

SMAI-M20-L25: Nonlinear methods: SVM, Kernels and MLP

C. V. Jawahar

IIIT Hyderabad

October 12, 2020

Class Review

$$[\phi(x_1) \dots \phi(x_n)] \quad (x) \rightarrow \phi(x)$$

Questions is based on the brief review of Kernels you had seen at:
<https://www.dropbox.com/s/qryziuo3u143q5e/KERNEL-REVIEW.pdf?dl=0> (shared in the class last week).

(see specially Section 4 for today)

$$[x_1 \dots x_n]$$

- 1 Relationship between Data Matrix and Kernel Matrix.
- 2 How computations such as (i) centering (ii) distance (iii) normalization can be done in the feature space ($\phi()$).

$$\underline{\underline{K}} = [\phi(x) - \phi(y)]^T [\phi(x) - \phi(y)]$$

pegsos

$w^T x$ @ params

~~$w^T \phi(x)$ / many para~~ X

$\sum_{i=1}^N \phi(x_i) k(x, x_i)$ \approx N parameter

\bar{e} :

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
- **Other Topics:**
 - Linear Regression
 - Probabilistic View, Bayesian View, MLE
 - Gradient Descent: Stochastic and Batch GD
 - Loss Functions and Optimization
 - Eigen Vector based optimization
 - Neuron model, Single Layer Perceptrons
 - Kernel Functions and Kernel Matrix

$$K(x, y) = \underline{\underline{(x^T \cdot y)}}$$

This Lecture:

$q(w, x)$ $f(x)$ nonlinear

1 Kernel SVMs:

- What happens during training?
- What happens during testing?

Find α

α , SVs and
order

2 MLP:

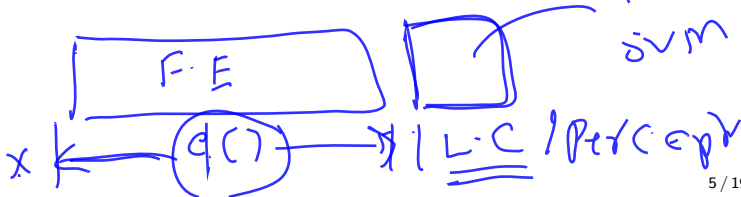
- Introduction to a typical MLP architectures
- Appreciating "Deep MLP" (aka Deep Neural Network) as feature transformation followed by a classification.

F.C

Questions? Comments?

D-T
svm

MLP:



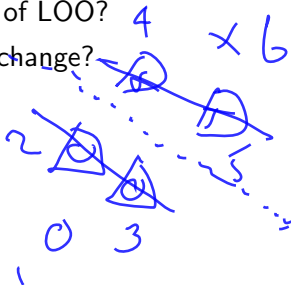
Discussions Point -I

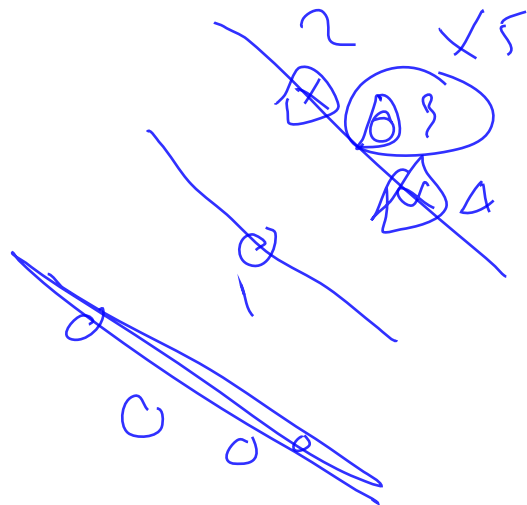
"Leave One Out"

A typical SVM training leaves you with the set of Support Vectors (SV) and the corresponding α s (for all others α s are zero).

- 1 Comment: "The ratio $\frac{|SV|}{N}$ gives us an idea of the error rate." Why? Argue?
- 2 Is this an upper bound or lower bound of LOO?
- 3 If we change the kernel, does the SVs change?

$$\frac{4}{6}$$





$x_1 \dots x_N$

$$\boxed{\begin{matrix} \{x_1, \dots\} = \text{SV} \\ \{\alpha_i\} \end{matrix}}$$

$$\begin{aligned} & \frac{\omega^T x}{100} \quad o(d) = 1 \\ & \sum_{i=1} \alpha_i y_i k(x_i, x) \\ & \quad \quad \quad \downarrow \\ & \quad \quad \quad (x_N^T x)^2 \\ & \quad \quad \quad o(d) \end{aligned}$$

Discussions Point -II

Starting from the Lagrangian

$$L(\mathbf{w}, b, \alpha, \xi) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{2} \sum_{i=1}^N \xi_i^2 - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i]$$

i.e., Derive the dual function for L2 SVM

(Write on a paper, complete and submit later.)

$$\frac{\partial L}{\partial \xi_1} = 0$$

$$\frac{C}{2} \cdot 2 \xi_1 - \alpha_1 = 0$$
$$C \xi_1 - \alpha_1 = 0$$

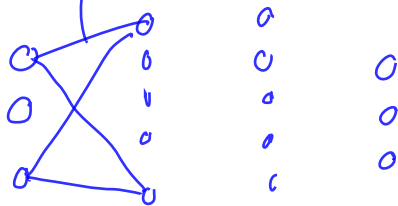
$$C \varepsilon_i - d_i = 0 \quad \forall i$$

$$\begin{aligned} \max_d \quad & \sum_{i=1}^N d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j \\ & - \frac{1}{2} \sum_{i=1}^N \frac{d_i^2}{C} \end{aligned}$$

Discussion Point - III



An MLP has 3 inputs, two hidden layers each of 5 neurons, and three output neurons. No bias anywhere. Hidden neurons have sigmoid activations. Output neurons have linear activations. How many parameters are there to learn?



how many w?

$$3 \times 5 + 5 \times 5 + 5 \times 3 = \underline{\underline{55}}$$

What Next:? (next few)

- ① SVMs and Kernels (winding up)
- ② MLP and Backpropagation