

## Taller 8

Métodos Computacionales para Políticas Públicas - Urosario

Entrega: viernes 18-oct-2019 11:59 PM

**Julian Santiago Ramirez**	
jullians.ramirez@urosario.edu.co	

### Instrucciones:

- Guarde una copia de este *Jupyter/ Notebook* en su computador, idealmente en una carpeta destinada al material del curso.
- Modifique el nombre del archivo del *notebook*, agregando al final un guión inferior y su nombre y apellido, separados estos por otro guión inferior. Por ejemplo, *mi notebook se llamaría: mlopq\_taller8\_santiago\_matalana*
- Marque el *notebook* con su nombre y e-mail en el bloque verde arriba. Reemplace el texto ["Su nombre acá"] con su nombre y apellido. Similar para su email
- Desarrolle la totalidad del taller sobre este *notebook*, insertando las celdas que sea necesario debajo de cada pregunta. Haga buen uso de las celdas para código y de las celdas tipo *markdown* según el caso.
- Decuade salvar periódicamente sus avances.
- Cuando termine el taller:
  1. Descargúelo en PDF. Si tiene algún problema con la conversión, descargúelo en HTML.
  2. Suba todos los archivos a su repositorio en GitHub, en una carpeta destinada exclusivamente para este taller, antes de la fecha y hora límites.

### 1. [1 punto]

Usando expresiones regulares extraiga en una lista todos los números presentes en el siguiente objeto de Python:

```
ob1 = "JEFF BEZOS, the founder of Amazon, has reached a divorce settlement with his wife, MacKenzie. Mr Bezos will keep all the shares in the Washington Post and Blue Origin, a space-exploration firm, as well as 75% of the couple's Amazon stock. Mrs Bezos will retain a 4% stake in the tech giant, worth nearly $36bn, which is likely to make her the third-richest woman alive when the divorce is finalised."
```

```
In [1]: ### Importamos La librería que nos permite usar expresiones regulares ###
import re
```

```
In [3]: ### variable ob1
ob1 = "JEFF BEZOS, the founder of Amazon, has reached a divorce settlement with his wife, MacKenzie. Mr Bezos will keep all the shares in the Washington Post and Blue Origin, a space-exploration firm, as well as 75% of the couple's Amazon stock. Mrs Bezos will retain a 4% stake in the tech giant, worth nearly $36bn, which is likely to make her the third-richest woman alive when the divorce is finalised."

numeros=re.findall('(\\d+)',ob1)

print("Números en el texto: ",numeros)

Números en el texto:  ['75', '4', '36']
```

### 2. [1 punto]

Usando expresiones regulares ahora extraiga de *ob1* sólo los números que correspondan a porcentajes.

```
In [5]: porcentajes=re.findall('(\\d+)',ob1)

print("Números con porcentaje en el texto: ",porcentajes)

Números con porcentaje en el texto:  ['75%', '4%']
```

### 3. [2 puntos]

Usando expresiones regulares, escriba una función de Python que reciba una fecha en formato **Marzo 7, 2019** y retorne la fecha en formato **2019-07-03**

```
In [18]: def conversion_fecha(fecha):
def_fecha = []
meses = ['Enero', 'Febrero', 'Marzo', 'Abril', 'Mayo', 'Junio', 'Julio',
          'Agosto', 'Septiembre', 'Octubre', 'Noviembre', 'Diciembre']
### ahora de la fecha queremos obtener el día y el año
# Primero busquemos los numeros
numeros = re.search('([\\d+]),([\\d+])',fecha)
fec=numeros.groups()
def_fecha.append(fec[1])
## El primer dato de fec es el día y el segundo es el año, por lo tanto preguntamos si el día es igual a 23
if fec[0] == '23':
    dia_actualizado = '23'
    def_fecha.append(dia_actualizado)

## Hasta el momento hemos agregado el año y el día, falta el mes
mes=0
for i in range(0,len(meses)):
    if meses[i]=="Octubre":
        mes=i+1
        break
## Agregamos el mes a la fecha definitiva
def_fecha.append(mes)
total=str(def_fecha[0])+ '-' +str(def_fecha[1])+ '-' +str(def_fecha[2])
return total

## variable que guarda la fecha ###
fecha = 'Octubre 23, 2019'
nueva = conversion_fecha(fecha)

print("Antiguo formato: ", fecha)
print("Nuevo formato: ", nueva)

Antiguo formato: Octubre 23, 2019
Nuevo formato: 2019-23-10
```

### 4. [3 puntos]

o02 es un string que reúne una lista de clases en una universidad. Use expresiones regulares para extraer los códigos de cada una de las clases. Ejemplo: El código de la clase **COMPSCI 143 (Spring 2018): Machine Learning** es 143.

```
In [ ]: ## Guardamos ob2 ##
ob2 = "COMPSCI 270 (Spring 2018): Introduction to Artificial Intelligence. COMPSCI 590.2 (Fall 2018): Computational Microeconomics. Game Theory, Social Choice, and Mechanism Design. COMPSCI 223 (Spring 2018): Computational Microeconomics. COMPSCI 570 (Fall 2017): Artificial Intelligence. COMPSCI 570 (Fall 2017): Artificial Intelligence. COMPSCI 590.3 (Fall 2017) / 590.1 (Spring 2018): Ethics and AI. COMPSCI 590.2 (Spring 2017): Computation, Information, and Learning in Market Design. COMPSCI 590.4 (Spring 2016): Computational Microeconomics. Game Theory, Social Choice, and Mechanism Design. COMPSCI 290.4/590.4 (Spring 2015): Crowdsourcing Societal Tradeoffs. COMPSCI 590.4 (Spring 2014): Computational Microeconomics. Game Theory, Social Choice, and Mechanism Design. COMPSCI 570 (Fall 2014): Artificial Intelligence. COMPSCI 590.2 (Spring 2014): Computational Microeconomics. Game Theory, Social Choice, and Mechanism Design. COMPSCI 590.1 (Fall 2013): Computational Microeconomics. Game Theory, Social Choice, and Mechanism Design. COMPSCI 590.4 (Spring 2013): Crowdsourcing Societal Tradeoffs. COMPSCI 296.1 (Fall 2010): Linear and Integer Programming. COMPSCI 173 (Spring 2010): Computational Microeconomics. COMPSCI 196.1/296.1 (Fall 2009): Computational Microeconomics. Game Theory, Social Choice, and Mechanism Design. COMPSCI 170 (Spring 2009): Introduction to Artificial Intelligence. COMPSCI 270 (Fall 2008): Artificial Intelligence. COMPSCI 196/296.2 (Spring 2008): Linear and Integer Programming. COMPSCI 196.2 (Fall 2007): Introduction to Computational Economics. COMPSCI 296.3 (Spring 2007): Topics in Computational Economics. COMPSCI 296.2 (Fall 2006): Computational Game Theory and Mechanism Design."
```

```
In [ ]: ## Observamos que hay dos clases de código: el primero un numero normal y el segundo se puede ver de esta forma: 290.4/590.4
```

```
In [181]: ## Buscamos los números que acompañan La palabra COMPSCI

## Codigos tipo 1
codigo_tipo1=re.findall(' (\\d+)',ob2)

## Codigos tipo 2
codigo_tipo2=re.findall(' (\\d+)',ob2)

## Supongo que los codigos tipo 2 son das clases diferentes y por ende debo separarlos
nuevo_codigo=[]
```

```
for cod in codigo_tipo2:
    tupla=re.search('(\\d+)',(\\d+)',cod).groups()
    #agregamos ambos codigos
    nuevo_codigo.append(tupla[0])
    nuevo_codigo.append(tupla[1])

total=codigo_tipo1+nuevo_codigo
```

```
print("Códigos:", total)
#Code=re.findall(' (\\d+)',ob2)

Códigos: ['178', '598.2', '223', '578', '598.3', '590.1', '590.2', '598.4', '578', '590.4', '590.1', '173', '296.1', '296.1', '296.1', '178', '278', '196.2', '296.3', '296.2', '290.4', '590.4', '196.1', '196', '296.2']
```

### 5. [5 puntos]

o03 es un string que reúne una lista de publicaciones. Use expresiones regulares para extraer todos los *Journals* en los cuales el autor ha publicado. Ejemplo: El paper **Bail, CA. "The configuration of symbolic boundaries against immigrants in Europe."** *American Sociological Review* **73.1 (January 1, 2008): 37-59. Full Text fue publicado en el *Journal American Sociological Review***

```
In [ ]: ## Variable ob3
ob3 = "Bail, CA, Argyle, LP, Brown, TW, Bumpus, JP, Chen, H, Hunzaker, MBF, Lee, J, Mann, M, Merhout, F, and Volfovsky, A. "Exposure to opposing views on social media can increase political polarization." Proceedings of the National Academy of Sciences of the United States of America 115.37 (September 2018): 9216-9221. Full Text Open Access Copy. In: "Bail, CA, Merhout, F, and Din, P. "Using Internet search data to examine the relationship between anti-Muslim and pro-ISIS sentiment in U.S. counties." Science Advances 4.6 (June 6, 2018): eaas0948-null. Full Text Open Access Copy. In: "Bail, CA, Brown, TW, and Mann, M. "Channeling Hearts and Minds: Advocacy Organizations, Cognitive-Emotional Currents, and Public Conversation." American Sociological Review 82.6 (December 1, 2017): 1188-1213. Full Text. In: "Bail, CA. "Taming Big Data: Using App Technology to Study Organizational Behavior on Social Media." Sociological Methods and Research 46.2 (March 1, 2017): 189-217. Full Text. In: "McDonnell, TE, Bail, CA, and Tavorly, I. "A Theory of Resonance." Sociological Theory 35.1 (March 1, 2017): 1-14. Full Text. In: "Bail, CA. "Combining natural language processing and network analyses to examine how advocacy organizations stimulate conversation on social media." Proceedings of the National Academy of Sciences of the United States of America 113.42 (October 2016): 11823-11828. Full Text. In: "Bail, CA. "Emotional Feedback and the Viral Spread of Social Media Messages About Autism Spectrum Disorders." American journal of public health 106.7 (July 2016): 1173-1180. Full Text. In: "Bail, CA. "The public life of secrets: Deception, disclosure, and discursive framing in the policy process." Sociological Theory 33.2 (January 1, 2015): 97-124. Full Text. In: "Bail, CA. "The cultural environment: Measuring culture with big data." Theory and Society 43.3 (January 1, 2014): 465-524. Full Text."
```

```
In [57]: ## Guardamos todos los Journals en esta lista
Journals = []

# Buscamos en el texto frases: palabra + fecha
frases = re.findall(' (\\w+ +\\d+)',ob3)

for i in range(0,len(frases)):
    jour = re.search('(\\w+)',frases[i]).group()
    Journals.append(jour)

# Modificamos el primer Journal
Journals[0] = 'National Academy of Sciences of the United States of America'
borrar=[]
for i in range(0,len(Journals)):
    if Journals[i]!='Sociological Theory':
        borrar.append(i)
del Journals[borrar[0]]
del Journals[borrar[1]-1]

print("Journals: ",Journals)
```

```
Journals: ['National Academy of Sciences of the United States of America', 'Science Advances', 'American Sociological Review', 'Sociological Methods and Research', 'Proceedings of the National Academy of Sciences of the United States of America', 'American journal of public health', 'Theory and Society']
```

### 6. [10 puntos]

Vamos a hacer "scraping" a esta página: <https://archive.ics.uci.edu/ml/datasets.php>, que contiene un listado de 468 bases de datos que hacen parte del repositorio de la Universidad de California, Irvine.

Su tarea consiste en crear un "Pandas dataframe" que contenga 468 filas (una por base de datos) y las siguientes columnas:

- Nombre de la base de datos
- Link a la base de datos
- Tipo de datos
- Tipo de tarea a resolver (default task)
- Tipo de las variables
- Número de observaciones
- Número de variables
- Año
- Descripción de la base (Pista: Utilice la opción list view: <https://archive.ics.uci.edu/ml/datasets.php?format=&task=&att=&numAtt=&numIns=&type=&sort=name&view=list>)

```
In [94]: import requests
from text import BeautifulSoup
import pandas as pd
```

```
In [68]: # Info pag
link = requests.get('https://archive.ics.uci.edu/ml/datasets.php').text
l = BeautifulSoup(link, 'lxml')
```

```
In [85]: ### Buscamos primero todas las líneas que inician con <b>
lineas = l.find_all('b')
str_lineas=str(lineas)

# En esas líneas buscamos unas con una característica específica
especificas= re.findall('<(b>+href="datasets/(\\w+)', str_lineas)
```

```
# Creamos una lista con todos los nombres
nombres = []
for i in range(0,len(especificas)):
    nom = re.search('>(\\w+)',especificas[i]).group()
    nombres.append(nom)

# Ahora de los nombres los quitamos el </b>
for i in range(0,len(especificas)):
    nom = re.search('>(.+)',especificas[i]).group()
    nombres[i]=nom

#print("Variables: ", nombres)
```

```
In [84]: ## datos característicos
# líneas que inician con <p>
carac = l.find_all('p')
# Característica de los datos, falta el ultimo
caracteristicas=re.findall('(href="datasets/(\\w+)',str(carac))['href="datasets/RiceLeafDiseases">Rice Leaf Diseases</p>
<p>, <p class="normal">Multivariate</p>, <p class="normal">Classification</p>, <p class="normal">Integer</p>, <p class="normal">120</p>, <p class="normal">120</p>, <p class="normal">120</p>']
```

```
# Tipos de datos
## Guardare todo en una Lista de Listas
tipos=[]
for i in range(0,len(especificas)):
    carac_es = re.search('<p>'+caracteristicas[i]).group()
    tipo=re.findall('<(p>+</p>)',carac_es)
    for j in range(0,6):
        pal1=re.search('<(\\w+)',tipo[j]).group()
        pal2=re.search('>(.+)',especificas[i]).group()
        tipos[j].append(pal2)

# x0ab indica missing
# Combinados por NA
for j in tipos:
    for i in range(0,len(especificas)):
        if j[i] == 'x0ab':
            j[i]='NA'
```

```
In [98]: # LINKS
l1 = []
for i in range(0,len(especificas)):
    pal1 = re.search('<(\\w+)',especificas[i]).group()
    pal2 = re.search('<(\\w+)',especificas[i]).group()
    l1.append(pal2)

l2 = []
for i in range(0,len(l1)):
    pal = 'https://archive.ics.uci.edu/ml/'+l1[i]
    l2.append(pal)
```

```
# Links parte III
l3 = []
for i in range(0,len(l1)):
    html = requests.get(l2[i]).text
    s_html = BeautifulSoup(html, 'lxml')
    pal_a = s_html.find_all('a')
    html_nue = re.findall('<(a href=".,/machine-learning-databases/(\\w+)',str(pal_a))
    if html_nue == []:
        l3.append('NA')
    else:
        l3.append(html_nue[0])
```

```
In [92]: l_def = []
for i in range(0,len(especificas)):
    if j[i] == "NA":
        l_def.append('NA')
    else:
        pal1 = re.search('<(\\w+)',l3[i]).group()
        pal2 = 'https://archive.ics.uci.edu/ml/'+pal1
        l_def.append(pal2)
```

```
In [99]: tabla = ["Base de datos": nombres,
                "Tipo de datos": tipos[0],
                "Tipo de tarea a resolver": tipos[1],
                "Tipo de variables": tipos[2],
                "Número de obs": tipos[3],
                "Número de variables": tipos[4],
                "Año": tipos[5],
                "Links": l_def]
definitiva = pd.DataFrame(tabla)
```

```
# Ordenando el data frame
definitiva["nombres mayu"] = definitiva["Base de datos"].str.upper()
definitiva.sort_values(["nombres mayu"], axis=0, ascending=True, inplace=True)
# Borramos la variable creada del data frame
del definitiva["nombres mayu"]
```

```
In [108]: # Descripción de la base
link = requests.get('https://archive.ics.uci.edu/ml/datasets.php?format=&task=&att=&numAtt=&numIns=&type=&sort=name&view=list')
l1list = link.text
des = BeautifulSoup(l1list, 'lxml')
des_re=re.findall('<(b>+href="datasets/(\\w+)',des)
for i in range(0,len(des)):
    try:
        pal1 = re.search('<(\\w+)',des[i]).group()
        pal2 = re.search('<(\\w+)',des[i]).group()
        description.append(re.sub('<\\n>', '\\n', pal2))
    except AttributeError:
        description.append('NA')
```

```
definitiva['Descripción base datos'] = description

# Resultado
definitiva
```

```
Out[108]:
```

<pre># Descripción de la base link = requests.get("https://archive.ics.uci.edu/ml/datasets.php?format=&amp;task=&amp;att=&amp;area=&amp;numAtt=&amp;numIns=&amp;type=&amp;sort=&amp;nameid&amp;view=list") link = link.text  s_link = BeautifulSoup(link, "lxml")  des = re.findall('&lt;div class="description"&gt;[^\&lt;]*&lt;/div&gt;', s_link)  description = [] for i in range(0, len(des)):     try:         pal1 = re.search('&lt;div class="description"&gt;[^\&lt;]*&lt;/div&gt;', des[i]).group(1)         pal2 = re.search('&lt;div class="description"&gt;[^\&lt;]*&lt;/div&gt;', pal1).group(1)         description.append(re.sub('&lt;br&gt;', '\n', pal2))     except AttributeError:         description.append('NA')  definitiva["Descripción base datos"] = description  # Resultado definitiva</pre>									
	Base de datos	Tipo de datos	Tipo de tarea a resolver	Tipo de variables	Número de observaciones	Número de variables	Año	Links	Descripción base datos
	424 GHz Indoor Channel Measurements	Multivariate	Classification	Real	7840	5	2018	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Measurement of the S21 consists of 10 sweeps...
237	3D Road Network (North Juland, Denmark)	Sequential, Text	Regression, Clustering	Real	434874	4	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	3D road network with highly accurate elevation...
301	AAAI 2013 Accepted Papers	Multivariate	Clustering	NA	150	5	2014	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	This data set compromises the metadata for the...
294	AAAI 2014 Accepted Papers	Multivariate	Clustering	NA	399	6	2014	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	This data set compromises the metadata for the...
0	Abalone	Multivariate	Classification	Categorical, Integer, Real	4177	8	1995	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Predict the age of abalone from physical measurements...
170	Absciss Acid Signaling Network	Multivariate	Causal-Discovery	Integer	300	43	2008	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The objective is to delineate the set of people...
427	Absenteeism at work	Multivariate, Time-Series	Classification, Clustering	Integer, Real	740	21	2018	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The database was created with records of absenteeism...
260	Activities of Daily Living (ADLs) Recognition from a Single Chest-Mounted Accelerometer	Sequential, Time-Series	Classification, Clustering	NA	2747	NA	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	This dataset comprises of a timeline information regarding individual activities...
278	Activity Recognition from Single Chest-Mounted Accelerometer	Univariate, Sequential, Time-Series	Classification, Clustering	Real	NA	NA	2014	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The dataset collects data from a wearable accelerometer...
351	Activity Recognition system based on Multisensor Data	Multivariate, Sequential, Time-Series	Classification	Real	42240	6	2016	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	This dataset contains temporal data from a Wireless Sensor Network...
411	Activity recognition with healthy older people	Sequential	Classification	Real	75128	9	2016	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Sequential motion data from 14 healthy older people...
181	Acute Inflammations	Multivariate	Classification	Categorical, Integer	120	6	2009	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The data was created by a medical expert as a part of a research project...
1	Adult	Multivariate	Classification	Categorical, Integer	48842	14	1996	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Predict whether income exceeds \$50K/yr based on demographic attributes...
371	Air quality	Multivariate, Time-Series	Regression	Real	9358	15	2016	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Contains the responses of a gas multisensor device...
345	Air Quality	Multivariate, Time-Series	Regression	Real	9358	15	2016	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Contains the responses of a gas multisensor device...
279	Airfoil Self-Noise	Multivariate	Regression	Real	1503	6	2014	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	NASA data set obtained from a series of aerodynamic experiments...
476	Alcohol QCM Sensor Dataset	Multivariate	Classification, Regression, Clustering	Real	125	8	2019	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Five different QCM gas sensors are used, and features are extracted from the data...
208	Amazon Access Samples	Time-Series, Domain-Theory	Regression, Clustering, Causal-Discovery	NA	30000	20000	2011	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Amazon's InfoSec is getting smarter about the data it collects...
207	Amazon Commerce reviews self-reported	Multivariate, Text, Domain-Theory	Classification	Real	1500	10000	2011	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The dataset is used for authorship identification...
2	Annealing	Multivariate	Classification	Categorical, Integer, Real	798	38	NA	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Steel annealing data
3	Anonymous Microsoft Web Data	NA	Recommender-Systems	Categorical	37711	294	1998	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Log of anonymous users of www.microsoft.com...
390	Anuran Calls (MFCCs)	Multivariate	Classification, Clustering	Real	7195	22	2017	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Acoustic features extracted from syllables of anuran calls...
358	Appliances energy prediction	Multivariate, Time-Series	Regression	Real	19735	29	2017	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Experimental data used to create regression models...
405	APS Failure at Scania Trucks	Multivariate	Classification	Integer, Real	60000	171	2017	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The dataset's positive class consists of components that have failed...
164	Arcene	Multivariate	Classification	Real	900	10000	2008	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	ARCENE's task is to distinguish cancer versus non-cancerous cells...
4	Arrhythmia	Multivariate	Classification	Categorical, Integer, Real	452	279	1998	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Distinguish between the presence and absence of arrhythmia...
5	Artificial Characters	Multivariate	Classification	Categorical, Integer, Real	8000	7	1992	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Dataset artificially generated by using first-order Markov chains...
6	Audiology (Original)	Multivariate	Classification	Categorical	226	NA	1987	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Nominal audiology dataset from Bayle
7	Audiology (Standardized)	Multivariate	Classification	Categorical	226	69	1992	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Standardized version of the original audiology dataset...
457	Audit Data	Multivariate	Classification	Real	777	18	2018	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Exhaustive one year non-confidential data in the form of a CSV file...
...	...	...	...	...	...	...	...	...	...
248	User Knowledge Modeling	Multivariate	Classification, Clustering	Integer	403	5	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	It is the real dataset about the students' knowledge...
257	USPTO Algorithm Challenge, run by NASA-Harvard	Domain-Theory	Classification	Integer	306	5	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Data used for USPTO Algorithm Competition. Contains patent data...
234	Vertebral Column	Multivariate	Classification	Real	310	6	2011	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Data set containing values for six biomechanical parameters...
206	Vicon Physical Action Dataset	Time-Series	Classification	Real	3000	27	2011	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The Physical Action Data Set includes 10 normal activities...
438	Victorian Era	Text	Classification	NA	93600	1000	2018	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	To create the largest authorship competition dataset...
139	Volcanoes on Venus - JARTool experiment	Image	Classification	NA	NA	NA	NA	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The JARTool project was a pioneering effort to study Venus...
190	Wearable Computing: Classification of Body Postures	Sequential	Classification	Integer, Real	165632	18	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	A dataset with 5 classes (sitting-down, standing, walking, running, and lying-down)...
363	Website Phishing	Multivariate	Classification	Integer	1353	10	2016	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	...
262	Weight Lifting Exercises monitored with Inertial Sensors	Multivariate	Classification	Real	39242	152	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Six young healthy subjects were asked to perform various weight lifting exercises...
447	WESAD (Wearable Stress and Affect Detection)	Multivariate, Time-Series	Classification, Regression	Real	63000000	12	2018	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	WESAD (Wearable Stress and Affect Detection) contains data from a wearable sensor...
280	Wholesale Customers	Multivariate	Classification, Clustering	Integer	440	8	2014	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	The data set refers to clients of a wholesale distributor...
321	wik4ME	Multivariate	Regression, Clustering, Causal-Discovery	NA	913	53	2015	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Survey of faculty members from two Spanish universities...
273	Wilt	Multivariate	Classification	NA	4889	6	2014	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	High-resolution Remote Sensing data set (QuickBird)...
106	Wine	Multivariate	Classification	Integer, Real	178	13	1991	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Predicting the Chemical Composition of Wine
162	Wine Quality	Multivariate	Classification, Regression	Real	4898	12	2009	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Two datasets are included, related to red and white wine...
406	Wireless Indoor Localization	Multivariate	Classification	Real	2000	7	2017	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Collected in indoor space by observing signal strength...
487	WISDM Smartphone and Smartwatch Activity and Behavior	Multivariate	Classification	Integer	120	NA	2019	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Contains accelerometer and gyroscope time-series data...
234	Yacht Hydrodynamics	Multivariate	Regression	Real	308	7	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Data set, used to predict the hydrodynamic performance of a yacht...
196	YearPredictionMSD	Multivariate	Regression	Real	515345	90	2011	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Prediction of the release year of a song from its acoustic features...
107	Yeast	Multivariate	Classification	Real	1484	8	1996	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	Predicting the Cellular Localization of Yeast Proteins
215	YouTube Comedy Slam Preference Data	Text	Classification	NA	1138562	3	2012	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	This dataset provides user vote data on which videos are funny...
258	YouTube Multiview Video Games Dataset	Multivariate, Text	Classification, Clustering	Integer	120000	1000000	2013	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	This dataset contains about 12k instances, each with a video game...
364	YouTube Spam Collection	Text	Classification	NA	1956	5	2017	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	It is a public set of comments collected for YouTube videos...
398	Z-Alzadeh Sani	NA	Classification	Integer, Real	303	56	2017	<a href="https://archive.ics.uci.edu/ml/machine-learning...">https://archive.ics.uci.edu/ml/machine-learning...</a>	...