

Predicting house prices using ML



Ask a home buyer to describe their dream house, and they probably won't begin with the height of the basement ceiling or the proximity to an east-west railroad. But this playground competition's dataset proves that much more influences price negotiations than the number of bedrooms or a white-picket fence.

With 79 explanatory variables describing (almost) every aspect of residential homes in Ames, Iowa, this competition challenges you to predict the final price of each home.

In the contemporary digital era, useful information of the society can be retrieved from a wide variety of sources and stored in the form of structured, unstructured and semi-structured formats. In the analysis of economic phenomena or social observations, advancement of innovative technology makes it possible to systematically extract the relevant information, transform them into complex data formats and structures, and then perform suitable analyses. Because of these new circumstances, traditional data processing and analytical tools may not be able to capture, process and analyse highly complex

information in the social and economic worlds. New techniques have been developed in response to the treatment of the colossal amount of available data.

Machine learning is one of the cutting edge techniques that can be used to identify, interpret, and analyse hugely complicated data structures and patterns . It allows consequential learning and improves model predictions with a systematic input of more recent data Contemporary research on machine learning is a subset of artificial intelligence (AI) that attempts to train computers with new knowledge through input of data, such as texts, images, numerical values, and so on, and support its interaction with other computer networks. Machine learning is about 'the science of getting computers to learn and act like humans do, and improve their learning over time in autonomous fashion, by feeding them data and information in the form of observations and real-world interactions'.

Machine learning can be broadly categorised into three types, namely supervised learning, unsupervised learning and semi-supervised learning. In essence, supervised machine learning algorithms intend to identify a function that is able to give accurate out-of-sample forecasts. For example, in property research, if an investigator intends to make forecast of housing prices y_i from its physical, neighbourhood and accessibility characteristics x_{ij} from a sample of n apartments, one can assume $L(\hat{y}_i, y_i)$ to be the prediction loss function. A machine learning algorithm will look for a function \hat{f} that produces lowest expected prediction loss $E(y_i, x_{ij})[L(\hat{f}(x_{ij}), y_i)]$ on the *test* data from the same distribution.... A model trained through supervised learning is said to be successful if it can make predictions within an acceptable level of accuracy. Examples of supervised learning include linear regression and support vector machine.

- **Creative feature engineering**

Advanced regression techniques like random forest and gradient boosting.

- **Fun with Real Estate Data**

Use Rmarkdown to learn advanced regression techniques like random forests and XGBoost

- **XGBoost with Parameter Tuning**

Implement LASSO regression to avoid multicollinearity

Includes linear regression, random forest, and XGBoost models as well

- **Ensemble Modeling: Stack Model Example**

Use “ensembling” to combine the predictions of several models

Includes GBM (gradient boosting machine), XGBoost.

• A Clear Example of Overfitting

Learn about the dreaded consequences of overfitting data

Comprehensive Data Exploration with Python

Understand how variables are distributed and how they interact

Apply different transformations before training machine learning models

• House Prices EDA

Learn to use visualization techniques to study missing data and distributions

Includes correlation heatmaps, pairplots, and t-SNE to help inform appropriate inputs to a linear model

A Study on Regression Applied to the Ames Dataset

Demonstrate effective tactics for feature engineering

Explore linear regression with different regularization methods including ridge, LASSO, and ElasticNet using scikit-learn

Regularized Linear Models

Build a basic linear model

Try more advanced algorithms including XGBoost

```
import numpy as np
import pandas as pd

%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
color = sns.color_palette()
sns.set_style('darkgrid')
import warnings
def ignore_warn(*args, **kwargs):
```

```

pass
warnings.warn = ignore_warn

from scipy import stats
from scipy.stats import norm, skew

pd.set_option('display.float_format', lambda x: '{:.3f}'.format(x))

from subprocess import check_output
print(check_output(["ls", "../input"]).decode("utf8"))

train = pd.read_csv('/kaggle/input/house-prices-advanced-regression-techniques/train.csv')
test = pd.read_csv('/kaggle/input/house-prices-advanced-regression-techniques/test.csv')
train.head(5)

```

Id	MS Sub Class	MSZ oning	Lot Frontage	Lot Area	Street	Alley	Lot Shape	Land Contour	Utilities	...	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr Sold	Sale Type	Sale Condition	Sale Price	
0	1	60	RL	65.0000	8450	Pave	NaN	Reg	Lvl	AlIPub	...	0	NaN	NaN	NaN	0	2	2008	WD	Normal	208500
1	2	20	RL	80.0000	9600	Pave	NaN	Reg	Lvl	AlIPub	...	0	NaN	NaN	NaN	0	5	2007	WD	Normal	181500

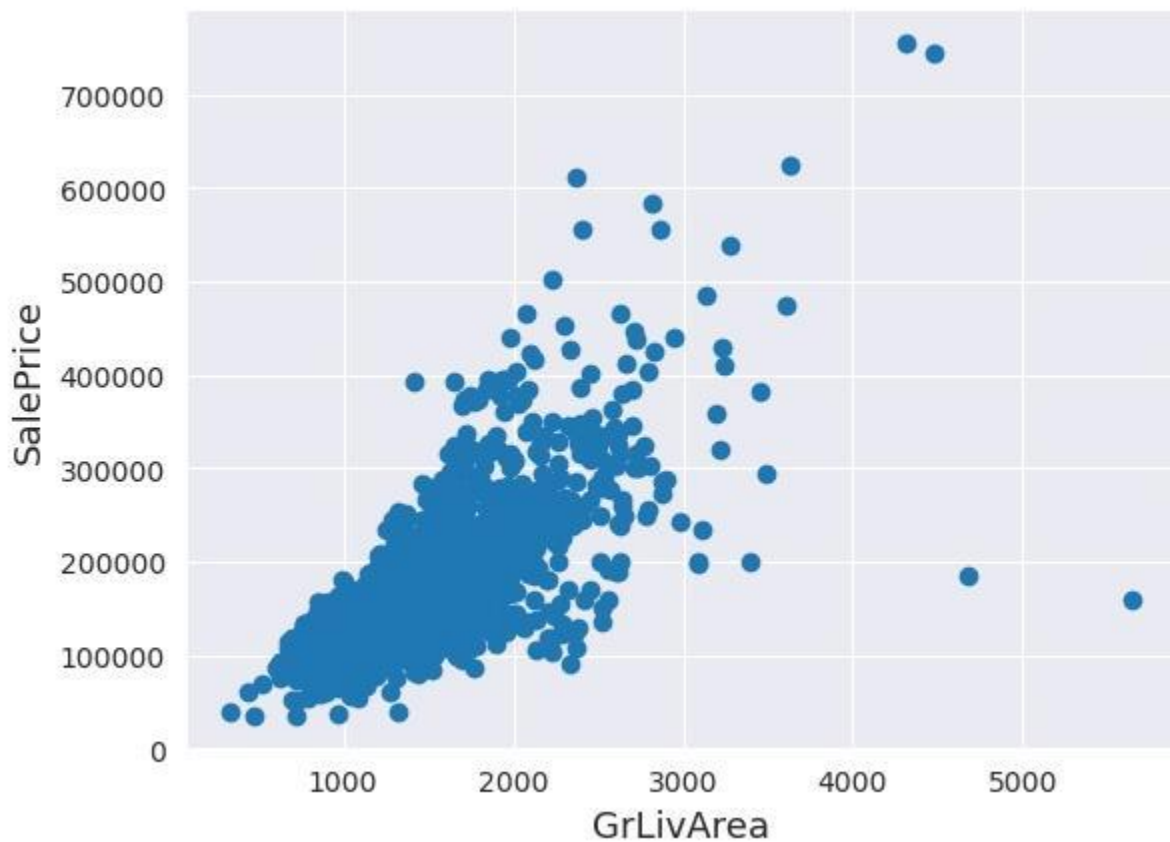
I d	MS Sub Class	M SZ oni ng	Lot Fro nta ge	L ot Ar e a	S t r e e t	A l l e y	Lo tS ha pe	Lan dC ont our	Ut ilit ie s	...	Po ol Ar ea	P o ol Q C	F e n c e	Mis cFe atu re	Mi sc V al	M o S ol d	Y r S ol d	Sa le Ty pe	Sal eCo nditi on	Sa le Pri ce	
2	3	60	RL	6 8. 0 0 0	1 1 2 5 0	P a v e	Na N	IR1	L vl	Al IP u b	...	0	Na N	Na N	Na N	0	9	20 08	WD	No rm al	2 2 3 5 0 0
3	4	70	RL	6 0. 0 0 0	9 5 5 0	P a v e	Na N	IR1	L vl	Al IP u b	...	0	Na N	Na N	Na N	0	2	20 06	WD	Ab no rm l	1 4 0 0 0 0
4	5	60	RL	8 4. 0 0 0	1 4 2 6 0	P a v e	Na N	IR1	L vl	Al IP u b	...	0	Na N	Na N	Na N	0	1 2	20 08	WD	No rm al	2 5 0 0 0 0

```
test.head(5)
```

I d	MS SubCl ass	M SZ oni ng	Lot Fron ta ge	L ot Ar e a	S t r e t	A l l e y	Lo t S ha pe	Lan d C ont our	U t i l i t i e s	...	Scr een Por ch	Po ol Ar ea	P o ol Q C	F e n c e	Mis c Fe atu re	Mi sc V al	M o S old	Y r S old	Sa le Ty pe	Sal e Co nditi on	
0	146 1	20	RH	8 0. 0 0 0	1 1 6 2 2	P a v e	Na N	Re g	L vl	A ll P ub	...	12 0	0	Na N	Mn Prv	Na N	0	6	20 10	WD	N or mal
1	146 2	20	RL	8 1. 0 0 0	1 4 2 6 7	P a v e	Na N	IR1	L vl	A ll P ub	...	0	0	Na N	Na N	G ar 2	1 2 5 0 0	6	20 10	WD	N or mal
2	146 3	60	RL	7 4. 0 0 0	1 3 8 3 0	P a v e	Na N	IR1	L vl	A ll P ub	...	0	0	Na N	Mn Prv	Na N	0	3	20 10	WD	N or mal
3	146 4	60	RL	7 8. 0 0 0	9 9 7 8	P a v e	Na N	IR1	L vl	A ll P ub	...	0	0	Na N	Na N	Na N	0	6	20 10	WD	N or mal
4	146 5	12 0	RL	4 3. 0 0 0	5 0 0 5	P a v e	Na N	IR1	H L S	A ll P ub	...	14 4	0	Na N	Na N	Na N	0	1	20 10	WD	N or mal

Outliers

```
fig, ax = plt.subplots()
ax.scatter(x = train['GrLivArea'], y = train['SalePrice'])
plt.ylabel('SalePrice', fontsize=13)
plt.xlabel('GrLivArea', fontsize=13)
plt.show()
```



```
train["SalePrice"] = np.log1p(train["SalePrice"])
sns.distplot(train['SalePrice'] , fit=norm);
```

```

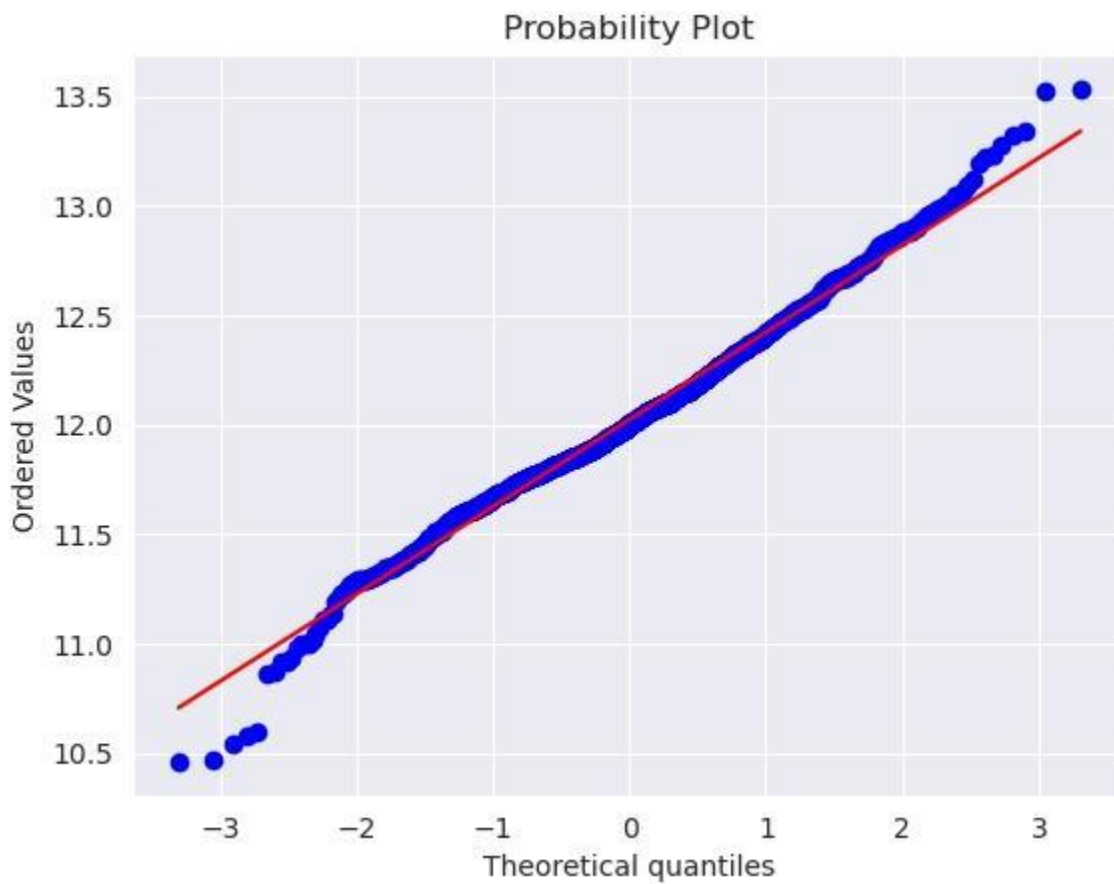
(mu, sigma) = norm.fit(train['SalePrice'])
print( '\n mu = {:.2f} and sigma = {:.2f}\n'.format(mu, sigma))

plt.legend(['Normal dist. ( $\mu$ = $ {:.2f} and  $\sigma$ = $ {:.2f} )'.format(mu, sigma)],
           loc='best')
plt.ylabel('Frequency')
plt.title('SalePrice distribution')

fig = plt.figure()
res = stats.probplot(train['SalePrice'], plot=plt)
plt.show()

```





The skew seems now corrected and the data appears more normally distributed...

CONCLUSION

Thus, the machine learning model using linear regression algorithm is very helpful in predicting the house prices for real estate customers