



UNIVERSIDAD  
**CATÓLICA**  
BOLIVIANA

DEPARTAMENTO DE INGENIERIA Y CIENCIAS EXACTAS  
**MACHINE LEARNING**

**INFORME DE PROYECTO  
MODELO PREDICTIVO:  
VOLUNTARIADO**

**Estudiante:**

**HENRRY ALBERTO CORONADO VILLCA**

**Fecha:** 06 de octubre de 2025

Santa Cruz – Bolivia

**IMPORTANTE**

El presente documento desarrollado tiene como objetivo presentar el desarrollo que se realizo en este proyecto, buscando documentar e informar el paso a paso de el producto. El resultado final y este documento es de dominio publico y no se buscara tener fines de lucro. El documento estará sujeto a cambios en el futuro por lo que se recomienda tener el ultimo versionado de este documento y su contenido.

Versión: 1.0

---

**Tabla de contenido**

1	INTRODUCCION .....	3
2	OBJETIVOS.....	3
2.1	OBJETIVO GENERAL.....	3
2.2	OBJETIVOS ESPECIFICOS .....	3
3	DESARROLLO DEL TRABAJO .....	4
3.1	Etapa 1: Desarrollo y Diagnóstico del Modelo.....	4
3.1.1	1. Preprocesamiento (División y Escalado) .....	4
3.1.2	2. Entrenamiento y Evaluación .....	4
3.1.3	3. Persistencia de Activos .....	4
3.2	Etapa 2: Operacionalización y Diagnóstico .....	4
3.2.1	1. Carga de Activos y Nueva Entrada.....	5
3.2.2	2. Preprocesamiento Operacional.....	5
3.2.3	3. Predicción y Reglas de Negocio.....	5
4	RESULTADOS OBTENIDOS.....	5
4.1	Coeficientes Aprendidos (Pesos del Modelo) .....	6
5	CONCLUSIONES.....	6
6	ANEXOS.....	7
6.1	Anexo 1: Training.....	7
6.1.1	Anexo1.1 .....	7
6.1.2	Anexo1.2 .....	8
6.1.3	Anexo1.3 .....	8
6.2	Anexo 2: Output Training.....	9
6.3	Anexo 3: Testing.....	9
6.3.1	Anexo3.1 .....	9
6.3.2	Anexo 3.2 .....	10
6.4	Anexo 4: Output Testing .....	10

## 1 INTRODUCCION

El presente informe documenta el desarrollo e implementación de un modelo de Machine Learning (ML) de Regresión Lineal Múltiple, cuyo objetivo principal es predecir la consistencia y el cumplimiento final de las horas de servicio de los voluntarios pastorales.

Tradicionalmente, la evaluación de riesgo se basa en la simple resta de horas acumuladas contra la meta. Este proyecto utiliza variables de patrón y capacidad (frecuencia, conflicto estructural, etc.) para ofrecer un diagnóstico predictivo avanzado a mitad del ciclo. El resultado es un sistema de apoyo que genera alertas tempranas y recomendaciones específicas para los voluntarios en riesgo de no alcanzar su meta al final del año, permitiendo a la Pastoral una gestión proactiva de la carga de trabajo.

## 2 OBJETIVOS

### 2.1 OBJETIVO GENERAL

Desarrollar un sistema de apoyo predictivo basado en Machine Learning para la gestión de voluntarios, capaz de pronosticar las Horas Totales Finales que cada voluntario acumulará, y utilizar esta predicción para emitir un diagnóstico de riesgo (Riesgo, Alerta, Normal) basado en la meta anual (N).

### 2.2 OBJETIVOS ESPECIFICOS

- Entrenar y evaluar un modelo de Regresión Lineal Múltiple que demuestre una capacidad de generalización (bajo RMSE y alto R2 en el conjunto de prueba).
- Implementar la persistencia para serializar el modelo entrenado y el escalador.
- Implementar reglas de negocio que traduzcan la predicción numérica ( $y$ -prediction) en un diagnóstico categórico práctico (Rojo, Amarillo, Verde).

### 3 DESARROLLO DEL TRABAJO

El trabajo se divide claramente en dos etapas para garantizar que el modelo no solo sea preciso, sino también funcional y usable en un entorno de producción (o en la herramienta de gestión diaria).

#### 3.1 Etapa 1: Desarrollo y Diagnóstico del Modelo

Esta etapa se centra en la construcción, entrenamiento y validación rigurosa del algoritmo.

##### 3.1.1 1. Preprocesamiento (División y Escalado)

División de Datos (60% Train / 20% Validation / 20% Test): Esta partición es la práctica estándar en ML. La separación del Test Set desde el inicio es crucial para garantizar que el error reportado (RMSETest) sea una estimación honesta de la performance del modelo en el mundo real (generalización).

Escalado de Características (StandardScaler): Las variables X fueron estandarizadas. El escalado asegura que todas las variables contribuyan de manera equitativa a la función de coste del optimizador Descenso de Gradiente Estocástico (SGD). El escalador se ajustó solo con el conjunto de entrenamiento para prevenir el Data Leakage.

##### 3.1.2 2. Entrenamiento y Evaluación

Modelo de Regresión Lineal (SGDRegressor): Se eligió esta técnica por su alta interpretabilidad. Los coeficientes son directamente analizables, permitiendo a la Pastoral entender la contribución de cada factor a la predicción final.

Métricas de Evaluación (RMSE y R2):

El RMSE (Raíz del Error Cuadrático Medio) se usó para medir la magnitud del error en las unidades de la variable objetivo (horas).

El R2 (Coeficiente de Determinación) se usó para cuantificar el porcentaje de la varianza en Y explicado por las variables X.

##### 3.1.3 3. Persistencia de Activos

Se utilizó la librería joblib para serializar el Modelo y el Escalador. Esto permite que el modelo se cargue en un entorno de aplicación sin necesidad de reentrenarlo, pasando a la Etapa 2.

#### 3.2 Etapa 2: Operacionalización y Diagnóstico

Esta etapa describe cómo la Pastoral utiliza los activos guardados en un escenario de aplicación.

### 3.2.1 1. Carga de Activos y Nueva Entrada

Los archivos .pkl del modelo y del escalador son cargados en memoria.

Se ingresan las métricas actuales del voluntario en un nuevo DataFrame.

### 3.2.2 2. Preprocesamiento Operacional

Uso del Escalador Persistido: El nuevo DataFrame se somete a la transformación (.transform()) utilizando el escalador cargado.

Justificación: Los datos de entrada deben estandarizarse exactamente con la misma media y desviación estándar utilizadas durante el entrenamiento del modelo para que la predicción sea válida.

### 3.2.3 3. Predicción y Reglas de Negocio

Predicción Numérica: El modelo predice un valor continuo: (Y-prediction) (horas totales proyectadas).

Aplicación de Reglas de Negocio: La (Y-prediction) se compara con la Meta Anual (N) y los umbrales predefinidos (85% de N y 100% de N) para emitir el diagnóstico final:

ROJO (Riesgo Estructural): Proyección bajo el 85%.

AMARILLO (Cuerda Floja): Proyección bajo la Meta (N).

VERDE (Cumplimiento): Proyección supera la Meta.

## 4 RESULTADOS OBTENIDOS

La evaluación del modelo entrenado arrojó los siguientes resultados en el conjunto de prueba (Test Set), lo que indica su capacidad de generalización:

Métrica	Valor Obtenido	Interpretación
RMSE (Raíz del Error Cuadrático Medio)	±7.96 Horas	El modelo predice las horas finales con un margen de error promedio de aproximadamente ±8 horas.
R2 (Coeficiente de Determinación)	0.8640	El 86.40% de la variación observada en las horas finales de los voluntarios es explicada por las variables X incluidas en el modelo.

#### 4.1 Coeficientes Aprendidos (Pesos del Modelo)

Los coeficientes (pesos) del modelo confirman el patrón conceptual diseñado, indicando la influencia de cada variable en la predicción de Horas Finales.

Variable	Peso (Coeficiente)	Interpretación
X1_Horas_Actuales	+17.65	Es el predictor más fuerte. Por cada desviación estándar en horas actuales, el resultado final aumenta significativamente.
X2_Frec_Semanal	+5.26	Alta frecuencia constante es un fuerte predictor positivo de cumplimiento final.
X3_Horas_Fallidas	-3.67	Las horas fallidas por conflicto actúan como un penalizador importante, indicando problemas estructurales.
X5_Disp_Neta_Restante.	+0.22	La capacidad restante de participar en actividades es un factor positivo, aunque menor que el ritmo actual.
X6_Antiguedad	+3.08	La experiencia (antigüedad) ofrece un beneficio predictivo positivo.
X4_Semanas_Restantes	≈0	Se excluye debido a que fue una variable constante en el conjunto de entrenamiento, siendo su efecto absorbido por el Intercepto

## 5 CONCLUSIONES

El modelo de Regresión Lineal optimizado con SGD es una herramienta predictiva robusta y eficiente para la Pastoral.

- Diagnóstico Proactivo: El sistema permite pasar de una gestión reactiva a una proactiva, identificando voluntarios que enfrentan un riesgo estructural. La separación en dos etapas asegura que la herramienta sea sostenible y fácil de actualizar.

- Validez Científica: El valor de R<sup>2</sup>≈0.89 demuestra la fuerte relación entre los factores de ritmo/conflicto y la consistencia final del voluntario.
- Usabilidad y Persistencia: La serialización con joblib y la clara distinción entre la fase de entrenamiento y la fase de uso (Etapa 1 y Etapa 2) hacen que el modelo sea inmediatamente operacionalizable en el sistema de gestión del voluntariado.

## 6 ANEXOS

### 6.1 Anexo 1: Training

#### 6.1.1 Anexo1.1

```
VoluntaryPredictor_ModelTraining.py > ...
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.linear_model import SGDRegressor
5 from sklearn.metrics import mean_squared_error
6 from sklearn.metrics import r2_score
7 import numpy as np
8 import joblib
9 import os
10
11 # Cargar el Dataset simulado
12 df = pd.read_csv('voluntariado_dataset.csv')
13
14 # Definir X (Variables Independientes) y Y (Variable Dependiente)
15 X = df[['X1_Horas_Actuales', 'X2_Frec_Semanal', 'X3_Horas_Fallidas',
16 |   |   'X4_Semanas_Restantes', 'X5_Disp_Neta_Restante', 'X6_Antiguedad']]
17 Y = df['Y_Horas_Totales_Finales']
18
19 # Separar el Test Set (20%) del resto (80%)
20 X_train_val, X_test, Y_train_val, Y_test = train_test_split(
21 |   X, Y, test_size=0.2, random_state=42
22 )
23
24 # Separar Validation Set (25% de X_train_val = 20% del total)
25 X_train, X_val, Y_train, Y_val = train_test_split(
26 |   X_train_val, Y_train_val, test_size=0.25, random_state=42
27 )
28
29 print(f"Total de muestras (100%): {len(df)}")
30 print(f"Muestras Train (60%): {len(X_train)}")
31 print(f"Muestras Validation (20%): {len(X_val)}")
32 print(f"Muestras Test (20%): {len(X_test)}")
33
```

### 6.1.2 Anexo1.2

```
👉 VoluntaryPredictor_ModelTraining.py > ...
34 # Instanciamos el StandardScaler para estandarizar las características
35 scaler = StandardScaler()
36
37 # Entrenar el Scaler
38 scaler.fit(X_train)
39
40 # 2. Aplicar la Transformación a TODOS los conjuntos
41 X_train_scaled = scaler.transform(X_train)
42 X_val_scaled = scaler.transform(X_val)
43 X_test_scaled = scaler.transform(X_test)
44
45 # Instanciar el SGDRegressor:
46 model = SGDRegressor(
47     loss='squared_error', # Función de coste para Regresión Lineal (Error Cuadrático Medio)
48     max_iter=1000,
49     tol=1e-3,
50     random_state=42
51 )
52
53 # Entrenamiento: Ajuste de los parámetros (pesos Theta) para encontrar el patrón.
54 model.fit(X_train_scaled, Y_train)
55
56 # Mostrar el patrón aprendido por el modelo (coeficientes)
57 print("\n--- Patrón Aprendido (Coeficientes Theta) ---")
58 feature_names = X.columns
59 for name, coef in zip(feature_names, model.coef_):
60     print(f"Peso de {name}: {coef:.2f}")
61
62 # Predicciones en los tres conjuntos
63 Y_train_pred = model.predict(X_train_scaled)
64 Y_val_pred = model.predict(X_val_scaled)
65 Y_test_pred = model.predict(X_test_scaled)
66
```

### 6.1.3 Anexo1.3

```
👉 VoluntaryPredictor_ModelTraining.py > ...
67 # Función para calcular el RMSE (Raíz del Error Cuadrático Medio)
68 def calculate_rmse(Y_true, Y_pred):
69     return np.sqrt(mean_squared_error(Y_true, Y_pred))
70
71 # Cálculo del RMSE para diagnóstico
72 rmse_train = calculate_rmse(Y_train, Y_train_pred)
73 rmse_val = calculate_rmse(Y_val, Y_val_pred)
74
75 print("\n--- Diagnóstico de Generalización ---")
76 print(f"RMSE (Train Set): {rmse_train:.2f} Horas")
77 print(f"RMSE (Validation Set): {rmse_val:.2f} Horas")
78
79 # Análisis de diagnóstico: Buscando Overfitting (Alta Varianza)
80 if rmse_val > rmse_train * 1.1: # Si el error de validación es > 10% mayor que el de train
81     print("⚠️ ¡ATENCIÓN! Alta Varianza (Overfitting). El modelo predice mejor sus datos de entrenamiento.")
82 else:
83     print("✅ Buena Generalización. Los errores son similares, indicando bajo Overfitting.")
84
85 # Evaluación Final (Generalización Real)
86 rmse_test = calculate_rmse(Y_test, Y_test_pred)
87 r2_test = r2_score(Y_test, Y_test_pred) # Cálculo del R^2
88
89 print("\n--- Evaluación Final (Generalización Real) ---")
90 print(f"Error promedio de predicción (RMSE): {rmse_test:.2f} horas.")
91 print(f"Coeficiente de Determinación (R^2): {r2_test:.4f}")
92
93 print(f"Interpretación: El {r2_test*100:.2f}% de la variación en las horas finales es explicada por las variables del modelo.")
94
95 # Crear la carpeta para guardar los activos
96 if not os.path.exists('model_assets'):
97     os.makedirs('model_assets')
98
99 joblib.dump(model, 'model_assets/voluntariado_regresor.pkl')
100 joblib.dump(scaler, 'model_assets/voluntariado_scaler.pkl')
101
102 print("\n✅ Persistencia exitosa: Modelo y Escalador guardados en 'model_assets/'")
```

## 6.2 Anexo 2: Output Training

```
(entVirtual) C:\Users\HenrryCoronado\Documents\UCB\Semestre_VI\Machine L\Practices\Proyecto>py VoluntaryPredictor_ModelTraining.py
Total de muestras (100%): 1000
Muestras Train (60%): 600
Muestras Validation (20%): 200
Muestras Test (20%): 200

--- Patrón Aprendido (Coeficientes Theta) ---
Peso de X1_Horas_Actuales : 17.65
Peso de X2_Frec_Semanal : 5.26
Peso de X3_Horas_Fallidas : -3.67
Peso de X4_Semanas_Restantes : 0.00
Peso de X5_Disponibilidad_Neta_Restante : 0.22
Peso de X6_Antiguedad : 3.08

--- Diagnóstico de Generalización ---
RMSE (Train Set): 7.87 Horas
RMSE (Validation Set): 7.78 Horas
 Buena Generalización. Los errores son similares, indicando bajo Overfitting.

--- Evaluación Final (Generalización Real) ---
Error promedio de predicción (RMSE): ±7.96 horas.
Coeficiente de Determinación (R²): 0.8640
Interpretación: El 86.40% de la variación en las horas finales es explicada por las variables del modelo.

 Persistencia exitosa: Modelo y Escalador guardados en 'model_assets'.
```

(entVirtual) C:\Users\HenrryCoronado\Documents\UCB\Semestre\_VI\Machine L\Practices\Proyecto>

## 6.3 Anexo 3: Testing

### 6.3.1 Anexo3.1

```
VoluntaryPredictor_UseModel.py > ...
1 import joblib
2 import pandas as pd
3 import numpy as np
4
5 HORAS_META_N = 100
6
7 # cargar modelo entrenado
8 loaded_model = joblib.load('model_assets/voluntariado_regressor.pkl')
9 loaded_scaler = joblib.load('model_assets/voluntariado_scaler.pkl')
10
11 # ejemplo
12 nuevos_datos_voluntario = pd.DataFrame([
13     'X1_Horas_Actuales': 30,
14     'X2_Frec_Semanal': 1.0,
15     'X3_Horas_Fallidas': 15,
16     'X4_Semanas_Restantes': 8,
17     'X5_Disponibilidad_Neta_Restante': 70, # Capacidad restante de actividades
18     'X6_Antiguedad': 1
19 ])
20
21 # usar el scaler cargado para transformar los nuevos datos
22 X_nuevo_escalado = loaded_scaler.transform(nuevos_datos_voluntario)
23
24 # generar predicción con los nuevos datos escalados
25 prediccion_horas_array = loaded_model.predict(X_nuevo_escalado)
26 y_prediccion = prediccion_horas_array[0]
27
28 print(f"Predicción Numérica del Modelo (y_prediccion): {y_prediccion:.1f} horas")
29
30 # definir umbrales para el diagnóstico
31 UMbral_Riesgo = HORAS_META_N * 0.85 # 85 horas
32 UMbral_Alerta = HORAS_META_N * 1.00 # 100 horas
33
34 diagnostico = {
35     'estado': 'Indefinido',
36     'color': 'gris',
37     'recomendacion_voluntario': '',
38     'carga_semestral': 0
39 }
40
```

### 6.3.2 Anexo 3.2

```

VoluntaryPredictor_UseModel.py > ...
40
41 # Lógica de Diagnóstico
42 if y_prediccion < UMBRAL_RIESGO:
43
44     # condición de riesgo (no cumple de ninguna manera)
45     horas_faltantes = HORAS_META_N - y_prediccion
46
47     diagnostico['estado'] = "RIESGO ESTRUCTURAL / ANOMALÍA"
48     diagnostico['color'] = "ROJO"
49     diagnostico['recomendacion_voluntario'] = (
50         f"Alerta inmediata: Se proyecta una falta de {horas_faltantes:.1f} horas. "
51         "Debe solicitar un plan de recuperación urgente y priorizar el cumplimiento de X5."
52     )
53     diagnostico['carga_semestre_2'] = horas_faltantes
54
55 elif y_prediccion < UMBRAL_ALERTA:
56
57     # condición de cuerda floja
58     diagnostico['estado'] = "CUERDA FLOJA / ALERTA"
59     diagnostico['color'] = "AMARILLO"
60     diagnostico['recomendacion_voluntario'] = (
61         "Atención: Se proyecta que quedará ligeramente por debajo de la meta. "
62         "Debe asistir a TODAS las actividades remanentes posibles para cumplir. El margen de error es mínimo."
63     )
64     diagnostico['carga_semestre_2'] = HORAS_META_N - y_prediccion
65
66 else:
67
68     # condición normal
69     horas_extra = y_prediccion - HORAS_META_N
70
71     diagnostico['estado'] = "OPERACIÓN NORMAL / CUMPLIMIENTO"
72     diagnostico['color'] = "VERDE"
73     diagnostico['recomendacion_voluntario'] = (
74         f"Felicitaciones: Se proyecta que superará la meta con {horas_extra:.1f} horas de colchón."
75     )
76     diagnostico['carga_semestre_2'] = 0
77
78     # mostrar diagnóstico
79     print("\n-----")
80     print(f"Diagnóstico ML para el Voluntario: {diagnostico['estado']} ")
81     print(f"Estado de Riesgo: {diagnostico['color']} ")
82     print("-----")
83     print(f"Recomendación para el Voluntario: {diagnostico['recomendacion_voluntario']} ")
84     print(f"Proyección de Carga para Semestre 2: {diagnostico['carga_semestre_2']:.1f} horas a arrastrar.\n")

```

### 6.4 Anexo 4: Output Testing

```

(entVirtual) C:\Users\HenrryCoronado\Documents\UCB\Semestre_VI\Machine_L\Practices\Project\entVirtual\scripts>cd ../..
(entVirtual) C:\Users\HenrryCoronado\Documents\UCB\Semestre_VI\Machine_L\Practices\Project>py VoluntaryPredictor_UseModel.py
Predicción Numérica del Modelo (y_prediccion): 60.2 horas
-----
| Diagnóstico ML para el Voluntario: RIESGO ESTRUCTURAL / ANOMALÍA
| Estado de Riesgo: ROJO
-----
Recomendación para el Voluntario: Alerta inmediata: Se proyecta una falta de 39.8 horas. Debe solicitar un plan de recuperación urgente y priorizar el cumplimiento de X5.
Proyección de Carga para Semestre 2: 39.8 horas a arrastrar.

(entVirtual) C:\Users\HenrryCoronado\Documents\UCB\Semestre_VI\Machine_L\Practices\Project>

```