

Notes on Turkish Software Industry: Developer Participation and Effort Estimation

Ekrem Kocaguneli
Lane Department of Computer Science and
Electrical Engineering
West Virginia University
Morgantown, WV 26505, USA
ekocagun@mix.wvu.edu

Bora Caglayan, Ayse Tosun
Ayse Bener
Computer Engineering Department
Bogazici University
34342 Bebek, Istanbul, TURKEY
bora.caglayan@boun.edu.tr,
ayse.tosun@boun.edu.tr,
bener@boun.edu.tr

ABSTRACT

Software effort estimation is used for the purpose of resource allocation and planning. Accurate estimates enable managers to accurately assess their available resources and distribute the load among programmers in a balanced manner. To assign the tasks and work-hours to developers, organizations use different project management tools. However, the actual participation of developers may be different from the values observed through project management tools.

In this research, we provide an overview of effort estimation activities of software companies in Turkey that specialize in various domains. For each company we analyze the developer participation through churn data of a representative project.

As a result of our analysis we observe that Turkish software industry is mature enough to appreciate the importance of software effort estimation, yet using ad-hoc methods. Furthermore, we observe that developer participation in different projects share common characteristics: More than 80% of the whole edits in code are performed by a very limited number of developers.

1. INTRODUCTION

Being able to estimate the development effort of a software project accurately is an important indicator of mature software development processes. Being able to plan and accurately manage the activities and resources associated with software development is also at least as important as accurate estimation. Software effort estimation as well as resource planning and monitoring has been addressed in various contexts. A considerable number of the studies have used publicly available datasets in their experimental settings. Results elicited from publicly available datasets have

provided a global picture of two issues: Software effort estimation and resource planning.

In our research, we tackle the two problems in a local context and provide a local picture of Turkish software industry. In terms of effort estimation we provide details about the estimation models that are used by companies in Turkish software industry. As for the resource planning, we focus on programmer related activities and evaluate programmer participation. Our research includes a wide range of companies from various domains: Banking, telecommunication as well as white good companies. When evaluating our results, we compare them to findings coming from open source projects. Our purpose in that is to see whether there are any similarities between development activities in open source and company settings.

As a result of our research we have seen that the importance of software effort estimation as well as resource planning is well understood in Turkish software industry. Companies have their own ad-hoc models of dealing with effort estimation as well as resource planning. For effort estimation companies use linear regression models, a form of expert judgment or a combination of the two. The programmer activities and workload distribution is monitored and managed through third party or in-house developed project management tools.

Although each problem is tackled by companies, there are some discomfort regarding both issues. The regression models used for software effort estimation are built on expert judgment. Estimation include features identified by experts and weights assigned to them also by experts. However, as the characteristics of projects change over time and more factors become influential on the software effort, human judgment becomes less effective and the linear models based on human expert judgment become obsolete. This fact also been pointed out by Jorgensen et. al. as the poor capability of humans to improve their expert judgment. [6]. Another discomfort regarding the monitoring developer participation through project management tools is that most of the time everyone completes the assigned work hours through the tool whereas some developers are believed to produce more than the others. The fact that some developers producing more than the others is not a new concept. Koch has conducted an

extensive case study in which he has questioned the manpower distribution on open source projects and has come up with the conclusion that commercial projects and open source projects share common characteristics [9]. In both commercial and open source projects, the manpower distribution and corresponding production per developer in terms of static code attributes like lines of code (LOC) show different characteristics. In Turkish software industry, we have seen similar results among the projects we have analyzed. Although workload is evenly distributed to developers, more than 80% of all source code development is performed by less than 20% of the developers.

Our datasets are composed of various projects coming from multiple software companies from different domains in Turkey. For our analysis we make use of churn analysis as well as static code attributes. Since we collect our data from private companies, the proprietary rights is a big concern. We analyze the churn data of all the projects in our research and give the necessary information regarding our analysis. However, project names, company names as well as the developer names will be kept anonymous due to proprietary rights.

The rest of the paper is organized as follows: In Section 2 we provide background information regarding software effort estimation modeling as well as case studies focusing on programmer participation, in Section 3 we give the details of our datasets. Then we continue with Section 4 in which we provide the methodology we adopted in this research. Section 5 provides the results we elicited from our analysis. In Section 6 we identify the possible threats to the validity of our research. Finally we conclude our research in Section 7.

2. BACKGROUND

In this research we present a ~~screenshot~~ of Turkish software industry in two parts: 1) ~~Adopted methods for software effort estimation~~ and 2) ~~in terms of developer participation in development activity~~. The first part is about our observations that we elicited during various projects with the industry. The second part is closer to a case study of developer participation in Turkish software industry and our conclusions will be based on churn analysis of various different projects coming from different software companies in Turkey. Therefore, in this section we provide some background information for both parts. In Section 2.1 we present some background information regarding software effort estimation and in Section 2.2 we present information regarding case studies.

2.1 Software Effort Estimation

We can categorize software effort estimation into two groups [14]: Expert judgment-based techniques and model-based techniques.

A very widely used technique in software effort estimation domain is expert judgment [5]. The employment of expert judgment methods may be either through well defined methods like Delphi [2] or via ad-hoc methods like formal meetings between experts in an organization. Although expert judgment methods are highly regarded in various settings, they have their own pitfalls. To begin with, the application

of expert judgment methods may create an atmosphere of competing interests. The competing interests may be between different departments of an organization or between experts with different levels of expertise. In the former case, the urgent needs of a particular department may result earlier deadlines for the project, which would in return affect the allocation of many resources and which eventually may challenge or even fail the project. In the latter case, estimations of a senior expert may be more dominant to those of a junior expert regardless of their accuracy. Secondly, human experts are evaluated as poor in terms of improving their estimation skills [6]. Lastly, an organization willing to completely rely on expert judgment methods for software effort estimation shall keep in mind that lack of good experts will result in imprecise estimations. Furthermore bad estimations would affect critical decisions regarding whether or not to start new projects or when to cancel challenged projects.

Unlike expert judgment-based techniques, model-based methods rely on parametric or algorithmic models. Parametric methods adapt an expert-proposed model to local data. Therefore, for the use of parametric methods collection of local data is a must. One of the most widely known example to parametric methods is the COCOMO method proposed by Boehm [3]. Algorithmic methods are useful when local data does not conform to the specifications of a parametric method. Algorithmic methods employ various algorithms to evaluate local data, build a model and make an estimation. Some well known examples to algorithmic methods are linear regression, case based reasoning systems, neural nets and model trees [8, 10, 12, 15].

The latter approach is useful in the case where local data does not conform to the specifications of the expert's method. A few examples of induced prediction systems are linear regression, neural nets, model trees and analogies [12, 15]. These methods have been mostly aided with certain types of patches such as noise removal, instance selection, feature selection or by feature reduction. Depending on the fundamental assumptions of each algorithm, a different patch can be selected.

Certain organizations may choose to use an expert judgment and some other organizations may prefer to go with a model-based approach, or alternatively a combination of the two in different settings. Regardless of the adopted method, the goal of any estimation model is to attain high estimation accuracy values.

2.2 Type of Research

According to two-dimension classification scheme proposed by Basili et. al. [1], our research falls into the class of blocked subject-project studies, which examine objects across a set of teams and a set of projects. To increase the granularity of this classification we can say that our research is a survey of effort estimation practices and developer participation in Turkish software industry. Any study can be a formal experiment, a case study or a survey depending how many teams and how many projects are considered in this study and what the scale of the study is [7]. Due to the nature of formal experiments, they are limited in size and need careful control. Case studies focus on what happens on a typical project and surveys on the other hand focus on what

is happening over various types of projects. In our research we deal with different companies and different projects that were developed with different languages. Furthermore, we are more of an observer in our research rather than experimenter. This is due to the fact that rather than being able to experiment the response of projects as well as project teams in response to a changing factor, we merely observed the effort estimation practices and developer participations across multiple projects and multiple companies. From our reading of Kitchenham et. al. [7] our research conforms to the definition of a survey.

Different surveys regarding both effort estimation practices and developer participations have been conducted in different settings. For example Jorgensen et. al. addresses the effort estimation practices and budget overruns in Norwegian software industry [13]. Koch questions effort modeling as well as developer participation over open source software projects in his extensive survey [9]. Although similar surveys have been conducted by different researchers in different contexts previously, to the best of our knowledge there is no previous survey evaluating the effort estimation practices and developer participation in the context of Turkish software industry.

3. DATA

In our research we analyze 6 different projects coming from different companies. Both in terms of their development languages and in terms of their magnitude they are very different from one another. Ideally the suggested and practiced strategy for selecting projects in surveys is random sampling [7, 13]. However, private organizations are not completely optimistic about sharing churn data and static code attribute information of all their projects. Therefore, in our research we have asked them to provide us with representative projects of their software development activities. This may have introduced a selection bias into our study and we will address this issue as threat to validity more broadly in Section 6.

The project names and their particular characteristics are provided in Figure 1. Due to proprietary rights, we will name the projects and developers with numbers. The properties of projects given Figure 1 are number of developers employed, total edited lines of code (*total edited LOC* = *added LOC* + *deleted LOC*), total number of commits that were performed during the lifecycle of the project and the language in which the project was developed. We see from Figure 1 that our sample set contains a wide variety of projects in terms of the features given in the same Figure. The smallest project in terms of size (Project1) employs only 3 developers and a total of 5579 LOC were edited for this project. The biggest project in terms of size (Project3) on the other hand employs 98 developers and includes a total of 1324956 edited LOC. Therefore, although we were unable to apply random selection in our research, the organization selected representative projects cover projects with quite different characteristics.

4. METHODOLOGY

In this section we provide the details about the methodology we adopted and present our research questions that guided us in this research.

4.1 Estimation Models

Our method for observing the estimation models employed in organizations is to conduct face to face meetings with domain experts in these organizations. The domain experts we have met are the people who are responsible for software effort estimation related activities or for the effort estimation model (if organization has any).

4.2 Metrics for Developer Participation

For developer participation we analyze churn data related to each project in our research. From churn data analysis we first extract the commit information of projects from source code control systems. Since different companies use different types of source control systems, this process takes some time. For each source control system, we write scripts to extract the commit information.

For this research we extract the commit information associated with each developer for each project. For each developer we extract the information of how many lines of code (LOC) added and deleted in each commit. Furthermore, we record the total number of commits performed by each developer. In this paper we will name the total of added and deleted lines of code as edited lines of code.

$$editedLOC = addedLOC + deletedLOC \quad (1)$$

After extracting the commit information related to each developer from the projects, we use two metrics to measure developer participation in projects. The first metric we use is the percentage edited LOC (*Perc.LOC*) for each developer in a project. Percentage edited LOC is the ratio of total edited LOC by a developer to the sum of edited LOC for all developers. The calculation of *Perc.LOC* for *developer_i* is given in Equation 2, where *n* corresponds to the total number of developers.

$$Perc.LOC = \frac{editedLOC_{developer_i} * 100}{\sum_{i=1}^n editedLOC_{developer_j}} \quad (2)$$

The second metric we use to evaluate the developer participation is the percentage of commits (*Perc.Commit*) performed by each developer during the development of a software project. Calculation of *Perc.Commit* is similar to that of *Perc.LOC*. *Perc.Commit* is the ratio of total number of commits performed by a single developer to the number of commits that was conducted by all the developers that worked in this project. The calculation of *Perc.Commit* is given in Equation 3, where *n* is the total number of developers in that project.

$$Perc.Commit = \frac{totalCommits_{developer_i} * 100}{\sum_{j=1}^n totalCommits_{developer_j}} \quad (3)$$

4.3 Research Questions

For guidance during our research we have formed three research questions. With these research questions our aim was to question how widely software effort estimation practices were employed in Turkish software industry and what was the implications of all the estimation and planning activities

Project	Total # of developers	Total edited LOC	Total # of commits	Dev. Language
Project1	3	5,579	72	Java
Project2	110	623,173	12,384	Java+JSP
Project3	97	1,324,956	23,403	SQL
Project4	7	7,034	212	Java
Project5	11	16,358	322	Java+JSP
Project6	19	68,550	2,387	C

Figure 1: 6 projects coming from different organizations. Projects have a wide diversity in terms of their size, developers they use and the languages they were developed in.

on the developer participation. The three research questions we formed are:

RQ1 How widely is effort estimation practices employed?

RQ2 Which effort estimation methods are used?

RQ3 How does estimation and planning activities reflect on developer participation?

5. RESULTS

In this section we present the results of our research. While presenting the results, we follow the research questions we have formed for this research.

RQ1 How widely is effort estimation practices employed?

Previously it was reported by Lederer et. al. in their survey that the software managers puts a high importance on the necessity of software effort estimation [11]. On a scale from 1 (min) to 5 (max) the average importance attributed to software effort estimation practice was 4.7 [11]. In our survey of Turkish industry we observed a similar tendency. However, rather than questioning the importance attributed to software effort estimation, we questioned whether the company conducts any estimation related activity or employs an estimation model. All the companies we questioned performs software estimation related activities. Of course the level of the activities change drastically, some companies have undergone into projects for developing their own algorithmic model, whereas some companies suffice with formal or informal meetings among domain experts. However, the fact that all companies perform estimation related activities at different levels show that all the companies are well aware of the software effort estimation concept. Therefore, we can say that effort estimation practices are widely employed, but the level of employment differs greatly from company to company.

RQ2 Which effort estimation methods are used?

Every company in our survey implicitly or explicitly employs a software effort estimation model. 2 of the companies employs a linear regression based model, in which attributes and the weights related to these attributes are defined by domain experts (the explicit cases). One of these two companies have recently switched to a machine learning based estimation method. The other companies employs expert judgment in their estimation activities. However, they do not strictly follow a Delphi like expert judgment estimation

method. Instead they follow their ad-hoc processes that entails a number of meetings between domain experts for discussing the cost of a project. This shows that mostly used estimation method is expert judgment that is probably due to the fact that it is the least effort requiring method. However, since the companies do not record the initial expert estimates of a project it is almost impossible to track the efficiency of this method. Although being the minority case, two companies use an algorithmic model and only of these two companies have undergone the effort of developing a complex estimation model that employs machine learning algorithms.

RQ3 How does estimation and planning activities reflect on developer participation?

After initial estimates of software effort, each company plans the allocation of their resources. For resource allocation and planning of effort, companies use various project management tools. Some of the companies prefer their in-house developed project management tools, whereas some prefer to use open source or outsourced tools. From our interviews with domain experts in each company, it is our understanding that managers try to distribute the total effort evenly to employees working in particular projects. The track keeping of planned efforts is again performed on the project management tools. Each employee to whom a task was assigned and a certain hours of effort is allocated is responsible to record the actually performed effort into the management tool. However, the consensus among the domain experts from different companies is that the actual effort values is nothing but a confirmation of the planned values. Therefore, the virtually even allocation of resources on the project management tools may be misleading. In that case using churn data may be helpful to get a better view of actual effort spent by developers. On the other hand, none of the companies we have conducted this survey uses churn data for observing the participation of developers to the project.

For each project from these companies, we extracted churn data from their source control system and analyzed the developer participation. We have seen that a big portion of the whole development activity is performed by a very limited number of developers. So the answer to our third research question is that reflection of estimation and planning activities on the developer participation is relatively weak. In other words, the effort of managers to evenly distribute the development effort among developers does not effectively occur in actual development environment. In the next paragraphs, we will evaluate the results for each project separately.

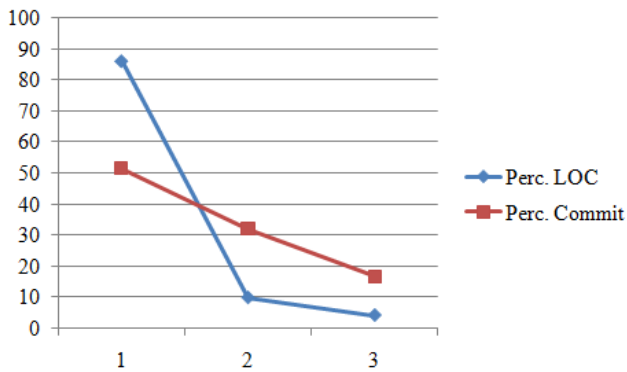


Figure 2: The *Perc.LOC* and *Perc.Commit* information for Project1. Developer 1 is responsible for more than 80% of the whole edited LOC.

In Figure 2 we see the *Perc.LOC* and *Perc.Commit* information regarding Project1. The x-axis corresponds to different developers whereas the y-axis is the percentage value associated with *Perc.LOC* and *Perc.Commit*. Project1 is the smallest project among 6 projects both in terms of size as well as in terms of employed developers. We can see that the majority of the edited LOC in context of this project was developed by a single developer. Developer 1 has edited more than 80% of the total edited LOC. The difference between developers in terms of *Perc.Commit* metric becomes less. However, this may be due to the fact that Developer 1 prefers to commit less often and the other two developers prefer to commit their changes more often.

Figure 3 provides the *Perc.LOC* and *Perc.Commit* values of Project2. Project2 is has the highest number of developers employed. A total of 110 developers have contributed to this project during its development lifecycle. The difference between developer participation in Project2 is less than that of Project1. However, this may be due to the effect of distributing a value of 100% to 110 developers. We see in Figure 3 that after 8 developers, the participation of individual developers goes below 2%. Furthermore, we can see from Figure 3 that the participation of more than half of the developers is less than 1%.

In Figure 4 we see the *Perc.LOC* and *Perc.Commit* values of Project3. Project3 is another densely populated project with 97 developers in total. The behaviour of developer participation of Project3 is similar to that of Project2, i.e. only a small group of developers have a participation of more than 2% and a large number of developers (after developer 31) is responsible for less than 1% of the total edited LOC. We also observe from Figure 4 that the behaviour of *Perc.LOC* is very similar to that of *Perc.Commit*. Although there are some irregularities between the two lines, this may be due to different commit habits of individual developers.

Figure 5 summarizes the developer participation of Project4. Project4 is a relatively small project when compared to Project2 and Project3. The participation pattern among 7 developers employed in Project4 is similar to previous projects. Only two developers account for more than 80%

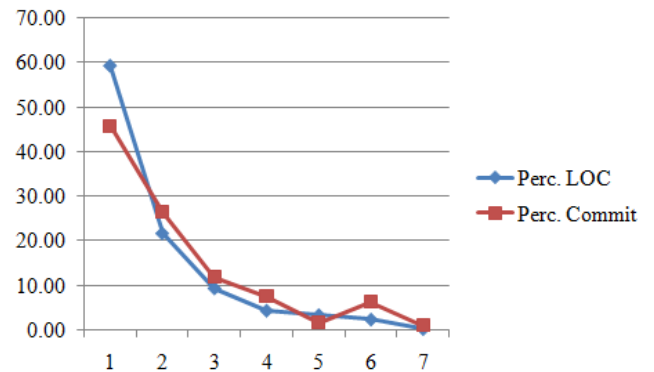


Figure 5: Project4 *Perc.LOC* and *Perc.Commit* information. Developer 1 and developer 2 has more than 80% of the *Perc.LOC*, whereas the rest is shared by other 5 developers.

of the whole edits and the remaining 5 developers account for less than 20% of the edits. Furthermore, like previous projects the *Perc.LOC* and *Perc.Commit* lines go hand in hand in Project4 as well.

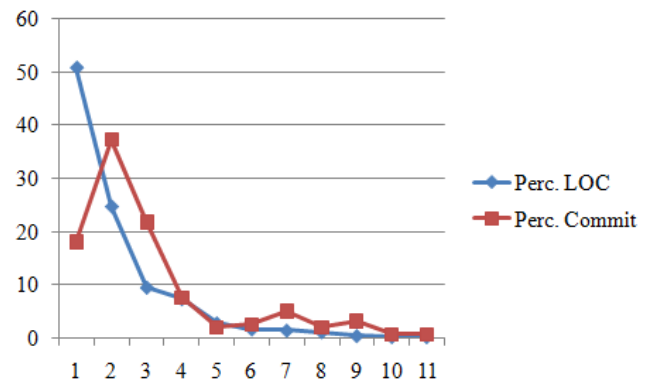


Figure 6: *Perc.LOC* and *Perc.Commit* values of Project5. Similar trend of previous projects continues here. The first three developers account for more than 80% of all edited LOC.

Perc.LOC and *Perc.Commit* information associated with Project5 are given in Figure 6. Project5 has 11 developers employed in the development phase. Among 11 developers the first three developers edits more than 80% of all the edited LOC. With Developer 5 the participation of developers in terms of *Perc.LOC* goes below 5% per developer. One suspicious point in Figure 6 is that there is a considerable difference between *Perc.LOC* and *Perc.Commit* values of Developer 1, who is the top developer. The likely explanation to that case could be that Developer 1 was reviewing the code of other programmers and some of the developers were committing through Developer 1. Although this scenario is very likely in open source communities, this is not the case for this particular company and project. After our inquiry with the company, we have learned that such a review and commit scenario does not exist for Project5.

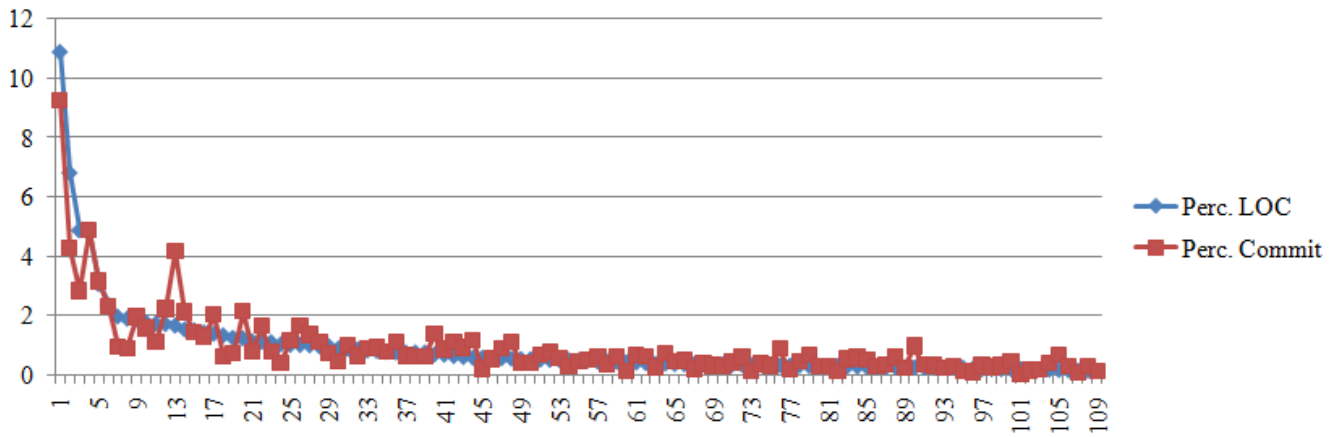


Figure 3: Project2 has the highest number of developers employed. Only 8 developers contribute more than 2% to the project. The participation of more than half of the developers is less than 1%.

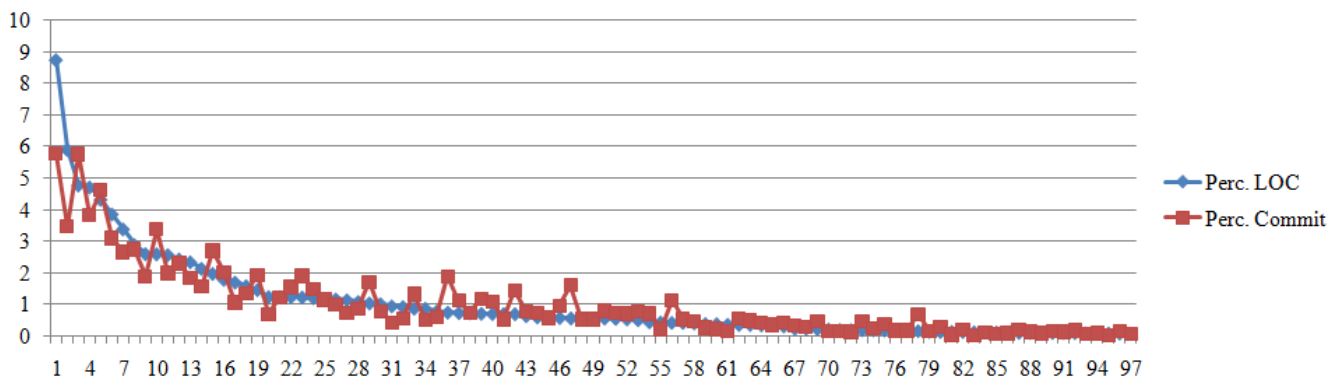


Figure 4: Project3 has 97 developers participating in the development activity. Majority of the developers have limited participation in the development activity. After developer 31, the participation of developers is less than 1%.

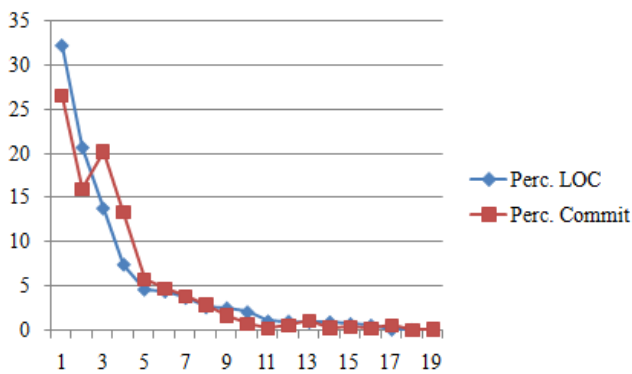


Figure 7: Project6 *Perc.LOC* and *Perc.Commit* values. Out of 19 developers, the first 2 developers account for more than 50% of the total edited LOC.

The last project in our research is Project6. Figure 7 provides the *Perc.LOC* and *Perc.Commit* information of Project6. In Project6 a total of 19 developers are employed. Among 19 developers the first two developers alone have the total of more than 50% of *Perc.LOC*. Starting with Developer 5, the individual participation of developers fall below 5% in terms of *Perc.LOC*. Furthermore, as in the case of all the previous projects the *Perc.Commit* line in Figure 7 behaves very similar to the line of *Perc.LOC*.

6. THREATS TO VALIDITY

The most obvious threat to our research is the selection process of the projects. Unlike widely used method of random selection of projects, our projects are selected by the experts of the organizations. Expert selection is based on the fact that projects shall be representative of the company. When considering the wide range of projects in our research we can consider the dataset as representative of the organizations. However, it cannot be proven that there is no selection bias coming from domain experts.

Another threat to validity is the selected metrics to evaluate the developer participation. Churn data stored in source control systems of open source projects and private organizations have some characteristic differences. Therefore, use of *Perc.LOC* and *Perc.Commit* may be questionable. However, edited LOC and commit count are successfully used in studies about open source projects [4,9], so we also based our metrics on edited LOC and commit count as well.

7. CONCLUSION

In this research we observed the effort estimation and its use in planning activities in the local perspective of Turkish software industry. We also addressed the implications of estimation and planning activities on developer participation in software development. Our dataset in this research included 6 projects coming from different companies that specialize in different domains. Our analysis concerning these 6 projects were performed in two phases. For effort estimation activities, we conducted personal meetings with domain experts to gain an insight into their estimation and planning processes. For the developer participation analysis on the other hand, we extracted churn data from source control systems of different companies and derived our metrics from the extracted churn data.

One implication of our study is that Turkish industry is well aware of the software effort estimation concept and appreciates its importance. All companies in our research implicitly or explicitly employ certain estimation methods. However, estimation practices are mostly far from having a scientific basis. Rather than following literature in effort estimation, organizations suffice with ad-hoc developed processes. Among the companies we worked with, only one company had a solid machine learning based effort estimation model. However, this is not a particular situation to Turkish software industry. It is a known fact that industry practices are usually more simplistic than the methods suggested in the literature.

Another implication of our research is about developer participation. The workload distribution that is being tracked by managers through project management tools does not reflect the actual participation of developers. In all 6 projects we analyzed in our research, we see that majority of the edited LOC is committed by a very limited number of developers. The fact that most of the code is developed by a limited number of developers is previously reported in open source domain [4,9]. In our research we see a very similar pattern in private software companies. Although software development in a private organization has very different characteristics in comparison to open source software development, they share common properties in terms of developer participation.

8. ACKNOWLEDGMENTS

This research has been conducted under the support of...



9. REFERENCES

- [1] V. R. Basili, R. W. Selby, and D. H. Hutchens. Experimentation in software engineering. *IEEE Trans. Softw. Eng.*, 12(7):733–743, 1986.
- [2] B. Boehm, C. Abts, and S. Chulani. Software development cost estimation approaches: A survey. *Annals of Software Engineering*, 10:177–205, 2000.
- [3] B. W. Boehm. *Software Engineering Economics*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1981.
- [4] B. Caglayan, A. Bener, and S. Koch. Merits of using repository metrics in defect prediction for open source projects. In *FLOSS '09: Proceedings of the 2009 ICSE Workshop on Emerging Trends in Free/Libre/Open Source Software Research and Development*, pages 31–36, Washington, DC, USA, 2009. IEEE Computer Society.
- [5] M. Jørgensen. A review of studies on expert estimation of software development effort. *Journal of Systems and Software*, 70:37–60, February 2004.
- [6] M. Jørgensen and T. Gruschke. The impact of lessons-learned sessions on effort estimation and uncertainty assessments. *IEEE Trans. Softw. Eng.*, 35(3):368–383, May-June 2009.
- [7] B. Kitchenham, L. Pickard, and S. L. Pfleeger. Case studies for method and tool evaluation. *IEEE Softw.*, 12(4):52–62, 1995.
- [8] E. Kocaguneli. Better methods for configuring case-based reasoning systems. Master's thesis, Bogazici University, 2010.
- [9] S. Koch. Effort modeling and programmer participation in open source software projects. *Information Economics and Policy*, 20(4):345 – 355, 2008. Empirical Issues in Open Source Software.
- [10] Y. Kultur, B. Turhan, and A. B. Bener. ENNA: software effort estimation using ensemble of neural networks with associative memory. In *SIGSOFT '08/FSE-16: Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of software engineering*, pages 330–338, New York, NY, USA, 2008.
- [11] A. L. Lederer and J. Prasad. Information systems software cost estimating: a current assessment. *J Inf technol.*, 8:22–33, 1993.
- [12] T. Menzies, Z. Chen, J. Hihn, and K. Lum. Selecting Best Practices for Effort Estimation. *IEEE Trans. Softw. Eng.*, 32:883–895, 2006.
- [13] K. Moløkken-Østfold, M. Jørgensen, S. S. Tanilkan, H. Gallis, A. C. Lien, and S. E. Hove. A survey on software estimation in the norwegian industry. *IEEE International Symposium on Software Metrics*, pages 208–219, 2004.
- [14] M. Shepperd. Software project economics: a roadmap. In *FOSE '07: 2007 Future of Software Engineering*, pages 304–315, 2007.
- [15] M. Shepperd and G. Kadoda. Comparing software prediction models using simulation. *IEEE Trans. Softw. Eng.*, pages 1014–1022, 2001.