



**الدليل الأساسي للذكاء الاصطناعي القابل  
للتفسير  
(Explainable AI (XAI))**

**إزالة الغموض عن الشفافية والثقة والمساءلة  
في الأنظمة الذكية**

**ياسر إسماعيل**

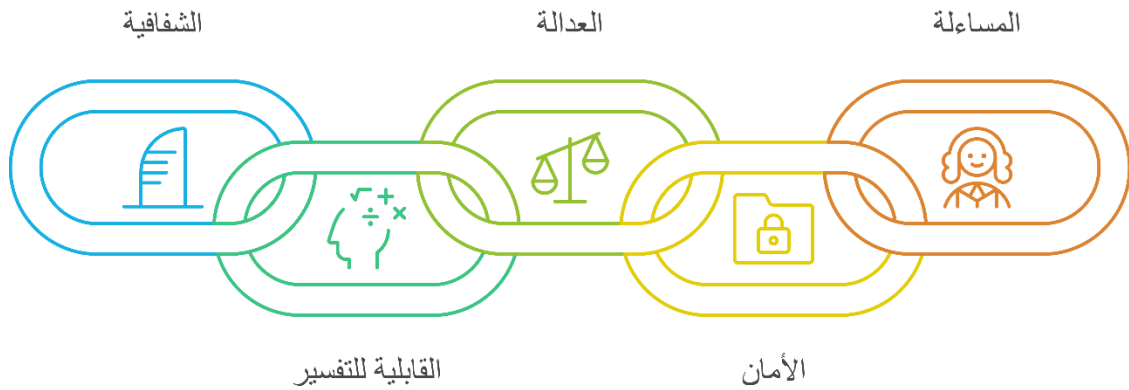


# الدليل الأساسي للذكاء الاصطناعي القابل للتفسير

## (XAI) Explainable AI

إزالة الغموض عن الشفافية والثقة والمساءلة في الأنظمة الذكية

الذكاء الاصطناعي القابل للتفسير



## محتويات الكتاب

مقدمة .....	4
الفصل الأول: مقدمة في الذكاء الاصطناعي القابل للتفسير .....	6
الفصل الثاني: أسس الذكاء الاصطناعي والتعلم الآلي .....	13
الفصل الثالث: التقنيات الأساسية للتفسير .....	19
الفصل الرابع: الأدوات والتقنيات لتطبيق الذكاء الاصطناعي القابل للتفسير .....	25
الفصل الخامس: تطبيقات الذكاء الاصطناعي القابل للتفسير بحسب القطاعات .....	30
الفصل السادس: الإبحار في تحديات الذكاء الاصطناعي القابل للتفسير .....	36
الفصل السابع: الاعتبارات الأخلاقية والأثر المجتمعي .....	39
الفصل الثامن: نظرة إلى الأمام—مستقبل الذكاء الاصطناعي القابل للتفسير .....	44
الفصل التاسع: التدريبات العملية والمشاريع الختامية .....	49
الملحق: موارد إضافية .....	56
قوائم ومراجع إضافية .....	58
عن المؤلف .....	60

## مقدمة

لقد تحرّر الذكاء الاصطناعي من قيود المختبرات البحثية وأصبح حاضرًا في مختلف جوانب حياتنا تقريبًا. فما كان يومًا مفهومًا بعيد المنال أضحت الآن قوة دافعة تقف خلف الأدوات والخدمات التي نعتمد عليها يوميًا. نراه يساعد الأطباء في تحديد الأسباب الجذرية لأعراض المرضى، ويوجّه البنوك في اتخاذ قرارات منح القروض، ويقدم المشورة للقضاة بشأن قضايا الكفالة، بل ويحفّز القادة في المؤسسات على اتخاذ قرارات توظيف أكثر حكمة. بعبارة أخرى، أصبح الذكاء الاصطناعي شريكًا فاعلاً، وإن كان غالبًا غير مرئي، في القرارات الإنسانية المصيرية.

ومع توسّع نطاق الذكاء الاصطناعي، ازدادت الحاجة الملحة لفهم الآليات التي يتخذ من خلالها قراراته. فلم يعد بمقدورنا اعتبار هذه النماذج بمثابة "صناديق سوداء" غامضة، لا سيما مع تعاضد تأثيرها على حياة البشر. كيف توصل خوارزمية ما إلى تشخيص طبي بعينه أو قرار انتمائي محدد؟ لماذا توصي باتباع نهج معين في قضية قانونية؟ إذا أردنا الحفاظ على الثقة، وضمان العدالة، والامتثال للمعايير التنظيمية المتطورة، فنحن بحاجة إلى أجوبة واضحة. وهنا يبرز دور الذكاء الاصطناعي القابل للتفسير. (XAI)

## الذكاء الاصطناعي القابل للتفسير



يرتكز الذكاء الاصطناعي القابل للتفسير على إزالة اللثام عن الآليات الداخلية لأنظمة الذكاء الاصطناعي. فهو يهدف إلى جعل المنطق الكامن وراء التنبؤات والقرارات الصادرة عن هذه الأنظمة مفهومًا وذا مغزى للبشر. إن القدرة على التفسير ليست إنجازًا تقنيًا وحسب، بل ضرورة أخلاقية واجتماعية واقتصادية. فالمؤسسات والعملاء والجهات الرقابية والجمهور العام جميعهم ينشدون الوضوح، ويسعون إلى فهم الكيفية التي تتشكل بها هذه الأنظمة متزايدة التأثير الناتج التي تمس رفاهيتنا واستقرارنا المالي وحياتنا الشخصية

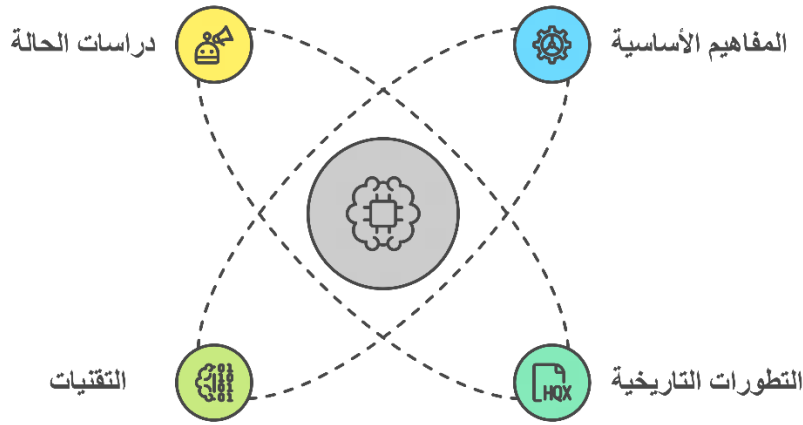
في "الدليل الأساسي للذكاء الاصطناعي القابل للتفسير"، سنسعى لإزالة الغموض عن هذا المجال الأسر. سنستعرض المفاهيم الجوهرية و سنتتبع التطورات التاريخية التي أوصلتنا إلى ما نحن عليه اليوم. سنتعرف على طيف واسع من التقنيات، بدءًا من الأساليب البسيطة التي تسهل تفسير النماذج الصغيرة، وصولاً إلى الأدوات المتقدمة التي تكشف عن

المنطق الكامن في الشبكات العصبية الأكثر تعقيدًا. وسنغوص في دراسات حالات متعددة، مستكشفين كيفية تطبيق الذكاء الاصطناعي القابل للتفسير

والاستفادة منه في المستشفيات والمصارف ومكاتب المحاماة على سبيل المثال

ستدرك أنه لا توجد صيغة موحدة لشرح "جيد" يلائم جميع السياقات. فالصناعات والحالات المختلفة والجمهور المتنوع لديها معاييرها الخاصة. ومع تقدّمنا في الفصول، سنوازن بين الدقة وقابلية التفسير، ونبحث في سبل إيصال الأفكار لغير المتخصصين بطريقة فاعلة، كما سنناقش الضغوط الأخلاقية والتنظيمية المتزايدة التي تطالب بمزيد من الشفافية والمساءلة

### إزالة الغموض عن الذكاء الاصطناعي القابل للتفسير

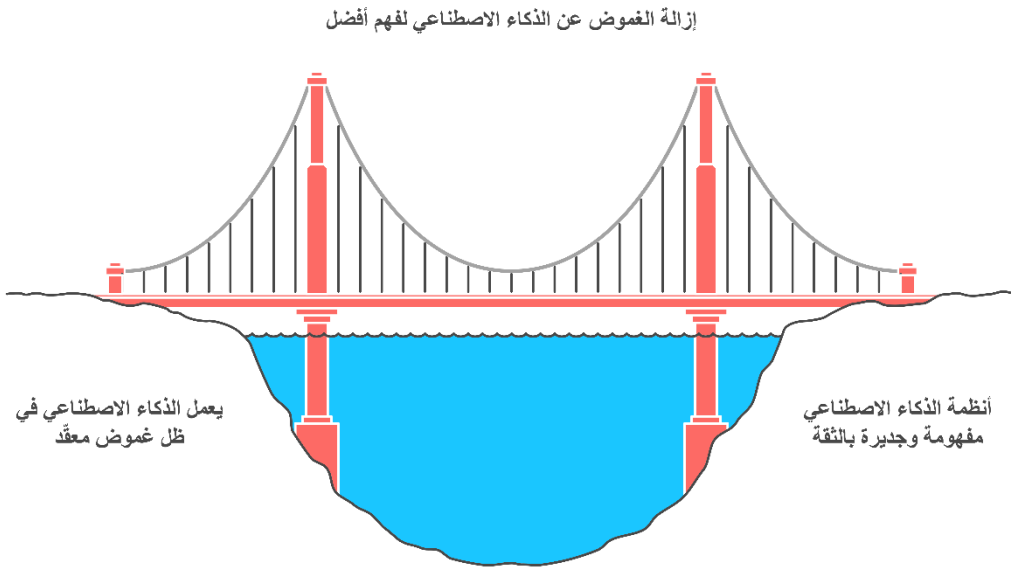


إن هذه الرحلة لا ترمي إلى توسيع معارفك فحسب، بل تهدف أيضًا إلى تزويدك بمهارات عملية. وبحلول الوقت الذي تصل فيه إلى نهاية هذا الدليل، ستكون قادرًا على دمج مفاهيم التفسير في مشاريعك وسياساتك، مدعومًا بأدوات تساعدك على فهم ما تقوم به نماذج الذكاء الاصطناعي القابل للتفسير لديك، ولماذا تتخذ قراراتها بتلك الصورة. ومع استمرار تطوّر عالم الذكاء الاصطناعي القابل للتفسير، فإن دورك سواء كنت تعمل على تطوير هذه الأنظمة، أو الإشراف عليها، أو التعايش معها سيترك بصمة واضحة في صياغة مستقبل يتحلّى بالمسؤولية ويضع الإنسان في صميم اهتماماته

ياسر إسماعيل

## الفصل الأول: مقدمة في الذكاء الاصطناعي القابل للتفسير

أصبح الذكاء الاصطناعي يتغلغل على نحو متزايد في عمليات اتخاذ القرارات التي تصوغ ملامح حياتنا اليومية. فنحن نشهده يؤثر في التشخيصات الطبية، ويرشد القرارات المالية، ويعزز توصيات أنظمة الاقتراح، بل ويساهم في توجيه الأحكام القضائية والسياسات الحكومية. ومع ذلك، ورغم هذا المستوى العالي من التطور، غالبًا ما يعمل الذكاء الاصطناعي خلف حجاب من التعقيد. وقد أثار هذا الوضع تساؤلًا جوهريًا: كيف نضمن أن تبقى أنظمة الذكاء الاصطناعي مفهومة، وجديرة بالثقة، ومتوافقة مع القيم الإنسانية؟ تكمن الإجابة في المجال الآخذ في التوسع، والمعروف باسم الذكاء الاصطناعي القابل للتفسير



### (XAI) تعريف الذكاء الاصطناعي القابل للتفسير

(XAI) يشمل الذكاء الاصطناعي القابل للتفسير

طيفًا واسعًا من التقنيات والمنهجيات والأطر التي تهدف إلى إلقاء الضوء على الآليات الداخلية لنماذج الذكاء الاصطناعي—لا سيما تلك التي يكتنفها الغموض. فالنماذج التقليدية للذكاء الاصطناعي، وبخاصة المعقدة منها كالشبكات العصبية العميقة، قد تحقق دقةً لافتةً، لكنها تعمل أحيانًا كصناديق سوداء. تقوم هذه النماذج بمعالجة كميات هائلة من البيانات المدخلة، ثم تنتج مخرجات—كالتنبؤات، أو التصنيفات، أو التوصيات—دون أن تقدم تفسيرًا واضحًا لكيفية وصولها إلى تلك النتائج

يمكن لنقص الشفافية هذا أن يسبب مشكلات كبيرة في الواقع العملي. تخيل طبيبًا يعتمد على أداة تشخيصية مدعومة بالذكاء الاصطناعي؛ فهو بحاجة إلى فهم المنطق الكامن وراء التشخيص المقترح قبل أن يأخذه بعين الاعتبار عند وضع الخطة

العلاجية. وبالمثل، يحتاج موظف القروض أو المراقب المالي إلى تبرير سبب رفض طلب ائتمان معين. وحتى العملاء يستحقون معرفة الأسباب التي دفعت النموذج لتقديم توصيات بمنتجات محدّدة أو إعلانات مستهدفة. بدون تفسيرات واضحة، تتآكل الثقة وتتضاءل الرغبة في الاستمرار باستخدام الذكاء الاصطناعي

تتصدى منهجيات الذكاء الاصطناعي القابل للتفسير لهذه التحديات عبر تقديم تفسيرات مفهومة للبشر. ويمكن أن تتراوح هذه التفسيرات بين قوائم بسيطة توضح أهمية كل خاصية ("ساهمت الخاصية (أ) بنسبة 30% في القرار، بينما ساهمت الخاصية (ب) بنسبة 20%") وصولاً إلى روايات أكثر تفصيلاً ("يُوصي هذا النموذج بمبلغ قرض أعلى أساساً بفضل تاريخ دخل المستفيد المستقر وانخفاض معدل استخدامه لبطاقة الائتمان خلال العامين الماضيين").

ويمكن لهذه التفسيرات أن تكون:

شاملة (عالمية)	محلية
تلخّص السلوك العام والمنطق الكامن وراء نموذج ما ككل	تشرح التنبؤات الفردية لحالات معينة أو لمجموعات صغيرة منها
<b>(Model-Specific) خاصة بالنموذج</b>	<b>(Model-Agnostic) محايدة للنموذج.</b>
مصمّمة خصيصاً لنوع معيّن من النماذج، مستفيدة من هيكلها الداخلي لتقديم تفسير أدق	تعمل مع أي نوع من النماذج لتوليد التفسيرات، بغض النظر عن الخوارزميات الأساسية

في نهاية المطاف، يهدف يشمل الذكاء الاصطناعي القابل للتفسير إلى ضمان قدرة جميع المعنيين—من الفرق التقنية والخبراء الميدانيين إلى المستهلكين والجهات التنظيمية—على الوصول إلى المعلومات التي يحتاجونها للوثوق بأنظمة الذكاء الاصطناعي وفحصها وإدارتها بمسؤولية

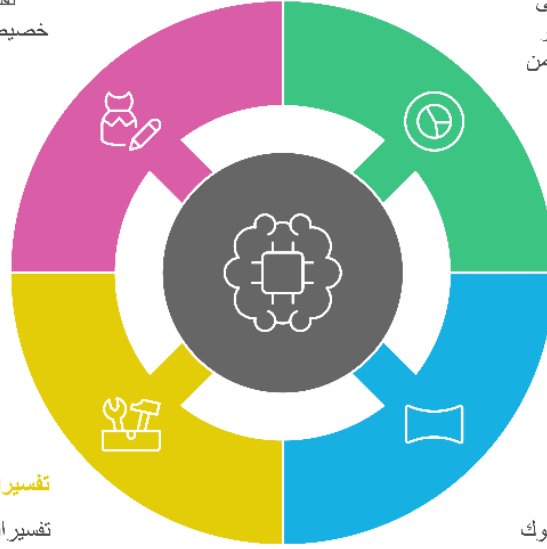
## نتائج الذكاء الاصطناعي القابل للتفسير.

### تفسيرات خاصة بالتمودج

تفسيرات مصممة خصيصاً لأنواع معينة من النماذج.

### تفسيرات محلية

تفسيرات تركز على التنبؤات الفردية أو مجموعات صغيرة من الحالات.



### تفسيرات محايدة للتمودج

تفسيرات تعمل مع أي نوع من نماذج الذكاء الاصطناعي.

### تفسيرات شاملة

تفسيرات تلخص السلوك العام للذكاء الاصطناعي.

## منظور تاريخي

إن السعي نحو الشفافية في مجال الذكاء الاصطناعي ليس وليد اللحظة. ففي البدايات، تم تطوير أنظمة خبيرة قائمة على القواعد (Rule-Based Systems) القواعد

كان المنطق فيها مشفراً بوضوح على هيئة قواعد "إذا-فإن". وكانت هذه الأنظمة قابلة للتفسير بطبيعتها: يمكن تتبع سلسلة القواعد لفهم كل قرار يتخذه النظام. على سبيل المثال، قد يقول نظام طبي يعتمد على القواعد: "إذا كان المريض يعاني من حمى وطفح جلدي، فاقترح اختباراً محدداً." ورغم شفافية هذه الأنظمة، إلا أنها كانت محدودة من حيث النطاق والقدرة على التكيف.

(ML) مع تقدّم الذكاء الاصطناعي، بدأت نماذج التعلّم الآلي

بالتعلّم التلقائي من البيانات، مما قلّل الحاجة إلى كتابة القواعد يدوياً. وبينما حسّنت هذه المقاربة الأداء وقابلية التوسّع بشكل كبير، إلا أنها جعلت عملية الاستدلال أكثر غموضاً. فقدّمت نماذج مثل الأشجار القرارية والنماذج الخطية مستوى ما من القابلية للتفسير، غير أن المجال شهد قفزة هائلة في التعقيد مع ظهور التعلّم العميق في العقد الثاني من القرن الحادي والعشرين



قدّمت نماذج التعلّم العميق—مثل الشبكات العصبية الالتفافية في التعرّف على الصور، والشبكات العودية والشبكات القائمة على المحوّلات في معالجة اللغة—دقة غير مسبوقّة. ومع ذلك، فإن تمثيلاتها الداخلية، التي غالبًا ما تتكون من ملايين أو مليارات المعاملات، ظلّت مبهمّة. ورغم تفوّق هذه النماذج في مهام ككشف الأورام في الصور الطبية أو الترجمة بين اللغات، إلا أن فهم سبب قيامها بذلك ظلّ أمرًا مرًا مرًا.

برزت أزمة “قابلية التفسير” هذه بشكل أوضح في السيناريوهات الحساسة. كيف يمكن لطبيب أن يثق في تشخيص اقترحه نموذج لا يمكن تفسير منطقته؟ كيف يمكن لمصرف الاعتماد على الموافقات الائتمانية المدفوعة بالذكاء الاصطناعي من دون التأكّد من عدالتها وامتثالها التنظيمي؟ وكيف يمكن لقاضي أو واضع سياسات أن يشعر بالراحة في اعتماد تقييمات خوارزمية للمخاطر دون قابلية للتفسير؟

استجابة لذلك، بدأ الباحثون في تطوير طرق لتوضيح هذه النماذج “الصندوقية السوداء”. وقدّمت التقنيات المبكرة مثل (التفسيرات النموذجية المحلية المحايدة) أسلوبًا لتقريب حدود قرار النموذج المعقّد حول حالة مفردة باستخدام LIME نماذج أبسط وأكثر قابلية للتفسير. سمح هذا الأسلوب للممارسين بإلقاء نظرة على سلوك النموذج على نطاق محلي، حتى وإن ظلّ غامضًا على المستوى العام.

(تفسيرات شابلي الإضافية) التي استندت إلى مفاهيم من نظرية الألعاب لتوزيع SHAP لاحقًا، جاءت ابتكارات مثل ، مستفيدًا من XAI المساهمات بإنصاف واتساق على كل خاصية من خصائص النموذج. وبمرور الوقت، توسّع مجال مجالًا متعدّد التخصصات XAI أفكار في الإحصاء وتفاعل الإنسان والحاسوب وعلم النفس والأخلاقيات. واليوم، يشكّل يتطوّر باستمرار لتلبية الطلب المتزايد على الشفافية في أنظمة الذكاء الاصطناعي المتزايدة التعقيد

### تطور شفافية الذكاء الاصطناعي



## (XAI) أثر الذكاء الاصطناعي القابل للتفسير

لا تقتصر أهمية الذكاء الاصطناعي القابل للتفسير على الجوانب التقنية فحسب، بل تمتد لتشمل أبعادًا تتعلق بالثقة، والامتثال التنظيمي، والإنصاف، وتحسين النماذج ذاتها

### 1. الثقة والاعتماد

لكي يصبح الذكاء الاصطناعي جزءًا لا يتجزأ من عمليات اتخاذ القرارات الحساسة، لا بد أن يتمتع المعنيون بالثقة في مخرجاته. تخيل مستشفى يطبق أداة ذكاء اصطناعي للتنبؤ بمخاطر وفاة المرضى. يرغب الأطباء والمرضى وعائلات المرضى في فهم سبب ارتفاع الخطر أو انخفاضه. إذا استند التفسير إلى عوامل طبية موثوقة، تزداد الثقة بالنظام. أما في حال غياب التفسيرات، فقد يسود الشك وتتعرض فرص تبني هذه التقنيات

#### مثال

يشعر فريق طبي بالارتياح عند الاعتماد على نموذج تنبؤ بالإنتان (تسمم الدم) مدعوم بالذكاء الاصطناعي عندما يتبين لهم أن منطقه التحليلي يستند إلى علامات إكلينيكية راسخة—مثل ارتفاع عدد كريات الدم البيضاء واضطرابات في معدل التنفس. بفضل هذه الشفافية، يثق الأطباء بالتحذيرات الصادرة عن النظام ويتدخلون مبكرًا، ما قد ينقذ حياة المرضى

### 2. الامتثال التنظيمي وإدارة المخاطر

"تبدي الهيئات التنظيمية حول العالم اهتمامًا متزايدًا بالمساءلة الخوارزمية. وقد فُسر "النظام العام لحماية البيانات في الاتحاد الأوروبي على أنه يتضمن "حق الحصول على تفسير" للقرارات الآلية، الأمر الذي يدفع (GDPR) المنظمات لجعل نماذج الذكاء الاصطناعي أكثر شفافية. وفي الولايات المتحدة، يدفع قانون تكافؤ الفرص الائتمانية واقتراحات تنظيمية أخرى باتجاه مزيد من التفسيرية في قرارات الإقراض

#### مثال

يتعين على مؤسسة مالية إثبات أن أداة تقييم الائتمان المدعومة بالذكاء الاصطناعي لا تميز ضد مجموعات لتفسير قرارات الائتمان، يمكن للمصرف إثبات للجهات التنظيمية أن تقييم العملاء SHAP معينة. باستخدام قيم يستند إلى معايير موضوعية ومسموح بها قانونيًا—مثل السجل الائتماني واستقرار الدخل—مما يقلل من مخاطر الغرامات أو الدعاوى القضائية

### 3. الإنصاف والأخلاقيات

تكتسب نماذج الذكاء الاصطناعي معارفها من بيانات تاريخية، مما قد يرسخ التحيزات القائمة ويعزز أوجه اللامساواة. وفي غياب التفسير، يصعب الكشف عما إذا كانت قرارات النموذج تضر بفئات معينة. تتيح التفسيرات مراجعة المنطق الداخلي للنماذج، ورصد أماكن وجود التحيزات الكامنة

#### مثال

تخيل نظام توظيف مدعوم بالذكاء الاصطناعي يبدو أنه يفضل مرشحين معينين. من خلال دراسة تفسيرات

النموذج، يكتشف خبراء الموارد البشرية أنه يولي وزناً كبيراً لكلمات مفتاحية تاريخياً مرتبطة بالمتقدمين الذكور. عند إدراك هذا التحيز الخفي، يعيدون تدريب النموذج على بيانات أكثر تنوعاً، ويحذفون الخصائص المنحازة، مما يؤدي في النهاية إلى تحسين العدالة

#### 4. تحسين أداء النماذج والتصحيح

لا تقتصر فوائد التفسيرات على المستخدمين النهائيين والجهات التنظيمية فحسب، بل تمتد أيضاً لفرق علوم البيانات والمطورين الذين يرغبون في تحسين نماذجهم. إذ يمكن للإطلاع على منطق النموذج أن يبرز الخصائص التي تتسبب في أخطاء مستمرة أو نتائج غير متوقعة، ما يوجّه جهود التطوير والتحسين

#### مثال

يطور فريق من المهندسين نموذجاً للتنبؤ بالصيانة الوقائية في بيئة صناعية. تكشف التفسيرات أن النموذج يولي اهتماماً مفرطاً بحساس كثير التشويش، مما ينتج عنه إنذارات خاطئة. عبر تعديل إجراءات معالجة البيانات أو حذف هذا الحساس من مجموعة الخصائص، يتمكّن الفريق من تحسين دقة النموذج واستقراره

#### ما وراء الحلول التقنية: ضرورة مجتمعية

على كونه مجموعة من الأدوات التقنية، بل يمثل تغييراً في طريقة تصميم الأنظمة الذكية ونشرها XAI لا يقتصر الـ وحوكمتها. ومع اندماج الذكاء الاصطناعي في النسيج الاجتماعي والاقتصادي والسياسي، تضمن القابلية للتفسير أن القرارات الآلية لا تتخذ في فراغ أخلاقي. فهي تتيح آلية لمساءلة استخدام الذكاء الاصطناعي، وتمكّن الأفراد من فهم النتائج المؤثرة على حياتهم، وتشجّع حواراً مجتمعياً أوسع حول متى ينبغي الوثوق بهذه الأنظمة وكيف

على سبيل المثال، تأمل في الاعتدال الآلي للمحتوى على منصات التواصل الاجتماعي. بدون تفسيرات، قد يشعر المستخدمون بأن خوارزمية غامضة قد حذفت منشوراتهم تعسفاً. أما النماذج القابلة للتفسير فيمكنها توضيح أن محتوى معيناً أزيل لانتهاكه إرشادات المجتمع المحددة، مما يعزّز فهم المستخدمين وتقبّلهم لهذه القرارات الآلية

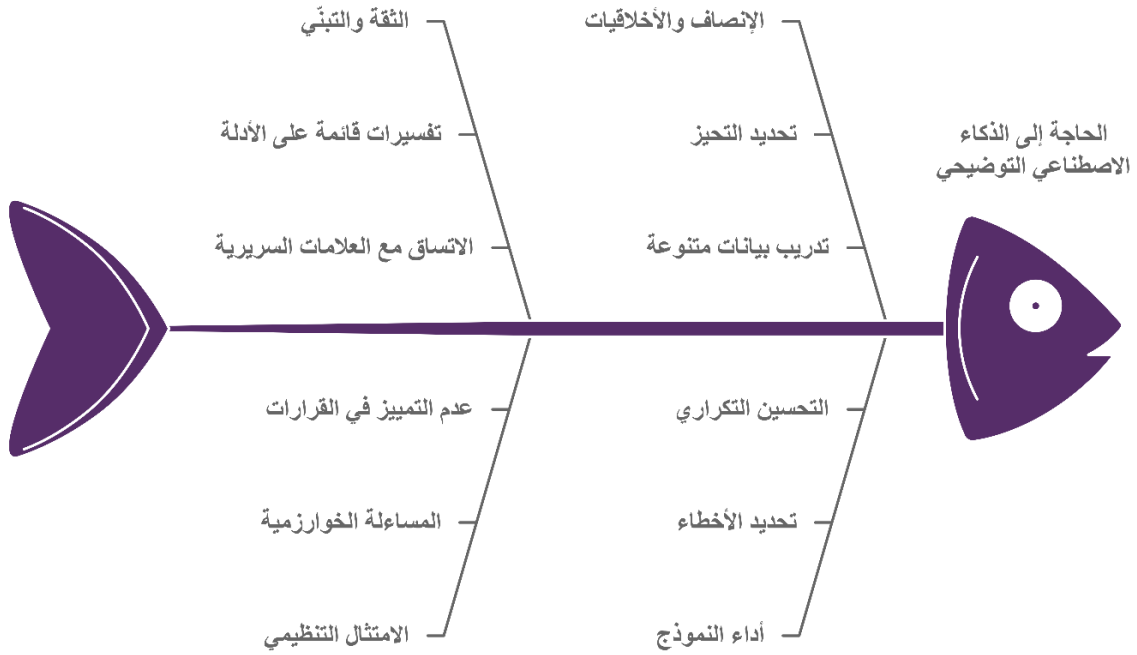
#### نظرة إلى الأمام

يمهّد هذا الفصل الطريق لما سيأتي. سنستعرض في الفصول اللاحقة الأسس الخاصة بالذكاء الاصطناعي والتعلّم الآلي، وندرس التقنيات الرئيسية للتفسير—بدءاً من النماذج القابلة للتفسير بطبيعتها ووصولاً إلى المناهج المعقّدة المحايدة ممارسة عملية. وسنستعرض تطبيقات متخصصة في XAI للنماذج—كما سننتعمق في الأدوات والتقنيات التي تجعل الـ قطاعات مختلفة، من الرعاية الصحية والتمويل إلى النظم القانونية والتصنيع، لتعرّف كيف يُحدث التفسير نقلة نوعية في هذه المجالات

سنناقش أيضاً التحديات المصاحبة للتفسير، بما في ذلك الموازنة بين تعقيد النموذج وشفافيته، والتكاليف الحسابية لإنتاج التفسيرات، والحاجة إلى تكيف التفسيرات مع شرائح مختلفة من الجمهور. وستتخلل مناقشتنا اعتبارات أخلاقية مثل الكشف عن التحيزات واحترام الخصوصية، إضافة إلى استعراض المشهد التنظيمي والسياسات المتطورة في هذا الميدان.

، ستكون أكثر قدرة على ابتكار وتقييم وحوكمة XAI من خلال فهم الجذور التاريخية والمناهج الحالية والأثر الواقعي للـ حلول ذكاء اصطناعي ليست قوية فحسب، بل مفهومة ومسؤولة ومتوافقة مع القيم الإنسانية

### (XAI) التأثير متعدد الجوانب للذكاء الاصطناعي التوضيحي



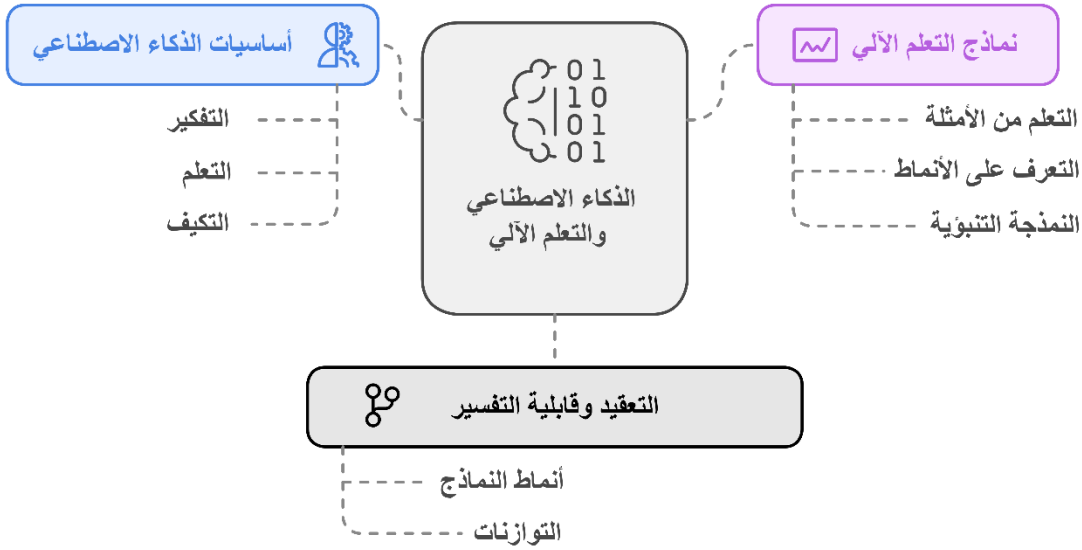
## الفصل الثاني: أسس الذكاء الاصطناعي والتعلم الآلي

لفهم كيفية جعل أنظمة الذكاء الاصطناعي أكثر شفافية وجديرة بالثقة، يجب أولاً استيعاب المبادئ الأساسية التي تحكمها. ، متناولاً الأنماط الأساسية للتعلم (ML) والتعلم الآلي (AI) يستكشف هذا الفصل المفاهيم الجوهرية للذكاء الاصطناعي ونماذج الذكاء الاصطناعي النموذجية والمقايضات بين التعقيد وقابلية التفسير. بوضع هذه الأسس، نمهّد الطريق لنقاشات أكثر عمقاً حول قابلية التفسير

### أساسيات الذكاء الاصطناعي والتعلم الآلي 2.1

يطمح الذكاء الاصطناعي إلى محاكاة القدرات المعرفية المرتبطة عادةً بالذكاء البشري—كالاستدلال والتعلم والتكيف مع الظروف الجديدة. ويُعدّ التعلم الآلي أحد أهم فروع الذكاء الاصطناعي، إذ يركّز على تمكين الخوارزميات من التعلّم مباشرة من البيانات، دون الحاجة إلى قواعد يدوية

بدلاً من تحديد قواعد صريحة، نقدّم للنماذج أمثلة. ومع مرور الوقت، تتوصّل هذه الأنظمة لاكتشاف الأنماط والعلاقات الكامنة، ثم تطبّق ما تعلّمته لإصدار تنبؤات حول بيانات جديدة. وقد دفع هذا النهج عجلة التقدّم في مجالات شتى، بدءاً من التوصيات الصحية المخصصة ووصولاً إلى السيارات ذاتية القيادة في المدن المزدهمة



## المناهج الرئيسية للتعلم

### • التعلم الموجّه (Supervised Learning):

في هذا النهج، نقدّم للنموذج أمثلة تحتوي على المدخلات والمخرجات الصحيحة (التصنيفات أو القيم). يتعلّم النموذج ربط المدخلات بالمخرجات، كتنبؤ أسعار المنازل أو تصنيف البريد الإلكتروني إلى بريد عادي أو مزعج.

#### مثال:

لنفترض أننا جمعنا بيانات عقارات معروفة الأسعار. بتدريب نموذج موجّه على مساحة المنزل وعدد الغرف والقرب من المدارس، يمكن للنموذج تقدير سعر بيع عقار جديد بدون معلومة مسبقة عن سعره.

### • التعلم غير الموجّه (Unsupervised Learning):

يتعامل التعلم غير الموجّه مع بيانات غير معنونة، حيث يحاول النموذج الكشف عن البنى الكامنة والأنماط. في كشف مجموعات طبيعية داخل البيانات (Clustering) الخفية. تساعد خوارزميات التجميع

#### مثال:

قد يكتشف تاجر تجزئة، عبر التحليل التجميعي، أن بعض العملاء يفضلون السلع المخفضة، بينما يفضل آخرون المنتجات الفاخرة باستمرار. تساعد هذه الرؤى في تحسين استراتيجيات التسويق دون وجود فئات مسبقة للتصنيف.

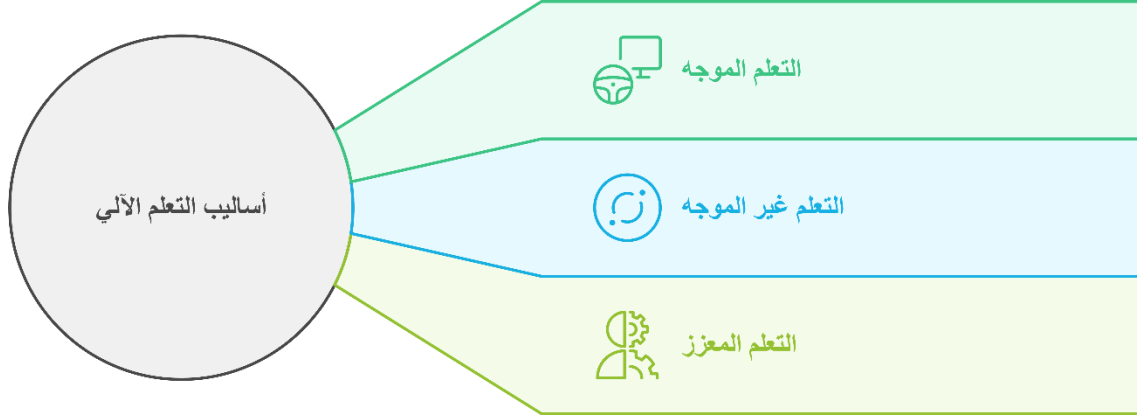
### • التعلم بالتعزيز (Reinforcement Learning):

مع بيئته ويتلقى مكافآت أو عقوبات بناءً على أفعاله. وبمرور الوقت، (Agent) في هذا النهج، يتفاعل وكيل يطور الوكيل سياسة عمل تهدف لتعظيم المكافآت التراكمية

#### مثال:

مكنسة روبوتية تتعلم تدريجياً المسار الأمثل للتنظيف من خلال التجربة والخطأ، مما يعزّز كفاءتها في استهلاك الطاقة وتغطية المساحات دون تعليمات بشرية مباشرة.

تتراوح نماذج التعلم الآلي من نماذج خطية بسيطة إلى نماذج عميقة بالغة التعقيد. يساعد فهم مكان نموذج معين على هذا الطيف في اختيار الإستراتيجيات المناسبة لجعله مفهومًا



## فئات النماذج وقابلية التفسير 2.2

:تختلف عائلات النماذج من حيث مدى سهولة تفسيرها

- (مثل الانحدار الخطي والانحدار اللوجستي) **النماذج الخطية**

تعتمد هذه النماذج على جمع موزون للخصائص المدخلة. توضح كل معاملات النموذج مدى تأثير خاصية معينة على المخرجات، مما يسهل تفسيرها.

**مثال:**

في نموذج تنبؤ بدرجات الطلاب، قد يظهر أن كل ساعة دراسة إضافية تضيف نقطتين للنتيجة المتوقعة، في حين أن تفويت حصة دراسية يقللها بنقطة واحدة، مما يسهل على المعلمين والطلاب فهم المنطق.

- **الأشجار القرارية والأنظمة القائمة على القواعد**

تشبه الأشجار القرارية قوائم قرارات شرطية، حيث يمكن تتبع المسار من الجذر إلى الورقة لفهم سلسلة الشروط التي أدت إلى التنبؤ.

**مثال:**

شجرة قرارية في المجال الطبي قد تقول: "إذا تجاوزت حمى المريض درجة معينة، افحص خلايا الدم البيضاء. إذا تجاوز عددها الحد، أوص باختبار محدد." يمكن للأطباء التحقق من منطق هذه الخطوات

- والانحدار المعزز (Random Forest) مثل الغابة العشوائية (**Ensemble Methods**) **طرق التجميع**

(Gradient Boosted Trees):

تدمج هذه الطرق نماذج أبسط لتحقيق دقة أعلى. ورغم متانتها، إلا أنها أقل وضوحًا من الشجرة القرارية المفردة. ومع ذلك، فإن تقنيات مثل ترتيب أهمية الخصائص قد تساعد في إظهار العوامل الأكثر تأثيرًا

**مثال:**

يستخدم مصرف نموذج "غابة عشوائية" للتنبؤ باتجاهات الأسهم. رغم صعوبة تفسير كل شجرة، إلا أن فحص أهمية الخصائص الإجمالية يبين أن المؤشرات الاقتصادية العالمية والتقارير المالية الأخيرة تشكل العوامل الأهم

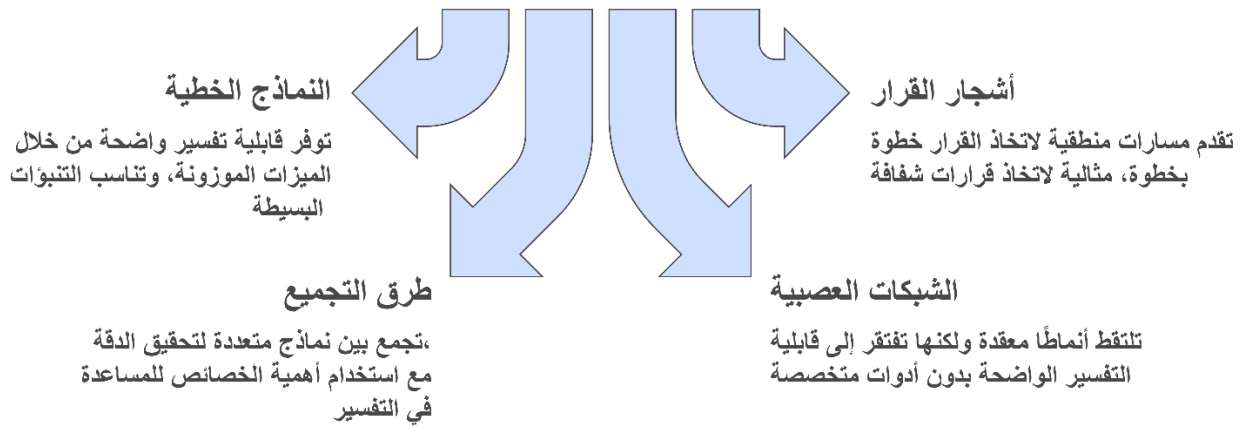
#### • الشبكات العصبية

تتألف الشبكات العميقة من طبقات عديدة، وتبرع في اكتشاف أنماط معقدة—مثل خصائص الصور الدقيقة أو الفروق اللغوية الخفية. إلا أن تمثيلاتها الداخلية غالبًا ما تكون غامضة. في غياب أدوات تفسير متخصصة، يصعب فهم سبب تنشيط عصبون معين أو كيفية تمييز الشبكة بين الفئات

**مثال:**

يبرع نموذج عصبي عميق في التعرف على أنواع الطيور من الصور، لكن تحديد المجموعة الدقيقة من سمات البكسل التي يستند إليها لاتخاذ القرار يظل صعبًا بدون أساليب تفسيرية إضافية

### أي فئة من النماذج تختار بناءً على قابلية التفسير والتطبيق؟





### أهمية الشفافية 2.3

تكتسب الشفافية أهمية قصوى مع تسليم نماذج الذكاء الاصطناعي مهام ذات تبعات خطيرة على الصحة والمعيشة والحريات الشخصية. خذ هذه الأبعاد بعين الاعتبار:

#### • تقييم الموثوقية:

يحتاج الأطباء، عند الاعتماد على أداة تشخيصية بالذكاء الاصطناعي، إلى الثقة بأن المنطق الذي يقود النموذج يتوافق مع المعرفة الطبية المقبولة. بدون تفسير، يبقى الأمر مجرد حدس.

#### • ضمان الإنصاف:

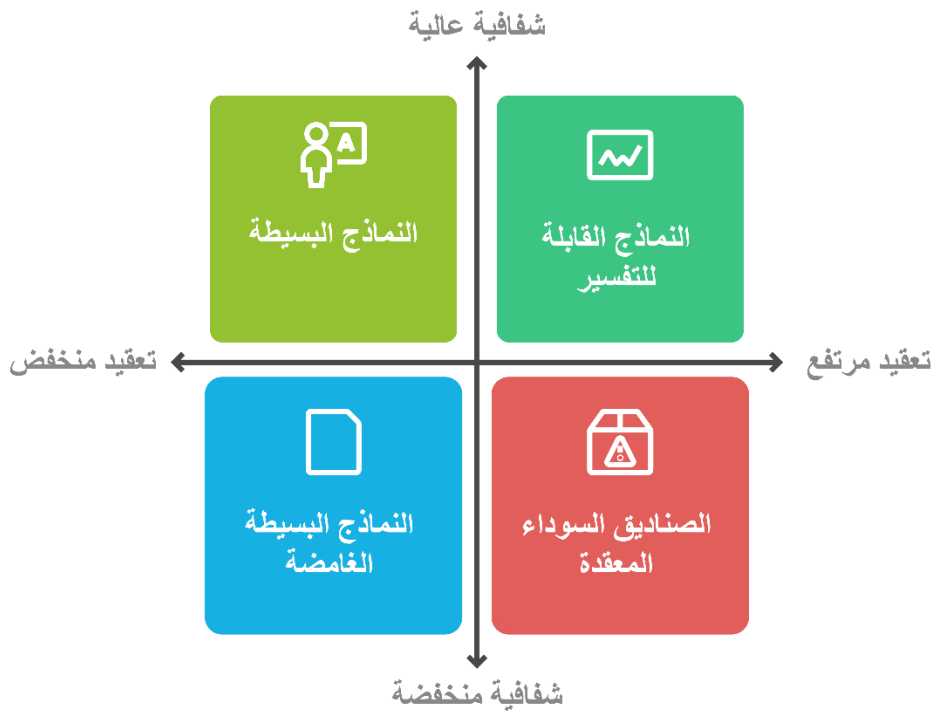
قد يؤدي نموذج توظيف غامض إلى تفضيل فئات معينة، مما يعيد إنتاج التحيزات التاريخية. بدون وضوح، يصعب على المؤسسات اكتشاف هذه الاختلالات ومعالجتها، مما يعرضها للمخاطر القانونية والأخلاقية.

#### • تعزيز الثقة العامة:

يسعى الأفراد الذين يخضعون لتقييم انتمائي آلي، أو تقييم أداء، إلى التأكد من أن القرارات تستند إلى منطق "عقلاني وغير متحيز". تضفي الشفافية مصداقية على العملية وتبديد الشكوك حول "الصندوق الأسود".

في سبيل الشفافية، يمكن اختيار نماذج يسهل تفسيرها بطبيعتها، أو استخدام استراتيجيات تفسيرية متخصصة لفهم نماذج أكثر تعقيداً. ويبقى اختيار النهج المناسب رهين السياق والمتطلبات التنظيمية وتوقعات أصحاب المصلحة.

### الموازنة بين الشفافية والتعقيد في نماذج الذكاء الاصطناعي



## حالة توضيحية: الشفافية في تقييم الائتمان 2.4

تخيل مؤسسة مالية تستخدم شبكة عصبية لتحديد أهلية القروض. رغم دقة النموذج، يظل المنطق غامضًا. بدون تفسير، سيفقد العملاء ثقتهم، وقد يشكون في تعسف القرارات. كما قد يطالب المشرعون بتفسيرات تبرهن عدم وجود تمييز غير مشروع

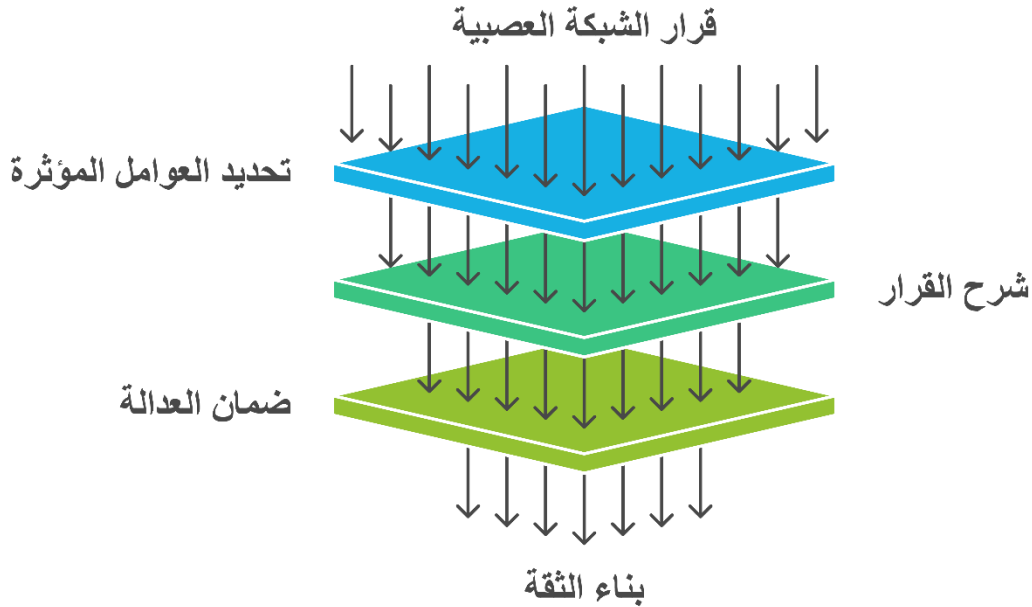
SHAP أو LIME باستخدام تقنيات مثل

، تستطيع المؤسسة إظهار الخصائص الأكثر تأثيرًا في القرار—مثل السجل الوظيفي غير المنتظم أو التخلف المتكرر عن سداد القروض. يتيح هذا الوضوح لموظف القروض شرح القرار بهدوء، مؤكّدًا التزام البنك بمعايير إقراض عادلة، مما يطمئن العملاء والجهات التنظيمية

### الخلاصة

استعرض هذا الفصل المبادئ الأساسية في الذكاء الاصطناعي والتعلم الآلي، مسلطًا الضوء على تنوع أنماط التعلم، ودرجات قابلية تفسير النماذج، وأهمية الشفافية. وباستخدام هذه الأسس، يمكننا الآن التعمق في جوهر التفسيرية—دراسة الأدوات والتقنيات وأفضل الممارسات التي تضمن فهم النماذج الأكثر تعقيدًا من قبل من يعتمدون على مخرجاتها.

## تعزيز الشفافية في تقييم الائتمان



## الفصل الثالث: التقنيات الأساسية للتفسير

لإيصال أنظمة الذكاء الاصطناعي إلى مستوى يمكن للبشر فهمه، لا بد من اعتماد أساليب تكشف عن أسس عمل هذه النماذج. تركز بعض الاستراتيجيات على بناء الوضوح في تصميم النموذج منذ البداية، بينما تعمل أخرى بأثر رجعي لتفسير المنطق في النماذج المعقدة التي تم تدريبها مسبقاً. يستعرض هذا الفصل التقنيات الرئيسية المعززة للتفسير، إلى جانب الطرق المتبعة لتقييم جودة وفائدة التفسيرات الناتجة

### 3.1 نماذج قابلة للتفسير بطبيعتها

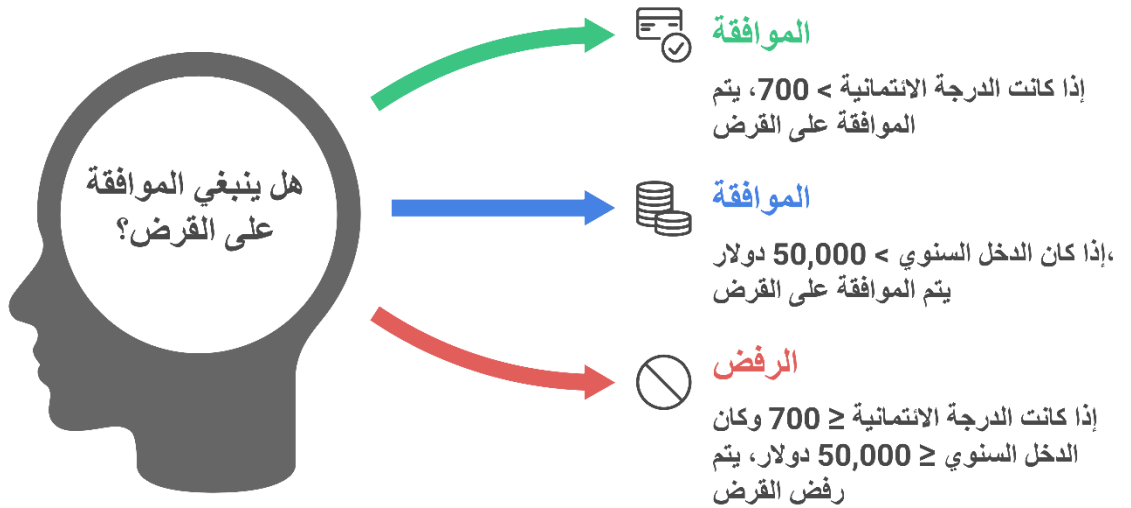
يتمثل أحد الطرق المباشرة نحو التفسير في اختيار خوارزميات وعائلات نماذج شفافة من الأساس. ورغم أن هذه النماذج قد لا تضاهي أحياناً قوة النماذج الأكثر تعقيداً، فإن وضوحها الفوري يجعلها لا تُقدَّر بثمن في السيناريوهات التي تتقدّم فيها الثقة والمساءلة على زيادة طفيفة في الدقة

### الأشجار القرارية

تعمل الشجرة القرارية كسلسلة منظمة من الأسئلة الشرطية، بحيث يقودنا تتبّع المسار من الجذر إلى الورقة لفهم المنطق الذي قاد إلى التنبؤ النهائي. ورغم أن كبر حجم الشجرة قد يقلل من وضوحها، إلا أن تقنيات التقليم وأدوات التصوّر تساعد في الحفاظ على قابليتها للفهم

### مثال

في عملية الموافقة على القروض، قد تقول الشجرة القرارية: "إذا تجاوزت درجة الائتمان 700، فوافق. وإلا، إذا تجاوز الدخل السنوي 50 ألف دولار، فوافق. وإلا، فرفض." هذه السردية مفهومة لموظفي القروض والمدققين



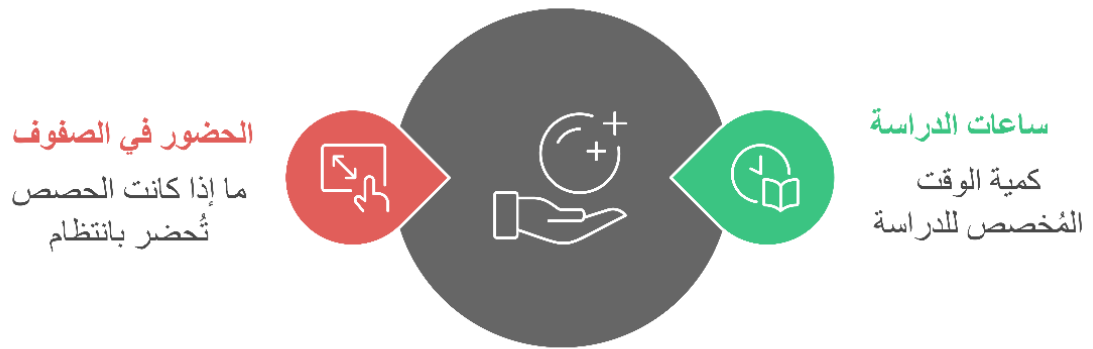
## الانحدار الخطي واللوجستي

تعتبر هذه النماذج عن التنبؤات على شكل مجموع موزون للخصائص. تشير كل معلمة (معامل) بوضوح إلى مدى تأثير خاصية ما على المخرجات، مما يجعل التفسير مباشرًا وسلسًا

### مثال

نموذج انحدار خطي يتنبأ بنتائج الطلاب قد يظهر أن كل ساعة دراسة إضافية ترفع الدرجة المتوقعة بنقطتين، فيما يقلل تفويت حصة دراسية النتيجة بنقطة واحدة. هذه العلاقة الواضحة تساعد المعلمين والطلاب على وضع خطط فعالة

## العوامل المؤثرة على التنبؤات



## (GAMs) النماذج الإضافية المعممة

توازن بين المرونة وقابلية التفسير. إذ تصف تأثير كل خاصية عبر دالة سلسلة، ثم تجمع هذه التأثيرات لإنتاج التنبؤ النهائي. تسمح هذه البنية غير الخطية، ولكن المضافة، بفهم إسهام كل خاصية بشكل مستقل

### مثال

للتنبؤ بخطر إعادة إدخال المرضى للمستشفى قد يعرض تأثير كل خاصية على شكل منحني GAM نموذج سهل القراءة. يمكن للطبيب رؤية أن ارتفاع مستوى الكوليسترول يرتبط بزيادة منتظمة في خطر العودة إلى المستشفى، مما يقدم أساسًا منطقيًا لإجراء فحوصات إضافية

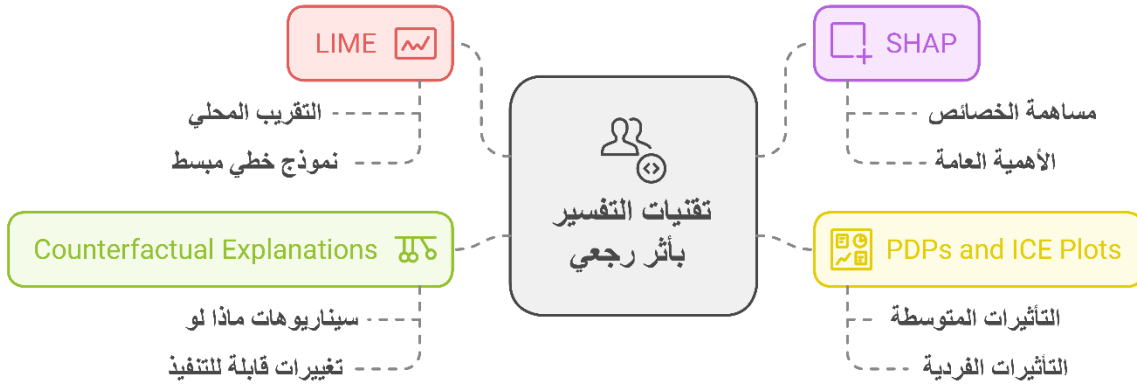
## المساهمات في التنبؤات باستخدام النماذج الإضافية المعممة (GAMs)



في القطاعات المنظمة مثل التمويل والرعاية الصحية، حيث تكتسب القدرة على تبرير القرارات أهمية قصوى، تعزز هذه النماذج ثقة المستخدمين. إلا أنه يجدر التنويه إلى وجود مفاضلة: فقد لا تصل دقة النماذج البسيطة إلى مستوى النماذج الأحدث والأكثر تعقيداً، الأمر الذي يطرح تساؤلات حول أولوياتنا: أيهما أولى، الفهم أم الأداء الأمثل؟

### 3.2 (Post-Hoc) تقنيات التفسير بأثر رجعي

عندما تكون النماذج الشفافة بطبيعتها غير مناسبة أو غير كافية من حيث الدقة، تتدخل تقنيات التفسير بأثر رجعي. هذه التقنيات توضح المنطق الكامن خلف النماذج المعقّدة—مثل الشبكات العصبية العميقة—دون تعديل بنيتها الداخلية. وبذلك تقدّم نافذة على تفكير النموذج لم تكن متوفرة في التصميم الأصلي



#### (التفسيرات المحلية المحايدة للنموذج) LIME

حدود قرار النموذج المعقّد حول حالة معينة بنموذج أبسط قابل للتفسير، غالبًا ما يكون نموذجًا خطيًا LIME تقرّب بسيطًا. ومن خلال فحص معاملات هذا النموذج البسيط، يمكننا استنتاج الخصائص التي أثّرت على التنبؤ المحدد

#### مثال

إذا رفضت شبكة عصبية طلب قرض، قد تُظهر

أن القرار استند محليًا إلى ارتفاع نسبة استخدام البطاقة الائتمانية والتاريخ الوظيفي المتقلّب للمتقدّم LIME

#### (تفسيرات شاب الإضافية) SHAP

تعتمد قيم SHAP

على مفاهيم من نظرية الألعاب، حيث تعامل الخصائص كلاعبين يساهمون في مخرجات النموذج. من خلال النظر في جميع مجموعات الخصائص الممكنة، تحدّد

حصة عادلة من التأثير لكل خاصية على التنبؤ النهائي. توفر هذه الإطارات تفسيرات محلية وعالمية SHAP

#### مثال

لنموذج تشخيص طبي أن إحدى الأعراض ترفع خطر المرض باستمرار عبر العديد SHAP قد يكشف تحليل من المرضى، بينما يقلل عامل آخر من هذا الخطر. هذا النمط يطمئن الأطباء بأن منطق النموذج يتوافق مع المعارف الطبية المعروفة

## (ICE) ومخططات التوقع الشرطي الفردي (PDP) مخططات الاعتماد الجزئي

PDP تعرض الـ

- فنتيح ICE كيف يؤثر تغيير خاصية واحدة في متوسط تنبؤات النموذج، مع تثبيت الخصائص الأخرى. أما مزيداً من التفصيل بعرض تأثير الخاصية على كل حالة فردية، ما يكشف عن اختلافات قد تختفي في المتوسط

### مثال

علاقة إيجابية بين زيادة الإنفاق على التسويق وزيادة المبيعات PDP في نموذج تسعير المنتجات، قد تُظهر الـ عن أن بعض فئات المنتجات لا تتبع هذا النمط، ما يساعد في وضع استراتيجيات ICE المتوقعة، بينما تكشف الـ تسويقية أكثر استهدافاً

### • (Counterfactual Explanations) التفسيرات المضادة

تسأل التفسيرات المضادة: "ما الذي ينبغي تغييره للحصول على نتيجة مختلفة؟" تساعد هذه المنهجية المستخدمين على فهم ما يلزم تعديله في مدخلاتهم للحصول على مخرجات مختلفة، وتوفّر تفسيرات قابلة للتنفيذ

### مثال:

إذا رُفض طلب قرض، قد تشير التفسيرات المضادة إلى أن زيادة الدخل السنوي بمقدار 5,000 دولار أو تخفيض نسبة استخدام البطاقة الائتمانية بمقدار 10% كان سيؤدي إلى الموافقة، مما يجعل التفسير عملياً

## تقييم جودة التفسيرات 3.3

ليست كل التفسيرات مفيدة أو دقيقة أو مفهومة. ومع تنوع أساليب التفسير، ينبغي وضع معايير لتقييم فعاليتها

### • (Fidelity) المصداقية:

هل يعكس التفسير منطق النموذج الحقيقي؟ إذا كانت التفسيرات مبسطة أكثر من اللازم أو تُضلل بشأن الآليات الأساسية، فإنها تفوّض الثقة

### • الاستقرار والمتانة:

إذا أدت تغييرات طفيفة في المدخلات إلى تفسيرات مختلفة تماماً، يشكك المستخدمون في موثوقية النموذج. فالتفسيرات المستقرة تعزز الثقة

### • قابلية الفهم:

لا جدوى من تفسير محكم إذا كان المستخدمون المستهدفون عاجزين عن فهمه. الوضوح والبساطة واستخدام لغة ومفاهيم ملائمة للمجال عوامل تضمن استفادة الجمهور من التفسير

### • (Actionability) الإمكانية التطبيقية:

أفضل التفسيرات تلك التي تمكّن من اتخاذ إجراءات ملموسة. فإذا علم الطبيب أن ارتفاع مستوى الكوليسترول عنصر حاسم، يمكنه طلب فحوصات إضافية. وإذا أدرك المصرف أن السجل الوظيفي يؤثر على القروض، قد يراجع معاييرها أو يقدّم برامج توعية مالية

، تتطور معايير ومقاييس موحدة ودراسات مقارنة وأبحاث متعلقة بتجربة المستخدم لتحسين طرق XAI مع نضوج مجال التفسير واختيار تلك التي تقدّم قيمة حقيقية.

### الخلاصة

استعرض هذا الفصل مجموعة واسعة من الاستراتيجيات لتعزيز قابلية تفسير النماذج، بدءًا من النماذج القابلة للتفسير الأكثر تعقيدًا والمحايدة لنوع النموذج. كما تطرّقنا إلى تقنيات مثل (Post-Hoc) بطبيعتها إلى أساليب ما بعد التدريب ، إضافة إلى أدوات كالمخططات التفاعلية والتفسيرات المضادة، مع التأكيد على ضرورة تقييم جودة SHAP و LIME التفسيرات وفق معايير الدقة والثبات والفهم والقابلية للتطبيق.

تساعد هذه الأدوات والمفاهيم الممارسين على التعامل بفاعلية مع تعقيد أنظمة الذكاء الاصطناعي الحديثة، والتأكد من أن "الصناديق السوداء" يمكن أن تصبح مفهومة وجديرة بالثقة ومسؤولة، ومتوافقة مع القيم الإنسانية.



## الفصل الرابع: الأدوات والتقنيات لتطبيق الذكاء الاصطناعي القابل للتفسير

بعد استعراض الأطر المفاهيمية والأساليب النظرية لجعل أنظمة الذكاء الاصطناعي أكثر شفافية، ننتقل الآن إلى الأدوات العملية التي تُمكن هذه المبادئ من أن تزدهر في البيئات الواقعية. فهناك منظومة متكاملة من المكتبات البرمجية والحزم المتكاملة وأفضل الممارسات التي تساعد الممارسين على إدماج قابلية التفسير في خطوط التطوير الخاصة بهم. ومن خلال الاستفادة من هذه الأدوات، يمكن للمؤسسات غرس مبدأ الشفافية في صلب مبادراتها في مجال الذكاء الاصطناعي

### الأدوات والتقنيات لتطبيق الذكاء الاصطناعي القابل للتفسير (XAI)

#### الأدوات المتكاملة

تقدم حلولاً شاملة تجمع بين أدوات متعددة للتكامل السلس

#### مكتبات البرمجيات

أساسية لتوفير الكود التأسيسي والوظائف اللازمة لـ XAI.



#### أفضل الممارسات

إرشادات واستراتيجيات تضمن التنفيذ الفعال والكفاء لـ XAI.

### XAI الأطر والمكتبات المتخصصة في 4.1

تتيح مجموعة قوية من المكتبات مفتوحة المصدر إنتاج التفسيرات وتصوّرها وتطوير نماذج أكثر قابلية للفهم بسهولة متزايدة. وقد حازت الأدوات التالية على شهرة واسعة بفضل تنوّعها وفعاليتها

#### • (بايثون) LIME:

توليد تفسيرات محلية لأي نموذج تنبؤ تقريباً. وبضعة أسطر من الشفرة تكفي لإنتاج تفصيل LIME تيسر يوضّح سبب تفصيل النموذج لمخرج معين في حالة محددة. تجعل طبيعتها المحايدة للنماذج واستخدامها البسيط منها خياراً شائعاً لإجراء تقييمات سريعة وتفسيرية

#### • (مكتبة بايثون) SHAP:

طرائق التفسير المتنوّعة ضمن إطار رياضي مستند إلى نظرية الألعاب التعاونية. سواء اعتمدت SHAP توحده واجهة SHAP أو استخدمت الشبكات العصبية العميقة، توفر XGBoost على أساليب قائمة على الأشجار مثل التي توضح دور كل خاصية محلياً وعالمياً (SHAP values) موحدة لحساب قيم المساهمة

- **InterpretML (من مايكروسوفت):**

تسمح (Post-Hoc) مجموعة شاملة من التقنيات القابلة للتفسير ومفسرات ما بعد التدريب InterpretML تقدم لوحات التحكم البصرية للعلماء وأصحاب المصلحة من غير التقنيين باستكشاف كيفية تشكل التنبؤات، ما يعزز التعاون والنقاش البناء داخل الفرق.

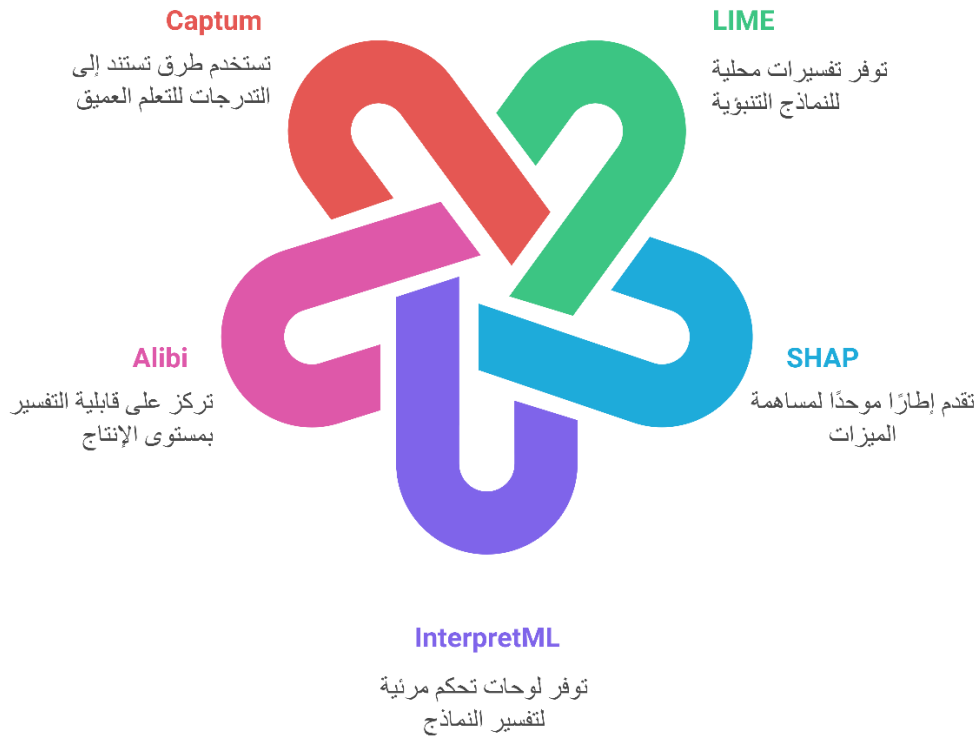
- **Alibi (من Seldon):**

على التفسيرية في البيانات الإنتاجية وقدرات المراقبة. يدعم مجموعة من أساليب التفسير—بدءًا من Alibi يركز التي تحدّد السمات المحورية في القرارات وصولاً إلى أساليب التفسير المضاد "Anchors" لتزويد المؤسسات بحلول قابلة للتوسع وملائمة للمتطلبات المؤسسية—(Counterfactual)

- **Captum (من Facebook AI Research):**

يتضمّن أساليب تستند إلى التدرجات PyTorch بسلاسة مع Captum مصمم لتطبيقات التعلم العميق، يندمج لمساعدة المطوّرين على فهم التمثيلات الداخلية التي تتعلّمها الشبكات العصبية، كاشفًا (Gradient-based) الخصائص المدخّلية المحفّزة لاستجابات أو مخرجات معيّنة

### نظرة عامة على أدوات الذكاء الاصطناعي القابل للتفسير (XAI)



## 4.2 XAI الخطوات العملية لدمج

يتطلب إدماج قابلية التفسير في بيئة إنتاجية تخطيطاً منهجياً وتنفيذاً مدروساً. ضع في الاعتبار المراحل التالية:

### 1. تحديد الأهداف:

حدّد بوضوح سبب حاجتك للتفسيرات. هل الهدف هو الامتثال التنظيمي؟ بناء الثقة لدى المستخدمين النهائيين؟ أم تحسين تصحيح النماذج وتطويرها داخلياً؟ يساعد تحديد الأهداف منذ البداية في اختيار الأدوات والأساليب المناسبة.

### 2. اختيار النموذج:

إذا كانت قابلية التفسير تفوق في أهميتها الأداء التنبؤي الخام، ابدأ بنماذج مفهومة بطبيعتها. أما إذا كنت بحاجة لأحدث درجات الدقة، فاستعد لاستخدام مفسّرات ما بعد التدريب على النماذج الأكثر تعقيداً. يؤثر نوع النموذج الأنسب XAI المختار على أساليب

### 3. دمج الأدوات:

في خط تطويرك. أجر تجارب على حالات اختبارية للتحقق من أن SHAP أو LIME ضمن مكتبات مثل التفسيرات تتوافق مع المعرفة الميدانية ومع توقعات أصحاب المصلحة. يُحوّل هذا الدمج المعرفة النظرية إلى مخرجات ملموسة.

### 4. التصور والتواصل:

SHAP تعزز الأدوات البصرية إمكانية استيعاب التفسيرات. تتيح مخططات أهمية الخصائص ومخططات تحويل المنطق المعقد إلى رسومات يسهل فهمها. تضمن لوحات (PDP) الملخصة ومخططات الاعتماد الجزئي القيادة المشتركة أن الفرق متعددة التخصصات تتشارك رؤية موحدة للنموذج وتفسيراته.

### 5. التكرار والتحسين:

كحال الدقة، تتطور جودة التفسيرات بمرور الوقت. احصل على ملاحظات من الخبراء الميدانيين—كالأطباء أو مسؤولي القروض أو محلي التسويق—وحسّن أساليبك بناءً على ذلك. كما تفيد الدراسات الميدانية لمعرفة أي التفسيرات تملك قيمة حقيقية للمستخدمين. عدّل وحسّن حتى تلبي التفسيرات احتياجات جميع الأطراف.

## 4.3 SHAP مثال تطبيقي: تفسير نموذج الموافقة على القروض باستخدام

(XGBoost) لتوضيح كيفية عمل هذه الأدوات عملياً، تخيل سيناريوًا درّبت فيه مصنعاً قائماً على الأشجار المعرزة للتنبؤ بالموافقة على القروض. تشمل البيانات خصائص مثل الدخل، ومدة الوظيفة، ونسبة استخدام البطاقة الائتمانية، والتخلف السابق عن السداد.

### الخطوة 1: تدريب النموذج

```
python
Copy code
import xgboost as xgb
```

```
model = xgb.XGBClassifier().fit(X_train, y_train)
```

- **SHAP الخطوة 2: تثبيت وتشغيل**

```
python
Copy code
pip install shap
import shap
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)
```

- **(Global Explanation) الخطوة 3: تفسير عالمي**

```
python
Copy code
shap.summary_plot(shap_values, X_test)
```

يعرض المخطط الملخص الخصائص مرتبة حسب أهميتها العالمية، ويوضح تأثيرها على دفع التنبؤ نحو الموافقة أو الرفض. إذا تبين أن "نسبة استخدام البطاقة الائتمانية" تدفع باستمرار نحو الرفض، يفهم أصحاب المصلحة على الفور أن ارتفاع هذه النسبة يمثل عامل مخاطرة رئيسيًا.

- **(Local Explanation) الخطوة 4: تفسير محلي**

```
python
Copy code
# فحص حالة أحد المتقدمين
shap.force_plot(explainer.expected_value, shap_values[0,:],
                X_test.iloc[0,:])
```

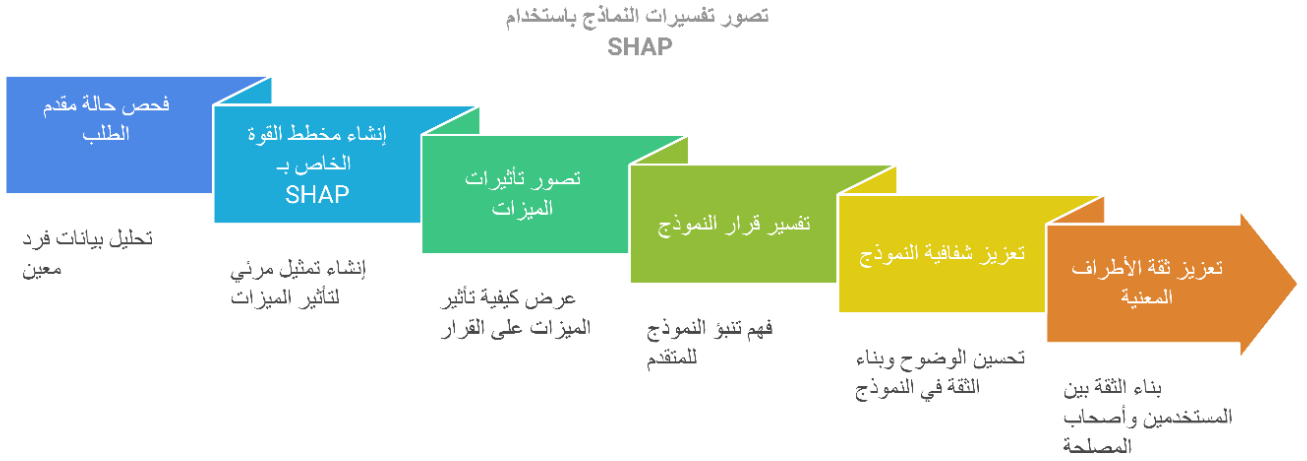
كيف أثر كل عامل على قرار النموذج. إذا رُفض القرض، فإن انخفاض الدخل "Force Plot" بالنسبة لحالة فردية، يبين وارتفاع نسبة استخدام البطاقة قد يكونان السببين الرئيسيين، في حين ساهم سجل وظيفي مستقر في تخفيف الأثر السلبي جزئيًا. يقدّم هذا المنظور المفصّل لموظفي القروض أساسًا لتبرير القرار والتواصل بوضوح مع المتقدم.

عبر هذه الخطوات، يتحوّل نموذج تنبؤي مبهم في البداية إلى نموذج شفاف. يمكن لأصحاب المصلحة رؤية دوافع كل قرار، وتحديد مجالات التحسين، والثقة بالنظام بدرجة أكبر.

## الخلاصة

إلى روى قابلة للتطبيق. من مكتبات XAI استعرض هذا الفصل مجموعة الأدوات العملية التي تحول الوعد النظري للـ ، تتوافر اليوم ثروة من الموارد لإدماج Alibi و InterpretML إلى منصات شاملة مثل SHAP و LIME مشهورة مثل

التفسيرية في خطوط العمل. من خلال تحديد الأهداف، واختيار النماذج المناسبة، والاستفادة من أدوات التصور، والاستمرار في تحسين الأساليب، يمكنك ضمان أن تظل أنظمة الذكاء الاصطناعي ليس دقيقة وحسب، بل ومفهومة وجديرة بالثقة، ما يعزز اتخاذ قرارات مستنيرة ويمكن المؤسسة من المضي قدمًا بثقة



## الفصل الخامس: تطبيقات الذكاء الاصطناعي القابل للتفسير بحسب القطاعات

حدود الأمانة التقنية، إذ يؤثر مباشرة على الطريقة التي يثق بها (XAI) يتجاوز الذكاء الاصطناعي القابل للتفسير المحترفون في مختلف المجالات بالأنظمة الذكية، وكيف يتبنونها ويستفيدون منها. ورغم أن الأدوات والمبادئ الأساسية للـ XAI تبقى ثابتة عمومًا، فإن أهميتها العملية ونقاط التركيز تختلف باختلاف القطاع. ففي بعض السياقات، تحافظ قابلية XAI التفسير على الأرواح وتدعم المعايير الأخلاقية، بينما في سياقات أخرى، تعزز الامتثال التنظيمي وثقة العملاء والكفاءة التشغيلية

### دور الذكاء الاصطناعي القابل للتفسير عبر الصناعات (XAI)



## 5.1 الرعاية الصحية

في الطب، قد تمثل الدقة بدون وضوح خطورة حقيقية. يحتاج الأطباء والممرضون والإداريون إلى فهم منطقي واضح لتوصيات النموذج التشخيصي المدعوم بالذكاء الاصطناعي. فشفافية النموذج لا تعزز الثقة فحسب، بل تساعد أيضًا في كشف الأخطاء وصقل استراتيجيات العلاج.

### • مثال:

يمكن لأخصائي الأشعة الذي يستعين بمصنف ذكاء اصطناعي لتحليل صور الرنين المغناطيسي الاستفادة من التي تبرز المناطق المشبوهة. إذا أشار النظام إلى منطقة معينة باعتبارها دليلاً (Heatmaps) خرائط الحرارة على وجود ورم، يمكن لأخصائي الأشعة تأكيد أو دحض هذا التركيز. عبر إظهار المنطق، يكمل الذكاء الاصطناعي الخبرة البشرية بدلاً من أن يحل محلها.

### • الفوائد:

تعزيز الثقة والالتزام بالتوصيات، تسريع القرارات التشخيصية، تقليل الحالات المرضية التي لا تُكتشف، واستنباط علاقات جديدة تفيد الأبحاث المستقبلية.

## 5.2 التمويل والمصارف

في المجال المالي، يطالب المنظمون والمدققون والمستهلكون بمعرفة سبب إصدار أنظمة الذكاء الاصطناعي لقرارات معينة، سواء كانت في الإقراض أو الاستثمار أو كشف الاحتيال. يطمئن النموذج الشفاف العملاء إلى أن القرارات تستند إلى معايير موضوعية، وليست عشوائية أو تمييزية.

### • مثال:

لتفسير رفض قرض معين. إذا أوضح النموذج أن SHAP يمكن لمصرف يقيم المخاطر الائتمانية استخدام قيم انخفاض الدرجة الائتمانية وتعدّد حالات التخلف عن السداد تفوق تأثير الدخل الثابت والمدخرات، فإن العميل يفهم المنطق. توفّر هذه الشفافية ردعًا عن الاتهامات بالمعاملة غير العادلة، وتلبي التوقعات التنظيمية حول ضرورة التفسير.

### • الفوائد:

الامتثال للقوانين (مثل قانون تكافؤ الفرص الائتمانية)، تقليل المخاطر القانونية، تحسين علاقات العملاء، إدارة أفضل للمخاطر، وتعزيز السمعة بالعدالة والإنصاف.

## 5.3 الاعتبارات القانونية والأخلاقية

في السياقات القانونية—مثل التنبؤ بالجريمة، وتحليل العقود، والتوصيات في الأحكام القضائية—تتطلب جسامه القرارات تدقيقًا صارمًا. يضمن الذكاء الاصطناعي القابل للتفسير عدم تعرض الأفراد لأحكام خوارزمية غير مفهومة أو متحيزة، تؤثر على حقوقهم وحياتهم الأساسية.

### • مثال:

قد يصنّف خوارزمية تنبؤ بالجريمة أحد الأحياء على أنه "عالي الخطورة". يكشف التفسير الشفاف عن أن

معدلات الاعتقال التاريخية والمؤشرات الاجتماعية الاقتصادية أثرت في هذا التقييم. يتيح هذا الاعتراف بالمحركات الأساسية للنموذج إجراء نقاش حول التحيز المنهجي وإمكانية الحاجة إلى مراجعة البيانات أو السياسات.

• **الفوائد:**

الحفاظ على العدالة، تجنّب الأنماط التمييزية، التوافق مع المعايير القانونية، وتعزيز ثقة الجمهور بالأحكام القضائية والإدارية.

#### 5.4 (Industry 4.0) التصنيع والثورة الصناعية الرابعة

تعتمد العمليات الصناعية على الذكاء الاصطناعي للتنبؤ بأعطال المعدات، وتحسين سلاسل الإمداد، وتقليل فترات التوقف. تساعد قابلية التفسير المديرين على فهم سبب توقع النموذج لفشل ما أو التوصية بتدخل معيّن، مما يحسّن خطط الصيانة وتوزيع الموارد.

• **مثال:**

، قد يكتشف مدير مصنع أن ارتفاع وتيرة الاهتزاز في جهاز (PDP) عند فحص مخططات الاعتماد الجزئي معين مؤشر رئيسي لقرب تعطل الماكينة. عبر هذه المعرفة، يمكن لفرق الصيانة التدخل مبكرًا، مما يوفر الوقت والتكاليف.

• **الفوائد:**

تقليل فترات التوقف غير المخطط لها، خفض تكاليف الصيانة، تحسين السلامة، واستخدام أكثر كفاءة للموارد البشرية والمادية.

#### 5.5 إدارة المشتريات وسلاسل التوريد

تشكل المشتريات واللوجستيات العمود الفقري لكثير من الأعمال، مما يضمن تدفق المواد الخام والمكونات والسلع بسلاسة. ومع اتخاذ نماذج الذكاء الاصطناعي قرارات حول استراتيجيات الطلب، واختيار الموردين، ومستويات المخزون، فإن الشفافية توضح كيفية التوصل لهذه القرارات، مما يساعد خبراء سلاسل التوريد على الثقة بتوجيهات النموذج.

• **مثال:**

قد يوصي أداة تحسين المخزون المدعومة بالذكاء الاصطناعي بزيادة مخزون مكون معيّن قبل موسم الأعياد. عبر فحص تفسيرات النموذج—مثل أوقات التوريد، والارتفاع التاريخي في الطلب، وموثوقية المورد—يكتسب مسؤولو المشتريات الثقة بأن هذه التوقعات ليست عشوائية. كما يمكنهم رصد ما إذا كان النموذج يعتمد بشكل مفرط على بيانات قديمة أو يتجاهل مؤشرات حديثة لجودة المورد.

• **الفوائد:**

سلاسل توريد أكثر مرونة، تقليل النقص أو الفائض في المخزون، تقليل مخاطر المشتريات، تحسّن استراتيجيات التفاوض مع الموردين، واتخاذ قرارات شراء مستنيرة تعتمد على البيانات.



## الخدمات اللوجستية والنقل 5.6

مع تزايد تعقيد شبكات اللوجستيات العالمية وأنظمة النقل، يساعد الذكاء الاصطناعي في تحسين المسارات وخطط السعة وجدول التسليم. لكن مديري اللوجستيات يحتاجون لفهم سبب اقتراح الخوارزمية لمسار معين أو تحديد الاختناقات

### • مثال:

قد يقترح نموذج لتحسين المسارات في شركة شحن إعادة توجيه سفن الشحن بعيداً عن مرفأ معين خلال موسم الأمطار. عبر تفسير ذلك بأن القرار يستند إلى توقعات الطقس، وبيانات التأخير التاريخية، وتقارير الازدحام، يمكن لفريق اللوجستيات التحقق من صحة الاستراتيجية والاستعداد لحالات الطوارئ

### • الفوائد:

أوقات تسليم أكثر موثوقية، تخطيط مسارات أكثر فعالية من حيث التكلفة، استباق الاضطرابات، تحسين الالتزام بالمعايير البيئية (مثلاً تجنب الازدحام لتقليل الانبعاثات)، وتعزيز جودة الخدمة

## التجزئة والتسويق 5.7

في قطاع التجزئة، تؤثر قابلية التفسير على تجربة العملاء والتخطيط الاستراتيجي. ورغم أن المخاطر قد لا تكون جسيمة كما في الرعاية الصحية أو القانون، فإن الشفافية تظل مهمة. يقدر المستهلكون فهم سبب ظهور منتجات معينة في توصياتهم، مما يعزز شعورهم بأن النظام يحترم تفضيلاتهم

### • مثال:

قد يوضح نظام توصية في موقع تجارة إلكترونية كيف أثرت مشتريات العميل السابقة، وأنماط التقييم، والاتجاهات الموسمية في اقتراحاته. يعزز هذا الاطمئنان ولاء العملاء ويشجعهم على الشراء المتكرر، إذ يدركون أن التوصيات نابعة من رؤية حقيقية وليست محاولات للتلاعب

### • الفوائد:

تعزيز ثقة العلامة التجارية، تحسين تفاعل المستخدمين، اتخاذ قرارات مستندة إلى البيانات بشأن تشكيل تشكيلة المنتجات، وتقوية قدرة الاستجابة للتغيرات في ديناميات السوق

## الطاقة والمرافق العامة 5.8

تعتمد شركات الطاقة والمرافق على الذكاء الاصطناعي للتنبؤ بالطلب وتحقيق توازن الأحمال ودعم مبادرات الاستدامة. توضح النماذج القابلة للتفسير كيف تشكل أنماط الطقس والاستخدام الصناعي واللوائح السياسية التوقعات

### • مثال:

قد يتنبأ نموذج لاستهلاك الكهرباء بارتفاع في الطلب، استناداً إلى موجة حر مرتقبة ونمو سكاني في مناطق معينة. عبر شرح هذا المنطق، يمكن لصنّاع السياسات أو مخططي البنية التحتية الاستعداد لمواجهة التحديات وتحسين استثمار الموارد

#### • الفوائد:

استقرار أكبر للشبكة الكهربائية، تقليل مخاطر انقطاع التيار، قرارات استثمارية مدروسة في البنية التحتية، تحسين الالتزام بالإرشادات البيئية، وزيادة ثقة الجمهور في إدارة الموارد.

### 5.9 الاتصالات وخدمات تكنولوجيا المعلومات

يستخدم مشغلو الشبكات ومقدمو خدمات تكنولوجيا المعلومات الذكاء الاصطناعي لتوقع الازدحام وتحسين تخصيص النطاق الترددي وتعزيز الأمن السيبراني. تضمن الشفافية في هذه القرارات أن تظل جودة الخدمة وبيانات العملاء على رأس الأولويات.

#### • مثال:

قد يحدد نموذج في شركة اتصالات فترات ازدحام الشبكة، موضحًا أن بعض ساعات الذروة أو نقاط الضعف في الأجهزة هي الأسباب الرئيسية. عبر تفسير هذا المنطق، يمكن للمهندسين تعزيز سعة الشبكة قبل حدوث المشكلات، ويفهم العملاء سبب أي تغييرات في جودة الخدمة أو التوصيات.

#### • الفوائد:

جودة خدمة أعلى، قدرة استباقية على حل المشكلات، استراتيجيات محسنة لأمن المعلومات، ورضا المستخدمين المحسن.

### 5.10 التعليم والتعلم عبر الإنترنت

تعتمد منصات التعلم عبر الإنترنت ونظم التدريس المدعومة بالذكاء الاصطناعي على التوصيات التكيفية. تبرز النماذج القابلة للتفسير سبب اقتراح خطط دراسية أو تمارين معينة، ما يساعد المعلمين على تحسين المناهج ويقود الطلاب نحو المناطق التي تحتاج إلى تطوير.

#### • مثال:

قد يقترح نظام تعليمي عبر الإنترنت إعادة دراسة مفهوم معين لطالب محدد. عبر توضيح أن الطالب واجه صعوبة مسبقاً مع مفهوم مشابه وحقق تحسناً بعد مراجعته، يفهم كل من المعلم والطالب المنطق وراء هذه التوصية.

#### • الفوائد:

تجارب تعليمية أكثر خصوصية، تحسين النتائج الأكاديمية، تقديم رؤى قابلة للتنفيذ حول تقدم الطالب، وتعزيز التفاعل بين المتعلمين والمواد التعليمية.

### خلاصة التطبيقات القطاعية

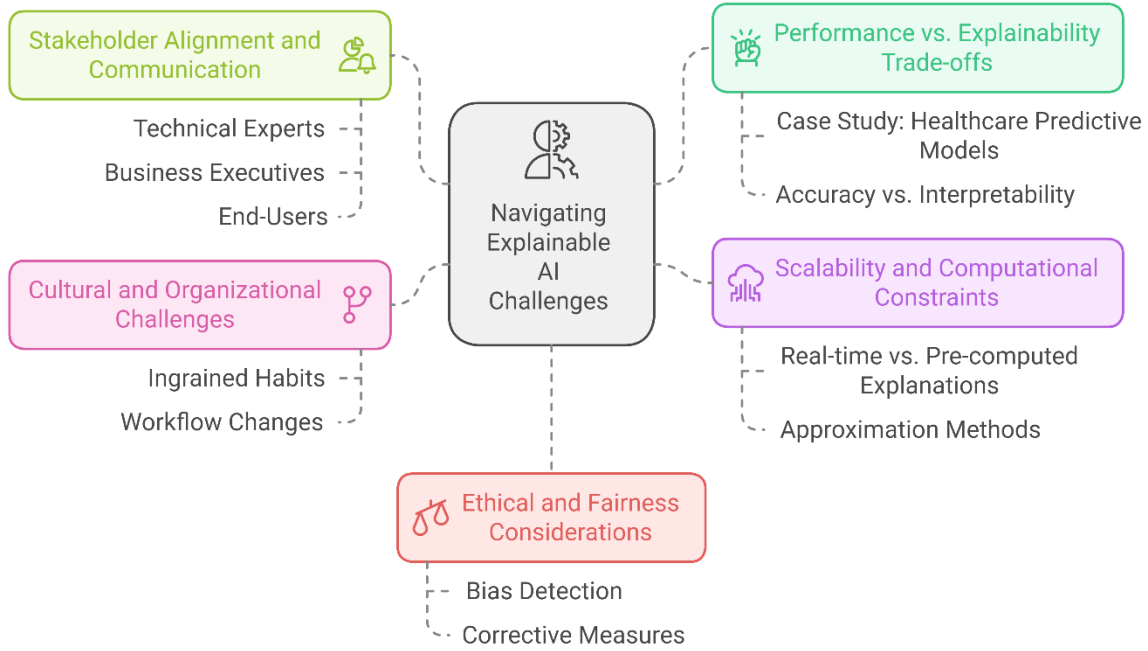
من إنقاذ الأرواح في المستشفيات إلى ضمان العدالة في المحاكم، ومن تحسين قرارات المشتريات إلى صقل استراتيجيات التسويق، يتكيف الذكاء الاصطناعي القابل للتفسير مع التحديات والأولويات الفريدة لكل قطاع. ورغم ثبات الأساليب والتفسيرات المضادة والنماذج القابلة للتفسير بطبيعتها—تختلف دوافع الاهتمام SHAP وLIME الأساسية—مثل

بالتفسير. ففي الرعاية الصحية والتمويل، تهيمن الحاجة إلى الامتثال والثقة. وفي اللوجستيات والتصنيع، تقود الكفاءة والوضوح التشغيلي الاهتمام، بينما يركز قطاعا التجزئة والتسويق على راحة المستهلك وثقة العلامة التجارية.

المكيف حسب السياق، تعزز المنظمات موثوقية أنظمتها، وتشجع الابتكار المسؤول، وتلتزم بالمعايير XAI عبر تبني التنظيمية والأخلاقية. وبذلك تُمكن صنّاع القرار—مهما كان مجالهم—من تحقيق أهدافهم بثقة أكبر، معتمدين على النماذج الذكية القابلة للتفسير

## الفصل السادس: الإبحار في تحديات الذكاء الاصطناعي القابل للتفسير

رغم ما يبشّر به الذكاء الاصطناعي القابل للتفسير من آفاق واعدة، فإن تطبيقه عمليًا لا يخلو من العقبات. يمثل تحقيق التوازن بين قابلية التفسير والأداء، والتعامل مع القيود الحسابية، ومواءمة أسلوب التواصل مع أطراف مختلفة، وتحفيز التغيير المؤسسي، وضمان العدالة—كل ذلك تحديات جسيمة. إن تجاوز هذه العقبات يحوّل التفسيرية من هدف مثالي إلى حقيقة مستدامة



### الموازنة بين الأداء وقابلية التفسير 6.1

كثيرًا ما تحقق النماذج ذات الدقة التنبؤية العالية أداءً متميزًا بفضل تعقيدها، إلا أنّ هذا التعقيد يحوّل منطق اتخاذ القرار إلى صندوق أسود غير مفهوم. هنا يتجلى السؤال: هل تستحق زيادة طفيفة في الدقة التضحية بقابلية التفسير والثقة؟

في كثير من المجالات المنظمة أو الحساسة، تميل الكفة لصالح الشفافية. فالنموذج الأقل دقة بنسبة بسيطة لكنه أكثر قابلية للفهم قد يشكل الفرق بين صنع قرارات واعية وثقة، والاعتماد غير المريح على تنبؤات مبهمّة.

#### • دراسة حالة:

تخيّل مستشفى يسعى إلى التنبؤ بإعادة إدخال المرضى. رغم أن شبكة عصبية عميقة قد تحقق دقة استثنائية، يفضل الأطباء ربما نموذج انحدار لوجستي يضحى ببضعة نقاط من الدقة مقابل معاملات واضحة ومفهومة. تمكن هذه الشفافية المتخصصين من تحديد عوامل الخطر الأساسية—مثل مدة الإقامة في المستشفى أو قيم مخبرية معيّنة—وصياغة تدخلات علاجية بثقة. عمليًا، تتفوق قابلية الفهم على إغراء التحسن الهامشي في الدقة.

## قابلية التوسع والقيود الحسابية 6.2

قد تكون عمليات التفسير مكلفة حسابياً، خاصةً مع النماذج المعقدة والبيانات الضخمة. تتطلب بعض تقنيات التفسير، مثل SHAP، تقييم النموذج عدة مرات لكل حالة. وفي حين قد يكون ذلك مقبولاً في بيئة بحثية، يصبح الأمر مرهقاً عند تطبيقه SHAP في الزمن الحقيقي أو على نطاق واسع.

على المؤسسات تحديد مدى عمق التفسيرات المطلوبة وتواترها. فهل تحتاج فقط لوضع حالات تمثيلية لفهم سلوك النموذج ككل، أم تحتاج لتفسيرات شبه فورية لكل تنبؤ؟ في كثير من الأحيان، يمكن اعتماد استراتيجيات مثل حساب التفسيرات الشاملة مسبقاً، أو استخدام طرق تقريبية، أو التركيز على الحالات الأكثر أهمية، لإيجاد توازن بين العمق والكفاءة.

## توافق أصحاب المصلحة والتواصل 6.3

تخدم التفسيرات جمهوراً متنوعاً، لكل منه أولوياته:

### • علماء البيانات والمهندسون

يحتاج الخبراء التقنيون إلى تفسيرات دقيقة ومرتكزة إلى أسس رياضية لتصحيح الأخطاء وتحسين ميزات النموذج. يقدرون الرسوم البيانية التفصيلية والنسب الإحصائية والتفاصيل على مستوى الشفرة.

### • الإداريون والمنظمون

يهتم المدبرون والمسؤولون بوصول سريع وموجز لأسباب النتائج دون الخوض في التفاصيل التقنية. وبالمثل، يرغب المشرفون في وضوح يضمن الامتثال والإنصاف، دون الحاجة إلى دقائق الخوارزميات.

### • المستخدمون النهائيون والعملاء

المستفيدون من قرارات الذكاء الاصطناعي—مثل طالب قرض—يحتاجون تفسيرات بلغة واضحة. بدل معادلة طويلة، يكتفون بجملة بسيطة: "رُفض القرض بسبب ارتفاع نسبة استخدامك لبطاقة الائتمان"، ليتمكنوا من اتخاذ إجراءات تصحيحية.

يضمن تكييف التفسيرات لكل فئة وصول المعلومات بوضوح وفعالية. ويمكن لأدوات كلوحات التحكم التفاعلية، والتفسيرات المتدرجة التي تكشف تفاصيل أكثر عند الحاجة، والملخصات النصية بلغة بسيطة، ردم الفجوة بين المنطق التقني العميق والسرور المفهوم للمستخدم.

## التحديات الثقافية والتنظيمية 6.4

قد تعيق العادات الراسخة التحول نحو الشفافية. فالفرق التي اعتادت الاعتماد على نماذج الصندوق الأسود قد تقاوم التحول نحو التفسير، خوفاً من عبء عمل إضافي أو انخفاض محتمل في الأداء. إقناعهم يتطلب إظهار أن التفسير يعزز صنع القرار، ويقوي سمعة العلامة التجارية، ويقلل المخاطر.

أيضاً إعادة تفكير في التدفقات العملية. فعلى غرار اختبارات الجودة للشفرة، ينبغي التحقق من صحة XAI يتطلب تبني التفسيرات. وقد يستدعي ذلك أدوات متخصصة—مثل "مهندسي التفسير" أو "أخصائيي حوكمة النماذج"—يعملون على ضمان جودة التفسيرات والالتزام بالمعايير الأخلاقية والأنظمة ذات الصلة.

## 6.5 الاعتبارات الأخلاقية وتحقيق العدالة

رغم ما تقدمه التفسيرية من قوة في كشف التحيز والظلم، فهي ليست حلاً سحرياً. فإذا أظهرت التفسيرات أن المنطق الداخلي للنموذج ينطوي على أنماط تمييزية، يبقى على المعنيين اتخاذ خطوات تصحيحية. يكشف فهم مصدر التحيز—سواء كان من البيانات التدريبية المنحازة أو الخصائص غير الملائمة أو بنية النموذج—أوجه القصور ويقود لجهود الإصلاح.

على سبيل المثال، قد توضح التفسيرات أن خوارزمية التوظيف تعاقب فئات سكانية معينة نتيجة اختلالات تاريخية في البيانات. بالاستعانة بهذه المعرفة، يمكن للمؤسسة إعادة تدريب النموذج على بيانات أكثر تمثيلاً، وحذف الخصائص المنحازة، أو فرض قيود تعزز العدالة. يوضّح هذا التفاعل كيف يتكامل التفسير والإنصاف لتعزيز نتائج أكثر عدلاً.

## 6.6 رسم المسار نحو الأمام

يتطلّب اجتياز هذه التحديات استراتيجيات شمولية تدرك التنازلات وتختلف حسب السياق. قد تقبل بعض المؤسسات بتضحية بسيطة في الدقة مقابل قابلية تفسير أعلى وامتثال تنظيمي أفضل؛ بينما قد يستثمر البعض الآخر في تقنيات تفسير متطورة. وبُنِي تحثية حسابية للحفاظ على الأداء مع تعزيز الشفافية.

في النهاية، فإن مكاسب تحقيق توازن مناسب عميقة. بتحويل الذكاء الاصطناعي من صندوق أسود غامض إلى مستشار موثوق، يفهم أصحاب المصلحة التنبؤات، ويثق المشرّعون بالعملية، ويشعر المستخدمون بالاحترام، ويكتسب علماء البيانات رؤى جديدة لمزيد من الابتكار.

وبذلك، تتجاوز التفسيرية مكانتها كميزة تقنية إضافية لتصبح ركيزة أساسية لنشر الذكاء الاصطناعي بصورة أخلاقية وفعالة ومستدامة.

## الفصل السابع: الاعتبارات الأخلاقية والأثر المجتمعي

مع اندماج أنظمة الذكاء الاصطناعي في نسيج الحياة اليومية، لا يمكن تجاهل تبعاتها الأخلاقية. وفي هذا السياق، يقف في صميم النقاش، إذ يمثل جسراً حيويًا بين القدرات التقنية والمسؤولية (XAI) الذكاء الاصطناعي القابل للتفسير قدرة المجتمعات على مساءلة الخوارزميات، XAI الأخلاقية. فمن خلال جعل القرارات الآلية أكثر شفافية، يضمن ومواجهة التحيزات، ومواءمة التكنولوجيا مع القيم الإنسانية المشتركة

### أخلاقيات الذكاء الاصطناعي والمجتمع



## بناء أطر أخلاقية للذكاء الاصطناعي 7.1

يعد وضع أطر أخلاقية للذكاء الاصطناعي عملية مستمرة للتفاوض والتحسين والتنفيذ. فقد تبنت العديد من المؤسسات مدونات أخلاقيات تشمل مبادئ مثل العدالة والمساءلة والشفافية والخصوصية. غالبًا ما تُترجم هذه المبادئ إلى متطلبات عملية، مثل ضرورة اعتماد نماذج قابلة للتفسير، ومنح الأفراد المتأثرين قرارات آلية حق الحصول على تفسيرات ذات مغزى، ووضع إجراءات للتدقيق في النتائج الظالمة وتصحيحها.

يلاحظ تزايد إنشاء مجالس مراجعة داخلية—تضم خبراء أخلاقيات، ومختصين في المجال، وممثلين عن المجتمع، ومتخصصين تقنيين—لمراجعة النماذج قبل نشرها. لا تقتصر معايير هذه المجالس على أداء النموذج فحسب، بل تمتد لتشمل سلامة التفسيرات المقدمة. وبذلك تضمن أنظمة الذكاء الاصطناعي عند نشرها مستوىً أدنى من المعايير الأخلاقية. كما تعترف الحكومات والمنظمات الدولية بأهمية قابلية التفسير، وتدرجها ضمن اللوائح والإرشادات الوليدة في هذا المجال.

## خصوصية البيانات والأمن 7.2

يتقاطع مفهوم التفسير مع مساحات حساسة عندما يكشف تفاصيل ينبغي أن تبقى سرية. فقد يؤدي تفسير آليات عمل نموذج طبي إلى كشف معلومات صحية محمية. وبالمثل، قد تتضمن التفسيرات في المجال المالي إفصاحًا عن أسرار تجارية أو هياكل تسعير الموردين.

، (k-Anonymity) وتقنيات إخفاء الهوية (Differential Privacy) تساعد مناهج مثل الخصوصية التفاضلية والقواميس المضبوطة، في تحقيق توازن بين الشفافية والسرية. ومن خلال إضافة "ضوضاء" إلى البيانات أو عرض رؤى مجمعة بدلًا من تفاصيل على مستوى الأفراد، تستطيع المؤسسات تقديم تفسيرات ذات مغزى دون المساس بخصوصية البيانات أو المنطق التجاري المحمي. والهدف هو الحفاظ على مستوى من قابلية التفسير يمكن الأطراف المعنية دون أن يعرض البيانات الحساسة أو المعارف التجارية للخطر.

## الأثر المجتمعي ورصد التحيزات 7.3

غالبًا ما يرث الذكاء الاصطناعي التحيزات المضمنة في البيانات التاريخية. وفي غياب الشفافية، تبقى هذه التحيزات مخفية فهم العوامل التي يعتمد عليها النموذج، مما يمكن الأطراف XAI وتعيد إنتاج حالات اللامساواة القائمة. تتيح أساليب المعنية من تحديد الأنماط التمييزية. وبفضل هذه المعرفة، يمكن للمؤسسات اتخاذ إجراءات تصحيحية—كإعادة تدريب النماذج على بيانات أكثر تمثيلاً، أو إزالة الخصائص المنحازة، أو فرض قيود تضمن الإنصاف.

### دراسة حالة: الحد من التحيز في التوظيف

نشرت شركة تكنولوجية كبيرة أداة لفرز السير الذاتية فوجدت أنها تميز ضد المتقدمات الإناث. أظهرت قيم أن بعض الكلمات المفتاحية—التي تنتشر إحصائيًا أكثر في السير الذاتية للذكور—تعزز التقييم. بإعادة SHAP تدريب النموذج على بيانات أكثر إنصافًا وفرض قيود أخلاقية، أعادت الشركة نشر نموذج أكثر عدالة. ما كان في البداية XAI لهذا التحول أن يتم لولا الوضوح الذي وفره.



#### 7.4 الفهم العام والنقاش المجتمعي

مع قيام أنظمة الذكاء الاصطناعي بتحديد مسارات التوظيف والعلاج الطبي والموافقة الانتمائية، بل وحتى الإعلانات الانتخابية، لا بد للمواطنين من فهم كيفية تشكيل هذه النماذج لمصائرهم. تعزز التفسيرية "الثقافة الذكائية" وتساعد العامة في إدراك أن هذه الأنظمة ليست سحرية أو معصومة. ويؤدي ارتفاع مستوى الوعي والقدرة على النقد البناء إلى نقاشات عامة أكثر نضجًا، ما يفضي إلى رقابة أقوى وسياسات أفضل، وأطر عمل تعكس القيم الجماعية.

تلعب الحكومات ومجموعات المناصرة والمؤسسات التعليمية دورًا محوريًا هنا. إذ يمكن أن تساعد ورش العمل العامة، والدورات التدريبية عبر الإنترنت، والأفلام الوثائقية، والقصص الإعلامية في تبديد غموض الذكاء الاصطناعي. عندما يفهم الناس آليات استخلاص الخوارزميات للاستنتاجات، يصبح بإمكانهم المطالبة بممارسات أفضل، ومساءلة صنّاع القرار، والمساهمة في صياغة مستقبل رقمي أكثر عدلاً.

#### 7.5 الاعتبارات البيئية

يستهلك تدريب النماذج المعقدة وتشغيلها طاقة كبيرة، مما يساهم في البصمة البيئية للذكاء الاصطناعي. ورغم أن أساليب التفسير قد تزيد من التكاليف الحسابية، إلا أنها يمكن أن تخدم الاستدامة بشكل غير مباشر. فمن خلال إيضاح الخصائص الأكثر تأثيرًا، قد تكشف التفسيرات عن فرص لتبسيط النماذج دون التضحية بالدقة. وتتطلب النماذج الأصغر والأكثر تركيزًا موارد حوسبة أقل، فتخفّض استهلاك الطاقة وتقلل انبعاثات الغازات الدفيئة.

يبرز هنا النقاء مفهوم التفسير مع الوعي البيئي، حيث يتعدى الذكاء الاصطناعي الأخلاقي نطاق العدالة والشفافية ليشمل حرصنا على موارد الكوكب. فعبر الموازنة بين الدقة وقابلية التفسير والكفاءة، نقرب من أنظمة تكنولوجية مستدامة حقًا.

#### 7.6 المنظور العالمي والتباينات الثقافية

تتوافق التفسيرية مع منظومات القيم والنظم القانونية المتباينة عالميًا. ففي بعض المجتمعات، تستند توقعات الشفافية إلى تاريخ طويل من الدفاع عن حقوق المستهلكين والمساءلة الحكومية. وفي مجتمعات أخرى، قد تشكل مفاهيم مثل الرفاهية "الجماعية أو سيادة البيانات محددات لما ينبغي أن يكون عليه" التفسير.

مع المعايير المحلية والقوانين السائدة ألا يؤدي تبني الذكاء الاصطناعي إلى تعميق أوجه XAI يضمن تكييف أساليب اللامساواة العالمية. فبعض المناطق قد تفضل تفسيرات بسيطة وسردية تتوافق مع أساليب التواصل المحلية، بينما تفضل أخرى مؤشرات كمية. احترام هذه الاختلافات يتيح للذكاء الاصطناعي أن يكون قوة خير، ويعزز الثقة والتماسك الاجتماعي بدلًا من الانقسام أو الصراع.

#### 7.7 المشهد التنظيمي الدولي

في (GDPR) "تحتل قابلية التفسير مكانة بارزة في الأطر التنظيمية المتطورة. فقد فسّر "النظام العام لحماية البيانات المقترح (AI Act) "الاتحاد الأوروبي باعتباره يمنح الأفراد حقًا في التفسير، ويؤكد مشروع "قانون الذكاء الاصطناعي في الاتحاد الأوروبي على المساءلة والشفافية. كما تتناول مناطق أخرى—مثل الولايات المتحدة بمقترحاتها الفدرالية

الوليدة حول الذكاء الاصطناعي، أو مبادئ حوكمة الذكاء الاصطناعي في الصين—كيفية ترسيخ قابلية التفسير في سياسات قابلة للتنفيذ.

ومع تبلور هذه اللوائح، على المؤسسات التنبؤ بمتطلباتها. قد يستلزم الامتثال توثيق كيفية إنتاج التفسيرات، وضمان قابلية تدقيق النماذج، وتوفير آليات اعتراض للأفراد الذين يعتقدون أن قرارًا مدفوعًا بالذكاء الاصطناعي أضرّ بهم. ويؤكد التفاعل بين اللوائح المحلية وسلاسل التوريد العالمية أو الأعمال متعددة الجنسيات على تعقيد تحقيق حلول ترضي الجميع.

### 7.8 التعاون متعدد التخصصات

لا يمكن معالجة الأبعاد الأخلاقية والاجتماعية للذكاء الاصطناعي حصراً بواسطة التقنيين. إذ يتعين على خبراء الأخلاقيات والقانون وعلم الاجتماع وعلم النفس والنشطاء وقادة المجتمع العمل جنباً إلى جنب مع المهندسين وعلماء البيانات. يضمن هذا النهج متعدد التخصصات أن تتوافق التفسيرات مع الاهتمامات البشرية الحقيقية، وأن تستجيب التدخلات للأفراد الأكثر تأثراً.

على سبيل المثال، قد يوضح خبير قانوني كيفية تقديم منطوق النموذج بما يتوافق مع الإجراءات القانونية الواجبة، في حين XAI يشرح عالم اجتماع كيفية تأثير التفسيرات على تصورات الرأي العام في سياق ثقافي محدد. بهذه الطريقة، يصبح مشروعاً تعاونياً، يمزج خبرات من مجالات متنوعة لإنتاج نتائج قوية وشاملة وقابلة للتكيف.

### 7.9 معالجة الفجوة الرقمية

الواعدة، لا تحظى جميع المجتمعات بنفس الفرص للوصول إلى الموارد التعليمية والبنية التحتية اللازمة XAI رغم آفاق للتفاعل معه. ففي المناطق ذات الموارد المحدودة، أو حيث اعتماد التقنية الرقمية ضعيف، قد تتفاقم اللامساواة القائمة بسبب الأنظمة الغامضة. بدلاً من ذلك، يمكن للتفسيرية المتاحة، والسياقية، والمكيفة مع ظروف المجتمعات المختلفة أن تضمن ألا تكون هذه المجموعات تحت رحمة خوارزميات غير مرئية تتخذ قرارات حاسمة بشأن توزيع الموارد والفرص التعليمية أو الرعاية الصحية.

ومنهجياته—مع تقديم تفسيرات بلغات متعددة، واستخدام استعارات ثقافية ملائمة، وتوفير XAI يضمن توفير أدوات خيارات استخدام دون اتصال أو ذات نطاق ترددي منخفض—جسر الفجوة الرقمية. هذه الديمقراطية للتفسيرية تتيح للمجتمعات المهمشة المساواة والتحسين والاستفادة من أنظمة الذكاء الاصطناعي، مما يعزز الشمول والإنصاف على نطاق عالمي.

### 7.10 التطور المجتمعي على المدى الطويل

مع تطوّر أنظمة الذكاء الاصطناعي، ستظل التفسيرية محور النقاشات حول مستقبل العمل والحوكمة والنظام الاجتماعي. فالأنظمة الشفافة تعزز الثقة حتى مع انتشار اتخاذ القرارات المعتمد على الذكاء الاصطناعي. وعلى المدى البعيد، قد تطوّر المجتمعات معايير قياسية وأفضل الممارسات ومدونات سلوك مهنية مكرّسة للتفسير. وقد تظهر شركات تدقيق متخصصة في المناهج الأساسية عبر تخصصات متعددة XAI للتحقق من جودة التفسيرات، في حين تدمج الجامعات مبادئ

يشير هذا التصور طويل الأمد إلى مستقبل تصبح فيه التفسيرية عنصرًا تأسيسيًا في تصميم ونشر الذكاء الاصطناعي، لا استباقيًا، توجه المجتمعات تطوّر الذكاء الاصطناعي نحو مسار أخلاقي ومسؤول، XAI خيارًا اختياريًا. ومن خلال تبني لضمان بقاء التقنية أداة للتقدم الجماعي، لا مصدرًا لتقييد القدرات الإنسانية

### الخلاصة

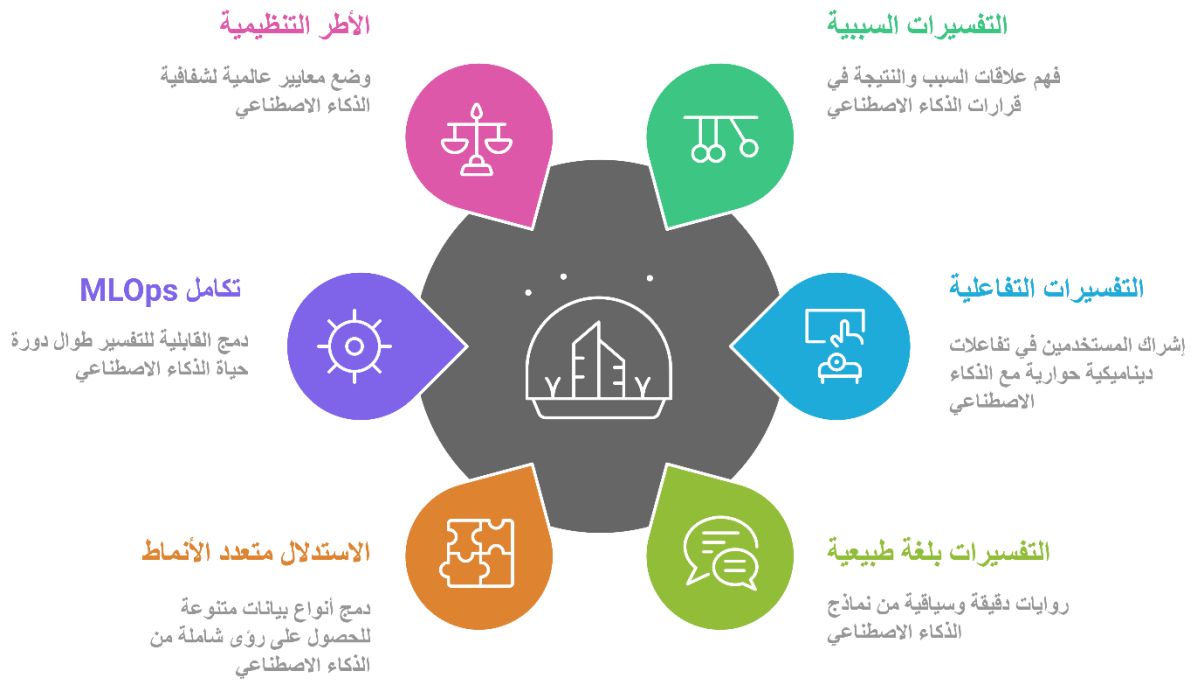
يقف الذكاء الاصطناعي القابل للتفسير عند نقطة تقاطع حاسمة حيث تلتقي الإمكانيات التقنية بالضرورة الأخلاقية. عبر بناء أطر أخلاقية، وحماية الخصوصية، ورصد التحيزات وتصحيحها، وإشراك الجمهور، ومراعاة القيود البيئية، واحترام إسهام أنظمة الذكاء الاصطناعي بشكل إيجابي في رفاهية المجتمعات XAI التنوع الثقافي، يضمن

وفي عالم يتزايد فيه تأثير الذكاء الاصطناعي على مصائر الأفراد وبنى الاقتصاد والتوازنات الجيوسياسية، تصبح الالتزام بالتفسير مسؤولية أخلاقية. ومن خلال التعاون متعدد التخصصات، والتنظيم الواعي، والحوار المجتمعي المستمر، يمكننا توجيه الذكاء الاصطناعي نحو مستقبل تحكمه العدالة والإنصاف واحترام الكرامة الإنسانية

## الفصل الثامن: نظرة إلى الأمام—مستقبل الذكاء الاصطناعي القابل للتفسير

مهيباً للتوسع في التعقيد والتأثير. ومع تطوّر أساليب التفسير، سنشهد (XAI) يبدو مشهد الذكاء الاصطناعي القابل للتفسير مواجهاً أكبر مع الأولويات العملية في العالم الواقعي، وازدياد الرقابة التنظيمية صرامة، واندماجاً أعمق في العمليات اليومية للمؤسسات والشركات. وبعيداً عن كونه ميزة إضافية اختيارية، ستتحول التفسيرية إلى مبدأ أساسي يوجّه آليات تصميم الأنظمة الذكية ونشرها وحوكمتها

### مستقبل الذكاء الاصطناعي القابل للتفسير (XAI)



## 8.1 التوجّهات الناشئة

عند حدٍ معيّن؛ فالبحوث والابتكارات المستمرة توسّع حدود ما يمكن أن تحقّقه التفسيرات XAI لا يقف مجال

### • التفسيرات السببية

بينما تركّز الأساليب الحالية غالبًا على إبراز العلاقات الارتباطية، قد تتيح النماذج المستقبلية تحديد العلاقات السببية بوضوح. لن يقتصر الأمر على معرفة "ما" أثر في القرار، بل "لماذا" هو صحيح ضمن إطار السببية. سيطلق هذا النهج مزيدًا من الرؤى القابلة للتطبيق، ليتمكّن صناع السياسات والأطباء والمخطّطون الاستراتيجيون من تنفيذ تدخلات دقيقة بثقة أكبر.

### • التفسيرات التفاعلية

ستصبح التفسيرات أكثر حوارية. بدلًا من الاكتفاء بمشاهدة رسوم ساكنة أو ملخصات نصية، قد يتمكن المستخدمون من "سؤال" النموذج أسئلة إضافية. فقد يسأل معلم يستخدم منصة تعليمية مدعومة بالذكاء الاصطناعي: "كيف سيتغيّر أداء الطلاب إذا خصصنا أسبوعًا إضافيًا لتدريس الكسور؟" فيجيب النموذج، مما يساعد المعلمين على تعديل خطط الدروس ديناميكيًا.

### • التفسيرات باللغة الطبيعية

بفضل تطوّر تقنيات توليد اللغة الطبيعية، ستمكّن النماذج من تقديم تفسيرات بلغات بشرية سلسلة وسياقية. تخيل نظامًا طبيًا يشرح منطقه كما لو كان طبيبًا كبيرًا يوجّه طبيبًا مبتدئًا، موضّحًا كيف تفاعلت الأعراض والنتائج المخبرية والتاريخ الطبي لتوجيه التشخيص الموصى به.

### • التفسير متعدد الوسائط وعبر المجالات

مع معالجة الذكاء الاصطناعي للنصوص والصور والصوت وبيانات الحساسات معًا، ستتطلب التفسيرات دمج أنماط بيانات متعددة في ملخصات مترابطة وسهلة الفهم. على سبيل المثال، يمكن تفسير منطوق اتخاذ القرار لدى سيارة ذاتية القيادة عبر إبراز العناصر المرئية التي دفعتها للفرملة، مع شرح لفظي: "أبطأت السرعة لأنني رصدت أحد المارة يعبر الطريق".

## 8.2 MLOps الاندماج في دورة حياة النماذج وعمليات الـ

(الذي يوحد تطوير أنظمة التعلم الآلي ونشرها وصيانتها) التفسيرية في كل مرحلة MLOps سيتبنّى مجال الـ

### • جمع البيانات ومعالجتها مسبقًا

ستساعد الأدوات التي تكشف التحيزات وتفسر الشذوذات في البيانات المهندسين على ضمان توازن مجموعات التدريب وتمثيليتها وخلوها من التحيز قبل البدء في النمذجة.

### • تدريب النماذج والتحقق منها

واختيار الخصائص وبنية النموذج. (Hyperparameters) أثناء التطوير، ترشد التفسيرات ضبط المعاملات يمكن للمهندسين استبعاد الخصائص غير المفيدة، والتحقق من توقعات الخبراء الميدانيين، والتأكد من أن التحسينات في الدقة لا تأتي على حساب الثقة أو العدالة.

### • النشر والمراقبة

بعد إطلاق النماذج، تتيح التفسيرات المستمرة معرفة تأثير التغيّرات الواقعية—مثل تحوّل الخصائص

الديموغرافية أو ظروف السوق أو القيود التنظيمية—على سلوك النموذج. عندما تكشف التفسيرات أن خاصية كانت فعالة في السابق لم تعد ذات مغزى، قد يكون الوقت مناسباً لإعادة تدريب النموذج أو إعادة ضبطه

#### • الإيقاف والاستبدال

في النهاية، تصبح جميع النماذج قديمة. تُبرر التفسيرات متى ولماذا ينبغي التخلص من نموذج ما. فعندما يلاحظ أصحاب المصلحة اعتماد النموذج على أنماط عفا عليها الزمن، يمكنهم استبداله بمسؤولية، وبحجج تستند إلى الأدلة.

### المشهد التنظيمي وتأثير السياسات 8.3

مع تزايد وعي المشرّعين بقوة القرارات المعتمدة على الذكاء الاصطناعي ومخاطرها، ستكثر القوانين المتعلقة بقابلية التفسير:

#### • التنسيق الدولي

على المدى الطويل، قد تظهر معايير عالمية للتفسير، مما يقلل التشتت بين الأنظمة القضائية المختلفة. سيؤدي ذلك إلى تبسيط الامتثال وخلق منافسة عادلة بين الشركات متعددة الجنسيات

#### • خدمات التصديق والتدقيق

قد نشهد ظهور شركات تدقيق متخصصة للتحقق من جودة التفسيرات ودقتها. قد تصدر هذه الجهات شهادات شبيهة بالتدقيق المالي، مما يطمئن العملاء والجهات التنظيمية بأن النماذج تلبّي معايير عالية للشفافية

#### • إرشادات متخصصة بالصناعات

قد تطوّر كل صناعة—مثل الرعاية الصحية والتمويل والطاقة والنقل—أفضل الممارسات الخاصة بها ونظم تقييم للتفسير. على سبيل المثال، قد يضع اتحاد مالي معايير تحدد مستوى التفاصيل والزمن الذي يجب أن يحصل فيه مقدّم طلب القرض على التفسير

### التعاون مع الخبراء الميدانيين 8.4

على حوار مثمر بين مختصي الذكاء الاصطناعي والخبراء في المجال XAI يعتمد نجاح

#### • التشاركية في إنشاء التفسيرات

بدلاً من تخمين مهندسي الذكاء الاصطناعي احتياجات المستخدمين من التفسيرات، يعمل الخبراء الميدانيون—كالأطباء والقضاة والمحللين الماليين—مباشرة على تصميم واجهات التفسير. يضمن هذا النهج التكراري دقة التفسيرات وملاءمتها

#### • أنماط تفسيرات مخصصة

في الطب، قد يرغب الطبيب بمراجع علمية مرفقة بمنطق النموذج. وفي القانون، قد يفضل القاضي تفسيرات مؤطرة في سياق السوابق القضائية. أما في التصنيع، فقد يقدر المهندسون تفاصيل على مستوى المستشعرات. من خلال مواءمة التفسيرات مع منطق المجال، يمكن دمج رؤى الذكاء الاصطناعي بسهولة أكبر في عمليات صنع القرار.

## نحو مجتمع متمكّن من فهم الذكاء الاصطناعي 8.5

يبقى تمكين الجمهور من فهم التساؤلات حول قرارات الذكاء الاصطناعي أمرًا ضروريًا

### • الدمج في المناهج التعليمية

ستدمج المدارس والجامعات مفاهيم التفسيرية في المناهج القياسية. ستصبح مقررات أخلاقيات الذكاء الاصطناعي والثقافة الرقمية وشفافية النماذج شائعة مثل أساسيات علوم الحاسوب. وسيتخرج الطلاب إلى سوق العمل وهم يتوقعون—لا يطالبون فقط—أن تكون أدوات الذكاء الاصطناعي مصحوبة بتفسير مفهم

### • منصات المعرفة العامة

قد تطوّر المتاحف ومراكز العلوم والمنصات الإلكترونية معارض ودروس تفاعلية تشرح كيفية "تفكير" الأنظمة الذكية. تصوّر منشأة عامة حيث يستطيع الزوار ضبط مدخلات النموذج ورؤية كيفية تغير التفسيرات فورًا. هذا التعلم التشاركي ينمي مواطنين واعين

### • المشاركة المدنية

مع ازدياد وعي الجمهور بالذكاء الاصطناعي، ستتعاظم المطالب المجتمعية بالشفافية. قد يكافئ الناخبون السياسيين الذين يدعمون الذكاء الاصطناعي المسؤول، وقد يفضل المستهلكون الشركات المعروفة بنماذج شفافة. وبمرور الوقت، سيعزز الضغط الشعبي التفسيرية كمعيار اجتماعي وحمية اقتصادية

## البحث متعدد التخصصات وأدوار جديدة 8.6

على علماء البيانات وحدهم. سيساهم باحثون من العلوم الإنسانية والاجتماعية والقانون والتصميم XAI لا يقتصر مستقبل والتواصل في تحديد ماهية "التفسير الجيد" وكيف تشكل التفسيرات الإدراك العام والمعايير المرجعية للتفسير

### • التفاعل بين الإنسان والحاسوب وتصميم التجارب

سيسمخ خبراء التفاعل وتجربة المستخدم واجهات تفسيرية بديهية تمكّن المستخدمين من استكشاف منطق النموذج بسهولة. تتيح واجهات التطبيقات ولوحات التحكم المصممة بعناية تقليل تعقيد العمليات المنطقية للنماذج وتحويلها إلى تجارب تفاعلية

### • علم النفس والمعرفة الإدراكية

سيساعد فهم طريقة تفسير المستخدمين للتفسيرات وثقتهم بها في توجيه تطوير أساليب تفسير تتوافق مع الإدراك البشري. وستساعد رؤى علم النفس المعرفي على اختيار تقنيات تفسيرية تقلل الالتباس والعبء الذهني

### • التكيف الثقافي واللغوي

مع تقديم الأنظمة الذكية لجمهور عالمي متنوع، يجب تكيف التفسيرات لتلائم الفروق اللغوية والقيم الثقافية. سيتطلب ذلك ترجمة المصطلحات التقنية إلى لغة بسيطة، واستخدام استعارات ثقافية مألوفة، الأمر الذي سيمثل مجالًا بحثيًا وتطبيقيًا نشطًا

## الأنظمة الذكية تفسّر بعضها البعض 8.7

في عالم قد تتواصل فيه النماذج الذكية مع بعضها البعض—متعاونة في إدارة سلاسل التوريد، أو التنسيق في أنظمة

المروء، أو التفاوض في الأسواق المالية—قد تتطوّر أطر التفسير لمساعدة هذه النماذج على فهم وتفسير منطق بعضها البعض. هذا "التفسير بين الوكلاء الذكيين" سيضمن استقرار النظم المعقدة وقدرة البشر على التدخل عند الضرورة

## 8.8 الاستدامة وإدارة الموارد

مع مواجهة التحديات البيئية العالمية، يمكن للتفسيرية دعم جهود الاستدامة

### • تحسين استهلاك الطاقة

عبر الكشف عن الخصائص أو النماذج الفرعية الأكثر استهلاكًا للطاقة، توجه التفسيرات المطورين لتبسيط البنى وتقليل التعقيد غير الضروري، واختيار أساليب استدلال موفرة للطاقة، مما يقلل البصمة الكربونية للذكاء الاصطناعي.

### • تحليلات دورة الحياة

المستقبلية تحليلات لدورة حياة استخدام الطاقة في النموذج. توضيح MLOps قد تتضمن خطوات عمليات ال- أنماط استهلاك الطاقة أو مكاسب الكفاءة يمكن أن يساعد الأطراف المعنية في تحديد متى وكيف يعتمدون حلول الذكاء الاصطناعي الأكثر صداقة للبيئة.

## 8.9 خاتمة حول مستقبل XAI

يشير مسار التفسيرية إلى عالم لن يكون فيه الذكاء الاصطناعي الشفاف مفهومًا مترقًا، بل ضرورة حتمية. مع تقدم أساليب ، ستقدم رؤى سببية أعمق، وتفاعلات حوارية ولغوية طبيعية، وستندمج بسلاسة في دورات حياة النماذج ضمن XAI سيضع الإطار التنظيمي معايير أكثر صرامة، في حين يضمن التعاون متعدد MLOps معايير تتطوّر لها صناعة ال- التخصصات والنطاقات أن تكون التفسيرات متينة تقنيًا وملائمة للسياقات المتنوعة

في هذا المستقبل، سيكفل الطلب العام على الفهم الكيفية التي تُصمم بها أدوات الذكاء الاصطناعي وتُنشر وتُدرس. ومع ازدياد وعي الناس بالذكاء الاصطناعي، سيصرّون على تكنولوجيا موثوقة وقابلة للمساءلة. والنتيجة النهائية تحوّل ثقافي: تصبح التفسيرية جزءًا لا يتجزأ من بناء الثقة وضمان الامتثال الأخلاقي وتشجيع الابتكار المسؤول

هي قصة تحسن مستمر. وبمرور الوقت، ستصبح النماذج القابلة للتفسير هي القاعدة، موجّهة XAI باختصار، قصة طريقتنا في التعامل مع الأنظمة الذكية في كل المجالات—من الرعاية الصحية والتمويل إلى التعليم واستراتيجيات المناخ. ومع نضج الأساليب وتأقلم المجتمعات، سيصبح شرح قرارات الذكاء الاصطناعي أمرًا طبيعيًا ومتوقعًا ولا غنى عنه، تمامًا كما نتحقّق من أوراق اعتماد أي خبير بشري نعتمد عليه



## الفصل التاسع: التدريبات العملية والمشاريع الختامية

الخطوة الأولى فقط. فالتمكن الحقيقي يأتي عبر التجربة (XAI) يمثل فهم النظرية وراء الذكاء الاصطناعي القابل للتفسير العملية—تجربة الأساليب، وحل المشكلات، وصقل التقنيات حتى تتوافق مع احتياجات الأطراف المعنية. يقدم هذا الفصل مجموعة واسعة من التمارين والمشاريع، بدءاً من المهام الأساسية ووصولاً إلى سيناريوهات معقدة تتناول مجالات ، ومساعدتك على تجنب الأخطاء الشائعة، وتطوير قدرة XAI متخصصة. تهدف هذه الأنشطة إلى ترسيخ فهمك لمبادئ الحكم لاختيار الاستراتيجيات التفسيرية الأنسب للسياق المحدد

### مسارات التعلم في الشرح



## تدريبات أساسية على التفسير 9.1

### التمرين 1: تفسير نموذج خطي بسيط

#### • المهمة:

اختر مجموعة بيانات عامة، مثل مجموعة بيانات أسعار المنازل في بوسطن، ودرب نموذج انحدار خطي للتنبؤ بأسعار المنازل اعتماداً على خصائص مثل عدد الغرف والمسافة إلى مراكز العمل ومعدلات الجريمة المحلية.

#### • الهدف:

افحص معاملات النموذج وفَسِّر معناها. على سبيل المثال، إذا كانت معامل "عدد غرف النوم" موجباً، فهل يتوافق ذلك مع الفكرة الشائعة بأن زيادة عدد الغرف تزيد عادةً من قيمة المنزل؟

#### • خطوات إضافية:

- تصوّر فترات الثقة للمعاملات، وفكّر في مدى استقرار هذه التفسيرات عبر تجارب متعددة بتقسيمات بيانات مختلفة.
- قارن التفسيرات مع المعرفة الميدانية—مثلاً راجع تقارير عقارية محلية، أو ناقش النتائج مع محترف في مجال العقارات.

#### • التأمل:

- هل كانت المعاملات بديهية؟
- هل تصرفت أي خاصية بشكل غير متوقع؟
- إذا بدت النتائج غير منطقية، تحقق من جودة البيانات أو تداخل الخصائص أو صغر حجم البيانات.

### التمرين 2: التحقق من التفسيرات المحلية مقابل التفسيرات العالمية

#### • المهمة:

استخدم نموذج شجرة قرارية بسيطة على مجموعة بيانات صغيرة (مثل التنبؤ بنجاح طالب في دورة دراسية)، ثم فسّر بنية الشجرة.

#### • الهدف:

حدّد القواعد العالمية بقراءة الشجرة كاملة، ثم ركّز على مسار قرار واحد لتفسيره محلياً. قارن بين المنطق العالمي للشجرة والمنطق المحلي لمسار واحد.

#### • خطوات إضافية:

- قلمّ الشجرة لترى إن كان نموذج أبسط يوفر تفسيرات أوضح دون فقدان كبير في الدقة.
- اطلب رأي خبير ميداني (مثل معلم) لمعرفة ما إذا كان المنطق يتوافق مع التوقعات.

#### • التأمل:

- هل توافق الهيكلية العالمية للشجرة القرارية مع القرارات المحلية؟
- هل بدت بعض القواعد شديدة التفصيل أو متداخلة؟

## 9.2 (Post-Hoc) أساليب التفسير بعد التدريب

### للمناذج المعقدة LIME التمرين 3: استخدام

#### • المهمة:

للتنبؤ بالتخلف عن السداد الائتماني باستخدام مجموعة بيانات (Random Forest) درّب نموذج غابة عشوائية لتفسير تنبؤ حالة واحدة LIME مالية. طَبِّقْ

#### • الهدف:

"حدّد أهم الخصائص التي دفعت التنبؤ نحو "تخلف" أو "عدم تخلف".

#### • خطوات إضافية:

- هل يؤثر زيادة أو تقليل عدد LIME جَرَب استخدام أعداد مختلفة من الخصائص في تفسير الخصائص المعروضة على وضوح التفسير؟
- افحص حالات متعددة من شرائح مختلفة من السكان (مثل أصحاب الدخل المرتفع مقابل أصحاب الدخل متناصفة LIME المنخفض) لترى إن كانت تفسيرات

#### • التأمل:

- مع الحدس أو المعرفة الميدانية؟ LIME هل تتوافق تفسيرات
- LIME لا تستطيع subtle إذا بدت التفسيرات غريبة، ففكّر فيما إذا كان النموذج يستغل تفاعلات تقربها بنموذج خطي بسيط

### للحصول على رؤى عالمية ومحلية SHAP التمرين 4: تطبيق

#### • المهمة:

على مجموعة بيانات طبية، مثل التنبؤ بإعادة (Gradient Boosted Trees) درّب نموذج أشجار متدرجة لكل من مجموعتي التدريب والاختبار SHAP إدخال المرضى للمستشفى المصابين بأمراض مزمنة. احسب قيم

#### • الهدف:

لتحديد الخصائص المهمة عالمياً، ثم أنشئ مخططات القوة (Summary Plot) ملخص SHAP أنشئ مخطط لعدد من الحالات الفردية (Force Plots)

#### • خطوات إضافية:

- قارن أهمية الخصائص العالمية مع الإرشادات الطبية المعروفة. هل الخصائص الأهم متوافقة مع المعرفة الإكلينيكية؟
- لرؤية كيف يتغير تأثير كل خاصية عبر نطاق قيمها SHAP أنشئ مخططات اعتماد
- إذا أمكن، استشر طبيباً للتحقق مما إذا كانت هذه الأنماط تتوافق مع الخبرة الإكلينيكية

#### • التأمل:

- هل توجد مجموعات مرضى تختلف تفسيراتهم المحلية بشكل كبير عن الأنماط العالمية؟
- هل جعلك ذلك تشكك في تمثيلية البيانات أو قابلية تعميم النموذج؟

### لتقديم توصيات قابلة للتطبيق (Counterfactual) التفسيرات المضادة 9.3

#### التمرين 5: التفسيرات المضادة في التعليم

##### • المهمة:

طوّر نموذج تصنيف (انحدار لوجستي أو شجرة بسيطة) للتنبؤ بنجاح الطالب في دورة دراسية. أنشئ تفسيرات "مضادة تظهر الحد الأدنى من التغييرات المطلوبة لتحويل تنبؤ "فشل" إلى "نجاح

##### • الهدف:

حدّد خطوات قابلة للتنفيذ، مثل زيادة ساعات الدراسة، تحسين الحضور، أو طلب المساعدة في موضوعات محددة.

##### • خطوات إضافية:

- قدّم هذه التفسيرات المضادة لمعلم واطلب منه تقييم مدى واقعية هذه التدخلات
- جرّب فرض قيود: ماذا لو لم يستطع الطالب تحسين نسبة حضوره بسبب التزامات عمل؟ هل توجد تفسيرات مضادة بديلة؟

##### • التأمل:

- هل اقترحت التفسيرات المضادة تدخلات قابلة للتطبيق، أم كانت غير واقعية (مثل "زيادة الدخل بمقدار 10,000 دولار" لطالب)؟
- فكّر في العنصر الإنساني: هل سيجد الطالب هذه الاقتراحات مشجّعة أم محبطة؟

#### التمرين 6: توليد التفسيرات المضادة في التمويل أو الرعاية الصحية

##### • المهمة:

بالنسبة لسيناريو قبول القروض أو التشخيص الطبي، أنشئ تفسيرات مضادة تظهر كيف يمكن لتغيير طفيف في الخصائص (مثل تخفيض ضغط الدم أو تحسين نسبة استخدام البطاقة الائتمانية) أن يغيّر قرار النموذج

##### • الهدف:

قيّم ما إذا كانت هذه التفسيرات المضادة مناسبة أخلاقياً ومفيدة. على سبيل المثال، في الرعاية الصحية، قد يكون اقتراح خفض الكوليسترول قابلاً للتنفيذ لكن قد يعتمد على عوامل اجتماعية اقتصادية

##### • التأمل:

- هل التغييرات المقترحة عادلة في ضوء ظروف الفرد؟
- هل قد تشير بعض التفسيرات المضادة عن غير قصد إلى خصائص حساسة؟

#### بناء لوحات قيادة وواجهات سهلة الاستخدام 9.4

#### التمرين 7: لوحة تفاعلية للتفسير

### • المهمة:

، أنشئ واجهة تفاعلية InterpretML أو لوحة تحكم Streamlit أو Plotly Dash باستخدام أدوات مثل تعرض تنبؤات النموذج وأهمية الخصائص عالمياً، إضافة إلى تفسيرات محلية. أضف منزلقات أو قوائم منسدلة لتغيير مدخلات الخصائص ورؤية كيف تتغير التنبؤات والتفسيرات

### • الهدف:

شارك هذه اللوحة مع شخص غير تقني (مثل مدير أو معلم) واجمع ملاحظاته

### • خطوات إضافية:

- أضف معلومات إرشادية ونوافذ توضيح وأدلة تعليمية لمساعدة المستخدمين لأول مرة
- نفذ مرشحات تتيح للمستخدمين التركيز على مقاطع بيانات معينة أو مخرجات محددة ("أرني فقط المرضى فوق 60 عاماً")

### • التأمل:

- هل وجد المستخدمون اللوحة بديهية أم مربكة؟
- ما التحسينات التي اقترحوها—رسوم أبسط، مصطلحات أقل تقنية، تفسيرات أكثر سردية؟

## مشاريع ختامية متخصصة 9.5

### مشروع الرعاية الصحية

### • المهمة:

، ثم SHAP وLIME للتنبؤ بإعادة إدخال المرضى للمستشفى. أنشئ تفسيرات (مثل XGBoost) درّب نموذج طور سيناريو تفسيري مضاد لاقتراح تدخلات (مثل فحوصات إضافية أو التزام بالأدوية)

### • التفاعل مع أصحاب المصلحة:

قدّم النتائج للعاملين في القطاع الطبي واجمع ملاحظاتهم حول الوضوح والدقة والملاءمة الإكلينيكية

### • التأمل:

- هل زادت ثقة الطاقم الطبي بمنطق النموذج بعد رؤية التفسيرات؟
- كيف ساعدت ملاحظاتهم في فهم ما هو "تفسير جيد" في مجال الرعاية الصحية؟

### مشروع التمويل

### • المهمة:

أنشئ نموذجًا لتقييم الائتمان وقدم تفسيرات عالمية ومحلية للقروض المقبولة والمرفوضة. استخدم هذه التفسيرات لصياغة تقرير امتثال يلبي المعايير التنظيمية

### • التفاعل مع أصحاب المصلحة:

اعرض التقرير على مسؤولي الامتثال أو مدققي الحسابات الماليين. اسألهم إن كانت التفسيرات تلبّي متطلبات الشفافية والإنصاف

#### • التأمل:

- هل وجد فريق الامتثال التفسيرات كافية للمراجعات التنظيمية؟
- هل اكتشفت تحيزات أو أنماط غير متوقعة بحاجة لمعالجة؟

#### مشروع قانوني:

#### • المهمة:

LIME أو SHAP درّب نموذجًا لتصنيف الوثائق القانونية (عقود، براءات اختراع، أحكام قضائية). استخدم لتسليط الضوء على المقاطع النصية الدافعة للتصنيف

#### • التفاعل مع أصحاب المصلحة:

قدّم هذه المقاطع المميزة لخبراء قانونيين واسألهم إن كان منطق النموذج يتوافق مع طريقتهم في تفسير مثل هذه الوثائق.

#### • التأمل:

- هل زادت ثقة الخبراء القانونيين بعملية تصنيف النموذج بعد رؤية المقاطع المسلّط عليها الضوء؟
- هل كشفت هذه التجربة عن مصاعب خاصة بالمجال، مثل تركيز النموذج على عبارات غير ذات صلة؟

#### مجالات إضافية:

#### • المشتريات:

على نموذج يتنبأ بموثوقية الموردين. أظهر لموظفي المشتريات الخصائص المؤثرة XAI طبّق أساليب (معدلات التسليم في الوقت المحدد، القدرة الإنتاجية) في التوصيات. عدّل البيانات أو إعدادات النموذج استنادًا إلى ملاحظاتهم

#### • الخدمات اللوجستية:

أشّرح قرارات تحسين المسارات لمسؤول لوجستي. إذا اقترح النموذج مسار شحن جديد، وضّح العوامل (مثل أنماط الطقس والازدحام التاريخي) التي شكّلت اختياره. اطلب ملاحظات حول ما إذا كانت هذه الرؤى تحسّن التخطيط وتخصيص الموارد

#### • التسويق:

بالنسبة لحملة تسويقية تستهدف شرائح عملاء مختلفة، استخدم التفسيرات للتحقق من سبب اختيار شرائح معينة. اطلب من محلي التسويق تقييم ما إذا كان المنطق يتوافق مع أهداف العلامة التجارية والمبادئ الأخلاقية

### 9.6 المراجعة والتقييم

بعد إتمام هذه التمارين والمشاريع، فكّر في الآتي

#### • الفعالية:

هل ساعدتك هذه الأنشطة على فهم منطق نماذجك بشكل أفضل؟ هل رصدت تحيزات أو مشكلات في جودة البيانات أو تفاعلات غير متوقعة؟

#### • سهولة الاستخدام:

أي أساليب تفسيرية لاقت قبولا أكبر لدى أصحاب المصلحة؟ هل أربكت بعض الأساليب المستخدمين؟ هل سهّلت اللوحات التفاعلية والأدوات الحوارية الحوار الواضح؟

#### • القابلية للتنفيذ:

هل دفعتك التفسيرات إلى تغيير قرارات النمذجة، أو جمع البيانات، أو الاستراتيجيات الميدانية؟ على سبيل المثال، هل أعدت تدريب نموذج بعد اكتشاف خاصية مزعجة؟ هل عدل أصحاب المصلحة بروتوكولات اتخاذ القرار؟

#### • قابلية التوسع والصيانة:

المستمرة. هل يمكنك الحفاظ عليها MLOps فُكر في كيفية دمج هذه الأساليب التفسيرية في ممارسات الـ بمرور الوقت، أم تتطلب إعادة ضبط متكررة وتدقيقات مستمرة؟

### تطوير الحدس وأفضل الممارسات 9.7

مع تكرار المحاولات والمراجعات، ستلاحظ أنماطًا

- النماذج الأبسط أسهل تفسيرًا ولكن قد تفتقر لأحدث درجات الدقة.
- النماذج المعقدة تحقق أداءً استثنائيًا لكنها تتطلب أساليب تفسيرية دقيقة ومركزة.
- بعض الأساليب ممتازة في تلخيص السلوك العام، بينما يتألق البعض الآخر في تقديم رؤى محلية أو تفسيرات مضادة قابلة للتنفيذ.

بمرور الوقت، سيُرشدك هذا الخبرة نحو اختيار الأساليب الأنسب لمجالك وخصائص بياناتك وأولويات أصحاب المصلحة. سيمنحك الحدس المتنامي القدرة على صياغة تفسيرات لا توضح منطق النموذج فحسب، بل تمكن المستخدمين وتدعم القرارات المستنيرة وتعزز الثقة في الأنظمة الذكية.

#### الخلاصة

تحول الممارسة العملية المعرفة النظرية إلى خبرة تطبيقية. عبر تجربة أساليب التفسير المختلفة، وعرض النتائج على تضمن هذه XAI خبراء المجال، والتأمل في الملاحظات، ستكتسب المهارة والثقة اللازمة للتعامل مع المشهد المعقد للقاعدة العملية أنه عندما تمضي قدمًا، لن تفهم آليات عمل نماذج الذكاء الاصطناعي فحسب، بل ستوجهها أيضًا نحو عمليات نشر شفافة، عادلة، عالية التأثير.

## الملحق: موارد إضافية

### قراءات مكملة

- Molnar, C. (2019). *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*.  
يقدّم كتاب مولنار مرجعًا شاملاً حول المفاهيم الأساسية والأساليب وأفضل الممارسات في مجال التعلّم الآلي القابل للتفسير. يُعد هذا الكتاب مدخلاً سهلاً للوصول ودليلاً مرجعياً، ويشمل أمثلة برمجية ومرئيات تفاعلية.  
متوفر عبر الإنترنت: <https://christophm.github.io/interpretable-ml-book/>
- Lipton, Z. C. (2018). “The Mythos of Model Interpretability.” *Queue*, 16(3), 31–57.  
يفحص هذا البحث مصطلح "التفسيرية" بشكل نقدي، ويناقش معانيه المختلفة ويؤكد الحاجة إلى تعريفات ومقاييس دقيقة. يوفر أساساً مفاهيمياً للباحثين والممارسين الذين يواجهون تحديات في التفسير.  
النسخة الأولية: <https://arxiv.org/abs/1606.03490>
- Rudin, C. (2019). “Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead.” *Nature Machine Intelligence*, 1(5), 206–215.  
تؤكد رادين على أنه في التطبيقات الحرجة—مثل الرعاية الصحية والعدالة الجنائية—يفضّل استخدام نماذج قابلة للتفسير بطبيعتها بدلاً من النماذج المعقدة التي تتطلب تفسيرات لاحقة، داعيةً إلى تغيير أولويات البحث في هذا المجال.  
الناشر: <https://www.nature.com/articles/s42256-019-0048-x>
- Doshi-Velez, F., & Kim, B. (2017). “Towards A Rigorous Science of Interpretable Machine Learning.”  
arXiv: <https://arxiv.org/abs/1702.08608>  
يدعو هذا العمل التأسيسي إلى وضع تعريفات أكثر رسمية وأطر تقييم صارمة للتفسيرية، ويحث على تطوير معايير قياسية ومنهجيات موحدة لتقييم أساليب التفسير ومقارنتها.  
مؤتمر FAccT (Fairness, Accountability, and Transparency):  
الموقع: <https://facctconference.org/>  
مؤتمراً رائداً متعدد التخصصات يعرض أحدث البحوث حول العدالة والمساءلة والشفافية في FAccT يعد الأنظمة الخوارزمية. تغطي الأوراق البحثية مجموعة واسعة من الموضوعات، بما في ذلك الأطر الأخلاقية والمنظورات التنظيمية وأساليب التفسير المبتكرة.
- Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness and Machine Learning: Limitations and Opportunities*.  
الكتاب متوفر عبر الإنترنت: <https://fairmlbook.org/>



رغم التركيز على الإنصاف، يوقّر هذا المصدر سياقًا حول كيفية مساهمة التفسيرية في كشف التحيزات وضمان الاستخدام الأخلاقي للأنظمة الذكية.

- (2022) "مخطط البيت الأبيض الأمريكي لـ"شريعة حقوق الذكاء الاصطناعي

الوثيقة: <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>

يضع هذا الإطار مبادئ لتصميم ونشر أنظمة الذكاء الاصطناعي بما يحترم حقوق الإنسان، بما في ذلك الشفافية والتفسيرية كعناصر أساسية

#### موارد ويب ودروس تعليمية

- LIME توثيق:

**GitHub:** <https://github.com/marcotcr/lime>

على نماذج مختلفة، مع نصائح حول ضبط المعاملات LIME يشمل أمثلة ودفاتر تجريبية وتوجيهات لتطبيق وتفسير النتائج

- SHAP توثيق:

**GitHub:** <https://github.com/slundberg/shap>

في تصنيف المهام SHAP ، ومعرض تصوّرات يوضح قدرات Jupyter وثنائق شاملة مع أمثلة برمجية ودفاتر والانحدار والتعلّم العميق

- InterpretML (من مايكروسوفت):

**GitHub:** <https://github.com/interpretml/interpret>

يقدم لوحة تحكم تفاعلية وعدداً من المفسّرات المحايدة للنماذج والمتخصصة بالنماذج. يحتوي المستودع على تعليمات مفصلة ومجموعات بيانات أمثلة وأمثلة تكامل

- Captum (من Facebook AI Research):

الموقع: <https://captum.ai/>

يوفر الموقع دروساً مرئية ومرجعية PyTorch على تفسير الشبكات العصبية المبنية على Captum يركز لواجهة برمجة التطبيقات وأمثلة تعرض طرق الإسناد القائمة على التدرجات وتقنيات التصور

- Alibi (من Seldon):

**GitHub:** <https://github.com/SeldonIO/alibi>

وتدابير أهمية الخصائص. يحتوي Anchors يقدم مجموعة من أساليب التفسير تشمل التفسيرات المضادة والمستودع على دفاتر وأدلة للدمج في البيانات التشغيلية

- Fairlearn:

**GitHub:** <https://github.com/fairlearn/fairlearn>

تداخل الإنصاف مع التفسيرية والمقاييس Fairlearn رغم تركيزه على الإنصاف، تتناول الوثائق والأدوات في اللاحقة للشرح

- OpenAI Cookbook:

GitHub: <https://github.com/openai/openai-cookbook>

، يناقش هذا المستودع أحيانًا مفاهيم التفسيرية وتقنيات التعامل مع النماذج OpenAI رغم تركيزه على نماذج اللغوية المعقدة، ما يقدم رؤية حول آفاق التفسير المستقبلي.

## قوائم ومراجع إضافية

- Awesome-Explainable-AI (GitHub):

<https://github.com/wangyongjie-ntu/Awesome-Explainable-AI>

، يُحدثها المجتمع باستمرار لتعكس أحدث التطورات XAI قائمة تضم موارد وأوراقًا وأدوات ودروسًا حول

- Distill Pub:

<https://distill.pub/>

يقدم مقالات تفاعلية وغنية بالمرئيات لتبسيط مفاهيم التعلم الآلي. رغم أنه ليس مخصصًا للتفسيرية فقط، إلا أن يُظهر كيف يمكن للتفسيرات المرئية تسهيل فهم المنطق المعقد Distill أسلوب

## المصطلحات

- (XAI) الذكاء الاصطناعي القابل للتفسير:

تقنيات وأدوات تجعل قرارات نماذج الذكاء الاصطناعي مفهومة للبشر، غالبًا عبر مرئيات أو تقييم أهمية الخصائص أو نماذج بديلة مبسطة

- التفسير المحلي:

تفسير يركز على حالة فردية أو مجموعة صغيرة من الحالات. يوضح كيفية وصول المدخلات المحددة إلى مخرجات معينة، بدلًا من تلخيص منطق النموذج ككل

- التفسير العالمي:

تفسير يلتقط السلوك العام للنموذج عبر مجموعة البيانات الكاملة. يعطي نظرة حول الخصائص الأكثر تأثيرًا. ويصف تأثيرها على التنبؤات بالمتوسط

- (Model-Agnostic) تفسير محايد للنموذج:

طريقة للتفسير يمكن تطبيقها على أي نوع من النماذج دون الحاجة للوصول إلى بارامترات النموذج الداخلية. من أشهر هذه الطرق SHAP وLIME تعتبر

- (Feature Importance) أهمية الخصائص:

مقياس يشير إلى التأثير النسبي للمدخلات على تنبؤات النموذج. يمكن احتساب الأهمية عالميًا (عبر مجموعة البيانات) أو محليًا (لتنبؤ واحد)

- (Counterfactual Explanation) التفسير المضاد:

سيناريو افتراضي يوضح كيف سيؤدي تغيير بعض المدخلات إلى تغيير قرار النموذج. تساعد هذه التفسيرات على تحديد الإجراءات القابلة للتنفيذ لتحقيق النتائج المرجوة

- **Anchors:**

قواعد محايدة للنموذج تشرح التنبؤات الفردية عبر تحديد شروط الخصائص التي "ترسي" التنبؤ. تقدم بيانات شرطية إذا-فان بسيطة وعالية الدقة، مما يوفر شكلاً بديهيًا للتفسير المحلي Anchors

- **(PDP) مخطط الاعتماد الجزئي:**

تصور يوضح العلاقة بين خاصية معينة والمخرج المتوقع، مع تثبيت الخصائص الأخرى. يساعد على فهم إذا ما كان تأثير الخاصية خطيًا أو رتيبًا أو أكثر تعقيدًا

- **(ICE) مخطط التوقع الشرطي الفردي:**

ولكنه يركز على الحالات الفردية بدلًا من المتوسطات. يكشف التباينات داخل المجموعات PDP مشابه لـ الفرعية، مما يساعد على اكتشاف التفاعلات والفئات المتميزة

- **SHAP (SHapley Additive exPlanations):**

على الخصائص. توفر هذه (SHAP قيم) إطار يعتمد على نظرية الألعاب التعاونية لتوزيع درجات المساهمة القيم تفسيرًا عادلًا ومتسقًا محليًا وعالميًا، بخصائص نظرية مرغوبة

- **LIME (التفسيرات المحلية المحايدة للنموذج):**

تقنية تقرب النموذج محليًا حول حالة محددة بنموذج أبسط قابل للتفسير (غالبًا خطي)، ما يمكن من فهم تأثير كل خاصية على التنبؤ الفردي

---

## المراجع

مذكورة في نهاية الملحق، وتشمل أعمال مولنار وليبتون وريببيرو وغيرهم من الباحثين، إضافة إلى وثائق إرشادية (ولوائح دولية وإطارات تنظيمية)

## عن المؤلف

ياسر إسماعيل هو محترف متميز بخبرته المتعددة المجالات التي تجمع بين الذكاء الاصطناعي، المشتريات الاستراتيجية، وإدارة المشاريع. حاصل على درجات علمية متقدمة في علوم الحاسب وإدارة الأعمال و المحاسبه ، بالإضافة إلى شهادة مدير مشاريع مما يجعله حلقة وصل فريدة بين الجوانب التقنية والتشغيلية لتحقيق تحول مؤسسي ناجح محترف

طوال مسيرته المهنية، قاد ياسر مؤسسات في مجالات متنوعة لتطبيق حلول الذكاء الاصطناعي القابلة للتفسير والمبنية على أسس أخلاقية. وقد عززت خبرته الواسعة في إدارة المشتريات وسلاسل الإمداد، التي تجلت بوضوح في منصبه الأخير كمدير للمشتريات والعقود في قطاع النفط ، نهجه العملي في ضمان دمج الذكاء الاصطناعي بسلاسة ضمن معايير التوريد العالمية وأطر الامتثال وأفضل الممارسات. هذا النهج يدعم اتخاذ قرارات ذكية قائمة على البيانات، تعزز الكفاءة، وتحافظ على ثقة أصحاب المصلحة، وتراعي المسؤوليات القانونية والاجتماعية

تتضمن الخلفية القيادية لياسر توجيه فرق الذكاء الاصطناعي داخل الهياكل المؤسسية، والمساهمة في صياغة سياسات تنظيمية للتقنيات الناشئة، وقيادة مبادرات التحول الرقمي. يعتمد نهجه على الجمع بين التحليلات المتقدمة، النمذجة التنبؤية، والرؤى المستندة إلى البيانات مع منهجيات المشتريات الاستراتيجية لإنشاء سلاسل إمداد مرنة ومبنية على القيمة. ومن خلال تعزيز الحوار المستمر بين الخبراء التقنيين وصناع القرار التنفيذيين، يضمن ياسر تحويل استثمارات الذكاء الاصطناعي إلى مكاسب تشغيلية ملموسة

بالإضافة إلى مساهماته المهنية المباشرة، يكرس ياسر جهوده للريادة الفكرية والتعليم. يشارك معرفته من خلال التدريس، المحاضرات، والنشر، بهدف تمكين قادة الأعمال، وصناع السياسات، والتقنيين من تسخير إمكانيات الذكاء الاصطناعي بشكل مسؤول. تطلع ياسر يتمثل في مستقبل تتسم فيه الأنظمة الذكية بالشفافية، الإنصاف، والالتزام بالقيم الإنسانية، لتصبح أدوات متكاملة لتعزيز الكفاءة التشغيلية والتميز التنافسي

## الذكاء الاصطناعي القابل للتفسير" (XAI)

ليس مجرد تطور تقني، بل هو حجر الأساس لبناء الثقة في الأنظمة الذكية. في عصر تؤثر فيه القرارات المدفوعة بالذكاء الاصطناعي على حياة الأفراد والأعمال والمجتمعات، يصبح ضمان الشفافية والمساءلة أمراً بالغ الأهمية.



يتيح لنا الذكاء الاصطناعي القابل للتفسير سد الفجوة بين الخوارزميات المعقدة وفهم البشر، مما يعزز الابتكار بشكل أخلاقي وشامل ومتوافق مع المعايير العالمية."

— ياسر إسماعيل