

# SMAI-M20-L30:Neural Network Learning

C. V. Jawahar

IIIT Hyderabad

October 26, 2020

# Announcements

- Quiz 2: Tentatively for next week. (say wed.)
- Refresh your programming and Jupyter notebooks.

# Class Review

SVM and Kernels: Revisiting the illustrative problems we already solved.

- if samples scale in magnitude what happens to the solutions?
- if more samples come around a sample, what happens to the number of support vectors and  $\alpha$ ?
- ExOR and Parity problems
- Does addition of new samples lead to new SVs?



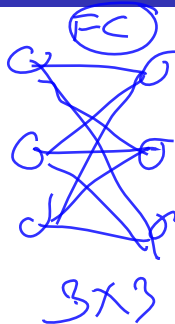


# Recap:

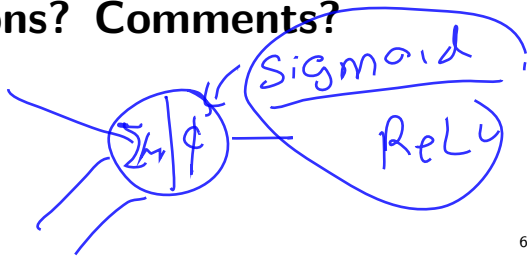
- **Supervised Learning:** Formulation, Conceptual Issues, Concerns etc.
- **Classifiers:** (i) Nearest Neighbour, (ii) Notion of a Linear Classifier (iii) Perceptrons (iv) Bayesian Optimal Classifier (v) Logistic Regression (vi) Multiclass classification architectures (v) SVMs (hard margin, soft margin, kernel)
- **Dimensionality Reduction and Applications:** (i) Feature Selection and Extraction (ii) PCA (iii) LDA (iv) Eigen face
- **Matrix Factorization and Applications:** (i) SVD, (ii) Eigen Decomposition (iii) Matrix Completion (iv) LSI (v) Recommendations
- **Neural Network Architectures and Learning** (i) Neuron model, Single Layer Perceptrons (ii) SLP (iii) MLP. (iv) Backpropagation
- **Other Topics:**
  - Linear Regression
  - Probabilistic View, Bayesian View, MLE
  - Gradient Descent: Stochastic and Batch GD
  - Loss Functions and Optimization
  - Eigen Vector based optimization
  - Kernel Functions and Kernel Matrix

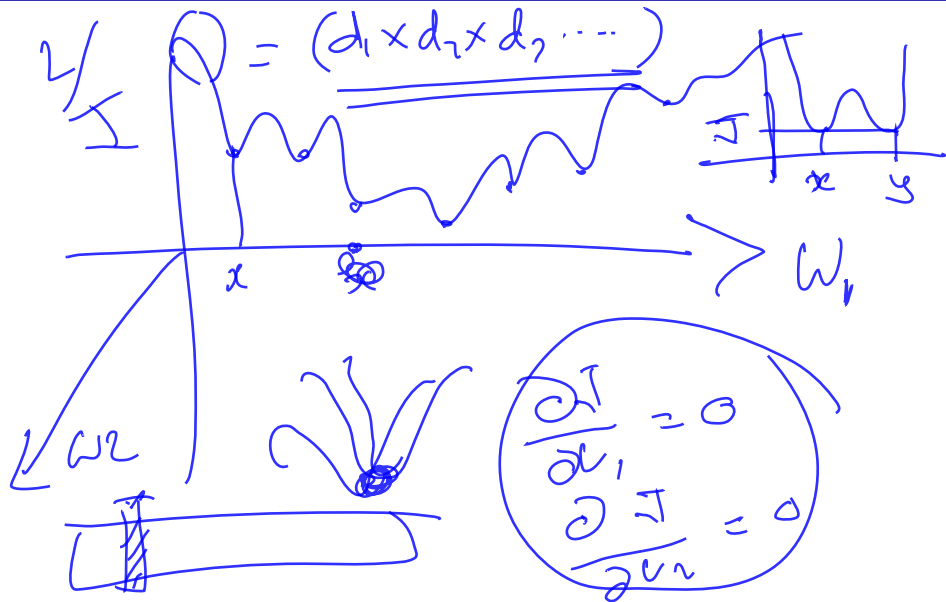
# This Lecture:

- ① Activations in Neural Networks
  - ① Sigmoid and Tanh
  - ② ReLU and L-ReLU
- ② Learning in NNs
  - ① Local minima and non-convexity
  - ② Plateaus and saddle points
  - ③ Vanishing and Exploding Gradients

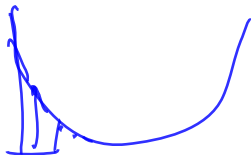


**Questions? Comments?**





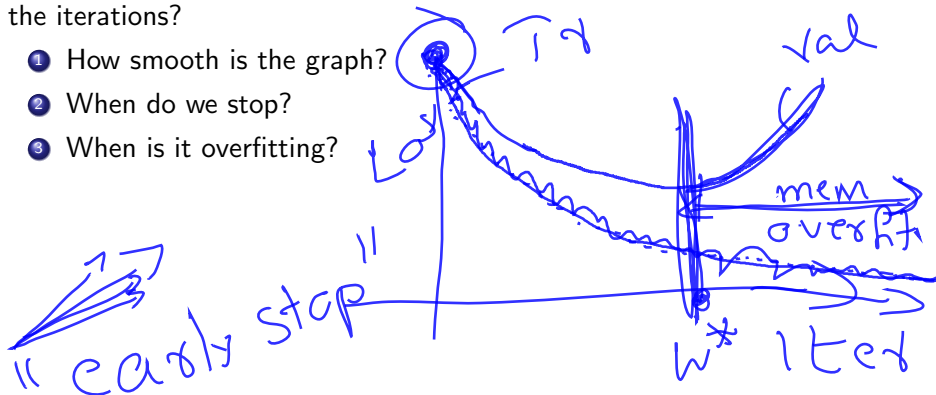
# Discussions Point - I



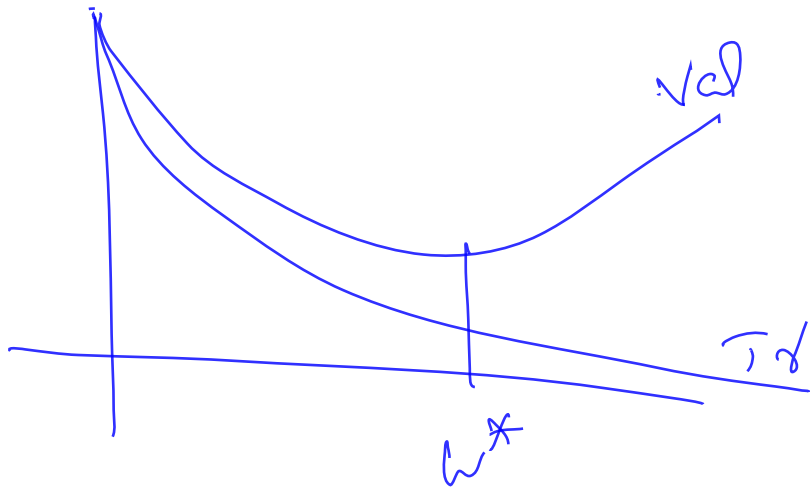
$$\frac{1}{n} \sum_{i=1}^n$$

What happens to the NN Loss graphs on “Train” and “Val/Test” during the iterations?

- 1 How smooth is the graph?
- 2 When do we stop?
- 3 When is it overfitting?











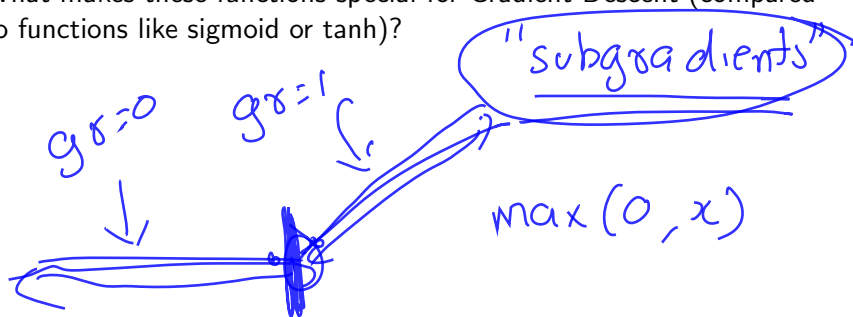


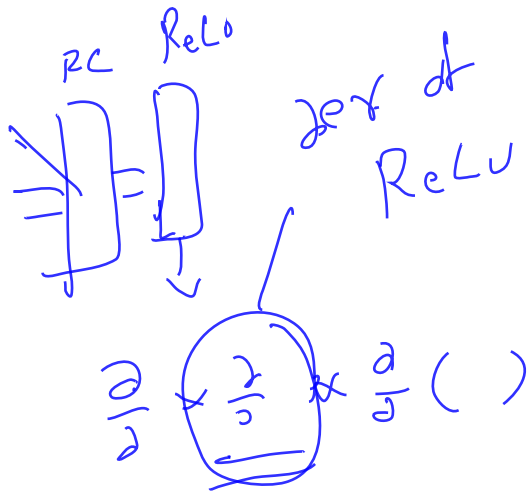
## Discussions Point -II



We know the functions like ReLu and Hinge Loss used in NN/ML.

- What makes these functions special for Gradient Descent (compared to functions like sigmoid or tanh)?







## Discussion Point - III

(Advanced; Out of Syllabus!!) We know that a new kernel can be defined in terms of existing kernels:

$$\sum_{i=1}^K \alpha_i \kappa(\cdot, \cdot)$$

Then why don't we formulate the overall learning problem in SVM, including that of learning these  $\alpha_i$

- 1 Discuss why it is a good idea?
- 2 How do we use it for “fusing” different features?
- 3 Why do we limit to  $\sum$ ?

See some of the works relevant<sup>1</sup> and <sup>2</sup>. Read later.

---

<sup>1</sup><http://manikvarma.org/pubs/varma07c.pdf>

<sup>2</sup><https://cvit.iiit.ac.in/images/ConferencePapers/2009/Rakesh09More.pdf>







## Discussion Point - IV

(Advanced; Out of Syllabus!!)

We know that linear SVMs are superefficient (compared to K-SVMs). Can we find a  $\phi()$  corresponding to a Kernel and solve the problem as

$$\mathbf{w}^T \phi(x)$$

Indeed, this may become difficult for many kernels (eg. RBFs). **why?** Can we find a finite dimensional approximation of  $\phi()$ ? How does it help in speeding up SVM with no major reduction in accuracy? read <sup>3</sup> and <sup>4</sup> later.

- 1 Discuss why it is a good idea?
- 2 Suggest an application where speed matters (eg. in the reference is that of object detection).

---

<sup>3</sup><https://cvit.iiit.ac.in/images/ConferencePapers/2010/Sreekanth10Generalized.pdf>

<sup>4</sup><https://www.robots.ox.ac.uk/vgg/publications/2011/Vedaldi11/vedaldi11.pdf>





# What Next:?

- ① NN Architectures and NN Learning
- ② Programming for Deep Learning.