

Zelfstudieopdracht - R, RStudio en RMarkdown

In deze zelfstudieopgave raak je bekend met R en R Markdown. Je hoeft de opdracht niet in te leveren.

1. Van start gaan

Alvorens te beginnen, moet je een RStudio Cloud account hebben gemaakt, of R en RStudio gedownload en geïnstalleerd hebben op je eigen computer. Zie Brightspace voor instructies.

Open nu het bestand `Introductie_R.Rmd` in RStudio, als je dat nog niet gedaan had. Als je vervolgens op “Knit PDF” in de balk boven de file klikt (misschien staat er aanvankelijk “Knit HMTL”; je kunt “PDF” kiezen met de muis), zet RStudio de `.Rmd` om in een `.pdf` file. Dit is precies de file die je nu leest.

2. Voorbeelden

Ter introductie een aantal voorbeelden. Voor de meeste bekende verdelingen is in R alvast de *dichtheidsfunctie* f , de *cumulatieve verdelingsfunctie* $x \mapsto F(x) = \mathbb{P}(X \leq x)$, en de *kwantiel functie* $\alpha \mapsto F^{-1}(\alpha)$ geïmplementeerd. Het betreffende commando begint respectievelijk met een `d`, `p` of `q`. We illustreren dit nu voor een normale verdeling met parameters $\mu = 2$ en $\sigma = 3$. Het eerste argument is altijd het argument van de functie, gevolgd door de parameters. Hieronder staat wat R code in grijze blokken, zogeheten ‘chunks’.

```
dnorm(0.5, 2, 3)
```

geeft de dichtheid in $x = 0.5$, en is dus gelijk aan

$$\frac{1}{\sqrt{2\pi}3^2} e^{-\frac{(0.5-2)^2}{2 \cdot 3^2}}.$$

Inderdaad vinden we dat

```
dnorm(0.5, 2, 3)-1/sqrt(2*pi*3^2)*exp(-(0.5-2)^2/(2*3^2))
```

```
## [1] 0
```

```
pnorm(0.5, 2, 3)
```

geeft $\mathbb{P}(X \leq 0.5)$ voor $X \sim \mathcal{N}(2, 9)$.

```
pnorm(0.5, 2, 3, lower.tail = FALSE)
```

geeft $\mathbb{P}(X \geq 0.5)$, wat gelijk is aan $1 - \mathbb{P}(X \leq 0.5)$.

```
qnorm(0.5, 2, 3)
```

geeft de waarde van x zodanig dat $\mathbb{P}(X \leq x) = 0.5$. In dit geval is deze waarde gelijk aan de verwachtingswaarde, dus twee.

Je kunt een steekproef genereren door `r` voor de naam van de verdeling te zetten, bijvoorbeeld

```
x <- runif(10,0,1)
```

maakt een steekproef met tien trekkingen uniform uit het interval $[0, 1]$, of

```
x <- rnorm(100, 2, 3)
```

maakt een steekproef met honderd trekkingen uit een normale verdeling met verwachting $\mu = 2$ en standaarddeviatie $\sigma = 3$. De steekproef bevindt zich nu in de vector `x`, en je kan de waarden zien door te typen

```
x
```

```
## [1] 4.80799107 7.32706805 3.85181002 3.47921843 3.65845455 4.03855585
## [7] 5.66998761 4.75129936 2.27009072 3.04375429 2.75397796 3.82493953
## [13] 0.24668452 5.82979019 5.66511103 -2.47634871 4.88882701 9.57463334
## [19] -0.24654457 -0.89053789 2.56086554 1.61766816 6.39481792 1.66718895
## [25] 7.76172448 5.87058564 5.56300339 1.15348975 1.98718823 1.72181738
## [31] 6.10855702 -3.57776246 -2.43479668 -3.30102298 4.96408563 -1.83746937
## [37] 3.30400907 4.60786943 3.40118317 3.34903535 4.92598638 -0.41208452
## [43] 4.73821008 -0.52275856 -0.26098306 3.93098777 -3.75837159 1.51510659
## [49] 2.28667952 0.07370292 -0.04951375 5.47232615 2.26717825 2.24775771
## [55] -3.88493541 -4.10555973 -2.40856218 5.22043551 4.07284134 -0.53272774
## [61] 3.98923155 0.48014647 -0.73518242 -1.24730685 3.49031471 2.58427422
## [67] 2.86096105 9.00202475 -2.81555068 1.37406561 2.43895290 4.58157094
## [73] -6.25568878 3.51099441 -0.79072606 1.96939352 6.75190869 2.06184837
## [79] 8.75545006 1.85339434 1.81048213 -1.33529272 -4.38373399 3.13233844
## [85] 4.59423829 6.18167006 -2.58379287 4.25455194 3.80754836 0.13835391
## [91] 2.54645711 4.03310341 -1.47370438 0.57059225 5.32282519 6.16530021
## [97] 5.26363112 6.71899558 5.41031527 -1.51491556
```

wat gegevens verkrijgen door te typen

```
summary(x)
```

```
##      Min.   1st Qu.   Median     Mean 3rd Qu.     Max.
## -6.25569 -0.09877  2.66913  2.36286  4.76547  9.57463
```

en een histogram maken door te typen

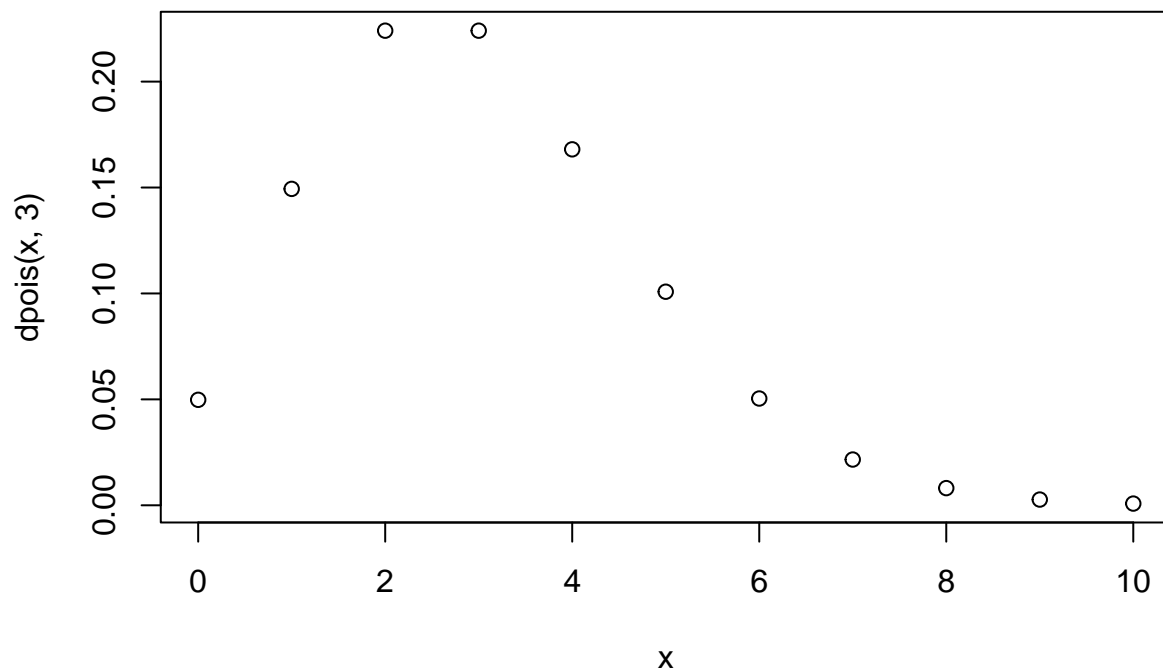
```
hist(x)
```

Met R kun je ook plots van functies maken. Je kunt bijvoorbeeld eerst een vector `x` maken die de waarden $0, 1, \dots, 10$ bevat en dan een plot van de Poisson dichtheid met parameter $\lambda = 3$ met

```
x <- seq(0, 10, by = 1)
class(x)
```

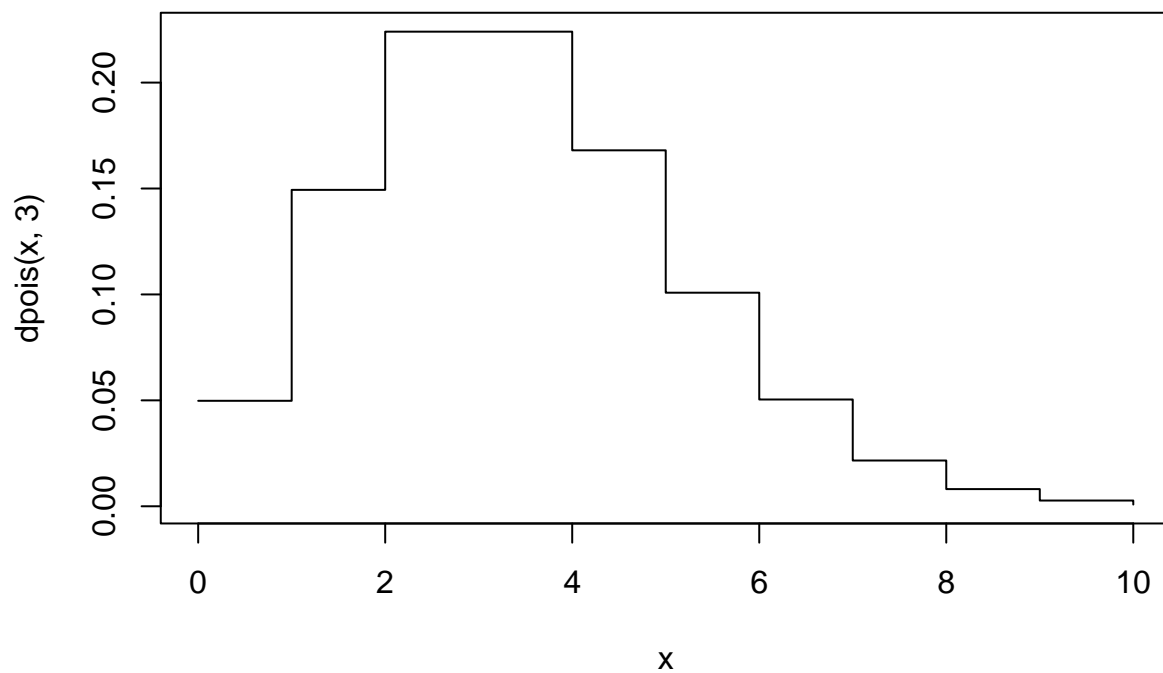
```
## [1] "numeric"
```

```
plot(x, dpois(x, 3), type = "p")
```



Of misschien is deze mooier?

```
plot(x, dpois(x,3), type = "s")
```



Tot slot laten we zien hoe je een barplot kunt maken. Eerst maken we een steekproef met 100 waarnemingen uit een binomiale verdeling.

```
binom <- rbinom(100, size = 5, p = 0.6)
binom
```

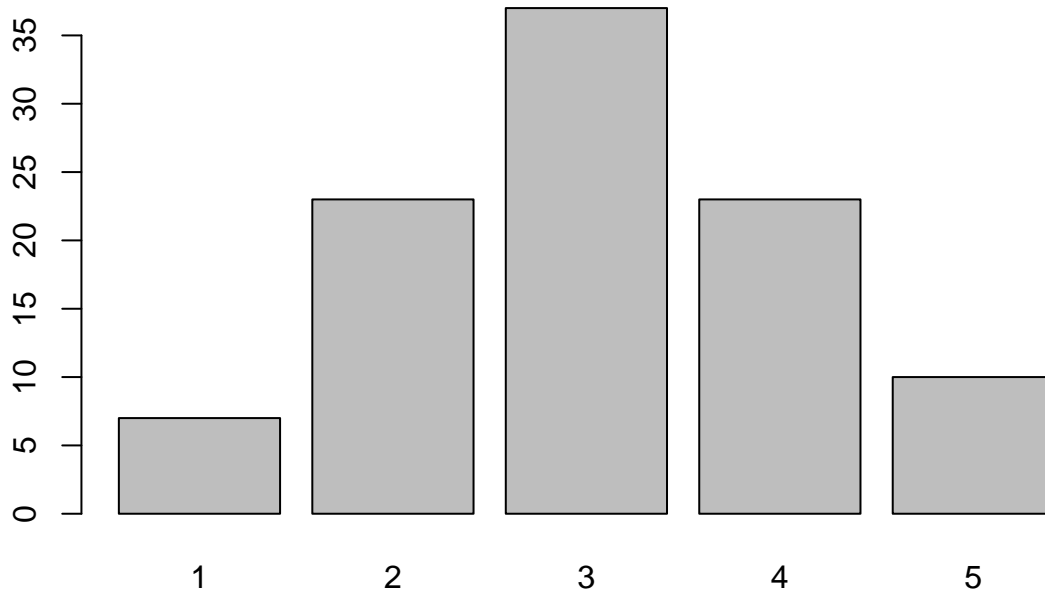
```
## [1] 1 3 4 3 4 4 2 3 4 2 2 4 3 3 3 4 3 3 5 3 3 1 1 3 5 2 3 3 2 3 4 3 4 5 5 3 4
## [38] 5 3 3 3 4 4 2 4 4 4 3 2 1 3 2 5 2 4 4 5 2 3 4 2 3 3 3 2 3 3 2 3 3 1 2 3 5
## [75] 4 2 3 2 4 4 4 5 2 2 3 3 2 5 2 4 2 3 1 1 2 3 3 4 3 2
```

```
table(binom)
```

```
## binom  
## 1 2 3 4 5  
## 7 23 37 23 10
```

De tabel geeft aan hoe vaak elke uitkomst is geobserveerd. We visualiseren dit met een barplot.

```
barplot(table(binom))
```



3. Vragen

Typ je antwoorden op onderstaande vragen onder elke vraag in de .Rmd-file. Niet alle benodigde informatie staat in deze opdracht gegeven, voor een aantal dingen zul je op internet op zoek moeten. Je zult zien dat er veel online hulp voor R beschikbaar is. Daarnaast kan het nuttig zijn om de documentatie te bekijken van R functies die je gaat gebruiken. Typ bijvoorbeeld

```
help("hist")
```

of

```
?hist
```

om meer over `hist()` te lezen.

Vraag 1 De toevoeging ‘eval=FALSE’ die in enkele van bovenstaande chunks staat, is een zogeheten *chunk option*. Wat is het effect van het toevoegen van ‘eval=FALSE’ tussen de chunk options?

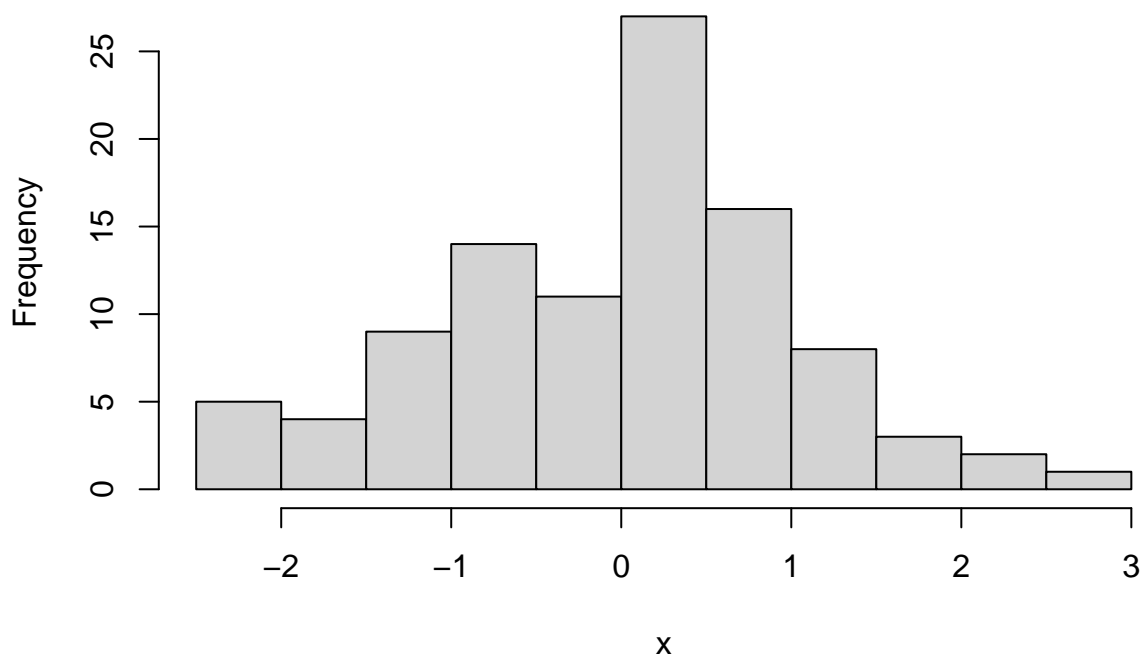
Antwoord: Het zorgt ervoor dat de chunks niet uitgevoerd worden bij het construeren van de PDF, zo kunnen we ze zien in de PDF (en in de source editor).

Vraag 2 Zoek op welke chunk options er verder nog zijn, en beschrijf er twee die je handig lijken.

Antwoord: Echo zorgt voor text wel/niet in het document (de PDF) en wel in de output. Met Result kan je zelfs iets wel in de PDF, maar niet in de output hebben (zeer handig!).

Vraag 3 Voeg hieronder een nieuwe chunk in, maak een histogram van 500 trekkingen uit een standaard normale verdeling, en zorg dat in het geknitte document de R-code niet zichtbaar is, maar het histogram wel.

Histogram of x



Vraag 4 Maak drie vectoren met 100 trekkingen uit een normale verdeling met verwachtingswaarde gelijk aan respectievelijk 1, 2, en 3, en standaarddeviatie ook gelijk aan respectievelijk 1, 2 en 3. Bereken het gemiddelde en de variantie van elk van de vectoren.

Vraag 5 Wat is de theoretische verdeling van de som van de drie vectoren uit de vorige vraag?

Vraag 6 Bereken nu met R het gemiddelde en de variantie van de som van de drie vectoren uit vraag 5.

Vraag 7 Maak een vector van lengte 50, gevuld met alleen maar enen. Gebruik hiervoor `rep()`.

Vraag 8 Maak een vector `x` van lengte 50, gevuld met de getallen 1 t/m 50, en een vector `y` die gelijk is aan tweemaal `x`. Maak met `plot()` een scatterplot van `x` en `y`.

Vraag 9 Als vraag 8, maar trek nu een lijn door de punten heen.

Vraag 10 Bereken de correlatie tussen `x` en `y` met behulp van R.

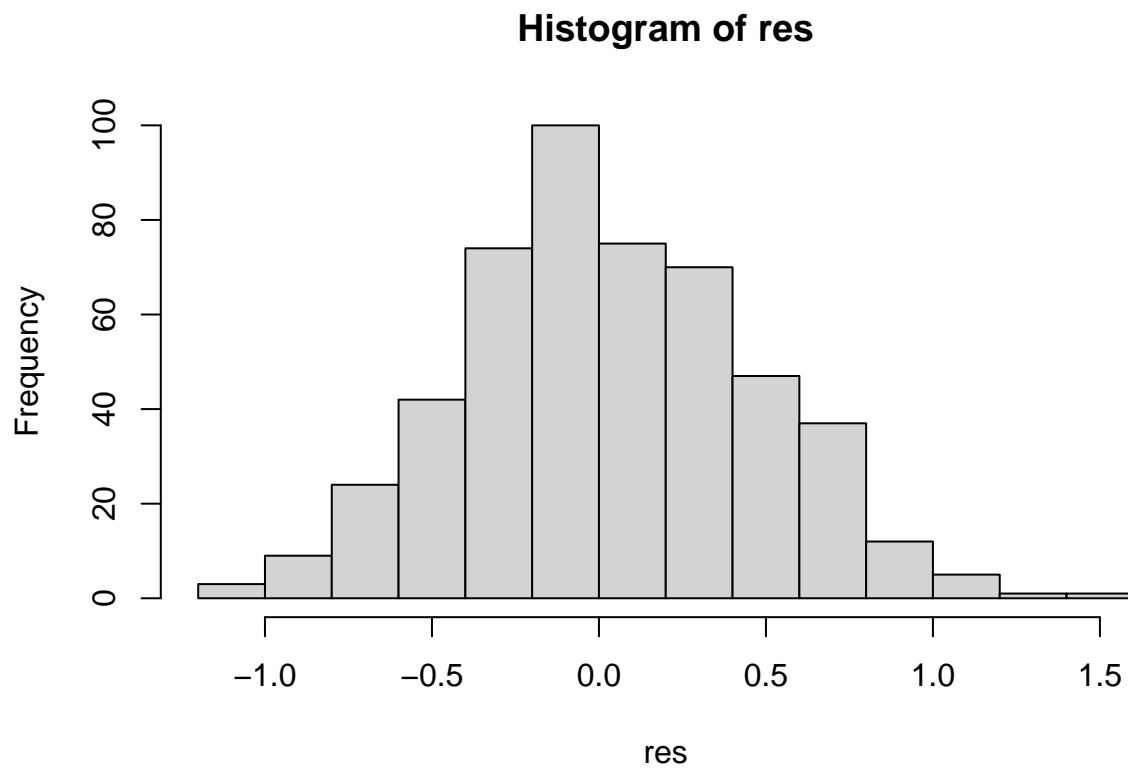
Vraag 11 Als vraag 8, maar tel nu een vector standaard normaal verdeelde variabelen op bij `y`. Bereken tevens de correlatie tussen `x` en de nieuwe vector. Waarom is het antwoord anders dan bij vraag 9?

Vraag 12 Beschrijf in woorden wat er in onderstaande code gebeurt.

```
N <- 500
n <- 5
res <- rep(0, N)

for(i in 1:N){
  vector <- rnorm(n)
  res[i] <- mean(vector)
}

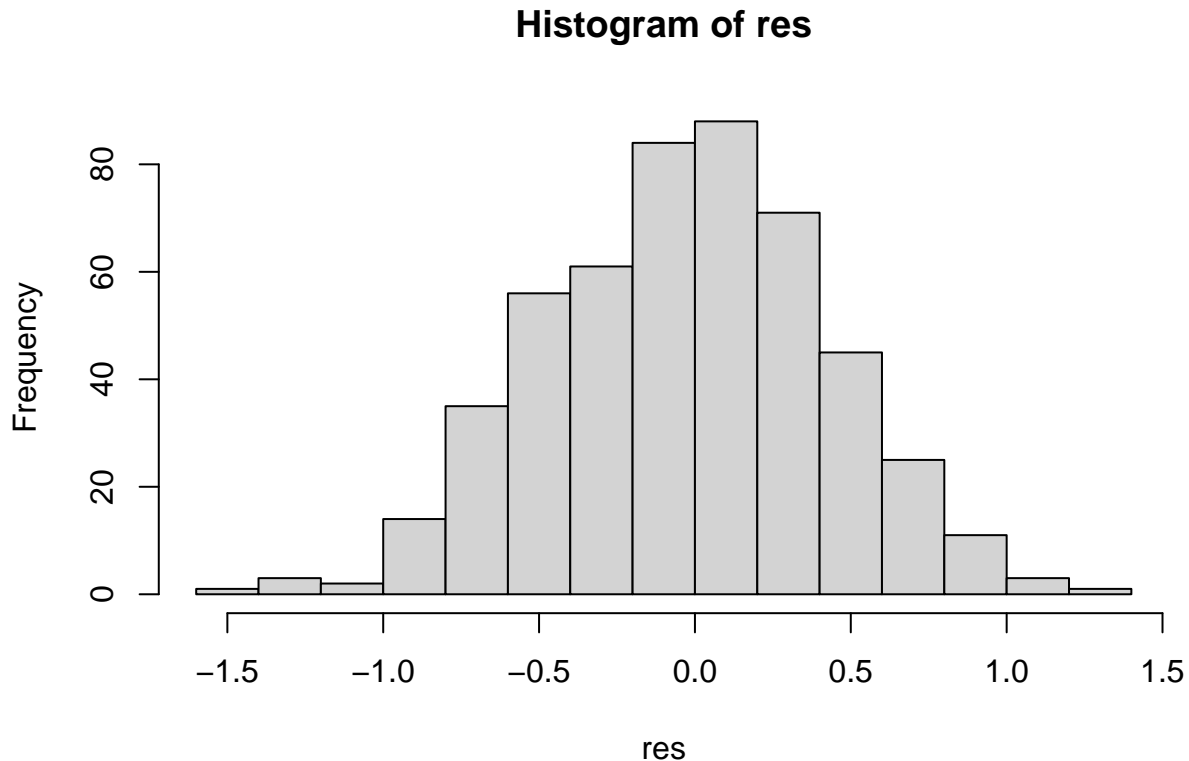
hist(res)
```



NB Een andere manier om hetzelfde te bereiken is:

```
N <- 500
n <- 5

obs <- matrix(rnorm(N*n), nrow = N)
res <- apply(obs, 1, mean)
hist(res)
```



Vraag 13 Maak het histogram uit bovenstaande vraag opnieuw, nu zo dat het histogram precies 20 bins heeft.

Vraag 14 Beschrijf wat in `pnorm()` het effect is van de keuze `lower.tail = TRUE` of juist `lower.tail = FALSE`.

NB In de help is te zien dat `lower.tail = TRUE` de 'default' is. Als je die optie wil, hoef je dat niet expliciet te vermelden wanneer je de functie aanroept.

Vraag 15 Zij $q_{0.3}$ de waarde zodanig dat de kans dat een standaard normaal verdeelde variabele kleiner dan $q_{0.3}$ is, gelijk is aan 0.3. Gebruik `qnorm()` om de waarde $q_{0.3}$ te vinden.

Vraag 16 We maken een vector met 20 uniform verdeelde waarnemingen.

```
obs <- runif(20)
obs
```

```
## [1] 0.9730347 0.9823020 0.2460618 0.5753167 0.6238989 0.8007011 0.3894284
## [8] 0.2826525 0.5032696 0.7499731 0.6125626 0.2836830 0.5929346 0.8181316
## [15] 0.6225257 0.4668239 0.7493313 0.7505331 0.3392878 0.3698955
```

Bekijk wat de volgende functies doen:

```
min(obs)
```

```
## [1] 0.2460618
```

```
max(obs)
```

```
## [1] 0.982302
```

```
which.min(obs)
```

```
## [1] 3
```

```
which.max(obs)
```

```
## [1] 2
```

```
obs[which.min(obs)]
```

```
## [1] 0.2460618
```

```
obs[which.max(obs)]
```

```
## [1] 0.982302
```

Gebruik nu de functie `which()` om de indices te vinden van de waarnemingen die groter dan 0.5 zijn. Gebruik deze indices vervolgens om de waarnemingen die groter dan 0.5 zijn te selecteren.