

Data Mining and Knowledge Discovery

Practice notes – 24.11.2009



Numeric prediction and descriptive DM

Data Mining and Knowledge Discovery

Knowledge Discovery and Knowledge Management in e-Science



Petra Kralj Novak
Petra.Kralj.Novak@ijs.si

Practice, 2009/11/24



Practice plan

- 2009/11/10: Predictive data mining
 - Decision trees
 - Naive 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


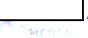



Numeric prediction

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





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


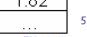



Example

- 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

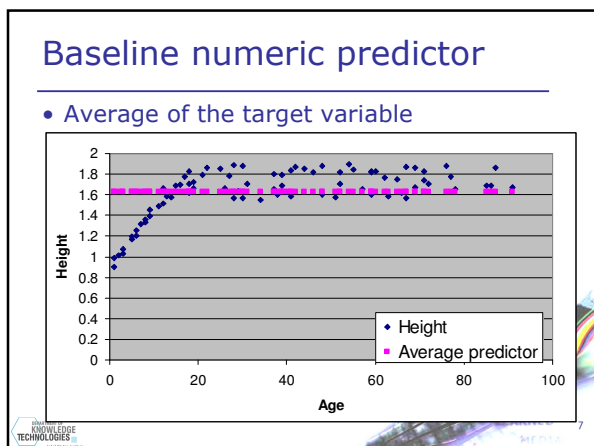
Age	Height
2	0.85
10	1.4
35	1.7
70	1.6

Data Mining and Knowledge Discovery

Practice notes – 24.11.2009

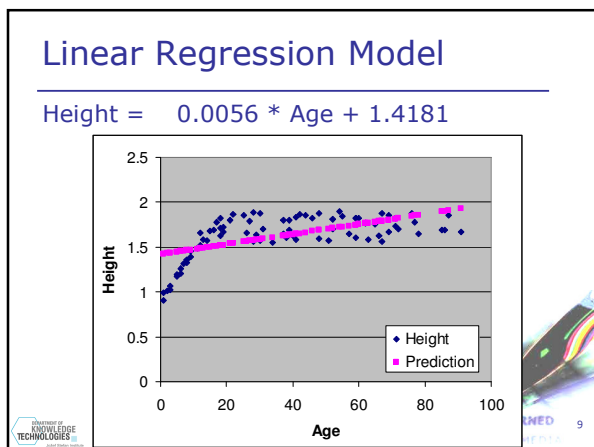
Numeric prediction and descriptive DM



Baseline predictor: prediction

Average of the target variable is 1.63

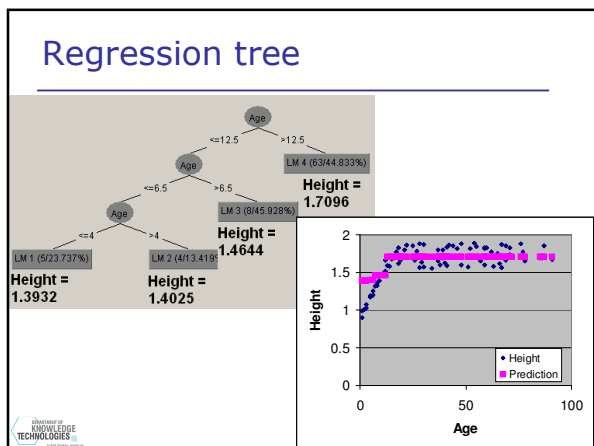
Age	Height	Baseline
2	0.85	
10	1.4	
35	1.7	
70	1.6	



Linear Regression: 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	



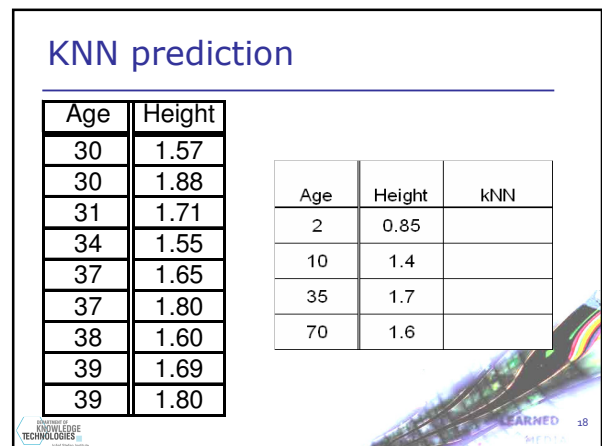
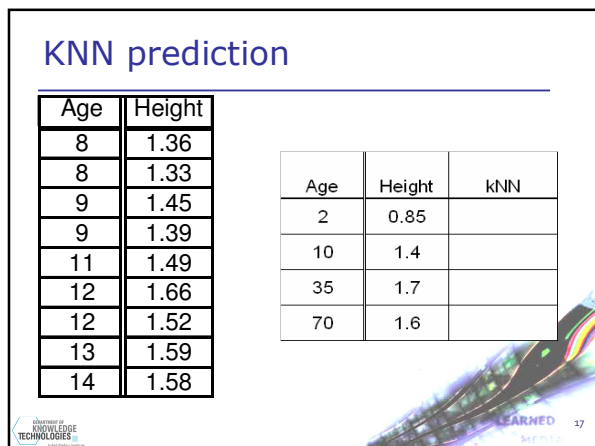
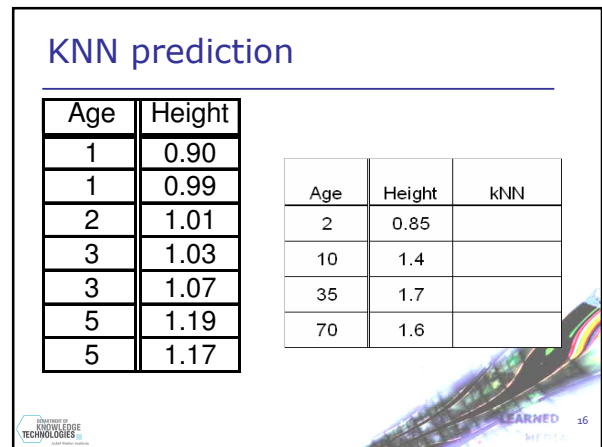
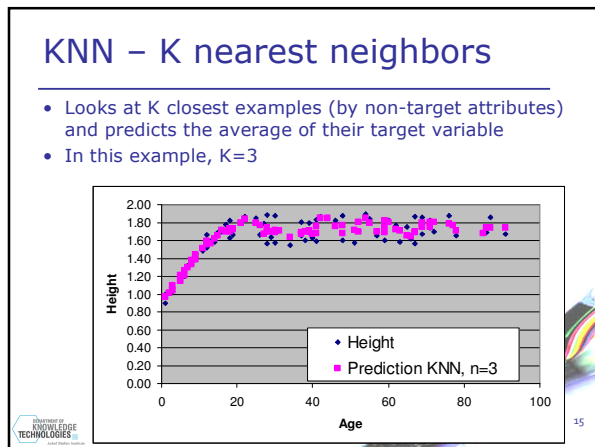
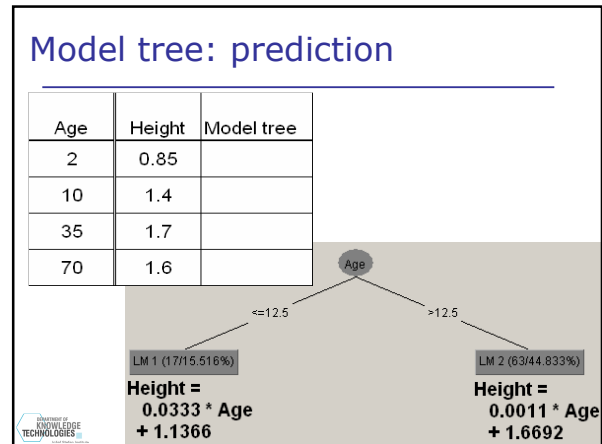
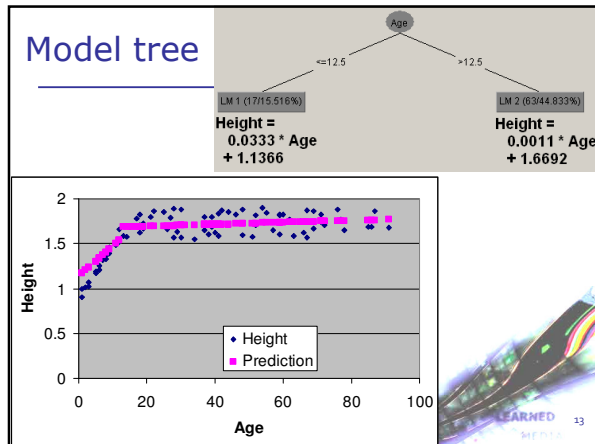
Regression tree: prediction

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

Data Mining and Knowledge Discovery

Practice notes – 24.11.2009

Numeric prediction and descriptive DM



Data Mining and Knowledge Discovery

Practice notes – 24.11.2009

Numeric prediction and descriptive DM

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

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 Naïve 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.

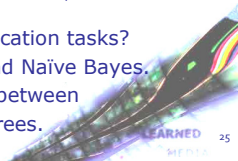
Data Mining and Knowledge Discovery

Practice notes – 24.11.2009

Numeric prediction and descriptive DM

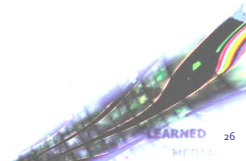
Numeric prediction discussion

- 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.



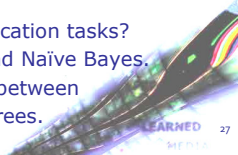
Can KNN be used for classification tasks?

- 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



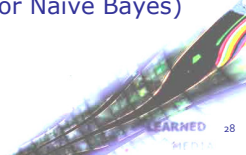
Numeric prediction discussion

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



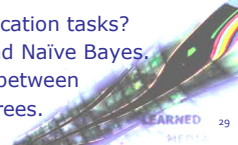
Similarities between KNN and Naïve Bayes.

- 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)

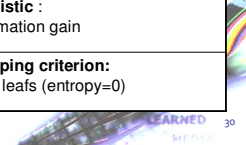


Numeric prediction discussion

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



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




Data Mining and Knowledge Discovery

Practice notes – 24.11.2009

Numeric prediction and descriptive DM

Association Rules



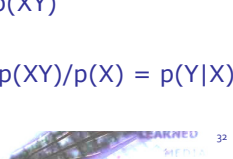
31

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:**

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

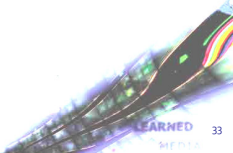


32

Association rules - algorithm

- generate frequent itemsets with a minimum support constraint
- generate rules from frequent itemsets with a minimum confidence constraint

* Data are in a transaction database

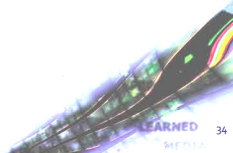


33

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

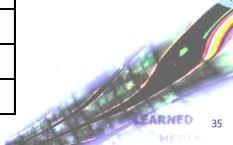


34

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	

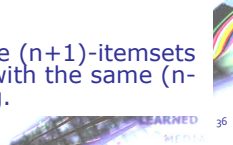


35

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.



36

Data Mining and Knowledge Discovery

Practice notes – 24.11.2009

Numeric prediction and descriptive DM

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

Support(X) = $\#X / \#D$ P(X)
 Support(X→Y) = $\#XY / \#D$ P(XY)
 Confidence(X→Y) = $\#XY / \#X$ P(Y|X)

Lift(X→Y) = Support(X→Y) / (Support(X)*Support(Y))

Leverage(X→Y) = Support(X→Y) - Support(X)*Support(Y)

Conviction(X → Y) = 1-Support(Y)/(1-Confidence(X→Y))

Quality of association rules

Support(X) = $\#X / \#D$ P(X)
 Support(X→Y) = $\#XY / \#D$ P(XY)
 Confidence(X→Y) = $\#XY / \#X$ P(Y|X)

Lift(X→Y) = Support(X→Y) / (Support(X)*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.

Leverage(X→Y) = Support(X→Y) - Support(X)*Support(Y)
 Similar to lift, difference instead of ratio.

Conviction(X → Y) = 1-Support(Y)/(1-Confidence(X→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, ¬A, B, ¬B
- Compare decision trees and association rules regarding handling an attribute like "PersonID". What about attributes that have many values (eg. Month of year)