

# Introduktion til Statistik

## Forelæsning 3: Kontinuerte fordelinger

Peder Bacher

DTU Compute, Dynamiske Systemer  
Bygning 303B, Rum 009  
Danmarks Tekniske Universitet  
2800 Lyngby – Danmark  
e-mail: pbac@dtu.dk

Forår 2017

## Kapitel 2: Kontinuerte fordelinger

### Grundlæggende koncepter:

- Tæthedsfunktion:  $f(x)$  (*pdf*)
- Fordelingsfunktion:  $F(x) = P(X \leq x)$  (*cdf*)
- Middelværdi ( $\mu$ ) og varians ( $\sigma^2$ )
- Regneregler for stokastiske variable

### Specifikke fordelinger:

- Normal
- Log-Normal
- Uniform
- Eksponential
- $t$
- $\chi^2$  (*Chi-i-anden*)
- $F$

## Chapter 2: Continuous Distributions

### General concepts:

- Density function:  $f(x)$  (*pdf*)
- Distribution:  $F(x) = P(X \leq x)$  (*cdf*)
- Mean ( $\mu$ ) and variance ( $\sigma^2$ )
- Calculation rules for random variables

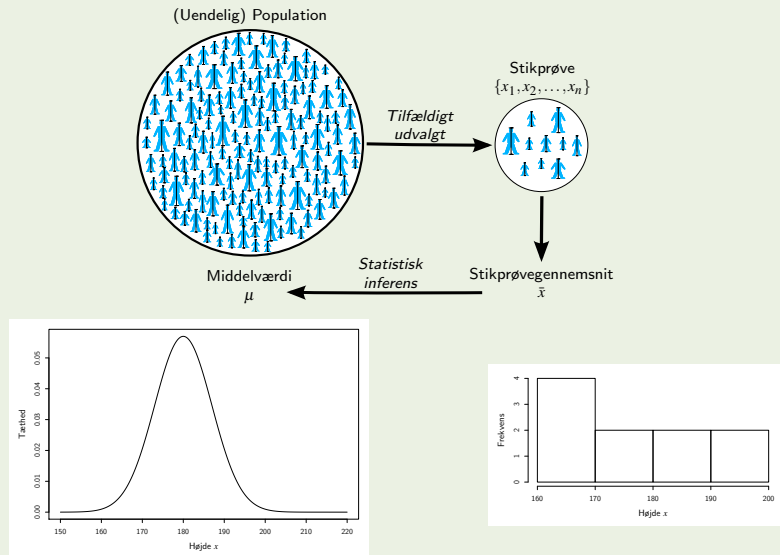
### Specific distributions:

- Normal
- Log-Normal
- Uniform
- Exponential
- $t$
- $\chi^2$  (*Chi-square*)
- $F$

## Oversigt

- 1 Kontinuerte Stokastiske variable og fordelinger
  - Tæthedsfunktion
  - Fordelingsfunktion
  - Middelværdi af en kontinuert stokastisk variabel
  - Varians af en kontinuert stokastisk variabel
- 2 Konkrete Statistiske fordelinger
  - Kontinuerte fordelinger i R
  - Uniform fordeling
  - Normalfordelingen
  - Log-Normal fordelingen
- 3 Eksponentialfordelingen
- 4 Regneregler for middelværdi og varians

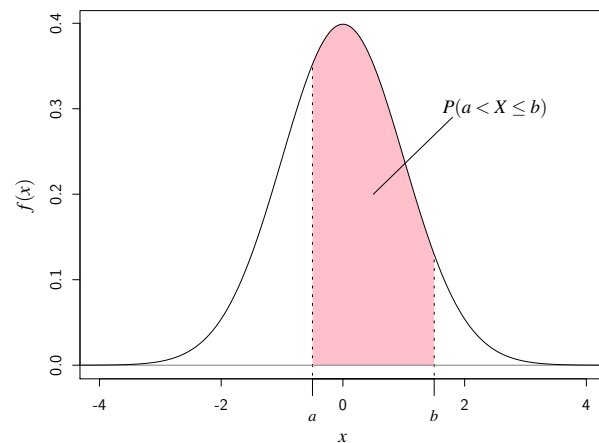
## Example: Population and distribution



## Tæthedsfunktion (probability density function (pdf))

- **Tæthedsfunktionen** for en stokastisk variabel betegnes ved  $f(x)$
- $f(x)$  siger noget om hyppigheden af udfaldet  $x$  for den stokastiske variabel  $X$
- For kontinuerte variable svarer tætheden ikke til sandsynligheden, dvs.  $f(x) \neq P(X = x)$
- Et godt plot af  $f(x)$  er et histogram (kontinuert)

## Tæthedsfunktion for en kontinuert variabel



## Tæthedsfunktion for en kontinuert variabel

- Der gælder:
    - Ingen negative værdier
- $$f(x) \geq 0 \quad \text{for alle mulige } x$$
- Areal under kurven er een

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

## Fordelingsfunktion (distribution function eller cumulative density function (cdf))

- **Fordelingsfunktion** for en kontinuert stokastisk variabel betegnes ved

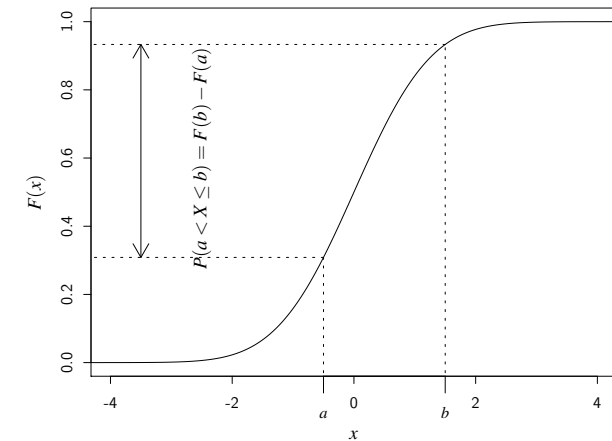
$$F(x)$$

- **Fordelingsfunktionen** svarer til den kumulerede tæthedsfunktion ved

$$F(x) = P(X \leq x)$$

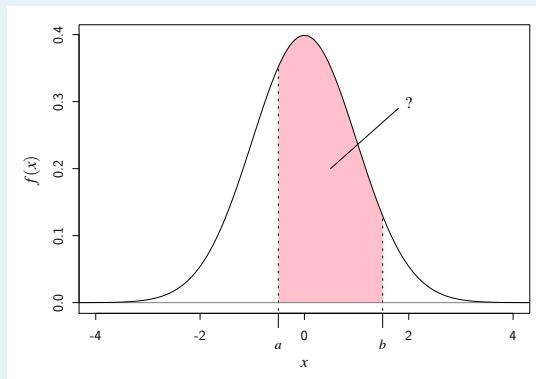
$$F(x) = \int_{-\infty}^x f(u) du$$

$$f(x) = F'(x)$$



## Fordelingsfunktion (distribution function eller cumulative density function (cdf))

### Spørgsmål om sandsynligheder (socrative.com, room: PBAC)

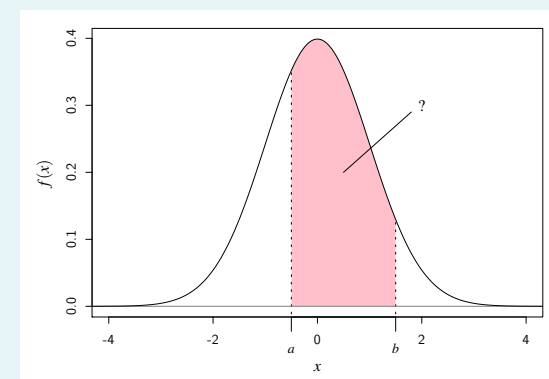


Hvilken sandsynlighed (areal) er markeret?

A:  $\int_{-\infty}^b f(x) dx$     B:  $1 - \int_a^b f(x) dx$     C:  $\int_a^b f(x) dx$     D:  $1 - \int_a^{\infty} f(x) dx$

Svar C:  $\int_a^b f(x) dx$

### Spørgsmål om sandsynligheder (socrative.com, room: PBAC)



Hvordan kan vi nemmest udregne det markerede areal?

A:  $\int_a^b f(x) dx$     B:  $\int_a^b F(x) dx$     C:  $f(b) - f(a)$     D:  $F(b) - F(a)$

Svar D:  $F(b) - F(a)$  (vi gør det i R med (normalfordelt): `pnorm(b) - pnorm(a)`)

## Middelværdi (mean) af en kontinuert stokastisk variabel

**Middelværdien** af en kontinuert stokastisk variabel

$$\mu = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

Sammenlign med den diskrete definition:  $\mu = \sum_{\text{alle } x} x \cdot f(x)$

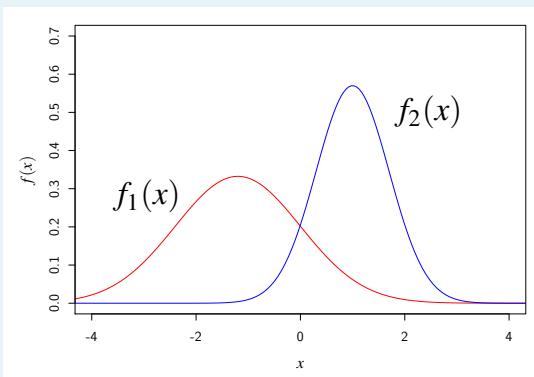
## Varians af en kontinuert stokastisk variabel

**Variansen** af en kontinuert stokastisk variabel:

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx$$

Sammenlign med den diskrete definition:  $\sigma^2 = \sum_{\text{alle } x} (x - \mu)^2 \cdot f(x)$

## Spørgsmål om middelværdi (socrative.com, room: PBAC)

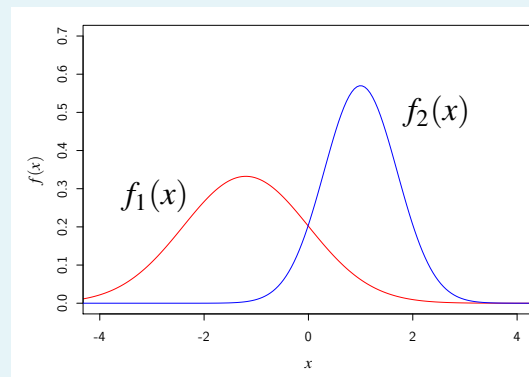


Hvilken pdf har størst middelværdi (begge er symmetriske)?

A:  $\mu_1 < \mu_2$     B:  $\mu_1 > \mu_2$     C:  $\mu_1 = \mu_2$     D: Kan ikke afgøres

Svar A:  $\mu_1 < \mu_2$  (umiddelbart). Svar D, også fint, da man ikke kan se hvad der er udenfor plottet.

## Spørgsmål om spredning (socrative.com, room: PBAC)



Hvilken pdf har størst spredning (begge er symmetriske)?

A:  $\sigma_1 < \sigma_2$     B:  $\sigma_1 > \sigma_2$     C:  $\sigma_1 = \sigma_2$     D: Kan ikke afgøres

Svar B:  $\sigma_1 > \sigma_2$  (umiddelbart). Svar D, også fint, da man ikke kan se hvad der er udenfor plottet.

## Konkrete statistiske fordelinger

Der findes en række statistiske fordelinger, som kan bruges til at beskrive og analysere forskellige problemstillinger med

- Følgende kontinuerte fordelinger:
  - Uniform fordeling
  - Normalfordelingen
  - Log-normalfordelingen
  - Eksponentialfordelingen

## Kontinuerte fordelinger i R

R	Betegnelse
<code>norm</code>	Normalfordelingen
<code>unif</code>	Uniform fordeling
<code>lnorm</code>	Log-normalfordelingen
<code>exp</code>	Eksponentialfordelingen

- `d` Tæthedsfunktion  $f(x)$  (probability density function).
- `p` Fordelingsfunktion  $F(x)$  (cumulative distribution function).
- `q` Fraktil (quantile) i fordeling.
- `r` Tilfældige tal fra fordelingen.

## Uniform fordeling

Skrivemåde:

$X \sim U(\alpha, \beta)$  (Læses:  $X$  følger en uniform fordeling med parametre  $\alpha$  og  $\beta$ )

Tæthedsfunktion:

$$f(x) = \frac{1}{\beta - \alpha}$$

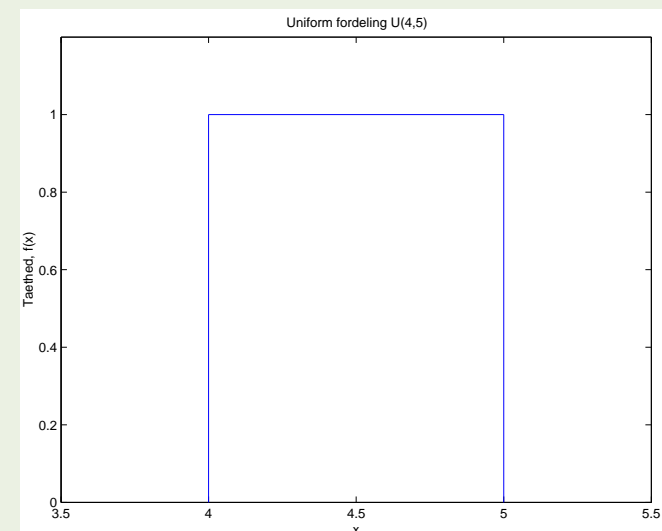
Middelværdi:

$$\mu = \frac{\alpha + \beta}{2}$$

Varians:

$$\sigma^2 = \frac{1}{12}(\beta - \alpha)^2$$

## Eksempel: Uniform fordeling



## Spørgsmål: Uniform fordeling (socrative.com, room: PBAC)

Medarbejdere på en arbejdsplads ankommer mellem klokken 8:00 og 8:30. Det antages, at ankomsttiden kan beskrives ved en uniform fordeling.

Hvad er sandsynligheden for at en tilfældig udvalgt medarbejder ankommer mellem 8:20 og 8:30?

A: 1/2    B: 1/6    C: 1/3    D: 0

Svar C:  $10/30=1/3$

`punif(30,0,30)-punif(20,0,30)`

[1] 0.33

## Spørgsmål: Uniform fordeling (socrative.com, room: PBAC)

Medarbejdere på en arbejdsplads ankommer mellem klokken 8:00 og 8:30. Det antages, at ankomsttiden kan beskrives ved en uniform fordeling.

Hvad er sandsynligheden for at en tilfældig udvalgt medarbejder ankommer efter 8:30?

A: 1/2    B: 1/6    C: 1/3    D: 0

Svar:  $P(X > 30) = 0$

`1-punif(30,0,30)`

[1] 0

## Normalfordelingen

Skrivemåde:

$$X \sim N(\mu, \sigma^2)$$

Tæthedsfunktion:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

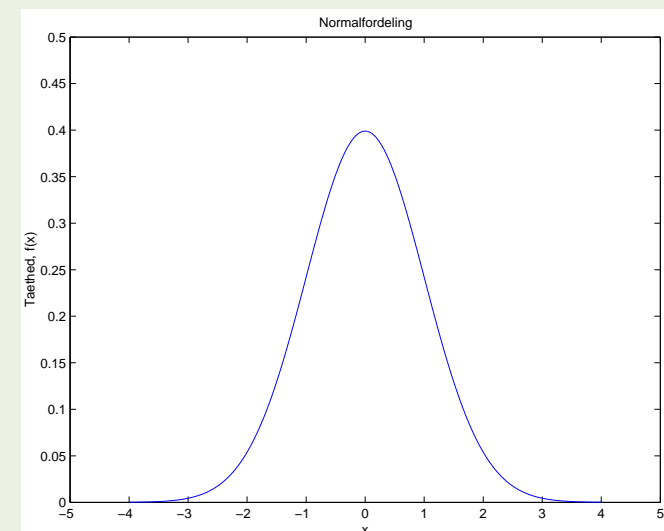
Middelværdi:

$$\mu = \mu$$

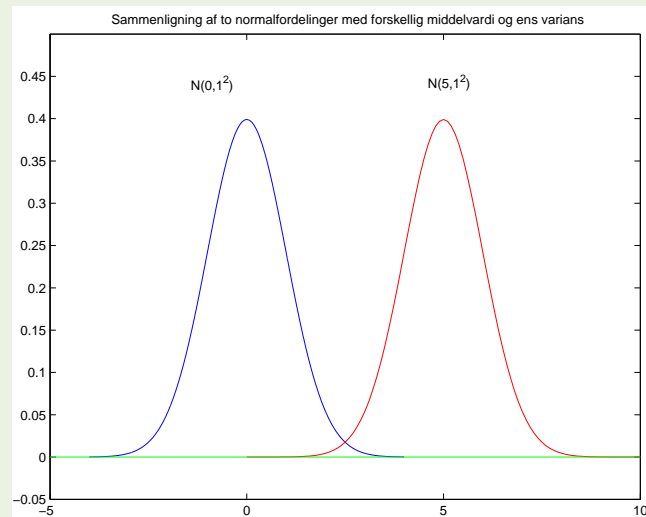
Varians:

$$\sigma^2 = \sigma^2$$

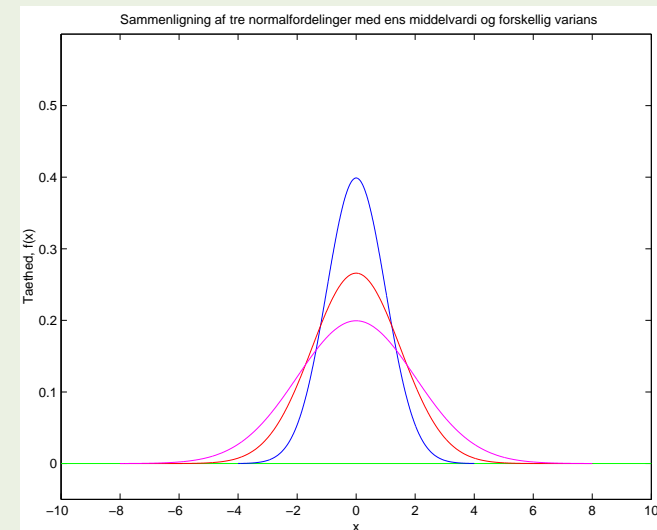
## Eksempel: Normalfordelingen



## Eksempel: Normalfordelingen



## Eksempel: Normalfordelingen



## Eksempel: Normalfordeling, sandsynligheder

## Fordeling af vægt af rugbrød:

Antag at vægten af et rugbrød fra en produktionslinie kan beskrives med en normalfordeling

$$X \sim N(500, 10^2)$$

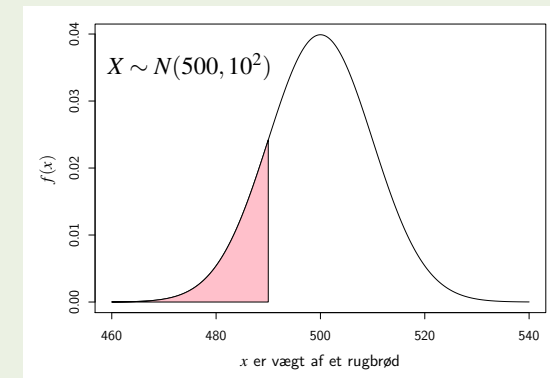
dvs. middelværdi  $\mu = 500$  gram og standardafvigelse  $\sigma = 10$  gram.

Vi vil måle vægten af ét tilfældigt udvalgt brød.

## Spørgsmål:

- 1: Hvad er sandsynligheden for at brødet vejer under 490 g?
- 2: Hvad er sandsynligheden for at brødet vejer mere en  $\pm 20$  g fra 500 g?

## Eksempel: Normalfordeling, spørgsmål 1

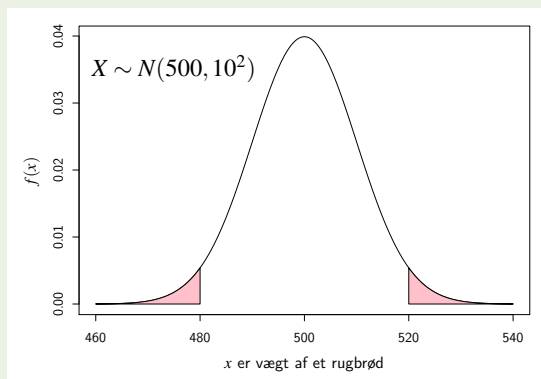


1: Hvad er sandsynligheden for at brødet vejer under 490 g?

Svar:  $P(X \leq 490) = F(490) = 0.16$

```
pnorm(490, mean=500, sd=10)
```

## Eksempel: Normalfordeling, spørgsmål 2

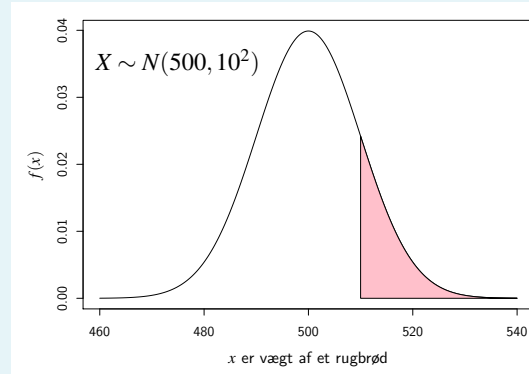


1: Hvad er sandsynligheden for at brødet vejer mere end  $\pm 20$  g fra 500 g?

Svar:  $P(X \leq 480 \vee X > 520) = 2 \cdot P(X \leq 480) = 2 \cdot F(480) = 0.046$

```
2 * pnorm(480, mean=500, sd=10)
```

## Spørgsmål: Sandsynlighed i normalfordeling

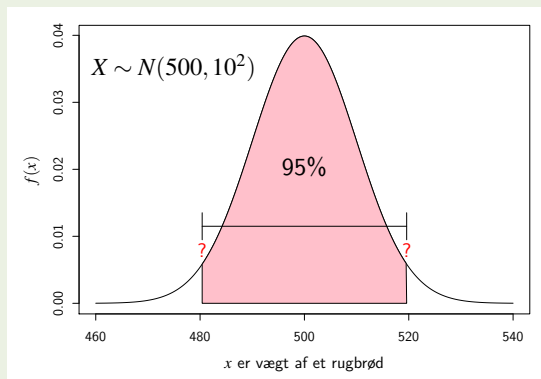


Hvad er sandsynligheden for at rugbrødet vejer over 510 g lig med?

A:  $F(510)$     B:  $1 - F(490)$     C:  $1 - F(520)$     D:  $1 - F(510)$

Svar:  $P(X > 510) = 1 - P(X \leq 510) = 0.16$

## Eksempel: Normalfordeling fraktiler



"Omvendt spørgsmål": Hvilket interval dækker 95% af rugbrødene?

```
qnorm(c(0.025, 0.975), mean=500, sd=10)
```

[1] 480.4 519.6

## Standard normalfordelingen

En standard normalfordeling

$$Z \sim N(0, 1^2)$$

En normalfordeling med middelværdi 0 og varians 1.

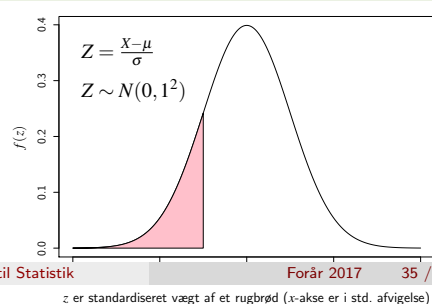
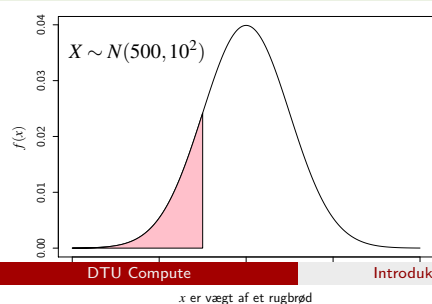
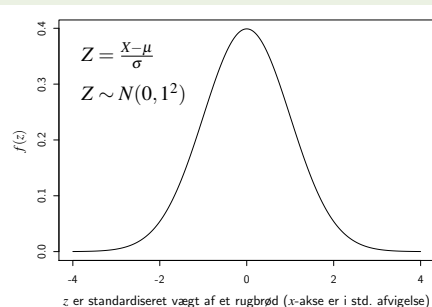
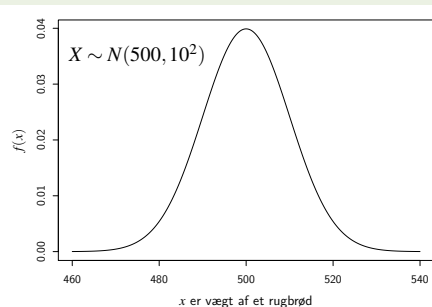
Standardisering

En vilkårlig normalfordelt variabel  $X \sim N(\mu, \sigma^2)$  kan standardiseres ved at beregne

$$Z = \frac{X - \mu}{\sigma}$$



## Eksempel: Standard Normalfordeling



1: *Hvad er sandsynligheden for at brødet veier under 490 gram?*

## Log-Normalfordelingen

Skrivemåde:

$X \sim LN(\alpha, \beta^2)$  (Hvis  $X$  følger log-normal så følger  $\ln(X)$  normal)

Tæthedsfunktion:

$$f(x) = \begin{cases} \frac{1}{x\sqrt{2\pi}\beta} e^{-(\ln(x)-\alpha)^2/2\beta^2} & x > 0, \beta > 0 \\ 0 & \text{ellers} \end{cases}$$

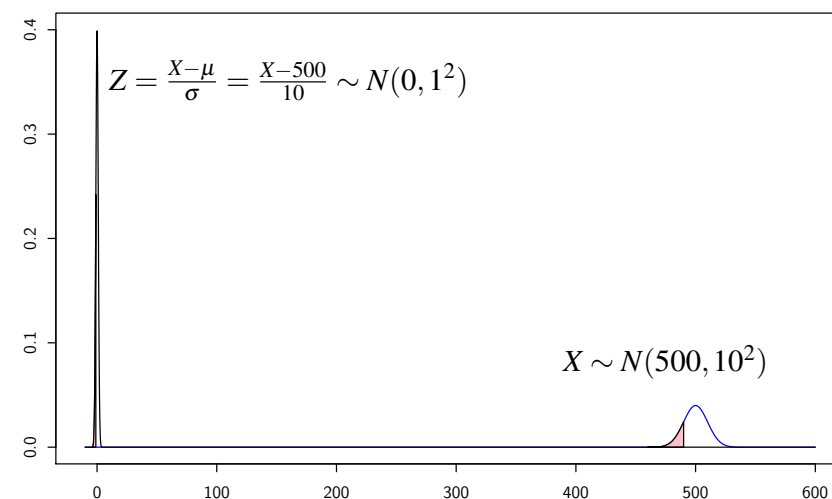
Middelværdi:

$$\mu = e^{\alpha + \beta^2/2}$$

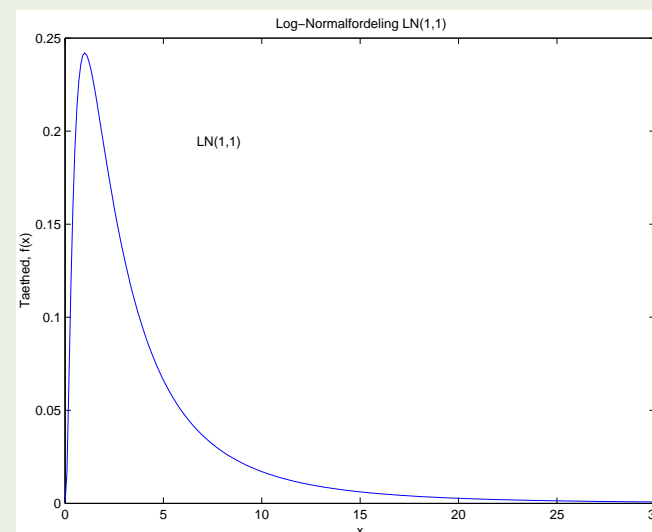
Varsians:

$$\sigma^2 = e^{2\alpha + \beta^2} (e^{\beta^2} - 1)$$

## Eksempel: Transformation til standard normalfordeling



## Eksempel: Log-normal fordelingen



## Log-normal fordelingen

### Lognormal og Normalfordelingen:

En log-normal fordelt variabel  $Y \sim LN(\alpha, \beta^2)$ , kan transformeres til en standard normalfordelt variabel  $Z$  ved

$$Z = \frac{\ln(Y) - \alpha}{\beta}$$

dvs.

$$Z \sim N(0, 1^2)$$

## Eksponentialfordelingen

### Skrivemåde:

$$X \sim Exp(\lambda)$$

### Tæthedsfunktionen

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & \text{ellers} \end{cases}$$

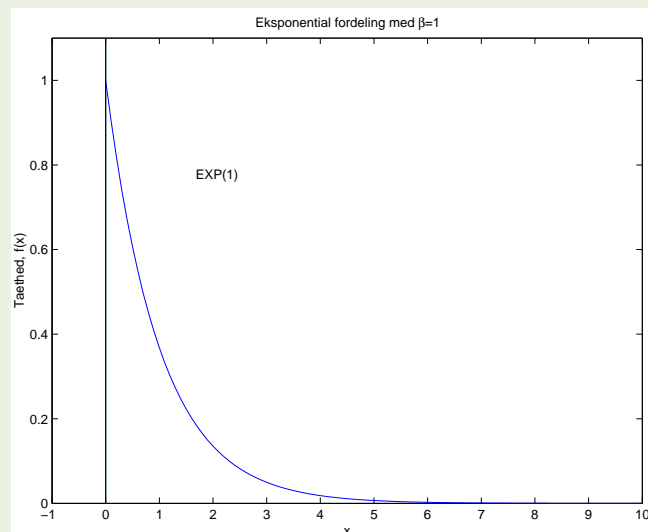
### Middelværdi

$$\mu = \frac{1}{\lambda}$$

### Varians

$$\sigma^2 = \frac{1}{\lambda^2}$$

## Eksempel: Eksponentialfordelingen



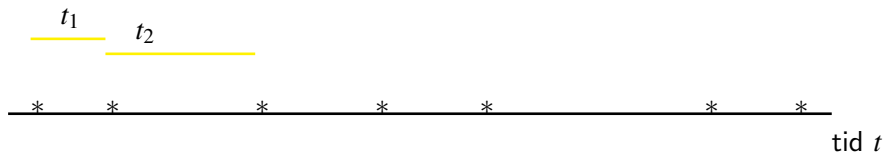
## Eksponentialfordelingen

- Eksponentialfordelingen er et special tilfælde af Gammafordelingen
- Eksponentialfordelingen anvendes f.eks. til at beskrive levetider og ventetider
- Eksponentialfordelingen kan bruges til at beskrive (vente)tiden mellem hændelser i poissonproces

## Sammenhæng mellem eksponential- og poissonfordelingen

Poisson: Diskrete hændelser pr. enhed

Eksponential: Kontinuert afstand mellem hændelser



## Regneregler for lineær funktion af et $X$

Hvis:

- $X$  er en stokastisk variabel
- Vi antager at  $a$  og  $b$  er konstanter

Da gælder (gælder BÅDE kontinuert og diskret):

Middelværdi-regel:

$$E(aX + b) = aE(X) + b$$

Varians-regel:

$$V(aX + b) = a^2 V(X)$$

## Eksempel: Eksponentielfordeling

Kø-model - poissonproces

Tiden mellem kundeankomster på et posthus er eksponentialfordelt med middelværdi  $\mu = 2$  minutter.

Spørgsmål:

En kunde er netop ankommet. Beregn sandsynligheden for at der ikke kommer flere kunder indefor en periode på 2 minutter vha. poissonfordelingen

Svar:

Eksponentialfordeling:  $X_{\text{exp}} \sim \text{Exp}(\lambda)$  with  $\lambda = \frac{1}{\mu} = \frac{1}{2} \frac{1}{\text{min}}$ , find  $P(X_{\text{exp}} > 2)$ .  
Brug Poissonfordeling:  $\lambda_{2\text{min}} = 1$ , find  $P(X_{\text{pois}} = 0)$ .

```
1-pexp(q=2, rate=1/2); dpois(x=0, lambda=1)
```

[1] 0.37 [1] 0.37

## Eksempel: Regneregler for lineær funktion af et $X$

$X$  er en stokastisk variabel

En stokastisk variabel  $X$  har middelværdi 4 og varians 6.

Spørgsmål:

Beregn middelværdi og varians for  $Y = -3X + 2$

Svar:

$$E(Y) = -3E(X) + 2 = -3 \cdot 4 + 2 = -10$$

$$V(Y) = (-3)^2 V(X) = 9 \cdot 6 = 54$$

## Regneregler for lineær funktion af flere $X$ 'er

Hvis:

- $X_1, \dots, X_n$  er stokastiske variable

Da gælder (når de er uafhængige) (gælder BÅDE kontinuert og diskret):

Middelværdi-regel:

$$E(a_1X_1 + a_2X_2 + \dots + a_nX_n) = a_1E(X_1) + a_2E(X_2) + \dots + a_nE(X_n)$$

Varians-regel:

$$V(a_1X_1 + a_2X_2 + \dots + a_nX_n) = a_1^2V(X_1) + \dots + a_n^2V(X_n)$$

## Eksempel: Regneregler for lineær funktion af flere $X$ 'er

### Flypassager-planlægning

Vægten af en passagerer på fly på en strækning antages normalfordelt  $X \sim N(70, 10^2)$ .

Et fly, der kan tage 55 passagerer, må max. lastes med 4000 kg (kun passageres vægt betragtes som last).

### Spørgsmål:

Beregn sandsynligheden for at flyet bliver overlastet.

Hvad er den samlede passagervægt  $Y$  på en afgang?

A:  $Y = 55 \cdot X$     B:  $Y = \sum_{i=1}^{55} X_i$     C:  $Y = 55 + X$     D: Ej A, B eller C

Svar B:  $Y = \sum_{i=1}^{55} X_i$ , det er summen af 55 forskellige passagerer.

## Eksempel: Regneregler 3

Hvad er den samlede passagervægt  $Y$  på en afgang?

$Y = \sum_{i=1}^{55} X_i$ , hvor  $X_i \sim N(70, 10^2)$

Middelværdi og varians for  $Y$ :

$$E(Y) = \sum_{i=1}^{55} E(X_i) = \sum_{i=1}^{55} 70 = 55 \cdot 70 = 3850$$

$$V(Y) = \sum_{i=1}^{55} V(X_i) = \sum_{i=1}^{55} 100 = 55 \cdot 100 = 5500$$

Bruger normalfordeling for  $Y$ :

```
1-pnorm(4000, mean = 3850, sd = sqrt(5500))
```

```
[1] 0.022
```

## Eksempel: Regneregler 3 - FORKERT ANALYSE

Hvad er  $Y$ ?

I hvert fald IKKE:  $Y = 55 \cdot X$  !!!!!

Middelværdi og varians for  $Y$ :

$$E(Y) = 55 \cdot 70 = 3850$$

$$V(Y) = 55^2 V(X) = 55^2 \cdot 100 = 550^2 = 302500$$

Bruger normalfordeling for  $Y$ :

```
1-pnorm(4000, mean = 3850, sd = 550)
```

```
[1] 0.39
```

Konsekvens af forkert beregning:

MANGE spildte penge for flyselskabet!!!

## Oversigt

- 1 Kontinuerte Stokastiske variable og fordelinger
  - Tæthedsfunktion
  - Fordelingsfunktion
  - Middelværdi af en kontinuert stokastisk variabel
  - Varians af en kontinuert stokastisk variabel
- 2 Konkrete Statistiske fordelinger
  - Kontinuerte fordelinger i R
  - Uniform fordeling
  - Normalfordelingen
  - Log-Normal fordelingen
- 3 Eksponentialfordelingen
- 4 Regneregler for middelværdi og varians