

intro R

jmm

Monday, August 03, 2015

R de reproducibilidad

Hoy en día, los científicos deben estar capacitados para “sacarle el jugo” a distintas fuentes de información. Parte del trabajo consiste en procesar bases de datos, resumir, graficar, y ajustar modelos a datos. Trabajar con herramientas tipo Excel no solo es limitante (y poco elegante) sino que además atenta contra la “reproducibilidad” de la investigación ya que normalmente no nos queda una documentación detallada de cómo se llega a un resultado o un gráfico en particular (hay mucho movimiento y “click” de mouse de por medio). La reproducibilidad es (o debería ser) un estándar mínimo para el trabajo científico. Si bien, el hecho de que algo sea reproducible no quiere decir que sea correcto, al menos lo hace más fácil de evaluar, de revisar y de re-utilizar. Un paso importante para lograr la reproducibilidad es hacer todo con instrucciones escritas (*scripts*) en lugar de con el mouse.

En resumen, para trabajar en ciencia hay que saber un mínimo de programación. Si tenemos que elegir en qué lenguaje programar, R o Python parecen las mejores opciones a la fecha, pero R es actualmente el más popular entre los ecólogos y es el que vamos a usar en este curso. R es una implementación de dominio público del lenguaje estadístico S que ha ganado popularidad por varios motivos. Se puede descargar de <https://cran.r-project.org/> y funciona en PC, Mac y Linux. Además de usarse para análisis estadísticos, R se puede programar para hacer casi cualquier cosa. Tiene programación basada en objetos y programación funcional. Hay una creciente comunidad de usuarios, incluyendo a algunos grosos...

R se distribuye con una interfase gráfica rudimentaria, pero hay varios motivos para trabajar con R a través de RStudio. Entre otras cosas, RStudio hace posible conectarse con repositorios online como GitHub o Bitbucket para control de versiones y crear documentos (como éste que fue escrito en RStudio) que incluyen código de R, figuras, ecuaciones, etc. Aprendiendo a usar R y a documentar nuestro trabajo con las herramientas de RStudio podremos trabajar mejor y de manera más eficiente. Al principio no parece, pero es **mucho mejor** trabajar con comandos (*scripts*) y lograr que nuestra investigación sea reproducible. De lo contrario y parafraseando a Walter Sobchak (1998), seremos “fuckin’ amateurs”...

Cosas Básicas de R

R funciona como una calculadora:

```
2 + 2
```

```
## [1] 4
```

```
2 * 3
```

```
## [1] 6
```

```
3^2
```

```
## [1] 9
```

```
2/3
```

```
## [1] 0.6666667
```

```
exp(1)
```

```
## [1] 2.718282
```

```
log(10)
```

```
## [1] 2.302585
```

```
log10(10)
```

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

R resuelve operaciones lógicas:

```
2 > 1
```

```
## [1] TRUE
```

```
2 == 1
```

```
## [1] FALSE
```

```
2 <= 2
```

```
## [1] TRUE
```

Un aspecto básico que tenemos que dominar es que R trabaja con vectores, matrices, listas y “dataframes”. Empecemos generando un vector:

```
n <- c(2, 6, 18, 54, 162, 486)
```

El comando `c` concatena los valores entre comas. La flecha `<-` asigna lo que está a la derecha al ‘objeto’ de la izquierda. En vez de `<-` podemos usar `=` pero los puristas de R prefieren `<-`.

Para ver los valores del vector en la ventana de comando escribimos su nombre en la ventana de comando:

```
n
```

```
## [1] 2 6 18 54 162 486
```

Un tipo especial de vector son las secuencias. Para generar una secuencia de valores podemos hacer:

```
yr <- 2009:2014
```

Hagamos de cuenta que los vectores `n` y `yr` corresponden a censos de una especie invasora. Si queremos ver el tamaño poblacional en el año 2010 hacemos:

```
n[2]
```

```
## [1] 6
```

O podemos hacer:

```
n[yr==2010]
```

```
## [1] 6
```

O usar la función `which`:

```
idx <- which(yr == 2010)
n[idx]
```

```
## [1] 6
```

Si queremos ver el tamaño poblacional desde el 2011 al 2014

```
n[3:6]
```

```
## [1] 18 54 162 486
```

O bien

```
n[yr >= 2011]
```

```
## [1] 18 54 162 486
```

o con `which`:

```
idx <- which(yr > 2010)
n[idx]
```

```
## [1] 18 54 162 486
```

Algo interesante para ver en un set de datos como este es el “factor reproductivo” o “tasa finita de incremento la población”, que no es otra cosa que $n(t+1)/n(t)$. Por ejemplo, para el 2010 la tasa de crecimiento fue de: $n[2]/n[1] = 3$.

Para no repetir las cuentas año por año podemos hacer un *for loop* pero antes necesitamos generar un vector para guardar los resultados:

```
tasa <- numeric(length(n) - 1)
```

Aquí combinamos dos funciones de R: `numeric` y `length`. Para ver qué hace ‘`numeric`’ podemos hacer `?numeric`. Poniendo un signo de pregunta y luego el nombre de una función de R accedemos a la documentación de esa función (que a veces no es de mucha ayuda...). Una mejor forma de encontrar ayuda e incluso información acerca de funciones que no conocemos es visitando a <http://rseek.org/>

Volviendo a nuestro *loop*, lo que queremos es repetir una operación (o varias) e ir guardando los resultados:

```
for(i in 2:length(n)){  
  tasa[i-1] <- n[i]/n[i-1]  
}
```

El comando `for` aquí itera en valores de `i` la cuenta $n(t+1)/n(t)$. Fijense que `i` está definido desde 2 hasta el largo del vector `n`. Usamos el valor de `i` para cada iteración para ver qué valores de `n` usar en la cuenta y para indicar en qué posición del vector `tasa` guardamos el resultado.

En este caso podríamos obtener los mismos valores haciendo:

```
for(i in 1:(length(n)-1)){  
  tasa[i] <- n[i+1]/n[i]  
}
```

Es importante notar que `1:(length(n)-1)` no es lo mismo que `1:length(n)-1`. Compruébenlo...

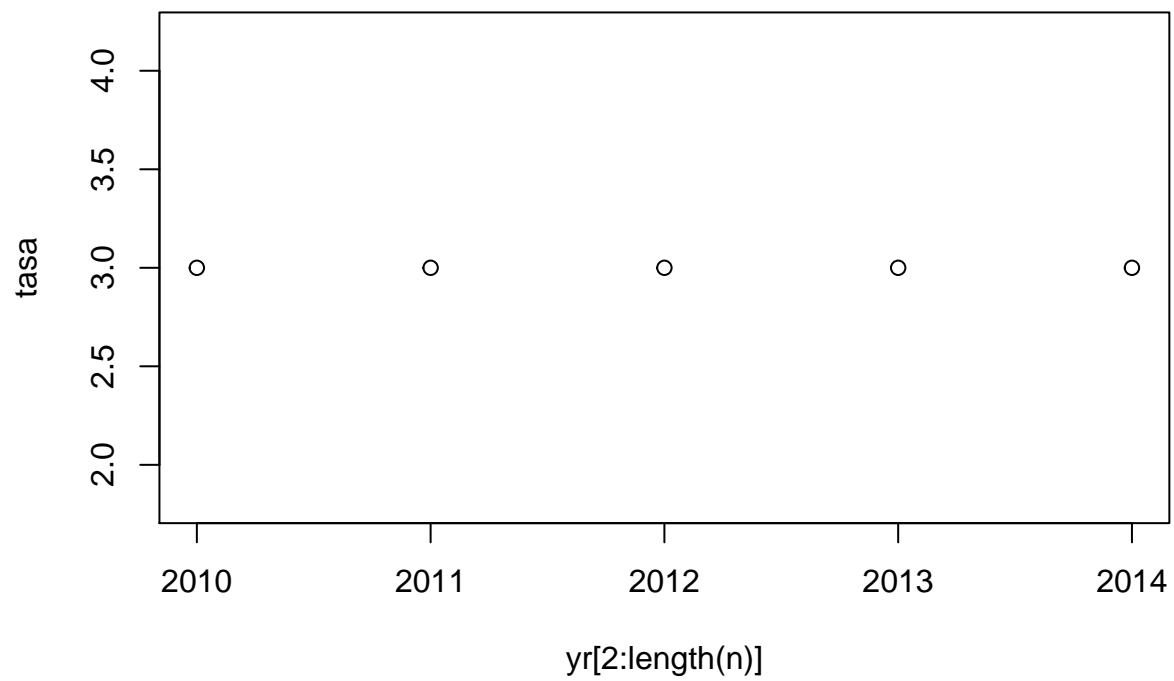
Para ver qué salió de estas cuentas escribimos

```
tasa
```

```
## [1] 3 3 3 3 3
```

Para graficar la tasa de incremento en cada año podemos usar

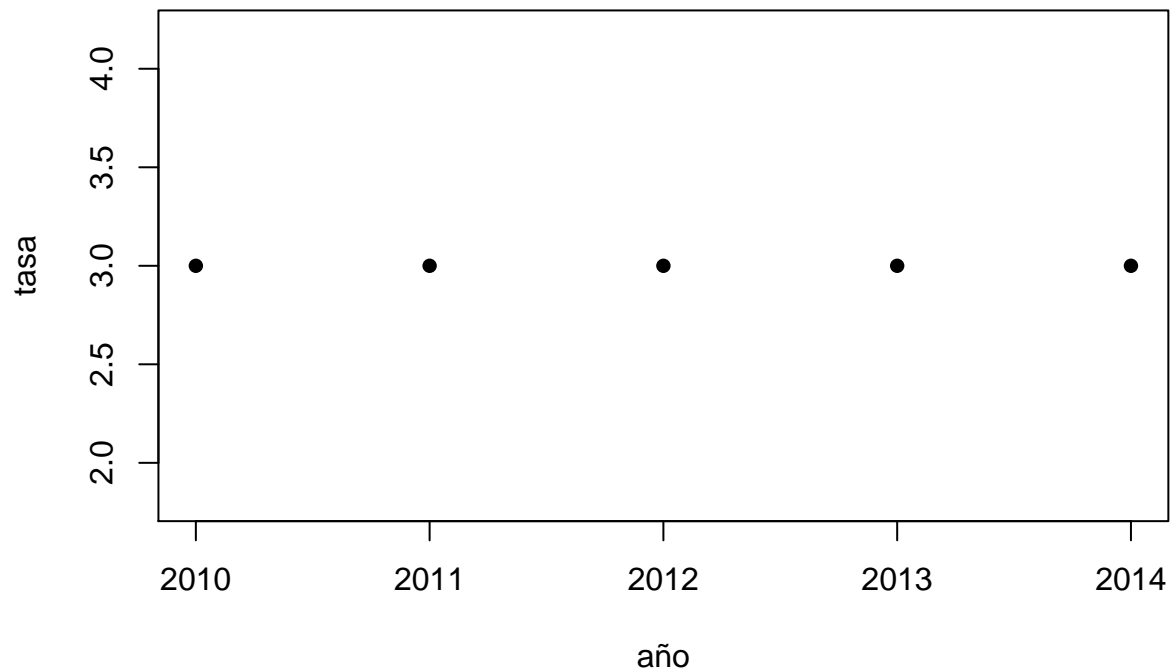
```
plot(yr[2:length(n)], tasa)
```



Podemos “embellecer” un poco el gráfico usando `main`, `xlab`, y `pch`

```
plot(yr[2:length(n)], tasa, main="Tasa de Crecimiento", xlab="año", pch=16)
```

Tasa de Crecimiento



Ahora que aprendimos a usar un *for loop*, podemos ver que en este caso no lo necesitábamos! De hecho, una muy buena característica de R es que permite hacer operaciones sobre vectores, listas, etc. como veremos oportunamente. De todas formas, es importante aprender a manejar los *loops*.

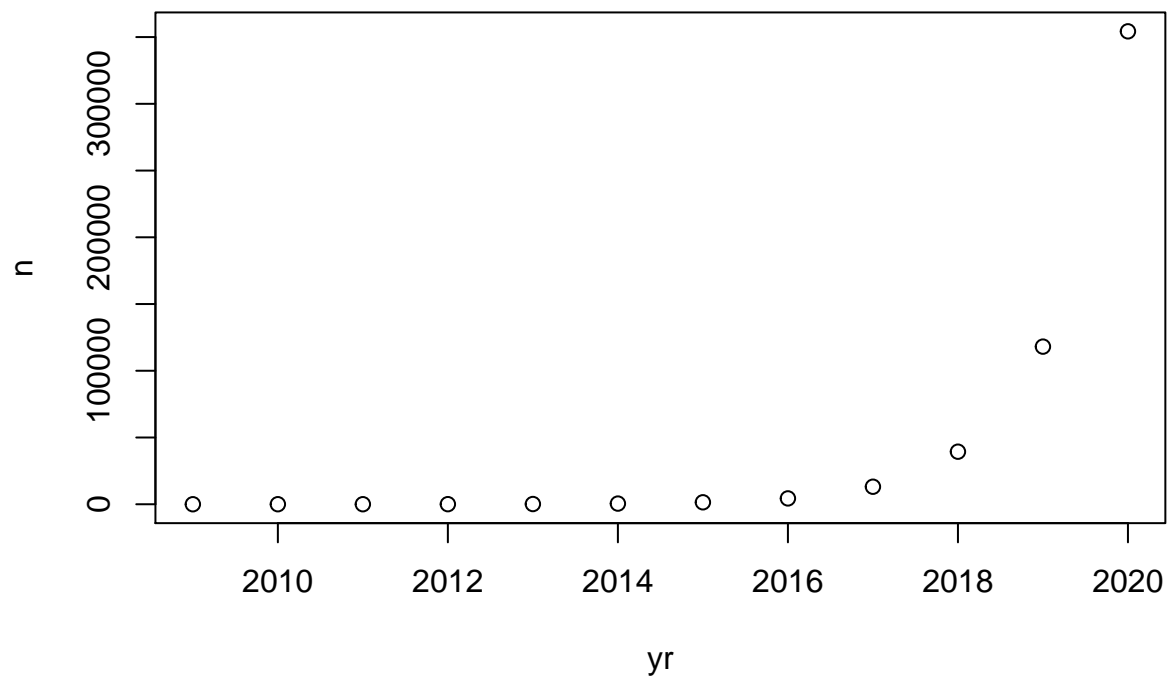
```
tasa <- n[2:6]/n[1:5]
```

En cualquier caso, vemos que la tasa de incremento poblacional es constante e igual a 3. Podemos usar ese valor para proyectar el tamaño poblacional en el futuro:

```
lambda <- 3
yr <- c(yr, 2015:2020)
n <- c(n, numeric(6))

for(i in 7:length(n)){
  n[i] <- n[i-1] * lambda
}

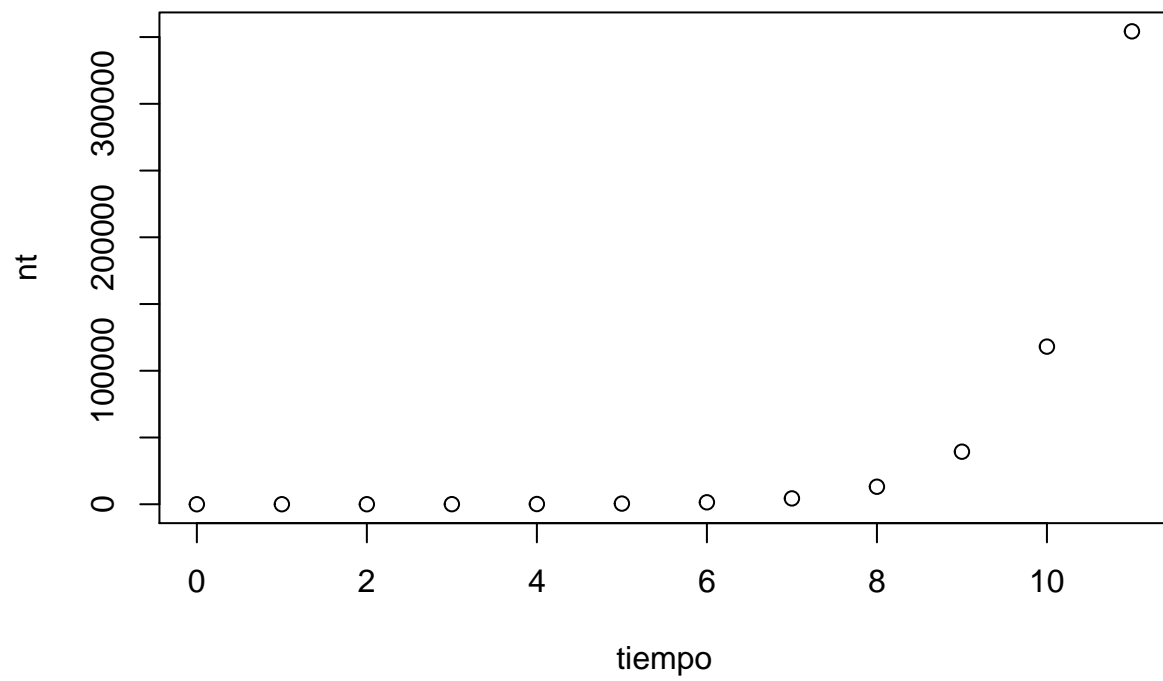
plot(yr,n)
```



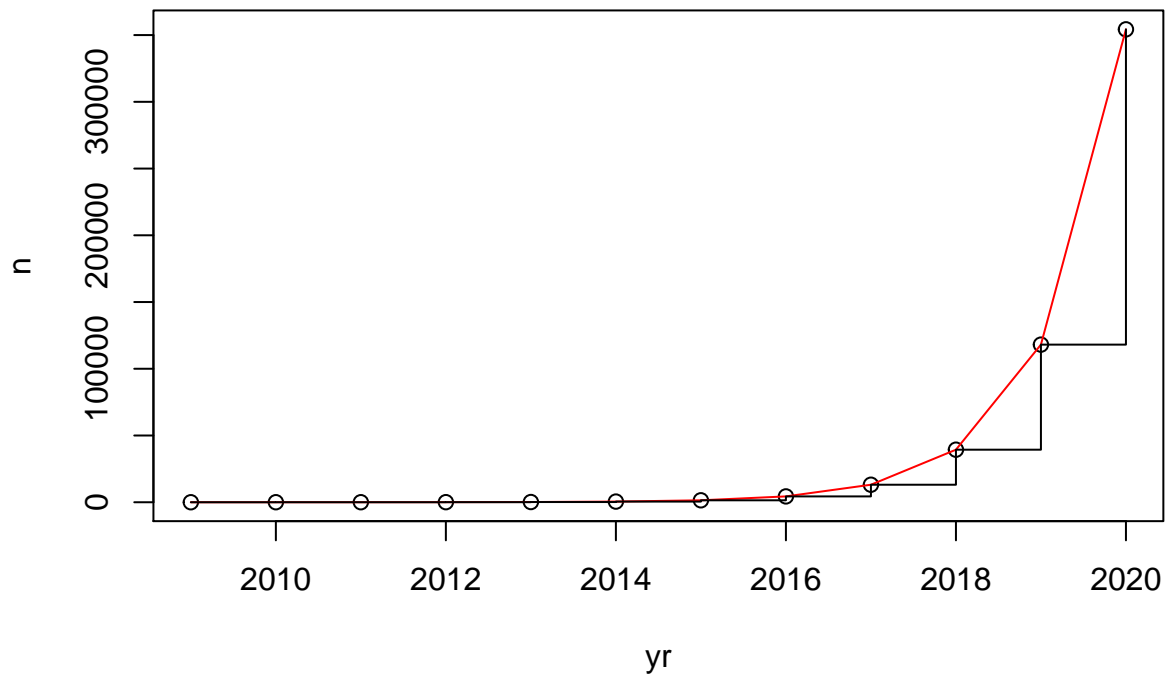
Lo mismo podríamos hacer definiendo condiciones iniciales y usando la fórmula $n(t) = n(0)\lambda^t$.

```
n0 <- 2
tiempo <- 0:11
nt <- n0 * lambda^tiempo

plot(tiempo, nt)
```



```
# o podemos graficar las dos trayectorias poblacionales  
plot(yr,n)  
lines(tiempo+yr[1], nt, col=2)  
lines(tiempo+yr[1], nt, type="s")
```

```
# 'lines' agrega lineas a un gr?fico existente
# 'col' es una opci?n de gr?ficos para definir colores
```

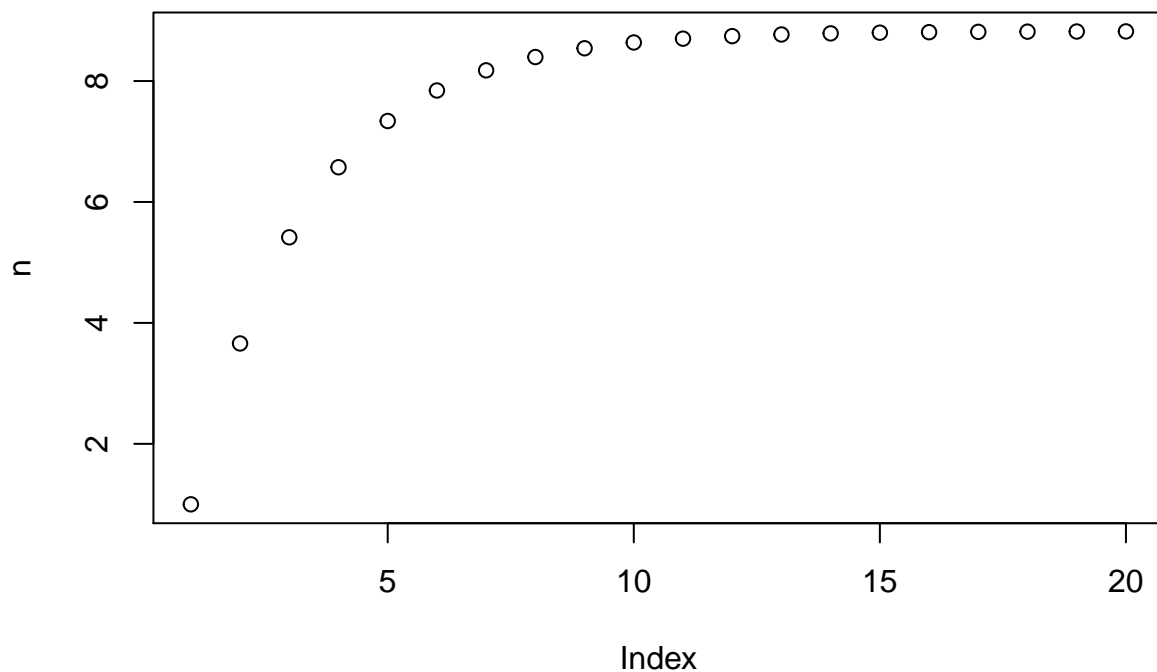
Ejercicios:

1. definir la escala de iteraci3n del modelo, revisar las unidades de variables y par3metros
2. cambiar n_0 y λ para ver qu3 pasa con la poblaci3n
3. graficar la trayectoria de la poblaci3n como $\log(n)$
4. Vimos un modelo de depredaci3n de monito del monte por gatos dom3sticos (diapositiva en margen superior derecho de la p3gina 4 en intro1.pdf). El modelo es $n(t+1) = n(t)(1-d)(1+b) + m$ +(a) definir la escala de iteraci3n del modelo, revisar las unidades de variables y par3metros asumiendo $d=0.4$, $b=0.1$ y $m=3$, +(b) modelar y graficar la trayectoria poblacional por 20 unidades de tiempo y con $n(0) = 1$ +(c) contrastar los resultados de (a) con un modelo en el cual el ciclo de vida es: censo, nacimientos, inmigraci3n, depredaci3n (usar los mismos valores para $d, b, m, n(0)$)

```

d=0.4
b=0.1
m=3
n <- numeric(20)
n[1] <- 1
for(t in 1:19){
  n[t+1] <- n[t]*(1-d)*(1+b)+m
}
plot(n)

```

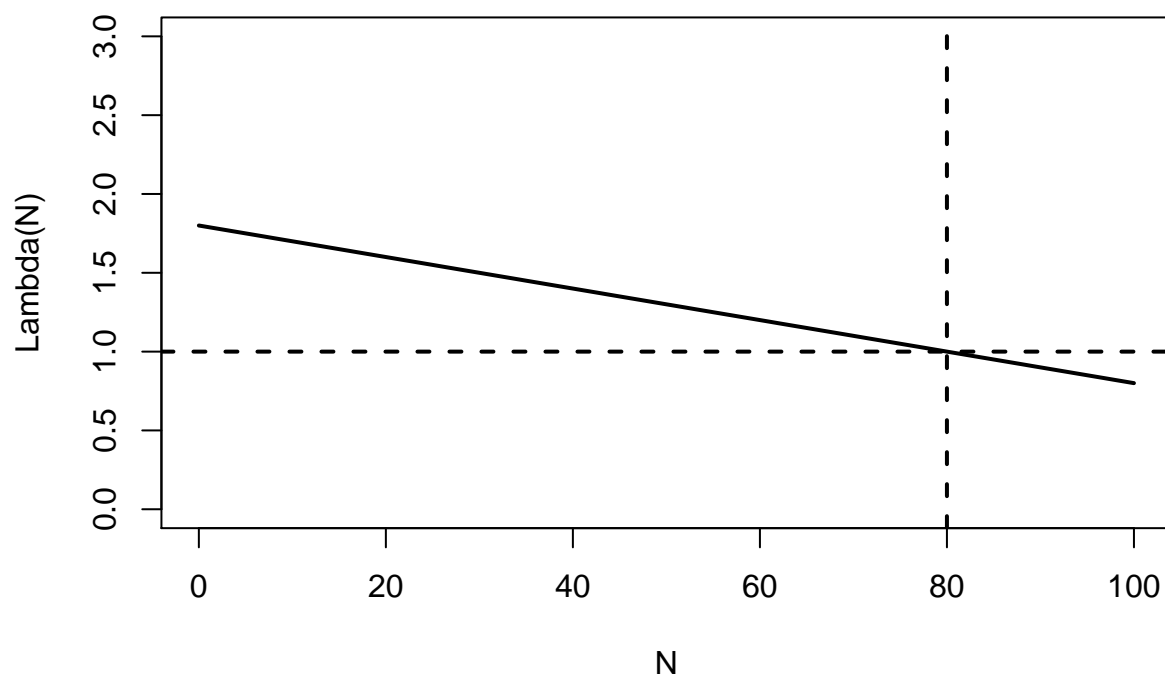


Es poco creíble que la población siga y siga creciendo como $n(t) = n(0)\lambda^t$. Una opción es pensar que la tasa de incremento poblacional decrece con el tamaño poblacional (modelo logístico).

```

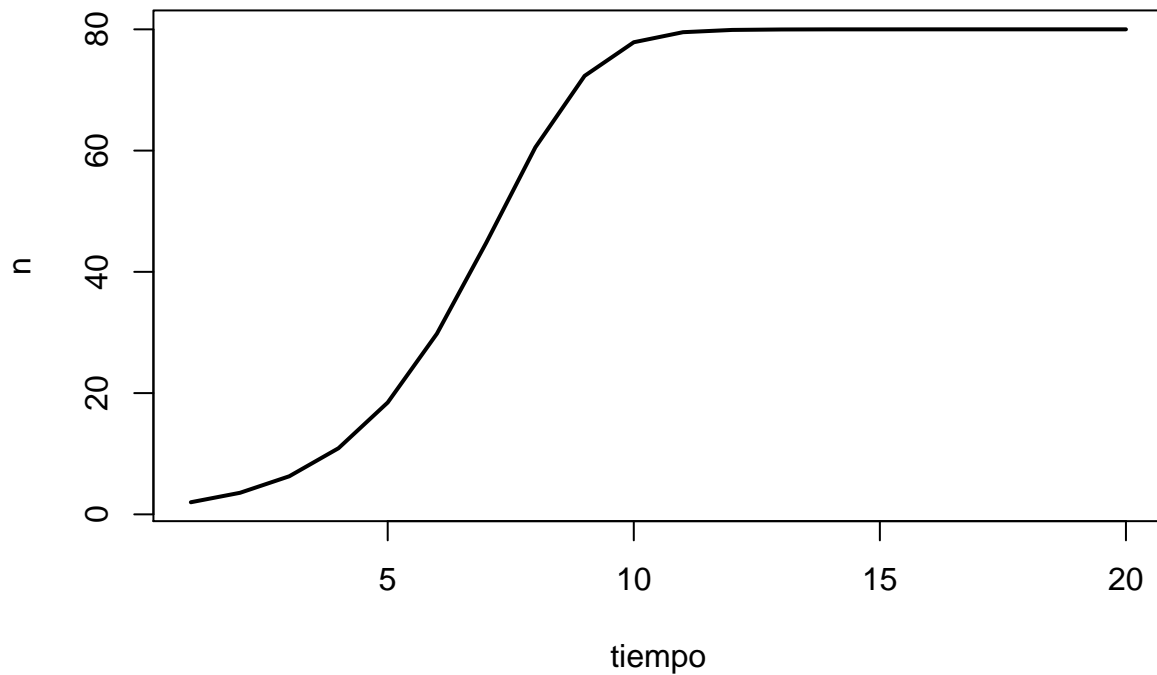
lambda <- 1.8
N <- 0:100
K <- 80 # capacidad de carga
rd <- lambda - 1 # rd es la tasa de crecimiento poblacional per cápita
plot(N, lambda - rd/K * N, type="l", ylab="Lambda(N)", ylim=c(0,3), lwd=2)
abline(h=1, lty=2, lwd=2)
abline(v=K, lty=2, lwd=2)

```



```
# simulemos una trayectoria
n0 <- 2
n <- numeric(20)
n[1] <- n0
for(t in 1:(length(n)-1)){
  n[t+1] <- n[t] + rd*n[t]*(1 - n[t]/K)
}

plot(1:length(n), n, type="l", lwd=2, xlab="tiempo")
```



Escribiendo Funciones

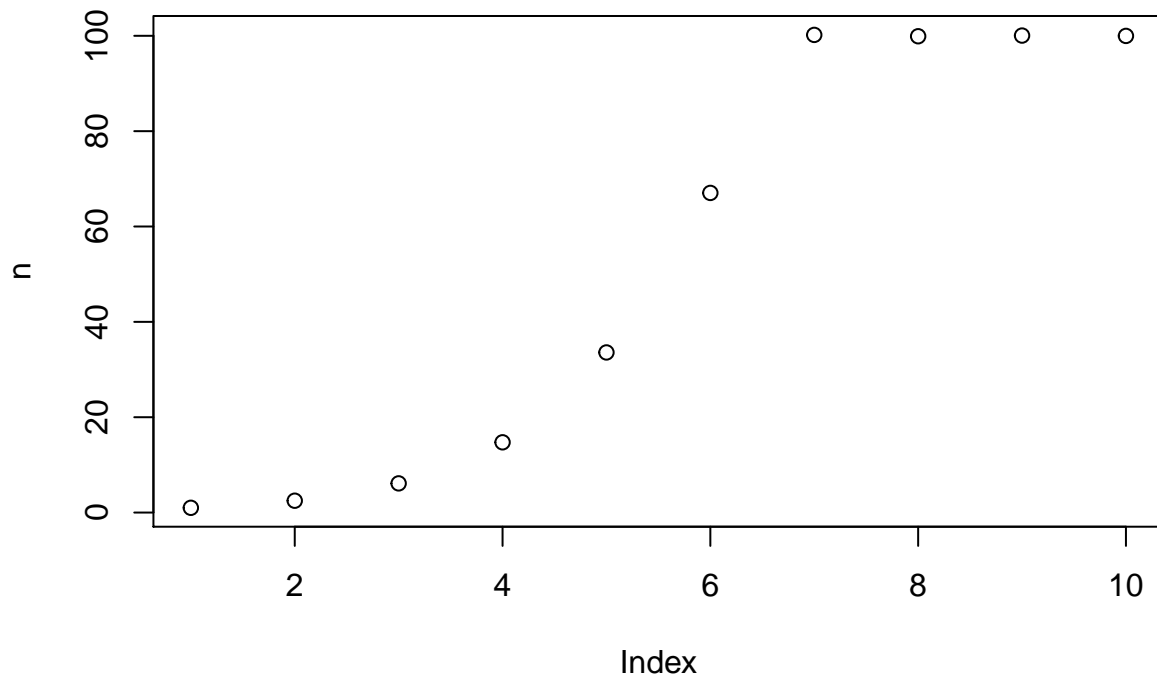
Un aspecto clave de la programación es escribir funciones. Una función es un conjunto de instrucciones que toma “inputs” o “argumentos”, usa esos inputs para calcular otros valores que devuelve como resultado. Por ejemplo, podemos escribir una función para el modelo logístico:

```
logistic <- function(n0,rd,K,nyr=10){  
  n <- numeric(nyr)  
  n[1] <- n0  
  for(t in 1:(length(n)-1)){  
    n[t+1] <- n[t] + rd*n[t]*(1 - n[t]/K)  
  }  
  return(n)  
}
```

Esta función tiene como argumentos a n_0 , rd , K y el número de años (nyr) que queremos simular.

A menos que le digamos otra cosa, R respeta el orden en que definimos los inputs de la función. Entonces, si queremos simular la trayectoria de una población que sigue un modelo logístico con $n_0=1$, $rd=1.5$, y $K=100$ por 10 generaciones, podemos hacer:

```
n <- logistic(1,1.5,100)
#y graficar
plot(n)
```



En este caso, no le dijimos a la función que queríamos $\text{nyr}=10$ porque está definido por defecto.

En realidad no es muy buena idea pasarle valores de esta manera a las funciones. Nos conviene ser más explícitos:

```
n <- logistic(n0=1, rd=1.5, K=100, nyr=20)
```

Incluso podemos cambiar el orden de los argumentos si es que los llamamos por su nombre:

```
n <- logistic(nyr=20, rd=1.5, n0=2, K=100)
```

Otras veces conviene definir los valores de los inputs aparte. Por ejemplo:

```
n0 <- 2
rd <- 1.5
K <- 100
nyears <- 20 # no necesariamente usamos el mismo nombre
n <- logistic(nyr=nyears, rd=rd, n0=n0, K=K)
```

Ejercicios:

(5) qué pasa si cambiamos $n(0)$?

(6) y si cambiamos rd ?

Ricker (1954), cambió la relación lineal entre λ y n por $\exp(-a \cdot n)$ (7) simular trayectorias poblacionales con el modelo de Ricker y compararlo con el logístico