

# Machine Learning

## Week 9: Support Vector Machines

Maheesan Niranjana

School of Electronics and Computer Science  
University of Southampton

Autumn Semester 2017/18

## Bias and Variance in Estimation

Decompose (generalization) error into two terms

- We are interested in training a model and testing it on unseen data
  - Quantify generalization over the space of data
  - Dataset we have is just one realization of the underlying process.

- Truth:  $f(\mathbf{x})$ , Estimated function:  $f(\mathbf{x}|\mathbf{w})$

- Generalization Error:  $E_g = \left\langle \{f(\mathbf{x}) - f(\mathbf{x}|\mathbf{w})\}^2 \right\rangle_{\mathcal{D}, \mathbf{x}}$

Over all space  $\mathbf{x}$  and datasets  $\mathcal{D}$  ( $\langle \cdot \rangle$  denotes expectation.)

Algebra: add and subtract a term and expand out...

$$\left\langle \{f(\mathbf{x}) - \langle f(\mathbf{x}|\mathbf{w}) \rangle_{\mathcal{D}} + \langle f(\mathbf{x}|\mathbf{w}) \rangle_{\mathcal{D}} - f(\mathbf{x}|\mathbf{w})\}^2 \right\rangle_{\mathcal{D}, \mathbf{x}}$$

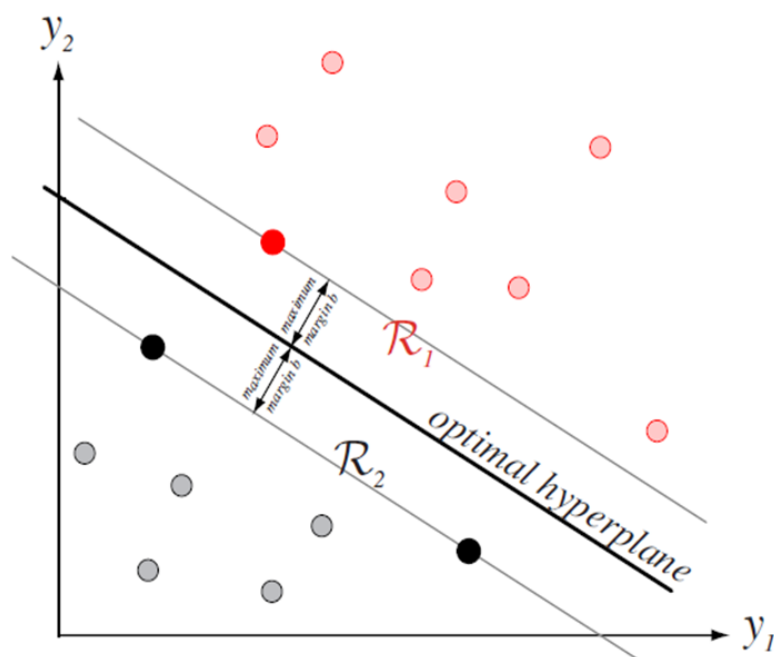
$$B = \left\langle \{f(\mathbf{x}) - \langle f(\mathbf{x}|\mathbf{w}) \rangle_{\mathcal{D}}\}^2 \right\rangle_{\mathbf{x}}$$

$$V = \left\langle \{f(\mathbf{x}|\mathbf{w}) - \langle f(\mathbf{x}|\mathbf{w}) \rangle_{\mathcal{D}}\}^2 \right\rangle_{\mathcal{D}, \mathbf{x}}$$

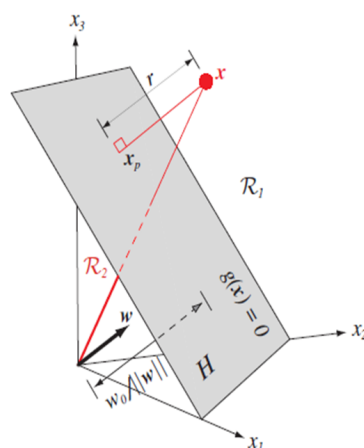
$$C = 2 \langle \{f(\mathbf{x}|\mathbf{w}) - \langle f(\mathbf{x}|\mathbf{w}) \rangle_{\mathcal{D}}\} \{f(\mathbf{x}) - \langle f(\mathbf{x}|\mathbf{w}) \rangle_{\mathcal{D}}\} \rangle_{\mathcal{D}, \mathbf{x}}$$

Bias, Variance and a term that reduces to zero (when averaged).

# Perceptron Classification and Margin



## Margin



( $b$  in formula is  $w_0$  in figure)

- Hyperplane:  $\mathbf{w}^t \mathbf{x} + b = 0$  See Lab 2

- Data:

$$\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N, \mathbf{x}_n \in \mathcal{R}^d, y_n \in \{-1, +1\}$$

- Learning problem:

$$y_n [\mathbf{w}^t \mathbf{x}_n + b] \geq 1, n = 1, \dots, N$$

- Distance from data  $\mathbf{x}_n$  to a hyperplane  $(\mathbf{w}, b)$ :

$$d(\mathbf{w}, b, \mathbf{x}_n) = \frac{|\mathbf{w}^t \mathbf{x}_n + b|}{\|\mathbf{w}\|}$$

- The margin – distance between data closest to the hyperplane on either side

$$\begin{aligned}\rho(\mathbf{w}, b) &= \min_{\mathbf{x}_n: y_n = -1} d(\mathbf{w}, b, \mathbf{x}_n) + \min_{\mathbf{x}_n: y_n = +1} d(\mathbf{w}, b, \mathbf{x}_n) \\ &= \min_{\mathbf{x}_n: y_n = -1} \frac{|\mathbf{w}^t \mathbf{x}_n + b|}{\|\mathbf{w}\|} + \min_{\mathbf{x}_n: y_n = +1} \frac{|\mathbf{w}^t \mathbf{x}_n + b|}{\|\mathbf{w}\|} \\ &= \frac{1}{\|\mathbf{w}\|} \left( \min_{\mathbf{x}_n: y_n = -1} |\mathbf{w}^t \mathbf{x}_n + b| + \min_{\mathbf{x}_n: y_n = +1} |\mathbf{w}^t \mathbf{x}_n + b| \right) \\ &= \frac{2}{\|\mathbf{w}\|}\end{aligned}$$

## Lagrangian for SVM Classification

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \alpha_n (y_n [\mathbf{w}^t \mathbf{x}_n + b] - 1), \quad \alpha_n \geq 0$$

- Setting  $\frac{\partial \mathcal{L}}{\partial b}$  to zero, gives  $\sum_{n=1}^N \alpha_n y_n = 0$
- Setting  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$  to zero, gives  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$
- Note: the unknown weights are computed as a weighted sum of the training examples; do you see a similarity to the perceptron algorithm?
- Substitute to get the dual problem

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j + \sum_{k=1}^N \alpha_k$$

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j - \sum_{k=1}^N \alpha_k$$

subject to  $\alpha_n \geq 0$  and  $\sum_{n=1}^N \alpha_n y_n = 0$

- Quadratic programming

MATLAB> help quadprog

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^t \mathbf{H} \mathbf{x} + \mathbf{f}^t \mathbf{x}$$

Subject to

$$\begin{aligned} \mathbf{A} \mathbf{x} &\leq \mathbf{b} \\ \mathbf{A}_{\text{eq}} \mathbf{x} &= \mathbf{b}_{\text{eq}} \\ \mathbf{lb} &\leq \mathbf{x} \leq \mathbf{ub} \end{aligned}$$

MATLAB> `x = quadprog( H,f,A,b,Aeq,beq,lb,ub );`

## Calculating the Bias Term

- Constraints  $\alpha_n \geq 0$ ; Parameters  $\mathbf{w} = \sum_{n=1}^N y_n \alpha_n \mathbf{x}_n$
- Non-zero  $\alpha_n$ 's correspond to Support Vectors
- For any of these support vectors ( $\mathbf{x}_s$ ):  $y_s [\mathbf{w}^t \mathbf{x}_s + b] = 1$ ; we can compute the bias term  $b$  from this.

$$y_s \left[ \sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s + b \right] = 1$$

$$y_s^2 \left( \sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s + b \right) = y_s$$

Note :  $y_s^2 = 1$ ; Hence  $b = y_s - \sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s$

- In practice, instead of using any one support vector, use we average:

$$b = \frac{1}{N_s} \sum_{s \in S} \left( y_s - \sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s \right)$$

# Non Separable Data

- Allow some slack:  $y_n(\mathbf{w}^t \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Some examples near the boundary need not achieve  $\pm 1$ , determined automatically.
- 

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

$$\text{subject to } y_n(\mathbf{w}^t \mathbf{x}_n + b) - 1 + \xi_n \geq 0, \forall n$$

- The Lagrangian for this problem is:

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n [y_n(\mathbf{w}^t \mathbf{x}_n + b) - 1 + \xi_n] - \sum_{n=1}^N \mu_n \xi_n$$

- We need the  $\alpha_n$ 's for classification constraints and  $\mu_n$ 's for positivity of slack variables.

## Nonseparable case (cont'd)

Differentiate with respect to  $\mathbf{w}$ ,  $b$  and  $\xi_n$  and equate to zero:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \implies \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \implies \sum_{n=1}^N \alpha_n y_n = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \implies C = \alpha_n + \mu_n$$

- Substitute to get (note  $\mu_n \leq 0 \implies \alpha_n \leq C$ ):

$$\max_{\alpha} \left[ \sum_{n=1}^N \alpha_n - \frac{1}{2} \alpha^t \mathbf{H} \alpha \right]$$

$$\text{subject to } 0 \leq \alpha_n \leq C \quad \forall n \quad \text{and} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

## Nonseparable case (cont'd)

$$\max_{\alpha} \left\{ \mathbf{1}^t \alpha - \frac{1}{2} \alpha^t \mathbf{H} \alpha \right\} \quad \text{subject to } 0 \leq \alpha_n \leq C \text{ and } \sum_{n=1}^N y_n \alpha_n = 0.$$

$$\min_{\alpha} \left\{ \frac{1}{2} \alpha^t \mathbf{H} \alpha - \mathbf{1}^t \alpha \right\} \quad \text{subject to } 0 \leq \alpha_n \leq C \text{ and } \sum_{n=1}^N y_n \alpha_n = 0.$$

MATLAB quadprog:

$$\min_x \left\{ \frac{1}{2} \mathbf{x}^t \mathbf{H} \mathbf{x} - \mathbf{f}^t \mathbf{x} \right\}, \quad \{\mathbf{A}, b\}, \quad \{\mathbf{A}_{eq}, b_{eq}\}, \quad \{lb, ub\}$$

quadprog( H, -ones(N,1), [ ], [ ], y', 0, 0, C )

## Nonlinear SVM Classifiers

- Matrix in QP problem:  $H_{ij} = y_i y_j \mathbf{x}_i^t \mathbf{x}_j$
- Generalize the scalar product  $\mathbf{x}_i^t \mathbf{x}_j$  to  $K(\mathbf{x}_i, \mathbf{x}_j)$
- $K(.,.)$  is called *kernel*. Under some conditions

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^t \Phi(\mathbf{x}_j)$$

- Kernels in input space map to dot products in a transformed (high dimensional) space
- Maximum margin solution in the  $\Phi$  space (a.k.a. “The Kernel Trick”).
- But without explicitly mapping the data onto that space!
- Example kernels
  - Radial Basis Function:  $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-|\mathbf{x}_i - \mathbf{x}_j| / \sigma^2)$
  - Polynomial:  $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^t \mathbf{x}_j + a)^b$