

SELAMAT DATANG DI TUGAS LEONARDO !!

Saya menggunakan Python karna pernah bersih-bersih string bareng temen di Python ehe.

► Packages

Karna Packages begitu berharga

[] ↵ 1 cell hidden

▼ Memasukkan Data

```
1 # Memasukkan Data
2 data = pd.read_csv('/content/raw data.csv')
3 # Melihat sekilas Data
4 data.head(5)
```

	Unnamed: 0	nama	harga	seller	terjual	tanggal_ambil_data
0	1	Bertolli Extra Virgin Olive Oil / Minyak Zaitu...	Rp85.000	Food Republic	Terjual 1.779 Produk	2020-09-04
1	2	Original Extra Virgin Olive Oil (EVOO) Casa Di...	Rp60.000	Yubyre	Terjual 1.624 Produk	2020-09-04
2	3	EVOO Baby Olive Oil Casa Di Oliva Olivia For K...	Rp58.000	Papamama Babyshop	Terjual 2.161 Produk	2020-09-04
3	4	Casa di Oliva - Extra Virgin Olive Oil for Kid...	Rp59.000	Chubby Baby Shop	Terjual 1.696 Produk	2020-09-04

Hemm... Kalau diliat-liat datanya kotor juga ya. Tapi sebelum dikomentari lebih jauh akan dilihat terlebih dahulu ukuran datanya

```
1 # Melihat Ukuran Data
2 print('Banyak Baris =', data.shape[0], 'Banyak Kolom =', data.shape[1])
```

Banyak Baris = 3259 Banyak Kolom = 6

```
1 # Melihat Struktur Data
2 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3259 entries, 0 to 3258
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Unnamed: 0             3259 non-null   int64
1   nama                   3259 non-null   object
2   harga                  3259 non-null   object
3   seller                 3259 non-null   object
4   terjual                3259 non-null   object
5   tanggal_ambil_data     3259 non-null   object
dtypes: int64(1), object(5)
memory usage: 152.9+ KB
```

Setelah dilihat-lihat ternyata datanya lengkap (tidak ada yang kosong) dan ukuran datanya lumayan besar. Berikutnya data akan dimanipulasi supaya bisa diambil sari/insight dari data tersebut.

Analisis Awal

Sekilas data tersebut memiliki kolom yang tidak perlu sehingga dapat dibuang.

Perhatikan juga bahwa kolom harga dan terjual memiliki tipe data object. Ini merupakan suatu masalah karna tidak bisa diolah.

Jika dilihat sekilas pada kolom nama dan seller, nama-nama terdiri dari string yang bervariasi (ada yang kapital ada yang tidak). Perhatikan bahwa pada python atau bahasa apapun huruf kapital dan tidak adalah 2 variabel yang berbeda. Sehingga perlu diubah menjadi kapital semua atau huruf kecil semua.

Kembali lihat kolom nama, terdapat 2 buah hal menarik yang dapat kita ambil, yakni merk dan volume dari jenis minyak.

Dapat juga dibuat suatu kolom baru yaitu Revenue yang diperoleh dari perkalian antara kolom terjual dengan kolom harga.

▼ Bersih - Bersih Biar Bisa Syukuran

Langkah penanganan Pertama akan dibuang kolom yang tidak perlu, menyamakan string nama dan seller (membuat semua string kapital atau tidak) dan string yang tidak membantu seperti Rp, terjual dan produk. Dan memastikan tanggal ambil

```
1 # Membuang Kolom Unnamed
2 data = data.drop('Unnamed: 0',axis=1)
3 data.head()
```

	nama	harga	seller	terjual	tanggal_ambil_data
0	Bertolli Extra Virgin Olive Oil / Minyak Zaitu...	Rp85.000	Food Republic	Terjual 1.779 Produk	2020-09-04

```
1 # Mengubah semua nama string menjadi huruf kecil
2 data['nama'] = data['nama'].str.lower()
3 data['seller'] = data['seller'].str.lower()
```

```
1 # Membuat data Harga Menjadi Integer
2 harga = data['harga'].str.split('Rp')
3 harga = harga.str[1]
4 harga = harga.str.replace('.', '')
5 harga = harga.astype(int)
6 data['harga'] = harga
7 data
```

	nama	harga	seller	terjual	tanggal_ambil_data
0	bertolli extra virgin olive oil / minyak zaitu...	85000	food republic	Terjual 1.779 Produk	2020-09-04
1	original extra virgin olive oil (evoo) casa di...	60000	yubyre	Terjual 1.624 Produk	2020-09-04
2	evoo baby olive oil casa di oliva olivia for k...	58000	papamama babyshop	Terjual 2.161 Produk	2020-09-04
3	casa di oliva - extra virgin olive oil for kid...	59000	chubby baby shop	Terjual 1.696 Produk	2020-09-04
4	best price casa di oliva extra virgin olive oi...	59900	whiz world	Terjual 3.115 Produk	2020-09-04
...
3254	minyak zaitun olive oil pak	70000	tokobundatitin	Terjual 34	2020-09-14

```
1 # Membuat data terjual Menjadi Integer
2 terjual =data['terjual'].str.split()
3 terjual = terjual.str[1]
4 terjual = terjual.str.replace('.', '')
5 terjual = terjual.astype(int)
6 data['terjual'] = terjual
7 data
```

	nama	harga	seller	terjual	tanggal_ambil_data
0	bertolli extra virgin olive oil / minyak zaitu...	85000	food republic	1779	2020-09-04
1	original extra virgin olive oil (evoo) casa di...	60000	yubyre	1624	2020-09-04
2	evoo baby olive oil casa di oliva olivia for k...	58000	papamama babyshop	2161	2020-09-04
3	casa di oliva - extra virgin olive oil for kid...	59000	chubby baby shop	1696	2020-09-04
4	best price casa di oliva extra virgin olive oi...	59900	whiz world	3115	2020-09-04
...

```

1 # Utak Atik Tanggal Biar Seru Ehe
2 tanggal = data['tanggal_ambil_data'].str.split('-')
3 tahun = tanggal.str[0]
4 bulan = tanggal.str[1]
5 hari = tanggal.str[2]
6 data['tahun'] = tahun.astype(str)
7 data['bulan'] = bulan.astype(str)
8 data['hari'] = hari.astype(str)

```

Merk Semua Olive Oil

```

1 # Karna di matematika kalau ga nguli ga asyik ye ga
2 utik = data['nama']
3 b = [0]*len(data)
4 merk = ['fillipo','filippo ','bertol','casa','borges','yummy bites', 'pak haji',
5         'mueloliva','zaituna','pomace','olio luglio','oilum','ayudya','olitalia',
6         'luglio', 'cobram', 'sasso','rafael salgado','riche', 'orkide',
7         'arbequina','nigella', 'mubarok','rs','syekh ali jaber','rumman',
8         'forever arctic','costa','selva','acropolis','coosur','bragg','al arobi',
9         'al afiat','ammara','coreysa','minyak kelapa voc','medina','co.e',
10        'la rambla','doug','ghuroba','deep olive oil','mumtaz','balsari',
11        'ecozest','amir','geofoods tartufo','orillia','naroop','afra','sesa',
12        'orilia canola','az zaitun','desert miracle','sufi','jadied','jadeed',
13        'espanola','altivolia','primarasa']
14 for i in range (len(merk)):
15     a = utik.str.find(merk[i])
16     for j in range(len(a)):
17         if a[j] != -1 :
18             b[j] = merk[i]
19     for k in range(len(a)):
20         if b[k] == 0:
21             b[k] = 'lainnya'
22 # Karena ada beberapa nama produk yang 'spesial'
23 benar = ['bertolli','costa d oro','filippo','al ghuroba','le riche','al amir','jade
24 for l in range (len(b)):
25     if b[l] == 'bertol':
26         b[l]= benar[0]

```

```

27 elif b[1] == 'costa':
28     b[1]=benar[1]
29 elif b[1] == 'fillipo':
30     b[1] = benar[2]
31 elif b[1] == 'ghuroba':
32     b[1] = benar[3]
33 elif b[1] == 'riche':
34     b[1] =benar[4]
35 elif b[1] == 'amir':
36     b[1] =benar[5]
37 elif b[1] == 'jadied':
38     b[1] =benar[6]
39 data['merk'] = b

1 # Menemukan volum dari nama
2 c = ['tidak tersedia']*len(data)
3 ukuran = ['250', '500', '1','325']
4
5 for i in range (len(ukuran)):
6     a = utik.str.find(ukuran[i])
7     for j in range(len(a)):
8         if a[j] != -1 :
9             c[j] = ukuran[i]
10         if c[j] == '1' :
11             c[j] = '1000'
12 data['ukuran'] = c

```

```

1 # Membuat kolom revenue
2 data['revenue'] = data['terjual']*data['harga']

```

Setelah bermalam-malam dan bergelas-gelas kopi yang telah diminum akhirnya datanya bersih juga. Perhatikan output dari cell dibawah ini

```

1 # Rekap
2 print('Banyaknya Seller ',len(data['seller'].unique()))
3 print('Banyaknya Barang ',len(data['nama'].unique()))
4 print('Banyaknya Merk ',len(data['merk'].unique()))
5 print('Range pengumpulan Data', data['tanggal_ambil_data'].unique())

Banyaknya Seller  211
Banyaknya Barang  416
Banyaknya Merk    57
Range pengumpulan Data ['2020-09-04' '2020-09-05' '2020-09-07' '2020-09-09' '2020-09-11' '2020-09-14']

```

Yang menarik dari data adalah, terdapat suatu produk yang bukan olive oil atau sejenisnya. Yaitu prima rasa. Produk tersebut adalah sambel untuk makan bukan untuk memasak. Jadi... tidak begitu signifikan jika diubang karna setelah diperiksa produk prima rasa ada 14 dari 3259 data. Jadi dibiarkan saja

▼ Pertanyaan Pertama

Pada rentang waktu tersebut, siapa yang menjadi market leader?

Yang menjadi Market Leader adalah Seller yang paling banyak melakukan penjualan. Sehingga langkah pertama adalah menghitung jumlah banyak terjual dari seller

```
1 market_leader = data.terjual.groupby(data['seller'],axis=0).sum().sort_values(ascen
2 market_leader.head(1)
```

```
seller
food republic    40954
Name: terjual, dtype: int64
```

▼ Pertanyaan Kedua

Brand minyak mana yang memiliki sales value terbesar saat big sale 9.9?

Pertama akan dikelompokkan pada waktu yang ditentukan dan akan dikelompokkan berdasarkan merk dan dijumlahkan jumlah terjualnya.

```
1 big_sale = data[data['tanggal_ambil_data']=='2020-09-09']
2 big_sale['revenue'] = big_sale['harga']*big_sale['terjual']
3 minyak_paling_laris = big_sale.groupby('merk',axis=0)
4 minyak_paling_laris['revenue'].max().sort_values(ascending=False).head(1)
```

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:2: SettingWithCopyWa
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <https://pandas.pydata.org/pandas-docs/stabl>

```
merk
bertolli    209455000
Name: revenue, dtype: int64
```

▼ Pertanyaan Ketiga

Bayangkan kalian bekerja di salah satu brand minyak tersebut, dari data ini hal apa yang bisa Anda sarankan ke tim sales / marketing?

Misalkan merk yang akan ditinjau adalah Bertolli.

Pertama akan dilihat terlebih dulu karakteristik dari data.

```

1 bertolli = data[data.merk=='bertolli']
2 print('Banyak Seller dari Merk Bertolli',len(bertolli.seller.unique()))
3 bertolli.seller.unique()

Banyak Seller dari Merk Bertolli 35
array(['food republic', 'ceva', 'ss suppliers f&b jakarta', 'bayininja',
      'bestcmart', 'momogo id', 'w-jaya store', 'warung mas budi',
      'theonlinegrocer', 'foodsupply.co', 'w2fitthealthy', 'herbs & co',
      'househerbal', 'genki plant', 'house of organix', 'niconicoshop',
      'kurakushop', 'serbasherbi', 'pd laris', 'the spices house',
      'secondition', 'jfm26', 'ampunmurahnya', 'greenara.id',
      'philocoffee', '888 seasoning', 'grandstand', 'provision master',
      'khas jaya nusantara', 'toko 3 rasa', "danieltyo's shop",
      "chef'schoice", 'choconola', 'disgro', 'dapoer canoli'],
      dtype=object)

```

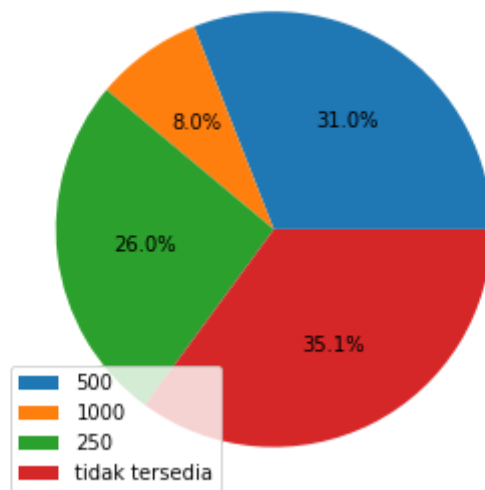
Berikutnya akan dianalisis ukuran dari merk bertolli yang dijual.

```

1 label_analisis = bertolli.ukuran.unique()
2 analisis = bertolli.groupby('ukuran')
3 plt.figure(figsize=(5,5))
4 plt.pie(analisis.revenue.sum(), autopct='%1.1f%%')
5 plt.legend(label_analisis,loc='lower left')
6 plt.title('Perbandingan Ukuran Minyak yang Paling Banyak dibeli pada Merk Bertolli')
7 plt.show()

```

Perbandingan Ukuran Minyak yang Paling Banyak dibeli pada Merk Bertolli



Perhatikan bahwa banyak dari produk Bertolli yang terjual masih kurang informasi mengenai ukuran dari produk yang dibeli.

Berikutnya akan dilihat pertumbuhan dari revenue Bertolli

```

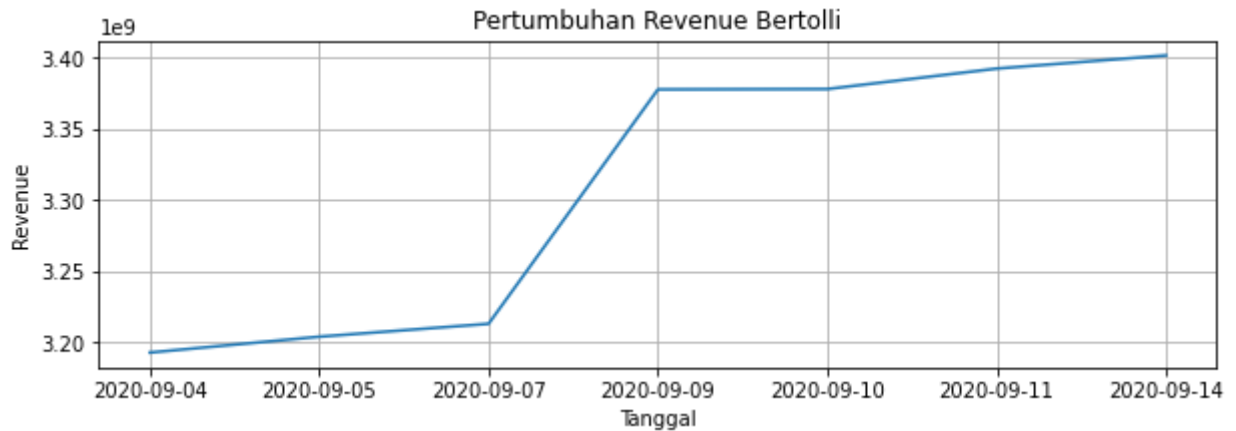
1 bertolli = bertolli.groupby('tanggal_ambil_data')
2 plt.figure(figsize=(10, 3))
3 plt.plot(bertolli['revenue'].sum())

```

```

4 plt.axis()
5 plt.title('Pertumbuhan Revenue Bertolli')
6 plt.xlabel('Tanggal')
7 plt.ylabel('Revenue')
8 plt.grid(True)
9 plt.show()

```



Setelah melihat data dari merk Bertolli akan dilihat data secara keseluruhan. Kita akan tinjau dari segi revenuenya.

```

1 merk = ['bertolli','casa','borges','lainnya']
2 for i in range(len(b)):
3     if b[i] != merk[0]:
4         if b[i] != merk[1]:
5             if b[i] !=merk[2]:
6                 b[i]=merk[3]
7 data13 =data
8 data13['merk'] = b

```

```

1 data_acuan = data13.groupby('merk')
2 labels = data_acuan['merk'].unique()
3
4 plt.figure(figsize=(5,5))
5 plt.pie(data_acuan['revenue'].sum().sort_values(ascending=False), autopct='%1.1f%%'
6 plt.legend(labels,loc='lower left')
7 plt.title('Perbandingan Reveneue antara 4 merk teratas')
8 plt.show()

```


Perbandingan Reveneue antara 4 merk teratas

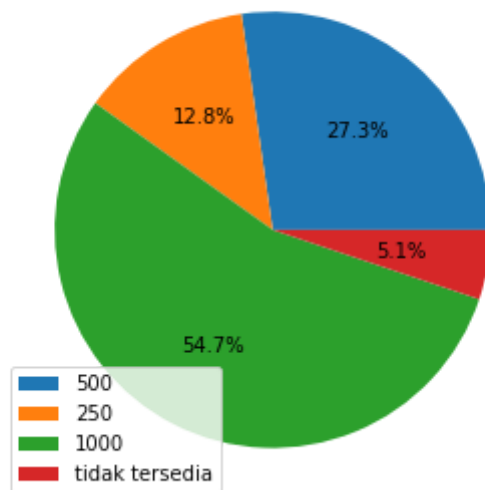


Ternyata merk Bertolli berbeda tipis dengan Borges. Akan dilihat data borges



```
1 borges = data[data.merk=='borges']
2 label_borges = borges.ukuran.unique()
3 analisis2 = borges.groupby('ukuran')
4 plt.figure(figsize=(5,5))
5 plt.pie(analisis2.revenue.sum(), autopct='%1.1f%%')
6 plt.legend(label_borges,loc='lower left')
7 plt.title('Perbandingan Ukuran Minyak yang Paling Banyak dibeli pada Merk Borges')
8 plt.show()
```

Perbandingan Ukuran Minyak yang Paling Banyak dibeli pada Merk Borges



Perhatikan bahwa banyak dari produk Borges yang terjual paling banyak adalah 1000ml dan perhatikan bahwa merk yang tidak memiliki ukuran hanya 5.1%.

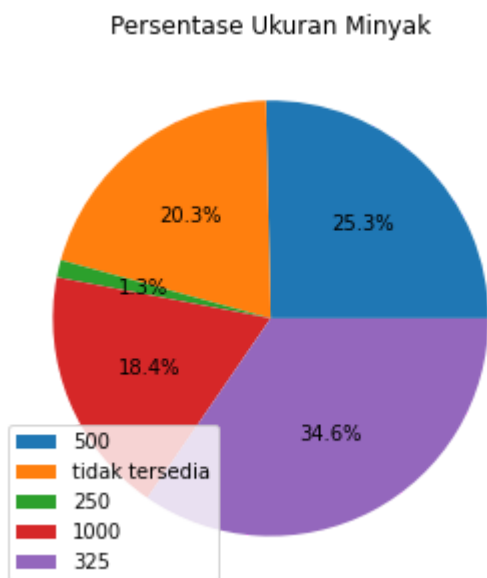
```
1 borges = borges.groupby('tanggal_ambil_data')
2 plt.figure(figsize=(10, 3))
3 plt.plot(borges['revenue'].sum())
4 plt.axis()
5 plt.title('Pertumbuhan Revenue Borges')
6 plt.xlabel('Tanggal')
7 plt.ylabel('Revenue')
8 plt.grid(True)
9 plt.show()
```



Perhatikan kembali bahwa merk Borges mengalami penurunan yang signifikan pada tanggal 9 September 2020. Hal ini berkebalikan pada merk borges, karena secara umum revenue borges terus naik. Kemungkinan besar, merk borges mengalami penurunan karena merk borges adalah barang inferior (Karena kemungkinan ada kecenderungan bahwa merk borges tidak menarik untuk dibeli pada event tertentu)

Berikutnya akan dicoba untuk dicocokkan antara analisis sementara dengan data yang lebih holistik

```
1 data_ukuran = pd.get_dummies(data['ukuran'])
2 labels1 = data['ukuran'].unique()
3 plt.figure(figsize=(5,5))
4 plt.pie(data_ukuran.sum(),autopct='%1.1f%%')
5 plt.legend(labels1,loc='lower left')
6 plt.title('Persentase Ukuran Minyak')
7 plt.show()
```



```
1 q = data[data['ukuran']=='325']
2 print('Rataan harga dari minyak ukuran 325ml', q.harga.mean())
3 q.harga.unique()
```

```
Rataan harga dari minyak ukuran 325ml 47083.333333333336
array([43500, 42000, 48000, 58000, 48500, 42500])
```

```
1 plt.figure(figsize=(10, 3))
2 q = q.groupby('tanggal_ambil_data')
```

```

3 plt.plot(q['revenue'].sum())
4 plt.axis()
5 plt.title('Pertumbuhan Revenue Untuk Minyak yang Berukuran 325ml')
6 plt.xlabel('Tanggal')
7 plt.ylabel('Revenue')
8 plt.grid(True)
9 plt.show()

```



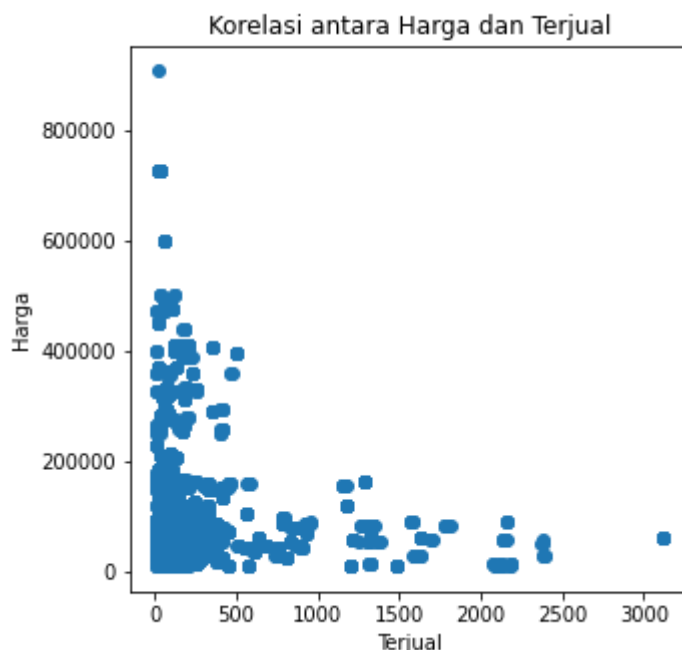
Perhatikan bahwa ukuran yang paling banyak laku adalah ukuran 325 ml dengan harga di kisaran 47000 rupiah hal ini berbeda dengan merk bertolli dan merk borges yang penjualan terbanyaknya bukan pada kedua merk tersebut.

Berikutnya akan dilihat apakah terdapat korelasi antara harga dan jumlah barang yang dijual

```

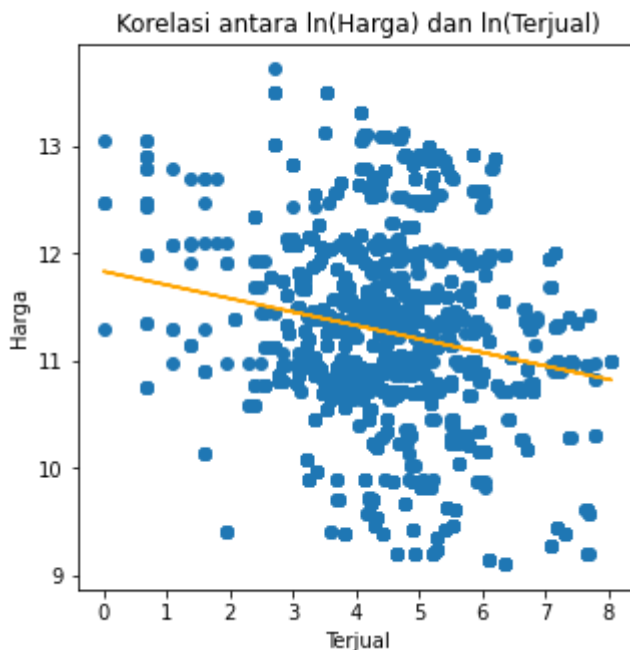
1 plt.figure(figsize=(5,5))
2 plt.scatter(x=data['terjual'],y=data['harga'])
3 plt.title('Korelasi antara Harga dan Terjual')
4 plt.ylabel('Harga')
5 plt.xlabel('Terjual')
6 plt.show()

```



Ternyata tidak ada korelasi antara harga dan produk yang terjual. Tapi jika data kita transformasikan dengan menggunakan log natural diperoleh.

```
1 m, b = np.polyfit(np.log(data['terjual']), np.log(data['harga']), 1)
2 plt.figure(figsize=(5,5))
3 plt.scatter(x=np.log(data['terjual']),y=np.log(data['harga']))
4 plt.plot(np.log(data['terjual']),m*np.log(data['terjual'])+b, color='orange')
5 plt.title('Korelasi antara ln(Harga) dan ln(Terjual)')
6 plt.ylabel('Harga')
7 plt.xlabel('Terjual')
8 plt.show()
```



Perhatikan bahwa sedikit terlihat ada korelasi negatif pada $\ln(\text{harga})$ terhadap $\ln(\text{terjual})$

Rekomendasi

Lebih gencar dengan penjualan dengan ukuran 325 ml.

▼ LEBIH LANJUT

Ke depannya ada beberapa hal yang bisa dilakukan, salah satunya adalah implementasi Machine Learning untuk menentukan merk dari barang-barang yang tidak memiliki merk maupun ukuran. Idennya adalah menggunakan kemiripan nama ataupun menggunakan nama seller maupun harga yang ditetapkan. Jika memenuhi parameter yang telah disampaikan maka dapat membuat data menjadi lebih baik untuk diolah.

Demikian tugas ini saya buat, saya merasa saya masih kurang mengerti dan kurang paham

```

1 #f = data1.nama[data1.merk=='lainnya']
2 #f.iloc[125]
3 #merk = data['nama'].str.replace('oil','')
4 #merk = merk.str.replace('virgin','')
5 #merk = merk.str.replace('minyak','')
6 #merk = merk.str.replace('olive','')
7 #
8 #merk = data['nama'].str.replace('best','')
9 #merk = merk.str.replace('price','')
10 #merk = merk.str.split(n=0)
11 #merk = merk.str[0]
12 #merk
13 #data['merk'] = merk
14 #data['merk'].unique()
15 #f = f.str.replace('oil','')
16 #f = f.str.replace('virgin','')
17 #f = f.str.replace('minyak','')
18 #f = f.str.replace('olive','')
19 #f = f.str.replace('zaitun','')
20 #f.unique()
21 #len(f)
22 #q = data['nama'].str.find('ayudya')
23 #for i in range (len(q)):
24 #   if q[i] != -1:
25 #       q[i] == 1
26 #for j in range(len(q)):
27 #   if q[j] == -1:
28 #       q.pop(j)
29 #q
30 #data['tanggal_ambil_data'] = pd.to_datetime(data['tanggal_ambil_data'],format='%Y-%m-%d')
31 #data['nama'].iloc[450]

```

