

# Master en Big Data. Fundamentos Matemáticos del Análisis de Datos (FMAD).

## Tarea 2

Departamento de Matemática Aplicada

Curso 2021-22. Última actualización: 2021-09-17

### Instrucciones preliminares

- Empieza abriendo el proyecto de RStudio correspondiente a tu repositorio personal de la asignatura.
- En todas las tareas tendrás que repetir un proceso como el descrito en la sección *Repite los pasos Creando un fichero Rmarkdown para esta práctica* de la *Práctica00*. Puedes releer la sección *Practicando la entrega de las Tareas* de esa misma práctica para recordar el procedimiento de entrega.

### Ejercicio 1. Simulando variables aleatorias discretas.

**Apartado 1:** La variable aleatoria discreta  $X1$  tiene esta tabla de densidad de probabilidad (es la variable que se usa como ejemplo en la Sesión ):

valor de $X1$	0	1	2	3
Probabilidad de ese valor $P(X = x_i)$	$\frac{64}{125}$	$\frac{48}{125}$	$\frac{12}{125}$	$\frac{1}{125}$

Calcula la media y la varianza teóricas de esta variable.

**Apartado 2:** Combina `sample` con `replicate` para simular cien mil muestras de tamaño 10 de esta variable  $X1$ . Estudia la distribución de las medias muestrales como hemos hecho en ejemplos previos, ilustrando con gráficas la distribución de esas medias muestrales. Cambia después el tamaño de la muestra a 30 y repite el análisis.

**Apartado 3:** La variable aleatoria discreta  $X2$  tiene esta tabla de densidad de probabilidad:

valor de $X2$	0	1	2
Probabilidad de ese valor $P(X = x_i)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$

Suponemos que  $X1$  y  $X2$  son independientes. ¿Qué valores puede tomar la suma  $X1 + X2$ ? ¿Cuál es su tabla de probabilidad?

**Apartado 4:** Calcula la media teórica de la suma  $X1 + X2$ . Después usa `sample` y `replicate` para simular cien mil *valores* de esta variable suma. Calcula la media de esos valores. *Advertencia:* no es el mismo tipo de análisis que hemos hecho en el segundo apartado.

### Ejercicio 2. Datos limpios

- Descarga el fichero de este enlace

<https://gist.githubusercontent.com/fernandosanseguno/471b4887737cfcec7e9cf28631f2e21e/raw/b3944599d02df494f5903740db5acac9da35bc6f/testResults.csv>

- Este fichero contiene las notas de los alumnos de una clase, que hicieron dos tests cada semana durante cinco semanas. La tabla de datos no cumple los principios de *tidy data* que hemos visto en clase. Tu tarea en este ejercicio es explicar por qué no se cumplen y obtener una tabla de datos limpios con la misma información usando *tidyR*.

**Indicación:** lee la ayuda de la función `separate` de *tidyR*.

### Ejercicio 3. Lectura de R4DS.

Continuando con nuestra *lectura conjunta* de este libro, si revisas el índice verás que hemos cubierto (holgadamente en algún caso) el contenido de los Capítulos 6, 8, 9, 10 y 11. Todos esos Capítulos son relativamente ligeros. Por eso esta semana conviene detenerse un poco en la lectura de los Capítulos 7 y 12, que son los más densos en información. Y como motivación os proponemos un par de ejercicios, uno por cada uno de esos capítulos.

- Haz el ejercicio 2 de la Sección 7.5.1.1 de R4DS. Las ideas de esa sección son importantes para nuestro trabajo de las próximas sesiones.
- Haz el ejercicio 4 de la Sección 12.6.1 de R4DS. ¡Aprovecha el código previo de esa sección para trabajar con datos limpios!