

Apprentissage statistique et Big Data, focus sur l'algorithme online-EM

Olivier Cappé

Laboratoire Traitement et Communication de l'Information
CNRS, Télécom ParisTech, 75013 Paris

8 octobre 2013

1 Apprentissage statistique

2 Big Data

3 Online EM

L'apprentissage automatique

Apprentissage automatique (*machine learning*)

Apprendre à effectuer des tâches à partir d'exemples

De façon supervisée Avec des exemples de données et de résultats souhaités

↪ Classification

De façon non supervisée Avec uniquement des exemples de données

↪ Clustering

De façon séquentielle En traitant un flux de données de façon causale

↪ Prédiction séquentielle

Par extension, problèmes dont le traitement présente des analogies avec ce qui précède

↪ Régression en grande dimension

Quelques idées sur l'apprentissage statistique

Je m'intéresse aux **approches statistiques** de l'apprentissage



- 1 Comme en statistique usuelle j'utilise un **modèle probabiliste des données**

par exemple, $Y = f(x)\beta + \epsilon$, avec ϵ aléatoire

- 2 Mais ce qui m'intéresse c'est **prédire** | **classifier** | **reconstruire** plus qu'estimer

prédire Y' par $f(x')\hat{\beta}$ pour un nouvel x'
plus que $\hat{\beta}$ en lui même

Quelques idées sur l'apprentissage statistique ...

- 3 D'autant plus que je crois que le **modèle** utilisé est **instrumental** (probablement faux)

on doit avoir $Y = \mathcal{F}(x) + \epsilon$, avec \mathcal{F} assez compliquée

- 4 De ce fait, j'ai intérêt à faire **croître la complexité du modèle avec le nombre de données disponibles**

si j'observe Y_1, \dots, Y_n , j'utilise plutôt le modèle

$$Y_i = \sum_{j=1}^{k_n} f_j(x_i) \beta_j + \epsilon_i, \text{ où } k_n \uparrow \text{ avec } n$$

- 5 Le **choix / sélection de modèle** est crucial

choisir \hat{k}_n pour que $\sum_{j=1}^{\hat{k}_n} f_j(x') \hat{\beta}_j$ soit le plus proche de Y'

Quelques idées sur l'apprentissage statistique

6 L'optimisation est un outil central

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|Y - x\beta\|^2 \text{ plutôt que } \hat{\beta} = (x^T x)^{-1} x^T Y$$

$$\text{de façon à penser à } \underset{\beta: \|\beta\|_1 \leq \lambda}{\operatorname{argmin}} \|Y - x\beta\|^2$$

qu'il va falloir **résoudre numériquement** . . .

7 Comme les modèles sont faux j'ai intérêt à les **combiner**

J'estime $\hat{\beta}_1, \dots, \hat{\beta}_k$ séparément dans les
 modèles $Y = f_1(x)\beta_1 + \epsilon, \dots, Y = f_k(x)\beta_k + \epsilon,$
 j'utilise l'ensemble $\{f_1(x')\hat{\beta}_1, \dots, f_k(x')\hat{\beta}_k\}$ pour prédire Y'

Quelques idées sur l'apprentissage statistique

- 8 L'approche bayésienne permet de retrouver beaucoup des méthodes précédentes ($k_n \uparrow n$, $\|\beta\|_1 \leq \lambda$, $\{f_1(x')\hat{\beta}_1, \dots, f_k(x')\hat{\beta}_k\}$) ainsi que certaines de leurs variantes 

J'imagine que β (voire k) est une variable aléatoire

- 1 Apprentissage statistique
- 2 Big Data
 - Quelques idées générales
 - Quelques idées de modélisation
- 3 Online EM

- Il est trop tard pour s'interroger sur l'opportunité du Big* Data...
- Néanmoins son potentiel reste assez largement spéculatif

Search

[Journée de rencontre Big Data et décisionnel](#)
www.digitalplace.fr/fr/.../item/...bigdata-decisionnel Cached
 Accueil DigitalPlace Actualités Les Événements Les événements DigitalPlace Journée de rencontre **Big Data** ... Université de **Toulouse** ... **Olivier CAPPÉ**, CNRS ...

[Journée de rencontre Big Data, 8 oct. Toulouse](#)
perso.math.univ-toulouse.fr/bigdata Cached
 ... FREMIT ainsi que le CMI SID organise le mardi 8 octobre une journée de rencontre entre l'Université de **Toulouse** ... **Olivier Cappé**, ... et **Big Data**, focus ...

[Aurélien Garivier - Institut de Mathématiques de Toulouse](#)
www.math.univ-toulouse.fr/~agarivie/?q=node/34 Cached
 professional web page . Search form. Search

[Midi-Pyrénées Innovation » Journée de rencontre BIG DATA ...](#)
www.mp-i.fr/2013/09/journee-de-rencontre-big-data-et... Cached
 Journée de rencontre **BIG DATA** et ... ainsi que le CMI SID organisent une journée de rencontre entre l'Université de **Toulouse** ... 9h00 **Olivier Cappé**, ...

* Suffisamment grand pour espérer répondre à des questions complexes

Un phénomène qui impacte aussi la science

Big Data

Web, réseaux sociaux, commerce en ligne, systèmes d'information d'entreprise. . .

mais aussi

Sciences du vivant, sciences du climat, physique des particules, sciences humaines, . . .



The screenshot shows the homepage of the journal 'Database: The Journal of Biological Databases and Curation'. The header features the journal title in large white letters on a dark blue background. Below the header is a navigation bar with links for 'ABOUT THIS JOURNAL', 'CONTACT THIS JOURNAL', 'SUBSCRIPTIONS', 'CURRENT ISSUE', and 'ARCHIVE'. A breadcrumb trail indicates the current page is 'Database' under 'Life Sciences'. The main content area is titled 'READ THIS JOURNAL' and includes a welcome message, a 'Fully Open Access Journal' badge, and several links: 'View Current Content', 'Browse the Archive', 'Biocuration Virtual Issue', 'Biomart Virtual Issue', and 'Now Indexed in PubMed Central'. A paragraph discusses the challenges of big data in biological research, and a final paragraph states the journal's mission as an open access platform for database research and biocuration.

DATABASE The Journal of Biological Databases and Curation

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE

[Oxford Journals](#) > [Life Sciences](#) > [Database](#)

READ THIS JOURNAL

DATABASE
The Journal of Biological Databases and Curation

Welcome to Database: *The Journal of Biological Databases and Curation*

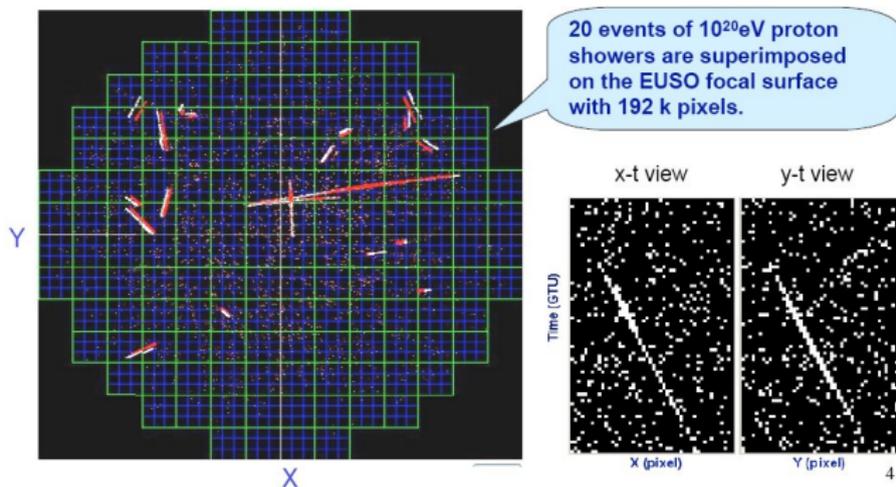
A Fully Open Access Journal

- [View Current Content](#)
- [Browse the Archive](#)
- [Biocuration Virtual Issue](#)
- [Biomart Virtual Issue](#)
- [Now Indexed in PubMed Central](#)

Huge volumes of primary data are archived in numerous open-access databases, and with new generation technologies becoming more common in laboratories, large datasets will become even more prevalent. The archiving, curation, analysis and interpretation of all of these data are a challenge. Database development and biocuration are at the forefront of the endeavor to make sense of this mounting deluge of data.

Database: The Journal of Biological Databases and Curation provides an open access platform for the presentation of novel ideas in database research and biocuration, and aims to help strengthen the bridge between database developers, curators, and users.

Vers une “science des données” ?



Expérience JEM-EUSO Capteur de $2 \cdot 10^5$ pixels / $2.5 \mu\text{s}$ pour détecter une centaine de particules de très haute énergie par jour (source Balázs Kégl, LAL/LRI)

↪ détecteur obtenu par des méthodes d'apprentissage statistique (boosting) à partir de données préliminaires et de simulations

Une demande très forte sur l'enseignement

Qui mêle des compétences en **mathématiques appliquées** (statistiques, optimisation) et en **informatique**

	A		B	
	A1	A2	B1	B2
P1	Concepts Fondamentaux de la Sécurité (30H, 3ECTS) INF721	Statistique (30H, 3ECTS) MDI 220	Bases de Données (30H, 3ECTS) INF225	Economie de l'Internet et des Données Personnelles (20H, 3ECTS) SES720
P2	L'Ecosystème du Big Data (20H, 3ECTS) SES721	Visualisation d'Information (30H, 3ECTS) INF229	Bases de Données Avancées (60H, 5ECTS) INF345	
P3	Systèmes Répartis (60H, 5ECTS) INF346		Machine-Learning (60H, 5ECTS) MDI 343	
P4	Machine-Learning Avancé (60H, 5ECTS) INFMDI341		Données du Web (60H, 5ECTS) INF344	

L'exemple du Master spécialisé Big Data à Télécom ParisTech (Stephan Cléménçon)

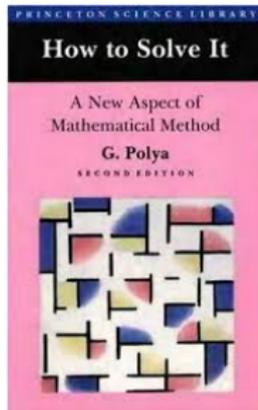
Modélisation et Big Data

Il est nécessaire (et utile) de **modéliser** le Big Data



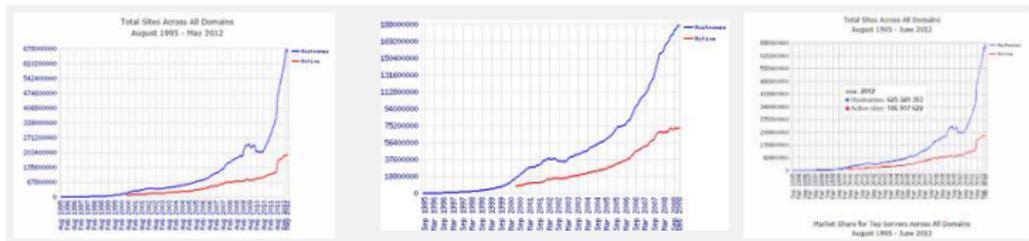
*If you cannot solve the proposed problem, try to solve first some related problem. **Could you imagine a more accessible related problem?***

George Pólya,
How to Solve It (1945)



Grand n , grand p

Plus de données

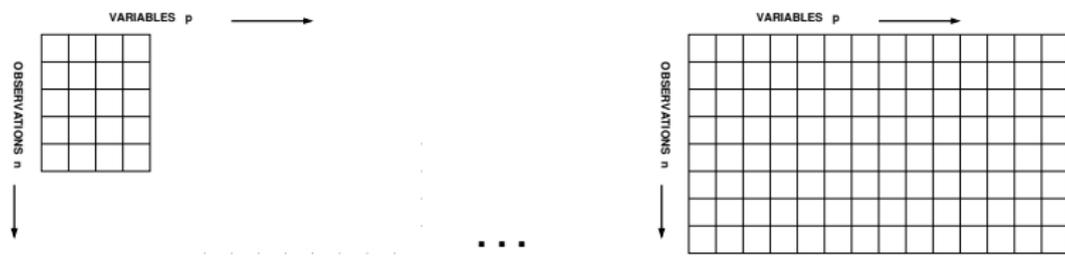


mais aussi plus de statistiques!



Grand n , grand p

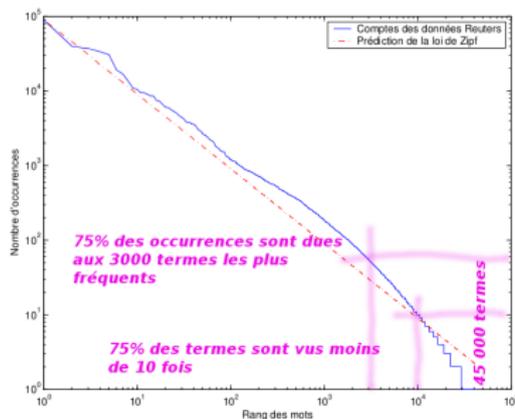
On imagine une **séquence d'expériences** dans lesquelles les nombres d'observations et de variables explicatives augmentent simultanément



Pour modéliser certaines situations (expérience “Microarray” en génomique), il est possible de considérer aussi des cas où $p > n$

“Lois” de Zipf

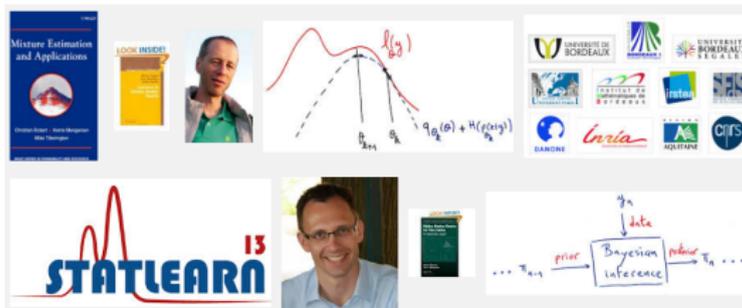
Rare n'est pas synonyme d'inintéressant



- Comment échantillonner les données ?
- Peut on construire (avec une complexité acceptable) des résumés exhaustifs des données ?

Contrôle du taux de faux positifs

Q: olivier cappe big data



Q: olivier cappe big data toulouse



Compromis performance / robustesse / calcul

Le volume de données nécessite de réduire la complexité des traitements

Idées : Algorithmes sous optimaux, utilisation de résumés des données. . .

- Régression sous-échantillonnée en observations et en prédicteurs
- Régressions marginales

Le calcul n'est plus une abstraction

Contraintes : structure des ressources de calcul, accès aux données

Algorithmes *online*, distribués

1 Apprentissage statistique

2 Big Data

3 Online EM

- L'algorithme EM
- L'algorithme online EM
- Applications
- Références

Modèles à données latentes

Un **modèle à données latentes** est un modèle statistique étendu $\{p_\theta(x, y)\}_{\theta \in \Theta}$ où **seul Y peut être observé**

- Les estimés θ_n du paramètre ne peuvent dépendre que des observations Y_1, \dots, Y_n
- Le modèle $\{f_\theta(y)\}_{\theta \in \Theta}$ est défini par marginalisation :
$$f_\theta(y) = \int p_\theta(x, y) dx$$

Les données $\{Y_t\}_{t \geq 1}$ sont supposées indépendantes et distribuées sous la loi marginale π (et $\pi \notin \{f_\theta\}_{\theta \in \Theta}$)

Exemple : Modèle de mélange

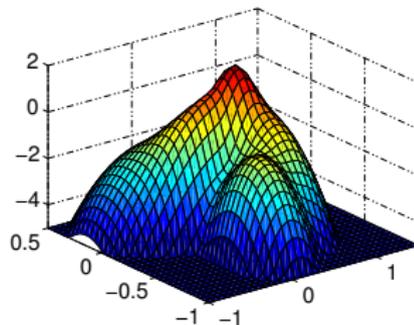
Densité de mélange

$$f(y) = \sum_{i=1}^m \alpha_i f_i(y)$$

Interprétation en terme de données latentes

$$P(X_t = i) = \alpha_i$$

$$Y_t | X_t = i \sim f_i(y)$$



Souvent utilisé pour réaliser une version probabiliste du *clustering* (traitement de la parole, du langage naturel. . .)

L'algorithme EM usuel vise à déterminer numériquement l'estimateur du **maximum de vraisemblance**

$$\theta_n = \arg \max_{\theta} \sum_{t=1}^n \log f_{\theta}(Y_t)$$

Expectation-Maximisation (Dempster, Laird & Rubin, 1977)

A l'étape k , étant donné l'estimé courant θ_n^k du paramètre

Etape E Calculer

$$q_{n, \theta_n^k}(\theta) = \sum_{t=1}^n \mathbb{E}_{\theta_n^k} [\log p_{\theta}(X_t, Y_t) | Y_t]$$

Etape M Mettre à jour l'estimé du paramètre

$$\theta_n^{k+1} = \arg \max_{\theta \in \Theta} q_{n, \theta_n^k}(\theta)$$

Principe

- 1 C'est un algorithme MM (maximisation d'un minorant) de par l'inégalité de Jensen

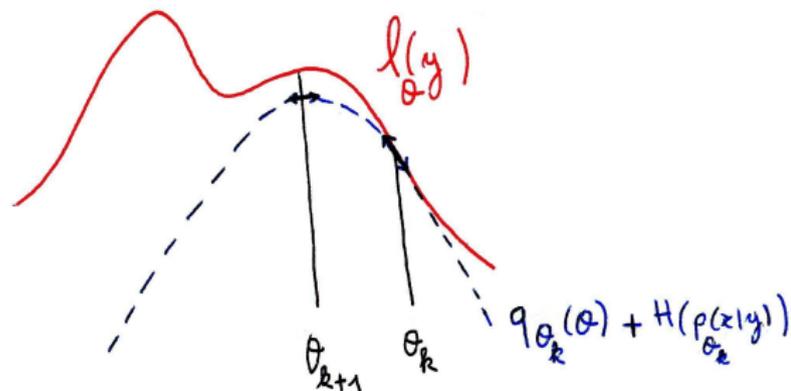
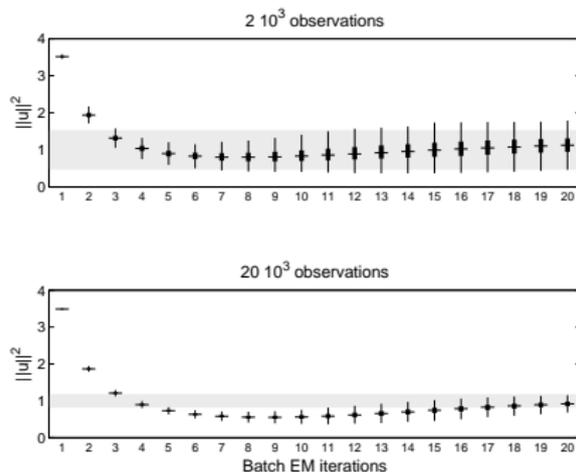


Figure: La quantité intermédiaire de l'EM minore localement la log-vraisemblance

- 2 Qui ne peut s'arrêter qu'en un point stationnaire de la log-vraisemblance (du fait de la relation dite de Fisher)

L'algorithme EM n'a pas un comportement satisfaisant lorsque n est grand

Le coût de calcul (et de stockage) de l'algorithme EM (pour un nombre fixé d'itération) est proportionnel à n pour un **résultat qui ne s'améliore pas quand $n \uparrow \infty$**



Algorithme online EM [C & Moulines, 2009–2011]

Pour remédier à ces défauts, on cherche une **variante de l'algorithme EM**

- qui à **chaque nouvelle observation Y_n** , met à jour l'estimé θ_n
- avec un coût de calcul (et de stockage) fixe par observation
- peut être interrompue à tout moment (et n'a pas besoin de connaître au préalable le nombre d'observations)
- est robuste vis à vis de contraintes d'accès au données
- se comporte comme le maximum de vraisemblance lorsque n est grand (mesuré via un TLC)

Ingrédient 1. Modèle de la famille exponentielle

Pour que la récursion de l'EM soit explicite on doit supposer

Un modèle de la famille exponentielle

$$p_{\theta}(x, y) = \exp(s(x, y)\psi(\theta) - A(\theta))$$

où $s(x, y)$ est une **statistique exhaustive** pour la loi p_{θ}
pour lequel le maximum de vraisemblance est explicite

$$S \mapsto \bar{\theta}(S) = \arg \max_{\theta} S\psi(\theta) - A(\theta)$$

Ingrédient 1. EM dans la famille exponentielle

La k -ème itération de l'algorithme EM

Etape E

$$S_n^{k+1} = \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\theta_n^k} [s(X_t, Y_t) | Y_t]$$

Etape M

$$\theta_n^{k+1} = \bar{\theta} \left(S_n^{k+1} \right)$$

peut être décrite uniquement via la **statistique exhaustive**

$$S_n^{k+1} = \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\bar{\theta}(S_n^k)} [s(X_t, Y_t) | Y_t]$$

Ingrédient 2. L'algorithme EM limite

En faisant tendre n vers l'infini on obtient

Version limite de l'algorithme EM

Mise à jour de la statistique exhaustive

$$S^k = E_{\pi} \left(E_{\bar{\theta}(S^{k-1})} [s(X_1, Y_1) | Y_1] \right)$$

Mise à jour du paramètre

$$\theta^k = \bar{\theta} \{ E_{\pi} (E_{\theta^{k-1}} [s(X_1, Y_1) | Y_1]) \}$$

Avec les mêmes arguments que pour l'algorithme EM, on montre que

- 1 la divergence de Kullback-Leibler $D(\pi | f_{\theta^k})$ décroît de façon monotone avec k
- 2 θ_k converge vers $\{\theta : \nabla_{\theta} D(\pi | f_{\theta}) = 0\}$

Ingrédient 3. L'approximation stochastique

- On recherche les solutions de

$$\mathbb{E}_\pi \left(\mathbb{E}_{\bar{\theta}(S)} [s(X_1, Y_1) | Y_1] \right) - S = 0$$

- En voyant $\mathbb{E}_{\bar{\theta}(S)} [s(X_n, Y_n) | Y_n]$ comme une observation bruitée de $\mathbb{E}_\pi \left(\mathbb{E}_{\bar{\theta}(S)} [s(X_1, Y_1) | Y_1] \right)$, on reconnaît un cas d'utilisation de **l'approximation stochastique** (ou algorithme de **Robbins-Monro**) :

$$S_n = S_{n-1} + \gamma_n \left(\mathbb{E}_{\bar{\theta}(S_{n-1})} [s(X_n, Y_n) | Y_n] - S_{n-1} \right)$$

où (γ_n) est une séquence de pas décroissants

L'algorithme

Online EM

Etape E via l'approximation stochastique

$$S_n = (1 - \gamma_n)S_{n-1} + \gamma_n \mathbf{E}_{\theta_{n-1}} [s(X_n, Y_n) | Y_n]$$

Etape M

$$\theta_n = \bar{\theta}(S_n)$$

Natural language processing

Online EM for Unsupervised Models, Percy Liang & Dan Klein, *NAACL Conference*, 2009

Context Text processing using word representation ($\approx 10k$ words) from large document corpora ($\approx 100k$ documents)

Tasks Tagging, classification, alignment (variant of mixture or HMM with multinomial observations)

In this application, the use of **mini-batch** blocking was found useful:

- Apply the proposed algorithm considering $Y_{mk+1}, Y_{mk+2} \dots Y_{m(k+1)}$ as one observation (with m of the order of a few k documents)

Online EM for Unsupervised Models

Percy Liang Dan Klein
Computer Science Division, EECS Department
University of California at Berkeley
Berkeley, CA 94720
{pliang,klein}@eecs.berkeley.edu

Abstract

The EM algorithm plays an important role in unsupervised induction, but its convergence often takes more iterations than is desirable. In this paper, we show that iterative variants of EM can converge significantly faster and (2) can even find local optima that are more robust than those found by plain EM. We support these findings on four unsupervised tasks: part-of-speech tagging, document classification, word segmentation, and word alignment.

1. Introduction

In unsupervised NLP tasks such as tagging, part-of-speech assignment, coreference, and named entity recognition, one wishes to induce latent linguistic structure from raw text. Probabilistic modeling has emerged as a dominant paradigm for these problems, and the EM algorithm has been a driving force for progress on such a simple and intuitive system. In many of these tasks, EM can converge significantly faster than an unoptimized parallel EM implementation, and it requires fewer iterations to converge than the well-known Baum-

(1966) of examples. Online algorithms have the potential to speed up learning by making updates more frequently. However, these updates can be seen as noisy approximations to the full batch update, and this noise can in fact impede learning.

This tradeoff between speed and stability is familiar to online algorithms for convex supervised learning problems—e.g., Perceptron, MIRA, stochastic gradient, etc. Unsupervised learning adds two additional issues: (1) Since the EM objective is non-convex, we often get convergence to different local optima of varying quality; and (2) we evaluate the accuracy metrics which are not directly correlated with the EM likelihood objective (Liang and Klein, 2009). We will see that these issues can limit the learning results.

In Section 4, we present a thorough investigation of online EM, mostly focusing on supervised tasks, and demonstrate incremental EM. For supervised tasks, we find that choosing a good regularizer can be as important as choosing a good learning rate. We also quantify relative improvements. We conclude with a discussion of the implications of our findings for supervised EM and for unsupervised learning.

Large-scale probabilistic sequence alignment

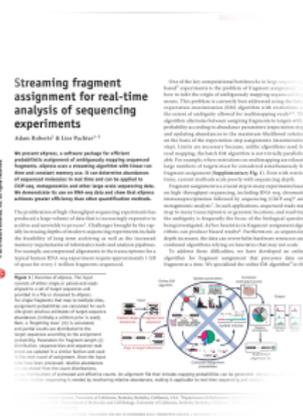
Streaming fragment assignment for real-time analysis of sequencing experiments,
Adam Roberts & Lior Pachter, *Nature*, 2013

Context High-throughput sequencing experiments in biology

Task Learn parameters of probabilistic models for sequence alignments (finite-valued sequences with insertion/deletion/substitution probabilities) from **10M sequence fragments** (reads) with **10k target sequences** (transcripts)

Note that the actual goal is the **retrospective** probabilistic sequence assignment (E-step)

Use a single pass of online EM through all the data, followed by a few more batch EM steps



- Roberts, A. & Pachter, L. (2013). Streaming fragment assignment for real-time analysis of sequencing experiments. *Nature*.
- Cappé, O. (2011). Online EM algorithm for hidden Markov models. *J. Comput. Graph. Statist.*
- Cappé, O. (2011). Online Expectation-Maximisation. In *Mixtures*, Wiley.
- Rohde, D. & Cappé, O. (2011). Online maximum-likelihood estimation for latent factor models. *IEEE Statistical Signal Processing Workshop*.
- Liang, P. & Klein, D. (2009). Online EM for Unsupervised Models. *NAACL Conference*.
- Cappé, O. & Moulines, E. (2009). On-line expectation-maximization algorithm for latent data models. *J. Roy. Statist. Soc. B*.
- Sato, M. & Ishii, S. (2000). On-line EM algorithm for the normalized Gaussian network. *Neural Computation*.
- Sato, M. (2000). Convergence of on-line EM algorithm. *International Conference on Neural Information Processing*.
- Neal, R. M. & Hinton, G. E. (1999). A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, MIT Press.
- Titterton, D. M. (1984). Recursive parameter estimation using incomplete data. *J. Roy. Statist. Soc. B*.

1 Apprentissage statistique

2 Big Data

- Quelques idées générales
- Quelques idées de modélisation

3 Online EM

- L'algorithme EM
- L'algorithme online EM
- Applications
- Références

Merci de votre attention