU account homepage on the bottom right.

## In the quiz

- ❖ Closed-book quiz
- ❖ Calculator - non-programmable is allowed  
**No Computer's calculator.**
- ❖ You may need 1 sheet scratch paper
- ❖ No bathroom breaks.

General concepts of machine learning and deep learning

No essay question.

Multilayer Neural Network

Optimization, Regularization, NumPy and PyTorch

# **Quick Review**

# A common training process for neural networks

1. Initialize the parameters
2. Choose an optimization algorithm
3. Repeat these steps:
  1. Forward propagate an input
  2. Compute the cost function
  3. Compute the gradients of the cost with respect to parameters using backpropagation
  4. Update each parameter using the gradients, according to the optimization algorithm

# Momentum

The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions.

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

$$\theta_t = \theta_{t-1} - v_t$$

Momentum term  $\gamma$  is usually set to 0.9 or a similar value.

# Nesterov accelerated gradient (NAG)

First make a big jump in the direction of the previous accumulated gradient.

Then measure the gradient where you end up.

Finally accumulate the gradient and make a correction.

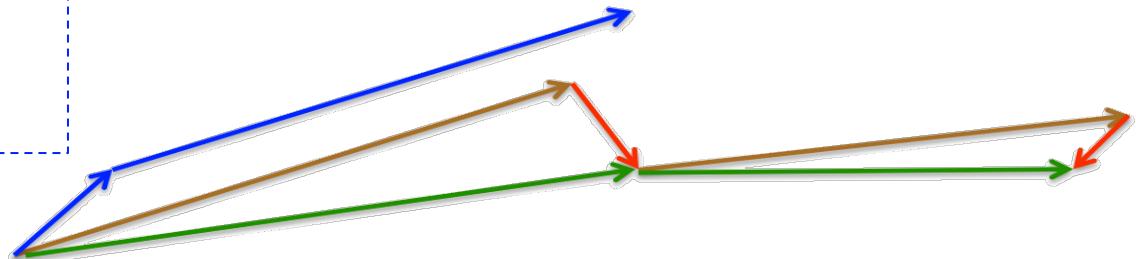
$$v_t = \gamma v_{t-1} + \theta - \theta - v_t$$

$$\eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$

Finally accumulate the gradient and make a correction

First make a big jump

Then measure the gradient where you end up



blue vectors = standard momentum

brown vector = jump, red vector = correction,  
green vector = accumulated gradient ,

- **Adagrad**  $\theta_t = [\theta_{t,1}, \theta_{t,2}, \dots, \theta_{t,i}, \dots, \theta_{t,d}]$

- **Original gradient descent:** The same for all parameters

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot \nabla_{\theta} J(\theta_i)$$

Perform an update for all parameters using the same learning rate.

- **Adagrad:** Adapt the learning rate to the parameters based on the past gradients that have been computed for  $\theta_i$ .

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot \nabla_{\theta} J(\theta_i)$$

$G_t \in R^{d \times d}$  is a diagonal matrix where each diagonal element is the sum of the squares of the gradients w.r.t.  $\theta_i$  up to time step t.

## Adadelta – 1 Accumulate over window

- Storing fixed previous squared gradients cannot accumulate to infinity and instead becomes a local estimate using recent gradients.
- Adadelta implements it as an exponentially **decaying average** of all past squared gradients.

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2$$

$$RMS[g]_t = \sqrt{E[g^2]_t + \epsilon}$$

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t$$

## Adadelta – 2 Correct units with Hessian approximation

$$H^{-1}g = \Delta\theta$$

$$\Delta\theta = \frac{\frac{\partial J}{\partial\theta}}{\frac{\partial^2 J}{\partial\theta^2}} \Rightarrow H^{-1} \propto \frac{1}{\frac{\partial^2 J}{\partial\theta^2}} \propto \frac{\Delta\theta}{\frac{\partial J}{\partial\theta}} \propto \frac{\Delta\theta}{g}$$

We assume the curvature is locally smooth and approximate  $\Delta\theta_t$  by computing the exponentially decaying RMS of  $\Delta\theta$ .

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t \xrightarrow{\text{red arrow}} \Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$$

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

# RMSprop

- RMSprop and Adadelta have both been developed independently around the same time.
- RMSprop in fact is identical to the first update vector of Adadelta.

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2$$

Hinton suggests  $\gamma$  to be set to 0.9, while a good default value for the learning rate  $\eta$  is 0.001.

# Adam

Combine the advantages of **Adgrad** and **RMSprop**.

- Adgrad: adaptive learning rate for different parameters.

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} m_t$$

- RMSprop: exponentially decaying average.

- Store an exponentially decaying average of **past squared gradients**  $v_t$  .

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

Keep an exponentially decaying average of **past gradients**  $m_t$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

# Initialization

All zero initialization:

Every neuron computes the same output.



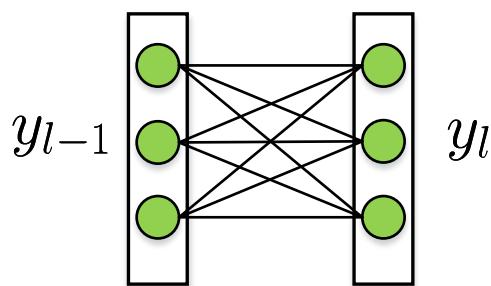
They will also compute the same gradients.



They will undergo the exact same parameter updates.



There will be no difference between all the neurons.



$$y_l = W y_{l-1} + b, \quad \frac{\partial \mathcal{L}}{\partial y_{l-1}} = W^\top \frac{\partial \mathcal{L}}{\partial y_l}$$

if  $W = 0$ , then  $\frac{\partial \mathcal{L}}{\partial y_{l-1}} = 0$

then,  $\frac{\partial \mathcal{L}}{\partial y_m} = 0, \quad m = 0, \dots, l-1$ .

Fills the input *Tensor* with values according to the method described in *Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification* - He, K. et al. (2015), using a uniform distribution. The resulting tensor will have values sampled from  $\mathcal{U}(-\text{bound}, \text{bound})$  where

$$\text{bound} = \boxed{\text{gain}} \times \sqrt{\frac{3}{\text{fan\_mode}}}$$

$$\mathcal{U}\left(-\sqrt{\frac{6}{n}}, \sqrt{\frac{6}{n}}\right)$$

Also known as He initialization.

### Parameters

- **tensor** – an n-dimensional `torch.Tensor`
- **a** – the negative slope of the rectifier used after this layer (only)
- **with 'leaky\_relu'** (used) –
- **mode** – either `'fan_in'` (default) or `'fan_out'`. Choosing `'fan_in'` preserves the magnitude of the weights in the forward pass. Choosing `'fan_out'` preserves the magnitude backwards pass.
- **nonlinearity** – the non-linear function (`nn.functional` name), recommended to use only with `'leaky_relu'` (default).

`torch.nn.init.calculate_gain(nonlinearity, param=None)`  
Return the recommended gain value for the given nonlinearity function. The values

nonlinearity	gain
Linear / Identity	1
Conv{1,2,3}D	1
Sigmoid	1
Tanh	$\frac{5}{3}$
ReLU	$\sqrt{2}$
Leaky Relu	$\sqrt{\frac{2}{1+\text{negative\_slope}^2}}$

### Examples

```
>>> w = torch.empty(3, 5)
>>> nn.init.kaiming_uniform_(w, mode='fan_in', nonlinearity='relu')
```

[https://pytorch.org/docs/stable/nn.init.html?highlight=kaiming\\_uniform#torch.nn.init.kaiming\\_uniform\\_](https://pytorch.org/docs/stable/nn.init.html?highlight=kaiming_uniform#torch.nn.init.kaiming_uniform_)

## How to find appropriate initialization values

1. The *mean* of the activations should be zero.
2. The *variance* of the activations should stay the same across every layer.

Under these two assumptions, the backpropagated gradient signal should not be multiplied by values too small or too large in any layer. It should travel to the input layer without exploding or vanishing.

## Initialization (Xavier) – [Xavier Glorot, Yoshua Bengio 2010]

We worked on activations computed during the forward propagation, and get

$$Var(W^{[l]}) = \frac{1}{n^{[l-1]}}$$

The same result can be derived for the backpropagated gradients:

$$Var(W^{[l]}) = \frac{1}{n^{[l]}}$$

Xavier initialization would either initialize the weights as

xavier\_normal

$$\mathcal{N}(0, \frac{2}{n^{[l-1]} + n^{[l]}})$$

xavier\_uniform

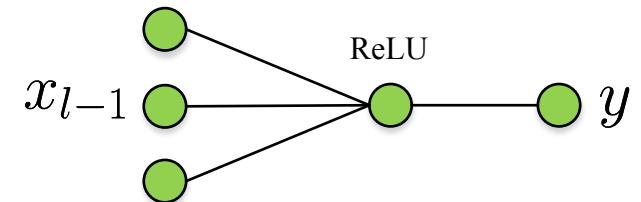
$$U\left(-\frac{\sqrt{6}}{\sqrt{n^{[l-1]}+n^{[l]}}}, \frac{\sqrt{6}}{\sqrt{n^{[l-1]}+n^{[l]}}}\right)$$

## Initialization (He) -- [Kaiming He, Xiangyu Zhang, et al. 2015]

If we use ReLU as the activation function,  $n$  is the number of inputs, then for layer  $l$  we have:

$$\text{Var}(y_l) = (n\text{Var}(w)) E[x_{l-1}^2]$$

$$x_{l-1} = \max(0, y_{l-1})$$



If  $y_{l-1}$  has zero mean, and has a symmetric distribution around zero:

$$E[x_{l-1}^2] = E[(\max(0, y_{l-1}))^2] = \frac{1}{2}\text{Var}[y_{l-1}]$$

So we have:  $\text{Var}(w) = \frac{2}{n}$

$$\text{Var}(x) = E(x^2) - (E(x))^2$$

`kaiming_normal`  $\mathcal{N}(0, \frac{2}{n})$

`kaiming_uniform`  $u(-\sqrt{\frac{6}{n}}, \sqrt{\frac{6}{n}})$

# Regularizations

# What is regularization?

In general: any method to prevent overfitting or help the optimization.

Regression using polynomials,

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \cdots + a_Mx^M + \epsilon$$

$$t = \sin(2\pi x) + \epsilon$$

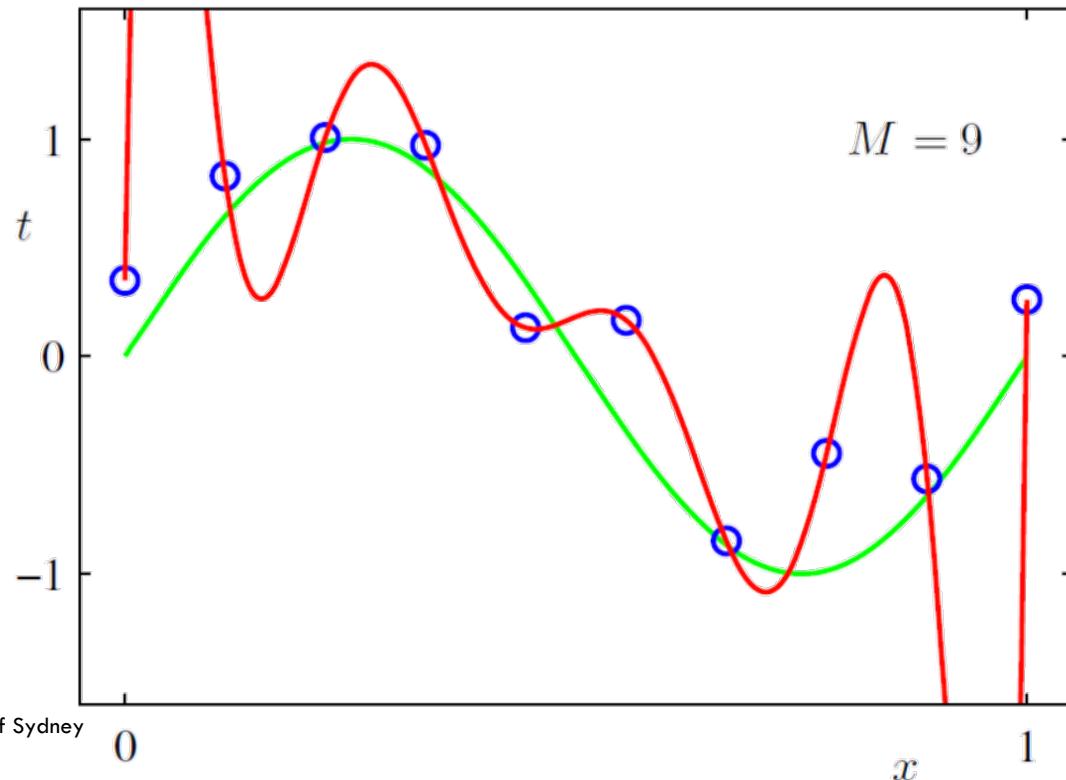


Figure from *Machine Learning and Pattern Recognition*, Bishop

# Overfitting

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \cdots + a_Mx^M + \epsilon$$

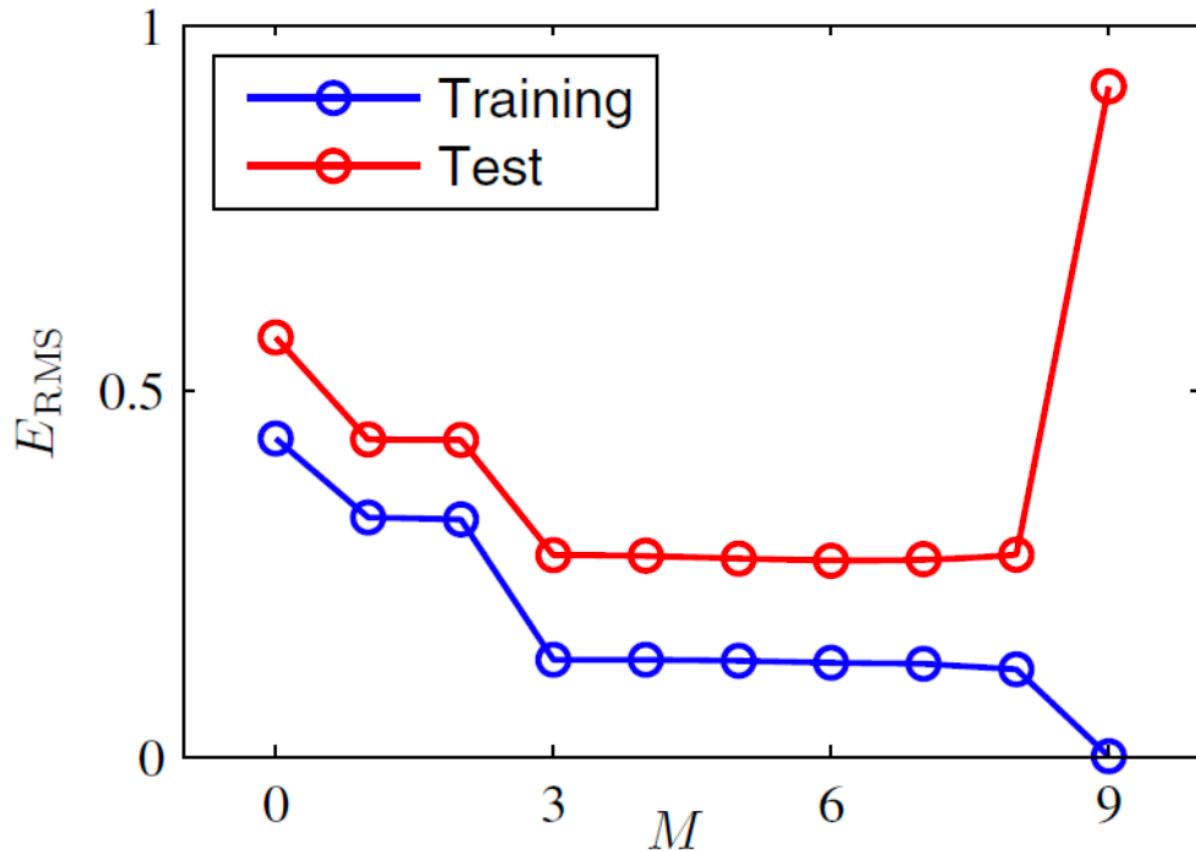


Figure from *Machine Learning and Pattern Recognition*, Bishop

## Prevent overfitting

- Larger data set helps.
  - Throwing away useless hypotheses also helps.
- 
- ✓ Classical regularization: some principal ways to constrain hypotheses.
  - ✓ Other types of regularization: data augmentation, early stopping, etc.

## Regularization as hard constraint

Training objective:

$$\min_{\theta} \hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta, x_i, y_i)$$
$$s.t. \quad R(\theta) \leq r$$

Example:  $\ell_2$  regularization

$$\min_{\theta} \hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta, x_i, y_i)$$
$$s.t. \quad \|\theta\|_2^2 \leq r$$

## Regularization as soft constraint

The hard-constraint optimization is equivalent to soft-constraint

$$\min_{\theta} \hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta, x_i, y_i) + \lambda^* R(\theta)$$

for some hyper parameter  $\lambda^* > 0$ .

Example:  $\ell_2$  regularization

$$\min_{\theta} \hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta, x_i, y_i) + \lambda^* \|\theta\|_2^2$$

# Regularization as Bayesian prior

- Bayesian view: everything is a distribution
- Prior over the hypotheses:  $p(\theta)$
- Posterior over the hypotheses:  $p(\theta|\{x_i, y_i\})$
- Likelihood:  $p(\{x_i, y_i\}|\theta)$
- Bayesian rule:

$$p(\theta|\{x_i, y_i\}) = \frac{p(\theta)p(\{x_i, y_i\}|\theta)}{p(\{x_i, y_i\})}$$

# Regularization as Bayesian prior

- Bayesian rule:

$$p(\theta | \{x_i, y_i\}) = \frac{p(\theta)p(\{x_i, y_i\}|\theta)}{p(\{x_i, y_i\})}$$

- Maximum A Posteriori (MAP):

$$\max_{\theta} \log p(\theta | \{x_i, y_i\}) = \max_{\theta} \underbrace{\log p(\theta)}_{\text{Regularization}} + \underbrace{\log p(\{x_i, y_i\}|\theta)}_{\text{MLE loss}}$$

# Some examples

# Weight Decay (Researchers' view)

- Limiting the growth of the weights in the network.
- A term is added to the original loss function, penalizing large weights:

$$\min_{\theta} \hat{L}_R(\theta) = \hat{L}(\theta) + \frac{\alpha}{2} \|\theta\|_2^2$$

- Gradient of regularized objective

$$\nabla \hat{L}_R(\theta) = \nabla \hat{L}(\theta) + \alpha \theta$$

- Gradient descent update

$$\theta \leftarrow \theta - \eta \nabla \hat{L}_R(\theta) = \theta - \eta \nabla \hat{L}(\theta) - \eta \alpha \theta = (1 - \eta \alpha) \theta - \eta \nabla \hat{L}(\theta)$$

# Weight Decay (Engineers' view)

- The L2 regularization technique for neural networks was worked out by researchers in the 1990s.
- At the same time, **engineers**, working independently from researchers, noticed that if you simply decrease the value of each weight on each training iteration, you get an improved trained model that isn't as likely to be overfitted.

$$\theta \leftarrow \theta - \eta \nabla \hat{L}_R(\theta) \quad // \text{ update}$$

$$\theta \leftarrow \theta * 0.98 \quad // \text{ or } \theta = \theta - (0.02 * \theta)$$

- The L2 approach has a solid underlying theory but is complicated to implement. The weight decay approach “just works” but is simple to implement.

[https://pytorch.org/docs/stable/\\_modules/torch/optim/adam.html#Adam](https://pytorch.org/docs/stable/_modules/torch/optim/adam.html#Adam)

```
import math
import torch
from .optimizer import Optimizer

[docs]class Adam(Optimizer):
    r"""Implements Adam algorithm.

    It has been proposed in 'Adam: A Method for Stochastic Optimization' .

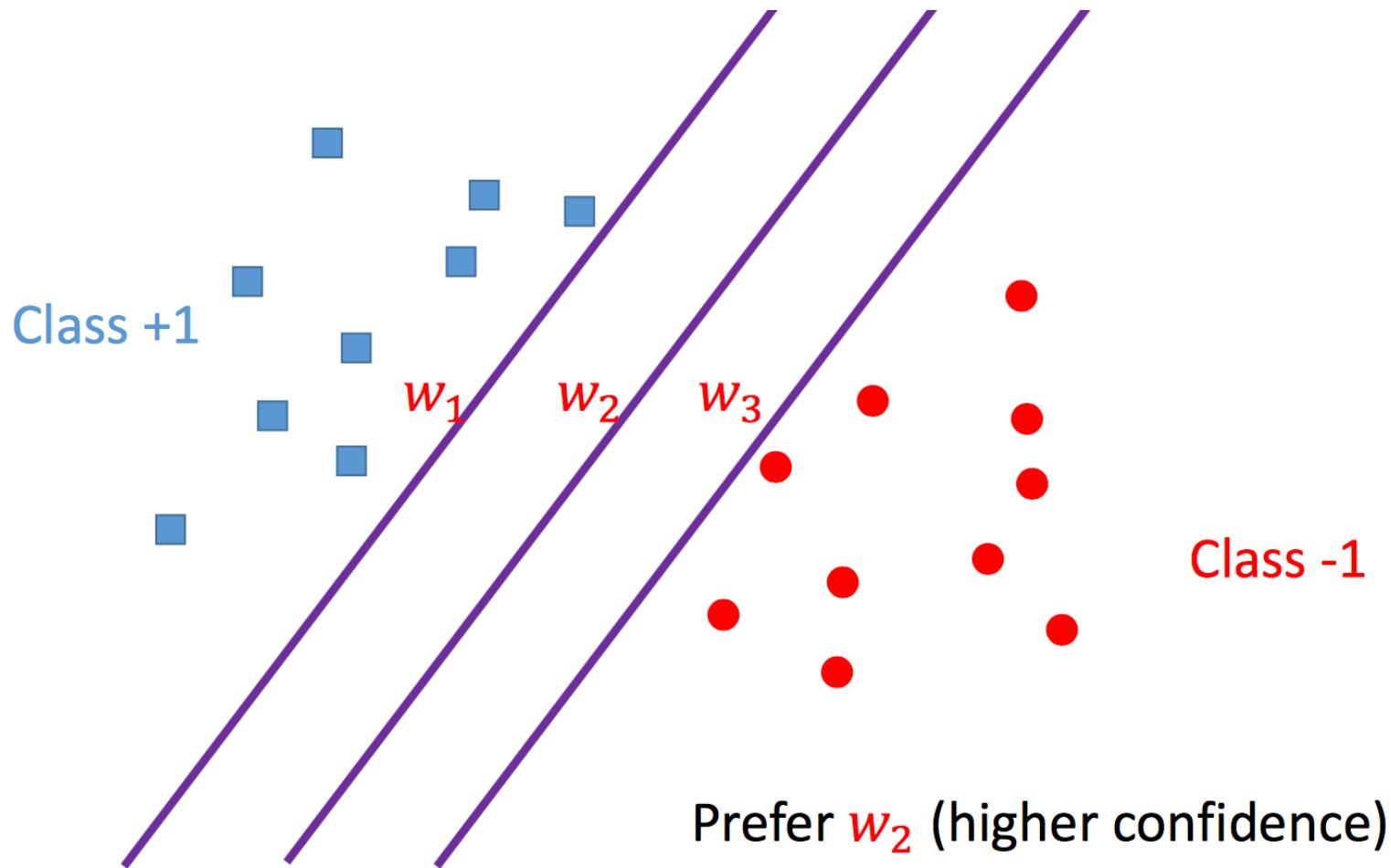
    Arguments:
        params (iterable): iterable of parameters to optimize or dicts defining
            parameter groups
        lr (float, optional): learning rate (default: 1e-3)
        betas (Tuple[float, float], optional): coefficients used for computing
            running averages of gradient and its square (default: (0.9, 0.999))
        eps (float, optional): term added to the denominator to improve
            numerical stability (default: 1e-8)
        weight_decay (float, optional): weight decay (L2 penalty) (default: 0)
        amsgrad (boolean, optional): whether to use the AMSGrad variant of this
            algorithm from the paper 'On the Convergence of Adam and Beyond'
            (default: False)

    .. _Adam\: A Method for Stochastic Optimization:
        https://arxiv.org/abs/1412.6980
    .. _On the Convergence of Adam and Beyond:
        https://openreview.net/forum?id=ryQu7f-RZ
    ...
```

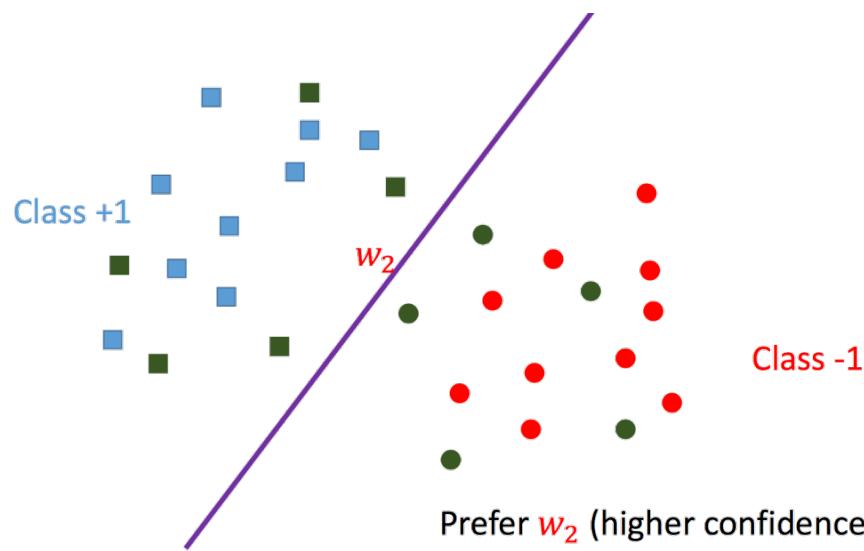
## Other types of regularizations

- Robustness to noise
  - Noise to the input
  - Noise to the weights
- Data augmentation
- Early stopping
- Dropout

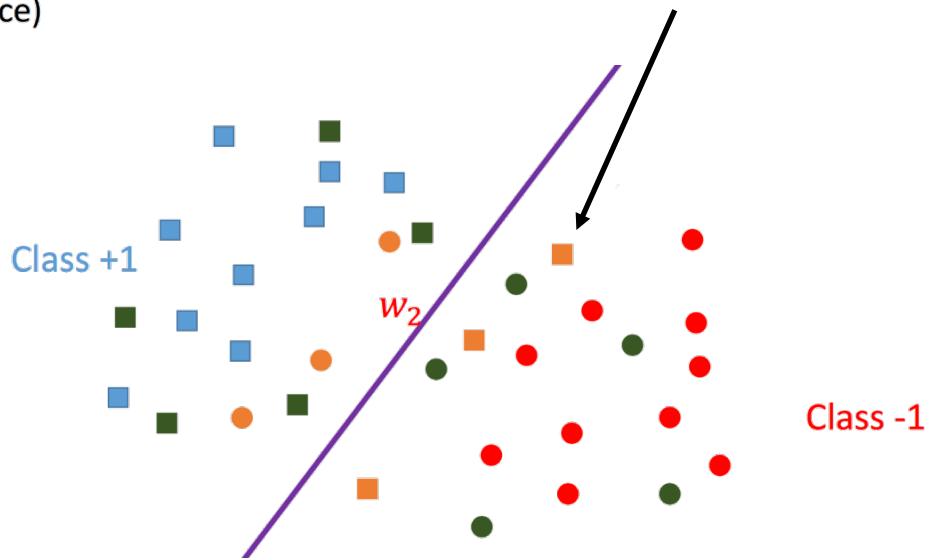
# Multiple optimal solutions.



## Add noise to the input



Too much noise leads to data points cross the boundary.



## Equivalence to weight decay

- Suppose the hypothesis is  $f(x) = w^T x$ , noise is  $\epsilon \sim N(0, \lambda I)$
- After adding noise, the loss is

$$\begin{aligned} L(f) &= \mathbb{E}_{x,y,\epsilon} [f(x + \epsilon) - y]^2 = \mathbb{E}_{x,y,\epsilon} [f(x) + w^T \epsilon - y]^2 \\ &= \mathbb{E}_{x,y} [f(x) - y]^2 + 2\mathbb{E}_{x,y,\epsilon} [w^T \epsilon (f(x) - y)] + \mathbb{E}_\epsilon [w^T \epsilon]^2 \\ L(f) &= \mathbb{E}_{x,y} [f(x) - y]^2 + \lambda \|w\|^2 \end{aligned}$$

## Add noise to the weights

- For the loss on each data point, add a noise term to the weights before calculating the prediction

$$\epsilon \sim N(0, \lambda I), \quad w' = w + \epsilon$$

- Prediction:  $f_{w'}(x)$  instead of  $f_w(x)$
- Loss becomes

$$L(f) = \mathbb{E}_{x,y,\epsilon} [f_{w+\epsilon}(x) - y]^2$$

## Add noise to the weights

- Loss becomes

$$L(f) = \mathbb{E}_{x,y,\epsilon} [f_{w+\epsilon}(x) - y]^2$$

- To simplify, use Taylor expansion

$$f_{w+\epsilon}(x) \approx f_w(x) + \epsilon^T \nabla f(x) + \frac{\epsilon^T \nabla^2 f(x) \epsilon}{2}$$

- Plug in

$$L(f) \approx \mathbb{E}_{x,y} [f_w(x) - y]^2 + \lambda \mathbb{E}_{x,y} \|\nabla f(x)\|^2 + \mathcal{O}(\lambda^2)$$

# Data augmentation

Adding noise to the input: a special kind of augmentation.

Horizontal Flip



Crop



Rotate



Figure from *Image Classification with Pyramid Representation and Rotated Data Augmentation on Torch 7*, by Keven Wang

The first very successful CNN on the ImageNet dataset.

---

# ImageNet Classification with Deep Convolutional Neural Networks

---

**Alex Krizhevsky**

University of Toronto

kriz@cs.utoronto.ca

**Ilya Sutskever**

University of Toronto

ilya@cs.utoronto.ca

**Geoffrey E. Hinton**

University of Toronto

hinton@cs.utoronto.ca

## Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

## 4 Reducing Overfitting

Our neural network architecture has 60 million parameters. Although the 1000 classes of ILSVRC make each training example impose 10 bits of constraint on the mapping from image to label, this turns out to be insufficient to learn so many parameters without considerable overfitting. Below, we describe the two primary ways in which we combat overfitting.

### 4.1 Data Augmentation

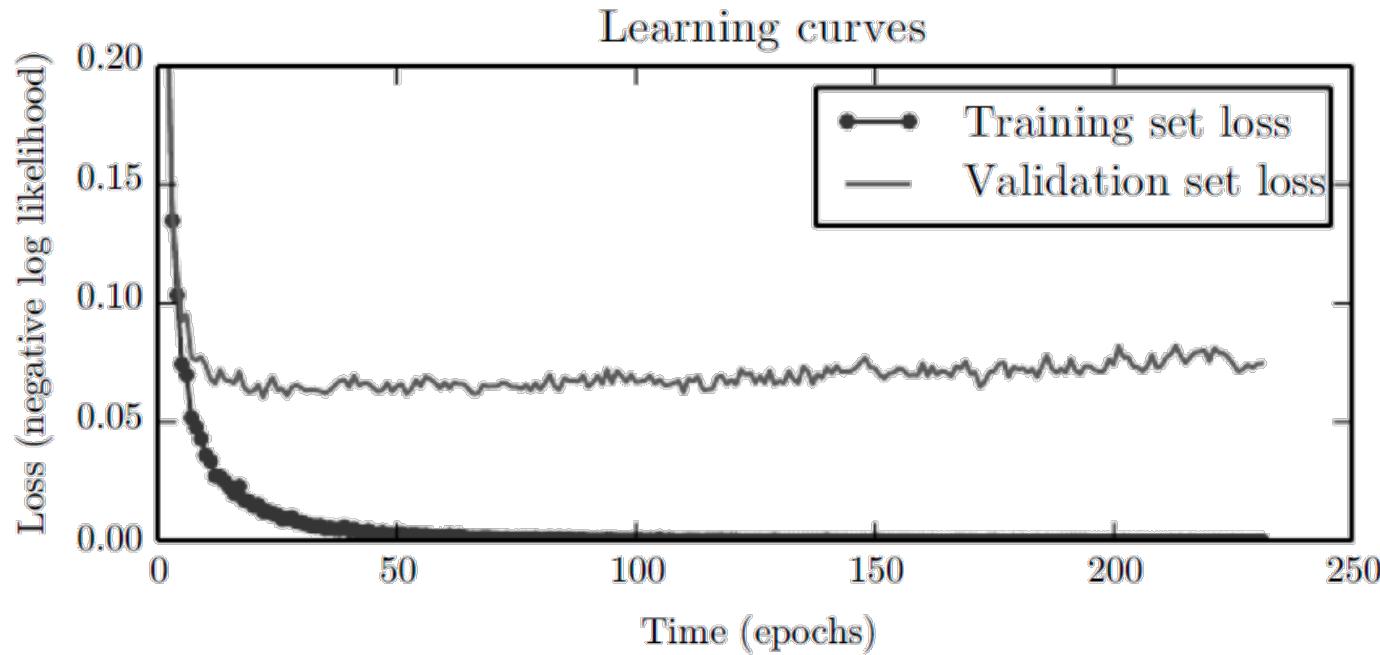
The easiest and most common method to reduce overfitting on image data is to artificially enlarge the dataset using label-preserving transformations (e.g., [25, 4, 5]). We employ two distinct forms of data augmentation, both of which allow transformed images to be produced from the original images with very little computation, so the transformed images do not need to be stored on disk. In our implementation, the transformed images are generated in Python code on the CPU while the GPU is training on the previous batch of images. So these data augmentation schemes are, in effect, computationally free.

### 4.2 Dropout

Combining the predictions of many different models is a very successful way to reduce test errors [1, 3], but it appears to be too expensive for big neural networks that already take several days to train. There is, however, a very efficient version of model combination that only costs about a factor of two during training. The recently-introduced technique, called “dropout” [10], consists of setting to zero the output of each hidden neuron with probability 0.5. The neurons which are “dropped out” in this way do not contribute to the forward pass and do not participate in back-

# Early stopping

- Idea: don't train the network to too small training error.

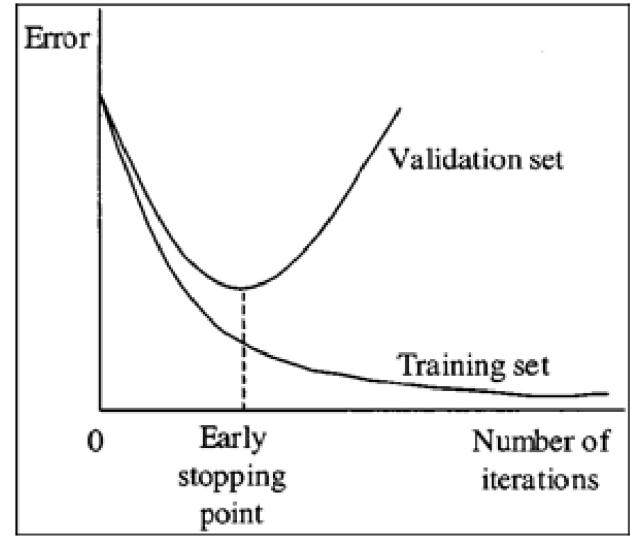


- When training, also output validation error
- Every time validation error improved, store a copy of the weights
- When validation error not improved for some time, stop
- Return the copy of the weights stored

# Early stopping

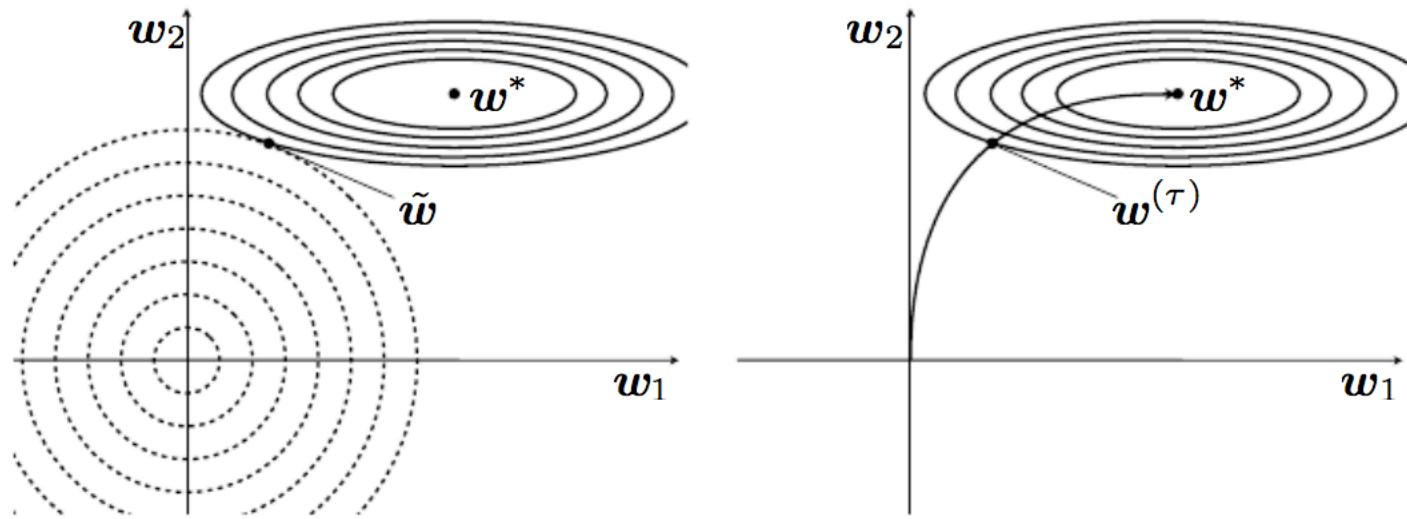
- hyperparameter selection: training step is the hyperparameter
- Advantage
  - Efficient: along with training; only store an extra copy of weights
  - Simple: no change to the model/algorithm
- Disadvantage: need validation data

# Early stopping



- Strategy to get rid of the disadvantage
  - After early stopping of the first run, train a second run and reuse validation data
- How to reuse validation data
  1. Start fresh, train with both training data and validation data up to the previous number of epochs
  2. Start from the weights in the first run, train with both training data and validation data until the validation loss < the training loss at the early stopping point

# Early stopping as a regularizer



# **Dropout**

# Dropout in PyTorch

<https://github.com/pytorch/pytorch/blob/10c4b98ade8349d841518d22f19a653a939e260c/torch/nn/modules/dropout.py#L36>

```
class Dropout(_DropoutNd):
    r"""During training, randomly zeroes some of the elements of the input
    tensor with probability :attr:`p` using samples from a Bernoulli
    distribution. Each channel will be zeroed out independently on every forward
    call.

    This has proven to be an effective technique for regularization and
    preventing the co-adaptation of neurons as described in the paper
    `Improving neural networks by preventing co-adaptation of feature
    detectors`_.

    Furthermore, the outputs are scaled by a factor of :math:`\frac{1}{1-p}` during
    training. This means that during evaluation the module simply computes an
    identity function.

    Args:
        p: probability of an element to be zeroed. Default: 0.5
        inplace: If set to ``True``, will do this operation in-place. Default: ``False``

    Shape:
        - Input: :math:`(*)`. Input can be of any shape
        - Output: :math:`(*)`. Output is of the same shape as input

    Examples::

        >>> m = nn.Dropout(p=0.2)
        >>> input = torch.randn(20, 16)
        >>> output = m(input)

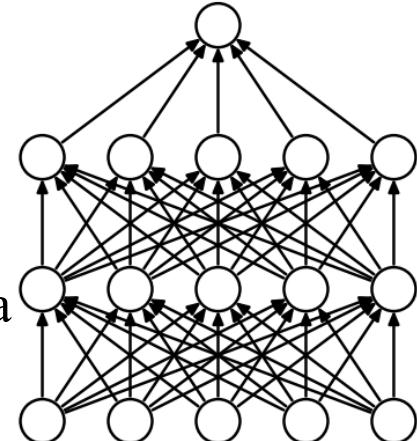
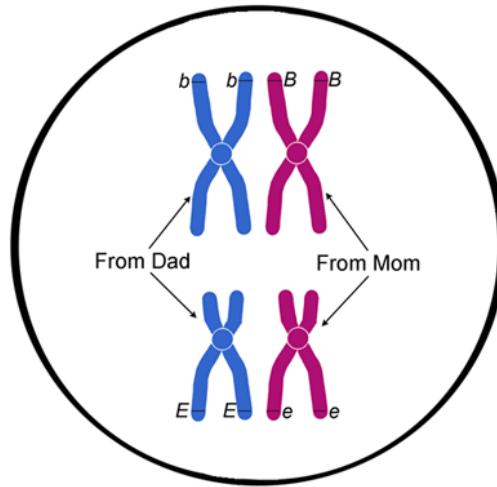
        .. _Improving neural networks by preventing co-adaptation of feature
           detectors: https://arxiv.org/abs/1207.0580
        ....
```

# Dropout

It is a simple but very effective technique that could alleviate overfitting in training phase.

The inspiration for Dropout (Hinton et al., 2013), came from the role of sex in evolution.

- Genes work well with another small random set of genes.
- Similarly, dropout suggests that each unit should work with a random sample of other units.



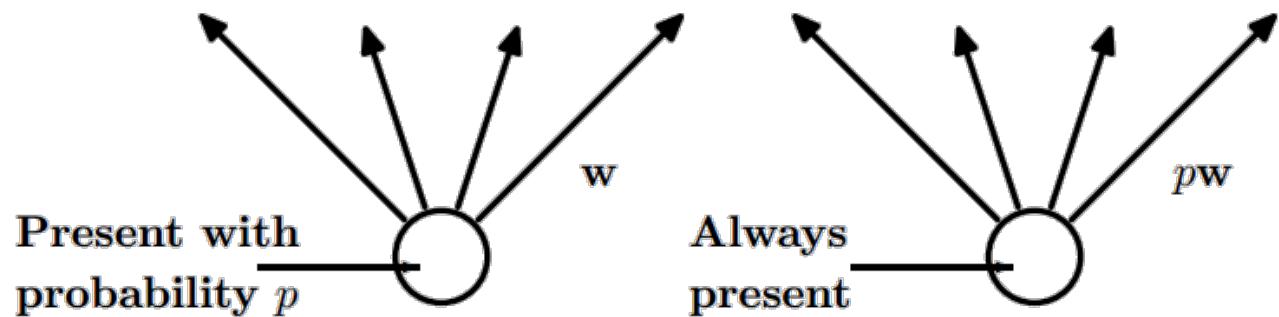
# Dropout

- At training (each iteration):

Each unit is retained with a probability  $p$ .

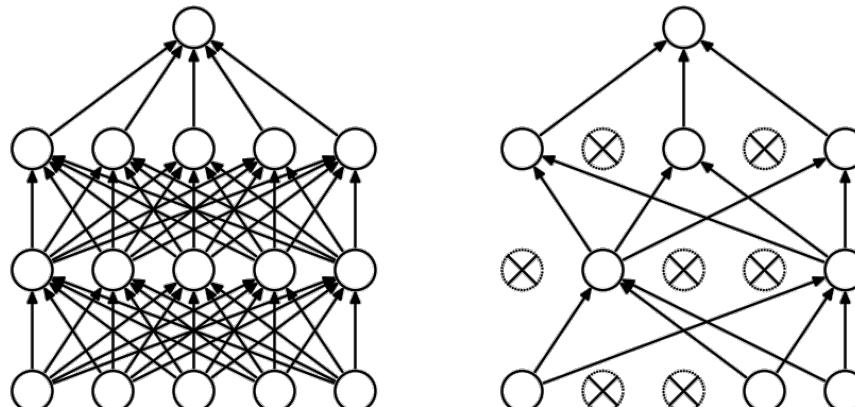
Expectation ?

- At test:



The network is used as a whole.

The weights are scaled-down by a factor of  $p$  (e.g. 0.5).

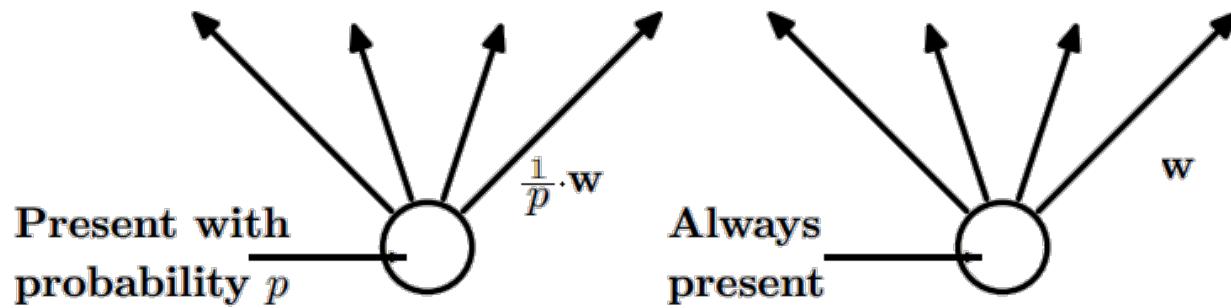


# Inverted Dropout

- At training (each iteration):

Each unit is retained with a probability  $p$ .

- ✓ Weights are scaled-up by a factor of  $1/p$ .



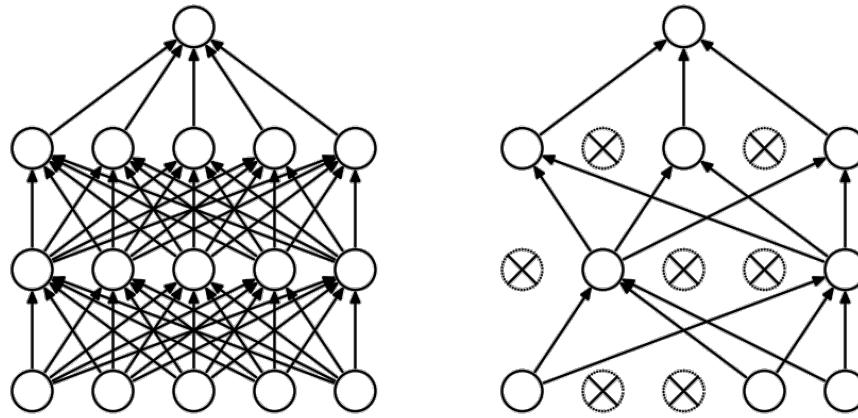
- At test:

The network is used as a whole.

- ✓ No scaling is applied.

# Dropout

In practice, dropout trains  $2^n$  networks ( $n$  – number of units).



Dropout is a technique which deals with overfitting by combining the predictions of many different large neural nets at test time.

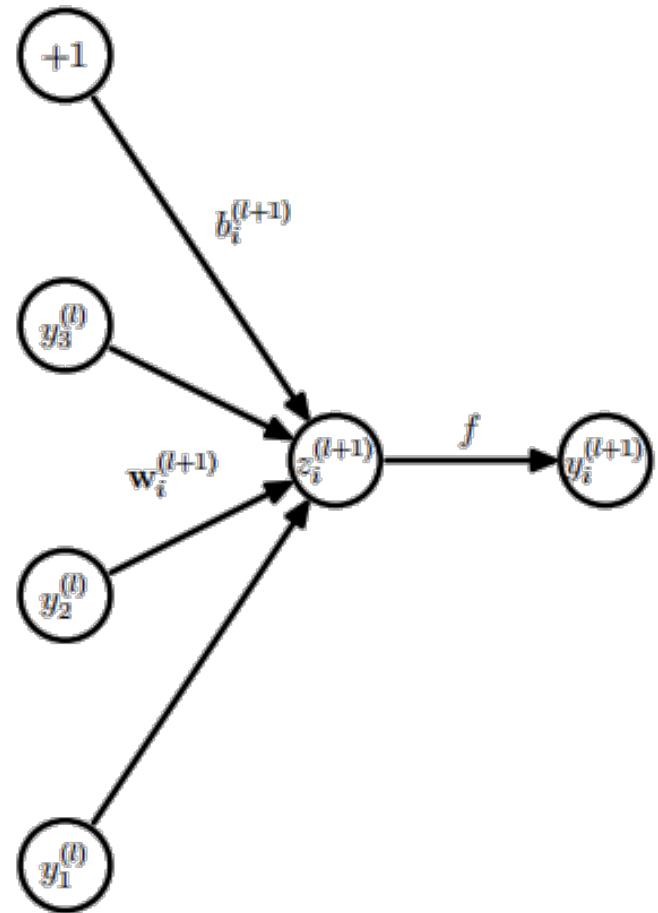
$$y^{l+1} = \sum_M p(M) f \left( (M * W)y^l \right) \approx f \left( \sum_M p(M)(M * W)y^l \right) = f(pW y^l)$$

# Dropout

The feed-forward operation of a standard neural network is

$$z_i^{(l+1)} = \mathbf{w}_i^{(l+1)} \mathbf{y}^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$

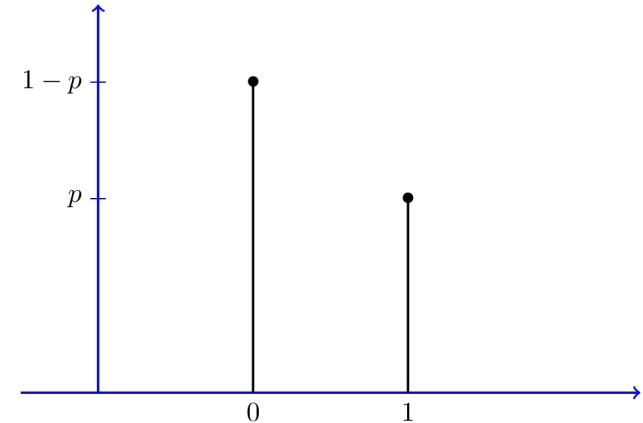


# Dropout

With dropout:

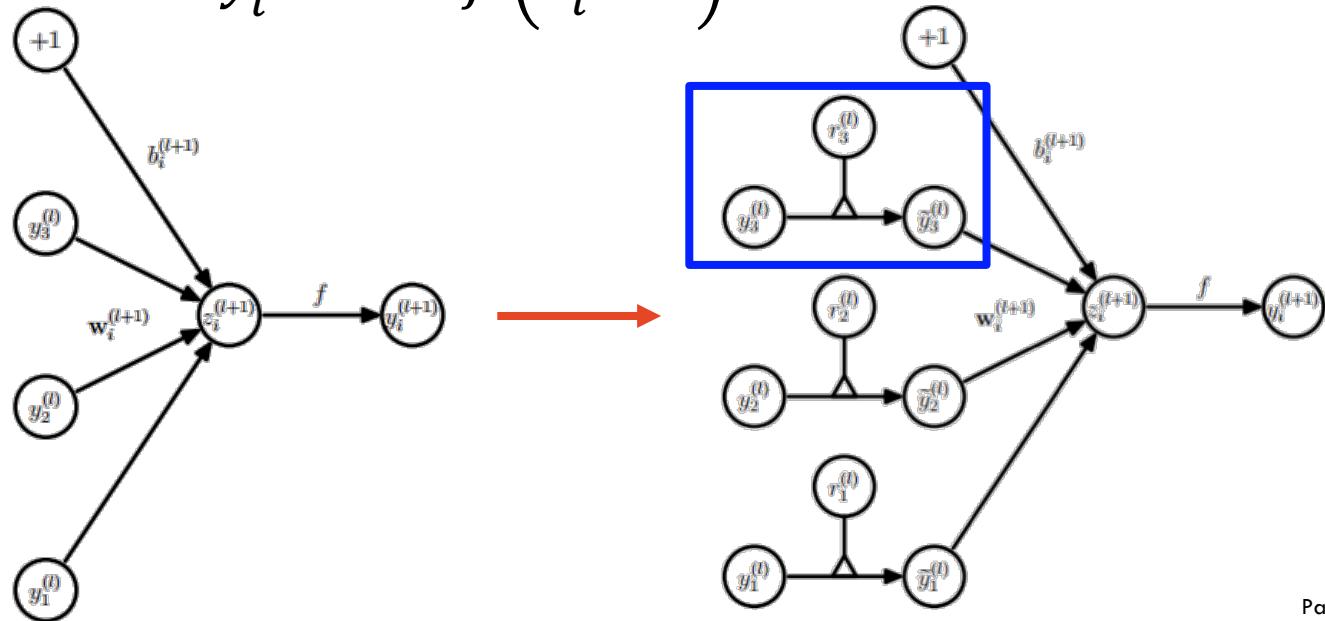
$$r_i^{(l)} \sim \text{Bernoulli}(p)$$

$$\tilde{\mathbf{y}}^{(l)} = \mathbf{r}^{(l)} * \mathbf{y}^{(l)}$$



$$z_i^{(l+1)} = \mathbf{w}_i^{(l+1)} \tilde{\mathbf{y}}^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$



# Weight Decay

- Limiting the growth of the weights in the network.
- A term is added to the original loss function, penalizing large weights:

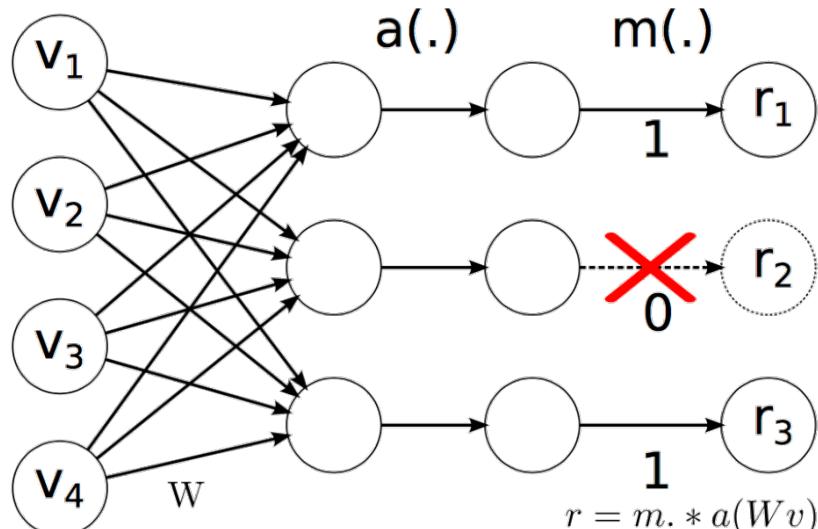
$$\mathcal{L}_{new}(\mathbf{w}) = \mathcal{L}_{old}(\mathbf{w}) + \frac{1}{2}\lambda\|\mathbf{w}\|_2^2$$

Dropout has more advantages over weight decay (Helmbold et al., 2016):

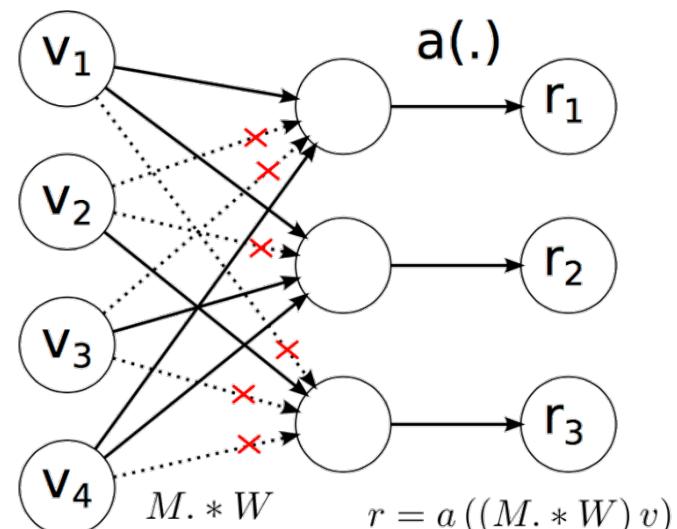
- Dropout is scale-free: dropout does not penalize the use of large weights when needed.
- Dropout is invariant to parameter scaling: dropout is unaffected if weights in a certain layer are scaled up by a constant  $c$ , and the weights in another layer are scaled down by a constant  $c$ .

# DropConnect

- DropConnect (Yann LeCun et al., 2013) generalizes dropout.
- Randomly drop connections in network with probability  $1 - p$ .
- As in Dropout,  $p = 0.5$  usually gives the best results.



DropOut Network



DropConnect Network

# DropConnect

A single neuron before activation function:

$$z = (M.* W)v$$

Approximate  $z$  by a **Gaussian distribution**:

$$E_M[z] = pWv$$

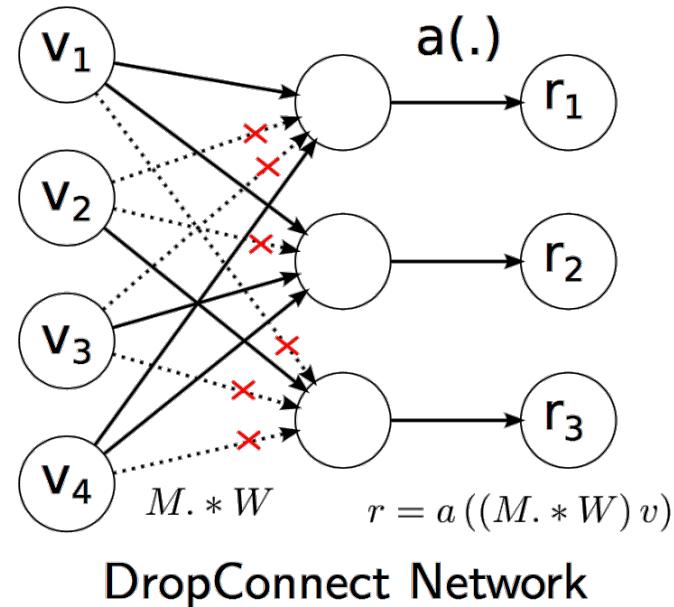
$$V_M[z] = p(1 - p)(W.* W)(v.* v)$$

At test:

Draw samples of  $z$  according to the Gaussian distribution.

Feed the samples into the activation function ( $f(z)$ ).

Average.



# **Batch Normalization**

## Feature scaling

- In machine learning algorithms, the functions involved in the optimization process are sensitive to normalization
  - For example: Distance between two points by the Euclidean distance. If one of the features has a broad range of values, the distance will be governed by this particular feature.
  - After, normalization, each feature contributes approximately proportionately to the final distance.
- In general, gradient descent converges much faster with feature scaling than without it.
- Good practice for numerical stability for numerical calculations, and to avoid ill-conditioning when solving systems of equations.

# Classical normalizations

Two methods are usually used for rescaling or normalizing data:

- Scaling data all numeric variables to the range [0,1]. One possible formula is given below:

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- To have zero mean and unit variance:

$$x_{new} = \frac{x - \mu}{\sigma}$$

- In the NN community this is call **Whitening**

# Whitening transformation



Fig.1: 10 raw sample images from each of 10 different classes of CIFAR10 namely airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck

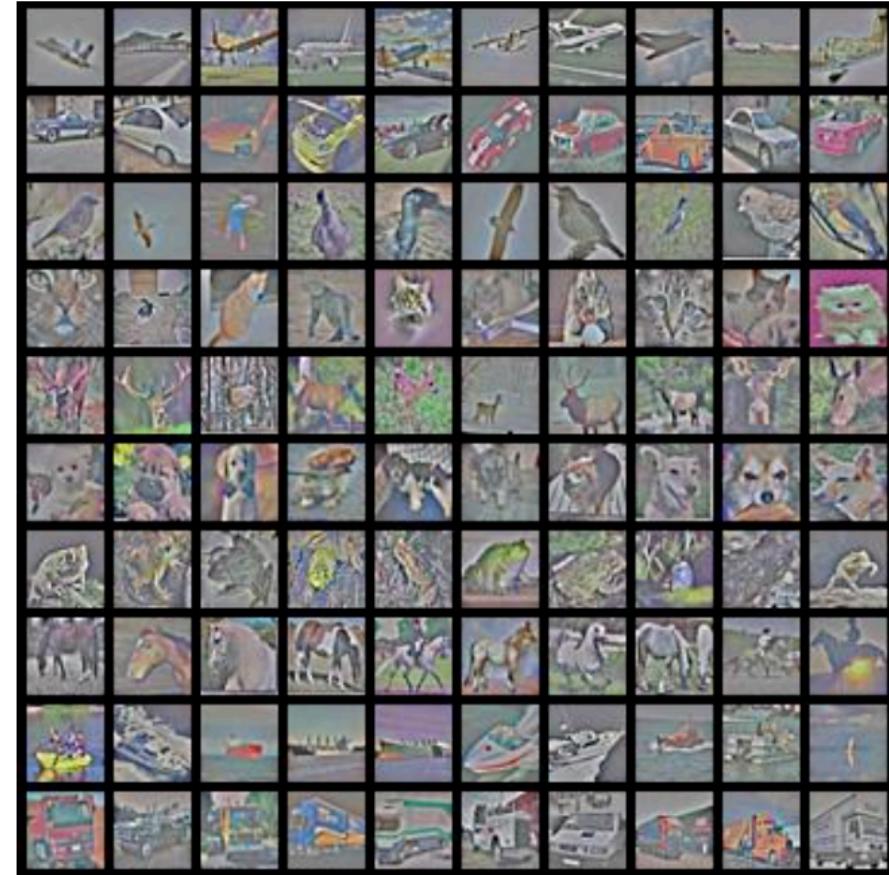
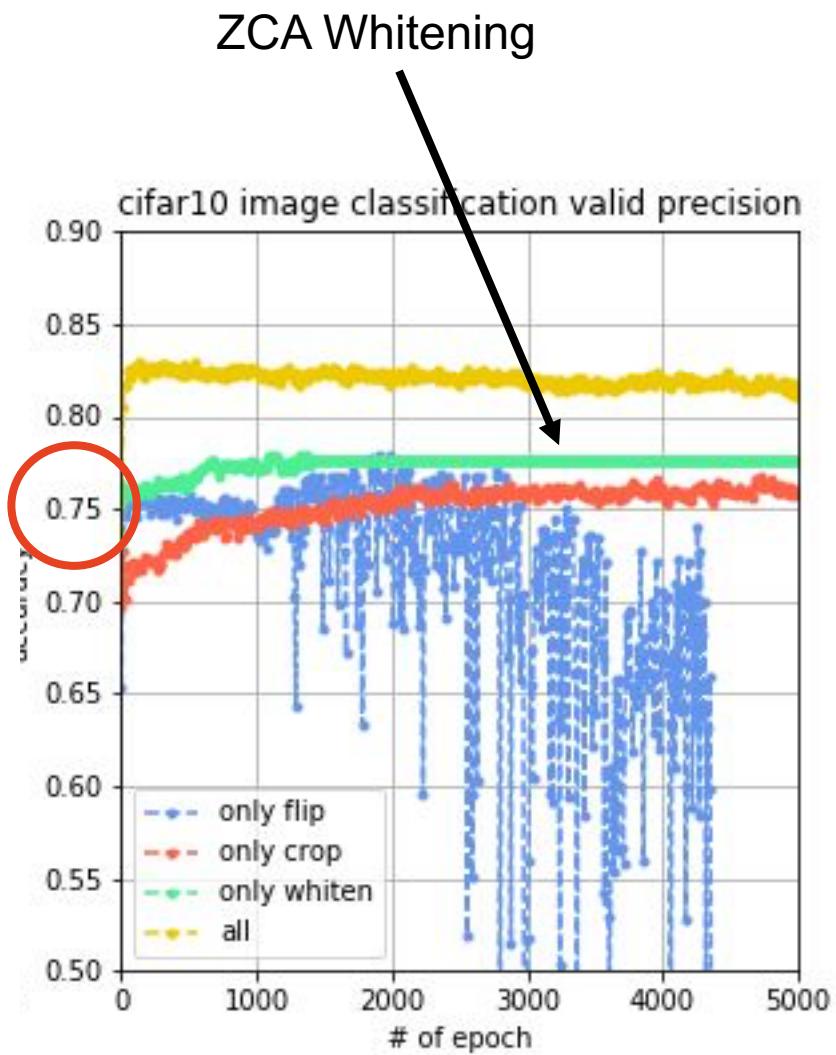
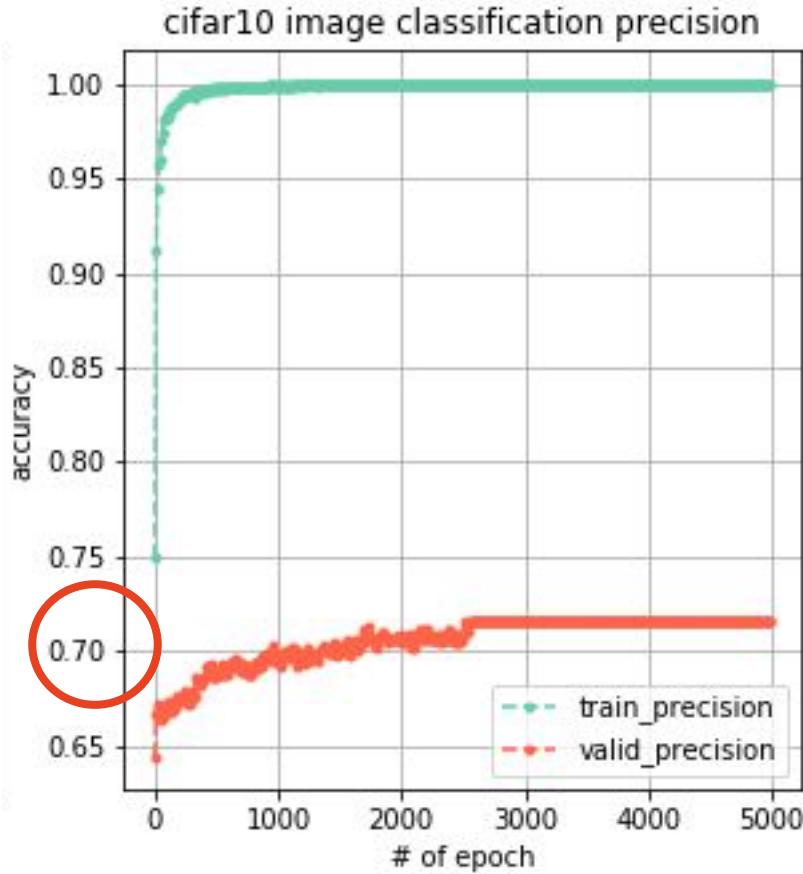


Fig.2:ZCA-whitened images of sample images of Fig. 1 with  $\epsilon=0.1$

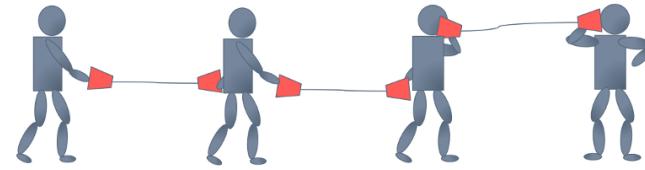
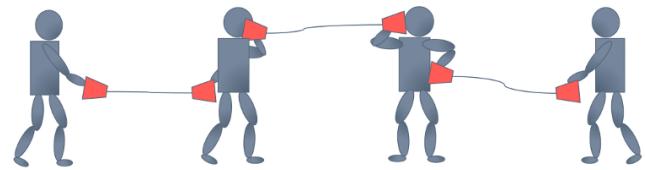
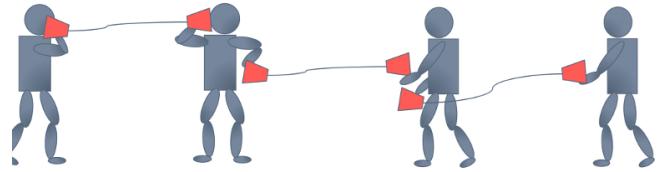
The ZCA transformation makes the edges of the objects more prominent.

# Validation Performance



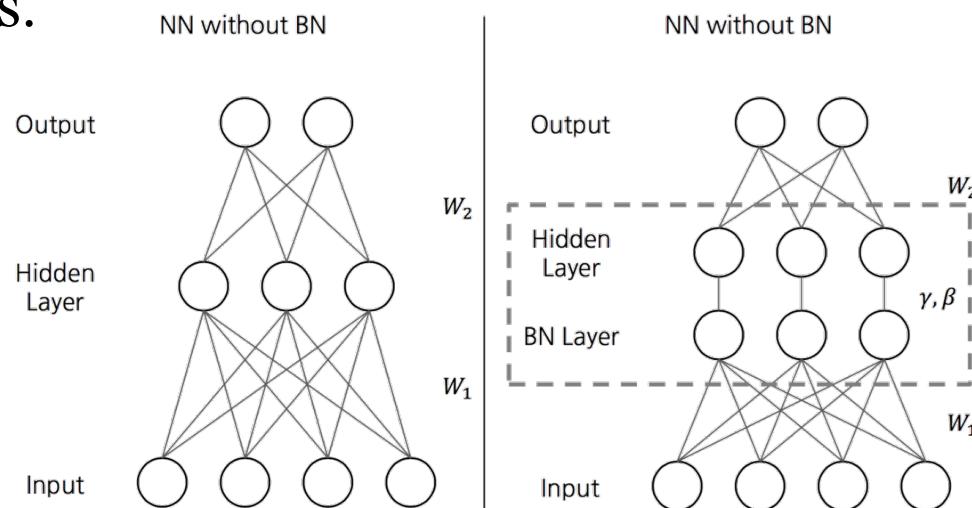
# Internal covariate shift:

- The cup game example
- The problems are entirely systemic and due entirely to faulty red cups.
- The situation is analogous to forward propagation
- First layer parameters change and so the distribution of the input to your second layer changes.



# Batch Normalization (BN)

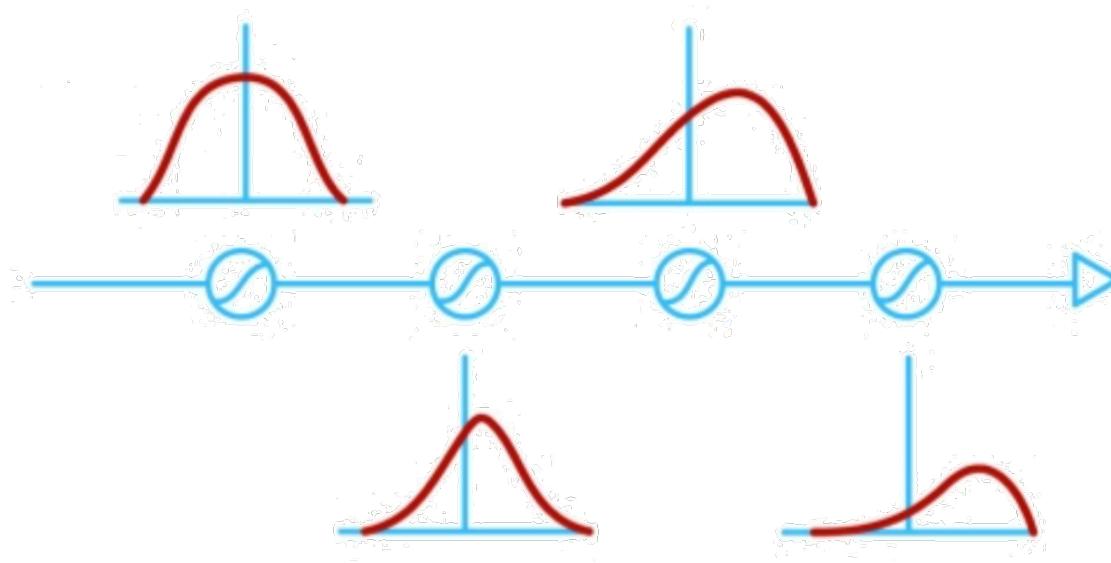
- Batch Normalization (BN) is a normalization method/layer for neural networks.
- Usually inputs to neural networks are normalized to either the range of [0, 1] or [-1, 1] or to mean=0 and variance=1
- BN essentially performs Whitening to the intermediate layers of the networks.



# Batch normalization

## Why it is good?

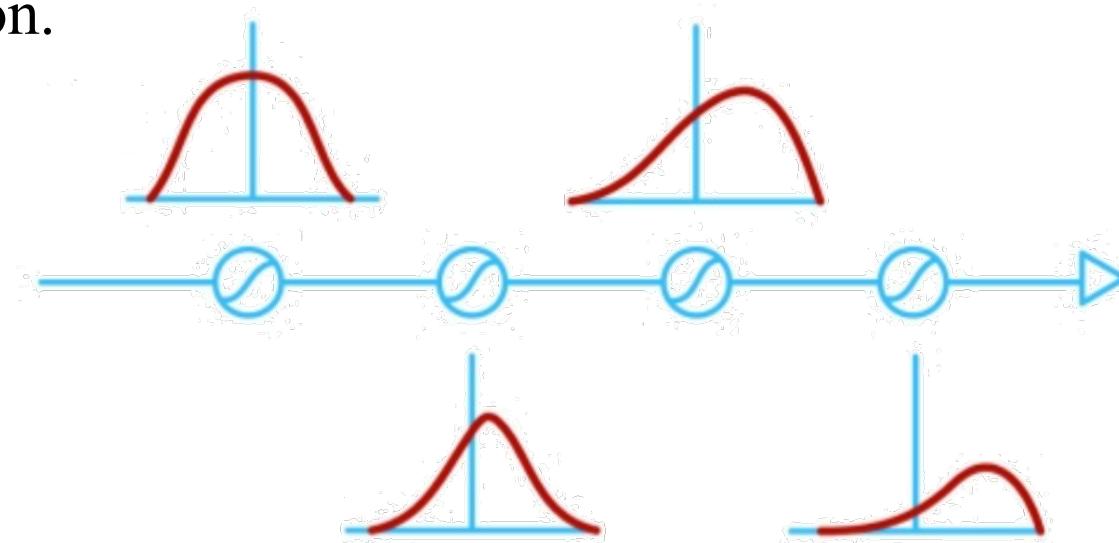
- BN reduces *Covariate Shift*. That is the change in distribution of activation of a component. By using BN, each neuron's activation (sigmoid) becomes (more or less) a Gaussian distribution, i.e. its usually not active, sometimes a bit active, rare very active.



# Batch normalization

## Why it is good?

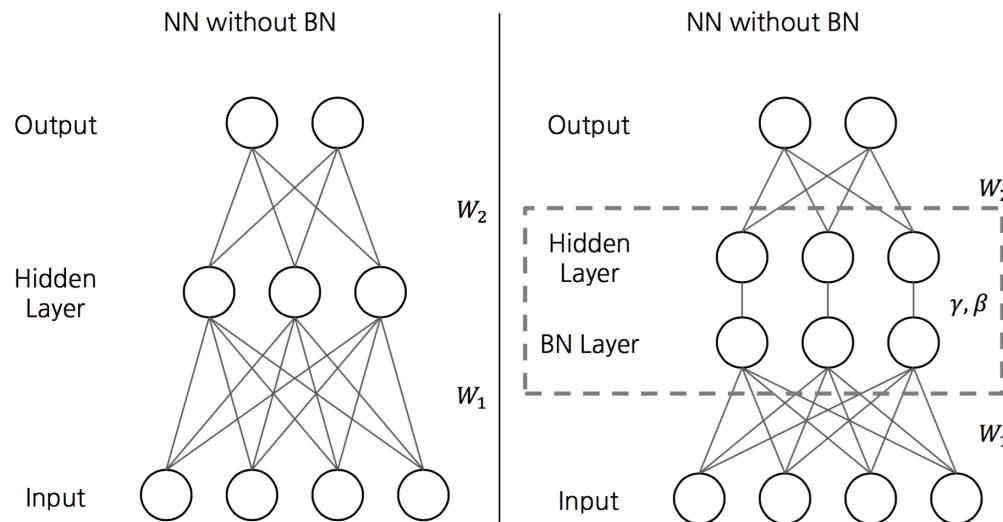
- Covariate Shift is undesirable, because the later layers have to keep adapting to the change of the type of distribution.
- BN reduces effects of exploding and vanishing gradients, because every becomes roughly normal distributed. Without BN, low activations of one layer can lead to lower activations in the next layer, and then even lower ones in the next layer and so on.



# Batch normalization layer

We introduce, for each activation  $x^{(k)}$ , a pair of parameters  $\gamma^{(k)}, \beta^{(k)}$ , which scale and shift the normalized value:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$



A new layer is added so the gradient can “see” the normalization and make adjustments if needed.

# The Batch Transformation: formally from the paper.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

$\epsilon$  is a constant added to the mini-batch variance for numerical stability.

$\gamma$  and  $\beta$  are learned to keep what the layer can represent.

# The full algorithm as proposed in the paper

**Input:** Network  $N$  with trainable parameters  $\Theta$ ;  
 subset of activations  $\{x^{(k)}\}_{k=1}^K$

**Output:** Batch-normalized network for inference,  $N_{\text{BN}}^{\text{inf}}$

- 1:  $N_{\text{BN}}^{\text{tr}} \leftarrow N$  // Training BN network
- 2: **for**  $k = 1 \dots K$  **do**
- 3:     Add transformation  $y^{(k)} = \text{BN}_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$  to  $N_{\text{BN}}^{\text{tr}}$  (Alg. 1) ←
- 4:     Modify each layer in  $N_{\text{BN}}^{\text{tr}}$  with input  $x^{(k)}$  to take  $y^{(k)}$  instead
- 5: **end for**
- 6: Train  $N_{\text{BN}}^{\text{tr}}$  to optimize the parameters  $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 7:  $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$  // Inference BN network with frozen // parameters
- 8: **for**  $k = 1 \dots K$  **do**
- 9:     // For clarity,  $x \equiv x^{(k)}$ ,  $\gamma \equiv \gamma^{(k)}$ ,  $\mu_B \equiv \mu_B^{(k)}$ , etc.
- 10:     Process multiple training mini-batches  $\mathcal{B}$ , each of size  $m$ , and average over them:  

$$\mathbb{E}[x] \leftarrow \mathbb{E}_{\mathcal{B}}[\mu_{\mathcal{B}}]$$

$$\text{Var}[x] \leftarrow \frac{m}{m-1} \mathbb{E}_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$
- 11:     In  $N_{\text{BN}}^{\text{inf}}$ , replace the transform  $y = \text{BN}_{\gamma, \beta}(x)$  with  

$$y = \frac{\gamma}{\sqrt{\text{Var}[x]+\epsilon}} \cdot x + \left(\beta - \frac{\gamma \mathbb{E}[x]}{\sqrt{\text{Var}[x]+\epsilon}}\right)$$
- 12: **end for**

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;

Parameters to be learned:  $\gamma, \beta$

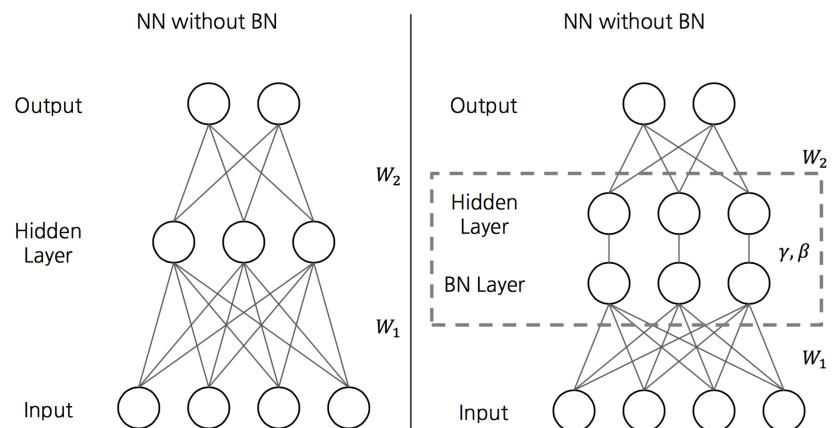
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$



**Algorithm 2:** Training a Batch-Normalized Network

# Batch normalization: Benefits in practice

- BN reduces training times. (less Covariate Shift, less exploding/vanishing gradients.)
- BN reduces demand for regularization, e.g. dropout or L2 norm.
  - Because the means and variances are calculated over batches and therefore every normalized value depends on the current batch, i.e. the network can no longer just memorize values and their correct answers.
- BN allows higher learning rates. (Because of less danger of exploding/vanishing gradients.)
- BN enables training with saturating nonlinearities in deep networks, e.g. sigmoid. (Because the normalization prevents them from getting stuck in saturating ranges, e.g. very high/low values for sigmoid.)

# The BN transformation is scalar invariant

- Batch Normalization also makes training more resilient to the parameter scale.

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

$$\text{BN}(W\mathbf{u}) = \text{BN}((aW)\mathbf{u})$$

# The BN transformation is scalar invariant

- Normally, large learning rates may increase the scale of layer parameters, which then amplify the gradient during backpropagation and lead to the model explosion
- However, with Batch Normalization, backpropagation through a layer is unaffected by the scale of its parameters.

$$\frac{\partial \text{BN}((aW)u)}{\partial u} = \frac{\partial \text{BN}(Wu)}{\partial u}$$

$$\frac{\partial \text{BN}((aW)u)}{\partial (aW)} = \frac{1}{a} \cdot \frac{\partial \text{BN}(Wu)}{\partial W}$$

Batch Normalization will stabilize the parameter growth.

## Normalization layers

### BatchNorm1d

**CLASS** `torch.nn.BatchNorm1d(num_features, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)`

[SOURCE]

Applies Batch Normalization over a 2D or 3D input (a mini-batch of 1D inputs with optional additional channel dimension) as described in the paper [Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift](#).

$$y = \frac{x - E[x]}{\sqrt{Var[x] + \epsilon}} * \gamma + \beta$$

The mean and standard-deviation are calculated per-dimension over the mini-batches and  $\gamma$  and  $\beta$  are learnable parameter vectors of size C (where C is the input size). By default, the elements of  $\gamma$  are set to 1 and the elements of  $\beta$  are set to 0.

Also by default, during training this layer keeps running estimates of its computed mean and variance, which are then used for normalization during evaluation. The running estimates are kept with a default `momentum` of 0.1.

If `track_running_stats` is set to `False`, this layer then does not keep running estimates, and batch statistics are instead used during evaluation time as well.

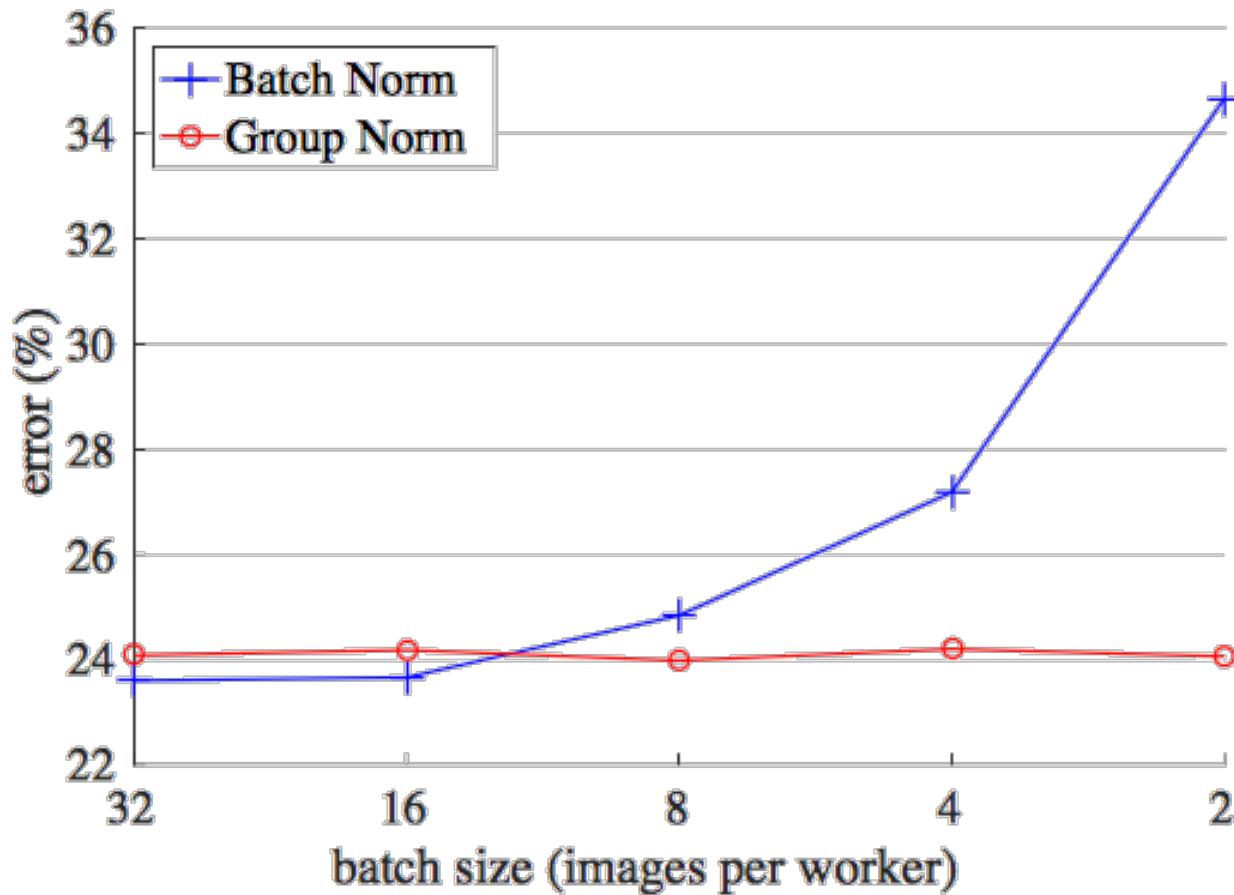
<p>torch.nn</p> <p>Parameters</p> <p>+ Containers</p> <p>+ Convolution layers</p> <p>+ Pooling layers</p> <p>+ Padding layers</p> <p>+ Non-linear activations (weighted sum)</p> <p>+ Non-linear activations (other)</p>	<p>- Normalization layers</p> <p>BatchNorm1d</p> <p>BatchNorm2d</p> <p>BatchNorm3d</p> <p>GroupNorm</p> <p>SyncBatchNorm</p> <p>InstanceNorm1d</p> <p>InstanceNorm2d</p> <p>InstanceNorm3d</p> <p>LayerNorm</p> <p>LocalResponseNorm</p>
	<p>+ Recurrent layers</p> <p>+ Transformer layers</p> <p>+ Linear layers</p> <p>+ Dropout layers</p>

# **Beyond Batch Normalization**

In batch normalization, we measure the mean/variance of the **mini batch**.

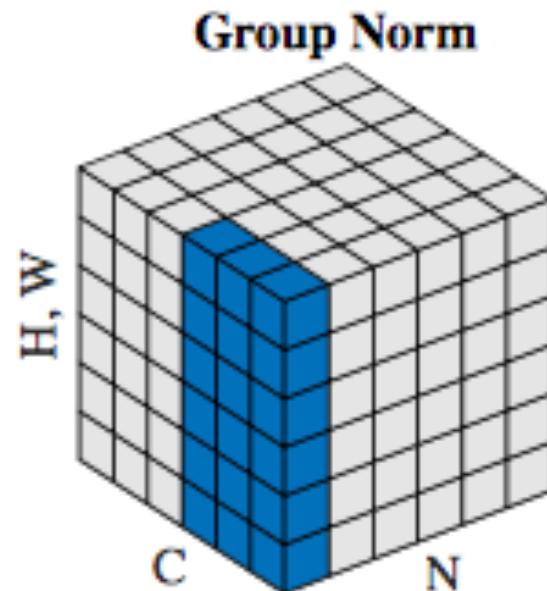
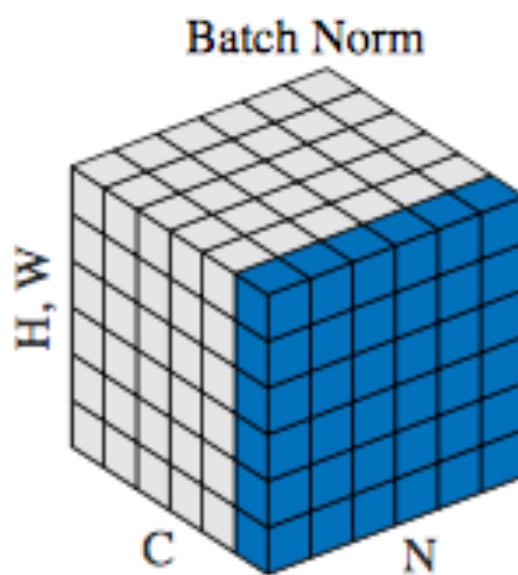
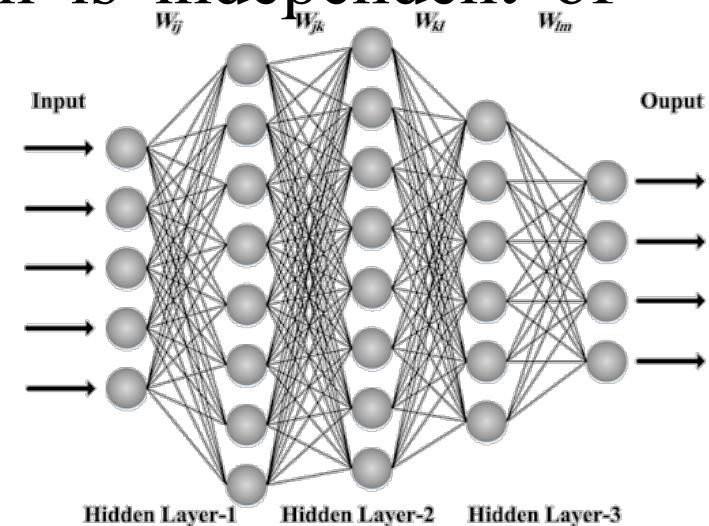
- ❖ Small batch size?
  - When the batch size is 1?
  - Noisy estimation
- ❖ Inconsistency at training and testing time?
  - At inference time, the mean and variance are pre-computed from the training set.
  - The pre-computed statistics may also change when the target data distribution changes
- ❖ Recurrent connections in an RNN?
  - Activations of each time step have different statistics.
  - Complicated and computational expensive

*BN's error increases rapidly when the batch size becomes smaller, caused by inaccurate batch statistics estimation.*

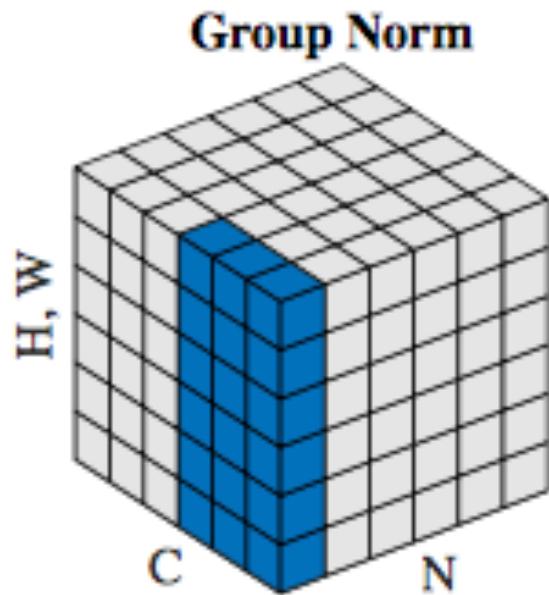
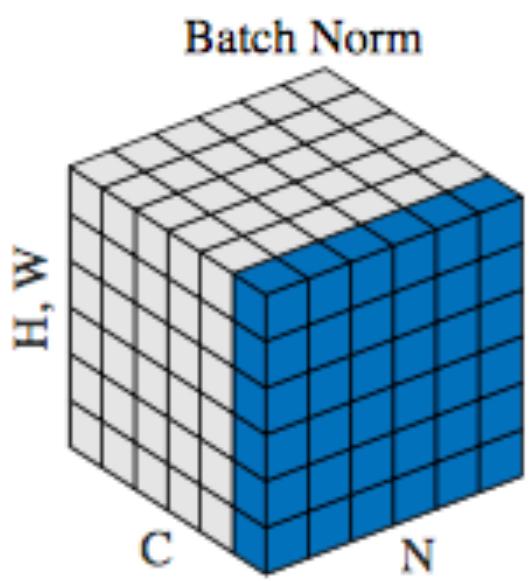


ImageNet classification error vs. batch sizes. This is a ResNet-50 model trained in the ImageNet training set using 8 workers (GPUs), evaluated in the validation set.

**Group Normalization** divides the channels into groups and computes within each group the mean and variance for normalization. GN's computation is independent of batch sizes.



# Group Normalization

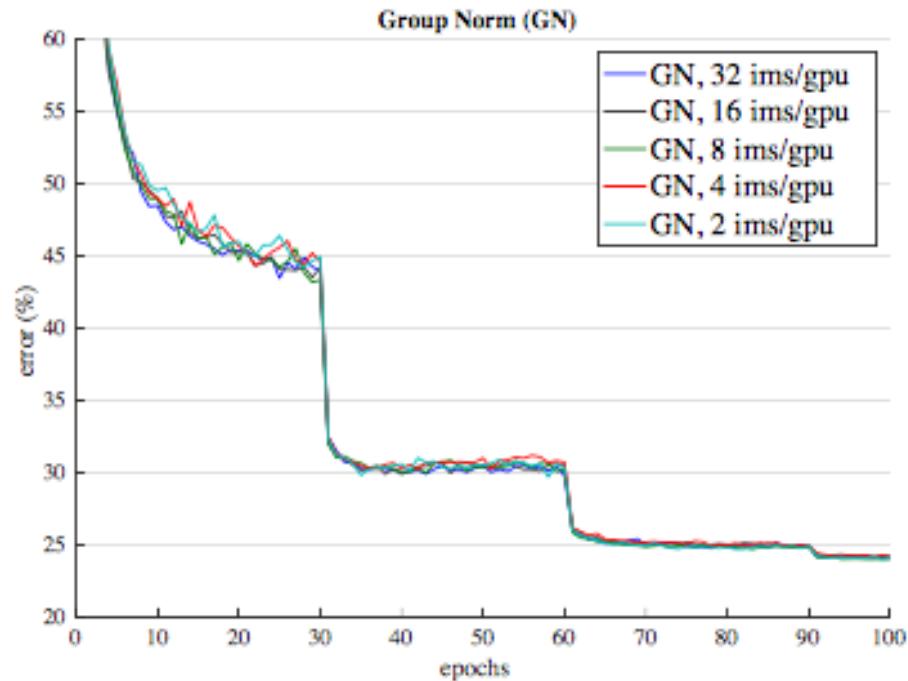
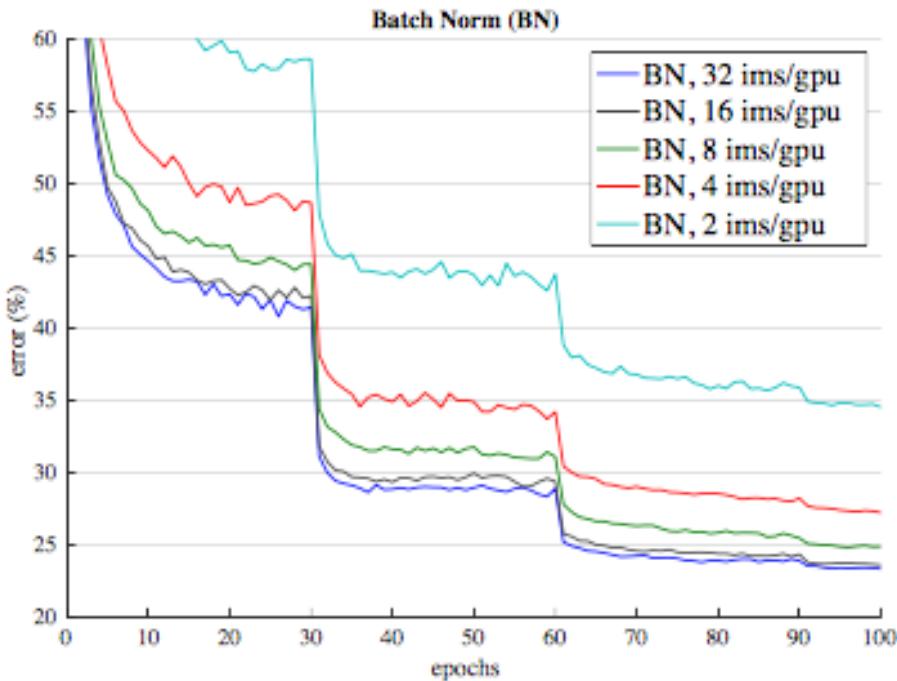


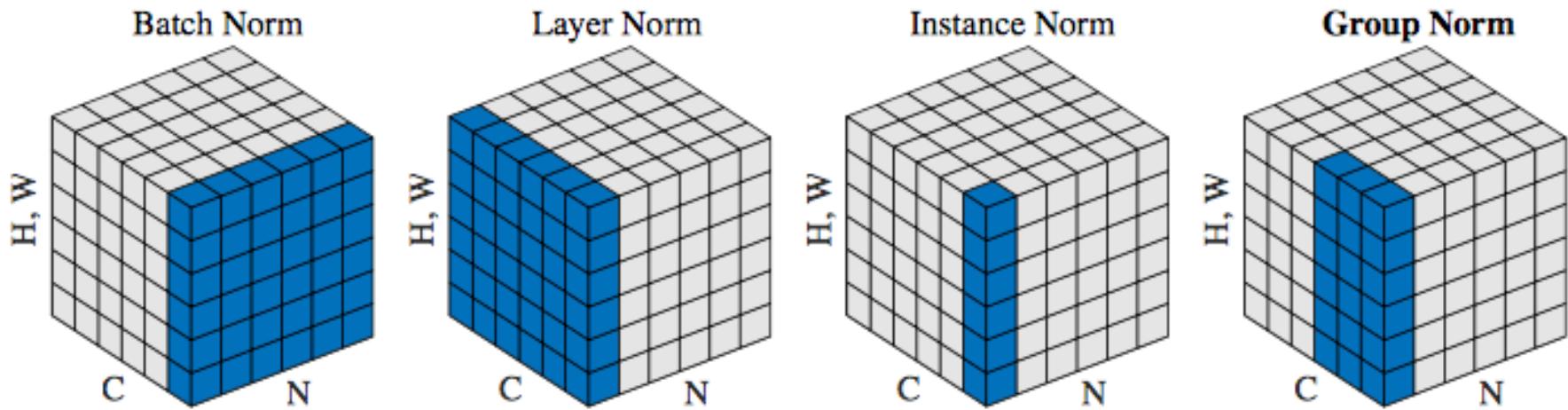
$$\mu_i = \frac{1}{m} \sum_{k \in \mathcal{S}_i} x_k,$$

$$\sigma_i = \sqrt{\frac{1}{m} \sum_{k \in \mathcal{S}_i} (x_k - \mu_i)^2 + \epsilon},$$

The pixels in the same group are normalized together by the same  $\mu$  and  $\sigma$ . GN also learns the per-channel  $\gamma$  and  $\beta$ .

# Sensitivity to batch sizes:





S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In ICML, 2015.

J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv:1607.06450*, 2016.

D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv:1607.08022*, 2016.

Wu Y, He K. Group normalization. Proceedings of the European Conference on Computer Vision (ECCV). 2018: 3-19.

# Thanks!

# **Some advanced discussions**

# Why does BatchNorm work?

Reducing Internal Covariate Shift (ICS) by normalizing activations

[When training deep models, the input distribution of each layer changes over time.] The change in the distributions of layers' inputs presents a problem because the layers need to continuously adapt to the new distribution.

[Ioffe, Szegedy 2015]

**But:** Is that really what happens?

# How Does Batch Normalization Help Optimization?

**Shibani Santurkar\***  
MIT  
[shibani@mit.edu](mailto:shibani@mit.edu)

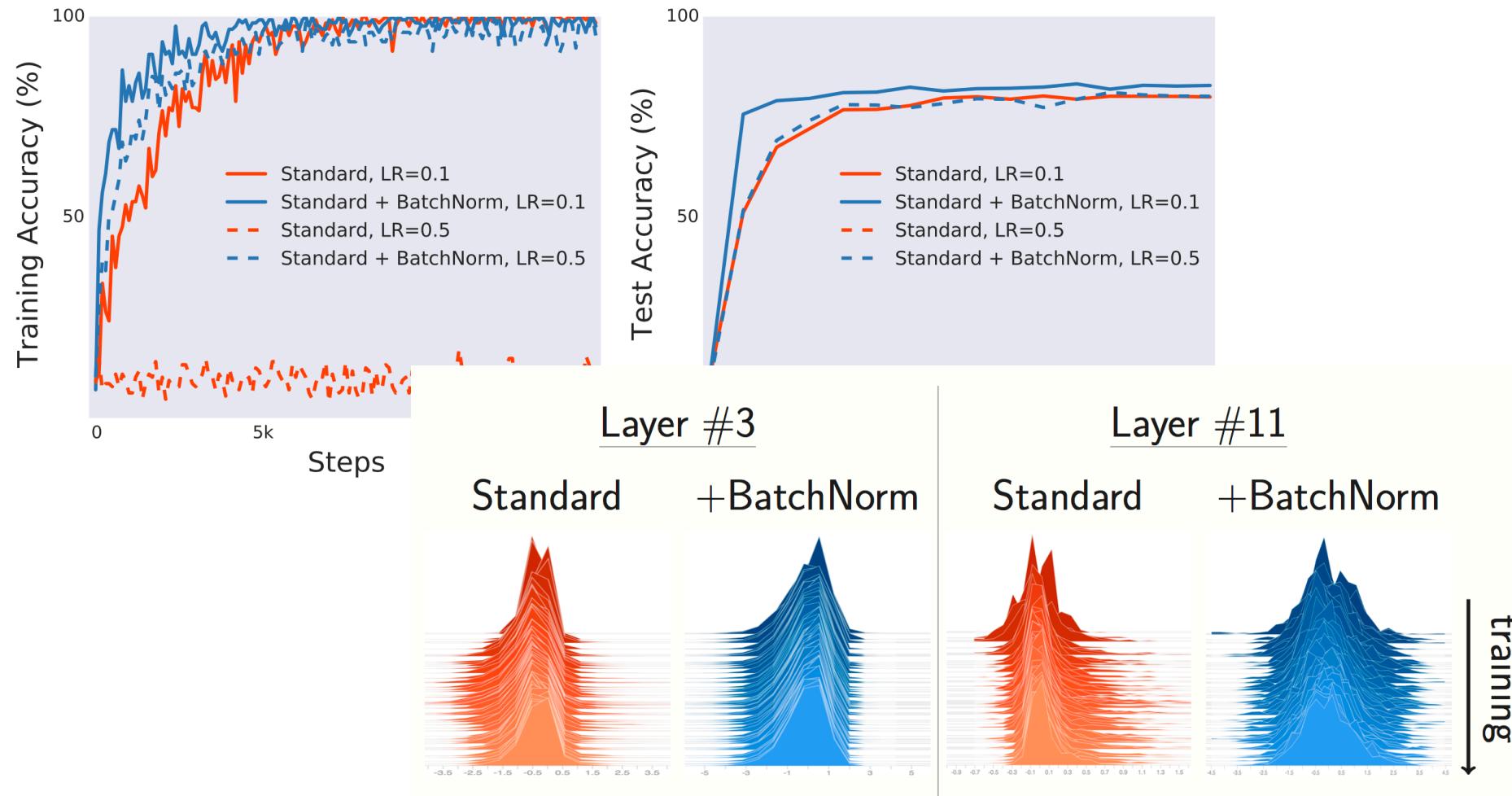
**Dimitris Tsipras\***  
MIT  
[tsipras@mit.edu](mailto:tsipras@mit.edu)

**Andrew Ilyas\***  
MIT  
[ailyas@mit.edu](mailto:ailyas@mit.edu)

**Aleksander Mądry**  
MIT  
[madry@mit.edu](mailto:madry@mit.edu)

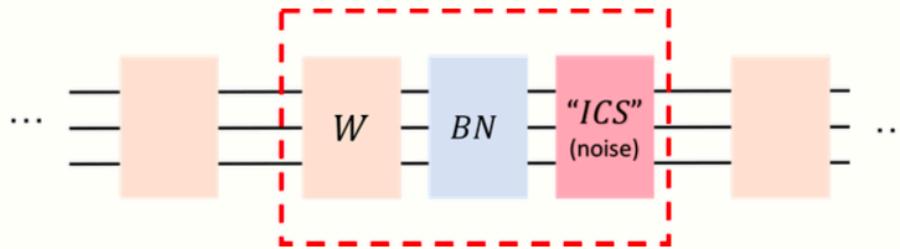
## Abstract

Batch Normalization (BatchNorm) is a widely adopted technique that enables faster and more stable training of deep neural networks (DNNs). Despite its pervasiveness, the exact reasons for BatchNorm’s effectiveness are still poorly understood. The popular belief is that this effectiveness stems from controlling the change of the layers’ input distributions during training to reduce the so-called “internal covariate shift”. In this work, we demonstrate that such distributional stability of layer inputs has little to do with the success of BatchNorm. Instead, we uncover a more fundamental impact of BatchNorm on the training process: it makes the optimization landscape significantly smoother. This smoothness induces a more predictive and stable behavior of the gradients, allowing for faster training.

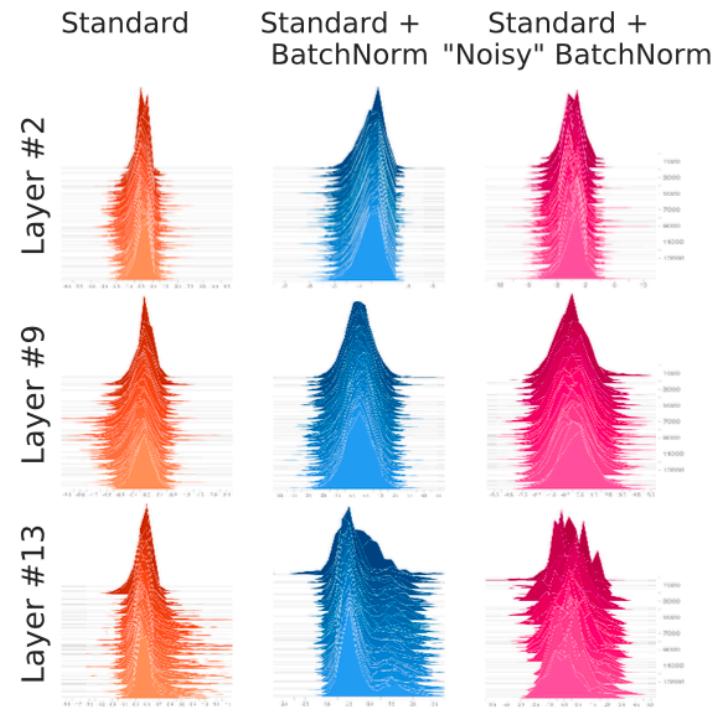
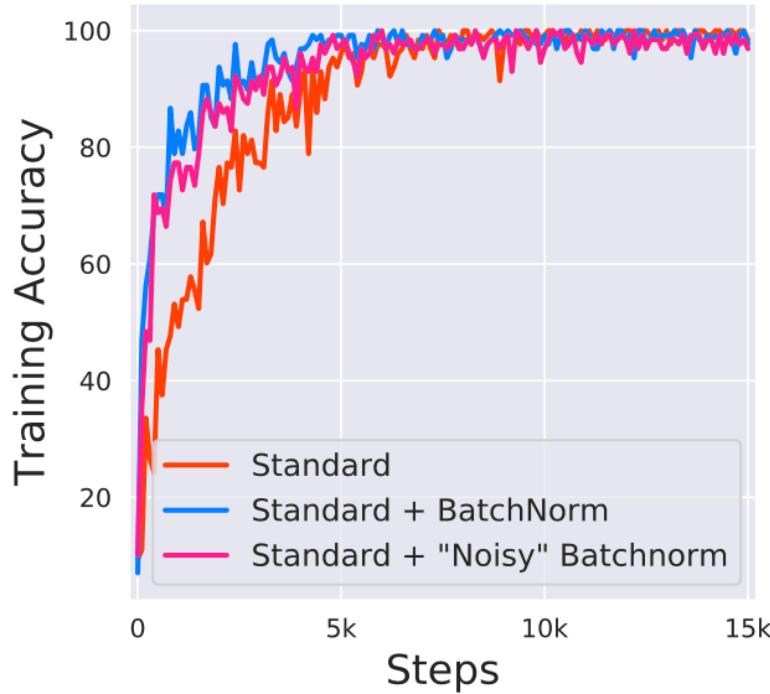


Even though the gap between the performance of the BatchNorm and non-BatchNorm networks is clear, the difference in the evolution of layer input distributions seems to be much less pronounced.

# What if we introduce additional (artificial) ICS?

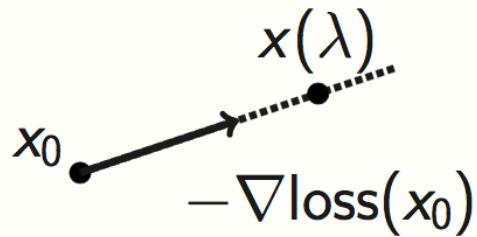


**Specifically:** We add **time-varying** noise (with **non-zero mean**) to the **outputs** of BatchNorm layers



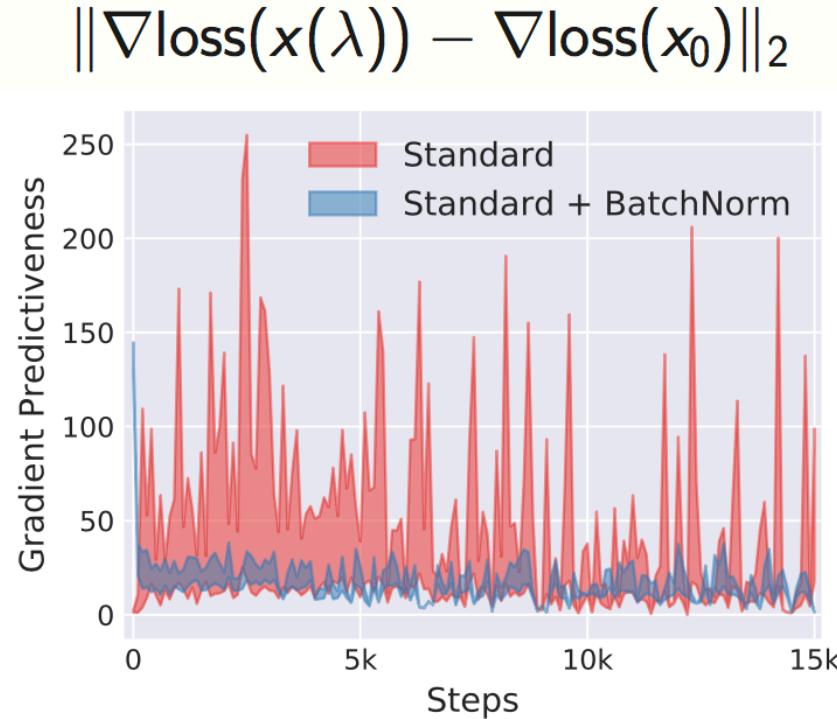
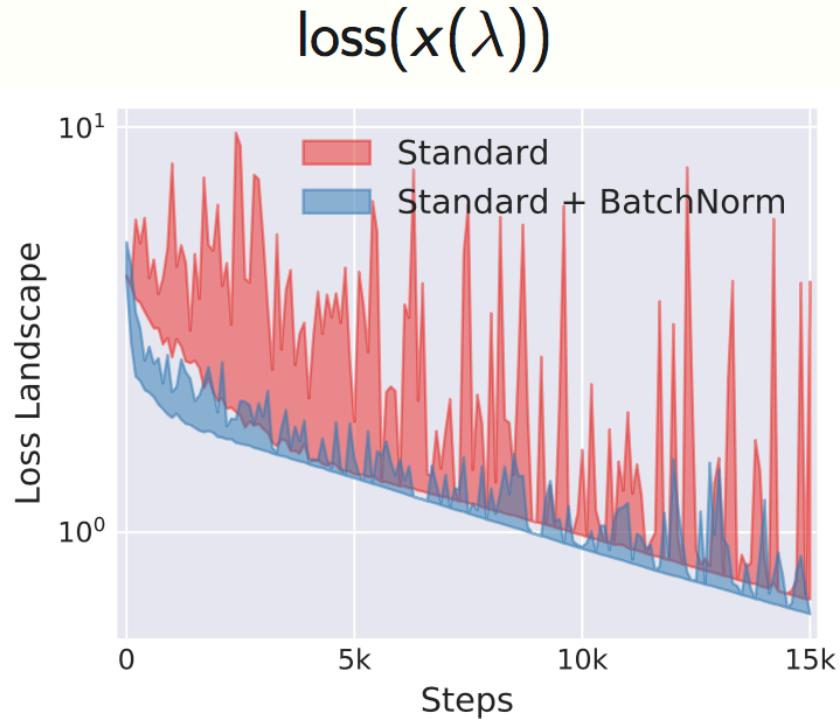
# Roots of BatchNorm's success

**Our approach:** Examine the loss and gradient landscape

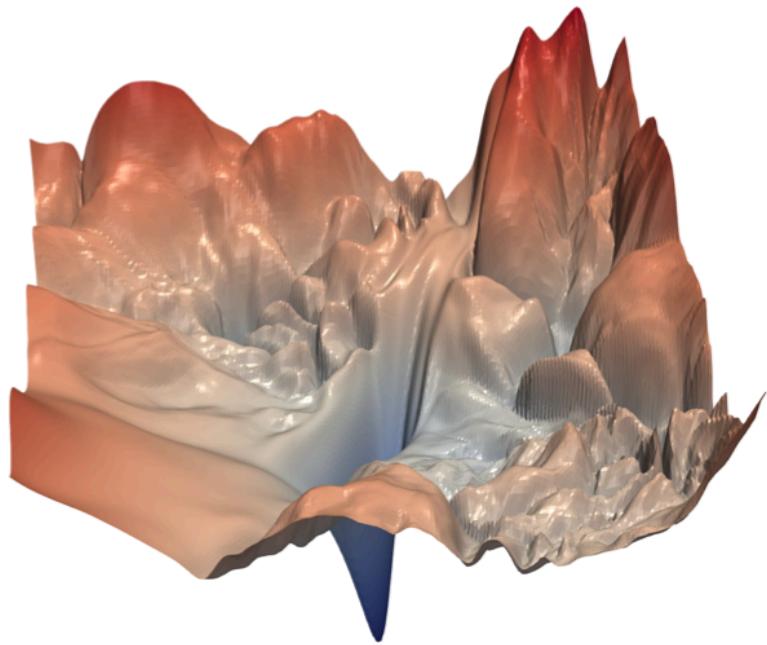


$$x(\lambda) = x_0 - \lambda \nabla \text{loss}(x_0)$$

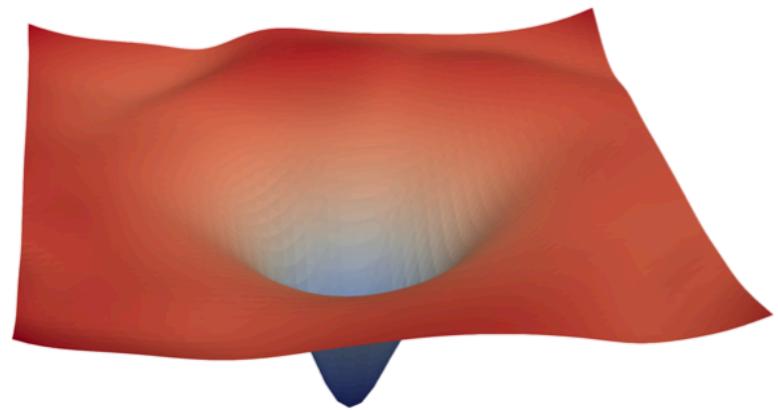
**Specifically:** Measure variation of loss and gradient over  $\lambda$



# The smoothing effect of BatchNorm



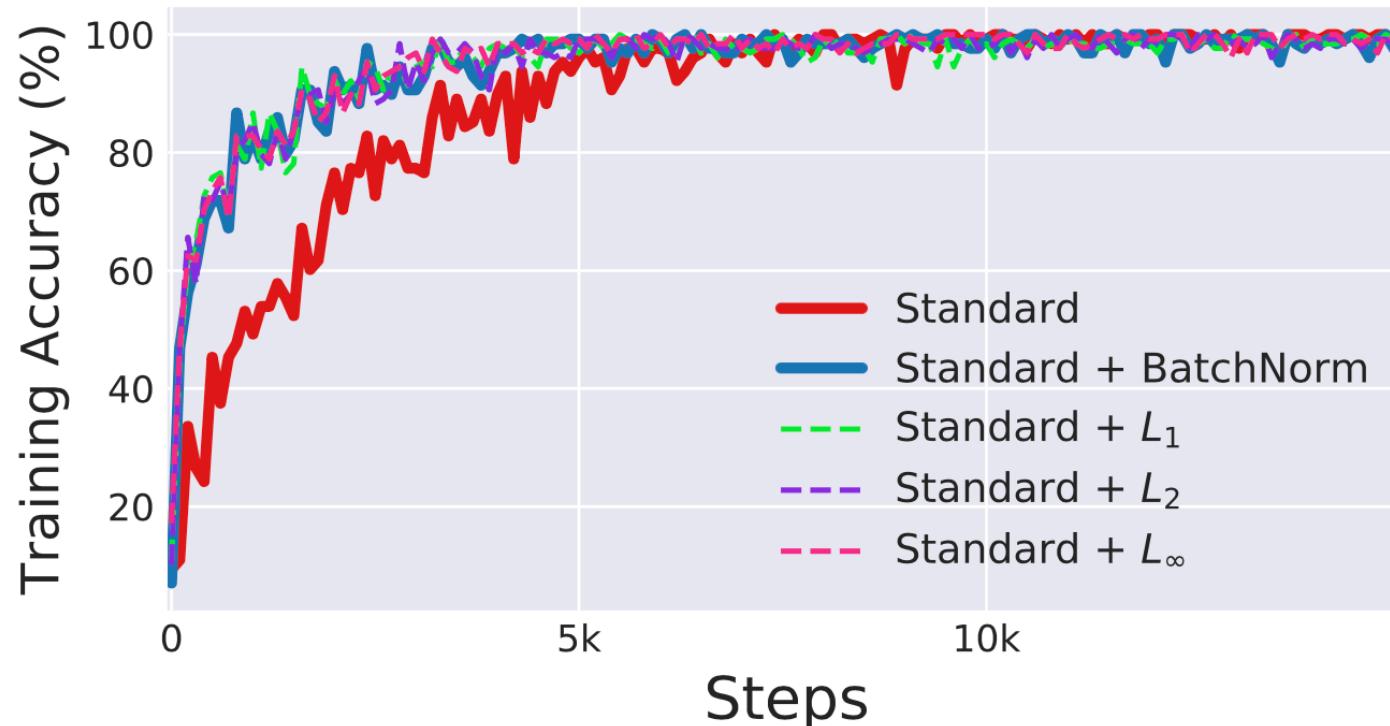
(a) without skip connections



(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

# Is BatchNorm the best (only?) way to smoothen the landscape?



(a) VGG