(Classification Models) مقارنة نماذج التصنيف

جدول المقارنة الرئيسي

الخاصية	Logistic Regression	Random Forest	XGBoost	LightGBM	CatBoost	SVM	Decision Tree
التعقيد الحسابي	O(nd)	O(n log n × trees)	O(n log n × rounds)	O(n log n × rounds)	O(n log n × rounds)	O(n²d)	O(n log n)
قابلية التفسير	عالية جداً	متوسطة	متوسطة	متوسطة	متوسطة	منخفضة	عالية جداً
الأداء مع البيانات الكبيرة	ممتاز	جيد جدآ	ممتاز	ممتاز جداً	ممتاز	ضعیف	جيد
مقاومة Overfitting	جيدة	ممتازة	جيدة جداً	جيدة جداً	ممتازة	جيدة	ضعيفة جداً
التعامل مع البيانات المفقودة	يحتاج معالجة	جيد	جيد	جيد	ممتاز	يحتاج معالجة	ممتاز
حساسية للمعايرة	متوسطة	منخفضة	متوسطة	منخفضة	منخفضة جداً	عالية جداً	منخفضة
التعامل مع الفئات غير المتوازنة	ضعیف	جيد	جيد جدآ	جيد جداً	جيد جدآ	متوسط	ضعیف
مناسب للبيانات عالية الأبعاد	جيد	ممتاز	ممتاز	ممتاز	ممتاز	ممتاز	ضعیف
سرعة التدريب	سريع جداً	متوسط	سريع	سريع جداً	متوسط	بطيء	سريع جداً

التفاصيل الفنية

1. **Logistic Regression**

- للتنبؤ بالاحتمالات sigmoid function المبدأ: يستخدم
- الافتراضات: العلاقة الخطية بين المتغيرات •
- **المخرجات**: احتمالات + قرار ثنائي •
- :المعاملات المهمة
 - قوة التنظيم (أصغر = تنظيم أكبر): C
 - (solver): طريقة الحل ('liblinear', 'lbfgs')

2. Random Forest

- المبدأ: مجموعة من أشجار القرار + تصويت جماعي •
- ويتعامل مع البيانات المعقدة Overfitting القوة: يقلل •
- :المعاملات المهمة
 - (n_estimators): (1000-100) عدد الأشجار
 - (max_depth): عمق الشجرة
 - (min_samples_split): أقل عينات للتقسيم

2. Random Forest

- المبدأ: مجموعة من أشجار القرار + تصويت جماعي •
- ويتعامل مع البيانات المعقدة Overfitting القوة: يقلل •
- :المعاملات المهمة
 - عدد الأشجار (100-100) (n_estimators
 - عمق الشجرة :(max_depth)
 - أقل عينات للتقسيم :(min_samples_split •

3. XGBoost (Extreme Gradient Boosting)

- مُحسن + تقنيات متقدمة Gradient Boosting :المبدأ •
- overfitting القوة: أداء ممتاز، سرعة عالية، مقاوم لل

:المعاملات المهمة

- (n_estimators): عدد الـ boosting rounds
- (learning_rate): (0.3-0.01) معدل التعلم
- عمق الأشجار (10-3) (max_depth):
- (subsample): نسبة العينات لكل round

4. / LightGBM

- المبدأ: Gradient Boosting المبدأ:
- يستهلك ذاكرة أقل ،XGBoost **القوة**: أسرع من •

:المعاملات المهمة

- عدد الأوراق (100-31) (num_leaves
- (learning_rate) معدل التعلم
- (feature_fraction): نسبة المتغيرات المستخدمة
- (bagging_fraction): نسبة العينات

5. @ CatBoost

- المبدأ: Gradient Boosting المبدأ: صُحسن للمتغيرات الفئوية
- تلقائياً categorical features **القوة**: يتعامل مع الـ •

:المعاملات المهمة

- (iterations): عدد ال boosting rounds
- (learning_rate) معدل التعلم
- عمق الأشجار :(depth

• (cat_features): المتغيرات الفئوية

6. X Support Vector Machine (SVM)

- **المبدأ**: يجد الحدود المثلى بين الفئات •
- **القوة**: ممتاز مع البيانات عالية الأبعاد والعلاقات المعقدة ●
- :المعاملات المهمة
 - قوة التنظيم:(C) •
 - (kernel): نوع الـ kernel ('linear', 'rbf', 'poly')
 - (gamma): تأثير كل نقطة تدريب

7. A Decision Tree

- (if-else) المبدأ: شجرة من القواعد البسيطة
- **القوة**: سهل الفهم والتفسير •
- بسهولة Overfitting :**نقاط الضعف** •
- :المعاملات المهمة
 - (max_depth): عمق الشجرة
 - أقل عينات في الورقة :(min_samples_leaf) •
 - (ˈginiˈ, ˈentropyˈ) معيار التقسيم :

8. ****** K-Nearest Neighbors (KNN)**

- جيران K المبدأ: يصنف بناءً على أقرب •
- **القوة**: بسيط وفعال مع البيانات المحلية •
- :المعاملات المهمة
 - عدد الجيران (5-15 عادة) : n_neighbors
 - (ˈuniformˈ, ˈdistanceˈ) طريقة الوزن :
 - (metric): مقياس المسافة ('euclidean', 'manhattan')

9. I Naive Bayes

- المبدأ: يفترض استقلالية المتغيرات + قانون بايز •
- القوة: سريع جداً، جيد مع النصوص
- :الأنواع
 - (GaussianNB): للبيانات المستمرة
 - (MultinomialNB): للبيانات المنفصلة
 - (BernoulliNB): للبيانات الثنائية

10. Neural Networks

- المبدأ: شبكة من العقد المترابطة •
- القوة: يتعلم العلاقات المعقدة جداً
- :المعاملات المهمة
 - (hidden_layer_sizes): حجم الطبقات المخفية
 - (relu', 'tanh') دالة التفعيل (relu', 'tanh')
 - (learning_rate): معدل التعلم

الفئة الجديدة: Gradient Boosting Models 💉

وأقوى النماذج للمسابقات والإنتاج 👱

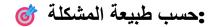
- الأشهر والأكثر استقراراً :1. XGBoost
- 2. **LightGBM**: الأسرع والأقل استهلاكاً للذاكرة
- الأفضل مع البيانات الفئوية :3. CatBoost

💪 مميزات Gradient Boosting:

- دقة عالية جداً (يفوز في أغلب المسابقات) •
- يتعامل مع أي نوع بيانات

- نسبياً) overfitting مقاوم لل
- يعطي feature importance

متى نستخدم كل نموذج؟



:للمشاكل البسيطة والسريعة

• Logistic Regression, Naive Bayes

:للدقة العالية

• Random Forest, SVM, Neural Networks

:للتفسير والفهم

• Decision Tree, Logistic Regression

:للبيانات الكبيرة

• Logistic Regression, Naive Bayes, Random Forest

وحسب نوع البيانات

نوع البيانات	الأفضل	البديل
نصوص	Naive Bayes	Logistic Regression
صور	Neural Networks	SVM
بيانات جدولية	Random Forest	Logistic Regression
بيانات زمنية	Neural Networks	Random Forest
بيانات فئوية	Decision Tree	Naive Bayes

:حسب الأولوية 🔸

السرعة أهم شيء:

- 1. Naive Bayes
- 2. Logistic Regression
- 3. KNN

:الدقة أهم شيء

- 1. Random Forest
- 2. Neural Networks
- 3. SVM

:التفسير أهم شيء

- 1. Decision Tree
- 2. Logistic Regression
- 3. Naive Bayes

نصائح للاختيار العملي

استراتيجية البداية السريعة 💉

- 1. بادأ بـ Logistic Regression (baseline سريع)
- (أداء جيد عادة) Random Forest جرب .2
- 3. أو **Neural Networks** إذا النتايج مش كافية، جرب

وقواعد ذهبية 💡

- Random Forest مشكلة جديدة؟ ابدأ بـ
- Decision Tree أو KNN بيانات أقل من **1000؟**
- Decision Tree أو Decision Tree
- بیانات نصیة؟ Naive Bayes
- Neural Networks مشكلة معقدة جداً؟

وتحذيرات مهمة 🔔

• KNN: بطيء جداً مع البيانات الكبيرة

• **Decision Tree**: يحفظ البيانات بسهولة (overfitting)

• SVM: يحتاج معايرة دقيقة جداً

• Neural Networks: يحتاج بيانات كتير ووقت تدريب طويل