

Support Vector Machine

A SVM is a discriminative classifier formally defined by a separating hyperplane. Given labeled training data, the algorithm outputs an optimal hyperplane which categorizes new examples.

Basics of Machine Learning Series

Index

Optimization Objective

The support vector machine objective can be seen as a modification to the cost of logistic regression. Consider the sigmoid function, given as,

where

The cost function of logistic regression as in the post [Logistic Regression Model](#) is given by,

Each training instance contributes to the cost function the following term,

So when $y = 1$, the contributed term is $\frac{1}{2} \ln(1 + e^{\frac{1}{C}})$, which can be seen in the plot below. The cost function of SVM, denoted as $J(w)$, is a modification the former and a close approximation.

Similarly, when $y = 0$, the contributed term is $\frac{1}{2} \ln(1 + e^{-\frac{1}{C}})$, which can be seen in the plot below. The cost function of SVM, denoted as $J(w)$, is a modification the former and a close approximation.

Regularized version of $J(w)$ can from the post [Regularized Logistic Regression](#) can be rewritten as,

In order to come up with the cost function for the SVM, $J(w)$ is modified by replacing the corresponding cost terms, which gives,

Following the conventions of SVM the following modifications are made to the cost in $J(w)$, which effectively is a change in notation but not the underlying logic,

- removing $\frac{1}{2}$ does not affect the minimization logic at all as the minima of a function is not changed by the linear scaling.
- change the form of parameterization from $\frac{1}{C}$ to C where it can be intuitively thought that C is a regularization parameter.

After applying the above changes, $J(w)$ gives,

The SVM hypothesis does not predict probability, instead gives hard class labels.

Large Margin Intuition

According to $J(w)$ and the plots of the cost function as shown in the image above, the following are two desirable states for SVM,

- if C is large, then $J(w)$ (not just $J(w)$)
- if C is small, then $J(w)$ (not just $J(w)$)

Let C in $J(w)$ be a large value. Consequently, in order to minimize the cost, the corresponding term $\frac{1}{C}$ must be close to 0.

Hence, in order to minimize the cost function, when $y = 1$, $\frac{1}{C}$ should be 0, and similarly, when $y = 0$, $\frac{1}{C}$ should be 0. And thus, from the plots in Fig.3, it is clear that it can only be fulfilled by the two states listed above.

Following the above intuition, the cost function can be written as,

subject to constraints,

What this basically leads to is the selection of a decision boundary that tries to maximize the margin from the support vectors as shown in the plot below. This maximization of the margin as seen for decision boundary A increases the robustness over decision boundaries with lesser margins like B. And it is this property of the SVMs that attributes the name **large margin classifier** to it.

As discussed in the section above, the effect of C can be considered as reciprocal of regularization parameter λ . This is more clear from Fig.5. A single outlier, can make the model choose the decision boundary with smaller margin if the value of C is large. A small value of C ensures that the outliers are overlooked and best approximation of large margin boundary is determined.

Mathematical Background

Vector Inner Product: Consider two vectors, \vec{a} and \vec{b} , given by,

Then, the **inner product** or the **dot product** is defined as $\vec{a} \cdot \vec{b}$.

Norm of a vector, \vec{a} , denoted as $\|\vec{a}\|$ is the euclidean length of the vector given by the pythagoras theorem as,

The inner product can also be defined as,

where $\vec{a} \cdot \vec{b}$ can be described as the projection of vector \vec{a} onto vector \vec{b} which can be either positive or negative signed based on the angle θ between the vectors as shown in the image below.

SVM Decision Boundary: From $J(w)$, the optimization statement can be written as,

subject to constraints,

Let n and m , i.e. number of features is 2 for simplicity, then \vec{a} can be written as,

Using \vec{a} , \vec{b} can be written as,

The plot of $J(w)$ can be seen below,

Hence, using \vec{a} and \vec{b} , the optimization objective in $J(w)$ and the constraints in $J(w)$ are written as,

subject to constraints,

where $\vec{a} \cdot \vec{b}$ is the projection of \vec{a} onto vector \vec{b} .

Consider two decision boundaries, A and B, and their respective perpendicular parameters, \vec{a} and \vec{b} , as shown in the plot below. As a consequence of choosing \vec{a} for simplification, all the corresponding decision boundaries pass through the origin.

Based on the two training examples of either class chosen, close to the boundaries, it can be seen that the magnitude of projection is more in case of \vec{a} than \vec{b} . This basically tells that it would be possible to choose smaller values of $\vec{a} \cdot \vec{b}$ and satisfy $\vec{a} \cdot \vec{b} \geq 1$ if the value of projection is bigger and hence, the decision boundary, B is more favourable to the optimization objective.

Why is decision boundary perpendicular to the ?

Consider two points \vec{a} and \vec{b} on the decision boundary given by,

Since the two points are on the line, they must satisfy $\vec{a} \cdot \vec{b} = 1$. Substitution leads to the following,

Subtracting $\vec{a} \cdot \vec{b}$ from $\vec{a} \cdot \vec{b}$,

Since \vec{a} and \vec{b} lie on the line, the vector \vec{a} is on the line too. Following the property of orthogonal vectors, is possible only if \vec{a} is orthogonal (perpendicular) to \vec{b} , and hence perpendicular to the decision boundary.

Kernels

When dealing with non-linear decision boundaries, a learning method like logistic regression relies on high order polynomial features to find a complex decision boundary and fit the dataset, i.e. predict y if,

where $\vec{a} \cdot \vec{b}$.

A natural question that arises is if there are choices of different features than in 7? A SVM does this by picking points in the space called **landmarks** and defining functions called **similarity** corresponding to the landmarks.

Say, there are three landmarks defined, \vec{a}_1, \vec{a}_2 and \vec{a}_3 as shown in the plot above, the for any given \vec{x} , \vec{a}_1, \vec{a}_2 and \vec{a}_3 are defined as follows,

Here, the similarity function is mathematically termed a **kernel**. The specific kernel used in $J(w)$ is called the γ kernel. γ kernels are sometimes also denoted as $\frac{1}{2C}$.

Consider \vec{a} from $J(w)$. If there exists \vec{a}_1 close to landmark \vec{a}_1 , then $\vec{a} \cdot \vec{a}_1$ and hence $\frac{1}{2C}$ will be a far from the landmark, \vec{a}_1 will be a larger value and hence exponential fall will cause $\frac{1}{2C}$ to be small. So effectively the choice of landmarks has helped in increasing the number of features \vec{a} had from 2 to 3, which can be helpful in discrimination.

For a gaussian kernel, the value of γ defines the spread of the normal distribution. If γ is small, the spread will be narrower and when its large the spread will be wider.

Also, the intuition is clear about how landmarks help in generating the new features. Along with the values of parameter, \vec{a} and \vec{b} , various different decision boundaries can be achieved.

How to choose optimal landmarks?

In a complex machine learning problem it would be advantageous to choose a lot more landmarks. This is generally achieved by choosing landmarks at the point of the training examples, i.e. landmarks equal to the number of training examples are chosen, ending up in n if there are n training examples. This translates to the fact that each feature is a measure of how close is an instance to the existing points of the class, leading to generation of new feature vectors.

For SVM training, given training examples $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$, features are computed, and \vec{a} is defined as follows

The training objective from $J(w)$ is modified as follows,

In this case, \vec{a} in $J(w)$ by the virtue of procedure used to choose \vec{a} .

The regularization term in $J(w)$ can be written as $\frac{1}{2C}$. But in practice most SVM libraries, instead, which can be considered a scaled version is used as it gives certain optimization benefits and scaling to bigger training sets, which will be taken up at a later point in maybe another post.

While the kernels idea can be applied to other algorithms like logistic regression, the computational tricks that apply to SVMs do not generalize as well to other algorithms.

Hence, SVMs and Kernels tend to go particularly well together.

Bias/Variance

Since \vec{a} ,

- Large C : Low bias, High Variance
- Small C : High bias, Low Variance

Regarding γ ,

- Large γ : High Bias, Low Variance (Features vary more smoothly)
- Small γ : Low Bias, High Variance (Features vary less smoothly)

Choice of Kernels

- Linear Kernel:** is equivalent to a no kernel setting giving a standard linear classifier given by,

Linear kernels are used when the number of training data is less but the number of features in the training data is huge.

- Gaussian Kernel:** Make a choice of γ to adjust the bias/variance trade-off.

Gaussian kernels are generally used when the number of training data is huge and the number of features are small.

Feature scaling is important when using SVM, especially Gaussian Kernels, because if the ranges vary a lot then the similarity feature would be dominated by features with higher range of values.

All the kernels used for SVM, must satisfy Mercer's Theorem, to make sure that SVM optimizations do not diverge.

Some other kernels known to be used with SVMs are:

- Polynomial kernels,
- Esoteric kernels, like string kernel, chi-square kernel, histogram intersection kernel, ..

Multi-Class Classification

- Most SVM libraries have multi-class classification.
- Alternatively, one may use one-vs-all technique to train different SVMs and pick class with largest

Logistic Regression vs SVM

- if C is large relative to γ , use logistic regression or SVM with linear kernel, like if $C \gg \gamma$
- if C is small and γ is intermediate, use SVM with gaussian kernel, like if $C \approx \gamma$
- if C is small and γ is large, create/add more features, then use logistic regression or SVM with no kernel, as with huge datasets SVMs struggle with gaussian kernels, like if $C \approx \gamma \gg 1$

Logistic Regression and SVM without a kernel (with linear kernel) generally give very similar. A neural network would work well on these training data too, but would be slower to train.

Also, the optimization problem of SVM is a convex problem, so the issue of getting stuck in local minima is non-existent for SVMs.

REFERENCES:

Machine Learning: Coursera - Optimization Objective
Machine Learning: Coursera - Large Margin Intuition
Machine Learning: Coursera - Mathematics of Large Margin Classification
Machine Learning: Coursera - Kernel I
Machine Learning: Coursera - Kernel II
Machine Learning: Coursera - Using An SVM
Quora - Why is theta perpendicular to the decision boundary?
Introduction to support vector machines