



# La science des données et l'avenir de la surveillance financière

Communication du surintendant Jeremy Rudin à l'occasion du 11<sup>e</sup> colloque sur les services bancaires et les finances en Asie

San Francisco, Californie  
Le 25 juin 2018

Le texte prononcé fait foi

Pour obtenir de plus amples renseignements :

Kaitlin Sabourin  
Communications et consultations  
[kaitlin.sabourin@osfi-bsif.gc.ca](mailto:kaitlin.sabourin@osfi-bsif.gc.ca)  
[www.osfi-bsif.gc.ca](http://www.osfi-bsif.gc.ca)



BSIF  
OSFI

Canada 

Communication du surintendant Jeremy Rudin<sup>1</sup>  
Bureau du surintendant des institutions financières (BSIF)  
à l'occasion du  
11<sup>e</sup> colloque sur les services bancaires et les finances en Asie  
San Francisco, Californie  
Le 25 juin 2018

---

## **Introduction**

Bonjour, et merci de votre aimable présentation. C'est avec plaisir que je participe au 11<sup>e</sup> colloque sur les services bancaires et les finances en Asie. Je tiens à remercier la Federal Reserve Bank de San Francisco et la Monetary Authority of Singapore de m'avoir invité.

Pour ce volet du colloque, les organisateurs nous ont demandé de déterminer si les organismes de surveillance font une promotion adéquate de leurs outils et de leurs capacités permettant de suivre l'évolution du système financier. Et ce, pour une bonne raison.

Nous avons tous vu comment le changement technologique transforme le secteur des services financiers. À mon avis, le plus important de ces changements est l'explosion remarquable de notre capacité à recueillir et à analyser des données. Les ensembles de données disponibles aujourd'hui sont si vastes que leur analyse exige de nouvelles approches qui exploitent la croissance continue de la puissance computationnelle grâce à l'utilisation de l'apprentissage machine et de l'intelligence artificielle de façon plus générale.

Lorsque les gens qualifient ce nouveau champ de « science des données », ce n'est pas du simple marketing : les techniques d'analyse appliquées à des ensembles massifs de données témoignent d'une rupture avec les statistiques et l'économétrie traditionnelles.<sup>2</sup>

Vous savez déjà que le secteur des services financiers a embrassé la science des données et qu'il y plonge plus avant. Les institutions financières utilisent ces techniques pour la souscription d'assurance, l'octroi du crédit, la lutte contre le recyclage des produits de la criminalité et dans bien d'autres domaines.

C'est un sujet très important pour les surveillants, mais ce n'est pas l'objet de mon propos aujourd'hui. J'aborderai non pas comment la science des données changera ce que font les *institutions financières*, mais plutôt comment elle changera ce que font leurs *autorités de contrôle*.<sup>3</sup>

---

<sup>1</sup> Je remercie Solon Angel, Jamil Abou Saleh, Andrew Kriegler et bon nombre de mes collègues du BSIF de leurs remarques et idées. Ils ne pourraient toutefois pas être tenus responsables d'erreurs potentielles dans le présent texte.

<sup>2</sup> Pour obtenir une comparaison entre l'apprentissage machine et l'économétrie appliquée, voir l'article de Sendhil Mullainathan et de Jann Spiess intitulé « Machine Learning: An Applied Econometric Approach », dans le *Journal of Economic Perspectives*, publié au printemps 2017 (en anglais seulement).

<sup>3</sup> Pour explorer les idées de l'un des pionniers de la science des données (la Monetary Authority of Singapore), consulter le discours principal de David Hardoon intitulé « Data Science and Machine Learning in Practice », qu'il a prononcé à l'occasion de la 7<sup>e</sup> conférence annuelle de l'institut Sim Kee Boon sur les progrès réalisés dans le domaine de la science des données et son incidence commerciale, le 26 mai 2017 (en anglais seulement) : <http://www.mas.gov.sg/News-and-Publications/Speeches-and-Monetary-Policy-Statements/Speeches/2017/Data-Science-and-Machine-Learning-in-Practice.aspx>

## ***La science des données et son incidence sur les autorités de contrôle***

À mon avis, la science des données aura une incidence majeure sur la façon dont nous surveillons les institutions financières.

Mon message à mes collègues surveillants est le suivant : la science des données est trop importante pour être confiée uniquement aux scientifiques des données. En tant que surveillants, nous devons comprendre ce que la science des données peut et ne peut pas faire pour nous aider. Si nous ne comprenons pas comment utiliser judicieusement la science des données, nous ne tirerons pas tous les avantages de ces puissantes techniques. De plus, nous pourrions détourner notre attention de certains risques importants, une erreur dont nous serions les seuls responsables.

Afin de vous convaincre que la science des données est trop importante pour être confiée uniquement aux scientifiques des données, je vais examiner comment elle peut s'appliquer au travail de deux surveillants différents. La première surveillante est chargée de la surveillance de la conduite des institutions et, en particulier, elle détecte et poursuit en justice les cas de délit d'initié illicites. Je parlerai de « elle » pour désigner cette personne. Le second surveillant est chargé de la surveillance prudentielle de grandes banques complexes. Ce surveillant, ce sera « moi » ou « je ».

### ***La science des données et la surveillance de la conduite des institutions : déceler le délit d'initié***

Parlons d'abord de notre collègue chargée de surveiller la conduite.

En tant que surveillante de la conduite, elle a plusieurs objectifs. Elle doit notamment décourager, voire prévenir, le délit d'initié illicite. Elle doit également déceler les cas de délit d'initié illicite antérieurs pour pouvoir les poursuivre en justice.

Elle dispose d'une grande quantité de données, y compris des données de très haute fréquence concernant les opérations sur le marché portant sur un large éventail d'instruments financiers et remontant à plusieurs années. Elle a des raisons de croire que son ensemble de données renferme des signes qui donnent à penser qu'il pourrait y avoir des cas de délit d'initié illicite. Si seulement elle savait où chercher...

Elle se tourne donc vers la science des données pour obtenir de l'aide.<sup>4</sup> Que peut-elle faire?

Supposons qu'elle ait déjà cerné et poursuivi en justice un certain nombre de cas de délit d'initié durant sa carrière. Elle utilise cette information pour étiqueter les opérations d'initié connues dans son ensemble de données. Elle programme ensuite un algorithme d'apprentissage machine pour examiner l'ensemble des données afin de trouver les opérations dont le profil correspond le mieux à ces opérations d'initié étiquetées. L'algorithme peut ensuite rechercher ce profil dans le reste des données.

---

<sup>4</sup> Il existe déjà plusieurs applications dans ce domaine de la science des données. Voir, par exemple, l'article intitulé « SEC's advanced analytics helps detect even the smallest illicit market activity », de Reuters, publié le 30 juin 2017, à l'adresse que voici : <https://www.reuters.com/article/bc-finreg-data-analytics/secs-advanced-data-analytics-helps-detect-even-the-smallest-illicit-market-activity-idUSKBN19L28C> et l'article de James Langton intitulé « IROC to strengthen market surveillance with new technology », dans l'*Investment Executive*, publié le 13 juillet 2017, à l'adresse que voici : <https://www.investmentexecutive.com/news/from-the-regulators/iroc-to-strengthen-market-surveillance-with-new-technology/> (les deux articles sont en anglais seulement).

Ce qu'elle obtient, c'est une liste des opérations les plus intéressantes qu'elle pourra étudier pour vérifier s'il s'agit bel et bien d'opérations d'initié illicites comme le montre l'algorithme. Que va-t-elle trouver lorsqu'elle enquêtera sur ces cas?

Si elle n'a qu'un petit nombre d'opérations d'initié avérées, elle constatera probablement que les prédictions de l'algorithme ne sont pas entièrement fiables. Paradoxalement, le problème est qu'il y a trop de données par rapport au nombre d'opérations d'initié prouvées. L'algorithme peut examiner une liste presque interminable de caractéristiques : la rentabilité des opérations; le nombre de personnes effectuant des opérations semblables; le nombre de fois que ces personnes ont effectué des opérations; la rapidité avec laquelle ces opérations ont été effectuées; les opérations qu'effectuaient ces mêmes personnes sur les dérivés connexes, et ainsi de suite. Tout ensemble d'opérations est susceptible d'avoir un certain nombre de caractéristiques en commun parmi cette liste extrêmement longue, ne serait-ce que par hasard. Un algorithme d'apprentissage machine trouvera bon nombre, sinon la totalité, de ces points communs. Certains de ces indicateurs dévoilent bel et bien des cas de délit d'initié, tandis que d'autres ne sont que des coïncidences.

Ce problème, appelé « surajustement », bien connu dans le domaine de l'économétrie, est courant dans ce genre d'exercice d'apprentissage machine, et il est plus grave lorsque le nombre d'exemples de l'événement que nous tentons de détecter est relativement restreint.<sup>5</sup>

Le surajustement est un problème, mais ce n'est pas un problème critique dans son cas. Elle n'utilise pas la science des données pour conclure qu'une opération donnée constituait une opération d'initié. Elle utilise la science des données afin de repérer les opérations suspectes pour un examen plus approfondi. Une fois qu'elle dispose d'une liste de pistes prometteuses, elle peut recueillir d'autres éléments de preuve et décider quelles opérations suspectes recensées par l'algorithme peuvent être poursuivies, et lesquelles ne peuvent l'être. En effet, elle a l'obligation de recueillir d'autres preuves; elle ne peut sanctionner quelqu'un pour délit d'initié sur la seule foi de la science des données.

### ***La science des données et la surveillance prudentielle : contributions possibles***

À mon tour, maintenant.

En tant que surveillant prudentiel, mon travail consiste à ramener à un niveau acceptable la probabilité de faillite des grandes banques complexes du pays. Comment faire?

La surveillance prudentielle comporte de nombreux aspects. En règle générale, les surveillants prudentiels d'institutions bancaires passent en revue ou examinent leurs pratiques de gestion du risque, cernent les lacunes et voient à ce que ces dernières soient corrigées. Nous établissons également des exigences de fonds propres et de liquidité qui vont au-delà des exigences réglementaires minimales.

---

<sup>5</sup> Il y a plusieurs façons de corriger le surajustement, et diverses façons de l'éviter. Vous trouverez une explication judicieuse et intuitive du surajustement de grandes quantités de données dans l'article de Vincent Spruyt intitulé « The Curse of Dimensionality in Classification », à l'adresse que voici : <http://www.visiondummy.com/2014/04/curse-dimensionality-affect-classification/> (en anglais seulement).

À l'appui de ces activités, nous cherchons des indices montrant que les banques prennent trop de risques ou que certaines de leurs activités sont trop risquées. Nous pouvons ensuite intervenir pour ramener le risque à un niveau acceptable.

Mes collègues et moi faisons déjà un travail considérable d'analyse de données. À mesure que nous jumelons une puissance informatique accrue avec des données plus granulaires, nous pouvons accroître la vitesse, l'exactitude et le niveau de détail de notre analyse des données existantes. C'est un progrès; peut-être ce sera même un progrès considérable.

Mais ce n'est pas le but ultime que je cherche à atteindre. Je dois croire qu'il existe dans les données des indicateurs plus puissants, plus informatifs et plus utiles qui attendent d'être découverts. Je vais essayer d'utiliser la science des données pour les trouver.

L'ensemble de données dont j'ai besoin pour travailler provient d'un seul pays; il se rapporte donc à un nombre limité de grandes banques complexes. Il contient toutefois beaucoup de détails couvrant plusieurs années sur chacune de ces banques.

Si je suis dans la même situation que beaucoup de surveillants prudentiels, mon ensemble de données ne fera état que de quelques faillites de grandes banques complexes. Pourquoi? Parce qu'il n'y a pas tant de grandes banques complexes que cela et qu'elles ne font pas souvent faillite. À l'instar de ma collègue chargée de surveiller la conduite, mes résultats seront aussi susceptibles d'être surajustés.

Je ne peux cependant pas adopter la même approche qu'elle pour atténuer le problème de surajustement. Je ne cherche pas à déceler quelque chose qui s'est déjà produit; je cherche à prédire un avenir qui demeure non réalisé. Je n'ai donc aucune autre preuve qui puisse confirmer ou réfuter de façon concluante les cas relevés par l'algorithme.

Dans mon cas, c'est trop de demander de chercher des preuves concluantes. Je pourrais plutôt décider d'accorder une attention particulière aux banques choisies par l'algorithme. C'est faisable, mais je n'ai que quelques grandes banques complexes que j'ai déjà étroitement à l'œil, alors ce n'est pas très une solution très intéressante.

Il serait préférable d'utiliser les prédictions de l'algorithme pour déterminer où je devrais examiner de plus près. Ici, j'aurai deux problèmes auxquels ma collègue ne sera pas confrontée.

Premièrement, il n'existe probablement pas d'interprétation claire ou intuitive des prédictions faites par l'algorithme. L'algorithme choisira les banques qu'il considère comme les plus à risque à l'aide d'une combinaison d'indicateurs pouvant couvrir plusieurs activités bancaires et comportant divers délais de faillite. Cela ne me dira ni où chercher, ni quoi chercher. Je devrais m'y attendre : les algorithmes d'apprentissage machine ne sont pas configurés pour expliquer le bien-fondé de leurs prédictions.<sup>6</sup>

---

<sup>6</sup> Plusieurs professionnels du domaine font déjà des efforts pour corriger ce défaut bien connu. Consulter, par exemple, l'article de Lars Hulstaert intitulé « Interpreting machine learning models », à l'adresse que voici : <https://towardsdatascience.com/interpretability-in-machine-learning-70c30694a05f> (en anglais seulement).

Deuxièmement, comme je suis un surveillant d'un organisme de réglementation prudentielle, je ne suis pas dans le domaine de la prédiction, mais bien dans celui de la prévention. L'algorithme a peut-être effectivement trouvé les meilleurs *prédicteurs* d'une éventuelle faillite. Ce que j'ai besoin de connaître, ce sont les *causes* les plus probables d'une éventuelle faillite. Ce sont peut-être deux choses différentes. Il ne sert à rien de prendre des mesures de surveillance afin de modifier quelque chose dont on sait qu'il permet de prédire les faillites des banques si ce « quelque chose » ne contribue pas également à la faillite. En clair, faire taire le coq n'empêchera pas le soleil de se lever.

Ma collègue n'a pas ce problème. Une fois qu'elle a décelé une opération suspecte, elle est en mesure de confirmer si elle peut la poursuivre ou non. Tout ce qui l'oriente vers des opérations prometteuses à étudier cachées parmi les millions d'opérations dans son ensemble de données lui est utile. Peu lui importe que la machine ait trouvé le feu, ou seulement la fumée.

### ***Le problème de la boucle de rétroaction : l'incidence de nos pratiques sur les données***

Nous avons toutefois un autre problème en commun. Nous devons tous deux composer avec le fait que notre façon de surveiller aura une incidence sur les données que nous utilisons.

Prenons cette fois un exemple prudentiel pour illustrer cette question. Supposons que, selon nos données, une banque exerçant une certaine activité soit susceptible de faire éventuellement faillite. Puisque nous œuvrons dans le domaine de la prévention, nous pourrions décider d'interdire cette pratique à moins que la banque ne prenne des mesures approuvées pour compenser le risque. Supposons en outre que cela contribue à prévenir les faillites attribuables à cette pratique risquée. Cette nouvelle pratique de surveillance sera intégrée aux données.

Au fur et à mesure que nous recueillons des données sous notre nouveau régime, nous observons que la pratique risquée n'annonce plus la faillite éventuelle de la banque qui l'exerce. Notre pratique de surveillance fera en sorte que les données occulteront les risques sous-jacents.

Notre pratique de surveillance masque peut-être les risques réels d'une autre façon, peut-être plus pernicieuse celle-là. Supposons que nous interdisions par erreur une pratique qui n'est pas vraiment risquée. Les données ne révéleront jamais cette erreur, car nous ne permettrons pas que la pratique apparaisse dans les données, prouvant ainsi son innocence.<sup>7</sup>

Ma collègue est également confrontée à ce problème. Ceux qui sont tentés d'effectuer des opérations d'initié pourraient décoder les indicateurs qu'elle utilise pour cerner ces dernières. Ils pourraient ensuite apprendre à effectuer des opérations d'initié que ces indicateurs ne décèleront pas. Elle devra donc éviter de tomber dans la complaisance si le nombre de cas signalés par son algorithme diminue. Cette baisse pourrait indiquer qu'elle cerne mieux les opérations d'initié et que cela a un effet dissuasif. Ou alors, cela pourrait aussi vouloir dire que les criminels deviennent plus habiles pour éviter la détection.

---

<sup>7</sup> Le risque qu'un algorithme d'apprentissage machine ne puisse recenser et corriger ses erreurs est l'un des thèmes du livre de Cathy O'Neil intitulé *Weapons of Math Destruction*, publié en 2016 (en anglais seulement).

## ***Exploiter le pouvoir de l'apprentissage machine non surveillé***

Jusqu'à maintenant, j'ai utilisé la première personne en parlant du surveillant prudentiel dans mes exemples, mais ce n'est pas tout à fait exact. Je ne peux pas essayer l'approche décrite ci-dessus pour le moment, car mon ensemble de données ne fait état d'aucune faillite de grandes banques complexes.

La science des données peut-elle faire quelque chose pour moi? Certainement. Lorsque nous ne sommes pas en mesure d'étiqueter les cas d'intérêt dans notre ensemble de données, nous pouvons adopter une approche au nom quelque peu menaçant : l'apprentissage machine non surveillé. En mode d'apprentissage non surveillé, l'algorithme recherche des anomalies dans les données. L'algorithme ne sait pas que ces cas posent problème; il sait seulement qu'ils se démarquent d'une certaine façon. Cela dit, ces cas peuvent être de bons candidats pour une enquête plus approfondie.

Par exemple, si ma collègue n'avait pas relevé de cas de délit d'initié illicites confirmés, elle aurait pu recourir à l'apprentissage non surveillé afin de repérer les anomalies en vue d'une enquête plus approfondie.

Que vais-je faire avec mes cas d'anomalie? Je peux certainement les examiner plus avant pour voir s'ils laissent entrevoir des activités à risque excessif. Contrairement à ma collègue, je n'ai aucun moyen de prouver de façon concluante que les anomalies trouvées sont excessivement risquées; après tout, mon ensemble de données ne fait état d'aucune faillite. Je peux plutôt me servir de mon jugement et de mon expérience en tant que surveillant prudentiel pour tenter de déterminer quelles anomalies, le cas échéant, signalent une prise de risque excessive. Ce faisant, je dois être conscient de deux écueils potentiels.

Tout d'abord, après avoir investi temps et efforts pour trouver des anomalies, je serai enclin à trouver des raisons de croire que les anomalies constatées témoignent effectivement de la prise de risques excessifs. Le biais de confirmation, c'est-à-dire la tendance à penser que vous avez trouvé ce que vous vous attendez à trouver, est difficile à surmonter.

Si nous succombons à ce biais maintenant, nous ne serons pas en mesure de le corriger plus tard. En décourageant les banques de faire tout ce qui crée l'anomalie, nous veillerons à ce que l'anomalie ne réapparaisse pas dans les données. Comme on l'a vu, ces pratiques n'auront donc pas l'occasion de prouver leur innocence.

Deuxièmement, à mesure que j'éliminerai progressivement les anomalies dans les données, je rendrai mes grandes banques complexes plus semblables les unes aux autres. C'est peut-être une bonne chose si je supprime des activités excessivement risquées. Voilà le genre de similarité que nous cherchons. Mais ce pourrait être mauvais si cela réduit la diversification à l'échelle des banques, rendant toutes les banques vulnérables aux mêmes chocs.

Il y a peut-être une façon moins ambitieuse, mais potentiellement plus fiable, d'appliquer l'apprentissage non surveillé à un ensemble de données n'indiquant aucune faillite bancaire. Je pourrais faire porter l'apprentissage non surveillé sur les relevés réglementaires produits par les banques, puis vérifier si les anomalies signalées par l'algorithme témoignent de relevés erronés ou

frauduleux.<sup>8</sup> C'est quelque chose que je peux vérifier de façon indépendante, mais qui requiert un certain effort. Et c'est certainement utile. La production de relevés réglementaires invariablement incorrects ou, pire encore, frauduleux devrait certainement attirer l'attention sur la banque en question.<sup>9</sup>

## **Conclusion**

Au cours des dernières minutes, nous avons tout juste effleuré le sujet; il y a beaucoup d'autres cas que nous pourrions considérer, et d'autres techniques qui sont maintenant mises à notre disposition grâce à la science des données. Je veux m'arrêter avant d'avoir atteint les limites de votre patience, et je dois m'arrêter avant d'avoir atteint celles de mes connaissances.

Je terminerai donc en vous soumettant une réflexion. Les professions comme la surveillance prudentielle, celles qui reposent fortement sur le jugement et l'expérience, ont tendance à résister fortement à l'idée que les machines et les algorithmes peuvent améliorer ce que l'expertise humaine peut accomplir seule.<sup>10</sup> C'est une erreur. Et nous ne pouvons pas nous permettre de faire cette erreur.

La science des données sera importante pour la surveillance financière. Très importante. Elle sera si importante qu'elle ne doit pas être l'apanage des scientifiques des données.

Je vous remercie.

---

<sup>8</sup> Pour explorer l'application de cette approche dans le secteur de l'assurance, consulter le document de Sutapat Thiprungsri et de Miklos A. Vasarhelyi intitulé « Cluster Analysis for Anomaly Detection in Accounting Data: An Audit Approach », dans l'*International Journal of Digital Accounting Research*, volume 11, 2011 (en anglais seulement). L'idée d'utiliser les anomalies mathématiques pour déceler la fraude précède le développement de l'apprentissage machine. Consulter, par exemple, le document de Cindy Durtschi, de William Hillison et de Carl Pacini intitulé « The Effective Use of Benford's Law to Assist in Detecting Fraud in Accounting Data », publié dans le *Journal of Forensic Accounting*, volume 5, 2004, qui raconte l'histoire de cette méthode.

<sup>9</sup> La Banque d'Angleterre met l'apprentissage non surveillé à l'essai dans quelques contextes. Consulter la page Web à l'adresse que voici : <https://www.bankofengland.co.uk/research/fintech/proof-of-concept>.

<sup>10</sup> Consulter, par exemple, le livre de Daniel Kahneman intitulé *Système 1 / Système 2 : Les deux vitesses de la pensée*, publié en 2011, surtout les chapitres 21 et 22.