

Resumen

Este análisis de datos tiene como objetivo identificar las condiciones que hacen que una persona especialista en analizar datos tenga un mejor sueldo. Para encontrar respuesta las principales preguntas planteadas, se utilizó una base de datos de Kaggle con información de personas especialistas en análisis de datos alrededor del mundo.

Introducción

Este análisis se enfoca en buscar respuesta por medio de visualización de datos a las siguientes preguntas: ¿En qué países se ofrecen mejores salarios?, ¿Se han incrementado los salarios a lo largo del tiempo? y ¿Influye el nivel de experiencia en el salario?

▼ 1) EXPLORACIÓN DE LA BASE DE DATOS

A) Accede a la base de datos de Data Science Jobs Salaries: [Aquí Descargar Aquí](#).

```
1 import pandas as pd
2 import plotly.express as px
3 import plotly.graph_objects as go
4
5 ds = pd.read_csv('ds_salaries.csv')
```

B) Explora las variables y familiarízate con su significado. La página de Data

▼ Science Jobs Salaries (Enlaces a un sitio externo.) de Kaggle te puede ser de utilidad.

```
1 #Identifica la cantidad de datos y variables presentes.
2 print("-----Columnas en data frame-----")
3 print(ds.columns)

-----Columnas en data frame-----
Index(['Unnamed: 0', 'work_year', 'experience_level', 'employment_type',
      'job_title', 'salary', 'salary_currency', 'salary_in_usd',
      'employee_residence', 'remote_ratio', 'company_location',
      'company_size'],
      dtype='object')
```

```

1 #Clasifica las variables de acuerdo a su tipo y escala de medición.
2 print("-----Tipos de variables-----")
3 print(ds.dtypes)
4
5 v_cuali = ['experience_level', 'employment_type', 'job_title', 'salary_currency', 'employee_residence']
6 v_cuant = ['salary', 'salary_in_usd']
7

```

```

-----Tipos de variables-----
Unnamed: 0          int64
work_year          int64
experience_level    object
employment_type     object
job_title          object
salary             int64
salary_currency     object
salary_in_usd      int64
employee_residence object
remote_ratio       int64
company_location   object
company_size       object
dtype: object

```

```

1 df_cuali = ds[['experience_level', 'employment_type', 'job_title', 'salary_currency', 'employee_residence']]
2 df_cuanti = ds[['salary', 'salary_in_usd']]

```

▼ C) Exploración de la base de datos

1. Calcula medidas estadísticas

* Variables cuantitativas

- Medidas de tendencia central: promedio, media, mediana y moda de los datos.
- Medidas de dispersión: rango: máximo - mínimo, varianza, desviación estándar.

```

1 # Medidas de tendencia central: promedio, media, mediana y moda de los datos.
2
3 mean = df_cuanti.mean()
4 median = df_cuanti.median()
5 mode = df_cuanti.mode()
6
7 print("---MEDIA---")
8 print(mean)
9 print("-----")
10 print("---MEDIANA---")
11 print(median)
12 print("-----")

```

```

13 print("---MODA---")
14 print(mode)
15 print("-----")
16

```

```

---MEDIA---
salary          324000.062603
salary_in_usd    112297.869852
dtype: float64
-----

```

```

---MEDIANA---
salary          115000.0
salary_in_usd    101570.0
dtype: float64
-----

```

```

---MODA---
   salary  salary_in_usd
0   80000      100000.0
1  100000             NaN
-----

```

```

1 #Medidas de dispersión: rango: máximo - mínimo, varianza, desviación estándar.
2
3 rango = df_cuanti.max() - df_cuanti.min()
4 varianza = df_cuanti.var()
5 des_est = df_cuanti.std()
6
7 print("---RANGO---")
8 print(rango)
9 print("-----")
10 print("---VARIANZA---")
11 print(varianza)
12 print("-----")
13 print("---DESVIACIÓN ESTANDAR---")
14 print(des_est)
15 print("-----")
16

```

```

---RANGO---
salary          30396000
salary_in_usd    597141
dtype: int64
-----

```

```

---VARIANZA---
salary          2.385040e+12
salary_in_usd    5.034933e+09
dtype: float64
-----

```

```

---DESVIACIÓN ESTANDAR---
salary          1.544357e+06
salary_in_usd    7.095726e+04
dtype: float64
-----

```

* Variables cualitativas

- Tabla de distribución de frecuencia
- Moda

```
1 # Tabla de distribución de frecuencia
2 cols = df_cuali.columns
3 for i in range(1, cols.size):
4     print(df_cuali.groupby(cols[i]).agg(frequency=(cols[i], "count")))
5
```

employment_type	frequency
CT	5
FL	4
FT	588
PT	10

job_title	frequency
3D Computer Vision Researcher	1
AI Scientist	7
Analytics Engineer	4
Applied Data Scientist	5
Applied Machine Learning Scientist	4
BI Data Analyst	6
Big Data Architect	1
Big Data Engineer	8
Business Data Analyst	5
Cloud Data Engineer	2
Computer Vision Engineer	6
Computer Vision Software Engineer	3
Data Analyst	97
Data Analytics Engineer	4
Data Analytics Lead	1
Data Analytics Manager	7
Data Architect	11
Data Engineer	132
Data Engineering Manager	5
Data Science Consultant	7
Data Science Engineer	3
Data Science Manager	12
Data Scientist	143
Data Specialist	1
Director of Data Engineering	2
Director of Data Science	7
ETL Developer	2
Finance Data Analyst	1
Financial Data Analyst	2
Head of Data	5
Head of Data Science	4
Head of Machine Learning	1

Lead Data Analyst	3
Lead Data Engineer	6
Lead Data Scientist	3
Lead Machine Learning Engineer	1
ML Engineer	6
Machine Learning Developer	3
Machine Learning Engineer	41
Machine Learning Infrastructure Engineer	3
Machine Learning Manager	1
Machine Learning Scientist	8
Marketing Data Analyst	1
NLP Engineer	1
Principal Data Analyst	2
Principal Data Engineer	3
Principal Data Scientist	7
Product Data Analyst	2
Research Scientist	16
Staff Data Scientist	1

```
1 #MODA
2 df_cuali.mode()
```

experience_level	employment_type	job_title	salary_currency	employee_residence	c
SE	FT	Data	USD	US	

2) Explora los datos usando herramientas de visualización

* Variables cuantitativas:

- Medidas de posición: cuartiles, outlier (valores atípicos), boxplots
- Análisis de distribución de los datos (Histogramas). Identificar si tiene forma simétrica o así

```
1 #Medidas de posición: cuartiles, outlier (valores atípicos), boxplots
2 import plotly.express as px
3
4 print("---CUARTILES---")
5 print(df_cuanti.quantile([.25, .50, .75]))
6 print("---OUTLIER---")
7 q1 = df_cuanti.quantile(.25)
8 q3 = df_cuanti.quantile(.75)
9 iqr = q3 - q1
10 print(q3)
11 print("---BOXPLOTS---")
12 fig = px.box(df_cuanti, y= v_cuant)
13 fig.show()
14
```

```
---CUARTILES---
      salary  salary_in_usd
0.25   70000.0       62726.0
0.50  115000.0      101570.0
0.75  165000.0      150000.0
---OUTLIER---
salary           165000.0
salary_in_usd    150000.0
Name: 0.75, dtype: float64
---BOXPLOTS---
```

* Variables categóricas

- Distribución de los datos (diagramas de barras, diagramas de pastel)

```
1 cols = df_cuali.columns
2 for i in range(1, cols.size):
3     aux = df_cuali.groupby(cols[i]).agg(
4         frequency=(cols[i], "count"))
5
6     fig = px.bar(aux, x='frequency')
```

```
7     fig.show()  
8
```


3) Identifica problemas de calidad de datos (registros duplicados, valores faltantes, outliers, etc).

```
1 from scipy.stats import zscore
2 #Valores faltantes?
3 print(ds.isnull().values.any())
4 #Valores nulos?
5 print(ds.isna().values.any())
6
```

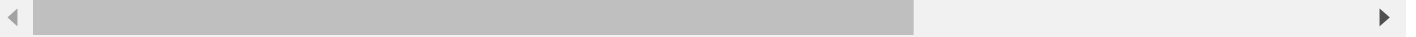
```
False
False
```

```
1 # Delete outliers
2 df_cuanti = df_cuanti[(zscore(df_cuanti) < 3).all(axis=1)]
3 print("---BOXPLOTS---")
4 fig = px.box(df_cuanti, y=v_cuant)
5 fig.show()
6
```

```
---BOXPLOTS---
```

D) Preparación de los datos: 1. Selecciona el conjunto de datos a utilizar.

- * Decide qué conjunto de datos se utilizará. Identifica variables objetivo. En caso necesario, ex
- * Maneja datos categóricos: transforma a datos numéricos si es necesario.
- * En caso de necesidad de recorte de datos (atípicos, faltantes, duplicados, etc), explica el mot
- * Maneja apropiadamente datos atípicos.



```
1 v_selec = ds[['experience_level', 'company_location', 'work_year', 'salary_in_usd']]
```

2. Transforma los datos en caso necesario.

- * Revisa si es necesario discretizar los datos
- * Revisa si es necesario escalar y normalizar los datos
- * Construye atributos si es conveniente

2) ANALIZA LOS DATOS Y CONTESTA TUS PREGUNTAS GUÍA

```
1
2 fig = px.histogram(v_selec, x='company_location', y='salary_in_usd', color="company_locati
3 fig.show()
```

```
1 fig = px.histogram(v_selec, x='work_year', y='salary_in_usd', color="work_year",  
2                     title='¿Se han incrementado los salarios a lo largo del tiempo?')  
3 fig.show()  
4
```

```
1 fig = px.histogram(v_selec, x='experience_level', y='salary_in_usd',  
2                     color="experience_level", title='¿Influye el nivel de experiencia en el  
3 fig.show()  
4
```

[Productos pagados de Colab](#) - [Cancela los contratos aquí](#)

