

Christophe Brasseur

Enjeux et usages du big data

2^e édition



Direction éditoriale : Emmanuel Leclerc
Édition : Céline Poiteaux
Fabrication : Estelle Perez
Composition : Nord Compo, Villeneuve-d'Ascq
Image de couverture : iconimage – Fotolia.com

© 2016, Lavoisier, Paris
ISBN : 978-2-7462-4758-1

Collection
Information numérique
Traitement, interprétation, communication
dirigée par Olivier Rioul

Professeur, Télécom ParisTech,
Université Paris-Saclay, Paris

Comité éditorial :

Gérard Blanchet, Professeur émérite, Télécom ParisTech, Université Paris-Saclay, Paris.

Isabelle Bloch, Professeur, Télécom ParisTech, Université Paris-Saclay, Paris.

Valérie Fernandez, Professeur, Télécom ParisTech, Université Paris-Saclay, Paris.

Benoît Geller, Professeur, ENSTA ParisTech, Université Paris-Saclay, Paris.

Pour ma femme Christine et mes enfants Arthur et Capucine.

Table des matières

Introduction : les mégadonnées au cœur de la révolution digitale	1
--	---

Partie 1 — Origines et enjeux du big data

Chapitre 1

Qu'est-ce que le big data ?

1. Essor des réseaux sociaux, des objets connectés et du <i>cloud</i>	7
<hr/>	
2. Définition et caractéristiques du big data	10
<hr/>	
2.1. Définition du big data	10
2.2. Caractéristiques du big data	12
2.3. Différents types de données	14
2.4. Big data et décisionnel	17
<hr/>	
3. Qui est concerné par le big data ?	18
<hr/>	
3.1. Secteurs économiques	18
3.2. Métiers concernés	20

Chapitre 2

Enjeux du big data et du digital

1. Quand les données créent de la valeur	25
1.1. Un levier de compétitivité incontournable	25
1.2. Une source de progrès scientifiques et humains	27
1.3. Gains engendrés par le big data	29
2. Big data et transformation digitale	30
2.1. Pourquoi la transformation digitale est-elle inéluctable ?	30
2.2. Le big data au cœur de l'expérience client.....	31
2.3. Comment aborder la transformation digitale ?.....	33
2.4. Analyse du positionnement actuel.....	33
2.5. Définition du modèle cible.....	34
2.6. Mise en œuvre de la transformation digitale.....	34
3. Un défi technologique majeur	35
4. Big data et qualité des données	37
4.1. Quelques exemples courants de défauts de qualité des données	37
4.2. La gouvernance des données au cœur de la stratégie big data.....	39
5. Big data et protection des données personnelles	40
5.1. Législation en matière de protection des données personnelles	41
5.2. Les bonnes pratiques à adopter.....	43
6. Ouverture des données publiques	44
7. Prévisions de croissance du marché big data	46

Chapitre 3

Exemples d'usage du big data

1. Quand le big data apporte des réponses au marketing	49
1.1. Analyse des conversations sur les réseaux sociaux.....	49
1.2. Personnalisation de l'offre et communication ciblée.....	52
1.3. Innovation participative	53
2. Vers une gestion intelligente et responsable de l'énergie	54
2.1. Enjeux et mutation des réseaux électriques	55
2.2. Caractéristiques des réseaux électriques intelligents.....	56
2.3. Un défi pour la gestion des données volumineuses	57
2.4. Projets de <i>smart grids</i> en cours	58
3. Le métier de l'assurance réinventé	59
4. Les mégadonnées pour améliorer la santé	60
4.1. Recherche médicale	61
4.2. Prévention des épidémies.....	62
4.3. Vers une médecine mobile et personnalisée	63
5. Un levier de transformation des services publics	64
5.1. Prérequis à la transformation digitale de l'État.....	65
5.2. Lutte contre la fraude.....	65
5.3. Lutte contre le terrorisme	66
5.4. Personnalisation des services publics.....	67
6. Du journalisme traditionnel <i>au data journalism</i> (journalisme de données)	68
6.1. Collecte des données.....	69
6.2. Transformation des données en informations.....	69

6.3. Visualisation des données	70
6.4. Quelques spécialistes du journalisme de données	71
7. Et tant d'autres domaines d'application...	72

Partie 2 — Technologies et méthodes du big data

Chapitre 4 Technologies du big data

1. Problématique technique	77
1.1. Limites des bases de données classiques.....	77
1.2. De nouvelles exigences techniques	80
2. Solutions techniques du big data	81
2.1. Des solutions majoritairement <i>open source</i>	81
2.2. Hadoop et son écosystème	83
2.3. Apache Spark et sa solution en mémoire	86
2.4. De nouveaux types de bases de données.....	87
2.5. Outils d'analyse des mégadonnées.....	90
3. La rencontre du big data et du <i>cloud computing</i>	92
3.1. Qu'est-ce que le <i>cloud computing</i> ?	92
3.2. Avantages et inconvénients du <i>cloud</i>	93
3.3. Le <i>cloud computing</i> , facilitateur du big data	94
4. Big data et web sémantique	95
4.1. Web sémantique	95
4.2. Intérêt du web sémantique pour le big data.....	96

Chapitre 5

Méthodes et techniques d'analyse du big data

1. Analyse des données massives	99
1.1. Qu'est-ce que l'analyse des données ?	99
1.2. Méthodes d'analyse de référence	102
2. Principales techniques d'analyse des données massives	105
2.1. <i>Data mining</i>	105
2.2. Apprentissage automatique (<i>machine learning</i>) et informatique cognitive (<i>cognitive computing</i>)	107
2.3. Analyse des réseaux sociaux	109
2.4. Test A/B	110
2.5. <i>Crowdsourcing</i>	111
2.6. Géomarketing	112
2.7. Analyse des séries temporelles.....	113
3. Techniques de visualisation	114
3.1. Objectif de la visualisation des données	114
3.2. Principes de la visualisation	115
3.3. Quelques exemples de visualisation de données.....	117

Partie 3 — Comment tirer parti du big data ?

Chapitre 6

Compétences et ressources humaines

1. Le virage du digital et de l'analytique	121
1.1. Manager avec l'analytique	121

1.2. De nouvelles opportunités pour les DSI	122
1.3. Compétences et organisation	123
2. Montée en puissance des <i>data scientists</i>	126
2.1. Les compétences fondamentales du <i>data scientist</i>	126
2.2. Qualités indispensables.....	127
2.3. Comment devient-on <i>data scientist</i> ?.....	128
2.4. Compétences IT requises	130

Chapitre 7

Gérer un projet big data

1. Caractéristiques d'un projet big data	133
1.1. Qu'est-ce qu'un projet big data ?.....	133
1.2. Une organisation projet pluridisciplinaire	134
1.3. Une approche progressive et itérative	134
2. Quel retour sur investissement pour le big data ?	136
2.1. Définition du retour sur investissement.....	136
2.2. Limites du ROI	137
2.3. Le ROI des projets big data.....	137
3. Étapes d'un projet big data	138
3.1. Phase préparatoire.....	138
3.2. Analyse et réalisation.....	139
3.3. Mise en œuvre et exploitation.....	140
4. Risques à prendre en compte	140
4.1. Degré de maturité de l'entreprise	140

4.2. Des technologies relativement jeunes.....	141
4.3. Comment limiter les risques ?.....	141
Conclusion : en route vers les <i>smart data</i>	143
Bibliographie.....	147
Index.....	149

Introduction

Les mégadonnées au cœur de la révolution digitale

L'utilité des systèmes d'information dans les entreprises n'est plus à démontrer. Pratiquement aucune entreprise ne peut aujourd'hui s'en passer. L'activité quotidienne des banques, des assurances, des sociétés industrielles, des entreprises de distribution ou des services publics dépend du bon fonctionnement du système d'information. Celui-ci facilite et optimise les processus métier de l'entreprise, et pratiquement toutes les fonctions (comptabilité, contrôle de gestion, marketing, vente, production, achats, ressources humaines, qualité, maintenance, recherche) sont concernées. On oublie parfois que la matière première du SI est constituée de données devenues indispensables à la vie et au développement des organisations. Il est clair que les données sont un facteur clé de succès et les entreprises qui en font bon usage gagnent en compétitivité. Des données fiables permettent notamment de prendre de bonnes décisions.

Depuis quelques années, on assiste à une prise de conscience de plus en plus importante de l'avantage compétitif que l'on peut obtenir avec des données de qualité totalement maîtrisées. Il reste encore du chemin à parcourir, mais ce mouvement vers la reconnaissance de la valeur des *data* est une bonne nouvelle. Les initiatives des entreprises pour améliorer la qualité des données notamment *via* des référentiels et des solutions MDM (*master data management*) sont des signes positifs de cette évolution. Et il était temps !

Car l'histoire des données est loin d'être terminée. Un autre phénomène relativement récent est venu se greffer aux problématiques de qualité de données : il s'agit de l'explosion du volume de données appelée plus communément « big data », terme anglais traduit depuis peu en français par celui de « mégadonnées ». Le développement

spectaculaire d'internet, et en particulier celui du web social (réseaux sociaux, blogs), la généralisation des intranets, l'essor de la mobilité avec les téléphones mobiles et les tablettes, la multiplication des capteurs et des objets connectés provoquent une avalanche de données à laquelle les organisations doivent faire face.

En conséquence, les entreprises sont submergées de données dont les volumes se chiffrent désormais en dizaines de téraoctets et parfois même en pétaoctets. Une bonne partie de ces nouvelles sources d'informations ne sont pas structurées, et les fournisseurs de bases de données et de solutions d'intégration se sont adaptés à cette nouvelle situation.

Ces dernières années, un énorme *buzz* autour du big data s'est développé à tel point qu'un certain nombre d'analystes n'y ont vu qu'un phénomène marketing de plus, visant à favoriser les ventes des fournisseurs de technologie. Les protagonistes des systèmes d'information, à commencer par les entreprises utilisatrices, s'aperçoivent à présent que le phénomène est bien réel. Et les enjeux sont de taille, car le big data n'est pas qu'une question technique de volumétrie et de stockage. Il constitue au contraire l'opportunité de comprendre le contenu de ces nouvelles sources et d'en tirer profit. La valeur des données non structurées issues des réseaux sociaux en est le parfait exemple. Quel service marketing n'a pas rêvé de connaître avec précision les sentiments des consommateurs à propos de ses produits ou de ses marques ?

Encore récemment, les moyens d'exploiter de telles volumétries de données, structurées ou non, étaient inexistants. La technologie est aujourd'hui disponible, les bases de données ont évolué et les solutions techniques dédiées à l'exploitation des données massives sont opérationnelles. La *business intelligence* laisse peu à peu place au *business analytics* et aux *algorithmes d'intelligence artificielle* tels que le *machine learning*. Ces solutions s'appuient sur des modèles mathématiques et des algorithmes pointus permettant entre autres de faire de l'analyse prédictive. Grâce à ces progrès remarquables, l'avalanche de données se transforme en valeur ajoutée, et les entreprises peuvent ainsi accroître leurs performances, être à la fois plus proactives et plus compétitives. En pratique, les entreprises sont de plus en plus nombreuses à mettre en œuvre des programmes de transformation digitale incluant un important volet big data.

Le big data ne se réduit pas à une problématique unique ; il se décline au contraire en différentes variantes selon le métier et la situation particulière de chaque entreprise. Les techniques doivent donc être choisies et intégrées en fonction de la situation. Les projets big data ont cependant des caractéristiques communes dans leur approche et leur méthode de mise en œuvre. Le présent ouvrage développe aussi cet aspect essentiel à la réussite des projets.

1. Objectifs de l'ouvrage

L'ouvrage a pour objectifs d'expliquer les enjeux du big data, de présenter différents cas d'usage et de détailler les méthodes et techniques utilisées. Il présente aussi les principes de mise en œuvre d'un projet big data avec les recommandations associées.

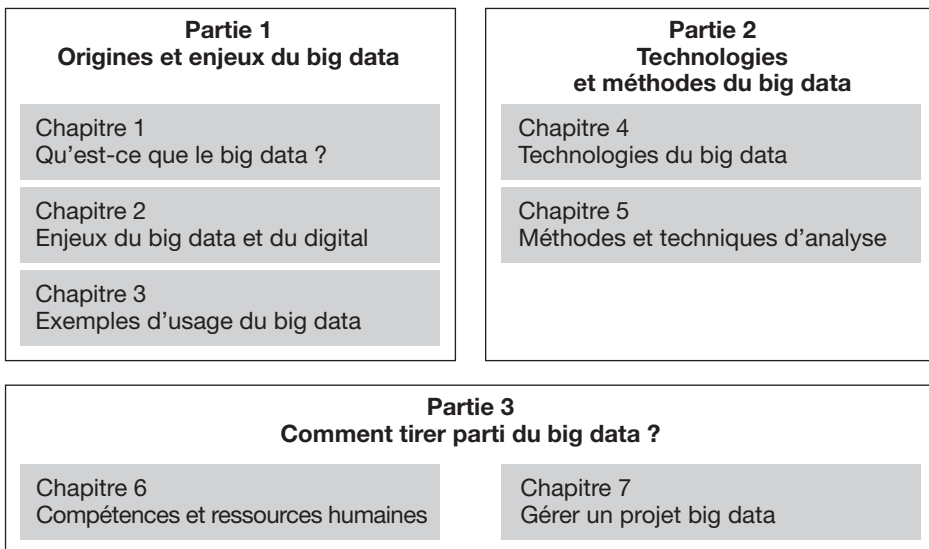
Cette seconde édition prend en compte les principales évolutions du big data, et approfondit les nouveaux enjeux numériques parmi lesquels figurent l'expérience client, les réseaux sociaux, les objets connectés, le *cloud* et l'analytique. Les exemples et les cas ont notamment été renouvelés pour coller aux préoccupations actuelles des entreprises. Enfin, l'ouvrage met l'accent sur les méthodes, les techniques et les ressources nécessaires pour permettre aux entreprises d'entrer avec succès dans l'ère de l'information à grande échelle.

Le livre s'adresse aux dirigeants d'entreprise, aux managers opérationnels et aux professionnels des systèmes d'information (responsables informatiques, directeurs et chefs de projets, consultants, ingénieurs). Les entreprises ne peuvent plus ignorer le phénomène du big data et il est essentiel que leurs responsables soient au fait des enjeux, des usages et des pratiques.

Le livre s'adresse aussi aux chercheurs et aux scientifiques confrontés à des volumétries de données de plus en plus importantes dans le cadre de leurs travaux. Outre le monde de l'entreprise, le big data concerne en effet des domaines divers tels que la recherche biomédicale, la météorologie, l'éducation, la criminologie et d'une façon plus générale la science, dont les applications et/ou les expérimentations nécessitent de plus en plus de données. Ainsi, l'informatique cognitive (*cognitive computing*) est un domaine prometteur s'appuyant sur des données volumineuses.

Cet ouvrage s'adresse aux étudiants des écoles de commerce, d'ingénieurs et des universités soucieux de comprendre les problématiques du big data auxquelles ils seront confrontés, de près ou de plus loin.

2. Organisation de l'ouvrage



La première partie, « Origines et enjeux du big data », permet de comprendre la problématique du big data :

- le chapitre 1, « Qu'est-ce que le big data ? », rappelle le rôle stratégique des données et définit ce qu'est le big data. Les différents types de données sont présentés, de même que les différents acteurs concernés ;
- le chapitre 2 présente les enjeux du big data et du digital. La création de valeur est au cœur de ce chapitre qui s'intéresse à la transformation digitale des organisations et à l'expérience client. Le chapitre traite aussi les sujets connexes que sont la qualité des données, la confidentialité et l'ouverture des données publiques. Le chapitre se termine par l'analyse des perspectives du big data ;
- le chapitre 3, « Exemples d'usage du big data », propose plusieurs cas pratiques. Les possibilités du big data sont illustrées au travers d'exemples concrets appliqués au marketing, à la santé, à l'énergie, aux services publics, à l'assurance et au *data journalism*.

La seconde partie présente les technologies et les méthodes utilisées pour traiter et analyser les mégadonnées :

- le chapitre 4 décrit « Les technologies du big data ». La problématique technique et les différentes solutions du marché sont présentées. Un focus spécifique sur les technologies Hadoop et Spark est réalisé. Les principaux fournisseurs de solutions sont également listés dans ce chapitre ;
- le chapitre 5, « Les méthodes et techniques d'analyse du big data », dresse un panorama des principales méthodes et techniques du *business analytics*. Les techniques de visualisation sont également développées dans ce chapitre.

La troisième partie fournit les clés permettant de tirer parti du big data :

- le chapitre 6, « Compétences et ressources humaines », s'intéresse aux ressources et profils nécessaires d'une part pour réaliser un projet big data et, d'autre part, pour utiliser les modèles et outils mis en œuvre. L'organisation des ressources est également abordée dans ce chapitre ;
- le chapitre 7, « Gérer un projet big data », précise les caractéristiques propres à ce type de projet, et définit la démarche de mise en œuvre. Des recommandations et les risques à prendre en compte complètent ce chapitre pour amener les projets vers le succès.

Le développement spectaculaire d'internet, des réseaux sociaux, de la technologie mobile et des objets connectés provoque une croissance exponentielle des données à laquelle toutes les entreprises sont confrontées : c'est le phénomène du big data.

Cette nouvelle édition de *Enjeux et usages du big data* prend en compte les principales évolutions du big data et approfondit les nouveaux enjeux numériques parmi lesquels figurent l'analytique, l'expérience client, le cloud, les réseaux sociaux et les objets connectés. Les exemples et les cas ont été entièrement renouvelés afin d'être au plus près des préoccupations actuelles des entreprises. L'ouvrage met également l'accent sur les méthodes, les techniques et les ressources nécessaires pour permettre aux entreprises d'entrer avec succès dans l'ère de l'information à grande échelle.

Cet ouvrage concis et didactique s'adresse aux dirigeants d'entreprise, aux managers opérationnels et aux professionnels des systèmes d'information, ainsi qu'aux étudiants des écoles de commerce, d'ingénieurs et des universités, soucieux de comprendre les problématiques et les applications du big data.

Christophe Brasseur, *ingénieur diplômé de l'ESTP et de l'IAE, travaille depuis plus de 25 ans dans les systèmes d'information. Il a participé à de nombreux projets internationaux dans des postes de conseil et de management, principalement dans le secteur des services publics, de l'industrie et de l'énergie. Il est actuellement Senior Manager chez Capgemini.*

