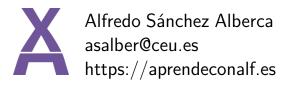
# Pácticas de Aprendizaje Automático con Julia





## Tabla de contenidos

Ρı	efacio	)	3	
	Lice	ncia	3	
1	Intro	oducción	5	
	1.1	El REPL de Julia	6	
	1.2	El gestor de paquetes de Julia	6	
	1.3	Operadores aritméticos	7	
	1.4	Operadores de comparación	7	
	1.5	Operadores booleanos	8	
	1.6	Funciones de redondeo	8	
	1.7	Funciones de división	9	
	1.8	Funciones para el signo y el valor absoluto	.0	
	1.9	Raíces, exponenciales y logaritmos	.1	
	1.10	Funciones trigonométricas	.2	
	1.11	Funciones trigonométricas inversas	.3	
	1.12	Precedencia de operadores	.4	
	1.13	Definición de variables	.5	
2	Preprocesamiento de datos			
	2.1	Ejercicios Resueltos	6	
	2.2	Ejercicios propuestos	0	
3	Regi	resión 3	3	
	3.1	Fiercicios Resueltos	3	

## **Prefacio**

¡Bienvenido a Prácticas de Aprendizaje Automático con Julia!

Este libro presenta una recopilación de prácticas de Aprendizaje Automático (Machine Learning) con el lenguaje de programación Julia.

No es un libro para aprender a programar con Julia, ya que solo enseña el uso del lenguaje y de algunos de sus paquetes para implementar los algoritmos más comunes de Aprendizaje Automático. Para quienes estén interesados en aprender a programar en este Julia, os recomiendo leer este manual de Julia.

#### Licencia

Esta obra está bajo una licencia Reconocimiento – No comercial – Compartir bajo la misma licencia 3.0 España de Creative Commons. Para ver una copia de esta licencia, visite <a href="https://creativecommons.org/licenses/by-nc-sa/3.0/es/">https://creativecommons.org/licenses/by-nc-sa/3.0/es/</a>.

Con esta licencia eres libre de:

- Copiar, distribuir y mostrar este trabajo.
- Realizar modificaciones de este trabajo.

Bajo las siguientes condiciones:

- Reconocimiento. Debe reconocer los créditos de la obra de la manera especificada por el autor o el licenciador (pero no de una manera que sugiera que tiene su apoyo o apoyan el uso que hace de su obra).
- No comercial. No puede utilizar esta obra para fines comerciales.
- Compartir bajo la misma licencia. Si altera o transforma esta obra, o genera una obra derivada, sólo puede distribuir la obra generada bajo una licencia idéntica a ésta.

Al reutilizar o distribuir la obra, tiene que dejar bien claro los términos de la licencia de esta obra.

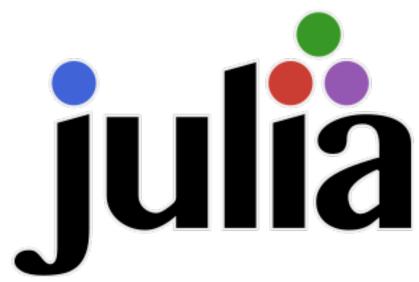
Estas condiciones pueden no aplicarse si se obtiene el permiso del titular de los derechos de autor.

Nada en esta licencia menoscaba o restringe los derechos morales del autor.

## 1 Introducción

La gran potencia de cálculo alcanzada por los ordenadores en las últimas décadas ha convertido a los mismos en poderosas herramientas al servicio de todas aquellas disciplinas que, como las matemáticas, requieren cálculos largos y complejos.

Julia es un lenguaje de programación especialmente orientado al cálculo numérico y el análisis de datos. Julia permite además realizar cálculos simbólicos y dispone de una gran biblioteca de paquetes con aplicaciones en muy diversas áreas de las Matemáticas como Cálculo, Álgebra, Geometría, Matemática Discreta o Estadística.



La ventaja de Julia frente a otros programas habituales de cálculo como Mathematica, MATLAB o Sage radica en su potencia de cálculo y su velocidad (equiparable al lenguaje C), lo que lo hace ideal para manejar grandes volúmenes de datos o realizar tareas que requieran largos y complejos cálculos. Además, es software libre por lo que resulta ideal para introducirlo en el aula como soporte computacional para los modelos matemáticos sin coste alguno.

En el siguiente enlace se explica el procedimiento de instalación de Julia.

Existen también varios entornos de desarrollo online que permiten ejecutar código en Julia sin necesidad de instalarlo en nuestro ordenador, como por ejemplo Replit, Cocalc o Codeanywhere.

El objetivo de esta práctica es introducir al alumno en la utilización de este lenguaje, enseñándole a realizar las operaciones básicas más habituales en Cálculo.

#### 1.1 El REPL de Julia

Para arrancar el REPL^(REPL es el acrónimo de Read, Evaluate, Print and Loop, que describe el funcionamiento del compilador de Julia) de julia basta con abrir una terminal y teclear julia.

#### 1.2 El gestor de paquetes de Julia

Julia viene con varios paquetes básicos preinstalados, como por ejemplo el paquete LinearAlgebra que define funciones básicas del Álgebra Lineal, pero en estas prácticas utilizaremos otros muchos paquetes que añaden más funcionalidades que no vienen instalados por defecto y tendremos que instalarlos aparte. Julia tiene un potente gestor de paquetes que facilita la búsqueda, instalación, actualización y eliminación de paquetes.

Por defecto el gestor de paquetes utiliza el repositorio de paquetes oficial pero se pueden instalar paquetes de otros repositorios.

Para entrar en el modo de gestión de paquetes hay que teclear ]. Esto produce un cambio en el *prompt* del REPL de Julia.

Los comandos más habituales son:

- add p: Instala el paquete p en el entorno activo de Julia.
- update: Actualiza los paquetes del entorno activo de Julia.
- status: Muestra los paquetes instalados y sus versiones en el entorno activo de Julia.

• remove p: Elimina el paquete p del entorno activo de Julia.

```
i Ejemplo

Para instalar el paquete SymPy para cálculo simbólico basta con teclear add Sympy.

(@v1.7) pkg> add SymPy

Updating registry at `~/.julia/registries/General.toml`

Resolving package versions...

Updating `~/.julia/environments/v1.7/Project.toml`

[24249f21] + SymPy v1.1.6

Updating `~/.julia/environments/v1.7/Manifest.toml`

[3709ef60] + CommonEq v0.2.0

[38540f10] + CommonSolve v0.2.1

[438e738f] + PyCall v1.93.1

[24249f21] + SymPy v1.1.6
```

#### 1.3 Operadores aritméticos.

El uso más simple de Julia es la realización de operaciones aritméticas como en una calculadora. En Julia se utilizan los siguientes operadores.

Operador	Descripción
x + y	Suma
х - у	Resta
x * y	Producto
х / у	División
x ÷ y	Cociente división entera
х % у	Resto división entera
x ^ y	Potencia

## 1.4 Operadores de comparación

Operador	Descripción
==	Igualdad
!=,	Desigualdad
<	Menor que
<=,	Menor o igual que

Operador	Descripción
>	Mayor que
>=,	Mayor o igual que

## 1.5 Operadores booleanos

Operador	Descripción
!x	Negación
x && y	Conjunción (y)
x    y	Disyunción (o)

Existen también un montón de funciones predefinidas habituales en Cálculo.

## 1.6 Funciones de redondeo

Función	Descripción
round(x)	Devuelve el entero más próximo a x
<pre>round(x, digits = n)</pre>	Devuelve al valor más próximo a ${\tt x}$ con ${\tt n}$
	decimales
floor(x)	Redondea x al próximo entero menor
ceil(x)	Redondea x al próximo entero mayor
trunc(x)	Devuelve la parte entera de ${\tt x}$

```
i Ejemplo
julia> round(2.7)
3.0
julia> floor(2.7)
2.0
julia> floor(-2.7)
-3.0
julia> ceil(2.7)
3.0
julia> ceil(-2.7)
-2.0
julia> trunc(2.7)
2.0
julia> trunc(-2.7)
-2.0
julia> round(2.5)
2.0
julia> round(2.786, digits = 2)
```

#### 1.7 Funciones de división

Función	Descripción
div(x,y), x÷y	Cociente de la división entera
fld(x,y)	Cociente de la división entera redondeado hacia abajo
<pre>cld(x,y)</pre>	Cociente de la división entera redondeado hacia arriba
rem(x,y), x%y	Resto de la división entera. Se cumple $x == div(x,y)*y +$
	rem(x,y)
mod(x,y)	Módulo con respecto a y. Se cumple $x == fld(x,y)*y + mod(x,y)$
gcd(x,y)	Máximo común divisor positivo de x, y,
lcm(x,y)	Mínimo común múltiplo positivo de x, y,

```
i Ejemplo

julia> div(5,3)
1

julia> cld(5,3)
2

julia> 5%3
2

julia> -5%3
-2

julia> mod(5,3)
2

julia> mod(-5,3)
1

julia> gcd(12,18)
6

julia> lcm(12,18)
36
```

## 1.8 Funciones para el signo y el valor absoluto

Función	Descripción
abs(x)	Valor absoluto de x
sign(x)	Devuelve -1 si $\mathbf{x}$ es positivo, -1 si es negativo y 0 si es 0.

```
i Ejemplo

julia> abs(2.5)
2.5

julia> abs(-2.5)
2.5

julia> sign(-2.5)
-1.0

julia> sign(0)
0

julia> sign(2.5)
1.0
```

## 1.9 Raíces, exponenciales y logaritmos

Función	Descripción
$sqrt(x), \sqrt{x}$	Raíz cuadrada de x
cbrt(x), x	Raíz cúbica de $\mathbf{x}$
exp(x)	Exponencial de $\mathbf{x}$
log(x)	Logaritmo neperiano de $\mathbf{x}$
log(b,x)	Logaritmo en base ${\tt b}$ de ${\tt x}$
log2(x)	Logaritmo en base 2 de $\mathbf{x}$
log10(x)	Logaritmo en base 10 de x

```
i Ejemplo
julia> sqrt(4)
2.0
julia> cbrt(27)
3.0
julia> exp(1)
2.718281828459045
julia> exp(-Inf)
0.0
julia> log(1)
0.0
julia> log(0)
-Inf
julia > log(-1)
ERROR: DomainError with -1.0:
log will only return a complex result if called with a complex argument.
julia > log(-1+0im)
0.0 + 3.141592653589793im
julia > log2(2^3)
3.0
```

## 1.10 Funciones trigonométricas

Función	Descripción
hypot(x,y)	Hipotenusa del triángulo rectángulo con catetos ${\tt x}$ e ${\tt y}$
sin(x)	Seno del ángulo $x$ en radianes
sind(x)	Seno del ángulo $x$ en grados
cos(x)	Coseno del ángulo ${\tt x}$ en radianes
cosd(x)	Coseno del ángulo x en grados
tan(x)	Tangente del ángulo ${\bf x}$ en radianes

Función	Descripción
tand(x)	Tangente del ángulo x en grados
sec(x)	Secante del ángulo $x$ en radianes
csc(x)	Cosecante del ángulo ${\bf x}$ en radianes
cot(x)	Cotangente del ángulo ${\bf x}$ en radianes

```
i Ejemplo
julia> sin(/2)
1.0
julia> cos(/2)
6.123233995736766e-17
julia> cosd(90)
0.0
julia> tan(/4)
0.99999999999999
julia> tand(45)
1.0
julia> tan(/2)
1.633123935319537e16
julia> tand(90)
Inf
julia> \sin(/4)^2 + \cos(/4)^2
1.0
```

## 1.11 Funciones trigonométricas inversas

Función	Descripción
asin(x)	Arcoseno (inversa del seno) de x en radianes
asind(x)	Arcoseno (inversa del seno) de x en grados
acos(x)	Arcocoseno (inversa del coseno) de x en radianes
acosd(x)	Arcocoseno (inversa del coseno) de ${\tt x}$ en grados

Función	Descripción
atan(x)	Arcotangente (inversa de la tangente) de x en radianes
atand(x)	Arcotangente (inversa de la tangente) de ${\bf x}$ en grados
asec(x)	Arcosecante (inversa de la secante) de x en radianes
acsc(x)	Arcocosecante (inversa de la cosecante) de x en radianes
acot(x)	Arcocotangente (inversa de la cotangente) de ${\tt x}$ en radianes

```
i Ejemplo

julia> asin(1)
1.5707963267948966

julia> asind(1)
90.0

julia> acos(-1)
3.141592653589793

julia> atan(1)
0.7853981633974483

julia> atand(tan(/4))
45.0
```

## 1.12 Precedencia de operadores

A la hora de evaluar una expresión aritmética, Julia evalúa los operadores según el siguiente orden de prioridad (de mayor a menor prioridad).

Categoría Operadores	Asociatividad
Funciones exp, log, sin, etc.	
Exponenciación	Derecha
Unarios + - √	Derecha
Fracciones //	Izquierda
Multiplicación/ % & \ ÷	Izquierda
Adición + -	Izquierda
Comparaciones >= <= == != !==	
Asignaciones += -= *= /= //= ^= ÷= %=  = &=	Derecha

Cuando se quiera evaluar un operador con menor prioridad antes que otro con mayor prioridad, hay que utilizar paréntesis.

```
i Ejemplo

julia> 1 + 4 ^ 2 / 2 - 3
6.0

julia> (1 + 4 ^ 2) / 2 - 3
5.5

julia> (1 + 4) ^ 2 / 2 - 3
9.5

julia> 1 + 4 ^ 2 / (2 - 3)
-15.0

julia> (1 + 4 ^ 2) / (2 - 3)
-17.0
```

#### 1.13 Definición de variables

Para definir variables se pueden utilizar cualquier carácter Unicode. Los nombres de las variables pueden contener más de una letra y, en tal caso, pueden usarse también números, pero siempre debe comenzar por una letra. Así, para Julia, la expresión xy, no se interpreta como el producto de la variable x por la variable y, sino como la variable xy. Además, se distingue entre mayúsculas y minúsculas, así que no es lo mismo xy que xy.

## 2 Preprocesamiento de datos

Esta práctica contiene ejercicios que muestran como preprocesar un conjunto de datos con Julia. El preprocesamiento de datos es una tarea fundamental en la construcción de modelos de aprendizaje automático que consiste en la limpieza, transformación y preparación de los datos para que puedan alimentar el proceso de entrenamiento de los modelos, así como para la evaluación de su rendimiento. El preprocesamiento de datos incluye tareas como

- Limpieza de datos.
- Imputación de valores perdidos.
- Recodificación de variables.
- Creación de nuevas variables.
- Transformación de variables.
- Selección de variables.
- Fusión de datos.
- Reestructuración del conjunto de datos.
- División del conjunto de datos en subconjuntos de entrenamiento y prueba.

#### 2.1 Ejercicios Resueltos

Para la realización de esta práctica se requieren los siguientes paquetes:

```
using CSV # Para la lectura de archivos CSV.
using DataFrames # Para el manejo de datos tabulares.
using PrettyTables # Para mostrar tablas formateadas.
using Plots # Para el dibujo de gráficas.
using Makie # Para obtener gráficos interactivos.
```

Ejercicio 2.1. La siguiente tabla contiene los ingresos y gastos de una empresa durante el primer trimestre del año.

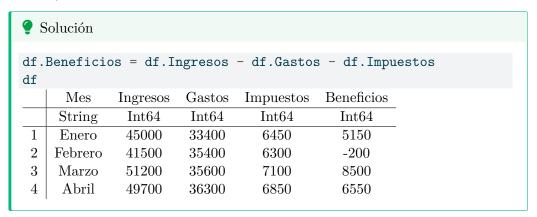
a. Crear un data frame con los datos de la tabla.

#### i Ayuda

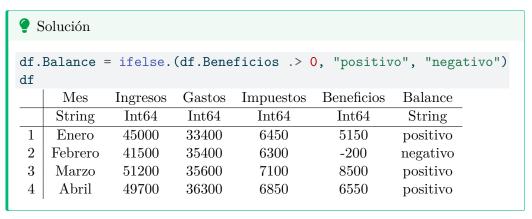
Utilizar la función <code>DataFrame</code> del paquete <code>DataFrames</code> para partir el rango de valores en intervalos y asociar a cada intervalo una categoría.

```
🅊 Solución
using DataFrames
df = DataFrame(
    Mes = ["Enero", "Febrero", "Marzo", "Abril"],
    Ingresos = [45000, 41500, 51200, 49700],
    Gastos = [33400, 35400, 35600, 36300],
    Impuestos = [6450, 6300, 7100, 6850]
    )
      Mes
             Ingresos
                       Gastos
                               Impuestos
              Int64
                       Int64
                                  Int64
     String
              45000
                       33400
                                  6450
 1
     Enero
 2
    Febrero
              41500
                       35400
                                  6300
 3
    Marzo
              51200
                       35600
                                  7100
 4
              49700
                       36300
                                  6850
     Abril
```

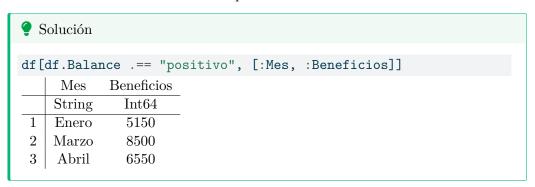
b. Crear una nueva columna con los beneficios de cada mes (ingresos - gastos - impuestos).



c. Crear una nueva columna con el factor Balance con dos posibles categorías: positivo si ha habido beneficios y negativo si ha habido pérdidas.



d. Filtrar el conjunto de datos para quedarse con los nombres de los meses y los beneficios de los meses con balance positivo.



Ejercicio 2.2. El fichero colesterol.csv contiene información de una muestra de pacientes donde se han medido la edad, el sexo, el peso, la altura y el nivel de colesterol, además de su nombre.

a. Crear un data frame con los datos de todos los pacientes del estudio a partir del fichero colesterol.csv.

#### i Ayuda

Utilizar la función CSV.read del paquete CSV para partir el rango de valores en intervalos y asociar a cada intervalo una categoría.



	nombre	edad	sexo	peso	altura	colesterol
	String	Int64	String1	Float64?	Float64	Float64?
1	José Luis Martínez Izquierdo	18	Н	85.0	1.79	182.0
2	Rosa Díaz Díaz	32	${f M}$	65.0	1.73	232.0
3	Javier García Sánchez	24	${ m H}$	missing	1.81	191.0
4	Carmen López Pinzón	35	${f M}$	65.0	1.7	200.0
5	Marisa López Collado	46	${ m M}$	51.0	1.58	148.0
6	Antonio Ruiz Cruz	68	${ m H}$	66.0	1.74	249.0
7	Antonio Fernández Ocaña	51	${ m H}$	62.0	1.72	276.0
8	Pilar Martín González	22	${f M}$	60.0	1.66	missing
9	Pedro Gálvez Tenorio	35	${ m H}$	90.0	1.94	241.0
10	Santiago Reillo Manzano	46	${ m H}$	75.0	1.85	280.0
11	Macarena Álvarez Luna	53	$\mathbf{M}$	55.0	1.62	262.0
12	José María de la Guía Sanz	58	${ m H}$	78.0	1.87	198.0
13	Miguel Angel Cuadrado Gutiérrez	27	${ m H}$	109.0	1.98	210.0
14	Carolina Rubio Moreno	20	${\bf M}$	61.0	1.77	194.0
	•					

b. Crear una nueva columna con el índice de masa corporal, usando la siguiente fórmula

$$\mathrm{IMC} = \frac{\mathrm{Peso}\ (\mathrm{kg})}{\mathrm{Altura}\ (\mathrm{cm})^2}$$

```
    Solución

df.imc = df.peso ./ (df.altura .^ 2)
df
```

	nombre	edad	sexo	peso	altura	colesterol	
	String	Int64	String1	Float64?	Float64	Float64?	
1	José Luis Martínez Izquierdo	18	Н	85.0	1.79	182.0	
2	Rosa Díaz Díaz	32	$\mathbf{M}$	65.0	1.73	232.0	•••
3	Javier García Sánchez	24	Η	missing	1.81	191.0	
4	Carmen López Pinzón	35	$\mathbf{M}$	65.0	1.7	200.0	
5	Marisa López Collado	46	$\mathbf{M}$	51.0	1.58	148.0	
6	Antonio Ruiz Cruz	68	Η	66.0	1.74	249.0	
7	Antonio Fernández Ocaña	51	Н	62.0	1.72	276.0	
8	Pilar Martín González	22	$\mathbf{M}$	60.0	1.66	missing	
9	Pedro Gálvez Tenorio	35	Η	90.0	1.94	241.0	
10	Santiago Reillo Manzano	46	Η	75.0	1.85	280.0	
11	Macarena Álvarez Luna	53	M	55.0	1.62	262.0	
12	José María de la Guía Sanz	58	Н	78.0	1.87	198.0	
13	Miguel Angel Cuadrado Gutiérrez	27	Н	109.0	1.98	210.0	
14	Carolina Rubio Moreno	20	${\bf M}$	61.0	1.77	194.0	

c. Crear una nueva columna con la variable obesidad recodificando la columna imc en las siguientes categorías.

Rango IMC	Categoría
Menor de 18.5	Bajo peso
De 18.5 a 24.5	Saludable
$\mathrm{De}\ 24.5\ \mathrm{a}\ 30$	Sobrepeso
Mayor de 30	Obeso

#### i Ayuda

Utilizar la función cut del paquete CategoricalArrays para partir el rango de valores en intervalos y asociar a cada intervalo una categoría.

		$\operatorname{nombre}$	edad	sexo	peso	altura	colesterol	
		String	Int64	String1	Float64?	Float64	Float64?	
-	1	José Luis Martínez Izquierdo	18	Н	85.0	1.79	182.0	
	2	Rosa Díaz Díaz	32	$\mathbf{M}$	65.0	1.73	232.0	
	3	Javier García Sánchez	24	Η	missing	1.81	191.0	
	4	Carmen López Pinzón	35	M	65.0	1.7	200.0	
	5	Marisa López Collado	46	$\mathbf{M}$	51.0	1.58	148.0	
	6	Antonio Ruiz Cruz	68	Η	66.0	1.74	249.0	
	7	Antonio Fernández Ocaña	51	Η	62.0	1.72	276.0	
	8	Pilar Martín González	22	$\mathbf{M}$	60.0	1.66	missing	
	9	Pedro Gálvez Tenorio	35	Η	90.0	1.94	241.0	
	10	Santiago Reillo Manzano	46	Η	75.0	1.85	280.0	
	11	Macarena Álvarez Luna	53	M	55.0	1.62	262.0	
	12	José María de la Guía Sanz	58	Η	78.0	1.87	198.0	
	13	Miguel Angel Cuadrado Gutiérrez	27	Н	109.0	1.98	210.0	
	14	Carolina Rubio Moreno	20	M	61.0	1.77	194.0	
		ı						

d. Seleccionar las columnas nombre, sexo y edad.

• So	lución						
df[:	df[:, [:nombre, :sexo, :edad]]						
	nombre	sexo	edad				
	String	String1	Int64				
1	José Luis Martínez Izquierdo	Н	18				
2	Rosa Díaz Díaz	${ m M}$	32				
3	Javier García Sánchez	${ m H}$	24				
4	Carmen López Pinzón	${f M}$	35				
5	Marisa López Collado	${ m M}$	46				
6	Antonio Ruiz Cruz	${ m H}$	68				
7	Antonio Fernández Ocaña	${ m H}$	51				
8	Pilar Martín González	${ m M}$	22				
9	Pedro Gálvez Tenorio	${ m H}$	35				
10	Santiago Reillo Manzano	${ m H}$	46				
11	Macarena Álvarez Luna	${ m M}$	53				
12	José María de la Guía Sanz	${ m H}$	58				
13	Miguel Angel Cuadrado Gutiérrez	${ m H}$	27				
14	Carolina Rubio Moreno	${\bf M}$	20				

e. Anonimizar los datos eliminando la columna nombre.

## i Ayuda

Utilizar la función select del paquete DataFrames para seleccionar las columnas deseadas y eliminar las columnas no deseadas. Existe también la función select! que modifica el data frame original eliminando las columnas no seleccionadas.

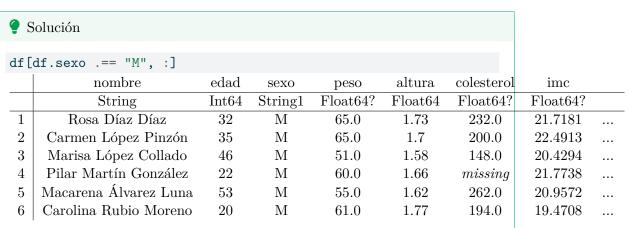
• Sc	lución						
sele	ct(df,	Not(:nor	nbre))				
	edad	sexo	peso	altura	colesterol	imc	obesidad
	Int64	String1	Float64?	Float64	Float64?	Float64?	Cat?
1	18	Н	85.0	1.79	182.0	26.5285	Sobrepeso
2	32	${ m M}$	65.0	1.73	232.0	21.7181	Saludable
3	24	Η	missing	1.81	191.0	missing	missing
4	35	${ m M}$	65.0	1.7	200.0	22.4913	Saludable
5	46	${ m M}$	51.0	1.58	148.0	20.4294	Saludable
6	68	H	66.0	1.74	249.0	21.7994	Saludable
7	51	H	62.0	1.72	276.0	20.9573	Saludable
8	22	${ m M}$	60.0	1.66	missing	21.7738	Saludable
9	35	H	90.0	1.94	241.0	23.9133	Saludable
10	46	Η	75.0	1.85	280.0	21.9138	Saludable
11	53	${ m M}$	55.0	1.62	262.0	20.9572	Saludable
12	58	Η	78.0	1.87	198.0	22.3055	Saludable
13	27	Η	109.0	1.98	210.0	27.8033	Sobrepeso
14	20	M	61.0	1.77	194.0	19.4708	Saludable

f. Reordenar las columnas poniendo la columna sexo antes que la columna edad.

```
Solución
select(df, Cols(:sexo, :edad, Not(:sexo, :edad)))
```

1		1 1	1		1,	1 , 1	
	sexo	edad	nombre	peso	altura	colesterol	
	String1	Int64	String	Float64?	Float64	Float64?	
1	Н	18	José Luis Martínez Izquierdo	85.0	1.79	182.0	
2	M	32	Rosa Díaz Díaz	65.0	1.73	232.0	
3	Η	24	Javier García Sánchez	missing	1.81	191.0	
4	${ m M}$	35	Carmen López Pinzón	65.0	1.7	200.0	
5	${ m M}$	46	Marisa López Collado	51.0	1.58	148.0	
6	Η	68	Antonio Ruiz Cruz	66.0	1.74	249.0	
7	Η	51	Antonio Fernández Ocaña	62.0	1.72	276.0	
8	${ m M}$	22	Pilar Martín González	60.0	1.66	missing	
9	Η	35	Pedro Gálvez Tenorio	90.0	1.94	241.0	
10	Η	46	Santiago Reillo Manzano	75.0	1.85	280.0	
11	M	53	Macarena Álvarez Luna	55.0	1.62	262.0	
12	${ m H}$	58	José María de la Guía Sanz	78.0	1.87	198.0	
13	${ m H}$	27	Miguel Angel Cuadrado Gutiérrez	109.0	1.98	210.0	
14	M	20	Carolina Rubio Moreno	61.0	1.77	194.0	

g. Filtrar el data frame para quedarse con las mujeres.



h. Filtrar el data frame para quedarse con los hombres mayores de 30 años.

```
    Solución

df[(df.sexo .== "H") .& (df.edad .> 30), :]
```

	nombre	$\operatorname{edad}$	sexo	peso	altura	colesterol	imc	
	String	Int64	String1	Float64?	Float64	Float 64?	Float64?	
1	Antonio Ruiz Cruz	68	Н	66.0	1.74	249.0	21.7994	
2	Antonio Fernández Ocaña	51	Η	62.0	1.72	276.0	20.9573	
3	Pedro Gálvez Tenorio	35	Η	90.0	1.94	241.0	23.9133	
4	Santiago Reillo Manzano	46	Η	75.0	1.85	280.0	21.9138	
5	José María de la Guía Sanz	58	Н	78.0	1.87	198.0	22.3055	
'	!							

i. Filtrar el data frame para quedarse con las filas sin valores perdidos.

#### i Ayuda

Utilizar la función dropmissing del paquete DataFrames para eliminar las filas con valores perdidos.

Solución     Solución									
dropmissing(df)									
nombre	$\operatorname{edad}$	sexo	peso	altura	colesterol	imc			
String	Int64	String1	Float64	Float64	Float64	Floate			
José Luis Martínez Izquierdo	18	Н	85.0	1.79	182.0	26.528			
Rosa Díaz Díaz	32	${ m M}$	65.0	1.73	232.0	21.718			
Carmen López Pinzón	35	${ m M}$	65.0	1.7	200.0	22.491			
Marisa López Collado	46	${ m M}$	51.0	1.58	148.0	20.429			
Antonio Ruiz Cruz	68	Η	66.0	1.74	249.0	21.799			
Antonio Fernández Ocaña	51	Н	62.0	1.72	276.0	20.957			
Pedro Gálvez Tenorio	35	Н	90.0	1.94	241.0	23.913			
Santiago Reillo Manzano	46	Н	75.0	1.85	280.0	21.913			
Macarena Álvarez Luna	53	${ m M}$	55.0	1.62	262.0	20.957			
José María de la Guía Sanz	58	${ m H}$	78.0	1.87	198.0	22.305			
Miguel Angel Cuadrado Gutiérrez	27	${\rm H}$	109.0	1.98	210.0	27.803			
Carolina Rubio Moreno	20	${\bf M}$	61.0	1.77	194.0	19.470			
	nombre String José Luis Martínez Izquierdo Rosa Díaz Díaz Carmen López Pinzón Marisa López Collado Antonio Ruiz Cruz Antonio Fernández Ocaña Pedro Gálvez Tenorio Santiago Reillo Manzano Macarena Álvarez Luna José María de la Guía Sanz Miguel Angel Cuadrado Gutiérrez	missing(df)  nombre edad  String Int64  José Luis Martínez Izquierdo 18  Rosa Díaz Díaz 32  Carmen López Pinzón 35  Marisa López Collado 46  Antonio Ruiz Cruz 68  Antonio Fernández Ocaña 51  Pedro Gálvez Tenorio 35  Santiago Reillo Manzano 46  Macarena Álvarez Luna 53  José María de la Guía Sanz 58  Miguel Angel Cuadrado Gutiérrez 27	missing(df)  String Int64 String1  José Luis Martínez Izquierdo 18 H Rosa Díaz Díaz 32 M Carmen López Pinzón 35 M Marisa López Collado 46 M Antonio Ruiz Cruz 68 H Antonio Fernández Ocaña 51 H Pedro Gálvez Tenorio 35 H Santiago Reillo Manzano 46 H Macarena Álvarez Luna 53 M José María de la Guía Sanz 58 H Miguel Angel Cuadrado Gutiérrez 27 H	missing(df)         edad         sexo         peso           String         Int64         String1         Float64           José Luis Martínez Izquierdo         18         H         85.0           Rosa Díaz Díaz         32         M         65.0           Carmen López Pinzón         35         M         65.0           Marisa López Collado         46         M         51.0           Antonio Ruiz Cruz         68         H         66.0           Antonio Fernández Ocaña         51         H         62.0           Pedro Gálvez Tenorio         35         H         90.0           Santiago Reillo Manzano         46         H         75.0           Macarena Álvarez Luna         53         M         55.0           José María de la Guía Sanz         58         H         78.0           Miguel Angel Cuadrado Gutiérrez         27         H         109.0	missing(df)         edad         sexo         peso         altura           String         Int64         String1         Float64         Float64           José Luis Martínez Izquierdo         18         H         85.0         1.79           Rosa Díaz Díaz         32         M         65.0         1.73           Carmen López Pinzón         35         M         65.0         1.7           Marisa López Collado         46         M         51.0         1.58           Antonio Ruiz Cruz         68         H         66.0         1.74           Antonio Fernández Ocaña         51         H         62.0         1.72           Pedro Gálvez Tenorio         35         H         90.0         1.94           Santiago Reillo Manzano         46         H         75.0         1.85           Macarena Álvarez Luna         53         M         55.0         1.62           José María de la Guía Sanz         58         H         78.0         1.87           Miguel Angel Cuadrado Gutiérrez         27         H         109.0         1.98	missing(df)         edad         sexo         peso         altura         colesterol           String         Int64         String1         Float64         Float64         Float64           José Luis Martínez Izquierdo         18         H         85.0         1.79         182.0           Rosa Díaz Díaz         32         M         65.0         1.73         232.0           Carmen López Pinzón         35         M         65.0         1.7         200.0           Marisa López Collado         46         M         51.0         1.58         148.0           Antonio Ruiz Cruz         68         H         66.0         1.74         249.0           Antonio Fernández Ocaña         51         H         62.0         1.72         276.0           Pedro Gálvez Tenorio         35         H         90.0         1.94         241.0           Santiago Reillo Manzano         46         H         75.0         1.85         280.0           Macarena Álvarez Luna         53         M         55.0         1.62         262.0           José María de la Guía Sanz         58         H         78.0         1.87         198.0           Miguel Angel Cuadrado Gutiérrez         27			

j. Filtrar el data frame para eliminar las filas con datos perdidos en la columna colesterol.

#### i Ayuda

Utilizar la función ismissing en la condición del filtro.

#### Solución

df[.!ismissing.(df.colesterol), :]

	$\operatorname{nombre}$	edad	sexo	peso	altura	colesterol	
	$\operatorname{String}$	Int64	String1	Float64?	Float64	Float64?	
1	José Luis Martínez Izquierdo	18	Н	85.0	1.79	182.0	
2	Rosa Díaz Díaz	32	${ m M}$	65.0	1.73	232.0	
3	Javier García Sánchez	24	Η	missing	1.81	191.0	
$_4$	Carmen López Pinzón	35	${f M}$	65.0	1.7	200.0	
5	Marisa López Collado	46	${ m M}$	51.0	1.58	148.0	
6	Antonio Ruiz Cruz	68	Η	66.0	1.74	249.0	
7	Antonio Fernández Ocaña	51	Η	62.0	1.72	276.0	
8	Pedro Gálvez Tenorio	35	Η	90.0	1.94	241.0	
9	Santiago Reillo Manzano	46	Η	75.0	1.85	280.0	
10	Macarena Álvarez Luna	53	${ m M}$	55.0	1.62	262.0	
11	José María de la Guía Sanz	58	${ m H}$	78.0	1.87	198.0	
12	Miguel Angel Cuadrado Gutiérrez	27	${ m H}$	109.0	1.98	210.0	
13	Carolina Rubio Moreno	20	${ m M}$	61.0	1.77	194.0	
'							

k. Imputar los valores perdidos en la columna colesterol con la media de los valores no perdidos.

#### i Ayuda

Utilizar la función coalesce para reemplazar los valores perdidos por otros valores.

## Solución

```
using Statistics
media_colesterol = mean(skipmissing(df.colesterol))
df.colesterol = coalesce.(df.colesterol, media_colesterol)
df
```

		nombre	edad	sexo	peso	altura	colesterol	
		String	Int64	String1	Float64?	Float64	Float64	
	1	José Luis Martínez Izquierdo	18	Н	85.0	1.79	182.0	
:	2	Rosa Díaz Díaz	32	${\bf M}$	65.0	1.73	232.0	
;	3	Javier García Sánchez	24	Η	missing	1.81	191.0	
	4	Carmen López Pinzón	35	${ m M}$	65.0	1.7	200.0	
	5	Marisa López Collado	46	${ m M}$	51.0	1.58	148.0	
(	6	Antonio Ruiz Cruz	68	Η	66.0	1.74	249.0	
,	7	Antonio Fernández Ocaña	51	Η	62.0	1.72	276.0	
	8	Pilar Martín González		${\bf M}$	60.0	1.66	220.231	
!	9	Pedro Gálvez Tenorio		Η	90.0	1.94	241.0	
1	0	Santiago Reillo Manzano	46	Η	75.0	1.85	280.0	
1	1	Macarena Álvarez Luna	53	M	55.0	1.62	262.0	
1:	2	José María de la Guía Sanz	58	Н	78.0	1.87	198.0	
13	3	Miguel Angel Cuadrado Gutiérrez		Н	109.0	1.98	210.0	
1	4	Carolina Rubio Moreno	20	${\bf M}$	61.0	1.77	194.0	
		•						

l. Ordenar el data frame según la columna nombre.

## i Ayuda

Utilizar la función **sort** para ordenar las filas del data frame según los valores de una o varias columnas. Utilizar el parámetro **rev** para especificar mediante un vector de booleanos si el orden es ascendente o descendente.

## Solución

sort(df, :nombre)

		1	edad			1.	1 , 1	
١.		nombre		sexo	peso	altura	colesterol	
		String	Int64	String1	Float64?	Float64	Float64	
	1	Antonio Fernández Ocaña	51	Η	62.0	1.72	276.0	
	2	Antonio Ruiz Cruz	68	H	66.0	1.74	249.0	
	3	Carmen López Pinzón	35	${ m M}$	65.0	1.7	200.0	
	4	Carolina Rubio Moreno	20	$\mathbf{M}$	61.0	1.77	194.0	
	5	Javier García Sánchez	24	H	missing	1.81	191.0	
	6	José Luis Martínez Izquierdo	18	H	85.0	1.79	182.0	
	7	José María de la Guía Sanz	58	H	78.0	1.87	198.0	
	8	Macarena Álvarez Luna	53	${ m M}$	55.0	1.62	262.0	
	9	Marisa López Collado	46	$\mathbf{M}$	51.0	1.58	148.0	
	10	Miguel Angel Cuadrado Gutiérrez	27	Η	109.0	1.98	210.0	
	11	Pedro Gálvez Tenorio	35	Η	90.0	1.94	241.0	
	12	Pilar Martín González	22	$\mathbf{M}$	60.0	1.66	220.231	
	13	Rosa Díaz Díaz	32	${ m M}$	65.0	1.73	232.0	
	14	Santiago Reillo Manzano	46	Н	75.0	1.85	280.0	
	,							

m. Ordenar el data frame ascendentemente por la columna sexo y descendentemente por la columna edad.

Solución								
sort								
	nombre	colesterol						
	String	Int64	String1	Float64?	Float64	Float64		
1	Antonio Ruiz Cruz	68	Н	66.0	1.74	249.0		
2	José María de la Guía Sanz	58	Η	78.0	1.87	198.0		
3	Antonio Fernández Ocaña	51	Η	62.0	1.72	276.0		
4	Santiago Reillo Manzano	46	${ m H}$	75.0	1.85	280.0		
5	Pedro Gálvez Tenorio	35	${ m H}$	90.0	1.94	241.0		
6	Miguel Angel Cuadrado Gutiérrez	27	${ m H}$	109.0	1.98	210.0		
7	Javier García Sánchez	24	${ m H}$	missing	1.81	191.0		
8	José Luis Martínez Izquierdo	18	${ m H}$	85.0	1.79	182.0		
9	Macarena Álvarez Luna	53	${ m M}$	55.0	1.62	262.0		
10	Marisa López Collado	46	${ m M}$	51.0	1.58	148.0		
11	Carmen López Pinzón	35	${ m M}$	65.0	1.7	200.0		
12	Rosa Díaz Díaz	32	${\rm M}$	65.0	1.73	232.0		
13	Pilar Martín González	22	${\rm M}$	60.0	1.66	220.231		
14	Carolina Rubio Moreno	20	${ m M}$	61.0	1.77	194.0		

Ejercicio 2.3. El fichero notas-curso2.csv contiene información de las notas de los alumnos de un curso.

a. Crear un data frame con los datos de los alumnos del curso a partir del fichero notas-curso2.csv.

	g CSV, Da CSV.read		notas-cu	rso2.csv'	', DataFr	ame; miss:	ingstring=	"NA")	
	sexo	turno	grupo	trabaja	notaA	notaB	notaC	notaD	notaE
	String7	String7	String1	String1	Float64	Float64?	Float64?	Float64?	Float64
1	Mujer	Tarde	С	N	5.2	6.3	3.4	2.3	2.0
2	Hombre	Mañana	A	N	5.7	5.7	4.2	3.5	2.7
3	Hombre	Mañana	В	N	8.3	8.8	8.8	8.0	5.5
4	Hombre	Mañana	В	N	6.1	6.8	4.0	3.5	2.2
5	Hombre	Mañana	A	N	6.2	9.0	5.0	4.4	3.7
6	Hombre	Mañana	A	$\mathbf{S}$	8.6	8.9	9.5	8.4	3.9
7	Mujer	Mañana	A	N	6.7	7.9	5.6	4.8	4.2
8	Mujer	Tarde	$\mathbf{C}$	$\mathbf{S}$	4.1	5.2	1.7	0.3	1.0
9	Hombre	Tarde	$\mathbf{C}$	N	5.0	5.0	3.3	2.7	6.0
10	Hombre	Tarde	$\mathbf{C}$	N	5.3	6.3	4.8	3.6	2.3
11	Mujer	Mañana	A	N	7.8	missing	6.5	6.7	2.8
12	Hombre	Mañana	A	N	6.5	8.0	5.0	3.2	3.3
13	Hombre	Mañana	В	N	6.6	7.6	5.3	4.0	1.0
4	Hombre	Mañana	В	N	6.2	6.7	5.3	4.7	4.7
5	Hombre	Mañana	В	N	5.2	4.1	5.8	5.0	1.9
6	Hombre	Mañana	В	$\mathbf{S}$	8.7	missing	7.6	6.3	9.3
7	Mujer	Mañana	В	N	6.7	6.3	6.8	5.3	2.8
18	Mujer	Tarde	$\mathbf{C}$	N	3.1	4.8	2.7	1.8	1.0
19	Hombre	Mañana	В	N	7.1	10.0	5.9	4.9	2.5
20	Hombre	Mañana	В	$\mathbf{S}$	5.3	4.8	2.8	0.9	4.2
21	Hombre	Mañana	A	N	4.4	6.1	2.9	1.9	2.4
22	Mujer	Tarde	$\mathbf{C}$	$\mathbf{S}$	5.7	6.4	4.2	3.3	1.0
23	Mujer	Tarde	$\mathbf{C}$	$\mathbf{S}$	4.5	3.9	3.5	2.6	2.2
24	Mujer	Tarde	$\mathbf{C}$	N	4.9	4.4	5.2	3.7	0.4
25	Hombre	Tarde	$\mathbf{C}$	N	5.1	4.5	6.4	5.6	0.5
26	Hombre	Mañana	В	N	3.5	3.9	4.0	3.7	1.7
27	Mujer	Mañana	A	N	7.3	6.8	6.1	5.4	1.8
28	Hombre	Mañana	В	N	6.2	7.9	3.5	1.9	3.1
29	Mujer	Mañana	A	N	4.3	7.4	2.9	1.7	0.2
30	Hombre	Mañana	В	N	7.9	7.6	7.0	6.1	0.4
							•••		

b. Obtener el número de datos perdidos en cada columna.

#### Solución describe(df)[:, [:variable, :nmissing]] nmissing variable Symbol Int64 1 0 sexo2 turno0 3 0 grupo 4 0 trabaja 5 0 notaA 6 notaB 5 7 notaC1 2 8 notaD 9 2 notaE

c. Recodificar la variable grupo en una colección de columnas binarias.

#### i Ayuda

Utilizar la función onehotbatch del paquete OneHotArrays para recodificar una variable categórica en una colección de columnas binarias.

```
    Solución

using OneHotArrays
codificacion = permutedims(onehotbatch(df.grupo, unique(df.grupo)))
hcat(df, DataFrame(codificacion, :auto))
```

l									
	sexo	turno	grupo	trabaja	notaA	notaB	notaC	notaD	notaE
	String7	String7	String1	String1	Float64	Float64?	Float64?	Float64?	Float64
1	Mujer	Tarde	С	N	5.2	6.3	3.4	2.3	2.0
2	Hombre	Mañana	A	N	5.7	5.7	4.2	3.5	2.7
3	Hombre	Mañana	В	N	8.3	8.8	8.8	8.0	5.5
4	Hombre	Mañana	В	N	6.1	6.8	4.0	3.5	2.2
5	Hombre	Mañana	A	N	6.2	9.0	5.0	4.4	3.7
6	Hombre	Mañana	A	$\mathbf{S}$	8.6	8.9	9.5	8.4	3.9
7	Mujer	Mañana	A	N	6.7	7.9	5.6	4.8	4.2
8	Mujer	Tarde	$\mathbf{C}$	$\mathbf{S}$	4.1	5.2	1.7	0.3	1.0
9	Hombre	Tarde	$\mathbf{C}$	N	5.0	5.0	3.3	2.7	6.0
10	Hombre	Tarde	$\mathbf{C}$	N	5.3	6.3	4.8	3.6	2.3
11	Mujer	Mañana	A	N	7.8	missing	6.5	6.7	2.8
12	Hombre	Mañana	A	N	6.5	8.0	5.0	3.2	3.3
13	Hombre	Mañana	В	N	6.6	7.6	5.3	4.0	1.0
14	Hombre	Mañana	В	N	6.2	6.7	5.3	4.7	4.7
15	Hombre	Mañana	В	N	5.2	4.1	5.8	5.0	1.9
16	Hombre	Mañana	В	$\mathbf{S}$	8.7	missing	7.6	6.3	9.3
17	Mujer	Mañana	В	N	6.7	6.3	6.8	5.3	2.8
18	Mujer	Tarde	$\mathbf{C}$	N	3.1	4.8	2.7	1.8	1.0
19	Hombre	Mañana	В	N	7.1	10.0	5.9	4.9	2.5
20	Hombre	Mañana	В	$\mathbf{S}$	5.3	4.8	2.8	0.9	4.2
21	Hombre	Mañana	A	N	4.4	6.1	2.9	1.9	2.4
22	Mujer	Tarde	$\mathbf{C}$	$\mathbf{S}$	5.7	6.4	4.2	3.3	1.0
23	Mujer	Tarde	$\mathbf{C}$	$\mathbf{S}$	4.5	3.9	3.5	2.6	2.2
24	Mujer	Tarde	$\mathbf{C}$	N	4.9	4.4	5.2	3.7	0.4
25	Hombre	Tarde	$\mathbf{C}$	N	5.1	4.5	6.4	5.6	0.5
26	Hombre	Mañana	В	N	3.5	3.9	4.0	3.7	1.7
27	Mujer	Mañana	A	N	7.3	6.8	6.1	5.4	1.8
28	Hombre	Mañana	В	N	6.2	7.9	3.5	1.9	3.1
29	Mujer	Mañana	A	N	4.3	7.4	2.9	1.7	0.2
30	Hombre	Mañana	В	N	7.9	7.6	7.0	6.1	0.4
						•••	•••	•••	

## 2.2 Ejercicios propuestos

**Ejercicio 2.4.** Calcular el décimo término de la sucesión  $\left(\frac{3n^2+n}{6n^2-1}\right)_{n=1}^{\infty}$ .

<sup>\*</sup>Hint: \*

#### Introducir hasta 5 decimales

Ejercicio 2.5. Los ficheros vinos-blancos.xls y vinos-tintos.csv contienen información sobre las características de vinos blancos y tintos portugueses de la denominación "Vinho Verde". Las variables almacenadas en estos archivos son las siguientes:

		Tipo
Variable	Descripción	(unidades)
tipo	Tipo de vino	Factor
		(blanco, tinto)
meses.barrica	Mesesde envejecimiento en barrica	Numérica(meses)
acided.fija	Cantidadde ácidotartárico	Numérica(g/dm3)
acided.volatil	Cantidad de ácido acético	Numérica(g/dm3)
acido.citrico	Cantidad de ácidocítrico	Numérica(g/dm3)
azucar.residual	Cantidad de azúcarremanente después de	Numérica(g/dm3)
	la fermentación	
cloruro.sodico	Cantidad de clorurosódico	Numérica(g/dm3)
dioxido.azufre.libre	Cantidad de dióxido de azufreen formalibre	Numérica(mg/dm3)
dioxido.azufre.total	Cantidadde dióxido de azufretotal en	Numérica(mg/dm3)
	forma libre o ligada	
densidad	Densidad	Numérica(g/cm3)
ph	pН	Numérica(0-
		14)
sulfatos	Cantidadde sulfato de potasio	$Num{\'e}rica(g/dm3)$
alcohol	Porcentajede contenidode alcohol	Numérica(0-
		100)
calidad	Calificación otorgada porun panel de	Numérica(0-
	expertos	10)

- a. Crear un data frame con los datos de los vinos blancos partir del fichero de Excel vinos-blancos.xlsx.
- b. Crear un data frame con los datos de los vinos tintos partir del fichero csv vinos-tintos.csv.
- c. Fusionar los datos de los vinos blancos y tintos en un nuevo data frame.
- d. Convertir el tipo de vino en un factor.
- e. Imputar los valores perdidos del alcohol con la media de los valores no perdidos para cada tipo de vino.
- f. Crear un factor Envejecimiento recodificando la variable meses.barrica en las siguientes categorías.

Rango en meses	Categoría
Menos de 3	Joven
Entre 3 y 12	Crianza
Entre 12 y 18	Reserva
Más de 18	Gran reserva

g. Crear un factor <code>Dulzor</code> recodificando la variable <code>azucar.residual</code> en las siguientes categorías.

Rango azúcar	Categoría
Menos de 4	Seco
Más de 4 y menos de 12	Semiseco
Más de 12 y menos de 45	Semidulce
Más de 45	Dulce

- h. Filtrar el conjunto de datos para quedarse con los vinos Reserva o Gran Reserva con una calidad superior a 7 y ordenar el data frame por calidad de forma descendente.
- i. ¿Cuántos vinos blancos con un contenido en alcohol superior al 12% y una calidad superior a 8 hay en el conjunto de datos?

## 3 Regresión

Los modelos de aprendizaje basados en regresión son modelos bastante simples que pueden utilizarse para predecir variables cuantitativas (regresión lineal) o cualitativas (regresión logística). Esta práctica contiene ejercicios que muestran como construir modelos de aprendizaje de regresión lineal y regresión logística con Julia.

#### 3.1 Ejercicios Resueltos

Para la realización de esta práctica se requieren los siguientes paquetes:

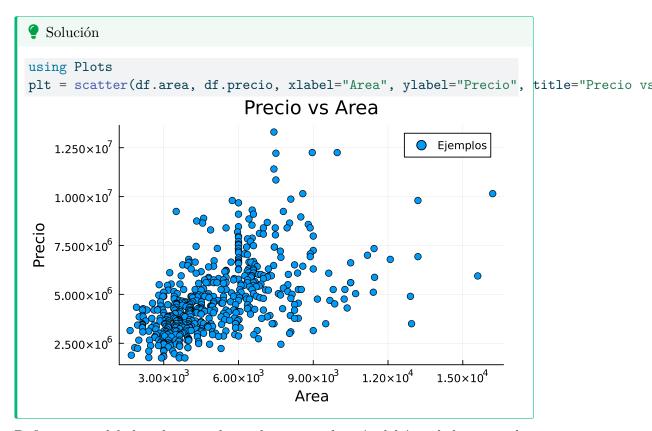
```
using CSV # Para la lectura de archivos CSV.
using DataFrames # Para el manejo de datos tabulares.
using PrettyTables # Para mostrar tablas formateadas.
using Plots # Para el dibujo de gráficas.
using GLMakie # Para obtener gráficos interactivos.
```

Ejercicio 3.1. El conjunto de datos viviendas.csv contiene información sobre el precio de venta de viviendas en una ciudad.

a. Cargar los datos del archivo viviendas.csv en un data frame.

```
🅊 Solución
using CSV, DataFrames
df = CSV.read("datos/viviendas.csv", DataFrame)
first(df, 5)
      precio
                         dormitorios
                                       baños
                                                habitaciones
                                                               calleprincipal
                                                                               huespedes
                                                                                             sotano
       Int64
                 Int64
                            Int64
                                        Int64
                                                    Int64
                                                                   String3
                                                                                 String3
                                                                                             String3
                                          2
                                                      3
     13300000
                 7420
                              4
                                                                     \sin
                                                                                    no
                                                                                               no
                                                      4
 ^{2}
    12250000
                 8960
                              4
                                          4
                                                                     \sin
                                                                                    no
                                                                                               no
                                          2
                                                      2
 3
     12250000
                 9960
                              3
                                                                     \sin
                                                                                    no
                                                                                                \sin
                                          2
                                                      2
 4
     12215000
                 7500
                              4
                                                                     si
                                                                                    no
                                                                                                si
                                                      2
     11410000
                 7420
                              4
                                          1
                                                                     \sin
                                                                                    si
                                                                                                si
```

b. Dibujar un diagrama de dispersión entre el precio y el area de las viviendas.



c. Definir un modelo lineal que explique el precio en función del área de las viviendas.

## i Ayuda

Un modelo lineal tiene encuación  $y = \theta_1 + \theta_2 x$ .

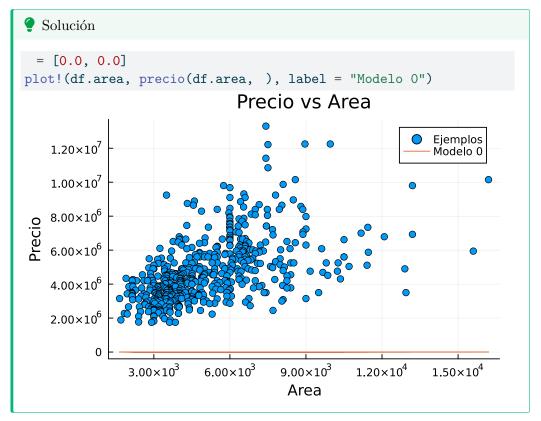
```
    Solución

precio(area, ) = [1] .+ [2] * area

precio (generic function with 1 method)

Observa que la función precio está vectorizada, lo que significa que puede recibir un vector de áreas y devolver un vector de precios.
```

d. Inicializar los parámetros del modelo lineal con valores nulos y dibujar el modelo sobre el diagrama de dispersión.



e. Definir una función de costo para el modelo lineal y evaluar el coste para el modelo lineal construido con los parámetros iniciales. A la vista del coste obtenido, ¿cómo de bueno es el modelo?

#### i Ayuda

La función de coste para un modelo lineal es el error cuadrático medio.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

donde  $h_{\theta}$  es el modelo,  $h_{\theta}(x^{(i)})$  es la predicción del modelo para el ejemplo *i*-ésimo,  $y^{(i)}$  es el valor real observado para el ejemplo *i*-ésimo, y m es el número de ejemplos.

#### Solución

```
function coste(, X, Y)
    m = length(Y)
    return sum((precio(X, ) .- Y).^2) / (2 * m)
end

coste(, df.area, df.precio)
```

#### 1.3106916364659266e13

La función de coste nos da una medida de lo lejos que están las predicciones del modelo de los valores reales observados. En este caso, el coste es muy alto, lo que indica que el modelo no es bueno.

f. ¿En qué dirección debemos modificar los parámetros del modelo para mejorar el modelo?



Para minimizar la función de coste, debemos modificar los parámetros del modelo en la dirección opuesta al gradiente de la función de coste, ya que el gradiente de una función indica la dirección de mayor crecimiento de la función.

g. Crear una función para modificar los pesos del modelo lineal mediante el algoritmo del gradiente descendente, y aplicarla a los parámetros actuales tomando una tasa de aprendizaje de  $10^{-8}$ . ¿Cómo han cambiado los parámetros del modelo? Dibujar el modelo actualizado sobre el diagrama de dispersión. ¿Cómo ha cambiado el coste?

#### i Ayuda

El algoritmo del gradiente descendente actualiza los parámetros del modelo de acuerdo a la siguiente regla:

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

donde  $\alpha$  es la tasa de aprendizaje y  $\frac{\partial J(\theta)}{\partial \theta_j}$  es la derivada parcial de la función de coste con respecto al parámetro  $\theta_j$ .

#### Solución

```
function gradiente_descendente!(, X, Y, )
    # Calculamos el número de ejemplos
    m = length(Y)
    # Actualizamos el término independiente del modelo lineal.
    [1] -= * sum(precio(X, ) - Y) / m
    # Actualizamos la pendiente del modelo lineal.
    [2] -= * sum((precio(X, ) - Y) .* X) / m
    return
end
```

#### gradiente\_descendente! (generic function with 1 method)

Aplicamos la función a los parámetros del modelo actual y mostramos los nuevos parámetros.

```
gradiente_descendente!( , df.area, df.precio, 1e-8)
```

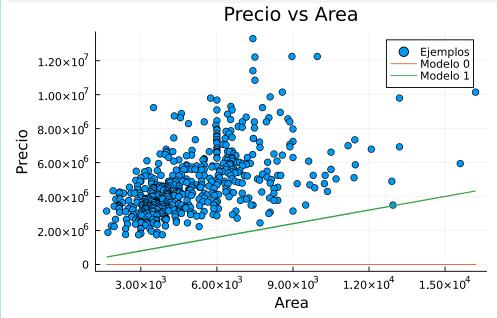
#### 2-element Vector{Float64}:

0.04766729247706422

267.22919804579385

Dibujamos el nuevo modelo.

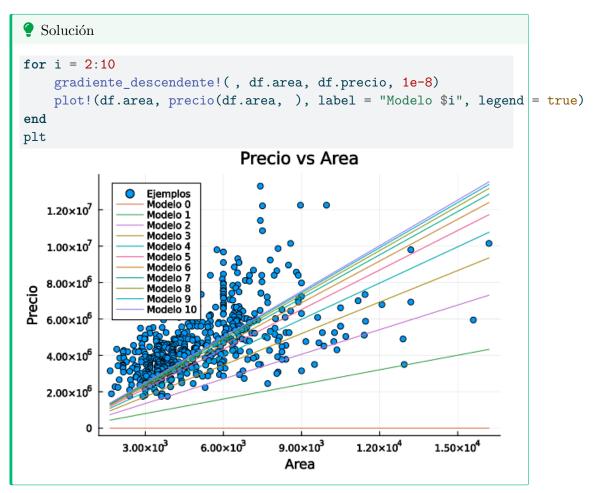




Se observa que ahora la recta está más cerca de la nube de puntos, por lo que el modelo ha mejorado. Calculamos el coste del nuevo modelo.

```
coste(, df.area, df.precio)
7.080823787113201e12
```

h. Repetir el proceso de actualización de los parámetros del modelo mediante el algoritmo del gradiente descendente durante 9 iteraciones más y dibujar los modelos actualizados.



i. Dibujar un gráfico con la evolución del coste del modelo a lo largo de las iteraciones. ¿Cómo se comporta el coste a lo largo de las iteraciones?

```
Solución
costes = Float64[]
for i = 1:10
    gradiente_descendente!(, df.area, df.precio, 1e-8)
    push!(costes, coste(, df.area, df.precio))
end
costes
10-element Vector{Float64}:
4.230808760870044e12
2.882906194020343e12
2.2454213686913755e12
 1.9439256128790886e12
 1.8013344680594421e12
 1.7338965877160208e12
 1.7020021263374993e12
 1.6869177748236997e12
 1.6797836937723748e12
 1.6764096595632322e12
El coste del modelo disminuye en cada iteración, lo que indica que el modelo
está mejorando. Esto se debe a que el algoritmo del gradiente descendente
modifica los parámetros del modelo en la dirección que minimiza la función
```

j. ¿Hasta qué iteración habrá que llegar para conseguir un reducción del coste menor de un 0.0001%?

de coste.

```
② Solución

= [0.0, 0.0]
costes = [0, coste(, df.area, df.precio)]
i = 1
while abs(costes[end] - costes[end-1]) / costes[end-1] > 0.000001
i += 1
gradiente_descendente!(, df.area, df.precio, 1e-8)
push!(costes, coste(, df.area, df.precio))
end
i

23
En este caso, el algoritmo del gradiente descendente converge en 1000 iteraciones.
```

k. ¿Qué sucede si se utiliza una tasa de aprendizaje  $\alpha=0.0001?$  ¿Cómo afecta al coste y a la convergencia del modelo?

```
Solución
 = [0.0, 0.0]
costes = [coste( , df.area, df.precio)]
for i = 1:10
    gradiente_descendente!(, df.area, df.precio, 0.0001)
    push!(costes, coste(, df.area, df.precio))
end
costes
11-element Vector{Float64}:
 1.3106916364659266e13
 1.114133369099188e20
 1.0856750832581238e27
 1.05794371802143e34
 1.0309206941949286e41
 1.004587918634273e48
 9.789277603492545e54
 9.539230386975057e61
 9.29557011881276e68
 9.058133657380397e75
 8.826762028174244e82
Si la tasa de aprendizaje es demasiado grande, el algoritmo del gradiente
descendente puede no converger y el coste puede oscilar en lugar de disminuir.
En este caso, el coste aumenta en cada iteración, lo que indica que la tasa de
```

aprendizaje es demasiado grande.