

# Apuntes

201<sub>200</sub>¡Hola! Como tu profesor experto en **Análisis de Datos Masivos**, he analizado exhaustivamente todas las transcripciones de las sesiones. He destilado el contenido eliminando el ruido administrativo para entregarte esta **Guía Maestra de Estudio**.

Esta guía no sigue el orden cronológico de las clases, sino el **flujo lógico de un proyecto de datos**, priorizando lo que será evaluado en el examen (interpretación y programación en R).

---



## GUÍA MAESTRA DE ESTUDIO: Análisis de Datos con R

### 1. PREPARACIÓN Y LIMPIEZA DE DATOS (El cimiento)

Antes de modelar, debes entender y limpiar los datos. Si los datos son basura, el modelo será basura (“Garbage in, Garbage out”).

#### Conceptos Clave

- **Tipos de Variables:**

- **Numéricas (int, num):** Para cálculos matemáticos (medias, correlaciones).
  - **Factor (factor):** Crucial para variables categóricas (texto con opciones limitadas, ej: “Sí/No”, “Rojo/Verde”). R asigna un código interno a cada categoría.
  - **Valores Nulos (NA):** Deben tratarse antes de modelar. Opciones: Eliminar la fila (si son pocos datos) o imputar (reemplazar por la media o cero).

#### Código Recurrente

```
# Carga de datos
library(readxl)
datos <- read_excel("ruta/archivo.xlsx") # Para Excel
```

```

datos <- read.csv("ruta/archivo.csv")      # Para CSV (más
                                             # estándar)

# Exploración inicial (¡OBBLIGATORIO!)
dim(datos)       # Dimensiones (filas, columnas)
str(datos)        # Estructura y tipos de datos
summary(datos)   # Resumen estadístico (Media, Mediana,
                  # Cuantiles)

# Tratamiento de Nulos
colSums(is.na(datos))          # Cuenta nulos por columna
datos_limpios <- na.omit(datos) # Elimina filas con nulos
                               # (Recomendado para examen si hay pocos)

# Conversión a Factor (Vital para Clasificación)
datos$columna_texto <- as.factor(datos$columna_texto)

```

---

## 2. ANÁLISIS DESCRIPTIVO Y RELACIONAL

### Métricas Clave

- **Media vs. Mediana:** Si son muy diferentes, indica presencia de **Outliers** (valores anómalos). La mediana es más robusta a outliers.
- **Cuantiles:**
  - *1er Cuantil:* El 25% de los datos es menor o igual a este valor.
  - *3er Cuantil:* El 75% de los datos es menor o igual a este valor.
- **Covarianza:** Indica el **signo** de la relación (positiva o negativa), pero no la fuerza.
- **Correlación (Pearson):** Indica **fuerza** y signo. Va de -1 a 1.
  - 0: Independientes.
  - 0 a 0.5: Débil.
  - 0.5 a 1: Fuerte.

▪ **Nota:** Correlación no implica causalidad.

### Código Recurrente

```

# Solo con variables numéricas
cor(datos_num) # Matriz de correlación
cov(datos_num) # Matriz de covarianza

# Selección de numéricas (con dplyr)
library(dplyr)
datos_num <- select_if(datos,
                        is.numeric)

```

---

### **3. APRENDIZAJE SUPERVISADO (Predicción y Clasificación)**

**Concepto:** Tienes una columna “objetivo” (solución) y datos históricos. Buscas predecir ese valor. Requiere dividir en **Entrenamiento (Training)** y **Testeo (Test)**.

#### **A. Regresión Lineal (lm)**

- **Objetivo:** Predecir un número continuo (ej. Precio, Ventas).
- **Interpretación:**
  - **P-valor ( $\text{Pr}(>|t|)$ ):** Si es **< 0.05**, la variable es **significativa** (importante). Si es mayor, se puede despreciar.
  - **R-cuadrado ( $R^2$ ):** Qué porcentaje de la variabilidad explica el modelo (cerca de 1 es bueno).
  - **Coeficiente:** Cuánto aumenta la variable objetivo si aumenta en 1 la variable explicativa.

#### **B. Clasificación (Logística y Árboles)**

- **Objetivo:** Predecir una categoría (Sí/No, 0/1).
- **Regresión Logística (glm):** Solo para 2 categorías (Binomial). Da una probabilidad. Corte habitual en 0.5.
- **Árboles de Decisión (rpart):** Reglas visuales. Sirve para más de 2 categorías.
- **Validación (Matriz de Confusión):**
  - **Accuracy:** % de aciertos totales.
  - **Sensibilidad:** Capacidad de detectar Positivos (ej. enfermos).
  - **Especificidad:** Capacidad de detectar Negativos (ej. sanos).

## Código Recurrente

```
# División Train/
# Test
# (Crucial)
library(caret)
set.seed(123) # Para
# reproducibilidad
indice <-
  createDataPartition(datos$V1,
  p=0.8, list=FALSE)
train <-
  datos[indice, ]
test <- datos[-indice, ]

# Regresión Lineal
modelo_lm <- lm(Y
~ ., data
= train)
summary(modelo_lm)
# Para ver
# p-valores y
# R2

# Regresión
# Logística
modelo_glm <-
  glm(Y ~ .,
  data =
  train,
  family =
  "binomial")
prediccion_prob <-
  predict(modelo_glm,
  newdata = test,
  type = "response")
prediccion_clase <-
  ifelse(prediccion_prob
  > 0.5, 1, 0) #
# Convertir a 0/1

# Árbol de Decisión
library(rpart);
library(rpart.plot)
modelo_arbol <-
  rpart(Y
~ ., data
= train,
method =
"class")
rpart.plot(modelo_arbol)
# Dibujo del
árbol
```

```
# Matriz de  
Confusión  
confusionMatrix(as.factor(prediccion)  
as.factor(test$Y))
```

---

## 4.

# APRENDIZAJE NO SUPERVISADO (Agrupación)

**Concepto:** No hay columna objetivo (no hay solución previa). Buscas patrones o grupos. **NO** se divide en Train/Test.

## Clustering (K-Means)

- **Objetivo:** Agrupar datos por similitud.
- **Elección de K (Número de grupos):**
  - **Método del Codo:** Donde la curva dobla.
- **Función nb clust:**  
Regla de la mayoría (Recomendada).
- **Interpretación:** Usar un Árbol de Decisión sobre el resultado del clúster para entender las características de cada

grupo (el  
“truco”).

## Código Recurrente

```
library(factoextra);
  library(NbClust)
# Calcular número óptimo
NbClust(datos_num,
  min.nc=2,
  max.nc=8,
  method="kmeans")

# Ejecutar K-Means
modelo_km <-
  kmeans(datos_num,
  centers = 2) # Si elegimos 2
datos$cluster <-
  modelo_km$cluster
# Guardar resultado
```

---

## 5. SERIES TEMPORALES

**Concepto:**  
Datos ordenados cronológicamente. Importante definir la frecuencia.

### Componente s y Predicción

- Componentes:  
Tendencia, Ciclo, Estacionalidad, Ruido.

- **Modelo:**  
Usamos `auto.arim` a para que ajuste automática mente.
- **Predicción:**  
Usamos `forecast`.

## Código Recurrente

```
library(forecast)
# Crear objeto serie temporal (ej. mensual frecuencia 12)
serie <- ts(datos$valor, start=c(2000,1), frequency=12)

# Modelo y Predicción
modelo_arima <- auto.arima(serie)
prediccion <- forecast(modelo_arima, h=12) # Predecir 12 periodos
plot(prediccion)
```

---



## LIBRERÍAS ESENCIALES

Librería	Función Principal
----------	-------------------

<code>readxl</code>	<code>read_excel</code>
---------------------	-------------------------

<code>readr</code>	<code>read_csv</code>
--------------------	-----------------------

---

Librería	Función Principal
----------	-------------------

---

<b>dplyr</b>	select, filter, %>%
--------------	---------------------

<b>caret</b>	createDataPartition confusionMatrix
--------------	--

<b>rpart</b>	rpart
--------------	-------

<b>rpart.plot</b>	rpart.plot
-------------------	------------

<b>factoextra</b>	fviz_nbclust
-------------------	--------------

<b>NbClust</b>	NbClust
----------------	---------

<b>forecast</b>	auto.arima, forecast
-----------------	----------------------

<b>ggplot2</b>	ggplot
----------------	--------

---

---



# PREGUNTAS DE EXAMEN PROBABLES

(Basadas en el énfasis repetitivo del profesor durante las sesiones)

## 1. Interpretación de Regresión:

- “Dado el siguiente *summary*, ¿es significativa la variable *precio*? ”

- Res puesta: Mirar el *p-value*. Si es < 0.05, SÍ es significativa. Si tiene asteriscos

(\*\*  
\*),  
es  
mu  
y  
sign  
ific  
ativ  
a.

1. **Dif**  
**ere**  
**nci**  
**a**  
**Sup**  
**ervi**  
**sad**  
**o**  
**vs.**  
**No**  
**Sup**  
**ervi**  
**sad**  
**o:**

■ “  
j  
Q  
u  
é  
a  
l  
g  
o  
r  
i  
t  
m  
o  
u  
s  
a  
r  
í  
a  
s  
p  
a  
r  
a  
s  
e  
g  
m  
e  
n

*t  
a  
r  
c  
l  
i  
e  
n  
t  
e  
s  
s  
i  
n  
c  
o  
n  
o  
c  
e  
r  
g  
r  
u  
p  
o  
s  
p  
r  
e  
v  
i  
o  
s  
?*

*ζ  
Y*

*p  
a  
r  
a  
p  
r  
e  
d  
e  
c  
i  
r  
s  
i  
u  
n  
c*

*l  
i  
e  
n  
t  
e  
a  
b  
a  
n  
d  
o  
n  
a  
r  
á  
(  
*S*  
*i*  
/  
*N*  
*o*  
)  
?  
,,*

▪ **R**  
e  
s  
p  
u  
e  
s  
t  
a  
:  
S  
e  
g  
m  
e  
n  
t  
a  
r  
=  
C  
l  
u  
s  
t  
e  
r  
i  
n

g  
(  
N  
o  
s  
u  
p  
e  
r  
v  
i  
s  
a  
d  
o  
)  
.A  
b  
a  
n  
d  
o  
n  
o  
=C  
l  
a  
s  
i  
f  
i  
c  
a  
c  
i  
ó  
n  
/  
L  
o  
g  
í  
s  
t  
i  
c  
a  
(  
S  
u  
p  
e

r  
v  
i  
s  
a  
d  
o  
)

. 1. M  
a  
t  
r  
i  
z  
d  
e  
C  
o  
n  
f  
u  
s  
i  
ó  
n  
:

■ “

I  
n  
t  
e  
r  
p  
r  
e  
t  
a  
l  
a  
S  
e  
n  
s  
i  
b  
i  
l  
i  
d  
a  
d  
o  
c

*a  
l  
c  
u  
l  
a  
e  
l  
A  
c  
c  
u  
r  
a  
c  
y  
.,"*

**▪ R**  
**e**  
**s**  
**p**  
**u**  
**e**  
**s**  
**t**  
**a**  
**:**  
**A**  
**c**  
**c**  
**u**  
**r**  
**a**  
**c**  
**y**  
**=**  
**(**  
**A**  
**c**  
**i**  
**e**  
**r**  
**t**  
**o**  
**s**  
**/**  
**T**  
**o**  
**t**  
**a**  
**l**  
**.**

S  
e  
n  
s  
i  
b  
i  
l  
i  
d  
a  
d  
= C  
a  
p  
a  
c  
i  
d  
a  
d  
d  
e  
d  
e  
t  
e  
c  
t  
a  
r  
p  
o  
s  
i  
t  
i  
v  
o  
s  
(  
r  
e  
c  
o  
r  
d  
a  
r  
e  
j  
e  
m

p  
l  
o  
T  
e  
s  
t  
C  
O  
V  
I  
D  
)

1. E  
l  
e  
c  
c  
i  
ó  
n  
d  
e  
K  
e  
n  
C  
l  
u  
s  
t  
e  
r  
i  
n  
g  
:

“  
S  
e  
g  
ú  
n  
l  
a  
s  
a  
l  
i  
d  
a  
d  
e

*N  
b  
C  
l  
u  
s  
t,  
f  
o  
r  
m  
a  
r  
?  
”*















