

## Bootstrapping the Support Vector Machine

Thomas Goerttler<sup>1,2,3</sup>, Christian Koopmann<sup>1,2,3</sup>, Patricia Craja<sup>1,2,3</sup>

<sup>1</sup>Humboldt University of Berlin, Unter den Linden 6, 10099 Berlin

<sup>2</sup>Free University Berlin, Kaiserswerther Str. 16-18, 14195 Berlin

<sup>3</sup>Technische Universitt Berlin, 10623 Berlin

thomas.goerttler@gmail.com

c.k.e.koopmann@gmail.com

Patricia.craja@gmx.de

## Abstract

Support Vector Machines are one of the most successful methods of Machine Learning. After training they can predict for each point a possible class. Unfortunately there is not uncertainty captured, that the point is really in this class. We found out a way to calculate variance for a support vector machine. We use bootstrapped samples and take this into account. The variance can be calculated parallelized.

## Introduction

Support Vector Machines are one of the most successful methods of Machine Learning. By reducing non-linear complex decisions problems to linear problems through application of the kernel-Trick, they represent a computationally efficient way to tackle these problems.

svm uses hyperplane space feature to separate data. The dimension of feature space is controlled by choice of kernel function. Common kernels are the linear kernel, the gaussian kernel (rbf) and the polynomial kernel.

SVMs based on certain kernels (e.g. Gaussian RBF Kernel) are non parametric methods. Since the distribution of the underlying data is generally unknown so is the finite sample distribution of these methods. There has been considerable research on the asymptotic distribution of SVMs, which have been shown to be asymptotically normally distributed under certain conditions. An alternative idea to estimate these distributions is using Efrons empirical bootstrap. The idea behind this method is to repeatedly draw samples with replacement from the full data according to the empirical distribution function of the data. Through the repeated calculation of the statistic of interest one can get an estimate of its distribution. For the SVM this estimate has been shown to be consistent under relatively mild conditions. ..

## Background

What is the status in the science world right now. Important has been **bootstrap, distances and parallelization**

## Problem

There is no way to calculate uncertainty (variance) of the prediction as we have no assumptions about the distribution. We found out that we can sample variance by bootstrapping. Furthermore we found out that they are tuning parameter, data-distribution option and svm attributes that influence this variance.

Predictions are only binary variables and therefore might be identical across all bootstrap samples. Therefore we use minimal distance of predictionpoint to decision boundary as real valued substitute.

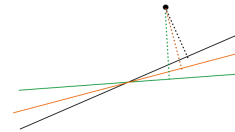
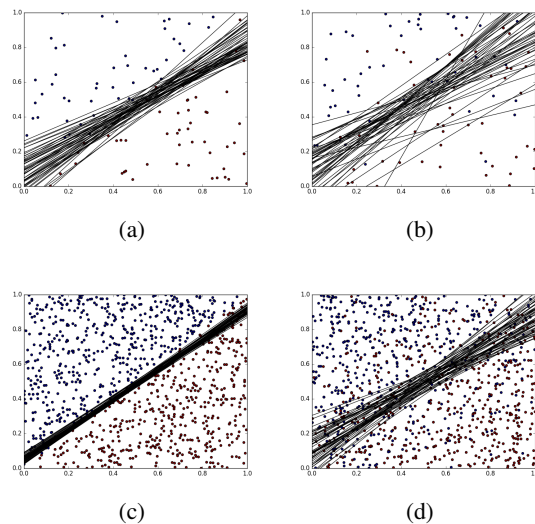


Figure 1: Really awesome trick of distance

## Implementation

[illegible]

[illegible]

```
def bootstrap_svm(trainings_data,
prediction_data, kernel, C, gamma = "
auto", degree = 3, processes = 10,
replications = 1000):

    input_parameters = SVM_Input(
        trainings_data, prediction_data,
        kernel, C, gamma, degree)
    data = input_parameters
    real_svm = do_svm(data)

    ### Do bootstrapping
    PROCESSES = processes
    REPLICATIONS = replications
    pool = Pool(processes = PROCESSES)
    results = pool.map(single_sample_and_svm
        , [data] * REPLICATIONS)

    ### Calculate the Variance of the
    Support Vector Machine

    points_information = Points_Information(
        results)

    variance_of_svm_probabilites =
        calculate_variance_of_svm(
            points_information.probabilites)
    variance_of_svm_distance_to_hyperplane =
        calculate_variance_of_svm(
            points_information.distances)

    result = Bootstrap_Result([real_svm[0],
        real_svm[1], real_svm[2],
        variance_of_svm_probabilites,
        variance_of_svm_distance_to_hyperplane
        ], real_svm[0].n_support_ )

    return(result)
```

erful! Kernels are powerful! Kernels are powerful! Kernels  
are powerful! Kernels are powerful! Kernels are powerful!  
Kernels are powerful! Kernels are powerful! Kernels are  
powerful! Kernels are powerful! Kernels are powerful! Ker-  
nels are powerful! Kernels are powerful! Kernels are pow-  
erful! Kernels are powerful! Kernels are powerful! Kernels  
are powerful! Kernels are powerful! Kernels are powerful!  
Kernels are powerful!

There we three indicators we found out. The C Parameter, Balance and number of support vector. Once we used linear, once we used gaussian.

[illegible]

```

1 def get_data_normally_distributed(coefs,
2     errorCoef, intercept, size):
3     inputs = []
4     error = errorCoef*rd.standard_normal(size)
5     y = error + intercept
6     for i in range(len(coefs)):
7         inputs = inputs + [rd.standard_normal(
8             size)]
9         y = y+coefs[i]*inputs[i]
10    y = sign(y)
11    inputs = list(zip(*inputs))
12    return([y,inputs])

```

Figure 4: Implementation of Datasimulation normally distributed

```

1 def get_data_centroid_distributed(coefs,
2   locations, errorCoef, size, intercept,
3   distance, xdistribution = "normal", parl
4   = 0, par2 = 1):
5
6     X = []
7     dimension = len(list(zip(*locations)))
8     for i in range(dimension):
9         if(xdistribution == "normal"):
10             X = X + [rd.normal(parl, par2, size)]
11         elif(xdistribution == "uniform"):
12             X = X + [rd.uniform(parl, par2, size)]
13         else:
14             print("Please_choose_supported_
15                 Distribution")
16             return None
17     X = list(zip(*X))
18     distances = []
19     error = errorCoef*rd.standard_normal(size)
20     y = error + intercept
21     for i in range(size):
22         newDistance = pdist([X[i]]+locations,
23                             distance)
24         newDistance = newDistance[:len(locations)
25                                   )]
26         inverseDistance = power(newDistance, -1)
27         y[i] = y[i] + dot(coefs, inverseDistance)
28     distances = distances + [newDistance]
29     y = sign(y)
30     #distances = list(zip(*distances))
31     return([y,X])

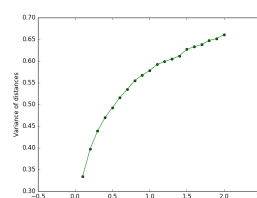
```

Figure 5: Implementation of Datasimulation normally distributed

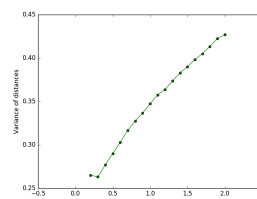
are powerful! Kernels are powerful! Kernels are powerful!  
Kernels are powerful! Kernels are powerful! Kernels are  
powerful! Kernels are powerful! Kernels are powerful! Ker-  
nels are powerful! Kernels are powerful! Kernels are pow-  
erful! Kernels are powerful! Kernels are powerful! Kernels  
are powerful! Kernels are powerful! Kernels are powerful!  
Kernels are powerful! Kernels are powerful! Kernels are

[illegible]

### Influence of the tuning parameter C

[illegible]

(a)



(b)

Figure 6: The results for changing Balances linear Datasize  
100 Rep| 1000 Datasets 20

### Influence of the balance of the training dataset

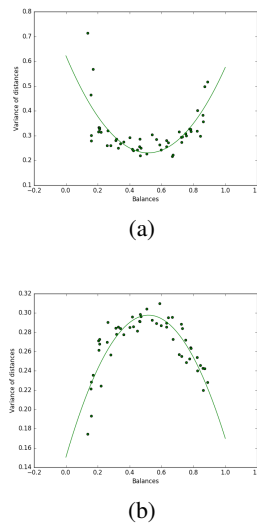
[illegible]

Figure 7: The results for changing Balances linear Datasize  
500 Repl 1000

### Influence of the svm attributes

[illegible]

## Conclusion

Nice to work with. I also works whit real datasets

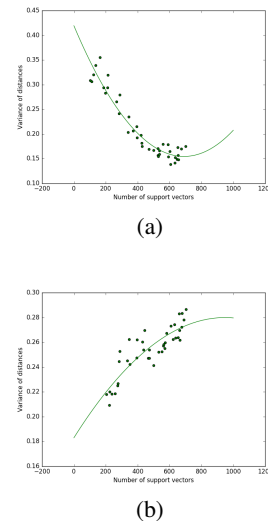


Figure 8: The results for changing number of support vectors Datasize 1000 Repl 500

### Influence of different aspects on variance

C-Parameter: Positive but decreasing influence on variance. Balance of Data: Positive Quadratic influence for Linear and negative quadratic influence for Gaussian SVM. Both with extremum around 0.5 (perfect balance). Number of support vectors: Positive influence for Gaussian and negative influence for linear SVMs

## Outlook

Analyse influence of dimensionality on variances. Analyse Data simulated from more "exotic" distributions. Analyse influence of other tuning parameters e.g. "Gamma" of rbf-kernel, degree of polynomial kernel etc.

## References