

**Course:** Data Mining  
**Student:** Ekaterina Eremina  
**Topic:** HW2

## Part\_1: written assignment

### Exercise\_1

(a)  $p(C1) = 5$ ,  $n(C2) = 4$ .

$I(C1, C2) = 0,9911$

	$p_i$	$n_i$	$I(p_i, n_i)$
Height = Tall	3	2	0,9710
Height = Short	2	1	0,9183
Height = Medium	0	1	0
Hair = Blond	2	2	1
Hair = Dark	3	1	0,8113
Hair = Red	0	1	0
Eye = Brown	3	0	0
Eye = Blue	2	4	0,9183

$$E(\text{Height}) = 5/9 * I(3,2) +$$

$$3/9 * I(2,1) + 1/9 * I(0,1) = 0,8455$$

$$E(\text{Hair}) = 4/9 * I(2,2) + 4/9 * I(3,1) +$$

$$1/9 * I(0,1) = 0,8050$$

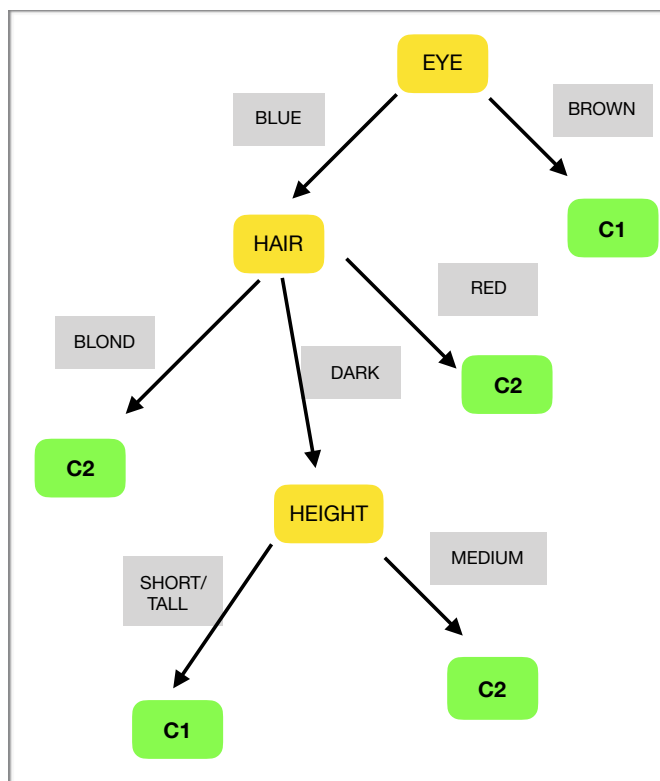
$$E(\text{Eye}) = 3/9 * I(3,0) + 6/9 * I(2,4) = 0,6122$$

$$\text{Gain}(\text{Height}) = 0,9911 - E(\text{Height}) = 0,1456$$

$$\text{Gain}(\text{Hair}) = 0,9911 - E(\text{Hair}) = 0,1861$$

$$\text{Gain}(\text{Eye}) = 0,9911 - E(\text{Eye}) = 0,3789$$

(b)



## Exercise\_2

$Z = \{\text{Brown, Blond, Short}\} \rightarrow C1$

Let's count probabilities for naive Bayesian classifier.

$$P(C1) = 5/9$$

$$P(C2) = 4/9$$

$$P(\text{Brown}|C1) = 3/5$$

$$P(\text{Brown}|C2) = 0$$

$$P(\text{Blond}|C1) = 2/5$$

$$P(\text{Blond}|C2) = 2/4$$

$$P(\text{Short}|C1) = 2/5$$

$$P(\text{Short}|C2) = 1/4$$

$$P(Z|C1) = 2/5 * 2/5 * 3/5 = 12/125$$

$$P(Z|C2) = 1/4 * 2/4 * 0 = 0$$

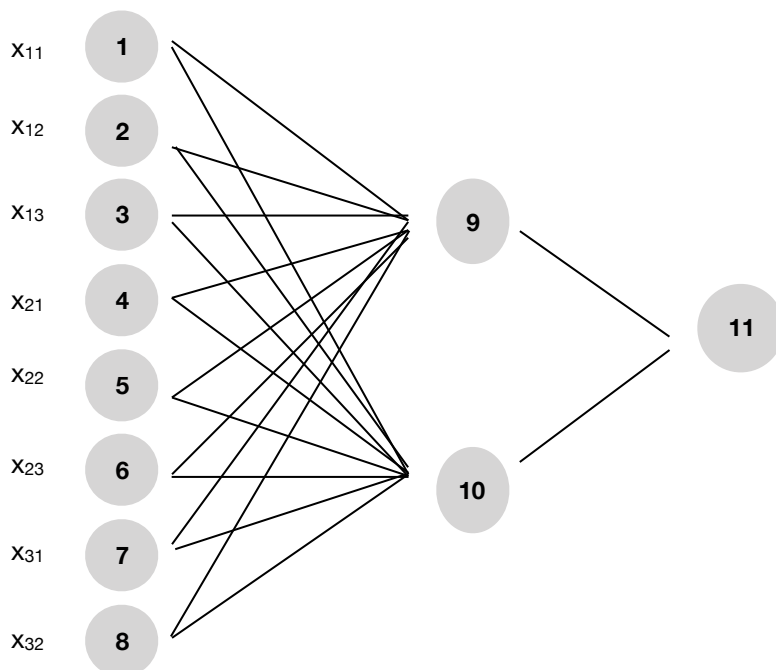
$$P(Z|C1) * P(C1) = 12/125 * 5/9 = 0,053$$

$$P(Z|C2) * P(C2) = 0$$

Result will be the same:  $= \{\text{Brown, Blond, Short}\} \rightarrow C1$ .

## Exercise\_3

(a)



(b) Training = {Medium, Blond, Blue -> C1}

X11	X12	X13	X21	X22	X23	X31	X32
tall	short	medium	blond	dark	red	brown	blue
0	0	1	1	0	0	0	1
W1,9	W1,10	W2,9	W2,10	W3,9	W3,10	W4,9	W4,10
0,2	0,2	-0,1	0,4	0,3	-0,1	-0,2	0,1
W5,9	W5,10	W6,9	W6,10	W7,9	W7,10	W8,9	W8,10
0,2	0,2	0,1	-0,1	0,3	0,2	0,1	0,2
W9,11	W10,11	$\theta_9$	$\theta_{10}$	$\theta_{11}$			
-0,3	0,1	-0,2	0,2	0,3			

Node, j	Input, I <sub>j</sub>	Output, O <sub>j</sub>
9	$0+0+0,3-0,2+0+0+0+0,1-0,2 = 0$	$1/(1+e^0) = 0,5$
10	$0+0-0,1+0,1+0+0+0+0,1+0,2 = 0,3$	$1/(1+e^{(-0,3)}) = 0,5744$
11	$(-0,3)*0,5+0,1*0,5744+0,1 = 0,00744$	$1/(1+e^{(-0,00744)}) = 0,502$

Node, j	Error, Err <sub>j</sub>
11	$0,502*(1-0,502)*(1-0,502) = 0,1245$
10	$0,5744*(1-0,5744)*0,125*0,1 = 0,0031$
9	$0,5*(1-0,5)*0,125*(-0,3) = -0,0094$

W1,10	$0,2 + 0,9*(0,0031)*0 = 0,2$	W6,10	$0,1 + 0,9*(0,0031)*0 = 0,1$
W1,9	$0,2 + 0,9*(-0,0094)*0 = 0,2$	W7,9	$0,3 + 0,9*(-0,0094)*0 = 0,3$
W2,10	$-0,1 + 0,9*(0,0031)*0 = -0,1$	W7,10	$0,3 + 0,9*(0,0031)*0 = 0,3$
W2,9	$-0,1 + 0,9*(-0,0094)*0 = -0,1$	W8,9	$0,1 + 0,9*(-0,0094)*1 = 0,092$
W3,10	$0,3 + 0,9*(0,0031)*1 = 0,303$	W8,10	$0,1 + 0,9*(0,0031)*1 = 0,103$
W3,9	$0,3 + 0,9*(-0,0094)*1 = 0,292$	W9,11	$-0,3+0,9*(0,1245)*0,5 = -0,244$
W4,9	$-0,2 + 0,9*(-0,0094)*1 = -0,209$	W10,11	$0,1+0,9*(0,1245)*0,5744 = 0,164$
W4,10	$-0,2 + 0,9*(0,0031)*1 = -0,197$	$\theta_9$	$-0,2+0,9*(-0,0094) = -0,209$
W5,9	$0,2 + 0,9*(-0,0094)*0 = 0,2$	$\theta_{10}$	$0,2+0,9*(0,0031) = 0,203$
W5,10	$0,2 + 0,9*(0,0031)*0 = 0,2$	$\theta_{11}$	$0,3+0,9*(0,1245) = 0,412$
W6,9	$0,1 + 0,9*(-0,0094)*0 = 0,1$		

## Exercise\_4

(a) We have clusters centroids  $A_1$ ,  $B_1$  and  $C_1$ . Let's calculate Euclidian distances.

	$A_1$	$B_1$	$C_1$
$A_2$	7,348	9,899	4,123
$A_3$	6,403	10,05	7,874
$B_2$	3,742	2,449	10,817
$B_3$	5,745	9,644	11,045
$C_2$	3,742	5,831	11,874
$C_3$	5,477	10	9,644
$C_4$	4,583	8,775	8,367

After the first round execution we will have three clusters  $\{A_1A_3B_3C_2C_3C_4\}$ ,  $\{B_1B_2\}$  and  $\{A_2C_1\}$ . Cluster centers are  $(4,5; 4,5; 6,83)$ ,  $(1,5; 2; 1,5)$  and  $(10,5; 2; 2)$  respectively.

(b)

	$(4,5; 4,5; 6,83)$	$(1,5; 2; 1,5)$	$(10,5; 2; 2)$
$A_1$	3,138	4,301	7,159
$A_2$	7,337	9,028	3,041
$A_3$	3,54	8,86	9,552
$B_1$	7,684	1,225	9,605
$B_2$	5,642	1,225	8,559
$B_3$	3,035	8,631	11,011
$C_1$	9,264	11,811	7,018
$C_2$	3,632	4,95	10,5
$C_3$	5,703	9,354	5,315
$C_4$	1,59	7,649	8,441

Finally, we have three clusters:  $\{A_1A_3B_3C_2C_3C_4\}$ ,  $\{B_1B_2\}$  and  $\{A_2C_1\}$ .

## Part\_2: Lab

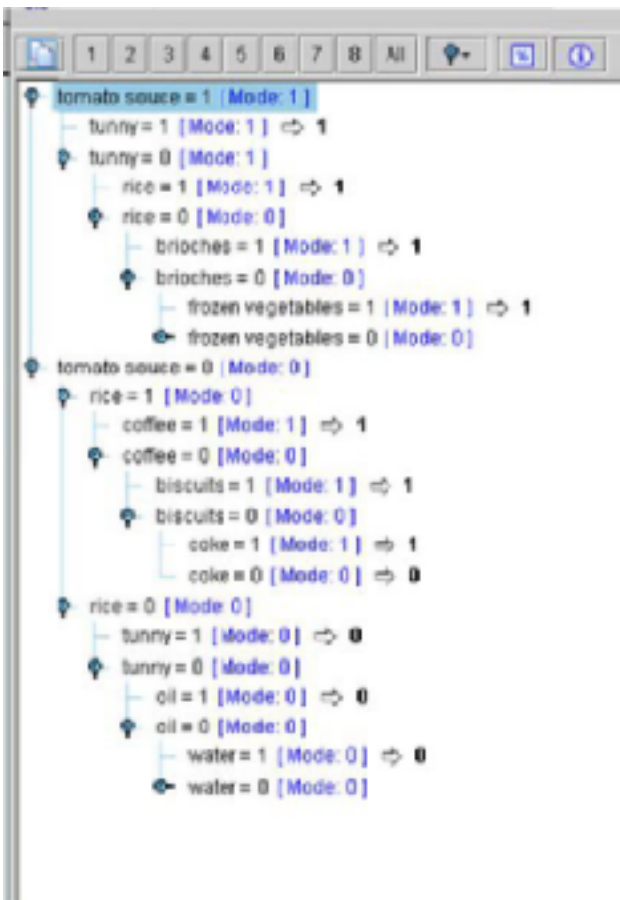
### Exercise\_1

(1) Decision tree can be seen on the next page.

(2)

	milk	water	biscuits	coffee	croissants	yoghurt	frozen vegetables	tunny	beer	tomato sauce	coke	rice	juices	crackers	oil	frozen fish	ice cream	mozzarella	baked meat	ICG-pasta	ICG-pasta
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.731
2	0	0	1	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0.740
3	1	1	1	0	0	1	1	1	1	0	1	0	0	0	0	0	1	0	0	0	0.610
4	1	1	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0.736
5	1	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0.610
6	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0.776
7	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0.760
8	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.729
9	1	0	1	0	0	1	1	1	0	0	1	0	0	0	0	0	1	0	0	0	0.610
10	1	0	0	1	0	1	0	0	0	0	0	1	1	0	1	0	0	0	1	0	0.600
11	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.740
12	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0.731
13	1	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0.729
14	1	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0.729
15	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0.680
16	1	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0.726
17	0	1	0	0	0	0	1	0	0	0	1	1	0	0	0	0	0	0	0	1	0.664
18	1	0	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	1	0	0	0.680
19	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0.731
20	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0.662

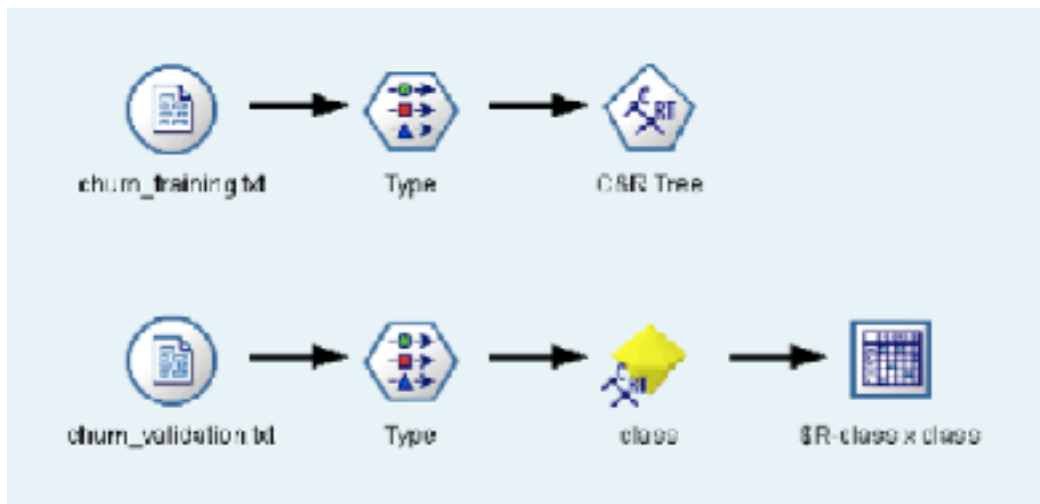
(3)





## Exercise\_2

(1)



Matrix of class by \$C-class #1

File Edit Generate

\$C-class

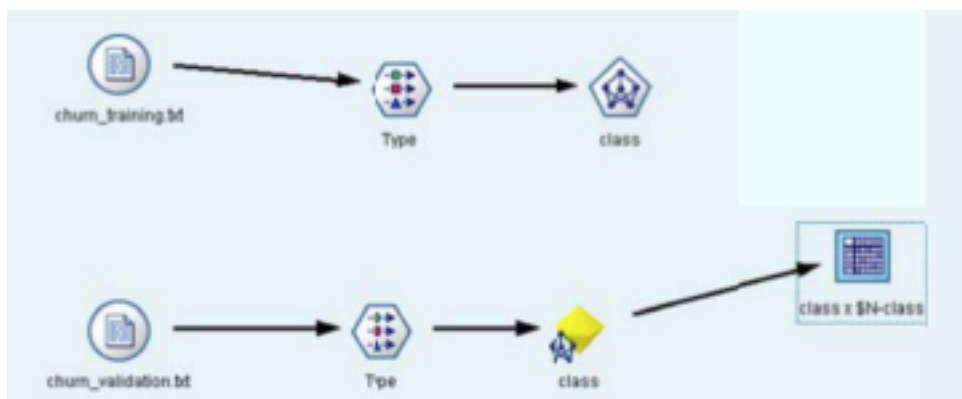
class	0	1	Total
0	981	9	990
1	22	21	43
Total	1003	30	1033

Print the output

Cells contain: cross-tabulation of fields

Matrix Appearance Annotations

(2)



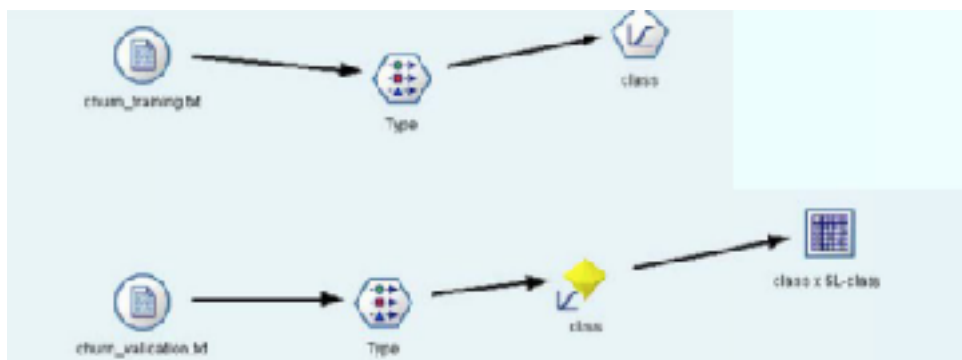
Matrix of class by \$N-class

\$N-class			
class	0	1	Total
0	968	22	990
1	10	33	43
Total	978	55	1033

Cells contain: cross-tabulation of fields

Matrix | Appearance | Annotations

(3)





Matrix of class by \$L-class #2

\$L-class			
class	0	1	Total
0	957	33	990
1	32	11	43
Total	989	44	1033

Cells contain: cross-tabulation of fields

Matrix Appearance Annotations

(4)

	sensitivity	specificity	precision	accuracy
decision_tree	0,4884	0,9909	0,7000	0,9700
neural_network	0,7674	0,9778	0,6000	0,9690
logistic_regression	0,2558	0,9667	0,2500	0,9371

After we calculated all measures for evaluation, let's choose the best method. Accuracy of decision\_tree and neural\_network are almost the same. Next we check sensitivity and specificity - neural\_network predict positive values much better, than decision\_tree, while difference in negative values prediction is not that dramatic. It allows to conclude, that neural\_network is better.