



COLLEGE OF ENGINEERING AND MANAGEMENT PUNNAPRA

Under CAPE | Government of Kerala

Issue No. / Date: 01/26.07.2019		TUTORIAL -3	CEMP/ISO/UG/CSE/653 I																																								
Revn. No. / Date:00																																											
Programme: B.TECH		Branch: COMPUTER SCIENCE & ENGINEERING																																									
Semester: S8		Academic Year: 2022-2023																																									
Subject Code: CST 466	L-T-P- Credits: 2-1-0-3	Subject Name: DATA MINING																																									
Date:	Hour:	Topic: Data Preprocessing																																									
Objective																																											
To understand the basic ideas of Steps in preprocessing																																											
Questions																																											
<p>1. Suppose that the data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.</p> <p>(a) What is the mean of the data? What is the median?</p> <p>(b) Use smoothing by bin means to smooth the data, using a bin depth of 3. Illustrate your steps. Comment on the effect of this technique for the given data.</p> <p>(c) How might you determine outliers in the data?</p> <p>(d) What other methods are there for data smoothing?</p> <p>(e) Use min-max normalization to transform the value 35 for age onto the range [0:0;1:0].</p> <p>(f) Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years.</p> <p>(g) Use normalization by decimal scaling to transform the value 35 for age.</p> <p>(h) Comment on which method of normalization you would prefer to use for the given data, giving reasons as to why.</p> <p>2. Suppose a hospital tested the age and body fat data for 18 randomly selected adults with the following result:</p> <table><tr><td>age</td><td>23</td><td>23</td><td>27</td><td>27</td><td>39</td><td>41</td><td>47</td><td>49</td><td>50</td></tr><tr><td>%fat</td><td>9.5</td><td>26.5</td><td>7.8</td><td>17.8</td><td>31.4</td><td>25.9</td><td>27.4</td><td>27.2</td><td>31.2</td></tr><tr><td>age</td><td>52</td><td>54</td><td>54</td><td>56</td><td>57</td><td>58</td><td>58</td><td>60</td><td>61</td></tr><tr><td>%fat</td><td>34.6</td><td>42.5</td><td>28.8</td><td>33.4</td><td>30.2</td><td>34.1</td><td>32.9</td><td>41.2</td><td>35.7</td></tr></table> <p>a. Calculate the mean, median, and standard deviation of age and %fat.</p> <p>b. Normalize the two variables based on z-score normalization.</p> <p>c. Calculate the correlation coefficient (Pearson's product moment coefficient). Are these two variables positively or negatively correlated?</p>				age	23	23	27	27	39	41	47	49	50	%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2	age	52	54	54	56	57	58	58	60	61	%fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7
age	23	23		27	27	39	41	47	49	50																																	
%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2																																		
age	52	54	54	56	57	58	58	60	61																																		
%fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7																																		
1. 1. Dunham M H, "Data Mining: Introductory and Advanced Topics", Pearson Education, New Delhi, 2003.																																											