

Hoofdstuk 1: Basis en Datasets

Functies en Methodes (Datasets)

Functie/Methode	Omschrijving	Voorbeeld	Output (voorbeeld)
<code>read_csv()</code>	Leest een CSV-bestand in een DataFrame	<code>titanic = pd.read_csv('https://raw.githubusercontent.com/DataRepo2019/Data-files/master/titanic.csv')</code>	DataFrame met Titanic-data geladen
<code>head()</code>	Toont de eerste rijen van een DataFrame	<code>titanic.head()</code>	Eerste 5 rijen van DataFrame
<code>tail()</code>	Toont de laatste rijen van een DataFrame	<code>print(not_null_df.tail(3))</code>	Laatste 3 rijen
<code>info()</code>	Toont informatie over kolommen en datatypes	<code>titanic.info()</code>	Overzicht van kolomnamen, niet-null counts, datatypes
<code>describe()</code>	Geeft statistische samenvatting	<code>print(titanic.Survived.describe())</code>	Count, mean, std, min, max enz.
<code>count()</code>	Telt niet-NA waarden per kolom	<code>print(titanic.count())</code>	Aantal niet-NA per kolom
<code>sum()</code>	Som van waarden	<code>print(not_null_df.sum())</code>	Som per kolom of rij
<code>mean()</code>	Gemiddelde van waarden	<code>avg_age = titanic['Age'].mean()</code>	Eén getal (gemiddelde leeftijd)
<code>unique()</code>	Unieke waarden in een kolom	<code>print(titanic.Embarked.unique())</code>	Array van unieke waarden
<code>value_counts()</code>	Telt hoe vaak unieke waarden voorkomen	<code>print(titanic.dtypes.value_counts())</code>	Tabel met telling per uniek type
<code>astype()</code>	Zet kolom om naar een ander datatype	<code>titanic.Survived = titanic.Survived.astype('category')</code>	Kolomtype aangepast
<code>drop()</code>	Verwijdert kolommen of rijen	<code>titanic.drop("PassengerId", axis="columns")</code>	DataFrame zonder opgegeven kolom

Functie/Methode	Omschrijving	Voorbeeld	Output (voorbeeld)
<code>dropna()</code>	Verwijdert rijen met NA-waarden	<code>cleaned = titanic.dropna()</code>	DataFrame zonder NA-rijen
<code>fillna()</code>	Vult NA-waarden op	<code>titanic = titanic.fillna(value={'Age' : avg_age})</code>	DataFrame met ingevulde waarden
<code>set_index()</code>	Zet een kolom als index	<code>titanic.set_index(['PassengerId'])</code>	DataFrame met nieuwe index
<code>map()</code>	Past een functie toe op een Series	<code>titanic['Sex'] = titanic['Sex'].map(mapping_dict)</code>	Kolom met aangepaste waarden
<code>query()</code>	Filtret data met een query-string	<code>titanic.query("(Sex=='male') and (Age < 18)")</code>	Gefilterde DataFrame
<code>round()</code>	Rondt waarden af	<code>print(f"(Rounded) Average age of passengers: {round(avg_age)}")</code>	Afgerond getal
<code>notnull()</code>	Controleert of waarden niet null zijn	<code>not_null_df = titanic.notnull()</code>	DataFrame met True/False waarden

Extra (niet direct pandas maar gezien in notebook)

Functie	Omschrijving	Voorbeeld	Output (voorbeeld)
<code>print()</code>	Print tekst naar console	<code>print(f"Number of rows: {len(titanic)}")</code>	"Number of rows: 891" (bijv.)
<code>len()</code>	Geeft lengte van object	<code>print(f"Number of rows: {len(titanic)}")</code>	Getal met aantal rijen
<code>and</code>	Logische EN	<code>titanic.query("(Sex=='male') and (Age < 18)")</code>	Gefilterde DataFrame

Visualisaties (Grafieken)

Naam Grafiek	Functie	Omschrijving	Voorbeeld	Output (voorbeeld)
Countplot	<code>countplot()</code>	Toont telling van categorieën	<code>sns.countplot(data=titanic, x='Embarked');</code>	Balkgrafiek met tellingen per categorie

Hoofdstuk 2: Univariate Analysis

Functies en Methodes

Functie/Methode	Omschrijving	Voorbeeld	Output (voorbeeld)
<code>load_dataset()</code>	Laadt voorbeelddataset (Seaborn)	<code>tips = sns.load_dataset("tips")</code>	DataFrame met tips-data
<code>head()</code>	Toont eerste rijen van DataFrame	<code>tips.head()</code>	Eerste 5 rijen
<code>mean()</code>	Gemiddelde	<code>print(f"Mean: {tips.tip.mean()}")</code>	Gemiddelde waarde

Functie/Methode	Omschrijving	Voorbeeld	Output (voorbeeld)
median()	Mediaan	print(f"Median: {tips.tip.median()}")	Mediaanwaarde
mode()	Modus	tips.mode()	Waarde(n) die het vaakst voorkomen
min()	Minimumwaarde	print(f"Minimum: {tips.tip.min()}")	Kleinste waarde
max()	Maximumwaarde	print(f"Maximum: {tips.tip.max()}")	Grootste waarde
var()	Variantie	print(f"Variance: {tips.tip.var()}")	Variantie (n-1 in noemer)
std()	Standaardafwijking	print(f"Standard deviation: {tips.tip.std()}")	Std-dev
quantile()	Percentielen	print(f"Percentiles {percentiles}\n{tips.tip.quantile(percentiles)}")	Waarden op percentielen
describe()	Statistische samenvatting	tips.day.describe()	Count, mean, std, min, max enz.
kurtosis()	Scheefheid (spitsheid)	print(f"Kurtosis: {tips.tip.kurtosis()}")	Kurtosis-waarde
skew()	Scheefheid	print(f"Skewness: {tips.tip.skew()}")	Skewness-waarde
iqr()	Interkwartielafstand	print(f"Inter Quartile Range: {stats.iqr(tips.tip)}")	IQR-waarde
array()	Maakt numpy-array	a = np.array([4, 8, 6, 5, 3, 2, 8, 9, 2, 5])	Numpy-array
arange()	Maakt range-array	population = np.arange(0, 101)	Array van 0 tot 100
choice()	Willekeurige selectie uit array	sample = np.random.choice(population, size=sample_size)	Array met random waarden
sum()	Som van waarden	mean = sum(x) / n	Totale som
len()	Lengte van object	n = len(x)	Getal met lengte
sqrt()	Vierkantswortel	return np.sqrt(pop_var(x))	Wortelwaarde

Visualisaties (Grafieken)

Naam Grafiek	Functie	Omschrijving	Voorbeeld	Output (voorbeeld)
Displot	displot()	Histogram met verdelingscurve	sns.displot(data=tips, x='tip');	Histogram
KDE Plot	kdeplot()	Kernel Density Estimate plot	sns.kdeplot(data=tips, x='tip');	Gladde verdelingslijn
Boxplot	boxplot()	Boxplot	sns.boxplot(data=tips, x='tip');	Boxplot met mediaan, IQR enz.
Violinplot	violinplot()	Violinplot	sns.violinplot(data=tips, x='tip');	Violinvormige plot
Catplot (count)	catplot(kind='count')	Categorische tellingen plot	sns.catplot(data=tips, x='day', kind='count');	Balkgrafiek met tellingen

Hoofdstuk 3: Probability

Functies en Methodes

Functie/Methode	Omschrijving	Voorbeeld	Output (voorbeeld)
-----------------	--------------	-----------	--------------------

Functie/Methode	Omschrijving	Voorbeeld	Output (voorbeeld)
<code>arange()</code>	Maakt een array met stappen	<code>x = np.arange(4)</code>	Array [0,1,2,3]
<code>linspace()</code>	Verdeelt bereik in evenredige stappen	<code>x = np.linspace(mu - 4 * sigma, mu + 4 * sigma, num=201)</code>	Array met 201 waarden
<code>randint()</code>	Genereert willekeurige gehele getallen	<code>random.randint(1,6)</code>	Getal tussen 1-6
<code>normal()</code>	Genereert normale verdeling	<code>observations = np.random.normal(loc=m, scale=s, size=n)</code>	Array met normaal verdeelde waarden
<code>full()</code>	Array vullen met constante waarde	<code>y = np.full(6, 1/6)</code>	Array [1/6,1/6,...]
<code>append()</code>	Voegt waarde toe aan lijst	<code>l_game_1.append(number_of_times_won / index)</code>	Gewijzigde lijst
<code>range()</code>	Genereert reeks getallen	<code>for index in range(1, number_of_games + 1):</code>	1 tot number_of_games
<code>pdf()</code>	Probability density function (scipy.stats)	<code>y = stats.norm.pdf(x, mu, sigma)</code>	Waarden van de kansdichtheid
<code>cdf()</code>	Cumulative density function (scipy.stats)	<code>stats.norm.cdf(1.62, loc=0, scale=1)</code>	Kanswaarde
<code>sf()</code>	Survival function (1 - CDF) (scipy.stats)	<code>stats.norm.sf(1.62, loc=0, scale=1)</code>	Complementaire kanswaarde
<code>isf()</code>	Inverse survival function (scipy.stats)	<code>stats.norm.isf(0.05, loc=0, scale=1)</code>	Waarde bij gegeven overlevingskans
<code>interval()</code>	Bereken betrouwbaarheidsinterval	<code>stats.norm.interval(confidence, loc=m, scale=s/vn)</code>	Tuple met (lower, upper) grens
<code>sqrt()</code>	Vierkantswortel	<code>lo = m - z * s / np.sqrt(n)</code>	Wortelwaarde
<code>subplots()</code>	Maakt figuur en subplot	<code>fig, tplot = plt.subplots(1, 1)</code>	Lege figuur met as
<code>plot()</code>	Lijnplot tekenen	<code>tplot.plot(x, stats.norm.pdf(x, 0, 1))</code>	Lijngrafiek
<code>legend()</code>	Voeg legenda toe aan plot	<code>tplot.legend(loc='best')</code>	Legenda op plot

Visualisaties (Grafieken)

Naam Grafiek/Plot	Functie	Omschrijving	Voorbeeld	Output (voorbeeld)
Lijnplot (seaborn)	<code>lineplot()</code>	Lijnplot voor trends	<code>sns.lineplot(ax=ax, x='number of games', y='game 1 die', data=data)</code>	Lijngrafiek
Staafdiagram (bar)	<code>bar()</code>	Staafdiagram	<code>ax.bar(x, y, 0.35)</code>	Staafdiagram
Histogram	<code>histplot()</code>	Histogram met KDE	<code>sns.histplot(observations, kde=True)</code>	Histogram met verdelingslijn
Vulgebied	<code>fill_between()</code>	Kleurt gebied onder curve	<code>plt.fill_between(dist_x, 0, dist_y, where=dist_x <= x, color='lightblue')</code>	Ingekleurde grafiek
Verticale lijn	<code>axvline()</code>	Verticale lijn in grafiek	<code>plt.axvline(75, color='green')</code>	Verticale lijn

Naam Grafiek/Plot	Functie	Omschrijving	Voorbeeld	Output (voorbeeld)
Subplots	<code>subplots()</code>	Meerdere grafieken naast elkaar	<code>f, ax = plt.subplots(1, 1)</code>	Lege subplot
Figure	<code>figure()</code>	Nieuwe figuur starten	<code>fig = plt.figure()</code>	Lege figuur
Add Axes	<code>add_axes()</code>	Handmatig assen toevoegen	<code>ax = fig.add_axes([0,0,1,1])</code>	Assen toegevoegd
Show Plot	<code>show()</code>	Toon de plot	<code>plt.show()</code>	Toont plot op scherm
Plot (lijn)	<code>plot()</code>	Lijnplot (matplotlib)	<code>plt.plot(x, stats.norm.pdf(x, 0, 1))</code>	Lijngrafiek
X-as label	<code>xlabel()</code>	Zet label op x-as	<code>plt.xlabel('number of games')</code>	Label op x-as
Y-as label	<code>ylabel()</code>	Zet label op y-as	<code>plt.ylabel('win percentage')</code>	Label op y-as
Set X Label	<code>set_xlabel()</code>	Zet label op x-as via Axes	<code>ax.set_xlabel('x')</code>	Label op x-as
Set Y Label	<code>set_ylabel()</code>	Zet label op y-as via Axes	<code>ax.set_ylabel('P(X = x)')</code>	Label op y-as
Titel toevoegen	<code>title()</code>	Voeg titel toe aan plot	<code>plt.title('Standard Normal Distribution')</code>	Titel boven plot
X-limieten	<code>xlim()</code>	Beperk x-as bereik	<code>plt.xlim((0,20))</code>	X-as limiet
Y-limieten	<code>ylim()</code>	Beperk y-as bereik	<code>plt.ylim((0,1))</code>	Y-as limiet
Legenda toevoegen	<code>legend()</code>	Voeg legenda toe	<code>plt.legend()</code>	Legenda op plot

Symbolische Notaties

Symbol	Naam	Betekenis
\cap	Doorsnede	A én B (snijpunt van gebeurtenissen)
\cup	Unie	A of B (alles wat in A of B zit)
A^c	Complement	Alles wat niet in A zit
		Voorwaardelijke kans
μ (mu)	Gemiddelde (verwachtingswaarde)	Gemiddelde van verdeling
σ (sigma)	Standaardafwijking	Spreiding rond het gemiddelde
α (alpha)	Significantieniveau	Kans buiten betrouwbaarheidsinterval (bijv. 0,05)

Kansregels en Formules

Begrip	Betekenis	Formule
Kans	Kans dat gebeurtenis A gebeurt	$P(A)$
Complementregel	Kans dat A niet gebeurt	$P(A^c) = 1 - P(A)$
Productregel (onafhankelijk)	Kans op A én B als onafhankelijk	$P(A \cap B) = P(A) * P(B)$
Productregel (afhankelijk)	Kans op A én B als afhankelijk	$P(A \cap B) = P(A) * P(B)$

Begrip	Betekenis	Formule
Somregel (disjunct)	Kans op A of B (disjunct)	$P(A \cup B) = P(A) + P(B)$
Somregel (niet-disjunct)	Kans op A of B	$P(A \cup B) = P(A) + P(B) - P(A \cap B)$
CDF (cumulatief)	Kans dat X kleiner of gelijk aan x	$P(X \leq x) = \text{CDF}(x)$
SF (survival)	Kans dat X groter is dan x	$P(X > x) = \text{SF}(x) = 1 - \text{CDF}(x)$
Z-score (kritieke waarde)	Kritieke waarde uit standaardnormale verdeling	$z = \text{isf}(\alpha/2)$
Standaardfout (SE)	Standaardafwijking van steekproefgemiddelde	$SE = s / \sqrt{n}$
Bereken z-score voor interval	Stap 1 voor betrouwbaarheidsinterval	$z = \text{stats.norm.isf}(\alpha/2)$
Bereken SE voor interval	Stap 2 voor betrouwbaarheidsinterval	$SE = s / \sqrt{n}$
Bereken onder- en bovengrens	Stap 3 voor betrouwbaarheidsinterval	$lo = m - z \cdot SE; hi = m + z \cdot SE$
α (alpha)	Kans buiten interval (significantieniveau)	$\alpha = 1 - \text{confidence}$

Formule (handmatig):

- Bereken z-score: $z = \text{stats.norm.isf}(\alpha/2)$
- Bereken standaardfout: $SE = s / \sqrt{n}$
- Ondergrens: $lo = m - z \cdot SE$
- Bovengrens: $hi = m + z \cdot SE$

Uitleg: (Standaard)normale verdeling

De **standaardnormale verdeling** is een speciale normale verdeling met gemiddeld $\mu = 0$ en standaardafwijking $\sigma = 1$. Het is een klokvormige symmetrische curve die wordt gebruikt om veel natuurlijke verschijnselen te modelleren.

Eigenschappen:

- Symmetrisch rond het gemiddelde $\mu = 0$
- 68% van de waarden ligt binnen 1 standaardafwijking (σ)
- 95% van de waarden ligt binnen 2 standaardafwijkingen (σ)
- 99,7% van de waarden ligt binnen 3 standaardafwijkingen (σ)

Python-voorbeeld:

```
x = np.linspace(-4, 4, 100)
y = stats.norm.pdf(x, 0, 1)
plt.plot(x, y)
plt.title("Standaardnormale verdeling")
plt.show()
```

Uitleg: Wat is een z-score?

Een **z-score** geeft aan hoeveel standaardafwijkingen een waarde van het gemiddelde af ligt.

Formule:

$$z = (x - \mu) / \sigma$$

Wat doet het? Het zet elke waarde uit een normale verdeling om naar de standaardnormale schaal, zodat je kanswaarden kunt opzoeken met bijvoorbeeld `cdf()`, `sf()` of `isf()`.

Voorbeeld: Stel je hebt een waarde $x = 80$, met een gemiddelde $\mu = 70$ en een standaardafwijking $\sigma = 5$, dan geldt:

$$z = (80 - 70) / 5 = 10 / 5 = 2$$

Dit betekent dat de waarde 2 standaardafwijkingen boven het gemiddelde ligt.

Wat is α (alpha)? De **α -waarde** (significantieniveau) geeft aan hoeveel kans je buiten een betrouwbaarheidsinterval laat. Bijvoorbeeld:

- Bij 95% betrouwbaarheidsinterval: $\alpha = 1 - 0,95 = 0,05$
- Bij 99% betrouwbaarheidsinterval: $\alpha = 1 - 0,99 = 0,01$

Omdat de interval symmetrisch verdeeld is, gebruik je vaak **$\alpha/2$** voor de bovenste of onderste kant.

Python-voorbeeld:

```
confidence = 0.95
m = 50
s = 10
n = 100
z = stats.norm.isf((1 - confidence) / 2)
lo = m - z * s / np.sqrt(n)
hi = m + z * s / np.sqrt(n)
print(f"Confidence interval: [{lo:.2f}, {hi:.2f}]")
```

Hypothesis Testing: Overzicht en Voorwaarden

Wanneer gebruik je welke test?

- **z-test** → als je een grote steekproef hebt ($n \geq 30$), de standaardafwijking σ van de populatie kent, en de data normaal verdeeld zijn.
- **t-test** → als je een kleine steekproef hebt ($n < 30$) of σ onbekend is.

Stappenplan Hypothesis Testing

1 Formuleer de hypotheses

- H_0 : nulhypothese (bijvoorbeeld: er is geen verschil)
- H_1 : alternatieve hypothese (bijvoorbeeld: er is wél verschil)

2 Kies significantieniveau

- Meestal $\alpha = 0,05$ (5%)

3 Bereken de teststatistiek

- $z = (\bar{x} - \mu) / (\sigma / \sqrt{n}) \rightarrow$ bij z-test
- $t = (\bar{x} - \mu) / (s / \sqrt{n}) \rightarrow$ bij t-test

4 Bepaal het kritieke gebied

- Right-tailed: kijk naar bovenkant ($H_1: \mu > \mu_0$)
- Left-tailed: kijk naar onderkant ($H_1: \mu < \mu_0$)
- Two-tailed: kijk naar beide kanten ($H_1: \mu \neq \mu_0$)

5 Maak een beslissing

- Als p-waarde $\leq \alpha \rightarrow$ verwerp H_0
- Als p-waarde $> \alpha \rightarrow$ behoud H_0

Python formules voor variabelen right tailed z test

Symbool	uitleg	Python formule
n	grootte van steekproef	<code>len(array)</code>
μ (mu)	gemiddelde (niet van steekproef)	/
m_sample	gemiddelde (steekproef)	<code>np.mean(array)</code>
p	/	<code>stats.norm.sf(m_sample, loc=mu, scale=sigma/np.sqrt(n))</code>

Symbool	uitleg	Python formule
g	/	<code>stats.norm.isf(1-alpha, loc=mu, scale=(sigma/ np.sqrt(n)))</code>

Python formules voor variabelen left tailed z test

Symbool	uitleg	Python formule
n	grootte van steekproef	<code>len(array)</code>
μ (mu)	gemiddelde (niet van steekproef)	/
m_sample	gemiddelde (steekproef)	<code>np.mean(array)</code>
p	/	<code>stats.norm.cdf(m_sample, loc=mu, scale=sigma/np.sqrt(n))</code>
g	/	<code>stats.norm.isf(1-alpha, loc=mu, scale=(sigma/ np.sqrt(n)))</code>

Python formules voor variabelen two-tailed z test

Symbool	uitleg	Python formule
n	grootte van steekproef	<code>len(array)</code>
μ (mu)	gemiddelde (niet van steekproef)	/
m_sample	gemiddelde (steekproef)	<code>np.mean(array)</code>
p	/	<code>stats.norm.sf(m_sample, loc=mu, scale=sigma/np.sqrt(n))</code>
g1	linkerkant	<code>mu - stats.norm.isf(1-alpha, loc=mu, scale=(sigma/ np.sqrt(n)))</code>
g2	rechterkant	<code>mu + stats.norm.isf(1-alpha, loc=mu, scale=(sigma/ np.sqrt(n)))</code>

Voorbeeld: Right-tailed z-test

```
n = 50
x_bar = 105
mu = 100
sigma = 10
z = (x_bar - mu) / (sigma / np.sqrt(n))
p_value = stats.norm.sf(z)
```

Voorbeeld: Left-tailed z-test

```
n = 50
x_bar = 95
mu = 100
sigma = 10
z = (x_bar - mu) / (sigma / np.sqrt(n))
p_value = stats.norm.cdf(z)
```

Voorbeeld: Two-tailed z-test

```
n = 50
x_bar = 105
mu = 100
sigma = 10
z = (x_bar - mu) / (sigma / np.sqrt(n))
p_value = 2 * stats.norm.sf(abs(z))
```


Voorbeeld: Right-tailed t-test

```
n = 15
x_bar = 105
mu = 100
s = 10
t = (x_bar - mu) / (s / np.sqrt(n))
p_value = stats.t.sf(t, df=n-1)
```

Voorbeeld: Left-tailed t-test

```
n = 15
x_bar = 95
mu = 100
s = 10
t = (x_bar - mu) / (s / np.sqrt(n))
p_value = stats.t.cdf(t, df=n-1)
```

Voorbeeld: Two-tailed t-test (met `ttest_1samp`)

```
sample = np.array([98, 102, 100, 105, 97])
mu = 100
t_stat, p_value = stats.ttest_1samp(sample, popmean=mu)
```

Symbolen en Formules

Symbool	Betekenis
H_0	Nulhypothese
H_1	Alternatieve hypothese
α	Significantieniveau (bijv. 0,05)
p	p-waarde, kans op verkregen resultaat onder H_0
z	z-score bij z-test
t	t-score bij t-test
μ	Populatiegemiddelde
\bar{x}	Steekproefgemiddelde
σ	Populatiestandaardafwijking
s	Steekproefstandaardafwijking
n	Steekproefgrootte

Wat is een Type I en Type II fout?

- **Type I fout (α):** je verworpt H_0 terwijl die waar is → vals alarm.
- **Type II fout (β):** je behoudt H_0 terwijl H_1 waar is → je mist een effect.

Gevolgen:

- Type I fout: je denkt dat er effect is, maar dat is er niet.
- Type II fout: er is wél effect, maar je ziet het niet.

Hoofdstuk 4: Chi-Squared & Associaties

Welke test gebruik je wanneer?

Onafhankelijke Variabele (Independent)	Afhankelijke Variabele (Dependent)	Test/Metric
Kwalitatief	Kwalitatief	χ^2 -test, Cramér's V
Kwalitatief	Kwantitatief	two-sample t-test, Cohen's d
Kwantitatief	Kwantitatief	Regressie, correlatie

Uitleg per test

- **χ^2 -test:** Verband tussen twee categorische variabelen.
- **Goodness-of-Fit test:** Vergelijkt verdeling van één categorische variabele met een verwachte verdeling.
- **Cramér's V:** Sterkte van het verband tussen categorische variabelen.
- **Two-sample t-test:** Vergelijkt gemiddelden tussen twee onafhankelijke groepen.
- **Gepaalde t-test:** Vergelijkt gemiddelden tussen twee gekoppelde metingen (voor/na).
- **Cohen's d:** Effectgrootte van verschil tussen twee groepen.
- **Regressie (least squares):** Voorspelt waarden van een afhankelijke kwantitatieve variabele uit een onafhankelijke kwantitatieve variabele.
- **Covariantie:** Meet hoe twee kwantitatieve variabelen samen variëren.
- **Correlatie (Pearson):** Meet de sterkte en richting van een lineair verband tussen twee kwantitatieve variabelen.
- **Coëfficiënt van determinatie (R^2):** Geeft aan hoeveel van de variantie verklaard wordt door het model.

Wat doet de χ^2 -test?

De **χ^2 -test** controleert of er een verband is tussen twee categorische variabelen door de geobserveerde frequenties in een kruistabel te vergelijken met de verwachte frequenties als er géén verband zou zijn.

Wat doet de Goodness-of-Fit test?

De **Goodness-of-Fit test** controleert of de verdeling van één categorische variabele overeenkomt met een verwachte verdeling (bijvoorbeeld een uniforme verdeling of een andere theoretische verdeling).

Beste grafieken (met voorbeeld)

Grafiektype	Omschrijving	Voorbeeldcode
Mozaïekdiagram	Visualiseert frequenties tussen categorieën	<code>mosaic(data, ['col1', 'col2'])</code>
Gestapelde staafdiagram (met percentages)	Vergelijkt proporties per categorie	<code>pd.crosstab(df['col1'], df['col2'], normalize='index').plot(kind='bar', stacked=True)</code>

Functies voor χ^2 , Goodness-of-Fit en Cramér's V

Functie/Methode	Omschrijving	Voorbeeld
<code>pd.crosstab()</code>	Maakt een kruistabel (contingency table)	<code>table = pd.crosstab(df['col1'], df['col2'])</code>
<code>chi2_contingency()</code>	Voert de χ^2 -test uit op een kruistabel	<code>chi2, p, dof, expected = stats.chi2_contingency(table)</code>
<code>chisquare()</code>	Voert de goodness-of-fit test uit	<code>chi2, p = stats.chisquare(observed, f_exp=expected)</code>
Handmatige Cramér's V	Berekent Cramér's V	zie functie hieronder

Handmatige berekening van Cramér's V

```
dof = min(observed.shape) - 1
cramers_v = np.sqrt(chi_squared / (n * dof))
```

```
print(cramers_v)
```

Cramér's V	Interpretation
0	No association
0.1	Weak association
0.25	Moderate association
0.50	Strong association
0.75	Very strong association
1	Complete association

Stappenplan χ^2 -test (contingentie)

1 Formuleer de hypotheses

- H_0 : Er is geen verband tussen de variabelen.
- H_1 : Er is wel een verband tussen de variabelen.

2 Maak een kruistabel

- Gebruik bijvoorbeeld `pd.crosstab()`.

3 Voer de test uit

- Gebruik `stats.chi2_contingency()`.

4 Bekijk de resultaten

- χ^2 -waarde, p-waarde, vrijheidsgraden, expected counts.

5 Maak een beslissing

- Als $p \leq \alpha \rightarrow$ verwerp H_0 .
- Als $p > \alpha \rightarrow$ behoud H_0 .

Stappenplan Goodness-of-Fit test

1 Formuleer de hypotheses

- H_0 : De geobserveerde verdeling komt overeen met de verwachte.
- H_1 : De verdeling wijkt af van de verwachte.

2 Bereken de teststatistiek

- Gebruik `stats.chisquare(observed, f_exp=expected)`.

3 Maak een beslissing

- Als $p \leq \alpha \rightarrow$ verwerp H_0 .
- Als $p > \alpha \rightarrow$ behoud H_0 .

Voorbeeld χ^2 Goodness-of-Fit test in Python

```
import numpy as np
from scipy import stats

observed = np.array([30, 50, 20])
expected = np.array([33.3, 33.3, 33.3])
chi2, p = stats.chisquare(observed, f_exp=expected)
print(f"Chi²: {chi2:.2f}, p-waarde: {p:.4f}")
```

Symbolen en Formules

Symbool	Betekenis
χ^2	Chi-kwadraat statistiek
p	p-waarde, kans onder nulhypothese
dof	Degrees of freedom (vrijheidsgraden)
Cramér's V	Sterkte van het verband tussen variabelen

Hoofdstuk 5: t-tests & Effectgrootte

Beste grafieken (met voorbeeld)

Grafiektype	Omschrijving	Voorbeeldcode
Density plot	Toont verdeling van een continue variabele	<code>sns.kdeplot(data, x='value', hue='group')</code>
Staafdiagram met foutbalken	Vergelijkt gemiddelden met onzekerheid (error bars)	<code>sns.barplot(data=df, x='group', y='score', ci='sd')</code>

Functies voor t-tests en Cohen's d

Functie/Methode	Omschrijving	Voorbeeld
<code>ttest_ind()</code>	t-test voor twee onafhankelijke steekproeven	<code>stats.ttest_ind(group1, group2)</code>
<code>ttest_rel()</code>	t-test voor gepaarde steekproeven	<code>stats.ttest_rel(before, after)</code>
<code>cohen_d()</code> (zelf gedefinieerd)	Berekent Cohen's d (effectgrootte) handmatig	Zie functie hieronder

Formule van Cohen's d

```
def cohen_d(a, b):
    na = len(a)
    nb = len(b)
    pooled_sd = np.sqrt(((na - 1) * np.var(a, ddof=1) + (nb - 1) * np.var(b, ddof=1)) / (na + nb - 2))
    return (np.mean(b) - np.mean(a)) / pooled_sd
```

Interpretatie van Cohen's d (Effect Size)

d	Effectgrootte
0,01	Zeer klein
0,20	Klein
0,50	Gemiddeld
0,80	Groot
1,20	Zeer groot
2,00	Enorm

Stappenplan t-test

1 Formuleer de hypotheses

- H_0 : Er is geen verschil tussen de groepen.
- H_1 : Er is wel een verschil tussen de groepen.

2 Kies de juiste test

- Onafhankelijke groepen → `ttest_ind()`
- Gepaalde metingen (voor/na) → `ttest_rel()`

3 Bereken de teststatistiek

- t-waarde en p-waarde met `stats.ttest_ind()` of `stats.ttest_rel()`

4 Bereken de effectgrootte (optioneel)

- Cohen's d met de handmatige `cohen_d()` functie

5 Maak een beslissing

- Als $p \leq \alpha$ → verwerp H_0 .
- Als $p > \alpha$ → behoud H_0 .

Voorbeeld: t-test voor twee onafhankelijke steekproeven

```
from scipy import stats
t_stat, p_value = stats.ttest_ind(group1, group2)
print(f"t: {t_stat:.3f}, p: {p_value:.4f}")
```

Voorbeeld: t-test voor gepaarde steekproeven

```
from scipy import stats
t_stat, p_value = stats.ttest_rel(before, after)
print(f"t: {t_stat:.3f}, p: {p_value:.4f}")
```

Voorbeeld: Cohen's d

```
d = cohen_d(group1, group2)
print(f"Cohen's d: {d:.3f}")
```

Symbolen en Formules

Symbool	Betekenis
t	t-statistiek
p	p-waarde, kans onder nulhypothese
μ	Populatiegemiddelde
\bar{x}	Steekproefgemiddelde
s	Steekproefstandaardafwijking
n	Steekproefgrootte
Cohen's d	Effectgrootte, verschil in standaarddeviaties

Hoofdstuk 6: Regressie, Covariantie & Correlatie

Beste grafieken (met voorbeeld)

Grafiektype	Omschrijving	Voorbeeldcode
Rel plot	Scatterplot met regressielijn of trends	<code>sns.relplot(data=df, x='x', y='y', kind='scatter')</code>

Grafiektype	Omschrijving	Voorbeeldcode
Regressieplot	Lijnfit met betrouwbaarheidsbanden	<code>sns.regplot(x='x', y='y', data=df)</code> of <code>sns.lmplot(x='x', y='y', data=df)</code>

Methode van de kleinste kwadraten (Method of Least Squares)

De kleinste-kwadratenmethode vindt de lijn (of regressievergelijking) die de som van de kwadraten van de residuen (afstanden tussen werkelijke waarden en voorspelde waarden) minimaliseert.

Formules:

- Regressielijn: $y = \beta_0 + \beta_1 \times x$
- Hellingscoëfficiënt (β_1): $\beta_1 = \text{Cov}(X, Y) / \text{Var}(X)$
- Intercept (β_0): $\beta_0 = \bar{y} - \beta_1 \times \bar{x}$

Voorbeeld β_0 en β_1 berekenen

```
import numpy as np
x = np.array([1, 2, 3, 4, 5])
y = np.array([2, 4, 5, 4, 5])

beta_1 = np.cov(x, y, ddof=1)[0, 1] / np.var(x, ddof=1)
beta_0 = np.mean(y) - beta_1 * np.mean(x)

print(f" $\beta_1$  (helling): {beta_1:.2f}")
print(f" $\beta_0$  (intercept): {beta_0:.2f}")
```

Covariantie

Begrip	Omschrijving	Voorbeeldcode
Covariantie	Meet de mate waarin twee variabelen samen variëren; positief = samen omhoog/omlaag, negatief = tegengesteld	<code>np.cov(x, y)[0, 1]</code>

Correlatie

- Meet de sterkte en richting van een lineair verband tussen twee variabelen.
- Waarden tussen -1 en 1; dicht bij 0 = geen lineair verband.

Functie en voorbeeld:

```
corr = np.corrcoef(x, y)[0, 1]
print(f"Correlatie: {corr:.2f}")
```

Coëfficiënt van determinatie (R^2)

- Geeft het percentage verklaarde variantie aan (hoe goed past het model?).
- $R^2 = r^2$, waar r de correlatiecoëfficiënt is.

Voorbeeld:

```
r = np.corrcoef(x, y)[0, 1]
r_squared = r ** 2
print(f" $R^2$ : {r_squared:.2f}")
```

R	R^2	Verklaarde variantie	Lineair verband
---	-------	----------------------	-----------------

R	R ²	Verklaarde variantie	Lineair verband
< 0,3	< 0,1	< 10%	zeer zwak
0,3 - 0,5	0,1 - 0,25	10% - 25%	zwak
0,5 - 0,7	0,25 - 0,5	25% - 50%	matig
0,7 - 0,85	0,5 - 0,75	50% - 75%	sterk
0,85 - 0,95	0,75 - 0,9	75% - 90%	zeer sterk
> 0,95	> 0,9	> 90%	uitzonderlijk sterk

Tabel met gebruikte termen

Term	Betekenis
Regressielijn	De beste lijn door de data (voorspelling)
Intercept (β_0)	Snijpunt met de y-as
Hellingscoëfficiënt (β_1)	Hoeveel y verandert als x met 1 stijgt
Covariantie	Samenhang in variatie tussen x en y
Correlatie (R)	Sterkte en richting lineair verband
Coëfficiënt van determinatie (R^2)	Hoeveel variantie door model verklaard wordt
Residuen	Verschil tussen geobserveerde en voorspelde waarden

Hoofdstuk 7: Time Series Analyse

Belangrijkste componenten

Component	Definitie
Level	Het gemiddelde of basisniveau van de reeks over tijd
Trend	Langetermijnbeweging (stijgend of dalend patroon)
Seasonal	Terugkerende patronen of fluctuaties over vaste perioden

Wanneer gebruik je welke methode?

Methode	Wanneer gebruiken?	Voorbeeldcode
Simple Moving Average (SMA)	Als je kortetermijnschommelingen wilt gladstrijken	<code>series.rolling(window=3).mean()</code>
Simple Exponential Smoothing (SES)	Voor tijdreeksen met alleen level, zonder trend of seizoen	<code>SimpleExpSmoothing(series).fit()</code>
Double Exponential Smoothing (DES)	Voor tijdreeksen met level én trend, maar zonder seizoen	<code>ExponentialSmoothing(series, trend='add').fit()</code>
Triple Exponential Smoothing (TES, Holt-Winters)	Voor tijdreeksen met level, trend én seizoen (add => mul om van additief naar multiplicatief verband)	<code>ExponentialSmoothing(series, trend='add', seasonal='add', seasonal_periods=12).fit()</code> * add = schommelingen zijn even groot mul = schommelingen worden groter

Uitgewerkte voorbeelden

Simple Moving Average (SMA)

--

```
import pandas as pd
sma = series.rolling(window=3).mean()
print(sma.head())
```

Dit berekent het gemiddelde van telkens 3 opeenvolgende waarden om ruis weg te filteren.

Simple Exponential Smoothing (SES)

```
from statsmodels.tsa.holtwinters import SimpleExpSmoothing
model = SimpleExpSmoothing(series).fit()
pred = model.predict(start=0, end=len(series)-1)
print(pred.head())
```

Gebruik SES als er geen trend of seizoen is, enkel een stabiel gemiddelde.

Double Exponential Smoothing (DES)

```
from statsmodels.tsa.holtwinters import ExponentialSmoothing
model = ExponentialSmoothing(series, trend='add').fit()
pred = model.predict(start=0, end=len(series)-1)
print(pred.head())
```

Gebruik DES als de data een duidelijke stijgende of dalende trend vertoont.

Triple Exponential Smoothing (TES / Holt-Winters)

```
model = ExponentialSmoothing(series, trend='add', seasonal='add', seasonal_periods=12).fit()
pred = model.predict(start=0, end=len(series)-1)
print(pred.head())
```

Gebruik TES als je ook seizoensinvloeden wilt meenemen (bijvoorbeeld maandelijkse data met jaarcyclus).

Evaluatiematen

Maat	Definitie	Voorbeeldcode
MAE	Gemiddelde absolute fout tussen voorspelde en werkelijke waarden	<pre>from sklearn.metrics import mean_absolute_error mae = mean_absolute_error(y_true, y_pred)</pre>
MSE	Gemiddelde van de kwadraten van de fouten (straft grotere fouten sterker)	<pre>from sklearn.metrics import mean_squared_error mse = mean_squared_error(y_true, y_pred)</pre>

Decomposing a Time Series

Decompositie splitst een tijdreeks op in zijn **level**, **trend**, **seasonal** en **residual** componenten om patronen beter te begrijpen.

Voorbeeld decompositie

```
from statsmodels.tsa.seasonal import seasonal_decompose
import matplotlib.pyplot as plt

result = seasonal_decompose(series, model='additive', period=12)
result.plot()
plt.show()
```


Samengevat

- **SMA**: gladstrijken van korte termijn fluctuaties.
- **SES**: voor stabiele reeksen zonder trend/seizoen.
- **DES**: voor reeksen met trend.
- **TES**: voor reeksen met trend én seizoenspatronen.
- **MAE/MSE**: beoordelen hoe nauwkeurig je voorspellingen zijn.
- **Decompositie**: inzicht krijgen in de bouwstenen van je tijdreeks.