

DRC®

# دليل لأفضل الممارسات في تحليل البيانات الاستكشافي (EDA)



# المحتويات

- 3 المقدمة
- 4 ما هو تحليل البيانات الاستكشافي؟
- 5 لماذا يعد تحليل البيانات الاستكشافي ضروريًا لمشاريع الذكاء الاصطناعي والتعلم الآلي؟
- 6 قيمة تحليل البيانات الاستكشافي للأعمال
- 7 تبعات ومخاطر إهمال تحليل البيانات الاستكشافي
- 8 أفضل الممارسات لتحليل البيانات الاستكشافي الفعال
- 15 الأدوات والأطر لتحليل البيانات الاستكشافي
- 17 دراسة حالة
- 20 الاستنتاجات النهائية

## المقدمة

في عالمنا اليوم المليء بالبيانات، تتمتع الشركات بإمكانية الوصول إلى كميات هائلة من المعلومات. ولكن الفرق بين الشركات الناجحة وتلك التي تواجه صعوبات يكمن في قدرتها على تحويل هذه البيانات إلى رؤى عملية. تحليل البيانات الاستكشافي (EDA) له دور رئيسي في تحقيق هذا التحول.



# ما هو تحليل البيانات الاستكشافي (EDA)؟

تحليل البيانات الاستكشافي (EDA) هو الخطوة الأساسية الأولى في أي مشروع يعتمد على الذكاء الاصطناعي (AI)، التعلم الآلي (ML)، وعلم البيانات. يتضمن هذا التحليل استكشافاً دقيقاً وتصوراً شاملاً وتلخيصاً للخصائص الرئيسية لمجموعة البيانات. الهدف الرئيسي من تحليل البيانات الاستكشافي هو الحصول على فهم متعمق للبيانات والكشف عن الأنماط، الاتجاهات، والعلاقات التي قد تكون غير واضحة في البداية. يوجه هذا الفهم كل خطوة في عملية التعلم الآلي، بدءاً من معالجة البيانات وهندسة الميزات وصولاً إلى بناء النماذج وتحليل النتائج. يعتبر تحليل البيانات الاستكشافي أساسياً لضمان سلامة وموثوقية القرارات المستندة إلى البيانات. من خلال تحديد مشكلات جودة البيانات مثل القيم المفقودة والقيم المتطرفة والتناقضات، يضمن تحليل البيانات الاستكشافي أن تكون البيانات نقية و مناسبة لإنشاء النماذج، مما يؤدي في النهاية إلى نتائج ذكاء اصطناعي وتعلم آلي أكثر دقة وموثوقية.



## لماذا يعد تحليل البيانات الاستكشافي (EDA) ضروريًا لمشاريع الذكاء الاصطناعي والتعلم الآلي؟

تحليل البيانات الاستكشافي (EDA) هو خطوة أساسية في مشاريع الذكاء الاصطناعي (AI) والتعلم الآلي (ML) لضمان نجاح وموثوقية النماذج. تعتمد جودة نماذج التعلم الآلي على جودة البيانات التي تُدرَّب عليها. يساعد تحليل البيانات الاستكشافي في ضمان جودة البيانات وتحديد التحيزات المحتملة التي قد تؤثر على نتائج نماذج الذكاء الاصطناعي والتعلم الآلي.

يساعد تحليل البيانات الاستكشافي في تحديد القيم المتطرفة، والقيم المفقودة، والتناقضات في البيانات، والتي يمكن أن تؤثر سلبًا على أداء النموذج. تنقية هذه المشكلات يضمن أن البيانات المستخدمة للتدريب دقيقة وموثوقة. بالإضافة إلى ذلك، يوجه تحليل البيانات الاستكشافي خطوات التحضير الضرورية مثل التطبيع والتجسيم والتحويل والتي تعتبر حاسمة لضمان أن البيانات في شكل مناسب لتدريب النموذج. يعد اكتشاف التحيزات جانبًا أساسيًا آخر من تحليل البيانات الاستكشافي، حيث يساعد في تحديد التحيزات التي قد تؤدي إلى نتائج مشوهة في النماذج. من خلال فهم توزيع البيانات عبر المجموعات المختلفة، يمكن تحديد ومعالجة مصادر التحيز المحتملة، مما يضمن نماذج ذكاء اصطناعي وتعلم آلي عادلة

وأخلاقية. يعزز تحليل البيانات الاستكشافي أداء النماذج من خلال تقديم رؤى حول الميزات الأكثر أهمية وفهم العلاقات والتفاعلات بين المتغيرات، مما يسهل هندسة الميزات واختيارها. يساعد هذا الإجراء في تقليل الإفراط في التخصيص عن طريق تحديد وإزالة الميزات الزائدة أو غير الضرورية، والتي تؤدي إلى نماذج أكثر عمومية. إضافة إلى ذلك، يساهم تحليل البيانات الاستكشافي في تفسير وشفافية نماذج الذكاء الاصطناعي والتعلم الآلي من خلال توفير فهم واضح للبيانات المستخدمة، وهو أمر ضروري لشرح قرارات النموذج لأصحاب المصلحة وضمان المسؤولية. البيانات النقية والمفهومة بوضوح هي الأساس لبناء نماذج ذكاء اصطناعي وتعلم آلي موثوقة ويمكن الاعتماد عليها.

# قيمة تحليل البيانات الاستكشافي (EDA) للأعمال

يُمكن تحليل البيانات الاستكشافي (EDA) الشركات من اتخاذ قرارات مبنية على البيانات. من خلال كشف الأنماط والاتجاهات المخفية والمعرفة الصحيحة في المجال، يمكن للشركات اتخاذ قرارات استراتيجية مدعومة بالأدلة وليست مبنية على الحدس. يمكن للشركات تحقيق فهم أعمق لجمهورها المستهدف من خلال تحليل بيانات استكشافي يركز على العملاء. من خلال إعطاء الأولوية لتحليل البيانات الاستكشافي، يمكن للشركات اكتساب ميزة تنافسية وتحسين العمليات وتعزيز فهمها للسوق والعملاء.

## قرارات مستنيرة:

يحول تحليل البيانات الاستكشافي (EDA) البيانات الخام إلى رؤى واضحة ومفهومة، مما يتيح للشركات اتخاذ قرارات قائمة على الأدلة التجريبية. ويقلل ذلك من الاعتماد على التخمين والحدس.

## تحديد الفرص:

من خلال تحليل البيانات الاستكشافي (EDA)، يمكن للشركات تحديد فرص جديدة في السوق والتوجهات الناشئة ومجالات للابتكار. يتيح هذا النهج الاستباقي للشركات البقاء في مقدمة المنافسين والاستفادة من تحركات السوق.

## التحليلات التنبؤية:

يمكن للشركات الاستفادة من تحليل البيانات الاستكشافي (EDA) لبناء نماذج تنبؤية تتوقع الاتجاهات والنتائج المستقبلية. يوفر هذا النهج رؤية استباقية لإدارة المخاطر والتخطيط الاستراتيجي المستنير.

## رؤى متمحورة حول العميل:

يتيح تحليل البيانات الاستكشافي (EDA) للشركات فهمًا أعمق لجمهورها المستهدف. من خلال تحليل بيانات العملاء، يمكن للشركات اكتشاف التفضيلات والسلوكيات ونقاط الضعف مما يؤدي إلى استراتيجيات تسويقية أكثر تخصيصًا وفعالية.

## مراقبة الأداء:

يتيح تحليل البيانات الاستكشافي (EDA) للشركات مراقبة وتقييم الأداء التشغيلي بشكل مستمر، مما يساعد في تحديد أوجه القصور ومجالات التحسين.

من خلال الاستفادة من تحليل البيانات الاستكشافي (EDA)، يمكن للشركات عبر مختلف الصناعات تحويل البيانات الخام إلى رؤى قابلة للتنفيذ، مما يؤدي إلى نتائج أعمال أفضل وميزة تنافسية مستدامة. إن الاستثمار في تحليل البيانات الاستكشافي لا يعزز العمليات الحالية فحسب، بل يضع الشركات في موقع أفضل لتحقيق النجاح في المستقبل في عالم يركز بشكل متزايد على البيانات.

## 1. تدهور النموذج:

النماذج التي تعتمد على التعلم الآلي والمدرّبة على بيانات لم يتم استكشافها بشكل كافٍ أو غير مفهومة جيداً تكون عرضة للأداء الضعيف، مما يؤدي إلى تنبؤات واتخاذ قرارات غير موثوقة.

## 2. رؤى غير دقيقة:

بدون استكشاف شامل، قد تكون القرارات مبنية على معلومات غير كاملة أو متحيزة، مما يؤدي إلى قرارات غير صائبة.

## 3. فرص ضائعة:

قد يؤدي إهمال تحليل البيانات الاستكشافي (EDA) إلى عدم استغلال القيمة الموجودة في البيانات، مما يؤدي إلى فقدان فرص الابتكار أو التحسين أو الميزة التنافسية.

## 4. الجهود المكررة:

قد يؤدي استكشاف البيانات غير الكافي إلى جهود زائدة أو غير ضرورية في جمع البيانات أو معالجتها مسبقاً أو نمذجتها، مما يؤدي إلى إهدار الموارد وزيادة تكاليف المشروع.

## 5. فقدان ثقة أصحاب المصلحة:

قد تؤدي التحليلات غير الدقيقة أو النماذج غير الموثوقة الناتجة عن إهمال تحليل البيانات الاستكشافي (EDA) إلى اهتزاز ثقة أصحاب المصلحة في عمليات صنع القرار القائمة على البيانات، مما قد يؤدي إلى تشويه سمعة المنظمة ومصداقيتها.

المخاطر المرتبطة بإهمال تحليل البيانات الاستكشافي (EDA) تمتد إلى ما هو أبعد من مجرد فشل المشروع إلى اتخاذ قرارات غير دقيقة وزيادة التكاليف وتضرر السمعة والمسؤوليات القانونية وفقدان القدرة التنافسية وتفويت فرص الابتكار. لتقليل هذه المخاطر، يجب أن تركز المؤسسات على الاستكشاف الكامل للبيانات وفهمها بشكل شامل طوال دورة حياة علم البيانات.

# تبعات ومخاطر تجاهل تحليل البيانات الاستكشافي (EDA)

تحليل البيانات الاستكشافي (EDA) هو مرحلة محورية في أي مشروع للتعلم الآلي (ML) أو علم البيانات، ويتضمن فحصاً دقيقاً لمجموعة البيانات بهدف استخراج رؤى ذات معنى واكتشاف الأنماط الخفية وتحضير البيانات للتحليل المستقبلي. تجاهل تحليل البيانات الاستكشافي يمكن أن يؤدي إلى نتائج سلبية، بما في ذلك فشل المشروع واستنتاجات غير صحيحة وتكاليف باهظة في مختلف القطاعات.



## المعرفة في المجال

### التفسير الصحيح للبيانات:

استخدم الخبرة التخصصية لتفسير البيانات بطريقة صحيحة. يساعد فهم السياق التجاري و مصدر البيانات في تقديم رؤى حول أهمية و تأثير المتغيرات المختلفة.

### أمثلة من العالم الواقعي:

في مجال الرعاية الصحية، المعرفة التخصصية ضرورية. على سبيل المثال، إذا لم يأخذ نموذج يتنبأ بنتائج المرضى في الاعتبار ميزات ضرورية مثل الأمراض المصاحبة أو تاريخ الأدوية بسبب نقص المعرفة التخصصية، فقد يؤدي ذلك إلى تنبؤات غير صحيحة ويؤثر سلبًا على رعاية المرضى.



## ملاءمة الميزات

### تحديد المتغيرات المهمة:

حدد الميزات الأكثر أهمية لتحليلك بناءً على المعرفة التخصصية. يساعد هذا في تركيز الجهود على المتغيرات الأكثر تأثيرًا.



## فهم البيانات

فهم البيانات هو عملية التعرف بعمق على مجموعة البيانات الخاصة بك قبل البدء في أي تحليل أو نمذجة. يتضمن ذلك استكشاف البيانات من خلال الإحصاءات الإجمالية والاستفادة من المعرفة التخصصية وتحديد الميزات الأكثر صلة بتحليلك. يضمن الفهم الصحيح للبيانات أن يكون تحليلك دقيقًا وذا معنى وقابلًا للتنفيذ.



## نظرة عامة عالية المستوى

### الاتجاه المركزي والتشتت:

ابدأ بالإحصاءات الإجمالية مثل المتوسط، الوسيط، المنوال، الانحراف المعياري، والمدى. تقدم هذه المقاييس لمحة سريعة عن الاتجاه المركزي و التشتت والتوزيع العام لبياناتك.

### التحليل الوصفي:

استخدم الإحصاءات الوصفية لتلخيص ووصف الخصائص الرئيسية لبياناتك. تساعد هذه الخطوة في تحديد الأنماط والاتجاهات والشذوذ في البيانات.

# أفضل الممارسات لتحليل بيانات استكشافي الفعال

يتضمن تحليل البيانات الاستكشافي الفعال اتباع طريقة منهجية لفهم البيانات من خلال الإحصاءات الإجمالية واستكشاف المتغيرات وتنقية البيانات لمعالجة التناقضات والأخطاء، والأهم من ذلك تمثيل البيانات باستخدام الرسومات البيانية والمخططات للكشف عن العلاقات والاتجاهات.

## أمثلة من قطاعات صناعية مختلفة:

**الرعاية الصحية:** في مجموعة بيانات الرعاية الصحية، يعتبر معرفة أهمية المتغيرات مثل عمر المريض وتاريخه الطبي ونتائج الفحوصات المخبرية أمراً أساسياً. على سبيل المثال، إذا كانت ميزة مثل معلومات حساسية المريض مفقودة من نموذج يتنبأ بفعالية الأدوية، فقد يؤدي ذلك إلى ردود فعل سلبية خطيرة، مما قد يهدد حياة المرضى.

**التكنولوجيا المالية (Fintech):** في القطاع المالي، يعتبر فهم الميزات مثل درجات الائتمان وسجلات المعاملات ومستويات الدخل أمراً بالغ الأهمية. بدون المعرفة التخصصية، قد يتجاهل نموذج يتنبأ بتعثر القروض عوامل هامة مثل التغييرات الأخيرة في الحالة الوظيفية أو الاتجاهات الصناعية، مما يؤدي إلى تقييمات غير دقيقة للمخاطر وخسائر مالية.

من خلال تلخيص الإحصاءات الإجمالية والاستفادة من المعرفة التخصصية، يمكنك ضمان تمثيل بياناتك بدقة وأن تكون ذات صلة بأهدافك. هذا النهج الشامل لا يعزز فقط جودة تحليلك، بل يقلل أيضاً من المخاطر ويؤدي إلى اتخاذ قرارات أكثر استنارة.

## تنقية البيانات



تنقية البيانات هي خطوة أساسية في مرحلة المعالجة المسبقة للبيانات، وتتضمن تحديد وتصحيح (أو إزالة) الأخطاء والتناقضات في البيانات. تعتبر هذه الخطوة ضرورية لضمان سلامة وجودة البيانات قبل إجراء أي تحليل أو بناء للنماذج. تنقية البيانات الفعال يعزز من موثوقية النتائج ويحسن أداء خوارزميات التعلم الآلي.

## التعامل مع القيم المفقودة



**المتوسط، الوسيط، والمنوال:** بالنسبة للبيانات الرقمية، يمكن تعويض القيم المفقودة باستخدام المتوسط أو الوسيط أو المنوال. هذه الطرق بسيطة ولكنها قد لا تكون مناسبة دائماً إذا كانت البيانات غير موزعة بشكل متماثل.

**التقنيات المتقدمة:** استخدم حلول متقدمة مثل خوارزمية K-Nearest Neighbors (KNN)

والتي يمكنها التنبؤ بالقيم المفقودة بناءً على قيم المتجاورة، أو بالتعويض المتعدد والذي يقوم بنمذجة كل متغير ذو قيم مفقودة كدالة معتمدة على المتغيرات الأخرى.

## اكتشاف القيم المتطرفة والأخطاء



**الطرق الإحصائية:** استخدم Z-scores لتحديد النقاط التي تقع على بعد عدة انحرافات معيارية عن المتوسط، أو استخدم النطاق الربعي (IQR) لاكتشاف القيم المتطرفة عن طريق تحديد القيم التي تقع خارج النطاق الربعي أعلى من الربع الثالث وأسفل الربع الأول بمقدار 1.5 من كل جانب.

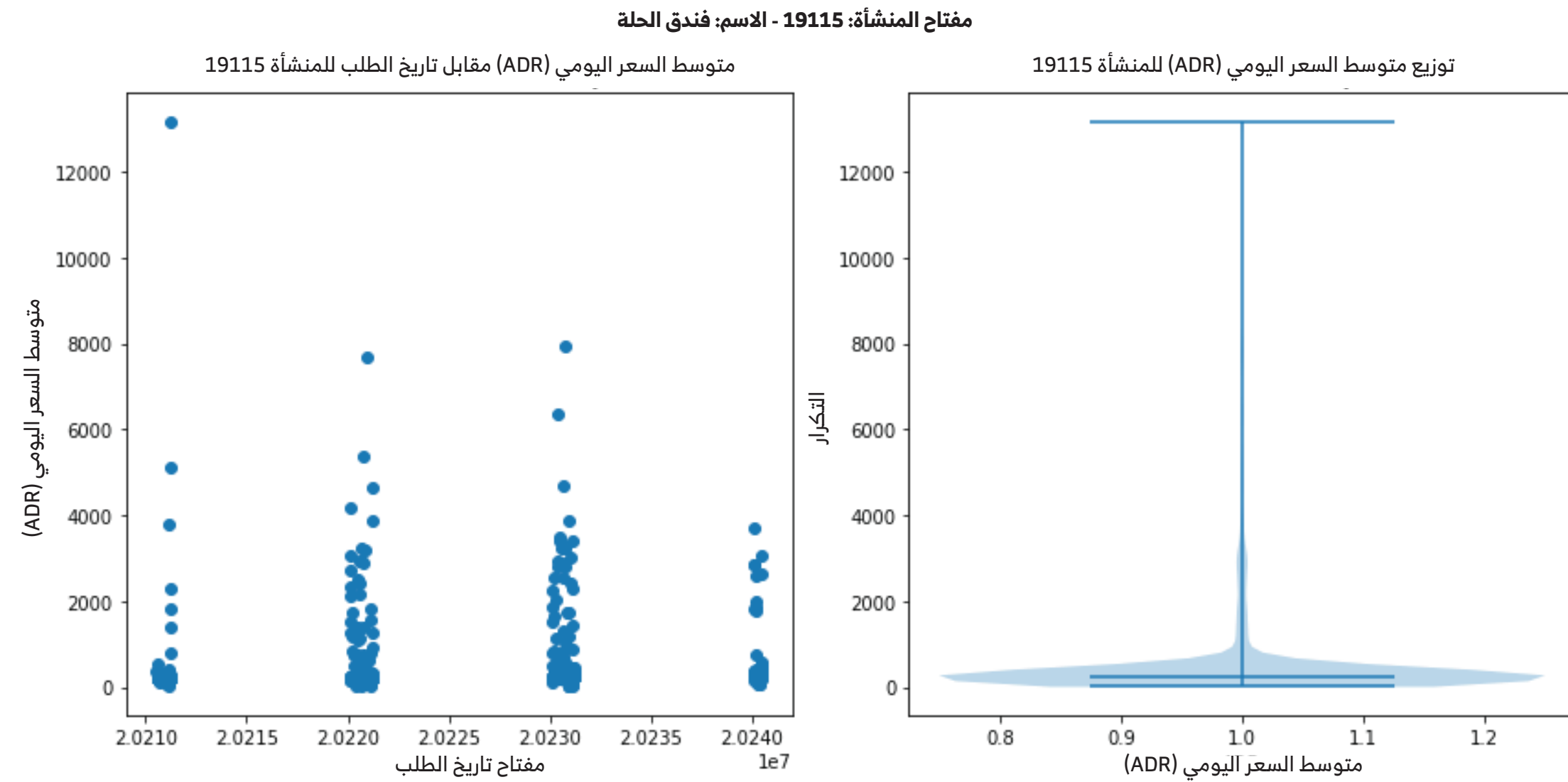
**أدوات تمثيل البيانات:** تعتبر رسومات الصناديق البيانية ورسومات النقاط المشتتة فعالة في تمثيل القيم المتطرفة وفهم تأثيرها على توزيع البيانات.



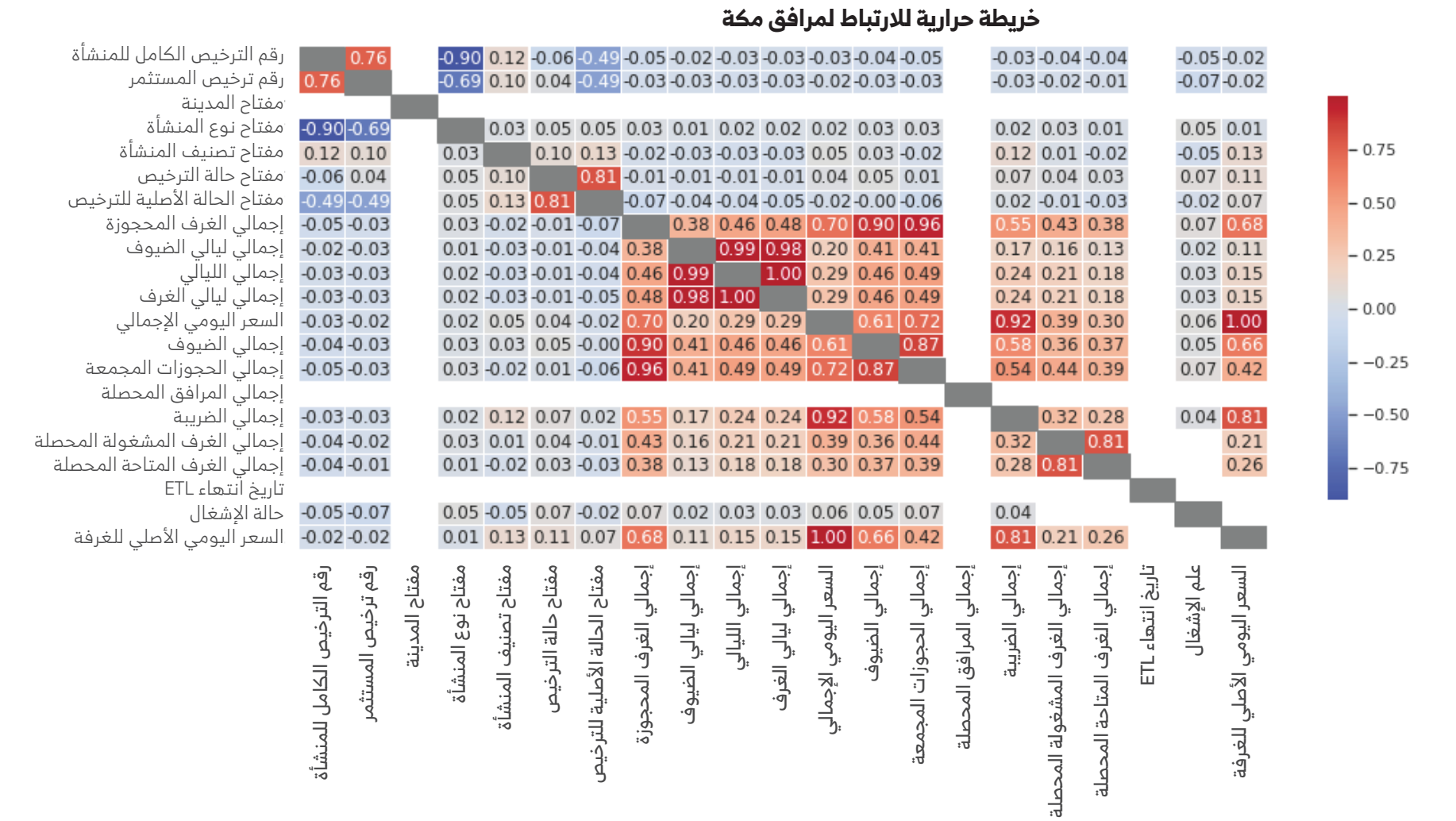
## معالجة القيم المتطرفة والأخطاء

**الإزالة أو التحويل:** قرر ما إذا كنت ستزيل القيم المتطرفة، أو تحولها (على سبيل المثال، باستخدام التحويل اللوغاريتمي أو طريقة وينسور)، أو تحدد سقفًا لها لتقليل تأثيرها. يجب أن يستند هذا القرار إلى طبيعة البيانات وأهداف التحليل.

**تصحيح الأخطاء:** قم بتصحيح أخطاء إدخال البيانات و تطرفها. قد يتضمن ذلك التحقق من النقاط غير الصحيحة وتحديثها أو إعادة إدخال البيانات من مصادر موثوقة.



الشكل 2. تصور بيانات المنشآت قبل معالجة القيم المتطرفة



الشكل 1. خريطة حرارية للارتباط لمرافق مكة



## معالجة اختلال توازن الفئات

### الزيادة في أخذ العينات (Oversampling):

قم بزيادة عدد العينات في فئة الأقلية عن طريق تكرار العينات أو توليد عينات جديدة (مثل تقنية الإفراط في أخذ العينات لفئة الأقلية الاصطناعية (SMOTE).

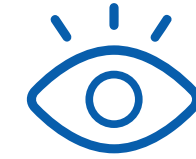
### النقص في أخذ العينات (Under-sampling):

تقليل عدد العينات في الفئة الأغلبية لتحقيق التوازن في مجموعة البيانات. يمكن أن يكون ذلك فعالاً ولكنه قد يؤدي إلى فقدان معلومات مهمة من الفئة الأغلبية.

### طرق الدمج:

قم بالجمع بين الزيادة في أخذ العينات والتقليل منها لتحقيق مجموعة بيانات متوازنة دون فقدان كبير للمعلومات.

من خلال تطبيق أفضل ممارسات تنقية البيانات هذه، يمكنك تحسين جودة مجموعة البيانات بشكل كبير، مما يؤدي إلى نتائج تحليلية ونماذج تنبؤية أكثر موثوقية وصحة.



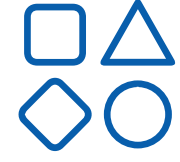
## تمثيل البيانات

تمثيل البيانات هي عملية عرض المعلومات والبيانات عن طريق الرسوم البيانية. من خلال استخدام العناصر المرئية مثل المخططات والرسوم البيانية والخرائط، توفر أدوات تمثيل البيانات طريقة ميسرة لرؤية وفهم الاتجاهات والقيم المتطرفة والأنماط في البيانات. يستعرض هذا القسم أهمية تمثيل البيانات، وكيف يمكن استخدامه بفعالية، وأنواع التصورات المختلفة التي يمكن توظيفها.

**الوضوح:** التمثيلات المرئية تجعل البيانات المعقدة أكثر قابلية للفهم والوصول إليها، مما يسهل تحديد الأنماط والرؤى التي قد يتم تجاهلها عند التعامل مع البيانات الخام.

**الكفاءة:** يتيح التمثيل المرئي للبيانات الكبيرة الفهم السريع لها، مما يساعد المحللين وأصحاب المصلحة على اتخاذ قرارات أسرع وأكثر استنارة.





## أنواع رسومات بيانية محددة

**Area Charts:** اتجاهات البيانات التراكمية: إظهار الاتجاهات عبر الزمن.

**Bubble Charts:** تحليل متعدد المتغيرات: تحليل العلاقات بين المتغيرات.

**Tree Maps:** البيانات الهرمية: تصور العلاقات بين الأجزاء والكل.

**Scatter Plots:** تحليل الارتباط: استكشاف العلاقات بين المتغيرات.

**Histograms:** تحليل التوزيع: تصور الانتشار والاتجاه المركزي لمتغير واحد.

**Violin Plots:** مقارنة التوزيع: مقارنة توزيع البيانات عبر مجموعات مختلفة.

**Heatmaps:** مصفوفة الارتباط: تصور العلاقات بين الميزات.

**Bar Charts:** مقارنة البيانات الفئوية: مقارنة الفئات المختلفة.

**Line Charts:** تحليل الاتجاهات: تحديد الاتجاهات عبر الزمن.

**Box Plots:** الإحصاءات الملخصة: اكتشاف القيم المتطرفة ومقارنة التوزيعات.

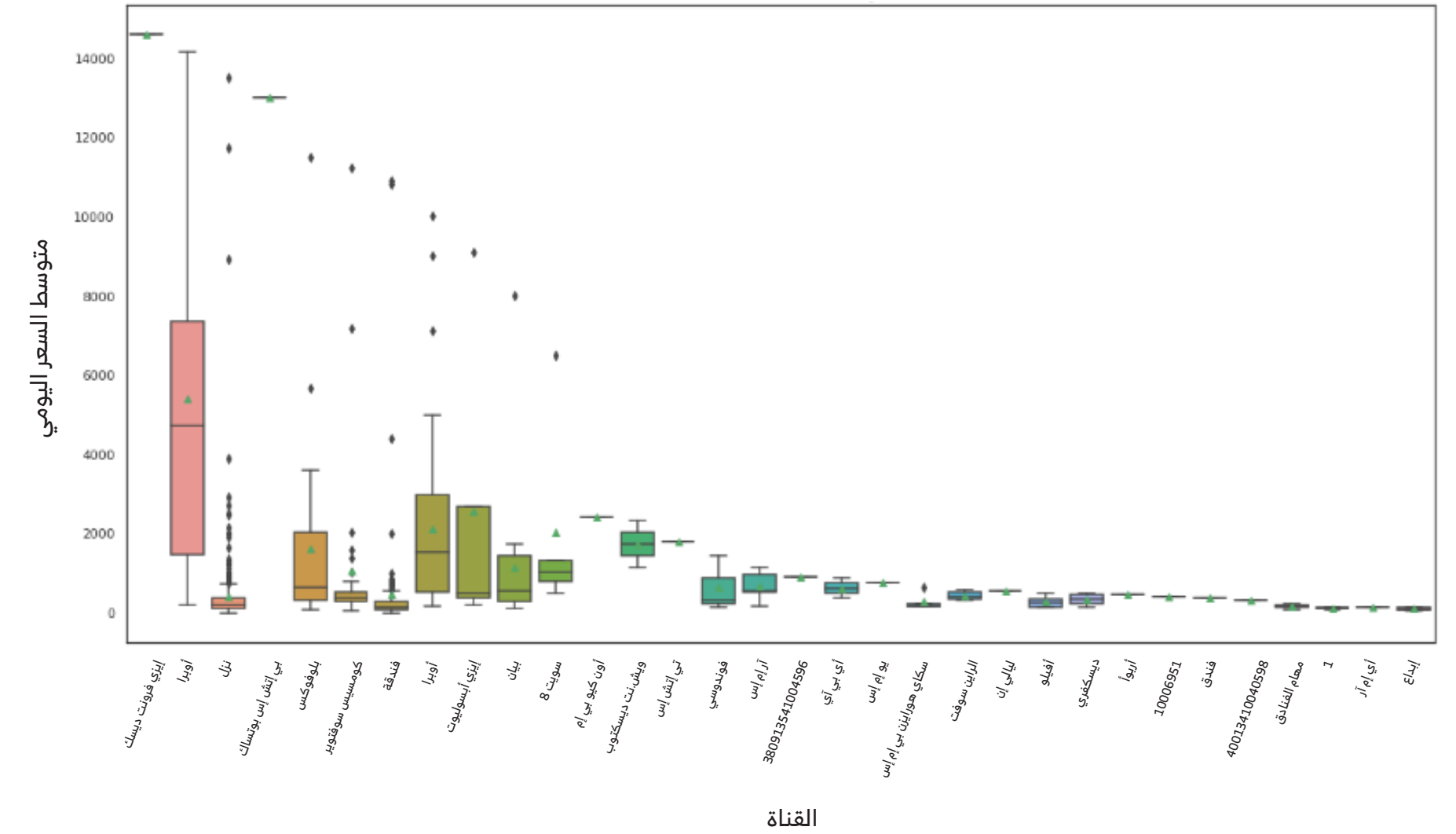
**Pie Charts:** تمثيل النسب: عرض نسب مكونات الكل.

من خلال اختيار واستخدام تقنيات تصور البيانات المناسبة، يمكنك الكشف عن القصة المخفية داخل بياناتك، مما يعزز التواصل الواضح واتخاذ القرارات المستنيرة.

**تفاعل أصحاب المصلحة:** تساعد التصورات المرئية المصممة جيداً في إيصال النتائج بشكل فعال إلى أصحاب المصلحة الذين قد لا يمتلكون خبرة تقنية، مما يسهل اتخاذ قرارات أفضل وتطوير استراتيجيات محسنة.

**سرد القصص:** يمكن للتصورات المرئية أن تروي قصة مقنعة باستخدام البيانات، مع تسليط الضوء على الرؤى والاتجاهات الرئيسية بطريقة واضحة وجذابة.

توزيع القناة حسب متوسط السعر اليومي (ADR)



الشكل 3. تصور توزيع متوسط السعر اليومي (ADR) عبر القنوات لمكة المكرمة

# أدوات وأطر تحليل البيانات الاستكشافي (EDA)



## الأدوات القائمة على البرمجة

**Python's pandas**: ضرورية لمعالجة البيانات وتحليلها باستخدام هياكل البيانات (DataFrames).

**NumPy**: هي برامج فعالة لإجراء العمليات الحسابية الرقمية وعمليات الجبر الخطي.

**SciPy**: يعد ساي باي إمتدادا ل نم باي (NumPy) ويتضمن عمليات للتحسين والتكامل والإحصاءات.

**Matplotlib and Seaborn**: هي مكتبات أساسية لإنشاء رسومات بيانية إحصائية ثابتة وجذابة.

**Plotly**: هي أداة لإنشاء رسومات تفاعلية وجذابة للوحات التحكم والتصورات على الويب.

**Jupyter Notebook**: هي أداة لإنشاء ومشاركة مستندات تحتوي على كود حي، ورسومات بيانية، ونصوص.



## واجهات سهلة الاستخدام

**Excel**: هو برنامج أساسي لتحليل البيانات وتصويرها باستخدام واجهة سهلة الاستخدام.

**Tableau**: هو برنامج قوي يقدم قدرات تفاعلية لتصوير البيانات.

**Microsoft Power BI**: هو برنامج شامل لتحليلات الأعمال مع تصورات تفاعلية.

**Google Data Studio**: هو برنامج لإنشاء لوحات تحكم وتقارير قابلة للمشاركة، مع دمج لمصادر البيانات المختلفة.



## أدوات متقدمة لتحليل البيانات الاستكشافي الفعال

**D3.js**: هو مكتبة جافا سكريبت لإنشاء تصورات ويب ديناميكية وتفاعلية.

**R Shiny**: هو أداة لبناء تطبيقات ويب تفاعلية مباشرة باستخدام لغة آر (R).

**Orange**: هو أداة مفتوحة المصدر لتحليل البيانات الاستكشافي واستكشاف البيانات التفاعلي.

**KNIME**: هو منصة لتحليل البيانات باستخدام خطوط بيانات وحدوية.

**RapidMiner**: هو بيئة متكاملة لتحضير البيانات والتعلم الآلي والتحليلات التنبؤية.

**QlikView**: هو أداة ذكاء الأعمال لتصوير البيانات التفاعلي واكتشاف البيانات.

**Databricks**: هو منصة تحليلات موحدة تتكامل مع أباتشي سبارك (Apache Spark) لمعالجة البيانات على نطاق واسع.



# دراسة حالة

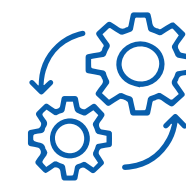
اليومي (ADR) والمتغيرات الأخرى لتحديد العوامل المؤثرة.

## اكتشاف القيم المتطرفة ومعالجتها

- حدد القيم المتطرفة في بيانات متوسط السعر اليومي (ADR) باستخدام الأساليب الإحصائية مثل زيسكور (Z-Score) أو المدى الرباعي (IQR).
- تصور القيم المتطرفة باستخدام الرسومات الصندوقية أو رسومات النقط المشتتة لفهم توزيعها.
- طبق تقنيات معالجة القيم المتطرفة مثل مثل إزالة البيانات المتطرفة وتحويلها.

## تحليل حسب المدينة

- حلل اتجاهات متوسط السعر اليومي (ADR) عبر المدن المختلفة لتحديد الاختلافات الجغرافية.
- استكشف العوامل المساهمة في ارتفاع قيم متوسط السعر اليومي (ADR) في مدن معينة من خلال التحليل الإحصائي والتصور.



## عملية تحليل البيانات الاستكشافي (EDA)

### تحميل البيانات وتحضيرها

- استخدم مكتبات Python مثل Pandas لتحميل البيانات ومعالجتها مسبقًا.
- تعامل مع القيم المفقودة، والقيم المتطرفة، والتناقضات في البيانات لضمان سلامة البيانات.
- عالج القيم المفقودة بتقنيات مثل التعويض بالمتوسط أو الوسيط أو القيم الأكثر تكرارًا.
- استفد من داتايكو (Dataiku) لتحميل البيانات وتسهيل عمليات التحضير الأولية للبيانات.

### الإحصاءات الوصفية والتصور

- اوجد الإحصاءات الإجمالية باستخدام Pandas لفهم توزيع البيانات والاتجاه المركزي.
- تصور توزيع متوسط السعر اليومي (ADR) عبر جميع المنشآت والمدن باستخدام Matplotlib أو Seaborn.
- استكشف العلاقات بين متوسط السعر

## استراتيجيات السياحة المدفوعة بالبيانات: استكشاف اتجاهات متوسط السعر اليومي (ADR) من خلال تحليل البيانات الاستكشافي (EDA)



### المقدمة:

تحليل البيانات الاستكشافي (EDA) هو خطوة ضرورية في كشف الرؤى والأنماط داخل مجموعات البيانات، وخاصة في إدارة الإيرادات للمنشآت السياحية. تتناول هذه الدراسة التطبيقية استخدام أطر البرمجة في بايثون (Python)، بما في ذلك داتايكو (Dataiku)، لإجراء تحليل بيانات استكشافي متعمق على بيانات متوسط السعر اليومي (ADR) عبر مدن ومنشآت مختلفة. هدفنا هو عرض نهج شامل لتحليل البيانات الاستكشافي، مع التركيز على اكتشاف القيم المتطرفة، وتحليل الاتجاهات، والعوامل المؤثرة في تباين متوسط السعر اليومي (ADR).

### مجموعة البيانات:

تتضمن مجموعة البيانات معلومات عن متوسط السعر اليومي (ADR) للمنشآت السياحية عبر مدن متعددة، وتشمل مؤشرات مثل إجمالي السعر اليومي (TOTAL\_DAILY\_RATE) وإجمالي الغرف المحجوزة (TOTAL\_COLLECTED\_BOOKED\_ROOMS) وتفاصيل المنشآت. توفر هذه البيانات رؤى حول اتجاهات متوسط السعر اليومي، القيم المتطرفة، والعوامل المؤثرة على الإيرادات عبر منشآت ومدن



## الخاتمة

يعد تحليل البيانات الاستكشافي (EDA) الشامل أمراً ضرورياً لتحسين استراتيجيات إدارة الإيرادات في قطاع السياحة. من خلال التعمق في اتجاهات متوسط السعر اليومي (ADR)، والقيم المتطرفة، والعوامل المحورية، يمكن لأصحاب المصلحة الاستفادة من الرؤى المستندة إلى البيانات لتعزيز الربحية ومواجهة التحديات المتغيرة باستمرار في القطاع السياحي.

- تحقق من الأنماط الزمنية والاتجاهات الموسمية التي تؤثر على تباين متوسط السعر اليومي (ADR) في كل مدينة.

### رؤى محددة على مستوى المنشآت

- تعمق في بيانات المنشآت لفهم أسباب ارتفاع قيم متوسط السعر اليومي (ADR) في بعض المنشآت.
- حدد نقاط البيع الفريدة، أو المرافق، أو استراتيجيات التسويق التي تساهم في تعظيم الإيرادات.
- تحقق من القيم المتطرفة على مستوى المنشأة واستكشف الأسباب المحتملة لحدوثها.

### قيم متوسط السعر اليومي (ADR) المرتفعة كقيم متطرفة

- استكشف الأسباب وراء القيم المرتفعة لمتوسط السعر اليومي (ADR) التي تعتبر قيمًا متطرفة.
- تحقق من العوامل مثل الأحداث الخاصة، أو المواسم المزدحمة، أو العروض المميزة التي تساهم في ارتفاع متوسط السعر اليومي (ADR) بشكل استثنائي.
- حلل القنوات التي تبلغ عن هذه القيم المتطرفة لفهم توزيعها ومصدرها.



# الاستنتاجات النهائية

يعد تحليل البيانات الاستكشافي (EDA) الشامل الركيزة الأساسية لاتخاذ القرارات المستندة إلى البيانات، حيث يمكن المؤسسات من اكتشاف رؤى مخفية وتحقيق نتائج مؤثرة.

من خلال تبني أفضل الممارسات والاستفادة من الأدوات المتقدمة، يمكن للمحللين تعظيم قيمة بياناتهم، مما يعزز الابتكار، ويحسن الكفاءة التشغيلية، ويكسب ميزة تنافسية في المشهد القائم على البيانات اليوم.

لا يضمن تحليل البيانات الاستكشافي (EDA) المتعمق موثوقية ونزاهة الرؤى التحليلية فحسب، بل يساهم أيضاً في تعزيز ثقافة اتخاذ القرارات المستنيرة والتوقعات الاستراتيجية.

ومع استمرار المؤسسات في التعامل مع تعقيدات بيئة الأعمال الحديثة، يصبح الاستثمار في ممارسات تحليل البيانات الاستكشافي (EDA)

أمراً استراتيجياً لدفع النمو المستدام، وتقليل المخاطر، وتقديم قيمة غير مسبوقة لأصحاب المصلحة.

