# Selective Bayesian Classifier: Feature Selection for the Naïve Bayesian Classifier Using Decision Trees

Chotirat "Ann" Ratanamahatana, Dimitrios Gunopulos
*Department of Computer Science,*
*University of California, Riverside, USA.*

## Abstract

It is known that Naïve Bayesian classifier (NB) works very well on some domains, and poorly on some. The performance of NB suffers in domains that involve correlated features. C4.5 decision trees, on the other hand, typically perform better than the Naïve Bayesian algorithm on such domains. This paper describes a Selective Bayesian classifier (SBC) that simply uses only those features that C4.5 would use in its decision tree when learning a small example of a training set, a combination of the two different natures of classifiers. Experiments conducted on eleven datasets indicate that SBC performs reliably better than NB on all domains, and SBC outperforms C4.5 on many datasets of which C4.5 outperform NB. SBC also can eliminate, on most cases, more than half of the original attributes, which can greatly reduce the size of the training and test data, as well as the running time. Further, the SBC algorithm typically learns faster than both C4.5 and NB, needing fewer training examples to reach high accuracy of classification.

## 1 Introduction

Two of the most widely used and successful methods of classification are C4.5 decision trees [9] and Naïve Bayesian learning (NB) [2]. While C4.5 constructs decision trees by using features to try and split the training set in to positive and negative examples until it achieves high accuracy on the training set, NB represents each class with a probabilistic summary, and finds the most likely class for each example it is asked to classify.

Severalresearchershaveemphasizedon theissueofredundantattributes,as wellasadvantagesoffeatureselectionfortheNaïveBayesianClassifier,notonly forin ductionlearning.Pazzani[8 ]exploresthemethodsofjoiningtwo(or more)relatedattributesintoanewcompoundattribute wheretheattribute dependenciesarepresent.Anothermethod,BoostingonNaïveBayesian classifier[3]hasbeenexperimentedbyapplyingseriesofclassifierstothe problemandpayingmoreattentiontotheexamplesmisclassifiedbyits predecessor.Ho wever,itwasshownthatitfailsonaverageinasetofnatural domain[7].AugmentedBayesianClassifiers[5]isanotherapproachwhere NaïveBayesisaugmentedbytheadditionofcorrelationarcsbetweenattributes. LangleyandSage[6],ontheother hand,useawrapperapproachforthesubset selectiontoonlyselectrelevantfeaturesforNB.

IthasbeenshownthatNaïveBayesianclassifierisextremelyeffectivein practiceanddifficulttosystematicallyimproveupon[1].Inthispaper,weshow thatitispossibletoreliablyimprovethisclassifierbyusingafeatureselection method.NaïveBayescansufferfromoversensitivitytoredundantand/or irrelevantattributes.Iftwoormoreattributesarehighlycorrelated,theyreceive toomuchweig htinthefinaldecisionastowhichclassanexamplebelongsto. Thisleadstoadeclineinaccuracyofpredictionindomainswithcorrelated features.C4.5doesnotsufferfromthisproblembecauseiftwoattributesare correlated,itwillnotbepossi bletousebothofthemtosplitthetrainingset, sincethiswouldleadtoexactlythesamesplit,whichmakesnodifferencetothe existingtree.ThisisoneofthemainreasonsC4.5performsbetterthanNBon domainswithcorrelatedattributes.

Weco njecturethattheperformanceofNBimprovesifitusesonlythose featuresthatC4.5usedinconstructingitsdecisiontree.Thismethodoffeature selectionwouldalsoperformwellandlearnquickly,thatis,itwouldneedfewer trainingexamplestorea chhighclassificationaccuracy.

Wepresentexperimentalevidencethatthismethodoffeatureselectionleads toimprovedperformanceoftheNaïveBayesianClassifier,especiallyinthe domainswhereNaïveBayesperformsnotaswellasC4.5.Weanalyzet he behavioronten domainsfromtheUCIrepository. Theexperimentalresults justifyourexpec tation.WealsotestedSBConanothersufficientlylarge syntheticdatasetandouralgorithmappearedtoscalenicely.OurSelective BayesianClassifieralwaysoutperformsNBandperformsaswellas,orbetter thanC4.5onalmostallthedomains.

## 2Naï veBayesianClassifier

### 2.1DescriptionandProblems

TheNaïveBayesianclassifierisastraightforwardandfrequentlyusedmethod forsupervisedlearning.Itprovidesaflexiblewayfordealingwithanynumber ofattributesorclasses,andisbasedon probabilitytheory (Bayes'rule) .Itisthe asymptoticallyfastestlearningalgorithmthatexaminesallitstraininginput.It hasbeendemonstratedtoperformsurprisinglywellinaverywidevarietyof

problemsinspiteofthesimplisticnatureofthe          model.Furthermore,small amountsofbaddata,or"noise,"donotperturbtheresultsbymuch.

However,therearetwocentralassumptionsinNaïveBayesianclassification. First,theclassificationassumesthattheelementsofeachclasscanbeassigned onprobabilitymeasurement,andthatthemeasurementissufficienttoclassifythe elementintoexactlyoneclass.Thisassumptionentailsthattheclassescanbe differentiatedonlybymeansoftheattributevalues.Thedependenceonthistype ofdiffe rentiationisrelatedtotheideaoflinearseparability;therefore,Naïve BayesianclassificationmaynoteasilylearnorpredictcomplicatedBoolean relations.

Theotherassumptionisthatgivenaparticularclassmembership,the probabilitiesofpartic ularattributeshavingparticularvaluesareindependentof eachother.However,thisassumptionisoftenviolatedinreality.

Aplausibleassumptionofindependenceiscomputationallyproblematic. Thisisbestdescribedbyredundantattributes.Ifw eposittwoindependent features,andathird   whichisredundant(i.e.perfectlycorrelated)with      thefirst , the firstattributewillhave    twiceasmuchinfluenceontheexpressionasthe secondhas,whichisastrengthnotreflectedinreality.Theincre      asedstrengthof thefirstattributeincreasesthepossibilityofunwantedbiasintheclassification. Evenwiththisindependenceassumption,HandandYuillustratedthatNaïve Bayesianclassificationstillworkswellinpractice[4].However,thispape         r showsthatifthoseredundantattributesareeliminated,theperformanceofNaïve Bayesianclassifiercansignificantlyincrease.

## 3C4.5DecisionTrees

Decisiontreesareoneofthemostpopularmethodsusedforinductiveinference. Theyarerobustf ornoisydataandcapableoflearningdisjunctiveexpressions.A decisiontreeisak -arytreewhereeachoftheinternalnodesspecifiesateston someattributesfromtheinputfeaturesetusedtorepresentthedata.Eachbranch descendingfromanode   correspondstooneofthepossiblevaluesofthefeature specifiedatthatnode.Andeachtestresultsinbranches,whichrepresent differentoutcomesofthetest.

Thealgorithmstartswiththeentiresetoftuplesinthetrainingset,selectsthe *best*attributethatyieldsmaximuminformationforclassification,andgeneratesa testnodeforthisattribute.Then,topdowninductionofdecisiontreesdivides thecurrentsetoftuplesaccordingtotheirvaluesofthecurrenttestattribute. Classifier generationstops,ifalltuplesinasubsetbelongtothesameclass,orif itisnotworthtoproceedwithanadditionalseparationintofurthersubsets,i.e.if furtherattributetestsyieldonlyinformationforclassificationbelowapre              - specifiedthres hold.

Decisiontreealgorithmusesanentropy              -basedmeasureknownas "informationgain"asaheuristicforselectingtheattributethatwillbestsplitthe trainingdataintoseparateclasses.Itsalgorithmcomputestheinformationgain ofeachattribu te,andineachround,theonewiththehighestinformationgain willbechosenasthetestattributeforthegivensetoftrainingdata.Awell              - chosensplitpointshouldhelpinsplittingthedatatothebestpossibleextent.

After all, a main criterion in the greedy decision tree approach is to build *shorter* trees. The best split point can be easily evaluated by considering each unique value for that feature in the given data as a possible split point and calculating the associated information gain.

A simple decision tree algorithm only selects one decision tree given an example set, though there may be many different trees consistent with the data. The information gain measure (implemented in ID3 decision trees) is biased in that it tends to prefer attributes with many values rather than those with few values. C4.5 suppresses this bias by using an alternative measure called *Information Gain Ratio*, which considers the probability of each attribute value. This removes the bias of information gain towards features with many values.

### 3.1 Tree Pruning

C4.5 builds a tree so that most of the training examples are classified correctly. Though this approach is correct when there is no noise, accuracy for unseen data might degrade in cases where there is a lot of noise associated with the training examples and/or the number of training examples is very small. To alleviate this problem, C4.5 uses the post-pruning method. This approach allows C4.5 to grow a complete decision tree first, and then *post-prune* the tree. It tries to shorten the tree in order to overcome overfitting. This generally involves removal of some of the nodes or subtrees from the original decision tree. Its goal is to improve (by pruning) the accuracy on the unseen set of examples. As a result, C4.5 achieves further elimination of features through pruning. It uses rule-postpruning to remove some of the *insignificant* nodes (and hence, some *not so relevant* features) from the tree.

## 4 Selective Bayesian Classifier

Our purpose is to improve the performance of the Naïve Bayesian classifier by removing redundant and/or irrelevant attributes from the dataset, and only choosing those that are most informative in classification task, according to the decision tree constructed by C4.5.

### 4.1 Description

As described in section 3, the features that C4.5 selected in constructing its decision tree are likely to be the ones that are most descriptive in terms of the classifier, in spite of the fact that a tree structure inherently incorporates dependencies among attributes, while Naïve Bayes works on a conditional independence assumption. C4.5 will naturally construct a tree that does not have an overly complicated branching structure if it does not have too many examples that need to be learned. As the number of training examples increases, the attributes that are considered will usually be the ones that are not correlated. This is mainly because C4.5 will use only one of a set of correlated features for making good splits in training set. However, sometimes many of the branches

mayreflectnoiseoroutliers(overfitting)inthetrainingdata."Treepruning" procedureinC4.5attemptstoidentifyandremovethoseleastreliablebranches, withthegoalofimp    rovingclassificationaccuracyonunseendata.Evenafter pruning,iftheresultdecisiontreeisstilltoo    *deep*orgrownintotoomanylevels, ouralgorithmonlypicksattributescontainedinthefirstfewlevelsofthetreeas themostrepresentativeat    tributes.Thisissupportedbythefactthatbythe selectionofattributesthatsplitthedatainthebestpossiblewayateverynode, C4.5willtrytoensurethatitencountersaleafattheveryearliestpossiblepoint, i.e.itpreferstoconstructsho    rtertrees.Andbyitsalgorithm,C4.5willfindtrees thathaveattributeswithhigherinformationgainnearertotheroot.We conjecturethatthissimplemethodoffeatureselectionwouldhelpNaïve Bayesianclassifierperformwellandlearnquickly,t    hatis,itwouldneedfewer trainingexamplestoreachhighclassificationaccuracy.

**4.2Algorithm**

+---------------------------------------------------------------+
| 1.  Shuffletheoriginaldata.                                   |
| 2.  Take10%fromtheori    ginaldataastrainingdata.             |
| 3.  RunC4.5ondatafromstep2.                                   |
| 4.  Selectasetofattributesthatappearonlyinthefirst           |
|     3levelsofthesimplifieddecisiontreeas    *relevant*       |
|     features.                                                 |
| 5.  Repeat10times(step1    -4)                                |
| 6.  Unionthesetsofattributesobtainedfromall10                |
|     rounds.                                                    |
| 7.  RunNaïveBayesianclassifieronthetrainingand               |
|     testdatausing    *only*thefinalfeaturesselectedinstep    |
+---------------------------------------------------------------+

Figure 1.Select    iveBayesianClassifierAlgorithm:
FeatureSelectionUsingC4.5

Figure1showsthealgorithmfortheSelectiveBayesianclassifier.Wefirst shufflethetrainingdataanduse10%ofthattorunC4.5on.Thisist            omakesure thatallthesubsamplesarenotbiasedtowardanyparticularclasses.Wefind 10%ofthetrainingtobeagoodsizeforfeatureselectionprocess          . Oncewerun C4.5andobtainthedecisiontree,weonlypickattributesthatonlyappearinthe first    3levelsofthedecisiontreeasthemostrelevantfeatures.Wehypothesize thatifafeatureinthedeeperlevelsonanyoneexecutionofC4.5is                *relevant* enough,itwillfinallyrisesupandappearinoneofthetoplevelsofthetreein someotherexe    cutionsofC4.5.Itisimportanttonotethatinthe10different iterations,C4.5maygiveslightlydifferentdecisiontrees,i.e.itusesdifferent attributestoproducedecisiontreefordifferenttrainingsets,evenwhenthe numberoftrainingexample    sisthesameacrossthesetrainingsets.Weunionall theattributesfromeachrun,andfinally,runtheNaïveBayesianclassifieronthe trainingandtestdatausingonlythosefeaturesselectedinthepreviousstep.

# 5 Experimental Evaluation

## 5.1 Th eDatasets

Weused10datas etsfromtheUCIrepository andonesyntheticdataset,shownin Table1.TheSyntheticData,createdwithGaussiandistribution,contains 1,200,000instanceswith20attributesand2classes.Wechose10datasetsfrom theUCIdatabases,5of whichNaïveBayes outperformsC4.5andtheother5of whichC4.5outperformsNaïveBayes .

Table 1.Descriptionsofdomainsused

| Dataset | #Attributes | #Classes | #Instances |
|---|---|---|---|
| Ecoli | 8 | 8 | 336 |
| GermanCredit | 20 | 2 | 1,000 |
| KrVsKp | 37 | 2 | 3,198 |
| Monk | 6 | 2 | 554 |
| Mushroom | 22 | 2 | 8,124 |
| Pima | 8 | 2 | 768 |
| Promoter | 57 | 2 | 106 |
| Soybean | 35 | 19 | 307 |
| Wisconsin | 9 | 2 | 699 |
| Vote | 16 | 2 | 435 |
| SyntheticData | 20 | 2 | 1,200,000 |

## 5.2 ExperimentalDesign

1. Eachdatasetisshuffledrando mly.
2. Produce *disjoint* trainingandtest setsasfollows.
   10%trainingand90%testdata
   20%trainingand80%testdata
   … … … …
   90%trainingand10%testdata
   99%trainingand1%testdata

3. Foreachsetoftrainingandtest data,run
   - NaïveBayesianClassifie r (NBC)
   - C4.5,and
   - SelectiveBayesian Classifier(SBC)
4. Repeat15times

Theclassifieraccuracyisdeterminedby *RandomSubsampling* method. The overallaccuracyestimat eisthemeanoftheaccuraciesobtainedfromall iterations.Thiswillgiveusinformationaboutboththelearningrates,aswellas theasymptoticaccuracyofthelearningalgorithmsused.

## 5.3 ExperimentalResults

Theresultsconfirmtheinitialhypotheses. ItisclearthatSBC doesimprove NBC'sperformance inalldomains ,andit does learnfasterthanbothC4.5a nd NBConallthedataset,i.e.withsmallnumberoftrainingdata(10%),the predictionaccuracyforSBCishigher.
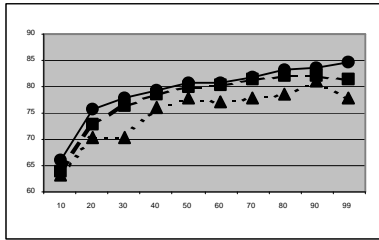
Figure2     –11depictthelearningc     urvesforthe10UCIdatasets.



Figure 2: Ecoli.336instances,8     attrib,
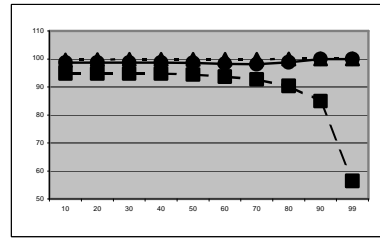8classes, 4 SBCattrib  .



Figure 6: Mushroom.8,124instances,   22
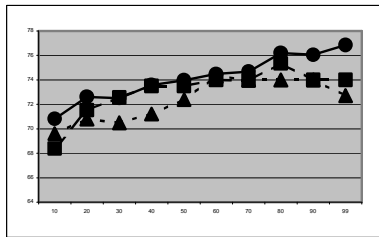attrib,2classes,   6 SBCattrib .



Figure 3:Germ  an.1,000instances,20
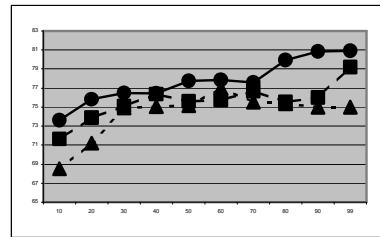attrib,2classes,   6 SBCattrib .



Figure 7: Pima.768instances,8     attrib,
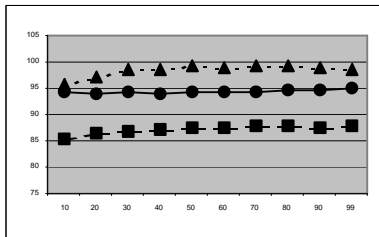2classes,  5 SBCattrib .



Figure 4:KrV  sKp.3,198instances,37
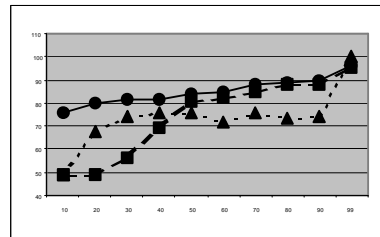attrib,2classes,   4 SBCattrib .



Figure 8:Pro  moter.106instances,57
attrib,2classes,   5 SBCattrib .



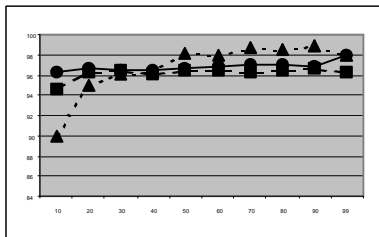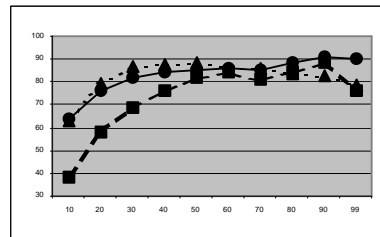Figure 5:Monk.554instances,6
attrib,2classes,  4 SBCattrib  .



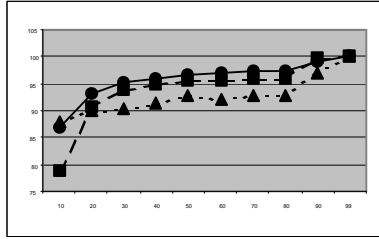Figure 9:Soybean.307in     stances,35
attrib,19classes,   12 SBCattrib .

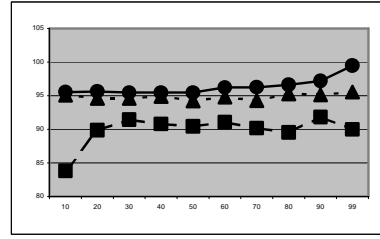Figure 10:Wisconsin.699   instances,9
attrib,2classes,   4 SBCattrib .

Figure 11:Vote.435instances,16   attrib,
2classes,  3 SBCattrib .

**TheX  -axisshowsthetrainingdata(%),and          theY  -axisshowsthe
accuracyontestdata.   SBCisrepresentedby   ● withasolidline  . NBCis
representedby  ■  withabigdashline  .  And C4.5isrepresentedby  ▲  with a
smalldashline  .**

NotethatalltheC4.5accuracyconsideredinthisexperimentisbasedonthe
simplifieddecisiontree(withpruning).Thisaccu          racyisusuallyhigheronthe
unseendata,comparingwiththeaccuracybasedonunpruneddecisiontrees.
     ToseeaclearerpictureontheSBCperformance,table2showstheresultsfor
NBC,C4.5,andSBC    algorithmsusing80%ofthedatafortraini       ngand20%for
testing.Thefiguresshowninboldreflectthewinningmethodoneachdataset.
ThelasttwocolumnsshowtheimprovementofSBC        overNBCandC4.5.

Table2.Accuracyofeachmethodusing5      -foldcross -validation(15iterations)

| Dataset | NBC | C4.5 | SBC | SBCvsNBC | SBCvsC4.5 |
|---------|-----|------|-----|----------|-----------|
| Ecoli | 81.99 | 78.65 | **83.27** | +1.6% | +5.9% |
| German | 75.35 | 74.00 | **76.21** | +1.1% | +3.0% |
| KrVsKp | 87.81 | **99.12** | 94.69 | +7.8% | -4.5% |
| Monk | 96.16 | **98.46** | 97.47 | +1.4% | -1.0% |
| Mushroom | 90.37 | **99.80** | 98.85 | +9.4% | -1.0% |
| Pima | 75.03 | 75.35 | **79.94** | +6.5% | +6.1% |
| Promoter | 87.66 | 66.67 | **88.72** | +1.2% | +33.1% |
| Soybean | 84.02 | 83.20 | **88.27** | +5.1% | +6.1% |
| Wisconsin | 95.78 | 92.63 | **97.38** | +1.7% | +5.1% |
| Vote | 89.54 | 95.29 | **96.61** | +7.9% | +1.4% |

Fromtable2,itisapparentthatSBCoutperformstheoriginalNBCin
EVERYdomain  ,givingtheaccuracyimprovementupto9.4%.SBCalso
outperformsC4.5inalmostallthedomain,givingtheaccuracyimprovementup
to33.1%.Eventhough,SBCcannotbeatC4.5insomecases,itstillgivesquite
bigimprovementovertheNaïveBayes(7.8      %,1.4%,and9.4%).
     Ourexperimentalresultsdemonstratethat C4.5doespickgoodfeaturesforits
decisiontree(especiallyonesthatarenearertotheroot),whichinturn
asymptoticallyimprovestheaccuracyoftheNaïveBayesianalgorithm,when

*only*t hosefeaturesareusedinthelearningprocess.Table3showsthenumberof
featuresselectedforSelectiveBayesianclassifier.Onalmostallthedatasets,
surprisinglymorethanhalfoftheoriginalattributeswereeliminated.30%or
lessofallattrib    utes *selected*wereshowninbold,whichmeansthatwecan
actuallypaynoattentiontomorethan70%oftheoriginaldataandstillachieve
highaccuracyinclassification.

Table3.Numberoffeaturesselected

| Dataset | #Attributes | #ofAttributes selected |
|---------|-------------|------------------------|
| Ecoli | 8 | 4 |
| GermanCredit | 20 | **6** |
| KrVsKp | 37 | **4** |
| Monk | 6 | 4 |
| Mushroom | 22 | **6** |
| Pima | 8 | 5 |
| Promoter | 57 | **5** |
| Soybean | 35 | 12 |
| Wisconsin | 9 | 4 |
| Vote | 16 | **3** |
| SyntheticData | 20 | 12 |

Forspeedupandscalabilityissues,weranSBConalargesyntheticdatajust
toseehowfast   itcanlearn.TherunningtimeforSBConoursyntheticdatagive
**1.14**and   **4.24**speedupovertheoriginalNBCandC4.5,respectively.Notethat
weonlyused    **2,000**instancesoutofthetotalof1,200,000instancesforC4.5
featureselectionprocess,whic    hmadeitaveryquickoperation.Hence,in
practice,ifthedatasetislargeenough,wecanevensamplemuchlessthan10%
ofdataforthefeatureselectionprocess.Thenumberofattributesselectedby
SBCwas **12**outofthetotalof20attributes.Tab        le4illustratesthemeanelapsed
time(userandsystemtime)foreachclassifieronthissyntheticdata,using
1,000,000instancesfortrainingand200,000instancesfortestdata.

Table4.MeanElapsedtimeforSyntheticDataset(sec)

| NBC | C4.5 | SBC |
|-----|------|-----|
| 37.546 | 139.5 | 32.912 |

TherunningtimesofbothSBCandNBCaremuchlessthanthatofC4.5
becauseBayesianclassifieronlyneedstogothroughthewholetrainingdata
once.Theyarealsospaceefficientbecausetheybuildupafrequencytablein
sizeofth  eproductofthenumberofattributes,numberofclassvalues,andthe
numberofvaluesperattribute    .SBC,comparingtoNBC,learnsfasterbecause
fewerattributesareinvolvedinlearning.However,itisobviousthatmostofthe
timespentinboth      algorithmswasonI/O,readingthetrainingdata.That
explainswhySBCtimedidnotreducemuchfromNBCtime.Ifthereexistsa

veryfastwayofremovingunwantedfeaturesfromaverylargedataset,SBC
wouldonlyneed25.746secondsandgive31.4%im      provementoverNBC.

## 6 Conclusion

Asimplemethod    toimprove  Naïve Bayesianlearning  thatusesC4.5decision
treestoselectfeatureshasbeendescribed.          Theempiricalevidenceshowsthat
thismethodisveryfastand      surprisinglysuccessful,giventheverydifferent
naturesofthetwoclassificationmethods.ThisSelectiveBayesianclassifieris
asymptoticallyatleastasaccurateasthebetterofC4.5andNaïveBayeson
almostall   thedomainsonwhichtheexperimentsw        ereperformed.Further,it
learnsfasterthanbothC4.5andNBoneachofthesedomains.

    ThisworksuggeststhatC4.5decisiontreessystematicallyselectgood
featuresforNaïveBayesianclassifiertouse.Webelievethereasonsarethat
C4.5doesnot   useredundantattributesinconstructingdecisiontrees,sincethey
cannotgeneratedifferentsplitsoftrainingdata.Whenfewtrainingexamplesare
available,C4.5usesthemostrelevantfeaturesitcanfind.Thehighaccuracy          of
SBCachieveswithfewtra    iningexamplesisindicativeofthefactthatusingthese
featuresforprobabilisticinduction    leadstohigheraccuracy    produced ineachof
thedomainswehaveexamined.

## References

[1] Domingos,P.andPazzani,M.OntheOptimalityoftheSimplieBayesianClassifier
    underZero -OneLoss.KulwerAcademicPublishers,Boston.

[2] Duda,R.O.andHart,P.E.(1973).PatternClassificationandSceneAnalysis.New
    York,NY:WileyandSons.

[3] Elkan,C.Bo   ostingandNaïveBayesianLearning.TechnicalReportNo.CS97          -557,
    DepartmentofComputerScienceandEngineering,UniversityofCalifornia,San
    Diego,Spetember1997.

[4] Hand,D.andYu,K.(2001)Idiot'sBayes        —NotSoStupidAfterAll?International
    StatisticalReview(2001),69,pp.385   -398.

[5] Keogh,E.andPazzani,M.LearningAugmentedBayesianClassifiers:Acomparisonof
    distribution-basedandclassification    -basedapproaches.Uncertainty99,7        [th]Int'l
    WorkshoponAIandStatistics,Ft.Lauderdale,Florida,    225-230.

[6] Langley,P.andSage,S.InductionofSelectiveBayesianClassifiers.Proceedingsofthe
    TenthConferenceonUncertaintyinArtificialIntelligence(1994).Seattle,WA:
    MorganKaufmann

[7] Ming,K.andZheng,Z.ImprovingthePerformanceofBoosting        forNaïveBayesian
    Classification.InProceedingsofthePAKDD    -99,pp.296 -305,Beijing,China.

[8] Pazzani,M.(1996).ConstructiveInductionofCartesianProductAttributes.
    Information,StatisticsandInductioninScience.Melbourne,Australia.

[9] Quinlan,J.R.(1993).C4.5:ProgramsforMachineLearning     .,CA:MorganKaufmann.