

SUPPORT VECTOR MACHINES

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

OPTIMAL SEPARATING HYPERPLANES

A main initiative in early computer science was to find separating hyperplanes among groups of data

(Rosenblatt (1958) with the perceptron algorithm)

The issue is that if there is a separating hyperplane, there is an infinite number

An optimal separating hyperplane can be generated by finding support points and bisecting them.

(Sometimes optimal separating hyperplanes are called maximum margin classifiers)

BASIC LINEAR GEOMETRY

A hyperplane in \mathbb{R}^p is given by

$$\mathcal{H} = \{X \in \mathbb{R}^p : h(X) = \beta_0 + \beta^\top X = 0\}$$

(Usually it is assumed that $\|\beta\|_2 = 1$)

1. The vector β is **normal** to \mathcal{H}

(To see this, let $X, X' \in \mathcal{H}$. Then $\beta^\top (X - X') = 0$)

2. **IMPORTANT:** For any point $X \in \mathbb{R}^p$, the (signed) length of its orthogonal complement to \mathcal{H} is $h(X)$

SUPPORT VECTOR MACHINES (SVM)

Let $Y_i \in \{-1, 1\}$

(It is common with SVMs to code Y this way. With logistic regression, Y is commonly phrased as $\{0, 1\}$ due to the connection with Bernoulli trials)

We will generalize this to supervisors with more than 2 levels at the end

A classification rule induced by a hyperplane is

$$g(X) = \text{sgn}(X^\top \beta + \beta_0)$$

SEPARATING HYPERPLANES

Our classification rule is based on a hyperplane \mathcal{H}

$$g(X) = \text{sgn}(X^\top \beta + \beta_0)$$

A **correct** classification is one such that $h(X)Y > 0$ and $g(X)Y > 0$

(Why?)

The larger the quantity $Yh(X)$, the more “sure” the classification

(**REMINDER:** The signed distance to \mathcal{H} is $h(X)$)

Under classical **separability**, we can find a function such that $Y_i h(X_i) > 0$

(That is, makes perfect training classifications via g)

OPTIMAL SEPARATING HYPERPLANE

This idea can be encoded in the following **convex program**

$$\begin{aligned} & \max_{\beta_0, \beta} M \text{ subject to} \\ & Y_i h(X_i) \geq M \text{ for each } i \text{ and } \|\beta\|_2 = 1 \end{aligned}$$

INTUITION:

- We know that $Y_i h(X_i) > 0 \Rightarrow g(X_i) = Y_i$. Hence, larger $Y_i h(X_i) \Rightarrow$ “more” correct classification
- For “more” to have any meaning, we need to **normalize** β , thus the other constraint

OPTIMAL SEPARATING HYPERPLANE

Let's take the original program:

$$\max_{\beta_0, \beta} M \text{ subject to}$$

$$Y_i h(X_i) \geq M \text{ for each } i \text{ and } \|\beta\|_2 = 1$$

and **rewrite** it as

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|_2^2 \text{ subject to}$$

$$Y_i h(X_i) \geq 1 \text{ for each } i$$

(Replace $Y_i h(X_i) \geq M$ with $\frac{1}{\|\beta\|_2} Y_i h(X_i) \geq M$, which redefines β_0)

This is still a **convex** optimization program: quadratic criterion, linear inequality constraints

OPTIMAL SEPARATING HYPERPLANE

Again, we can convert this **constrained** optimization problem into the **Lagrangian** form

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|_2^2 - \sum_{i=1}^n \alpha_i [Y_i (X_i^\top \beta + \beta_0) - 1]$$

In contrast to the lasso problem, there are now n Lagrangian parameters $\alpha_1, \dots, \alpha_n$

(There are n constraints, after all)

Before continuing, let's discuss a general method of theoretical optimization

Duality and KKT conditions

DUALITY

Suppose we wish to solve the **primal** problem

$$\min_{\beta} \ell(\beta) \text{ s.t. } c_j(\beta) \leq 0 \text{ for } j = 1, \dots, p$$

(Here, c is meant to convey constraint)

We define the **Lagrangian** as

$$L(\beta, \alpha) = \ell(\beta) + \sum_{j=1}^p \alpha_j c_j(\beta)$$

with $\alpha_j \geq 0$

Note that $L(\beta, \alpha) \leq \ell(\beta)$ at each feasible β and $\alpha \geq 0$

LAGRANGE DUALITY

Let C denote the primal feasible set and ℓ^* be the optimal primal value

Then,

$$\ell^* \geq \min_{\beta \in C} L(\beta, \alpha) \geq \min_{\beta} L(\beta, \alpha) := \Lambda(\alpha)$$

is the **Lagrange dual function**

Hence,

$$\Lambda^* := \max_{\alpha \geq 0} \Lambda(\alpha) \leq \ell^*$$

provides a lower bound known as **weak duality**

In some important cases, it turns out **strong duality** holds

$$\Lambda^* = \ell^*$$

STRONG DUALITY

The general sufficient condition for strong duality is known as
SLATER'S CONDITION

- Primal is convex

(Note that the dual is always convex (in the sense of a maximum concave problem) even if primal is not)

- There is at least one strictly feasible $\beta \in \mathcal{C}$

(This is an extremely weak condition)

IMPLICATION: We can use the **dual gap** $= \ell(\beta) - \Lambda(\alpha)$ as a stopping criterion for iterative approaches:

$$\ell(\beta) - \ell^* \leq \ell(\beta) - \Lambda(\alpha) \stackrel{\text{check}}{\leq} \text{tolerance}$$

KARUSH-KUHN-TUCKER

The **KKT** conditions are



$$0 \in \partial \ell|_{\beta} + \sum_{j=1}^p \alpha_j \partial c_j|_{\beta}$$

(This is just the classic stationary point for subgradients condition)

- **COMPLEMENTARY SLACKNESS:** $\alpha_j c_j(\beta) = 0$
- **PRIMAL FEASIBLE:** $\beta \in \mathcal{C}$
- **DUAL FEASIBLE:** $\alpha_j \geq 0$

These conditions are always sufficient for finding an optimum and are necessary under strong duality

(That is, β^* and α^* are primal/dual solutions if and only if β^* and α^* satisfy the KKT conditions)

This gives rise to the **alternating direction method of multipliers (ADMM)** and **strong/safe rules**

Back to optimal separating hyperplane

OPTIMAL SEPARATING HYPERPLANE

$$\frac{1}{2} \|\beta\|_2^2 - \sum_{i=1}^n \alpha_i [Y_i (X_i^\top \beta + \beta_0) - 1]$$

Derivatives with respect to β and β_0 :

- $\beta = \sum_{i=1}^n \alpha_i Y_i X_i$
- $0 = \sum_{i=1}^n \alpha_i Y_i$

Substituting into the **Lagrangian**:

$$\text{Wolfe dual} = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k Y_i Y_k X_i^\top X_k$$

(this is all subject to $\alpha_i \geq 0$)

We want to **maximize** Wolfe dual

OPTIMAL SEPARATING HYPERPLANE

Looking at the complementary slackness

$$\alpha_i [1 - Y_i h(X_i)] = 0 \text{ for all } i$$

(The product of Lagrangian parameters and inequality constraint equals 0)

This implies either:

- $\alpha_i = 0$, which happens if the constraint $Y_i h(X_i) > 1$
(That is, when the constraint is non binding)
- $\alpha_i > 0$, which happens if the constraint $Y_i h(X_i) = 1$
(That is, when the constraint is binding)

OPTIMAL SEPARATING HYPERPLANE

Taking this relationship

$$\alpha_i[Y_i h(X_i) - 1] = 0$$

we see that, for $i = 1, \dots, n$,

- The points (X_i, Y_i) such that $\alpha_i > 0$ are **support vectors**
- The points (X_i, Y_i) such that $\alpha_i = 0$ are **irrelevant** for classification

A point (X_i, Y_i) is a support vector if it defines the **margin** of the optimal separating hyperplane

Support vector classifier

SUPPORT VECTOR CLASSIFIER

Of course, we can't realistically assume that the data are linearly separated (even in a transformed space)

In this case, the previous program has no **feasible** solution

We need to introduce **slack** variables, ξ , that allow for overlap among the classes

These slack variables allow for us to encode **training misclassification** into the optimization problem

SUPPORT VECTOR CLASSIFIER

$$\begin{aligned} & \max_{\beta_0, \beta, \xi_1, \dots, \xi_n} M \text{ subject to} \\ & Y_i h(X_i) \geq M \underbrace{(1 - \xi_i), \xi_i \geq 0, \sum \xi_i \leq t}_{\text{new}}, \text{ for each } i \end{aligned}$$

Note that

- t is a **tuning parameter**. The literature usually refers to t as a **budget**
(Think: lasso)
- The separable case corresponds to $t = 0$

SUPPORT VECTOR CLASSIFIER

We can rewrite the problem again:

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|_2^2 \quad \text{subject to}$$

$$Y_i h(X_i) \geq 1 - \underbrace{\xi_i}_{\text{new}}, \xi_i \geq 0, \sum \xi_i \leq t, \text{ for each } i$$

(Convex optimization program: quadratic criterion, linear inequality constraints.)

Converting $\sum \xi_i \leq t$ to the Lagrangian:

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|_2^2 + \lambda \sum \xi_i \quad \text{subject to}$$

$$Y_i h(X_i) \geq 1 - \xi_i, \xi_i \geq 0, \text{ for each } i$$

(Think: lasso. $\lambda \sum \xi_i + \xi_i \geq 0 \Rightarrow \lambda \|\xi\|_1$)

SVMs: SLACK VARIABLES

The **slack variables** give us insight into the problem

- If $\xi_i = 0$, then that observation is on **correct** the side of the **margin**
- If $\xi_i \in (0, 1]$, then that observation is on the **incorrect** side of the **margin**, but still correctly classified
- If $\xi_i > 1$, then that observation is **incorrectly** classified

SUPPORT VECTOR CLASSIFIER

Continuing to convert constraints to **Lagrangian**

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|_2^2 + \lambda \sum \xi_i - \underbrace{\sum_{i=1}^n \alpha_i [Y_i (X_i^\top \beta + \beta_0) - (1 - \xi_i)] - \sum_{i=1}^n \gamma_i \xi_i}_{\text{remaining constraints}}$$

Necessary conditions (taking derivatives)

- $\beta = \sum_{i=1}^n \alpha_i Y_i X_i$
- $0 = \sum_{i=1}^n \alpha_i Y_i$
- $\alpha_i = \lambda - \gamma_i$

(As well as positivity constraints: $\lambda, \alpha_i, \gamma_i \geq 0$)

SUPPORT VECTOR CLASSIFIER

Substituting, we require the **Wolfe dual**

This, combined with the **KKT** conditions uniquely characterize the solution:

$$\max_{\alpha \text{ subject to: KKT + Wolfe dual}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{i'=1}^n \alpha_i \alpha_{i'} Y_i Y_{i'} X_i^\top X_{i'}$$

(See Chapter 12.2.1 in “Elements of Statistical Learning”)

Note: the necessary conditions $\beta = \sum_{i=1}^n \alpha_i Y_i X_i$ imply estimators of the form

- $\hat{\beta} = \sum_{i=1}^n \hat{\alpha}_i Y_i X_i$
- $\hat{\beta}^\top X = \sum_{i=1}^n \hat{\alpha}_i Y_i X_i^\top X$

SUPPORT VECTOR CLASSIFIER

We can think of t as a **budget** for the problem

If $t = 0$, then there is **no budget** and we won't tolerate any margin violations

If $t > 0$, then no more than $\lfloor t \rfloor$ observations can be misclassified

A larger t then leads to larger **margins**

(we allow more margin violations)

SUPPORT VECTOR CLASSIFIER

FURTHER INTUITION:

Like the optimal hyperplane, only observations that violate the margin determine \mathcal{H}

A large t allows for many violations, hence many observations factor into the fit

A small t means only a few observations do

Hence, t calibrates a bias/variance trade-off, as expected

In practice, t gets selected via cross-validation

SUPPORT VECTOR CLASSIFIER

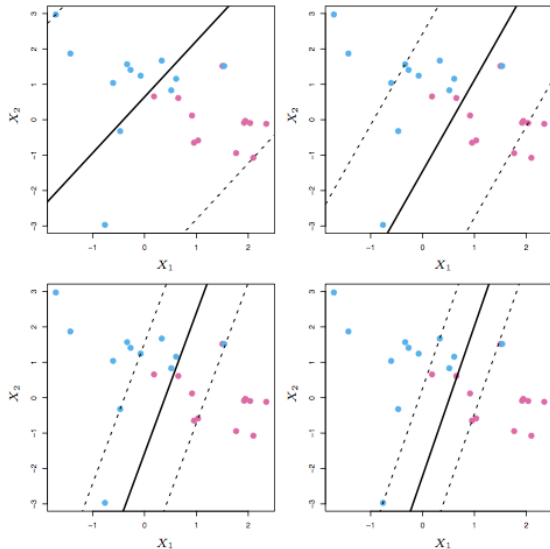


Figure 9.7 in ISL

SUPPORT VECTOR CLASSIFIER IN R

A common package to use is `e1071`

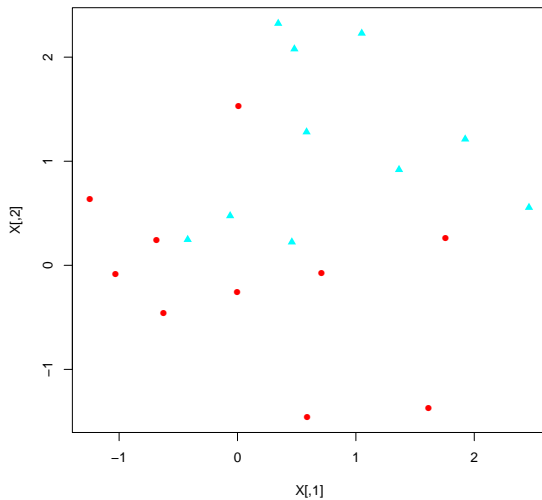
```
X = matrix(rnorm(20*2),ncol=2)
Y = c(rep(-1,10),rep(1,10))
X[Y == 1,] = X[Y == 1,] + 1
```

```
col = rep(0,length(Y))
col[Y == -1] = rainbow(2)[1]
col[Y == 1] = rainbow(2)[2]
```

```
pch = rep(0,length(Y))
pch[Y == -1] = 16
pch[Y == 1] = 17
```

```
plot(X,col=col,pch=pch)
```

SUPPORT VECTOR CLASSIFIER IN \mathbb{R}



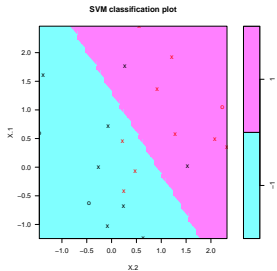
SUPPORT VECTOR CLASSIFIER IN R

```
library(e1071)
dat  =data.frame(X=X, Y=as.factor(Y))
svmfit=svm(Y~., data=dat, kernel="linear", cost=cost)
```

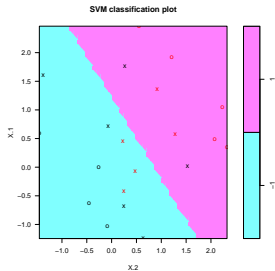
IMPORTANT: Their definition of cost is the **Lagrangian** version, which we defined as λ

Hence, a **small cost** means a large t and a **wider** margin

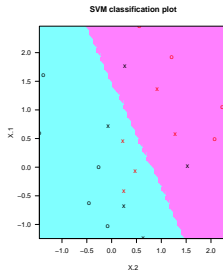
SUPPORT VECTOR CLASSIFIER IN \mathbb{R}^2



cost = .1



cost = 1

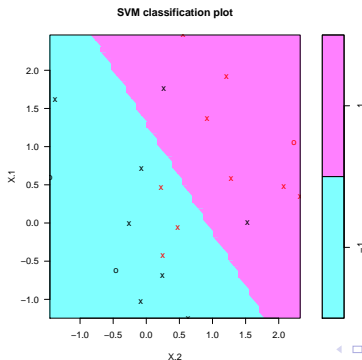


cost = 10

SUPPORT VECTOR CLASSIFIER IN R

```
tune.out = tune(svm,Y~.,data=dat,kernel="linear",  
               ranges=list(cost=c(0.001, 0.01, 0.1, 1,5,10,100)))  
best.model = tune.out$best.model
```

Note that `best.model` is an `svm` object:



NEXT TIME: KERNEL METHODS

INTUITION: Many methods have linear decision boundaries

We know that sometimes this isn't sufficient to represent data

EXAMPLE: Sometimes we need to include a polynomial effect or a log transform in multiple regression

Sometimes, a **linear** boundary, but in a different space makes all the difference..