

---

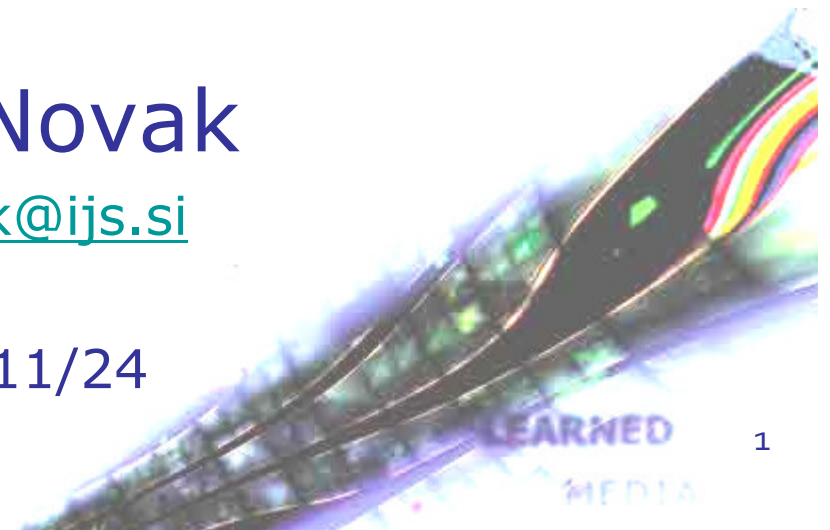
# Data Mining and Knowledge Discovery

## Knowledge Discovery and Knowledge Management in e-Science

Petra Kralj Novak

[Petra.Kralj.Novak@ijs.si](mailto:Petra.Kralj.Novak@ijs.si)

Practice, 2009/11/24



# Practice plan

---

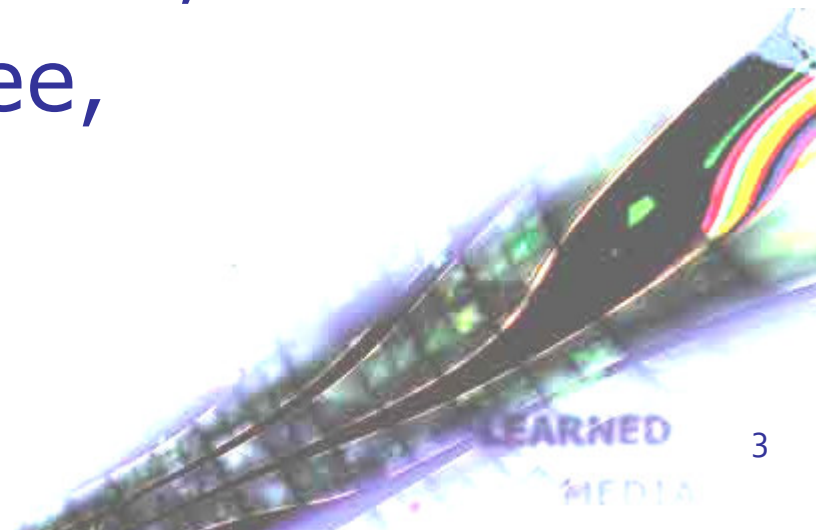
- 2009/11/10: Predictive data mining
  - Decision trees
  - Naïve Bayes classifier
  - Evaluating classifiers (separate test set, cross validation, confusion matrix, classification accuracy)
  - Predictive data mining in Weka
- 2009/11/24: Numeric prediction and descriptive data mining
  - Numeric prediction
  - Association rules
  - Regression models and evaluation in Weka
  - Descriptive data mining in Weka
  - Discussion about seminars and exam
- 2009/12/8: Written exam
- 2010/1/26: Seminar proposal presentations
- 2009/3/1: deadline for data mining papers (written seminar)
- 2009/3/3: Data mining seminar presentations



---

# Numeric prediction

Baseline,  
Linear Regression,  
Regression tree,  
Model Tree,  
KNN

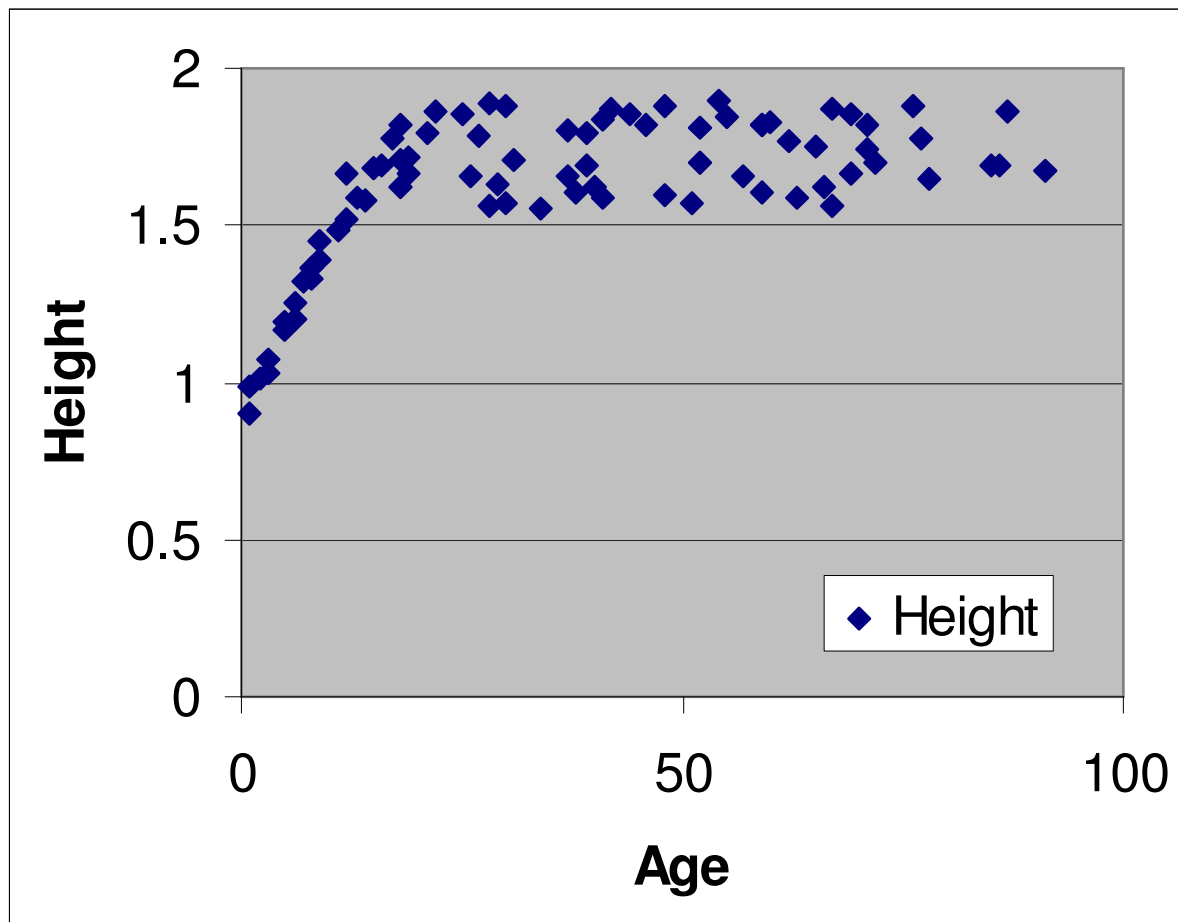


<b>Numeric prediction</b>	<b>Classification</b>
<b>Data:</b> attribute-value description	
<b>Target variable:</b> Continuous	<b>Target variable:</b> Categorical (nominal)
<b>Evaluation:</b> cross validation, separate test set, ...	
<b>Error:</b> MSE, MAE, RMSE, ...	<b>Error:</b> 1-accuracy
<b>Algorithms:</b> Linear regression, regression trees,...	<b>Algorithms:</b> Decision trees, Naïve Bayes, ...
<b>Baseline predictor:</b> Mean of the target variable	<b>Baseline predictor:</b> Majority class

# Example

---

- data about 80 people:  
Age and Height

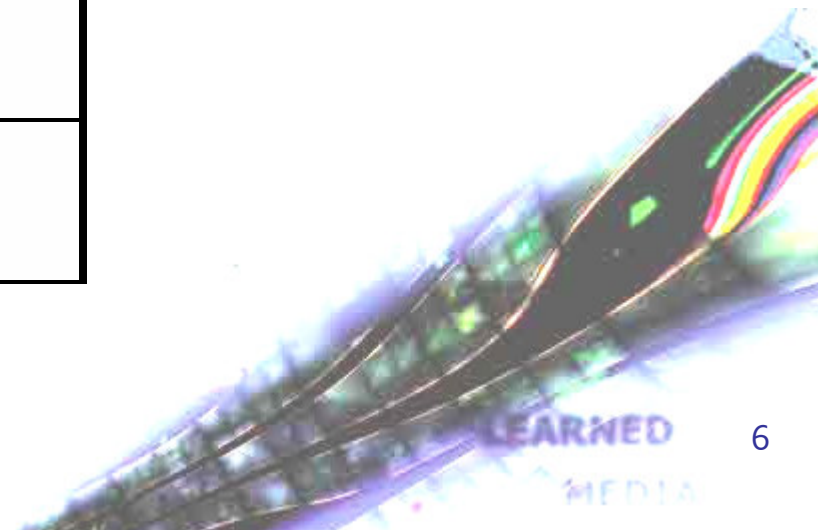


Age	Height
3	1.03
5	1.19
6	1.26
9	1.39
15	1.69
19	1.67
22	1.86
25	1.85
41	1.59
48	1.60
54	1.90
71	1.82
...	...

# Test set

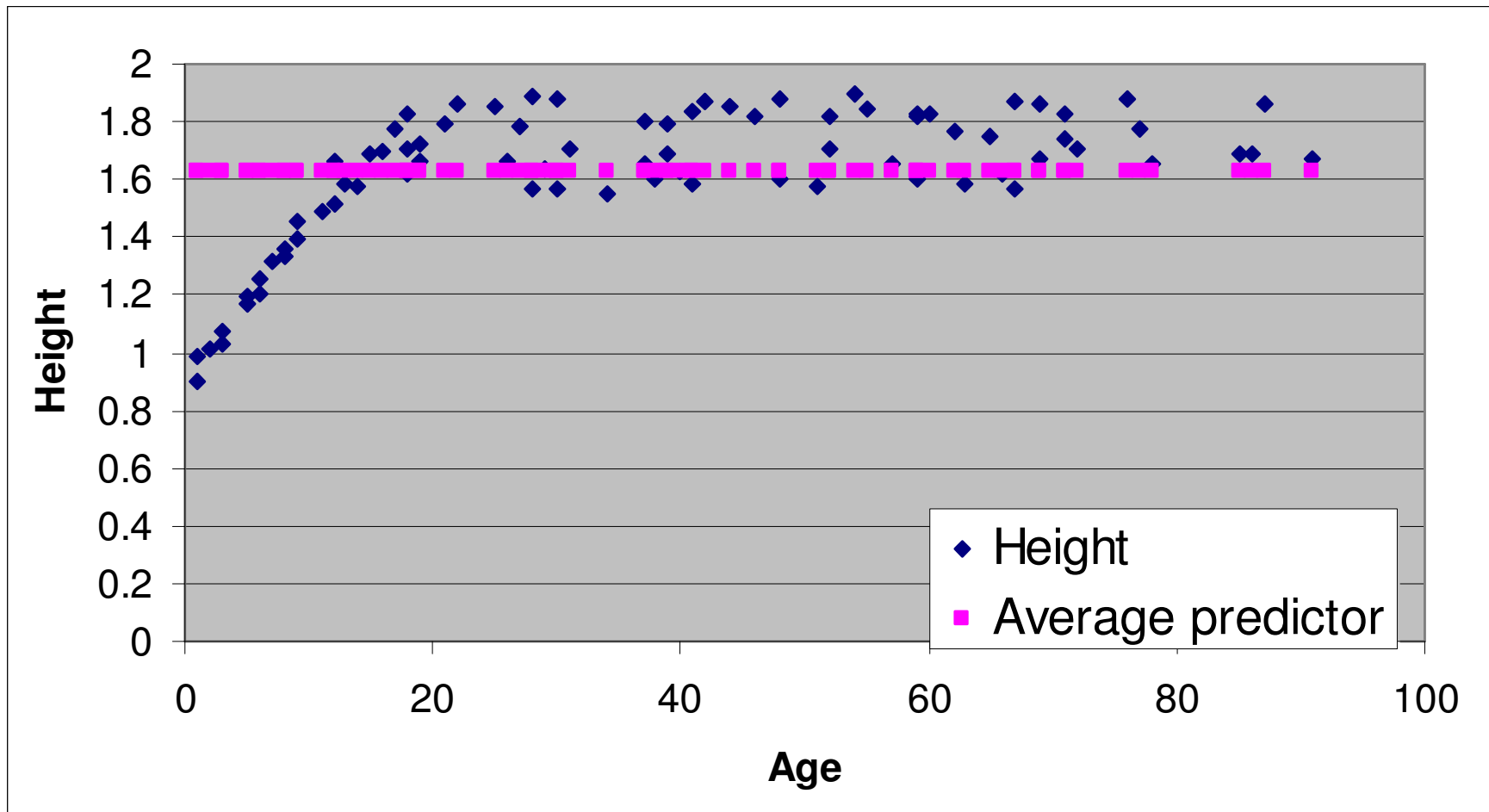
---

Age	Height
2	0.85
10	1.4
35	1.7
70	1.6



# Baseline numeric predictor

- Average of the target variable

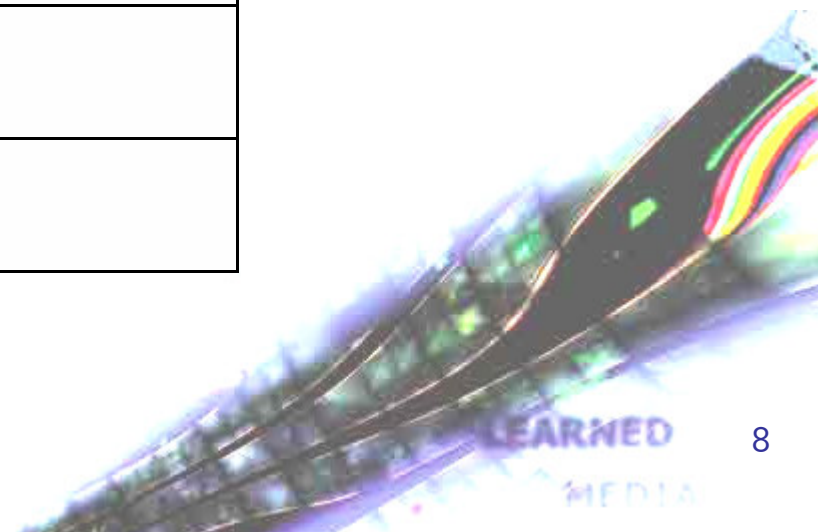


# Baseline predictor: prediction

---

Average of the target variable is 1.63

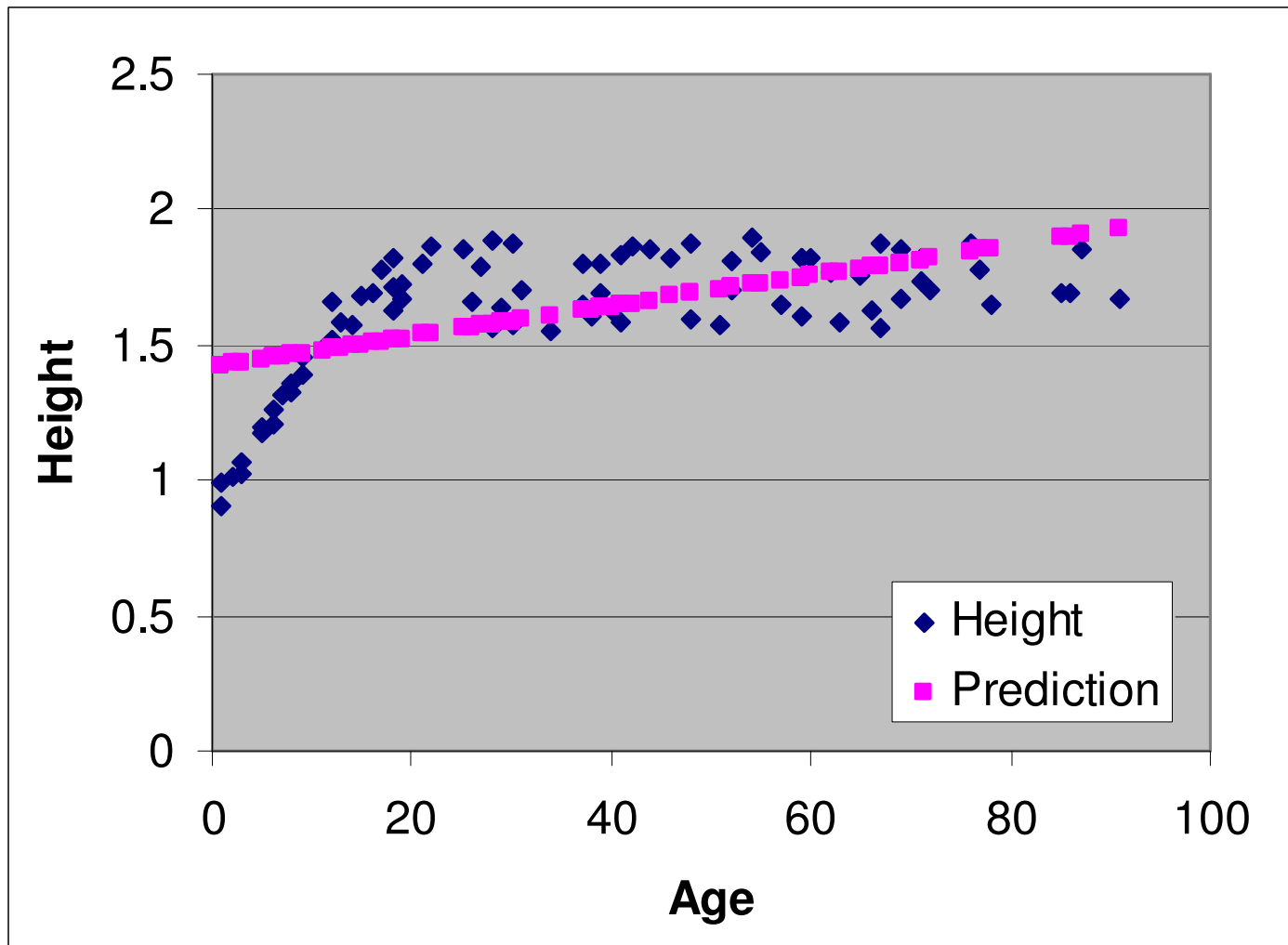
Age	Height	Baseline
2	0.85	
10	1.4	
35	1.7	
70	1.6	





# Linear Regression Model

$$\text{Height} = 0.0056 * \text{Age} + 1.4181$$

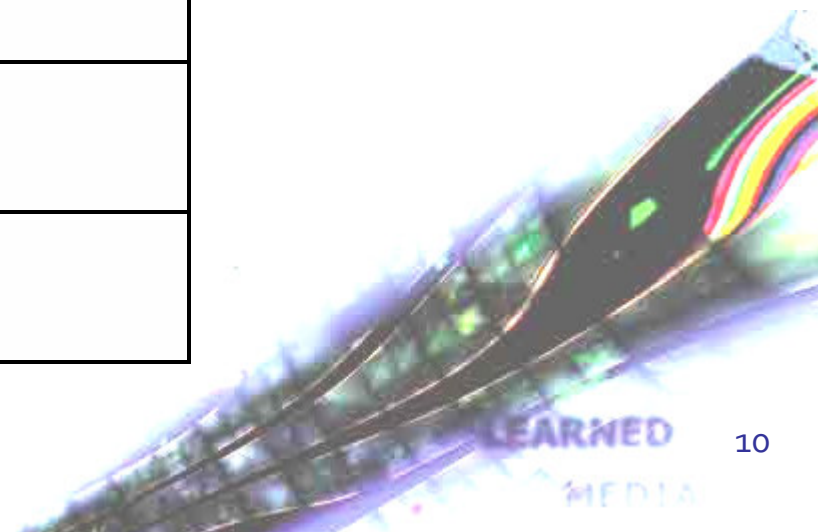


# Linear Regression: prediction

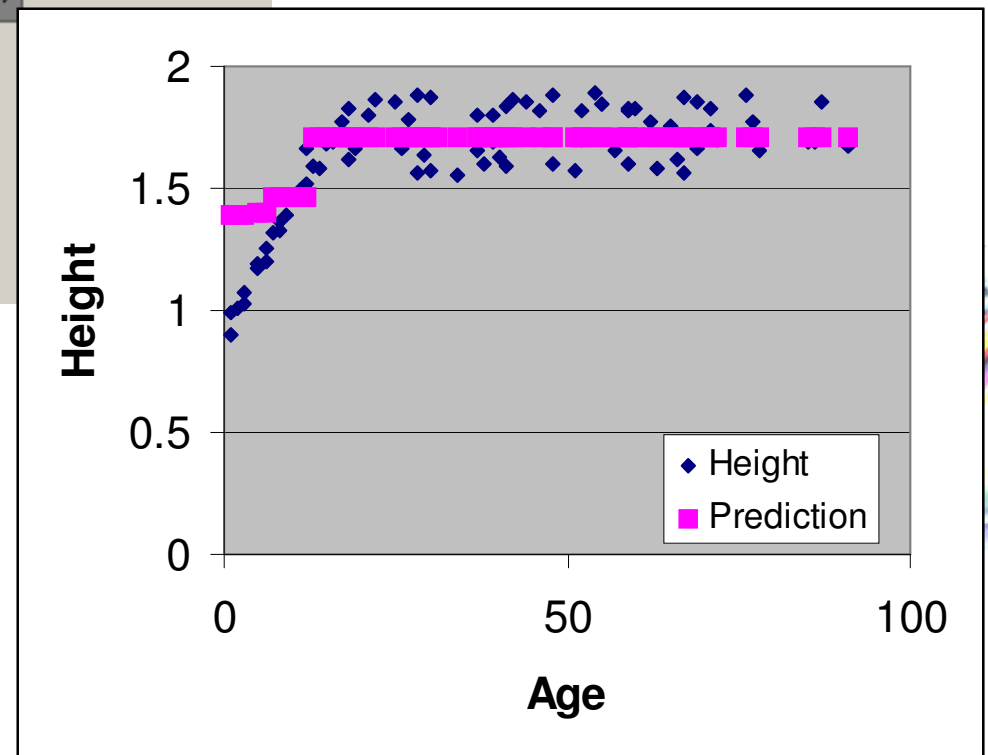
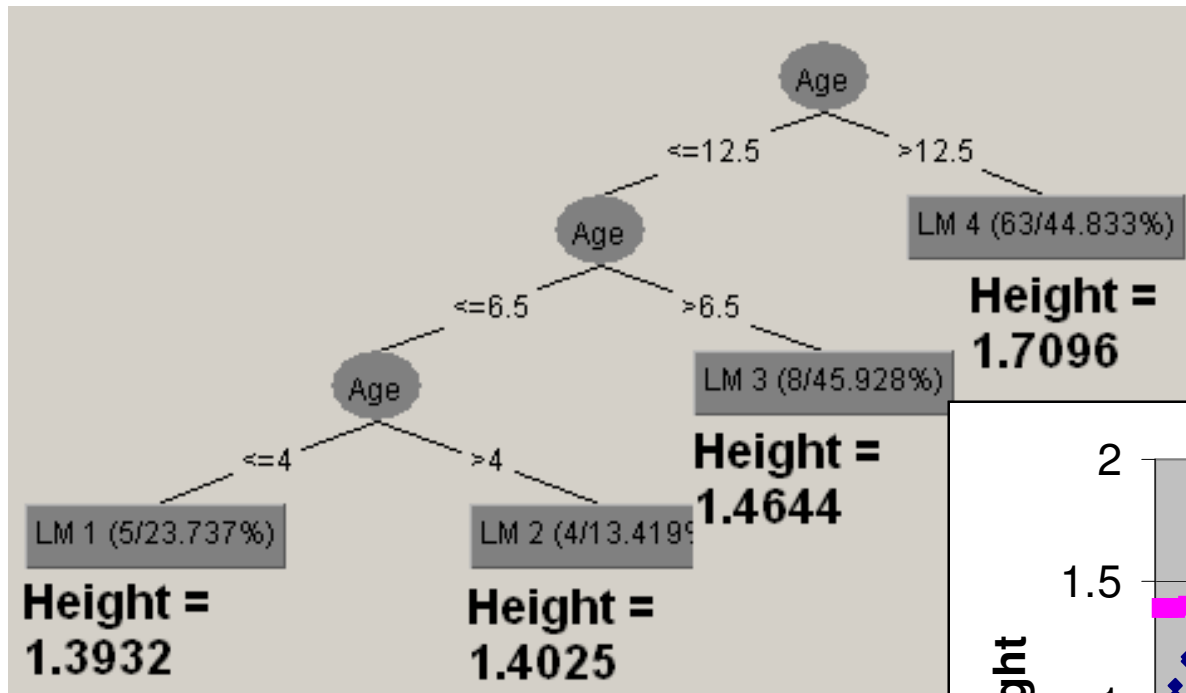
---

$$\text{Height} = 0.0056 * \text{Age} + 1.4181$$

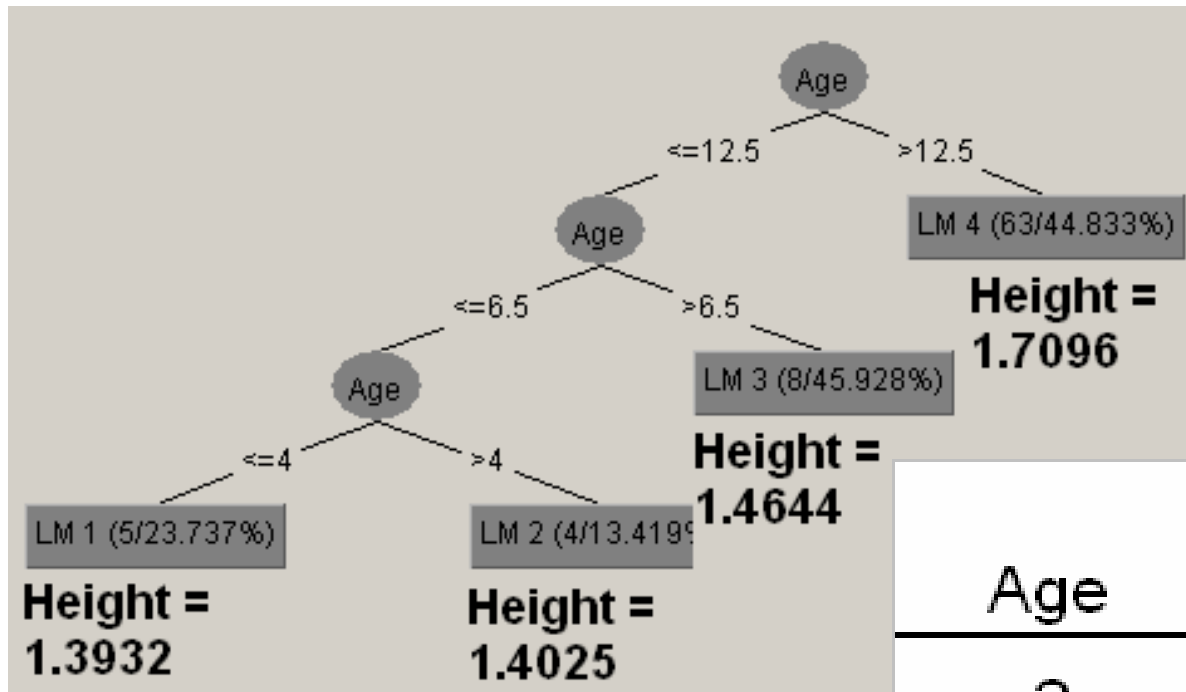
Age	Height	Linear regression
2	0.85	
10	1.4	
35	1.7	
70	1.6	



# Regression tree

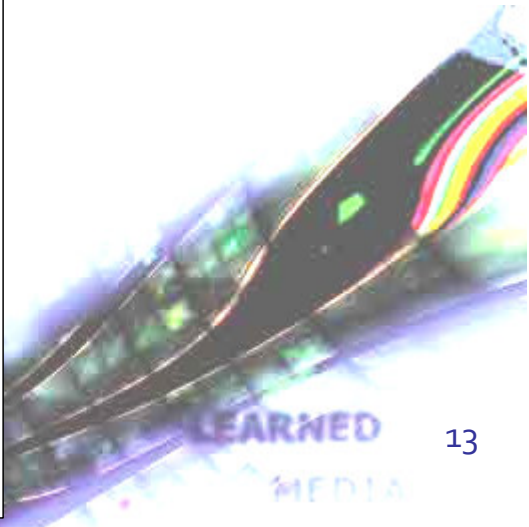
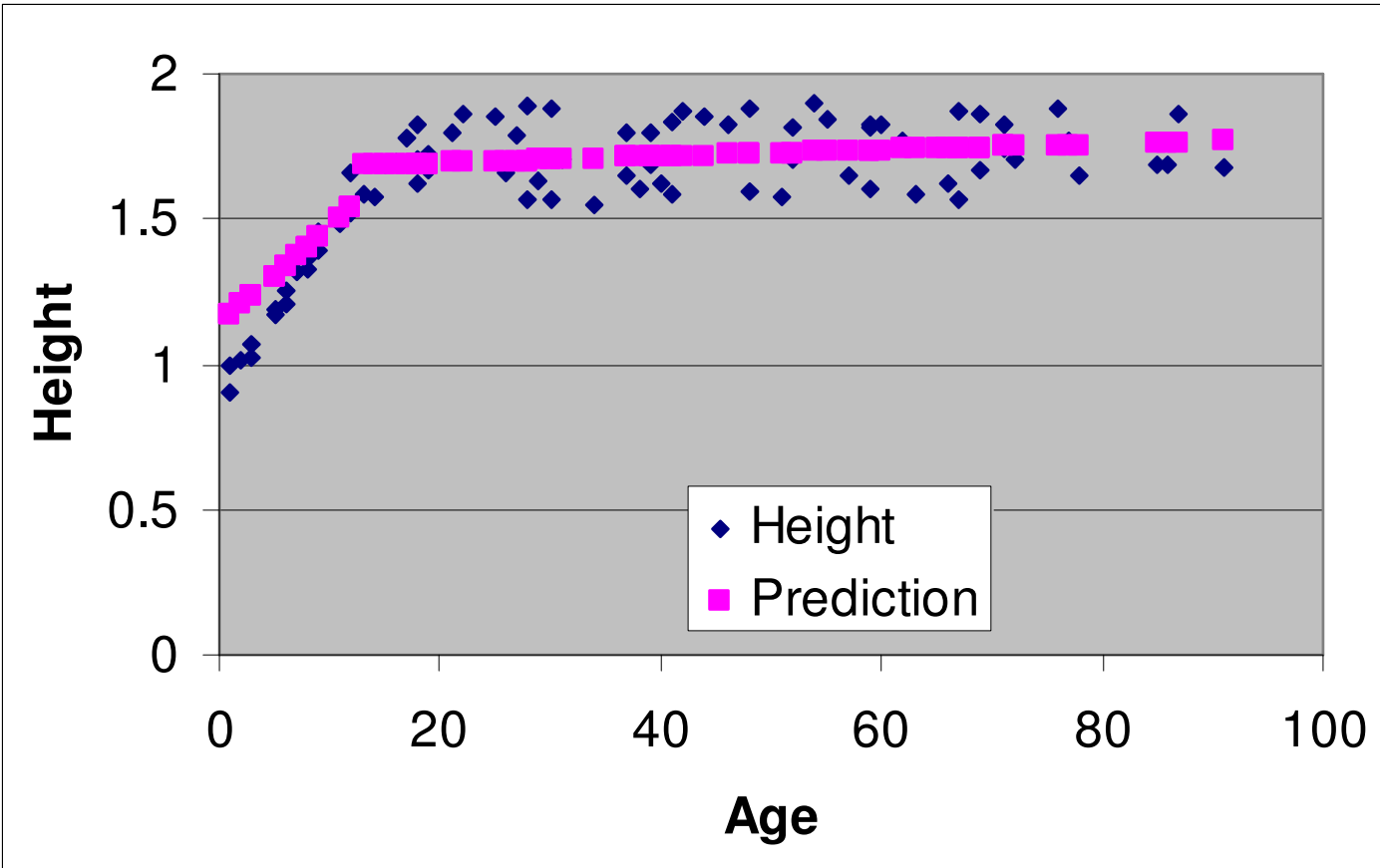


# Regression tree: prediction



Age	Height	Regression tree
2	0.85	
10	1.4	
35	1.7	
70	1.6	

# Model tree



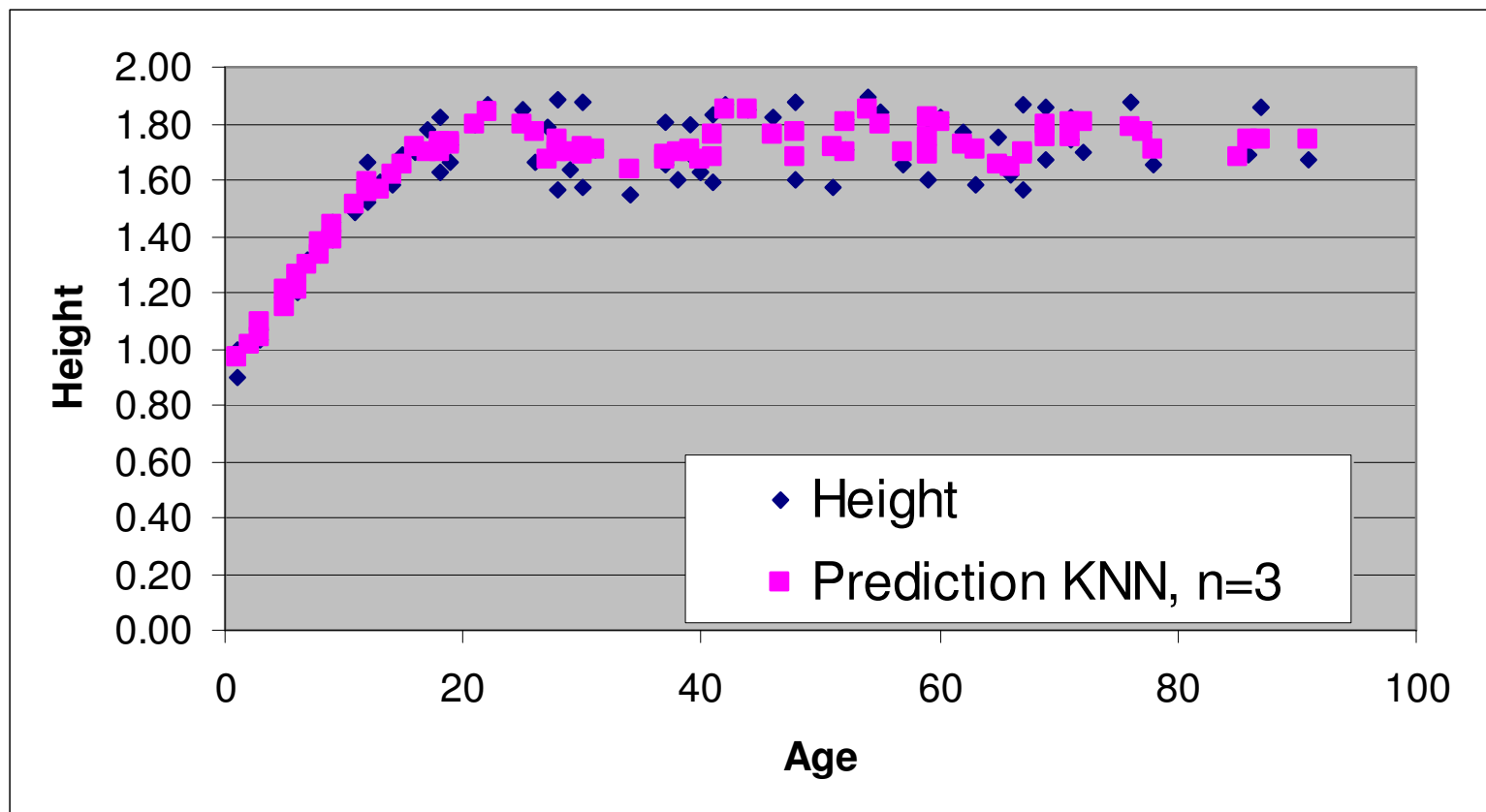
# Model tree: prediction

Age	Height	Model tree
2	0.85	
10	1.4	
35	1.7	
70	1.6	



# KNN – K nearest neighbors

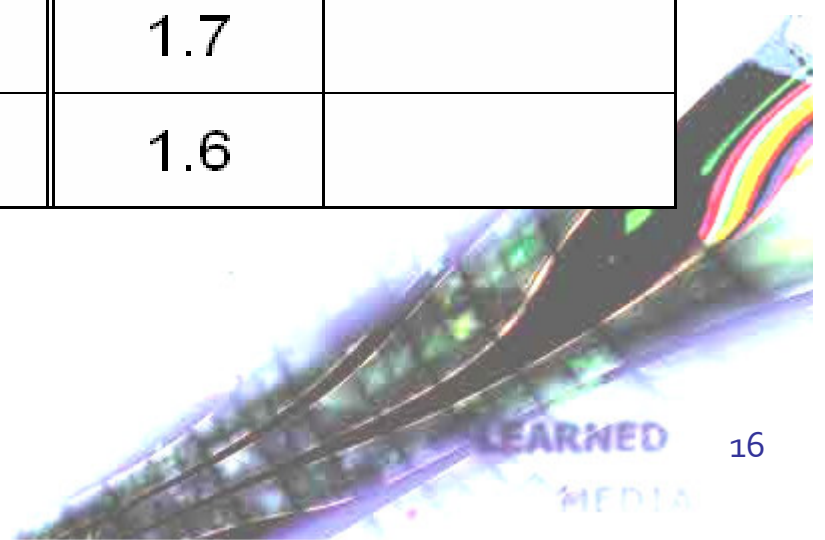
- Looks at K closest examples (by non-target attributes) and predicts the average of their target variable
- In this example,  $K=3$



# KNN prediction

Age	Height
1	0.90
1	0.99
2	1.01
3	1.03
3	1.07
5	1.19
5	1.17

Age	Height	kNN
2	0.85	
10	1.4	
35	1.7	
70	1.6	

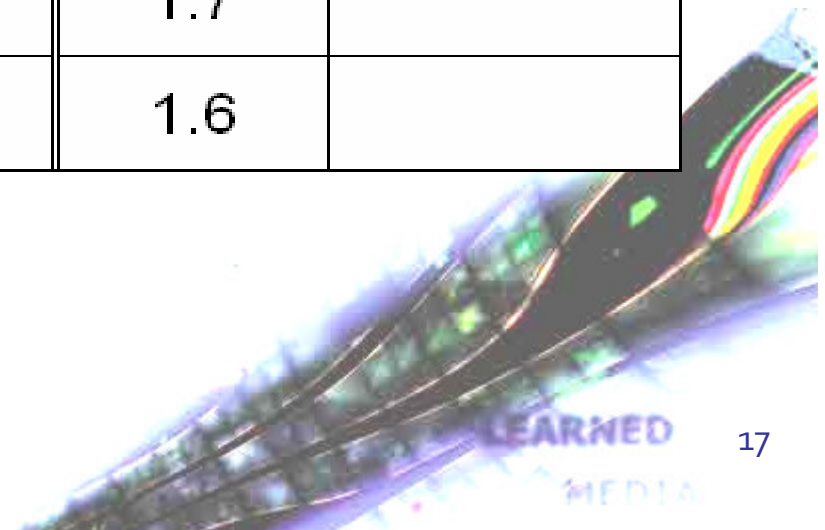




# KNN prediction

Age	Height
8	1.36
8	1.33
9	1.45
9	1.39
11	1.49
12	1.66
12	1.52
13	1.59
14	1.58

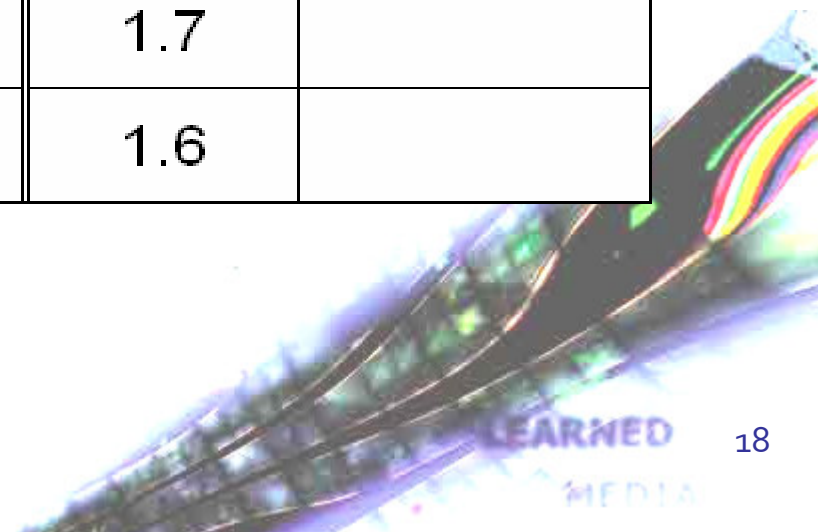
Age	Height	kNN
2	0.85	
10	1.4	
35	1.7	
70	1.6	



# KNN prediction

Age	Height
30	1.57
30	1.88
31	1.71
34	1.55
37	1.65
37	1.80
38	1.60
39	1.69
39	1.80

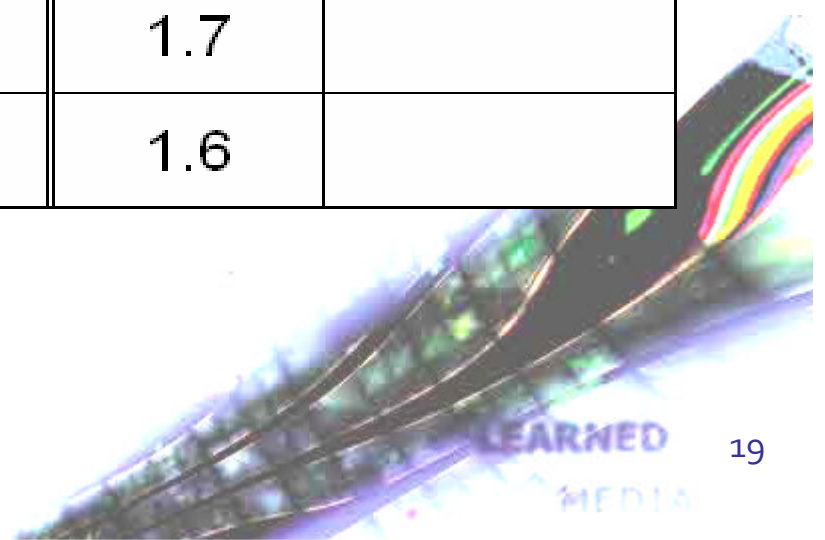
Age	Height	kNN
2	0.85	
10	1.4	
35	1.7	
70	1.6	



# KNN prediction

Age	Height
67	1.56
67	1.87
69	1.67
69	1.86
71	1.74
71	1.82
72	1.70
76	1.88

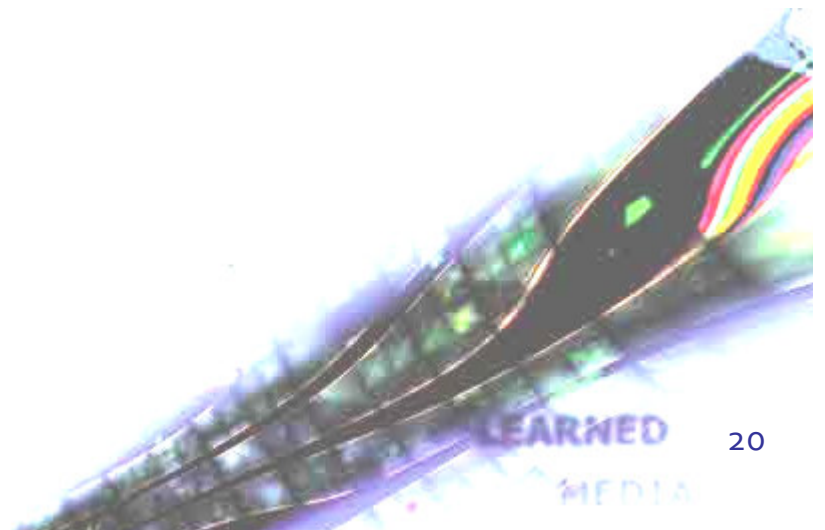
Age	Height	kNN
2	0.85	
10	1.4	
35	1.7	
70	1.6	



# Which predictor is the best?

---

Age	Height	Baseline	Linear regression	Regression tree	Model tree	kNN
2	0.85	1.63	1.43	1.39	1.20	1.00
10	1.4	1.63	1.47	1.46	1.47	1.44
35	1.7	1.63	1.61	1.71	1.71	1.67
70	1.6	1.63	1.81	1.71	1.75	1.77



# Evaluating numeric prediction

Performance measure	Formula
mean-squared error	$\frac{(p_1 - a_1)^2 + \dots + (p_n - a_n)^2}{n}$
root mean-squared error	$\sqrt{\frac{(p_1 - a_1)^2 + \dots + (p_n - a_n)^2}{n}}$
mean absolute error	$\frac{ p_1 - a_1  + \dots +  p_n - a_n }{n}$
relative squared error	$\frac{(p_1 - a_1)^2 + \dots + (p_n - a_n)^2}{(a_1 - \bar{a})^2 + \dots + (a_n - \bar{a})^2}, \text{ where } \bar{a} = \frac{1}{n} \sum_i a_i$
root relative squared error	$\sqrt{\frac{(p_1 - a_1)^2 + \dots + (p_n - a_n)^2}{(a_1 - \bar{a})^2 + \dots + (a_n - \bar{a})^2}}$
relative absolute error	$\frac{ p_1 - a_1  + \dots +  p_n - a_n }{ a_1 - \bar{a}  + \dots +  a_n - \bar{a} }$
correlation coefficient	$\frac{S_{PA}}{\sqrt{S_p S_A}}, \text{ where } S_{PA} = \frac{\sum_i (p_i - \bar{p})(a_i - \bar{a})}{n-1},$ $S_p = \frac{\sum_i (p_i - \bar{p})^2}{n-1}, \text{ and } S_A = \frac{\sum_i (a_i - \bar{a})^2}{n-1}$

# Numeric prediction discussion

---

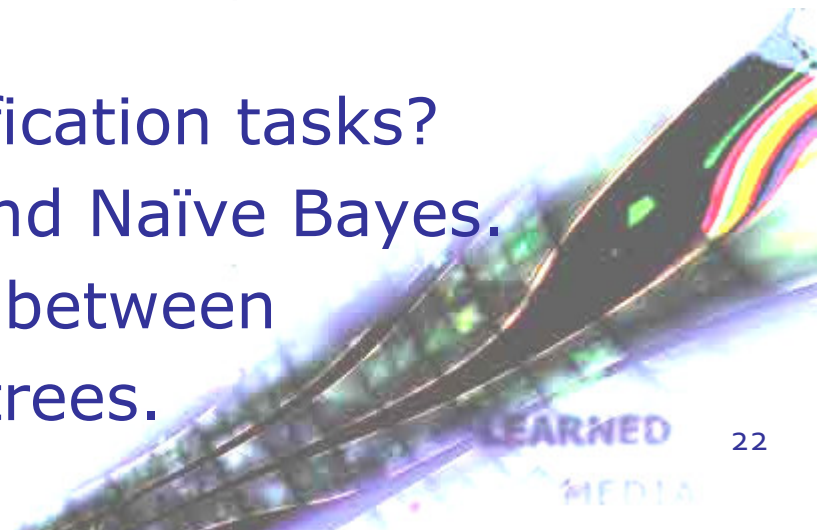
- Consider a dataset with a target variable with five possible values:

1. non sufficient
2. sufficient
3. good
4. very good
5. excellent



- Is this a classification or a numeric prediction problem?
- What if such a variable is an attribute, is it nominal or numeric?

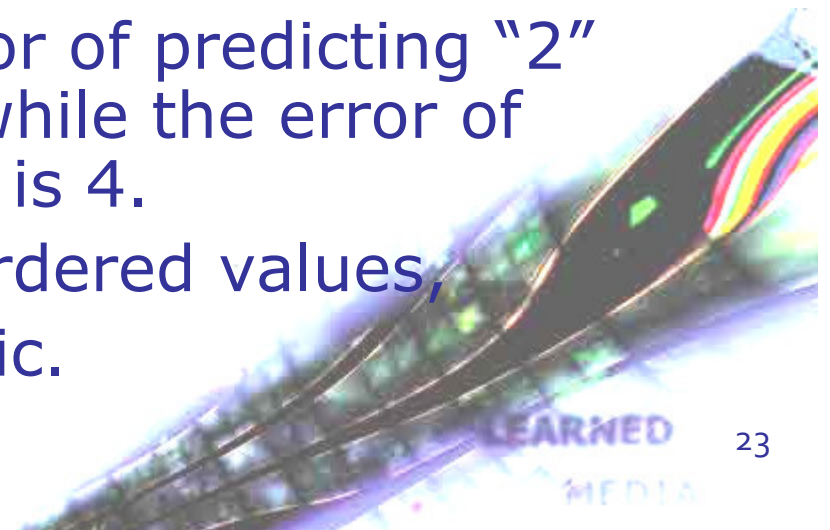
- Can KNN be used for classification tasks?
- Similarities between KNN and Naïve Bayes.
- Similarities and differences between decision trees and regression trees.



# Classification or a numeric prediction problem?

---

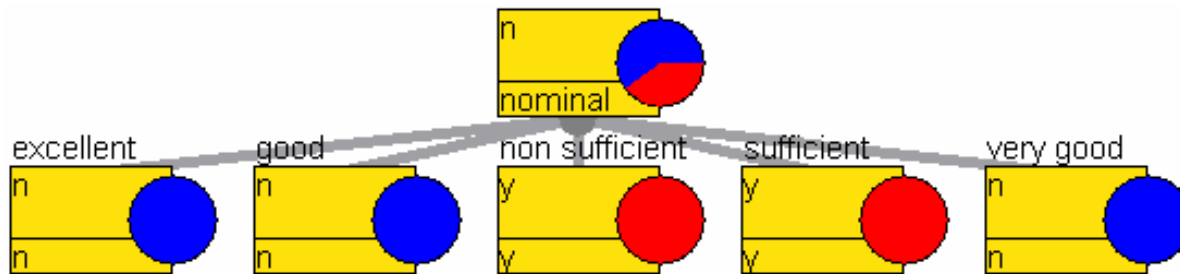
- Target variable with five possible values:
  - 1.non sufficient
  - 2.sufficient
  - 3.good
  - 4.very good
  - 5.excellent
- Classification: the **misclassification cost** is the same if "non sufficient" is classified as "sufficient" or if it is classified as "very good"
- Numeric prediction: The error of predicting "2" when it should be "1" is 1, while the error of predicting "5" instead of "1" is 4.
- If we have a variable with ordered values, it should be considered numeric.



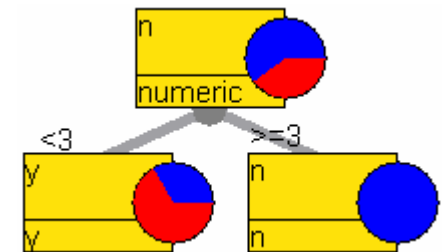
# Nominal or numeric attribute?

- A variable with five possible values:
  - 1.non sufficient
  - 2.sufficient
  - 3.good
  - 4.very good
  - 5.excellent

Nominal:



Numeric:



- If we have a variable with **ordered** values, it should be considered numeric.



# Numeric prediction discussion

---

- Consider a dataset with a target variable with five possible values:
  1. non sufficient
  2. sufficient
  3. good
  4. very good
  5. excellent
  - Is this a classification or a numeric prediction problem?
  - What if such a variable is an attribute, is it nominal or numeric?

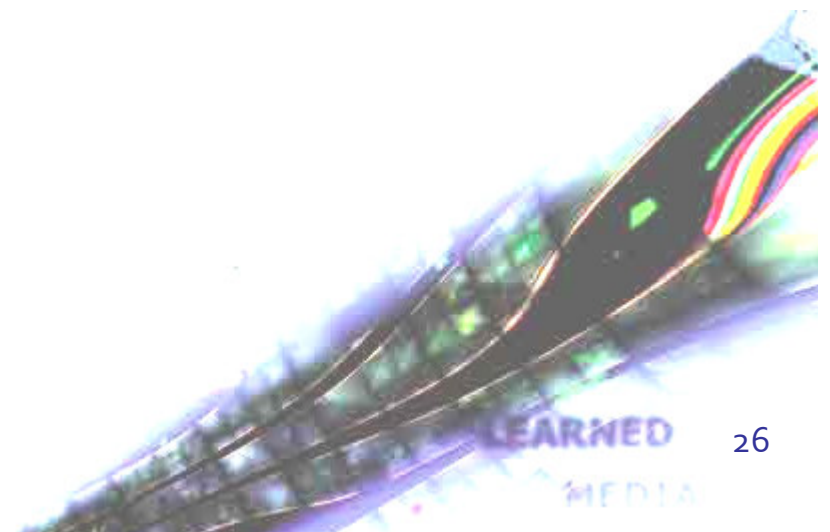
- • Can KNN be used for classification tasks?
- Similarities between KNN and Naïve Bayes.
- Similarities and differences between decision trees and regression trees.



# Can KNN be used for classification tasks?

---

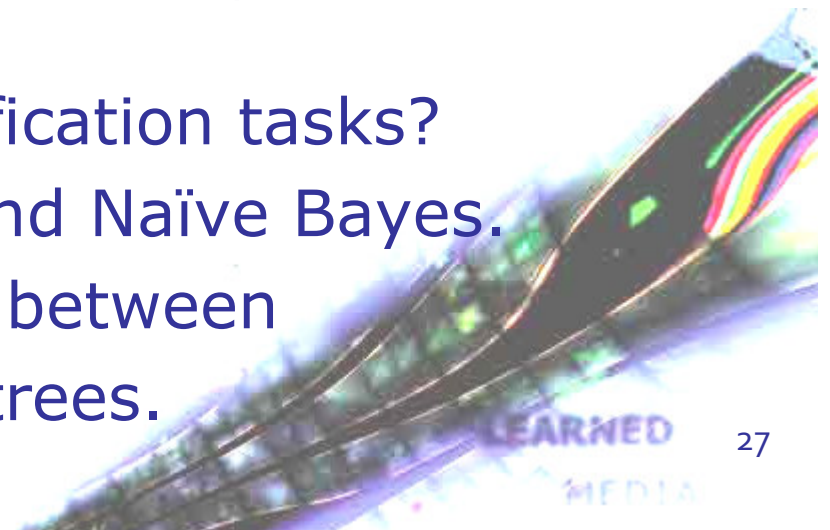
- **YES.**
- In numeric prediction tasks, the average of the neighborhood is computed
- In classification tasks, the distribution of the classes in the neighborhood is computed



# Numeric prediction discussion

---

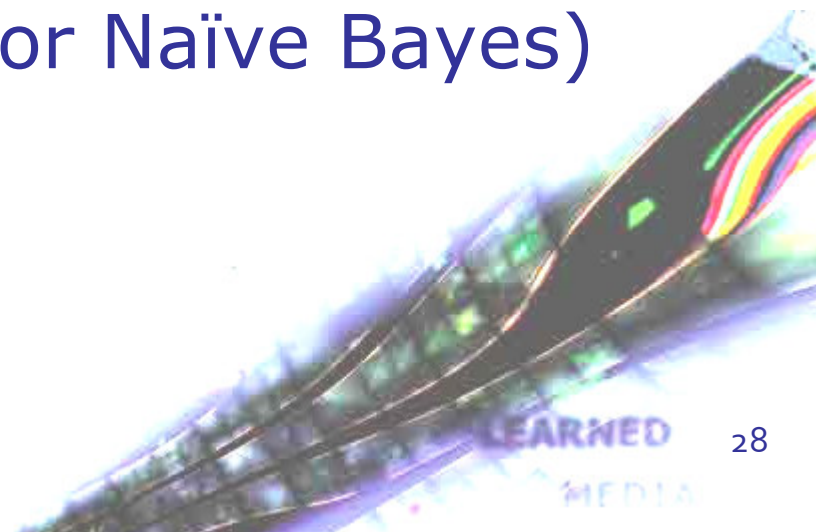
- Consider a dataset with a target variable with five possible values:
  1. non sufficient
  2. sufficient
  3. good
  4. very good
  5. excellent
  - Is this a classification or a numeric prediction problem?
  - What if such a variable is an attribute, is it nominal or numeric?
- Can KNN be used for classification tasks?
- • Similarities between KNN and Naïve Bayes.
- Similarities and differences between decision trees and regression trees.



# Similarities between KNN and Naïve Bayes.

---

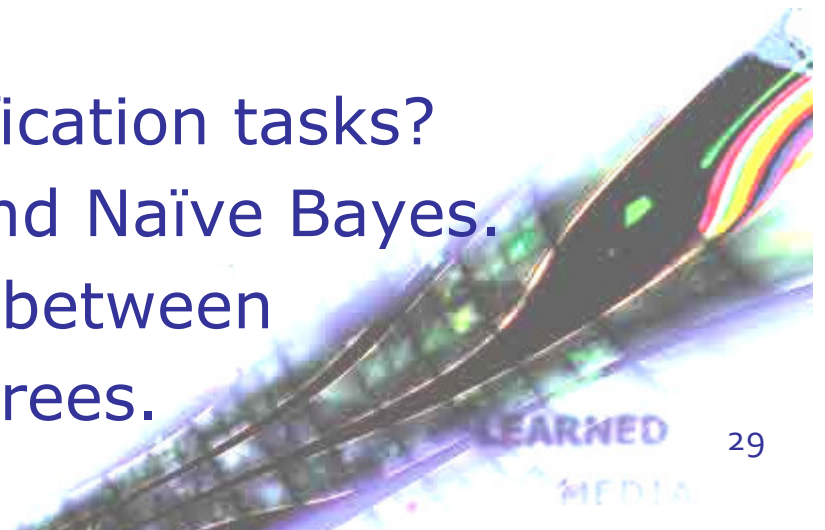
- Both are “**black box**” models, which do not give the insight into the data.
- Both are “**lazy classifiers**”: they do not build a model in the training phase and use it for predicting, but they need the data when predicting the value for a new example (partially true for Naïve Bayes)



# Numeric prediction discussion

---

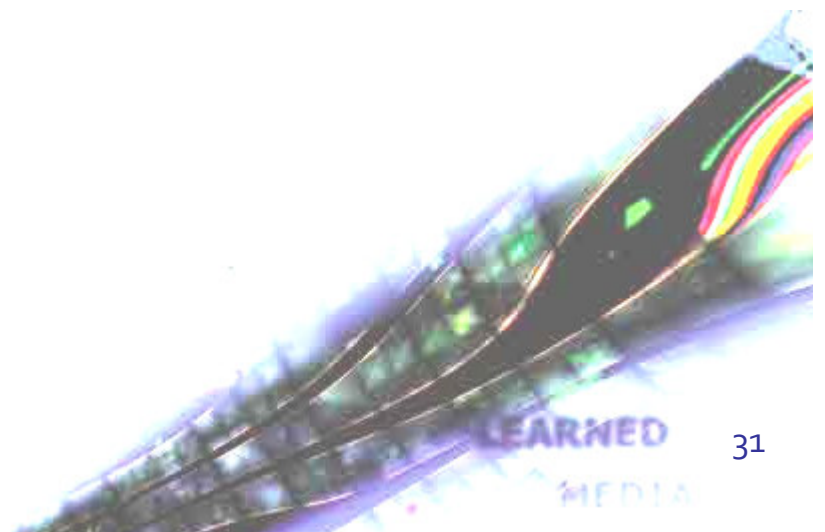
- Consider a dataset with a target variable with five possible values:
  1. non sufficient
  2. sufficient
  3. good
  4. very good
  5. excellent
  - Is this a classification or a numeric prediction problem?
  - What if such a variable is an attribute, is it nominal or numeric?
- Can KNN be used for classification tasks?
- Similarities between KNN and Naïve Bayes.
- • Similarities and differences between decision trees and regression trees.



Regression trees	Decision trees
<b>Data:</b> attribute-value description	
<b>Target variable:</b> Continuous	<b>Target variable:</b> Categorical (nominal)
<b>Evaluation:</b> cross validation, separate test set, ...	
<b>Error:</b> MSE, MAE, RMSE, ...	<b>Error:</b> 1-accuracy
<b>Algorithm:</b> Top down induction, shortsighted method	
<b>Heuristic:</b> Standard deviation	<b>Heuristic :</b> Information gain
<b>Stopping criterion:</b> Standard deviation < threshold	<b>Stopping criterion:</b> Pure leafs (entropy=0)

---

# Association Rules



# Association rules

---

- Rules  $X \rightarrow Y$ ,  $X, Y$  conjunction of items
- Task: Find **all** association rules that satisfy minimum support and minimum confidence constraints
- **Support:**  
$$\text{Sup}(X \rightarrow Y) = \#XY/\#D \cong p(XY)$$
- **Confidence:**  
$$\text{Conf}(X \rightarrow Y) = \#XY/\#X \cong p(XY)/p(X) = p(Y|X)$$

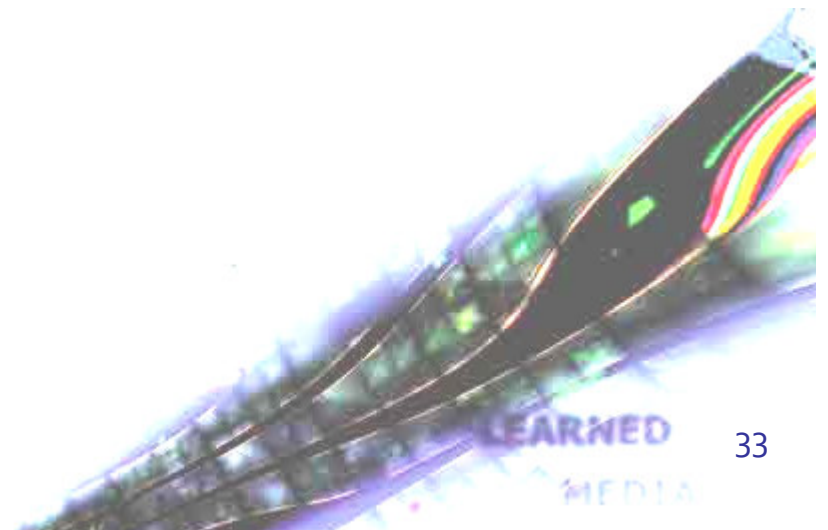


# Association rules - algorithm

---

1. generate frequent itemsets with a minimum support constraint
2. generate rules from frequent itemsets with a minimum confidence constraint

\* Data are in a transaction database

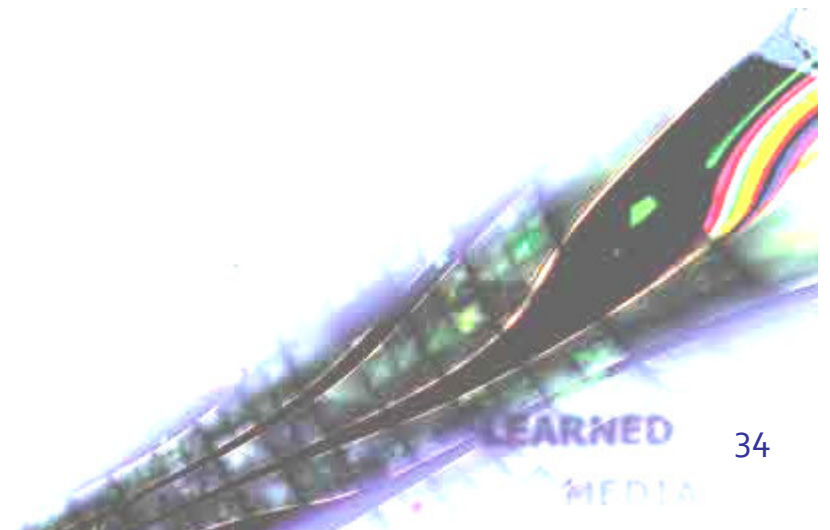


# Association rules – transaction database

---

Items: **A**=apple, **B**=banana,  
**C**=coca-cola, **D**=doughnut

- Client 1 bought: A, B, C, D
- Client 2 bought: B, C
- Client 3 bought: B, D
- Client 4 bought: A, C
- Client 5 bought: A, B, D
- Client 6 bought: A, B, C

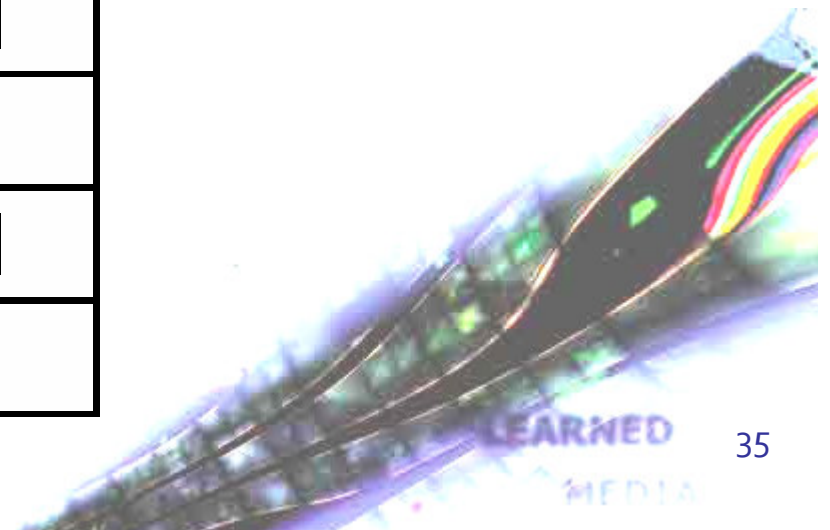


# Frequent itemsets

---

- Generate frequent itemsets with support at least 2/6

<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>
<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>
	<b>1</b>	<b>1</b>	
	<b>1</b>		<b>1</b>
<b>1</b>		<b>1</b>	
<b>1</b>	<b>1</b>		<b>1</b>
<b>1</b>	<b>1</b>	<b>1</b>	



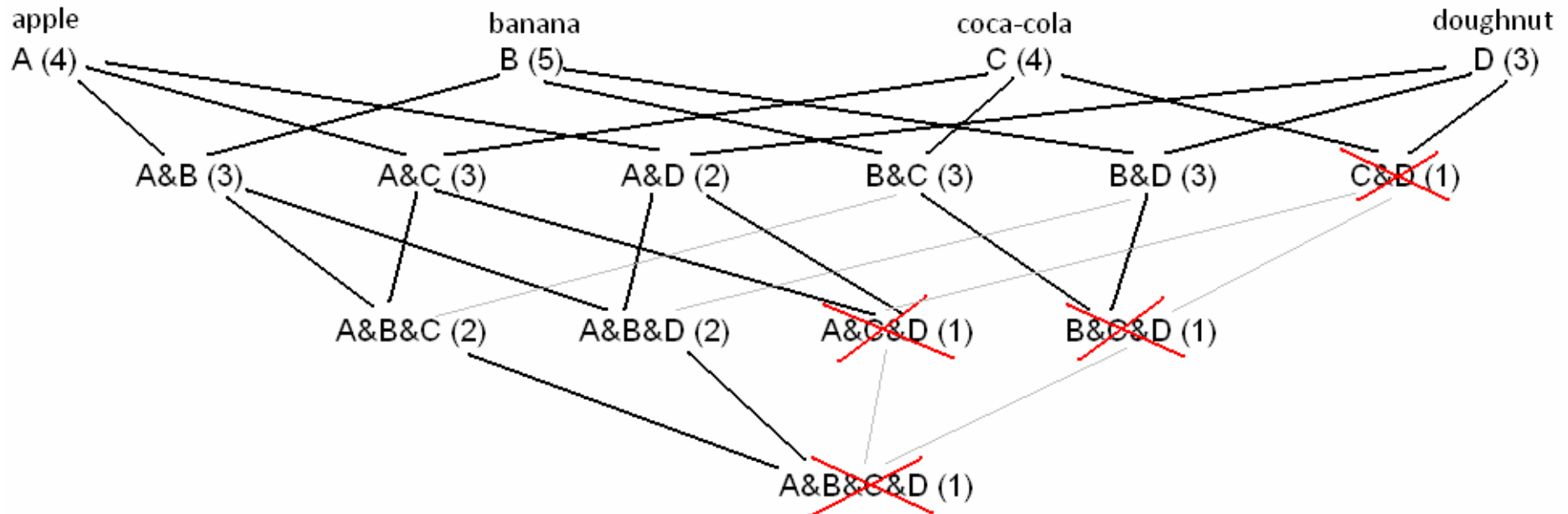
# Frequent itemsets algorithm

---

Items in an itemset should be sorted alphabetically.

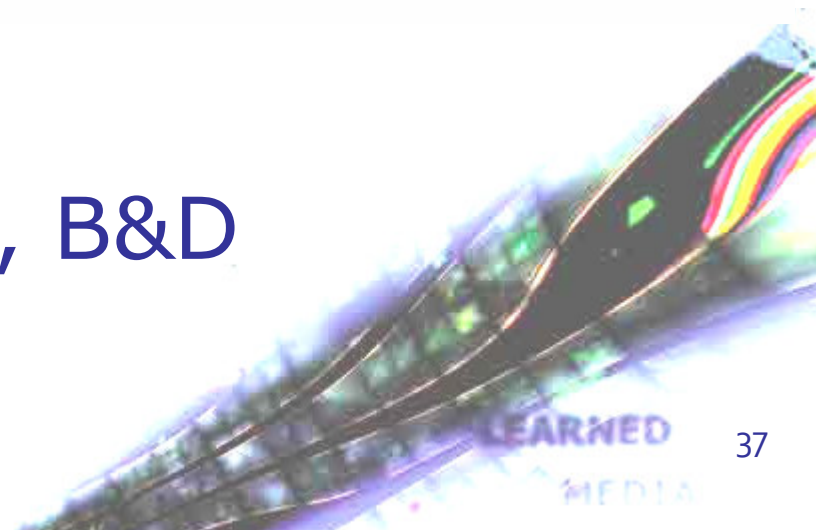
- Generate all 1-itemsets with the given minimum support.
- Use 1-itemsets to generate 2-itemsets with the given minimum support.
- From 2-itemsets generate 3-itemsets with the given minimum support as unions of 2-itemsets with the same item at the beginning.
- ...
- From  $n$ -itemsets generate  $(n+1)$ -itemsets as unions of  $n$ -itemsets with the same  $(n-1)$  items at the beginning.

# Frequent itemsets lattice



Frequent itemsets:

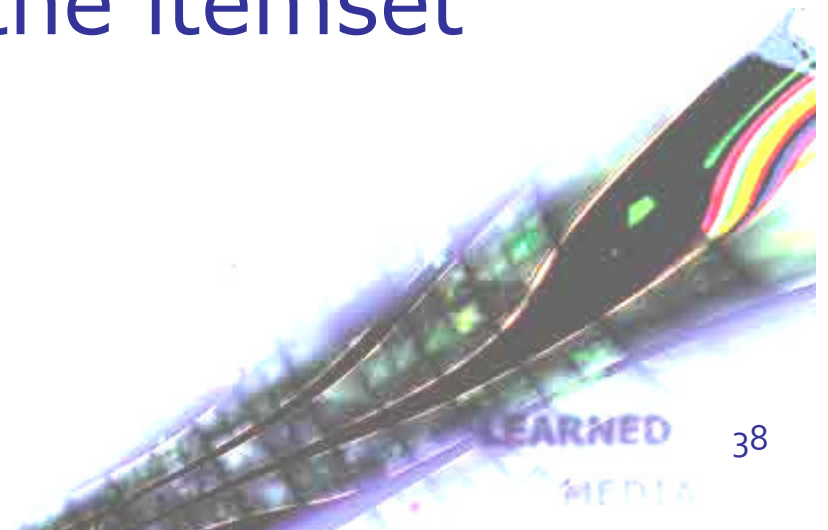
- A&B, A&C, A&D, B&C, B&D
- A&B&C, A&B&D



# Rules from itemsets

---

- A&B is a frequent itemset with support 3/6
- Two possible rules
  - $A \rightarrow B$  confidence =  $\#(A\&B)/\#A = 3/4$
  - $B \rightarrow A$  confidence =  $\#(A\&B)/\#B = 3/5$
- All the counts are in the itemset lattice!



# Quality of association rules

---

$$\text{Support}(X) = \#X / \#D \quad \dots\dots\dots P(X)$$

$$\text{Support}(X \rightarrow Y) = \text{Support}(XY) = \#XY / \#D \quad \dots\dots\dots P(XY)$$

$$\text{Confidence}(X \rightarrow Y) = \#XY / \#X \quad \dots\dots\dots P(Y|X)$$

---

$$\text{Lift}(X \rightarrow Y) = \text{Support}(X \rightarrow Y) / (\text{Support}(X) * \text{Support}(Y))$$

$$\text{Leverage}(X \rightarrow Y) = \text{Support}(X \rightarrow Y) - \text{Support}(X) * \text{Support}(Y)$$

$$\text{Conviction}(X \rightarrow Y) = 1 - \text{Support}(Y) / (1 - \text{Confidence}(X \rightarrow Y))$$



# Quality of association rules

---

$$\text{Support}(X) = \#X / \#D \quad \dots\dots\dots P(X)$$

$$\text{Support}(X \rightarrow Y) = \text{Support}(XY) = \#XY / \#D \quad \dots\dots\dots P(XY)$$

$$\text{Confidence}(X \rightarrow Y) = \#XY / \#X \quad \dots\dots\dots P(Y|X)$$

---

$$\text{Lift}(X \rightarrow Y) = \text{Support}(X \rightarrow Y) / (\text{Support}(X) * \text{Support}(Y))$$

How many more times the items in X and Y occur together than it would be expected if the itemsets were statistically independent.

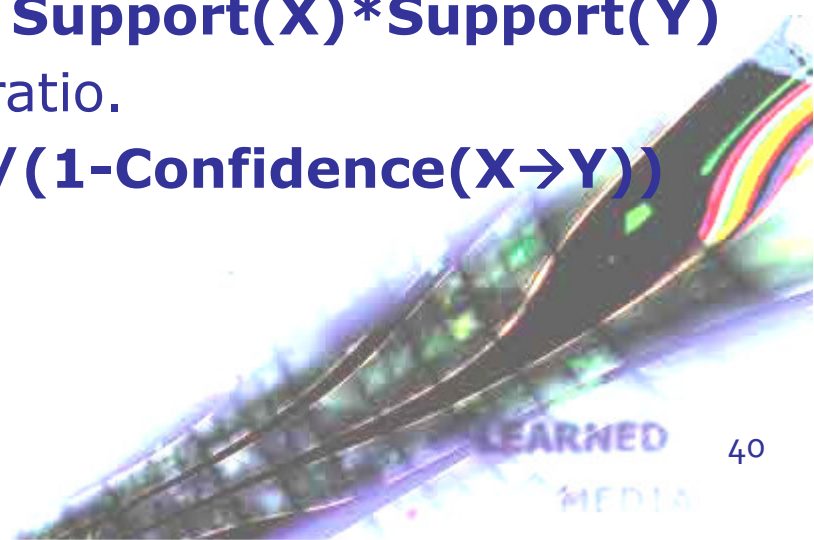
$$\text{Leverage}(X \rightarrow Y) = \text{Support}(X \rightarrow Y) - \text{Support}(X) * \text{Support}(Y)$$

Similar to lift, difference instead of ratio.

$$\text{Conviction}(X \rightarrow Y) = 1 - \text{Support}(Y) / (1 - \text{Confidence}(X \rightarrow Y))$$

Degree of implication of a rule.

Sensitive to rule direction.





# Discussion

---

- Transformation of an attribute-value dataset to a transaction dataset.
- What would be the association rules for a dataset with two items A and B, each of them with support 80% and appearing in the same transactions as rarely as possible?
  - minSupport = 50%, min conf = 70%
  - minSupport = 20%, min conf = 70%
- What if we had 4 items: A,  $\neg A$ , B,  $\neg B$
- Compare decision trees and association rules regarding handling an attribute like "PersonID". What about attributes that have many values (eg. Month of year)

A	B
Green	White
Green	White
Green	Blue
Green	Blue
Green	Blue
Green	Blue
Green	Blue
Green	Blue
White	Blue
White	Blue