

# Cólicos en Caballos

¿SE PUEDE DETERMINAR LA SUPERVIVENCIA?

---



Alumno: Juan Polo  
Tutor: Francisco Di Palma  
Profesor: Jorge Ruiz

**CODER HOUSE**

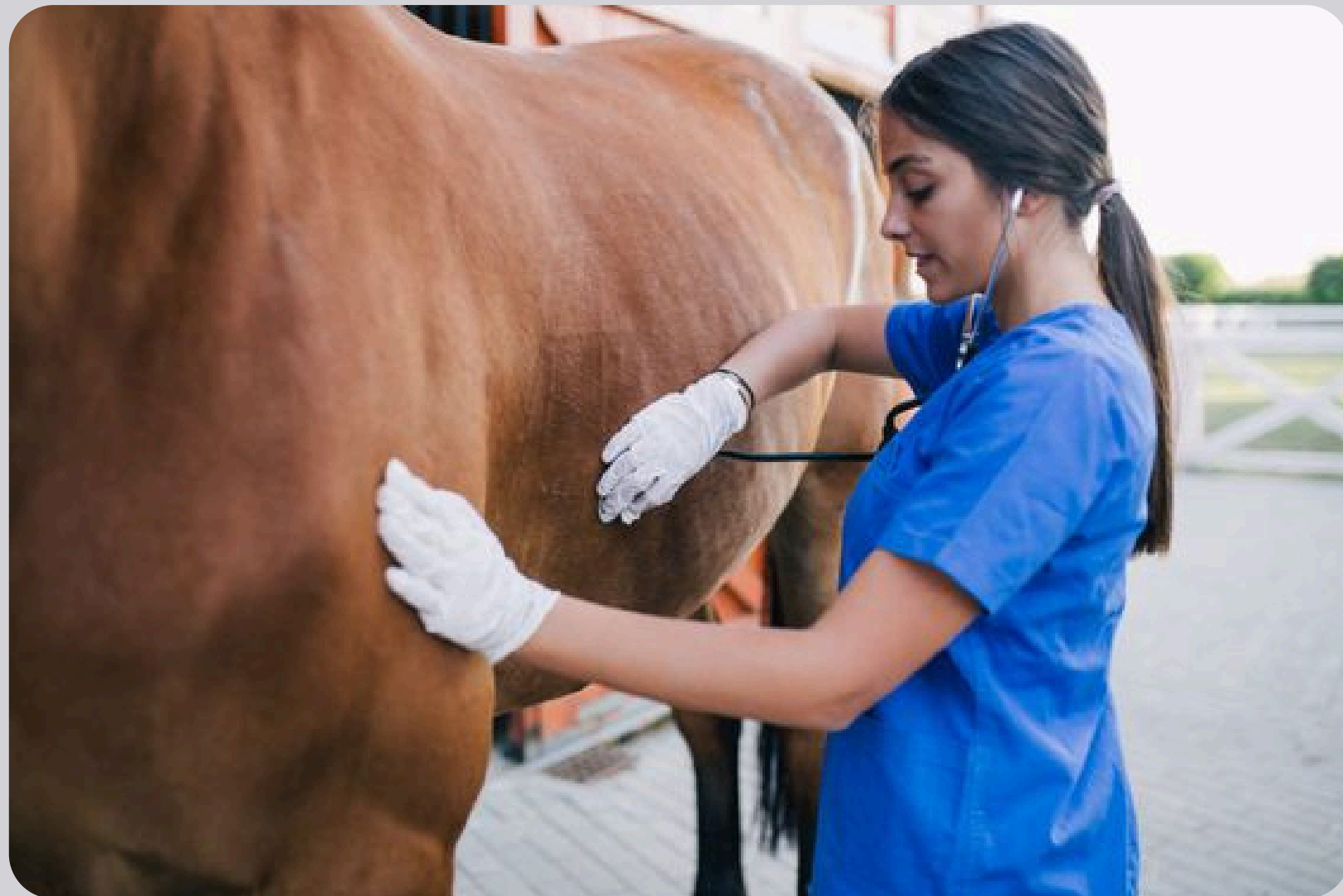
2024

# Cólicos en caballos

Conjunto de síntomas asociados con el dolor abdominal. Es una condición común y puede variar desde un malestar leve hasta una emergencia médica grave que requiere intervención inmediata.

## Sobre los datos

Fueron extraídos de kaggle y agregamientos sintéticos. Contiene 28 variables (numéricas continuas, discretas y categóricas) y 1299 filas. Variable *target* columna 'resultado'.



## OBJETIVO

*Predecir si un caballo vive o muere según sus condiciones médicas.*

## PROPUESTA

Desarrollar un modelo de **clasificación** para estimar la supervivencia de los caballos basado en un conjunto de variables relevantes. Se utilizará técnicas de machine learning (*Random Forest, Regresión Logística*) para predecir la supervivencia de los caballos. El modelo será entrenado, validado y ajustado para maximizar su rendimiento, facilitando la toma de decisiones críticas en contextos veterinarios o de manejo de animales.

# Limpieza y transformacion de los datos

- Concatenacion de los datos para formar un solo dataframe.
- Renombre de columnas para facilitar su entendimiento.
- Reemplazo de valores *nulls* por la moda en variables categoricas y la mediana en variables numéricas.
- Descripción de variables numéricas y categoricas.
- Los valores de la columna 'número\_hospital' no fueron considerados ID y por eso no se efectuo la eliminación de sus valores duplicados.





# ANÁLISIS Y VISUALIZACIÓN DE DATOS

- Análisis univariados de distribuciones y comportamientos de variables numéricas y categoricas.
- Visualización de *outliers* con gráficos de *boxsplot* y *violinplot*.
- Visualización de correlación de variables numéricas en gráfico de calor, *heatmap*.
- Transformación de columna numérica de 'pulso' en una columna categorica 'categoria\_pulso'.
- **IQR** de columna 'pulso' y creación de *features* 'umbral\_inferior' y 'umbral\_superior'.
- Visualización de gráficos multivaridos de columnas categoricas.
- Frecuencia de distribución de variable *target*.





# ANÁLISIS ESTADÍSTICOS

- **Test Chi Cuadrado:** El resultado nos revela 98 filas de  $p\text{-valor} < 0.05$  que rechazan la hipótesis nula, lo que sugiere que hay una relación significativa entre las variables.
- **Coeficiente Pearson:** Los resultados no nos muestran que haya una relación lineal significativa entre las variables numéricas continuas.
- **Isolation Forest:** Utilizado para buscar anomalías, con una contaminación de un 0,05 nos muestra 65 filas que podrían ser datos anómalos. Se tomaron solo los datos sin anomalías (*df\_horse\_sin\_anomalías*), para probar luego de ser encodeado, en mi modelo de *Random Forest*.
- **Coeficiente Biserial:** Los resultados nos muestran que algunas variables rechazan la hipótesis nula y muestran una relación entre la variable dicotómica ('resultado') y algunas variables numéricas continuas.
- **MCA:** Se usó para observar cuáles eran las variables más influyentes en los componentes que creó el análisis.



# MODELOS DE MACHINE LEARNING

1° vistazo de los modelos. Sin utilizar tecnicas de mejoramiento de hiperparámetros. Previamente se encodearon los datos con *OHE* para su utilización en *ML*. Tabien fue probado el dataframe sin anomalias proporcionado por IsolateForest, el cual no presento mejoras en las métricas.

## Random Forest



Accuracy



Recall



Precisión



F1-score

## Regresión Logística



Accuracy



Recall



Precisión



F1-score



# OPTIMIZACION DE MODELOS

- **Curva ROC:** El AUC (Randomforest=0.89, RegresionLogística=0.54) nos indica que el modelo que mejor rendimiento tiene es el de RandomForest
- **RandomOverSampler/SMOTE/TomekLinks/SMOTETomek:** Se utilizo con el modelo de Random Forest para resamplear las muestras, pero ninguno tuvo mejoras significativas en el modelo.
- **XGBoots:** El modelo es secuencial y trabaja tratando de minorizar errores, los resultados de las metricas son: accuracy: , precision: , recall: , f1-score: .
- **Bootstrap:** Su utilizacion en el rendimiento no fue significativa
- **Stratified K-Fold con parametros optimizados de GridSearchCV:** Se utilizo Stratified K-Fold para que el modelo se entrenado correctamente en equilibrio de las clases, y cargue los parametros en GridSercghCV para ver que pruebas consiguian el mejor paramentro, los resultados de las iteraciones nos dieron que en el modelo de Random Forest (El mejor de entre las dos pruebas) los siguientes parametros: max\_depth: 15, n\_estimators: 150.







# Conclusion

La medida F1-score nos hace un balance entre la precisión y el recall. Me quedo con esta medida mas armonica sobre las clases de mi modelo y seria importante conseguir mas datos de campo para trabajar mejor. Hay mucho por hacer y seguro con mas datos se puede mejorar muchos mas la calidad del trabajo.

# Utilizacion

Va dirigido a academicos veterinarios que esten interesados en la agilizar y mejorar su trabajo y en la implementacion de nuevas tecnologias para la mejora de sus procesos de trabajo.

# Muchas gracias por su atencion!

A sido un placer poder exponer un trabajo y aportar a temas de esta indole.