

الرموز ومفاهيم أساسية

□ **الفرضية (Hypothesis)** - الفرضية، ويرمز لها بـ h_θ ، هي النموذج الذي نختاره. إذا كان لدينا المدخل $x^{(i)}$ ، فإن المخرج الذي سيتوقعه النموذج هو $h_\theta(x^{(i)})$.

□ **دالة الخسارة (Loss function)** - دالة الخسارة هي الدالة $L: (z, y) \in \mathbb{R} \times Y \mapsto L(z, y) \in \mathbb{R}$ التي تأخذ كمداخلات القيمة المتوقعة z والقيمة الحقيقية y وتعطينا الاختلاف بينهما. الجدول التالي يحتوي على بعض دوال الخسارة الشائعة:

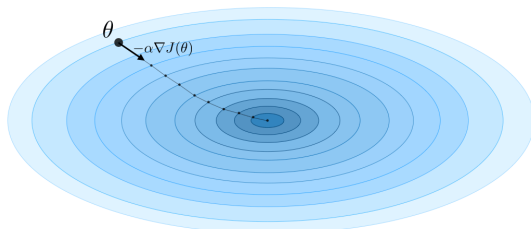
خطأ أصغر تربيع (Least squared error)	خسارة لوجستية (Logistic loss)	خسارة مفصلية (Hinge loss)	الانتروبيا التقاطعية (Cross-entropy)
$\frac{1}{2}(y - z)^2$	$\log(1 + \exp(-yz))$	$\max(0, 1 - yz)$	$-[y \log(z) + (1 - y) \log(1 - z)]$
الانحدار الخطي (Linear regression)	الانحدار اللوجستي (Logistic regression)	آلة المتجهات الداعمة (SVM)	الشبكات العصبية (Neural Network)

□ **دالة التكلفة (Cost function)** - دالة التكلفة J تستخدم عادة لتقييم أداء نموذج ما، ويتم تعريفها مع دالة الخسارة L كالتالي:

$$J(\theta) = \sum_{i=1}^m L(h_\theta(x^{(i)}), y^{(i)})$$

□ **النزول الاشتقاقي (Gradient descent)** - لنعرّف معدل التعلم $\alpha \in \mathbb{R}$ ، يمكن تعريف القانون الذي يتم تحديث خوارزمية النزول الاشتقاقي من خلاله باستخدام معدل التعلم ودالة التكلفة J كالتالي:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$



مرجع سريع للتعلم المُوجّه

افشين عميدى و شروين عميدى

١٤ ربيع الثاني، ١٤٤١

تمت الترجمة بواسطة فارس القنيعير. تمت المراجعة بواسطة زيد اليافعي.

مقدمة للتعلم المُوجّه

إذا كان لدينا مجموعة من نقاط البيانات $\{x^{(1)}, \dots, x^{(m)}\}$ مرتبطة بمجموعة مخرجات $\{y^{(1)}, \dots, y^{(m)}\}$ ، نريد أن نبني مُصنّف يتعلم كيف يتوقع y من x .

□ **نوع التوقع** - أنواع نماذج التوقع المختلفة موضحة في الجدول التالي:

الانحدار (Regression)	التصنيف (Classification)
المُخرَج	مستمر
أمثلة	انحدار خطي (Linear regression) آلة المتجهات الداعمة (SVM) بايز البسيط (Naive Bayes)

□ **نوع النموذج** - أنواع النماذج المختلفة موضحة في الجدول التالي:

نموذج توليدي (Generative)	نموذج تمييزي (Discriminative)
الهدف	التقدير المباشر لـ $P(y x)$
ماذا يتعلم	حدود القرار
توضيح	
أمثلة	الانحدار (Regression)، آلة المتجهات الداعمة (SVM)
GDA، بايز البسيط (Naive Bayes)	

التصنيف والانحدار اللوجستي

□ دالة سيجمويد (Sigmoid) – دالة سيجمويد g ، وتعرف كذلك بالدالة اللوجستية، تعرّف كالتالي:

$$\forall z \in \mathbb{R}, \quad g(z) = \frac{1}{1 + e^{-z}} \in]0,1[$$

□ الانحدار اللوجستي (Logistic regression) – نفترض هنا أن $y|x; \theta \sim \text{Bernoulli}(\phi)$. فيكون لدينا:

$$\phi = p(y = 1|x; \theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

ملاحظة: ليس هناك حل رياضي مغلق للانحدار اللوجستي.

□ انحدار سوفت ماكس (Softmax) – ويطلق عليه الانحدار اللوجستي متعدد الأصناف (multiclass logistic regression). يستخدم لتعميم الانحدار اللوجستي إذا كان لدينا أكثر من صنفين. في العرف يتم تعيين $\theta_K = 0$. بحيث تجعل مُدخل بيرنولي ϕ_i (Bernoulli) لكل فئة i يساوي:

$$\phi_i = \frac{\exp(\theta_i^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

النماذج الخطية العامة (GLM - Generalized Linear Models)

□ العائلة الأسية (Exponential family) – يطلق على صف من التوزيعات (distributions) بأنها تنتمي إلى العائلة الأسية إذا كان يمكن كتابتها بواسطة مُدخل قانوني (η canonical parameter)، إحصاء كافي (sufficient statistic) $T(y)$ ، ودالة تجزئة لوغاريتمية $a(\eta)$ ، كالتالي:

$$p(y; \eta) = b(y) \exp(\eta T(y) - a(\eta))$$

ملاحظة: كثيراً ما سيكون $y = T(y)$. كذلك فإن $\exp(-a(\eta))$ يمكن أن تفسر كمدخل تسوية (normalization) للتأكد من أن الاحتمالات يكون حاصل جمعها يساوي واحد.

تم تلخيص أكثر التوزيعات الأسية استخداماً في الجدول التالي:

التوزيع	η	$T(y)$	$a(\eta)$	$b(y)$
برنولي (Bernoulli)	$\log\left(\frac{\phi}{1-\phi}\right)$	y	$\log(1 + \exp(\eta))$	1
جاوسي (Gaussian)	μ	y	$\frac{\eta^2}{2}$	$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$
بواسون (Poisson)	$\log(\lambda)$	y	e^η	$\frac{1}{y!}$
هندسي (Geometric)	$\log(1 - \phi)$	y	$\log\left(\frac{e^\eta}{1 - e^\eta}\right)$	1

ملاحظة: في النزول الاشتقاقي العشوائي (Stochastic gradient descent (SGD)) يتم تحديث المُعاملات (parameters) بناءً على كل عينة تدريب على حدة، بينما في النزول الاشتقاقي الحزمي (batch gradient descent) يتم تحديثها باستخدام حُرْم من عينات التدريب.

□ الأرجحية (Likelihood) – تستخدم أرجحية النموذج $L(\theta)$ ، حيث أن θ هي المُدخلات، للبحث عن المُدخلات θ الأحسن عن طريق تعظيم (maximizing) الأرجحية. عملياً يتم استخدام الأرجحية اللوغاريتمية (log-likelihood) $\ell(\theta) = \log(L(\theta))$ حيث أنها أسهل في التحسين (optimize). فيكون لدينا:

$$\theta^{\text{opt}} = \arg \max_{\theta} L(\theta)$$

□ خوارزمية نيوتن (Newton's algorithm) – خوارزمية نيوتن هي طريقة حسابية للعثور على θ بحيث يكون $\ell'(\theta) = 0$. قاعدة التحديث للخوارزمية كالتالي:

$$\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

ملاحظة: هناك خوارزمية أعم وهي متعددة الأبعاد (multidimensional)، يطلق عليها خوارزمية نيوتن-رافسون (Newton-Raphson)، ويتم تحديثها عبر القانون التالي:

$$\theta \leftarrow \theta - \left(\nabla_{\theta}^2 \ell(\theta)\right)^{-1} \nabla_{\theta} \ell(\theta)$$

الانحدار الخطي (Linear regression)

هنا نفترض أن $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$

□ المعادلة الطبيعية/الناظرية (Normal) – إذا كان لدينا المصفوفة X ، القيمة θ التي تقلل من دالة التكلفة يمكن حلها رياضياً بشكل مغلق (closed-form) عن طريق:

$$\theta = (X^T X)^{-1} X^T y$$

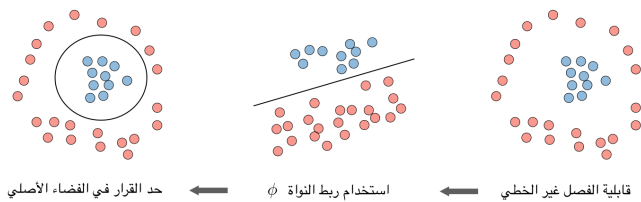
□ خوارزمية أصغر معدل تربيع LMS – إذا كان لدينا معدل التعلم α ، فإن قانون التحديث لخوارزمية أصغر معدل تربيع (Least Mean Squares (LMS)) لمجموعة بيانات m عينة، ويطلق عليه قانون تعلم ويدرو-هوف (Widrow-Hoff)، كالتالي:

$$\forall j, \quad \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] x_j^{(i)}$$

ملاحظة: قانون التحديث هذا يعتبر حالة خاصة من النزول الاشتقاقي (Gradient descent).

□ الانحدار الموزون محلياً (LWR) – الانحدار الموزون محلياً (Locally Weighted Regression)، ويعرف بـ LWR، هو نوع من الانحدار الخطي يزن كل عينة تدريب أثناء حساب دالة التكلفة باستخدام $w^{(i)}(x)$ ، التي يمكن تعريفها باستخدام المُدخل (parameter) $\tau \in \mathbb{R}$ كالتالي:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$



ملاحظة: نقول أننا نستخدم "حيلة النواة" (kernel trick) لحساب دالة التكلفة عند استخدام النواة لأننا في الحقيقة لا نحتاج أن نعرف التحويل الصريح ϕ ، الذي يكون في الغالب شديد التعقيد. ولكن، نحتاج أن فقط أن نحسب القيم $K(x, z)$.

□ **اللاغرانجي (Lagrangian)** - يتم تعريف اللاگرانجي $\mathcal{L}(w, b)$ على النحو التالي:

$$\mathcal{L}(w, b) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

ملاحظة: المعاملات β_i (coefficients) يطلق عليها مضروبات لاغرانج (Lagrange multipliers).

التعلم التوليدي (Generative Learning)

النموذج التوليدي في البداية يحاول أن يتعلم كيف تم توليد البيانات عن طريق تقدير $P(x|y)$ ، التي يمكن حينها استخدامها لتقدير $P(y|x)$ باستخدام قانون بايز (Bayes' rule).

تحليل التمايز الجاوسي (Gaussian Discriminant Analysis)

□ **الإطار** - تحليل التمايز الجاوسي يفترض أن y و $x|y = 0$ و $x|y = 1$ بحيث يكونوا كالتالي:

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \text{و} \quad x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$$

□ **التقدير** - الجدول التالي يلخص التقديرات التي يمكننا التوصل لها عند تعظيم الأرجحية (likelihood):

$\hat{\Sigma}$	$\hat{\mu}_j \quad (j = 0, 1)$	$\hat{\phi}$
$\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$	$\frac{\sum_{i=1}^m 1_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m 1_{\{y^{(i)}=j\}}}$	$\frac{1}{m} \sum_{i=1}^m 1_{\{y^{(i)}=1\}}$

بايز البسيط (Naive Bayes)

□ **الافتراض** - يفترض نموذج بايز البسيط أن جميع الخصائص لكل عينة بيانات مستقلة (independent):

$$P(x|y) = P(x_1, x_2, \dots | y) = P(x_1|y)P(x_2|y)\dots = \prod_{i=1}^n P(x_i|y)$$

□ **افتراضات GLMs** - تهدف النماذج الخطية العامة (GLM) إلى توقع المتغير العشوائي y كدالة لـ $x \in \mathbb{R}^{n+1}$ وتستند إلى ثلاثة افتراضات:

$$(1) \quad y|x; \theta \sim \text{ExpFamily}(\eta) \quad (2) \quad h_\theta(x) = E[y|x; \theta] \quad (3) \quad \eta = \theta^T x$$

ملاحظة: أصغر تربيع (least squares) الاعتيادي و الانحدار اللوجستي يعتبران من الحالات الخاصة للنماذج الخطية العامة.

آلة المتجهات الداعمة (Support Vector Machines)

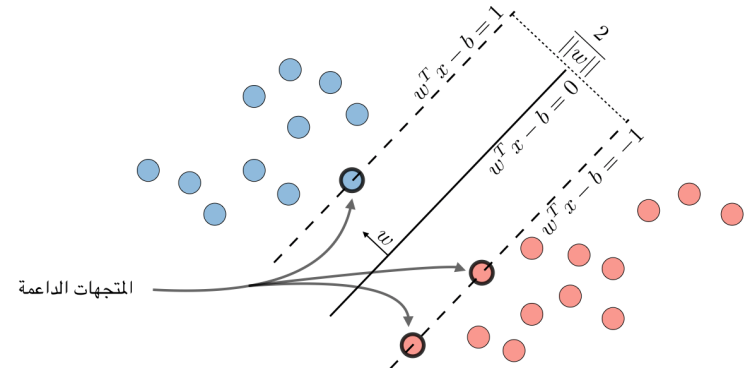
تهدف آلة المتجهات الداعمة (SVM) إلى العثور على الخط الذي يعظم أصغر مسافة إليه:

□ **مُصنّف الهامش الأحسن (Optimal margin classifier)** - يعرف مُصنّف الهامش الأحسن h كالتالي:

$$h(x) = \text{sign}(w^T x - b)$$

حيث $(w, b) \in \mathbb{R}^n \times \mathbb{R}$ هو الحل لمشكلة التحسين (optimization) التالية:

$$\min \frac{1}{2} \|w\|^2 \quad \text{بحيث أن} \quad y^{(i)}(w^T x^{(i)} - b) \geq 1$$



ملاحظة: يتم تعريف الخط بهذه المعادلة $w^T x - b = 0$.

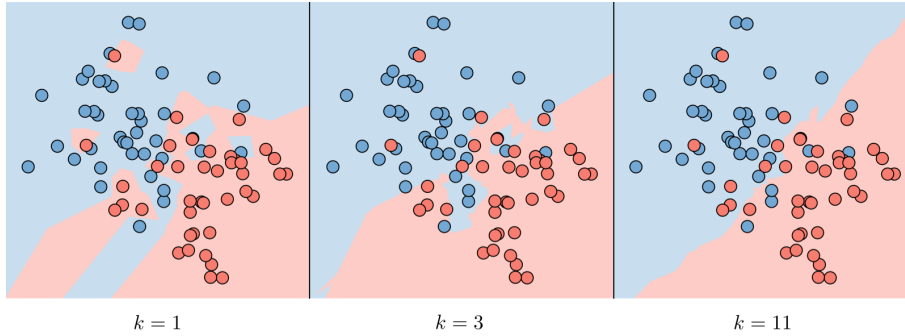
□ **الخسارة المفصلية (Hinge loss)** - تستخدم الخسارة المفصلية في حل SVM ويعرف على النحو التالي:

$$L(z, y) = [1 - yz]_+ = \max(0, 1 - yz)$$

□ **النواة (Kernel)** - إذا كان لدينا دالة ربط الخصائص ϕ ، يمكننا تعريف النواة K كالتالي:

$$K(x, z) = \phi(x)^T \phi(z)$$

عملياً، يمكن أن تُعرّف الدالة K عن طريق المعادلة $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$ ، ويطلق عليها النواة الجاوسية (Gaussian kernel)، وهي تستخدم بكثرة.



الحل - تعظيم الأرجحية اللوغاريتمية (log-likelihood) يعطينا الحلول التالية إذا كان $k \in \{0,1\}, l \in [1,L]$:

$$P(y = k) = \frac{1}{m} \times \#\{j | y^{(j)} = k\} \quad \text{و} \quad P(x_i = l | y = k) = \frac{\#\{j | y^{(j)} = k \text{ و } x_i^{(j)} = l\}}{\#\{j | y^{(j)} = k\}}$$

ملاحظة: بايز البسيط يستخدم بشكل واسع لتصنيف النصوص واكتشاف البريد الإلكتروني المزجج.

الطرق الشجرية (tree-based) والتجميعية (ensemble)

هذه الطرق يمكن استخدامها لكل من مشاكل الانحدار (regression) والتصنيف (classification).

التصنيف والانحدار الشجري (CART) - والاسم الشائع له أشجار القرار (decision trees)، يمكن أن يمثل كأشجار ثنائية (binary trees). من المزايا لهذه الطريقة إمكانية تفسيرها بسهولة.

الغابة العشوائية (Random forest) - هي أحد الطرق الشجرية التي تستخدم عدداً كبيراً من أشجار القرار مبنية باستخدام مجموعة عشوائية من الخصائص. بخلاف شجرة القرار البسيطة لا يمكن تفسير النموذج بسهولة، ولكن أدائها العالي جعلها أحد الخوارزميات المشهورة.

ملاحظة: أشجار القرار نوع من الخوارزميات التجميعية (ensemble).

التعزيز (Boosting) - فكرة خوارزميات التعزيز هي دمج عدة خوارزميات تعلم ضعيفة لتكوين نموذج قوي. الطرق الأساسية ملخصة في الجدول التالي:

التعزيز التكراري (Adaptive boosting)	التعزيز الاشتقاقي (Gradient boosting)
- يتم التركيز على مواطن الخطأ لتحسين النتيجة في الخطوة التالية. "Adaboost"	- يتم تدريب خوارزميات التعلم الضعيفة على الأخطاء المتبقية.

طرق أخرى غير بارامترية (non-parametric)

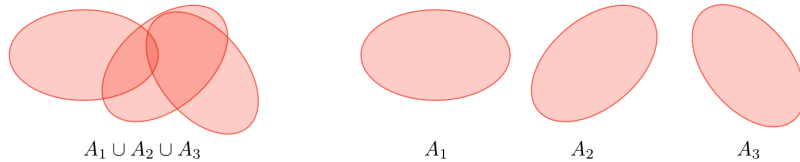
خوارزمية أقرب الجيران (k-nearest neighbors) - تعتبر خوارزمية أقرب الجيران، وتعرف بـ k -NN، طريقة غير بارامترية، حيث يتم تحديد نتيجة عينة من البيانات من خلال عدد k من البيانات المجاورة في مجموعة التدريب. ويمكن استخدامها للتصنيف والانحدار.

ملاحظة: كلما زاد المدخل k ، كلما زاد الانحياز (bias)، وكلما نقص k ، زاد التباين (variance).

نظرية التعلم

حد الاتحاد (Union bound) - لنجعل A_1, \dots, A_k تمثل k حدث. فيكون لدينا:

$$P(A_1 \cup \dots \cup A_k) \leq P(A_1) + \dots + P(A_k)$$



مراجعة هوفدينج (Hoeffding) - لنجعل Z_1, \dots, Z_m تمثل m متغير مستقل وموزعة بشكل مماثل (iid) مأخوذة من توزيع برنولي (Bernoulli distribution) ذا مدخل ϕ . لنجعل $\hat{\phi}$ متوسط العينة (sample mean) و $\gamma > 0$ ثابت. فيكون لدينا:

$$P(|\hat{\phi} - \phi| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$

ملاحظة: هذه المتراجحة تعرف كذلك بحد تشرنوف (Chernoff bound).

خطأ التدريب - ليكن لدينا المُصنّف h ، يمكن تعريف خطأ التدريب $\epsilon(h)$ ، ويعرف كذلك بالخطر التجريبي أو الخطأ التجريبي، كالتالي:

$$\hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m 1_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

تقريباً صحيح احتمالياً (Probably Approximately Correct (PAC)) - هو إطار يتم من خلاله إثبات العديد من نظريات التعلم، ويحتوي على الافتراضات التالية:

- مجموعتي التدريب والاختبار يتبعان نفس التوزيع.
- عينات التدريب تؤخذ بشكل مستقل.

□ **مجموعة تكسيرية (Shattering Set)** - إذا كان لدينا المجموعة $S = \{x^{(1)}, \dots, x^{(d)}\}$ ، ومجموعة مُصنّفات \mathcal{H} ، نقول أن \mathcal{H} تكسر S (H shatters S) إذا كان لكل مجموعة علامات $\{y^{(1)}, \dots, y^{(d)}\}$ (labels) لدينا:

$$\exists h \in \mathcal{H}, \quad \forall i \in [1, d], \quad h(x^{(i)}) = y^{(i)}$$

□ **مبرهنة الحد الأعلى (Upper bound theorem)** - لنجعل \mathcal{H} فئة فرضية محدودة (finite hypothesis class) بحيث $|\mathcal{H}| = k$ ، و δ وحجم العينة m ثابتين. حينها سيكون لدينا، مع احتمال على الأقل $1 - \delta$ ، التالي:

$$\epsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}} \epsilon(h) \right) + 2 \sqrt{\frac{1}{2m} \log \left(\frac{2k}{\delta} \right)}$$

□ **بُعد فابنيك - تشرفونيكس (Vapnik-Chervonenkis - VC)** لفئة فرضية غير محدودة (infinite hypothesis class) \mathcal{H} ، ويرمز له بـ $VC(\mathcal{H})$ ، هو حجم أكبر مجموعة (set) التي تم تكسيورها بواسطة \mathcal{H} (shattered by \mathcal{H}). ملاحظة: بُعد فابنيك تشرفونيكس VC \mathcal{H} = {مجموعة التصنيفات الخطية في بُعدين} يساوي ٣.



□ **مبرهنة فابنيك (Vapnik theorem)** - ليكن لدينا \mathcal{H} ، مع $VC(\mathcal{H}) = d$ وعدد عيّنات التدريب m . سيكون لدينا، مع احتمال على الأقل $1 - \delta$ ، التالي:

$$\epsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}} \epsilon(h) \right) + O \left(\sqrt{\frac{d}{m} \log \left(\frac{m}{d} \right)} + \frac{1}{m} \log \left(\frac{1}{\delta} \right) \right)$$