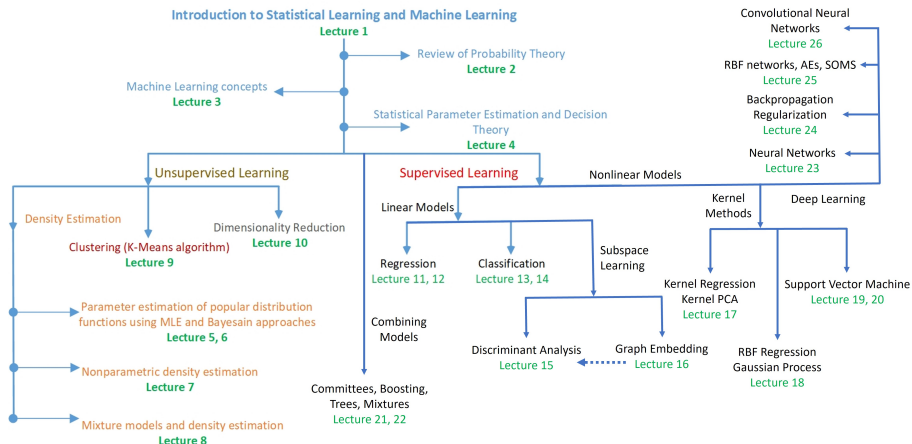


Statistical Learning and Machine Learning

Lecture 20 - Sparse Kernel Machines

October 13, 2021

Course overview and where do we stand



Support Vector Machine: Linearly separable case

Binary linear decision function defined on $\phi(x)$:

$$y(x) = w^T \phi(x) + b \quad (1)$$

where b is the bias (off-set from the origin) parameter.

Given a set of data points x_n , $n = 1, \dots, N$, the corresponding class labels $t_n \in \{-1, 1\}$ and parameters w and b classifying correctly all data points:

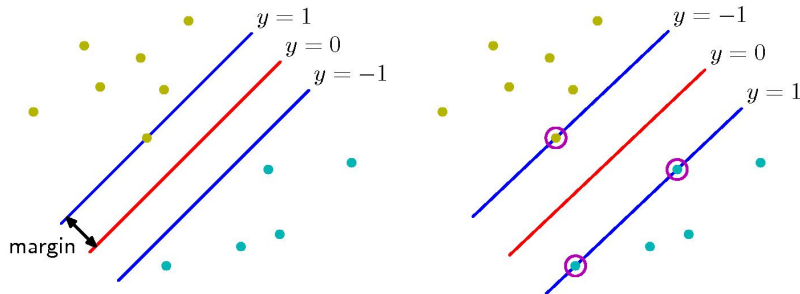
- For $x_n \in \mathcal{C}_1$: $y(x) > 0$
- For $x_n \in \mathcal{C}_2$: $y(x) < 0$
- We can unify the two above cases using:

$$t_n y(x) > 0, \quad n = 1, \dots, N. \quad (2)$$

For linearly separable classes, there may be multiple combinations of w and b satisfying the above conditions.

Support Vector Machine: Linearly separable case

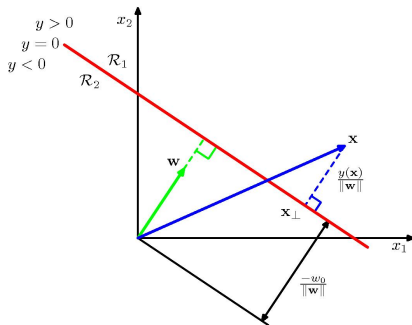
SVM selects the w and b maximizing the *margin* between the two classes.



The margin is determined by a subset of the data points, known as *support vectors*. SVM defines the decision function using only the support vectors (*sparse kernel machine*).

Support Vector Machine: Linearly separable case

Reminder: The perpendicular distance of a point x from a hyperplane defined by $y(x) = 0$ is given by $|y(x)|/\|w\|$.



Thus the distance of a point x_n to the decision hyperplane is:

$$\frac{t_n y(x_n)}{\|w\|} = \frac{t_n (w^T \phi(x_n) + b)}{\|w\|} \quad (3)$$

Support Vector Machine: Linearly separable case

The margin is defined by the perpendicular distance to the closest point x_n . Thus:

$$\arg \max_{w,b} \left\{ \frac{1}{\|w\|} \min_n \left[t_n (w^T \phi(x_n) + b) \right] \right\} \quad (4)$$

To solve the above problem, we use the following trick:

- We observe that by rescaling $w \rightarrow kw$ and $b \rightarrow kb$ the decision function $t_n y(x_n) / \|w\|$ does not change
- we (implicitly) use a k such that for the closest data point to the decision hyperplane x_j :

$$t_j (w^T \phi(x_j) + b) = 1 \quad (5)$$

- then for all data points:

$$t_n (w^T \phi(x_n) + b) \geq 1, \quad n = 1, \dots, N \quad (6)$$

Support Vector Machine: Linearly separable case

The decision function maximizing the margin ($\|w\|^{-1}$) can be obtained by optimizing for:

$$\arg \min_{w,b} \frac{1}{2} \|w\|^2 \quad (7)$$

subject to the constraints:

$$t_n \left(w^T \phi(x_n) + b \right) \geq 1, \quad n = 1, \dots, N \quad (8)$$

To solve the problem above, we introduce Lagrange multipliers $\alpha_n \geq 0$:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N \alpha_n \left[t_n \left(w^T \phi(x_n) + b \right) - 1 \right] \quad (9)$$

where $\alpha = [\alpha_1, \dots, \alpha_n]^T$.

Support Vector Machine: Linearly separable case

Setting the derivatives of $L(w, b, \alpha)$ w.r.t. w and b to zero we obtain:

$$w = \sum_{n=1}^N \alpha_n t_n \phi(x_n) \quad (10)$$

$$0 = \sum_{n=1}^N \alpha_n t_n \quad (11)$$

Eliminating w and b from $L(w, b, \alpha)$ gives the *dual representation* of the problem, in which we maximize:

$$\tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \kappa(x_n, x_m) \quad (12)$$

subject to constraints $\sum_n \alpha_n t_n = 0$ and $\alpha_n \geq 0$, $n = 1, \dots, N$.

The kernel is $\kappa(x_n, x_m) = \phi(x_n)^T \phi(x_m)$.

Support Vector Machine: Linearly separable case

After solving for $\tilde{L}(\alpha)$ (quadratic problem w.r.t. α), b is given by:

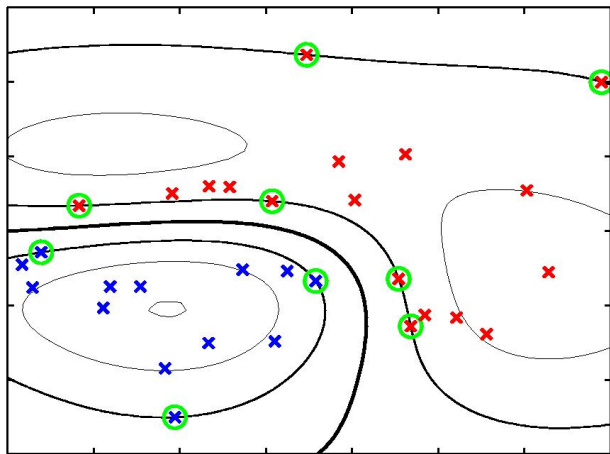
$$b = \frac{1}{N_S} \sum_{n \in \mathcal{S}} \left(t_n - \sum_{m \in \mathcal{S}} \alpha_m t_m \kappa(x_n, x_m) \right) \quad (13)$$

where \mathcal{S} is the index set of support vectors and N_S is the number of support vectors.

To classify a new data point x_* we evaluate the sign of:

$$y(x_*) = \sum_{n=1}^N \alpha_n t_n \kappa(x_n, x_*) + b. \quad (14)$$

Support Vector Machine: Non-linearly separable case



Contours of constant $y(x)$ of an SVM using RBF kernel function

Support Vector Machine: Non-linearly separable case

The decision function maximizing the margin and allowing some data points to be mis-classified can be obtained by optimizing for:

$$\arg \min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n \quad (15)$$

subject to the constraints:

$$t_n (w^T \phi(x_n) + b) \geq 1 - \xi_n, \quad n = 1, \dots, N \quad (16)$$

$$\xi_n \geq 0, \quad n = 1, \dots, N \quad (17)$$

where ξ_n , $n = 1, \dots, N$ are the *slack variables* and $C > 0$ is a hyper-parameter controlling the importance of the two terms.

For a correctly classified data point x_j the corresponding slack variable is $\xi_j = 0$.

Support Vector Machine: Non-linearly separable case

To solve the problem above, we introduce Lagrange multipliers $\alpha_n \geq 0$ and $\mu_n \geq 0$:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n \left[t_n (w^T \phi(x_n) + b) - 1 + \xi_n \right] - \sum_{n=1}^N \mu_n \xi_n \quad (18)$$

where $\alpha = [\alpha_1, \dots, \alpha_n]$.

Setting the derivatives of $L(w, b, \alpha)$ w.r.t. w , b and ξ_n , $n = 1, \dots, N$:

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{n=1}^N \alpha_n t_n \phi(x_n) \quad (19)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n t_n = 0 \quad (20)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow \alpha_n = C - \mu_n \quad (21)$$

Support Vector Machine: Non-linearly separable case

We introduce the results obtained by the derivatives above in $L(w, b, \alpha)$ and we obtain the dual Lagrangian:

$$\tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \kappa(x_n, x_m) \quad (22)$$

subject to the constraints $\sum_n \alpha_n t_n = 0$ and $0 \leq \alpha_n \leq C$, $n = 1, \dots, N$ (*box constraints*).

After solving for $\tilde{L}(\alpha)$ (quadratic problem w.r.t. α), b is given by:

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} \left(t_n - \sum_{m \in \mathcal{S}} \alpha_m t_m \kappa(x_n, x_m) \right) \quad (23)$$

where \mathcal{S} is the set of support vectors, \mathcal{M} is the set of data satisfying the box constraints and $N_{\mathcal{S}}$, $N_{\mathcal{M}}$ are the corresponding set sizes.

Support Vector Machine: Non-linearly separable case

To classify a new data point x_* we evaluate the sign of:

$$y(x_*) = \sum_{n=1}^N \alpha_n t_n \kappa(x_n, x_*) + b. \quad (24)$$

Support Vector Machine: $K > 2$

For $K > 2$, we can formulate K we can define an optimization problem optimizing K decision functions of the form:

$$y_k(x) = w_k^T x + b_{k0} \quad (25)$$

and then classify a new sample x_* using:

$$y(x_*) = \max_k y_k(x_*) \quad (26)$$

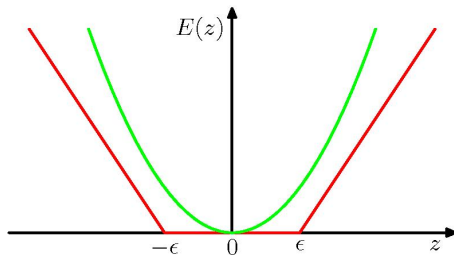
Otherwise, we can combine binary SVM classifiers using the following combination schemes:

- One-versus-rest: Define K binary SVM classifiers, each discriminating between one class the the rest
- One-versus-one: Define $K(K - 1)/2$ binary SVM classifiers, one per each class pair

Support Vector Regression

To define SVR we use the ϵ -sensitive error function:

$$E_{\epsilon}(y(x) - t) = \begin{cases} 0, & \text{if } |y(x) - t| < \epsilon \\ |y(x) - t| - \epsilon, & \text{otherwise} \end{cases} \quad (27)$$



Red: ϵ -sensitive error function

Green: quadratic error function

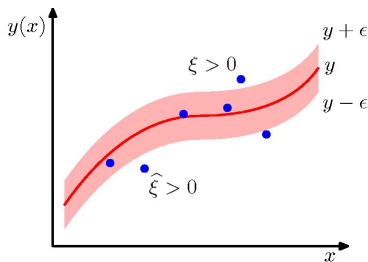
Support Vector Regression

We minimize for:

$$\mathcal{J} = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N E_{\epsilon}(y(x_n) - t_n) \quad (28)$$

where $y(x) = w^T \phi(x) + b$.

To account for errors, we can introduce slack variables (we need two for each data point ξ_n and $\hat{\xi}_n$, $n = 1, \dots, N$).



Support Vector Regression

We minimize for:

$$\mathcal{J} = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) \quad (29)$$

subject to the constraints:

$$t_n \leq y(\mathbf{x}_n) + \epsilon + \xi_n, \quad n = 1, \dots, N \quad (30)$$

$$t_n \geq y(\mathbf{x}_n) - \epsilon - \hat{\xi}_n, \quad n = 1, \dots, N \quad (31)$$

$$\xi_n \geq 0, \quad n = 1, \dots, N \quad (32)$$

$$\hat{\xi}_n \geq 0, \quad n = 1, \dots, N \quad (33)$$

Support Vector Regression

We introduce Lagrange multipliers $\alpha_n \geq 0$, $\hat{\alpha}_n \geq 0$, $\mu_n \geq 0$ and $\hat{\mu}_n \geq 0$:

$$\begin{aligned} L = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) \\ & - \sum_{n=1}^N \alpha_n (\epsilon + \xi_n + y(\mathbf{x}_n) - t_n) \\ & - \sum_{n=1}^N \hat{\alpha}_n (\epsilon + \hat{\xi}_n - y(\mathbf{x}_n) + t_n) \end{aligned} \quad (34)$$

Support Vector Regression

Setting the derivatives of L w.r.t. w , b , ξ_n and $\hat{\xi}_n$, $n = 1, \dots, N$ equal to zeros:

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{n=1}^N (\alpha_n - \hat{\alpha}_n) \phi(x_n) \quad (35)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N (\alpha_n - \hat{\alpha}_n) = 0 \quad (36)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow \alpha_n + \mu_n = C \quad (37)$$

$$\frac{\partial L}{\partial \hat{\xi}_n} = 0 \Rightarrow \hat{\alpha}_n + \hat{\mu}_n = C \quad (38)$$

Support Vector Regression

Eliminating w , b , ξ_n and $\hat{\xi}_n$, $n = 1, \dots, N$ from L we get the dual problem maximizing:

$$\begin{aligned}\tilde{L}(\alpha, \hat{\alpha}) = & -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (\alpha_n - \hat{\alpha}_n)(\alpha_m - \hat{\alpha}_m) \kappa(x_n, x_m) \\ & - \epsilon \sum_{n=1}^N (\alpha_n + \hat{\alpha}_n) + \sum_{n=1}^N (\alpha_n - \hat{\alpha}_n) t_n\end{aligned}\quad (39)$$

subject to the constraints $0 \leq \alpha_n \leq C$, $0 \leq \hat{\alpha}_n \leq C$ and $\sum_n (\alpha_n - \hat{\alpha}_n) = 0$.

The parameter b is calculated by:

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} \left(t_n - \epsilon - \sum_{m=1}^N (\alpha_m - \hat{\alpha}_m) \kappa(x_n, x_m) \right) \quad (40)$$

where \mathcal{M} is the set of data points for which $0 < \alpha_n < C$ or $0 < \hat{\alpha}_n < C$ and $N_{\mathcal{M}}$ is the size of \mathcal{M} .

Support Vector Regression

A new data point \mathbf{x}_* is evaluated by:

$$y(\mathbf{x}_*) = \sum_{n=1}^N (\alpha_n - \hat{\alpha}_n) \kappa(\mathbf{x}_n, \mathbf{x}_*) + b. \quad (41)$$