



---

# DATA SCIENCE

---

PROYECTO FINAL - CODERHOUSE



GOZZARINO NAHUEL

AGOSTO 2023



## INDICE

CONTEXTO COMERCIAL.....	3
OBJETIVOS DEL MODELO.....	3
PREGUNTAS/HIPOTESIS.....	3
EXPLORATORY DATA ANALYSIS (EDA).....	6
RESPUESTAS A LAS PREGUNTAS/HIPOTESIS.....	9
ENCODING .....	9
ELECCION DE MODELOS.....	10
SECUENCIAL FEATURE SELECTION (SFS).....	12
RANDOM FOREST REGRESSOR.....	13
COMPARACION DE MODELOS.....	13
K FOLD CROSS VALIDATION.....	15
OPTIMIZACION DE HIPERPARAMETROS.....	15
CONCLUSIONES GENERALES.....	17



## ABSTRACT.

Este conjunto de datos contiene información del precio de ventas de casas en el condado de King, Washington durante los años 2014 y 2015. Puede ser útil para diversos fines. Por ejemplo, un agente inmobiliario podría utilizar esta técnica para estimar el precio de una casa basándose en sus características, lo que lo puede ayudar a determinar el precio de venta adecuado para una propiedad en particular. También podría ser utilizado por un comprador o vendedor de viviendas para comprender mejor el mercado inmobiliario y tomar decisiones con un criterio objetivo.

Además, los datos y las técnicas de análisis de este conjunto de datos también podrían ser utilizados por profesionales de bienes raíces, urbanistas y planificadores urbanos para comprender mejor los precios de las viviendas en una determinada área geográfica. Esto podría ayudar en la planificación de la expansión urbana y el desarrollo de nuevas viviendas en una zona determinada. En resumen, puede ser útil para aquellos que están interesados en el mercado inmobiliario y quieren utilizar los datos para tomar decisiones informadas en la compra, venta o desarrollo de propiedades.

## CONTEXTO COMERCIAL.

El siguiente trabajo podría ser utilizado por diversas empresas y organizaciones que operan en el mercado inmobiliario, como agentes inmobiliarios, desarrolladores de viviendas, planificadores urbanos y empresas de bienes raíces.

Actuaremos como un científico de datos en una gran inmobiliaria de EEUU que trabaja en conjunto con el area de ventas. Se analizarán los precios de las casas en california como así también las características de las mismas.

## CONTEXTO ANALITICO.

Se enfocará el analisis en la identificación de los factores que tienen el mayor impacto en el precio de venta de las casas, como el número de habitaciones, la ubicación, el tamaño del terreno, entre otros. Esto puede ayudar a los actores del mercado inmobiliario a tomar decisiones informadas sobre la compra, venta y desarrollo de propiedades en el condado de King, Washington.

El problema principal, de la mano de nuestro objetivo principal será **construir 3 modelo de aprendizaje supervisado y elegir el que nos permita predecir mejor cuanto será el precio de venta de una determinada casa**, de acuerdo a las características de la misma.

## OBJETIVOS DEL MODELO.

En este caso, el objetivo será comunicar los resultados del análisis de datos de los precios de las viviendas en King. Nos centraremos en las variables más importantes que influyen en el precio de venta. Analizaremos si existe algun tipo de relación de una variable respecto a otra. Si este tipo de dependencia existe, veremos de que forma se da esta relación.

## PREGUNTAS/HIPOTESIS.

- ¿De qué categoría son las casas que más se venden? ¿Las que menos se venden?
- ¿Las casas con una buena vista son las más costosas?
- ¿Cuáles son las casas más económicas? ¿Que características tienen?

En base a estas preguntas, intentaremos descubrir aquellas variables que mayor impacto tengan en el precio de una casa.

## DESCRIPCION DE LOS DATOS.

El conjunto de datos que se utilizará contiene información del precio de ventas de casas en el condado de King, Washington durante los años 2014 y 2015.

El mismo contiene 21 columnas:

- ID: Correlativo numérico que distingue cada una de las casas vendidas.
- Date: Fecha en la cual fue vendida la casa.
- Price: Precio de venta de la casa (en dólares norteamericanos).
- Bedrooms: Número de cuartos en la casa.
- Bathrooms: Número de baños disponibles en la casa, donde un valor de 0.5 representa un baño con inodoro pero sin ducha.
- Sqft\_living: Número de pies cuadrados del espacio habitable de la casa.
- Sqft\_lot: Pies cuadrados del espacio total del terreno donde se ubica la casa.
- Floors: Cantidad de pisos en la casa.
- Waterfront: Variable que indica la presencia o no de vista al mar en la casa.
- View: Índice del 0 al 4 que indica que tan buena es la vista de la propiedad.
- Condition: Índice del 1 al 5 para calificar la condición actual de la casa.
- Grade: Índice del 1 al 13, el cuál califica el nivel de calidad de construcción de la casa.
- Sqft\_above: Cantidad de pies cuadrados del espacio interior de la casa. que está sobre el nivel del suelo.
- Sqft\_basement: Los pies cuadrados del espacio interior de la casa. que está por debajo del nivel del suelo.
- Yr\_built: El año en fue construida la casa.
- Yr\_renovated: Año de la última renovación de la casa.
- Zipcode: Código postal del área donde se encuentra la casa.
- Lat: Latitud de la ubicación de la casa.
- Long: Longitud de la ubicación de la casa.
- Sqft\_living15: Los pies cuadrados de espacio habitable de la casa interior para los 15 vecinos más cercanos.
- Sqft\_lot15: Los metros cuadrados de los terrenos de los 15 vecinos más cercanos.

A continuación, se muestra una tabla resumen de todas las variables, con información de nulos, tipos de datos, cantidad de valores únicos y un ejemplo de cada variable.

Column	Type	Non-Null	Nulls	Unique	Example
id	int64	21613	0	21436	7129300520
date	object	21613	0	372	20141013T000000
price	float64	21613	0	4028	221900.0
bedrooms	int64	21613	0	12	3
bathrooms	float64	21613	0	29	1.0
sqft_living	int64	21613	0	1038	1180
sqft_lot	int64	21613	0	9782	5650
floors	float64	21613	0	6	1.0
waterfront	int64	21613	0	1	0
view	int64	21613	0	4	0
condition	int64	21613	0	5	3
grade	int64	21613	0	12	7
sqft_above	int64	21613	0	946	1180
sqft_basement	int64	21613	0	305	0
yr_built	int64	21613	0	116	1955
yr_renovated	int64	21613	0	69	0
zipcode	int64	21613	0	70	98178
lat	float64	21613	0	5034	47.5112
long	float64	21613	0	752	-122.257
sqft_living15	int64	21613	0	777	1340
sqft_lot15	int64	21613	0	8689	5650

Debido a que no realizaremos ningún análisis enfocado en fechas, la columna "date" fue eliminada del dataset. Del mismo modo, la columna "id" no aporta información relevante, por lo que también fue descartada de nuestro análisis. Resulta interesante destacar que no tenemos valores nulos ni tampoco filas duplicadas.

## EXPLORATORY DATA ANALYSIS (EDA).

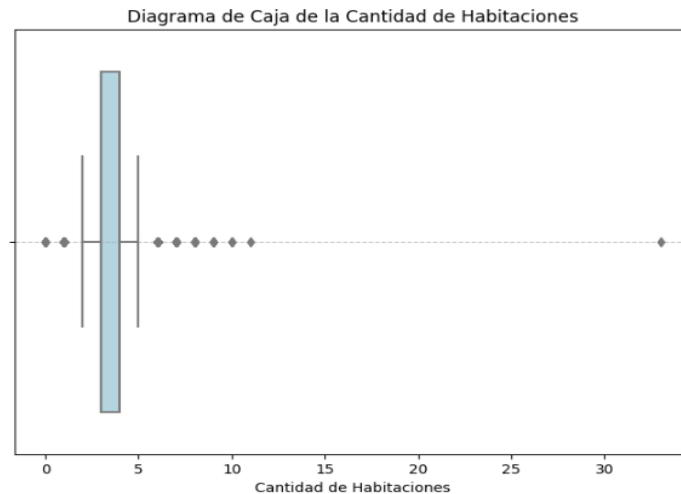
	count	mean	std	min	25%	50%	75%	max
price	21613.00	540088.14	367127.20	75000.00	321950.00	450000.00	645000.00	7700000.00
bedrooms	21613.00	3.37	0.93	0.00	3.00	3.00	4.00	33.00
bathrooms	21613.00	1.75	0.73	0.00	1.00	2.00	2.00	8.00
sqft_living	21613.00	2079.90	918.44	290.00	1427.00	1910.00	2550.00	13540.00
sqft_lot	21613.00	15106.97	41420.51	520.00	5040.00	7618.00	10688.00	1651359.00
floors	21613.00	1.45	0.55	1.00	1.00	1.00	2.00	3.00
waterfront	21613.00	0.01	0.09	0.00	0.00	0.00	0.00	1.00
view	21613.00	0.23	0.77	0.00	0.00	0.00	0.00	4.00
condition	21613.00	3.41	0.65	1.00	3.00	3.00	4.00	5.00
grade	21613.00	7.66	1.18	1.00	7.00	7.00	8.00	13.00
sqft_above	21613.00	1788.39	828.09	290.00	1190.00	1560.00	2210.00	9410.00
sqft_basement	21613.00	291.51	442.58	0.00	0.00	0.00	560.00	4820.00
yr_built	21613.00	1971.01	29.37	1900.00	1951.00	1975.00	1997.00	2015.00
yr_renovated	21613.00	84.40	401.68	0.00	0.00	0.00	0.00	2015.00
zipcode	21613.00	98077.94	53.51	98001.00	98033.00	98065.00	98118.00	98199.00
lat	21613.00	47.56	0.14	47.16	47.47	47.57	47.68	47.78
long	21613.00	-122.21	0.14	-122.52	-122.33	-122.23	-122.12	-121.31
sqft_living15	21613.00	1986.55	685.39	399.00	1490.00	1840.00	2360.00	6210.00
sqft_lot15	21613.00	12768.46	27304.18	651.00	5100.00	7620.00	10083.00	871200.00

El rango de precios es muy variable. Tenemos casas cuyo precio de venta es \$75.000 USD, mientras que la más costosa supera los 7 millones de dólares. Cabe destacar que el 75% de las casas no superan los \$645.000 USD, por lo que existe una minoría de casas que presentan precios muy elevados respecto a la gran mayoría.

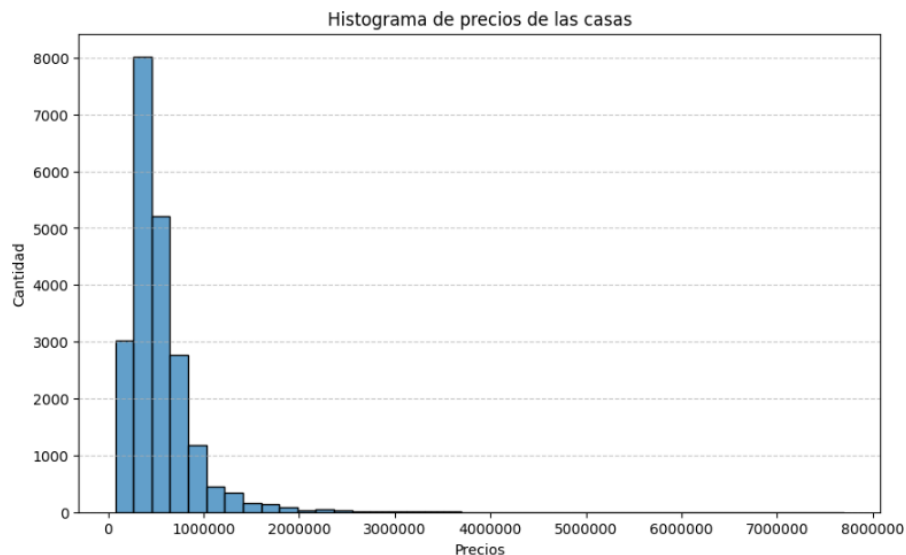
En promedio, el precio de venta ronda los \$540.000 USD, las casas cuentan con 3 habitaciones, 2 baños y 1 piso.

**OBSERVACIÓN 1:** Hay casas sin habitaciones y/o sin baño. Supondremos que mínimamente, una casa debe contar con un baño. Aceptaremos valores ceros en la cantidad de habitaciones para incluir aquellas viviendas de tipo "monoambiente". Por lo tanto, para la cantidad de baños, reemplazaremos los valores ceros con la mediana de dicha columna.

**OBSERVACION 2:** Hay por lo menos una casa donde los datos indican que posee 33 habitaciones. Esta cantidad se considerará como un outlier y se eliminará del dataset original.



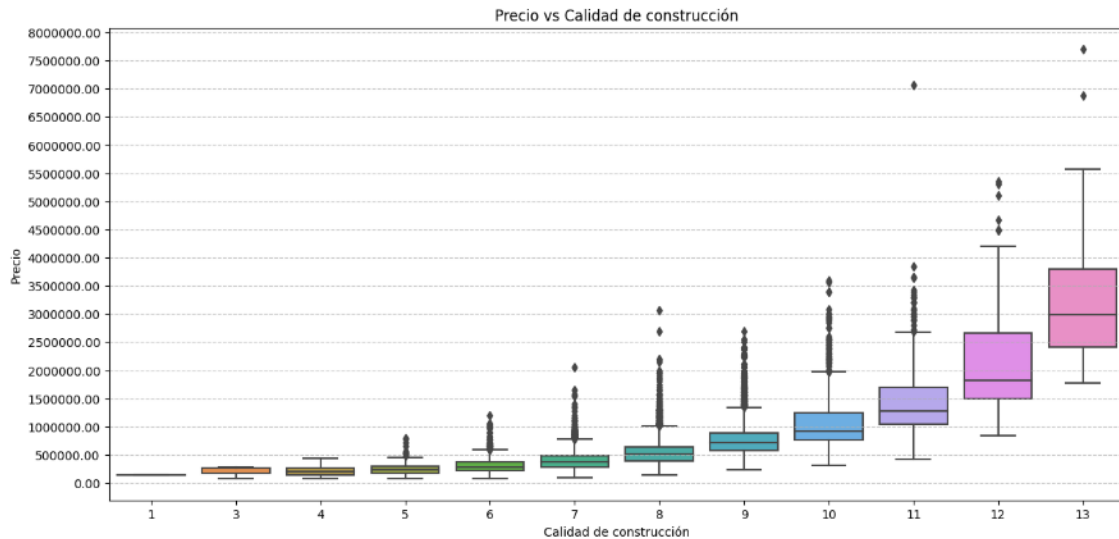
Tenemos 62 casas con una cantidad mayor a 6 habitaciones. Las cuales, a simple vista, considerando el precio y la cantidad de baños, pueden tratarse de casas de gran tamaño. Por lo que se mantendrán en el análisis.



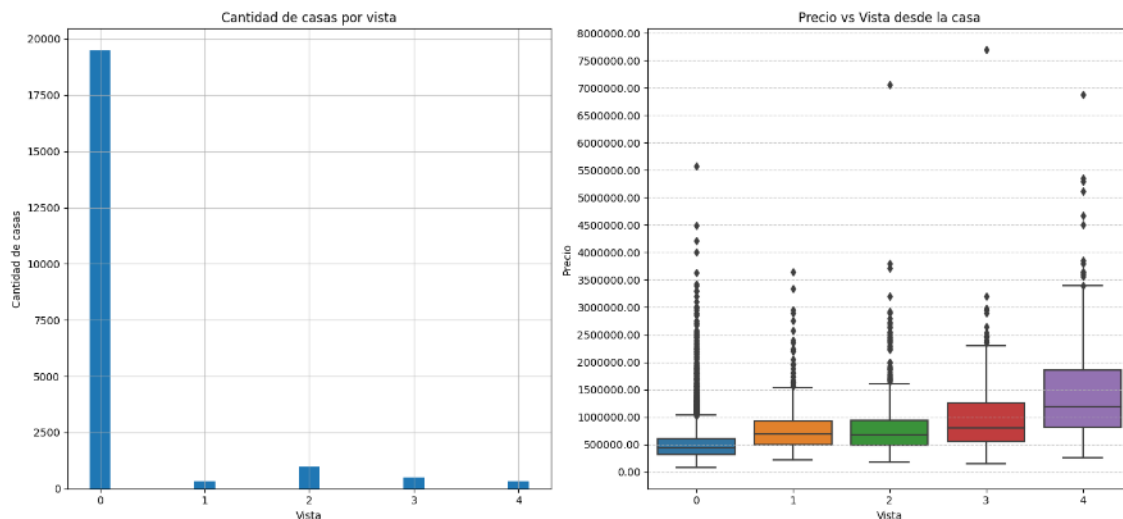
En el histograma de precios, los datos están sesgados a la izquierda. Es decir, la gran mayoría de las casas que se venden, se venden a menos de un millón de dólares. Mientras que hay una proporción menor de casas que se venden por encima del millón de dólares.

Si analizamos la variable precio en relación a otras, podemos llegar a las siguientes afirmaciones:

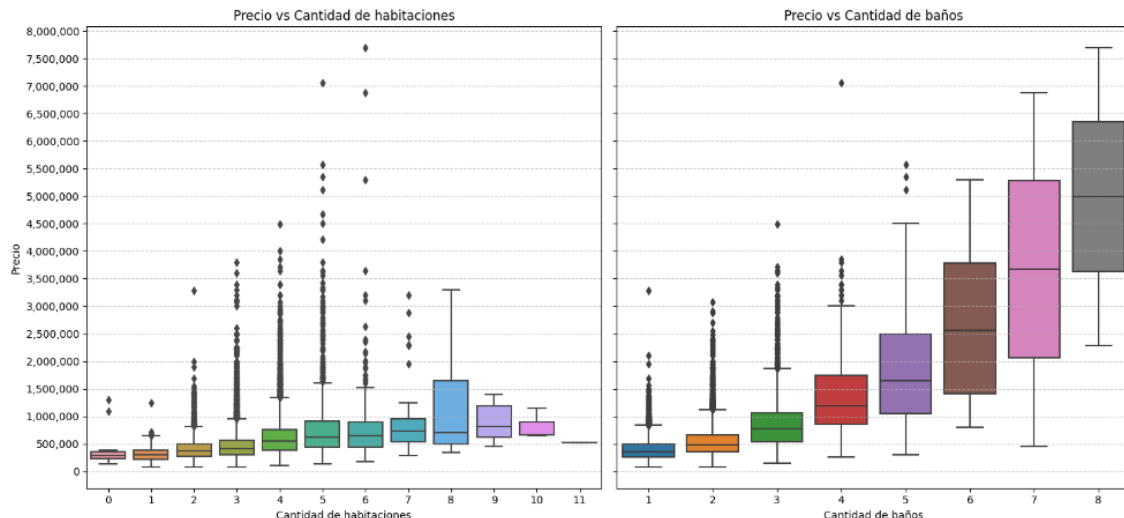




Teniendo en cuenta que "grade" es un índice del 1 al 13 que califica el nivel de calidad de construcción de la casa, podemos observar que, respecto a la calidad, el precio mantiene su media hasta una calidad grado 7. Luego, comienza a aumentar levemente la media a medida que la calidad aumenta. Se observan valores por encima del límite superior de cada boxplot. Esto no necesariamente implica que sean outliers, ya que independientemente de la calidad de construcción, estas casas podrían tener características únicas y privilegiadas, lo que justificaría los precios altos para cada nivel.



La gran mayoría de las casas están clasificadas con una categoría de vista "0". Respecto al precio, las pocas casas cuya vista es mayor a 2, poseen una media en el precio más elevada que en el resto de las casas. Esto debe surgir por la ubicación privilegiada. Que a pesar de la calidad de construcción que tenga, le aumenta su precio de venta.



Para la cantidad de habitaciones y la cantidad de baños, sucede algo parecido a la calidad de la construcción. Si bien se cumple la correlación positiva del precio con ambas variables, la media del precio aumenta de manera más exponencial respecto a la cantidad de baños.

## RESPUESTAS A LAS PREGUNTAS/HIPOTESIS.

La gran mayoría de las casas vendidas se concentran en la vista de calidad "0". Esto es debido a que la categoría cero posee en promedio el precio de venta mas bajo.

En las vistas 1,2 y 3 la media del precio es muy parecida, mientras que para la vista número 4 se observa un claro aumento en la media del precio de venta. Por lo tanto, podemos afirmar y validar que la vista de una casa de King tiene una correlación positiva (0,40) respecto al precio. A su vez, podemos afirmar que las cantidades de casas que se venden con una vista privilegiada son menores a las cantidades de casa vendidas con vista calidad 0,1,2 y 3.

La categoría de construcción de cada casa se representa mediante un índice del 1 al 13. Donde 1 es el nivel más bajo en cuanto a categoría y 13 es la categoría mas alta.

La mayoría de las casas se concentran en las categorías siete, ocho y nueve. Respecto al precio de venta, se puede apreciar un pequeño crecimiento en la mediana del precio de venta en las casas con mayor categoría de construcción, sin embargo, hay ciertas categorías que engloban muy poca cantidad de casas. Si bien en un principio resultaba lógico pensar que las casas más baratas serían las que peores vistas tengan y las que estaban construidas con materiales de no tan buena categoría, el EDA permitió confirmar esto.

## ENCODING

No se realizó encoding ya que en nuestro dataset no tenemos variables categóricas que necesiten ser codificadas.

## TECNICAS DE SELECCIÓN Y REDUCCION DE FEATURES.

Con el objetivo de reducir la dimensionalidad del dataset, se profundizó en:

- PCA (Principal Component Analysis)
- Secuencial Feature Selector - Cross Validation

PCA para reducir la dimensionalidad y transformar las variables originales en un nuevo conjunto de variables no correlacionadas llamadas componentes principales.

Sequential Feature Selector (SFS) con Cross Validation para encontrar el subconjunto óptimo de variables predictoras.

## ELECCION DE MODELOS.

Se crearon 2 modelos de **regresión lineal**. El primero será entrenado con las componentes principales de PCA, mientras que al segundo lo entrenaremos con el conjunto optimo seleccionado a través de SF. Por ultimo, crearemos un tercer modelo de Random Forest Regressor.

Los 3 modelos serán evaluados y se validará aquel con mejores métricas.

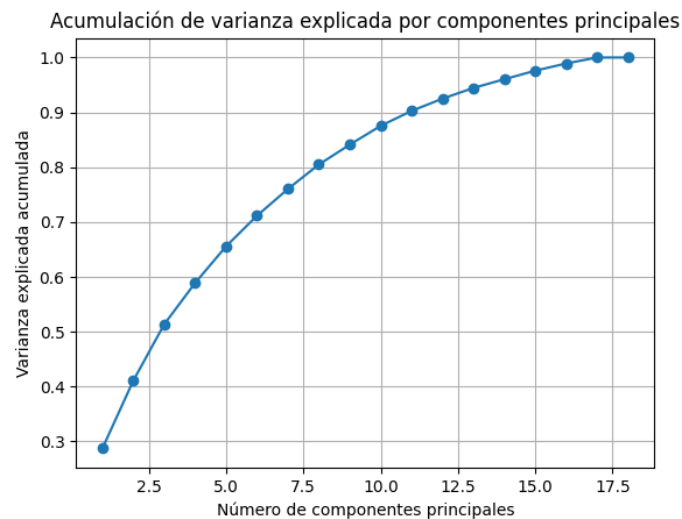
Nuestra variable dependiente en cada modelo será el precio de venta de las casas.

## PRINCIPAL COMPONENT ANALYSIS (PCA).

Luego de aplicar PCA, se observó que, con 9 componentes principales, obtenemos un 90% de porcentaje de varianza explicada. Por lo tanto, con el objetivo de reducción de dimensionalidad y a la vez de conservación de información relevante, en una primera instancia, se tomaron las primeras 9 componentes para evaluar el rendimiento de nuestro primer modelo. Sin embargo, se realizaron una serie de pruebas considerando diferentes cantidades de componentes principales con el objetivo de optimizar los valores con los cuales evaluamos el modelo (MSE, RMSE, MAE y R2).

14 componentes principales fue la cantidad que optimizó las métricas en nuestro modelo de regresión lineal. Cabe destacar que con 18 PCA los valores de MSE y R2 mejoraban muy levemente. Por lo que se decidió optar por 14 con el objetivo de trabajar con menor cantidad de features y reducir la dimensionalidad de nuestro dataset original.

El modelo entrenado con este conjunto de variables fue llamado “regressor\_pca”



## SECUENCIAL FEATURE SELECTION (SFS).

Este método se realizó con el objetivo de obtener un subconjunto óptimo de características (variables predictoras) a partir de nuestro dataset original y así mejorar la eficiencia computacional y disminuir los requerimientos de memoria. Se aplicó la técnica Sequential Feature Selector (SFS) con Cross Validation. Esto permitió reducir el riesgo de sobreajuste al eliminar características irrelevantes y reducir la cantidad de características utilizadas en el modelo.

El valor de  $R^2$  fue el utilizado para seleccionar gradualmente el conjunto óptimo de características que mejoran el rendimiento del modelo.

Posteriormente, se compararon los resultados con los valores obtenidos de las predicciones realizadas utilizando las primeras 14 componentes principales.

En el método SFS se utilizó un rango de 2 a 18 con el objetivo de evaluar el modelo de regresión lineal con todas las combinaciones posibles que nos arroje SFS. Es decir, se evaluó el modelo con 2 features obtenidas a través de SFS. Luego con 3, y así sucesivamente hasta testear el modelo con las 18 features disponibles en nuestro dataset original.

Luego de analizar las iteraciones, las métricas alcanzaron sus mejores valores a partir del index 11. Luego, desde la iteración 11 a la 16, los valores no presentaron mejoras significativas. Por lo que asumimos que el conjunto de features con las que nuestro modelo performa de mejor manera son:

- Bedrooms
- Bathrooms
- sqft\_living
- waterfront
- view
- condition
- grade
- sqft\_basement
- yr\_built
- yr\_renovated
- zipcode
- lat
- long

El modelo entrenado con este conjunto de variables fue llamado “regressor\_sfs”

## RANDOM FOREST REGRESSOR.

Para el modelo de random forest, el modelo se ejecutó en primer lugar con sus hiperparametros configurados por defecto.

El modelo entrenado fue llamado “regressor\_rf”

## COMPARACION DE MODELOS.

Las métricas con la que se evaluaron los 3 modelos arrojaron los siguientes resultados:

	Modelo	Error Cuadrático Medio (MSE)	Raíz del Error Cuadrático Medio (RMSE)	Error Absoluto (MAE)	Coefficiente de Determinación (R2)
0	regressor_pca	47632432137.065	218248.556	130142.285	0.683
1	regressor_sfs	45098130404.422	212363.204	127476.317	0.700
2	RandomForestRegressor	21392969164.105	146263.356	72825.476	0.857

Comparando los modelos de regresión lineal, el modelo regressor\_sfs parece tener un mejor rendimiento en comparación con el modelo regressor\_pca para predecir los precios de las casas.

Considerando que el MSE mide el promedio de los errores al cuadrado entre las predicciones del modelo y los valores reales, y que un MSE más bajo indica un mejor ajuste del modelo a los datos, regressor\_sfs posee un mejor MSE.

En ambos modelos, el coeficiente de determinación posee valores muy parecidos. Esto indica que el 70% de la variabilidad de los precios de las casas se explica a través de ambos modelos. El modelo tiene cierta capacidad para explicar y predecir los precios de las casas pero también significa que aproximadamente el 30% de la variabilidad en los datos no se explica por el modelo. Esto implica que todavía hay factores adicionales o fuentes de variabilidad que no se están teniendo en cuenta y que podrían influir en los precios de las casas.

En cuanto al RMSE y el Error absoluto, son métricas que nos indican qué tan cerca están las predicciones obtenidas de los valores reales. Según las métricas observadas y comparando ambos modelos de regresión lineal, en el regressor\_sfs las predicciones están más cerca de los valores reales.

Sin embargo, **RandomForestRegressor obtuvo valores superadores en todas las métricas.** Se observan valores inferiores de MSE y MAE. Lo que indica que se generaron predicciones más precisas y cercanas a los valores reales y el R2 es el más alto de todos los modelos (86%), por lo que el modelo explica una mayor cantidad de variabilidad en los datos.

**Aclaración:** Cabe destacar que normalmente se exigen valores de MAE, MSE o RMSE muy cercano a cero, pero no siempre es así, ya que en nuestro dataset los precios de las casas están cercanas al millón de dólares y la mayoría por encima de los cientos de miles ( $10^6$ ,  $10^7$ ). Por lo tanto, tener un MAE en el orden de los 72.800 no es tan descabellado.



## K FOLD CROSS VALIDATION.

Si bien el modelo regressor\_rf parece performar bien en los datos actuales, aplicar la validación cruzada permitió realizar una validación más sólida y confiable. Con esta técnica se intentó mitigar los riesgos de sobreajuste (overfitting) y subajuste (underfitting), dos problemas comunes en el desarrollo de modelos.

Los scores R2 de cada iteración de la validación cruzada arrojaron valores por encima de 0.86 y al ser un valor cercano a 1, esto indicó un buen ajuste del modelo a los datos.

El promedio de todos los R2 fue de 0.87. Lo cual indica que el modelo generalizó bien y tuvo un buen rendimiento en datos no vistos.

Dicho esto, **validamos este modelo y lo consideramos como el mejor de los 3 modelos creados.**

## OPTIMIZACION DE HIPERPARAMETROS.

Una vez validado nuestro modelo, se aplicó la técnica de optimización Randomized Search.

Principalmente para encontrar la combinación de hiperparámetros que maximice el rendimiento del modelo y para mejorar la eficiencia computacional.

El resultado de la búsqueda aleatoria fue:

- n\_estimators: 300
- min\_samples\_split: 2
- min\_samples\_leaf: 1
- max\_features: "sqrt"
- max\_depth: 40

Posterior a esto, se creó un último modelo optimizado y se realizaron predicciones con estos nuevos hiperparametros. Los resultados obtenidos fueron los siguientes:

Error cuadrático medio (MSE) : 21.411.069.166,22

Raíz del error cuadrático medio (RMSE): 146.325,21

Error absoluto: 74.519,15

Coefficiente de determinación (R2): 0,86



## CONCLUSIONES ACERCA DE LOS MODELOS.

El modelo RandomForestRegressor con hiperparámetros optimizados arrojó métricas muy parecidas a los valores obtenidos previamente sin estos hiperparámetros. Esto puede deberse a que el modelo puede estar alcanzando sus límites de rendimiento o debido a variaciones aleatorias en los datos o el proceso de entrenamiento.

RandomForestRegressor obtuvo valores superadores en todas las métricas. Tanto para el modelo optimizado como para el modelo sin hiperparámetros optimizados. Respecto a los modelos de regresión lineal se observaron valores inferiores de MSE y MAE. Lo que indica que se generaron predicciones más precisas y cercanas a los valores reales y el R2 fue el más alto de todos los modelos (86%), por lo que el modelo explica una mayor cantidad de variabilidad en los datos.

**Se validó este modelo como el mejor de los modelos entrenados.**

A continuación, se visualizan las 10 mejoras predicciones alcanzadas con nuestro modelo.

	Valor Actual	Valor Predicho	Error Absoluto (MAE)	Error Cuadrático (MSE)
1061	1035480.000000	1035467.740000	12.260000	150.307600
640	378000.000000	378063.200000	63.200000	3994.240000
928	359000.000000	359078.500000	78.500000	6162.250000
4691	423000.000000	423086.400000	86.400000	7464.960000
20735	275000.000000	274909.600000	90.400000	8172.160000
8552	315000.000000	315107.400000	107.400000	11534.760000
17654	450000.000000	450110.000000	110.000000	12100.000000
5282	515000.000000	515131.260000	131.260000	17229.187600
6544	532500.000000	532346.166667	153.833333	23664.694444
17901	219000.000000	218839.500000	160.500000	25760.250000

Más allá del modelo elegido y de los buenos resultados obtenidos, resultó interesante tener en cuenta las variables obtenidas con secuencial feature selector, para **conocer el negocio y entender el mercado**. Ya que gran parte del precio de una casa viene representado por estas variables.

Es fundamental considerar el tamaño de la propiedad en conjunto con la cantidad de baños y habitaciones de cada casa ('bedrooms', 'bathrooms') Los inversores deben buscar propiedades de un tamaño adecuado para el mercado objetivo, teniendo en cuenta las preferencias de los compradores potenciales.

Los inversores deben considerar la presencia de frente de agua en la propiedad, ya que esto puede aumentar significativamente el valor y la demanda.

Se debe prestar atención a la vista. La variable 'view' también tiene una gran influencia en el precio de las casas. Los inversores deben buscar propiedades con vistas impresionantes, como vistas al mar, a la montaña, a un lago, etc.

Se debe considerar la calidad de la propiedad. La variable 'grade' se refiere a la calidad general de la propiedad, por lo que es importante que los inversores se centren en propiedades de alta calidad para maximizar su potencial de reventa.

La variable 'lat' tiene una gran influencia en el precio de las casas. Por lo tanto, es importante que los inversores se centren en áreas de King County que tengan una ubicación privilegiada, con buenas conexiones de transporte, acceso a servicios y comodidades cercanas, y atractivos paisajes.

Resulta necesario comprender la edad de la propiedad, el año en el que fue renovada, su latitud y longitud. Los inversores deben tener en cuenta la edad de la propiedad al evaluar su potencial de inversión. Una propiedad antigua y con una mala ubicación puede requerir más reparaciones y actualizaciones, lo que puede afectar su rentabilidad.

## CONCLUSIONES GENERALES.

Se cumplieron todos los objetivos propuestos al principio de este trabajo final integrador. Los contenidos teóricos fueron llevados a la práctica y se profundizó en cada uno de los diferentes temas desarrollados a lo largo del curso dictado.

El trabajo final fue fundamental para tener una perspectiva clara de lo que implica trabajar en la actualidad como Data Scientist.

Se logró cumplir con lo solicitado en la consigna, e incluso se obtuvieron muy buenos resultados.