



## Analisis Prediksi Pendapatan Penduduk dengan Metode K-Nearest Neighbor, Decision Tree, Naive Bayes, Ensemble Methods, dan Linear Regression

Eri Mardiani<sup>1</sup>, Nur Rahmansyah<sup>2</sup>, Endah Tri Esti Handayani<sup>3</sup>, Sari Ningsih<sup>4</sup>, Deny Hidayatullah<sup>5</sup>,  
Dhieka Avrilia Lantana<sup>6</sup>, Yuni Latifah<sup>7</sup>, Alica Dwi Fahira<sup>8</sup>, Keysha Belynda Tyva Panggabean<sup>9</sup>,  
Imelta Natalia Ginting<sup>10</sup>

(1), (3), (4), (5), (6) Universitas Nasional, Indonesia

(2) Program Studi Animasi Politeknik Negeri Media Kreatif, Indonesia

(7), (8), (9), (10) UPN Veteran Jakarta, Indonesia

Email: [erimardiani1@gmail.com](mailto:erimardiani1@gmail.com) 

### Abstrak

Data mining bermula dari peningkatan data yang cukup pesat dilihat dari segi volume serta variasi data yang dihasilkan oleh berbagai sumber, dan jumlahnya yang sangat besar, serta kompleksitasnya data hingga pembuatannya yang cepat. Dengan data bisa menghasilkan prediksi yang membantu pemerintah dalam mengambil keputusan dan kebijakan di masa mendatang. Selain itu prediksi dapat membantu pemerintah dalam perencanaan kegiatan yang akan dilakukan untuk mencapai tujuan, karena prediksi ini dapat memberikan output terbaik sehingga diharapkan resiko kesalahan yang disebabkan oleh kesalahan perencanaan dapat ditekan seminimal mungkin. Prediksi biasanya digunakan untuk menemukan informasi dari sejumlah data yang besar sehingga diperlukan data mining. Data mining dapat digunakan untuk menggali informasi dari data yang besar sehingga didapatkan informasi yang dapat digunakan dalam memprediksi sesuatu. Dalam data mining terdapat banyak teknik dalam pengerjaannya, untuk menemukan pola atau informasi yang tersembunyi diantaranya adalah Klasterisasi (clustering), Regresi (regression), Asosiasi (association), dan Klasifikasi (classification)

Kata Kunci : *Data mining, income*

### Abstract

Copyright © Eri Mardiani, Nur Rahmansyah, Endah Tri Esti Handayani, Sari Ningsih, Deny Hidayatullah, Dhieka Avrilia Lantana, Yuni Latifah, Alica Dwi Fahira, Keysha Belynda Tyva Panggabean, Imelta Natalia Ginting

Data mining begins with a fairly rapid increase in data in terms of the volume and variety of data produced by various sources, and the numbers are very large, as well as the complexity of the data and its rapid creation. Data can produce predictions that help the government in making decisions and policies in the future. Apart from that, predictions can help the government in planning activities that will be carried out to achieve goals, because these predictions can provide the best output so it is hoped that the risk of errors caused by planning errors can be reduced to a minimum. Predictions are usually used to find information from large amounts of data, so data mining is needed. Data mining can be used to dig up information from large amounts of data to obtain information that can be used to predict something. In data mining, there are many techniques used to find hidden patterns or information, including clustering, regression, association and classification.

*Keywords: Keywords contain basic ideas or concepts that represent the field under study; The number of keywords is between 3-5 Phrases and are sorted alphabetically*

## PENDAHULUAN

Di era digital saat ini, ekonomi merupakan salah satu industri terpenting dan berkembang sangat cepat. Hal ini didukung oleh perkembangan teknologi yang semakin pesat di seluruh dunia. Itulah sebabnya setiap orang yang berkecimpung dalam bisnis bergerak dan saling bersaing untuk menjadi yang terbaik di bidangnya.[10]

Penelitian ini mempunyai tujuan agar dapat membandingkan metode K-NN, Naive Bayes, Decision Tree, Ensemble Methods, dan Linear Regression yang dilakukan untuk klasifikasi apakah seorang penduduk memiliki pendapatan lebih besar atau sama dengan \$50.000 USD per tahun atau tidak. Sedangkan aplikasi yang digunakan adalah aplikasi data mining orange yang merupakan aplikasi data mining open source yang terbukti dapat membantu dalam hal penganalisaan data. Untuk melakukan hal tersebut kita akan menunjukkan prosesnya mulai dari akuisisi data sampai prediksi.[1]

Data Mining menggunakan Teknik Klasifikasi dengan 5 Model Algoritma[2]

### 1. K-Nearest Neighbor (k-NN) [3]

Algoritma K-Nearest Neighbor (k-NN) merupakan sebuah model algoritma yang digunakan untuk melakukan klasifikasi terhadap objek yang berdasarkan dari data pembelajaran yang jaraknya paling dekat dengan objek tersebut. Ataupun dapat dipahami juga bahwa k- nearest neighbor adalah salah satu algoritma yang paling sederhana dan banyak. Titik data akan diklasifikasikan berdasarkan kesamaan kelompok tertentu dari titik data lain yang berdekatan. Sehingga, algoritma ini akan memberikan hasil yang kompetitif.

### 2. Naïve Bayes[5]

Salah satu metode data mining ialah klasifikasi Naive Bayes. Naive Bayes Classifier adalah metode klasifikasi yang berakar pada teorema Bayes. Naive bayes merupakan metode

pengklasifikasian berdasarkan probabilitas sederhana dan dirancang agar dapat dipergunakan dengan asumsi antar variabel penjelas saling bebas (independen). Pada algoritma ini pembelajaran lebih ditekankan pada pengestimasi probabilitas. Tujuan dari metode Naïve Bayes adalah untuk menemukan probabilitas ketika kita mengetahui probabilitas tertentu lainnya. Hasil dari perhitungan data mining menggunakan metode klasifikasi Naïve Bayes akan makin berguna jika penyajiannya menarik dan dapat dipahami dengan baik oleh penerima data.

### 3. Decision Tree[9]

Decision tree adalah algoritma machine learning yang menggunakan seperangkat aturan untuk membuat keputusan dengan struktur seperti pohon yang memodelkan kemungkinan hasil, biaya sumber daya, utilitas dan kemungkinan konsekuensi atau resiko. Konsepnya adalah dengan cara menyajikan algoritma dengan pernyataan bersyarat, yang meliputi cabang untuk mewakili langkah-langkah pengambilan keputusan yang dapat mengarah pada hasil yang menguntungkan. Klasifikasi ini menggunakan observasi pada node untuk menemukan target pada leaves. Decision Tree merupakan salah satu metode klasifikasi yang paling populer karena mudah untuk diinterpretasi oleh manusia dengan kemampuannya untuk mem-break down proses pengambilan keputusan yang kompleks menjadi lebih simple.[6]

### 4. Ensemble Method

Ensemble Method adalah algoritma dalam pembelajaran mesin (machine learning) dimana algoritma ini sebagai pencarian solusi prediksi terbaik dibandingkan dengan algoritma yang lain karena metode ensemble ini menggunakan beberapa algoritma pembelajaran untuk pencapaian solusi prediksi yang lebih baik daripada algoritma yang bisa diperoleh dari salah satu pembelajaran algoritma konstituen saja. Tidak seperti ensemble statistika dalam mekanika statistika biasanya selalu tak terbatas. Ensemble Pembelajaran hanya terdiri dari seperangkat model alternatif yang bersifat terbatas, namun biasanya memungkinkan untuk menjadi lebih banyak lagi struktur fleksibel yang ada diantara alternatif model itu sendiri. Evaluasi prediksi dari ensemble biasanya memerlukan banyak komputasi daripada evaluasi prediksi model tunggal (single model).

### 5. Linear Regression[4]

Linear regression merupakan salah satu algoritma yang memodelkan suatu persamaan untuk menghitung estimasi. Pada metode ini bertujuan untuk mencari pola pada nilai

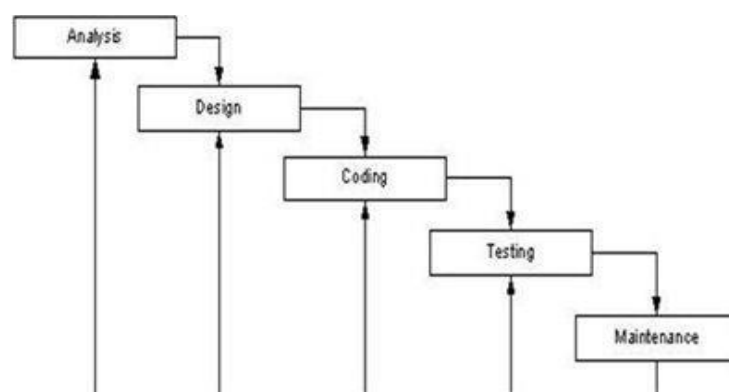
numerik,

sehingga data yang dibutuhkan berupa data numerik agar dapat diolah dengan model algoritma ini. Linear regression memprediksi nilai data yang tidak diketahui dengan menggunakan nilai data lain yang terkait dan diketahui. Secara matematis memodelkan variabel yang tidak diketahui atau tergantung dan variabel yang dikenal atau independen sebagai persamaan linier..

## METODE PENELITIAN

Data yang dikumpulkan untuk penelitian ini yaitu primer dan sekunder. Teknik pengumpulan data primer berasal dari studi literatur seperti buku maupun jurnal. Adapun Teknik pengumpulan data sekunder pada penelitian ini diperoleh dari media online maupun sumber-sumber lainnya. Studi pustaka akan digunakan oleh peneliti guna menjabarkan dan menganalisis data yang sudah dikumpulkan yang memiliki kaitan dengan penelitian ini .[7]

Selain metode kualitatif , kami juga menggunakan metode waterfall karena dibutuhkan pendekatan yang sistematis dan berurutan dalam membangun sistem. Aliran cascading adalah bahwa sistem ditulis secara berurutan. Sistem yang dihasilkan menghasilkan kualitas yang sangat baik karena implementasinya bertahap dan tidak terfokus pada fase tertentu [12]



Gambar 1. Metode Waterfall

Pada tahap pengembangan sebuah sistem, metode ini bersifat sistematis, serta rangkaian dalam sistem perancangan perangkat lunak dilakukan secara urut. Metode waterfall atau disebut juga metode air terjun merupakan metode yang sering dipakai untuk melakukan pembaharuan pada sistem aplikasi yang sedang berjalan. Terdapat beberapa tahapan pada metode ini, yaitu analisis aplikasi, desain aplikasi, implementasi aplikasi, pengujian aplikasi, dan pemeliharaan [8][11]

# HASIL DAN PEMBAHASAN

Hasil

A. K-Nearest Neighbor

Implementasi menggunakan K-Nearest Neighbor (KNN)

The screenshot shows a table with columns: knn, error, salary, native-country, age, workclass, triaget, education, education-num, marital-status, occupation, relationship, race, sex. The table contains 34 rows of data. Below the table, there are performance metrics: Model AUC CA F1 Precision Recall, with values: knn 0.944 0.882 0.879 0.879 0.882.

Gambar 2 Tabel Prediction

Berdasarkan hasil penelitian menggunakan Predictions dapat dilihat tabel nilai prediksi dan nilai aktualnya. Jika dilihat secara manual oleh tabel pada Predictions, nilai prediksi menggunakan metode K-NN merupakan hampir keseluruhan nilainya mendeskripsikan dengan benar kepada nilai aktualnya. Hal tersebut didukung dengan nilai K-NN yang berada di bawah, nilai K-NN memperoleh nilai AUC sebesar 94,4% yang berarti sangat baik karena hampir mendekati 100%, nilai CA sebesar 88,2%, nilai F1 sebesar 87,9%, nilai Precisions sebesar 87,9%, dan nilai Recall sebesar 88,2%.

		Predicted		
		<=50K	>50K	Σ
Actual	<=50K	24465	1512	25977
	>50K	2523	5690	8213
Σ		26988	7202	34190

Gambar 3 Confusion KNN

Selanjutnya, penilaian metode K-NN dapat dilihat dengan evaluasi menggunakan Confusion Matrix. Pada gambar diatas dijelaskan bahwa untuk variabel <=50K menunjukkan hasil prediksi dengan metoda K-NN memiliki nilai benar sebanyak 24,465 sedangkan nilai salah sebanyak 2.523. Kemudian, untuk variabel >50K menunjukkan hasil prediksi dengan metode K-NN memiliki nilai benar sebanyak 5.690, sedangkan nilai salah sebanyak 1.512.

Hasil penelitian dengan Confusion Matrix juga dapat dilihat menggunakan ikon Data Table, untuk melihat secara manual dan lebih jelas, manakah nilai prediksi dari metode K-NN yang benar dan manakah nilai prediksi dari metode K-NN yang salah. Untuk hasil Data Table ditampilkan seperti gambar yang tertera di bawah.

id	salary	native-country	knn	knn (<=50K)	knn (>50K)	age	workclass	privgt	education	education-num	marital-status	occupation	relation
1	<=50K	United-States	<=50K	1	0	18	Private	423024	HS-grad	9	Never-married	Other-service	Not-in-fan
2	<=50K	United-States	<=50K	1	0	17	Private	178953	12th	8	Never-married	Sales	Own-child
3	<=50K	United-States	<=50K	1	0	25	Local-gov	348986	HS-grad	9	Never-married	Handlers-clean...	Other-rela
4	<=50K	United-States	<=50K	1	0	20	Private	218215	Some-college	10	Never-married	Sales	Own-child
5	<=50K	Puerto-Rico	<=50K	1	0	47	Private	248025	HS-grad	9	Never-married	Machine-op-ins...	Unmarried
6	<=50K	United-States	>50K	0.4	0.6	33	Private	399531	Bachelors	13	Married-civ-sp...	Craft-repair	Husband
7	<=50K	United-States	<=50K	0.6	0.4	38	Private	200220	HS-grad	9	Married-civ-sp...	Craft-repair	Husband
8	<=50K	Mexico	<=50K	1	0	21	Private	329530	11th	7	Married-civ-sp...	Craft-repair	Husband
9	<=50K	United-States	<=50K	0.8	0.2	43	Private	282155	Assoc-acdm	12	Divorced	Prof-specialty	Not-in-fan
10	<=50K	United-States	<=50K	1	0	55	Private	202220	HS-grad	9	Married-civ-sp...	Other-service	Wife
11	>50K	United-States	>50K	0	1	46	Private	129607	Bachelors	13	Married-civ-sp...	Sales	Husband
12	>50K	United-States	<=50K	0.8	0.2	34	Private	261799	Assoc-voc	11	Married-civ-sp...	Adm-clerical	Husband
13	<=50K	United-States	<=50K	1	0	40	?	246862	Bachelors	13	Widowed	?	Not-in-fan
14	<=50K	United-States	<=50K	0.6	0.4	50	Private	173754	HS-grad	9	Married-civ-sp...	Craft-repair	Husband
15	<=50K	United-States	<=50K	1	0	29	Private	176727	Some-college	10	Never-married	Craft-repair	Not-in-fan
16	<=50K	United-States	<=50K	1	0	62	Private	24515	9th	5	Married-civ-sp...	Exec-managerial	Husband
17	<=50K	United-States	<=50K	1	0	56	Private	158776	11th	7	Married-civ-sp...	Sales	Husband
18	>50K	United-States	>50K	0.4	0.6	40	Federal-gov	90737	Some-college	10	Married-civ-sp...	Adm-clerical	Husband
19	<=50K	United-States	<=50K	1	0	74	?	340599	9th	5	Married-civ-sp...	?	Husband
20	<=50K	United-States	<=50K	0.8	0.2	21	Private	305874	Some-college	10	Married-civ-sp...	Craft-repair	Husband
21	<=50K	United-States	<=50K	1	0	31	Private	74501	HS-grad	9	Divorced	Other-service	Unmarried
22	>50K	United-States	>50K	0	1	36	Private	295706	Masters	14	Married-civ-sp...	Exec-managerial	Husband
23	<=50K	United-States	<=50K	1	0	36	Private	108293	HS-grad	9	Widowed	Other-service	Unmarried
24	>50K	United-States	>50K	0.2	0.8	64	State-gov	111795	Bachelors	13	Married-civ-sp...	Craft-repair	Husband
25	>50K	United-States	>50K	0.8	0.2	32	Local-gov	250585	Bachelors	13	Never-married	Prof-specialty	Not-in-fan
26	<=50K	United-States	<=50K	0.8	0.2	55	Local-gov	159028	HS-grad	9	Married-civ-sp...	Transport-mov...	Husband
27	>50K	United-States	<=50K	1	0	33	Private	179758	HS-grad	9	Married-civ-sp...	Sales	Wife
28	>50K	United-States	<=50K	0.6	0.4	52	Self-emp-not-inc	34973	HS-grad	9	Married-civ-sp...	Farming-fishing	Husband
29	>50K	United-States	<=50K	0.6	0.4	45	Self-emp-inc	204196	Bachelors	13	Divorced	Exec-managerial	Unmarried
30	<=50K	United-States	<=50K	0.8	0.2	47	Private	47247	Some-college	10	Married-civ-sp...	Adm-clerical	Wife
31	<=50K	United-States	<=50K	1	0	27	Private	150080	Bachelors	13	Never-married	Exec-managerial	Own-child
32	<=50K	United-States	<=50K	1	0	19	Private	376540	HS-grad	9	Never-married	Adm-clerical	Not-in-fan
33	>50K	United-States	>50K	0.2	0.8	52	Private	99736	Masters	14	Divorced	Prof-specialty	Unmarried
34	<=50K	United-States	<=50K	1	0	26	Private	166301	Bachelors	13	Never-married	Tech-support	Not-in-fan
35	<=50K	United-States	<=50K	1	0	50	Private	321770	HS-grad	9	Divorced	Adm-clerical	Not-in-fan
36	<=50K	United-States	<=50K	1	0	40	Private	106698	Assoc-acdm	12	Divorced	Transport-mov...	Unmarried
37	<=50K	United-States	<=50K	1	0	45	Self-emp-not-inc	32172	Some-college	10	Never-married	Farming-fishing	Not-in-fan
38	<=50K	United-States	<=50K	1	0	25	Private	243410	HS-grad	9	Never-married	Other-service	Not-in-fan

Gambar 4 Hasil Data Table

## B. Naive Bayes

Model	AUC	CA	F1	Precision	Recall
Naive Bayes	0.903	0.824	0.831	0.847	0.824

Gambar 5 Data Tabel Naive Bayes

Berdasarkan hasil penelitian menggunakan Test and Score dapat dilihat pada gambar diatas, nilai yang diperoleh menggunakan metode Naive Bayes. Nilai Naive Bayes memperoleh nilai AUC sebesar 90,3% yang berarti sangat baik karena hampir mendekati 100%, nilai CA sebesar 82,4%, nilai F1 sebesar 83,1%, nilai Precisions sebesar 84,7%, dan nilai Recall sebesar 82,4%.

		Predicted		
		<=50K	>50K	Σ
Actual	<=50K	31129	6026	37155
	>50K	2555	9132	11687
Σ		33684	15158	48842

Gambar 6 Confusion Naïve Bayes

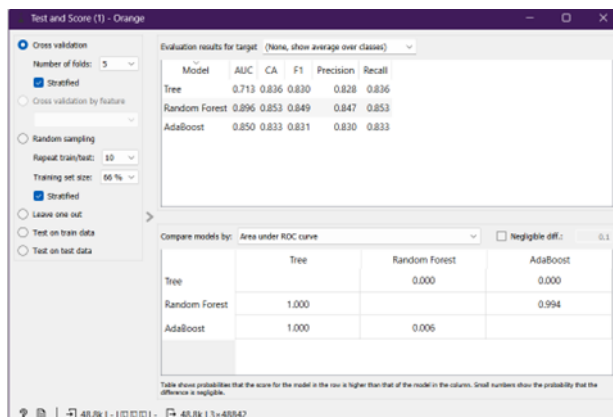
Selanjutnya, penilaian metode Naive Bayes dapat dilihat dengan evaluasi menggunakan Confusion Matrix. Pada gambar diatas dijelaskan bahwa untuk variabel <=50K menunjukkan hasil prediksi dengan metode Naive Bayes memiliki nilai benar sebanyak 31.129, sedangkan nilai salah sebanyak 2555. Kemudian, untuk variabel >50K menunjukkan hasil prediksi dengan metode Naive Bayes memiliki nilai benar sebanyak 9.132, sedangkan nilai salah sebanyak 6.026.

	y	Naive Bayes	Naive Bayes (>50K)	Naive Bayes (<=50K)	Fold
1	<=50K	<=50K	0.00908346	0.990917	1
2	<=50K	<=50K	0.162727	0.837273	1
3	<=50K	<=50K	0.0771687	0.922831	1
4	<=50K	<=50K	0.490398	0.509602	1
5	<=50K	<=50K	9.98283e-05	0.9999	1
6	<=50K	<=50K	8.18219e-08	1	1
7	<=50K	<=50K	0.10121	0.89879	1
8	<=50K	<=50K	0.128342	0.871658	1
9	<=50K	<=50K	6.6426e-05	0.999934	1
10	<=50K	>50K	0.888352	0.111648	1
11	<=50K	<=50K	0.0358846	0.964115	1
12	>50K	>50K	0.983005	0.0169948	1
13	<=50K	<=50K	0.000424583	0.999575	1
14	>50K	>50K	0.579867	0.420133	1
15	<=50K	<=50K	0.0812535	0.918746	1
16	>50K	>50K	0.740311	0.259689	1
17	<=50K	<=50K	1.87571e-05	0.999981	1

Gambar 6 Hasil Data Table Naïve Bayes

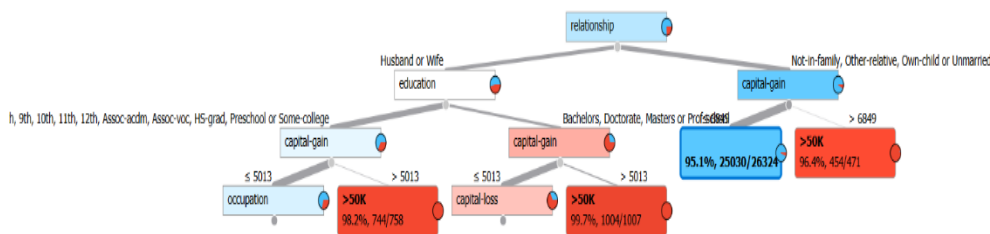
Hasil penelitian dengan Test and Score juga dapat dilihat menggunakan ikon Data Table yang tertera pada gambar diatas, untuk melihat secara manual dan lebih jelas, manakah nilai prediksi dari metode Naive Bayes yang benar dan manakah nilai prediksi dari metode Naive Bayes yang salah.

### c. Decision Tree dan Ensemble Tree



Gambar 7 Test And Score

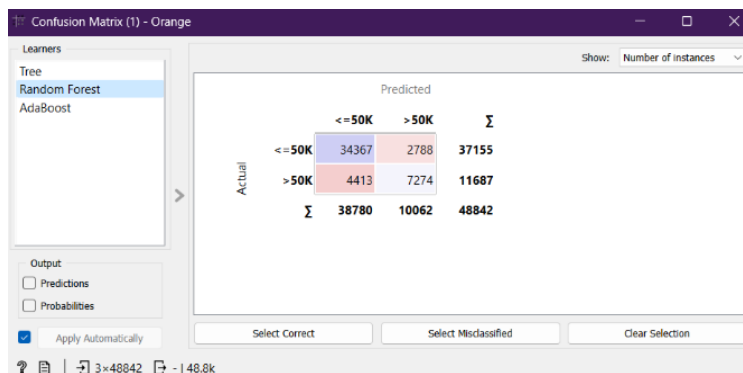
Berdasarkan hasil penelitian menggunakan Test and Score dapat dilihat pada gambar diatas, nilai yang diperoleh menggunakan metode Tree, Random Forest, dan AdaBoost. Untuk nilai Tree memperoleh nilai AUC sebesar 71,3% yang berarti sangat baik karena hampir mendekati 100%, nilai CA sebesar 83,6%, nilai F1 sebesar 83%, nilai Precisions sebesar 82,8%, dan nilai Recall sebesar 83,6%. Untuk nilai Random Forest memperoleh nilai AUC sebesar 89,6% yang berarti sangat baik karena hampir mendekati 100%, nilai CA sebesar 85,3%, nilai F1 sebesar 84,9%, nilai Precisions sebesar 84,7%, dan nilai Recall sebesar 85,3%. Untuk nilai AdaBoost memperoleh nilai AUC sebesar 85% yang berarti sangat baik karena hampir mendekati 100%, nilai CA sebesar 83,3%, nilai F1 sebesar 83,1%, nilai Precisions sebesar 83%, dan nilai Recall sebesar 83,3%. Dari ketiga metode, disimpulkan bahwa nilai metode terbaik untuk prediksi adalah metode Random Forest.



Gambar 8 Hasil Tree Viewer

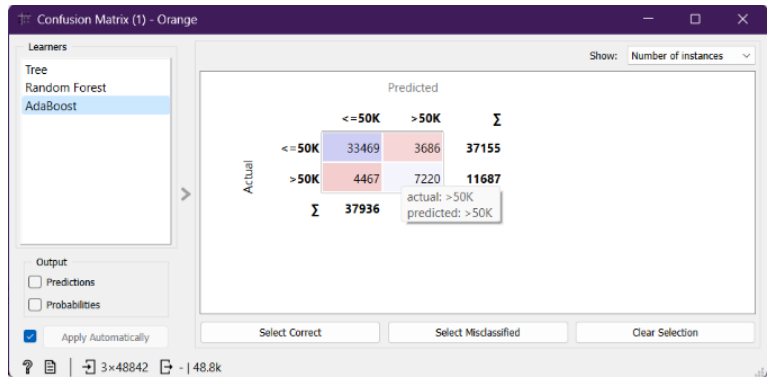
Untuk melihat hasil evaluasi Tree bisa digunakan Tree Viewer, pada gambar dijelaskan bahwa variabel-variabel terbagi ke dalam beberapa klasifikasi. Penjelasan dalam gambar diatas terdiri dari beberapa cabang. Cabang Utama yang menjelaskan mengenai hubungan

atau relationship, memiliki dua cabang lainnya yang pertama mengenai status pernikahan yang berpendidikan dan yang masih single dengan memiliki pendapatan sendiri. Dalam status pernikahan terdapat dua cabang mengenai beberapa lama waktu menempuh pendidikan dan status pernikahan yang sudah mempunyai gelar pendidikan (bachelor, master, etc). Sedangkan untuk pribadi yang masih single dan mendapatkan pendapatan sendiri memiliki dua cabang dimana nilai yang kurang dari 6849 sebesar 95,1% dan yang lebih dari 6849 sebesar 96,4%. Lalu, untuk status pernikahan yang sedang mengambil gelar memiliki dua cabang yaitu kurang dari 5013 dengan hasil occupation (pekerjaan) dan yang lebih dari 5013 dengan hasil pendapatan lebih dari \$50.000 memiliki hasil sebesar 98,2%. Lalu untuk status pernikahan yang sudah mendapatkan gelar memiliki dua cabang yang pertama kurang dari 5013 menghasilkan capital-loss dan yang lebih dari 5013 dengan hasil pendapatan lebih dari \$50.000 mendapatkan nilai sebesar 99.7%.



Gambar 9 Confusion Tree Viewer

Selanjutnya, penilaian metode Tree dapat dilihat dengan evaluasi menggunakan Confusion Matrix. Pada gambar diatas dijelaskan bahwa untuk variabel <=50K menunjukkan hasil prediksi dengan metode Tree memiliki nilai benar sebanyak 34.084, sedangkan nilai salah sebanyak 4.959. Kemudian, untuk variabel >50K menunjukkan hasil prediksi dengan metode Tree memiliki nilai benar sebanyak 6.728, sedangkan nilai salah sebanyak 3.071.



Gambar 10 Confusion Random Forest

Untuk penilaian metode Random Forest dapat dilihat dengan evaluasi menggunakan Confusion Matrix. Pada gambar diatas dijelaskan bahwa untuk variabel <=50K menunjukkan hasil prediksi dengan metode Random Forest memiliki nilai benar sebanyak 34.367, sedangkan nilai salah sebanyak 4.413. Kemudian, untuk variabel >50K menunjukkan hasil prediksi dengan metode Random Forest memiliki nilai benar sebanyak 7.274, sedangkan nilai salah sebanyak 2.728.

Untuk penilaian metode AdaBoost dapat dilihat dengan evaluasi menggunakan Confusion Matrix. Pada gambar diatas dijelaskan bahwa untuk variabel <=50K menunjukkan hasil prediksi dengan metode AdaBoost memiliki nilai benar sebanyak 33.469, sedangkan nilai salah sebanyak 4.467. Kemudian, untuk variabel >50K menunjukkan hasil prediksi dengan metode AdaBoost memiliki nilai benar sebanyak 7.220, sedangkan nilai salah sebanyak 3.686.

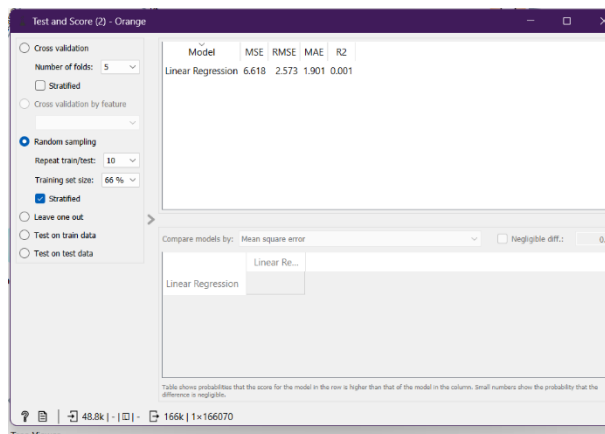
Hasil penelitian dengan Test and Score juga dapat dilihat menggunakan

Instance	Salary	Native-Country	Tree	Random Forest	AdaBoost	Tree (<=50K)	Tree (>50K)	Random Forest (<=50K)	Random Forest (>50K)	AdaBoost (<=50K)	AdaBoost (>50K)	F1
1	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	0.766667	0.233333	0.973557	0.0264432	1
2	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.981144	0.0188558	1
3	>50K	United-States	>50K	<=50K	>50K	0	0	0.625556	0.374444	0.199018	0.810382	1
4	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	0.818667	0.181333	0.973963	0.0270367	1
5	>50K	United-States	>50K	>50K	>50K	0.951464	0.0483363	0.4875	0.5125	0.00626589	0.993734	1
6	<=50K	United-States	>50K	>50K	>50K	0	1	0.125	0.875	0.193317	0.806683	1
7	<=50K	United-States	>50K	<=50K	>50K	0	1	0.545	0.455	0.499004	0.500996	1
8	<=50K	United-States	<=50K	<=50K	<=50K	0.75	0.25	0.775	0.225	0.89482	0.10518	1
9	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.989738	0.0102619	1
10	<=50K	Mexico	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.996548	0.00345206	1
11	<=50K	United-States	<=50K	>50K	<=50K	1	0	0.466667	0.533333	0.504003	0.495997	1
12	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.997866	0.00213443	1
13	<=50K	Cuba	>50K	<=50K	<=50K	0	1	0.741667	0.258333	0.008307	0.991693	1
14	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	0.983333	0.016667	0.988187	0.0118128	1
15	>50K	United-States	>50K	>50K	>50K	0.75	0.25	0.288333	0.711667	0.193682	0.806318	1
16	<=50K	England	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.986812	0.0131884	1
17	<=50K	United-States	>50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.98267	0.0173312	1
18	<=50K	United-States	>50K	<=50K	<=50K	0.0266667	0.9733333	0.55	0.45	0.946091	0.0539091	1
19	<=50K	United-States	<=50K	<=50K	<=50K	1	0	0.91	0.09	0.877867	0.122133	1
20	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.999721	0.000279138	1
21	>50K	United-States	<=50K	<=50K	>50K	0.666667	0.333333	0.528333	0.471667	0.499241	0.500759	1
22	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.975617	0.0243823	1
23	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.988187	0.0118128	1
24	>50K	United-States	<=50K	<=50K	<=50K	1	0	0.878667	0.121333	0.011557	0.988443	1
25	<=50K	United-States	>50K	>50K	>50K	0	1	0.138909	0.861091	0.189672	0.810328	1
26	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	0.93	0.07	0.987192	0.0128081	1
27	>50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	0.894048	0.105952	0.973725	0.0262752	1
28	<=50K	United-States	<=50K	<=50K	<=50K	0.970297	0.029703	1	0	0.816791	0.183209	1
29	<=50K	Mexico	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.988187	0.0118128	1
30	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.999737	0.000262742	1
31	>50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.995006	0.00499399	1
32	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.989816	0.0101838	1
33	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.976463	0.023537	1
34	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.990945	0.00905455	1
35	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.998932	0.00106064	1
36	<=50K	England	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.999054	0.000945053	1
37	<=50K	Mexico	<=50K	<=50K	<=50K	0.951464	0.0483363	1	0	0.988187	0.0118128	1
38	<=50K	United-States	<=50K	<=50K	<=50K	0.951464	0.0483363	0.801429	0.198571	0.809774	0.190226	1

Gambar 11 Data Table

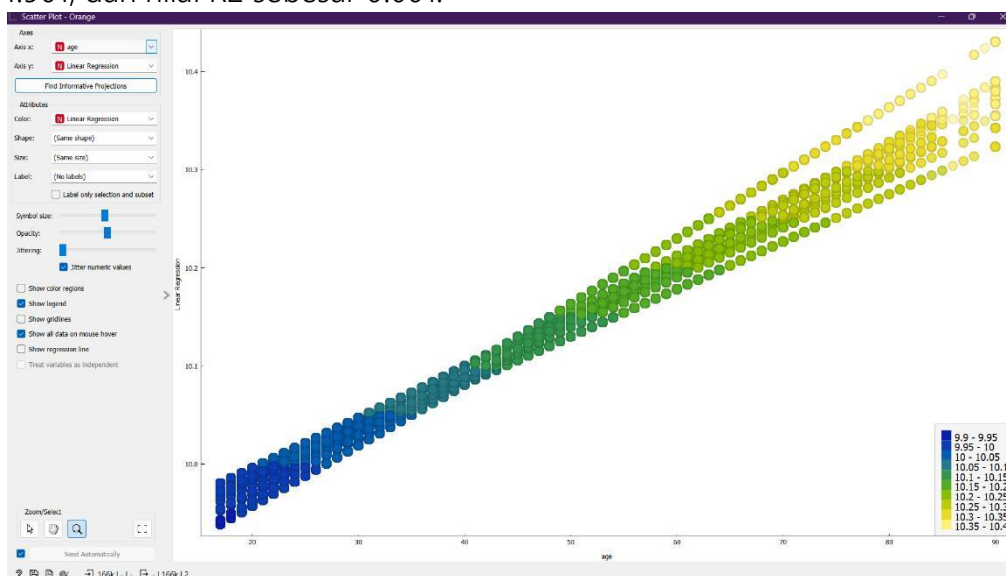
ikon Data Table yang tertera pada gambar diatas, untuk melihat secara manual dan lebih jelas, manakah nilai prediksi dari ketiga metode yang benar dan manakah nilai prediksi dari ketiga metode yang salah.

#### D. Linear Regression



Gambar 12 Test and Score

Berdasarkan hasil penelitian menggunakan Test and Score dapat dilihat pada gambar diatas, nilai yang diperoleh menggunakan metode Linear Regression. Nilai Linear Regression memperoleh nilai MSE sebesar 6.618 yang berarti bahwa pada saat nanti kita melakukan prediksi akan terdapat error sebesar 6.618, nilai MRSE sebesar 2.573, nilai MAE sebesar 1.901, dan nilai R2 sebesar 0.001.



Gambar 3 Scatter Plotter

Untuk penilaian metode Linear Regression dapat dilihat dengan evaluasi menggunakan Scatter Plotter. Pada gambar diatas dijelaskan bahwa jika yang menjauhi garis, maka nilai prediksi terdapat kesalahan.

## SIMPULAN

Prediksi atau peramalan (forecasting) adalah suatu perhitungan untuk meramalkan keadaan di masa mendatang melalui pengujian keadaan di masa lalu. Salah satu dari kegunaan prediksi adalah untuk membantu pemerintah dalam mengambil keputusan dan kebijakan di masa mendatang. Selain itu prediksi dapat membantu pemerintah dalam perencanaan kegiatan yang akan dilakukan untuk mencapai tujuan, karena prediksi ini dapat memberikan output terbaik sehingga diharapkan resiko kesalahan yang disebabkan oleh kesalahan perencanaan dapat ditekan seminimal mungkin. Prediksi biasanya digunakan untuk menemukan informasi dari sejumlah data yang besar sehingga diperlukan data mining. Data mining merupakan bidang dari beberapa bidang keilmuan yang menyatukan teknik dari pembelajaran mesin, pengenalan pola, statistic, database dan visualisasi untuk penanganan permasalahan pengambilan informasi dari penyimpanan database yang besar. Data mining dapat digunakan untuk menggali informasi dari data yang besar sehingga didapatkan informasi yang dapat digunakan dalam memprediksi sesuatu. Dalam data mining terdapat banyak teknik dalam pengerjaannya, untuk menemukan pola atau informasi yang tersembunyi diantaranya adalah Klasterisasi (clustering), Regresi (regression), Asosiasi (association), dan Klasifikasi (classification).

## DAFTAR PUSTAKA

- [1] <https://www.kaggle.com>
- [2] FH Pratama, A Triayudi, E Mardiani.(2022). Data Mining K-Medoids Dan K-Means Untuk Pengelompokan Potensi Produksi Kelapa Sawit di Indonesia. JIPI (Jurnal Ilmiah Penelitian dan Pembelajaran Informatika) 7 (4), 1294-1310
- [3] Indriyanti, Ichsan, Nurul., Fatah, Haerul ,. Wahyuni,. Tri., Ermawati, Erni. Implementasi Orange Data Mining Untuk Prediksi Harga Bitcoin. JURNAL RESPONSIF, Vol. 4 No.2 Agustus 2022, pp. 118~125
- [4] Indriyawati.Henny , Khoirudin.2019. PENERAPAN METODE REGRESI LINIER DALAM KOHERENSI PENGOLAHAN DATA BAHAN BAKU TIANDRA STORE GUNA MENINGKATKAN MUTU PRODUKSI. Proceeding SINTAK 2019 Universitas Semarang <https://www.unisbank.ac.id/ojs/index.php/sintak/article/view/7603>
- [5] MA Djamaludin, A Triayudi, E Mardiani.(2022) Analisis Sentimen Tweet KRI Nanggala 402 di Twitter menggunakan Metode Naïve Bayes Classifier. Jurnal JTIK(Jurnal Teknologi Informasi dan Komunikasi ) 6 (2) 2022 pp 2580-1643
- [6] Mardiani, E., Hayati, N., Purnama, M. I. W., Ningsih, S., Darusalam, U., Handayani, T. E.,

- & Rahmansyah, N. (2023). *Kumpulan Latihan VB. Net*. Elex Media Komputindo.
- [7] Mardiani, E., Rahmansyah, N., & Kurniati, I. (2023). Website Design at SDN Cipete Utara 07. *SITEKIN: Jurnal Sains, Teknologi Dan Industri*, 20(2), 891–898.
- [8] Mardiani, E., Rahmansyah, N., Kurniawan, H., & Sensuse, D. I. (2016). *Kumpulan Latihan SQL*. Elex Media Komputindo.
- [9] Mardiani, E., Rahmansyah, N., Ningsih, S., Lantana, D. A., Wirawan, A. S. P., Wijaya, S. A., & Putri, D. N. (2023). Komparasi Metode Knn, Naive Bayes, Decision Tree, Ensemble, Linear Regression Terhadap Analisis Performa Pelajar Sma. *Innovative: Journal Of Social Science Research*, 3(2), 13880–13892.
- [10] Mardiani, E., & Ramadhan, F. A. (2023). Design Information System Sales of Nuts and Bolts at PT. Catur Naga Steelindo. *SITEKIN: Jurnal Sains, Teknologi Dan Industri*, 20(2), 729–735.
- [11] Matondang, Nurhafifah, Mardiani, E., Wahyudi, Praptiningsih, & Saebani, A. (2019). *Aplikasi Komputer*.
- [12] Rahmansyah, N., Mulyani, D., Mardiani, E., & Rahman, A. (2022). Perancangan Sistem Transaksi Berbasis Web pada UKM Pangkas Rambut Tasik. *Jurnal Sistem Informasi Bisnis (JUNSIBI)*, 3(1), 22–31.