



Universidad Nacional  
Autónoma de México



Facultad  
de  
Ciencias

# Integrations distinct sources databases to improve the estimation models

October 14th, 2021, Virtual IT CONFIDENCE  
CONFERENCE



**ISBSG**

**IT Confidence 2021**

**Dr. Francisco Valdés-Souto**

Associate Professor  
Mathematics Department of  
Science Faculty at  
National Autonomous University of Mexico (UNAM)  
[fvaldes@Ciencias.unam.mx](mailto:fvaldes@Ciencias.unam.mx)

*COSMIC President*

*Mexican Software Metrics Association (AMMS), Founder*



01

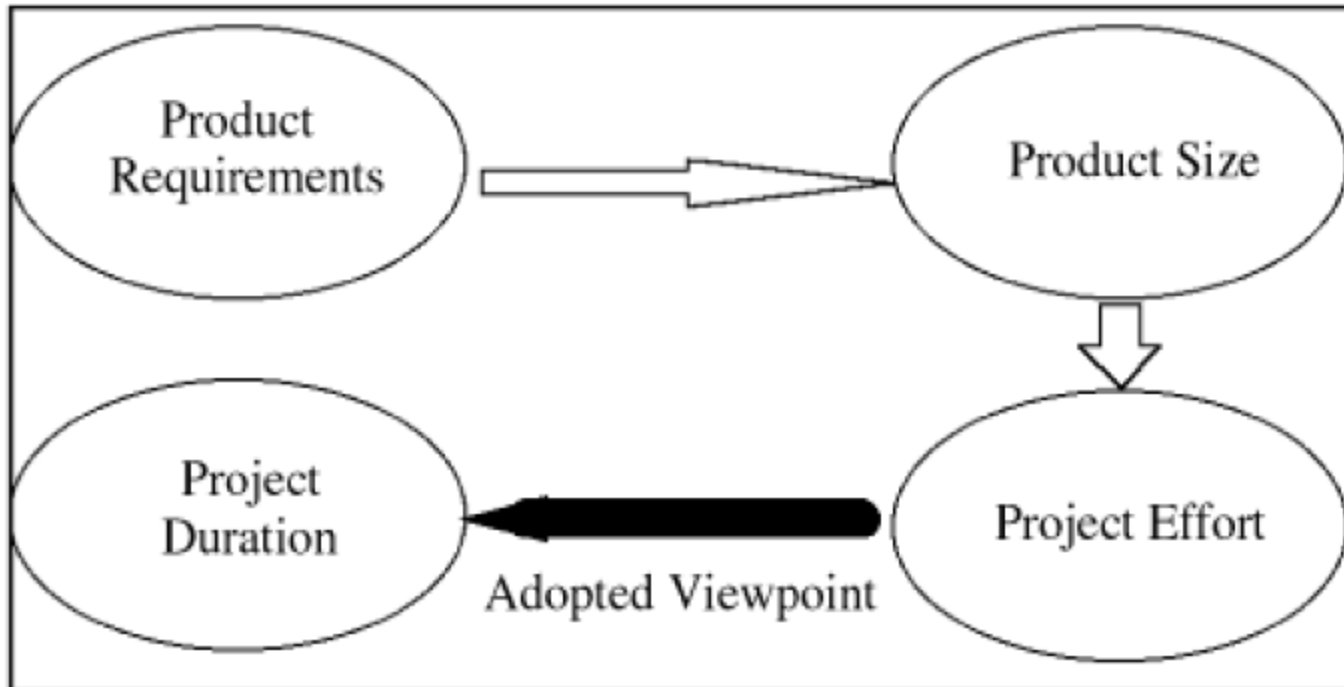
**Problems with datasets in industry**

02

**Case study explanation**

03

**3 Case study results**

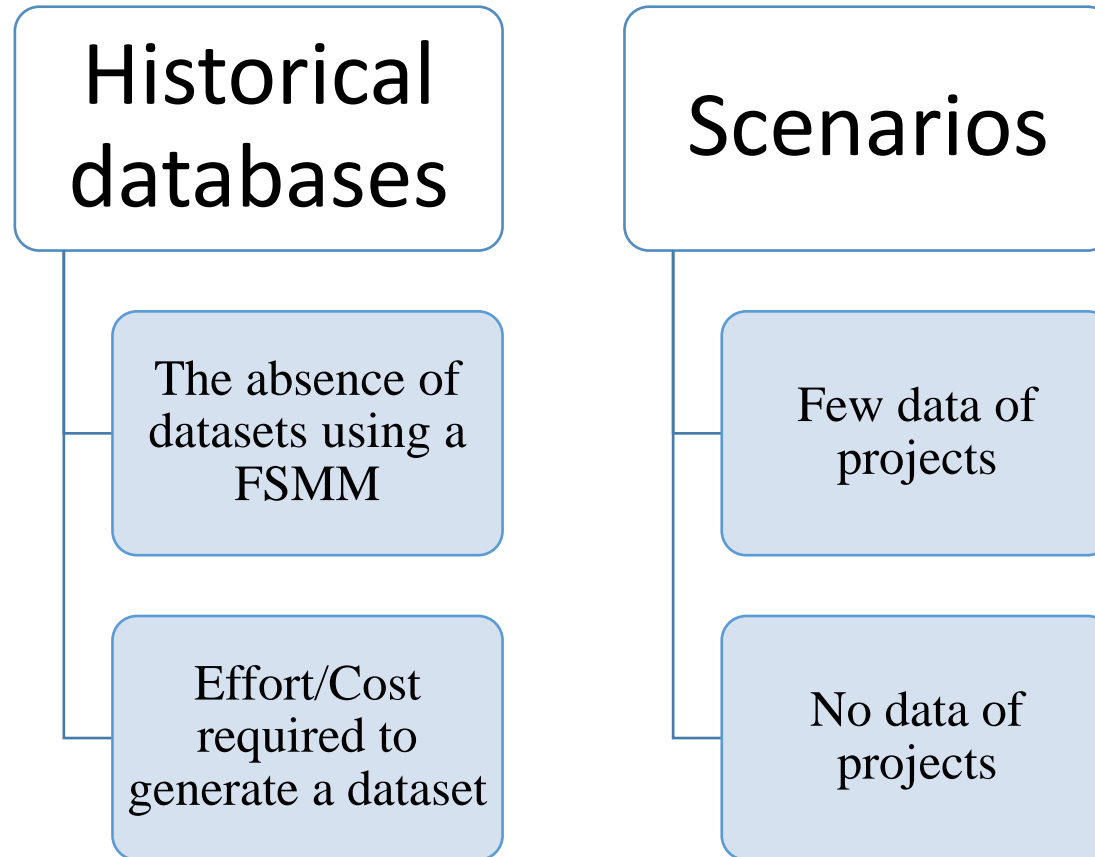


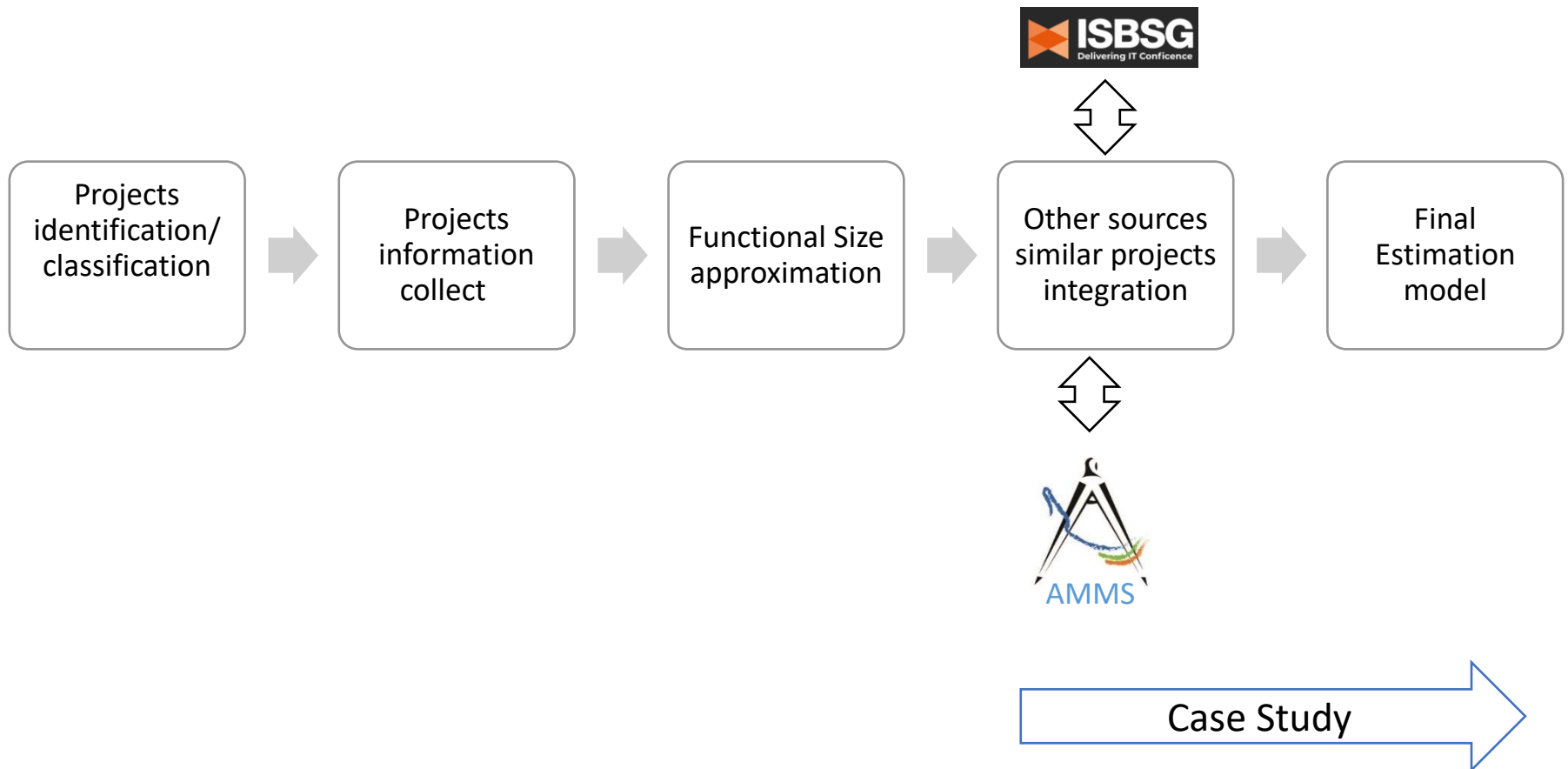
*Bourque, 2007*

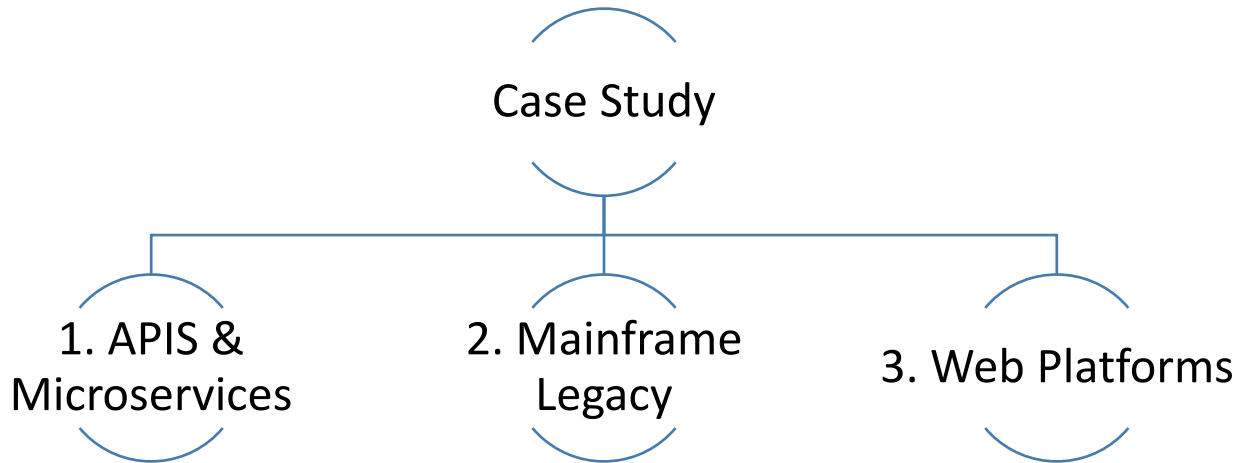
Several authors identify the measurement of the size of the piece of software as a relevant factor in the precision of the estimate (Linda, 2006) (Koch,2009) (De Lucia, 2005) (Hill, 2000)












- Morgenshtern pointed out that “Algorithmic models **need historic data**, and many organizations do not have this information. Additionally, collecting such data may be both **expensive and time consuming**.” (Morgenshtern, 2007).
- The majority of the estimation models developed are dependent on the **representativeness of the samples** (databases) used.
- Abran mentions that “most of the so-called estimation models in the literature are productivity models.” They **represent the past behavior for a specific organization** to develop software projects, representing the relationship across the two variables, usually the functional size as an independent variable and the effort or cost as dependent variables. (Abran 2015)
- In order to generate estimation models, **the researchers** have used databases documented on the **basis of the past completed projects they participated in**, usually, this **information is not available** to all the persons or is **difficult to acquire** or has elements that **do not make sense** for all the database’s users.

- Jørgensen et al. [1], in a systematic review of estimation studies, found “that there are good reasons to claim that the **availability of a data set is more indicative for its use than its representativeness or other properties**”.
- Braga et al. [3] mention they do not found “**any reliable information about the way in which the projects included in a dataset were obtained,**”
- Carbonera et al. [7] analyze the number of data points in the datasets and classify in **high quality** (more than 15 points) **medium quality** (10 to 15 data points and ) or **low quality** (less than 10 data points), where it is possible to observe the **lack of datasets with a high number of data points, a main statistical principle.**
- Carbonera et al. find that **students or researchers are the most common participants** in the primary studies about effort estimation (91.67%). However, the presence of professionals is fundamental to produce realistic findings.







|   | 1. APIS & Microservices   | 2. Mainframe Legacy  | 3. Web Platforms  |
|---|---|--|---|
| CLIENT  |    |    |    |
|  ISBSG<br>Delivering IT Confidence |  |  |  |
|  AMMS                              |  |  |  |



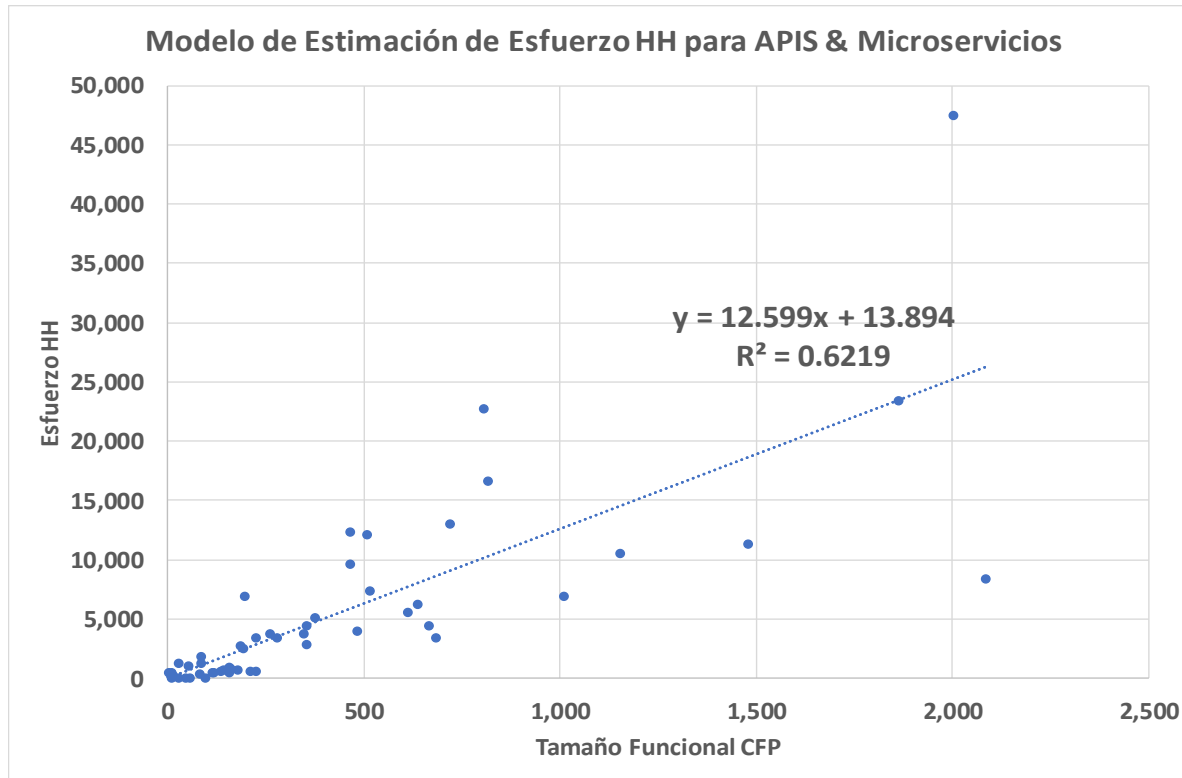
# Case Study 1. APIS & Microservices

| SOURCE       | Sample Size | %              |
|--------------|-------------|----------------|
| CLIENT       | 8           | 14.04%         |
| ISBSG        | 15          | 26.31%         |
| IMDS         | 34          | 59.65%         |
| <b>TOTAL</b> | <b>57</b>   | <b>100.00%</b> |

| SOURCE       | COSMIC Functional Size (CFP) | %           |
|--------------|------------------------------|-------------|
| CLIENT       | 2,418.7                      | 11.01%      |
| ISBSG        | 3,873.0                      | 17.63%      |
| IMDS         | 15,674.6                     | 71.36%      |
| <b>TOTAL</b> | <b>21,966.4</b>              | <b>100%</b> |

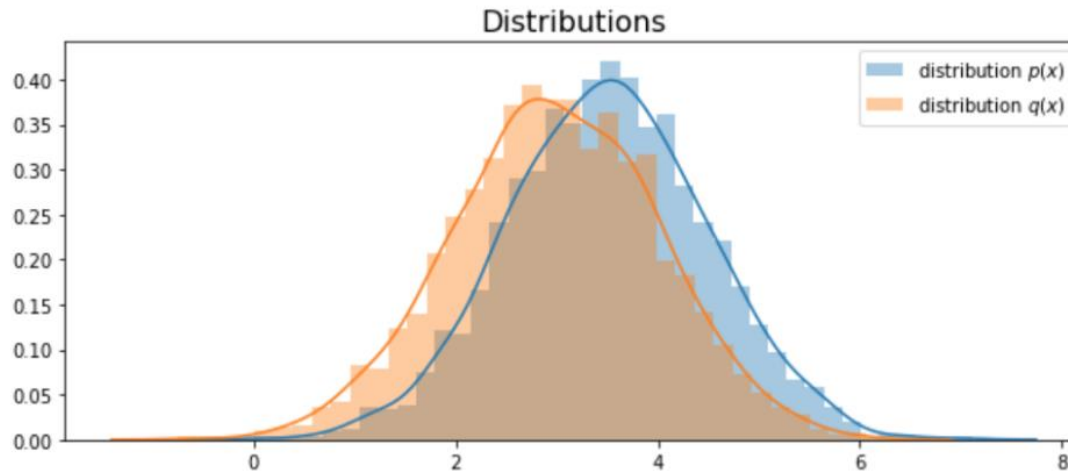
| Productivity  | PDR   |
|---|---|
| CFP/WH  | WH/CFP  |
| Productivity represents how many CFPs are implemented per work-hour | The PDR represents how many WH are required per CFP |

| SOURCE | Mín  | P10  | P25  | Media n | P75  | P90  | Máx   | Media | DesvE st |
|--------|------|------|------|---------|------|------|-------|-------|----------|
| CLIENT | 12.3 | ---- | 12.5 | 14.2    | 20.9 | ---- | 35.5  | 17.7  | 8.0      |
| ISBSG  | 0.2  | 0.3  | 2.3  | 3.9     | 4.7  | 12.9 | 23.7  | 4.6   | 5.6      |
| IMDS   | 1.4  | 4.5  | 8.0  | 11.7    | 20.5 | 33.1 | 143.2 | 18.1  | 24.1     |
| TOTAL  | 0.2  | 2.2  | 4.6  | 10.1    | 18.2 | 26.8 | 143.2 | 14.5  | 19.9     |



Could be integrated the three databases considering statistical foundations to get a high number of datapoints?

The integration make sense and it is valid?



- The **Kruskal-Wallis test**, also known as the H test, is the **non-parametric alternative to the one-way ANOVA test** for unpaired data.
- It is considered an **extension of the Mann-Whitney test for more than two groups**. It is therefore a test that uses ranges to contrast the hypothesis that  $k$  samples have been obtained from the same population.
- The Kruskal-Wallis test **contrasts whether the different samples are equidistributed and therefore belong to the same distribution (population)**. In a simple way, the Kruskal-Wallis test compares the medians.

The test was developed considering the Product Delivery Rate (PDR), using SPSS ver 25 in Spanish.

Hypothesis:

- $H_0: Med_1 = Med_2 = \dots = Med_k$
- $H_1: Med_i \neq Med_j$  for at least one pair  $(i, j)$

|   |              |
|---|--------------|
| <b>N</b>  | <b>57</b>    |
| <b>Grados de libertad (número de agrupaciones -1)</b> | <b>2</b>     |
| <b>Sig. Asintótica (p-value)</b>                      | <b>0.000</b> |

We can say that, since the **p-value** (Sig. Asymptot.) is **lower than 0.05**, then **the null hypothesis ( $H_0$ ) is rejected** and it is concluded that with a **significance level of 5%**, there is a significant difference in at least one of the PDR distributions of the databases

To determine **which databases have different distributions**, it is necessary to perform an analysis using **the Kruskal-Wallis test in pairs**, adjusting the resulting p-value considering the number of tests, **this correction is known as the Bonferroni correction** for various tests.

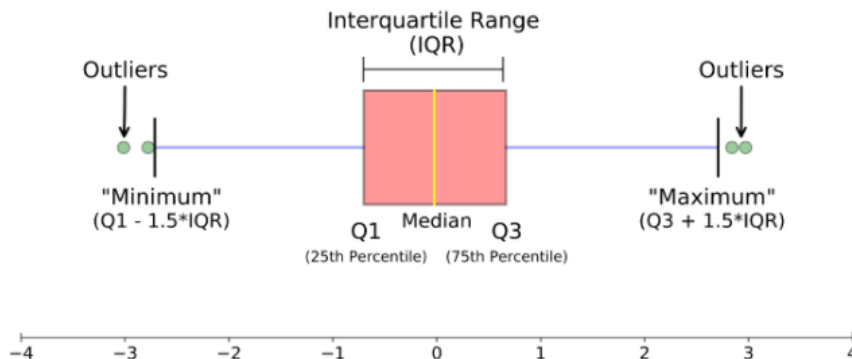
| Pareja         | Sig. Asintótica (p-value) | Sig. Asintótica (p-value) Ajustada [33] |
|----------------|---------------------------|---|
| ISBSG – IMDS   | 0.000                     | 0.000                                   |
| ISBSG – CLIENT | 0.000                     | 0.000                                   |
| IMDS – CLIENT  | <b>0.292</b>              | <b>0.876</b>                            |

We observe that **the IMDS - CLIENT pair** is the only one that has an adjusted p-value  $(0.876) > 0.05$ ; from which we conclude that the **IMDS and CLIENT databases have the same distribution and could be integrated**.



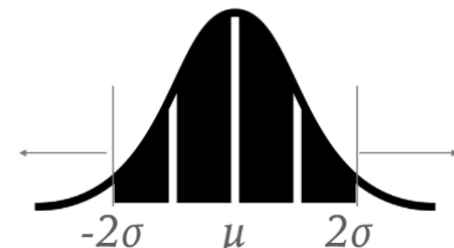
## Inter quartile distance:

The method most taught academically for its simplicity and results is the **Tukey test**, which takes as a reference **the difference between the first quartile Q1 and the third quartile Q3, or the interquartile range IQR (Q3-Q1).**

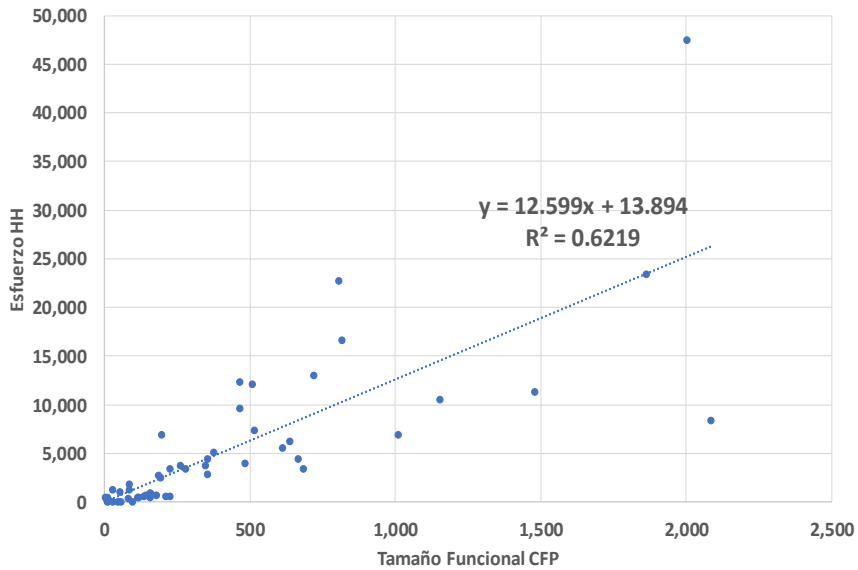


## Using the sample standard deviation:

Another method used is to take the sample mean as a reference and all those points that are outside the interval of two standard deviations around the mean can be considered atypical; This method allows defining if we want to use 1, 2 or 3 standard deviations to be more or less flexible.

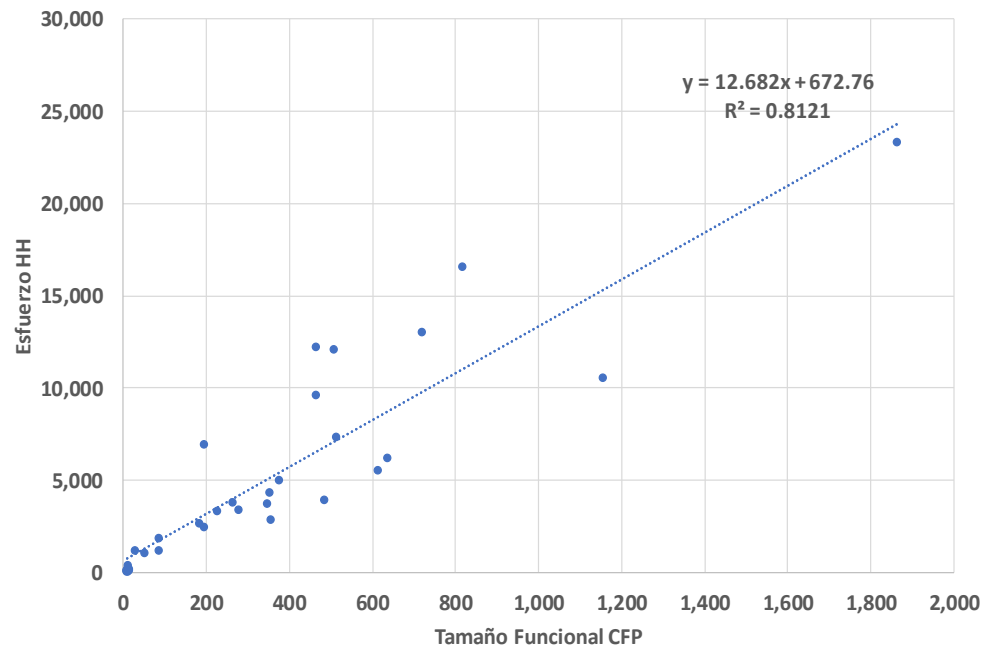


Modelo de Estimación de Esfuerzo HH para APIS & Microservicios



|                 |           |
|-----------------|-----------|
| <b>CLIENT</b>   | <b>8</b>  |
| <b>IMDS</b>     | <b>34</b> |
|                 | <b>42</b> |
| <b>Outliers</b> | <b>9</b>  |
| <b>Total</b>    | <b>33</b> |

Modelo de Estimación de Esfuerzo HH para APIS & Microservicios



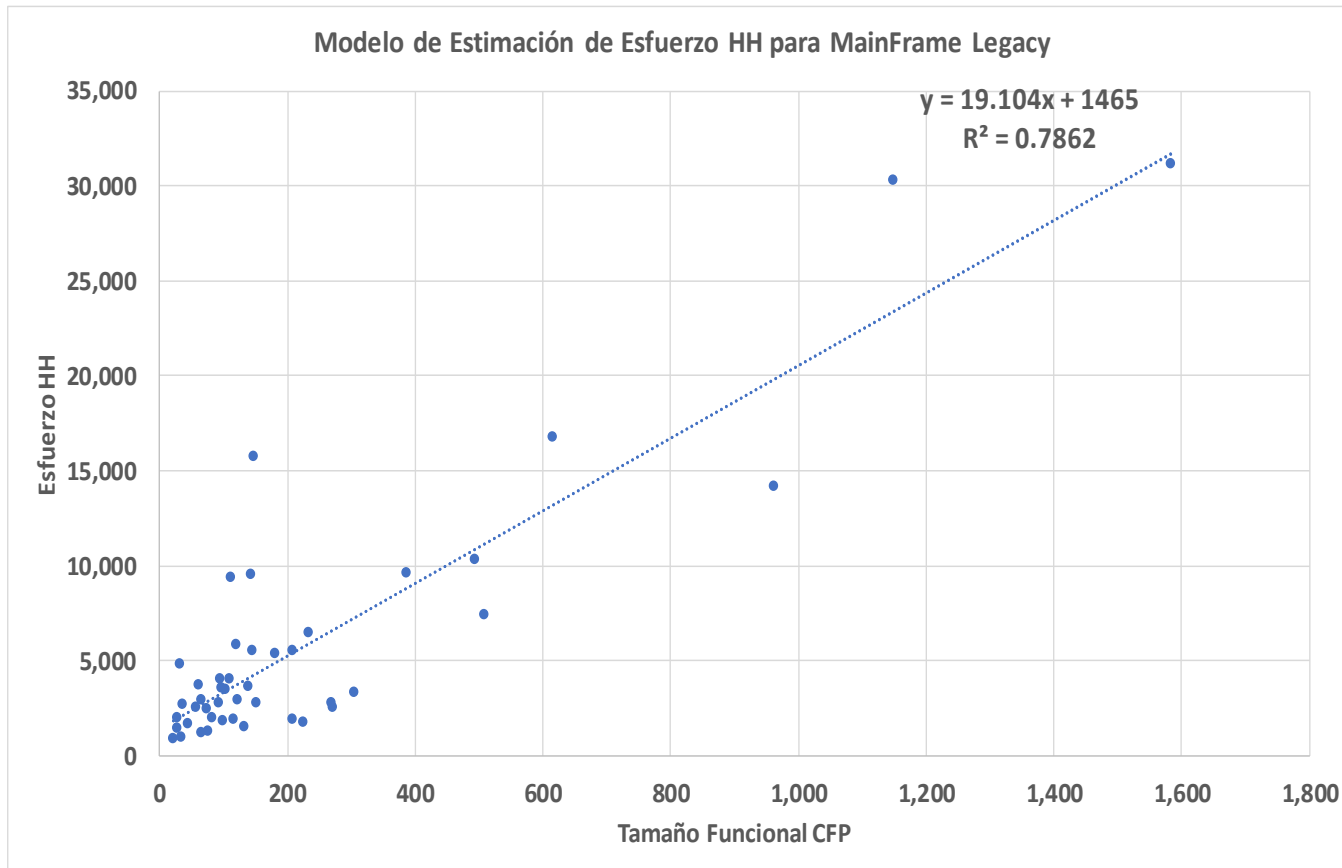
## Case Study 2. Mainframe Legacy

| SOURCE       | Sample Size | %             |
|--------------|-------------|---------------|
| CLIENT       | 10          | 22.2%         |
| ISBSG        | 35          | 77.8%         |
| IMDS         | 0           | 0.0%          |
| <b>TOTAL</b> | <b>45</b>   | <b>100.0%</b> |

| SOURCE       | COSMIC Functional Size (CFP) | %           |
|--------------|------------------------------|-------------|
| CLIENT       | 914.0                        | 8.98%       |
| ISBSG        | 9,267.0                      | 91.02%      |
| <b>TOTAL</b> | <b>10,181.0</b>              | <b>100%</b> |

| Productivity<br>CFP/WH  | PDR<br>WH/CFP                                       |
|---|---|
| Productivity represents how many CFPs are implemented per work-hour | The PDR represents how many WH are required per CFP |

| SOURCE | Mín  | P10  | P25  | Media<br>n | P75  | P90   | Máx   | Media | DesvE<br>st |
|--------|------|------|------|------------|------|-------|-------|-------|-------------|
| CLIENT | 10.5 | 11.3 | 19.4 | 38.4       | 53.3 | 145.3 | 152.9 | 47.5  | 41.4        |
| ISBSG  | 8.0  | 10.6 | 18.6 | 26.8       | 43.2 | 70.7  | 107.4 | 33.8  | 23.0        |
| TOTAL  | 8.0  | 10.9 | 18.9 | 28.2       | 44.5 | 76.7  | 152.9 | 36.9  | 28.2        |



Could be integrated the three databases considering statistical foundations to get a high number of datapoints?

The integration make sense and it is valid?

The test was developed considering the Product Delivery Rate (PDR), using SPSS ver 25 in Spanish.

Hypothesis:

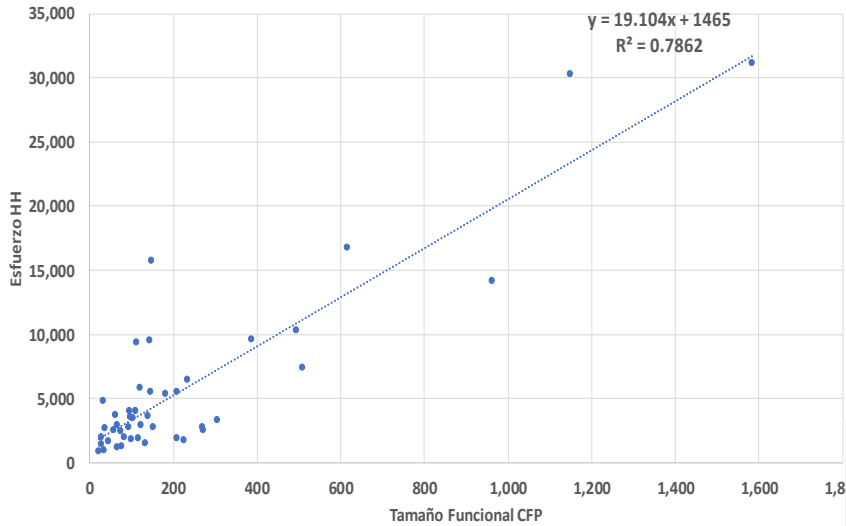
- $H_0: Med1 = Med2 = \dots = Medk$
- $H_1: Medi \neq Medj$  for at least one pair (i, j)

|   |              |
|---|--------------|
| <b>N</b>  | <b>45</b>    |
| <b>Grados de libertad<br/>(número de<br/>agrupaciones -1)</b> | <b>1</b>     |
| <b>Sig. Asintótica (p-<br/>value)</b>                         | <b>0.275</b> |

We can say that, since the **p-value** (Sig. Asymptot.) is **higher than 0.05**, then **the null hypothesis (H0) is accepted** and it is concluded that there is **NO** significant difference in the distributions of the CLIENT and ISBSG databases



Modelo de Estimación de Esfuerzo HH para MainFrame Legacy



|               |           |
|---------------|-----------|
| <b>CLIENT</b> | <b>10</b> |
|---------------|-----------|

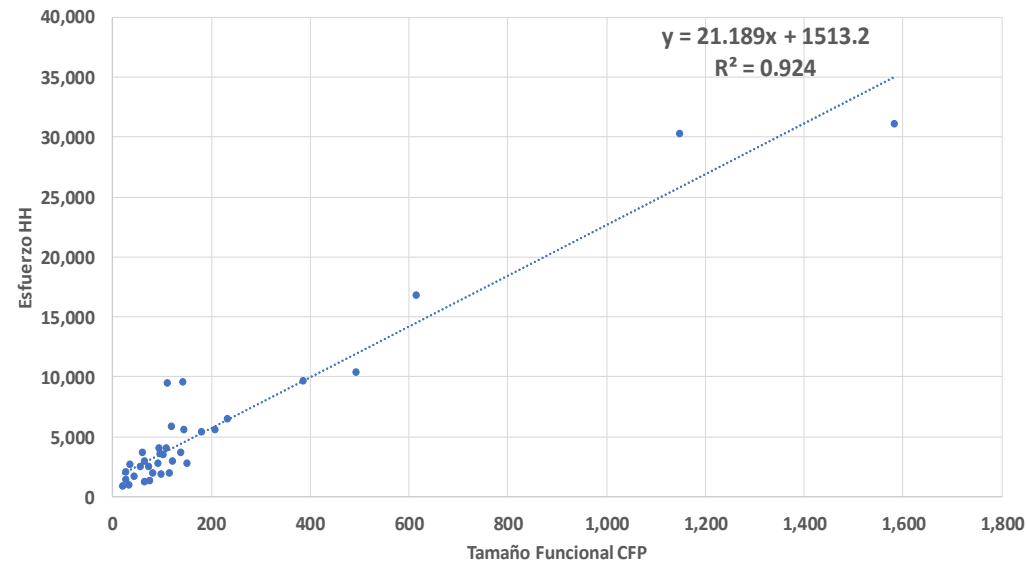
|              |           |
|--------------|-----------|
| <b>ISBSG</b> | <b>35</b> |
|--------------|-----------|

|  |           |
|--|-----------|
|  | <b>45</b> |
|--|-----------|

|                 |           |
|-----------------|-----------|
| <b>Outliers</b> | <b>10</b> |
|-----------------|-----------|

|              |           |
|--------------|-----------|
| <b>Total</b> | <b>35</b> |
|--------------|-----------|

Modelo de Estimación de Esfuerzo HH para MainFrame Legacy



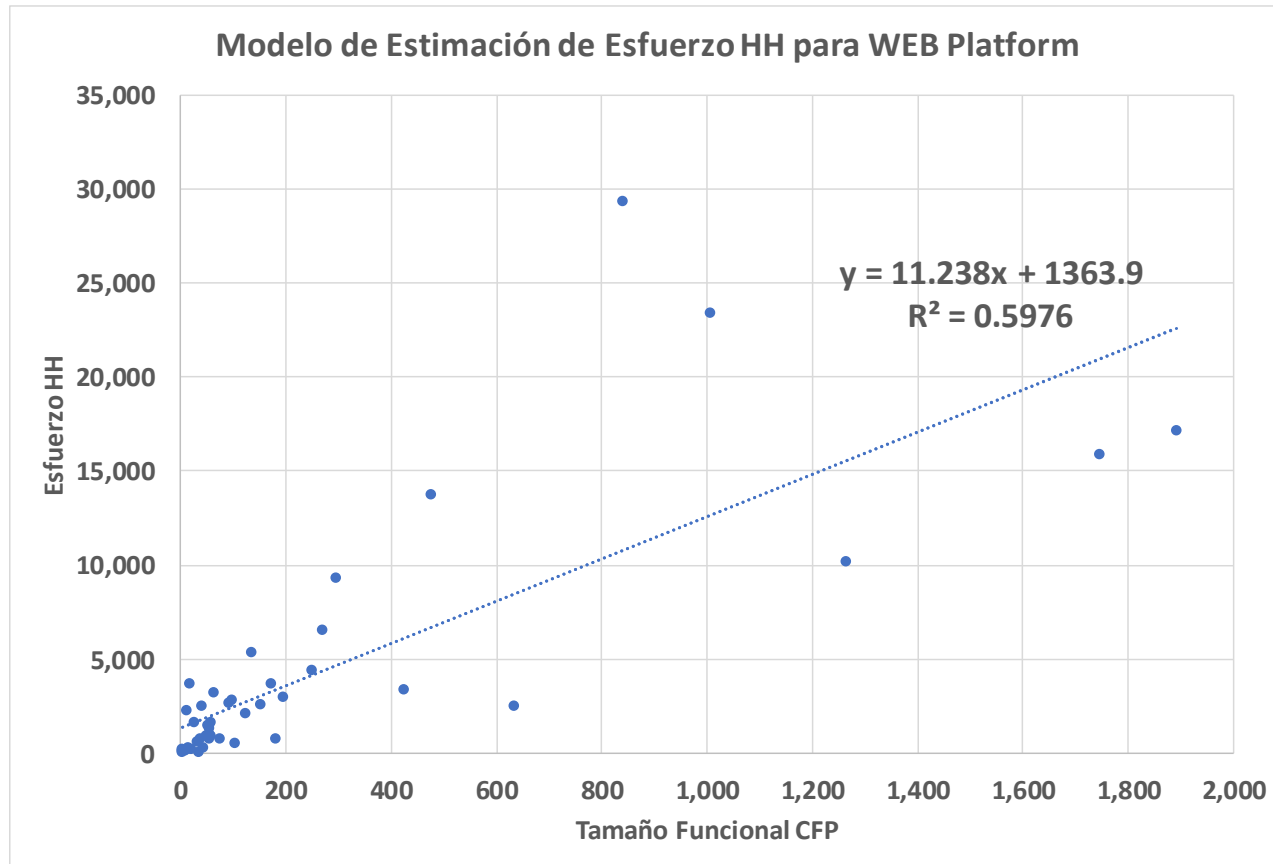
## Case Study 3. Web Platforms

| SOURCE       | Sample Size | %             |
|--------------|-------------|---------------|
| CLIENT       | 0           | 0.0%          |
| ISBSG        | 20          | 44.4%         |
| IMDS         | 25          | 55.6%         |
| <b>TOTAL</b> | <b>45</b>   | <b>100.0%</b> |

| SOURCE       | COSMIC Functional Size (CFP) | %           |
|--------------|------------------------------|-------------|
| CLIENT       | 0.0                          | 0.00%       |
| ISBSG        | 3,721.0                      | 33.31%      |
| IMDS         | 7,449.1                      | 66.69%      |
| <b>TOTAL</b> | <b>11,170.1</b>              | <b>100%</b> |

| Productivity<br>CFP/WH  | PDR<br>WH/CFP                                       |
|---|---|
| Productivity represents how many CFPs are implemented per work-hour | The PDR represents how many WH are required per CFP |

| SOURCE | Mín | P10  | P25  | Media<br>n | P75  | P90   | Máx   | Media | DesvE<br>st |
|--------|-----|------|------|------------|------|-------|-------|-------|-------------|
| ISBSG  | 4.0 | 13.9 | 18.3 | 24.8       | 48.5 | 217.2 | 257.8 | 50.4  | 68.9        |
| IMDS   | 2.8 | 5.1  | 8.6  | 19.2       | 29.1 | 46.8  | 131.2 | 24.2  | 25.9        |
| TOTAL  | 2.8 | 6.3  | 14.5 | 21.6       | 30.9 | 66.4  | 257.8 | 35.8  | 50.9        |



Could be integrated the three databases considering statistical foundations to get a high number of datapoints?

The integration make sense and it is valid?

The test was developed considering the Product Delivery Rate (PDR), using SPSS ver 25 in Spanish.

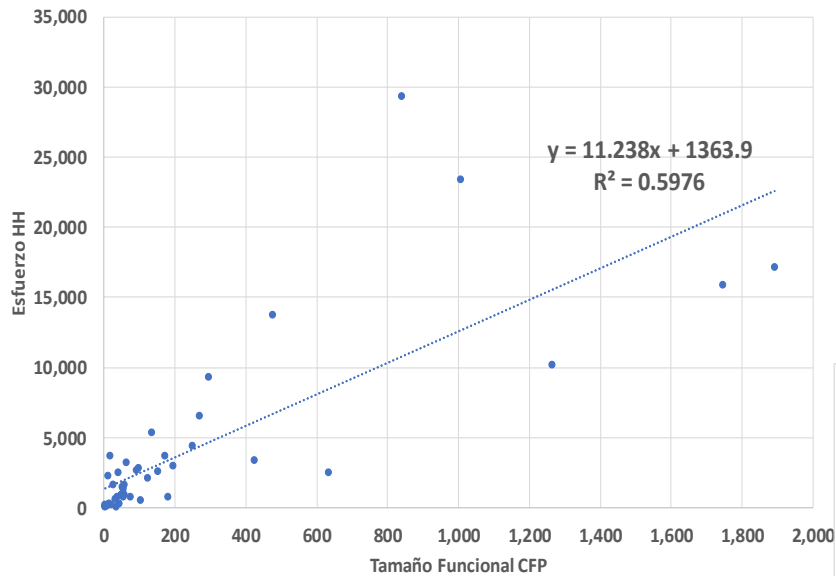
Hypothesis:

- $H_0: Med1 = Med2 = \dots = Medk$
- $H_1: Medi \neq Medj$  for at least one pair (i, j)

|   |              |
|---|--------------|
| <b>N</b>  | <b>45</b>    |
| <b>Grados de libertad (número de agrupaciones -1)</b> | <b>1</b>     |
| <b>Sig. Asintótica (p-value)</b>                      | <b>0.054</b> |

We can say that, since the **p-value** (Sig. Asymptot.) is **higher than 0.05**, then **the null hypothesis (H0) is accepted** and it is concluded that there is **NO** significant difference in the distributions of the IMDS and ISBSG databases

Modelo de Estimación de Esfuerzo HH para WEB Platform



|             |           |
|-------------|-----------|
| <b>IMDS</b> | <b>25</b> |
|-------------|-----------|

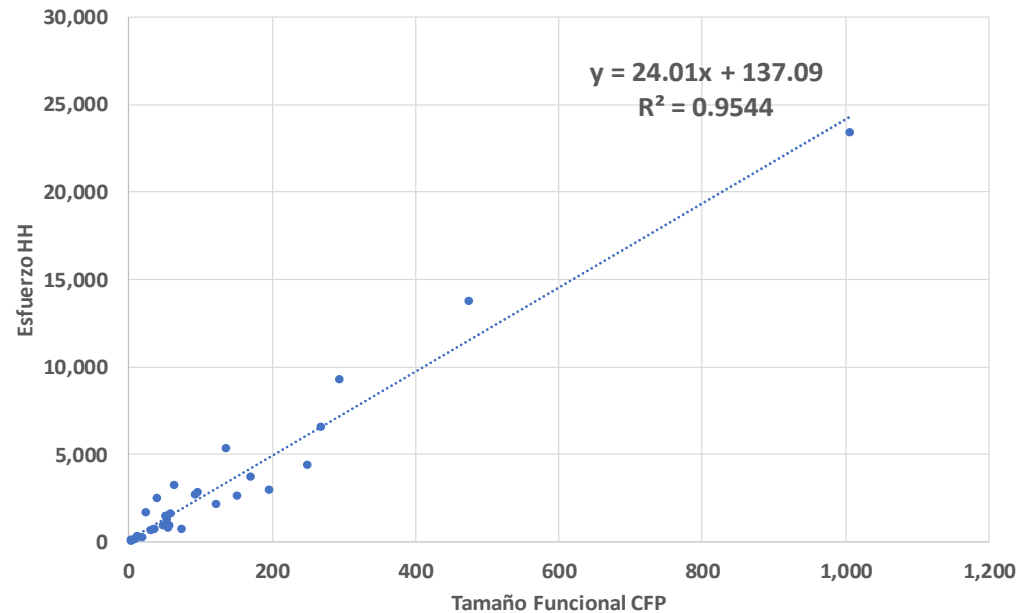
|              |           |
|--------------|-----------|
| <b>ISBSG</b> | <b>20</b> |
|--------------|-----------|

|  |           |
|--|-----------|
|  | <b>45</b> |
|--|-----------|

|                 |           |
|-----------------|-----------|
| <b>Outliers</b> | <b>13</b> |
|-----------------|-----------|

|              |           |
|--------------|-----------|
| <b>Total</b> | <b>32</b> |
|--------------|-----------|

Modelo de Estimación de Esfuerzo HH para WEB Platform





*The use of the proposed procedure has made it possible to improve the estimation models in the Mexican industry from the integration of different databases, considering statistical foundations to validate the integration of different data sources.*



# Questions?

**Dr. Francisco Valdés-Souto**

Associate Professor  
Department of Mathematics,  
Science Faculty,  
National Autonomous University of Mexico (UNAM)  
*[fvaldes@Ciencias.unam.mx](mailto:fvaldes@Ciencias.unam.mx)*



*COSMIC President*