

dr Filip Wójcik
Wroclaw University
of Economics and Business

Extreme Gradient Boosting

H2O POWERED MACHINE LEARNING

FILIP WÓJCIK

SENIOR DATA SCIENTIST AT CREDIT SUISSE, UNIVERSITY LECTURER

FILIP.WOJCIK@OUTLOOK.COM

Agenda

1. Machine learning madness
 - Core machine learning technologies
 - Machine learning overview
2. Is deep learning an ultimate algorithm (hint: no!)?
 - Neural networks vs the rest of the world
 - Top machine learning algorithms
3. Xgboost overview
 - Decision trees primer
 - Boosted trees - powerful extension to decision trees
4. H2O framework – scalable machine learning tool
 - H2O architecture
 - H2O and big data
 - DEMO!

Machine learning madness



Machine learning madness

- Data science & machine learning are gaining popularity in recent years
- Rapid boost of data science & analytical software
- Big data influence – larger volumes of data can be processed now, without any problems
- Machine learning everywhere!



Hi. I'm Cortana.



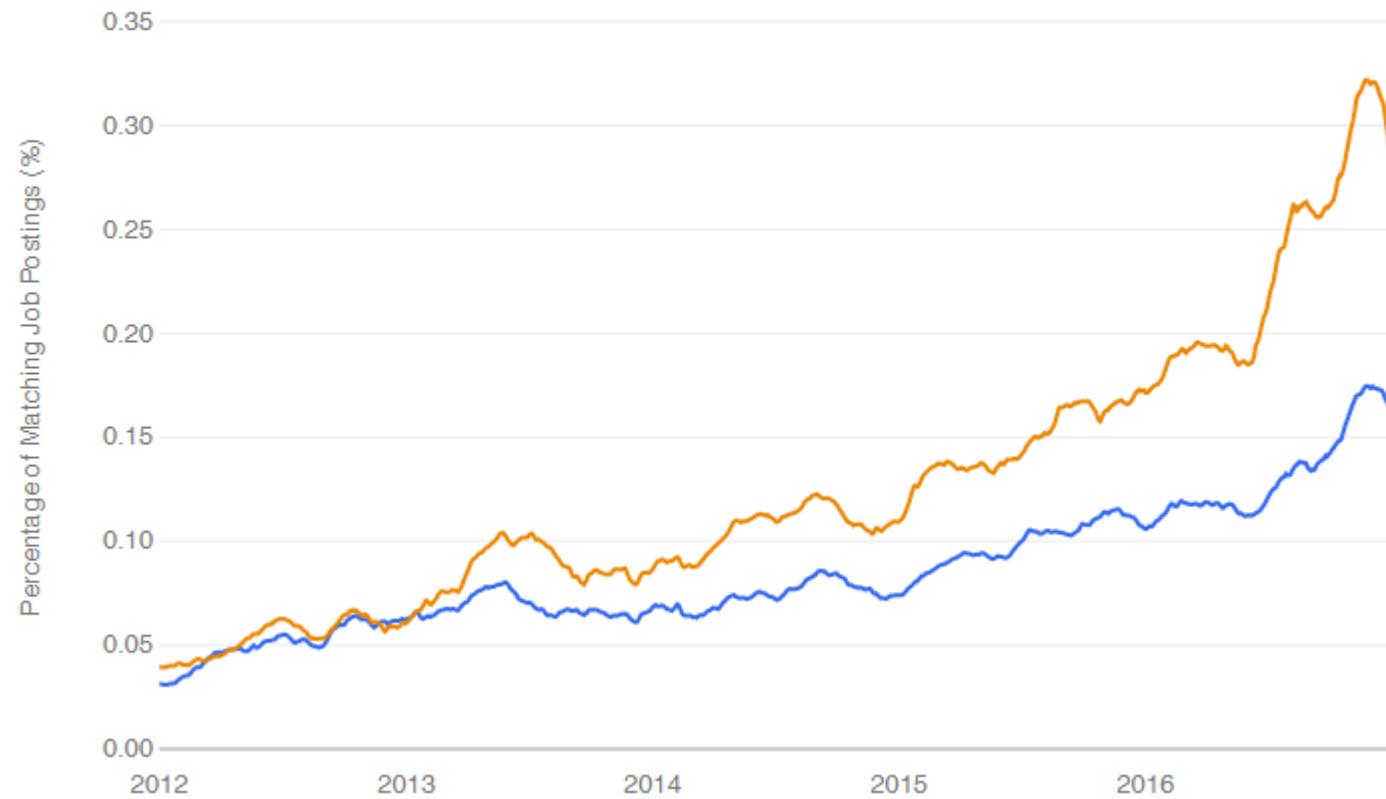
Machine learning madness

Machine learning job positions
count

R jobs



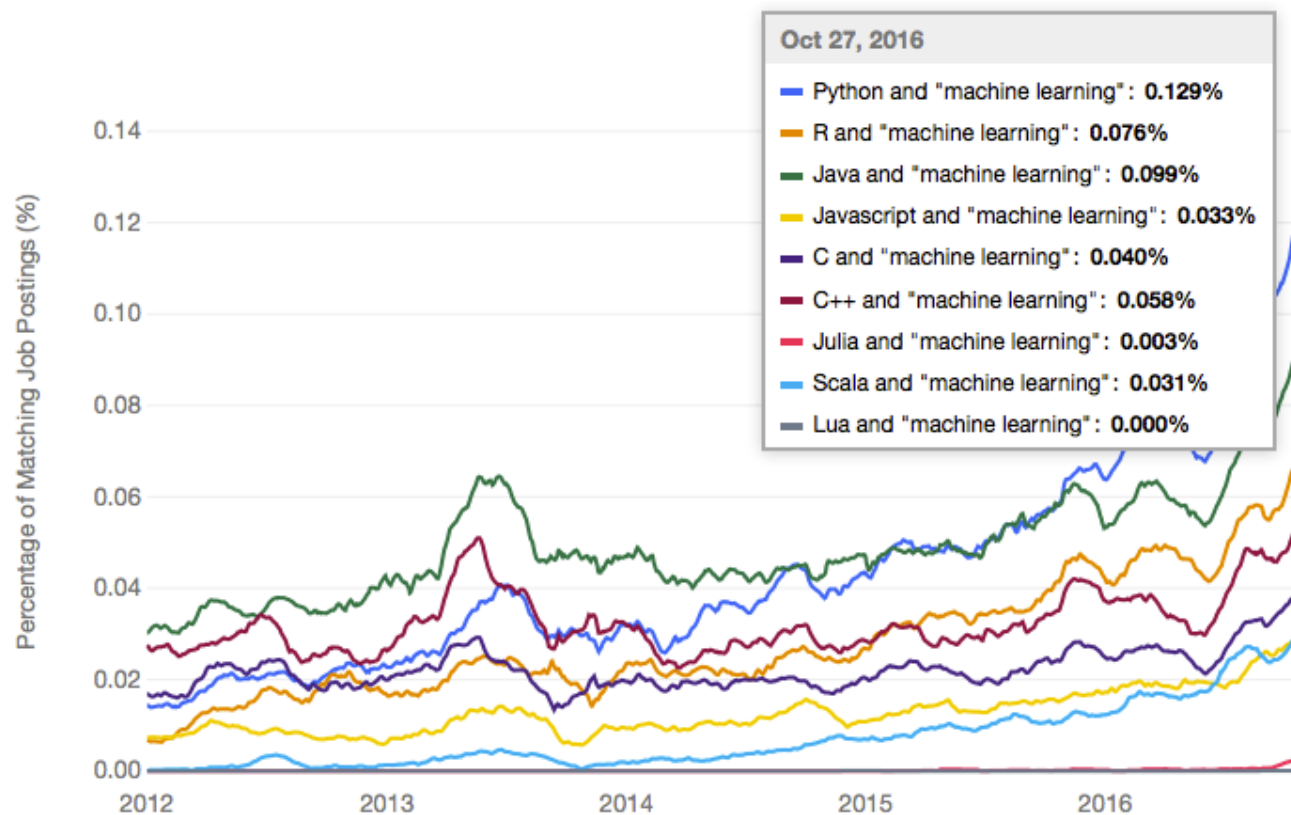
Python jobs



Source: r4stat, The Popularity of Data Science Software by Robert A. Muenchen, <http://r4stats.com/articles/popularity/>

Machine learning madness

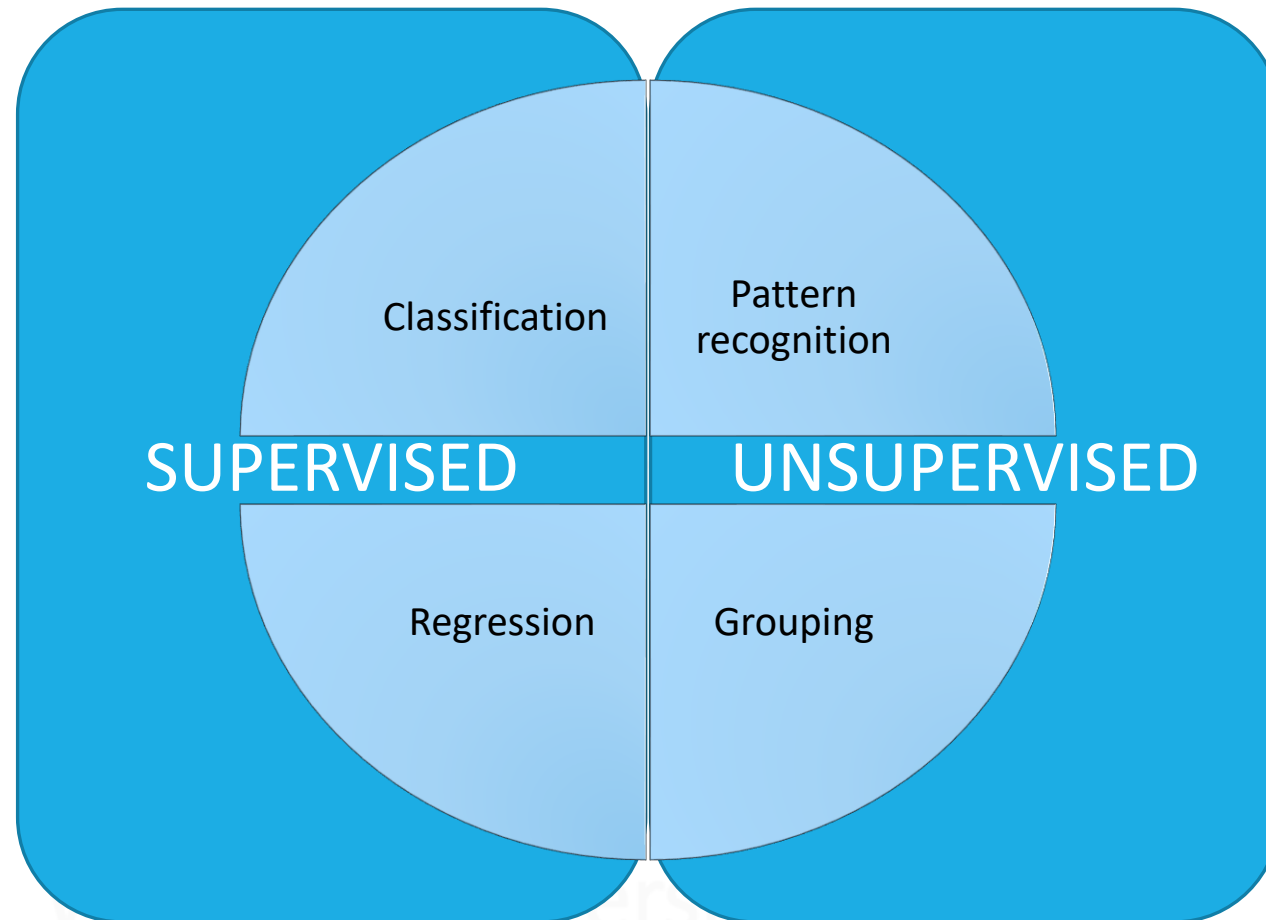
Machine learning core technologies



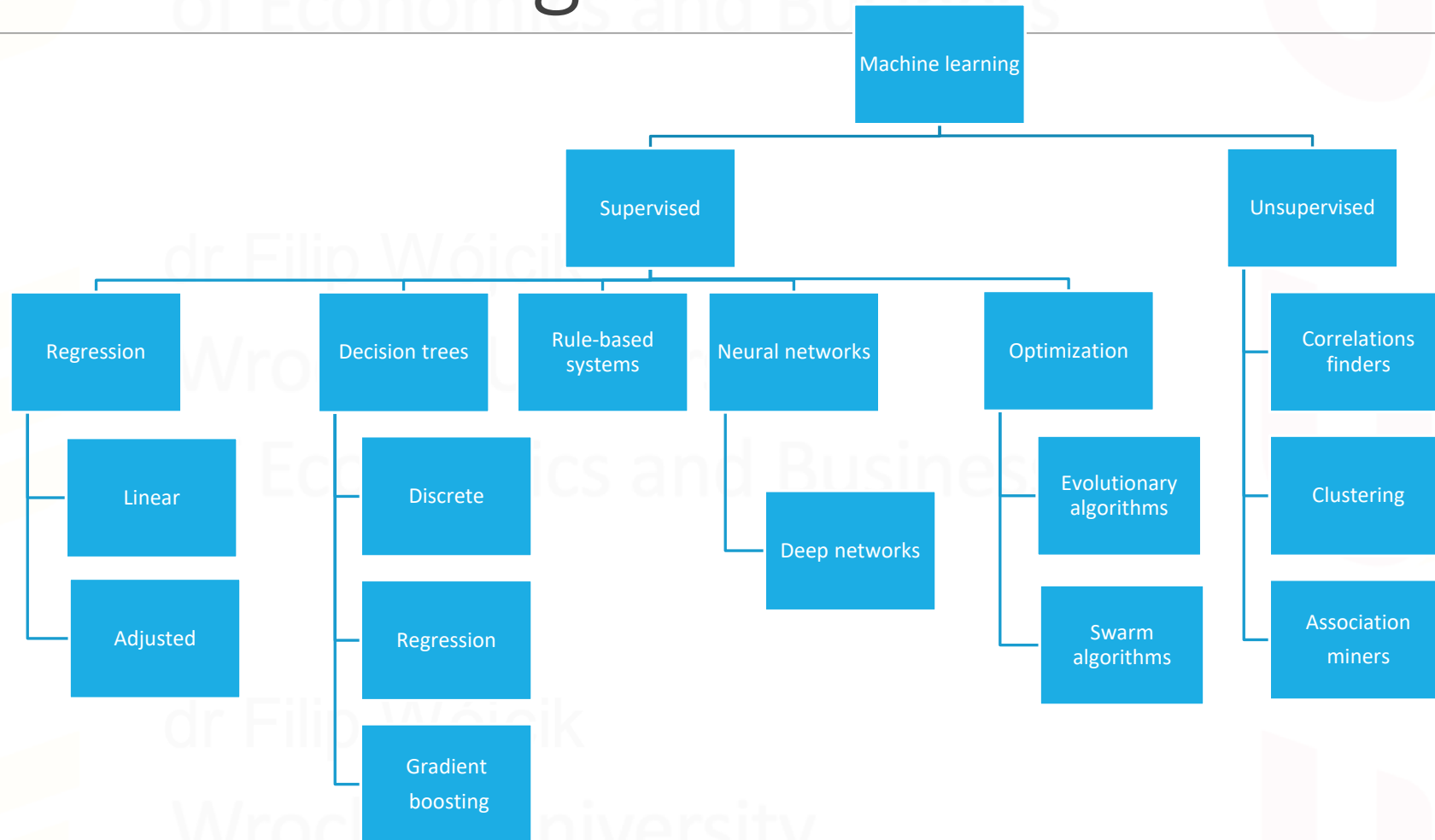
Source: KDNuggets, Jean-Francois Puget, 2017 The Most Popular Language For Machine Learning and Data Science Is..., <http://www.kdnuggets.com/2017/01/most-popular-language-machine-learning-data-science.html>



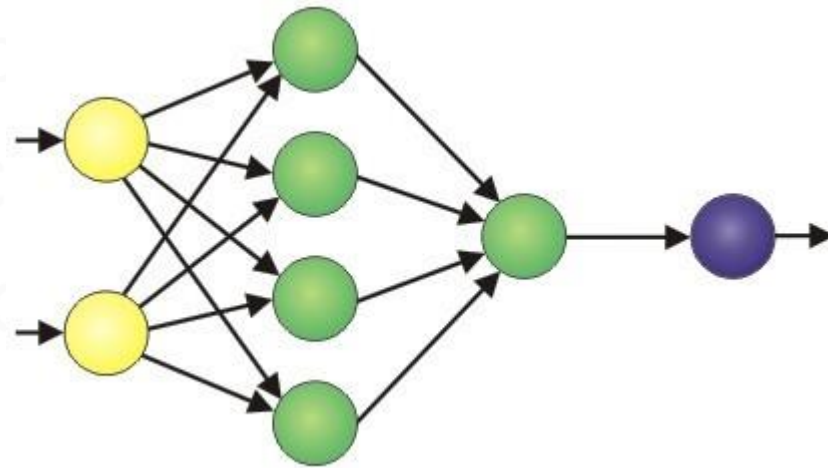
Machine learning madness



Machine learning madness



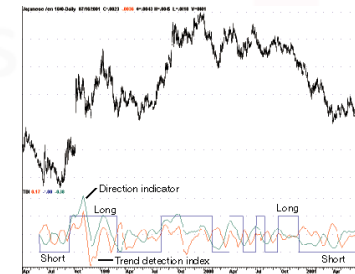
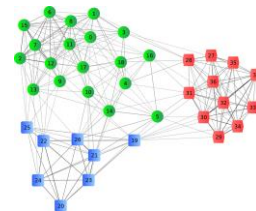
Is deep learning an ultimate algorithm?



Is deep learning an ultimate algorithm?

- Deep learning seems to be the most hyped ML algorithm nowadays
- Many analysts try to use it to solve ALL kinds of problems
- In reality deep learning is great tool for:

- ✓ Pattern recognition in images
- ✓ Motion detection
- ✓ Sentiment analysis
- ✓ Classification in fuzzy contexts
- ✓ Trend identification



Is deep learning an ultimate algorithm?

No free lunch theorem!

No learning algorithm has an inherent superiority over other learning algorithms for all problems.

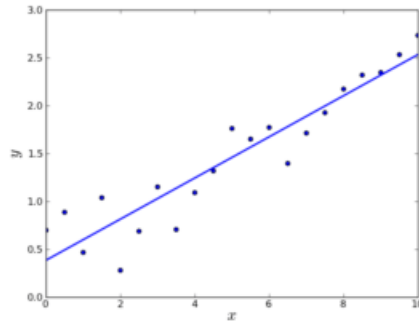
Wolpert, D.H., Macready, W.G. (1997), "No Free Lunch Theorems for Optimization"



Is deep learning an ultimate algorithm?

Top multi-purpose machine learning algorithms

1. Linear regression



2. Logistic regression

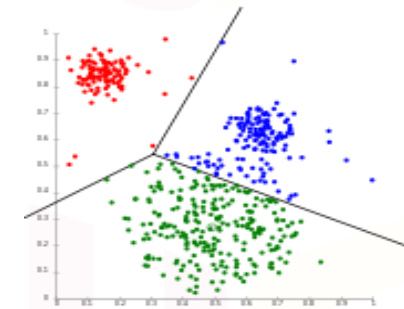
3. Decision trees



4. Random forests

5. K-means clustering

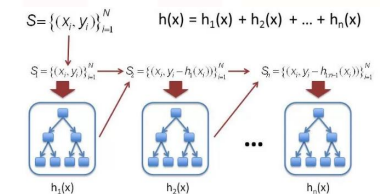
6. Support Vector Machines



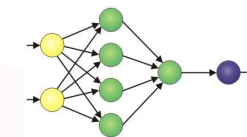
7. Gradient tree boosting

8. Naive Bayes

9. Neural networks



$$P(C|X) = \frac{P(X|C) P(C)}{P(X)}$$



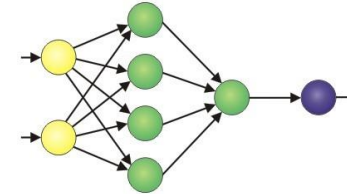
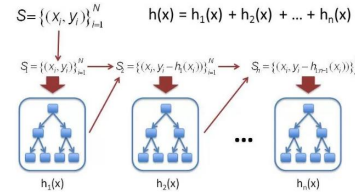
Source:

UpX, Namartha Peddi, 10 most popular machine learning algorithms, <https://upxacademy.com/10-popular-machine-learning-algorithms/>

Data Science Central, Top 10 most popular machine learning algorithms, <http://www.datasciencecentral.com/profiles/blogs/top-10-machine-learning-algorithms>



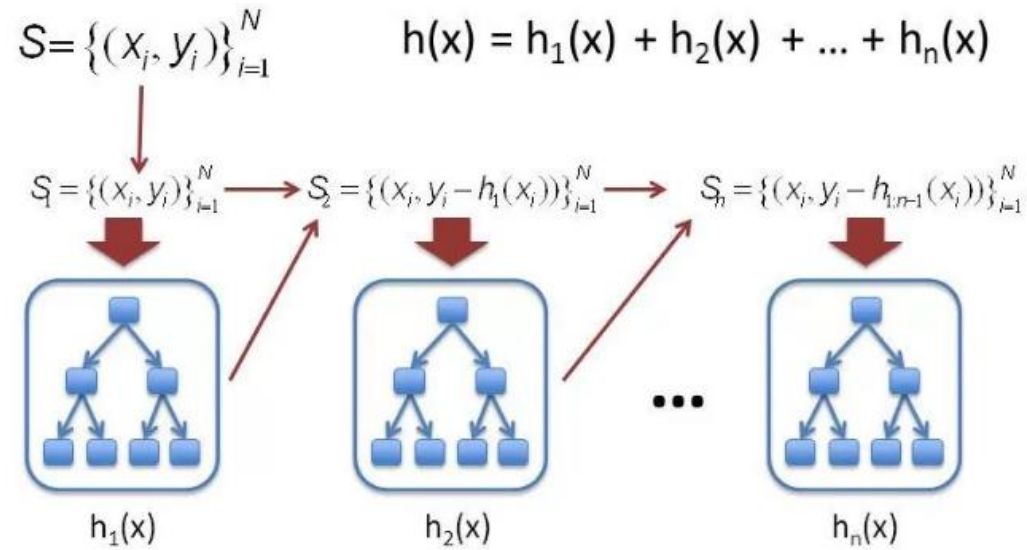
Is deep learning an ultimate algorithm?



	Xgboost	Deep learning / ANN
Building complexity	Fast to design & train	Require careful tuning
Data types	Tabular/structured data	Unstructured data/pictures/speech/etc.
Volumes of data	High to moderate	Small to INSANE
Mathematical explanation	Gradient function	Black - box



Xgboost – an overview



Xgboost – an overview

- Based on decision trees & random forests
- Using boosting procedure – learning on previous mistakes:
 - Putting more emphasis on wrongly classified examples
 - Training new classifiers one-by-one, with instruction to „watch out” for particular errors
- Can be used to both tasks:
 - Classification
 - Regression
- Both tasks are based on assigning numerical score, which corresponds to decision certainty level (classification) and/or raw numeric result

Xgboost – an overview

client	hotel	addons	money_spent	offer
business	Hilton	trip	40,000	deluxe
business	Hilton	full board	38,000	deluxe
business	Hilton	trip	40,000	deluxe
middle class	Meta	none	800	basic
middle class	Meta	meal	900	basic
manager	Meta	spa	1,500	premium

Value	Count	%
Deluxe	3	0.5
Basic	2	0.333
Premium	1	0.16666



client	hotel	addons	money_spent	offer
business	Hilton	trip	40,000	deluxe
business	Hilton	full board	38,000	deluxe
business	Hilton	trip	40,000	deluxe
middle class	Meta	none	800	basic
middle class	Meta	meal	900	basic
manager	Meta	spa	1,500	premium

Client == business?

True

False

hotel	addons	money_spent	offer
Hilton	trip	40,000	deluxe
Hilton	full board	38,000	deluxe
Hilton	trip	40,000	deluxe

hotel	addons	money_spent	offer
Meta	none	800	basic
Meta	meal	900	basic
Meta	spa	1,500	premium



Xgboost – an overview

- Multiple decision trees, **each learning from mistakes of its predecessors**
- Each tree receives **slightly randomized dataset** (different columns, resampled rows), to get rid of noise
- Steps:

1. Fit initial model - simple prediction, e.g. mean: $F_1(x) = \hat{y}$

2. Calculate error magnitude for each data point: $h_1(x) = y - F_1(x)$

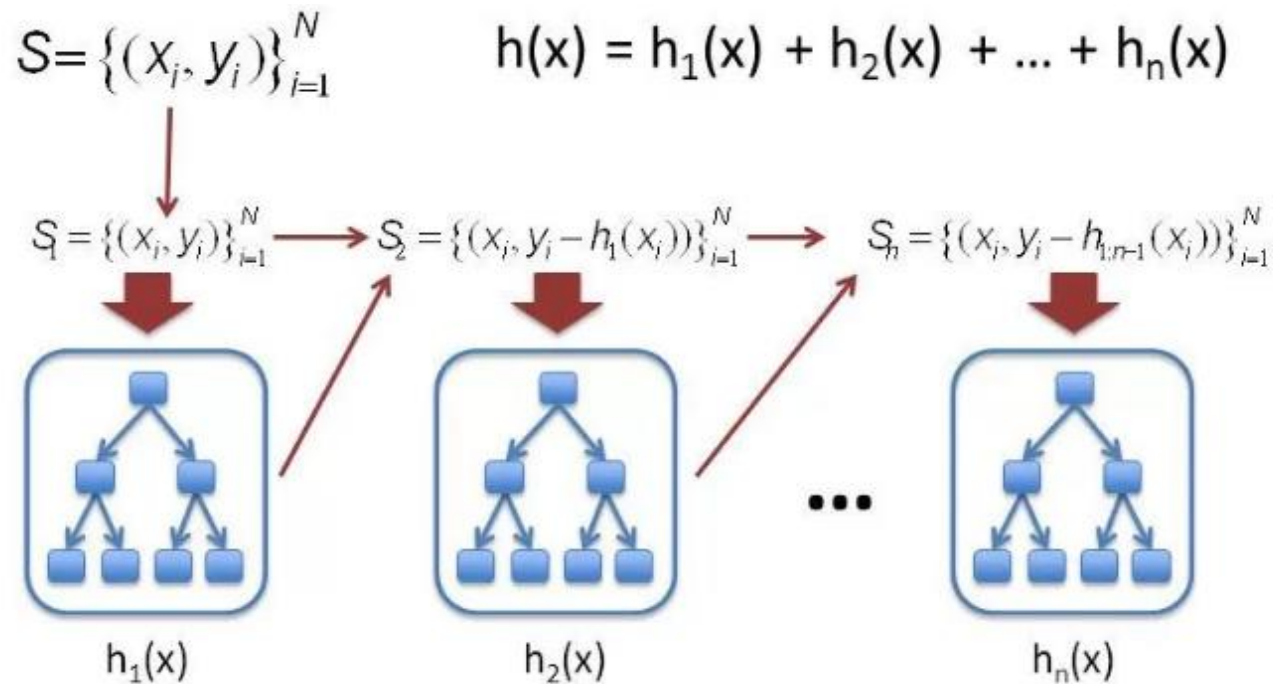
3. Create a new model, which will correct errors of its predecessor: $F_2(x) = F_1(x) + h_1(x)$

4. Continue until error rates are small enough or until reaching tree limit

$$F(x) = F_1(x) \rightarrow F_2(x) \rightarrow \dots \rightarrow F_n(x)$$



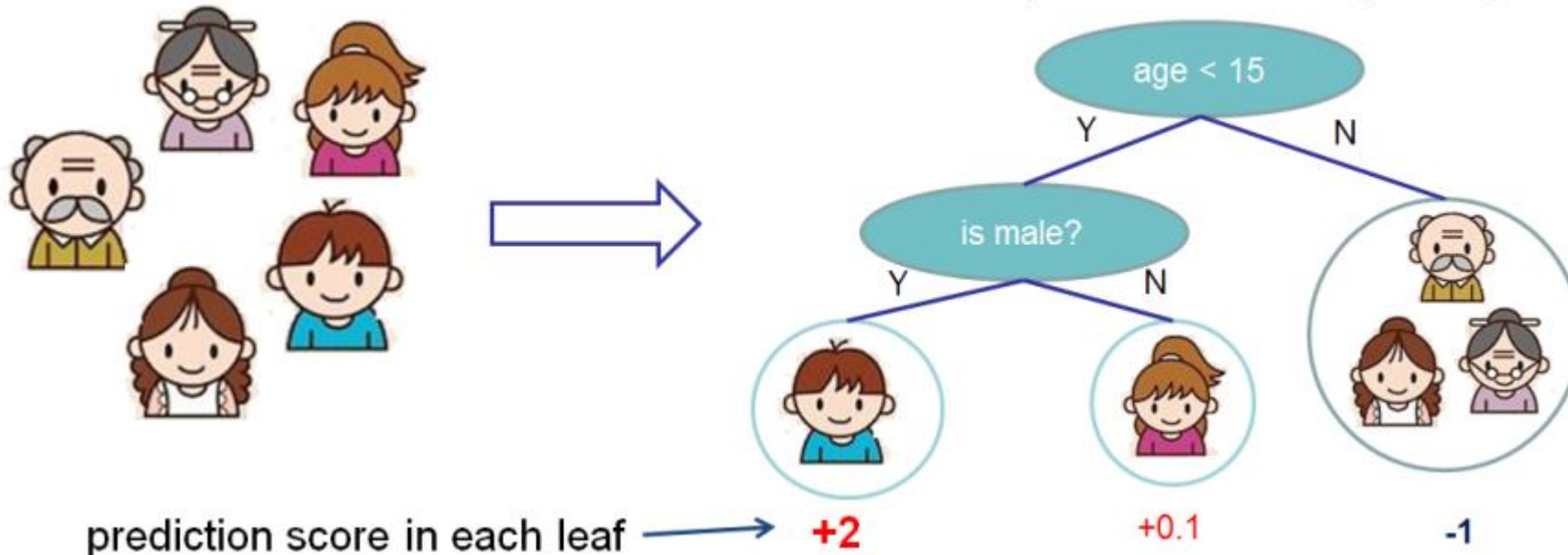
Xgboost – an overview



Xgboost – an overview

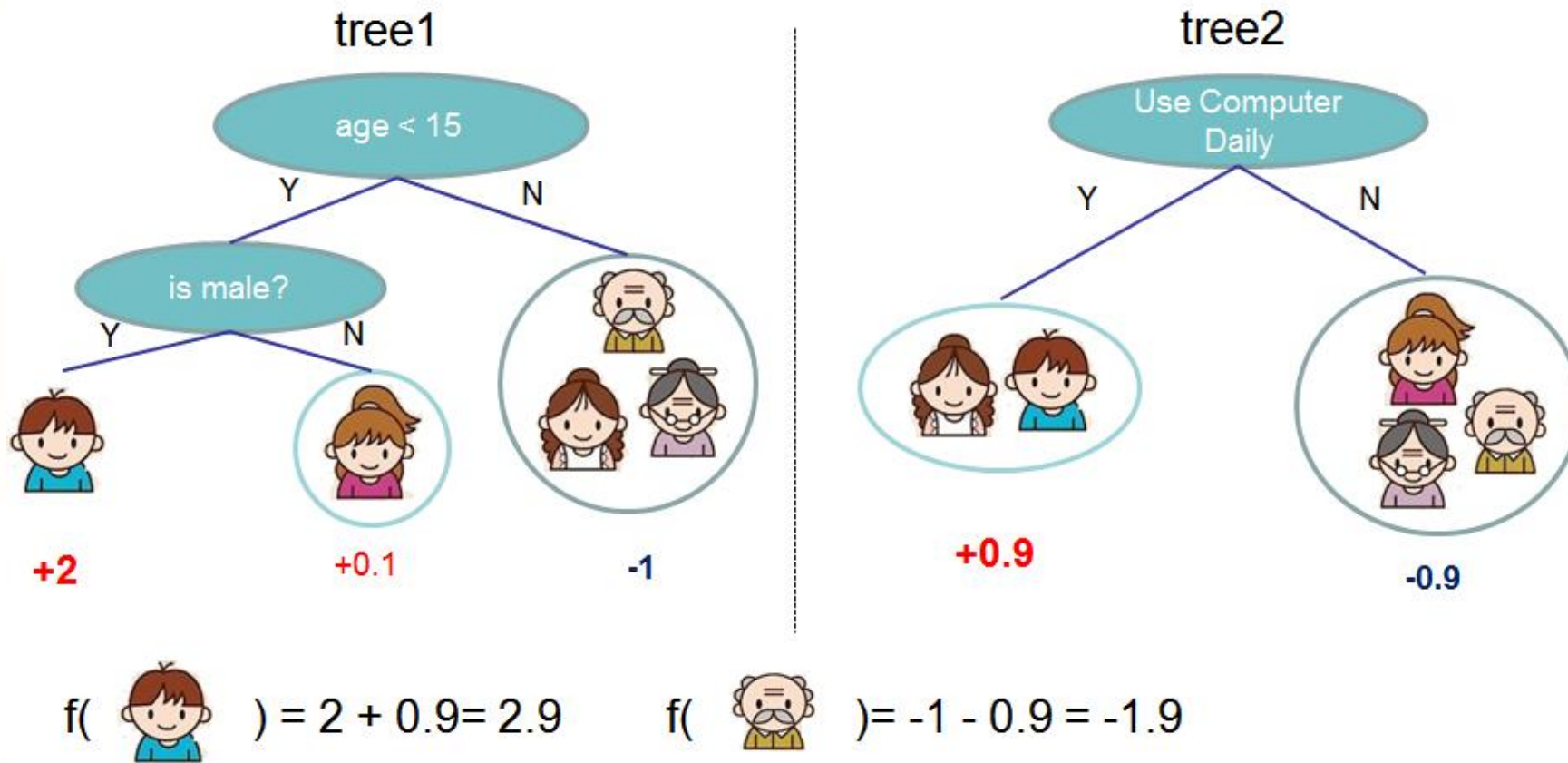
Input: age, gender, occupation, ...

Does the person like computer games



Source: official Xgboost documentation. <http://xgboost.readthedocs.io/en/latest/model.html>

Xgboost – an overview



Source: official Xgboost documentation. <http://xgboost.readthedocs.io/en/latest/model.html>

Xgboost – an overview

- Additional performance tuning tricks:
 - Using random set of columns for every tree – prevents overfitting
 - Using random set of examples for every tree – prevents overfitting
 - Using different tree depths – prevents overfitting 😊



Achievement unlocked

Survived math details part

H2O.AI and XGBoost

The logo for H2O.ai is centered on a yellow square. It consists of the text "H2O.ai" in a bold, black, sans-serif font. The "2" is a subscript, and the ".ai" is in a smaller font size than the "H2O".

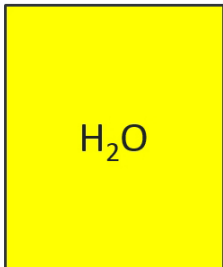
H₂O.ai

H2O.AI and XGBoost

- H2O AI is
 - Open-source
 - Fast
 - Scalable
 - In-memory

processing engine, equipped with predefined set of machine learning models
- Big-data-ready and optimized
 - Special data structures (hex) – SPARSE MATRICES HANDLING
 - Highly compressed
 - Lazy transformations (like in Apache Spark)
 - Immutable, distributed structures

H2O.AI and XGBoost



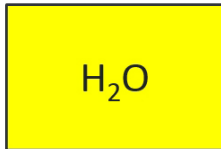
Standalone



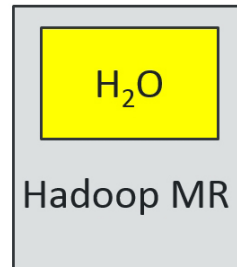
Scala



python™



Over YARN



H₂O in MR

Cluster

- Multi-node clusters
- Data distributed across nodes
- No limit on number of nodes 😊

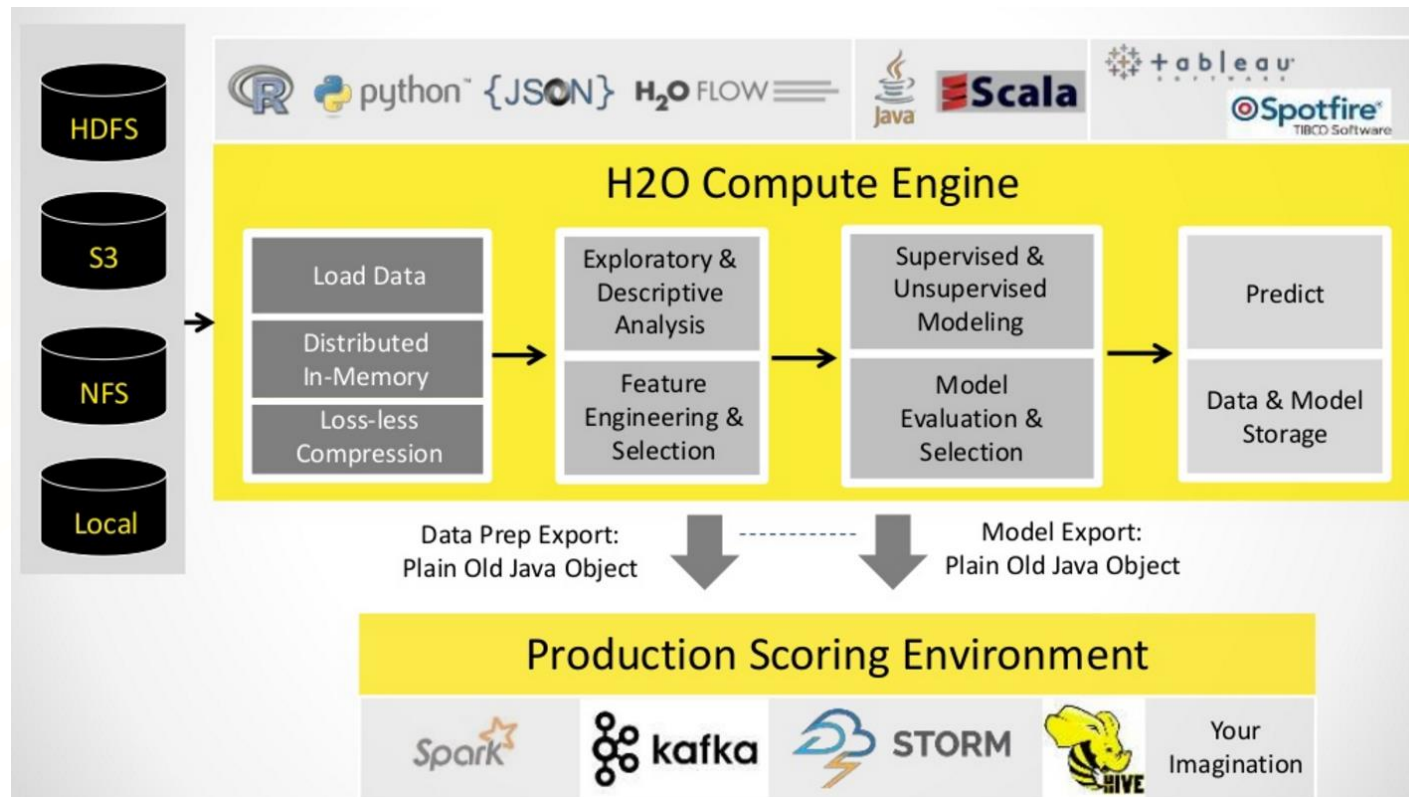
Key-value
pairs data

- Datasets referenced by key
- Columnary datastore
- SQL-like structures

Multi-
language

- Native implementation in Java
- Interfaces for Python/Scala/R

H2O.AI and XGBoost



Source: official H2O documentation, <http://docs.h2o.ai/h2o/latest-stable/h2o-docs/index.html>



H2O.AI and XGBoost



DEMO!

Recap

1. XGBoost is great algorithm for classification and predictive modelling:
 - Based on decision trees (basic estimator) and random forest (a lot of randomized trees)
 - Using gradient boosting algorithm – subsequent trees can learn on mistakes
2. Good alternative to Neural Networks and other algorithms
 - Fast training
 - Scalability
 - Accuracy
3. H2O is a scalable framework that comes with fast and accurate XGBoost implementation
4. H2O is:
 1. Big-data ready
 2. Integrating with Spark
 3. Open source
 4. Multi-language
 5. Easy to use for end-users and non-developers

THANK YOU!
