

Projektopgave 1

“Sandsynlighedsregning og statistik”

Michael Andersen – michael@diku.dk
Henrik Jensen – henrikjensen@gmail.com
Ulrik Bonde – bonde@diku.dk

30. november 2009

Opgave 1

I bilag A.1 er kildekoden til vores **R**-program inkluderet. Funktionen `assignment1()` løser den første opgave. Vi bruger et *for-loop* til at gå gennem en vektor, hvori vi har defineret de tre antalsparametre. Koden er i øvrigt flot indenteret og selv-dokumenterende.

Centralt i funktionen `assignment1()` har vi metoden fra **R** kaldet `dbinom(x, n, p)`. Denne metode er defineret som binomialfordelingens sandsynlighedsfunktion som er givet ved

$$p(x) = \binom{n}{x} p^x (1-p)^{n-x} \quad (1)$$

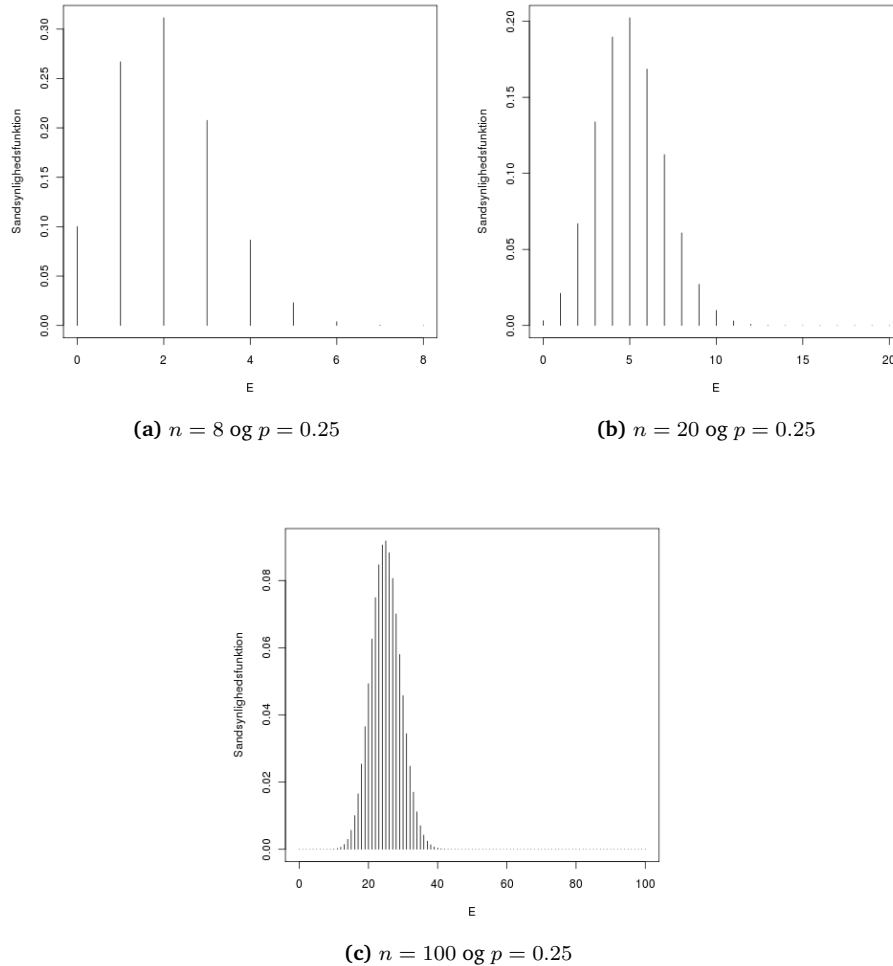
Hvis x er en vektor vil **R** returnere en vektor med resultatet fra ligning 1. Vi kan da plote resultatet med metoden `plot()` i **R**.

I programmets løkke sættes for hver iteration et nyt filnavn til plottet. Vi eksporterer plottet til et png-billede. I figur 1 ses resultatet fra programmet.

Opgave 2

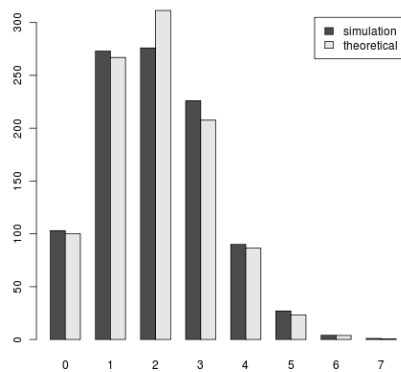
På samme måde som ovenstående opgave, bruger vi et *for-loop* til at simulere 1000 udfald af binomialfordelinger med antalsparameter 8, 20 og 100 og sandsynlighedsparameter $\frac{1}{4}$. Metoden kan findes i kildekoden og er kaldt `assignment2_3_5()` da de resterende programmeringsopgaver afhænger af de simuleringer vi kommer frem til.

En simulering laves ved at bruge metoden `rbinom(n, size, prob)` i **R**, hvor n er antallet af observationer, *size* er antalsparameteren og *prob* er sandsynligheden for et gunstigt udfald.

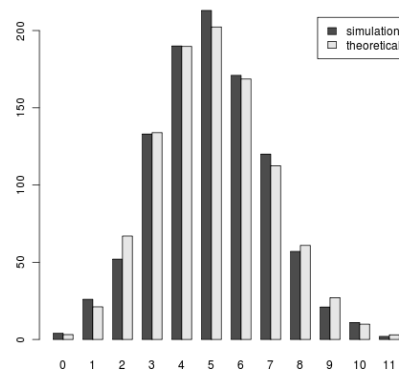


Figur 1: Resultater fra binomialfordelinger med antalsparameter 8, 20 og 100 og sandsynlighedsparameter $\frac{1}{4}$

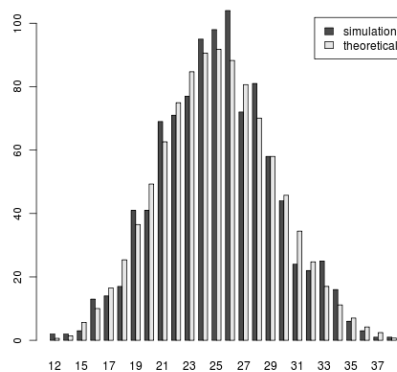
Vi ønsker nu at plote de observerede værdier for at se dem i forhold til de teoretiske værdier. Disse kan ses i figur 2. Det ses at de simulerede værdier ligger tæt op ad de teoretiske værdier. Værdierne er ikke helt ens, men de holder sig i nærheden. Plottene viser kun de observerede værdier, hvilket kan forklare symmetrien i disse plots. Ved at sammeligne med figur 1 ses det at der kun er gjort meget små observationer af de udfald med lille sandsynlighed. Det mest usædvanlige plot ses i figur 2c, hvor $x = 24$, $x = 25$ og $x = 26$ har en lidt større observation end den teoretiske. En nærmere analyse med et konfidensinterval er ikke blevet foretaget for at fastslå om de virkelig er bemærkelsesværdige. De resterende observationer viser ikke denne tendens.



(a) $n = 8$ og $p = 0.25$



(b) $n = 20$ og $p = 0.25$



(c) $n = 100$ og $p = 0.25$

Figur 2: Barplot af simulerede værdier og de teoretiske værdier.

Opgave 3

Metoden `assignment2_3_5()` printer ved kørsel gennemsnittet og den empiriske varians for de simulerede kørsler ud. Dette gøres ved metoder `mean()` og `var()` i **R**. Endvidere bruges metoden `paste()` til at manipulere tekststrengene. Vi får følgende udskrift:

```
[1] Mean for simulation with size 8 is: 2.029
[1] Var for simulation with size 8 is: 1.60376276276276
[1]
[1] Mean for simulation with size 20 is: 4.997
[1] Var for simulation with size 20 is: 3.55654754754755
[1]
```

[1] Mean for simulation with size 100 is: 25.087
 [1] Var for simulation with size 100 is: 17.6550860860861

Middelværdien for en binomialfordeling \mathbf{X} med antalsparameter n og sandsynlighedsparameter p er givet som $E(\mathbf{X}) = np$ og variansen er givet ved $\text{Var}(\mathbf{X}) = np(1 - p)$. For $\mathbf{X} \sim \text{bin}(n, p)$, hvor $n = 8$ og $p = 0.25$, får vi da

$$E(\mathbf{X}) = np = 8 \cdot 0.25 = 2 \quad (2)$$

$$\text{Var}(\mathbf{X}) = np(1 - p) = 8 \cdot 0.25(1 - 0.25) = 1.5 \quad (3)$$

For $\mathbf{X} \sim \text{bin}(n, p)$, hvor $n = 20$ og $p = 0.25$, får vi

$$E(\mathbf{X}) = 20 \cdot 0.25 = 5 \quad (4)$$

$$\text{Var}(\mathbf{X}) = 20 \cdot 0.25(1 - 0.25) = 3.75 \quad (5)$$

Endelig for $\mathbf{X} \sim \text{bin}(n, p)$, hvor $n = 100$ og $p = 0.25$, får vi

$$E(\mathbf{X}) = 100 \cdot 0.25 = 25 \quad (6)$$

$$\text{Var}(\mathbf{X}) = 100 \cdot 0.25(1 - 0.25) = 18.75 \quad (7)$$

Det ses at middelværdien for binomialfordelingerne ligger meget tæt på det beregnede gennemsnit for simulationerne. Den empiriske varians ligger ligeledes tæt på variansen af binomialfordelingen, men dog med større procentvis afvigelse end middelværdi mod gennemsnit.

Opgave 4

a)

Vi betragter k uafhængige stokastiske variable med antalsparameter n og sandsynlighedsparameter p , dette vil sige

$$\mathbf{X}_i \sim \text{bin}(n, p), \text{ for } i = 1, \dots, k$$

Der skal nu findes middelværdien $E(\bar{\mathbf{X}})$, hvor $\bar{\mathbf{X}} = \frac{1}{k} \sum_{i=1}^k \mathbf{X}_i$, da \mathbf{X}_i er binomialfordelt med antalsparameter n og sandsynlighedsparameter p er $E(\mathbf{X}_i) = np$. Derfor følger udregningen

$$E\bar{\mathbf{X}} = E\left(\frac{1}{k} \sum_{i=1}^k \mathbf{X}_i\right)$$

$$= \frac{1}{k} \sum_{i=1}^k E(\mathbf{X}_i)$$

$$= \frac{1}{k} knp$$

$$= np$$

Vi skal nu finde det samme for variansen, dette gøres påfølgende måde, vi ved at $\text{Var}(\bar{\mathbf{X}}) = \text{E}(\bar{\mathbf{X}}^2) - \text{E}(\bar{\mathbf{X}})^2$ og der vides at $\text{Var}(\mathbf{X}_i) = np(1-p)$.

$$\text{Var}(\bar{\mathbf{X}}) = \text{Var}\left(\frac{1}{k} \sum_{i=1}^k \mathbf{X}_i\right)$$

$$= \frac{1}{k^2} \sum_{i=1}^k \text{Var}(\mathbf{X}_i)$$

$$= \frac{1}{k^2} \sum_{i=1}^k np(1-p)$$

Her benyttes sætning 3.7.13.

$$= \frac{1}{k^2} knp(1-p)$$

Derfor er

$$\text{Var}(\bar{\mathbf{X}}) = \frac{np(1-p)}{k}$$

b)

Den empiriske varians, kan findes ved at kigge på opgave 3.26, og beskrives som

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})^2$$

Da vi betragter k stokastiske variable fremfor n , har vi altså at

$$s^2 = \frac{1}{k-1} \sum_{i=1}^k (\mathbf{X}_i - \bar{\mathbf{X}})^2$$

Vi skal nu finde middelværdien af dette, hvilket vil sige at vi skal kigge på

$$\begin{aligned} \text{E}(s^2) &= \text{E}\left(\frac{1}{k-1} \sum_{i=1}^k \mathbf{X}_i - \bar{\mathbf{X}}\right)^2 \\ &= \frac{1}{k-1} \sum_{i=1}^k \text{E}(\mathbf{X}_i - \bar{\mathbf{X}})^2 \\ &= \frac{1}{k-1} \sum_{i=1}^k \text{E}(\mathbf{X}_i^2 + \bar{\mathbf{X}}^2 - 2\mathbf{X}_i\bar{\mathbf{X}}) \end{aligned}$$

Vi kan ved at bruge sætning 3.7.6 og 3.7.4, dele ovenstående op og tage middel værdien af hvert enkelt led, således at vi har

$$= \frac{1}{k-1} \sum_{i=1}^k \left(E(\mathbf{X}_i^2) + E(\bar{\mathbf{X}}^2) - 2E(\mathbf{X}_i \bar{\mathbf{X}}) \right)$$

Vi betragter nu kun det første led, altså $E(\mathbf{X}_i^2)$. Til dette ved vi at der i henhold til 3.7.9 gælder at

$$\text{Var}(\mathbf{X}) = E(\mathbf{X}^2) - (E(\mathbf{X}))^2$$

Her isolerer vi $E(\mathbf{X}^2) = \text{Var}(\mathbf{X}) + (E(\mathbf{X}))^2$ fra første del af opgave 4 ved vi at $\text{Var}(\mathbf{X}_i) = np(1-p)$ og $(E(\mathbf{X}_i))^2 = n^2 p^2$ Derfor har vi at

$$E(\mathbf{X}_i^2) = np(1-p) + n^2 p^2$$

Det samme gør sig næsten gældende for ledet $E(\bar{\mathbf{X}}^2)$, her har vi fra første del af opgaven at

$$\text{Var}(\bar{\mathbf{X}}) = \frac{np(1-p)}{k} \text{ og } (E(\bar{\mathbf{X}}))^2 = n^2 p^2$$

Derfor har vi at

$$E(\bar{\mathbf{X}}^2) = \frac{np(1-p)}{k} + n^2 p^2$$

Der betragtes nu det sidste led som har formen

$$\begin{aligned} & E(\mathbf{X}_i \bar{\mathbf{X}}) \\ &= E\left(\mathbf{X}_i \frac{1}{k} \sum_{j=1}^k \mathbf{X}_j\right) \\ &= \frac{1}{k} \sum_{j=1}^k E(\mathbf{X}_i \mathbf{X}_j) \end{aligned}$$

Her kan vi i henhold til sætning 3.7.7 gøre følgende, da $\mathbf{X}_i \perp \mathbf{X}_j$

$$= \frac{1}{k} \left(\sum_{\substack{j=1 \\ j \neq i}}^k (E\mathbf{X}_i E\mathbf{X}_j) + E(\mathbf{X}_i^2) \right)$$

Vi tager nu resultaterne fra første halv del af opgaven og bruger her.

$$\begin{aligned} &= \frac{1}{k} \left(\sum_{\substack{j=1 \\ j \neq i}}^k (n^2 p^2 + np(1-p) + n^2 p^2) \right) \\ &= \frac{1}{k} ((k-1)n^2 p^2 + np(1-p) + n^2 p^2) \end{aligned}$$

Disse resultater indsættes nu i vores oprindelige formel

$$\begin{aligned}
& \frac{1}{k-1} \sum_{i=1}^k \left(E(\mathbf{X}_i^2) + E(\bar{\mathbf{X}}^2) - 2E(\mathbf{X}_i \bar{\mathbf{X}}) \right) \\
&= \frac{1}{k-1} \sum_{i=1}^k \left(np(1-p) + n^2 p^2 + \frac{np(1-p)}{k} + n^2 p^2 - 2 \frac{1}{k} ((k-1)n^2 p^2 + np(1-p) + n^2 p^2) \right) \\
&= \frac{1}{k-1} ((knp(1-p) + kn^2 p^2 + np(1-p) + kn^2 p^2) - 2((k-1)n^2 p^2 + np(1-p) + n^2 p^2)) \\
&= \frac{1}{k-1} (knp(1-p) + kn^2 p^2 + np(1-p) + kn^2 p^2 - 2kn^2 p^2 + 2n^2 p^2 - 2np(1-p) - 2n^2 p^2) \\
&= \frac{1}{k-1} (knp(1-p) - np(1-p)) \\
&= \frac{1}{k-1} ((k-1)np(1-p)) \\
&= np(1-p)
\end{aligned}$$

Opgave 5

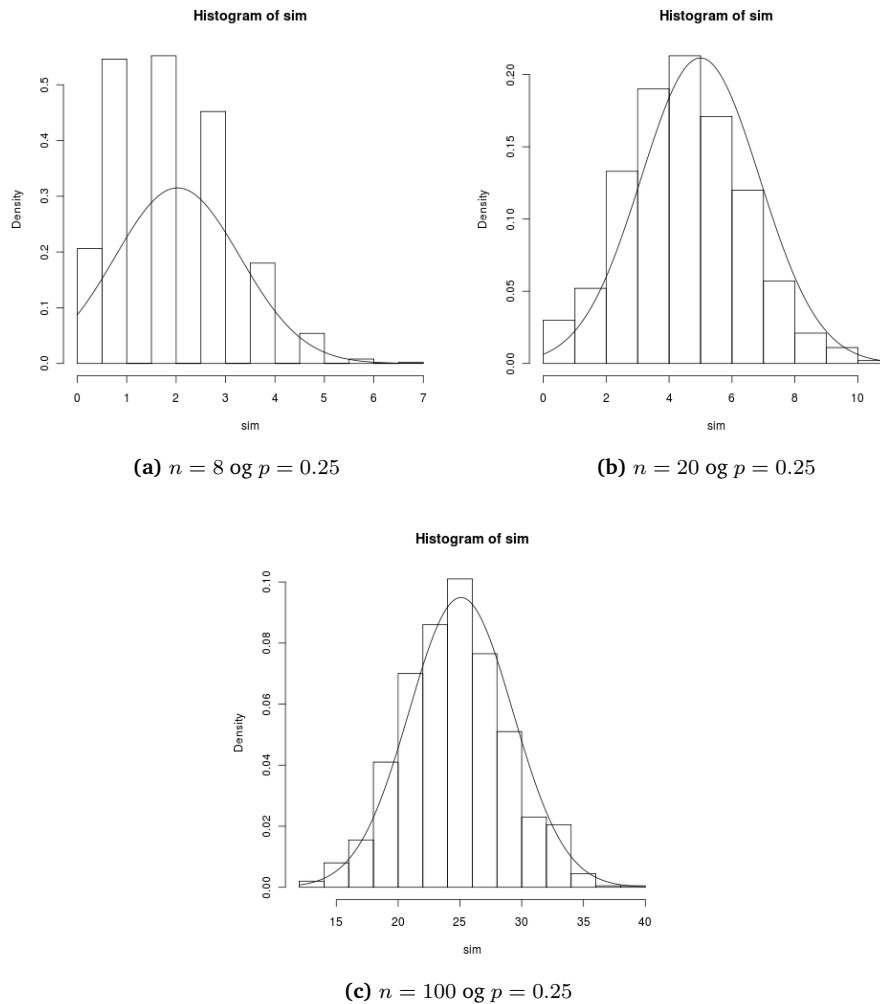
Vi får at vide, at binomialfordelingen tilnærmer sig normalfordelingen for store værdier af n . Dette ses i figur 3. For $n = 8$ ses det at de observerede værdier fra simulationen ikke følger normalfordelingen. Ved $n = 20$ ligger vi tættere på normalfordelingen. For $n = 100$ ligger observationerne stort set på niveau med normalfordelingen, men vi lægger også her mærke til at der i intervallet $[24, 26)$ er et højt antal observationer, som vi så i opgave 2. Også intervallet $[30, 32)$ ligger udenfor normalfordelingen, men her er der et færre antal observationer.

Opgave 6

Da der ligger et uendeligt antal handsker på stranden antager den stokastiske variable \mathbf{X} værdier i udfaldsrummet $E = \{0, \dots, n\}$, $n \rightarrow \infty$. Det oplyses at der er lige mange højre og venstre handsker, derfor er sandsynlighedsparameteren $p = 0.5$

I det konkrete tilfælde med $n = 14$, har vi altså

$$\mathbf{X} \sim \text{bin}(14, 0.5)$$



Figur 3: Histogram og normalfordeling

Opgave 7

Der skal beregnes sandsynligheden for at få netop 10 venstrehandsker og 4 højrehandsker. I henhold til den opstillede statistiske model. Dette gøres ved at bruge formelen $\binom{n}{k} p^k (p-1)^{n-k}$.

Derfor har vi at

$$\binom{14}{10} 0.5^{10} (1 - 0.5)^{14-10} = \binom{14}{10} 0.5^{14} = 0.061$$

Sandsynligheden for udfaldet 10 venstrehandsker og 4 højrehandsker er derfor ca. 6%.

Opgave 8

Vi skal finde sandsynligheden for at finde 10 eller flere handsker til samme hånd i den nævnte model, hvor der er lige mange højre- og venstrehandsker.

Det vil sige at vi er interesseret i sandsynligheden for

$$P(\mathbf{X} \leq 4, \mathbf{X} \geq 10) = P(\mathbf{X} \leq 4) + P(\mathbf{X} \geq 10) = P(\mathbf{X} \leq 4) + (1 - P(\mathbf{X} \leq 9)) = 2 \cdot P(\mathbf{X} \leq 4)$$

Det sidste kan lade sig gøre da sandsynlighedsfunktionen er symmetrisk omkring middelværdien 7

$$P(\mathbf{X} \leq 4) + (1 - P(\mathbf{X} \leq 9)) = 0.0898 + 0.0898 = 0.18$$

Så i henhold til den opstillede statistiske model er der 18% sandsynlighed for at finde 10 eller flere af samme type handske.

Opgave 9

Vi skal nu kommentere på hvor vidt det årvågne strandvandrepars fund er bemærkelsesværdigt eller ej. For at undersøge dette laver vi et 95%-konfidensinterval for p , hvis 0.18 er inkluderet i dette interval er fundet ikke bemærkelsesværdigt. Formlen der skal benyttes er som følger

$$p \pm 1.96 \cdot s(p), \text{ hvor } s(p) = \sqrt{\frac{p(1-p)}{n}}$$

Vi har $p = 0.5$ og $n = 14$, udregningen er derfor som følger.

$$0.5 \pm 1.96 \cdot 0.134 = (0.238, 0.762)$$

Da

$$0.18 \notin (0.238, 0.762)$$

må fundet af de 10 venstrehandsker være bemærkelsesværdigt.

A Source code

A.1 main.r

```
1  #!/usr/bin/env Rscript
2
3  # R statistics for the course "Sandsynlighedsregning og statistik"
4  # at the Department of Mathematical Sciences, University of Copenhagen
5
6  #####
7  #
8  # Project 1
9  #
10 #####
11
12 # Assignment 1
13 assignment1 <- function() {
14   # We define n as an array with the values 8, 20 and 100
15   n <- c(8, 20, 100);
16   prob <- 0.25;
17
18   # Traverse the n-array
19   for (size in n) {
20     # Define E
21     E <- 0:size;
22
23     # Define the probability function
24     sshfunk <- dbinom(E, size, prob);
25
26     # Set the filename for the plot and plot the graph
27     name <- paste(size, "_plot_1.png", sep="");
28     png(name);
29     plot(E, sshfunk, "h", ylab="Sandsynlighedsfunktion");
30     dev.off();
31   }
32 }
33
34 #####
35
36 # Assignment 2, 3, and 5
37 assignment2_3_5 <- function() {
38   # Define the simulation using a random binom
39   n <- 1000;
40   size <- c(8, 20, 100);
41   prob <- 0.25;
42
43   for (s in size) {
44     # Do a simulation and count the observations of each occurrence
45     sim <- rbinom(n, s, prob);
46     simulation <- table(sim);
47
48     # Assignment 3
49     # Print mean and var
50     print(paste("Mean for simulation with size ",s, " is: ", mean(sim), sep=""));
51     print(paste("Var for simulation with size ",s, " is: ", var(sim), sep=""));
52     print("");
53
54     # Get the values corresponding to an occurrence
55     values <- as.numeric(names(simulation));
56
57     theoretical <- 1000*dbinom(values, s, prob);
58
59     # Combine the simulation and the theoretical vectors
60     comb <- rbind(simulation, theoretical);
61
62     # Assignment 2
63     # Set up the filename for the plot, then draw it
64     name <- paste(s, "_sim-theo_plot_2.png", sep="");
65     png(name);
66     barplot(comb, beside=TRUE, legend.text=TRUE);
67
68     # Assignment 5
```

```

69         # Set up the filename for the histogram, draw and close pipe
70         histname <- paste(s, "_hist_5.png", sep="");
71         png(histname);
72         hist(sim, prob=TRUE);
73         curve(dnorm(x, mean = mean(sim), sd = sd(sim)), add=TRUE);
74         dev.off();
75     }
76 }
77
78 #####
79 # Main
80
81 assignment1 ();
82 assignment2_3_5();

```