



Matematisk Statistik

SF1910 Tillämpad statistik, HT 2016  
Laboration 1 för CSAMHS, CLGYM-TEMI

## Introduktion

Detta är handledningen till Laboration 1, ta med en en utskriven kopia av den till laborationen. Läs handledningen två gånger. Försäkra dig om att du förstår hur de MATLAB-kommandon som finns i den bifogade koden fungerar. Laborationen bedöms som godkänd eller ej godkänd. För att få delta i laborationen skall svar på förberedelseuppgifter kunna redovisas **individuellt**. Arbete i grupp är tillåtet (och uppmuntras) med **högst två** personer per grupp. Godkänd laboration ger 3 poäng på ordinarie tentamenstillfälle.

## Förberedelseuppgifter

1. Definiera begreppen sannolikhet, fördelningsfunktion, och täthetsfunktion. Skriv upp sambandet mellan dem.
2. Den stokastiska variabeln  $X$  har täthetsfunktion

$$f_X(x) = \lambda e^{-\frac{x}{\lambda}} + \frac{\lambda}{x}, \quad x \in [1, 10] \quad (1)$$

för ett specifikt  $\lambda$ .

- (a) Bestäm (med ett närmevärde)  $\lambda$  så att  $f_X$  blir en täthetsfunktion.

**Svar:** .....

- (b) Bestäm fördelningsfunktionen för  $X$ .

**Svar:** .....

- (c) Bestäm sannolikheten att  $X$  är mindre än 7. Använd  $\lambda = 0.4267$ .

**Svar:** .....

- (d) Bestäm  $E[X]$ .

**Svar:** .....

3. Låt  $U$  vara likformigt fördelad över intervallet  $[0, 2\pi]$ , beräkna

(a)  $E[\cos(U)]$

**Svar:** .....

(b)  $E[\sin(U)^2]$

**Svar:** .....

4. Redogör för Stora talens lag.

5. Redogör för Centrala gränsvärdessatsen.

6. Läs sid. 272-273 om bootstrap i läroboken.

## Syfte och vidare introduktion

Börja laborationen med att ladda ner följande filer från kurshemsidan.

- `plot_mvnpdf.m`
- `hist_density.m`
- `birth.dat`
- `birth.txt` - beskrivning av datat `birth.dat`

Se till att de ligger i den mapp du kommer att arbeta i. För att kontrollera att du har lagt filerna rätt, skriv `ls` och se om filerna ovan listas.

Du kan skriva dina kommandon direkt i MATLAB-prompten men det är absolut att föredra att arbeta i editorn. Om den inte är öppen så kan du öppna den och skapa ett nytt dokument genom att skriva `edit lab1.m`. Koden som ges nedan är skriven i celler. En ny cell påbörjas genom att skriva två procenttecken. **Ctrl+Enter** exekverar innehållet i en cell.

Temat för den här datorlaborationen är simulering. Sannolikhetsteoridelen av kursen handlar om hur man genom beräkningar kan ta fram olika storheter som sannolikheter, väntevärden osv, för en given stokastisk modellen. För mer komplicerade system är det ibland inte alls möjligt att göra exakta beräkningar, eller så är de så tidskrävande att man avstår. I sådana sammanhang kan simulering vara ett alternativ. Simulering innebär att man med hjälp av en dator drar 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 enkla problem, men grundprinciperna går att använda på långt mer komplicerade problem som vi inte kan lösa med enkla beräkningar.

**Problem 0 - Beräkna sannolikheter**

Läs help för funktionerna `binocdf`, `binopdf`, `normcdf`, `normpdf`, `expcdf` och `exppdf`. Observera att MATLABs `exprnd` har väntevärdet  $\mu$  som parameter i motsats till [2] som har  $1/\mu$  som parameter.

Låt  $X_1$  vara  $\text{Bin}(10, 0.3)$ ,  $X_2 \in N(5, 3)$ ,  $X_3 \in \text{Exp}(7)$  och bestäm (med hjälp av funktionerna ovan) för  $k = 1, 2, 3$

1.  $P(X_k \leq 3)$

**Svar:** .....

2.  $P(X_k > 7)$

**Svar:** .....

3.  $P(3 < X_k \leq 4)$

**Svar:** .....

**Problem 1 - Täthetsfunktioner**

Plotta täthetsfunktionen för en exponentialfördelad stokastisk variabel med väntevärde  $\mu$ .

```

1 %% Problem 1: exp-pdf
2     dx = 0.1;
3     x = 0:dx:15;           % Skapar en vektor med dx som inkrement
4     mu = 1;
5     y = exppdf(x, mu);    % exponential-fordelningen
6     plot(x, y)

```

Gör nu samma sak för täthetsfunktionen i förberedelseuppgift 2.

```

1 %% Problem 1: lambda-plot
2     lambda = 0.4267;
3     f=(lambda*exp(-x/lambda)+lambda./x).*(x >= 1 & x <= 10);
4     plot(x, f)

```

Diskutera skillnaden mellan fördelningarna.

**Kommentar:** .....

.....

## Problem 2 - Multivariat Normalfördelning

Täthetsfunktionen för den multivariata normalfördelningen ritas upp av funktionen `plot_mvnpdf`. Undersök hur funktionen fungerar och testa med lite olika parametervärden. Parametrarna `mux` och `muy` kan anta alla reella värden, parametrarna `sigmax` och `sigmay` kan anta alla positiva värden och parametern `rho` kan anta alla värden på intervallet  $[-1, 1]$ . Observera att plotfönstret i funktionen `plot_mvnpdf` är fixt, så för parametervärden som är av storleksordningen tio eller större, så kommer merparten av täthetsfunktionen att hamna utanför plotfönstret.

```
1 %% Problem 2: Multivariat normal
2     mux = 0; muy = -2; sigmax = 1; sigmay = 4; rho = 0.7;
3     plot_mvnpdf(mux, muy, sigmax, sigmay, rho)
```

Hur påverkar olika parametervärden utseendet på plotten?

**Kommentar:** .....

.....

## Problem 3 - Simulering av slumpstal

I denna övning så genererar vi ett stort antal slumpstal, ritar upp deras histogram och plottar slutligen den sanna täthetsfunktionen ovanpå .

```
1 %% Problem 3: Simulering av slumpstal
2     mu = 10;
3     N = 1e4;
4     y = exprnd(mu, N, 1); % Genererar N exp-slumpstal
5     hist_density(y); % Skapar ett normaliserat histogram
6     t = linspace(0, 100, N/10); % Vektor med N/10 punkter
7     hold on
8     plot(t, exppdf(t, mu), 'r') % 'r' betyder rod linje
9     hold off
```

Upprepa simuleringarna och studera hur histogrammet förändras. Hur förhåller sig histogrammet till den röda linjen och hur förklaras variationen kring denna linje?

**Kommentar:** .....

.....

### Problem 4 - Stora talens lag, Monte Carlo och CGS

I detta avsnitt kommer vi igen att generera stora mängder slumpstal, nu för att beräkna väntevärden och sannolikheter. Antag att man vill veta väntevärdet på en tärning, det är inte svårt att beräkna för hand men det går också att kasta tärningen många gånger och sedan räkna ut medelvärdet på dessa kast. Om  $X_1, X_2, \dots, X_n$  är likafördelade med väntevärde  $\mu$  så gäller enligt stora talens lag i [2] att

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| < \varepsilon\right) \rightarrow 1 \quad (2)$$

för varje  $\varepsilon > 0$  när  $n \rightarrow \infty$ . Sannolikheten att skillnaden mellan medelvärdet och det sanna väntevärdet är mindre än  $\varepsilon$  går alltså mot ett när antalet observationer går mot oändligheten. Att använda detta för att beräkna väntevärden kallas för Monte Carlo-metoder.

Idén bakom Monte Carlo är gammal och har funnits inom matematiken åtminstone sedan 1700-talet, men synen kom att förändras under andra halvan av 1900-talet då det på allvar blev möjligt att utföra stora beräkningar. Under 1940-talet utvecklade Stanislaw Ulam och John von Neumann metoder för att göra dessa "tärningskast" med hjälp av dator enligt [1]. Arbetet var kopplat till Manhattanprojektet vars syfte var att ta fram den första atombomben. Metoden namngavs efter casinot Monte Carlo i Monaco.

#### Illustration av Stora talens lag

Koden nedan simulerar exponentialfördelade stokastiska variabler och plottar medelvärdet efter hand. Plotten som visas är alltså din skattning av medelvärdet upp till och med variabel nummer k. Om du tröttnar på att vänta, så kan du kommentera bort `pause`-raden.

```
1 %% Problem 4: Stora talens lag
2 mu = 0.5;
3 M = 500;
4 X = exprnd(mu, M, 1);
5 plot(ones(M, 1)*mu, 'r-.')
6 hold on
7 for k = 1:M
8     plot(k, mean(X(1:k)), '.')
9     if k == 1
10        legend('Sant \mu', 'Skattning av \mu')
11    end
12    xlabel(num2str(k)), pause(0.001)
13 end
14 hold off
```

Ser det ut som förväntat?

**Svar:** .....

### Illustration av Centrala Gränsvärdesatsen

Koden nedan simulerar exponentialfördelade slumpstal och summerar sedan dessa. Studera koden och förklara vad  $N$  representerar.

**Svar:** .....

```
1 %% Problem 4: CGS
2     M = 1e3;
3     N = 4;
4     mu = 5;
5     X = exprnd(mu, M, N);
6     S = cumsum(X, 2);
7     for k = 1:N
8         hist(S(:, k), 30)
9         xlabel(num2str(k))
10        pause(0.1)
11     end
```

Justera  $N$ , vad händer när du ökar respektive minskar värdet? Varför?

**Kommentar:** .....

.....

Vid vilket  $N$  ser det ut som att det inte gör någon skillnad att öka  $N$ ?

**Svar:** .....

Vilken fördelning verkar summorna ha? Varför har de denna fördelning?

**Svar:** .....

### Väntevärden

Uppgiften är nu att beräkna samma väntevärden som i förberedelseuppgift 3. Det första väntevärdet kan du beräkna enligt nedan.

```
1 %% Problem 4: Monte Carlo
2     N = 1e5;
3     U = rand(N, 1)*2*pi;
4     mean(sin(U).^2);
```

Kontrollera att resultatet stämmer med vad du förväntade dig i båda fallen. Låt nu  $X$  och  $Y$  vara oberoende stokastiska variabler där  $X \in Exp(4)$  och  $Y \in N(0, 1)$ . Använd Monte Carlo-metoden för att beräkna  $E[(e^X)^{\cos(Y)}]$ .

Svar: .....

### Problem 5 - Deskriptiv statistik

Vi skall nu gå vidare till att studera skillnaden mellan väntevärden i två populationer, t ex skillnaden i födelsevikt för barn vars mammor röker respektive inte röker under graviditeten. (Om ni vill kan ni ta två andra populationer, och/eller andra variabler att studera!).

I filen `birth.txt` ser man att kolonn 20 i `birth.txt` innehåller rökvanor, och att värdena 1 och 2 betyder att mamman inte röker under graviditeten, medan värdet 3 betyder att hon gör det. Ni kan skapa två variabler `x` och `y` för födelsevikter hörande till icke-rökande respektive rökande mammor enligt

```
>> x = birth(birth(:, 20) < 3, 3);  
>> y = birth(birth(:, 20) == 3, 3);
```

Vad som händer här är att `birth(:, 20) < 3` returnerar en vektor av "sant" och "falskt", och att bara de rader av kolonn 3 (födelsevikterna) i `birth` för vilka jämförelsen är sann, väljs ut.

Använd koden nedan för att visuellt inspektera datat.

```
1 %% Problem 5: Deskriptiv statistik  
2 load lab2data/birth.dat  
3 x = birth(birth(:, 20) < 3, 3);  
4 y = birth(birth(:, 20) == 3, 3);  
5 subplot(2,2,1)  
6 boxplot(x)  
7 axis([0 2 500 5000])  
8 subplot(2,2,2)  
9 boxplot(y)  
10 axis([0 2 500 5000])  
11 subplot(2,2,3:4)  
12 ksdensity(x)  
13 hold on  
14 [fy, ty] = ksdensity(y);  
15 plot(ty, fy, 'r')  
16 hold off
```

Vad betyder plotarna? Vilka slutsatser kan ni dra?

Svar: .....

## Problem 6 - Bootstrap (gör i mån av tid)

Namnet bootstrap syftar till metaforen att dra sig upp ur ett knivig situation genom att ta tag i sina stövlskaft. Ett klassiskt exempel är historien om Baron von Münchhausen i vilken han ska ha räddat sig och sin häst ur ett träsk genom att dra ur de båda genom att lyfta sig själv i håret. Detta förfarande beskriver idén bakom den statistiska varianten av metoden mycket väl: Man har observerat en begränsad mängd data och man vill bilda sig en uppfattning om vad som hade hänt om man hade haft fler observationer. I vårt fall innebär det att vi gör upprepade slumpvisa dragningar från vårt observerade data med funktionen `randsample`.

Vi ska börja med att titta på ett simulerat exempel där vi vet vad som *borde* hända och fortsätter sedan i nästa laboration med att titta på riktiga data, vilket är det som vi huvudsakligen är intresserade av.

### Simulering av summerade exponentialfördelade variabler

Om  $X_1, X_2, \dots, X_n$  alla är oberoende och exponentialfördelade med intensitet  $\lambda$  så gäller

$$Y_n = \sum_{k=1}^n X_k \in \Gamma(n, \frac{1}{\lambda}). \quad (3)$$

Med andra ord är summan av de exponentialfördelade variablerna gammafördelad. (Specialfallet när  $\Gamma$ -fördelningen har ett positivt heltal  $n$  som första parameter kallas mer specifikt för *Erlang-fördelning*.)

Vi ska nu simulera  $M$  stycken summor av typen (3), med  $n$  noterat som `n_sum` i koden nedan, och studera ett histogram för dem. Börja med att studera koden nedan. Läs `help randsample` och testa hur funktionen fungerar. Kommandot `reshape` används nedan för att omforma matrisen `x` till en vektor så att den ska kunna användas i `randsample`. Notera att vi använder väntevärdet  $\mu = \frac{1}{\lambda}$  i koden nedan.



```
1 %% Problem 6: Bootstrap - Simulering
2   % lambda = 1/5;
3   mu = 5;
4   M = 1e3;
5   n_sum = 5;
6   X = exprnd(mu, M, n_sum);
7   g = sum(X, 2);
8   subplot(211)
9   hist_density(g)
10  hold on
11  t = 0:0.01:mu*10;
12  plot(t, gampdf(t, n_sum, mu), 'r')
13  hold off
14  B = 1e3;   % Antal bootstrapreplikat
15  totalNoSamples = M*n_sum;
16  X = reshape(X, totalNoSamples, 1); % Gör om X till vektor
17  yBoot = zeros(B, 1);
18  for j = 1:B
19      sampleDraws = X(randsample(totalNoSamples, n_sum, 1));
20      yBoot(j) = sum(sampleDraws);
21  end
22  subplot(212)
23  hist_density(yBoot)
24  hold on
25  plot(t, gampdf(t, n_sum, mu), 'r')
26  hold off
```

Förklara vad de två raderna i for-loopen gör!

**Kommentar:** .....

.....

Justera B, M, respektive mu. Vad händer? Varför?

**Kommentar:** .....

.....

## Referenser

- [1] Eckhardt, Roger (1987) Stan Ulam, John von Neumann and the Monte Carlo, Method *Los Alamos Sci.*, Vol **15**, p. 131-43.
- [2] Blom, G., Enger, J., Englund, G., Grandell, J., och Holst, L., (2005). Sannolikhets teori och statistikteori med tillämpningar.