



Matematisk Statistik

SF1911 Statistik för bioteknik: HT 2016 lp2  
Lab 0 för CBIOT

## 1 Förbedelser

*Denna lab är till skillnad från de andra inte poänggivande utan är till för den som vill bekanta sig med MATLAB. Fokusera på att lära dig att använda MATLAB och se till att du förstår de kommandon som du använder.*

**Förberedelser:** Innan du går till laborationen, läs igenom den här handledningen, repetera exponentialfördelningen.

Temat för den här datorlaborationen är *simulering*. Sannolighetsteoridelen av kursen handlar om hur man genom beräkningar kan ta fram olika storheter som sannolikheter, väntevärden osv., för en given stokastisk modell. För mer komplicerade system är det ibland inte alls möjligt att göra exakta beräkningar, eller så är det så tidskrävande att man avstår.

I sådana sammanhang kan simulering vara ett alternativ. Simulering innebär att man med hjälp av en dator simulerar ett antal replikeringar av det stokastiska systemet, och sedan använder t.ex. medelvärden eller empiriska kvantiler (mer om det nedan) för att uppskatta de storheter man söker. I den här laborationen skall vi göra detta för några enklare problem men grundprinciperna går att använda på långt mer komplicerade problem som vi inte kan lösa med enkla beräkningar.

I det allra enklaste fallet kan det vara fråga om att uppskatta väntevärdet för en fördelning. Antag att vi har en fördelningsfunktion  $F$ , och låt  $X$  vara en stokastisk variabel med denna fördelningsfunktion. Antag också att vi vill uppskatta tillhörande väntevärde,  $\mu$  säg. Om vi nu drar observationer  $x_1, x_2, \dots, x_n$  från  $F$ , kan vi uppskatta  $\mu$  med hjälp av

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k. \quad (1)$$

Att detta är en rimlig uppskattning följer av stora talens lag. Om vi vill uppskatta  $F(a) = \mathbb{P}(X \leq a)$  för något tal  $a$ , så kan vi göra det

genom att räkna ut hur stor andel av de simulerade observationerna som är  $\leq a$ . Vi kan skriva detta som

$$\hat{F}(a) = \frac{\text{antal } x_k \text{ som är } \leq a}{n} = \frac{1}{n} \sum \mathbb{I}(x_k \leq a); \quad (2)$$

här är  $\mathbb{I}$  en s.k. *indikatorfunktion*, som tar värdet 1 om villkoret inom parentes är uppfyllt, annars 0. Alltså är  $\mathbb{I}(x_k \leq a)$  lika med 1 precis för de  $k$  sådana att  $x_k \leq a$ , så att summan ovan räknar antalet index  $k$  som uppfyller villkoret.

## 2 Introduktion till MATLAB

I samtliga laborationer i den här kursen kommer *MATLAB*, som är ett interaktivt program för numeriska beräkningar, att användas. Det är också ett programmeringsspråk. MATLAB finns på de flesta datorer på KTH, och till fördelarna med programmet hör att det ser i stort likadant ut oberoende av på vilken sorts dator man kör det. Om du vill ha mer att läsa om MATLAB så finns det olika handledningar att ladda ned och köpa.

Börja med att logga in på ditt vanliga konto. Starta sedan MATLAB genom att klicka på ikonen. Programmet avslutas med kommandot `exit`. Till att börja med kan man tänka på MATLAB som en avancerad räknedosa som beräknar uttryck. Man skriver in vad man vill ha gjort och MATLAB svarar.

```
>> 3*11.5 + 2.3^2/4
ans =
    35.8225
```

Variabler tilldelas värden med tecknet `=` och finns sedan kvar i minnet. Prova att tilldela några variabler värden.

```
>> a = 1;
>> b = sqrt(36);
>> width = 3.89;
>> who
```

Your variables are:

```
a          b          width
```

Kommandot `who` visar alltså vilka variabler som finns i minnet. En variabel, t.ex. `b`, kan raderas ur minnet med `clear b`. Lägg ett semikolon, `;`, till efter en kommandorad, skrivs resultatet inte ut på skärmen. MATLAB kan också hantera vektorer och matriser, och de hanteras precis lika enkelt som ovan.

```
>> x = [1 3 7]
```

```
x =
```

```
1 3 7
```

```
>> y = [2 1 8]'
```

```
y =
```

```
2
```

```
1
```

```
8
```

```
>> z = [1 2 ; 3 4]'
```

```
z =
```

```
1 3
```

```
2 4
```

```
>> w = rand(1, 4)
```

```
w =
```

```
[0.2190 0.0470 0.678 0.6793]
```

```
>>
```

Tecknet ' betyder som synes transponat och semikolon används för att skilja rader åt i matriser. Funktionen `rand(m, n)` ger en  $m \times n$ -matris med slumpantal som är likformigtördelade mellan 0 och 1. Notationen för algebraiska uttryck är den vanliga, men kom ihåg att multiplikationstecknet \* tolkas som matrismultiplikation. Elementvis multiplikation mellan tvåmatriser A och B av samma dimensioner skrivs `A.*B`. I MATLAB finns alla vanliga funktioner inbyggda, t.ex.

```
exp log sin asin cos acos tan atan
```

Observera att `log` är den naturliga logaritmen. Prova att plotta en funktion, t.ex. genom följande kommandon.

```
>> x = 0.5:0.1:2
```

```
>> help log
```

```
>> y = log(x)
>> plot(x,y)
```

Den första raden tilldelar  $x$  en vektor som löper från 0.5 till 2 i steg om 0.1. Hjälpfunktionen `help` ger information om en funktion. Skriver du bara `help` visas en lista med tillgängliga funktioner, sorterade efter funktionspaket (ett sådant paket kallas en *toolbox*).

I MATLAB finns det en stor mängd funktioner som har att göra med sannolikhetssteori och statistik. Se `help stats`. Titta speciellt under rubriken *Random Number Generators*, som kommer i början på den långa listan av funktioner.

### 3 Väntevärde av exponentialfördelning

MATLABs funktion för att simulera exponentialfördelade stokastiska variabler heter `exprnd`. Använd gärna MATLABs hjälpkommando `help` för att ta reda på precis hur funktionen `exprnd` fungerar, dvs. skriv `help exprnd!` Funktionen kan också användas för att simulera vektorer (eller matriser) av oberoende exponentialfördelade variabler. T.ex. ger

```
>> n = 1000;
>> x = exprnd(2.5, n, 1);
```

en  $n \times 1$ -vektor av värden från exponentialfördelningen med väntevärde 2.5. Antag att vi nu inte vet att denna fördelning har väntevärdet just 2.5, och att vi vill uppskatta detta från våra simulerade data. Det kan vi göra genom att beräkna medelvärdet.

```
>> mean(x)
```

Hur bra blir din uppskattning? Pröva att göra om simuleringen och medelvärdesberäkningen några gånger. Pröva också olika värden på  $n$ !

### 4 Svanssannolikheter för exponentialfördelning

Vi ska nu använda simulering för att uppskatta svanssannolikheten

$$\mathbb{P}(X > x) = 1 - F_X(x) \quad (3)$$

för en exponentialfördelning. Syntaxen `(x>5)` i MATLAB ger en vektor av samma storlek som  $x$ , men där ett element i vektorn är 1 eller 0 beroende på om motsvarande element i  $x$  uppfyller villkoret  $> 5$  eller inte.

Du kan alltså simulera en vektor  $x$  av data som ovan, och sedan skriva

```
>> mean(x>5)
```

för att beräkna skattningen  $1 - \hat{F}(5)$ . Vad får du för värde? Vad är den sanna svanssannolikheten? Prova med olika  $n$ , och olika svanssannolikheter (dvs byt ut 5 mot något annat)!