

Grado en Inteligencia Artificial
PMAAD

Advanced Profiling

Prof. Dante Conti
Prof. Sergi Ramirez

Créditos Prof. Tomás Aluja

Profiling. What is the profile of a group of individuals?

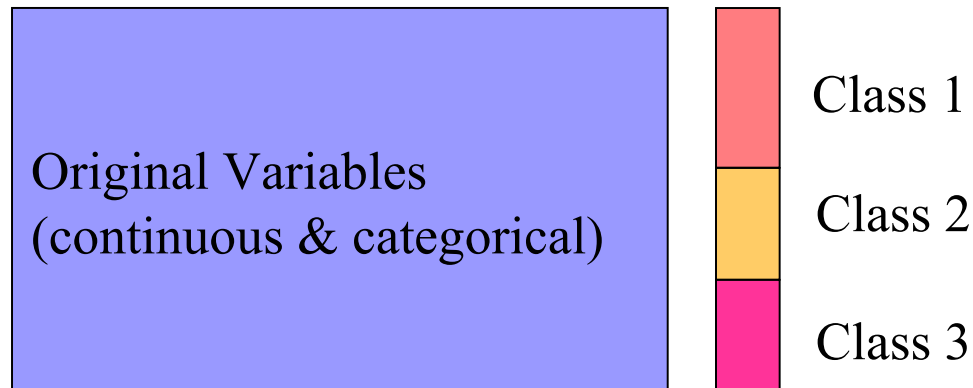
i.e.: a class resulting from a clustering algorithm, or a modality of a categorical variable. What is the profile of A buyers?

Profiling is finding the significant characteristics which make the group of individuals different than the whole set of individuals.

→ *Automatic detection of significant deviations*

- Differential characterisation among classes

Categorical variable or a partition



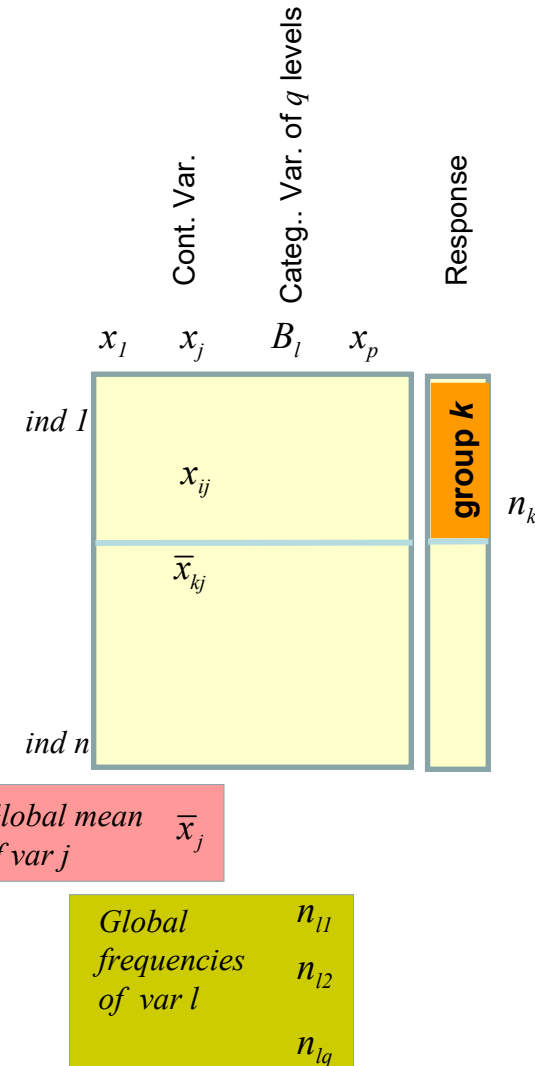
- Statistical characterization

We have a group of individuals defined by a level of a categorical variable.

- We take as response variable the variable identifying the groups that we want to find their profile.
- The explanatory variables are all the others and can be either categorical or continuous.

Problem: For every group of individuals detect

1. which modalities (of the categorical explanatory variables) deviate significantly from what were expected and,
2. which continuous explanatory variables, deviate significantly from what were expected.



Hypothesis test with continuous variables

groups	means	counts
1	\bar{x}_1	n_1
\vdots	\vdots	\vdots
q	\bar{x}_q	n_q

Global \bar{x} n

$$H_0 : \mu_k = \mu \quad k = 1, \dots, q$$

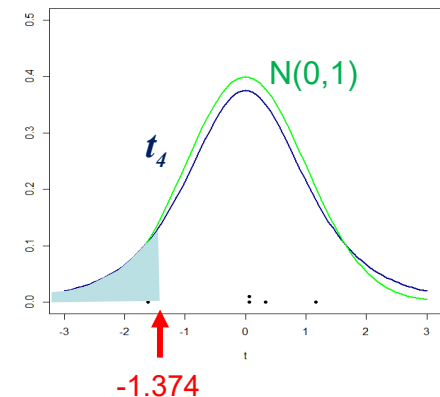
Test statistic: Difference between the mean in group k and the global mean

Student's t

$$t = \frac{\bar{x}_k - \bar{x}}{\sqrt{(1 - \frac{n_k}{n}) \frac{s^2}{n_k}}} : t_{n-1}$$



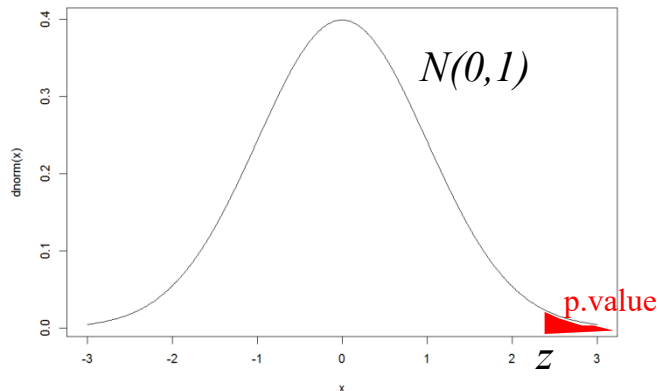
William Gosset "Student",
English, 1876, 1937



Rank the continuous variables by p.value (ascending)

	1... j ... J	
1	\vdots	
k	$\cdots n_{kj} \cdots$	n_k
q	\vdots	
	n_j	

Test statistic: Difference between proportion of modality j in group k and proportion of modality j in whole data



$$H_0 : p_{j \cdot k} = p_j \quad k = 1, \dots, q; j = 1, \dots, J$$

Assumption of normality of proportions:

$$\frac{n_{kj}}{n_k} : N\left(p_j = \frac{n_j}{n}, \left(1 - \frac{n_k}{n}\right) \frac{p_j(1 - p_j)}{n_k}\right)$$

$$Z = \frac{\frac{n_{kj}}{n_k} - \frac{n_j}{n}}{\sqrt{\left(1 - \frac{n_k}{n}\right) \left(\frac{p_j(1 - p_j)}{n_k}\right)}} : N(0,1)$$

Rank the levels of the categorical explanatory variables by *p.value* (ascending)

- For each group to profile,
 1. rank the modalities of the categorical explanatory variables according their p-value (ascending).
 2. Likewise, rank the continuous variables according their p-value
- Select the most significant modalities and continuous variables by a threshold (0.05, 0.01, ..) defined a priori.
(what matters is the ordering, actual significance depends on the number of individuals)

Interpreting the classes of expenses

```
> library(FactoMineR)
> catdes(data, num.var, proba = 0.05, row.w = NULL)
```

```
> catdes(cbind(as.factor(c1),despeses),1)
```

```
$category$`2` Cla/Mod Mod/Cla Global p.value v.test
STATUS=3      75      100 33.33333 0.01818182 2.361894
```

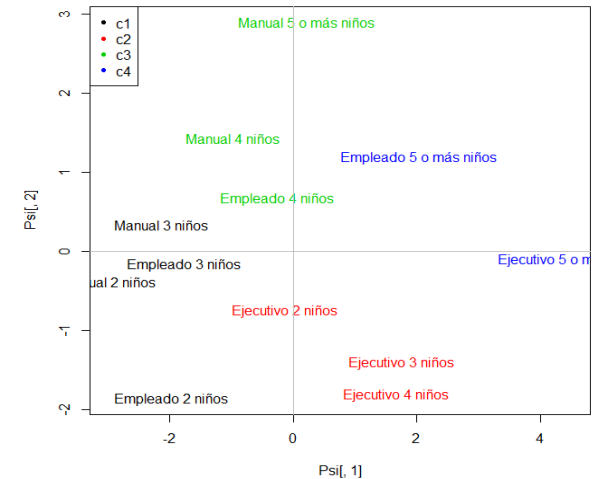
```
$category$`4` Cla/Mod Mod/Cla Global p.value v.test
CHILDREN=5    66.66667      100      25 0.04545455 2.000424
```

```
$quanti$`1` v.test Mean in category Overall mean sd in category Overall sd p.value
fruits      -2.006717      369.75      505.0000      22.89514      158.0638 0.04477978
bread       -2.112728      354.25      446.6667      44.53299      102.5860 0.03462408
meat        -2.437108     1493.00     1886.7500      33.73426      378.9022 0.01480525
vegetables  -2.492476      539.50      732.0000      67.18817      181.1261 0.01268558
poultry     -2.497153      548.75      803.1667      15.48184      238.9361 0.01251951
```

```
$quanti$`2` v.test Mean in category Overall mean sd in category Overall sd p.value
poultry     2.175824     1074.667      803.1667     104.4169     238.9361 0.02956844
```

```
$quanti$`4` v.test Mean in category Overall mean sd in category Overall sd p.value
vegetables   2.571343      1046.0      732.00      51.0      181.1261 0.01013051
milk         2.397322      539.5      358.25      21.5      112.1406 0.01651540
fruits       1.994058      717.5      505.00     169.5      158.0638 0.04614569
```

Clustering of expenses in 4 classes





Sir Ronald A. Fisher
English, 1890-1962

To be used for feature selection

levels	means	counts
1	\bar{x}_1	n_1
\vdots	\vdots	\vdots
q	\bar{x}_q	n_q

$$H_0 : \mu_1 = \dots = \mu_q = \mu$$

\bar{x}, s^2 global mean and global variance

$$F = \frac{S_B^2 / (q-1)}{S_W^2 / (n-q)} : F_{q-1, n-q}$$

Ranking of variables by p.values (ascending)

$$S_W^2 = \sum_{k=1}^q \sum_{i=1}^{n_k} (x_{ki} - \bar{x}_k)^2$$

$$S_B^2 = \sum_{k=1}^q n_k (\bar{x}_k - \bar{x})^2$$

Equivalent statistic:

$$\eta^2 = \frac{S_B^2}{S_W^2 + S_B^2}$$

Dependence of a partition respect to a categ. Var.



Karl Pearson
English, 1857, 1936

To be used for feature selection

	1...	j	...	J
1		\vdots		
k	...	n_{kj}	...	n_k
q		\vdots		
		n_j		

$$H_0 : n_{kj} = np_k p_j \quad \forall kj$$

$$X^2 = \sum_{k=1}^p \sum_{j=1}^q \frac{(n_{kj} - \frac{n_k n_j}{n})^2}{\frac{n_k n_j}{n}} : \chi^2_{(p-1)(q-1)}$$

Ranking of variables by p.values (ascending)

Improvement when Interpreting Clusters Results

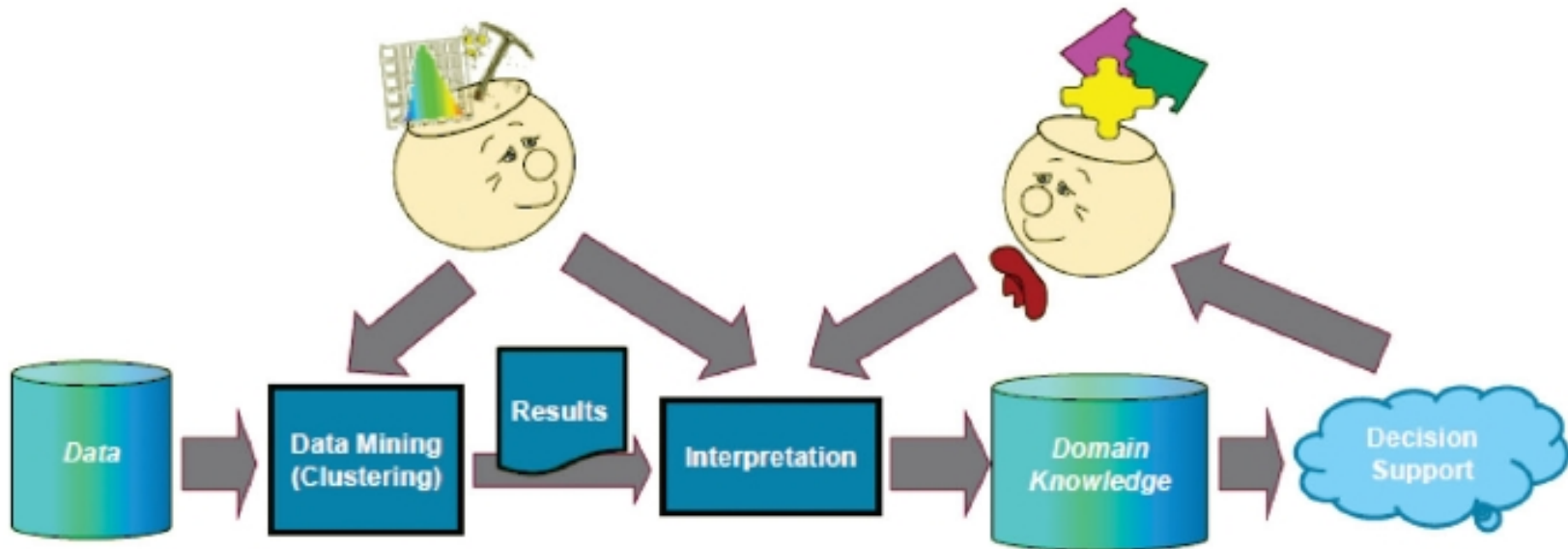
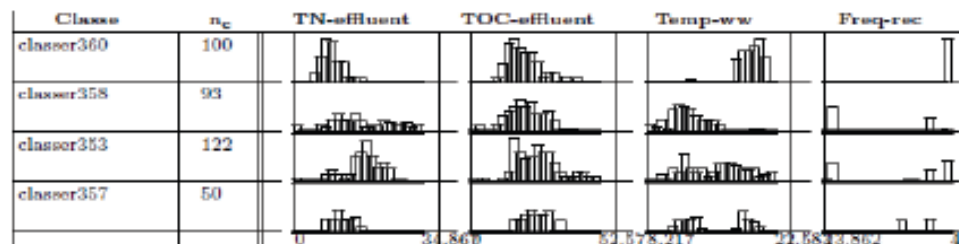
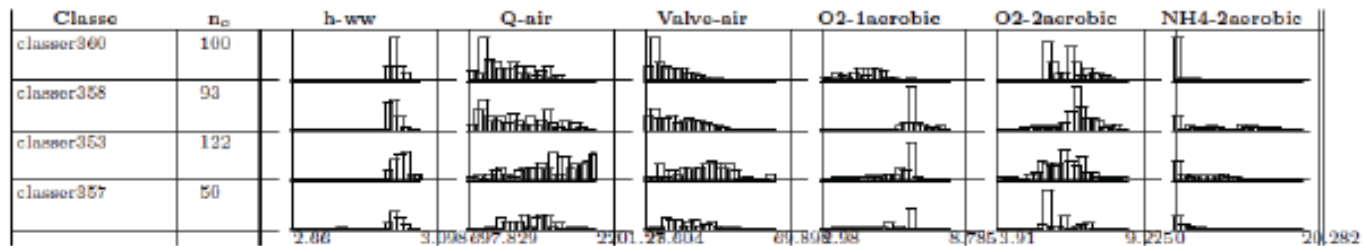
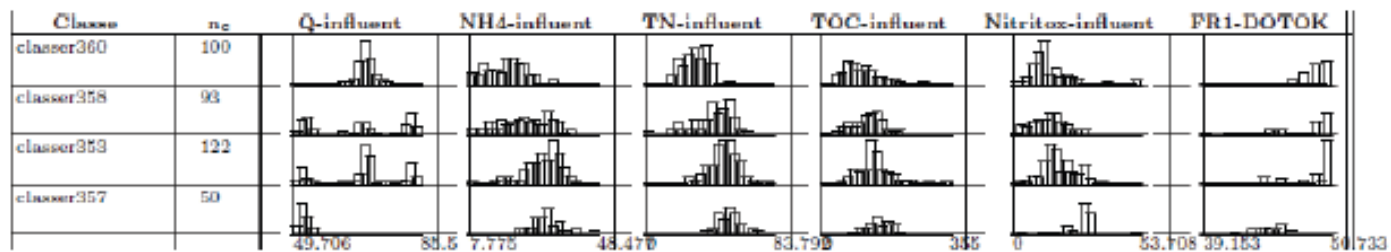
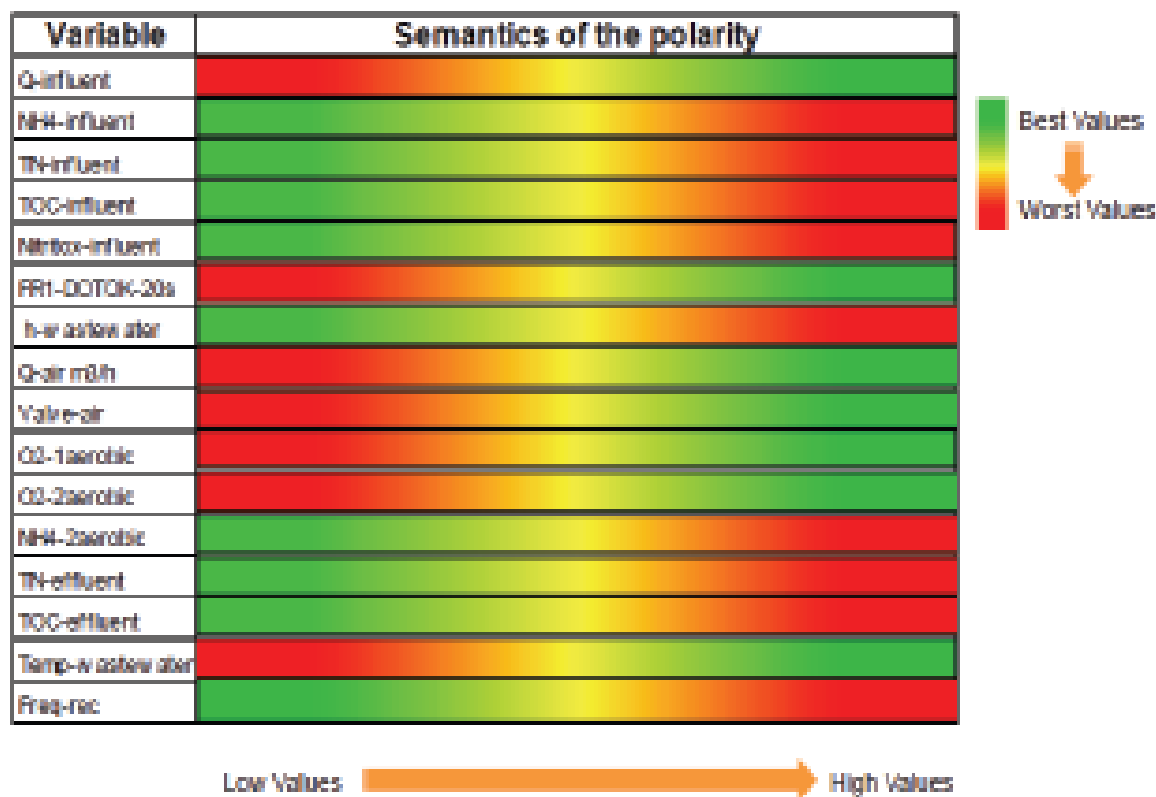


Figure 1. Interpretation support tools bridging the gap between DM and decision support.

CPG → Class Panel Graph



Semantics and Polarity




TLP → Traffic Light Panel

The analyst has to read the CPG, mark the characteristic values of the variables for the different classes and assign qualitative levels to those values. The dominant level of a variable in a given class can be determined in two ways:

- Identify the qualitative level of the mean or median of the variable in the class,
- Identify the qualitative level of the mode of the variable in the class.

When the variables have only 3 qualitative levels assigned, it is very interesting to assign the colours of a traffic light to those levels (red for the bad or negative value, yellow for the medium or neutral value and green for the good or positive value), being the *bad* values the higher or lower values of the variable depending on the variable's semantics. It is

		VARIABLES						
Class	nc	CANTV	EDC	MVZB	BVP	IGBC	BOVESPA	DJIA
C357	90	Yellow	Yellow	Yellow	Yellow	Yellow	Red	Yellow
C356	123	Yellow	Yellow	Yellow	Yellow	Yellow	Green	Green
C359	33	Green	Green	Green	Green	Green	Green	Green
C346	16	Red	Red	Red	Red	Red	Yellow	Yellow
C347	102	Yellow	Yellow	Yellow	Red	Yellow	Green	Green



Best

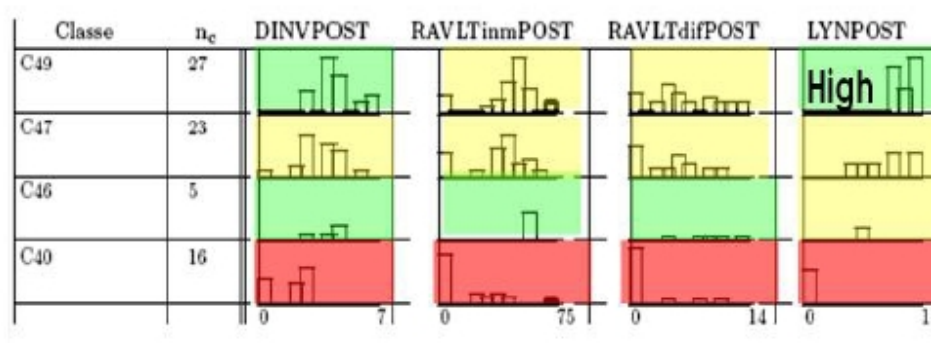
Neut

Worse

Figure 2. Traffic light panel (colour online only).

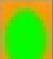


TLP → Traffic Light Panel (How to get it?)

1. Perform the Class panel graph (CPG) of all the variables versus the discovered classes
2. Calculate the basic statistics per class
3. Using materials from steps 1 and 2 identify variables or combination of variables with specific ranges of values in a class that distinguish the class from the others.
4. Assign qualitative levels to the variables implied in step 3 by detecting the area where the mass of the distribution is placed.
5. Perform a TLP for the variables, using the qualitative values assigned in step 4
6. Show the Traffic lights panel to the expert and ask him to select a label for the class. The expert is conceptualizing the class in this step, on the basis of the Traffic lights panel
7. Perform significance tests assessing relevance of differences for the variables implied in step 6, for instance: ANOVA, Kruskal-Wallis or χ^2 independence tests depending on the item.












TLP → Improvements (Grouping variables)

			ATENCIÓN			MEMORIA				FUNCIONES EJECUTIVAS				
	Class	r_p	DDP	TMTA	PALCOL	DIN	RAVLT _{mem}	RAVLT _{vis}	DIN	TMTB	INTERF	CATEG	PERSEV	PMB
POST	C19	17	Good	Good	Mode	Good	Mode	Mode	Good	Good	Good	Good	Good	Mode
	C17	18	Good	Good	Mode	Mode	Mode	Mode	Mode	Good	Good	Mode	Good	Mode
	C46	5	Good	Mode	Mode	Good	Good	Good	Mode	Mode	Mode	Good	Good	Mode
	C40	16	Mode	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad
DIF	C19	17	Bad	Bad	Bad	Mode	Bad	Bad	Mode	Mode	Bad	Mode	Mode	Mode
	C17	18	Good	Good	Good	Mode	Mode	Mode	Good	Good	Good	Good	Good	Good
	C46	5	Good	Bad	Bad	Good	Good	Good	Good	Bad	Bad	Mode	Mode	Bad
	C40	16	Mode	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad	Bad

 Good
 Mode
 Bad

aTLP →
annotated Traffic Light Panel

$$VC_k|C = \frac{s_{X_k|C}}{\bar{X}_k|C}.$$

Low VC			
Intermediate VC			
High VC			

aTLP →
annotated Traffic Light Panel

Class		nc	influent						Aerobic Tanks 1 & 2 - Anoxic Tank 2						Effluent			Other	VC < 0.30				
			Q	NH4	TN	TOC	Ni Tri tox	FR1- DO TOK	k- ww	Q- air	Val- ve air	O2-1 aero- bic	O2-2 aero- bic	NH4 aero- bic	TN	TOC	Temp ww		Frec rec	0.30<VC<0.90			
																				VC>0.90			
C360	100																						
C358	93																						
C353	122																						
C357	50																						

Class		nc	influent						Aerobic Tanks 1 & 2 - Anoxic Tank 2						Effluent			Other	VC < 0.10				
			Q	NH4	TN	TOC	Ni Tri tox	FR1- DO TOK	k- ww	Q- air	Val- ve air	O2-1 aero- bic	O2-2 aero- bic	NH4 aero- bic	TN	TOC	Temp ww		Frec rec	0.10<VC<0.50			
																				VC>0.50			
C360	100																						
C358	93																						
C353	122																						
C357	50																						