

# Study on Residual Defect Prediction using Multiple Technologies

Wanjiang. Han

School Of Software Engineering, Beijing University of Posts and Telecommunication, Beijing, China  
Email: hanwanjiang@bupt.edu.cn

Lixin. Jiang

Department of Emergency Response, China Earthquake Networks Center, Beijing, China  
Email: jl@seis.ac.cn

Tianbo. Lu, Xiaoyan. Zhang, Sun Yi

School Of Software Engineering, Beijing University of Posts and Telecommunication, Beijing, China  
Email: { lutb, xiaoyan, sunyisse }@bupt.edu.cn

**Abstract**—Finding defects in a software system is not easy. Effective detection of software defects is an important activity of software development process. In this paper, we propose an approach to predict residual defects, which applies machine learning algorithms (classifiers) and defect distribution model. This approach includes two steps. Firstly, use machine learning Algorithms to get defect classification table, then confirm the defect distribution trend referring to several distribution models. Experiment results on a GUI project show that the approach can effectively improve the accuracy of defect prediction and be used for test planning and implementation.

**Index Terms**—residual defect prediction, defect distribution model, software defect classification, defect trend, classifiers

## I. INTRODUCTION

With the rapid development of information technology, the application of computer software is more and more widely. High efficient and security software system is highly dependent on the software reliability, and the software defects have become the underlying causes which lead to the systems error, failure, collapse or even the disaster. Software defect prediction is a very important research topic in software engineering. It is based on the defect records in historical data to predict the defects in the future. It helps software project planning and process management.

The growing demand for higher operational efficiency and safety in industrial processes has resulted in a huge interest in defect detection techniques. Engineering researchers and practitioners remain concerned with accurate defect prediction when building systems. Defect prediction has both safety and economic benefits in technical systems by preventing future failures and further improves process maintenance schedules. Lack of adequate tools to evaluate and estimate the cost for a software system failure is one of the main challenges in

systems engineering.

For a software development project, it is highly desirable to reduce software defects. A variety of information products emerged and the requirements of product quality are also increasing. By analyzing and studying defects, building defect model can help to find solutions for defects quickly and efficiently, which provides a good exploration to improve the quality of product and will help to forecast the reliability of the software and improve software reliability[1].

Software defects include detection defects and residual defects. Detection defects are the defects which were detected before software release. Residual defects are the defects which were detected after software release. The purpose of residual defects' prediction is to keep the software defects number under the acceptable level in testing times[2][3][4].

During testing process, there will be a large accumulation of defect data, it is necessary to merge the similar defect to the same type of defect for unify solution easily. And defect model can predict residual defects. Defects can be effectively managed and predicted by defect classification and prediction model[5].

This paper discusses the residual defect prediction based on more techniques by a large number of defect data from testing. Firstly, This paper describes how to use classifiers to classify software defects, then gives the distribution of residual defects. Finally, the residual defects can be predicted. It provides a favorable recommendation for product development and design.

## II. RELATED CLASSIFIERS

In recent years, the use of machine learning algorithms (classifiers) has proven to be of great value in solving a variety of problems in software engineering including software defect prediction[6][7].

The classifier is an hypothesis about the true classification function that is learned from, or fitted to,

training data. The classifier is then tested on test data. In supervised learning, especially for multivariate data, a classification function  $y = f(x)$  predicts one (or more) output attribute(s) or dependent variable(s) given the values of the input attributes[8].

This section will discuss three classifiers: Association rules, Decision tree and k-Nearest Neighbour[8][9][10].

#### A. Association Rules

Association Rule, Proposed by Agrawal et al in 1993, is an important data mining model studied extensively by the database and data mining community. Assume all data are categorical. It is the task of discovering association rules that occur frequently in a given transaction data set[11].

An association rule is a pattern that states when X occurs, Y occurs with certain probability.

Its task is to find certain relationships among a set of data (item set) in the database.

Association rules mining is interested in finding frequent rules that define relations between unrelated frequent items in databases, and it has two main measurements: support and confidence values. Confidence values are measurements of rule's strength, while support value corresponds to statistical significance.

An association rule states that an item or group of items implies the presence of another item with some probability. Association rules reflect the relationship among data or is a study whether the generation of a data can speculate the generation of another data. Data association reflect the relationship in Database. The association degree can be expressed through support and confidence. Thus, finding the relationship among data in the transaction database is the mining purposes of association rule.

Apriori algorithm has been shown as a classical association rule in mining algorithms that have been cited so far, which is used to find frequent item sets in a database and to generate Association Rules from the frequent item sets. Algorithm is how to regenerate all the frequent itemsets. It uses an iterative approach layer by layer, applying K-Itemset to search for (K+1)-Itemset[12].

#### B. Decision Tree

Decision tree (DT) induction is one of the simplest and yet most successful forms of supervised learning algorithm. It has been extensively pursued and studied in many areas such as statistics and ML [13] for the purposes of classification and prediction.

DTs are non-parametric (no assumptions about the data are made) and a useful means of representing the logic embodied in software routines. A decision tree takes as input a case or example described by a set of attribute values, and outputs a Boolean or multi-valued decision.

A classification tree, as opposed to a regression tree means that the response variable is qualitative rather than quantitative. In the classification case, when the response variable takes value in a set of previously defined classes the node is assigned to the class which represents the highest proportion of observations. Whereas, in the regression case, the value assigned to cases in a given

terminal node is the mean of the response variable values associated with the observations belonging to a given node. Note that in both cases, this assignment is probabilistic, in the sense that a measure of error is associated with it. Clustering trees just group instances in leaves.

#### C. k-Nearest Neighbour

One of the most venerable algorithms in machine learning is the nearest neighbour (NN). Nearest-neighbour methods are sometimes referred to as memory-based reasoning or instance-based learning (IBL) or case-based learning (CBL) techniques and have been used for classification tasks. They essentially work by assigning to an unclassified sample point the classification of the nearest of a set of previously classified points.

A further non-parametric procedure of this form is the k-nearest neighbour (k-NN) approach. To classify an unknown pattern, the k-NN approach looks at a collection of the k nearest points and uses a "voting" mechanism to select between them, instead of looking at the single nearest point and classifying according to that with ties broken at random. If there are ties for the  $k^{th}$  nearest observations, all candidates are included in the vote[14].

### III. RELATED DEFECT DISTRIBUTION MODEL

Software defect prediction technology includes lots of models. Here, We will discuss several models used in this paper.

#### A. Rayleigh Distribution Model

Rayleigh model is a common reliability model. It can forecast the defect distribution during software life cycle. This kind of model is based on the Weibull statistical distribution. Weibull distributed reliability analysis is widely used in different field. The probability density function of Weibull distributed tag end is gradually converge to 0, but never equal to 0. The experts of Trachtenberg and IBM have researched on the software project faults and they found that the distribution accords with Rayleigh distribution model. The probability distribution density function of Rayleigh model is as in (1).

$$f(t) = 2K(t/c^2)e^{-(t/c)^2} \quad (1),$$

the cumulative distribution function is as in (2)

$$F(t) = K(1 - e^{-(t/c)^2}) \quad (2),$$

where K is the total faults. T stands for time, c is an

constant and  $c = \sqrt{2}t_m$ , where  $t_m$  is the time when  $f(t)$  reached maximum, and  $F(t_m) / K$

approximately equal to 0.4. Therefore, we can estimate the total number of defects at a certain time as well as the specific Rayleigh distribution parameter. Thus we can simplify the counting process. It is easy to control the quality of the enterprise performance goal by using

Rayleigh model. It is necessary to take measures to correct it when the process appeared abnormal[15][16].

*B. Exponential Distributed Model*

Exponential model accords with defect distribution for the testing phase, especially the acceptance testing phase. The basic principle is that the defect occurs at this stage or the failure mode is a good indication of the reliability of the product. Exponential model, also known as reliability growth model, is divided into fault/failure count model and the failure interval model. The defect probability distribution's density function of the exponential model is as in (3).

$$f(t) = K(\lambda e^{-\lambda t}) \tag{3}$$

And the defects cumulative distribution function is as in (4).

$$F(t) = K(1 - e^{-\lambda t}) \tag{4}$$

Where t is time, K is the total number of defects,  $\lambda$  is the defect detection rate or failure occur rates. Exponential model is the most simple and most important model in software reliability research, and also is a basis of many other reliability growth model. Its basic assumptions are:

- Each defect leading to failure found in test phase is equal;
- The defect repair time is negligible;
- Each defect can be perfectly repaired and the repair will do not introduce new defects, and so on.

Defects exponential probability density function and cumulative distribution function are shown in Fig. 1 and Fig. 2[15].

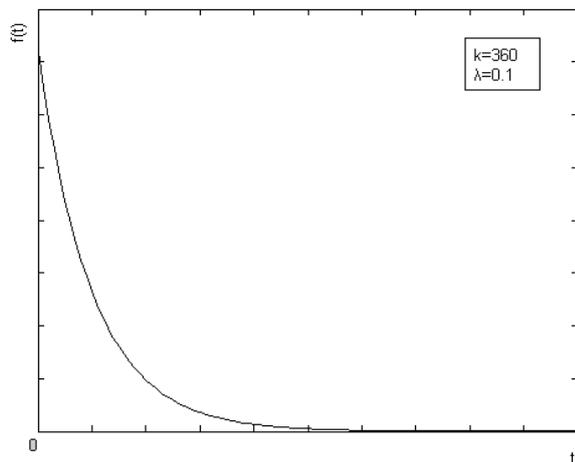


Figure 1. The defect exponential distribution probability density function

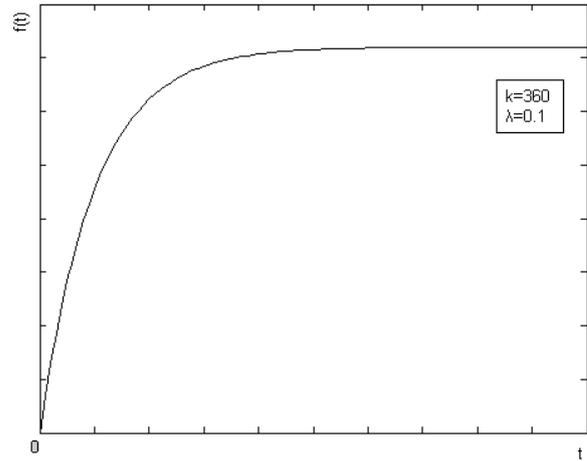


Figure 2. The defects exponential cumulative distribution function

*C. S-curve Distributed Model*

Yamada and other people proposed a test period not only includes detecting the defects but also includes isolating defects. When there comes a defect, we need to find the reason why defects are disabled. Therefore, there exists a time delay for finding a defect until we report it. The accumulated time delayed defects accord with S curve distribution and this is called delayed S-curve model which is also a reliable increment model. S-curve model meet the requirement of nonhomogeneous Poisson process. Its CDF is as in (5).

$$F(t) = K(1 - (1 + \lambda t)e^{-\lambda t}) \tag{5}$$

where t represents time; k represents total fault number;  $\lambda$  represents the possibility of detecting a fault. The PDF is as in (6).

$$f(t) = K\lambda^2 t e^{-\lambda t} \tag{6}$$

In 1984, Oghba[20] proposed another S-curve distributed model called transformed S model. This model considers the detected faults are interdependent and the more defects are detected the more faults will be detected later. The accumulated fault distributed function (CDF) is as in (7).

$$F(t) = K(1 - e^{-\lambda t}) / (1 + \phi e^{-\lambda t}) \tag{7}$$

where t represents time; k represents total fault number;  $\lambda$  represents the possibility of detecting a defect. Its PDF is as in (8) [15].

$$f(t) = K\lambda e^{-\lambda t} K(1 + \phi) / (1 + \phi e^{-\lambda t})^2 \tag{8}$$

*D. The lognormal distribution model*

Lognormal model which can predict the defect distribution in the entire life cycle of software is a most commonly used reliability model. Lognormal model is based on the Weibull statistical distribution, and Weibull distribution is widely used statistical distribution family

in different areas of a reliability analysis, which has a large number of statistical data to support [15][17]. Set software defects in the program block, B1, B2, ..., Bm, are m blocks of the program; b1, b2, ..., bn are n branches of the program. The probability of each branch being executed obeys normal distribution, the probability of each block being executed can be expressed as in (9).

$$P(B) = \prod_{j=1}^n P(b_{i,j}) \tag{9}$$

Since  $P(b_{i,j})$  obeyed normal distribution, according to the central limit theorem,  $P(B_i)$  obeys the lognormal distribution. The lognormal distribution of defects exposure rate  $\lambda$  can be expressed as in (10).

$$P(\lambda) = \frac{N}{\lambda \delta \sqrt{2\pi}} e^{-\frac{(\ln \lambda - \mu)^2}{2\delta^2}} \tag{10}$$

Where N is the initial number of defects;  $\lambda$  as a variable;  $\mu$  is the average of the logarithm of  $\lambda$ ;  $\delta^2$  is the variance of the logarithm of  $\lambda$ . Defects lognormal probability density function and cumulative distribution functions are shown in Fig. 3 and Fig. 4.

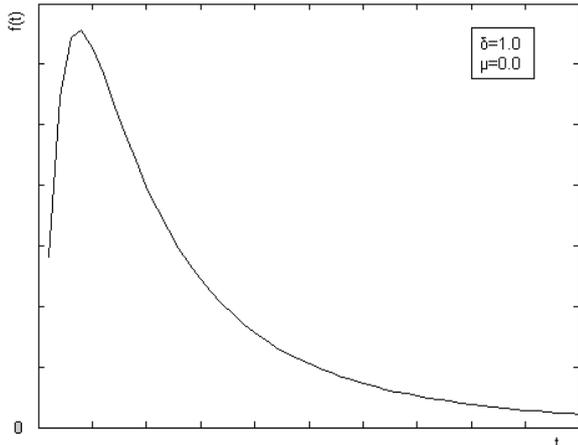


Figure 3. Defects probability density function

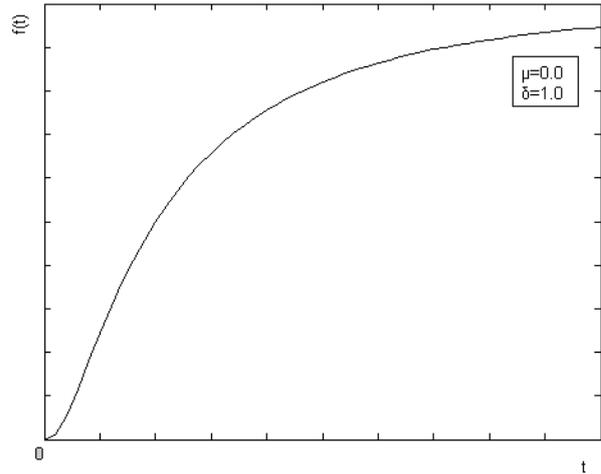


Figure 4. The number of defects cumulative distribution function

E. Bayesian Belief Networks

Bayesian Belief Networks (BBNs) is a graphical network. The graph is made up of nodes and arcs where the nodes represent uncertain variables and the arcs represent the causal relationships between variables.

There are many advantages of using BBNs. The benefits of using BBN include:

- specification of complex relationships using conditional probability statements;
- use of "what - if? Analysis and forecasting of effects of process changes;
- explicit modeling of "ignorance" and uncertainty in estimates, forecasting with missing data.

There are three stages to set up software residual defects prediction model using BBNs : (1) Setting variables set and field, (2) Constructing Network topology, (3) ascertainment local probability distribution. Conditional probability expresses the relation of nodes and their parents' nodes. Its transcendental probability is conditional probability which has no parents node. These three steps are implemented interactively, not serially [15].

IV. DEFECT CLASSIFICATION USING CLASSIFIERS

Through a large number of test data, defect distribution can be summed up. We applied association rules of data mining, and studied the relative defects, which can analyze efficiently and prevent product defects, then improve product quality.

An association rules-based method has the strength of not requiring prior knowledge on the system. The association rules approach compares with existing data mining methods but has the slight advantage of detecting some anomalies which (otherwise) could have been overlooked by conventional approaches.

In order to show the steps of residual defect prediction, this paper uses a typical GUI defect data to analyze the classification and defects distribution. The data is shown in Table I.

TABLE I.  
GUI DEFECT CLASSIFICATION

Term	Defect classification
1	Interface display problem
2	Interface indication problem
3	Wrong character problem
4	Punctuation format problem
5	Inconformity semantic expression
6	Unreadable codes
7	Messy versions
8	Inconformity pictures
9	Inappropriate learning problem
10	Sequencing problem
11	Boundary problem

From Table I, we use the association rules-based technology to get a defect association relationship, shown as Table II.

Relative analysis informs that a happening of one defect may incur another or many other defects occur. For example, the problem of ‘interface displaying’ is associated with ‘interface indication problems’, ‘consistency problems’ and ‘boundary problems’.

After the threshold values of minimum support and minimum confidence support are given, frequent item sets and association rule will be generated by applying Apriori arithmetic. We focus on the strong association rule which is the rule satisfying the minimum support and minimum confidence[18][19][20].

TABLE II GUI DEFECT ASSOCIATION

Defect types	Defect associations
Interface display problem	Interface indication problem, consistency problem, boundary problem
Interface indication problem	Interface display problem, consistency problem, boundary problem, messy codes
Messy Versions	Interface display problem
Graphical inconsistency	Interface indication problem, consistency problem
Learnability problem	Interface indication problem, consistency problem, boundary problem
Boundary problem	Interface display problem, Interface indication problem, Consistency problem
Consistency problem	Interface display problem, Interface indication problem, Boundary problem, Graphical inconsistency

By this relationship, we can predict more defects.

IV. DEFECT PREDICTION USING DITRIBUTION MODEL

After getting the defect relationship, we can study the defects distribution, then we can predict the residual defects. In order to obtain accurate distribution of software defects, this paper uses many typical project defect data to analyze

the distribution of defects. Table III shows the average data of several projects.

TABLE III DEFECT DATA

Date	Defects
1	6
2	7
3	9
4	11
5	13
6	15
7	16
8	19
9	22
10	25
11	28
12	31
13	35
14	40
15	46
16	51

We draw the defect trend picture using data from TABLE III, as shown in Fig 5. We can infer from the trend of the data fitting curve that the distribution of software defects is more consistent with the lognormal distribution. So we can select lognormal distribution model to predict residual defects.

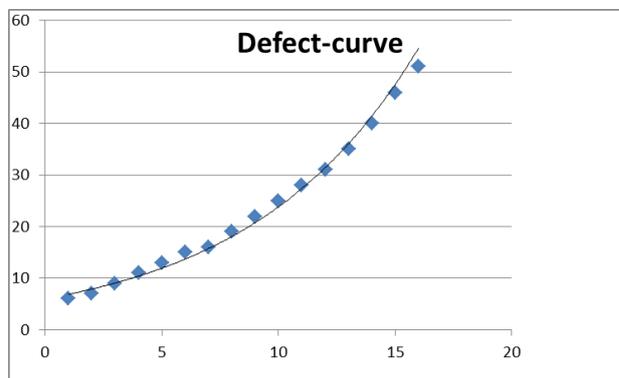


Figure 5. Defect trend fitting picture

We used the least square method to make the regression analysis about the parameters of the model and in the last we have achieved the expression of the model, shown as in (8).

$$F(t) = 123(1 - e^{-0.0325t}) \tag{8}$$

Then, using this model to predict the project defects during the test phase. Fig. 6 shows the difference between

estimated data and actual data. Obviously, the two curves are basically identical.

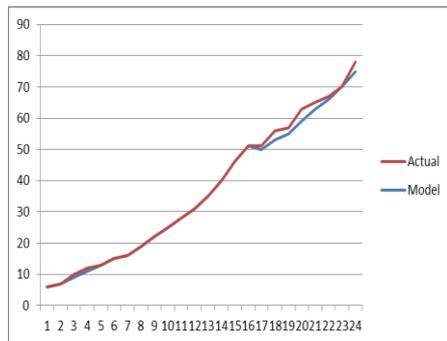


Figure 5. The comparison between the predicted and actual data

We can evaluate the accuracy of estimation by MRE (Magnitude of Relative Error). Namely,  $MRE = |\text{Actual} - \text{Predicted}| / \text{Actual}$ . The MMRE (Mean of the Magnitude of Relative Error) is also a useful evaluation tool, Its Equation is as in (9).

$$MMRE = \frac{1}{n} \sum_{i=1}^n MRE_i \quad (9)$$

For Fig. 5, the MMRE is 8.6%. The calculated MMRE of similar projects is close to 8.9%.

## V. CONCLUSIONS

This paper analyzed the related technologies about classifiers and distribution model. From the representative collected software defects data of GUI projects, the paper used several classifier algorithms to get defect classification table, then applied mathematical methods to show that the distribution of this kind of software project defects is consistent with the lognormal distribution better. If we can find out which distribution the software defects obeyed in accordance with the defects classification, we can use the fault injection method to simulate software fault, and study the accelerated test method under certain defects distribution, which can effectively improve the software test coverage, reduce test time, reduce cost of test.

Further study will continuously work on specializing projects and keep on researching other prediction model as well.

## ACKNOWLEDGMENT

This work was supported in part by the National Natural Science Foundation of China (Grant No. 61170273).

## REFERENCES

- [1] GAFFNEY J E. Estimating the number of faults in code [J]. IEEE Transactions on Software engineering, 1984, 10(6):459- 464
- [2] Fenton N, Neil M. A Critique of Software Defect Prediction Models. IEEE Transaction on Software Engineering, 1999, 25(5):675- 689.

- [3] May J. Fault prediction for software[M] Oxford: Oxford University Press, 1995
- [4] C. Wohlin, and P Runeson, Defect Content Estimations from Review Data, In Proceedings of the Twentieth International Conference on Software Engineering, IEEE Computer Society Press, pp.400 - 409, 1998.
- [5] Fenton N, Krause P, Neilm M. A probabilistic Model for Software Defect Prediction.
- [6] L. Breiman, Bagging predictors. Machine Learning, 26 (2) (1996a), pp. 123-140, 1996a.
- [7] L. Breiman, J. Friedman, R. Olshen and C. Stone. (1984). Classification and Regression Trees, Wadsworth, 1984.
- [8] E. Bauer and R. Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting and variants. Machine Learning, 36, 1/2, pp. 105-139, 1999.
- [9] T. Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. Machine Learning, 40 (2), pp. 139-158, 2000.
- [10] Bhakesispho Twala. Software Faults Prediction Using Multiple Classifiers, [www.cnki.net](http://www.cnki.net).
- [11] C.J.C. Burges. A Tutorial on Support Vector Machines for Pattern Recognition. Data Mining and Knowledge Discovery. Kluwer Academic Publishers, 1998.
- [12] Wanjiang Han, Tianbo Lu and Sun Yi. *Research on the Problem Model of GUI based on Knowledge Discovery in Database*. Proceedings of the 2013 International Conference on Software Engineering and Computer Science. 2013:5-9.
- [13] J.R. Quinlan. C.4.5: Programs for machine learning. Los Altos, California: Morgan Kauffman Publishers, INC., 1993.
- [14] P. Domingos. On the optimality of the simple Bayesian classifier under zero-one loss. Machine Learning, 29, pp. 103-130, 1997.
- [15] Ouyang Xin, Shuai Chun-Yan. The research of lognormal probability distribution of software defects, *2010 International Conference on Circuit and Signal Processing*.
- [16] Mullen R. The lognormal distribution of software failure rates application to software reliability growth modeling[c]. Proc of the 9th international symposium on software reliability engineering Washington DC: IEEE computer society, 1998:134.
- [17] Mullen R. The lognormal distribution of software failure rates: origin and evidence[c]. Proc of the 9th international symposium on software reliability engineering Washington DC: IEEE computer society, 1998:124.
- [18] Wanjiang Han, "Research of the Defect Model Based on Similarity and Association Rule" unpublished.
- [19] Wanjiang Han, "Research on Software Fault Distribution for Web Application" unpublished
- [20] Li C, Fan M (2004) Generating association rules based on threaded frequent pattern tree. *Comput Eng Appl*. 2004; 4:188-192 (in Chinese).



**Wan-Jiang HAN** was born in HeiLongJiang province, China, 1967. She received her Bachelor Degree in Computer Science from Hei Long Jiang University in 1989 and her Master Degree in Automation from Harbin Institute of Technology in 1992.

She is an assistant professor in School Of Software Engineering, Beijing University of Posts and Telecommunication, China. Her technical interests include software project management and software process improvement.



**Li-Xin Jiang** was born in HeiLongJiang province, China, 1966. He received his Bachelor Degree and Master Degree in physical geography from Beijing University in 1989 and 1992.

He is a professor in the department of Emergency Response of China Earthquake Networks Center. His technical interests include software cost estimation and Emergency Response software development.

**Xiao-Yan Zhang** was born in Shandong Province, China, 1973. She received her Master Degree in Computer Application in 1997 and her PH.D Degree in Communication and information system from Beijing University of Posts and Telecommunication, China, in 2011.



She is an Associate professor in School of Software Engineering, Beijing University of Posts and Telecommunications, Beijing, China. Her technical interests include software cost estimation and software process improvement.



**Tian-Bo Lu** was born in Guizhou Province, China, 1977. He received his Master Degree in computer science from Wuhan University in 2003 and his PH.D Degree in computer science from the Institute of Computing Technology of the Chinese Academy of Sciences in 2006.

He is an Associate professor in School of Software Engineering, Beijing University of Posts and Telecommunications, China. His technical interests include information and network security, trusted software and P2P computing.