

Banking Marketing Targets

Dokumen
Laporan Final
Project (Stage 2)

By: Group 3 DS Batch 21 aka **Jump-start**



Descriptive Statistic

- A. Apakah ada kolom dengan tipe data kurang sesuai, atau nama kolom dan isinya kurang sesuai? Pada kolom `job` ada value 'admin.' yang mana seharusnya penulisan tidak perlu menggunakan tanda titik (.) seperti pekerjaan lainnya. Selebihnya, semua tipe data sudah sesuai.**

```
data = pd.read_csv('train.csv', sep = ';')
df = data.rename(columns={'y': 'subscribed'})
df.sample(5)
```

	age	job	marital	education	default	balance	housing	loan	contact	day	month
33026	26	admin.	single	secondary	no	1375	yes	no	cellular	17	aug
19393	49	blue-collar	married	secondary	no	141	no	no	cellular	6	aug
16347	36	management	divorced	tertiary	no	14930	no	no	cellular	23	jan
2820	42	admin.	married	secondary	no	2656	yes	no	unknown	14	mar
41162	89	retired	married	tertiary	no	553	no	no	telephone	19	aug

- B. Apakah ada kolom yang memiliki nilai kosong? Jika ada, apa saja? Tidak ada kolom dengan nilai kosong.**

```
[6] df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45211 entries, 0 to 45210
Data columns (total 18 columns):
#   Column      Non-Null Count  Dtype
---  -
0    age         45211 non-null  int64
1    job         45211 non-null  object
2    marital     45211 non-null  object
3    education   45211 non-null  object
4    default     45211 non-null  object
5    balance     45211 non-null  int64
6    housing     45211 non-null  object
7    loan        45211 non-null  object
8    contact     45211 non-null  object
9    day         45211 non-null  int64
10   month       45211 non-null  object
11   duration    45211 non-null  int64
12   campaign    45211 non-null  int64
13   pdays       45211 non-null  int64
14   previous    45211 non-null  int64
15   poutcome    45211 non-null  object
16   y           45211 non-null  object
17   y2          45211 non-null  int64
dtypes: int64(8), object(10)
memory usage: 6.2+ MB
```


C. Apakah ada kolom yang memiliki nilai summary agak aneh? (min/mean/median/max/unique/top/freq)

- Pada kolom **previous** nampak issue pada nilai maksimalnya, dimana salah satu customer dihubungi pada campaign sebelumnya sebanyak 275 kali. Kemungkinan akan di drop pada saat pre-processing.
- kolom **balance**, **duration**, dan **pdays** tampak right-skewed (median < mean).

✓ 0s df.describe()

	age	balance	day	duration	campaign	pdays	previous
count	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000
mean	40.936210	1362.272058	15.806419	258.163080	2.763841	40.197828	0.580323
std	10.618762	3044.765829	8.322476	257.527812	3.098021	100.128746	2.303441
min	18.000000	-8019.000000	1.000000	0.000000	1.000000	-1.000000	0.000000
25%	33.000000	72.000000	8.000000	103.000000	1.000000	-1.000000	0.000000
50%	39.000000	448.000000	16.000000	180.000000	2.000000	-1.000000	0.000000
75%	48.000000	1428.000000	21.000000	319.000000	3.000000	-1.000000	0.000000
max	95.000000	102127.000000	31.000000	4918.000000	63.000000	871.000000	275.000000

Keanehan ditemukan
pada summary data
numerical

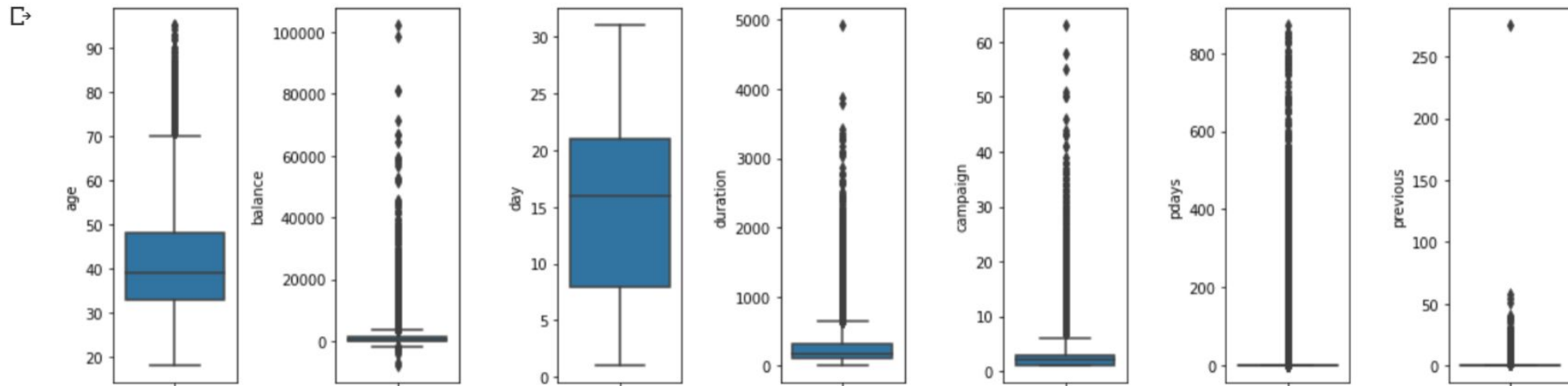
✓ 0s df[cats].describe()

	marital	education	default	housing	loan	contact	month	poutcome	y
count	45211	45211	45211	45211	45211	45211	45211	45211	45211
unique	3	4	2	2	2	3	12	4	2
top	married	secondary	no	yes	no	cellular	may	unknown	no
freq	27214	23202	44396	25130	37967	29285	13766	36959	39922

- Data didominasi oleh customer yang sudah menikah (**marital**) dan/atau tidak memiliki tunggakan (**default**) ataupun pinjaman (**loan**).
- lebih dari 75% data customer tidak diketahui hasil/output dari campaign sebelumnya (**poutcome**).

Univariate Analysis

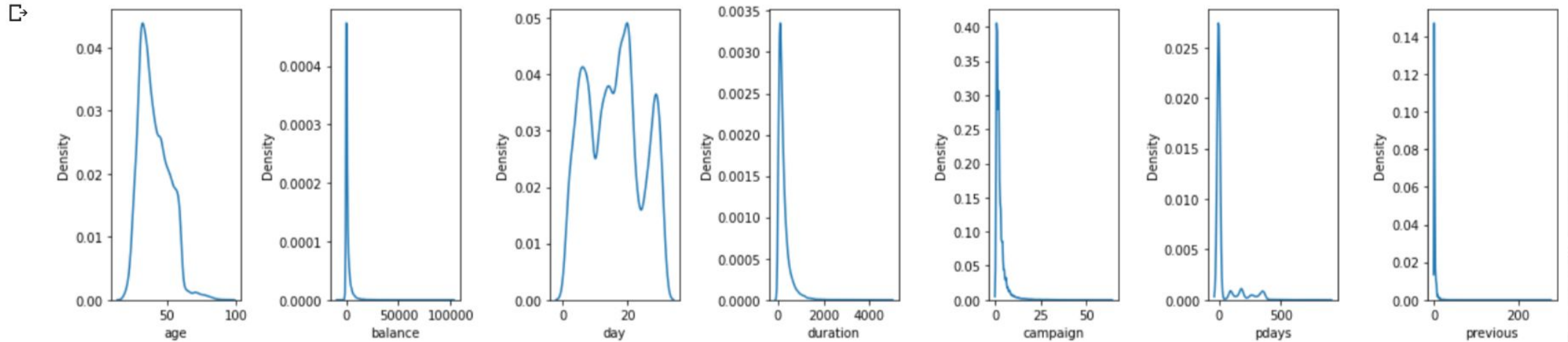
```
✓ 1s ▶ plt.figure(figsize=(15, 4))
for i in range(0, len(nums)):
    plt.subplot(1, len(nums), i+1)
    sns.boxplot(y=df[nums[i]], orient='v')
plt.tight_layout()
```



- Pada **previous** dan **duration** terdapat outlier yang berbeda sangat signifikan.

✓
4s

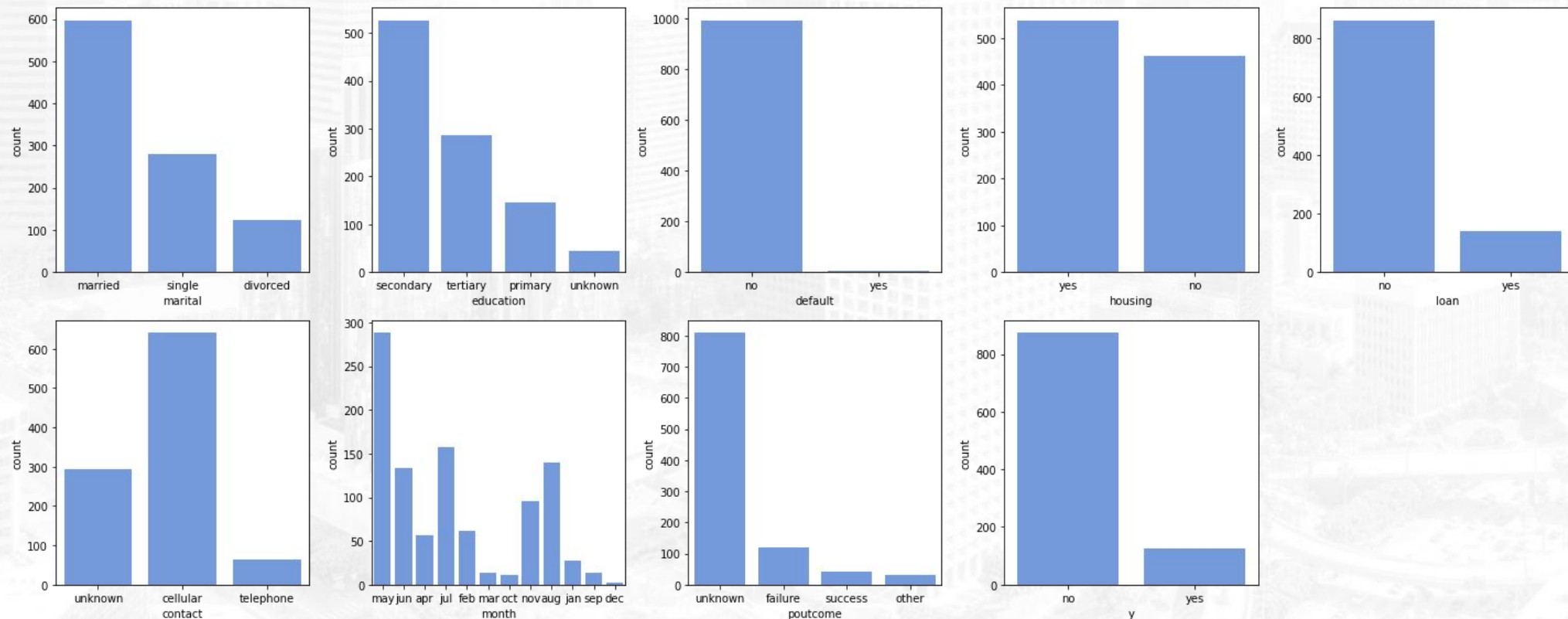
```
plt.figure(figsize=(17, 4))
for i in range(0, len(nums)):
    plt.subplot(1, len(nums), i+1)
    sns.kdeplot(x=df[nums[i]])
    plt.tight_layout()
```



- hampir semua fitur numerical, distribusinya (sangat) right-skewed.

Sedangkan untuk data kategorikal, dapat terlihat hampir seluruh fitur data memiliki ketimpangan, kecuali **housing**.

```
✓ 3s
cat_sample = df.sample(1000, random_state=42)
plt.figure(figsize=(20, 8))
for i in range(0, len(cats)):
    plt.subplot(2, 5, i+1)
    sns.countplot(x=cat_sample[cats[i]], color = 'cornflowerblue')
plt.tight_layout()
```



Multivariate Analysis

Encoding **y** untuk untuk kebutuhan cek korelasi fitur dengan target output.

```
✓ [9] def segment(x):
    0s     if x['y'] == 'yes':
        segment = 1
    else:
        segment = 0
    return segment
```

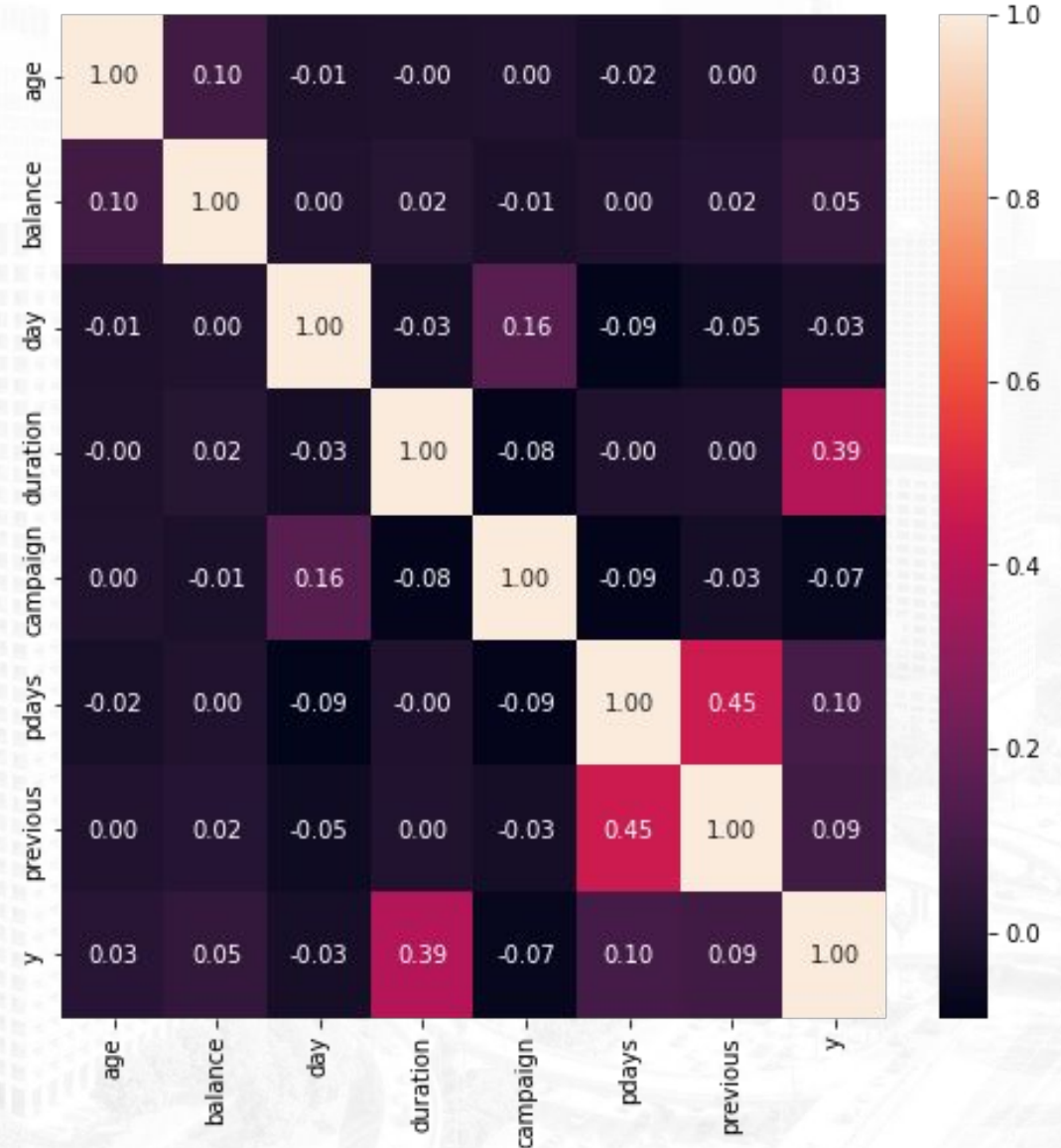
```
✓ [10] df['y'] = df.apply(lambda x: segment(x), axis= 1)
    1s df.sample(10)
```

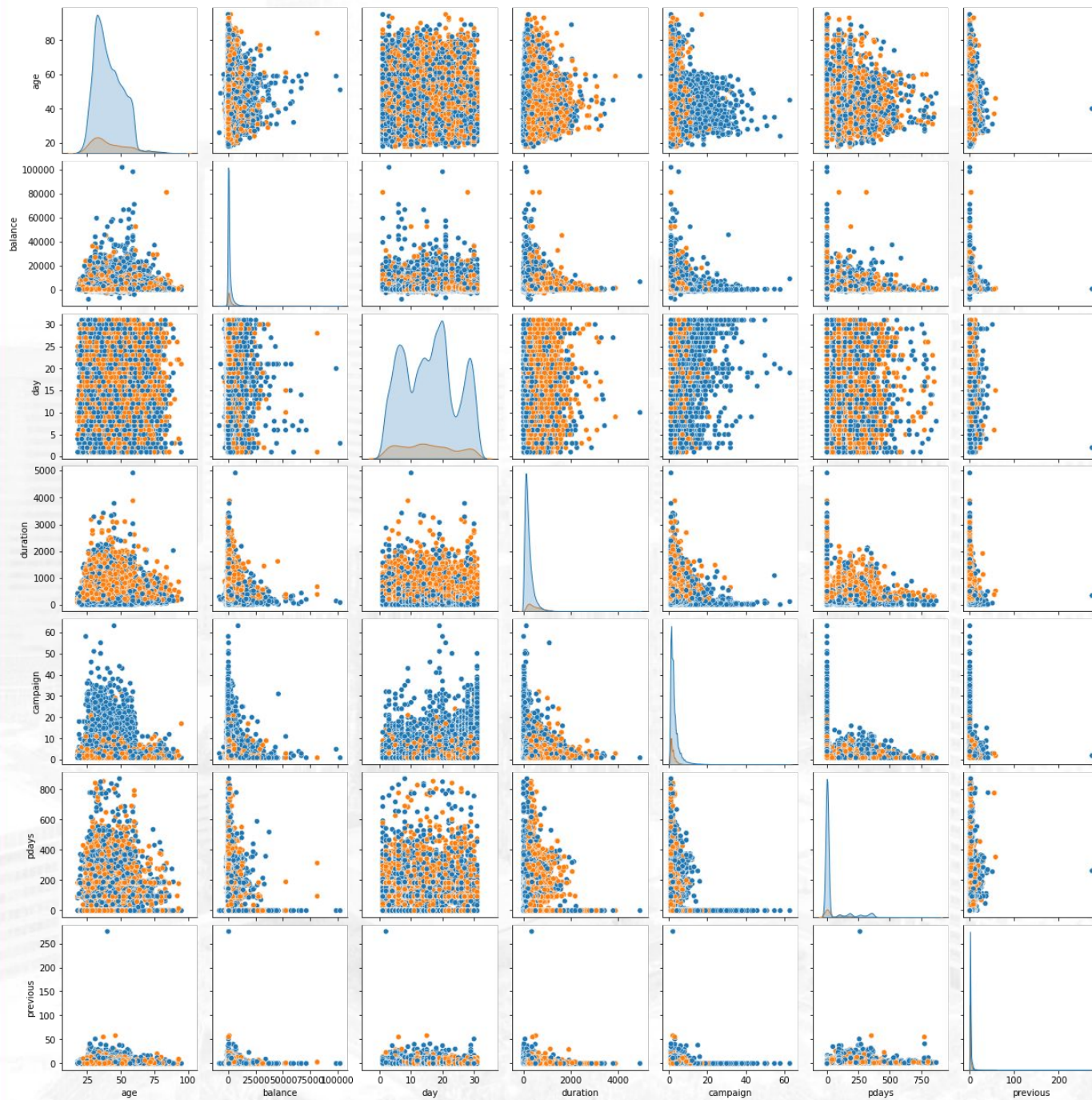
	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
13195	57	technician	married	primary	no	4442	no	no	cellular	8	jul	97	6	-1	0	unknown	0
20666	37	technician	single	secondary	no	3665	no	no	cellular	12	aug	664	3	-1	0	unknown	1
7142	38	retired	single	secondary	no	44	no	no	unknown	29	may	148	1	-1	0	unknown	0
592	41	admin.	divorced	primary	no	4070	yes	no	unknown	6	may	140	2	-1	0	unknown	0
15482	31	entrepreneur	single	secondary	no	379	yes	no	cellular	18	jul	570	2	-1	0	unknown	0
8457	38	services	married	secondary	no	823	yes	no	unknown	3	jun	132	5	-1	0	unknown	0
37178	39	management	married	tertiary	no	141	yes	no	cellular	13	may	788	2	331	6	other	0
40357	59	self-employed	married	tertiary	no	185	no	no	cellular	22	jun	177	5	138	1	failure	0
36342	46	blue-collar	married	secondary	no	-27	yes	no	cellular	11	may	254	1	-1	0	unknown	0

Korelasi Data Numerik

Dari correlation heatmap di atas dapat disimpulkan bahwa:

- **y** memiliki potensi korelasi yang tinggi dengan *duration* (strong potential correlation)
- **y** juga memiliki korelasi yang lemah dengan **previous** dan **pdays** (decent potential feature)
- **campaign** memiliki korelasi positif dengan **day**
- Sedangkan korelasi **campaign** dengan **age** sangat lemah, kemungkinan bukan fitur yang potensial (decent potential feature)
- ada korelasi antar **previous** dengan **pdays**, namun tidak cukup kuat untuk dikatakan redundant

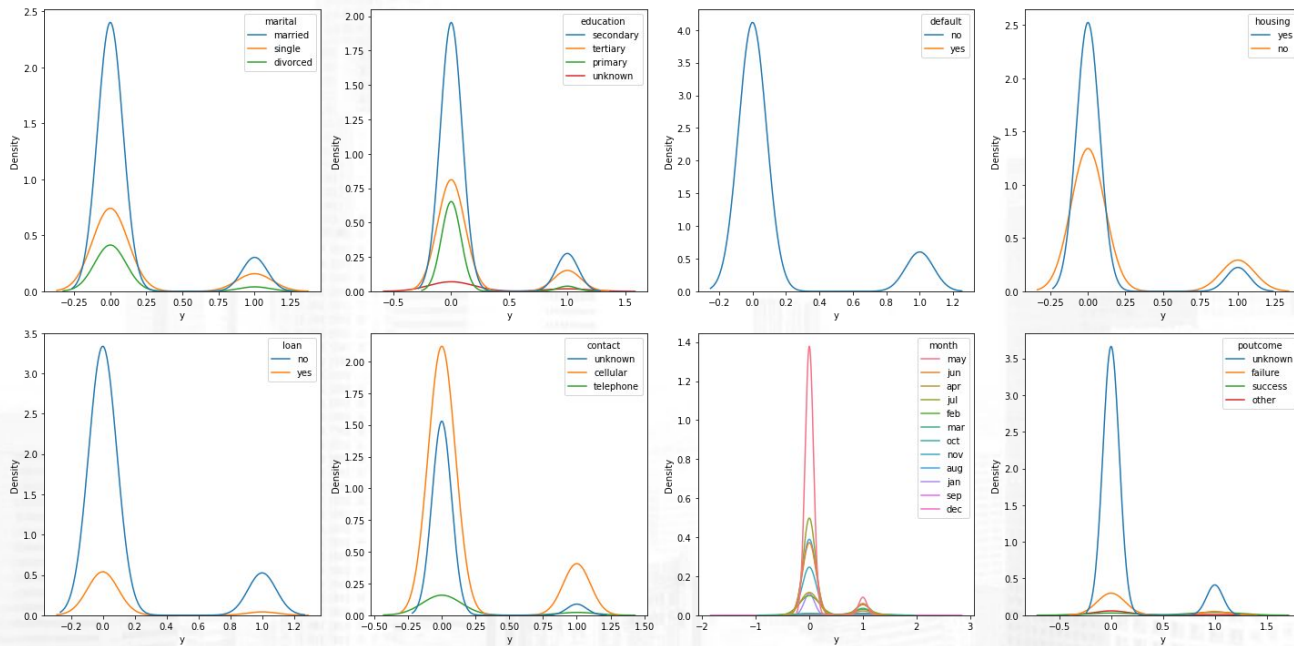




belum ada fitur yg memiliki korelasi linear yg cukup kuat.

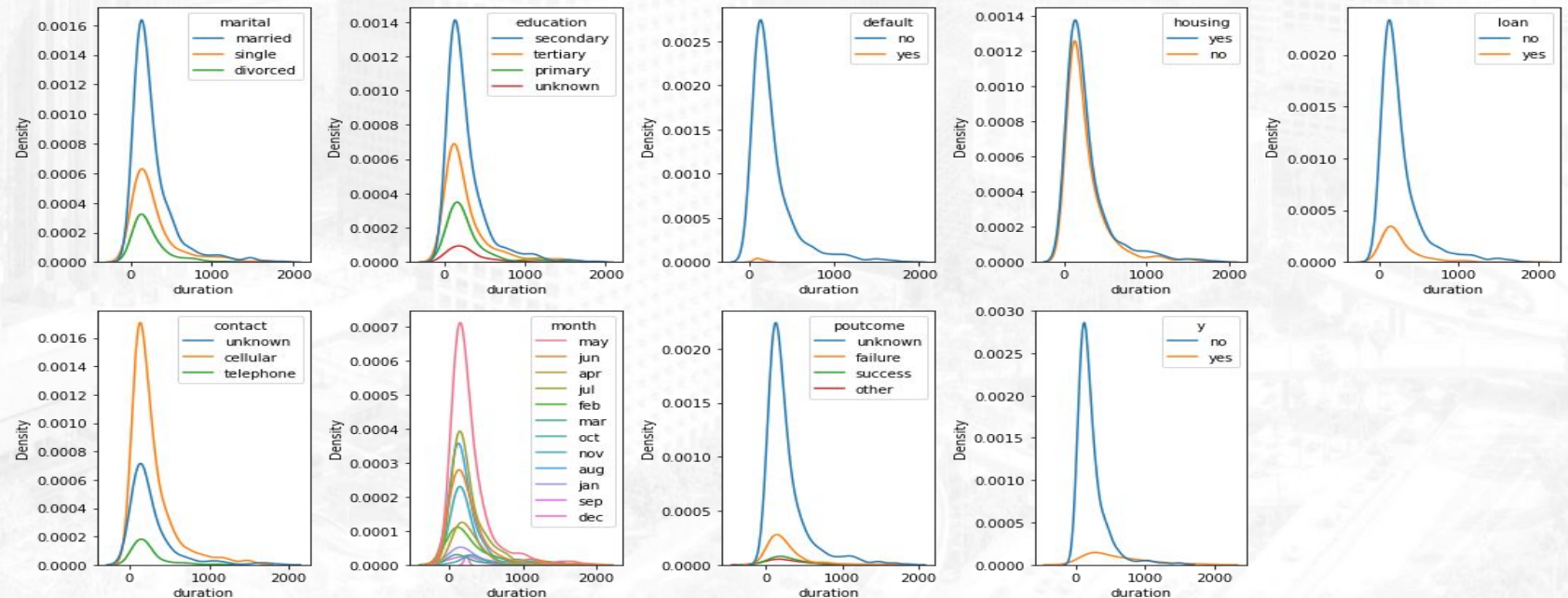
Korelasi Data Kategorik terhadap y

Belum nampak ada korelasi yg kuat pada data kategorik terhadap output target.



Korelasi data kategorik terhadap **duration**

namun customer yang melakukan subscribed deposito ($y = \text{'yes'}$) cenderung memiliki durasi telepon yg lebih lama dibandingkan yang tidak subscribed.



Business Insight

- Dilihat pada heatmap, terdapat korelasi yg cukup tinggi antara target output (**y**) dengan **duration**. Jadi, untuk meningkatkan jumlah nasabah yang melakukan deposito dengan menambahkan durasi telepon kepada nasabah. Tetapi perlu diperhatikan juga semakin lama durasi telepon semakin besar biaya yang diperlukan.
- Pada **previous** dan **pdays** walaupun memiliki korelasi tapi sepertinya tidak ada kausalitas.
- Terdapat korelasi **campaign** terhadap **day**. Sehingga dapat dipilih hari-hari tertentu yang memiliki tingkat keberhasilan campaign yang tinggi, agar campaign dapat lebih efektif.

0s

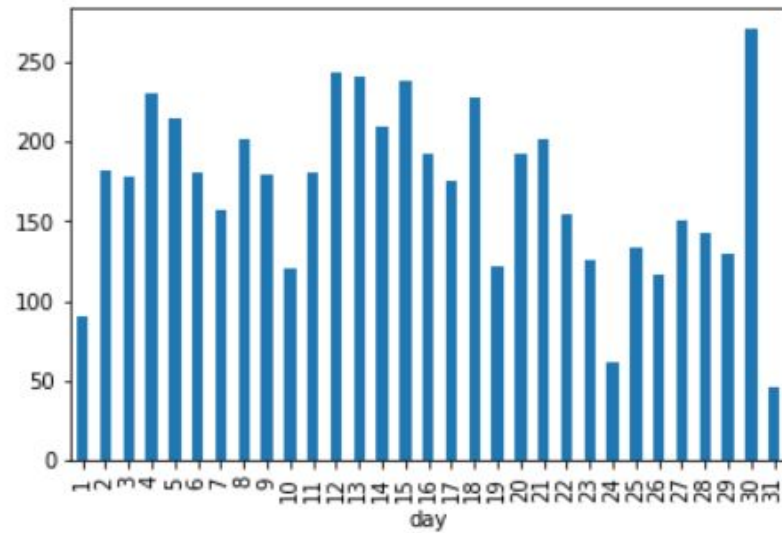


```
dy = df[df['y'] == 1]

dy_cpg = dy.groupby('day')['campaign'].count()

dy_cpg.plot(x='day', y='campaign', kind='bar')
plt.title("Campaign vs Day", loc='left', y=1.03, fontsize=15, weight='bold')
plt.show()
```

Campaign vs Day



campaign yg menghasilkan output positif (customer subscribed to deposit) jarang terjadi pada akhir bulan.

Langkah Meningkatkan Model

- a. Encode Data Yes dan No pada **y** menjadi boolean agar bisa lebih mudah melihat
- b. Ambil insight dari korelasi antara **y** dengan data kategorik serta korelasi **duration** dengan data kategorik.
- c. Handling outlier menggunakan metode manual filtering.
- d. Melakukan box-cox transformation pada **balance**, **duration**, **campaign**, dan **previous** kemudian lakukan normalisasi data.
- e.

Data Cleansing

A. Handle missing values

```
[ ] 1 df.isnull().sum()
```

```
age      0
job      0
marital  0
education 0
default  0
balance  0
housing  0
loan     0
contact  0
day      0
month    0
duration 0
campaign 0
pdays   0
previous 0
poutcome 0
y        0
dtype: int64
```

Data tidak ada yang null

B. Handle duplicated data

```
[ ] 1 df.duplicated().sum()
```

```
0
```

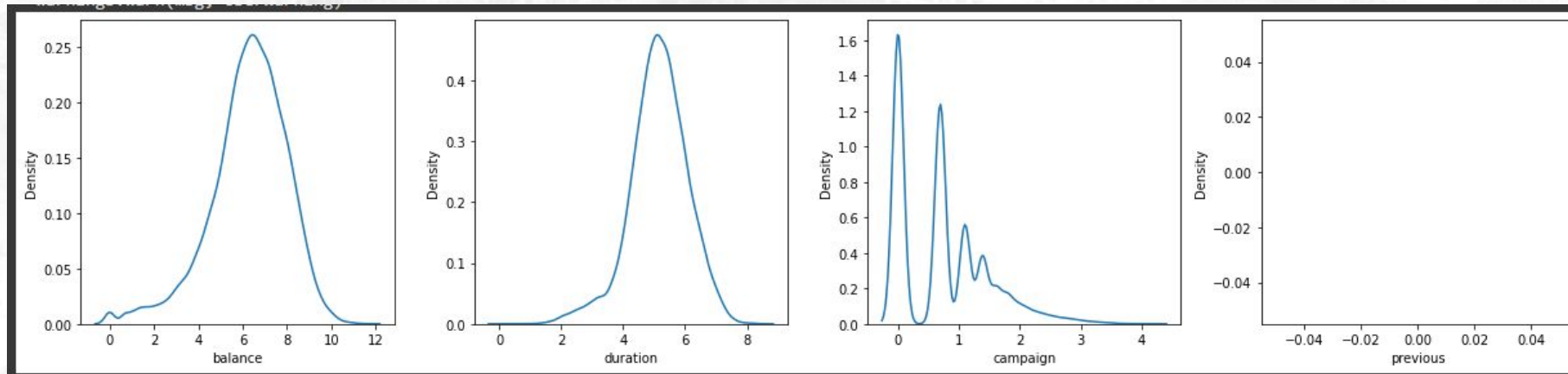
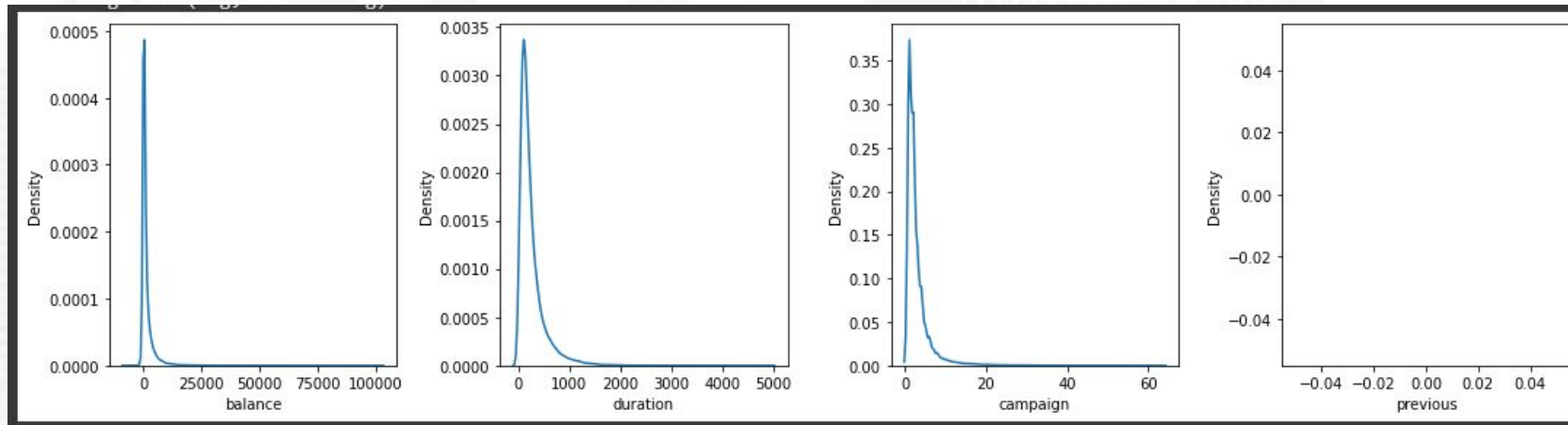
Data tidak ada yang duplikat

C. Handle outliers

```
1 from sklearn.preprocessing import MinMaxScaler, StandardScaler
2
3 df['age_norm'] = MinMaxScaler().fit_transform(df['age'].values.reshape(len(df), 1))
4 df['day_norm'] = MinMaxScaler().fit_transform(df['day'].values.reshape(len(df), 1))
5 df['campaign_norm'] = MinMaxScaler().fit_transform(df['campaign'].values.reshape(len(df), 1))
6 df['pdays_norm'] = MinMaxScaler().fit_transform(df['pdays'].values.reshape(len(df), 1))
7
8
9 df.describe()
```

	age	balance	day	duration	campaign	pdays	previous	y	log_balance	log_duration	age_norm	day_norm	campaign_norm	pdays_norm
count	36954.000000	36954.000000	36954.000000	36954.000000	36954.000000	36954.0	36954.0	36954.000000	3.367600e+04	3.695400e+04	36954.000000	36954.000000	36954.000000	36954.0
mean	40.932430	1318.788846	16.145424	257.726119	2.921957	-1.0	0.0	0.091573	-inf	-inf	0.297824	0.504847	0.030999	0.0
std	10.430218	3039.557077	8.372554	262.256406	3.325791	0.0	0.0	0.288427	NaN	NaN	0.135457	0.279085	0.053642	0.0
min	18.000000	-8019.000000	1.000000	0.000000	1.000000	-1.0	0.0	0.000000	-inf	-inf	0.000000	0.000000	0.000000	0.0
25%	33.000000	55.000000	9.000000	101.000000	1.000000	-1.0	0.0	0.000000	4.867534e+00	4.615121e+00	0.194805	0.266667	0.000000	0.0
50%	39.000000	414.000000	17.000000	177.000000	2.000000	-1.0	0.0	0.000000	6.240276e+00	5.176150e+00	0.272727	0.533333	0.016129	0.0
75%	49.000000	1358.000000	22.000000	318.000000	3.000000	-1.0	0.0	0.000000	7.333023e+00	5.762051e+00	0.402597	0.700000	0.032258	0.0
max	95.000000	102127.000000	31.000000	4918.000000	63.000000	-1.0	0.0	1.000000	1.153397e+01	8.500657e+00	1.000000	1.000000	1.000000	0.0

D. Feature transformation



	age	balance	day	duration	campaign	pdays	previous	y
count	36954.000000	36954.000000	36954.000000	36954.000000	36954.000000	36954.0	36954.0	36954.000000
mean	40.932430	1318.788846	16.145424	257.726119	2.921957	-1.0	0.0	0.091573
std	10.430218	3039.557077	8.372554	262.256406	3.325791	0.0	0.0	0.288427
min	18.000000	-8019.000000	1.000000	0.000000	1.000000	-1.0	0.0	0.000000
25%	33.000000	55.000000	9.000000	101.000000	1.000000	-1.0	0.0	0.000000
50%	39.000000	414.000000	17.000000	177.000000	2.000000	-1.0	0.0	0.000000
75%	49.000000	1358.000000	22.000000	318.000000	3.000000	-1.0	0.0	0.000000
max	95.000000	102127.000000	31.000000	4918.000000	63.000000	-1.0	0.0	1.000000

	age	balance	day	duration	campaign	pdays	previous	y	log_balance	log_duration
count	36954.000000	36954.000000	36954.000000	36954.000000	36954.000000	36954.0	36954.0	36954.000000	3.367600e+04	3.695400e+04
mean	40.932430	1318.788846	16.145424	257.726119	2.921957	-1.0	0.0	0.091573	-inf	-inf
std	10.430218	3039.557077	8.372554	262.256406	3.325791	0.0	0.0	0.288427	NaN	NaN
min	18.000000	-8019.000000	1.000000	0.000000	1.000000	-1.0	0.0	0.000000	-inf	-inf
25%	33.000000	55.000000	9.000000	101.000000	1.000000	-1.0	0.0	0.000000	4.867534e+00	4.615121e+00
50%	39.000000	414.000000	17.000000	177.000000	2.000000	-1.0	0.0	0.000000	6.240276e+00	5.176150e+00
75%	49.000000	1358.000000	22.000000	318.000000	3.000000	-1.0	0.0	0.000000	7.333023e+00	5.762051e+00
max	95.000000	102127.000000	31.000000	4918.000000	63.000000	-1.0	0.0	1.000000	1.153397e+01	8.500657e+00

E. Feature encoding

	age	job	education	default	balance	housing	loan	day	month	duration	...	day_norm	campaign_norm	pdays_norm	marital_divorced	marital_married	marital_single	contact_cellular	contact_telephone	contact_unknown	poutcome_unknown
0	58	management	3	0	2143	1	0	5	5	261	...	0.133333	0.0	0.0	0	1	0	0	0	1	1
1	44	technician	2	0	29	1	0	5	5	151	...	0.133333	0.0	0.0	0	0	1	0	0	1	1
2	33	entrepreneur	2	0	2	1	1	5	5	76	...	0.133333	0.0	0.0	0	1	0	0	0	1	1
3	47	blue-collar	0	0	1506	1	0	5	5	92	...	0.133333	0.0	0.0	0	1	0	0	0	1	1
4	33	unknown	0	0	1	0	0	5	5	198	...	0.133333	0.0	0.0	0	0	1	0	0	1	1

5 rows × 27 columns

F. Handle class imbalance



2. Feature Engineering

A. Feature selection

Beberapa feature yang di drop menggunakan korelasi

```
1 abs(df.corr()['y'])[abs(df.corr()['y'])>0.05].drop('y').index.tolist()

['housing',
 'duration',
 'campaign',
 'log_balance',
 'log_duration',
 'campaign_norm',
 'marital_married',
 'marital_single',
 'contact_cellular',
 'contact_unknown']
```

B. Feature extraction

duration_minute : Mengubah duration pada satuan detik menjadi menit agar lebih mudah dipahami

```

1 cost_perminute = 0.01
2 df3['cost'] = (df['duration'] * cost_perminute).round(2)
3 df3.head()

```

	age	job	education	default	balance	housing	loan	day	month	duration	...	pdays_norm	marital_divorced	marital_married	marital_single	contact_cellular	contact_telephone	contact_unknown	poutcome_unknown	duration_minute
0	58	management	3	0	2143	1	0	5	5	261	...	0.0	0	1	0	0	0	1	1	4.4
1	44	technician	2	0	29	1	0	5	5	151	...	0.0	0	0	1	0	0	1	1	2.5
2	33	entrepreneur	2	0	2	1	1	5	5	76	...	0.0	0	1	0	0	0	1	1	1.3
3	47	blue-collar	0	0	1506	1	0	5	5	92	...	0.0	0	1	0	0	0	1	1	1.5
4	33	unknown	0	0	1	0	0	5	5	198	...	0.0	0	0	1	0	0	1	1	3.3

5 rows x 29 columns

C. 4 feature tambahan :

- `duration_minute` : Mengubah duration pada satuan detik menjadi menit agar lebih mudah dipahami
- `cost` : Menambahkan fitur tambahan cost yaitu biaya tambahan dengan tarif 0.01 euro per menit, sehingga menjadi bisa menjadi pertimbangan perusahaan untuk memperhatikan duration
- `group_balance` : Mengkategorikan nasabah sesuai dengan balancenya
- `group_age` : Mengkategorikan nasabah sesuai dengan umurnya, agar dapat diperhatikan usia produktif untuk bekerja

<code>duration_minute</code>	<code>cost</code>	<code>group_balance</code>	<code>group_age</code>
4.4	2.61	High	Adults
2.5	1.51	Low	Adults
1.3	0.76	Low	Adults
1.5	0.92	High	Adults
3.3	1.98	Low	Adults