Lecture Slides for

INTRODUCTION TO Machine Learning 2nd Edition



ETHEM ALPAYDIN, modified by Leonardo Bobadilla and some parts from http://www.cs.tau.ac.il/~apartzin/MachineLearning/ and

www.cs.princeton.edu/courses/archive/fall01/cs302 /notes/11.../EM.ppt © The MIT Press, 2010 alpaydin@boun.edu.tr

http://www.cmpe.boun.edu.tr/~ethem/i2m

Outline

Previous class Ch 8: Decision Trees This class: Ch 9: Decision Trees

CHAPTER 9: Decision Trees

Slides from: Blaž Zupan and Ivan Bratko magix.fri.uni-lj.si/predavanja/uisp



Decision Tree



Classification Trees

- What is the good split function?
- Use Impurity measure
- Assume N_m training samples reach node m
- N_m^i of N_m belong to class C_i , with $\sum_i N_m^i = N_m$.

•
$$\hat{P}(C_i|\mathbf{x},m) \equiv p_m^i = \frac{N_m^i}{N_m}$$

• Node m is pure if for all classes either 0 or 1

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Entropy

- Measure amount of uncertainty on a scale from 0 to 1
- Example: 2 events
- If p1=p2=0.5, entropy is 1 which is maximum uncertainty
- If p₂=1=1-p₁, entropy is 0, which is no uncertainty (Eq 9.3)

$$\mathcal{I}_m = -\sum_{i=1}^{\kappa} p_m^i \log_2 p_m^i$$

Entropy



Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Best Split

- Node is impure, need to split more
- Have several split criteria (coordinates), have to choose optimal
- Minimize impurity (uncertainty) after split
- Stop when impurity is small enough
 - Zero stop impurity=>complex tree with large variance
 - Larger stop impurity=>small tress but Based on E Alpavdin 2004 Introduction to Machine Learning © The MIT Press (V1.1) large bias

Best Split

- Impurity after split: N_{mj} of N_m take branch *j*.
- Nⁱ_{mj} belong to C_i

$$\hat{P}(C_{i} | \mathbf{x}, m, j) \equiv p_{mj}^{i} = \frac{N_{mi}^{i}}{N_{mj}}$$

• (Eq 9.8)
$$\mathbf{I'}_{m} = -\sum_{j=1}^{n} \frac{N_{mj}}{N_{m}} \sum_{i=1}^{K} p_{mj}^{i} \log_{2} p_{mj}^{i}$$

- Find the variable and split that min impurity
 - among all variables
 - split positions for numeric variables

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)





I(yellow) = 0.0 bits

 $Gain(Color) = I - I_{res}(Color) = 0.940 - 0.694 = 0.246 \ bits$



Regression Trees

- Value not a label in a leaf nodes
- Need other impurity measure
- Use Average Error

$$b_{m}(x) = \begin{cases} 1 & \text{if } x \in X_{m} : x \text{ reaches node } m \\ 0 & \text{otherwise} \end{cases}$$
$$E_{m} = \frac{1}{N_{m}} \sum_{t} (r^{t} - g_{m})^{2} b_{m}(x^{t}) \qquad g_{m} = \frac{\sum_{t} b_{m}(x^{t}) r^{t}}{\sum_{t} b_{m}(x^{t})}$$

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Regression Trees

• After splitting:

 $b_{mj}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \mathsf{X}_{mj} : \mathbf{x} \text{ reaches node } m \text{ and branch } j \\ 0 & \text{otherwise} \end{cases}$

$$E'_{m} = \frac{1}{N_{m}} \sum_{j} \sum_{t} \left(r^{t} - g_{mj} \right)^{2} b_{mj} \left(\mathbf{x}^{t} \right) \qquad g_{mj} = \frac{\sum_{t} b_{mj} \left(\mathbf{x}^{t} \right) r^{t}}{\sum_{t} b_{mj} \left(\mathbf{x}^{t} \right)}$$

Lecture Notes for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)



Pruning Trees

- Number of data instances reach a node is small
 - Less then 5% of training data
 - Don't want to split further regardless of impurity
- Remove subtrees for better generalization
 - Prepruning: Early stopping
 - Postpruning: Grow the whole tree then prune subtrees
 - Set aside pruning set
 - Make sure pruning does not significantly increase error

Decision Trees and Feature Extraction

- Univariate Tree uses only certain variable
- Some variables might not get used
- Features closer to the root have greater importance

Interpretability

- Conditions that are simple to understand
- Path from the root =>one conjunction of test
- All paths can be defined using set of IF_THEN rules
 - Form a rule base
- Percentage of training data covered by the rule
 - Rule support
- Tool for a Knowledge Extraction
- Can be verified by experts Nachine Learning © The MIT Press (V1.1)

Rule Extraction from Trees



- R1: IF (age>38.5) AND (years-in-job>2.5) THEN y = 0.8
- R2: IF (age>38.5) AND (years-in-job \leq 2.5) THEN y = 0.6
- R3: IF (age \leq 38.5) AND (job-type='A') THEN y = 0.4
- R4: IF (age \leq 38.5) AND (job-type='B') THEN y = 0.3
- R5: IF (age \leq 38.5) AND (job-type='C') THEN y = 0.2

Learning Rules

- Rule induction is similar to tree induction but
 - tree induction is breadth-first,
 - rule induction is depth-first; one rule at a time
- Rule set contains rules; rules are conjunctions of terms
 - A rule covers an example if all terms of the rule evaluate to true for the example.
 - A rule is said to cover an example if the example satisfies all the conditions of the rule.
- Sequential covering:
 - Generate rules one at a time until all positive examples are covered

Rule-Based Classifier

- Classify records by using a collection of "if... then..." rules
- Rule: (Condition) $\rightarrow y$
 - where
 - Condition is a conjunctions of attributes
 - y is the class label
 - Examples of classification rules:
 - (Blood Type=Warm) \land (Lay Eggs=Yes) \rightarrow Birds
 - (Taxable Income < 50K) \land (Refund=Yes) \rightarrow Evade=No

Rule-based Classifier (Example)

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
human	warm	yes	no	no	mammals
python	cold	no	no	no	reptiles
salmon	cold	no	no	yes	fishes
whale	warm	yes	no	yes	mammals
frog	cold	no	no	sometimes	amphibians
komodo	cold	no	no	no	reptiles
bat	warm	yes	yes	no	mammals
pigeon	warm	no	yes	no	birds
cat	warm	yes	no	no	mammals
leopard shark	cold	yes	no	yes	fishes
turtle	cold	no	no	sometimes	reptiles
penguin	warm	no	no	sometimes	birds
porcupine	warm	yes	no	no	mammals
eel	cold	no	no	yes	fishes
salamander	cold	no	no	sometimes	amphibians
gila monster	cold	no	no	no	reptiles
platypus	warm	no	no	no	mammals
owl	warm	no	yes	no	birds
dolphin	warm	yes	no	yes	mammals
eagle	warm	no	yes	no	birds

R1: (Give Birth = no) \land (Can Fly = yes) \rightarrow Birds

- R2: (Give Birth = no) \land (Live in Water = yes) \rightarrow Fishes
- R3: (Give Birth = yes) \land (Blood Type = warm) \rightarrow Mammals

R4: (Give Birth = no) \land (Can Fly = no) \rightarrow Reptiles

R5: (Live in Water = sometimes) \rightarrow Amphibians

Application of Rule-Based Classifier

 A rule is said to cover an example if the example satisfies all the conditions of the rule.

R1: (Give Birth = no) \land (Can Fly = yes) \rightarrow Birds

R2: (Give Birth = no) \land (Live in Water = yes) \rightarrow Fishes

R3: (Give Birth = yes) \land (Blood Type = warm) \rightarrow Mammals

R4: (Give Birth = no) \land (Can Fly = no) \rightarrow Reptiles

R5: (Live in Water = sometimes) \rightarrow Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
hawk	warm	no	yes	no	?
grizzly bear	warm	yes	no	no	?

The rule R1 covers a hawk => Bird

The rule R3 covers the grizzly bear => Mammal

From: Luke Huan 2006

Example of Sequential Covering



From: Luke Huan 2006

Example of Sequential Covering...



From: Luke Huan 2006

Ripper Algorithm

- There are two kinds of loop in Ripper algorithm (Cohen, 1995):
 - Outer loop : adding one rule at a time to the rule base
 - Inner loop : adding one condition at a time to the current rule
 - Conditions are added to the rule to maximize an information gain measure.
 - Conditions are added to the rule until it covers no negative example.

Ripper Algorithm

- In Ripper, conditions are added to the rule to
 - Maximize an information gain measure

$$Gain(R',R) = s \cdot (\log_2 \frac{N'}{N'} - \log_2 \frac{N}{N})$$

- R : the original rule
- *R*': the candidate rule after adding a condition
- N(N'): the number of instances that are covered by R(R')
- N_+ (N'_+): the number of true positives in R (R')
- *s* : the number of true positives in *R* and *R*' (after adding the condition)
- Until it covers no negative example

Stopping Criterion and Rule Pruning

- Stopping criterion
 - Compute the gain
 - If gain is not significant, discard the new rule
- Rule Pruning
 - Similar to post-pruning of decision trees
 - Reduced Error Pruning:
 - Remove one of the conjuncts in the rule
 - Compare error rate on validation set before and after pruning
 - If error improves, prune the conjunct

Summary of Direct Method

Grow a single rule

- Remove Instances from rule
- Prune the rule (if necessary)
- Add rule to Current Rule Set



Direct Method: RIPPER

- For 2-class problem, choose one of the classes as positive class, and the other as negative class
 - Learn rules for positive class
 - Negative class will be default class
- For multi-class problem
 - Order the classes according to increasing class prevalence (fraction of instances that belong to a particular class)
 - Learn the rule set for smallest class first, treat the rest as negative class
 - Repeat with next smallest class as positive class

Multivariate Trees



Based on for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Likelihood- vs. Discriminant-based Classification

 Likelihood-based: Assume a model for p(x|C_i), use Bayes' rule to calculate P(C_i|x)

 $g_i(\boldsymbol{x}) = \log P(C_i|\boldsymbol{x})$

- Discriminant-based: Assume a model for g_i(**x**|Φ_i); no density estimation
 - Estimating the boundaries is enough; no need to accurately estimate the densities inside the boundaries
 - Learning is the optimization of the model parameters Φ_i to maximize the classification accuracy on a given labeled training set.

Linear Discriminant

• Linear discriminant:

$$g_i(\mathbf{x} | \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0} = \sum_{j=1}^d \mathbf{w}_{ij} \mathbf{x}_j + \mathbf{w}_{i0}$$

Advantages:

- Simple: O(d) space/computation
- Easy to understand:
 - The final output is weighted sum of several factors.

• Accurate:

- When p(x|C_i) are Gaussian with a shared covariance matrix, the optimal discriminant is linear.
- The linear discriminant can be used even when this assumption does not hold.

Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e $\ensuremath{\mathbb{C}}$ The MIT Press (V1.0)

Generalized Linear Model

• Quadratic discriminant:

 $g_i(\mathbf{x} | \mathbf{W}_i, \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$

• Higher-order (product) terms:

$$Z_1 = X_1, Z_2 = X_2, Z_3 = X_1^2, Z_4 = X_2^2, Z_5 = X_1X_2$$

Map from **x** to **z** using nonlinear basis functions and use a linear discriminant in **z**-space

$$g_{i}(\mathbf{x}) = \sum_{j=1}^{k} w_{j} \phi_{ij}(\mathbf{x})$$
$$\phi_{ij}(\mathbf{x}) : \text{ basis functions}$$

Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e $\$ The MIT Press (V1.0)

Two Classes



Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e $\ensuremath{\mathbb{C}}$ The MIT Press (V1.0)



 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \mathbf{w}_0$

 w_0 : determines the location of the hyper-plane with respect to the origin

w : determines its orientation

Lecture Notes for E Alpaydin 2010 Introduction to Machine Learning 2e © The MIT Press (V1.0)

Multiple Classes



 $g_i(\mathbf{x} | \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$ Choose C_i if $g_i(\mathbf{x}) = \max_{j=1}^{\kappa} g_j(\mathbf{x})$

> H_i : the hyperplane separate the examples of C_i from the examples of all other classes

Classes are linearly separable

Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e \odot The MIT Press (V1.0)

Pairwise Separation



Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e $\ensuremath{\mathbb{C}}$ The MIT Press (V1.0)