

# 9 Support Vector Machines

Andrew Liang

11/23/2020

## Notes

- classification technique
- considered one of the best “out of the box” classifiers

### Maximal Margin Classifier

- first we need to define *hyperplane*
  - $p$  dimensional space
  - flat affine (doesn't need to pass through origin) subspace of a dimension  $p - 1$
  - the mathematical definition is:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p = 0$$

for parameters  $\beta_0, \beta_1, \beta_2$

- for any  $X = (X_1, X_2, \dots, X_p)^T$  that holds the equation above, defines a hyperplane
- if  $X$  does NOT satisfy the equation above, then:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p > 0$$

OR

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p < 0$$

which tells us that  $X$  must lie on either side of the hyperplane

### Classification Using a Separating Hyperplane

- suppose we have a  $n \times p$  data matrix  $\mathbf{X}$  consisting of  $n$  training observations in  $p$ -dimensional space:

$$x_1 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{1p} \end{bmatrix}, \dots, x_n = \begin{bmatrix} x_{n1} \\ \vdots \\ x_{np} \end{bmatrix}$$

and also suppose these observations fall into two distinct classes:  $y_1, \dots, y_n \in \{-1, 1\}$ .

We are also given a test observation, a  $p$ -vector of observed features  $x^* = (x_1^*, \dots, x_p^*)^T$ . Our goal is to classify test observation using its feature measurements through a concept of *separating hyperplane*.

Not only can we classify observations based on the sign of  $f(x^*)$ , but also the *magnitude* of it as well.

- if magnitude is large, then we know that the observation  $x^*$  lies far from the hyperplane, and vice versa
- separating hyperplane leads to a linear decision boundary

## Maximal Margin Classifier

- if the data can be perfectly separated using a hyperplane, then there exists an infinite number of such hyperplanes, as it can be shifted or rotated without coming into contact any of the observations. Must decide which of the infinite separating hyperplanes to use
- leads to the *maximal margin hyperplane* (AKA *optimal separating hyperplane*)
  - this is basically the hyperplane that is furthest from all the training observations
    - \* determined by the perpendicular distance from each data point to the hyperplane, called the *margin*
    - \* thus we try to maximize this distance, hence the name *maximal margin classifier*
    - \* can lead to overfitting when  $p$  is large
- the points that are closest and equidistant to this maximal margin hyperplane are called *support vectors*
  - they are vectors in  $p$  dimensions
  - thus the hyperplane are “supported” by these points and only depend on them, and not on any other observation
    - \* a movement in any of these support vectors would move the hyperplane as well

## Construction of Maximal Margin Classifier

Consider  $n$  training observations  $x_1, \dots, x_n \in \mathbb{R}^p$  associated with class labels  $y_1, \dots, y_n \in \{-1, 1\}$ . We try to solve the optimization:

$$\max_{\beta_0, \beta_1, \dots, \beta_p, M} M$$

subject to:

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

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

- the second constraint ensures that each observation is on the correct side of hyperplane

**Non-separable Case** \* however, if no hyperplane exists, then maximal margin classifier won't exist

## Support Vector Classifiers (AKA soft margin classifier)

- we now consider a classifier that does *not* perfectly separate the classes, but *most* of it, this gives a couple of advantages:
  - greater robustness to individual observations
  - better classification of *most* of the training observations

## Construction of Support Vector Classifiers

$$\max_{\beta_0, \dots, \beta_p, \epsilon_1, \dots, \epsilon_p, M} M$$

subject to:

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

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

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

Where C is a non-negative tuning parameter.

- $\epsilon_1, \dots, \epsilon_n$  are *slack variables* that allow individual observations to be on wrong side of the margin or the hyperplane
  - if  $\epsilon_i > 0$ , then the  $i$ th observation is on wrong side of the margin
  - if  $\epsilon_i > 1$ , then it is on the wrong side of the hyperplane
- $C$  bounds the sums of  $\epsilon_i$ 's so it determines number and severity of the violations
  - if  $C = 0$  then there is no budget for violations
  - for  $C > 0$ , no more than  $C$  observations can be on wrong side of hyperplane
  - generally chosen via CV
- observations that lie on the margin or violate the margin will affect the hyperplane (known as *support vectors*), and so observations that lie on the correct side of the margin does not affect the support vector classifier
- meaning that it is robust to observations far away from the hyperplane
- different from LDA classifications where it depends on *all* of the observations within each class
- similar to logistic where it is robust to observations far from decision boundary

## Support Vector Machines

### Classification with Non-linear Decision Boundaries

- can create non-linear boundaries by expanding feature space (ie enabling quadratic and cubic terms)

Rather than fitting support vector classifier using  $p$  features, we can fit a support vector classifier using  $2p$  features:

$$X_1, X_1^2, \dots, X_p, X_p^2$$

and so our maximization problem would be:

$$\max_{\beta_0, \beta_{11}, \beta_{12}, \dots, \beta_{p1}, \beta_{p2}, \epsilon_1, \dots, \epsilon_p, M} M$$

subject to:

$$y_i(\beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2) \geq M(1 - \epsilon_i)$$

$$\epsilon_i \geq 0, \quad \sum_{i=1}^n \epsilon_i \leq C, \quad \sum_{j=1}^p \sum_{k=1}^2 \beta_{jk}^2 = 1$$

\* in this case, the hyperplane is non-linear because it is a quadratic polynomial within the original feature space

## Support Vector Machines (SVMs)

- extension of support vector classifier and uses *kernels* to enlarge feature space
  - allows for efficient computational approach from using kernels
- calculation involves taking *inner products* (dot products) of observations. Inner product of two observations  $x_1, x_{i'}$  is given by:

$$\langle x_i, x_{i'} \rangle = \sum_{j=1}^p x_{ij} x_{i'j}$$

- and so the linear support vector classifier can be shown as:

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$$

where there are  $n$  parameters  $\alpha_i$ ,  $i = 1, \dots, n$ , one per training observation

- to estimate these parameters  $\alpha_1, \dots, \alpha_n$  and  $\beta_0$ , we just need  $\binom{n}{2}$  inner products of  $\langle x_i, x_{i'} \rangle$  between all pairs of observations

It turns out that  $\alpha_i$  is nonzero for only the *support points*. And so if  $S$  is the collection of indices for these support points, we can rewrite:

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle$$

which typically involves fewer terms than the previous equation

We replace the inner product calculation with a *generalization* of the inner product in the form:

$$K(x_i, x_{i'})$$

where  $K$  is some function that is call the *kernel*

- the kernel is a function that quantifies the similarity of two observations, for example:

$$K(x_i, x_{i'}) = \sum_{j=1}^p x_{ij} x_{i'j}$$

is just a support vector classifier (linear). The linear kernel essentially quantifies the similarity of a pair of observations using Pearson correlation. One could also replace the linear kernel with:

$$K(x_i, x_{i'}) = (1 + \sum_{j=1}^p x_{ij}x_{i'j})^d$$

AKA a *polynomial kernel* of degree  $d$ , where  $d$  is a positive integer. With  $d > 1$ , the support vector classifier algorithm leads to a more flexible decision boundary

- does this by fitting a support vector classifier in a higher-dimensional space involving polynomials of degree  $d$ , rather than in the original feature space
- the process of combining a support vector classifier and a non-linear kernel is called a *support vector machine*
- the function has the form:

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

Another popular choice is the *radial kernel*, which is:

$$K(x_i, x_{i'}) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2)$$

where  $\gamma$  is a positive constant (essentially a tuning parameter)

- if given test observation  $x^* = (x_1^* \dots x_p^*)^T$  is far from the training observation  $x_i$  in terms of Euclidean distance, then  $\sum_{j=1}^p (x_j^* - x_{ij})^2$  will be very large, and thus  $K(x_i, x_{i'}) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2)$  will be small. Thus, training observations far from  $x^*$  will not play a role in the predicted class label for  $x^*$ 
  - means that the radial kernel has a very *local* behavior
- advantage of using kernels over expanding feature space is that it's more computationally feasible
  - only need to compute  $K(x_i, x_{i'})$  for all  $\binom{n}{2}$  distinct pairs  $i, i'$

## SVMs with More than Two Classes

### One-Verses-One Classification (all-pairs)

- suppose we want to classify using SVMs when there are  $K > 2$  classes
- this method constructs  $\binom{K}{2}$  SVMs, each of which compares a pair of classes
- classify a test observation using each of the  $\binom{K}{2}$  classifiers, and tally the number of times the test obs is assigned to each of the  $K$  classes
- final classification is performed by assigning the obs to class which it was most frequently assigned in these  $\binom{K}{2}$  pairwise classifications
- usually used for smaller number of classes

## One-Verses-All Classification

- We fit  $K$  SVMs, each time comparing one of the  $K$  classes to the remaining  $K - 1$  classes
- let  $\beta_{0k}, \beta_{1k}, \dots, \beta_{pk}$  denote parameters that result from fitting an SVM comparing the  $k$ th class to the others
- let  $x^*$  denote test observation, we then assign observation to the class for which  $\beta_{0k} + \beta_{1k}x_1^* + \dots + \beta_{pk}x_p^*$  is largest (classifier that yields the highest margin)
- usually used for a large amount of classes
  - amounts to highest confidence that the test obs belongs to the  $k$ th class rather than any other classes

## Which Classifier to Use?

- when classes are (nearly) separable, SVM does better than Logistic, and so does LDA
  - if not, then LR (with ridge/lasso penalty) and SVM are very similar
- if you want to estimate probabilities, LR should be the choice
- for nonlinear boundaries, SVMs are popular
  - however, using kernels with LR and LDA can work as well, but computations more expensive
- drawback with SVMs is that it doesn't really select features well like the lasso, all features are used