# Choosing independent variables

- The main idea of variables selection is to reduce the number of independent variables.
- The goal is to identify the independent variables that are decently correlated with dependent variable and possibly not correlated among themselves.
- Otherwise the independent variables might be of linear relationship which may seriously damage the model.

# Choosing independent variables

Three popular methods of choosing independent variables are:

- Hellwig's method
- Graphs analysis method
- Correlation matrix method.

Three steps:

- 1. Number of combinations: 2<sup>m</sup>-1
- 2. Individual capacity of every independent variable in the combination:

$$h_{kj} = \frac{r_{0j}^2}{\sum_{i \in I_k} |r_{ij}|}$$

3. Integral capacity of information for every combination:  $H_k = \sum h_{ki}$ 

 Number of combinations
In Hellwig's method the number of combinations is provided by the formula 2<sup>m</sup>-1 where *m* is the number of independent variables.

2. Individual capacity of each independent variable in the combination is given by the formula:



where:

 $h_{kj}$  – individual capacity of information for *j*-th variable in *k*-th combination

2. Individual capacity of each independent variable in the combination is given by the formula:



where:

r<sub>0j</sub> – correlation coefficient between *j*-th variable (independent) and dependent variable

2. Individual capacity of each independent variable in the combination is given by the formula:



where:

 $r_{ij}$  – correlation coefficient between *i*-th and *j*-th variable (both independent)

2. Individual capacity of each independent variable in the combination is given by the formula:



where:

 $I_k$  – the set of numbers of variables in *k*-th combination

3. Integral capacity of information for every combination

The next step is to calculate  $H_k$  – integral capacity of information for each combination as the sum of individual capacities of information within each combination:

$$H_k = \sum h_{kj}$$

- Q: HOW TO CHOOSE INDEPENDENT VARIABLES?
- A: LOOK AT INTEGRAL CAPACITIES OF INFORMATION. THE GREATEST Hk MEANS THAT VARIABLES FROM THIS COMBINATION SHOULD BE INCLUDED IN THE MODEL.

• Let's choose independent variables, using Hellwig's method.

S	8	1	2 A A A A A A A A A A A A A A A A A A A
Y	X1	X2	Х3
12	1	5	8
14	1	6	7
17	3	6	6
20	2	8	5
25	4	6	6
30	4	9	5
36	6	9	5

• First we need to have vector and matrix of correlation coefficients.

Correlation coefficients between every independent variable X1, X2 and X3 and dependent variable Y are provided in vector  $R_0$ .

	0,9467	<b>r</b> <sub>yx1</sub> ( <b>r</b> <sub>01</sub> )
R₀	0,8374	<b>r</b> yx2 <b>(r</b> 02 <b>)</b>
	-0,7895	<b>r</b> yx3 ( <b>r</b> 03)

- First we need to have vector and matrix of correlation coefficients.
- Correlation matrix R includes correlation coefficients between independent variables.

		1	0,671	-0,712			<b>r</b> ×1×1 ( <b>r</b> 11)	<b>r</b> <sub>×1×2</sub> ( <b>r</b> <sub>12</sub> )	<b>r</b> ×1×3 ( <b>r</b> 13)
8	R	0,671	1	-0,884		R	<b>r</b> <sub>x2x1</sub> ( <b>r</b> <sub>21</sub> )	<b>r</b> <sub>x2x2</sub> ( <b>r</b> <sub>22</sub> )	<b>r</b> x2x3 ( <b>r</b> 23)
		-0,712	-0,884	1	l		<b>r</b> <sub>x3x1</sub> ( <b>r</b> <sub>31</sub> )	<b>r</b> x3x2 ( <b>r</b> 32)	<b>r</b> x3x3 ( <b>r</b> 33)

1. Number of combinations We have 3 independent variables X1, X2 and X3. Thus we may have  $2^m-1 = 2^3-1 = 8-1=7$ combinations of independent variables.

$$\{X_1\} \ \{X_2\} \ \{X_3\} \ \{X_1, X_2\} \ \{X_1, X_3\} \ \{X_2, X_3\} \ \{X_1, X_2, X_3\} \ \{X_1, X_2, X_3\} \ \{X_2, X_3\} \ \{X_3, X_3, X_3\} \ \{X_3,$$

2. Individual capacity of independent variable in the combination 1

Combination 1

 $h_{11} = \frac{0,9467^2}{1} = 0,8962 = H_1$ {X1} 0,671 0,9467 -0,712 1 0,8374 0,671 -0,884 Ro R 1 -0,7895 -0,712 -0,884 1

2. Individual capacity of independent variable in the combination 2

{X2} 
$$h_{22} = \frac{0.8374^2}{1} = 0.7013 = H_2$$



2. Individual capacity of independent variable in the combination 3

{X3} 
$$h_{33} = \frac{(-0,7895)^2}{1} = 0,6234 = H_3$$

2	0,9467		1	0,671	-0,712
Ro	0,8374	R	0,671	1	-0,884
(	-0,7895		-0,712	-0,884	

2. Individual capacity of every independent variable in the combination 4

{X1, X2} 
$$h_{41} = \frac{0,9467^2}{1+0,671} = 0,5364; \quad h_{42} = \frac{0,8374^2}{1+0,671} = 0,4197$$



*2. Individual capacity of independent variables in the combination 5* 

{X1, X3} 
$$h_{51} = \frac{0,9467^2}{1+\left|-0,712\right|} = 0,5236; \quad h_{53} = \frac{(-0,7895)^2}{1+\left|-0,712\right|} = 0,3642$$



*2. Individual capacity of every independent variables in the combination 6* 

{X2, X3} 
$$h_{62} = \frac{0,8374^2}{1+|-0,884|} = 0,3723; \quad h_{63} = \frac{(-0,7895)^2}{1+|-0,884|} = 0,3309$$



#### Combination 7

{X1, X2, X3}



3. Integral capacity of information for each combination

Combination number	Independent variables for combination	Individual capacity of information <b>h</b>		Integral capacity of information <b>H</b>	
1	{X <sub>1</sub> }	<b>h</b> <sub>11</sub>	0,8962	<b>H</b> <sub>1</sub>	0,89621
2	{X <sub>2</sub> }	<b>h</b> <sub>22</sub>	0,7013	$H_2$	0,7013
3	{X <sub>3</sub> }	<b>h</b> 33	0,6234	H <sub>3</sub>	0,62338
Λ	(Y. Y.)	<b>h</b> 41	0,5364	H <sub>4</sub> 0,	0 05612
4	<b>{</b> ^1, ^2 <b>}</b>	<b>h</b> 42	0,4197		0,35012
5	(Y. Y.)	<b>h</b> 51	0,5236	ц.	0 99796
5	{ <b>^</b> 1, <b>^</b> 3}	<b>h</b> 53	0,3642	Π5	0,00700
6	(Y. Y.)	<b>h</b> 62	0,3723	н	0 70346
v	{^2, ^3}	<b>h</b> 63	0,3309	116	0,70510
7		<b>h</b> <sub>71</sub>	0,3762		
	$\{X_1, X_2, X_3\}$	<b>h</b> <sub>72</sub>	0,2745	$H_7$	0,89089
		<b>h</b> <sub>73</sub>	0,2402		

The greatest integral capacity is for combination C4.
Independent variables
X1, X2 - will be included in model.

Three steps

- 1. Calculating r\*
- 2. Modification of correlation matrix
- 3. Drawing the graph

- *Q: HOW TO CHOOSE INDEPENDENT VARIABLES?*
- A: LOOK AT THE GRAPHS. THE NUMBER OF **GROUPS MEANS THE NUMBER OF VARIABLES** INCLUDED IN THE MODEL. IF THERE'S SEPARATED (ISOLATED) VARIABLE, YOU SHOULD INCLUDE IT IN THE MODEL. FROM EACH GROUP, THE VARIABLE WITH THE GREATEST NUMBER OF LINKS SHOULD BE INCLUDED IN MODEL. IF THERE'S TWO VARIABLES WITH THE GREATEST NUMBER OF LINKS, YOU SHOULD TAKE THE VARIABLE WHICH IS MORE STRONGLY CORRELATED WITH DEPENDENT VARIABLE.

### 1. Calculating r\*

We start with calculating critical value of r\* using the formula:

$$r^* = \sqrt{\frac{t_\alpha^2}{n-2+t_\alpha^2}}$$

where  $t_{\alpha}$  is provided in the table of t-Student distribution at the significance level  $\alpha$  and the degrees of freedom *n-2* (sometimes r\* can be given, so there's no need to calculate it).

### 2. Modification of correlation matrix

The correlation coefficients for which  $|r_{ij}| \le r^*$  are statistically irrelevant and we replace them with nulls in correlation matrix.

- 3. Drawing the graph
- Using modified correlation matrix we draw the graphs with bulbs representing the variables and the links representing correlation coefficients of statistical significance.

Let's have an example (the same one as for Hellwig's method, n=7)

	0,9467		1	0,671	-0,712
Ro	0,8374	R	0,671	1	-0,884
	-0,7895		-0,712	-0,884	1

1. Calculating r\* (n=7,  $t_{\alpha, n-2} = t_{0,05,5} = 2,571$ )

$$r^* = \sqrt{\frac{t_{\alpha}^2}{n-2+t_{\alpha}^2}} = \sqrt{\frac{2,571^2}{5+2,571^2}} = \sqrt{\frac{6,61}{11,61}} = \sqrt{0,569337} = 0,7545$$

### 2. Modification of correlation matrix

 $|r^*| \le 0,7545$ 

	1	0,671	-0,772
R	0,871	1	-0,884
	-0,712	-0,884	1

	1	0	0
R*	0	1	-0,884
	0	-0,884	1

### 3. Drawing the graph

	1	0	0
R*	0	1	-0,884
	0	-0,884	1



Conclusion: Model will consist of X1 (as isolated variable) and x2 (cause is more strongly correlated with dependent variable – you may check it in  $R_0$  vector).

	0,9467
Ro	0,8374
	-0,7895
	31

### Correlation matrix method

### 1. Calculate r\*

We start with calculating critical value of r\* using the formula:

$$r^* = \sqrt{\frac{t_\alpha^2}{n-2+t_\alpha^2}}$$

where  $t_{\alpha}$  is provided in the table of t-Student distribution at the significance level  $\alpha$  and the degrees of freedom *n-2* (sometimes r\* can be given, so there's no need to calculate it). 2. To eliminate Xi variables weakly correlated withY

$$\left|r_{ij}\right| \leq r^*$$

3. To choose  $X_s$  where  $|r_s| = \max\{|r_i|\}$ 

[Xs is the best source of information]

4. To eliminate Xi variables strongly correlated with Xs

$$\left|r_{si}\right| > r^*$$

Let's have an example (the same one as for Hellwig's method and graph analysis metod, n=7)

	0,9467		1	0,671	-0,712
Ro	0,8374	R	0,671	1	-0,884
	-0,7895		-0,712	-0,884	1

1. Calculating r\* (n=7,  $t_{\alpha, n-2} = t_{0,05,5} = 2,571$ )

$$r^* = \sqrt{\frac{t_{\alpha}^2}{n - 2 + t_{\alpha}^2}} = \sqrt{\frac{2,571^2}{5 + 2,571^2}} = \sqrt{\frac{6,61}{11,61}} = \sqrt{0,569337} = 0,7545$$

2. To eliminate Xi variables weakly correlated withY

$$|r_{ij}| \le r^*$$
  $r^* = 0,7545$ 



None of the variables will be eliminated

### 3. To choose $X_s$ where $|r_s| = \max\{|r_i|\}$

	0,9467		1	0,671	-0,712
Ro	0,8374	R	0,671	1	-0,884
	-0,7895		-0,712	-0,884	1

4. To eliminate Xi variables strongly correlated with Xs

$$|r_{si}| > r^*$$
  $r^* = 0,7545$ 

	0,9467		1	0,671	-0,712	
Ro	0,8374	R	0,671	1	-0,884	
	-0,7895		-0,712	-0,884	1	

None of the variables will be eliminated. X1, X2, X3 will be included in model. In this example level of significane can be changed – this will give us different results (you may check it if you want).

### DON'T EXPECT TO GET THE SAME RESULTS FROM THESE THREE METHODS...