

DS4102PC: Machine Learning

Mr. G Lachiram Assistant Professor Department of IT



Maisammaguda (V), Kompally - 500100, Secunderabad, Telangana State, India

UGC - Autonomous Institute Accredited by NBA & NAAC with 'A' Grade Approved by AICTE Permanently affiliated to JNTUH

References



TextBooks:

1. TomM.Mitchell,MachineLearning,IndiaEdition2013,McGraw HillEducation.

Reference Books:

- 1. TrevorHastie,RobertTibshirani,JeromeFriedman,hTheElements of Statistical Learning, 2nd edition, springer series in statistics.
- 2. EthemAlpaydın,Introductiontomachinelearning,secondedition,MITpress.



Prerequisites

For Machine Learning Course were commend that students meet the following prerequisites:

- Basic programmingskills(in Python)
- Algorithm design
- Basics of probability & statistics



Content

- Unit–1 Introduction, Concept Learning, Decision Tree
- Unit-2 LearningArtificialNeuralNetworks-1,ArtificialNeural Networks-2, Evaluating Hypothesis,
- Unit–3 BayesianLearning,Computationallearningtheory, Instance Based Learning,
- Unit-4 GeneticAlgorithms,LearningSetsofRules,
 - Reinforcement Learning
- Unit–5 AnalyticalLearning-1,AnalyticalLearning-2,CombiningInductiveand Analytical Learning



UNIT-1





Machine Learning Introduction

Ever since computers were invented, we have wondered whether they might be made to learn.

If we could understand how to program them to learn-to improve automatically with experience-the impact would be dramatic.

Imagine computers learning from medical records which treatments are most effective for new diseases

Houseslearningfromexperiencetooptimizeenergycostsbasedontheparticularus age patterns of their occupants.

Personal software assistants learning the evolving interests of their users in order to highlight especially relevant stories from the online morning newspaper



Example of Successful Learning

Learning to recognize spoken words
Learning to drive an autonomous vehicle
Learning to classify new astronomical structures
Learning to play world-class backgammon



WhyisMachineLearningImportant?

- Sometaskscannotbedefinedwell,exceptbyexamples(e.g., recognizing people).
- Relationshipsandcorrelationscanbehiddenwithinlargeamountsof data.Machine Learning/DataMining maybe able to find these relationships.
- Humandesignersoftenproducemachinesthatdonotworkaswellas desired in the environments in which they are used.
- Theamountofknowledgeavailableaboutcertaintasksmightbetoolarge for explicit encoding by humans (e.g., medical diagnostic).
- Environmentschangeovertime.
- Newknowledgeabouttasksisconstantlybeingdiscoveredbyhumans.It may be difficult to continuously re-designsystems "by hand".

T.Aparna,AssistantProfessor,CSE,NRCM



AreasofInfluenceforMachineLearning

- *Statistics:*Howbesttousesamplesdrawnfromunknownprobabilitydistributionsto help decidefrom which distributionsomenewsampleisdrawn?
- *Brain Models*: Non-linear elements with weighted inputs (Artificial NeuralNetworks)havebeensuggested assimple models of biological neurons.
- *AdaptiveControlTheory*:Howtodealwithcontrollingaprocesshavingunknown parameters that must be estimated during operation?
- *Psychology*:Howtomodelhumanperformanceonvariouslearningtasks?
- *ArtificialIntelligence*:Howtowritealgorithmstoacquiretheknowledgehumansare able to acquire, atleast, as well as humans?
- *EvolutionaryModels*:Howtomodelcertainaspectsofbiologicalevolutionto improve the performance of computer programs?



MachineLearning:ADefinition

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured byP,improveswithexperienceE.



Why"Learn"?

Learningissusedwhen:

- Humanexpertise doesnot exist(navigatingonMars)
- Humansare unabletoexplaintheirexpertise(speech recognition)
- Solution changes intime(routingona computernetwork)
- Solutionneedstobe adapted toparticular cases(userbiometrics)



Well-PosedLearningProblem

Definition: A computer program is said to learn from experience E with respect to someclassoftasksT andperformancemeasureP,ifitsperformance attasksinT,as measured by P, improves with experience E.

Tohave a well-defined learning problem, three features needs to be identified:

- 1. The classoftasks
- 2. Themeasureofperformancetobe improved
- 3. The source of experience



CheckersGame



Computer Science and Engineering ,NRCM



GameBasics

- Checkers is played by two players. Each player begins the game with 12 coloreddiscs.(Onesetofpiecesisblackandtheotherred.)Eachplayerplaceshis or her pieces on the 12 dark squares closest to him or her. Black moves first. Players then alternate moves.
- Theboard consistsof64squares, alternating between 32 dark and 32 light squares.
- Itispositionedsothateachplayerhasalightsquareontherightsidecornerclosest to him or her.
- A player wins the game when the opponent cannot make a move. In most cases, thisisbecausealloftheopponent'spieceshavebeencaptured,butitcouldalsobe because all of his pieces are blocked in.

RulesoftheGame



- Moves are allowed only on the dark squares, so pieces always move diagonally. Single pieces are always limited to forward moves (toward the opponent).
- Apiecemakinganon-capturingmove(notinvolvingajump)maymoveonlyone square.
- Apiecemakingacapturingmove(ajump)leapsoveroneoftheopponent'spieces, landinginastraightdiagonallineontheotherside.Onlyonepiecemaybecaptured in a single jump; however, multiple jumps are allowed during a single turn.
- When apiece iscaptured, itisremoved from the board.
- If a player is a ble to make a capture, there is no option; the jump must be made.
- If morethanone capture is available, the player is free to choose which ever heor she prefers.



RulesoftheGameCont.

- Whenapiecereachesthefurthestrowfromtheplayerwhocontrolsthatpiece,itis crownedandbecomesaking.Oneofthepieceswhichhadbeencapturedisplaced on top of the king so that it is twice as high as a single piece.
- Kingsarelimitedtomovingdiagonallybutmaymovebothforwardandbackward. (Remember that single pieces, i.e. non-kings, are always limited to forward moves.)
- Kings may combine jumps in several directions, forward and backward, on the same turn. Single pieces may shift direction diagonally during a multiple captureturn, but must always jump forward (toward the opponent).



Well-DefinedLearningProblem

Acheckerslearning problem:

- TaskT:playingcheckers
- Performancemeasure P:percent of games wonagainstopponents
- Training experienceE:playing practicegames against itself

Ahandwritingrecognition learning problem:

- TaskT:recognizingandclassifyinghandwrittenwordswithinimages
- PerformancemeasureP:percentofwordscorrectlyclassified
- TrainingexperienceE:adatabaseofhandwrittenwordswith given classifications

Arobotdriving learningproblem:



- TaskT:drivingonpublicfour-lanehighwaysusingvisionsensors
- Performance measureP:averagedistancetravelledbeforeanerror(asjudgedby human overseer)
- TrainingexperienceE:asequence of images and steering commands recorded While observing a human driver



DesigningaLearningSystem

- 1. ChoosingtheTrainingExperience
- 2. ChoosingtheTargetFunction
- 3. ChoosingaRepresentationfortheTarget Function
- 4. ChoosingaFunctionApproximationAlgorithm
 - 1. Estimating training values
 - 2. Adjusting theweights
- 5. TheFinalDesign



1.ChoosingtheTrainingExperience

- Thefirstdesignchoiceistochoosethetypeoftrainingexperiencefromwhich the system will learn.
- Thetypeoftrainingexperienceavailablecanhaveasignificantimpact on success or failure of the learner.

Therearethreeattributeswhichimpactonsuccessorfailureofthe learner

- 1. Whetherthetrainingexperienceprovidesdirectorindirectfeedbackregarding the choices made by the performance system.
- 2. Thedegreetowhichthelearnercontrols thesequence of training examples
- 3. Howwellitrepresentsthedistributionofexamplesoverwhichthefinal system performance mustbe measured.

1. Whetherthetrainingexperienceprovidesdirectorindirectfeedbackregarding the choices madeby

Forexample, incheckersgame:

- Inlearningto playcheckers, the systemmight learn from direct training examples consisting of individual *Checkersboardstates and the correct move for each*.
- Indirect trainingexamplesconsistingof the *movesequences and final outcomes of various games played*.
- Theinformationaboutthecorrectnessofspecificmovesearlyinthegamemustbeinferredindirectly from the fact that the game was eventually won or lost.
- Herethelearnerfacesanadditionalproblem of <u>creditassignment</u>, or determining the degree to which each move in the sequence deserves creditor blame for the final outcome.
- Creditassignment canbeaparticularly difficult problem because the game can be lost even when early moves are optimal, if the seare followed later by poor moves.
- Hence, learning from direct training feedback is typically easier than learning from indirect feedback.

2. Asecond important attribute of the training experience is the degree to which the learner controls is the sequence of training examples

Forexample, incheckersgame:

- The learner might depends on the teacher to select informative board states and to provide the correct move foreach.
- Alternatively,thelearnermightitselfproposeboardstatesthatitfindsparticularlyconfusingandaskthe teacher for the correct move.
- The learner may have complete control over both the board states and (indirect) training classifications, as it does when it learns by playing against itself with no teacher present.
- Noticeinthislastcasethelearnermaychoosebetweenexperimentingwithnovelboardstatesthatithasnotyet considered, or honing its skill by playing minor variations of lines of play it currently finds mostpromising.

3. A third attribute of the training experience is how well it represents the

distribution of examples overwhich the final system performance must be measured **•**

Learning ismostreliablewhenthe trainingexamplesfollowadistribution similartothatoffuture test examples.

Forexample, incheckersgame:

- In checkers learning scenario, the performance metric P is the percent of games the system wins in the worldtournament.
- If its training experience E consists only of games played against itself, there is an danger that this training experiencemight not be fully representative of the distribution of situations overwhich it will later be tested. Forexample,thelearnermightneverencountercertaincrucialboardstatesthatareverylikelytobeplayedby the human checkers champion.
- It is necessary to learn from a distribution of examples that is somewhat different from those on which the finalsystemwillbeevaluated.Suchsituationsareproblematicbecausemasteryofonedistributionofexamples will not necessary lead to strong performance over some other distribution.





3.ChoosingtheTargetFunction

The next design choice is to determine exactly what type of knowledge will be learned and how this will be used by the performance program.

- Let's begin with a checkers-playing program that can generate the legal moves from any boardstate.
- The program needs only to learn how to choose the best move from among these legalmoves. This learning task is representative of a large class of tasks for which the legal moves that define some large search space are known a priori, but for which the best search strategy is not known.



Given this setting where we must learn to choose among the legal moves, the most obvious choiceforthetypeofinformation to be learned is a program, or function, that chooses the best move for any given board state.

1. Let ChooseMovebethetarget functionand thenotationis

ChooseMove: $B \longrightarrow M$

Which indicate that this function accepts as input any board from the set of legal boardstates Band produces as output somemove from the set of legal moves M.

ChooseMove is an choice for the target function in checkers example, but this function will turn out to be very difficult to learn given the kind of indirect training experience available to our system



2. Analternative targetfunction is an*evaluationfunction* thatassignsa *numericalscore* to any given boardstate LetthetargetfunctionVandthenotation

 $V:B \longrightarrow R$

Which denote that V maps any legal board state from the set B to some real value

We intend for this target function V to assign higher scores to better board states. If thesystemcansuccessfullylearnsuchatargetfunctionV,itcaneasilyuseittoselect the best move from any current board position.



Let usdefine the target value V(b) for an arbitrary board state bin B, as follows:

- 1. If bisa finalboardstate that is won, then V(b)=100
- 2. If bisa finalboard state that islost, then V(b)=-100
- 3. If b is a finalboard state that is drawn, then V(b)=0
- 4. If b is a nota final state in the game, then V(b) = V(b'),

Whereb'isthebestfinalboardstatethatcanbeachievedstartingfromband playing optimally until the end of the game



3. ChoosingaRepresentationforthe TargetFunction

Letuschooseasimplerepresentation-foranygivenboardstate,thefunctionwill be calculated as a linear combination of the following board features:

XL: the number of black pieces on the boardx2: the number of red pieces on the boardx3:thenumberofblackkingsonthe boardx4: the number of red kings on the board

X5:thenumberofblackpiecesthreatenedbyred(i.e.,whichcan be captured on red's next turn)

x6: the number ofred pieces threatened byblack



 $\hat{V}(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_6 x_6$

Where,

- wothroughw6arenumericalcoefficients,orweights,tobechosenbythe learning algorithm.
- Learnedvaluesfortheweights w1throughw6will determinetherelative Importanceofthe variousboardfeaturesindetermining the valueoftheboard
 - Theweightwo willprovide anadditive constant to the board value

Partialdesignofacheckerslearningprogram:



- TaskT:playingcheckers
- PerformancemeasureP:percentofgameswonintheworldtournament
- TrainingexperienceE:gamesplayedagainstitself
- Targetfunction:V:Board $\longrightarrow R$
- Targetfunctionrepresentation

 $\hat{V}(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_6 x_6$

The first three items above correspond to the specification of the learning task, whereasthefinaltwoitemsconstitutedesignchoicesfortheimplementation of the learning program.



4. ChoosingaFunctionApproximation Algorithm

- Inordertolearnthetargetfunctionwerequireasetoftrainingexamples,each describing a specific board state band the training value V_{train}(b) for b.
- Eachtraining example is an ordered pair of the form (b, V_{train}(b)).
- For instance, the following training example describes a board state b in whichblackhaswonthegame(notex₂=0indicatesthatredhasnoremaining pieces) and for which the target function valueV_{train}(b) is therefore +100.

 $((x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0), +100)$



FunctionApproximationProcedure

- 1. Derivetrainingexamplesfromtheindirecttrainingexperienceavailableto the learner
- 2. Adjuststheweightswitobestfitthesetraining examples



Where, visthe learner'scurrentapproximation to Successor(b)denotesthenextboardstatefollowingforwhichitisagainthe program's turn to move

Rulefor estimatingtrainingvalues

 $V_{train}(b) \leftarrow \forall (Successor(b))$



Specify the learning algorithm for choosing the weights with the set of training examples $\{(b, V_{train}(b))\}$

Afirststepistodefinewhat wemeanbythebestfit tothetrainingdata.

• Onecommonapproach istodefinethebest hypothesis,orsetofweights,as that which minimizes the squared error E between the training values and the values predicted by the hypothesis.

$$E \equiv \sum_{\langle b, V_{train}(b) \rangle \in \ training \ examples} (V_{train}(b) - \hat{V}(b))^2$$

• Severalalgorithmsareknownforfindingweightsofalinearfunction that minimizeE.



Onesuchalgorithmiscalledthe*leastmeansquares*,or*LMStrainingrule*.Foreach observed training example it adjusts the weights asmall amount in the direction that reduces the error on this training example

LMS weightupdate rule:-For each training example (b,V_{train}(b))Usethecurrentweightstocalculatev(b) Foreachweightw_i,updateit as

 $w_i \leftarrow w_i + \eta(Vtrain(b) - v(b))\hat{x}_i$



g.,0.1)thatmoderatesthe



- When the error (Vtrain(b)-V(b))isz ero, no weights are changed.
 - When (Vtrain(b)-V(b))ispositive(i.e.,whenV(b)isby advised polyative of V(b); reducing the error.
 - $If the value of some feature x_i is zero, then its weight is not altered regardless of the$
- error, so that the only weights updated are those whose features actually occur on the training example board.
Thefinaldesignofcheckerslearningsystemcanbedescribedbyfourdistinct program modules that represent the central components in many learning systems

5. The Final Design







1. The Performance System is the module that must solve the given performance taskby using the learned target function(s).

Ittakesaninstanceofanewproblem(newgame)asinputandproducesatraceofits solution (game history) as output.

In checkers game, the strategy used by the Performance System to select its next move at each step is determined by the learned \forall evaluation function. Therefore, we expect its performance to improve as this evaluation function becomes increasingly accurate.

2. *The Critic* takes as input the history or trace of the game and produces as output a setoftrainingexamplesofthetargetfunction. Asshowninthediagram, eachtraining example in this case corresponds to some game state in the trace, along with an estimateV_{train} of the target function value for this example.



3. TheGeneralizer takes a sinput the training examples and produces an output hypothes is that is its estimate of the target function. It generalizes from the specific training examples, hypothesizing a general function that covers these examples and other cases beyond the training examples. Inour example, the Generalizer corresponds to the LMS algorithm, and the output hypothesis is the function vdescribed by the learned weights wo, ..., W6.

*4. TheExperimentGenerator*takesasinputthecurrenthypothesisandoutputsanew problem(i.e.,initial board state) for the Performance System to explore. Its roleisto pick new practice problems that will maximize the learning rate of the over all system.

Inourexample,theExperimentGeneratoralwaysproposesthesameinitialgame board to begin a new game.





PerspectivesofMachineLearning



Perspectiveofmachinelearninginvolvessearchingverylarge space of possible hypothesis to determine one that

Best fits the observed data and any prior knowledge heldby learner.



Issuesin MachineLearning

- What algorithms exist for learning general target functions from specific training examples? In what settings will particular algorithms converge to the desired function, given sufficient training data?Which algorithms perform best for which types of problems and representations?
- Howmuchtrainingdataissufficient?Whatgeneralboundscanbefound torelate the confidence in learned hypotheses to the amount of training experience and the character of the learner's hypothesis space?
- When and how can prior knowledge held by the learner guide the process of generalizingfromexamples?Canpriorknowledgebehelpfulevenwhenitisonly approximately correct?



- What is thebest strategyforchoosing a useful next training experience, and how does the choice of this strategy alter the complexity of the learning problem?
- What is the best way to reduce the learning task to one or more function approximation problems? Put another way, what specific functions should the system attempt tolearn? Can this process itself be automated?
- Howcanthelearnerautomaticallyalteritsrepresentationtoimproveitsabilityto represent and learn the target function?



ConceptLearning

- Learning involves acquiring general concepts from specific training examples. Example: People continually learn general concepts or categories such as "bird,""car,""situationsinwhichIshouldstudymoreinordertopasstheexam,"etc.
- Each such concept can be viewed as describing some subset of objects or events defined over a larger set
- Alternatively, each concept can be though to fasaBoolean-valued function defined over this larger set. (Example: A function defined over all animals, whose value is true for birds and false for other animals).

Conceptlearning-InferringaBoolean-valuedfunctionfromtrainingexamplesof its input andoutput



AConceptLearningTask

Consider the example task of learning the target concept "DaysonwhichmyfriendAldoenjoyshisfavoritewatersport."

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

Table-Describesasetof example days, recacharepresented byasetofattributes

47

Whathypothes isrepresentation is provided tothelearner?



Let'sconsiderasimplerepresentationinwhicheachhypothesisconsistsofaconjunction of constraints on the instance attributes.

Leteachhypothesisbeavectorofsixconstraints,specifyingthevaluesofthesix attributes Sky, AirTemp, Humidity,Wind,Water, and Forecast.

Foreachattribute, the hypothesis will either

- Indicatebya"?'thatanyvalueisacceptableforthis attribute,
- Specifyasingle required value(e.g.,Warm)forthe attribute,or
- Indicatebya"Φ"that novalue is acceptable



The hypothesis that PERSON enjoys his favo rite sport only on cold days with highhumidity (independent of the values of the other attributes) is represented by theexpression

(?, Cold, High,?,?,?)

Themostgeneralhypothesis-thateverydayisapositive example-isrepresented by (?,?,?,?,?,?)

The most specific possible hypothesis-that day is a positive example-is norepresented by

Notation



Theset of itemsover which the concept is defined is called the set of *instances*, which we denote by X.

*Example:*Xisthesetofall possibledays,eachrepresentedbytheattributes:Sky, AirTemp, Humidity,Wind,Water,and Forecast

Theconceptorfunctiontobelearnediscalledthe*targetconcept*,whichwedenote by c. ccan be anyBooleanvalued function defined overtheinstancesXc:X{O, 1}

Example: The target conceptcorresponds to the value of the attribute *EnjoySport* (i.e., c(x)=1if*EnjoySport*=Yes, and c(x)=0if*EnjoySport*=No).



- Instancesforwhichc(x)=1arecalledpositiveexamples,ormembersofthe target concept.
- Instancesforwhichc(x)=0arecallednegativeexamples,ornon-membersofthe target concept.
- Theorderedpair(x,c(x))todescribe the training example consisting of the instancex and its target concept value c(x).
- **D**todenote the set of available training examples
- The symbol *H* to denote the set of all possible hypotheses that the learner mayconsiderregardingtheidentityofthetargetconcept.Eachhypothesis *h* in *H*represents a Boolean-valued function defined over X

h:X $= \{0,1\}$

• The goal of the learner is to find a hypothesis such that h(x)=c(x) for all x in X.



Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes



• Given:

- Instances X: Possible days, each described by the attributes
 - Sky (with possible values Sunny, Cloudy, and Rainy),
 - AirTemp (with values Warm and Cold),
 - Humidity (with values Normal and High),
 - Wind (with values Strong and Weak),
 - Water (with values Warm and Cool), and
 - Forecast (with values Same and Change).
- Hypotheses H: Each hypothesis is described by a conjunction of constraints on the attributes Sky, AirTemp, Humidity, Wind, Water, and Forecast. The constraints may be "?" (any value is acceptable), "Ø" (no value is acceptable), or a specific value.
- Target concept c: $EnjoySport: X \rightarrow \{0, 1\}$
- Training examples D: Positive and negative examples of the target function (see Table 2.1).
- Determine:
 - A hypothesis h in H such that h(x) = c(x) for all x in X.

TABLE The EnjoySport concept learning task.

Any hypothesis found to approximate the target function well over a sufficiently large set of training examples will also approximate the target function well overother unobserved examples.

TheInductiveLearningHypothesis



ConceptlearningasSearch

- Conceptlearningcanbeviewedasthetaskofsearchingthroughalargespace of hypotheses implicitly defined by the hypothesis representation.
- Thegoalofthissearchistofindthehypothesisthatbestfitsthetraining examples.

Example,theinstancesXandhypothesesHintheEnjoySportlearningtask. The attribute Sky has three possible values, and AirTemp,Humidity

, *Wind, Water Forecast* each have two possible values, the instances pace X contains exactly

- 3.2.2.2.2=96 Distinctinstances
- 5.4.4.4.4=5120syntacticallydistinct hypotheses withinH.

Everyhypothesiscontainingoneormore" Φ "symbolsrepresentstheemptysetof instances; that is, it classifies every instance as *negative*.

1+(4.3.3.3.3.3)=973. Se m antically distinct h ypotheses

•



Computer Science and Engineering ,NRCM



General-to-SpecificOrderingofHypotheses

• Consider the two hypotheses

h₁= (Sunny,?,?,Strong,?,?) h₂= (Sunny,?,?,?,?)

- Consider these tso finst ances that are classified positive by hand by h2.
- h₂imposesfewerconstraintsontheinstance, itclassifies more instancesaspositive. So, any instance classified positive by h₁ will also be classified positiveby h₂.Therefore, h₂is more generalthan h₁.



General-to-SpecificOrderingofHypotheses

 Givenhypotheses hjandhk,hjismore-general-thanorequaldohkifandonlyifanyinstancethatsatisfieshkalsosatisfieshi

Definition: Let h_j and h_k be Booleanvalued functions defined over X. Then h_j is more general-than-or-equaltoh_k(written $h_j \ge h_k$) if and only if

$$(\forall x \in X)[(h_k(x) = 1) \rightarrow (h_j(x) = 1)]$$







 $x_1 = <$ Sunny, Warm, High, Strong, Cool, Same> $x_2 = <$ Sunny, Warm, High, Light, Warm, Same> $\begin{aligned} h_1 &= <Sunny, ?, ?, Strong, ?, ?> \\ h_2 &= <Sunny, ?, ?, ?, ?, ?> \\ h_3 &= <Sunny, ?, ?, ?, Cool, ?> \end{aligned}$

- In the figure, the box on the leftrepresents these tX of all instances, the box on the right the set H of all hypotheses.
- Eachhypothesiscorrespondstosomes ubsetofXthesubsetofinstancesthatitclassifies positive.
- The arrows connecting hypothesesrepresent the *more general than*relation,withthearrowpointingt

owardthelessgeneralhypothesis.

 Notethesubsetofinstancescharacteri zedbyh2subsumesthesubset characterizedbyh1,henceh2is moregeneral-thanh1



FIND-S:FindingaMaximallySpecific Hypothesis

FIND-SAlgorithm

- 1. Initialize h to the most specific hypothesis in H
- 2. Foreachpositivetraininginstancex
 - Foreachattributeconstraintainh
 - Iftheconstraintaissatisfiedbyx
 - Thendonothing

Elsereplace*ai*n*h*bythenextmoregeneralconstraintthatissatisfiedby *x*

3. Outputhypothesish



Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

The first step of FIND-Sistoinitialize h to the most specific hypothesis in Hh- $(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset)$

$x_1 = < SunnyWarmNormalStrongWarmSame>, +$



Observing the first training example, it is clear that our hypothesis is too specific. Inparticular, none of the " \emptyset " constraints in h are satisfied by this example, so each isreplaced by the nextmore general constraint that fits the example

$h_1 = < Sunny WarmNormalStrong WarmSame >$

Thishisstillveryspecific; it asserts that all linst ances are negative except for the single positive training example

*x*₂=*<Sunny*, *Warm*, *High*, *Strong*, *Warm*, *Same*>, +

The second training example forces the algorithm to further generalize h, this timesubstituting a "?' in place of any attribute value in h that is not satisfied by the newexample

h₂=<SunnyWarm?StrongWarmSame>



x3=<Rainy,Cold,High,Strong,Warm,Change>,-

Uponencounteringthethirdtrainingthealgorithmmakesnochangetoh.TheFIND-Salgorithm simply ignores everynegative example.

h3=<SunnyWarm?StrongWarmSame>

x4=<SunnyWarmHighStrongCoolChange>,+

Thefourthexampleleadstoafurthergeneralizationofh

h4=<SunnyWarm?Strong??>



NRCM your roots to success...

ThekeypropertyoftheFIND-Salgorithmis

- FIND-S is guaranteed to output the most specific hypothesis within H that isconsistent with the positive training examples
- FIND-S algorithm's final hypothesis will also be consistent with the negative examples provided the correct target concept is contained in H, and provided the training examples are correct.

T.Rupa Rani, AssistantProfessor, CSE, NRCM



UnansweredbyFIND-S

1. Hasthe

hear wrspherered to the specefter governess pt?

- 3. Arethetrainingexamplesconsistent?
- 4. What if there are several maximally specific consistent hypotheses?

T.Rupa Rani, AssistantProfessor, CSE, NRCM



VersionSpaceandCANDIDATEE LIMINATIONAlgorithm

ThekeyideaintheCANDIDATE-

ELIMINATIONalgorithmistooutputadescriptionofthesetofall*hypothesesconsistentwiththetr* ainingexamples

Representation

• **Definition:** Ahypothesishis**consistent** with a set of training examples D if and only if h(x) = c(x) for each example (x, c(x)) in D.

 $Consistent(h,D) \equiv (\forall \langle x, c(x) \rangle \in D) h(x) = c(x))$

Notedifferencebetweendefinitionsof consistent and satisfies

- Anexamplexissaidto*satisfy*hypothesis*h*when*h*(*x*)=1,regardless of whether *x* is a positive or negative example of the target concept.
- an example x is said to consistent with hypothesis h if h(x) = c(x)

T.Rupa Rani, AssistantProfessor, CSE, NRCM

VersionSpa ce



 $\label{eq:constraint} A representation of the set of all hypotheses which are consistent with D$

Definition: The version space, denoted $VS_{H,D}$ with respect to hypothesisspace H and training examples D, is the subset of hypotheses from H consistent with the training examples in D

 $VS_{H,D} = \{h \in H | Consistent(h,D)\}$





TheLIST-THEN-ELIMINATE Algorithm

TheLIST-THEN-ELIMINATEalgorithmfirstinitializestheversionspacetocontainall hypotheses in H and then eliminates any hypothesis found inconsistent with anytraining example.

DeepakD, Asst. Prof., Dept. of CSE, Canara Engg. College

TheLIST-THEN-ELIMINATE





- 1. VersionSpacecalistcontainingeveryhypothesisinH
- 2. Foreachtrainingexample,(x,c(x))

remove from *VersionSpace* any hypothesis *h* for which $h(x) \neq c(x)$

3. OutputthelistofhypothesesinVersionSpace

TheLIST-THEN-ELIMINATEAlgorithm

- List-Then-Eliminateworksinprinciple, solong as version space is finite.
- However, since it requires exhaustive enumeration of all hypotheses in practice it is not feasible.

AMoreCompactRepresentationforVersio nSpaces



- The version space is represented by its most general and least general members.
- Thesemembersformgeneralandspecificboundarysetsthatdelimittheversionspace within the partially ordered hypothesis space.



- Aversionspacewithitsgeneral and specific boundary sets
- The version space includes allsix hypotheses shown here, butcanberepresentedmoresimpl yby SandG.
- Arrowsindicateinstanceof the*more-general-than* relation.Thisistheversion spaceforthe *Enjoysport*conceptlearning
- problemandtrainingexamples describedinbelowtable


$G = \{g \in H | Consistent(g, D) \land (\neg \exists g' \in H) [(g' >_g g) \land Consistent(g', D)] \}$

Definition:The**specificboundary**S,withrespecttohypothesisspaceHandtraining data *D*, is the set of minimally general (i.e., maximally specific) members ofHconsistentwithD.

 $S = \{s \in H | Consistent(s, D) \land (\neg \exists s' \in H) [(s >_g s') \land Consistent(s', D)]\}$



VersionSpacerepresentationth eorem



Theorem: Let X be an arbitrary set of instances and Let H be a set of Boolean-valued hypotheses defined over X. Let $c : X \rightarrow \{O, 1\}$ be an arbitrary target conceptdefined over X, and let D be an arbitrary set of training examples $\{(x, c(x))\}$. For allX, H,c, and D such that Sand G are well defined,

 $VS_{H,D} = \{h \in H | (\exists s \in S) (\exists g \in G) (g \geq_g h \geq_g s)\}$





ToProve:

1. Everyheatisfyingtherighthandsideoftheaboveexpressionisin*VS*_{*H*,*D*} 2. Everymemberof*VS*_{*H*,*D*}satisfiestheright-handsideoftheexpression

Sketchofproof:

1. letg,h,sbearbitrarymembersofG,H,Srespectivelywith $g \ge_g h \ge_g s$

By the definition of *S*, **s** must be satisfied by all positive examples in D. Because $h \ge_{gS}$, h must also be satisfied by all positive examples in D.

Bythedefinition of G, gcannot be satisfied by any negative example in D. Because h is satisfied by all positive examples in Dandbynone gative examples in D, his consistent with D, and therefore his a member of $VS_{H,D}$

2. It can be proven by assuming some hin $VS_{H,D}$, that does not satisfy the right-hand side of the expression, then showing that this leads to an inconsistency

TheCANDIDATE-ELIMINATIONLearningAlgorithm

The CANDIDATE-ELIMINTION algorithm computes the *version space* containingall hypotheses from H that are consistent with an observed sequence of trainingexamples.

Initialize G to the set of maximally general hypotheses in HInitialize S to the set of maximally specific hypotheses in HForeachtraining example d, do

- Ifdisapositiveexample
 - RemovefromGanyhypothesisinconsistentwithd
 - ForeachhypothesissinSthatisnotconsistentwithd
 - RemovesfromS
 - AddtoSallminimalgeneralizationshofssuchthat
 - hisconsistent with d, and some member of Gismoregeneral than h
 - Remove from Sanyhypothesis that is more general than another hypothesis in S
- Ifdisanegativeexample
 - RemovefromSanyhypothesisinconsistentwithd
 - ForeachhypothesisginGthatisnotconsistentwithd
 - RemovegfromG
 - AddtoGallminimalspecializationshofgsuchthat
 - hisconsistent with d, and some member of Sismore specific than h
 - $\bullet \ Remove from Gany hypothesis that is less general than another hypothesis in G$



AnIllustrativeExam



esesin H.



S₀





(Sunny,Warm,Normal,Strong,Warm,Same)+







$\langle Sunny, Warm, High, Strong, Warm, Same \rangle +$









{Rainy,Cold,High,Strong,Warm,Change}-

unny, Warm, **S**₂,**S**₃ rong,Warm,Same>





\langle Sunny,Warm,High,Strong,CoolChange \rangle +





The final version space for the *EnjoySport* concept learning problem and



InductiveBias

Thefundamentalquestionsforinductiveinference

- Whatifthetargetconceptisnotcontained in the hypothesis space?
- Canweavoidthisdifficultybyusing ahypothesisspacethatincludeseverypossiblehypothesis?
- Howdoesthesizeofthishypothesisspaceinfluencetheabilityofthealgorithmtogeneralizet ounobservedinstances?
- Howdoesthesizeofthehypothesisspaceinfluencethenumberoftrainingexamples thatmustbeobserved?





Y

Y

N

Effectofincompletehypothesi sspace

PrecedingalgorithmsworkiftargetfunctionisinH WillgenerallynotworkiftargetfunctionnotinH

Considerfollowingexampleswhichrepresenttargetfunction "sky=sunnyorsky=cloudy": {SunnyWarmNormalStrongCoolChange} {CloudyWarmNormalStrongCoolChange} {RainyWarmNormalStrongCoolChange}

If apply Candidate Elimination algorithm as before, endup with empty Version Space After f

irsttwotrainingexample

 $S = \langle ?WarmNormalStrongCoolChange \rangle$

Newhypothesisisoverlygeneralanditcoversthethirdnegativetrainingexample!Our

AnUnbiasedLear

ner



Incompletehypothesisspace

- IfcnotinH,thenconsidergeneralizingrepresentationofHtocontainc
- ThesizeoftheinstancespaceXofdaysdescribedbythesixavailableattributesis
 96.ThenumberofdistinctsubsetsthatcanbedefinedoverasetXcontaining|X| elements(i.e., thesizeofthepowersetofX)is2|X|
- Recallthatthereare96instancesin*EnjoySport*;hencethereare296possiblehypothesesinfullsp aceH
- CandothisbyusingfullpropositionalcalculuswithAND,OR,NOT
- HenceHdefinedonlybyconjunctionsofattributesisbiased(containingonly973h's)



- Let us reformulate the *Enjoysport*learning task in an unbiased way by defining a newhypothesisspace*H*'thatcanrepresenteverysubsetofinstances;thatis,letH'correspondtoth epowersetofX.
- Onewayto definesuchan H'istoallowarbitrarydisjunctions,conjunctions,andnegationsofourearlierhypothes es.

Forinstance, the target concept ''Sky=SunnyorSky=Cloudy ''could then be described as (Sunny, ?, ?, ?, ?, ?)V(Cloudy, ?, ?, ?, ?)

Definition:



Consider a concept learning algorithm L for the set of instances X.

- LetcbeanarbitraryconceptdefinedoverX
- LetD_c={(x,c(x))}beanarbitrarysetoftrainingexamplesofc.
- $\bullet \ Let L(x_i, D_c) denote the classification assigned to the instance x_i by Lafter training on the data D_c.$
- TheinductivebiasofLisany minimalsetofassertionsBsuch thatforanytargetconceptcandcorrespondingtrainingexamplesDc

 $(\forall \langle x_i \in X) [(B \land D_c \land x_i) \mid L(x_i, D_c)]$



characterizinginductivesystems

by theirinductive biasallowsmodellingthembytheirequivalentdeductivesystems. This provides away to compare inductive systems according to their policies for generalizing beyond the observed training data

DECISIONTREELEARNING



DECISIONTREEREPRESENTATION



FIGURE: Α decision for tree theconcept PlayTennis.Anexam pleisclassified by sortingit through the tree totheappropriateleaf node, then returningtheclassific ationassociated with thisleaf



- Decision trees classify instances by sorting them down the tree from the root tosome leaf node, which provides the classification of the instance.
- Each node in the tree specifies a test of some attribute of the instance, and eachbranch descending from that node corresponds to one of the possible values forthisattribute.
- Aninstanceisclassifiedbystartingattherootnodeofthetree,testingtheattribute specified by this node, then moving down the tree branch correspondingtothevalueoftheattributeinthegivenexample.Thisprocessisthenrepea tedforthe subtree rootedat the new node.



- Decisiontreesrepresentadisjunctionofconjunctionsofconstraintsontheattrib utevalues of instances.
- Eachpathfromthetreeroottoaleafcorrespondstoaconjunctionofattributetests,an d the treeitself toa disjunction of these conjunctions

Forexample,

The decision treeshown in a bove figure corresponds to the expression (Outlook = Sunny AHumidity = Normal) (Outlook=Overcast) (Outlook=RainAWind=Weak)





APPROPRIATEPROBLEMSFOR DECISIONTREELEARNING

Decisiontreelearningisgenerallybestsuitedtoproblemswiththefollowingcharacteristics:

- *1. Instancesarerepresentedbyattribute-valuepairs*—Instancesaredescribedbyafixed set of attributes and their values
- 2. The target function has discrete output values The decision tree assigns aBoolean classification (e.g., yes or no) to each example. Decision tree methodseasilyextendtolearningfunctions withmore than two possible output values.
- 3. Disjunctivedescriptionsmayberequired



- **4.** *Thetrainingdatamaycontainerrors*—Decisiontreelearningmethodsarerobust to errors, both errors in classifications of the training examples and errorsinthe attribute values that describe these examples.
- 5. *Thetrainingdatamaycontainmissingattributevalues*–Decisiontreemethodscanbe used evenwhen some training exampleshaveunknown values
- Decision tree learning has been applied to problems such as learning to classify*medical patients by their disease, equipment malfunctions by their cause,* and*loan applicants by their likelihood of defaulting on payments.*
- Suchproblems, inwhich the task is to classify examples into one of a discrete set of possible categories, are often referred to as *classification problems*.

THEBASICDECISIONTREE LEARNING

• Mostalgorithmsthathavebeendevelopedforlearningdecisiontreesarevariations on a core algorithm that employs a top-down, greedy search through thespaceofpossibledecisiontrees. This approach is exemplified by the ID3 algorithm and its successor C4.5



WhatistheID3algorithm?

- ID3standsforIterativeDichotomiser3
- ID3isaprecursortotheC4.5Algorithm.
- TheID3algorithmwasinventedbyRossQuinlanin1975
- Usedtogenerateadecisiontreefromagivendatasetbyemployingatopdown,greedysearch, to test each attribute at every nodeof thetree.
- Theresultingtreeisusedtoclassifyfuturesamples.



ID3algorithm

ID3(Examples, Target_attribute, Attributes)

Examples are the training examples. Target_attribute is the attribute whose value is tobe predicted by the tree. Attributes is a list of other attributes that may be tested by the learned decision tree.Returns a correctly classifies the given Examples.

- CreateaRootnodeforthetree
- IfallExamplesarepositive,Returnthesingle-nodetreeRoot,withlabel=+
- IfallExamplesarenegative,Returnthesingle-nodetreeRoot,withlabel=-
- IfAttributesisempty,ReturnthesinglenodetreeRoot,withlabel=mostcommonvalueofTarget_attributeinExamples



- OtherwiseBegin

 - The decision attribute for Root \leftarrow A
 - Foreachpossiblevalue, *v*_i, of A,
 - Addanewtreebranchbelow*Root*,correspondingtothetestA=*v*_i
 - Let*Examplesvi*, bethesubsetofExamplesthathavevaluev_iforA
 - If *Examplesvi*, is empty
 - Thenbelowthisnewbranchaddaleafnodewithlabel=mostcommonvalueofTarget_attri buteinExamples
 - Elsebelowthisnewbranchaddthesubtree ID3(*Examples_{vi}*,Targe_tattribute,Attributes-{A}))
- End
- ReturnRoot

 * The best attribute is the one with highest information gain



WhichAttributeIstheBestClassifier?

- ThecentralchoiceintheID3algorithmisselectingwhichattributetotestat eachnode in thetree.
- Astatisticalpropertycalled*informationgain*thatmeasureshowwellagivenattributese paratesthetraining examples according to their target classification.
- ID3uses*informationgain*measuretoselectamongthecandidateattributesateach step while growingthe tree.

ENTROPYMEASURESHOMOGENEITYOFEXAMPLES



- Todefineinformationgain, we begin by defining a measure called entropy. *Entropymeasures the impurity of a collection of examples*.
- GivenacollectionS, containing positive and negative examples of some target conce pt, the entropy of Srelative to this Boolean classification is

$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Where,

p+istheproportionofpositiveexamplesinS *p*-istheproportionofnegativeexamplesinS.

Example:Entropy



• Suppose S is a collection of 14 examples of some boolean concept, including 9positive and 5 negative examples. Then the entropy of S relative to this booleanclassificationis

$$Entropy([9+, 5-]) = -(9/14) \log_2(9/14) - (5/14) \log_2(5/14)$$
$$= 0.940$$



- Theentropyis0ifallmembersofSbelongtothesameclass
- Theentropyis1whenthecollectioncontainsanequalnumberofpositiveandnegativ e examples
- If the collection contains unequal numbers of positive and negative examples, the entro pyis between 0 and 1





FIGURE The entropy function relative to a boolean classification, as the proportion, p_{\oplus} , of positive examples varies between 0 and 1.

INFORMATIONGAINMEASURESTHEEXPECTEDREDU

- *Informationgain*, is the expected reduction in entropy caused by partitioning the examples according to this attribute.
- Theinformationgain,Gain(S,A)ofanattributeA,relativetoacollectionofexamplesS, is defined as

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$



Example:Informationgain

Let, Values(Wind)={Weak, Strong} S = [9+,5-] $S_{Weak} = [6+,2-]$ $S_{Strong} = [3+,3-]$

Informationgainofattribute*Wind*:

 $Gain(S,Wind) = Entropy(S) - \frac{8}{14}Entropy(S_{Weak}) - \frac{6}{14}Entropy(S_{Strong})$ = 0.94-(8/14)*0.811-(6/14)*1.00 = 0.048



AnIllustrativeExample

- ToillustratetheoperationofID3,considerthelearningtask representedbythetrainingexamples of below table.
- Herethetargetattribute*PlayTennis*,which canhave values*yes* or*no*fordifferentdays.
- Consider the first step through the algorithm, in which the top most node of the decision tree is created.



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No
Theinformationgainvaluesforallfourattributesare



- Gain(S,Outlook) = 0.246
- Gain(S,Humidity) = 0.151
- Gain(S,Wind) =0.048
- Gain(S,Temperature) =0.029
- Accordingtotheinformationgainmeasure,the*Outlook*attributeprovidesthebest prediction of the target attribute, *PlayTennis*, over the training examples.Therefore,*Outlook*isselectedasthedecisionattributefortherootnode, andbranchesarecreatedbelowtherootforeachofitspossiblevaluesi.e.,



men annoure snoura de restea nere:

DeepakD,Asst.Prof.,Dept.ofCSE,CanaraEngg.College

 $S_{sunnv} = \{D1, D2, D8, D9, D11\}$

 $Gain (S_{sunny}, Humidity) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$ $Gain (S_{sunny}, Temperature) = .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$ $Gain (S_{sunny}, Wind) = .970 - (2/5) 1.0 - (3/5) .918 = .019$

 $S_{Rain} = \{D4, D5, D6, D10, D14\}$

 $Gain(S_{Rain}, Humidity) = 0.970 - (2/5)1.0 - (3/5)0.917 = 0.019$ $Gain(S_{Rain}, Temperature) = 0.970 - (0/5)0.0 - (3/5)0.918 - (2/5)1.0 = 0.019$ $Gain(S_{Rain}, Wind) = 0.970 - (3/5)0.0 - (2/5)0.0 = 0.970$





HYPOTHESISSPACESEARCHINDECISIONTREEL EARNING



- ID3 can be characterized as searching a space of hypotheses for one that fits thetraining examples.
- ThehypothesisspacesearchedbyID3isthesetofpossibledecisiontrees.
- ID3 performs a *simple-to complex, hill-climbing search* through this hypothesisspace, beginning with the emptytree, then considering progressively more ela borate hypotheses insearch of a decision tree that correctly classifies the training data



Figure:

• HypothesisspacesearchbyID3.



• ID3searchesthroughthespaceofpossible decision trees from simplest toincreasinglycomplex,guidedbytheinfo rmationgain heuristic



1. ID3's hypothesis space of all deision trees is a *complete* space of finite discrete-valued functions, relative to the available attributes. Because every finite discrete-valuedfunction can be represented by some decision tree

• ID3avoidsoneofthemajorrisksofmethodsthats*earchincompletehypothesisspaces* :thatthehypothesisspacemightnotcontainthetargetfunction.

NRCM your roots to success...

2. ID3maintains*onlyasinglecurrenthypothesis*asitsearchesthroughthespaceof decision trees.

Forexample, with the earlier version space candidate elimination method, which maint ains the set of *all* hypotheses consistent with the available training examples.

Bydeterminingonlyasinglehypothesis,ID3losesthecapabilitiesthatfollowfromexplicitly the total all consistent hypotheses.

For example, it does not have the ability to determine how many alternativedecisiontreesareconsistent with the available training data, or to pose new instance queries that

optimally resolve among the secompeting hypotheses



3. ID3 in its pure form performs *no backtracking in its search*. Once it selectoria anattribute to test at a particular level in the tree, it never backtracks to reconsider thischoice.

• In the case of **ID3**, a locally optimal solution corresponds to the decision tree itselects along the single search path it explores. However, this locally optimal solution may be less desirable than trees that would have been encountered along adifferentbranch of the search.

4. ID*3usesalltrainingexamplesateachstep* in the search to make statistically based decis ions regarding how to refine ts current hypothesis.

- Oneadvantageofusingstatisticalpropertiesofalltheexamplesisthattheresultingsearc h ismuch *less sensitiveto errors* inindividualtraining examples.
- **ID3**canbeeasilyextendedtohandlenoisytrainingdataby modifyingitsterminationcriteriontoaccepthypotheses thatimperfectlyfitthetrainingdata.

INDUCTIVEBIASINDECISIONTREELEARNING



Inductive bias is the set of assumptions that, together with the training data, deductively justify the classifications assigned by the learner to future instances

Givenacollectionoftrainingexamples, there are typically many decision trees consistent with these examples. Which of these decision trees does ID3 choose?

ID3searchstrategy

- (a) selects infavour of shorter trees overlonger ones
- (b) selectstreesthatplacetheattributes with highestinformation gain closest to the root.

$\label{eq:approximate} Approximate inductive bias of ID3: {\it Shorter trees are preferred overlarger trees} and {\it Shorter trees} and {\it Shorter trees} are preferred overlarger trees} are preferred overlarger trees} and {\it Shorter trees} are preferred overlarger trees} ar$



- Consideranalgorithmthatbeginswiththeemptytreeandsearches*breadthfirst* throughprogressivelymorecomplextrees.
- Firstconsideringalltreesofdepth1,thenalltreesofdepth2,etc.
- Once itfinds a decision tree consistent with thetraining data, it returns thesmallestconsistenttreeatthatsearchdepth(e.g.,thetreewiththefewestnodes).
- Letuscallthisbreadth-firstsearchalgorithmBFS-ID3.
- BFS-ID3findsashortestdecisiontreeandthusexhibitsthebias"shortertrees arepreferredover longer trees.



- ID3 can be viewed as an efficient approximation to BFS-ID3, using a greedyheuristic search to attempt to find the shortest tree without conducting the entirebreadth-firstsearchthrough thehypothesis space.
- Because ID3 uses the information gain heuristic and a hill climbing strategy, itexhibits more complexbiasthan BFS-ID3.
- In particular, it does not always find the shortest consistent tree, and it is biased to favourtrees that place attributes with high information gain closest to the root.

RestrictionBiasesandPreferenceBiases



DifferencebetweenthetypesofinductivebiasexhibitedbyID3andbytheCANDIDATE-ELIMINATIONAlgorithm.

<u>ID3</u>

- ID3searchesacompletehypothesisspace
- Itsearchesincompletelythroughthisspace,fromsimpletocomplexhypotheses,untilitstermin ation condition ismet
- Its inductive bias is solely a consequence of the ordering of hypotheses by its search strategy. Its hypothesis space introduces no additional bias

CANDIDATE-ELIMINATIONAlgorithm

- The version space CANDIDATE-ELIMINATION Algorithms earches an incomplete hypothesis space
- Itsearchesthisspacecompletely, finding every hypothesis consistent with the training data.
- Its inductive bias is solely a consequence of the expressive power of its hypothesisrepresentation. Its search strategy introduces no additional bias

RestrictionBiasesandPreferenceBiases



- The inductive bias of ID3 is a *preference* for certain hypotheses over others (e.g., preference for shorter hypotheses over larger hypotheses), with no hard restrictiononthehypothesesthatcanbeeventuallyenumerated. This form of bias is called a *preference bias* or *a search bias*.
- The bias of the CANDIDATE ELIMINATION algorithm is in the form of a*categorical* restriction on the set of hypotheses considered. This form of bias istypicallycalled a*restriction bias* or *alanguage bias*.

Whichtypeofinductivebiasispreferredinordertogeneralizebeyondthetrainingdata.apreferenc e bistrom

- A preference bias is more desirable than a restriction bias, because it allows thelearner to work within a complete hypothesis space that is assured to contain theunknown target function.
- In contrast, a restriction bias that strictly limits the set of potential hypotheses isgenerally less desirable, because it introduces the possibility of excluding the unknown target function altogether.

Occam'srazor



Occam's razor: is the problem-solving principle that the simplest solution tends to bethe right one. When presented with competing hypotheses to solve a problem, oneshouldselect the solution with the fewest assumptions.

Occam'srazor:"Preferthesimplesthypothesisthatfitsthedata".

WhyPreferShortHypotheses?



Argumentinfavour:

Fewershorthypothesesthanlongones:

- Shorthypothesesfitsthetrainingdatawhicharelesslikelytobecoincident
- Longerhypothesesfitsthetrainingdatamightbecoincident.

Manycomplexhypothesesthatfitthecurrenttrainingdatabutfailtogeneralizecorrec tly tosubsequent data.

Argumentopposed:



- There are few small trees, and our priori chance of finding one consistent with anarbitrary set of data is therefore small. The difficulty here is that there are verymanysmallsetsofhypothesesthatonecandefine*butunderstoodbyfewerlearner*.
- The size of a hypothesis is determined by the representation used *internally* by thelearner.Occam'srazorwillproduce*twodifferenthypothesesfromthesametrainingex ampleswhenitisappliedbytwolearners*,bothjustifyingtheircontradictory conclusions by Occam's razor. On this basis we might be tempted torejectOccam's razor altogether.



ISSUESINDECISIONTREELEARNING

1. AvoidingOverfittingtheData

Reduced error pruningRulepost-pruning

- 2. IncorporatingContinuous-ValuedAttributes
- 3. AlternativeMeasuresforSelectingAttributes
- 4. HandlingTrainingExampleswithMissingAttributeValues
- 5. HandlingAttributeswithDifferingCosts

1. AvoidingOverfittingtheData



- The ID3 algorithm grows each branch of the tree just deeply enough to perfectlyclassify the training examples but it can lead to difficulties when there is noise inthe data, or when the number of training examples is too small to produce are presentative sample of the true target function. This algorithm can produce trees that *overfit* the training examples.
- **Definition-Overfit:** Given a hypothesis space H, a hypothesis $h \in$ Hissaid to overfit the training data if there exists some alternative hypothesis $h' \in$ H, such that h has smaller error than h' over the training examples, but h' has a smaller error than h over the entire distribution of instances.



• Thebelowfigureillustratestheimpactofoverfittinginatypicalapplicationofdecisiontreelearnin g.



• The horizontalaxis of this

plotindicates the total number of nodes in the decision tree, as the tree is being constructed. The **vertical axis** indicates the accuracy of predictions made by the tree.

- The **solidline** shows the accuracy of the decision tree over the training examples. The **broken line** shows accuracy meas ure dover an independent set of test example
- The accuracy of the training examples increases monotonically as the tree is grown. The accuracy measured over the independent test examples first increases, then decreases.



Howcanitbepossiblefortreehtofitthetrainingexamplesbetterthanh',butforittoperformmorepoorl y over subsequentexamples?

- 1. Overfittingcanoccurwhenthetrainingexamplescontainrandomerrorsornoise
- $2.\ When small numbers of examples are associated with leaf nodes.$

NoisyTrainingExample

Example15:<Sunny,Hot,Normal,Strong,->

- Exampleisnoisybecausethecorrectlabelis+
- Previouslyconstructedtreemisclassifiesit





Criterionusedtodeterminethecorrectfinaltreesize

- Useaseparatesetofexamples, distinctfromthetrainingexamples, to evaluate the utility of post-pruning nodes from the tree
- Usealltheavailabledatafortraining, butapplyastatisticaltest to estimate whether expanding pruning) a particular node is likely to produce an improvement beyond thetrainingset
- Usemeasureofthecomplexityforencodingthetrainingexamplesandthedecisiontree,h alting growth of the tree when this encoding size is minimized. This approach is theMinimumDescription Length

called

MDL-

Minimize:size(tree)+size(misclassifications(tree))

Reduced-ErrorPruning



- *Reduced-error pruning*, is to consider each of the decision nodes in the tree to becandidates for pruning
- *Pruning* a decision node consists of removing the subtree rooted at that node, making it a leaf node, and assigning it the most common classification of thetraining examples affiliated with that node
- Nodes are removed only if the resulting pruned tree performs no worse thantheoriginalover the validation set.
- Reduced error pruning has the effect that any leaf node added due to coincidentalregularities in the training set is likely to be pruned because these same coincidences are unlikely to occur in the validation set

Theimpactofreduced-errorpruningontheaccuracyofthedecisiontreeisillustratedinbelowfigure





- The additional line in figure shows *accuracy over the test examples* as the tree is pruned. Whenpruning begins, the tree is at its maximum size and lowest accuracy over the test set. As pruningproceeds,the numberofnodes is reduced and accuracy overthetest set increases.
- The available data has been split into three subsets: the training examples, the validation examplesused for pruning the tree, and a set of test examples used to provide an unbiased estimate of accuracy overfuture unseen examples. The plotshows accuracy over the training and test sets.

Summary



In general, the problem of estimating confidence intervals isapproachedbyidentifyingtheparametertobeestimated($error_bh$) and anestimator($error_sh$)forthisquantity.

Because the estimator is a random variable it can be characterised by the probability distribution that governs its value.

upbiasedestima₂tor,theobservedvalueoftheestimatorislikelyto vary fromone Experimentto another.





The variance of the distribution governing the estimate is likely to orcharacterises how widely this

Confidence intervals can then be calculated by determining the interval that contains the desired probability mass under this distribution.

'Acause of estimation error is the variance in the estimate. Even with an



Prodigy and Explanation Based Learning

Prodigy defines a set of target concepts to learn, e.g., which operator given the current state takes you to the goal state?

An example of a rule learned by Prodigy in the blockstacking problem is:

IF One subgoal to be solved is On(x,y) ANDOne subgoal to be solved is On(y,z)THEN Solve the subgoal On(y,z) before On(x,y)



Prodigy and Explanation Based Learning

The rationale behind the rule is that it would avoid a conflict when stacking blocks.

Prodigy learns by first encountering a conflict, then explaining the reason for the conflict and creating a rule like the one above.

Experiments show an improvement in efficiency by a factor of two to four.



Problems with EBL

- ✓ The number of control rules that must be learned is very large.
- ✓ If the control rules are many, much time will be spent looking for the best rule. Utility analysis is used to determine what rules to keep and what rules to forget.

Prodigy:

328 possible rules \longrightarrow 69 pass test \longrightarrow 19 were retained



Problems with EBL

✓ Another problem with EBL is that it is sometimes intractable to create an explanation for the target concept.

For example, in chess, learning a concept like: "states for which operator A leads to a solution" The search here grows exponentially.



Summary

- Different from inductive learning, analytical learning looks for a hypothesis that fit the background knowledge and covers the training examples.
- Explanation based learning is one kind of analytical learning that divides into three steps:
 - a. Explain the target value for the current example
 - b. Analyze the explanation (generalize)
 - c. Refine the hypothesis



Summary

- Prolog-EBG constructs intermediate features after analyzing examples.
- Explanation based learning can be used to find search control rules.
- In all cases we depend on a perfect domain theory.



부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr

Artificial Intelligence Laboratory

Machine Learning Chapter 12. Combining Inductive and Analytical Learning

Tom M. Mitchell



Inductive and Analytical Learning

Inductive learning

- Hypothesis fits data
- Statistical inference
- Requires little prior knowledge
- Syntactic inductive bias

Analytical learning

- Hypothesis fits domain the
- Deductive inference
- Learns from scarce data
- Bias is domain theory



부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr

Artificial Intelligence Laboratory

What We Would Like

Inductive learning

Analytical learning

Plentiful data No prior knowledge Perfect prior knowledge Scarce data

General purpose learning method:

- No domain theory \rightarrow learn as well as inductive methods
- Perfect domain theory \rightarrow learn as well as Prolog-EBG
- Accomodate arbitrary and unknown errors in domain theory
- Accomodate arbitrary and unknown errors in training data


부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr

Artificial Intelligence Laboratory

Domain theory:

Cup ← Stable, Liftable, Open Vessel Stable ← BottomIsFlat Liftable ← Graspable, Light Graspable ← HasHandle Open Vessel ← HasConcavity, ConcavityPointsUp

Training examples:







KBANN

- KBANN (data D, domain theory B)
 1. Create a feedforward network h equivalent to B
 - 2. Use BACKPROP to tune h to t D









Creating Network Equivalent to Domain Theory

Create one unit per horn clause rule (i.e., an AND unit)

- Connect unit inputs to corresponding clause antecedents
- For each non-negated antecedent, corresponding input weight w ← W, where W is some constant
- For each negated antecedent, input weight $w \leftarrow -W$
- Threshold weight w₀ ← -(n-.5)W, where n is number of non-negated antecedents

Finally, add many additional connections with near-zero weights

Liftable \leftarrow *Graspable*, \neg *Heavy*



borame.cs.pusan.ac.kr

Result of refining the network







Artificial Intelligence Laboratory

KBANN Results

Classifying promoter regions in DNA leave one out testing:

- Backpropagation : error rate 8/106
- **KBANN: 4/106**

Similar improvements on other classification, control tasks.



🙆 부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr

Artificial Intelligence Laboratory

Hypothesis space search in KBANN





11



Key idea:

- Previously learned approximate domain theory
- Domain theory represented by collection of neural networks
- Learn target function as another neural network



부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr

Artificial Intelligence Laboratory





Modified Objective for Gradient Descent

$$E = \sum_{i} \left[(f(x_i) - \hat{f}(x_i))^2 + \mu_i \sum_{j} \left(\frac{\partial A(x)}{\partial x^j} - \frac{\partial \hat{f}(x)}{\partial x^j} \right)_{(x=x_i)}^2 \right]$$

where

$$\mu_i \equiv 1 - \frac{|A(x_i) - f(x_i)|}{c}$$

- f(x) is target function
- $\hat{f}(x)$ is neural net approximation to f(x)
- A(x) is domain theory approximation to f(x)





Artificial Intelligence Laboratory







· 부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr

Artificial Intelligence Laboratory



Hypothesis Space Search in EBNN

Hypothesis Space



부산대학교 인공지능 연구실 borame.cs.pusan.ac.kr



Search in FOCL





THANK YOU