

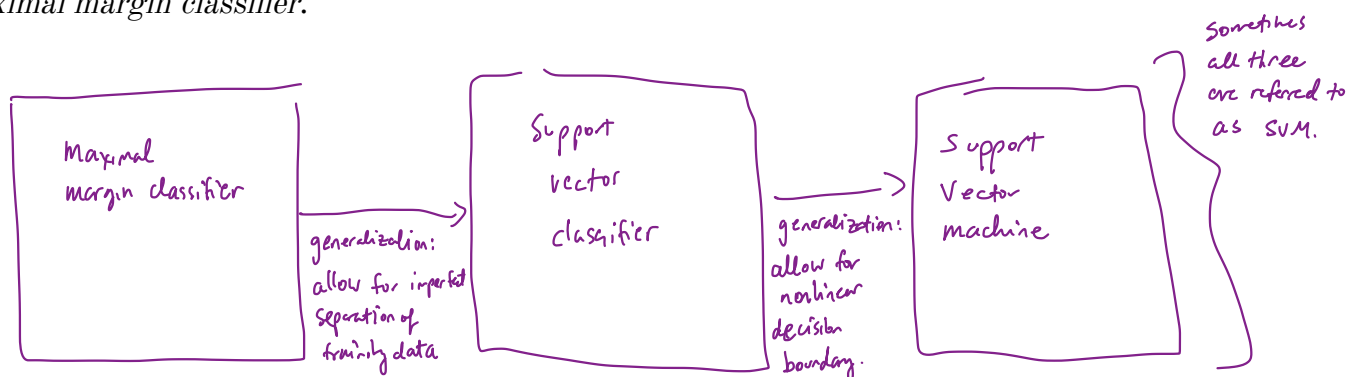
Chapter 9: Support Vector Machines

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 well in variety of settings

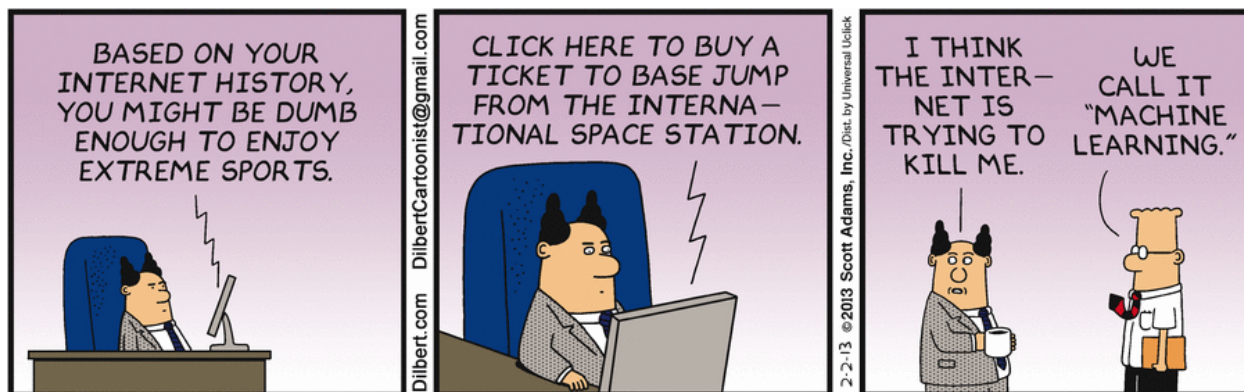
often considered one of the best "out of the box" classifiers.

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.

↳ categorical response w/ 2 levels



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

1 Maximal Margin Classifier

→ based on a hyperplane separator

→ extension of euclidean space.

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

eg. in 2D, a hyperplane is a flat 1 dim subspace - a line.

in 3D, a hyperplane is a flat 2 dim subspace - a plane.

in $p > 3$ dimension, hyperplane is harder to conceptualize, but still a flat $p-1$ dim. subspace.

The mathematical definition of a hyperplane is quite simple,

In 2D, a hyperplane is by $\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$.

i.e., any $X = (X_1, X_2)$ for which this equation holds, lies on the hyperplane.

Note this is just equation for a line.

This can be easily extended to the p -dimensional setting.

$\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p = 0$ define a p -dim hyperplane.

i.e., any $X = (X_1, \dots, X_p)$ for which this equation holds lies on the hyperplane.

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

If $\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p > 0$ then $X = (X_1, \dots, X_p)$ lies on one side of the hyperplane and

If $\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p < 0$ then X lies on the other side of the hyperplane.

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

$$\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

1.1 Classifier Using a Separating Hyperplane

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

$$\underline{x}_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1p} \end{pmatrix}, \dots, \underline{x}_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{np} \end{pmatrix}$$

↑
training observations.

and that these observations fall into two classes.

$$y_1, \dots, y_n \in \{-1, 1\}$$

where -1 represents one class
 1 other class.

We also have a test observation.

p -vector of observed features

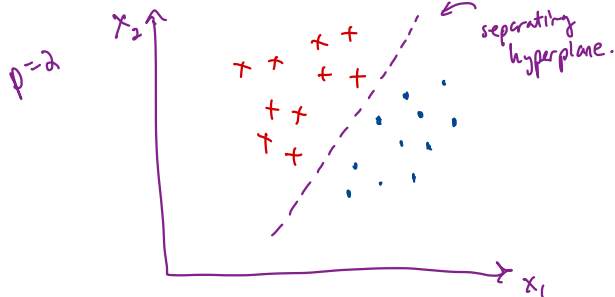
$$\underline{x}^* = (x_1^*, \dots, x_p^*)^T$$

Our Goal: Develop a classifier based on training data that will correctly classify the test observation based on feature measurements.

We have already used many approaches:
 logistic regression (penalized version, i.e. LASSO or ridge)
 LDA
 classification trees
 Bagging
 Random forests
 Boosting, etc.

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

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} > 0 \quad \text{if } y_i = 1 \quad \text{and}$$

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} < 0 \quad \text{if } y_i = -1$$



$$y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) > 0 \quad \forall i = 1, \dots, n.$$

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

A test observation is assigned to a class depending on which side of the hyperplane it is located.

That is, we classify the test observation x^* based on the sign of

$$f(x^*) = \beta_0 + \beta_1 x_1^* + \dots + \beta_p x_p^*.$$

If $f(x^*) > 0$ assign x^* to class 1

$f(x^*) < 0$ assign x^* to class -1

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

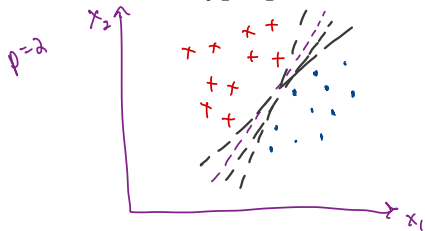
If $f(x^*)$ is far from zero (large magnitude) means x^* lies far from the hyperplane
 \Rightarrow we can be confident about our class assignment for x^* .

If $f(x^*)$ is close to zero (small magnitude) it is located near the hyperplane
 \Rightarrow we are less confident about the class assignment for x^* .

Note: a classifier based on separating hyperplane leads to a linear decision boundary.

1.2 Maximal Margin Classifier

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



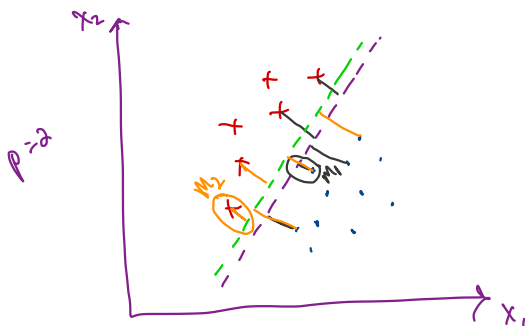
A given separating hyperplane can be shifted a tiny bit or rotated without coming into contact w/ any training observations

⇒ Which one to use for our classifier?

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 to a given separating hyperplane.
- the smallest distance is known as the "margin"

The maximal margin hyperplane is the one w/ the largest margin, i.e. farthest from all training points.



$$M_2 > M_1$$

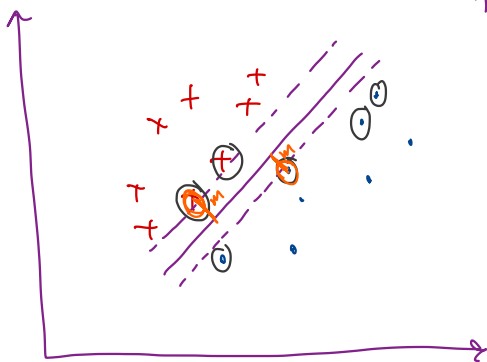
⇒ larger margin

⇒ 2nd hyperplane is preferred.

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

Hopefully a large margin on training data will lead to a large margin on test data
⇒ classify test data correctly

When p is large, we can see over fitting.



the two equidistant points from the maximal margin hyperplane are known as support vectors because they are p -dim vectors that "support" the hyperplane.

i.e. if these points move, the maximal margin hyperplane would move as well.

↙ a small # of points.

NOTE: The maximal margin hyperplane only depends on the support vectors!
the rest of the points can move and it doesn't matter.

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.

$$x_1, \dots, x_n \in \mathbb{R}^p$$

$$y_i \rightarrow y_n \in \{-1, 1\}$$

The maximal margin hyperplane is the solution to the optimization problem

① maximize M ^{margin}

$$\beta_0, \beta_1, \dots, \beta_p, M$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1 \quad \textcircled{2}$$

$$\textcircled{3} y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M \quad \forall i=1, \dots, n.$$

③ means each observation will be on the correct side of the hyperplane ($M \geq 0$) with some cushion (if $M > 0$).

② ensures $y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})$ is perp. distance to hyperplane, and ③ means the point is at least M distance away \Rightarrow M is the margin.

① chooses $\beta_0, \beta_1, \dots, \beta_p, M$ to maximize the margin
 \Rightarrow maximal margin hyperplane!

This problem can be solved efficiently, but the details are outside the scope of this course.

\hookrightarrow we'll talk a little bit more later.

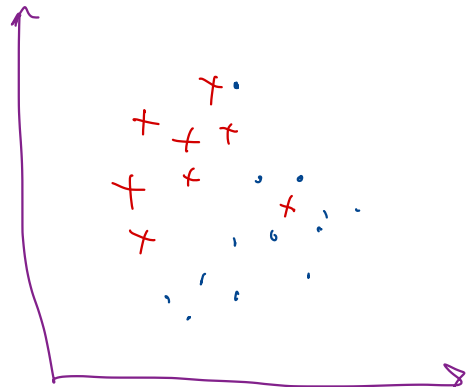
What happens when no separating hyperplane exists?

\Rightarrow no maximal margin hyperplane!

We can develop a hyperplane that almost separates the classes

\uparrow

a "soft margin"



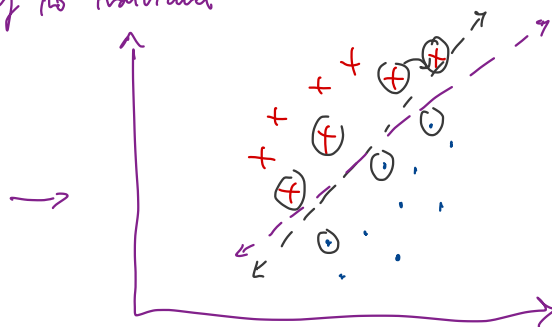
2 Support Vector Classifiers

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 based on a separating hyperplane will necessarily perfectly classify all training observations.

↳ This can lead to oversensitivity to individual observations.

shifting / adding one data point can result in dramatic change in hyperplane (and margin).



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

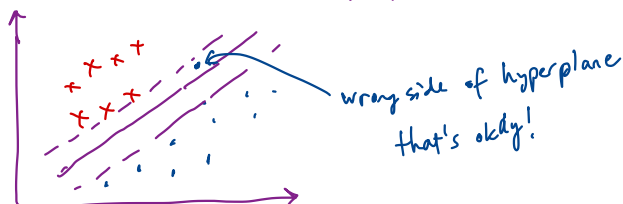
- greater robustness to individual observations
- better classification of most of observations.

i.e. could be worthwhile to misclassify a few observations to do a better job classifying the rest (also test observations).

↳ sometimes called "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.

↳ when there is no separating hyperplane this is inevitable.



The support vector classifier classifies 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 ← margin

$\beta_0, \beta_1, \dots, \beta_p, \varepsilon_1, \dots, \varepsilon_n, M$

subject to

$$\sum_{j=1}^p \beta_j^2 = 1$$

$$y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M (1 - \varepsilon_i)$$

$$\varepsilon_i \geq 0, \quad \sum_{i=1}^n \varepsilon_i \leq C$$

↗

"slack variables"

allow observations to be on the wrong 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^* based on sign of $f(x^*) = \beta_0 + \beta_1 x_1^* + \dots + \beta_p x_p^*$

ε_i - tells us where observations lie relative to hyperplane & margin.

if $\varepsilon_i = 0 \Rightarrow$ obs on correct side of margin.

if $\varepsilon_i > 0 \Rightarrow$ obs on wrong side of margin (violated margin)

if $\varepsilon_i > 1 \Rightarrow$ obs on wrong side of hyperplane.

C - tuning parameter, bounds sum of ε_i 's \Rightarrow determines # and severity of violations we will allow
think of C as a budget for amount of violations

if $C = 0 \Rightarrow$ no budget for violations $\Rightarrow \varepsilon_1 = \dots = \varepsilon_n = 0 \Rightarrow$ support vector classifier = maximal margin classifier (if exists).

if $C > 0 \Rightarrow$ no more than C observations can be on the wrong side of the hyperplane.

because $\varepsilon_i > 1$ and $\sum_{i=1}^n \varepsilon_i \leq C$.

small $C \Rightarrow$ narrow margin, large $C \Rightarrow$ wider margins and allows for more violations

C controls bias-variance trade-off.

choose
 C by
CV

The optimization problem has a very interesting property.

only observations on the margin or violate the margin (or hyperplane) affect the hyperplane \Rightarrow the classifier!

i.e. observations that lie strictly on the correct side of the margin don't affect the classifier!

Observations that lie directly on the margin or on the wrong side of the margin ^{or hyperplane} are called support vectors.

these observations do affect 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 large \Rightarrow margin is wide, many observations violate the margin \Rightarrow many support vectors
 i.e. many observations used to determine the hyperplane
 \Rightarrow low variance but potentially high bias.

When C small \Rightarrow fewer support vectors
 \Rightarrow 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 behavior of other classifier methods.

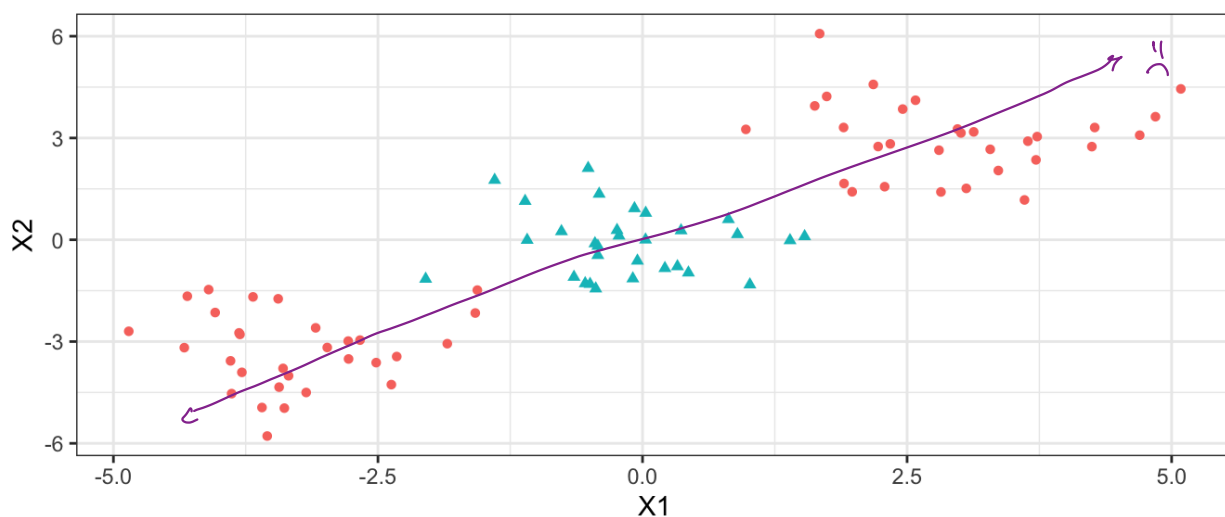
e.g. LOA depends on the mean of all observations within each class + within class covariance matrix.

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!

Sometimes we have nonlinear boundaries:

How to draw a line separating?
won't work well.



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

QDA, bagging, RF, boosting trees

logistic regression w/ polynomial features, etc.

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

e.g. adding quadratic or cubic terms

instead of fitting SV classifier w/ X_1, \dots, X_p

could use $X_1, \dots, X_p, X_1^2, \dots, X_p^2$ etc.

Then our optimization problem would become

Maximize M

$\beta_0, \beta_{11}, \beta_{12}, \dots, \beta_{1p}, \beta_{21}, \dots, \beta_{2p}, \varepsilon_1, \dots, \varepsilon_n, M$

Subject $\sum_{j=1}^p \sum_{k=1}^2 \beta_{kj} = 1$

$$y_i \left(\beta_0 + \sum_{j=1}^p \beta_{1j} X_{ij} + \sum_{j=1}^p \beta_{2j} X_{ij}^2 \right) \geq M(1 - \varepsilon_i)$$

$$\sum_{i=1}^n \varepsilon_i \leq C$$

quadratic polynomial leads to nonlinear boundary.

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.

→ using "kernels"

Want to enlarge feature 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).

$$\text{inner product } \langle a, b \rangle = \sum_{i=1}^r a_i b_i$$

$$\text{inner product of 2 obs.: } \langle x_i, x_j \rangle = \sum_{j=1}^p x_{ij} x_{kj}$$

It can be shown that

- The linear support vector classifier can be represented as

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$$
 , $\alpha_i, i=1, \dots, n$ additional parameters.
- To estimate $\alpha_1, \dots, \alpha_n$ and β_0 need $\binom{n}{2} = \frac{n(n-1)}{2}$ inner products between all pairs of training observations

- α_i nonzero only for support vectors in the solution!
 - ↳ typically less than n .
 - ↳ only need inner products w/ support vectors (not observations).

$$\Rightarrow \text{rewrite } f(x) = \beta_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle, \quad S = \text{indices of support vectors}$$

Now suppose every time the inner product shows up in the SVM representation above, we replaced it with a generalization.

$$\curvearrowright \langle x, x_i \rangle.$$

Kernel: $K(x_i, x_j)$ (some function).

↳ a function that quantifies similarity of two observations

e.g. $K(x_i, x_j) = \sum_{j=1}^p x_{ij} x_{kj}$ results in support vector classifier "linear kernel" b/c linear boundary.

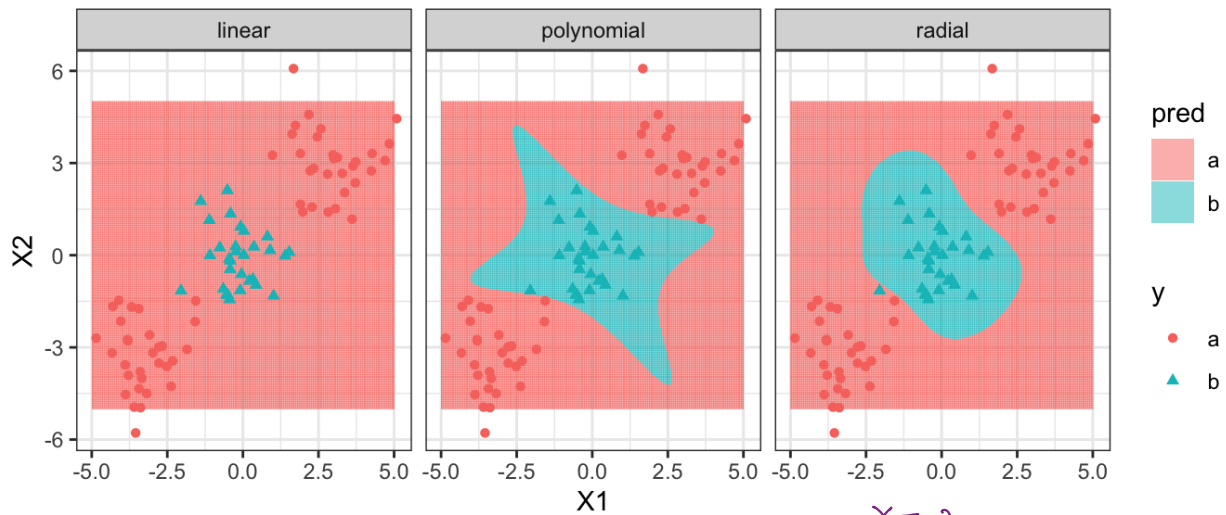
$$\left\{ \begin{array}{l} K(x_i, x_j) = \left(1 + \sum_{j=1}^p x_{ij} x_{kj} \right)^d \quad d \leftarrow \text{pos. integer} \\ K(x_i, x_j) = \exp \left(-\gamma \sum_{j=1}^p (x_{ij} - x_{kj})^2 \right) \end{array} \right. \quad \begin{array}{l} \text{"polynomial kernel"} \\ \text{"radial kernel"} \end{array}$$

↑
pos. constant.

"support vector machine"

Computation of SV classifier... idea.

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i K(x, x_i).$$



↗
When polynomial kernel w/
 $d=1$
same as linear.

$d=4$
= fitting support
vector classifier in
a higher dimension
feature space w/
polynomial of degree d .

$\gamma = 2$
 $K(x_i, x_j) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{ij}^*)^2)$
If a test observation i^* is far from a training
observation, then
 x_i^* $\sum_{j=1}^p (x_{ij}^* - x_{ij})^2$
will be large.

$\Rightarrow K(x^*, x_i)$ will be very small.

\Rightarrow will play a ^{very} small role in $f(x^*)$.

\Rightarrow training observations far from x^* will
play little role in prediction of class
of x^* .

Why use a kernel instead of enlarging the feature space using functions of features?

Computational - Only need to compute $K(x_i, x_j)$ at $\binom{n}{2}$ distinct pairs i, j

- don't have to explicitly work in enlarged space (may be too large to compute
hyperplane).

= radial kernel - feature space 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 extend SVMs to the more general case with some arbitrary number of classes?

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

Two popular options:

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

One-Versus-One

- ① Construct $\binom{K}{2}$ SVMs each comparing a pair of classes.
- ② Classify a test observation using each of the $\binom{K}{2}$ SVMs
- ③ Assign test observation to class it was most frequently assigned.

One-Versus-All Let x^* be a test observation.

- ① Fit K SVMs comparing each class to remaining $K-1$ classes.
- ② assign x^* to the class for which $\beta_{0k} + \beta_{1k}x_1^* + \dots + \beta_{pk}x_p^*$ is the largest

(results in high level of confidence test observation belongs to the k^{th} class over any other).