# (Regression Models) مقارنة نماذج الانحدار

# جدول المقارنة الرئيسي

الخاصية	Linear Regression	Ridge	Lasso	Elastic Net	Random Forest	XGBoost
التعقيد الحسابي	O(nd²)	O(nd²)	O(nd²)	O(nd²)	O(n log n × trees)	O(n log n × rounds)
قابلية التفسير	عالية جداً	عالية	عالية	عالية	متوسطة	متوسطة
الأداء مع البيانات الكبيرة	ممتاز	ممتاز	ممتاز	ممتاز	جيد جداً	ممتاز
مقاومة Overfitting	ضعيفة	جيدة جداً	جيدة جداً	ممتازة	ممتازة	جيدة جداً
التعامل مع البيانات المفقودة	يحتاج معالجة	يحتاج معالجة	يحتاج معالجة	يحتاج معالجة	جيد	جيد
حساسية للمعايرة	منخفضة	متوسطة	متوسطة	متوسطة	منخفضة	متوسطة
Feature Selection	И	И	نعم (تلقائي)	نعم (متوازن)	نعم (importance)	نعم (importance)
التعامل مع Multicollinearity	ضعيف جداً	ممتاز	جيد	ممتاز	جيد	جيد
مناسب للبيانات عالية الأبعاد	ضعیف	جيد	ممتاز	ممتاز	ممتاز	ممتاز
سرعة التدريب	سريع جداً	سريع جدآ	سريع جدآ	سريع جدآ	متوسط	سريع

# التفاصيل الفنية

# 1. \ Linear Regression

- المبدأ: يفترض علاقة خطية بين المتغيرات والهدف •
- المعادلة:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_n x_n$
- الافتراضات: خطية، عدم ارتباط الأخطاء، توزيع طبيعي للأخطاء •
- المخرجات: معاملات واضحة المعنى •

# 2. Ridge Regression (L2 Regularization)

- المبدأ Linear regression + عاملة
- المعادلة:  $\left( Loss + \alpha \times \Sigma(\beta_i^2) \right)$
- القوة: يقلل المعاملات الكبيرة، يحل مشكلة
- :المعاملات المهمة
  - قوة التنظيم (1.0-0.1) (alpha

# 3. P Lasso Regression (L1 Regularization)

- المبدأ عاملة + Linear regression :المبدأ
- المعادلة:  $(Loss + \alpha \times \Sigma | \beta_i |)$
- (تلقائي feature selection) **القوة**: يصفر المعاملات غير المهمة •
- :المعاملات المهمة
  - (alpha): (10-0.01) قوة التنظيم

#### 4. 6 Elastic Net

- e Lasso و Ridge المبدأ: يجمع بين
- المعادلة:  $\left( Loss + \alpha_1 \times \Sigma |\beta_i| + \alpha_2 \times \Sigma (\beta_i^2) \right)$
- والاستقرار feature selection **القوة**: يوازن بين •

#### :المعاملات المهمة

- قوة التنظيم الإجمالية :(alpha
- L2 (0-1) إلى L1 نسبة :(11\_ratio) •

# 5. Random Forest Regressor

- المبدأ: مجموعة من أشجار الانحدار + متوسط النتائج •
- overfitting القوة: يتعامل مع العلاقات غير الخطية، مقاوم للـ
- :المعاملات المهمة
  - عدد الأشجار (100-100): (n\_estimators)•
  - عدد المتغيرات لكل انقسام :(max\_features •

# 6. 🖋 XGBoost Regressor

- المبدأ : Gradient boosting المبدأ
- القوة: دقة عالية جداً، يتعامل مع أي نوع علاقة •
- :المعاملات المهمة
  - (n\_estimators): عدد الـ rounds (100-1000)
  - (learning\_rate): (0.3-0.01) معدل التعلم
  - عمق الأشجار (10-3) (max\_depth)

# 7. / LightGBM Regressor

- المبدأ: Gradient boosting المبدأ
- نفس الدقة تقريباً ،XGBoost **القوة**: أسرع من ●
- مناسب لـ: البيانات الضخمة والإنتاج •
- :المعاملات المهمة
  - عدد الأوراق (31-300) (num\_leaves •
  - (learning\_rate): (0.3-0.01) معدل التعلم

#### 8. **@** CatBoost Regressor

- otategorical features مع معالجة متقدمة للـ Gradient boosting المبدأ
- hyperparameter tuning القوة: أقل حاجة لـ •
- :المعاملات المهمة
  - (iterations): عدد الـ rounds (100-1000)
  - عمق الأشجار (10-4): (depth)

# 9. X Support Vector Regression (SVR)

- margin يحتوي أكبر عدد من النقاط ضمن hyperplane المبدأ: يجد
- **القوة**: فعال مع البيانات عالية الأبعاد •
- :المعاملات المهمة
  - قوة التنظيم :C
  - (kernel): نوع الـ kernel ('linear', 'rbf')
  - (epsilon): عرض الـ margin

#### 10. A Decision Tree Regressor

- المبدأ: شجرة من القواعد البسيطة للتنبؤ بالقيم
- القوة: سهل الفهم والتصور
- نقاط الضعف overfitting:نقاط الضعف

# 11. **\*\*\*** K-Nearest Neighbors Regressor

- ◄يران K المبدأ: يتنبأ بناءً على متوسط أقرب
- القوة: بسيط، يتعامل مع الأنماط المحلية •
- نقاط الضعف: بطيء مع البيانات الكبيرة •

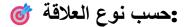
## 12. Neural Networks (MLPRegressor)

المبدأ: شبكة عصبية متعددة الطبقات

**القوة**: يتعلم أي علاقة معقدة •

يحتاج: بيانات كثيرة ووقت تدريب طويل •

# متی نستخدم کل نموذج؟



#### :علاقة خطية بسيطة

- Linear Regression (baseline)
- Ridge (مع multicollinearity)
- Lasso (مع متغيرات كثيرة)

#### علاقة غير خطية:

- XGBoost/LightGBM (الأفضل)
- Random Forest
- Neural Networks

#### بيانات فئوية كثيرة:

- CatBoost (الأمثل)
- XGBoost/LightGBM

## وحسب حجم البيانات 📊

حجم البيانات	الأفضل	البديل
< 1K	Linear/Ridge/Lasso	Decision Tree
1K-10K	Random Forest	XGBoost
10K-100K	XGBoost	LightGBM

حجم البيانات	الأفضل	البديل
100K-1M	LightGBM	XGBoost
> 1M	LightGBM	Linear Regression

# حسب الهدف 🏫

# :فهم العلاقات والتأثيرات

- 1. Linear Regression
- 2. Ridge/Lasso
- 3. Decision Tree

# :أقصى دقة ممكنة

- 1. XGBoost
- 2. **LightGBM**
- 3. CatBoost

# :سرعة في الإنتاج

- 1. **LightGBM**
- 2. Linear Regression
- 3. **XGBoost**

### :مع بيانات معقدة

- 1. Gradient Boosting Family
- 2. Neural Networks
- 3. Random Forest

# نصائح للاختيار العملي

# واستراتيجية البداية السريعة

- 1. **Linear Regression** (baseline سريع)
- (أداء قوي) **LightGBM** أو 2. XGBoost
- Neural Networks → إذا النتائج مش كافية

# وللمشاريع الحقيقية

- تطویر سریع: LightGBM
- إنتاج مستقر: XGBoost
- تفسیر مطلوب: Ridge/Lasso
- **Ensemble :دقة قصوى**

#### تحذيرات مهمة 🛕

- Linear models: تفترض علاقة خطية
- Tree models: حساسة لل outliers
- KNN: curse of dimensionality
- Neural Networks: يحتاج بيانات كثيرة جداً

# الاختيار الذهبي لمعظم المشاكل 👱

- (فهم البيانات) Ridge Regression ابدأ بـ 1
- (أداء قوي) **LightGBM** جرب .2
- 3. مع hyperparameter tuning مع hyperparameter tuning

# مقاييس التقييم المناسبة

النموذج	المقياس الأفضل	السبب
Linear Models	$R^2 + RMSE$	سهل التفسير

النموذج	المقياس الأفضل	السبب	
Tree Models	MAE + RMSE	outliers مقاوم للـ	
<b>Boosting Models</b>	Custom metrics	مرونة في التحسين	
Neural Networks MSE + validation curves		overfitting مراقبة	
4		<b>▶</b>	

# M Recordsتطبيق عملي لـ 2

# للبيانات بحجم 2 مليون:

الخيار الأول (سرعة + دقة) 1. LightGBM:

2. **Linear Models**: السريع baseline

إذا الوقت مش مشكلة :3. **XGBoost** 

4. **X KNN/SVR**: مش عملي خالص

### :مقترح Pipeline

ightarrow EDA ightarrow StandardScaler ightarrow LightGBM ightarrow تقییم النتائج