Big Data &Al Point Of View

Christophe Burgaud - PhD

Novembre, 2020



AGENDA

- ·L'enjeu de la transformation digitate chez nos clients
- · Retour sur les cas d'usage
- Concepts et definitions
- Marche, Technologies et Architecture



Transformation digitale et cas d'usage



The digital age has changed the way we Live, Play, Learn and Work...





Netflix provides personalized recommendations

Waze provides a personalized driving experience

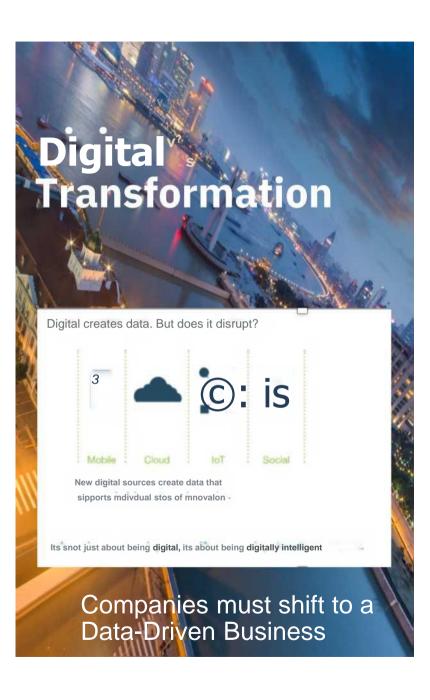
Uber delivers food that you like and is the right temperature

Self driving cars react to changing conditions

But ALL based on DATA and AI

Some Very successfully

Some not as successful



75%

of large enterprises will have digital transformation at the center of corporate strategy within two years

81%

of companies do not yet understand the data required for AI

>85%

of enterprise IT organizations will commit to multi-cloud architectures

La data au creur de la transformation digitale

- Laurent Mignon, President de la BPCE (Banque Populaire Caisse d'Epargne), 29 janvier 2019
 - « Tout depend de la maniere dont vous allez utiliser la donnee » reprend le dirigeant.
 - « La technologie vous pouvez l'appeler Machine Learning, Deep Learning, intelligence artificielle, etc. le veritable probleme est comment bien utiliser la donnee dont nous disposons pour ameliorer notre capacite a mieux servir nos clients. C'est le premier enjeu » insiste-t-il.
 - « Si nous reussissons cela, je pense que nous transformerons en profondeur notre entreprise. De maniere implicite, l'usage de la donnee etait realise dans la tete de nos conseillers en agence. Nous devons elargir l'acces a la donnee » dit-il.

http://www.larevuedudigital.com/la-banque-bpce-en-attente-de-son-data-lake-pour-passer-a-la-vitesse-superieure/

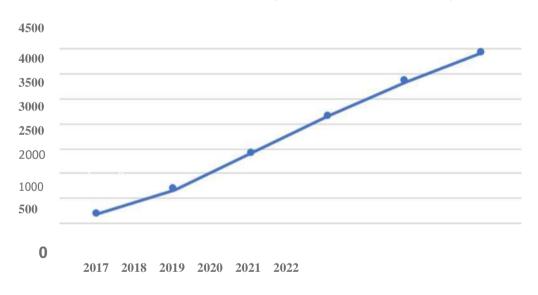


Business Value of Data & Al

Some main sources of AI business value:

S Customer Experience S Cost Reduction S New Revenue/New Services

Al-Derived Business Value (Billions of U.S. Dollars)



Gartner, April 2018 (https://www.gartner.com/newsroom/id/3872933)



Hyper or Radical Personalization





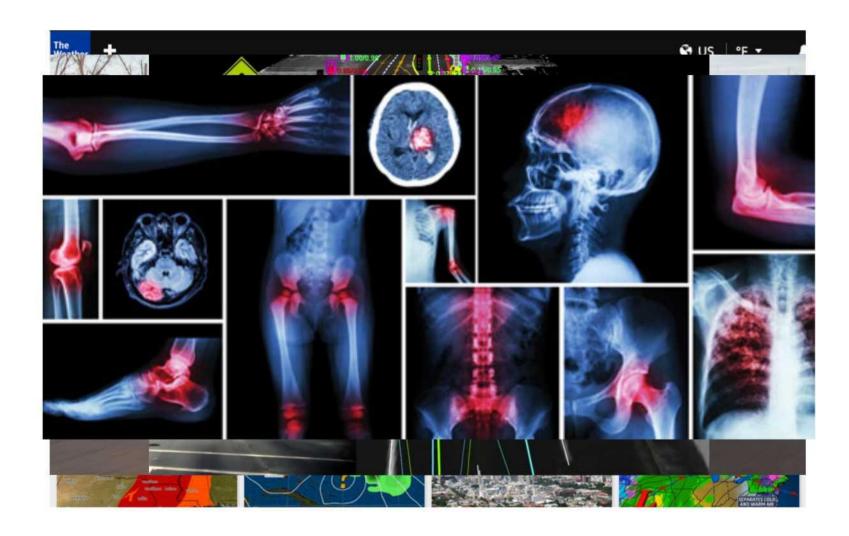
Resource Allocation & Strategic Planning





Predictions and Classifications







Industries: All

Smarter Healthcare



























Combine

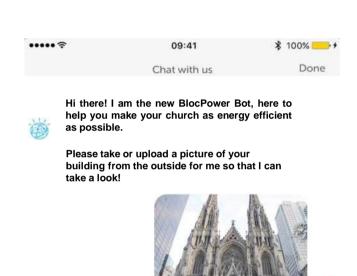
Building Model to Predict Energy Consumption of Buildings.





Combine

Chat Bot to estimate energy cost from an image of building.



Great, thanks for that picture! Looks like your building is made of stone and has large windows

I estimate your building has a high energy usage intensity (EUI), with a 97.01% probability

Great, thanks for that picture! Looks like your building is made of stone and has large windows.



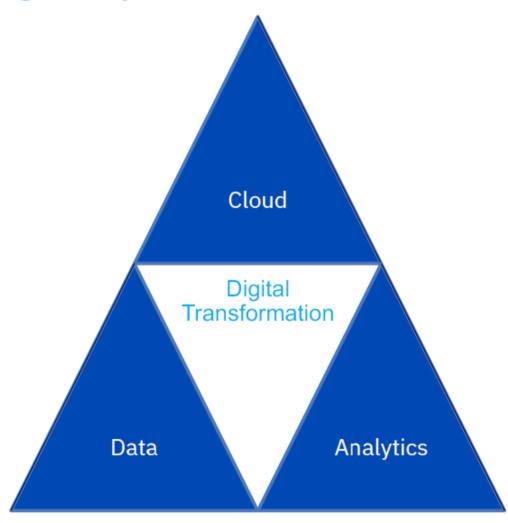
How old is this church building?



Buzzwords and challenges

- Al, machine learning, cloud, self-service, data governance, etc...there is no shortage of buzzwords in data today.
- Every organization is seeking to outpace their competition by leveraging data to drive differentiation for their business.
- To win this race, companies are building up data science teams, investing in faster/more scalable cloud data platforms and utilizing the growing variety of publicly available datasets and algorithms.
- How do you stay ahead of what's next and help drive the successful adoption of new technology and processes within your organization?

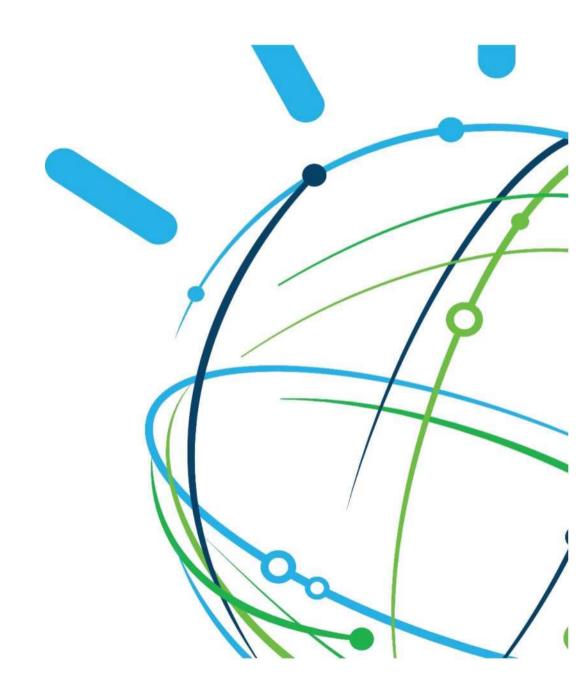
Transformation digitale: points clés



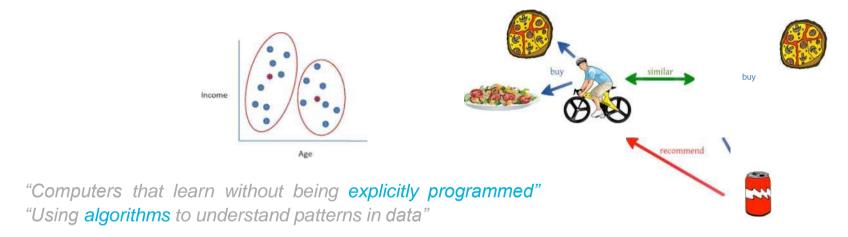
Data and AI Forum / © 2019 IBM Corporation 15

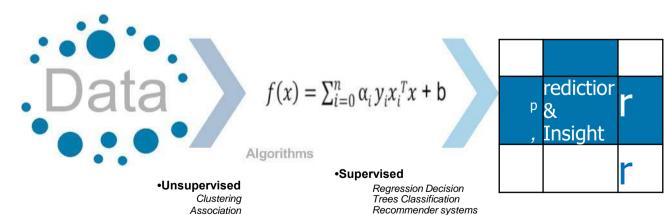


Concepts et definitions



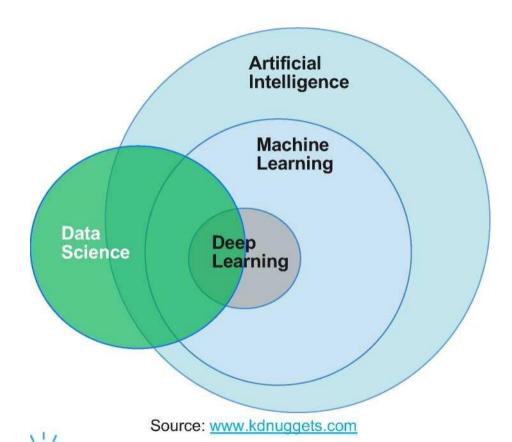
Machine Learning... What is it?







Artificial Intelligence



All Machine Learning is Al but not all Al is Machine Learning.

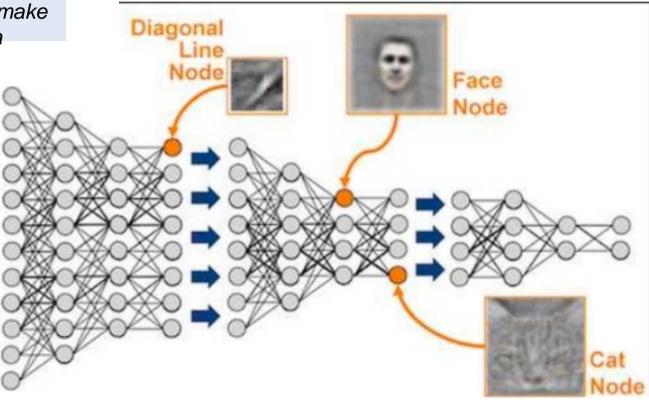
It's the fastest growing part of Al and shows the most promise.

Deep Learning is a subset of Machine Learning and focuses even more narrowly on a subset of Machine Learning techniques and requires "thought".

Deep Learning Example

Technical Definition:

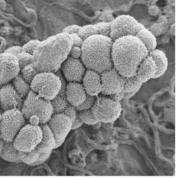
Deep learning is a class of machine learning algorithms in the form of a neural network that uses a cascade of layers (tiers) of processing units to extract features from data and make predictive guesses about new data





Deep Learning Use Cases











INTERNET & CLOUD

Image Classification
Speech Recognition
Language Translation
Language Processing
Sentiment Analysis
Recommendation

MEDICINE & BIOLOGY

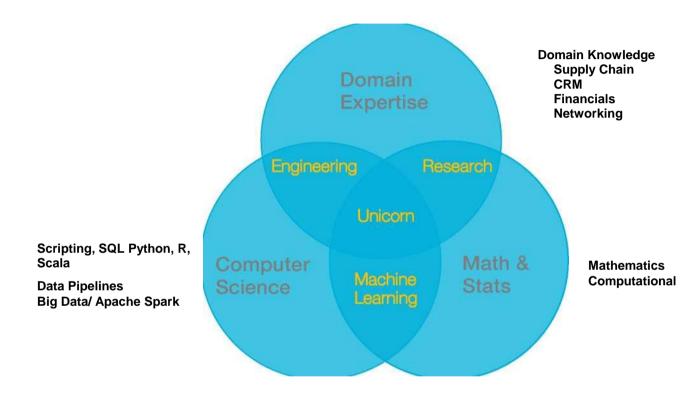
Cancer Cell Detection Diabetic Grading Drug Discovery MEDIA & ENTERTAINMENT

Video Captioning Video Search Real Time Translation SECURITY & DEFENSE

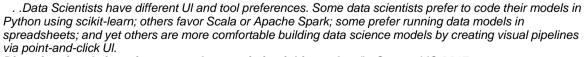
Face Detection Video Surveillance Satellite Imagery AUTONOMOUS MACHINES

Pedestrian Detection Lane Tracking Recognize Traffic Sign

What makes a "Data Scientist"?



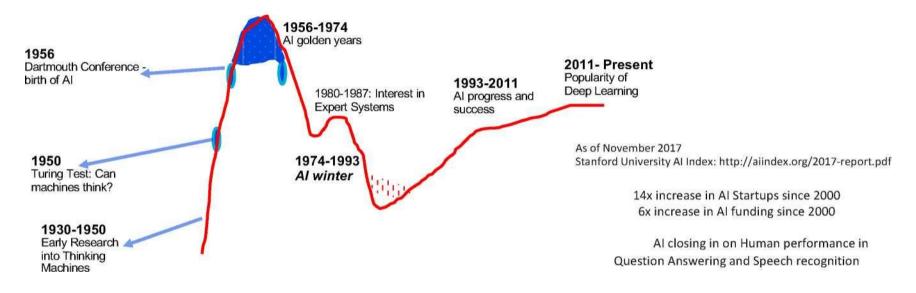
Data Science Projects Require multiple Skills



Diversity of tools is an important characteristic of this market. " - Gartner MQ 2017



Al and Machine Learning have been around for some time, so why the hype now?



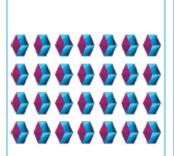
Big Data: Structured and unstructured data, as well as its storage, has grown exponentially decade. Data is also more accessible.

Distributed processing: The Machine Learning algorithms and models take advantage of new technologies: Spark, GPUs and Big Data to train models with large data sets. The new challenge is to industrialize and democratize AI in applications



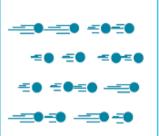


- Fit-for-purpose data architectures are necessary to accommodate the specialized data needs of the business, and individuals within its organizations.
- Gone are the days of a single, structured, data-at-rest architecture.



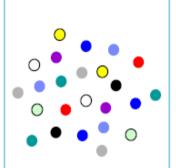
Data at rest

Terabytes to Zettabytes of data to process



Data in motion

Event-based data, streaming data, milliseconds to seconds to respond



Data in many forms

Structured, unstructured, text, documents, images, speech, video



Data in doubt

Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations

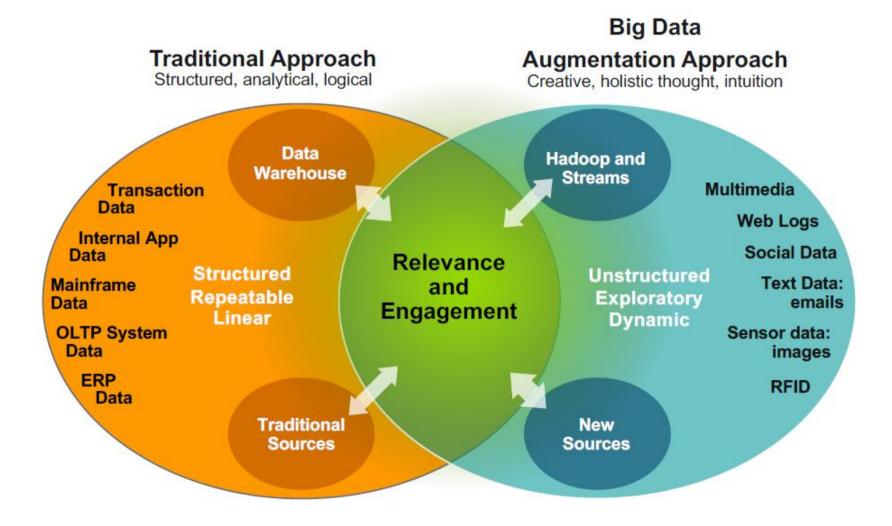


Une Velle Véfinition du Vig Vata... les 11 V

"Vast, Volumes of Vigorously, Verified, Vexingly Variable Verbose yet Valuable Visualized high Velocity 23 Data"... I just would add "... creating Value" :-) Décideo



Big Data Approaches Help Accumulate Context







Market and Technology



Marketplace Observations and hot topics

■ "Big Data" is not any more Hadoop technology only: architectures rely more and more on Object Storage

"Big Data 2018: Cloud storage becomes the de facto data lake"



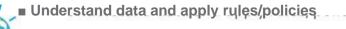


Spark is the layer for massive distributed processing



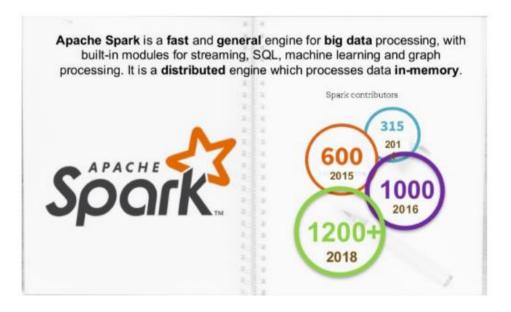
- Data Science and Machine Leaning industrialization
 - Gain insight from more data exploration on more data (volumes, internal, external,...)
 - Challenge is to deploy and monitor predictive and prescriptive modeling
- Requirement to interoperate with Enterprise tools/application that understand SQL:
 - BI tools like Cognos, Tableau, QlikView, ...
- Data Governance is a hot topic
 - **Sensitive Data (RGPD)**
 - Security





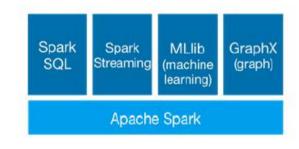
What is Spark?





Spark is an in-memory framework

Supports general workloads as well as streaming, interactive queries and machine learning providing performance gains



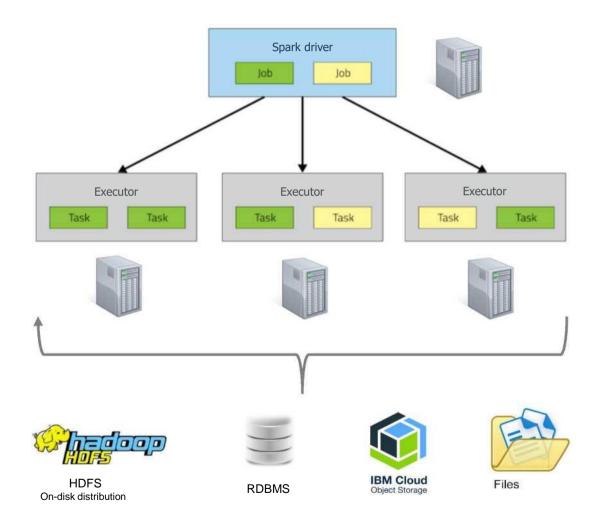


Why Spark?





Spark DataFrames In-memory distribution





Key Reasons for the Interest in Spark

Performant



In-memory architecture greatly reduces disk I/O

Productive



Concise and expressive syntax, especially compared to prior approaches

Single programming model across a range of use cases and steps in data lifecycle

Integrated with common programming languages - Java, Python, Scala, R

New tools continually reduce skill barrier for access (e.g. SQL for analysts)

Leverages existing investments



Works well within existing Hadoop ecosystem

Improves with age



Large and growing community of contributors continuously improve full analytics stack and extend capabilities



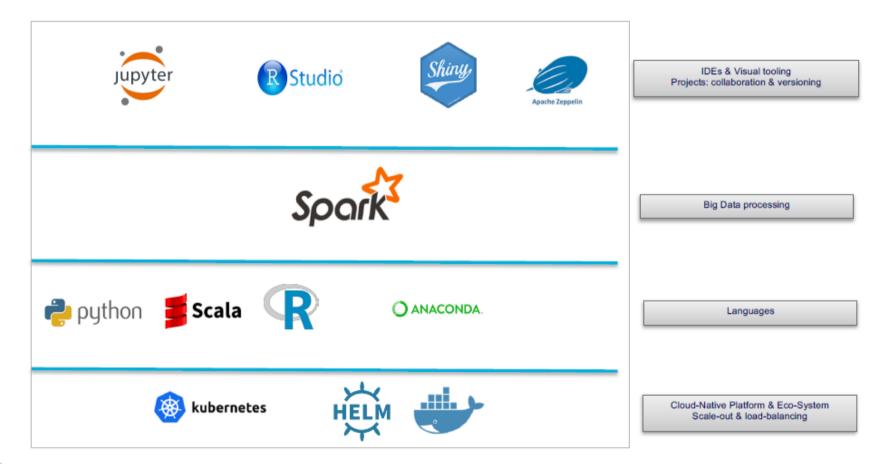
What Spark Is Not!

- Not only for Hadoop Spark can work within Hadoop (especially with HDFS), but Spark is a standalone system
- Not a data store Spark attaches to other data stores but does not provide its own
- Not only for machine learning Spark includes machine learning and does it very well, but it can handle much broader tasks equally well
 - Distributed File System
 - Data preparation
 - Streaming processing
 - SQL engine
 - Graph Engine
 - Distributed R
- Not a language!!!

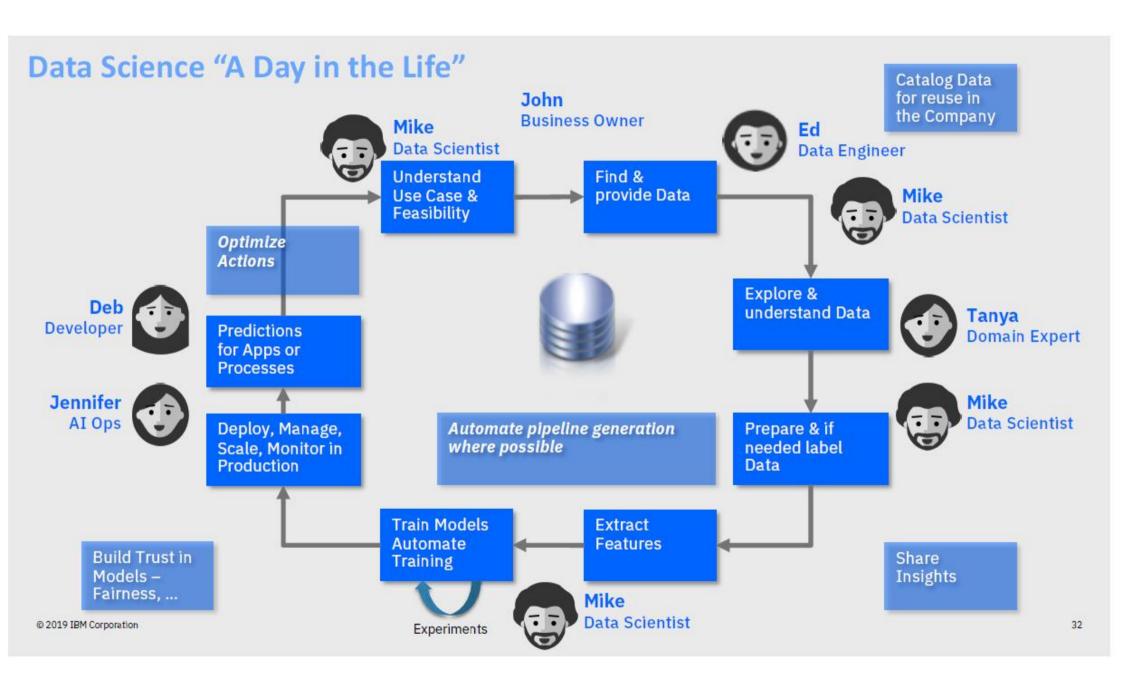




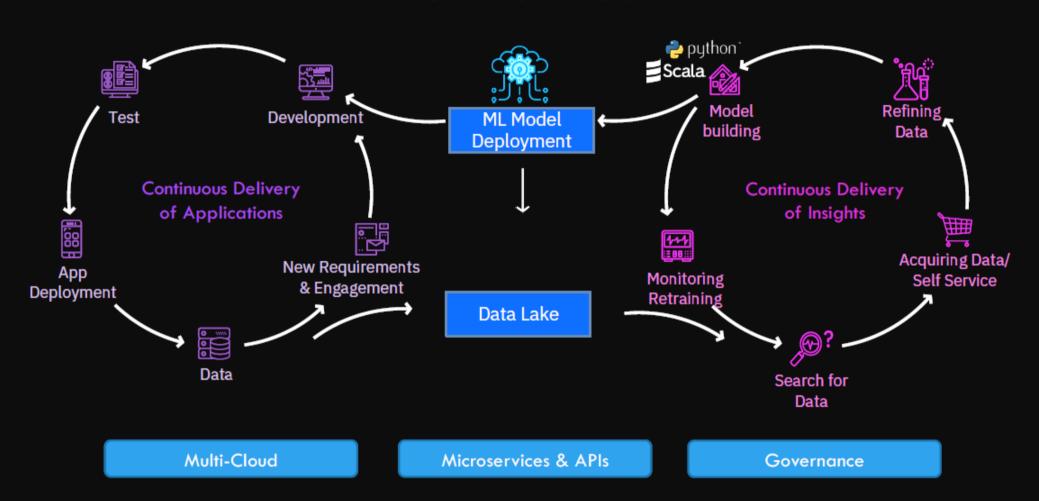
Data Science and Big Data technologies convergence







Delivering Insight Applications



Organizations lack the resources to manage and monitor AI in production.

60%

of companies see **regulatory constraints** as a barrier to implementing AI.

IBM IBV AI 2018

63%

cite availability of **technical skills** as a challenge to implementation.

- IBM IBV AI 2018

Without expensive Data Science resources handholding multiple AI models in a production application:

- No way to validate if AI models are compliant with regulations and will achieve expected business outcomes before deploying
- Difficult to track and measure indicators of business success in production
- Resource intensive and unreliable processes for ongoing business monitoring and compliance
- Impossible for business users to feedback subtle domain knowledge into model lifecycle

Data & Al Work Together

Data makes AI feasible & necessary.





Al makes Data valuable.

The Data Challenge

- **0** Which data should we use to train ML models?
- **0** Do we have the right data to train an ML model? The data that are easy to get may not be the most informative.
- **0** Are they available and easy to access? What about the volumes?
- **0** Are they structured? Unstructured?
- 0 Do we have labeled data and are those labels valid?
- **0** How do we make sure that data can be used without violating any compliance and regulation requirements?

The Al challenge

- **0** Understand the business use case
- **0** Choosing the right Machine Learning Algorithm

0 Supervised Learning 0

Unsupervised Learning

0 Industrialize it in Applications

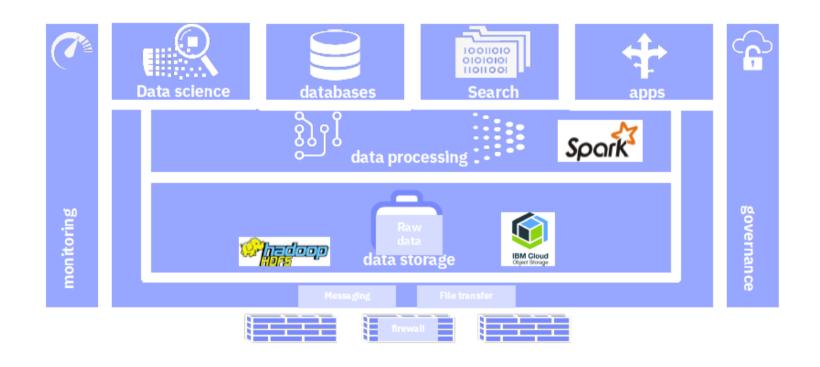
35



How Smart Is AI ? Can I Use It For My Problem Domain ?
IT'S ONLY AS GOOD AS THE DATA



Data Lake Architecture

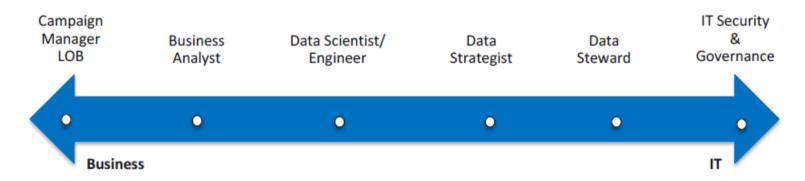






Who benefits from a Data Lake?





- LOB users, business analysts and data scientists can easily find the information they need without extensive IT involvement
- Data strategists and data stewards can make information available to users in an organized and well-governed manner
- IT security and governance teams can be assured that information is governed according to well-defined organizational and regulatory policies





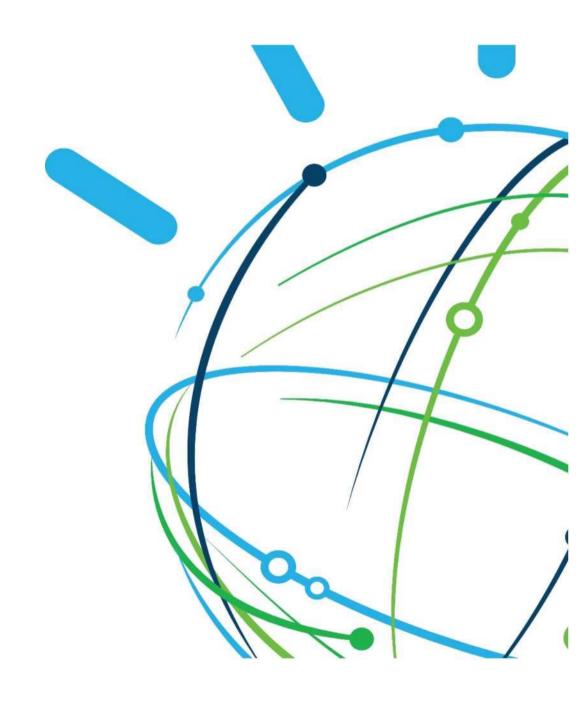
Quelques principes d'architecture

- Governance:
 - **◆◆◆**Share the same understanding on data: Data Catalog
- User Experience :
 - **♦**Collaborative : Share the assets (data, modeles, dashboards, flows, ...)
 - Managed Services
- Deploy Anywhere:
 - Cloud native architecture
 - Cloud agnostic & multi-cloud any vendor cloud or data center
 - Coherent, efficient, and scalable data & analytics services





Skills & Tips



5 tips for machine learning success

Contribute to Open Source

- Use Github and learn from some tutorials.
- Scan Apache's open source projects and choose one to explore. Play around with Jupyter notebooks, and reach out to other contributors. Then, when contributing to open source becomes more importantin your own dev community, you'll be ready.

Participate to Challenges and stay active in the community

- Kaggle, Hackathon

Keep learning Python

If you're already using Python for machine learning, you might be frustrated by the lack of support, especially for large data sets. But many data scientists love Python, so more open source projects are improving their support. You'll be able to do more as time goes on.

Pay attention to data gravity

- It almost always makes sense to move computation to the data rather than moving data to the computation. Same for tools. Bring tools to the data instead of bringing data to the tools.
- Simple advice, but it's hard to follow when you need to set up a data management architecture quickly or cheaply. If the people managing your data have security or data integrity concerns, spendtime making your case. In the long term, it makes a huge difference, especially with machine learning, where computation is getting more and more intense.

Know your data:

- Maybe this seems obvious, but check the data quality; You'll probably find and correct the obvious ones, but don't get overconfident.
- Data quality issues are everywhere.

Get started with Cloud

- · Sooner or later, you'll probably need to work with the cloud, so plan in advance.
- · Watson Studio is a good example to do it





Do not Stop Educating







https://www.kdnuggets.com/



https://www.reddit.eom/r/MachineLearning/





https://cognitiveclass.ai/

Free

courses, free Spaterials: MOOCs on Big Data, Data Science; Spark, Hadoop, R, SQL, NoSQL

Cloud-based sandbox for exercises





Suivre IBM: jn # O ^ @ D

Rechercher parmi les offres juniors

Rechercher parmi les offres experiments

