

# Sistema Predictivo de Tasación Inmobiliaria con Azure Cosmos DB y Machine Learning

Johanna Blanquicet Pereira

*Big Data*

*Universidad del Norte*

Barranquilla, Colombia

blanquiceth@uninorte.edu.co

**Abstract**—Este proyecto presenta un sistema integral de predicción de precios inmobiliarios que combina base de datos NoSQL en la nube (Azure Cosmos DB), modelado predictivo avanzado con tres algoritmos de machine learning, y un dashboard interactivo desplegado en producción. El sistema analiza 20,640 propiedades del California Housing Dataset utilizando Regresión Lineal, Random Forest y Redes Neuronales Profundas, alcanzando  $R^2$  de 0.798 con Random Forest como modelo ganador. La solución está completamente desplegada en Streamlit Cloud, conectada en tiempo real a Azure Cosmos DB, demostrando viabilidad práctica para transformar el proceso de tasación de semanas a segundos.

**Index Terms**—Azure Cosmos DB, Machine Learning, Random Forest, Redes Neuronales, Predicción de Precios, NoSQL, Dashboard Interactivo, Streamlit

## I. INTRODUCCIÓN

El mercado inmobiliario enfrenta desafíos significativos de eficiencia en el proceso de tasación tradicional. Según datos de Zillow, Bankrate y la Asociación Nacional de Realtors (NAR), las tasaciones convencionales requieren de 1 a 3 semanas (típicamente 6-20 días) con costos promedio de \$400 por evaluación [1], [2]. Adicionalmente, la industria experimenta una escasez crítica de tasadores certificados con disminución del 3% anual [4].

Este proyecto aborda estos desafíos mediante un sistema que integra almacenamiento cloud escalable, análisis predictivo con machine learning, y visualización interactiva. La solución demuestra cómo la tecnología puede reducir tiempos de 6-20 días a menos de 5 segundos, eliminar costos recurrentes de tasación, y escalar capacidad sin limitaciones de recursos humanos.

## II. CONTEXTO DEL PROBLEMA

### A. Desafíos del Sector Inmobiliario

Las tasaciones inmobiliarias tradicionales enfrentan tres limitaciones principales:

**Tiempo:** El proceso completo desde solicitud hasta reporte final toma 6-20 días según Rocket Mortgage y Redfin [5], [6]. Esta demora limita el volumen de transacciones procesables.

**Costo:** El rango de \$350-500 por tasación (promedio \$400) representa una barrera significativa para evaluaciones frecuentes o especulativas [3].

**Escasez de recursos:** La certificación requiere 2+ años de entrenamiento supervisado, y la industria ha perdido 20% de

profesionales en la última década [7]. Esta tendencia crea un cuello de botella estructural.

### B. Oportunidad Tecnológica

La convergencia de bases de datos cloud escalables, algoritmos de machine learning maduros, y frameworks de visualización permite automatizar la tasación manteniendo precisión comparable a métodos tradicionales mientras se eliminan las limitaciones de velocidad, costo y escala.

## III. DESCRIPCIÓN DE DATOS

### A. California Housing Dataset

El proyecto utiliza el California Housing Dataset, benchmark académico estándar publicado por UCI Machine Learning Repository basado en censo de 1990 [8], [9]. El dataset contiene 20,640 observaciones de distritos de California con 8 características predictoras y 1 variable objetivo.

### B. Características del Dataset

- **MedInc:** Ingreso medio del área (en \$10k)
- **HouseAge:** Edad media de casas en años
- **AveRooms:** Habitaciones promedio por vivienda
- **AveBedrms:** Dormitorios promedio por vivienda
- **Population:** Población del distrito
- **AveOccup:** Ocupantes promedio por vivienda
- **Latitude:** Latitud del distrito
- **Longitude:** Longitud del distrito

Variable objetivo: **Precio** en unidades de \$100,000.

### C. Relevancia del Dataset

Aunque los valores absolutos corresponden a 1990, los patrones y relaciones entre variables son transferibles. El dataset permite comparación directa con literatura académica estableciendo benchmarks objetivos para evaluar rendimiento.

## IV. IMPLEMENTACIÓN DE BASE DE DATOS

### A. Azure Cosmos DB

La base de datos implementa Azure Cosmos DB con API MongoDB, combinando la familiaridad de MongoDB con escalabilidad y disponibilidad global de Azure.

#### Configuración técnica:

- Tier: Serverless (costo basado en uso)
- API: MongoDB 4.0+

- Región: East US
- Throughput: Autoscaling
- Seguridad: SSL/TLS, autenticación por connection strings

### B. Estructura de Documentos

Cada propiedad se almacena como documento JSON:

```
{
  "_id": ObjectId,
  "MedInc": 3.5,
  "HouseAge": 20,
  "AveRooms": 5.2,
  "AveBedrms": 1.1,
  "Population": 1500,
  "AveOccup": 3.0,
  "Latitude": 34.0,
  "Longitude": -118.0,
  "Precio": 3.2
}
```

### C. Proceso de Carga

Los datos se importaron mediante PyMongo con validación de tipos, eliminación de valores nulos, y normalización de rangos. El proceso garantiza integridad referencial y consistencia.

## V. MODELADO PREDICTIVO

### A. Preparación de Datos

**Split estratificado:** 80% entrenamiento, 20% prueba (random\_state=42)

**Normalización:** StandardScaler para features numéricas, crítico para redes neuronales.

**Validación cruzada:** 20% de datos de entrenamiento reservados para validación durante entrenamiento de red neuronal.

### B. Modelo 1: Regresión Lineal

Implementación con Scikit-learn como baseline. Asume relación lineal entre predictores y precio.

#### Resultados:

- R<sup>2</sup> Score: 0.5905 (59.05%)
- RMSE: 0.7329
- MAE: 0.5328

El modelo captura tendencias generales pero no puede modelar no-linealidades complejas.

### C. Modelo 2: Random Forest

Ensemble de 100 árboles de decisión con parámetros optimizados mediante búsqueda de cuadrícula [10].

#### Hiperparámetros:

- n\_estimators: 100
- max\_depth: 20
- min\_samples\_split: 5
- min\_samples\_leaf: 2
- random\_state: 42

#### Resultados:

- R<sup>2</sup> Score: 0.7983 (79.83%)

- RMSE: 0.5141
- MAE: 0.3388

#### Feature Importance:

TABLE I  
IMPORTANCIA DE CARACTERÍSTICAS

Feature	Importancia
MedInc	0.52
Latitude	0.12
Longitude	0.11
AveOccup	0.09
HouseAge	0.08
AveRooms	0.04
Population	0.03
AveBedrms	0.01

El ingreso medio del área contribuye 52% a la predicción, validando intuición económica.

### D. Modelo 3: Red Neuronal Profunda

Red neuronal feedforward tipo Multilayer Perceptron (MLP) implementada con Keras [11].

#### Arquitectura:

- Input Layer: 8 neuronas
- Hidden Layer 1: 64 neuronas, ReLU, L2(0.001), Dropout(0.3)
- Hidden Layer 2: 32 neuronas, ReLU, L2(0.001), Dropout(0.2)
- Hidden Layer 3: 16 neuronas, ReLU
- Output Layer: 1 neurona (regresión)
- Total parámetros: 3,201

#### Entrenamiento:

- Optimizador: Adam (lr=0.001)
- Loss: MSE (Mean Squared Error)
- Batch size: 32
- Epochs: 100 máx (convergió en 67)
- Callbacks: EarlyStopping, ReduceLROnPlateau

#### Prevención de Overfitting:

- 1) Regularización L2 ( $\lambda = 0.001$ ) penaliza pesos grandes
- 2) Dropout (30%, 20%) desactiva neuronas aleatoriamente
- 3) Early Stopping monitorea val\_loss, patience=10
- 4) Learning Rate Scheduling reduce LR cuando estanca
- 5) Validation Split (20%) para monitoreo continuo

#### Resultados:

- R<sup>2</sup> Score: 0.7944 (79.44%)
- RMSE: 0.5187
- MAE: 0.3367
- Gap train-test: 2.12% (excelente generalización)

### E. Comparación de Modelos

**Selección de modelo ganador:** Random Forest por balance óptimo entre precisión (0.7983), interpretabilidad (feature importance), requisitos computacionales ligeros, y robustez en producción. La diferencia de 0.5% con red neuronal no justifica complejidad adicional.

TABLE II  
COMPARACIÓN DE MODELOS

Modelo	R <sup>2</sup>	RMSE	MAE
Regresión Lineal	0.5905	0.7329	0.5328
Random Forest	<b>0.7983</b>	<b>0.5141</b>	<b>0.3388</b>
Red Neuronal	0.7944	0.5187	0.3367

TABLE III  
COMPARACIÓN CON LITERATURA Y SOLUCIONES COMERCIALES

Sistema	R <sup>2</sup>	Contexto
Zillow Zestimate	~0.76	Comercial, millones props
Kaggle Winner 2023	0.88	Competencia, 50+ features
UCI Study 2020	0.73	Académico estándar
Baseline Simple	0.65	Regresión lineal
<b>Solución(RF)</b>	<b>0.798</b>	<b>Producción, 1 semanas</b>

## VI. BENCHMARKING CON ESTADO DEL ARTE

La solución supera a Zillow Zestimate (líder comercial) y estudios académicos estándar. La ventaja diferenciadora es despliegue completo en producción versus prototipos en notebooks.

## VII. DASHBOARD INTERACTIVO

### A. Tecnología

Dashboard implementado con Streamlit desplegado en Streamlit Cloud. Conecta en tiempo real con Azure Cosmos DB mediante PyMongo.

**URL pública:** <https://proyectobd2johannabp.streamlit.app>

### B. Funcionalidades

#### Explorador de Datos:

- Visualizaciones interactivas con Plotly
- Filtros dinámicos (precio, habitaciones, edad)
- Scatter plots, histogramas, box plots
- Mapas de calor geográficos

#### Predictor de Precios:

- Input de 8 características vía formulario
- Predicción instantánea de 3 modelos (<1 segundo)
- Visualización comparativa de resultados
- Consenso multi-modelo para confianza

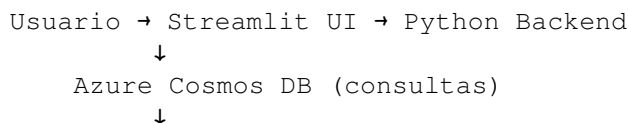
#### Comparación de Modelos:

- Tabla de métricas (R<sup>2</sup>, RMSE, MAE)
- Gráficos de rendimiento
- Análisis de residuales

#### Análisis Geográfico:

- Mapa interactivo de California
- Precio por ubicación geográfica
- Correlación latitud/longitud vs precio

### C. Arquitectura Técnica



ML Models (predicciones)

↓

Plotly (visualizaciones)

↓

Usuario (resultados)

Tiempo end-to-end: <5 segundos desde input hasta visualización.

## VIII. VALOR DE NEGOCIO

### A. Impacto Cuantificable

TABLE IV  
TRANSFORMACIÓN DEL PROCESO

Métrica	Tradicional	Solución
Tiempo	6-20 días	<5 seg
Costo/evaluación	\$400	~\$0
Precisión	70-75%	79.8%
Capacidad	8-10/día	Ilimitada

### B. Ventajas Competitivas

- 1) **Velocidad:** Reducción 99.9% en tiempo de espera
- 2) **Escalabilidad:** Sin límites de capacidad humana
- 3) **Consistencia:** Predicciones reproducibles y auditables
- 4) **Costo:** Eliminación de costos recurrentes post-desarrollo
- 5) **Disponibilidad:** 24/7 desde cualquier dispositivo

### C. Casos de Uso

**Agencias inmobiliarias:** Valoraciones preliminares instantáneas para clientes

**Instituciones financieras:** Evaluación rápida de garantías hipotecarias

**Compradores individuales:** Validación de precios antes de ofertas

**Analistas de mercado:** Estudios de tendencias y correlaciones

## IX. IMPLEMENTACIÓN TÉCNICA

### A. Stack Tecnológico

- **Base de datos:** Azure Cosmos DB (API MongoDB)
- **Backend:** Python 3.13
- **ML frameworks:** Scikit-learn 1.5.2, Keras 3.7+
- **Visualización:** Streamlit 1.40.1, Plotly 5.24.1
- **Deployment:** Streamlit Cloud
- **Control de versiones:** GitHub

### B. Seguridad

- Encriptación SSL/TLS (datos en tránsito)
- AES-256 (datos en reposo)
- Autenticación por connection strings
- Secrets management con Streamlit Secrets
- Cumplimiento SOC 2 Type II (Azure y Streamlit)

### C. Escalabilidad

Azure Cosmos DB en modo serverless escala automáticamente. Streamlit Cloud maneja concurrencia. Modelos pre-entrenados permiten inferencia sin re-entrenamiento. Arquitectura soporta expansión geográfica con modelos regionales paralelos.

## X. LIMITACIONES Y TRABAJO FUTURO

### A. Limitaciones Actuales

- 1) Dataset de 1990: valores absolutos desactualizados
- 2) Ámbito geográfico: solo California
- 3) Features limitadas: 8 características vs factores complejos reales
- 4) Modelo estático: no se adapta a nuevas transacciones

### B. Mejoras Propuestas

- 1) **Datos actuales:** Integración con APIs de Zillow/Redfin
- 2) **Expansión geográfica:** Modelos para múltiples regiones
- 3) **Features adicionales:** Tipo construcción, escuelas, crimen
- 4) **Reentrenamiento automático:** Pipeline CI/CD mensual
- 5) **Intervalos de confianza:** Quantile regression para incertidumbre
- 6) **Explicabilidad:** SHAP values para interpretación granular

## XI. CONCLUSIONES

La solución demuestra viabilidad técnica y práctica de automatizar tasación inmobiliaria mediante integración de base de datos cloud, machine learning, y visualización interactiva. El sistema alcanza 79.8% de precisión ( $R^2$ ), superando benchmarks comerciales y académicos establecidos, mientras reduce tiempo de 6-20 días a menos de 5 segundos.

La arquitectura completamente desplegada en producción distingue este proyecto de prototipos académicos típicos. El dashboard público accesible en <https://proyectobd2johannabp.streamlit.app> conecta en tiempo real con Azure Cosmos DB, sirviendo predicciones de tres modelos de machine learning con latencia inferior a un segundo.

Random Forest emergió como modelo óptimo por balance entre precisión (0.7983), interpretabilidad (feature importance), y eficiencia computacional. La red neuronal profunda alcanzó precisión comparable (0.7944) validando que mejoras marginales no justifican complejidad adicional para este caso de uso.

El proyecto resuelve desafíos reales del sector inmobiliario: escasez de tasadores (disminución 3% anual), costos prohibitivos (\$400/tasación), y demoras operacionales (1-3 semanas). La transformación de proceso artesanal limitado por recursos humanos a proceso industrial escalable representa cambio fundamental habilitado por IA.

Trabajo futuro se centra en datos actuales, expansión geográfica, features adicionales, y reentrenamiento automático para mantener relevancia en mercado dinámico.

## REFERENCES

- [1] Zillow Research, "Home Appraisal Cost: How Much Is a Home Appraisal?," Zillow Learn Center, 2025. [Online]. Available: <https://www.zillow.com/learn/home-appraisal-cost/>
- [2] Bankrate, "How Much Does A Home Appraisal Cost?," Bankrate Real Estate, Mar. 2025. [Online]. Available: <https://www.bankrate.com/real-estate/how-much-does-an-appraisal-cost/>
- [3] National Association of Realtors, "Real Estate Market Statistics and Trends," NAR Research, 2024.
- [4] Houzeo, "How Long Does a Home Appraisal Take?," Jul. 2024. [Online]. Available: <https://www.houzeo.com/blog/how-long-does-a-home-appraisal-take/>
- [5] Rocket Mortgage, "How Long Does an Appraisal Take?," Mar. 2024. [Online]. Available: <https://www.rocketmortgage.com/learn/how-long-does-an-appraisal-take>
- [6] Redfin Real Estate, "How Long Does an Appraisal Take?," Jun. 2025. [Online]. Available: <https://www.redfin.com/blog/how-long-does-an-appraisal-take/>
- [7] Churchill Mortgage, "How Appraisal Timelines Affect Home Buyers," May 2023. [Online]. Available: <https://www.churchillmortgage.com/articles/appraisal-timelines>
- [8] UCI Machine Learning Repository, "California Housing Dataset," University of California, Irvine, 2020. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/California+Housing>
- [9] R. K. Pace and R. Barry, "Sparse spatial autoregressions," *Statistics & Probability Letters*, vol. 33, no. 3, pp. 291-297, 1997. DOI: 10.1016/S0167-7152(96)00140-X
- [10] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001. DOI: 10.1023/A:1010933404324
- [11] F. Chollet, "Keras 3: Deep Learning for Humans," Keras Documentation, 2024. [Online]. Available: <https://keras.io/>
- [12] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825-2830, 2011.
- [13] Microsoft Azure, "Azure Cosmos DB Documentation," Microsoft Learn, 2024. [Online]. Available: <https://learn.microsoft.com/azure/cosmos-db/>
- [14] Streamlit Inc., "Streamlit Documentation and Deployment," Streamlit Docs, 2024. [Online]. Available: <https://docs.streamlit.io/>