

Classification Analysis of Waiting Period for Telkom University Alumni to Get Jobs Using Decision Tree and Support Vector Machine

Annisa Miranda¹, Kemas Muslim Lhaksana^{2,*}

School of Computing, Telkom University, Bandung, Indonesia

Email: ¹annisamiranda@student.telkomuniversity.ac.id, ^{2,*}kemasmuslim@telkomuniversity.ac.id

Email Penulis Korespondensi: kemasmuslim@telkomuniversity.ac.id

Submitted: 25/07/2022; Accepted: 18/08/2022; Published: 30/09/2022

Abstract— Tracer analysis is one of the ways to increase a university's accreditation. Tracer studies, also known as graduate surveys, are beneficial for enhancing learning and developing university curricula. The period it takes graduates to secure employment is a measure of their quality. The sooner graduates obtain a job, the higher their perceived quality. Conversely, if it takes graduates longer to find employment, their quality is deemed lower. To gain new knowledge from the tracer study dataset regarding the relationship between university contribution and alumni capability in the job market, in this study, data mining techniques are used to determine what factors influence the length of time it takes college graduates to find employment. This classification model contains a total of 2288 data instances from the categorical type of dataset. The features are selected using chi-square. Two classification algorithms, Decision Tree and Support Vector Machine, are compared for the best model. This study also used hyperparameter tuning to improve accuracy. The results show decision tree produces higher accuracy compared to the support vector machine. The accuracy obtained from the decision tree model is 55.02% and increased to 65.06% after hyperparameter tuning. Meanwhile, the support vector machine brought an accuracy of 60.40% and increased to 62.15% after hyperparameter tuning. Factors that affect the classification of the alumni waiting period in getting a job in this study are sex, faculty of the study field, department of the study field, study period, company specification, company category, and work location.

Key Words: Tracer Study; Classification; Chi2; Decision Tree; Support Vector Machine

1. INTRODUCTION

University is a final-level educational institution that is able to prepare and create graduates who can compete and are ready to enter the world of work. Fierce competition in getting work cannot be avoided, and as a result, challenges arise in the field of education that must be faced by universities. The ability of graduates to be able to compete and get jobs is a measure for universities in producing good quality graduates.

The quality of graduates is measured by the period it takes to get a job. The work obtained by college graduates measures success in tertiary education in the fields of education. Every college has a different system for creating the quality of graduates. The Director General of Learning and Student Affairs requires every university to conduct a Tracer Study every year to fulfill the needs of accreditation data and improvement of the learning process and curriculum development in the university [1].

Tracer Study is a research on the situation of alumni in job search, work situations, and the use of competencies obtained during their education at university. Universities in developed countries evaluate themselves by conducting tracer studies. Graduate Tracer Studies (GTS) are crucial to Higher Education Institutions because they allow them to adjust to changes in society, particularly the expectations of current and future employers, thorough evaluation, and continuous revision of their curriculum [2]. One of the ways to get information, knowledge, and patterns from tracer study is data mining.

Data mining is the process of extracting patterns from data using specialized algorithms. As information can be acquired from the large quantity of data created and gathered nowadays, this branch of data science helps the decision-making process in an information system [3]. One of the analytical methods in data mining is classification. Classification is a machine learning technique that categorizes data into a certain number of classes by predicting class labels or categories for new data [4]. Classification is included in supervised learning, which can estimate the mapping function of an input variable and predict its output variable.

Our research on the waiting period classification of Telkom University Alumni to get a job is based on several related studies. One of the former studies used one of the decision tree algorithms, the C4.5 algorithm. The study results indicate that the C4.5 model achieves the highest performance when using forward selection features. The accuracy obtained in this study is 80.37% and 79.56% of precision. This study also shows that some features have a more significant influence on classification waiting time, such as the city of work, class year, graduation year, GPA, organizational history, expertise certificate, and type of company they work for [5].

Other study was conducted using various machine learning methods such as Decision Tree, Support Vector Machine, and Neural Network. The study results indicate that the Support Vector Machine produces the highest accuracy compared to other models by 66.097%. The author decided to do parameter tuning for the Decision Tree model because of its short time to build. From 66.0651%, the accuracy of the decision tree model can be increased to 66.1824% after parameter tuning. Age, industrial internship, and faculty contain the most information and have the most effects on the final class out of all the criteria considered [6].

Another study uses Naïve Bayes algorithm to classify alumni waiting period to get a job. The main goal of this research is to make a prediction model for Budi Luhur Secretary Academy alumni's waiting periods when getting their



first job. This study produces 90.90% accuracy level by using 199 training data and 22 testing data. One of the features used in this research besides waiting time to get a job is GPA or index performance for every semester [7].

Similar studies on alumni waiting period classification using Naïve Bayes algorithm divided the sample data into several combinations with a composition of 9/10 as training data and 1/9 as testing data. The results show that from 1240 test data, only 603 data can be properly classified with an accuracy rate of 48.629% [8]. Several previous studies on the classification of the alumni waiting period in getting the first job after graduation has shown that it is necessary to use machine learning to build a classification model that predicts valuable factors that might be helpful to classifying alumni fluency in getting a job.

This study was conducted to build a classification model system using Decision Tree and Support Vector Machine. The model aims to classify the data of Telkom University alumni's waiting period to get a job by analyzing their performance throughout their active study years. This research was also conducted to analyze the effect of using several combinations of sample data. Due to the categorical type of the dataset, the Chi-Square Test was used to select features for this study. This study also implements hyperparameter tuning to improve accuracy. In a former study [6], hyperparameter tuning was proven to improve the accuracy of the decision tree. In addition, we also analyze the factors that may affect alumni fluency in getting a job.

2. RESEARCH METHODOLOGY

2.1 Research Flow

In this study, the development of waiting time classification system for alumni seeking employment based on the tracer study dataset using Decision Tree and Support Vector Machine was carried out in several stages, as illustrated by the flowchart diagram in Figure 1.

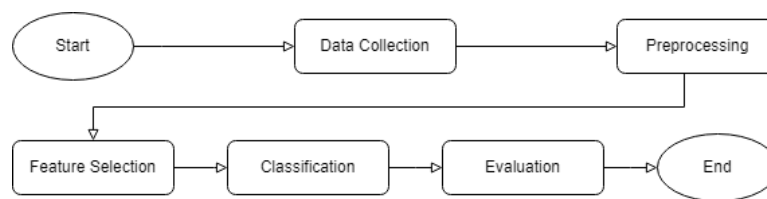


Figure 1. Research Process

The system model is made from raw data that goes through preprocessing, then the data is divided into training data and test data. Furthermore, the data goes through the feature selection stage. The selected features then go through the classification stage using a decision tree and support vector machine. The results of the classification are then evaluated to see the performance of each model built.

2.1.1 Data Collection

The data used to develop the alumni waiting period classification system to get a job comes from a survey of Telkom University alumni between 2015 and 2019. Table 1 provides a more detailed explanation of the data attributes.

Table 1. Tracer Study Dataset Description

Attributes	Description	Attributes	Description
student_id	Student's identification number	computer_skills	Computer skills
Sex	Sex	critical_thinking	Critical thinking skills
department	Departement of study field	research_skills	Research skills
faculty	Faculty of study field	learning_skills	Learning skills
class_year	Entrance year	communication_skills	Communication skills
graduate_year	Graduation year	pressure_tolerance_skills	Working under pressure
study_period	Student study period	time_management_skills	Time management skills
organization_exp	Organizational experience	self_employed_skills	Work independently
gpa	Graduate Point Average	team_work_skills	Team work skills
job_hunt_time	Times taken to get a job	problem_solving_skills	Problem solving skills
job_search	Way to get a job	negotiation_skills	Negotiation skills
company_type	Type of company	analytical_skills	Analytical skills
company_specification	Specification of company	tolerance_skills	Tolerance ability
company_category	Category of company	adaptability_skills	Adaptability skills
work_location	Work location	loyalty_integrity_skills	Loyalty and integrity
study_work_fit	Suitability of study and work	work_diverse_culture	Different culture adapt
applications_submitted	Number of applications sent	leadership_skills	Leadership skills

applications_received	Applications' responses number	responsibility_skills	Responsibility skills
education_level	Level education's background	initiative_skills	Conducting initiative skills
inside_study_competence	Expertise in study field	project_management_skills	Project management skills
outside_study_competence	Expertise outside study field	presenting_idea_skills	Presenting idea skills
general_knowledge	General knowledge	writing_skills	Writing document skills
internet_skills	Internet skills		

2.1.2 Preprocessing

There are two steps in this preprocessing, including data cleaning and data transformation. Several ways to clean the data in data cleaning include removing duplicate records, data imputation for missing values using the mode for numerical and categorical data, and dropping outliers. In data transformation, categorical attributes with non-ordinal associations are encoded numerically using the one-hot encoder, whereas categorical attributes with ordinal associations are encoded numerically using the label encoder.

Label encoder is part of Python's Scikit-learn library and converts categorical data into a computer-readable form so that machine learning models can understand feature correlations. One-hot encoding is a method for avoiding the misinterpretation of correlations between independent variables by emphasizing the presence of feature variables. Each categorical value is transformed into a new column, and the label values are converted to a digital form (1 or 0) [9]. The target column is obtained from the results of the classification of the alumni waiting period in getting a job. The dataset target will be classified into two classes, namely alumni who go well in getting a job and alumni who do not go well in getting a job. The average waiting period for alumni to get a job is two months, obtained from the job_hunt_time attribute mode. Alumni who get a job for ≤ 2 months will be classified as alumni who go well in getting a job, whereas alumni who get a job for > 2 months will be classified as alumni who do not go well in getting a job.

2.1.3 Feature Selection

In the building model, feature selection is used to discrete irrelevant attributes. It enables selecting the best and most valuable attributes when creating a model. Classification algorithms that use related attributes improve prediction accuracy, reduce research time, and simplify concepts [6]. The feature selection method used in this study is Chi-Square Test because the chi-square test works on categorical variables. The chi-square test compares the observed and predicted distributions. In this study, the p-value will be used to find the relationship between attributes in the dataset. P-value represents the probability that the observed value deviates from the expected value. Calculating the chi-square test score requires testing both the Null Hypothesis and the Alternate Hypothesis. This feature selection is based on the following hypothesis [10]:

$$H_0 : X_{obs} = X_{exp} \tag{1}$$

$$H_1 : X_{obs} > X_{exp} \text{ OR } X_{obs} < X_{exp} \tag{2}$$

Where H_0 is the null hypothesis, H_1 is the alternate hypothesis, X_{obs} is observed data, and X_{exp} is expected data. If the Null hypothesis is accepted, which means that the p-value is higher than alpha (α) that is 0.05, the feature will be removed. In meantime, the feature will pass the selection and be used in the model if the p-value is lower than 0.05.

2.1.4 Classification

Decision Tree and Support Vector Machine is used to build the classification model in classifying alumni waiting period to get a job, whether they go well in getting a job or don't go well in getting a job. Several classification procedures are performed, consisting of data train and data test distribution, model creation, model training, and prediction of each class on each model created. The more specific classification processes are illustrated in figure 2.

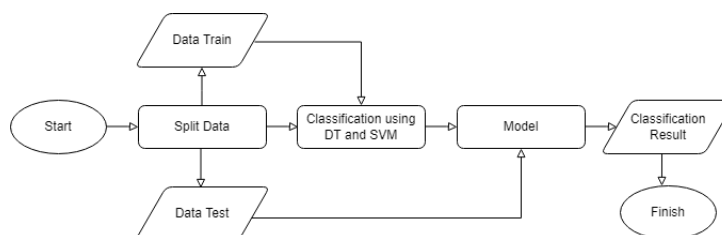


Figure 2. Classification Process

The decision tree is used to train algorithms by using choices such as leaves and branches. Regression and classification tasks with continuous and categorical values can be performed using a decision tree [11]. In 1972, Messenger and Mandell proposed the first classification tree algorithm, THeta Automatic Interaction Detection or THAID. Several decision tree algorithms, including ID3, C4.5, C5.0, and CART. In this study, we proposed the

Classification and Regression Tree (CART). Leo Breiman et al.(1984) introduced the CART algorithm, which refers to decision tree algorithms for classification and regression modeling. The CART algorithm generates binary trees, also known as Hierarchical Optimal Discriminate Analysis (HODA). The term binary implies that a decision tree node can only be divided into two groups. For attribute selection, CART uses the Gini index as a measure of impurity. The node's records are split based on the attribute with the highest reduction in impurity [12]. Equation 2 shows how the Gini index is defined.

$$\text{Gini}(D) = 1 - \sum_{j=1}^n p_j^2 \quad (3)$$

Where D is the data set containing n samples, that p_j is the relative probability that the sample of category j appears in D. The Gini index separates most categories from other nodes. Smaller gini values indicate a more uneven sample category distribution. If the subset generated by the splitting point has higher category purity, it's easier to distinguish between categories [13].

Support Vector Machine is a classifier-building approach that uses hyperplanes in n-dimensional space to separate various class occurrences [14]. The support vector machine was first introduced by Vladimir Vapnik and Alexey Chervonenkis in 1963. The kernel is a technique in machine learning that is used to tackle non-linear problems using linear classifiers and is engaged in converting linear non-separable input into linearly separable data [15]. In particular, support vector machine algorithms with universal kernels that can approximate any continuous function with any level of accuracy are universally consistent because they use a method for minimizing structural risk [16]. Several support vector machine kernels include linear, Radial Basic Function (RBF), and polynomial. In this study, we proposed an RBF kernel. Equation 3 shows how the RBF kernel is defined.

$$K(x, y) = \exp\left(\frac{-||x - y||^2}{2\sigma^2}\right) \quad (4)$$

Where σ is the variance and the parameter and $-||x - y||^2$ is the squared euclidean distance between two points x and y.

2.1.5 Evaluation

As part of this study, we determined the performance of the system by evaluating its level of accuracy and precision. The number of classes correctly classified is represented by accuracy. The accuracy obtained of True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) from the confusion matrix. The confusion matrix is a table commonly used to describe the performance of a classification model on a set of test data with known values. Each row in a confusion matrix represents an actual class, while each column represents a predicted class [17].

Precision describes the degree of congruence between the requested data and the model's prediction results. The confusion matrix is a square matrix that depicts the actual and predicted class [18]. Recall describes the model's ability to retrieve information. F1-Score is a weighted comparison of the average precision and recall used to evaluate the performance of classifiers. Accuracy, precision, recall, and f1-score are defined in the following equation.

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN} \quad (5)$$

$$\text{Precision} = \frac{TP}{TP+FP} \quad (6)$$

$$\text{Recall} = \frac{TP}{TP+FN} \quad (7)$$

$$F1 - \text{Score} = \frac{2(\text{precision} * \text{recall})}{\text{precision} + \text{recall}} \quad (8)$$

3. RESULT AND DISCUSSION

3.1 Dataset Information

Several attributes have missing values, according to Table 1. All attributes with missing values are replaced in accordance with the condition in the research method except gpa and job_hunt_time. Missing values in gpa and job_hunt_time will be dropped to avoid biased evaluation. The problem of missing values in attributes other than gpa and job_hunt_time will be fixed manually filling in the values using the mode that corresponds to each attribute in question. Gpa will be converted into a categorical attribute by categorizing the value of gpa. Job_hunt_time will be transformed into the target label of the data set.

Handling outliers on the job_hunt_time is necessary to maintain the distribution of the data. The ranges obtained after handling outliers in job_hunt_time are <-10 to >20. The results of feature selection using chi-square

revealed several attributes that are related to or dependent on the target class after being encoded by a one-hot encoder. More explanations are in table 2.

Table 2. Dependent Attributes

Attributes	P-Value	Attributes	P-Value
company_specification_consultant	0.0415	company_category_local	0.0074
department_accountancy	0.0378	study_period_6yrs	0.0074
department_communication	0.0318	sex_female	0.0071
faculty_ict_business	0.0317	department_business_management	0.0062
department_electro_engineering	0.0250	work_location_kepri	0.0060
company_specification_mining	0.0250	department_computer_engineering	0.0038
job_search	0.0228	company_category_multinational	0.0029
class_year_2012	0.0177	company_specification_digital	0.0025
organization_art	0.0174	faculty_business_communication	0.0015
department_business_adm	0.0173	faculty_informatics	0.0002
work_location_bali	0.0139	faculty_business_economy	0.0002
study_period_2yrs	0.0134	department_informatics	0.0000
class_year_2013	0.0129		

Handling outliers on the job_hunt_time is necessary to maintain the distribution of the data. The ranges obtained after handling outliers in job_hunt_time are <-10 to >20. The results of feature selection using chi-square revealed several attributes that are related to or dependent on the target class after being encoded by a one-hot encoder. More explanations are in table 2.

3.2 Baseline

The first scenario is to divide the data set into different combinations of train data and test data. Data splitting combinations include 70:30 proportion, 80:20 proportion, and 90:10 proportion with combinations of features selected using chi-square into 80 features, 90 features, and 100 features. These combinations of split data and the number of features use the decision tree and support vector machine with both use default parameters which are no maximum depth for the decision tree and cost equal to 1.0 for support vector machine.

Figure 3 illustrates the outcome of a 70:30 proportion of data split and 80 features, 90 features, and 100 features. The decision tree produces 55.02% accuracy and 56.76% precision using 80 features, 54.14% accuracy, and 57.65% precision using 90 features, and 54.00% accuracy and 57.79% precision using 100 features. Support vector machine produces 60.40% accuracy and 61.63% precision using 80 features, 59.24% accuracy, and 60.77% precision using 90 features, and 59.24% accuracy and 61.08% precision using 100 features.

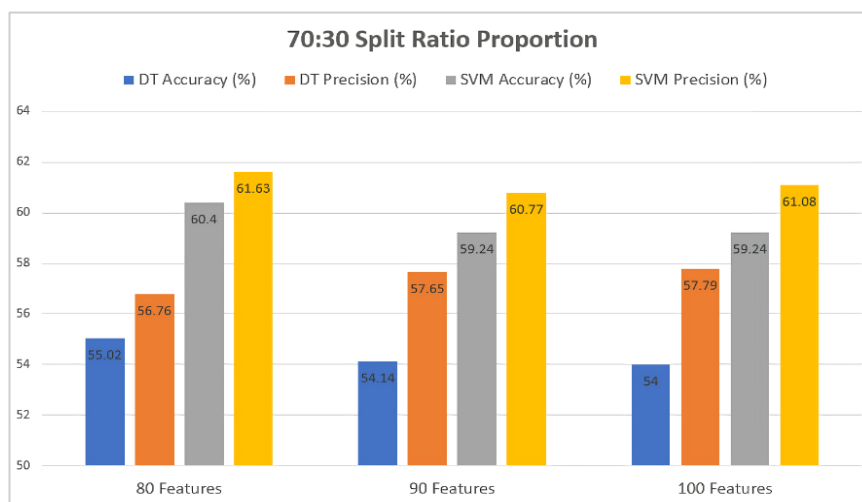


Figure 3. 70:30 Data Split Proportion Result

Figure 4 illustrates the outcome of an 80:20 proportion of data split and 80 features, 90 features, and 100 features. The decision tree produces 53.93% accuracy and 57.87% precision using 80 features, 53.93% accuracy, and 56.85% precision using 90 features, and 53.93% accuracy and 57.55% precision using 100 features. Support vector machine produces 57.64% accuracy and 57.42% precision using 80 features, 57.64% accuracy, and 57.38% precision using 90 features, and 57.64% accuracy and 57.18% precision using 100 features.

Figure 5 illustrates the outcome of a 90:10 proportion of data split and 80 features, 90 features, and 100 features. The decision tree produces 55.02% accuracy and 57.59% precision using 80 features, 48.47% accuracy, and 53.19% precision using 90 features, and 51.76% accuracy and 56.14% precision using 100 features. Support vector machine



produces 56.76% accuracy and 55.67% precision using 80 features, 56.33% accuracy, and 55.67% precision using 90 features, and 56.76% accuracy and 55.80% precision using 100 features.

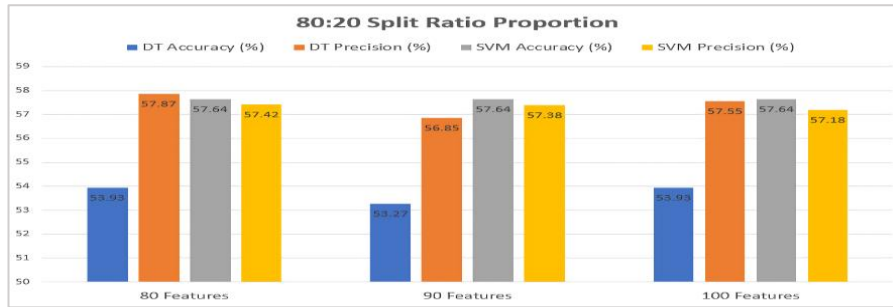


Figure 4. 80:20 Data Split Proportion Result

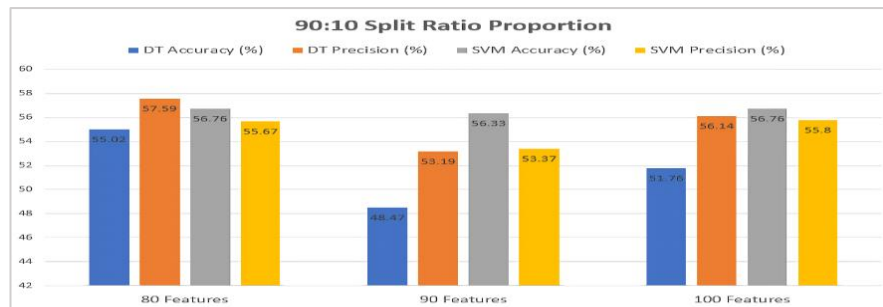


Figure 5. 90:10 Data Split Proportion Result

Based on the results, the decision tree model generated the best accuracy and precision in a 90:10 proportion of data split and 80 features with 55.02% of accuracy and 57.09% of precision. Support vector machine model generated the best accuracy and precision in 70:30 proportion of data split and 80 features with 60.04% of accuracy and 61.63% of precision. These results will be the baseline for the next test scenario.

3.3 Hyperparameter Tuning

The second test scenario involves hyperparameter tuning by varying the maximum depth for the decision tree model and cost for the support vector machine. Hyperparameter tuning for both models is to determine whether the maximum depth and cost value have a significant impact on performances. In the decision tree, values of the maximum depth to be used in the decision tree is with a range of 1 to 30. Whereas the values of the cost to be used in the support vector machine are in a range of 1.0 to 2.9. The result of the decision tree hyperparameter tuning can be seen in Figure 6 and Figure 7.

Figure 6 depicts the results of decision tree accuracy using various values of the maximum parameter depth of the tree against the data train and data test. As can be seen, increasing the maximum depth values improves the accuracy of train data but reduces the accuracy of test data. Starting from the maximum depth value of 6, the accuracy gap between the train data and test data is wider and shrinks back when the maximum depth value is 8. Then the distance appears again at the maximum depth is 9. This result is known as overfitting, in which the model predicts the train data nearly perfectly but fails to predict the test data. We performed pre-pruning by halting the tree's growth in advance. We determined that the optimal maximum depth value is 8. The accuracy of the train data is 67.21%, and the test data is 65.06% based on this value. Hyperparameter tuning improved the accuracy by 10.04% compared to the baseline. Next is a visualization of support vector machine results using hyperparameter tuning. A more detailed explanation can be seen in Figure 7.

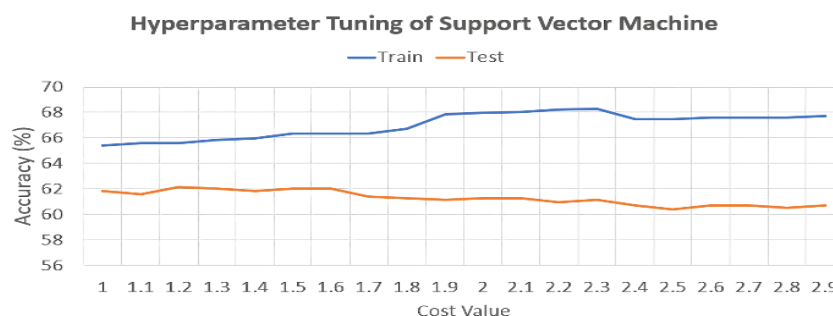


Figure 6. Decision Tree Hyperparameter Tuning Result

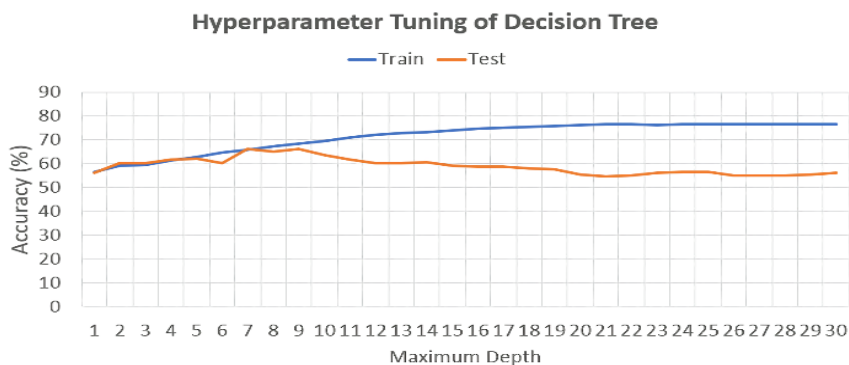


Figure 7. Support vector Machine Hyperparameter Tuning Result

Figure 7 depicts the results of support vector machine accuracy using various values of the parameter of the cost against the train data and test data. As can be seen, a big gap occurs between the accuracy of the train data and test data when the cost is 1.1 and 1.3 to 2.9. The accuracy of the train data continues to increase when the accuracy of the test data decreases as the cost value increases. The smallest gap between the accuracy of train data and test data is when the cost is 1.2. Using hyperparameter tuning, the accuracy was boosted by 1.75% compared to the baseline. Based on these results, we decided that the best cost value as a parameter in the support vector machine model that was built was 1.2.

3.4 Model Comparison

The third scenario compares the accuracy results of the decision tree model and support vector machine model. Table 2 displays the comparative results.

Table 3. The Comparison of Decision Tree and Support Vector Machine Results

Model	Accuracy (%)	Precision (%)
Decision Tree (Baseline)	55.02	57.59
Decision Tree (Hyperparameter Tuning)	65.06	64.24
Support Vector Machine (Baseline)	60.40	61.63
Support Vector Machine (Hyperparameter Tuning)	62.15	63.38

As seen from the table above, the decision tree model in the baseline produces a lower accuracy level than the support vector machine in the baseline. But, with hyperparameter tuning, the decision tree model can produce a higher accuracy level than the support vector machine model, with an accuracy level of 65.06%. A previous accuracy value of 55.02% at the decision tree baseline increased by 10.04% after the use of hyperparameter tuning. Although not statistically significant, the accuracy of the support vector machine increased by 1.75%, from 60.40% at baseline to 62.15% after hyperparameter tuning. Because the attributes in the dataset used are categorical type, these results are supported by the results of research by Kumari, Kumar, Prakash, and Divya that decision trees outperform support vector machines for categorical and collinear information [19].

4. CONCLUSION

Considering the results and discussion, classification modeling using a decision tree and support vector machine is proven to be able to classify the waiting period of alumni in getting a job. The combination of split data and the number of features are also proven to have a great impact on classifying data. Combination of 90:10 proportion of data split and 80 features produces highest accuracy in baseline scenario using the decision tree, whereas 70:30 proportion of data split and 80 features produces highest accuracy in baseline using the support vector machine. The result showed that the decision tree performed better compared to the support vector machine with 65.06% over 62.15%. Using hyperparameter tuning is proven to improve the accuracy, evidenced by an increase in accuracy of 10.04% in the decision tree and 1.75% in the support vector machine. Tracer study dataset has also proven usable to classify the fluency of alumni waiting period in getting a job after graduation. This is due to the fact that the dataset utilized is a categorical data type, performing chi-square test in feature selection is proven to be able to find factors that affect the fluency of alumni in getting a job. P-value is used to measure the association between attributes and the target column. A significance level of $p < 0.01$ is considered to be powerful evidence against the null hypothesis. According to table 2, some of the most influential factors in classifying alumni waiting period to get a job are alumni from business and economy faculty, alumni from informatics faculty, alumni from business and communication faculty, alumni from the informatics department, alumni from the computer engineering department, alumni from business management department, alumni with female gender, alumni with six years of the study period, alumni who get a job in Kepulauan

Riau, alumni who get a job in a company with digital specification, and alumni who get a job in a company with local and multinational category. For further research, utilizing tracer study datasets with more employability-related characteristics for students is a feasible option.

REFERENCES

- [1] “Pelaksanaan Tracer Study di Perguruan Tinggi - Website LLDIKTI Wilayah V.” <https://lldikti5.kemdikbud.go.id/home/detailpost/pelaksanaan-tracer-study-di-perguruan-tinggi> (accessed Jul. 12, 2022).
- [2] A. C. Albina and L. P. Sumagaysay, “Employability tracer study of Information Technology Education graduates from a state university in the Philippines,” *Social Sciences & Humanities Open*, vol. 2, no. 1, p. 100055, Jan. 2020, doi: 10.1016/J.SSAHO.2020.100055.
- [3] P. M. Seeger, Z. Yahouni, and G. Alpan, “Literature review on using data mining in production planning and scheduling within the context of cyber physical systems,” *J Ind Inf Integr*, vol. 28, p. 100371, Jul. 2022, doi: 10.1016/J.JII.2022.100371.
- [4] S. Priya, R. Priyatharshini, R. Shruthi, V. Pooja, and R. S. Swarna, “Early detection of Parkinson’s disease using data mining techniques from multimodal clinical data,” *Advanced Machine Vision Paradigms for Medical Image Analysis*, pp. 213–228, Jan. 2021, doi: 10.1016/B978-0-12-819295-5.00008-1.
- [5] F. Rezkika, B. N. Sari, and A. Susilo, “Klasifikasi Masa Tunggu Alumni Untuk Mendapatkan Pekerjaan Menggunakan Algoritma C4 . 5,” *Progresif: Jurnal Ilmiah Komputer*, vol. Vol. 17, pp. 95–106, 2021.
- [6] Z. Othman, S. W. Shan, I. Yusoff, and C. P. Kee, “Classification techniques for predicting graduate employability,” *International Journal on Advanced Science, Engineering and Information Technology*, vol. 8, no. 4–2, pp. 1712–1720, 2018, doi: 10.18517/ijaseit.8.4-2.6832.
- [7] R. Amalia and A. Wibowo, “Prediction of the Waiting Time Period for Getting a Job Using the Naive Bayes Algorithm,” *Research Journal of Advanced Engineering and irjaes.com*, 2020. [Online]. Available: <http://irjaes.com/wp-content/uploads/2020/10/IRJAES-V5N2P219Y20.pdf>
- [8] I. M. B. Adnyana, “Implementasi Naïve Bayes Untuk Memprediksi Waktu Tunggu Alumni Dalam Memperoleh Pekerjaan,” *Seminar Nasional Teknologi Komputer & Sains (SAINTEKS)*, vol. 1, no. 1, pp. 131–134, 2020, [Online]. Available: <http://prosiding.seminar-id.com/index.php/sainteks/article/view/418>
- [9] T. Al-shehari and R. A. Alsowail, “An Insider Data Leakage Detection Using One-Hot Encoding, Synthetic Minority Oversampling and Machine Learning Techniques,” *Entropy* 2021, Vol. 23, Page 1258, vol. 23, no. 10, p. 1258, Sep. 2021, doi: 10.3390/E23101258.
- [10] K. F. Weaver, “An introduction to statistical analysis in research : with applications in the biological and life sciences,” p. 594.
- [11] N. Peker and C. Kubat, “Application of Chi-square discretization algorithms to ensemble classification methods,” *Expert Systems with Applications*, vol. 185, p. 115540, Dec. 2021, doi: 10.1016/J.ESWA.2021.115540.
- [12] Information Resources Management Association, “Cognitive analytics : concepts, methodologies, tools, and applications”.
- [13] C. L. Lin and C. L. Fan, “Evaluation of CART, CHAID, and QUEST algorithms: a case study of construction defects in Taiwan,” *Journal of Asian Architecture and Building Engineering*, vol. 18, no. 6, pp. 539–553, Nov. 2019, doi: 10.1080/13467581.2019.1696203.
- [14] V. N. Gudivada, “Data Analytics: Fundamentals,” *Data Analytics for Intelligent Transportation Systems*, pp. 31–67, Jan. 2017, doi: 10.1016/B978-0-12-809715-1.00002-X.
- [15] G. Battineni, N. Chintalapudi, and F. Amenta, “Machine learning in medicine: Performance calculation of dementia prediction by support vector machines (SVM),” *Informatics in Medicine Unlocked*, vol. 16, p. 100200, Jan. 2019, doi: 10.1016/J.IMU.2019.100200.
- [16] V. Vapnik and R. Izmailov, “Knowledge transfer in SVM and neural networks,” *Annals of Mathematics and Artificial Intelligence* 2017 81:1, vol. 81, no. 1, pp. 3–19, Feb. 2017, doi: 10.1007/S10472-017-9538-X.
- [17] S. Smys, R. Bešťák, R. Palanisamy, and I. Kotuliak, “Computer networks and inventive communication technologies : proceedings of Fourth ICCNCT 2021”.
- [18] O. Caelen, “A Bayesian interpretation of the confusion matrix,” *Annals of Mathematics and Artificial Intelligence*, vol. 81, no. 3–4, pp. 429–450, Dec. 2017, doi: 10.1007/s10472-017-9564-8.
- [19] A. D. Kumari, J. P. Kumar, V. S. Prakash, and K. S. Divya, “Supervised Learning Algorithms : A Comparison,” vol. 1, no. 1, pp. 1–12, 2020.