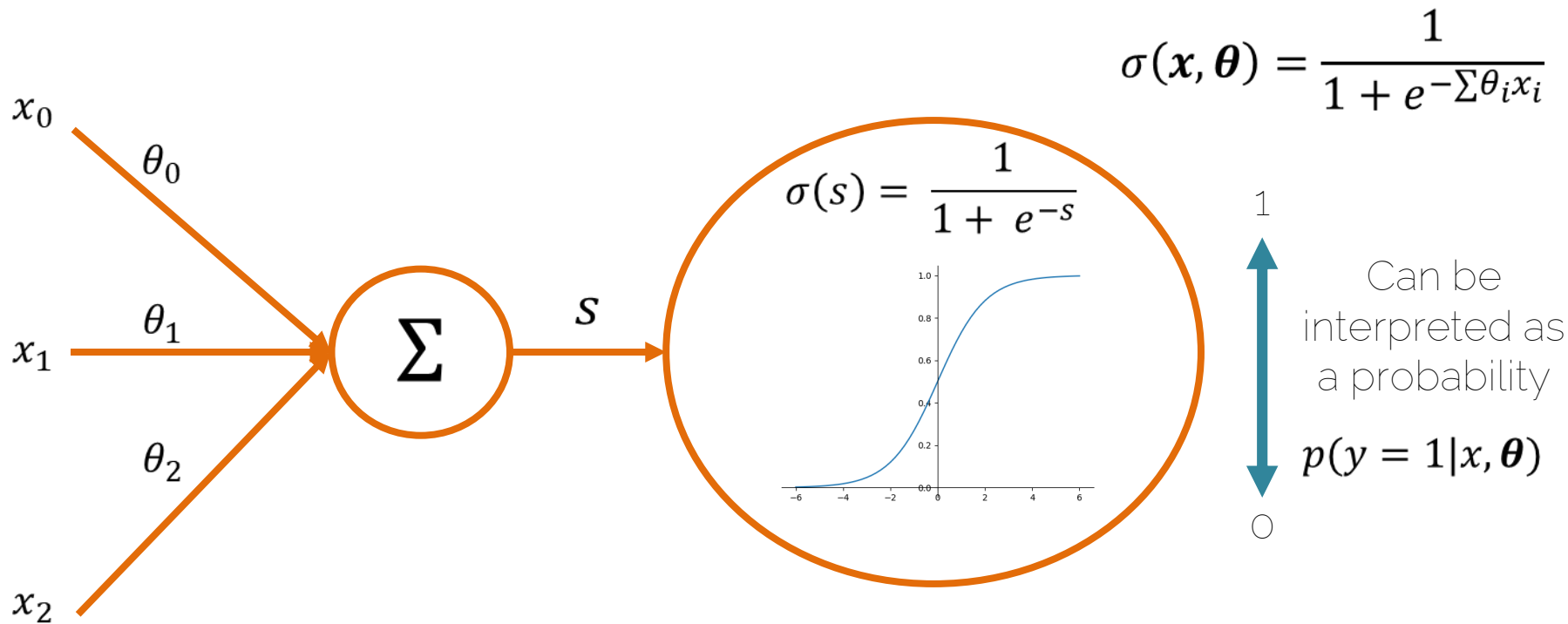


Data Augmentation and Advanced Regularization

Lecture 7 Recap

Binary Classification: Sigmoid



Multiclass Classification: Softmax

- Softmax

$$p(y_i | \mathbf{x}_i, \Theta) = \frac{e^{s_{y_i}}}{\sum_{k=1}^C e^{s_k}} = \frac{e^{\mathbf{x}_i \boldsymbol{\theta}_{y_i}}}{\sum_{k=1}^C e^{\mathbf{x}_i \boldsymbol{\theta}_k}}$$

Probability of the true class

Exp

normalize

training pairs $[\mathbf{x}_i; y_i]$,
 $\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{1, 2 \dots C\}$
 y_i : label (true class)

Parameters:

$$\Theta = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_C]$$

C : number of classes

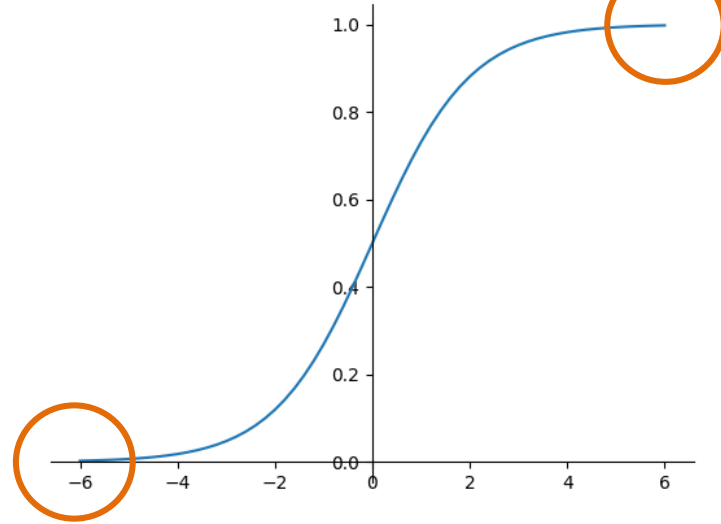
s : score of the class

1. Exponential operation: make sure probability > 0
2. Normalization: make sure probabilities sum up to 1.

Sigmoid Activation

Forward

$$\sigma(s) = \frac{1}{1 + e^{-s}}$$



X Saturated neurons kill the gradient flow

$$\cancel{\frac{\partial L}{\partial w}} = \frac{\partial s}{\partial w} \frac{\partial L}{\partial s}$$

$$\cancel{\frac{\partial L}{\partial s}} = \frac{\partial \sigma}{\partial s} \frac{\partial L}{\partial \sigma}$$

$$\cancel{\frac{\partial \sigma}{\partial s}}$$

$$\frac{\partial L}{\partial \sigma}$$

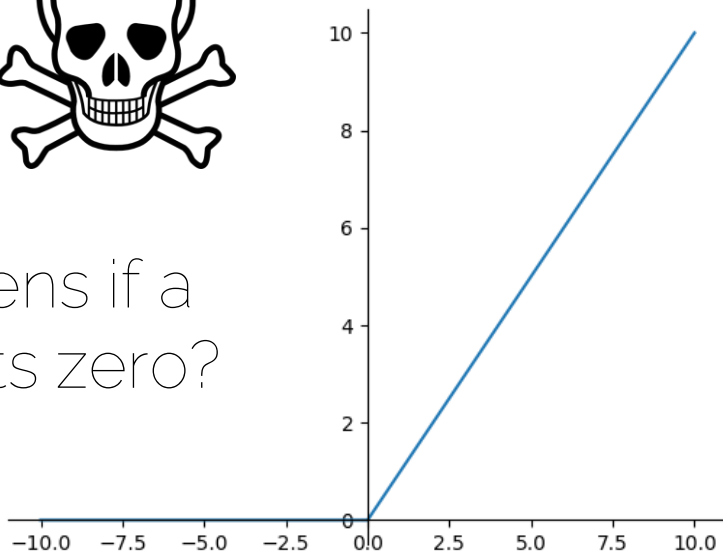
Rectified Linear Units (ReLU)



Dead ReLU



What happens if a ReLU outputs zero?



Large and consistent gradients



Fast convergence



Does not saturate

[Krizhevsky et al. NeurIPS 2012] ImageNet Classification with Deep Convolutional Neural Networks

Xavier/Kaiming Initialization

- How to ensure the variance of the output is the same as the input?

$$\underbrace{(n\text{Var}(w)\text{Var}(x))}_{= 1}$$

$$\text{Var}(w) = \frac{1}{n}$$

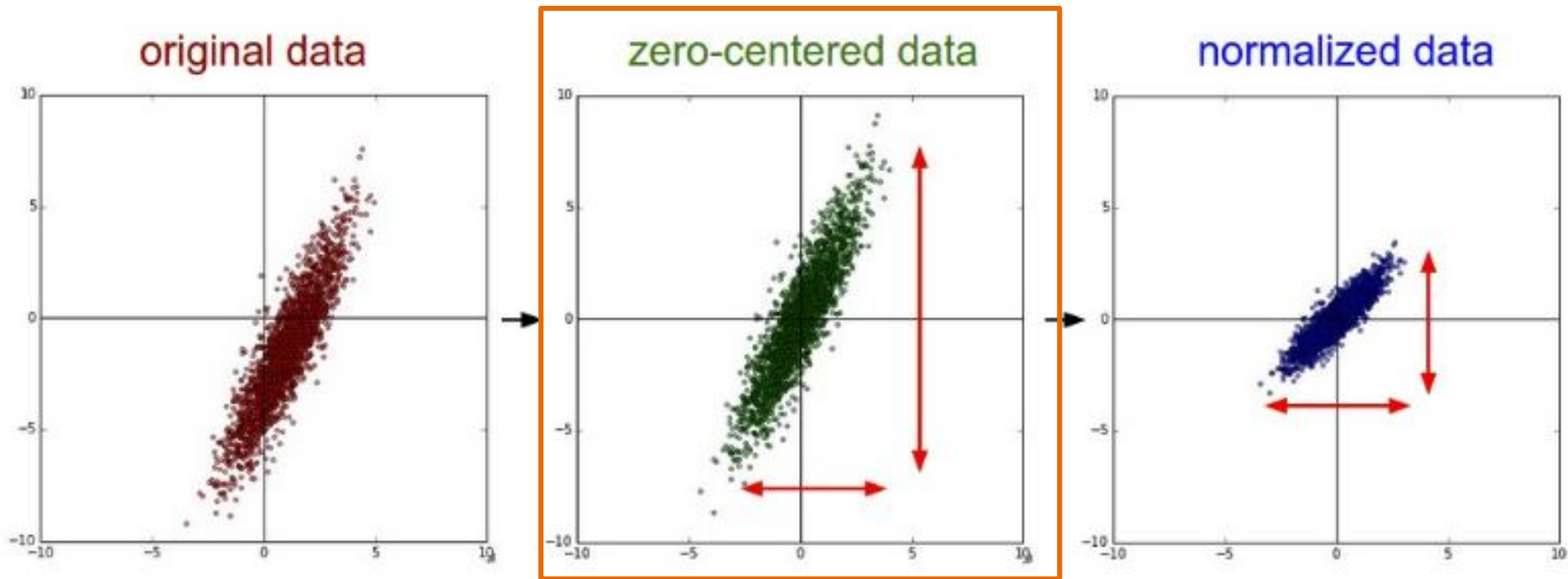
ReLU Kills half of the activations
-> adjust var by a factor of 2

$$\text{Var}(w) = \frac{2}{n}$$

Lecture 8

Data Augmentation

Data Pre-Processing



For images subtract the mean image (AlexNet) or per-channel mean (VGG-Net)

Data Augmentation

- A classifier has to be invariant to a wide variety of transformations

All

Images

Videos

News

Shopping

More

Settings

Tools

SafeSearch ▼



Cute



And Kittens



Clipart



Drawing



Cute Baby



White Cats And Kittens



Pose

Appearance

Illumination

Data Augmentation

- A classifier has to be invariant to a wide variety of transformations
- Helping the classifier: synthesize data simulating plausible transformations

Data Augmentation

a. No augmentation (= 1 image)



b. Flip augmentation (= 2 images)

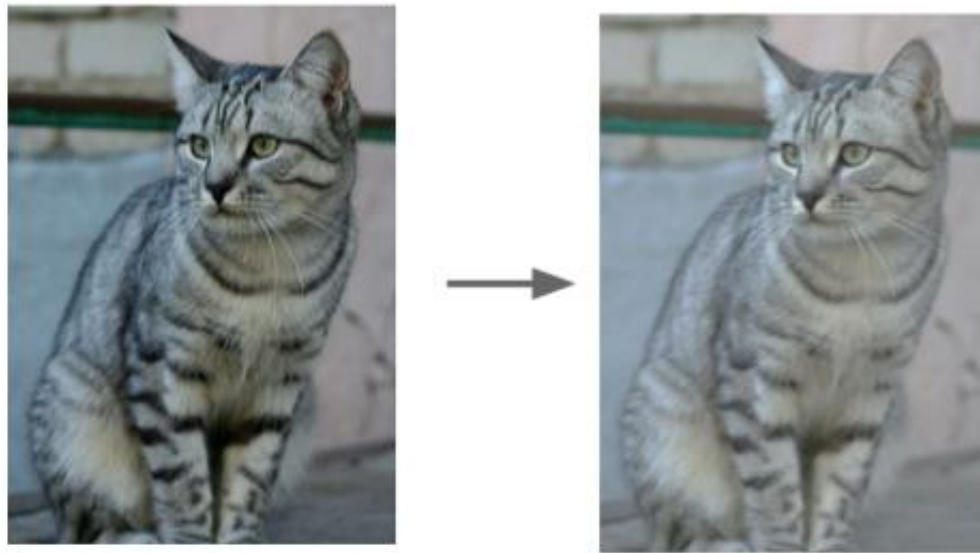


c. Crop+Flip augmentation (= 10 images)



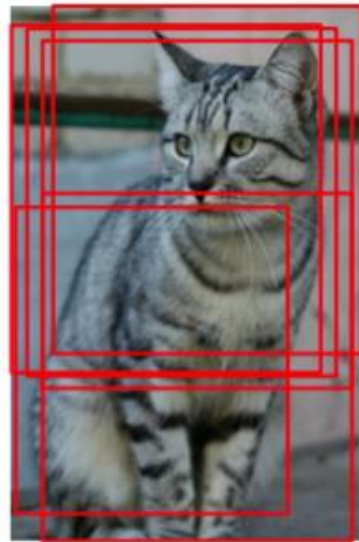
Data Augmentation: Brightness

- Random brightness and contrast changes

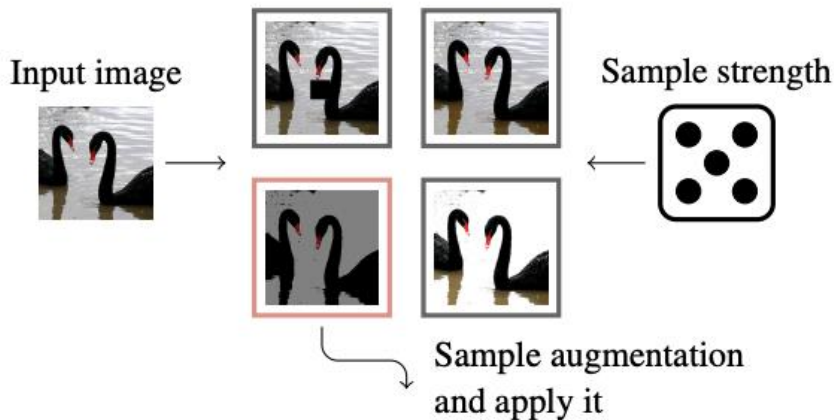
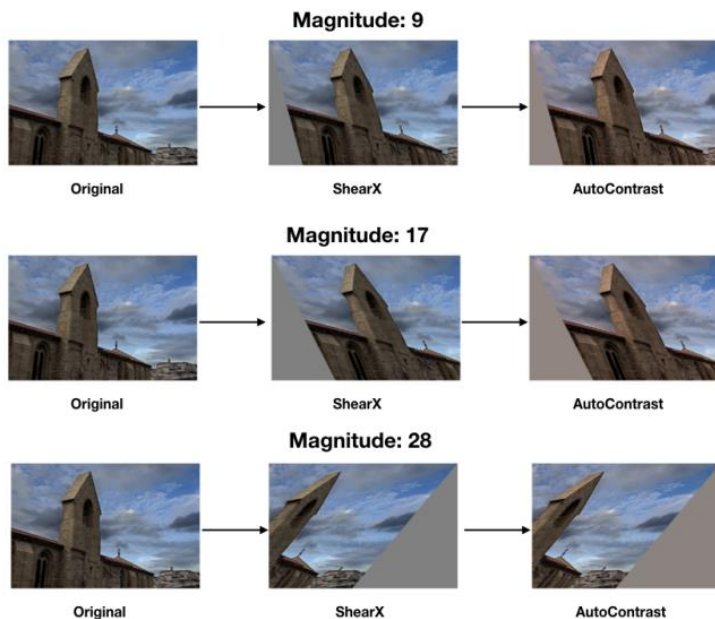


Data Augmentation: Random Crops

- Training: random crops
 - Pick a random L in $[256, 480]$
 - Resize training image, short side L
 - Randomly sample crops of 224×224
- Testing: fixed set of crops
 - Resize image at N scales
 - 10 fixed crops of 224×224 : (4 corners + 1 center) \times 2 flips



Data Augmentation: Advanced



Algorithm 1 TrivialAugment Procedure

- 1: **procedure** TA(x : image)
 - 2: Sample an augmentation a from \mathcal{A}
 - 3: Sample a strength m from $\{0, \dots, 30\}$
 - 4: Return $a(x, m)$
 - 5: **end procedure**
-

Cubuk et al., RandAugment, CVPRW 2020

Muller et al., Trivial Augment, ICCV 2021

Data Augmentation

- When comparing two networks make sure to use the same data augmentation!
- Consider data augmentation a part of your network design

Augmentation – Practical Considerations

- Augmentations should not distort the labels (e.g., '6' vs '9')
- Memory vs speed: on-the-fly vs pre-computed
- Test-time augmentation: generated multiple augmentations of an input image and aggregate model predictions (more robustness)

Advanced Regularization

L2 regularization, also (wrongly) called weight decay

- L2 regularization

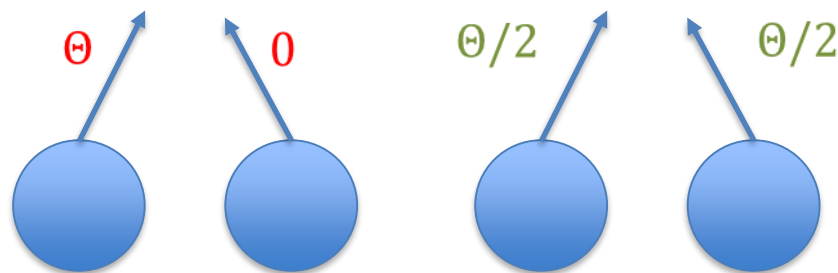
$$\Theta_{k+1} = \Theta_k - \epsilon \nabla_{\Theta} f(\Theta_k) - \lambda \Theta_k$$

Learning rate

Gradient

Gradient of L2-regularization

- Penalizes large weights
- Improves generalization



L2 regularization, also (wrongly) called weight decay

- Weight decay regularization

$$\Theta_{k+1} = (1 - \lambda)\Theta_k - \alpha \nabla_{\Theta} f(\Theta_k)$$

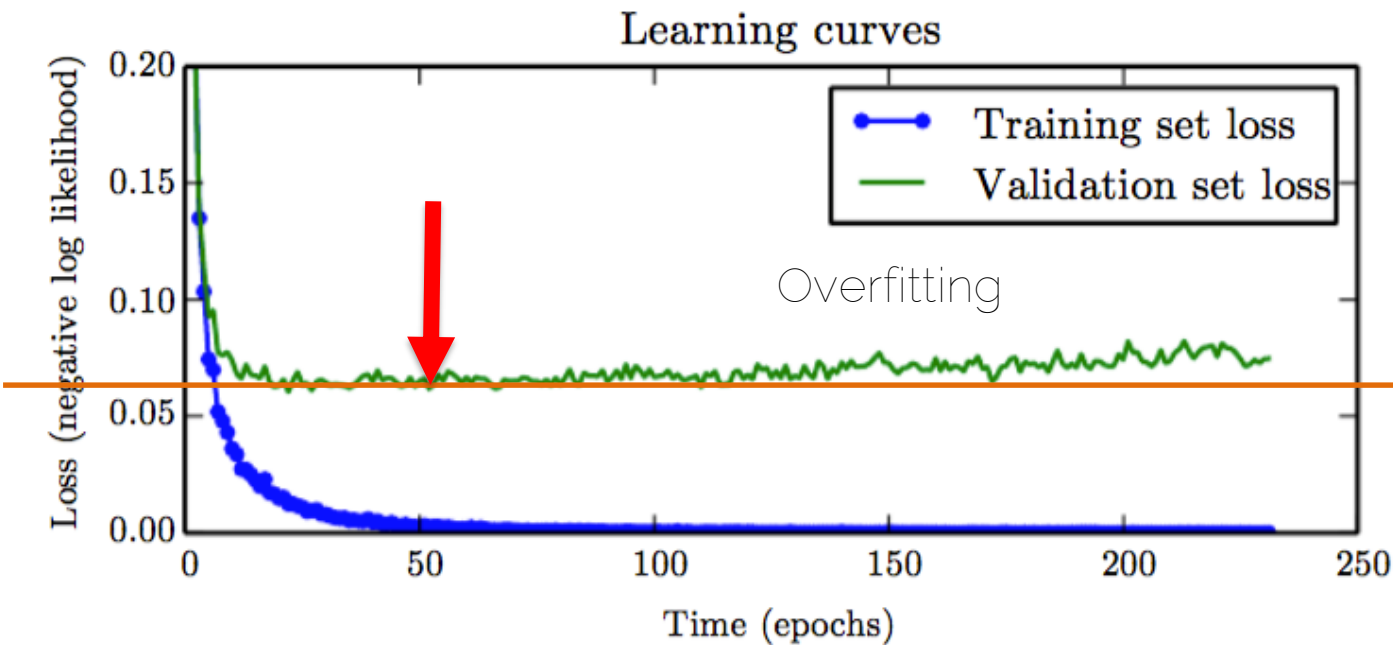
Learning rate of weight
decay

Learning rate of the
optimizer

- Equivalent to L2 regularization in GD, but not in Adam.

Loshchilov and Hutter, Decoupled Weight Decay
Regularization, ICLR 2019

Early Stopping

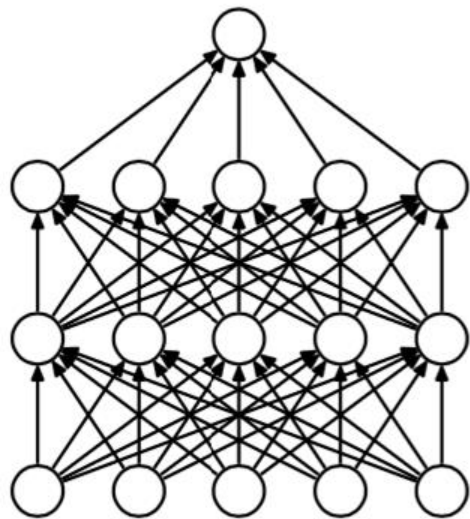


Bagging and Ensemble Methods

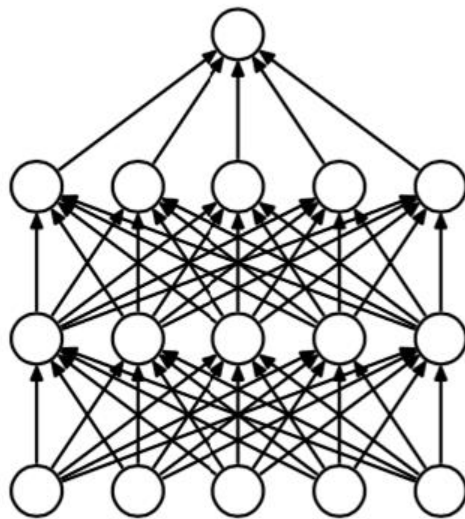
- Train multiple models and average their results
- E.g., use a different algorithm for optimization or change the objective function / loss function.
- If errors are uncorrelated, the expected combined error will decrease linearly with the ensemble size

Bagging and Ensemble Methods

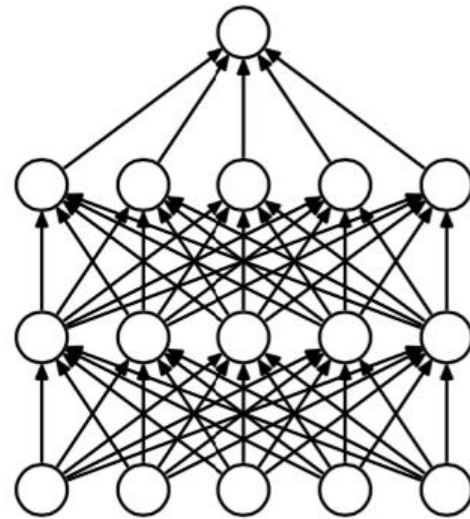
- Bagging: uses k different datasets (or SGD/init noise)



Training Set 1



Training Set 2



Training Set 3

Image Source: [Srivastava et al., JMLR'14] Dropout

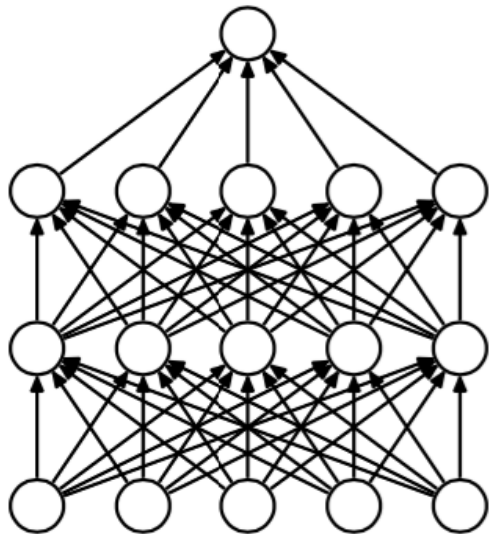
Ensembling Variants

- Avoid training multiple different models
- Different checkpoints as ensemble members
- Ensemble via subnetworks
 - Train one big network that acts as an ensemble
 - E.g., multiple inputs -> multiple outputs (MIMO)
 - Single shared network that acts as ensemble (different inputs)

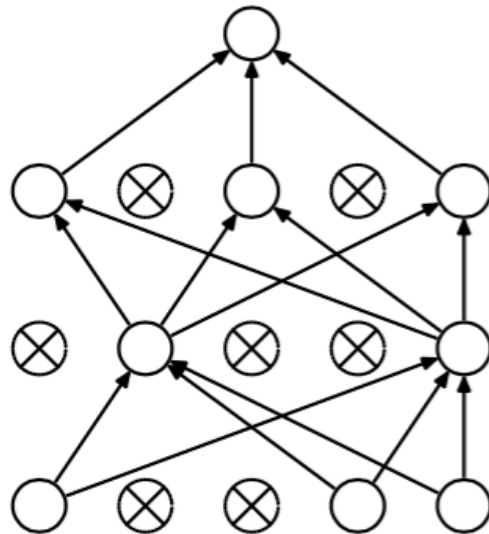
Dropout

Dropout

- Disable a random set of neurons (typically 50%)



(a) Standard Neural Net

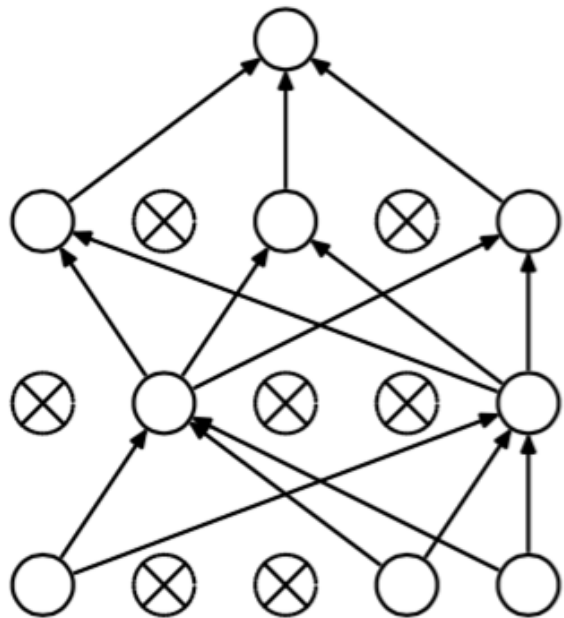


(b) After applying dropout.

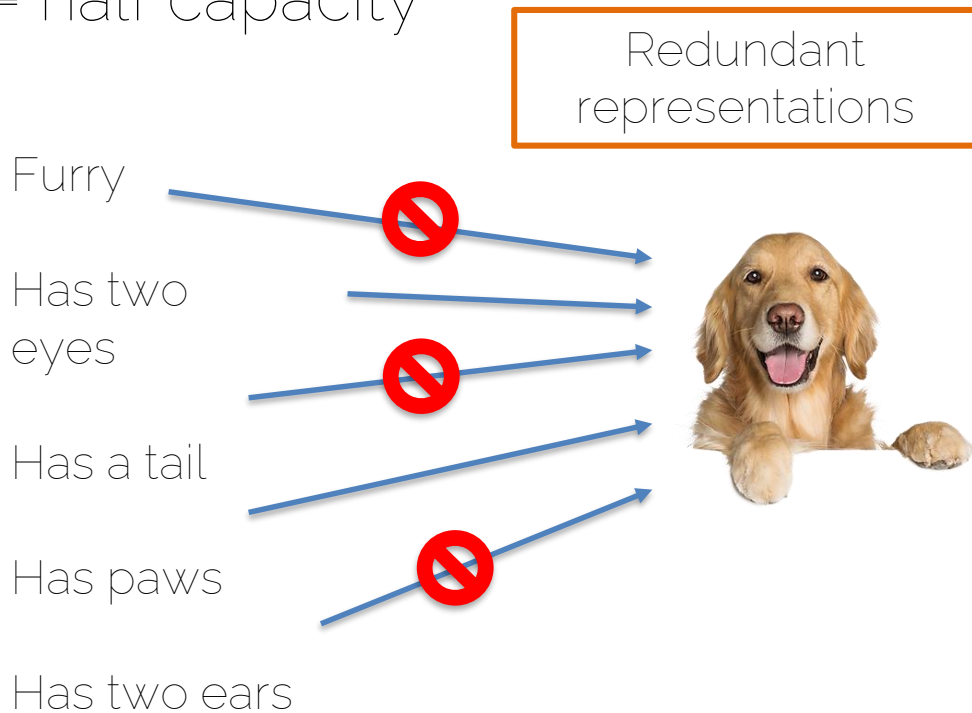
Forward ↑

Dropout: Intuition

- Using half the network = half capacity



(b) After applying dropout.

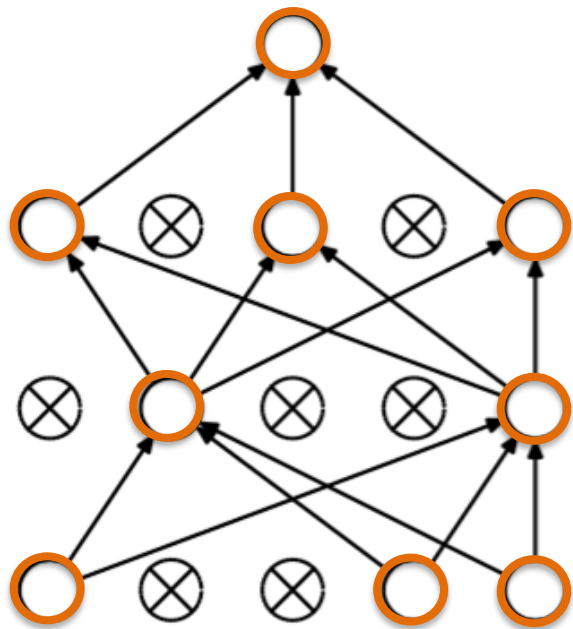


Dropout: Intuition

- Using half the network = half capacity
 - Redundant representations
 - Base your scores on more features
- Consider it as a model ensemble

Dropout: Intuition

- Two models in one



(b) After applying dropout.



Model 1



Model 2



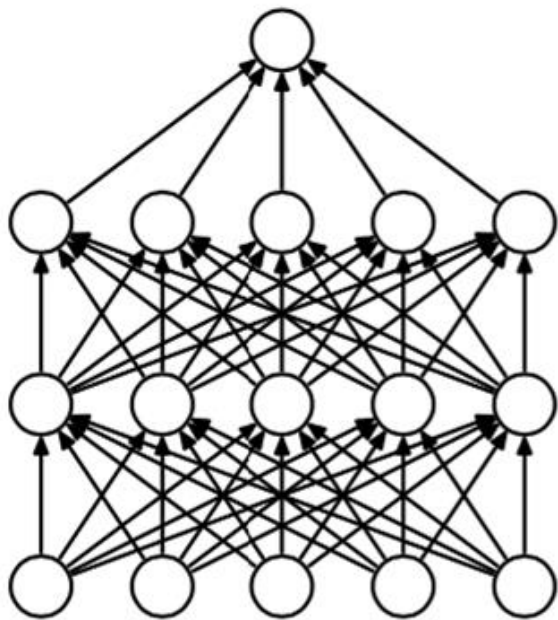
Dropout: Intuition

- Using half the network = half capacity
 - Redundant representations
 - Base your scores on more features
- Consider it as two models in one
 - Training a large ensemble of models, each on different set of data (mini-batch) and with SHARED parameters

Reducing co-adaptation between neurons

Dropout: Test Time

- All neurons are “turned on” – no dropout



Conditions at train and test time are not the same

PyTorch: `model.train()` and `model.eval()`

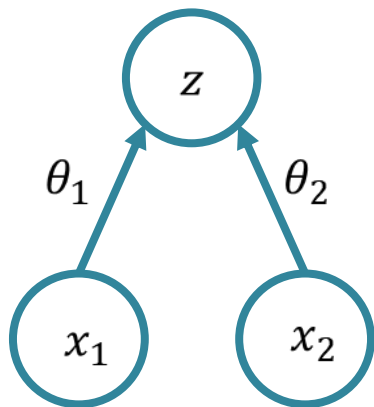
Dropout: Test Time

- Test:

$$z = (\theta_1 x_1 + \theta_2 x_2) \cdot p$$

Dropout
probability
 $p = 0.5$

- Train:



Weight scaling
inference rule

$$\begin{aligned} E[z] &= \frac{1}{4} (\theta_1 0 + \theta_2 0 \\ &\quad + \theta_1 x_1 + \theta_2 0 \\ &\quad + \theta_1 0 + \theta_2 x_2 \\ &\quad + \theta_1 x_1 + \theta_2 x_2) \\ &= \frac{1}{2} (\theta_1 x_1 + \theta_2 x_2) \end{aligned}$$

Dropout: Before

- Efficient bagging method with parameter sharing
- Try it!
- Dropout reduces the effective capacity of a model, but needs more training time
- Efficient regularization method, can be used with L2

Dropout: Nowadays

- Usually does not work well when combined with batch-norm.
- Training takes a bit longer, usually 1.5x
- But, can be used for uncertainty estimation.
- Monte Carlo dropout (Yarin Gal and Zoubin Ghahramani series of papers).

Monte Carlo Dropout

- Neural networks are massively overconfident.
- We can use dropout to make the softmax probabilities more calibrated.
- Training: use dropout with a low p (0.1 or 0.2).
- Inference, run the same image multiple times (25-100), and average the results.

Gal et al., Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference, ICLRW 2015

Gal and Ghahramani, Dropout as a Bayesian approximation, ICML 2016

Gal et al., Deep Bayesian Active Learning with Image Data, ICML 2017

Gal, Uncertainty in Deep Learning, PhD thesis 2017

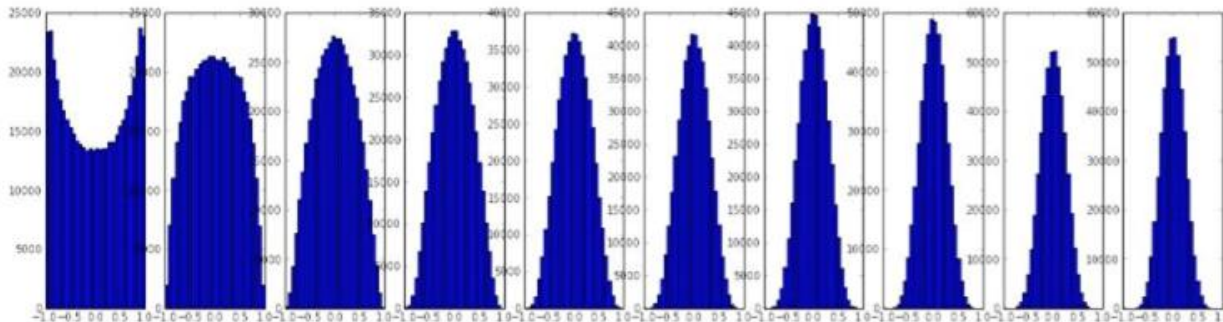
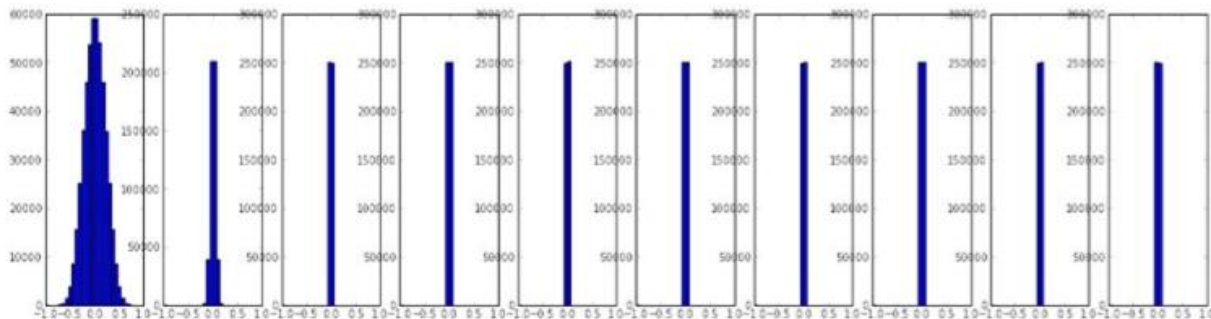
Batch Normalization: Reducing Internal Covariate Shift

Batch Normalization: Reducing Internal Covariate Shift

What is internal covariate shift, by the way?

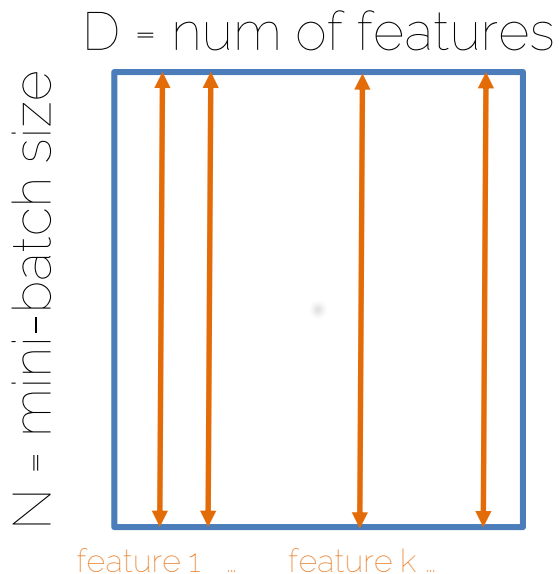
Our Goal

- All we want is that our activations do not die out



Batch Normalization

- Wish: Unit Gaussian activations (in our example)
- Solution: let's do it



Mean of your mini-batch
examples over feature k

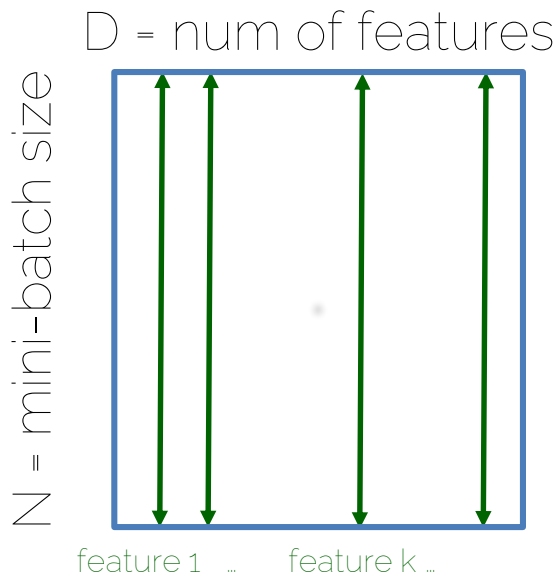
$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - E[\mathbf{x}^{(k)}]}{\sqrt{\text{Var}[\mathbf{x}^{(k)}]}}$$

An orange arrow points from the text 'Mean of your mini-batch examples over feature k' to the term $E[\mathbf{x}^{(k)}]$ in the numerator of the equation.

[Ioffe and Szegedy, PMLR'15] Batch Normalization

Batch Normalization

- In each dimension of the features, you have a unit gaussian (in our example)



Mean of your mini-batch examples over feature k

Unit Gaussian

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - E[\mathbf{x}^{(k)}]}{\sqrt{\text{Var}[\mathbf{x}^{(k)}]}}$$

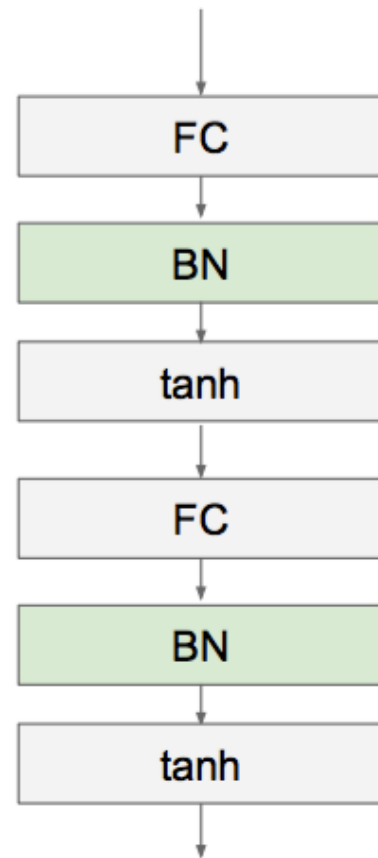
[Ioffe and Szegedy, PMLR'15] Batch Normalization

Batch Normalization

- In each dimension of the features, you have a unit gaussian (in our example)
- For NN in general, BN normalizes the mean and variance of the inputs to your activation functions

BN Layer

- A layer to be applied after Fully Connected (or Convolutional) layers and before non-linear activation functions



[Ioffe and Szegedy, PMLR'15] Batch Normalization

Batch Normalization

- 1. Normalize

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - E[\mathbf{x}^{(k)}]}{\sqrt{Var[\mathbf{x}^{(k)}]}}$$

← Differentiable function so we can backprop through it...

- 2. Allow the network to change the range

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

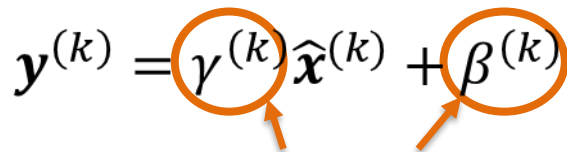
← These parameters will be optimized during backprop

Batch Normalization

- 1. Normalize

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - E[\mathbf{x}^{(k)}]}{\sqrt{Var[\mathbf{x}^{(k)}]}}$$

- 2. Allow the network to change the range

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


backprop

The network *can* learn to undo the normalization

$$\gamma^{(k)} = \sqrt{Var[\mathbf{x}^{(k)}]}$$

$$\beta^{(k)} = E[\mathbf{x}^{(k)}]$$

Batch Normalization

- Ok to treat dimensions separately?
Shown empirically that even if features are not correlated, convergence is still faster with this method

BN: Train vs Test

- Train time: mean and variance is taken over the mini-batch

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mathbf{E}[\mathbf{x}^{(k)}]}{\sqrt{\text{Var}[\mathbf{x}^{(k)}]}}$$

- Test-time: what happens if we can just process one image at a time?
 - No chance to compute a meaningful mean and variance

BN: Train vs Test

Training: Compute mean and variance from mini-batch 1,2,3 ...

Testing: Compute mean and variance by running an exponentially weighted averaged across training mini-batches:

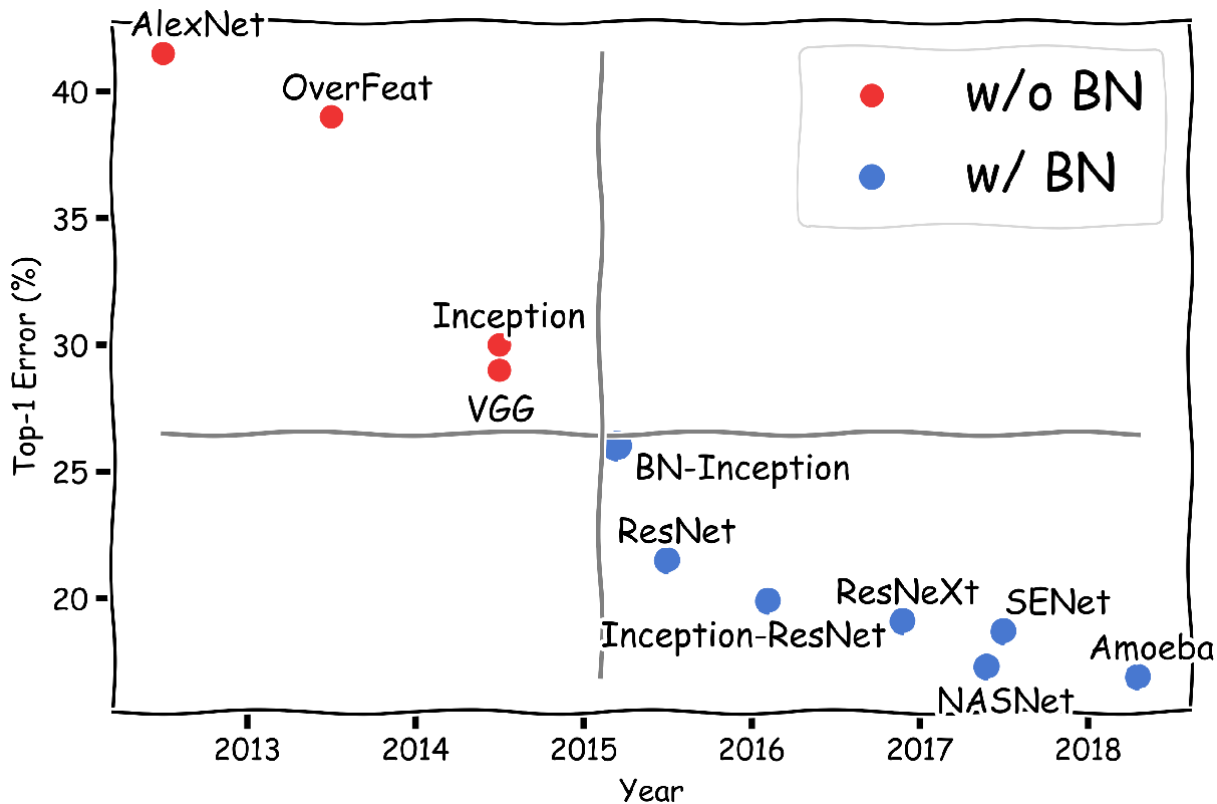
$$\begin{aligned} Var_{running} &= \beta_m * Var_{running} + (1 - \beta_m) * Var_{minibatch} \\ \mu_{running} &= \beta_m * \mu_{running} + (1 - \beta_m) * \mu_{minibatch} \end{aligned}$$

β_m : momentum (hyperparameter)

BN: What do you get?

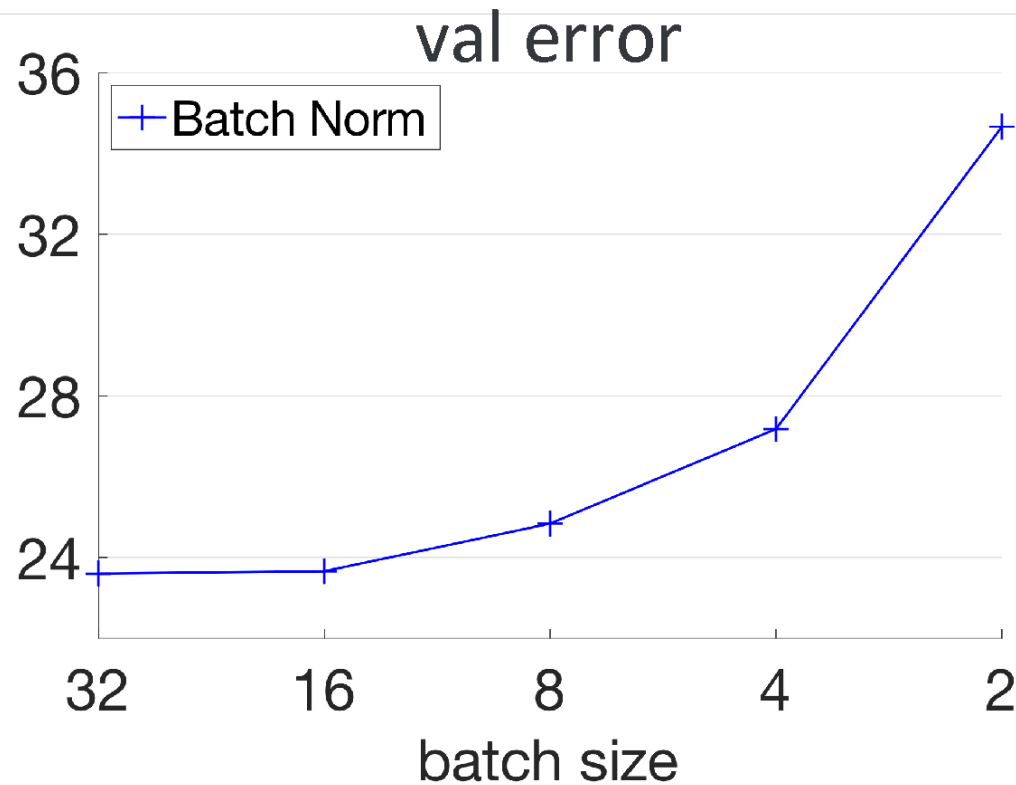
- Very deep nets are much easier to train, more stable gradients
- A much larger range of hyperparameters works similarly when using BN

BN: A Milestone



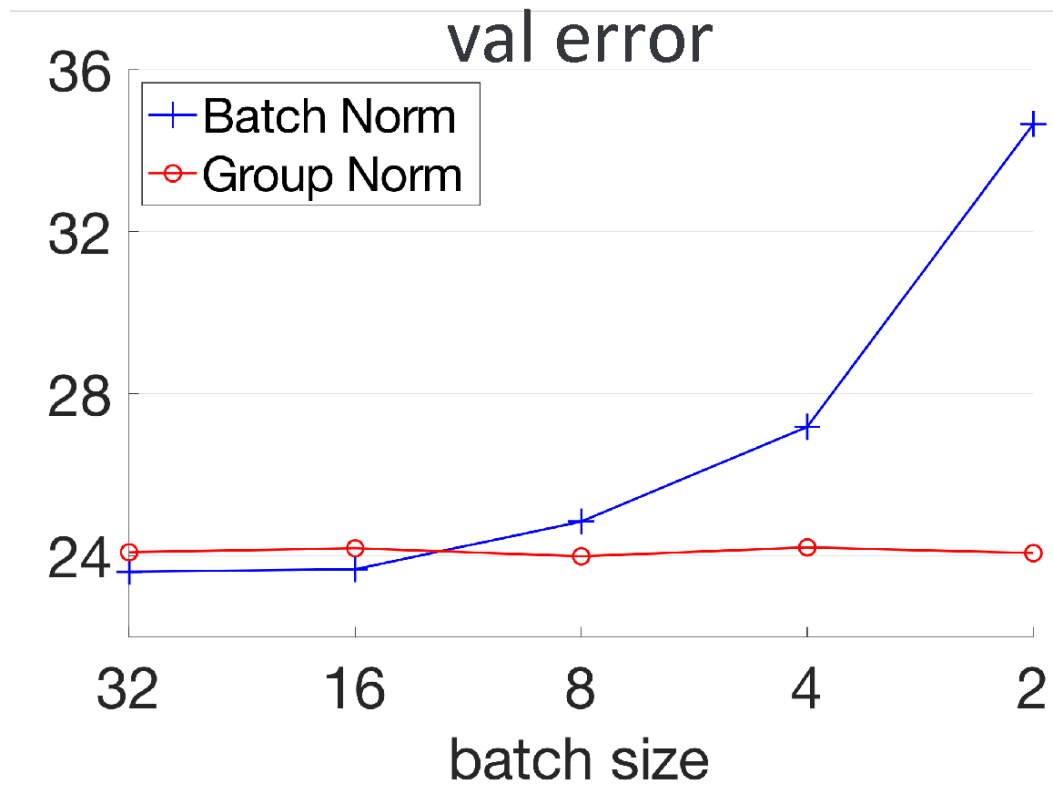
[Wu and He, ECCV'18] Group Normalization

BN: Drawbacks



[Wu and He, ECCV'18] Group Normalization

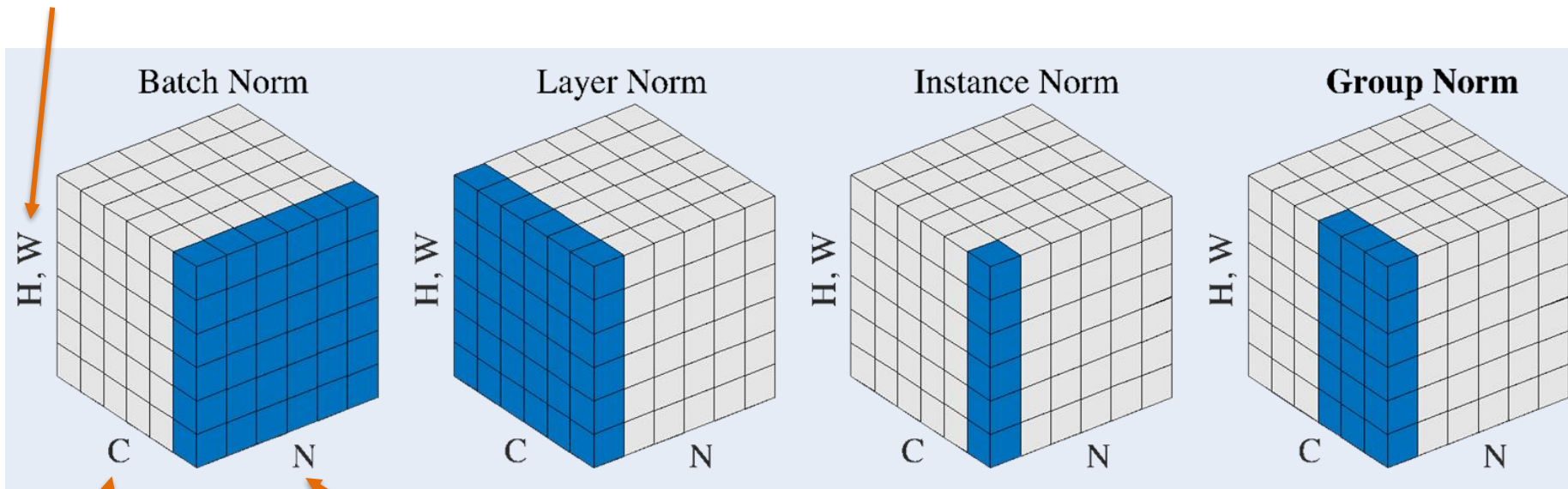
Other Normalizations



[Wu and He, ECCV'18] Group Normalization

Other Normalizations

Image size



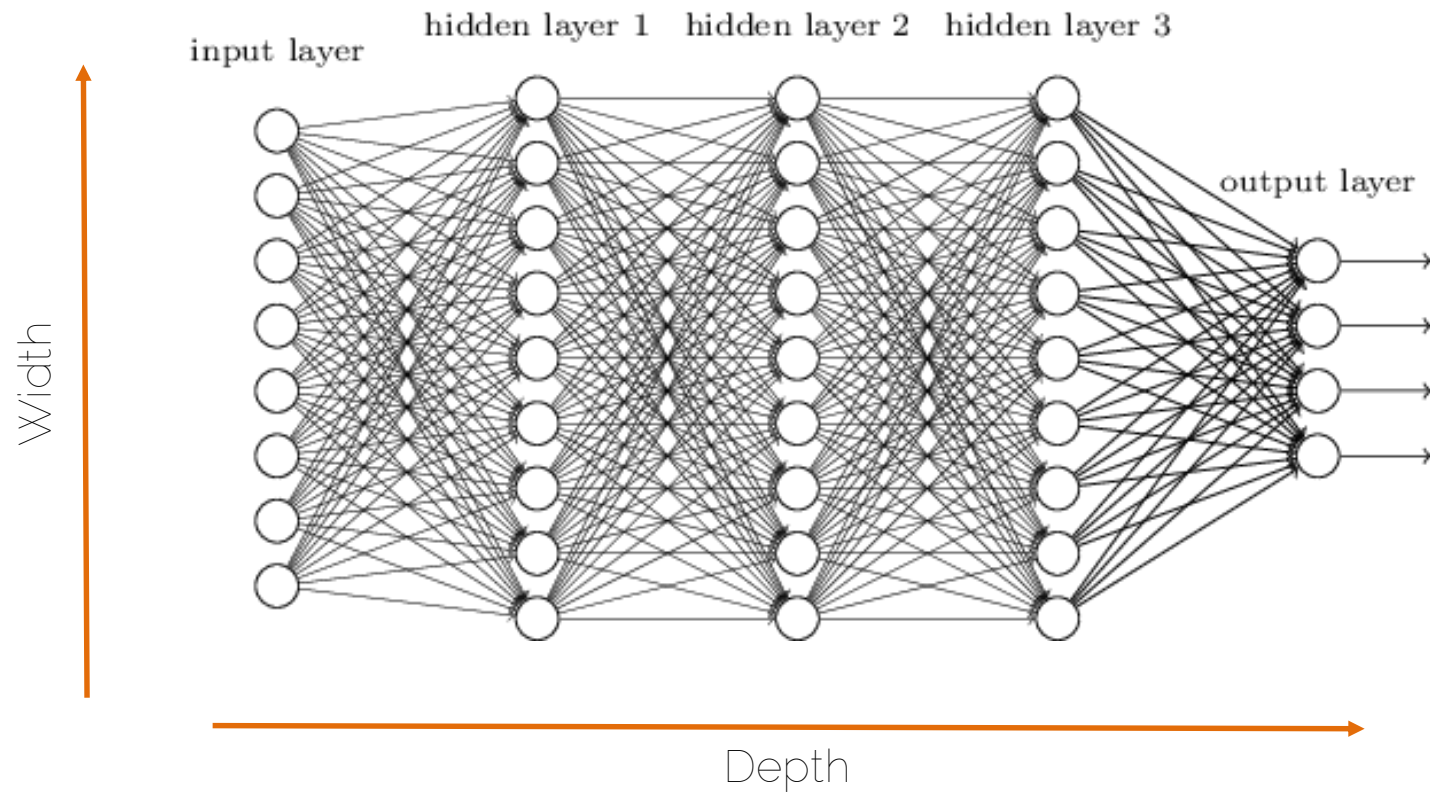
Number of elements in the batch

Number of channels

[Wu and He, ECCV'18] Group Normalization

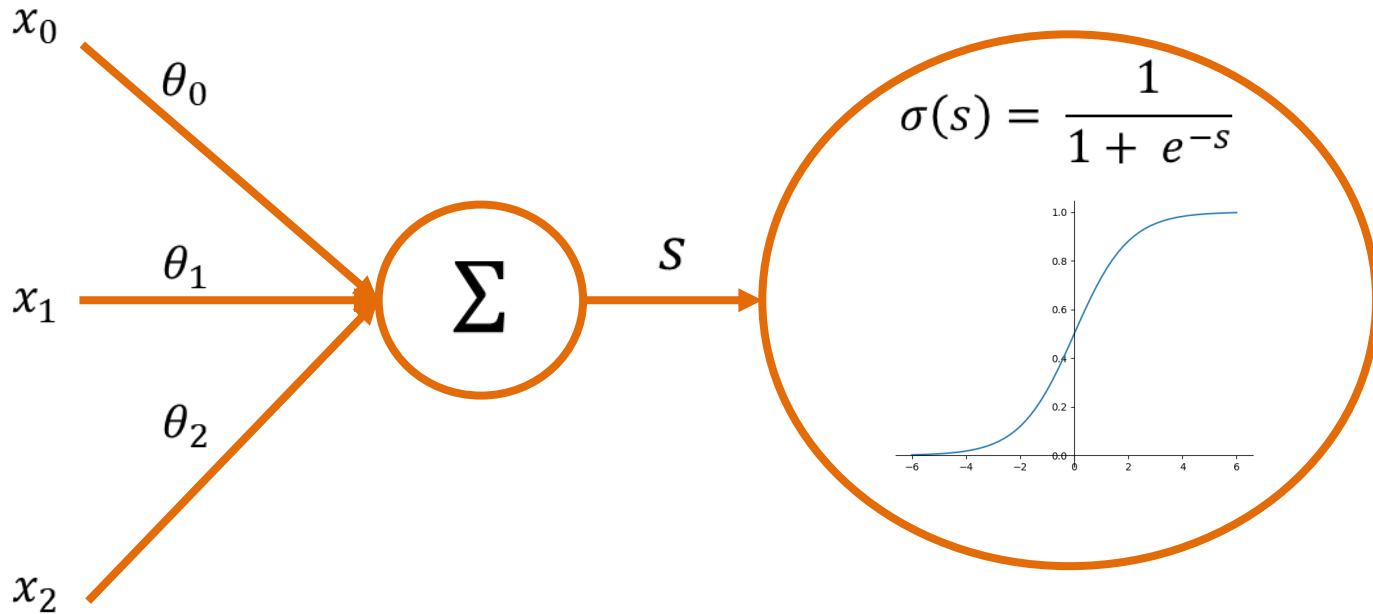
What We Know

What do we know so far?



What do we know so far?

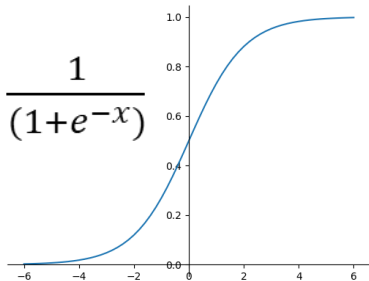
Concept of a 'Neuron'



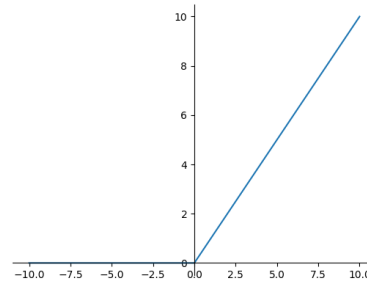
What do we know so far?

Activation Functions (non-linearities)

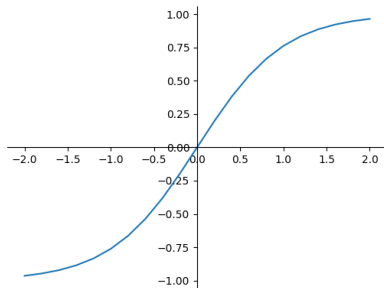
- Sigmoid: $\sigma(x) = \frac{1}{(1+e^{-x})}$



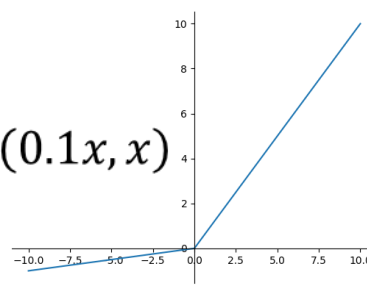
- ReLU: $\max(0, x)$



- TanH: $\tanh(x)$

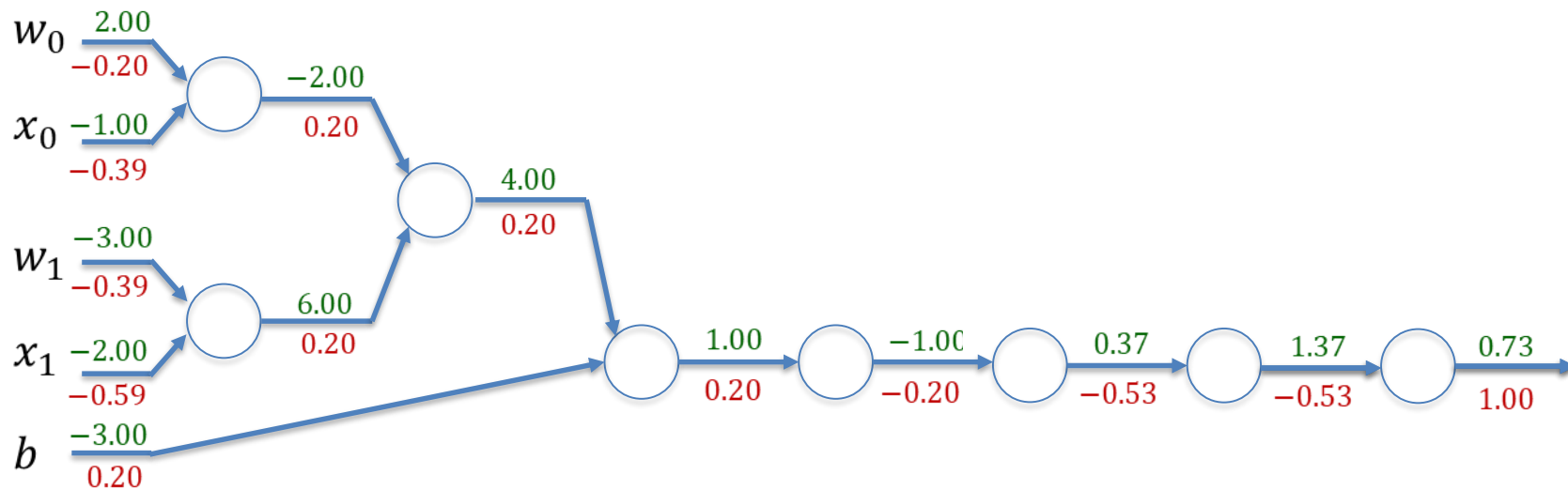


- Leaky ReLU: $\max(0.1x, x)$



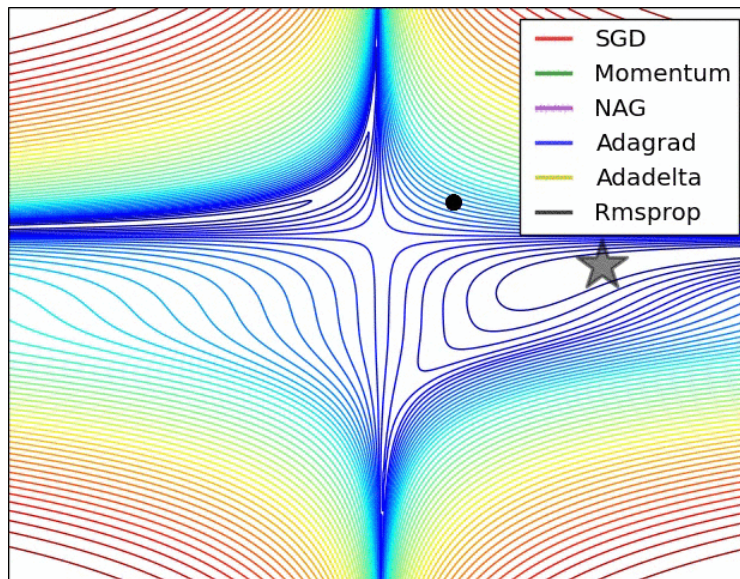
What do we know so far?

Backpropagation



What do we know so far?

SGD Variations (Momentum, etc.)



What do we know so far?

Data Augmentation

a. No augmentation (= 1 image)



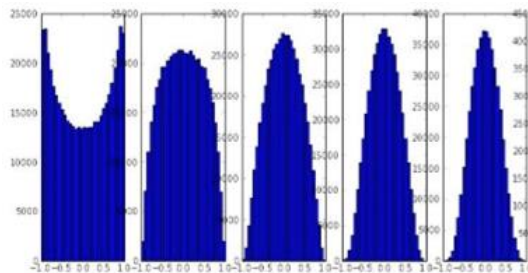
b. Flip augmentation (= 2 images)



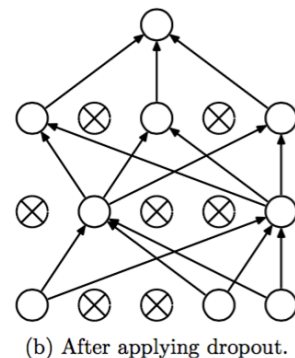
Batch-Norm

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - E[\mathbf{x}^{(k)}]}{\sqrt{\text{Var}[\mathbf{x}^{(k)}]}}$$

Weight Initialization (e.g., Kaiming)



Dropout



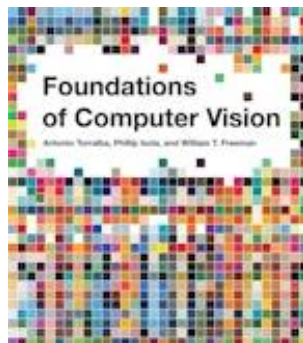
Weight Regularization
e.g., L^2 -reg: $R^2(\mathbf{W}) = \sum_{i=1}^N w_i^2$

Why not simply more layers?

- Neural nets with at least one hidden layer are universal function approximators.
- But generalization is another issue.
- Why not just go deeper and get better?
 - No structure!!
 - It is just brute force!
 - Optimization becomes hard
 - Performance plateaus / drops!
- We need more! More means CNNs, RNNs, and Transformers.

Useful References (Recently Released)

- Foundations of Computer Vision (2024; Torralba, Isola, Freeman)
 - Foundational concepts of computer vision with a machine learning perspective
 - Free online at: <https://visionbook.mit.edu/>



References

- Goodfellow et al. "Deep Learning" (2016),
 - Chapter 6: Deep Feedforward Networks
- Bishop "Pattern Recognition and Machine Learning" (2006),
 - Chapter 5.5: Regularization in Network Nets
- <http://cs231n.github.io/neural-networks-1/>
- <http://cs231n.github.io/neural-networks-2/>
- <http://cs231n.github.io/neural-networks-3/>

See you next week!