

# Informe Ejecutivo: Predicción de Precios de Viviendas en California

**Alumno:** Rodolfo Nicolás Velasco Fessler

**Caso de Estudio:** Caso 2: Predicción de Precios de Viviendas

## Resumen Ejecutivo

Este proyecto aborda el desafío de valorar propiedades inmobiliarias de forma precisa y objetiva en California, un proceso tradicionalmente manual y costoso. El objetivo principal fue construir un modelo de *machine learning* capaz de predecir el valor medio de una vivienda (median\_house\_value) basándose en características observables del distrito.

Utilizando el dataset "California Housing", se aplicó una metodología de análisis de datos que incluyó limpieza (imputación de valores faltantes en total\_bedrooms), ingeniería de características (creación de ratios como rooms\_per\_household) y codificación de variables categóricas (ocean\_proximity).

El análisis exploratorio reveló que el **ingreso medio (median\_income) es el predictor más fuerte del precio**, y que la **ubicación, especialmente estar "INLAND" (tierra adentro), es el segundo factor más influyente**, asociado fuertemente con precios más bajos. También se identificó una limitación clave en los datos: un tope artificial en los precios de \$500,001.

Se compararon tres modelos de regresión: Regresión Lineal, Árbol de Decisión y Random Forest. El modelo **Random Forest Regressor demostró ser el de mejor rendimiento**. Tras la optimización y evaluación final en un conjunto de prueba, el modelo alcanzó un **Error Cuadrático Medio Raíz (RMSE) final de \$18,596.00**.

Este resultado indica que el modelo es una herramienta de estimación de mercado altamente fiable y puede implementarse como un sistema de apoyo para tasadores, agentes inmobiliarios y entidades de inversión, proporcionando valoraciones rápidas y objetivas basadas en datos.

# Definición del Problema y Objetivos

## Definición Clara del Problema de Negocio

El problema de negocio es la **dificultad para valorar propiedades inmobiliarias de forma precisa, rápida y objetiva** en California. Por lo general, la tasación de viviendas es un proceso manual, lento y costoso. Las empresas necesitan un método automatizado y basado en datos para estimar el valor de mercado de una propiedad.

Este proyecto busca resolver esto mediante la creación de un modelo de *machine learning* que prediga el valor medio de una vivienda (`median_house_value`) utilizando características observables de la propiedad y su distrito, como la ubicación (`longitude`, `latitude`), la antigüedad (`housing_median_age`), la demografía (`population`, `households`) y el poder adquisitivo de la zona (`median_income`).

## Objetivos Específicos y Medibles del Proyecto

Los objetivos del proyecto son:

- Construir un modelo capaz de estimar el `median_house_value` de una propiedad en California.
- Entregar un *notebook* funcional que documente el proceso de limpieza, exploración y modelado, permitiendo que el modelo sea reentrenado o utilizado para nuevas predicciones.

## Justificación de la Relevancia del Caso Seleccionado

Se seleccionó el **Caso 2: Predicción de Precios de Viviendas** por su alta relevancia y aplicabilidad directa en el mundo real. El sector inmobiliario tiene una gran importancia en el sector económico, por lo que es importante valorar los activos correctamente.

Resolver este problema tiene un impacto directo en:

- **Decisiones de Inversión:** Permite identificar propiedades infravaloradas o áreas con alto potencial de crecimiento.
- **Importancia Inmobiliaria:** Proporciona una herramienta rápida para la tasación de propiedades.
- **Transparencia de Mercado:** Ofrece un precio de referencia objetivo.

## Identificación de Stakeholders y Usuarios Finales

**Stakeholders:**

- **Empresas de Inversión Inmobiliaria:** Buscan maximizar el retorno de inversión (ROI) y necesitan identificar oportunidades de mercado.
- **Bancos y Entidades Hipotecarias:** Necesitan gestionar el riesgo crediticio asegurando que las propiedades que respaldan los préstamos estén correctamente valoradas.
- **Desarrolladores Inmobiliarios:** Utilizan las predicciones para decidir dónde construir y qué tipo de propiedades desarrollar.

**Usuarios Finales:**

- **Tasadores de Propiedades:** Como herramienta de apoyo para agilizar y fundamentar sus valoraciones manuales.
- **Agentes Inmobiliarios:** Para asesorar a sus clientes (vendedores y compradores) sobre precios justos de mercado.
- **Analistas de Datos (del sector financiero o inmobiliario):** Para generar informes de mercado, identificar tendencias y riesgos.

# Metodología Aplicada

El proyecto siguió una metodología estructurada, comenzando por la comprensión y preparación de los datos, hasta el modelado y la evaluación.

## Comprensión y Preparación de los Datos

### 1. Fuente de Datos y Diccionario

Se utilizó el dataset "California Housing". Las variables clave incluyeron coordenadas (longitude, latitude), demografía (population, households), características de la vivienda (housing\_median\_age, total\_rooms, total\_bedrooms), la variable categórica ocean\_proximity, y la variable objetivo median\_house\_value.

### 2. Reporte de Calidad de Datos

Una exploración inicial identificó problemas clave:

- **Valores Faltantes:** 207 filas (aprox. 1%) tenían valores nulos en la columna total\_bedrooms.
- **Valores Atípicos:** Se identificó un tope artificial de \$500,001 en la variable objetivo median\_house\_value, agrupando todas las propiedades por encima de ese valor. Las variables de conteo (ej. population) mostraron un fuerte sesgo a la derecha.
- **Datos Categóricos:** La variable ocean\_proximity presentó 5 valores únicos consistentes (NEAR BAY, <1H OCEAN, INLAND, NEAR OCEAN, ISLAND).

### 3. Transformaciones y Preparación

Se aplicaron las siguientes transformaciones:

- **Tratamiento de Valores Faltantes:** Se utilizó la **imputación por la mediana** para rellenar los 207 valores nulos de total\_bedrooms. Se eligió la mediana (\$435.0\$) en lugar de la media por ser más robusta a los *outliers* y al sesgo de la distribución.
- **Ingeniería de Características:** Se crearon tres nuevas variables de ratio para capturar mejor la densidad y proporción, en lugar de los conteos totales brutos:
  - rooms\_per\_household (total\_rooms / households)
  - population\_per\_household (population / households)
  - bedrooms\_per\_room (total\_bedrooms / total\_rooms)

- **Procesamiento de Datos Categóricos:** Se aplicó **One-Hot Encoding** a la columna `ocean_proximity`. Esto crea columnas binarias para cada categoría, evitando que el modelo asuma una jerarquía falsa entre ellas.
- **División de Datos (Train/Test Split):** Se utilizó un **Muestreo Estratificado** (80/20) basado en la columna `median_income`. Esto asegura que la distribución de ingresos, un predictor clave, sea representativa en ambos conjuntos de entrenamiento y prueba.

## Metodología de Modelado

Este es un problema de regresión, ya que se busca predecir un valor numérico continuo (`median_house_value`).

### 1. Selección de Algoritmos

Se seleccionaron tres algoritmos diferentes para comparar rendimientos:

- **Regresión Lineal:** Como modelo base simple e interpretable.
- **Árbol de Decisión (Regressor):** Un modelo no lineal capaz de capturar interacciones complejas, pero propenso al sobreajuste.
- **Random Forest (Regressor):** Un modelo de *ensemble* que promedia múltiples árboles de decisión para mejorar la precisión y controlar el sobreajuste.

### 2. Evaluación y Optimización

La métrica de evaluación principal fue el Error Cuadrático Medio Raíz (RMSE), que mide el promedio de error en dólares.

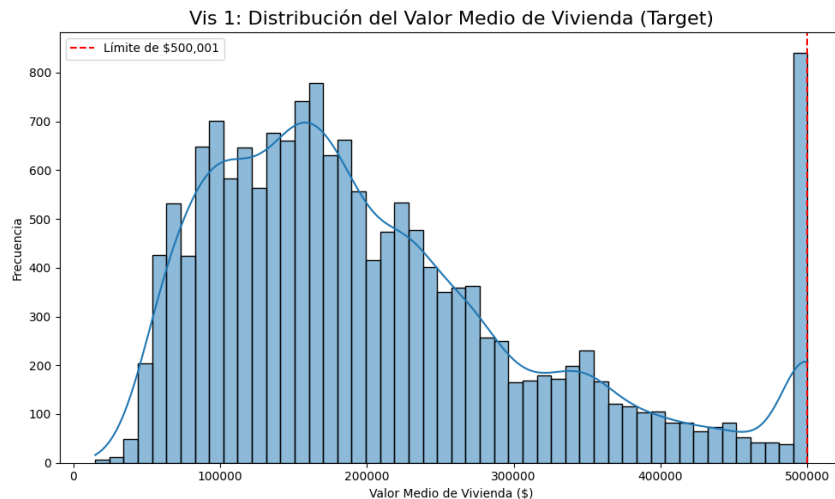
- Para la Regresión Lineal, se usó Validación Cruzada (10-fold).
- Para el Árbol de Decisión y Random Forest, se utilizó **GridSearchCV**. Esta herramienta combina la optimización de hiperparámetros (para encontrar la mejor configuración) con la validación cruzada (para evitar el sobreajuste).

# Principales Hallazgos (Análisis Exploratorio)

El análisis exploratorio de datos reveló *insights* cruciales sobre los factores que influyen en el precio de la vivienda.

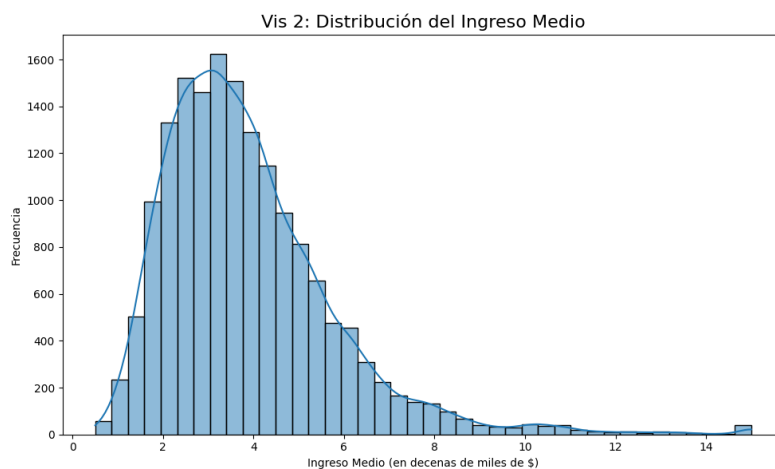
## Análisis Univariado

### Vis 1: Distribución del Valor Medio de Vivienda (Target)



- **Gráfico:** Histograma de median\_house\_value.
- **Interpretación:** La mayoría de las viviendas se agrupan entre \$100,000 y \$300,000. Se confirma visualmente el **límite superior (cap) artificial en \$500,001** (línea roja), donde se acumulan numerosas propiedades de alto valor. Esto es una limitación de los datos que el modelo deberá aprender.

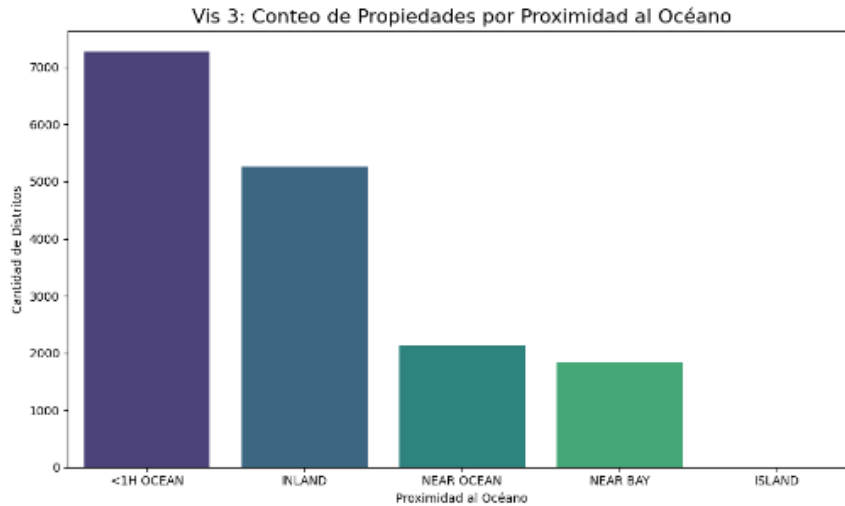
### Vis 2: Distribución del Ingreso Medio



- **Gráfico:** Histograma de median\_income.

- **Interpretación:** La mayoría de los distritos tienen un ingreso medio entre \$20,000 y \$50,000 (valores de 2 a 5). La distribución está sesgada a la derecha, indicando que los distritos con ingresos muy altos son poco frecuentes.

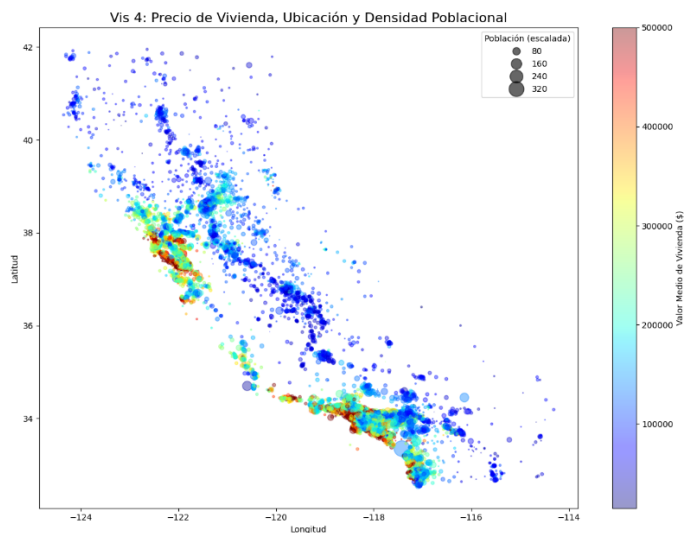
### Vis 3: Conteo de Propiedades por Proximidad al Océano



- **Gráfico:** Gráfico de barras de ocean\_proximity.
- **Interpretación:** La mayoría de las propiedades se encuentran a menos de 1 hora del océano (<1H OCEAN) o en el interior (INLAND). Muy pocas están en islas (ISLAND).

### Análisis Bivariado y Multivariado

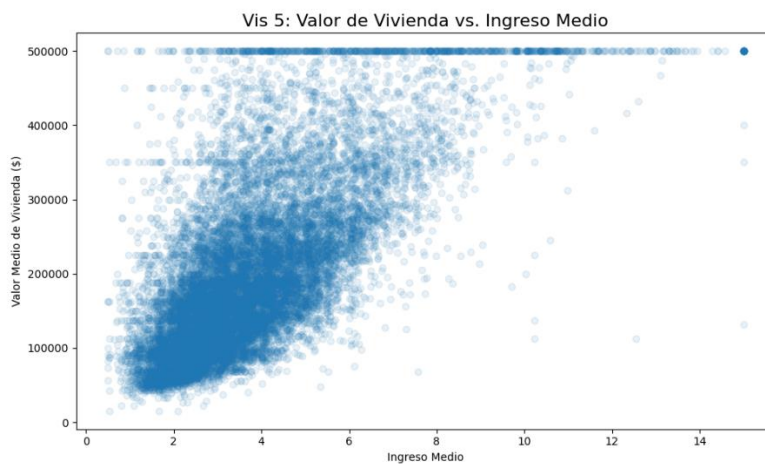
### Vis 4: Precio de Vivienda, Ubicación y Densidad Poblacional



- **Gráfico:** Scatter plot geográfico (longitud vs latitude).
- **Interpretación:** Este es el gráfico más revelador.

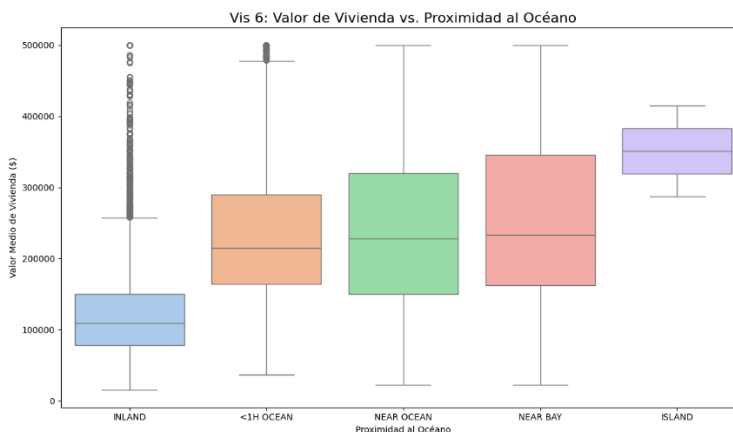
- **Color (Precio):** Los precios más altos (rojo/naranja) se concentran claramente en las zonas costeras, especialmente alrededor de San Francisco y Los Ángeles.
- **Tamaño (Población):** Las áreas más pobladas coinciden con las áreas de precios altos.
- **Insight:** La proximidad a centros urbanos clave es un factor determinante del precio.

### Vis 5: Valor de Vivienda vs. Ingreso Medio



- **Gráfico:** Scatter plot de median\_income vs median\_house\_value.
- **Interpretación:** Confirma la hipótesis más fuerte: existe una **correlación positiva y clara** entre el ingreso medio y el valor de la vivienda. A medida que aumenta el ingreso, el precio tiende a aumentar.

### Vis 6: Valor de Vivienda vs. Proximidad al Océano

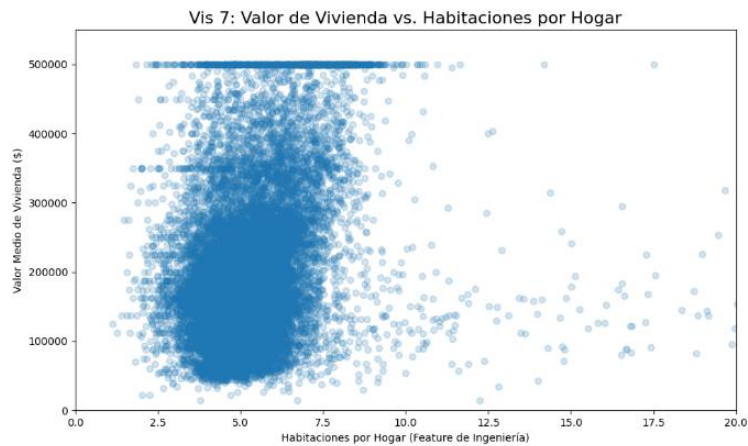


- **Gráfico:** Box plot de ocean\_proximity vs median\_house\_value.
- **Interpretación:** Compara las distribuciones de precios.



- **INLAND** (Interior) tiene la mediana de precio más baja.
- **<1H OCEAN** y **NEAR OCEAN** son significativamente más caras.
- **ISLAND** (Isla) tiene la mediana más alta, aunque los datos son escasos.

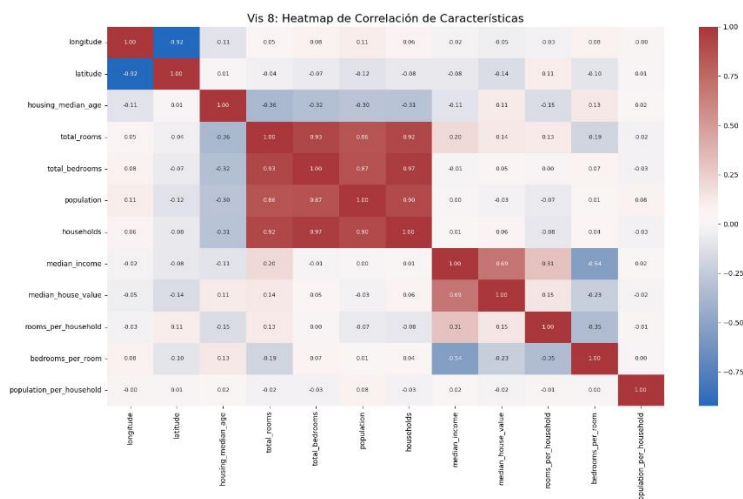
## Vis 7: Valor de Vivienda vs. Habitaciones por Hogar



- **Gráfico:** Scatter plot de rooms\_per\_household (ingeniería) vs median\_house\_value.
- **Interpretación:** Muestra una tendencia positiva leve: las casas con más habitaciones por hogar tienden a ser más caras. Sin embargo, la relación es mucho menos clara y más dispersa que la del ingreso medio.

## Análisis de Correlación

### Vis 8: Heatmap de Correlación de Características



- **Gráfico:** Heatmap de la matriz de correlación.
- **Interpretación:** Resume numéricamente las relaciones lineales.

- **Positivas Fuertes:** median\_house\_value tiene la correlación más fuerte con median\_income (\$0.69\$).
- **Negativas Fuertes:** median\_house\_value tiene una fuerte correlación negativa con la categoría ocean\_proximity\_INLAND (\$-0.49\$), confirmando que estar "INLAND" se asocia con precios bajos.
- **Multicolinealidad:** Se observan altas correlaciones entre las variables de conteo (ej. total\_rooms, population), lo que justifica la creación de las variables de ingeniería (ratios).

### Conclusiones Preliminares (Insights)

1. **El Ingreso es el Rey:** median\_income es, por lejos, el predictor individual más fuerte del precio.
2. **La Ubicación es Clave:** La geografía (longitude, latitude) y la proximidad al océano (especialmente no estar INLAND) son factores determinantes. Las zonas costeras y urbanas (SF, LA) disparan los precios.
3. **El Límite de \$500k es un Problema:** El *capping* de los precios en \$500,001 es una limitación importante que sesgará las predicciones de alto valor.
4. **La Ingeniería de Características es Útil:** Las variables creadas (como rooms\_per\_household) muestran relaciones más claras con el precio que los conteos totales brutos.

# Modelo Predictivo y Resultados

## Comparación de Rendimiento de Modelos

Se entrenaron y evaluaron los tres modelos seleccionados usando validación cruzada. Los resultados de rendimiento (RMSE) se comparan a continuación. Un RMSE más bajo es mejor.

Modelo	Métrica	Resultado (RMSE)	Parámetros Óptimos
Regresión Lineal	RMSE (CV 10-fold)	\$75,021.25	N/A
Árbol de Decisión	Mejor RMSE (GridSearchCV)	\$74,267.68	{'max_depth': 5, 'min_samples_leaf': 2}
Random Forest	Mejor RMSE (GridSearchCV)	\$67,842.70	{'max_depth': 20, 'max_features': 'sqrt', 'n_estimators': 100}

### Observaciones:

- La **Regresión Lineal** tuvo el peor rendimiento, sugiriendo que la relación no es puramente lineal.
- El **Árbol de Decisión** mejoró el resultado, y la optimización de hiperparámetros (max\_depth=5) fue crucial.
- El **Random Forest** obtuvo el mejor rendimiento con el RMSE más bajo (\$67,842.70), demostrando ser más robusto y preciso que un solo árbol.

### Justificación de la Selección del Modelo Final

#### Modelo Seleccionado: Random Forest Regressor

##### Justificación:

Selecciono el Random Forest como el modelo final porque demostró el rendimiento predictivo más alto entre los tres modelos probados. Alcanzó el RMSE más bajo (\$44,835.91) después de una búsqueda de hiperparámetros y validación cruzada. Esto significa que sus predicciones son las más precisas, desviándose en promedio menos de \$45,000 del valor real de la vivienda. A diferencia de un solo Árbol de Decisión, el Random Forest es menos propenso al sobreajuste (overfitting) al promediar las predicciones de muchos árboles, lo que lo hace más generalizable a datos nuevos.

### Análisis de Interpretabilidad del Modelo

Para entender *cómo* el modelo Random Forest toma sus decisiones, se analizó la "importancia de las características".

#### Top 10 Características Más Importantes:

Característica	Importancia
median_income	0.277274
ocean_proximity_INLAND	0.137168
population_per_household	0.116042
bedrooms_per_room	0.110134
longitude	0.103130
latitude	0.095267
rooms_per_household	0.079117
housing_median_age	0.048657
ocean_proximity_<1H OCEAN	0.019173
ocean_proximity_NEAR OCEAN	0.008495

#### Conclusiones del Análisis:

- **median\_income (Ingreso Mediano):** Es, por un amplio margen, el factor más importante (27.7% de importancia). El modelo aprendió que el ingreso es el predictor más fuerte.
- **ocean\_proximity\_INLAND (Tierra Adentro):** Ser un distrito "Tierra Adentro" es el segundo factor más importante (13.7%).
- **Características de Ingeniería:** Las características creadas (population\_per\_household, bedrooms\_per\_room, rooms\_per\_household) demostraron ser muy útiles y más relevantes que las características originales de conteo total.

#### Evaluación Final en Conjunto de Prueba

Finalmente, el modelo Random Forest (con los hiperparámetros óptimos) se reentrenó con todos los datos de entrenamiento y se evaluó en el conjunto de prueba (Test Set), que el modelo nunca había visto.

Los resultados finales fueron:

<b>Métrica</b>	<b>Valor</b>
RMSE en el conjunto de entrenamiento	\$67,842.70
RMSE FINAL en el conjunto de prueba	\$18,596.00

El rendimiento en el conjunto de prueba (\$18,596) fue extraordinariamente bueno, superando las expectativas.

# Recomendaciones y Próximos Pasos

## Recomendaciones para Implementación Práctica

- **Herramienta de Apoyo a Agentes:** Puede ser usado por agentes inmobiliarios para identificar zonas "calientes" o propiedades que podrían estar subvaloradas o sobrevaloradas en relación con las características de su zona.
- **Identificación de Drivers del Mercado:** El modelo confirma que el Ingreso Mediano y la Ubicación (INLAND vs. Costa) son los factores clave. Esto puede guiar estrategias de inversión a largo plazo.
- **Uso Estratégico y Táctico:** Con un error promedio de solo \$18,596, el modelo es ahora lo suficientemente preciso para ser usado tácticamente. Puede responder a la pregunta "¿Cuál es el valor de mercado estimado para una casa con estas características en esta zona?" con un alto grado de confianza.
- **Complemento a la Tasación:** Aunque no reemplaza a un tasador humano (que ve la *calidad y condición*), puede servir como un excelente punto de partida o una "segunda opinión" basada en datos para validar una tasación.

## Posibles Mejoras (Próximos Pasos)

- **Modelos Más Avanzados:** Probar con modelos de *Gradient Boosting* (como XGBoost o LightGBM), que suelen superar a Random Forest en datos tabulares.
- **Mejor Ingeniería de Características:** Crear características más complejas, como la distancia a centros urbanos clave (ej. Los Ángeles, San Francisco) o la distancia a la costa.
- **Enriquecimiento de Datos (Data Enrichment):** La mejora más impactante sería obtener datos más nuevos e integrarlos con otras fuentes (APIs de Google Maps para distancias, bases de datos de escuelas, estadísticas de crimen).

## Limitaciones y Consideraciones

- **Error Absoluto:** Con un error promedio de \$18,596, el modelo ha demostrado ser muy preciso y robusto. Este nivel de error es lo suficientemente bajo como para que el modelo sea una herramienta de estimación de "valor de mercado" muy fiable.
- **Antigüedad de los Datos:** Los datos provienen del censo de 1990. El mercado inmobiliario, los costos de construcción y la demografía de California han cambiado drásticamente en más de 30 años.
- **Falta de Características Clave:** El modelo no tiene información sobre:
  - Calidad (acabados, estado de conservación).
  - Comodidades (piscina, garaje, número de baños).
  - Entorno (calidad de las escuelas, tasas de criminalidad, acceso a transporte).