



إسم المادة: التنقيب عن البيانات

إسم المحاضر: م. خليل محمد

الأكاديمية العربية الدولية – منصة أعد

مفهوم التنقيب عن البيانات

التنقيب عن البيانات، أو ما يُعرف بـ **Data Mining**، هو عملية تحليل مجموعات كبيرة من البيانات لاكتشاف أنماط، علاقات، واتجاهات جديدة ومفيدة لم تكن معروفة من قبل. يُعتبر التنقيب عن البيانات جزءاً من مجال أوسع يُعرف باكتشاف المعرفة في قواعد البيانات.

المفاهيم الأساسية في التنقيب عن البيانات:

1. **استخراج الأنماط:** في هذه المرحلة، يتم تحليل البيانات لاستخراج الأنماط المتكررة، مثل الترابطات بين المتغيرات.
2. **التصنيف:** التصنيف هو عملية تخصيص العناصر إلى فئات محددة بناءً على ميزات معينة.
3. **التنظيم(Clustering):** على عكس التصنيف الذي يعتمد على الفئات المعروفة مسبقاً، يُستخدم التنظيم لتجميع البيانات بناءً على التشابه بينها.
4. **الارتباطات:** تهدف إلى اكتشاف العلاقات بين المتغيرات المختلفة في قاعدة البيانات.
5. **التنبؤ:** يستخدم التنبؤ النماذج المستخلصة من البيانات الحالية للتنبؤ بقيم المتغيرات المستقبلية.
6. **تحليل التتابع(Sequence Analysis):** يركز على اكتشاف الأنماط في البيانات التي تحدث بترتيب معين.

الخطوات الأساسية في التقيب عن البيانات

1. **جمع البيانات:** الخطوة الأولى هي جمع البيانات من مصادر مختلفة. قد تشمل هذه البيانات بيانات رقمية، نصية، أو صور.
2. **تنظيف البيانات:** بعد جمع البيانات، تأتي مرحلة تنظيفها من الأخطاء والتكرارات.
3. **التحليل الاستكشافي:** قبل بدء عملية التقيب، يُجرى تحليل استكشافي للبيانات لفهم خصائصها والتعرف على الأنماط الأولية وال العلاقات المحتملة بين المتغيرات.
4. **تطبيق التقنيات:** هنا، يتم تطبيق تقنيات التقيب عن البيانات مثل التجميع(Clustering) ، التصنيف(Classification) ، والتوقع(Prediction) لاكتشاف الأنماط والاتجاهات.
5. **التقييم والتفسير:** بعد تطبيق التقنيات، تُقيّم النتائج للتحقق من صحتها وقيمتها.
6. **عرض النتائج:** الخطوة الأخيرة هي تقديم النتائج بطريقة سهلة الفهم، من خلال الرسوم البيانية والتقارير، لتمكين الأطراف المعنية من استخدام المعلومات المكتشفة في اتخاذ القرارات الاستراتيجية.

أهمية التقريب عن البيانات

- 1 . اكتشاف الأنماط والاتجاهات: من خلال تطبيق تقنيات التقريب عن البيانات مثل التجميع (Clustering) والتصنيف (Classification) ، يمكن اكتشاف أنماط غير مرئية في البيانات.
- 2 . تحسين اتخاذ القرارات: من خلال استخراج المعلومات القيمة، يمكن التقريب عن البيانات من تقديم رؤى دقيقة ومدعومة بالبيانات، مما يعزز القدرة على اتخاذ قرارات مبنية على معلومات موثوقة.
- 3 . التنبؤ بالأحداث المستقبلية: يمكن للتقريب عن البيانات استخدام تقنيات التوقع (Prediction) للتنبؤ بالأحداث المستقبلية بناءً على الأنماط الحالية. على سبيل المثال، يمكن استخدامه للتنبؤ بمبيعات المستقبل بناءً على بيانات المبيعات التاريخية.
- 4 . تحسين العمليات والكافاءات: يمكن استخدام التقريب عن البيانات لتحليل وتحسين العمليات الداخلية في المؤسسات. على سبيل المثال، يمكن تحليل بيانات الإنتاج لتحسين كفاءة خطوط الإنتاج.
- 5 . اكتشاف الفرص التجارية: يمكن للتقريب عن البيانات الكشف عن فرص تجارية جديدة من خلال تحليل بيانات السوق والاتجاهات الحالية.

التطبيقات العملية للتذكير عن البيانات

- **أمثلة على التطبيقات:** التسويق (استهداف العملاء)، الصحة (توقع الأمراض)، والأمن (الكشف عن الاحتيال).
- 1. **التسويق: استهداف العملاء:** يستخدم التذكير عن البيانات في التسويق لتحليل سلوك العملاء وتوقع احتياجاتهم المستقبلية. من خلال تحليل بيانات المبيعات، وتاريخ الشراء، وسلوك التصفح، يمكن للشركات تحديد الأنماط والاتجاهات. هذا يسمح للشركات بتطوير استراتيجيات تسويقية موجهة، مثل تخصيص العروض الترويجية وتقديم توصيات مخصصة.
- 2. **الصحة: توقع الأمراض:** في مجال الصحة، يمكن للتذكير عن البيانات أن يساعد في توقع انتشار الأمراض وتحليل الاتجاهات الصحية. من خلال تحليل سجلات المرضى، وبيانات التشخيص، والتقارير الطبية، يمكن تحديد عوامل الخطر والتنبؤ بالأمراض المحتملة.
- 3. **الأمن: الكشف عن الاحتيال:** في الأمن، يستخدم التذكير عن البيانات للكشف عن الأنشطة المشبوهة ومنع الاحتيال. يمكن تحليل سجلات المعاملات، وسلوك المستخدمين، والنشاطات غير الطبيعية لاكتشاف الأنماط المشبوهة وتحليلها.

الأساليب الأساسية للتقييّب عن البيانات - التحليل الوصفي

تعريف: هو عملية تهدف إلى فهم البيانات من خلال تلخيص الأنماط والاتجاهات فيها دون محاولة التنبؤ أو تعميم النتائج على مجموعات أخرى.

أهداف التحليل الوصفي:

1. **تلخيص البيانات:** تقديم ملخص بسيط للبيانات الكبيرة، مما يسهل فهمها وتحليلها. يشمل ذلك حساب مقاييس أساسية مثل الوسط الحسابي، الوسيط، والمنوال.
2. **التعرف على الأنماط:** اكتشاف الأنماط أو الاتجاهات الأساسية التي قد لا تكون واضحة من خلال النظر في البيانات الخام فقط.
3. **تحليل التشتت:** قياس مدى تباين البيانات، الذي يمكن أن يوفر رؤى حول مدى تجانس أو تباين القيم.

التقنيات المستخدمة في التحليل الوصفي:

1. **الإحصاء الوصفي:**
 - **الوسط الحسابي:** هو المعدل البسيط لجميع القيم في مجموعة بيانات. يستخدم لفهم النقطة المركزية لتوزيع البيانات.
 - **الوسيط:** القيمة التي تقع في منتصف مجموعة البيانات عندما يتم ترتيبها. يستخدم عندما تكون هناك قيم شاذة تؤثر على الوسط الحسابي.
 - **المنوال:** القيمة التي تتكرر بشكل أكثر تكراراً في مجموعة البيانات. يوفر معلومات حول القيمة الأكثر شيوعاً.

الأساليب الأساسية للتقييّب عن البيانات - التحليل الوصفي

2. التحليل البياني:

- الرسوم البيانية: مثل الرسوم البيانية الشريطية، والرسوم البيانية الدائرية، والرسوم البيانية الخطية، تُستخدم لتصوير البيانات بصرياً. تساعد هذه الرسوم في تقديم فكرة واضحة وسريعة عن توزيع البيانات والاتجاهات الرئيسية.
- الرسوم البيانية الصندوقية: تُستخدم لتوضيح التوزيع الكمي للبيانات، بما في ذلك تحديد القيم الشاذة والتباين.

3. مقاييس التشتت:

- الانحراف المعياري: يقيس مدى تباين البيانات حول الوسط الحسابي. كلما كان الانحراف المعياري أكبر، زادت درجة التباين في البيانات.
- التباين: هو مربع الانحراف المعياري ويعكس مدى تباين البيانات.
- المدى: الفرق بين أعلى وأدنى قيمة في مجموعة البيانات. يُستخدم لتحديد نطاق القيم التي تشملها البيانات.

الأساليب الأساسية للتقييّب عن البيانات - التحليل الوصفي

1. تحليل توزيع البيانات:

- التوزيع التكراري: يستخدم لتصنيف البيانات إلى فئات محددة وحساب عدد القيم التي تقع في كل فئة. هذا يمكن أن يساعد في فهم كيف يتم توزيع البيانات بشكل عام.
- الرسم البياني التوزيعي: **Histogram** يعرض توزيع البيانات عبر نطاقات مختلفة، مما يتيح رؤية كيف تتواءم القيم في مجموعة البيانات.

2. الأنماط العامة:

- التحليل الزمني: يستخدم لمراقبة كيفية تغير القيم بمرور الوقت. هذا يمكن أن يتضمن دراسة الاتجاهات العامة أو الأنماط الموسمية.
- الأنماط الموسمية: يتناول الأنماط التي تتكرر بشكل دوري. يمكن تحليل البيانات لتحديد تأثيرات موسم معين أو تغييرات دورية في الأنشطة.

3. التعرف على القيم الشاذة:

- الرسم البياني الصندوقي: **Box Plot** يستخدم لتحديد القيم الشاذة أو الخارجة عن المألوف في مجموعة البيانات. يُظهر هذا الرسم نطاق القيم ويحدد القيم التي تقع خارج النطاق المعتاد.
- مقاييس التشتت: مثل المدى، والانحراف المعياري، التي تساعده في قياس مدى تباين البيانات حول المتوسط.

تقنيات التزقيب عن البيانات – الاستدلال الإحصائي

التقنيات الرئيسية في الاستدلال الإحصائي:

1. **التنبؤ:** تشمل بناء نماذج إحصائية للتنبؤ بالنتائج المستقبلية بناءً على الأنماط الموجودة في البيانات التاريخية. يمكن استخدام نماذج مثل الانحدار الخطي أو الانحدار اللوجستي لتحليل العلاقات بين المتغيرات وتوقع القيم المستقبلية.
2. **الاختبارات الإحصائية:** تُستخدم لاختبار صحة افتراضات حول البيانات. تتضمن هذه الاختبارات مقارنة البيانات الفعلية مع التوقعات لمعرفة ما إذا كانت هناك دلائل على أن الاختلافات غير ناتجة عن الصدفة.
3. **التحليل الاستكشافي للبيانات (EDA):** يتضمن استخدام مقاييس مثل الوسط الحسابي والانحراف المعياري لفهم توزيع البيانات وتحديد الأنماط الرئيسية.
4. **التزقيب عن الأنماط: استخراج الأنماط:** يشمل تحديد الأنماط والعلاقات المخفية في البيانات التي قد لا تكون واضحة من خلال التحليل البسيط. يمكن أن تشمل هذه الأنماط الترابطات بين المتغيرات، الأنماط الزمنية، أو الأنماط الجغرافية.

تقنيات التزقيب عن البيانات – التجميع Clustering

أنواع خوارزميات التجميع

1. التجميع غير المراقب (Unsupervised Clustering):

- **K-Means:** يقسم البيانات إلى k مجموعة (عناقيد) بناءً على متوسط نقاط البيانات داخل كل مجموعة.
 - القيود: يتطلب تحديد عدد العناقيد k مسبقاً، وقد لا يعمل بشكل جيد مع العناقيد غير المتجانسة.
- **Hierarchical Clustering:** ينشئ شجرة هرمية من العناقيد عبر مرحلتين: التكوين من الأسفل إلى الأعلى (أو العكس) لتجميع البيانات بناءً على المسافة أو التشابه.
 - القيود: يمكن أن يكون بطيناً عند التعامل معمجموعات بيانات كبيرة.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise):

- **المفهوم:** يحدد العناقيد بناءً على كثافة النقاط في الفضاء، حيث تجمع النقاط ذات الكثافة العالية في مجموعة واحدة وتعتبر النقاط ذات الكثافة المنخفضة كضوضاء.
- القيود: يتطلب تحديد معلمات الكثافة، قد يكون صعباً للتنفيذ في بيانات عالية الأبعاد.

2. التجميع المراقب (Supervised Clustering):

يستخدم المعلومات المصنفة (الملصقات) لتوجيه عملية التجميع بحيث تكون العناقيد الناتجة ذات معنى بناءً على المعرفة المسبقة.

- القيود: يتطلب بيانات معلمة مسبقاً، وقد يكون أقل فعالية إذا كانت البيانات قليلة أو غير ممثلة بشكل جيد.

تقنيات التزقيب عن البيانات – التصنيف Classification

أنواع خوارزميات التصنيف

1. الانحدار اللوجستي(**Logistic Regression**) : يستخدم لتصنيف البيانات إلى فئتين أو أكثر بناءً على علاقة خطية بين المتغيرات المستقلة والاحتمالية التابعة للانتماء إلى فئة معينة.

- القيود: قد يكون محدوداً في معالجة علاقات غير خطية بين المتغيرات.

2. شجرة القرار(**Decision Tree**) : ينشئ نموذجاً شجرياً يستخدم مجموعة من القواعد البسيطة لتصنيف البيانات بناءً على قيم الخصائص.

- القيود: يمكن أن يكون عرضة للتجاوز (Overfitting) عندما يكون النموذج معقداً جداً.

3. دعم الآلات الشعاعية(**Support Vector Machine - SVM**) : يحاول إيجاد أفضل خط أو مستوى يفصل بين الفئات المختلفة في مجموعة البيانات.

- القيود: قد يكون بطيئاً في التدريب علىمجموعات البيانات الكبيرة جداً.

4. الشبكات العصبية(**Neural Networks**) : تستخدم شبكة من العقد (الخلايا العصبية) لمعالجة المعلومات واكتساب أنماط معقدة في البيانات. تتكون الشبكة العصبية من طبقات إدخال، وطبقات مخفية، وطبقة إخراج.

- القيود: تحتاج إلى كمية كبيرة من البيانات والتدريب، وتعتبر أكثر تعقيداً في التنفيذ والفهم.

الأساليب المتقدمة في التنقيب عن البيانات - الخوارزميات

خوارزميات التجميع Clustering

K-Means Clustering: .1

- المفهوم: خوارزمية تجمع البيانات إلى K مجموعة (Cluster) بناءً على مدى تشابهها. تبدأ الخوارزمية باختيار K مراكز عشوائية (Centroids) ثم تقوم بتكرار عملية تعيين كل عنصر إلى أقرب مركز وتحريك المراكز بناءً على العناصر المخصصة لها، حتى تقارب المراكز بشكل ثابت.
- الميزات: بسيطة وفعالة في التعامل معمجموعات بيانات كبيرة. سريعة في التنفيذ عند العمل مع بيانات كبيرة.
- القيود: تحتاج إلى تحديد عدد المجموعات K مسبقاً، وقد تؤدي إلى نتائج غير جيدة إذا كانت البيانات تحتوي على مجموعات غير متساوية أو غير كروية.

Hierarchical Clustering: .2

- المفهوم: يتضمن بناء شجرة (Dendrogram) تُظهر العلاقات بين العناصر أو المجموعات. هناك نوعان من التجميع الهرمي: التجميع الهرمي التصاعدي (Agglomerative) الذي يبدأ بكل عنصر في مجموعة منفصلة ثم يدمج المجموعات تدريجياً، والتجميع الهرمي التنازلي (Divisive) الذي يبدأ بمجموعة واحدة كبيرة ثم يقسمها تدريجياً.
- الميزات: لا يتطلب تحديد عدد المجموعات مسبقاً ويقدم تصوراً مفصلاً عن العلاقة بين العناصر.
- القيود: قد يكون بطئاً في التعامل معمجموعات بيانات كبيرة، ويمكن أن يكون حساساً للضجيج في البيانات.

الأساليب المتقدمة في التنبؤ عن البيانات - الخوارزميات

خوارزميات التصنيف Classification

Decision Trees: .1

- المفهوم: تستخدم شجرة القرار لتصنيف البيانات بناءً على سلسلة من الأسئلة البسيطة التي تقسم البيانات إلى فئات مختلفة. تبدأ الشجرة بجذر (Root) وتفرع إلى عقد (Nodes) وأوراق (Leaves) تمثل الفئات النهائية.
- الميزات: سهلة الفهم والتفسير، وتوفر رؤية واضحة حول كيفية اتخاذ القرارات. يمكن التعامل مع البيانات المفقودة بشكل جيد.
- القيود: عرضة للتجاوز (Overfitting) إذا كانت الشجرة عميقه جداً وتحتوي على تفاصيل دقيقة جداً.

Support Vector Machines (SVM): .2

- المفهوم: تستخدم SVM لإنشاء حدود (Hyperplanes) تفصل بين الفئات المختلفة في الفضاء متعدد الأبعاد. الهدف هو إيجاد الحدود التي توفر أقصى فصل بين الفئات.
- الميزات: فعالة في التعامل مع البيانات ذات الأبعاد العالية وتعمل بشكل جيد مع البيانات غير الخطية عند استخدام kernel functions.
- القيود: قد تكون بطيئة في التدريب على مجموعات بيانات كبيرة، وتحتاج إلى ضبط معلمات الـ kernel بعناية.

الأساليب المتقدمة في التنبؤ عن البيانات - تحليل الارتباطات

التحليل الارتباطي: يركز على قياس قوة واتجاه العلاقة بين متغيرين. يساعد هذا التحليل في فهم كيفية تغير المتغيرات المرتبطة عندما يتغير أحدها.

أنواع الارتباطات

1. الارتباط الإيجابي:

التفصير: يشير إلى أن المتغيرات تتحرك في نفس الاتجاه.
المفهوم: يحدث عندما يزيد المتغير الثاني مع زيادة المتغير الأول.

2. الارتباط السلبي:

التفصير: يشير إلى أن المتغيرات تتحرك في اتجاهين متعاكسين.
المفهوم: يحدث عندما ينخفض المتغير الثاني مع زيادة المتغير الأول.

3. الارتباط الصفيري:

التفصير: لا يوجد ارتباط بين المتغيرين؛ أي أن تغييرات في أحد المتغيرات لا تؤثر على المتغير الآخر.
المفهوم: لا يوجد ارتباط بين المتغيرين؛ أي أن تغييرات في أحد المتغيرات لا تؤثر على المتغير الآخر.

الأساليب المتقدمة في التنقيب عن البيانات - تحليل الارتباطات

كيفية قياس الارتباط

1. معامل الارتباط (Correlation Coefficient):

المفهوم: هو رقم يتراوح بين -1 و $+1$ ويعكس مدى قوة العلاقة بين المتغيرات.

أنواع المعاملات:

معامل الارتباط بيرسون (Pearson Correlation Coefficient): يقيس العلاقة الخطية بين المتغيرات.

معامل الارتباط سبيرمان (Spearman's Rank Correlation Coefficient): يستخدم لقياس العلاقة بين المتغيرات عندما تكون البيانات غير خطية أو ذات توزيع غير طبيعي.

2. مصفوفة الارتباط:

المفهوم: هي جدول يعرض معامل الارتباط بين كل زوج من المتغيرات في مجموعة بيانات. يساعد في فحص العلاقات بين جميع المتغيرات بشكل جماعي.

التفسير: يمكن أن تكون المصفوفة مفيدة لاكتشاف الأنماط والعلاقات المعقدة بين عدة متغيرات.

الأساليب المتقدمة في التقريب عن البيانات - تحليل البيانات الكبيرة

تقنيات التعامل مع البيانات الكبيرة

1. التخزين والتوزيع:

- أنظمة الملفات الموزعة : مثل (HDFS) Hadoop Distributed File System التي تقسم البيانات الكبيرة إلى أجزاء أصغر وتوزعها عبر عدة خوادم، مما يساعد في إدارة حجم البيانات الكبير وتحسين الأداء.
- التخزين السحابي : خدمات مثل Amazon S3 و Google Cloud Storage توفر تخزينًا مرئيًّا وقابلًا للتوسيع للبيانات الكبيرة.

2. التحليل المتوازي:

- إطار MapReduce: عمل يستخدم لتقسيم معالجة البيانات الكبيرة إلى مهام أصغر يتم تنفيذها بشكل متوازي عبر عدة خوادم. يستخدم لتسرير عملية تحليل البيانات.
- إطار Apache Spark: عمل أكثر تطوراً من MapReduce ، يوفر معالجة البيانات في الذاكرة (in-memory) لتحسين السرعة والأداء.

الأساليب المتقدمة في التنبؤ عن البيانات - تحليل البيانات الكبيرة

3. التنبؤ عن البيانات:

- التجميع: استخدام خوارزميات مثل K-Means لتجميع البيانات الكبيرة إلى مجموعات ذات خصائص مشابهة، مما يساعد في تحليل الأنماط.
- التصنيف: استخدام خوارزميات مثل Random Forest و Decision Trees لتصنيف البيانات الكبيرة بناءً على معايير محددة.

4. تحليل البيانات في الوقت الحقيقي:

- تدفق البيانات: استخدام أدوات مثل Apache Kafka لمعالجة وتحليل البيانات المتداقة في الوقت الحقيقي، مما يساعد في اتخاذ القرارات الفورية.
- تحليل البيانات الحية: أدوات مثل Apache Flink و Apache Storm تستخدم لمعالجة وتحليل تدفق البيانات بشكل فوري وفعال.

أدوات التنقيب عن البيانات: KNIME، Weka، RapidMiner

هو منصة مفتوحة المصدر للتنقيب عن البيانات وتحليلها. تُستخدم بشكل واسع لتصميم، تنفيذ، وتقدير حلول التحليل المعقدة.

الميزات الرئيسية:

- **تطوير النماذج:** يوفر مجموعة متنوعة من خوارزميات التنقيب عن البيانات مثل التجميع، التصنيف، والانحدار.
- **التكامل:** يدعم التكامل مع قواعد بيانات متعددة وأنظمة تحليل بيانات كبيرة.
- **تحليل البيانات:** يتضمن أدوات لتحليل البيانات والنماذج والتصور.
- **التوسيع:** يمكن توسيع وظائف RapidMiner عبر إضافات وبرامج مدمجة.

استخدامات شائعة:

- **التسويق:** تحليل سلوك العملاء وتوقع الاتجاهات المستقبلية. **الصحة:** تطوير نماذج تنبؤية لتشخيص الأمراض.

أدوات التقيب عن البيانات: Weka

Weka .2 هو مجموعة أدوات مفتوحة المصدر مصممة لأغراض التقيب عن البيانات والتحليل. تُستخدم بشكل شائع في الأبحاث الأكاديمية والتجارية.

الميزات الرئيسية:

- **التنوع:** يقدم Weka مجموعة واسعة من الخوارزميات لتصنيف البيانات، التجميع، التقدير، والبحث عن القواعد.
- **الواجهة:** يحتوي على واجهة رسومية (GUI) لسهولة استخدام الأدوات بدون الحاجة للبرمجة.
- **التحليل البياني:** يوفر أدوات لتصور البيانات والنتائج.
- **تحديث البيانات:** يدعم التحديث المستمر للبيانات وتحليلها.

استخدامات شائعة:

- **التحليل التنبؤي:** إنشاء نماذج للتنبؤ بالنتائج المستقبلية بناءً على البيانات التاريخية.
تحليل سلوك المستهلك: تصنيف وتحليل أنماط الشراء.

أدوات التقيب عن البيانات: KNIME

.3 : هو منصة مفتوحة المصدر تُستخدم في التقيب عن البيانات، التحليل، والتصور. يتميز بواجهة رسومية تفاعلية تسهل تصميم سير العمل وتحليل البيانات.

الميزات الرئيسية:

- التكامل: يدعم التكامل مع مجموعة واسعة من الأدوات والأنظمة مثل Hadoop وSpark.
- سير العمل: يتيح إنشاء سير عمل مرن يمكن تعديله بسهولة ليتناسب مع احتياجات التحليل.
- التوسيع: يمكن توسيع وظائف KNIME عبر إضافات ووحدات جديدة.
- تحليل البيانات الكبيرة: يوفر أدوات لمعالجة وتحليل البيانات الكبيرة بفعالية.

استخدامات شائعة:

- التحليل المالي: تحليل البيانات المالية وإعداد التقارير.
- البحث العلمي: تحليل مجموعات بيانات كبيرة من الأبحاث والتجارب.

التطبيقات العملية والدراسات الحالة - تحليل سلوك العملاء

دراسة حالة : Amazon

- الشركة Amazon : أحد أكبر المتاجر الإلكترونية في العالم.
- التحدي : تحسين تجربة العملاء وزيادة المبيعات من خلال فهم سلوك العملاء بشكل أفضل.
- التطبيق : استخدمت Amazon التقىب عن البيانات لتحليل سلوك الشراء الخاص بالعملاء، والتعرف على الأنماط في بيانات المعاملات.
- التقنيات المستخدمة : خوارزميات التجميع (Clustering) والتوصية (Recommendation Systems) لتحليل بيانات التصفح والشراء.
- النتائج:
 - التوصيات الشخصية : أدت التحليلات إلى تحسين نظام التوصية، مما زاد من مبيعات المنتجات ذات الصلة.
 - استهداف الحملات الإعلانية : تحسين استراتيجيات التسويق من خلال استهداف العملاء بناءً على اهتماماتهم وسلوكياتهم.
- الأثر : ساعدت هذه التحليلات Amazon في زيادة رضا العملاء وتحقيق نمو ملحوظ في المبيعات.

التطبيقات العملية والدراسات الحالة - التنبؤ بالأمراض

دراسة حالة : Mayo Clinic

- المنظمة Mayo Clinic :، واحدة من أبرز المؤسسات الطبية في الولايات المتحدة.
- التحدي :تحسین دقة التنبؤ بالأمراض المزمنة مثل السكري وأمراض القلب.
- التطبيق :استخدمت Mayo Clinic تقنيات التقسيب عن البيانات لتحليل سجلات المرضى الطبية، بما في ذلك المعلومات الوراثية والعوامل البيئية.
- التقنيات المستخدمة :خوارزميات التصنيف (Classification) والتجميع (Clustering) لتحديد الأنماط المرتبطة بالأمراض.
- النتائج:
 - التنبؤ المبكر :ساعدت التحليلات في تحسين القدرة على التنبؤ بالإصابة بالأمراض، مما يسمح بالتدخل المبكر.
 - التخصيص الشخصي :تحسين استراتيجيات العلاج من خلال تخصيص العلاج بناءً على البيانات الفردية للمريض.
- الأثر :ساعدت هذه التحليلات في تحسين الرعاية الصحية وتقليل تكاليف العلاج من خلال التدخل المبكر.

التطبيقات العملية والدراسات الحالة - الكشف عن الاحتيال

دراسة حالة : Visa

- الشركة Visa :، إحدى أكبر شركات الدفع الإلكتروني في العالم.
- التحدي : الكشف عن المعاملات الاحتيالية في الوقت الفعلي لحماية العملاء.
- التطبيق : استخدمت Visa تقنيات التنقيب عن البيانات لتحليل معاملات بطاقة الائتمان واكتشاف الأنماط غير الطبيعية التي قد تشير إلى الاحتيال.
- التقنيات المستخدمة : تحليل الارتباطات (Association Analysis) واستخدام خوارزميات الانحراف (Anomaly Detection) لرصد الأنشطة المشبوهة.

النتائج:

- الكشف المبكر : تحسين القدرة على اكتشاف المعاملات الاحتيالية بسرعة أكبر.
 - تحسين الأمان : تقليل الخسائر المالية وحماية العملاء من الاحتيال.
- الأثر : ساعدت هذه التحليلات في تعزيز أمان المعاملات الإلكترونية وحماية العملاء من الاحتيال.

التحديات في التنبؤ عن البيانات – التحديات التقنية

١. البيانات غير النظيفة

أ. **البيانات المفقودة: المشكلة**: البيانات المفقودة هي مشكلة شائعة قد تؤثر على جودة التحليل ونتائجـه. يمكن أن تكون البيانات مفقودة بسبب أخطاء في جمع البيانات أو مشاكل في الإدخال.

- **الحلول: ملء القيم المفقودة**: استخدم تقنيات مثل الإحصاء الوصفي أو النماذج التنبؤية لملء القيم المفقودة.
- **إزالة السجلات غير المكتملة**: إذا كانت النسبة المئوية للبيانات المفقودة صغيرة، يمكن إزالة السجلات غير المكتملة لتجنب تأثيرها على التحليل.
- **الاستدلال الذكي**: استخدام طرق مثل التنبؤ بالاستدلال لاستكمال البيانات المفقودة بناءً على الأنماط في البيانات المتاحة.

ب. **البيانات غير الدقيقة: المشكلة**: قد تحتوي البيانات على أخطاء أو معلومات غير دقيقة بسبب مشاكل في إدخال البيانات أو مصادر غير موثوقة.

- **الحلول:**
- **التنظيف المسبق**: تطبيق تقنيات تنظيف البيانات مثل تصحيح الأخطاء وإزالة القيم الشاذة.
- **التحقق من المصادر**: تأكد من صحة البيانات من خلال التحقق من المصادر ومقارنتها ببيانات أخرى موثوقة.
- **التحقق التكراري**: استخدام طرق التتحقق التكراري لتأكيد دقة البيانات من خلال مقارنة النتائج مع البيانات السابقة أو البيانات المحسوبة.

التحديات في التزريب عن البيانات – التحديات التقنية

2. التحليل في الوقت الفعلي

أ. تحديات معالجة البيانات الكبيرة: المشكلة: تحليل البيانات في الوقت الفعلي يتطلب معالجة كميات ضخمة من البيانات بسرعة وبدون تأخير، وهو ما يمثل تحدياً كبيراً.

• الحلول:

- التقنيات الموزعة: استخدام تقنيات الحوسبة الموزعة مثل Apache Hadoop وApache Spark لمعالجة البيانات بشكل أسرع عبر مجموعة من الخوادم.
- التحليل التدريجي: تطبيق استراتيجيات التحليل التدريجي لمعالجة البيانات على دفعات صغيرة بدلاً من معالجة كل البيانات دفعة واحدة.

ج. التحديات المستمرة: المشكلة: التعامل مع التحديات المستمرة في البيانات يمكن أن يعقد عملية التحليل و يؤدي إلى تأخير في النتائج.

• الحلول:

- التحليل المتكامل: تطبيق نماذج التحليل المتكامل التي تعامل مع البيانات المتداقة وتقوم بتحديث النتائج في الوقت الفعلي.
- التخزين المؤقت: استخدام تقنيات التخزين المؤقت لضمان توافر البيانات الفورية وتحسين أداء التحليل.

التحديات الأخلاقية في التزبيب عن البيانات - الخصوصية

أ. جمع البيانات الشخصية: المشكلة: جمع البيانات الشخصية يمكن أن ينتهك خصوصية الأفراد إذا لم يتم بشكل شفاف أو مع الحصول على موافقة مسبقة.

• **الحلول:**

◦ **موافقة المستهلك:** التأكد من الحصول على موافقة واضحة ومستنيرة من الأفراد قبل جمع بياناتهم.

◦ **إشعار الشفافية:** تقديم إشعارات واضحة حول كيفية جمع البيانات واستخدامها.

◦ **تقليل البيانات:** جمع أقل قدر ممكن من البيانات الازمة لتحقيق الهدف المحدد، وتجنب جمع بيانات غير ضرورية.

ب. استخدام البيانات لأغراض غير معنفة: المشكلة: استخدام البيانات لأغراض لم يُصرح بها من قبل يمكن أن يؤدي إلى إساءة استخدام البيانات وانتهاك حقوق الأفراد.

• **الحلول:**

◦ **الالتزام بالغرض:** استخدام البيانات فقط للأغراض التي تم جمعها من أجلها.

◦ **السياسات الداخلية:** وضع سياسات واضحة تحدد كيفية استخدام البيانات وتتضمن التزام جميع الأطراف بالقوانين والأخلاقيات.

التحديات الأخلاقية في التزوير عن البيانات - الأمان

أ. حماية البيانات الحساسة: المشكلة: تعرض البيانات الحساسة لمخاطر الأمان مثل الاختراقات أو التسريبات يمكن أن يؤدي إلى فقدان الثقة وأضرار كبيرة للأفراد.

• **الحلول:**

- التشفير: استخدام تقنيات التشفير لحماية البيانات المخزنة والمنقولة.

- التحكم بالوصول: تطبيق ضوابط وصول قوية لضمان أن الأفراد المصرح لهم فقط يمكنهم الوصول إلى البيانات الحساسة.

- تحديث الأنظمة: الحفاظ على تحديث الأنظمة وتطبيق تصحيحات الأمان بانتظام لحماية البيانات من التهديدات الجديدة.

ب. إدارة البيانات بعد انتهاء استخدامها: المشكلة: عدم معالجة البيانات بشكل صحيح بعد انتهاء استخدامها يمكن أن يؤدي إلى حفظ البيانات بشكل غير آمن أو استخدامها بشكل غير صحيح.

• **الحلول:**

- الإتلاف الآمن: ضمان إتلاف البيانات بشكل آمن بعد انتهاء فترة استخدامها أو عندما لم تعد ضرورية.

- مراجعات الأمان: إجراء مراجعات دورية لضمان أن البيانات التي لم تعد مطلوبة قد تم إزالتها بشكل آمن.

التحديات الأخلاقية في التنقيب عن البيانات - التلاعب في البيانات

أ. التلاعب في النتائج: المشكلة: قد يحدث تلاعب في البيانات أو النتائج لتوجيه النتائج بطريقة معينة، مما يضلل المستفيدين.

• **الحلول:**

◦ **الشفافية:** ضمان أن تكون العمليات التحليلية شفافة وقابلة للمراجعة.

◦ **المراجعة الداخلية:** تنفيذ مراجعات داخلية للتأكد من نزاهة العمليات التحليلية ونتائجها.

ب. التحييز في النماذج: المشكلة: يمكن أن يتسبب التحييز في النماذج في تمييز غير عادل أو تأثيرات سلبية على مجموعات معينة من الأفراد.

• **الحلول:**

◦ **التدقيق للتأكد من التوازن:** مراجعة النماذج للتأكد من عدم وجود تحيز وعدم عدالة النتائج.

◦ **التحليل العادل:** استخدام تقنيات تحليلات للتأكد من أن النتائج لا تميز ضد أي مجموعة.

دعوة لمزيد من البحث والتعلم

التنقيب عن البيانات هو مجال يتطور بسرعة ويقدم إمكانيات هائلة. أدعو الطلاب إلى:

- استكشاف المزيد : التعمق في دراسة تقنيات وأدوات التنقيب عن البيانات من خلال الدورات التدريبية المتقدمة، ورش العمل، والمقالات العلمية.
 - التطبيق العملي : تطبيق ما تعلموه على مشاريع حقيقة أو بيانات مفتوحة لتعزيز فهمهم وتطوير مهاراتهم العملية.
 - التحديث المستمر : متابعة أحدث الاتجاهات والتقنيات في مجال التنقيب عن البيانات لضمان الاستفادة من التطورات الجديدة وتحقيق التميز في هذا المجال.
- بمواصلة البحث والتعلم، يمكن للطلاب أن يصبحوا خبراء في التنقيب عن البيانات ويساهموا في حل مشاكل معقدة وتحقيق النجاح في مجالاتهم المستقبلية.

أمل أن تكونوا قد حققتم الفائدة