

Support Vector Machines

AW

Lecture Overview

1 Recap

2 SVM

3 Soft Margin SVM

Margins

Concept of a **margin**: The margin of a hyperplane with respect to a training set = $2 \times$ minimal distance between a point in the training set and the hyperplane.

Definition

D is a distribution over $\mathbb{R}^d \times \{\pm 1\}$. We say that D is **separable** with a (γ, ρ) -margin if there exists hyperplane $[\bar{w}^* : b^*]$ with $\|w\| = 1$ such that with probability 1 over the choice of $(\bar{x}, y) \sim D$ we have that $y(\bar{w}^* \bullet \bar{x} + b^*) \geq \gamma$ and $\|\bar{x}\| \leq \rho$.

Margins

Concept of a **margin**: The margin of a hyperplane with respect to a training set = $2 \times$ minimal distance between a point in the training set and the hyperplane.

Definition

D is a distribution over $\mathbb{R}^d \times \{\pm 1\}$. We say that D is **separable** with a (γ, ρ) -margin if there exists hyperplane $[\bar{w}^* : b^*]$ with $\|w\| = 1$ such that with probability 1 over the choice of $(\bar{x}, y) \sim D$ we have that $y(\bar{w}^* \bullet \bar{x} + b^*) \geq \gamma$ and $\|\bar{x}\| \leq \rho$.

Theorem (Optimal Hyperplane)

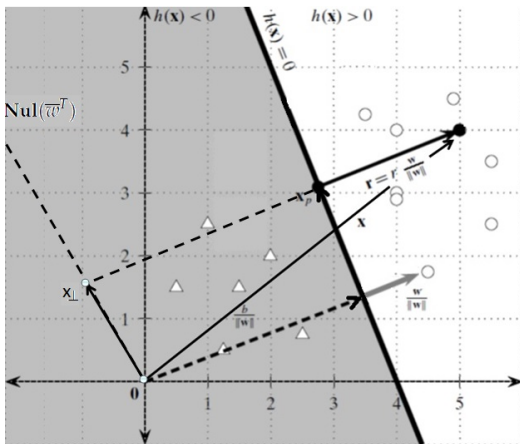
D is a distribution over $\mathbb{R}^d \times \{\pm 1\}$ that is separable with a (γ, ρ) -margin. Then, with probability of at least $1 - \delta$ over the choice of a training set of size m , there is a hyperplane that has upper bound on error $\sqrt{\frac{(\rho/\gamma)^2}{m}} + \sqrt{\frac{2 \log(2/\delta)}{m}}$

Our goal is to find this hyperplane.

Distance of a Point to the Hyperplane

What is the distance of \bar{x} from h along \bar{w} ? It is the difference of the following terms:

- The length of projection of \bar{x} onto \bar{w} where the projection vector is $\frac{\bar{x} \bullet \bar{w}}{\|\bar{w}\|^2} \bar{w} = \frac{\bar{x} \bullet \bar{w}}{\|\bar{w}\|} \hat{w}$ so length is $\frac{\bar{x} \bullet \bar{w}}{\|\bar{w}\|}$
- The length of translation vector = $\frac{-b}{\|\bar{w}\|}$



This value is negative when $\bar{w} \bullet \bar{x} < -b$ (class -1) and positive when $\bar{w} \bullet \bar{x} > -b$ (class 1). But distance is always positive so we need to take absolute value

$$d(\bar{x}, h) = \left| \frac{\bar{x} \bullet \bar{w}}{\|\bar{w}\|} - \frac{-b}{\|\bar{w}\|} \right| = \frac{y(\bar{x} \bullet \bar{w} + b)}{\|\bar{w}\|}$$

Support Vectors

Distance example: Suppose separating plane is $h = \left[\begin{pmatrix} 4 \\ 3 \end{pmatrix}^T : 20 \right]$.

Then for a point $\bar{x} = \left(\begin{pmatrix} 3 \\ 4 \end{pmatrix}, 1 \right)$ distance from h is

$$d(\bar{x}, h) = \frac{1 \left([4 \ 3] \begin{pmatrix} 3 \\ 4 \end{pmatrix} - 20 \right)}{\sqrt{4^2 + 3^2}} = \frac{4}{5}$$

Over all the m points in training set $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, we define the **margin** d^* of the linear classifier as the minimum distance of a point in S from the separating hyperplane h , i.e.

$$d^* = \min_{\bar{x}_i \in S} \left\{ \frac{y_i(\bar{x}_i \bullet \bar{w} + b)}{\|\bar{w}\|} \right\}$$

Note that $d^* \neq 0$, since h is a (strictly) separating hyperplane. All the vectors in S that are at margin distance d^* from h are the **support vectors** for the hyperplane h .

Support Vectors

Over all the m points in training set $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, we define the **margin** d^* of the linear classifier as the minimum distance of a point in S from the separating hyperplane h , i.e.

$$d^* = \min_{\bar{x}_i \in S} \left\{ \frac{y_i(\bar{x}_i \bullet \bar{w} + b)}{\|\bar{w}\|} \right\}$$

Note that $d^* \neq 0$, since h is a (strictly) separating hyperplane. All the vectors in S that are at margin distance d^* from h are the **support vectors** for the hyperplane h .

Example continued: Let $S =$

$$\left\{ \left(\begin{pmatrix} 0 \\ 1 \end{pmatrix}, -1 \right); \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, -1 \right); \left(\begin{pmatrix} 2.5 \\ 2 \end{pmatrix}, -1 \right); \left(\begin{pmatrix} 3 \\ 4 \end{pmatrix}, 1 \right); \left(\begin{pmatrix} 1.5 \\ 6 \end{pmatrix}, 1 \right) \right\}$$

We can check that $h = \left[\begin{pmatrix} 4 \\ 3 \end{pmatrix}^T : 20 \right]$ is a separating hyperplane for S . What are support vectors of this hyperplane?

Support Vectors

Over all the m points in training set $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, we define the **margin** d^* of the linear classifier as the minimum distance of a point in S from the separating hyperplane h , i.e.

$$d^* = \min_{\bar{x}_i \in S} \left\{ \frac{y_i(\bar{x}_i \bullet \bar{w} + b)}{\|\bar{w}\|} \right\}$$

Note that $d^* \neq 0$, since h is a (strictly) separating hyperplane. All the vectors in S that are at margin distance d^* from h are the **support vectors** for the hyperplane h .

Example continued: Let $S =$

$$\left\{ \left(\begin{pmatrix} 0 \\ 1 \end{pmatrix}, -1 \right); \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, -1 \right); \left(\begin{pmatrix} 2.5 \\ 2 \end{pmatrix}, -1 \right); \left(\begin{pmatrix} 3 \\ 4 \end{pmatrix}, 1 \right); \left(\begin{pmatrix} 1.5 \\ 6 \end{pmatrix}, 1 \right) \right\}$$

We can check that $h = \left[\begin{pmatrix} 4 \\ 3 \end{pmatrix}^T : 20 \right]$ is a separating hyperplane for S .

Easy to check that $d^* = \frac{4}{5}$ for h and that support vectors are $\left(\begin{pmatrix} 2.5 \\ 2 \end{pmatrix}, -1 \right)$, $\left(\begin{pmatrix} 1.5 \\ 6 \end{pmatrix}, 1 \right)$ and $\left(\begin{pmatrix} 3 \\ 4 \end{pmatrix}, 1 \right)$.

Lecture Overview

1 Recap

2 SVM

3 Soft Margin SVM

Canonical Hyperplane

- For hyperplane $h = [\bar{w} : b]$ multiplying on both sides by a scalar s yields an equivalent hyperplane representation $h = [s\bar{w} : sb]$. To obtain the **canonical hyperplane**, set s so that $d^* \|\bar{w}\| = 1$.
- So for any support vector (\bar{x}^*, y^*) holds $sy^*(\bar{w} \bullet \bar{x}^* + b) = 1$. Thus $s = \frac{1}{y^*(\bar{w} \bullet \bar{x}^* + b)}$.

Example continued: For all support vectors $\left(\begin{pmatrix} 2.5 \\ 2 \end{pmatrix}, -1\right)$,

$\left(\begin{pmatrix} 1.5 \\ 6 \end{pmatrix}, 1\right)$ and $\left(\begin{pmatrix} 3 \\ 4 \end{pmatrix}, 1\right)$ of the hyperplane

$h = \left[\begin{pmatrix} 4 \\ 3 \end{pmatrix}^T : 20\right]$ holds $s = \frac{1}{y^*(\bar{w} \bullet \bar{x}^* + b)} = \frac{1}{4}$. So canonical

representation of the hyperplane is $h = \left[\begin{pmatrix} 1 \\ 3/4 \end{pmatrix}^T : 5\right]$ and

$$d^* = \frac{1}{5}.$$

Canonical Hyperplane

- For hyperplane $h = [\bar{w} : b]$ multiplying on both sides by a scalar s yields an equivalent hyperplane representation $h = [s\bar{w} : sb]$. To obtain the **canonical hyperplane**, set s so that $d^* \|\bar{w}\| = 1$.
- So for any support vector (\bar{x}^*, y^*) holds $sy^*(\bar{w} \bullet \bar{x}^* + b) = 1$. Thus
$$s = \frac{1}{y^*(\bar{w} \bullet \bar{x}^* + b)}.$$
- Given the canonical hyperplane $h = [\bar{w} : b]$
 - for each support vector (\bar{x}^*, y^*) holds $y^*(\bar{w} \bullet \bar{x}^* + b) = 1$
 - for any point $(\bar{x}_i, y_i) \in S$ that is not a support vector we have $y_i(\bar{w} \bullet \bar{x}_i + b) > 1$, because, it is farther from the hyperplane than a support vector

Thus for all $(\bar{x}_i, y_i) \in S$ holds $y_i(\bar{w} \bullet \bar{x}_i + b) \geq 1$

Maximum Margin Hyperplane

Let \mathcal{H} be class of linear hyperplanes. For canonical hyperplane h distance of a support vector from hyperplane is $d^* = \frac{1}{\|\bar{w}\|}$, so geometric margin of a hyperplane is at least $2d^* = \frac{2}{\|\bar{w}\|}$.

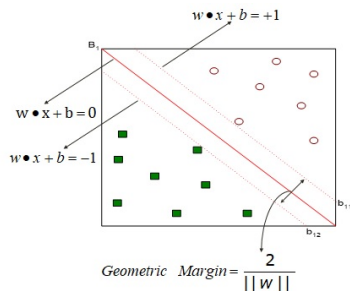
So we need

$$h^* = \arg \max_{h_S \in \mathcal{H}} \{d_{h_S}^*\} = \arg \max_{\bar{w}, b} \left\{ \frac{2}{\|\bar{w}\|} \right\}$$

under constraints $y_i(\bar{w} \bullet \bar{x}_i + b) \geq 1$ for all $(\bar{x}_i, y_i) \in S$.

Instead of maximizing margin we can minimize its inverse, and then instead of norm (to get rid of square roots) we can minimize the square of norm, the resulting $\arg \min$ will be the same as original $\arg \max$. So we need to solve

$$\begin{aligned} & \min_{\bar{w}, b} \left\{ \frac{\|\bar{w}\|^2}{2} \right\} \\ \text{subject to } & y_i(\bar{w} \bullet \bar{x}_i + b) \geq 1 \quad \text{for all } (\bar{x}_i, y_i) \in S \end{aligned}$$



Solving SVM Optimization Problem

As before introducing Lagrange multipliers λ_i for each constraint the objective function becomes:

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b} \left\{ \frac{\|\bar{w}\|^2}{2} - \sum_{i=1}^m \lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) \right\}$$

We need stationary point solution that also satisfy KKT (Karush-Kuhn-Tacker) conditions (as constraints are inequalities).

Simply put we need $\frac{\partial L}{\partial \bar{w}} = 0$ and $\frac{\partial L}{\partial b} = 0$ at the same time

$\lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) = 0$ and $\lambda_i \geq 0$ for all $0 \leq i \leq m$.

$$\frac{\partial L}{\partial \bar{w}} = \bar{w} - \sum_{i=1}^m \lambda_i y_i \bar{x}_i = 0 \quad \text{or} \quad \bar{w} = \sum_{i=1}^m \lambda_i y_i \bar{x}_i$$

and

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \lambda_i y_i = 0$$

Notice that

$$\sum_{i=1}^m \lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) = \underbrace{\bar{w} \bullet \sum_{i=1}^m \lambda_i y_i \bar{x}_i}_{\bar{w}} + \underbrace{b \sum_{i=1}^m \lambda_i y_i}_0 - \sum_{i=1}^m \lambda_i$$

Solving SVM Optimization Problem

As before introducing Lagrange multipliers λ_i for each constraint the objective function becomes:

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b} \left\{ \frac{\|\bar{w}\|^2}{2} - \sum_{i=1}^m \lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) \right\}$$

We need stationary point solution that also satisfy KKT (Karush-Kuhn-Tacker) conditions (as constraints are inequalities). Simply put we need $\frac{\partial L}{\partial \bar{w}} = 0$ and $\frac{\partial L}{\partial b} = 0$ at the same time $\lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) = 0$ and $\lambda_i \geq 0$ for all $0 \leq i \leq m$.

$$\frac{\partial L}{\partial \bar{w}} = \bar{w} - \sum_{i=1}^m \lambda_i y_i \bar{x}_i = 0 \quad \text{or} \quad \bar{w} = \sum_{i=1}^m \lambda_i y_i \bar{x}_i$$

and

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \lambda_i y_i = 0$$

Notice that $\sum_{i=1}^m \lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) = \bar{w} \bullet \bar{w} - \sum_{i=1}^m \lambda_i$

Solving SVM Optimization Problem

As before introducing Lagrange multipliers λ_i for each constraint the objective function becomes:

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b} \left\{ \frac{\|\bar{w}\|^2}{2} - \sum_{i=1}^m \lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) \right\}$$

We need stationary point solution that also satisfy KKT (Karush-Kuhn-Tacker) conditions (as constraints are inequalities). Simply put we need $\frac{\partial L}{\partial \bar{w}} = 0$ and $\frac{\partial L}{\partial b} = 0$ at the same time $\lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1) = 0$ and $\lambda_i \geq 0$ for all $0 \leq i \leq m$.

$$\frac{\partial L}{\partial \bar{w}} = \bar{w} - \sum_{i=1}^m \lambda_i y_i \bar{x}_i = 0 \quad \text{or} \quad \bar{w} = \sum_{i=1}^m \lambda_i y_i \bar{x}_i$$

and

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \lambda_i y_i = 0$$

So when pugged into $\min L$ becomes

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b} \left\{ \frac{\|\bar{w}\|^2}{2} - \|\bar{w}\|^2 + \sum_{i=1}^m \lambda_i \right\}$$

Solving SVM - continued

So we have

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b} \left\{ \frac{-\|\bar{w}\|^2}{2} + \sum_{i=1}^m \lambda_i \right\}$$

where $\bar{w} = \sum_{i=1}^m \lambda_i y_i \bar{x}_i$. Thus substituting, and noticing that $\max_{\bar{\lambda}}$ under constraints $\frac{\partial L}{\partial b} = \sum_{i=1}^m \lambda_i y_i = 0$ for all i , gives us $\min_{\bar{w}, b}$ under constraints $y_i(\bar{w} \bullet \bar{x}_i + b) \geq 1$ for all i we obtain dual quadratic program with *linear* constraints

$$\min_{\bar{w}, b} L = \max_{\bar{\lambda}} \left\{ \sum_{i=1}^m \lambda_i - \frac{\sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j}{2} \right\}$$

subject to

$$\begin{aligned} \sum_{i=1}^m \lambda_i y_i &= 0 \\ \lambda_i &\geq 0 \quad \text{for } 1 \leq i \leq m \end{aligned}$$

where $\bar{\lambda} = (\lambda_1, \dots, \lambda_m)^T$.

This program can be solved by standard methods (e.g. gradient ascent), which gives optimal values $\lambda_1, \dots, \lambda_m$.

Finding \bar{w} and b

By KKT conditions we have that for all λ_i holds

$\lambda_i(y_i(\bar{w} \bullet \bar{x}_i + b) - 1) = 0$ one of the two possibilities hold

- $\lambda_i > 0$ and then $y_i(\bar{w} \bullet \bar{x}_i + b) = 1$, i.e. \bar{x}_i is a support vector
- $\lambda_i = 0$ and $y_i(\bar{w} \bullet \bar{x}_i + b) > 1$, i.e. \bar{x}_i is not a support vector

In other words, since $\bar{w} = \sum_{i=1}^m \lambda_i y_i \bar{x}_i$, orthogonal vector of a hyperplane is a combination of support vectors. It is computed as

$$\bar{w} = \sum_{i:\lambda_i>0} \lambda_i y_i \bar{x}_i$$

Now as before we can find projections of each support vector \bar{x}_i onto \bar{w} and average these projections to get a separation point, but rather than normalizing \bar{w} , we can compute these from $y_i(\bar{w} \bullet \bar{x}_i + b) = 1$ as $b_i = y_i - \bar{w} \bullet \bar{x}_i$ and then take average

$$b = \frac{1}{|\{i : \lambda_i > 0\}|} \sum_{\lambda_i > 0} (y_i - \bar{w} \bullet \bar{x}_i).$$

Error again

Since margins are the same for all supporting vectors instead of solving $y_i(\bar{w} \bullet \bar{x}_i + b) = 1$ for all support vectors and averaging, we could solve it for one vector \bar{x}_i in class -1 and one vector \bar{x}_j in class 1 and take $b = \frac{b_i + b_j}{2}$.

Let distribution D from be separable with probability 1 over the choice of data $(\bar{x}, y) \sim D$ with some margin γ .

- Margin γ of optimal (over all data) hyperplane is at most the same as learned margin $|b - b_i|$ (where b_i projection coordinate of any support vector)
- Any computed from training set S hyperplane other than SVM-hyperplane $[\bar{w} : b]$ already has smaller margin than $[\bar{w} : b]$

Hyperplanes have finite VC-dimensions, so by Fundamental theorem they are ERM-agnostically learnable. Since D is separable and if $m > m(\epsilon, \delta)$ our hyperplane is optimal. If there is radius ρ for which $\|\bar{x}\| \leq \rho$ for all datapoints \bar{x} then by Optimal Hyperplane theorem $[\bar{w} : b]$ has the error bound

$$\sqrt{\frac{(\rho/\gamma)^2}{m}} + \sqrt{\frac{2 \log(2/\delta)}{m}}.$$

SVM in R

```
library(e1071);library(mlbench)
data(BreastCancer)
BC<-BreastCancer[!rowSums(is.na(BreastCancer)),-1]
#remove id column and rows with na valuesSVM
#can't handle it
dtrain<-sample(1:nrow(BC),2/3*nrow(BC),F)
#randomly select record #'s for training.
# remaining -dtrain numbers used for test
BC_model<- svm(Class ~ .,BC[dtrain,],
                type='C',kernel='linear')
#kernel linear is a vector (line);
#can be another curve
summary(BC_model)
pred<-predict(BC_model,BC[dtrain,-10]).
table(pred, BC[dtrain, ]$Class)
table(predicted=predict(BC_model,BC[-dtrain,-10]),
       true=BC[-dtrain, ]$Class)
```

Lecture Overview

1 Recap

2 SVM

3 Soft Margin SVM

Soft Margins

If data that is linearly separable there is no \overline{w}, b that satisfy constraints of maximization problem! No solution will be returned. How to modify optimization problem to handle inseparable data?

Key Idea: use of slack parameters. *Intuition:* suppose $[\overline{w}, b]$ reliably separates 99.9(9)% of data points with a large margin. But there is couple data points that end up on the wrong side of the hyperplane, perhaps these are noisy.

- Suppose we can 'move' these points across the hyperplane to the correct side
- We pay for the 'move' of a point, small 'cost', so that a total cost is acceptable if we are not moving a lot of points around

Soft Margin SVM explained

- We formalize the notion of moving to the right side by introducing one slack variable for each training example:

$$y_i(\bar{w} \bullet \bar{x}_i + b) + \xi_i \geq 1 \text{ where } \xi_i \geq 0$$

- We formalize the notion of cost by penalizing oneself for having to use slack in the objective. The smaller the slack the less we pay:

$$\min_{\bar{w}, b, \bar{\xi}} \left\{ \frac{\|\bar{w}\|^2}{2} + C \sum_{i=1}^m (\xi_i)^k \right\}$$

C and k are parameters of **Soft margin SVM**:

- k defines **loss** function $f(\bar{\xi}) \sum_{i=1}^m (\xi_i)^k$ - total cost that we pay for violating constraints; $k = 1$ called **hinge** loss, $k = 2$ quadratic loss
- C is a regularization constant that controls the trade-off between maximizing the margin (minimizing $\frac{\|\bar{w}\|^2}{2}$) and minimizing loss (minimizing total error = distance away from hyperplane $\sum_{i=1}^m (\xi_i)^k$)

Solving Soft Margin SVM w/Hinge Loss

Soft Margin SVM with hinge loss optimization problem:

$$\begin{aligned} & \min_{\bar{w}, b, \bar{x}_i} \left\{ \frac{\|\bar{w}\|^2}{2} + C \sum_{i=1}^m \xi_i \right\} \\ \text{subject to} \quad & y_i(\bar{w} \bullet \bar{x}_i + b) + \xi_i \geq 1 && \text{for all } (\bar{x}_i, y_i) \in S \\ & \xi_i \geq 0 && \text{for all } 1 \leq i \leq m \end{aligned}$$

Introducing Lagrange multipliers λ_i for margin constraints and β_i for non-negativity of slack variables constraints we get

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b, \bar{\xi}} \left\{ \frac{\|\bar{w}\|^2}{2} - \sum_{i=1}^m \lambda_i (y_i(\bar{w} \bullet \bar{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^m (\beta_i - C) \xi_i \right\}$$

with KKT conditions being $\lambda_i(y_i(\bar{w} \bullet \bar{x}_i + b) - 1 + \xi_i) = 0$ and $\beta_i \xi_i = 0$.
Equating partial derivatives wrt \bar{w} , b and all ξ_i to 0 we get:

$$\begin{aligned} \bar{w} &= \sum_{i=1}^m \lambda_i y_i \bar{x}_i && \text{as before} \\ \sum_{i=1}^m \lambda_i y_i &= 0 && \text{as before} \\ \text{and} \quad C - \lambda_i - \beta_i &= 0 && \text{for all } 1 \leq i \leq m \end{aligned}$$

Soft Margin SVM w/Hinge Loss - continued

Substituting

$$\bar{w} = \sum_{i=1}^m \lambda_i y_i \bar{x}_i \quad \text{as before}$$

$$\sum_{i=1}^m \lambda_i y_i = 0 \quad \text{as before}$$

and

$$C - \lambda_i - \beta_i = 0 \quad \text{for all } 1 \leq i \leq m$$

back into

$$\min_{\bar{w}, b} L = \min_{\bar{w}, b, \xi} \left\{ \frac{\|\bar{w}\|^2}{2} - \sum_{i=1}^m \lambda_i (y_i (\bar{w} \bullet \bar{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^m (\beta_i - C) \xi_i \right\}$$

we obtain dual program that is almost the same as in separable case except to the range of λ_i

$$\min_{\bar{w}, b} L = \max_{\lambda} \left\{ \sum_{i=1}^m \lambda_i - \frac{\sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j}{2} \right\}$$

subject to

$$\sum_{i=1}^m \lambda_i y_i = 0$$

$$C \geq \lambda_i \geq 0 \quad \text{for } 1 \leq i \leq m$$

The range of λ_i is determined by constraint $C = \lambda_i + \beta_i$ and the fact that $\beta_i \geq 0$

Soft Margin SVM w/Hinge Loss - continued

Solving

$$\begin{aligned} \min_{\bar{w}, b} L &= \max_{\bar{\lambda}} \left\{ \sum_{i=1}^m \lambda_i - \frac{\sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j}{2} \right\} \\ \text{subject to} & \quad \sum_{i=1}^m \lambda_i y_i = 0 \\ & \quad C \geq \lambda_i \geq 0 \quad \text{for } 1 \leq i \leq m \end{aligned}$$

gives optimal values $\lambda_1, \dots, \lambda_m$. As before support vectors are those for which $\lambda_i > 0$. So \bar{x}_i for which holds $y_i(\bar{w} \bullet \bar{x}_i + b) + \xi_i = 1$ are support vectors (due to KKT conditions). They now include all points that are on the margin (which have zero slack $\xi_i = 0$), as well as all points with positive slack ($\xi_i > 0$)!

As before orthogonal vector of the hyperplane is $\bar{w} = \sum_{\lambda_i > 0} \lambda_i y_i \bar{x}_i$. Since $C = \lambda_i + \beta_i$ and $\beta_i \xi_i = 0$ (KKT condition) we get that those \bar{x}_i for which $0 < \lambda_i < C$ are exactly on the margins ($\xi_i = 0$ while those for which $\lambda_i = C$ are behind hyperplane. To compute bias b we take points on the margins and as before compute $b_i = y_i - \bar{w} \bullet \bar{x}_i$ and then take

Reading

TSKK (main textbook) Section 4.9

Zaki, Meira, Sections 21.1, 21.2, 21.3.1.