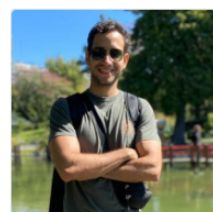


# **Predicción de costo de venta de una propiedad**

# RESUMEN DEL TRABAJO

## Perfil Kaggle

https://www.kaggle.com/amilcarrodriguez



**Amilcar Rodriguez**

[Add organization](#)

Electrical Engineer

Santiago, Santiago Metropolitan Region, Chile

Joined 5 months ago · last seen in the past day

[in](#)

[Home](#) [Competitions \(1\)](#) [Datasets](#) [Code \(1\)](#) [Discussion](#) [Followers](#) [Notifications](#) [Account](#)


## Integrantes

Amilcar Rodriguez	<a href="mailto:joserba91@gmail.com">joserba91@gmail.com</a>
-------------------	--

## Resultados Kaggle

Overview **Data** Code Models Discussion Leaderboard Rules Team

Submissions [Submit Predictions](#) ...

3845	Diamantakiou François		0.18806	1	1d
3846	FelipeAraujo		0.18819	4	1mo
3847	Y.Yamakawa		0.18820	1	2mo
3848	Robert Grady Williams		0.18835	13	13d
3849	Vetle OyeOpheim		0.18837	3	1mo
3850	Sebastian Cuya		0.18849	1	2mo
3851	<b>Amilcar Rodriguez</b>		0.18876	11	2h
<div> <b>Your Best Entry!</b> Your most recent submission scored 0.18876, which is an improvement of your previous score of 0.18893. Great job!</div> <div><a href="#">Tweet this</a></div>					
3852	Madeline Ginsberg		0.18882	8	13d
3853	goldbabyerim		0.18891	2	14d
3854	SacredDeer		0.18900	4	2mo

# DESCRIPCIÓN DEL PROBLEMA

---

Debido a los requerimientos actuales de la compañía y la reciente creación del departamento de ciencia de datos, ha surgido la necesidad de construir una herramienta de cálculo para tasar potenciales propiedades de los clientes que puedan ingresar a nuestro sitio web; el objetivo entonces apunta en ayudar a nuestros clientes mejorando la experiencia y uso del sitio web que desemboque en más visitas y más propiedades que estén a cargo de nuestra compañía (inmobiliaria).

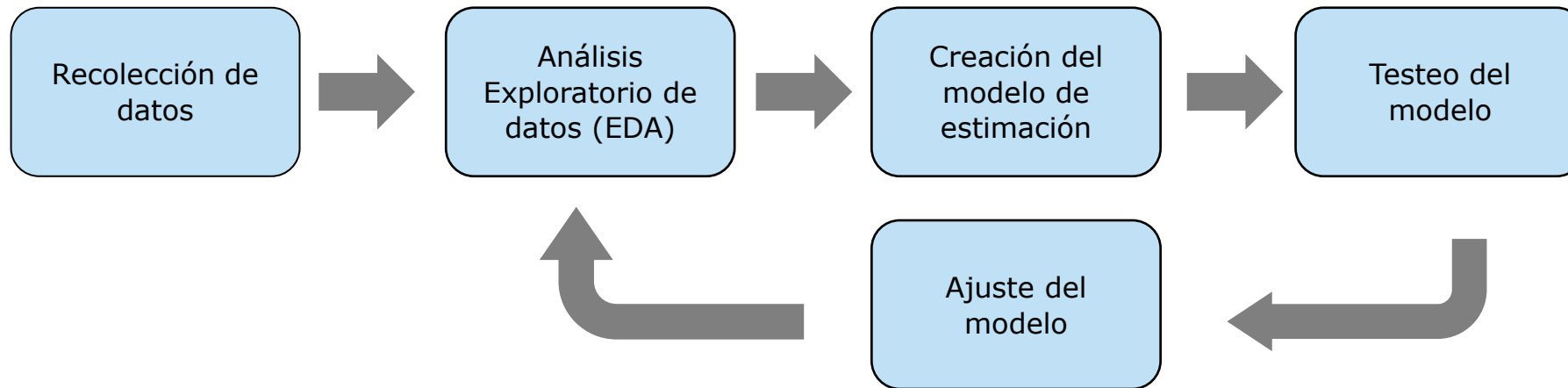
Esta herramienta interactiva debe ser capaz de estimar el costo de una propiedad según diversas características que el cliente pueda suministrar, por tal motivo, como parte de un plan piloto, el departamento de experiencia de cliente y ciencia de datos ha iniciado un análisis para construir un modelo de predicción que a priori permita estimar el costo de una vivienda y así determinar características claves, modelos a utilizar y posibles desempeños que generen una primera visión de solución a esta necesidad.

A continuación, se describirán ciertos elementos de análisis y metodológicos que desembocarán en un pliego de resultados que esclarezcan los próximos pasos a seguir para el asentamiento de este plan piloto.

# METODOLOGÍA

---

A nivel metodológico para este plan piloto se han considerado las siguientes etapas:



Cada etapa será explicada de forma resumida a continuación.

# METODOLOGÍA

A nivel metodológico para este plan piloto se han considerado las siguientes etapas:

Etapa	Procedimiento
Recolección de datos	<ul style="list-style-type: none"><li>• Datos referenciales de Ames, Iowa (Fuente: Kaggle)</li><li>• Sets de datos de Entrenamiento y Pruebas</li><li>• Volumen de datos 79 variables (categóricas y numéricas) x 1460 filas</li></ul>
Análisis Exploratorio de datos (EDA)	<ul style="list-style-type: none"><li>• Análisis de características numéricas y categóricas</li><li>• Revisión de medidas estadísticas dentro de la data</li><li>• Imputación de valores nulos</li><li>• Tratamiento de Valores perdidos (outliers)</li><li>• Análisis de correlación de variables</li><li>• Selección de características Claves</li></ul>

# METODOLOGÍA

A nivel metodológico para este plan piloto se han considerado las siguientes etapas:

Etapa	Procedimiento
<div data-bbox="364 482 693 675">Creación del modelo de estimación</div>	<ul style="list-style-type: none"><li>• Codificación de características (one-hot encoding)</li><li>• Split de data (data de entrenamiento y testeo del modelo)</li><li>• Escalamiento de los datos</li><li>• Modelado (Regresión Lineal Lasso)</li><li>• Obtención de métricas de desempeño del modelo</li></ul>
<div data-bbox="364 803 693 996">Testeo del modelo</div>	<ul style="list-style-type: none"><li>• Ajuste de data de testeo en base a data de entrenamiento</li><li>• Codificación de características (one-hot encoding)</li><li>• Ejecución de predicciones</li><li>• Impresión de resultados y análisis de estos</li></ul>
<div data-bbox="364 1132 693 1325">Ajuste del modelo</div>	<ul style="list-style-type: none"><li>• Ajuste de coeficientes del modelo</li><li>• Selección de más o menos características</li><li>• Validación de métodos para tratamientos de Outliers y N/A</li><li>• Definición de modelo final</li></ul>

# RESULTADOS

---

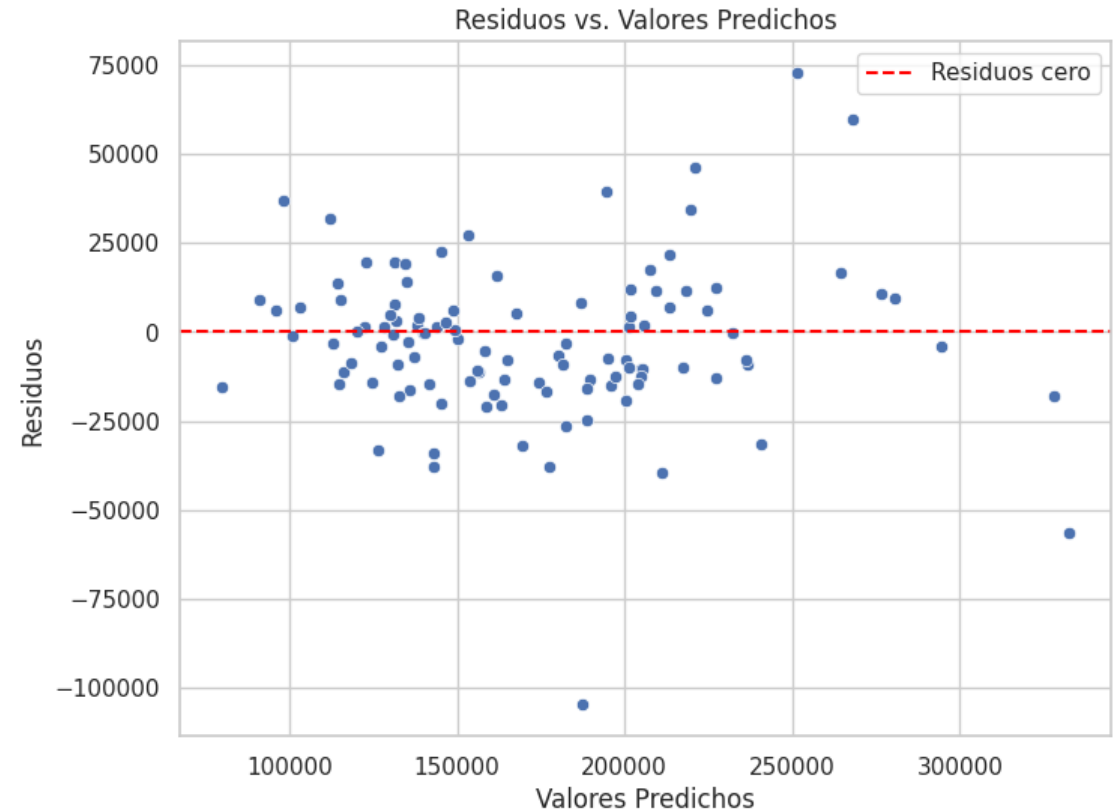
Al tratarse de una predicción cuyo valor objetivo es el precio de venta de una propiedad según ciertas características, concluimos:

- El modelo utilizado ha sido la regresión Lasso debido a la característica del problema, donde se requiere predecir una variable objetivo numérica como el precio de venta (SalePrice) de la propiedad. La ventaja de esta regresión es que permite la penalización de coeficientes de algunas variables a 0 que se traduce en una mejor selección de características, ventajoso frente a más de 70 características iniciales.
- La regresión Lasso es menos sensible al sobre ajuste por lo que la convierte en un buen modelo para estimar este tipo de escenarios.
- Las características numéricas con mayor incidencia en la predicción son: 'YearRemodAdd', 'TotalBsmtSF', 'FullBath', '1stFlrSF', 'TotRmsAbvGrd', 'GrLivArea', 'GarageCars', 'OverallQual', 'MasVnrArea', 'GarageYrBlt', 'YearBuilt', 'GarageArea'.
- Las características categóricas con mayor incidencia en la predicción son: 'MSZoning', 'Street', 'LotShape', 'LandContour', 'LotConfig', 'LandSlope', 'Neighborhood', 'Condition1', 'BldgType', 'HouseStyle', 'RoofStyle', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'Heating', 'HeatingQC', 'CentralAir', 'Electrical', 'KitchenQual', 'Functional', 'FireplaceQu', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageCond', 'PavedDrive', 'SaleType', 'SaleCondition'.

# RESULTADOS

Al tratarse de una predicción cuyo valor objetivo es el precio de venta de una propiedad según ciertas características, concluimos:

- El resultado de las predicciones es prometedor y como podemos ver en el siguiente gráfico, los residuos a 0 (coincidencias) se encuentran distribuidos uniformemente dentro de la nube de estimación, lo que conlleva a pensar que el modelo está representando adecuadamente la variabilidad de los datos.
- Respecto al indicador  $R^2$  concluimos que existe un 83% de la variabilidad representada por el modelo en la predicción del precio de venta.





# RESULTADOS

---

Al tratarse de una predicción cuyo valor objetivo es el precio de venta de una propiedad según ciertas características, concluimos:

- Respecto al indicador RMSE concluimos que las predicciones del modelo tienen un error de alrededor de \$21975 en comparación con los valores reales, que para en términos de un rango de SalePrice de ejemplo (Train) que va desde 34.900\$ a 755.000\$ representa un 5.8%.

Fuente de datos	MSE	R <sup>2</sup>	RMSE
Train (80%)	1.264730e+08	0.957033	11246.021868
Traing (20%)	4.829266e+08	0.836646	21975.589981

- Podemos agregar que se si bien se ha hecho un trabajo de feature engineering, tratando de identificar variables claves, tratamiento de outliers y valores nulos se sugiere una mayor cantidad de iteraciones con el objetivo de obtener mejores resultados frente a escenarios reales, y más iteraciones para reducir (en la medida de lo posible el número de variables).
- También se sugiere la implementación de modelos de ensamblado en bagging para contrastar los resultados previos y en caso de ser necesario llegar a un modelo más robusto.