

# Data Summary

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
df = pd.read_csv('/work/diabetes_prediction_dataset.csv')
```

## Reporte de Datos

El conjunto de datos esta compuesto por las siguientes columnas

['gender', 'age', 'hypertension', 'heart\_disease', 'smoking\_history', 'bmi', 'HbA1c\_level', 'blood\_glucose\_level', 'diabetes'] Donde diabetes es la variable objetivo y el resto son las características

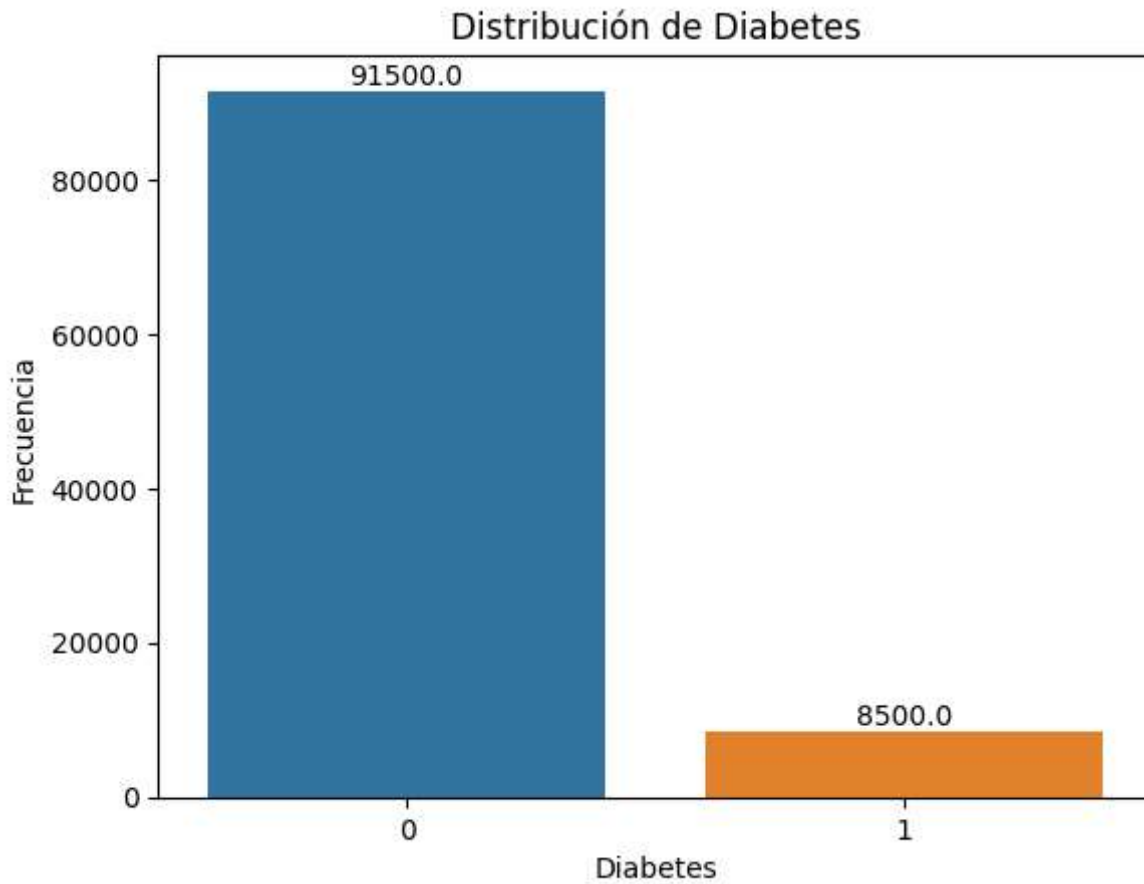
Por otra parte, se observa que es un df desbalanceado y no hay existencia de datos nulos

```
sns.countplot(x='diabetes', data=df, order=df['diabetes'].value_counts().index)

for p in plt.gca().patches:
    plt.gca().annotate(f'{p.get_height()}', (p.get_x() + p.get_width() / 2, p.get_height()),
                      ha='center', va='bottom')

plt.xlabel('Diabetes')
plt.ylabel('Frecuencia')
plt.title('Distribución de Diabetes')

plt.show()
```



## Correlacion e Importancia de Variables

Se observa que las variables que tienen una correlacion con respecto a diabetes son:

1. blood\_glucose\_level
2. HbA1c\_level
3. age
4. bmi
5. hypertension
6. heart\_disease

```
heatmap = sns.heatmap(df.corr(), vmin=-1, vmax=1, annot=True)
heatmap.set_title('Correlation Heatmap', fontdict={'fontsize':12}, pad=12);
```

