# Improving Estimation Models by Merging Independent Data Sources

[1] *F. Valdés-Souto, ORCID: 0000-0001-6736-0666 <fvaldes@ciencias.unam.mx>*
[2] *J. Valeriano-Assem, ORCID: 0009-0008-6473-1271 <jorge.valeriano@spingere.com.mx>*
[1] *National Autonomous University of Mexico Science Faculty,*
*CDMX, México.*
[2] *SPINGERE, CDMX, México.*

**Abstract.** Software cost/effort estimation has been a key research topic for over six decades due to its industry impact. Despite numerous models, regression-based approaches dominate the literature. Challenges include insufficient datasets with enough data points and arbitrary integration of different source databases. This study proposes using the Kruskal-Wallis test to validate the integration of distinct source databases, aiming to avoid mixing unrelated data, increase data points, and enhance estimation models. A case study was conducted with data from an international company's Mexico office, which provides software development for "Microservices and APIs." Data from 2020 were analyzed. The estimation model's quality improved significantly. MMRE decreased by 25.4% (from 78.6% to 53.2%), standard deviation dropped by 97.2% (from 149.7% to 52.5%), and the Pred (25%) indicator rose by 3.2 percentage points. The number of data points increased, and linear regression constraints were met. The Kruskal-Wallis test effectively improved the estimation models by validating database integration.

**Ключевые слова:** linear regression model; software estimation; effort estimation; cost estimation; functional size; COSMIC method; Kruskal-Wallis.

**For citation:** Valdés-Souto F., Valeriano-Assem J. Improving estimation models by merging independent data sources. Trudy ISP RAN/Proc. ISP RAS, vol. 36, issue 6, 2024. pp. 7-18. DOI: 10.15514/ISPRAS-2024-36(6)- 1.

# Совершенствование моделей оценки путем объединения независимых источников данных

[1] *Ф. Вальдес-Соуто, ORCID: 0000-0001-6736-0666 <fvaldes@ciencias.unam.mx>*
[2] *Х. Валериано-Ассем, ORCID: 0009-0008-6473-1271 <jorge.valeriano@spingere.com.mx>*

[1] *Национальный автономный университет Мексики, факультет Науки,*
*Мехико, Мексика.*
[2] *SPINGERE, Мехико, Мексика.*

**Аннотация.** Оценка затрат/усилий на программное обеспечение является ключевой темой исследований более шести десятилетий из-за его влияния на отрасль. Несмотря на многочисленные модели, подходы, основанные на регрессии, доминируют в литературе. Проблемы включают в себя недостаточные наборы данных с достаточным количеством точек данных и произвольную интеграцию различных исходных баз данных. В этом исследовании предлагается использовать тест Крускала-Уоллиса для проверки интеграции отдельных исходных баз данных с целью избежать смешивания несвязанных данных, увеличения точек данных и улучшения моделей оценки. Было проведено тематическое исследование с данными из офиса международной компании в Мексике, который обеспечивает разработку программного обеспечения для «микросервисов и API». Были проанализированы данные за 2020 год. Качество модели оценки значительно улучшилось. MMRE снизился на 25,4% (с 78,6% до 53,2%), стандартное отклонение снизилось на 97,2% (с 149,7% до 52,5%), а показатель Pred (25%) вырос на 3,2 процентных пункта. Количество точек данных увеличилось, и были соблюдены ограничения линейной регрессии. Тест Крускала-Уоллиса эффективно улучшил модели оценки, подтвердив интеграцию базы данных.

**Ключевые слова:** модель линейной регрессии; оценка программного обеспечения; оценка затрат; оценка стоимости; функциональный размер; метод COSMIC; тест Крускала-Уоллиса.

## 1. Introduction

Software cost and effort estimation is crucial for planning, budgeting, and project success in software engineering. Regression-based estimation approaches are common, but the literature highlights challenges such as small datasets and the arbitrary combination of different source databases without proper validation. These issues are prevalent in both academia and industry, where small datasets are more common than expected. Estimation has potential for significant contributions to software engineering, particularly through leveraging statistical methods. This paper proposes a method for using established statistical techniques to validate the integration of distinct databases, thereby improving estimation models by increasing the number of data points. The paper outlines background information, the proposed procedure for data integration and validation, and the improvements observed in the estimation model, concluding with a discussion of the findings [1-4].

## 2. Background

### 2.1 Parametric Software Estimation

Software estimation, which began in the 1950s, has been crucial to the success of development projects, influencing budgeting and planning. Over more than 70 years, various techniques and classifications of estimation methods have been developed. However, many challenges and unanswered questions remain in software estimation research. A key factor in estimation accuracy is the measurement of software size [5]. Today, functional size is the only software feature that can be consistently quantified, emphasizing its importance. Every estimation model is closely tied to the method used to measure the input variables that produce the estimate.

### 2.1.1 Database Conformation for Parametric Estimation

When creating an estimation model, it's crucial to integrate a reference database from past projects. This database helps identify correlations between variables, with functional size being primary. However, regression-based models often face replication issues, as noted by various authors [1-4, 6-10]. These models are typically based on previous projects, but if a new project differs significantly from those in the database, the model's forecasts may be inaccurate. Challenges in accessing and interpreting data, as well as the limitations of available datasets, further complicate the process.

One major problem is the insufficient number of data points in datasets, which is critical for statistical reliability. Carbonera et al. classified datasets based on the number of points: high quality (more than 15), medium quality (10-15), and low quality (fewer than 10). The central limit theorem suggests that at least 30 data points are needed for each variable to approximate a normal distribution effectively.

Collecting 30 similar projects is often difficult, as noted by Morgenshtern et al., who highlight the cost and time involved in gathering historical data. Furthermore, combining data from different sources without proper evaluation can compromise its utility.

Organizations like the International Software Benchmarking Standards Group (ISBSG) and the Mexican Software Metrics Association (AMMS) maintain databases of past projects. The AMMS dataset, which includes data from real Mexican industry projects, shares similar features with the ISBSG dataset. Addressing the issue of limited data points may require developing techniques to integrate distinct databases using statistical methods, thereby improving the reliability of estimation models.

### 2.1.2 Estimation Models Performance Comparison

The performance of estimation models is assessed by applying quality criteria to evaluate their accuracy. Discrepancies between estimated and actual values are measured using criteria like Mean Magnitude of Relative Error (MMRE), Standard Deviation of MRE (SDMRE), Prediction level (PRED), Median Magnitude of Relative Error (MdMRE), and Mean Absolute Residual (MAR). Researchers have analyzed these techniques and identified various concerns regarding their effectiveness and reliability [10].

## 2.2 Kruskal-Wallis test

The Kruskal-Wallis [11, 12, 13] test is a nonparametric method used to compare the distributions of independent groups, serving as an alternative to one-way ANOVA when assumptions of normality and homogeneity are violated or when data are ordinal. Introduced by William Kruskal and Wilson Wallis in 1952, it ranks data from all groups and calculates a test statistic H based on these ranks. A higher H value suggests more evidence against the null hypothesis of no difference among distributions. H follows a chi-square distribution with $k-1$ degrees of freedom under the null hypothesis, where $k$ is the number of groups. The test assesses whether the sample rank distributions differ significantly, indicating differences in population medians. If H exceeds the critical chi-square value, it implies significant differences among groups, leading to rejection of the null hypothesis. It is widely used in experimental and observational studies.

## 2.3 Outliers identification using Tukey test

The Tukey [12, 13] test, developed by John Tukey, is used to identify outliers in a dataset. It involves calculating the interquartile range (IQR) by subtracting the first quartile (Q1) from the third quartile (Q3). A threshold, typically 1.5 times the IQR, is set to flag outliers. Observations falling below Q1 minus the threshold or above Q3 plus the threshold are considered potential outliers. The test is robust against moderate deviations from normality and is effective for skewed or non-normally distributed data.

## 3. Case study: integrating distinct sources databases

This section presents a summary of a case study carried out at an international company with a Mexico office, referred to as COMPANY for confidentiality purposes. The office offers software development services to a financial institution, with data gathered in 2020.

The case study comprises three steps, but this paper will focus only on the final two:

1. **Project Identification/Classification and Functional Size Approximation:** COMPANY carried out this step to determine the types of projects that required estimation. They selected projects from a technological tower labeled "Microservices and APIs". Using the integrated information, we applied the EPCU approximation method [14] to measure the FURs of each project using COSMIC (ISO/IEC 19761).

2. **Incorporation of Additional Projects from Other Sources:** Since the COMPANY's provided projects were insufficient to develop a reliable estimation model, we sought out similar projects related to Microservices or API development in the ISBSG and AMMS databases. A total of forty-nine (49) projects were identified: 15 from the ISBSG database and 34 from the AMMS database.

3. **Constructing the Final Estimation Model:** To develop the final estimation model, we utilized the Kruskal-Wallis test to compare the distributions of independent groups and assess the feasibility of integration. We then followed the steps outlined by Valdés-Souto et al. in [7-9] to build and refine the estimation model.

The COMPANY and alternative sources (ISBSG, AMMS) are crucial in our project characterization process. They provide the essential data needed to select projects with similar characteristics, allowing us to compare size and effort. The effort was measured using COSMIC (ISO/IEC 19761), which serves as the fundamental metric for our project characterization. In Table 1.a), column 1 lists the acronym of the source from which the projects were obtained, column 2 shows the number of projects in the sample, and column 3 indicates the proportion of each group relative to the total number of projects. All the projects involved microservices or API development. Table 1.b) presents in column 1 the source acronym, in column 2 the functional size in CFP per group, and in column 3 the proportion of size per group relative to the total functional size in the sample.

*Table 1. a) Sample size by source, "Microservices and APIS" projects. b) Total functional size by source.*

| SOURCE | Sample Size | % |
|---|---|---|
| COMPANY | 8 | 14.0% |
| ISBSG | 15 | 26.3% |
| AMMS | 34 | 59.7% |
| **Total** | **57** | **100.0%** |

a)

| SOURCE | COSMIC Functional Size (CFP) | % |
|---|---|---|
| COMPANY | 2418.7 | 11.0% |
| ISBSG | 3873 | 17.6% |
| AMMS | 15674.6 | 71.4% |
| **TOTAL** | **21966.3** | **100.0%** |

b)

Based on the tables above, it can be seen that the AMMS database is a major contributor, representing 71.4% of the total functional size and 59.7% of the total projects. The ISBSG database is the second largest contributor in terms of project quantity and size, accounting for 26.3% of the total projects and 17.4% of the total size. Data from the COMPANY had the smallest contribution in both size and quantity. Due to the central limit theorem, the number of projects from the

COMPANY is insufficient to create a significant estimation model using only the initially provided data from the COMPANY.

Originally, the estimation model that could be developed using only the COMPANY's data is shown in Fig. 1, where the x-axis represents CFP and the y-axis represents effort. Although this model achieved an $R^2$ of over 77%, the limited amount of data prevents meaningful extrapolation of the conclusions.

The model with three datasets is shown in Fig. 2. In this case the model presents a R2 of 62% that it is lower than 77% of the initial model. The main concern is whether the additional dataset has a different distribution or if its mean significantly varies from the mean of the previous data. If so, the impact could be significant, and it might not be a good idea to integrate the data.
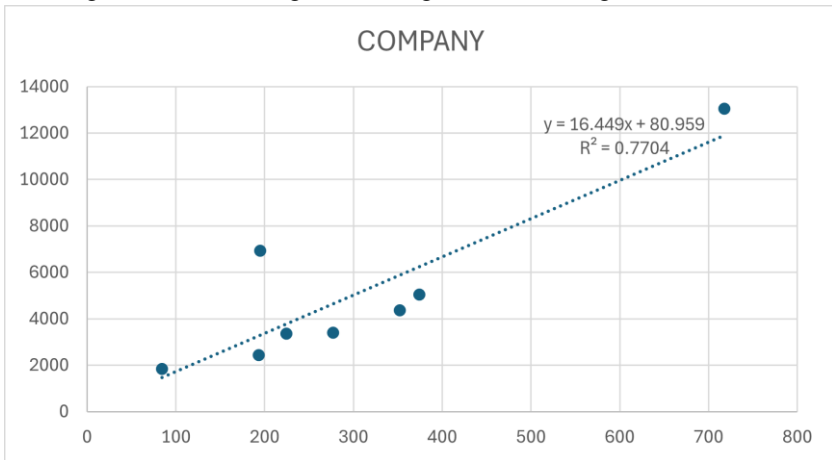


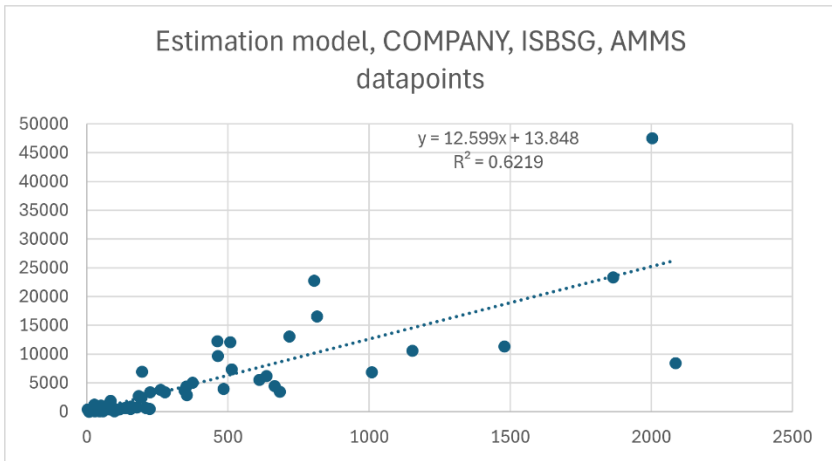*Fig. 1. COMPANY estimation model.*



*Fig. 2. COMPANY, ISBSG, AMMS estimation model.*

The next step was to compare the distributions of independent groups to evaluate whether the integration was feasible with solid statistical foundations. Specifically, we assessed whether the distributions of the three databases (COMPANY, ISBSG, AMMS) for the PDR variable (HH/CFP) are the same or different. This evaluation helps determine whether it is appropriate to combine the three databases into a single database and build estimation models. Since the project samples come from different databases, they are considered independent samples in statistics. In this case, there are three samples. To assess these, we used a nonparametric test called the Kruskal-Wallis test [11],

[13], which allows us to conclude whether the distributions of the three samples are equal or different. The null hypothesis (H0) is that there is no significant difference between the distributions of the COMPANY, ISBSG, and AMMS databases. The alternative hypothesis (H1) is that at least one of the distributions from the COMPANY, ISBSG, or AMMS databases is significantly different. The significance level required is α=0.05. If the test's p-value is greater than or equal to 0.05, then H0 is correct; otherwise, if it is less than 0.05, H1 is correct. The Kruskal-Wallis test was performed using SPSS® version 25 in Spanish.

Table 2 summarizes the results of the Kruskal-Wallis test for the COMPANY, ISBSG, and AMMS databases. The p-value is less than 0.05 (LINE 3); therefore, the null hypothesis (H0) is rejected. Consequently, we conclude that there is a significant difference in at least one of the distributions of the COMPANY, ISBSG, and AMMS databases, as stated by the alternative hypothesis (H1). To determine which databases have different distributions, it is necessary to perform pairwise comparisons using the Kruskal-Wallis test, adjusting the resulting p-values to account for the number of tests. This adjustment is known as the Bonferroni correction [13]. Table 3 displays the results for each pair of datasets analyzed. The AMMS – COMPANY pair is the only one with an adjusted p-value (0.6171) greater than 0.01667 (0.05/3). From this, we conclude that the AMMS and COMPANY databases have the same distribution, while the ISBSG database has a different distribution. Therefore, it is only possible to integrate the COMPANY and AMMS datasets to build the estimation model, resulting in a total of 42 datapoints (COMPANY (8), AMMS (34)).

*Table 2. Kruskal-Wallis test results for three datasets.*

| N | 57 |
|---|---|
| Degrees of freedom (Number of sets -1) | 2 |
| Asymp.sig. (p-value) | 0.00001696 |

*Table 3. Kruskal-Wallis test results by couple of datasets.*

| Pair | Asymp.sig. (p-value) | Asymp.sig. (p-value) with Bonferroni correction |
|---|---|---|
| ISBSG – AMMS | 0.00001578 | 0.00004735 |
| ISBSG – COMPANY | 0.0006197 | 0.001859 |
| AMMS – COMPANY | 0.2057 | 0.6171 |

Once the integration validation is performed, we have the final dataset (COMPANY + AMMS) to develop an estimation model directly. The results are shown in Fig. 3. The generated estimation model is **y = 8.6672x + 1586.9**, with a determination coefficient **$R^2$ = 0.5388**. However, it is necessary to develop a linear regression model validation and diagnostics. The Normal probability graph and the Residuals graph were obtained using an Excel add-in to analyze the regression model, as shown in Fig. 4.

Fig. 4 shows evidence against normality, as the points do not follow the identity line in the normal probability graph. Additionally, in the residuals graph, the variance of the residuals increases with the fitted values, showing a systematic pattern and indicating non-constant variance, which means the data do not exhibit homoscedasticity. To correct the model's assumptions, we applied a logarithmic transformation to the functional size and effort variables and built a new estimation model. Refer to Fig. 5, where the x-axis represents Log(CFP) and the y-axis represents Log(effort). The new estimation model is Log(y) = 0.9326 Log(x) + 2.8916, with a coefficient of determination $R^2$ = 0.8339. We conducted validation and diagnostics using the transformed data, with results

12

shown in Fig. 6. The plot of fitted values against residuals indicates constant variance, as the dots do not display patterns, demonstrating homoscedasticity. Consequently, the estimation model in Fig. 5 meets the statistical principles of normality and homoscedasticity, making the linear regression model appropriate for this dataset.
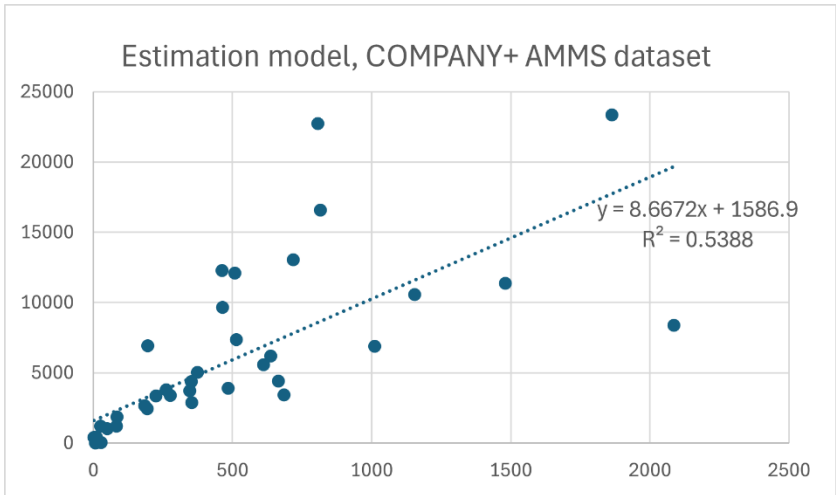


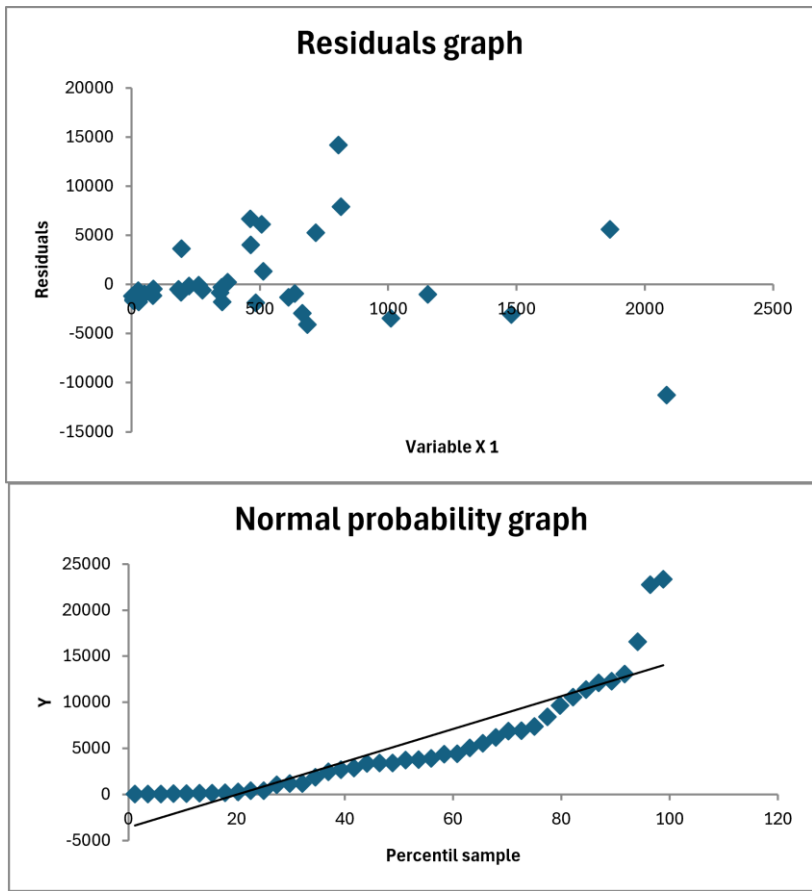*Fig. 3. Initial estimation model AMMS-COMPANY dataset.*



*Fig. 4. Graph for validation and diagnostics AMMS-COMPANY dataset.*
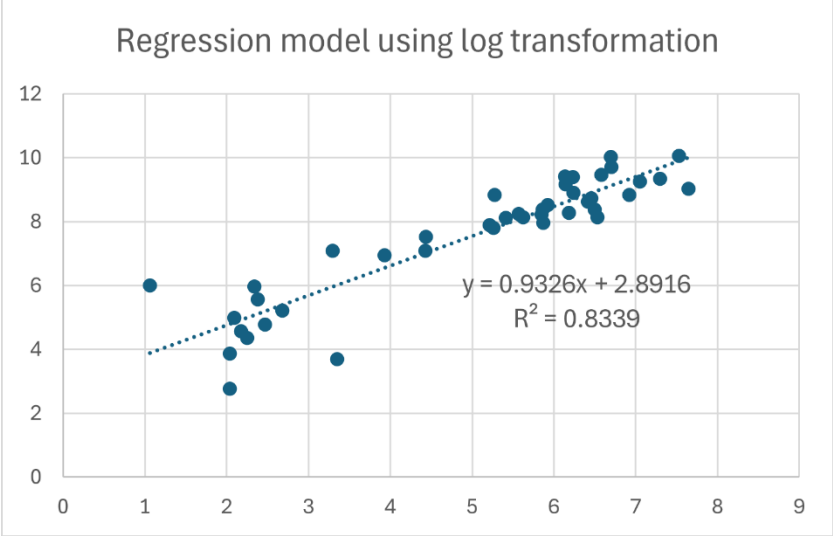
*Fig. 5. Estimation model AMMS-COMPANY dataset using logarithm transformation.*
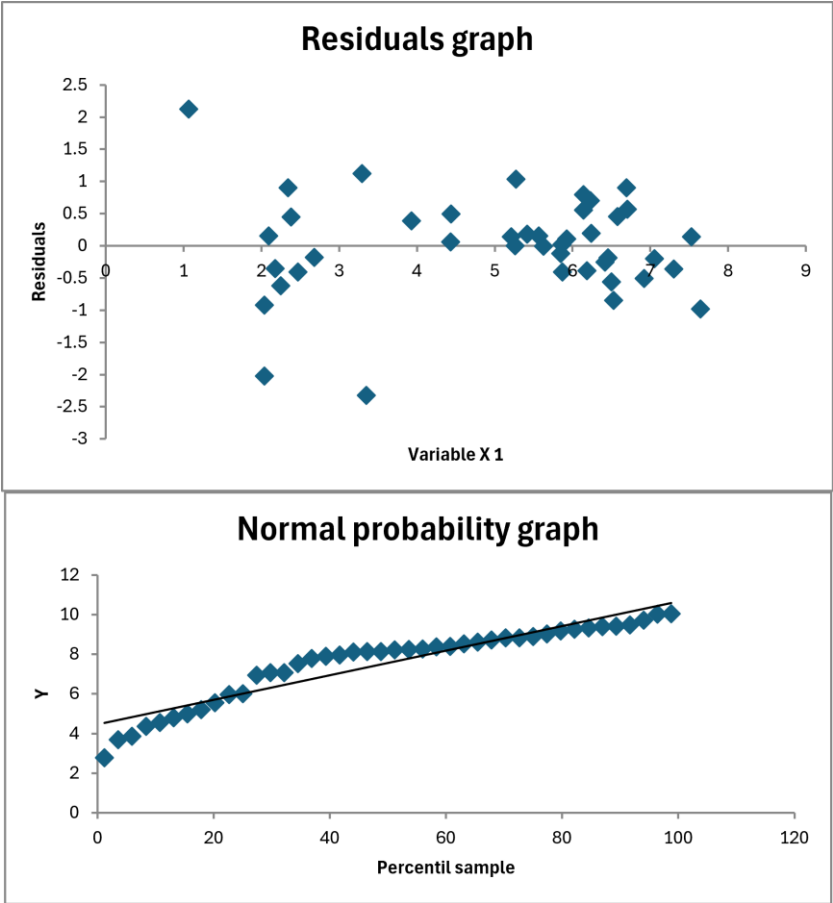


*Fig. 6. Graph for validation and diagnostics AMMS-COMPANY dataset with logarithmic transformation.*

After that, we search for outliers using the Tukey test, finding four (4) outliers as shown in Fig. 7. After removing the outliers, a new estimation model was obtained, the result is:

$$\text{Log(y)} = 0.9377 \, \text{Log(x)} + 2.8996$$

with a Determination coefficient $R^2 = 0.9023$. The model uses logarithmic variables. To apply it to the actual variables, we need to eliminate the logarithmic transformation using the inverse operation (Euler's number, $e$), resulting in the final model: $\mathbf{y = x^{0.\,9377} * e^{2.8996}}$.

Table 4 presents the quality criteria for the developed estimation models. The best model is the last one, achieved after applying validations and diagnostics, then performing a transformation and removing the outliers.
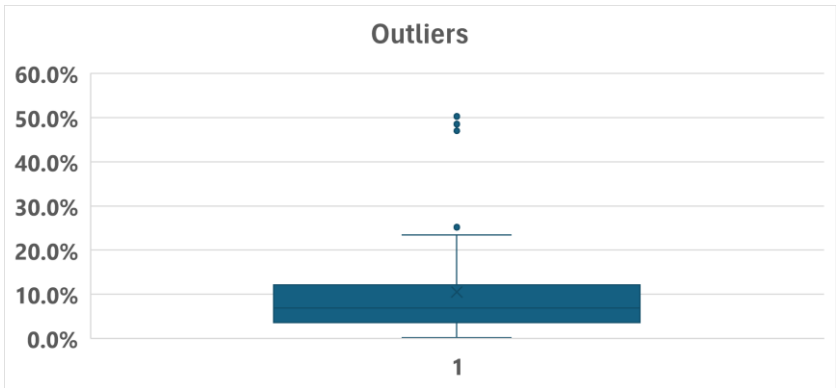


*Fig. 7. Outliers AMMS-COMPANY dataset with logarithmic transformation.*

*Table 4. Quality criteria comparison for estimation models.*

|  | y = 16.449x + 80.959<br>R² = 0.7704 | y = 8.6672x + 1586.9<br>R² = 0.5388 | y= x^{0.9422}* e ^{3.0156} |
|---|---|---|---|
| N | 8 | 42 | 38 |
| MMRE | 27.8% | 78.6% | 53.2% |
| SDMRE | 13.4% | 149.7% | 52.5% |
| Pred (25%) | 50.0% | 31.0% | 34.2% |
| Enough data | NO | YES | YES |

## 4. Analysis

In the case study presented, the COMPANY under study had only eight (8) data points. Two additional datasets were considered: the ISBSG dataset with fifteen (15) data points and the AMMS dataset with thirty-four (34) data points.

However, the ISBSG data points were rejected by the Kruskal-Wallis test, resulting in a final dataset with forty-two (42) data points. After removing outliers, the dataset contained thirty-eight (38) data points.

The results obtained are related to MMRE that was reduced by 25.4% (from 78.6% to 53.2%), the standard deviation was reduced by 97.2% (from 149.7% to 52.5%), and the Pred (25%) indicator increased by 3.2 percentage points. Notably, the number of data points was significantly increased, from 8 to 38 (475%), bolstering the robustness of our findings.

## 5. Conclusions

Software cost/effort estimation has been a key research issue for over 60 years, with regression-based methods being widely used. However, issues have arisen related to dataset conformity,

including insufficient data points and arbitrary combining of different sources. This study presents a real case applying statistical methods, specifically the Kruskal-Wallis test, to determine if data from different sources can be integrated without compromising dataset integrity. The integration, validated through this analysis, allows for a larger and more significant dataset, improving the estimation models.

The case study demonstrated that an estimation model generated from validated integrated datasets outperformed one created from unvalidated sources. This underscores the importance of validation and diagnostic analysis in integration efforts. The study aimed to establish a formal methodology for creating reliable estimation models from diverse data sources, addressing a common issue in both industry and academia. While the proposed approach showed promise, it is crucial to apply statistical principles correctly; otherwise, the models might be ineffective. The methodology, developed and applied in the study, represents a significant advance in addressing the problem of small dataset sizes in software estimation.

# References

[1]. M. Jørgensen and M. Shepperd, "A systematic review of software development cost estimation studies," IEEE Trans. Softw. Eng., vol. 33, no. 1, pp. 33–53, 2007, doi: 10.1109/TSE.2007.256943.

[2]. P. L. Braga, A. L. I. Oliveira, and S. R. L. Meira, "Software Effort Estimation using Machine Learning Techniques with Robust Confidence Intervals," in 7th International Conference on Hybrid Intelligent Systems, 2007, no. October 2007. doi: 10.1109/his.2007.56.

[3]. C. E. Carbonera, K. Farias, and V. Bischoff, "Software development effort estimation: A systematic mapping study," IET Res. Journals, vol. 14, no. 4, pp. 1–14, 2020, doi: 10.1049/iet-sen.2018.5334.

[4]. A. Abran, Software Project Estimation: The Fundamentals for Providing High Quality Information to Decision Makers, 1st ed. John Wiley & Sons, 2015.

[5]. A. Abran, Software Metrics and Software Metrology. Hoboken, New Jersey: John Wiley & Sons, 2010.

[6]. N. Kinoshita, A. Monden, M. Tshunoda, and Z. Yucel, "Predictability classification for software effort estimation," in Proceedings - 2018 IEEE/ACIS 3rd International Conference on Big Data, Cloud Computing, Data Science and Engineering, BCD 2018, 2018, no. 1, pp. 43–48. doi: 10.1109/BCD2018.2018.00015.

[7]. F. Valdés-Souto, "Validation of supplier estimates using cosmic method," CEURInternational Work. Softw. Meas. Int. Conf. Softw. Process Prod. Meas. (IWSM Mensura 2019), vol. 2476, pp. 15–30, 2019.

[8]. F. Valdés-Souto and L. Naranjo-Albarrán, "Improving the Software Estimation Models Based on Functional Size through Validation of the Assumptions behind the Linear Regression and the Use of the Confidence Intervals When the Reference Database Presents a Wedge-Shape Form," Program. Comput. Softw., vol. 47, no. 8, pp. 673–693, 2021, doi: 10.1134/S0361768821080259.

[9]. F. Valdés-Souto, "Creating an Estimation Model from Functional Size Approximation Using the EPCU Approximation Approach for COSMIC (ISO 19761)," in Software Engineering: Methods, Modeling and Teaching, Volume 4, Editorial., C. Mario, Z. Jaramillo, C. Elena, D. Vanegas, and W. P. Charry, Eds. Bogotá, Colombia, 2017, p. 468.

[10]. L. Lavazza, "Accuracy Evaluation of Model-based COSMIC Functional Size Estimation," in ICSEA 2017: The Twelfth International Conference on Software Engineering Advances, 2017, no. c, pp. 67–72.

[11]. W. A. Kruskal, W. H., & Wallis, "Use of Ranks in One-Criterion Variance Analysis," J. Am. Stat. Assoc., vol. 47, no. 260, pp. 583–621, 1952, doi: https://doi.org/10.1080/01621459.1952.10483441.

[12]. J. W. T. W. Tukey, Exploratory Data Analysis, 1st ed. Addison & Wesley, 1977.

[13]. R. R. Wilcox, Introduction to Robust Estimation and Hypothesis Testing, 4th Editio. Academic Press, 2016. [Online]. Available: https://shop.elsevier.com/books/introduction-to-robust-estimation-and-hypothesis-testing/wilcox/978-0-12-804733-0

[14]. A. Abran et al., "Early Software Sizing with COSMIC: Experts Guide," vol. 2020, no. May. Common Software Measurement International Consortium (COSMIC), pp. 1–67, 2020. doi: 10.13140/RG.2.1.4195.0567.

## Информация об авторах / Information about authors

Франсиско ВАЛЬДЕС-СОУТО имеет степень PhD в области инженерии программного обеспечения по специальности "Измерение и оценка программного обеспечения" в Высшей технологической школе (ETS) в Канаде, две магистерские степени в Мексике и Франции. Президент COSMIC. Доцент факультета наук Национального автономного университета Мексики. Основатель Мексиканской ассоциации метрик программного обеспечения (AMMS). Имеет более 25 лет опыта в разработке критически важного программного обеспечения. К настоящему времени опубликовал более 50 научных работ, включая статьи в индексированных журналах, трудах научных конференций, книгах и главах книг. Является главным промоутером проекта изучения формальных метрик программного обеспечения в Мексике, продвигая COSMIC (ISO/IEC 19761) в качестве национального стандарта. Член Национальной системы исследователей (SNI). Область научных интересов: измерение и оценка программного обеспечения, применяемого для управления проектами программного обеспечения, управление тематикой, производительностью и экономикой разработок программного обеспечения.

Francisco VALDÉS-SOUTO has a PhD degree in Software Engineering with a specialty in Software Measurement and Estimation at the École de Technologie Supérieure (ETS) in Canada, two master's degrees in Mexico and France. President of COSMIC. Associate Professor of the Faculty of Sciences of the National Autonomous University of Mexico (UNAM). Founder of the Mexican Association of Software Metrics (AMMS). More than 25 years of experience in critical software development. He currently has more than 50 publications including articles in Indexed Journals, Proceedings, books and book chapters. He is the main promoter of the topic of formal software metrics in Mexico, promoting COSMIC (ISO/IEC 19761) as a National Standard. Member of the National System of Researchers (SNI). Research interests: software measurement and estimation applied to software project management, scope management, productivity and economics in software projects.

Хорхе ВАЛЕРИАНО-АССЕМ – Магистр компьютерных наук и инженерии в Национальном автономном университете Мексики, специалист-консультант по формальному измерению и оценке программного обеспечения с 2016 года. Сфера научных интересов: метрики программного обеспечения (COSMIC), модели оценки программного обеспечения, модели валидации программного обеспечения, оценка функциональных и нефункциональных требований, оценка эффективности проектов разработки программного обеспечения на основе метрик программного обеспечения, оценка качества программных продуктов.

Jorge VALERIANO-ASSEM – master in Computer Science and Engineering from the National Autonomous University of Mexico, specialist consultant in formal software measurement and estimation since 2016. Areas of interest: Software metrics (COSMIC), Software estimation models, Software Validation Models, Estimation of Functional and Non-Functional Requirements, Evaluation of the Performance of Software Development Projects aligned to Software Metrics, Evaluation of the Quality of the Software Development Product.