

1.- Question or problem definition.

We will try to identify using the data presented in the exercise database whether there is a direct relationship between the work performed by each of the bank's customers and their account balance in order to identify whether there is a group that is more attractive to the bank because of the banking security they offer.

2.- Acquire training and testing data.

The data is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls.

3.- Wrangle, prepare, cleanse the data.

In order to be able to answer this question we will extract from the database the work column where the type of work performed by the client appears and the balance column which indicates what has happened in the last year in the client's bank account, either positive or negative.

We will then divide this information into thirteen columns, one for each type of work available, and with the values of all the available balances.

4.- Analyze, identify patterns, and explore the data.

the first step in analyzing the data was to find out some statistical data on the different jobs to see if there were any differences.

<i>admin.</i>		<i>blue-collar</i>		<i>entrepreneur</i>		<i>housemaid</i>		<i>management</i>	
Media	1226.736	Media	1085.162	Media	1645.125	Media	2083.804	Media	1766.929
Error típico	108.4067	Error típico	66.33322	Error típico	342.654	Error típico	435.0217	Error típico	104.9747
Mediana	430	Mediana	408.5	Mediana	365.5	Mediana	296.5	Mediana	577
Moda	0	Moda	0	Moda	0	Moda	0	Moda	0
Desviación estándar	2370.119	Desviación estándar	2040.218	Desviación estándar	4441.304	Desviación estándar	4603.837	Desviación estándar	3267.733
Varianza de la muestra	5617465	Varianza de la muestra	4162490	Varianza de la muestra	19725178	Varianza de la muestra	21195312	Varianza de la muestra	10678079
Curtosis	31.25918	Curtosis	18.09018	Curtosis	45.66067	Curtosis	13.79546	Curtosis	17.04516
Coefficiente de asimetría	4.748747	Coefficiente de asimetría	3.750573	Coefficiente de asimetría	6.049316	Coefficiente de asimetría	3.538975	Coefficiente de asimetría	3.642617
Rango	23138	Rango	17753	Rango	44127	Rango	27724	Rango	29105
Mínimo	-967	Mínimo	-1400	Mínimo	-2082	Mínimo	-759	Mínimo	-1746
Máximo	22171	Máximo	16353	Máximo	42045	Máximo	26965	Máximo	27359
Suma	586380	Suma	1026563	Suma	276381	Suma	233386	Suma	1712154
Cuenta	478	Cuenta	946	Cuenta	168	Cuenta	112	Cuenta	969
<i>retired</i>		<i>self-employed</i>		<i>services</i>		<i>student</i>		<i>technician</i>	
Media	2319.191	Media	1392.41	Media	1103.957	Media	1543.821	Media	1330.996
Error típico	385.4989	Error típico	183.3003	Error típico	119.7439	Error típico	281.4887	Error típico	94.91109
Mediana	672.5	Mediana	483	Mediana	288	Mediana	422.5	Mediana	434.5
Moda	0	Moda	0	Moda	0	Moda	0	Moda	0
Desviación estándar	5846.38	Desviación estándar	2479.641	Desviación estándar	2445.24	Desviación estándar	2579.887	Desviación estándar	2630.253
Varianza de la muestra	34180158	Varianza de la muestra	6148619	Varianza de la muestra	5979199	Varianza de la muestra	6655815	Varianza de la muestra	6918233
Curtosis	86.43204	Curtosis	12.65537	Curtosis	42.60244	Curtosis	5.134617	Curtosis	30.61018
Coefficiente de asimetría	8.071526	Coefficiente de asimetría	3.183122	Coefficiente de asimetría	5.402962	Coefficiente de asimetría	2.368667	Coefficiente de asimetría	4.494781
Rango	72394	Rango	19743	Rango	27596	Rango	11785	Rango	29413
Mínimo	-1206	Mínimo	-3313	Mínimo	-1202	Mínimo	-230	Mínimo	-1680
Máximo	71188	Máximo	16430	Máximo	26394	Máximo	11555	Máximo	27733
Suma	533414	Suma	254811	Suma	460350	Suma	129681	Suma	1022205
Cuenta	230	Cuenta	183	Cuenta	417	Cuenta	84	Cuenta	768
				<i>unemployed</i>					
				Media	1089.422				
				Error típico	149.5767				
				Mediana	473.5				
				Moda	0				
				Desviación estándar	1692.268				
				Varianza de la muestra	2863770				
				Curtosis	6.674153				
				Coefficiente de asimetría	2.350104				
				Rango	9891				
				Mínimo	-872				
				Máximo	9019				
				Suma	139446				
				Cuenta	128				

Analyzing these first statistical data, we realize that there are notable differences in the balances if we group the sample into jobs.

We know that, on average, retirees and household employees are the people with the highest balance, but they are also the jobs in which there is the highest variance.

5.- Model

To verify that the differences between the different papers are significant and that there is a considerable difference depending on the paper, we performed an ANOVA analysis.

SUMMARY

<i>Groups</i>	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>
admin.	477	564209	1182.828	4705775
blue-collar	945	1010210	1069.005	3919703
entrepreneur	167	234336	1403.21	9952900
housemaid	111	206421	1859.649	15709349
management	968	1684795	1740.491	10011117
retired	229	462226	2018.454	13436982
self-employed	182	238381	1309.786	4926393
services	416	433956	1043.163	4448730
student	83	118126	1423.205	5500018
technician	767	994472	1296.574	6016070
unemployed	127	130427	1026.984	2383535

ANOVA

<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
Between Groups	4.22E+08	10	42215888	6.243729	1.47E-09	1.83282
Within Groups	3.02E+10	4461	6761327			
Total	3.06E+10	4471				

With the ANOVA analysis done, we can prove that with a 95% certainty there are notable differences in the means of the variables or jobs we are evaluating.