

# SVMs Part 1: Maximum Margin Classifiers

Stephen Vardeman  
Analytics Iowa LLC  
ISU Statistics and IMSE

# Support vector classifiers-linearly separable cases

Consider a 2-class classification problem. For notational convenience, we'll suppose that output  $y$  takes values in  $\mathcal{G} = \{-1, 1\}$ . We further develop linear classification methodology.

For  $\boldsymbol{\beta} \in \mathbb{R}^p$  and  $\beta_0 \in \mathbb{R}$  we'll consider the form

$$g(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta} + \beta_0 \quad (1)$$

and a theoretical predictor/classifier

$$f(\mathbf{x}) = \text{sign}(g(\mathbf{x})) \quad (2)$$

We will approach the problem of choosing  $\boldsymbol{\beta}$  and  $\beta_0$  to in some sense provide a maximal cushion around a hyperplane separating between  $\mathbf{x}_i$  with corresponding  $y_i = -1$  and  $\mathbf{x}_i$  with corresponding  $y_i = 1$ .

# Optimization problem

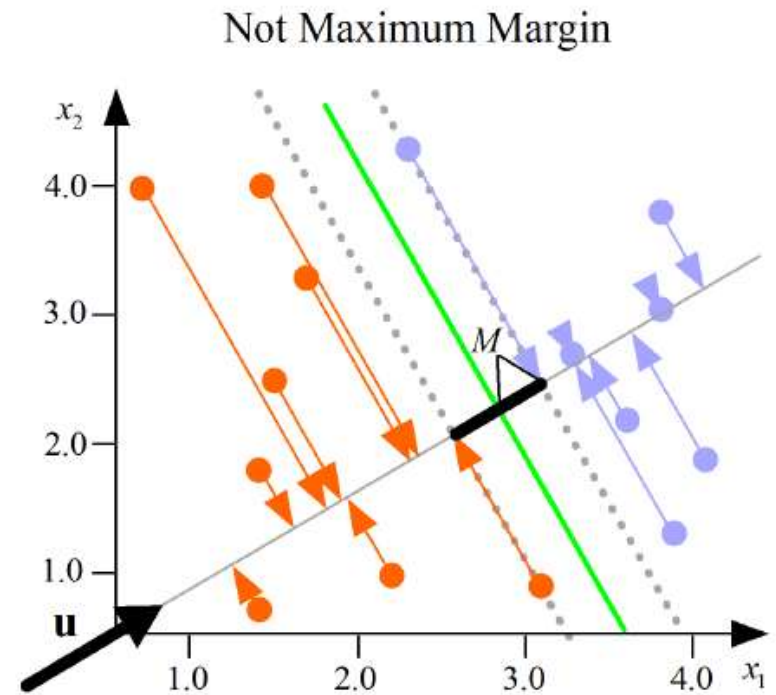
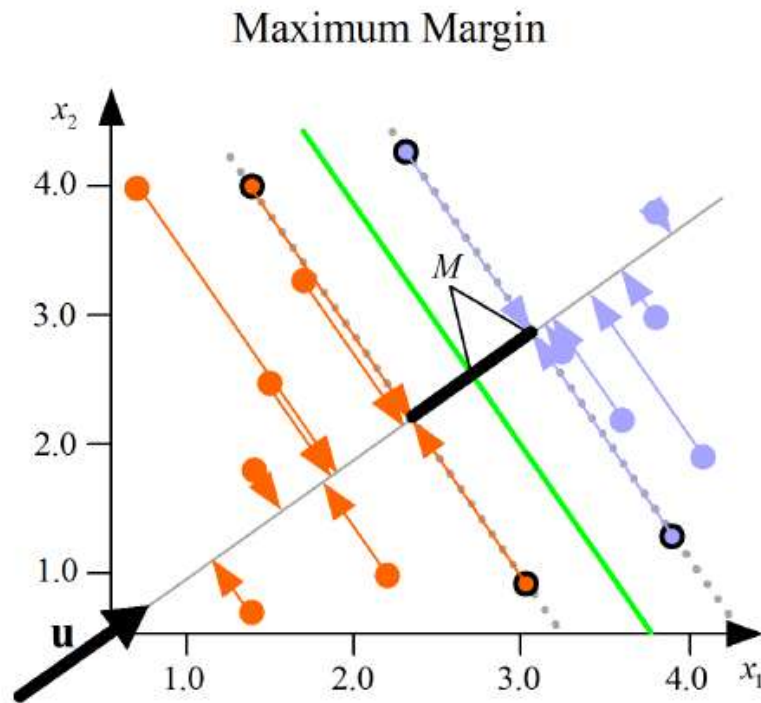
In the case that there is a classifier of form (2) with 0 training error rate, we consider the optimization problem

$$\begin{array}{ll} \text{maximize} & M \\ \mathbf{u} \text{ with } \|\mathbf{u}\| = 1 & \text{subject to } y_i (\mathbf{x}'_i \mathbf{u} + \beta_0) \geq M \quad \forall i \\ \text{and } \beta_0 \in \mathcal{R} & \end{array} \quad (3)$$

This can be thought of in terms of choosing a unit vector  $\mathbf{u}$  (or direction) in  $\mathcal{R}^p$  so that upon projecting the training input vectors  $\mathbf{x}_i$  onto the subspace of multiples of  $\mathbf{u}$  there is maximum separation between convex hull of projections of the  $\mathbf{x}_i$  with  $y_i = -1$  and the convex hull of projections of  $\mathbf{x}_i$  with corresponding  $y_i = 1$ . (The sign on  $\mathbf{u}$  is chosen to give the latter larger  $\mathbf{x}'_i \mathbf{u}$  than the former.) The resulting classifier might be termed a **maximum margin classifier**.

# Maximum margin goal

Here is a  $p = 2$  cartoon illustrating the basic "maximum margin" goal



# Properties of a solution

Notice that if  $\mathbf{u}$  and  $\beta_0$  solve maximization problem (3), the margin is

$$M = \frac{1}{2} \left( \min_{\substack{\mathbf{x}_i \text{ with} \\ y_i = 1}} \mathbf{x}'_i \mathbf{u} - \max_{\substack{\mathbf{x}_i \text{ with} \\ y_i = -1}} \mathbf{x}'_i \mathbf{u} \right)$$

and the constant that makes the voting function (1) take the value 0 on the separating hyperplane is

$$\beta_0 = -\frac{1}{2} \left( \min_{\substack{\mathbf{x}_i \text{ with} \\ y_i = 1}} \mathbf{x}'_i \mathbf{u} + \max_{\substack{\mathbf{x}_i \text{ with} \\ y_i = -1}} \mathbf{x}'_i \mathbf{u} \right)$$

## Second version of the optimization problem

For purposes of applying standard optimization theory and software, it is useful to reformulate the basic problem (3) several ways. First, note that (3) may be rewritten as

$$\begin{array}{l} \text{maximize} \\ \mathbf{u} \text{ with } \|\mathbf{u}\| = 1 \\ \text{and } \beta_0 \in \mathfrak{R} \end{array} M \quad \text{subject to } y_i \left( \mathbf{x}'_i \left( \frac{\mathbf{u}}{M} \right) + \frac{\beta_0}{M} \right) \geq 1 \quad \forall i \quad (4)$$

Then if we let

$$\boldsymbol{\beta} = \frac{\mathbf{u}}{M}$$

it's the case that

$$\|\boldsymbol{\beta}\| = \frac{1}{M} \quad \text{or} \quad M = \frac{1}{\|\boldsymbol{\beta}\|}$$

## Third version of the optimization problem

So (4) can be rewritten

$$\begin{aligned} & \underset{\beta \in \mathbb{R}^p}{\text{minimize}} && \frac{1}{2} \|\beta\|^2 && \text{subject to } y_i (\mathbf{x}'_i \beta + \beta_0) \geq 1 \quad \forall i && (5) \\ & \text{and } \beta_0 \in \mathbb{R} \end{aligned}$$

This formulation is that of a convex (quadratic criterion, linear inequality constraints) optimization problem for which there exists standard theory and algorithms.

The so-called primal functional corresponding to (5) is (for  $\alpha \in \mathbb{R}^N$ )

$$F_P(\beta, \beta_0, \alpha) \equiv \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^N \alpha_i (y_i (\mathbf{x}'_i \beta + \beta_0) - 1) \quad \text{for } \alpha \geq \mathbf{0}$$

To solve (5), one may for each  $\alpha \geq \mathbf{0}$  choose  $(\beta(\alpha), \beta_0(\alpha))$  to minimize  $F_P(\cdot, \cdot, \alpha)$  and then choose  $\alpha \geq \mathbf{0}$  to *maximize*  $F_P(\beta(\alpha), \beta_0(\alpha), \alpha)$ .

# Karush-Kuhn-Tucker conditions

The "Karush-Kuhn-Tucker conditions" are necessary and sufficient for solution of this optimization problem.

These are the *gradient conditions*

$$\frac{\partial F_P(\boldsymbol{\beta}, \beta_0, \boldsymbol{\alpha})}{\partial \beta_0} = - \sum_{i=1}^N \alpha_i y_i = 0 \quad (6)$$

and

$$\frac{\partial F_P(\boldsymbol{\beta}, \beta_0, \boldsymbol{\alpha})}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = \mathbf{0} \quad (7)$$

the *feasibility conditions*

$$y_i (\mathbf{x}_i' \boldsymbol{\beta} + \beta_0) - 1 \geq 0 \quad \forall i \quad (8)$$



# Karush-Kuhn-Tucker conditions cont.

the *non-negativity conditions*

$$\alpha \geq \mathbf{0}$$

and the *orthogonality conditions*

$$\alpha_i (y_i (\mathbf{x}_i' \boldsymbol{\beta} + \beta_0) - 1) = 0 \quad \forall i \quad (9)$$

Now (6) and (7) are

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad \boldsymbol{\beta} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \equiv \boldsymbol{\beta}(\boldsymbol{\alpha}) \quad (10)$$

and plugging these into  $F_P(\boldsymbol{\beta}, \beta_0, \boldsymbol{\alpha})$  gives a function of  $\boldsymbol{\alpha}$  only.

# Standard optimization theory

That is, write

$$\begin{aligned}F_D(\boldsymbol{\alpha}) &\equiv \frac{1}{2} \|\boldsymbol{\beta}(\boldsymbol{\alpha})\|^2 - \sum_{i=1}^N \alpha_i (y_i \mathbf{x}'_i \boldsymbol{\beta}(\boldsymbol{\alpha}) - 1) \\&= \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}'_i \mathbf{x}_j - \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}'_i \mathbf{x}_j + \sum_i \alpha_i \\&= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}'_i \mathbf{x}_j \\&= \mathbf{1}' \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}' \mathbf{H} \boldsymbol{\alpha}\end{aligned}$$

for

$$\mathbf{H}_{N \times N} = (y_i y_j \mathbf{x}'_i \mathbf{x}_j)$$

# Standard optimization theory

Then the "dual" problem is

$$\text{maximize } \mathbf{1}'\alpha - \frac{1}{2}\alpha'\mathbf{H}\alpha \quad \text{subject to } \alpha \geq \mathbf{0} \quad \text{and} \quad \alpha'\mathbf{y} = 0$$

and apparently this problem is easily solved.

# Properties of the solution

Now condition (9) implies that if  $\alpha_i^{\text{opt}} > 0$

$$y_i (\mathbf{x}'_i \boldsymbol{\beta} (\boldsymbol{\alpha}^{\text{opt}}) + \beta_0 (\boldsymbol{\alpha}^{\text{opt}})) = 1$$

so that

1. by (8) the corresponding  $\mathbf{x}_i$  has minimum  $\mathbf{x}'_i \boldsymbol{\beta} (\boldsymbol{\alpha}^{\text{opt}})$  for training vectors with  $y_i = 1$  or maximum  $\mathbf{x}'_i \boldsymbol{\beta} (\boldsymbol{\alpha}^{\text{opt}})$  for training vectors with  $y_i = -1$  (so that  $\mathbf{x}_i$  is a **support vector** for the "slab" of thickness  $2M$  around a separating hyperplane),
2.  $\beta_0 (\boldsymbol{\alpha}^{\text{opt}})$  may be determined using the corresponding  $\mathbf{x}_i$  from

$$y_i \beta_0 (\boldsymbol{\alpha}^{\text{opt}}) = 1 - y_i \mathbf{x}'_i \boldsymbol{\beta} (\boldsymbol{\alpha}^{\text{opt}}) \quad \text{i.e.} \quad \beta_0 (\boldsymbol{\alpha}^{\text{opt}}) = y_i - \mathbf{x}'_i \boldsymbol{\beta} (\boldsymbol{\alpha}^{\text{opt}})$$

(apparently for reasons of numerical stability it is common practice to average values  $y_i - \mathbf{x}'_i \boldsymbol{\beta} (\boldsymbol{\alpha}^{\text{opt}})$  for support vectors in order to evaluate  $\beta_0 (\boldsymbol{\alpha}^{\text{opt}})$ ) and,

## Properties of the solution cont.

3.

$$\begin{aligned} 1 &= y_i \beta_0 (\boldsymbol{\alpha}^{\text{opt}}) + y_i \left( \sum_{j=1}^N \alpha_j^{\text{opt}} y_j \mathbf{x}_j \right)' \mathbf{x}_i \\ &= y_i \beta_0 (\boldsymbol{\alpha}^{\text{opt}}) + \sum_{j=1}^N \alpha_j^{\text{opt}} y_j y_i \mathbf{x}_j' \mathbf{x}_i \end{aligned}$$

The fact (10) that  $\boldsymbol{\beta}(\boldsymbol{\alpha}) \equiv \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$  implies that only the training cases with  $\alpha_i > 0$  (typically corresponding to a relatively few support vectors) determine the nature of the solution to this optimization problem.

## Properties of the solution cont.

Further, for  $\mathcal{SV}$  the set indices of support vectors in the problem,

$$\begin{aligned}\|\boldsymbol{\beta}(\boldsymbol{\alpha}^{\text{opt}})\|^2 &= \sum_{i \in \mathcal{SV}} \sum_{j \in \mathcal{SV}} \alpha_i^{\text{opt}} \alpha_j^{\text{opt}} y_i y_j \mathbf{x}'_i \mathbf{x}_j \\ &= \sum_{i \in \mathcal{SV}} \alpha_i^{\text{opt}} \sum_{j \in \mathcal{SV}} \alpha_j^{\text{opt}} y_i y_j \mathbf{x}'_j \mathbf{x}_i \\ &= \sum_{i \in \mathcal{SV}} \alpha_i^{\text{opt}} (1 - y_i \beta_0(\boldsymbol{\alpha}^{\text{opt}})) \\ &= \sum_{i \in \mathcal{SV}} \alpha_i^{\text{opt}}\end{aligned}$$

the next to last of these following from 3. above, and the last following from the gradient condition (6). Then the **margin** for this problem is simply computed in terms of the  $\alpha_i^{\text{opt}}$ s as

$$M = \frac{1}{\|\boldsymbol{\beta}(\boldsymbol{\alpha}^{\text{opt}})\|} = \frac{1}{\sqrt{\sum_{i \in \mathcal{SV}} \alpha_i^{\text{opt}}}}$$