

ECE C147/C247, Winter 2022

Department of Electrical and Computer Engineering
University of California, Los Angeles

Midterm Review

Prof. J.C. Kao
TAs: T. Monsoor, T. Wang, P. Lu, Y. Li

3. Training neural networks (Pan)**Covered topics**

1. Gradient problems
2. Weight initialization
3. Batch normalization
4. L2 and L1 regularization
5. Dataset augmentation
6. Multitask and transfer learning
7. Ensemble methods
8. Dropout
9. Other techniques discussed in the lectures

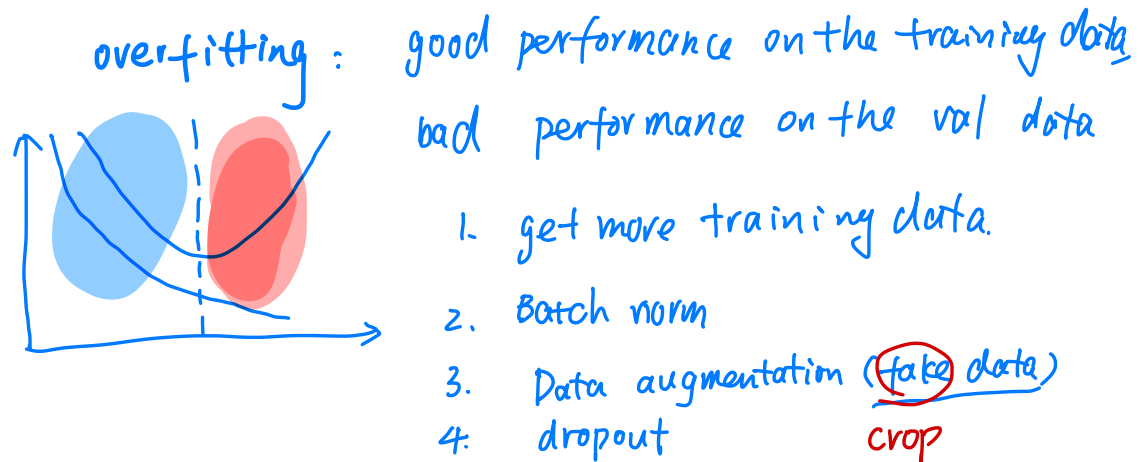
Review materials

1. Lecture slides (formal notes)
2. Homework
3. Discussion handouts
4. Textbook (optional)

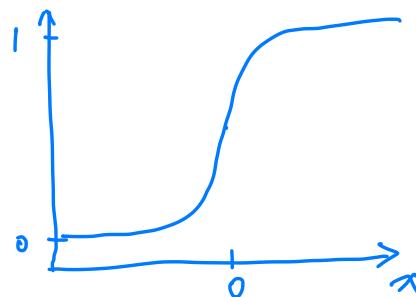
Problem examples

1. Which of the following techniques can be used to reduce overfitting? **Select all that apply.**

- ✓ A. Batch normalization
- ✓ B. Data augmentation
- ✓ C. Dropout
- ✗ D. Use less training data
- ✗ E. None of the above



2. In homework 3, assume you are implementing a fully connected neural network using the sigmoid activation function. Is this a good idea to initialize the weights with large positive numbers? Explain why or why not.

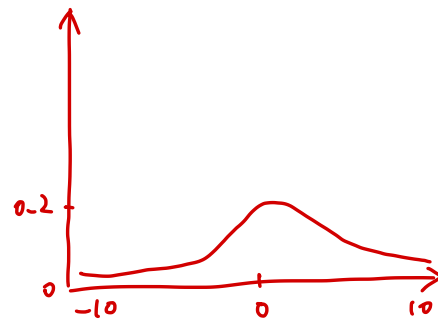


$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$f(w) = \sigma(w^T x + b)$$

$$z = w^T x + b,$$

$$\frac{\partial f(w)}{\partial w} = \frac{\partial z}{\partial w} \frac{\partial f(z)}{\partial z} = \frac{\sigma(z)(1-\sigma(z))}{0} \times \rightarrow 0$$



$$\frac{d\sigma(x)}{dx} = \sigma(x)(1-\sigma(x))$$

x : large

$$\sigma(x) \rightarrow 0 \text{ or } 1$$

$$\frac{d\sigma(x)}{dx} \rightarrow 0.$$

vanishing gradient = zero gradients = no learning