

Laporan Kerja Individu

Final Project - Data Analytics

Program Zenius Studi Independen Bersertifikat - Angkatan 4

Nama Lengkap: Zulfa Nabilah Nurvitasari

Nomor ID live class: 099

Nomor Kelompok: 30

Mentor: Erwin Fernanda

Deskripsi peran

- Melakukan *cleansing* dan *preparation* pada dataset *credit card balance* dan *previous application*
- Melakukan *cleansing* dan *preparation* pada dataset *application train* (Semua anggota melakukan *cleansing* dan *preparation*. Oleh karena itu, dari semua bagian *cleansing* dan *preparation* data yang telah saya lakukan, yang digunakan ialah bagian *data encoding*).
- Melakukan pemodelan dengan menggunakan *Supervised Learning* menggunakan KNN, *logistic regression*, dan *random forest*
- Membuat dashboard dengan Looker

Lampiran Hasil Kerja

- *Cleansing* dan *Preparation* pada dataset *credit card balance* dan *previous application*:
<https://colab.research.google.com/drive/1wyx2ypaaOnYMmLx5s3dTkE0JF76isL6y?usp=sharing>
- *Cleansing*, *Preparation*, dan *modeling* pada dataset *application train*:
https://colab.research.google.com/drive/1ki0chl6qm01JWFk0_sUftZ_e372jSsmc?usp=sharing
- Dashboard: <https://lookerstudio.google.com/reporting/b9b058f8-df07-4f8f-ac71-63d985fb2455>

Proses Kerja

No	Hari, Tanggal	Pekerjaan	Keterangan
1	Rabu, 31 Mei	<i>Cleansing dan preparation pada dataset credit card balance.</i>	Saya melakukan data understanding seperti melihat jumlah kolom dan baris pada data, melihat kolom pada dataset, melihat tipe data, melihat jumlah missing value di setiap kolom, dan melihat frekuensi setiap <i>unique values</i> pada data.
2	Kamis, 1 Juni	<i>Cleansing dan preparation pada dataset credit card balance.</i>	Saya melakukan <i>missing value handling</i> pada data bertipe numerik dengan mengisi <i>missing value</i> menggunakan nilai median karena tidak berdistribusi normal. Selain itu, saya juga melakukan <i>data encoding</i> dengan menggunakan <i>frequency encoding</i> .
3	Jumat, 2 Juni	<i>Cleansing dan preparation pada dataset previous application.</i>	Saya melakukan data understanding seperti melihat jumlah kolom dan baris pada data, melihat kolom pada dataset, melihat tipe data, melihat jumlah missing value di setiap kolom, dan melihat frekuensi setiap <i>unique values</i> pada data.
4	Sabtu, 3 Juni	<i>Cleansing dan preparation pada dataset previous application.</i>	Saya melakukan <i>missing value handling</i> pada data bertipe numerik dan kategorik. Saya menghapus baris dengan kolom yang memiliki missing value kurang dari 10 karena jumlah tidak terlalu besar. Selain itu, saya juga menghapus kolom dengan missing value lebih dari 1000000. Untuk data bertipe numerik yang memiliki <i>missing value</i> di bawah 1000000 saya isi dengan mean untuk kolom yang memiliki distribusi normal. Sedangkan, <i>missing value</i> pada kolom-kolom yang tidak memiliki distribusi normal, saya isi dengan median. Untuk data bertipe

			kategorik, <i>missing value</i> saya isi dengan modus.
5	Minggu, 4 Juni	<i>Cleansing dan preparation pada dataset previous application.</i>	Saya melakukan data encoding pada kolom-kolom bertipe data kategori. Untuk kolom-kolom yang memiliki <i>unique value</i> kurang dari 5, saya melakukan <i>label encoding</i> , sedangkan untuk kolom-kolom yang memiliki <i>unique value</i> lebih dari 5 dilakukan <i>frequency encoding</i> .
6	Senin, 5 Juni	<i>Cleansing dan preparation pada dataset application train</i>	Saya melakukan data understanding seperti melihat jumlah kolom dan baris pada data, melihat tipe data, melihat statistik deskriptif pada data, melihat jumlah missing value di setiap kolom, melihat jumlah duplikat, melihat frekuensi setiap <i>unique values</i> pada data, dan juga melakukan pengecekan keseimbangan data. Selain itu, saya juga melakukan unusual value handling dengan mengganti salah satu unique value pada kolom bertipe data kategori dengan modus.
7	Selasa, 6 Juni	<i>Cleansing dan preparation pada dataset application train</i>	Saya melakukan <i>missing value handling</i> pada data bertipe numerik dan kategorik. Saya menghapus baris dengan kolom yang memiliki missing value kurang dari 1000 karena jumlah tidak terlalu besar. Selain itu, saya juga menghapus kolom dengan missing value lebih dari 70%. Untuk data bertipe numerik yang memiliki <i>missing value</i> di bawah 70% saya isi dengan mean untuk kolom-kolom yang memiliki distribusi normal. Sedangkan, <i>missing value</i> pada kolom-kolom yang tidak memiliki distribusi normal, saya isi dengan median.

			Untuk data bertipe kategorik yang memiliki <i>missing value</i> di atas 50% saya isi dengan 'unk', sedangkan untuk kolom yang memiliki <i>missing value</i> di bawah 50% saya isi dengan modus.
8	Rabu, 7 Juni	<i>Cleansing dan preparation pada dataset application train</i>	Saya melakukan data encoding pada kolom-kolom bertipe data kategori. Untuk kolom-kolom yang memiliki <i>unique value</i> kurang dari 5, saya melakukan <i>label encoding</i> , sedangkan untuk kolom-kolom yang memiliki <i>unique value</i> lebih dari 5 dilakukan <i>frequency encoding</i> .
9	Jumat, 9 Juni	<i>Modeling</i>	Sebelum melakukan modeling, saya mencari nilai korelasi antar semua kolom dengan kolom TARGET terlebih dahulu. Setelah itu, saya mengambil kolom-kolom yang memiliki nilai korelasi dengan TARGET sebesar ≥ 0.01 dan sebesar ≤ -0.01 . Selanjutnya, saya membagi data menjadi X (Independen) dan Y (Dependen) dan melakukan standarisasi pada data X karena setiap kolomnya memiliki satuan yang berbeda. Setelah itu, saya melakukan <i>split train test pada data</i> . Tahap selanjutnya, yaitu melakukan <i>undersampling</i> dan <i>oversampling</i> pada data X_train dan Y_train karena ada ketidakseimbangan antar data tersebut. Setelah itu, saya melakukan modeling dengan menggunakan <i>logistic regression</i> . Setelah melakukan modeling, saya menghitung metric evaluasi untuk mengetahui performa pada model.
10	Sabtu, 10 Juni	<i>Modeling</i>	Saya melakukan modeling dengan menggunakan Random Forest dan KNN pada data yang telah di sampling dengan menggunakan <i>undersampling</i> . Setelah melakukan modeling, saya menghitung metric

			evaluasi untuk mengetahui performa pada model.
11	Minggu, 11 Juni	<i>Modeling</i>	Saya melakukan modeling dengan menggunakan Random Forest dan KNN pada data yang telah di sampling dengan menggunakan <i>oversampling</i> . Setelah melakukan modeling, saya menghitung metric evaluasi untuk mengetahui performa pada model.
12	Rabu, 14 Juni	Membuat Dashboard	Membuat pie chart client by target, score card jumlah client dengan target 1 dan 0, membuat dua tabel jumlah client dengan target 1 dan 0 berdasarkan <i>education type</i> dan <i>organization type</i> , membuat dua barchart jumlah client dengan target 1 dan 0 berdasarkan <i>contract type</i> dan gender.
13	Kamis, 15 Juni	Membuat Dashboard	Menampilkan tabel metric evaluation yang terdiri dari precision, recall, f1-score, dan support. Selain itu, saya juga menampilkan kurva ROC dan nilai ROC AUC pada model <i>logistic regression</i> tanpa <i>sampling</i> , dengan SMOTETomek, dan <i>over sampling</i> serta pada model random forest tanpa <i>sampling</i> , dengan SMOTETomek, dan <i>over sampling</i> .

Persetujuan Anggota

Laporan ini telah disetujui oleh:

- ☒ Christ Jordan Baeha
- ☒ Zulfa Nabilah Nurvitasari
- ☒ Dinny Meilinda Sari
- ☐ Sania Salsabila Agustin
- ☐ Fiqih Imanul Haq