

```
In [1]: import pandas as pd  
import matplotlib.pyplot as plt  
import numpy as np
```

## Classification

```
In [2]: df_class = pd.read_csv("../data/classification.csv")
```

```
In [3]: df_class.shape
```

```
Out[3]: (1500, 16)
```

```
In [4]: df_class.dtypes
```

```
Out[4]: X1      float64  
X2      float64  
X3      float64  
X4      float64  
X5      float64  
X6      float64  
X7      float64  
X8      float64  
X9      float64  
X10     float64  
X11     float64  
X12     float64  
X13     float64  
X14     float64  
X15     float64  
Y       int64  
dtype: object
```

```
In [5]: df_class.head()
```

```
Out[5]:      X1        X2        X3        X4        X5        X6        X7  
0   -2.388741  6.221087  3.442447  1.273807  0.912272  8.908027  8.441999  -2.6  
1   -6.012792 -9.884413 -1.590610  4.999943  0.247758 -1.197048 -10.939272  1.5  
2    2.270829 -8.849332 -6.619179 -2.861520 -6.720253  5.715418  6.493857 -4.4  
3   -7.092421 -10.254081 -0.907321  3.712683 -0.567676  0.254027 -10.135377 -0.4  
4   -2.246293  7.617936  3.580218  2.412760  3.881735  8.096439  8.372886 -4.6
```

```
In [6]: df_class.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1500 entries, 0 to 1499
Data columns (total 16 columns):
 #   Column  Non-Null Count  Dtype  
--- 
 0   X1      1500 non-null   float64
 1   X2      1497 non-null   float64
 2   X3      1500 non-null   float64
 3   X4      1497 non-null   float64
 4   X5      1499 non-null   float64
 5   X6      1500 non-null   float64
 6   X7      1499 non-null   float64
 7   X8      1499 non-null   float64
 8   X9      1500 non-null   float64
 9   X10     1498 non-null   float64
 10  X11     1498 non-null   float64
 11  X12     1499 non-null   float64
 12  X13     1499 non-null   float64
 13  X14     1500 non-null   float64
 14  X15     1500 non-null   float64
 15  Y       1500 non-null   int64  
dtypes: float64(15), int64(1)
memory usage: 187.6 KB
```

In [7]: `df_class.describe()`

	X1	X2	X3	X4	X5	X
<b>count</b>	1500.000000	1497.000000	1500.000000	1497.000000	1499.000000	1500.000000
<b>mean</b>	-1.883910	-4.075099	-1.382158	0.920043	-2.471415	4.634421
<b>std</b>	3.834396	7.897076	5.051711	3.539492	3.470882	3.955870
<b>min</b>	-9.028149	-13.311122	-10.556390	-6.713349	-10.050322	-3.461190
<b>25%</b>	-5.695571	-9.919480	-7.157683	-2.919740	-6.301859	0.371100
<b>50%</b>	-1.986221	-8.899808	-0.788744	2.032352	-1.602947	5.129050
<b>75%</b>	1.998095	6.259801	3.629647	3.778093	0.507027	8.334600
<b>max</b>	6.040819	11.162159	7.478247	7.514823	3.915434	12.017230

In [8]: `missing_summary = df_class.isnull().sum().to_frame(name="n_missing")  
missing_summary["pct_missing"] = missing_summary["n_missing"] / len(df_class)  
print("\nResumen de valores faltantes:")  
missing_summary`

Resumen de valores faltantes:

Out[8]:

	n_missing	pct_missing
X1	0	0.000000
X2	3	0.200000
X3	0	0.000000
X4	3	0.200000
X5	1	0.066667
X6	0	0.000000
X7	1	0.066667
X8	1	0.066667
X9	0	0.000000
X10	2	0.133333
X11	2	0.133333
X12	1	0.066667
X13	1	0.066667
X14	0	0.000000
X15	0	0.000000
Y	0	0.000000

In [9]: df\_class.shape

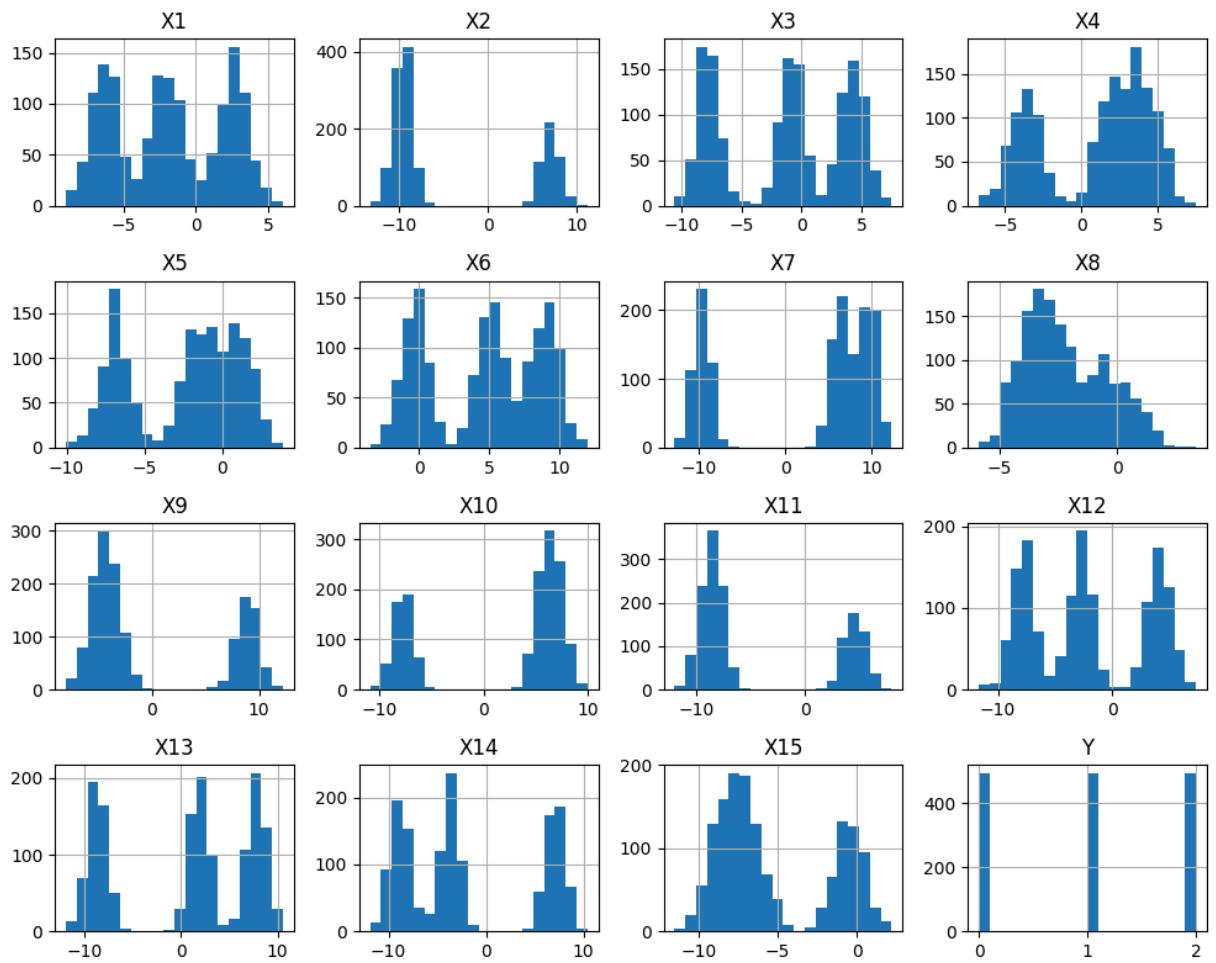
Out[9]: (1500, 16)

In [10]: df\_class = df\_class.dropna()

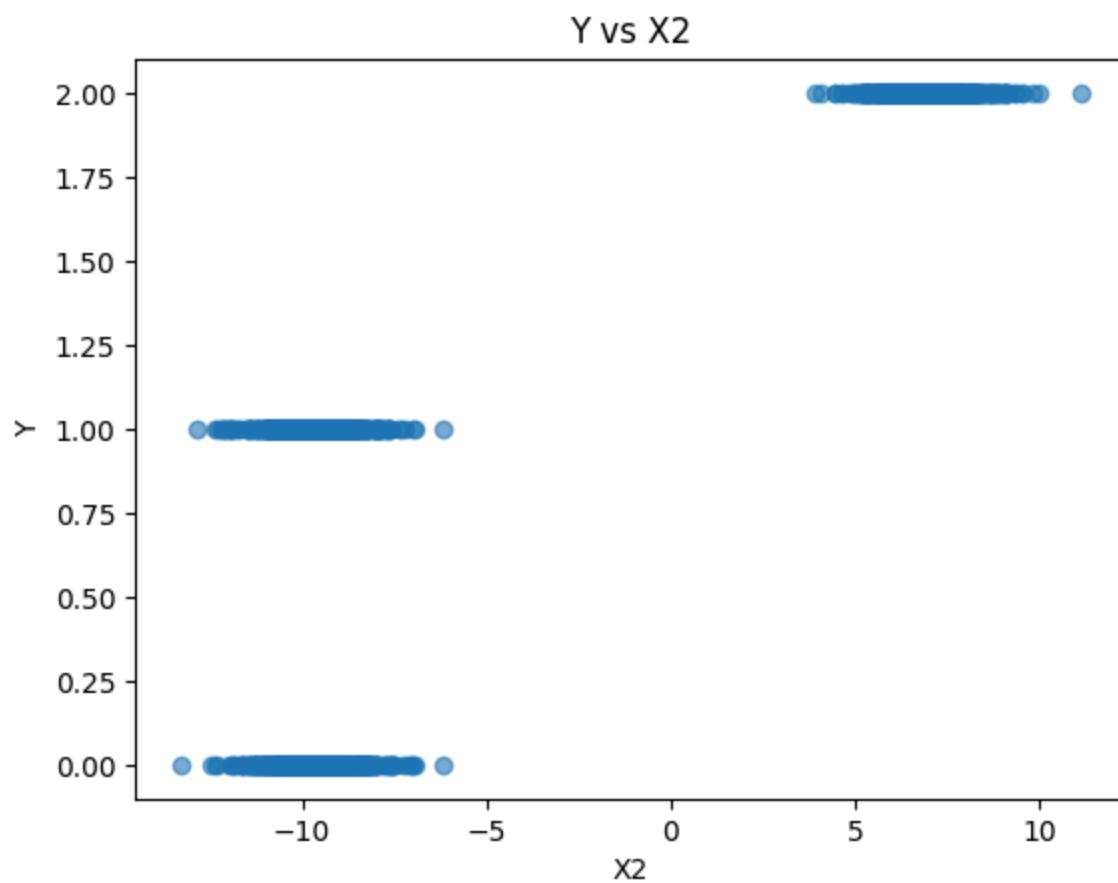
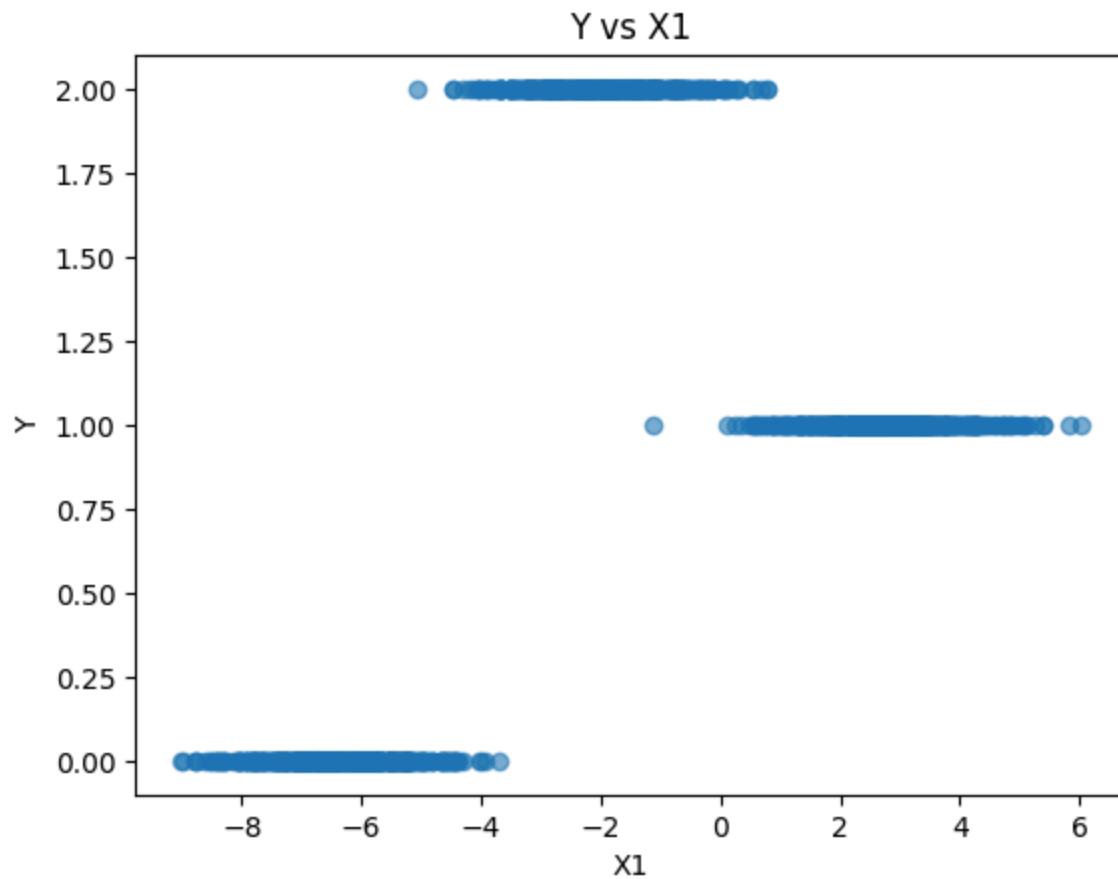
In [11]: df\_class.shape

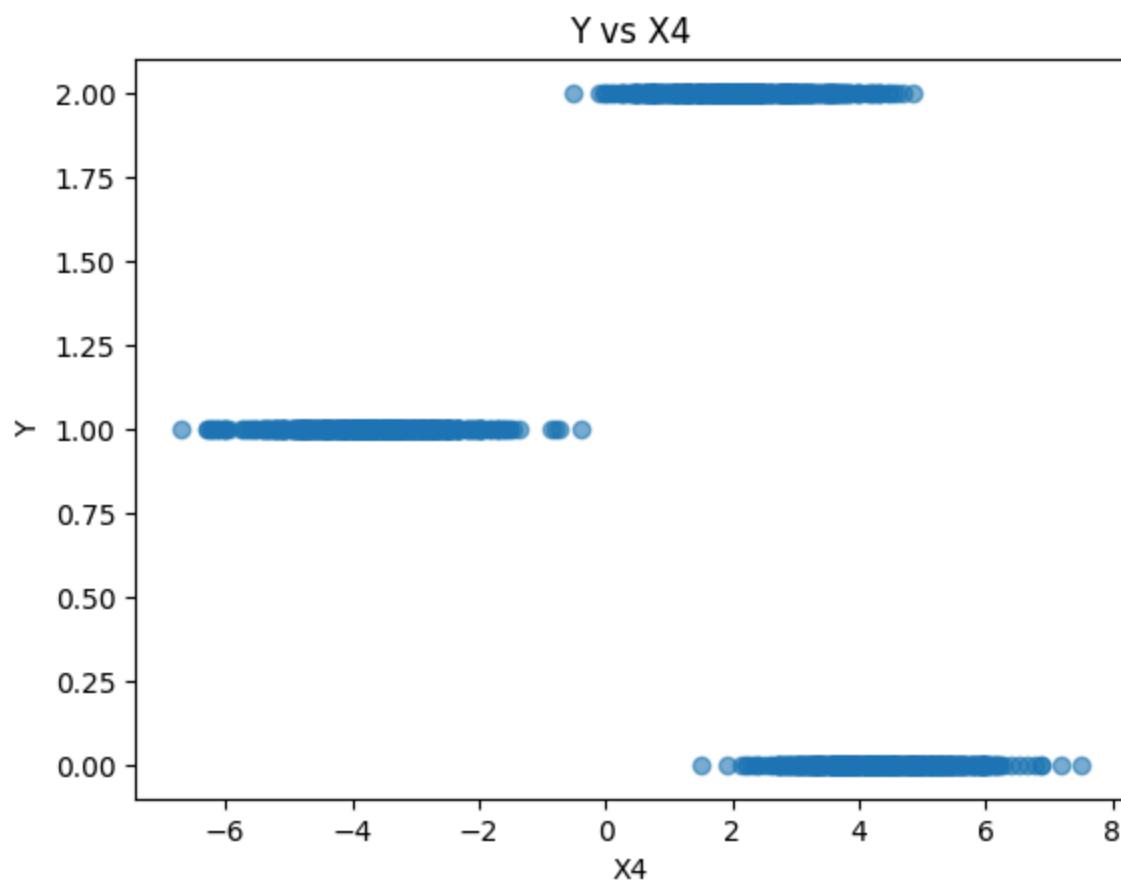
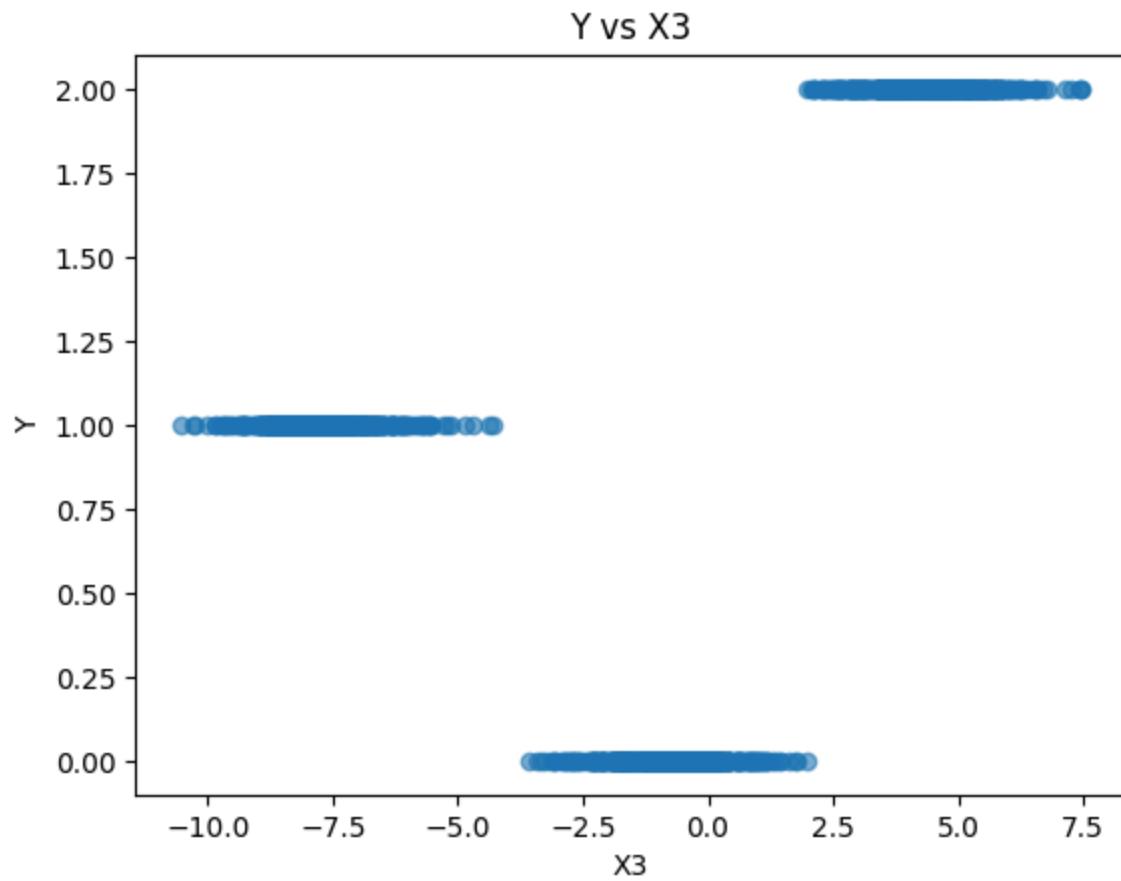
Out[11]: (1485, 16)

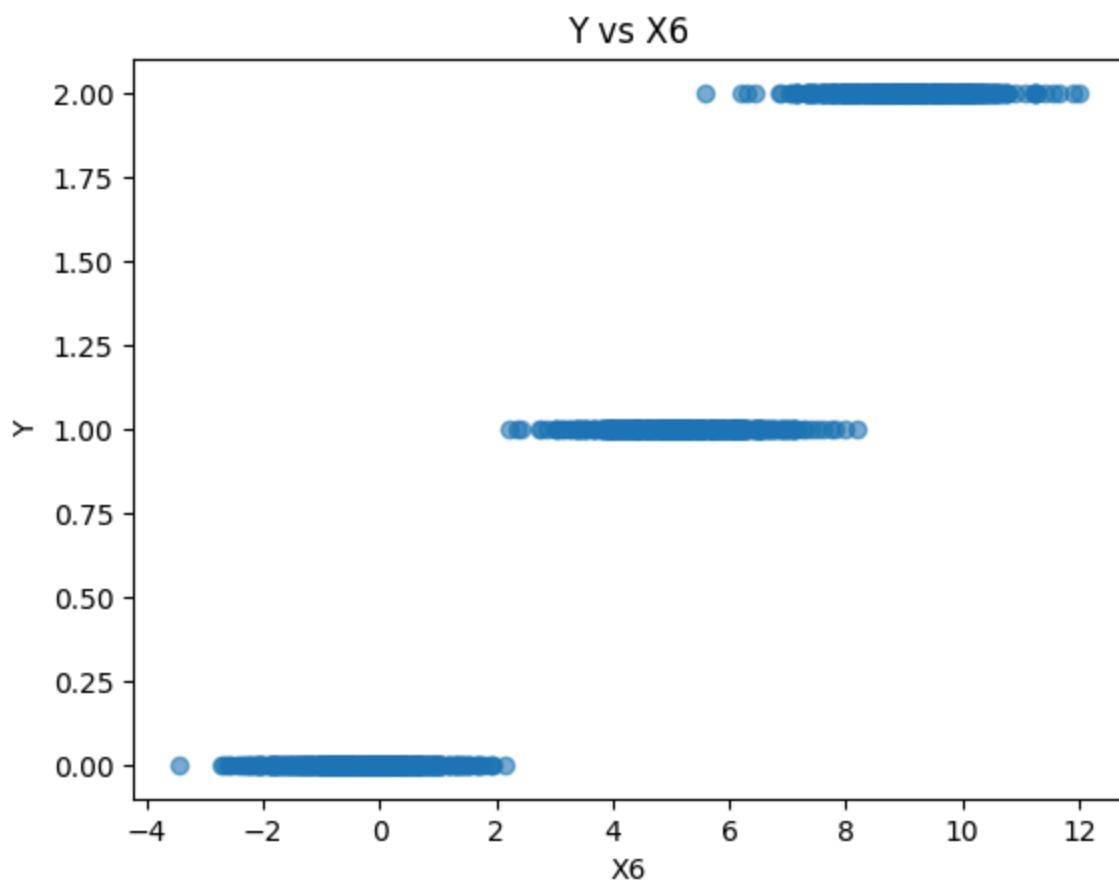
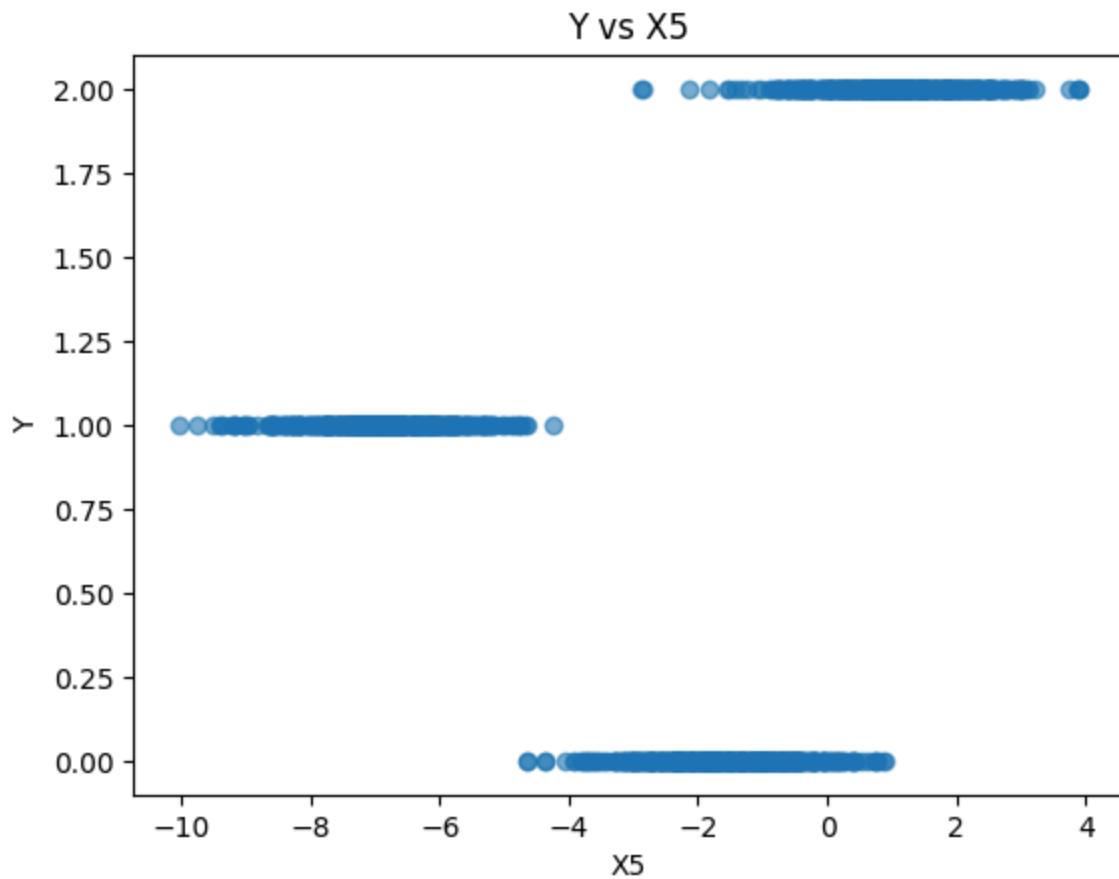
In [12]: num\_cols\_class = df\_class.select\_dtypes(include=[np.number]).columns  
df\_class[num\_cols\_class].hist(bins=20, figsize=(10, 8))  
plt.tight\_layout()  
plt.show()

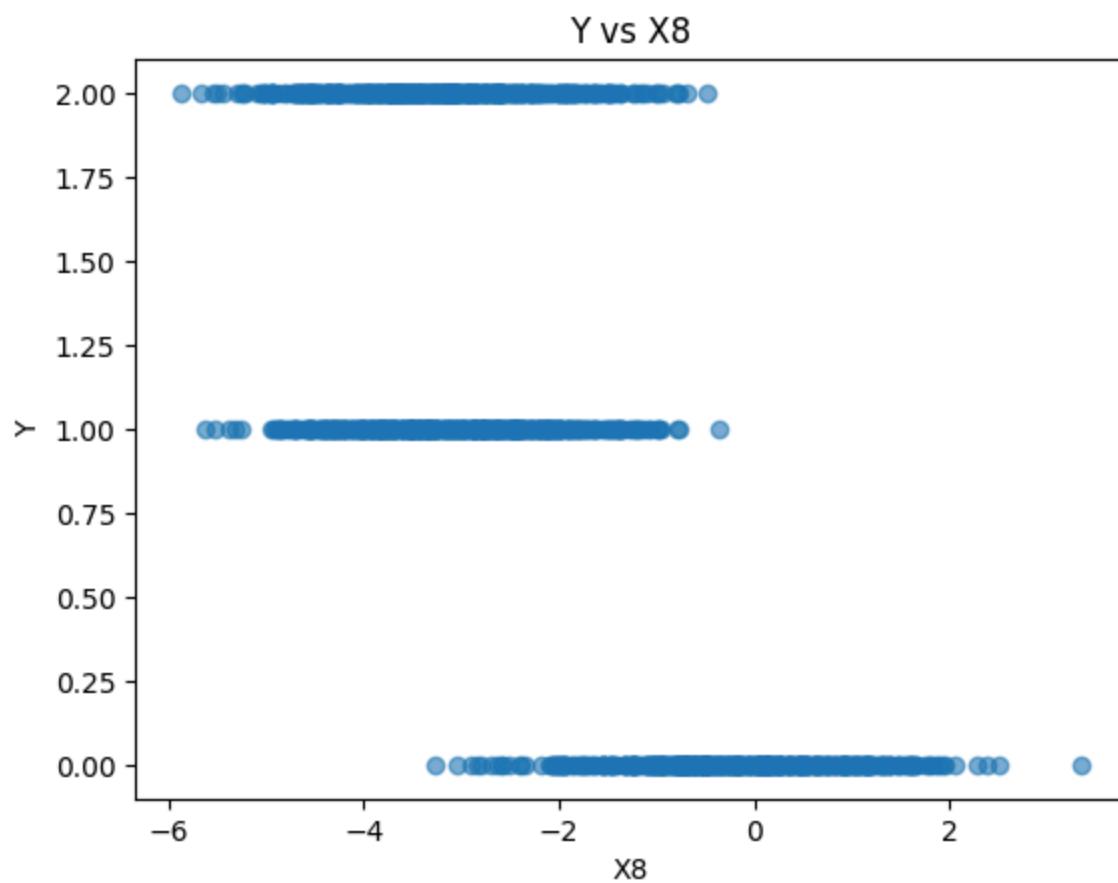
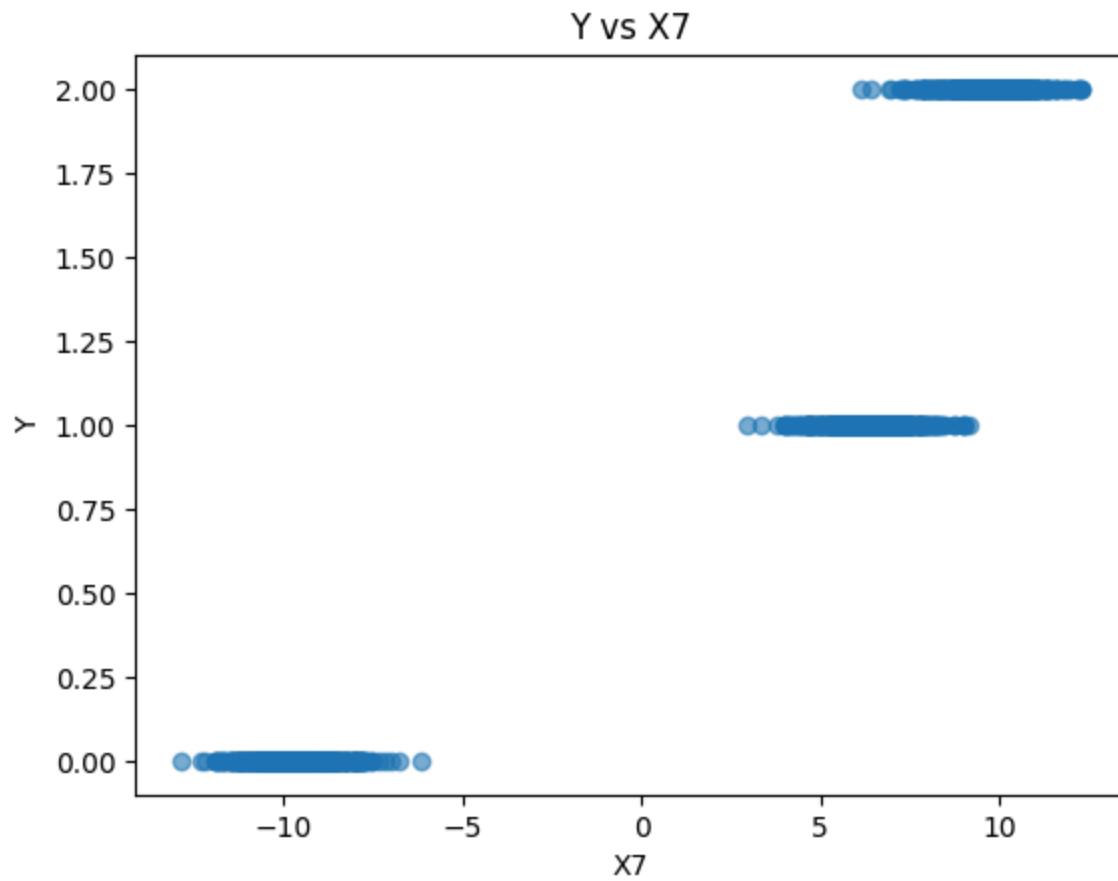


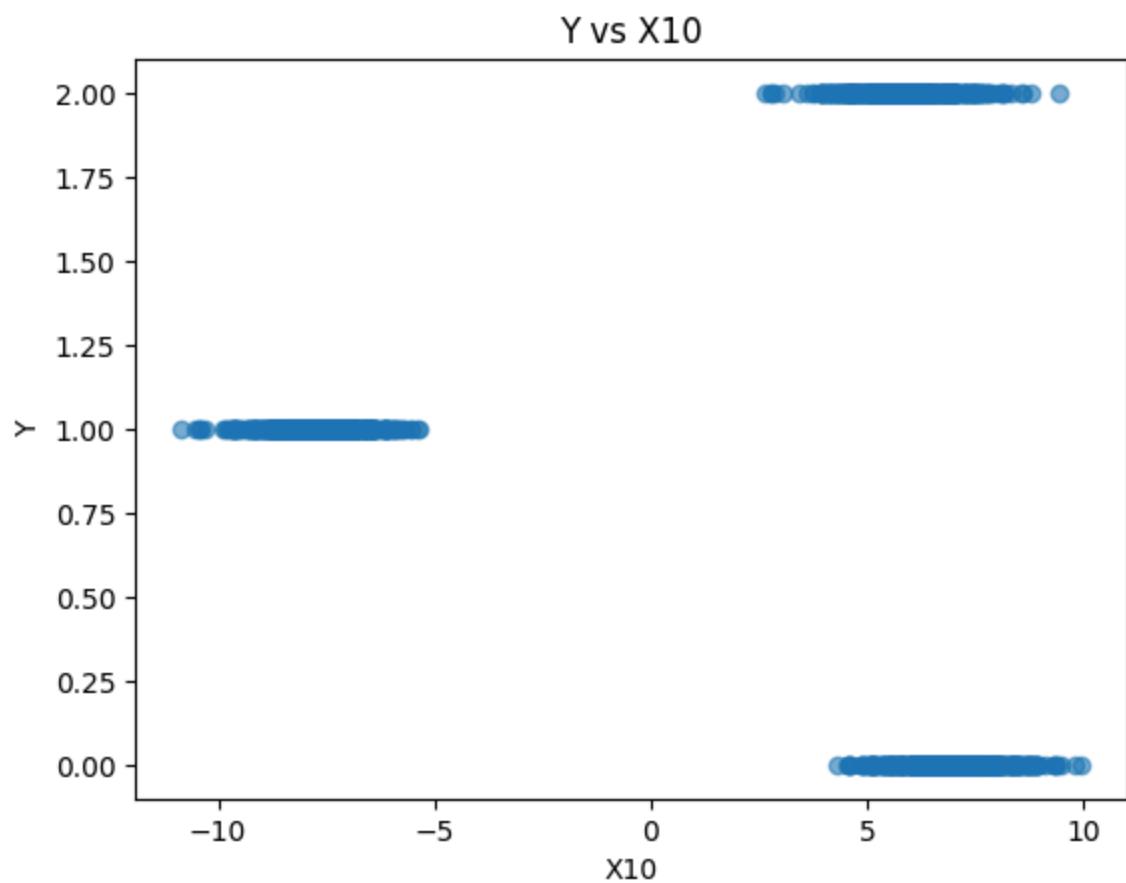
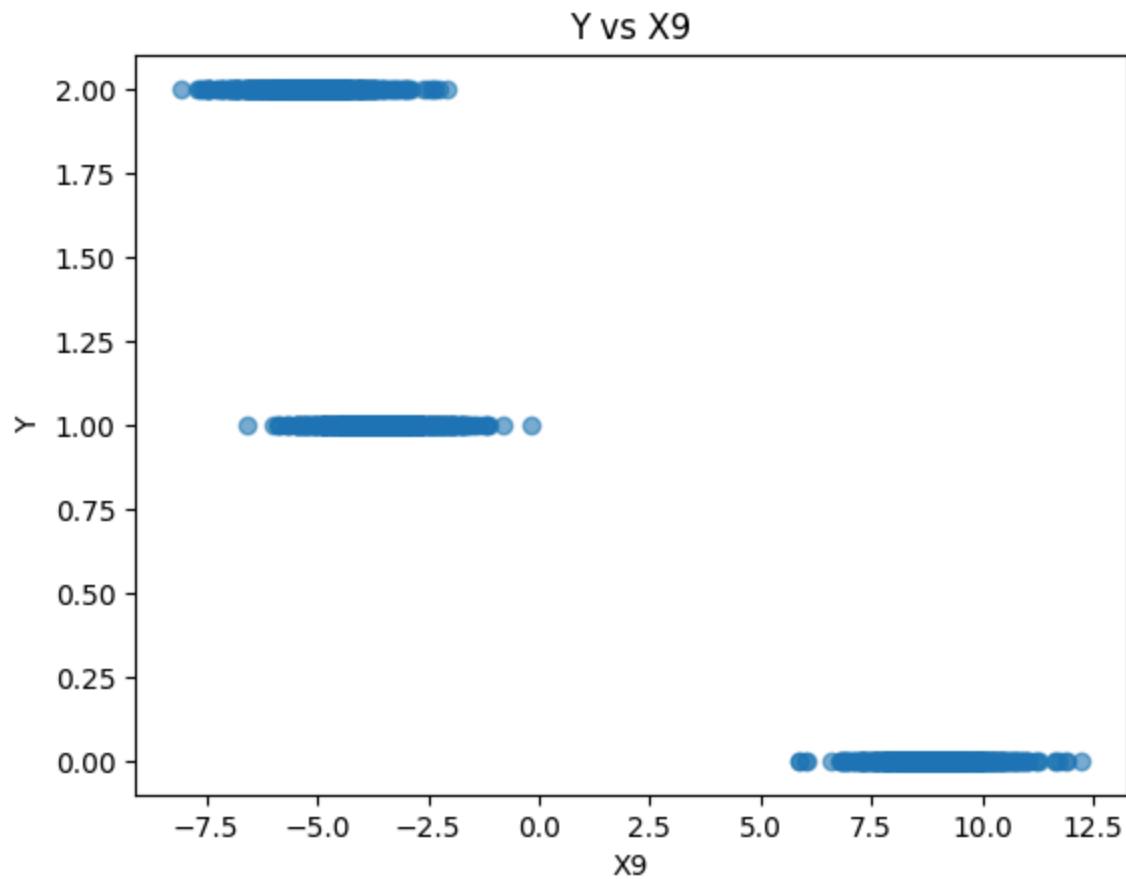
```
In [13]: for col in [c for c in num_cols_class.to_list() if c != "Y"]:
    plt.figure()
    plt.scatter(df_class[col], df_class["Y"], alpha=0.6)
    plt.title(f"Y vs {col}")
    plt.xlabel(col)
    plt.ylabel("Y")
    plt.show()
```

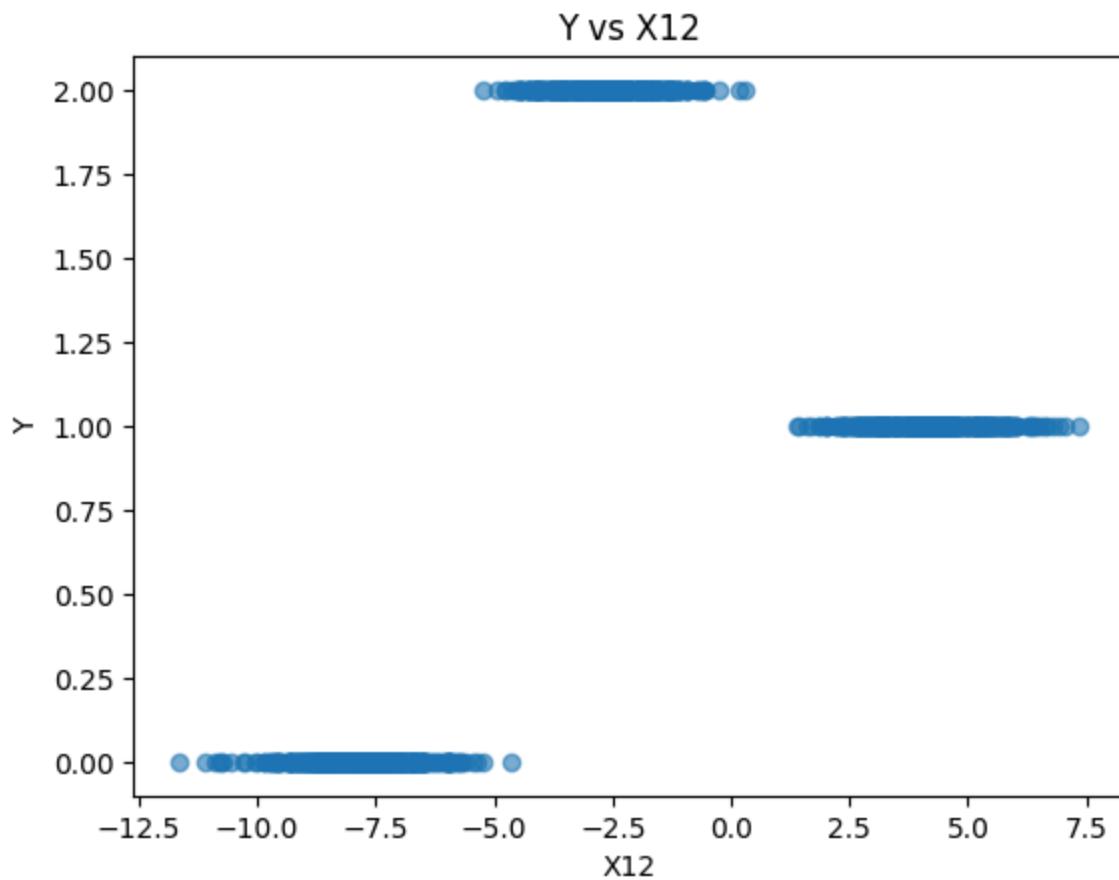
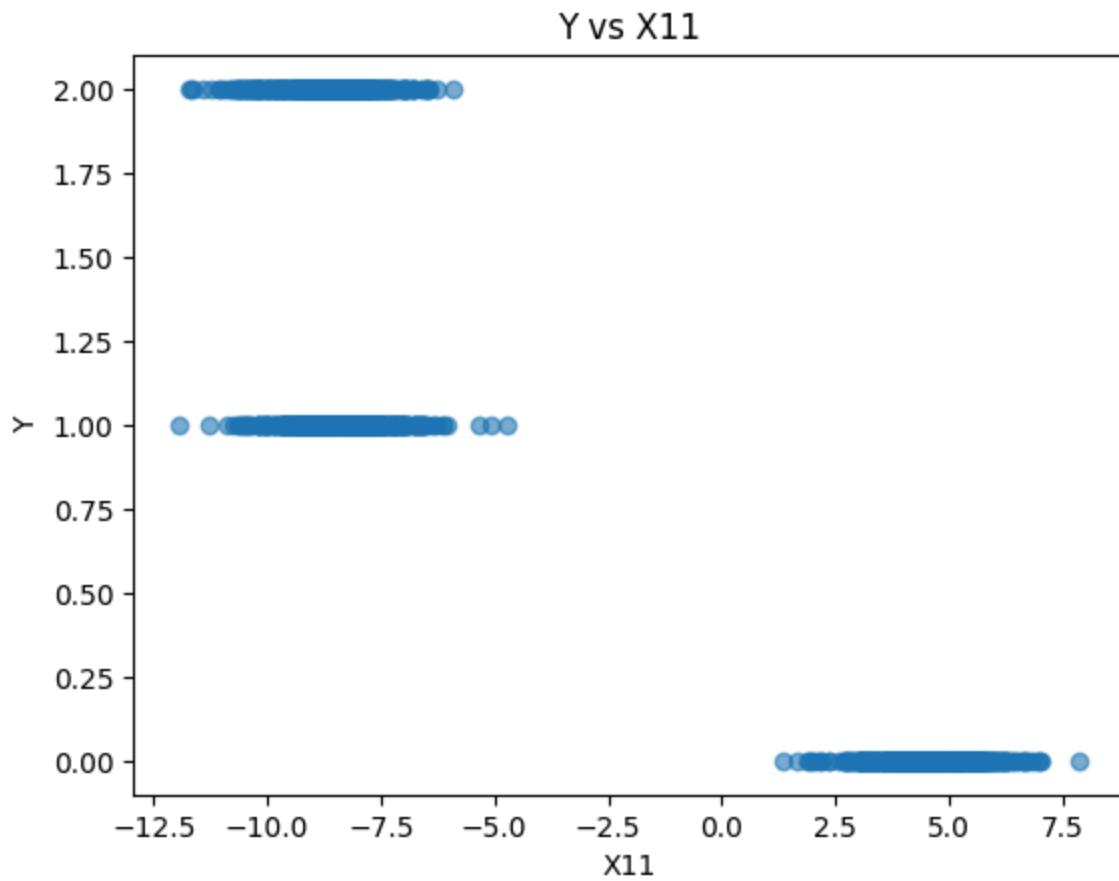


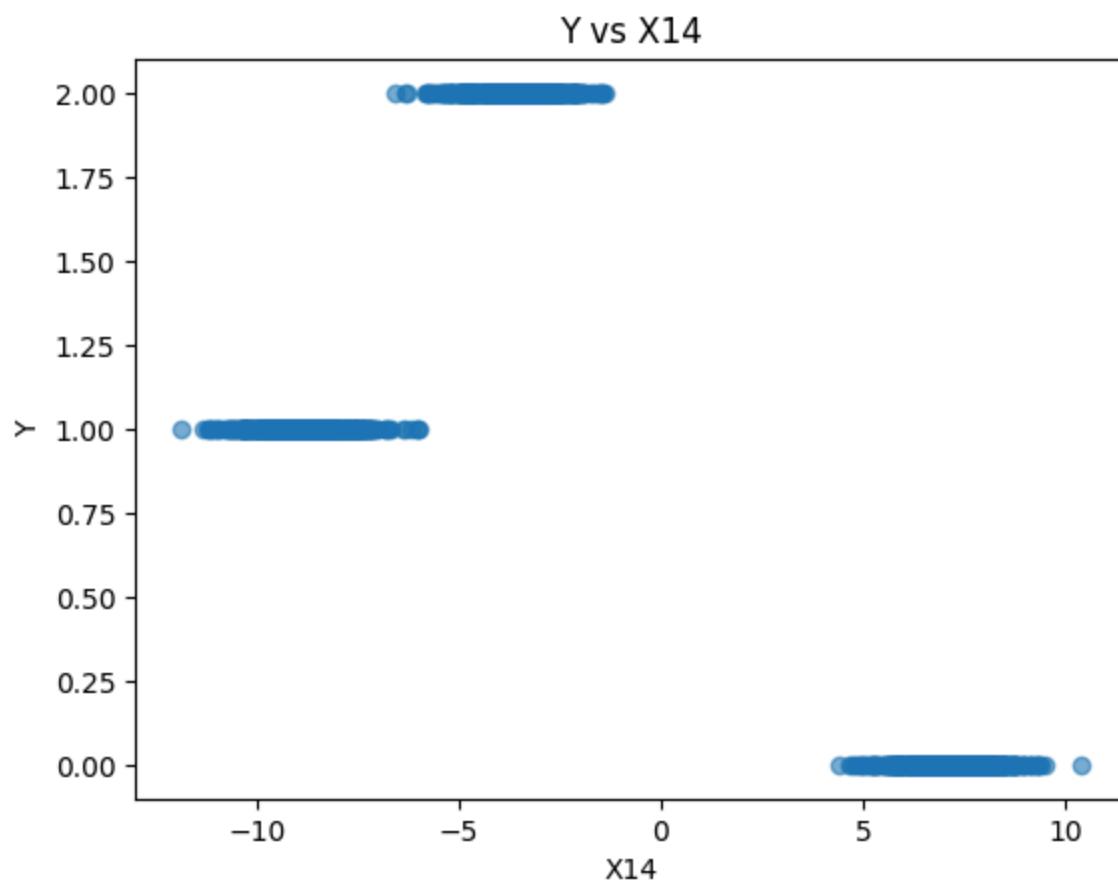
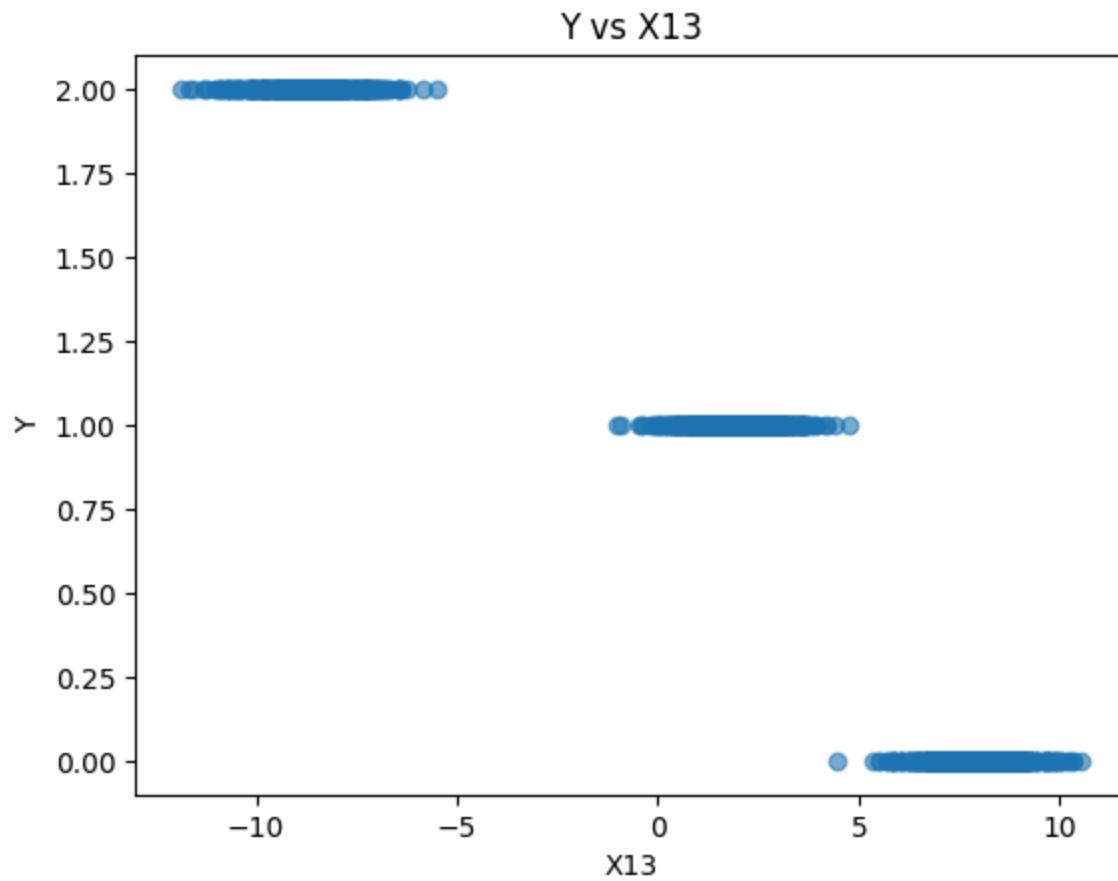


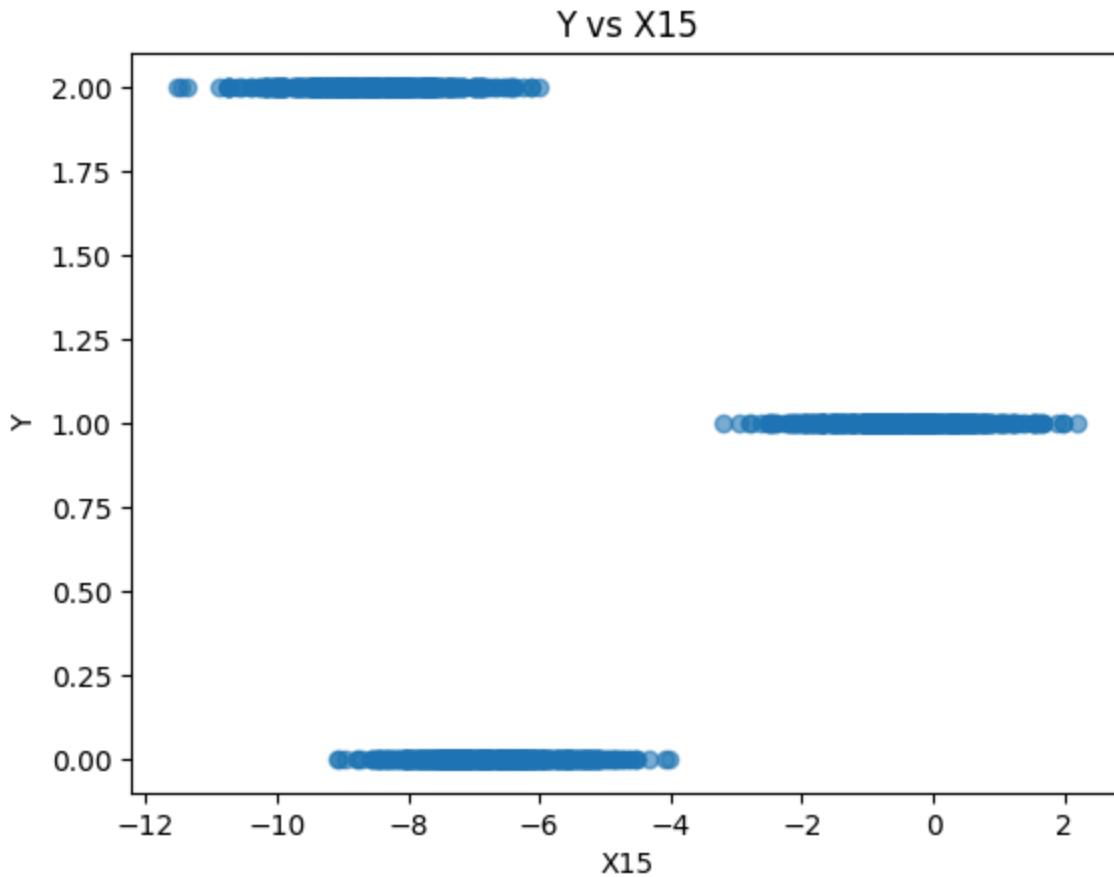










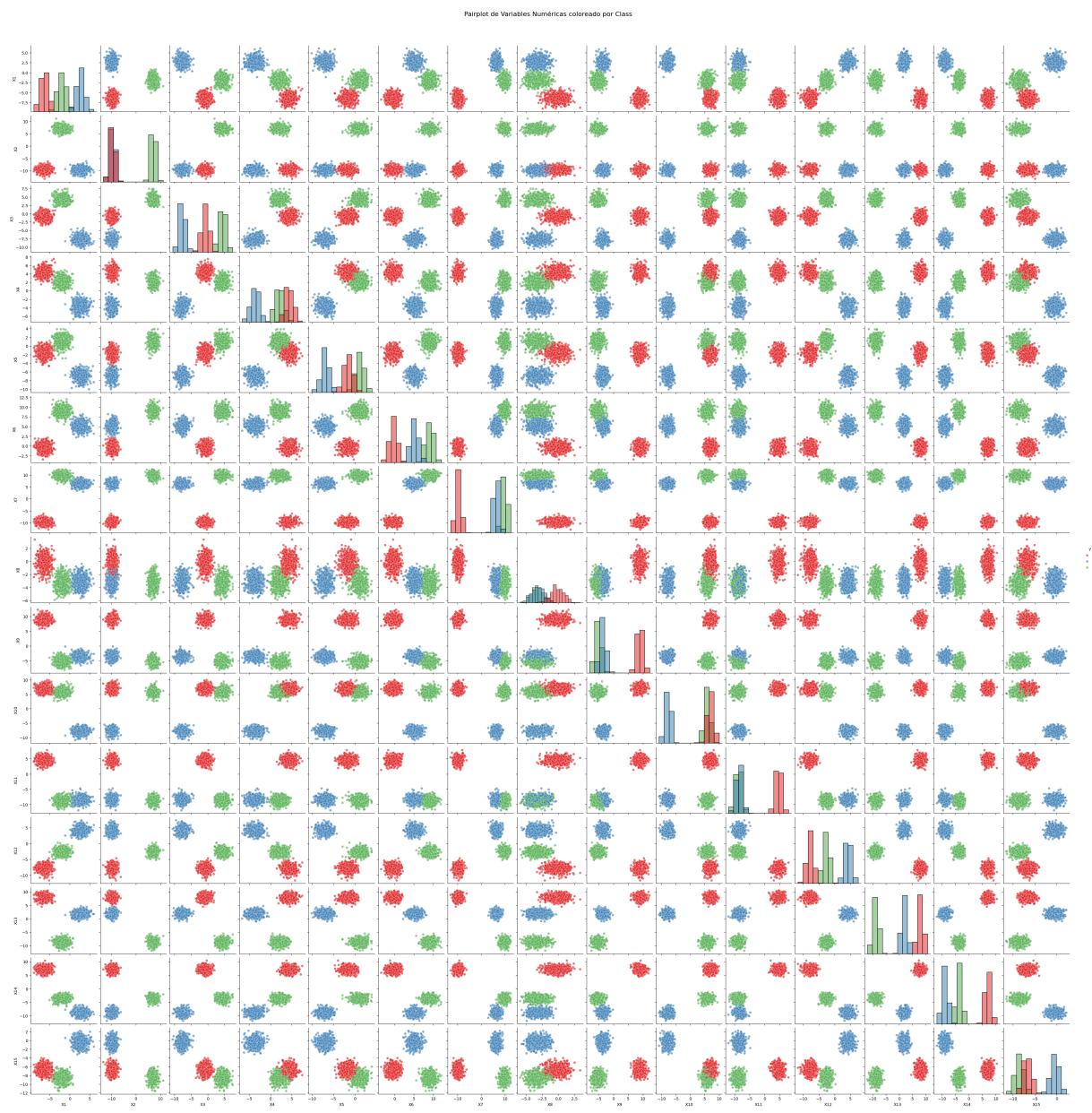


```
In [14]: from pandas.plotting import scatter_matrix

numeric_cols_class = num_cols_class.drop("Y")
class_vals = df_class["Y"]

import seaborn as sns

sns.pairplot(
    df_class[numeric_cols_class.tolist() + ["Y"]],
    hue="Y",
    diag_kind="hist",
    palette="Set1",
    plot_kws={"alpha": 0.6, "marker": "o"},
)
plt.suptitle("Pairplot de Variables Numéricas coloreado por Class", y=1.02,
plt.show()
```



```
In [15]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

corr = df_class[numeric_cols_class.tolist() + ["Y"]].corr()

sns.set_theme(style="white")

mask = np.triu(np.ones_like(corr, dtype=bool))

n = corr.shape[0]
fig, ax = plt.subplots(figsize=(1.2 * n, 1.2 * n))

sns.heatmap(
    corr,
    mask=mask,
    annot=True,
    fmt=".2f",
    
```

```

        cmap="RdBu_r",
        center=0,
        linewidths=1.0,
        square=True,
        annot_kws={"size": 8},
        cbar_kws={"shrink": 0.6, "label": "Correlación"},
        vmin=-1,
        vmax=1,
        ax=ax,
    )

ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha="right", fontsize=10)
ax.set_yticklabels(ax.get_yticklabels(), rotation=0, fontsize=10)

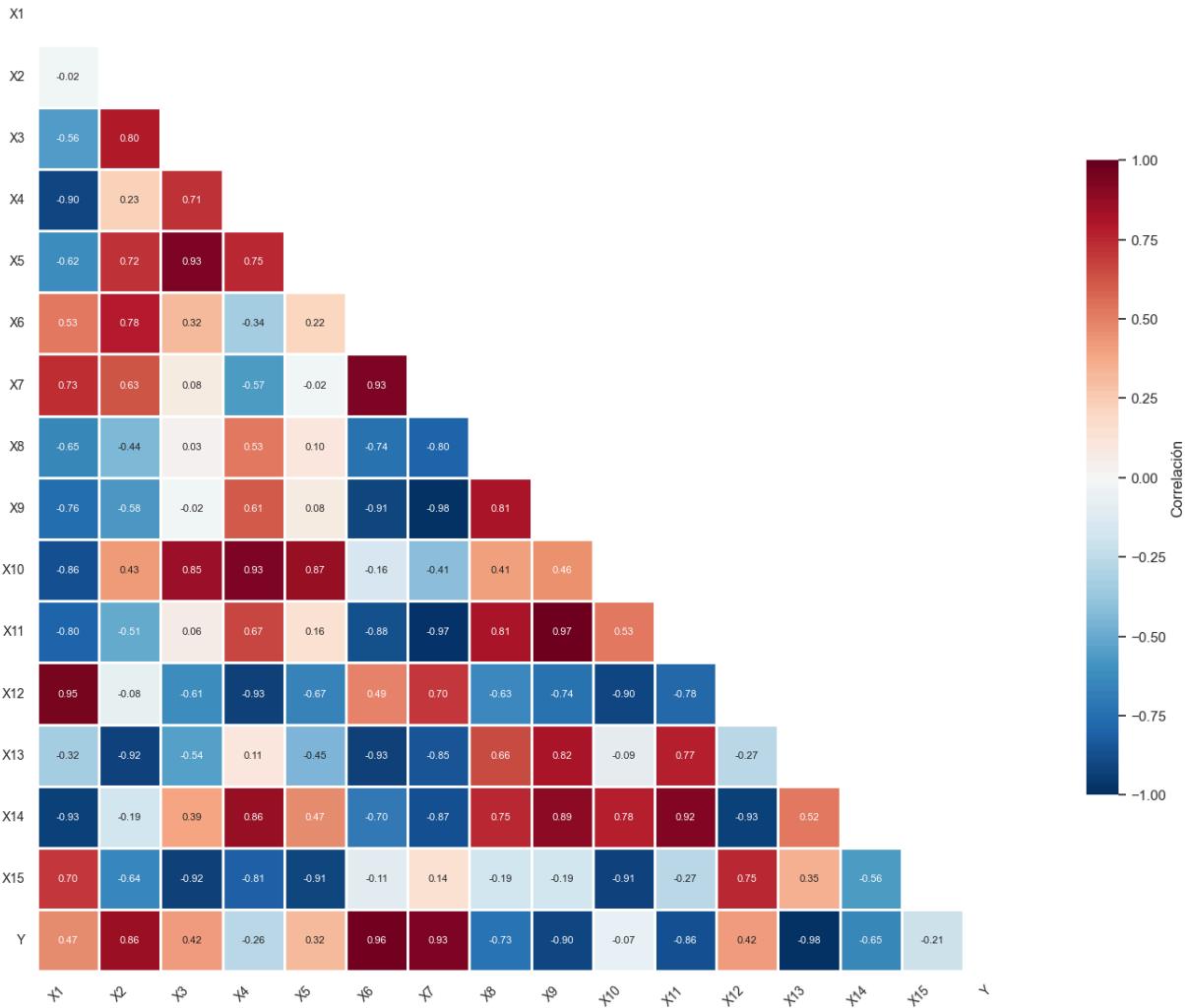
plt.subplots_adjust(bottom=0.3, left=0.2, top=0.9)

ax.set_title("Mapa de Calor de Correlaciones", fontsize=18, pad=20)

plt.show()

```

Mapa de Calor de Correlaciones



# Regression

```
In [16]: df_reg = pd.read_csv("../data/regression.csv")
```

```
In [17]: df_reg.shape
```

```
Out[17]: (1500, 13)
```

```
In [18]: df_reg.dtypes
```

```
Out[18]: X1      float64  
X2      float64  
X3      float64  
X4      float64  
X5      float64  
X6      float64  
X7      float64  
X8      float64  
X9      float64  
X10     float64  
X11     float64  
X12     float64  
Y       float64  
dtype: object
```

```
In [19]: df_reg.head()
```

```
Out[19]:
```

	X1	X2	X3	X4	X5	X6	
0	5468.737967	1.703822	-8613.627601	0.284719	-1.585462	-0.554724	-1.0230
1	-13104.449210	-1.019682	23061.174510	-0.820498	-0.265240	2.099421	-0.3163
2	6204.272157	0.734018	133965.765100	-0.710429	0.574379	0.511245	0.2981
3	23979.672370	-1.777481	45514.777250	-0.148987	-0.006399	0.464724	0.3568
4	-2390.952073	-2.418137	47102.194950	0.384587	0.359033	0.178320	-1.0630

```
In [20]: df_reg.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1500 entries, 0 to 1499
Data columns (total 13 columns):
 #   Column  Non-Null Count  Dtype  
--- 
 0   X1      1499 non-null   float64
 1   X2      1500 non-null   float64
 2   X3      1495 non-null   float64
 3   X4      1497 non-null   float64
 4   X5      1495 non-null   float64
 5   X6      1498 non-null   float64
 6   X7      1499 non-null   float64
 7   X8      1498 non-null   float64
 8   X9      1500 non-null   float64
 9   X10     1496 non-null   float64
 10  X11     1498 non-null   float64
 11  X12     1499 non-null   float64
 12  Y       1499 non-null   float64
dtypes: float64(13)
memory usage: 152.5 KB
```

In [21]: `df_reg.describe()`

	X1	X2	X3	X4	X5
<b>count</b>	1499.000000	1500.000000	1495.000000	1497.000000	1495.000000
<b>mean</b>	-298.970046	-0.000182	2505.655649	0.019644	-0.037554
<b>std</b>	12664.656465	1.036299	95364.185732	1.017798	0.975190
<b>min</b>	-39838.521840	-3.710679	-336700.253400	-2.796599	-3.282985
<b>25%</b>	-9433.010398	-0.677534	-60485.305915	-0.694018	-0.698821
<b>50%</b>	-496.180129	0.021941	1748.071119	0.027749	-0.034482
<b>75%</b>	8318.518644	0.705290	70803.472755	0.728685	0.615946
<b>max</b>	44114.014510	3.156200	312030.703900	3.344649	3.515152

In [22]: `missing_summary = df_reg.isnull().sum().to_frame(name="n_missing")  
missing_summary["pct_missing"] = missing_summary["n_missing"] / len(df_reg)  
print("\nResumen de valores faltantes:")  
missing_summary`

Resumen de valores faltantes:

Out[22]:

	n_missing	pct_missing
X1	1	0.066667
X2	0	0.000000
X3	5	0.333333
X4	3	0.200000
X5	5	0.333333
X6	2	0.133333
X7	1	0.066667
X8	2	0.133333
X9	0	0.000000
X10	4	0.266667
X11	2	0.133333
X12	1	0.066667
Y	1	0.066667

In [23]: df\_reg.shape

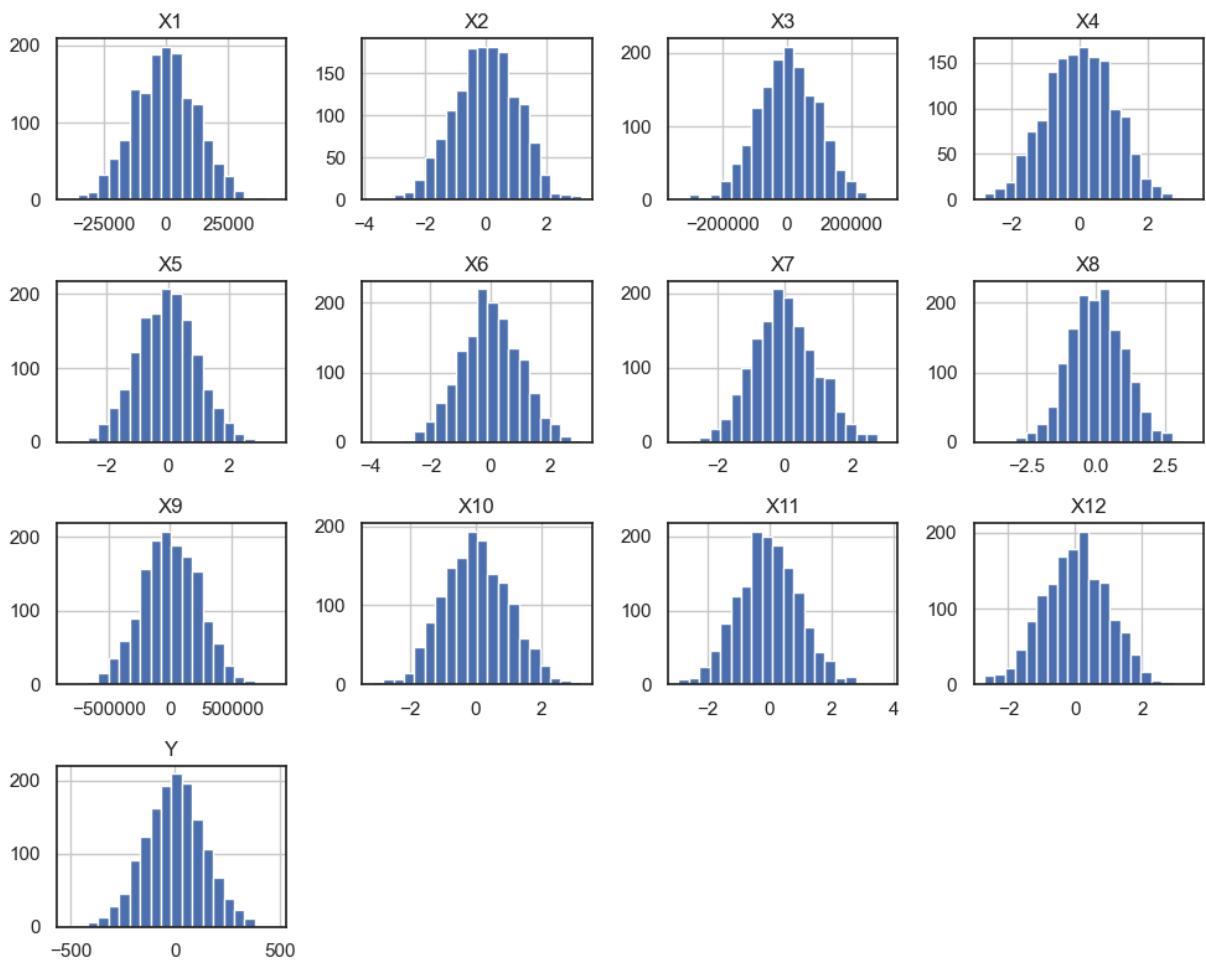
Out[23]: (1500, 13)

In [24]: df\_reg = df\_reg.dropna()

In [25]: df\_reg.shape

Out[25]: (1473, 13)

In [26]: num\_cols\_reg = df\_reg.select\_dtypes(include=[np.number]).columns  
df\_reg[num\_cols\_reg].hist(bins=20, figsize=(10, 8))  
plt.tight\_layout()  
plt.show()



```
In [27]: from pandas.plotting import scatter_matrix

numeric_cols_reg = num_cols_reg.drop("Y")
reg_vals = df_reg["Y"]

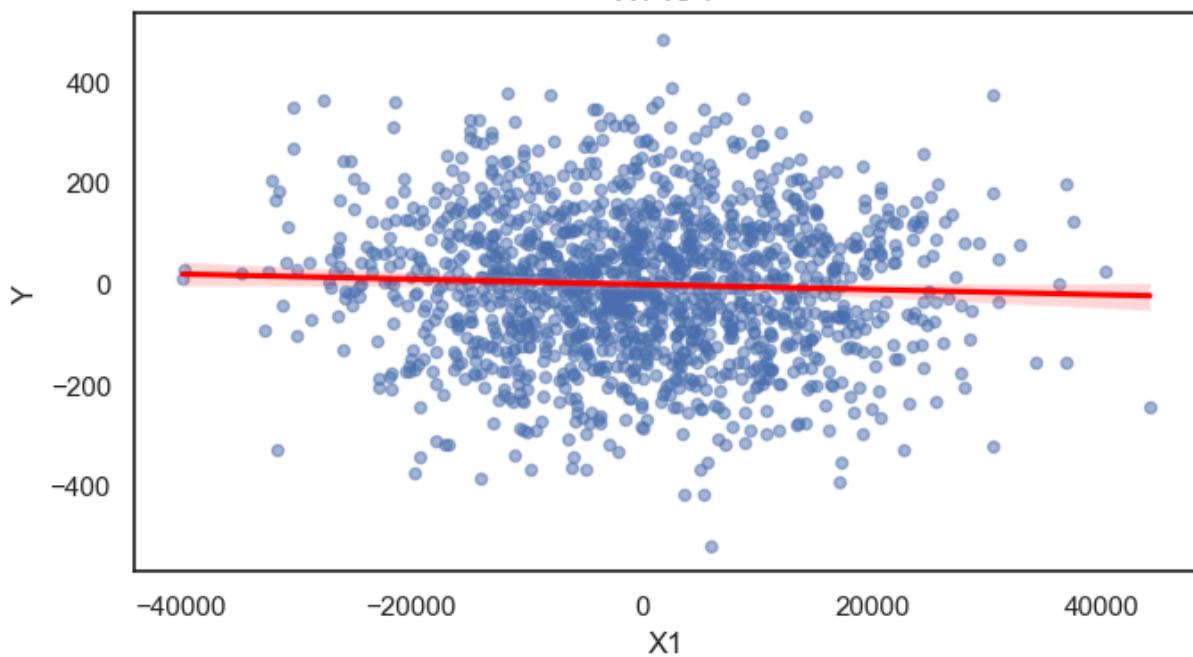
import matplotlib.pyplot as plt
import seaborn as sns

n_feats = len(numeric_cols_reg)
fig, axes = plt.subplots(n_feats, 1, figsize=(7, 4 * n_feats))

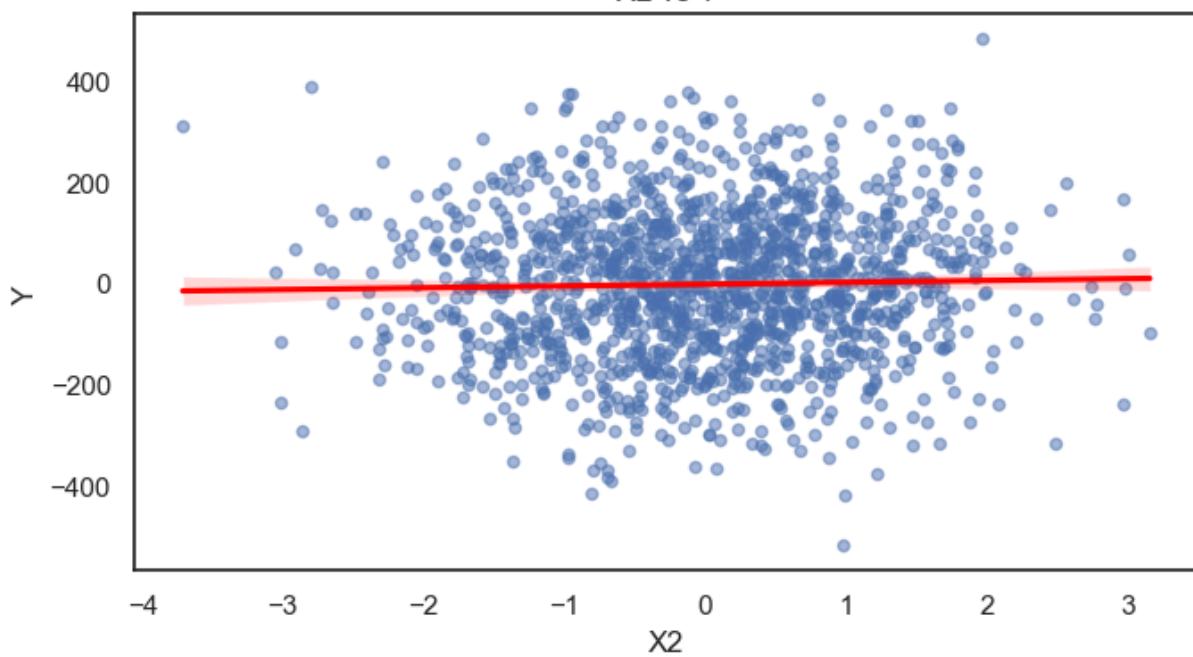
for ax, feat in zip(axes, numeric_cols_reg):
    sns.regplot(
        x=feat,
        y="Y",
        data=df_reg,
        ax=ax,
        scatter_kws={"alpha": 0.5, "s": 20},
        line_kws={"color": "red"},
    )
    ax.set_title(f"{feat} vs Y")
    ax.set_xlabel(feat)
    ax.set_ylabel("Y")

plt.tight_layout()
plt.show()
```

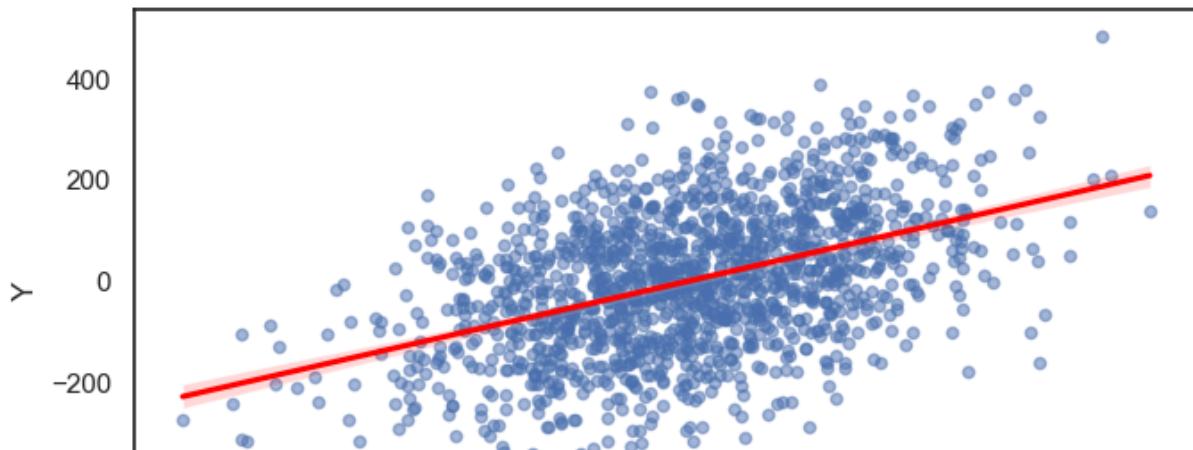
X1 vs Y



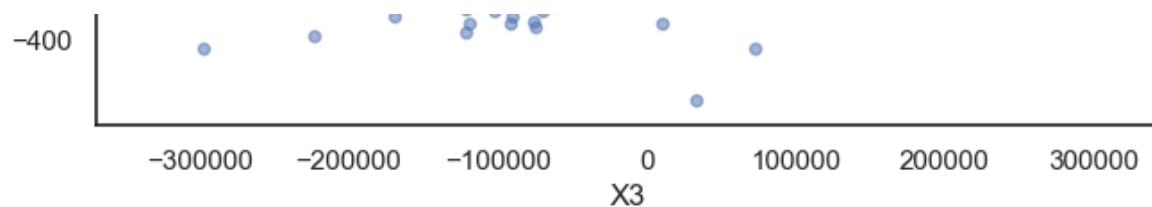
X2 vs Y



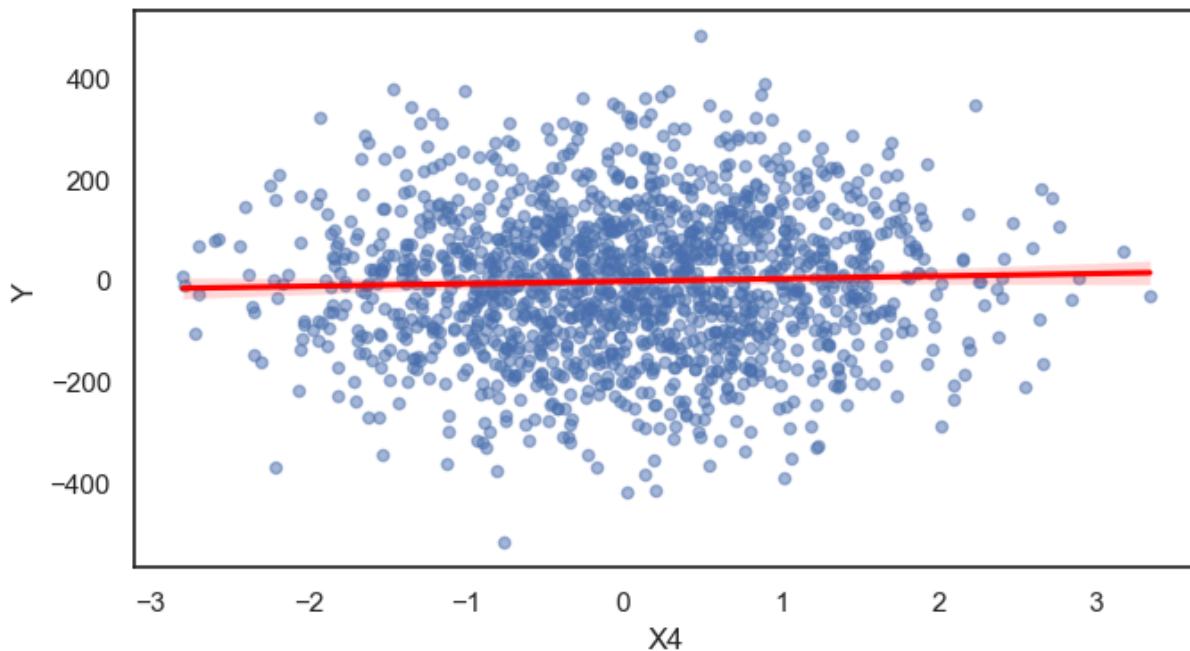
X3 vs Y



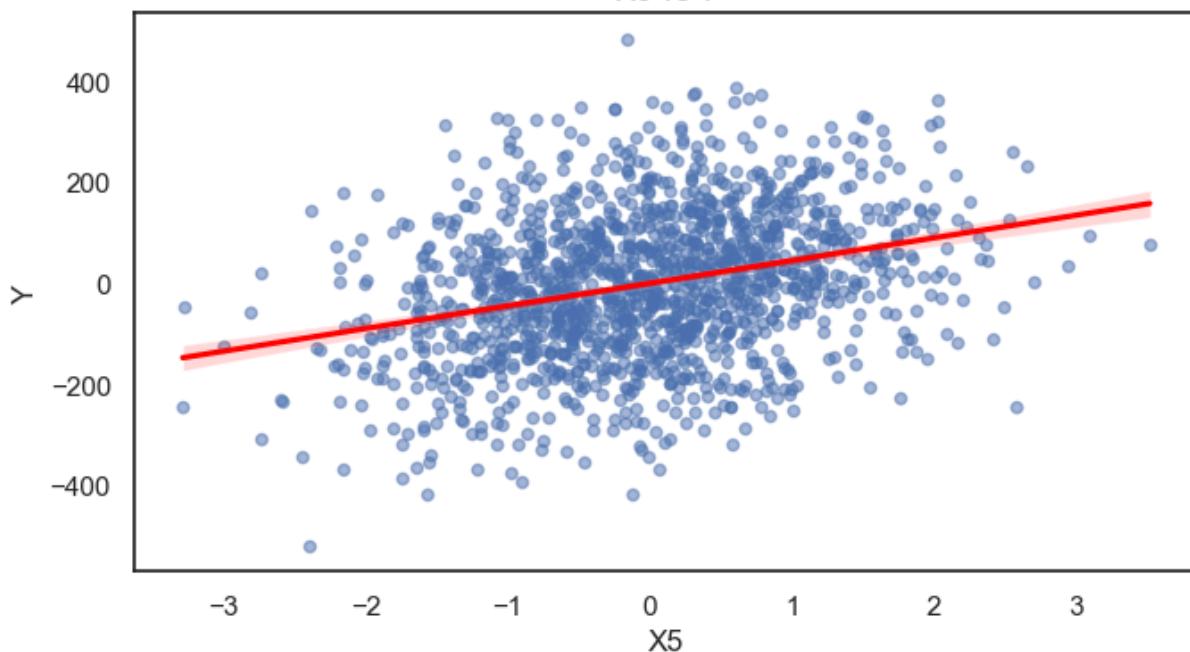
## EDA



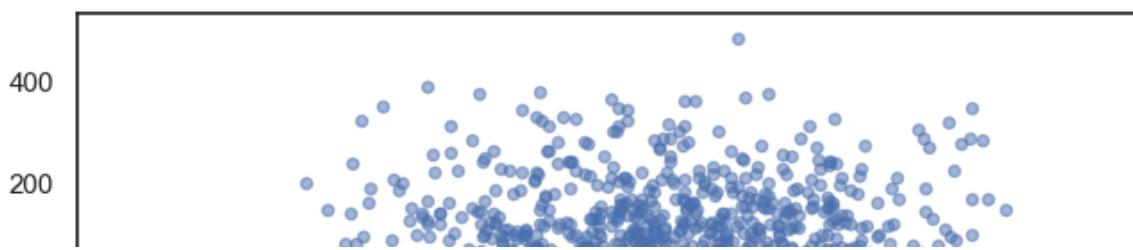
X4 vs Y

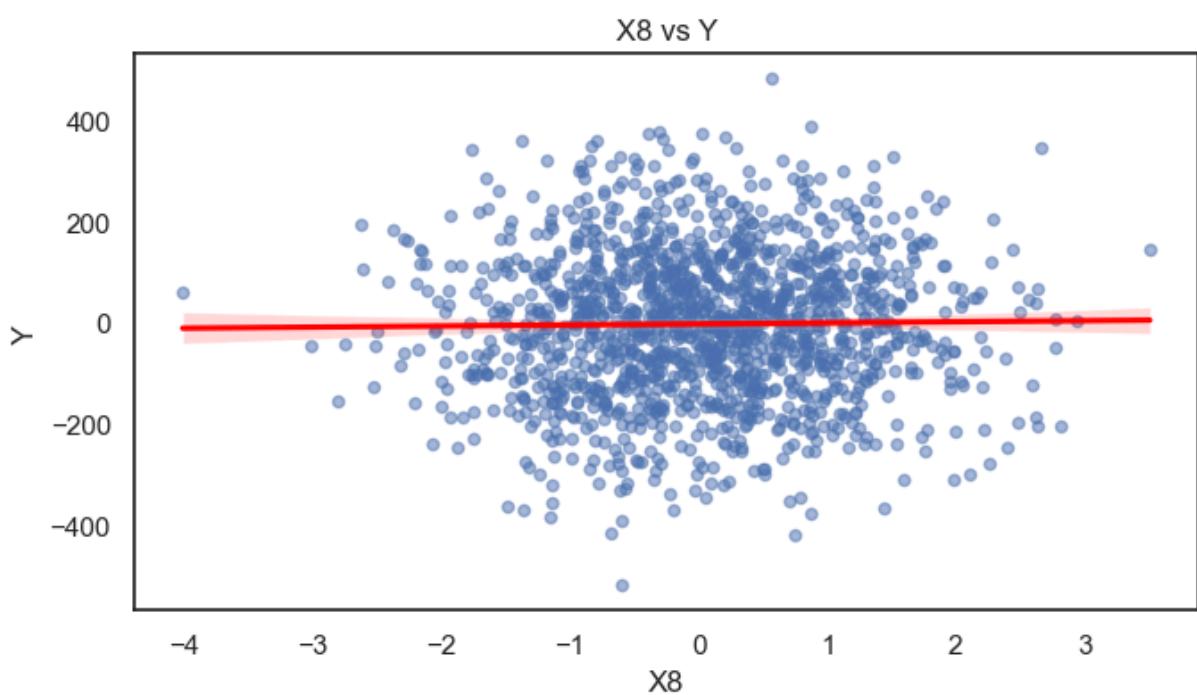
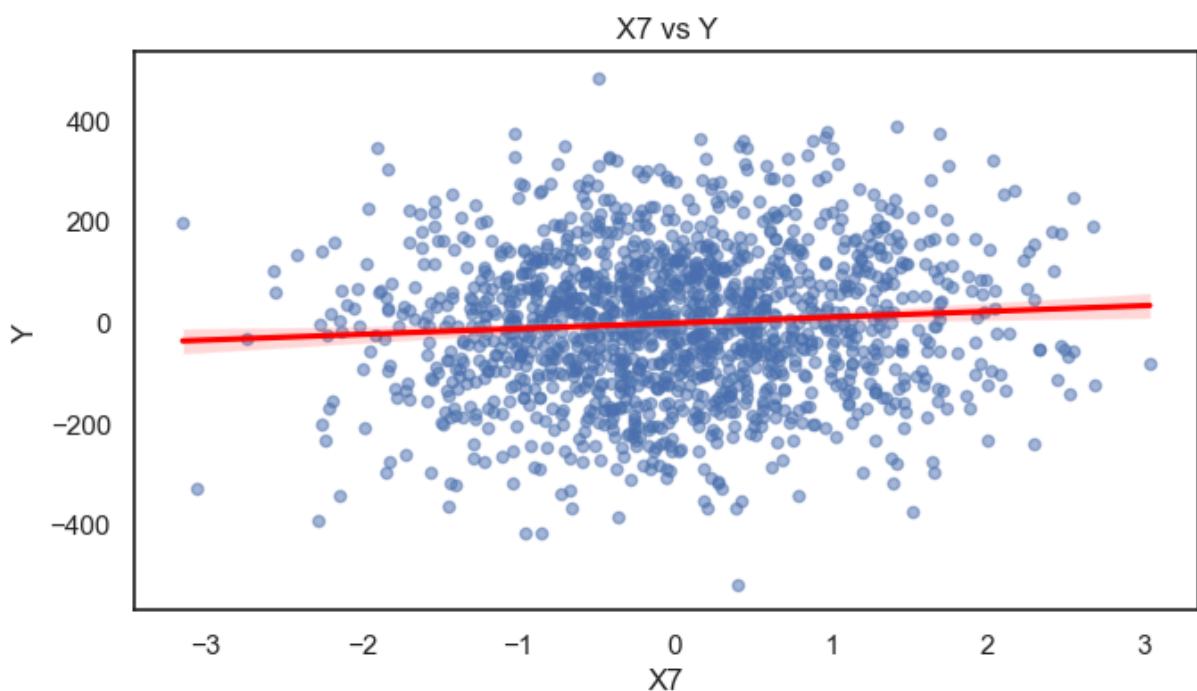
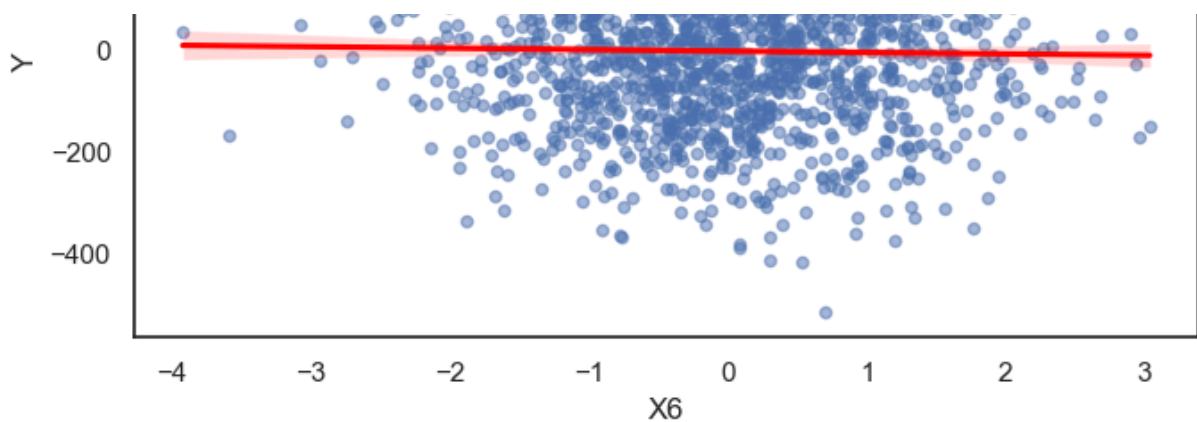


X5 vs Y

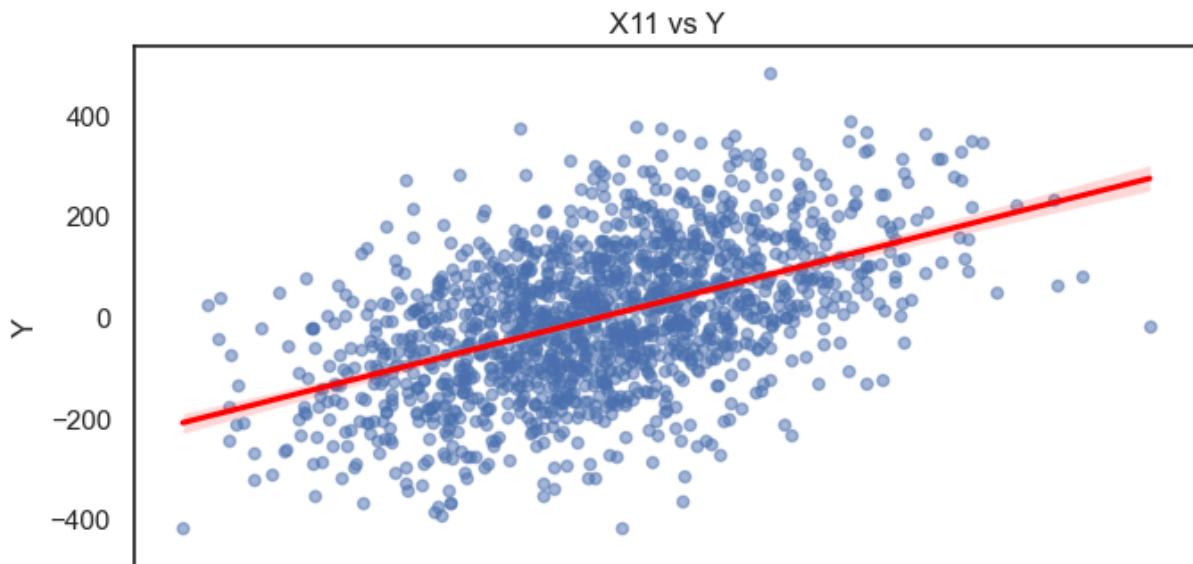
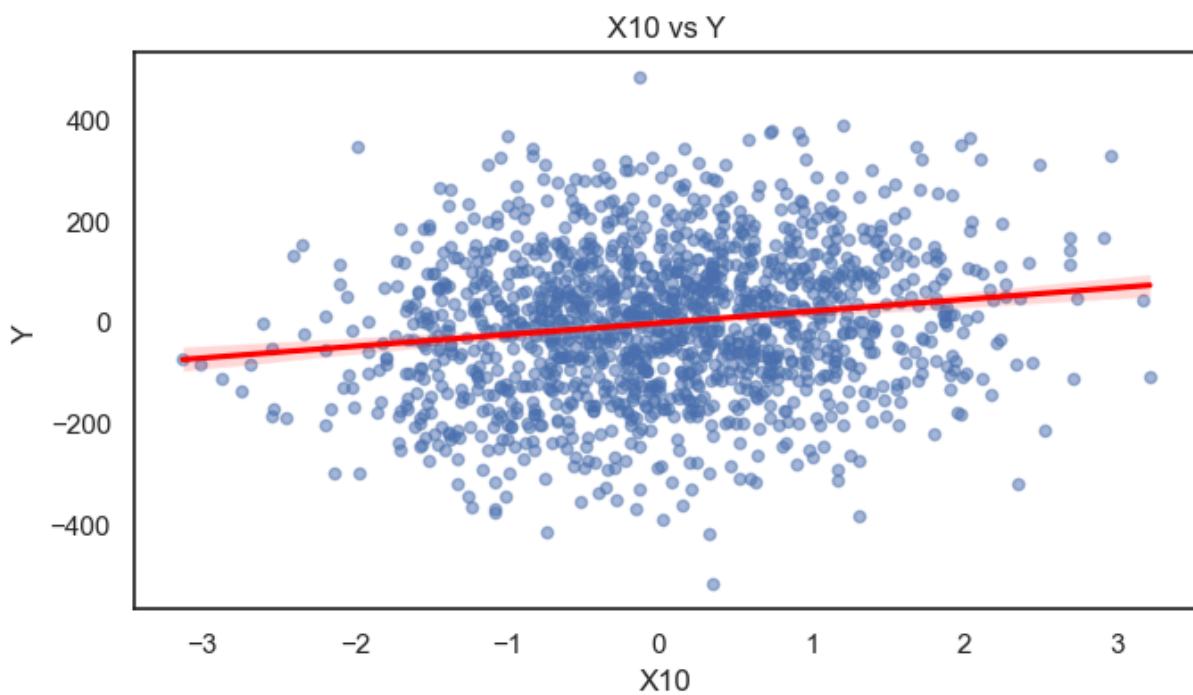
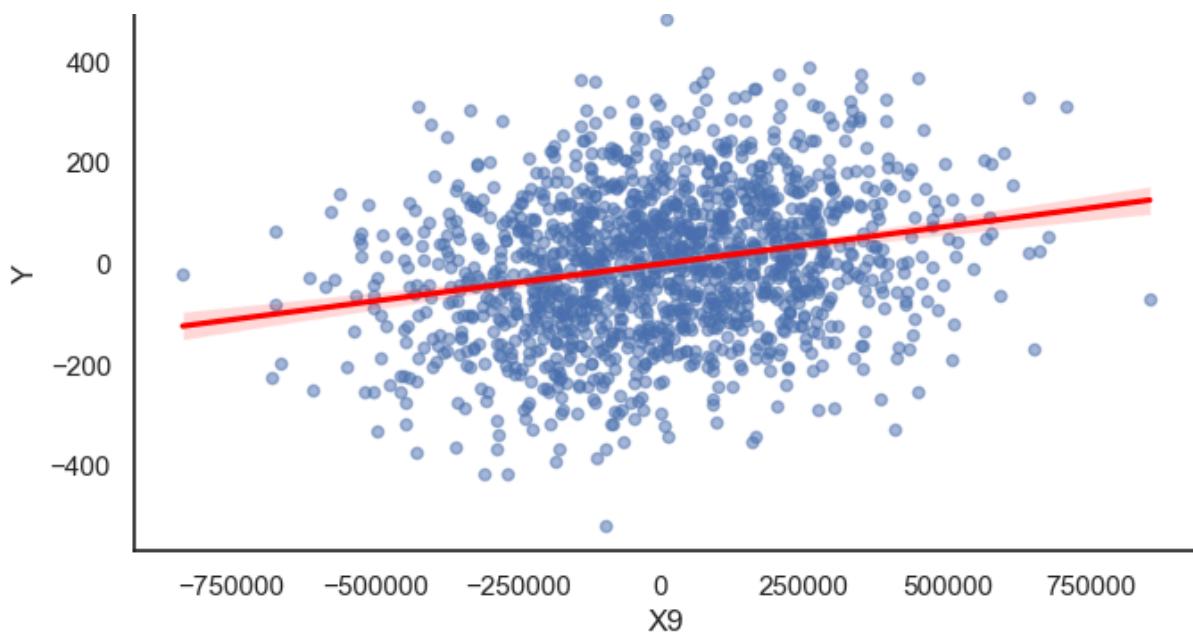


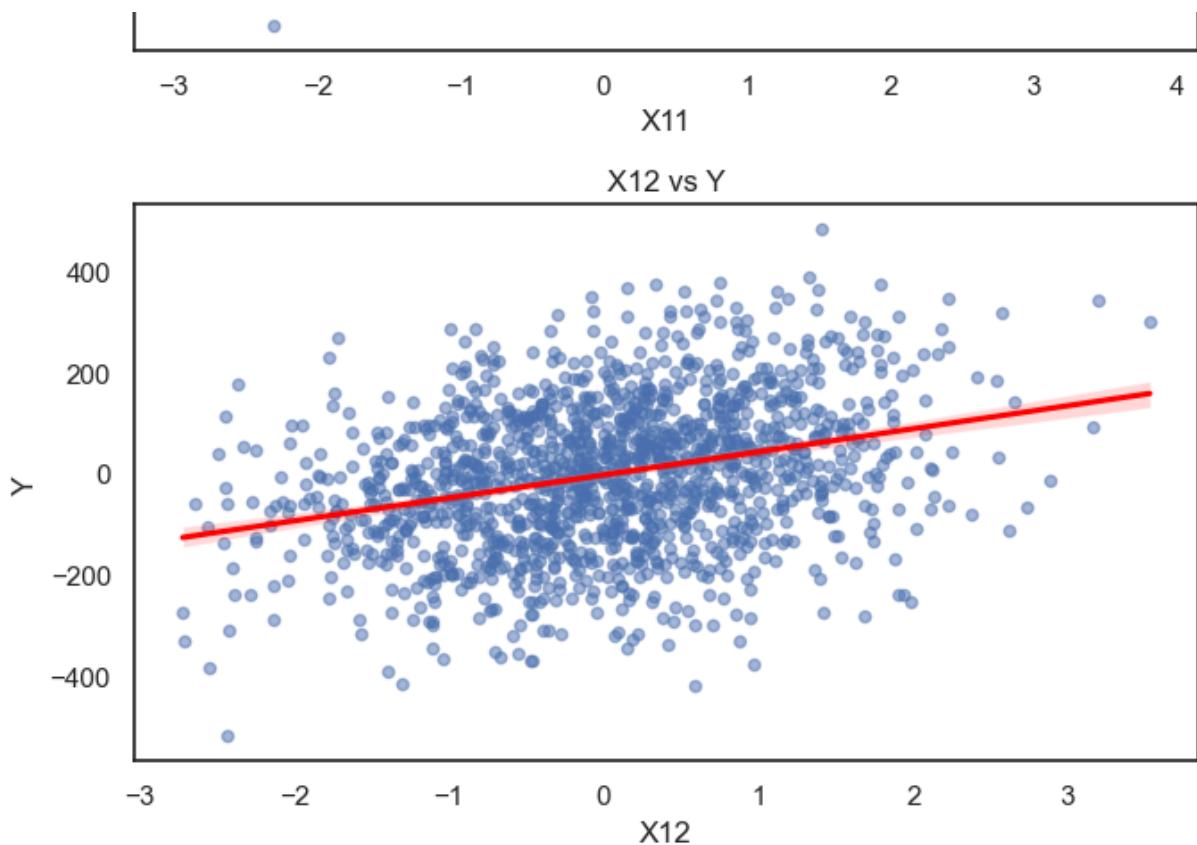
X6 vs Y





X9 vs Y





```
In [28]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

corr = df_reg[numeric_cols_reg.tolist() + ["Y"]].corr()

sns.set_theme(style="white")

mask = np.triu(np.ones_like(corr, dtype=bool))

n = corr.shape[0]
fig, ax = plt.subplots(figsize=(1.2 * n, 1.2 * n))

sns.heatmap(
    corr,
    mask=mask,
    annot=True,
    fmt=".2f",
    cmap="RdBu_r",
    center=0,
    linewidths=1.0,
    square=True,
    annot_kws={"size": 8},
    cbar_kws={"shrink": 0.6, "label": "Correlación"},
    vmin=-1,
    vmax=1,
    ax=ax,
)
```

```

ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha="right", fontsize=1)
ax.set_yticklabels(ax.get_yticklabels(), rotation=0, fontsize=10)

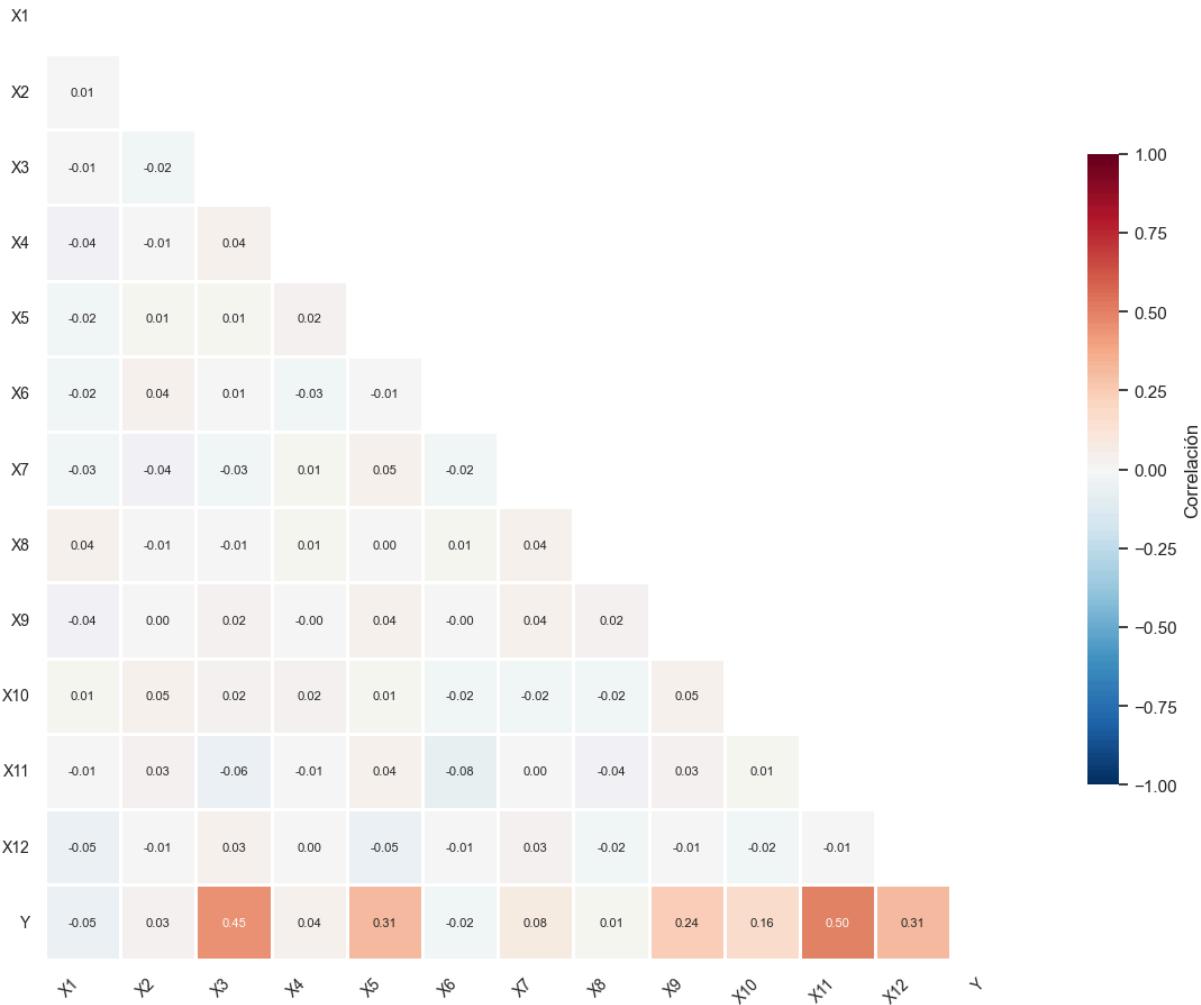
plt.subplots_adjust(bottom=0.3, left=0.2, top=0.9)

ax.set_title("Mapa de Calor de Correlaciones", fontsize=18, pad=20)

plt.show()

```

Mapa de Calor de Correlaciones



## Comments

### Classification

- El dataset tiene: 1.500 observaciones y 16 columnas (15 numéricas y 1 categórica de clase Y). Las variables X1–X15 son float y Y es int con valores {0,1,2}
- Menos del 1% de datos faltantes en 8 de las 15 variables (máx. 3 filas en X2/X4), representando sólo ~15 muestras; tras eliminar nulos quedan 1,485 filas.

- Percentiles 25, 50 y 75 en 0,1,2 nos indican clases equilibradas (~500 ejemplos por clase).
- Medias centradas cerca de cero, desviaciones estándar en rango 3–8.
- Casi todas las  $X_i$  presentan tres picos (formas multimodales) bien definidos en sus histogramas, lo que podrían ser agrupamientos alineados con las tres clases.
- Variables como  $X_1, X_{13}$  o  $X_{14}$  muestran 3 grupos marcados.
- Los scatterplots de  $Y$  vs  $X_1, X_8, X_{13}$  nos confirman agrupamientos clasificables por rangos de cada variable.
- En el pairplot, la mayoría de pares de variables muestran tres nubes separadas.
- En el mapa de calor de correlaciones se aprecian múltiples pares con valores muy fuertes  $|\rho| > 0.9$ . La variable  $Y$  se correlaciona positivamente con  $X_6 (+0.96)$ ,  $X_7 (+0.93)$  y negativamente con  $X_{13} (-0.98)$ .
- **Multicolinealidad:** Alta redundancia entre características (varias parejas  $>0.9$ ), se podrían usar técnicas de selección o reducción de dimensión para evitar sobreajuste.

## Regression

- El dataset contiene 1 473 observaciones y 13 columnas — 12 predictoras ( $X_1-X_{12}$ , float) y la variable objetivo  $Y$  (float).
- Menos del 2% de datos faltantes repartidos en varias  $X_i$ ; tras eliminar nulos se conservan las 1 473 filas (pérdida  $\approx 1,8\%$ ).
- Varias de las variables parecieran seguir distribuciones normales centradas en 0.
- $X_1, X_3, X_9$  presentan rangos mucho más grandes y colas pesadas.
- Desviaciones estándar distintas y distantes en  $X_1$  ( $\sim 12\ 665$ ),  $X_3$  ( $\sim 95\ 364$ ) y  $X_9$  ( $\sim 233\ 212$ ) frente a  $\approx 1$  en el resto, podría considerarse escalado para aquellos modelos en los que impacte esto.
- No existe correlación con  $Y$ , la mayoría de  $X_i$  muestra  $\rho \approx 0$ .
- Los puntos del scatter plot no siguen un patrón claro y están muy dispersos alrededor de la línea roja (recta de regresión).
- **Multicolinealidad mínima:** coeficientes  $X_i-X_j$  muy bajos ( $|\rho| < 0.06$ ), por lo que las variables son prácticamente independientes.

In [ ]: