# Chapter 9: Support Vector Machines

The *support vector machine* is an approach for classification that was developed in the computer science community in the 1990s and has grown in popularity.

```
SVMs perform will in a variety of suthtys
often considered one of the bast "out of the box" classifiers.
```

The support vector machine is a generalization of a simple and intuitive classifier called the *maximal margin classifier*.



Support vector machines are intended for binary classification, but there are extensions for more than two classes.

cuteborial response v/ only 2 classes.



Credit: https://dilbert.com/strip/2013-02-02

### 1 Maximal Margin Classifier

> oppension of endiden space.

In *p*-dimensional space, a *hyperplane* is a flat affine subspace of dimension p-1.

c.y. In 2 dimensions, a hyperplane is a flat 1 dimensional subspace - a line. In 3 dimensions, a hyperplane is a flat 2 dimensional subspace - a plane

In p > 3 demensions, horder the conceptualize, but still a flat p-1 dim. subspace. The mathematical definition of a hyperplane is quite simple,

In 2 dimensions, a hyperplane is defined by  $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$ i.e. any X = (X, X2) for which this equation holds lies on the hyperplace.

Note this is just the equation of a line. This can be easily extended to the *p*-dimensional setting.

B. + B. X, + ... + B. XD = 0 defines a p-din hyperplane.

i.e. any x = (x1,-, xp) for which this equation holds lies on the hyperplane.

We can think of a hyperplane as dividing *p*-dimensional space into two halves.

If: Bo + B1 X1 + ... + BpXp >0 then X=(X1,...,Xp) lies on one side of the hyperplace. Bo + B, X, + --+ PpXp <0 then X lies on the other side of the hyperplace

You can determine which side of the hyperplane by just determining the sign of

#### 1.1 Classificaton Using a Separating Hyperplane

Suppose that we have a  $n \times p$  data matrix  $\boldsymbol{X}$  that consists of n training observations in p-dimensional space.

$$\underline{\mathcal{X}}_{i} = \begin{pmatrix} \mathcal{X}_{i} \\ \vdots \\ \mathcal{X}_{ip} \end{pmatrix} \qquad j = -j \quad \underline{\mathcal{X}}_{ip} = \begin{pmatrix} \mathcal{X}_{ij} \\ \vdots \\ \mathcal{X}_{np} \end{pmatrix}$$
freahuly observations.

and that these observations fall into two classes.

We also have a test observation.

$$p$$
-vector of observed features:  
 $\underline{\mathcal{X}}^{*} = (\underline{\mathcal{X}}^{*}_{1,3}, -3, \underline{\mathcal{X}}^{*}_{p})^{T}$ 

Our Goal: Develop a classifier based on training dotes but will uncitly classify the test observation based on feature measurements.

We will see a new approach using a separating hyperplane

Suppose it is possible to construct a hyperplane that separates the training observations perfectly according to their class labels.



Then a separating hyperplane has the property that

$$\beta_{0} + \beta_{i} x_{i_{1}} + \dots + \beta_{p} x_{i_{p}} > 0 \quad i \neq \forall i = 1 \quad \text{and} \quad f_{r} \quad i = 1, \dots, n$$

$$\beta_{0} + \beta_{i} x_{i_{1}} + \dots + \beta_{p} x_{i_{p}} \neq 0 \quad i \neq \forall i = 1, \dots, n$$

$$\forall j \quad \left( \beta_{0} + \beta_{p} x_{i_{1}} + \dots + \beta_{p} x_{i_{p}} \right) = 0 \quad f_{r} \quad i = 1, \dots, n.$$

If a separating hyperplane exists, we can use it to construct a very natural classifier:

a fest dosnation is assigned to a class depending on which side of the hyperplace it is backed.

That is, we classify the test observation  $x^*$  based on the sign of  $f(x^*) = \beta_0 + \beta_1 x_1^* + \cdots + \beta_p x_p^*$ .

- If f(x\*) > 0 assign x\* to class 1
- IF floc\*) <0 assign oc\* to class -1.

We can also use the magnitude of  $f(x^*)$ .

- If  $f(\mathbb{Z}^{\#})$  is is far from zero, this means some lies far from the hyperplane we can be confident about our class assignment for zem
- If f (20t) is close to Zero, if is located near the hyperplane >> we are less sure about class assignment.

#### 1.2 Maximal Margin Classifier

If our data cab be perfectly separated using a hyperplane, then there will exist an infinite number of such hyperplanes.



A natural choice for which hyperplane to use is the *maximal margin hyperplane* (aka the *optimal separating hyperplane*), which is the hyperplane that is farthest from the training observations.

- We compute the perpendicular distance from each observation the a given supervised hyperplane - the smallest distance is called the margin.

The maximal margin hyperplane is the hyperplane of the largest margin, i.e. farthest from all trainity points,



We can then classify a test observation based on which side of the maximal margin hyperplane it lies – this is the *maximal margin classifier*.

- When p is large, ourfitting can occur.

1

We now need to consider the task of constructing the maximal margin hyperplane based on a set of n training observations and associated class labels.

ション・・・ショーモル アン・・・ショーモシーリ、13.

The maximal margin hyperplane is the solution to the optimization problem:

- (3) means each observation in training data set will be on the correct side of the hyperplane (MZO) with some cushion if (MZO).
- (2) ensures yi (Bo+Bizcii+--+Bpzip) is perp. distance to the hyperplane and (3) means the point z; is otherst Manuary =7 defines M as the Margin.
- (1) chooses Bos--sfp, M the maximize the margin.

=> maximal margin hyperplane!

This problem can be solved efficiently, but the details are outside the scope of this course.  $\downarrow_{\mathcal{P}}$  we'll talk more later...

What happens when no separating hyperplane exists?

=> no maximal margin hyperplane!

We can durelop a hyperplane that <u>almost</u> Separates the classes L> a "soft margin"



# 2 Support Vector Classifiers

It's not always possible to separate training observations by a hyperplane. In fact, even if we can use a hyperplane to perfectly separate our training observations, it may not be desirable.

A classifier band on a persently separating (manimul Mappin) hyperplane can lead to oversus inity to individual observations (high variability).



7

We might be willing to consider a classifier based on a hyperplane that does *not perfectly* separate the two classes in the interest of

- · greater robustness to individual observations
- · proper classificition of most of the trachily observations.
- i.e. it might be worthville to missdanity a few observations in training data set to do abitter job classifying a few fest data set.

" soft margin classifier"

The *support vector classifier* does this by finding the largest possible margin between classes, but allowing some points to be on the "wrong" side of the margin, or even on the "wrong" side of the hyperplane.

The support vector classifier glassifies a test observation depending on which side of the hyperplane it lies. The hyperplane is chosen to correctly separate **most** of the training observations.

Solution to The following optimization problem: maximize M βο,βο,-,βρ, ε..., En, M Subject to  $\sum_{i=1}^{p} \beta_{i}^{2} = 1$  $\eta_i \left( \beta_0 + \beta_i x_{i1} + \dots + \beta_p x_{ip} \right) = M \left( 1 - \xi_i \right)$  $\mathcal{E}_i \ge 0$ ,  $\sum_{i=1}^{\infty} \mathbb{E}_i \le C$   $\mathcal{I}$  nonnegative tuning parameter  $\mathcal{I}_i$  slack variables" (budget for how wrong we are willing to be in training dota). allow observations to be on The wary side of margin ( or hyperplane)

Once we have solved this optimization problem, we classify  $x^*$  as before by determining which side of the hyperplane it lies.

classify x\* Land on sigh of f(x\*)= Bo+B,x1\*+ ...+ Bp x \*.

 $\epsilon_i$  - tells us where the observation lies relative to hyperplane and margin.

If  $2i=0 \Rightarrow obs.$  on conect side of Margin.

- If  $\xi_1 = 0 \implies obs$  on wrong side of margin If  $\xi_1 = 1 \implies obs$  on wrong side of hyperplace.
- C tuning porameter, bounds the sum of Ei's => determines # and severity of violations we will allow think of C as a budget for the amount of volations.
  - If  $C=0 \Rightarrow$  no budget for violations  $\Rightarrow E_1 = ... = E_n = 0 \Rightarrow$  SU classifier = maximal margin classifier.

If  $C = 0 \Rightarrow$  no more than C ois. can be on the wrong side of the hyperplane because  $\varepsilon_i^{(2)}$  and  $\hat{\Sigma} = \hat{\varepsilon}_i^{(2)} \leq C$ .

Jmall C => narrow margins large C => wider morgin, allow for more iriolations trade-off. cross-volidation!

The optimization problem has a very interesting property.

```
only obserations on the Morgin or vidade the Margin (or hyperplane) affect the hyperplane!

> The classifier.
```

i.e. observations that lie on correct side of magin do not affect the support vator classifier!

Observations that lie directly on the margin or on the wrong side of the margin are called *support vectors*.

```
These observations do effect the classifier.
```

The fact that only support vectors affect the classifier is in line with our assertion that C controls the bias-variance tradeoff.

When C small => Server support vectors => low bias but high variance.

Because the support vector classifier's decision rule is based only on a potentially small subset of the training observations means that it is robust to the behavior of observations far away from the hyperplane.

```
distinct from Librarior of other methods
eng. LDA depicts on mean of observations in each class.
```

## **3** Support Vector Machines

The support vector classifier is a natural approach for classification in the two-class setting... if the decision boundary is linear!



We've seen ways to handle non-linear classification boundaries before.

```
neahinear basis function + Logistic regression, KNN, QDA
```

In the case of the support vector classifier, we could address the problem of possible nonlinear boundaries between classes by enlarging the feature space.

```
e.g. add guadratic or cubic terms
instead of filting SV dassition w/ X10-02Xp
could use X10-02Xp, X<sup>2</sup>, ..., Xp, etc.
```

Then our optimization problem would become

could consider higher order polynomials or other functions.

The support vector machine allows us to enlarge the feature space used by the support classifier in a way that leads to efficient computation.  $\rightarrow using "kernels"$ 

Want to enlarge fature space to have non-linear boundary.

It turns out that the solution to the support vector classification optimization problem involves only *inner products* of the observations (instead of the observations themselves).

inner product 
$$\langle \underline{a}, \underline{b} \rangle = \sum_{i=1}^{\infty} a_i b_i$$
  
inner product of 2 obs.  $\langle \underline{x}_i, \underline{x}_i \rangle = \frac{\underline{b}}{\underline{b}_{i+1}} \underline{x}_i, \underline{x}_i$ 

It can be shown that

- The solution to (lincer) support vector classifier on be written as  $f(x) = \beta_0 + \sum_{i=1}^{n} \alpha_i < \sum_{j=1}^{n} \alpha_{j,j} > \sum_{i=1}^{n} \alpha_{i,j} = 1, ..., n \text{ additional parameters.}$
- To estimate  $a_{(2-1)}a_n$  and fo und  $\binom{n}{2}$  inner products between all pairs of training observations.
- O(i) Monzero only for support vectors in the solution!
   ⇒ typically list than n.
   ⇒ rewrite f(x) = Bo + Σαi < x, xi?</li>
   B = indias & support vectors.

Now suppose every time the inner product shows up in the SVM representation above, we replaced it with a generalization.  $\sum \langle \mathcal{Z}, \mathcal{Z}, \mathcal{Z} \rangle$ 

Kernel: K(Zi,Zir). some function. Les a function that quantities similarity of two observations.

e.q. 
$$K(\underline{x}_{i},\underline{x}_{i}) = \sum_{j=1}^{r} \mathbb{I}_{ij} \mathbb{I}_{i'j}$$
 results in support vector classifier "kincar kernel"  

$$\begin{cases}
K(\underline{x}_{i},\underline{x}_{i'}) = (1 + \sum_{j=1}^{r} \mathbb{I}_{ij} \mathbb{I}_{i'j})^{d} \leftarrow pos. integer \\
K(\underline{x}_{i},\underline{x}_{i'}) = \exp(1 + \sum_{j=1}^{r} \mathbb{I}_{i'j} \mathbb{I}_{i'j})^{d} \\
K(\underline{x}_{i},\underline{x}_{i'}) = \exp(-\mathcal{E}\sum_{j=1}^{r} (\mathcal{R}_{ij} - \mathbb{I}_{i'j})^{2}) \\
pos. construct
\end{cases}$$
vadial kernel"

Compartation whore the support den.



L7 radial kernel - enlarged featurespace is infinite dimensional!

### 4 SVMs with More than Two Classes

So far we have been limited to the case of binary classification. How can we exted SVMs to the more general case with some arbitrary number of classes?

This is not that clear. There is no one way to do this.

#### Two optims!

Suppose we would like to perform classification using SVMs and there are K > 2 classes.

#### **One-Versus-One**

(1) Construct (<sup>K</sup><sub>2</sub>) SVMs each comparing a pair of classes.
 (2) classify a test observation using each of the (<sup>K</sup><sub>2</sub>) <u>SVMs</u>
 (3) Assign test observation to class it was most frequently predicted.

One-Versus-All lit 2 the a fest observation.

- 1) Fit K SUMs imparing each dass to remaining K-1 classes.
- (2) assign JC\* to the class for which Box + BIR 20t + .-+ Bpx 20pt is largest.

results in high level of confidence fest observation belongs to ken dass our any other.