# Heterogeneous Cross Project Defect Predication Based on Canonical Correlation Analysis

[1]Prakash More, [2]Snehal Patil, [3]Prof. Rahul Kapse

Department of Information Technology,

Pillai HOC College of Engineering and Technology, India

*Abstract:*   Generally, there was lack of training data at the early stage of software testing that limits the efficiency of prediction in practice. Cross project defect prediction uses the data from other projects. In most of the previous efforts assumed that the cross-project defect data have the same metrics set which means that the metrics used, and size of metrics set are same in the data of projects. However, in real scenarios, this assumption may not hold. Here we consider the scenario as heterogeneous cross project defect prediction which add joint features space for associating cross project data and a novel support vector machine algorithm which incorporates the correlation transfer information into classifier designed for prediction. Heterogeneous Cross-Project Defect Prediction (HCPDP) has entire different metrics. In this project, we are going to present CCA+ and the within project defect prediction (WPDP) result.

**Keywords- canonical correlation analysis, defect prediction, heterogeneous metrics.**

## 1. INTRODUCTION

Software defect is an error, mistakes, failure, bug, fault, flaws in a program or in our system that's creates an incorrect or unexpected result, so that problem overcome by the software defect prediction models. There are many models of software defect prediction that are used for quality improvement, system testing, maintenance resource planning and software insurance. Software Defect Prediction (SDP) is one of the most assisting activities of the Testing Phase. It identifies the modules that are defect prone and require extensive testing. This way the testing resources can be used efficiently without violating the constraints. Though software defect prediction is very helpful in testing, it's not always easy to predict the defective modules. There are various issues that hinder the smooth performance as well as use of the Defect Prediction models.

In WPDP (Within Project Defect Prediction), which trained prediction model from historical data to detect the defect proneness of new software modules within the same project

Canonical correlation analysis is a method for exploring the relationships between two multivariate sets of variables (vectors), all measured on the same individual. Canonical correlation terminology makes an important distinction between the words variables and variates. The term variables are reserved for referring to the original variables being analyzed.

The term variates are used to refer to variables that are constructed as weighted averages of the original variables. Thus, a set of Y variates is constructed from the original Y variables. Likewise, a set of X variates is constructed from the original X Variable Suppose you have given a group of students two tests of ten questions each and wish to determine the overall correlation between these two tests. Canonical correlation finds a weighted average of the questions from the first test and correlates this with a weighted average of the questions from the second test. The weights are constructed to maximize the correlation between these two averages.

Heterogeneous Cross-Project Defect Prediction (HCDP) In some systems collecting the same metrics set is challenging as their return in different languages. When collecting metrics for new project it was difficult to obtain the tool license because of this situation, publicly available defect dataset is widely used in defect prediction usually have heterogeneous metrics sets. Therefore CPDP (Cross-Project Defect Prediction) is introduced. In CPDP, the existing systems having same target and source data are being checked with the current system, identify the defect, resolve them and provide the output. But when the source and target projects are different programmer cannot expect correct prediction performance output. Therefore, we are going to implement HCPDP (Heterogeneous Cross-Project Defect Prediction).
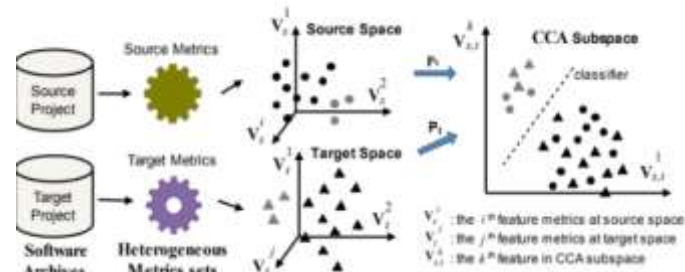
Fig1. The overview of Heterogeneous Cross Project Prediction

## 2.RELATED WORK

Jincheng Nam, Wei Fu, Sunghun kim ,Tim Menzies and Lin Tan proposes Software measurements are the terms used to portray programming projects. Regularly utilized programming measurements for imperfection expectation are unpredictability measurements, (for example, lines of code, Halstead measurements, McCabe's cyclomatic multifaceted nature, and CK measurements) and process measurements. When learning imperfection models, names demonstrate whether the source code is carriage or clean for double arrangement Most proposed deformity forecast models have been assessed on "within-project" deformity forecast (WPDP) settings. As appeared in Figure 1a, in WPDP, each case speaking to a source code record or capacity comprises of programming metric esteems what's more, is marked as carriage or clean.

Ming Cheng ,Hongyan Wan , Meeting Yuan ,Min Jiang proposes Programming imperfection expectation for the most part fabricates models from intra-project information. Absence of preparing information at the beginning time of programming testing limits the productivity of expectation practically speaking. In this manner analysts proposed cross-project deformity forecast utilizing the information from different projects. Most past endeavors accepted the cross-project imperfection information have similar measurements set which implies the measurements utilized and size of measurements set are same in the information of projects. Nonetheless, in genuine situations, this supposition may not hold. Moreover, programming imperfection datasets have the class irregularity issue expanding the trouble for the student to foresee surrenders. In this paper, we progress accepted connection examination for determining a joint component space for partner cross- project information and propose a novel help vector machine calculation which consolidates the relationship exchange data into classifier outline for cross-project expectation. In addition, we take diverse misclassification costs into thought to make the characterization slanting to group a module as a deficient one, mitigating the effect of imbalanced information. Examinations on open heterogeneous datasets from distinctive projects demonstrate that our strategy is more viable, contrasted with best in class strategies.

## 3.PROBLEM STATEMENT

In Within Project Defect prediction (WPDP) and Cross project defect prediction (CPDP) the existing models having same target and source data as well as metrics used for to predict defects. However, we are focusing to analyze the results of defect predication for heterogeneous data and metrics.

## 4.CORRELATION BASED ON CCA+

In light of the acquired UMR for heterogeneous information, we utilize CCA method to decide a typical portrayal (e.g. a joint subspace) for highlights extricated from source and target ventures, so the model prepared in the source venture can be connected to recognize the

test modules in the objective venture. CCA learns two projection vectors $P_s \in R^{ds}$ and $P_t \in R^{dt}$, which boost the accompanying direct relationship coefficient. ρ:

$$\max_{\mathbf{p}_s, \mathbf{p}_t} \rho = \frac{\mathbf{p}_s^T \Sigma_{st} \mathbf{p}_t}{\sqrt{\mathbf{p}_s^T \Sigma_{ss} \mathbf{p}_s} \sqrt{\mathbf{p}_t^T \Sigma_{tt} \mathbf{p}_t}}$$

Where $\sum_{ss}$ and $\sum_{tt}$ represent the within-project covariance matrices of $\overline{X}s$ and $\overline{X}t$ respectively, while Σst = Σts represents the cross-project covariance matrix of $\overline{X}s$ and $\overline{X}t$. Σss , Σtt and Σst are separately defined as

$$\Sigma_{ss} = \frac{1}{N} \sum_{i=1}^{N} (\overline{\mathbf{x}}_s^i - \mathbf{m}_s)(\overline{\mathbf{x}}_s^i - \mathbf{m}_s)^T$$

$$\Sigma_{tt} = \frac{1}{M} \sum_{i=1}^{M} (\overline{\mathbf{x}}_t^i - \mathbf{m}_t)(\overline{\mathbf{x}}_t^i - \mathbf{m}_t)^T$$

$$\Sigma_{st} = \frac{1}{NM} \sum_{i=1}^{N} \sum_{j=1}^{M} (\overline{\mathbf{x}}_s^i - \mathbf{m}_s)(\overline{\mathbf{x}}_t^j - \mathbf{m}_t)^T$$
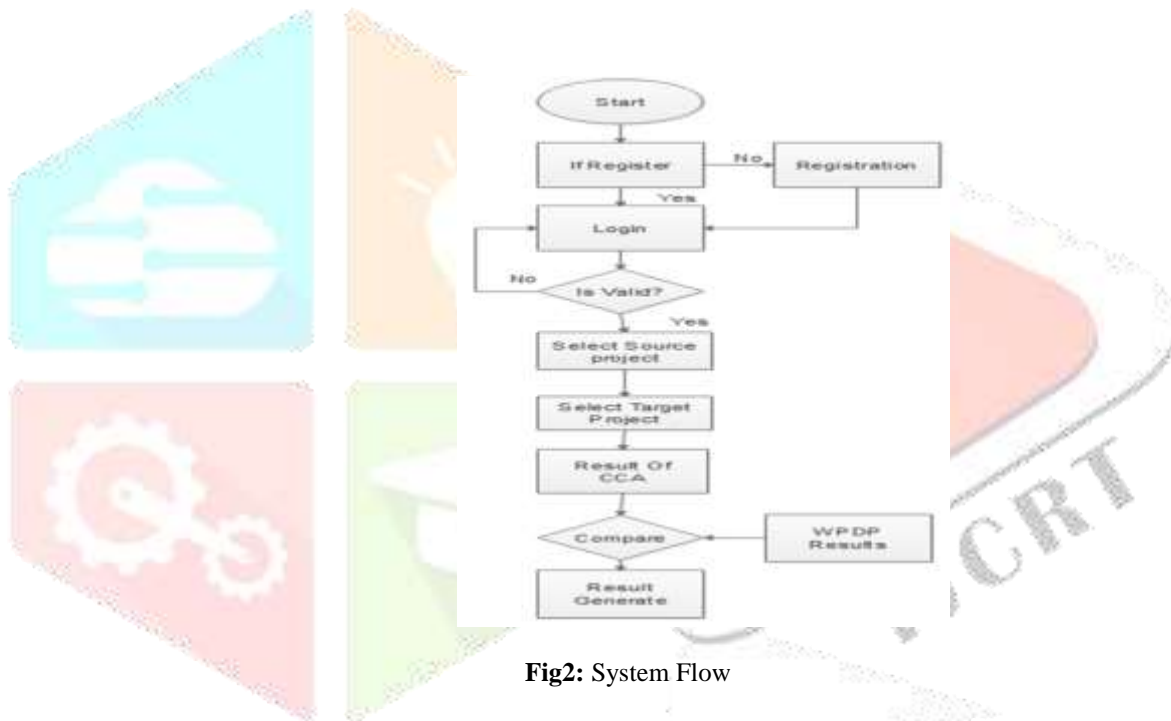


**Fig2:** System Flow

## 5. PROPOSED SYSTEM

 In Within Project Defect Predication (WPDP) historical data are us to predict the Defect in the project the proposed system has to use Within Project Defect Predication Results as references these Referenced results are compared with the HCPDC along with the CCA and Predict the Results. Canonical correlation analysis is a method for exploring the relationships between two multivariate sets of variables (vectors), all measured on the same individual
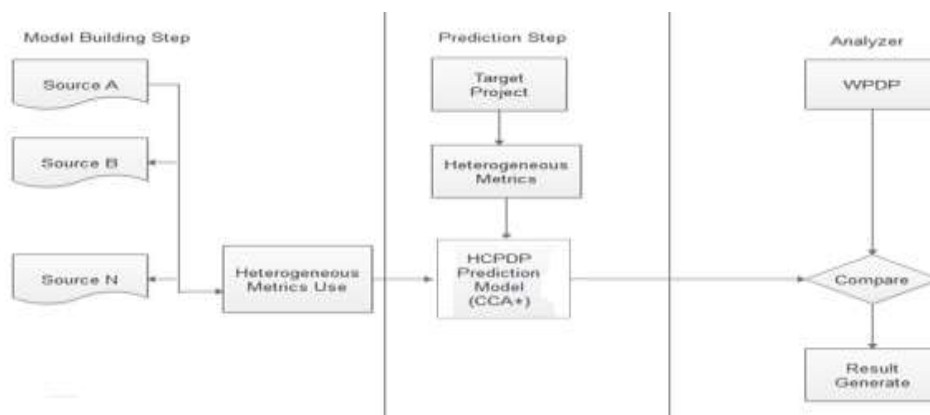
## 6. SYSTEM ARCHITECTURE

**Fig3:** System Architecture

## 7. CONCLUSION AND FUTURE WORK

Cross-venture programming deformity expectation assumes an imperative part in enhancing the nature of a product item if there should be an occurrence of activities without adequate authentic information. Be that as it may, it is hard to lead with heterogeneous measurements set. Also, programming imperfection datasets have the class-irregularity trademark. Without considering this issue, the adequacy of programming deformity expectation would be enormously decreasing. In this paper, we tended to these two essential issues at the same time and proposed a novel cost-touchy connection exchange bolster vector machine technique for heterogeneous deformity expectation. Exploratory outcomes on the open source ventures from various gatherings demonstrated that our strategy is achievable and yields promising outcomes.

For the future work, we will present other complex class unevenness learning systems in the heterogeneous cross project imperfection expectation, and we will assess our approach in more heterogeneous deformity datasets

### REFERENCES

1] X.Y .Jing, F.Wu,X.Dong, F.Qi and B.Xu, "Heterogeneous cross-company defect predication by unified matric representation and CCA based transfer learning ," In Proceeding of the 10th Joint Meeting on Foundations of Software Engineering 2015 ,pp.496-507.

[2] Jaechang Nam, Wei Fu, Sunghun kim ,Tim Menzies and Lin Tan Heterogenous Defect Prediction,IEEE transaction on Software engineering

[3] Ming Cheng ,Hongyan Wan , Meeting Yuan ,Min Jiang .Heterogeneous Defect Prediction via Exploiting Correlation Subspace The 28th International conference on Software Engineering

**References Sites**

[1]https://www.researchgate.net/publication/305972951_Heterogeneous_Defect_Prediction_via_Exploiting_CorrelationSubspace

[2]http://www.worldscientific.com/doi/abs/10.1142/S0218194016710017

[3]https://www.slideshare.net/mobile/hunkim/heterogeneous-defect-prediction-esecfse-2015.