

10-ALGORITMA

Yang Diperlukan Data Scientists

Oleh:

[Prof ir Rudy C Tarumingkeng, PhD](#)

Guru Besar Manajemen, NUP: 9903252922

[Sekolah Pascasarjana, IPB-University](#)

RUDYCT e-PRESS

rudyct75@gmail.com

Bogor, Indonesia

26 Desember 2024

Berikut adalah **10 algoritma machine learning yang harus ketahui Data Scientist**, dengan penjelasan konsep, cara kerja, serta kasus penggunaannya

1. Linear Regression

Konsep: Linear Regression adalah algoritma **supervised learning** yang digunakan untuk memodelkan hubungan antara variabel independen (*features*) dan variabel dependen (*target*). Algoritma ini mencari garis lurus terbaik yang meminimalkan kesalahan kuadrat antara nilai prediksi dan nilai sebenarnya.

Cara Kerja:

- Linear Regression menggunakan rumus $y=mx+b$, di mana m adalah gradien (kemiringan) dan b adalah intercept.
- Algoritma ini memanfaatkan metode **Least Squares** untuk mengurangi error total.

Kasus Penggunaan:

- Memodelkan harga rumah berdasarkan ukuran dan lokasi.
 - Memprediksi penjualan berdasarkan data iklan.
-

2. Logistic Regression

Konsep: Meskipun namanya "regresi," Logistic Regression adalah algoritma **klasifikasi** yang digunakan untuk memprediksi probabilitas suatu peristiwa. Hasilnya diklasifikasikan ke dalam kategori (misalnya, 0 atau 1) menggunakan fungsi sigmoid.

Cara Kerja:

- Logistic Regression memanfaatkan fungsi sigmoid: $f(x) = \frac{1}{1+e^{-x}}$, yang mengubah nilai kontinu menjadi probabilitas antara 0 dan 1.
- Prediksi dibuat berdasarkan ambang batas (biasanya 0.5).

Kasus Penggunaan:

- Memprediksi apakah pelanggan akan membeli produk (ya/tidak).
 - Mendeteksi email sebagai spam atau bukan spam.
-

3. Decision Tree

Konsep: Decision Tree adalah algoritma **supervised learning** yang membagi dataset ke dalam subgrup berdasarkan fitur tertentu. Hasilnya adalah struktur seperti pohon dengan simpul keputusan dan simpul daun.

Cara Kerja:

- Algoritma ini menggunakan metrik seperti **Gini Impurity** atau **Entropy** untuk memutuskan bagaimana membagi data.
- Proses pembagian dilakukan hingga semua data dalam simpul memiliki karakteristik serupa (pure).

Kasus Penggunaan:

- Memutuskan apakah nasabah layak mendapatkan pinjaman.
 - Menentukan segmentasi pasar.
-

4. Random Forest

Konsep: Random Forest adalah ensemble learning yang menggabungkan banyak Decision Tree untuk meningkatkan akurasi dan mengurangi overfitting. Setiap pohon dilatih pada subset data yang berbeda (bagging).

Cara Kerja:

- Algoritma memilih subset data dan fitur secara acak untuk membangun banyak Decision Tree.
- Hasil akhir ditentukan berdasarkan rata-rata (untuk regresi) atau voting mayoritas (untuk klasifikasi).

Kasus Penggunaan:

- Memprediksi penyakit berdasarkan gejala.
 - Analisis churn pelanggan.
-

5. Support Vector Machine (SVM)

Konsep: SVM adalah algoritma **supervised learning** yang digunakan untuk klasifikasi dan regresi. Algoritma ini mencari *hyperplane* terbaik yang memisahkan data dalam ruang berdimensi tinggi.

Cara Kerja:

- SVM menggunakan margin terbesar antara dua kelas untuk menentukan *hyperplane* optimal.
- Untuk data yang tidak dapat dipisahkan secara linear, SVM menggunakan kernel trick untuk memetakan data ke dimensi yang lebih tinggi.

Kasus Penggunaan:

- Klasifikasi teks, seperti analisis sentimen.
 - Pengenalan pola wajah.
-

6. K-Nearest Neighbors (KNN)

Konsep: KNN adalah algoritma sederhana untuk **klasifikasi** dan **regresi**. Algoritma ini menentukan kelas suatu titik data berdasarkan mayoritas tetangga terdekatnya.

Cara Kerja:

- Algoritma menghitung jarak antara titik data baru dan semua titik data dalam dataset (misalnya, jarak Euclidean).
- Kelas ditentukan berdasarkan mayoritas dari kkk tetangga terdekat.

Kasus Penggunaan:

- Sistem rekomendasi produk.
 - Pengenalan tulisan tangan.
-

7. Naive Bayes

Konsep: Naive Bayes adalah algoritma **probabilistik** berdasarkan Teorema Bayes. Algoritma ini mengasumsikan bahwa semua fitur bersifat independen, meskipun pada kenyataannya mungkin tidak.

Cara Kerja:

- Menggunakan formula $P(C|X) = \frac{P(X|C)P(C)}{P(X)}$, di mana $P(C|X)$ adalah probabilitas kelas C diberikan fitur X .
- Keputusan dibuat berdasarkan probabilitas tertinggi.

Kasus Penggunaan:

- Deteksi email spam.
 - Klasifikasi teks, seperti analisis sentimen.
-

8. K-Means Clustering

Konsep: K-Means adalah algoritma **unsupervised learning** yang membagi data menjadi kkk cluster berdasarkan kesamaan data. Tujuannya adalah meminimalkan jarak intra-cluster.

Cara Kerja:

- Menentukan k centroid awal secara acak.
- Setiap titik data diberi label berdasarkan kedekatannya dengan centroid.
- Centroid diperbarui hingga konvergensi.

Kasus Penggunaan:

- Segmentasi pelanggan dalam pemasaran.
 - Analisis pola dalam data tanpa label.
-

9. Principal Component Analysis (PCA)

Konsep: PCA adalah algoritma **unsupervised learning** yang digunakan untuk **reduksi dimensi**. Tujuannya adalah mengurangi jumlah fitur dalam dataset sambil mempertahankan varians maksimum.

Cara Kerja:

- PCA menghitung komponen utama dengan menggunakan dekomposisi matriks (*eigenvector* dan *eigenvalue*).
- Data diproyeksikan ke ruang berdimensi rendah dengan mempertahankan informasi sebanyak mungkin.

Kasus Penggunaan:

- Visualisasi data berdimensi tinggi.
 - Pra-pemrosesan data dalam model machine learning.
-

10. Gradient Boosting (XGBoost, LightGBM)

Konsep: Gradient Boosting adalah algoritma ensemble learning yang membangun model secara iteratif dengan menyesuaikan kesalahan model sebelumnya. Contoh populer adalah **XGBoost** dan **LightGBM**.

Cara Kerja:

- Model pertama dibuat, dan kesalahan dihitung.

- Model berikutnya dilatih untuk memperbaiki kesalahan model sebelumnya.
- Proses ini berulang hingga model konvergen atau mencapai batas iterasi.

Kasus Penggunaan:

- Kompetisi machine learning (Kaggle).
 - Prediksi kredit macet.
-

Kesimpulan

Sebagai seorang Data Scientist, memahami algoritma-algoritma di atas adalah fundamental. Setiap algoritma memiliki kekuatan dan kelemahan tergantung pada jenis data dan masalah yang ingin dipecahkan. Linear Regression mungkin sederhana tetapi efektif untuk regresi dasar, sedangkan algoritma seperti Gradient Boosting menjadi pilihan dalam model prediktif yang kompleks. Memahami kapan dan bagaimana menggunakan algoritma ini akan membantu Anda memecahkan masalah data dengan lebih efektif dan efisien.

Berikut adalah penjelasan **detail masing-masing algoritma** yang telah disebutkan sebelumnya:

1. Linear Regression

- **Cara Kerja:**

- Linear Regression memodelkan hubungan linier antara variabel independen (X) dan variabel dependen (Y).
- Fungsi yang digunakan: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$, di mana:
 - β : koefisien regresi,
 - ϵ : error (kesalahan residu).
- Model mencoba meminimalkan Mean Squared Error (MSE), yaitu $\frac{1}{n} \sum (y_i - \hat{y}_i)^2$.

- **Kelebihan:**

- Sederhana dan mudah diinterpretasikan.
- Cocok untuk dataset kecil dengan hubungan linier.

- **Kekurangan:**

- Tidak efektif untuk hubungan non-linear.
- Sensitif terhadap outlier.

2. Logistic Regression

- **Cara Kerja:**

- Digunakan untuk klasifikasi biner (0 atau 1).
- Fungsi logit ($g(z)$) adalah $\frac{1}{1+e^{-z}}$, di mana $z = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n$.
- Model memprediksi probabilitas, kemudian menggunakan ambang batas (biasanya 0.5) untuk menentukan kelas.

- **Kelebihan:**

- Mudah diimplementasikan.
- Interpretasi probabilitas.

- **Kekurangan:**

- Tidak cocok untuk hubungan non-linear kecuali dengan transformasi fitur.
- Asumsi independensi fitur bisa jadi tidak realistis.

3. Decision Tree

- **Cara Kerja:**

- Mulai dari root node.
- Algoritma memutuskan atribut mana yang membagi data menjadi subset paling homogen menggunakan metrik seperti:

- Gini Impurity: $G = 1 - \sum_{i=1}^k p_i^2$,

- Entropy: $H = - \sum_{i=1}^k p_i \log(p_i)$.

-
- Pohon berkembang hingga mencapai kedalaman tertentu atau simpul menjadi homogen (pure).

- **Kelebihan:**

- Mudah dipahami dan diinterpretasikan.
- Tidak memerlukan pra-pemrosesan data (misalnya, normalisasi).

- **Kekurangan:**

- Rentan terhadap overfitting (jika tidak dipangkas).
- Sensitif terhadap perubahan kecil dalam data.

4. Random Forest

- **Cara Kerja:**

- Ensemble learning: Menggunakan banyak Decision Tree.

- Setiap pohon dilatih pada subset data dan subset fitur yang dipilih secara acak (bagging).
 - Voting mayoritas (klasifikasi) atau rata-rata prediksi (regresi) digunakan untuk menentukan hasil akhir.
 - **Kelebihan:**
 - Mengurangi overfitting dibanding Decision Tree.
 - Robust terhadap outlier.
 - **Kekurangan:**
 - Lebih lambat dibanding Decision Tree tunggal.
 - Kurang interpretatif dibanding Decision Tree.
-

5. Support Vector Machine (SVM)

- **Cara Kerja:**
 - SVM mencoba memisahkan kelas dengan margin terbesar antara data dari dua kelas.
 - Untuk data yang tidak dapat dipisahkan secara linier, **kernel trick** (linear, polynomial, radial basis function/RBF) digunakan untuk memetakan data ke dimensi yang lebih tinggi.
 - **Kelebihan:**
 - Cocok untuk dataset berdimensi tinggi.
 - Efektif untuk dataset kecil dan data tidak linier.
 - **Kekurangan:**
 - Lambat untuk dataset besar.
 - Sulit memilih kernel yang tepat.
-

6. K-Nearest Neighbors (KNN)

- **Cara Kerja:**

- Hitung jarak (misalnya Euclidean) antara data baru dengan semua data yang ada.
 - Pilih kkk tetangga terdekat.
 - Klasifikasi dilakukan berdasarkan mayoritas kelas tetangga.
 - **Kelebihan:**
 - Mudah diimplementasikan.
 - Non-parametrik, sehingga cocok untuk pola yang kompleks.
 - **Kekurangan:**
 - Memerlukan banyak memori.
 - Sensitif terhadap nilai kkk dan outlier.
-

7. Naive Bayes

- **Cara Kerja:**
- Menggunakan Teorema Bayes:
$$P(C|X) = \frac{P(X|C)P(C)}{P(X)}$$
 - Asumsi: Semua fitur bersifat independen (*naive assumption*).
 - Probabilitas posterior dihitung untuk setiap kelas, dan kelas dengan probabilitas tertinggi dipilih.
- **Kelebihan:**
 - Cepat dan efisien untuk dataset besar.
 - Cocok untuk data teks (misalnya analisis sentimen).
- **Kekurangan:**
 - Asumsi independensi jarang berlaku di dunia nyata.

- Tidak cocok untuk data numerik dengan korelasi tinggi.

8. K-Means Clustering

• Cara Kerja:

- Tentukan jumlah cluster (k).
- Inisialisasi centroid secara acak.
- Setiap titik data diberi label berdasarkan kedekatannya dengan centroid.
- Centroid diperbarui sebagai rata-rata dari titik dalam cluster.
- Proses berlanjut hingga centroid tidak berubah (konvergensi).

• Kelebihan:

- Cepat untuk dataset besar.
- Mudah diimplementasikan.

• Kekurangan:

- Perlu menentukan k sebelumnya.
- Sensitif terhadap outlier.

9. Principal Component Analysis (PCA)

• Cara Kerja:

- PCA menemukan *principal components* (kombinasi linier fitur) yang memaksimalkan varians dalam data.
- Komponen dihitung menggunakan **eigenvalue** dan **eigenvector** dari matriks kovarians.
- Data diproyeksikan ke ruang berdimensi lebih rendah.

- **Kelebihan:**
 - Mengurangi dimensi data tanpa kehilangan informasi penting.
 - Membantu visualisasi dataset berdimensi tinggi.
 - **Kekurangan:**
 - Informasi sulit diinterpretasikan setelah reduksi.
 - Tidak cocok untuk data dengan hubungan non-linear.
-

10. Gradient Boosting (XGBoost, LightGBM)

- **Cara Kerja:**
 - Gradient Boosting membangun model secara bertahap dengan menyesuaikan kesalahan model sebelumnya.
 - Setiap model baru dilatih untuk meminimalkan **residuals** (sisa kesalahan).
 - Contoh: **XGBoost** dan **LightGBM** adalah implementasi canggih dari Gradient Boosting dengan fitur seperti penanganan data hilang dan paralelisasi.
 - **Kelebihan:**
 - Akurasi tinggi untuk kompetisi machine learning.
 - Menangani data tidak seimbang dengan baik.
 - **Kekurangan:**
 - Membutuhkan tuning parameter yang hati-hati.
 - Konsumsi memori tinggi.
-

Penutup

Setiap algoritma memiliki karakteristik unik, kekuatan, dan kelemahan. Sebagai seorang Data Scientist, memahami bagaimana algoritma bekerja dan memilih algoritma yang tepat

berdasarkan sifat data dan masalah adalah keterampilan penting. Mempraktikkan algoritma ini pada berbagai jenis data akan memperdalam pemahaman Anda dan meningkatkan keahlian Anda dalam machine learning.

Tambahan untuk Linear Regression

- **Overfitting dan Regularization:**
 - Untuk mencegah overfitting, teknik seperti Ridge Regression (penalti $L2$) dan Lasso Regression (penalti $L1$) sering digunakan.

Ridge Regression mengecilkan semua koefisien secara proporsional, sementara Lasso Regression dapat mengurangi beberapa koefisien hingga nol, menghasilkan model yang lebih sederhana.

- **Asumsi Penting:**
 - Hubungan linier antara variabel independen dan dependen.
 - Residual memiliki distribusi normal dan varian konstan (*homoscedasticity*).

Tambahan untuk Logistic Regression

- **Multiclass Logistic Regression:**
 - Untuk masalah klasifikasi dengan lebih dari dua kelas, Logistic Regression dapat diperluas menggunakan metode seperti **One-vs-Rest (OvR)** atau **Softmax Regression**.
- **Feature Scaling:**
 - Logistic Regression sensitif terhadap skala fitur, sehingga normalisasi atau standarisasi sering diperlukan.

Tambahan untuk Decision Tree

- **Overfitting:**
 - Decision Tree rentan terhadap overfitting jika dibiarkan tumbuh tanpa batas. Teknik **pruning** (pemangkasan) digunakan untuk mengurangi ukuran pohon dan meningkatkan generalisasi.
 - **Feature Importance:**
 - Decision Tree dapat digunakan untuk menentukan fitur mana yang paling penting berdasarkan informasi yang diberikan setiap fitur dalam membagi data.
-

Tambahan untuk Random Forest

- **Out-of-Bag (OOB) Error:**
 - Random Forest secara otomatis menghitung error dengan menggunakan data yang tidak termasuk dalam subset pelatihan setiap pohon, sehingga tidak memerlukan set validasi terpisah.
 - **Hyperparameter Tuning:**
 - Beberapa parameter penting yang perlu disetel, seperti jumlah pohon (*n_estimators*), kedalaman maksimum (*max_depth*), dan jumlah fitur maksimum (*max_features*).
-

Tambahan untuk SVM

- **Linear vs Non-Linear SVM:**
 - Linear SVM digunakan ketika data dapat dipisahkan dengan garis lurus, sedangkan Non-Linear SVM

menggunakan kernel trick untuk menangani data yang lebih kompleks.

- **Regularization Parameter (C):**
 - Parameter C mengontrol keseimbangan antara margin besar dan kesalahan klasifikasi. Nilai C kecil menghasilkan margin besar (lebih toleran terhadap kesalahan).
-

Tambahan untuk K-Nearest Neighbors (KNN)

- **Peningkatan KNN:**
 - **Weighted KNN:** Tetangga yang lebih dekat diberi bobot lebih tinggi.
 - **Ball Tree/KD Tree:** Struktur data untuk mempercepat pencarian tetangga terdekat dalam dataset besar.
 - **Skala Data:**
 - Karena KNN berbasis jarak, skala fitur sangat memengaruhi hasil. Normalisasi sangat penting.
-

Tambahan untuk Naive Bayes

- **Varian Naive Bayes:**
 - **Gaussian Naive Bayes:** Cocok untuk data kontinu yang diasumsikan mengikuti distribusi normal.
 - **Multinomial Naive Bayes:** Digunakan untuk data diskrit, seperti jumlah kata dalam teks.
 - **Bernoulli Naive Bayes:** Digunakan untuk fitur biner.
- **Manfaat untuk Data Teks:**

- Sangat efektif untuk klasifikasi teks karena mampu menangani data sparsity (misalnya, dokumen yang direpresentasikan sebagai bag-of-words).
-

Tambahan untuk K-Means Clustering

- **Pemilihan Jumlah Cluster (k):**
 - Metode **Elbow** digunakan untuk menentukan kkk optimal dengan melihat grafik antara jumlah cluster dan *inertia* (dalam-cluster sum of squares).
 - **K-Means++:**
 - Teknik inisialisasi centroid untuk memastikan hasil yang lebih stabil dan akurat dibanding inisialisasi acak.
 - **Kekurangan Tambahan:**
 - K-Means hanya bekerja dengan data berbasis Euclidean dan tidak efektif untuk cluster berbentuk non-sferis.
-

Tambahan untuk PCA

- **Komponen yang Dipilih:**
 - Pemilihan jumlah komponen utama didasarkan pada persentase varians total yang ingin dipertahankan (misalnya, 95% dari varians total).
 - **Kelebihan untuk Data Besar:**
 - PCA sering digunakan untuk mempercepat algoritma machine learning dengan mengurangi jumlah fitur tanpa kehilangan informasi yang signifikan.
-

Tambahan untuk Gradient Boosting

- **Varian Gradient Boosting:**

- **XGBoost:** Implementasi yang dioptimalkan dengan dukungan paralelisasi dan penanganan data yang hilang.
 - **LightGBM:** Dirancang untuk dataset besar dengan fitur seperti *leaf-wise tree growth* yang lebih efisien.
 - **CatBoost:** Mendukung fitur kategori tanpa perlu encoding manual.
 - **Tuning Hyperparameter:**
 - **Learning Rate** (η): Mengontrol ukuran langkah dalam setiap iterasi. Nilai kecil menghasilkan pelatihan yang lambat tetapi lebih akurat.
 - **Max Depth:** Kedalaman maksimum pohon menentukan kompleksitas model.
-

Poin Tambahan Umum

1. Pra-Pemrosesan Data:

- Algoritma seperti Logistic Regression, SVM, dan KNN sangat bergantung pada skala data. Fitur harus dinormalisasi atau distandardisasi sebelum digunakan.
- Missing values perlu ditangani, terutama untuk algoritma ensemble seperti Random Forest dan Gradient Boosting.

2. Evaluasi Model:

- **K-Fold Cross Validation:** Teknik untuk mengurangi overfitting dan mendapatkan evaluasi model yang lebih handal.
- **Metode Evaluasi:**

- Klasifikasi: Akurasi, Precision, Recall, F1-Score, ROC-AUC.
- Regresi: MSE, MAE, R^2 -Score.

3.

Interpretasi Model:

- Model seperti Decision Tree, Random Forest, dan Logistic Regression lebih mudah diinterpretasikan dibanding SVM atau Gradient Boosting.
- Alat seperti SHAP (SHapley Additive exPlanations) dapat digunakan untuk interpretasi model kompleks.

4. Kombinasi Algoritma:

- Teknik seperti **Stacking** atau **Blending** dapat digunakan untuk menggabungkan kelebihan dari beberapa algoritma untuk meningkatkan akurasi.

Kesimpulan Akhir

Pemahaman mendalam tentang algoritma-algoritma ini memungkinkan seorang Data Scientist untuk memilih alat yang paling sesuai dengan karakteristik data dan tujuan analisis. Dengan memahami tambahan poin seperti regularisasi, tuning hyperparameter, dan penggunaan varian yang lebih canggih, Anda dapat meningkatkan efisiensi dan efektivitas model machine learning yang Anda bangun. Mempraktikkan algoritma ini dengan dataset nyata akan semakin memperkuat pemahaman dan keahlian Anda.

1. ChatGPT 4o (2024). Kopilot Artikel ini. Tanggal akses: 25 Desember 2024. Akun penulis.
<https://chatgpt.com/c/676b7c7b-d7b8-8013-b1c8-114a0f46108f>

Buku Utama

1. **Bishop, C. M.** (2006). *Pattern Recognition and Machine Learning*. Springer.
 - Buku klasik yang mencakup penjelasan teori dan implementasi algoritma seperti Linear Regression, Logistic Regression, dan SVM.
2. **Goodfellow, I., Bengio, Y., & Courville, A.** (2016). *Deep Learning*. MIT Press.
 - Buku ini memiliki bab awal yang mencakup banyak konsep dasar machine learning, termasuk Gradient Boosting.
3. **Hastie, T., Tibshirani, R., & Friedman, J.** (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
 - Referensi komprehensif untuk algoritma seperti PCA, Random Forest, dan K-Means.
4. **Murphy, K. P.** (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
 - Buku yang menjelaskan konsep Naive Bayes, Logistic Regression, dan teori probabilitas lainnya.
5. **Friedman, J., Hastie, T., & Tibshirani, R.** (2001). *The Elements of Statistical Learning*. Springer.
 - Buku lanjutan untuk memahami teori ensemble methods seperti Gradient Boosting.

Dokumentasi dan Artikel Online

1. Scikit-learn Documentation

- Scikit-learn adalah pustaka Python populer untuk machine learning yang menyediakan implementasi semua algoritma di atas dengan dokumentasi rinci.
 - <https://scikit-learn.org/stable/>

2. XGBoost Documentation

- Panduan resmi untuk implementasi XGBoost.
 - <https://xgboost.readthedocs.io/>

3. LightGBM Documentation

- Panduan implementasi untuk LightGBM, varian Gradient Boosting yang lebih efisien.
 - <https://lightgbm.readthedocs.io/>

4. Stanford Machine Learning Course

- Kursus gratis yang diajarkan oleh Andrew Ng melalui Coursera, mencakup algoritma seperti Linear Regression, Logistic Regression, dan SVM.
 - <https://www.coursera.org/learn/machine-learning>

Jurnal Akademik

1. **Ho, T. K.** (1995). "Random Decision Forests." *Proceedings of 3rd International Conference on Document Analysis and Recognition*. IEEE.
 - Penelitian awal yang memperkenalkan Random Forest sebagai metode ensemble.
2. **Cortes, C., & Vapnik, V.** (1995). "Support-Vector Networks." *Machine Learning*, 20(3), 273–297.
 - Artikel seminal yang memperkenalkan Support Vector Machines (SVM).
3. **Breiman, L.** (2001). "Random Forests." *Machine Learning*, 45(1), 5–32.
 - Penjelasan teoretis dan aplikasi Random Forest.
4. **Pearson, K.** (1901). "On Lines and Planes of Closest Fit to Systems of Points in Space." *Philosophical Magazine*, 2(11), 559–572.

- Artikel yang memperkenalkan PCA untuk pertama kalinya.
5. **Quinlan, J. R.** (1986). "Induction of Decision Trees." *Machine Learning*, 1(1), 81–106.
- Penelitian klasik yang mendasari Decision Tree.
-

Website dan Blog Referensi

1. Machine Learning Mastery by Jason Brownlee

- Blog yang menyediakan panduan implementasi untuk hampir semua algoritma machine learning.
- <https://machinelearningmastery.com/>

2. Towards Data Science (Medium Blog)

- Blog ini berisi artikel praktis tentang penggunaan algoritma machine learning dengan contoh kode Python.
- <https://towardsdatascience.com/>

3. Analytics Vidhya

- Platform untuk belajar machine learning melalui tutorial dan studi kasus nyata.
- <https://www.analyticsvidhya.com/>

4. **ChatGPT 4o (2024)**. Kopilot Artikel ini. Tanggal akses: 25 Desember 2024. Akun penulis.
<https://chatgpt.com/c/676b7c7b-d7b8-8013-b1c8-114a0f46108f>
-

Kursus Video

1. YouTube Channels

- **StatQuest with Josh Starmer:** Penjelasan algoritma machine learning yang sederhana dan visual.
 - <https://www.youtube.com/user/joshstarmer>
- **Kaggle Learn:** Kursus singkat untuk algoritma seperti K-Means, Random Forest, dan XGBoost.
 - <https://www.kaggle.com/learn>

2. Fast.ai

- Platform pembelajaran machine learning dan deep learning.
- <https://www.fast.ai/>

Software dan Tools

1. Python Libraries:

- **Scikit-learn:** Algoritma machine learning untuk analisis data.
- **XGBoost:** Pustaka untuk Gradient Boosting.
- **LightGBM:** Untuk menangani dataset besar dengan ensemble methods.

2. R Packages:

- **caret:** Framework untuk membangun model prediktif.
- **randomForest:** Paket untuk implementasi Random Forest.

3. MATLAB:

- MATLAB memiliki toolbox lengkap untuk algoritma seperti PCA, SVM, dan Clustering.
- <https://www.mathworks.com/products/machine-learning.html>

