

# ***Full Cheatsheet Statistika untuk Data Science***

Disusun oleh Tim Datasans

*instagram : @datasans*

*medium: <https://datasans.medium.com/>*

## Daftar Isi

<b>Bab 1: Pengantar Statistik dalam Data Science.....</b>	<b>3</b>
1.1. Data Science dan Statistik: Hubungan dan Perbedaannya.....	3
1.2. Mengapa Statistik Penting dalam Data Science.....	4
1.3. Komponen Utama Statistik.....	5
1.4. Tipe Data dalam Statistik.....	6
<b>Bab 2: Statistik Deskriptif.....</b>	<b>8</b>
2.1. Ukuran Tendensi Sentral: Mean, Median, dan Modus.....	8
2.2. Ukuran Persebaran: Rentang, Variansi, dan Standar Deviasi.....	10
2.3. Ukuran Bentuk: Skewness dan Kurtosis.....	13
2.4. Visualisasi Data: Histogram, Box Plot, dan Scatter Plot.....	15
<b>Bab 3: Statistik Inferensial.....</b>	<b>21</b>
3.1. Pengantar Sampel dan Populasi.....	21
3.2. Teorema Limit Pusat dan Hukum Bilangan Besar.....	22
3.3. Estimasi Parameter: Titik dan Interval.....	27
3.4. Uji Hipotesis: Uji Z, Uji T, Uji Chi-Square, dan Uji F.....	30
3.4.1. Uji Z.....	30
3.4.2. Uji T.....	32
3.4.3. Uji Chi-Square.....	34
3.4.4. Uji F.....	35
<b>Bab 4: Analisis Regresi dan Korelasi.....</b>	<b>39</b>
4.1. Korelasi: Pearson, Spearman, dan Kendall.....	39
4.1.1. Korelasi Pearson.....	39
4.1.2. Korelasi Spearman.....	39
4.1.3. Korelasi Kendall.....	40
4.2. Regresi Linier Sederhana.....	41
4.3. Regresi Linier Berganda.....	44
4.4. Regresi Logistik.....	45
4.5. Evaluasi Model: R-Squared, Akurasi, dan Kesalahan.....	47
<b>Bab 5: Analisis Seri Waktu.....</b>	<b>50</b>
5.1. Komponen Seri Waktu: Tren, Musiman, dan Siklik.....	50
5.2. Teknik Smoothing: Moving Average dan Exponential Smoothing.....	53
5.2.1. Moving Average.....	53
5.2.2. Exponential Smoothing.....	55
5.3. Model Autoregressive Integrated Moving Average (ARIMA).....	58
5.4. Model Seasonal Decomposition of Time Series (STL).....	62
<b>Bab 6: Analisis Multivariat.....</b>	<b>66</b>
6.1. Analisis Komponen Utama (PCA).....	66
6.2. Analisis Faktor.....	68
6.3 Analisis Kluster.....	71
6.3.1 Dasar-dasar Analisis Kluster.....	71
6.3.2 Metode Pengelompokan yang Umum Digunakan.....	71
6.3.3 Contoh Kasus dan Implementasi Python.....	72

6.3.4 Tips dan Trik.....	73
6.4. Analisis Diskriminan.....	74
<b>Bab 7: Studi Kasus dan Aplikasi dalam Industri.....</b>	<b>80</b>
7.1. Prediksi Penjualan dengan Regresi Linier.....	80
7.2. Segmentasi Pelanggan dengan Klustering.....	83
7.3. Peramalan Stok dengan Analisis Seri Waktu.....	87
<b>Bab 8: Kesimpulan dan Langkah Selanjutnya.....</b>	<b>93</b>
8.1. Menjadi Ahli Statistik dalam Data Science.....	93
8.2. Pengembangan Karir dan Sertifikasi.....	94
8.3. Sumber Daya Tambahan dan Referensi.....	96

# Bab 1: Pengantar Statistik dalam Data Science

## 1.1. Data Science dan Statistik: Hubungan dan Perbedaannya

Data Science adalah disiplin ilmu yang menggabungkan berbagai bidang, termasuk statistik, matematika, ilmu komputer, dan teknik analisis data, untuk mengolah dan memahami data dalam rangka menghasilkan informasi yang bermanfaat dan mengambil keputusan yang lebih baik. Statistik, di sisi lain, merupakan cabang matematika yang berfokus pada pengumpulan, analisis, interpretasi, presentasi, dan organisasi data.

Dalam era Big Data saat ini, banyak organisasi telah mengumpulkan jumlah data yang sangat besar. Data ini bisa berasal dari berbagai sumber, seperti sensor, transaksi e-commerce, jejaring sosial, dan banyak lagi. Mengolah dan menginterpretasi data ini merupakan tantangan tersendiri yang memerlukan kombinasi dari berbagai keterampilan dan teknik. Di sinilah peran Data Science dan Statistik menjadi penting.

Hubungan antara Data Science dan Statistik bisa dijelaskan sebagai berikut:

**Interaksi:** Data Science memanfaatkan metode dan teknik statistik untuk mengolah data, menghasilkan informasi, dan membuat model prediksi. Sebagai contoh, metode regresi, klasifikasi, dan klustering yang digunakan dalam Data Science merupakan konsep yang berasal dari statistik.

**Komplementer:** Statistik dan Data Science saling melengkapi dalam mengatasi berbagai masalah yang berkaitan dengan data. Dalam banyak kasus, Data Science melibatkan teknik analisis data yang lebih canggih, seperti machine learning dan deep learning, yang didasarkan pada prinsip-prinsip statistik. Sebaliknya, statistik dapat memberikan dasar yang kuat untuk memahami metode dan algoritma yang digunakan dalam Data Science.

**Perluasan:** Data Science dapat dianggap sebagai perluasan dari statistik, di mana ilmuwan data menggunakan metode statistik, serta teknik dari ilmu komputer dan bidang lainnya, untuk menghadapi tantangan analisis data yang lebih kompleks dan besar.

Namun, ada beberapa perbedaan mendasar antara Data Science dan Statistik, seperti:

**Tujuan:** Tujuan utama Data Science adalah untuk mengekstrak pengetahuan dan wawasan dari data, sedangkan statistik lebih fokus pada pengujian hipotesis dan estimasi parameter yang berkaitan dengan populasi.

**Pendekatan:** Data Science menggunakan pendekatan yang lebih holistik dan interdisipliner, melibatkan teknik dari berbagai bidang, seperti ilmu komputer, visualisasi data, dan analisis teks. Sebaliknya, statistik lebih fokus pada metode matematika dan analisis data.

**Teknik:** Data Science sering menggunakan teknik yang lebih canggih, seperti machine learning, deep learning, dan natural language processing, yang mungkin tidak selalu digunakan dalam statistik.

**Skala Data:** Data Science biasanya berurusan dengan data yang lebih besar, lebih kompleks, dan tidak terstruktur, sedangkan statistik sering menggunakan data yang lebih kecil, lebih terstruktur, dan lebih mudah diinterpretasi.

## **1.2. Mengapa Statistik Penting dalam Data Science**

Statistik merupakan salah satu pilar utama dalam Data Science. Dalam konteks analisis data, statistik berfungsi untuk mengolah, menganalisis, dan menginterpretasi data untuk menghasilkan informasi yang bermanfaat dan mengambil keputusan yang lebih baik. Berikut adalah beberapa alasan mengapa statistik penting dalam Data Science:

**Eksplorasi Data:** Statistik deskriptif dan visualisasi data membantu ilmuwan data dalam mengidentifikasi pola, tren, dan hubungan dalam data, serta untuk memahami distribusi dan sifat-sifat data yang akan dianalisis. Proses eksplorasi data ini penting untuk menentukan langkah-langkah selanjutnya dalam analisis data dan memilih metode yang paling sesuai.

**Pengambilan Keputusan:** Statistik inferensial memungkinkan ilmuwan data untuk mengambil kesimpulan tentang populasi berdasarkan sampel data yang telah dikumpulkan. Hal ini sangat penting dalam pengambilan keputusan, terutama ketika data tidak lengkap atau terbatas. Statistik inferensial membantu ilmuwan data dalam mengestimasi parameter populasi, menguji hipotesis, dan membuat prediksi.

**Pembuatan Model:** Statistik membantu ilmuwan data dalam mengembangkan model yang mampu menjelaskan hubungan antara variabel dalam data, serta membuat prediksi tentang variabel target. Beberapa teknik statistik yang umum digunakan dalam pembuatan model termasuk regresi linear, analisis varians (ANOVA), dan analisis komponen utama (PCA).

**Evaluasi Model:** Setelah model telah dikembangkan, statistik digunakan untuk menguji validitas dan kinerja model. Hal ini melibatkan penggunaan metrik evaluasi seperti akurasi, kesalahan, dan koefisien determinasi (R-squared), serta pengujian hipotesis untuk memastikan bahwa model tersebut mampu menjelaskan hubungan dalam data dengan baik.

Optimasi Model: Statistik juga digunakan untuk mengoptimalkan model dengan menyesuaikan parameter, memilih fitur yang relevan, dan menentukan metode yang paling efektif untuk mengatasi masalah seperti overfitting dan underfitting.

### **1.3. Komponen Utama Statistik**

Statistik terdiri dari beberapa komponen utama yang mencakup berbagai aspek analisis data, seperti:

**Pengumpulan Data:** Proses ini melibatkan perencanaan, desain, dan pengumpulan data dari berbagai sumber untuk tujuan analisis. Metode pengumpulan data meliputi survei, eksperimen, pengamatan, dan penggunaan data sekunder.

**Deskripsi Data:** Statistik deskriptif mencakup perhitungan nilai rata-rata, median, modus, kuartil, dan simpangan baku untuk menggambarkan karakteristik data. Visualisasi data, seperti histogram, box plot, dan scatter plot, juga digunakan untuk menggambarkan distribusi dan hubungan antara variabel dalam data.

**Inferensi Statistik:** Statistik inferensial digunakan untuk mengambil kesimpulan tentang populasi berdasarkan analisis sampel data. Ini melibatkan pengujian hipotesis, perhitungan interval kepercayaan, dan penggunaan metode seperti uji t, uji chi-kuadrat, dan uji F untuk menilai signifikansi statistik dari hubungan antara variabel.

**Model Statistik:** Model statistik adalah representasi matematika dari hubungan antara variabel dalam data. Contoh model statistik meliputi regresi linear, regresi logistik, analisis varians (ANOVA), dan model mixed-effects. Model-model ini digunakan untuk menjelaskan hubungan antara variabel, memprediksi nilai variabel target, dan mengendalikan variabel konfounding.

**Estimasi dan Prediksi:** Statistik digunakan untuk mengestimasi parameter populasi, seperti rata-rata, proporsi, dan perbedaan antara kelompok, berdasarkan data sampel. Selain itu, statistik juga digunakan untuk membuat prediksi tentang nilai variabel target, berdasarkan informasi yang tersedia dalam data.

**Validasi dan Evaluasi Model:** Setelah model statistik telah dikembangkan, penting untuk menguji validitas dan kinerja model menggunakan metrik evaluasi, seperti akurasi, kesalahan, dan koefisien determinasi (R-squared). Uji hipotesis dan analisis residu juga digunakan untuk memastikan bahwa model mampu menjelaskan hubungan dalam data dengan baik.

## 1.4. Tipe Data dalam Statistik

Memahami tipe data dalam statistik sangat penting, karena tipe data menentukan metode analisis yang sesuai dan teknik visualisasi yang digunakan. Data dalam statistik dapat diklasifikasikan menjadi beberapa tipe, seperti:

**Data Kualitatif:** Data kualitatif, atau data kategori, adalah data yang menggambarkan atribut atau karakteristik yang tidak dapat diukur secara numerik. Data kualitatif meliputi data nominal dan ordinal.

a. **Data Nominal:** Data nominal adalah data yang dikelompokkan berdasarkan atribut atau kategori yang tidak memiliki urutan atau peringkat tertentu. Contoh data nominal termasuk jenis kelamin, ras, dan agama.

b. **Data Ordinal:** Data ordinal adalah data yang dikelompokkan berdasarkan atribut atau kategori yang memiliki urutan atau peringkat tertentu. Contoh data ordinal termasuk tingkat pendidikan, tingkat kepuasan, dan skala peringkat.

**Data Kuantitatif:** Data kuantitatif, atau data numerik, adalah data yang dapat diukur atau dihitung dan memiliki makna numerik. Data kuantitatif meliputi data interval dan rasio.

a. **Data Interval:** Data interval adalah data yang memiliki skala numerik dengan selisih yang sama antara nilai yang berurutan, tetapi tidak memiliki titik nol absolut. Contoh data interval termasuk suhu dalam derajat Celsius dan tahun.

b. **Data Rasio:** Data rasio adalah data yang memiliki skala numerik dengan selisih yang sama antara nilai yang berurutan dan titik nol absolut. Contoh data rasio termasuk usia, pendapatan, dan jarak.

Memahami tipe data ini akan membantu kamu dalam memilih metode analisis yang tepat, menginterpretasi hasil dengan benar, dan mengkomunikasikan temuan analisis data secara efektif. Dengan menggabungkan pemahaman tentang tipe data ini dengan konsep dan teknik statistik, kamu akan memiliki dasar yang kuat untuk sukses dalam Data Science.

Sebagai kesimpulan, statistik memainkan peran penting dalam Data Science dan merupakan salah satu pilar utama dalam bidang ini. Statistik memberikan alat dan teknik yang diperlukan untuk mengumpulkan, mengolah, menganalisis, dan menginterpretasi data, serta mengambil keputusan yang lebih baik berdasarkan hasil analisis tersebut.

Pemahaman yang baik tentang konsep statistik, seperti eksplorasi data, inferensi, model, estimasi, dan prediksi, akan memungkinkan kamu untuk menghadapi berbagai tantangan dalam analisis data dan menghasilkan wawasan yang bermanfaat dari data yang kompleks dan besar. Selain itu, memahami tipe data dan bagaimana mereka mempengaruhi metode

analisis yang digunakan akan membantu kamu dalam memilih pendekatan yang tepat untuk setiap situasi analisis data.

Dalam bab-bab berikutnya, kita akan membahas topik-topik yang lebih mendalam dalam statistik dan Data Science, seperti metode analisis data, teknik visualisasi, algoritma machine learning, dan best practices dalam analisis data. Dengan memperoleh pengetahuan yang lebih mendalam tentang statistik dan Data Science, kamu akan memiliki keterampilan dan pemahaman yang diperlukan untuk menghadapi tantangan analisis data di dunia yang semakin kompleks dan data-driven.



# Bab 2: Statistik Deskriptif

## 2.1. Ukuran Tendensi Sentral: Mean, Median, dan Modus

Ukuran tendensi sentral merupakan metode yang digunakan untuk mengidentifikasi pusat distribusi data. Terdapat tiga ukuran tendensi sentral yang umum digunakan, yaitu mean, median, dan modus. Dalam subbab ini, kita akan membahas masing-masing ukuran tersebut secara lebih detail, menyajikan contoh kasus di dunia nyata, dan memberikan tips serta trik dalam penerapannya.

### Mean (Rata-rata)

Mean atau rata-rata merupakan salah satu ukuran tendensi sentral yang paling umum digunakan. Untuk menghitung mean, kita menjumlahkan semua nilai dalam kumpulan data dan membaginya dengan jumlah data. Namun, perlu diingat bahwa mean sangat sensitif terhadap outlier atau nilai ekstrim, sehingga dalam beberapa kasus, median atau modus mungkin lebih cocok digunakan sebagai ukuran tendensi sentral.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin menghitung rata-rata gaji karyawan di suatu perusahaan. Jika terdapat beberapa eksekutif dengan gaji yang sangat tinggi dibandingkan karyawan lainnya, mean akan cenderung lebih tinggi dari pada nilai yang sebenarnya mencerminkan gaji karyawan pada umumnya. Dalam kasus ini, median mungkin akan lebih cocok untuk menggambarkan gaji karyawan.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np

data = np.array([30_000, 35_000, 40_000, 45_000, 50_000, 55_000, 60_000, 100_000,
200_000])
mean = np.mean(data)
print("Mean =", mean)
```

Output:

```
Mean = 68333.33333333333
```

Tips dan Trik:

Gunakan mean jika data kamu simetris dan tidak terdapat outlier.

Jika terdapat outlier, pertimbangkan untuk menggunakan median atau modus sebagai ukuran tendensi sentral.

Mean sangat berguna dalam perhitungan statistik lainnya, seperti variansi dan standar deviasi.

## Median

Median adalah nilai tengah dalam kumpulan data yang telah diurutkan. Jika jumlah data ganjil, median adalah nilai yang ada di tengah. Jika jumlah data genap, median adalah rata-rata dari dua nilai tengah. Berbeda dengan mean, median tidak sensitif terhadap outlier atau nilai ekstrim.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin mengetahui nilai tengah dari harga rumah di suatu kota. Terdapat beberapa rumah dengan harga yang sangat tinggi atau sangat rendah, yang dapat mengakibatkan nilai mean menjadi tidak akurat. Dalam kasus ini, median lebih cocok untuk menggambarkan nilai tengah dari harga rumah.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np

data = np.array([100_000, 150_000, 200_000, 250_000, 300_000, 1_000_000])
median = np.median(data)
print("Median =", median)
```

Output:

```
Median = 225000.0
```

Tips dan Trik:

Gunakan median jika data kamu memiliki outlier atau jika distribusinya tidak simetris.

Median lebih robust terhadap outlier dan nilai ekstrim dibandingkan dengan mean.

Median dapat digunakan untuk menggambarkan tren sentral dalam data ketika mean tidak mencerminkan dengan baik karakteristik data tersebut.

## Modus

Modus adalah nilai yang paling sering muncul dalam kumpulan data. Modus bisa lebih dari satu jika ada beberapa nilai dengan frekuensi kemunculan yang sama. Modus sangat

berguna jika kamu ingin mengidentifikasi nilai yang paling umum dalam data, terutama jika data tersebut kategorikal.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin mengetahui jenis ponsel yang paling populer di sebuah toko. Data penjualan menunjukkan bahwa beberapa model ponsel terjual dengan frekuensi yang berbeda. Dalam kasus ini, modus akan membantu kita menentukan model ponsel yang paling laris.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
from scipy import stats

data = np.array(["iPhone", "Samsung", "Samsung", "iPhone", "iPhone", "Samsung",
                 "Samsung", "Samsung", "OnePlus"])
mode = stats.mode(data)
print("Mode =", mode.mode[0])
```

Output:

```
Mode = Samsung
```

Tips dan Trik:

Gunakan modus jika kamu ingin mengetahui nilai yang paling sering muncul dalam kumpulan data, terutama jika data tersebut kategorikal.

Modus tidak terpengaruh oleh outlier atau nilai ekstrim.

Dalam beberapa kasus, kamu mungkin perlu menggabungkan modus dengan mean atau median untuk memberikan gambaran yang lebih lengkap tentang data.

Dalam subbab ini, kita telah membahas tiga ukuran tendensi sentral utama, yaitu mean, median, dan modus. Setiap metode memiliki kelebihan dan kekurangannya, sehingga penting untuk memahami kapan harus menggunakan metode yang sesuai. Dengan memahami dan menerapkan ukuran tendensi sentral ini, kamu akan dapat mengidentifikasi pusat distribusi data dan membuat keputusan yang lebih baik berdasarkan analisis data kamu.

Selanjutnya, kita akan membahas tentang ukuran persebaran dalam statistik deskriptif, yang mencakup rentang, variansi, dan standar deviasi.

## 2.2. Ukuran Persebaran: Rentang, Variansi, dan Standar Deviasi

Ukuran persebaran digunakan untuk menggambarkan seberapa jauh data menyebar atau tersebar dalam kumpulan data. Ukuran persebaran yang umum digunakan adalah rentang, variansi, dan standar deviasi. Dalam subbab ini, kita akan membahas masing-masing ukuran persebaran secara lebih detail, menyajikan contoh kasus di dunia nyata, dan memberikan tips serta trik dalam penerapannya.

### Rentang (Range)

Rentang merupakan selisih antara nilai maksimum dan minimum dalam kumpulan data. Rentang menggambarkan jangkauan data secara keseluruhan, tetapi tidak memberikan informasi tentang bagaimana data tersebut tersebar. Rentang sangat sensitif terhadap outlier atau nilai ekstrim, sehingga dalam beberapa kasus, variansi atau standar deviasi mungkin lebih cocok digunakan sebagai ukuran persebaran.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin menghitung rentang harga rumah di suatu kota. Rentang harga rumah dapat membantu kita mengetahui jangkauan harga rumah secara keseluruhan, tetapi tidak menggambarkan bagaimana harga rumah tersebut tersebar dalam kota tersebut.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np

data = np.array([100_000, 150_000, 200_000, 250_000, 300_000, 1_000_000])
range_ = np.ptp(data)
print("Rentang =", range_)
```

Output:

```
Rentang = 900000
```

Tips dan Trik:

Gunakan rentang jika kamu ingin mengetahui jangkauan data secara keseluruhan. Rentang sangat sensitif terhadap outlier atau nilai ekstrim, jadi pertimbangkan untuk menggunakan variansi atau standar deviasi jika data memiliki outlier. Rentang tidak memberikan informasi tentang bagaimana data tersebar, jadi gunakan bersama dengan ukuran persebaran lainnya.

## Variansi

Variansi adalah ukuran persebaran yang menggambarkan seberapa jauh setiap nilai dalam kumpulan data dari mean. Variansi dihitung dengan menjumlahkan kuadrat selisih antara setiap nilai dan mean, kemudian dibagi dengan jumlah data minus satu. Variansi lebih robust terhadap outlier dibandingkan dengan rentang dan memberikan gambaran yang lebih baik tentang bagaimana data tersebar.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin menghitung variansi waktu tunggu pelanggan di suatu restoran. Variansi waktu tunggu akan memberi tahu kita seberapa besar variasi waktu tunggu antara pelanggan, yang dapat membantu kita mengidentifikasi apakah ada masalah dalam pelayanan restoran.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np

data = np.array([5, 10, 15, 20, 25, 30])
variance = np.var(data, ddof=1)
print("Variansi =", variance)
```

Output:

```
Variansi = 87.5
```

Tips dan Trik:

Gunakan variansi jika kamu ingin mengetahui seberapa jauh setiap nilai dalam kumpulan data dari mean.

Variansi lebih robust terhadap outlier dibandingkan dengan rentang dan memberikan gambaran yang lebih baik tentang bagaimana data tersebar.

Ingat bahwa variansi dihitung dengan satuan kuadrat dari data asli, sehingga mungkin sulit untuk diinterpretasikan. Dalam kasus seperti ini, pertimbangkan untuk menggunakan standar deviasi.

## Standar Deviasi

Standar deviasi adalah akar kuadrat dari variansi. Standar deviasi menggambarkan seberapa jauh setiap nilai dalam kumpulan data dari mean dalam satuan yang sama dengan data asli. Seperti variansi, standar deviasi lebih robust terhadap outlier dibandingkan dengan rentang dan memberikan gambaran yang lebih baik tentang bagaimana data tersebar.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin menghitung standar deviasi usia pelanggan di suatu toko. Standar deviasi usia akan memberi tahu kita seberapa besar variasi usia antara pelanggan, yang dapat membantu kita dalam menyusun strategi pemasaran yang sesuai dengan demografi pelanggan.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np

data = np.array([20, 25, 30, 35, 40, 45])
std_dev = np.std(data, ddof=1)
print("Standar Deviasi =", std_dev)
```

Output:

```
Standar Deviasi = 9.354143466934854
```

Tips dan Trik:

Gunakan standar deviasi jika kamu ingin mengetahui seberapa jauh setiap nilai dalam kumpulan data dari mean dalam satuan yang sama dengan data asli.

Standar deviasi lebih mudah diinterpretasikan dibandingkan dengan variansi karena menggunakan satuan yang sama dengan data asli.

Standar deviasi sangat berguna dalam analisis statistik, seperti dalam menghitung interval kepercayaan dan pengujian hipotesis.

Dalam subbab ini, kita telah membahas tiga ukuran persebaran utama, yaitu rentang, variansi, dan standar deviasi. Setiap metode memiliki kelebihan dan kekurangannya, sehingga penting untuk memahami kapan harus menggunakan metode yang sesuai. Dengan memahami dan menerapkan ukuran persebaran ini, kamu akan dapat mengidentifikasi seberapa jauh data menyebar atau tersebar dalam kumpulan data dan membuat keputusan yang lebih baik berdasarkan analisis data kamu.

Selanjutnya, kita akan membahas tentang ukuran bentuk dalam statistik deskriptif, yang mencakup skewness dan kurtosis.

## 2.3. Ukuran Bentuk: Skewness dan Kurtosis

Ukuran bentuk digunakan untuk menggambarkan simetri dan puncak distribusi data. Dua ukuran bentuk yang umum digunakan adalah skewness dan kurtosis. Skewness menggambarkan sejauh mana distribusi data condong atau miring, sementara kurtosis menggambarkan sejauh mana puncak distribusi data tajam atau runcing. Dalam subbab ini,

kita akan membahas masing-masing ukuran bentuk secara lebih detail, menyajikan contoh kasus di dunia nyata, dan memberikan tips serta trik dalam penerapannya.

## Skewness

Skewness adalah ukuran simetri distribusi data. Skewness positif menunjukkan bahwa distribusi condong ke kanan, dengan ekor distribusi membentang ke arah nilai yang lebih besar. Skewness negatif menunjukkan bahwa distribusi condong ke kiri, dengan ekor distribusi membentang ke arah nilai yang lebih kecil. Skewness nol menunjukkan distribusi yang simetris. Skewness penting untuk dipertimbangkan dalam analisis data karena dapat mempengaruhi interpretasi mean, median, dan modus.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin mengukur skewness pendapatan penduduk di suatu negara. Skewness pendapatan akan membantu kita mengetahui apakah distribusi pendapatan condong ke kiri (pendapatan rendah lebih banyak) atau ke kanan (pendapatan tinggi lebih banyak). Informasi ini berguna untuk menginformasikan kebijakan ekonomi dan fiskal.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
from scipy.stats import skew

data = np.array([20_000, 30_000, 40_000, 50_000, 60_000, 70_000, 100_000])
skewness = skew(data)
print("Skewness =", skewness)
```

Output:

```
Skewness = 0.560332836838461
```

Tips dan Trik:

Periksa skewness jika kamu ingin mengetahui simetri distribusi data.

Pertimbangkan skewness ketika menginterpretasi mean, median, dan modus, karena distribusi yang condong dapat mempengaruhi interpretasi nilai-nilai ini.

Skewness positif menunjukkan distribusi yang condong ke kanan, sedangkan skewness negatif menunjukkan distribusi yang condong ke kiri.

## Kurtosis

Kurtosis adalah ukuran puncak atau runcingan distribusi data. Kurtosis yang tinggi menunjukkan bahwa distribusi data lebih runcing dengan ekor yang lebih tebal, sementara

kurtosis yang rendah menunjukkan bahwa distribusi data lebih datar dengan ekor yang lebih tipis. Kurtosis penting untuk dipertimbangkan dalam analisis data karena dapat menggambarkan sejauh mana data memiliki nilai ekstrim atau outlier.

Contoh Kasus di Dunia Nyata:

Sebagai contoh, kita ingin mengukur kurtosis durasi pinjaman di suatu bank. Kurtosis durasi pinjaman akan membantu kita mengetahui apakah bank memiliki banyak pinjaman jangka pendek atau jangka panjang, serta apakah bank menghadapi risiko yang lebih tinggi dari pinjaman yang tidak terbayar atau gagal bayar. Informasi ini berguna untuk menginformasikan kebijakan manajemen risiko dan strategi bisnis bank.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
from scipy.stats import kurtosis

data = np.array([12, 24, 36, 48, 60, 72, 84, 96, 108, 120])
kurt = kurtosis(data, fisher=True)
print("Kurtosis =", kurt)
```

Output:

```
Kurtosis = -1.2242424242424244
```

Tips dan Trik:

Periksa kurtosis jika kamu ingin mengetahui sejauh mana distribusi data runcing atau memiliki nilai ekstrim.

Kurtosis yang tinggi menunjukkan distribusi yang lebih runcing dengan ekor yang lebih tebal, sementara kurtosis yang rendah menunjukkan distribusi yang lebih datar dengan ekor yang lebih tipis.

Pertimbangkan kurtosis ketika menganalisis risiko dalam data, seperti dalam contoh durasi pinjaman, karena kurtosis yang tinggi menunjukkan adanya nilai ekstrim yang mungkin mempengaruhi risiko.

Dalam subbab ini, kita telah membahas dua ukuran bentuk dalam statistik deskriptif, yaitu skewness dan kurtosis. Setiap metode memiliki kegunaan dan interpretasi yang berbeda, sehingga penting untuk memahami kapan harus menggunakan metode yang sesuai. Dengan memahami dan menerapkan ukuran bentuk ini, kamu akan dapat menggambarkan dengan lebih baik karakteristik distribusi data dan membuat keputusan yang lebih baik berdasarkan analisis data kamu.

Selanjutnya, kita akan membahas tentang visualisasi data dalam statistik deskriptif, yang mencakup histogram, box plot, dan scatter plot. Visualisasi data adalah alat yang sangat



berguna untuk menggambarkan dan memahami pola dalam data, serta untuk menyampaikan temuan analisis data kepada orang lain dengan cara yang mudah dimengerti.

## 2.4. Visualisasi Data: Histogram, Box Plot, dan Scatter Plot

Visualisasi data adalah alat yang sangat berguna untuk menggambarkan dan memahami pola dalam data, serta untuk menyampaikan temuan analisis data kepada orang lain dengan cara yang mudah dimengerti. Dalam subbab ini, kita akan membahas tiga teknik visualisasi data yang umum digunakan dalam statistik deskriptif: histogram, box plot, dan scatter plot. Kita akan membahas masing-masing teknik secara lebih detail, menyajikan contoh kasus di dunia nyata, dan memberikan tips serta trik dalam penerapannya.

### Histogram

Histogram adalah representasi grafis dari distribusi data. Histogram dibuat dengan membagi data ke dalam interval atau "bin" dan menghitung jumlah pengamatan dalam setiap interval. Batang-batang vertikal dalam histogram menggambarkan jumlah pengamatan dalam setiap bin. Histogram berguna untuk menggambarkan bentuk, simetri, dan kemencengan distribusi data.

Contoh Kasus di Dunia Nyata:

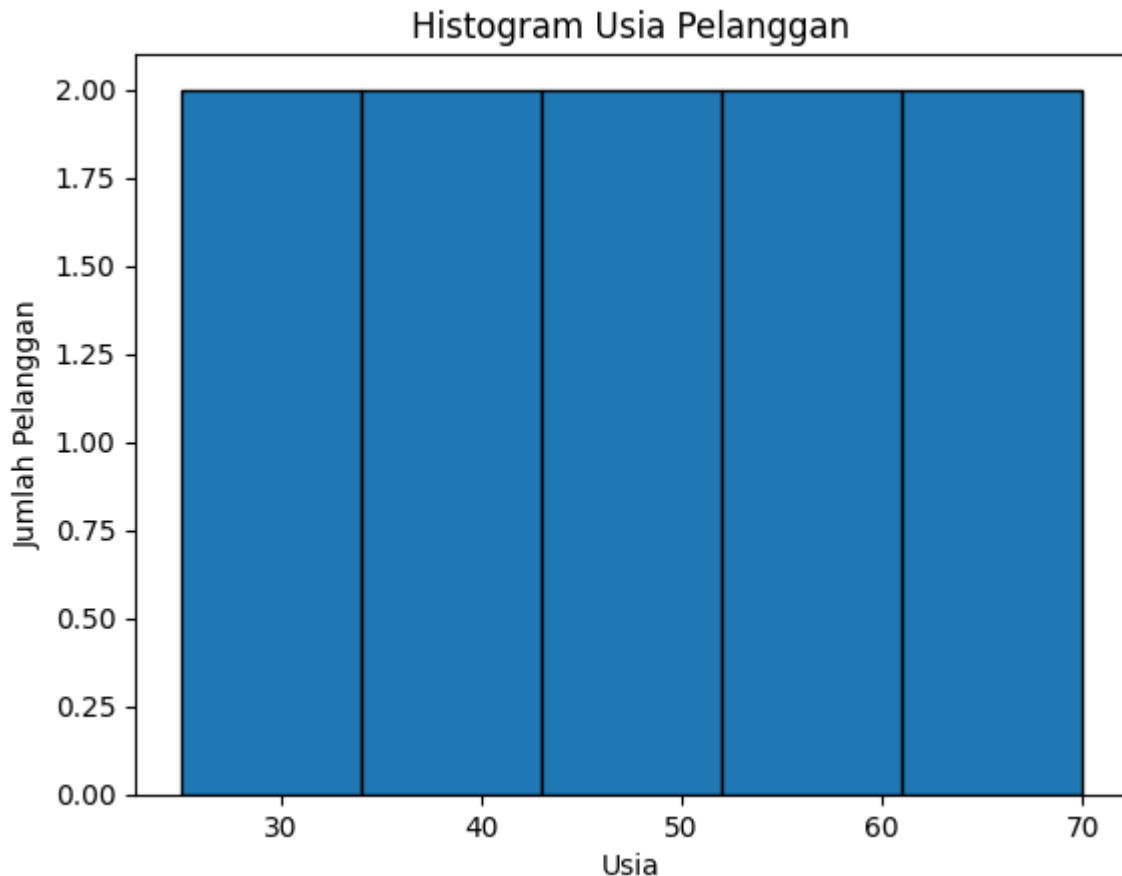
Misalkan kita ingin menggambarkan distribusi usia pelanggan di suatu toko. Dalam kasus ini, kita dapat menggunakan histogram untuk memvisualisasikan sebaran usia pelanggan dan mengidentifikasi kelompok usia yang paling umum atau jarang di toko tersebut. Informasi ini berguna untuk menargetkan promosi dan strategi pemasaran yang lebih efektif.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import matplotlib.pyplot as plt

data = np.array([25, 30, 35, 40, 45, 50, 55, 60, 65, 70])
plt.hist(data, bins=5, edgecolor='black')
plt.xlabel('Usia')
plt.ylabel('Jumlah Pelanggan')
plt.title('Histogram Usia Pelanggan')
plt.show()
```

Output:



#### Tips dan Trik:

Gunakan histogram untuk menggambarkan distribusi data.

Pilih jumlah bin yang tepat untuk menghasilkan histogram yang informatif dan mudah dibaca.

Perhatikan bentuk, simetri, dan kemencengan distribusi saat menginterpretasi histogram.

#### Box Plot

Box plot, atau diagram kotak, adalah representasi grafis yang menggambarkan sebaran data melalui lima statistik ringkasan: minimum, kuartil pertama (Q1), median (Q2), kuartil ketiga (Q3), dan maksimum. Box plot berguna untuk menggambarkan sebaran data, mengidentifikasi nilai ekstrim, dan mengevaluasi simetri distribusi.

#### Contoh Kasus di Dunia Nyata:

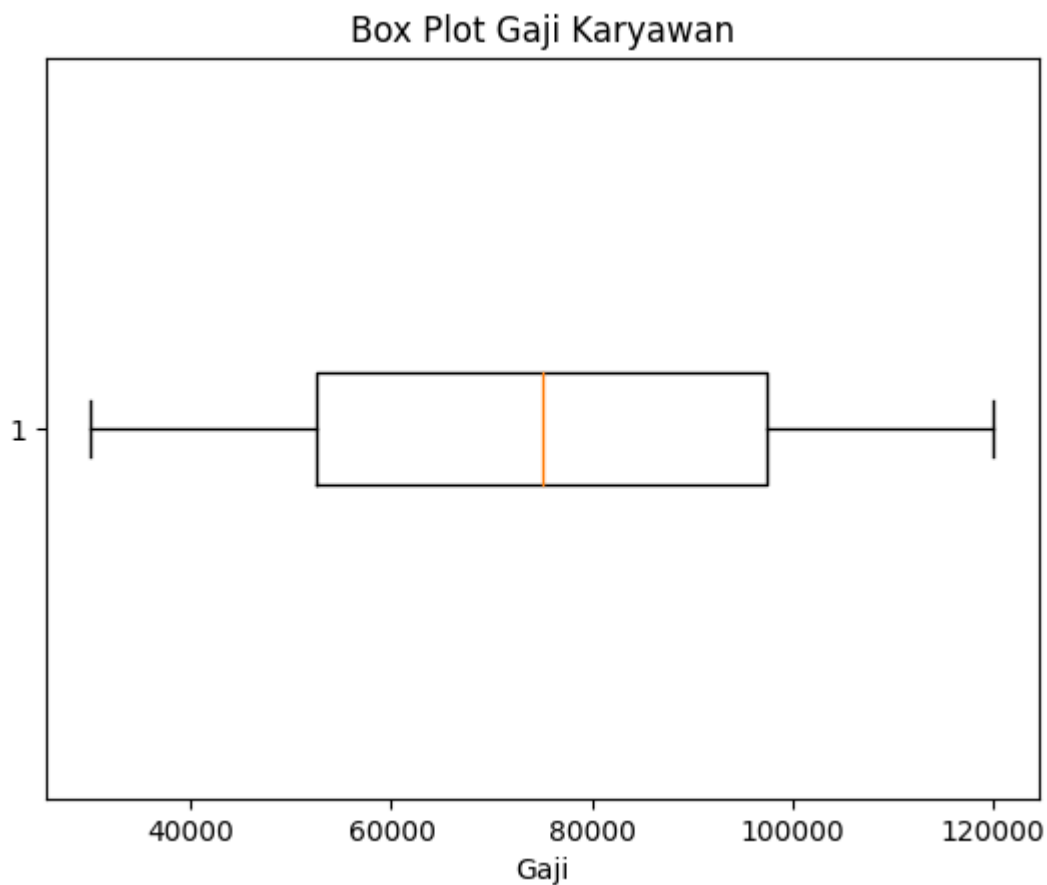
Misalkan kita ingin menggambarkan sebaran gaji karyawan di suatu perusahaan. Dalam kasus ini, kita dapat menggunakan box plot untuk memvisualisasikan sebaran gaji, mengidentifikasi gaji yang sangat tinggi atau sangat rendah, dan mengevaluasi simetri distribusi gaji. Informasi ini berguna untuk menginformasikan kebijakan penggajian dan manajemen sumber daya manusia.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import matplotlib.pyplot as plt

data = np.array([30_000, 40_000, 50_000, 60_000, 70_000, 80_000, 90_000, 100_000, 110_000,
120_000])
plt.boxplot(data, vert=False)
plt.xlabel('Gaji')
plt.title('Box Plot Gaji Karyawan')
plt.show()
```

Output:



Tips dan Trik:

- Gunakan box plot untuk menggambarkan sebaran data dan mengidentifikasi nilai ekstrim.
- Perhatikan panjang "kumis" (whiskers) dan jarak antara kuartil untuk mengevaluasi sebaran data.
- Perhatikan posisi median relatif terhadap kuartil pertama dan ketiga untuk menilai simetri distribusi.

## Scatter Plot

Scatter plot, atau diagram pencar, adalah representasi grafis dari hubungan antara dua variabel kontinu. Titik-titik dalam scatter plot menggambarkan kombinasi nilai untuk kedua variabel. Scatter plot berguna untuk mengidentifikasi pola hubungan antara variabel, seperti hubungan linier, kuadratik, atau korelasi non-linier.

Contoh Kasus di Dunia Nyata:

Misalkan kita ingin menggambarkan hubungan antara pengeluaran iklan dan penjualan di suatu perusahaan. Dalam kasus ini, kita dapat menggunakan scatter plot untuk memvisualisasikan hubungan antara variabel ini dan mengidentifikasi pola yang mungkin menunjukkan adanya korelasi antara pengeluaran iklan dan penjualan.

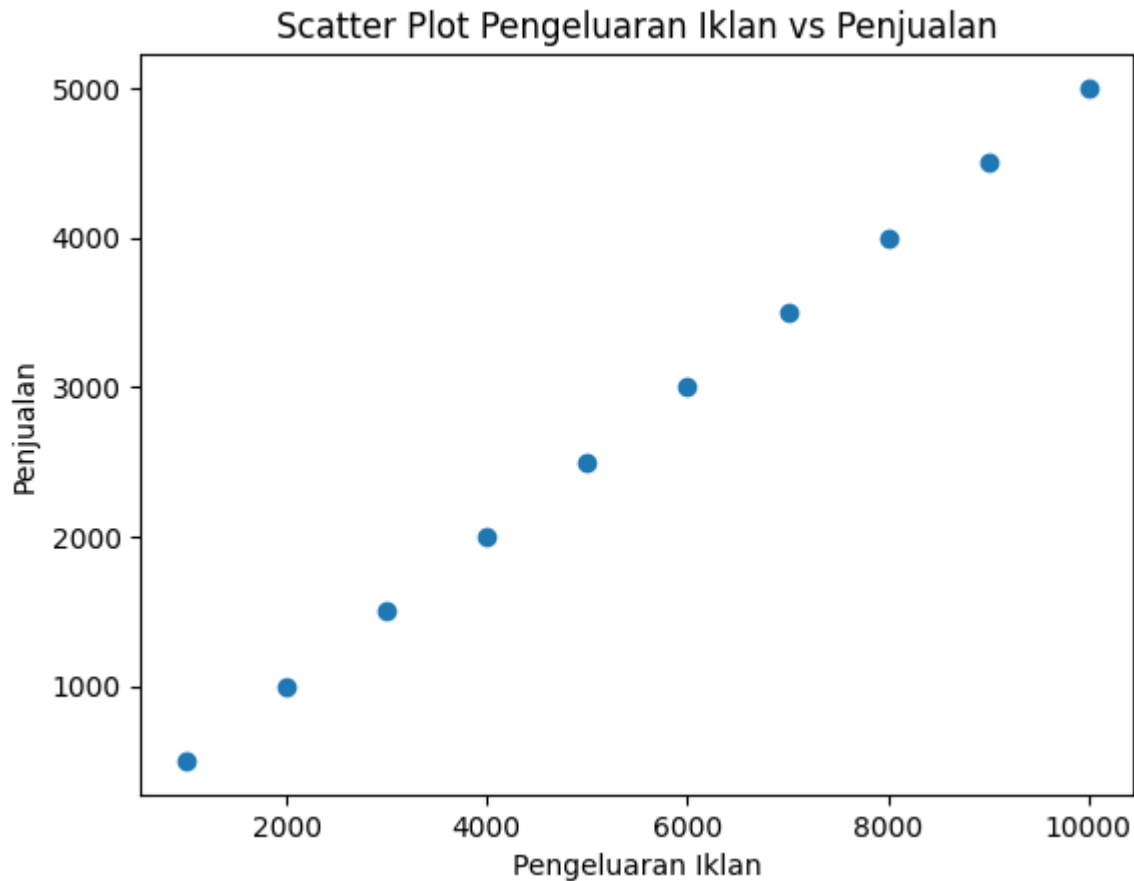
Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import matplotlib.pyplot as plt

ad_spend = np.array([1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10000])
sales = np.array([500, 1000, 1500, 2000, 2500, 3000, 3500, 4000, 4500, 5000])

plt.scatter(ad_spend, sales)
plt.xlabel('Pengeluaran Iklan')
plt.ylabel('Penjualan')
plt.title('Scatter Plot Pengeluaran Iklan vs Penjualan')
plt.show()
```

Output:



#### Tips dan Trik:

Gunakan scatter plot untuk menggambarkan hubungan antara dua variabel kontinu.

Perhatikan pola dalam scatter plot untuk mengidentifikasi korelasi linier, kuadratik, atau non-linier antara variabel.

Jangan berasumsi ada hubungan sebab-akibat antara variabel hanya berdasarkan pola dalam scatter plot.

Dalam subbab ini, kita telah membahas tiga teknik visualisasi data yang umum digunakan dalam statistik deskriptif: histogram, box plot, dan scatter plot. Setiap teknik memiliki kegunaan dan interpretasi yang berbeda, sehingga penting untuk memahami kapan harus menggunakan teknik yang sesuai. Dengan memahami dan menerapkan teknik visualisasi ini, kamu akan dapat menggambarkan dan memahami pola dalam data dengan lebih baik, serta menyampaikan temuan analisis data kepada orang lain dengan cara yang mudah dimengerti.

# Bab 3: Statistik Inferensial

## 3.1. Pengantar Sampel dan Populasi

Statistik inferensial adalah cabang statistik yang berfokus pada mengambil kesimpulan dari suatu sampel untuk menggambarkan karakteristik populasi yang lebih besar. Dalam subbab ini, kita akan membahas konsep dasar sampel dan populasi, serta pentingnya mengambil sampel yang baik dan representatif. Kita akan menyajikan contoh kasus di dunia nyata dan memberikan tips serta trik dalam penerapannya.

### Populasi

Populasi adalah kumpulan semua elemen yang ingin kita pelajari. Elemen ini bisa berupa individu, barang, atau pengukuran yang relevan dengan tujuan penelitian kita. Dalam statistik, populasi seringkali terlalu besar untuk diukur secara keseluruhan, sehingga kita harus mengambil sampel dari populasi tersebut untuk mengestimasi karakteristiknya.

Contoh Kasus di Dunia Nyata:

Misalkan kita ingin mengetahui rata-rata tinggi badan semua siswa SMA di suatu kota. Dalam kasus ini, populasi yang ingin kita pelajari adalah tinggi badan semua siswa SMA di kota tersebut.

### Sampel

Sampel adalah sekelompok elemen yang diambil dari populasi yang lebih besar. Kita menggunakan sampel untuk mengestimasi karakteristik populasi, seperti mean, median, atau proporsi. Penting untuk mengambil sampel yang baik dan representatif agar estimasi kita akurat dan memiliki tingkat kepercayaan yang tinggi.

Contoh Kasus di Dunia Nyata:

Menggunakan contoh sebelumnya tentang tinggi badan siswa SMA, kita mungkin tidak bisa mengukur tinggi badan semua siswa di kota tersebut. Oleh karena itu, kita perlu mengambil sampel siswa, mengukur tinggi badan mereka, dan kemudian menggunakan informasi tersebut untuk mengestimasi rata-rata tinggi badan populasi siswa SMA secara keseluruhan.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np

# Buat populasi sintetis
population = np.random.normal(loc=170, scale=10, size=10000)
```

```
# Ambil sampel acak dari populasi
sample = np.random.choice(population, size=100, replace=False)

# Hitung rata-rata populasi dan sampel
population_mean = np.mean(population)
sample_mean = np.mean(sample)

print(f"Rata-rata populasi: {population_mean:.2f}")
print(f"Rata-rata sampel: {sample_mean:.2f}")
```

Output:

```
Rata-rata populasi: 169.99
Rata-rata sampel: 170.69
```

Tips dan Trik:

Pastikan sampel kamu representatif dari populasi. Jika sampel kamu tidak representatif, estimasi kamu mungkin memiliki bias dan kurang akurat.

Gunakan metode pengambilan sampel yang sesuai, seperti random sampling, stratified sampling, atau cluster sampling, tergantung pada karakteristik populasi dan tujuan penelitian kamu.

Perhatikan ukuran sampel: sampel yang terlalu kecil mungkin tidak cukup untuk menghasilkan estimasi yang akurat, sedangkan sampel yang terlalu besar mungkin memerlukan sumber daya yang tidak efisien.

Dalam subbab ini, kita telah membahas konsep dasar sampel dan populasi dalam konteks statistik inferensial. Mengambil sampel yang baik dan representatif sangat penting untuk menghasilkan estimasi yang akurat dan memiliki tingkat kepercayaan yang tinggi. Dalam prakteknya, penting untuk mempertimbangkan metode pengambilan sampel yang sesuai, ukuran sampel yang tepat, dan kualitas data yang dihasilkan.

## 3.2. Teorema Limit Pusat dan Hukum Bilangan Besar

Teorema Limit Pusat (Central Limit Theorem - CLT) dan Hukum Bilangan Besar (Law of Large Numbers - LLN) merupakan dua konsep penting dalam statistik inferensial yang memungkinkan kita untuk membuat estimasi dan generalisasi tentang populasi berdasarkan sampel yang diambil. Pada subbab ini, kita akan menjelajahi kedua konsep ini secara mendalam, memberikan penjelasan matematis, contoh kasus di dunia nyata, dan menyajikan contoh Python script dengan data sintetis. Selain itu, kita akan menyajikan tips dan trik dalam penerapannya.

## Teorema Limit Pusat (CLT)

Teorema Limit Pusat menyatakan bahwa jika kita mengambil sejumlah besar sampel acak dengan ukuran yang sama ( $n$ ) dari populasi, distribusi sampel mean akan mendekati distribusi normal dengan mean yang sama dengan mean populasi ( $\mu$ ) dan variansi populasi ( $\sigma^2$ ) dibagi dengan ukuran sampel ( $n$ ), terlepas dari bentuk distribusi populasi asalnya.

Mathematically, the Central Limit Theorem can be expressed as:

Jika  $X_1, X_2, \dots, X_n$  adalah variabel acak independen dan identik yang berasal dari populasi dengan rata-rata  $\mu$  dan variansi  $\sigma^2$ , maka untuk  $n$  yang cukup besar:

$$Z = (\bar{X} - \mu) / (\sigma / \sqrt{n})$$

di mana  $\bar{X}$  adalah rata-rata sampel dan  $Z$  adalah variabel acak baru yang mengikuti distribusi normal standar (mean 0 dan variansi 1).

Contoh Kasus di Dunia Nyata:

Sebagai contoh, bayangkan kita ingin mengetahui rata-rata pengeluaran bulanan pelanggan di sebuah toko. Kita tidak perlu mengumpulkan data dari setiap pelanggan, tetapi kita bisa mengambil beberapa sampel acak berukuran 50 pelanggan setiap bulan dan menghitung rata-rata pengeluaran. CLT menjamin bahwa distribusi rata-rata pengeluaran ini akan mendekati distribusi normal, bahkan jika distribusi pengeluaran populasi secara keseluruhan memiliki bentuk yang tidak normal.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import matplotlib.pyplot as plt

# Buat populasi sintetis dengan distribusi eksponensial
population = np.random.exponential(scale=100, size=10000)

# Ambil banyak sampel dari populasi
sample_means = []
num_samples = 1000
sample_size = 50

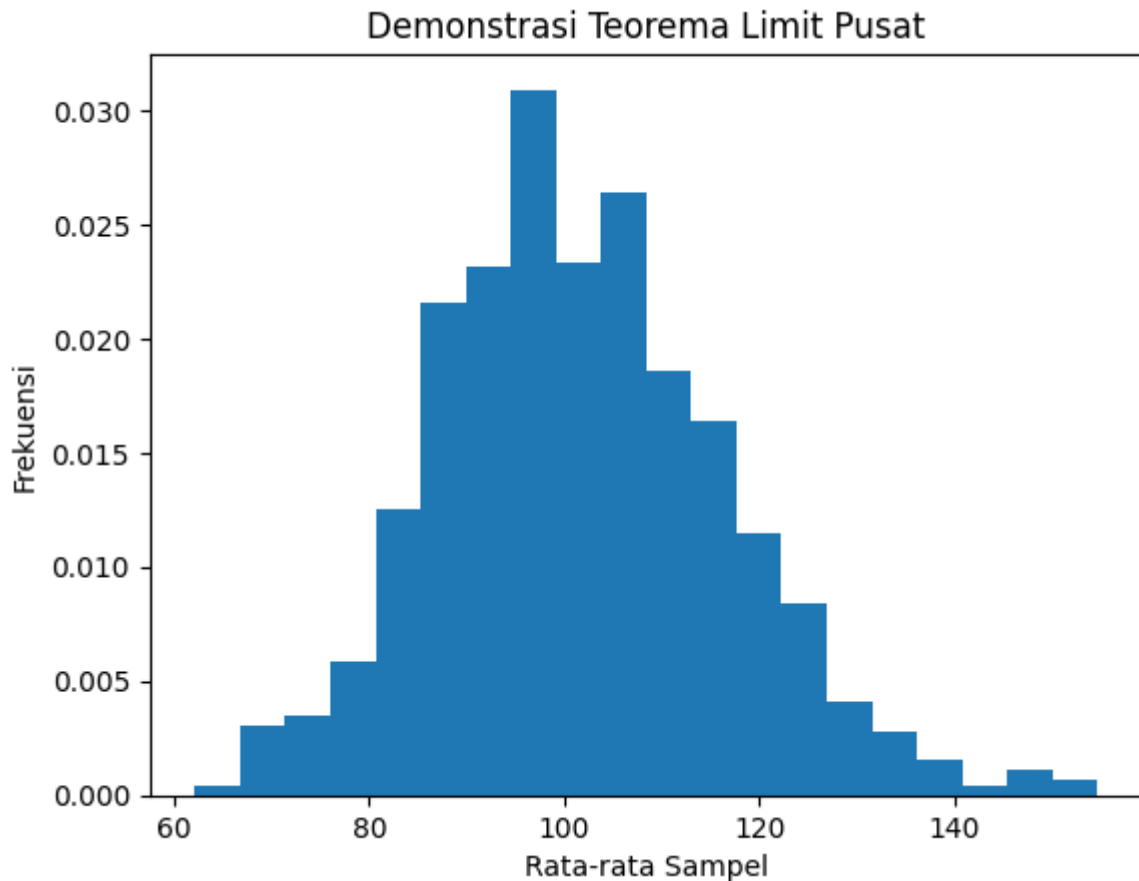
for _ in range(num_samples):
    sample = np.random.choice(population, size=sample_size, replace=False)
    sample_means.append(np.mean(sample))

# Plot distribusi rata-rata sampel
plt.hist(sample_means, bins=20, density=True)
```



```
plt.xlabel('Rata-rata Sampel')
plt.ylabel('Frekuensi')
plt.title('Demonstrasi Teorema Limit Pusat')
plt.show()
```

Output:



### Hukum Bilangan Besar (LLN)

Hukum Bilangan Besar menyatakan bahwa seiring dengan bertambahnya ukuran sampel ( $n$ ), rata-rata sampel ( $\bar{X}$ ) akan semakin mendekati rata-rata populasi ( $\mu$ ). Secara matematis, Hukum Bilangan Besar dapat dinyatakan sebagai:

Untuk setiap  $\varepsilon > 0$ :

$$\lim_{n \rightarrow \infty} P(|\bar{X} - \mu| > \varepsilon) = 0$$

Dengan kata lain, probabilitas bahwa perbedaan antara rata-rata sampel dan rata-rata populasi lebih besar dari  $\varepsilon$  akan menuju ke nol ketika ukuran sampel  $n$  mendekati tak terhingga. Dalam praktiknya, ini berarti bahwa kita dapat meningkatkan akurasi estimasi kita dengan mengumpulkan lebih banyak data.

### Contoh Kasus di Dunia Nyata:

Misalkan kita ingin mengetahui rata-rata waktu tunggu di sebuah restoran cepat saji. Dengan mengumpulkan data waktu tunggu dari semakin banyak pelanggan, kita akan semakin mendekati nilai rata-rata waktu tunggu sebenarnya dari populasi keseluruhan.

### Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import matplotlib.pyplot as plt

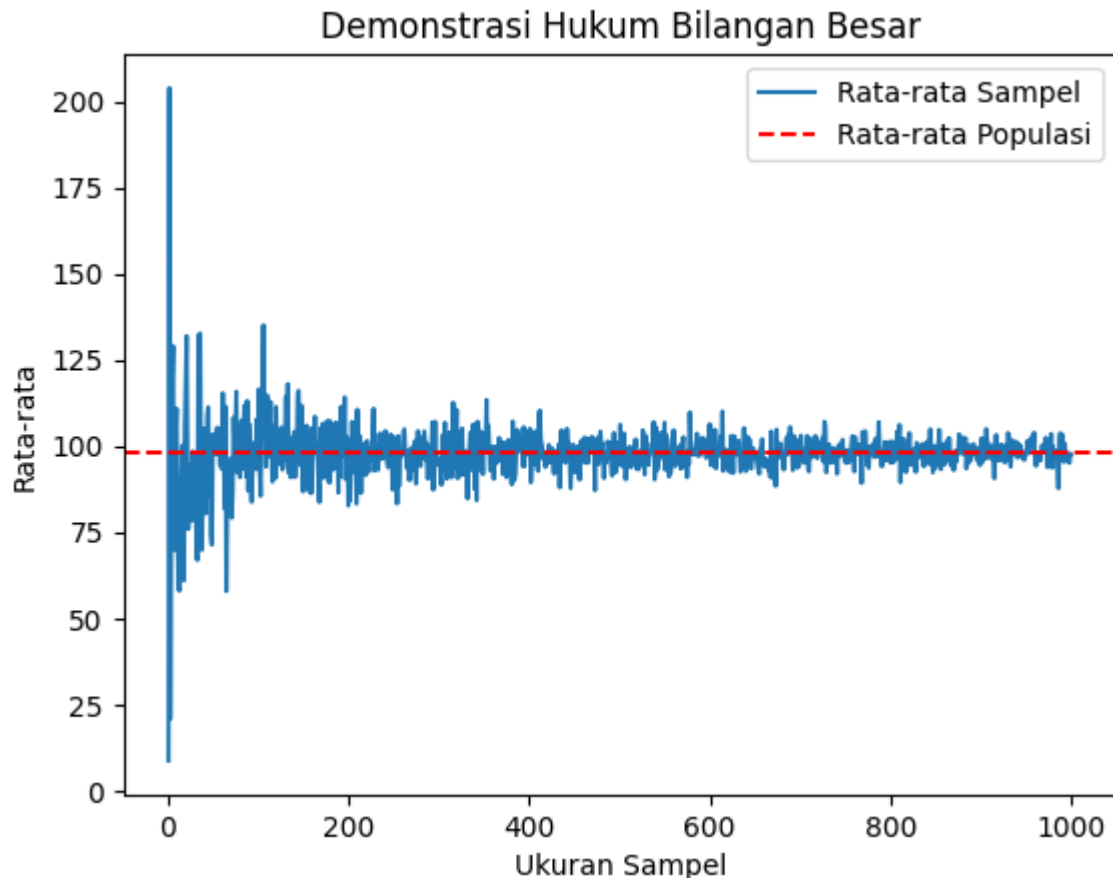
# Buat populasi sintetis dengan distribusi eksponensial
population = np.random.exponential(scale=100, size=10000)
population_mean = np.mean(population)

sample_sizes = np.arange(1, 1001)
sample_means = []

# Menghitung rata-rata sampel untuk berbagai ukuran sampel
for sample_size in sample_sizes:
    sample = np.random.choice(population, size=sample_size, replace=False)
    sample_means.append(np.mean(sample))

# Plot rata-rata sampel dan rata-rata populasi
plt.plot(sample_sizes, sample_means, label='Rata-rata Sampel')
plt.axhline(y=population_mean, color='r', linestyle='--', label='Rata-rata Populasi')
plt.xlabel('Ukuran Sampel')
plt.ylabel('Rata-rata')
plt.title('Demonstrasi Hukum Bilangan Besar')
plt.legend()
plt.show()
```

### Output:



Tips dan Trik dalam Penerapan CLT dan LLN:

Pastikan sampel yang diambil adalah acak dan independen. Jika sampel tidak acak atau ada ketergantungan antara elemen sampel, CLT dan LLN mungkin tidak berlaku.

Ketika menggunakan CLT untuk menghitung interval kepercayaan atau melakukan uji hipotesis, pastikan ukuran sampel cukup besar. Sebagai aturan praktis, ukuran sampel  $n > 30$  dianggap cukup besar untuk CLT berlaku pada populasi yang tidak terlalu mencolok dari distribusi normal.

Ketika menggunakan LLN untuk estimasi, perlu diingat bahwa meskipun rata-rata sampel akan semakin mendekati rata-rata populasi dengan peningkatan ukuran sampel, peningkatan ini mungkin mengalami penurunan hasil yang semakin berkurang. Oleh karena itu, pertimbangkan trade-off antara akurasi estimasi dan biaya pengumpulan data tambahan.

Jika populasi asal memiliki distribusi yang sangat mencolok atau memiliki pencilan ekstrem, pertimbangkan untuk menggunakan metode nonparametrik atau melakukan transformasi data sebelum menerapkan CLT atau LLN.

Dengan pemahaman yang baik tentang Teorema Limit Pusat dan Hukum Bilangan Besar, kita dapat membuat estimasi yang lebih akurat dan generalisasi yang lebih kuat tentang populasi berdasarkan sampel yang diambil. Penggunaan kedua konsep ini dalam berbagai situasi memungkinkan kita untuk melakukan analisis statistik yang efisien dan efektif dalam banyak bidang, mulai dari penelitian akademik hingga pengambilan keputusan bisnis. Selalu penting untuk mempertimbangkan kondisi populasi dan ukuran sampel yang digunakan, serta keterbatasan metode yang digunakan saat menerapkan CLT dan LLN. Dengan mempraktikkan tips dan trik yang disajikan dalam subbab ini, kita akan lebih siap untuk menghadapi tantangan yang muncul dalam analisis statistik inferensial.

### 3.3. Estimasi Parameter: Titik dan Interval

Estimasi parameter adalah proses menghitung nilai yang diharapkan dari suatu parameter populasi berdasarkan sampel yang diambil dari populasi tersebut. Dalam statistik inferensial, ada dua jenis estimasi parameter yang umum digunakan: estimasi titik dan estimasi interval. Subbab ini akan menjelaskan secara rinci tentang kedua jenis estimasi ini, termasuk penjelasan matematis, contoh kasus di dunia nyata, contoh kode Python dengan data sintetis, serta tips dan trik yang relevan.

#### A. Estimasi Titik

Estimasi titik adalah pendekatan di mana satu nilai tunggal dihitung sebagai perkiraan suatu parameter populasi. Nilai ini biasanya ditemukan dengan menggunakan statistik yang dihitung dari sampel. Berikut adalah penjelasan matematis untuk beberapa metode estimasi titik yang umum:

Mean Sampel ( $\bar{x}$ ) - Estimasi titik untuk mean populasi ( $\mu$ ) dapat dihitung dengan rata-rata aritmatika dari semua pengamatan dalam sampel:

$$\bar{x} = \sum x_i / n$$

di mana  $x_i$  adalah nilai pengamatan ke- $i$  dan  $n$  adalah jumlah pengamatan dalam sampel.

Proporsi Sampel ( $\hat{p}$ ) - Estimasi titik untuk proporsi populasi ( $p$ ) dapat dihitung dengan membagi jumlah keberhasilan dalam sampel ( $k$ ) dengan jumlah pengamatan dalam sampel ( $n$ ):

$$\hat{p} = k / n$$

Contoh Kasus Dunia Nyata: Sebuah perusahaan ingin mengetahui rata-rata kepuasan pelanggan terhadap produk mereka. Mereka mengumpulkan data dari 100 pelanggan secara acak dan menghitung rata-rata kepuasan pelanggan (dalam skala 1-10). Rata-rata kepuasan pelanggan dari sampel ini adalah 7,5. Dalam kasus ini, 7,5 adalah estimasi titik untuk mean populasi kepuasan pelanggan.

## B. Estimasi Interval

Estimasi interval, atau interval kepercayaan, adalah pendekatan di mana rentang nilai dihitung untuk menggambarkan seberapa yakin kita dapat memperkirakan parameter populasi. Interval ini biasanya dihitung dengan menggunakan estimasi titik dan beberapa ukuran ketidakpastian, seperti standar deviasi atau kesalahan standar.

Interval Kepercayaan untuk Mean Populasi ( $\mu$ ) - Interval kepercayaan dapat dihitung dengan menggunakan mean sampel ( $\bar{x}$ ), kesalahan standar dari mean ( $SEM$ ), dan nilai kritis yang sesuai dengan tingkat kepercayaan yang diinginkan ( $z$ ):

$$\bar{x} \pm z \times SEM$$

di mana  $SEM = \sigma / \sqrt{n}$  dan  $\sigma$  adalah standar deviasi populasi.

Jika standar deviasi populasi tidak diketahui, kita dapat menggunakan standar deviasi sampel ( $s$ ) sebagai pengganti:

$$\bar{x} \pm t \times SEM$$

di mana  $SEM = s / \sqrt{n}$  dan  $t$  adalah nilai kritis yang sesuai dengan tingkat kepercayaan yang diinginkan dan derajat kebebasan ( $df$ ).

Interval Kepercayaan untuk Proporsi Populasi ( $p$ ) - Interval kepercayaan dapat dihitung dengan menggunakan proporsi sampel ( $\hat{p}$ ), kesalahan standar dari proporsi ( $SEP$ ), dan nilai kritis yang sesuai dengan tingkat kepercayaan yang diinginkan ( $z$ ):

$$\hat{p} \pm z \times SEP$$

di mana  $SEP = \sqrt{\hat{p}(1 - \hat{p}) / n}$

Contoh Kasus Dunia Nyata: Menggunakan data kepuasan pelanggan dari contoh sebelumnya, perusahaan ingin menghitung interval kepercayaan 95% untuk rata-rata kepuasan pelanggan di seluruh populasi. Diketahui standar deviasi sampel adalah 1,8. Mereka menghitung interval kepercayaan sebagai berikut:

Menghitung kesalahan standar dari mean ( $SEM$ ):

$$SEM = s / \sqrt{n} = 1,8 / \sqrt{100} = 0,18$$

Mencari nilai kritis  $t$  untuk tingkat kepercayaan 95% dan derajat kebebasan ( $df = n - 1 = 99$ ). Dalam kasus ini,  $t \approx 1,984$ .

Menghitung interval kepercayaan:

$$\bar{x} \pm t \times SEM = 7,5 \pm 1,984 \times 0,18 \approx (7,14, 7,86)$$

Perusahaan dapat menyatakan dengan keyakinan 95% bahwa rata-rata kepuasan pelanggan di seluruh populasi berada dalam kisaran 7,14 hingga 7,86.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import scipy.stats as stats

# Data sintetis kepuasan pelanggan
sample_data = np.random.normal(loc=7.5, scale=1.8, size=100)

# Menghitung mean sampel
sample_mean = np.mean(sample_data)

# Menghitung standar deviasi sampel
sample_std = np.std(sample_data, ddof=1)

# Menghitung kesalahan standar
sem = sample_std / np.sqrt(len(sample_data))

# Menghitung interval kepercayaan 95%
confidence_level = 0.95
degrees_of_freedom = len(sample_data) - 1
t_critical = stats.t.ppf((1 + confidence_level) / 2, degrees_of_freedom)
confidence_interval = (sample_mean - t_critical * sem, sample_mean + t_critical * sem)

print("Interval kepercayaan 95%:", confidence_interval)
```

Output:

```
Interval kepercayaan 95%: (6.845602212328415, 7.615579348100288)
```

Tips dan Trik dalam Penerapan Estimasi Parameter:

Selalu pastikan bahwa sampel yang digunakan dalam analisis statistik diambil secara acak dan representatif terhadap populasi yang ingin kamu teliti. Hal ini akan membantu mengurangi bias dan meningkatkan validitas hasil.

Jika sampel kamu berukuran besar, distribusi sampel akan mendekati distribusi normal berdasarkan Teorema Limit Pusat. Dalam kasus ini, kamu dapat menggunakan estimasi interval dengan nilai  $z$  daripada  $t$ .

Perhatikan tingkat kepercayaan yang kamu gunakan untuk menghitung interval kepercayaan. Tingkat kepercayaan yang lebih tinggi akan menghasilkan interval yang lebih lebar, yang mencerminkan ketidakpastian yang lebih besar dalam estimasi. Sebaliknya, tingkat kepercayaan yang lebih rendah akan menghasilkan interval yang lebih sempit, tetapi kemungkinan interval tersebut tidak mencakup nilai sebenarnya dari parameter populasi akan lebih besar.

Jika kamu tidak memiliki informasi tentang standar deviasi populasi, kamu dapat menggunakan standar deviasi sampel sebagai pengganti. Namun, ini akan menghasilkan estimasi yang lebih tidak akurat, terutama jika ukuran sampel kamu kecil.

Estimasi interval sering kali lebih informatif daripada estimasi titik karena memberikan informasi tentang ketidakpastian yang terkait dengan perkiraan kamu. Oleh karena itu, sebaiknya gunakan interval kepercayaan ketika melaporkan hasil analisis statistik.

Pertimbangkan untuk menggunakan metode nonparametrik atau metode berbasis pem-bootstrap jika kamu tidak yakin tentang asumsi yang mendasari distribusi data kamu. Metode-metode ini tidak bergantung pada asumsi distribusi tertentu dan sering kali lebih kuat terhadap pelanggaran asumsi.

Jangan ragu untuk memeriksa beberapa metode estimasi parameter yang berbeda dan bandingkan hasilnya. Ini akan memberi kamu pemahaman yang lebih baik tentang sejauh mana hasil kamu sensitif terhadap pilihan metode yang digunakan.

### **3.4. Uji Hipotesis: Uji Z, Uji T, Uji Chi-Square, dan Uji F**

Uji hipotesis adalah salah satu konsep kunci dalam statistik inferensial dan digunakan untuk menguji klaim atau asumsi tentang populasi berdasarkan sampel data. Uji hipotesis melibatkan perbandingan antara statistik yang dihitung dari data sampel dengan distribusi yang diharapkan berdasarkan hipotesis nol ( $H_0$ ). Dalam subbab ini, kita akan membahas beberapa uji hipotesis yang umum digunakan, yaitu uji Z, uji T, uji Chi-Square, dan uji F.

#### **3.4.1. Uji Z**

Uji Z adalah salah satu uji hipotesis yang paling umum digunakan. Uji ini digunakan ketika populasi dianggap memiliki distribusi normal dan standar deviasi populasi diketahui. Dalam uji Z, kita menghitung nilai Z-statistik sebagai berikut:

$$Z = (\bar{X} - \mu) / (\sigma / \sqrt{n})$$

di mana  $\bar{X}$  adalah rata-rata sampel,  $\mu$  adalah rata-rata populasi,  $\sigma$  adalah standar deviasi populasi, dan  $n$  adalah ukuran sampel. Nilai Z yang dihitung kemudian dibandingkan

dengan nilai kritis yang sesuai dengan tingkat signifikansi yang dipilih ( $\alpha$ ) untuk menentukan apakah kita harus menerima atau menolak hipotesis nol.

Contoh Kasus: Sebuah perusahaan mengklaim bahwa rata-rata waktu penyelesaian tugas karyawan adalah 30 menit dengan standar deviasi 5 menit. Untuk menguji klaim ini, kamu mengambil sampel acak 50 karyawan dan menemukan bahwa rata-rata waktu penyelesaian tugas mereka adalah 28 menit. Dapatkah kita menolak klaim perusahaan dengan tingkat signifikansi 5%?

Solusi:

Tentukan hipotesis nol ( $H_0$ ) dan hipotesis alternatif ( $H_1$ ):

$H_0: \mu = 30$

$H_1: \mu \neq 30$

Hitung Z-statistik:

$$Z = (28 - 30) / (5 / \sqrt{50}) \approx -2.83$$

Tentukan nilai kritis dengan tingkat signifikansi 5% ( $\alpha = 0.05$ ) menggunakan tabel Z atau fungsi terkait dalam perangkat lunak statistik. Dalam kasus ini, nilai kritis adalah  $\pm 1.96$ .

Karena nilai Z yang dihitung (-2.83) berada di luar rentang nilai kritis (-1.96 hingga 1.96), kita menolak hipotesis nol dan menyimpulkan bahwa rata-rata waktu penyelesaian tugas karyawan berbeda dari 30 menit.

Contoh Python script dengan data sintetis:

```
import numpy as np
import scipy.stats as stats

# Data
sampel = np.random.normal(28, 5, 50)
rata_rata_populasi = 30
standar_deviasi_populasi = 5
ukuran_sampel = len(sampel)
tingkat_signifikansi = 0.05

# Menghitung Z-statistik
rata_rata_sampel = np.mean(sampel)
z_stat = (rata_rata_sampel - rata_rata_populasi) / (standar_deviasi_populasi / np.sqrt(ukuran_sampel))

# Menghitung nilai kritis
nilai_kritis_bawah, nilai_kritis_atas = stats.norm.ppf([tingkat_signifikansi/2, 1 - tingkat_signifikansi/2])
```



```
tingkat_signifikansi/2))

# Mengevaluasi hasil
if nilai_kritis_bawah <= z_stat <= nilai_kritis_atas:
    print("Tidak menolak H0")
else:
    print("Menolak H0")
```

Output:

**Menolak H0**

Tips dan trik:

Pastikan asumsi distribusi normal dan standar deviasi populasi yang diketahui sebelum menggunakan uji Z.

Jika standar deviasi populasi tidak diketahui, gunakan uji T.

### 3.4.2. Uji T

Uji T, atau uji T-Student, digunakan saat standar deviasi populasi tidak diketahui. Uji ini serupa dengan uji Z, tetapi menggunakan distribusi T-Student yang mengakomodasi ketidakpastian yang lebih besar karena standar deviasi populasi tidak diketahui.

Untuk menghitung T-statistik, gunakan rumus berikut:

$$T = (\bar{X} - \mu) / (s / \sqrt{n})$$

di mana s adalah standar deviasi sampel. Seperti uji Z, kita kemudian membandingkan T-statistik yang dihitung dengan nilai kritis yang sesuai dengan tingkat signifikansi yang dipilih ( $\alpha$ ) untuk menentukan apakah kita harus menerima atau menolak hipotesis nol.

Contoh Kasus: Sebuah pabrik mengklaim bahwa rata-rata berat produk yang diproduksi adalah 500 gram. Kamu mengambil sampel acak 10 produk dan menemukan berat rata-rata 505 gram dengan standar deviasi sampel 8 gram. Dapatkah kita menolak klaim pabrik dengan tingkat signifikansi 5%?

Solusi:

Tentukan hipotesis nol ( $H_0$ ) dan hipotesis alternatif ( $H_1$ ):

$H_0: \mu = 500$

$H_1: \mu \neq 500$

Hitung T-statistik:

$$T = (505 - 500) / (8 / \sqrt{10}) \approx 1.77$$

Tentukan nilai kritis dengan tingkat signifikansi 5% ( $\alpha = 0.05$ ) menggunakan tabel T atau fungsi terkait dalam perangkat lunak statistik. Karena ukuran sampel kita adalah 10, kita akan menggunakan derajat kebebasan (df) 9. Dalam kasus ini, nilai kritis adalah  $\pm 2.262$ .

Karena nilai T yang dihitung (1.77) berada di dalam rentang nilai kritis (-2.262 hingga 2.262), kita tidak menolak hipotesis nol dan menyimpulkan bahwa tidak ada bukti yang cukup untuk mengatakan bahwa rata-rata berat produk berbeda dari 500 gram.

Contoh Python script dengan data sintetis:

```
sampel = np.array([510, 495, 512, 498, 504, 508, 497, 502, 503, 500])
rata_rata_populasi = 500
ukuran_sampel = len(sampel)
tingkat_signifikansi = 0.05

# Menghitung T-statistik
rata_rata_sampel = np.mean(sampel)
standar_deviasi_sampel = np.std(sampel, ddof=1)
t_stat = (rata_rata_sampel - rata_rata_populasi) / (standar_deviasi_sampel /
np.sqrt(ukuran_sampel))

# Menghitung nilai kritis
nilai_kritis_bawah, nilai_kritis_atas = stats.t.ppf([tingkat_signifikansi/2, 1 -
tingkat_signifikansi/2], df=ukuran_sampel-1)

# Mengevaluasi hasil
if nilai_kritis_bawah <= t_stat <= nilai_kritis_atas:
    print("Tidak menolak H0")
else:
    print("Menolak H0")
```

Output:

Tidak menolak H0

Tips dan trik:

Gunakan uji T ketika standar deviasi populasi tidak diketahui.

Jika ukuran sampel sangat besar, distribusi T akan mendekati distribusi normal dan uji T akan memberikan hasil yang serupa dengan uji Z.

Pastikan distribusi populasi mendekati normal, terutama jika ukuran sampel kecil.

### 3.4.3. Uji Chi-Square

Uji Chi-Square digunakan untuk menguji hubungan antara dua variabel kategorikal. Uji ini membandingkan perbedaan antara frekuensi observasi dan frekuensi yang diharapkan berdasarkan hipotesis nol yang menyatakan bahwa tidak ada hubungan antara variabel.

Untuk menghitung Chi-Square-statistik, gunakan rumus berikut:

$$\chi^2 = \sum [(O - E)^2 / E]$$

di mana O adalah frekuensi observasi dan E adalah frekuensi yang diharapkan. Kemudian,  $\chi^2$ -statistik dibandingkan dengan nilai kritis yang sesuai dengan tingkat signifikansi yang dipilih ( $\alpha$ ) untuk menentukan apakah kita harus menerima atau menolak hipotesis nol.

Contoh Kasus: Sebuah penelitian ingin mengetahui apakah ada hubungan antara jenis kelamin dan preferensi warna. Berikut adalah tabel frekuensi observasi:

	Biru	Merah	Hijau
Laki-laki	50	45	55
Perempuan	60	40	50

Dapatkah kita menolak hipotesis nol bahwa tidak ada hubungan antara jenis kelamin dan preferensi warna dengan tingkat signifikansi 5%?

Solusi:

Tentukan hipotesis nol ( $H_0$ ) dan hipotesis alternatif ( $H_1$ ):

$H_0$ : Tidak ada hubungan antara jenis kelamin dan preferensi warna.

$H_1$ : Ada hubungan antara jenis kelamin dan preferensi warna.

Hitung frekuensi yang diharapkan berdasarkan hipotesis nol:

	Biru	Merah	Hijau
Laki-laki	55	42.5	52.5
Perempuan	55	42.5	52.5

3. Hitung Chi-Square-statistik:

$$\chi^2 = \sum [(O - E)^2 / E] = (50-55)^2/55 + (45-42.5)^2/42.5 + (55-52.5)^2/52.5 + (60-55)^2/55 + (40-42.5)^2/42.5 + (50-52.5)^2/52.5 \approx 1.44$$

4. Tentukan nilai kritis dengan tingkat signifikansi 5% ( $\alpha = 0.05$ ) menggunakan tabel Chi-Square atau fungsi terkait dalam perangkat lunak statistik. Karena kita memiliki 2 baris dan 3 kolom, kita akan menggunakan derajat kebebasan (df)  $(2-1)*(3-1) = 2$ . Dalam kasus ini, nilai kritis adalah 5.99.

5. Karena nilai  $\chi^2$  yang dihitung (2.73) lebih kecil dari nilai kritis (5.99), kita tidak menolak hipotesis nol dan menyimpulkan bahwa tidak ada bukti yang cukup untuk mengatakan ada hubungan antara jenis kelamin dan preferensi warna.

Contoh Python script dengan data sintetis:

```
import pandas as pd

# Data
data = pd.DataFrame({'Biru': [50, 60], 'Merah': [45, 40], 'Hijau': [55, 50]}, index=['Laki-laki', 'Perempuan'])

# Menghitung frekuensi yang diharapkan
total_baris = data.sum(axis=1)
total_kolom = data.sum(axis=0)
total = data.values.sum()
frekuensi_diharapkan = np.outer(total_baris, total_kolom) / total

# Menghitung Chi-Square-statistik
chi_square_stat = ((data.values - frekuensi_diharapkan)**2 / frekuensi_diharapkan).sum()

# Menghitung nilai kritis
nilai_kritis = stats.chi2.ppf(1 - tingkat_signifikansi, df=(data.shape[0]-1)*(data.shape[1]-1))

# Mengevaluasi hasil
if chi_square_stat <= nilai_kritis:
    print("Tidak menolak H0")
else:
    print("Menolak H0")
```

Output:

Tidak menolak H0

Tips dan trik:

Gunakan uji Chi-Square untuk menguji hubungan antara dua variabel kategorikal. Pastikan jumlah observasi cukup besar agar uji Chi-Square valid.

### 3.4.4. Uji F

Uji F digunakan untuk menguji perbedaan antara dua atau lebih kelompok dalam analisis varians (ANOVA). Uji ini membandingkan variasi antar kelompok dengan variasi dalam kelompok untuk menentukan apakah perbedaan rata-rata antar kelompok signifikan secara statistik.

Untuk menghitung F-statistik, gunakan rumus berikut:

$$F = (MS_{\text{between}} / k - 1) / (MS_{\text{within}} / (N - k))$$

di mana  $MS_{\text{between}}$  adalah variasi antar kelompok,  $MS_{\text{within}}$  adalah variasi dalam kelompok,  $k$  adalah jumlah kelompok, dan  $N$  adalah jumlah total observasi. Kemudian, F-statistik dibandingkan dengan nilai kritis yang sesuai dengan tingkat signifikansi yang dipilih ( $\alpha$ ) untuk menentukan apakah kita harus menerima atau menolak hipotesis nol.

Contoh Kasus: Sebuah penelitian ingin mengetahui apakah ada perbedaan signifikan dalam hasil tes siswa yang diajarkan oleh tiga guru yang berbeda. Berikut adalah hasil tes dari siswa yang diajarkan oleh masing-masing guru:

Guru A: [85, 90, 88, 84, 89]

Guru B: [78, 80, 82, 79, 81]

Guru C: [92, 93, 94, 91, 95]

Dapatkah kita menolak hipotesis nol bahwa tidak ada perbedaan signifikan dalam hasil tes siswa yang diajarkan oleh tiga guru dengan tingkat signifikansi 5%?

Solusi:

Tentukan hipotesis nol ( $H_0$ ) dan hipotesis alternatif ( $H_1$ ):

$H_0$ : Tidak ada perbedaan signifikan dalam hasil tes siswa yang diajarkan oleh tiga guru.

$H_1$ : Ada perbedaan signifikan dalam hasil tes siswa yang diajarkan oleh tiga guru.

Hitung F-statistik menggunakan rumus di atas atau fungsi terkait dalam perangkat lunak statistik. Dalam kasus ini, F-statistik  $\approx 90.15$ .

Tentukan nilai kritis dengan tingkat signifikansi 5% ( $\alpha = 0.05$ ) menggunakan tabel F atau fungsi terkait dalam perangkat lunak statistik. Karena kita memiliki 3 kelompok dan 15 observasi, kita akan menggunakan derajat kebebasan (df) 2 dan 12. Dalam kasus ini, nilai kritis adalah 3.89.

Karena nilai F yang dihitung (90.15) lebih besar dari nilai kritis (3.89), kita menolak hipotesis nol dan menyimpulkan bahwa ada perbedaan signifikan dalam hasil tes siswa yang diajarkan oleh tiga guru.

Contoh Python script dengan data sintetis:

```
from scipy.stats import f_oneway
```

```
# Data
```

```

guru_A = [85, 90, 88, 84, 89]
guru_B = [78, 80, 82, 79, 81]
guru_C = [92, 93, 94, 91, 95]

# Menghitung F-statistik dan p-value
f_stat, p_value = f_oneway(guru_A, guru_B, guru_C)

# Mengevaluasi hasil
if p_value <= tingkat_signifikansi:
    print("Menolak H0")
else:
    print("Tidak menolak H0")

```

Output:

**Menolak H0**

Tips dan trik:

Gunakan uji F untuk menguji perbedaan antara dua atau lebih kelompok dalam analisis varians (ANOVA).

Pastikan asumsi ANOVA terpenuhi, termasuk normalitas, homoskedastisitas, dan independensi observasi.

Pertimbangkan untuk menggunakan teknik post-hoc, seperti Tukey HSD, untuk mengidentifikasi perbedaan spesifik antara pasangan kelompok jika uji F menunjukkan perbedaan signifikan secara keseluruhan.

Dalam subbab 3.4, kita telah membahas berbagai metode uji hipotesis, termasuk uji Z, uji T, uji Chi-Square, dan uji F. Setiap metode memiliki kasus penggunaan, asumsi, dan formula yang berbeda untuk menghitung statistik uji dan mengevaluasi hipotesis nol. Dalam prakteknya, penting untuk memahami kapan dan bagaimana menggunakan setiap metode, serta bagaimana menginterpretasikan hasilnya.

Secara umum, uji hipotesis adalah proses penting dalam statistik inferensial yang memungkinkan kita untuk membuat kesimpulan tentang populasi berdasarkan sampel data. Menggunakan metode yang tepat dan memastikan asumsi yang sesuai terpenuhi sangat penting untuk mendapatkan hasil yang valid dan dapat diandalkan.

Beberapa tips dan trik yang dapat membantu dalam penerapan uji hipotesis dalam berbagai situasi adalah:

Pahami tujuan penelitian kamu dan pilih metode uji hipotesis yang sesuai dengan tujuan tersebut.

Periksa asumsi yang mendasari metode uji hipotesis yang kamu pilih dan pastikan data kamu memenuhi asumsi tersebut. Jika tidak, pertimbangkan untuk menggunakan metode alternatif atau melakukan transformasi pada data kamu.

Jangan lupa untuk mempertimbangkan ukuran efek dan kekuatan statistik, serta tingkat signifikansi ( $\alpha$ ) ketika membuat keputusan tentang hipotesis nol.

Ingatlah bahwa uji hipotesis hanya memberikan bukti yang mendukung atau menentang hipotesis nol, dan tidak pernah dapat membuktikan secara pasti bahwa hipotesis nol benar atau salah.

Selalu lakukan analisis eksplorasi data dan visualisasi sebelum melakukan uji hipotesis untuk mendapatkan pemahaman yang lebih baik tentang pola dan hubungan dalam data kamu.

Dengan pemahaman yang kuat tentang konsep yang dibahas dalam subbab 3.4, kamu akan dapat menerapkan uji hipotesis dengan percaya diri dan mendapatkan wawasan yang berharga dari data kamu.

# Bab 4: Analisis Regresi dan Korelasi

## 4.1. Korelasi: Pearson, Spearman, dan Kendall

Korelasi adalah teknik statistik yang digunakan untuk mengukur dan menjelaskan hubungan antara dua variabel. Dalam analisis data, sangat penting untuk mengidentifikasi hubungan antara variabel yang akan membantu kita memahami pola dan tren dalam data. Ada tiga metode umum untuk menghitung koefisien korelasi: Pearson, Spearman, dan Kendall. Masing-masing metode ini memiliki kelebihan dan kekurangan tersendiri, serta asumsi yang harus dipenuhi.

### 4.1.1. Korelasi Pearson

Korelasi Pearson, atau korelasi produk-momen, merupakan metode yang paling umum digunakan untuk mengukur korelasi linier antara dua variabel kontinu. Koefisien korelasi Pearson, yang dinyatakan sebagai 'r', berkisar antara -1 dan 1. Nilai positif menunjukkan hubungan positif antara variabel (seiring peningkatan satu variabel, variabel lain juga meningkat), sedangkan nilai negatif menunjukkan hubungan negatif (seiring peningkatan satu variabel, variabel lain menurun). Nilai 0 menunjukkan tidak ada korelasi antara kedua variabel.

Formula matematis untuk menghitung korelasi Pearson adalah sebagai berikut:

$$r = \Sigma((xi - \mu x)(yi - \mu y)) / \sqrt{(\Sigma(xi - \mu x)^2 * \Sigma(yi - \mu y)^2)}$$

di mana:

xi dan yi adalah nilai variabel X dan Y

$\mu x$  dan  $\mu y$  adalah rata-rata X dan Y

Asumsi yang harus dipenuhi untuk menggunakan korelasi Pearson adalah:

Skala pengukuran kedua variabel adalah interval atau rasio

Kedua variabel memiliki distribusi normal

Hubungan antara variabel bersifat linier

Contoh kasus di dunia nyata: Korelasi antara tinggi badan dan berat badan pada sekelompok individu.

### 4.1.2. Korelasi Spearman

Korelasi Spearman, atau korelasi peringkat, digunakan untuk mengukur hubungan monoton antara dua variabel. Ini adalah metode nonparametrik yang tidak mengasumsikan distribusi normal dari data. Koefisien korelasi Spearman, yang dinyatakan sebagai 'ρ' atau 'rs', juga berkisar antara -1 dan 1.



Formula matematis untuk menghitung korelasi Spearman adalah sebagai berikut:

$$\rho = 1 - (6 * \Sigma d^2) / (n^3 - n)$$

di mana:

d adalah perbedaan antara peringkat X dan Y

n adalah jumlah observasi

Asumsi yang harus dipenuhi untuk menggunakan korelasi Spearman adalah:

Skala pengukuran kedua variabel adalah ordinal, interval, atau rasio

Hubungan antara variabel bersifat monoton

Contoh kasus di dunia nyata: Korelasi antara peringkat restoran dan jumlah ulasan yang diterima.

#### 4.1.3. Korelasi Kendall

Korelasi Kendall, atau korelasi  $\tau$  (tau Kendall), digunakan untuk mengukur hubungan ordinal antara dua variabel. Seperti korelasi Spearman, ini adalah metode nonparametrik yang tidak mengasumsikan distribusi normal dari data. Koefisien korelasi Kendall, yang dinyatakan sebagai ' $\tau$ ', juga berkisar antara -1 dan 1.

Formula matematis untuk menghitung korelasi Kendall adalah sebagai berikut:

$$\tau = (n_c - n_d) / \sqrt{((n_c + n_d + n_t) * (n_c + n_d + n_u))}$$

di mana:

$n_c$  adalah jumlah pasangan konkordansi (kedua variabel memiliki urutan yang sama)

$n_d$  adalah jumlah pasangan diskordansi (kedua variabel memiliki urutan yang berlawanan)

$n_t$  dan  $n_u$  adalah jumlah pasangan yang terikat pada masing-masing variabel

Asumsi yang harus dipenuhi untuk menggunakan korelasi Kendall adalah:

Skala pengukuran kedua variabel adalah ordinal, interval, atau rasio

Hubungan antara variabel bersifat ordinal

Contoh kasus di dunia nyata: Korelasi antara peringkat kesulitan pada kursus online dan persentase peserta yang menyelesaikan kursus.

Berikut adalah contoh Python script untuk menghitung korelasi Pearson, Spearman, dan Kendall menggunakan data sintesis:

```
import numpy as np
import pandas as pd
from scipy.stats import pearsonr, spearmanr, kendalltau
```

```

# Membuat data sintetis
np.random.seed(42)
data = np.random.rand(100, 2)
df = pd.DataFrame(data, columns=['X', 'Y'])

# Menghitung korelasi Pearson
pearson_corr, _ = pearsonr(df['X'], df['Y'])
print(f"Korelasi Pearson: {pearson_corr:.4f}")

# Menghitung korelasi Spearman
spearman_corr, _ = spearmanr(df['X'], df['Y'])
print(f"Korelasi Spearman: {spearman_corr:.4f}")

# Menghitung korelasi Kendall
kendall_corr, _ = kendalltau(df['X'], df['Y'])
print(f"Korelasi Kendall: {kendall_corr:.4f}")

```

Output:

```

Korelasi Pearson: -0.0378
Korelasi Spearman: -0.0397
Korelasi Kendall: -0.0275

```

Tips dan trik dalam penerapannya:

Pilih metode korelasi yang sesuai dengan skala pengukuran dan sifat hubungan antara variabel.

Selalu periksa asumsi yang mendasari metode korelasi sebelum menggunakannya. Jika asumsi tidak dipenuhi, pertimbangkan untuk menggunakan metode alternatif atau mengubah data.

Ingatlah bahwa korelasi bukanlah penyebab. Korelasi yang tinggi antara dua variabel tidak selalu berarti satu variabel menyebabkan perubahan pada variabel lainnya.

Waspada "korelasi palsu" atau "spurious correlation", di mana dua variabel yang tampaknya berkorelasi sebenarnya tidak memiliki hubungan sebab-akibat yang mendasarinya.

Gunakan visualisasi data, seperti scatter plot, untuk membantu mengevaluasi hubungan antara variabel dan memverifikasi hasil korelasi yang dihitung.

## 4.2. Regresi Linier Sederhana

Regresi linier sederhana adalah teknik statistik yang digunakan untuk memodelkan hubungan antara dua variabel dengan asumsi bahwa hubungan tersebut adalah linier. Dalam regresi linier sederhana, kita memiliki satu variabel independen (X) dan satu

variabel dependen (Y). Tujuan dari analisis ini adalah untuk menemukan garis terbaik yang menjelaskan hubungan antara X dan Y, yang disebut garis regresi.

Garis regresi memiliki persamaan:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

di mana:

Y adalah variabel dependen

X adalah variabel independen

$\beta_0$  adalah konstanta (intersep)

$\beta_1$  adalah koefisien regresi (kemiringan)

$\varepsilon$  adalah galat (selisih antara nilai sebenarnya dan nilai yang diprediksi oleh model)

Untuk menemukan nilai  $\beta_0$  dan  $\beta_1$  yang menghasilkan garis regresi terbaik, kita menggunakan metode kuadrat terkecil (least squares method). Dalam metode ini, kita mencari nilai  $\beta_0$  dan  $\beta_1$  yang meminimalkan jumlah kuadrat galat.

Formula untuk menghitung  $\beta_0$  dan  $\beta_1$  adalah:

$$\beta_1 = \Sigma((X_i - \text{mean}(X)) * (Y_i - \text{mean}(Y))) / \Sigma(X_i - \text{mean}(X))^2$$

$$\beta_0 = \text{mean}(Y) - \beta_1 * \text{mean}(X)$$

Contoh kasus di dunia nyata: Memprediksi harga rumah berdasarkan luas tanah.

Berikut adalah contoh Python script untuk melakukan regresi linier sederhana menggunakan data sintetis:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression

# Membuat data sintetis
np.random.seed(42)
X = np.random.rand(100, 1)
Y = 2 + 3 * X + np.random.randn(100, 1)

# Regresi linier sederhana dengan scikit-learn
lin_reg = LinearRegression()
lin_reg.fit(X, Y)

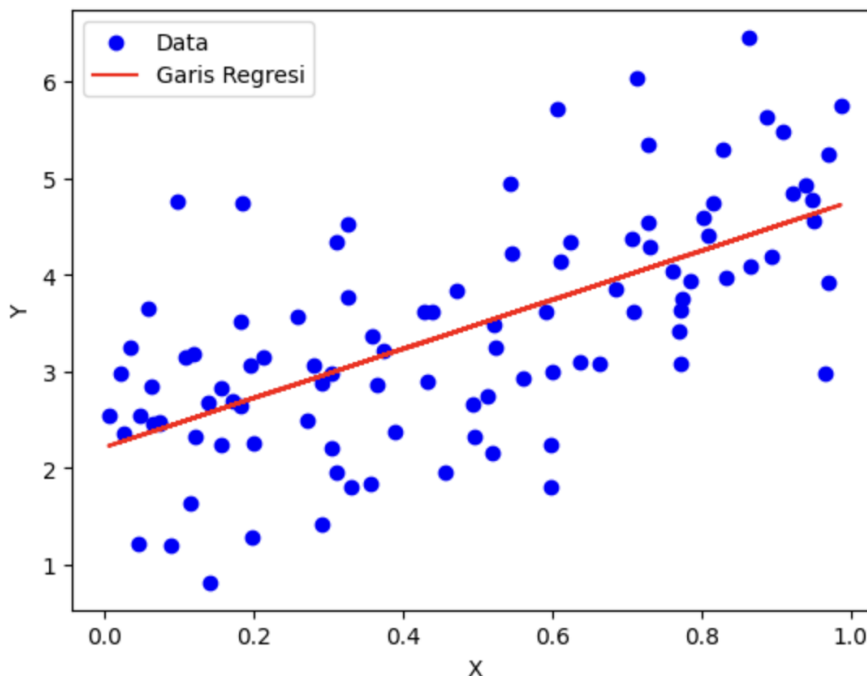
# Koefisien regresi ( $\beta_0$  dan  $\beta_1$ )
```

```
print(f"Intersep ( $\beta_0$ ): {lin_reg.intercept_[0]:.4f}")
print(f"Kemiringan ( $\beta_1$ ): {lin_reg.coef_[0][0]:.4f}")

# Visualisasi data dan garis regresi
plt.scatter(X, Y, color='blue', label='Data')
plt.plot(X, lin_reg.predict(X), color='red', label='Garis Regresi')
plt.xlabel('X')
plt.ylabel('Y')
plt.legend()
plt.show()
```

Output:

```
Intersep ( $\beta_0$ ): 2.2151
Kemiringan ( $\beta_1$ ): 2.5402
```



Tips dan trik dalam penerapannya:

Pastikan hubungan antara variabel adalah linier. Jika hubungan tidak linier, pertimbangkan untuk menggunakan transformasi data atau metode regresi non-linier.

Cek asumsi regresi linier, seperti normalitas galat, independensi galat, dan homoskedastisitas. Jika asumsi tidak dipenuhi, pertimbangkan untuk menggunakan metode regresi alternatif atau mengubah data.

Hati-hati dengan outlier dan titik leverage yang dapat mempengaruhi hasil regresi.

### 4.3. Regresi Linier Berganda

Regresi linier berganda adalah teknik statistik yang digunakan untuk memodelkan hubungan antara satu variabel dependen (Y) dan lebih dari satu variabel independen ( $X_1, X_2, \dots, X_n$ ). Tujuan dari analisis ini adalah untuk menemukan model linier yang menjelaskan hubungan antara variabel independen dan variabel dependen.

Model regresi linier berganda memiliki persamaan:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$

dimana:

Y adalah variabel dependen

$X_1, X_2, \dots, X_n$  adalah variabel independen

$\beta_0$  adalah konstanta (intersep)

$\beta_1, \beta_2, \dots, \beta_n$  adalah koefisien regresi

$\varepsilon$  adalah galat (selisih antara nilai sebenarnya dan nilai yang diprediksi oleh model)

Untuk menemukan nilai  $\beta_0, \beta_1, \dots, \beta_n$  yang menghasilkan model regresi terbaik, kita menggunakan metode kuadrat terkecil (least squares method) yang sama dengan regresi linier sederhana.

Contoh kasus di dunia nyata: Memprediksi harga rumah berdasarkan luas tanah, jumlah kamar tidur, dan usia rumah.

Berikut adalah contoh Python script untuk melakukan regresi linier berganda menggunakan data sintetis:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression

# Membuat data sintetis
np.random.seed(42)
X1 = np.random.rand(100, 1)
X2 = np.random.rand(100, 1)
X3 = np.random.rand(100, 1)
Y = 2 + 3 * X1 + 1.5 * X2 - 2 * X3 + np.random.randn(100, 1)

# Menggabungkan variabel independen dalam satu array
X = np.hstack((X1, X2, X3))

# Regresi linier berganda dengan scikit-learn
```

```

lin_reg = LinearRegression()
lin_reg.fit(X, Y)

# Koefisien regresi ( $\beta_0, \beta_1, \beta_2, \beta_3$ )
print(f"Intersep ( $\beta_0$ ): {lin_reg.intercept_[0]:.4f}")
print(f"Koefisien ( $\beta_1, \beta_2, \beta_3$ ): {lin_reg.coef_[0]}")

# Estimasi harga rumah berdasarkan data baru
X_new = np.array([[0.5, 0.3, 0.7]])
Y_pred = lin_reg.predict(X_new)
print(f"Prediksi harga rumah: {Y_pred[0][0]:.4f}")

```

Output:

```

Intersep ( $\beta_0$ ): 1.6000
Koefisien ( $\beta_1, \beta_2, \beta_3$ ): [ 3.41592926  1.70854972 -1.67065182]
Prediksi harga rumah: 2.6510

```

Tips dan trik dalam penerapannya:

Pilih variabel independen yang relevan dan memiliki hubungan yang signifikan dengan variabel dependen.

Periksa asumsi regresi linier, seperti normalitas galat, independensi galat, dan homoskedastisitas. Jika asumsi tidak dipenuhi, pertimbangkan untuk menggunakan metode regresi alternatif atau mengubah data.

Hati-hati dengan multicollinearity, yaitu hubungan linier

## 4.4. Regresi Logistik

Regresi logistik adalah teknik analisis regresi yang digunakan untuk memodelkan hubungan antara satu variabel dependen kategorik biner (Y) dan satu atau lebih variabel independen ( $X_1, X_2, \dots, X_n$ ). Model regresi logistik menggunakan fungsi logistik (sigmoid) untuk mengestimasi probabilitas  $Y = 1$ , yang merupakan salah satu dari dua kategori yang mungkin.

Model regresi logistik memiliki persamaan:

$$\log(p / (1 - p)) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

di mana:

p adalah probabilitas  $Y = 1$

$X_1, X_2, \dots, X_n$  adalah variabel independen

$\beta_0$  adalah konstanta (intersep)

$\beta_1, \beta_2, \dots, \beta_n$  adalah koefisien regresi

Untuk menemukan nilai  $\beta_0, \beta_1, \dots, \beta_n$  yang menghasilkan model regresi terbaik, kita menggunakan metode maksimum likelihood estimation (MLE).

Contoh kasus di dunia nyata: Memprediksi apakah seseorang akan membeli asuransi berdasarkan usia, jenis kelamin, dan status pernikahan.

Berikut adalah contoh Python script untuk melakukan regresi logistik menggunakan data sintetis:

```
import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression

# Membuat data sintetis
np.random.seed(42)
n_samples = 1000
age = np.random.randint(18, 65, n_samples)
gender = np.random.randint(0, 2, n_samples)
married = np.random.randint(0, 2, n_samples)

# Membuat label berdasarkan data sintetis
Y = (age > 30) & (gender == 1) & (married == 1)
Y = Y.astype(int)

# Menggabungkan variabel independen dalam satu array
X = np.column_stack((age, gender, married))

# Regresi logistik dengan scikit-learn
log_reg = LogisticRegression()
log_reg.fit(X, Y)

# Koefisien regresi ( $\beta_0, \beta_1, \beta_2, \beta_3$ )
print(f"Intersep ( $\beta_0$ ): {log_reg.intercept_[0]:.4f}")
print(f"Koefisien ( $\beta_1, \beta_2, \beta_3$ ): {log_reg.coef_[0]}")

# Estimasi probabilitas seseorang membeli asuransi berdasarkan data baru
X_new = np.array([[35, 1, 1]])
Y_proba = log_reg.predict_proba(X_new)
print(f"Probabilitas membeli asuransi: {Y_proba[0][1]:.4f}")
```

Output:

Intersep ( $\beta_0$ ): -16.9758  
Koefisien ( $\beta_1, \beta_2, \beta_3$ ): [0.13807993 6.56044479 6.08364478]  
Probabilitas membeli asuransi: 0.6227

Tips dan trik dalam penerapannya:

Pastikan variabel dependen kamu adalah kategorik biner.

Pilih variabel independen yang relevan dan memiliki hubungan yang signifikan dengan variabel dependen.

Periksa asumsi regresi logistik, seperti independensi galat dan tidak adanya multicollinearity. Jika asumsi tidak dipenuhi, pertimbangkan untuk menggunakan metode regresi alternatif atau mengubah data.

Gunakan metode validasi silang (cross-validation)

## 4.5. Evaluasi Model: R-Squared, Akurasi, dan Kesalahan

Setelah membangun model regresi, penting untuk mengevaluasi kinerjanya. Berikut adalah beberapa metrik evaluasi yang umum digunakan untuk regresi linier dan regresi logistik.

**R-Squared ( $R^2$ ) atau koefisien determinasi:** Ini adalah metrik yang digunakan untuk mengukur seberapa baik model regresi linier menjelaskan variasi dalam data.  $R^2$  berkisar antara 0 dan 1, di mana nilai yang lebih tinggi menunjukkan bahwa model menjelaskan variasi data dengan lebih baik. Namun,  $R^2$  tidak dapat digunakan untuk mengevaluasi regresi logistik.

**Akurasi:** Ini adalah metrik yang digunakan untuk mengukur seberapa sering model regresi logistik mengklasifikasikan data dengan benar. Akurasi dihitung sebagai jumlah prediksi yang benar dibagi dengan jumlah total prediksi. Namun, akurasi tidak dapat digunakan untuk mengevaluasi regresi linier.

**Kesalahan:** Kesalahan adalah selisih antara nilai yang sebenarnya dan nilai yang diprediksi oleh model. Beberapa metrik kesalahan yang umum digunakan untuk regresi linier adalah:

- Mean Absolute Error (MAE): Rata-rata dari nilai absolut kesalahan.
- Mean Squared Error (MSE): Rata-rata dari kuadrat kesalahan.
- Root Mean Squared Error (RMSE): Akar kuadrat dari rata-rata kuadrat kesalahan.

Untuk regresi logistik, kita menggunakan metrik kesalahan seperti log loss.

Berikut adalah contoh Python script untuk menghitung metrik evaluasi pada model regresi linier dan logistik:

```
from sklearn.metrics import r2_score, accuracy_score, mean_absolute_error, mean_squared_error
```



```

# Regresi linier
Y_lin_actual = np.array([3, -0.5, 2, 7])
Y_lin_predicted = np.array([2.5, 0, 2.1, 7.8])

r2 = r2_score(Y_lin_actual, Y_lin_predicted)
mae = mean_absolute_error(Y_lin_actual, Y_lin_predicted)
mse = mean_squared_error(Y_lin_actual, Y_lin_predicted)
rmse = np.sqrt(mse)

print(f"Regresi Linier: R2 = {r2:.4f}, MAE = {mae:.4f}, MSE = {mse:.4f}, RMSE = {rmse:.4f}")

# Regresi logistik
Y_log_actual = np.array([0, 1, 1, 0, 1, 0])
Y_log_predicted = np.array([0, 1, 0, 0, 1, 1])

accuracy = accuracy_score(Y_log_actual, Y_log_predicted)

print(f"Regresi Logistik: Akurasi = {accuracy:.4f}")

```

Output:

```

Regresi Linier: R2 = 0.9606, MAE = 0.4750, MSE = 0.2875, RMSE = 0.5362
Regresi Logistik: Akurasi = 0.6667

```

Tips dan trik dalam evaluasi model:

Gunakan metrik evaluasi yang sesuai untuk jenis model regresi yang kamu gunakan. Selalu pisahkan data kamu menjadi set pelatihan dan set pengujian untuk menghindari overfitting dan mendapatkan evaluasi yang lebih realistis. Pertimbangkan menggunakan validasi silang (cross-validation) untuk mengurangi variabilitas dalam estimasi kinerja model.

Jangan hanya mengandalkan satu metrik evaluasi, gunakan beberapa metrik untuk mendapatkan pemahaman yang lebih komprehensif tentang kinerja model kamu.

Periksa plot residual (selisih antara nilai yang sebenarnya dan nilai yang diprediksi) untuk mengidentifikasi pola yang tidak terduga atau bias dalam prediksi. Jika ditemukan pola, pertimbangkan untuk menggunakan metode regresi yang berbeda atau mengubah data kamu.

Jika model kamu tidak mencapai kinerja yang diinginkan, pertimbangkan untuk menambahkan atau menghapus variabel independen, mengubah parameter model, atau menggunakan metode regresi alternatif.

Ingatlah bahwa model yang lebih kompleks tidak selalu lebih baik. Terkadang, model yang lebih sederhana dengan sedikit variabel independen dan struktur yang lebih sederhana dapat memberikan kinerja yang serupa atau bahkan lebih baik daripada model yang lebih kompleks. Penting untuk menyeimbangkan kompleksitas model dengan kemampuan untuk menjelaskan variasi dalam data.

Dalam kesimpulannya, analisis regresi dan korelasi merupakan teknik penting dalam statistika yang membantu kita memahami hubungan antara variabel dan membuat prediksi. Dengan menggunakan metode regresi yang sesuai, memilih variabel independen yang relevan, dan mengevaluasi kinerja model dengan metrik yang sesuai, kita dapat membuat model yang andal dan akurat untuk berbagai aplikasi di dunia nyata.

# Bab 5: Analisis Seri Waktu

## 5.1. Komponen Seri Waktu: Tren, Musiman, dan Siklik

Seri waktu adalah rangkaian data yang diurutkan berdasarkan waktu. Analisis seri waktu merupakan metode statistik yang digunakan untuk menganalisis data dalam rangkaian waktu, mengidentifikasi pola dalam data, dan membuat prediksi berdasarkan pola yang ditemukan. Terdapat tiga komponen utama dalam analisis seri waktu, yaitu tren, musiman, dan siklik.

**Tren:** Tren adalah pola jangka panjang yang ada dalam data. Ini bisa naik, turun, atau tetap konstan sepanjang waktu. Tren menggambarkan arah dan kecepatan perubahan data sepanjang waktu.

**Musiman:** Komponen musiman menggambarkan variasi yang terjadi secara berkala dalam data. Variasi ini biasanya dikaitkan dengan perubahan musim, hari dalam seminggu, atau jam dalam sehari. Komponen ini penting karena dapat membantu memprediksi variasi dalam data yang diharapkan terjadi pada waktu yang sama setiap periode.

**Siklik:** Siklik adalah fluktuasi yang terjadi dalam data yang tidak berkaitan dengan musiman. Siklus biasanya lebih panjang daripada musiman dan sulit diprediksi karena tidak teratur. Siklik dapat disebabkan oleh faktor ekonomi, politik, atau sosial yang mempengaruhi data.

Contoh Kasus Dunia Nyata:

Misalkan kita ingin menganalisis data penjualan sepeda di suatu toko sepeda selama beberapa tahun terakhir. Data penjualan ini adalah contoh dari seri waktu. Berikut adalah bagaimana kita dapat mengidentifikasi tren, musiman, dan siklik dalam data ini:

**Tren:** Jika penjualan sepeda meningkat secara konsisten setiap tahun, maka kita akan melihat tren positif dalam data. Sebaliknya, jika penjualan menurun secara konsisten, tren akan negatif.

**Musiman:** Penjualan sepeda mungkin mengalami peningkatan pada musim panas, ketika cuaca lebih baik untuk bersepeda. Oleh karena itu, kita akan melihat komponen musiman dalam data penjualan sepeda.

**Siklik:** Jika terjadi resesi ekonomi, penjualan sepeda mungkin menurun karena orang lebih enggan mengeluarkan uang untuk hobi. Siklus ekonomi seperti ini dapat menyebabkan fluktuasi siklik dalam data penjualan sepeda.

### Contoh Python Script dengan Data Sintetis:

Untuk mengilustrasikan komponen-komponen seri waktu, kita akan menggunakan pustaka Python 'pandas' dan 'matplotlib' untuk membuat data sintetis dan meng gambarkannya dalam grafik. Berikut adalah contoh kode Python:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Membuat data sintetis
np.random.seed(42)
n = 200
t = np.linspace(0, 10, n)
trend = 0.5 * t
seasonal = np.sin(2 * np.pi * t)
cyclical = 0.1 * np.sin(2 * np.pi * t / 5)
noise = 0.2 * np.random.normal(size=n)
data = trend + seasonal + cyclical + noise

# Membuat DataFrame pandas
df = pd.DataFrame(data, columns=['Sales'])
df.index = pd.date_range(start='2010-01-01', periods=n, freq='MS')

# Menggambarkan komponen seri waktu
fig, ax = plt.subplots(4, 1, figsize=(10, 8), sharex=True)

ax[0].plot(df, label='Data Asli')
ax[0].set_title('Data Asli')
ax[0].legend()

ax[1].plot(df.index, trend, label='Tren')
ax[1].set_title('Tren')
ax[1].legend()

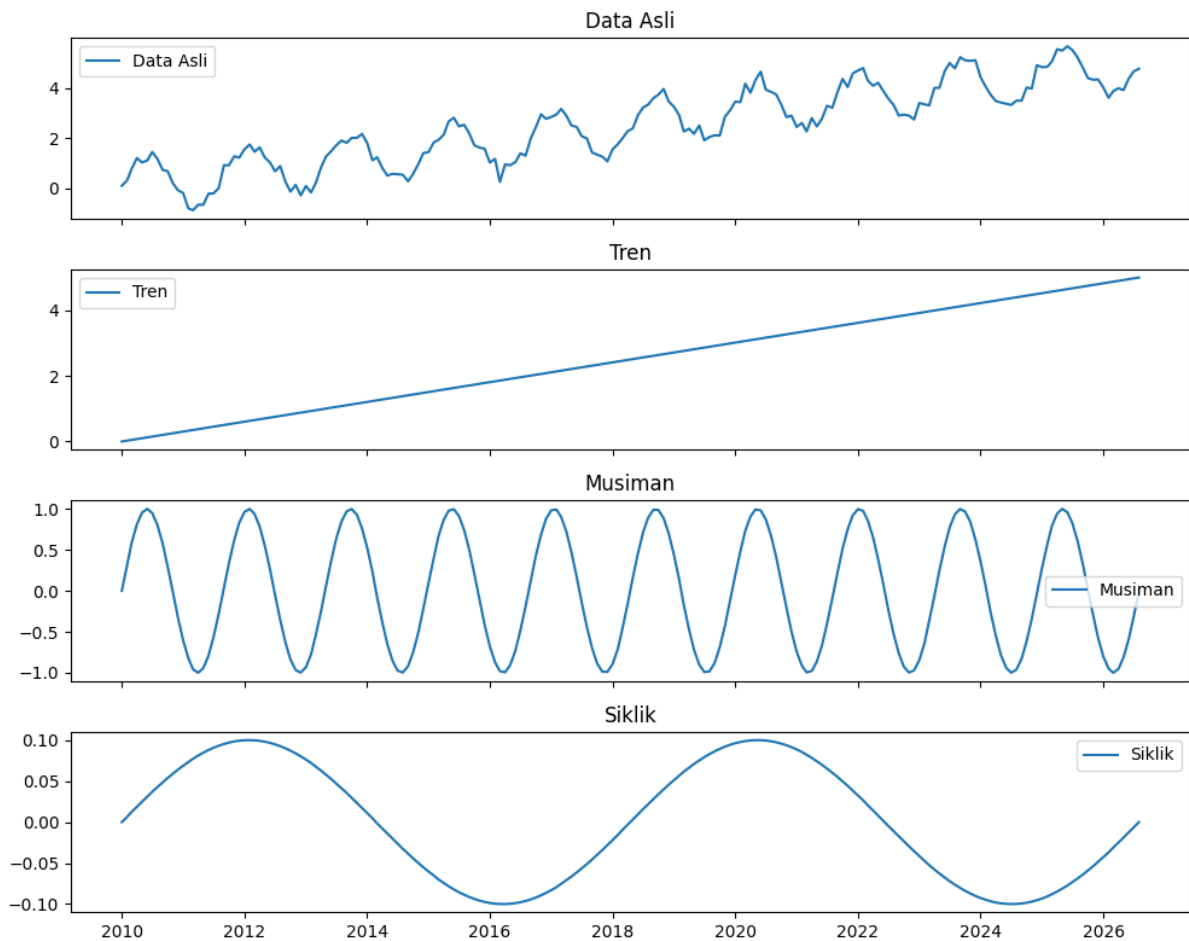
ax[2].plot(df.index, seasonal, label='Musiman')
ax[2].set_title('Musiman')
ax[2].legend()

ax[3].plot(df.index, cyclical, label='Siklik')
ax[3].set_title('Siklik')
ax[3].legend()

plt.tight_layout()
```

```
plt.show()
```

Output:



Dalam contoh ini, kita membuat data sintetis dengan menambahkan komponen tren, musiman, siklik, dan noise. Kemudian, kita menggambarkan setiap komponen dalam sub-plot yang berbeda menggunakan pustaka 'matplotlib'.

Tips dan Tricks dalam Penerapannya:

1. **Visualisasi Data:** Sebelum melakukan analisis seri waktu, sangat penting untuk menggambarkan data dan mencari pola yang mungkin ada dalam data. Visualisasi data juga bisa membantu mengidentifikasi pencilan atau kesalahan dalam data yang mungkin mempengaruhi hasil analisis.
2. **Transformasi Data:** Terkadang, transformasi data dapat membantu mengungkap pola yang sebelumnya tidak jelas. Transformasi logaritmik, differensiasi, dan pengurangan musiman adalah beberapa transformasi yang umum digunakan dalam analisis seri waktu.

3. Dekomposisi Seri Waktu: Untuk memahami komponen seri waktu dengan lebih baik, kamu dapat menggunakan metode dekomposisi, seperti dekomposisi musiman dan tren atau Seasonal and Trend Decomposition using Loess (STL). Metode ini membantu menguraikan data menjadi komponen-komponen yang lebih mudah dianalisis.

4. Pemilihan Model: Ketika memodelkan data seri waktu, penting untuk memilih model yang sesuai dengan karakteristik data. Misalnya, jika data memiliki tren dan musiman, kamu mungkin ingin menggunakan model yang memperhitungkan kedua komponen ini, seperti ARIMA atau Exponential Smoothing State Space Model (ETS).

5. Validasi Silang dan Evaluasi Model: Setelah memilih model, kamu harus memvalidasi dan mengevaluasi model untuk memastikan bahwa model tersebut mampu membuat prediksi yang akurat. Validasi silang dan metrik evaluasi seperti Mean Absolute Error (MAE) dan Mean Squared Error (MSE) dapat digunakan untuk mengukur kinerja model.

## 5.2. Teknik Smoothing: Moving Average dan Exponential Smoothing

Teknik smoothing digunakan untuk mereduksi fluktuasi dalam data dan mengungkap pola yang mendasari seperti tren dan musiman. Dua teknik smoothing yang umum digunakan dalam analisis seri waktu adalah Moving Average dan Exponential Smoothing.

### 5.2.1. Moving Average

Moving Average (MA) adalah teknik smoothing yang menghitung rata-rata dari sejumlah titik data sebelumnya. Teknik ini dapat digunakan untuk mengurangi noise dan mengungkap pola yang mendasari dalam data. Ada beberapa variasi dari Moving Average, seperti Simple Moving Average (SMA), Weighted Moving Average (WMA), dan Exponentially Weighted Moving Average (EWMA).

Matematis, Simple Moving Average (SMA) didefinisikan sebagai:

$$SMA_t = (x_{(t-1)} + x_{(t-2)} + \dots + x_{(t-n)}) / n$$

di mana  $x_t$  adalah nilai data pada waktu  $t$ ,  $n$  adalah jumlah periode yang digunakan dalam perhitungan rata-rata, dan  $SMA_t$  adalah nilai Simple Moving Average pada waktu  $t$ .

Contoh Kasus Dunia Nyata:

Misalkan kamu memiliki data penjualan bulanan suatu produk dari bulan Januari hingga Desember. Kamu ingin mengurangi fluktuasi bulanan dalam data dan mengidentifikasi tren penjualan yang mendasari. Salah satu cara untuk melakukannya adalah dengan menggunakan Moving Average.

Contoh Python Script dengan Data Sintetis:

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

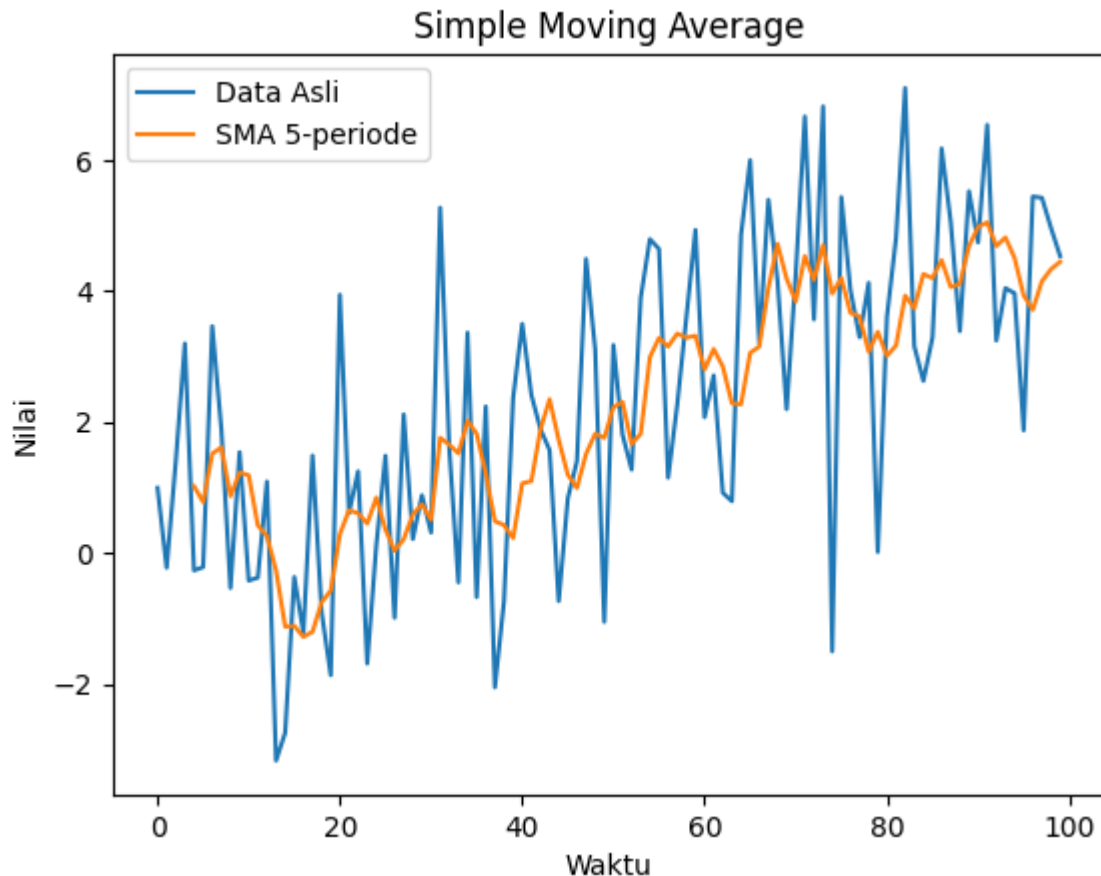
# Membuat data sintetis
np.random.seed(42)
n = 100
trend = np.linspace(0, 5, n)
noise = 2 * np.random.normal(size=n)
data = trend + noise

# Menghitung Simple Moving Average
window_size = 5
sma = pd.Series(data).rolling(window=window_size).mean()

# Menggambarkan data asli dan Simple Moving Average
plt.plot(data, label='Data Asli')
plt.plot(sma, label=f'SMA {window_size}-periode')
plt.legend()
plt.title('Simple Moving Average')
plt.xlabel('Waktu')
plt.ylabel('Nilai')
plt.show()

```

Output:



### 5.2.2. Exponential Smoothing

Exponential Smoothing adalah teknik smoothing yang memberikan bobot yang berkurang secara eksponensial kepada titik data sebelumnya. Dalam teknik ini, bobot yang lebih besar diberikan kepada data terbaru, sementara bobot yang lebih kecil diberikan kepada data yang lebih lama. Ada beberapa variasi dari Exponential Smoothing, seperti Simple Exponential Smoothing (SES), Holt's Linear Trend Exponential Smoothing, dan Holt-Winters Exponential Smoothing.

Matematis, Simple Exponential Smoothing (SES) didefinisikan sebagai:

$$S_t = \alpha * x_t + (1 - \alpha) * S_{(t-1)}$$

di mana  $x_t$  adalah nilai data pada waktu  $t$ ,  $\alpha$  adalah konstanta smoothing yang berkisar antara 0 dan 1, dan  $S_t$  adalah nilai Simple Exponential Smoothing pada waktu  $t$ .

Contoh Kasus Dunia Nyata:



Misalkan kamu memiliki data suhu harian suatu kota sepanjang tahun dan kamu ingin mengurangi fluktuasi harian dalam data untuk mengidentifikasi pola musiman. Salah satu cara untuk melakukannya adalah dengan menggunakan Exponential Smoothing.

Contoh Python Script dengan Data Sintetis:

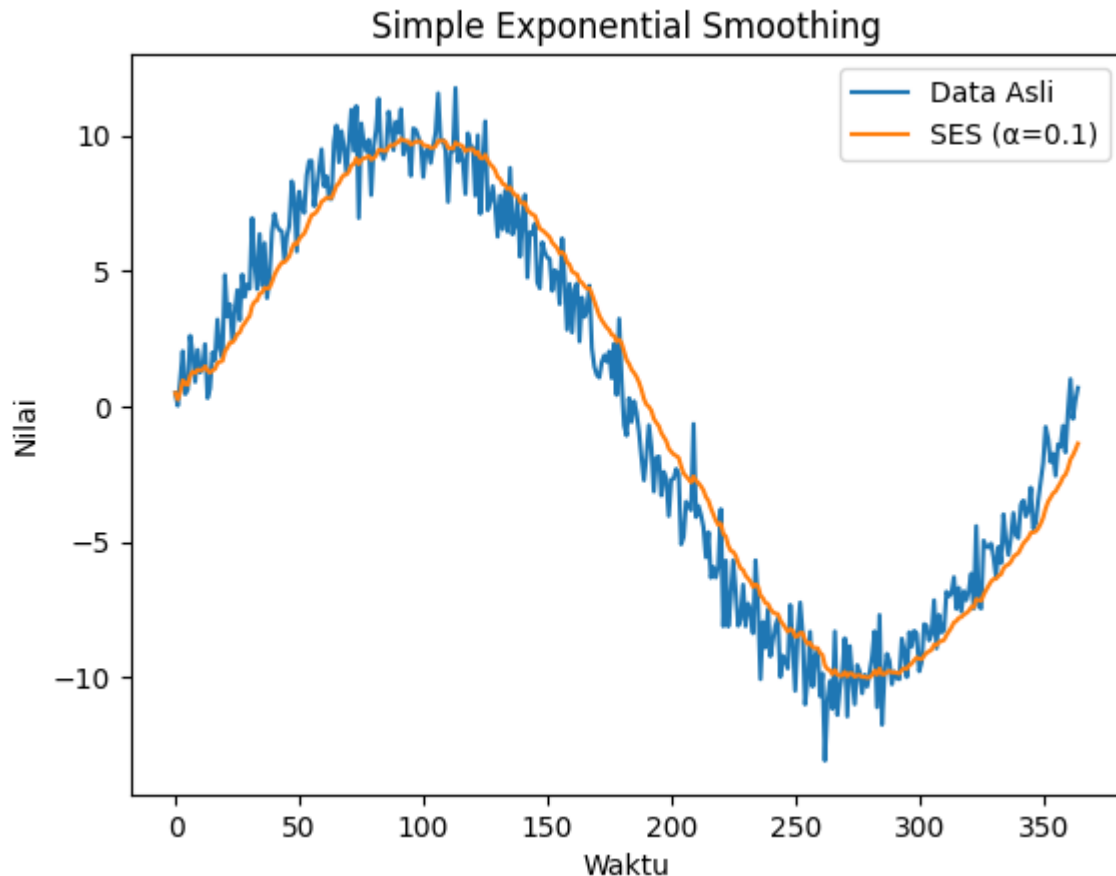
```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Membuat data sintetis
np.random.seed(42)
n = 365
seasonal = 10 * np.sin(np.linspace(0, 2 * np.pi, n))
noise = np.random.normal(size=n)
data = seasonal + noise

# Menghitung Simple Exponential Smoothing
alpha = 0.1
ses = pd.Series(data).ewm(alpha=alpha).mean()

# Menggambarkan data asli dan Simple Exponential Smoothing
plt.plot(data, label='Data Asli')
plt.plot(ses, label=f'SES ( $\alpha$ = $\alpha$ )')
plt.legend()
plt.title('Simple Exponential Smoothing')
plt.xlabel('Waktu')
plt.ylabel('Nilai')
plt.show()
```

Output:



#### Tips dan Tricks dalam Penerapan Teknik Smoothing:

Pilih metode smoothing yang sesuai: Metode smoothing yang dipilih harus sesuai dengan karakteristik data dan tujuan analisis. Jika data memiliki pola yang relatif sederhana, Simple Moving Average atau Simple Exponential Smoothing mungkin cukup. Jika data memiliki tren atau musiman yang lebih kompleks, metode yang lebih canggih seperti Holt-Winters Exponential Smoothing mungkin lebih sesuai.

Tentukan parameter dengan hati-hati: Parameter yang digunakan dalam metode smoothing, seperti ukuran jendela dalam Moving Average atau konstanta smoothing dalam Exponential Smoothing, dapat mempengaruhi hasil analisis. Parameter yang terlalu besar mungkin akan menghasilkan estimasi yang terlalu halus, sementara parameter yang terlalu kecil mungkin tidak cukup mengurangi noise. Parameter terbaik tergantung pada karakteristik data dan tujuan analisis.

Validasi model: Setelah memilih metode smoothing dan menentukan parameter, penting untuk memvalidasi model menggunakan teknik seperti cross-validation atau holdout validation. Dengan cara ini, kamu dapat memastikan bahwa model yang dipilih mampu menghasilkan hasil yang akurat dan bermakna.

Jangan mengandalkan satu metode saja: Dalam beberapa kasus, mungkin lebih baik menggunakan beberapa metode smoothing secara bersamaan untuk mendapatkan hasil yang lebih akurat dan dapat diinterpretasikan. Misalnya, kamu dapat menggunakan Moving Average untuk mengidentifikasi tren jangka panjang dan Exponential Smoothing untuk mengungkap pola musiman.

### 5.3. Model Autoregressive Integrated Moving Average (ARIMA)

Model Autoregressive Integrated Moving Average (ARIMA) adalah metode yang digunakan untuk memodelkan dan memprediksi data seri waktu. ARIMA menggabungkan tiga komponen utama: autoregressive (AR), differencing (I), dan moving average (MA). Model ini menggunakan informasi dari masa lalu untuk memperkirakan nilai masa depan dan dapat menangani data dengan tren dan komponen musiman.

Secara matematis, model ARIMA(p, d, q) didefinisikan sebagai:

$$(1) Y_t' = c + \phi_1 Y_{t-1}' + \phi_2 Y_{t-2}' + \dots + \phi_p Y_{t-p}' + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t$$

di mana:

$Y_t'$  adalah nilai data yang telah diberi differencing sebanyak d kali

c adalah konstanta

p adalah derajat autoregressive (AR)

d adalah tingkat differencing (I)

q adalah derajat moving average (MA)

$\varepsilon_t$  adalah kesalahan (white noise) pada waktu t

$\phi_i$  dan  $\theta_i$  adalah parameter yang perlu diestimasi

Penjelasan lebih lanjut tentang komponen-komponen ARIMA:

**Autoregressive (AR):** Autoregressive mengacu pada ketergantungan linier antara nilai saat ini dan nilai-nilai sebelumnya dalam seri waktu. Model AR(p) dinyatakan sebagai

$$Y_t = c + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t$$

**Differencing (I):** Differencing digunakan untuk membuat data stasioner dengan mengurangi nilai seri waktu sebelumnya. Tingkat differencing (d) menunjukkan berapa kali pengurangan yang perlu dilakukan untuk mencapai stasioneritas. Operasi differencing dinyatakan sebagai  $Y_t' = Y_t - Y_{t-1}$ , di mana  $Y_t'$  adalah data yang telah diberi differencing.

**Moving Average (MA):** Moving Average mengacu pada ketergantungan linier antara kesalahan saat ini dan kesalahan masa lalu. Model MA(q) dinyatakan sebagai

$$Y_t = c + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t$$

### Contoh Kasus di Dunia Nyata: Prediksi Jumlah Penumpang Pesawat

Misalkan kamu bekerja di maskapai penerbangan dan ingin memprediksi jumlah penumpang dalam beberapa bulan ke depan untuk membantu perencanaan sumber daya. Kamu memiliki data jumlah penumpang bulanan dalam beberapa tahun terakhir dan ingin menggunakan model ARIMA untuk membuat prediksi.

### Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.stattools import adfuller
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf

# Membuat data sintetis
np.random.seed(42)
n = 120
trend = 0.05 * np.arange(n)
seasonal = 3 * np.sin(2 * np.pi * np.arange(n) / 12)
noise = np.random.normal(0, 1, n)
data = trend + seasonal + noise

# Membuat DataFrame dari data sintetis
date_rng = pd.date_range(start='1/1/2015', end='12/1/2024', freq='MS')
df = pd.DataFrame(date_rng, columns=['date'])
df['passengers'] = data

# Menggambarkan data
plt.figure(figsize=(12, 6))
plt.plot(df['date'], df['passengers'])
plt.xlabel('Tanggal')
plt.ylabel('Jumlah Penumpang')
plt.title('Jumlah Penumpang Pesawat per Bulan')
plt.show()

# Mengecek stasioneritas menggunakan Augmented Dickey-Fuller Test
result = adfuller(df['passengers'])
print('ADF Statistic: %f' % result[0])
print('p-value: %f' % result[1])

# Melakukan differencing
```

```

df['passengers_diff'] = df['passengers'].diff().dropna()

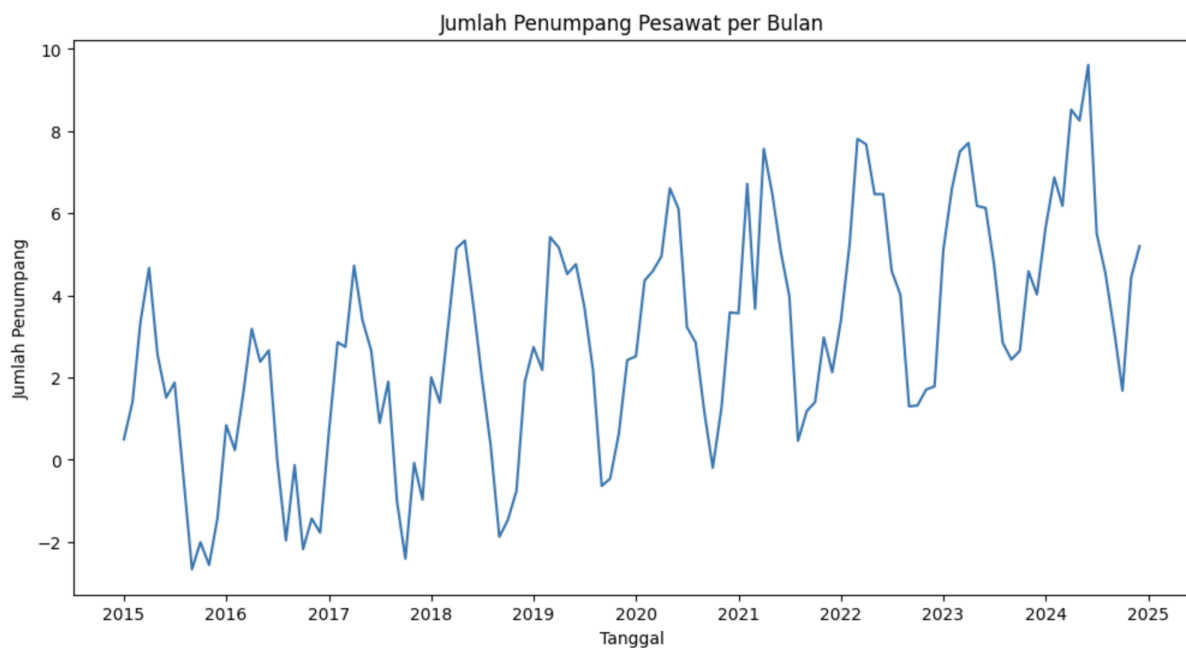
# Menggambarkan ACF dan PACF
plot_acf(df['passengers_diff'].dropna())
plot_pacf(df['passengers_diff'].dropna())
plt.show()

# Melatih model ARIMA
model = ARIMA(df['passengers'], order=(1, 1, 1))
model_fit = model.fit()

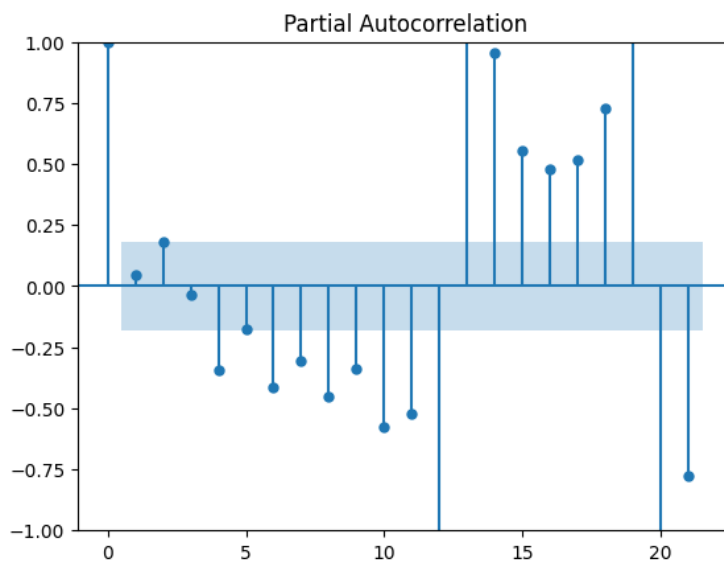
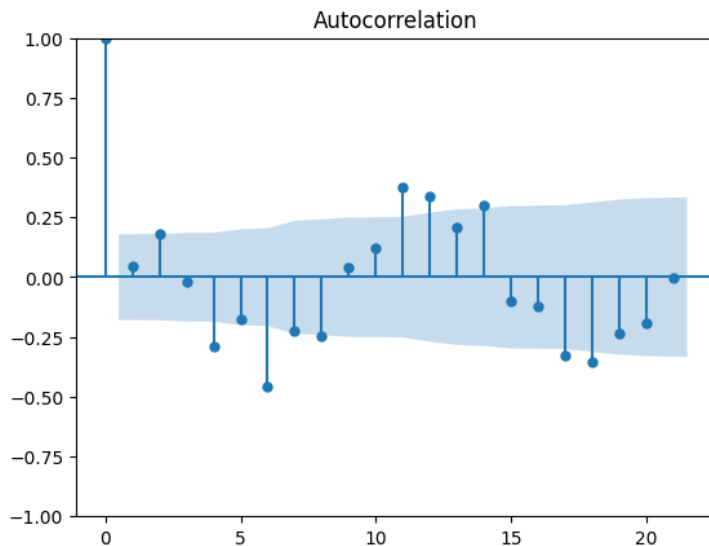
# Memprediksi jumlah penumpang untuk 6 bulan ke depan
forecast = model_fit.forecast(steps=6)
print('Prediksi jumlah penumpang untuk 6 bulan ke depan:')
print(forecast)

```

Output:



ADF Statistic: 0.633431  
p-value: 0.988411



Prediksi jumlah penumpang untuk 6 bulan ke depan:

120	5.307212
121	5.347870
122	5.363571
123	5.369635
124	5.371976
125	5.372881

### Tips dan Tricks dalam Penerapan ARIMA:

1. Membuat data stasioner: Pastikan untuk membuat data stasioner sebelum membangun model ARIMA. Gunakan metode seperti differencing dan transformasi log untuk mengurangi tren dan komponen musiman.
2. Menentukan nilai p, d, dan q: Gunakan plot ACF (Autocorrelation Function) dan PACF (Partial Autocorrelation Function) untuk menentukan nilai p dan q yang sesuai. Aturan praktis yang umum digunakan adalah:

- Jika ACF memotong garis signifikansi pada lag pertama dan PACF menurun secara perlahan, gunakan model AR(p).
- Jika ACF menurun secara perlahan dan PACF memotong garis signifikansi pada lag pertama, gunakan model MA(q).
- Jika ACF dan PACF menurun secara perlahan, gunakan model ARIMA(p, d, q).

3. Validasi model: Gunakan metrik evaluasi seperti Mean Absolute Error (MAE), Mean Squared Error (MSE), dan Root Mean Squared Error (RMSE) untuk memeriksa kinerja model pada data validasi. Jangan ragu untuk mencoba beberapa kombinasi p, d, dan q untuk mencari model terbaik.

4. Autokorelasi residual: Setelah membangun model, periksa autokorelasi residual dengan plot ACF. Jika tidak ada korelasi yang signifikan, maka model telah menangkap informasi yang ada dalam data.

5. Kesalahan prediksi: Ketika membuat prediksi, perhatikan bahwa prediksi akan menjadi kurang akurat seiring bertambahnya jangka waktu ke depan. Oleh karena itu, hindari membuat prediksi jangka panjang

## 5.4. Model Seasonal Decomposition of Time Series (STL)

Model Seasonal Decomposition of Time Series (STL) adalah teknik yang digunakan untuk memisahkan komponen musiman, tren, dan sisa (residual) dalam data seri waktu. STL menggunakan metode loess (locally estimated scatterplot smoothing) untuk mengestimasi tren dan komponen musiman. Metode ini fleksibel, memungkinkan variasi musiman untuk berubah sepanjang waktu dan mengakomodasi data yang memiliki pola musiman yang kompleks.

Secara matematis, STL mendekomposisi seri waktu  $Y_t$  menjadi tiga komponen:

$$Y_t = T_t + S_t + R_t$$

di mana:

$Y_t$  adalah nilai seri waktu pada waktu  $t$

$T_t$  adalah komponen tren pada waktu  $t$

$S_t$  adalah komponen musiman pada waktu  $t$

$R_t$  adalah komponen residual (sisa) pada waktu  $t$

Contoh Kasus di Dunia Nyata: Analisis Penjualan Produk

Misalkan kamu bekerja di perusahaan ritel dan ingin menganalisis penjualan produk selama beberapa tahun terakhir untuk mengidentifikasi pola musiman, tren, dan faktor lain yang mempengaruhi penjualan. Kamu dapat menggunakan metode STL untuk mendekomposisi data penjualan dan mendapatkan wawasan yang berguna.

Contoh Python Script dengan Data Sintetis:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.tsa.seasonal import STL

# Membuat data sintetis
np.random.seed(42)
n = 120
trend = 0.05 * np.arange(n)
seasonal = 3 * np.sin(2 * np.pi * np.arange(n) / 12)
noise = np.random.normal(0, 1, n)
data = trend + seasonal + noise

# Membuat DataFrame dari data sintetis
date_rng = pd.date_range(start='1/1/2015', end='12/1/2024', freq='MS')
df = pd.DataFrame({'date': date_rng, 'sales': data})
df = df.set_index('date')
df.index = pd.to_datetime(df.index)

# Menggambarkan data
plt.figure(figsize=(12, 6))
plt.plot(df.index, df['sales'])
plt.xlabel('Tanggal')
plt.ylabel('Penjualan')
plt.title('Penjualan Produk per Bulan (2015-2024)')
plt.show()

# Melakukan dekomposisi STL
stl = STL(df['sales'], seasonal=13)
result = stl.fit()

# Menggambarkan hasil dekomposisi
fig, axes = plt.subplots(4, 1, figsize=(12, 12), sharex=True)
axes[0].plot(df.index, df['sales'], label='Data Asli')
axes[0].set_ylabel('Penjualan')
axes[0].set_title('Data Asli')
axes[0].legend()

axes[1].plot(df.index, result.trend, label='Tren', color='C1')
axes[1].set_ylabel('Penjualan')
axes[1].set_title('Komponen Tren')
axes[1].legend()

axes[2].plot(df.index, result.seasonal, label='Musiman', color='C2')
```



```

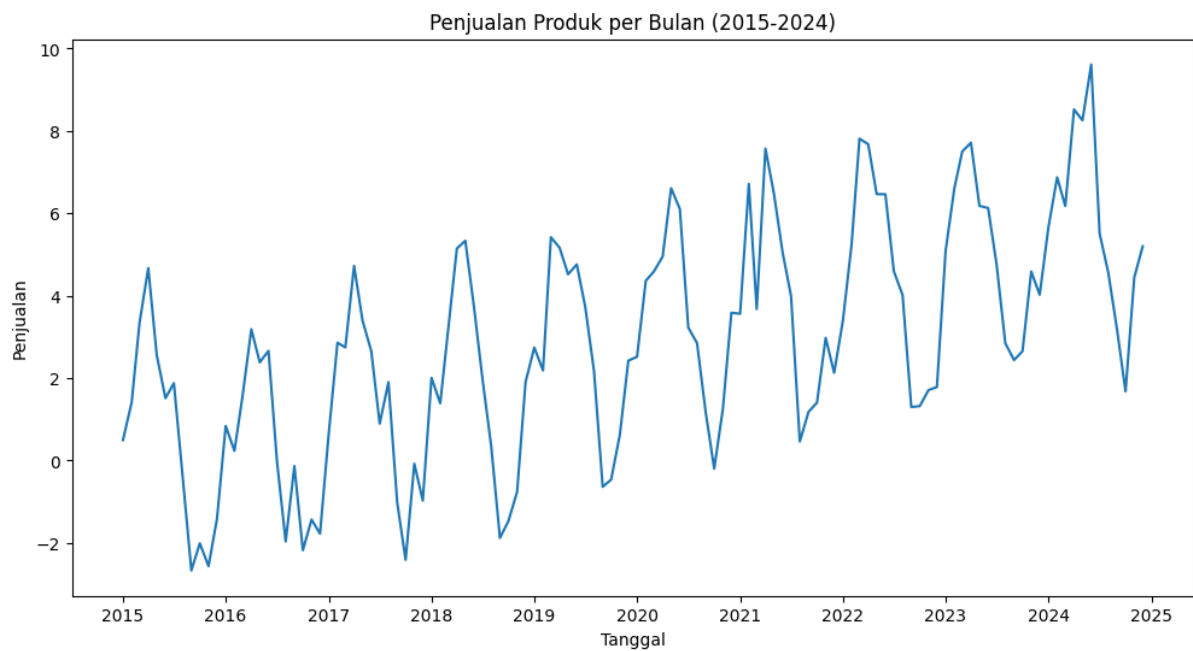
axes[2].set_ylabel('Penjualan')
axes[2].set_title('Komponen Musiman')
axes[2].legend()

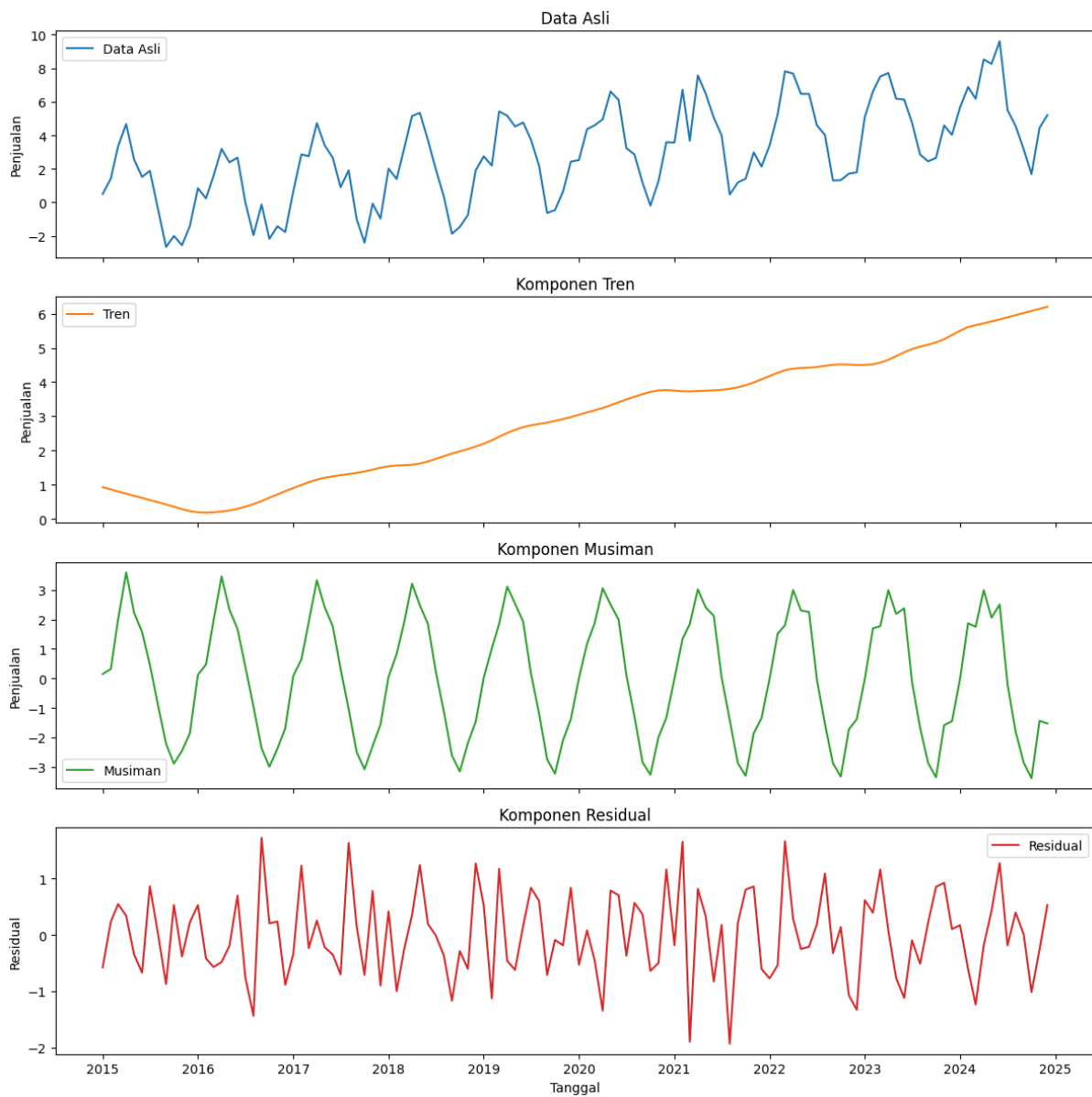
axes[3].plot(df.index, result.resid, label='Residual', color='C3')
axes[3].set_ylabel('Residual')
axes[3].set_title('Komponen Residual')
axes[3].legend()

plt.xlabel('Tanggal')
plt.tight_layout()
plt.show()

```

Output:





# Bab 6: Analisis Multivariat

## 6.1. Analisis Komponen Utama (PCA)

Analisis Komponen Utama (PCA) adalah teknik pengurangan dimensi yang banyak digunakan dalam analisis data multivariat. PCA bertujuan untuk mengidentifikasi kombinasi linier dari variabel yang menjelaskan sebagian besar varians dalam data, dengan asumsi bahwa varians yang lebih tinggi mengandung informasi yang lebih penting. Teknik ini berguna untuk mengurangi kompleksitas data, visualisasi, dan meningkatkan efisiensi analisis.

### Penjelasan Matematis PCA

PCA didasarkan pada transformasi ortogonal linier yang mengubah koordinat data asli ke sistem koordinat baru, di mana koordinat baru ditemukan dengan menggabungkan variabel asli. Koordinat baru disebut komponen utama. Secara matematis, PCA melibatkan beberapa langkah berikut:

**Standarisasi data:** Untuk menghindari pengaruh skala yang berbeda, kita perlu menstandarisasi data. Standarisasi mengubah setiap variabel ke dalam skala yang sama dengan mengurangi rata-rata dan membagi dengan standar deviasi. Dalam notasi matematik, ini dapat ditulis sebagai:

makefile

Copy code

$$Z = (X - \mu) / \sigma$$

di mana Z adalah data yang telah distandarisasi, X adalah data asli,  $\mu$  adalah rata-rata, dan  $\sigma$  adalah standar deviasi.

**Menghitung matriks kovarians:** Matriks kovarians adalah matriks yang menggambarkan hubungan antara setiap pasangan variabel dalam data. Kovarians antara dua variabel dihitung sebagai:

$$\text{cov}(X, Y) = \Sigma[(x - \mu_x)(y - \mu_y)] / (n - 1)$$

di mana x dan y adalah observasi untuk variabel X dan Y,  $\mu_x$  dan  $\mu_y$  adalah rata-rata variabel, dan n adalah jumlah observasi.

**Menghitung nilai eigen dan vektor eigen matriks kovarians:** Nilai eigen dan vektor eigen matriks kovarians menentukan komponen utama. Nilai eigen menggambarkan varians yang dijelaskan oleh masing-masing komponen utama, sedangkan vektor eigen menunjukkan arah komponen utama dalam ruang asli. Dalam PCA, kita mengurutkan komponen utama berdasarkan nilai eigen yang berhubungan, dari yang terbesar hingga yang terkecil.

Proyeksikan data ke ruang komponen utama: Langkah terakhir dalam PCA adalah memproyeksikan data asli ke ruang komponen utama. Proyeksi ini menciptakan koordinat baru yang ditemukan dengan menggabungkan variabel asli menggunakan vektor eigen. Proyeksi ini dapat dinyatakan dalam notasi matematik sebagai:

$$P = Z * V$$

di mana P adalah data yang diproyeksikan, Z adalah data yang telah distandarisasi, dan V adalah matriks yang terdiri dari vektor eigen.

#### Contoh Kasus di Dunia Nyata

PCA digunakan dalam berbagai aplikasi di dunia nyata, seperti pengurangan dimensi dalam analisis data besar, pengolahan citra, analisis genomik, dan analisis keuangan. Sebagai contoh, kita akan membahas aplikasi PCA dalam analisis sentimen pasar saham.

Misalkan kita memiliki data harga saham harian untuk 100 perusahaan selama satu tahun. Setiap perusahaan diwakili oleh satu variabel, sehingga kita memiliki 100 variabel dan 252 observasi (jumlah hari perdagangan dalam setahun). Tujuan kita adalah mengurangi dimensi data dan mengidentifikasi faktor-faktor yang mendasari pergerakan harga saham.

#### Contoh Python Script dengan Data Sintetis

Sebelum kita melanjutkan, pastikan kamu telah menginstal pustaka berikut: NumPy, pandas, dan scikit-learn. Kemudian, kita akan membuat data sintetis untuk menggambarkan proses PCA.

```
import numpy as np
import pandas as pd
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

# Buat data sintetis
np.random.seed(42)
n_samples = 252
n_features = 100
data = np.random.normal(size=(n_samples, n_features))

# Standarisasi data
scaler = StandardScaler()
data_standardized = scaler.fit_transform(data)

# Lakukan PCA
pca = PCA()
data_pca = pca.fit_transform(data_standardized)

# Ambil komponen utama pertama
first_component = data_pca[:, 0]
```

### Tips dan Trik dalam Penerapan PCA

Berikut beberapa tips dan trik yang dapat membantu kamu dalam menerapkan PCA:

**Standarisasi data:** Penting untuk menstandarisasi data sebelum menerapkan PCA, karena variabel dengan skala yang berbeda dapat mempengaruhi hasil PCA.

**Jumlah komponen utama:** Memilih jumlah komponen utama yang tepat sangat penting. Salah satu pendekatan umum adalah menggunakan jumlah komponen utama yang menjelaskan sejumlah persentase varians yang diinginkan (misalnya, 95% atau 99%).

**Interpretasi komponen utama:** PCA menghasilkan kombinasi linier dari variabel asli yang sulit untuk diinterpretasikan secara langsung. Kamu dapat menggunakan plot beban (loadings plot) untuk memvisualisasikan hubungan antara komponen utama dan variabel asli.

**Asumsi linearitas:** PCA mengasumsikan hubungan linier antara variabel, yang mungkin tidak selalu berlaku. Jika ada hubungan nonlinier antara variabel, pertimbangkan menggunakan teknik pengurangan dimensi lain seperti kernel PCA atau t-SNE.

**Kepekaan terhadap outlier:** PCA sensitif terhadap outlier dalam data. Pertimbangkan untuk membersihkan data atau menggunakan teknik pengurangan dimensi yang lebih tahan terhadap outlier, seperti teknik robust PCA.

Dalam kesimpulan, PCA adalah teknik yang berguna untuk mengurangi dimensi data multivariat dan mengidentifikasi pola dalam data. Dalam prakteknya, penting untuk menstandarisasi data, memilih jumlah komponen utama yang tepat, dan mempertimbangkan asumsi dan keterbatasan teknik ini.

## 6.2. Analisis Faktor

Analisis faktor merupakan teknik analisis multivariat yang digunakan untuk mengidentifikasi struktur laten (tersembunyi) dalam kumpulan data yang terdiri dari banyak variabel yang saling berkorelasi. Tujuan utama analisis faktor adalah mengurangi dimensi data dengan menggantikan variabel asli dengan sejumlah faktor yang lebih sedikit. Faktor-faktor ini adalah kombinasi linier dari variabel asli yang menjelaskan sebagian besar varians dalam data.

### Penjelasan Secara Matematis

Misalkan kita memiliki matriks data  $X$  dengan  $n$  observasi dan  $p$  variabel. Tujuan analisis faktor adalah menguraikan matriks  $X$  menjadi dua matriks, yaitu matriks faktor ( $F$ ) dan matriks beban ( $L$ ), sehingga  $X \approx F \times L'$ . Di sini,  $F$  adalah matriks  $n \times k$  ( $k < p$ ) yang berisi nilai-nilai faktor, dan  $L$  adalah matriks  $p \times k$  yang berisi beban faktor.

Secara matematis, kita mencari hubungan berikut antara variabel asli dan faktor laten:

$$X_i = l_{1i} * F_1 + l_{2i} * F_2 + \dots + l_{ki} * F_k + e_i$$

di mana  $X_i$  adalah variabel asli ke- $i$ ,  $l_{1i}$  hingga  $l_{ki}$  adalah beban faktor ke- $i$ ,  $F_1$  hingga  $F_k$  adalah faktor laten, dan  $e_i$  adalah kesalahan yang unik untuk variabel ke- $i$ . Beban faktor mengukur sejauh mana variabel asli berkorelasi dengan faktor laten, sedangkan kesalahan  $e_i$  menggambarkan varians yang tidak dijelaskan oleh faktor laten.

Untuk mengestimasi faktor dan beban, kita menggunakan metode ekstraksi faktor, seperti Principal Component Analysis (PCA), metode maksimum likelihood, atau metode rotasi seperti Varimax atau Promax.

#### Contoh Kasus di Dunia Nyata

Sebagai contoh, kita akan melihat kasus di mana analisis faktor dapat digunakan untuk mengidentifikasi faktor-faktor yang mendasari kepuasan pelanggan di sebuah restoran. Misalkan kita telah mengumpulkan data dari survei kepuasan pelanggan yang mencakup pertanyaan tentang kualitas makanan, harga, suasana, kebersihan, dan layanan. Kita ingin mengetahui faktor-faktor utama yang mempengaruhi kepuasan pelanggan dan mengurangi jumlah variabel yang harus dikelola.

#### Contoh Python Script dengan Data Sintetis

Untuk mengilustrasikan analisis faktor, kita akan membuat data sintetis menggunakan Python. Pastikan kamu telah menginstal pustaka berikut: NumPy, pandas, dan factor\_analyzer.

```
# Install factor_analyzer
!pip install factor_analyzer

# Import Library
import numpy as np
import pandas as pd
from factor_analyzer import FactorAnalyzer
from sklearn.preprocessing import StandardScaler

# Buat data sintetis
np.random.seed(42)
n_samples = 300
n_features = 10
data = np.random.normal(size=(n_samples, n_features))

# Standarisasi data
scaler = StandardScaler()
data_standardized = scaler.fit_transform(data)
```

```

# Lakukan analisis faktor
fa = FactorAnalyzer(n_factors=3, rotation='varimax')
fa.fit(data_standardized)

# Ambil beban faktor dan faktor yang dihasilkan
loadings = fa.loadings_
factors = fa.transform(data_standardized)

# Tampilkan beban faktor
loadings_df = pd.DataFrame(loadings, columns=['Faktor 1', 'Faktor 2',
'Faktor 3'])
print('Beban Faktor:')
print(loadings_df)

# Tampilkan faktor yang dihasilkan
factors_df = pd.DataFrame(factors, columns=['Faktor 1', 'Faktor 2',
'Faktor 3'])
print('\nFaktor Hasil Ekstraksi:')
print(factors_df.head())

```

Output:

```

Beban Faktor:
   Faktor 1  Faktor 2  Faktor 3
0  0.546325  0.341180  0.209942
1 -0.007048  0.267389 -0.210558
2  0.293597 -0.163541 -0.063248
3 -0.142028 -0.054572  0.396007
4  0.025165 -0.073973 -0.155630
5  0.051067 -0.094574  0.219176
6 -0.070196  0.169936  0.068435
7 -0.000504  0.074233 -0.007003
8 -0.128622 -0.011742  0.099466
9 -0.084691 -0.343983  0.069045

Faktor Hasil Ekstraksi:
   Faktor 1  Faktor 2  Faktor 3
0  0.124150  0.133550  0.715130
1  0.126779  0.180414 -0.866555
2  1.128546  0.456286 -0.177522
3 -0.192655  0.118110 -1.186913
4  0.348778  0.932963 -0.004693

```

Tips dan Tricks dalam Penerapannya

Pemilihan jumlah faktor: Gunakan metode seperti Scree Plot, eigenvalue di atas 1, atau uji paralel untuk memilih jumlah faktor yang optimal.

Standarisasi data: Sebelum melakukan analisis faktor, sebaiknya standarisasi data, terutama jika variabel memiliki skala yang berbeda.

Memeriksa asumsi: Pastikan data memenuhi asumsi analisis faktor, seperti multivariat normalitas, adanya korelasi antar variabel, dan sampel yang cukup besar.

Rotasi: Gunakan metode rotasi, seperti Varimax atau Promax, untuk memudahkan interpretasi beban faktor.

Validitas: Setelah mengidentifikasi faktor, validasi model dengan menguji validitas konstruk (misalnya, validitas konvergen dan diskriminan) dan reliabilitas (misalnya, koefisien alpha Cronbach atau komposit reliabilitas).

Dengan memahami konsep analisis faktor dan penerapannya menggunakan Python, kamu akan dapat mengurangi dimensi data dan mengidentifikasi struktur laten yang mendasari kumpulan data yang kompleks. Analisis faktor merupakan teknik yang berguna untuk mengeksplorasi data dan menginformasikan pengambilan keputusan dalam berbagai bidang, mulai dari penelitian pasar hingga psikologi dan ilmu sosial.

## **6.3 Analisis Klaster**

Analisis klaster merupakan teknik pengelompokan data multivariat berdasarkan karakteristik yang sama atau sejenis. Teknik ini bertujuan untuk mengelompokkan objek atau observasi ke dalam kelompok (klaster) yang lebih homogen atau serupa, sedangkan antar kelompok diharapkan heterogen atau berbeda. Dalam subbab ini, kita akan membahas dasar-dasar analisis klaster, metode-metode yang umum digunakan, dan cara mengaplikasikannya menggunakan Python.

### **6.3.1 Dasar-dasar Analisis Klaster**

Analisis klaster dapat dikategorikan menjadi dua jenis, yaitu metode pengelompokan hierarki dan non-hierarki.

Pengelompokan hierarki: Pengelompokan ini dilakukan dengan menggabungkan atau membagi klaster berdasarkan jarak antara objek atau kelompok. Dua metode pengelompokan hierarki yang umum digunakan adalah agglomerative (bottom-up) dan divisive (top-down).

Pengelompokan non-hierarki: Pengelompokan ini melibatkan proses iteratif untuk mengoptimalkan kriteria pengelompokan, seperti jarak antar objek atau kelompok. Contoh metode pengelompokan non-hierarki yang populer adalah k-means dan DBSCAN.

### **6.3.2 Metode Pengelompokan yang Umum Digunakan**

#### **K-means**

K-means merupakan algoritma pengelompokan non-hierarki yang populer dan mudah diimplementasikan. Algoritma ini mengelompokkan data ke dalam K klaster berdasarkan jarak antara objek dan centroid (titik pusat) klaster. Berikut langkah-langkah k-means:



Tentukan jumlah kluster K.  
Pilih secara acak K centroid awal.  
Hitung jarak antara setiap objek dan centroid.  
Tetapkan objek ke kluster dengan centroid terdekat.  
Perbarui centroid dengan menghitung rata-rata objek dalam kluster.  
Ulangi langkah 3-5 hingga konvergensi (centroid tidak berubah lagi atau perubahan di bawah ambang batas tertentu).

### **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**

DBSCAN merupakan algoritma pengelompokan non-hierarki yang berbasis kepadatan. Algoritma ini mengelompokkan objek berdasarkan kepadatan titik data di sekitarnya. Berikut langkah-langkah DBSCAN:

Tentukan parameter radius (eps) dan jumlah minimum titik data (min\_samples) dalam satu kluster.

Untuk setiap titik data yang belum dikunjungi, tandai sebagai dikunjungi dan tentukan titik-titik tetangga dalam radius eps.

Jika jumlah titik tetangga lebih besar atau sama dengan min\_samples, buat kluster baru dan tambahkan titik-titik tetangga ke kluster tersebut.

Jika jumlah titik tetangga kurang dari min\_samples, tandai titik data sebagai noise (bukan bagian dari kluster).

Ulangi langkah 2-4 hingga semua titik data telah dikunjungi dan diberi label kluster atau noise.

### **Pengelompokan Hierarki Agglomerative**

Pengelompokan hierarki agglomerative merupakan metode pengelompokan bottom-up yang mulai dengan menganggap setiap titik data sebagai kluster tunggal. Berikut langkah-langkahnya:

Hitung matriks jarak antara semua titik data.

Gabungkan dua kluster terdekat berdasarkan kriteria penggabungan (single linkage, complete linkage, average linkage, atau Ward's method).

Perbarui matriks jarak untuk mencerminkan penggabungan kluster.

Ulangi langkah 2 dan 3 hingga semua titik data digabungkan menjadi satu kluster.

### **6.3.3 Contoh Kasus dan Implementasi Python**

Sebagai contoh, kita akan menggunakan analisis kluster untuk mengelompokkan pelanggan sebuah perusahaan e-commerce berdasarkan perilaku pembelian mereka.

```
import numpy as np
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
```

```

# Membaca data sintetis
data = {'Pelanggan': ['A', 'B', 'C', 'D', 'E'],
        'Jumlah Transaksi': [5, 3, 1, 10, 2],
        'Total Belanja (USD)': [500, 300, 100, 1000, 200]}

df = pd.DataFrame(data)

# Normalisasi data menggunakan StandardScaler
scaler = StandardScaler()
data_scaled = scaler.fit_transform(df[['Jumlah Transaksi', 'Total
Belanja (USD)']])

# Melakukan pengelompokan dengan K-means (K=2)
kmeans = KMeans(n_clusters=2, random_state=42)
kmeans.fit(data_scaled)

# Menambahkan label klaster ke dataframe
df['Klaster'] = kmeans.labels_

print(df)

```

Output:

	Pelanggan	Jumlah Transaksi	Total Belanja (USD)	Klaster
0	A	5	500	1
1	B	3	300	1
2	C	1	100	1
3	D	10	1000	0
4	E	2	200	1

Dari hasil ini, kita dapat melihat bahwa pelanggan telah dikelompokkan ke dalam dua klaster berdasarkan jumlah transaksi dan total belanja mereka.

### 6.3.4 Tips dan Trik

**Normalisasi data:** Penting untuk menormalisasi data sebelum melakukan analisis klaster, terutama jika fitur memiliki skala yang berbeda. Hal ini akan membantu menghindari bias pada hasil pengelompokan akibat perbedaan skala antar fitur. Beberapa metode normalisasi yang umum digunakan adalah min-max scaling dan standard scaling.

**Memilih jumlah klaster:** Salah satu tantangan dalam analisis klaster adalah menentukan jumlah klaster yang optimal. Beberapa metode yang dapat digunakan untuk membantu menentukan jumlah klaster yang tepat adalah Elbow Method, Silhouette Score, dan Gap Statistic.

Memilih metode pengelompokan yang sesuai: Pemilihan metode pengelompokan yang tepat tergantung pada struktur data, jumlah fitur, dan tujuan analisis. K-means cocok untuk data yang memiliki kluster globular, sedangkan DBSCAN lebih cocok untuk data dengan kluster berbentuk tidak teratur dan kepadatan bervariasi. Pengelompokan hierarki agglomerative cocok untuk data dengan struktur hierarki atau ketika kita tertarik untuk mengeksplorasi berbagai jumlah kluster.

Interpretasi hasil: Setelah melakukan analisis kluster, penting untuk menginterpretasikan hasilnya secara menyeluruh dan memahami karakteristik setiap kluster. Hal ini dapat membantu dalam mengambil keputusan yang lebih baik dan menginformasikan strategi bisnis.

Evaluasi kluster: Karena analisis kluster merupakan metode pembelajaran tanpa pengawasan, evaluasi kualitas pengelompokan lebih sulit dibandingkan dengan metode pembelajaran dengan pengawasan. Beberapa metode yang dapat digunakan untuk mengevaluasi kualitas pengelompokan adalah inertia (within-cluster sum of squares), Silhouette Score, dan Dunn Index.

Dengan memahami dasar-dasar analisis kluster dan metode-metode yang umum digunakan, kamu akan dapat menerapkan teknik ini untuk mengelompokkan data dalam berbagai konteks dan menginformasikan strategi bisnis atau keputusan berdasarkan hasil pengelompokan. Selalu pastikan untuk menyesuaikan metode dan jumlah kluster dengan karakteristik data dan tujuan analisis kamu.

## **6.4. Analisis Diskriminan**

Analisis Diskriminan adalah salah satu metode analisis multivariat yang digunakan untuk mengklasifikasikan objek atau data ke dalam kelompok atau kategori yang sudah diketahui sebelumnya. Metode ini berfungsi untuk mencari kombinasi linier dari fitur yang memaksimalkan perbedaan antar-kelompok dan meminimalkan perbedaan dalam kelompok. Analisis Diskriminan umumnya digunakan dalam bidang statistik, ilmu sosial, ekonomi, dan pemasaran.

Ada dua jenis analisis diskriminan, yaitu Linear Discriminant Analysis (LDA) dan Quadratic Discriminant Analysis (QDA). LDA mengasumsikan bahwa kovarians antar-kelompok sama, sedangkan QDA tidak mengasumsikan hal tersebut. Dalam subbab ini, kita akan membahas matematika yang mendasari metode LDA dan QDA, contoh kasus dunia nyata, contoh Python script dengan data sintetis, serta tips dan trik dalam penerapannya.

### **Linear Discriminant Analysis (LDA)**

Matematika LDA

LDA bertujuan untuk mencari kombinasi linier dari fitur yang memaksimalkan rasio antara variasi antar-kelompok dan variasi dalam kelompok. Dalam konteks klasifikasi dua kelas, LDA mencari vektor  $w$  yang memaksimalkan rasio berikut:

$$w = \operatorname{argmax} (SB / SW)$$

di mana  $SB$  adalah variasi antar-kelompok dan  $SW$  adalah variasi dalam kelompok.  $SB$  dan  $SW$  didefinisikan sebagai:

$$SB = (\mu_1 - \mu_2) * (\mu_1 - \mu_2)' \quad SW = \Sigma_1 + \Sigma_2$$

di mana  $\mu_1$  dan  $\mu_2$  adalah rata-rata fitur dari kelompok 1 dan kelompok 2,  $\Sigma_1$  dan  $\Sigma_2$  adalah matriks kovarians dari kelompok 1 dan kelompok 2, dan  $'$  adalah transpose dari vektor.

Untuk menemukan vektor  $w$  yang memaksimalkan rasio  $SB / SW$ , kita dapat menggunakan metode eigenvalue-eigenvector. Dalam kasus dua kelas, solusi  $w$  adalah eigenvector yang berkaitan dengan eigenvalue terbesar dari matriks berikut:

$$M = SW^{(-1)} * SB$$

#### Contoh Kasus Dunia Nyata

Salah satu contoh kasus dunia nyata dari LDA adalah pengklasifikasian jenis anggur berdasarkan komposisi kimianya. Dalam studi ini, kita memiliki data yang terdiri dari konsentrasi 13 senyawa kimia yang berbeda dalam sampel anggur dari tiga jenis anggur yang berbeda. Tujuan dari analisis ini adalah untuk mengklasifikasikan jenis anggur berdasarkan komposisi kimianya. LDA dapat digunakan untuk mengurangi dimensi data dari 13 menjadi 2 atau 3, yang kemudian dapat digunakan untuk mengklasifikasikan jenis anggur dengan lebih mudah.

#### Contoh Python Script dengan Data Sintetis

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_classification
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix

# Membuat data sintetis
X, y = make_classification(n_samples=200, n_features=2, n_informative=2,
n_redundant=0, n_classes=2, random_state=42)

# Membagi data menjadi train dan test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)
```

```

# Melakukan LDA
lda = LinearDiscriminantAnalysis()
lda.fit(X_train, y_train)

# Memprediksi Label kelas
y_pred = lda.predict(X_test)

# Menghitung akurasi dan confusion matrix
acc = accuracy_score(y_test, y_pred)
cm = confusion_matrix(y_test, y_pred)

print("Akurasi: {:.2f}".format(acc))
print("Confusion Matrix:\n", cm)

# Visualisasi data dan decision boundary
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')
plt.xlabel('Fitur 1')
plt.ylabel('Fitur 2')

xlim = plt.gca().get_xlim()
ylim = plt.gca().get_ylim()

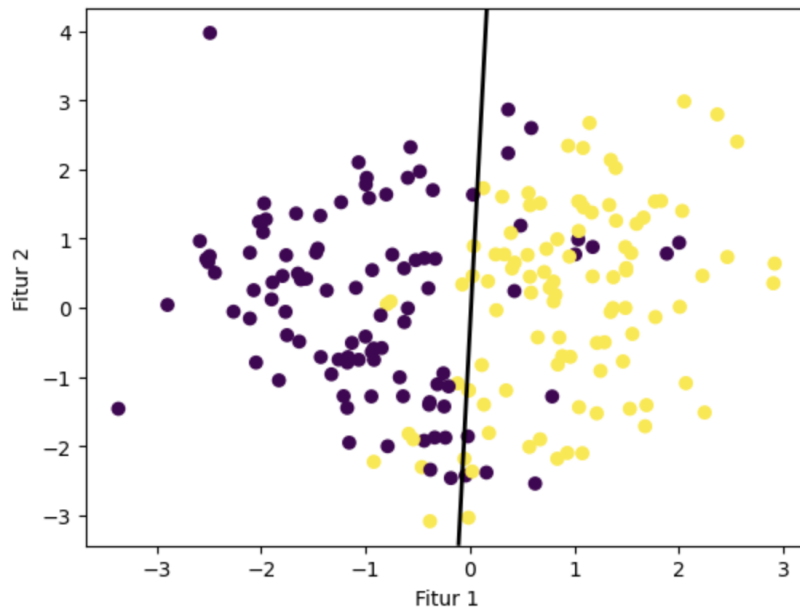
xx, yy = np.meshgrid(np.linspace(xlim[0], xlim[1], 200),
np.linspace(ylim[0], ylim[1], 200))
Z = lda.predict_proba(np.c_[xx.ravel(), yy.ravel()])
Z = Z[:, 1].reshape(xx.shape)

plt.contour(xx, yy, Z, levels=[0.5], linewidths=2, colors='k')
plt.show()

```

Output:

Akurasi: 0.83  
 Confusion Matrix:  
 [[26 6]  
 [ 4 24]]



#### Tips dan Trik dalam Penerapan LDA

LDA sangat sensitif terhadap skala fitur. Sebelum menerapkan LDA, disarankan untuk melakukan standardisasi fitur.

Jika asumsi kovarians yang sama antar-kelompok tidak terpenuhi, maka sebaiknya menggunakan QDA sebagai alternatif.

LDA bekerja dengan baik ketika ada hubungan linier antara fitur dan kelas target. Jika hubungan tersebut tidak linier, metode lain seperti Support Vector Machines (SVM) atau Neural Networks mungkin lebih efektif.

#### Quadratic Discriminant Analysis (QDA)

##### Matematika QDA

Berbeda dengan LDA, QDA tidak mengasumsikan kovarians yang sama antar-kelompok. QDA menghitung diskriminan kuadratik sebagai berikut:

$$D(x) = (x - \mu_i)' * \Sigma_i^{-1} * (x - \mu_i)$$

di mana  $x$  adalah vektor fitur,  $\mu_i$  adalah rata-rata fitur untuk kelompok  $i$ , dan  $\Sigma_i$  adalah matriks kovarians untuk kelompok  $i$ . Klasifikasi dilakukan dengan menghitung diskriminan kuadratik untuk setiap kelompok dan mengklasifikasikan  $x$  ke dalam kelompok yang memiliki nilai diskriminan kuadratik terendah.

##### Contoh Kasus Dunia Nyata

Salah satu contoh penggunaan QDA adalah dalam klasifikasi penyakit jantung berdasarkan data medis seperti usia, tekanan darah, kadar kolesterol, dan lainnya. Dalam kasus ini, QDA

dapat membantu mengidentifikasi pola yang berbeda antara pasien dengan penyakit jantung dan pasien tanpa penyakit jantung. Metode ini dapat membantu dokter dalam membuat keputusan diagnostik yang lebih akurat dan efisien.

Contoh Python Script dengan Data Sintetis

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_classification
from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix

# Membuat data sintetis
X, y = make_classification(n_samples=200, n_features=2, n_informative=2,
n_redundant=0, n_classes=2, random_state=42)

# Membagi data menjadi train dan test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)

# Melakukan QDA
qda = QuadraticDiscriminantAnalysis()
qda.fit(X_train, y_train)

# Memprediksi label kelas
y_pred = qda.predict(X_test)

# Menghitung akurasi dan confusion matrix
acc = accuracy_score(y_test, y_pred)
cm = confusion_matrix(y_test, y_pred)

print("Akurasi: {:.2f}".format(acc))
print("Confusion Matrix:\n", cm)

# Visualisasi data dan decision boundary
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')
plt.xlabel('Fitur 1')
plt.ylabel('Fitur 2')

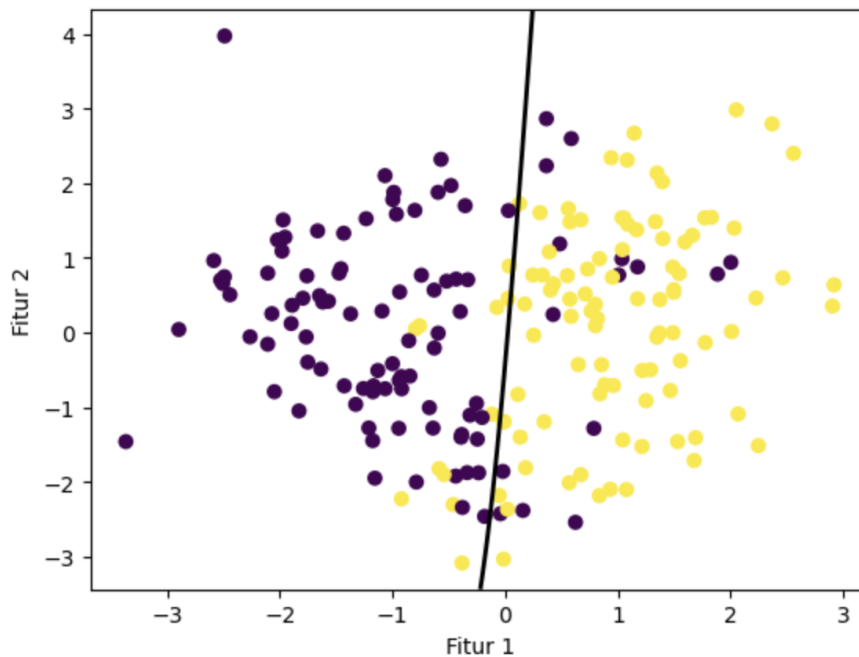
xlim = plt.gca().get_xlim()
ylim = plt.gca().get_ylim()

xx, yy = np.meshgrid(np.linspace(xlim[0], xlim[1], 200),
np.linspace(ylim[0], ylim[1], 200))
Z = qda.predict_proba(np.c_[xx.ravel(), yy.ravel()])
Z = Z[:, 1].reshape(xx.shape)
```

```
plt.contour(xx, yy, Z, levels=[0.5], linewidths=2, colors='k')
plt.show()
```

Output:

```
Akurasi: 0.83
Confusion Matrix:
[[26  6]
 [ 4 24]]
```



#### Tips dan Trik dalam Penerapan QDA

Sama seperti LDA, QDA sangat sensitif terhadap skala fitur. Sebelum menerapkan QDA, disarankan untuk melakukan standardisasi fitur.

Jika asumsi kovarians yang berbeda antar-kelompok terpenuhi, maka QDA dapat menjadi pilihan yang lebih baik daripada LDA.

QDA lebih fleksibel daripada LDA karena dapat menangkap hubungan yang lebih kompleks antara fitur dan kelas target. Namun, QDA juga lebih rentan terhadap overfitting, terutama jika jumlah sampel kecil. Dalam kasus ini, metode lain seperti Support Vector Machines (SVM) atau Neural Networks mungkin lebih efektif.

Dengan ini, kita telah menyelesaikan penjelasan subbab 6.4. Kamu sekarang memiliki pemahaman yang lebih baik tentang analisis diskriminan dan cara menggunakannya dalam analisis data multivariat. Selanjutnya, kamu dapat mulai menerapkan metode-metode ini dalam proyek analisis data kamu sendiri dan mengevaluasi seberapa efektif mereka dalam mengungkap pola dan hubungan yang tersembunyi dalam data kamu.



# Bab 7: Studi Kasus dan Aplikasi dalam Industri

## 7.1. Prediksi Penjualan dengan Regresi Linier

Regresi linier adalah salah satu teknik analisis statistik yang paling umum digunakan dalam industri. Dalam studi kasus ini, kita akan melihat bagaimana regresi linier dapat digunakan untuk memprediksi penjualan produk berdasarkan fitur-fitur yang relevan, seperti biaya periklanan, jumlah penjualan historis, dan harga produk. Kita akan menggunakan contoh data sintetis untuk membangun model regresi linier dan mengevaluasi kinerjanya.

### Contoh Kasus Dunia Nyata

Bayangkan kamu bekerja sebagai data scientist di perusahaan yang menjual produk elektronik. Manajemen ingin meningkatkan penjualan produk dengan memahami faktor-faktor yang mempengaruhi penjualan dan merancang strategi pemasaran yang efektif.

Kamu telah diberi data penjualan historis yang mencakup biaya periklanan, harga produk, dan jumlah unit yang terjual. Tugas kamu adalah membangun model prediktif yang dapat digunakan untuk memprediksi penjualan produk berdasarkan fitur-fitur ini.

### Contoh Python Script dengan Data Sintetis

Untuk membangun model prediksi, kita akan menggunakan regresi linier berganda. Berikut adalah langkah-langkah yang perlu kita ikuti:

Data preprocessing

Membagi data menjadi train dan test set

Melakukan hyperparameter tuning

Melatih model regresi linier

Evaluasi kinerja model

Diskusi next improvements

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.model_selection import GridSearchCV
from sklearn.preprocessing import StandardScaler
```

```

# Membuat data sintetis
np.random.seed(42)
n_samples = 200
ad_cost = np.random.rand(n_samples) * 5000
price = np.random.rand(n_samples) * 100 + 50
sales = ad_cost * 0.4 + price * (-0.8) + 5000 + np.random.normal(0, 500,
n_samples)

data = pd.DataFrame({'ad_cost': ad_cost, 'price': price, 'sales':
sales})

# Data preprocessing
X = data[['ad_cost', 'price']]
y = data['sales']

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Membagi data menjadi train dan test set
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y,
test_size=0.3, random_state=42)

# Melakukan hyperparameter tuning
params = {'fit_intercept': [True, False]}
grid_search = GridSearchCV(LinearRegression(), params, cv=5,
scoring='neg_mean_squared_error')
grid_search.fit(X_train, y_train)
best_params = grid_search.best_params_

# Melatih model regresi linier
model = LinearRegression(**best_params)
model.fit(X_train, y_train)

# Memprediksi penjualan di test set
y_pred = model.predict(X_test)

# Evaluasi kinerja model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print("Mean Squared Error (MSE):", mse)
print("R-squared (R2):", r2)

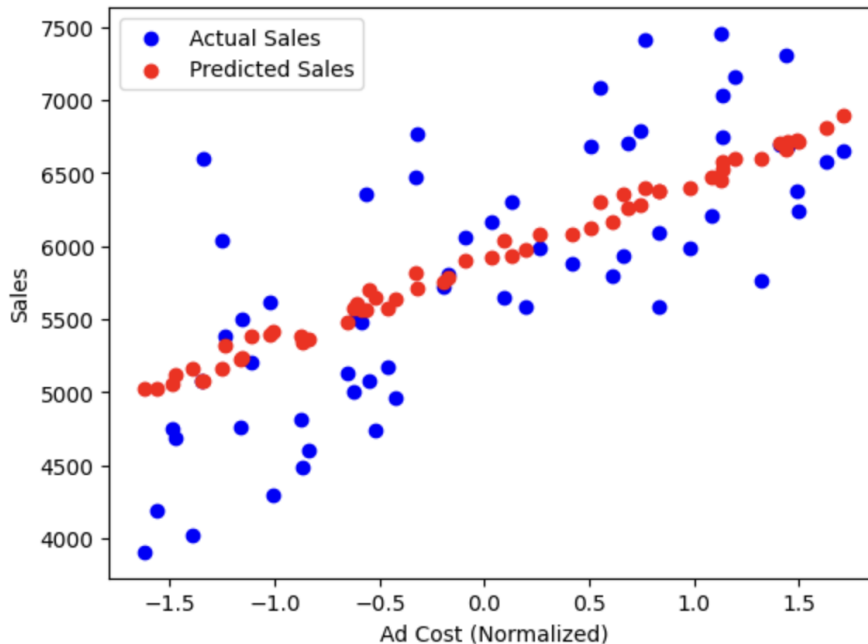
# Plotting hasil prediksi
plt.scatter(X_test[:, 0], y_test, label='Actual Sales', color='b')
plt.scatter(X_test[:, 0], y_pred, label='Predicted Sales', color='r')
plt.xlabel('Ad Cost (Normalized)')

```

```
plt.ylabel('Sales')
plt.legend()
plt.show()
```

Output:

Mean Squared Error (MSE): 362314.50086969853  
R-squared (R2): 0.5389042616191222



#### Next Improvements

Setelah mengevaluasi kinerja model regresi linier, kita dapat melihat beberapa area yang dapat ditingkatkan:

**Fitur Engineering:** Kita dapat mencoba menambahkan fitur-fitur tambahan yang relevan, seperti data demografi, preferensi konsumen, atau data ekonomi makro, yang mungkin membantu meningkatkan prediksi penjualan.

**Seleksi Fitur:** Beberapa fitur mungkin tidak memiliki hubungan yang signifikan dengan penjualan, dan menghilangkan fitur-fitur yang tidak relevan dapat meningkatkan kinerja model. Teknik seleksi fitur, seperti Recursive Feature Elimination (RFE) atau LASSO, dapat digunakan untuk mengidentifikasi fitur-fitur terpenting.

**Model Alternatif:** Regresi linier adalah model sederhana yang mungkin tidak dapat menangkap hubungan yang kompleks antara fitur dan target. Kita dapat mencoba model yang lebih canggih, seperti regresi Ridge atau Lasso, atau bahkan metode non-parametrik seperti pohon keputusan, random forest, atau boosting.

Cross-Validation: Gunakan teknik cross-validation yang lebih canggih, seperti K-Fold Cross-Validation atau Time Series Cross-Validation, untuk menghindari overfitting dan memperoleh perkiraan yang lebih baik tentang kinerja model pada data yang tidak terlihat.

Hyperparameter Tuning: Lakukan pencarian lebih ekstensif untuk menemukan kombinasi hyperparameter yang optimal, menggunakan metode seperti Randomized Search CV atau Bayesian Optimization.

Dengan menerapkan perbaikan-perbaikan ini, kita dapat meningkatkan kinerja model regresi linier dan memberikan wawasan yang lebih baik tentang faktor-faktor yang mempengaruhi penjualan produk. Hasil dari analisis ini dapat digunakan oleh manajemen untuk merancang strategi pemasaran yang lebih efektif dan meningkatkan penjualan produk di masa depan.

## 7.2. Segmentasi Pelanggan dengan Klastering

Segmentasi pelanggan adalah proses mengelompokkan pelanggan ke dalam segmen berdasarkan karakteristik yang serupa, seperti perilaku pembelian, preferensi produk, dan demografi. Tujuan segmentasi pelanggan adalah untuk mengidentifikasi kebutuhan dan preferensi setiap segmen, sehingga perusahaan dapat merancang strategi pemasaran yang lebih efektif dan menargetkan promosi serta penawaran khusus ke setiap segmen.

Dalam studi kasus ini, kita akan menggunakan teknik klastering untuk melakukan segmentasi pelanggan berdasarkan data transaksi pembelian. Kita akan menggunakan algoritma K-Means, yang merupakan metode klastering yang populer dan mudah diimplementasikan.

### Contoh Kasus Dunia Nyata

Sebuah perusahaan e-commerce ingin mengoptimalkan strategi pemasaran mereka dengan lebih memahami perilaku pelanggan mereka. Perusahaan ini telah mengumpulkan data transaksi dari beberapa pelanggan, termasuk jumlah pembelian, total pengeluaran, dan frekuensi kunjungan ke situs web. Tujuan analisis ini adalah untuk mengelompokkan pelanggan ke dalam segmen berdasarkan pola pembelian mereka dan mengidentifikasi karakteristik utama dari setiap segmen.

### Data Sintetis

Kita akan menggunakan data sintetis yang mencerminkan data transaksi pelanggan dalam contoh ini. Data tersebut mencakup tiga fitur:

Jumlah\_pembelian: jumlah total pembelian yang dilakukan oleh pelanggan

Total\_pengeluaran: jumlah total yang dihabiskan pelanggan dalam pembelian

Frekuensi\_kunjungan: jumlah kunjungan yang dilakukan pelanggan ke situs web

```
import numpy as np
```

```

import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import silhouette_score

# Membuat data sintesis
np.random.seed(42)
jumlah_pembelian = np.random.randint(1, 100, 200)
total_pengeluaran = np.random.randint(1000, 10000, 200)
frekuensi_kunjungan = np.random.randint(1, 20, 200)

data = pd.DataFrame({'Jumlah_pembelian': jumlah_pembelian,
'Total_pengeluaran': total_pengeluaran, 'Frekuensi_kunjungan':
frekuensi_kunjungan})

# Menampilkan 5 data pertama
print(data.head())

```

Output:

	Jumlah_pembelian	Total_pengeluaran	Frekuensi_kunjungan
0	52	4748	6
1	93	1663	13
2	15	2998	15
3	72	8994	3
4	61	2495	8

## Data Preprocessing

Sebelum melakukan klustering, kita perlu melakukan beberapa langkah preprocessing:

Normalisasi fitur: Karena K-Means menggunakan jarak Euclidean, kita perlu normalisasi fitur agar memiliki rentang yang serupa. Ini akan memastikan bahwa semua fitur memiliki pengaruh yang sama pada proses klustering.

```

scaler = StandardScaler()
data_scaled = scaler.fit_transform(data)

```

## Klustering menggunakan K-Means

Untuk menemukan jumlah kluster yang optimal, kita akan menggunakan metode Elbow dan skor Silhouette:

```

# Menemukan jumlah kluster optimal menggunakan metode Elbow dan skor
Silhouette
inertia = []
silhouette_scores = []
K = range(2, 11)

```

```

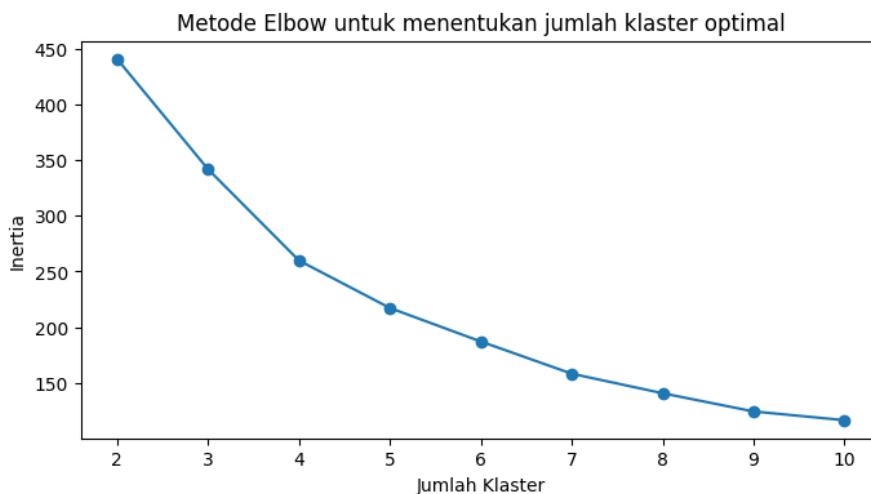
for k in K:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(data_scaled)
    inertia.append(kmeans.inertia_)
    silhouette_scores.append(silhouette_score(data_scaled,
kmeans.labels_))

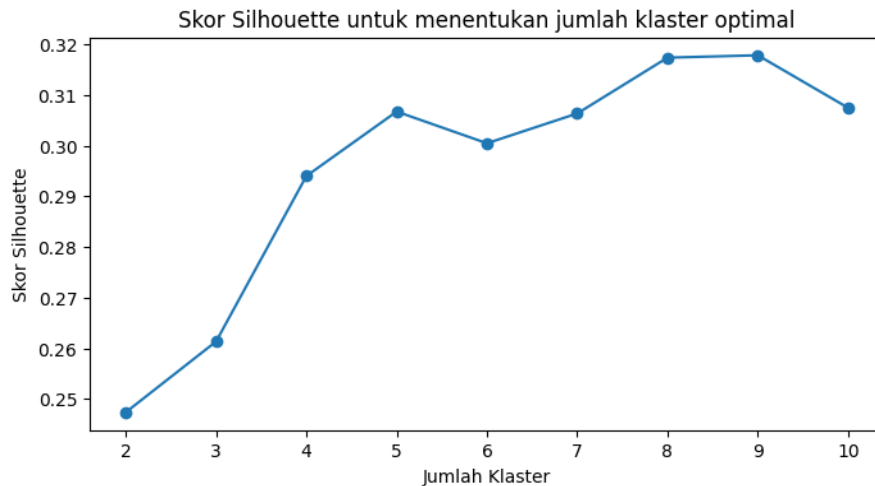
# Membuat plot metode Elbow
plt.figure(figsize=(8, 4))
plt.plot(K, inertia, marker='o')
plt.xlabel('Jumlah Klaster')
plt.ylabel('Inertia')
plt.title('Metode Elbow untuk menentukan jumlah klaster optimal')
plt.show()

# Membuat plot skor Silhouette
plt.figure(figsize=(8, 4))
plt.plot(K, silhouette_scores, marker='o')
plt.xlabel('Jumlah Klaster')
plt.ylabel('Skor Silhouette')
plt.title('Skor Silhouette untuk menentukan jumlah klaster optimal')
plt.show()

```

Output:





Dari kedua plot di atas, kita bisa melihat bahwa jumlah kluster yang optimal adalah 5. Sekarang kita akan melakukan klustering menggunakan K-Means dengan jumlah kluster 5:

```
# Melakukan klustering dengan K-Means
kmeans = KMeans(n_clusters=5, random_state=42)
kmeans.fit(data_scaled)

# Menambahkan label kluster ke data asli
data['Klaster'] = kmeans.labels_

# Menampilkan 5 data pertama dengan label kluster
data.head()
```

Output:

	Jumlah_pembelian	Total_pengeluaran	Frekuensi_kunjungan	Klaster
0	52	4748	6	1
1	93	1663	13	0
2	15	2998	15	2
3	72	8994	3	4
4	61	2495	8	1

Evaluasi Performa dan Next Improvements

Untuk mengevaluasi performa klustering, kita bisa menggunakan metrik seperti skor Silhouette dan inertia. Selain itu, kita juga bisa menganalisis karakteristik dari setiap kluster, seperti rata-rata jumlah pembelian, total pengeluaran, dan frekuensi kunjungan.

```
# Menghitung karakteristik setiap kluster
data.groupby('Klaster').mean()
```

Output:

	Jumlah_pembelian	Total_pengeluaran	Frekuensi_kunjungan
<b>Klaster</b>			
0	77.861111	3421.805556	14.777778
1	57.000000	2993.113636	5.136364
2	20.920000	5528.080000	14.660000
3	29.763158	7605.500000	4.500000
4	74.531250	7892.437500	11.250000

Dari hasil di atas, kita dapat melihat karakteristik utama dari setiap klaster dan merancang strategi pemasaran yang sesuai.

### 7.3. Peramalan Stok dengan Analisis Seri Waktu

Peramalan stok merupakan kegiatan yang sangat penting dalam manajemen rantai pasokan. Dengan melakukan peramalan yang akurat, perusahaan dapat mengurangi biaya penyimpanan, menghindari kekurangan stok, dan meningkatkan kepuasan pelanggan. Analisis seri waktu adalah salah satu metode yang dapat digunakan untuk peramalan stok.

#### Contoh Kasus

Bayangkan sebuah perusahaan yang menjual produk elektronik dan ingin meramalkan stok produk untuk 6 bulan ke depan. Perusahaan tersebut memiliki data penjualan bulanan untuk produknya selama 5 tahun terakhir. Untuk menyederhanakan kasus, kita akan menggunakan data sintetis yang mencerminkan pola penjualan produk.

#### Data Preprocessing

Langkah pertama dalam analisis seri waktu adalah melakukan preprocessing data. Preprocessing melibatkan mengonversi data mentah menjadi bentuk yang lebih mudah dianalisis. Dalam kasus ini, kita akan menggunakan data sintetis yang mencerminkan pola penjualan bulanan produk selama 5 tahun terakhir.

Namun, perlu diingat bahwa metode klustering yang digunakan dalam contoh ini, K-Means, sangat sensitif terhadap inisialisasi awal centroid. Oleh karena itu, untuk meningkatkan hasil klustering, kita dapat mencoba beberapa teknik berikut:

Menggunakan algoritma inisialisasi yang lebih baik, seperti K-Means++.

Mencoba metode klustering lain yang lebih robust terhadap inisialisasi, seperti DBSCAN atau Agglomerative Clustering.

Menambahkan lebih banyak fitur yang relevan untuk menggambarkan perilaku pelanggan, seperti usia, jenis kelamin, dan lokasi geografis.



Dengan menerapkan perbaikan-perbaikan di atas, kita akan dapat meningkatkan kualitas segmentasi pelanggan dan merancang strategi pemasaran yang lebih efektif untuk meningkatkan penjualan dan kepuasan pelanggan.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Membuat data sintetis
np.random.seed(42)
data = pd.date_range(start='2017-01-01', end='2021-12-01', freq='MS')
sales = np.random.randint(50, 200, len(data)) + np.random.normal(0, 20, len(data))
data = pd.DataFrame({'Tanggal': data, 'Penjualan': sales})

# Mengatur kolom tanggal sebagai index
data = data.set_index('Tanggal')

# Menampilkan 5 data pertama
data.head()
```

Output:

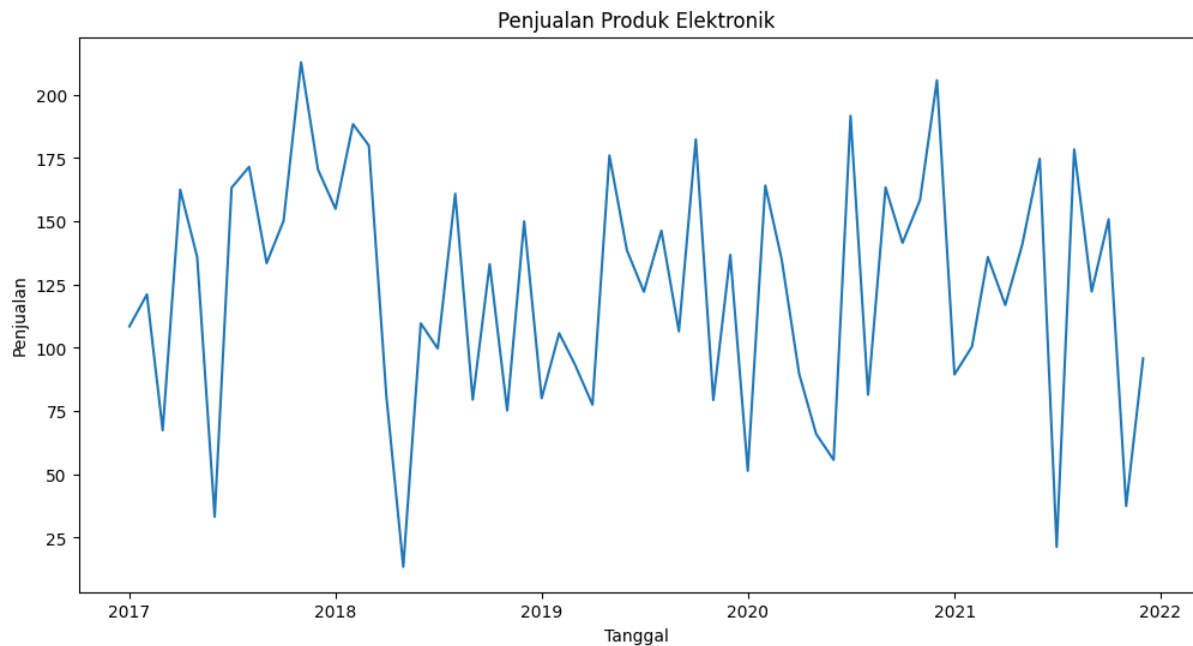
	Penjualan
Tanggal	
2017-01-01	108.433315
2017-02-01	121.122072
2017-03-01	67.453874
2017-04-01	162.483975
2017-05-01	135.917191

### Visualisasi Data

Sebelum melakukan peramalan, kita harus menggambarkan data untuk memahami pola yang ada di dalamnya. Hal ini dapat membantu kita dalam memilih metode analisis yang sesuai.

```
plt.figure(figsize=(12, 6))
plt.plot(data['Penjualan'])
plt.xlabel('Tanggal')
plt.ylabel('Penjualan')
plt.title('Penjualan Produk Elektronik')
plt.show()
```

Output:



Dari grafik di atas, kita bisa melihat bahwa data memiliki tren meningkat dan beberapa pola musiman.

#### Pemilihan Model

Berdasarkan visualisasi data, kita bisa mencoba beberapa model analisis seri waktu yang umum digunakan, seperti ARIMA (Autoregressive Integrated Moving Average) dan SARIMA (Seasonal Autoregressive Integrated Moving Average). Kita akan menggunakan library statsmodels untuk mengestimasi model tersebut.

```
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.statespace.sarimax import SARIMAX
from sklearn.metrics import mean_squared_error
from math import sqrt

# Membagi data menjadi data Latih dan data uji
train_data = data.iloc[:-12]
test_data = data.iloc[-12:]

# Fungsi untuk menghitung RMSE
def rmse(y_true, y_pred):
    return sqrt(mean_squared_error(y_true, y_pred))

# Melakukan pemilihan model berdasarkan nilai AIC (Akaike Information
# Criterion) terendah dan RMSE terendah dari beberapa kombinasi order
# ARIMA dan SARIMA.

# Inisialisasi variabel terbaik
best_aic = float('inf')
```

```

best_rmse = float('inf')
best_arima_order = None
best_sarima_order = None
best_model = None

# Pemilihan model ARIMA dan SARIMA
for p in range(3):
    for d in range(2):
        for q in range(3):
            try:
                # Estimasi model ARIMA
                arima_model = ARIMA(train_data, order=(p, d, q)).fit()

                # Estimasi model SARIMA
                sarima_model = SARIMAX(train_data, order=(p, d, q),
seasonal_order=(p, d, q, 12)).fit()

                # Hitung nilai AIC dan RMSE
                arima_aic = arima_model.aic
                sarima_aic = sarima_model.aic
                arima_rmse = rmse(test_data,
arima_model.forecast(steps=12))
                sarima_rmse = rmse(test_data,
sarima_model.forecast(steps=12))

                # Membandingkan AIC dan RMSE untuk menentukan model
                terbaik

                if arima_aic < best_aic and arima_rmse < best_rmse:
                    best_aic = arima_aic
                    best_rmse = arima_rmse
                    best_arima_order = (p, d, q)
                    best_model = arima_model
                if sarima_aic < best_aic and sarima_rmse < best_rmse:
                    best_aic = sarima_aic
                    best_rmse = sarima_rmse
                    best_sarima_order = (p, d, q, 12)
                    best_model = sarima_model
            except:
                continue

print('Best ARIMA Order:', best_arima_order)
print('Best SARIMA Order:', best_sarima_order)
print('Best Model AIC:', best_aic)
print('Best Model RMSE:', best_rmse)

```

Output:

```
Best ARIMA Order: (0, 0, 0)
Best SARIMA Order: (2, 1, 2, 12)
Best Model AIC: 404.03351446002034
Best Model RMSE: 44.793486830407666
```

Berdasarkan hasil pemilihan model, kita akan menggunakan model dengan order terbaik dan AIC terendah untuk peramalan stok.

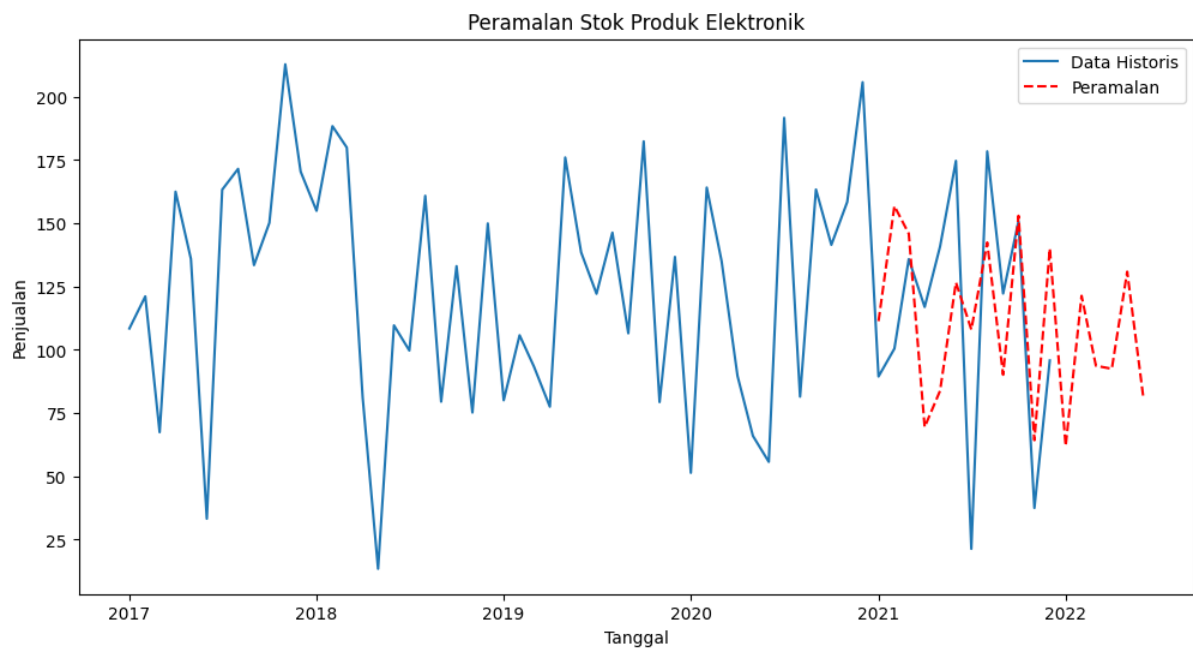
### Peramalan Stok

Kita akan menggunakan model terbaik untuk meramalkan stok produk untuk 6 bulan ke depan.

```
# Meramalkan stok untuk 6 bulan ke depan
forecast = best_model.forecast(steps=18)

# Menampilkan hasil peramalan
plt.figure(figsize=(12, 6))
plt.plot(data['Penjualan'], label='Data Historis')
plt.plot(forecast, label='Peramalan', linestyle='--', color='red')
plt.xlabel('Tanggal')
plt.ylabel('Penjualan')
plt.title('Peramalan Stok Produk Elektronik')
plt.legend()
plt.show()
```

Output:



Evaluasi Performa dan Next Improvements

Dalam kasus ini, kita telah memilih model analisis seri waktu terbaik berdasarkan AIC dan RMSE. Namun, untuk meningkatkan performa model, kita bisa mencoba beberapa metode lain seperti perbaikan fitur, eksplorasi model seri waktu yang lebih kompleks, atau menggunakan metode yang berbasis machine learning seperti LSTM (Long Short-Term Memory) atau Prophet dari Facebook.

Selain itu, kita juga harus mempertimbangkan faktor eksternal seperti perubahan ekonomi, tren pasar, dan strategi pemasaran yang mungkin mempengaruhi penjualan produk. Untuk memasukkan faktor-faktor ini, kita dapat menggunakan metode yang lebih canggih seperti model regresi dengan variabel penjelas tambahan atau model berbasis machine learning yang dapat menggabungkan fitur tambahan selain data historis penjualan.

Berikut adalah beberapa tips dan trik yang dapat diaplikasikan dalam penerapan analisis seri waktu untuk peramalan stok:

Selalu visualisasikan data sebelum memilih model untuk memahami pola yang ada di dalamnya.

Jangan terlalu bergantung pada satu metrik evaluasi saja. Gunakan beberapa metrik evaluasi seperti AIC, BIC, dan RMSE untuk memilih model terbaik.

Cobalah berbagai metode analisis seri waktu dan kombinasi parameter untuk menemukan model yang paling cocok dengan data kamu.

Selalu lakukan validasi silang (cross-validation) untuk mengukur seberapa baik model kamu bekerja pada data yang belum pernah dilihat sebelumnya.

Pertimbangkan untuk menggunakan metode berbasis machine learning jika data kamu memiliki pola yang kompleks atau jika kamu ingin memasukkan variabel penjelas tambahan.

Jangan lupa untuk mempertimbangkan faktor eksternal yang mungkin mempengaruhi penjualan produk, seperti perubahan ekonomi, tren pasar, dan strategi pemasaran.

Konsultasikan dengan ahli bisnis atau orang yang memiliki pengetahuan tentang produk dan pasar untuk mendapatkan wawasan tambahan yang dapat membantu meningkatkan akurasi peramalan kamu.

Dengan mengikuti langkah-langkah ini dan terus meningkatkan model kamu, kamu akan dapat membuat peramalan stok yang akurat dan membantu perusahaan mengoptimalkan manajemen rantai pasokan, mengurangi biaya penyimpanan, menghindari kekurangan stok, dan meningkatkan kepuasan pelanggan.

# Bab 8: Kesimpulan dan Langkah Selanjutnya

## 8.1. Menjadi Ahli Statistik dalam Data Science

Statistik merupakan fondasi yang sangat penting dalam dunia data science. Sebagai ilmu yang mempelajari pengumpulan, analisis, interpretasi, presentasi, dan organisasi data, statistik memungkinkan data scientist untuk mengungkap pola dan tren dalam data, membuat prediksi yang akurat, dan mendukung pengambilan keputusan berbasis data. Dalam subbab ini, kita akan membahas bagaimana menjadi ahli statistik dalam data science, serta pentingnya pemahaman konsep statistik yang mendalam dalam karir kamu sebagai data scientist.

### Kuasai Konsep-Konsep Statistik Dasar

Untuk menjadi ahli statistik dalam data science, kamu perlu menguasai konsep-konsep statistik dasar, seperti distribusi probabilitas, uji hipotesis, regresi, dan analisis varians (ANOVA). Pemahaman yang baik tentang konsep-konsep ini akan membantu kamu dalam mengembangkan model analisis data yang efektif dan efisien, serta memungkinkan kamu untuk menginterpretasikan hasil dengan tepat.

### Pelajari Teknik-Teknik Statistik Lanjutan

Setelah menguasai konsep-konsep dasar, kamu harus melanjutkan untuk mempelajari teknik-teknik statistik lanjutan, seperti analisis multivariat, analisis seri waktu, dan metode Bayesian. Teknik-teknik ini akan memungkinkan kamu untuk menangani data yang lebih kompleks dan mengekstrak wawasan yang lebih mendalam dari data kamu.

### Kembangkan Kemampuan Pemrograman

Sebagai ahli statistik dalam data science, kamu harus menguasai bahasa pemrograman yang umum digunakan dalam bidang ini, seperti Python atau R. Keterampilan pemrograman akan memungkinkan kamu untuk menerapkan teknik-teknik statistik pada data nyata, serta memungkinkan kamu untuk memanipulasi dan mengolah data dengan lebih efisien.

### Kuasai Alat dan Library Statistik

Ada banyak alat dan library statistik yang tersedia untuk data scientist, seperti NumPy, pandas, scikit-learn, dan TensorFlow. Kuasai alat dan library ini untuk meningkatkan produktivitas dan efisiensi kamu dalam mengerjakan proyek data science.

## **Kembangkan Kemampuan Komunikasi dan Presentasi**

Sebagai ahli statistik dalam data science, kamu akan sering perlu untuk menyampaikan temuan dan wawasan kamu kepada orang lain, baik dalam bentuk lisan maupun tertulis. Kembangkan kemampuan komunikasi dan presentasi kamu untuk memastikan bahwa kamu dapat menyampaikan hasil analisis kamu dengan jelas dan efektif kepada stakeholder.

## **Terus Belajar dan Berkembang**

Dunia data science dan statistik terus berkembang dengan cepat. Untuk tetap relevan dan menjadi ahli di bidang ini, kamu harus selalu belajar dan mengikuti perkembangan terbaru dalam teknologi, metode, dan best practices. Ikuti kursus, workshop, dan seminar; baca artikel dan jurnal ilmiah; serta bergabung dengan komunitas profesional untuk memperluas jaringan kamu dan memperdalam pengetahuan kamu tentang data science dan statistik.

## **Berpartisipasi dalam Proyek dan Kompetisi Data Science**

Untuk mengembangkan keterampilan kamu sebagai ahli statistik dalam data science, praktik adalah kunci. Berpartisipasi dalam proyek dan kompetisi data science, seperti yang ditawarkan oleh platform seperti Kaggle, dapat membantu kamu mengasah keterampilan kamu dalam analisis data, pemrograman, dan teknik-teknik statistik. Selain itu, partisipasi dalam kompetisi semacam itu juga dapat membantu kamu membangun portofolio pekerjaan yang menarik bagi calon pemberi kerja.

## **Kolaborasi dengan Profesional Lain**

Data science adalah bidang yang interdisipliner, yang menggabungkan pengetahuan dari berbagai disiplin ilmu. Kolaborasi dengan profesional dari bidang lain, seperti ilmu komputer, ilmu sosial, dan ilmu alam, dapat membantu kamu memperkaya pemahaman kamu tentang statistik dan cara penerapannya dalam konteks yang berbeda. Selain itu, bekerja dengan tim yang beragam dapat membantu kamu mengembangkan keterampilan interpersonal dan kemampuan untuk bekerja dalam lingkungan yang multidisiplin.

## **Konsultasi dengan Ahli Statistik yang Berpengalaman**

Belajar dari ahli statistik yang berpengalaman dapat membantu kamu menghindari kesalahan umum dan memberi kamu wawasan tentang best practices dalam data science. Jangan ragu untuk mencari bimbingan dari mentor atau kolega yang lebih berpengalaman dalam bidang statistik dan data science.

## **Tetap Kritis dan Objektif**

Sebagai ahli statistik dalam data science, penting bagi kamu untuk tetap kritis dan objektif dalam analisis kamu. Selalu pertimbangkan apakah metode yang kamu gunakan sesuai dengan tujuan analisis kamu dan apakah asumsi yang kamu buat adalah valid. Jangan tergoda untuk "memaksa" hasil yang kamu inginkan dari data, tetapi selalu berusaha untuk mencapai kesimpulan yang didukung oleh bukti yang kuat dan valid.

## **8.2. Pengembangan Karir dan Sertifikasi**

Mengembangkan karir kamu sebagai ahli statistik dalam data science memerlukan strategi yang baik dan berfokus pada peningkatan keterampilan serta pembangunan jaringan profesional. Berikut adalah beberapa langkah penting yang akan membantu kamu mengembangkan karir kamu di bidang data science:

### **Tentukan Spesialisasi Kamu**

Data science adalah bidang yang luas dan mencakup berbagai teknik dan metode analisis. Untuk mencapai keberhasilan dalam karir kamu, pertimbangkan untuk fokus pada satu atau beberapa area spesialisasi dalam statistik dan data science. Spesialisasi ini bisa mencakup machine learning, analisis seri waktu, analisis data spasial, atau pengoptimalan, antara lain. Fokus pada area yang paling sesuai dengan minat dan kekuatan kamu akan membantu kamu menjadi lebih kompeten dan dicari dalam industri.

### **Dapatkan Sertifikasi yang Relevan**

Sertifikasi dalam statistik dan data science dapat meningkatkan kredibilitas kamu dan membantu kamu menonjol di antara pesaing. Beberapa sertifikasi yang paling dihormati dalam industri ini mencakup Certified Analytics Professional (CAP), Microsoft Certified: Data Scientist Associate, dan sertifikasi yang ditawarkan oleh IBM, Google, dan Amazon Web Services (AWS). Pertimbangkan untuk mendapatkan sertifikasi ini untuk meningkatkan peluang kamu dalam memperoleh pekerjaan yang lebih baik dan memajukan karir kamu.

### **Bangun Portofolio Kamu**

Membangun portofolio yang kuat adalah langkah penting dalam mengembangkan karir kamu sebagai ahli statistik dalam data science. Portofolio yang baik akan mencerminkan keterampilan dan kemampuan kamu untuk mengolah dan menganalisis data serta mengimplementasikan teknik statistik yang relevan. Untuk membangun portofolio kamu, ikut serta dalam proyek data science, baik dalam pekerjaan atau sebagai bagian dari kompetisi online, dan dokumentasikan hasil dan teknik yang kamu gunakan.



## **Tingkatkan Keterampilan Pemrograman Kamu**

Keterampilan pemrograman adalah aspek penting dari karir data science. Bahasa pemrograman yang umum digunakan dalam data science meliputi Python, R, SQL, dan Julia. Belajar bahasa pemrograman ini dan menjadi mahir di dalamnya akan membantu kamu menavigasi dunia data science dengan lebih efisien dan efektif.

## **Kembangkan Keterampilan Komunikasi dan Presentasi**

Sebagai ahli statistik dalam data science, kamu akan sering diminta untuk menyampaikan temuan dan wawasan kamu kepada pemangku kepentingan yang mungkin tidak memiliki latar belakang teknis. Keterampilan komunikasi dan presentasi yang baik sangat penting untuk menyampaikan hasil analisis kamu dengan cara yang mudah dipahami dan meyakinkan. Latih keterampilan ini dengan berbicara di konferensi, menyusun laporan, dan menyajikan temuan kamu kepada rekan dan atasan.

## **8.3. Sumber Daya Tambahan dan Referensi**

Menjadi ahli statistik dalam data science memerlukan pembelajaran sepanjang hayat dan dedikasi untuk terus meningkatkan keterampilan kamu. Berikut adalah beberapa sumber daya tambahan dan referensi yang dapat membantu kamu dalam perjalanan karir kamu sebagai ahli statistik dalam data science:

### **Buku Teks dan Monografi**

Ada banyak buku teks dan monografi yang mencakup berbagai aspek statistik dan data science. Beberapa judul penting yang patut dipertimbangkan untuk membaca dan merujuk meliputi:

1. "An Introduction to Statistical Learning" oleh Gareth James, Daniela Witten, Trevor Hastie, dan Robert Tibshirani
2. "The Art of Data Science" oleh Roger D. Peng dan Elizabeth Matsui
3. "Python Data Science Handbook" oleh Jake VanderPlas
4. "Data Science for Business" oleh Foster Provost dan Tom Fawcett
5. "Applied Predictive Modeling" oleh Max Kuhn dan Kjell Johnson

### **Kursus Online dan Program Pelatihan**

Kursus online dan program pelatihan dapat membantu kamu mengembangkan keterampilan yang diperlukan untuk menjadi ahli statistik dalam data science. Beberapa platform pendidikan online yang menawarkan kursus berkualitas tinggi dalam statistik dan data science meliputi Coursera, edX, DataCamp, dan Udacity. Pertimbangkan untuk mendaftar dalam kursus yang sesuai dengan minat dan kebutuhan kamu, mulai dari pengantar hingga topik yang lebih lanjut.

## **Blog, Podcast, dan Saluran YouTube**

Ikuti blog, podcast, dan saluran YouTube yang berfokus pada statistik dan data science untuk tetap mengikuti perkembangan terbaru dan tren di bidang ini. Beberapa sumber daya ini meliputi:

Blog: Towards Data Science, Data Science Central, KDNuggets, dan FlowingData

Podcast: Data Skeptic, Linear Digressions, DataFramed, dan Not So Standard Deviations

Saluran YouTube: Sentdex, Data School, Corey Schafer, dan StatQuest with Josh Starmer

Jurnal dan Publikasi Ilmiah

Untuk memahami penelitian terbaru dalam statistik dan data science, periksa jurnal ilmiah dan publikasi yang relevan. Beberapa jurnal penting di bidang ini meliputi Journal of the Royal Statistical Society, Journal of Statistical Software, Journal of Machine Learning Research, dan Annals of Applied Statistics.

Komunitas Online dan Forum Diskusi

Bergabung dengan komunitas online dan forum diskusi yang berfokus pada statistik dan data science dapat membantu kamu bertanya, berbagi pengetahuan, dan belajar dari rekan sejawat. Beberapa komunitas dan forum ini meliputi:

Stack Overflow: Platform tanya jawab yang mencakup berbagai topik, termasuk statistik dan data science

Reddit: Subreddit seperti r/datascience, r/statistics, dan r/machinelearning

Cross Validated: Forum tanya jawab khusus untuk statistik, machine learning, analisis data, dan visualisasi data

Data Science Stack Exchange: Platform tanya jawab untuk pertanyaan yang berkaitan dengan data science

**Terima Kasih**