# PROCESAMIENTO

✗ NO DATOS FALTANTES

✓ DUPLICADOS (-24)

TRUNCAR DECIMALES

| | Gender | Age | Height | Weight | family_history_with_overweight | FAVC | FCVC | NCP | CAEC | SMOKE | CH2O | SCC | FAF | TUE | CALC | MTRANS | NObeyesdad |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | 21.0 | 1.620000 | 64.000000 | yes | no | 2.0 | 3.0 | Sometimes | no | 2.000000 | no | 0.000000 | 1.000000 | no | Public_Transportation | Normal_Weight |
| 1 | Female | 21.0 | 1.520000 | 56.000000 | yes | no | 3.0 | 3.0 | Sometimes | yes | 3.000000 | yes | 3.000000 | 0.000000 | Sometimes | Public_Transportation | Normal_Weight |
| 2 | Male | 23.0 | 1.800000 | 77.000000 | yes | no | 2.0 | 3.0 | Sometimes | no | 2.000000 | no | 2.000000 | 1.000000 | Frequently | Public_Transportation | Normal_Weight |
| 3 | Male | 27.0 | 1.800000 | 87.000000 | no | no | 3.0 | 3.0 | Sometimes | no | 2.000000 | no | 2.000000 | 0.000000 | Frequently | Walking | Overweight_Level_I |
| 4 | Male | 22.0 | 1.780000 | 89.800000 | no | no | 2.0 | 1.0 | Sometimes | no | 2.000000 | no | 0.000000 | 0.000000 | Sometimes | Public_Transportation | Overweight_Level_II |
| ... | ... | ... | ... | ... | | ... | ... | ... | ... | ... | ... | | ... | ... | ... | ... | ... |
| 2106 | Female | 21.0 | 1.710730 | 131.408528 | yes | yes | 3.0 | 3.0 | Sometimes | no | 1.728139 | no | 1.676269 | 0.906247 | Sometimes | Public_Transportation | Obesity_Type_III |
| 2107 | Female | 22.0 | 1.748584 | 133.742943 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.005130 | no | 1.341390 | 0.599270 | Sometimes | Public_Transportation | Obesity_Type_III |
| 2108 | Female | 23.0 | 1.752206 | 133.689352 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.054193 | no | 1.414209 | 0.646288 | Sometimes | Public_Transportation | Obesity_Type_III |
| 2109 | Female | 24.0 | 1.739450 | 133.346641 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.852339 | no | 1.139107 | 0.586035 | Sometimes | Public_Transportation | Obesity_Type_III |
| 2110 | Female | 24.0 | 1.738836 | 133.472641 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.863513 | no | 1.026452 | 0.714137 | Sometimes | Public_Transportation | Obesity_Type_III |

→ INCONSISTENTES

CONSUMO VERDURAS, CANTIDAD COMIDAS, CONSUMO AGUA, ACTIVIDAD FISICA, USO DISPOSITIVOS ELECTRONICOS

# DUMMY

## EXPLICATIVAS

### NIVEL DE PESO

| | |
|---|---|
| **BAJO PESO** | **0** |
| NORMAL | 0 |
| **EXCESO DE PESO** | 0 |
| **OBESIDAD I** | **1** |
| **OBESIDAD II** | **1** |
| **OBESIDAD III** | **1** |

### HISTORIAL FAM.

| | |
|---|---|
| **SI** | **1** |
| NO | 0 |

### FUMA

| | |
|---|---|
| **SI** | **1** |
| NO | 0 |

### EDAD

| |
|---|
| Numerico |

### CALORIAS

| | |
|---|---|
| **SI** | **1** |
| NO | 0 |

### GENERO

| | |
|---|---|
| **FEMENINO** | **0** |
| MASCULINO | 1 |

### PESO

| |
|---|
| Numerico |

## VERDURAS

| |
|---|
| NUNCA |
| ALGUNAS VECES |
| SIEMPRE |

## COMIDAS PRINCIPALES

| |
|---|
| 1 - 2 |
| 3 |
| + 3 |

## AGUA

| |
|---|
| -1L |
| 1L - 2L |
| + 2L |

## ENTRE COMIDAS

| |
|---|
| NO |
| ALGUNAS VECES |
| FRECUENTEMENTE |
| SIEMPRE |

## ALCOHOL

| |
|---|
| NO |
| FRECUENTEMENTE |
| SIEMPRE |

## MEDIO TRANSPORTE

| |
|---|
| AUTOMOVIL |
| MOTO |
| TRANSPORTE PUBLICO |
| CAMINAR Y BICICLETA |

# SELECCION A PRIORI
## PRUEBA F Y CHI-CUADRADO

X = base[['family_history_with_overweight', 'FAVC', 'FCVC','NCP'          X = base[['Weight','Height','Age']]
,'SMOKE' ,'CH2O' ,'FAF' ,'TUE']]

```
Variables seleccionadas:
family_history_with_overweight        int64
FAF                                 float64
TUE                                 float64
dtype: object


Número original de variables:  8
Número de variables seleccionadas:  3
```

```
Variables seleccionadas:
Weight     float64
Age        float64
dtype: object


Número original de variables:  3
Número de variables seleccionadas:  2
```

X=base[['Weight','FCVC', 'NCP', 'FAF']]

```
Variables seleccionadas:
Weight     float64
FAF        float64
dtype: object


Número original de variables:  4
Número de variables seleccionadas:  2
```

# ALGORITMOS

## HIPERPARAMETROS

## LASSO

- ALPHAS : 115 ***
- CROOS VALIDATION : 9
- TOLERANCIAS : 3 ***
- SOLVER : newton-cg', 'lbfgs', 'liblinear ***
- CANTIDAD DE MODELOS : 9315

| | alphas | tol | solver | validacion | recall | roc |
|---|---|---|---|---|---|---|
| 8073 | 1.6 | 0.000100 | newton-cg | 8 | 0.958546 | 0.988462 |
| 8074 | 1.6 | 0.000100 | lbfgs | 8 | 0.958546 | 0.988462 |
| 8076 | 1.6 | 0.000010 | newton-cg | 8 | 0.958546 | 0.988462 |
| 8077 | 1.6 | 0.000010 | lbfgs | 8 | 0.958546 | 0.988462 |
| 8079 | 1.6 | 0.000001 | newton-cg | 8 | 0.958546 | 0.988462 |
| 8080 | 1.6 | 0.000001 | lbfgs | 8 | 0.958546 | 0.988462 |

SELECCION DE PREDICTORES
MENOS INFLUYENTE

## ARBOL DECISION

- PROFUNDIDAD : 80
- CROOS VALIDATION : 9
- SPLIT : BEST, RANDOM
- MINIMO DE MUESTRAS POR HOJA : 9
- CANTIDAD DE MODELOS : 12960

ROC
ACCURACY

| | n_en_nodos | profundidad | validacion | divisor | recall | roc |
|---|---|---|---|---|---|---|
| 80 | 4 | 5 | 10 | best | 0.982692 | 0.980977 |
| 10192 | 14 | 7 | 6 | best | 0.949306 | 0.983404 |

# Modelo

$$\log(odds) = 2.07 + 6.75*Family + 0.91*FAVC + 1.83*Smoke$$
$$-2.18*Gender + 0.40*CAEC\_A,$$
$$-1.24*CAEC\_f - 0.75*CAEC\_S, 0.25*CAEC\_N ....$$
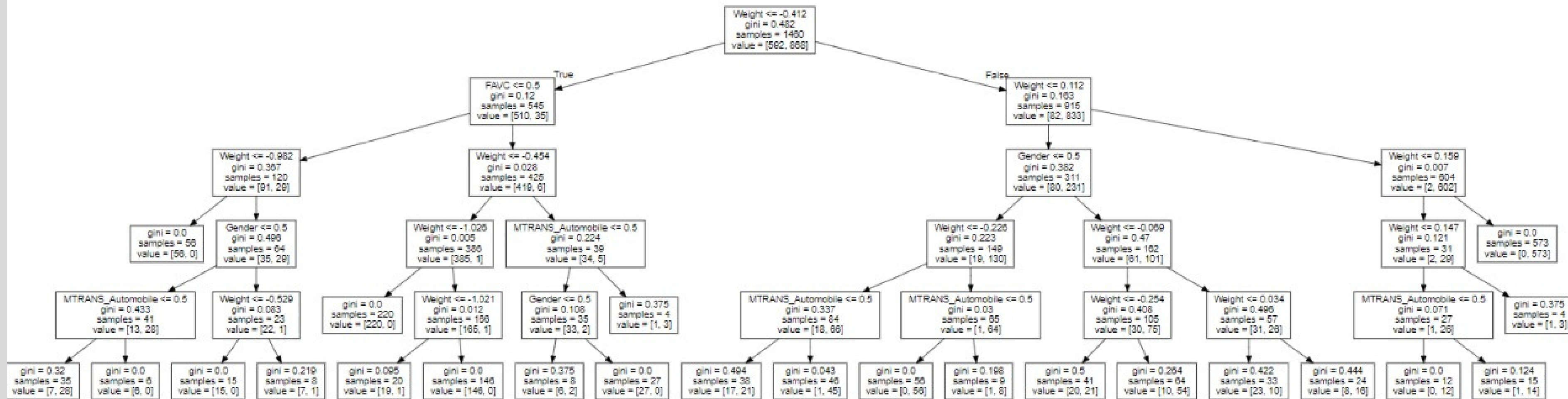
# SELECCION

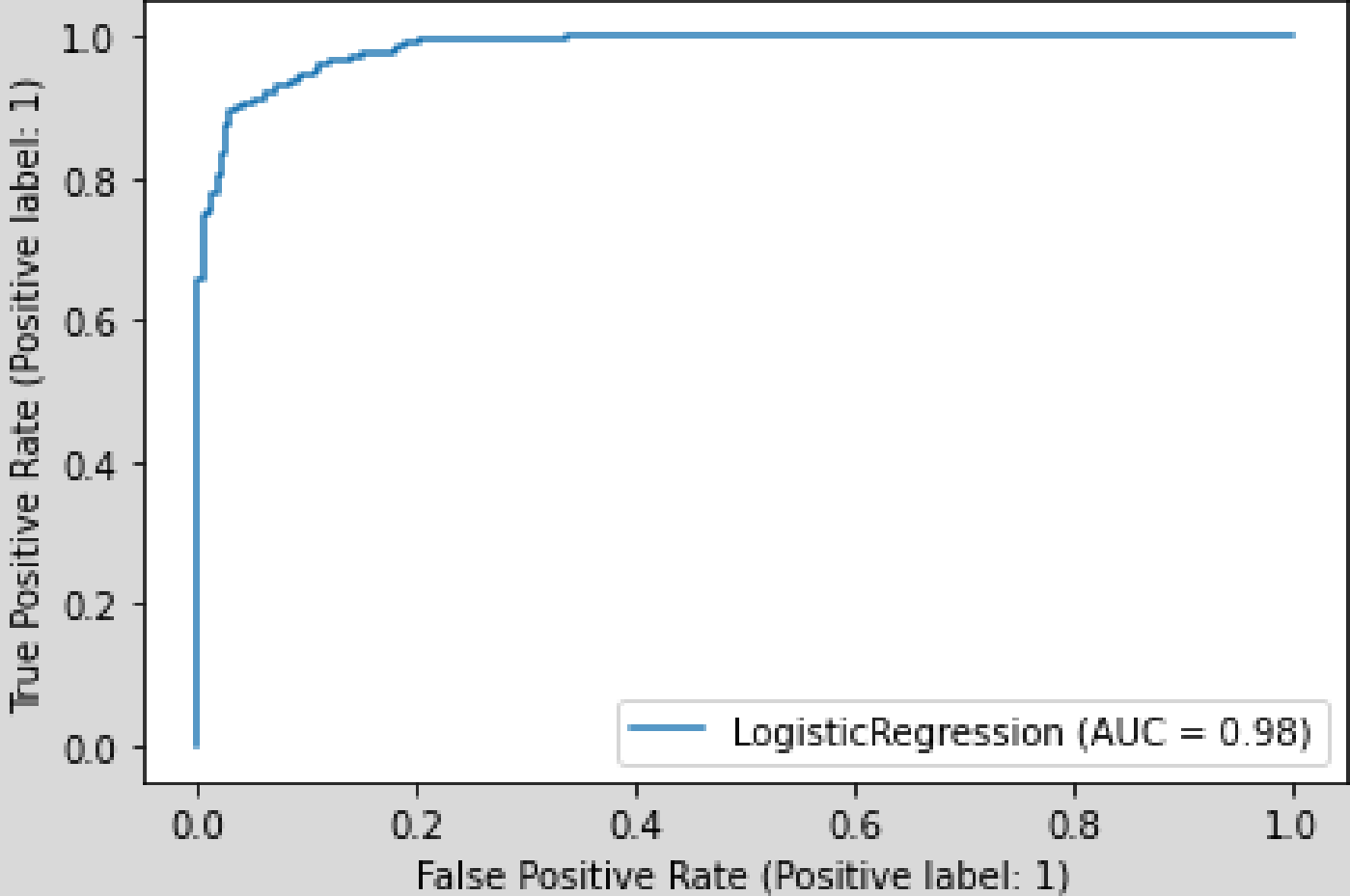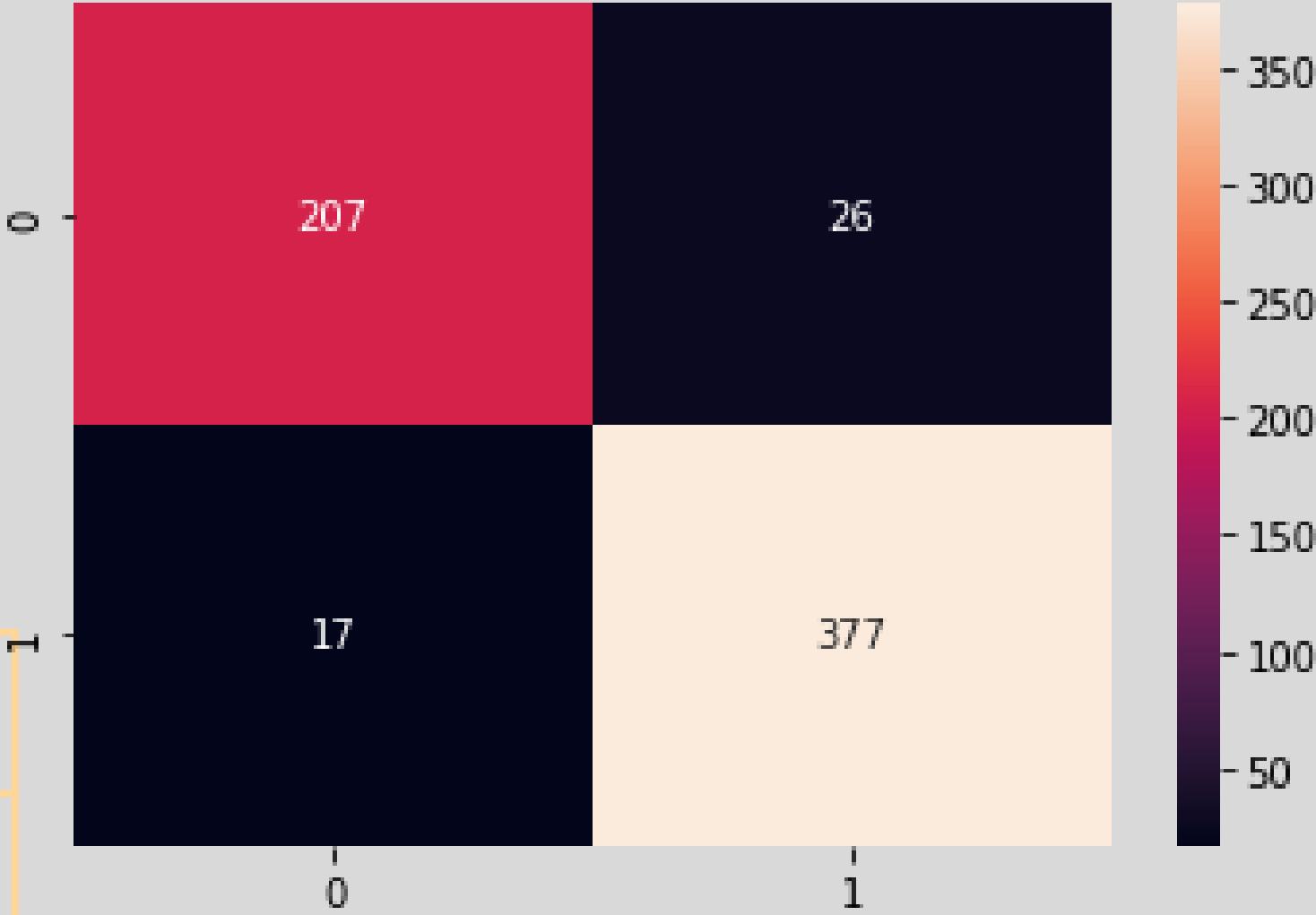## MEDIANTE ARBOL DE DECISION

SELECCION HACIA ADELANTE : 4 VARIABLES QUE MAXIMICEN ROC

```
Variables originales dataset: ['Weight' 'Age' 'family_history_with_overweight' 'FAVC' 'SMOKE' 'Gender'
 'CAEC_Always' 'CAEC_Frequently' 'CAEC_Sometimes' 'CAEC_no'
 'MTRANS_Automobile' 'MTRANS_Motorbike' 'MTRANS_Public_Transportation'
 'MTRANS_Walking_and_bike' 'CALC_Frequently' 'CALC_Sometimes' 'CALC_no'
 'CH2O' 'FAF' 'TUE']
Variables seleccionadas : ['Weight' 'FAVC' 'Gender' 'MTRANS_Automobile']
```

# ANÁLISIS GRÁFICOS

# SELECCION

## MEDIANTE REGRESION LOGISTICA L2

## SELECCION HACIA ATRAS

```
Variables originales dataset: ['Weight' 'Age' 'family_history_with_overweight' 'FAVC' 'SMOKE' 'Gender'
 'CAEC_Always' 'CAEC_Frequently' 'CAEC_Sometimes' 'CAEC_no'
 'MTRANS_Automobile' 'MTRANS_Motorbike' 'MTRANS_Public_Transportation'
 'MTRANS_Walking_and_bike' 'CALC_Frequently' 'CALC_Sometimes' 'CALC_no'
 'CH2O' 'FAF' 'TUE']
Variables seleccionadas : ['Weight' 'Age' 'family_history_with_overweight' 'FAVC' 'Gender' 'CALC_no']
```
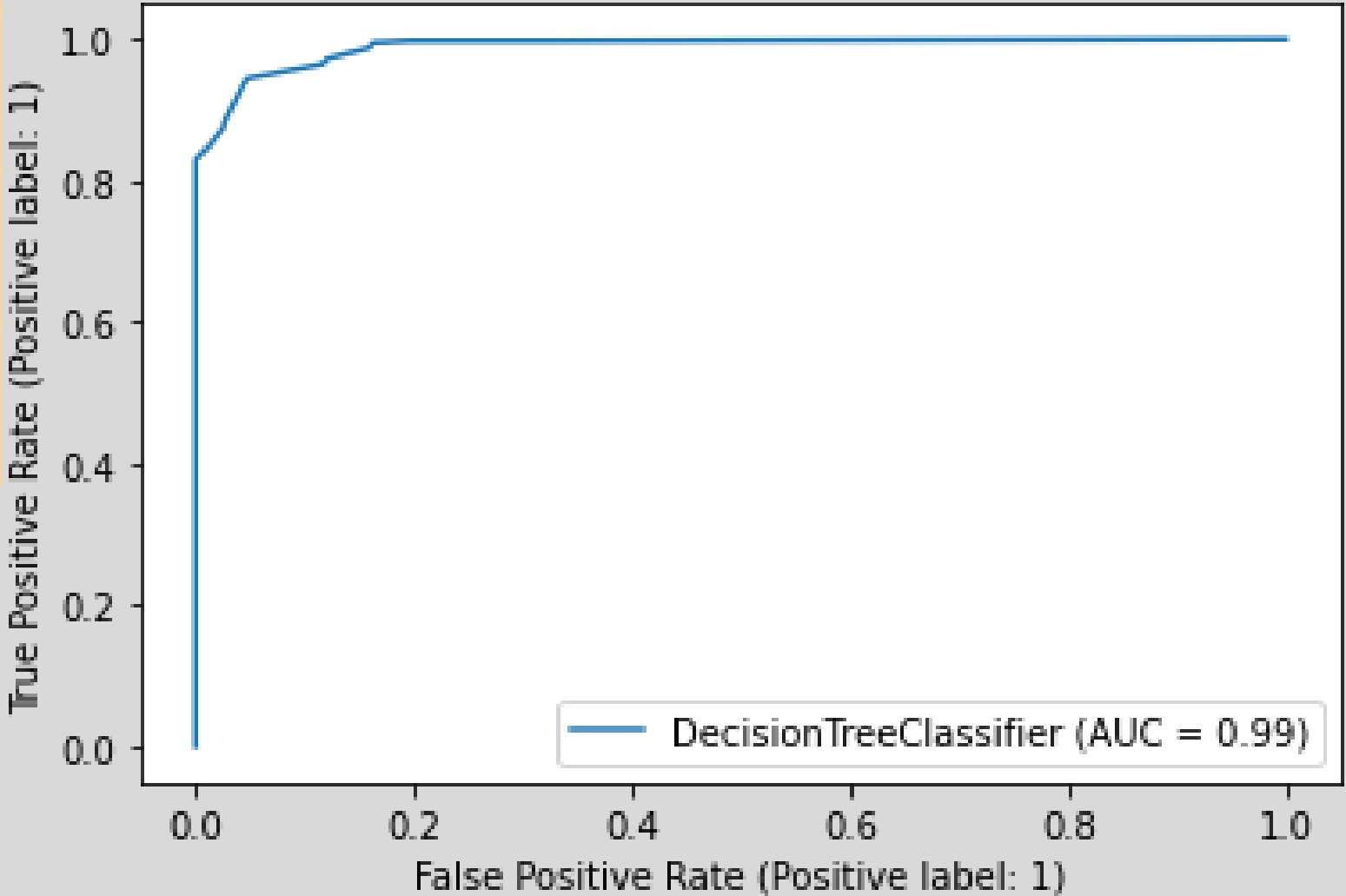
# EVALUACION MODELOS INDIVIDUALES

**LASSO**

| ACCURACY | 0.931 |
|----------|-------|
| ROC_AUC | 0.929 |
| RECALL | 0.956 |
| FBETA 2 | 0.952 |

# EVALUACION MODELOS INDIVIDUALES

## ARBOL DECISION

| | |
|---|---|
| ACCURACY | 0.934 |
| ROC_AUC | 0.935 |
| RECALL | 0.964 |
| FBETA 2 | 0.958 |

# BIBLIOGRAFIA

- https://archive.ics.uci.edu/ml/index.php
- https://archive.ics.uci.edu/ml/datasets/Estimation+of+obesity+levels+based+on+eating+habits+and+physical+condition+