



## DEPARTAMENTO DE INGENIERÍA DE SISTEMAS Y COMPUTACIÓN

### Inteligencia de negocios 202220 – Laboratorio 4

PROFESORA: Haydemar Nuñez

| Nombres         | Apellidos        | Código    | Login       |
|-----------------|------------------|-----------|-------------|
| María Sofía     | Álvarez López    | 201729031 | ms.alvarezl |
| Brenda Catalina | Barahona Pinilla | 201812721 | bc.barahona |
| Alvaro Daniel   | Plata Márquez    | 201820098 | ad.plata    |

### Informe de laboratorio #4

El objetivo de este laboratorio es reforzar el conocimiento adquirido en la construcción de Pipelines (y, preferiblemente -como se realizó-, implementar transformaciones personalizadas), exportar un modelo machine learning utilizando pickle y construir un API para montar el modelo en producción y realizar predicciones mediante peticiones HTTP. Esto último fue logrado tanto local como remotamente, en una instancia EC2 de AWS, cuya IP estática es <http://3.228.160.169/>.

#### Escenarios de prueba para el API y análisis:

Para probar los escenarios, puede encontrar cada uno de los Json en la carpeta "Test". En los escenarios se pondrán los pantallazos de la prueba y el resultado que dio. Además de puntualizar el nombre de la prueba (del archivo).

En este proyecto, se plantearon 6 escenarios. Se buscó que fuese un espectro tal que arrojara predicciones coherentes, incoherentes y erróneas. Los escenarios previamente descritos se muestran a continuación:

- **R<sup>2</sup> coherente:**  
Nombre del archivo: R2Coherente1

|  |               |
|--|---------------|
| POST /r2 Get R2  |               |
| Parameters   |               |
| No parameters  |               |
| Request body <small>required</small>   |               |
| <pre>{   "data": {     "data": [       {         "unnamed_0": "s",         "adult_mortality": 16.0,         "infant_deaths": 11.0,         "alcohol": 5.88,         "percentage_expenditure": 547.210141,         "hepatitis_B": 98.0,         "measles": 6071,         "hmi": 26.8,         "under_five_deaths": 12.0,         "polio": 99.0,         "total_expenditure": 4.11,         "diphtheria": 99.0,         "hiv_aids": 0.3,         "gdp": 4212.549200,         "population": 66881867.0,         "thinness 10-19 years": 8.3       }     ]   } }</pre> |               |
| Code   | Details       |
| 200  | Response body |
| <pre>{   "r^2": 0.9242702975164372 }</pre>   |               |

Para este escenario calculamos en Jupyter el  $r^2$  y este es el mismo que el  $r^2$  que obtuvimos con la API. Esta fue la petición realizada en el módulo de FASTAPI. Los resultados de la petición en Postman se presentan a continuación:

POST

http://3.228.160.169/r2

Params

Authorization

Headers (8)

Body

Pre-request Script

Tests

Settings

none

form-data

x-www-form-urlencoded

raw

binary

GraphQL

JSON

```

1  {
2    "data": {
3      "data": [
4        {
5          "unnamed_0": 1,
6          "adult_mortality": 16.0,
7          "infant_deaths": 11.0,
8          "alcohol": 5.88,
9          "percentage_expenditure": 547.210141,
10         "hepatitis_B": 98.0,
11         "measles": 6071,
12         "bmi": 26.8,
13         "under_five_deaths": 12.0,
14         "polio": 99.0,
15         "total_expenditure": 4.11,
16         "diphtheria": 99.0,
17         "hiv_aids": 0.3,
18         "gdp": 4212.549200,
19         "population": 66881867.0,
20         "thinness_10_19_years": 8.3,
21         "thinness_5_9_years": 8.5,
22         "income_composition_of_resources": 0.706,
23         "schooling": 13
24       },
25       {
26         "unnamed_0": 21,
27         "adult_mortality": 393.0,
28         "infant_deaths": 2,
29         "alcohol": 5.01,
30         "percentage_expenditure": 426.785566,
31         "hepatitis_B": 94.0

```

Body

Cookies

Headers (5)

Test Results

Pretty

Raw

Preview

Visualize

JSON

```

1  {
2    "r^2": 0.9242702975164372
3  }

```

La predicción tiene sentido puesto que los datos utilizados fueron seleccionados cuidadosamente tal que siguieran las mismas tendencias que algunos de los suministrados en el archivo de entrenamiento con el que se entrenó el modelo.

## Nombre del archivo: PrediccionCoherente1

|   |  |
|---|--|
| POST /predict Make Predictions  |  |
| Parameters  |  |
| No parameters   |  |
| Request body required   |  |
| <pre>{  "data": [    {      "unnamed_0": "s",      "adult_mortality": 16,      "infant_deaths": 11,      "alcohol": 5.88,      "percentage_expenditure": 547.218141,      "hepatitis_B": 98.0,      "measles": 6871,      "bmi": 26.8,      "under_five_deaths": 12,      "polio": 99.0,      "total_expenditure": 4.11,      "diphtheria": 99.0,      "hiv_aids": 0.3,      "gdp": 4212.549288,      "population": 66881867.0,      "thinness_10_19_years": 8.3,      "thinness_5_9_years": 8.5    }  ]}</pre> |  |
| Code  | Details  |
| 200   | <div>Response body</div> <pre>{  "predict": "[74.94561162835979, 61.731706271446886]"}</pre> |

En esta predicción se puede ver que se ponen a prueba dos casos, y podemos ver que el life expectancy que resulta es muy cercano al que es en realidad. Para el primer caso, el valor de la variable de interés es de 74.94, cuando en realidad debió ser de 73.7. Para la segunda, la predicción dio 61.73, siendo en realizada 59.2. Esta fue la petición realizada en el módulo de FASTAPI. Los resultados de la petición en Postman se presentan a continuación:

POST http://3.228.160.169/predict

Params Authorization Headers (8) **Body** Pre-request Script Tests Settings

none form-data x-www-form-urlencoded raw binary GraphQL JSON

```
1 {
2   "data": [
3     {
4       "unnamed_0": 200,
5       "adult_mortality": 16,
6       "infant_deaths": 11,
7       "alcohol": 5.88,
8       "percentage_expenditure": 547.210141,
9       "hepatitis_B": 98.0,
10      "measles": 6071,
11      "bmi": 26.8,
12      "under_five_deaths": 12,
13      "polio": 99.0,
14      "total_expenditure": 4.11,
15      "diphtheria": 99.0,
16      "hiv_aids": 0.3,
17      "gdp": 4212.549200,
18      "population": 66881867.0,
19      "thinness_10_19_years": 8.3,
20      "thinness_5_9_years": 8.5,
21      "income_composition_of_resources": 0.706,
22      "schooling": 13.0
23    },
24    {
25      "unnamed_0": 100,
26      "adult_mortality": 393.0,
27      "infant_deaths": 2,
28      "alcohol": 5.01,
29      "percentage_expenditure": 426.785566,
30      "hepatitis_B": 94,
31      "measles": 184.

```

Body Cookies Headers (5) Test Results

Pretty Raw Preview Visualize JSON

```
1 {
2   "predict": "[74.94561162835979, 61.731706271446086]"
3 }
```

Note que los resultados son iguales a los obtenidos con FastAPI.

- **Predicción incoherente:**

**Nombre del archivo:** PrediccionIncoherente1

| Parameters   |   |
|--|---|
| No parameters  |   |
| Request body <span>required</span>   |   |
| <pre>{   "data": [     {       "unnamed_0": 500,       "adult_mortality": 160,       "infant_deaths": 1100,       "alcohol": 35.88,       "percentage_expenditure": 947.210141,       "hepatitis_B": 9.0,       "measles": 0,       "bmi": 68.6,       "under_five_deaths": 12,       "polio": 99.0,       "total_expenditure": 4.11,       "diphtheria": 99.0,       "hiv_aids": 0.3,       "gdp": 2284.549200,       "population": 600.0,       "thinness_10_19_years": 100,       "thinness_5_9_years": 100     }   ] }</pre> |   |
| Code   | Details   |
| 200  | Response body <pre>{   "predict": "[8680.796497568856]" }</pre> |

Podemos ver que varios de estos datos no tienen sentido, entre estos la cantidad de población y la cantidad de muertes de niños no es coherente, de acuerdo con las suposiciones e información del negocio. Adicional a esto la delgadez en niños de 10 a 19 años y 5 a 9 años es alta. Sin embargo, se predice que la esperanza de vida es 8680.8 años, lo cual es imposible para cualquier país del mundo.

La realidad sobre esta predicción tan elevada recae en que el income composition of resources, que de acuerdo con el análisis cualitativo es la de mayor peso, está teniendo un valor de 693 (ver json), cuando debería tener un valor entre 0 y 1. Esto termina entonces generando un valor de expectativa de vida supremamente grande.

Note que los resultados son los mismos realizando la petición en Postman:

POST ⌵ http://3.228.160.169/predict

Params Authorization Headers (8) Body ● Pre-request Script Test

● none ● form-data ● x-www-form-urlencoded ● raw ● binary ● Gr

```
1  {
2    "data": [
3      {
4        "unnamed_0": 500,
5        "adult_mortality": 160,
6        "infant_deaths": 1100,
7        "alcohol": 35.88,
8        "percentage_expenditure": 947.210141,
9        "hepatitis_B": 9.0,
10       "measles": 0,
11       "bmi": 68.6,
12       "under_five_deaths": 12,
13       "polio": 99.0,
14       "total_expenditure": 4.11,
15       "diphtheria": 99.0,
16       "hiv_aids": 0.3,
17       "gdp": 2284.549200,
18       "population": 600.0,
19       "thinness_10_19_years": 100,
20       "thinness_5_9_years": 100,
21       "income_composition_of_resources": 693,
22       "schooling": 14.6
23     }
24   ]
25 }
```

Body Cookies Headers (5) Test Results

Pretty

Raw

Preview

Visualize

JSON ⌵



```
1  {
2    "predict": "[8680.796497568856]"
3  }
```

## Nombre del archivo: PrediccionIncoherenteMedia2

POST /predict Make Predictions

Parameters

No parameters

Request body required

```
{  "data": [    {      "unnamed_0": 78,      "adult_mortality": 1625,      "infant_deaths": 31,      "alcohol": 381,      "percentage_expenditure": 2945742256,      "hepatitis_B": 655,      "measles": 2433,      "bmi": 377,      "under_five_deaths": 43,      "polio": 820,      "total_expenditure": 449,      "diphtheria": 819,      "hiv_aids": 16,      "gdp": 1479315375,      "population": 114609342,      "thinness_10_19_years": 47,      "thinness_5_9_years": 48    }  ]}
```

Server response

| Code | Details   |
|------|---|
| 200  | <div>Response body<pre>{  "predict": "[6964.773366704773]"}</pre></div> |

Decidimos realizar una predicción usando los valores promedio de cada una de las variables del conjunto de entrenamiento. Note que los valores extremos (outliers) pueden llegar a afectar mucho la predicción. En particular, de nuevo, es income composition of resources la que más despista al modelo. Al tener una media de 597, cuando realmente debería estar en 1, hace que la predicción sobre la expectativa de vida sea mucho mayor. Una forma de mitigar este problema sería que, por ejemplo, a todos los países que tengan ese valor tan elevado, se les remplace esta variable por el valor que se tiene, dividido entre 1000. Esto será implementado en las mejoras al pipeline.

Con esto, se ve que la media termina representando datos incoherentes para un país. Esto también se ve reflejado en la variable de interés, ya que la esperanza de vida es de 6964.7773, lo cual no tiene sentido.

Note que los resultados en Postman son iguales:



POST

▼

http://3.228.160.169/predict

Params

Authorization

Headers (8)

Body ●

Pre-request Script

Test

● none

● form-data

● x-www-form-urlencoded

● raw

● binary

● Gr

1

{

2

"data": [

3

{

4

"unnamed\_0": 78,

5

"adult\_mortality": 1625,

6

"infant\_deaths": 31,

7

"alcohol": 381,

8

"percentage\_expenditure": 2945742256,

9

"hepatitis\_B": 655,

10

"measles": 2433,

11

"bmi": 377,

12

"under\_five\_deaths": 43,

13

"polio": 820,

14

"total\_expenditure": 449,

15

"diphtheria": 819,

16

"hiv\_aids": 16,

17

"gdp": 1479315375,

18

"population": 114609342,

19

"thinness\_10\_19\_years": 47,

20

"thinness\_5\_9\_years": 48,

21

"income\_composition\_of\_resources": 547,

22

"schooling": 115

23

}

24

]

25

}

Body

Cookies

Headers (5)

Test Results

Pretty

Raw

Preview

Visualize

JSON ▼

≡ ↻

1

{

2

"predict": "[6964.773366704773]"

3

}

### Nombre del archivo: R2Incoherencia1

En el último escenario que probamos, quisimos ver qué sucedía con el  $r^2$  si tomábamos datos incoherentes (por ejemplo, los usados para los anteriores dos escenarios). A continuación, el resultado:

**POST** /r2 GetR2

**Parameters**

No parameters

**Request body** required

```
{
  "data": {
    "data": [
      {
        "unnamed_0": 21,
        "adult_mortality": 151.0,
        "infant_deaths": 0.0,
        "alcohol": 2,
        "percentage_expenditure": 569.2953589,
        "hepatitis_B": 11,
        "measles": 0,
        "tbi": 76.5,
        "under_five_deaths": 0.0,
        "polio": 75.6,
        "total_expenditure": 5,
        "diphtheria": 10,
        "hiv_aids": 80,
        "gdp": 3000,
        "population": 200,
        "thinness 10-19 years": 3.
      }
    ]
  }
}
```

| Code | Details  |
|------|--|
| 200  | <b>Response body</b> <pre>{   "r^2": -.4956438.933716847 }</pre> |

Podemos ver que, en este escenario, el  $r^2$  toma un valor negativo que no tiene coherencia, esto se debe a que la predicción hecha y esperada es muy diferente. Este valor no tiene sentido este puede tener un valor máximo de 1 y, para representar el ajuste de los datos al modelo, debe ser positivo.

La razón por la cual esto ocurre es que el valor promedio del income composition of resources para ambos valores enviados es dos órdenes de magnitud mayor que lo esperado (está alrededor de 500, cuando debería ser 1 – como máximo -). Para mejorar este comportamiento, en el pipeline conviene dividir todos aquellos valores mayores a 1 entre 1000. Esto dará un resultado más certero con respecto a la variable de expectativa de vida.

Note que los resultados en postman son los mismos que los obtenidos previamente:

POST    http://3.228.160.169/r2

Params    Authorization    Headers (8)    **Body**    Pre-request Script    Test

☐ none    ☐ form-data    ☐ x-www-form-urlencoded    ☒ raw    ☐ binary    ☐ Gr

```
1  {
2    "data": {
3      "data": [
4        {
5          "unnamed_0": 21,
6          "adult_mortality": 151.0,
7          "infant_deaths": 0.0,
8          "alcohol": 2,
9          "percentage_expenditure": 569.2953509,
10         "hepatitis_B": 11,
11         "measles": 0,
12         "bmi": 76.5,
13         "under_five_deaths": 0.0,
14         "polio": 75.6,
15         "total_expenditure": 5,
16         "diphtheria": 10,
17         "hiv_aids": 80,
18         "gdp": 3000,
19         "population": 200,
20         "thinness_10_19_years": 3,
21         "thinness_5_9_years": 3,
22         "income_composition_of_resources": 856,
23         "schooling": 25
24       },
25     }
```

**Body**    Cookies    Headers (5)    Test Results

Pretty    Raw    Preview    Visualize    JSON   

```
1  {
2    "r^2": -4956438.933716047
3  }
```

- **Hay fallas**

Nombre del archivo: FallaUnnamed

| POST   | /predict                    | Make Predictions |
|--|-----------------------------|------------------|
| Parameters   |                             |                  |
| No parameters  |                             |                  |
| Request body <small>required</small>   |                             |                  |
| <pre>{   "data": [     {       "adult_mortality": 151.0,       "infant_deaths": 0,       "alcohol": 1.8,       "percentage_expenditure": 423.2953589,       "hepatitis_B": 9.0,       "measles": 0,       "bmi": 68.6,       "under_five_deaths": 0.0,       "polio": 91.0,       "total_expenditure": 4.87,       "diphtheria": 9.0,       "hiv_aids": 0.1,       "gdp": 2284.37858,       "population": 146.0,       "thinness_10_19_years": 0.1,       "thinness_5_9_years": 0.1,       "income_correlation_of_required": 693     }   ] }</pre> |                             |                  |
| Code   | Details                     |                  |
| 422  | Error: Unprocessable Entity |                  |
| Response body  |                             |                  |
| <pre>{   "detail": [     {       "loc": [         "body",         "data",         0,         "unnamed_0"       ],       "msg": "field required",       "type": "value_error.missing"     }   ] }</pre>   |                             |                  |

Podemos ver que ocurre un error, al realizar una prueba sin la columna "Unnamed". Esto ocurre porque en el modelo teníamos en los datos esta columna, que, en el pipeline eliminamos para hacer el modelo final de predicción. Con lo anterior generamos un problema ya que requerimos de esta columna para que el pipeline y por consiguiente el modelo funcione. Debemos definir entonces, en la clase DataModel.py, un atributo correspondiente a Unnamed\_0 (i.e. la columna Unnamed: 0 con la que venía el archivo de entrenamiento), y enviar el JSON con este campo. En realidad, como Unnamed\_0 no es un parámetro del modelo, no tiene sentido incluirlo. Por lo tanto, quitar este parámetro sería parte de la ejecución del bono (como se verá más adelante).

Note que los resultados obtenidos en Postman son idénticos:

POST http://3.228.160.169/predict

Params Authorization Headers (8) **Body** Pre-request Script Test

● none ● form-data ● x-www-form-urlencoded ● raw ● binary ● Gr

```
1 {
2   "data": [
3     {
4       "adult_mortality": 151.0,
5       "infant_deaths": 0,
6       "alcohol": 1.8,
7       "percentage_expenditure": 423.2953509,
8       "hepatitis_B": 9.0,
9       "measles": 0,
10      "bmi": 68.6,
11      "under_five_deaths": 0.0,
12      "polio": 91.0,
13      "total_expenditure": 4.87,
14      "diphtheria": 9.0,
15      "hiv_aids": 0.1,
16      "gdp": 2284.37858,
17      "population": 146.0,
18      "thinness_10_19_years": 0.1,
19      "thinness_5_9_years": 0.1,
20      "income_composition_of_resources": 693,
21      "schooling": 14.6
22    },
23    {
24      "unnamed_0": 22,
25      "adult_mortality": 153.0,
```

Body Cookies Headers (5) Test Results

Pretty Raw Preview Visualize JSON

```
1 {
2   "detail": [
3     {
4       "loc": [
5         "body",
6         "data",
7         0,
8         "unnamed_0"
9       ],
10      "msg": "field required",
11      "type": "value_error.missing"
12    }
13  ]
14 }
```

Como se puede apreciar, se obtiene un error de tipo: `value_error.missing`, correspondiente a la carencia del campo `unnamed_0` en el primer elemento del json.

### **Estrategia por desarrollar sobre el software para mitigar incoherencias en el resultado y fallas en el sistema:**

Una posible estrategia es mejorar el pipeline, en esta forma podríamos tener en cuenta la posibilidad de que no sea agregado la columna "Unnamed", por lo que, el modelo debería correr exitosamente con o sin esta columna.

Por otro lado, podríamos dar un mejor manejo a los valores nulos, nosotros los tratamos con la media, pero se podría probar alternativas como la eliminación de estos datos o la imputación con otro valor significativo que no sea la media. En este caso, proponemos dividir todos aquellos valores que sean mayores que 1, entre 1000, para evitar este error.

El intento de ambas mejoras en la implementación puede encontrarlo en el repositorio <https://github.com/sofiaalvarezlopez/BI-Lab4-Mejoras> .

### **Conclusiones:**

En este laboratorio, pudo implementarse satisfactoriamente un pipeline con transformaciones personalizadas, que fue desplegada tanto local como remotamente (en un servidor de AWS).

Hay que tener en cuenta que los datos solo funcionan para países vecinos a los Alpes y que sean de años cercanos, para que tengan estadísticas similares, ya que, debido a factores macro y micro, los índices o variables de estos países pudieron haber cambiado drásticamente, por lo que el modelo dejaría de predecir adecuadamente la expectativa de vida. Es por esta razón que los datos incoherentes (como, por ejemplo, el valor de `income composition of resources`) afectan tanto (y negativamente) la predicción del modelo.

### **Bono:**

#### **1. Construir transformaciones personalizadas e incluirlas en el pipeline y garantizar que el proceso completo es correcto.**

El pipeline que usamos para la realización de este laboratorio tiene transformaciones personalizadas, estas clases se pueden ver en el documento `clases.py`, en la carpeta Notebook.

Es importante notar que, al comienzo, las transformaciones personalizadas se implementaban directamente en el notebook. No obstante, por errores de serialización, estas fueron puestas en el archivo clases.py.

**2. Desplegar la API en un servidor gratuito como Heroku para que pueda prestar servicio a cualquiera haciendo uso de una URL.**

Desplegamos la API en AWS, en este logramos que la IP sea estática. Puede acceder en : <http://3.228.160.169/docs>

**3. Implementar la estrategia para mitigación de errores identificados en los escenarios y documentados por ustedes en el documento de entrega.**

Remítase al repositorio: <https://github.com/sofiaalvarezlopez/BI-Lab4-Mejoras>