



# CS F425: Deep Learning

# 11

## Neural Weight Initialization



**Dr. Kamlesh Tiwari**  
Assistant Professor, Department of CSIS,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Feb 14, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

### Deep Learning Approach

DL is constructing network (neural obviously) of parameterized functional modules and training them from examples using gradient based optimization

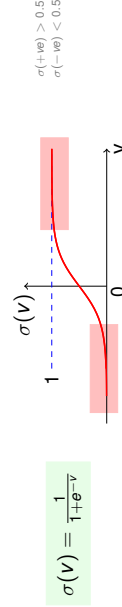
- Yun LeCun

There are many interesting questions

- How many Layers?
- How many units in each layer?
- What should be the learning rate?
- How to initialize the parameters?
- What is right activation function? etc.

### Exploding and Vanishing Gradient

- If weights are **too small**, variance of the input signal starts diminishing as it passes through layers.
- If weights are **too large**, variance get huge amplification.



We want variance of input and output should not differ much (atleast initially)

### One Hot Encoding

How should we encode the output?

- Assume there are five class **cat**, **dog**, **bus**, **plant**, **apple**

Class	Code-1	Code-2
cat	001	10000
dog	010	01000
bus	011	00100
plant	100	00010
apple	101	00001

Consider Code-1:

- Average of cat and bus?
- Determine similarity in  $sim(cat, dog)$  and  $sim(cat, bus)$
- Does the loss for **plant** and **apple** is same?

Issues leading to Code-2 (One Hot Encoding):

- Ordinal relationship should not be present
- Similarity between classes to be avoided (make them orthogonal)
- Loss should not favour (be biased) to any class

### Hyperparameter Tuning

It is an iterative process.



- Where is the nob? Underfitting/Overfitting
- Bad everywhere (high bias), vs on test data (high variance)
- Go towards more complex model, or increase data/regularization, (reduce complexity)

### Xavier Weight Initialization

- Random initial  $W$ , likely tend to exploding and vanishing gradient
- Initialization with zero weights makes things linear (as all neurons become symmetric).
- Standard is to make weights centered around 0 with small variance<sup>1</sup>  $k/S_{l-1}$  where previous layer has  $S_{l-1}$  neurons<sup>2</sup>
- Weights are sampled on normal distribution with variance  $k/S_{l-1}$   $k = 2$  for ReLU and  $k = 1$  for sigmoid/tanh
- This is called **xavier** initialization
- Another possibility is to use variance as  $k/(S_{l-1} + S_l)$  or  $k/(S_{l-1} S_l)$

Apply **gradient clipping** if needed.

<sup>1</sup>variance is average of the squared differences from the mean

<sup>2</sup>Standard deviation becomes  $\sqrt{k/S_{l-1}}$

## Xavier Initialization

Thank You!

Assume data is i.i.d. from normal distribution with  $\mu = 0, \sigma = 1$

- As  $\text{var}[w_i x_i] = E[x_i]^2 \text{var}[w_i] + E[w_i]^2 \text{var}[x_i] + \text{var}[x_i] \text{var}[w_i]$   
since  $E[x_i] = \mu = 0$  therefore  $\text{var}[w_i x_i] = \text{var}[x_i] \text{var}[w_i]$

- Consider a neuron  $y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$

$$\begin{aligned} \text{var}[y] &= \text{var}[w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b] \\ &= \text{var}[w_1] \text{var}[x_1] + \text{var}[w_2] \text{var}[x_2] + \dots + \text{var}[w_n] \text{var}[x_n] + \text{var}[b] \\ &= \text{var}[w_1] \text{var}[x_1] + \text{var}[w_2] \text{var}[x_2] + \dots + \text{var}[w_n] \text{var}[x_n] + 0 \\ &= S_{l-1} \cdot \text{var}[w_l] \text{var}[x_l] \quad n = S_{l-1} \end{aligned}$$

To have same variance in **input** and **output**  $S_{l-1} \cdot \text{var}[w_l] = 1$  or

$$\text{var}[w_l] = 1 / S_{l-1}$$

Thank you very much for your attention!