

المقارنة بين النماذج التقليدية والحديثة للتقيب في البيانات  
أمل احمد طلعت مدرس الاحصاء بكلية التجارة بنات الأزهر

مقدمة:

يتناول هذا البحث بالدراسة المداخل المختلفة للتقيب في البيانات. ويتعلق ذلك بالنماذج المفترضة لحل المشاكل، وأساليب التقدير المستخدمة، وطرق اختبار صلاحية النماذج والمقارنة بينها لاختيار أكفاءها.

ويمكن تصنيف طرق التقيب في البيانات في نوعين أساسيين: الطرق المعلممية، والطرق الامعممية بالإضافة للطرق نصف المعلممية. لذلك فسوف يأتي هذا البحث في أربعة أجزاء: يُخصص الأول منها للجانب النظري؛ ويحتوي على تصميم البحث ومشكلته والدراسات السابقة. ويُخصص الثاني للطرق المعلممية؛ ويحتوي على نموذج الانحدار الخطي، وتحليل المكونات الرئيسية، والتحليل العاملية، وتحليل التمايز، وتحليل التمازن التمييزي، والتحليل العنودي، ونماذج البروبت واللوجيست. ويُخصص الثالث للطرق الامعممية؛ ويحتوي على العملية الهرمية التحليلية، وطريقة أقرب الجيران، والنماذج الجمعية المعممة، والبرمجة الرياضية، ودوال الانحدار المقسمة المتنوائمة، وشجرة القرارت، والشبكات العصبية، والانحدار الامعممي. كما يُخصص الفصل الرابع وهو الأخير للطرق نصف المعلممية؛ ويحتوي على الانحدار نصف المعلممي.

ويهدف البحث إلى التعريف بأساليب التقيب في البيانات وأنواعها ، كما يهدف إلى استعراض المجالات المختلفة التي استُخدمت فيها تلك الأساليب بنجاح. وقد اعتمد البحث على دراسة العديد من البحوث التي تتضمن أساليب مختلفة للتقيب في البيانات للتعريف بهذه الأساليب وطرق تقديرها وتطبيقاته وشروط استخدامها بشكل مبسط، ومقارنة بين أفضلية أساليب التقيب والأساليب التقليدية.

## ٢-١ مشكلة البحث وأهميته

التقيب في البيانات مجال مهم للحديث عنه في هذا البحث، وتحليل البيانات إحدى مراحله. وتحتوي برمجياته على أساليب حديثة، ولكن استخدام هذه البرامج وبالتالي هذه الأساليب في البحوث العربية ما زال في أضيق الحدود بسبب ندرة منشوراتها باللغة العربية وبالتالي صعوبة فهمها. لذا، فقد أحجم معظم الباحثين عن تلك الأساليب في توفيق العلاقة بين المتغير التابع والمتغيرات المستقلة مستعينين بنموذج الانحدار الخطي المتعدد لسهولة فهمه واستخدامه. غير أن التطبيقات الحديثة أثبتت ضعف مصداقية نموذج الانحدار الخطي المتعدد في توفيق معظم المشاكل المعاصرة التي تتسم باللاخطية وجود تفاعلات بين المتغيرات بفعل مجموعات البيانات الكبيرة.

### ٣- الدراسات السابقة

تعاملت بعض دراسات التقيب في البيانات مع عملية التقيب في البيانات [11][13:17]. وبينت الخطوات التي يجب أن تسير عليها المنشأة بغية اكتشاف المعرفة والأنماط الهامة التي لم تكن معروفة من قبل، كما أشار البعض منها لأساليب التقيب في البيانات و مجالات استخدامها. وقد اهتمت بعض الدراسات بحصر وتقديم تعريفات لهذه الأساليب، وركز البعض الآخر على دراسة الجوانب النظرية [20:32][6:10] والخوارزميات [35][3:5] التي تستخدمها تلك الأساليب في عملية التقدير، كما اهتم البعض الآخر بتقديم البرمجيات [34] وأدوات المساعدة لهذه الأساليب. وقد تبين من الدراسات التي اهتمت بتطبيقات أساليب التقيب في البيانات [21][23] أن أهم التطبيقات الناجحة كانت: تحليل سلة السوق، وتقدير الجدارة الائتمانية، وتحليل أسواق الأسهم، واكتشاف الغش والأعطال، وتشخيص الأمراض.

### ٤- أسئلة البحث

يجيب البحث عن الأسئلة التالية:

ما هي أساليب التقيب في البيانات؟

ومتى تُستخدم؟ وما هي التطبيقات الناجحة التي استُخدمت فيها تلك الأساليب؟

وقد تم عمل ذلك لكل من: نموذج الانحدار الخطي، وتحليل المكونات الرئيسية، والتحليل العاملی، وتحليل التمايز، وتحليل التناظر التمييزي، وتحليل العنقودي، ونمذج البرويت واللوجيست، والعملية الهرمية التحليلية، والأنظمة الخبرية، وطريقة أقرب الجيران، والنمذج الجمعية المعممة، والبرمجة الرياضية، ودوال الانحدار المقسمة المتوازنة، وشجرة القرارت، والشبكات العصبية، والانحدار الامامي، والانحدار نصف المعلمی.

### ٢. الأساليب المعلمية

#### ١- الانحدار الخطي المتعدد MLR

يُعد أسلوب الانحدار الخطي المتعدد Multiple Linear Regression أقدم وأشهر الأساليب الإحصائية التي استُخدمت في حقل التقيب في البيانات. وهو نموذج تنبؤي يلجأ إليه المحللون عند الرغبة في تقييم العلاقة السببية بين أحد المتغيرات الكمية وعدة متغيرات أخرى. ويطلق على المتغير الذي نريد تفسير التغيير فيه أو التنبؤ بقيمه في المستقبل عدة أسماء: المتغير التابع dependent variable، متغير الاستجابة response، أو المتغير المفسر explained variable ويأخذ الرمز  $y$ . كما يطلق على المتغيرات الأخرى اسم: المتغيرات المستقلة independent vr's، المتغيرات المفسرة explanatory vr's

المتباينات predictors، السمات features أو covariates وتأخذ الرمز  $x_i$ ؛ حيث  $i = 1, 2, \dots, n$  وهو ما يشير إلى المشاهدات بينما  $(1, 2, \dots, p) = l$  وهو ما يشير إلى المتغيرات. وأأخذ نموذج الانحدار الخطي المتعدد (بصيغة المصفوفات) الشكل التالي:

$$\begin{matrix} Y &= X & B + l \\ n \times 1 & n \times p & p \times 1 & n \times 1 \end{matrix} \quad (1)$$

ويتم التوصل لشكله التحليلي بتقدير متوجه المعالم  $B$  بإحدى طرق التقدير. وتعتبر طريقة المربيعات الصغرى LS أشهر هذه الطرق، حيث يتم اختيار مستوى يصغر مجموع مربعات الباقي  $l$ . ويمكن فحص جودة توفيق النموذج من خلال الأدوات التشخيصية برسم الباقي مقابل القيم المقدرة من خط الانحدار، ثم النظر إلى الشكل الناتج. فإذا كان الانحدار صادقاً، فإن قيم المتغير التابع يجب أن تتوزع حول الخط المقدر عشوائياً بدون أن تتشكل أي اتجاه عام واضح. كما يمكن فحص جودة توفيق نموذج الانحدار بالاعتماد على مؤشر تشخيصي يُعرف باسم معامل التحديد  $R^2$  الذي يأخذ قيمة تتراوح بين الصفر والواحد، إذ كلما اقتربت قيمته من الواحد، كلما دل ذلك على إمكانية التنبؤ بقيم  $l$  بشكل أصدق اعتماداً على العلاقة التي تربطها بقيم  $x_i$ . وأخيراً، يتم اختبار المعنوية الإجمالية للنموذج باستخدام اختبار F، وختبار المعنوية الجزئية للمتغيرات المستقلة باستخدام اختبار t.

إذا كان لدينا مجموعة بيانات واحدة (متغير كمي واحد وعدة متغيرات مستقلة) وكانت الأخيرة لا تعتمد على بعضها (بمعنى عدم وجود ازدواج خطى متعدد multicollinearity)، يمكن تطبيق الانحدار الخطي المتعدد بأمان. أما في حالة وجود مجموعتي بيانات (مجموعة للمتغيرات المستقلة ومجموعة للمتغيرات التابع) أو أكثر (مجموعة للمتغيرات المستقلة وعدة مجموعات للمتغيرات التابع) أو كان هناك ازدواج خطى في حالة مجموعة البيانات الواحدة، فإن الانحدار الخطي لا يصلح ويمكن تطبيق أحد الأساليب التالية:

## 2-2 تحليل المكونات الرئيسية PCA

لا بد - عند التعامل إحصائياً مع أي مشكلة - من التعبير عنها بما يسمى بمصطلحات التقريب في البيانات بجدول البيانات. وجدول البيانات data table هو عبارة عن مصفوفة من الدرجة  $p \times n$ ، تشير فيه الصور  $n$  إلى القياسات التي أخذتها وحدات المعاينة في  $p$  من المتغيرات الخاضعة للدراسة.

ويهدف تحليل المكونات الرئيسية Principal Components Analysis إلى ضغط جدول البيانات في ظل القياسات المرتبطة والتعبير عنه بمجموعة جديدة من المتغيرات غير المرتبطة (المتعامدة) وهو ما يُعرف باختزال الأبعاد. وعندئذ، يقال أن المتغيرات الجديدة تعتمد على السياق context أو أنها المكونات الرئيسية

أو العوامل factors أو المتجهات المميزة principal components أو المتجهات eigenvectors أو التحamil singular vectors أو loadings. كما تمثل أيضاً كل وحدة (صف) بمجموعة من الدرجات scores تناظر تقديرها في المكونات.

ويبدأ تحليل المكونات الرئيسية بمعايرة جميع المتغيرات ثم حساب مصفوفة التغایر  $S$ . وتطبق عملية تكرارية تهدف للتوصل إلى  $k$  من المكونات الرئيسية حيث  $p > k$ . وتبدأ هذه العملية بالحصول على المكون الرئيس الأول الذي يصف جميع المتغيرات الموجودة، أي الحصول على متوجه المعاملات (الأوزان)

$$a_1 = (a_{11}, a_{21}, \dots, a_{p1})'$$

الناتج عن حل مشكلة تعظيم التباين في  $Y_1$  :

$$\max \text{var}(Y_1) = \max(a_1' S a_1)$$

باستخدام مضاعفات لاجرانج في ظل القيد  $a_1' a_1 = 1$ . ثم الحصول على المكون الثاني، وهكذا حتى الحصول على المكون رقم  $k$ ، أي الحصول على متوجه المعاملات (الأوزان)

$$a_k = (a_{1k}, a_{2k}, \dots, a_{pk})'$$

الناتج عن حل مشكلة تعظيم التباين في  $Y_k$  :

$$\max \text{var}(Y_k) = \max(a_k' S a_k)$$

باستخدام مضاعفات لاجرانج في ظل القيود  $a_k' a_k = 1, a_1' a_1 = a_2' a_2 = \dots = a_{k-1}' a_{k-1} = 0$ .

ويتم رسم أرقام المكونات الرئيسية على المحور الأفقي مقابل القيم المميزة لها على المحور الرئيسي، وهو ما يُعرف برسم الأحجار scree plot، ويختار المكون ذو أقصى ارتفاع.

### 3-2 التحليل العائلي FA

يُستخدم التحليل العائلي Factor analysis في اختزال الأبعاد، فهو يختصر المتغيرات من عدد أكبر إلى عدد أقل من العوامل عند نمذجة البيانات. ويختار FA مجموعة فرعية من المتغيرات من مجموعة أكبر استناداً إلى أعلى الارتباطات بين المتغيرات الأصلية مع عوامل المكونات الرئيسية. وبعد ذلك مدخل لعلاج الأزواج الخطية المتعدد عند توفيق نموذج الانحدار المتعدد لأن مجموعة العوامل الناتجة تكون متغيرات غير مرتبطة. لذلك فإن التحليل العائلي يُستخدم في بناء ما يسمى بنماذج المتغيرات المستترة latent variables، وهي المتغيرات غير المشاهدة التي لا يوجد لها قياسات مسجلة وإنما هي مستبنتة من متغيرات أخرى مشاهدة (من خلال نموذج رياضي) لها قياسات مسجلة. كما يُستخدم أيضاً لاكتشاف الهيكل في العلاقات بين المتغيرات، وهو ما يُعرف باسم تصنيف المتغيرات classify variables.

- وينقسم التحليل العاملی إلى نوعین: التحلیل العاملی الاستکشافی Exploratory factor analysis والتحليل العاملی التوكیدی Confirmatory factor analysis
- فالتحليل العاملی الاستکشافی EFA: هو الذي یبحث في طبيعة أبنية العلاقات المؤثرة على المتغيرات التابعة (أي هياکل النماذج أو أشكالها البنائية).
  - والتحليل العاملی التوكیدی CFA: یختبر أي من هذه الهياکل یؤثر على المتغيرات التابعة عند التنبؤ.
  - ويتعلق التحلیل العاملی بسابقه تحلیل المكونات الرئیسیة، لكنهما ليسا شيئاً واحداً. إذ یستخدم FA أسالیب نمذجة الانحدار لاختبار حدود الخطأ، في حين أن PCA هو مجرد أسلوب إحصائي وصفی.

## ٢-٤ تحلیل التمايز DA

یُستخدم تحلیل التمايز Discriminant Analysis -في علوم الإحصاء والتعرف على الأنماط pattern recognition وتعلیم الآلة machine learning- لإیجاد التولیفة الخطیة من المتغيرات المستقلة الكمية التي تمیز أو تفصل فئتين أو أكثر من الأحداث (متغير تابع تصنیفی). وبمعنى آخر، فإن DA هو طریقة لتصنیف القياسات في مجموعتين أو أكثر. فالغرض الرئیس من DA هو التنبؤ بما یسمی ببعضیة المجموعة group membership استناداً إلى تولیفة خطیة من المتغيرات الكمية. ویبدأ الأسلوب بمجموعة مشاهدات ذات قیم معلومة وذات مجموعات معروفة، وینتهي بنموذج یسمح بالتنبؤ ببعضیة المجموعة بمعلمیة المتغيرات المستقلة الكمية فقط. والغرض الثاني لتحليل التمايز هو فهم مجموعة البيانات بالفحص الدقيق لنموذج التنبؤ لأخذ فكرة عن العلاقة بين عضویة المجموعة والمتغيرات المستقلة المستخدمة في التنبؤ بتلك العضویة.

على سبيل المثال، فإن لجنة القبول بالجامعة قد تقسم خريجيها إلى مجموعتين: الطلاب الذين أنهوا البرنامج في خمس سنوات أو أقل، والطلاب خلاف ذلك. ويمكن استخدام DA للتنبؤ بالاستكمال الناجح لبرنامج الدراسة للطلاب الجدد على أساس درجاتهم في اختبار القدرات GRE score ومعدلهم التراکمي في الثانوية undergraduate grade point average. ويعطی فحص نموذج التنبؤ فكرة عن مدى مساهمة كل متغير (بمفرده وبالاشتراك مع المتغيرات الأخرى) في إكمال أو عدم إكمال البرنامج.

ويتشابه DA مع كل من تحلیل التباين ANOVA وتحلیل الانحدار RA، اللذان یعبران أيضاً عن المتغير التابع بتولیفة من المتغيرات المستقلة. غير أن المتغير التابع في الأسلوبین الآخرين یشترط أن يكون كمیاً، على عکس الحال في DA الذي یكون فيه تصنیفیاً. كما یقترب الانحدار اللوجستی logistic regression والانحدار الاحتمالي probit regression أيضاً بشدة من DA، إذ یفسر الكل متغير

تصنيفي ما. غير أن الأسلوبين الأولين يفضلان في التطبيقات التي لا تفترض أن المتغيرات المستقلة تتبع التوزيع الطبيعي، وهو الفرض الأساسي الذي يُبني عليه DA.

ويتشابه أيضًا DA مع كل من تحليل المكونات الرئيسية PCA والتحليل العاملی FA في أن الكل يبحث عن التوليفات الخطية للمتغيرات التي تعطي أصل تفسير للبيانات. وإذا كان DA يحاول نمذجة الفرق بين فئات البيانات بصرامة، فإن PCA لا يأخذ في حسابه أي فرق في الفئات، كما يبني FA توليفات المتغيرات على الفروق بدلاً من التشابه. كما يختلف DA عن FA في أنه ليس أسلوب تداخل interdependence technique يتم فيه التمييز بين المتغيرات المستقلة والمتغير التابع.

وأخيرًا، فإن DA يُطبق عندما تأخذ المتغيرات المستقلة قياسات كمية مستمرة. أما عندما نتعامل مع متغيرات مستقلة تصنيفية، فإن الأسلوب المكافئ يكون تحليل التمازن التمييزي Discriminant Analysis.

## 2-5 تحليل التمازن التمييزي DCA

كما يشير الاسم، فإن تحليل التمازن التمييزي هو امتداد لكل من تحليل التمايز DA وتحليل التمازن CA. ويهدف DCA (مثل DA) إلى تصنیف المشاهدات في مجموعات معرفة مسبقاً (ومثل CA) في أنه يستخدم مع المتغيرات الاسمية. إن الفكرة الأساسية وراء DCA هي تمثيل كل مجموعة بإجمالي مشاهداتها وإجراء CA بسيط على المجموعات عن طريق مصفوفة المتغيرات. ويتم التنبؤ بالمشاهدات الأصلية وتخصيص كل مشاهدة متوقعة في المجموعة الأقرب. ويمكن استخدام المقارنة بين التصنيفين القبلي والبعدي priori and the a posteriori classifications باستخدام أساليب التحقق المبدل من الصحة cross-validation techniques.

## 2-6 التحليل العنودي Cluster Analysis

يُعد التحليل العنودي من أشهر الطرق الوصفية (الاستكشافية) للتقييم في البيانات، وهو منهج لتجمیع grouping مجموعة معينة من المشاهدات. فإذا تكونت مصفوفة البيانات من  $n$  من المشاهدات (الحالات أو الصفوف) و  $p$  من المتغيرات (الحقول أو الأعمدة)، فإن هدف التحليل العنودي يكون عقدة أو تصنیف المشاهدات في مجموعات متاجنسة (مت Manson) داخلياً internal cohesion وغير متاجنسة من مجموعة إلى أخرى (منفصلة خارجياً external separation). ويفسر ذلك على أنه اختزال للأبعاد في الفضاء  $R^p$ ، ولكن ليس بنفس طريقة المكونات الرئيسية. إذ يقوم التحليل العنودي بالاختزال الرأسي بتجمیع

المشاهدات  $n$  في  $g$  من المجموعات الفرعية (حيث تكون  $n < g$ )، بينما يقوم تحليل المكونات الرئيسية بتحويل المتغيرات الأصلية  $p$  إلى  $k$  من المتغيرات الجديدة (حيث يكون  $k < p$ ). ويمكن تكوين التجمعات groupings أو التقسيمات partitions أو العناقيد clusters بنوعين من الطرق:

- الطرق الهرمية hierarchical methods: ويتم فيها تقدير عدد العناقيد بإجراء أسلوب التعاقب succession بدءً من  $n$  (وهي الحالة الأسطى التي تُعامل فيها كل مشاهدة على أنها مجموعة منفصلة) حتى 1 (كل المشاهدات تنتهي لمجموعة واحدة).
- الطرق غير الهرمية non-hierarchical methods: ويكون فيها عدد العناقيد معروف مسبقاً.

## 2-7 نماذج البروبيت واللوجيت

إن نموذج البروبيت أو نموذج الوحدة الاحتمالي probit model هو نوع خاص من الانحدار يكون فيه المتغير التابع من النوع التصنيفي ( الثنائي binary ) و يأخذ قيمتين فقط، النجاح ويشار إليه بالرمز 1 والفشل ويشار إليه بالرمز 0 . ومثال ذلك: متزوج وغير متزوج، ناجح وراسب، يفضل ولا يفضل، الإجابة بنعم أو لا، وجود أو غياب صفة معينة ... إلخ.  
ويأخذ نموذج البروبيت الشكل التالي:

$$\Pr(Y = 1 | X) = \Phi(X'\beta),$$

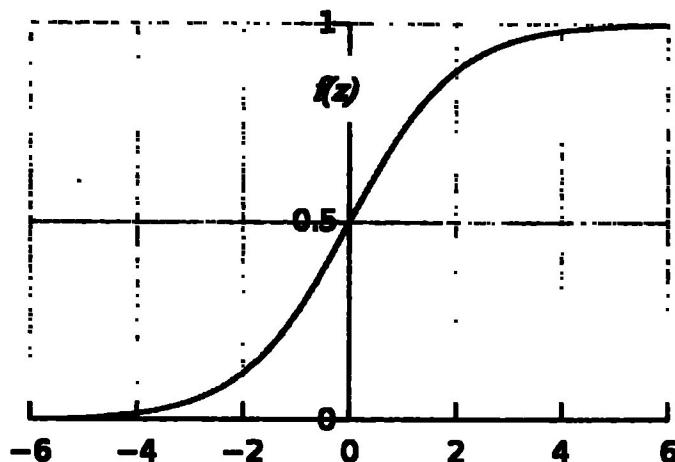
حيث يشير الرمز  $\Pr$  إلى الاحتمال، والرمز  $\Phi$  إلى دالة التوزيع المتجمعة للتوزيع المعتاد المعياري، والرمز  $\beta$  إلى المعالم المقدرة باستخدام طريقة الإمكان الأكبر التقليدية.

أما نموذج اللوجيت logit model فهو كسابقه أسلوب أحادي/متعدد المتغيرات يسمح بتقدير احتمال وقوع/عدم وقوع حدث ما لمتغير تابع ثانٍ، ولكنه يأخذ الشكل التالي:

$$y = \exp(b_0 + b_1*x_1 + \dots + b_n*x_n) / \{1 + \exp(b_0 + b_1*x_1 + \dots + b_n*x_n)\}$$

ويُعد نموذج الانحدار اللوجستي LRM مثلاً لهذا النوع من النماذج. وتأخذ فيه الدالة اللوجستية دائمًا مثلاً الاحتمالات - قيمًا تتراوح بين الصفر والواحد، وتعُرف بالنموذج والشكل التاليين:

$$f(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$



شكل (1): مثال للدالة اللوجستية

ويشير المتغير  $z$  (الذي يمكن أن يأخذ أي قيم عددي) إلى مدخلات الدالة، بينما تتحصر قيم المخرجات  $f(z)$  بين الصفر والواحد. ويمثل المتغير  $z$  التعرض لمجموعة ما من المتغيرات المستقلة، بينما تمثل  $f(z)$  احتمال الناتج المقابل في ظل قيم المتغيرات المفسرة. ويقيس المتغير  $z$  المساهمة الكلية لجميع المتغيرات المستقلة المستخدمة في النموذج ويُطلق عليه اسم logit، ويُعرف بالمعادلة:

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k,$$

ويُعد الانحدار اللوجستي طريقة مفيدة في وصف العلاقة بين متغير مستقل أو أكثر (مثل العمر، والسن، إلخ) ومتغير استجابة ثانوي يأخذ قيمتين فقط (النجاح أو الفشل).

### 3. الأساليب اللامعنية

#### 3-1 العملية الهرمية التحليلية AHP

تُعرَّف العملية الهرمية التحليلية Analytical hierarchy process بأنها أسلوب تَقْرُّع structured technique لتنظيم وتحليل القرارات المعقدة. وهي عبارة عن مزيج من علم الرياضيات وعلم النفس قدم على يد Thomas L. Saaty في السبعينات للتوصُّل لأفضل قرار من بين البديلات المتاحة<sup>[1]</sup>. وقد لاقت هذه العملية استحساناً كبيراً من جانب صانعي القرار وقبولاً واسعاً في المجالات الحكومية والتعليم والصحة والصناعة ومنشآت الأعمال. فبدلاً من أن تفترض المنشأة بأن قرارها "صحيح"، فإن AHP تساعد على العثور على أفضل قرار يتوافق مع هدفها وفهمها للمشكلة. وتتوفر AHP إطاراً شاملًا وعقلانية لبناء مشكلة القرار، ولتمثيل وقياس عناصرها، ولربط هذه العناصر لتحقيق الأهداف العامة، ولتقييم الحلول البديلة.

وأول ما يفعله مستخدمو AHP هو تفكيك مشكلة القرار بشكل هرمي إلى مشاكل فرعية يمكن فهمها بسهولة أكثر، بحيث يمكن تحليل كل منها بشكل مستقل. ويمكن أن تتعلق عناصر التسلسل

الهرمي بأي جانب من جوانب المشكلة، سواء كانت تلك العناصر ملموسة أو غير ملموسة، وسواء قيست بدقة أو قُدرت بشكل تقريري، وسواء فهمت جيداً أو بشكل ضعيف.

وبمجرد بناء الهرم، يقوم صانعو القرار بتقييم عناصره المختلفة بمقارنة كل عنصر بالعناصر الأخرى اثنين في كل مرة من حيث تأثيرها على العنصر الذي يعلوها في التقسيم الهرمي. وعند عمل المقارنات، يمكن لصانع القرار أن يستخدم بيانات واقعية عن العناصر، كما يمكنه أيضاً استخدام حكمه عن الأهمية النسبية للعناصر. وهكذا، فإن جوهر AHP يعتمد على استخدام الأحكام الشخصية إلى جانب المعلومات الأساسية في إجراء التقييمات.

وتحول AHP هذه التقييمات إلى قيم عدديّة<sup>[2]</sup> يمكن معالجتها ومقارنتها على المدى الكامل للمشكلة. ويُشتق الوزن العددي أو ما يُعرف بالأولوية priority بالنسبة لكل عنصر من التسلسل الهرمي، مما يسمح بالمقارنة مع العناصر المتعددة وغير القابلة للقياس في كثير من الأحيان مع بعضها البعض بطريقة عقلانية ومتسقة. وتميز هذه القدرة AHP عن غيرها من أساليب صنع القرار.

وفي الخطوة الأخيرة من العملية، يتم حساب الأولويات العددية لكل بديل من بدائل القرار. وتمثل هذه الأرقام القدرة النسبية للبدائل في تحقيق الهدف المقرر.

### 3-2 الأنظمة الخبرية ES

يُعرف النظام الخبير expert system (في مجال الذكاء الاصطناعي artificial intelligence) بأنه نظام حاسبي يحاكي قدرة الخبرة البشرية في صناعة القرار<sup>[25]</sup>. وتحصم الأنظمة الخبرية لحل المشاكل المعقدة عن طريق المنطق المكتسب من المعرفة، أي بطريقة الخبرير وليس بإتباع أسلوب المطور كما هو الحال في البرمجة الاتفاقية<sup>[38][9]</sup>. وقد قدمت أول الأنظمة الخبرية في السبعينيات ثم انتشرت بعد ذلك في الثمانينيات<sup>[10]</sup>.

ويكون النظام الخبير من قسمين: الأول ثابت مستقل عن النظام هو محرك الاستنتاج the inference engine، والثاني متغير يمثل قاعدة المعرفة the knowledge base. وفي الثمانينيات ظهر قسم ثالث يسمح بالاتصال بالمستخدمين هو واجهة الحوار<sup>[27]</sup>.

وقد تم تصميم النظم الخبرية لتسهيل المهام في مجالات المحاسبة، والقانون، والطب، التحكم في العمليات، والخدمات المالية، والإنتاج، والموارد البشرية. لذلك فقد ساندتها تطبيقات كثيرة في مجالات تشخيص الأعطال، والتشخيص الطبي، ودعم القرارت في الأنظمة المعقدة، والرقابة على العمليات، والبرامج التعليمية، وإدارة المعرفة.

### 3- طريقة أقرب الجيران k-NN

تُعد طريقة أقرب الجيران Nearest Neighbors من الطرق المبنية على الذاكرة، بمعنى أنها لا تتطلب (بمصطلاحات التقريب في البيانات) أي تدريب (توفيق نموذج للبيانات) على خلاف الطرق الإحصائية الأخرى. و تستند  $k$ -NN على فكرة بدائية تلخص في أن المشاهدات القريبة يجب أن تقع في نفس الفئة. فهي أسلوب تصنيف يقرر في أي فئة سنضع الحالة الجديدة بفحص عدد ما ( $k$ ) في معظم الحالات المشابهة أو الجيران. ويلجأ المحل لهذه الطريقة عند عمل التحاليل المقارنة باستخدام أساليب اختزال البيانات.

ويتطلب تطبيق الطريقة<sup>[7]</sup> معايرة جميع المتغيرات وحساب المسافات الإقليدية بين كل زوج من المشاهدات. ويمكن تصنيف المشاهدة الجديدة بوحدة من 4 طرق هي: طريقة Papadakis التي تسمى أحياناً Genesis (وهي مبنية على حساب البوافي ثم استخدام طريقة تحليل التغير)، وطريقة الارتباط ( وهي مبنية على استخدام الارتباط بين كل زوج من المشاهدات من خلال المربعات الصغرى المعممة بشرط معلومية هيكل الارتباط)، وطريقة Wilkinson، وطريقة تمهيد المربعات الصغرى.

ومن أهم تطبيقات  $k$ -NN: التعرف على الأنماط، والتصنيف الإحصائي، والتحليل العنقودي، واسترجاع المحتوى المبني على الصور من قواعد البيانات، والتسوق عبر الإنترنت.

### 4- النماذج المعممة المضافة GAMs

إن النماذج المعممة المضافة generalized additive models هي<sup>[20]</sup> أحد مداخل الانحدار اللامعجمي في حالة تعدد المتغيرات المستقلة. وقد قدم هذا الأسلوب<sup>[23]</sup> في التسعينيات على يد Trevor Hastie and Rob Tibshirani. ويخلط هذا النموذج بين خصائص النماذج الخطية المعممة generalized linear models والنماذج المضافة additive models. وإذا كان النموذج الخطي الجمعي يأخذ الشكل التالي:

$$Y = b_0 + b_1 * X_1 + \dots + b_m * X_m$$

فإن GAM يأخذ الشكل المختلف التالي:

$$g(\mathbf{E}(Y)) = \beta_0 + f_1(x_1) + f_2(x_2) + \cdots + f_m(x_m).$$

ويُلاحظ بمقارنة النموذجين أن GAMs أخذت من النموذج المتعدد حفاظها على الشكل الجمعي، غير أنها استبدلت الحدود البسيطة في المعادلة الخطية (أي  $b_i * X_i$ ) بالدوال ( $f_i(X_i)$ ) وهي دوال لامعلمية للمتغيرات المفسرة  $X$ . وبعبارة أخرى، فإن GAMs تقدر دوال لامعلمية غير محددة لكل متغير مستقل بدلاً من المعامل للتوصيل لأفضل تنبؤ لقيم المتغير التابع.

ويمكن توفيق الدوال  $f_i(X_i)$  باستخدام<sup>[28]</sup> أحد ممهدات الشكل الانتشاري التالية:

- 1) شرائح التمهيد المكعب cubic smoothing spline وهي متوفرة في برنامج SAS
- 2) أسلوب LOESS وهو أيضاً متاح في البرنامج السابق
- 3) ممهد النواة Kernel smoother وهو متاح في برنامج STATA
- 4) الشرائح الرقيقة thin-plate splines التي تسمح بوجود تفاعل بين المتغيرات المستقلة وهو متاح في برنامجي SAS و R.

وأخيراً، تُعد GAMs مفيدة في الحالات التالية<sup>[29]</sup>: 1) إذا كان شكل العلاقة بين المتغيرات شديد التعقيد بشكل يصعب معه توفيق نموذج خطى تقليدي أو أي من النماذج غير الخطية 2) إذا لم يتوفر سبب مسبق لاستخدام نموذج معين 3) إذا كنا نريد أن تقترح البيانات الشكل الدالى المناسب. ويعنى ذلك أن تلك النماذج تناسب معظم التطبيقات الحديثة التي تحتوى على عدد كبير من المتغيرات بينها تفاعلات ممكنة في ظل أحجام البيانات الكبيرة مثل أسواق الأسهم.

### 3-5 البرمجة الرياضية MP

تشير البرمجة الرياضية<sup>[30]</sup> mathematical programming (في كل من الرياضيات وعلم الإدارة وعلوم الحاسب) إلى الأمثلية optimization؛ أي عملية اختيار أفضل الحلول من بين عدة بدائل متاحة في ظل مجموعة من القيود. وت تكون مشكلة الأمثلية في شكلها البسيط من تعظيم أو تصغير دالة حقيقة باختيار قيم المتغيرات الهامة من بين مجموعة من المتغيرات وحساب قيمة دالة الهدف. ويسمح تعليم مشكلة الأمثلية بوجود تشكيلة متنوعة من دوال الهدف وأنواع مختلفة من النطاقات.

ويتيح زر البرامج الإضافية<sup>[31]</sup> Add-in في برنامج Excel بناء نماذج البرمجة الرياضية وحلها باستخدام حل المشاكل Solver Add-in Math Programming add-in، يتم إضافة سطور أوامر لكل من: البرمجة الخطية والكسرية linear and integer programming، والبرمجة غير الخطية nonlinear programming، وشبكات الأعمال network programming، ومشكلات النقل transportation.

### 3-6 دوال الانحدار المقسمة المتوازنة متعددة المتغيرات MARS

تُعد دوال الانحدار المقسمة المتوازنة متعددة المتغيرات<sup>[32][18]</sup> Multivariate Adaptive Regression Splines شكلًا من أشكال تحليل الانحدار. وهي أسلوب انحدار لامعجمي ينمذج اللاحظية والتفاعلات، ويبني النماذج بالشكل التالي:

$$\hat{f}(x) = \sum_{i=1}^k c_i B_i(x)$$

وهو عبارة عن مجموع لدوال الأساس  $B_i(x)$  المرجحة بالمعاملات الثابتة  $c_i$ . وتأخذ كل دالة أساس basis أحد الأشكال الثلاث التالية: 1) الثابت 1 وهو ما يسمح بظهور حد التقاطع intercept function في النموذج 2) دالة مفصلية hinge function على الشكل  $\max(0, x - const)$  أو الشكل  $\max(0, const - x)$  (3) حاصل ضرب دالتين knots المتغيرات وقيم العقد  $\max(0, const - x)$ ، حيث تختار MARS حيث تختار knots تلقائياً مفصليتين أو أكثر.

ويمكن توفيق نموذج MARS على مرحلتين؛ بنفس المنهج المستخدم في التقسيم المتكرر recursive عند توفيق شجرة القرارات partitioning:

1) المرور للأمام the forward pass: ويبداً بنموذج يحتوي على حد التقاطع فقط (متوسط قيم المتغير التابع) ثم إضافة زوج من دوال الأساس إلى النموذج في كل مرة إلى أن نصل لأقصى اخترال في الخطأ المعيّر عنه بمجموع مربعات الباقي. غير أن التوفيق الأمامي عادةً ما يبني نموذج ذو جودة توفيق فوقية overfit (نموذج ذو جودة توفيق جيدة بالنسبة للبيانات المستخدمة في بنائه، غير أن أدائه التنبؤي بالنسبة للبيانات الجديدة يكون ضعيف).

2) المرور للخلف the backward pass: وهو استكمال للمرحلة السابقة للتغلب على مشكلة التوفيق الفوقي (تحسين القدرة التنبؤية) بتقليم prunes النموذج عن طريق حذف حدوده واحداً تلو الآخر، حيث يتم حذف الحد الأقل تأثيراً في كل خطوة إلى أن يتم الوصول إلى أفضل نموذج فرعي. ويقارن أداء النماذج الفرعية باستخدام طريقة التحقق من الصحة المقاطعة المعممة Generalized cross validation (GCV) لاختيار أفضل نموذج فرعي؛ حيث تشير القيمة الأقل لـ GCV لنموذج أفضل. وتحسب GCV بالصيغة التالية:

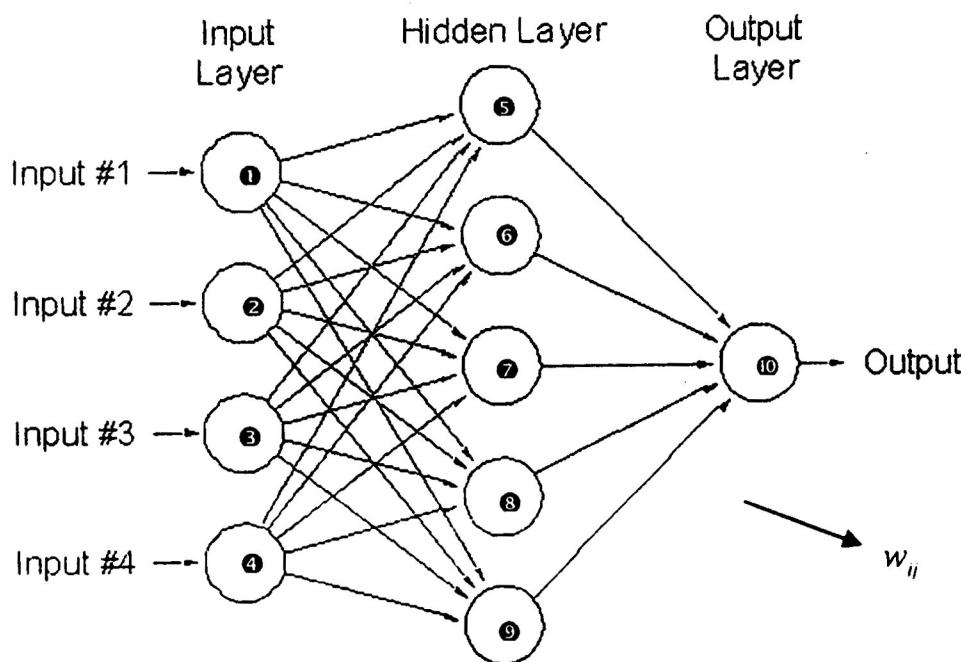
$$GCV = RSS / (N * (1 - Effective\ Number\ Of\ Parameters / N)^2)$$

### 3-8 الشبكات العصبية NN

تُستخدم الشبكات العصبية Neural Networks لتحقيق العديد من الأغراض الوصفية والتنبؤية عند التقريب في البيانات<sup>[20]</sup>. وقد نشأت NN في مجال تعليم الآلة Machine Learning في محاولة لتقليد الوظائف العصبية للمخ البشري من خلال توليفة من العناصر الحاسبية البسيطة (الخلايا العصبية Neurons) في نظام متداخل للغاية. وتتمتع NN بأهمية خاصة لأنها تقدم نمذجة عالية الكفاءة للمشاكل المعقدة (التي تحتوي على مئات المتغيرات المسنقة والعديد من التفاعلات ومتغير تابع أو أكثر) بطريقة

لامعلمية من قواعد البيانات الكبيرة. كما يمكن استخدامها في حل مشاكل التصنيف ومشاكل الانحدار سواء كانت البيانات مكتملة أو مبتورة.

مكوناتها: يوضح شكل (2) أن الشبكة العصبية تتكون من مجموعة من الوحدات الحاسبية الأولية (تعرف باسم الخلايا العصبية متصلة بما يليها من خلال روابط مرجحة). وتمثل كل خلية دائرة، وتأخذ رقمًا طبيعياً (من 1 : 10 في هذا المثال). كما تمثل الروابط بأسمهم وتأخذ الرمز  $w_{ij}$ ، حيث يشير الدليل  $i$  إلى رقم العقدة التي ينطلق منها السهم ويشير الدليل  $j$  إلى رقم العقدة التي ينتهي إليها. وتنظم هذه الوحدات في طبقات Layers بحيث تتصل كل خلية (في طبقة ما) بجميع خلايا الطبقة السابقة واللاحقة. وتبدأ الشبكة بطبقة المدخلات Input Layer (من 1 : 4 في هذا المثال) التي تنتظر كل عقدة فيها أحد المتغيرات المستقلة. وتتصل كل عقدة في طبقة المدخلات بجميع عقد الطبقة الخفية (من 5 : 9 في هذا المثال)، وربما تتصل عقد الطبقة الخفية بجميع عقد طبقة خفية أخرى (غير موضح على الرسم). وتنتهي الطبقات بطبقة المخرجات Output Layer (رقم 10 في هذا المثال) وهي عقدة (أو أكثر) تمثل المتغير التابع (أو المتغيرات التابع) وهي النقاء للأسماء الخارجية من آخر طبقة خفية.



شكل (2): نموذج لشبكة عصبية بسيطة

ويُحسب الوزن  $w_{ij}$  بمجموع حواصل ضرب الأوزان الداخلة على العقدة التي ينطلق منها في قيم العقد التي تتطرق منها تلك الأوزان. وكمثال، فإن قيمة الوزن الرابط بين الطبقة 7 والطبقة 10 هو :

$$w_{7,10} = w_{17} * \text{value of node 1} + w_{27} * \text{value of node 2} + w_{37} * \text{value of node 3} + w_{47} * \text{value of node 4}$$

ويمكن أن يُنظر إلى كل عقدة على أنها متغير مستقل (العقد من 1 : 4)، أو على أنها توليفة (تفاعل) من المتغيرات المستقلة (العقد من 5 : 10). فالعقدة 10 هي توليفة غير خطية للقيم في العقد من 1 : 4 بسبب وجود دالة التنشيط (القيم المجمعة في عقد الطبقة الخفية). وجدير بالذكر أنه إذا كانت دالة التنشيط خطية ولا توجد طبقة خفية، فإن الشبكة العصبية تُختزل إلى الانحدار الخطى. بينما تُختزل الشبكة العصبية إلى الانحدار اللوجيستى في ظل دوال تنشيط غير خطية ذات شكل معين.

### الإمكانية **Potential**

تعبر الأوزان في الشبكة العصبية (كما في النموذج البيولوجي) عن معاملات قابلة للتعديل استجابة للإشارات التي تسافر في الشبكة بحسب خوارزمية تعلم مناسبة وقيمة فاصلة **Threshold** (تعرف أيضاً باسم التحيز Bias) تشبه حد التقاطع في نموذج الانحدار. فالخلية  $j$  تأخذ القيمة الفاصلة  $\theta_j$  وتنstem إشارات داخلة  $[x_0, x_1, \dots, x_n] = x$  من الوحدات (الخلايا/العقد) المتصلة بها من الطبقة السابقة. وتقترب كل إشارة بوزن معين  $[w_0, w_1, \dots, w_n] = w$ .

وتتم دراسة الإشارات الداخلية وأوزانها والقيمة الفاصلة لكل خلية من خلال ما يسمى بدالة التوليف **Combination Function**. وتنتتج دالة التوليف (كل خلية) قيمة واحدة تسمى **الإمكانية** (أو الداخل الصافى Net Input). وتقوم دالة التنشيط Activation Function بتحويل **الإمكانية** إلى إشارة خارجة.

وتكون دالة التوليف عادةً خطية، لذلك فإن **الإمكانية**  $p_j$  تكون مجموع انحرافات قيم الخلايا السابقة  $x_i$  المرجحة بالأوزان الخارجية منها  $w_{ij}$  عن القيمة الفاصلة  $\theta_j$ ، وهو ما يُعبر عنه رمزاً كالتالي:

$$p_j = \sum_{i=1}^n (x_i w_{ij} - \theta_j) = \sum_{i=0}^n x_i w_{ij}$$

حيث  $x_0 = 1$ ،  $w_0 = -\theta_j$ . ويمكن الحصول على الإشارة الخارجية للخلية  $j$  (أي  $y_j$ ) بتطبيق دالة التنشيط على **الإمكانية**  $p_j$  لتعطى:

$$y_j = f(x, w_j) = f(p_j) = f(\sum_{i=0}^n x_i w_{ij})$$

### أنواع دالة التنشيط:

هناك طرق كثيرة لتنشيط الخلايا في الشبكة العصبية. ومن أشهرها: الطريقة الخطية، والطريقة المجزأة Piecewise، والطريقة الإيسية Sigmoidal، وطريقة أقصى تمديد **Softmax** .

**1) دالة التنشيط الخطية:** تُعرف دالة التنشيط الخطية بالصيغة التالية:

$$f(p_j) = \alpha + \beta p_j,$$

حيث تنتهي الإمكانية  $p$  لمجموعة الأعداد الحقيقة، و  $\alpha, \beta$  ثوابت. وعندما يتطلب النموذج أن يكون مخرج الخلية مساوً تماماً لمستوى تشيشتها (الإمكانية)، نضع  $\alpha = \beta = 1$  وتحول الدالة الخطية إلى ما يسمى بدالة الوحدة. يلاحظ التشابه القوي بين دالة التشيش الخطية ونموذج الانحدار الخطى البسيط، إذ يمكن النظر للأخير على أنه نوع بسيط من الشبكات العصبية.

(2) دالة التشيش المجزأة: تُعرف دالة التشيش الخطية بالصيغة التالية:

$$f(p_j) = \begin{cases} \alpha & p_j \geq \theta_j \\ \beta & p_j < \theta_j \end{cases}$$

ويوضح أن تأخذ قيمتين فقط بحسب تجاوز الإمكانية لقيمة الفاصلة من عدمه. وعندما تكون  $\alpha = 1, \beta = 0, \theta_j = 0$ ، تكون أمام حالة خاصة من التشيش المجزأ تُعرف باسم دالة تشيش الإشارة Sign Activation Function التي تأخذ القيمة 1 إذا كانت الإمكانية موجبة والقيمة 0 بخلاف ذلك.

(3) دالة التشيش الإيسية: أي التي تأخذ شكل حرف S، وهي الأكثر استخداماً في التطبيقات العملية. وتُنتج هذه الدالة قيمة موجبة فقط في الفترة  $[0, 1]$ . ويرجع شيوخ استخدامها إلى أنها غير خطية وإلى قابليتها للفهم وللتفاضل بسهولة. وتُعرف بالصيغة التالية:

$$f(p_j) = \frac{1}{1 + e^{-\alpha p_j}}$$

حيث تشير  $\alpha$  إلى معلمة موجبة تنظم ميل الدالة.

(4) دالة أقصى تمديد: تُستخدم في تطبيق Normalize مخرجات العقد المختلفة التي يوجد بينها علاقة. فإذا كانت الشبكة تحتوي على  $g$  من العقد بمخرجات عددها  $v$  (حيث  $v = 1, 2, \dots, g$ )، فإن دالة أقصى تمديد التي تُطبع  $v$  (تجعل مجموعها 1) تكون:

$$\text{soft max}(v_j) = \frac{e^{v_j}}{\sum_{j=1}^g e^{v_j}}$$

وستخدم هذه الدالة في حل مشاكل التصنيف المراقب Supervised Classification Problems عندما يأخذ المتغير التابع عدد  $g$  من المستويات.

### طرق التدريب [3] : Training Methods

يُقصد بالتدريب (بمفاهيم الشبكات العصبية) تعليم الشبكة كيف تنجز مهمة ما، وهو بلغة الإحصائيين الطريقة المستخدمة في تقدير أوزان الشبكة (المعالم المجهولة). ويمكن تدريب أو تعليم الشبكة بعدة طرق من أشهرها وأوسعاً انتشاراً؛ طريقة الإثمار الخلفي Backpropagation التي تبحث في تحديث أوزان الشبكة

بتضييف دالة الخطأ في فضاء الأوزان باستخدام عدة خوارزميات، من أشهرها: خوارزمية الهبوط المتدرج [12]، خوارزمية التدرج المقارن [33][8]، وخوارزمية Quasi-Newton، وخوارزمية Levenberg-Marquardt [35]، والخوارزميات الجينية [19]. Genetic Algorithms

### أنواع الشبكات العصبية:

يمكن تصنيف الشبكات العصبية بحسب عدد طبقاتها في نوعين: شبكة الفواهم ذوي الطبقة الواحدة Multi-Layer Perceptrons، وشبكة الفواهم متعددة الطبقات Feedforward Networks. كما يمكن تصنيف الشبكات العصبية بحسب اتجاه تدفق المعلومات إلى: شبكات التغذية الأمامية Feedforward Networks؛ حيث تتحرك فيها المعلومات من طبقة إلى طبقة التالية للأمام فقط دون السماح لها بالعودة للخلف، وشبكات التغذية الخلفية Feedback Networks؛ حيث تتحرك فيها المعلومات من طبقة إلى طبقة التالية للأمام مع السماح لها بالعودة للخلف إلى الطبقات السابقة.

### تطبيقات الشبكات العصبية:

تكون الشبكات العصبية قابلة للتطبيق في المشاكل السببية حين توجد علاقة معقدة بين عدة متغيرات مستقلة (مفسرة/ مبنية/ مدخلات) واحد أو أكثر من المتغيرات التابعه (مفسر/ متباً به/ مخرجات)، ويصعب التعبير عن تلك العلاقة بالمدخلات التقليدية كالارتباط والانحدار والاختلاف بين المجموعات. ومن أمثلة المشاكل التي طُبقت فيها الشبكات العصبية بنجاح [34][6]: الكشف عن الظواهر الطبيعية، و التنبؤ بسوق الأسهم، وتقييم الجدارة الائتمانية لطالبي القروض.

### 9-3 الانحدار اللامعملي NR

إذا قرر الباحث مثلاً استخدام كثيرة حدود تكعيبية على الشكل [29]:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon$$

في توفيق نموذج الانحدار الذي يربط بين المتغير  $y$  والمتغير  $x$ ، فإن ذلك يتشرط صحة الشكل الرياضي المفترض في التعبير عن البيانات. ويقال أن النموذج معلمي لأنّه يعتمد على المعالم  $\beta_1, \beta_2, \beta_3$ . أما عندما لا تتوافر معلومات كافية لصنع فرض مثل هذا، أو عند الرغبة في مجرد افتراض أن:

$$y = f(x) + \epsilon$$

في ظل فرض التمهيد العادي [يأن  $(x), f(x), f'(x), f''(x)$  كلها مستمرة] وتقدير  $f(x)$  من البيانات، فإننا نستخدم الانحدار اللامعملي nonparametric regression.

فالانحدار اللامعملي  $NR$  إذن<sup>[37]</sup> هو شكل من تحليل الانحدار لا يأخذ فيه المتغير المستقل شكل محدد، ولكنه يُبنى من المعلومات المشتقة من البيانات. لذلك فإن  $NR$  يتطلب حجم عينة أكبر من الحجم اللازم لحساب الانحدار المعملي لأن البيانات هي التي تقترح هيكل النموذج وتقديرات المعامل. ويُقدر  $NR$  من خلال<sup>[23][26]</sup> دوال الشرائح الممهدة المكعبية cubic smoothing spline التي تأخذ الشكل التالي:

$$RSS(f, \lambda) = \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_0^b (f''(t))^2 dt$$

وهو ما يعني تقدير مجموع مربعات الباقي من جميع الدوال الممكنة  $(x)f$  في ظل مشتقين مستمرتين. وتشير  $\lambda$  إلى معلمة التمهيد المثبتة، ويشير الحد الأول من الطرف الأيمن قرب البيانات، ويقيس الحد الثاني مدى انحصار الدالة، وتحدد  $\lambda$  المفاضلة بين الحدين. فإذا كانت  $\lambda = 0$ ، فإن  $f$  يمكن أن تكون أي دالة تستكمel البيانات. أما إذا كانت  $\lambda = \infty$ ، يتم توفيق  $f$  بخط المربيعات الصغرى المستقيم بسبب عدم الاستفادة من المشتقة الثانية.

#### 4. الأساليب نصف المعلمية

##### 1-4 الانحدار نصف المعملي SPR:

يُعد الانحدار نصف المعملي Semiparametric Regression توليفة من الانحدارين المعملي واللامعملي. وهو يستخدم إذا كان النموذج اللامعملي الكامل لا يعبر بشكل جيد عن البيانات و/أو إذا أراد الباحث استخدام نموذج معملي لكنه لا يعرف بالضبط شكله الدالي بالنسبة لمجموعة فرعية من المتغيرات المفسرة أو إذا كانت كثافة الأخطاء غير معروفة. وحيث أن SPR تحتوي على مركبة معلمية، فإنها تعتمد على فروض معلمية؛ وبالتالي فإنها تكون معرضة لمشكلتين مهمتين: خطأ التحديد misspecified (اختيار شكل رياضي خاطئ و/أو القصور في إدخال المتغيرات المعتبرة عن المشكلة)، وعدم الاتساق inconsistent (عدم تمركز توزيع المقدرات بالقرب من القيمة الحقيقية للمعلمة المقدرة) كما في النماذج المعلمية الكاملة.

ويوجد العديد من الطرق لتقدير نماذج SPR، أشهرها:

1] **النموذج الخطى الجزئى Partially Linear Model** المعرف بالشكل التالي:

$$Y_i = X_i\beta + g(Z_i) + u_i, \quad i = 1, \dots, n,$$

حيث يشير  $Y_i$  إلى المتغير التابع، وكل من  $X_i, Z_i$  إلى متوجهى المتغيرات المستقلة من الدرجة  $1 \times p$ ، و  $\beta$  إلى متوجه المعامل من الدرجة  $1 \times p$ ، و  $g \in R^q$ . ويعرف متوجه المعامل  $\beta$  الجزء المعلمي من

النموذج، بينما تُعرف الدالة المجهولة ( $Z_i$ ) والجزء الامامي منه ويتم تقديرها بأي طريقة انحدار لامعممية مناسبة.

[2] نموذج الرقم المفرد<sup>[24]</sup> Ichimura's method single index model ويأخذ الشكل التالي:

$$Y = g(X'\beta_0) + u,$$

وتقدير فيه المعلم  $\beta_0$  باستخدام طريقة المربيعات الصغرى غير الخطية لتصغير الدالة:

$$\sum_{i=1} (Y_i - g(X'_i\beta))^2.$$

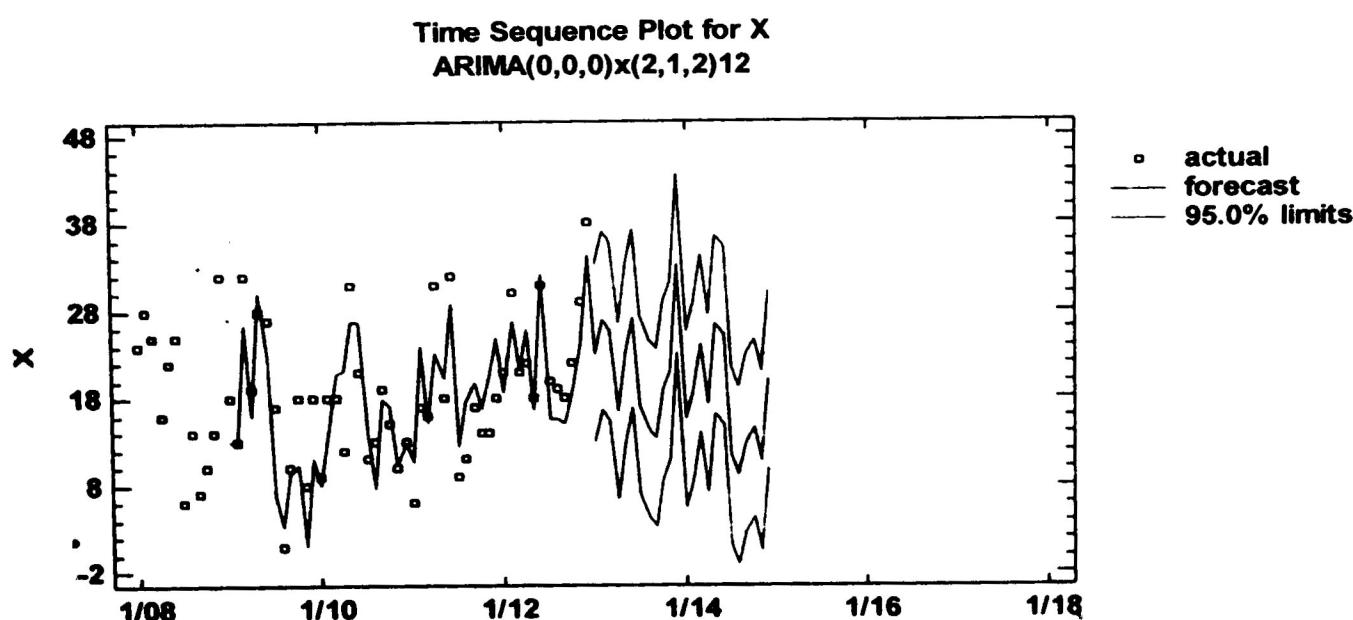
[3] نماذج المعامل الممهد أو المتغير<sup>[22]</sup> Smooth coefficient\varying coefficient models التي تُعرف بالصيغة التالية:

$$Y_i = \alpha(Z_i) + X'_i\beta(Z_i) + u_i = (1 + X'_i) \begin{pmatrix} \alpha(Z_i) \\ \beta(Z_i) \end{pmatrix} + u_i = W'_i\gamma(Z_i) + u_i,$$

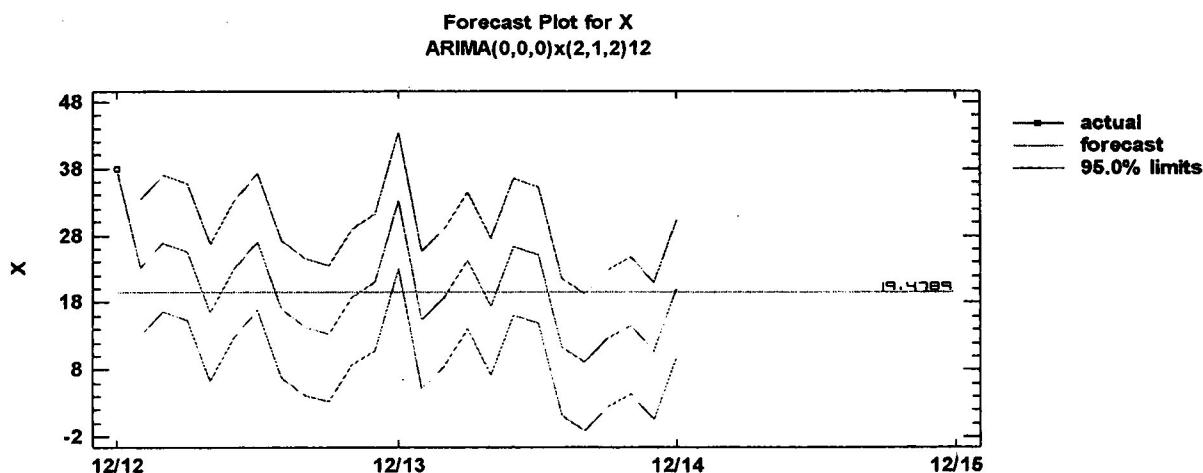
حيث يشير  $X_i$  إلى متوجه من الدرجة  $1 \times k$  و  $(z)$  إلى متوجه من الدوال الممهدة غير المحددة في  $z$ .

#### 5- تطبيق للمقارنة بين النماذج:

تم استخدام السلسلة الشهرية لبيانات القضايا المعروضة على المحكمة الدستورية العليا في مصر خلال الفترة 2008-2013 (المحكمة الدستورية العليا، إدارة المعلومات و البيانات) لتوفيق عدة نماذج حسب طبيعة البيانات المستخدمة والتي توضح وجود أثر موسمي والذي بلغ اقصاه شهر يونيو (611) وأدناء شهر أغسطس (809)، والتي يوضحها الشكل التالي:



شكل (3) أعداد القضايا المعروضة على المحكمة الدستورية خلال الفترة 2008-2013



شكل (4) الدليل الموسمي لعدد القضايا خلال الفترة 2008-2013

في البداية تم توفيق 18 نموذج (تقليدي واريما) موضحة بالجدول التالي، وكان أفضلها وفقاً للمعايير المحددة هو نموذج اريما الموسمي ARIMA(2,1,0)x(1,1,2)12 حيث بلغ معامل الارتباط الذاتي الأول 0.634

جدول (1) مقارنة لمؤشرات الجودة للنماذج التقليدية واريما عدد القضايا 2008-2013

Model	RMSE	MAE	MAPE
Random walk	9.06997	6.77431	62.3605
Random walk with drift = 0.0104027	9.16648	6.77481	62.3812
Constant mean = 18.8607	7.08446	5.35525	49.0912
Linear trend = $-68.2119 + 0.119852 t$	7.10294	5.2331	46.7143
Quadratic trend = $6693.6 + -18.5055 t + 0.0128185 t^2$	6.58787	4.78467	41.0339
Exponential trend = $\exp(-2.00416 + 0.006681 t)$	7.11337	5.11647	42.5103
S-curve trend = $\exp(7.54147 + -3406.72 /t)$	7.12695	5.12222	42.6378
Simple moving average of 2 terms	7.60982	5.44652	51.2748
Simple exponential smoothing with alpha = 0.2303	6.7967	5.04166	47.5017
Brown's linear exp. smoothing with alpha = 0.1006	6.84085	5.03122	47.4461
Holt's linear exp. smoothing with alpha = 0.1247 and beta = 0.1451	6.90477	4.87344	42.7063
Brown's quadratic exp. smoothing with alpha = 0.0636	6.89329	5.07301	48.5641
Winter's exp. smoothing with alpha = 0.2086, beta = 0.0834, gamma = 0.1744	7.12028	5.28484	58.1989
ARIMA(0,0,0)x(2,1,2)12	4.76095	3.85657	28.3341
<b>ARIMA(2,1,0)x(2,1,2)12</b>	<b>4.72003</b>	<b>3.64274</b>	<b>26.3177</b>
ARIMA(0,0,1)x(2,1,2)12	4.8306	3.88196	27.5077
ARIMA(1,0,0)x(2,1,2)12	4.85863	3.90788	28.9779
ARIMA(1,0,1)x(2,1,2)12	4.80716	3.72139	28.5118

وبعد ذلك تم استخدام أفضل نموذج من السابقة ومقارنته مع أربعة نماذج حديثة لعدد 60 شهر، ويظهر جدول (1) مقاييس الجودة للنموذج المقترن والتي تؤكد جودة ملائمة النموذج لتوفيق للبيانات باستخدام المعايير الإحصائية لقياس قدرة النموذج على التنبؤ واستخدمت المعايير التالية (العباسي، 2011، 2003 :

1- جذر متوسط مربع الخطأ (RMSE)

2- المتوسط النسبي للخطأ المطلق (MAPE)

3- متوسط القيمة المطلقة للخطأ (MAE)

4- معامل ثيل (T.C)

5- معامل التحديد ( $R^2$ )

6- مؤشر الدلالة (TS)

7- معامل الارتباط الذاتي الأول ( $p_1$ ).

وأثبتت المعايير المستخدمة للحكم على أفضلية النموذج، ولنتائج الباقي وعشوائيتها أن نموذج الشبكات العصبية يعد الأفضل والأكثر ملائمة للبيانات المستخدمة خلال الفترة 2008-2013.

جدول (2) مقارنة لمؤشرات الجودة للنماذج الخمس المستخدمة للتوفيق لاعداد القضايا المعروضة على المحكمة

الدستورية العليا 2008-2013

نموذج	جذر متوسط مربع الخطأ (RMSE)	المتوسط النسبي للخطأ المطلق (MAPE)	متوسط القيمة المطلقة للخطأ (MAE)	معامل ثيل Theil	معامل التحديد $R^2$	مؤشر الدلالة TS	الارتباط الذاتي الأول $p_1$
الانحدار المتعدد	32.490	76.399	-23.804	-0.836	14.20%	4.03	0.1831
الانحدار ال بواسني	-32.610	70.439	-23.921	-0.838	17.36%	4.01	0.2213
السلسل الزمنية (اريما)	-35.310	33.157	-26.139	-0.895	71.60%	2.39	0.2493
الشبكات العصبية	-36.119	23.685	-27.229	-0.905	85.53%	3.49	-0.0714

## ٦. الخلاصة:

مما سبق يتضح:

- إن استخدام نموذج الشبكات العصبية في التنبؤ، ورسم الخطط سواء الطويلة الأجل والقصيرة الأجل لما يتميز به هذا النموذج من سرعة ودقة في البيانات أكثر منه في الأساليب الإحصائية التقليدية.
- من خلال التطبيق لكل من النماذج الإحصائية التقليدية والشبكات العصبية الاصطناعية ANN يتبين لنا أن الشبكات العصبية قد تميزت عن الأساليب الإحصائية التقليدية بأن لديها منهجية في عدم الاعتماد على الخطية في البيانات.
- أن الشبكات العصبية الاصطناعية أكثر دقة وكفاءة في التنبؤ عن الأساليب الإحصائية التقليدية حيث وصلت الشبكات لمعدل مرتفع وعالي من الدقة مع بقاء أفضليتها في التنبؤ للسلسل الزمنية الطويلة والتي لا يوجد بها اثر واضح للموسمية او الارتباط الذاتي.
- يجب على كل من يقوم بدراسة يتطلب فيها نظرة مستقبلية أن يقوم باستخدام الشبكات العصبية وأن يتم تحليلها باستخدام الأساليب الإحصائية الحديثة، وذلك لتحقيق الاستفادة القصوى منها حيث أن الشبكات لديها السرعة والدقة.
- وجد أن الشبكات العصبية تتتفوق على النماذج التقليدية بدرجة ملحوظة، ويعنى آخر ونظراً لمنهجية الشبكات العصبية في اعتمادها على غير الخطية فإن أداؤها أفضل مقارنة بالنماذج التقليدية ، وينتج أيضاً أنه يمكن تطبيق الشبكات العصبية بنجاح في التنبؤ بالسلسل الزمنية الشهرية الطويلة والتي تتسم بالموسمية أو الارتباط الذاتي.

في النهاية تكون قد حققنا هدف البحث وهو التعريف بأساليب التنبؤ بالبيانات وأنواعها، كما استعرضنا المجالات المختلفة التي استُخدمت فيها تلك الأساليب بنجاح. وقد اعتمد البحث على دراسة العديد من البحوث التي تتضمن أساليب مختلفة للتنبؤ بالبيانات والتعريف بها وتطبيقاتها وشروط استخدامها بشكل مبسط، ومقارنة بين أفضلية أساليب التنبؤ الحديثة وأساليب التقليدية، وأنه يوضح أن أسلوب الشبكات العصبية يعد أفضل النماذج المستخدمة مقارنة بالنماذج التقليدية لبناء نموذج لعدد القضايا المعروضة على المحكمة الدستورية العليا في مصر خلال الفترة 2008-2013.

## المراجع

١. العباسى، عبد الحميد محمد (2012)، قوة العمل الحكومية الكويتية: الواقع والعوامل المؤثرة خلال الفترة ١٩٩٣-٢٠١١، المجلة الإحصائية المصرية، معهد الدراسات والبحوث الإحصائية - القاهرة - مصر، مجلد (٥٦) العدد (٢)، ديسمبر ٢٠١٢ ص (٣٠ - ٤٦).
٢. العباسى، عبد الحميد محمد (٢٠١٠)، التحليل الحديث للسلسل الزمنية باستخدام Eviwes، معهد الدراسات والبحوث الإحصائية- القاهرة.
٣. العباسى، عبد الحميد محمد (٢٠٠٤)، "المقارنة بين استخدام الشبكات العصبية وساريما للتنبؤ بأعداد الوفيات الشهرية الناتجة عن حوادث المرور بالكويت" ، المجلة العربية للعلوم الإدارية ، الكويت ، مجلد (٣) العدد (١١)، ص (٣٣٣ - ٣٥٩).
٤. الوزير، رزق السيد وسمري، حاتم عبد الواحد (٢٠١٢)، "اساليب التقريب في البيانات: الطرق المعلمية واللامعلمية" معهد الدراسات والبحوث الإحصائية-المجلة المصرية للسكان وتنظيم الاسرة، مجلد ٤٥ العدد ديسمبر ٢٠١٢، ص (٦٥-٨٥).

- [1] Analytic Hierarchy Process  
[http://en.wikipedia.org/wiki/Analytic\\_Hierarchy\\_Process](http://en.wikipedia.org/wiki/Analytic_Hierarchy_Process)
- [2] Analytic Hierarchy Process (AHP) Tutorial  
<http://www.cs.toronto.edu/~sme/CSC340F/slides/tutorial-prioritization.pdf>
- [3] Backpropagation  
<http://en.wikipedia.org/wiki/Backpropagation>
- [4] C4.5 algorithm  
[http://en.wikipedia.org/wiki/C4.5\\_algorithm](http://en.wikipedia.org/wiki/C4.5_algorithm)
- [5] CHAID  
<http://en.wikipedia.org/wiki/CHAID>
- [6] Christos Stergiou and Dimitrios Siganos. Neural Network  
[http://www.doc.ic.ac.uk/~nd/surprise\\_96/journal/vol4/cs11/report.html#Neural\\_Networks\\_in\\_Practice](http://www.doc.ic.ac.uk/~nd/surprise_96/journal/vol4/cs11/report.html#Neural_Networks_in_Practice)
- [7] CLPFA (2008). Nearest Neighbours Model: Methodology Note and Instructions  
[http://www.cipfastats.net/default\\_view.asp?content\\_ref=2748](http://www.cipfastats.net/default_view.asp?content_ref=2748)
- [8] Conjugate gradient method  
[http://en.wikipedia.org/wiki/Conjugate\\_gradient\\_method](http://en.wikipedia.org/wiki/Conjugate_gradient_method)
- [9] Conventional programming  
[http://www.pcmag.com/encyclopedia\\_term/0,2542,t=conventional+programming&i=40325,00.asp](http://www.pcmag.com/encyclopedia_term/0,2542,t=conventional+programming&i=40325,00.asp)
- [10] Cornelius T. Leondes (2002). Expert systems: the technology of knowledge management and decision making for the 21st century, *Academic Press*, pp. 1-22.
- [11] CRISP-DM (2003), CRoss Industry Standard Process for Data Mining  
<http://www.crisp-dm.org>.
- [12] Delta rule (gradient descent)  
[http://en.wikipedia.org/wiki/Delta\\_rule](http://en.wikipedia.org/wiki/Delta_rule)
- [13] Fayyad; U.M., Piatetsky-Shapiro; G., Smyth; P. and Uthurusamy; R (eds) (1996a), Advances in Knowledge Discovery and Data Mining, *AAAI Press*.
- [14] Fayyad; U.M., Piatetsky-Shapiro; G., and Smyth; P. (1996b), From Data Mining to Knowledge Discovery: An Overview. In Fayyad; U.M., Piatetsky-Shapiro; G., Smyth; P. and Uthurusamy; R (eds), Advances in Knowledge Discovery and Data Mining , *AI, DDM, AAAI/MIT Press*, pp. 1-34.
- [15] Fayyad; U.M., Piatetsky-Shapiro; G., and Smyth; P. (1996c), The KDD process for extracting useful knowledge from volumes of data, Communications of the ACM, 39 (11), pp. 27-34.
- [16] Fayyad; U.M., Piatetsky-Shapiro; G., and Smyth; P. (1996d), Knowledge Discovery and Data Mining: Towards a unifying framework, *AI, DDM, AAAI/MIT Press*, pp. 82-88.
- [17] Fayyad; U.M., Piatetsky-Shapiro; G., and Smyth; P. (1996e), From data mining to knowledge discovery in databases, *AI Magazine*, 17, (3), pp. 37-54.
- [18] Friedman, J. H. (1991). "Multivariate Adaptive Regression Splines" *Annals of Statistics*, 19 (1): 1-67.  
[doi:10.1214/aos/1176347963](https://doi.org/10.1214/aos/1176347963). [MR1091842](#). [Zbl 0765.62064](#).
- [19] Genetic algorithm

- [http://en.wikipedia.org/wiki/Genetic\\_algorithm](http://en.wikipedia.org/wiki/Genetic_algorithm)
- [20] Generalized additive model  
[http://en.wikipedia.org/wiki/Generalized\\_additive\\_model](http://en.wikipedia.org/wiki/Generalized_additive_model)
- [21] Giudici; P. (2003), Applied Data Mining: Statistical Methods for Business and Industry, *John Wiley & Sons Ltd.*
- [22] Hastie; T., Tibshirani; R. (1993). "Varying-Coefficient Models" *Journal of the Royal Statistical Society, Series B*, 55, pp. 757–796.
- [23] Hastie; T., Tibshirani; R., Friedman; J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2<sup>nd</sup> Edition, *Springer Series in Statistics*.
- [24] Ichimura, H. (1993). "Semiparametric Least Squares (SLS) and Weighted SLS Estimation of Single Index Models" *Journal of Econometrics*, 58, pp. 71–120. doi:10.1016/0304-4076(93)90114-K.
- [25] Jackson, Peter (1998). Introduction to Expert Systems, 3<sup>rd</sup> ed., *Addison Wesley*, p. 2, ISBN 978-0-201-87686-4.
- [26] John Fox (2002). Nonparametric regression  
<http://cran.r-project.org/doc/contrib/Fox-Companion/appendix-nonparametric-regression.pdf>
- [27] Koch; C. G., Isle; B. A., Butler; A. W. (1988). "Intelligent user interface for expert systems applied to power plant maintenance and troubleshooting" *IEEE Transactions on Energy Conversion*, PP. 3- 71.
- [28] Mark E. Irwin (2005). Generalized Additive Models, *Harvard University*  
<http://www.markirwin.net/stat135/Lecture/Lecture34.pdf>
- [29] Mark E. Irwin (2005). Non Parametric Regression, *Harvard University*  
<http://www.markirwin.net/stat135/Lecture/Lecture33.pdf>
- [30] Mathematical optimization  
[http://en.wikipedia.org/wiki/Mathematical\\_optimization](http://en.wikipedia.org/wiki/Mathematical_optimization)
- [31] Mathematical Programming  
<http://www.me.utexas.edu/~jensen/ORMM/frontpage/pdf/mathprog.pdf>
- [32] Multivariate adaptive regression splines  
[http://en.wikipedia.org/wiki/Multivariate\\_adaptive\\_regression\\_splines](http://en.wikipedia.org/wiki/Multivariate_adaptive_regression_splines)
- [33] Multilayer Perceptron Neural Networks  
<http://www.dtreg.com/mlfn.htm>
- [34] Neural Network Software *For researchers, data mining experts and predictive analysts*  
<http://www.alyuda.com/products/neurointelligence/neural-network-applications.htm>
- [35] Neural Network Toolbox, Levenberg-Marquardt (trainlm)  
[http://www.caspur.it/risorse/softappl/doc/matlab\\_help/toolbox/nnet/backpr11.html](http://www.caspur.it/risorse/softappl/doc/matlab_help/toolbox/nnet/backpr11.html)
- [36] Newton's Telecom Dictionary (2010), Harry Newton, CMP Books,  
<http://www.cmpbooks.com>.
- [37] Nonparametric regression  
[http://en.wikipedia.org/wiki/Nonparametric\\_regression](http://en.wikipedia.org/wiki/Nonparametric_regression)
- [38] Nwigbo Stella and Agbo Okechukwu Chuks (2011). "Exert system: a catalyst in educational development in Nigeria," *Proceedings of the 1st International Technology, Education and Environment Conference, (c) African Society for Scientific Research (ASSR)*  
<http://www.harmars.com/admin/pics/261.pdf>