

# El Científico de Datos



Is your Job sexy?

**ANALYZING the ANALYZERS**

ANALYTICS

INFORMATION

TIBCO®

SHARED

trend Science

Asset

Networks

capture volume

size

business

Large

INTERNET

# BIG DATA

cloud

SOFTWAR

CLUSTER

management

data

visualization

Databases

needed

amount

quantity

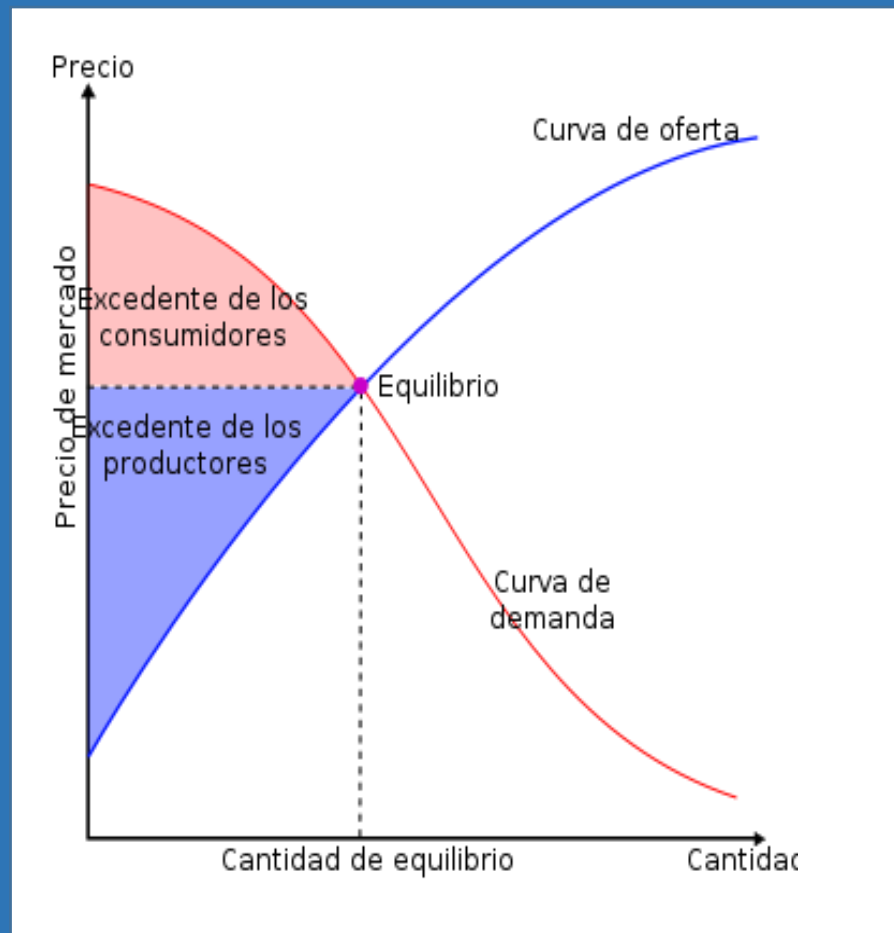
research

complex

Is your job sexy?



ANALYZING the ANALYZERS



## Científico de datos: profesión con un futuro BIG

por [Óscar Méndez](#)

Madrid, España 14 de octubre del 2015 [big data](#), [Big Data Spain](#), [empleo](#) 1 comentario

Actualmente estamos en la tercera generación de IT. Resumiendo mucho, la primera fueron los grandes servidores hosts, la segunda los sistemas cliente/servidor con clusters de dos o pocas



# GERENCIA

EXPO.ENTIC.16 26 de OCTUBRE 2016

INICIO NOTICIAS EVENTOS ARCHIVO CLASIFICADOS SUSCRIPCIONES AVISADORES PERFIL DEL MEDIO CONTACTO

Búsqueda

Miércoles 19 de Octubre de 2016 • Dólar= \$667,76 • UF=\$26.241,21 • UTM=\$45.999

GERENCIA

Ti en el Sector Salud

Edición Octubre 2016

Revisela aquí

LANIX ERP

www.lanixerp.cl

Recomendar

Twitter

Comentar

### DATA SCIENTIST O CIENTÍFICO DE DATOS

## Cómo ganar dinero con las habilidades adecuadas

El volumen de datos crece y también la necesidad de quien analice toda esa información. El "McKinsey Global Institute" estima que unas 500 mil personas serán empleadas como científicos de datos para fines de esta década.

# cibersur

DOSSIER | ENCUENTROS | CIBERSUR TV

Portada Empresa Internet Comunicaciones I+D+i Software Hardware i-Administ

cibersur.com e-Sociedad

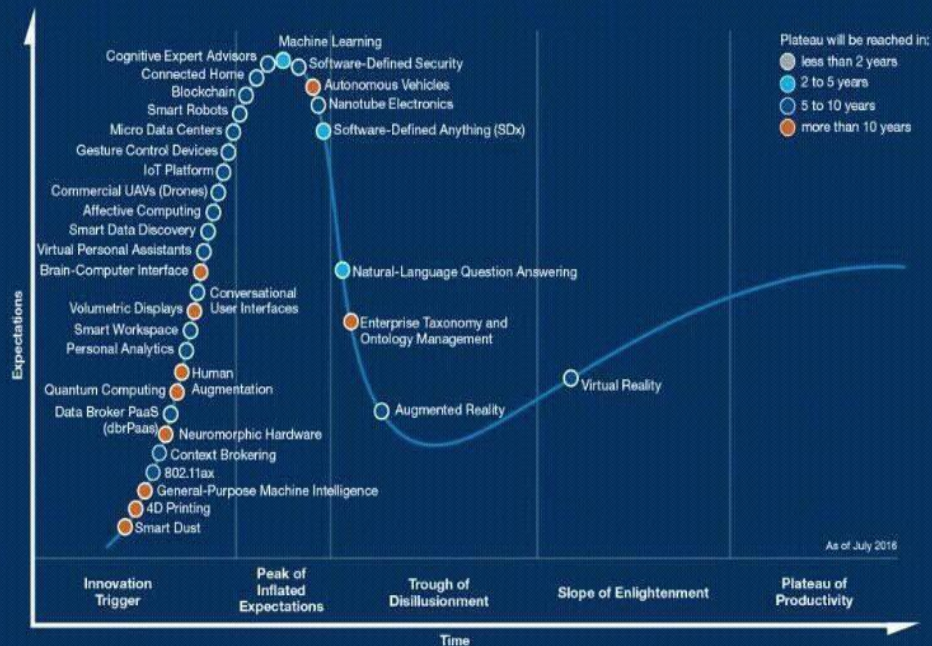
E-RRHH

## El científico de datos, entre los de mayor demanda de empleo

Nadie duda de la creciente importancia de los científicos de datos (data scientists) en las organizaciones y se estima una demanda mínima de 250.000 de estos profesionales a medio plazo en España, según MBIT School, único centro de formación español dedicado exclusivamente a Business Intelligence y Big Dat.

Según [Burtch Works](#),  
el 32% de los científicos de datos en activo vienen  
del mundo de las matemáticas y la estadística,  
el 19% de la ingeniería informática  
el 16% de otras ingenierías.

## Emerging Technology Hype Cycle for 2016



Source: Gartner  
© 2016 Gartner, Inc. and/or its affiliates. All rights reserved.

Gartner

	Mínimo	Máximo
<b>España</b>	250,000	500,000
<b>Perú</b>	36,567	73,134



# 01. ¿QUÉ ES UN CIENTÍFICO DE DATOS?

*Contexto*

## INGENIERO

Solucionar el problema del negocio.



## IT / DESARROLLADOR

Maneja la Infraestructura.



## ESTADÍSTICO

Construye el Modelo.



C i e n c i a   d e   D a t o s



# 01. ¿QUÉ ES UN CIENTÍFICO DE DATOS?

*Contexto*



## INGENIERO

Solucionar el problema del negocio.



## IT / DESARROLLADOR

Maneja la Infraestructura.



## ESTADÍSTICO

Valida el Modelo.



C i e n c i a   d e   D a t o s

## 02. ANALICE Y ACTÚE EN 'MOMENTOS CRÍTICOS DEL NEGOCIO'

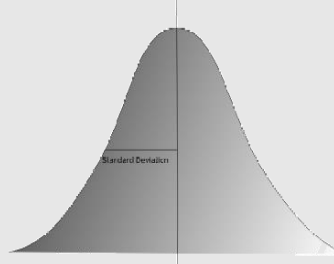
En tiempo real



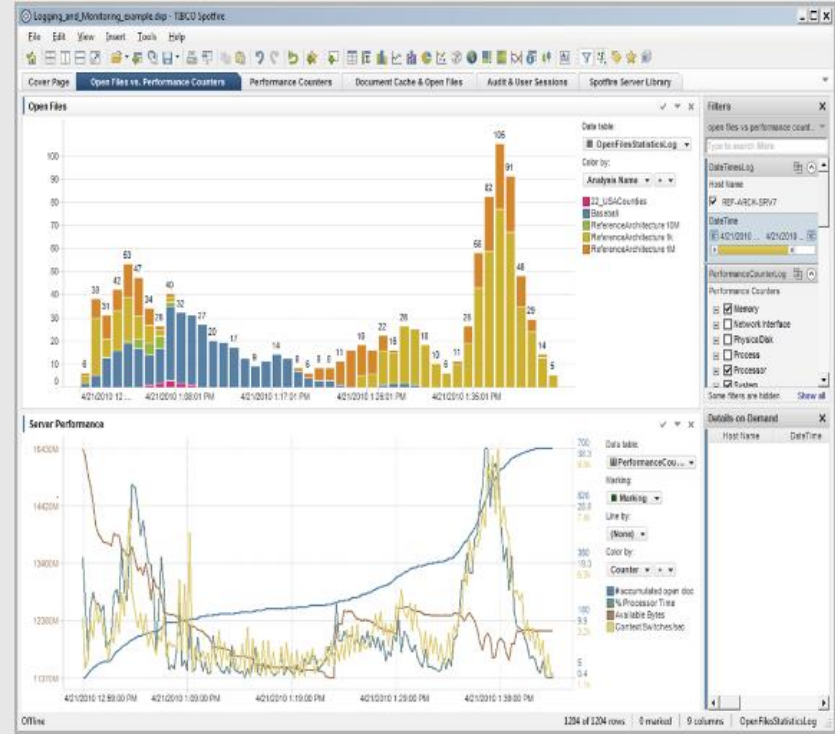
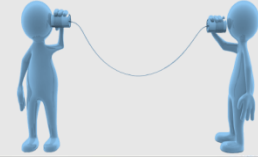
C i e n c i a   d e   D a t o s

## 02. ANALICE Y ACTÚE EN 'MOMENTOS CRÍTICOS DEL NEGOCIO'

En tiempo real



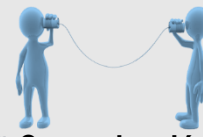
```
def add5(x):  
    return x+5  
  
def dotwrite(ast):  
    nodename = getNodeName()  
    label=symbol.sym_name.get(int(ast[0]),ast[0])  
    print ' %s [%label="%s" % (nodename, label),  
    if isinstance(ast[1], str):  
        if ast[1].strip():  
            print '%s' % ast[1]  
        else:  
            print ']'  
    else:  
        print ']'  
        children = []  
        for n, child in enumerate(ast[1:]):  
            children.append(dotwrite(child))  
        print ' %s -> {' % nodename,  
        for name in children:  
            print '%s' % name,
```



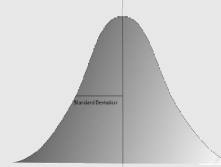


```
def add5(x):
    return x+5

def dotwrite(ast):
    nodename = getNodename()
    label=symbol.sym_name.get(int(ast[0]),ast[0])
    print ' %s [%label="%s" % (nodename, label),
    if isinstance(ast[1], str):
        if ast[1].strip():
            print '="%s";' % ast[1]
        else:
            print '""'
    else:
        print '""';'
        children = []
        for n, child in enumerate(ast[1:]):
            children.append(dotwrite(child))
        print ' %s -> {' % nodename,
        for name in children:
            print ' %s' % name,
```



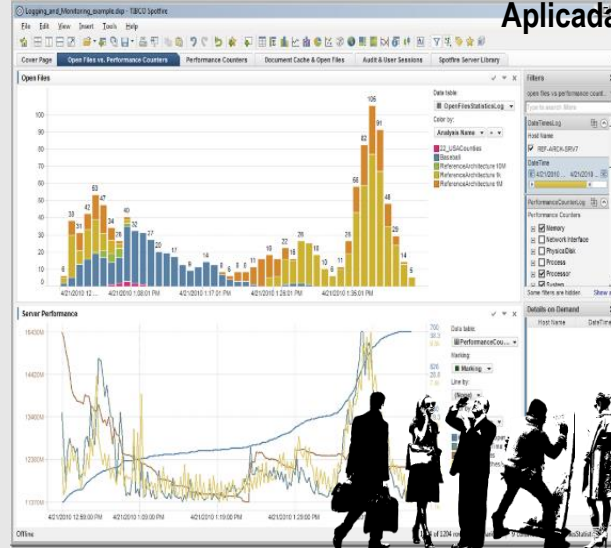
Comunicación



Programación

Estadística

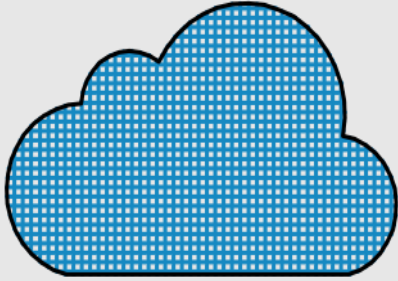
Aplicada



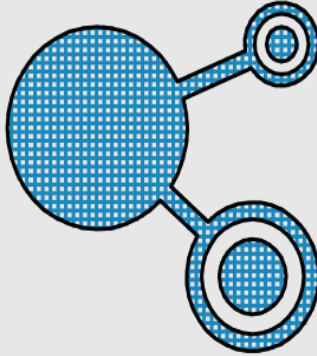
Visualización de Datos Experiencia en el Campo

## 03. PROCESO EN EL ANÁLISIS DE DATOS

### Cinco Pasos



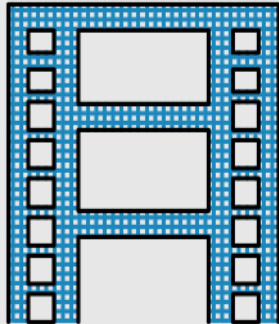
1/ Adquisición de Datos



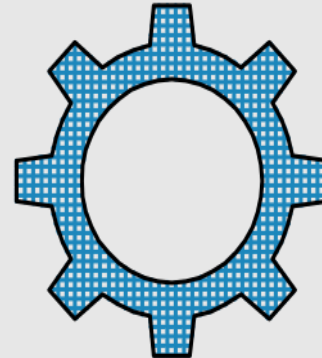
2/ Minería de Datos  
(Data Munging)



3/ Exploración de Datos



4/ Generar Modelos



5/ Validación de Modelos



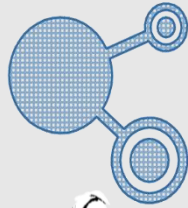
# 03. PROCESO EN EL ANÁLISIS DE DATOS

## Cinco Pasos

1 / Adquisición  
de Datos



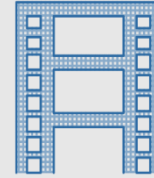
2 / Minería de Datos  
(Data Wrangling)



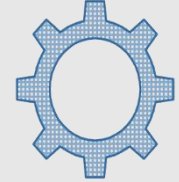
3 / Exploración  
de Datos



4 / Generar  
Modelos



5 / Validación  
de Modelos

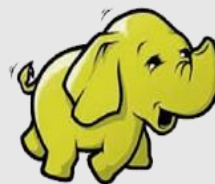


C i e n c i a d e D a t o s

## 04. ADQUISICIÓN DE DATOS

Paso Uno

UNO / Adquisición  
de Datos



SAP® Certified  
Integration with SAP HANA®



cloudera



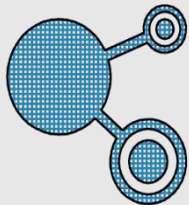
Databricks Certified Application

C i e n c i a d e D a t o s

# 05. MINERÍA DE DATOS

## Paso Dos

DOS / Minería de Datos  
(Data Wrangling)



## Data Wrangling o Data Mining

Minería de Datos: es un campo de la estadística y las ciencias de la computación referido al proceso que intenta descubrir patrones en grandes volúmenes de conjuntos de datos.



## munge

/mʌndʒ/ ☐ IPA ☐ Syllables

Word Origin

verb (used with object), verb (used without object), **munged, munging**. *Computer Slang*.

1. to manipulate (raw data), especially to convert (data) from one format to another:

"the munging of HTML content."

⌵ Hide transformations (no transformation steps added)

Transformations:

- Calculate and replace column
- Calculate and replace column
- Calculate new column
- Change column names
- Change data types
- Data function**
- Exclude columns
- Normalization
- Pivot
- Unpivot
- Help

Add...

Preview

Edit...

Remove

OK

Manage Relations...

C i e n c i a d e D a t o s

## ESTRATEGIA DE PRECIOS

Categoría	Productos con Stock	Productos con Matches						Productos con Matches con Stock			Ratio
		Totales		Sin Stock		Con Stock		Estoy más barato	Precio Igualado	Estoy más caro	
Tv y Video	388	291	75%	71	18%	220	57%	46 / 21%	81 / 37%	93 / 42%	1,03
Electrohogar	978	786	80%	237	24%	549	56%	257 / 47%	118 / 21%	174 / 32%	0,97
Cómputo	349	244	70%	107	31%	137	39%	40 / 29%	32 / 23%	65 / 47%	1,10
Más	263	42	16%	31	12%	11	4%	3 / 27%	4 / 36%	4 / 36%	1,15
Entretenimiento	770	514	67%	172	22%	342	44%	118 / 35%	103 / 30%	121 / 35%	1,01
Muebles	625	79	13%	30	5%	49	8%	19 / 39%	22 / 45%	8 / 16%	0,95
CALZADO Y RELOJES	1355	525	39%	143	11%	382	28%	38 / 56 / 15%	288 / 75%		1,29
Belleza y Accesorios	719	335	47%	77	11%	258	36%	49 / 19%	164 / 64%	45 / 17%	0,99
Ripley Home	3724	857	23%	199	5%	658	18%	162 / 25%	402 / 61%	94 / 14%	0,99
Celulares	381	244	64%	67	18%	177	46%	41 / 23%	38 / 21%	98 / 55%	1,16
Deporte	1452	499	34%	227	16%	272	19%	67 / 25%	133 / 49%	72 / 28%	1,06
Infantil	1938	946	49%	299	15%	647	33%	146 / 23%	197 / 30%	304 / 47%	1,13

# 05. EJEMPLO MINERÍA DE DATOS

Paso Dos

DOS / Minería de  
Datos  
(Data Wrangling)



PARA MAÑANA!  
???



CATEGOR Y	Q3-14	Q4-14	Q1-15	Q2-15	Q3-15	Q4-15	Q1-16	Q2-16	Well-Id
Total Cost in	2.00	2	3.00	4	2	2	3	6	Well 1
SMM	0.60	4	4.20	8	12	15	15	20	Well 1
Total Days	5.00	25	30.00	40	40	50	45	65	Well 1
Avg ABC per	6.00	50	40.00	60	100	130	120	150	Well 1
Operator	50.00	500	400.00	600	700	1000	500	600	Well 1
# of ABC	1.00	6	1.00	5	2	2	3	6	Well 2
Total ABC	0.60	4	4.20	8	12	15	15	20	Well 2
Total Cost in	5.00	25	30.00	40	40	50	45	65	Well 2

C i e n c i a   d e   D a t o s

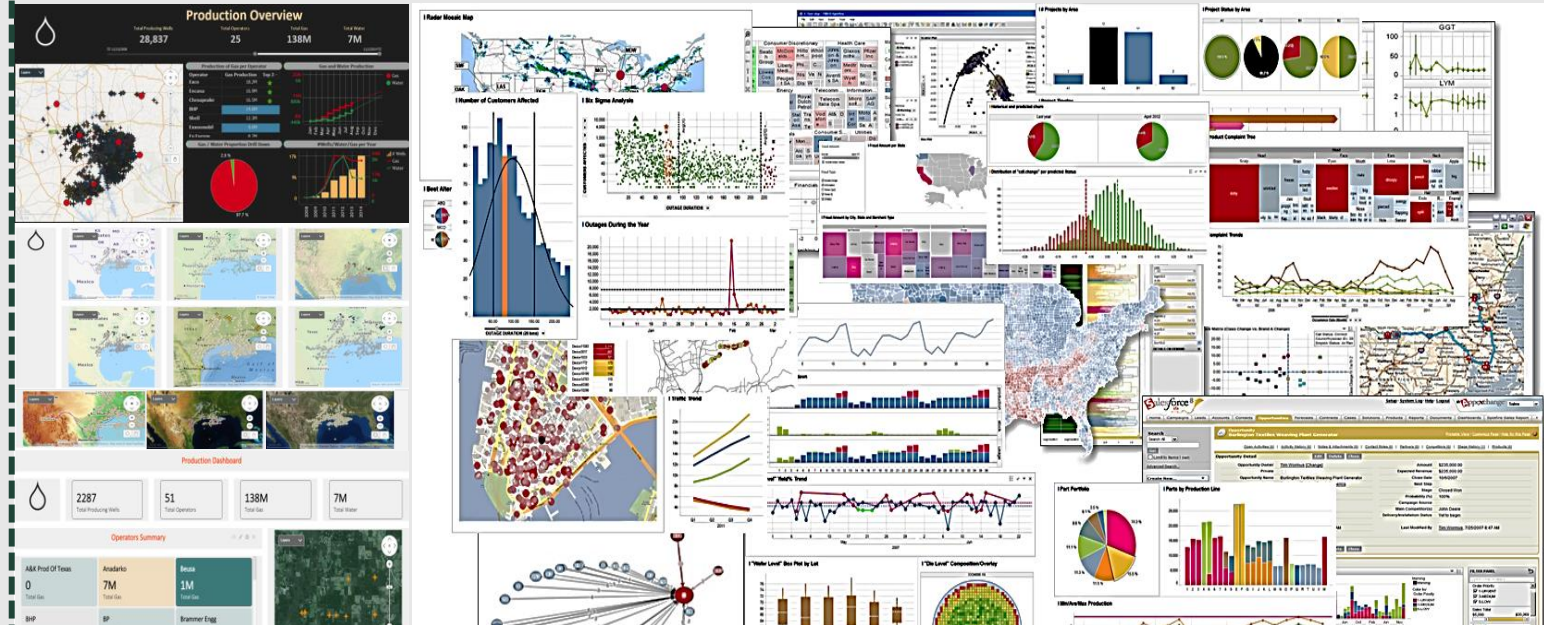


# 06. EXPLORACIÓN DE DATOS

Paso Tres



## TRES / Exploración de Datos



# C i e n c i a d e D a t o s



## 06. EXPLORACIÓN DE DATOS

Paso Tres



### ENGINEERING STATISTICS

H A N D B O O K

Welcome! The goal of this handbook is to help scientists and engineers incorporate statistical methods in their work as efficiently as possible.

Análisis exploratorio de datos (EDA) es un enfoque / filosofía para el análisis de datos que emplea una variedad de técnicas (principalmente GRÁFICA) para maximizar la penetración de datos; descubrir la trayectoria subyacente; extraer variables importantes; detectar valores atípicos y anomalías; probar los supuestos subyacentes; desarrollar modelos parsimoniosos; y determinar los valores óptimos de los factores.

John W. Tukey

### EXPLORATORY DATA ANALYSIS



“El mayor valor de una imagen es cuando se nos obliga a darnos cuenta de lo que nunca esperábamos ver”.

# 06. EXPLORACIÓN DE DATOS

Paso Tres

## TRES / Exploración de Datos



## EXPLORACIÓN DE DATOS

¿Cómo voy a visualizar los datos?



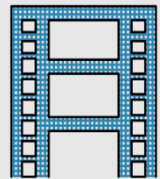
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
054	187	180	180	180	180	180	142	135	129	131	143	154	154	140	132	138	138	137	143	145
071	178	159	144	142	144	151	159	165	179	193	197	193	193	197	200	200	197	192	191	193
072	175	153	144	142	144	151	159	165	179	193	197	193	194	198	203	202	199	194	192	195
073	170	150	145	143	144	148	163	167	178	189	194	192	195	201	203	203	200	195	194	196
083	134	142	152	154	174	177	177	176	179	183	185	188	192	196	193	193	194	194	195	195
090	217	216	215	213	212	210	209	208	206	203	201	202	205	208	207	207	209	213	212	208
095	200	200	200	199	199	198	197	194	191	195	198	197	190	183	182	179	178	181	178	174
098	190	190	182	170	165	166	169	172	170	162	161	166	167	165	164	173	183	187	188	190
099	139	139	132	151	143	145	157	157	163	171	177	178	175	174	173	177	180	181	194	188
304	189	191	192	190	189	189	191	189	188	191	193	193	194	196	193	193	194	194	193	193
307	189	191	193	191	190	190	193	192	193	196	196	193	191	193	197	195	192	190	189	191
322	198	192	190	189	189	186	187	177	157	139	138	144	147	143	140	143	144	144	149	157
358	211	205	198	197	191	178	164	155	143	132	128	130	134	137	135	135	138	144	149	149
424	202	204	206	207	208	210	209	210	210	209	207	206	206	207	210	209	207	207	207	205
437	196	190	190	189	189	194	190	193	200	203	205	204	204	203	200	200	199	199	202	204
453	200	193	192	194	193	190	191	193	192	189	192	198	193	184	164	152	141	139	139	144
454	201	195	193	194	194	190	190	191	190	190	194	196	188	176	158	145	135	136	144	153
571	216	215	214	215	216	216	215	215	215	216	217	217	216	216	216	217	216	214	209	204
577	184	185	187	188	189	190	190	189	195	199	203	205	208	206	206	204	203	203	207	207

C i e n c i a   d e   D a t o s

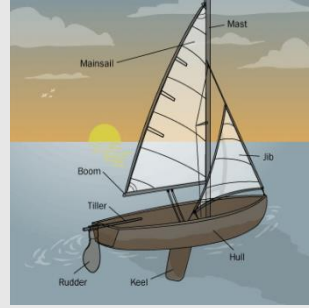
# 07. GENERAR MODELOS

Paso Cuatro

## CUATRO / Generar Modelos



DATA  
RELATIONSHIPS



	Numerical Vs Numerical
	Numerical Vs Categorical
	Categorical Vs Categorical



Lineal
Spearman R
Anova
Kruskal-Wallis
Chi-square



$$y = \beta_0 + \beta_1 x + \epsilon$$

C i e n c i a   d e   D a t o s

# El Científico de Datos



Is your Job sexy?

**Aldo Canales Bernal**  
**[acanales@abab.com.pe](mailto:acanales@abab.com.pe)**