

# TBDA01: Laboration 2

Måns Magnusson, Mattias Villani

20 september 2017

---

## Instruktioner

- Laborationen ska göras **två och två**
  - Labben ska vara en **PDF-rapport** med kod, analys och grafer. I rapporten ska följande ingå:
    - Båda studenternas namn och LiU-id.
    - Laborationsnummer
    - Uppgifterna ni besvarar (ex. som rubriker).
  - Ett tips är att använda R markdown. [Här](#) finns en R markdownmall att utgå ifrån. Antingen kan ni skapa PDF direkt från denna mall i R-Studio (med TeX) eller så skapar ni ett Word/HTML dokument som sedan skrivs ut till PDF.
  - Deadlinen för laborationen framgår av [kurshemsidan](#).
  - Laborationsrapporten skickas in via LISAM.
  - För att få godkänt på laborationen krävs godkänt på **75%** av uppgifterna.
-

# Innehåll

<b>1</b>	<b>Introduktion till R</b>	<b>3</b>
<b>2</b>	<b>Laboration</b>	<b>4</b>
2.1	Likelihoodfunktionen . . . . .	4
2.2	Punktskattningar från populationsmodell . . . . .	5
2.3	Estimatorers samplingfördelning . . . . .	6
2.4	Numerisk maximum likelihood estimation . . . . .	7

# Kapitel 1

## Introduktion till R

R är ett programmeringspråk för statistisk programmering som påminner mycket om Matlab. R bygger på öppen källkod och kan laddas ned [här](#). R-Studio är en mycket populär IDE för R (som också påminner mycket om Matlab). Denna IDE finns att tillgå [här](#). I R-Studio finns funktionalitet för literate programming med R markdown implementerat för att kombinera R kod med markdownsyntax. På detta sätt är det enkelt att generera rapporter med både text, grafik och kod. Det är R:s motsvarighet till Python Notebook.

För en ingång till R från andra språk kan onlineboken *Advanced R* rekommenderas som finns [här](#). Kapitlen *Data structures*, *Subsetting* och *Functions* bör ge en snabb introduktion.

Även boken *The art of R programming* av Norman Matloff kan vara till hjälp som referenslitteratur. Boken finns [här](#).

### Videomaterial

- För en introduktion till syntaxen i R se Google developers R videomaterial [här](#).
- Mer (detaljerat) videomaterial av Roger Peng finns att tillgå [här](#).
- För att visualisera med basgrafiken finns [följande](#) introduktionsvideo.
- För mer komplicerad grafik rekommenderas *ggplot2*-paketet. En introduktionsvideo finns [här](#).
- En introduktion till R markdown finns [här](#).

### Cheatsheets

- *R reference card v.2* av Matt Baggot med vanliga funktioner i R finns att tillgå [här](#).
- *R markdown cheatsheet* av R-Studio med tips för R markdown finns att tillgå [här](#).

## Kapitel 2

# Laboration

I denna laboration kommer vi gå djupare in på statistisk inferens. D.v.s. givet att vi observerat data, vilka slutsatser kan vi dra om parametrarna i våra sannolikhetsmodeller. När vi gör statistisk inferens kommer vi (oftast) att behöva göra **antaganden** om från vilken modell vårt datamaterial kommer. Det är därför nyttigt att ha med sig följande citat av George Box.

”All models are wrong, but some are useful.” - George E. P. Box

**Obs!** De flesta beräkningar i denna laboration är avrundade. Du kan mycket väl få något annorlunda resultat (men som avrundat bör ge samma resultat som i exemplen).

### 2.1 Likelihoodfunktionen

Vi ska nu visualisera likelihoodfunktionen. Nedan är 10 respektive 100 dragningar från en  $\text{Gamma}(\alpha = 4, \beta = 1)$ , vilket i detta exempel är vår ”sanna” fördelning. Observera att vi sätter slumpfröet till 4711. Det gör att vi kan replikera dessa dragningar exakt.

```
> set.seed(4711)
> x1 <- rgamma(n = 10, shape = 4, scale = 1)
> x1

[1] 7.7323 6.5334 6.0959 2.7802 1.3565 5.1966 1.9279 5.4122 8.5859 6.4369

> x2 <- rgamma(n = 100, shape = 4, scale = 1)
> x2

[1] 4.79331 6.40788 1.25250 7.06992 2.67750 1.74628 5.69628 4.95335
[9] 3.64830 5.56426 3.02685 5.63416 4.69548 3.87418 4.72661 4.91524
[17] 8.29348 6.57035 7.10917 4.31000 6.82118 3.66030 1.03559 9.27426
[25] 2.56627 3.50256 3.08015 4.86450 4.30884 6.73306 2.70627 4.57985
[33] 2.67979 2.27968 2.00768 1.40964 4.99229 5.48001 4.51633 5.14000
[41] 3.77294 3.45746 2.85032 2.46253 3.36590 3.49688 2.41453 5.07472
[49] 3.81785 3.65560 5.87285 4.31166 9.40850 8.03080 3.48122 5.37745
[57] 7.74533 8.56121 0.71804 4.34067 3.02849 5.75616 1.44207 3.84525
[65] 5.10882 2.07482 3.58257 0.94767 2.44326 4.97860 3.52383 1.69256
[73] 9.03576 3.61321 2.45790 2.94222 4.97442 1.56972 2.01078 7.01523
[81] 6.40548 1.39056 2.86646 1.20835 1.59489 4.04473 7.87768 2.75054
[89] 5.82472 3.32521 2.56622 2.55995 5.58725 4.63006 4.74396 3.84933
[97] 1.64613 4.71465 2.28516 3.08871
```

Vi ska nu använda log-likelihoodfunktionen. För att härleda en log-likelihoodfunktion utgår vi ofta från vår sannolikhetsfördelnings täthetsfunktion.

Nedan är ett exempel på log-likelihoodfunktionen för gammafördelade variabler och en vektor  $\mathbf{x}$  med data. Vi antar här att observationerna är **iid** (independent and identically distributed).

$$l(\alpha, \beta | \mathbf{x}) = \ln \left( \prod_{i=1}^n \frac{\beta^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\beta x_i} \right) \quad (2.1)$$

$$= n \cdot [\alpha \cdot \ln \beta - \ln \Gamma(\alpha)] + (\alpha - 1) \sum_{i=1}^n \ln(x_i) - \beta \sum_{i=1}^n x_i \quad (2.2)$$

### Uppgift 1 Likelihoodfunktioner

**a)** Skapa en funktion i R du kallar `llgamma(x, alpha, beta)` som tar datapunkterna i en vektor `x`, värden för  $\alpha$  och  $\beta$  och returnerar log-likelihoodvärdet för parametrarna givet det observerade data `x`. Inkludera funktionen i rapporten.

**Tips!**  $\ln \Gamma(x)$  heter `lgamma()` och  $\ln$  heter `log()` i R.

Nedan är ett exempel på värden funktionen ska returnera.

```
> llgamma(x = x1, alpha = 2, beta = 2)
[1] -75.19
```

**b)** Beräkna och visualisera loglikelihood värden för  $\mathbf{x}_1$  och  $\mathbf{x}_2$  som simulerades ovan då  $\alpha = 4$ . Du behöver beräkna log-likelihooden för olika värden på  $\beta$  på intervallet  $(0, 3]$ , förslagsvis med en stegstorlek på 0.01. För vilket av de upprepade värdena för  $\beta$  får du det maximala värdet på log-likelihoodfunktionen? Visualisera log-likelihoodfunktionen för de olika värdena på  $\beta$ .

**c)** Upprepa det du gjort i **b)** ovan men anta nu att  $\beta = 1$  och visualisera på ett liknande sätt  $\alpha$  på intervallet  $(0, 10]$  för  $\mathbf{x}_1$  och  $\mathbf{x}_2$ . Vilket värde på  $\alpha$  ger det maximala värdet för log-likelihood-funktionen?

**d)** Härled nu log-likelihoodfunktionen för en vektor `y` med normalfördelade datapunkter och inkludera härledningen i rapporten.

Implementera log-likelihoodfunktionen som en funktion i R som du kalla `llnorm(x, mu, sigma2)`. Funktionen ska kunna ta en vektor med datapunkter `x`, samt värden på  $\mu$  och  $\sigma$ .

Inkludera funktionen i rapporten.

Nedan är ett exempel på hur funktionen ska fungera.

```
> llnormal(x = x1, mu = 2, sigma2 = 1)
[1] -87.257
```

**Tips!**  $\pi$  är inkluderat som `pi` i R.

**e)** Visualisera på samma sätt som i **b)** och **c)** ovan log-likelihood funktionen för  $\mu$  på intervallet  $[0, 10]$  för de (gammafördelade) dragningarna  $\mathbf{x}_1$  och  $\mathbf{x}_2$  ovan. Vi antar att  $\sigma^2 = 1$ . Visualisera gammafördelnings täthetsfunktion (med  $\alpha$  och  $\beta$  som maximerade log likelihooden i **b)** och **c)**) och normalfördelningens täthetsfunktion (med  $\mu$  som maximerade log likelihood funktionen i **d)**). Visualisera också datamaterialet (som ett histogram med `histogram()` **Tips!** Se labb 1). Vilken modell, Normalfördelningen eller Gammafördelningen, tycker du passar datamaterialet bäst?

## 2.2 Punktskattningar från populationsmodell

Vi har sett ovan att likelihoodfunktionen är en funktion som innehåller information om de parametrar vi är intresserade av. Vi vill då uppskatta dessa parametrar baserat på den data vi har samlat in. En av de vanligaste metoderna för att göra detta är med vad som brukar kallas **maximum likelihood estimation** (MLE), som namnet antyder handlar det om att maximera likelihoodfunktionen vi studerade ovan. Är vi intresserade av att försöka uppskatta parametrarna i vår modell med ett enda värde per parameter är vi intresserad av en **punktskattning** av modellens parametrar. Det är också det vi ska studera i detta avsnitt.

För gammafördelningens parameter  $\beta$  går det att analytiskt härleda maximum för loglikelihoodfunktionen med avseende på parametern  $\beta$ . Detta kan göras på traditionellt sett genom att derivera log likelihoodfunktionen, sätta derivatan till 0 och lösa ut  $\beta$ . En punktskattning anges ofta med ett litet "tak" över den parameter vi är intresserad av för att indikera att det är en skattning av parametern  $\beta$ .

$$\hat{\beta}_{MLE} = n\alpha \left( \sum_{i=1}^n x_i \right)^{-1} \quad (2.3)$$

Funktionen för att beräkna vår punktskattning för ett antal datapunkter är vad som kallas en **estimator**. Det värde funktionen returnerar kallas för **estimat** eller **skattning**.

### Uppgift 2 Punktskattning med MLE i en gammafördelning

Implementera estimatoren ovan som en funktion du kallar `gamma_beta_mle(x, alpha)` med parametrarna `x` (data) och `alpha`. Nedan är ett exempel på hur funktionen ska fungera:

```
> gamma_beta_mle(x = x1, alpha=2)
```

```
[1] 0.38419
```

Skatta  $\hat{\beta}_{MLE}$  med denna estimator för  $\mathbf{x}_1$  och  $\mathbf{x}_2$  ovan med  $\alpha = 4$ . Vad är dina slutsatser?

På ett liknande sätt är det möjligt att skatta parametrarna  $\mu$  och  $\sigma^2$  i en normalfördelning.

### Uppgift 3 Punktskattning med MLE i en normalfördelning

a) MLE-estimatorn för  $\mu$  och  $\sigma^2$  i en normalfördelning finns [här](#). Implementera dessa estimatorer som `norm_mu_mle(x)` och `norm_sigma2_mle(x)` med `x` (data) som argument. Se nedan för ett exempel på hur funktionerna ska fungera.

```
> test_x <- 1:10
```

```
> norm_mu_mle(x = test_x)
```

```
[1] 5.5
```

```
> norm_sigma2_mle(x = test_x)
```

```
[1] 8.25
```

b) Gör 10 och 10000 dragningar från följande fördelning. Sätt slumpfröet till 42 innan du gör dina dragningar på samma sätt som  $\mathbf{x}_1$  och  $\mathbf{x}_2$  ovan.

$$\mathbf{y} \sim \mathcal{N}(\mu = 10, \sigma^2 = 4)$$

Använd sedan dina två estimatorer för att först skatta  $\mu$  och  $\sigma^2$  sedan. Vad är skillnaden mellan 10 och 10000 dragningar? Vad beror detta på (**Tips!** Föregående laboration)?

**Obs!** Tänk på att `rnorm()` är parametriserad med  $\sigma$ , inte  $\sigma^2$ .

## 2.3 Estimatorers samplingfördelning

Vi ska nu studera egenskaperna hos dessa estimatorer genom simulering, detta är en vanlig metod för att utvärdera hur väl olika **estimatorer** fungerar. Vi ska nu studera **fördelningen** för de estimatorer vi har implementerat ovan.

#### Uppgift 4 Samplingfördelningen för $\hat{\beta}_{MLE}$ , $\hat{\mu}_{MLE}$ och $\hat{\sigma}^2_{MLE}$

Som ett första steg ska vi studera fördelningen för de estimatorer vi har implementerat. Vi gör detta genom att upprepa skattningen för ett antal olika utfall från våra fördelningear och se hur våra estimatorer varierar när våra data varierar. Detta kallas estimatorns **samplingfördelning**. Implementera pseudokoden nedan för att beräkna 2000 MLE-skattningar.

```
for i in 1 to 2000 do
  x1 = draw Gamma(n = 10, alpha= 4, beta = 1)
  x2 = draw Gamma(n = 10000, alpha= 4, beta = 1)
  beta1_mle[i] = gamma_beta_mle(x = x1, alpha = 4)
  beta2_mle[i] = gamma_beta_mle(x = x2, alpha = 4)

  y1 = draw Normal(n = 10, mu = 10, sigma2 = 4)
  y2 = draw Normal(n = 10000, mu = 10, sigma2 = 4)
  mu1[i] = norm_mu_mle(x = y1)
  mu2[i] = norm_mu_mle(x = y2)
  sigma1[i] = norm_sigma2_mle(x = y1)
  sigma2[i] = norm_sigma2_mle(x = y2)
```

a) Visualisera samplingfördelningarna för  $\hat{\beta}_{MLE}$ ,  $\hat{\mu}_{MLE}$  och  $\hat{\sigma}^2_{MLE}$  då  $n = 10$  och då  $n = 10000$  i ett histogram. Vad är dina slutsatser? **Tips!** `hist()`

## 2.4 Numerisk maximum likelihood estimation

De exempel vi arbetat med såhär långt har varit "snälla" situationer där det finns analytiskt härledda estimatorer. Men som vi sett ovan (exempelvis för  $\alpha$  i gammafördelningen) är det långt ifrån självklart att det finns analytiska resultat för att göra en MLE. I lite mer komplicerade fall vill vi därför istället använda oss av numeriska optimerare för att göra vår skattningar. I R finns funktionen `optim()` som är en bra numerisk optimerare som kan användas för MLE oavsett om vi känner likelihoodens gradient och hessian eller inte.

Nedan är ett exempel på hur `optim()` kan användas för att numeriskt finna MLE skattningar av  $\mu$  och  $\sigma^2$  i exemplet ovan. `optim()` kräver att parametrarna som ska optimeras finns i en vektor, likaså finner `optim()` ett minimum varför vi behöver multiplicera log likelihooden med -1. Det ger följande funktion:

```
> llnormal2(par = c(2, 1), x = x1)
```

```
[1] 87.257
```

Med denna funktion kan vi sedan använda `optim()` för att numeriskt finna våra maximum likelihoodskattningar. Jämför med resultaten vi fick ovan med de analytiska lösningarna. Argumentet `par` i `optim()` är de initiala värdena där optimeraren startar. Tänk på att ge initieringsvärden som är definierade i funktionen. För att undvika att `optim()` ska söka i exakt 0 kan man ange `.Machine$double.eps` som lägsta värde (d.v.s. lägsta positiva numeriska värde).

```
> opt_res <- optim(par = c(0,1), fn = llnormal2, x=test_x, method="L-BFGS-B", lower = c(-Inf, .Machine$double.eps))
> opt_res$par
```

```
[1] 5.50 8.25
```



### Uppgift 5 Log-likelihoodfunktionen för betafördelningen

**a)** Härled (eller leta reda på) log-likelihoodfunktionen för Betafördelningen och implementera den som en funktion i R som kan optimeras med `optim()`. Nedan är ett exempel på log-likelihoodfunktionen (multiplicerad med -1) som kan användas i `optim()`.

```
> llbeta(par = c(2, 2), x = c(0.01, 0.5, 0.99))
```

```
[1] 5.2415
```

**b)** Simulera 100 dragningar från en  $\text{Beta}(\alpha = 0.2, \beta = 2)$  och visualisera dragningarna med ett histogram. **Tips!** `hist()`

**c)** Använd `optim()` för att baserat på dessa dragningar och log-likelihoodfunktionen uppskatta parametrarna  $\alpha$  och  $\beta$  i betafördelningen. Tänk på att  $\alpha$  och  $\beta$  är definierade på  $\mathbb{R}^+$  när du anger intervallen och startvärdena till `optim()`.

**d)** Upprepa **b)** och **c)** 2000 ggr (d.v.s. ta fram samplingfördelningen för  $\hat{\alpha}_{MLE}$  och  $\hat{\beta}_{MLE}$ ). Visualisera dessa samplingfördelningar i ett histogram.