

Algorithms for Data Science: Support Vector Machines

Héctor Corrada Bravo

University of Maryland, College Park, USA CMSC644: 2020-05-10



State-of-the-art classification and regression method

Flexible and efficient framework to learn classifers.

Nice geometric interpretation of how they are trained (based maximum margin arguments).

Can be estimated over *similarities* between observations (more on this later) rather than standard data in tabular form.

E.g., applications where string similarities, or network similarities are readily available.





SVMs follow the "predictor space partition" framework

Training data: $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

- \mathbf{x}_i is a vector of p predictor values for ith observation,
- y_i is the class label (we're going to use +1 and -1)

Training data: $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

- \mathbf{x}_i is a vector of p predictor values for ith observation,
- y_i is the class label (we're going to use +1 and -1)

Build a classifier by defining a *discriminative* function such that

$$w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip} > 0 ext{ if } y_i = 1$$

and

$$w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip} < 0 \, ext{if} y_i = -1$$

Points where the *discriminative* function equals 0 form a *hyper-plane* (i.e., a line in 2D)

$$\{x\,:\,w_0+w_1x_1+\dots+w_px_p=0\}$$



Hyper-plane partitions the predictor space into two sets on each side of the line.

Denote **w** is the vector (w_1, w_2, \ldots, w_p)

Hyper-plane partitions the predictor space into two sets on each side of the line.

Denote **w** is the vector (w_1, w_2, \ldots, w_p)

Restrict estimates to those for which $\mathbf{w'w} = \|\mathbf{w}\|^2 = 1$

Hyper-plane partitions the predictor space into two sets on each side of the line.

Denote **w** is the vector (w_1, w_2, \ldots, w_p)

Restrict estimates to those for which $\mathbf{w'w} = \|\mathbf{w}\|^2 = 1$

Then, the *signed* distance of any point x to the decision boundary L is $w_0 + \mathbf{w}' \mathbf{x}$.

With this we can easily describe the two partitions as

$$L^+=\{\mathbf{x}:w_0+\mathbf{w}'\mathbf{x}>0\},$$

$$L^-=\{\mathbf{x}: w_0+\mathbf{w}'\mathbf{x}<0\}$$

The \mathbf{w} we want as an estimate is one that separates the training data as perfectly as possible.

The \mathbf{w} we want as an estimate is one that separates the training data as perfectly as possible.

Describe this requirement as

$$y_i(w_0+\mathbf{w}'\mathbf{x}_i)>0, i=1,\ldots,N$$

Algorithm to find vector \mathbf{w} that satisfies the separation requirement as much as possible.

Penalize \mathbf{w} by how far into the wrong side misclassified points are:

$$D(w_0,\mathbf{w}) = -\sum_{i\in\mathcal{M}} y_i(w_0+\mathbf{w}'\mathbf{x}_i)$$

 \mathcal{M} : set of points misclassified by \mathbf{w} (on the wrong side of the hyperplane).

Estimate \mathbf{w} by minimizing D.

Assuming ${\mathcal M}$ is fixed, the gradient of D is

$$abla_{\mathbf{w}} D(w_0, \mathbf{w}) = -\sum_{i \in \mathcal{M}} y_i \mathbf{x}_i$$

and

$$rac{\partial D(w_0,\mathbf{w})}{\partial w_0} = -\sum_{i\in\mathcal{M}}y_i$$

Rosenblatt's algorithm uses *stochastic gradient descent*.

- Initialize parameters w_0 and ${f w}$
- Cycle through training points i, if it is misclassified, update parameters as

$$\mathbf{w} \leftarrow \mathbf{w} +
ho y_i \mathbf{x}_i$$

and

$$w_0 \leftarrow w_0 + \rho y_i$$

- Stop when converged (or get tired of waiting)
- (\$rho\$ is a learning rate parameter)

15 / 73

Update Rule:

 $\mathbf{w} \leftarrow \mathbf{w} +
ho y_i x_i$

Learning rate parameter ρ is used to control how much we update \mathbf{w} in each step.

Update Rule:

 $\mathbf{w} \leftarrow \mathbf{w} +
ho y_i x_i$

Learning rate parameter ρ is used to control how much we update ${\bf w}$ in each step.

This is the gradient descent algorithm that forms the basis of neural networks and deep learning.

There are a few problems with this algorithm:

If there exists w_0 and w that separates the training points perfectly,

There are a few problems with this algorithm:

If there exists w_0 and w that separates the training points perfectly,

then there are an infinite number of w_0 and \mathbf{w}_s that also separate the data perfectly

Algorithm will converge in a finite number of steps if the training data is separable

Algorithm will converge in a finite number of steps if the training data is separable

However, the number of finite steps can be *very* large

When the training data is *not* separable, the algorithm will not converge.

Support Vector Machines (SVMs) are designed to directly address these problems.



A central concept in SVMs that we did not see in logistic regression is **the margin**: the distance between the separating plane and its nearest datapoints.

When the data are separable, SVMs will choose the single optimal \mathbf{w} that *maximizes* the distance between the decision boundary and the closest point in each class.

When the data are separable, SVMs will choose the single optimal \mathbf{w} that *maximizes* the distance between the decision boundary and the closest point in each class.

Why is this a good idea?

SVMs are designed from three *key insights*:

1. Look for the maximum margin hyper-plane

- 2. Only depend on pair-wise "similarities" of observations
- 3. Only depend on a subset of observations (support vectors)

Let's see these in turn.

Goal: find the hyper-plane that separates training data with largest margin.

This will tend to *generalize* better since new observations have room to fall within margin and still be classified correctly.

This can be cast as *optimization* problem:

 $egin{aligned} \max_{w_0,\mathbf{w}} M \ & ext{ s. t.} |\mathbf{w}|^2 = 1 \ & ext{ } y_i(w_0 + \mathbf{w}'\mathbf{x}_i) \geq M \, orall i \end{aligned}$

Rewrite optimization problem setting $M = 1/||\mathbf{w}||^2$ and using a little bit of algebra:

$$egin{aligned} \min_{w_0,\mathbf{w}}&rac{1}{2}|\mathbf{w}|^2\ ext{s. t.} y_i(w_0+\mathbf{w}'\mathbf{x}_i) \geq 1\,orall i \end{aligned}$$

$$egin{aligned} \min_{w_0,\mathbf{w}}&rac{1}{2}|\mathbf{w}|^2\ ext{s. t.} y_i(w_0+\mathbf{w}'\mathbf{x}_i)\geq 1\,orall i \end{aligned}$$

This is a *constrained* optimization problem

Minimize the norm of \mathbf{w} under the constraint that it classifies every observation correctly.

Recall the standard constrained optimization problem we encountered before

 $egin{array}{lll} \min_x & f_0(x) \ ext{s.t.} & f_i(x) \leq 0 \; i=1,\ldots,m \ & h_i(x)=0 \; i=1,\ldots,p \end{array}$

Recall the standard constrained optimization problem we encountered before

$$egin{array}{lll} \min_x & f_0(x) \ ext{s.t.} & f_i(x) \leq 0 \; i=1,\ldots,m \ & h_i(x)=0 \; i=1,\ldots,p \end{array}$$

And associated Lagrangian

$$L(x,\lambda,
u)=f_0(x)+\sum_{i=1}^m\lambda_if_i(x)+\sum_{i=1}^p
u_ih_i(x)$$

Let's define the *dual* equation

$$g(\lambda,
u) = \min_x L(x,\lambda,
u)$$

Let's define the *dual* equation

$$g(\lambda,
u) = \min_x L(x,\lambda,
u)$$

And the dual problem

 $egin{array}{lll} \max_{\lambda,
u} & g(\lambda,
u) \ {
m s.\,t.} & \lambda_i \geq 0 \; i=1,\dots,m \end{array}$
Constrained Optimization

Duality gap: let \tilde{x} be primal optimal, then

 $g(\lambda,
u)\leq f_0(ilde{x})$

Constrained Optimization

Duality gap: let \tilde{x} be primal optimal, then

 $g(\lambda,
u)\leq f_0(ilde{x})$

If $ilde{x}$ is *primal optimal* and $(ilde{\lambda}, ilde{
u})$ are *dual optimal* then

$$g(ilde{\lambda}, ilde{
u})=f_0(ilde{x})$$

Constrained Optimization

Let's restate the optimality conditions (Karush-Kuhn-Tucker)

$f_i(\tilde{x}) \leq 0$	(primal feasible)		
$h_i(ilde{x})=0$			
$ ilde{\lambda}_i \geq 0$	$({ m dual\ feasible})$		
$ ilde{\lambda_i} f_i(ilde{x}) = 0$	(complementarity)		
$ abla_x L(ilde{x}, ilde{\lambda}, ilde{ u})=0$	$({ m saddle\ point})$		

Maximum-margin hyper-planes

What does the Lagrangian look like for the maximum-margin hyper-plane problem?

$$L(w_0, \mathbf{w}, lpha) = rac{1}{2} |\mathbf{w}|^2 + \sum_i -lpha_i [y_i(w_0 + \mathbf{w}' \mathbf{x}_i) - 1)]$$

Maximum-margin hyper-planes

What does the Lagrangian look like for the maximum-margin hyper-plane problem?

$$L(w_0,\mathbf{w},lpha)=rac{1}{2}|\mathbf{w}|^2+\sum_i -lpha_i[y_i(w_0+\mathbf{w'x}_i)-1)]$$

You can then find *dual function* by solving $\min_x L(x,\lambda,
u)$

Maximum-margin hyper-planes

In the maximum-margin hyper-plane case, the equivalent constrained maximization problem (the *dual* problem) is:

$$egin{aligned} \max_lpha \sum_{i=1}^N lpha_i - rac{1}{2} \sum_{i=1}^N \sum_{k=1}^N lpha_i lpha_k y_i y_k \mathbf{x}_i' \mathbf{x}_k \ ext{ s. t.} lpha_i \geq 0 \, orall i \end{aligned}$$

Maximum margin hyper-planes

$$egin{aligned} \max_lpha \sum_{i=1}^N lpha_i - rac{1}{2} \sum_{i=1}^N \sum_{k=1}^N lpha_i lpha_k y_i y_k \mathbf{x}_i' \mathbf{x}_k \ ext{ s. t.} lpha_i \geq 0 \, orall i \end{aligned}$$

This quadratic optimization problem is usually easier to optimize than the original problem (notice there is only positivity constraints on α).

Maximum margin hyper-planes

$$egin{aligned} \max_lpha \sum_{i=1}^N lpha_i - rac{1}{2} \sum_{i=1}^N \sum_{k=1}^N lpha_i lpha_k y_i y_k \mathbf{x}_i' \mathbf{x}_k \ ext{ s. t.} lpha_i \geq 0 \, orall i \end{aligned}$$

This quadratic optimization problem is usually easier to optimize than the original problem (notice there is only positivity constraints on α).

We can still recover original parameters w_0 and \mathbf{w} from solution α .

Key insight no. 2: **SVMs only depend on pairwise "similarity" functions of observations**

$$egin{aligned} \max_lpha \sum_{i=1}^N lpha_i - rac{1}{2} \sum_{i=1}^N \sum_{k=1}^N lpha_i lpha_k y_i y_k \mathbf{x}_i' \mathbf{x}_k \ ext{ s. t.} lpha_i = 0 \, orall i \end{aligned}$$

Only inner products between observations are required as opposed to the observations themselves.

Also, we can write the *discriminant* function in equivalent form

$$f(x) = w_0 + \sum_{i=1}^n y_i lpha_i \mathbf{x}' \mathbf{x}_i$$

Also, we can write the *discriminant* function in equivalent form

$$f(x) = w_0 + \sum_{i=1}^n y_i lpha_i \mathbf{x}' \mathbf{x}_i$$

Geometrically, we can think of the inner product between observations as a "similarity" measure.

Also, we can write the *discriminant* function in equivalent form

$$f(x) = w_0 + \sum_{i=1}^n y_i lpha_i \mathbf{x}' \mathbf{x}_i$$

Geometrically, we can think of the inner product between observations as a "similarity" measure.

Therefore, we can fit these models with other measures that works as "similarities".

Key insight no. 3: **SVMs only depend on a subset of observations** (support vectors)

Optimial solutions ${\bf w}$ and α must satisfy the following condition:

$$lpha_i[y_i(w_0+\mathbf{w}'\mathbf{x}_i)-1]=0\,orall i.$$

$$lpha_i[y_i(w_0+\mathbf{w}'\mathbf{x}_i)-1]=0\,orall i.$$

Case 1: $\alpha_i > 0$, then the signed distance between observation x_i and the decision boundary is 1.

This means that observation \mathbf{x}_i is *on the margin*

$$lpha_i[y_i(w_0+\mathbf{w}'\mathbf{x}_i)-1]=0\,orall i.$$

Case 2: $y_i(w_0 + \mathbf{w}'\mathbf{x}_i) > 1$, then observation x_i is not on the margin and $\alpha_i = 0$.

To define the discriminant function in terms of α s we only need observations that are *on the margin*,

i.e., those for which $\alpha_i > 0$.

To define the discriminant function in terms of α s we only need observations that are *on the margin*,

i.e., those for which $\alpha_i > 0$.

These are called *support vectors*.

To define the discriminant function in terms of α s we only need observations that are *on the margin*,

i.e., those for which $\alpha_i > 0$.

These are called *support vectors*.

Also implies we only need Support Vectors to make predictions.

The method we have discussed so far runs into an important complication:

What if there is no separating hyperplane?.



The solution is to penalize observations on the **wrong side of the margin** by introducing *slack variables* to the optimization problem.

$$egin{aligned} \min_{w_0,\mathbf{x},\xi} C\sum_{i=1}^N \xi_i + rac{1}{2}\|\mathbf{w}\|^2 \ ext{s. t} \ y_i(w_0+\mathbf{w}'\mathbf{x}_i) \geq 1-\xi_i \,orall i \ \xi_i \geq 0 \,orall i \end{aligned}$$

$$egin{aligned} \min_{w_0,\mathbf{w},\xi} C\sum_{i=1}^N \xi_i + rac{1}{2}\|\mathbf{w}\|^2 \ ext{s. t} \ y_i(w_0+\mathbf{w}'\mathbf{x}_i) \geq 1-\xi_i \,orall i \ \xi_i \geq 0 \,orall i \end{aligned}$$

C is a parameter that tradeoffs the width of the margin vs. the penalty on observations on the *wrong* side of the margin.

$$egin{aligned} \min_{w_0,\mathbf{w},\xi} C\sum_{i=1}^N \xi_i + rac{1}{2}\|\mathbf{w}\|^2 \ ext{s. t} \ y_i(w_0+\mathbf{w}'\mathbf{x}_i) \geq 1-\xi_i \,orall i \ \xi_i \geq 0 \,orall i \end{aligned}$$

C is a parameter that tradeoffs the width of the margin vs. the penalty on observations on the *wrong* side of the margin.

This is a "data fit + model complexity" learning objective.

C is a hyperparameter to be selected by the user or via cross-validation model selection methods.



An elegant result is that this formulation doesn't change the dual problem we saw before very much:

$$egin{aligned} \max_lpha \ &\sum_{i=1}^N lpha_i - rac{1}{2} \sum_{i=1}^N \sum_{k=1}^N lpha_i lpha_k y_i y_k x_i' x_k \ & ext{ s. t. } 0 \leq lpha_i \leq C \, orall i \end{aligned}$$

Only need support vectors, where $\alpha_i > 0$ to define the discriminant function and make predictions.

Only need support vectors, where $\alpha_i > 0$ to define the discriminant function and make predictions.

The smaller the cost parameter C, the learned SVM will have fewer support vectors.

Only need support vectors, where $\alpha_i > 0$ to define the discriminant function and make predictions.

The smaller the cost parameter C, the learned SVM will have fewer support vectors.

Think of the number of support vectors as a rough measure of the *complexity* of the SVM obtained.



What to do when we need nonlinear partitions of predictor space to get a classifier?

We can define the SVM discriminant function in terms of inner products of observations.

We can generalize inner product using "kernel" functions that provide something like an inner product:

$$f(\mathbf{x}) = w_0 + \sum_{i=1}^n y_i lpha_i k(\mathbf{x}, \mathbf{x}_i)$$

But, what is k? Let's consider two examples.

• Polynomial kernel: $k(\mathbf{x},\mathbf{x}_i) = 1 + \langle \mathbf{x},\mathbf{x}_i
angle^d$



• RBF (radial) kernel: $k(\mathbf{x},\mathbf{x}_i) = \exp\{-\gamma\sum_{j=1}^p (x_j-x_{ij})^2\}$



The optimization problem is very similar

$$egin{aligned} \max_lpha \ &\sum_{i=1}^N lpha_i - rac{1}{2}\sum_{i=1}^N \sum_{k=1}^N lpha_i lpha_k y_i y_k k(\mathbf{x}_i,\mathbf{x}_k) \ & ext{ s. t. } 0 \leq lpha_i \leq C \, orall i \end{aligned}$$

Let's try fitting SVMs to a credit card default dataset. Let's start with a linear SVM (where k is the inner product).

Here we are fitting three different SVMs resulting from using three different values of cost parameter C.

cost	number_svs	train_error	test_error
1e-02	350	3.44	3.22
1e+00	352	3.44	3.22
1e+02	354	3.44	3.22

Let's try now a *non-linear* SVM by using a radial kernel.

Notice now that we have two parameters to provide to the fitting function: cost parameter C and parameter γ of the radial kernel function.

cost	gamma	number_svs	train_error	test_error
0.01	0.01	344	3.44	3.22
1.00	0.01	348	3.44	3.22
10.00	0.01	347	3.44	3.22
0.01	1.00	392	3.44	3.22
1.00	1.00	426	2.82	2.58
10.00	1.00	382	2.64	2.66
0.01	10.00	491	3.44	3.22
1.00	10.00	1226	2.56	2.96

71/73
Support Vector Machines

Different algorithms depending on data size

- Massive number of examples with few predictors, train with stochastic gradient descent on the primal problem
- Moderate number of examples, use quadratic optimization with kernel functions
- For quadratic version, can subset observations that could be support vectors

Support Vector Machines

State-of-the-art for many applications

RBF kernels usually work well, but tuning γ properly is very important

Very elegant formulation

Kernel trick gives a lot of flexibility