



# Fundamentos de Análisis de Datos

## **Unidad 7**

Predicción y clasificación con  
técnicas numéricas

# Referencias

- García, J., Molina, J., Berlanga, A., Patricio, M., Bustamante, Á., & Padilla, W. (2018). Ciencia de datos: técnicas analíticas y aprendizaje estadísticos en un enfoque práctico (1st ed.). Alfaomega Altaria.  
<https://visorweb.utpl.edu.ec/reader/ciencia-de-datos-tecnicas-analiticas-y-aprendizaje-estadisticos-en-un-enfoque-practico?location=118> (apartado 4, pág 117)

# Agenda

- Predicción y clasificación con técnicas numéricas
  - Concepto
  - Técnicas
  - Pasos
- Regresión logística
- Árboles de decisión
- Random Forest (bosques aleatorios)
- Clasificación bayesiana

# Predicción y clasificación con técnicas numéricas

- El análisis predictivo utiliza diversas técnicas y modelos estadísticos para predecir resultados futuros y tomar decisiones estratégicas. Dos de los tipos más comunes de análisis predictivo son la **clasificación** y la **regresión**.
- **Clasificación:** asigna observaciones a categorías específicas. Se basa en algoritmos que analizan las características y propiedades de los datos disponibles para realizar predicciones precisas.
- **Regresión:** se utiliza para predecir valores numéricos o continuos. A diferencia del análisis de clasificación, que asigna categorías, el análisis de regresión estima el valor de una variable dependiente en función de una o más variables independientes.

# Predicción y clasificación con técnicas numéricas

Algunas de las técnicas más utilizadas para predicción y clasificación son:

Técnica	Usos
<b>Regresión Lineal:</b> asume una relación lineal entre la variable dependiente y una variable independiente.	Previsión de ventas, precios de viviendas, análisis de tendencias económicas
<b>Regresión Lineal Múltiple:</b> extensión de la regresión lineal que utiliza múltiples variables independientes.	Predicción de ingresos, análisis financiero
<b>Regresión Logística:</b> utilizada para modelar una variable dependiente binaria	Diagnóstico médico, predicción de fraude, análisis de comportamiento del consumidor.
<b>Árboles de decisión:</b> Modelos que dividen los datos en ramas basadas en reglas de decisión	Clasificación de clientes, análisis de riesgo, diagnóstico médico.
<b>Bosques Aleatorios (Random Forests):</b> Combina múltiples árboles de decisión para mejorar la precisión y reducir el sobreajuste.	Clasificación de texto, análisis de imágenes, predicción de ventas
<b>Máquinas de Soporte Vectorial (SVM):</b> Encuentran el hiperplano que maximiza el margen entre clases	Clasificación de imágenes, detección de fraude, análisis de texto
<b>K-Vecinos Más Cercanos (K-NN):</b> Clasifica una observación basada en las etiquetas de sus vecinos más cercanos	Reconocimiento de patrones, análisis de imagen, clasificación de texto
<b>Series Temporales (ARIMA):</b> modelos que capturan la autocorrelación en los datos de series temporales	Previsión de demanda, análisis financiero, meteorología
<b>Regresión Logística Multinomial:</b> Extensión de la regresión logística para variables dependientes categóricas con más de dos niveles	Análisis de mercado, clasificación de documentos, diagnóstico médico
<b>Redes neuronales:</b> Modelos inspirados en la estructura del cerebro humano, utilizados para capturar patrones complejos	Reconocimiento de voz, imágenes, predicción de series temporales complejas
<b>Naive Bayes:</b> Se basa en el teorema de Bayes para calcular la probabilidad de pertenencia a una clase	Clasificación de texto y análisis de sentimientos.

# Predicción y clasificación con técnicas numéricas

- Estas técnicas numéricas de predicción y clasificación son herramientas poderosas en el análisis de datos y el aprendizaje automático.
- La elección de la técnica adecuada depende del tipo de problema, la naturaleza de los datos y los objetivos del análisis.
- Es común probar varias técnicas y comparar sus rendimientos utilizando métricas de evaluación para seleccionar el modelo más adecuado.

UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA

# Regresión logística

- **Propósito:** Modelar la relación entre una variable dependiente BINARIA y múltiples variables independientes.
- **Descripción:** Es un tipo de modelo lineal generalizado (GLM) que se utiliza para clasificación binaria. A diferencia de la regresión lineal, que predice un valor continuo, la regresión logística predice una probabilidad que luego se mapea a una de dos posibles categorías.
- **Aplicaciones:** En cualquier ámbito en el que necesitemos predecir la ocurrencia o no de un evento (Ej: deserción de un estudiante, otorgar un crédito, compra de un artículo por un cliente, presencia de una enfermedad)
- **Condiciones:**
  - Aplica para la predicción de **variables categóricas BINARIAS**
  - Si los predictores (variables independientes) son variables categóricas, se las debe **factorizar**.

# Regresión logística

- Entre dos categorías A y B, la regresión logística estima la probabilidad de que una observación pertenezca a la 2da categoría. Si la probabilidad es mayor a 0.5 se asume que pertenece a B.
- **Ejemplo** (dataset “**titanic.csv**”):
  - **Objetivo:** Mediante regresión logística construir un modelo que permita predecir el status de supervivencia de cada pasajero, a partir de las variables independientes: *Pclass*, *Sex*, *Age*, *SibSp*, *Parch*, *Fare*, y *Embarked*.
  - **Librería a utilizar:** **caret**
    - Incluye un conjunto de funciones que facilitan el uso de métodos complejos de clasificación y regresión. Para regresión logística se usa la función **glm()**



# Regresión logística

```
7 # Instalar y cargar las librerías necesarias
8 install.packages("caret")
9 library(dplyr)
10 library(ggplot2)
11 library(caret)
12
13 # Cargar y explorar datos
14 df <- read.csv("titanic.csv")
15 head(df)
16 names(df)
17 str(df)
18
19 # seleccionar variables relevantes y eliminar valores faltantes
20 df_clean <- df %>%
21   select(Survived, Pclass, Sex, Age, SibSp, Parch, Fare, Embarked)
22
23 sapply(df_clean, function(x) mean(is.na(x) | !nzchar(x)))
24
25 df_clean <- df_clean %>%
26   filter(!is.na(Age) & nzchar(Embarked))
27
28 # Convertir variables categóricas a factores
29 df_clean$Survived <- as.factor(df_clean$Survived)
30 df_clean$Pclass <- as.factor(df_clean$Pclass)
31 df_clean$Sex <- as.factor(df_clean$Sex)
32 df_clean$Embarked <- as.factor(df_clean$Embarked)
33 str(df_clean)
```

# Regresión logística

```
35 # Dividir los Datos en Conjuntos de Entrenamiento y de Prueba
36 # 80% entrenamiento, 20% prueba, y generar un vector con
37 # los índices de las obseraciones seleccionadas
38 set.seed(88)
39 train_index <- createDataPartition(df_clean$Survived,
40                                   p = 0.8,
41                                   list = FALSE)
42 train_data <- df_clean[train_index, ] #Datos de entrenamiento
43 test_data <- df_clean[-train_index, ] #Datos de prueba
44
45 # Ajustar el Modelo de Regresión Logística
46 # Variables dependiente (a predecir): Survived
47 model <- glm(Survived ~ Pclass + Sex + Age + SibSp +
48             Parch + Fare + Embarked,
49             data = train_data,
50             family = binomial)
51
52 # Resumen del modelo
53 summary(model)
54 model$coefficients
```

Para que se ajuste a un modelo de regresión logística (tipo de variable binaria)

Coefficientes  $\beta$  (logaritmo de la razón de probabilidades) del modelo. Indica el sentido del cambio de valor de la variable dependiente, por cada aumento del valor del predictor, mientras el resto de valores se mantiene constante

(Intercept)	Pclass2	Pclass3	Sexmale	Age	SibSp	Parch
4.2769382253	-1.1208987508	-2.4797172142	-2.5819978513	-0.0432734526	-0.3731713769	0.0462173295
Fare	EmbarkedQ	EmbarkedS				
0.0001232415	-0.6217113936	-0.2232999545				

UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA

# Regresión logística

```
Call:
glm(formula = Survived ~ Pclass + Sex + Age + SibSp + Parch +
     Fare + Embarked, family = binomial, data = train_data)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	4.2769382	0.5813403	7.357	1.88e-13	***
Pclass2	-1.1208988	0.3658516	-3.064	0.00219	**
Pclass3	-2.4797172	0.3781737	-6.557	5.49e-11	***
Sexmale	-2.5819979	0.2488561	-10.375	< 2e-16	***
Age	-0.0432735	0.0090663	-4.773	1.82e-06	***
SibSp	-0.3731714	0.1416637	-2.634	0.00843	**
Parch	0.0462173	0.1466590	0.315	0.75266	
Fare	0.0001232	0.0029058	0.042	0.96617	
EmbarkedQ	-0.6217114	0.6500200	-0.956	0.33885	
EmbarkedS	-0.2233000	0.3141234	-0.711	0.47717	

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 770.64 on 570 degrees of freedom  
Residual deviance: 508.74 on 561 degrees of freedom  
AIC: 528.74

Number of Fisher Scoring iterations: 5

Medida del coeficiente  $\beta$  (logaritmo de la razón de probabilidades) asociado a la variable dependiente cuando todos los predictores con cero

P-value < 0.05 => la incidencia de la variable en el modelo de regresión es significativa

Coefficientes del modelo. Indica el sentido del cambio de valor de la variable dependiente, por cada aumento del valor del predictor, mientras el resto de valores se mantiene constante

Criterio de información de Akaike: mide la calidad del modelo, sirve para comparar modelos, mientras más bajo, mejor se ajusta el modelo.

# Regresión logística

```
56 # EVALUAR EL MODELO
57 # Predecir en el conjunto de prueba
58 predictions <- predict(model, test_data, type = "response")
59
60 # Convertir probabilidades a etiquetas binarias
61 predicted_classes <- ifelse(predictions > 0.5, 1, 0)
62
63 # Matriz de confusión
64 conf_matrix <- confusionMatrix(factor(predicted_classes),
65                                test_data$Survived,
66                                positive = "1")
67
68 # Ver la matriz de confusión
69 # y las métricas de evaluación
70 conf_matrix
71
```

Devuelve valores en una escala de probabilidades (probabilidad de que se cumpla la hipótesis alternativa)(que Survived = 1)

Exactitud: Proporción de predicciones correctas.

## Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	73	16
1	11	41

Accuracy : 0.8085  
95% CI : (0.7338, 0.8699)

No Information Rate : 0.5957  
P-value [Acc > NIR] : 5.647e-08

Kappa : 0.5968

McNemar's Test P-Value : 0.4414

Sensitivity : 0.7193  
Specificity : 0.8690

'Positive' Class : 1

Verdaderos negativos

Falsos positivos

	Reference	
Prediction	0	1
0	73	16
1	11	41

Falsos negativos

Verdaderos positivos

# Regresión logística

## Confusion Matrix and Statistics

Reference  
Prediction 0 1  
0 73 16  
1 11 41

Accuracy : 0.8085  
95% CI : (0.7338, 0.8699)

No Information Rate : 0.5957  
P-Value [Acc > NIR] : 5.647e-08

Kappa : 0.5968

McNemar's Test P-Value : 0.4414

Sensitivity : 0.7193  
Specificity : 0.8690  
Pos Pred Value : 0.7885  
Neg Pred Value : 0.8202  
Prevalence : 0.4043  
Detection Rate : 0.2908  
Detection Prevalence : 0.3688  
Balanced Accuracy : 0.7942

'Positive' Class : 1

Exactitud: Proporción de predicciones correctas.

Intervalo de confianza. La precisión estará entre 73.38% y 86.99% el 95% de las veces

Probabilidad de llegar a los mismos resultados por el azar

Proporción de verdaderos positivos correctamente detectados (41/57)

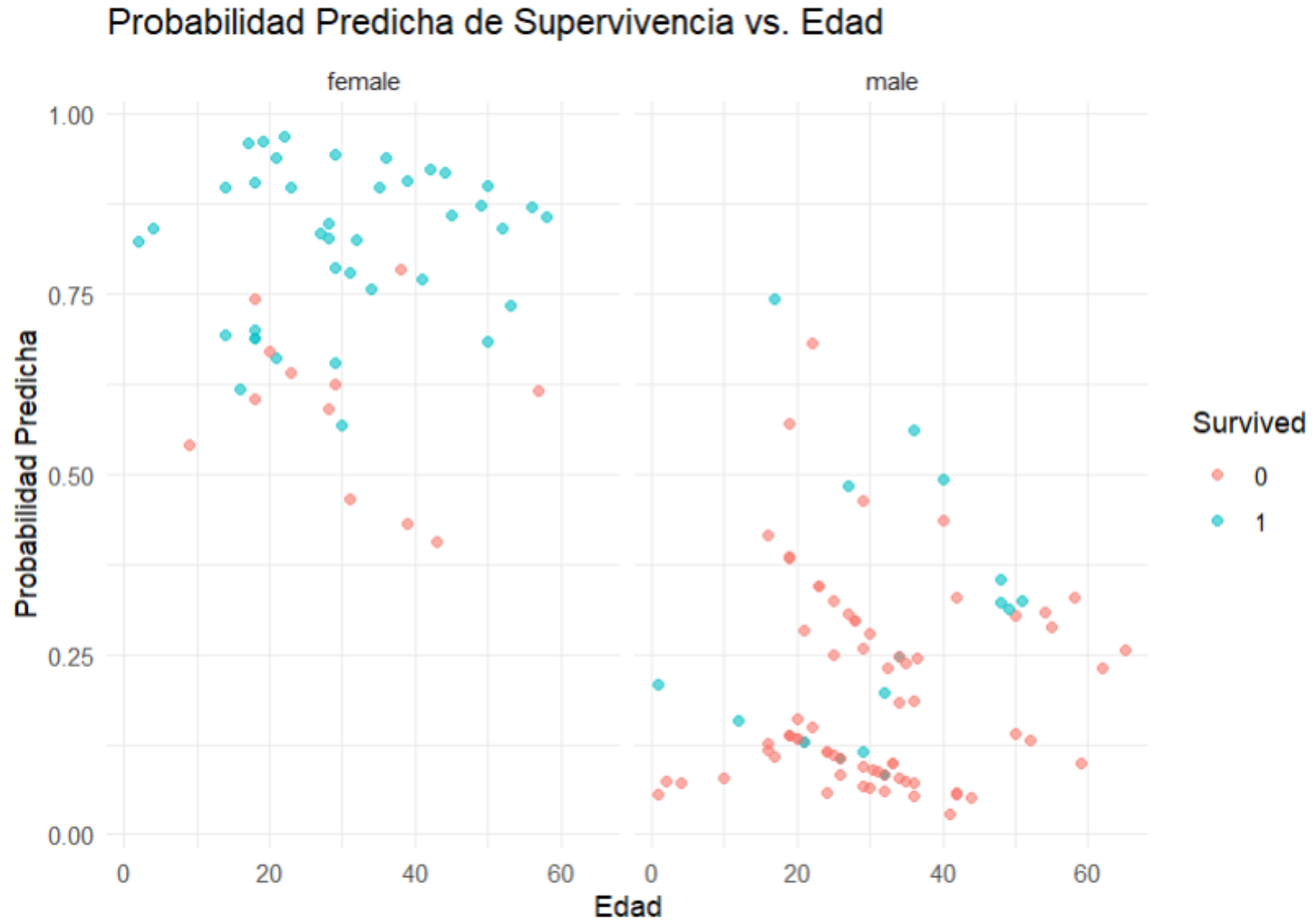
Proporción de verdaderos negativos correctamente identificados (73/84)

# Regresión logística

```
82 # visualizar los Resultados
83
84 # Añadir las predicciones al conjunto de prueba
85 test_data$predicted_prob <- predictions
86
87 # Gráfico de las probabilidades predichas vs. Edad,
88 # separado por Sexo y Supervivencia
89 ggplot(test_data, aes(x = Age, y = predicted_prob,
90                       color = Survived)) +
91   geom_point(alpha = 0.6) +
92   facet_wrap(~ Sex) +
93   labs(title = "Probabilidad Predicha de Supervivencia vs. Edad",
94        x = "Edad",
95        y = "Probabilidad Predicha") +
96   theme_minimal()
```



# Regresión logística



# Ejercicio 1

- Desarrollar tarea en clase sobre Regresión Logística





# Árboles de decisión

- **Propósito:** Predecir el valor de una variable objetivo (dependiente) basándose en varias variables independiente. Son utilizados tanto para problemas de clasificación como de regresión
- **Descripción:** Es una estructura de árbol donde los nodos internos representan decisiones basadas en las variables independientes, y los nodos terminales (hojas) representan los resultados finales (predicciones).
- **Aplicaciones:** Clasificación de enfermedades a partir de síntomas, segmentación de clientes, pronóstico de ventas, predicción de precios, decisiones de crédito etc.
- **Condiciones:**
  - Aplica para la predicción de **variables categóricas** o **continuas**
  - Si los predictores (variables independientes) son variables categóricas, se las debe **factorizar**.

# Árboles de decisión

- Se debe propender a que el árbol resultante guarde equilibrio entre precisión y complejidad. Esto es, el árbol no debe estar:
  - **Sobreajustado**: aumentar complejidad sin mejora significativa de precisión.
  - **Subajustado**: aumentar simplicidad disminuyendo drásticamente la precisión.
- Si es necesario se debe “podar” el árbol para mejorar el modelo.
- **Ejemplo** (dataset “**titanic.csv**”):
  - **Objetivo**: Mediante árboles de decisión construir un modelo que permita predecir el status de supervivencia de cada pasajero, a partir de las variables independientes: *Pclass*, *Sex*, *Age*, *SibSp*, *Parch*, *Fare*, y *Embarked*.
  - **Librería a utilizar**: **rpart** y **rpart.plot**

# Árboles de decisión

## Clasificación

```
11 # Instalar y cargar las librerías necesarias
12 install.packages("rpart")
13 install.packages("rpart.plot")
14 library(dplyr)
15 library(caret)
16 library(rpart)
17 library(rpart.plot)
18
19 # Cargar y explorar datos
20 df <- read.csv("titanic.csv")
21 head(df)
22 names(df)
23 str(df)
24
25 # Seleccionar variables relevantes y eliminar valores faltantes
26 df_clean <- df %>%
27   select(Survived, Pclass, Sex, Age, SibSp, Parch, Fare, Embarked)
28
29 sapply(df_clean, function(x) mean(is.na(x) | !nzchar(x)))
30
31 df_clean <- df_clean %>%
32   filter(!is.na(Age) & nzchar(Embarked))
33
34 # Convertir variables categóricas a factores
35 df_clean$Survived <- as.factor(df_clean$Survived)
36 df_clean$Pclass <- as.factor(df_clean$Pclass)
37 df_clean$Sex <- as.factor(df_clean$Sex)
38 df_clean$Embarked <- as.factor(df_clean$Embarked)
39 str(df_clean)
```

# Árboles de decisión

## Clasificación

```
41 # Dividir los datos de entrenamiento (80%) y de Prueba (20%)
42 set.seed(88)
43 train_index <- createDataPartition(df_clean$Survived,
44                                   p = 0.8,
45                                   list = FALSE)
46 train_data <- df_clean[train_index, ] #Datos de entrenamiento
47 test_data <- df_clean[-train_index, ] #Datos de prueba
48
49 # Ajustar el Modelo de Árbol de Decisión
50 model <- rpart(Survived ~ Pclass + Sex + Age + SibSp + Parch +
51               Fare + Embarked,
52               data = train_data,
53               method = "class")
54 summary(model)
55
56 # Visualizar el árbol de decisión
57 rpart.plot(model,extra = 104)
58
59 # EVALUAR EL MODELO
60 # Predicciones en el conjunto de prueba
61 predictions <- predict(model, test_data, type = "class")
62
63 # Matriz de confusión
64 conf_matrix <- confusionMatrix(predictions,
65                                test_data$Survived,
66                                positive = "1")
67 print(conf_matrix)
```

Indica que es  
clasificación y no  
regresión

# Árboles de decisión

## Clasificación

```
Call:
rpart(formula = Survived ~ Pclass + Sex + Age + SibSp + Parch +
      Fare + Embarked, data = train_data, method = "class")
n= 571
```

	CP	nsplit	rel error	xerror	xstd
1	0.45887446	0	1.0000000	1.0000000	0.05077096
2	0.03030303	1	0.5411255	0.5411255	0.04277523
3	0.02597403	2	0.5108225	0.5757576	0.04372530
4	0.02164502	3	0.4848485	0.5367965	0.04265153
5	0.01731602	4	0.4632035	0.5281385	0.04240070
6	0.01298701	5	0.4458874	0.5238095	0.04227357
7	0.01000000	7	0.4199134	0.5411255	0.04277523

Variable importance

Sex	Pclass	Fare	Age	Parch	SibSp	Embarked
50	15	13	11	7	3	1

Node number 1: 571 observations, complexity param=0.4588745  
predicted class=0 expected loss=0.4045534 P(node) =1  
class counts: 340 231  
probabilities: 0.595 0.405  
left son=2 (363 obs) right son=3 (208 obs)  
Primary splits:  
Sex splits as RL, improve=80.276740, (0 missing)  
Pclass splits as RRL, improve=30.558450, (0 missing)  
Fare < 52.2771 to the left, improve=30.223080, (0 missing)  
Embarked splits as RLL, improve=13.516790, (0 missing)  
Age < 5.5 to the right, improve= 8.729757, (0 missing)  
Surrogate splits:  
Parch < 0.5 to the left, agree=0.669, adj=0.091, (0 split)  
Fare < 64.17915 to the left, agree=0.664, adj=0.077, (0 split)  
Age < 15.5 to the right, agree=0.639, adj=0.010, (0 split)

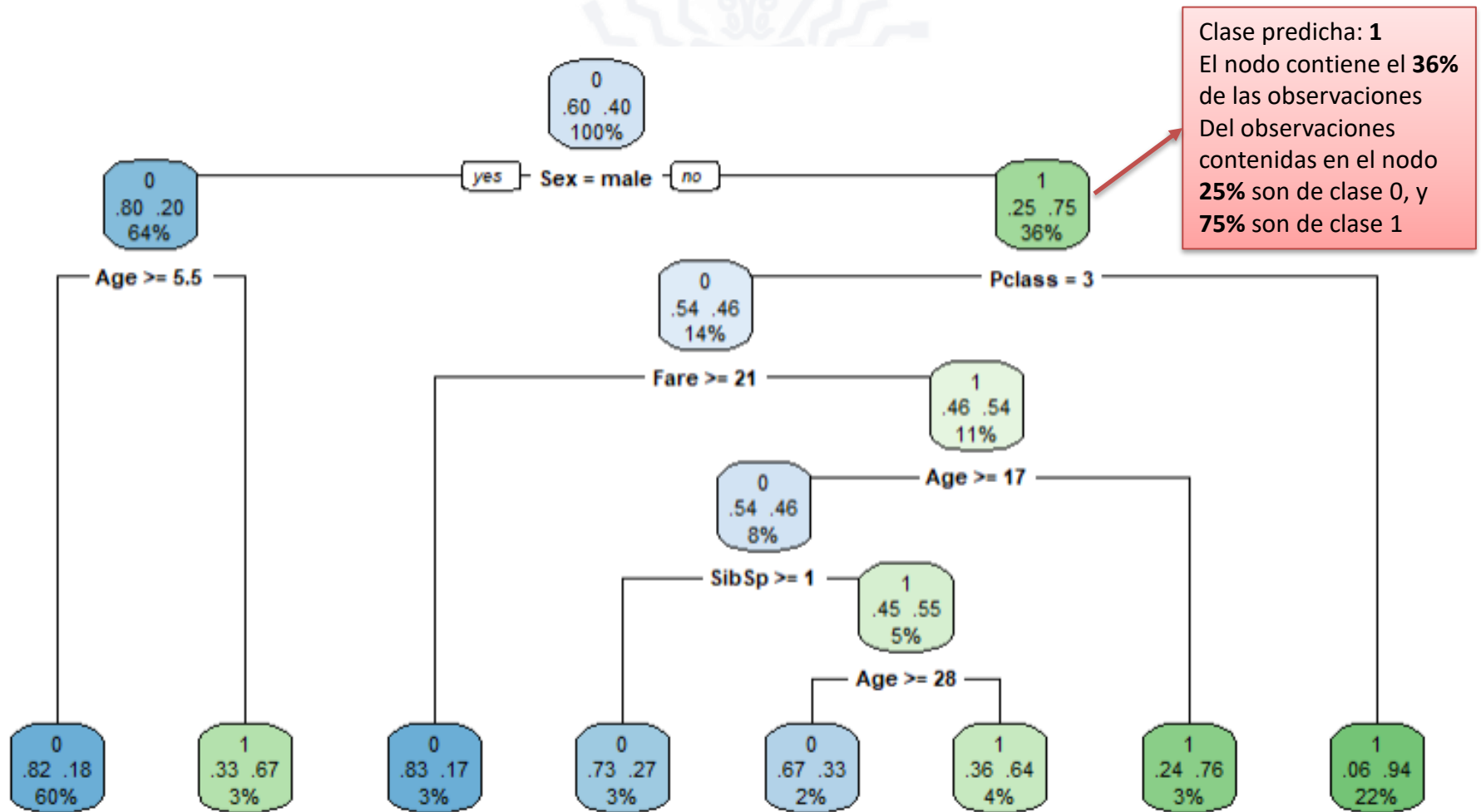
Parámetro de complejidad.  
Costo de añadir complejidad al modelo. Mientras más bajo, el árbol es más grande. Ayuda a prevenir el sobreajuste.

Error de validación cruzada comparado con el error del árbol sin ninguna división. El valor óptimo de CP es el que minimiza este error. Es donde hay mayor equilibrio entre complejidad y precisión

Importancia relativa de cada variable predictora en el modelo.

# Árboles de decisión

## Clasificación



# Árboles de decisión

- Respecto a la condiciones establecidas en cada nodo
  - El árbol elige la variable y el punto de corte (threshold) para dividir los datos, que reduce más impureza al dividir el nodo.
  - Impureza significa que hay una mezcla de clases dentro de un nodo. Cuanto más mezcladas estén mayor la impureza.
  - El algoritmo explora todos los puntos posibles de corte (basados en los valores que existen en los datos), y elige el que mejor separa las clases en ese nodo. No es un valor arbitrario: es el valor óptimo que reduce más la impureza en ese momento del árbol.



# Árboles de decisión

## Clasificación

Confusion Matrix and Statistics

Prediction \ Reference	0	1	
	0	1	
0	77	20	FN
1	7	37	VP

Accuracy : 0.8085  
95% CI : (0.7338, 0.8699)  
No Information Rate : 0.5957  
P-Value [Acc > NIR] : 5.647e-08  
Kappa : 0.5873  
McNemar's Test P-Value : 0.02092  
Sensitivity : 0.6491  
Specificity : 0.9167  
Pos Pred Value : 0.8409  
Neg Pred Value : 0.7938  
Prevalence : 0.4043  
Detection Rate : 0.2624  
Detection Prevalence : 0.3121  
Balanced Accuracy : 0.7829  
'Positive' Class : 1

Annotations: VN (True Negative) points to the 0,0 cell (77); FP (False Positive) points to the 0,1 cell (20); FN (False Negative) points to the 1,0 cell (7); VP (True Positive) points to the 1,1 cell (37).

**Precisión o exactitud:** Proporción de predicciones correctas

**Sensibilidad:** Proporción de verdaderos positivos sobre todos los casos reales positivos  $VP / (VP + FN)$ .  
En este caso: El modelo es capaz de identificar correctamente el 65% de las instancias positivas (supervivientes).

**Especificidad:** Proporción de verdaderos negativos sobre todos los casos reales negativos  $VN / (VN + FP)$ .  
En este caso: El modelo es capaz de identificar correctamente el 92% de las instancias negativas (fallecidos)



# Árboles de decisión

## Clasificación

```
62 # PODAR EL ÁRBOL BASADO EN EL VALOR ÓPTIMO DE CP
63
64 cp_optimo <- model$cptable[which.min(model$cptable[, "xerror"]), "CP"]
65
66 p_model <- prune(model, cp = cp_optimo)
67
68 summary(p_model)
69
70 # visualizar el árbol podado
71 rpart.plot(p_model)
72
73
74 # Evaluar el Modelo Podado
75 # Predicciones en el conjunto de prueba con el modelo podado
76 predictions <- predict(p_model, test_data, type = "class")
77
78 # Matriz de confusión
79 conf_matrix <- confusionMatrix(predictions, test_data$Survived)
80 print(conf_matrix)
```

# Árboles de decisión

## Clasificación

Call:

```
rpart(formula = Survived ~ Pclass + Sex + Age + SibSp + Parch +  
      Fare + Embarked, data = train_data, method = "class")  
n= 571
```

	CP	nsplit	rel error	xerror	xstd
1	0.45887446	0	1.0000000	1.0000000	0.05077096
2	0.03030303	1	0.5411255	0.5411255	0.04277523
3	0.02597403	2	0.5108225	0.5757576	0.04372530
4	0.02164502	3	0.4848485	0.5367965	0.04265153
5	0.01731602	4	0.4632035	0.5281385	0.04240070
6	0.01298701	5	0.4458874	0.5238095	0.04227357

Variable importance

Sex	Pclass	Fare	Age	Parch	SibSp	Embarked
51	15	13	10	7	2	1

Node number 1: 571 observations, complexity param=0.4588745  
predicted class=0 expected loss=0.4045534 P(node) =1  
class counts: 340 231  
probabilities: 0.595 0.405  
left son=2 (363 obs) right son=3 (208 obs)

Primary splits:

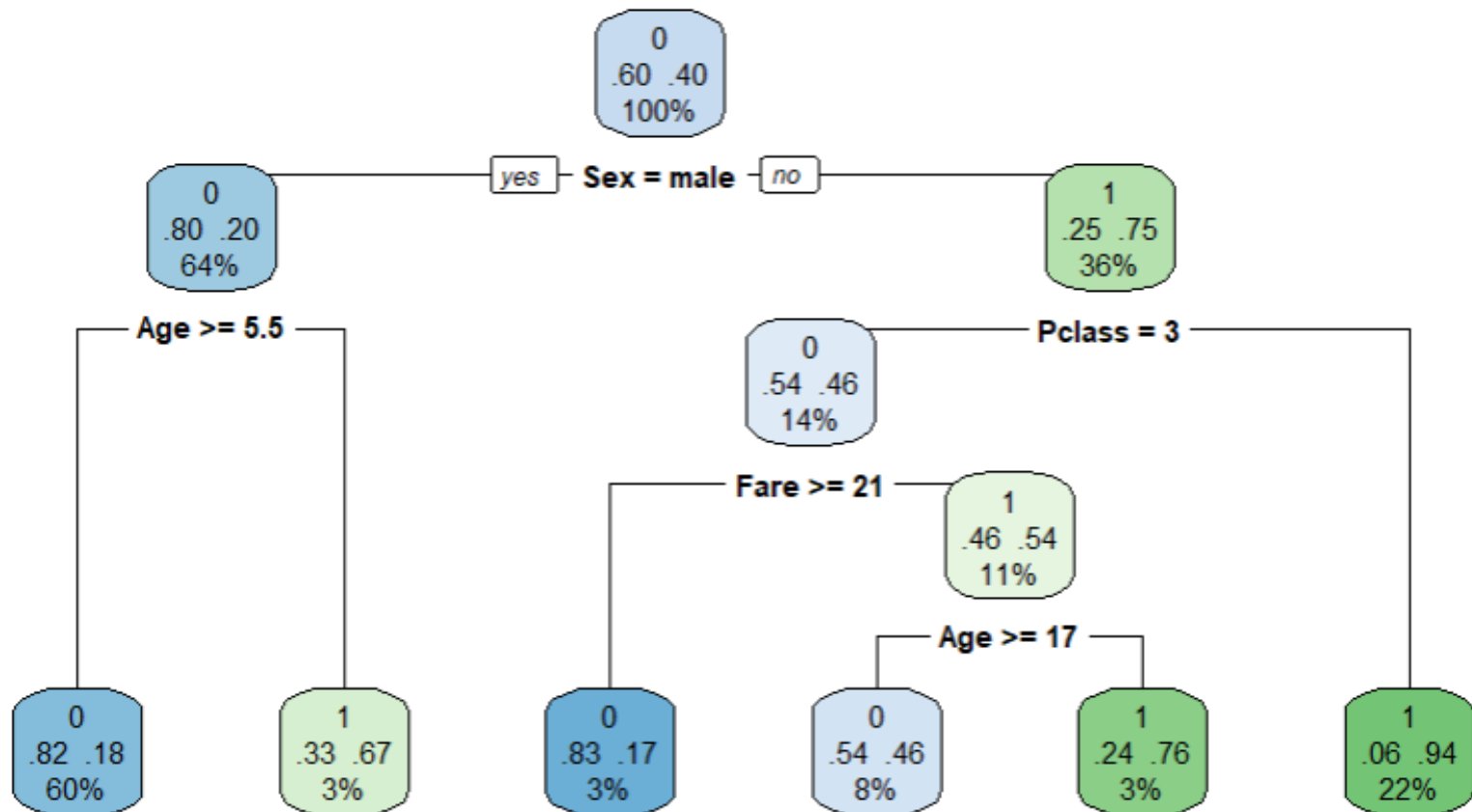
Sex	splits as	RL,	improve=80.276740, (0 missing)
Pclass	splits as	RRL,	improve=30.558450, (0 missing)
Fare	< 52.2771	to the left,	improve=30.223080, (0 missing)
Embarked	splits as	RLL,	improve=13.516790, (0 missing)
Age	< 5.5	to the right,	improve= 8.729757, (0 missing)

Surrogate splits:

Parch	< 0.5	to the left,	agree=0.669, adj=0.091, (0 split)
Fare	< 64.17915	to the left,	agree=0.664, adj=0.077, (0 split)
Age	< 15.5	to the right,	agree=0.639, adj=0.010, (0 split)

# Árboles de decisión

Clasificación



# Árboles de decisión

## Clasificación

### Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	80	23
1	4	34

Accuracy : 0.8085

95% CI : (0.7338, 0.8699)

No Information Rate : 0.5957

P-Value [Acc > NIR] : 5.647e-08

Kappa : 0.5799

McNemar's Test P-Value : 0.000532

Sensitivity : 0.5965

Specificity : 0.9524

Pos Pred Value : 0.8947

Neg Pred Value : 0.7767

Prevalence : 0.4043

Detection Rate : 0.2411

Detection Prevalence : 0.2695

Balanced Accuracy : 0.7744

'Positive' Class : 1

No se logra una mejora en la precisión pero si se logra reducir la complejidad del árbol

# Ejercicio 2

- Desarrollar la PARTE 1 de tarea en clase sobre Árboles de decisión.

