# AdaptiveBoostingTechniquesinHeterogeneousandSpatialDatabases

AleksandarLazarevic,ZoranObradovic
CenterforInformationScienceandTechnology,
TempleUniversity,Room303,WachmanHall(038    -24),
1805N.BroadSt.,Philadelphia,PA19122,USA,
aleks@ist.temple.edu,zoran@ist.temple.edu
phone:(215)204   –6265,fax:(215)204    -5082

**Abstract.** *Combiningmultipleclassifiersisaneffectivetechniqueforimprovingclassificationaccuracybyreducing thevariancethroughmanipulatingthetrain    ingdatadistributions.Inmanylarge    -scaledataanalysisproblems involvingheterogeneousdatabaseswithattributeinstability,however,standardboostingmethodsdonotimprove local classifiers(e.g.k  -nearestneighbors) duetotheirlowsensitivityto     dataperturbation . Here,weproposean adaptiveattributeboostingtechniquetocoalescemultiplelocalclassifierseachusingdifferentrelevantattribute information.Toreducethecomputationalcostsofk         -nearestneighbor(k  -NN)classifiers,anovelfas     tk -NN algorithmisdesigned.Weshowthattheproposedcombiningtechniqueisalsobeneficialwhenboostingglobal classifierslikeneuralnetworksanddecisiontrees.         Inaddition,amodificationoftheboostingmethodisdeveloped forheterogeneousspati   aldatabaseswithunstabledrivingattributesbydrawingspatialblocksofdataateach boostinground.Finally,whenheterogeneousdatasetscontainseveralhomogeneousdatadistributions,wepropose anewtechniqueofboostingspecializedclassifiers,wh        ereinsteadofasingleglobalclassifierforeachboosting round,therearespecializedclassifiersresponsibleforeachhomogeneousregion.Thenumberofregionsis identifiedthroughaclusteringalgorithmperformedateachboostingiteration.        Newboost ingmethodsappliedto syntheticspatialdataandreallifespatialdatashowimprovementsinpredictionaccuracyforbothlocalandglobal classifierswhenunstabledrivingattributesandheterogeneityarepresentinthedata.Inaddition,boosting specializedexpertssignificantlyreducesthenumberofiterationsneededforachievingthemaximalprediction accuracy.*

**Keywords:**Adaptiveattributeboosting,Spatialboosting,Clustering,Boostingspecializedexperts          ,Heterogeneous spatialdatabases

# 1.Intro duction

Incontemporarydatamining,manyrealworldknowledgediscoveryproblemsinvolvetheinvestigationof relationshipsbetweenattributesinheterogeneousdatasetswhererulesidentifiedamongtheobservedattributesin certainsubsetsdonotapply    elsewhere.Aheterogeneousdatasetcanbepartitionedintohomogeneoussubsetssuch thatlearningalocalmodelseparatelyoneachofthemresultsinimprovedoverallpredictionaccuracy.Inaddition, manylarge -scaledatasetsveryoftenexhibitattribu     teinstability,whichmeansthatthesetofrelevantattributesthat describesdataexamplesisnotthesamethroughtheentiredataspace.Thisisespeciallytrueinspatialdatabases, wheredifferentspatialregionsmayhavecompletelydifferentcharacte     ristics[18].

Itiswellknowninmachinelearningtheorythatacombinationofmanydifferentpredictorscanbeaneffective techniqueforimprovingpredictionaccuracy.Therearemanygeneralcombiningalgorithmssuchasbagging[5], boosting[9],orErr     orCorrectingOutputCodes(ECOC)[15]thatsignificantlyimproveglobalclassifierslike decisiontrees,rulelearners,andneuralnetworks.Thesealgorithmsmaymanipulatethetrainingpatternsusedby individualclassifiers(bagging,boosting)orthecl     asslabels(ECOC).Inmostcases,theimportanceofdifferent classifiersisthesameforallofthepatternswithinthedatasettowhichtheyareapplied.

Inordertoimprovetheglobalaccuracyofthewhole,anensembleofclassifiersmustbebothaccu     rateand diverse.Tomaketheensembleofclassifiersfor     heterogeneousdatabasesmoreaccurate,insteadof     applyingaglobal classificationmodelacrossentiredatasets,themodelsarevariedtobettermatchspecificneeds     ofthesubsetsthus improvingp redictioncapabilities[21].Insuchanapproach,thereisaspecializedclassificationexpertresponsible foreachregionwhichstronglydominatestheothersfromthepoolofspecializedexperts.

Ontheotherhand,diversityisrequiredtoensurethatall     theclassifiersdonotmakethesameerrors.Inorderto increasethediversityofcombinedclassifiersfor     heterogeneous spatial databaseswithattributeinstability,one cannotassumethatthesamesetofattributesisappropriateforeachsingleclassi     fier.Foreachtrainingsample, drawninabaggingorboostingiteration,adifferentsetofattributesisrelevantandthereforetheappropriate attributesetshouldbeusedforconstructingsingleclassifiersineveryiteration.Inaddition,theapplicat     ionof differentclassifiersonspatialdatabases,wherethedataarehighlyspatiallycorrelated,mayproducespatially correlatederrors[19].Insuchsituations,standardcombiningmethodsmightrequiredifferentschemesfor manipulatingthetrainingin     stancesinordertomaintainclassifierdiversity.

Inthispaper,weextendtheframeworkforconstructingmultipleclassifiersystemusingtheAdaBoostalgorithm [9].Inourapproach,wefirsttrytomaximizelocalspecificinformationforadrawnsample bychangingtheattribute representationusingattributeselection,attributeextractionandappropriateattributeweightingmethods[22]ateach boostingiteration.Second,inordertoexploitthespatialdataknowledge,amodificationoftheboostingmet hod appropriateforheterogeneousspatialdatabasesisproposed,where,ateachboostinground,spatialdatablocksare drawninsteadofsamplingsingleinstanceslikeinthestandardapproach.Finally,themaximalgainbyemphasizing localinformation,es peciallyforhighlyheterogeneousdatasets,wasachievedbyallowingtheweightsofthe differentweakclassifierstodependontheinput.Ratherthanhavingconstantweightsoftheclassifiersforalldata patterns(asinstandardapproaches),wealloww eightstobefunctionsovertheinputdomain.Inordertodetermine theseweights,ateachboostingiterationweidentifylocalregionshavingsimilarcharacteristicsusingaclustering algorithmandthenbuildspecializedclassificationexpertsoneachof theseregionswhichdescribetherelationship betweenthedatacharacteristicsandthetargetclass[18].Insteadofasingleclassifierbuiltonasampledrawnin eachboostingiteration,thereareseveralspecializedclassificationexpertsresponsible foreachofthelocalregions identifiedthroughtheclusteringprocess. Alldatapointsbelongingtothesameregionandhencetothesame classificationexpertwillhavethesameweightswhencombiningtheclassificationexperts.

Theinfluenceofallof theseadjustmentsisnotthesame,however,forlocalclassifiers[4](e.g.k –nearest neighbors,radialbasisfunctionnetworks)andglobalclassifiers(e.g.decisiontreesandartificialneuralnetworks).It isknownthatstandardcombiningmethodsdonot improvesimplelocalclassifiersduetocorrelatedpredictions acrosstheoutputsfrommultiplecombinedclassifiers[5,15].Weshowthat,byselectingdifferentattribute representationsforeachsample,predictionofcombinednearestneighboraswella sglobalclassifierscanbe considerablydecorrelated.Ourexperimentalresultsindicatethatsamplingspatialdatablocksduringboosting iterationsisbeneficialonlyforlocalbutnotforglobalclassifiers.Furthersignificantimprovementsinpredictio n accuracyobtainedbybuildingspecializedclassifiersresponsibleforlocalregionsshowthatthismethodseemstobe slightlymorebeneficialfork -nearestneighboralgorithmsthanforglobalclassifiers,althoughthetotalprediction accuracywassigni ficantlybetterwhencombiningglobalclassifiers.

Thenearestneighborclassifierisoftencriticizedforslowrun -timeperformanceandlargememoryrequirements, andusingmultiplenearestneighborclassifierscouldfurtherworsentheproblem.Therefore, weusedanovelfast methodfork -nearestneighborclassificationtospeeduptheboostingprocess.

In Section 2, we discuss current ensemble approaches and work related to specialized experts and changing attribute representations of combined classifiers. Section 3 describes the proposed methods and investigates their advantages and limitations. In Section 4, we evaluate the proposed methods on three synthetic and one real-life data set comparing it with standard boosting and other methods for dealing with heterogeneous spatial databases. Finally, section 5 concludes the paper and suggests further directions in current research.

## 2. Classifier Ensembles

### 2.1. Ensembles of Local Learning Algorithms

One of the oldest and simplest methods for performing general, non-parametric classification that belongs to the family of local learning algorithms [4] is a k-nearest neighbor classifier (k-NN) [7]. Despite its simplicity, the k-NN classifier can often provide similar accuracy to more sophisticated methods such as decision trees or neural networks. Its advantages include the ability to learn from a small set of examples, and to incrementally add new information at runtime.

Many general algorithms for combining multiple versions of a single classifier do not improve the k-NN classifier at all. For example, when experimenting with bagging, Breiman[5] found no difference in accuracy between the bagged k-NN classifier and the single model approach. Kong and Dietterich[15] also concluded that ECOC would not improve classifiers that use local information due to higher error correlation.

A popular alternative to bagging is boosting, which uses adaptive sampling of patterns to generate the ensemble. In boosting [9], the classifiers in the ensemble are trained serially, with the weights on the training instances set adaptively according to the performance of the previous classifiers. The main idea is that the classification algorithm should concentrate on the difficult instances. Boosting can generate more diverse ensembles than bagging does, due to its ability to manipulate the input distributions. However, it is not clear how one should apply boosting to the k-NN classifier for the following reasons: (1) boosting stops when a classifier obtains 100% accuracy on the training set, but this is always true for the k-NN classifier, (2) increasing the weight on a hard to classify instance does not help to correctly classify that instance ease each prototype can only help classify its neighbors, not itself. Freund and Schapire[9] applied a modified version of boosting to the k-NN classifier that worked around these problems by

limitingeachclassifiertoasmallnumberofprototypes.However,theirgoalwasnottoimproveaccuracy,butto improvespeedwhilemaintainingcurrentperfo    rmancelevels.

Althoughthereisalargebodyofresearchonmultiplemodelmethodsforclassification,verylittlespecifically dealswithcombiningk    -NNclassifiers.RicciandAha[31]appliedECOCtothek              -NNclassifier(NN    -ECOC). Normally,applyingEC   OCtok    -NNwouldnotworksincetheerrorswouldbecorrelatedacrossthebinarylearning problems.However,theyfoundthatapplyingattributeselectiontothetwo              -classproblemsdecorrelatederrorsif differentattributeswereselected.Unlikethisappro        ach,Bay'sMultipleFeatureSubsets(MFS)method[3]uses randomattributeswhencombiningindividualclassifiersbysimplevoting.Eachtimeapatternispresentedfor classification,anewrandomsubsetofattributesisselectedforeachclassifier.Heu              sedtwodifferentsampling functions:samplingwithreplacement,andsamplingwithoutreplacement.Eachofthek              -NNclassifiersusesthe samenumberofattributes.

Someresearchersdevelopedtechniquesforreducingmemoryrequirementsfork              -NNclassifiers    bytheir combining.Incombiningcondensednearestneighbor(CNN)classifiers[1],thesizeofeachclassifier'sprototype setisdrasticallyreducedinordertodestabilizethek              -NNclassifier.Bootstrapordisjointdatasetpartitioningwas usedincombi      nationwithCNNclassifierstoeditandreducetheprototypes.InVotingnearestneighbor subclassifiers[16],threesmallgroupsofexamplesareselectedsuchthateachk              -NNsubclassifier,whenusedon them,errsinadifferentpartoftheinstancespace.              Simplevotingmaythencorrectmanyfailuresofindividual subclassifiers.

## 2.2.EnsembleofGlobalLearningAlgorithms

Therehasbeenaverysignificantmovementduringthepastdecadetocombinethedecisionsofglobalclassifiers (e.g.decisiontrees,    neuralnetworks),andasignificantbodyofliteratureonthistopichasbeenproduced.All combiningmethodsareresultsoftwoparallellinesofstudy:(1)multipleclassifiersystemsthatattempttofindan optimalcombinationofthedecisionsfromag          ivensetofcarefullydesignedglobalclassifiers;and(2)specialized classifiersystemsthatbuild      mutuallycomplementaryclassificationexperts,eachresponsibleforaparticulardata subset,  andthenmergethemtogether.Althoughitisknownthatmulti              pleclassifiersystemsworkwellwithglobal classifierslikeneuralnetworks,therehavebeenseveralexperimentsinselectingdifferentattributesubsetsasan

attempttoforcetheclassifierstomakedifferentandhopefullyuncorrelatederrorswhenanal yzingheterogeneous databases.

FeatureBoost[26]isarecentlyproposedvariantofboostingwhereattributesareboostedratherthanexamples. Whilestandardboostingalgorithmsalterthedistributionbyemphasizingparticulartrainingexamples,FeatureBoo st altersthedistributionbyemphasizingparticularattributes.ThegoalofFeatureBoostistosearchforalternate hypothesesamongsttheattributes.Adistributionovertheattributesisupdatedateachboostingiterationby conductingasensitivityana lysisontheattributesusedbythemodellearnedinthecurrentiteration.Thedistribution isusedtoincreasetheemphasisonunusedattributesinthenextiterationinanattempttoproducedifferentsub - hypotheses.

Onlyafewmonthsearlier,aconside rablydifferentalgorithmexploringasimilarideaforanadaptiveattribute boostingtechniquewaspublished[19].Thetechniquecoalescesmultiplelocalclassifierseachusingdifferent relevantattributeinformation.Therelatedattributerepresentation ischangedthroughattributeselection,extraction andweightingprocessesperformedateachboostinground.Thismethodwasmainlymotivatedbythefactthat standardcombiningmethodsdonotimprove local classifiers(e.g.k -NN)duetotheirlowsensiti vitytodata perturbation,althoughthemethodwasalsousedwithglobalclassifierslikeneuralnetworks.

Inadditiontothepreviousmethod,therewereafewmoreexperimentsselectingdifferentattributesubsetsasan attempttoforcetheneuralnetwork classifierstomakedifferentandhopefullyuncorrelatederrors.Althoughthereis noguaranteethatusingdifferentattributesetswilldecorrelateerror,TumerandGhosh[35]foundthatwithneural networks,selectivelyremovingattributescoulddecorre lateerrors.Unfortunately,theerrorratesintheindividual classifiersincreased,andasaresulttherewaslittleornoimprovementintheensemble.Cherkauer[6]wasmore successful,andwasabletocombineneuralnetworksthatuseddifferenthandsel ectedattributestoachievehuman expertlevelperformanceinidentifyingvolcanoesfromimages.

Motivatedbytheproblemofhowtoavoidoverfittingasetoftrainingdatawhenusingdecisiontreesfor classification,Ho[12]proposeda" *decisionforest* ",anensembleofdecisiontreesconstructedsystematicallyby autonomouslyandpseudorandomlyselectingasmallnumberofdimensionsfromagivenattributespace.The decisionsofindividualtreesarecombinedbyaveragingtheconditionalprobabilityofeac hclassattheleaves.The methodmaintainshighaccuracyonthetrainingdataand,comparedwithsingletreeclassifiers,improvesonthe generalizationaccuracyasitgrowsincomplexity.

Opitz[25]hasinvestigatedthenotionofanensemblefeaturesele ctionwiththegoaloffindingasetofattribute subsetsthatwillpromotedisagreementamongthecomponentmembersoftheensemble.Ageneticalgorithm approachwasusedforsearchinganappropriatesetofattributesubsetsforensembles.First,aniniti alpopulationof classifiersiscreated,whereeachclassifierisgeneratedbyrandomlyselectingadifferentsubsetofattributes.Then, thenewcandidateclassifiersarecontinuallyproduced,byusingthegeneticoperatorsofcrossoverandmutationon theattributesubsets.Thealgorithmdefinestheoverallfitnessofanindividualtobethecombinationofaccuracyand diversity.

DynaBoost[24]isanextensionoftheAdaBoostalgorithmthatallowsaninput -dependentcombinationofthe basehypotheses.As eparateweaklearnerisusedfordeterminingtheinputdependentweightsofeachhypothesis. Theerrorfunctionminimizedbytheseadditionalweaklearnersisamargincostfunctionthatisalsominimizedby AdaBoost.Althoughtheweightsdependontheinp ut,thereisstillasinglehypothesisperiterationthatneedstobe combined.

Severalapproachesbelongingtospecializedclassifiersystemshavealsoappearedlately.Ourrecentapproach [21]isdesignedforanalysisofspatiallyheterogeneousdatabase s.Itfirstcluststhedatainthespaceofobserved attributes,withanobjectiveofidentifyingsimilarspatialregions.Thisisfollowedbylocalpredictionaimedat learningrelationshipsbetweendrivingattributesandthetargetattributeinsideeac hcluster.Themethodwasalso extendedforlearningwhenthedataaredistributedatmultiplesites.

Asimilarmethodisbasedonacombinationofclassifierselectionandfusionbyusingstatisticalinferenceto switchbetweenthesetwo[17].Selection isappliedinregionsoftheattributespacewhereoneclassifierstrongly dominatestheothersfromthepool(clustering -and-selectionstep),andfusionisappliedintheremainingregions. Decisiontemplates(DT)areadoptedforclassifierfusion,where allclassifiersaretrainedovertheentireattribute spaceandtherebyconsideredascompetitiveratherthancomplementary.

Someresearchersalsohavetriedtocombineboostingtechniqueswithbuildingsingleclassifiersinorderto improvepredictionin heterogeneousdatabases.Onesuchapproachisbasedonasupervisedlearningprocedure, whereoutputsofpredictorsaretrainedondifferentdistributionsfollowedbyadynamicclassifiercombination[2]. ThisalgorithmappliesprinciplesofbothboostingandMixtureofExperts[13]andshowshighperformanceon classificationorregressionproblems.Theproposedalgorithmmaybeconsideredeitherasaboost -wiseinitialized MixtureofExperts,orasavariantoftheBoostingalgorithm.Asavariantofthe MixtureofExperts,itcanbemade

appropriateforgeneralclassificationandregressionproblems,byinitializingthepartitionofthedatasettodifferent expertsinaboostinglikemanner.IfviewedasavariantoftheBoostingalgorithm,itusesadyn amicmodelfor combiningtheoutputsoftheclassifiers.

## 3.Methodology

### 3.1AdaptiveAttributeBoosting

TheadaptiveattributeboostingalgorithmwepresenthereisavariantoftheAdaBoost.M2procedure[9].The proposedalgorithm,showninFigure1,p roceedsinaseriesof $T$ rounds.Ineveryroundaweaklearningalgorithmis calledandpresentedwithadifferentdistribution $D_t$ alterednotonlybyemphasizingparticulartrainingexamples, butalsobyemphasizingparticularattributes.Thedistribution isupdatedtogivewrongclassificationshigher weightsthancorrectclassifications.Theentireweightedtrainingsetisgiventotheweaklearnertocomputethe weakhypothesis $h_t$.Attheend,thedifferenthypothesesarecombinedintoafinalhypothesi s $h_{fn}$.

Sinceateachboostingiteration $t$ wehavedifferenttrainingsamplesdrawnaccordingtothedistribution $D_t$ ,atthe beginningofthe"forloop"inFigure1wemodifythestandardalgorithmbyadding **step0** ,whereinwechoosea differentattribute representationforeachsample.Differentattributerepresentationsarerealizedthroughattribute selection,attributeextractionandattributeweightingprocessesthroughboostingiterations.Thisis anattemptto forceindividualclassifierstomakedif ferentandhopefullyuncorrelatederrors.

- Given:SetS$\{(x_1,y_1),\ldots,(x_m,y_m)\}$x$_i \in$X,withlabe lsy$_i \in$Y=$\{1,\ldots,k\}$
  - LetB=$\{(i,y):i \in \{1,2,3,4,\ldots m\},y \neq y_i\}$
  - Initializethedistribution $D_1$ overtheexamples,suchthat $D_1$(i)=1/m.
- Fort=1,2,3,4,… $T$
  - *0. Findrelevantfeatureinformationfordistribution$D_t$ usingsupervisedattributeselection*
  1. Trainweaklearnerusingdistribution $D_t$
  2. Computeweakhypothesis $h_t$:X $\times$Y $\to$[0,1]
  3. Computethepseudo -lossofhypothesish $_t$:

  $$\varepsilon_t = \frac{1}{2} \cdot \sum_{(i,y)\in B} D_t(i,y)(1 - h_t(x_i,y_i) + h_t(x_i,y))$$

  4. Set $\beta_t = \varepsilon_t/(1 - \varepsilon_t)$
  5. UpdateD $_t$: $D_{t+1}(i,y) = (D_t(i,y)/Z_t) \cdot \beta_t^{(1/2)\cdot(1-h_t(x_i,y)+h_t(x_i,y_i))}$
     where $Z_t$ isanormalizationconst antchosensuchthat $D_{t+1}$ isadistribution.
- Outputthefinalhypothesis: $h_{fn} = \arg\max_{y\in Y} \sum_{t=1}^{T} (\log \frac{1}{\beta_t}) \cdot h_t(x,y)$

Figure1.Theadaptiveattributeboostingwithperformingattributeselectionatstep0ineachboostingiteration

To eliminate irrelevant and highly correlated attributes, regression-based attribute selection is performed through performance feedback forward selection and backward elimination search techniques [22] based on linear regression mean square error (MSE) minimization. The $r$ most relevant attributes are selected according to the selection criterion at each round of boosting, and are used by the single classifiers. In addition, attribute extraction procedure is performed through Principal Components Analysis (PCA) [10]. Each of the single classifiers uses the same number of new transformed attributes. Another possibility is to choose an appropriate number of newly transformed attributes that will retain some predefined part of the variance.

The attribute weighting method for the proposed technique is used only for local classifiers (k-NN) and is based on a 1-layer feedforward neural network. First, we try to perform target value prediction for the drawn sample with a defined 1-layer feedforward neural network using all attributes. It turns out that this kind of neural network can discriminate relevant from irrelevant attributes. Therefore, the neural networks interconnection weights are taken as attribute weights for the k-NN classifier.

To further experiment with attribute stability properties, miscellaneous attribute selection algorithms [22] are applied to the entire training set and the most stable attributes are selected. These attributes are then used by the standard boosting method. When applying adaptive attribute boosting, in order to compare the most stable selected attributes, the attribute occurrence frequency is monitored at each boosting round. When attribute subsets selected through boosting rounds become stable, this is an indication to stop the boosting process.

### 3.1.1 Adaptive Attribute Boosting for k-NN Classifier

Nearest neighbors are stable to the data perturbation, so bagging and boosting generate poor k-NN ensembles. However, they are extremely sensitive to the attributes used. Our approach attempts to use this instability to generate a diverse set of local classifiers with uncorrelated errors. At each boosting round, we perform one of the methods for changing attribute representation, explained above, to determine a suitable attribute space for use in classification. When determining the least distant instances, we consider standard Euclidean distance and Mahalanobis distance.

To speed up the long-lasting boosting process, a fast k-NN classifier is proposed. For $n$ training examples and $d$ attributes our approach requires preprocessing which takes O($d \cdot n \cdot \log n$) steps to sort each attribute separately. However, this is performed only once, and we trade off this initial time for later speedups.

Initially, we form a hyper-rectangle with boundaries defined by the extreme values of the $k$ closest values for each attribute (Figure 2 – small dotted lines). If the number of training instances inside the identified hyper-rectangle is less than $k$, we compute the distances from the test point to all of $d \cdot k$ data points which correspond to the $k$ closest values for each of $d$ attributes, and sort them into a non-decreasing array $sx$. We take the nearest training example $cdp$ with the distance $dst_{min}$, and form a hypercube with boundaries defined by this minimum distance $dst_{min}$ (Figure 2 - larger dotted lines). If the hypercube does not contain enough data, i.e. $k$ training points, form the hypercube of a side $2 \cdot sx(k+1)$. Although this hypercube contains more than $k$ training examples, we need to find the one which contains the minimal number of training examples greater than $k$. Therefore, if needed, we search for a minimal hypercube by binary halving the index in the non-decreasing array $sx$. This can be executed at most $logk$ times, since we are reducing the size of the hypercube from $2 \cdot sx(k+1)$ to $2 \cdot sx(1)$. Therefore the total time complexity of our algorithm is $O(d \cdot logk \cdot logn)$, under the assumption that $n > d \cdot k$, which is always true in practical problems.
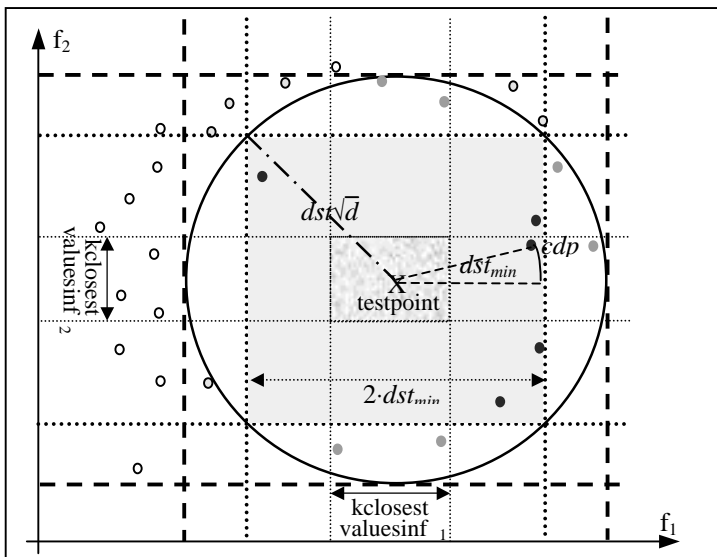


Figure 2. The used hyper-rectangle, hypersphere and hypercubes in the fast $k$-NN

If the number of training instances inside the identified hyper-rectangle (Figure 2 - small dotted lines) is larger than $k$, we also search for a minimal hypercube that contains at least $k$ and at most $2 \cdot k$ training instances inside that hypercube. This is accomplished by binary halving or by incrementing a side of the hypercube. After each modification of the hypercube's side, we compute the number of enclosed training instances and modify the hypercube accordingly. Analogously to the previous case, it can be shown that binary halving or incrementing the hypercube's side will not take more than $logk$ time, and therefore the total time complexity is still $O(d \cdot logk \cdot logn)$.

When we find a hypercube which contains the appropriate number of points, it is not necessary that all $k$ nearest neighbors are in the hypercube, since some of the closer training instances to the test point could be located in a hypersphere of identified radius $dst\sqrt{d}$ (Figure 2). Since there is no fast way to compute the number of instances inside the sphere without computing all the distances, we embed the hypersphere in a minimal hypercube (Figure 2 – dashed lines) and compute the number of training points inside this surrounding hypercube. The number of points inside the surrounding hypercube is much less than the total number of training instances and therefore speed up sour algorithm.

### 3.1.2 Adaptive Attribute Boosting for Global Classifiers

Although standard boosting can increase the prediction accuracy of global classifiers like neural networks [34] and decision trees [30], we change attribute representation to see if adaptive attribute boosting can further improve accuracy of an ensemble of global classifiers. The most stable attributes used in standard boosting of k-NN classifiers are also used here for the same purpose.

We train multilayer (2-layered) feed forward neural network classification models with the number of hidden neurons equal to the number of input attributes. The neural network classification models have the number of output nodes equal to the number of classes, where the predicted class is from the output with the largest response. We used two learning algorithms: resilient propagation [32] and Levenberg-Marquardt [11]. For a decision tree model, we used the ID3 learning algorithm [29] which employs the information gain criterion to choose which attribute to place at the root of each decision tree and subtree. After the trees are fully grown, a pruning phase replaces subtrees with leaves using the same predefined pruning factor for all trees.

### 3.2 Spatial Boosting

Spatial data represent a collection of attributes whose dependence is strongly related to spatial location; observations close to each other are more likely to be similar than observations widely separated in space. Explanatory attributes, as well as the target attribute in spatial datasets are very often highly spatially correlated. As a consequence, applying different classification techniques on such data is likely to produce errors that are also

spatiallycorrelated[27].Therefore,whenappliedtospatialdata, theboostingmethodmayrequiredifferent partitioningschemesthansimpleweightedselectionthatdoesnottakeintoaccountthespatialpropertiesofthedata.

Theproposedspatialboostingmethod(Figure3)startswithpartitioningthespatialdatas etintothespatialdata blocks(squaresofsizeMpoints $\times$Mpoints).Ratherthandrawing $n$datapointsaccordingtothedistribution $D_t$ (Figure1),theproposedmethoddrawsonly $\lfloor n/M^2 \rfloor$datapointsaccordingtothedistribution $P_t$(Figure3).Since each ofdrawnexamplesbelongsexactlytooneofthepartitionedspatialdatablocks,theproposedmethoddefines $\lfloor n/M^2 \rfloor$belongingspatialdatablocksandmergesthemintoasetusedforlearningaweakclassifier.Likeinstandard boosting,thedistribution $P_t$isalsoupdatedtogivewrongclassificationshigherweightsthancorrectclassifications, butduetospatialcorrelationofdata,attheendofeachboostingroundsimplemedianM $\times$Mfilteringisapplied overtheentiredatadistribution $P_t$.Usingthis approachwehopetoachievemoredecorrelatedclassifierswhose integrationcanfurtherimprovemodelgeneralizationcapabilitiesforspatialdata.Thespatialboostingtechniquewas appliedtobothlocal(k -NN)andglobal(neuralnetwork,decisiontrees) classifiers

---

- GivensetS$\{(x_1,y_1),\ldots,(x_m,y_m)\}$x$_i \in$X,withlabelsy$_i \in$Y=$\{1,\ldots,k\}$issplitinto $\lfloor n/M^2 \rfloor$squaresofsize
  MxMpoints.LetB=$\{(i,y):i \in \{1,2,3,4,\ldots m\}, y \neq y_i\}$
- Initializethedistribution $P_1$overtheexamples,suchthat $P_1$(i)=1/m.
- Fort=1,2,3,4,… $T$
  1. AccordingtodistributionP$_t$draw $\lfloor n/M^2 \rfloor$datapointsthatuniquelydeterminebelongingspatialdatablocks.
  2. Trainaweaklearneronasetcontainingal lbelongingspatialdatablocks.
  3. Computeweakhypothesis $h_t$:X $\times$Y $\to$[0,1]
  4. Computethepseudo -lossofhypothesish $_t$: $\mathcal{E}_t = \dfrac{1}{2} \cdot \sum_{(i,y) \in B} P_t(i,y)(1 - h_t(x_i, y_i) + h_t(x_i, y))$
  5. Set $\beta_t = \mathcal{E}_t/(1 - \mathcal{E}_t)$
  6. Update $P_t$: $P_{t+1}(i,y) = (P_t(i,y)/Z_t) \cdot \beta_t^{(1/2) \cdot (1 - h_t(x_i,y) + h_t(x_i,y_i))}$ where $Z_t$isanormalizationconstantchose n suchthat $D_{t+1}$isadistribution.ApplymedianM $\times$Mfilteringtothedistribution $P_t$.
- Outputthefinalhypothesis: $h_{fn} = \arg\max_{y \in Y} \sum_{t=1}^{T} (\log \dfrac{1}{\beta_t}) \cdot h_t(x,y)$

Figure3.Thespatialboostingalgorithmwithdrawingspatialdatablocksateachboostinground

## 3.3BoostingSpecializedClassifiers

Althoughpreviousboostingmodificationsimprovegeneralizabilityoffinalpredictors,itseemsthatin heterogeneousdatabaseswhereseveralmorehomogeneousregionsexistboostingdoesnotenhancetheprediction capabilitiesaswellasforhomogeneousdatabases[19].Insuchcasesitismoreusefultohaveseverallocalexperts

responsibleforeachregionof thedataset.Apossibleapproachtothisproblemistoclusterthedatafirstandthento

assignasingleclassifiertoeachdiscoveredcluster.Boostingspecializedclassifiers,describedinFigure4,modelsa

scenarioinwhichtherelativesignificance ofeachexpertadvisorisafunctionoftheattributesfromthespecific

inputpatterns.Thisextensionseemstobettermodelreallifesituationswhereparticularlycomplextasksaresplit

amongexperts,eachwithexpertiseinasmallspatialregion.

---

- Given:SetS=$\{(x_1,y_1),\ldots,(x_m,y_m)\}x_i \in X$,withlabelsy$_i \in Y=\{1,\ldots,k\}$
- LetB=$\{(i,y):i \in \{1,2,3,4,\ldots m\},y \neq y_i\}$
- Initializethedistribution $D_1$overtheexamples,suchthat $D_1(i)=1/m$.
- While(t< $T$ )or(globalaccuracyonset $S$startstodecrease)

1. Findrelevantattributeinformationfordistribution$D_t$usingunsupervisedwrapperapproacharound
   clusteringalgorithm.

2. Obtain $c$distributions $D_{t,j}$,j=1,… $c$andcorrespondingsets(clusters) $S_{t,j}=\{(x_{1,j},y_{1,j}),\ldots,(x_{m_j,j},y_{m_j,j})\}$

   $x_{i,j} \in X_j$, withlabelsy$_{i,j} \in Y_j=\{1,\ldots,k\}$byapplyingc lusteringwiththemostrelevantattributes identified
   instep1.LetB$_{j=} \{(i^j,y^j):i^j \in \{1,2,3,4,\ldots m^j\},y^j \neq y_i^j\}$

3. Forj=1… $c$(Foreachof $c$clusters)

   3.1.Findrelevantattributerepresentationforcluster s $S_{t,j}$usingsupervisedfeatureselection
   3.2.Trainweak learners $L_{t,j}$onthesets $S_{t,j}$, j=1,… $c$.
   3.3.Computeweakhypothesis $h_{t,j}$:X$_j \times$Y$_j \to$[0,1]
   3.4.Computeconvexhulls $H_{t,j}$foreachof $c$clusters $S_{t,j}$ fromtheentiresetS.
   3.5.Computethepse udo-lossofhypothesis $h_{t,j}$:

   $$\varepsilon_{t,j}= \frac{1}{2}\sum_{(i^j,y^j)\in B_j}D_{t,j}(i^j,y^j)(1-h_{t,j}(x_{i,j},y_{i,j})+h_{t,j}(x_{i,j},y^j))$$

   3.6.Set $\beta_{t,j}= \varepsilon_{t,j}/(1 - \varepsilon_{t,j})$
   3.7.Determineclustersontheentiretrainingsetaccordingtotheconvexhullmapping.Allpointsinside
        theconvexhull $H_{t,j}$belongtothe $j$-$th$cluster $T_{t,j}$from iteration $t$.

4. Mergeall $h_{t,j}$,j=1,...$c$ intoauniqueweakhypothesis $h_t$andall $\beta_{t,j}$,j=1,...$c$ intoaunique $\beta_t$according
   toconvexhullbelonging(examplefittinginthe $j$-$th$convexhullhasthehypothesis $h_{t,j}$andthevalue $\beta_{t,j}$).

5. UpdateD$_t$: $D_{t+1}(i, y) = (D_t(i, y)/Z_t)\cdot \beta_t(i, y)^{(1/2)\cdot(1+h_t(x_i,y_i)-h_t(x_i,y))}$

   where $Z_t$isanormalizationconstantchosensuchthat $D_{t+1}$isadistribution.

- Outputthefinalhypothesis: $h_{fn} = \arg\max_{y\in Y}\sum_{t=1}^{T}\bigcup_{j=1}^{c}(\log\frac{1}{\beta_{tj}(i^j,y^j)})\cdot h_{t,j}(x^j,y^j)$

Figure4. Theschemeforboostingspecializedclassifierswithperformingattributeselectionalgorithmwrapped
aroundclustering(step1)ineachboostingiteration.

Inthisworkasinmanyboostingalgorithms,thefinalcompositehypothesisisisconstructed asaweighted

combinationofbaseclassifiers.Thecoefficientsofthecombinationinthestandardboosting,however,donot

dependonthepositionofthepoint $x$whoselabelisofinterest.Theproposedboostingalgorithmachievesgreater

flexibilitybyb uildingclassifiersthatoperateonlyinspecializedregionsandhavelocalweights $\beta_t(x)$thatdependon

thepoint $x$wheretheyareapplied.

In order to partition the spatial dataset into these localized regions, two clustering algorithms are employed. The first is the standard $k$-means algorithm [14]. Here, dataset $S = \{(x_1, y_1), \ldots, (x_m, y_m)\}, x_i \in X$, is partitioned into $k$ clusters by finding $k$ points $\{m_j\}_{j=1}^{k}$ such that

$$\frac{1}{n} \sum_{x_i \in X} (\min_{j} d^2(x_i, m_j))$$

is minimized, where $d^2(x_i, m_j)$ usually denotes the Euclidiean distance between $x_i$ and $m_j$, although other distance measures can be used. The points $\{m_j\}_{j=1}^{k}$ are known as *cluster centroids* or cluster *means*.

The second clustering algorithm called DBSCAN relies on a density-based notion of clusters and was designed to discover clusters of an arbitrary shape [33]. The key idea of a density-based cluster is that for each point of a cluster its *Eps*-neighborhood for a given *Eps* > 0 has to contain at least a minimum number of points (*MinPts*), (i.e. the density in the *Eps*-neighborhood of points has to exceed some threshold), since the typical density of points inside clusters is considerably higher than outside of clusters. Unlike the cluster centroids in the $k$-means, here the centers of the clusters can be outside of the clusters due to their arbitrary shapes. Therefore, we define cluster medoids, the cluster *core* objects closest to the cluster centroids.

Since our boosting specialized experts involves clustering at step 1, there is a need to find a small subset of attributes that uncover "natural" groupings (clusters) from the data according to some criterion. For this purpose, we adopt the wrapper framework in unsupervised learning [8], where we apply the clustering algorithm to each attribute subset in the search space and then evaluate the attribute subset by a criterion function that utilizes the clustering result. If there are $d$ attributes in a dataset, an exhaustive search of the $2^d$ possible attribute subsets for the one that maximizes our selection criterion is computationally intractable. Therefore, in our experiments, fast sequential forward selection search is applied.

Like in [8] we also accept the scatter separability trace($S_w^{-1} S_b$) for attribute selection criterion, where $S_w$ is the within-class scatter matrix and $S_b$ is the between scatter matrix. $S_w$ measures the average covariance of each cluster and how scattered the samples are from their cluster medoids in the case of DBSCAN clustering, or from their cluster means in the case of k-means clustering. $S_b$ measures how the cluster means or medoids are distant from the total mean. Larger the value of the trace($S_w^{-1} S_b$) results in larger the normalized distance between the clusters and therefore in better cluster discrimination.

This procedure, performed at step 1 of every boosting iteration, results in $r$ most relevant attributes for clustering. Thus, for each round of boosting, there are different relevant attribute subsets that are responsible for distinguishing among homogeneous regions existing in a drawn sample. As a result of the clustering, applied to find those homogeneous regions, several distributions $D_{t,j}$ ($j=1,\ldots,c$) are obtained, where $c$ is the number of discovered clusters. For each of $c$ clusters $S_{t,j}$ discovered in the data sample, we identify its most relevant attributes, train a weak learner $L_{t,j}$ using a distribution $D_{t,j}$ and compute a weak hypothesis $h_{t,j}$. Furthermore, for every cluster $S_{t,j}$, we identify its convex hull $H_{t,j}$ in the attribute space used for clustering, and map these convex hulls to the entire training set in order to find the corresponding clusters $T_{t,j}$ (Figure 5) [20]. Data points inside the convex hull $H_{t,j}$ belong to cluster $T_{t,j}$, and data points outside the convex hulls are attached to the cluster containing the closest data pattern. Therefore, instead of a single global classifier constructed in every iteration by the standard boosting approach, there are $c$ classifiers $L_{t,j}$ and each of them is applied to the corresponding cluster $T_{t,j}$.
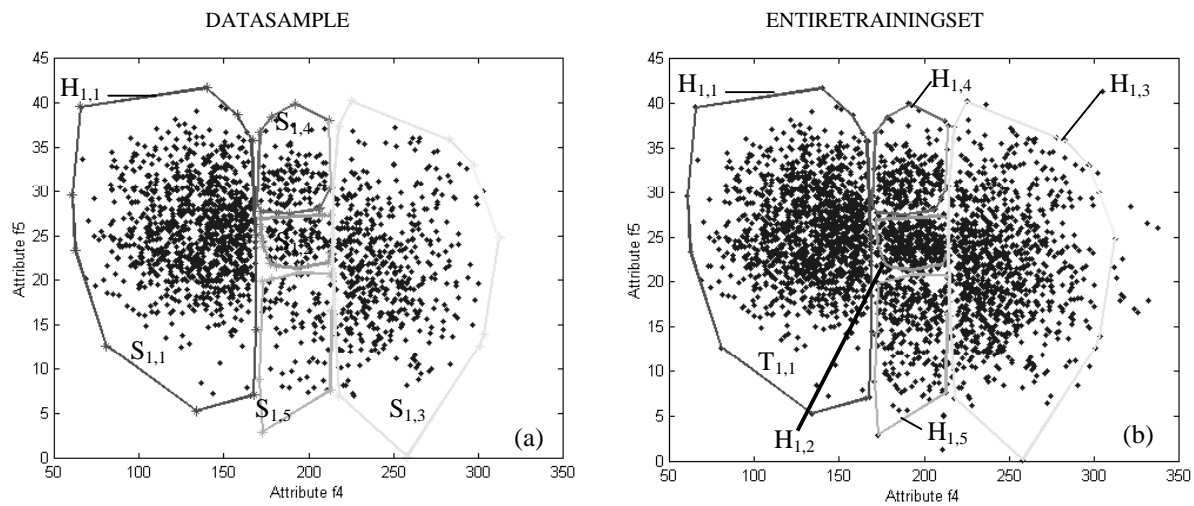


Figure 5. Mapping convex hulls $H_{1,j}$ of clusters $S_{1,j}$ discovered in the data sample to the entire training set in order to find corresponding clusters $T_{1,j}$. For example, all data points inside the contours of the convex hull $H_{1,1}$ (corresponding to the cluster $S_{1,1}$) belong to the new cluster $T_{1,1}$ identified on the entire training set.

In our boosting specialized classifiers, data points from different clusters have different pseudo-loss values and different parameter values $\beta_t$. For each cluster $T_{t,j}$, ($j=1,\ldots,c$) (Figure 5) defined with the convex hull $H_{t,j}$, there is a pseudo-loss $\varepsilon_{t,j}$ and the corresponding parameter $\beta_{t,j}$. Both the pseudo-loss value $\varepsilon_{t,j}$ and parameter $\beta_{t,j}$ are computed independently for each cluster $T_{t,j}$ where a particular classifier $L_{t,j}$ is responsible. Before updating the distribution $D_t$, the parameters $\beta_{t,j}$ for $c$ clusters are merged into a unique vector $\beta_t$ such that the $i$-th pattern from the dataset that belongs to the $j$-th cluster specified by the convex hull $H_{t,j}$, corresponds to the parameter $\beta_{t,j}$ at the $i$-th position in the

vector $\beta_t$. Analogously, the hypotheses $h_{t,j}$ are merged into a single hypothesis $h_t$. Since we merged $\beta_{t,j}$ into $\beta_t$ and $h_{t,j}$ into $h_t$, updating the distribution $D_t$ can be performed as in standard boosting. However, the local classifiers from each round are first applied to the corresponding clusters and integrated into a composite classifier responsible for that round. The composite classifiers are then combined into the final hypothesis using the AdaBoost.M2 algorithm.

When performing clustering during boosting iterations, it is possible that some of the discovered clusters have an insufficient number of data points needed for training a specialized classifier. This number of data patterns is defined as a function of the number of patterns in the entire training set. Several techniques for handling this scenario are considered.

The first technique denoted as *simple* halts the boosting process every time a cluster with an insufficient size is detected. When the boosting procedure is terminated, only the classifiers from the previous iterations are combined in order to create the final hypothesis $h_{fn}$. More sophisticated techniques do not stop the boosting process, but instead of training the specialized classifier on an insufficiently large cluster, they employ the specialized classifiers constructed in previous iterations. When an insufficiently large cluster is identified, its corresponding cluster from previous iterations is detected using the convex hull matching (Figure 5) and the model constructed on the corresponding cluster is applied to the cluster discovered in the current iteration. To determine this model, the most effective method (*best_local*) takes the classifier constructed in the iteration where the *local* prediction accuracy for the corresponding cluster was maximal. In two similar techniques, the *previous* method always takes the classifiers constructed on the corresponding cluster from the *previous* iteration, while the *best_global* technique uses the classification models from the iteration where the *global* prediction accuracy was maximal. In all of these sophisticated techniques, the boosting procedure ceases when the pre-specified number of iterations is reached or there is a significant drop in the prediction accuracy for the training set.

Furthermore, drawing spatial data blocks in boosting iterations, employed in the spatial boosting technique, is also integrated in boosting specialized classifiers.


## 4. Experimental Results

Our experiments were first performed on three synthetic datasets generated using our spatial data simulator [28] such that the distributions of generated data resembled the distributions of real life spatial data. All datasets had had

6561patternswithfiverelevant(f1,...f5)andfiveirrelevantattributes(f6,...,f10)and threeequalsizeclasses.The firstdatasetstemmedfromhomogeneousdistribution,whilethesecondonewasheterogeneouscontaining five homogeneousdatadistribu tions.In heterogeneousdataset,the attributes f4andf5 weresimulatedtoformfive clustersintheirattributespace(f4,f5)usingthetechniqueoffeatureagglomeration[28].Furthermore,insteadof usingasinglemodelforgeneratingthetargetattr ibuteontheentirespatialdataset,adifferentdatageneration process usingdifferentrelevantattributes wasappliedpereachcluster .Thedegreeofrelevancewasalsodifferent foreachdistribution.Thehistogramsofallfiveattributesforhomogeno usdatasetaswellasforheterogeneousdata setwithfivedistributionsareshowninFigures6aand6brespectively.Whenapplyingboostingspecialized classifiers,wealsoexperimentedwiththeheterogeneousdatasetwheretheoneofattributesrelevant forclustering wasmissingonlyduringclusteringprocess,whileallattributeswereavailablewhentrainingspecializedclassifiers.
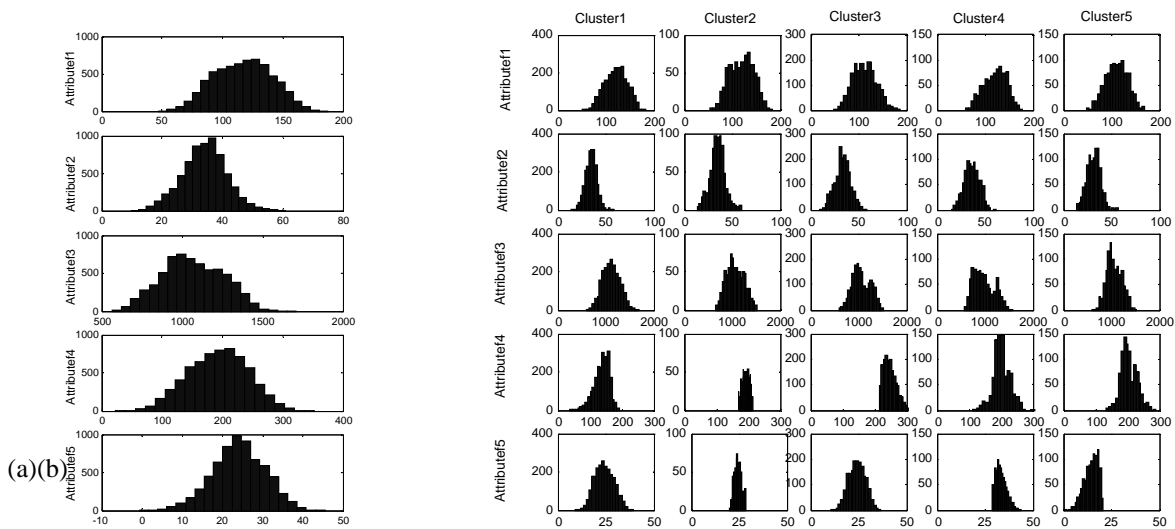


(a)(b)

Figure6. Histogramsofallfiverelevantat tributesfora)homogeneoussyntheticspatialdatasetb)heterogeneous syntheticspatialdatasetwithfiveclusters

Wealsoperformedexperimentsusingspatialdatafroma 220hafieldlocatednearPullman,WA. Allattributes wereinterpolatedtoa10x1 0mgridresultingin24,598patterns.ThePullmandatasetcontainedxandycoordinates (attributes1 -2),19soilandtopographicattributes(attributes3 -21)andthecorrespondingcropyield.

Forallperformedexperiments,syntheticandreallifedata setsweresplitintotrainingandtestdatasets.Theall reportedclassificationaccuracieswereachievedontestdatabyaveragingover10trialsoftheboostingalgorithm.

Forsyntheticdatasets,wefirstperformedstandardboostingandadaptiveattrib uteboosting(Figure7)for both local(k -NNclassifiers)andglobal(neuralnetworksanddecisiontrees)classifiers.Forthek -NNclassifier

experiments,thevalueof $k$ wassetusingcrossvalidationperformanceestimatesontheentiretrainingset.For

boostingneuralnetworkclassifiers,weusedthemodeldefinedinsection3.1.2.,and thebestpredictionaccuracies

wereachievedusing theLevenberq -Marquardtalgorithmfortrainingneuralnetworks.ForboostingID3decision

trees,weusedapost -pruningwithasmallconstantpruningfactorsuchthattheprunedtreesweresmallerthanthe

originalonesforapproximately20%.
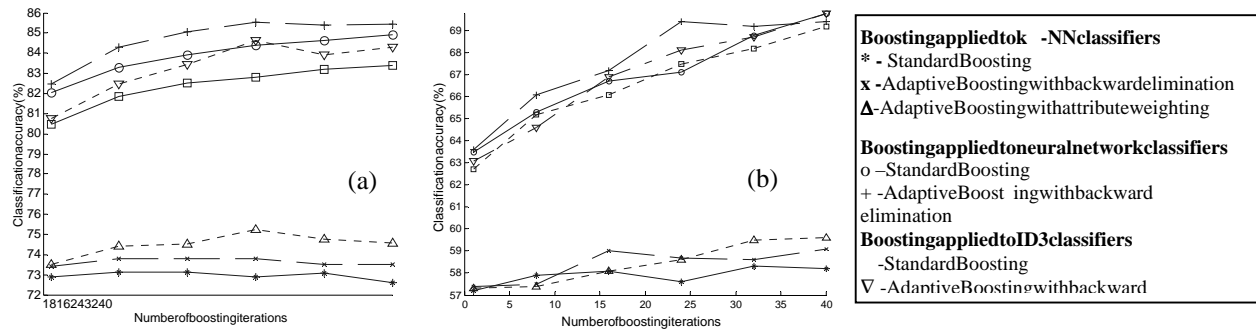


Figure7. Overallaveragedclassificationaccuracies(%)forthe3equal -sizeclassproblemson(a)homogeneous
synthetictes tdataset(b)heterogeneoussynthetictestdatasetwithfiveclustersdefinedby2of5relevantattributes.


AnalyzingthechartsinFigure7,itwasevidentthatthemethodofadaptiveattributeboostingappliedtolocal

andglobalclassifiersshowedo nlyminorimprovementsinpredictionaccuracyforbothsyntheticdatasets.Forthe

homogeneousdatasetthiswasbecausethereweredifferencesinrelevantattributesthroughthetrainingset,

whilefortheheterogeneousdatasetthiswasduetothefa ctthateachspatialregionnotonlyhaddifferentrelevant

attributesrelatedtoyieldclassbutalsoadifferentnumberofrelevantattributes.Insuchascenariowithuncertainty

regardingthenumberofrelevantattributesforeachregion,weneededto selectatleastthefourorfivemost

importantattributesateachboostinground,sinceselectingthethreemostrelevantattributesmaybeinsufficientfor

successfullearning.Sincethetotalnumberofrelevantattributesinthedatasetwasfiveaswel l,weselectedthefour

mostrelevantattributesforadaptiveattributeboosting,knowingthatforsomedrawnsampleswewouldlose

beneficialinformation.Duetothesefactsconcerningdeficientattributeinstability,t heselectedattributesduringthe

boostingiterationswerenotmonitored.

Inthestandardboostingmethod,weusedallfiverelevantattributesfromthedataset.Nevertheless,weobtained

similarclassificationaccuraciesforboththeadaptiveattributeboostingandthestandardboostingm ethod,but

adaptiveattributeboostingreachedthe"bounded"finalpredictionaccuracyinfewerboostingiterations.This

propertycouldbeusefulforreducingthetotalnumberoftheboostingrounds.Insteadofpost                -pruningtheboosted

classifiers[23]we  cantrytosettheappropriatenumberofboostingiterationsatthebeginningoftheprocedure.

   Applyingthespatialboostingmethodtoak               -NNclassifier,weachievedmuchbetterpredictionthanwiththe

adaptiveattributeboostingmethodsonak               -NNclass ifier(Table1).Furthermore,whenapplyingspatialboosting

withattributeselectionateachround,thepredictionaccuracywasincreasedslightlyasthesize(M)ofthespatial

blockwasincreased(Table1).Nosuchimprovementswerenoticedforspatial               boostingwithfixedattributesorwith

theattributeweightingmethod,andthereforetheclassificationaccuraciesforonlyM=5aregiven.

        Applyingspatialboostingonglobalclassifiers(neuralnetworksanddecisiontree)resultedinno

enhancementsin  classificationaccuracies.Moreover,forpurespatialboostingwithoutattributeselectionwe

obtainedslightlyworseclassificationaccuraciesthanusing"non               -spatial"boosting.Thisphenomenonwasdueto

spatialcorrelationofourattributes,whichmeans            thatdatapointsthatarecloseintheattributespaceareprobably

closeinrealspace,too.However,neuralnetworksordecisiontreesdonotconsiderspatiallocalinformationduring

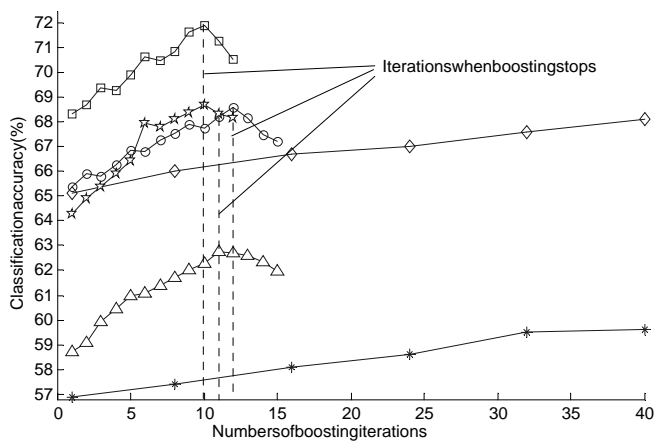thetraining,andunlikek   -NNdonotgainfromsamplingspatialdata            blocks.


Table1.   Overallaveragedclassificationaccuracy(%)ofspatialboostingforthe3equalsizeclassesonboth
synthetictestdatasetsusingk     -NNclassifiers.

| NumberofBoostingRounds | | Homogeneousdataset | | | | | Heterogeneousdatasetwith5clusters | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 8 | 16 | 24 | 32 | 40 | 8 | 16 | 24 | 32 | 40 |
| FixedAttributeSet(M=5) | | 79.1 | 79.6 | 80.1 | 80.7 | 80.6 | 65.6 | 65.5 | 65.8 | 66.0 | 66.1 |
| Backward Elimination | M=2 | 78.9 | 79.3 | 80.3 | 80.2 | 79.9 | 64.6 | 65.2 | 65.5 | 65.4 | 65.3 |
| | M=3 | 80.1 | 79.7 | 80.7 | 80.6 | 80.8 | 65.3 | 65.9 | 65.9 | 66.2 | 66.4 |
| | M=4 | 80.3 | 80.1 | 80.8 | 80.5 | 81.0 | 65.4 | 65.2 | 65.8 | 66.1 | 66.7 |
| | M=5 | 81.2 | 80.8 | 82.3 | 82.4 | 82.5 | 66.0 | 66.7 | 67.0 | 67.6 | 68.1 |
| AttributeWeighting(M=5) | | 79.4 | 78.8 | 80.1 | 80.7 | 80.3 | 64.2 | 64.7 | 65.4 | 66.3 | 65.9 |

   Whenperforming  boostingspecializedexperts(Table2,Figures8and              9)onheterogeneousdatasetwithall

attributes,insteadofperformingunsupervisedfeatureselectionaroundaclusteringalgorithmateachboosting

iteration,wealwaysappliedclusteringusingtheattributesf4andf5,sinceweknewthattheseattribute             sdetermine

homogeneousdistributions.Whenoneoftwoattributesresponsibleforclusteringwasmissing               ,weperformed

clusteringusingavailableclusteringattributeand            themostrelevantattributeobtainedthroughthefeatureselection

process. Inadditi on,wealwaysusedallfiverelevantattributesfortrainingspecializedclassifiers.Theexperiments

performedonhomogeneousdatasetshowedsimilarperformancelikeinheterogeneousdatawithmissingclustering

attributeandtheyarenotreportedhere.

Table2.Finalaveragedclassificationaccuracies(%)forthe3equalsizeclasses.Differentboostingalgorithmsare
appliedonbothsyntheticheterogeneoustestdatasets.

| Heterogeneousdatasets → | | | Setwithallrelevantattributes | | | Setwithmissingcluste ringattribute | | |
|---|---|---|---|---|---|---|---|---|
| Method | | | k-NN | NeuralNetwork | ID3 | k-NN | NeuralNetwork | ID3 |
| SingleClassifier | | | 57.3 | 61.0 ±2.2 | 63.3 | 57.3 | 61.0 ±2.2 | 63.3 |
| DBSCANClusteringwithsingle specializedclassifiers | | | 62.1 | 71.3 ±0.9 | 67.7 | 58.2 | 63.1 ±1.4 | 64.2 |
| StandardBoosting | | | 58.2 ±0.7 | 69.8 ±1.1 | 69.2 ±0.6 | 58.2 ±0.7 | 69.8 ±1.1 | 69.2 ±0.6 |
| AdaptiveAttributeBoosting | | | 59.1 ±0.6 | 69.4 ±1.1 | 69.8 ±0.6 | 59.1 ±0.6 | 69.4 ±1.1 | 69.8 ±0.6 |
| SpatialBoosting(M=5) | | | 68.1 ±0.9 | 69.1 ±1.2 | 68.2 ±0.07 | **68.1 ±0.9** | 69.1 ±1.2 | 68.2 ±0.07 |
| Boosting Specialized Expertswith Clustering | k-meansclustering | | **66.4** ±1.1 | **72.6** ±1.1 | **71.2** ±0.8 | 61.8 ±1.3 | 70.4 ±1.5 | 69.9 ±1.1 |
| | DBSCAN clustering | *simple* | 66.9 ±1.4 | 73.9 ±1.7 | 72.1 ±1.0 | 62.1 ±1.4 | 71.1 ±1.8 | 70.4 ±1.3 |
| | | *previous* | 67.4 ±1.3 | 74.4 ±1.5 | 72.8 ±1.2 | 63.3 ±1.5 | 71.3 ±1.9 | 70.5 ±1.3 |
| | | *best_global* | 67.9 ±1.3 | 74.9 ±1.4 | 73.4 ±1.1 | 62.4 ±1.4 | 71.6 ±1.5 | 70.8 ±1.1 |
| | | *best_local* | **68.6** ±1.1 | **76.6** ±1.2 | **74.5** ±0.9 | **62.7 ±1.3** | **71.9 ±1.4** | **71.1 ±1.2** |
| SpatialBoostingSpecialized Experts(DBSCAN+ *best_local*) | | | **71.9 ±1.0** | 76.4 ±1.3 | 74.4 ±1.0 | **68.6 ±1.1** | 71.4 ±1.5 | 70.8 ±1.3 |



**Boostingappliedtoheterogeneousdataset**
\* -AdaptiveAttributeBoosting
◊ -Drawingspatialblocks
o -Boosting specializedexperts with
**DBSCAN**clustering( *best_local*technique)
□ -SpatialBoosting specializedexpert swith
**DBSCAN**clustering( *best_local*)

**Boostingappliedtoheterogeneousdataset withmissingclusteringattribute**
Δ -Boosting specializedexperts with
**DBSCAN**clustering( *best_local*technique)
★-SpatialBoosting specializedexperts with
**DBSCAN**clustering ( *best_local*)

Figure8. Overallclassificationaccuracies(%)ofk -NNforthe3equal -sizeclassproblemsonheterogeneous synthetictestdatasetswith5relevantand5irreleva ntattributes.



(a)

(b)

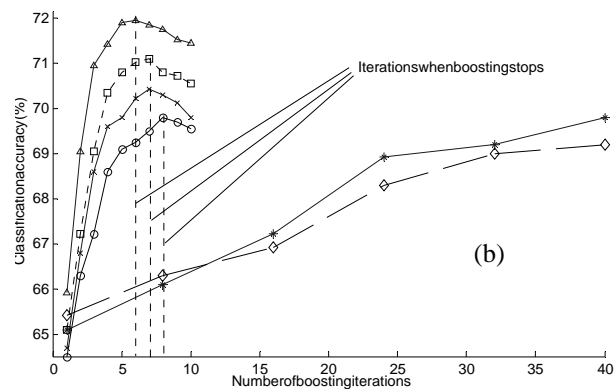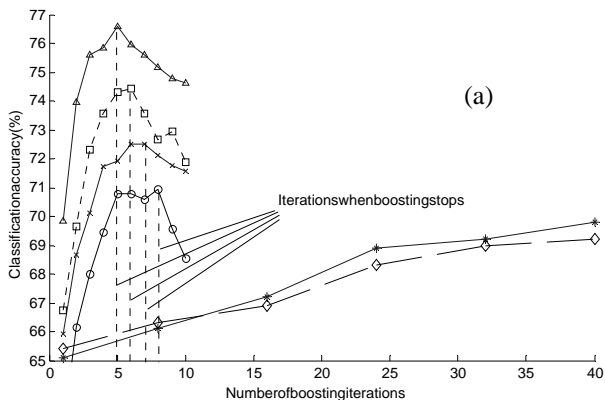Figure9.O verallclassificationaccuracies(%)forthe3equalsizeclassesforglobalpredictorsappliedon(a) heterogeneoussynthetictestdatasetwithallavailableattributes,(b)heterogeneoussynthetictestdatasetwith missingoneclusteringattribute.( \*-AdaptiveAttributeBoostingwithneuralnetworks,Boostingspecializedneural networkswith ✕ -k -meansclustering, Δ -DBSCANclustering(best_local), ◊ -AdaptiveAttributeBoostingwith ID3s,BoostingspecializedID3 classifierswith **o** -k -meansclustering, □-DBSCANclustering(best_local))

All methods of boosting specialized experts resulted in improved generalizations for all synthetic spatial data sets. However, improvements for heterogeneous dataset with all attributes (approximately 68–77%) were much more significant than for heterogeneous dataset with missing clustering attribute (approximately 63-72%) as compared to 57–63% obtained by single classifiers, specialized classifiers built on identified clusters, standard boosting and adaptive attribute boosting as shown at Table 2, Figures 8 and 9. Therefore, it is apparent that the prediction accuracy of all methods for boosting specialized experts directly depends on the quality of identified clusters during boosting iterations.

Boosting specialized experts is slightly more beneficial when boosting $k$-NN classifiers than global prediction models (Table 2), since the discovered clusters emphasize the local information, which is more helpful for local learning algorithms than for the global ones. Compared to the pure boosting specialized experts, the spatial boosting of global specialized classifiers again did not significantly affect the overall classification accuracy, while the influence of drawing spatial blocks when boosting specialized $k$-NN classifiers was reduced as compared to the improvements of pure spatial boosting over the standard and adaptive attribute boosting. This is due to the observed phenomenon that the smaller discovered clusters are not totally spatial, i.e. they contain scattered points in the spatial domain, and, in such cases, drawing spatial blocks does not help in reducing the total classification error.

It was also evident that the boosting of specialized experts required significantly fewer iterations in order to reach the maximal prediction accuracy. After prediction accuracy was maximized, the overall prediction accuracy on the training set, as well as the total classification accuracy on the test set, started to decline due to the fact that in the later iterations only data points that were difficult for learning were drawn. Therefore, there was not sufficient number of data examples in identified clusters needed for successful learning, and the prediction accuracy on these clusters began to deteriorate thus causing the drop of the total prediction accuracy too.

The data distribution of clusters discovered by applying DBSCAN clustering algorithm to heterogeneous dataset with all attributes was monitored at each boosting iteration (Figure 10). Unlike the previous adaptive attribute boosting method when around 30 boosting iterations were needed to achieve good generalization results, here typically only a few iterations (5-8 for global classification models and 8-12 for $k$-NN classifiers) were sufficient. As observed in Figure 10, data samples drawn in initial iterations (iteration 1) clearly included data points from all five clusters while samples drawn in later iterations (iterations 4, 5) contained a very small number of data points from the clusters where the prediction accuracy was good in previous iterations. As one of the criteria for stopping

boostingearly,westoptheboostingprocedurewhenthesizeofanyofthediscoveredclustersislessthansome

predefinednumber(usually  lessthan40).Anadditionalstoppingcriterionistoobservetheclassificationaccuracy

ontheentiretrainingsetandtostoptheprocedurewhenitstartstodecline.Figures8and9showtheiterationswhen

westoppedtheboostingprocedure.Although    inpracticethepredictionaccuracyonthetestsetdoesnotnecessarily

starttodropinthesameiteration,thisdifferenceisusuallywithintwoboostingiterationsanddoesnotsignificantly
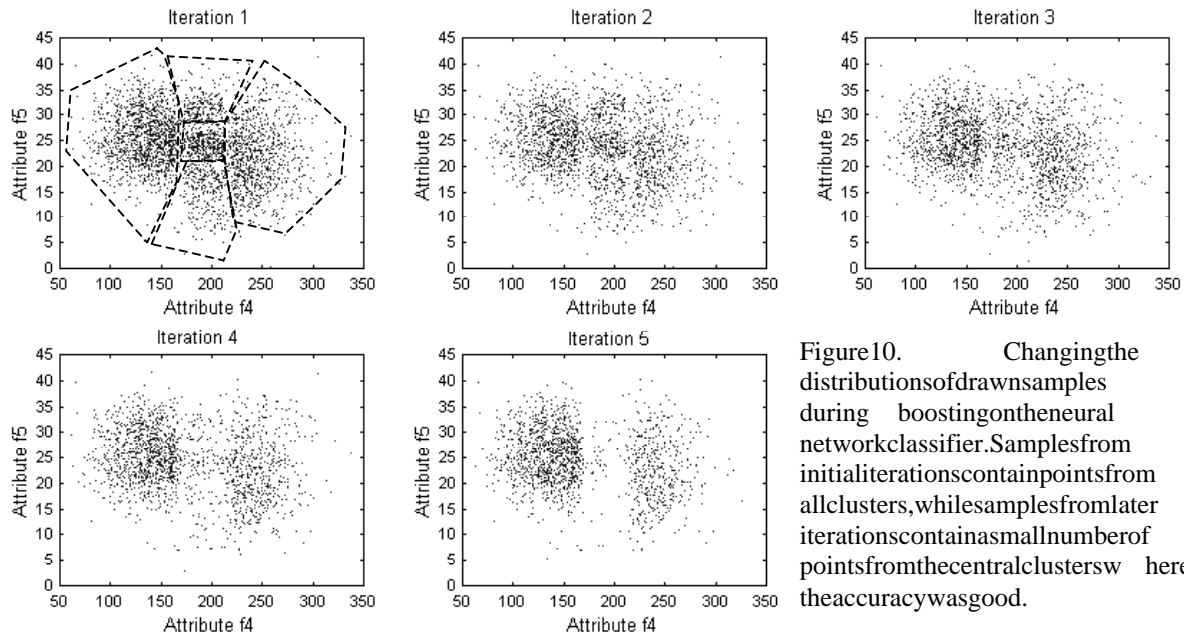
affectthetotalgeneralizabilityoftheproposedmethod.



Figure10.        Changingthe
distributionsofdrawnsamples
during    boostingontheneural
networkclassifier.Samplesfrom
initialiterationscontainpointsfrom
allclusters,whilesamplesfromlater
iterationscontainasmallnumberof
pointsfromthecentralclustersw  here
theaccuracywasgood.

Whenusingthek  -meansclusteringalgorithmduringtheboostingprocedure,wedidnotnoticethephenomenon

ofreducingthesizeofdiscoveredclustersandthereforewedidnotperformthemodificationsoftheproposed

algorithm.Inaddition,itwa  sevidentthatboostingspecializedexpertswhenusingthek        -meansclusteringalgorithm

wasnotassuccessfulasboostinglocalizedexpertswiththeDBSCANalgorithm,duetothebetterqualityofthe

clustersidentifiedbyDBSCANwhichwasdesignedtodisc        overspatialclustersofarbitraryshape.

Nevertheless,whenusingtheDBSCANalgorithmateachboostinground,the        *best_local*techniqueprovidedthe

bestpredictionaccuracy(Table2),whilethe        *simple*and  *previous*methodswerenotsignificantlybetterth        an

boostinglocalizedexpertswithk      -meansclustering.The        *simple*techniquefailedtoachieveimprovedprediction

results,becauseitdidnotreachenoughboostingiterationstodevelopthemostappropriateclassifiersforeach

clusterthatneededtobeco      mbined,whilethe  *previous*methodhadaboostingcyclethatwaslongenough,butdid

notcombineappropriatemodels.Finally,the        *best_global*and  *best_local*methodscombinedthemostaccurate

modelsforeachclustertakeninsomeoftheearlieriterations,        andhenceachievedthebestgeneralizability.

Experimentswithallproposedboostingmodificationswererepeatedforreallifespatialdata.Thegoalwasto
predict3equalsizeclassesofwheatyieldasafunctionofsoilandtopographicattributes.For              reallifedata(Pullman
dataset)16miscellaneousattributeselectionmethods(Table3)wereappliedonthetrainingdatasetinorderto
identifythefourmostrelevantattributesthatwereusedinthestandardboostingmethod.Histogramsforthesemost
stableattributes(4,7,9,20)areshowninFigure11.

Table3.Attributeselectionmethodsusedtoidentifythe4moststableattributesontraindataset.

| AttributeSelectionMethods | | | Selectedattributes |
|---|---|---|---|
| *Branch& Bound methods* | Probabilistic distance | Mahalanobisdistance | 7,9,11,20 |
| | | Bhatacharyadistance | 4,7,10,14 |
| | | Patrick-Fisherdistance | 13,17,20,21 |
| *Forward Selection methods* | Inter-class distance | Minkowski(order=1) | 7,9,10,11 |
| | | Minkowski(order=3) | 3,4,5,7 |
| | | Euclideandistance | 3,4 ,5,7 |
| | | Chebychevdistance | 3,4,5,7 |
| | Probabilistic distance | Bhatacharyadistance | 3,4,8,9 |
| | | Mahalanobisdistance | 7,9,11,20 |
| | | Divergencedistancemetric | 3,4,8,9 |
| | | Patrick-Fisherdistance | 13,16,20,21 |
| | MinimalErrorProbability,k  -NNwiths  ubstitution | | 4,7,11,19 |
| | Linearregressionperformancefeedback | | 5,9,7,18 |
| *Backward Elimination methods* | Probabilistic distance | Mahalanobisdistance | 7,9,11,20 |
| | | Bhatacharyadistance | 4,7,9,14 |
| | | Patrick-Fisherdistance | 13,17,20,21 |
| | Linearregressionperformancefeedback | | 7,9,11,20 |

Whenperformingattributeselectionduring         boostingonreallifedataset         ,thefourandfiveattributeswere
selectedandmonitoredandtheirfrequencywascomputed.         Thefrequencyofselectedattributesduringthe         boosting
rounds,whenboostingwasappliedtok         -NNclassifiers,neuralnetworkanddecisiontreeclassificationmodels,is
presentedinFigures12,13and14respectively.When         PCAwasusedwithboostingk         -NNclassifiers,projectionsto
fourdimensionsexp  lainedmostofthevarianceandtherewaslittleimprovementfromadditionaldimensions.For
theattributeweightingmethodinboostingk         -NNpredictors,weusedtheattributesynapticweightsbetweeninput
nodesandtheoutputnodeofa1         -layerneuralnetw  orkconstructedforeachdrawnsample.Whenboostingwas
appliedtoglobalclassifiers(neuralnetworkclassifiersanddecisiontrees),onlyattributeselectionproceduresfor
changingattributerepresentationwereconsidered.Theachievedclassificationa         ccuraciesforbothlocalandglobal
classifiersaregiveninTable4.
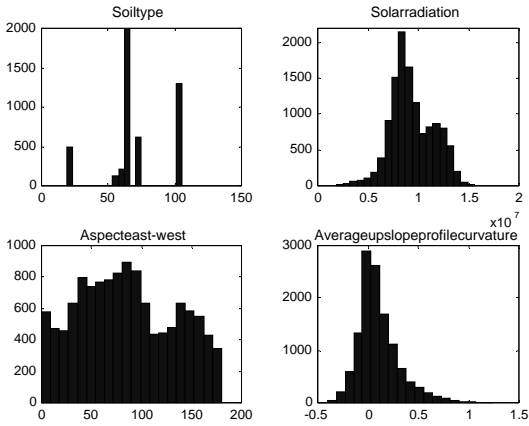
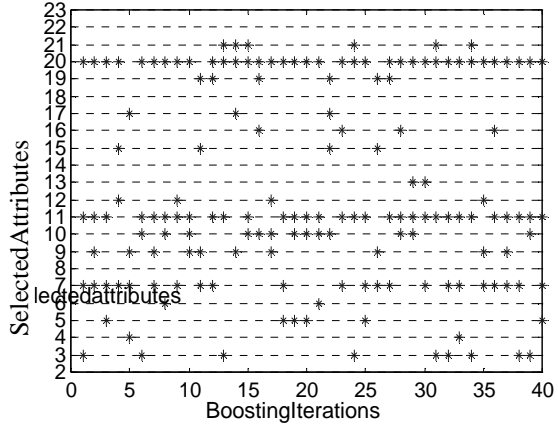Figure11. Histogramsof4mostrelevant attributesofreallifedataset



Figure12. Attributestabilityduringboostingonk -NN classifiers( *denotesthatattributeisselectedinboost - inground, -denotesthatattributeisnotselected)
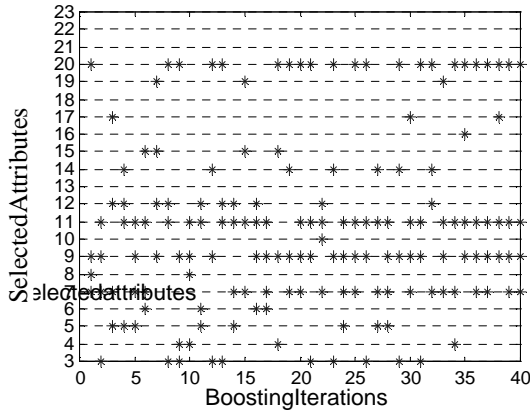


Figure13.Attributestabilityduring boostingon neuralnetworkwithLevenberq -Marquardtalgorithm
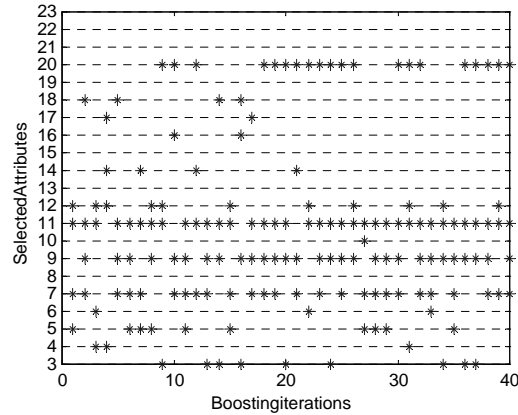


Figure14. Attributestabilityduring boostingonID3 decisiontreealgorithm

Table4. Comparativeanalysis ofoverallclassificationaccuracies(%)forthe3equal -sizeclassproblems onreallife testdatawit h 19soilandtopographicattributes.

| Number of Boosting Rounds | k-NNclassifier | | | | | Levenberg-Marquardt neuralnetworks | | ID3DecisionTrees | |
|---|---|---|---|---|---|---|---|---|---|
| | Standard Boosting | AdaptiveAttributeBoostingwith | | | | Standard Boosting | Backward Elimination | Standard Boosting | Backward Elimination |
| | | Forward Selection | Backward Elimination | PCA | Attribute Weighting | | | | |
| 8 | 38.2 | 40.9 | 38.5 | 42.4 | 43.0 | 43.6 | 47.5 | 43.3 | 46.9 |
| 16 | 39.5 | 41.3 | 38.8 | 42.4 | 43.9 | 44.1 | 47.8 | 43.7 | 47.3 |
| 24 | 38.8 | 41.9 | 42.1 | 44.5 | 44.8 | 44.8 | 48.3 | 44.3 | 47.8 |
| 32 | 38.5 | 41.8 | 43.5 | 45.1 | 46.1 | 45.5 | 48.8 | 45.0 | 48.2 |
| 40 | 39.3 | 42.1 | 42.8 | 43.4 | 44.3 | 44.9 | 48.5 | 45.2 | 48.4 |

ResultsfromTable4showthatthemethodsofadaptiveattributeboostingoutperformedthestandardboosting

techniqueforbothlocalandglobalclassifiers.Theresultsindicatethat30 boostingroundswereusuallysufficientto

maximizepredictionaccuracyandtosomewhatstabilizetheselectedattributesalthoughattributeselectionduring

boostingwaslessstablefork -NN(Figure12)thanforneuralnetworks(Figure13)ordecisiontr ees(Figure14).For

k-NN after approximately 30 boosting rounds the attributes became fairly stable with attributes 7, 11 and 20 obviously more stable than attributes 3 and 9, which also appeared in later iterations. The prediction accuracies when using k-NN classifier with Mahalanobis distance were worse than those using Euclidean distance, and are not reported here.

When boosting neural network classifiers we used models defined in section 3.1.2, and the best results were obtained using the applied backward elimination attribute selection and the Levenberq-Marquardt learning algorithm (Table 4). On the other hand, decision trees used all selected attributes for computing the splitting criterion, and after constructing they are pruned such that the number of nodes in pruned trees was reduced for 20%.

Classification accuracies of spatial boosting for k-NN classifiers on the real life dataset were again much better than without using spatial information and comparable to boosting neural networks and decision trees (Table 5). Here, the classification accuracy improvements from increasing the size (M) of the spatial blocks were less apparent than for synthetic spatial data probably due to the higher spatial correlation of the synthetic datasets.

Table 5. Overall classification accuracy (%) of spatial boosting for the 3 equal-size class problems for real life test data using k-NN classifiers.

| Number of Boosting Rounds | Spatial Boosting for k-NN with | | | | | |
|---|---|---|---|---|---|---|
| | Fixed Attribute Set | Backward Elimination Attribute Selection | | | | Attribute Weighting |
| | M=5 | M=2 | M=3 | M=4 | M=5 | M=5 |
| 8 | 46.4 | 45.8 | 47.7 | 48.1 | 47.8 | 45.2 |
| 16 | 46.6 | 46.2 | 47.6 | 48.1 | 47.7 | 45.6 |
| 24 | 46.7 | 46.7 | 47.9 | 48.2 | 48.2 | 45.8 |
| 32 | 46.9 | 46.9 | 48.3 | 48.4 | 47.9 | 46.3 |
| 40 | 47.0 | 47.2 | 48.3 | 47.9 | 47.8 | 45.9 |

When boosting specialized classifiers, all experiments were performed with the unsupervised wrapper procedure for identifying the most germane attributes for clustering and also with the supervised feature selection procedure for finding the most important attributes for each of the discovered clusters. In order to reduce the computational cost of the unsupervised wrapper approach, we did not identify more than three most appropriate attributes for clustering, since our previous experiments with clustering on the entire training set indicate that the best quality of clusters was obtained when using only two or three attributes [21]. The same experiments pointed out that modeling with four attributes results in the best prediction capability and therefore we were selecting only four attributes for constructing classifiers on discovered clusters. Figure 15 shows the overall classification accuracy when boosting k-

NNclassifiers,whiletheresultsinFigure16wereobtainedusingthe          Levenberq-Marquardtalgorithm for

optimizingneuraln  etworkparametersandusingtheprunedID3treeswitharelativelysmallpruningfactor.
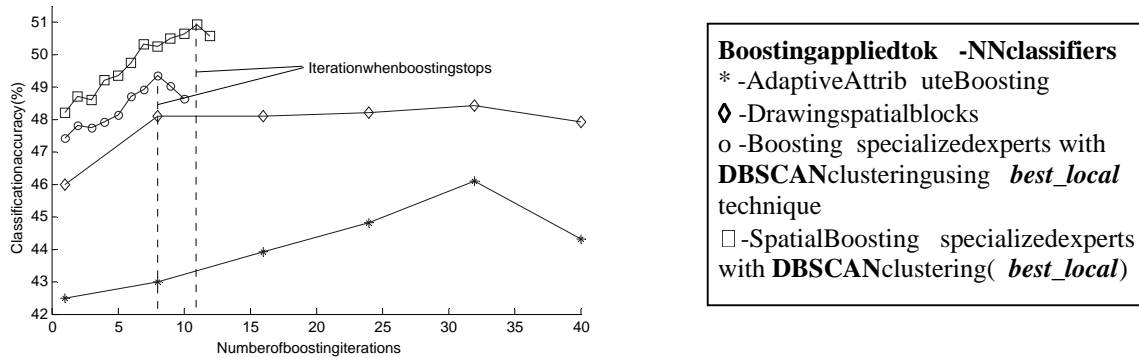


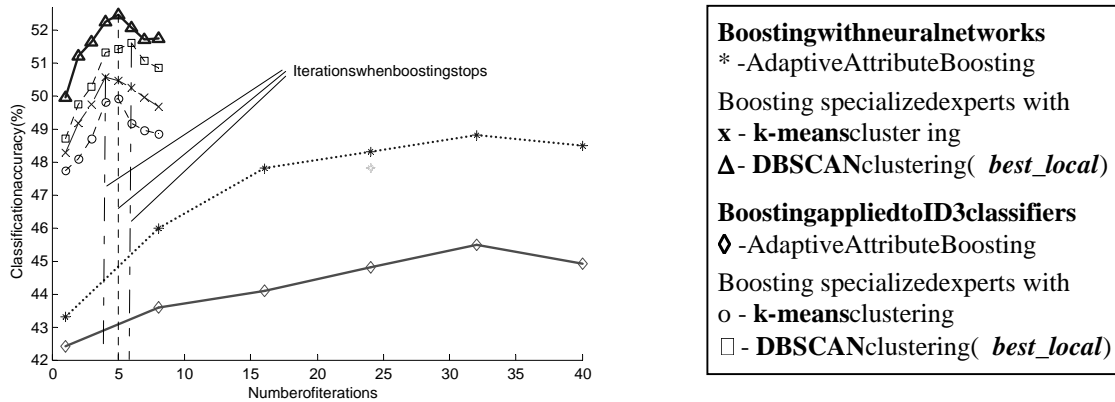Figure15.  Overallclassificationaccuraciesofk  -NNclassifiersforthe3    -classproblemsonreallifetestdata.



Figure16.  Overallclassification accuraciesofglobalpredictorsforthe3        -classproblemsonreallifetestdata.


    Boostingspecializedexpertsonareallifedatasetisnotassuperiortotheadaptiveattributeandspatialboosting

methodsasforthesyntheticheterogeneousdatasetwi          thallattributes.However,similarimprovementsinprediction

accuracywereachievedforsyntheticheterogeneousdatasetwithmissingclusteringattribute.            Thisindicatesthatin

reallifedata,itispossiblethereisalack          ofappropriatedrivingvaria   blesforexplainingthevariabilityofthetarget

attribute.  Thediscoveredspatialclustersinreallifedataarenotasdistinctasthespatialclustersinsynthetic

heterogeneousdatawithallattributes,butthehigherattributeinstabilitywasapparen            tlybeneficialforadaptive

attributeboosting.Unlikesyntheticheterogeneousdatasets,forreallifedatatheadditionaldiversityofconstructed

classifiersisachievedbyperformingunsupervisedattributeselectionandbydiscoveringclustersusingth            edifferent

attributes.Similartoexperimentsonsyntheticdata,the          *best_local*techniqueofboostinglocalizedexpertswasthe

mostsuccessfulamongalltheproposedmethods.

**5. Conclusions and Future Work**

Results from several spatial data sets indicate that the proposed techniques for combining multiple classifiers can result in significantly better predictions over existing classifier ensembles, especially for heterogeneous spatial data sets with attribute instabilities. First, this study provides evidence that by manipulating the attribute representation used by individual classifiers at each boosting round, classifiers could be more decorrelated thus leading to higher prediction accuracy. Second, our adaptive attribute boosting technique is more efficient than standard boosting, since a smaller number of iterations was sufficient to achieve the same final prediction accuracy. In addition, the attribute stability test served as a good indicator for properly stopping further boosting iterations. Third, the new boosting method proposed for spatial data showed promising results for k-NN classifiers making it competitive with powerful global classification models like neural networks and decision trees. Finally, boosting specialized experts with clustering performed at each boosting round further significantly improved both the prediction accuracy on highly heterogeneous databases and the efficiency of the algorithm by additional reducing the number of boosting iterations needed for achieving maximal prediction accuracy. However, for homogeneous data as well as for heterogeneous data sets with missing relevant attributes, the proposed method of boosting specialized classifiers showed only small improvements in achieved prediction accuracy.

Although boosting specialized experts required order of magnitude less boosting rounds to achieve the maximum prediction accuracy than the standard, adaptive attribute or spatial boosting, the number of constructed prediction models increases drastically through the iterations. This number depends on the number of discovered clusters and on the number of boosting rounds needed for making the final classifier. In our case, this drawback was alleviated by the fact that we were experimenting with small numbers of clusters and that only a few boosting iterations were sufficient to maximize prediction accuracy. Therefore, the memory needed for storing all prediction models is comparable or even less than for the standard boosting technique.

In addition to the prediction accuracy of the boosted specialized experts, the time required for building the model is also an important issue when developing a novel algorithm. Albeit the number of learned classifiers per iteration for the proposed method was much larger than for the standard boosting, the cluster data sets on which the classification models were built were smaller. The computation time for learning specialize experts was therefore comparable to learning the models on the entire training data. Hence, the total computation time depends only on the

numberofiterations,andismuchsmallerfortheproposedboostinglocalizedexpertsthanforstandardboostingor adaptiveattributeboosting.

Despitethefactthatthenewfastk -NNclassifiersignificantlyreducesthecomputationalreq uirements,anopen researchquestionistofurtherincreasethespeedofensemblesofk -NNclassifiersforhigh -dimensionaldata. Althoughtheperformedexperimentsprovideevidencethattheproposedapproachcanimprovepredictionsby ensemblesofbothlo calandglobalclassifiers,furtherworkisneededtoexaminetheadaptationofglobalclassifiers whenboostingspatialdata.Inordertousetheadvantagesfrombothlocalandnon -linearpredictionmodels,weare currentlyexperimentingwithamethodof boostingradialbasisfunctions.Inaddition,weareworkingtoextendthe methodtoregression -basedproblems.

**6.References**

1. Alpaydin,E.,VotingoverMultipleCondensedNearestNeighbors,in *LazyLearning* ,(D.Aha,ed.),Kluwer , 115-132,1997.

2. Avnimelech,R.,Intrator,N.,BoostingMixtureofExperts:AnEnsembleLearningScheme, *Neural Computation*,11:475 -490,1999.

3. Bay,S.,NearestNeighborClassificationfromMultipleFeatureSubsets, *IntelligentDataAnalysis* ,3(3):191 - 209, 1999.

4. Bottou,L.,Vapnik,V.,LocalLearningAlgorithms, *NeuralComputation* ,4(6):888 -900,1992.

5. Breiman,L.:Baggingpredictors, *MachineLearning* 24,123 -140,1996.

6. Cherkauer,K.,HumanExpert -levelPerformanceonaScientificImageAnalysisTaskbyaSy stemUsing CombinedArtificialNeuralNetworks,in *WorkingNotesoftheAAAIWorkshoponIntegratingMultiple LearnedModels* ,(P.Chan,ed.),15 -21,1996.

7. Dasarathy,B.V.,NearestNeighbor(NN)Norms:NNPatternClassificationTechniques, *IEEEComputer SocietyPress* ,388 -397,1991.

8. Dy,J.andBrodley,C.,FeatureSubsetSelectionandOrderIdentificationforUnsupervisedLearning,in *Proceedingsofthe17thInternationalConferenceonMachineLearning*, 247 -254,2000.

9. Freund,Y.,andSchapire,R.E.,Experi mentswithaNewBoostingAlgorithm,in *Proceedingsofthe13th InternationalConferenceonMachineLearning*, 325 -332,1996.

10. Fukunaga,K., *IntroductiontoStatisticalPatternRecognition*, AcademicPress,SanDiego,1990.

11. Hagan,M.,Menhaj,M.B.,Trainingf eedforwardnetworkswiththeMarquardtalgorithm. *IEEETransactions onNeuralNetworks* 5,989 -993,1994.

12. Ho,T.K.,TheRandomSubspaceMethodforConstructingDecisionForests, *IEEETransactionsonPattern AnalysisandMachineIntelligence*, 20(8),832 -844,1998.

13. JordanM.,JacobsR.,HierarchicalMixtureofExpertsandtheEMAlgorithm. *NeuralComputation*, 6(2):181 - 214,1994.

14. Kaufman,L.,Rousseeuw,P., *Findinggroupsindata:anintroductiontoclusteranalysis*, Willey,NewYork, 1990.

15. Kong,E.B.,Diett erich,T.G.,Error -CorrectingOutputCodingCorrectsBiasandVariance,In *Proceedingsof the12thNationalConferenceonArtificialIntelligence*, 725 -730,1996.

16. Kubat,M.,Cooperson,M.,VotingNearestNeighborSubclassifiers,in *Proceedingsofthe17th International ConferenceonMachineLearning*, 503 –510,2000.

17. Kuncheva,L.,Bezdek,J.,Duin,R.,DecisionTemplatesforMultipleClassifierFusion:AnExperimental Comparison, *PatternRecognition*, 34,299 -314,2001.

18. LazarevicA,XuX,FiezT,ObradovicZ .,Clu stering-Regression-OrderingStepsforKnowledgeDiscoveryin SpatialDatabases,in *ProceedingsofIEEE/INNSInternationalConferenceonNeuralNetworks*, No.345, Session8.1B,1999.

19. Lazarevic,A.,Fiez,T.,Obradovic,Z.,AdaptiveBoostingforSpatia lFunctionswithUnstableDriving Attributes,in *ProceedingsofPacific -AsiaConferenceonKnowledgeDiscoveryandDataMining*, 329 –340, 2000.

20. Lazarevic,A.,Pokrajac,D.,andObradovic,Z.,DistributedClusteringandLocalregressionforKnowledge DiscoveryinMultipleSpatialDatabases,in *Proceedingsof8thEuropeanSymposiumonArtificialNeural Networks*,129 -134,2000.

21. Lazarevic,A.andObradovic,Z.,KnowledgeDiscoveryinMultipleSpatialDatabases,inreview.

22. Liu,L.andMotoda,H. *FeatureSelection forKnowledgeDiscoveryandDataMining* ,KluwerAcademic Publishers,Boston,1998.

23. Margineantu,D.,andDietterich,T.,Pruningadaptiveboosting,in *Proceedingsofthe14thInternational ConferenceonMachineLearning* ,211 -218,1997.

24. Moerland,P.,Mayora z,E.,DynaBoost:CombiningBoostedHypothesesinaDynamicWay,IDIAPResearch Report99 -09,1999.

25. Opitz,D.,FeatureSelectionforEnsembles,in *Proceedingsof16thNationalConferenceonArtificial Intelligence*,379 -384,1999.

26. O'Sullivan,J.,Langford, J.,Caruna,R.,Blum,A.,FeatureBoost:AMeta -LearningAlgorithmthatImproves ModelRobustness,in *Proceedingsofthe17th InternationalConferenceonMachineLearning* ,703 -710,2000.

27. Pokrajac,D.,Obradovic,Z., CombiningRegressiveandAuto -Regressive ModelsforSpatio -Temporal Prediction,in *Proceedingsof17thInternationalMachineLearningWorkshoponSpatialKnowledge* ,2000.

28. PokrajacD,FiezT,ObradovicZ.,ASpatialDataSimulatorforAgricultureKnowledgeDiscovery Applications,inreview.

29. Quinlan,R.,InductionofDecisionTrees, *MachineLearning* ,1(1),81 –106,1986.

30. Quinlan,R.,Bagging,BoostingandC4.5,in *Proceedingsofthe13thNationalConferenceonArtificial Intelligence,*725 –730,1996.

31. Ricci,F.,andAha,D.W.,Error -CorrectingOutpu tCodesforLocalLearners,in *Proceedingsofthe10th EuropeanConferenceonMachineLearning* ,280 -291,1998.

32. Riedmiller,M.,Braun,H.,ADirectAdaptiveMethodforFasterBackpropagationLearning:TheRPROP Algorithm,in *ProceedingsoftheIEEEInternat ionalConferenceonNeuralNetworks* ,586 –591,1993.

33. SanderJ.,EsterM.,KriegelH -P,XuX.,Density -BasedClusteringinSpatialDatabases:TheAlgorithm GDBSCANanditsApplications, *DataMiningandKnowledgeDiscovery* ,2(2):169 -194,1998.

34. Schwenk,H.,Be ngio,Y.,BoostingNeuralNetworks, *NeuralComputation* ,12:1869 -1887,1999.

35. Tumer,K.,andGhosh,J.,ErrorCorrelationandErrorReductioninEnsembleClassifiers, *ConnectionScience* 8,385 -404,1996.