# EXTRACCIÓN DE DATOS EN *RECIBOS*

- ***Pilar Madariaga Lasala***
- Estudiante Máster Inteligencia Artificial en Unir
- Pasantia en Grupo Bancolombia

## OBJETIVO
*Extracción de características*

- *Servicio Rest en Python*
- **Input**: Imagen en base64
- **Output**:
  - Fecha compra
  - Establecimiento
  - Listado de producto y precio

POST ▼ http://demo-env.eba-r8em2rmc.us-east-1.elasticbeanstalk.com/extractFeatures

Params   Authorization   Headers (10)   Body ●   Pre-request Script   Tests   Settings

○ none   ○ form-data   ○ x-www-form-urlencoded   ● raw   ○ binary   ○ GraphQL   JSON ▼

```
1  {
2      "filename":"ticket.jpg",
3      "image": "iVBORw0KGgoAAAANSUhEUgAAAmkAAARdCAIAAACFKEIEAADGM01EQVR4nOydr7fkSJbfVXnmHDOzN0uUScyqoYkySc0iMyM
       fR4QQQQggZzOjWBSCEEEIeDGonIYQQEga1kxBCCAmD2kkIIYSEQe0khBBCCwqB2EkIIIWFQOwkhhJAwqJ2EEEJIGNROggghJAxqJyGE
       YSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kxBCCAmD2kkIIYSEQe0khBBCCwqB2EkIIIWFQOwkhhJAwqJ2EEEJIGNR
       IRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kxBCCAmD2kkIIYSEQe0khBBCCwqB2EkIIIWFQOwkhhJ.
       EEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kxBCCAmD2kkIIYSEQe0khBBCCwqB2E
       MaichhBASBrWTEEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kxBCCAmD2kkIIYSE
       BCSBjUTkIIISQMaichhBASBrWTEEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kxB
       DsJIYSQMKidhBBCSBjUTkIIISQMaichhBASBrWTEEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDG
       QsKgdhJCCCFhUDv/xrt3725dhLuA9QBYD4D1AFgPgPGUA3r17926/39
       +6GIQQQsgjQbuTEEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kxBCCAmD2kkIIYS
       BBCSBjUTkIIISQMaichhBASBrWTEEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkDGonIYQQEga1kx
       UDsJIYSQMKidhBBCSBjUTkIIISQMaichhBASBrWTEEIICYPaSQghhIRB7SSEEEELCoHYSQgghYVA7CSGEkDConYQQQkgY1E5CCCEkD
```

Body   Cookies (2)   Headers (7)   Test Results

Pretty   Raw   Preview   Visualize   JSON ▼

```
1  {
2      "date": "2018-12-02",
3      "items": {
4          "Bacon Cheese Fries ": "4.89",
5          "DBL ShackBurger ": "8.69",
6          "Shake Vanilla Shake ": "5.49"
7      },
8      "shop": "SHAKE SHACK"
9  }
```

# *Primeros pasos: pytesseract*

```python
1  import re
2  import cv2
3  import pytesseract
4  from pytesseract import Output
5  from google.colab.patches import cv2_imshow
6
7
8  img = cv2.imread('ticket.png')
9
10 d = pytesseract.image_to_data(img, output_type=Output.DICT)
11 keys = list(d.keys())
12
13 date_pattern = '^(0[1-9]|[12][0-9]|3[01])/(0[1-9]|1[012])/(19|20)\d\d$'
14
15 n_boxes = len(d['text'])
16 for i in range(n_boxes):
17     if int(d['conf'][i]) > 60:
18         if re.match(date_pattern, d['text'][i]):
19             (x, y, w, h) = (d['left'][i], d['top'][i], d['width'][i], d['height'][i])
20             img = cv2.rectangle(img, (x, y), (x + w, y + h), (0, 255, 0), 2)
21
22 cv2_imshow(img)
```

```
            SHAKE SHACK
       3790 Las Vegas Blvd South
Host: Lisa                    12/02/2018
166  WALTER                   11:44 AM
                              60042

DBL ShackBurger                   8.69
Bacon Cheese Fries                4.89
Shake                             5.49
  Vanilla Shake

Subtotal                         19.07
Tax                               1.57

To Stay Total          20.64

MasterCard #XXXXXXXXXXXX2825      20.64
   Auth:NGJJ58
```

```python
# regex
price_pattern = '\d+\.\d+$'
date_pattern = '^(0[1-9]|[12][0-9]|3[01])/(0[1-9]|1[012])/(19|20)\d\d$'
total_pattern = '(.*otal.*|.*Tax.*)$'

n_boxes = len(d['text'])
for i in range(n_boxes):
    if int(d['conf'][i]) > 60:
        # productos y precios
        if re.match(price_pattern, d['text'][i]) and d['block_num'][i]==2  and not re.match(total_pattern, d['text'][i-1]):
            k = lineaElementos(d,i)
            for j in range(0,k):
                (x, y, w, h) = (d['left'][i-k], d['top'][i-k], d['width'][i-k], d['height'][i-k])
                img = cv2.rectangle(img, (x, y), (x + w, y + h), (0, 255, 0), 2)

cv2_imshow(img)
```

```
            SHAKE SHACK
       3790 Las Vegas Blvd South
Host: Lisa                    12/02/2018
166  WALTER                   11:44 AM
                              60042

DBL ShackBurger                   8.69
Bacon Cheese Fries                4.89
Shake                             5.49
  Vanilla Shake

Subtotal                         19.07
Tax                               1.57

To Stay Total          20.64

MasterCard #XXXXXXXXXXXX2825      20.64
   Auth:NGJJ58
```

```
'establecimiento': ' SHAKE SHACK', 'fecha_compra': '12/02/2018', 'productos': {' DBL ShackBurger ': '8.69', ' Bacon Cheese Fries ': '4.89', ' Shake ': '5.49'}}
```

# DATASET
## *Recibos Compras (I)*



Duplicados

Intensidad

Formato Fecha

*SROIE 2019*  https://drive.google.com/drive/folders/1ShItNWXyiY1tFDM5W02bceHuJjyeeJl2

# DATASET
## Recibos Compras (II)



Formato Línea Producto

Varios productos en la misma línea

Sin Alinear

Sin fecha

**SROIE 2019** https://drive.google.com/drive/folders/1ShItNWXyiY1tFDM5W02bceHuJjyeeJl2

# *Herramientas*
## *AWS + Regex + libs*

- *Amazon Textract*
- *Amazon Comprehend*
- *Elastic Beanstalk*
- Dateparser
- Spacy
- Expresiones regulares



https://www.python.org/          https://dateparser.readthedocs.io/          https://spacy.io/

# Lógica del servicio

- 1. Decodificar la imagen en directorio /tmp

- 2. Procesar imagen con 2 servicios de Textract y 1 servicio de Comprehed. Salvar resultados en directorios /text /tables /entidades

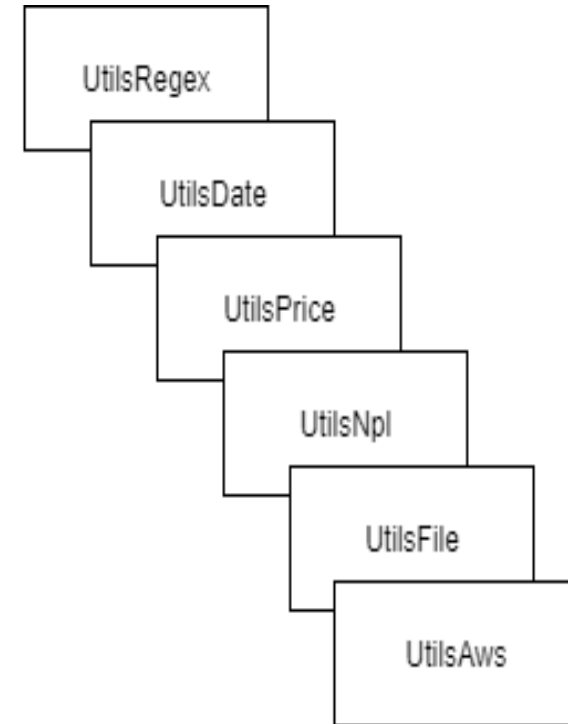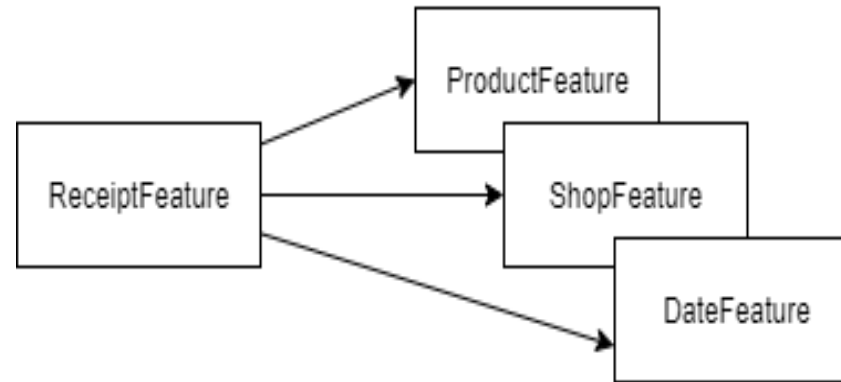- 3. Procesar los ficheros anteriores y extraer Features
  - Extraer Establecimiento (ShopFeature)
  - Extraer Listado de productos y precio (ProductFeature)
  - Extraer Fecha (DateFeature)
  - (Doble validación con expresiones regulares)

- 4. Eliminar ficheros e imagen de los directorios temporales

```python
def getFeatures(self, file):

    print("Procesando el fichero: {}".format(file))
    text = self.saveTextDetectedByTextract(self.tmpDirectory,file)
    entities=self.saveEntitiesDetectedByComprehed(self.entitiesDirectory,file,text)
    table = self.saveTableDetectedByTextract(self.tmpDirectory,file)


    text = self.readText(file)
    table = self.readTable(file)
    entities = self.readEntities(file)


    features = {}
    extractor = ShopFeature(entities,text)
    features['shop'] = extractor.process()
    extractor = DateFeature(entities)
    features['date']  = extractor.process()
    extractor = ProductFeature(table,text)
    features['items'] = extractor.process()
    return features
```

https://github.com/pilarcode/demo

# Componentes

- *Application.py*
- *Utils.py*
- *ReceiptFeature.py*
- *ProductFeature.py*
- *ShopFeature.py*
- *DateFeature.py*

# Tests *(I)*

```json
{
    "date": "2014-10-21",
    "items": {
        "1 Amer Kobe Filet ": "79.00",
        "1 Choc Pot De Crem ": "13.00",
        "1 Grlld Asparagus ": "13.00",
        "1 Kobe Beef Slider ": "18.00",
        "1 Lobster Tail ": "38.00",
        "1 New York ": "63.00",
        "1 Peppercorn ": "6.00",
        "1 Pepsi ": "6.00",
        "1 Pimm's Cup ": "14.00",
        "1 Potato Puree ": "12.00",
        "2 The Drifter @ 14.00 ": "28.00"
    },
    "shop": "Gordon Ramsay Steak"
}
```



```
Gordon Ramsay Steak
     Paris Las Vegas
10/21/2014              19:41
=============================
   Gordon Ramsay Steak
heck:  2042367    Table: T-35
erver: Lomberto   Guests: 2
erminal: 42
=============================
            Regular
1 Pimm's Cup            14.00
2 The Drifter          28.00
  @ 14.00
1 Kobe Beef Slider     18.00
1 New York             63.00
1 Amer Kobe Filet      79.00
1 Lobster Tail         38.00
1 Potato Puree         12.00
1 Grlld Asparagus      13.00
1 Peppercorn            6.00
1 Pepsi                 6.00
1 Choc Pot De Crem     13.00

Subtotal              290.00
     Tax               23.49
   Total              313.49
```

# Tests *(II)*



Se extrae la información aunque el ticket esta borroso.

```
{
    "date": "2019-07-11",
    "items": {
        "1 BOTTLED SODA ": "2.49",
        "1 S - TURKEY AND SWISS ": "9.99"
    },
    "shop": "MARKET FRESH KITCHEN"  👍
}
```

```json
{
    "date": "2017-02-18",
    "items": {
        "7 UP ": "3.25",
        "DIET SDA ": "3.75",
        "LING VONG ": "22.75",
        "PREM WHITE ": "10.75",
        "VONG RIPIENE ": "9.75"
    },
    "shop": "MARIANACCI'S"
}
```

| Column 1 | | Column 2 | | Column 3 |
|---|---|---|---|---|
| | | PREM WHITE | | 10.75 |
| | | PREM RED | | 10.75 |
| | | DIET SDA | | 3.75 |
| | | 7 UP | | 3.25 |
| | | VONG RIPIENE | | 9.75 |
| | | Chicken | | Parm21.75 |
| | | SCAMPI MILAN25.75 | | |
| | | LING VONG | | 22.75 |
| | | 1/2 PASTA | | SP16.75 |
| | | DIET SDA | | 3.75 |

*Faltan algunos productos...*

*Algunos casos son difíciles hasta para el propio servicio de AWS.*

# EXTRACCIÓN DE DATOS EN *RECIBOS*

SHAKE SHACK
3790 Las Vegas Blvd South
12/02/2018
Host: Lisa
11:44 AM
166 WALTER
60042

8.69
DBL ShackBurger
4.89
Bacon Cheese Fries
5.49
Shake
Vanilla Shake
19.07
Subtotal
1.57
Tax

To Stay Total    20.64

20.64

MasterCard #XXXXXXXXXXXX2825
Auth:NGJJ58

We wanna hear ya! Take our survey
for $5 off your next $20 App order.
http://bit.ly/shack-survey-1130

--- Check Closed ---

# ¡*Gracias!*

*¿Alguna pregunta?*

- ***Pilar Madariaga Lasala***
- Estudiante Máster Inteligencia Artificial en Unir
- Pasantia en Grupo Bancolombia