



CS F425: Deep Learning

12

Regularization For Neural NW



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

Feb 16, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

Recall Regularization for Logistic Regression

- Optimization minimizes the loss $\min_w J(w)$ by adjusting w where

$$J(w) = \frac{1}{m} \sum_{i=1}^m \text{Loss}(\hat{y}^{(i)}, y^{(i)})$$

- Regularization penalizes the large values of w by

$$J(w) = \frac{1}{m} \sum_{i=1}^m \text{Loss}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

where

$$\|w\|_2^2 = \sum_j w_j^2 = w^T w$$

λ being regularization parameter

Weight Decay: Regularization for NN (contd..)

- With regularization term $dw^{[l]} = (\text{from backpropogation}) + \frac{\lambda}{m} w^{[l]}$
- Therefore the update is modified to

$$w^{[l]} = w^{[l]} - \alpha \cdot ((\text{from backpropogation}) + \frac{\lambda}{m} w^{[l]})$$

Which is

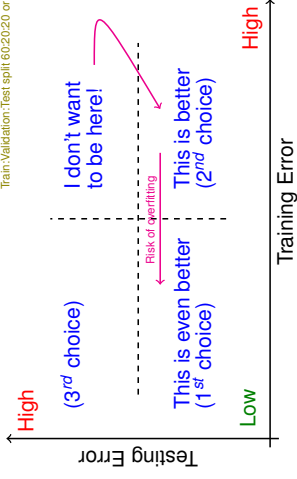
$$w^{[l]} = \left(1 - \frac{\alpha\lambda}{m}\right) w^{[l]} - \alpha \cdot (\text{from backpropogation})$$

- Due to the $(1 - \frac{\alpha\lambda}{m})$ factor, this update method is also called **weight decay**

Our objective here is to penalize the weight matrices being too large.

Which side do you want to be?

Train/Validation/Test split 60/20/20 or 99:1:1



Low training error comes with a risk of overfitting (high variance)

Regularization for Neural Network

- As there could be L layers, each having own set of parameters, so

$$J(w^{[1]}, w^{[2]}, \dots, w^{[L]}) = \frac{1}{m} \sum_{i=1}^m \text{Loss}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

where **frobenius** norm is

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} w_{ij}^2$$

- How you updated the parameter earlier? get $dw^{[l]} = (\text{from backpropogation})$ that is $\frac{\partial J}{\partial w^{[l]}}$ then

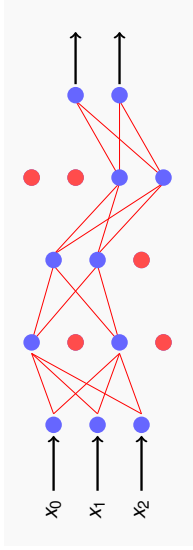
$$w^{[l]} = w^{[l]} - \alpha \cdot dw^{[l]}$$

Penalize the weight matrices from being too large

- With low weight, the connection get weakened so network has effectively less connections and become simpler.
- Another intuition is that, when weights are lower, the output of the units is also lower. Assume activation function be tanh small input values tends to produce linear output.
- So overall n/w tends to becomes linear (more biased).

Dropout Regularization

Shutdown some units **randomly**



- Drop probabilities for different layers may be different
- Now networks cannot rely on any specific connection. So have to give importance to others neuron also
- Issue is that cost function is now not well defined

Deep Learning (CS F.425)

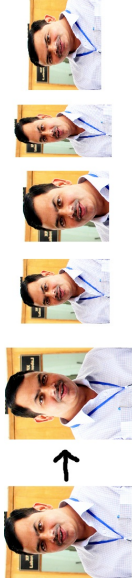
(Tu,Th,Fr,12PM) BITS-Pilani

Lecture-12 (Feb 16, 2023)

7/10

Other Regularization Methods

- Increase the data by **data augmentation**
- Flip, rotate, scale, translate, deform ...
- **Normalize** training examples → to **speed up your training**



Deep Learning (CS F.425)

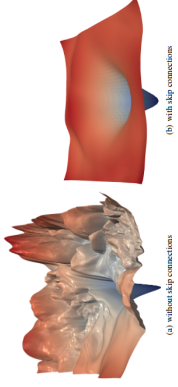
(Tu,Th,Fr,12PM) BITS-Pilani

Lecture-12 (Feb 16, 2023)

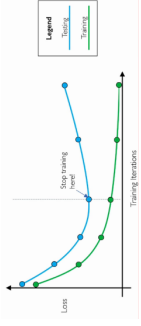
9/10

Regularization

- Dropout¹



- Early stopping



¹ Visualizing the Loss Landscape of Neural Nets, Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, Tom Goldstein, 2018

Deep Learning (CS F.425)

(Tu,Th,Fr,12PM) BITS-Pilani

Lecture-12 (Feb 16, 2023)

8/10

Thank You!

Thank you very much for your attention!

Deep Learning (CS F.425)

(Tu,Th,Fr,12PM) BITS-Pilani

Lecture-12 (Feb 16, 2023)

10/10