

## Implementation of Split Sampling for Decision Tree and K-Nearest Neighbor Algorithms in The DKI Jakarta Legislative Election

Bambang Wisnu Widagdo\*<sup>1</sup>, Muhammad Rizky Fadillah<sup>2</sup>, Mochamad Adhari Adiguna<sup>3</sup>, Sudarno Wiharjo<sup>4</sup>, Murni Handayani<sup>5</sup>

<sup>1,2,3,4,5</sup> Universitas Pamulang, Indonesia

Email: <sup>1</sup>dosen02092@unpam.ac.id, <sup>3</sup>dosen01864@unpam.ac.id, <sup>4</sup>dosen02710@unpam.ac.id

### Abstract

The 2009 legislative election was contested by 44 political parties, consisting of national and local parties. In the 2009 Legislative Election for DKI Jakarta, there were 2,268 candidates for the Regional House of Representatives (DPRD) from 44 parties competing for 94 seats in the DKI Jakarta Regional People's Representative Council. Data mining is a series of processes aimed at discovering added value in the form of information that has not been previously known manually from a database. The classification method can be used to predict the results of legislative elections. In this study, the author employs the Decision Tree and K-Nearest Neighbor classification algorithms. This research utilizes several data sampling techniques, namely Linear Sampling, Shuffled Sampling, and Stratified Sampling. The data split partition used in this study was 80% for training and 20% for testing. The software tool utilized was RapidMiner. The performance variables measured include Recall, Precision, and Accuracy. The results of this study indicate that, overall, Linear Split Sampling outperforms Shuffled Split Sampling and Stratified Split Sampling. For the Decision Tree algorithm, Linear Split Sampling achieved a Recall of 100%, Precision of 82.05%, and Accuracy of 98.46%. Shuffled Split Sampling recorded a Recall of 81.82%, Precision of 85.71%, and Accuracy of 98.46%. Meanwhile, Stratified Split Sampling obtained a Recall of 100%, Precision of 82.05%, and Accuracy of 97.80%. Meanwhile, for the K-Nearest Neighbor (KNN) algorithm, Linear Split Sampling achieved a Recall of 93.75%, Precision of 75%, and Accuracy of 97.36%. Shuffled Split Sampling recorded a Recall of 59.09%, Precision of 81.25%, and Accuracy of 97.36%. Stratified Split Sampling obtained a Recall of 31.58%, Precision of 85.71%, and Accuracy of 96.92%.

**Keywords:** *Data Mining, Decision Tree, K-Nearest Neighbor, Linear Sampling, Shuffled Sampling, Stratified Sampling*

## 1. INTRODUCTION

The legislative election is a fundamental democratic process through which citizens elect their representatives to legislative bodies at the national and regional levels. In Indonesia, the legislative election for the Special Capital Region of Jakarta (DKI Jakarta) holds significant importance due to the region's status as the nation's capital and economic center. The election determines the members of the Jakarta Regional People's Representative Council (DPRD DKI), who are responsible for legislating regional policies, overseeing the regional government, and representing the interests of Jakarta's residents.

In the 2009 legislative election, DKI Jakarta witnessed a highly competitive political atmosphere, with 44 political parties—both national and local—competing for a total of 94 seats in the DPRD DKI. A total of 2,268 candidates contested the election, reflecting the diverse political landscape and the active participation of various segments of society. The high level of competition underscores the complexity of predicting election outcomes, as voter preferences are influenced by multiple factors, including party platforms, candidate profiles, and socio-political issues.

Given the dynamic and multifaceted nature of elections in DKI Jakarta, data-driven approaches such as data mining and machine learning have become increasingly relevant for analyzing voter behavior and forecasting results. By leveraging historical electoral data, these methods can uncover hidden patterns and generate insights that are not easily observable through manual analysis, thereby contributing to more informed political strategies and public understanding.

General elections serve as a means of exercising the people's sovereignty in the Unitary State of

the Republic of Indonesia, based on Pancasila and the 1945 Constitution (Law of the Republic of Indonesia No. 10 of 2008). The objective of the general elections, as stated in Article 3 of Law No. 8 of 2012, is to elect members of the House of Representatives (DPR), Provincial Legislative Councils (DPRD Provinsi), and Regency/Municipal Legislative Councils (DPRD Kabupaten/Kota) within the Unitary State of the Republic of Indonesia, which is founded on Pancasila and the 1945 Constitution of the Republic of Indonesia.

The 2009 legislative election was contested by 44 political parties, consisting of both national and local parties. In the 2009 Legislative Election in DKI Jakarta, a total of 2,268 candidates from 44 parties competed for 94 seats in the DKI Jakarta Regional House of Representatives (DPRD).

Data mining is a series of processes used to extract added value in the form of information that was previously unknown manually from a database. Data mining is the process of extracting valuable knowledge, patterns, and insights from large volumes of data that may otherwise remain hidden. It integrates techniques from statistics, machine learning, database systems, and artificial intelligence to transform raw data into meaningful information for decision-making. In the era of big data, where organizations collect vast amounts of structured and unstructured data, data mining has become an essential tool for discovering trends, predicting outcomes, and supporting strategic planning.

The data mining process typically involves several stages, including data collection, preprocessing, transformation, pattern discovery, and evaluation. Various methods and algorithms—such as classification, clustering, regression, and association rule mining—are applied depending on the nature of the problem and the type of data being analyzed. Classification, in particular, is widely used for predictive tasks, where historical data is employed to build models capable of forecasting future events or behaviors.

Applications of data mining are broad and span multiple domains, including business analytics, healthcare, finance, manufacturing, education, and politics. In the context of elections, data mining can be used to analyze voting patterns, predict election results, and identify factors that influence voter preferences. By uncovering hidden relationships within electoral data, data mining provides a scientific foundation for understanding complex social phenomena and supports evidence-based decision-making.

The information is obtained by extracting and identifying significant or interesting patterns from the data within the database. Data mining is primarily used to discover knowledge from large databases, and is therefore often referred to as Knowledge Discovery in Databases (KDD) (Vulandari, 2017). One of the main functions of data mining is classification. Classification is the process of discovering a model or function that explains or distinguishes data concepts or classes, with the goal of predicting the class of an object whose label is known (Haskett, 2000). Algorithms that can be used in classification methods include Tree-based algorithms (CART, ID3, C4.5, Credal DT, Credal C4.5, Adaptive Credal C4.5), Naive Bayes (NB), K-Nearest Neighbor (KNN), Linear Discriminant Analysis (LDA), and Logistic Regression (LogR) (Wahono, 2020).

RapidMiner is a tool used in data mining processes. RapidMiner offers superior speed compared to the Weka data mining tool (Faid, M., et al., 2019). RapidMiner is a powerful, open-source data science platform widely used for data mining, machine learning, predictive analytics, and business intelligence applications. Originally developed as an academic research project in 2001 at the Technical University of Dortmund, Germany, RapidMiner has evolved into a comprehensive software environment that supports the entire data science lifecycle—from data preparation and visualization to model building, evaluation, and deployment.

One of RapidMiner's key strengths is its graphical user interface (GUI) based workflow design, which allows users to create analytical processes without the need for extensive programming knowledge. This makes the platform accessible to both beginners and experienced data scientists. The software provides a rich set of operators for tasks such as data preprocessing, transformation, modeling, validation, and visualization. Additionally, it supports a variety of machine learning algorithms, including Decision Tree, K-Nearest Neighbor (KNN), Support Vector Machines (SVM), Naïve Bayes, and neural networks.

RapidMiner also offers flexible data sampling and partitioning methods, such as Linear Sampling, Shuffled Sampling, and Stratified Sampling, which enable users to split datasets efficiently for training

and testing purposes. The platform integrates well with external databases, programming languages such as Python and R, and big data technologies, making it a versatile choice for diverse analytical tasks.

Split Data is an operator in RapidMiner that generates the desired number of subsets from a given ExampleSet. The ExampleSet is partitioned into several subsets based on specified relative sizes. The Split Data operator utilizes several sampling methods to build these subsets, including Linear Sampling, Shuffled Sampling, Stratified Sampling, and Automatic Sampling. To evaluate the performance of a classification algorithm, the Performance (Classification) operator can be used. This operator is designed to assess the statistical performance of classification tasks. It provides a list of performance metric values related to the classification task. (<https://docs.rapidminer.com/>)

Classification is one of the fundamental techniques in data mining and machine learning, aimed at predicting the categorical class labels of new instances based on patterns learned from historical data. The primary objective of classification is to develop a model, or classifier, that can accurately assign an input to one of several predefined categories. This process is widely used in predictive analytics, where labeled training data is analyzed to build rules or decision boundaries that can be applied to unseen data.

The classification process typically involves two main phases: training and testing. In the training phase, the algorithm learns from a dataset that contains both input features and known class labels. During the testing phase, the trained model is evaluated on a separate dataset to assess its performance, often using metrics such as accuracy, precision, recall, and F1-score.

Various algorithms can be used for classification, including Decision Tree, K-Nearest Neighbor (KNN), Naïve Bayes, Support Vector Machine (SVM), and ensemble methods like Random Forest and Gradient Boosting. Each algorithm has its own strengths and limitations, making the choice of method dependent on factors such as dataset size, feature types, noise levels, and computational requirements.

Classification has extensive applications across numerous fields, including medical diagnosis, fraud detection, customer segmentation, sentiment analysis, and election prediction. By leveraging classification techniques, analysts and researchers can transform raw data into actionable insights, supporting informed decision-making and strategic planning.

The classification method can be used to predict the results of legislative elections. In a study by Mafatikhu Habibi (2016), the Decision Tree method with the C4.5 algorithm was applied to predict the results of the legislative election for the Regional House of Representatives (DPRD) of Central Java Province. The prediction was based on entropy and gain calculations using data from 2014, which included 1,034 records from 10 electoral districts and 12 national political parties. The accuracy obtained, using a confusion matrix and attributes such as political party, candidate number, administrative city, valid votes per candidate, and valid votes per party, reached 94.68%. Another study by Mohammad Badrul (2015) concluded that the optimized model of K-Nearest Neighbor achieved an accuracy of 81.35% and an AUC value of 0.500, indicating an Excellent Classification diagnosis level. This data can later be used by parties requiring information related to general elections, particularly elections at the regional (Level I or Level II) level.

In the study conducted by Mohammad Badrul (2013), a cross-validation operator was used. In this research, Badrul employed the operator to evaluate a model using the Neural Network algorithm and a Neural Network algorithm based on Particle Swarm Optimization (PSO). In contrast, the study by Sofia Listyaningrum (2021) utilized the Split Validation operator to randomly divide the data into training and testing sets. In her research, Sofia applied enumeration and default split sampling. The study by Afrilio Franseda et al. (2020), which used an 80% : 20% split ratio, yielded an AUC value of 0.755, which falls under the category of fair classification. Meanwhile, the research by Eko et al. (2018) demonstrated that split validation using an 80% : 20% ratio resulted in higher accuracy compared to cross-validation. The highest accuracy achieved was 97.30%.

Based on the background of this study, the identified research problem is the absence of prior research that applies split sampling techniques to the Decision Tree and K-Nearest Neighbor algorithms in the context of the DKI Jakarta legislative election. From this problem identification, the scope of the study is defined as follows: The research utilizes data from the 2009 DKI Jakarta Legislative Election, which includes 2,268 candidates for the Regional House of Representatives (DPRD) from 44 political parties. The study compares the Decision Tree and K-Nearest Neighbor algorithms. The research

employs the data mining tool RapidMiner version 9.7. The performance variables used are Accuracy, Recall, and Precision. The study compares several types of split sampling operators: Linear Sampling, Shuffled Sampling, and Stratified Sampling. The research uses a data partition of 80% for training and 20% for testing. Based on the problem identification, the research question formulated in this study is: "How can split sampling techniques be applied to the Decision Tree and K-Nearest Neighbor algorithms in the context of the DKI Jakarta legislative election?"

## **2. RESEARCH METHOD**

The methodology consists of several stages: data collection, preprocessing, model building, evaluation, and comparison.

### **2.1. Data Collection**

The dataset used in this study is obtained from publicly available records relevant to the research domain—in this case, electoral data and associated attributes that may influence prediction outcomes. The dataset includes both numerical and categorical variables that serve as input features for the classification models.

### **2.2. Data Preprocessing**

Before model development, the raw dataset undergoes preprocessing to ensure quality and consistency. This step involves:

- **Data Cleaning:** Removing duplicate entries, handling missing values, and correcting inconsistencies.
- **Data Transformation:** Converting categorical attributes into numerical form using encoding techniques.
- **Normalization:** Scaling numerical values to a uniform range to prevent bias toward features with larger magnitudes.

### **2.3. Data Partitioning**

The dataset is split into training and testing subsets using three sampling techniques provided by RapidMiner:

- **Linear Sampling** – splits the data sequentially based on order.
- **Shuffled Sampling** – randomly shuffles data before splitting.
- **Stratified Sampling** – maintains the proportion of class labels in both training and testing sets. In this study, an 80:20 ratio is used, with 80% for training and 20% for testing.

### **2.4. Model Building in RapidMiner**

Two classification algorithms are implemented:

- **Decision Tree** – constructs a tree-based model to classify data based on feature splits.
- **K-Nearest Neighbor (KNN)** – classifies instances based on the majority class of the nearest neighbors in the feature space.

These models are configured and executed in RapidMiner's graphical process design interface, enabling reproducible workflows without extensive programming.

### **2.5. Model Evaluation**

The performance of each model is evaluated using three standard metrics:

- **Accuracy** – the proportion of correctly classified instances.
- **Precision** – the ratio of correctly predicted positive observations to the total predicted positives.
- **Recall** – the ratio of correctly predicted positive observations to all actual positives.

The evaluation is conducted for each sampling technique to determine its effect on model performance.

## 2.6. Comparative Analysis

Results from the different sampling methods and algorithms are compared to identify the optimal combination for achieving the highest predictive performance. This comparison is based on the evaluation metrics and supported by visualizations generated in RapidMiner.

Data mining is the process of discovering patterns or interesting information within selected data using specific techniques or methods. It is the process of analyzing data to identify patterns within a dataset. Data mining has the capability to analyze large datasets and transform them into meaningful patterns that are valuable for decision-makers (Hermawati, 2013). Data mining is a process that employs one or more machine learning techniques to analyze and automatically extract knowledge (Kusrini & Luthfi, 2009). It is an automatic process applied to existing data, particularly large datasets (Han & Kamber, 2006). According to Kusnawi (2007), the typical architecture of data mining consists of several main components, including: Database, data warehouse, or other information storage systems. Database or data warehouse server. Knowledge base. Data mining engine. Pattern evolution module. Graphical user interface. The stages of data mining require a systematic methodology, not only during the analysis phase but also when preparing the data and interpreting the results so they can become actionable insights or decisions. Therefore, data mining should be understood as a process with specific stages, including feedback from each stage to the previous one. Generally, the data mining process is interactive because it is not uncommon for the initial mining results to not meet the analyst's expectations, requiring a redesign of the process (Larose, 2005). Thus, data mining is a series of processes that can be divided into several stages. These stages are interactive, where the user is directly involved or interacts with the knowledge base.

General elections are one of the main pillars of a democracy. In a democratic country, elections serve as a crucial mechanism to select leaders who will represent the people in government, ranging from members of regional legislatures (DPRD Level II and Level I) to the national parliament (DPR RI) and the Regional Representative Council (DPD). According to Law No. 10 of 2008, Article 5 Paragraph 1, the electoral system used for legislative elections (DPR/DPRD) is a proportional representation system with an open list, while the DPD election is conducted through a multi-member district system. Based on the same law, participants in the elections for DPR, Provincial DPRD, and District/City DPRD are political parties, whereas candidates for the DPD are individuals (independent). Political parties participating in the election may nominate candidates up to 120 percent of the total number of contested seats in each electoral district. The nomination process must be democratic and transparent, and political parties are required to include at least 30% female representation in their list of candidates. Furthermore, political parties are mandated by law to submit a ranked list of candidates (numbered list) in order to compete for seats. The election system for choosing members of the DPR, Provincial DPRD, and District/City DPRD adheres to an open proportional representation model.

## 3. RESULTS AND DISCUSSION

The data obtained from the Jakarta General Elections Commission (KPUD Jakarta) is from the 2009 election, consisting of 2,268 records and comprising 10 variables or attributes. The variables used include: party number, party name, valid party votes, candidate number, gender, administrative city, electoral district, valid candidate votes, and number of seats obtained. The target variable is the election result. The algorithms applied in this study are Decision Tree and K-Nearest Neighbor (KNN), with a data split ratio of 80% for training and 20% for testing. The split sampling methods used are Linear Split Sampling, Shuffled Split Sampling, and Stratified Split Sampling. The performance metrics evaluated are Recall, Precision, and Accuracy.

The Decision Tree algorithm using Linear Split Sampling (LSS) achieved the highest Recall value among all split sampling types, reaching 100%. The Decision Tree algorithm using Shuffled Split Sampling (ShSS) produced the highest Precision value at 85.71%, outperforming other sampling types. Both Linear Split Sampling (LSS) and Shuffled Split Sampling (ShSS) achieved the same and highest Accuracy value of 98.46%, which is higher than that of Stratified Split Sampling (StSS) at 97.80%. The K-Nearest Neighbor (KNN) algorithm using Linear Split Sampling (LSS) achieved the highest Recall value among all sampling types, reaching 93.75%. The KNN algorithm using Stratified

Split Sampling (StSS) produced the highest Precision value at 85.71%, compared to other sampling types. Both Linear Split Sampling (LSS) and Shuffled Split Sampling (ShSS) achieved the same and highest Accuracy value of 97.36%, which is higher than that of Stratified Split Sampling (StSS) at 96.92%.

Table 1. Research Results tabel

Algorithm		Decision Tree			K-Nearest Neighbor (KNN)		
Enumerasi Split Sampling		80% : 20%			80% : 20%		
		LSS	ShSS	StSS	LSS	ShSS	StSS
Per	rec	100%	81.82%	63.16%	93.75%	59.09%	31.58%
for	prec	82.05%	85.71%	80%	75%	81.25%	85.71%
ma	Acc	98.46%	98.46%	97.80%	97.36%	97.36%	96.92%
nce							

LSS = Linear Split Sampling

ShSS = Shuffled Split Sampling

StSS = Stratified Split Sampling

Rec = Recall

Prec = Precision Acc = Accuracy

Overall, the Linear Split Sampling method outperforms the other sampling types, as it achieves the highest Recall and Accuracy values, with only a slight difference in Precision compared to the other methods.

#### 4. CONCLUSION

Based on the research conducted, it can be concluded that, overall, Linear Split Sampling performs better than Shuffled Split Sampling and Stratified Split Sampling. This is because the performance metrics of Linear Split Sampling are superior when applied to both the Decision Tree and K-Nearest Neighbor (KNN) algorithms compared to the other sampling types. In the implementation of the Decision Tree algorithm, the model's performance demonstrated consistent yet slightly varied results depending on the data splitting method used. With the Linear Split Sampling method, the algorithm achieved a perfect Recall of 100%, a Precision of 82.05%, and an Accuracy of 98.46%. These results indicate the model's excellent ability to correctly identify all positive cases, although there is still room for improvement in terms of precision. Using the Shuffled Split Sampling method, the model's Recall decreased to 81.82%, while Precision improved to 85.71%. Despite this change, the Accuracy remained stable at 98.46%. This suggests that the model became more selective in identifying positive cases, which resulted in fewer false positives but also a reduced detection of actual positives. On the other hand, with Stratified Split Sampling, the model once again achieved a perfect Recall of 100%, the same Precision of 82.05%, and a slightly lower Accuracy of 97.80%. In contrast, the K-Nearest Neighbor (KNN) algorithm exhibited more fluctuating performance across the three data splitting methods. With Linear Split Sampling, KNN achieved a Recall of 93.75%, Precision of 75%, and Accuracy of 97.36%. This indicates a reasonably good ability to detect positive cases, though the lower precision suggests a higher rate of false positives. Under the Shuffled Split Sampling method, the Recall dropped significantly to 59.09%, while Precision increased to 81.25%, with Accuracy remaining consistent at 97.36%. This reflects a trade-off where the model made fewer incorrect positive predictions but missed a substantial number of actual positive cases. The lowest performance was observed with Stratified Split Sampling, where Recall further declined to 31.58%, although Precision remained relatively high at 85.71%, and Accuracy slightly decreased to 96.92%. Overall, the Decision Tree algorithm demonstrated more stable and balanced performance across different data splitting techniques compared to the KNN algorithm, particularly in its ability to detect positive cases effectively (as reflected by higher Recall values). These results support the conclusion that Linear Split Sampling provides the most consistent and optimal performance for legislative election prediction in DKI Jakarta using Decision Tree and KNN algorithms.

Despite the promising results, this study is not without its limitations. Several aspects could be improved or explored further in future research to enhance the robustness and applicability of the findings. One potential direction is to incorporate a wider range of data mining classification methods. By comparing the performance of various algorithms, researchers can gain deeper insights and identify the most effective approach for similar datasets or problems.

Additionally, the dataset used in this study was relatively limited in size. Future research would benefit from expanding the dataset to include a larger and more diverse sample, which could improve the reliability and generalizability of the results.

Another promising avenue for future work is the development of a predictive application for election outcomes. By leveraging the algorithms and methods applied in this study, such an application could offer practical value in forecasting election results based on relevant data inputs.

## REFERENCES

- Faid, M., Syahputra, M., & Wahyudi, H. (2019). Perbandingan kinerja tool data mining Weka dan RapidMiner dalam algoritma klasifikasi. Universitas Nurul Jadid.
- Franseda, A., Subiyanto, & Sari, R. F. (2020). Integrasi metode Decision Tree dan SMOTE untuk klasifikasi data kecelakaan lalu lintas. *Jurnal Sistem dan Teknologi Informasi*, 8.
- Habibi, M. (2016). Prediksi hasil pemilihan umum legislatif DPRD Provinsi Jawa Tengah menggunakan metode decision tree dan algoritma C4.5 [Undergraduate thesis, Universitas Dian Nuswantoro Semarang].
- Han, J., & Kamber, M. (2006). *Data mining: Concepts and techniques* (2nd ed.). Morgan Kaufmann.
- Han, J., Kamber, M., & Pei, J. (2011). *Data mining: Concepts and techniques* (3rd ed.). Elsevier.
- Hermawati, F. A. (2013). *Data mining* (Edisi 1). Andi Offset.
- Kusnawi. (2007, November 24). Pengantar solusi data mining. Paper presented at Seminar Nasional Teknologi STMIK AMIKOM Yogyakarta, Yogyakarta, Indonesia.
- Kusrini, & Luthfi, E. T. (2009). *Algoritma data mining* (Edisi 1). Andi Offset.
- Larose, D. T. (2005). *Discovering knowledge in data: An introduction to data mining* (2nd ed.). Wiley Interscience.
- Listyaningrum, S. (2021). Penerapan data mining untuk analisis karakteristik DPT non-participate sebagai prediksi partisipan pemilu dengan menggunakan metode Naive Bayes Classifier [Undergraduate thesis, Universitas Dian Nuswantoro].
- Nanfack, O., Gaur, M., Almeida, F., Namoun, A., & Purohit, H. (2023). Decision trees: From efficient prediction to responsible AI. *Patterns*, 4(5), 100740. <https://doi.org/10.1016/j.patter.2023.100740>
- Nanfack, W., Gaur, M., Almeida, F., Namoun, A., & Purohit, H. (2023). Decision trees: From efficient prediction to responsible AI. *Patterns*, 4(5), 100740. <https://doi.org/10.1016/j.patter.2023.100740>
- RapidMiner Documentation. (n.d.). Split Data. Retrieved from [https://docs.rapidminer.com/latest/studio/operators/blending/examples/sampling/split\\_data.html](https://docs.rapidminer.com/latest/studio/operators/blending/examples/sampling/split_data.html)
- Republik Indonesia. (2008). Undang-undang No. 10 Tahun 2008 tentang Pemilihan Umum Anggota DPR, DPD, dan DPRD Pasal 5 Ayat 2. Lembaran Negara RI Tahun 2008. Sekretariat Negara.
- Republik Indonesia. (2008). Undang-undang No. 10 Tahun 2008 tentang Pemilihan Umum Anggota DPR, DPD, dan DPRD Pasal 12. Lembaran Negara RI Tahun 2008. Sekretariat Negara.
- Vitalaya, N. A. R. (2024). Perbandingan tipe sampling pada klasifikasi minat TIK [Master's thesis, UIN Syarif Hidayatullah Jakarta].

**Halaman Ini Dikosongkan**