

Coupon Collector's Problem With Unequal Probabilities Application to Ecology

Marko KACHAIKIN, Vincent RUNGE



05/07/2022

Contents

Introduction	1
From equal to linear probability models	1
The log-linear expectation in equality case	2
Estimation of N	3
A nice picture	3
Moment estimation	8

Introduction

In this work we study the coupon collector's problem in the more generic case of unequal occurrence probabilities. Our first goal consists in better analyzing the probability distribution (mean and variance) for the random variable of the time to completion for some particular cases. We found simple formulas in asymptotic regime and compare these results throughout an extensive simulation study.

Our secondary goal is the search for estimators for unknown coupon number when we stop the collection at some chosen step. Performances of proposed estimators are studied both theoretically and by simulations.

```
library(coupon)
```

From equal to linear probability models

We consider that the collection is made of N coupons with three possibles models.

1. the equal probability model:

$$p_i = \frac{1}{N}, \quad i = 1, \dots, N.$$

2. the linear probability model:

$$p_i = \frac{1}{N}, \quad i = 1, \dots, N,$$

with $\beta \in [0, \frac{2}{N(N-1)}]$.

3. the 2-probability model:

$$p_i = \frac{1}{N}, \quad i = 1, \dots, N.$$

The log-linear expectation in equality case

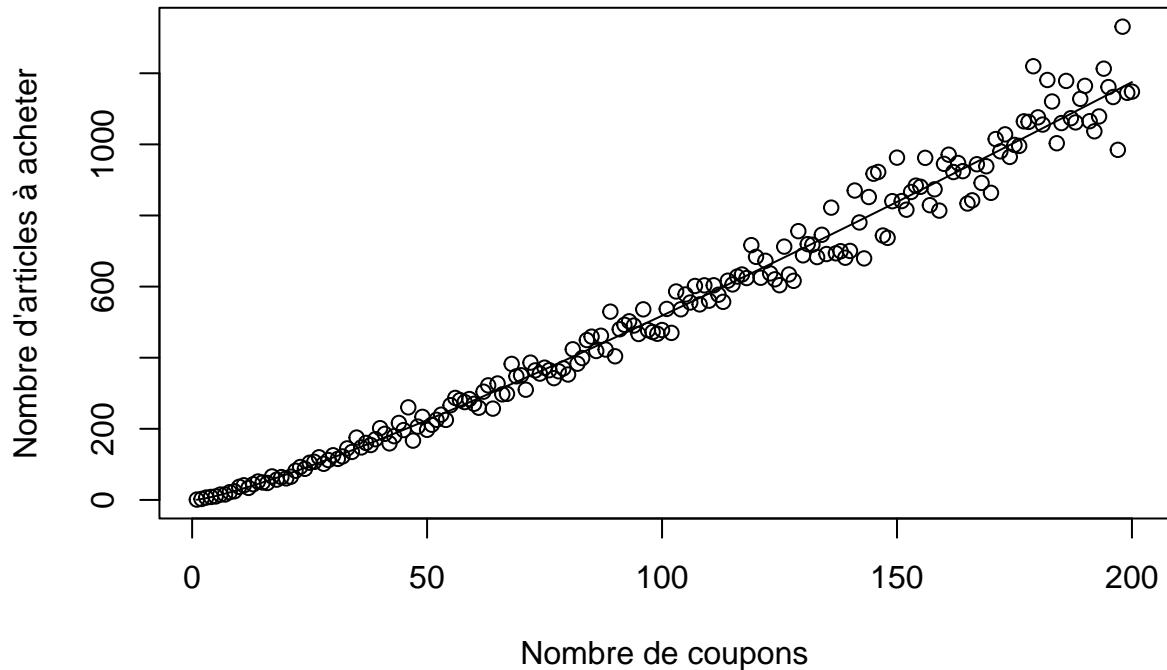
In case 1 we can easily verify the well-known asymptotic result:

$$\mathbb{E}[T] \approx N(\ln N + \gamma)$$

For $N \in \{1, \dots, 200\}$ we simulate 10^3 coupon problems. We show the detailed code using our package `coupon` available on GitHub (<https://github.com/vrunge/coupon>)

```
res_coupon <- function(n, nb_iterations = 10)
{
  mean(replicate(nb_iterations, simu.coupon(nbCoupons = n)))
}
res_theory <- function(n){n*(log(n) - digamma(1))}

N <- 200
moyenne_resultat1 <- sapply(1:N, res_coupon)
moyenne_theorique <- sapply(1:N, res_theory)
plot(moyenne_resultat1, xlab = "Nombre de coupons", ylab = "Nombre d'articles à acheter")
lines(moyenne_theorique)
```

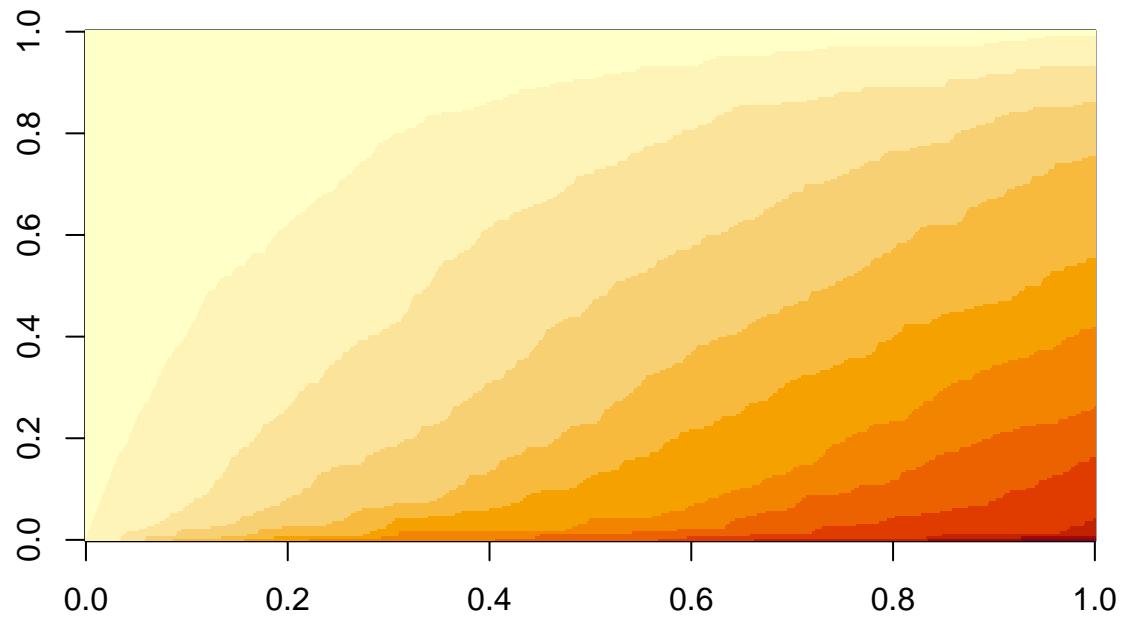


Estimation of N

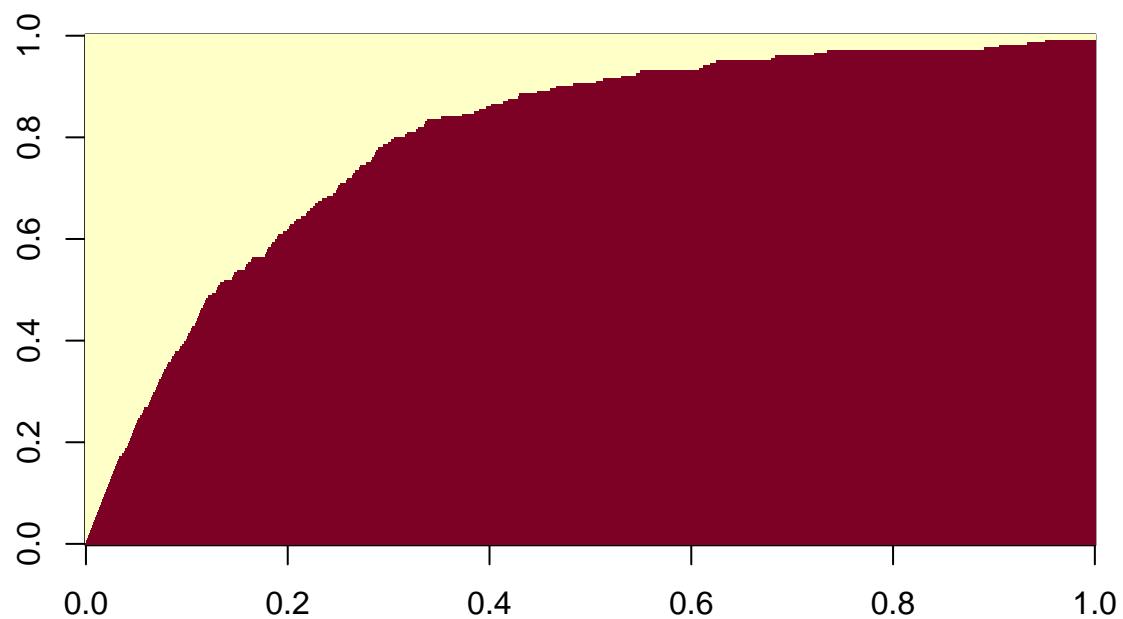
A nice picture

We give a few examples of dynamics in the coupon problem. The x-axis represents the time dynamics and the colors the number of observations for each coupon sorted from the highest occurrence (bottom) to the smallest (0 if collection not completed) (top)

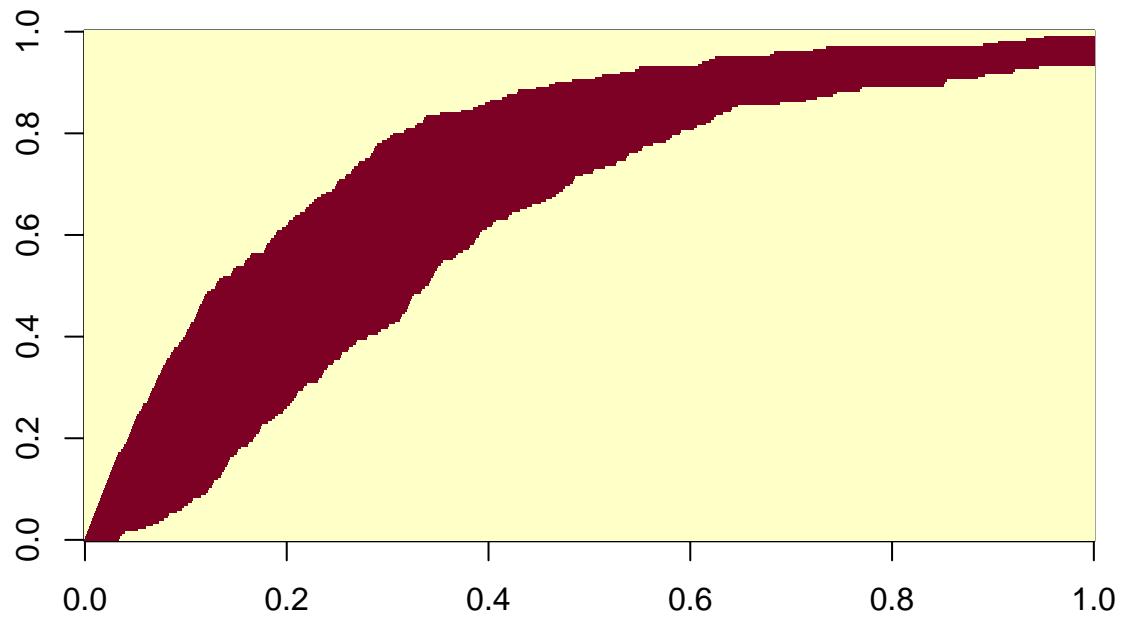
```
nbCoupons <- 200
myN <- 1000
image(res <- dynamicCollection(nbCoupons = nbCoupons, N = myN))
```



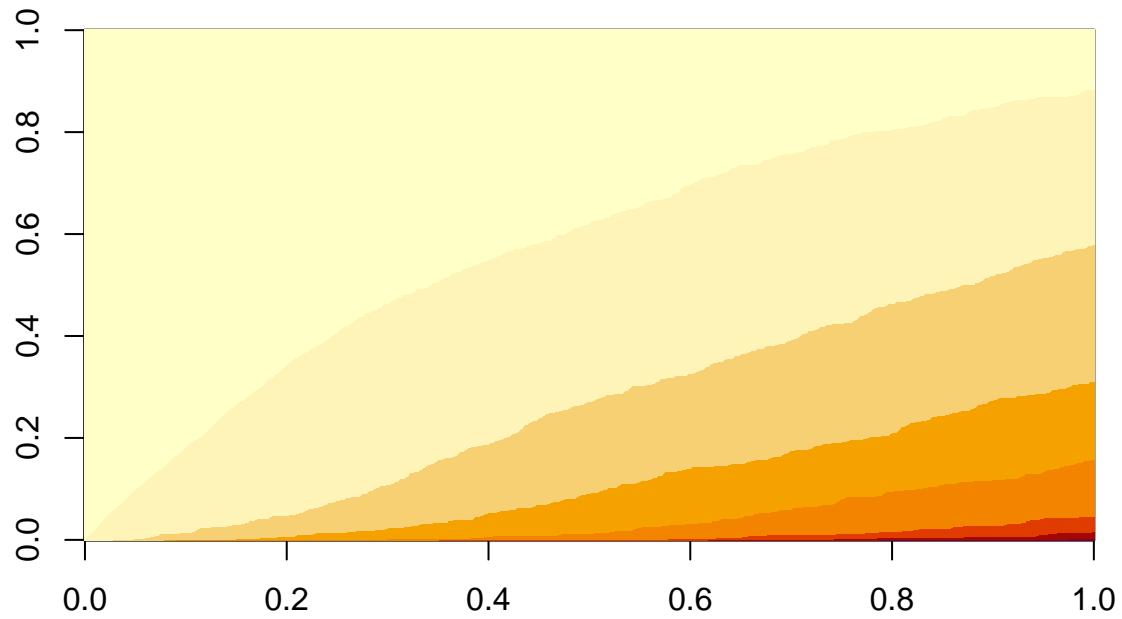
```
image(res>0)
```



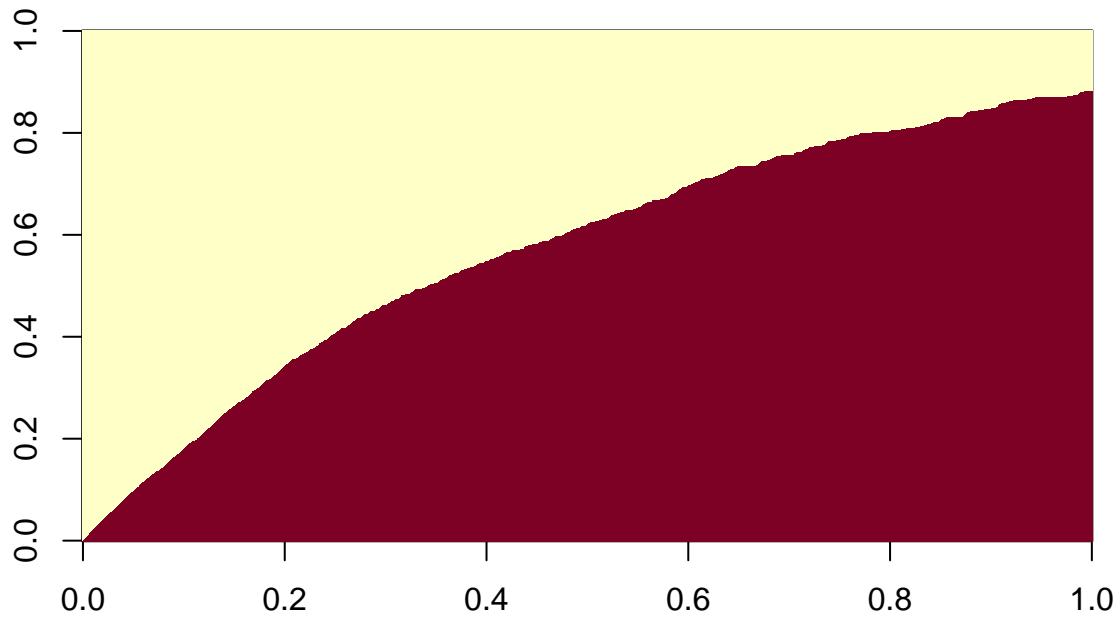
```
image(res==1)
```



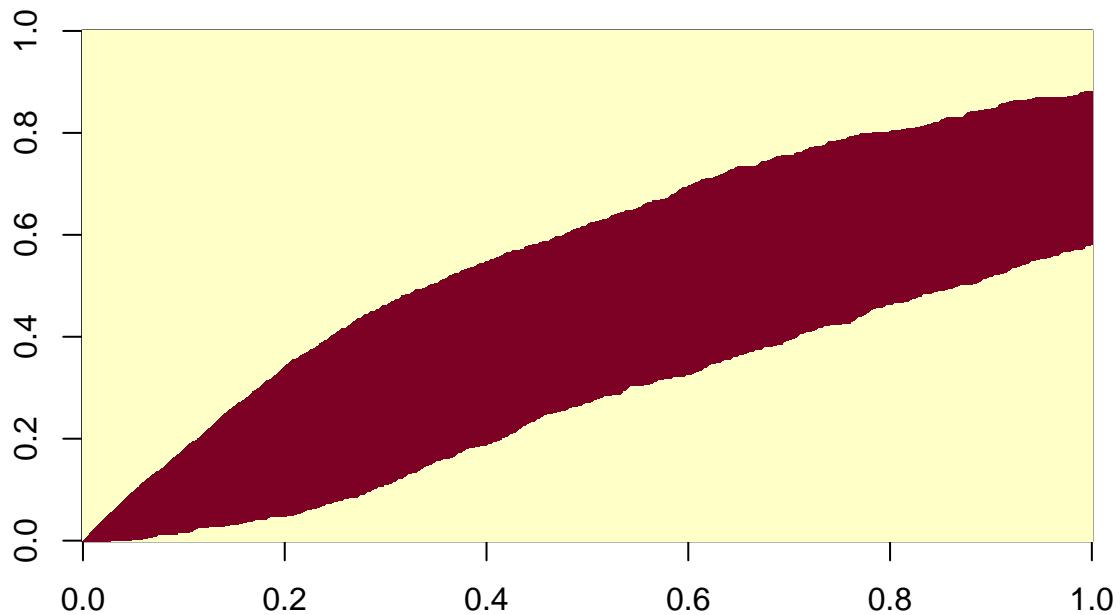
```
nbCoupons <- 500  
myN <- 1000  
image(res <- dynamicCollection(nbCoupons = nbCoupons, N = myN))
```



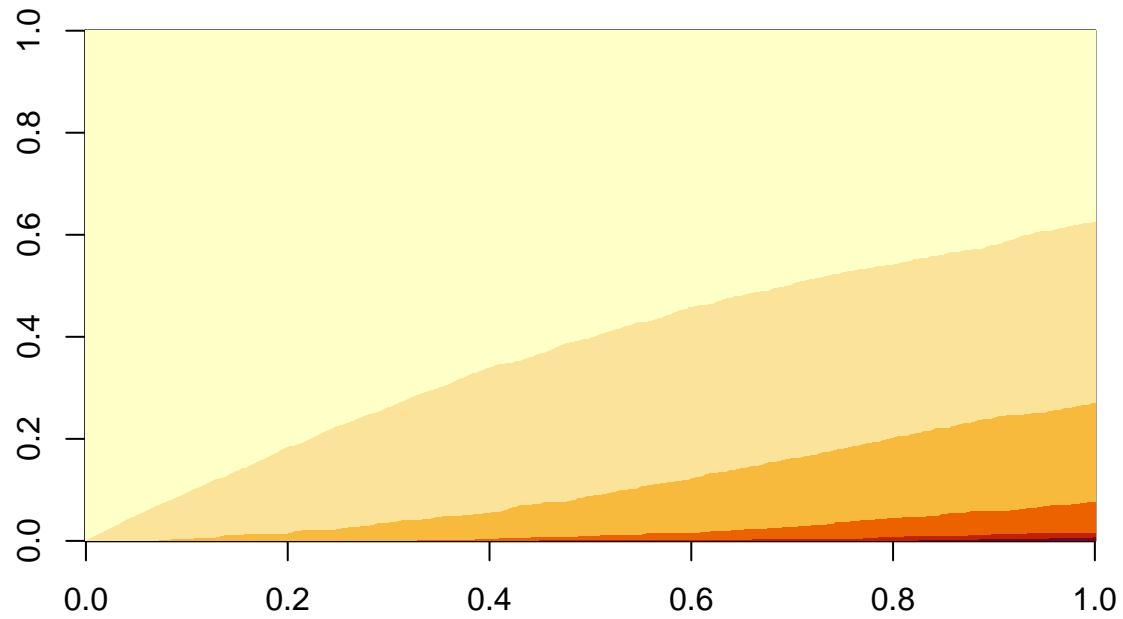
```
image(res>0)
```



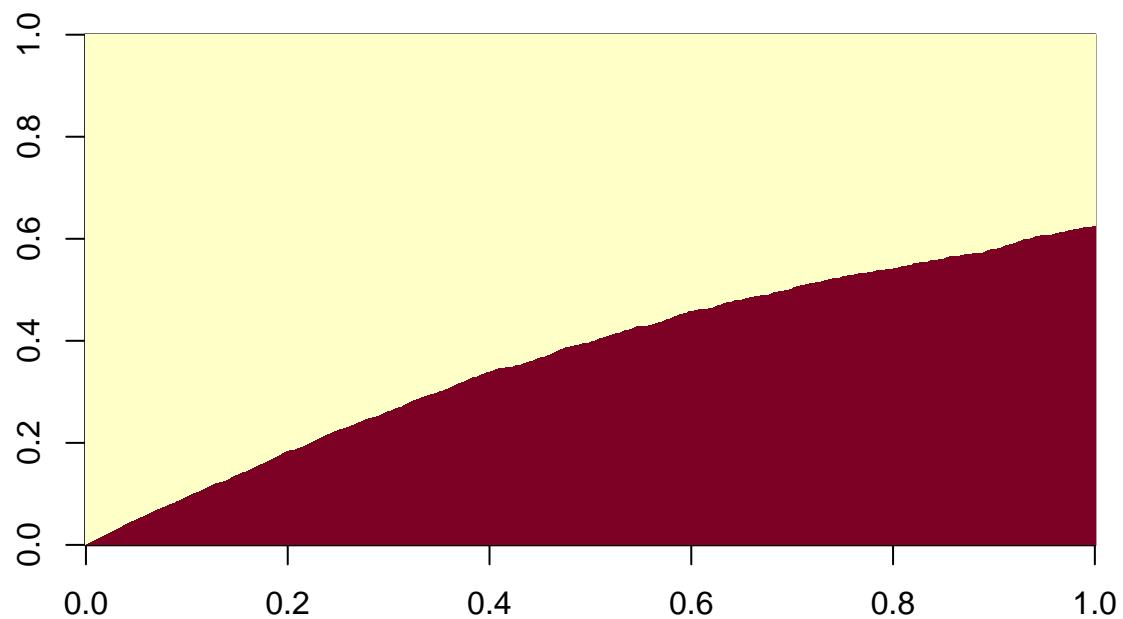
```
image(res==1)
```



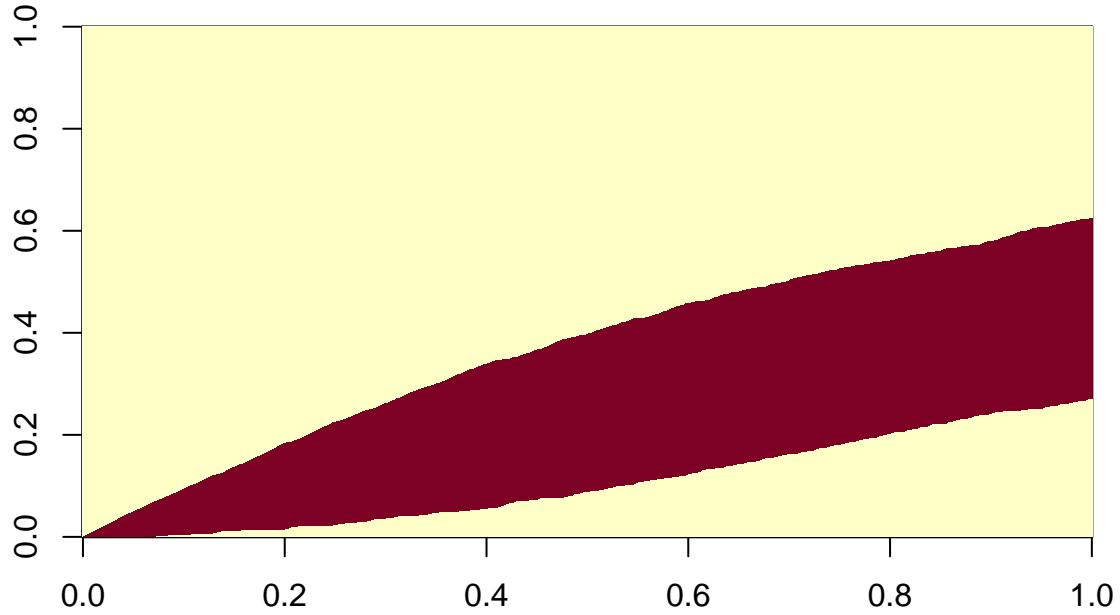
```
nbCoupons <- 1000  
myN <- 1000  
image(res <- dynamicCollection(nbCoupons = nbCoupons, N = myN))
```



image(res>0)



image(res==1)



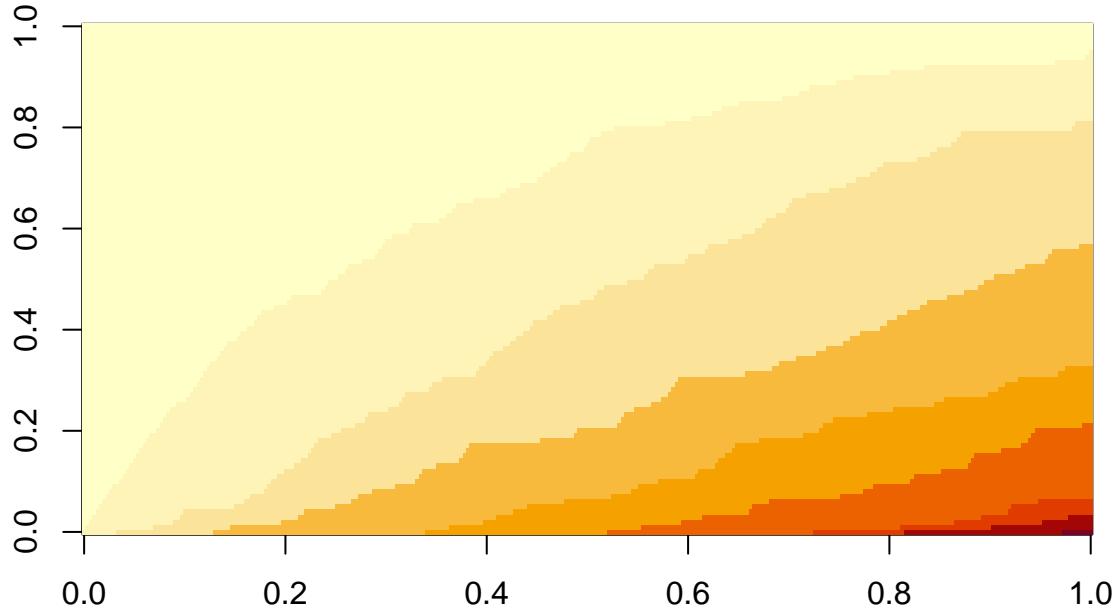
Moment estimation

The expected number of coupons at time t is given by formula

$$N \left(1 - \left(1 - \frac{1}{N} \right)^t \right)$$

We verify the closeness of the two curves (the expectation and the observed counts)

```
nbCoupons <- 100
myN <- 300
image(res <- dynamicCollection(nbCoupons = nbCoupons, N = myN))
```

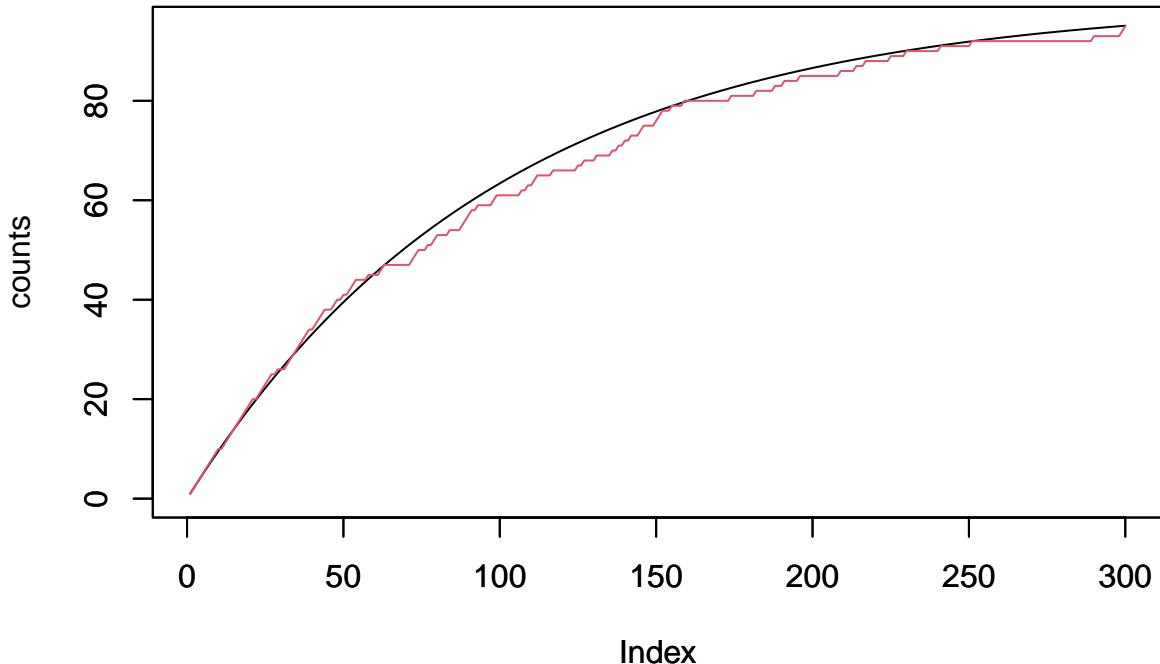


```

nb_time <- apply(res, 1, function(x) sum(x > 0))

curve <- nbCoupons*(1- (1-1/nbCoupons)^(1:myN))
ylimit <- max(c(nb_time, curve))
plot(curve, type = 'l', ylim = c(0,ylimit), ylab = "")
par(new = TRUE)
plot(nb_time, type = 'l', col = 2, ylim = c(0,ylimit), ylab = "counts")

```



Study of the quality of our estimator

```

res <- replicate(1000, estimatorExpectation(nbCoupons = 200,N = 300))
mean(unlist(res[1,]))

```

```
## [1] 200.617
```

```
sd(unlist(res[1,]))
```

```
## [1] 12.53728
```

```
mean(unlist(res[2,]))
```

```
## [1] 155.495
```

```
sd(unlist(res[2,]))
```

```
## [1] 4.520886
```

```
hist(unlist(res[1,]), breaks = 100)
```

Histogram of unlist(res[1,])

