

Олійник Богдан КА-83
Аналіз фінансово-економічних даних
Завдання 3
Дерева рішень

Розглянемо метод CRT

Training Tree

Result

bad
good

Node 0		
Category	%	n
bad	4,9	740
good	95,1	14210
Total	100,0	14950

Income_customer
Improvement=0,013

<= 700

> 700

Node 1		
Category	%	n
bad	38,3	314
good	61,7	506
Total	5,5	820

Interest rate_%
Improvement=0,000

Node 2		
Category	%	n
bad	3,0	426
good	97,0	13704
Total	94,5	14130

goal_credit
Improvement=0,000

<= 18

> 18

Node 3		
Category	%	n
bad	44,6	172
good	55,4	214
Total	2,6	386

Age
Improvement=0,000

Node 4		
Category	%	n
bad	32,7	142
good	67,3	292
Total	2,9	434

Costs_customer
Improvement=0,000

Node 5		
Category	%	n
bad	1,1	50
good	98,9	4505
Total	30,5	4555

flat; house

consumer_credit; overhaul; auto

Node 6		
Category	%	n
bad	3,9	376
good	96,1	9199
Total	64,0	9575

Credit_sum
Improvement=0,000

<= 38

> 38

<= 264

> 264

<= 5375

> 5375

Node 7		
Category	%	n
bad	35,5	59
good	64,5	107
Total	1,1	166

Node 8		
Category	%	n
bad	51,4	113
good	48,6	107
Total	1,5	220

Node 9		
Category	%	n
bad	26,2	60
good	73,8	169
Total	1,5	229

Node 10		
Category	%	n
bad	40,0	82
good	60,0	123
Total	1,4	205

Node 11		
Category	%	n
bad	2,7	118
good	97,3	4334
Total	29,8	4452

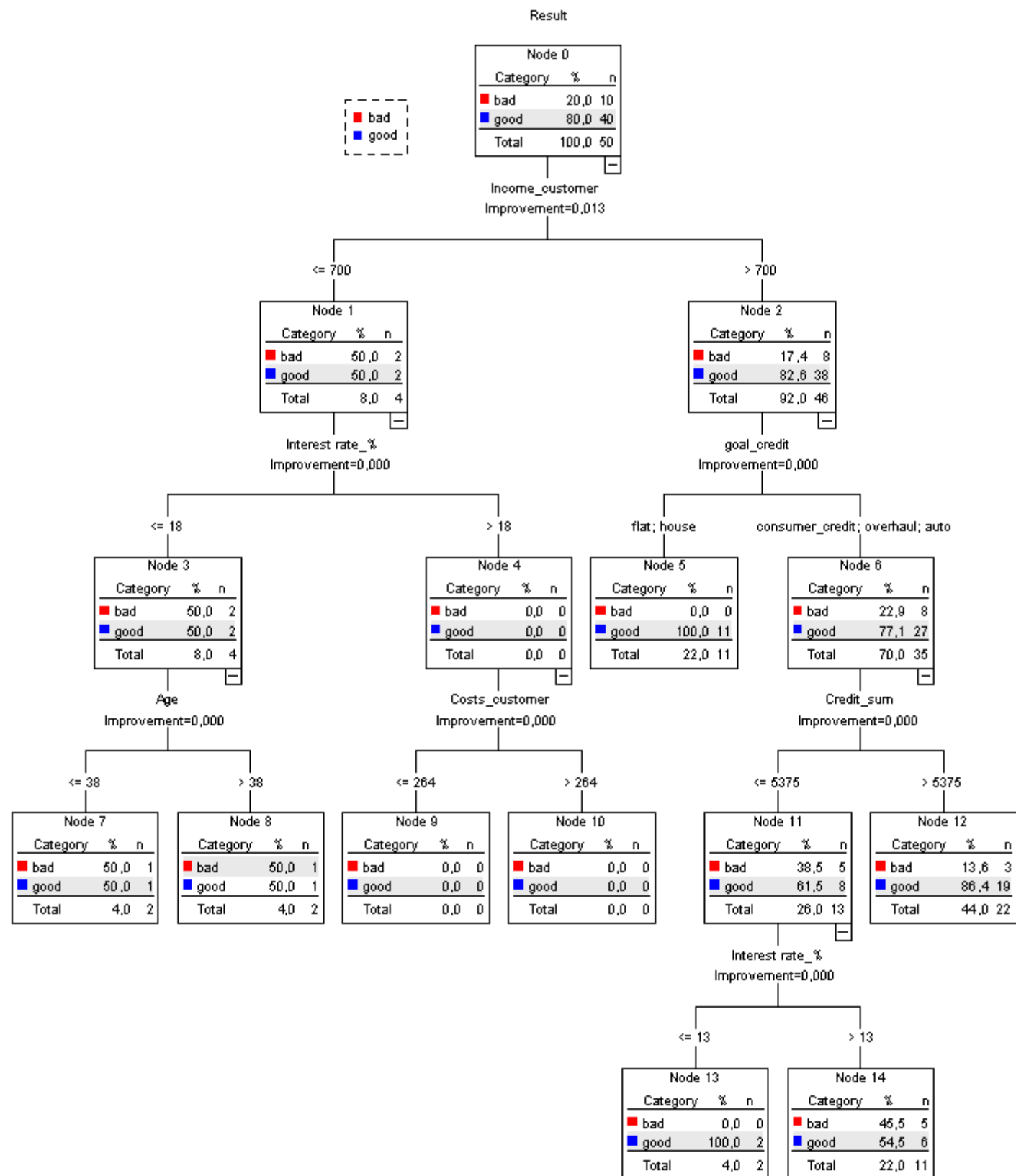
Interest rate_%
Improvement=0,000

<= 13

> 13

Node 13		
Category	%	n
bad	6,1	41
good	93,9	629
Total	4,5	670

Node 14		
Category	%	n
bad	2,0	77
good	98,0	3705
Total	25,3	3782



Test Tree

Sample	Estimate	Std. Error
Training	,049	,002
Test	,200	,057

Classification

		Predicted		
		bad	good	Percent Correct
Training	bad	113	627	15,3%
	good	107	14103	99,2%
	Overall Percentage	1,5%	98,5%	95,1%
Test	bad	1	9	10,0%
	good	1	39	97,5%
	Overall Percentage	4,0%	96,0%	80,0%

	Duration_of_stay_in_a_city	Marital_status	Children	Job_position	Tenure_with_current_employer	Term_of_existence_of_enterprise	Company_type	Number_of_employees_in_company	Income_customer	Costs_customer	goal_credit	Result	Label	PredictedProbability_1	PredictedProbability_2	var	var	
14945	12 120_9999	MARRIED	c2	AS	60_120	60_120	OT	100_9999	25000	1850	consumer_credit	good	1,00	0,05	0,95			
14946	12 120_9999	SINGLE	c0	AS	18_24	12_24	FS	100_9999	1012	265	consumer_credit	bad	1,00	0,06	0,94			
14947	12 120_9999	MARRIED	c1	AS	18_24	12_24	WB	6_15	2697	370	consumer_credit	bad	1,00	0,05	0,95			
14948	12 12_24	MARRIED	c0	SP	24_60	24_60	CI	31_50	2656	765	consumer_credit	good	1,00	0,05	0,95			
14949	15 120_9999	MARRIED	c1	MM	24_60	24_60	FS	31_50	8112	840	flat	good	1,00	0,01	0,99			
14950	16 120_9999	DIVORCED	c1	PE	24_60	60_120	FS	51_100	940	550	consumer_credit	good	1,00	0,02	0,98			
14951	24 120_9999	MARRIED	c1	SP	18_24	12_24	OT	51_100	2344	423	consumer_credit	good	.	0,02	0,98			
14952	16 60_120	SINGLE	c1	PE	24_60	24_60	FS	6_15	730	200	consumer_credit	good	.	0,05	0,95			
14953	24 120_9999	MARRIED	c0	TM	18_24	12_24	FS	6_15	1762	382	consumer_credit	good	.	0,05	0,95			
14954	13 60_120	MARRIED	c0	AS	24_60	60_120	FS	51_100	4049	170	flat	good	.	0,01	0,99			
14955	13 120_9999	MARRIED	c1	TM	60_120	60_120	FS	6_15	7165	310	auto	good	.	0,05	0,95			
14956	18 120_9999	MARRIED	c1	AS	18_24	12_24	OT	51_100	436	210	consumer_credit	bad	1,00	0,51	0,49			
14957	17 120_9999	DIVORCED	c0	SP	120_9999	120_9999	OT	6_15	5353	750	consumer_credit	good	.	0,05	0,95			
14958	16 120_9999	MARRIED	c0	PE	18_24	12_24	FS	51_100	621	290	consumer_credit	good	1,00	0,51	0,49			
14959	14 120_9999	SINGLE	c0	TM	24_60	24_60	WB	0_5	15333	4925	consumer_credit	good	.	0,05	0,95			
14960	15 120_9999	MARRIED	c1	TM	24_60	12_24	AP	0_5	2929	1146	consumer_credit	good	.	0,05	0,95			
14961	16 120_9999	MARRIED	c1	PE	24_60	24_60	FS	51_100	863	726	consumer_credit	bad	.	0,02	0,98			
14962	16 120_9999	DIVORCED	c0	MM	24_60	24_60	WB	100_9999	6559	1020	consumer_credit	good	.	0,05	0,95			
14963	13 120_9999	DIVORCED	c0	AS	24_60	24_60	HRC	100_9999	5405	2658	overhaul	bad	1,00	0,05	0,95			
14964	12 120_9999	MARRIED	c1	AS	120_9999	120_9999	WB	6_15	3299	1150	overhaul	bad	.	0,05	0,95			
14965	12 120_9999	SINGLE	c1	TM	24_60	24_60	FS	51_100	8000	760	consumer_credit	good	.	0,05	0,95			
14966	16 120_9999	SINGLE	c0	PE	18_24	0_12	FS	100_9999	1052	540	consumer_credit	good	.	0,02	0,98			
14967	16 120_9999	SINGLE	c0	TM	60_120	60_120	WB	0_5	13949	6468	flat	good	1,00	0,01	0,99			
14968	16 60_120	SINGLE	c0	PE	6_12	12_24	FS	100_9999	1020	400	consumer_credit	good	.	0,02	0,98			
14969	16 60_120	SINGLE	c0	AS	6_12	0_12	ASTL	51_100	1208	500	consumer_credit	good	.	0,05	0,95			
14970	16 120_9999	DIVORCED	c1	PE	6_12	12_24	FS	100_9999	520	332	consumer_credit							

Risk		
Sample	Estimate	Std. Error
Training	,050	,002
Test	,081	,045

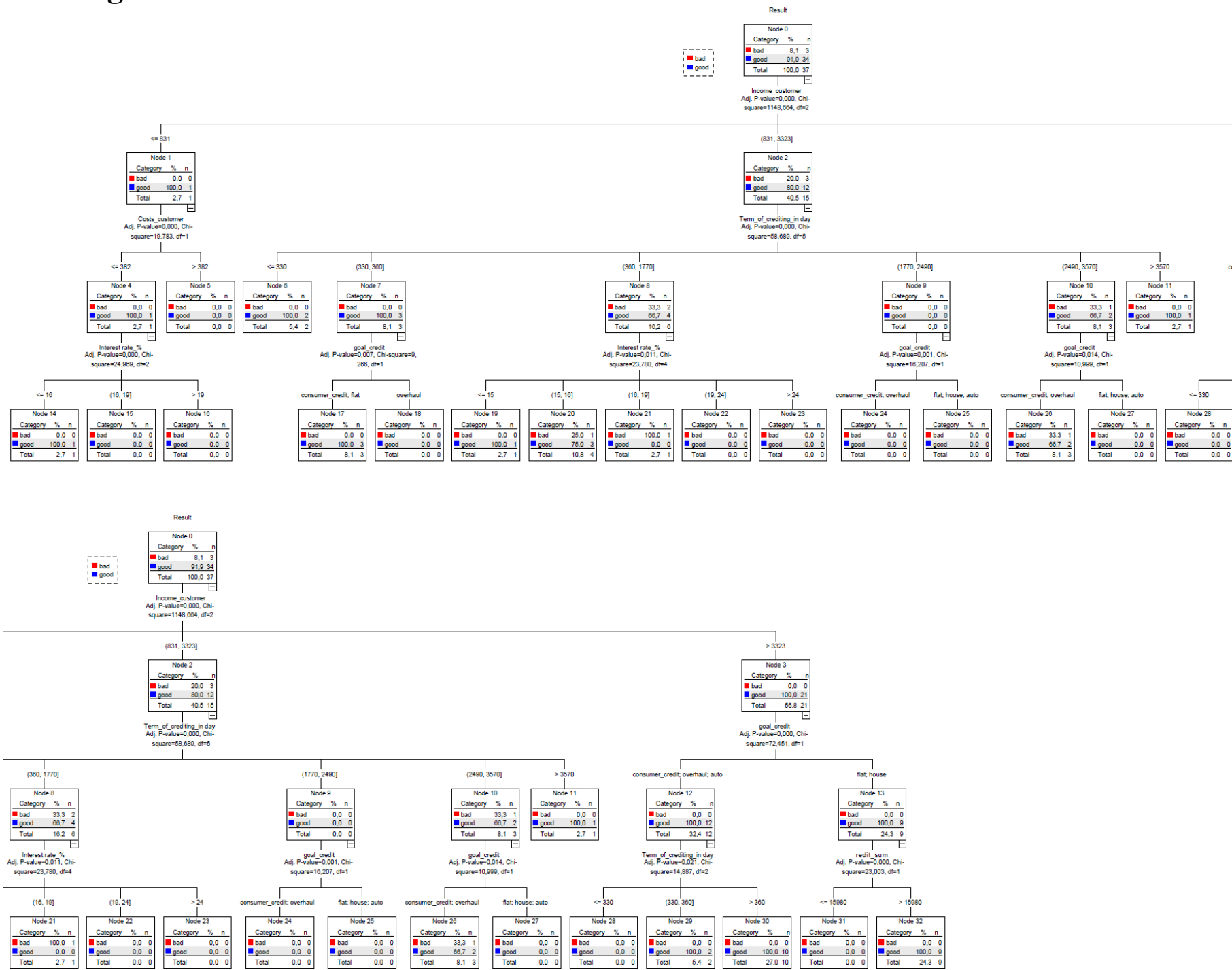
Growing Method: CRT
Dependent Variable: Result

Classification				
		Predicted		
		bad	good	Percent Correct
Training	bad	114	633	15,3%
	good	108	14108	99,2%
	Overall Percentage	1,5%	98,5%	95,0%
Test	bad	0	3	,0%
	good	0	34	100,0%
	Overall Percentage	,0%	100,0%	91,9%
Growing Method: CRT Dependent Variable: Result				

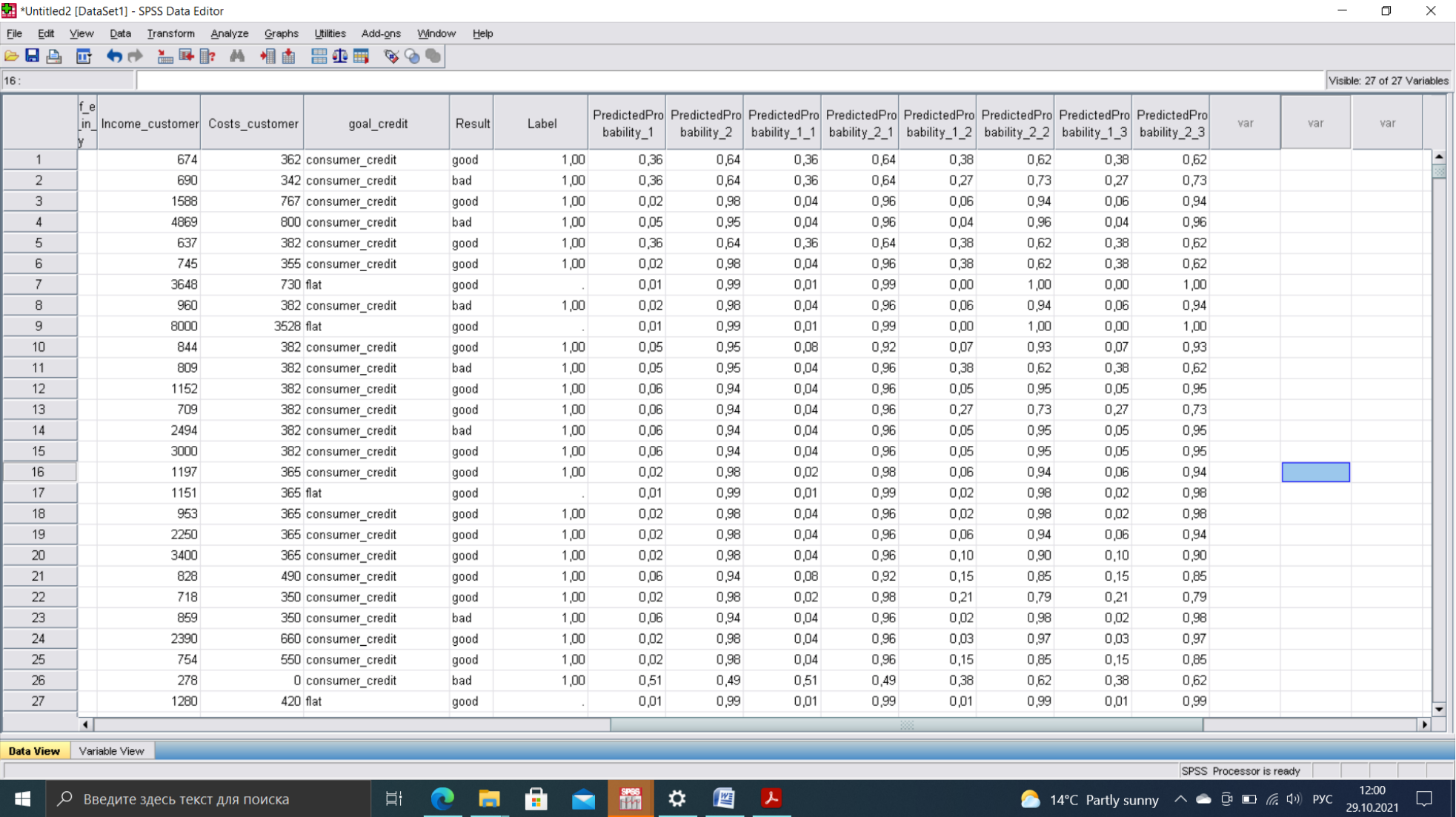
Як ми можемо бачити-це суттєво допомогло покращити якість нашого передбачення.

Розглянемо метод CHAID

Training Tree



Було прийнято рішення спробувати ще покращити результати: видаляємо з навчальної вибірки ті кейси, в яких вірогідність гарного предикту близька до 1



Risk		
Sample	Estimate	Std. Error
Training	,050	,002
Test	,070	,039

Growing Method: CHAID
Dependent Variable: Result

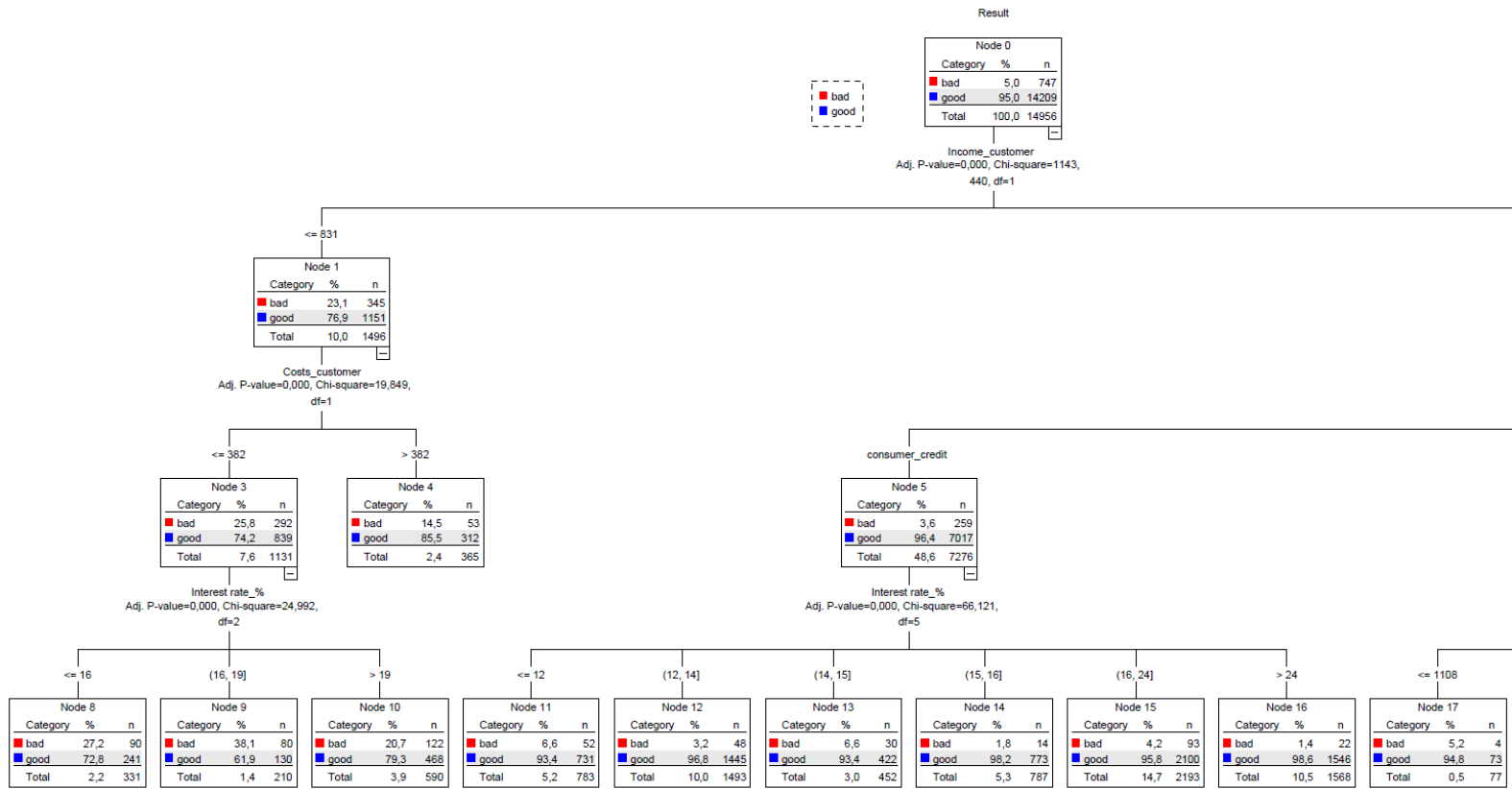
Classification				
Sample	Observed	Predicted		
		bad	good	Percent Correct
Training	bad	0	747	,0%
	good	0	14210	100,0%
	Overall Percentage	,0%	100,0%	95,0%
Test	bad	0	3	,0%
	good	0	40	100,0%
	Overall Percentage	,0%	100,0%	93,0%

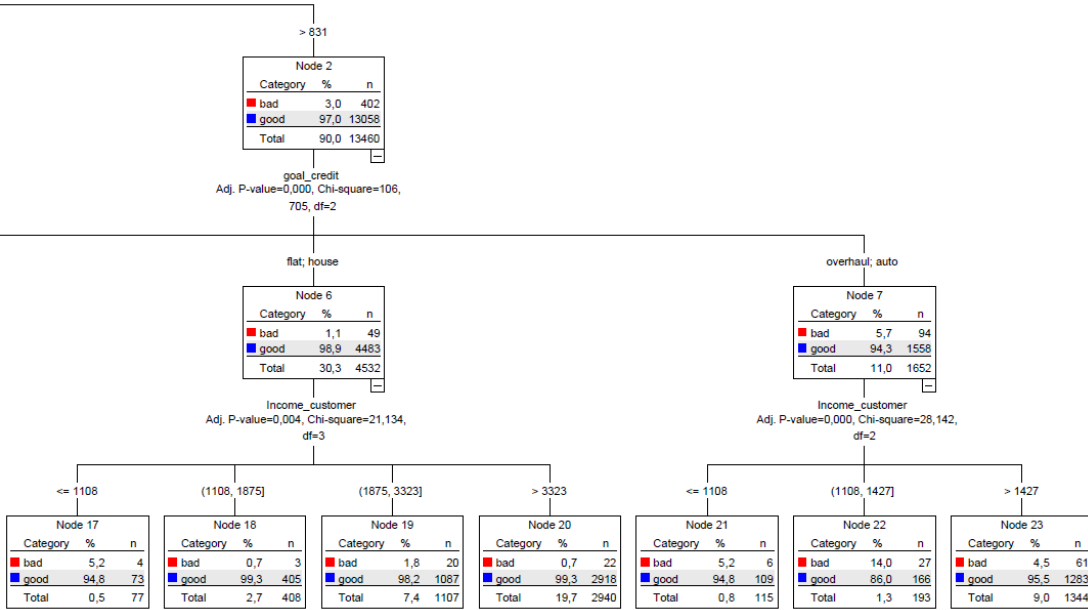
Growing Method: CHAID
Dependent Variable: Result

Якщо порівняти ці два методи, то можна сказати, що CHAID більш підходить до нашої моделі
Як ми бачимо, виділений нижче кейс у цьому методі дає результат набагато ближчий до істини.

	Result	Label	PredictedProbability_1	PredictedProbability_2	PredictedProbability_1_1	PredictedProbability_2_1	PredictedProbability_1_2	PredictedProbability_2_2	PredictedProbability_1_3	PredictedProbability_2_3	PredictedProbability_1_4	PredictedProbability_2_4
	bad	1,00	0,01	0,99	0,01	0,99	0,02	0,98	0,02	0,98	0,02	0,98
	good	1,00	0,01	0,99	0,01	0,99	0,00	1,00	0,00	1,00	0,00	1,00
	good	1,00	0,06	0,94	0,04	0,96	0,05	0,95	0,05	0,95	0,05	0,95
	good	1,00	0,05	0,95	0,04	0,96	0,04	0,96	0,04	0,96	0,04	0,96
	good	1,00	0,05	0,95	0,04	0,96	0,04	0,96	0,04	0,96	0,04	0,96
	good	1,00	0,01	0,99	0,01	0,99	0,00	1,00	0,00	1,00	0,00	1,00
	good	1,00	0,05	0,95	0,04	0,96	0,04	0,96	0,04	0,96	0,04	0,96
	bad	1,00	0,06	0,94	0,04	0,96	0,05	0,95	0,05	0,95	0,05	0,95
	bad	1,00	0,05	0,95	0,04	0,96	0,10	0,90	0,10	0,90	0,10	0,90
	good	1,00	0,05	0,95	0,08	0,92	0,10	0,90	0,10	0,90	0,10	0,90
	good	1,00	0,01	0,99	0,01	0,99	0,00	1,00	0,00	1,00	0,00	1,00
	good	1,00	0,02	0,98	0,04	0,96	0,01	0,99	0,01	0,99	0,01	0,99
	good	.	0,02	0,98	0,02	0,98	0,02	0,98	0,02	0,98	0,02	0,98
	good	.	0,05	0,95	0,04	0,96	0,27	0,73	0,27	0,73	0,27	0,73
	good	.	0,05	0,95	0,02	0,98	0,02	0,98	0,02	0,98	0,02	0,98
	good	.	0,01	0,99	0,01	0,99	0,00	1,00	0,00	1,00	0,00	1,00
	good	.	0,05	0,95	0,04	0,96	0,04	0,96	0,04	0,96	0,04	0,96
	bad	1,00	0,51	0,49	0,51	0,49	0,38	0,62	0,38	0,62	0,38	0,62
	good	.	0,05	0,95	0,04	0,96	0,04	0,96	0,04	0,96	0,04	0,96
	good	.	0,51	0,49	0,51	0,49	0,27	0,73	0,27	0,73	0,27	0,73
	good	.	0,05	0,95	0,04	0,96	0,04	0,96	0,04	0,96	0,04	0,96

Розглянемо метод CHAID-exhaustive Training Tree





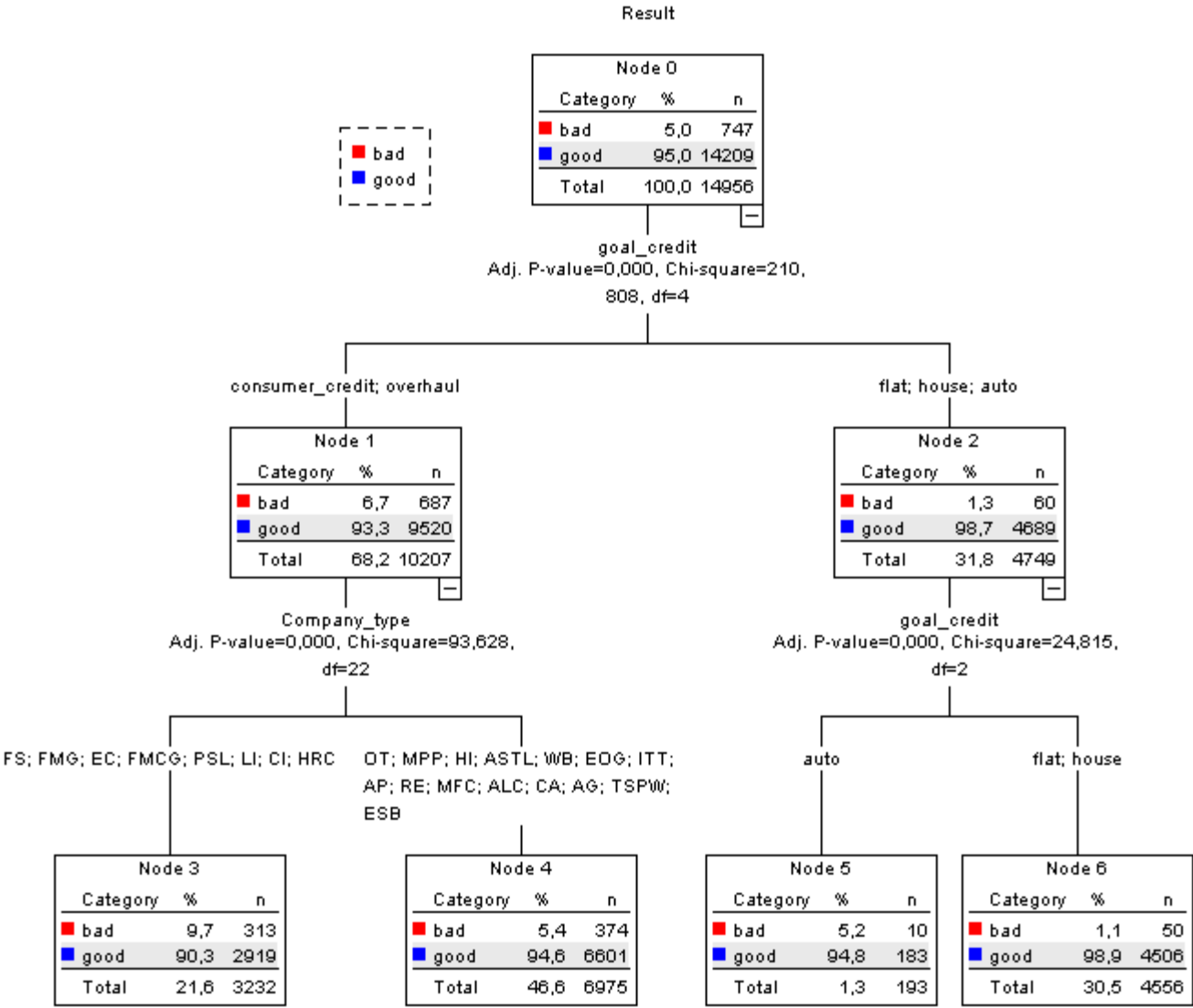
Classification				
Sample		Predicted		
		bad	good	Percent Correct
Training	bad	0	747	,0%
	good	0	14209	100,0%
	Overall Percentage	,0%	100,0%	95,0%
Test	bad	0	3	,0%
	good	0	41	100,0%
	Overall Percentage	,0%	100,0%	93,2%

Growing Method: EXHAUSTIVE CHAID
Dependent Variable: Result

Як ми можемо бачити цей метод дає ще більшу точність навіть без корегування вибірки.

Розглянемо метод QUEST

Training Tree



Risk		
Sample	Estimate	Std. Error
Training	,050	,002
Test	,068	,038

Growing Method: QUEST
Dependent Variable: Result

Classification				
		Predicted		
		bad	good	Percent Correct
Training	bad	0	747	,0%
	good	0	14209	100,0%
	Overall Percentage	,0%	100,0%	95,0%
Test	bad	0	3	,0%
	good	0	41	100,0%
	Overall Percentage	,0%	100,0%	93,2%

Growing Method: QUEST
Dependent Variable: Result

good	.	0,51	0,49	0,51	0,49	0,27	0,73	0,27	0,73	0,27	0,73	0,27	0,73	0,10	0,90		
bad	.	0,02	0,98	0,04	0,96	0,01	0,99	0,01	0,99	0,01	0,99	0,02	0,98	0,10	0,90		

Як ми можемо бачити наприкладі цих двох кейсів, метод QUEST дав нам найбільш ближчий результат до істини

ВИСНОВОК

Мною було реалізовано чотири методи дерев рішень на одній й тій самій вибірці. Спочатку я намагався розділити тестові та навчальні вибірки таким чином, щоб процент похибки був найменшим.

Реалізувавши всі ці методи, можна сказати, що всі вони не дають 100% результату, більш реалістичним методом був QUEST, але я вважаю, що якщо ми хочемо покращити наші результати предикту, то потрібно реалізовувати одразу декілька методів і на основі цих результатів аналізувати що нам підходить краще(або використовувати одразу декілька методів у комплексі)