Trabajo Práctico Final Diseño de una investigación basado en una base de datos

Materia: Bioestadística Alumnos:

- Axel Cesar Wood Niella
- Lucía Bernochi
- Priscilla Vanesa Tenas Vai



#### Elección de base de datos

Presentación del diseño de una investigación basado en una base de datos de interés que en este caso se trató de biomarcadores urinarios del cáncer de páncreas y posterior análisis estadístico de la misma.

## TABLE OF CONTENTS

**O** Introducción

Breve descripción del problema, objetivo y PICO.

**Métodos** 

Incluye: diseño del estudio, lugar donde se desarrollará el estudio y marco temporal, entre otros.

O3. Data Set

Base de datos de interés y tratamiento de la misma.

**04.** Resultados

Identificación de variables, análisis univariado y multivariado.

**O 5** Conclusiones

Interpretación de los resultados obtenidos.

01.

Introducción



#### Introducción

#### **Importancia**



La tasa de supervivencia a cinco años es inferior al 10%.

#### **Actualidad**



No existen biomarcadores específicos con la capacidad de generar una detección temprana.

#### Objetivo



Encontrar una prueba de diagnóstica capaz de identificar posibilidades de desarrollar la enfermedad.

## **Objetivo**

#### Creatinina

000

000

La elevación implica menor perfusión renal provocada en los casos graves de pancreatitis.

#### **REGIB**

Actúa como factor de crecimiento en la regeneración de los islotes pancreáticos.

#### **LYVEI**

000

. . .

Receptor que tiene un papel activo en la linfangiogénesis y remodelación endotelial.

#### TFFI

Factor involucrado en el desarrollo y la progresión de varios tipos de cáncer.

## Pregunta PICO

#### Población (P)

Hombres y mujeres con un rango etario entre 26 y 89 años con historias clínicas diferentes.

#### Intervención (I)

Medición de biomarcadores.



#### Comparación (C)

Pacientes con adenocarcinoma ductal pancreatico vs. sin.

#### Outcomes (O)

Presencia o no de cáncer pancreático.

# 02.

## Métodos



#### Diseño del estudio

#### Casos y controles

Baja prevalencia

#### **Hospital Boston EEUU**

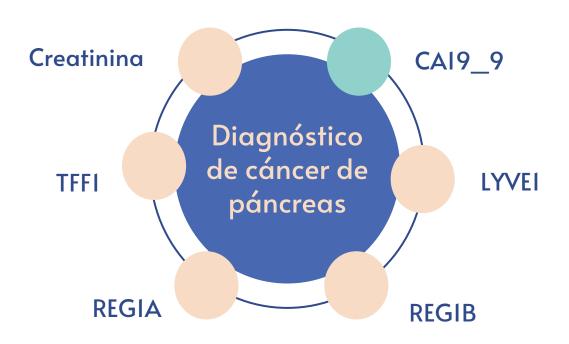
Previo a intervención quirúrgica o quimioterapéutico



#### Muestras

- Preservadas según procedimiento estándar.
- Sin afecciones pancreáticas
- Sin neoplasias malignas conocidas
- Adenocarcinoma ductal pancreatico
- Grupo benigno

## **Variables**



## Data sources / measurement



000



Grupos de 40

Muestras duplicadas para confirmación de resultados

#### Creatinina

ILab Aries del Laboratorio de Instrumentación

mmol/l

#### **REGIB**

TMB Substrate Set y Stop Solution de BioLegend

## FLUOstar Omega Microplate Reader

000

...

Límites de detección:

CA19\_9:0.3 U/ml

TFF1: 3.91 pg/ml

REG1B: 8 pg/ml

LYVE1: 56 pg/ml



## Sesgos

Sesgo de Selección Se usaron muestras de pacientes entre 26 y 89 años.

Las funciones renales de pacientes mayores de edad pueden no funcionar de la manera esperada

#### Tamaño del estudio

Determinado por la cantidad de personas que se sometieron una intervención quirúrgica o a un tratamiento quimioterapéutico para el tratamiento de cancer de pancreas durante 6 meses

#### Sexo

Se tomaron muestras tanto de hombres como mujeres

#### Edad

Se tomaron muestras de pacientes en el rango etario de interés

#### Muestras benignas

000

...

- PancreatitisCronica
- Casos de dolor Abdominal
- Enfermedades de vesícula biliar



#### Variables Cuantitativas

Variables Explicativas Fueron agrupadas dependiendo del diagnóstico del paciente. No se agrupo por edad.

Analisis Multivariado Las variables fueron agrupadas en dos grupos utilizando un valor umbral 03.

**Data Set** 





https://www.kaggle.com/johnjdavisiv/urinary-biomarkers-for-pancreatic-cancer

1

**Controles Saludables** 

3

Pacientes con adenocarcinoma ductal pancreático

2

Pacientes con afecciones pancreáticas no cancerosas, como pancreatitis crónica

Datos faltantes < 20%
En su lugar se reemplazó la media correspondiente a cada variable

000

04.

Resultados



•••

Tabla 0: Identificación de variables, su forma de medirlas y su distribución.

## Tabla 0: Identificación de variables, su forma de medirlas y su distribución.

Variable	Tipo
Sex	factor
Age	integer
Creatinine	numeric
LYVE1	numeric
REG1B	numeric
TFF1	numeric
Diagnosis	factor

## Variable: Sex

	No	Si
Fem	68	12
Masc	29	16

Variable explicatoria cualitativa nominal categórica que en consecuencia posee una distribución binomial.

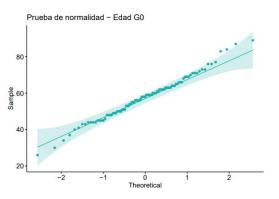
Comparación mediante TEST CHI CUADRADO.

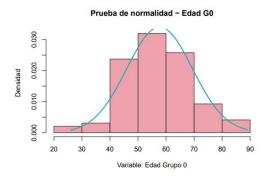
## Variables: Age, Creatinine, LYVEI, REGIB, TFFI

## En primer lugar: Análisis de normalidad.

```
##
## Shapiro-Wilk normality test
##
## data: x
## W = 0.9901, p-value = 0.6925
##
##
## Lilliefors (Kolmogorov-Smirnov) normality test
##
## data: x
## data: x
## D = 0.058762, p-value = 0.5633
```

Grupo 0



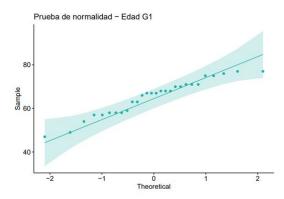


## Variables: Age, Creatinine, LYVEI, REGIB, TFFI

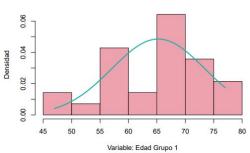
## En primer lugar: Análisis de normalidad.

```
##
## Shapiro-Wilk normality test
##
## data: x
## W = 0.9459 p-value = 0.1561
##
##
## Lilliefors (Kolmogorov-Smirnov) normality test
##
## data: x
## D = 0.1557, p-value = 0.08022
```

Grupo I







## Variables: Age, Creatinine, LYVEI, REGIB, TFFI

En segundo lugar: Comparación



## Distribución Normal

Para ambos grupos la distribución fue del tipo Normal.



0 0 0

## Test Z

Condiciones:

- Homocedasticidad
- → Muestra mayor a 30

## Tabla resumen

Variable	Respuesta o explicatoria	Tipo	Unidad	Dist. de Prob.	Test
Sex	Explicatoria	Cual. Dicotómica	-	Binomial	Chi Cuadrado
Age	Explicatoria	Cuant. cont.	Años	Normal	Test Z
Creatinine	Explicatoria	Cuant. cont.	mmol/l	No normal	Wilcoxon
LYVE1	Explicatoria	Cuant. cont.	pg/ml	No normal	Wilcoxon
REG1B	Explicatoria	Cuant. cont.	pg/ml	No normal	Wilcoxon
TFF1	Explicatoria	Cuant. cont.	pg/ml	No normal	Wilcoxon
Diagnosis	Respuesta	Cual. Dicotómica	-	Binomial	-



#### Tabla I: Analisis Univariado



H0: Se trata de variables independientes.

H1: No se trata de variables independientes



#### Test Z

 $H0: \mu Grupo0 - \mu Grupo1 = 0$ 

H1: μGrupo0 − μGrupo1 ≠ 0



 $H0: \theta Grupo0 - \theta Grupo1 = 0$ 

H1:  $\theta$ Grupo0 -  $\theta$ Grupo1  $\neq$  0

#### Tabla I: Analisis Univariado



#### Chi Cuadrado

HO: Se trata de variables independientes.

H1: No se trata de variables independientes



Test Z

 $H0: \mu Grupo0 - \mu Grupo1 = 0$ 

H1: µGrupo0 - µGrupo1 ≠ 0



 $H0: \theta Grupo0 - \theta Grupo1 = 0$ 

H1:  $\theta$ Grupo0 -  $\theta$ Grupo1  $\neq$  0

#### Variable: Sex

#### **TEST CHI CUADRADO**

```
>> chisq.test(tf, correct = TRUE) donde tf es la tabla de frecuencias
## Pearson's Chi-squared test with Yates' continuity correction
##
## data: tf
## X-squared = 5.8681, df = 1, p-value = 0.01542
Menor al n
```

Menor al nivel de significancia (5%)

Sex y el cancer de pancreas no son variables independientes

#### Tabla I: Analisis Univariado



#### Test Z

 $H0: \mu Grupo0 - \mu Grupo1 = 0$ 

H1: μGrupo0 - μGrupo1 ≠ 0



HO: Se trata de variables independientes.

H1: No se trata de variables independientes



H0:  $\theta$ Grupo0 –  $\theta$ Grupo1 = 0

H1:  $\theta$ Grupo0 −  $\theta$ Grupo1 ≠ 0

## Variable: Age

#### TEST DE LEVENE Y TEST Z

>> z.test(edad0,edad1, mu = 0, sigma.x=sigma.x, sigma.y =sigma.y,
conf.level = 0.95)



La diferencia entre las medias de edad del grupo 0 y del grupo 1 es distinta de cero

#### Tabla I: Analisis Univariado





H0: Se trata de variables independientes.

H1: No se trata de variables independientes



Test Z

 $H0: \mu Grupo0 - \mu Grupo1 = 0$ 

H1:  $\mu$ Grupo0 -  $\mu$ Grupo1  $\neq$  0

 $HO: \theta Grupo O - \theta Grupo I = O$ 

H1:  $\theta$ Grupo0 -  $\theta$ Grupo1  $\neq$  0

## Variables: Creatinine, LYVEI, REGIB, TFFI

```
## Wilcoxon rank sum test with continuity correction
##
## data: creatinina0 and creatinina1
## W = 1285.5, p-value = 0.6698
## alternative hypothesis: true location shift is not equal to 0
## 95 percent confidence interval:
## -0.2375452 0.1470490
## sample estimates:
## difference in location
## -0.03394606
```

**TEST DE WILCOXON** 

Mayor al nivel de significancia (5%)

La mediana de la creatinina del grupo 0 no difiere de la del grupo 1.

## Variables: Creatinine, LYVEI, REGIB, TFFI

```
TEST DE WILCOXON
```

Mediana -31.78906 pg/ml

95 %

• • (

(-31.79223; -31.61921) pg/ml.

La mediana de la LYVEI del grupo 0 difiere de la del grupo 1.

## Tabla resumen: Analisis univariado

Variable	Test	Hipotesis	p-valor	Inferencia
Sex	Chi Cuadrado	$H_0$ : Se trata de variables independientes	0.01542	Se rechaza $H_0$
Age	Test Z	$H_0: \mu_{Grupo0} - \mu_{Grupo1} = 0$	0.0001176	Se rechaza $H_0$
Creatinina	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	0.6698	Se acepta $H_0$
LYVE1	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	1.845e-10	Se rechaza $H_0$
REG1B	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	0.07389	Se acepta $H_0$
TFF1	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	5.53e-05	Se rechaza $H_0$

# Tabla 2: Análisis Multivariado

## Tabla resumen: Analisis univariado

Variable	Test	Hipotesis	p-valor	Inferencia
Sex	Chi Cuadrado	$H_0$ : Se trata de variables independientes	0.01542	Se rechaza $H_0$
Age	Test Z	$H_0: \mu_{Grupo0} - \mu_{Grupo1} = 0$	0.0001176	Se rechaza $H_0$
Creatinina	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	0.6698	Se acepta $H_0$
LYVE1	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	1.845e-10	Se rechaza $H_0$
REG1B	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	0.07389	Se acepta $H_0$
TFF1	Wilcoxon	$H_0: \theta_{Grupo0} - \theta_{Grupo1} = 0$	5.53e-05	Se rechaza $H_0$

## Variables dicotómicas

- Se tomaron las variables explicativas relevantes.
- Se las volvió variables dicotómicas utilizando un valor umbral (media)

LYVE1D		$\overline{\text{TF}}$	F1D
0	93	0	117
1	28	1	8

## Test de Chi cuadrado de Pearson

• 80% de las celdas deben tener una frecuencia esperada mayor o igual a 5 y ninguna frecuencia menor a 1

	Cáncer de pa	ncreas	
TFF1D	0	1	Row Total
0	79   72.168	18 24.832	97
1	14 20.832	14 7.168	28
Column Total	93	32	125

	Cáncer de pa	ncreas	
LYVE1D	0	1	Row Total
0	95	2 6. 208	97
1	22   26.208	6 1.792	28
Column Total	117	8	125

P valor = 0.001854415

• P valor = 0.001153486

## **Odds Ratio**

#### TFF1

#### Point estimates and 95% CIs:

Inc risk ratio	1.25 (1.03, 1.52)
Odds ratio	12.95 (2.45, 68.55)
Attrib risk *	19.37 (3.91, 34.83)
Attrib risk in population *	15.03 (-0.76, 30.82)
Attrib fraction in exposed (%)	19.77 (2.44, 34.03)
Attrib fraction in population (%)	16.06 (1.54, 28.43)

Test that OR = 1: chi2(1) = 13.604 Pr>chi2 = <0.001 Wald confidence limits
CI: confidence interval
\* Outcomes per 100 population units

Un valor de TFF1 alto incrementa en 12.95 las chances de tener cancer de pancreas IC(2.45,68.55)

#### LYVE1

#### Point estimates and 95% CIs:

Test that OR = 1: chi2(1) = 11.279 Pr>chi2 = <0.001 Wald confidence limits
CI: confidence interval

\* Outcomes per 100 population units

Un valor de LYVE1 alto incrementa en 4.39 las chances de tener cancer de pancreas IC(1.78,10.80)

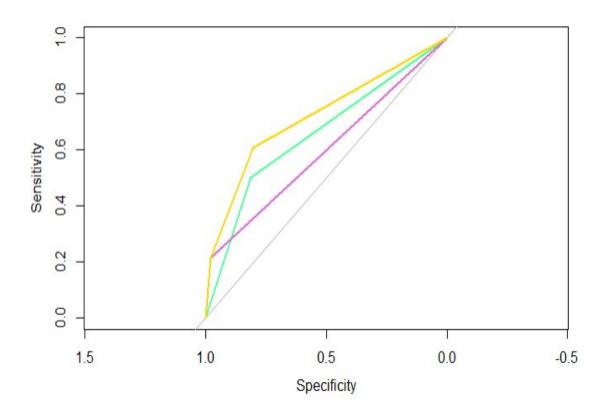
```
call:
glm(formula = diagnosis ~ TFF1D, family = "binomial", data = datafinal)
Deviance Residuals:
   Min
           1Q Median 3Q
                                  Max
-1.0727 -0.5712 -0.5712 -0.5712 1.9460
Coefficients:
          Estimate Std. Error z value Pr(>|z|)
TFF1D1 1.4791 0.4594 3.219 0.00128 **
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)
   Null deviance: 132.98 on 124 degrees of freedom
Residual deviance: 122.66 on 123 degrees of freedom
AIC: 126.66
Number of Fisher Scoring iterations: 4
```

```
call:
glm(formula = diagnosis ~ LYVE1D, family = "binomial", data = datafinal)
Deviance Residuals:
   Min 1Q Median 3Q Max
-1.6651 -0.6454 -0.6454 -0.6454 1.8282
coefficients:
          Estimate Std. Error z value Pr(>|z|)
LYVE1D1 2.5614 0.8501 3.013 0.00259 **
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '1
(Dispersion parameter for binomial family taken to be 1)
   Null deviance: 132.98 on 124 degrees of freedom
Residual deviance: 122.10 on 123 degrees of freedom
AIC: 126.1
Number of Fisher Scoring iterations: 4
```

```
call:
glm(formula = diagnosis ~ TFF1D + LYVE1D, family = "binomial",
   data = datafinal)
Deviance Residuals:
   Min
            10 Median 30
                                    Max
-2.0290 -0.5258 -0.5258 -0.5258 2.0235
Coefficients:
          Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.9090 0.3117 -6.125 9.07e-10 ***
TFF1D1 1.3998 0.4842 2.891 0.00384 **
LYVE1D1 2.4311 0.8866 2.742 0.00611 **
Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)
   Null deviance: 132.98 on 124 degrees of freedom
Residual deviance: 113.86 on 122 degrees of freedom
AIC: 119.86
Number of Fisher Scoring iterations: 4
```

## **CURVAS ROC**

- Morado LYVE1
- Verde TFF1
- Amarillo LYVE1 + TFF1



## Tabla resumen: Analisis Multivariado

	Variables Utilizadas	AUC
Modelo 1	LYVE1	0.5968
Modelo 2	TFF1	0.6572
Modelo 3	LYVE1 + TFF1	0.7204

05.

Conclusiones



## Conclusiones

. .



Las variables explicativas LYVE1 y TFF1 pueden ser utilizadas como biomarcadores para la detección temprana del cáncer de páncreas.



...

El modelo de regresión logística que provee mejores resultados es el que usa las dos variables (LYVE1 y TFF1).



. . .

En futuros trabajos utilizar otros valores umbral de modo de analizar si se consiguen mejores resultados.

Muchas Gracias!