

MODELO PREDICTIVO DE RIESGO DE INCIDENTES

Hackathon PUCE - UNACEM | Equipo GEOINNOVA

NOMBRE DEL MODELO

GEOINNOVA -
modelo_geoinnova_multitask.h5

TIPO DE MODELO

Deep Learning (LSTM) - Multi-Salida

OBJETIVO

Predicción de riesgo total y distribución de incidentes por categorías

EQUIPO RESPONSABLE

GEOINNOVA

RESPONSABLE TÉCNICO

Johnatan Guacho - Ricardo Carrion
(Desarrollo Python / API / Web)

COLABORACIÓN ANALÍTICA

Equipo Geoinnova - Integrantes Carrera de Estadística (Limpieza y análisis de variables)

1. PROPÓSITO Y ALCANCE DEL MODELO

El modelo predictivo desarrollado tiene como finalidad estimar el comportamiento futuro del riesgo de incidentes de seguridad industrial en función de registros históricos agregados temporalmente. Su objetivo principal es generar predicciones trimestrales para apoyar la toma de decisiones mediante indicadores anticipados.

La solución se diseñó con enfoque modular, permitiendo su integración con un sistema de visualización (dashboard) y con un servicio de consumo vía API.

2. VARIABLES DE ENTRADA (FEATURES)

VARIABLE	TIPO	DESCRIPCIÓN TÉCNICA
Año	Numérica	Representa el año del periodo analizado. Es escalado para modelado.
Trimestre	Numérica	Representa el trimestre (1 a 4) correspondiente al periodo.

VARIABLE	TIPO	DESCRIPCIÓN TÉCNICA	Página 1 de 1
Total de Reportes	Numérica	Total de incidentes reportados en el trimestre. Se normaliza a rango 0-1.	

3. VARIABLES OBJETIVO (SALIDAS DEL MODELO)

El modelo fue desarrollado bajo un enfoque multitarea (multi-output), generando múltiples predicciones simultáneas. Cada salida corresponde a una distribución probabilística sobre categorías, excepto el riesgo total que se modela como un valor continuo.

SALIDA	TIPO DE PREDICCIÓN	DESCRIPCIÓN
¿Qué?	Distribución (Softmax)	Predice la probabilidad de ocurrencia por tipo de incidente.
¿Quién?	Distribución (Softmax)	Predice la distribución de reportes según clasificación del personal.
¿Dónde?	Distribución (Softmax)	Predice la probabilidad de incidentes según sitio o ubicación operacional.
¿Cómo?	Distribución (Softmax)	Predice la probabilidad de incidentes por categoría de reporte o clasificación interna.
Riesgo Total	Continuo (Sigmoid)	Predictión de un score de riesgo en rango 0 a 1 para el periodo siguiente.

4. PREPARACIÓN Y NORMALIZACIÓN DE DATOS

Los datos fueron previamente depurados y agregados por trimestre. Posteriormente, se aplicaron transformaciones numéricas necesarias para entrenamiento:

PROCESO	DESCRIPCIÓN TÉCNICA
Conversión temporal	Se separó el campo año_trimestre en dos variables numéricas: año y trimestre .
Escalamiento	Se aplicó MinMaxScaler para normalizar el año, trimestre y total de reportes en rango 0 a 1.

PROCESO	DESCRIPCIÓN TÉCNICA	Página 1 de 1
Conversión a proporciones	Las variables categóricas agregadas (tipo, sitio, género y categoría) se transformaron a proporciones respecto al total trimestral, para evitar dependencia directa del volumen.	
Persistencia de scalers	Se almacenaron los escaladores generados en formato .pkl para reproducibilidad del pipeline.	

5. CONSTRUCCIÓN DE SECUENCIAS TEMPORALES (VENTANAS)

Para modelar dependencia temporal se construyeron secuencias tipo ventana deslizante. Cada muestra de entrenamiento contiene un conjunto de periodos consecutivos que permiten al modelo aprender patrones de comportamiento en el tiempo.

PARÁMETRO	VALOR	DESCRIPCIÓN
Window Size	3	Cada predicción utiliza 3 trimestres anteriores como entrada.
Input Shape	(3, 3)	Secuencia de 3 pasos temporales con 3 variables por paso.
Output Shape	Multi-output	Se predicen simultáneamente distribuciones y score total de riesgo.

6. ARQUITECTURA DEL MODELO

Se implementó un modelo basado en redes neuronales recurrentes tipo LSTM, debido a su capacidad para aprender patrones temporales en secuencias de datos históricos. La arquitectura fue definida como una red de entrada única con múltiples cabezales de salida (multi-head).

CAPA	CONFIGURACIÓN	FUNCIÓN
Input	(3,3)	Recibe secuencias temporales escaladas.
LSTM	32 unidades	Aprendizaje de patrones temporales en ventanas trimestrales.
Dense	32 neuronas ReLU	Transformación no lineal de representaciones internas.

CAPA	CONFIGURACIÓN	FUNCIÓN	Página 1 de 1
Dropout	0.15	Reducción de sobreajuste en dataset de tamaño reducido.	
Salidas Softmax	4 cabezales	Predicción de distribuciones probabilísticas por categoría.	
Salida Sigmoid	1 neurona	Score continuo de riesgo total entre 0 y 1.	

7. ENTRENAMIENTO Y CONFIGURACIÓN

PARÁMETRO	VALOR	DETALLE
Optimizador	Adam	Learning rate configurado para estabilidad del entrenamiento.
Learning Rate	0.001	Valor estándar ajustado para convergencia gradual.
Función de Pérdida	MSE	Error cuadrático medio aplicado en todas las salidas del modelo.
Epochs	Hasta 300	Entrenamiento controlado por EarlyStopping.
Validación	33%	Split train/test aplicado con random_state fijo.
EarlyStopping	Patience 25	Se detiene el entrenamiento si la pérdida no mejora.

8. EXPORTACIÓN Y PERSISTENCIA DEL MODELO

Para permitir replicabilidad y despliegue, se generaron archivos persistentes del pipeline completo:

ARCHIVO GENERADO	DESCRIPCIÓN TÉCNICA
modelo_geoinnova_multitask.h5	Modelo entrenado en formato HDF5 listo para despliegue.
X_seq.npy	Secuencias de entrada procesadas para pruebas y predicción.

ARCHIVO GENERADO	DESCRIPCIÓN TÉCNICA	Página 1 de 1
y_*.npy	Targets del entrenamiento por cada salida multitarea.	
scaler_time.pkl	Escalador MinMax para variables temporales (año y trimestre).	
scaler_total.pkl	Escalador MinMax para variable Total_reportes.	
label_map.json	Mapeo de etiquetas de salida para interpretación del modelo.	

9. PREDICCIÓN E INTERPRETACIÓN

La predicción final del modelo genera un score de riesgo y un conjunto de probabilidades por categoría. Para interpretación, se aplica un método de ranking Top-K que permite identificar las categorías más relevantes del siguiente periodo estimado.

ELEMENTO	DESCRIPCIÓN
Riesgo Total	Score entre 0 y 1, utilizado como indicador de riesgo general trimestral.
Top-3 por Categoría	Ranking de las tres categorías con mayor probabilidad para ¿Qué?, ¿Quién?, ¿Dónde? y ¿Cómo?.
Semáforo de Riesgo	Clasificación interpretativa (Bajo, Medio, Alto) basada en rangos del score total.

Declaración de Confidencialidad: Este documento describe únicamente el enfoque metodológico y técnico del modelo predictivo. Los datos originales utilizados para entrenamiento y validación son información confidencial y no se incluyen en este anexo debido a restricciones institucionales y acuerdos de privacidad establecidos durante el Hackathon.