

7 Apprentissage

Il y a trois grands avantages de modèles probabilistes :

- ★ **[apprentissage]** on peut ajuster les paramètres du modèle selon des principes mathématiques bien fondés (p.e., maximisation de vraisemblance)
- ★ **[confidence]** l'inférence est accompagnée de mesures statistiques de confiance (p.e., probabilité postérieure de colonnes d'un alignement)
- ★ **[choix]** on peut choisir l'architecture du modèle selon des mesures de *fit* statistique (p.e., pénalisation par description minimale)

7.1 Vraisemblance

Soit $X = (x_1, x_2, \dots, x_n)$ un jeu de données observées, ou un *échantillon*. Un modèle probabiliste M se décrit par son architecture (observation + variables aléatoires inobservées avec dépendances) et par la vecteur de ses paramètres θ . Le modèle définit la distribution d'observations. Typiquement, chaque x_i est tirée indépendamment de la même distribution $p(x|\theta)$. Ainsi, on définit la vraisemblance

$$L(\theta) = \mathbb{P}\left\{\text{on observe } X \mid \text{modèle } M \text{ avec paramètres } \theta\right\} = \prod_{i=1}^n p(x_i|\theta). \quad (7.1)$$

7.2 Vraisemblance et modèles graphiques

Le modèle probabiliste comprend un ensemble de variables aléatoires inobservées Y_1, \dots, Y_m et l'observation Y_0 même :

$$p(x|\theta) = \sum_{y_1, y_2, \dots, y_m} \mathbb{P}\left\{Y_0 = x \mid Y_1 = y_1, \dots, Y_m = y_m; \theta\right\} \cdot \mathbb{P}\left\{Y_1 = y_1, \dots, Y_m = y_m \mid \theta\right\}.$$

La somme suggère une temps exponentiel en m , mais on peut souvent exploiter les indépendances conditionnelles entre les Y_i pour accélérer le calcul. Considérons le graphe G de dépendances entre les variables aléatoires. Les nœuds de G correspondent aux variables $Y_i: i = 0, 1, \dots, m$. Les arcs de G sont orientés. L'ensemble des arcs entrants au nœud j , dénoté par $\text{pred}(j) = \{i: ij \in G\}$, définit l'indépendance conditionnelle :

$$\mathbb{P}\left\{Y_j \mid Y_i: i \neq j\right\} = \mathbb{P}\left\{Y_j \mid Y_i: i \in \text{pred}(j)\right\}.$$

Un *modèle graphique* est un modèle probabiliste dont le graphe de dépendances ne contient aucun cycle. Exemple : évolution d'un caractère dans un arbre phylogénétique. On associe une valeur inobservée avec chaque nœud ancestral. On observe le caractère aux feuilles. L'arbre montre les dépendances : $\text{pred}(j)$ contient le parent du nœud j . La vraisemblance se calcule rapidement dans l'ordre inverse du tri topologique.

[algorithme de Felsenstein pour calculer vraisemblance sur un arbre phylogénétique]

7.3 Maximisation de la vraisemblance

Une méthode standard d'ajuster les paramètres du modèle est de maximiser la vraisemblance. La valeur des paramètres θ_{ML} qui maximise $L(\theta)$ est un estimateur *consistent* sous quelques conditions génériques¹ sur

¹ Wald A. "Note on the consistency of the maximum likelihood estimate", *Annals of Mathematical Statistics*, **20** :595–601, 1949

la topologie de l'espace de paramètres : sur un échantillon généré par le même modèle, quand $n \rightarrow \infty$, θ_{ML} converge vers les vraies valeurs. Une des conditions très importantes est l'*identifiabilité* du modèle : $p(x|\theta)$ doit être une distribution différente pour toute valeur de θ .

Optimisation numérique. Le plus simple est de considérer la vraisemblance de l'équation (7.1) comme une fonction de la vecteur de paramètres. On maximise alors sa valeur par une méthode numérique standard comme la méthode de la bisection ou la méthode du gradient conjugué. Dans l'implantation du calcul, il faut payer attention aux issues de débordement numérique. Typiquement, on travaille avec la log-vraisemblance $\log L(\theta) = \sum_{i=1}^n \log p(x_i|\theta)$.

7.4 Espérance-maximisation

L'algorithme d'espérance-maximisation (EM) sert à optimiser la vraisemblance d'un modèle probabiliste. L'algorithme calcule une séquence d'estimateurs $\theta_0, \theta_1, \theta_2, \dots$ en une itération jusqu'à convergence. En chaque itération, on exécute un étape E («espérance») et un étape M («maximisation»). Pour simplifier la discussion, soit X l'observation aléatoire, et Y la vecteur de variables aléatoires inobservées. On observe $X = x$. En itération $i = 1, 2, \dots$, on fait :

★ **Étape E.** On définit la fonction

$$\begin{aligned} Q_i(\theta) &= \mathbb{E}_{Y|X, \theta_{i-1}} \log \mathbb{P}\{Y; X = x \mid \theta\} \\ &= \sum_y \mathbb{P}\{Y = y \mid X = x; \theta_{i-1}\} \log \mathbb{P}\{Y; X = x \mid \theta\}. \end{aligned} \quad (7.2)$$

★ **Étape M.** On maximise Q :

$$\theta_i \leftarrow \arg \max_{\theta} Q_i(\theta). \quad (7.3)$$

L'algorithme EM converge vers le maximum de la vraisemblance.

L'algorithme de Baum-Welch pour les HMM suit cette technique générique. Les variables inobservées sont les états. Supposons qu'on observe la séquence $\mathbf{x} = x_1 x_2 \dots x_n$, générée par la séquence d'états inconnue Q_1, \dots, Q_n . Soit $\gamma_i(q) = \mathbb{P}\{Q_i = q \mid \mathbf{x}\}$, la probabilité postérieure de l'état Q_i , ce qu'on calcule par l'algorithme forward-backward. Soit $e(x|q) = \mathbb{P}\{X_i = x \mid Q_i = q\}$ la probabilité d'émission du caractère x en état q . En équation (7.2), $e(x|q)$ paraît dans les termes :

$$Q(e) = \sum_{i=1}^n \gamma_i(q) \log e(x_i|q).$$

Pour un alphabet à $r < \infty$ lettres, on a $Q(e_1, \dots, e_r) = \sum_{a=1}^r n_a \log e_a$, où $n_a = \sum_{i: x_i=a} \gamma_i(q)$ est le nombre d'émissions a en état q , et $e_a = e(a|q)$ et la probabilité d'émission du caractère a . Or, cette expression est maximisée avec le choix $e_a = n_a / \sum_{j=1}^r n_j$ dans l'étape M, et c'est exactement la formule dans l'apprentissage Baum-Welch.

7.5 Algorithme des k -moyennes

L'algorithme des k -moyennes sert à partitionner un jeu de données x_1, x_2, \dots, x_n de d dimensions ($x_i \in \mathbb{R}^d$). Le modèle probabiliste sous-jacente comprend k centres $\mu_1, \mu_2, \dots, \mu_k$ inobservés. Pour générer un

W_(en)

W_(en)

seul point X observé, on choisit un des centres au hasard ($Y = 1, 2, \dots, k$), et on tire une valeur selon la distribution normale centrée à μ_Y . En étape E, on doit calculer

$$\rho_i(j) = \mathbb{P}\left\{Y_i = j \mid X_i = x_i; \mu_1, \dots, \mu_k\right\} \propto \exp\left(-\frac{(x_i - \mu_j)^2}{2}\right) :$$

c'est l'affectation «souple» de point i au centre j . En étape M, on trouve le choix de centres maximisant la fonction Q en (7.3) par $\mu_j \leftarrow \frac{\sum_i \rho_i(j)x_i}{\sum_i \rho_i(j)}$: c'est la vote des points pondérée selon l'affectation souple.

En une démarche alternative de groupage «dur» (évoquant la stratégie d'apprentissage Viterbi), on utilise $\rho_i(j) = 1$ pour le centre le plus proche ($(x_i - \mu_j)^2$ minimale) et $\rho_i(j) = 0$ pour tout autre centre. L'affectation de l'étape M est simplement $\mu_j \leftarrow \frac{\sum_i \{\rho_i(j)=1\} \cdot x_i}{\sum_i \{\rho_i(j)\}}$.

7.6 Bootstrapping

La maximisation de la vraisemblance nous fournit un seul choix optimal pour les paramètres. Parfois il est désirable d'étudier l'incertitude de valeurs de paramètres. Une approche générale consiste de simulations : on examine comment l'optimisation fonctionne sur un échantillon aléatoire. Le **bootstrapping** est applicable quand l'observation est iid. Étant donné la vraie observation x_1, x_2, \dots, x_n , on produit un échantillon simulé y_1, \dots, y_n de même taille. Pour chaque y_i on choisit un des x_j au hasard (avec répétitions permises). Ensuite, on performe la procédure d'optimisation de paramètres sur l'échantillon simulé. On répète la procédure plusieurs fois ($N = 1000, 10000, 100000, \dots$) et on collecte des statistiques sur la distribution de valeurs optimisées dans les échantillons. On construit une intervalle de confiance en excluant les valeurs les plus grandes et les plus petites (p.e., pour 95%, on exclut 25-25 parmi 1000).

En «bootstrap paramétrique», l'échantillon simulé est construit à partir du modèle directement. Notamment, on génère n observations iid, en utilisant la distribution impliée par le modèle. L'échantillon ainsi produit est soumis à la procédure normale d'apprentissage. On collecte des statistiques en répétant la procédure.

7.7 Markov Chain Monte Carlo

Idéalement, on voudrait avoir de l'information complète sur la distribution postérieure des paramètres : on veut calculer $\mathbb{P}\{\theta \mid \mathbf{x}\}$. Une telle formulation de la question est possible si on introduit la distribution *a priori* des paramètres $p(\theta)$. Par la loi de Bayes,

$$\mathbb{P}\{\theta \mid \mathbf{x}\} = \frac{\mathbb{P}\{\mathbf{x} \mid \theta\} \cdot p(\theta)}{\mathbb{P}\{\mathbf{x}\}} \propto L(\theta) \cdot p(\theta). \quad (7.4)$$

Notez qu'en général, on ne peut pas directement calculer le dénominateur $\mathbb{P}\{\mathbf{x}\}$ car il requiert l'intégrale de la vraisemblance pour tout choix de θ . L'approche de méthodes *Markov chain Monte Carlo* (MCMC) est de définir une chaîne de Markov dont les états sont les valeurs possibles de θ . L'espace des états est l'espace de paramètres. On définit les transitions entre états de telle façon que la chaîne converge vers la distribution postérieure de θ . Dans d'autres mots, on produit un *échantillonnage* de l'espace de paramètres. Il s'agit d'une collection $\theta_0, \theta_1, \theta_2, \dots, \theta_N$ (états simulés de la chaîne) telle que

$$\lim_{N \rightarrow \infty} \mathbb{P}\{\theta_i = \theta\} = \mathbb{P}\{\theta \mid \mathbf{x}\}.$$

Pour achever la convergence vers la distribution souhaitée, on choisit des probabilités de transition pour lesquelles la distribution stationnaire est exactement $\mathbb{P}\{\theta \mid \mathbf{x}\}$. Rappel : la distribution stationnaire π (vue

comme vecteur ligne) d'une chaîne de Markov avec matrice de transition \mathbf{M} est sa vecteur propre avec $\pi \cdot \mathbf{M} = \pi$. Si la chaîne est apériodique et irréductible (on peut accéder à tout état à partir de n'importe quel état), alors elle converge vers la distribution stationnaire. Étant une vecteur propre, on veut alors qu'avec $\pi(\theta) = \mathbb{P}\{\theta \mid \mathbf{x}\}$,

$$\sum_a \pi(a) \mathbf{M}[a, b] = \pi(b)$$

en choisissant les cases de la matrice \mathbf{M} .

```

MCMC // (algorithme de MCMC générique)
M1 initialiser  $\theta_0$  // (p.e., au hasard selon distribution  $\pi(\theta)$ )
M2 for  $i \leftarrow 1, 2, \dots$  do
M3     générer  $\theta$  aléatoire selon la distribution conditionnelle de transitions  $\mathbf{M}[\theta_i, \cdot]$ 
M4      $\theta_i \leftarrow \theta$  // (écrire valeur échantillonnée dans un fichier)

```

Algorithme Metropolis-Hastings. Dans l'algorithme de Metropolis-Hastings, on propose une nouvelle vecteur de paramètres θ' selon une règle quelconque de propositions $Q(\theta \rightarrow \theta')$. On calcule le *Hastings ratio*

$$\alpha = \frac{\mathbb{P}\{\theta' \mid \mathbf{x}\} \cdot Q(\theta' \rightarrow \theta)}{\mathbb{P}\{\theta \mid \mathbf{x}\} \cdot Q(\theta \rightarrow \theta')}.$$

Si $\alpha \geq 1$, on accepte ($\theta_i \leftarrow \theta'$), sinon on rejette ($\theta_i \leftarrow \theta_{i-1}$) la proposition.

```

MH // (algorithme Metropolis-Hastings)
H1 initialiser  $\theta_0$  // (p.e., au hasard selon distribution  $\pi(\theta)$ )
H2 for  $i \leftarrow 1, 2, \dots$  do
H3     générer  $\theta'$  aléatoire selon la distribution de propositions  $Q(\theta_{i-1} \rightarrow \theta)$ 
H4      $\alpha \leftarrow \frac{L(\theta') \cdot \pi(\theta') \cdot Q(\theta' \rightarrow \theta_{i-1})}{L(\theta_{i-1}) \cdot \pi(\theta_{i-1}) \cdot Q(\theta_{i-1} \rightarrow \theta')}$ 
H5     if  $\alpha \geq 1$  ou  $\text{rand}() < \alpha$  then  $\theta_i \leftarrow \theta'$  else  $\theta_i \leftarrow \theta_{i-1}$ 

```

Notez qu'en Ligne H4, il suffit de calculer les vraisemblances parce que le dénominateur normalisant en (7.4) est le même pour θ_{i-1} et θ' .