

732A99/TDDE01 Machine Learning
Lecture 3b Block 1: Support Vector Machines

Jose M. Peña
IDA, Linköping University, Sweden

Contents

- ▶ Support Vector Machines for Classification
- ▶ Support Vector Machines for Regression
- ▶ Summary

Literature

- ▶ Main source
 - ▶ Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006. Section 7.1.
- ▶ Additional source
 - ▶ Hastie, T., Tibshirani, R. and Friedman, J. *The Elements of Statistical Learning*. Springer, 2009. Sections 4.5 and 12.1-12.3.

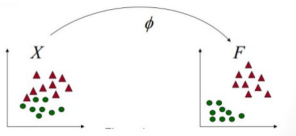
Support Vector Machines for Classification

- ▶ Consider binary classification with input space \mathbb{R}^D .
- ▶ Consider a training set $\{(\mathbf{x}_n, t_n)\}$ where $t_n \in \{-1, +1\}$.
- ▶ Consider using the linear model

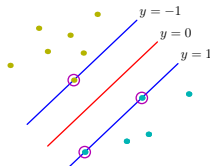
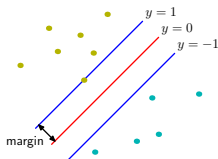
$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

so that a new point \mathbf{x} is classified according to the sign of $y(\mathbf{x})$.

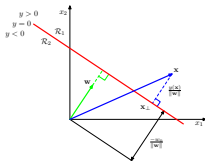
- ▶ Assume that the training set is linearly separable in the **feature space** (but not necessarily in the input space), i.e. $t_n y(\mathbf{x}_n) > 0$ for all n .



- ▶ Aim for the separating hyperplane that maximizes the **margin** (i.e. the smallest perpendicular distance from any point to the hyperplane) so as to minimize the generalization error.



Support Vector Machines for Classification



- ▶ The perpendicular distance from any point to the hyperplane is given by

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

and, thus, the margin of the hyperplane is given by

$$\min_n \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

- ▶ For any scalar κ , $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ and $\kappa y(\mathbf{x}) = \kappa \mathbf{w}^T \phi(\mathbf{x}) + \kappa b$ represent the same hyperplane. So, hereinafter we only consider rescaled hyperplanes where κ is such that $\min_n t_n (\kappa \mathbf{w}^T \phi(\mathbf{x}_n) + \kappa b) = 1$. For simplicity, we rename $\kappa \mathbf{w}$ and κb as \mathbf{w} and b . Then, the maximum margin separating hyperplane is given by

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|}$$

subject to $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$ for all n .

Support Vector Machines for Classification

- ▶ Then, the maximum margin separating hyperplane is given by

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$ for all n .

- ▶ To minimize the previous expression, we minimize

$$\frac{1}{2} \|\mathbf{w}\|^2 - \sum_n a_n (t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1)$$

where $a_n \geq 0$ are called Lagrange multipliers.

- ▶ Note that any stationary point of the Lagrangian function is a stationary point of the original function subject to the constraints. Moreover, the Lagrangian function is a quadratic function subject to linear inequality constraints. Then, it is concave, actually concave up because of the $+1/2$ and, thus, "easy" to minimize.
- ▶ Note that we are now minimizing with respect to \mathbf{w} and b , and maximizing with respect to a_n .
- ▶ Setting its derivatives with respect to \mathbf{w} and b to zero gives

$$\mathbf{w} = \sum_n a_n t_n \phi(\mathbf{x}_n)$$

$$0 = \sum_n a_n t_n$$

Support Vector Machines for Classification

- ▶ A new point \mathbf{x} is classified according to the sign of

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_n a_n t_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) + b = \sum_n a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b.$$

So, we have a **kernel method** !

- ▶ Replacing the previous expressions in the Lagrangian function gives the dual representation of the problem, in which we maximize

$$\sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m t_n t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = \sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0$ for all n , and $\sum_n a_n t_n = 0$.

- ▶ Again, this "easy" to maximize.
- ▶ Note that the dual representation makes use of the **kernel trick**, i.e. it allows working in a more convenient feature space without constructing it.

Support Vector Machines for Classification

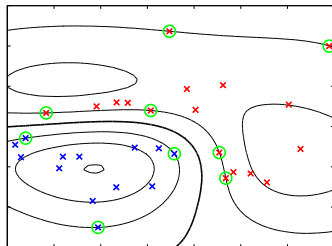
- ▶ When the Lagrangian function is maximized, the Karush-Kuhn-Tucker condition holds for all n :

$$a_n(t_n y(\mathbf{x}_n) - 1) = 0$$

- ▶ Then, $a_n > 0$ if and only if $t_n y(\mathbf{x}_n) = 1$. The points with $a_n > 0$ are called **support vectors** and they lie on the margin boundaries.
- ▶ A new point \mathbf{x} is classified according to the sign of

$$\begin{aligned} y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b &= \sum_n a_n t_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) + b = \sum_n a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \\ &= \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}, \mathbf{x}_m) + b \end{aligned}$$

where \mathcal{S} are the indexes of the support vectors. **Sparse** kernel method !



Support Vector Machines for Classification

- ▶ To find b , consider any support vector \mathbf{x}_n . Then,

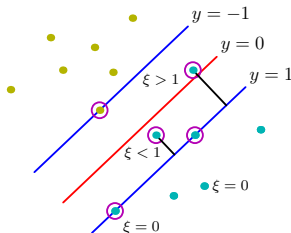
$$1 = t_n y(\mathbf{x}_n) = t_n \left(\sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right)$$

and multiplying both sides by t_n , we have that

$$b = t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

- ▶ We now drop the assumption of linear separability in the feature space, e.g. to avoid overfitting. We do so by introducing the **slack variables** $\xi_n \geq 0$ to penalize (almost-)misclassified points as

$$\xi_n = \begin{cases} 0 & \text{if } t_n y(\mathbf{x}_n) \geq 1 \\ |t_n - y(\mathbf{x}_n)| & \text{otherwise} \end{cases}$$

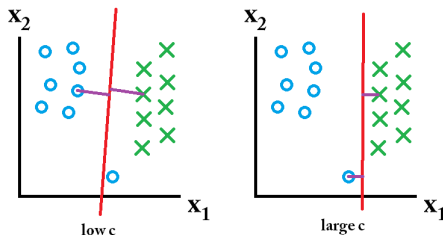


Support Vector Machines for Classification

- ▶ The optimal separating hyperplane is given by

$$\arg \min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \xi_n$$

subject to $t_n y(\mathbf{x}_n) \geq 1 - \xi_n$ and $\xi_n \geq 0$ for all n , and where $C > 0$ controls regularization. Its value can be decided by cross-validation. Note that the number of misclassified points is upper bounded by $\sum_n \xi_n$.



- ▶ To minimize the previous expression, we minimize with respect to \mathbf{w} , b , and ξ_n and maximize with respect to a_n

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \xi_n - \sum_n a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1 + \xi_n) - \sum_n \mu_n \xi_n$$

where $a_n \geq 0$ and $\mu_n \geq 0$ are Lagrange multipliers.

Support Vector Machines for Classification

- ▶ Setting its derivatives with respect to \mathbf{w} , b and ξ_n to zero gives

$$\mathbf{w} = \sum_n a_n t_n \phi(\mathbf{x}_n)$$

$$0 = \sum_n a_n t_n$$

$$a_n = C - \mu_n$$

- ▶ Replacing these in the Lagrangian function gives the dual representation of the problem, in which we maximize

$$\sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0$ and $a_n \leq C$ for all n , because $\mu_n \geq 0$.

- ▶ When the Lagrangian function is maximized, the Karush-Kuhn-Tucker conditions hold for all n :

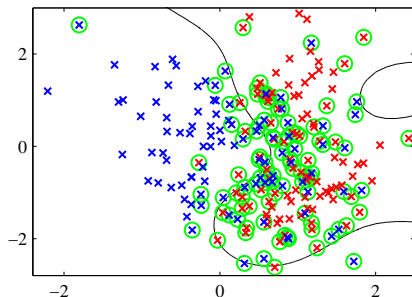
$$a_n(t_n y(\mathbf{x}_n) - 1 + \xi_n) = 0$$

$$\mu_n \xi_n = 0$$

- ▶ Then, $a_n > 0$ if and only if $t_n y(\mathbf{x}_n) = 1 - \xi_n$ for all n . The points with $a_n > 0$ are called support vectors and they lie
 - ▶ on the margin if $a_n < C$, because then $\mu_n > 0$ and thus $\xi_n = 0$, or
 - ▶ inside the margin (even on the wrong side of the decision boundary) if $a_n = C$, because then $\mu_n = 0$ and thus ξ_n is unconstrained.

Support Vector Machines for Classification

- ▶ Since the optimal \mathbf{w} takes the same form as in the linearly separable case, classifying a new point is done the same as before. Finding b is done the same as before by considering any support vector \mathbf{x}_n with $0 < a_n < C$.



- ▶ Not covered topics:
 - ▶ Classifying into more than two classes.
 - ▶ Returning class posterior probabilities.

Support Vector Machines for Regression

- ▶ Consider regressing an unidimensional continuous random variable on a D -dimensional continuous random variable.
- ▶ Consider a training set $\{(\mathbf{x}_n, t_n)\}$. Consider using the linear model

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

- ▶ To get a sparse solution, instead of minimizing the classical regularized error function

$$\frac{1}{2} \sum_n (y(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

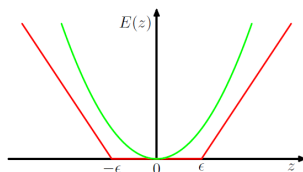
consider minimizing the **ϵ -insensitive** regularized error function

$$C \sum_n E_\epsilon(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

where $C > 0$ controls regularization and

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

Figure 7.6 Plot of an ϵ -insensitive error function (in red) in which the error increases linearly with distance beyond the insensitive region. Also shown for comparison is the quadratic error function (in green).



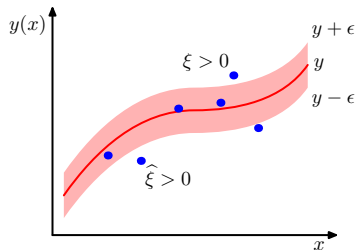
Support Vector Machines for Regression

- ▶ The values of C and ϵ can be decided by cross-validation.
- ▶ Consider the slack variables $\xi_n \geq 0$ and $\widehat{\xi}_n \geq 0$ such that

$$\xi_n = \begin{cases} t_n - y(\mathbf{x}_n) - \epsilon & \text{if } t_n > y(\mathbf{x}_n) + \epsilon \\ 0 & \text{otherwise} \end{cases}$$

and

$$\widehat{\xi}_n = \begin{cases} y(\mathbf{x}_n) - \epsilon - t_n & \text{if } t_n < y(\mathbf{x}_n) - \epsilon \\ 0 & \text{otherwise} \end{cases}$$



Support Vector Machines for Regression

- ▶ The optimal regression curve is given by

$$\arg \min_{\mathbf{w}, b, \{\xi_n\}, \{\widehat{\xi}_n\}} C \sum_n (\xi_n + \widehat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $\xi \geq 0$, $\widehat{\xi}_n \geq 0$, $t_n \leq y(\mathbf{x}_n) + \epsilon + \xi_n$ and $t_n \geq y(\mathbf{x}_n) - \epsilon - \widehat{\xi}_n$.

- ▶ To minimize the previous expression, we minimize with respect to \mathbf{w} , b , ξ_n and $\widehat{\xi}_n$ and maximize with respect to a_n

$$\begin{aligned} & C \sum_n (\xi_n + \widehat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_n (\mu_n \xi_n + \widehat{\mu}_n \widehat{\xi}_n) \\ & - \sum_n a_n (y(\mathbf{x}_n) + \epsilon + \xi_n - t_n) - \sum_n \widehat{a}_n (t_n - y(\mathbf{x}_n) + \epsilon + \widehat{\xi}_n) \end{aligned}$$

where $\mu_n \geq 0$, $\widehat{\mu}_n \geq 0$, $a_n \geq 0$ and $\widehat{a}_n \geq 0$ are Lagrange multipliers.

- ▶ Setting its derivatives with respect to \mathbf{w} , b , ξ_n and $\widehat{\xi}_n$ to zero gives

$$\mathbf{w} = \sum_n (a_n - \widehat{a}_n) \phi(\mathbf{x}_n)$$

$$0 = \sum_n (a_n - \widehat{a}_n)$$

$$C = a_n + \mu_n$$

$$C = \widehat{a}_n + \widehat{\mu}_n$$

- ▶ The prediction for a new point \mathbf{x} is made according to the **kernel method**

$$y(\mathbf{x}) = \sum (a_n - \widehat{a}_n) \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) + b = \sum (a_n - \widehat{a}_n) k(\mathbf{x}, \mathbf{x}_n) + b$$

Support Vector Machines for Regression

- ▶ Replacing the expressions above in the Lagrangian function gives the dual representation of the problem, in which we maximize

$$\frac{1}{2} \sum_n \sum_m (a_n - \widehat{a}_n)(a_m - \widehat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) - \epsilon \sum_n (a_n + \widehat{a}_n) + \sum_n (a_n - \widehat{a}_n) t_n$$

subject to $a_n \geq 0$ and $a_n \leq C$ for all n , because $\mu_n \geq 0$. Similarly for \widehat{a}_n .

- ▶ When the Lagrangian function is maximized, the Karush-Kuhn-Tucker conditions hold for all n :

$$a_n(y(\mathbf{x}_n) + \epsilon + \xi_n - t_n) = 0$$

$$\widehat{a}_n(t_n - y(\mathbf{x}_n) + \epsilon + \widehat{\xi}_n) = 0$$

$$\mu_n \xi_n = 0$$

$$\widehat{\mu}_n \widehat{\xi}_n = 0$$

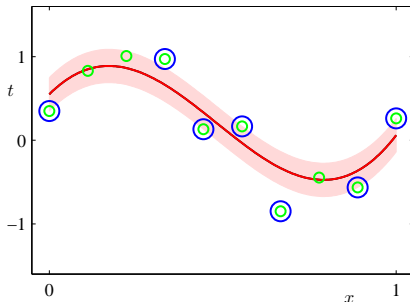
- ▶ Then, $a_n > 0$ if and only if $y(\mathbf{x}_n) + \epsilon + \xi_n - t_n = 0$, which implies that \mathbf{x}_n lies on or above the upper margin of the ϵ -tube. Similarly for $\widehat{a}_n > 0$.

Support Vector Machines for Regression

- The prediction for a new point \mathbf{x} is made according to

$$y(\mathbf{x}) = \sum_{m \in \mathcal{S}} (a_m - \widehat{a}_m) k(\mathbf{x}, \mathbf{x}_m) + b$$

where \mathcal{S} are the indexes of the support vectors. **Sparse** kernel method !



- To find b , consider any support vector \mathbf{x}_n with $0 < a_n < C$. Then, $\mu_n > 0$ and thus $\xi_n = 0$ and thus $0 = t_n - \epsilon - y(\mathbf{x}_n)$. Then,

$$b = t_n - \epsilon - \sum_{m \in \mathcal{S}} (a_m - \widehat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m)$$

Summary

- ▶ Kernel trick: It allows to work in the feature space without constructing it.
- ▶ Quadratic objective function: It allows to obtain the global optimum for a given kernel and C/ϵ (which are obtained by cross-validation).
- ▶ Sparse model: Only the support vectors are needed for classification/regression (compare with kernel models).