

• Recall NN: CNN

- activation map
- objective: learn the weights by: optimization
- mini-batch Stochastic Gradient Descent
- Loop:
  1. sample a batch of data
  2. forward prop it thru the graph (network) get loss
  3. backprop to calculate the gradients
  4. update parameters using gradients

- in this lecture:
  1. one-time setup: activation functions, preprocessing, weight initialization, gradient checking
  2. training dynamics: better learning process, parameters update, hyperparameter optimization
  3. evaluation model ensembles

1. activation functions

Activation Functions

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

**Leaky ReLU**  
 $\max(0, Lx)$

**tanh**  
 $\tanh(x)$

**ReLU**  
 $\max(0, x)$

**Maxout**  
 $\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**  
 $\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

- Sigmoid  $\sigma(x) = \frac{1}{1+e^{-x}}$
- squashes numbers to range [0, 1]
- historically popular since they have nice interpretation as a saturating "firing rate" of a neuron
- problem:
  1. saturated neurons "kill" the gradients

What happens when  $x = -10$ ?  
 What happens when  $x = 0$ ?  
 What happens when  $x = 10$ ?

2. sigmoid outputs are not zero-centered

Consider what happens when the input to a neuron (x) is always positive:

input =  $w_1 x_1, w_2 x_2$

What can we say about the gradients on  $w$ ? is always all positive or all negative.

Consider what happens when the input to a neuron is always positive...

What can we say about the gradients on  $w$ ? Always all positive or all negative. (this is also why you want zero-mean data!)

3. exp() is a bit computation expensive

• tanh

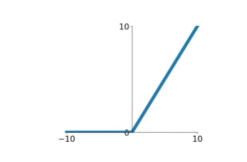
Activation Functions

**tanh(x)**

- Squashes numbers to range [-1, 1]
- zero centered (nice)
- still kills gradients when saturated :(

• ReLU

Activation Functions



ReLU (Rectified Linear Unit)

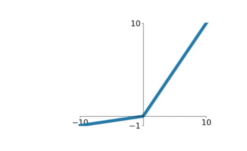


ReLU. The Rectified Linear Unit has become very popular in the last few years. It computes the function  $f(x) = \max(0, x)$ . In other words, the activation is simply thresholded at zero (see image above on the left). There are several pros and cons to using the ReLUs:

- (\*) It was found to greatly accelerate (e.g. a factor of 6 in Krizhevsky et al.) the convergence of stochastic gradient descent compared to the sigmoid/tanh functions. It is argued that this is due to its linear, non-saturating form.
- (\*) Compared to tanh/sigmoid neurons that involve expensive operations (exponentials, etc.), the ReLU can be implemented by simply thresholding a matrix of activations at zero.
- (\*) Unfortunately, ReLU units can be fragile during training and can "die". For example, a large gradient flowing through a ReLU neuron could cause the weights to update in such a way that the neuron will never activate on any datapoint again. If this happens, then the gradient flowing through the unit will forever be zero from that point on. That is, the ReLU units can irreversibly die during training since they can get knocked off of the data manifold. For example, you may find that as much as 40% of your network can be "dead" (i.e. neurons that never activate across the entire training dataset) if the learning rate is set too high. With a proper setting of the learning rate this is less frequently an issue.

• Leaky ReLU

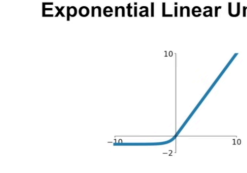
Activation Functions



**Leaky ReLU**  
 $f(x) = \max(0.01x, x)$

• ELU

Activation Functions



**Exponential Linear Units (ELU)**

$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$

• Maxout

Maxout "Neuron"

- Does not have the basic form of dot product → nonlinearity
  - Generalizes ReLU and Leaky ReLU
  - Linear Regime! Does not saturate! Does not die!
- $\max(w_1^T x + b_1, w_2^T x + b_2)$

Problem: doubles the number of parameters/neuron :(

TLDR:

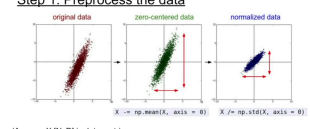
- △ use ReLU (learning rates matter!)
- △ try Leaky ReLU/Maxout/ELU
- △ try tanh yet don't expect much
- △ don't use sigmoid!!!

- Computes  $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid
- NOT zero-centered
- what is the gradient
  - $x = -10 \Rightarrow 0$
  - $x = 0 \hat{=} 0$
  - $x = 10 \Rightarrow 1$

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- will not "die".
- parametric rectifier (PReLU)
- $S(x) = \max(\alpha x, x)$

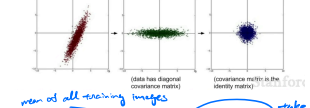
Computation requires exp()

• step 1. Preprocess the data



(Assume  $X$  [NxD] is data matrix, each example in a row)

Step 1: Preprocess the data



In practice, you may also see PCA and Whitening of the data

TLDR: In practice for images, center only

e.g. consider CIFAR-10 example with [32,32,3] images (mean image = [32,32,3] array)

- Subtract the mean image (e.g. AlexNet)
- Subtract per-channel mean (e.g. VGGNet) (mean along each channel = 3 numbers)

- weight initialization
  - Q: what will happen if  $\text{tr}(W) = 0$
  - A: neurons will learn the same thing
  - △ method 1: small random numbers (e.g.  $w \sim 0.1 \cdot \text{np.random.randn}(D)$ )
  - △ method 2: larger random numbers. → will saturate
  - △ method 3: Xavier initialization

• Batch Normalization

Batch Normalization (Ioffe and Szegedy, 2015)

"you want unit gaussian activations? just make them so."

consider a batch of activations at some layer. To make each dimension unit gaussian, apply:

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

this is a vanilla differentiable function... able to do backprop

Batch Normalization (Ioffe and Szegedy, 2015)

"you want unit gaussian activations? just make them so."

1. compute the empirical mean and variance independently for each dimension.

2. Normalize

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

Batch Normalization (Ioffe and Szegedy, 2015)

Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

Problem: do we necessarily want a unit gaussian input to a tanh layer?

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

Batch Normalization (Ioffe and Szegedy, 2015)

Normalize:

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

And then allow the network to squash the range if it wants to:

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

Note: the network can learn:  $\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$ ,  $\beta^{(k)} = E[x^{(k)}]$  to recover the identity mapping.

Batch Normalization (Ioffe and Szegedy, 2015)

Input: Values of  $x$  over a mini-batch:  $B = [x_1, \dots]$ . Parameters to be learned:  $\gamma, \beta$

Output:  $\hat{y} = \text{BN}_\gamma \beta(x)$

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

- Backstopping the training (each is hyperparameter optimization) itself (after lecture)
- cross-validation

- Training II

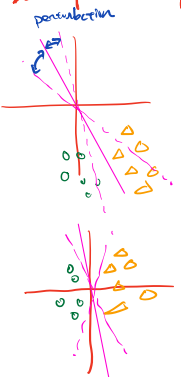
- recall:

① always / or mostly: use ReLU

② weight initialization:

- initialization too small (activations go to zero, gradients also zero, no learning)
- initialization too big: (activations saturate (for tanh) gradients zero, no learning)
- initialization just right (nice distribution of activations at all layers, learning proceeds nicely)

③ Data processing



without normalization:  
classification loss  
very sensitive to  
changes in weight  
matrix: very hard  
to optimize

with normalization:  
less sensitive  
better to optimize

Batch Normalization

Input:  $x : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Learnable params:  
 $\gamma, \beta : D$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Intermediates:  
 $\mu, \sigma : D$   
 $\hat{x} : N \times D$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

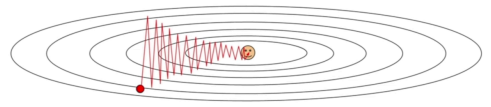
Output:  $y : N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

△ Faster optimization

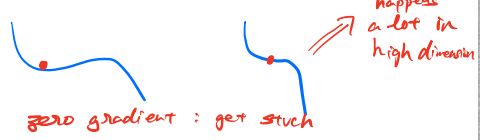
① Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?  
What does gradient descent do?  
Very slow progress along shallow dimension, jitter along steep direction



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large  
even more severe in high dimensions

② local minima // saddle point



②

Optimization: Problems with SGD

Our gradients come from minibatches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$

