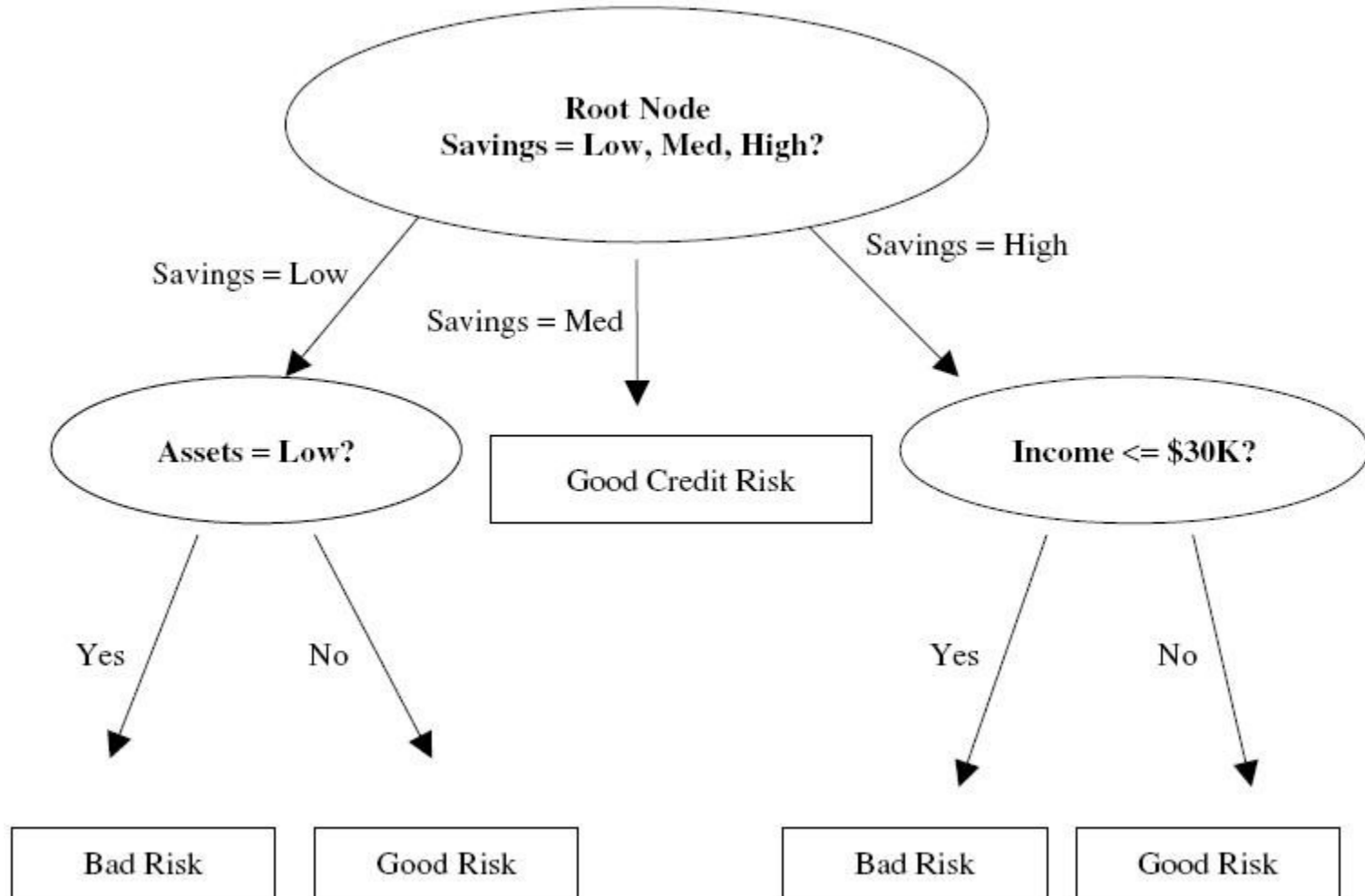


Stabila odlučivanja

Bojan Furlan

УНИВЕРЗИТЕТ У БЕОГРАДУ
ЕЛЕКТРОТЕХНИЧКИ ФАКУЛТЕТ

Stabla odlučivanja - cilj



Stabla odlučivanja – svojstva

- Kreiranje stabla odlučivanja (*Decision Trees* – DT) :
 - Ručno – na osnovu ekspertskog znanja
 - Automatsko – na osnovu trening skupa (ML)
- Zahtevi:
 - Trening skup sa vrednostima ciljnog atributa (nadgledano učenje)
 - Ciljni atribut mora biti diskretan (ili diskretizovan)

Stabla odlučivanja – svojstva

- Pitanje#1: Koji atribut uzeti za grananje?
 - Onaj koji deli početni čvor na što “čistije” potomke
 - U listovima treba da bude zastupljeno što više instanci iste klase
 - Ovo obezbeđuje klasifikaciju sa najvećom pouzdanošću.
 - Npr. za prethodni primer atribut ušteđevina (*savings*) je uzet jer najbolje deli početni trening skup.

Stabla odlučivanja – svojstva

- Pitanje#2: Kada stati sa grananjem?
 - I. Kada sve instance u čvoru pripadaju istoj klasi.
 - II. Kada su sva grananja iscrpljena.
- Algoritmi:
 - Classification and regression trees (CART) algoritam
 - C4.5 algoritam

CART algoritam

- CART stabla su binarna stabla i svaki čvor odlučivanja ima tačno dve grane.
- CART rekurzivno deli početni skup u podskupove sa istim vrednostima ciljnog atributa (iste klase).

$$\Phi(s|t) = 2P_L P_R \sum_{j=1}^{\# \text{ classes}} |P(j|t_L) - P(j|t_R)|$$

- $\Phi(s|t)$ je mera “povoljnosti” grananja za kandidat grananje s čvora t
- t_L i t_R su levi odnosno desni potomak čvora t
- Optimalno grananje je ono sa maksimalnom vrednosti $\Phi(s|t)$ za sva moguća grananja za čvor t .

CART algoritam – svojstva

$$\Phi(s|t) = 2P_L P_R \sum_{j=1}^{\# \text{ classes}} |P(j|t_L) - P(j|t_R)|$$

- Zastupljenost levog i desnog potomka:
 $P_L = |t_L|/|t|$ i $P_R = |t_R|/|t|$
- Zastupljenost klase j u levom i desnom potomku:
 - $P(j|t_L) = |t_L == j|/|t_L|$
 - $P(j|t_R) = |t_R == j|/|t_R|$

CART algoritam – svojstva

- $\Phi(s|t)$ raste kada obe komponente proizvoda rastu $2P_L P_R$ i $\sum_{j=1}^{\# \text{ classes}} |P(j|t_L) - P(j|t_R)|$
- 1. $2P_L P_R$ ima maksimalnu vrednost kada su oba potomka iste veličine (ista zastupljenost): $0.5*0.5 = 0.25$ ili $0.9*0.1 = 0.09$
- 2. $Q(s|t) = \sum_{j=1}^{\# \text{ classes}} |P(j|t_L) - P(j|t_R)|$
 - Maksimalna vrednost je kada sve instance čvora potomka su potpuno uniformne (čiste).
 - Teoretski maksimum za $Q(s|t)$ je k ,
gde je k broj klasa koje može uzimati ciljna promenljiva.

CART Primer

Customer	Savings	Assets	Income (\$1000s)	Credit Risk
1	Medium	High	75	Good
2	Low	Low	50	Bad
3	High	Medium	25	Bad
4	Medium	Medium	50	Good
5	Low	Medium	100	Good
6	High	High	25	Good
7	Low	Low	25	Bad
8	Medium	Medium	75	Good

Trening skup

Candidate Split	Left Child Node, t_L	Right Child Node, t_R
1	<i>Savings = low</i>	<i>Savings</i> \in { <i>medium, high</i> }
2	<i>Savings = medium</i>	<i>Savings</i> \in { <i>low, high</i> }
3	<i>Savings = high</i>	<i>Savings</i> \in { <i>low, medium</i> }
4	<i>Assets = low</i>	<i>Assets</i> \in { <i>medium, high</i> }
5	<i>Assets = medium</i>	<i>Assets</i> \in { <i>low, high</i> }
6	<i>Assets = high</i>	<i>Assets</i> \in { <i>low, medium</i> }
7	<i>Income</i> \leq \$25,000	<i>Income</i> $>$ \$25,000
8	<i>Income</i> \leq \$50,000	<i>Income</i> $>$ \$50,000
9	<i>Income</i> \leq \$75,000	<i>Income</i> $>$ \$75,000

Kandidat grananja za koreni čvor t

CART Primer

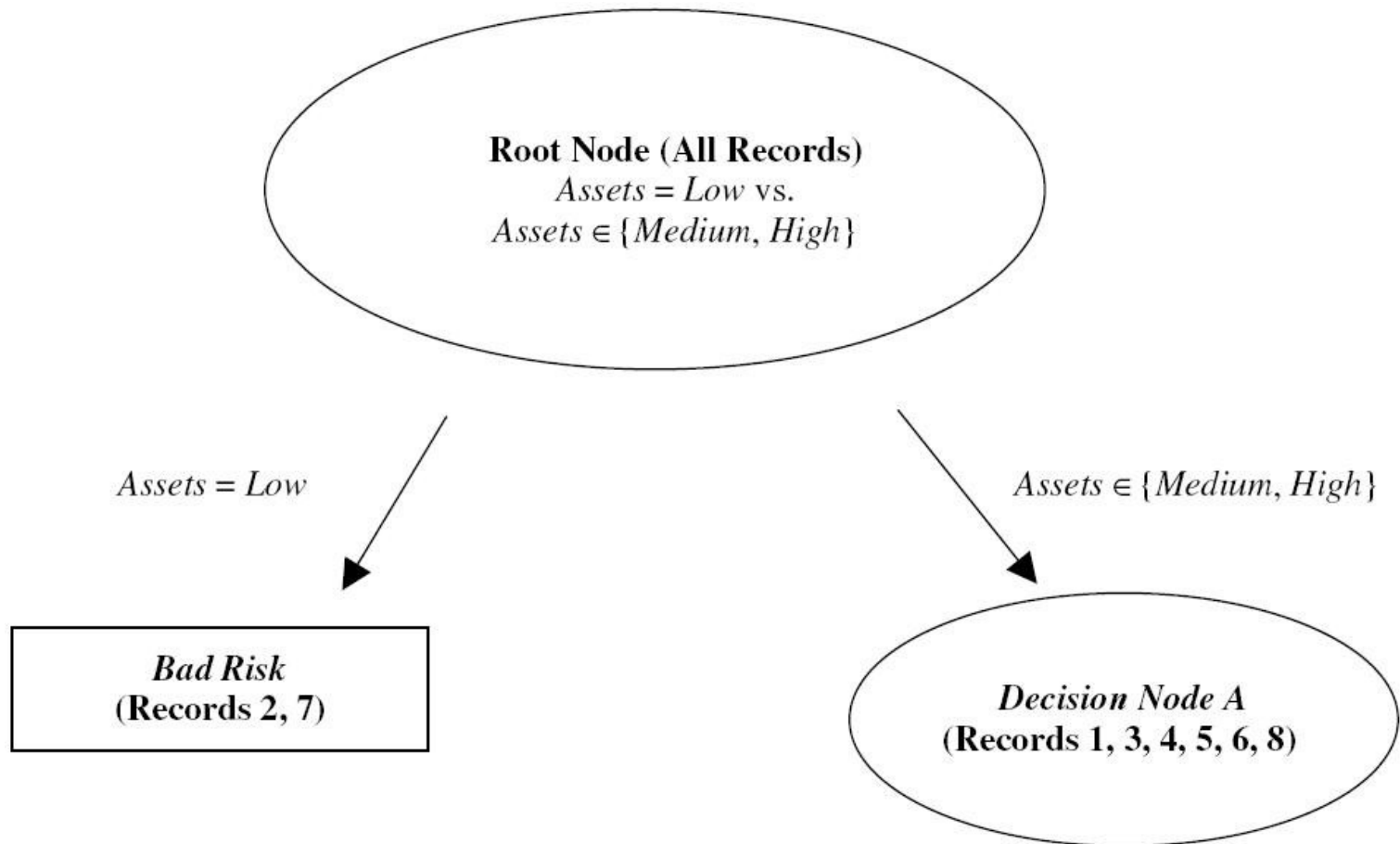
- Grananje 1. -> Savings == low (L-true, R-false)
 - Right:1,3,4,6,8
 - Left:2,5,7
- $P_R=5/8 = 0.625$ $P_L=3/8=0.375$ -> $2 * P_L P_R=15/64=0.46875$
- za $j(\text{klasu}) = \text{Bad}$
 - $P(\text{Bad}|t_R)= 1/5=0.2$
 - $P(\text{Bad}|t_L)= 2/3=0.67$
- za $j(\text{klasu}) = \text{Good}$
 - $P(\text{Good}|t_R)= 4/5 = 0.8$
 - $P(\text{Good}|t_L)= 1/3 = 0.33$
- $Q(s|t)= |0.67-0.2| + |0.8-0.33| = 0.934$

CART Primer

Split	P_L	P_R	$P(j t_L)$	$P(j t_R)$	$2P_L P_R$	$Q(s t)$	$\Phi(s t)$
1	0.375	0.625	G: .333 B: .667	G: .8 B: .2	0.46875	0.934	0.4378
2	0.375	0.625	G: 1 B: 0	G: 0.4 B: 0.6	0.46875	1.2	0.5625
3	0.25	0.75	G: 0.5 B: 0.5	G: 0.667 B: 0.333	0.375	0.334	0.1253
4	0.25	0.75	G: 0 B: 1	G: 0.833 B: 0.167	0.375	1.667	0.6248
5	0.5	0.5	G: 0.75 B: 0.25	G: 0.5 B: 0.5	0.5	0.5	0.25
6	0.25	0.75	G: 1 B: 0	G: 0.5 B: 0.5	0.375	1	0.375
7	0.375	0.625	G: 0.333 B: 0.667	G: 0.8 B: 0.2	0.46875	0.934	0.4378
8	0.625	0.375	G: 0.4 B: 0.6	G: 1 B: 0	0.46875	1.2	0.5625
9	0.875	0.125	G: 0.571 B: 0.429	G: 1 B: 0	0.21875	0.858	0.1877

Vrednosti $\Phi(s|t)$ za svako kandidat grananje za koreni čvor

CART Primer

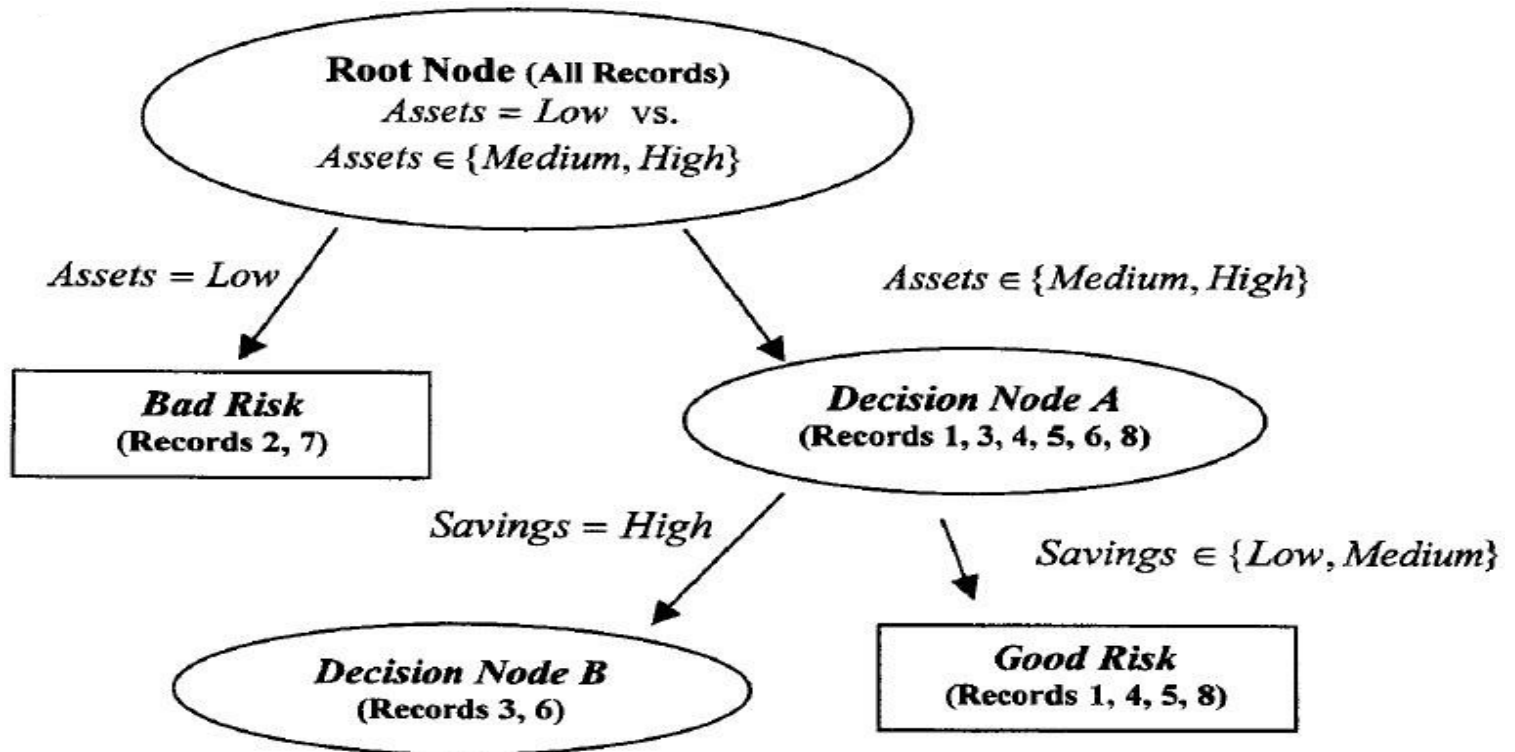


CART Primer

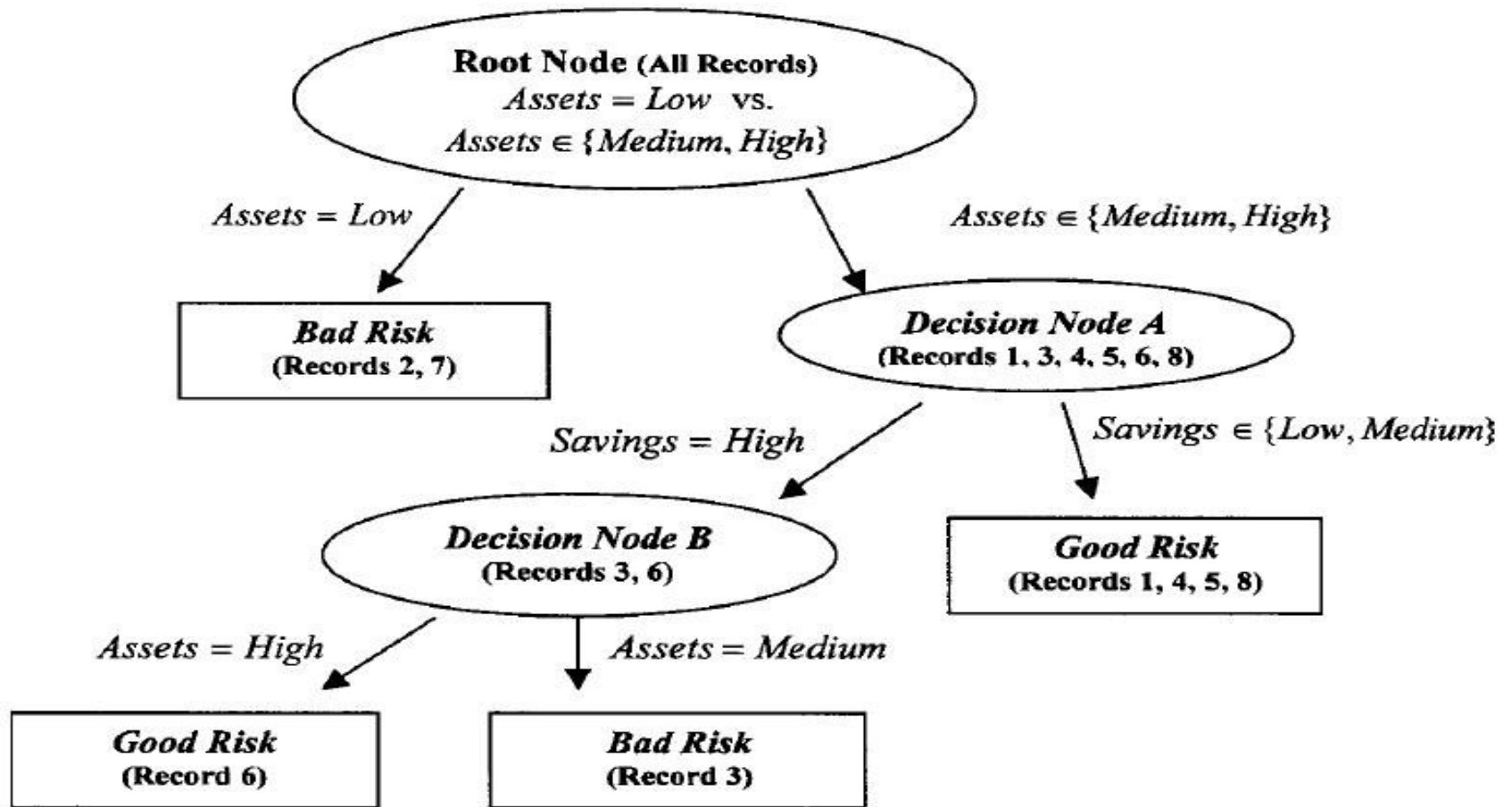
Split	P_L	P_R	$P(j t_L)$	$P(j t_R)$	$2P_L P_R$	$Q(s t)$	$\Phi(s t)$
1	0.167	0.833	G: 1 B: 0	G: .8 B: .2	0.2782	0.4	0.1112
2	0.5	0.5	G: 1 B: 0	G: 0.667 B: 0.333	0.5	0.6666	0.3333
3	0.333	0.667	G: 0.5 B: 0.5	G: 1 B: 0	0.4444	1	0.4444
5	0.667	0.333	G: 0.75 B: 0.25	G: 1 B: 0	0.4444	0.5	0.2222
6	0.333	0.667	G: 1 B: 0	G: 0.75 B: 0.25	0.4444	0.5	0.2222
7	0.333	0.667	G: 0.5 B: 0.5	G: 1 B: 0	0.4444	1	0.4444
8	0.5	0.5	G: 0.667 B: 0.333	G: 1 B: 0	0.5	0.6666	0.3333
9	0.167	0.833	G: 0.8 B: 0.2	G: 1 B: 0	0.2782	0.4	0.1112

Vrednosti $\Phi(s|t)$ za svako kandidat grananje za čvor A

CART Primer



CART Primer



Završetak

- Kada se iscrpe sva moguća grananja generisano je “puno” stablo.
- Ovakvo stablo ima najmanji stepen greške na trening skupu, ali često može dovesti do prepodešenosti modela.
- Stoga se stablo “odseca” kako bi se postigla generalnost modela.

Prečišćenost stabla

- Često svi listovi stabla nisu homogeni, što ostavlja određeni stepen greške prilikom klasifikacije.

Customer	Savings	Assets	Income	Credit Risk
004	High	Low	$\leq \$30,000$	Good
009	High	Low	$\leq \$30,000$	Good
027	High	Low	$\leq \$30,000$	Bad
031	High	Low	$\leq \$30,000$	Bad
104	High	Low	$\leq \$30,000$	Bad

Primer lista koji nije potpuno "čist"

- Dalja grananja nisu moguća, pa krajnje pravilo (odluka) ima pouzdanost od 60%.

C4.5 algoritam

- C4.5 u odnosu na CART nije ograničen na binarna grananja.
 - Generiše zasebnu granu za svaku vrednost kategoričkog atributa.
- C4.5 metod se zasniva na meri homogenosti.
- Svako kandidat grananje deli polazni skup T na nekoliko podskupova T_1, T_2, \dots, T_k .
- $entropy_reduction(S) = H(T) - H_S(T)$, gde je

$$H(X) = - \sum_j p_j \log_2(p_j)$$

C4.5 algoritam

- $H_S(T)$ je težinska suma entropija podskupova T_1, T_2, \dots, T_k i računa se kao:

$$H_S(T) = \sum_{i=1}^k P_i H_S(T_i)$$

- P_i predstavlja zastupljenost instanci u skupu i
- Za optimalno grananje uzima se ono sa najvećom vrednosti *entropy_reduction* (redukcijom entropije)

C4.5 primer

Customer	Savings	Assets	Income (\$1000s)	Credit Risk
1	Medium	High	75	Good
2	Low	Low	50	Bad
3	High	Medium	25	Bad
4	Medium	Medium	50	Good
5	Low	Medium	100	Good
6	High	High	25	Good
7	Low	Low	25	Bad
8	Medium	Medium	75	Good

Trening skup

Candidate Split	Child Nodes		
1	<i>Savings = low</i>	<i>Savings = medium</i>	<i>Savings = high</i>
2	<i>Assets = low</i>	<i>Assets = medium</i>	<i>Assets = high</i>
3	<i>Income ≤ \$25,000</i>		<i>Income > \$25,000</i>
4	<i>Income ≤ \$50,000</i>		<i>Income > \$50,000</i>
5	<i>Income ≤ \$75,000</i>		<i>Income > \$75,000</i>

Kandidat grananja za koreni čvor t

C4.5 primer

- 5/8 instanci je klasifikovano u klasu “Good”, a 3/8 u klasu “Bad”
- Početna entropija (pre grananja):

$$H(T) = - \sum_j p_j \log_2(p_j) = -\frac{5}{8} \log_2\left(\frac{5}{8}\right) - \frac{3}{8} \log_2\left(\frac{3}{8}\right) = 0.9544$$

- Za svako grananje se dobijena entropija poredi sa ovom vrednosti kako bi se odredilo koje grananje ima najveću redukciju u entropiji.

C4.5 primer

- Za kandidat grananje 1 (savings):

$$P_{\text{high}} = \frac{2}{8}, P_{\text{medium}} = \frac{3}{8}, P_{\text{low}} = \frac{3}{8}.$$

- Entropija za čvor savings = high

$$-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) = 1$$

- Entropija za čvor savings = medium

$$-\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - \frac{0}{3} \log_2 \left(\frac{0}{3} \right) = 0$$

- Entropija za čvor savings = low

$$-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) = 0.9183$$

C4.5 primer

- Težinska kombinacija ovih entropija:

$$H_{\text{savings}}(T) = \frac{2}{8}(1) + \frac{3}{8}(0) + \frac{3}{8}(0.9183) = 0.5944$$

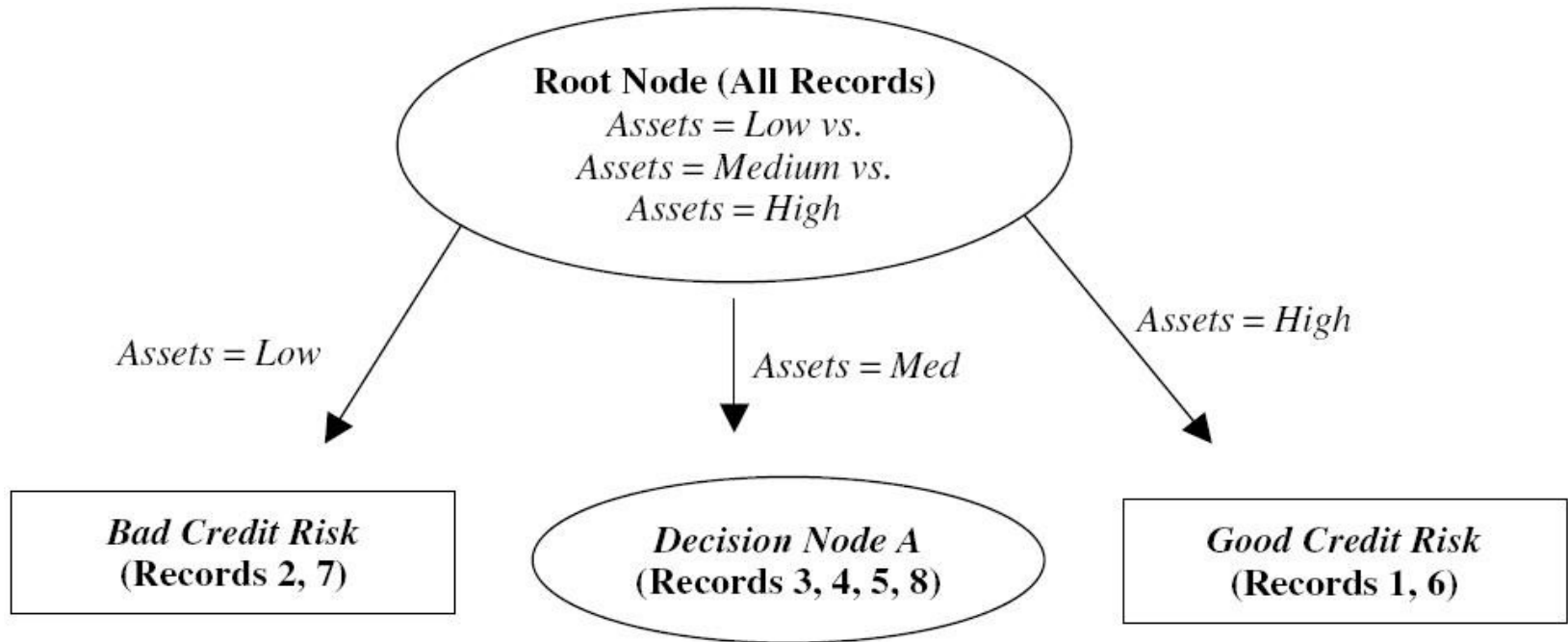
- Redukcija entropije za ovaj atribut je:

$$H(T) - H_{\text{savings}}(T) = 0.9544 - 0.5944 = 0.36$$

C4.5 primer

Candidate Split	Child Nodes	Information Gain (Entropy Reduction)
1	<i>Savings = low</i> <i>Savings = medium</i> <i>Savings = high</i>	0.36 bits
2	<i>Assets = low</i> <i>Assets = medium</i> <i>Assets = high</i>	0.5487 bits
3	<i>Income \leq \$25,000</i> <i>Income $>$ \$25,000</i>	0.1588 bits
4	<i>Income \leq \$50,000</i> <i>Income $>$ \$50,000</i>	0.3475 bits
5	<i>Income \leq \$75,000</i> <i>Income $>$ \$75,000</i>	0.0923 bits

C4.5 primer



C4.5 primer

Customer	Savings	Assets	Income (\$1000s)	Credit Risk
3	High	Medium	25	Bad
4	Medium	Medium	50	Good
5	Low	Medium	100	Good
8	Medium	Medium	75	Good

Instance čvora A

- Početna (pre grananja) entropija:

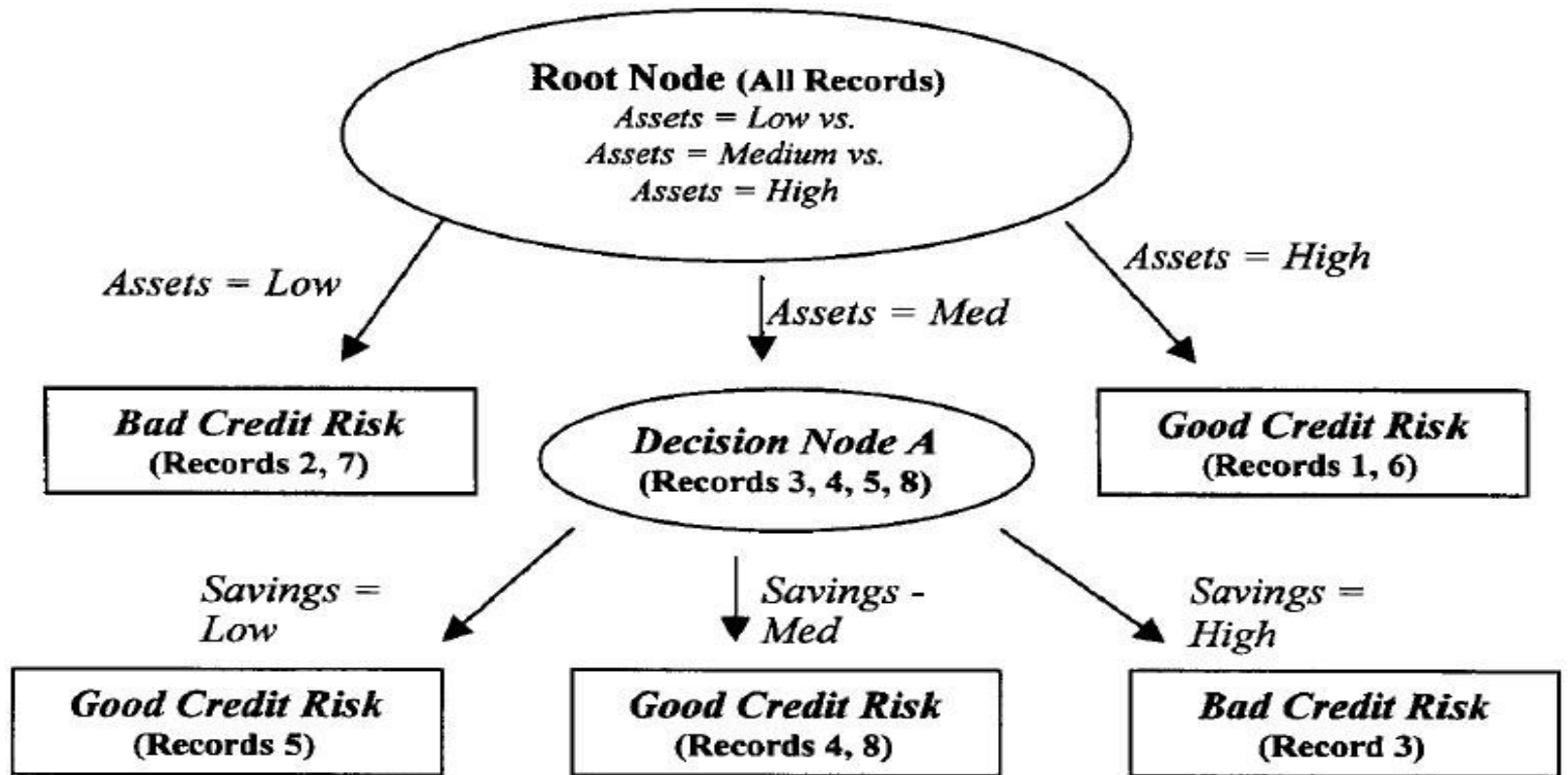
$$H(A) = - \sum_j p_j \log_2(p_j) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.8113$$

C4.5 primer

Candidate Split	Child Nodes		
1	<i>Savings = low</i>	<i>Savings = medium</i>	<i>Savings = high</i>
3	<i>Income \leq \$25,000</i>		<i>Income > \$25,000</i>
4	<i>Income \leq \$50,000</i>		<i>Income > \$50,000</i>
5	<i>Income \leq \$75,000</i>		<i>Income > \$75,000</i>

Kandidat grananja za čvor A

C4.5 primer



C4.5 stablo odlučivanja

Pravila

Antecedent	Consequent	Support	Confidence
<i>If assets = low</i>	<i>then bad credit risk.</i>	$\frac{2}{8}$	1.00
<i>If assets = high</i>	<i>then good credit risk.</i>	$\frac{2}{8}$	1.00
<i>If assets = medium and savings = low</i>	<i>then good credit risk.</i>	$\frac{1}{8}$	1.00
<i>If assets = medium and savings = medium</i>	<i>then good credit risk.</i>	$\frac{2}{8}$	1.00
<i>If assets = medium and savings = high</i>	<i>then bad credit risk.</i>	$\frac{1}{8}$	1.00

Pravila generisana na osnovu stabla odlučivanja

Reference

- *Discovering Knowledge in Data: An Introduction to Data Mining* (Wiley, 2005) Larose D.
- SQL Server Data mining tutorial (Lessons 1-6)
<http://technet.microsoft.com/en-us/library/ms167167.aspx>