

# Example Data Sets

- Contact Lens (symbolic)
- Weather (symbolic data)
- Weather ( numeric +symbolic)
- Iris (numeric; outcome:symbolic)
- CPU Perf.(numeric; outcome:numeric)
- Labor Negotiations (missing values)
- Soybean

# Contact Lens Data

age	spectacle prescription	astigmatism	tear production rate	recommendation lenses
young	myope	no	reduced	none
young	myope	no	normal	soft
young	myope	yes	reduced	none
young	myope	yes	normal	hard
young	hypermetrope	no	reduced	none
young	hypermetrope	no	normal	soft
young	hypermetrope	yes	reduced	none
young	hypermetrope	yes	normal	hard
pre-presbyopic	myope	no	reduced	none
pre-presbyopic	myope	no	normal	soft
pre-presbyopic	myope	yes	reduced	none
pre-presbyopic	myope	yes	normal	hard
pre-presbyopic	hypermetrope	no	reduced	none
pre-presbyopic	hypermetrope	no	normal	soft
pre-presbyopic	hypermetrope	yes	reduced	none
pre-presbyopic	hypermetrope	yes	normal	none
presbyopic	myope	no	reduced	none
presbyopic	myope	no	normal	none
presbyopic	myope	yes	reduced	none
presbyopic	myope	yes	normal	hard
presbyopic	hypermetrope	no	reduced	none
presbyopic	hypermetrope	no	normal	soft
presbyopic	hypermetrope	yes	reduced	none
presbyopic	hypermetrope	yes	normal	none

# Structural Patterns

- Part of structural description

```
If tear production rate = reduced then recommendation = none  
Otherwise, if age = young and astigmatic = no  
    then recommendation = soft
```

- Example is simplistic because all combinations of possible values are represented in table

# Structural Patterns

- In most learning situations, the set of examples given as input is far from complete
- Part of the job is to generalize to other, new examples

# Weather Data

outlook	temperature	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

# Weather Problem

- This creates 36 possible combinations ( $3 \times 3 \times 2 \times 2 = 36$ ), of which 14 are present in the set of examples

If outlook = sunny and humidity = high	then play = no
If outlook = rainy and windy = true	then play = no
If outlook = overcast	then play = yes
If humidity = normal	then play = yes
If none of the above	then play = yes

# Weather Data with Some Numeric Attributes

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

# Classification and Association Rules

## Rules

- **Classification Rules:** rules which predict the classification of the example in terms of whether to play or not

```
If outlook = sunny and humidity = >83,      then play = no
```

# Classification and Association Rules

## Rules

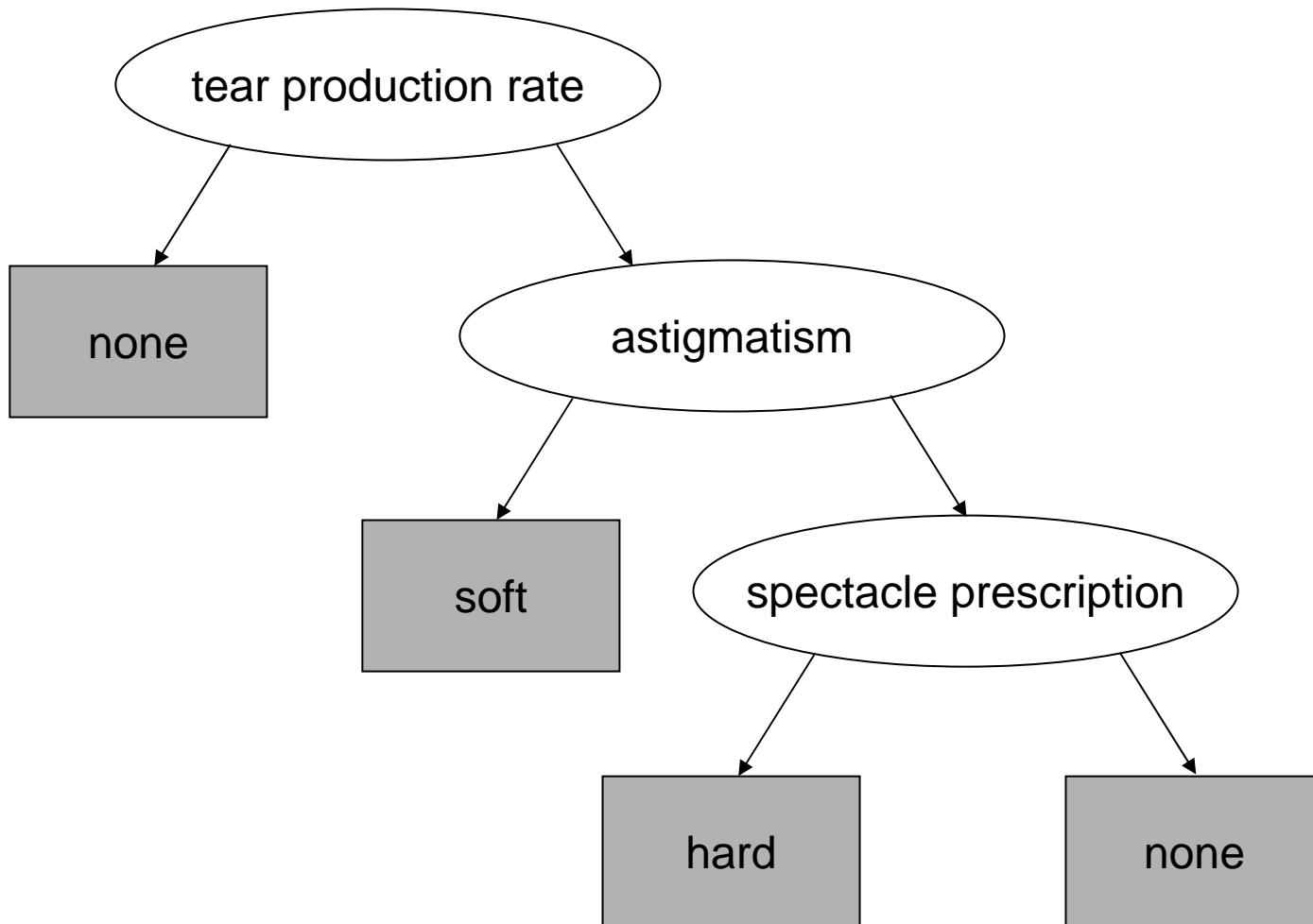
- Association Rules: rules which strongly associate different attribute values
- Association rules which derive from weather table

If temperature = cool	then humidity = normal
If humidity = normal and windy = false	then play = yes
If outlook = sunny and play = no	then humidity = high
If windy = false and play = no and humidity = high	then outlook = sunny

# Rules for Contact Lens Data

```
If tear production rate = reduced then recommendation = none
If age = young and astigmatic = no and
    tear production rate = normal then recommendation = soft
If age = pre-presbyopic and astigmatic = no and
    tear production rate = normal then recommendation = soft
If age = presbyopic and spectacle prescription = myope and
    astigmatic = no then recommendation = none
If spectacle prescription = hypermetrope and astigmatic = no and
    tear production rate = normal then recommendation = soft
If spectacle prescription = myope and astigmatic = yes and
    tear production rate = normal then recommendation = hard
If age = young and astigmatic = yes and
    tear production rate = normal then recommendation = hard
If age = pre-presbyopic and
    spectacle prescription = hypermetrope and astigmatic = yes
    then recommendation = none
If age = presbyopic and spectacle prescription = hypermetrope and
    astigmatic = yes then recommendation = none
```

# Decision Tree for Contact Lens Data



# Iris Data

	sepal length	sepal width	pedal length	pedal width	type
1	5.1	3.5	1.4	0.2	<i>Iris setosa</i>
2	4.9	3.0	1.4	0.2	<i>Iris setosa</i>
3	4.7	3.2	1.3	0.2	<i>Iris setosa</i>
4	4.6	3.1	1.5	0.2	<i>Iris setosa</i>
5	5.0	3.6	1.4	0.2	<i>Iris setosa</i>
...					
51	7.0	3.2	4.7	1.4	<i>Iris</i>
52	6.4	3.2	4.5	1.5	<i>Iris</i>
53	6.9	3.1	4.9	1.5	<i>Iris</i>
54	5.5	2.3	4.0	1.3	<i>Iris</i>
55	6.5	2.8	4.6	1.5	<i>Iris</i>
101	6.3	3.3	6.0	2.5	<i>Iris virginica</i>
102	5.8	2.7	5.1	1.9	<i>Iris virginica</i>
103	7.1	3.0	5.9	2.1	<i>Iris virginica</i>
104	6.3	2.9	5.6	1.8	<i>Iris virginica</i>
105	6.5	3.0	5.8	2.2	<i>Iris virginica</i>

# Iris Rules Learned

- If petal-length < 2.45 then Iris-setosa
- If sepal-width < 2.10 then Iris-versicolor
- If sepal-width < 2.45 and petal-length < 4.55 then Iris-versicolor
- ...

# CPU Performance Data

	cycle	main memory (Kb)		cache	channels		performance
	time (ns)	min	max	(Kb)	min	max	
	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	PRP
1	125	256	6000	256	16	128	198
2	29	8000	32000	32	8	32	269
3	29	8000	32000	32	8	32	220
4	29	8000	32000	32	8	32	172
5	29	8000	16000	32	8	16	132
...							
207	125	2000	8000	0	2	14	52
208	480	512	8000	32	0	0	67
209	480	1000	4000	0	0	0	45

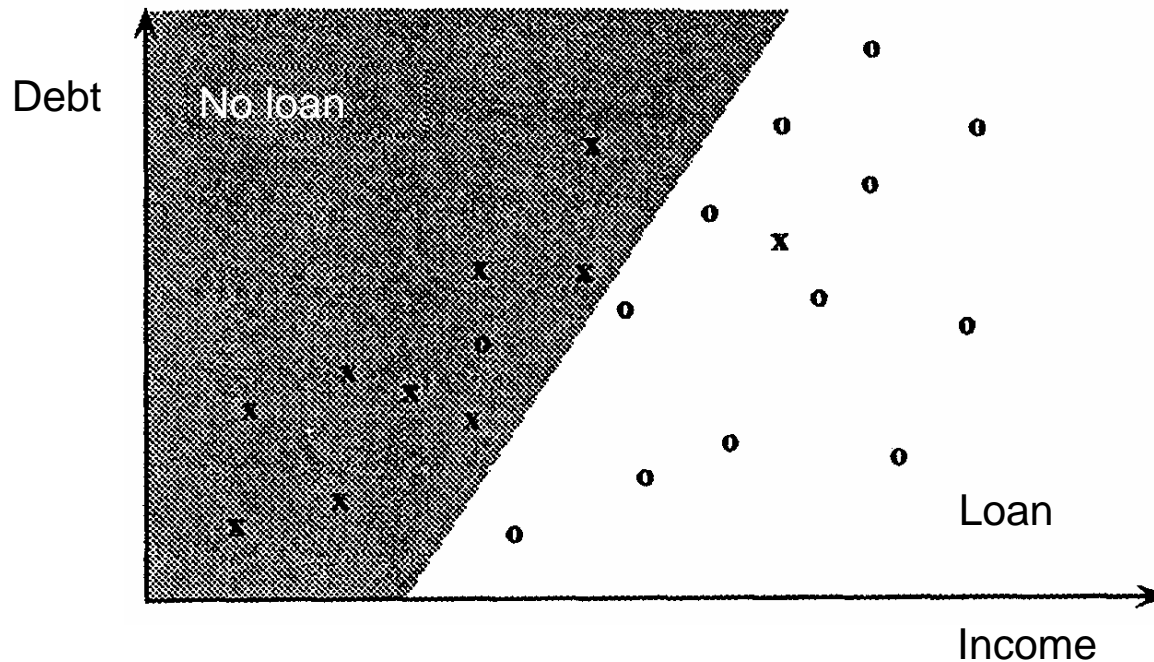
# CPU Performance

- Numerical Prediction: outcome as linear sum of weighted attributes
- Regression equation:
- $PRP = -55.9 + .049MYCT + . + 1.48CHMAX$
- Regression can discover linear relationships, not non-linear ones

# Labor Negotiations Data

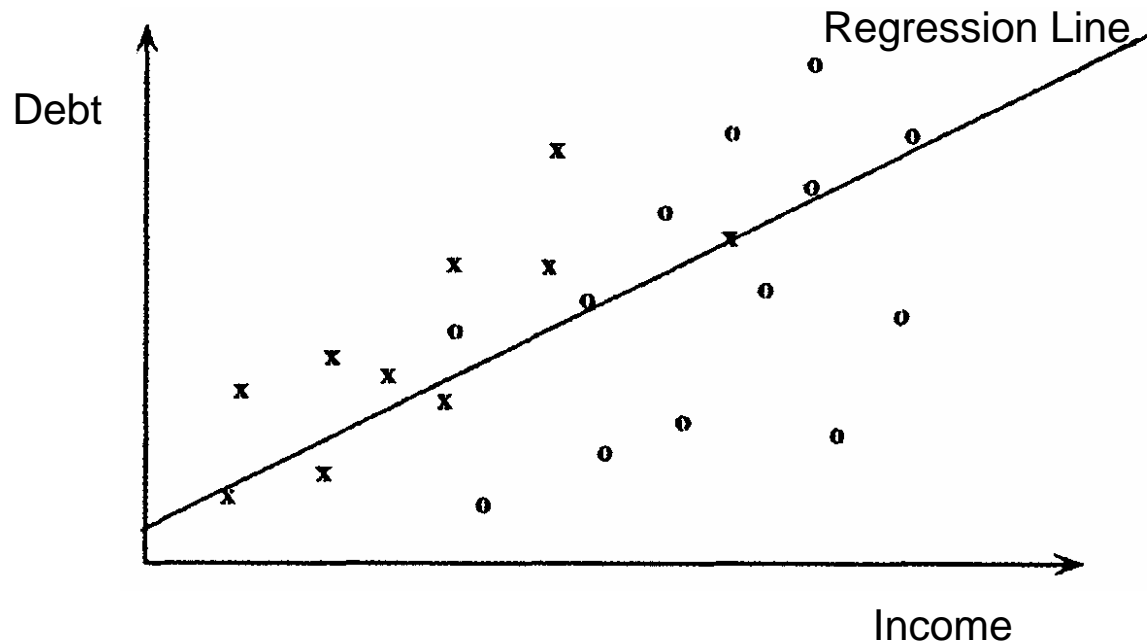
attribute	type	1	2	3	...	40
duration	(number of years)	1	2	3		2
wage increase first year	percentage	2%	4%	4.3%		4.5
wage increase second year	percentage	?	5%	4.4%		4.0
wage increase third year	percentage	?	?	?		?
cost of living adjustment	{none, tcf, tc}	none	tcf	?		none
working hours per week	(number of hours	28	35	38		40
pension	{none, ret-allw,	none	?	?		?
standby pay	percentage	?	13%	?		?
shift-work supplement	percentage	?	5%	4%		4
education allowance	{yes, no}	yes	?	?		?
statutory holidays	(number of days)	11	15	12		12
vacation	{below-avg, avg,	avg	gen	gen		avg
long-term disability	{yes, no}	no	?	?		yes
dental plan contribution	{none, half, full}	none	?	full		full
bereavement assistance	{yes, no}	no	?	?		yes
health plan contribution	{none, half, full}	none	?	full		half
acceptability of contract	{good, bad}	bad	good	good		good

# Classification



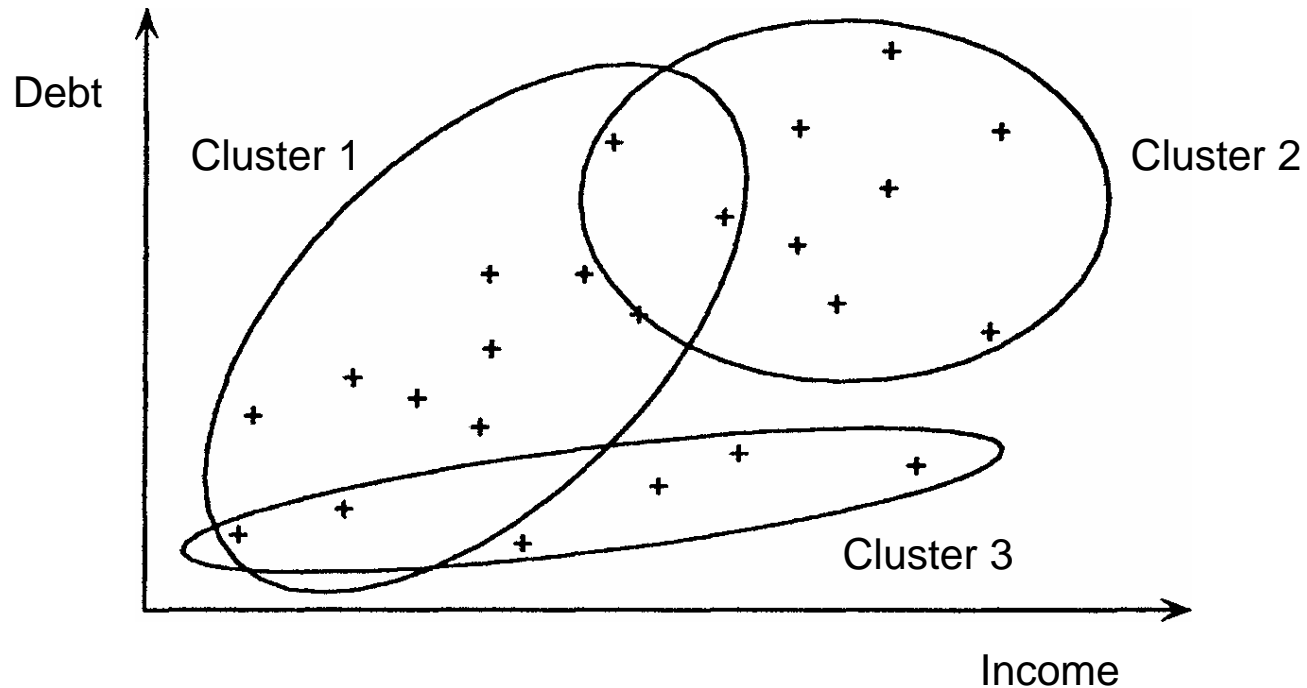
A simple linear classification boundary for the loan data set; shaded region denotes class “no loan”

# Linear Regression



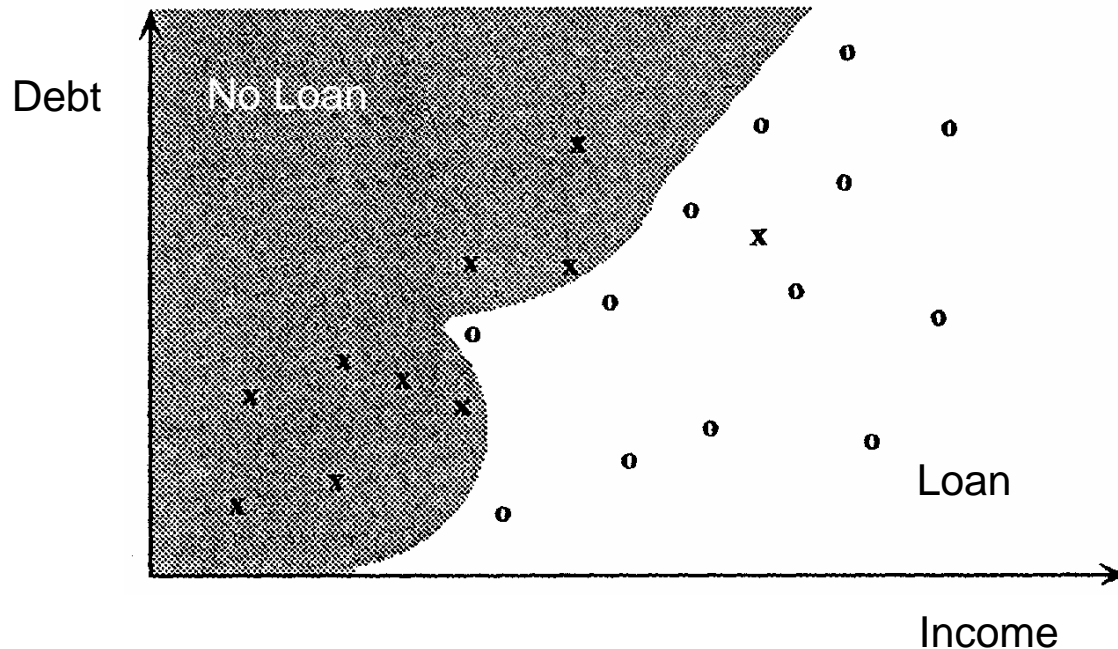
A simple linear regression for the loan data set

# Clustering



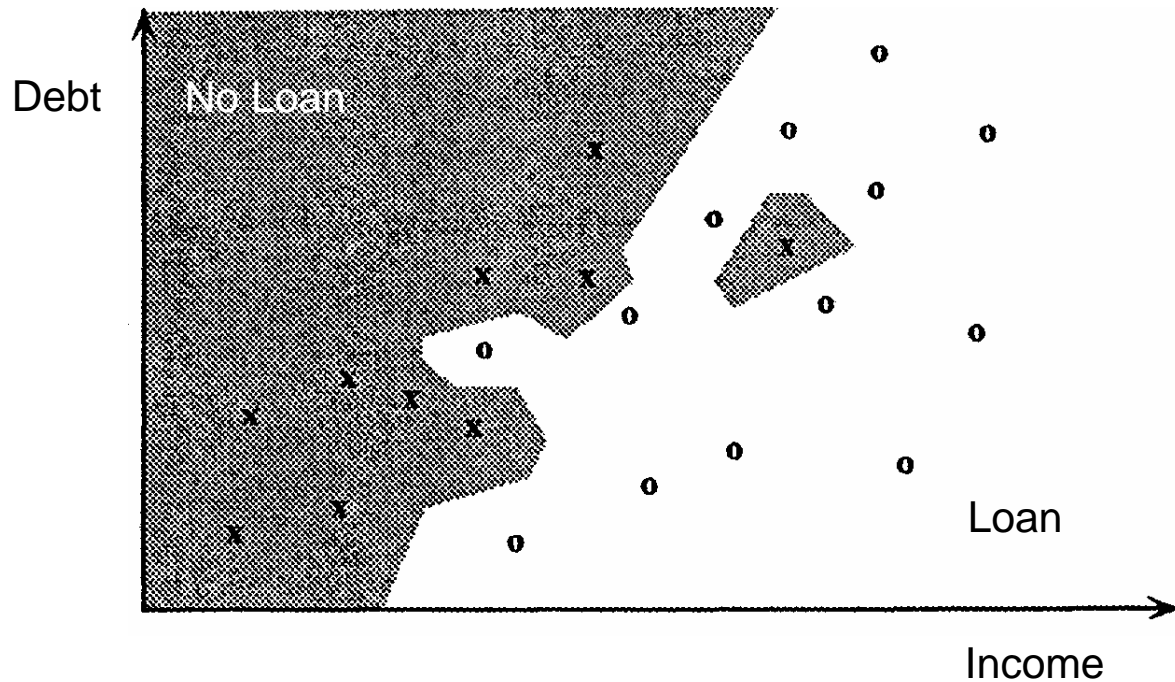
A simple clustering of the loan data set into 3 clusters; note that the original labels are replaced by '+'s

# Non-Linear Classification



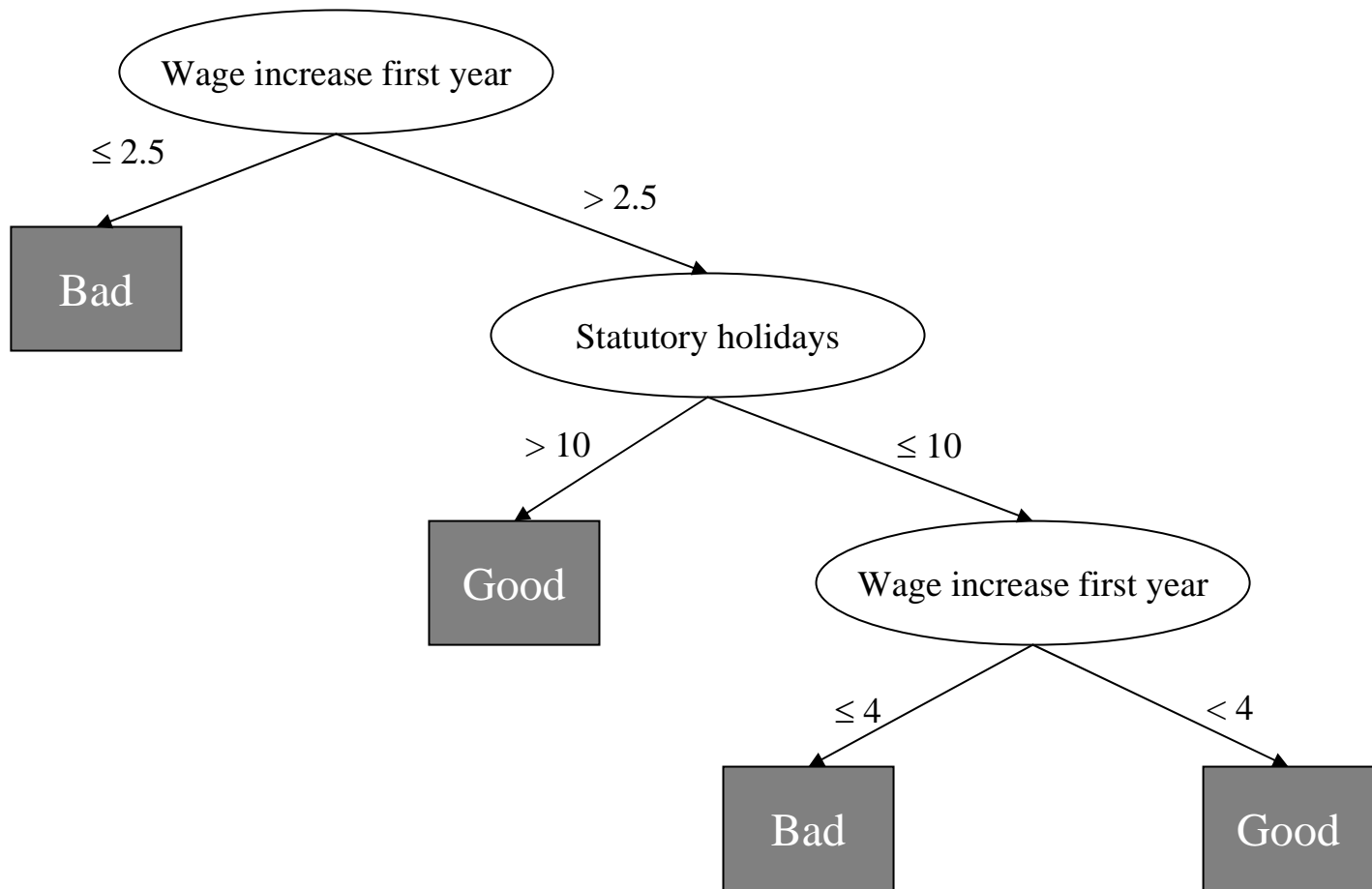
An example of classification boundaries learned by a non-linear classifier (such as a neural network) for the loan data set

# Nearest Neighbor Classifier

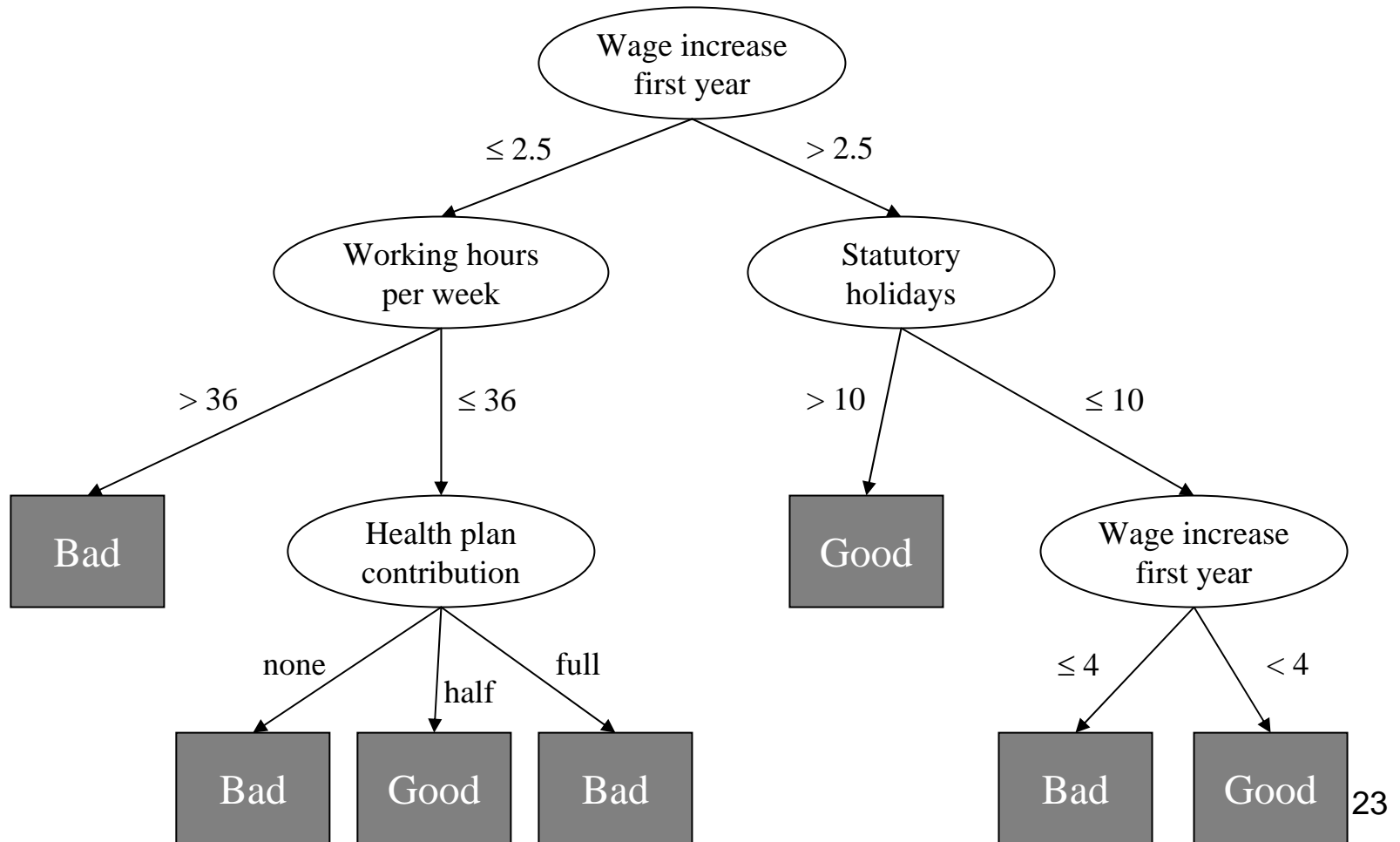


Classification boundaries for a nearest neighbor classifier for the loan data set

# Decision Trees for ...



# ... Labor Negotiations Data



# Soy Bean Data

	Attribute	Number of Values	Sample Value
<b>Environment</b>	time of occurrence	7	July
	precipitation	3	above normal
	temperature	3	normal
<b>Seed</b>	condition	2	normal
	mold growth	2	absent
	discoloration	2	absent
<b>Fruit</b>	condition of fruit pods	4	normal
<b>Leaves</b>	condition	2	abnormal
	yellow leaf spot halo	3	absent
	leaf spot margins	3	no data
<b>Stem</b>	condition	2	abnormal
	stem lodging	2	yes
	stem cankers	4	above the soil line
<b>Roots</b>	condition	3	normal
<b>Diagnosis</b>		19	diaporthe stem canker

# Two Example Rules

```
If      [leaf condition is normal and  
        stem condition is abnormal and  
        stem cankers is below soil line and  
        canker lesion color is brown]
```

```
then
```

```
    diagnosis is rhizoctonia root rot
```

```
If      [leaf malformation is absent and  
        stem condition is abnormal and  
        stem cankers is below soil line and  
        canker lesion color is brown]
```

```
then
```

```
    diagnosis is rhizoctonia root rot
```

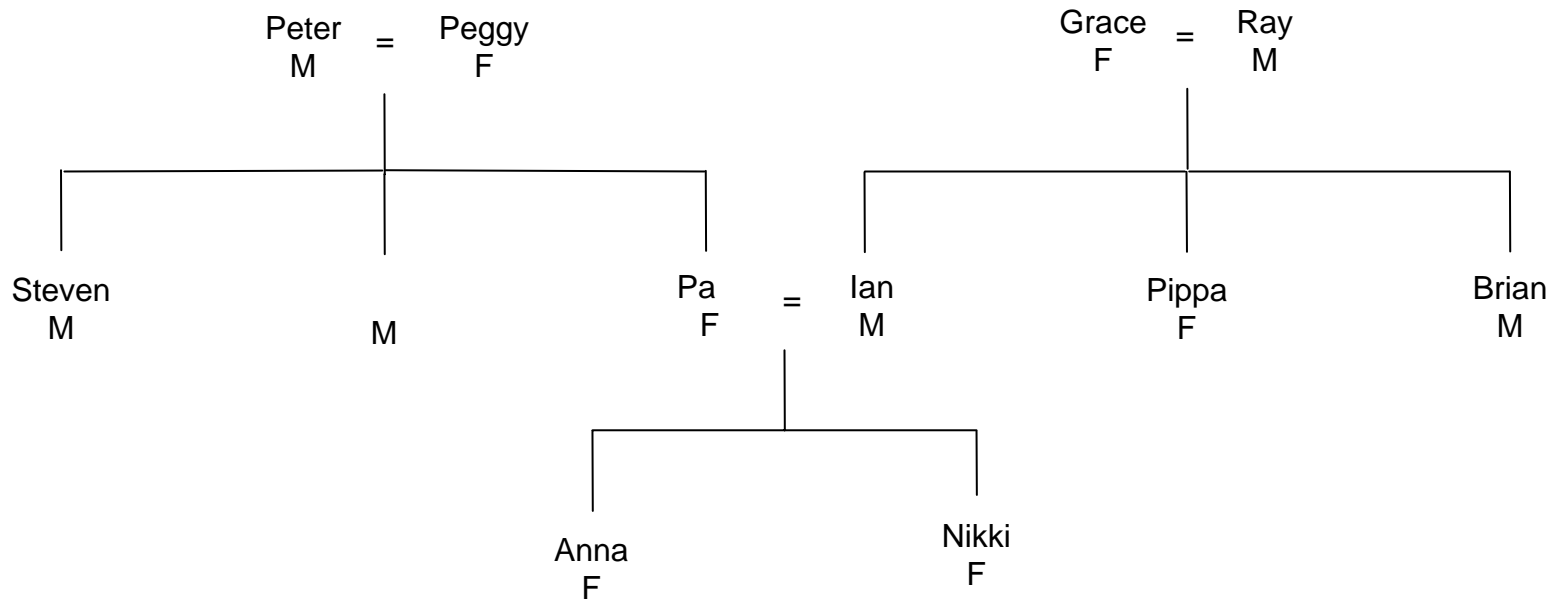
# Iris Data–Clustering Problem

	Sepal Length	Sepal Width	Petal Length	Petal Width
1	5.1	3.5	1.4	0.2
2	4.9	3	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
5	5	3.6	1.4	0.2
...				
51	7	3.2	4.7	1.4
52	6.4	3.2	4.5	1.5
53	6.9	3.1	4.9	1.5
54	5.5	2.3	4	1.3
55	6.5	2.8	4.6	1.5
...				
101	6.3	3.3	6	2.5
102	5.8	2.7	5.1	1.9
103	7.1	3	5.9	2.1
104	6.3	2.9	5.6	1.8
105	6.5	3	5.8	2.2

# Weather Data–Numeric Class

Outlook	Temperature	Humidity	Windy	Play-time
sunny	85	85	false	5
sunny	80	90	true	0
overcast	83	86	false	55
rainy	70	96	false	40
rainy	68	80	false	65
rainy	65	70	true	45
overcast	64	65	true	60
sunny	72	95	false	0
sunny	69	70	false	70
rainy	75	80	false	45
sunny	75	70	true	50
overcast	72	90	true	55
overcast	81	75	false	75
rainy	71	91	true	10

# Family Tree



# Family Tree

First Person	Second Person	Sister-of ?
Peter	Peggy	no
Peter	Steven	no
...	...	
Steven	Peter	no
Steven	Graham	no
Steven	Pam	yes
Steven	Grace	no
...	...	
Ian	Pippa	yes
...	...	
Anna	Nikki	yes
...	...	
Nikki	Anna	yes

First Person	Second Person	Sister-of ?
Steven	Pam	yes
Graham	Pam	yes
Ian	Pippa	yes
Brian	Pippa	yes
Anna	Nikki	yes
Nikki	Anna	yes
All the rest		no

# Family Tree As Table

Name	Gender	Parent1	Parent2
Peter	male	?	?
Peggy	female	?	?
Steven	male	Peter	Peggy
Graham	male	Peter	Peggy
Pam	female	Peter	Peggy
Ian	male	Grace	Ray

# Sister-of As Table

First Person				Second Person				Sister of?
Name	Gender	Parent1	Parent2	Name	Gender	Parent 1	Parent2	
Steven	male	Peter	Peggy	Pam	female	Peter	Peggy	yes
Graham	male	Peter	Peggy	Pam	female	Peter	Peggy	yes
Ian	male	Grace	Ray	Pippa	female	Grace	Ray	yes
Ian	male	Grace	Ray	Pippa	female	Grace	Ray	yes
Annna	female	Pam	Ian	Nikki	female	Pam	Ian	yes
Nikki	female	Pam	Ian	Anna	female	Pam	Ian	yes
All the rest								no

# Another Relationship As Table

First Person				Second Person				Ancestor of?
Name	Gender	Parent1	Parent2	Name	Gender	Parent 1	Parent2	
Peter	male	?	?	Steven	male	Peter	Peggy	yes
Peter	male	?	?	Pam	female	Peter	Peggy	yes
Peter	male	?	?	Anna	female	Pam	Ian	yes
Peter	male	?	?	Nikki	female	Pam	Ian	yes
Pam	female	Peter	Peggy	Nikki	female	Pam	Ian	yes
Grace	female	?	?	Ian	male	Grace	Ray	yes
Grace	female	?	?	Nikki	female	Pam	Ian	yes
Other examples here								yes
All the rest								no

# ARFF File for Weather Data

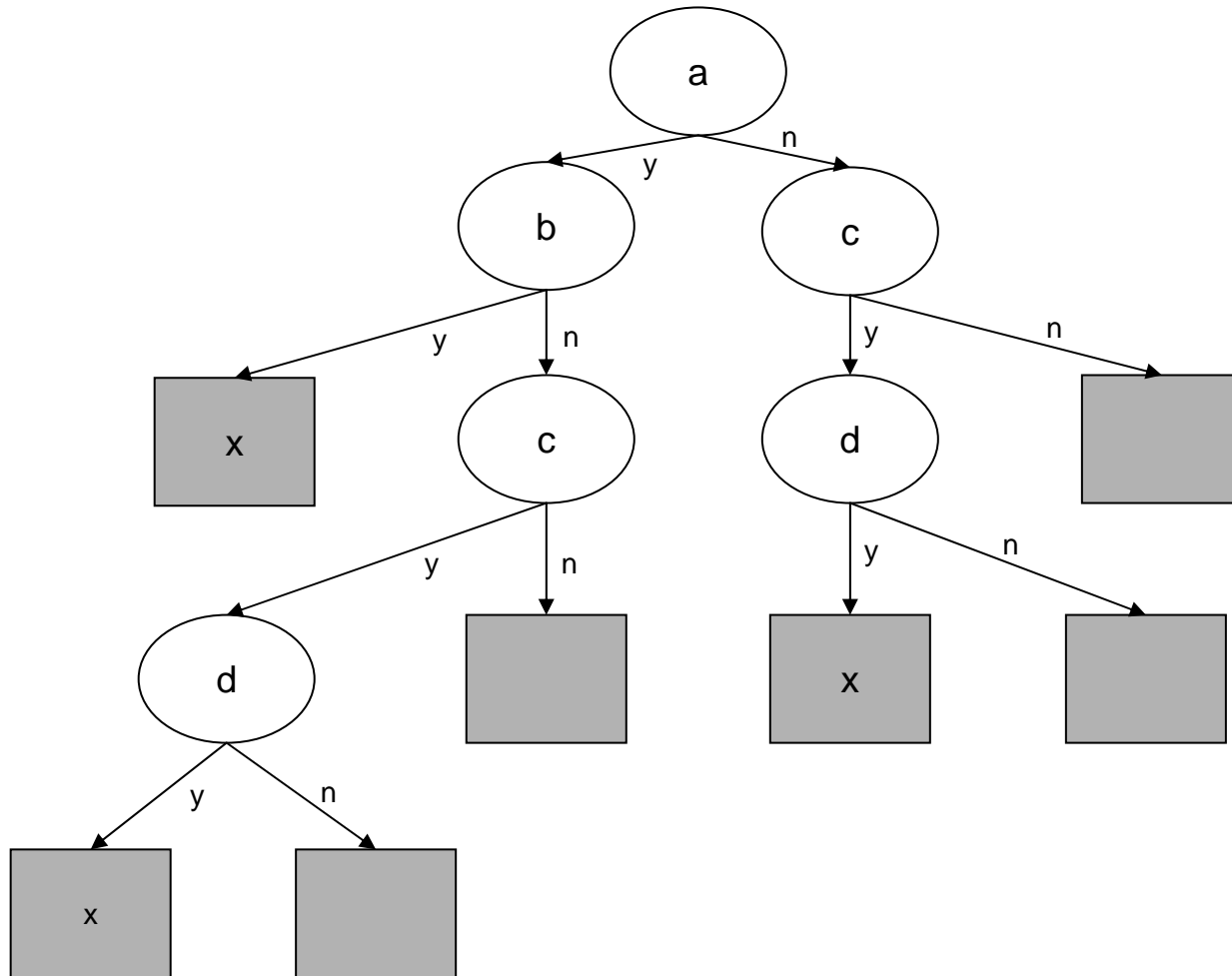
```
% ARFF file for the weather data
with some numeric features
%
@relation weather

@attribute outlook {sunny, overcast,
rainy}
@attribute temperature numeric
@attribute humidity numeric
@attribute windy {true, false}
@attribute play? {yes, no}

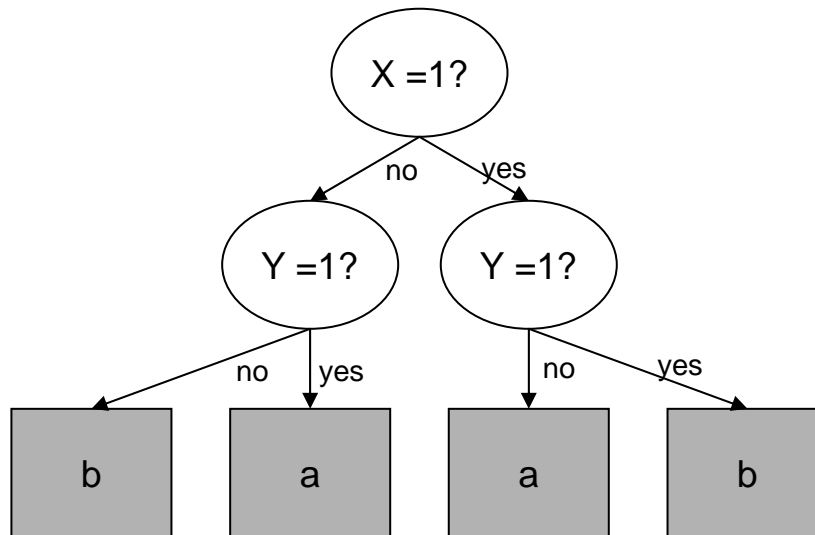
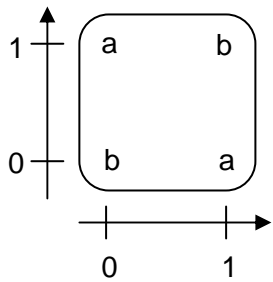
@data
%
%14 instances
%
sunny, 85, 85, false, no
sunny, 80, 90, true, no
overcast, 83, 86, false, yes
```

```
rainy, 70, 96, false, yes
rainy, 68, 80, false, yes
rainy, 65, 70, true, no
overcast, 64, 65, true, yes
sunny, 72, 95, false, no
sunny, 69, 70, false, yes
rainy, 75, 80, false, yes
sunny, 75, 70, true, yes
overcast, 72, 90, true, yes
overcast, 81, 75, false, yes
rainy, 71, 91, true, no
```

# Simple Disjunction



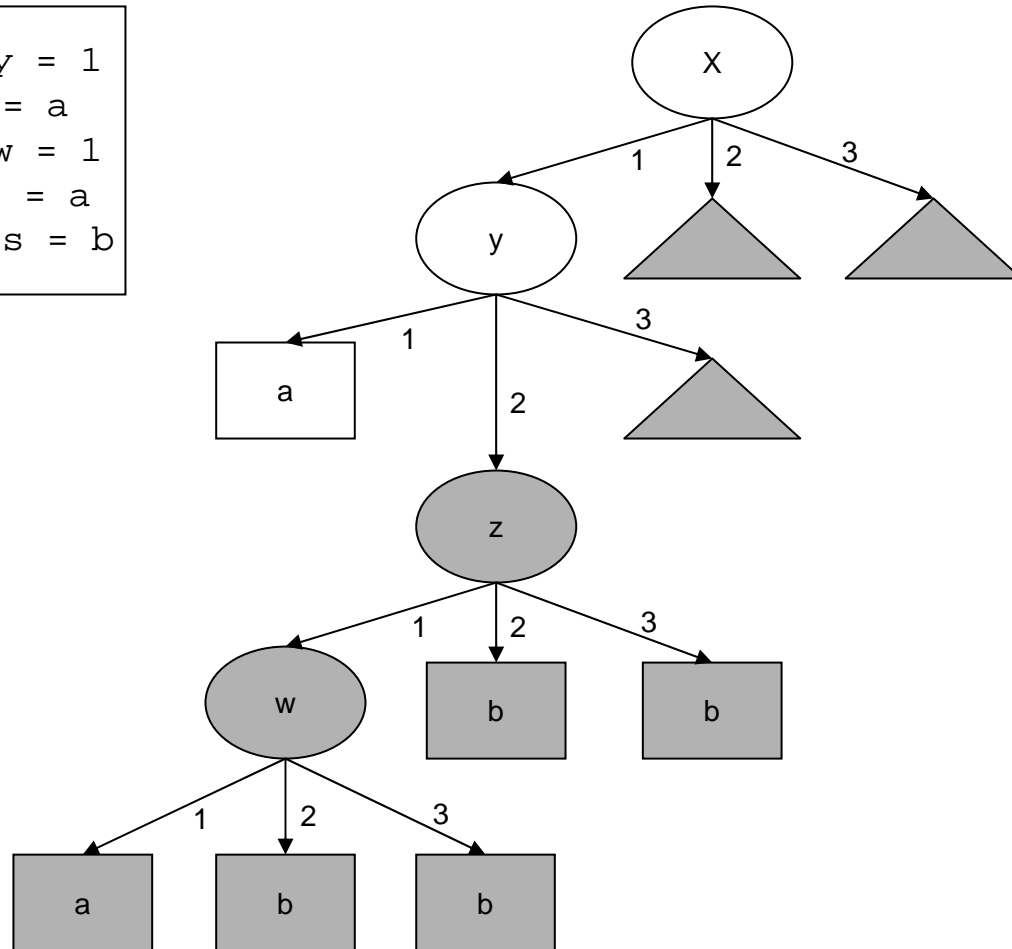
# Exclusive-Or Problem



If  $x = 1$  and  $y = 0$   
then class = a  
If  $x = 0$  and  $y = 1$   
then class = a  
If  $x = 0$  and  $y = 0$   
then class = b  
If  $x = 1$  and  $y = 1$   
then class = b

# Replicated Subtree

If  $x = 1$  and  $y = 1$   
then class = a  
If  $z = 0$  and  $w = 1$   
then class = a  
Otherwise class = b



# New Iris Flower

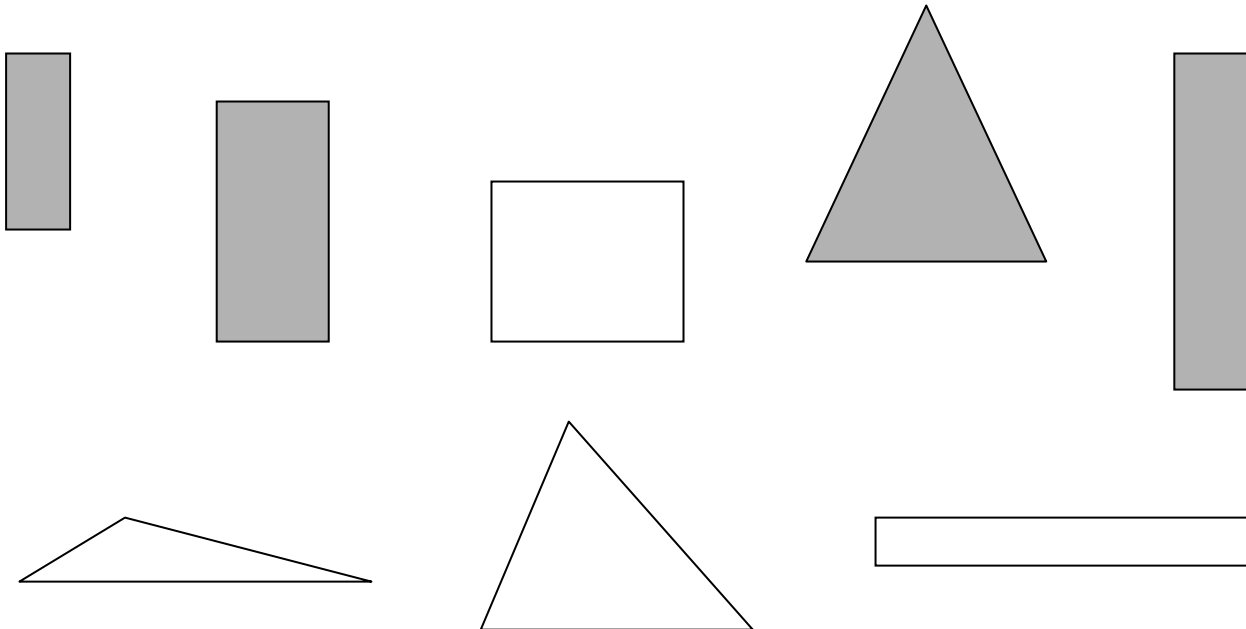
Sepal Length	Sepal Width	Petal Length	Petal Width	Type
5.1	3.5	2.6	0.2	?

# Rules for Iris Data

```
Default: Iris-setosa 1
except if petal-length ≥ 2.45 and petal-length < 5.355 2
    and petal-width < 1.75 3
then Iris-versicolor 4
    except if petal-length ≥ 4.95 and petal-width < 1.55 5
        then Iris-virginica 6
        else if sepal-length < 4.95 and sepal-width ≥ 2.45 7
            then Iris-virginica 8
    else if petal-length ≥ 3.35 9
        then Iris-virginica 10
            except if petal-length < 4.85 and sepal-length < 5.95 11
                then Iris-versicolor 12
```

# The *Shapes* Problem

Shaded: Standing  
Unshaded: Lying



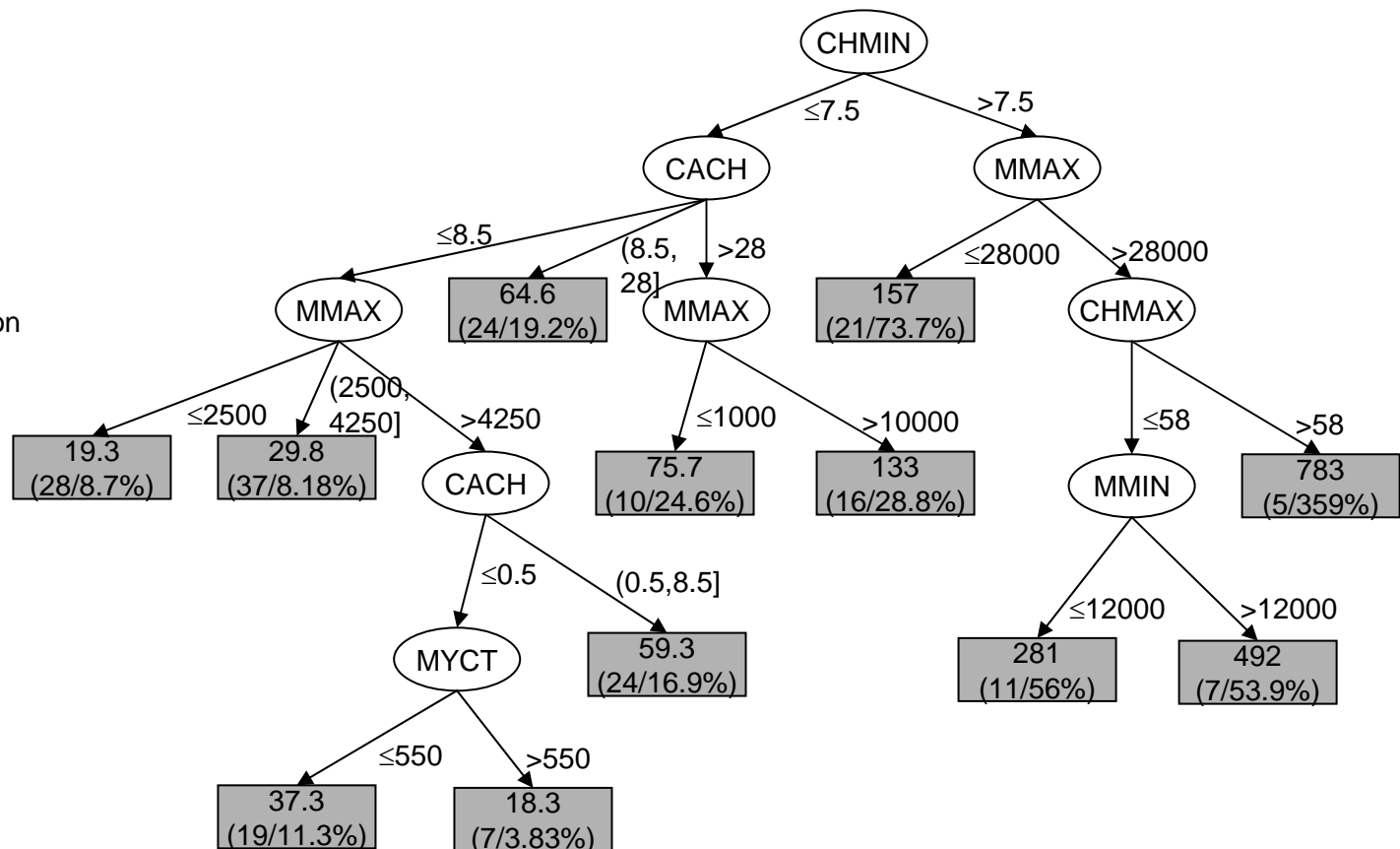
# Training Data for *Shapes* Problem

Width	Height	Sides	Class
2	4	4	standing
3	6	4	standing
4	3	4	lying
7	8	3	standing
7	6	3	lying
2	9	4	standing
9	1	4	lying
10	2	3	lying

# CPU Performance Data

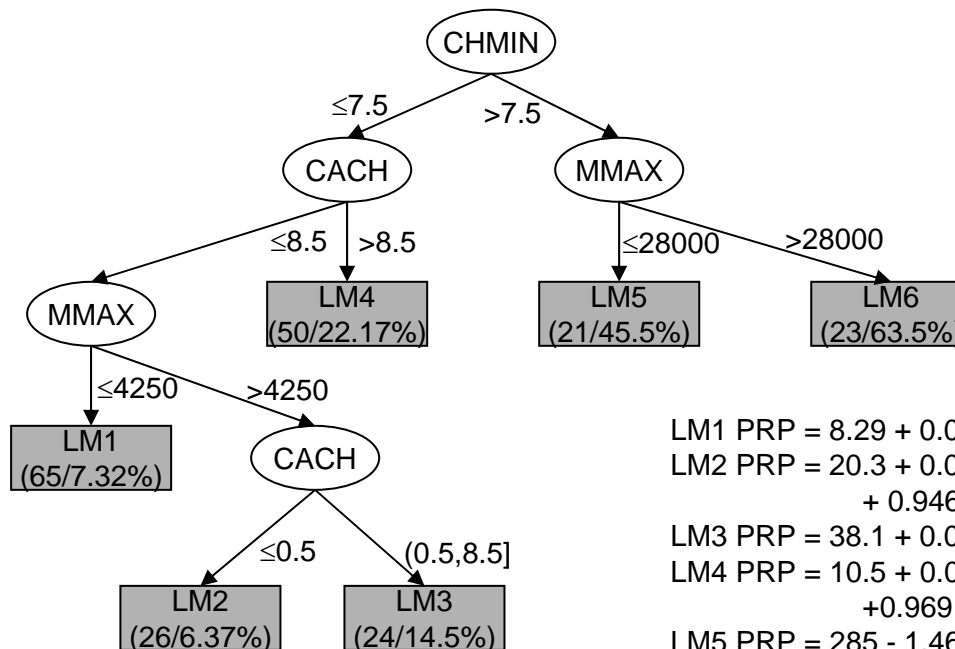
PRP =  
 -56.1  
 +0.049 MYCT  
 +0.015 MMIN  
 +0.006MMAX  
 +0.630CACH  
 -0.270CHMIN  
 +1.46 CHMAX

(a) linear regression



(b) regression tree

# CPU Performance Data



$$\text{LM1 PRP} = 8.29 + 0.004 \text{ MMAX} + 2.77 \text{ CHMIN}$$

$$\text{LM2 PRP} = 20.3 + 0.004 \text{ MMIN} - 3.99 \text{ CHMIN} + 0.946 \text{ CHMAX}$$

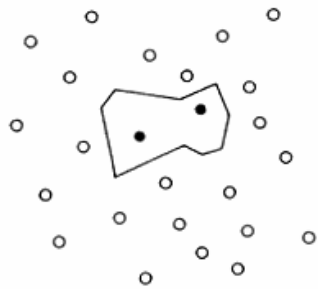
$$\text{LM3 PRP} = 38.1 + 0.012 \text{ MMIN}$$

$$\text{LM4 PRP} = 10.5 + 0.002 \text{ MMAX} + 0.698 \text{ CACH} + 0.969 \text{ CHMAX}$$

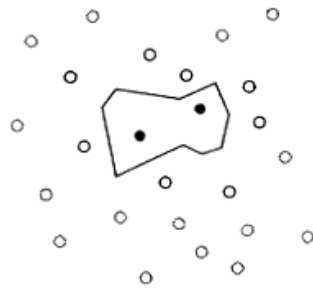
$$\text{LM5 PRP} = 285 - 1.46 \text{ MYCT} + 1.02 \text{ CACH} - 9.39 \text{ CHMIN}$$

$$\text{LM6 PRP} = -65.8 + 0.03 \text{ MMIN} - 2.94 \text{ CHMIN} = 4.98 \text{ CHMAX}$$

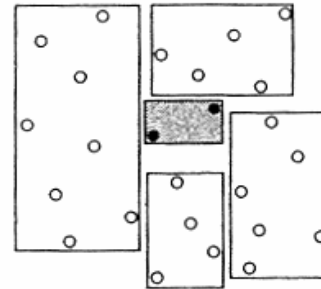
# Partitioning Instance Space



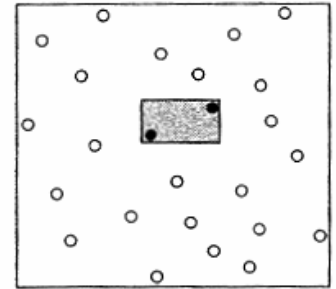
(a)



(b)

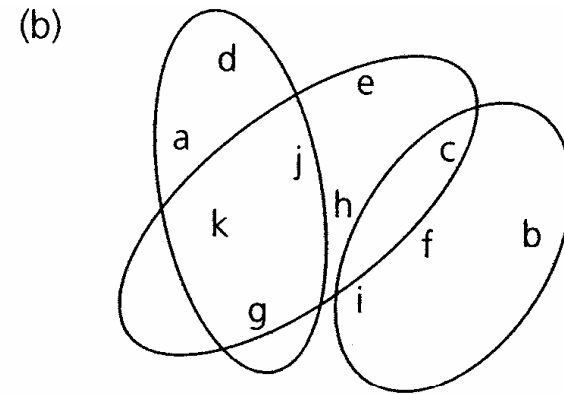
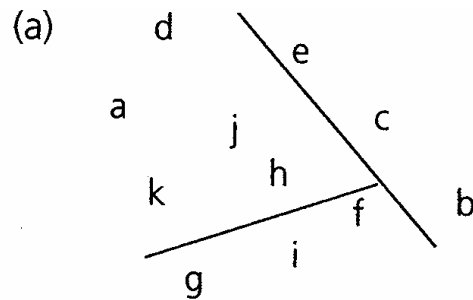


(c)



(d)

# Ways to Represent Clusters



(c)

	1	2	3
a	0.4	0.1	0.5
b	0.1	0.8	0.1
c	0.3	0.3	0.4
d	0.1	0.1	0.8
e	0.4	0.2	0.4
f	0.1	0.4	0.5
g	0.7	0.2	0.1
h	0.5	0.4	0.1
...			

