

## Práctica 2

### *Introducción al uso de redes bayesianas con R*

Diciembre de 2024



DECSAI

## Modelos Gráficos Probabilísticos

Diciembre de 2024

# Índice

<b>1. Introducción al guion</b>	<b>3</b>
<b>2. El problema a modelar: Encuesta sobre uso de medios de transporte</b>	<b>3</b>
<b>3. Creación del DAG de la red bayesiana</b>	<b>5</b>
<b>4. Dibujando el DAG de la red bayesiana</b>	<b>7</b>
<b>5. Especificación de las probabilidades</b>	<b>9</b>
<b>6. Creación de la red bayesiana</b>	<b>11</b>
<b>7. Inferencia en redes bayesianas</b>	<b>12</b>
7.1. Obtención de las independencias expresadas por el DAG .	12
7.2. Inferencia exacta . . . . .	13
7.3. Inferencia aproximada . . . . .	17

## 1. Introducción al guion

Este guion está basado en el libro *Bayesian Networks With Examples in R*. Marco Scutari y Jean-Baptiste Denis. CRC Press. 2015.

Para crear y manipular redes bayesianas, usaremos principalmente el paquete **bnlearn** (**B**ayesian **n**etwork **l**earning) de **R**. Si no lo tienes instalado, instálalo:

```
> install.packages("bnlearn")
```

Puedes encontrar la documentación sobre este paquete en **CRAN** en <http://cran.r-project.org/web/packages/bnlearn/index.html>.

La web asociada a este paquete es <http://www.bnlearn.com/>

Comienza cargando la librería **bnlearn**:

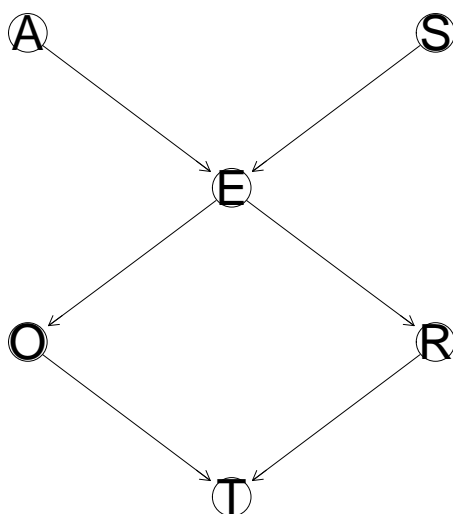
```
> library(bnlearn)
```

## 2. El problema a modelar: Encuesta sobre uso de medios de transporte

En este guion vamos a crear una red bayesiana para el problema siguiente. Se quiere hacer un estudio sobre los patrones de uso de diferentes medios de transporte por parte de una serie de personas, con el foco puesto en el coche y el tren. Para ello se realizará una encuesta. Para cada individuo, se examinarán seis variables:

- Edad ( $A$ ): puede tomar los valores *young* (menos de 30 años), *adult* (entre 30 y 60 años) y *old* (mayores de 60 años).
- Sexo ( $S$ ): puede ser  $M$  (hombre) o  $F$  (mujer).
- Nivel educativo ( $E$ ): el nivel de educación más alto alcanzado por el individuo, *high* (estudios superiores) o *uni* (universitarios).
- Ocupación ( $O$ ): indica si el individuo trabaja como empleado (*emp*) o como autónomo (*self*).
- Residencia ( $R$ ): el tamaño de la ciudad donde vive el individuo, *small* (ciudad pequeña) o *big* (ciudad grande).
- Travel ( $T$ ): el medio de transporte que suele usar el individuo, *car* (coche), *train* (tren) o *other* (otro).

Supondremos que este problema se puede modelar con el *Grafo Dirigido Acíclico* (DAG) que se muestra a continuación:



Suponemos las siguientes distribuciones de probabilidad:

$P(A)$	<i>young</i>	<i>adult</i>	<i>old</i>
	0,30	0,50	0,20

(1)

$P(S)$	<i>M</i>	<i>F</i>
	0,60	0,40

(2)

$P(E A, S)$	<i>high</i>	<i>uni</i>
<i>young, M</i>	0,75	0,25
<i>adult, M</i>	0,72	0,28
<i>old, M</i>	0,88	0,12
<i>young, F</i>	0,64	0,36
<i>adult, F</i>	0,70	0,30
<i>old, F</i>	0,90	0,10

(3)

$P(R E)$	<i>small</i>	<i>big</i>
<i>high</i>	0,25	0,75
<i>uni</i>	0,2	0,8

(4)

$P(O E)$	<i>emp</i>	<i>self</i>
<i>high</i>	0,96	0,04
<i>uni</i>	0,92	0,08

(5)

$P(T O, R)$	<i>car</i>	<i>train</i>	<i>other</i>
<i>small, emp</i>	0,48	0,42	0,10
<i>small, self</i>	0,56	0,36	0,08
<i>big, emp</i>	0,58	0,24	0,18
<i>big, self</i>	0,70	0,21	0,09

(6)

### 3. Creación del DAG de la red bayesiana

Crearemos el DAG de la red bayesiana, con un nodo por cada variable. Usamos la función `empty.graph` que permite generar un grafo dirigido acíclico vacío o aleatorio para un conjunto de nodos dado. Esta función crea el DAG como un objeto de la clase `bn`. En nuestro caso, crearemos un grafo con las variables *A*, *S*, *E*, *O*, *R* y *T*.

```
> dag <- empty.graph(nodes = c("A", "S", "E", "O",
                                "R", "T"))
```

El grafo anterior tiene por ahora un conjunto vacío de arcos. El DAG está almacenado en un objeto de la clase `bn`. Podemos imprimir en la consola el objeto que acabamos de crear:

```
> dag
```

El resultado sería el siguiente:

```
> dag
Random/Generated Bayesian network

model:
  [A] [S] [E] [O] [R] [T]
nodes:                                6
arcs:                                0
  undirected arcs:                    0
  directed arcs:                      0
average markov blanket size:          0.00
average neighbourhood size:           0.00
average branching factor:             0.00

generation algorithm:                 Empty
```

A continuación, añadiremos los arcos entre los nodos del DAG de la red bayesiana. Por ejemplo, el enlace entre las variables *A* y *E* lo añadimos con:

```
> dag <- set.arc(dag, from = "A", to = "E")
```

Añade tú el resto de arcos al DAG. Una vez hecho, si volvemos a imprimir nuestro DAG por la consola, obtendríamos:

```
> dag

Random/Generated Bayesian network

model:
```

```
[A] [S] [E|A:S] [O|E] [R|E] [T|O:R]
nodes:                                6
arcs:                                6
  undirected arcs:                    0
  directed arcs:                      6
average markov blanket size:          2.67
average neighbourhood size:           2.00
average branching factor:              1.00

generation algorithm:                  Empty
```

Como puede verse en la anterior salida, en el *model string* (modelo) existe un elemento por cada variable. Por ejemplo,  $[E|A:S]$  expresa que hay una dependencia directa de  $E$  con  $A$  y  $S$ , por lo que en el grafo tenemos los enlaces  $A \rightarrow E$  y  $S \rightarrow E$ . Sin embargo  $[A]$  expresa que no hay ningún arco apuntando hacia  $A$ .

Podemos usar la función `modelstring` para ver o modificar el modelo. Por ejemplo, para verlo usaríamos:

```
> modelstring(dag)
[1] "[A] [S] [E|A:S] [O|E] [R|E] [T|O:R]"
```

Podríamos haber definido los enlaces entre las variables con la función `modelstring`:

```
> modelstring(dag) = "[A] [S] [E|A:S] [O|E] [R|E] [T|O:R]"
```

Usando la función `model2network` podríamos haber creado también el DAG (variables y enlaces) especificando directamente el modelo, de la misma forma que con `modelstring`:

```
dag3 <- model2network("[A] [S] [E|A:S] [O|E] [R|E] [T|O:R]")
```

El DAG resultante del comando anterior (`dag3`), es idéntico al que hemos creado anteriormente (`dag`). Esto puede comprobarse con la función `all.equal`:

```
> all.equal(dag, dag3)
[1] TRUE
```

Podemos mostrar la lista de nodos del DAG de la red bayesiana con la función `nodes`:

```
> nodes(dag)
[1] "A" "S" "E" "O" "R" "T"
```

También podemos mostrar la lista de enlaces con la función `arcs`:

```
> arcs(dag)
      from to
[1,] "A"  "E"
[2,] "S"  "E"
[3,] "E"  "O"
[4,] "E"  "R"
[5,] "O"  "T"
[6,] "R"  "T"
```

Una forma alternativa para añadir los arcos a un DAG es mediante el uso de la función `arcs`. Por ejemplo, podemos crear un nuevo DAG con las mismas variables y arcos de la siguiente forma:

```
> dag2 <- empty.graph(nodes = c("A", "S", "E", "O",
                                "R", "T"))
> arc.set <- matrix(c("A", "E",
                      "S", "E",
                      "E", "O",
                      "E", "R",
                      "O", "T",
                      "R", "T"),
                    byrow = TRUE, ncol = 2,
                    dimnames = list(NULL, c("from", "to")))
> arcs(dag2) <- arc.set
```

Podemos comprobar si los dos DAGs creados hasta ahora son iguales con la función `all.equal`:

```
> all.equal(dag, dag2)
[1] TRUE
```

Las dos formas que hemos visto para añadir los arcos, garantizan que el DAG es acíclico. Por ello, si intentamos introducir un ciclo a propósito, se devolverá un error:

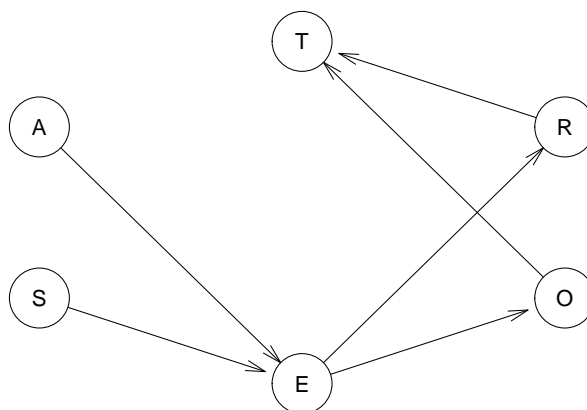
```
> try(set.arc(dag, from="T", to="E"))
Error in arc.operations(x = x, from = from, to = to,
  op = "set", check.cycles = check.cycles, :
  the resulting graph contains cycles.
```

## 4. Dibujando el DAG de la red bayesiana

Podemos mostrar un gráfico del DAG de la red bayesiana mediante la función `plot`:

```
> plot(dag)
```

El resultado tendría el siguiente aspecto:



Alternativamente podríamos usar el paquete **Rgraphviz** para mostrar el DAG. Este paquete no está en **CRAN** sino en **bioconductor** (<http://www.bioconductor.org/>). Podemos instalarlo de la siguiente forma:

1. Ejecutar el comando:

```
setRepositories()
```

y seleccionar el repositorio *BioC Software*.

2. Instalar el paquete **Rgraphviz** con el comando:

```
install.packages("Rgraphviz")
```

Si no funcionase lo anterior, otra alternativa para instalar **Rgraphviz** sería la siguiente:

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install(version = "3.14")
BiocManager::install(c("Rgraphviz"))
```

Si tuviésemos una versión de R anterior a la 3.5 entonces **Rgraphviz** se instalaría con:

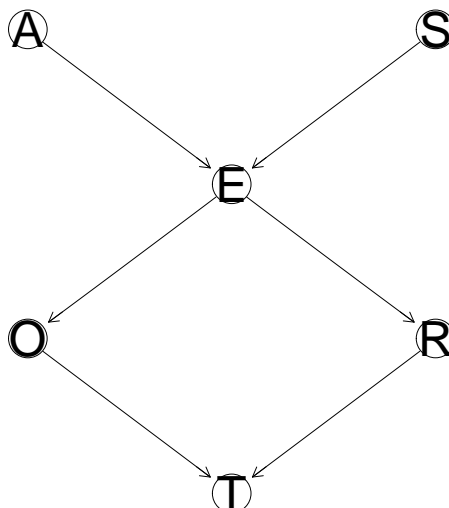
```
> source("http://bioconductor.org/biocLite.R")
> biocLite("Rgraphviz")
```



Una vez instalado **Rgraphviz**, mostramos nuestro DAG de la siguiente forma:

```
> graphviz.plot(dag)
```

El resultado tendría ahora el siguiente aspecto:



La función `graphviz.plot` tiene cuidado de colocar los nodos y arcos de forma que se minimice su solapamiento. Por defecto, los nodos se colocan de forma que los padres se dibujan encima de los hijos y la mayoría de los arcos apuntan hacia abajo. A esta forma de colocación (layout) se le llama *dot*. Podemos usar otros layouts como *fdp* y *circo*.

## 5. Especificación de las probabilidades

En la red bayesiana que estamos creando, todas las variables son discretas y por tanto cada una está definida para un número finito y no ordenado de *estados* (el dominio de esa variable). Vamos a definir los estados de las variables de la red bayesiana con los siguientes objetos:

```
> A.st <- c("young", "adult", "old")
> S.st <- c("M", "F")
> E.st <- c("high", "uni")
> O.st <- c("emp", "self")
> R.st <- c("small", "big")
> T.st <- c("car", "train", "other")
```

El DAG que hemos construido hasta ahora nos indica que la distribución de probabilidad conjunta se factoriza de la siguiente forma:

$$P(A, S, E, O, R, T) = P(A)P(S)P(E|A, S)P(O|E)P(R|E)P(T|O, R)$$

Por tanto, debemos definir las seis distribuciones de probabilidad condicional de la anterior factorización. Las variables *A* y *S* se definen con las

distribuciones de probabilidad unidimensionales (no tienen padres)  $P(A)$  y  $P(S)$ , las cuales podemos definir con un array:

```
> A.prob <- array(c(0.30, 0.50, 0.20), dim=3,
  dimnames = list(A = A.st))
```

Podemos mostrar ahora las probabilidades que acabamos de definir para  $P(A)$ :

```
> A.prob
A
young adult    old
   0.3    0.5   0.2
```

De igual forma, para  $P(S)$  haríamos:

```
> S.prob <- array(c(0.60, 0.40), dim=2,
+                  dimnames = list(S = S.st))
>
> S.prob
S
  M  F
0.6 0.4
```

Las distribuciones condicionales  $P(O|E)$  y  $P(R|E)$  se definen con:

```
> O.prob <- array(c(0.96, 0.04, 0.92, 0.08), dim=c(2,2),
+                  dimnames = list(O = O.st, E = E.st))
> O.prob
      E
O      high uni
emp  0.96 0.92
self 0.04 0.08
```

```
> R.prob <- array(c(0.25, 0.75, 0.20, 0.80), dim=c(2,2),
+                  dimnames = list(R = R.st, E = E.st))
> R.prob
      E
R      high uni
small 0.25 0.2
big   0.75 0.8
```

Para las distribuciones unidimensionales y bidimensionales, podemos usar también la función `matrix`. La sintaxis es casi idéntica a la de `array`; la única diferencia es que se debe especificar solo una dimensión (o el número de filas `nrow`, o el número de columnas `ncol`). Por ejemplo, para  $P(R|E)$  podríamos haber ejecutado:

```
> R.prob <- matrix(c(0.25, 0.75, 0.20, 0.80), ncol = 2,
  dimnames = list(R = R.st, E = E.st))
```

Las distribuciones  $P(E|A, S)$  y  $P(T|O, R)$  las definimos con:

```
> E.prob <- array(c(0.75, 0.25, 0.72, 0.28, 0.88, 0.12, 0.64,
  0.36, 0.70, 0.30, 0.90, 0.10), dim=c(2, 3, 2),
  dimnames = list(E = E.st, A = A.st, S = S.st))
```

```
> T.prob <- array(c(0.48, 0.42, 0.10, 0.56, 0.36, 0.08, 0.58,
  0.24, 0.18, 0.70, 0.21, 0.09), dim=c(3, 2, 2),
  dimnames = list(T = T.st, O = O.st, R = R.st))
```

## 6. Creación de la red bayesiana

Una vez definidas todas las distribuciones de probabilidad necesarias, construimos una lista (que llamaremos `cpt`) con ellas.

```
> cpt <- list(A=A.prob, S=S.prob, E=E.prob, O=O.prob,
  R=R.prob, T=T.prob)
```

Con esto, ya tenemos definido todo lo necesario para crear nuestra red bayesiana mediante un objeto de la clase `bn.fit` que llamaremos `bn`. Para ello, usaremos la función `custom.fit`:

```
> bn <- custom.fit(dag,cpt)
```

El número de parámetros de la red bayesiana puede ser calculado con la función `nparams`, que en este caso nos da el valor 21, como puede comprobarse observando las 6 tablas de distribuciones de probabilidad de la red. Por ejemplo, el número de parámetros en la tabla para  $P(A)$  sería de 2, ya que contiene tres probabilidades, pero solo 2 son independientes; la tercera puede obtenerse a partir de las otras ya que deben sumar 1.

```
> nparams(bn)
[1] 21
```

Los objetos de la clase `bn.fit` se usan para describir redes bayesianas en el paquete `bnlearn`. Estos objetos incluyen información sobre el DAG (padres e hijos de cada nodo) y las distribuciones de probabilidad locales (sus parámetros). Los objetos `bn.fit` pueden usarse como si fuesen objetos de la clase `bn` a la hora de investigar sus propiedades gráficas. Por ejemplo, podemos usar `arcs` con un objeto `bn.fit`:

```
> arcs(bn)
      from to
[1,] "A"  "E"
[2,] "S"  "E"
[3,] "E"  "O"
[4,] "E"  "R"
[5,] "O"  "T"
[6,] "R"  "T"
```

Podríamos usar de la misma forma las funciones `nodes`, `parents` y `children` de la clase `bn`.

Podemos también mostrar las tablas de probabilidad condicional del objeto `bn.fit` de la siguiente forma:

```
> bn$R

Parameters of node R (multinomial distribution)

Conditional probability table:

      E
R      high  uni
small 0.25 0.20
big   0.75 0.80
```

Podemos extraer una tabla de distribución de probabilidad condicional para su uso posterior con la función `coef`:

```
> R.cpt <- coef(bn$R)
```

Para mostrar en la consola todas las distribuciones de probabilidad condicional ejecutamos:

```
> bn
```

## 7. Inferencia en redes bayesianas

### 7.1. Obtención de las independencias expresadas por el DAG

En una red bayesiana, el DAG define el siguiente tipo de independencia entre variables. Dos conjuntos de variables  $X$  y  $Y$  son independientes dado un tercer conjunto  $Z$  si todos los caminos que conectan variables de  $X$  con  $Y$  están *bloqueados* o *separados* por las variables de  $Z$ . Para ello se aplica el criterio de *D-separación*. Las independencias expresadas gráficamente en el DAG con este criterio, se dan también como independencias probabilísticas entre las variables de la red bayesiana.

La función `dsep` permite saber si dos variables  $x$  e  $y$  del DAG están d-separadas dada otra variable  $z$  o un vector de variables  $z$ . Por ejemplo, para ver si  $S$  y  $R$  son independientes ( $S$  y  $R$  están d-separadas sin ninguna evidencia) ejecutamos:

```
> dsep(dag, x = "S", y = "R")
[1] FALSE
```

De igual forma, con  $O$  y  $R$ :

```
> dsep(dag, x = "O", y = "R")
[1] FALSE
```

La función `path.exists` permite ver si existe un camino dirigido de un nodo a otro. Por ejemplo:

```
> path.exists(dag, from = "S", to = "R")
[1] TRUE
```

El camino entre  $S$  y  $R$  queda bloqueado si condicionamos sobre  $E$  ( $S$  y  $R$  se vuelven independientes conocido el valor de  $E$ ) como puede verse con:

```
> dsep(dag, x = "S", y = "R", z = "E")
[1] TRUE
```

Esta independencia indica que:

$$P(S, R|E) = P(S|E) \cdot P(R|E)$$

Otro ejemplo,  $S$  y  $T$  son independientes conocido el valor de  $O$  y de  $R$ :

```
> dsep(dag, x="S", y="T", z=c("O", "R"))
[1] TRUE
```

Las variables  $O$  y  $R$  también están d-separadas dada la variable  $E$ :

```
> dsep(dag, x = "O", y = "R", z = "E")
[1] TRUE
```

Por otro lado, el condicionar sobre otras variables, puede hacer que variables que eran independientes marginalmente (sin evidencia) se vuelvan dependientes. Por ejemplo las variables  $A$  y  $S$  son independientes marginalmente, pero se vuelven dependientes al condicionar sobre  $E$ .

```
> dsep(dag, x = "A", y = "S")
[1] TRUE
> dsep(dag, x = "A", y = "S", z = "E")
[1] FALSE
```

## 7.2. Inferencia exacta

La inferencia exacta está implementada en el paquete **gRain** (**gR**aphical **m**odel **i**nference). Este paquete está en **bionconductor** y lo podemos instalar con:

```
install.packages("gRbase", dependencies=TRUE);
install.packages("gRain", dependencies=TRUE);
```

Alternativamente podríamos emplear el siguiente comando:

```
BiocManager::install(c("gRbase", "gRain"))
```

Si tenemos una versión de R anterior a la 3.5 usamos:

```
> source("http://bioconductor.org/biocLite.R")
> biocLite(c("graph", "RBGL", "Rgraphviz"))
> install.packages("gRain")
```

Finalmente, cargamos el paquete **gRain**:

```
> library(gRain)
```

El paquete **gRain** fue desarrollado por Søren Højsgaard. Podemos encontrar información sobre éste y otros de sus paquetes en <http://people.math.aau.dk/~sorenh/software/gR/>. También conviene echar un vistazo a la web <https://cran.r-project.org/web/views/GraphicalModels.html> que da una revisión general a los paquetes **CRAN** para modelos gráficos.

El algoritmo exacto que contiene el paquete **gRain** es el algoritmo de Lauritzen-Spiegelhalter que es un algoritmo basado en la construcción de un árbol de grupos (*junction tree*) y es muy similar a los algoritmos de Hugin y el de Shenoy-Shafer. Puedes ver una comparativa de estos algoritmos en el artículo *A Comparison of Lauritzen-Spiegelhalter, Hugin, and Shenoy-Shafer Architectures for Computing Marginals of Probability Distributions* disponible en <http://arxiv.org/pdf/1301.7394.pdf>.

Para construir el árbol de grupos con el paquete **gRain**, debemos convertir nuestra red bayesiana, que es de la clase `bn.fit`, en un objeto de la clase `grain`. Para ello usaremos la función `as.grain` del paquete **bn-learn**. Con el objeto `grain`, usaremos la función `compile` para construir el árbol de grupos y hacer la propagación:

```
> junction <- compile(as.grain(bn))
```

La función `querygrain` permite obtener la distribución de probabilidad a posteriori para una variable dada la evidencia. Por ejemplo, para obtener  $P(T|evidencia)$  usaremos:

```
> querygrain(junction, nodes="T")$T
      car      train      other
0.5618340 0.2808573 0.1573088
```

En la operación anterior, todavía no hemos introducido evidencia en la red bayesiana, por lo que nos ha mostrado  $P(T)$  para cada estado de  $T$ .

Una vez construido el árbol de grupos y calculadas las tablas de probabilidad a posteriori, podemos definir la evidencia disponible mediante la función `setEvidence`. Esto hará que automáticamente se propague dicha evidencia en el árbol de grupos. Por ejemplo, podemos definir la evidencia  $S = F$ :

```
> jsex <- setEvidence(junction, nodes="S", states="F")
```

Para obtener de nuevo  $P(T|evidencia)$  ( $P(T|S = F)$  en este caso) usaremos:

```
> querygrain(jsex,nodes="T")$T
      car      train      other
0.5620577 0.2806144 0.1573280
```

Puede verse que las probabilidades obtenidas para  $T$  con la evidencia  $S = F$ , no han cambiado sustancialmente respecto a cuando no había evidencia. Esto nos indica que las mujeres muestran más o menos las mismas preferencias hacia el uso de coche o tren que la población entrevistada como un todo.

Hagamos ahora una consulta para conocer cómo afecta el vivir en una ciudad pequeña al uso de coche o tren, o sea, queremos calcular  $P(T|R = \text{small})$ . La gente que trabaja en grandes ciudades a menudo vive en pueblos vecinos y viajan diariamente a sus lugares de trabajo, ya que los precios de las casas son más bajos a medida que te alejas hacia el campo. Esto fuerza al uso de tren o coche para viajar a la ciudad, ya que otros medios de transporte (bicicleta, metro, líneas de autobús, etc) no están disponibles o no son prácticas.

```
> jres <- setEvidence(junction, nodes = "R", states = "small")
> querygrain(jres, nodes = "T")$T
T
      car      train      other
0.48388675 0.41708494 0.09902831
```

Después de las anteriores consultas, tenemos:

	$P(T)$	$P(T S = F)$	$P(T R = small)$
car	0.5618340	0.5620577	0.48388675
train	0.2808573	0.2806144	0.41708494
other	0.1573088	0.1573280	0.09902831

Como podemos ver en la tabla anterior,  $P(T = other)$  pasa de 0,1573 (población general) a 0,099 (gente que vive en pequeñas ciudades), mientras que  $P(T = train)$  pasa de 0,2808 (población general) a 0,4170 (gente que vive en pequeñas ciudades). La probabilidad combinada de usar coche o tren ( $P(T = car \cup T = train)$ ) pasa de 0,8426 (población general) a 0,9009 (gente que vive en pequeñas ciudades).

La tabla anterior, también nos permite saber cual es la *explicación más probable* para la variable  $T$  en cada caso. Por ejemplo, para la gente que vive en pequeñas ciudades, el coche es el medio de transporte preferido (probabilidad igual a 0,48388675).

Es posible también obtener la probabilidad conjunta condicional dada la evidencia (probabilidad conjunta a posteriori) para más de una variable a la vez. Por ejemplo, para obtener  $P(S, T|E = high)$  haríamos:

```
> jedu <- setEvidence(junction, nodes = "E", states = "high")
> SxT.cpt <- querygrain(jedu, nodes = c("S", "T"),
+                        type = "joint")
> SxT.cpt
```

	T			
S	car	train	other	
M	0.3426644	0.1736599	0.09623271	
F	0.2167356	0.1098401	0.06086729	

Esta consulta, nos permite también ver que la *explicación más probable* para las dos variables  $S$  (sexo) y  $T$  (medio de transporte usado) conjuntamente cuando tienen un nivel de estudios *high* (alto) es que son hombres que usan el coche.

El argumento `type` en la función `querygrain` especifica el tipo de consulta que queremos. El valor por defecto es "marginal", que nos da la distribución marginal a posteriori para cada variable. Usando "marginal" en el anterior caso obtendríamos la marginal para  $S$  y la marginal para  $T$

```
> querygrain(jedu, nodes = c("S", "T"), type = "marginal")
$S
S
      M      F
0.612557 0.387443

$T
T
  car train other
0.5594 0.2835 0.1571
```

Otra posibilidad para el argumento `type` de la función `querygrain` es usar el valor `conditional`. Esto devuelve la distribución de probabilidad de la primera variable en `nodes` condicionado al resto de variables en `nodes` y a la evidencia especificada con `setEvidence`. Por ejemplo, para calcular  $P(S|T, E = high)$  ejecutamos:

```
> querygrain(jedu, nodes = c("S", "T"), type = "conditional")
T
S
  car train other
M 0.612557 0.612557 0.612557
F 0.387443 0.387443 0.387443
```

Los resultados obtenidos en esta consulta muestran que todas las probabilidades

$$P(S = M|T = t, E = high), \quad t \in \{car, train, other\}$$

son idénticas, independientemente del valor de  $T$ . Lo mismo ocurre cuando  $S$  es igual a  $F$ . En otras palabras, esto nos indica que:

$$P(S = M|T = t, E = high) = P(S = M|E = high)$$

y

$$P(S = F|T = t, E = high) = P(S = F|E = high)$$

Esto sugiere que  $S$  es independiente de  $T$  condicionado al valor de  $E$ ; conocer el sexo de una persona no es informativo sobre sus preferencias en el transporte si conocemos su nivel educativo. Esta conclusión podemos también extraerla con el criterio de la d-separación, puesto que  $S$  y  $T$  están d-separados por  $E$ .



```
> dsep(bn, x = "S", y = "T", z = "E")
[1] TRUE
```

### 7.3. Inferencia aproximada

El paquete **bnlearn** tiene implementados algoritmos aproximados de inferencia por Monte Carlo para redes bayesianas. Los algoritmos que tiene implementados son el de *muestreo lógico* (rejection sampling, logic sampling o muestreo por rechazo) y el de *ponderación por verosimilitud* (likelihood weighting). Para usar el algoritmo de *muestreo lógico* usaremos la función `cpquery`. Por ejemplo, para calcular la distribución conjunta a posteriori para  $S = M$  y  $T = car$  ejecutamos:

```
> cpquery(bn, event = (S == "M") & (T == "car"),
+          evidence = (E == "high"))
[1] 0.351592
```

Podemos ver que esta probabilidad es algo distinta a la exacta que obtuvimos anteriormente con `querygrain` (0,3426644). La calidad de la aproximación puede mejorarse usando el argumento `n` en `cpquery` para incrementar el número de muestras aleatorias que utilizamos. El valor por defecto es  $5000 \cdot nparams(bn)$ . Lo pondremos a un millón:

```
> cpquery(bn, event = (S == "M") & (T == "car"),
+          evidence = (E == "high"), n = 10^6)
[1] 0.343198
```

Podemos observar que el tiempo empleado es mayor que en el caso anterior. Como se vió en la clase de teoría, el muestreo lógico da malas estimaciones cuando la probabilidad de la evidencia tiene baja probabilidad, ya que este método rechaza las muestras que no son consistentes con la evidencia.

Una aproximación que suele funcionar mejor es el método de *ponderación por verosimilitud*, que funciona de forma que ninguna muestra es rechazada. Puede usarse con la función `cpquery` usando `method="lw"`. El método por defecto es el muestreo lógico (`method="ls"`). Por ejemplo, para la consulta anterior de  $P(S = M, T = car | E = high)$  ejecutamos:

```
> cpquery(bn, event = (S == "M") & (T == "car"),
+          evidence = list(E = "high"), method = "lw")
[1] 0.3422035
```

Podemos ver que el valor obtenido 0,3422035 es más cercano al exacto 0,3426644 que el obtenido que el muestreo lógico (0,343198) sin tener que generar  $10^6$  muestras aleatorias. En el anterior comando, podemos observar que la evidencia `evidence` se proporciona al algoritmo de ponderación por similitud como una lista de valores, una para cada variable de condicionamiento.

El algoritmo para muestreo lógico implementado en **bnlearn** permite hacer consultas más complejas. Por ejemplo la probabilidad de que un hombre viaje en coche dado que su edad es *young* (joven) y su nivel educativo es *uni* (universitario) o bien él es un *adult* (adulto) independientemente de su nivel educativo:

$$P(S = M, T = \text{car} | \{A = \text{young}, E = \text{uni}\} \cup \{A = \text{adult}\})$$

Esta consulta la haríamos con:

```
> cpquery(bn, event = (S == "M") & (T == "car"),
+           evidence = ((A == "young") & (E == "uni")) | (A == "adult"))
[1] 0.3336759
```

El algoritmo de ponderación por verosimilitud no está implementado para permitir hacer el anterior tipo de consulta.

La función `cpdist` del paquete **bnlearn** devuelve un *data frame* que contiene un conjunto de muestras consistentes con la evidencia (*evidence*) para las variables indicadas en *nodes*. Por ejemplo:

```
> SxT <- cpdist(bn, nodes = c("S", "T"),
+               evidence = (E == "high"))
> head(SxT)
  S      T
1 M   car
2 M train
3 M   car
4 M   car
5 M train
6 F other
```

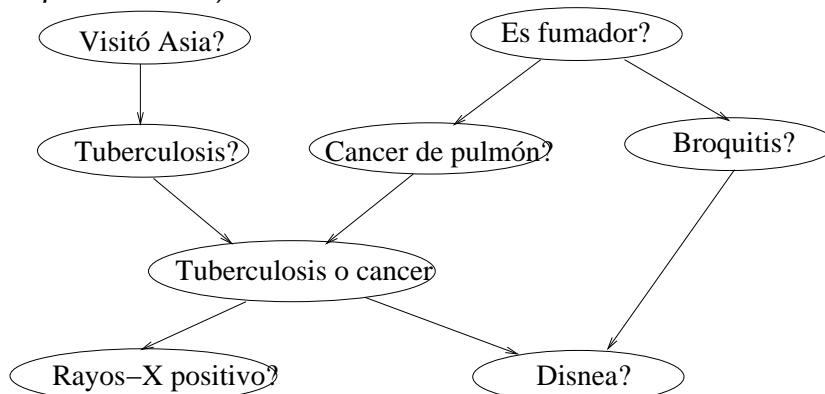
El data frame obtenido puede usarse ahora, por ejemplo, para obtener la probabilidad conjunta a posteriori de *S* y *T* ( $P(S, T | E = \text{high})$ ) usando la función `prop.table` de **R**:

```
> prop.table(table(SxT))
      T
S      car      train      other
M 0.35464961 0.16653344 0.09805489
F 0.21209699 0.10285105 0.06581402
```

La anterior consulta nos permite también obtener la *configuración más probable* para *S* y para *T* dada la evidencia  $E = \text{high}$ . De nuevo, la respuesta es que entre la gente con nivel de educación alto (*high*) la combinación más frecuente para *S* y *T* es *hombre* (*M*) y *coche* (*car*).

**Ejercicio 1** *Hacer un script en R para construir la red bayesiana (DAG y tablas de probabilidad) usada como ejemplo en esta práctica. El script debe contener también los comandos para dibujar la red bayesiana y para hacer las consultas (d-separación, probabilidad a posteriori) que se han explicado en la sección 7.3 usando el método de inferencia exacto y los dos aproximados.*

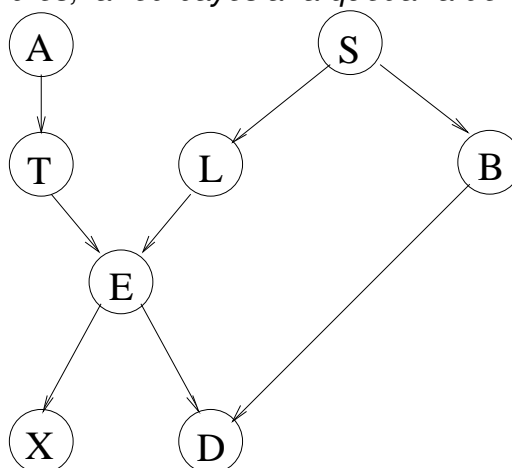
**Ejercicio 2** Construye un script en R para construir la red Asia (DAG y tablas de probabilidad).



Usaremos los siguientes nombres para las variables:

- Variable Visitó Asia:  $A$
- Variable Es fumador:  $S$
- Variable Tuberculosis:  $T$
- Variable Cáncer de pulmón:  $L$
- Variable Tuberculosis o cáncer de pulmón:  $E$
- Variable Bronquitis:  $B$
- Variable Rayos X positivo:  $X$
- Variable Disnea:  $D$

Con estos nombres, la red bayesiana quedaría de la siguiente forma:



Todas las variables tienen dos estados: yes y no.  
Los valores de probabilidad son los siguientes:

- $P(A=\text{yes}) = 0.1$
- $P(T=\text{yes})|A=\text{yes}) = 0.05$
- $P(T=\text{yes})|A=\text{no}) = 0.01$

- $P(S=\text{yes}) = 0.5$
- $P(L=\text{yes} \mid S=\text{yes}) = 0.1$
- $P(L=\text{yes} \mid S=\text{no}) = 0.01$
- $P(B=\text{yes} \mid S=\text{yes}) = 0.6$
- $P(B=\text{yes} \mid S=\text{no}) = 0.3$
- $P(E=\text{yes} \mid L=\text{yes}, T=\text{yes}) = 1$
- $P(E=\text{yes} \mid L=\text{yes}, T=\text{no}) = 1$
- $P(E=\text{yes} \mid L=\text{no}, T=\text{yes}) = 1$
- $P(E=\text{yes} \mid L=\text{no}, T=\text{no}) = 0$
- $P(X=\text{yes} \mid E=\text{yes}) = 0.98$
- $P(X=\text{yes} \mid E=\text{no}) = 0.05$
- $P(D=\text{yes} \mid E=\text{yes}, B=\text{yes}) = 0.9$
- $P(D=\text{yes} \mid E=\text{yes}, B=\text{no}) = 0.7$
- $P(D=\text{yes} \mid E=\text{no}, B=\text{yes}) = 0.8$
- $P(D=\text{yes} \mid E=\text{no}, B=\text{no}) = 0.1$

*El script debe también hacer consultas para obtener lo siguiente:*

1. *Tablas de probabilidad marginal de padecer tuberculosis, de padecer cáncer de pulmón y de padecer bronquitis. Usa el método exacto y los dos aproximados.*
2. *Tablas de probabilidad marginal a posteriori de padecer cada una de las tres enfermedades anteriores dado que se sabe que el paciente visitó Asia. Usa el método exacto y los dos aproximados.*
3. *Tablas de probabilidad marginal a posteriori de padecer cada una de las tres enfermedades anteriores dado que se sabe que el paciente visitó Asia, y tiene asma. Usa el método exacto y los dos aproximados.*
4. *Tabla de probabilidad conjunta de padecer cáncer de pulmón y bronquitis. Indica cual es la probabilidad de la configuración más probable para las variables de estas dos enfermedades.*
5. *Tabla de probabilidad conjunta de padecer cáncer de pulmón y bronquitis dado que el paciente no visitó Asia y el paciente es fumador. Indica cual es la probabilidad de la configuración más probable para las variables de estas dos enfermedades cuando el paciente no visitó Asia y el paciente es fumador.*

6. *Si no se conoce el valor de ninguna variable, ¿qué comando(s) podemos usar para saber si visitar Asia influye en tener cáncer de pulmón o no?*
7. *Sabiendo que el paciente es fumador, ¿qué comando(s) podemos usar para saber si dar positivo en una prueba de rayos X, puede influir en tener bronquitis?*
8. *La variable Tuberculosis o cáncer tiene asociada una distribución de probabilidad condicional que es una función determinista de los padres, concretamente una puerta or. Sabiendo que el paciente tiene tuberculosis, obtén las tablas de probabilidad marginal de padecer cáncer de pulmón y las tablas de probabilidad marginal a posteriori de padecer cáncer de pulmón si además se sabe que el paciente da positivo en la prueba de rayos X. ¿Puedes comentar si ha influido en la probabilidad de padecer cáncer de pulmón, el conocer que la prueba de rayos X ha sido positiva?*