



## Data Mining and Knowledge Discovery


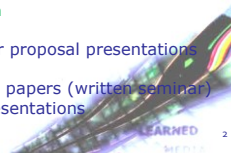
Petra Kralj Novak  
[Petra.Kralj.Novak@ijs.si](mailto:Petra.Kralj.Novak@ijs.si)

Practice, 2010/12/2


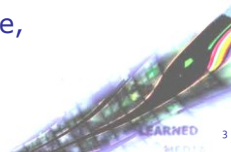
## Practice plan

- 2010/11/25: Predictive data mining
  - Decision trees
  - Naive Bayes classifier
  - Evaluating classifiers (separate test set, cross validation, confusion matrix, classification accuracy)
  - Predictive data mining in Weka
- 2010/12/2: Numeric prediction and descriptive data mining
  - Numeric prediction models
  - Association rules
  - Regression models and evaluation in Weka
  - Descriptive data mining in Weka
  - Discussion about seminars and exam
- 2010/12/16: Written exam, Seminar proposal presentations
- 2011/2/1: Deadline for data mining papers (written seminar)
- 2011/2/3: Data mining seminar presentations


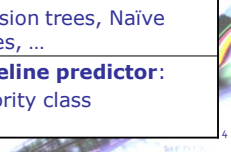



## Numeric prediction

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

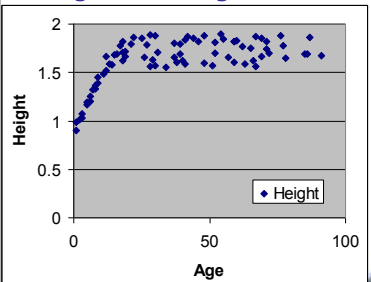



Numeric prediction	Classification
<b>Data:</b> attribute-value description	
<b>Target variable:</b> Continuous	<b>Target variable:</b> Categorical (nominal)
<b>Evaluation:</b> cross validation, separate test set, ...	
<b>Error:</b> MSE, MAE, RMSE, ...	<b>Error:</b> 1-accuracy
<b>Algorithms:</b> Linear 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


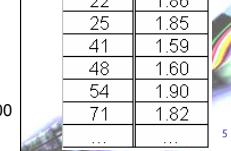



## Example

- 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
...	...

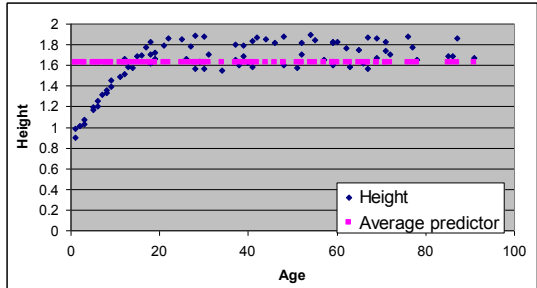
## Test set

Age	Height
2	0.85
10	1.4
35	1.7
70	1.6




## Baseline numeric predictor

- Average of the target variable



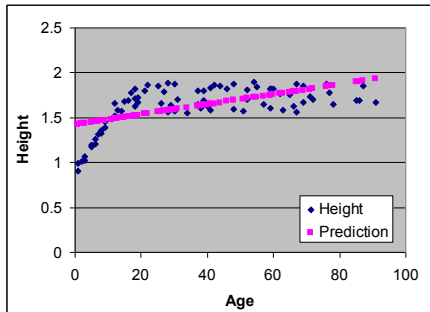
## Baseline predictor: prediction

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

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

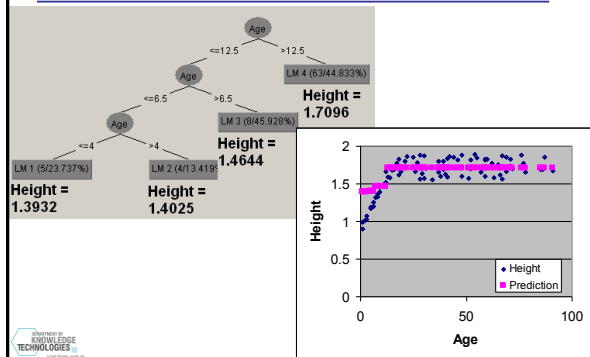


## Linear Regression: prediction

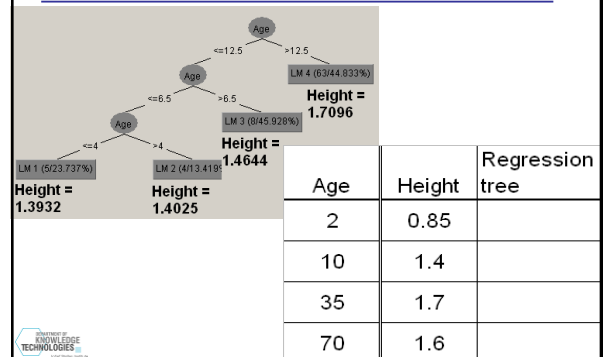
$$\text{Height} = 0.0056 * \text{Age} + 1.4181$$

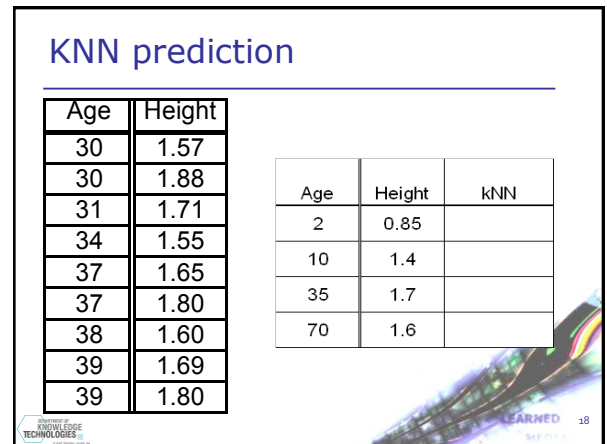
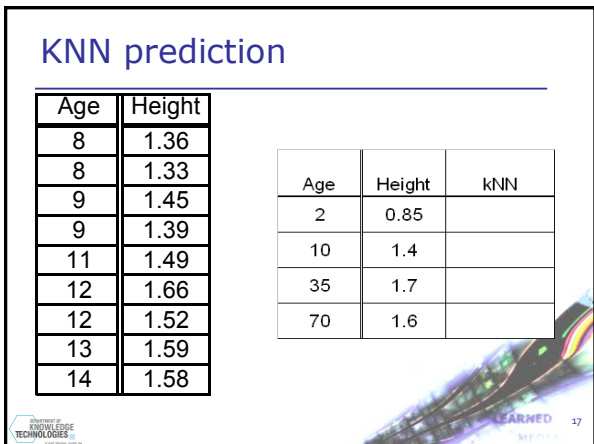
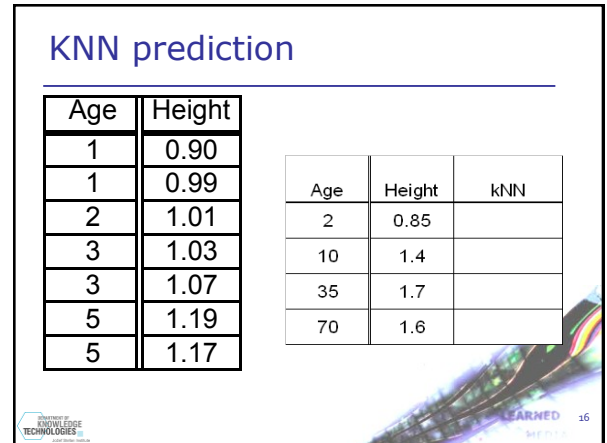
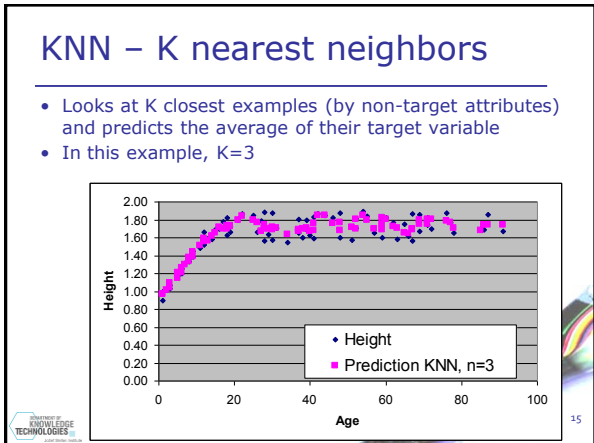
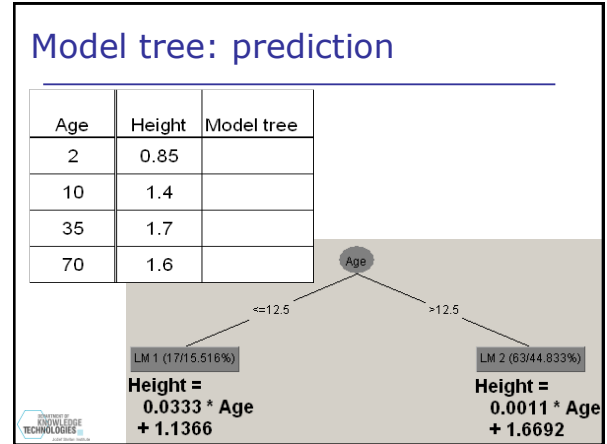
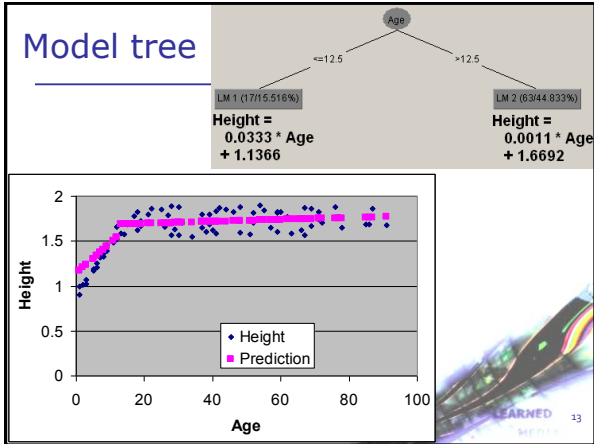
Age	Height	Linear regression
2	0.85	
10	1.4	
35	1.7	
70	1.6	

## Regression tree



## Regression tree: prediction





## 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 on 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	$\frac{(p_1 - a_1)^2 + \dots + (p_n - a_n)^2}{(a_1 - \bar{a})^2 + \dots + (a_n - \bar{a})^2}$ , where $\bar{a} = \frac{1}{n} \sum a_i$
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	$\frac{S_{p,a}}{\sqrt{S_p S_a}}$ , where $S_{p,a} = \sum_{i=1}^n (p_i - \bar{p})(a_i - \bar{a})$ , $S_p = \sum_{i=1}^n (p_i - \bar{p})^2$ , and $S_a = \sum_{i=1}^n (a_i - \bar{a})^2$

## Numeric prediction discussion

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

- non sufficient
- sufficient
- good
- very good
- 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 Naive Bayes.
- Similarities and differences between decision trees and regression trees.

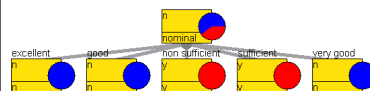
## Classification or a numeric prediction problem?

- Target variable with five possible values:
  - non sufficient
  - sufficient
  - good
  - very good
  - 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.

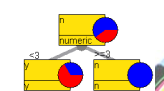
## Nominal or numeric attribute?

- A variable with five possible values:
  - non sufficient
  - sufficient
  - good
  - very good
  - 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:
  - non sufficient
  - sufficient
  - good
  - very good
  - 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.



25

## 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



26

## Numeric prediction discussion

- Consider a dataset with a target variable with five possible values:
  - non sufficient
  - sufficient
  - good
  - very good
  - 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.



27

## 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)



28

## Numeric prediction discussion

- Consider a dataset with a target variable with five possible values:
  - non sufficient
  - sufficient
  - good
  - very good
  - 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.



29

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 leaves (entropy=0)



30

## 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:**  

$$\text{Sup}(X \rightarrow Y) = \#XY/\#D \equiv p(XY)$$
- **Confidence:**  

$$\text{Conf}(X \rightarrow Y) = \#XY/\#X \equiv p(XY)/p(X) = p(Y|X)$$

## 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

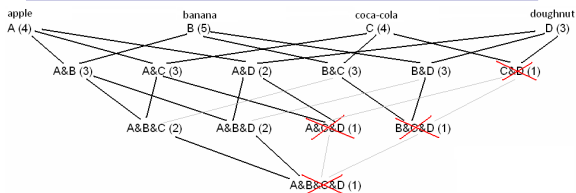
- Generate frequent itemsets with support at least 2/6

A	B	C	D
1	1	1	1
	1	1	
	1		1
1		1	
1	1		1
1	1	1	

## 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

$$\begin{aligned} \text{Support}(X) &= \#X / \#D && \dots\dots\dots P(X) \\ \text{Support}(X \rightarrow Y) &= \text{Support}(XY) = \#XY / \#D && \dots\dots\dots P(XY) \\ \text{Confidence}(X \rightarrow Y) &= \#XY / \#X && \dots\dots\dots P(Y|X) \end{aligned}$$

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

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

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

## Quality of association rules

$$\begin{aligned} \text{Support}(X) &= \#X / \#D && \dots\dots\dots P(X) \\ \text{Support}(X \rightarrow Y) &= \text{Support}(XY) = \#XY / \#D && \dots\dots\dots P(XY) \\ \text{Confidence}(X \rightarrow Y) &= \#XY / \#X && \dots\dots\dots P(Y|X) \end{aligned}$$

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

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
1	Green	Blue
2	Green	Blue
3	Green	Blue
4	Green	Blue
5	Green	Blue
6	Green	Blue
7	Green	Blue
8	Green	Blue
9	Green	Blue
10	Green	Blue
11	Green	Blue
12	Green	Blue
13	Green	Blue
14	Green	Blue
15	Green	Blue
16	Green	Blue
17	Green	Blue
18	Green	Blue
19	Green	Blue
20	Green	Blue