COMP0271: Inteligência Artificial Árvores de decisão



Subconjunto do censo 1990 US: 48.000 registros, 16 atributos [Kohavi 1995]

-		Ļ

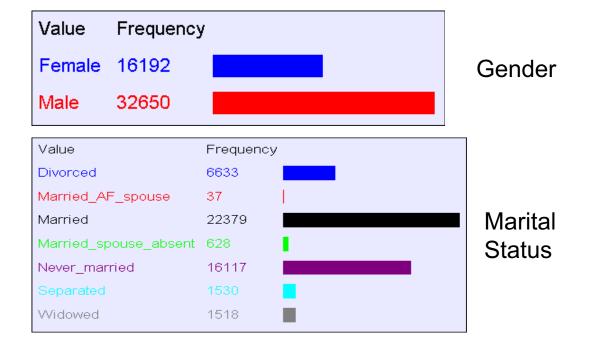
age	employment	education	edunum	marital		job	relation	race	gender	hours	country	wealth
39	State gov	Bachelors	13	Never mar		Adm cleric	Not in fam	White	Male	40	United Sta	noor
51	Self_emp_not_i			Married		Exec mana		White	Male		United_Sta	
	Private	HS_grad		Divorced			Not in fam		Male		United Sta	
	Private	11th		Married		Handlers c		Black	Male		United_Sta	_
	Private	Bachelors		Married		Prof specia		Black	Female		Cuba	poor
	Private	Masters		Married		Exec mana		White	Female		United Sta	-
	Private	9th		Married sp			Not in fam		Female		Jamaica	poor
52				Married_sp		Exec mana		White	Male		United Sta	-
	Private	Masters		Never mar	•••		Not in fam		Female		United_Sta	
	Private	Bachelors		Married		Exec mana		White	Male		United_Sta	
	Private	Some colle		Married	•••	Exec mana		Black	Male		United_Sta	
	State gov	Bachelors		Married	•••	Prof specia		Asian	Male			rich
											India	
	Private	Bachelors		Never_mar		_	Own_child	White	Female		United_Sta	
	Private	Assoc_acd		Never_mar	•••	Sales	Not_in_fam		Male		United_Sta	-
41		Assoc_voc		Married		Craft_repai		Asian	Male		*MissingVa	
	Private	7th_8th		Married	•••	Transport_i		Amer_India			Mexico	poor
	Self_emp_not_i			Never_mar			Own_child		Male		United_Sta	
	Private	HS_grad	9	Never_mar		Machine_o		White	Male		United_Sta	
	Private	11th		Married		Sales	Husband	White	Male		United_Sta	
44	Self_emp_not_i	Masters	14	Divorced		Exec_mana	Unmarried	White	Female	45	United_Sta	rich
41	Private	Doctorate	16	Married		Prof_specia	Husband	White	Male	60	United_Sta	rich
:	:	:	:	:	:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:	:	:	:	:	:

Atributos

age	edunum	race	hours_worked
employment	marital	gender	country
taxweighting	job	capitalgain	wealth
education	relation	capitalloss	agegroup

O que podemos fazer com essas informações?

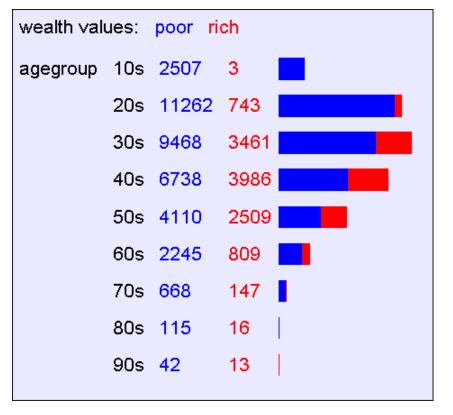
Montar Histogramas



O que podemos fazer com essas informações?

- Montar Tabelas de Contingência (1-d, 2-d, 3-d)
 - Escolha d atributos: a₁,a₂, ... a_d.
 - 2. Para cada combinação possível de valores: $a_1,=x_1, a_2,=x_2,... a_d,=x_d$, grave o quão frequentemente essa combinação ocorre

Tabela de Contingência 2-d



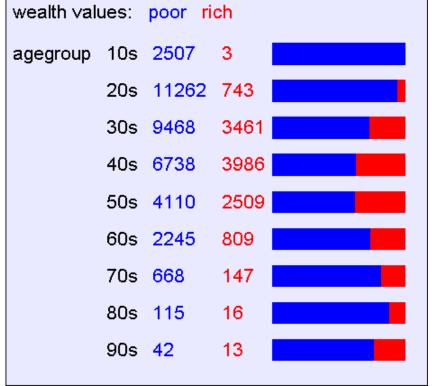
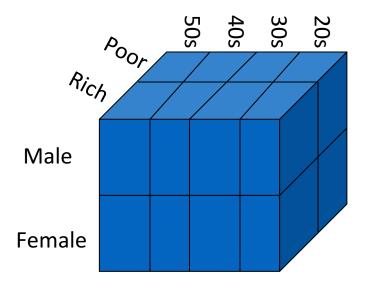


Tabela de Contingência 2-d

job valu	es:	Adm_clerical	Craft_	repair		Farm	ning_fis	hing	Ma	chine_c	op_ins	oct	Priv_ho	ouse_	serv	Protec	tive_s	erv Tech_support	
Missing	gValue	Armed_Forces	Exec_	manage	erial	Hand	dlers_c	leanei	rs Oth	ner_ser	vice		Prof_s	pecial	ty	Sales		Transport_movin	g
marital	Divorce	k	270	1192	0	679	890	90	197	434	762		795	121	664	239	254		
	Married_	_AF_spouse	5	6	0	4	3	1	1	1	5		4	1	5	0	1		
	Married		928	1495	7	3818	3600	869	724	1469	1088		3182	583	2491	609	1489		
	Married_	_spouse_absent	45	84	0	77	52	35	32	37	92		64	7	55	9	30		
	Never_n	narried	1242	2360	8	1301	1260	434	1029	872	2442		1849	237	1992	506	486		
	Separat	ed	97	224	0	160	126	23	63	123	275		145	23	146	48	56		
	Widowe	d	222	250	0	73	155	38	26	86	259		133	11	151	35	39		

Tabela de Contingência 3-d

Mais difícil de visualizar



Refletindo...

- Com 16 atributos,
 - quantas tabelas 1-d existem? 16
 - quantas 2-d? 16 * 15 / 2 = 120
 - 3-d? 16*15*14 / 6 = 560
- Com 100 atributos, quantas tabelas 3-d existem?
 - 100*99*98 / 6 = 161.700 !!!!!

Imagina ter que olhar/analisar isso manualmente???

Mineração dos dados

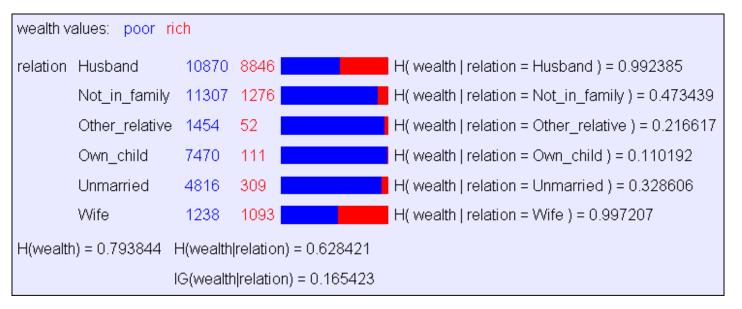
Automatização do processo de busca por padrões nos dados.

Que padrões são interessantes? Quais seriam meras ilusões? Como podem ser explorados?

usar Teoria da Informação (Shannon, 1948)

Pedindo ajuda à máquina...

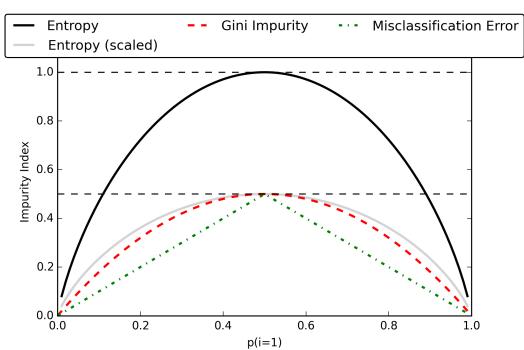
Considerando algo que você quer prever (ex: wealth), é mais fácil pedir ao computador para encontrar o atributo de maior **Ganho de Informação**



Ganho de informação e Impureza

$$IG(D_p,f) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N_p} I(D_j)$$
 (Impureza)

$$I_H(t) = -\sum_{i=1}^c p(i|t) \log_2 p(i|t)$$
 (Entropia)



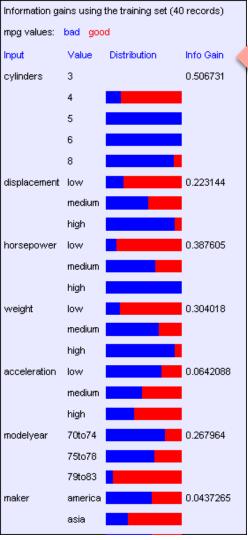
Árvore de Decisão

- Plano estruturado em árvore de um conjunto de atributos para testar a fim de se prever a saída.
- Qual atributo deve ser testado primeiro? Aquele que produzir maior Ganho de Informação
 - E aí, recorrentemente...

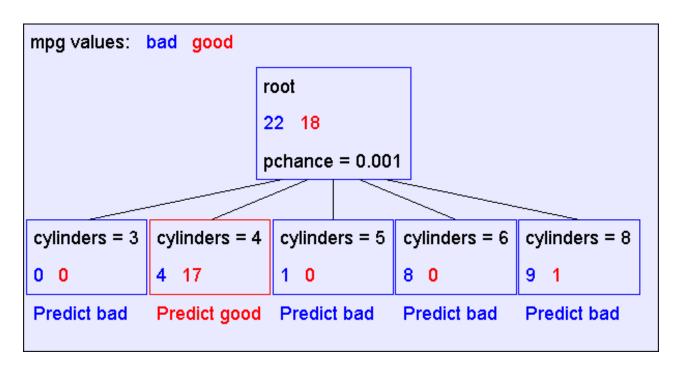
Outro dataset 40 exemplos

mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europe
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low	medium	low	medium	70to74	asia
bad	4	low	medium	low	low	70to74	asia
bad	8	high	high	high	low	75to78	america
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
bad	8	high	high	high	low	70to74	america
good	8	high	medium	high	high	79to83	america
bad	8	high	high	high	low	75to78	america
good	4	low	low	low	low	79to83	america
bad	6	medium	medium	medium	high	75to78	america
good	4	medium	low	low	low	79to83	america
good	4	low	low	medium	high	79to83	america
bad	8	high	high	high	low	70to74	america
good	4	low	medium	low	medium	75to78	europe
bad	5	medium	medium	medium	medium	75to78	europe

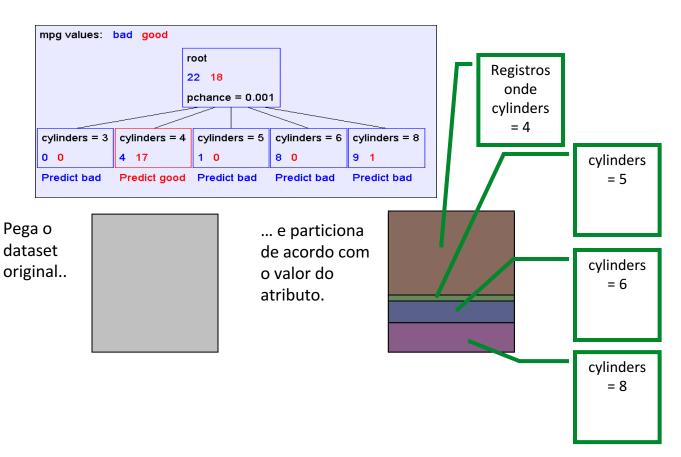
Considere prever MPG (milhas por galão)...



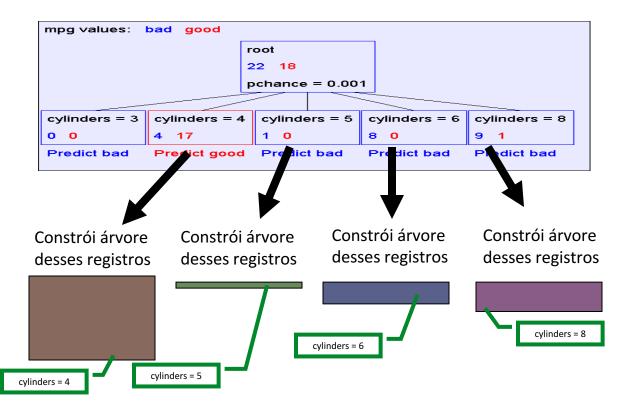
Um pedaço da árvore



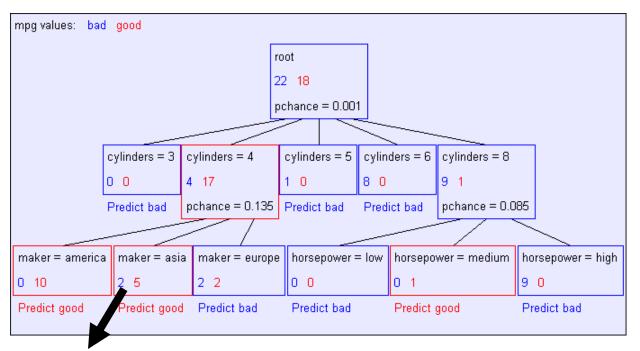
Passo recursivo



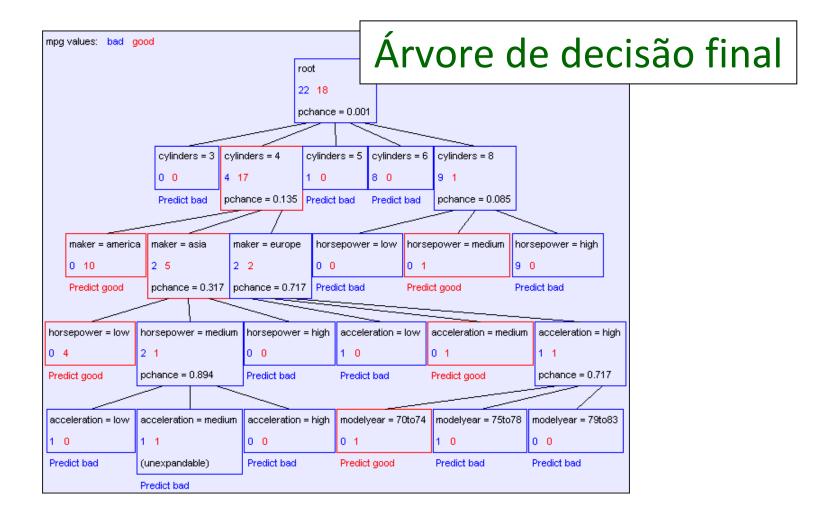
Passo recursivo

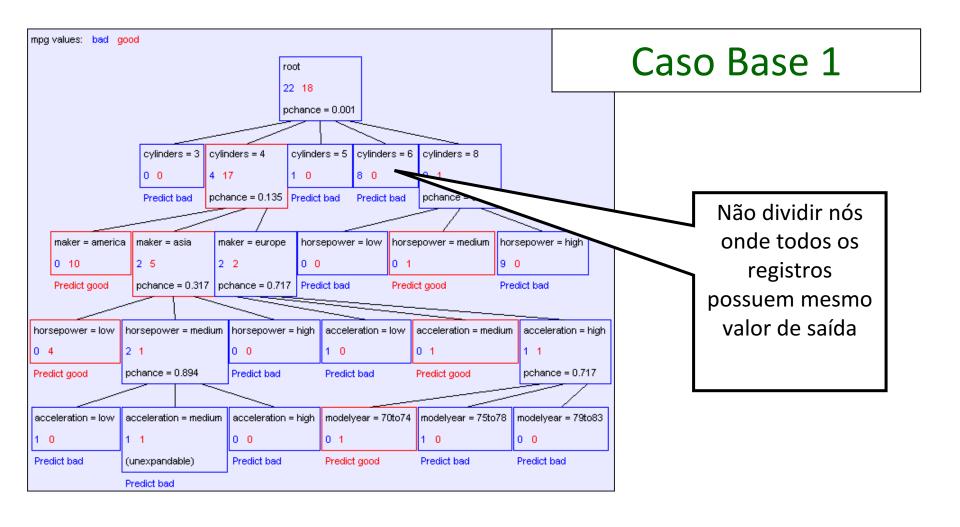


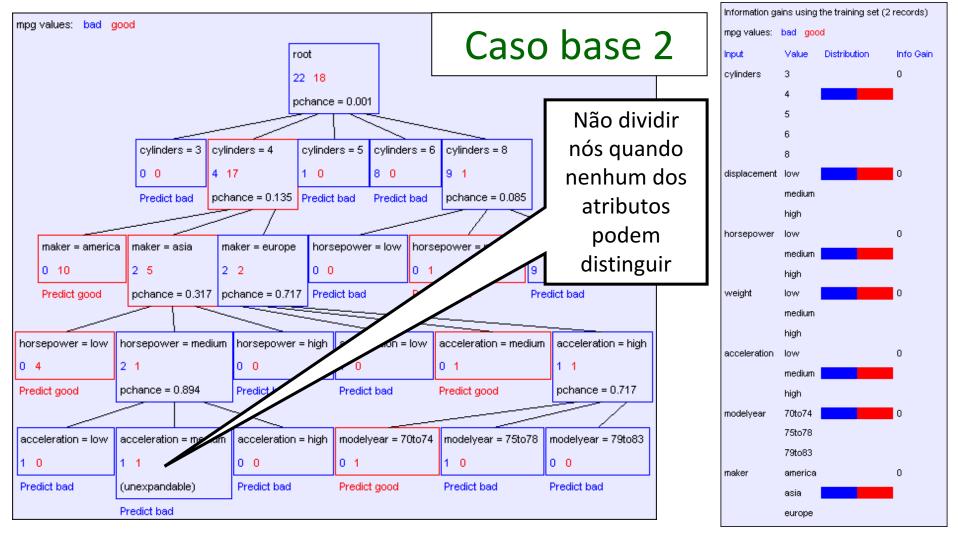
Segundo nível da árvore...



Constrói uma árvore a partir dos 7 (sete) registros em que cylinders = 4 e maker = Asia







Algoritmo sumarizado

BuildTree(DataSet,Output)

If all output values are the same in *DataSet*, return a leaf node that says "predict this unique output"

If all input values are the same, return a leaf node that says "predict the majority output"

Else find attribute *X* with highest Info Gain

Suppose *X* has n_X distinct values (i.e. X has arity n_X).

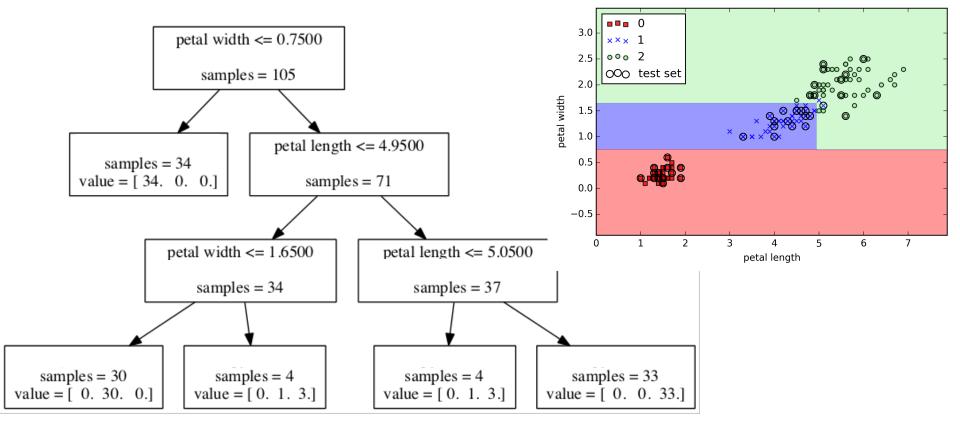
Create and return a non-leaf node with n_x children.

The *i*th child should be built by calling

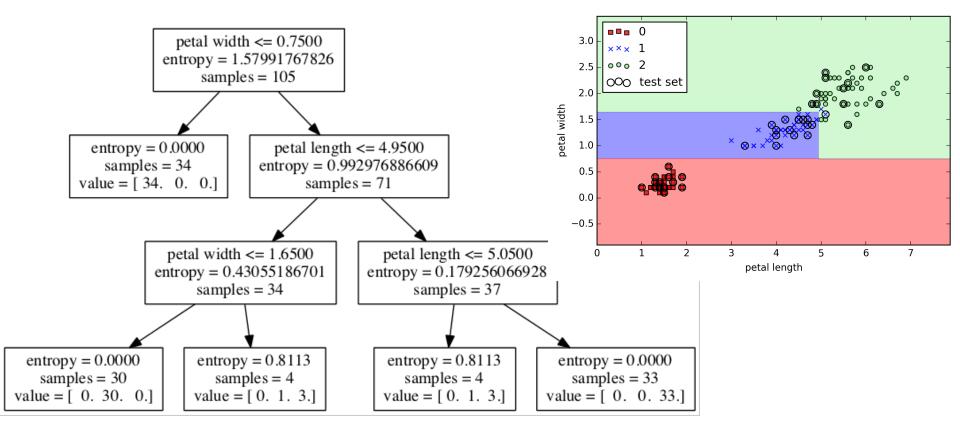
BuildTree(DS_i,Output)

(Where DS_i built consists of all those records in DataSet for which X = ith distinct value of X)

Dataset Iris



Dataset Iris



Exemplo brinquedo (livro): Devo jogar tenis?

Dia	Aspecto	Temp.	Umidade	Vento	Jogar Tênis
D1	Sol	Quente	Elevada	Fraco	Não
D2	Sol	Quente	Elevada	Forte	Não
D3	Nuvens	Quente	Elevada	Fraco	Sim
D4	Chuva	Ameno	Elevada	Fraco	Sim
D5	Chuva	Fresco	Normal	Fraco	Sim
D6	Chuva	Fresco	Normal	Forte	Não
D7	Nuvens	Fresco	Normal	Fraco	Sim
D8	Sol	Ameno	Elevada	Fraco	Não
D9	Sol	Fresco	Normal	Fraco	Sim
D10	Chuva	Ameno	Normal	Forte	Sim
D11	Sol	Ameno	Normal	Forte	Sim
D12	Nuvens	Ameno	Elevada	Forte	Sim
D13	Nuvens	Quente	Normal	Fraco	Sim
D14	Chuva	Ameno	Elevada	Forte	Não

O conteúdo de informação da árvore é calculado a partir das probabilidades das diferentes classificações.

```
p(sim) = 9/14; p(nao) = 5/14

I[JogarTenis] = -p(sim)*log_2(p(sim)) -p(nao)*log_2(p(nao))

= -(9/14)*log_2(9/14) -(5/14)*log_2(5/14)

= 0.9404858
```

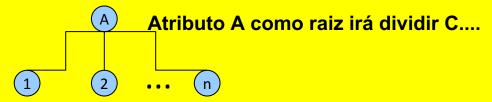
Ganho de Informação

Ganho
$$A = I C - \sum_{v \in valores A} \frac{|C_v|}{|C|} \cdot I C_v$$

Mede a redução na Entropia da (sub)árvore obtida pela divisão escolhida.

Objetivo: escolher a divisão que obtém maior redução (maior Ganho)

Considere um conjunto C de exemplos de treinamento



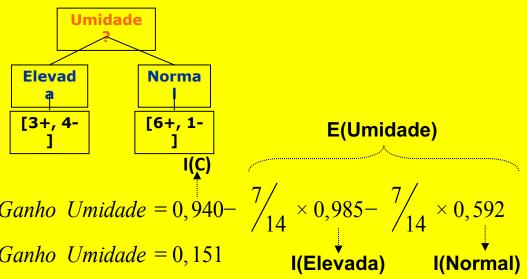
...em subconjuntos $\{C_1, C_2, ...C_n\}$

<u>Primeiro passo</u>: são analisados todos os atributos, começando pela Umidade.

$$C = [9, 5-]$$

$$I C = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0,940$$

).	Umidade	Vento	Jogar Tênis	
te	Elevada	Fraco	Não	
te	Elevada	Forte	Não	
te	Elevada	Fraco	Sim	
10	Elevada	Fraco	Sim	
:О	Normal	Fraco	Sim	
:0	Normal	Forte	Não	
:О	Normal	Fraco	Sim	
10	Elevada	Fraco	Não	
:0	Normal	Fraco	Sim	
10	Normal	Forte	Sim	
10	Normal	Forte	Sim	6
10	Elevada	Forte	Sim	U
te	Normal	Fraco	Sim	
10	Elevada	Forte	Não	G



Para o resto dos atributos....

$$C = [9+, 5-]$$

$$I = 0,940$$

$$Aspect$$

$$MAX$$

$$Ganho(Umidade) = 0,151$$

$$Ganho(Vento) = 0,048$$

$$Ganho(Aspecto) = 0,247$$

$$Ganho(Temp.) = 0,029$$

$$[2+, 3-]$$

$$[4+, 0-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

$$[3+, 2-]$$

Ganho (
$$Aspecto$$
) = 0,940 - $\binom{5}{14}$ × 0,971 - $\binom{4}{14}$ × 0,0 - $\binom{5}{14}$ × 0,971
Ganho ($Aspecto$) = 0,247

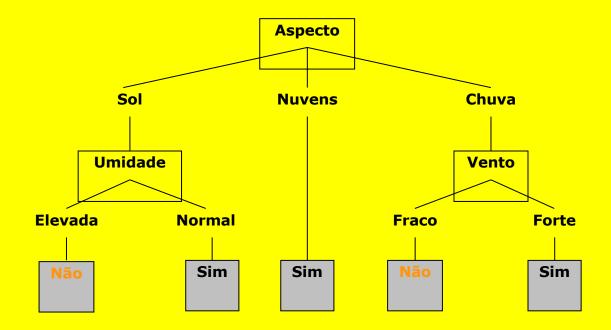
Dia	Aspecto	Temp.	Umidade	Vento	Jogar Tênis
D1	Sol	Quente	Elevada	Fraco	Não
D2	Sol	Quente	Elevada	Forte	Não
D3	Nuvens	Quente	Elevada	Fraco	Sim
D4	Chuva	Ameno	Elevada	Fraco	Sim
D5	Chuva	Fresco	Normal	Fraco	Sim
D6	Chuva	Fresco	Normal	Forte	Não
D7	Nuvens	Fresco	Normal	Fraco	Sim
D8	Sol	Ameno	Elevada	Fraco	Não
D9	Sol	Fresco	Normal	Fraco	Sim
D10	Chuva	Ameno	Normal	Forte	Sim
D11	Sol	Ameno	Normal	Forte	Sim
D12	Nuvens	Ameno	Elevada	Forte	Sim
D13	Nuvens	Quente	Normal	Fraco	Sim
D14	Chuva	Ameno	Flevada	Forte	Não

$$Ganho(C_{Sol}, Umidade) = 0.971 - \binom{3}{5} \times 0.0 - \binom{2}{5} \times 0.0 = 0.971$$

$$MAX Ganho(C_{Sol}, Temp.) = 0.971 - \binom{2}{5} \times 0.0 - \binom{2}{5} \times 1.0 - \binom{1}{5} \times 0.0 = 0.570$$

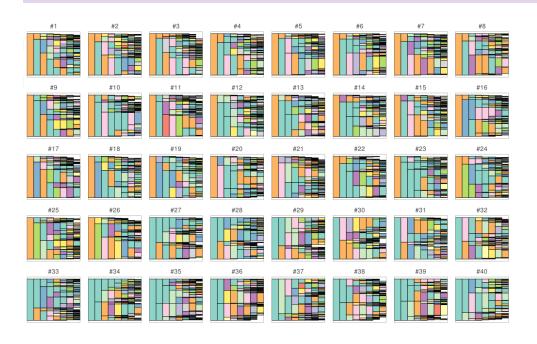
$$Ganho(C_{Sol}, Vento) = 0.971 - \binom{2}{5} \times 1.0 - \binom{3}{5} \times 0.918 = 0.019$$

$$= Ganho(C_{Sol}, Umidade)$$



∀ x Aspecto(x,Sol) ∧ Umidade(x, Normal) ⇒ JogarTenis(x)
 ∀ x Aspecto(x,Nuvens) ⇒ JogarTenis(x)
 ∀ x Aspecto(x,Chuva) ∧ Vento(x, Forte) ⇒ JogarTenis(x)

Árvores de decisão --> Random forests



(ensemble learning)

Inúmeras árvores de decisão

Random forests

- 1. Draw a random **bootstrap** sample of size *n* (randomly choose *n* samples from the training set with replacement).
- 2. Grow a decision tree from the bootstrap sample. At each node:
 - 1. Randomly select *d* features without replacement.
 - 2. Split the node using the feature that provides the best split according to the objective function, for instance, by maximizing the information gain.
- 3. Repeat the steps 1 to 2 k times.
- 4. Aggregate the prediction by each tree to assign the class label by majority vote.

Random forests



Majority vote

