

# DSAI\_Samenvatting

May 28, 2022

## 1 Datascience & AI

Samenvatting voor examen van AJ 2021-2022. Door ydm#1001.

---

## 2 Module 1 Basisbegrippen, steekproefonderzoek

---

### 2.1 Basisbegrippen

**Variabele** = algemene eigenschap object, kan objecten onderscheiden. **Waarde** = Specifieke eigenschap, interpretatie van var.

Variabele	Waarde
Gender	Man
Hoogte	180 cm
Funny	Neen.

#### 2.1.1 Meetniveaus

= variabele types. Bepalen beste analyse methode. (visualisatie, centrale tendens en spreiding, verband onderzoeken,...)

**Kwalitatief** = niet noodzakelijk numeriek. Beperkt aantal waardes.

**Nominaal**: categorieën zoals gender, ras, land, vorm,...

**Ordinaal**: Order, rank zoals militaire rank, onderwijsniveau,...

**Kwantitatief** = Numeriek met eenheid. Veel waardes die vaak uniek zijn.

**Interval**: Geen vast nulpunt => geen proporties. <sup>1</sup> (°C, °F)

**Ratio**: Absoluut nulpunt => wel proporties (bv afstand, energie, gewicht,...)

---

<sup>1</sup>20 m is 1/3de (~33%) langer dan 15 meter (wel proportie) <-> 20°C is niet 1/3de warmer dan 15 °C (geen proportie)

**Relaties tussen variabelen.** variabelen hebben en verband als hun waardes systematisch veranderen.

	Pepsi	Coca Cola	Total
Like	56	24	80
Dislike	14	6	20
Total	70	30	100

Totalen zijn *Marginale totalen*

Onderzoek vaak naar **oorzakelijk verband** (frustratie leidt tot agressie, ...).

*Oorzaak:* onafhankelijke variabele

*Verband:* Afhankelijke variabele

**Een verband tussen 2 variabelen zijn niet noodzakelijk oorzakelijk verband!**

## 2.2 Steekproef

**Populatie:** Volledige verzameling objecten/personen die je wilt onderzoeken

**Steekproef:** Deel van de populatie waarop metingen uitgevoerd worden.

In bepaalde gevallen is het resultaat van de steekproef toepasbaar op de volledige populatie.

Steekproefmethode: Bepalen populatie -> bepalen steekproefgrootte -> Kiezen van steekproefmethode (budget en tijd)

Hoe keuze maken voor steekproef?

**aselecte steekproef:** elk element van de populatie heeft evenveel kans om gekozen te worden.

**Niet aselecte steekproef:** De elementen van een sample zijn niet random gekozen. Objecten die makkelijker verkregen worden zijn waarschijnlijker om deel te nemen aan de steekproef. (convenience sampling genoemd).

Stratified to variables: populatie verdeeld op basis van een kenmerk (bijvoorbeeld leeftijd,...). (ook kan vgm bij dit voorbeeld alles /10 gedaan worden (zie slides voorbeeld) en is dit ook stratified)

Gender	<=18	]18,25]	]25,40]	>40	Totaal
Vrouw	500	1500	1000	250	3250
man	400	1200	800	160	2560
Totaal	900	2700	1800	410	5810

### 2.2.1 Fouten

	Steekproeffout	niet steekproeffout
Accidental	Puur toeval	Onjuist antwoord aangeduid

	Steekproeffout	niet steekproeffout
Systematisch	Online onderzoek: mensen zonder internet uitgesloten. Straat onderzoek: enkel die op dat moment daar aan het wandelen is Vrijwilligers onderzoek: enkel geïnteresseerde mensen	Slecht of niet gecalibreerd meetmateriaal Waarde beïnvloed door het feit dat je het meet. Antwoorders liegen (bv aantal sigaretten per dag)

## 2.3 Algemene imports.

```
[ ]: #imports
import numpy as np                # "Scientific computing"
import scipy.stats as stats       # Statistical tests

import pandas as pd               # Data Frame
from pandas.api.types import CategoricalDtype

import matplotlib.pyplot as plt   # Basic visualisation
from statsmodels.graphics.mosaicplot import mosaic # Mosaic diagram
import seaborn as sns             # Advanced data
    ↪ visualisation
import altair as alt              # Alternative visualisation
    ↪ system
import math
```

## 2.4 Python Module 1

```
[ ]: #Import data van een csv file
ais = pd.read_csv('../data/ais.csv')
#indien geen , maar bijvoorbeeld ; gebruikt dan is het
#pd.read_csv(fileLink, delimiter=';')

##Eerste aantal lijnen tonen
ais.head()

#Aantal rijen en kolommen in een dataset printen
print(f"Aantal rijen: {len(ais)}")
#Aantal kolommen
```

```

print(f"Aantal kolommen: {len(ais.columns)}")
#Algemene info over dataset.
ais.info()
#lijntje * printn
print("*"*50)
#Aantal kolommen per type
print(ais.dtypes.value_counts())

#kolom als index instellen
ais.set_index(['id'])

#Voor een kolom categorie als meetvariabele instellen
ais.sex = ais.sex.astype("category")

#Kan ook variabelen als ordinaal aanduiden met een ordening. Bijvoorbeeld als
→we voor sex zouden doen.
# Voorbeeld:
print(ais.sex.unique()) #uniek
sex_Type = CategoricalDtype(categories=['f','m'], ordered=True) #en ordenen
ais.sex= ais.sex.astype(sex_Type)
#een kolom beschrijven
print(ais.ferr.describe())

#SELECTEREN DATA
#Toon de tweede rij
ais.iloc[[1]]

#Toon rij 4 tot en met 6
ais.iloc[4:7]
#Toon KOLOM 6 tem 8: (ferr, bmi, ssf)
ais.iloc[:,5:8]
#Toon 1 variabelen (pcBfat)
ais['pcBfat']

#Toon alles van specifieke query (sport=netball)
ais.query("(sport=='Netball')")
#Toon specifieke kolom met specifieke query (kolom wt van sport=netball)
ais.query("(sport=='Netball')").wt

#Toon alles met een bmi>26
print("BMI ding")
print(ais[ais.bmi>26])

#Toon frequentie en dergelijke
bmiais = ais[ais.bmi>26]
sns.countplot(x=bmiais.sport, data=bmiais)

```

```
#Tel hoe vaak een bepaalde categorie voorkomt  
ais["Sport"].value_counts()
```

---

## 3 Module 2 Analyse van 1 variabele

---

### 3.1 Centrale tendens en spreiding

#### 3.1.1 Maten van centrale tendens

**Mean of Average** De *arithmetic mean* is de som van alle waarden gedeeld door het aantal waarden.

$$> \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

**Median** Sorteer alle waarden en neem het middelste (gemiddelde bij een oneven).

**Mode** de mode is de waarde die het meest voorkomt in een dataset.

#### 3.1.2 Maten van centrale spreiding

**Range** Absolute waarde van het verschil tussen het hoogste en laagste waarde.

**Quartielen** De quantielen van een gesorteerde set zijn 3 waarde die de set in 4 gelijke delen verdelen.  $Q_1$ ,  $Q_2$ ,  $Q_3$

**Variantie:** De variantie ( $S^2$  of  $\sigma^2$ ) is het gemiddelde (mean) van het kwadraat van het verschil van de waarden van de dataset en het gemiddelde (arithmetic mean).

$$> S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

**Standaard afwijking:** De standaard afwijking ( $S$  of  $\sigma$ ) is de wortel van de variantie

#### 3.1.3 Samenvatting Centrale tendens en spreiding

Meetniveau	Center	Sprijdingsmaat
Kwalitatief	Mode	-
Kwantitatief	Average/mean Median	Variantie, standaard afwijking, range, interkwartielafstand

#### 3.1.4 Samenvatting Symbolen

	Populatie	Steekproef
Aantal elementen	N	n
Gemiddelde (mean)	$\mu$	$\bar{x}$
Variantie	$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$	$S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$

	Populatie	Steekproef
Standaard deviatie	$\sigma$	S

## 3.2 Data visualisatie

### 3.2.1 grafiek type overzicht

Meetniveau	Grafiek type
Kwalitatief	Staafdiagram
Kwantitatief	Boxplot Histogram Density plot

#### Taart diagrammen

vermijd gebruiken van taart diagrammen. Hoeken vergelijken is moeilijker dan lengtes, onbruikbaar voor veel categorieën

**Tips** Assen labelen, duidelijke titel, eenheid, label die de grafiek verduidelijkt.

**Data distortion** = zorgt voor fout interpreteren.

## 3.3 Python Module 2

```
[ ]: #distributie van gevens voor sport (distribution)
sns.displot(data= ais["sport"])

#categorie plot voor sports
sns.catplot(data= ais, kind="count", x="sport")

#distribution met Kernel density estimate (soort van normaalverdeling achtige
↳ding te krijgen)
sns.displot(data=ais[ais.sex=="f"].ht, kde=True)

#dingen
rowers = ais[ais.sport == "Row"].ht
print(f"Mean: {rowers.mean()}")
print(f"Standard deviation: {rowers.std()}") # Pay attention: n-1 in the
↳denominator
print(f"Variance: {rowers.var()}") # Pay attention: n-1 in the
↳denominator
print(f"Skewness: {rowers.skew()}")
print(f"Kurtosis: {rowers.kurtosis()}")

# Median & co
```

```

print(f"Minimum:    {rowers.min()}")
print(f"Median:     {rowers.median()}")
print(f"Maximum:    {rowers.max()}")
percentiles = [0.0, 0.25, 0.5, 0.75, 1.0]
print("Percentiles", percentiles, "\n", rowers.quantile(percentiles))
print("Inter Quartile Range:", rowers.quantile(.75) - rowers.quantile(.25))
print(f"Range :     {rowers.max() - rowers.min()}")

```

---

## 4 Module 3.1 De centrale limietstelling, betrouwbaarheidsintervallen

---

### 4.1 Kansverdeling van een steekproef

#### 4.1.1 Kans

Kans is de relatieve frequentie van het voorkomen van een bepaald event (bij uitvoeren van groot aantal onafhankelijke experimenten)

- kansen zijn getallen aan een set toegewezen - Die sets zijn deel van een allesomvattende set, het *universum*  $\Omega$  - De nummers (kansen) toegewezen aan een set voldoen aan 3 basis regels (axiom van kans) om overeen te komen met hoe kansen werken

1. Kansen zijn niet negatief  $P(A) \geq 0$  voor elke  $A$ .
2. Het universum heeft een kans 1:  $P(\Omega) = 1$ .
3. Als  $A$  en  $B$  disjunct zijn ( $A \cap B = \emptyset$ ) dan geldt  $P(A \cup B) = P(A) + P(B)$  dit heet de somregel.

**Eigenschappen**

1. Complement regel: voor elke  $A$  geldt  $P(\bar{A}) = 1 - P(A)$  als  $\bar{A}$  voorstelt dat  $A$  niet voorkomt.
2. Het onmogelijke event is kans nul:  $P(\emptyset) = 0$
3. De algemene som regel:  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

**Onafhankelijke events** Een event is onafhankelijk als het voorkomen van dit event (of het weten dat dit voorkomt) de kans dat een ander event gebeurt niet beïnvloed. Wiskundig:  $P(A \cap B) = p(A)p(B)$

#### 4.1.2 Random variabele

Een random variable is een waarde toekennen aan verschillende gebeurtenissen. Bijvoorbeeld.

1 als je een J trekt uit een kaart spel, 2 bij een Q en 3 bij een K en tot slot 0 bij alle andere mogelijkheden

Kansverdeling functie (PDF) wiskundig.  $f_x(x) = P(X = x)$

Voorbeeld:

**Expectation of a R.V.** Verwachting van een random variabele is geschreven door  $\mu_x$  of  $E(X)$  en is gegeven door  $\mu_X = \sum_i x_i P(X = x_i) = \sum_i x_i f_x(x_i)$

**Variantie van een R.V.** De variantie van een random variabele is bepaald door  $\sigma^2 = \sum_i (x_i - \mu_x)^2 P(X = x_i) = \sum_i (x_i - \mu_x)^2 f_x(x_i)$ .

Standaard afwijking:  $\sigma_x = \sqrt{\sigma_x^2}$

## Continue random variabele

- Een continue R.V. neemt een ontelbaar oneindig aantal mogelijke waarden
- In dat geval niet logisch dat de kans van  $X = a$  exact te bekijken, omdat de kans altijd 0 is.
- Wat wel zin heeft is te bekijken wat de kans is van  $X=[a,b]$ .
- Deze kans kan gevonden worden door te integreren de PDF van random variabale.
- voorbeeld van Continue random value: lengte van mensen in een populatie...

De lengte van mensen volgt vaak ongeveer een **Normale verdeling**. De normale verdeling is een type van **Continue kansverdeling**. De formules voor de variantie en de verwachting zijn dezelfde als voor gewone R.V.'s maar dan met een integraal van -oneindig tot +oneindig.

### 4.1.3 Standaard normaal verdeling

$x$  en  $z$  hebben gelijkaardige positie in de gauss curve. Wat is de wiskundige relatie tussen  $x$  en  $z$ ?

$$x = \mu + z * \sigma \text{ and } z = \frac{x - \mu}{\sigma}$$

- $[-1, 1]$  ( $[\mu - \sigma, \mu + \sigma]$ ) bevat 68,3% van de populatie of kans
- $[-2, 2]$  ( $[\mu - 2\sigma, \mu + 2\sigma]$ ) bevat 95,4%
- $[-3, 3]$  ( $[\mu - 3\sigma, \mu + 3\sigma]$ ) bevat 99,7%

### 4.1.4 Exponentiele spreiding

is een andere veelgebruikte continue distributie. Dit gebeurt als er minder grote waarden zijn en meer kleine waarde. Bijvoorbeeld het bedrag dat klanten uitgeven volgt een exponentiele distributie. Er zijn meer mensen die kleine bedragen uitgeven dan mensen die veel uitgeven.

### 4.1.5 Continue uniforme spreiding

beschrijft een experiment waar een arbitreire uitkomst is tussen bepaalde grenzen. de density functie is constant omdat er voor elke waarde een even grote kans is dat ze voorkomt. Bijvoorbeeld een lift gaat altijd tussen de 10 en 15 seconden naar de 2de verdieping dan is de kans altijd 1 dat je binnen de 10 en 15 seconden op dat verdiep bent. (neem de trap is beter voor je gezondheid!).

## 4.2 Van steekproef naar populatie

### 4.2.1 De centrale limietstelling

Als de steekproefgrootte groot genoeg is dan zal de kansverdeling van het steekproefgemiddelde ongeveer een normale verdeling zijn, onafhankelijk van de kansverdeling van de onderliggende populatie.

Bekijk een random steekproef met  $n$  observaties uit een populatie met verwachte waarde  $\mu$  en standaard deviatie  $\sigma$ . Als  $n$  groot genoeg is dan zal de kansdichtheid van de steekproefgemiddelde  $\bar{x}$  ongeveer een normale verdeling met gemiddelde  $\mu_{\bar{x}}$  en een standaardafwijking  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$ . Hoe groter de steekproef hoe beter de kansverdeling van  $\bar{x}$  zal benaderen met verwachte waarde van de populatie,  $\mu$ .

### 4.2.2 Punt schatting

Een punt schatting van een populatie parameter is een formule of vergelijking die toestaat om een verwachte waarde te bepalen voor die parameter.



### 4.2.3 Confidence interval

Een confidence interval is een vergelijking of formule die toestaat een interval op te stellen die met een bepaalde zekerheid een parameter bevat.

Voor kleine steekproef is de centrale limietstelling niet geldig. In de plaats daarvan zeggen we als de populatie X een normale verdeling heeft en je hebt een kleine steekproef met  $\bar{x}$  en standaard afwijking s dan  $t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$  zal gedragen als een t-distributie met  $n-1$  graden van vrijheid.

## 4.3 Python Module 3.1

```
[ ]: #voorbeeld incl print met 3 cijfers na komma.

#kansberekening voor normale verdeling met
m = 0 #gemiddelde
s = 1 #standaard afwijking
print('P(Z<1.33)=%.3f'%stats.norm.cdf(1.33, loc=m, scale=s))
print('P(Z>1.33)=%.3f'% stats.norm.sf(1.33, loc=m,scale=s))
print('P(-1.35<Z<-0.10)=%.3f'% (stats.norm.cdf(-0.10, loc=m,scale=s)-stats.norm.
    ↪cdf(-1.35, loc=m,scale=s)))

## Probability density function (blauw) ingekleur
m=2.5
s= 1.5
dist_x = np.linspace(m - 4 * s, m + 4 * s, num=201)
dist_y = stats.norm.pdf(dist_x, m, s)
plt.plot(dist_x, dist_y)
plt.fill_between(dist_x, 0, dist_y, color='lightblue')
#zelfde oef de cdf (oranje)
dist_y_cdf = stats.norm.cdf(dist_x)
plt.plot(dist_x, dist_y_cdf)
#de area onder de pdf tussen 0.5 en 4
stats.norm.cdf(4, loc=m,scale=s)-stats.norm.cdf(0.5, loc=m,scale=s)

#Genereer random nummers die de standaard normaal verdeling volgen
n=25
observations = np.random.normal(loc=m, scale=s, size=n)
#print een histogram met kansdichtheids functie en theoretische kansdichtheids
sns.histplot(observations, kde=True)

#Bij standaardafwijking en gemiddelde van populatie naar s en m van een
    ↪steekprof te gaan (n= steekproefgrootte)
n=81
```

```
standaarddivPop= 36
s= standaarddivPop/math.sqrt(n) #s= standaarddiv van steekproef voor dat
↪ gemiddelde
```

---

## 5 Module 3.2 Hypothesetesten

---

### 5.1 Testprocedure

#### 5.1.1 Statistische hypothesetesten

- **Hypothesis:** idee dat nog moet bewezen worden: statement over een numerische waarde van een populatie parameter.
- **Hypothesetest:** Verificatie van het statement over de waarden van 1 of meerdere populatie parameters.
- **Null Hypothese ( $H_0$ ):** Basis hypothese, aan nemen dat die waar is
- **Alternatieve Hypothese ( $H_1, H_a$ ):** Conclusie als de null hypothese waarschijnlijk fout is.

#### 5.1.2 Elementen van een test procedure

- **Test statistiek:** De waarde die berekend word van een steekproef
- **Acceptatieregio:** De regio van waarden die de null hypothesis bevestigen
- **Kritieke regio/ Regio van afwijzing:** De regio van waarden die de null hypothesis verworpen.
- **Significantie niveau:** De waarschijnlijkheid van verwerpen van de null hypothese  $H_0$

#### 5.1.3 Test procedure

1. Formuleer beide hyptoheses ( $H_0, H_1$ )
2. Bepaal significantie niveau ( $\alpha$ )
3. Berkenen test statistiek
4. Bepalen kritieke regio van de kans waarde
5. Conclusies trekken.

### 5.2 Kans waarde (Probability value)

**p-waarde:** De p-waarde is de kans, als de null hyptohese waar is, een waarde voor de test statistiek te krijgen die minimaal zo extreem is als de geobserveerde waarde.

- p-waarde  $< \alpha \Rightarrow$  verwerpen  $H_0$ : de ontdekte waarde voor  $\bar{x}$  is te extreem.
- p-waarde  $\leq \alpha \Rightarrow H_0$  niet verwerpen: De ontdekte waarde voor  $\bar{x}$  kan nogsteeds uitgelegd worden door toeval.

### 5.3 Kritieke regio

De Kritieke regio is de verzameling van alle waarden van een test statistiek waarvoor de null hypothese verworpen kan worden.

Kijk naar kritieke waarde  $g$  waarvoor geldt:  $P(M > g) = \alpha$

Bepaal  $z_\alpha$  waarvoor geldt:  $P(Z > z_\alpha) = \alpha \Rightarrow g = \mu + z_\alpha * \frac{\sigma}{\sqrt{n}}$

- Links van  $G$ : Regio van acceptatie ( $H_0$  niet verwerpen)
- Rechts van  $G$ : Kritieke regio ( $H_0$  verwerpen)

### 5.3.1 Samengevat Testing procedures

Goal	Test met betrekking tot de waarde van de populatie gemiddelde $\mu$ gebruik makend van steekproef van $n$ onafhankelijke waarden.
Voorwaarde	De populatie heeft een random verdeling, $n$ is groot genoeg

Test type	Two-tailed	Left-tailed	Right-tailed
$H_0$	$\mu = \mu_0$	$\mu = \mu_0$	$\mu = \mu_0$
$H_1$	$\mu \neq \mu_0$	$\mu < \mu_0$	$\mu > \mu_0$
Critieke regio	$ \bar{x}  > g$	$\bar{x} < -g$	$\bar{x} > g$
Test statistiek	$z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$		

### 5.3.2 Voorwaardes voor z-test

- De steekproef moet aselekt zijn
- De steekproefgrootte moet groot genoeg zijn ( $n \geq 30$ )
- De test statistiek moet zich gedragen als een normale verdeling
- De standaarddeviatie van de populatie  $\sigma$  is gekend.

Soms zijn de voorwaardes niet voldaan en dan kan geen z-test

## 5.4 Voorbeelden

zie slides.

## 5.5 Student's t-test

Wat als voorwaardes voor z-test niet voldaan zijn? - steekproef niet groot genoeg - populatie standaard deviatie niet gekend Als de variabele normaal verdeeld zijn kunnen we de t-test gebruiken.

### 5.5.1 De t-test

Kritieke waarde bepalen:

$$g = \mu \pm t * \frac{s}{\sqrt{n}}$$

- t-waarde afgeleid uit de Student t-distributie gebaseerd op  $n-1$  vrijheidsgraden
- De waarde opzoeken door `t.isf` in Python
- Los van dit is de procedure gelijk aan de procedure voor de z-test.

## 5.6 Fouten in hypothese testen

	Rea	liteit
Conclusion	$H_0$ True	$H_1$ True
$H_0$ niet verworpen	correct conclusie	Type II fout (vals negatief)
$H_0$ verworpen	Type I fout (vals positief)	correct conclusie

$P(\text{type I error}) = \alpha$  (=significantie level)

$P(\text{type II error}) = \beta$

$\beta$  berekenen is niet triviaal maar als  $\alpha$  afneemt dan neemt  $\beta$  toe.

## 5.7 Python Module 3.2

```
[ ]: #voorbeeld met Z
m = 44 #gemiddelde populatie
s = 6.2 #standaard dev populatie
n = 72 #steekproefgrootte
m_intro = 46.2 #steekproef gemiddelde
s_intro = s/math.sqrt(n) #steekproef std?
a = 0.025 #alpha waarde, significantie niveau

#Plot de grafiek voor ja de normaal verdeling enz
dist_x = np.linspace(m-4*s_intro,m+4*s_intro, num=201)
dist_y = stats.norm.pdf(dist_x, m, s_intro)
plt.plot(dist_x, dist_y)

#Inkleuren het stuk dat beter is? ofja gun het te onderzoeken het H0 ding
plt.fill_between(dist_x, 0, dist_y, where=(dist_x>=m_intro), color='red')
#vertikale lijn daarvoor
plt.axvline(m_intro, color="green")

#p waarde berekenn
p_waarde = stats.norm.sf(m_intro, loc=m, scale=s_intro)
print("p waarde: %.4f"%p_waarde)
if(p_waarde < a):
    print("p < a: reject H0")
else:
    print("p > a: do not reject H0")
#kritieke regio bepalen
g_value = stats.norm.isf(0.025, loc = m, scale= s_intro)
print("Critical value g   %.3f" % g_value)
if (m_intro < g_value):
    print("sample mean = %.3f < g = %.3f: do not reject H0" % (m_intro,
    ↪g_value))
else:
```

```

    print("sample mean = %.3f > g = %.3f: reject H0" % (m_intro, g_value))

#voorbeeld met t
#H0: dat het niet significant groter is
prijs= [400, 350, 400, 500, 300, 350, 200,
        500, 200, 250, 250, 500, 350, 100] #is gun dataset steekproef
m = 300 #mag ni significant groetr zijn dan
a=0.05 #significantieniveau
n = len(prijs) #lengte sample
m_samp = np.mean(prijs) #gemiddelde samp
s = np.std(prijs, ddof=1) #std samp
sn = s/math.sqrt(n) #s/wortel(n)

print(f"mean is {m_samp}")
print(f"std is {s}")

#tekenen
dist_x = np.linspace(m-4*sn,m+4*sn, num=201)
dist_y = stats.t.pdf(dist_x, loc=m, scale=sn, df=n-1)
plt.plot(dist_x, dist_y)
#aanduiden welk deel "oke" is voor H0
plt.fill_between(dist_x, 0, dist_y, where=(dist_x<=m_samp), color='yellow')
plt.axvline(m_samp, color="blue")

#p waarde
p = stats.t.sf(m_samp, loc=m, scale = sn, df = n-1)
print("p:value %.5f"% p)
if(p<a):
    print("p<a:rejectH0")
else:
    print("p>a: niet reject H0")

#kritieke waarde (g)
g = stats.t.isf(a, loc=m, scale=sn, df=n-1)
print('critiek g ong= %.3f' % g)
if (m_samp < g):
    print("sample mean = %.3f < g = %.3f: do not reject H0" % (m_samp, g))
else:
    print("sample mean = %.3f > g = %.3f: reject H0" % (m_samp, g))

```

## 6 Samenvatting Symbolen

	Populatie	Steekproef
Aantal elementen	N	n

	Populatie	Steekproef
Gemiddelde (mean)	$\mu$	$\bar{x}$
Variantie	$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$	$S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$
Standaard deviatie	$\sigma$	S

	Symbol	Formule
Expectation van een random value	$\mu_x$ of $E(X)$	$\mu_X = \sum_i x_i P(X = x_i) = \sum_i x_i f_x(x_i)$
Variantie van een random value	$\sigma^2$	$\sum_i (x_i - \mu_x)^2 P(X = x_i) = \sum_i (x_i - \mu_x)^2 f_x(X_i)$
Standaardafwijking van een R.V.	$\sigma_x$	$\sqrt{\sigma_x^2}$

	Symbol	Tailed	Formule
Kritieke regio	g	Rechts	$g = \mu + z_\alpha * \frac{\sigma}{\sqrt{n}}$
		Links	$g = \mu - z_\alpha * \frac{\sigma}{\sqrt{n}}$
		2 zijdig:	$g = \mu \pm z_\alpha * \frac{\sigma}{\sqrt{n}}$

## 7 Samenvatting python code (zoals in slides in overzicht staan)

### 7.1 Normaalverdeling met mean m en standaard deviation s

Function stats.	Doel
<code>norm.pdf(x, loc=m, scale=s)</code>	kansdichtheid bij X
<code>norm.cdf(x, loc=m, scale=s)</code>	Links kans (left-tail) \$ P(X < x)\$
<code>norm.sf(x, loc=m, scale=s)</code>	Rechts kans (right-tail) \$ P(X > x)\$
<code>norm.isf(1-p, loc=m, scale=s)</code>	p% van observaties die lager verwacht werden dan het resultaat

### 7.2 Student t-verdeling

(df= degrees of freedom) |Function stats. | Betekenis | |`t.pdf(x, df=d)`| kansdichtheid voor x| |`t.cdf(x, df=d)`| Left-tail kans \$ P(X < x)\$| |`t.sf(x, df=d)`| Right-tail kans \$ P(X > x)\$| |`t.isf(1-p, df=d)`| p% van observaties is de verwachting lager als het resultaat|

### 7.3 Enkele functies

```
[ ]: def kansdichtheidNorm(x, mean, standardDiviatie):
    return stats.norm.pdf(x, loc=mean, scale=standardDiviatie)

def leftTailNorm(x, mean, standardDiviatie):
    return stats.norm.cdf(x, loc=mean, scale=standardDiviatie)
```

```

def rightTailNorm(x, mean, standardDiviatie):
    return stats.norm.sf(x, loc=mean, scale=standardDiviatie)

def zScore(x, mean, standardDiviatie):
    return (x-mean)/standardDiviatie

#Confidence interval Large sample (met z)
def confIntervallLarge(confidenceLevel, meanSample, standaardDivSample):
    alpha = 1-confidenceLevel
    p = 1-alpha/2
    zAlpha2Kans = stats.norm.isf(1-p)
    return [meanSample-zAlpha2Kans*(standaardDivSample),
    ↪meanSample+zAlpha2Kans*(standaardDivSample)]

#confidence interval small sample (met t)
def confIntervalSmall(confidenceLevel, meanSample, standaardDivSample,
    ↪sampleSize):
    alpha = 1-confidenceLevel
    p = 1-alpha/2
    zAlpha2Kans = stats.t.isf(1-p, df=(sampleSize-1))
    return [meanSample-zAlpha2Kans*(standaardDivSample/math.sqrt(sampleSize)),
    ↪meanSample+zAlpha2Kans*(standaardDivSample/math.sqrt(sampleSize))]

```