

- Lecture 6: Training Neural Networks – Part 1

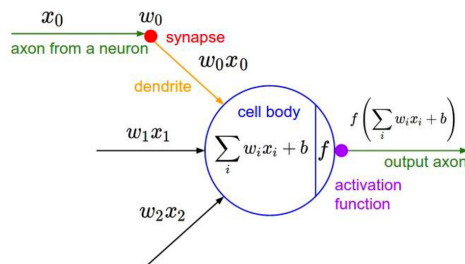
- Training Neural Network: Overview

Overview

1. **One time setup**
activation functions, preprocessing, weight initialization, regularization, gradient checking
2. **Training dynamics**
babysitting the learning process, parameter updates, hyperparameter optimization
3. **Evaluation**
model ensembles

- Activation Functions

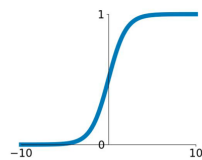
Activation Functions



- Makes neural network Non-Linear

Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$



Sigmoid

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

3 problems:

1. Saturated neurons "kill" the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit compute expensive

- Sigmoid Function

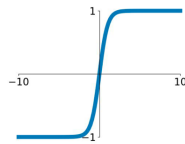
- ◆ Properties

- $(-\infty, \infty) \rightarrow [0, 1]$

- ◆ Problems

- Killed gradients of saturated neurons
- Not zero-centered output
- Exponential Calculation is Expensive

Activation Functions



$\tanh(x)$

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

[LeCun et al., 1991]

■ Tanh Function

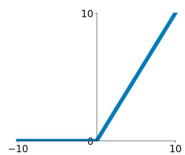
◆ Properties

- $(-\infty, \infty) \rightarrow (-1, 1)$
- Zero-Centered Output

◆ Problems

- Killed gradients of saturated neurons

Activation Functions



ReLU
(Rectified Linear Unit)

- 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 output
- An annoyance:

hint: what is the gradient when $x < 0$?

■ ReLU Function

◆ Properties

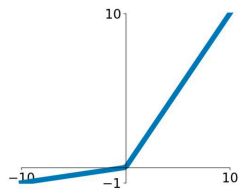
- $(-\infty, \infty) \rightarrow [0, \infty)$
- Does not saturate
- Computational-Efficiency

◆ Problems

- Gradients still die in $(-\infty, 0)$
- Not Zero-Centered Output

Activation Functions

[Mass et al., 2013]
[He et al., 2015]



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into α
(parameter)

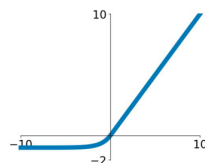
■ Leaky ReLU and PReLU

- ◆ Leaky ReLU: $(-\infty, 0)$ 영역에서 zero output 대신 negative slope 사용
- ◆ PReLU: Learned negative slope version of Leaky ReLU
- ◆ Properties
 - $(-\infty, \infty) \rightarrow (\infty, \infty)$
 - Gradients will not die in $(-\infty, 0)$

Activation Functions

[Clevert et al., 2015]

Exponential Linear Units (ELU)



- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

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

- **Computation requires $\exp()$**

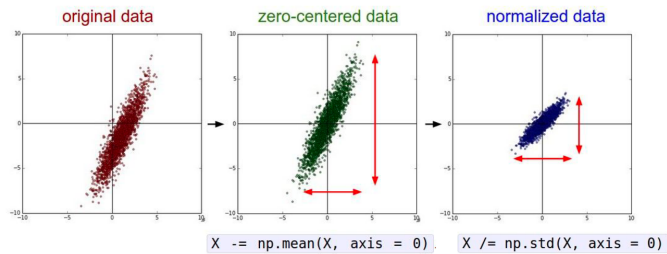
■ ELU Function

- ◆ Properties
 - Closer to Zero-Mean Outputs
- ◆ Problems
 - Costly Computation - Exponential

- Data Preprocessing

-

Step 1: Preprocess the data

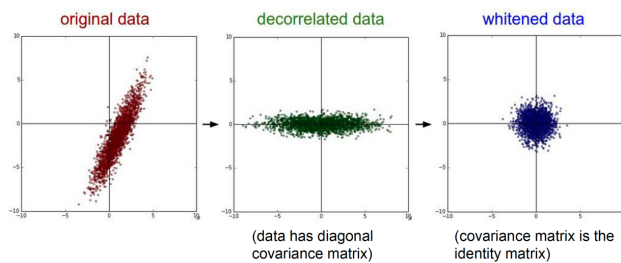


(Assume X [NxD] is data matrix,
each example in a row)

■ Zero-Centering and Normalizing

Step 1: Preprocess the data

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



■ PCA and Whitening

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:

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

this is a vanilla
differentiable function...

Batch Normalization

[Ioffe and Szegedy, 2015]

Normalize:

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

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

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

Note, the network can learn:

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

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

to recover the identity mapping.

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- 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

■ Batch Normalization