



Modern Machine Learning — Computing the Support Vector Classifier

Kenneth E. Barner

Department of Electrical and Computer Engineering
University of Delaware



Perfectly Separable Case: Recall that in this case, the optimization can be cast as:

$$\min_{\beta_0, \boldsymbol{\beta}} \|\boldsymbol{\beta}\| \text{ subject to } y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1, \forall i$$

Note that the constant can be arbitrarily scaled, so the optimization can be equivalently expressed as:

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|^2 \text{ subject to } y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1, \forall i$$

Lagrangian Approach: Consider a general optimization problem of the form

$$\begin{aligned} &\min_x f_0(x) \\ &\text{subject to } \begin{cases} f_i(x) \leq 0, i = 1, 2, \dots, n \\ h_i(x) = 0, i = 1, 2, \dots, p \end{cases} \end{aligned}$$

This optimization problem can be recast in a Lagrangian formulation as:

$$L_p = f_0(x) + \sum_{i=1}^n \lambda_i f_i(x) + \sum_{i=1}^p v_i h_i(x)$$

which is to be minimized with respect to x and where the constraints are in the form of additive penalties



For the perfectly separable case,

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|^2 \text{ subject to } y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1, \forall i$$

The equivalent Lagrange (primal) function to be minimized is:

$$L_p = \frac{1}{2} \|\boldsymbol{\beta}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - 1] \quad (*)$$

Setting the derivatives to zero yields: with respect to $\boldsymbol{\beta}$

$$\frac{\partial L_p}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T = 0 \implies \boldsymbol{\beta} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T$$

And with respect to β_0

$$\frac{\partial L_p}{\partial \beta_0} = \sum_{i=1}^n \alpha_i y_i = 0$$



Substituting these results into L_p , equation (*) in previous slide, yields the Wolfe dual

$$L_D = \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \mathbf{x}_i^T \mathbf{x}_k$$

subject to $\alpha_i \geq 0$ and $\sum_{i=1}^n \alpha_i y_i = 0$

Observations:

Optimization objective is to maximize L_D subject to the constraints

Optimization problem is convex

Solution must satisfy the Karush-Kuhn-Tucker conditions:

$$\boldsymbol{\beta} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T, \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad \forall i$$

and

$$\alpha_i [y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - 1] = 0 \quad \forall i$$



Aside: Recall our previous problem development

The $\|\beta\| = 1$ normalization was incorporated as follows:

$$\begin{aligned}\frac{1}{\|\beta\|} y_i (\mathbf{x}_i^T \beta + \beta_0) &\geq M \\ \Rightarrow y_i (\mathbf{x}_i^T \beta + \beta_0) &\geq M \|\beta\|\end{aligned}$$

Rescaling (β_0, β) so that $\|\beta\| = \frac{1}{M}$ yielded the optimization:

$$\min_{\beta_0, \beta \in \mathbb{R}^p} \|\beta\|$$

subject to

$$y_i (\mathbf{x}_i^T \beta + \beta_0) \geq 1, i = 1, \dots, n$$

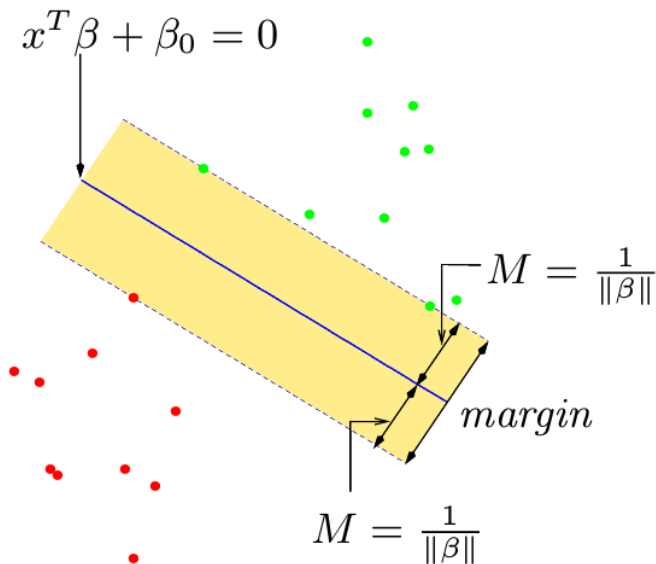
Recollections & Observations:

$\mathbf{x}_i^T \beta + \beta_0$ is the signed distance between \mathbf{x}_i and the decision boundary, $\mathbf{x}^T \beta + \beta_0 = 0$

$y_i \in \{-1, 1\}$ are the class labels

$y_i (\mathbf{x}_i^T \beta + \beta_0) \geq 1 \Rightarrow \mathbf{x}_i$ has a distance to the boundary $\geq M$

$y_i (\mathbf{x}_i^T \beta + \beta_0) = 1 \Rightarrow \mathbf{x}_i$ is on the margin boundary, i.e., has distance M





Resume Optimization Development: Recall the Wolfe dual

$$L_D = \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \mathbf{x}_i^T \mathbf{x}_k$$

that must be maximized and satisfy the Karush-Kuhn-Tucker conditions:

$$\boldsymbol{\beta} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T, \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad \forall i$$

and

$$\alpha_i [y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - 1] = 0 \quad \forall i$$

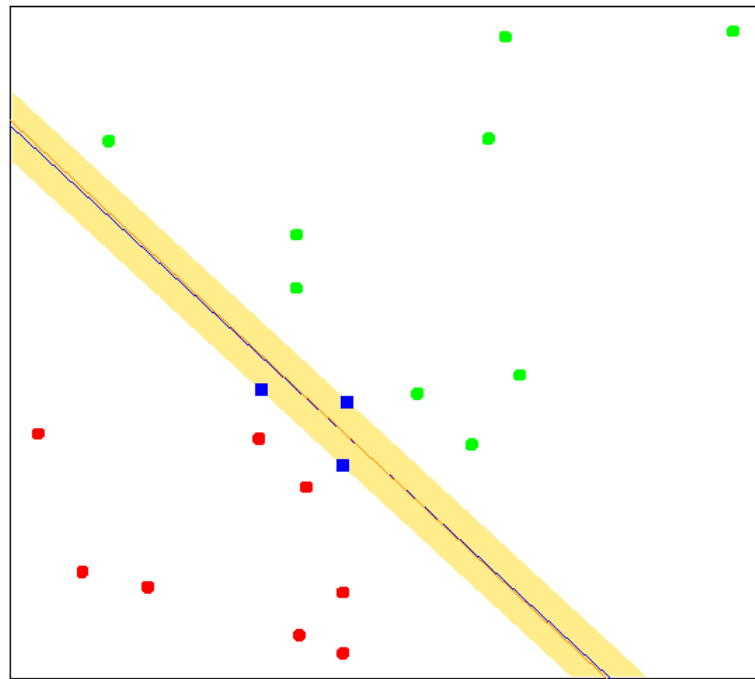
Observations:

- If $\alpha_i > 0$, then $y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) = 1 \Rightarrow \mathbf{x}_i$ is on the boundary of the slab (has distance M to the decision boundary)
- If $y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) > 1$, then \mathbf{x}_i is not on the boundary of the slab (has distance $> M$ to the decision boundary), and $\alpha_i = 0$

Interpretation: optimal $\boldsymbol{\beta}$,

$$\boldsymbol{\beta} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T$$

is defined in terms of a linear combination of **support vectors** \mathbf{x}_i —those points on the boundary of the slab (having $\alpha_i > 0$).



Binary Class Example: Optimal SVM hyperplane (blue line) and logistic regression determined hyperplane (red line) are shown. Yellow region shows the maximum margin separating the two classes. Data yields three support vectors, shown in blue.



Non-Separable Case: Recall that the non-separable case yields:

$$\min \|\beta\| \quad \text{subject to} \quad \begin{cases} y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1 - \xi_i, \forall i \\ \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq \text{constant} \end{cases}$$

Similarly to the separable case, this can be expressed as:

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \xi_i$$

subject to $\xi_i \geq 0$, and $y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1 - \xi_i, \forall i$

Note that $C = \infty$ corresponds to the separable case.

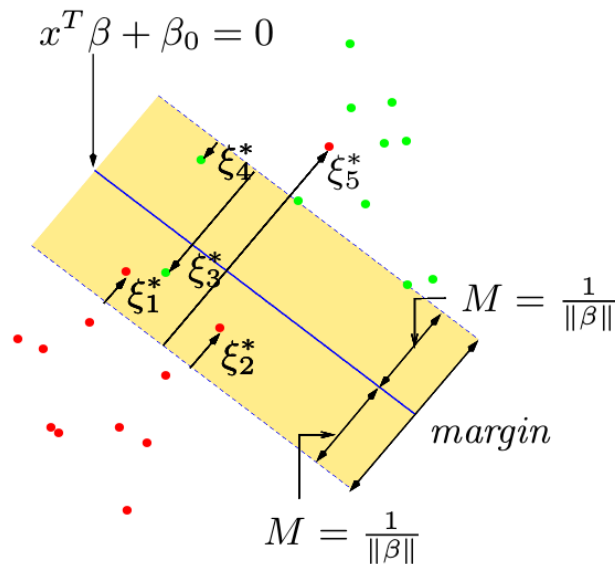
The Lagrange function equivalent representation is:

$$L_p = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i(\mathbf{x}_i^T \beta + \beta_0) - (1 - \xi_i)] - \sum_{i=1}^n \mu_i \xi_i$$

Where $\alpha_i, \mu_i, \xi_i \geq 0$. L_p is to be minimized with respect to β, β_0 and ξ_i .

Setting the derivatives with respect to each to zero yields:

$$\beta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T, \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i = C - \mu_i, \forall i$$





Substituting these results into L_D , for the non-separable case, yields the Wolfe dual

$$L_D = \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \mathbf{x}_i^T \mathbf{x}_k$$

subject to $0 \leq \alpha_i \leq C$ and $\sum_{i=1}^n \alpha_i y_i = 0$

Observations:

Solution must satisfy the Karush-Kuhn-Tucker conditions:

$$\boldsymbol{\beta} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i, \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i = C - \mu_i, \forall i$$

and

$$\alpha_i [y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - (1 - \xi_i)] = 0 \quad (*)$$

$$\mu_i \xi_i = 0$$

$$y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - (1 - \xi_i) \geq 0 \quad (**)$$

for $i = 1, \dots, n$.

Maximizing L_D , subject to the constraints, is a convex quadratic programming problem — standard quadratic programming algorithms can determine the solution

Thus

$$\hat{\boldsymbol{\beta}} = \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i$$

where $\hat{\alpha}_i$ terms are established by the quadratic programming solution maximizing solved determined by the quadratic programming algorithm maximizing L_D

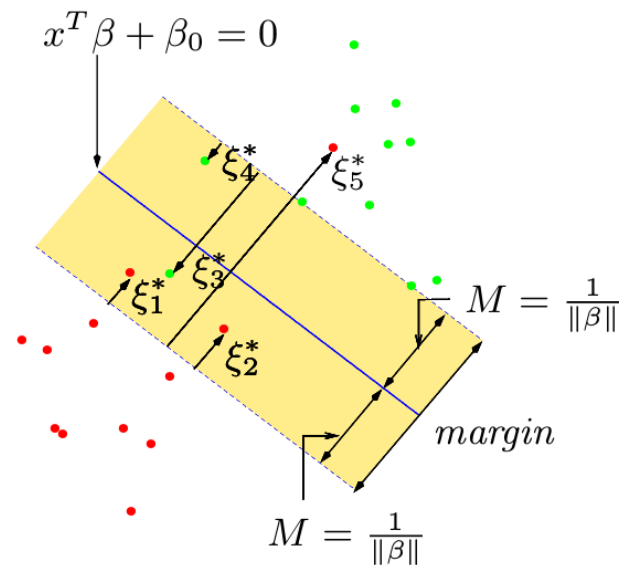
Observations:

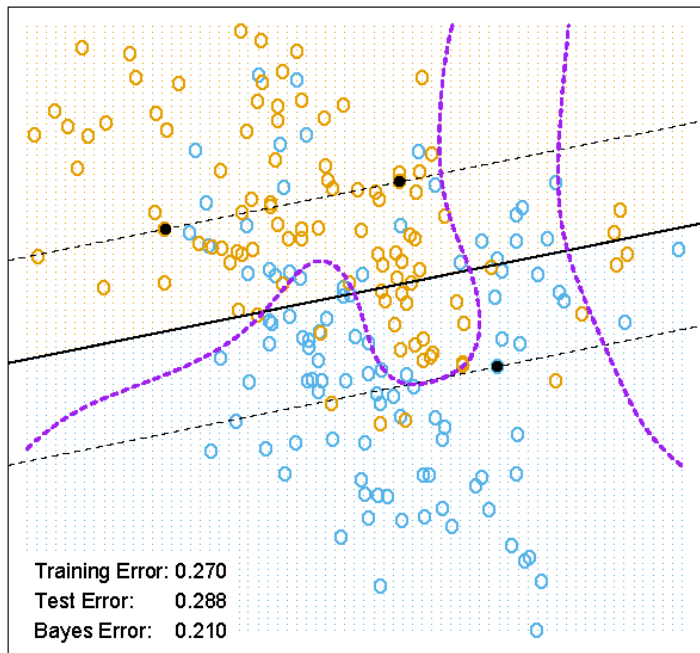
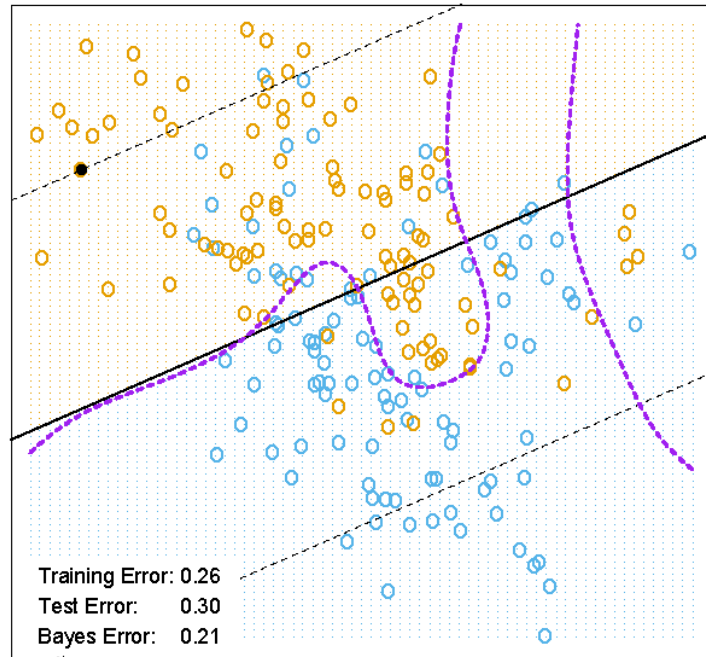
- Constraint $(*) \Rightarrow \alpha_i > 0$ only if $(**)$ is satisfied with equality, i.e.

$$y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - (1 - \xi_i) = 0 \quad (***)$$
- Those \mathbf{x}_i satisfying $(***)$ are the **support vectors**
- Support vectors are those on the wrong side of the boundary and those on the correct side of the boundary, but close to the boundary (within the margin)
- $\hat{\boldsymbol{\beta}}$ is determined as a linear combination of the support vectors
- The tuning parameter is the cost C

**Additional Observations:**

- The tuning parameter is the cost C
- Small $C \Rightarrow$ a larger margin, while a large $C \Rightarrow$ a smaller margin
- Large C focuses more on correctly classified points near the decision boundary (because margin is smaller)
- Small C involves correctly classified samples further away from the decision boundary (because margin is bigger)
- Incorrectly classified samples are always considered
- Optimal value for C can be estimated by cross-validation



 $C = 10000$  $C = 0.01$

Example: Mixture data with overlapping classes. Dashed lines indicate the margins, where $f(x) = x^T \beta + \beta_0 = \pm 1$. Support vectors ($\alpha_i > 0$) are those on the wrong side of the margin. Black dots are support vectors falling exactly on the margin ($\xi_i = 0, \alpha_i > 0$). Purple line is the Bayes decision boundary.

**Summary:**

- Linear **Support Vector Machines** determine the **optimal separating hyperplane** between two classes
- If the data is not separable, the SVM problem statement can be expressed as the Wolfe dual

$$L_D = \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \mathbf{x}_i^T \mathbf{x}_k$$

subject to $0 \leq \alpha_i \leq C$ and $\sum_{i=1}^n \alpha_i y_i = 0$

- Solution must satisfy the Karush-Kuhn-Tucker conditions:

$$\begin{aligned} \boldsymbol{\beta} &= \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i, & \sum_{i=1}^n \alpha_i y_i &= 0, & \alpha_i &= C - \mu_i, \forall i \\ \alpha_i [y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - (1 - \xi_i)] &= 0, & \mu_i \xi_i &= 0, & y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) - (1 - \xi_i) &\geq 0, \forall i \end{aligned}$$

- Maximizing L_D , subject to the constraints, is a convex quadratic programming problem — standard quadratic programming algorithms can determine the solution
- $\hat{\boldsymbol{\beta}}$ is determined as a linear combination of the **support vectors** — the vectors on the wrong side of the margin
- The tuning parameter C controls the margin; Small $C \Rightarrow$ a larger margin, while a large $C \Rightarrow$ a smaller margin
- For $f(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta} + \beta_0$, $f(\mathbf{x}) = 0$ defines the decision hyperplane boundary, $f(\mathbf{x}) = \pm 1$ define the hyperplane margins, and classification is determined as $G(\mathbf{x}) = \text{sgn}(f(\mathbf{x}))$ for binary (± 1) classes