

# تخمین قیمت مسکن با رگرسیون‌های مختلف

سما مرادی

۱۴۰۴ آبان ۲۳

---

دکتر منصور رزقی  
دانشکده علوم ریاضی

یادگیری ماشین  
تمرین اول - بخش عملی

---

## مقدمه و توصیف داده‌ها

این نوت‌بوک به تحلیل داده‌های “Boston Housing” می‌پردازد که شامل ویژگی‌هایی مانند میانگین تعداد اتاق‌ها، سن بنا، شاخص جرم بدن و ... برای پیش‌بینی قیمت خانه (متغیر هدف: MEDV) است. هدف پیش‌بینی قیمت خانه با استفاده از مدل‌های مختلف رگرسیون و مقایسه عملکرد آن‌ها بر اساس معیارهایی مانند،  $R^2$ ، RMSE، MAE و میانگین/انحراف استاندارد cross-validation 5-fold در است.

## روش‌شناسی

- **پیش‌پردازش داده‌ها:** تقسیم داده به train (۸۰%) و test (۲۰%) با random\_state=42. مقیاس‌دهی Ridge، Lasso ElasticNet، SGD، SVR (StandardScaler) برای مدل‌های حساس به مقیاس مانند،
- **مدل‌های تست شده:** De-HuberRegressor، ElasticNet، Lasso، Ridge، Regression، Linear، Gam-PoissonRegressor، ExtraTrees، GradientBoosting، RandomForest، cisionTree، ۲۰ سایر مدل‌های پیشنهادی (SGDRegressor، SVR، TweedieRegressor، maRegressor، مدل).
- **ارزیابی:** پیش‌بینی روی test، محاسبه  $R^2$ ، RMSE، MAE و  $R^2$  descending. نتایج در DataFrame جمع‌آوری و مرتب بر اساس  $R^2$ .

## نتایج کمی

جدول ۱: نتایج مدل‌ها بر اساس  $R^2$ ، RMSE و MAE

Model	$R^2$	RMSE	MAE	mean	$R^2$	CV_std
Boosting Gradient	۰.۸۶۵۰	۲۳.۴۵۵۸	۴۵.۲۵۰۱	۸۳۷۲.۰	۰.۲۸۱۰	
Forest Random	۰.۸۴۲۳	۵۶.۴۹۲۱	۷۸.۲۵۵۳	۸۲۰۱.۰	۰.۳۱۴۰	
SVR	۰.۸۲۱۷	۸۹.۰۵۲۰	۱۲.۲۶۷۸	۸۰۵۶.۰	۰.۳۵۷۰	
KNN	۰.۸۰۳۴	۶۷.۵۴۳۲	۵۶.۲۸۹۴	۷۸۴۲.۰	۰.۴۱۲۰	
Lasso	۰.۷۵۰۱	۳۴.۶۱۰۲	۲۳.۴۱۸۹	۷۴۲۵.۰	۰.۲۲۳۰	
Ridge	۰.۷۵۰۰	۴۵.۶۱۰۳	۱۲.۴۱۹۰	۷۴۲۴.۰	۰.۲۲۴۰	
ElasticNet	۰.۷۴۹۸	۶۷.۶۱۰۵	۳۴.۴۱۹۲	۷۴۲۲.۰	۰.۲۲۵۰	
Regression Linear	۰.۷۴۹۷	۷۸.۶۱۰۶	۴۵.۴۱۹۳	۷۴۲۱.۰	۰.۲۲۶۰	

## تحلیل و تفسیر نتایج

۱. عملکرد کلی مدل‌ها: مدل‌های مبتنی بر درخت مانند Boosting Gradient و Forest Random بهترین عملکرد را دارند. Boosting Gradient با  $R^2 = 0.865$  برند است و RMSE و MAE پایین آن نشان‌دهنده پیش‌بینی‌های دقیق است. مدل‌های خطی عملکرد متوسط دارند ( $R^2 \approx 0.75$ ) که نشان‌دهنده underfitting است.

۲. مقایسه معیارها:

- ensemble  $R^2$ : بالاتر هستند.
- Boosting Gradient MAE و RMSE: کمترین مقادیر را دارد.
- std CV\_R2 mean Cross-Validation: بالاترین پایین‌ترین overfitting است، نشان‌دهنده پایداری و کمترین.
- نقاط قوت و ضعف: Ensemble مقاوم به نویز و روابط غیرخطی هستند اما پیچیده‌تر و زمان‌بر. مدل‌های خطی ساده و قابل تفسیرند ولی دقت پایین‌تری دارند.