

Laporan Final Project Machine Learning

Prediksi Nilai Ujian Pada Sekolah Menengah ke Atas
Menggunakan Model RandomForest
Regressor, XGboost, MLP Tensorflow, dan Linear Regresi

Muhammad Risang - 50242310xx

Akhmad Rizqullah Ridlohi - 5024231037

Abraham Napitupulu - 5024231037

*Departemen Teknik Komputer
Institut Teknologi Sepuluh Nopember
Surabaya*

Penulis: Kelompok 6
20 Juni 2025

Daftar Isi

Dasar Teori

Pada bagian ini, akan dibahas landasan teoretis yang menjadi dasar pelaksanaan proyek. Penjelasan mencakup konsep fundamental machine learning, alur kerja yang diterapkan, metodologi pra-pemrosesan data, serta arsitektur dan cara kerja dari setiap model yang digunakan untuk memprediksi kinerja akademik siswa.

1 Machine Learning dan Supervised Learning

Machine learning adalah cabang dari kecerdasan buatan (AI) yang berfokus pada pengembangan algoritma yang memungkinkan komputer untuk belajar dari data dan membuat prediksi atau keputusan tanpa diprogram secara eksplisit. Dalam proyek ini, pendekatan yang digunakan adalah Supervised Learning (Pembelajaran Terarah).

Dalam Supervised Learning, model dapat belajar dari dataset yang telah memiliki label atau target yang diketahui. Model menganalisis serangkaian fitur input (variabel independen, X) dan hubungannya dengan variabel output (variabel dependen, y). Tujuan utamanya adalah untuk mempelajari fungsi pemetaan ($y=f(X)$) sehingga ketika data baru tanpa label diberikan, model dapat memprediksi outputnya dengan akurasi tinggi.

Solusi dari permasalahan yang ingin diselesaikan dalam proyek ini meliputi model regresi, yaitu salah satu bentuk Supervised Learning di mana variabel output (y) yang diprediksi bersifat kontinu atau numerik. Dalam konteks ini, model bertujuan untuk memprediksi nilai akhir siswa (G3) yang merupakan sebuah angka.

2 Alur Kerja Proyek Machine Learning

Untuk memastikan hasil yang sistematis dan dapat dipertanggungjawabkan, proyek ini mengikuti alur kerja standar dalam pengembangan model machine learning, yang terdiri dari beberapa tahapan utama:

1. **Eksplorasi Data (EDA):** Memahami karakteristik, distribusi, dan hubungan antar variabel dalam dataset melalui statistik deskriptif dan visualisasi data.
2. **Rekayasa Fitur (Feature Engineering):** Membuat fitur-fitur baru yang lebih informatif dari fitur yang sudah ada (misalnya, $grade_{avg,rev}$) untuk meningkatkan performa model.
3. **Pemodelan (Modeling):** Memilih, melatih, dan mengevaluasi beberapa algoritma yang berbeda untuk menemukan model dengan performa terbaik.
4. **Evaluasi Model:** Mengukur kinerja model terbaik menggunakan metrik kuantitatif seperti R-squared dan RMSE untuk menilai akurasi dan keandalannya.
5. **Prediksi (Inference):** Menggunakan model yang telah dilatih untuk membuat prediksi pada data baru yang belum pernah dilihat sebelumnya, seperti pada simulasi uji coba real-time.

3 Pra-pemrosesan Data

Data mentah jarang sekali bisa langsung digunakan oleh model. Oleh karena itu, diperlukan beberapa teknik pra-pemrosesan untuk membersihkan dan menstrukturkan data.

3.1 One-Hot Encoding

Sebagian besar model machine learning hanya dapat bekerja dengan data numerik. Fitur-fitur kategorikal dalam dataset ini (seperti sex, address, higher) yang berbentuk teks perlu diubah menjadi representasi numerik. One-Hot Encoding adalah teknik yang digunakan untuk tujuan ini, di mana setiap kategori unik dalam sebuah fitur diubah menjadi kolom biner baru (bernilai 0 atau 1).

3.2 Penskalaan Fitur (Feature Scaling)

Fitur-fitur numerik seringkali memiliki rentang nilai yang sangat berbeda (misalnya, age (15-22) vs absences (0-93)). Jika tidak diskalakan, fitur dengan rentang nilai yang lebih besar dapat mendominasi proses pembelajaran model secara tidak adil. StandardScaler digunakan untuk mentransformasi setiap fitur sehingga memiliki rata-rata 0 dan standar deviasi 1. Langkah ini sangat krusial untuk model yang sensitif terhadap skala seperti Regresi Linear dan MLP.

4 Model yang Digunakan

Pemilihan model dilakukan secara strategis untuk membandingkan pendekatan dari berbagai paradigma dengan tingkat kompleksitas yang berbeda.

4.1 Regresi Linear

Model ini dipilih sebagai baseline atau titik awal perbandingan. Regresi Linear bekerja dengan mengasumsikan adanya hubungan linear antara fitur-fitur input dan variabel target. Ia mencoba menemukan satu formula matematis terbaik untuk memetakan input ke output. Meskipun sederhana, model ini sangat berguna untuk mengukur seberapa baik masalah dapat diselesaikan dengan pendekatan linear dan untuk menilai relevansi fitur yang dipilih.

4.2 Random Forest

Random Forest adalah perwakilan dari model ensemble berbasis pohon keputusan. Cara kerjanya adalah dengan membangun ratusan pohon keputusan secara independen, di mana setiap pohon dilatih pada sampel data yang sedikit berbeda. Untuk membuat prediksi, model ini mengambil rata-rata dari semua prediksi pohon individual. Pendekatan "wisdom of the crowd" ini membuatnya sangat kuat dalam menangkap interaksi fitur yang kompleks, robust terhadap noise, dan cenderung tidak overfitting, sehingga sangat cocok untuk data tabular seperti pada proyek ini.

4.3 XGBoost

XGBoost (eXtreme Gradient Boosting) dipilih sebagai perwakilan dari metode gradient boosting. Berbeda dari Random Forest, XGBoost membangun pohon keputusan secara sekuensial. Setiap pohon baru yang dibuat secara spesifik bertugas untuk memperbaiki kesalahan prediksi dari pohon-pohon sebelumnya. Proses belajar iteratif yang fokus pada kesalahan ini membuat XGBoost menjadi salah satu model dengan performa tertinggi di industri untuk masalah regresi dan klasifikasi pada data tabular.

4.4 Multi-Layer Perceptron (MLP)

MLP dipilih untuk menjelajahi pendekatan deep learning. Model ini bekerja seperti jaringan otak buatan, terdiri dari lapisan-lapisan "neuron" yang saling terhubung. MLP mampu belajar dan mengenali pola-pola yang sangat

kompleks dan abstrak dari kombinasi semua fitur secara bersamaan. Model ini dipilih untuk menguji apakah arsitektur non-linear yang sangat fleksibel ini dapat menemukan pola yang terlewatkan oleh model berbasis pohon.

5 Metrik Evaluasi Model

Untuk mengukur dan membandingkan kinerja model secara objektif, dua metrik utama digunakan:

- **R-squared (R^2):** Metrik ini mengukur seberapa besar persentase variasi (perbedaan) dalam variabel target (G3) yang dapat dijelaskan oleh model. Nilai R^2 berkisar dari 0 hingga 1, di mana nilai yang lebih tinggi menunjukkan model yang lebih baik.
- **RMSE (Root Mean Squared Error):** Metrik ini mengukur rata-rata besaran kesalahan prediksi model dalam satuan yang sama dengan variabel target. Dalam proyek ini, RMSE menunjukkan rata-rata berapa poin nilai prediksi meleset dari nilai sebenarnya. Nilai RMSE yang lebih rendah menunjukkan model yang lebih akurat.

Hasil Praktikum

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Hasil

Kami melakukan uji coba model dengan skenario dataset yang memiliki total 33 fitur, dan fitur yang dipakai untuk memprediksi kinerja akademik siswa. Fitur-fitur tersebut meliputi nilai ujian sebelumnya (G1, G2, dan rata-ratanya), data demografi dan keluarga (usia, pendidikan orang tua, jenis kelamin, alamat), serta kebiasaan belajar dan sekolah (waktu belajar, kegagalan sebelumnya, absensi, akses internet, partisipasi kegiatan). Kami juga menyertakan aspek gaya hidup dan sosial siswa, seperti frekuensi bersosialisasi, status kesehatan, status hubungan romantis, dan konsumsi alkohol, dengan jumlah siswa total 395. Setelah melakukan uji coba didapatkan hasil sebagai berikut

6 Uji Coba dengan Data Siswa

```
=== LANGKAH 7: UJI COBA REAL-TIME ===

--- Data Siswa Baru ---
G1          15
G2          15
failures    0
higher      yes
age         16
Medu        4
Fedu        3
studytime   3
absences    2
goout       2
health      5
sex         M
address     U
internet     yes
romantic    no
freetime    4
Dalc        1
Walc        2
activities  yes
dtype: object
...
>>> Prediksi Nilai Akhir (G3) oleh MLP: 15.13
(Catatan: RMSE model ini ~2.66)
```

Gambar 1: Data Siswa Pertama


Di sini pada nilai aktualnya nilai G3 siswa adalah 15, dan saat dilakukan uji coba prediksi didapatkan setiap model seperti yang ada pada gambar berikut.

```
--- Prediksi Menggunakan Linear Regression (Simple) ---
>>> Prediksi Nilai Akhir (G3) oleh LinReg (Simple): 10.91
(Catatan: RMSE model ini ~4.67)

--- Prediksi Menggunakan Linear Regression (Lengkap) ---
>>> Prediksi Nilai Akhir (G3) oleh LinReg (Lengkap): 18.06
(Catatan: RMSE model ini ~2.24)

--- Prediksi Menggunakan RandomForest Regressor ---
>>> Prediksi Nilai Akhir (G3) oleh RandomForest: 15.23
(Catatan: RMSE model ini ~1.75)

--- Prediksi Menggunakan XGBoost Regressor ---
>>> Prediksi Nilai Akhir (G3) oleh XGBoost: 14.82
(Catatan: RMSE model ini ~1.82)

--- Prediksi Menggunakan MLP TensorFlow/Keras ---
1/1  0s 158ms/step
>>> Prediksi Nilai Akhir (G3) oleh MLP: 15.13
(Catatan: RMSE model ini ~2.66)
```

Gambar 2: Hasil Uji Coba Model Pertama

Dilanjutkan uji coba dengan data siswa sebagai berikut:

Dengan data siswa tersebut, kami mengasumsikan bahwa siswa tersebut memiliki nilai G3 sebesar 14,5 yang didapat dengan perkiraan secara intuitif menggunakan heatmap yang ada dan didapatkan hasil seperti pada gambar berikut.

7 Evaluasi Model

Berikut adalah visualisasi untuk RMSE (Root Mean Square Error) dan R-Squared, di mana nilai RMSE yang baik adalah mendekati 0 dan R-Squared mendekati 1.

Dari data tersebut RandomForest Regressor mendapatkan nilai RMSE yang lebih rendah daripada semua model yang ada, diikuti oleh XGBoost, Linear Regression dengan fitur yang lengkap, lalu MLP TensorFlow. Nilai RMSE menunjukkan berapa rentang error yang dihasilkan oleh data. Untuk R-Squared, RandomForest Regressor juga merupakan yang paling baik

```

...
=== LANGKAH 7: UJI COBA REAL-TIME ===

--- Data Siswa Baru ---
G1      15
G2      16
failures 0
higher  yes
age     17
Medu    4
Fedu    4
studytime 3
absences 2
goout   2
health  5
sex     M
address U
internet yes
romantic no
freetime 3
Dalc    1
Walc    1
activities yes
dtype: object

--- Prediksi Menggunakan Linear Regression (Simple) ---
>>> Prediksi Nilai Akhir (G3) oleh LinReg (Simple): 10.91
(Catatan: RMSE model ini ~4.67)

```

Gambar 3: Data Siswa Kedua

```

--- Prediksi Menggunakan Linear Regression (Simple) ---
>>> Prediksi Nilai Akhir (G3) oleh LinReg (Simple): 10.91
(Catatan: RMSE model ini ~4.67)

--- Prediksi Menggunakan Linear Regression (Lengkap) ---
>>> Prediksi Nilai Akhir (G3) oleh LinReg (Lengkap): 18.76
(Catatan: RMSE model ini ~2.24)

--- Prediksi Menggunakan RandomForest Regressor ---
>>> Prediksi Nilai Akhir (G3) oleh RandomForest: 15.67
(Catatan: RMSE model ini ~1.75)

--- Prediksi Menggunakan XGBoost Regressor ---
>>> Prediksi Nilai Akhir (G3) oleh XGBoost: 15.01
(Catatan: RMSE model ini ~1.82)

--- Prediksi Menggunakan MLP TensorFlow/Keras ---
1/1 [56ms/step]
>>> Prediksi Nilai Akhir (G3) oleh MLP: 16.23
(Catatan: RMSE model ini ~2.68)

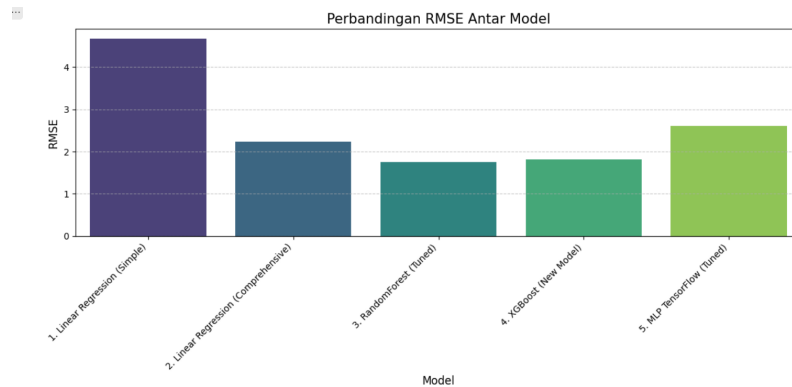
```

Gambar 4: Hasil Uji Coba Model Kedua

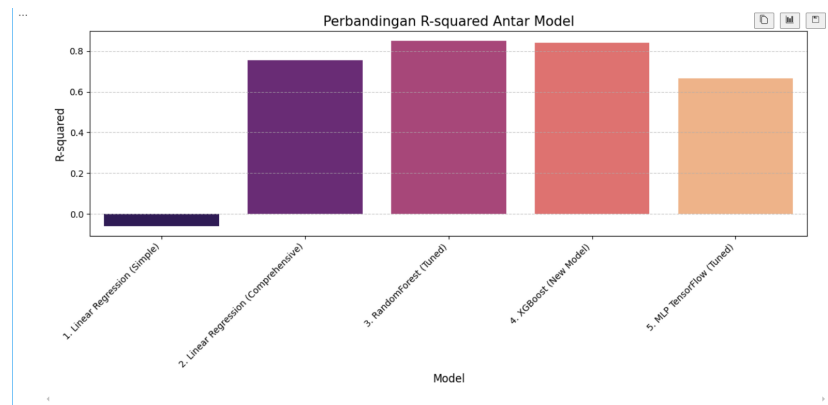
di antara model yang ada. R-Squared menunjukkan informasi mengenai bagaimana model tersebut menjelaskan variabilitas data terhadap model yang digunakan. R-Squared menggambarkan seberapa baik model sesuai dengan data yang digunakan.

Secara teknis, RandomForest Regressor mungkin lebih cocok untuk dataset yang relatif lebih kecil dibandingkan MLP TensorFlow dan XGBoost. RandomForest umumnya dianggap lebih tidak rumit dalam interpretasi dan implementasi dibandingkan dua model lainnya.

Dapat dilihat bahwa performa model RandomForest Regressor adalah model dengan prediksi yang paling akurat dibandingkan dengan model-model lain. Hal ini dapat terjadi dikarenakan RandomForest Regressor merupakan



Gambar 5: Visualisasi RMSE

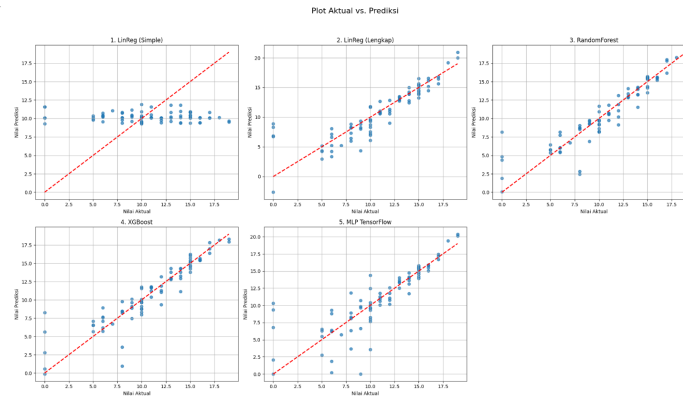


Gambar 6: Visualisasi R-Squared

model yang memang cocok untuk dataset yang lebih kecil daripada MLP TensorFlow dan XGBoost. Jika dibandingkan secara teknis, RandomForest merupakan model yang paling tidak rumit dibanding dua model tersebut. RandomForest Regressor bekerja dengan membuat decision tree secara paralel dengan memilih data secara acak, dan dari decision tree yang sudah dibuat secara paralel nilai prediksi dari setiap decision tree tersebut akan dirata-rata untuk mendapatkan hasilnya. XGBoost juga masih menggunakan decision tree, namun dengan cara memperbaiki kesalahan berulang kali dari decision tree sebelumnya. Untuk penjelasan simpelnya, MLP TensorFlow sendiri merupakan salah satu model jaringan saraf tiruan yang memiliki tiga komponen, yaitu *Input Layer*, *Hidden Layer*, dan *Output Layer* yang memiliki cara pemodelan yang lebih kompleks dibandingkan RandomForest Regressor dan XGBoost.

8 Visualisasi Scatter Plot

Kami juga melakukan visualisasi menggunakan scatter plot.



Gambar 7: Visualisasi Scatter Plot

Dari visualisasi scatter plot yang membandingkan nilai aktual dengan nilai prediksi untuk setiap model, beberapa wawasan penting dapat ditarik mengenai performa dan karakteristik setiap algoritma. Idealnya, titik-titik pada scatter plot harus terkonsentrasi rapat di sepanjang garis diagonal merah ($y=x$), menunjukkan bahwa nilai prediksi sangat mendekati nilai aktual.

Model RandomForest Regressor dan XGBoost menunjukkan performa yang paling baik dalam visualisasi ini. Titik-titik data untuk kedua model ini tampak paling padat dan terpusat di sekitar garis diagonal. Ini mengindikasikan bahwa baik RandomForest maupun XGBoost mampu menangkap hubungan kompleks dalam data dan menghasilkan prediksi yang konsisten dan akurat di berbagai rentang nilai G3. Kerapatan titik di sekitar garis menunjukkan bahwa error prediksi (perbedaan antara aktual dan prediksi) cenderung kecil dan terdistribusi merata.

Sebaliknya, model Linear Regression yang menggunakan fitur baseline menunjukkan sebaran titik yang lebih luas dan menyebar jauh dari garis diagonal. Ini terutama terlihat pada nilai-nilai G3 yang lebih ekstrem (sangat rendah atau sangat tinggi), di mana prediksi cenderung menyimpang lebih jauh dari nilai aktual. Sebaran yang lebih tersebar tidak merata mencerminkan keterbatasan model linear dalam menangkap hubungan non-linear atau interaksi kompleks antar fitur, yang mengakibatkan nilai RMSE yang lebih tinggi dibandingkan model berbasis pohon atau jaringan saraf.

Model MLP TensorFlow menunjukkan performa yang kuat, berada di antara Linear Regression dan RandomForest/XGBoost. Meskipun titik-titiknya tidak sepadat RandomForest atau XGBoost, mereka tetap menun-

jukkan konsentrasi yang baik di sekitar garis diagonal. Ini menegaskan bahwa MLP adalah model yang efektif dan mampu memberikan prediksi yang cukup akurat, meskipun dalam kasus ini RandomForest sedikit unggul.

Secara keseluruhan, scatter plot ini secara visual mengkonfirmasi metrik evaluasi yang disajikan sebelumnya. Model-model yang memiliki nilai R-Squared tinggi dan RMSE rendah (seperti RandomForest dan XGBoost) ditunjukkan dengan titik-titik yang sangat dekat dengan garis ideal, sementara model dengan metrik yang kurang baik menunjukkan sebaran yang lebih lebar. Visualisasi ini juga membantu mengidentifikasi rentang nilai di mana suatu model mungkin berkinerja kurang baik atau di mana terdapat outlier prediksi.

Diskusi

Hasil yang diperoleh dari uji model menggunakan dataset ini adalah model RandomForest Regressor yang terbaik pada skenario ini yaitu dataset yang memiliki banyak fitur namun jumlah data yang ada pada dataset hanya 395 dimana data ini relatif cukup kecil jika dibandingkan dengan lainnya. Pada dataset kami, kami memiliki keterbatasan untuk menguji model ini pada dataset yang lebih banyak sehingga dapat melihat potensial penuh dari XGboost,dan MLP TensorFlow. Rekomendasi dari kami penulis adalah untuk mencoba model ini pada dataset yang mempunyai data yang lebih banyak daripada dataset ini untuk mendapatkan potensi penuh dari XGboost,dan MLP TensorFlow.

Kesimpulan

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.