

BIG DATA

Analyse et valorisation de masses de données

Florent **BERANGER**

Directeur de projets Décisionnel

Smile

OPEN SOURCE SOLUTIONS

www.smile.fr • +33 (0)1 41 40 11 00 • contact@smile.fr
www.smile-oss.com • blog.smile.fr • twitter: @GroupeSmile

PREAMBULE

SMILE

Smile est une **société d'ingénieurs experts** dans la mise en œuvre de **solutions open source** et l'intégration de systèmes appuyés sur l'open source. Smile est membre de l'**APRIL**, l'association pour la promotion et la défense du logiciel libre, du **PLOSS** – le réseau des entreprises du Logiciel Libre en Ile-de-France et du **CNLL** – le conseil national du logiciel libre.

Smile compte près de 700 collaborateurs dans le monde, dont plus de 500 en France (février 2014), ce qui en fait *le premier intégrateur français et européen de solutions open source*.

Depuis 2000 environ, **Smile mène une action active de veille technologique** qui lui permet de découvrir les produits les plus prometteurs de l'open source, de les qualifier et de les évaluer, de manière à proposer à ses clients les produits les plus aboutis, les plus robustes et les plus pérennes.

Cette démarche a donné lieu à **toute une gamme de livres blancs** couvrant différents domaines d'application. La gestion de contenus (2004), les portails (2005), la business intelligence (2006), la virtualisation (2007), la gestion électronique de documents (2008), les PGIs/ERPs (2008), les VPN open source (2009), les Firewall et Contrôle de flux (2009), les Middleware orientés messages (2009), l'e-commerce et les Réseaux Sociaux d'Entreprise (2010), le Guide de l'open source et NoSQL (2011) et plus récemment, Mobile et Recensement et audit (2012). Chacun de **ces ouvrages présente une sélection des meilleures solutions open source** dans le domaine considéré, leurs qualités respectives, ainsi que des retours d'expérience opérationnels.

Au fur et à mesure que des solutions open source solides gagnent de nouveaux domaines, Smile sera présent pour proposer à ses clients d'en bénéficier sans risque. Smile apparaît dans le paysage informatique français comme **le prestataire intégrateur de choix pour accompagner** les plus grandes entreprises dans l'adoption des meilleures solutions open source.

Ces dernières années, Smile a également étendu la gamme des services proposés. Depuis 2005, un département consulting accompagne nos clients, tant dans les phases d'avant-projet, en recherche de solutions, qu'en accompagnement de projet. Depuis 2000, Smile dispose d'un studio graphique, devenu en 2007 Smile Digital – agence interactive, proposant outre la création graphique, une expertise e-marketing, éditoriale, et interfaces riches. Smile dispose aussi d'une agence spécialisée dans la TMA (support et l'exploitation des applications) et d'un centre de formation complet, Smile Training.

Enfin, Smile est implanté à Paris, Lille, Lyon, Grenoble, Nantes, Bordeaux, Marseille et Montpellier. Et présent également en Espagne, en Suisse, au Benelux, en Ukraine, au Maroc et en Côte d'Ivoire.

QUELQUES REFERENCES DE SMILE

SMILE est fier d'avoir contribué, au fil des années, aux plus grandes réalisations Web françaises et européennes. Vous trouverez ci-dessous quelques clients nous ayant adressé leur confiance.

Sites Internet

EMI Music, Salon de l'Agriculture, Mazars, Areva, Société Générale, Gîtes de France, Patrice Pichet, Groupama, Eco-Emballage, CFnews, CEA, Prisma Pub, Veolia, NRJ, JCDecaux, Spie, PSA, Boiron, Larousse, Dassault Systèmes, Action Contre la Faim, BNP Paribas, Air Pays de Loire, Forum des Images, IFP, BHV, ZeMedical, Gallimard, Cheval Mag, Afssaps, Bénéteau, Carrefour, AG2R La Mondiale, Groupe Bayard, Association de la Prévention Routière, Secours Catholique, Canson, BNP Paribas, Bouygues Telecom, CNIL, Eiffage, Poweo, Mazars, Kering...

Portails, Intranets et Systèmes d'Information

HEC, Bouygues Telecom, Prisma, Veolia, Arjowiggins, INA, Primagaz, Croix Rouge, Eurosport, Invivo, Faceo, Château de Versailles, Eurosport, Ipsos, VSC Technologies, Sanef, Explorimmo, Bureau Veritas, Région Centre, Dassault Systèmes, Fondation d'Auteuil, INRA, Gaz Electricité de Grenoble, Ville de Niort, Ville de Saint-Etienne, Ministère de la Culture, PagesJaunes Annonces, Feu Vert, Bouygues Immobilier, Biomérieux, Generali ...

E-Commerce

Krys, La Halle, Gibert Joseph, De Dietrich, Adenclassifieds, Macif, Furet du Nord, Gîtes de France, Camif Collectivités, GPdis, Projectif, ETS, Bain & Spa, Yves Rocher, Bouygues Immobilier, Nestlé, Stanhome, AVF Périmédical, CCI, Pompiers de France, Commissariat à l'Energie Atomique, Snowleader, Darjeeling, Le Bon Marché, VF Corporation, Histoire d'Or, MyEvian, Chantelle, Yamaha, Wesco...

ERP et Décisionnel

Veolia, La Poste, Christian Louboutin, Eveha, Sun'R, Home Ciné Solutions, Pub Audit, Effia, France 24, Publicis, iCasque, Nomadvantage, Gets, Nouvelles Frontières, Anevia, Jus de Fruits de Mooréa, Espace Loggia, Bureau Veritas, Skyrock, Lafarge, Cadremploi, Meilleurmobilite.com, Groupe Vinci, IEDOM (Banque de France), Carrefour, Jardiland, Trésorerie Générale du Maroc, Ville de Genève, ESCP, Sofia, Faiveley Transport, INRA, Deloitte, Yves Rocher, ETS, DGAC, Generalitat de Catalunya, Gilbert Joseph, Perouse Médical, V'Lille, Casden, Corsair...

Gestion documentaire

Generali, HEC, JCDecaux, Serimax, Pierre Audoin Consultant, Alstom Power services, NetasQ, CS informatique, SNCF - Direction du matériel, Mazars, EDF R&D, EDF Nucléaire, Conseil Régional du Centre, Leroy Merlin, Primagaz, Renault F1, INRIA, Ministère belge de la Communauté Française, APAVE, Conseil Général de Loire Atlantique, CNIL, Services du Premier Ministre...

Infrastructure et Hébergement

Agence Nationale pour les Chèques Vacances, Pierre Audoin Consultants, Rexel, Motor Presse, OSEO, Sport24, Eco-Emballage, Institut Mutualiste Montsouris, ETS, Ionis, Osmoz, SIDEL, Atel Hotels, Cadremploi, SETRAG, Institut Français du Pétrole, Mutualité Française, Orange, Bouygues Télécom, Fiducial, Ministère du Développement Durable...

Consultez nos références, en ligne, à l'adresse : <http://www.smile.fr/clients>.

SOMMAIRE

PREAMBULE..... 2

SMILE 2

QUELQUES REFERENCES DE SMILE..... 3

SOMMAIRE 5

EN RESUME..... 7

LE BIG DATA GENERATEUR D’OPPORTUNITES POUR LES ENTREPRISES ET COLLECTIVITES 7

UNE TENDANCE DE FOND POUR L’ANALYSE DE DONNEES MASSIVES 9

CHECKLIST D’UN PROJET DECISIONNEL BIG DATA 11

 CADRER LES OPPORTUNITES METIER 11

 CADRER L’ARCHITECTURE 11

CE LIVRE BLANC 12

SUJETS TRAITES 12

CONCEPTS ET DEFINITIONS 14

BIG DATA 14

ENTREPOT DE DONNEES OU DATAWAREHOUSE 14

STOCKAGE DISTRIBUE - NOSQL..... 15

 LIMITES DES SGBDR DANS LES ARCHITECTURES DISTRIBUEES 15

 PRINCIPES DE DISTRIBUTION ET DE REPLICATION DES DONNEES 16

 STRUCTURES DES BASES ET ORGANISATION DES DONNEES NOSQL 16

INTEGRATION ET TRAITEMENT (DISTRIBUE) DE DONNEES MASSIVES..... 19

 ETL 19

 FRAMEWORKS DE TRAITEMENTS DISTRIBUES - MAP-REDUCE 19

L’ANALYSE MULTIDIMENSIONNELLE OU OLAP..... 19

REQUETAGE AD-HOC EN LANGAGE NATUREL 20

DATA MINING..... 20

CAS D’USAGES 21

USAGES COUVERTS PAR LES SOLUTIONS BIG DATA POUR L’ANALYSE ET LA VALORISATION..... 21

MARKETING 21

 VUE A 360° DES CLIENTS ET ANALYSE DES COMPORTEMENTS DE CONSOMMATION 21

 E-COMMERCE..... 22

 RESSENTI SUR LES SERVICES, PRODUITS ET CONCEPTS 22

IMPLANTATION DE POINTS DE VENTE	22
LOGISTIQUE ET CHAINE D'APPROVISIONNEMENT	22
LE BIG DATA AU SERVICE DE LA TRAÇABILITE	22
LE BIG DATA FACTEUR D'OPTIMISATION DE LA CHAINE D'APPROVISIONNEMENT	23
TELECOMS	23
PANORAMA DES SOLUTIONS BIG DATA POUR LA BI	24
COMPOSANTS D'INTEGRATION ET DE TRAITEMENT DE DONNEES	25
SYNTHESE	25
HADOOP	26
ETL TALEND FOR BIG DATA.....	28
ETL PENTAHO DATA INTEGRATION	34
STOCKAGE DE MASSES DE DONNEES	38
SYNTHESE	38
FEDERATION DE DONNEES NoSQL DANS DES BASES RELATIONNELLES	38
MONGODB	40
ELASTICSEARCH	41
ANALYSER ET RESTITUER DES MASSES DE DONNEES.....	42
SYNTHESE	42
PENTAHO BUSINESS ANALYTICS	43
JASPERSOFT BI SUITE.....	45
SPAGOBİ	47
ELASTICSEARCH KIBANA	50
REMERCIEMENTS.....	52

EN RESUME

LE BIG DATA GENERATEUR D'OPPORTUNITES POUR LES ENTREPRISES ET COLLECTIVITES

Chaque jour, la quantité de données créées et manipulées ne cesse d'augmenter, et ce quel que soit le secteur d'activité concerné.

Cette croissance, exponentielle, est liée à :

- l'évolution du nombre d'utilisateurs des solutions IT
- l'évolution des périmètres couverts et des usages (mobile,...)
- la génération de données par des machines
- la finesse de l'information tracée
- la croissance des volumes opérationnels
- l'évolution de l'historique de données disponible.

Ces données sont issues de sources multiples :

RFID, compteurs d'énergie, opérations commerciales en volumes, transactions financières, blogs, réseaux de capteurs industriels, réseaux sociaux, téléphonie, indexation Internet, parcours de navigation GPS, détails d'appels en call center, e-commerce, dossiers médicaux, informatique embarquée, Internet des objets, données biologiques, textes de tickets ou mails, sondages,...

Ces masses de données apportent des opportunités d'analyses plus larges et plus fines ainsi que de nouveaux usages de l'information, qu'elle soit pleinement ou partiellement structurée à la source.

La question n'est plus "Le Big Data peut-il devenir un avantage concurrentiel pertinent ?" mais **"Comment pouvons-nous exploiter les possibilités offertes par ces solutions pour optimiser nos processus d'analyse et de prise de décision ?"**.

En effet, les masses de données constituent un matériau brut. Au delà de leur exploitabilité (pertinence, disponibilité et qualité), c'est la capacité à les transformer en analyse et en service qui apporte une valeur maximale.

“ Big Data Analyse et valorisation de masses de données ”



Le Big Data transforme progressivement les organisations autour de la valorisation de l'information. Avec la finesse d'information sur les opérations passées et de plus en plus d'informations prospectives, le Big Data va permettre l'éclosion de modèles prédictifs plus pertinents.

UNE TENDANCE DE FOND POUR L'ANALYSE DE DONNEES MASSIVES

Les systèmes de base de données relationnelles et les outils d'aide à la décision n'ont initialement pas été créés afin de manipuler une telle quantité et richesse de données, et il peut vite devenir compliqué et improductif pour les entreprises d'accéder à ces masses de données avec les outils classiques.

Cette nouvelle problématique a donné naissance aux systèmes de gestion de base de données appelés « NoSQL » (Not Only SQL), qui ont fait le choix d'abandonner certaines fonctionnalités des SGBD classiques au profit de la simplicité, la performance et de la capacité à monter en charge.

Des frameworks comme Hadoop ont également été créés et permettent le requêtage, l'analyse et la manipulation de ces données en masse.

Nous relevons que les principales solutions de Big Data sont Open Source. Ce contexte favorise leur vitesse de développement et de diffusion au sein des entreprises et collectivités.

Et ce à moindre coût par rapport à des solutions dont l'évolution de la capacité est verticale : coût des ressources matérielles, licences,...

Il est possible de mettre en place une solution décisionnelle Big Data complète uniquement basée sur des solutions Open Source sans coût de licence. Toutefois, des versions commerciales basées sur de l'Open Source apportent des facilités qui vont dans le sens de la productivité de mise en oeuvre et de l'exploitabilité des solutions avec des outils d'administration complémentaires notamment.

Beaucoup d'entreprises et de collectivités publiques utilisent déjà des solutions Big Data, souvent hébergées dans le cloud (ex : Google Analytics, réseaux sociaux, Salesforce,...). Les solutions Big Data ont fait leurs preuves et sont mûres pour un déploiement en production.

Les fonctionnalités de visualisation graphique (DataViz), pour illustrer des analyses portant sur des masses de données, et de datamining prennent avec le Big Data toute leur importance.

Techniquement, le format JSON (*JavaScript Object Notation*) émerge comme un standard d'échange et d'exploitation de données (massives), en complément du SQL, comme cela s'observe aussi sur les solutions web non décisionnelles.

Ce mouvement va de pair avec le développement des bibliothèques JavaScript de visualisation graphique avancées (d3.js,...) et des frameworks Javascript d'interactivité avec les données.

“ Big Data Analyse et valorisation de masses de données ”



Nous relevons aussi les possibilités de consolidation de données (massives) et hétérogènes “à la volée” en complément de l’entrepôt de données : la fédération de données.

**CHECKLIST D'UN PROJET
DÉCISIONNEL BIG DATA**

Au delà des principes et bonnes pratiques de mises en oeuvre de solutions IT, une vigilance sur les points suivants peut éviter des écueils lors du cadrage d'un projet Big Data :

Cadrer les opportunités métier

- identifier des leviers de gain d'exploitation de masses de données sur les activités de l'entreprise
- identifier le périmètre (légal, technique, historique) d'information disponible : SI interne, données fournies par des partenaires, OpenData, ...
- identifier le ou les cas d'utilisation résultant de l'adéquation entre les leviers de gain et le périmètre d'information disponible

Cadrer l'architecture

- définir une architecture flexible adaptée au(x) cas d'utilisation; il n'existe pas un modèle d'architecture Big Data idéal adapté à tous les usages
- valider la disponibilité et l'exploitabilité des données sources
- valider l'architecture (matérielle, réseau, applicative) par un test de montée en charge.

CE LIVRE BLANC

Cet ouvrage constitue le premier livre blanc de Smile sur le sujet. Nous espérons qu’il vous apportera l’information souhaitée et qu’il vous sera agréable à parcourir.

Comme les autres livres blancs publiés par Smile, cet ouvrage s’efforce de réunir :

- une approche générale de la thématique, ici : l’analyse et la valorisation de masses de données, ses concepts, ses champs d’application, ses besoins spécifiques.
- un recensement des meilleures solutions Open Source dans ce domaine.
- une présentation assez complète de ces solutions, de leurs forces, de leurs limites, de leur maturité et de leur aptitude à satisfaire des besoins opérationnels.

Cette étude, réalisée par notre équipe de consultants, a été fondée sur plusieurs années de travail de recherche et de premiers déploiements effectifs de solutions Big Data.

Cet ouvrage vient compléter livres blancs Business Intelligence et NoSQL.

Les marques et logos présents dans ce livre blanc sont la propriété des entreprises concernées.

SUJETS TRAITES

Ce livre blanc est concentré sur les solutions applicatives d’analyse et de valorisation de masses de données. D’autres aspects de l’exploitation des masses de données sont importants mais non décrits ici :

- **Qualité des données** : prendre en compte la qualité et le nettoyage des données, ainsi que la gestion du cycle de vie des données référentielles dans le scope du projet évite d’aboutir à une masse de données inexploitable. Des solutions de traitement, qualification et nettoyage automatique des données existent : fonctionnalités intégrées aux flux de données ETL, briques complémentaires telles DataQuality de Talend.
- **Infrastructures techniques** : les solutions Big Data nécessitent une architecture répartie. La composante système et réseaux est un facteur clé de performance et d’exploitabilité d’une solution Big Data.
- **Sécurité de l’information** : les aspects de sécurisation des accès et de gestion de l’intégrité des données sont importants pour la mise en oeuvre d’une solution pérenne.

“ Big Data Analyse et valorisation de masses de données ”



- **Respect de la vie privée** : les solutions Big Data peuvent apporter une puissance informative importante. Cette puissance doit respecter les libertés individuelles.
- **Solutions** : l'écosystème des solutions Big Data est riche et évolutif. Il nous serait difficile de détailler toutes les solutions. Nous nous sommes concentrés sur les solutions les plus pertinentes à l'heure actuelle.

CONCEPTS ET DEFINITIONS

BIG DATA

Le Big Data consiste en un/des ensemble(s) de données plus ou moins structurées qui deviennent tellement volumineux qu'ils sont difficiles à travailler avec des outils classiques de gestion de base de données.

En 2012, Gartner a posé les bases de la définition du Big Data, basée sur les 3V :

- Volume
- Vitesse
- Variété des données.

→ "Big data is high volume, high velocity, and/or high variety information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization."

Sans seuil ni repère, beaucoup de bases de données classiques peuvent prétendre à répondre à ces trois critères.

Dans le présent livre blanc, pour les usages d'analyse, nous compléterons pragmatiquement la combinaison des 3V avec une considération de volumétries en dizaines de millions d'enregistrements minimum.

ENTREPOT DE DONNEES OU DATAWAREHOUSE

L'entrepôt de données est une base de données qui concentre de l'information issue de différents systèmes d'information de l'entreprise, à des fins d'analyse et de reporting des activités et marchés.

STOCKAGE DISTRIBUE - NoSQL

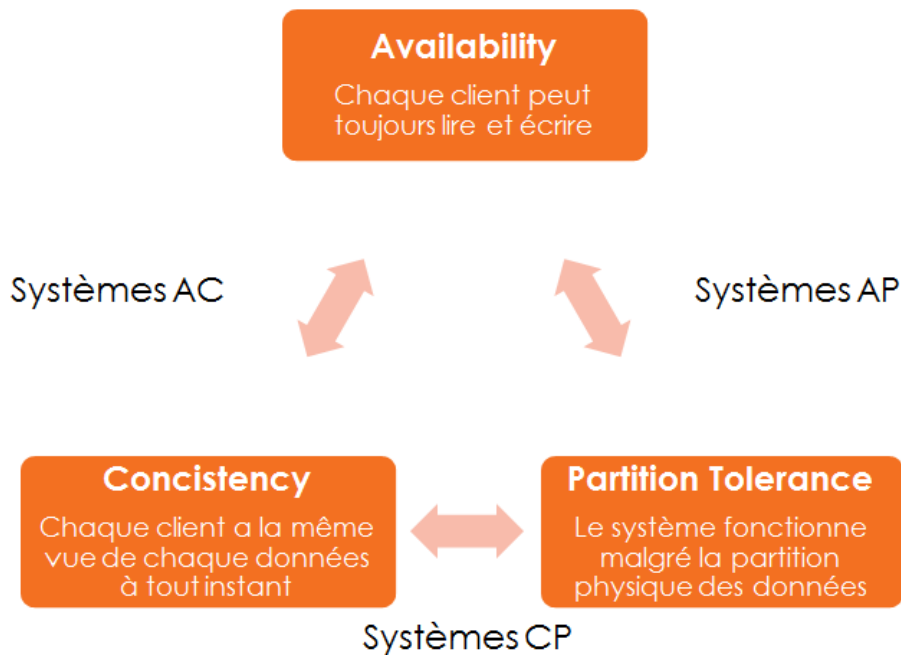
NoSQL, pour Not Only SQL désigne les systèmes de gestion de base de données qui ne s'appuient plus, du fait des volumétries et de la variété des données contenues, sur une architecture relationnelle et transactionnelle.

Ces systèmes privilégient la simplicité et l'évolutivité de la capacité via des architectures distribuées.

Limites des SGBDR dans les architectures distribuées

Outre leur modèle relationnel, la plupart des moteurs de SGBDs relationnels sont transactionnels ce qui leur impose le respect des contraintes Atomicity Consistency Isolation Durability, communément appelé par son acronyme ACID.

Théorème de CAP



Il est actuellement impossible d'obtenir ces trois propriétés en même temps dans un système distribué. Sur de nombreux SGBDR classiques, la réplication devient plus complexe avec de fortes volumétries et une forte vélocité des données.

Principes de distribution et de réplication des données

Les capacités de montée en charge des bases NoSQL reposent, au delà de leur simplicité, sur la distribution (sharding) et la réplication des données sur différents noeuds (cluster de quelques serveurs à plusieurs DataCenter).

Pour simplifier, une analogie peut être faite entre les mécanismes de partitionnements verticaux (sur plusieurs tables physiques de la même instance) de certains moteurs de bases de données relationnelles et la distribution horizontale (sur plusieurs serveurs) des données en NoSQL.

Les données peuvent également être répliquées, sur un principe analogue aux mécanismes de stockage en RAID, afin de garantir un haut niveau de service, même en cas de problème ou de maintenance d'un nœud du cluster.

Structures des bases et organisation des données NoSQL

Il existe plusieurs paradigmes au niveau des systèmes de stockage NoSQL :

Type documentaire

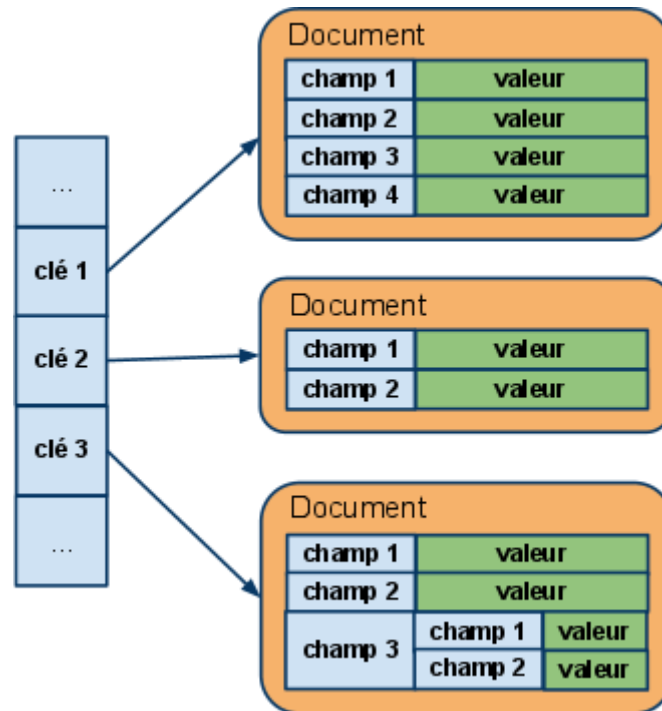
Les bases de données documentaires sont constituées de collections de documents. Les collections sont généralement assimilées à des tables d'un modèle relationnel.

Bien que les documents soient structurés, ces bases sont sans schéma de données prédéfini. Il n'est donc pas nécessaire de définir au préalable l'ensemble des champs utilisés dans un document. Les documents peuvent donc avoir une structure hétérogène au sein de la base.

Un document est composé de champs et de valeurs associées, ces dernières pouvant être requêtées. Les valeurs peuvent être, soit d'un type simple (entier, chaîne de caractère, date, ...), soit composées de plusieurs couples clé/valeur (imbrications nested sets).

Les structures de données sont donc très souples.

“ Big Data Analyse et valorisation de masses de données ”



La souplesse du modèle de données, les performances et les capacités de requêtage orientent l'usage des bases documentaires vers du stockage opérationnel de masse (ODS) dans un système décisionnel.

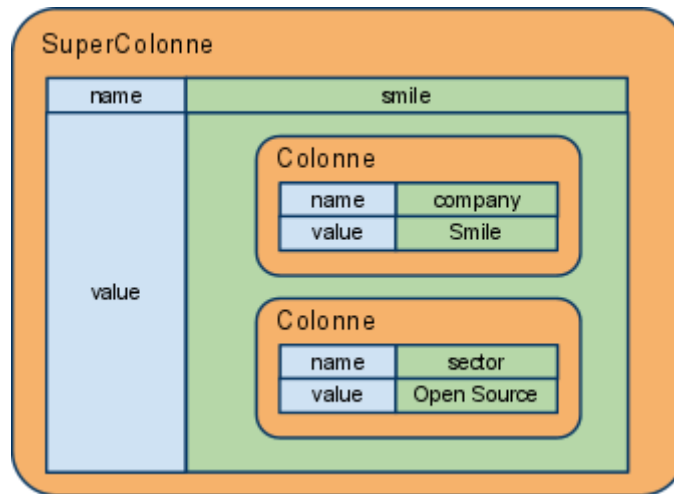
Type graphe

Au delà du moteur de stockage sous la forme d'une base documentaire, ce type de base propose également des relations entre objets. Ces derniers sont orientés et peuvent porter des propriétés.

Type orienté colonnes

La colonne représente l'entité de base de la structure de données. Chaque colonne d'un objet est défini par un couple clé / valeur. Une colonne contenant d'autres colonnes est nommée super-colonne.

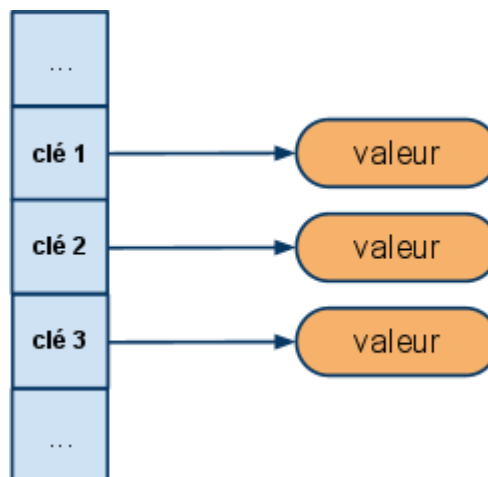
“ Big Data Analyse et valorisation de masses de données ”



Ces types de bases sont adaptés au stockage opérationnel de masse (ODS) et de source d’analyses massives dans un système décisionnel.

Type clé/valeur

Dans ce modèle, chaque objet/enregistrement est identifié par une clé unique. La structure de l’objet est libre.



Dans ce modèle on ne dispose généralement que des quatre opérations Create, Read, Update, Delete (CRUD) en utilisant la clé de l’enregistrement à manipuler.

Du fait des limites fonctionnelles d’accès aux données de ces types de base, nous ne leur voyons pas d’application décisionnelle.

**INTEGRATION ET
TRAITEMENT (DISTRIBUE) DE
DONNEES MASSIVES**

ETL

Afin d'alimenter un datawarehouse à partir des différentes sources de données ou de synchroniser en batch des données entre systèmes, on utilise une gamme d'outils appelés ETL, pour « Extract, Transform, Load ».

Comme le nom l'indique, ces outils permettent d'extraire des données à partir de différentes sources, de les transformer (rapprochement, format, dénomination, calculs), et de les charger dans la base de données cible, comme un datawarehouse dans le cas d'un projet décisionnel. L'ETL permet de masquer, grâce à une modélisation visuelle, la complexité de manipulations (réparties) des données (hétérogènes) au sein des traitements et ainsi d'en réduire fortement les coûts de développement et maintenance.

Un ETL est généralement composé d'un studio de modélisation des traitements ainsi que d'un ou plusieurs environnements d'exécution et des outils d'administration voire de visualisation de données suivant les versions.

Frameworks de traitements distribués - Map-Reduce

Modèle d'architecture portant sur la distribution et la répartition des traitements de données sur plusieurs noeuds d'une grappe de serveurs (cluster).

Dans l'étape Map, les données à traiter et traitements à effectuer sont répartis sur les noeuds de traitement.

Dans l'étape Reduce, les noeuds de traitements remontent leur résultat pour agrégation (il peut y avoir plusieurs niveaux de traitement).

**L'ANALYSE
MULTIDIMENSIONNELLE OU
OLAP**

L'analyse multidimensionnelle permet l'analyse de mesures suivant différents aspects métiers appelés dimensions ou axes d'analyse et ce, à plusieurs niveaux de regroupement.

Par exemple, la mesure de Montant HT d'une ligne de facture peut être agrégée par :

- jour → mois → trimestre → année
- produit → catégorie de produits → ligne de produits
- client → segment de client.

**REQUETAGE AD-HOC EN
LANGAGE NATUREL**

Le requêtage ad-hoc permet à des non informaticiens de construire visuellement des requêtes, en s'appuyant sur un dictionnaire d'informations en langage naturel (métadonnées) faisant abstraction du langage technique d'accès aux bases de données (SQL, JSON).

DATA MINING

Le data mining consiste à rechercher des informations statistiques utiles cachées dans un grand volume de données.

L'utilisateur est à la recherche d'une information statistique qu'il n'identifie pas encore.

CAS D'USAGES

USAGES COUVERTS PAR LES SOLUTIONS BIG DATA POUR L'ANALYSE ET LA VALORISATION

Il existe de nombreux cas d'usage des solutions de valorisation et d'analyse massive de données. Nous en avons détaillé quelques unes ci-dessous mais nous pouvons aussi citer l'analyse fine de processus, la recherche scientifique, les analyses politiques et sociales, l'analyse de données de capteurs sur les chaînes industrielles...

MARKETING

Le Big Data transforme en profondeur les métiers du marketing, avec les facilités suivantes :

Vue à 360° des clients et analyse des comportements de consommation

Une vue complète de chaque client nécessite la manipulation de larges ensembles de données:

- informations sur le client : stockées dans le SI, disponibles sur les réseaux sociaux publics
- comportements d'achat : détail des commandes, fréquence, canaux
- segmentation
- parcours / historique de la relation depuis la prospection
- niveau d'engagement; parrainage d'autres clients
- enquêtes de satisfaction
- expérience d'utilisation, utilisation des services après-vente.

La collecte et la consolidation de toutes ces données représente une tâche fastidieuse, rarement faite ou uniquement sur un petit panel de clients. Les solutions Big Data peuvent permettre d'automatiser cela et apporter les gains suivants :

- optimisation de l'adéquation des produits et services proposés
- affinage des ciblage et optimisation des communications avec chaque client : canal, message,...

E-commerce

Les principales solutions d'analyse d'audience web (pages visitées, recherches,...) du marché utilisent des solutions Big Data. Des solutions d'analyse Big Data complémentaires peuvent apporter un plus :

- analyse des critères et freins de transformation en fonction d'informations complémentaires aux mesures d'audience web
- corrélation avec les retours, livraisons et données financières
- analyse fine des interactions des utilisateurs avec le site e-commerce : Real User Monitoring.

Elles permettent également de faire le lien avec l'analyse d'achat :

- analyse du tunnel de vente
- analyse des comportements d'achat ou d'usage des clients afin d'optimiser leur expérience
- détection de fraudes
- les bases NoSQL documentaires sont particulièrement adaptées à l'entreposage et l'analyse de données souples et complexes, telles les caractéristiques de produits.

Ressenti sur les services, produits et concepts

Analyse de mots postés sur les réseaux sociaux publics.

Implantation de points de vente

La technologie Big Data offre la possibilité de corréliser des données de différentes natures et de différentes sources pour déterminer le meilleur emplacement pour un point de vente :

- OpenData
- données géographiques
- données socio-économiques
- informations disponibles sur le marché et la concurrence.

LOGISTIQUE ET CHAINE D'APPROVISIONNEMENT

Le Big Data au service de la traçabilité

Les solutions Big Data permettent une pleine traçabilité des opérations logistiques :

- mouvements de stock - RFID
- produits frais ou sensibles

“ Big Data Analyse et valorisation de masses de données ”

- suivi de flotte ou de colis, y compris lors de transport inter-modal

Ces solutions facilitent les opérations de suivi des voyages dans le temps : geo corrdoring, analyse des voyages et taux de rotation

Le Big Data facteur d'optimisation de la chaîne d'approvisionnement

La masse de données disponible sur tous les mouvements permet d'analyser et de piloter plus finement les processus logistiques et d'approvisionnement.

La richesse d'information permet de combiner les différents facteurs de qualité (délais, défauts, qualité de service,...) et économiques (prix d'achat, coût de possession et de stockage,...) dans les analyses.

Le Big Data permet d'intégrer plus facilement les données logistiques dans les informations du cycle de vie des objets (commande, logistique, exploitation, recyclage,...) et permet ainsi une vision à 360° autour de la fonction d'approvisionnement.

TELECOMS

Les télécoms génèrent des masses de données sur les flux transités. Le Big Data est une solution utile pour :

- l'analyse de capacité
- la segmentation des usagers et des comportements d'usage des réseaux
- la corrélation avec les processus de vente et de support
- la qualité de service de réseaux complexes, la corrélation avec les appels aux call center.

“ Big Data Analyse et valorisation de masses de données ”

PANORAMA DES SOLUTIONS BIG DATA POUR LA BI

Il est important de noter que les solutions, notamment d'intégration et de traitement ne sont pas concurrentes mais souvent complémentaires voire intégrées.

Par exemple :

- l'intégration de briques de traitement et requêtage Hadoop avec du stockage MongoDB ou Cassandra.
- plusieurs ETL peuvent s'appuyer sur les frameworks de traitement distribué Hadoop.

Il est intéressant de constater que les principales technologies Big Data sont initiées par des acteurs majeurs du Web tels Google, Facebook, Twitter, Yahoo puis passées sous licence libre ce qui leur permet un développement et une diffusion rapide.


Nous relevons également la forte présence de la fondation Apache dans ce domaine de solutions.

**COMPOSANTS D'INTEGRATION ET DE
TRAITEMENT DE DONNEES**

Synthèse

Type	Solution	Site web de la solution
Framework de traitement	Apache Flume	http://flume.apache.org
Framework de requête	Apache Hive	http://hive.apache.org
Framework de requête et traitement	Apache Pig	https://pig.apache.org
Framework de requête	Cloudera Impala	http://www.cloudera.com/content/cloudera/en/products-and-services/cdh/impala.html
ETL	Talend for Big Data	http://fr.talend.com/products/big-data
ETL	Pentaho Data Integration	http://www.pentaho.fr/explore/pentaho-data-integration
ESB	Mule ESB	http://www.mulesoft.org
Framework de traitement	Hadoop YARN & MapReduce	https://hadoop.apache.org
Framework de traitement	Storm	http://storm-project.net

Hadoop

Editeur : Fondation Apache Licence : Apache License V2 et commerciales (suivant la distribution et la version) Version actuelle : 2 (+ suivant les composants)	
--	---



Présentation

Hadoop est un ensemble de projets et d'outils Open source de la fondation Apache permettant de stocker et traiter massivement des données. Hadoop a été développé à l'origine par Facebook et Yahoo.

Il existe plusieurs distributions d'Hadoop, parmi lesquelles on distinguera les principales à l'heure actuelle : HortonWorks, Cloudera et MapR.

Framework de traitements parallélisés Map-Reduce

Hadoop Map-Reduce est un puissant framework Java de traitement de données massives. A noter que dans le cas de l'utilisation conjointe avec HDFS et HBase et suivant la configuration du cluster Hadoop, il est possible qu'une partie des traitements soient effectués au niveau des noeuds de stockage, afin de limiter les échanges de données massives entre noeuds du cluster.

HDFS : Hadoop Distributed File System

HDFS est un système de fichiers distribué sur des noeuds d'un cluster Hadoop. HDFS est adapté au stockage et la réplication de fichiers de grande taille (>256MB).

Hbase

HBase est une base de données NoSQL répartie en colonnes, inspirée de Google BigTable. La mise en oeuvre de HBase repose généralement sur un système de fichiers répartis HDFS.

Hive

Hadoop Hive permet d'exploiter des traitements MapReduce de manière analogue à une base de données. En effet, des connecteurs JDBC et ODBC pour Hive sont disponibles.

Oozie

Oozie est un moteur de workflow et de coordination de tâches Hadoop (Map-Reduce, Pig).

Mahout

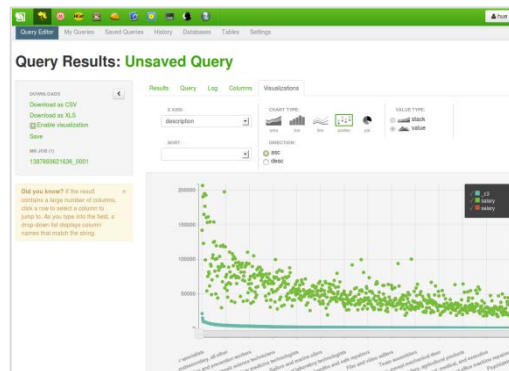
Mahout est une librairie Java qui permet d'implémenter différents algorithmes de data mining sur un cluster Hadoop.

Ces algorithmes sont développés à partir de MapReduce. Cependant, ils ne se limitent pas uniquement à Hadoop et certains fonctionnent sur d'autres environnements, dont non distribués.

Hue

Hue est un portail web d'exploitation de clusters Hadoop :

- requêtes Hive (Beeswax) :



- éditer, gérer et exécuter des traitements (jobs MapReduce et scripts Pig)


Usages et possibilités Big Data

L'ensemble Hadoop fournit plusieurs briques puissantes pour le décisionnel Big Data :

- l'entreposage de données opérationnelles (ODS HDFS ou Hbase) ou en entrepôt de données (Hbase et Hive).
- l'intégration et le traitement parallélisé de données (YARN/Map-Reduce, Pig)
- le requêtage et l'analyse de masses de données (Hive+YARN/Map-Reduce, Pig)
- le datamining (Mahout)

Notons que les principaux portails décisionnels Open Source intègrent directement un connecteur Hive pour une exploitation des données traitées dans un cluster Hadoop.

ETL Talend for Big Data

<p>Editeur : Talend Licences : Apache V2 et commerciale (suivant la version) Version actuelle : 5.4.1</p>	
---	---

Présentation

Éditeur et solutions

Talend est un éditeur basé en France (Talend SA) et en Californie (Talend Inc.). La société Talend, fondée en 2005, est soutenue dans son développement par des investisseurs tels Idinvest Partners (AGF Private Equity), Silver Lake Sumeru, Balderton Capital, Bpifrance et Iris Capital. Talend a réussi une levée de fonds de 40 millions de dollars fin 2013.

Talend offre un large éventail de solutions middleware répondant aux besoins de gestion de données et d'intégration d'applications, au travers une plateforme unifiée et flexible :

- l'intégration de données (ETL)
- la qualité de données (DQ)
- les architectures orientées services (ESB)
- la gestion de référentiels de données (MDM)
- la gestion de processus d'information (BPM).

Talend obtient une reconnaissance forte de la part des observateurs tel le Gartner (magic quadrants).

Les solutions sont disponibles en version communautaire (Talend Open Studio for Data Integration / Big Data) et en version commerciale avec des fonctionnalités supplémentaires et un support éditeur.

Les fonctionnalités ETL classiques de Talend sont présentées plus en détail dans le livre blanc BI (<http://www.smile.fr/Livres-blancs/Erp-et-decisionnel/Le-decisionnel-open-source>).

Talend et le Big Data

Talend propose depuis début 2012 une gamme de solutions Big Data, allant de la version Open Studio à la plateforme d'intégration massive de données (Talend Platform for Big Data).

Talend a établi des partenariats avec des acteurs majeurs du Big Data, notamment : Cloudera, EMC Greenplum, Google, HortonWorks, MapR.

Plus d'informations :

- <http://fr.talend.com/solutions/etl-analytics>

“ Big Data Analyse et valorisation de masses de données ”

- <http://www.talend.com/solutions/big-data>
- <http://fr.talend.com/products/platform-for-big-data>

Fonctionnalités

ETL Talend Open Studio for Big Data

Talend est un ETL de type « générateur de code », c'est-à-dire qu'il offre la capacité de créer graphiquement des processus (répartis) de manipulation et de transformation de données puis de générer l'exécutable correspondant sous forme de programme Java (et scripts Pig).

Ce programme peut ensuite être déployé sur un ou plusieurs serveur(s) d'exécution.

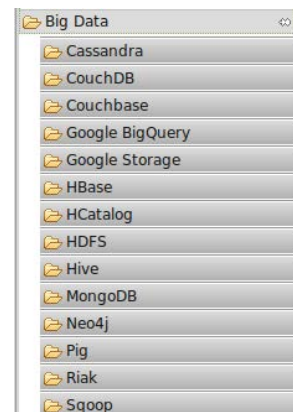
La modélisation des traitements se fait dans le Studio Talend, qui permet d'utiliser des connexions prédéfinies et les tâches de transformations pour collecter, transformer et charger les données par simple **glisser-déposer** dans l'espace de modélisation.

Palette de connecteurs Big Data

L'ETL Talend fournit nativement une large palette de connecteurs permettant de s'interfacer à la plupart des systèmes existants (bases de données, fichiers locaux ou distants, web services, annuaires,...).

Si l'ETL classique Talend peut se connecter aux principales bases NoSQL via des connecteurs communautaires ou APIs, la version Talend Open Studio for Big Data fournit nativement toute la flexibilité et les connecteurs d'intégration de masses de données, dont :

- les bases NoSQL : MongoDB, Apache Hadoop/Hive, Cassandra, Google BigQuery, Neo4j
- HDFS, HCatalog
- le chargement massif de bases NoSQL MongoDB et Cassandra ainsi qu'Apache Sqoop.



Composants de transformation

Les composants de transformation permettent entre autres :

- les multiplexages et jointures
- les filtrages (lignes, colonnes), le dédoublement
- l'exécution d'opérations sur des événements en base ou sur des fichiers
- les manipulations de fichiers locaux ou distants...

La liste des composants Talend est disponible à l'adresse suivante :

<http://www.talendforge.org/components/index.php>

La palette peut même être étendue grâce aux composants disponibles sur Talend Exchange ou du code Java spécifique.

Gestion des différents environnements d'exécution des traitements

L'ETL Talend gère des contextes d'exécution permettant d'externaliser l'ensemble des paramètres d'accès et variables d'exécution utilisés dans les composants / jobs.

Les utilisateurs peuvent ainsi configurer les paramètres à la volée lors de l'exécution ou utiliser des paramètres différents pour chaque contexte d'exécution : le développement, la recette et la production.

Intégration Hadoop

Génération de traitements répartis Pig :

Talend for BigData propose de générer des traitements (répartis) Hadoop Pig avec des composants graphiques prédéfinis.

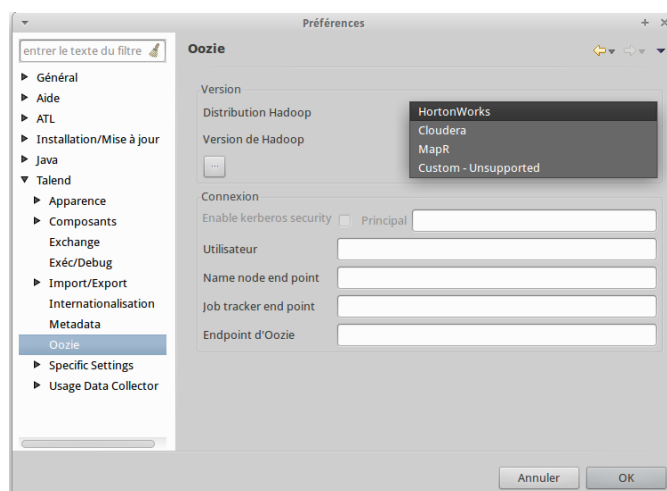
Il est également possible d'utiliser le mode ELT (Extract, Load and Transform) avec Hive pour répartir des traitements sur un cluster Hadoop.

De plus, le framework de traitement Hadoop YARN est intégré.

Coordination et intégration aux plateformes Hadoop :

Talend utilise Oozie pour la coordination des jobs sur un cluster Hadoop.

L'intégration est facilitée avec les outils des distributions HortonWorks, Cloudera et MapR :



Paramétrage de la connexion à Hadoop Oozie

“ Big Data Analyse et valorisation de masses de données ”

Talend Enterprise for Big Data

De manière analogue à Talend Enterprise for Data Integration pour l'ETL, cette version commerciale apporte notamment :

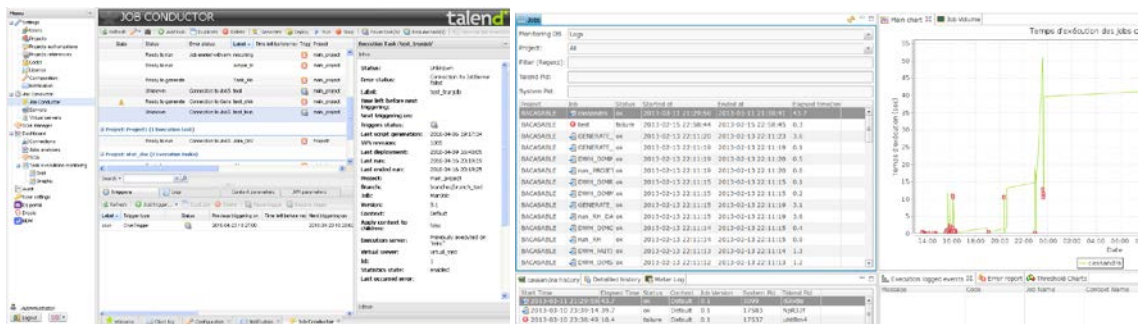
La gestion et le partage de métadonnées

- connexion aux bases des données (relationnelles ou NoSQL)
- connexion aux clusters Hadoop
- métadonnées de tables, fichiers,....
- analyse d'impacts.

La console Talend Administration Center

- gestion des référentiels des projets d'intégration, utilisateurs et droits associés
- ordonnancement des traitements (Job Conductor)
- console de monitoring AMC (Activity Monitoring Console) web
- gestion des reprises de traitements sur erreur d'exécution
- gestion des environnements d'exécution des traitements.

WWW.SMILE.FR



Job Conductor Talend - Activity Monitoring Console Talend

Autres fonctionnalités de productivités et d'exploitabilité

Cette version apporte également :

- le versionning des traitements
- la capacité de définir des points de reprise des traitements en cas d'erreur d'exécution
- un moteur de règles (Drools)
- joblets : morceaux de jobs réutilisables pour la factorisation des développements
- design de jobs à partir de templates
- visualisateur de données en sortie des composants
- change data capture

Jobs MapReduce

Cette version offre la possibilité de développer visuellement des traitements MapReduce, dont l'exécution peut se faire sur un cluster Hadoop.

L'exécution de jobs MapReduce depuis le studio offre un suivi d'avancement visuel de chaque étape map et reduce.

Talend Platform for Big Data

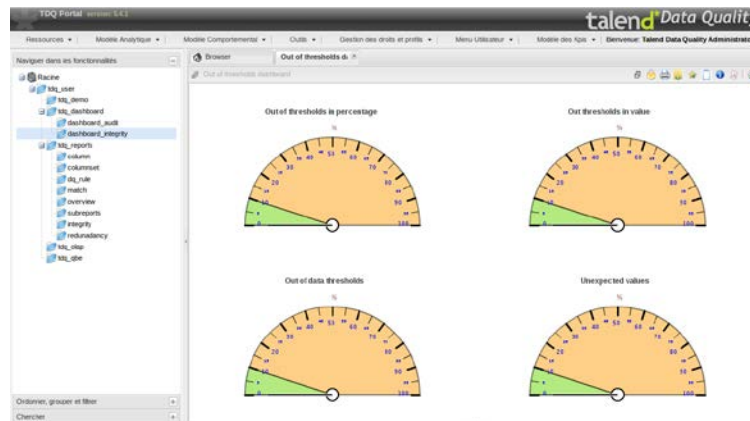
Cette version apporte notamment des fonctionnalités complémentaires et intégrées de qualité de données et de gestion de mapping complexes (XML, EDI) :

Profilage des données

Les analyses modélisées depuis le studio unifié, se font sur des sources, dont la définition peut être partagée avec les métadonnées définies au niveau de l'intégration.

L'outil produit des métriques sur le taux d'unicité, de remplissage, la conformité à un format, la diversité des formats ...

Des rapports, tableaux de bords et données requêtables peuvent être produits et publiés sur un portail décisionnel intégré (basé sur SpagoBI, présenté plus loin dans le document) à partir des analyses de données afin de piloter le processus d'amélioration de la qualité des données :



Composants de correction et enrichissement des données

Le studio de modélisation est enrichi de composants de traitement et correction supplémentaires de qualité des données :

- correction/enrichissement d'adresses postales via des services tiers QAS, Google
- rapprochements complexes en utilisant des technologies de logique floue
- création de tâches de correction manuelle des données.

“ Big Data Analyse et valorisation de masses de données ”

Workflow web de correction des données

La solution intègre la console web Data Stewardship avec la définition de workflows de correction et validation de données :


ID	Status	Type	Created By
1	resolved	manually	administrator
2	locked	Oracle	administrator
3	resolved	SAP_Excel_Oracle	administrator
4	new	TaskName	administrator
5	new	TaskName	administrator
6	new	TaskName	administrator

Liste des tâches de correction/validation de données

Column	Value	Customer	Finance
Name	Talent	Talent	Talent
Address	9, rue de pagès	9, rue page	9, rue de pagès
Phone	+33 1 46 25 06 00		+33 1 46 25 06 00
Zlrcode	92150	92150	92150

Détail d'une tâche de résolution de données

ETL Pentaho Data Integration

<p>Editeur : Pentaho Licence : Apache V2 et commerciale (suivant la version) Version actuelle : 5</p>	
---	---

Présentation

Editeur et solutions

Pentaho est un éditeur basé en Floride et en Californie, avec des bureaux en France. L'éditeur est un acteur impliqué de l'Open Source, qui a rallié dès le début des produits Open Source comme Kettle ou Mondrian et qui anime sa communauté.

Au delà de la solution d'intégration de données, Pentaho fournit aussi une solution complète d'analyse et d'exploitation décisionnelle des données : [Pentaho Business Analytics, présentés plus loin dans le document.](#)

Pentaho et le Big Data

Pentaho a établi des partenariats avec des acteurs majeurs du Big Data, notamment : MongoDB, HortonWorks, Cloudera, MapR et DataStax.

L'éditeur publie également un portail web dédié aux problématiques Big Data : <http://www.pentahobigdata.com>

Fonctionnalités

Pentaho Data Integration (PDI) est un ETL qui permet de concevoir et exécuter des opérations de manipulation et de transformation de données.

Grâce à un modèle graphique à base d'étapes, il est possible de créer dans le studio de modélisation (Spoon), sans programmation, des processus composés d'imports et d'exports de données, et de différentes opérations de transformation (conversions, jointures, application de filtres, ou même exécution de fonctions Javascript si besoin).

Les fonctionnalités ETL classiques de Pentaho Data Integration sont présentées plus en détail dans le livre blanc BI.

PDI Community Edition

L'ETL Pentaho Data Integration propose des connecteurs aux principales Bases NoSQL/Big Data telles Hadoop (HDFS, HBase, Hive et MapReduce), Cassandra, CouchDb, MongoDB,

“ Big Data Analyse et valorisation de masses de données ”

ElasticSearch ainsi qu’aux bases de données Amazon S3 et aux réseaux sociaux Twitter et Facebook.

Pour les traitements en masse, la connectivité avec Hadoop Map-Reduce et le moteur MongoDB Map-reduce sont intéressants, tout comme les capacités de répartition de charge des traitements ETL dans une configuration cluster de PDI.

En sus des composants et techniques dédiées aux technologies Big Data, il y a d'autres options de PDI qui permettent une meilleure gestion de grosses volumétries de données :

- lecture en parallèle de fichiers plats de grande taille tels que des fichiers de logs
- exécution concurrente de plusieurs copies d'une même étape d'une transformation avec distribution aléatoire en entrée des données en conséquence
- partitionnement, même option que la précédente avec une distribution plus intelligente des données à l'aide d'algorithmes proposés ou possibilité de développer des algorithmes de répartition spécifiques
- pour un environnement distribué, possibilité depuis la version 5.0 de faire du load balancing pour la distribution des données entre deux étapes d'une transformation.

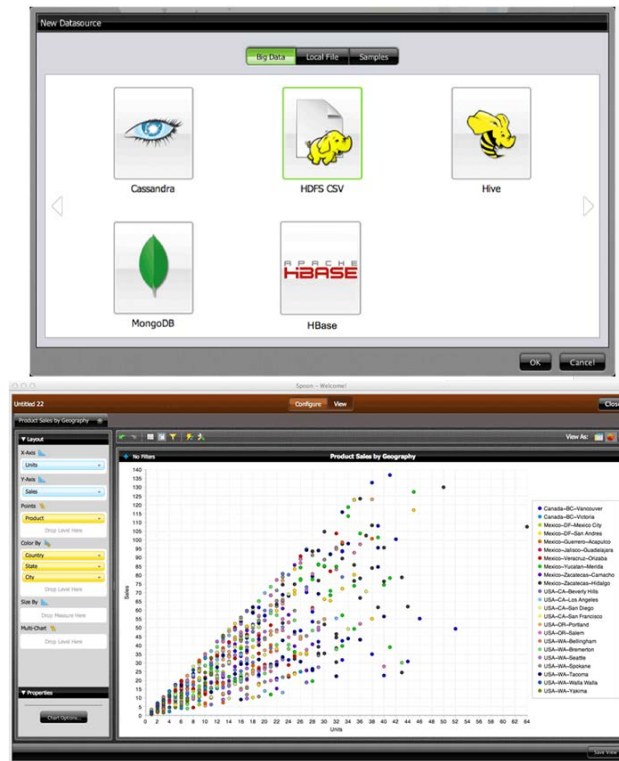
PDI Enterprise Edition

La version Enterprise apporte plusieurs outils pour plus de productivité dans la manipulation de données massives.

Les **possibilités de visualisation et d'analyse intégrées Instaview** sont utiles aux Data Scientists¹ pour développer rapidement des applications analytiques Big Data, en limitant les allers-retours entre outils :

¹ <http://blog.smile.fr/Pentaho-4-8-l-analyse-instantanee-et-interactive-des-donnees-mobiles-et-big-data>

“ Big Data Analyse et valorisation de masses de données ”



Perspective Instaview de Pentaho Data Integration Enterprise Edition

En effet, dans le cadre de la méthodologie AgileBI, cette perspective intégrée au studio de modélisation des traitements ETL permet d’analyser avec l’outil Analyzer Pentaho des données, Big Data ou non, issues des transformations et mises en cache dans une base MongoDB.

Fédération de données

La version Enterprise propose également des possibilités de fédération de données au travers d’un connecteur JDBC. Ce dernier permet de projeter une transformation PDI comme source de données JDBC : cela ouvre des perspectives intéressantes de connectivité et de restitutions en quasi temps réel sur des processus métiers distribués au niveau applicatif.

Cela permet également de faire une interface entre des technologies Big Data, NoSQL et certains outils de restitutions plutôt orientés SQL (workbench/Mondrian). Et ainsi, permet d’éviter dans certains cas une structure de stockage hybride (NoSQL / SQL).

Pentaho MapReduce

Pentaho MapReduce permet le développement de traitements MapReduce (mettant en oeuvre 1 transformation pour l’étape map et 1 transformation pour l’étape reduce) depuis le studio de modélisation des traitements ETL.

“ Big Data Analyse et valorisation de masses de données ”



Ils sont ensuite exécutables sur un cluster Hadoop.

Pentaho Predictive Analytics

En plus des méthodes d'analyse classiques (analyse d'événements passés et/ou présents), un des enjeux du Big Data notamment dans le domaine scientifique est de faire parler ces gros volumes de données pour de la prévision.

Weka est un projet data mining open source dont Pentaho est un acteur majeur, dans ce contexte de nombreux plugins sont disponibles par défaut ou non pour l'utilisation de certaines briques de Weka (Scoring, Knowledge Flow, ...) via Pentaho Data Integration.

Pour plus de précision sur les possibilités en termes de Data Mining via Pentaho, rendez-vous sur :

<http://wiki.pentaho.com/display/DATAMINING/Pentaho+Data+Mining+Community+Documentation>.

STOCKAGE DE MASSES DE DONNEES

Synthèse

Type	Solution	Site web de la solution
NoSQL Colonne	Apache Cassandra base de données répartie en Peer to Peer	http://cassandra.apache.org
NoSQL Colonne	Apache HBase Base de données du framework Hadoop Voir Hadoop pour sa description	http://hbase.apache.org
NoSQL Document	MongoDB	http://www.mongodb.org
NoSQL Document	ElasticSearch	http://www.elasticsearch.org
NoSQL Graph	Neo4j	http://www.neo4j.org

WWW.SMILE.FR

Fédération de données NoSQL dans des bases relationnelles

Plusieurs moteurs de bases de données relationnels permettent de fédérer des lacs de données massives NoSQL externes au sein de bases de données classiques.

Le modèle est ici d'utiliser un moteur de stockage NoSQL (réparti et qui reste accessible de manière autonome) au sein d'une base de données relationnelle pour son exploitation.

Citons par exemple le mécanisme de Foreign Data Wrapper de PostGreSQL ou le connecteur Cassandra de MariaDB.

Ces mécanismes offrent l'avantage d'intégrer facilement des données de bases NoSQL au sein d'un ODS ou un entrepôt de données de type base de données relationnelle et ainsi d'y accéder avec un langage SQL classique.

Par exemple, cela peut être une source de fait MongoDB à très forte volumétrie et vitesse intégrée de manière "transparente" à un ODS PostGreSQL.


“ Big Data Analyse et valorisation de masses de données ”



Par contre, il faut bien garder à l'esprit les limites de ce modèle :

- limites de performances techniques du moteur de la base de données fédératrice par rapport à un système de traitement réparti (agrégation de masses de données notamment)
- perte de performance due à l'intégration d'un système tiers
- mapping “rigide” des champs entre la base NoSQL et les tables virtuelles de la base de données fédératrice.

MongoDB

Type NoSQL : document Editeur : MongoDB Licences : GNU AGPL v3.0 et commerciale (suivant la version) Version actuelle : 2.4	
--	---



Présentation

MongoDB est une base de données NoSQL de type [document](#), la définition des données est très souple et chaque enregistrement a sa propre structure, dont les objets sont stockés au format JSON binaire (BSON).

Persistence

MongoDB permet de gérer la réplication et la répartition de données sur un ensemble de serveurs (cluster).

Connectivité, requêtage et traitement

L'avantage du format JSON est son utilisation native dans de nombreux langages de programmation, notamment le Javascript; la console MongoDB est d'ailleurs un interpréteur Javascript.

MongoDB fournit également des fonctions JavaScript de traitement réparti MongoDB Map-reduce.

Usages Big Data BI

MongoDB peut servir d'Operating Data Store.

Avec ses connecteurs disponibles au sein de la plupart des solutions BI OpenSource, MongoDB peut aussi servir d'entrepôt de données de masse à des fins de requêtage et de reporting.

L'analyse multidimensionnelle (OLAP) avec MongoDB nécessite actuellement l'emploi combiné d'un composant supplémentaire, tel :

- Hadoop Hive+Map-Reduce
- une fédération de données JDBC :
 - l'ETL Pentaho Data Integration avec son connecteur JDBC et du moteur Map-Reduce de MongoDB
 - Foreign Data Wrapper de PostgreSQL.

Conclusion

A l'heure où nous écrivons ces lignes, MongoDB est la base NoSQL la plus populaire d'après le site db-engines.com, bénéficiant d'une relative facilité de mise en oeuvre ainsi que d'un scope fonctionnel utile à l'entreposage opérationnel de masse de données.

ElasticSearch

Type NoSQL : document Editeur : ElasticSearch Licence : Apache V2 Version actuelle : 0.90	 elasticsearch.
--	---

Présentation

Sous le système de recherche d'ElasticSearch, propulsé par Apache Lucene, se cache un moteur de base de données NoSQL documentaire.

Persistence

ElasticSearch permet la mise en cluster pour la réplication et la répartition de données. A noter que les indexes (de recherche/requêtage) générés sont de type colonne.

Connectivité, requêtage et traitement

L'accès et la manipulation de données se fait simplement via l'API REST et le format JSON. Le moteur de requêtage propose des capacités d'agrégation et d'analyse, utile pour du requêtage décisionnel.

Usages Big Data BI

ElasticSearch peut servir d'Operating Data Store et à la mise en oeuvre de datamarts combinés avec des outils de restitution compatibles.

Conclusion


Cette solution est intéressante et prometteuse sur le plan technologique. Notons toutefois qu'elle est relativement jeune et encore peu intégrée aux portails décisionnels classiques, malgré une API très accessible.

**ANALYSER ET RESTITUER DES MASSES DE
DONNEES**

Synthèse

Type	Solution	Site web de la solution
Portail décisionnel complet	Pentaho Business Analytics	http://www.pentaho.fr
Portail décisionnel complet	JasperSoft BI Suite	http://www.jaspersoft.com/fr
Portail décisionnel complet	Spago BI	http://www.spagobi.org
Portail de tableaux de bord web	ElasticSearch Kibana	http://www.elasticsearch.org/overview/kibana
Portail décisionnel complet	Vanilla Platform	http://bpm-conseil.com

Pentaho Business Analytics

<p>Editeur : Pentaho Licences : GNU GPL V2 et commerciale (suivant la version) Version actuelle : 5</p>	
---	---

Présentation

Pentaho Business Analytics est un portail décisionnel qui permet la distribution d'outils d'analyse et requêtage en langage naturel ainsi que des documents décisionnels à un grand nombre de personnes par l'intermédiaire d'une interface web :

WWW.SMILE.FR



Page d'accueil de Pentaho Business Analytics

Pentaho est proposé en version communautaire et en version entreprise soumise à souscription annuelle, avec des modules supplémentaires (Pentaho Analyzer) ainsi qu'un support produit.

La communauté enrichit le portail en version communautaire sous forme de modules disponibles depuis le Pentaho MarketPlace, parmi lesquels l'interface d'analyse Saiku et les CTools qui ont le vent en poupe.

Fonctionnalités

Pentaho fournit un portail décisionnel complet, permettant aux utilisateurs finaux :

- l'analyse multidimensionnelle : Pentaho Analyzer, Saiku Analytics
- le requêtage ad-hoc : Interactive Report, Saiku Reporting, WAQR
- l'exploitation de tableaux de bords dynamiques (CTools).

Les capacités de répartition de charge (load balancing) entre plusieurs instances Pentaho Business Analytics sont intéressantes dans le cadre d'analyses en masses.

“ Big Data Analyse et valorisation de masses de données ”



Connectivité NoSQL et exploitation de données massives

Pentaho fournit nativement des connecteurs Big Data au niveau des connections du portail pour les sources NoSQL offrant une connectivité JDBC :

- Hive
- Impala
- connecteur JDBC générique.

A noter qu'il est également possible d'accéder à d'autres sources de données NoSQL au sein du portail en passant par de la fédération de données, en utilisant l'[ETL PDI](#) ou un [mécanisme de stockage externe d'une base relationnelle](#).

L'outil Pentaho Report Designer permet de plus d'élaborer et de publier des rapports à partir d'une source MongoDB.

JasperSoft BI Suite

Editeur : JasperSoft
Licences : GPL et commerciale (suivant la version)
Version actuelle : 5.5



Présentation

JasperSoft BI Suite est la plateforme décisionnelle de JasperSoft, société qui développe également le générateur d'états JasperReports, disponible depuis 2001.

La plateforme propose des fonctionnalités de reporting et d'analyse et est disponible sous deux licences : GPL et commerciale.

Fonctionnalités

JasperServer, en versions Professionnelle et Entreprise, offre des fonctionnalités supplémentaires par rapport à la version open source, limitée à la publication et la diffusion de rapports :

- outil de création de rapports ad-hoc en ligne (listes, graphiques ou tableaux croisés), accessible à tout utilisateur
- outil de composition de tableaux de bord.

WWW.SMILE.FR

The screenshot shows the JasperServer interface with a pivot table titled "Product Results by Store Type". The table displays sales data for various products across three store types: Deluxe Supermarket, Mid-Size Grocery, and Totaux. The data is organized by product name and includes a "Measures" column for sales figures.

Type de magasin	Deluxe Supermarket	Mid-Size Grocery	Totaux
Measures	Ventes du magasin 1998	Ventes du magasin 1998	Ventes du magasin 1998
ADJ Rosy Sunglasses	82,80	22,60	110,40
Akron City Map	52,20	13,92	66,12
Akron Eyeglass Screwdriver	29,92	12,32	42,24
American Beef Bologna	17,94	8,58	26,52
American Chicken Hot Dogs	55,14	7,56	62,70
American Cole Slaw	37,38	6,23	43,61
American Corned Beef	105,00	0,00	105,00
American Foot-Long Hot Dogs	29,86	8,56	38,42
American Low Fat Bologna	71,73	17,22	88,92
American Low Fat Cole Slaw	34,03	11,33	45,40
American Pimento Loaf	66,24	16,56	82,80
American Potato Salad	22,90	6,20	29,10
American Roasted Chicken	109,38	12,82	118,80
American Sliced Chicken	12,38	1,22	14,16
American Sliced Ham	49,68	44,16	93,84
American Sliced Turkey	44,38	22,18	66,52
American Turkey Hot Dogs	104,12	16,44	120,56

Module de requête ad-hoc de JasperServer

“ Big Data Analyse et valorisation de masses de données ”

Connectivité NoSQL et exploitation de données massives

JasperSoft BI fournit nativement, en versions commerciales Professionnel et Entreprise, un outil de requêtage et d'analyse ad-hoc qui permet une exploitation directe de sources de données NoSQL :

- MongoDB
- Hadoop via Hive

Un système de cache de données est présent, pour optimiser le temps de réponse des requêtes.

JasperSoft Studio fournit également une large palette de connecteurs au delà du JDBC classique pour le reporting et les tableaux de bord :

- MongoDB
- Hadoop via Hive
- Cassandra
- JSON.

Il existe aussi des connecteurs communautaires pour d'autres bases NoSQL, comme Google BigQuery ou Neo4j.

SpagoBI

Editeur : Engineering Group / OW2 Consortium
Licence : Mozilla Public License V2
Version actuelle : 4.1



Présentation

SpagoBI est une suite décisionnelle uniquement distribuée sous licence Open Source, développée par la société italienne Engineering Ingegneria Informatica au sein du consortium OW2.

Fonctionnalités

Afin de couvrir les différents besoins fonctionnels propres à la valorisation et l'analyse de données, SpagoBI propose une vingtaine de modules (ou « moteurs ») complémentaires, offrant des fonctionnalités de reporting/dashboarding, requêtage et analyse OLAP ad-hoc, geoBI, KPI et datamining :



Exemples de restitutions SpagoBI

Ces modules s'appuient sur un ensemble de projets Open Source phares, offrant ainsi une richesse de modules fonctionnels unique : l'ETL Talend, le moteur MOLAP Palo, les moteurs de reporting BIRT et Jasper, R et weka datamining.

“ Big Data Analyse et valorisation de masses de données ”



Modules de SpagoBI

WWW.SMILE.FR

Connectivité NoSQL et exploitation de données massives

Afin de répondre à la problématique du Big Data, SpagoBI développé de nouveaux connecteurs permettant le requêtage de bases de données NoSQL via des datasets :

- HBase: développement de requête HBQL, langage de requête Hbase, intégré nativement dans SpagoBI
- Hive: développement de requête HQL, langage de requête Hive, intégré nativement dans SpagoBI
- Impala: connecteur Cloudera Impala JDBC, récemment rendu disponible par Cloudera
- Cassandra: développement de requêtes CQL, langage de requête Cassandra

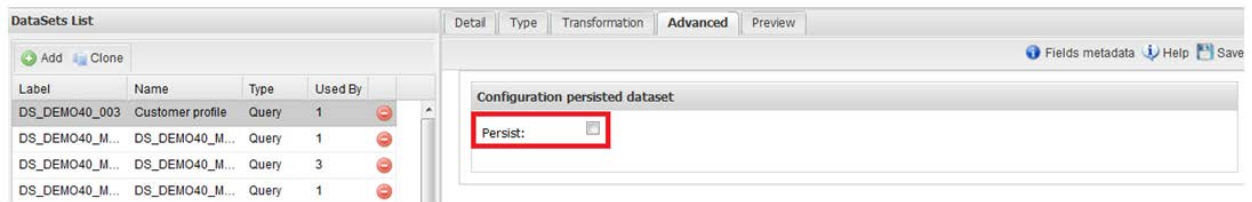
Label:	BIGDATA_DATASOURCE
Description:	BIGDATA_DATASOURCE
Dialect:	Default Dialect
Multischema:	Default Dialect
Read Only:	Oracle(any version) Oracle (Oracle 9i/10g) SQL Server
Write Default:	HSQL
Type:	MySql PostgreSQL
URL:	Ingres
User:	HBase QL Hive QL
Password:	DB2 AS400
Driver:	

Sélection du langage d'un connecteur

“ Big Data Analyse et valorisation de masses de données ”

Dans la version 4 de SpagoBI, la définition de dataset a évolué afin de permettre des temps de réponses plus courts sur les larges volumes de données :

- possibilité de planifier l'alimentation des datasets pour une restitution différée
- possibilité de définir des datasets persistants où les données sont stockées en cache.



Définition d'un dataset persistant

SpagoBI travaille actuellement à introduire les problématiques d'accès en temps réel ainsi qu'à la mise en place d'une couche sémantique sur les données Big Data.

ElasticSearch Kibana

Editeur : ElasticSearch
Licence : Apache V2
Version actuelle : 3m4



Présentation

Kibana est le module de dashboard d'ElasticSearch. Il permet d'associer la puissance du moteur de recherche d'ElasticSearch (des recherches complexes peuvent être faites pour filtrer les données pertinentes à l'analyse) aux modules de reporting classiques.

Cette solution est jeune : la première publication sur github date de début 2013. Toutefois, l'éditeur ElasticSearch propose un service de support en production pour ce composant. L'interface est entièrement écrite en javascript, avec les frameworks angular.js, bootstrap et jquery notamment. Un simple serveur web suffit donc à déployer la solution.

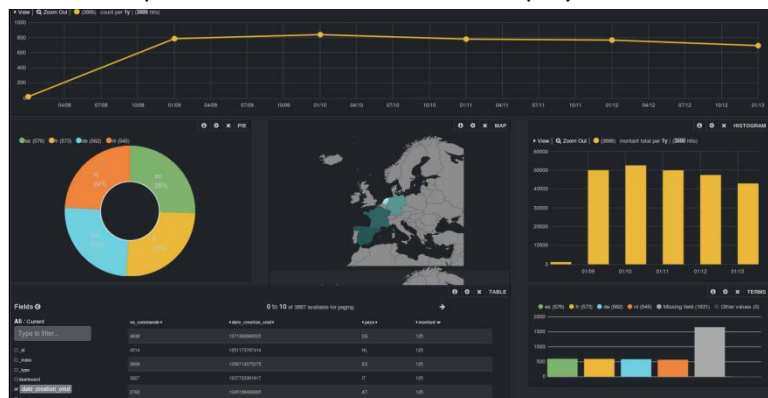


Tableau de bord Kibana

Fonctionnalités

L'usage unique de Kibana est la publication de tableaux de bords visuels, souples, hautement paramétrables par l'utilisateur final, grâce aux fonctionnalités de recherche et de filtrage offertes par ElasticSearch.

L'outil propose un rafraîchissement automatique, adapté à des problématiques de monitoring de processus en temps quasi réel.

Le design des tableaux de bord se fait via l'insertion de panels (graphiques, listes, tendances, cartographies,...) dans une structure de type tableau. Un tableau de bord peut ainsi être bâti en quelques minutes. Les panels communiquent entre eux : recherche, zoom,...

Notons toutefois que cette solution, jeune, ne permet pas encore de mise en forme complexe et les composants de restitution intégrables sont en nombre limité.

“ Big Data Analyse et valorisation de masses de données ”



Les tableaux de bord peuvent être enregistrés dans une base ElasticSearch afin d’être ré-exécutés et partagés.

L’accès à Kibana peut être protégé (authentification au niveau du virtualhost d’Apache).

Par contre cette solution ne permet pas encore de gérer complètement une bibliothèque de tableaux de bords (arborescence de tableaux de bords, droits d’accès aux tableaux de bord).

REMERCIEMENTS

Un grand remerciement à toutes les personnes ayant travaillé sur le livre blanc :

- Florent BERANGER,
- Elise BENZAGLOU,
- Laury GIRONDIN,
- Aurélien FOUCRET,
- Adrien FUTSCHIK,
- Pierre-Antoine MARC.

N'hésitez pas à nous transmettre vos avis et évaluations sur ce livre blanc.
Une seule adresse : contact@smile.fr