

CLASIFICACIÓN DE NIVELES DE OBESIDAD MEDIANTE MODELOS DE APRENDIZAJE AUTOMÁTICO: UN ENFOQUE PREDICTIVO BASADO EN DATOS

Simón Olivieri
Leandro Urturi
Camila Muñoz



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

**DIPLOMATURA EN
CIENCIAS DE DATOS**



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

INTRODUCCIÓN

100%

AUMENTO DE LA
OBESIDAD
1990 - 2022



IMC

≥ 25 **SOBREPESO**
 ≥ 30 **OBESIDAD**



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

METODOLOGÍA

OBJETIVOS

Desarrollar modelos de clasificación del nivel de obesidad basados en el Aprendizaje Automático utilizando variables características.

Comparar el rendimiento de los distintos modelos propuestos para esta clasificación.

DATASET

**Estimation of Obesity Levels Based On
Eating Habits and Physical Condition**

2111
REGISTROS

17
VARIABLES

23% **DATOS REALES**

77% **DATOS SINTÉTICOS**

E1
RECOPIACIÓN
DE DATOS



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

METODOLOGÍA

E2
PREPROCESAMIENTO
DE DATOS

ESPERADO

3 VARIABLES NUMÉRICAS

14 VARIABLES CATEGÓRICAS

REAL

8 NUMÉRICAS

9 CATEGÓRICAS

FINAL

17 NUMÉRICAS

Train_set
700/
1477
REGISTROS

Val_set
150/
317
REGISTROS

Test_set
150/
317
REGISTROS



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

METODOLOGÍA

E3

SELECCIÓN DE
CARACTERÍSTICAS
IMPORTANTES



RANDOM FOREST CLASSIFIER
SELECT FROM MODEL

E4

MODELADO

RANDOM FOREST

n_estimators

bootstrap

max_features

max_depth

HIPERPARÁMETROS



OPTUNA

Optuna



scikit
learn

RandomizerSearchCV

OPTIMIZADORES



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

METODOLOGÍA

E4
MODELADO

MODELO DEEP LEARNING

ARQUITECTURA DE LA RED

Modelo Red neuronal	Número de capas ocultas	cantidad de nodos de capas ocultas	Ecuación de Agregación	Ecuación de activación capas de entrada y ocultas	Ecuación de activación capa de salida	DropOut por capa
Modelo 1	2	64	Lineal	Relu	Softmax	No tiene
Modelo 2	2	64	Lineal	Relu	Softmax	0.2
Modelo 3	2	64	Lineal	Tanh	Softmax	0.2
Modelo 4	2	64	Lineal	Tanh	Softmax	0.2
Modelo 5	1	64	Lineal	Tanh	Softmax	No tiene



MODELO DEEP LEARNING

FUN. COSTO – OPTIMIZADORES- HIPERPARÁMETROS

Modelo Red neuronal	Función de coste	Optimizador				
		Optimizador	learning rate	beta	Regularización L2	amsgrad
Modelo 1	CrossEntropyLoss	Adam	0.001	Valor por defecto: (0.9, 0.999)	0	False
Modelo 2	CrossEntropyLoss	Adam	0.0005	Valor por defecto: (0.9, 0.999)	0	False
Modelo 3	CrossEntropyLoss	Adam	0.001	Valor por defecto: (0.9, 0.999)	0	False
Modelo 4	CrossEntropyLoss	Adam	0.0005	Valor por defecto: (0.9, 0.999)	0	False
Modelo 5	CrossEntropyLoss	Adam	0.0005	Valor por defecto: (0.9, 0.999)	0	False



E1
RECOPIACIÓN
DEDATOS

E2
PREPROCESAMIENTO
DEDATOS

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2111 entries, 0 to 2110
Data columns (total 17 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Gender                                2111 non-null   object
1   Age                                   2111 non-null   float64
2   Height                               2111 non-null   float64
3   Weight                               2111 non-null   float64
4   family_history_with_overweight       2111 non-null   object
5   FAVC                                  2111 non-null   object
6   FCVC                                  2111 non-null   float64
7   NCP                                   2111 non-null   float64
8   CAEC                                  2111 non-null   object
9   SMOKE                                 2111 non-null   object
10  CH2O                                  2111 non-null   float64
11  SCC                                   2111 non-null   object
12  FAF                                   2111 non-null   float64
13  TUE                                   2111 non-null   float64
14  CALC                                  2111 non-null   object
15  MTRANS                                2111 non-null   object
16  NObeyesdad                           2111 non-null   object
dtypes: float64(8), object(9)
memory usage: 280.5+ KB
```

RESULTADOS

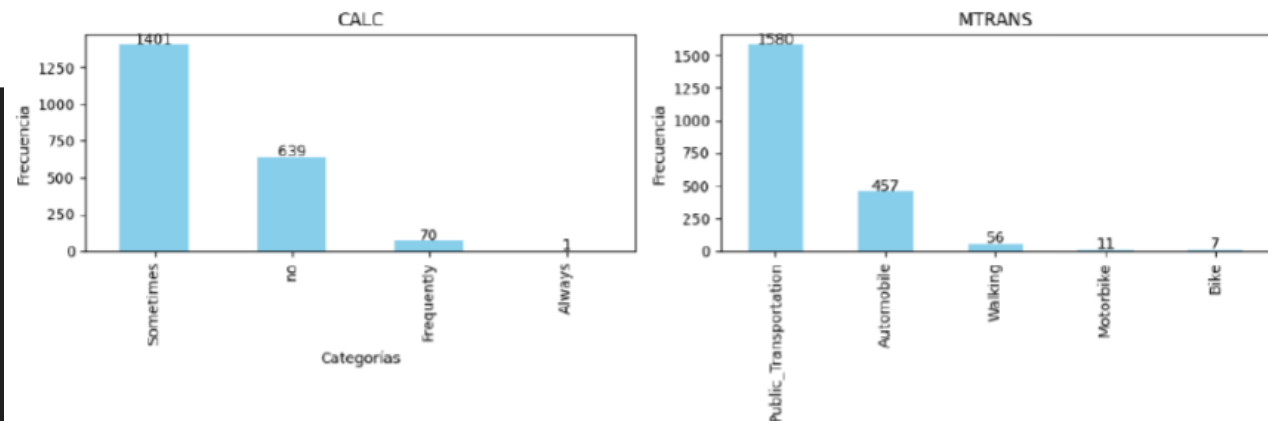


Figura 2. Variables MTRANS y CALC antes del preprocesamiento.

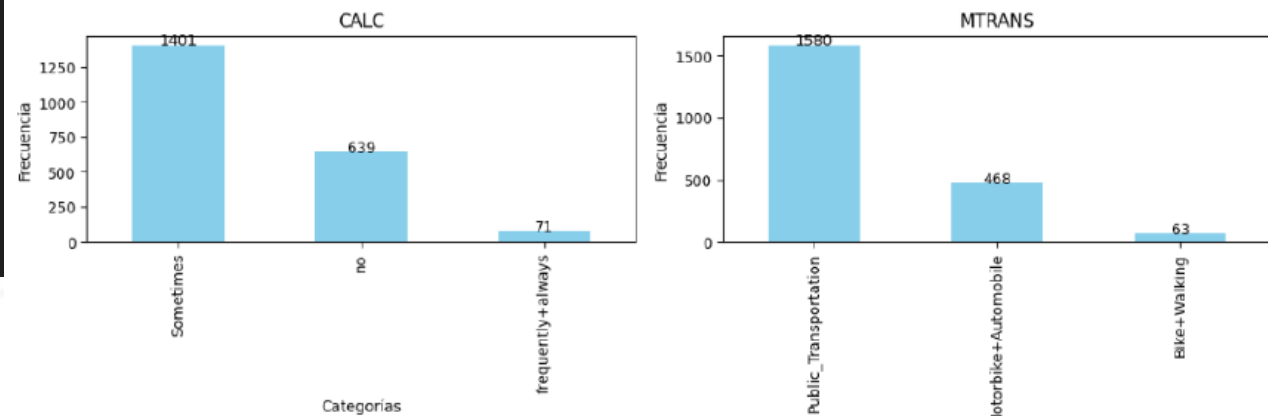


Figura 3. Variables MTRANS y CALC después del preprocesamiento.



E2

PREPROCESAMIENTO
DE DATOS

RESULTADOS

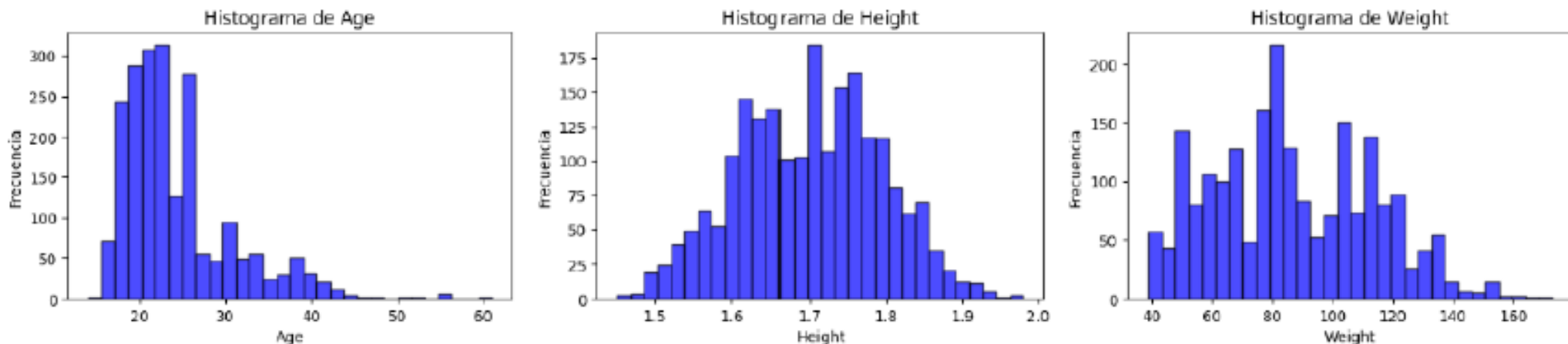


Figura 4. Histograma de las variables numéricas.

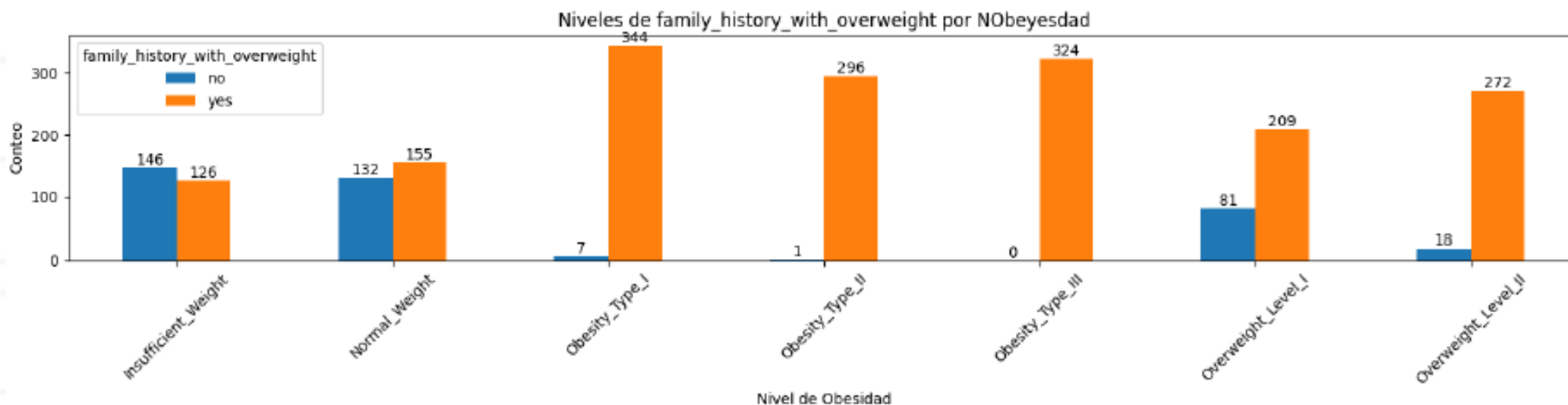


Figura 5. Relación entre la variable de historial familiar de obesidad y los niveles de obesidad.

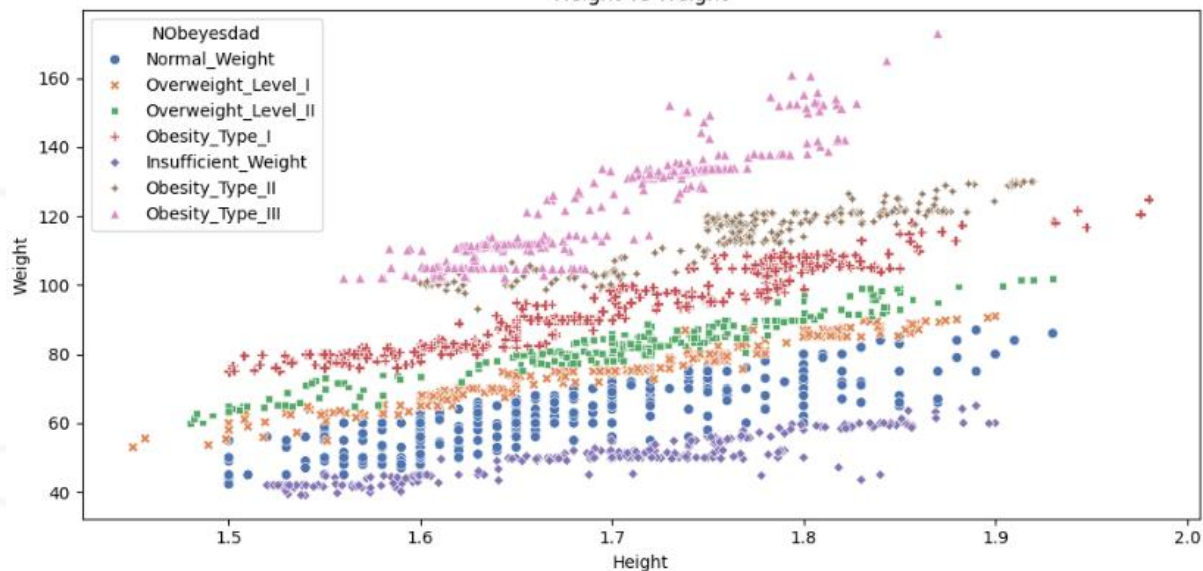


Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

E2

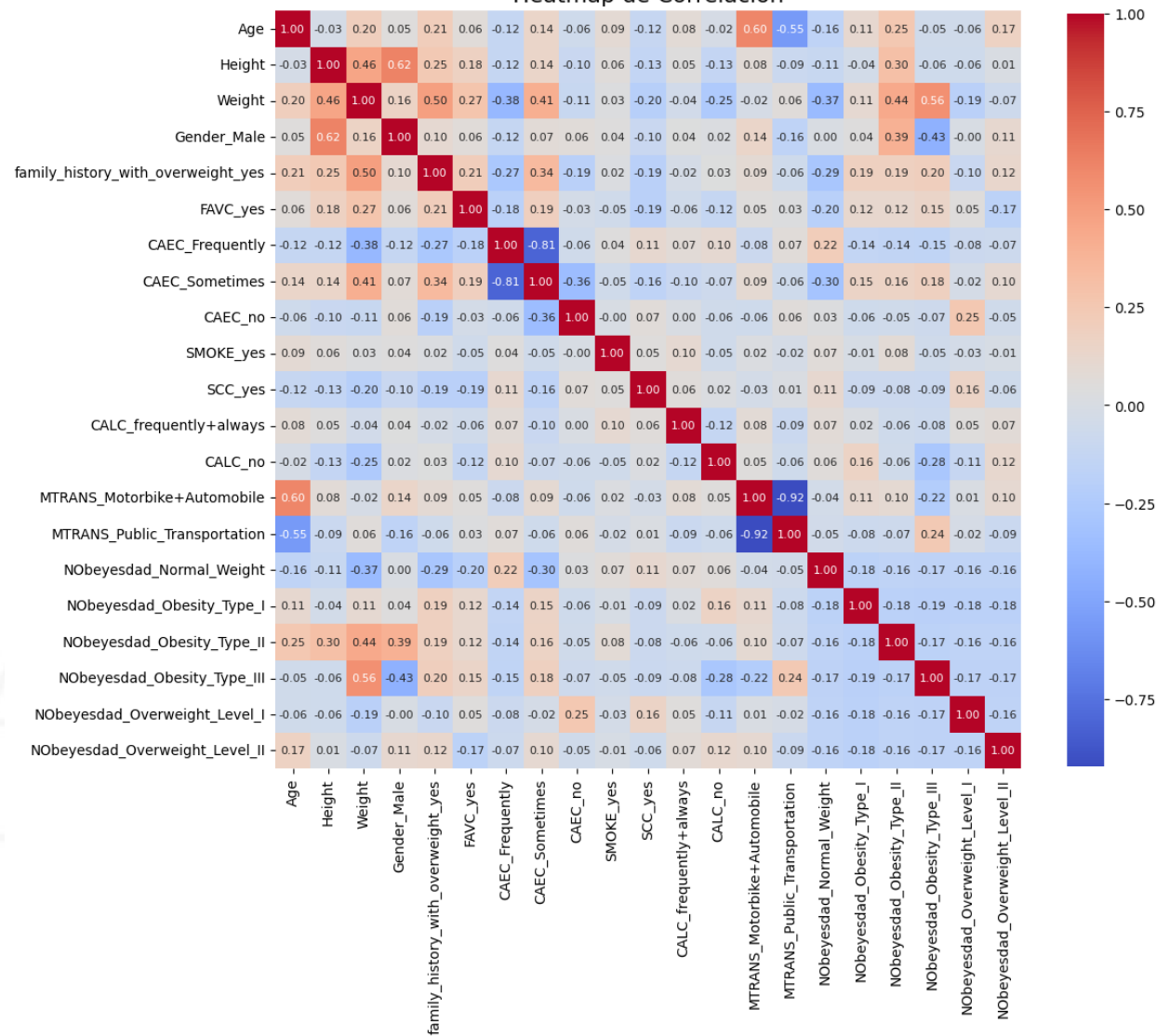
PREPROCESAMIENTO
DE DATOS

Height vs Weight



RESULTADOS

Heatmap de Correlación





RESULTADOS

E3 SELECCIÓN DE CARACTERÍSTICAS IMPORTANTES

SELECCIÓN DE CARACTERÍSTICAS DEL MODELO

Optimización de Hiperparámetros	max_features							
	1	2	3	4	5	6	7	8
RandomizedSearch CV	Weight	Height	Age	Gender_Male	family_history_with_overweight_yes	CALC_no	FAVC_yes	
Optuna	Weight	Height	Age	Gender_Male	family_history_with_overweight_yes	CALC_no	FAVC_yes	
Select model y RandomizedSearch CV	Weight	Height	Age	Gender_Male	family_history_with_overweight_yes	CALC_no	FAVC_yes	CAEC_Sometimes

E4 MODELADO

RESULTADOS DE LOS MODELOS CON RANDOMFORESTCLASSIFIER

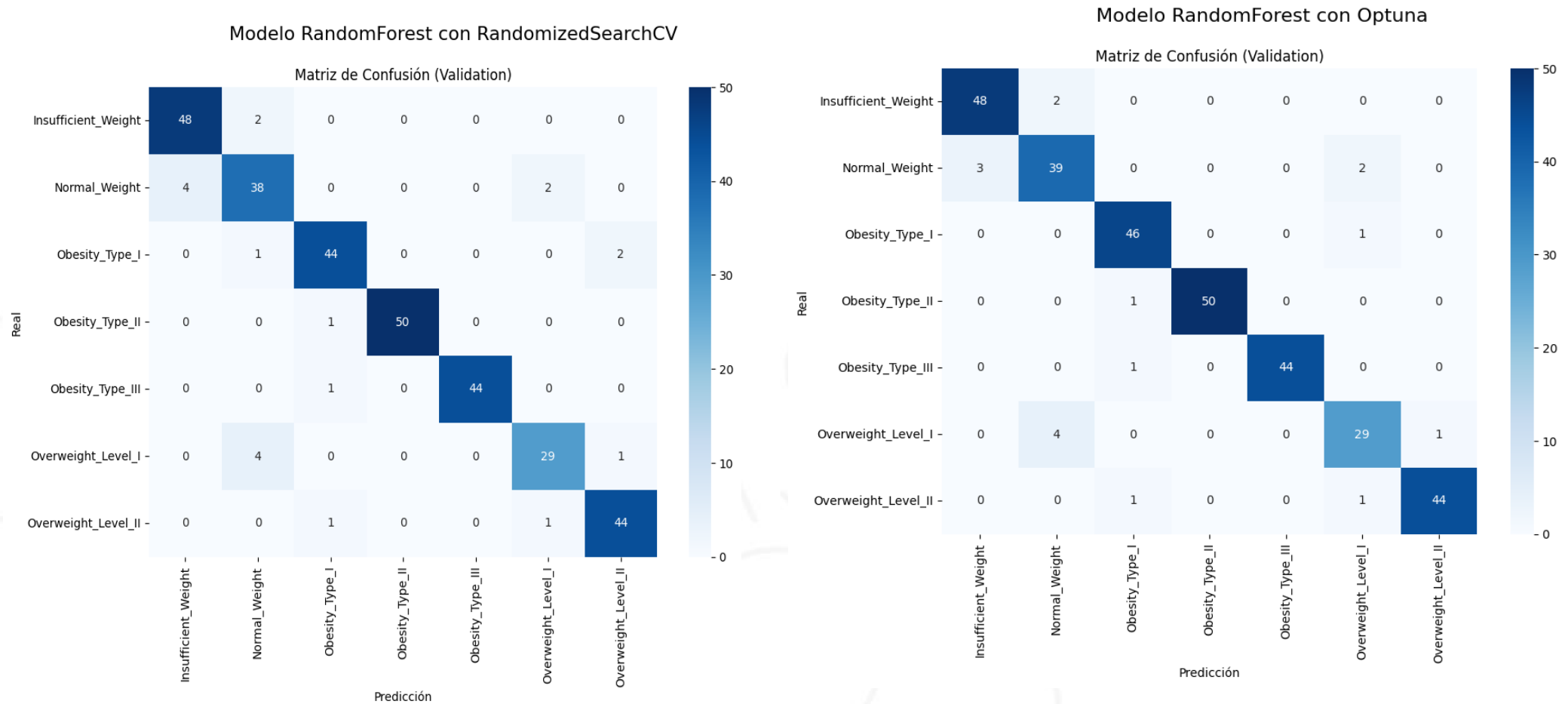
Optimización de Hiperparámetros	Mejores hiperparametros				Muestra de validación		Muestra de testeo	
	n_estimators	max_depth	bootstrap	max_features	Accuracy	F1-score	Accuracy	F1-score
RandomizedSearch CV	120	58	False	7	0.952	0.952	0.927	0.927
Optuna	203	84	False	7	0.952	0.95	0.946	0.946



RESULTADOS

E4
MODELADO

COMPARACIÓN ENTRE MATRICES DE CONFUSIÓN ENTRE AMBOS MODELOS





RESULTADOS

REDES NEURONALES

E4
MODELADO

Modelo Red neuronal	mini batch	epochs	Muestra de validación			Muestra de testeo	
			Acurracy	F1-score	Error	Acurracy	F1-score
Modelo 1	32	1250	0.959	0.958	0.30	0.981	0.981
Modelo 2	32	1250	0.962	0.961	0.30	0.974	0.974
Modelo 3	32	1250	0.962	0.962	0.12	0.981	0.981
Modelo 4	64	1250	0.962	0.962	0.09	0.977	0.978
Modelo 5	32	1250	0.971	0.971	0.12	0.971	0.971



MODELO 5



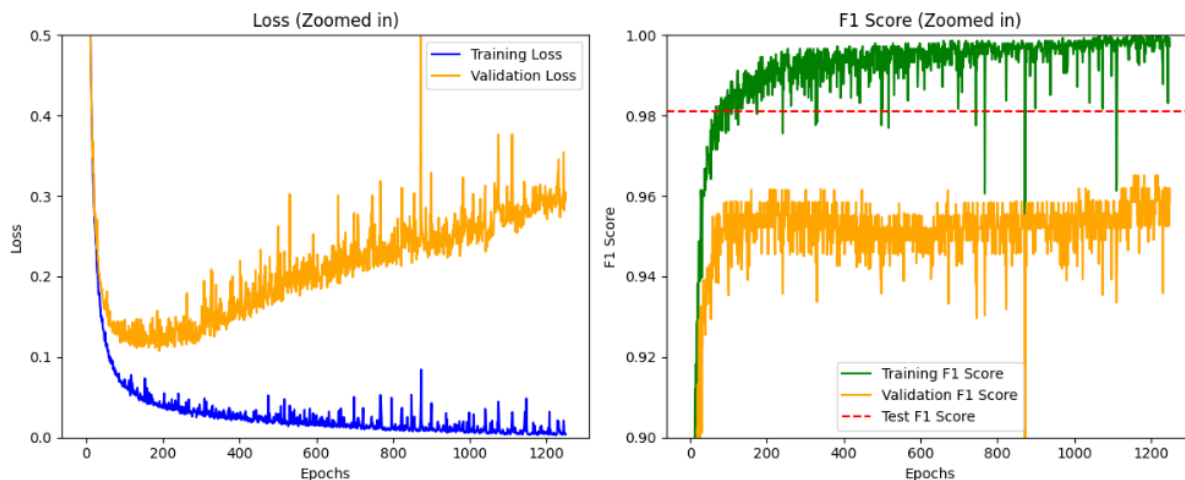


Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

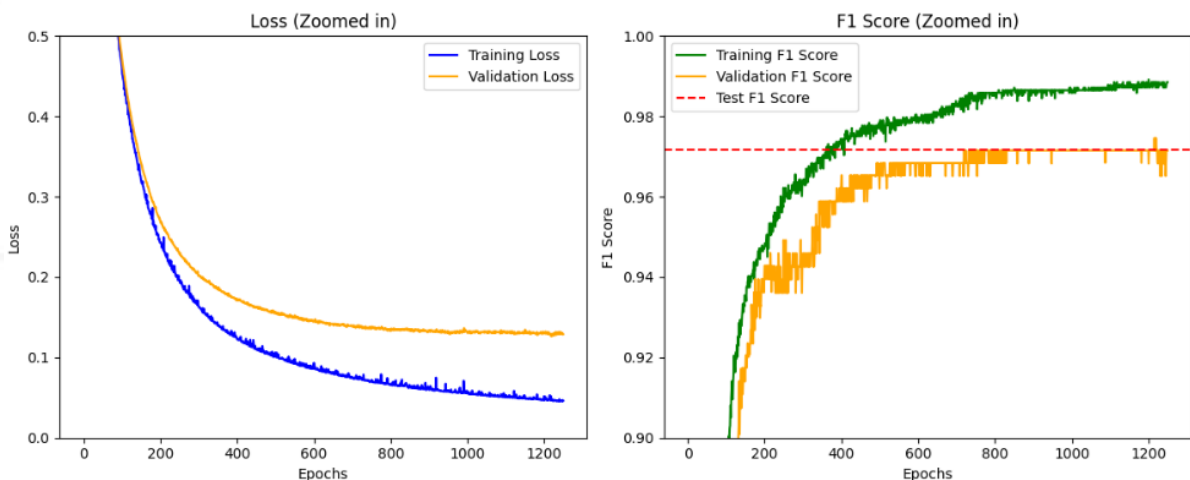
E4
MODELADO

RESULTADOS

MODELO 1



MODELO 5



MODELO 1

Matriz de Confusión (Validation)

Insufficient_Weight	49	1	0	0	0	0	0
Normal_Weight	2	40	0	0	0	2	0
Obesity_Type_I	0	0	46	1	0	0	0
Obesity_Type_II	0	0	0	51	0	0	0
Obesity_Type_III	0	0	0	1	44	0	0
Overweight_Level_I	0	3	0	0	0	29	2
Overweight_Level_II	0	0	0	0	0	1	45

Predicción

MODELO 5

Matriz de Confusión (Validation)

Insufficient_Weight	50	0	0	0	0	0	0
Normal_Weight	2	41	0	0	0	1	0
Obesity_Type_I	0	0	47	0	0	0	0
Obesity_Type_II	0	0	0	51	0	0	0
Obesity_Type_III	0	0	0	1	44	0	0
Overweight_Level_I	0	1	0	0	0	31	2
Overweight_Level_II	0	0	0	0	0	2	44

Predicción



Facultad de Ciencias Exactas
y Naturales y Agrimensura
**UNIVERSIDAD NACIONAL
DEL NORDESTE**

CONCLUSIONES

OPTUNA
RANDOM FOREST
94,6%
F1-score

MODELO 5
REDES NEURONALES
97,1%
F1-score