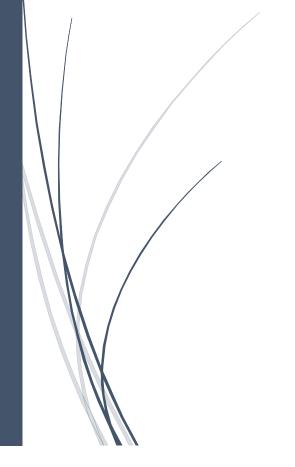
1/5/2025

Ensemble Learning

Ames Housing Dataset





Dr. F Keshavarze

Mohammad Mahdi Shafighi

SHAHED UNIVERSITY

فهرست مطالب

٣	<u>ٔ .</u> توضیح دادهها و پیشپردازشهای انجامشده :	1
٣	'. پیشپردازش دادهها شامل مراحل زیر بود:	<u> </u>
٣	پاکسازی دادهها:	<u>. 1</u>
٣	تحلیل دادههای اکتشافی:(EDA)	<u>. ۲</u>
٣	مهندسی ویژگیها:	<u>.</u> ٣
٤	پیشپر داز ش نهایی دادهها:	٤_
٤	ر. توضیح الگوریتمهای پایه و روشهایEnsemble	<u>٣</u>
٤	مدلهای پایه:	4
	روشهای Ensemble :	
٤	نتایج ارزیابی مدلها	٤
٥	. نتایج ارزیابی مدلها:	0
٥) نتایج	تحليل
٦	لدات برای بهبود بیشتر	پیشنه
٦	تنظیم هایپرپارامترها:	<u>.1</u>
٦	افز ایش و یژگیهای مشتقشده:	<u>.II.</u>
٦	ِ استفاده از Ensemble های پیچیدهتر:	<u>.III</u>
٦	بهبود دادهها:	IV
		<u>.V</u>
	ندى:	جمعين

گزارش تحلیل و نتایج پروژه Ames Housing Data

۱. توضیح داده ها و پیش پر داز شهای انجامشده:

دادههای استفادهشده شامل اطلاعات مرتبط با ویژگیهای خانهها و قیمت نهایی فروش آنها

ستون SalePrice در منطقه Ames Iowa است. ویژگیها به دو دسته کلی تقسیم میشوند:

ویژگیهای عددی مانندLotArea ، GrLivArea، وYearBuilt،

ویژگیهای غیر عددی مانندBldgType ، Neighborhood، ویژگیهای غیر عددی مانند

٢. پیشپردازش داده ها شامل مراحل زیر بود:

١. پاکسازي دادهها:

- حذف مقادیر تکراری
- پر کردن مقادیر گمشده در ستونهای عددی با مقدار میانه و در ستونهای غیر عددی با پرتکرارترین مقدار
 - حذف رکوردهایی که همچنان مقادیر گمشده داشتند.
 - ۲. تحلیل دادههای اکتشافی: (EDA)
 - بررسی توزیع متغیر هدف (SalePrice) با استفاده از هیستوگرام.
 - محاسبه ماتریس همبستگی ویژگیهای عددی و نمایش آن با Heatmap
 - ایجاد نمودار پراکندگی (scatter plot) برای بررسی رابطه GrLivArea و TotalBsmtSF با SalePrice

مهندسی ویژگیها:

• افزودن ویژگی جدید TotalLivingArea (مجموع مساحت زندگی در طبقات اصلی و زیرزمین)

٤. پیشیردازش نهایی دادهها:

- تقسیم دادهها به ویژگیها (X) و متغیر هدف(y)
- دستهبندی ستونها به ویژگیهای عددی و غیر عددی
- استفاده از استاندارد سازی برای ستونهای عددی و تبدیل OneHotEncoding برای ستونهای غیر عددی

۳. توضیح الگوریتمهای پایه و روشهای Ensemble

• مدلهای پایه:

Linear Regression :یک مدل خطی ساده که برای پیش بینی متغیرهای پیوسته استفاده می شود.

Decision Tree :الگوریتمی مبتنی بر درخت که تصمیمات را بر اساس بیشترین کاهش خطا (مانندMSE) اتخاذ می کند.

Random Forest :یک مدل Bagging که مجموعهای از درختهای تصمیم گیری را ایجاد می کند و نتایج را ترکیب می کند.

• روشهای Ensemble

Bagging (Random Forest) الگوریتمی که نمونه گیری Bootstrap انجام داده و درختهای تصمیم متعدد را ترکیب می کند.

Boosting (Gradient Boosting) الگوریتمی که مدلهای متوالی ایجاد کرده و خطای نمونههای قبلی را بهبود میدهد.

Stacking :ترکیب چندین مدل پایه (مانند Random Forest و Gradient Boosting) و آموزش یک مدل متا برای بهبود عملکرد.

٤. نتایج ارزیابی مدلها

در این پروژه، از معیارهای زیر برای ارزیابی عملکرد مدلها استفاده شد:

. Mean Squared Error (MSE) ميانگين مربع خطاها.

Root Mean Squared Error (RMSE): ریشه دوم میانگین مربع خطاها.

(اضافی شده) میزان واریانس توضیح داده شده توسط مدل. (R^2 Score

٥. نتایج ارزیابی مدلها:

مدلهایBoosting ، Bagging، و Stacking با استفاده از دادههای Ames Housing مورد ارزیابی قرار گرفتندکه نتایج از این قرار می باشد :

R ² Score	MSE	RMSE	MAE	مدل
0.9130	697,214,755.73	26,404.82	15,823.80	Bagging
0.9119	706,663,100.46	26,583.14	15,249.18	Boosting
0.9075	741,483,154.40	27,230.19	17,135.91	Stacking

تحليل نتايج

Bagging (Random Forest) .\

- RMSE. و MAE و معیارهای \circ
- این مدل توانایی بسیار خوبی در کاهش خطاهای پیشبینی دارد، به دلیل استفاده از چندین
 درخت تصمیم و ترکیب نتایج آنها.

Boosting (Gradient Boosting) . 7

- o عملکرد کمی ضعیفتر از Bagging در معیارهای RMSE و RSE و RMSE
- مدل Boosting به دلیل افزایش تدریجی و بهبود خطاهای مدلهای قبلی، مناسب
 پیشبینیهای پیچیده تر است.

Stacking . "

o عملکرد نسبتاً ضعیفتر در مقایسه با Bagging و

این مدل به دلیل پیچیدگی بیشتر و وابستگی به کیفیت مدلهای پایه، ممکن است در
 دادههای خاص کارایی پایین تری داشته باشد.

بیشنهادات برای بهبود بیشتر

تنظیم هاییرپارامترها:

استفاده از GridSearchCV یا RandomizedSearchCV برای یافتن بهترین مقادیر
 هایپرپارامترهای مدلها.

برای مثال:

- تعداد درختها (n_estimators) و عمق درختها (max_depth) در • Random Forest.
- نرخ یادگیری (learning_rate) و تعداد مراحل تقویت در learning_rate)

اا. افزایش ویژگیهای مشتقشده:

شناسایی و افزودن ویژگیهای جدید مانند تعامل بین متغیرها یا استفاده از تحلیلهای دامنه
 تخصصی (مانند نسبت مساحت زیرزمین به کل مساحت خانه).

ااا. استفاده از Ensemble های پیچیدهتر:

o ترکیب مدلهای دیگری مانند XGBoost یا LightGBM با مدلهای فعلی.

۱۷. بهبود دادهها:

o افزایش تعداد دادهها یا استفاده از تکنیکهایی مانند Data Augmentation برای کاهش خطاها.

۷. بررسی حذف ویژگیهای کماهمیت:

• ویژگیهایی که همبستگی کمتری با متغیر هدف دارند، میتوانند حذف شوند تا پیچیدگی مدل کاهش یابد.

جمعبندی:

مدل Bagging با MAE کمترین و Score بالاترین، بهترین عملکرد را در این تحلیل ارائه داد. مدلهای Boosting با Stacking بیچیدگی بیشتر، عملکرد قابل قبولی داشتند. برای بهبود بیشتر، تنظیم دقیق هایپرپارامترها و استخراج ویژگیهای جدید پیشنهاد می شود.