

PROYECTO

Presentación final aprendizaje estadístico

Presentado por:

Julio Avila
Guillermo Furlan

AGENDA

- 01** Información General del proyecto
- 02** Objetivos
- 03** Análisis exploratorio
- 04** Clasificación
- 05** Modelo de preicción
- 06** final



Contexto

Se tiene una base de datos que contiene información acerca de empleos relacionados a datos, se cuenta entre otras variables; el salario promedio, puesto, empresa donde se trabaja y lista de habilidades colocadas en el cv.

Objetivos

01

**ELABORAR UN MODELO PARA
PREDECIR EL RANGO SALARIAL
PARA UN PUESTO EN CIENCIA
DE DATOS**

02

**CLASIFICAR EL PERFIL DE
UNA PERSONA BASADO EN
SUS HABILIDADES**

ANÁLISIS EXPLORATORIO

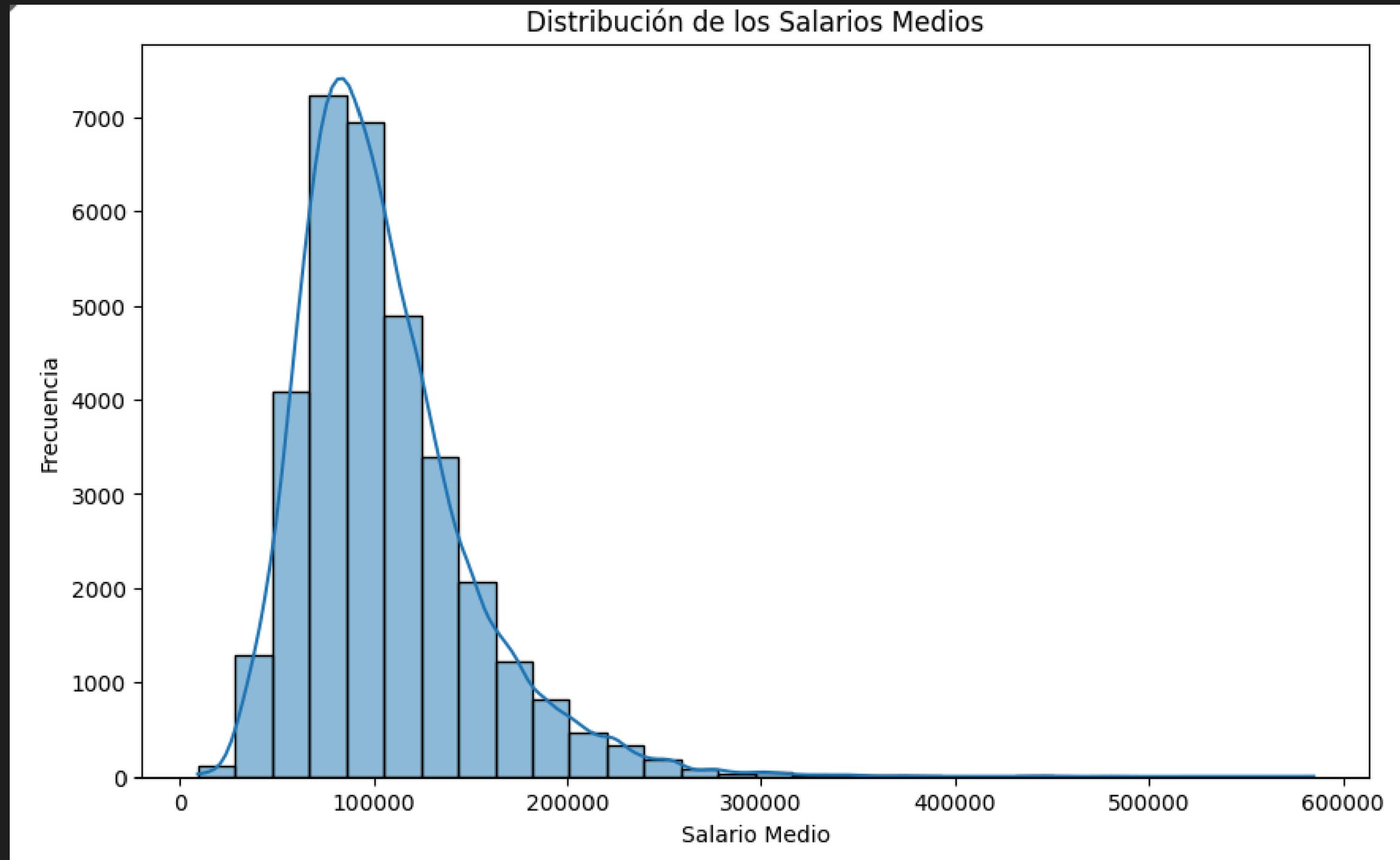


Resultado

#	Column	Non-Null	Count	Dtype
0	ID	33248	non-null	object
1	Job	33248	non-null	object
2	Jobs_Group	33248	non-null	object
3	Profile	12141	non-null	object
4	Remote	13929	non-null	object
5	Company	33239	non-null	object
6	Location	33235	non-null	object
7	City	29424	non-null	object
8	State	30136	non-null	object
9	Frecuency_Salary	33248	non-null	object
10	Mean_Salary	33248	non-null	float64
11	Skills	33248	non-null	object
12	Sector	26034	non-null	object
13	Sector_Group	26034	non-null	object
14	Revenue	14930	non-null	object
15	Employee	20449	non-null	object
16	Company_Score	24486	non-null	float64
17	Reviews	24486	non-null	float64
18	Director	12463	non-null	object
19	Director_Score	11324	non-null	float64
20	URL	17215	non-null	object



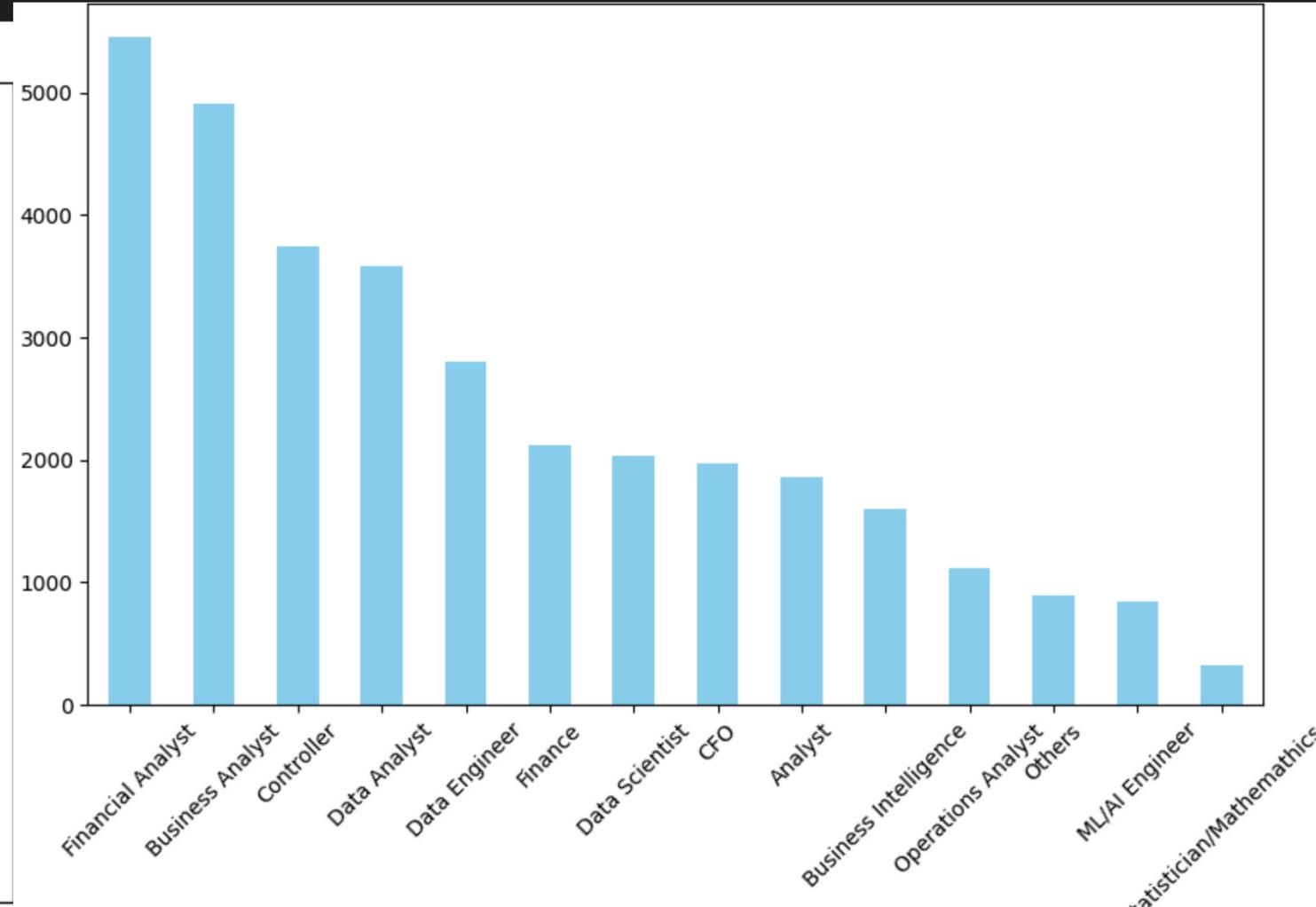
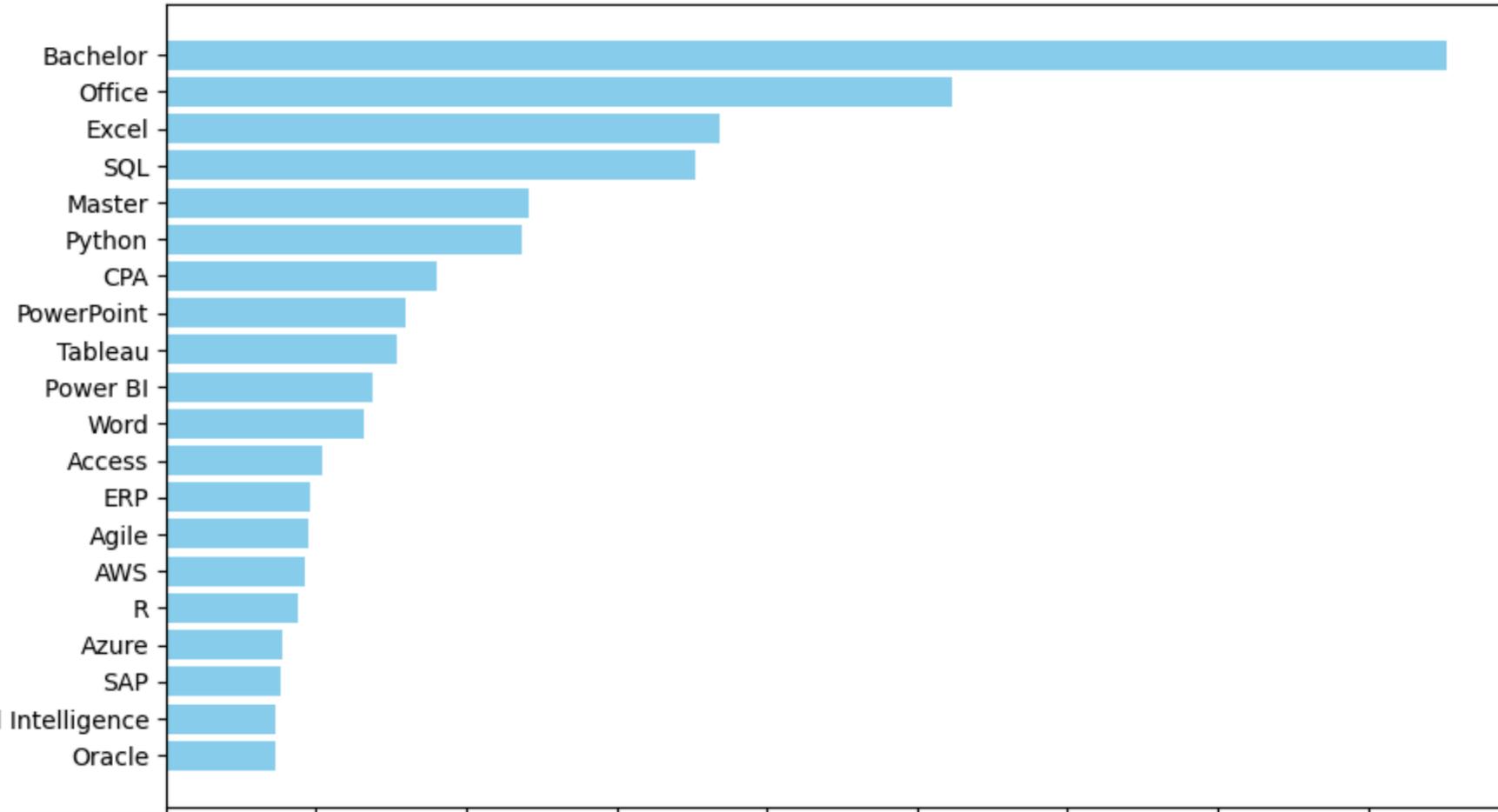
Resultado





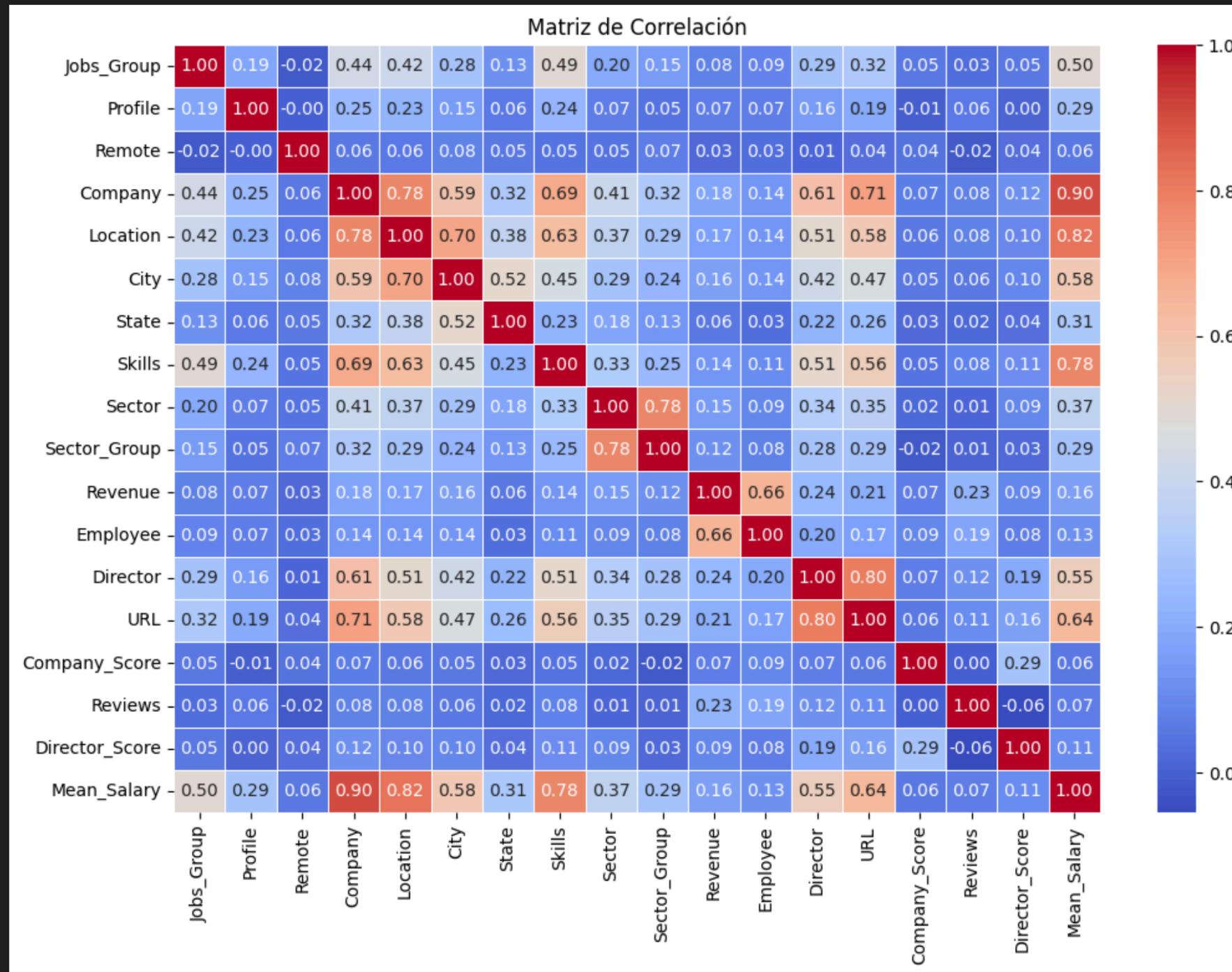
Resultado

Top 20 Elementos Más Frecuentes





Resultado



CLASIFICACIÓN

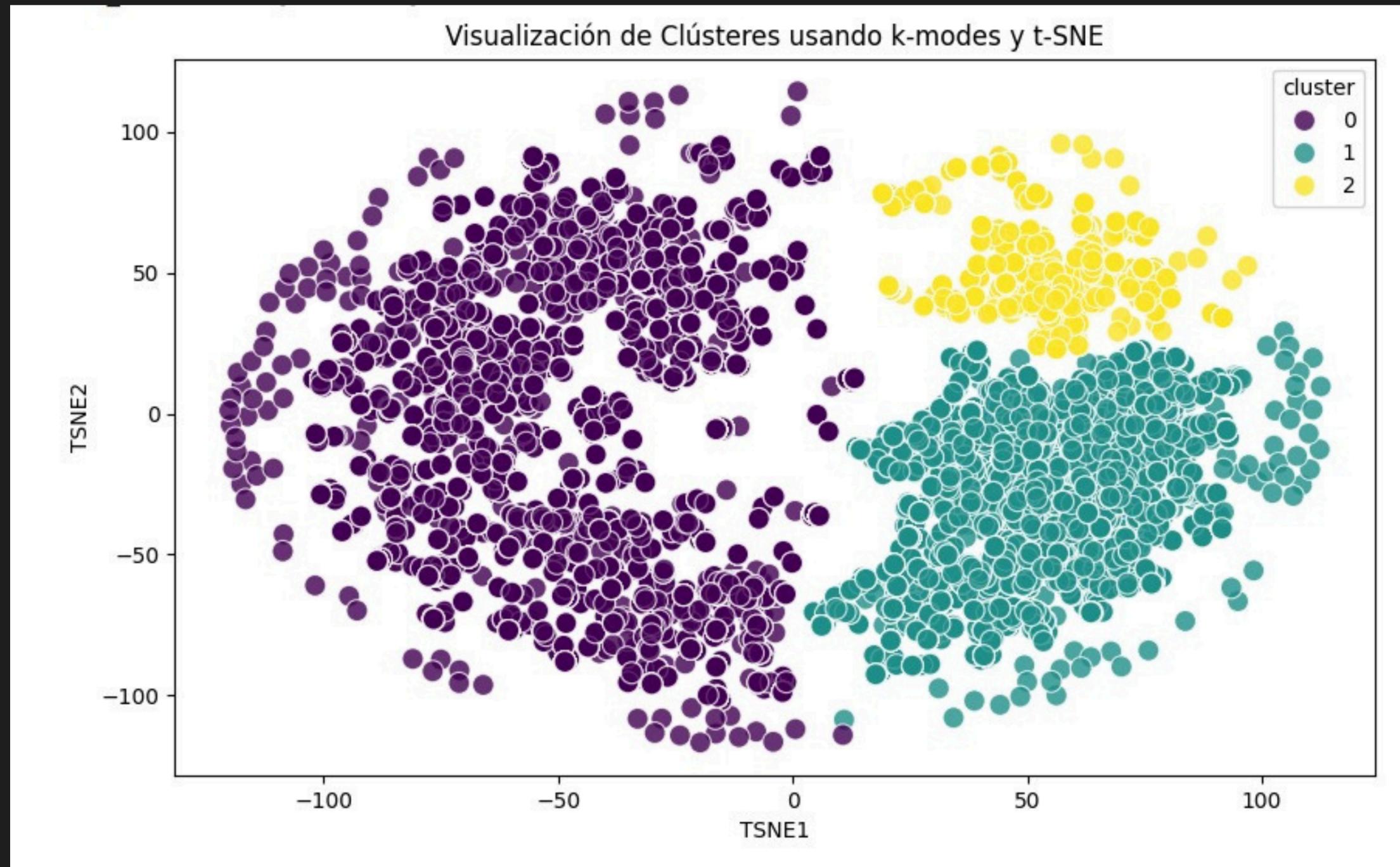
Técnicas usadas

- **K-MODES:** el algoritmo K-Modes es una variante del algoritmo K-Means diseñada específicamente para datos categóricos. Usando la métrica Hamming
- **TSNE:** t-distributed Stochastic Neighbor Embedding, es un algoritmo de reducción de dimensionalidad no lineal





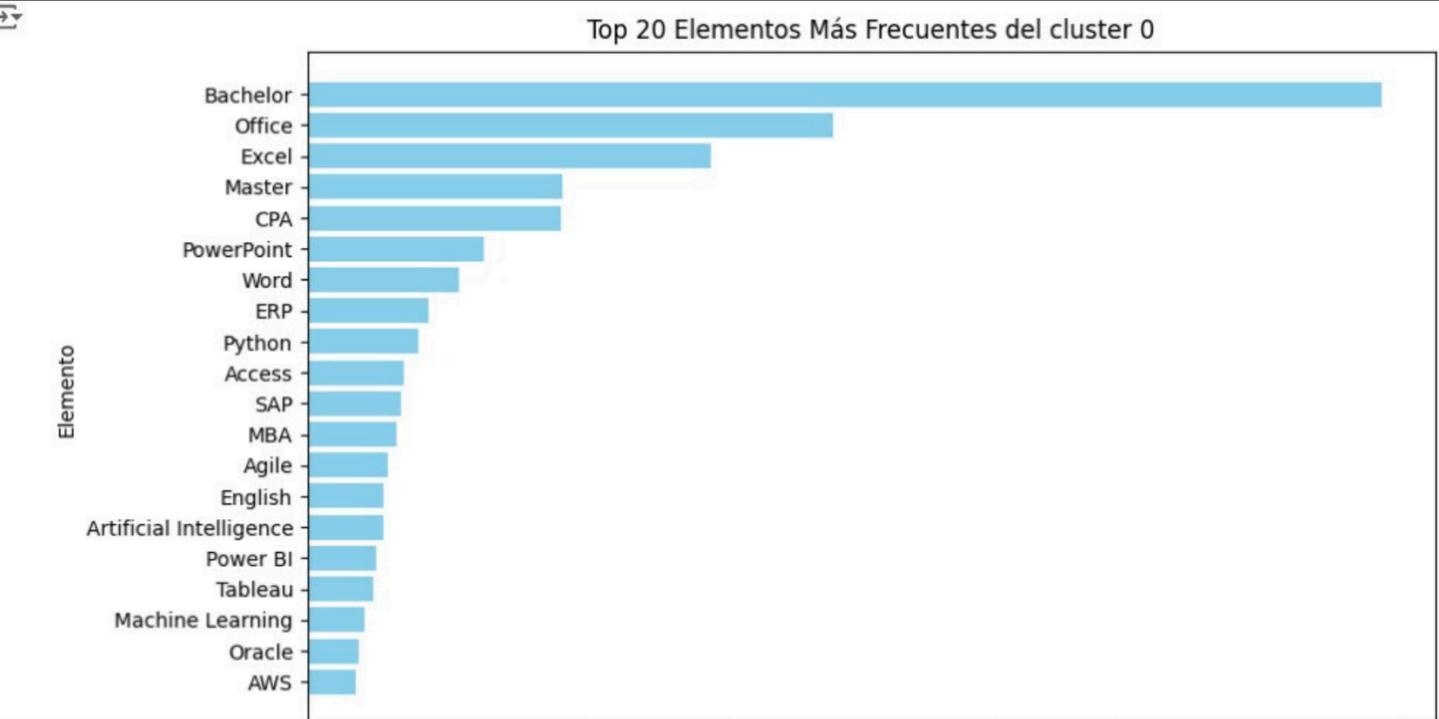
Resultado



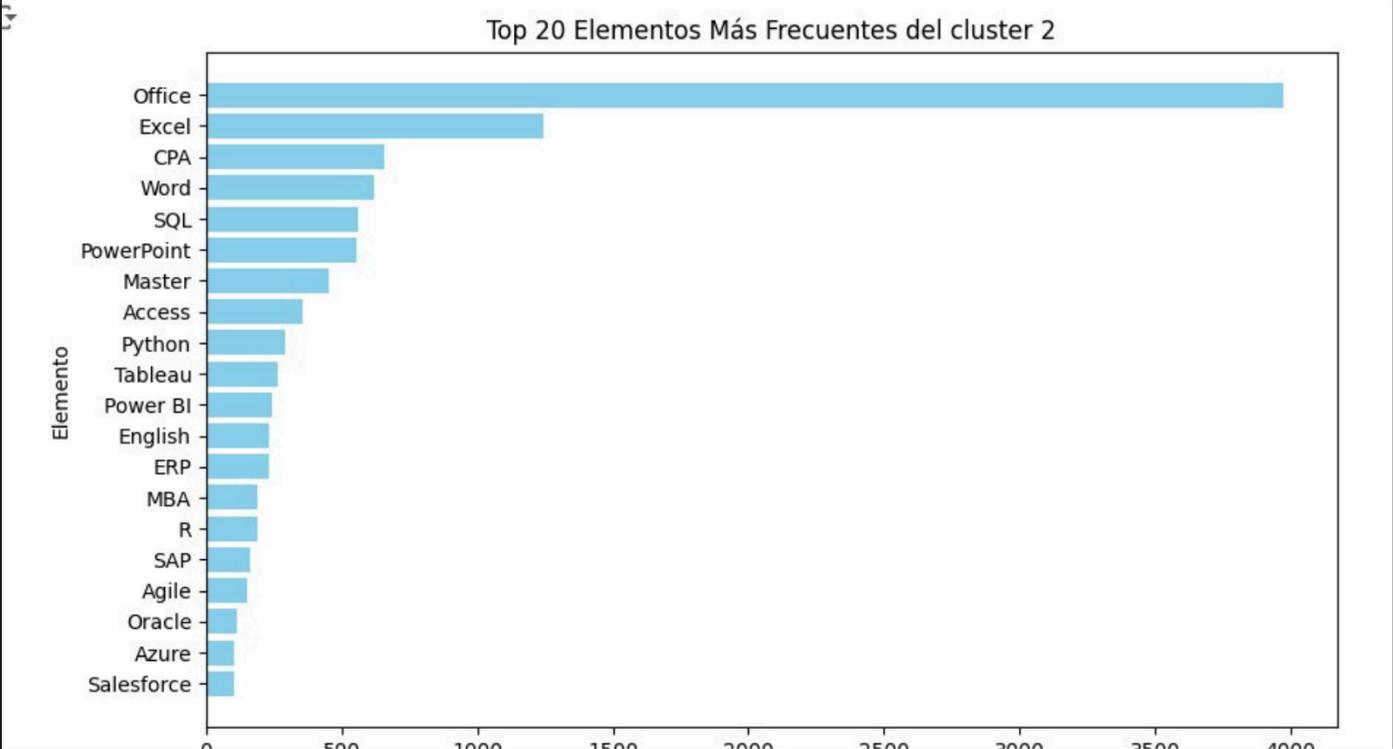
Habilidades



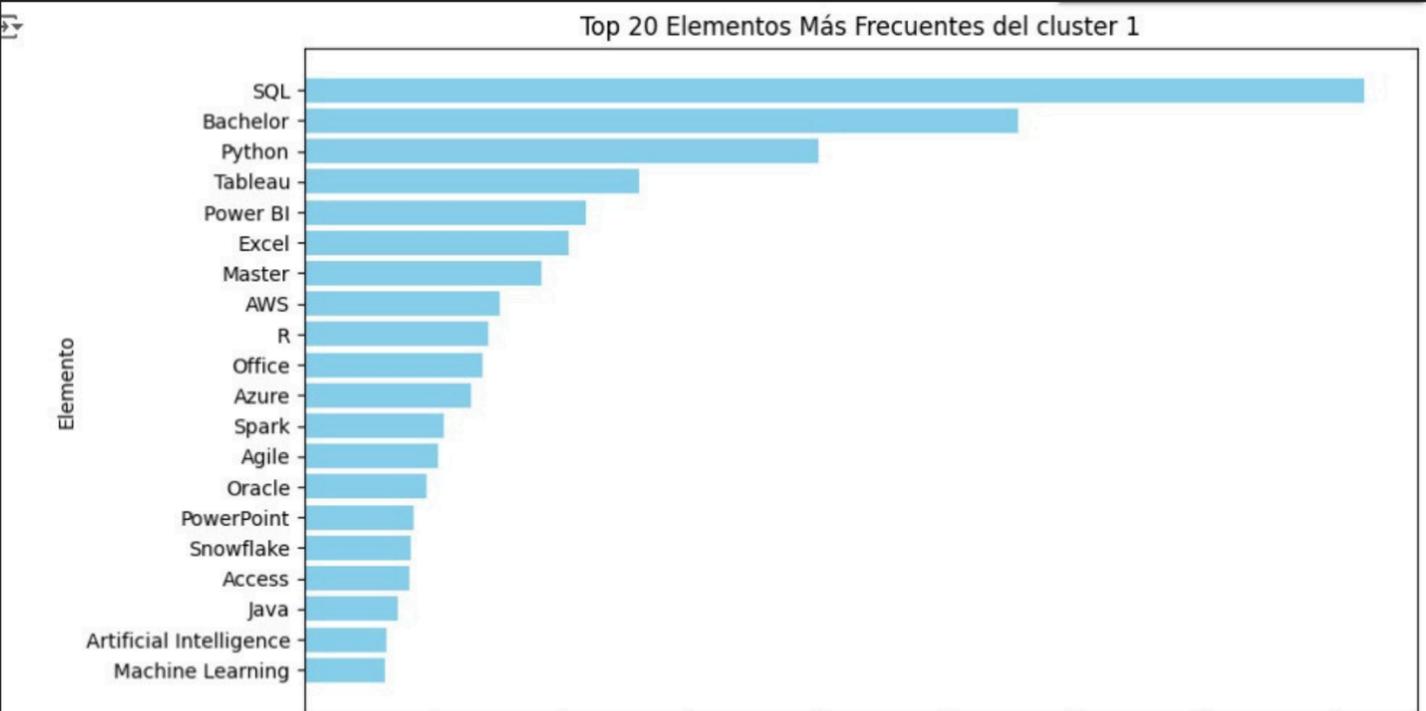
Top 20 Elementos Más Frecuentes del cluster 0



Top 20 Elementos Más Frecuentes del cluster 2

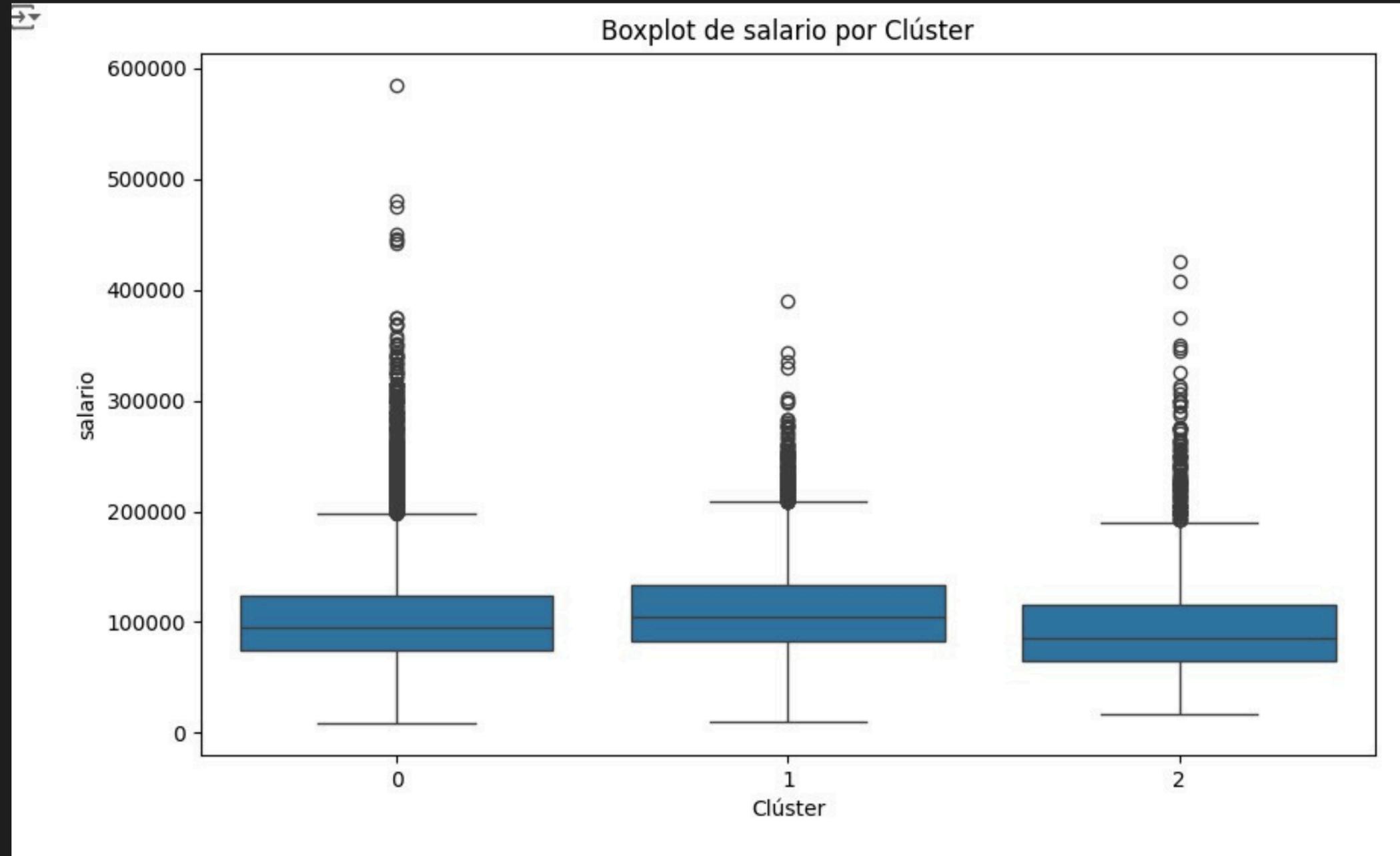


Top 20 Elementos Más Frecuentes del cluster 1



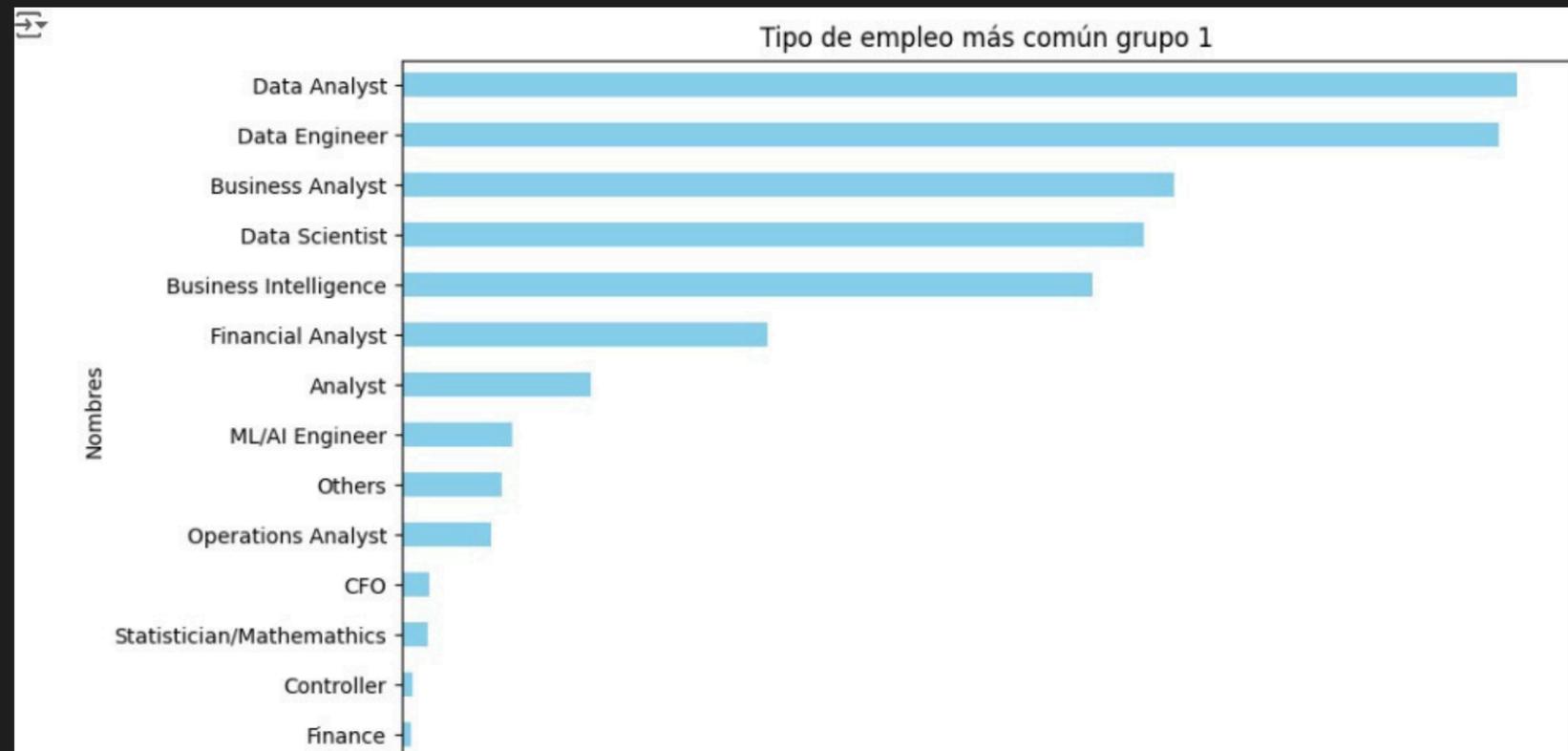
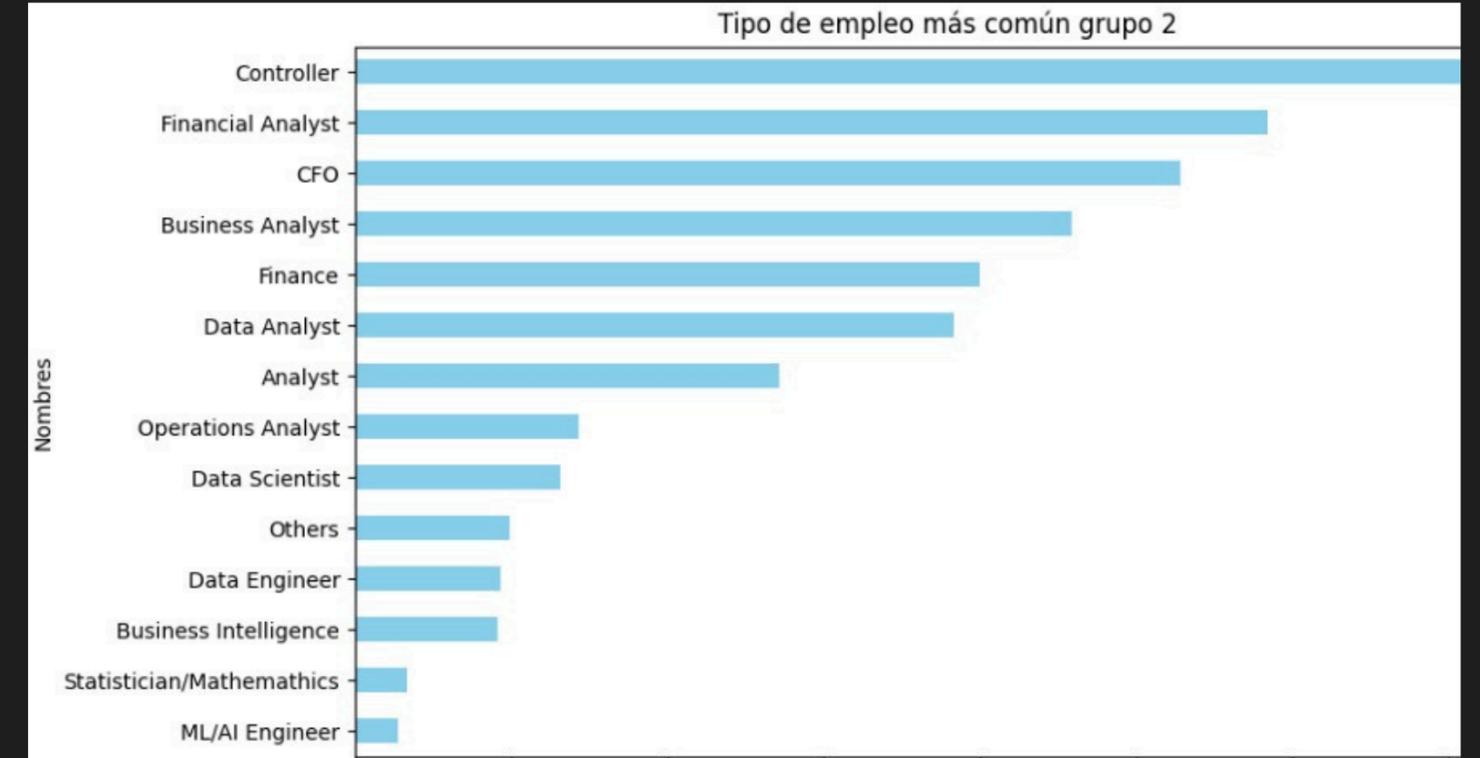
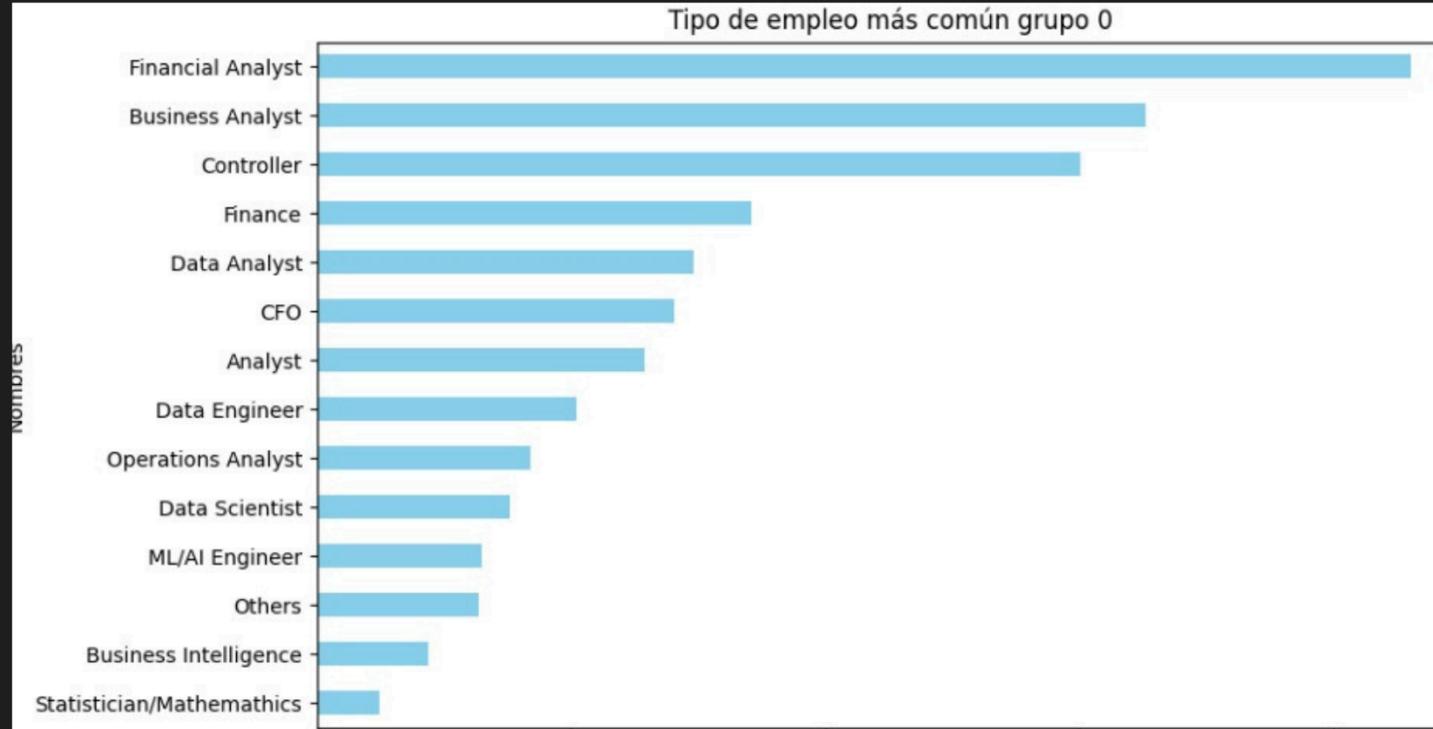


Salario Medio





Tipo de trabajo

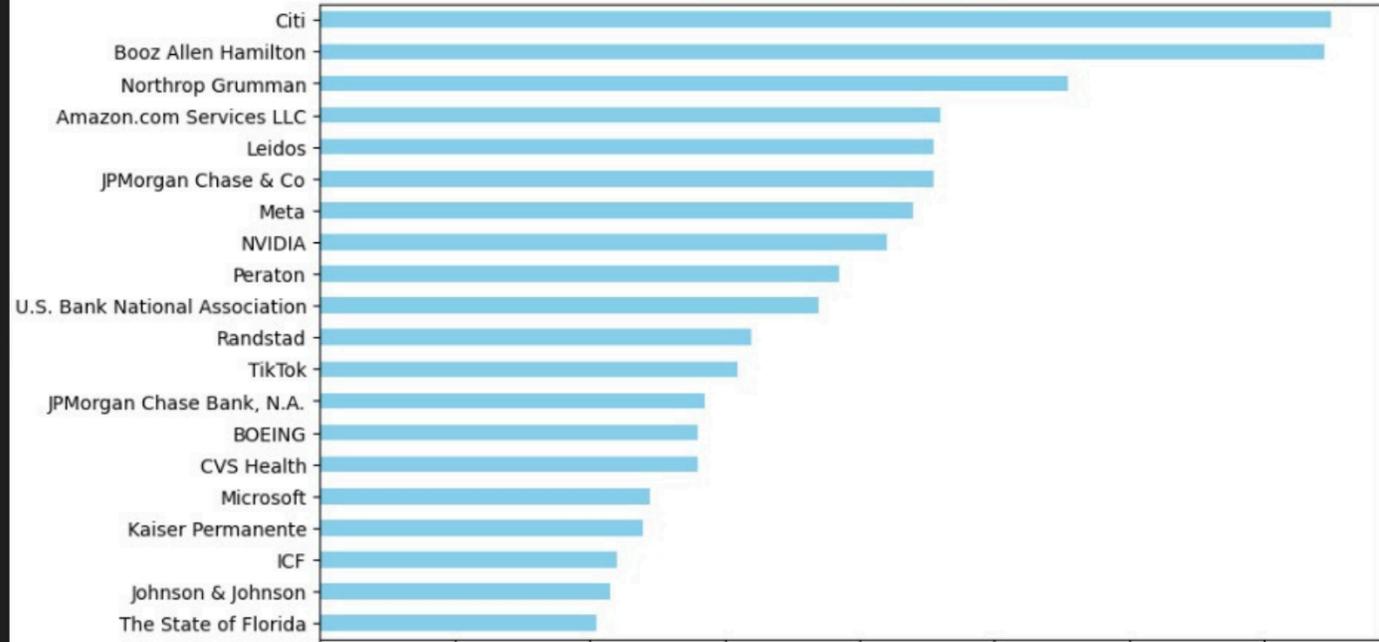




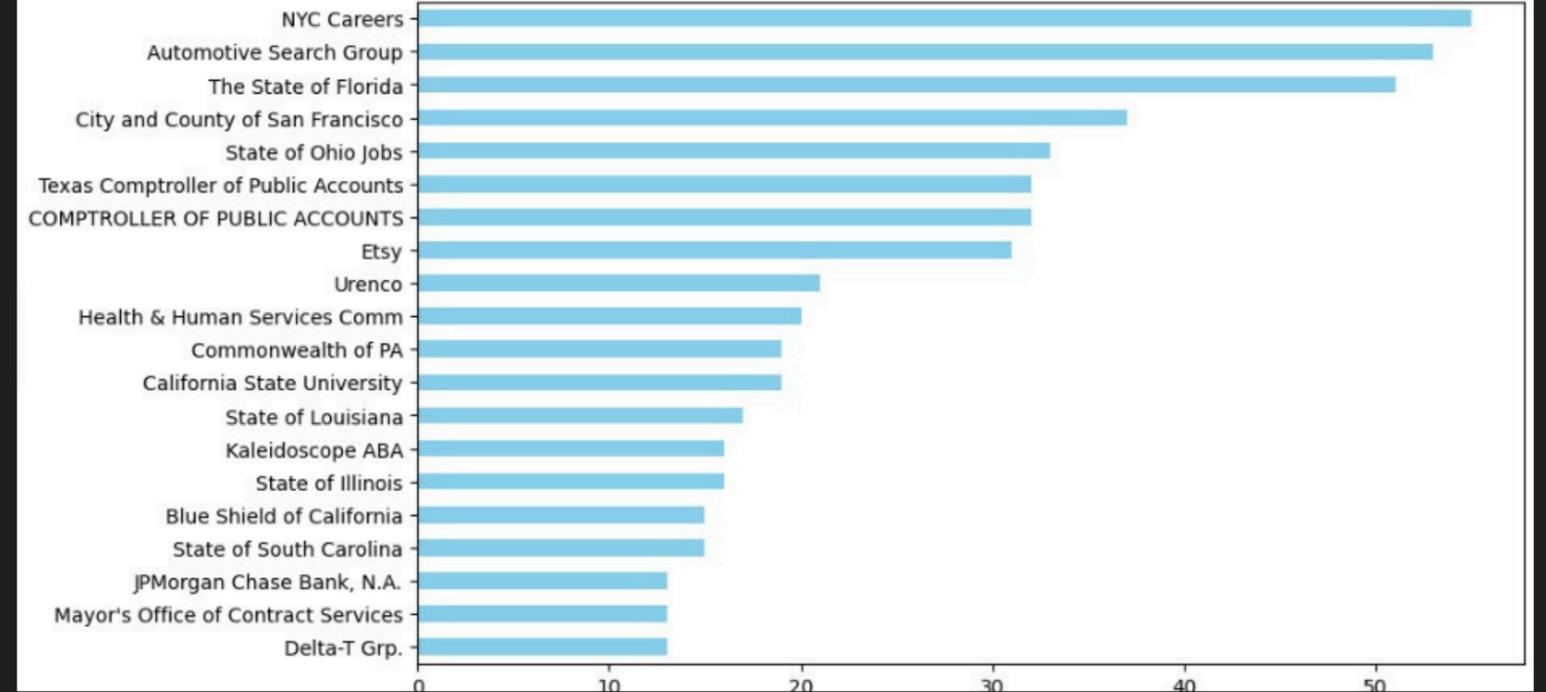
Empresas



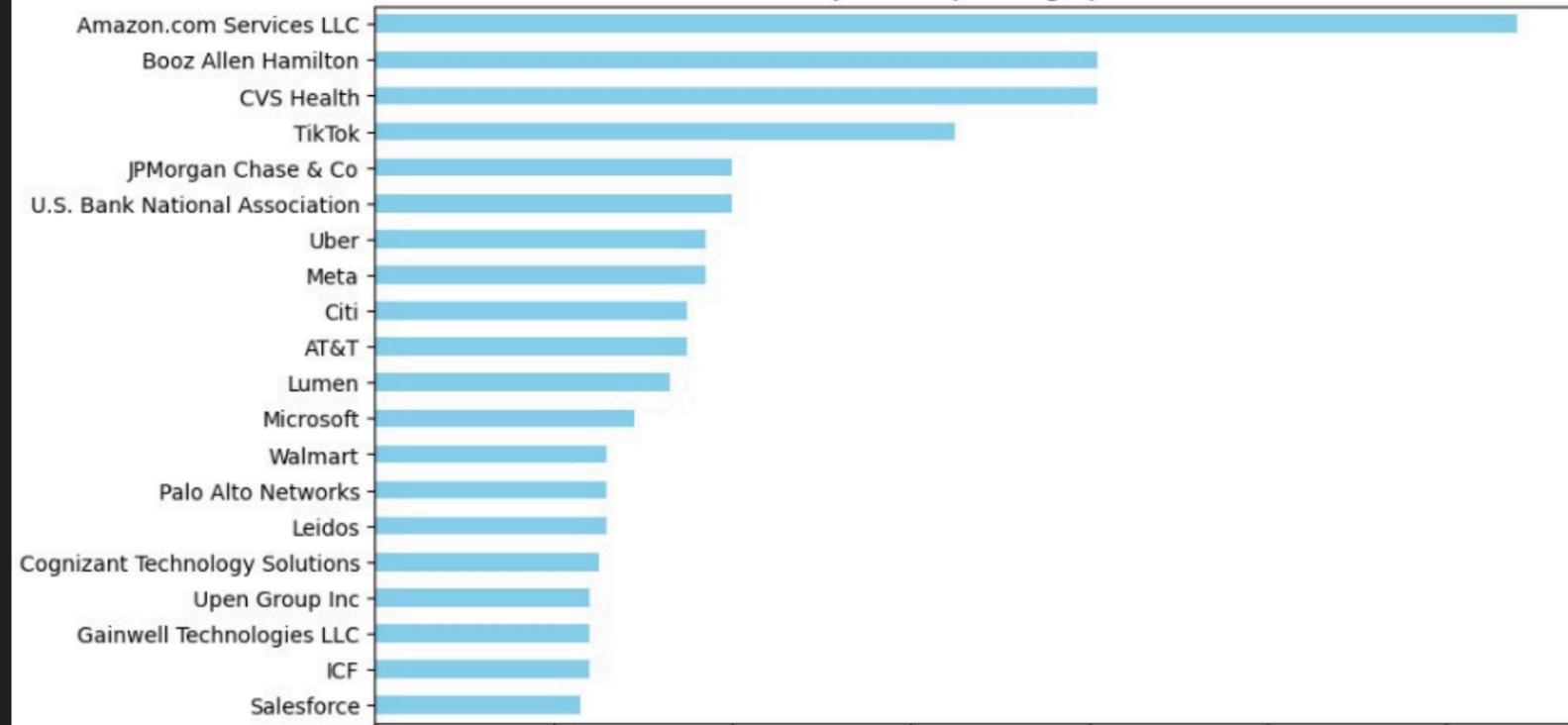
top 20 compañías grupo 0



top 20 compañías grupo 2



top 20 compañías grupo 1



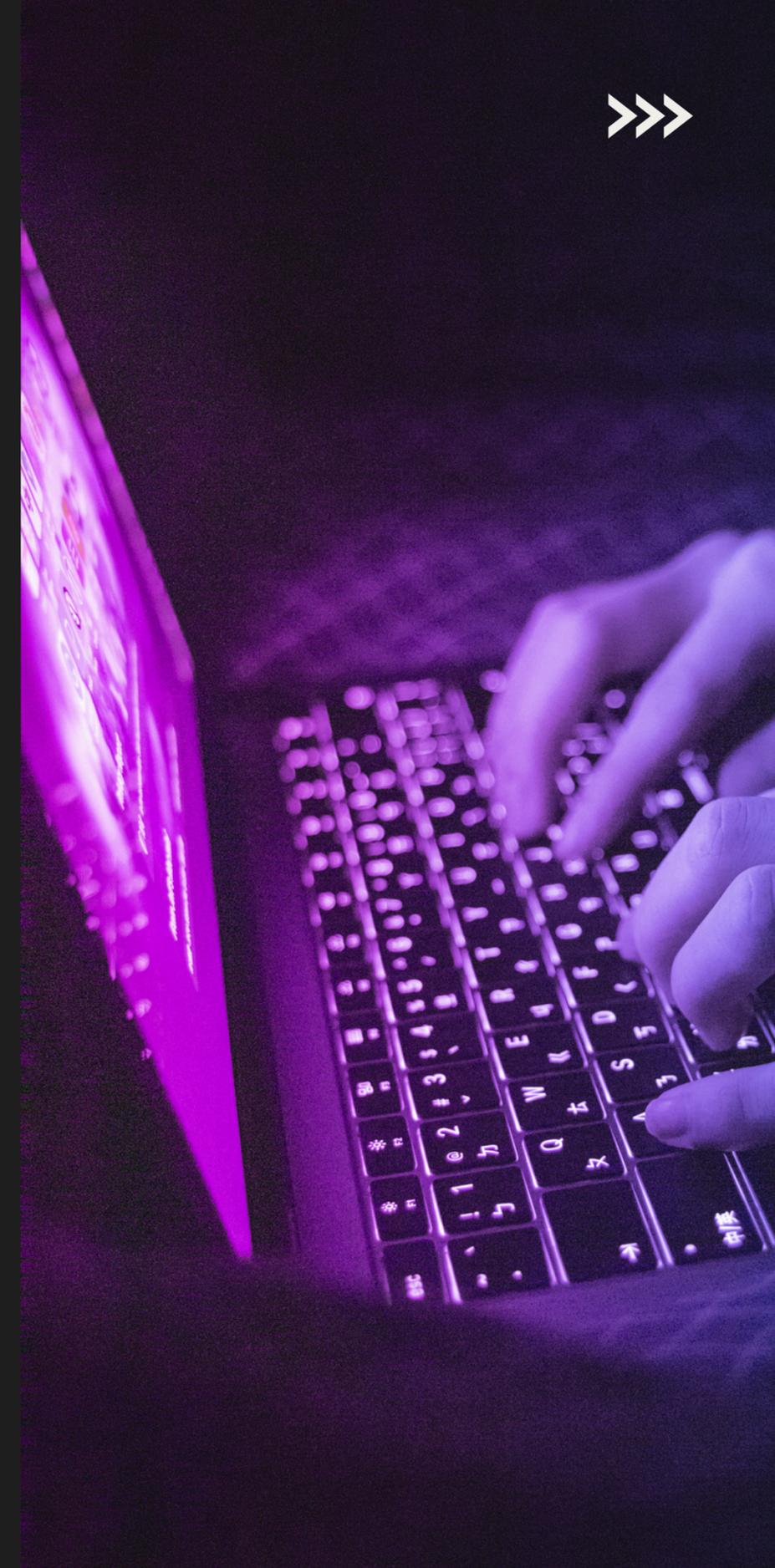
MODELO DE PREDICCIÓN



Modelos utilizados

Para realizar este proyecto, se hicieron las predicciones con varios diferentes modelos, los cuales fueron:

- Random Forest
- Decision Tree
- AdaBoost
- Gradient Boosting
- Linear Regression
- XGBoost
- SVM



Resultados

Después de implementar dichos modelos, se concluyó que el mejor modelo para este problema es el Random Forest, aunque la regresión lineal tuvo rendimiento similar.

```
Mejor modelo: Random Forest
Métricas:
RMSE: 28692.82
R2 Score: 0.55
```

```
Random Forest:
RMSE: 28692.82
R2 Score: 0.55
```

```
Decision Tree:
RMSE: 35982.70
R2 Score: 0.29
```

```
AdaBoost:
RMSE: 40026.16
R2 Score: 0.12
```

```
Gradient Boosting:
RMSE: 31792.15
R2 Score: 0.44
```

```
Linear Regression:
RMSE: 29701.97
R2 Score: 0.51
```

```
XGBoost:
RMSE: 31841.47
R2 Score: 0.44
```

```
SVM:
RMSE: 38789.30
R2 Score: 0.17
```



Adicionalmente, se quiso analizar si los resultados que se estaban prediciendo se encontraban en un rango aceptable para cada grupo de trabajo, por lo que se tomó la media de cada grupo y se hizo un rango con los cuartiles 1 y 3 para comparar si las predicciones estaban muy alejadas en un panorama más general y se obtuvieron resultados positivos.

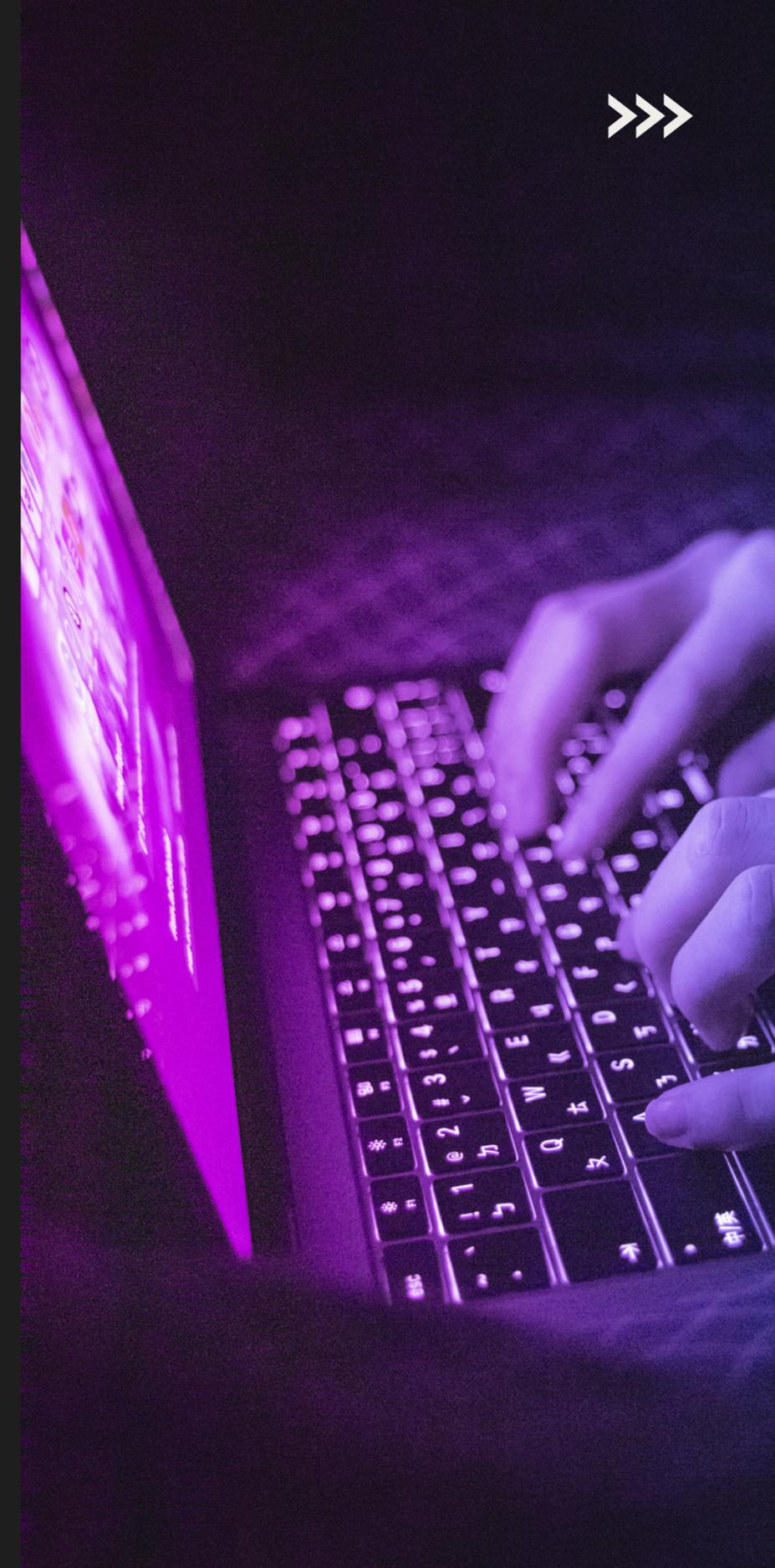
```
Matriz de conteo:
```

Jobs_Group	Within_Range	Out_of_Range	Lower_Salary	Upper_Salary	Percentage
Controller	388	124	85000.00	127500.00	75.78
Financial Analyst	442	230	71241.75	103500.00	65.77
Business Intelligence	144	52	87045.25	130000.00	73.47
Business Analyst	342	176	75000.00	110000.00	66.02
Data Engineer	238	106	108503.50	160000.00	69.19
Analyst	139	69	67500.00	98777.81	66.83
ML/AI Engineer	84	30	136000.00	203625.00	73.68
Data Analyst	271	137	71000.00	112500.00	66.42
CFO	202	61	101362.50	175000.00	76.81
Data Scientist	172	100	112500.88	171050.25	63.24
Others	64	27	85000.00	137401.50	70.33
Operations Analyst	89	40	68363.00	107500.00	68.99
Finance	189	87	72325.00	110000.00	68.48
Statistician/Mathemathics	28	17	80000.00	134500.00	62.22



Modelos utilizados

Debido a que los modelos funcionaban correctamente, se consideró que la precisión aún era muy baja y se buscó una implementación más potente. Por lo que se optó por un modelo stacked que tomó en cuenta el random forest y la regresión lineal



Resultados

Los resultados de las predicciones de este modelo mejoraron los modelos anteriores, obteniendo un R^2 score de 0.6

Modelo Stacked:

RMSE: 26803.81

R^2 Score: 0.60

CONCLUSIONES Y RECOMENDACIONES

Conclusiones

- ES POSIBLE HACER PREDICCIONES DE SALARIOS SEGÚN DIFERENTES CARACTERÍSTICAS DE LA PERSONA
- EN ESTADOS UNIDOS, HAY UNA CORRELACIÓN FUERTE ENTRE LA LOCACIÓN Y EL SALARIO

Recomendaciones

- UTILIZAR MAYOR CANTIDAD DE DATOS PARA PODER MEJORAR LA PRESICIÓN DE LOS MODELOS
- UTILIZAR UNA BASE DE DATOS CON MENOS VALORES FALTANTES, YA QUE AUNQUE SE MANEJEN ESOS CASOS DE FORMA ESPECIAL, PUEDEN TENER INCIDENCIA EN EL RESULTADO FINAL DE LOS MODELOS

**MUCHAS
GRACIAS**