

# Neural Network Optimization 1

CS 519: Deep Learning, Winter 2018

Fuxin Li

With materials from Zsolt Kira

# Backpropagation learning of a network

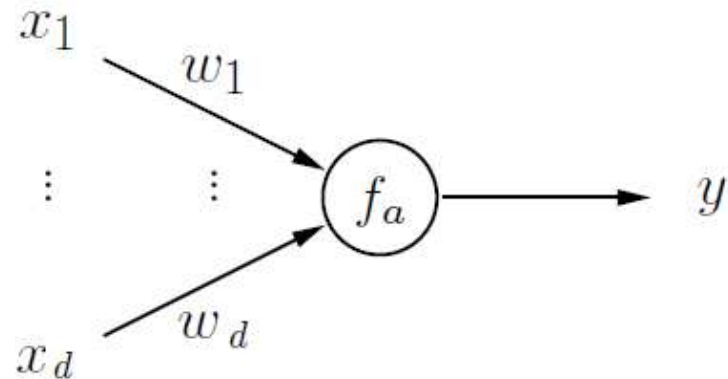
- The algorithm
  - 1. Compute a forward pass on the compute graph (DAG) from the input to all the outputs
  - 2. Backpropagate all the outputs back all the way to the input and collect all gradients  $\mathbf{G}$
  - 3.  $\mathbf{W} = \mathbf{W} - \alpha\mathbf{G}$  for all the weights in all layers

# Modules (Layers)

- Each layer can be seen as a module
- Given input, return
  - Output  $f_a(x)$
  - Network gradient  $\frac{\partial f_a}{\partial x}$
  - Gradient of module parameters  $\frac{\partial f_a}{\partial w_a}$
- During backprop, propagate/update
  - Backpropagated gradient  $\frac{\partial E}{\partial f_a}$

$$\frac{\partial E}{\partial \mathbf{W}_k} = \frac{\partial E}{\partial f_k} g(f_{k-1}(x)) = \frac{\partial E}{\partial f_{k+1}} \frac{\partial f_{k+1}}{\partial f_k} g(f_{k-1}(x))$$

$$\text{where } f_k(x) = \mathbf{w}_k^\top g(f_{k-1}(x)), f_0(x) = x$$



# The abundance of online layers

Core Layers	Convolutional Layers	Advanced Activations Layers
Layer	Convolution1D	LeakyReLU
Methods	Convolution2D	PReLU
Masking	MaxPooling1D	ELU
TimeDistributedMerge	AveragePooling1D	ParametricSoftplus
Merge	MaxPooling2D	ThresholdedLinear
Dropout	AveragePooling2D	ThresholdedReLU
Activation	UpSampling1D	<b>Normalization Layers</b>
Reshape	UpSampling2D	BatchNormalization
Permute	ZeroPadding1D	<b>Embedding Layers</b>
Flatten	ZeroPadding2D	Embedding
RepeatVector	<b>Recurrent Layers</b>	<b>Noise layers</b>
Dense	Recurrent	GaussianNoise
ActivityRegularization	SimpleRNN	GaussianDropout
TimeDistributedDense	GRU	
AutoEncoder	LSTM	
MaxoutDense		
Lambda		
LambdaMerge		
Siamese		
Highway		

# Learning Rates

- Gradient descent is only guaranteed to converge with *small enough* learning rates
  - So that's a sign you should decrease your learning rate if it explodes

- Example:

- $C(w) = \frac{1}{2}w^2$
- Learning rate of  $\alpha = 0.5$
- $\alpha = 1$
- $\alpha = 2$
- $\alpha = 3$

$$w_0 = 1$$

$$w_1 = 1 - 0.5 \cdot 1 \\ = 0.5$$

$$w_2 = 0.5 - 0.5 \cdot 0.5 \\ = 0.25$$

# Weight decay regularization

- Instead of using a normal step, add a

$$\mathbf{G} = \mathbf{G} + \lambda \mathbf{W}$$

- This corresponds to:

$$\min_{\mathbf{W}} \frac{1}{N} \sum_{i=1}^N l(f(x_i; \mathbf{W}), y_i) + \frac{1}{2} \lambda \|\mathbf{W}\|^2$$

- Early stopping as well!
  - Help generalization

$$W = W - 2G$$

$$W = W - 2G - \underbrace{2\lambda}_{\beta} W$$

# Momentum

$$(D_1 = -\alpha G_0) W_1 = W_0 - \alpha G_0$$

- Basic updating equation (with momentum):

$$\begin{aligned} D_0 &= 0 \\ D_{t+1} &= \mu D_t - \alpha G_t \\ W_{t+1} &= W_t + D_{t+1} \end{aligned}$$

- $\mu = 0.6 \sim 0.9$ , a lot of "inertia" in optimization
- Check the previous example with a momentum of 0.5

$$D_2 = -\alpha \mu G_0 - \alpha G_1$$

$$W_2 = W_1 - \alpha \mu G_0 - \alpha G_1$$

$$W_3 = W_2 - \alpha \mu^2 G_0 - \alpha \mu G_1 - \alpha G_2$$

$$\min \frac{1}{2} W^2$$

$$\alpha = 0.5$$

$$W_0 = 1$$

$$W_1 = 0.5$$

$$W_2 = -0.25$$

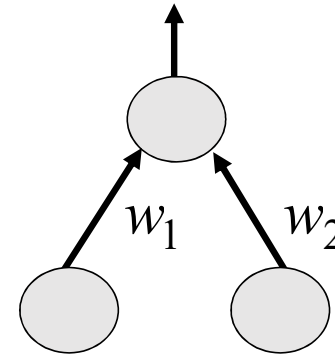
$$- \alpha G_2$$

# Normalization

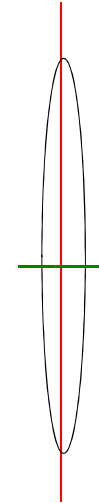
- Each component to 0 mean, 1 standard deviation
- For ease of L2 regularization + optimization convergence rates

0.1, 10  
0.1, -10

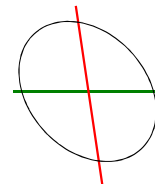
1, 1  
1, -1



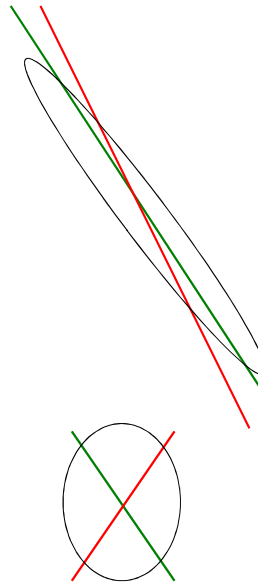
color indicates training case



101, 101  
101, 99



1, 1  
1, -1





# Computing the energy function and gradient

- Usual ERM energy function

$$\min_W E(f) = \sum_{i=1}^n L(f(x_i; W), y_i)$$
$$\nabla_W E = \sum_{i=1}^n \frac{\partial L(f(x_i; W), y_i)}{\partial W}$$

- One problem:
  - Very slow to compute when  $n$  is large
  - One gradient step takes long time!
  - Approximate?

# Stochastic Mini-batch Approximation

$$\min_W E(f) = \sum_{i=1}^n L(f(x_i; W), y_i)$$

$$\nabla_W E = \sum_{i=1}^n \frac{\partial L(f(x_i; W), y_i)}{\partial W}$$

$$\nabla_W \tilde{E} \approx \sum_{i \in N_m} \frac{\partial L(f(x_i; W), y_i)}{\partial W}$$

$$N_m \subset \{1, \dots, n\}$$

- Ensure the expectation is the same

$$\mathbb{E}(\nabla_W \tilde{E}) = \nabla_W E$$

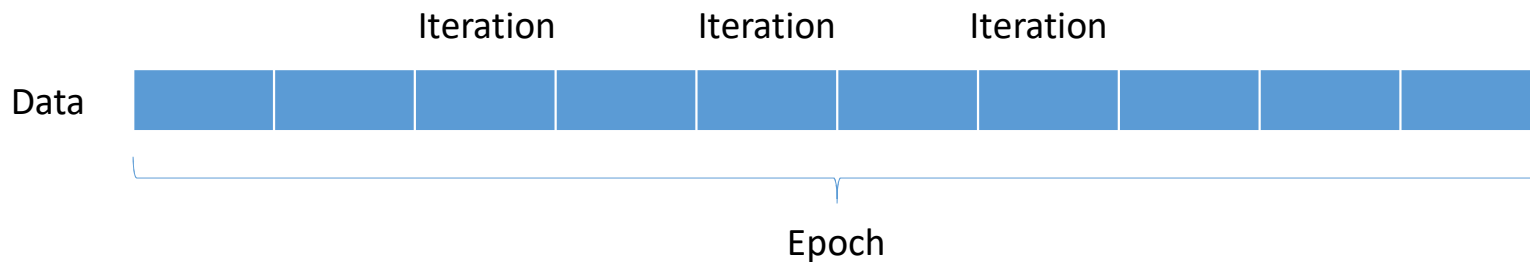
- Uniformly sample every time
  - Sample how many? 1 (SGD) – 256 (Mini-batch SGD)
  - Common mini-batch size is 32-256
  - In practice: dependent on GPU memory size

$$\mathbb{E} \left[ \frac{\partial L(f(x_i; W), y_i)}{\partial W} \right]$$

# In Practice

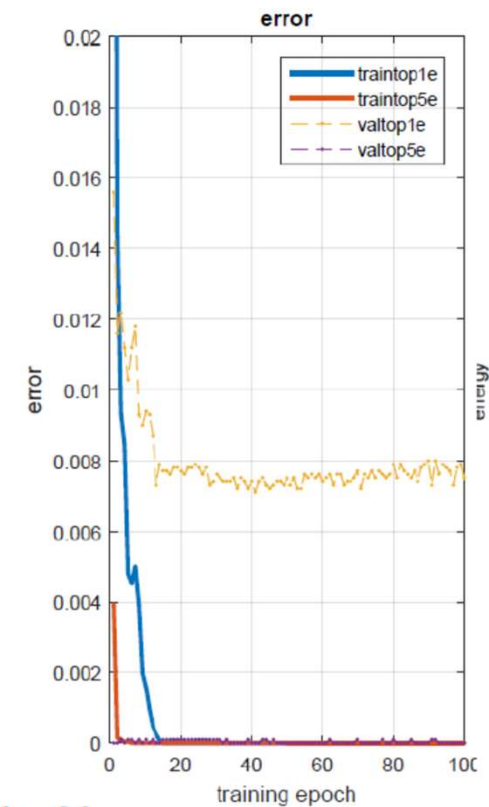
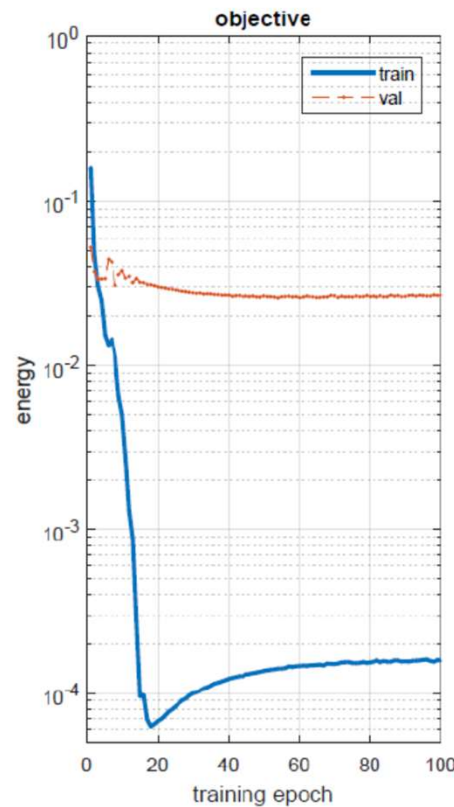
$(n-2 \ 3 \ 5 \ 7)$   
 $(8 \ 9 \ 2 \ 1)$

- Randomly re-arrange the input examples
- Use a fixed order on the input examples
- Define an *iteration* to be every time the gradient is computed
- An *epoch* to be every time that all the input examples is looped through once



# A practical run of training a neural network

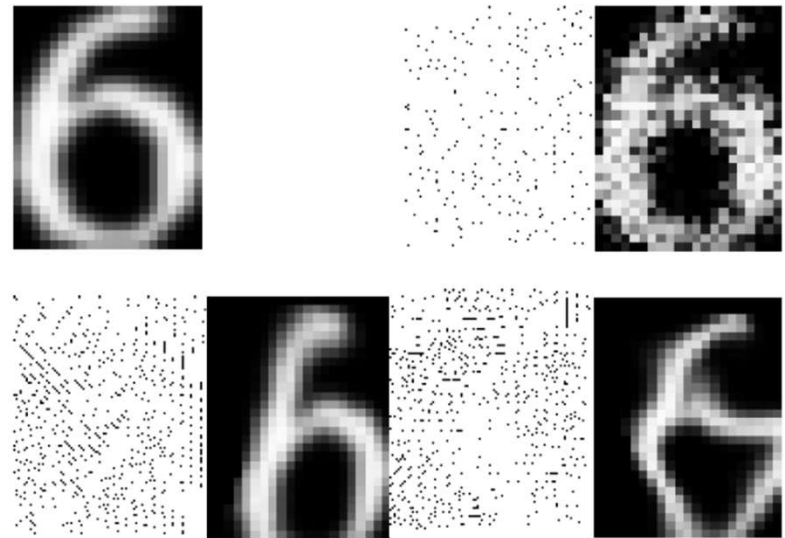
- Check:
  - Energy
  - Training error
  - Validation error



M = 64

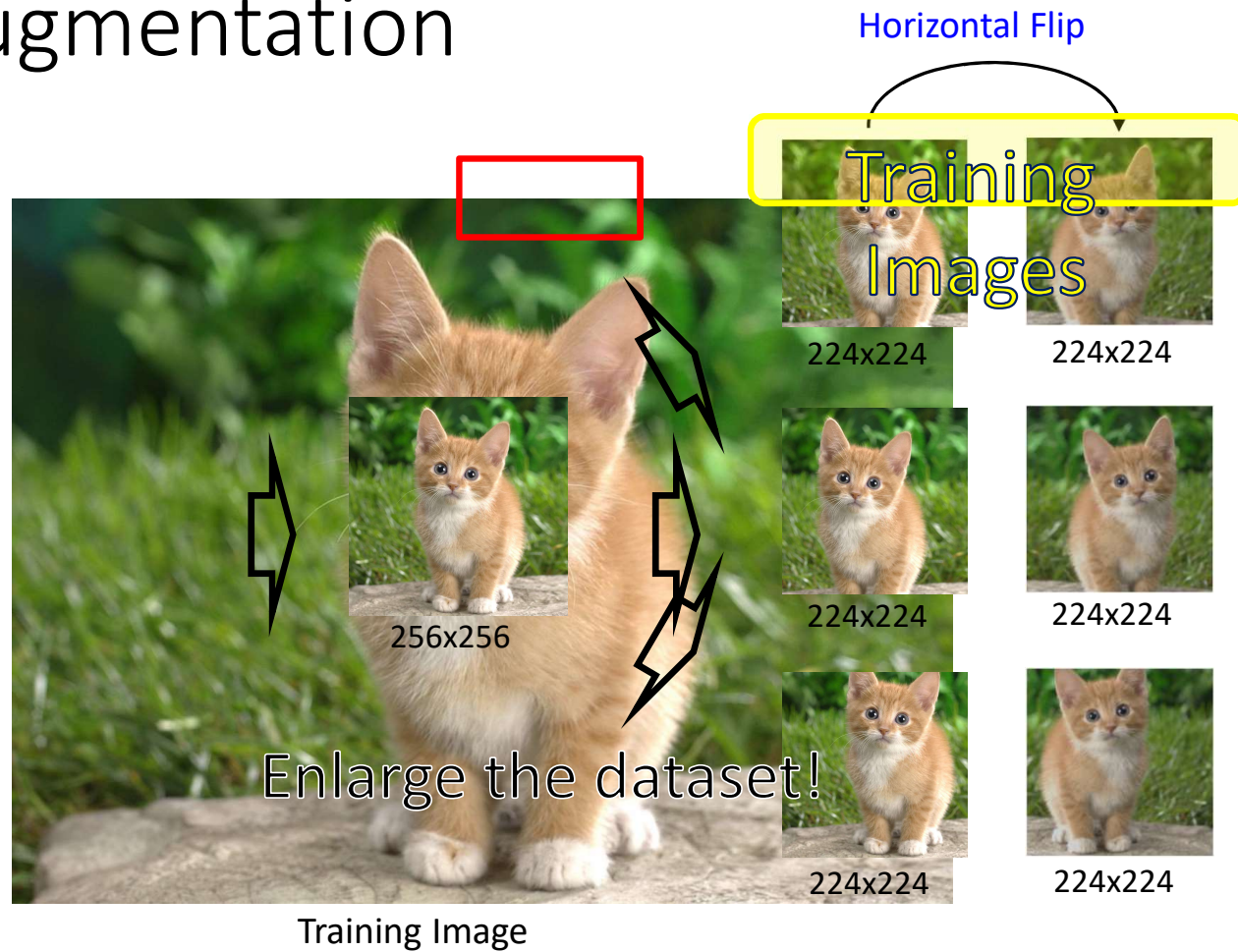
# Data Augmentation

- Create artificial data to increase the size of the dataset
- Example: Elastic deformations on MNIST



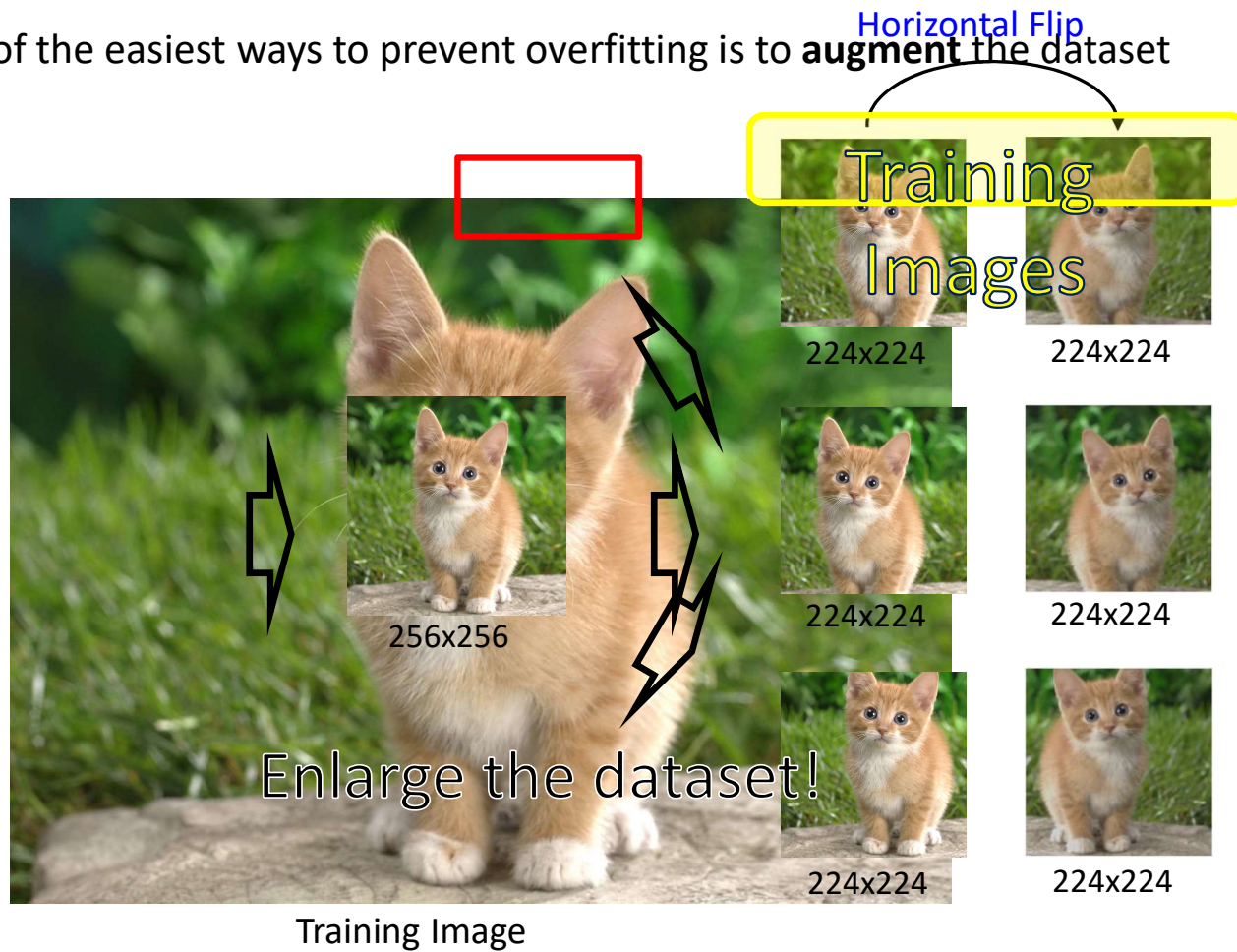
**Figure 2. Top left: Original image. Right and bottom: Pairs of displacement fields with various smoothing, and resulting images when displacement fields are applied to the original image.**

# Data Augmentation



# Data Augmentation

- One of the easiest ways to prevent overfitting is to **augment** the dataset





# CIFAR-10 dataset

- 60,000 images in 10 classes
  - 50,000 training
  - 10,000 test
- Designed to mimic MNIST
- 32x32
- Assignment (will post on Canvas with more explicity):
  - Write your own backpropagation NN to test on CIFAR-10

**airplane**



**automobile**



**bird**



**cat**



**deer**



**dog**



**frog**



**horse**



**ship**



**truck**

