# Chapter 9: Support Vector Machines

*categorical response y*
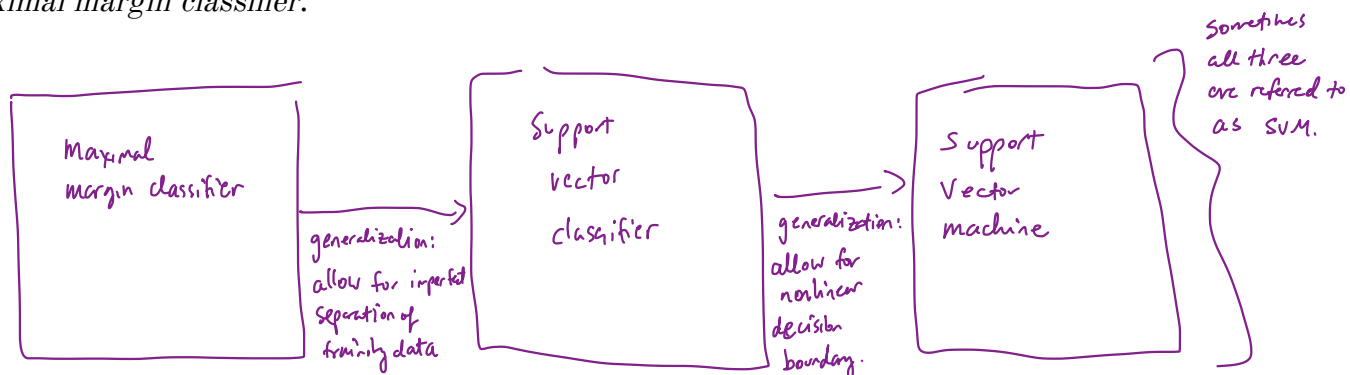
The *support vector machine* is an approach for classification that was developed in the computer science community in the 1990s and has grown in popularity.
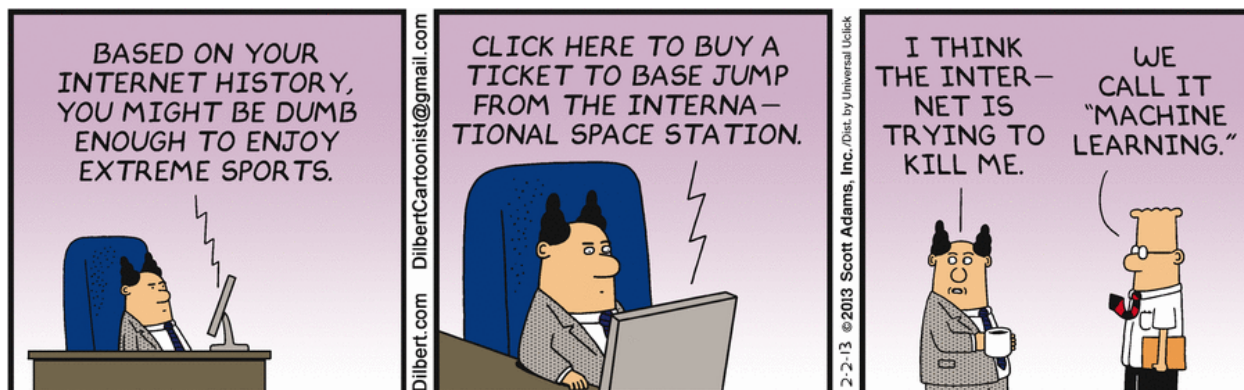
SVMs perform well in variety of settings

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*.

Sometimes all three are referred to as SVM.

Maximal margin classifier
→ generalization: allow for imperfect separation of training data
→ Support vector classifier
→ generalization: allow for nonlinear decision boundary.
→ Support Vector machine

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 $p - 1$.

e.g. 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 $\underline{\beta_0 + \beta_1 X_1 + \beta_2 X_2} = 0$.

parameters

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.

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

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

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

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

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

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

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

## 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.

$$x_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1p} \end{pmatrix}, \ldots, \underline{x}_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{np} \end{pmatrix}$$

training observations.

and that these observations fall into two classes.

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

where $-1$ represents one class

$1$ other class.

We also have a test observation.
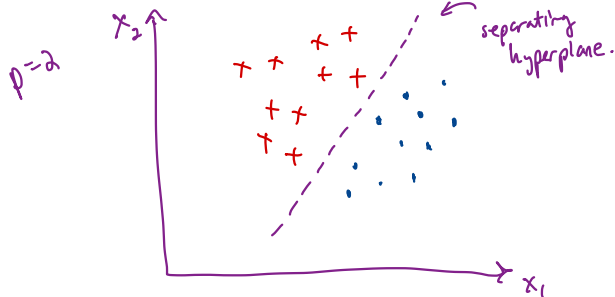
p-vector of observed features

$$x^* = (x_1^*, \ldots, x_p^*)^T$$

**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 <u>separating hyperplane</u>.

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

$p=2$



Then a separating hyperplane has the property that

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

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

$$\Longleftrightarrow$$

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

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

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

$f(x^*) < 0$ assign $x^*$ to class $-1$
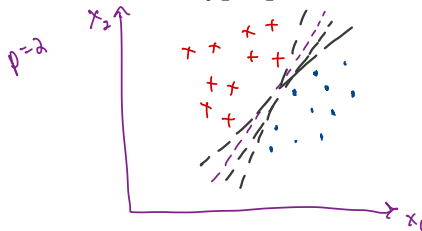
We can also use the magnitude of $f(x^*)$.

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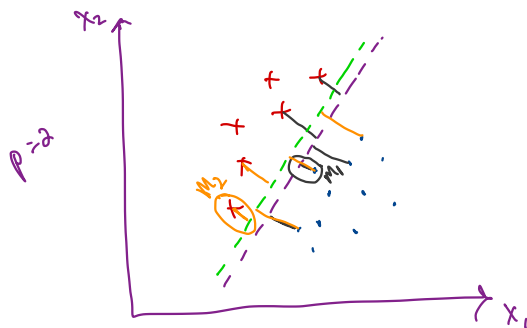
A given separating hyperplane can be shifted a tiny bit or rotated without coming into contact w/ any training observations

$\implies$ 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.
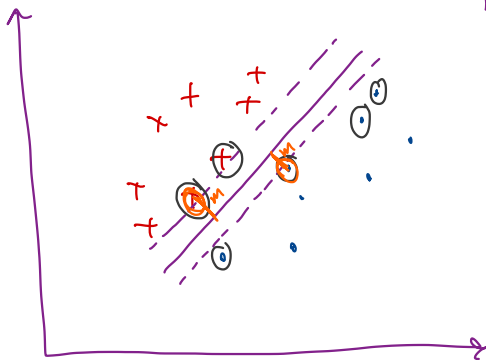
$M_2 > M_1$

$\implies$ larger margin

$\implies$ 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
$\implies$ classify test data correctly

When p is large, we can see overfitting.

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.

$$x_1, \ldots, x_n \in \mathbb{R}^p \qquad\qquad y_1, \ldots, y_n \in \{-1, 1\}.$$

The maximal margin hyperplane is the solution to the optimization problem

① $\underset{\beta_0, \beta_1, \ldots, \beta_p, M}{\text{maximize}} \; M$ ← margin

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

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

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

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

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

This problem can be solved efficiently, but the details are outside the scope of this course.
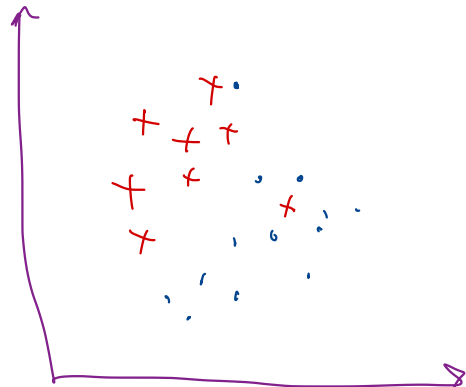
↳ we'll talk a little bit more later.

What happens when no separating hyperplane exists?

⟹ no maximal margin hyperplane!

We can develop a hyperplane that
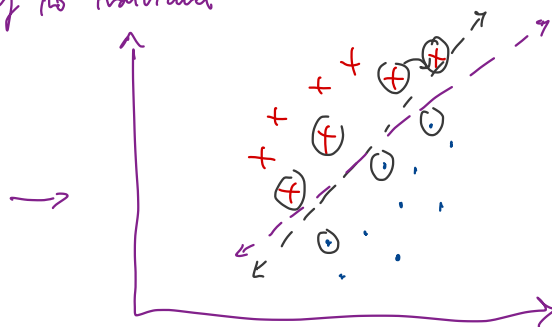almost separates the classes
↑
a "soft margin"

# 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 seprating 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
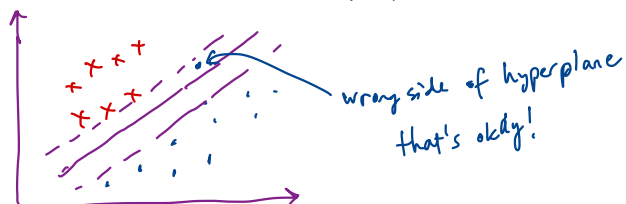
- greater robustness to individual observations

- better classification of __most__ of observations.

i.e. Could be worthwile 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.

wrong side of hyperplane that's okay!

7

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:

$$\text{maximize } M \leftarrow \text{margin}$$
$$\beta_0, \beta_1, \dots, \beta_p, \varepsilon_1, \dots, \varepsilon_n, M$$

subject to

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

$$y_i \left( \beta_0 + \beta_1 x_{i_1} + \dots + \beta_p x_{ip} \right) \geq M(1 - \varepsilon_i)$$

$$\varepsilon_i \geq 0, \quad \sum_{i=1}^{n} \varepsilon_i \leq C$$
$$\uparrow \qquad \qquad \curvearrowleft \text{ non negative tuning parameter}$$

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

$\epsilon_i$ — 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.

$\boxed{\text{choose } C \text{ by CV}}$ $C$ — tuning parameter, bounds sum of $\varepsilon_i$'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.

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 are called
*support vectors.*

or hyperplane

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.

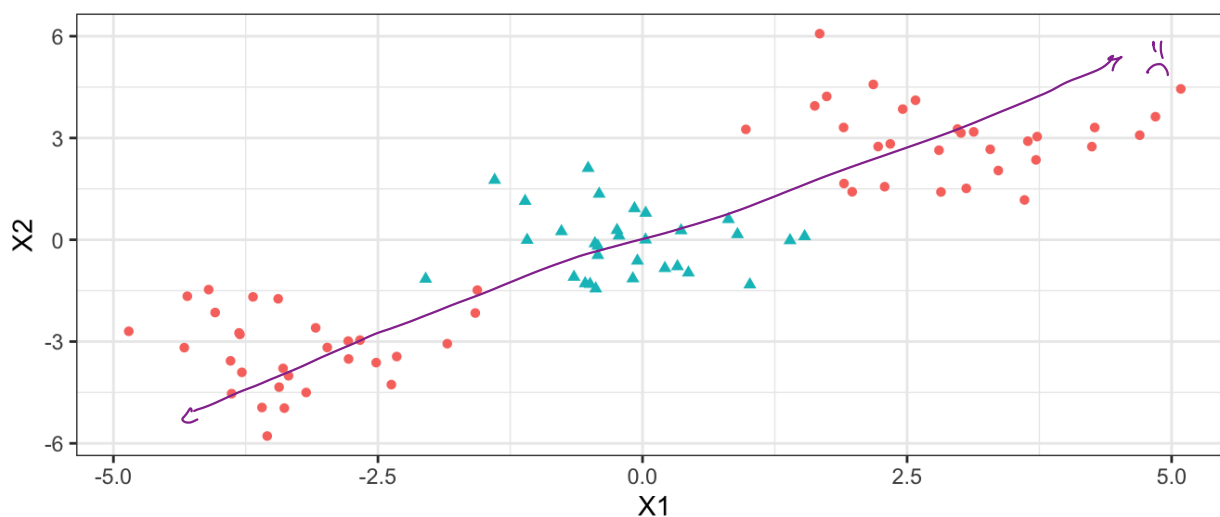distinct from behavior of other classifier methods.

e.g. LDA 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/ $X_1,...,X_p$*

*could use $X_1,...,X_p, X_1^2,...,X_p^2$ etc.*

Then our optimization problem would become

$$\text{Maximize } M \atop \beta_0, \beta_{11}, \beta_{12},...,\beta_{1p}, \beta_{21},...,\beta_{2p}, \varepsilon_1,...,\varepsilon_n, M$$

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

*quadratic polynomial leads to nonlinear boundary.*

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

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

*could consider higher order polynomials or other functions.*

*→ using "kernels"*

The *support vector machine* allows us to enlarge the feature space used by the support classifier in a way that leads to efficient computation.

*Want to enlarge feature space to have non-linear boundary.*

*Computation of SV classifier... idea.*

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).
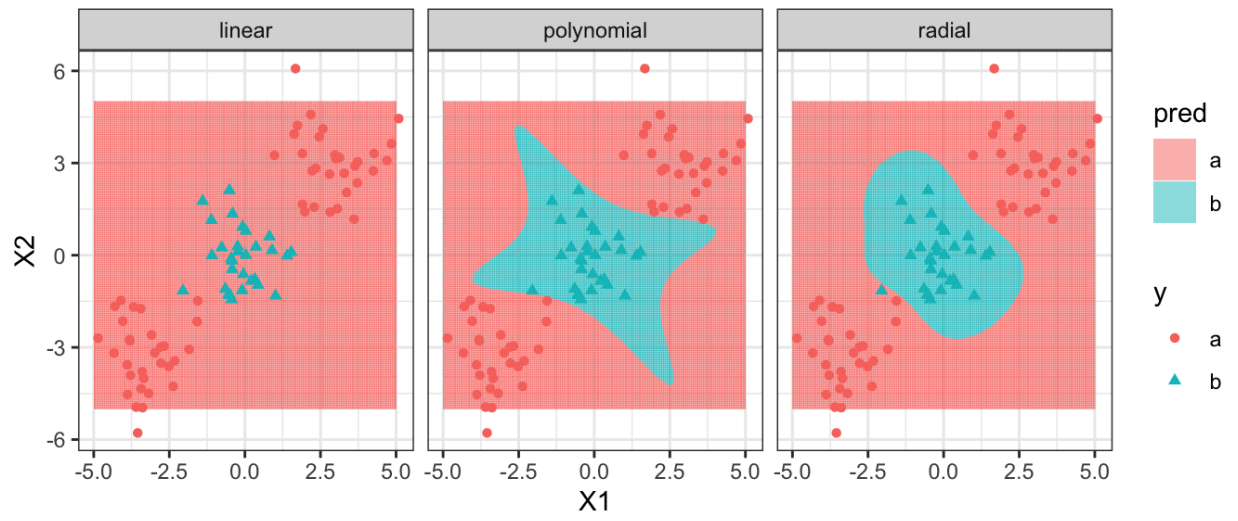
It can be shown that

- 

- 

- 

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

# 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?

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

**One-Versus-One**

**One-Versus-All**