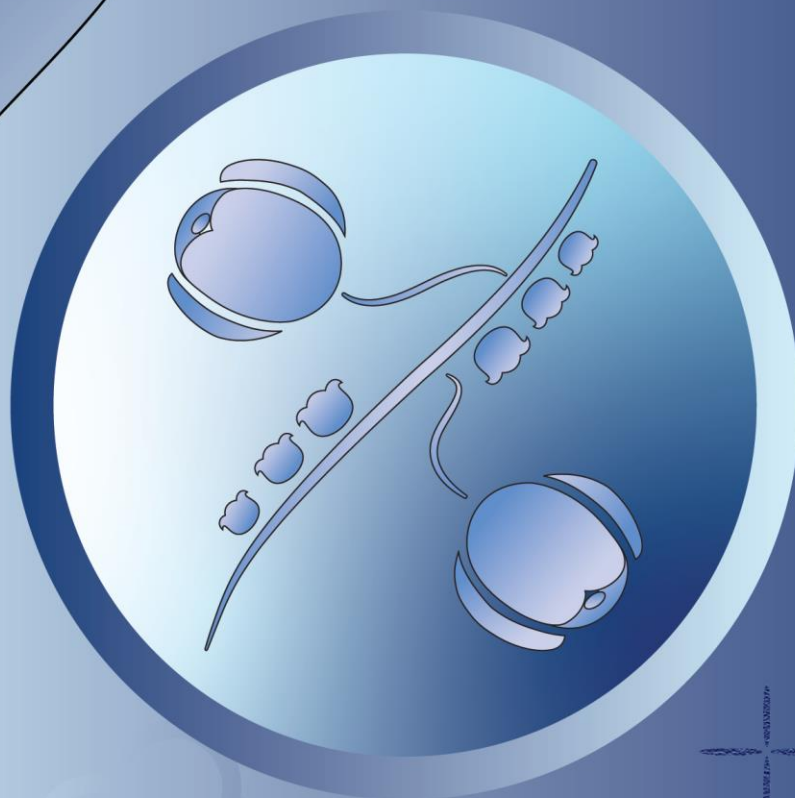




Ajang Pengenalan Statistika dan Festival Data #19



# Makalah



## DATAVERS

**NAMA TIM**

Aquilae

**NOMOR PESERTA**

DVS0000278

## 1. PENDAHULUAN

Manusia telah memanfaatkan mesin untuk memudahkan kehidupan sejak awal peradaban. Mesin sederhana seperti roda, poros dan pengerek telah digunakan oleh orang-orang Mesopotamia dan Mesir kuno yang eksis pada Zaman Perunggu [1]. Mesin termekanisasi pertama yang telah diciptakan manusia adalah Mekanisme Antikythera (Αντικύθηρα). Mekanisme berbasis roda gigi ini ditemukan di dalam satu bangkai kapal yang terdapat di pesisir pulau Yunani Antikythera pada tahun 1902 Masehi dan bertanggal kembali ke sekitar tahun 100-150 sebelum Masehi. Mesin ini digunakan untuk memprediksi posisi astronomis benda langit, misalnya untuk memprediksi kejadian gerhana [2].

Sepanjang sejarah peradaban manusia, berbagai jenis mesin telah ditemukan dan dikembangkan dengan kompleksitas yang terus meningkat. Hal ini dilakukan untuk menyelesaikan berbagai masalah baru yang muncul seiring dengan kompleksitas peradaban manusia. Mesin uap pertama kali ditemukan oleh Thomas Savery pada tahun 1698 Masehi, diikuti oleh Thomas Newcomen pada tahun 1712 Masehi [3]. Mesin ini kemudian dikembangkan lebih lanjut oleh James Watt pada pertengahan abad ke-18 Masehi yang menandakan kemajuan signifikan pada perkembangan mesin. Perkembangan ini menjadi pemicu dimulainya Revolusi Industri, yang pada akhirnya menciptakan mesin-mesin kompleks yang memungkinkan keberadaan pabrik-pabrik modern dan berbagai industri baru.

Seperti segala sesuatu yang diciptakan oleh manusia, semua mesin pada akhirnya akan mengalami kerusakan sehingga berhenti berfungsi. Kerusakan dapat terjadi akibat deteriorasi pada bagian-bagian yang membentuk suatu mesin tertentu, yang dapat mengakibatkan beberapa hal buruk untuk terjadi [4]. Kejadian seperti terganggunya proses produksi serta kegagalan dalam memenuhi permintaan pelanggan dapat terjadi sebagai akibat dari kerusakan mesin [5]. Hal ini yang menjadi alasan untuk melakukan penelitian guna mengetahui faktor-faktor yang berpotensi meningkatkan kerusakan mesin.

Berdasarkan penelitian yang telah dilakukan sebelumnya, deteriorasi mesin berpotensi untuk terjadi apabila terdapat bagian mesin yang cacat dan sistem perawatan yang kurang efektif [6]. Analisis lebih lanjut diperlukan untuk menemukan berbagai faktor spesifik yang dapat menyebabkan deteriorasi, yang pada akhirnya dapat mengakibatkan kerusakan suatu mesin. Dengan kata lain, dapat diketahui faktor-faktor apa saja yang berpotensi menyebabkan deteriorasi, yang pada akhirnya berakibat pada kerusakan suatu mesin tertentu. Selanjutnya, berdasarkan hasil analisis ini, model dibentuk dengan pendekatan sains data dan matematika untuk memprediksi kerusakan suatu mesin tertentu berdasarkan nilai-nilai yang terdapat pada berbagai variabel yang berhubungan dengan mesin tersebut.



## 2. LANDASAN TEORI

Karya ilmiah ini ditulis berdasarkan hasil penelitian ini yang telah dilakukan melalui pendekatan sains data dan matematis. Peneliti menggunakan pendekatan sains data untuk langkah pra-pemrosesan data guna membentuk data sesuai dengan keperluan penelitian ini. Setelah mengolah data, peneliti menggunakan pendekatan matematis untuk menguji dan menafsirkan hasil olah data tersebut. Segala tindakan yang dilakukan selama proses penelitian dapat diketahui melalui karya ilmiah ini.

### 2.1. Uji Kolmogorov-Smirnov Satu Sampel

Uji Kolmogorov-Smirnov (Колмогоров–Смирнов) atau uji KS Satu Sampel adalah uji statistika yang digunakan untuk mengidentifikasi distribusi peluang tertentu dari sampel tersebut. Hipotesis nol ( $H_0$ ) menyatakan bahwa sampel berasal dari distribusi tertentu, sedangkan hipotesis alternatif ( $H_a$ ) menyatakan sebaliknya [7]. Dalam penelitian ini, fitur yang tidak mendukung hipotesis nol atau tidak berasal dari distribusi normal akan lebih lanjut dianalisis dengan menggunakan Uji Kruskal Wallis dalam mengidentifikasi signifikansi hubungan dengan variabel dependen.

### 2.2. Uji Kruskal Wallis

Uji Kruskal Wallis adalah uji statistika non-parametrik yang digunakan untuk membandingkan tiga atau lebih kelompok data independen untuk mengetahui perbedaan statistik signifikan antar kelompok data ini. Uji ini adalah ekstensi dari Uji Mann-Whitney U yang membandingkan dua kelompok data yang berbeda. Uji ini dapat digunakan ketika asumsi uji ANOVA (*Analysis of variance*) tidak berhasil untuk dipenuhi, yaitu ketika data tidak berdistribusi normal atau data tidak bernilai ordinal [8].

### 2.3. Uji Chi-Square

Uji Chi-Square ( $\chi^2$ ) adalah uji statistika non-parametrik yang pada umumnya digunakan untuk menguji asosiasi antara variabel-variabel kategorikal. Hal ini terjadi karena uji ini mampu menentukan hubungan pada frekuensi-frekuensi yang diobservasi dan pada frekuensi-frekuensi harapan berdasarkan hipotesis tertentu [9]. Jenis Uji Chi-Square yang relevan dalam

penelitian ini adalah Uji Chi-Square untuk menentukan kebebasan atau independensi dalam menentukan asosiasi antara dua variabel kategorikal.

#### 2.4. *Variance Inflation Factor (VIF)*

*Variance Inflation Factor* atau umumnya disebut VIF adalah salah satu metode yang digunakan untuk mendeteksi multikolinearitas. Multikolinearitas menjadi masalah yang umum pada proses penelitian sains data, terutama ketika suatu variabel bebas dapat berkorelasi dengan variabel-variabel lainnya pada kelompok data tertentu. Masalah ini mengakibatkan prediksi yang tidak akurat serta menyebabkan misinterpretasi pada model yang digunakan. Metode VIF mampu mendeteksi multikolinearitas dengan cara mengukur seberapa banyak nilai variansi dari koefisien regresi meningkat karena korelasi antar variabel independen [10].

#### 2.5. *Principal Component Analysis (PCA)*

*Principal Component Analysis* atau disebut PCA, adalah metode statistika yang menerapkan transformasi ortogonal untuk mengubah kelompok variabel yang saling berkorelasi menjadi kelompok variabel yang tidak saling berkorelasi. Metode ini bekerja dengan cara memetakan data di ruang dimensi yang tinggi ke data di ruang dimensi yang rendah. PCA digunakan untuk mengurangi jumlah dimensi dari suatu kelompok data tertentu dengan mempertahankan hubungan antar variabel yang terdapat pada kelompok data tertentu. PCA akan dilakukan tanpa *input* informasi apapun mengenai variabel-variabel dependen yang menjadi target [11].

#### 2.6. Nilai *SHAP*

Nilai *Shapley Additive Explanations* atau umumnya disebut SHAP adalah salah satu *framework* yang dapat digunakan untuk menafsirkan hasil *output* dari berbagai model pembelajaran mesin [12]. SHAP adalah salah satu teknik untuk menjelaskan model dengan dasar teoritis yang baik [13]. SHAP mampu membantu peneliti dengan menyediakan informasi rinci mengenai peran setiap fitur terhadap prediksi yang dihasilkan oleh suatu model. SHAP dapat memperbaiki nilai-nilai signifikansi dengan melakukan analisis terhadap hubungan antara berbagai kovariat dengan *output* dari model.

## 2.7. Random Forest

*Decision Tree* terdiri dari satu *root node*, beberapa *internal node* dan beberapa *leaf node*. Setiap *tree* hanya akan memiliki satu *root node* sebagai tempat mulainya *tree* ini dalam melakukan percabangan untuk menentukan hasil. *Root node* akan selalu berada pada tingkat ke nol untuk setiap *tree*. Berjalan dari *root node*, setiap *tree* akan melakukan percabangan ke salah satu dari *internal node* yang ada pada tingkat selanjutnya. Ketika satu *internal node* sudah tercapai, hal yang sama akan dilakukan secara terus-menerus hingga *internal node* pada tingkat ke  $n-1$  sudah tercapai. Ketika salah satu *leaf node* pada tingkat  $n$  telah tercapai, *tree* akan menentukan hasil berdasarkan *leaf node* yang telah terpilih [14].

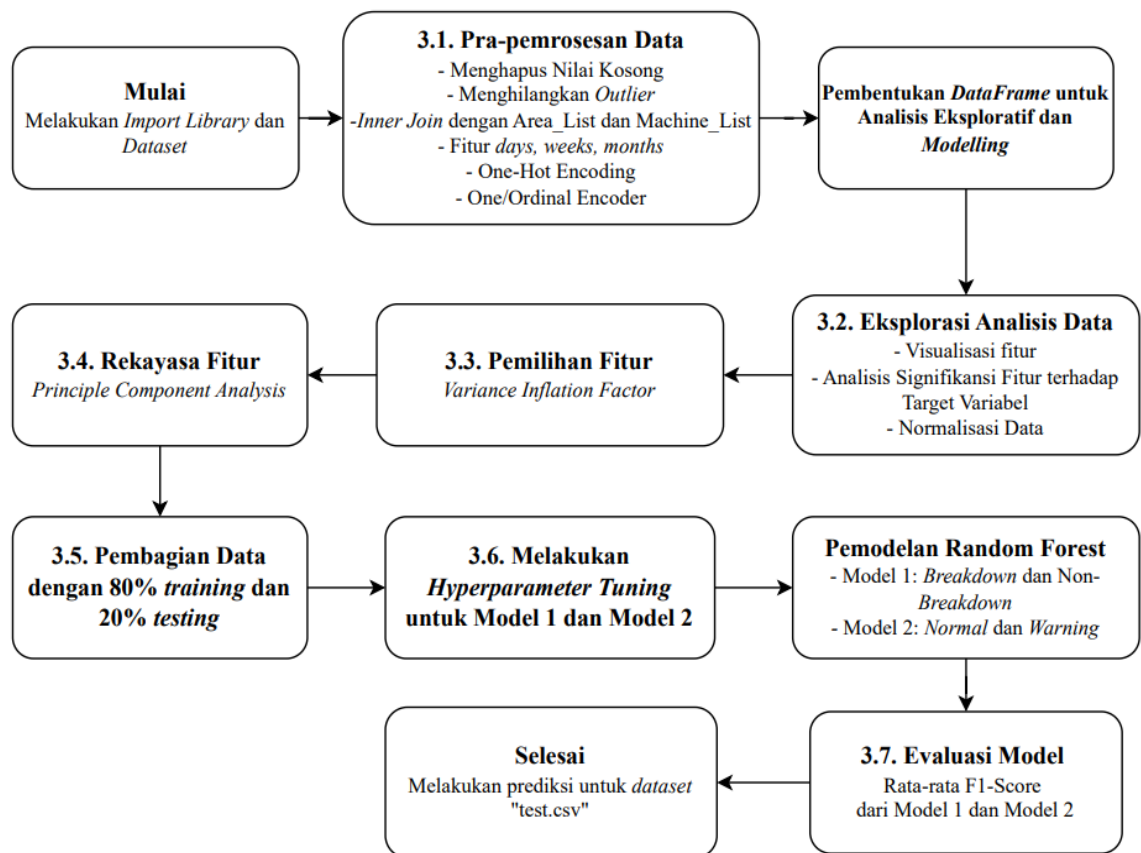
*Random Forest* adalah salah satu model pembelajaran mesin dan modifikasi dari algoritma *Decision Tree* dalam melakukan prediksi klasifikasi dan regresi. *Random Forest* bekerja dengan cara memanfaatkan *Decision Tree* individu dalam jumlah yang besar. Model ini akan mengambil bagian-bagian berbeda dari kelompok data secara acak untuk melatih setiap *tree* yang ada. Untuk mendapatkan hasil akhir, model ini akan mengambil nilai mayoritas atau rata-rata dari hasil proses setiap *tree* yang telah digabungkan [15].



### 3. PROSES ANALISIS

Proses analisis dalam penelitian ini mencakup beberapa tahap utama, yaitu pra-pemrosesan data, eksplorasi data, pemilihan fitur, rekayasa fitur, dan pemodelan. Seluruh analisis dilakukan menggunakan Jupyter Notebook Python sebagai platform utama. Alur lengkap dari proses analisis ini disajikan pada **Gambar 1**.

*Gambar 1. Diagram Alur Proses Analisis*



#### 3.1. Pra-pemrosesan Data

Proses pra-pemrosesan data dimulai dengan menghapus observasi atau baris yang memiliki nilai kosong pada fitur. Pendekatan ini dilakukan tanpa menggunakan imputasi data untuk menghindari potensi peningkatan *noise* yang memengaruhi kualitas analisis dan prediksi. Kemudian, data pencilan atau *outlier* diidentifikasi dan dihapus berdasarkan nilai absolut *Z-score* yang lebih besar atau sama dengan tiga. Selanjutnya, dilakukan *inner join* dengan *dataset machine\_list* dan *area\_list* untuk menambahkan fitur baru, yaitu

“Mesin,” “Country”, “Area”, dan “Priority.” Selain itu, fitur baru “days,” “weeks,” dan “months” dibuat dengan menghitung rentang waktu dari fitur “timestamp”, menggunakan tanggal 10 Januari 2025 sebagai batas acuan, yaitu hari dimulainya kompetisi ANAVA DataVers. Adapun fitur “StatusI” yang menandakan status “Breakdown” dan “Non-Breakdown”. Untuk data kategori, dilakukan *One Hot Encoding* pada fitur nominal dan *One Encoding* pada fitur ordinal atau fitur yang nilainya memiliki makna berurut atau *ranking*. Terakhir, tahap normalisasi data dilakukan khususnya setelah melakukan Eksplorasi Analisis Data untuk memastikan konsistensi skala data serta meningkatkan performa model.

### 3.2. Eksplorasi Analisis Data

Proses analisis melibatkan visualisasi data dan uji signifikansi antara setiap fitur dengan variabel “Status” untuk mengevaluasi asosiasi secara langsung dan mempermudah pemilihan fitur. Analisis ini dilakukan menggunakan Estimasi Parametrik untuk data fitur yang berdistribusi normal dan Estimasi Non-Parametrik untuk data yang tidak berdistribusi normal. Untuk data bertipe kategorikal, digunakan Uji Chi-Square, sedangkan untuk data numerik atau berskala kontinu digunakan Uji Kruskal-Wallis. Selain itu, pengujian normalitas data dilakukan dengan menggunakan Uji Kolmogorov-Smirnov Satu Sampel untuk memastikan distribusi dan normalitas data sebelum analisis lebih lanjut.

### 3.3. Pemilihan Fitur

Proses ini menggunakan *Variance Inflation Factor* (VIF) untuk menguji multikolinearitas antar fitur bertipe numerik atau kontinu. Hal ini bertujuan untuk mengidentifikasi dan menghapus fitur-fitur dengan nilai VIF lebih besar dari 5, karena fitur tersebut cenderung memiliki korelasi tinggi dengan fitur lainnya. Dengan membuang fitur-fitur tersebut, redundansi dalam *dataset* berkurang sehingga meningkatkan efisiensi model serta memastikan bahwa setiap fitur yang tersisa memberikan kontribusi secara identik terhadap prediksi model. Pendekatan ini juga membantu meningkatkan ketepatan dan interpretabilitas model.



### 3.4. Rekayasa Fitur

Untuk mengurangi dimensi *dataset*, penelitian ini mengimplementasikan *Principal Component Analysis* (PCA). Pendekatan ini mereduksi jumlah fitur dalam *dataset* menjadi  $n$  komponen utama yang ditentukan berdasarkan kontribusi kumulatif variansi data. Jumlah komponen utama dipilih sedemikian rupa sehingga menjelaskan setidaknya 80% dari total variansi data sebelumnya, tanpa adanya kehilangan informasi secara signifikan yang relevan untuk analisis.

### 3.5. Pembagian Data *Training* dan Data *Testing*

Penelitian ini menerapkan *stratified data splitting* untuk memastikan dan mempertahankan proporsi antar kelas atau label dalam data pelatihan (*training*) dan data pengujian (*testing*). Pembagian data dilakukan sebanyak dua tahap, dengan komposisi 80% data untuk pelatihan dan 20% untuk pengujian. Tahap pertama bertujuan untuk memprediksi atau membedakan data dengan label “*Breakdown*” sebagai kelas dominan (setelah melakukan *pre-processing*) dan label “*Non-Breakdown*” yang mencakup label “*Normal*” dan “*Warning*”. Pada tahap kedua, pembagian dari data label “*Non-Breakdown*” dilakukan untuk memprediksi atau membedakan data dengan label “*Normal*” dan “*Warning*” menggunakan himpunan data “*Non-Breakdown*” dari tahap sebelumnya.

### 3.6. Pemodelan dan *Hypertuning*

Proses prediksi ini dilakukan melalui dua tahap model. Model pertama digunakan untuk membedakan data berlabel “*Breakdown*” dan “*Non-Breakdown*”, sedangkan model kedua digunakan untuk membedakan data berlabel “*Normal*” dan “*Warning*”. Algoritma yang digunakan dalam penelitian ini adalah *Random Forest Classifier* yang memanfaatkan metode *ensemble learning* atau melibatkan beberapa prediktor (*decision trees*). Proses *hypertuning parameter* dilakukan dengan mencari kombinasi parameter paling optimal dari daftar yang telah ditentukan, yaitu jumlah pohon ( $n\_estimators$ ), kedalaman pohon ( $max\_depth$ ), jumlah data minimum pada suatu node supaya node tersebut dapat dibagi ( $min\_samples\_split$ ), jumlah data minimum pada masing-masing *leaf* supaya node dapat dibagi

(*min\_sample\_leaf*). Tujuan dari pendekatan ini adalah untuk memaksimalkan tingkat akurasi prediksi atau performa model.

### 3.7. Evaluasi Model

Evaluasi performa model dalam penelitian ini menggunakan *F1-Score* sebagai metrik utama karena dapat mengukur ketepatan prediksi pada data yang tidak seimbang dengan mempertimbangkan keseimbangan presisi dan sensitivitas. Dengan menggunakan *F1-Score*, evaluasi model berfokus pada kemampuan untuk meminimalkan kesalahan dalam prediksi antar kelas. Rata-rata *F1-Score* dari prediksi model pertama (label “*Breakdown*” dan “*Non-Breakdown*”) dan model kedua (label “*Normal*” dan “*Warning*”) digunakan untuk menjelaskan atau memberikan gambaran keseluruhan dari performa model dalam penelitian ini.

## 4. HASIL ANALISIS

### 4.1. Pra-pemrosesan Data

Jumlah data setelah menghapuskan data kosong adalah **163.078 observasi**, dengan jumlah data setelah penghapusan *outlier* adalah **133.703 observasi**. Setelah menerapkan *inner join* dengan *dataset machine\_list* dan *area\_list*, pembuatan fitur “*days*”, “*weeks*”, dan “*months*”, *One Encoding* pada fitur “*Power\_Backup*”, “*Priority*” dan “*Status*”, serta *One Hot Encoding* pada fitur “*Area*”, “*Country*”, dan “*Machine*”, *dataset* ini memiliki kolom-kolom tambahan sehingga secara keseluruhan terdapat **96 kolom**, termasuk variabel dependen. Contoh representasi *output* nilai-nilai data sebelum dinormalisasikan dan setelah dinormalisasikan disajikan pada **Lampiran**.

### 4.2. Eksplorasi Analisis Data

Berdasarkan visualisasi data (*bar chart* dan *histogram*) dari *dataset* yang telah dibersihkan, mayoritas mesin memiliki status “*Breakdown*”, yang menunjukkan ada proporsi kelas *imbalanced*. Tipe mesin yang paling dominan adalah “*Formax*”. Negara dengan jumlah produksi mesin terbanyak adalah China, sementara negara dengan produksi paling sedikit adalah Taiwan. Dari segi area, mesin paling banyak ditemukan di daerah JGJ dan BNTN, sedangkan area dengan jumlah mesin paling sedikit adalah SLJA dan SDA. Adapun usia mesin yang diproduksi bervariasi antara 0 hingga 24 bulan, dengan rentang usia lainnya antara 9 hingga 740 hari. Grafik visualisasi data dan deskripsi statistik lebih lanjut disajikan pada **Lampiran** (bagian 4.2.1 sampai 4.2.6).

Berdasarkan distribusi data melalui Uji Kolmogorov-Smirnov Satu Sampel, ditemukan bahwa semua fitur numerik tidak berdistribusi normal. Dari segi hubungan antar fitur dengan “*Status*”, hasil Uji Kruskal-Wallis menunjukkan bahwa fitur numerik yang memiliki asosiasi dengan “*Status*” adalah RPM-1 dengan selang kepercayaan 90% dengan *p-value* sebesar 0,05131, sedangkan dari hasil Uji Chi-Square, fitur kategorikal yang memiliki asosiasi dengan “*Status*” adalah nama mesin dan negara mesin yang diproduksi dengan masing-masing *p-value* sebesar 0,05423 untuk selang kepercayaan 90% dan



0,02998 untuk selang kepercayaan 95%. Karena hanya terdapat tiga fitur yang memiliki hubungan signifikan terhadap “*Status*”, mayoritas data berdasarkan “*Status*” tersebar secara acak pada fitur-fitur lainnya.

Selain dari hasil uji statistik yang menunjukkan sebagian besar fitur lainnya tidak memiliki asosiasi signifikan terhadap “*Status*”, secara visual juga memaparkan bahwa distribusi antar fitur berdasarkan label “*Status*” tidak menunjukkan perbedaan sebagaimana terlihat pada grafik di **Lampiran** (bagian 4.2.7). Hal ini menegaskan bahwa pola hubungan antara fitur dan “*Status*” perlu dieksplorasi lebih lanjut melalui pendekatan rekayasa fitur dan pemilihan fitur menggunakan metode VIF. Penelitian ini juga tidak menghapus beberapa fitur lainnya yang tidak signifikan terhadap “*Status*” supaya mempertahankan makna keseluruhan dari suatu observasi.

#### 4.3. Pemilihan Fitur melalui VIF

Berdasarkan hasil uji VIF, ditemukan bahwa fitur “*days*” dan “*weeks*” memiliki tingkat multikolinearitas tinggi terhadap fitur-fitur lainnya ( $VIF > 5$ ) sehingga *dataset* memiliki **94 kolom**. Oleh karena itu, kedua fitur tersebut tidak digunakan sebagai indikator prediksi. Hasil uji VIF setelah menghapus kedua fitur tersebut disajikan pada **Tabel 1** dan akan ditetapkan menjadi indikator prediksi.

*Tabel 1. Hasil VIF pada Fitur Numerik*

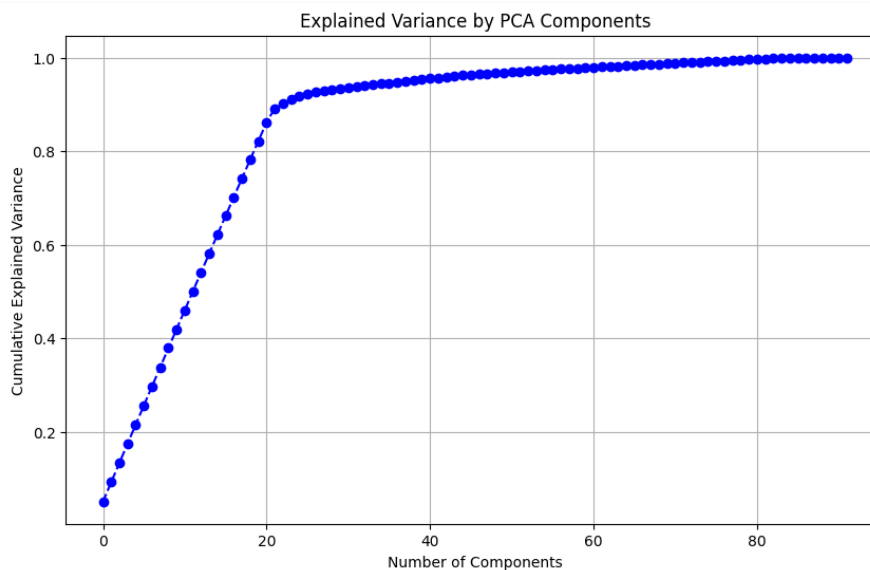
No.	Fitur Numerik	VIF Akhir
1	temperature_10H_max (°C)	1,074523
2	temperature_10H_min (°C)	1,074461
3	temperature-1	1,000195
4	temperature-2	1,000192
5	temperature-3	1,000210
6	apparent_temperature_max	1,000176
7	apparent_temperature_min	1,000177
8	humidity	1,000198
9	Voltage-L	1,000128
10	Voltage-R	1,000104

11	Voltage-M	1,000111
12	Current-M	1,000158
13	Current-R	1,000133
14	Current-T	1,000167
15	RPM	1,000209
16	RPM-1	1,000181
17	RPM-2	1,000120
18	RPM-3	1,000108
19	Vibration-1	1,000110
20	Vibration-2	1,000067
21	Power	1,000176
22	Months	1,000152

#### 4.4. Rekayasa Fitur melalui PCA

Berdasarkan variansi kumulatif, jumlah komponen minimum untuk merepresentasikan 80% dari data adalah 20 komponen, sehingga seluruh fitur akan ditransformasi dari 92 fitur menjadi 20 fitur berupa PC1 sampai PC20. Grafik peningkatan variansi kumulatif berdasarkan jumlah komponen disajikan pada **Gambar 4**, serta contoh representasi nilai-nilai disajikan pada **Gambar 5**.

*Gambar 4. Grafik Jumlah Komponen terhadap Variansi Kumulatif*



Gambar 5. Hasil Transformasi PCA Dataset (*pca\_df*)

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10
0	2.373799	-0.297787	-0.901328	-1.258437	0.647345	0.204878	-0.931182	-0.912180	-0.504509	0.171421
1	-0.180669	0.960748	0.726794	-0.209655	1.602245	1.960067	0.812967	-0.662528	0.734368	-2.306839
2	0.207715	0.983052	-0.832967	-0.613272	0.963920	-0.634680	2.690692	-1.209402	-0.037445	-0.897856
3	0.069907	1.416061	0.856010	0.837557	-0.345037	-0.992348	0.229582	-0.457033	-0.298375	1.409099
4	-0.462163	0.176828	1.035365	-0.670833	0.336980	-0.716499	-0.903791	0.797164	-0.825100	-0.953816

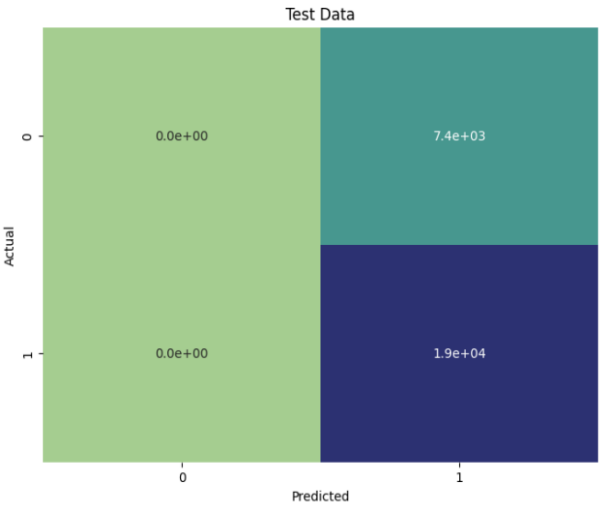
  

	PC11	PC12	PC13	PC14	PC15	PC16	PC17	PC18	PC19	PC20	Status
0	-0.731324	-0.725640	1.375104	1.139741	0.583154	1.186383	-0.708040	-0.068422	0.521952	-0.609302	2
1	1.042076	0.476950	0.005122	0.989776	0.301895	2.046232	0.090130	-0.334088	0.015959	-1.302072	0
2	-0.009786	1.284706	-1.314232	0.573685	0.291695	0.990441	0.341007	1.296999	-0.337940	0.146195	0
3	-0.493833	-0.494935	0.143667	0.285287	-0.898482	2.132934	-0.068479	-1.460945	-0.180344	-0.220704	0
4	-0.330280	-0.295634	-1.932337	0.181884	1.240509	-0.423963	0.381743	0.512540	0.371846	0.610062	0

#### 4.5. Evaluasi Performa Model Pertama (“Breakdown” dan “Non-Breakdown”)

Berdasarkan hasil *hyperparameter tuning*, parameter paling optimal adalah *n\_estimators* = 100, *max\_depth* = 10, *min\_samples\_split* = 2, *min\_samples\_leaf* = 1, dan *criterion* = ‘entropy’. Model ini memperoleh *F1-Score* sebesar **83,95%** untuk memprediksi data yang berlabel “Breakdown”. Hasil prediksi label pada data *testing* menunjukkan bahwa jumlah *true positive* adalah **19.344 observasi**, sedangkan jumlah *false positive* adalah **7.397 observasi**. Perbandingan antara data prediksi terhadap data sebenarnya disajikan pada **Gambar 6**.

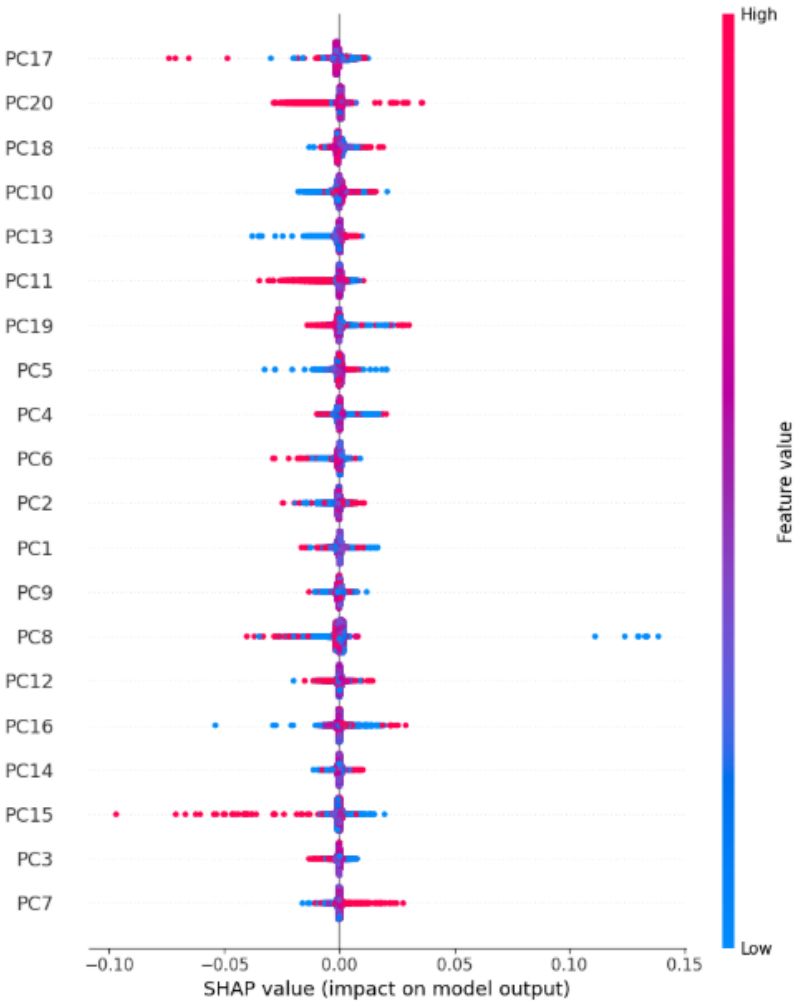
Gambar 6. Confusion Matrix Data Testing pada Model Pertama





Selanjutnya, penelitian ini juga memaparkan hubungan antar fitur dan variabel dependen (“*Breakdown*” dan “*Non-Breakdown*”) melalui visualisasi nilai SHAP, yang setiap peningkatan nilai SHAP meningkatkan kemungkinan data berlabel “*Breakdown*”. Pada **Gambar 7**, nilai pada *vertical axis* merujuk pada besaran nilai dari setiap fitur, dalam kata lain titik berwarna biru memiliki nilai fitur yang rendah, sedangkan titik berwarna merah memiliki nilai fitur yang tinggi.

Gambar 7. SHAP Value Data Testing pada Model 1



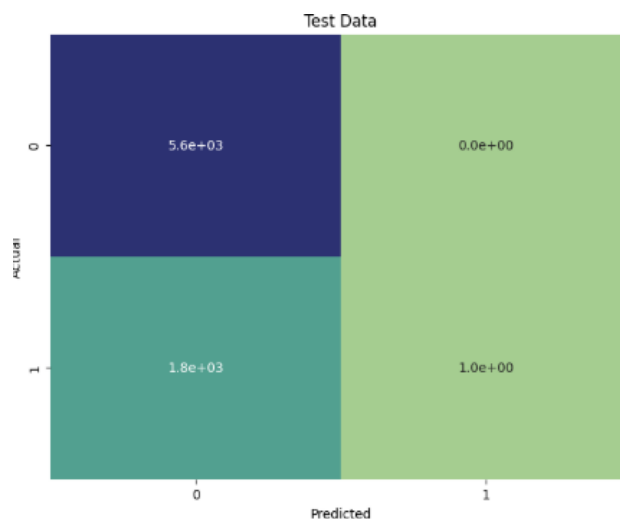
Berdasarkan nilai SHAP, algoritma dari model pertama mempelajari bahwa setiap kenaikan nilai pada fitur PC3, PC6, PC8, PC11, PC15 dan PC17 cenderung menurunkan kemungkinan bahwa mesin tersebut mengalami kerusakan (“*Breakdown*”), sedangkan setiap kenaikan nilai pada fitur PC2, PC7, PC10, PC13, PC14 dan PC16 cenderung meningkatkan kemungkinan

bahwa mesin tersebut mengalami kerusakan. Namun, fitur PC1, PC4, PC5, PC9, PC12, PC18, PC19 dan PC20 tidak menunjukkan hubungan terhadap “Status” karena distribusi nilai-nilai antar fitur tersebar pada  $x$ -axis negatif dan positif.

#### 4.6. Evaluasi Performa Model Kedua (“Normal” dan “Warning”)

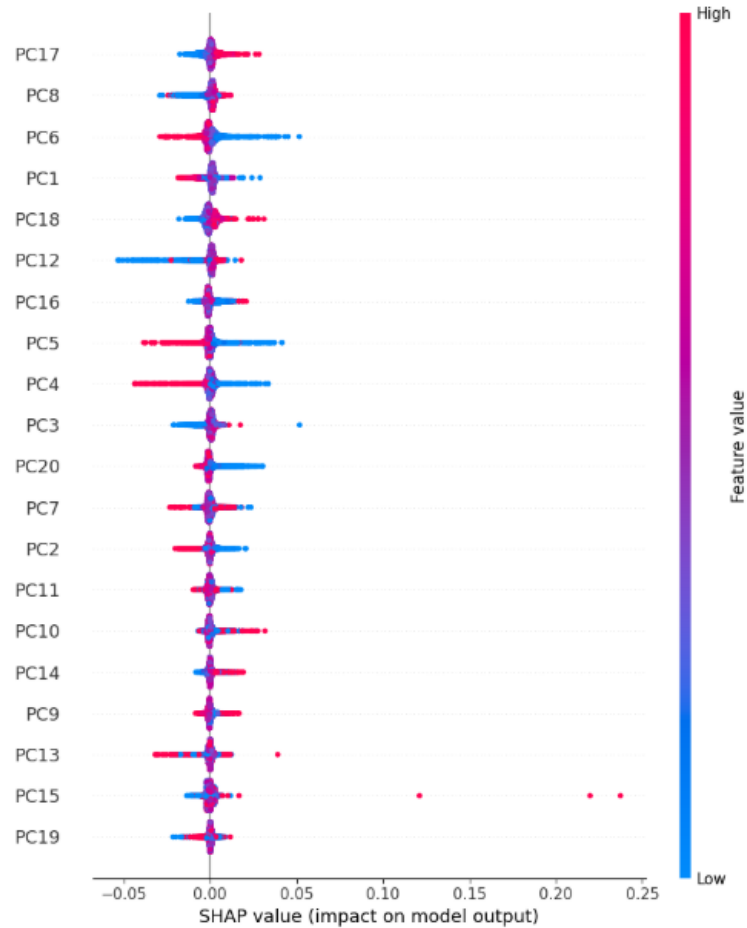
Berdasarkan hasil *hyperparameter tuning*, parameter paling optimal adalah  $n\_estimators = 100$ ,  $max\_depth = 10$ ,  $min\_samples\_split = 2$ ,  $min\_samples\_leaf = 1$ , dan  $criterion = 'entropy'$ . Model ini memperoleh **F1-Score** sebesar **0,11%** untuk memprediksi data yang berlabel “Warning” (label 1), serta **F1-Score** sebesar **86,08%** untuk memprediksi data yang berlabel “Normal” (label 0). Hasil prediksi menunjukkan bahwa jumlah data yang benar diklasifikasi adalah **5.589 observasi**, sedangkan jumlah data berlabel “Warning” yang telah diprediksi sebagai label “Normal” adalah **1.807 observasi**. Perbandingan antara data prediksi terhadap data sebenarnya disajikan pada **Gambar 8**.

Gambar 8. Confusion Matrix Data Testing pada Model Pertama



Selanjutnya, penelitian ini juga memaparkan hubungan antar fitur dan variabel dependen (“Normal” dan “Warning”) melalui visualisasi nilai SHAP, yang setiap peningkatan nilai SHAP meningkatkan kemungkinan data berlabel “Warning”. Pada **Gambar 9**, nilai pada *vertical axis* merujuk pada besaran nilai dari setiap fitur, dalam kata lain titik berwarna biru

memiliki nilai fitur yang rendah, sedangkan titik berwarna merah memiliki nilai fitur yang tinggi.



Berdasarkan nilai SHAP, algoritma dari model kedua mempelajari bahwa setiap kenaikan nilai pada fitur PC1, PC2, PC4, PC5, PC6, PC11 dan PC20 cenderung menurunkan kemungkinan bahwa mesin tersebut mengalami pertanda kerusakan (“Warning”), sedangkan setiap kenaikan nilai pada fitur PC8, PC10, PC12, PC14, PC15, PC16, PC17 dan PC18 cenderung meningkatkan kemungkinan bahwa mesin tersebut mengalami pertanda kerusakan. Selanjutnya, setiap penurunan nilai PC3 dan PC19 cenderung meningkatkan kemungkinan bahwa mesin terkategori sebagai “Normal”. Namun, fitur PC7, PC9 dan PC13 tidak menunjukkan hubungan terhadap “Status” karena distribusi nilai-nilai antar fitur tersebar pada  $x$ -axis negatif dan positif.

Secara keseluruhan, rata-rata  $F1$ -Score pada kedua model adalah **42.03%**.



## 5. KESIMPULAN DAN REKOMENDASI

Dalam penelitian ini, upaya deteksi deteriorisasi mesin diimplementasi dengan menggunakan algoritma *Random Forest* dengan parameter yang telah dioptimasi melalui *hyperparameter tuning*. Hasil prediksi menunjukkan bahwa secara keseluruhan performa model antara data pelatihan dan data pengujian mencapai *F1-Score* sebesar 42.03% karena mayoritas fitur tidak memiliki asosiasi secara statisitka dengan variabel dependen, ketidakseimbangan kelas, serta terdapat pengaruh nilai SHAP yang saling kontradiksi antar komponen (PCA) dalam hasil prediksi model.

Berdasarkan hasil analisis data melalui uji statistika, ditemukan bahwa variabel RPM-1, jenis mesin dan negara produksi mesin memiliki hubungan sigmifikan terhadap status kerusakan mesin. Dari sudut pandang *machine learning*, fitur-fitur yang telah ditransformasi ke dalam komponen PCA yang berkontribusi terhadap peningkatan potensi kerusakan mesin, yaitu label “Warning” dan “Breakdown”, adalah PC2, PC7, PC8, PC10, PC12, PC13, PC14, PC15, PC16, PC17, PC18.

Dengan adanya kontribusi dari penelitian ini dalam memberikan deteksi deteriorasi mesin sejak dini melalui pendekatan sains data dan matematis, diharapkan dapat mendukung sektor ekonomi dengan cara mengantisipasi kerusakan mesin, meningkatkan efisiensi produksi serta memilih mesin produksi yang lebih kredibel sehingga mengoptimalkan kinerja operasional. Untuk meningkatkan performa model dari penelitian ini, tindakan lanjut dapat dilakukan dengan cara mengoptimasikan model lebih lanjut untuk deteksi “Warning” dan “Normal” melalui algoitma *Computer Vision*, LSTM, atau *transformer* untuk mengolah data operasional mesin tambahan yang berbasis waktu dan gambar seperti data lingkungan operasional mesin. Selain itu, penelitian ini dapat dikembangkan dengan melibatkan fitur-fitur komponen mesin lain yang secara statistik signifikan terhadap status deteriorisasi mesin serta menganalisis relevansi fitur-fitur dengan hasil transformasi dalam bentuk PCA.

## DAFTAR PUSTAKA

- [1] J. Baldi and V. Roux, "The innovation of the potter's wheel: a comparative perspective between Mesopotamia and the southern Levant," *Levant*, vol. 48, no. 3, pp. 236–253, Sep. 2016, doi: 10.1080/00758914.2016.1230379.
- [2] A. N. Safronov, "Antikythera Mechanism and the Ancient World," *Journal of Archaeology*, vol. 2016, pp. 1–19, May 2016, doi: 10.1155/2016/8760513.
- [3] B. D. Solomon and K. Krishna, "The coming sustainable energy transition: History, strategies, and outlook," *Energy Policy*, vol. 39, no. 11, pp. 7422–7431, Nov. 2011, doi: 10.1016/j.enpol.2011.09.009.
- [4] G. C. Lin and D.-C. Gong, "On a production-inventory system of deteriorating items subject to random machine breakdowns with a fixed repair time," *Math Comput Model*, vol. 43, no. 7–8, pp. 920–932, Apr. 2006, doi: 10.1016/j.mcm.2005.12.013.
- [5] M. Al-Salamah and A. Abudari, "Production Lot Sizing and Process Targeting under Process Deterioration and Machine Breakdown Conditions," *Modelling and Simulation in Engineering*, vol. 2012, pp. 1–7, 2012, doi: 10.1155/2012/393495.
- [6] P. W. Tse and D. P. Atherton, "Prediction of Machine Deterioration Using Vibration Based Fault Trends and Recurrent Neural Networks," *J Vib Acoust*, vol. 121, no. 3, pp. 355–362, Jul. 1999, doi: 10.1115/1.2893988.
- [7] S. Facchinetti, "A Procedure to Find Exact Critical Values of Kolmogorov-Smirnov Test," *Statistica Applicata*, vol. 21, pp. 337–359, 2009.
- [8] P. E. McKight and J. Najab, "Kruskal-Wallis Test," in *The Corsini Encyclopedia of Psychology*, Wiley, 2010, pp. 1–1. doi: 10.1002/9780470479216.corpsy0491.

- [9] M. L. McHugh, "The Chi-square test of independence," *Biochem Med (Zagreb)*, pp. 143–149, 2013, doi: 10.11613/BM.2013.018.
- [10] R. M. O'brien, "A Caution Regarding Rules of Thumb for Variance Inflation Factors," *Qual Quant*, vol. 41, no. 5, pp. 673–690, Sep. 2007, doi: 10.1007/s11135-006-9018-6.
- [11] H. Abdi and L. J. Williams, "Principal component analysis," *WIREs Computational Statistics*, vol. 2, no. 4, pp. 433–459, Jul. 2010, doi: 10.1002/wics.101.
- [12] N. Basheer, B. Pranggono, S. Islam, S. Papastergiou, and H. Mouratidis, "Enhancing Malware Detection Through Machine Learning Using XAI with SHAP Framework," 2024, pp. 316–329. doi: 10.1007/978-3-031-63211-2\_24.
- [13] A. Janssen, M. Hoogendoorn, M. H. Cnossen, and R. A. A. Mathôt, "Application of <scp>SHAP</scp> values for inferring the optimal functional form of covariates in pharmacokinetic modeling," *CPT Pharmacometrics Syst Pharmacol*, vol. 11, no. 8, pp. 1100–1110, Aug. 2022, doi: 10.1002/psp4.12828.
- [14] J. R. Quinlan, "Learning decision tree classifiers," *ACM Comput Surv*, vol. 28, no. 1, pp. 71–72, Mar. 1996, doi: 10.1145/234313.234346.
- [15] L. Breiman, "Random Forests," *Mach Learn*, vol. 45, no. 1, pp. 5–32, 2001, doi: 10.1023/A:1010933404324.



[illegible]

### 4.1.2. Hasil Representasi *Output* Setelah Tahap Normalisasi

	temperature_10H_max (°C)	temperature_10H_min (°C)	temperature-1	temperature-2	temperature-3	apparent_temperature_max	apparent_temperature_min	humidity	Voltage-L
0	-2.642552	-0.731460	-0.324868	-0.594051	-0.789549	-0.442491	0.957295	0.424456	-1.051294
1	-0.376189	0.482965	-0.972137	-0.731278	1.459644	-1.522419	0.940198	-1.318928	-2.765264
2	-0.760331	0.460477	-0.312359	0.335031	-0.253346	-0.900067	-0.685081	-0.914389	-0.049909
3	-0.420057	0.293990	0.502498	-1.086542	-0.700825	0.177077	-1.119496	0.037994	-2.158361
4	0.144814	0.488463	1.009304	1.179576	-0.320343	-0.822367	0.761294	-0.980009	1.194674

	Voltage-L	Voltage-R	Voltage-M	Current-M	Current-R	Current-T	RPM	RPM-1	RPM-2	RPM-3	Vibration-1	Vibration-2
0	-1.051294	-1.054127	0.610881	1.052089	0.059707	-0.908450	1.000818	-0.832805	-0.690146	0.203373	-0.874794	-0.950238
1	-2.765264	0.339455	-0.368923	-1.909963	0.247120	0.279544	-0.388655	-1.156049	-0.150442	-0.627337	-1.672233	1.144025
2	-0.049909	1.658259	-0.856404	-0.999453	1.021696	1.603401	-1.223584	-1.323249	-0.563717	1.215342	-1.840971	-0.133002
3	-2.158361	0.669901	0.246062	0.472269	-0.812895	-0.761831	-0.874230	-0.412245	0.324907	1.368461	0.279016	-0.520018
4	1.194674	0.526778	0.777158	0.410290	1.037966	-1.404241	0.204774	1.084648	0.000501	1.178321	0.347006	1.314161

	Power	months	Power_Backup	Priority	BFMG	BGR	BKS	BLJA	BNTN	BPN	JGJ	KDR	KLT	KRWG	LMPG	MKS
0	-1.273974	-0.001312	0.333352	0.642119	0	0	0	0	1	0	0	0	0	0	0	0
1	-0.068371	0.282334	1.000000	0.219894	0	0	0	0	0	0	0	0	0	0	0	0
2	-0.976339	0.424158	1.000000	1.000000	0	0	0	0	0	0	0	0	0	0	0	0
3	1.042962	-0.994075	1.000000	0.219894	0	0	0	0	1	0	0	0	0	0	0	0
4	-0.851325	0.140511	1.000000	0.642119	0	0	0	0	0	0	0	0	0	0	0	0

	SKBM	SMGS	SRBY	SRG	TGR	JP	KR	TW	US	Forklift	Formax	Hitech-1	Hiwell	Innotech-1	Jawfeng	Mixer	Xiaojin	NL1	NL2	NL3	Novamax
0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

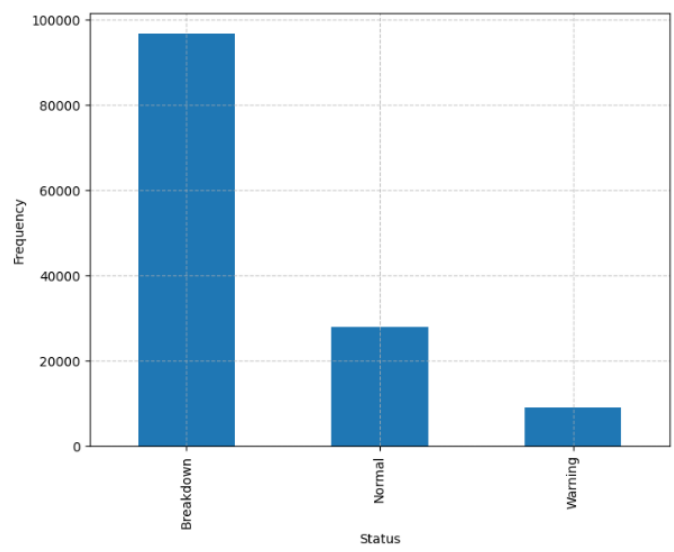
	Palette Jack	Palletizer	Palletizer-2	Palletizer-3	Palletizer-4	Palletizer-5	Palletizer-6	Plate Fomer Heidelberg	Plate Fomer Komori	Plate Fomer Revo	Plate Fomer Stork	Risco-TR130	Risco-TR200	Risco-TR300	Risco-TR400	Risco-TR500	Risco-TR600
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
3	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
4	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0

	Risco-TR700	Risco-TR800	Stuffer Linker 1	Stuffer Linker 2	Stuffer Linker 3	Stuffer Linker 4	Stuffer Linker 5	Stuffer Linker 6	Vacuum Filler	Vacuum Filler-2	Vacuum Filler-3	Vacuum Filler-4	Vacuum Filler-5
0	0	0	0	0	0	1	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0

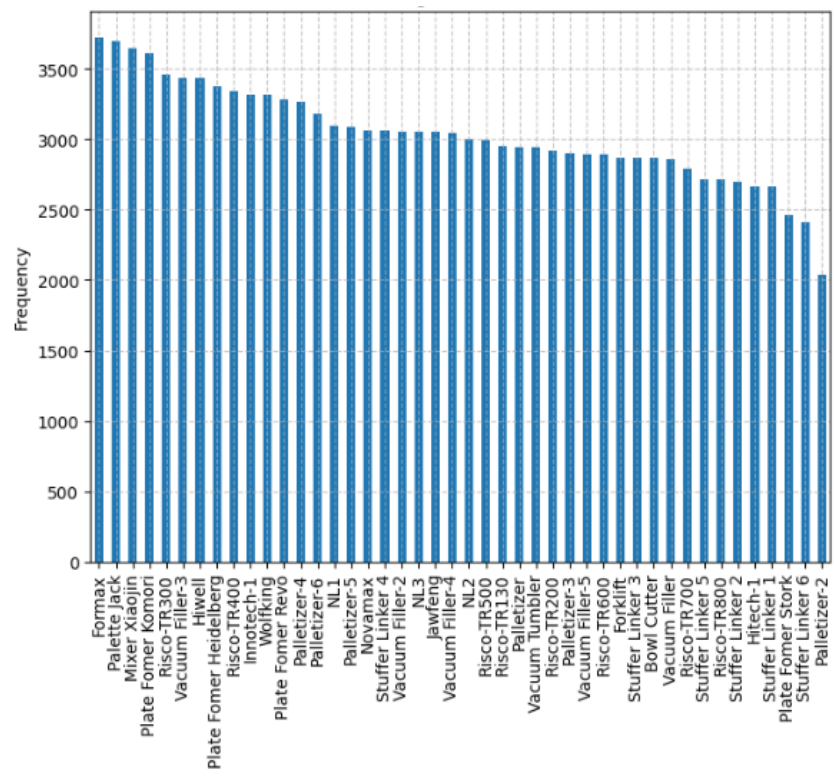


	Vacuum Tumbler	Wolfking	Status	Status1
0	0	0	2	True
1	0	0	0	False
2	0	0	0	False
3	0	0	0	False
4	0	0	0	False

#### 4.2.1. Visualisasi Distribusi Fitur “Status”

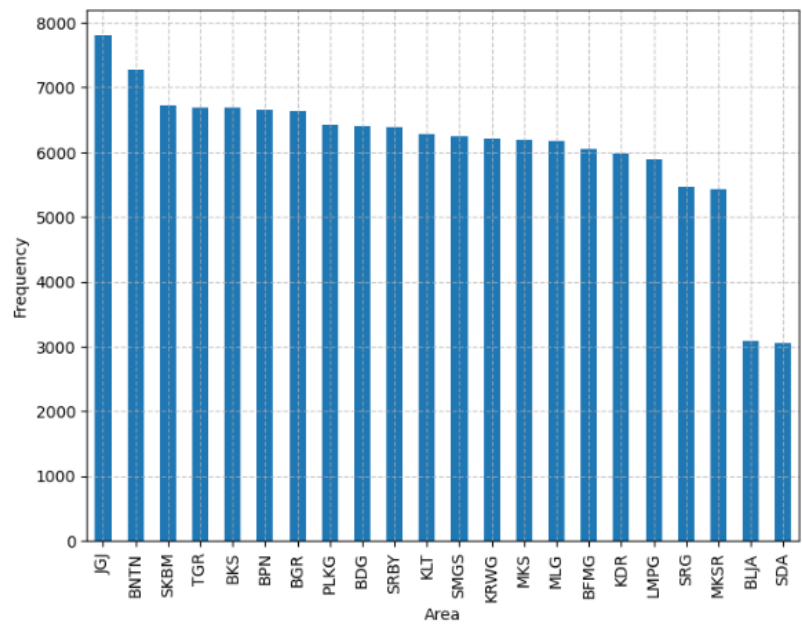


#### 4.2.2. Visualisasi Distribusi Fitur Tipe Mesin

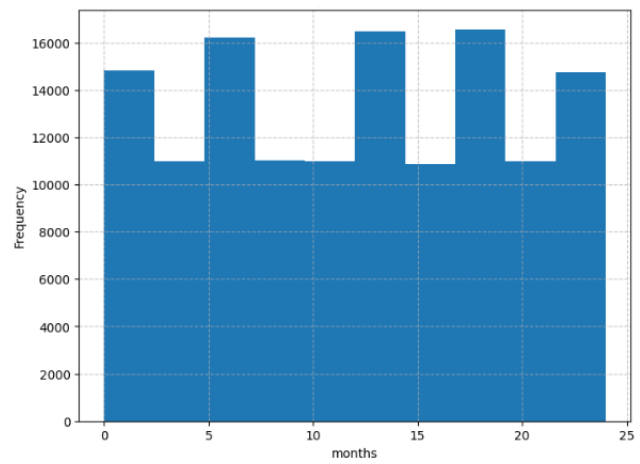




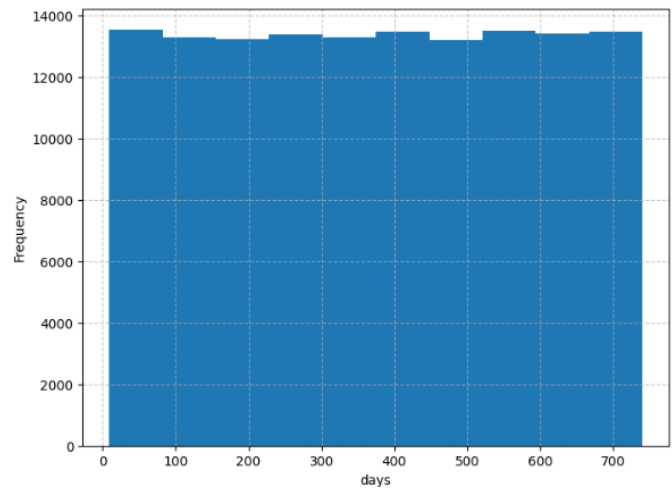
### 4.2.3. Visualisasi Distribusi Daerah Produksi Mesin



### 4.2.4. Visualisasi Distribusi Usia Mesin (Fitur “months”)



### 4.2.5. Visualisasi Distribusi Usia Mesin (Fitur “days”)



#### 4.2.6. Deskriptif Statistik *Dataset* yang Sudah Dibersihkan (*df\_cleanedEDA*)

	temperature_10H_max (°C)	temperature_10H_min (°C)	temperature-1	temperature-2	temperature-3	apparent_temperature_max	apparent_temperature_min
count	133703.000000	133703.000000	133703.000000	133703.000000	133703.000000	133703.000000	133703.000000
unique	NaN	NaN	NaN	NaN	NaN	NaN	NaN
top	NaN	NaN	NaN	NaN	NaN	NaN	NaN
freq	NaN	NaN	NaN	NaN	NaN	NaN	NaN
mean	684.858371	-14.176383	296.046655	28.327196	143.199543	28.407149	-28.381971
std	353.506250	66.058482	115.943062	15.687927	181.127327	15.679908	15.657071
min	-1259.392376	-285.688673	16.518964	0.132182	0.000189	0.012649	-76.437730
25%	474.672012	1.512306	210.058731	16.009633	30.032261	16.105883	-38.710648
50%	602.947525	12.976147	282.032671	26.334361	78.700597	26.410316	-26.349436
75%	826.837317	19.729531	368.395435	38.741708	180.303429	38.845174	-16.090835
max	2589.257322	81.005313	667.326116	76.400403	1237.911909	76.496701	-0.034805

humidity	Voltage-L	Voltage-R	...	Status	Breakdown Category	Mesin	Country Machine	Area	Priority	days	weeks	months	Status1
133703.000000	133703.000000	133703.000000	...	133703	133703	133703	133703	133703	133703	133703.000000	133703.000000	133703.000000	133703
NaN	NaN	NaN	...	3	3	44	5	22	3	NaN	NaN	NaN	2
NaN	NaN	NaN	...	Breakdown	Shutdown	Formax	CN	JGJ	Low	NaN	NaN	NaN	True
NaN	NaN	NaN	...	96719	44884	3727	33417	7803	50885	NaN	NaN	NaN	96719
37.753625	429.672211	287.947678	...	NaN	NaN	NaN	NaN	NaN	NaN	374.812465	53.113991	12.009252	NaN
7.630494	10.209004	19.647372	...	NaN	NaN	NaN	NaN	NaN	NaN	211.248672	30.180357	7.051053	NaN
22.432530	396.950366	253.720988	...	NaN	NaN	NaN	NaN	NaN	NaN	9.000000	1.000000	0.000000	NaN
31.672948	423.297694	272.149713	...	NaN	NaN	NaN	NaN	NaN	NaN	192.000000	27.000000	6.000000	NaN
36.348716	431.609765	284.195501	...	NaN	NaN	NaN	NaN	NaN	NaN	375.000000	53.000000	12.000000	NaN
42.521629	437.849858	300.176663	...	NaN	NaN	NaN	NaN	NaN	NaN	558.000000	79.000000	18.000000	NaN
62.316631	448.385732	350.935023	...	NaN	NaN	NaN	NaN	NaN	NaN	740.000000	105.000000	24.000000	NaN

#### 4.2.7. Visualisasi Distribusi Fitur-fitur Berdasarkan “Status”

