

ciência de dados
+ computação inteligente
==
ciência da computação
+ dados inteligentes

Wagner Arbex
wagner.arbex@embrapa.br
wagner.arbex@ufjf.edu.br



Computação científica...

- ... é também chamada *e-science*, *ciência computacional*, *ciência intensiva* etc. e não deve ser confundida com ciência da computação;
- ... é a aplicação de modelagem matemática e computacional em problemas científicos e tecnológicos;
- ... promove a computação massiva e/ou complexa, que são “novos” paradigmas da pesquisa científica;

Computação científica...

All models are wrong, but some are useful (George Box, 1976[8])


All models are wrong, and increasingly you can succeed without them (Peter Norvig, 2008)

Ciência de dados...

WIRED MAGAZINE: 16.07

SCIENCE : DISCOVERIES 

The End of Theory: The Data Deluge Makes the Scientific Method Obsolete

By Chris Anderson  06.23.08

"All models are wrong, but some are useful."

So proclaimed statistician George Box 30 years ago, and he was right. But what choice did we have? Only models, from cosmological equations to theories of human behavior, seemed to be able to consistently, if imperfectly, explain the world around us. Until now. Today companies like Google, which have grown up in an era of massively abundant data, don't have to settle for wrong models. Indeed, they don't have to settle for models at all.

Sixty years ago, digital computers made information readable. Twenty years ago, the Internet made it reachable. Ten years ago, the first search engine crawlers made it a single database. Now Google and like-minded companies are sifting through the most measured age in history, treating this massive corpus as a laboratory of the human condition. They are the children of the Petabyte Age.

Computação científica...

COMPUTAÇÃO CIENTÍFICA...

...is like youthful sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it ...

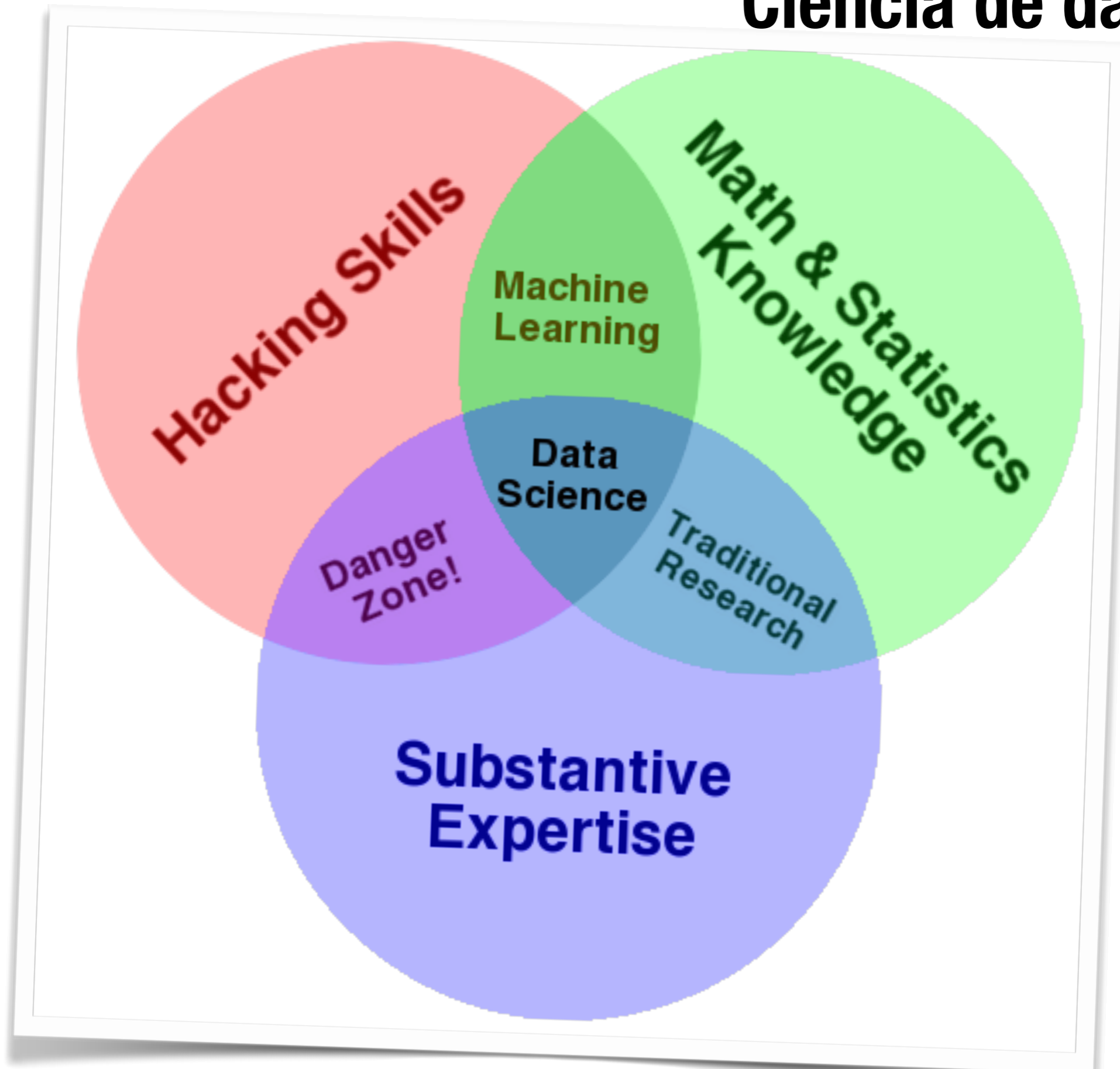
In a Data Deluge, Companies Seek to Fill a New Role

By Jessica Leber on May 22, 2013

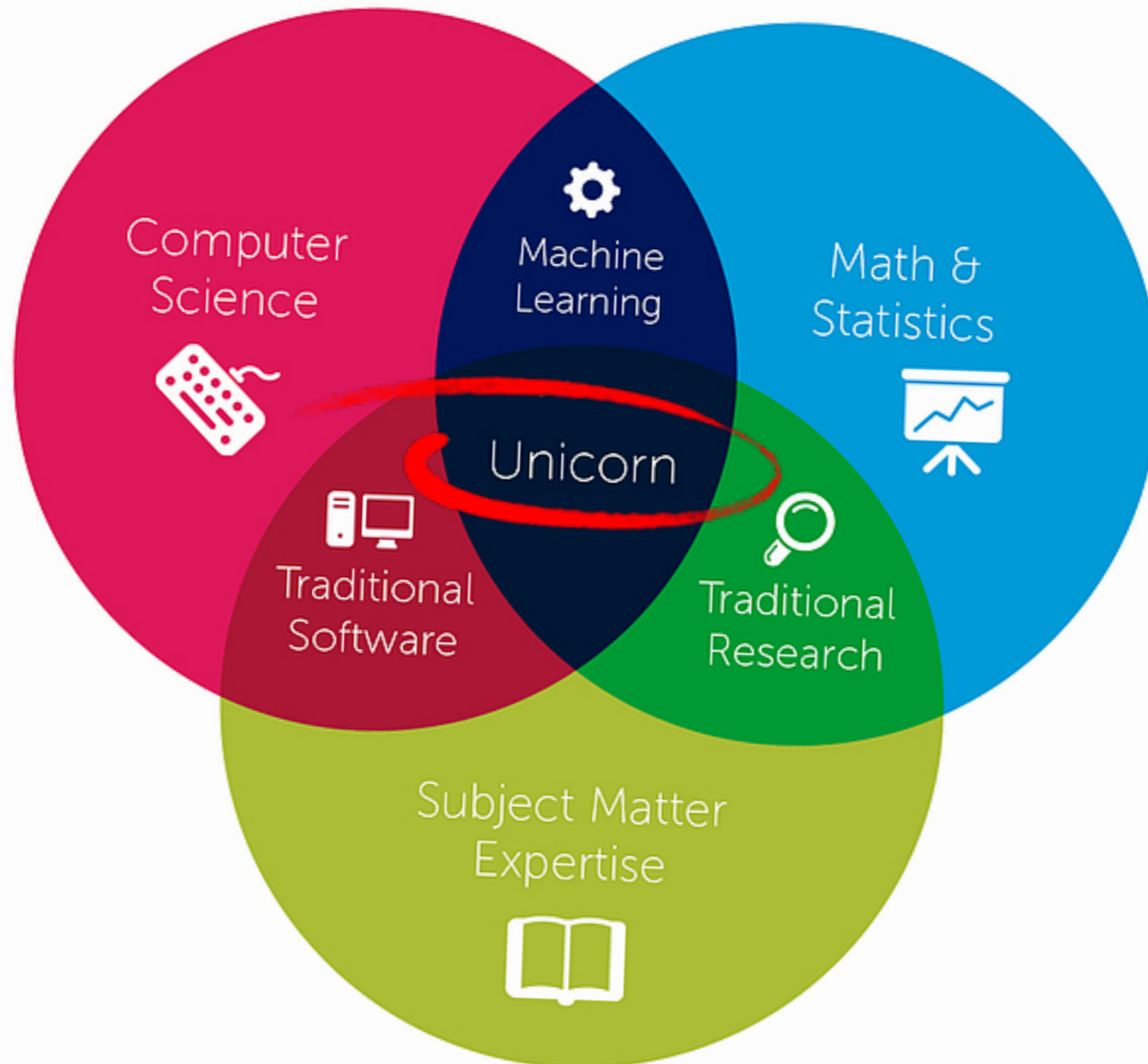
The job description "data scientist" didn't exist five years ago. No one advertised for an expert in data science, and you couldn't go to school to specialize in the field. Today, companies are fighting to recruit these specialists, courses on how to become one are popping up at many universities, and the *Harvard Business Review* even proclaimed that data scientist is [the "sexiest"](#) job of the 21st century.

Data scientists take huge amounts of data and attempt to pull useful information out. The job combines statistics and programming to identify sometimes subtle factors that can have a big impact on a company's bottom line, from whether a person will click on a certain type of ad to whether a new chemical will be toxic in the human body.

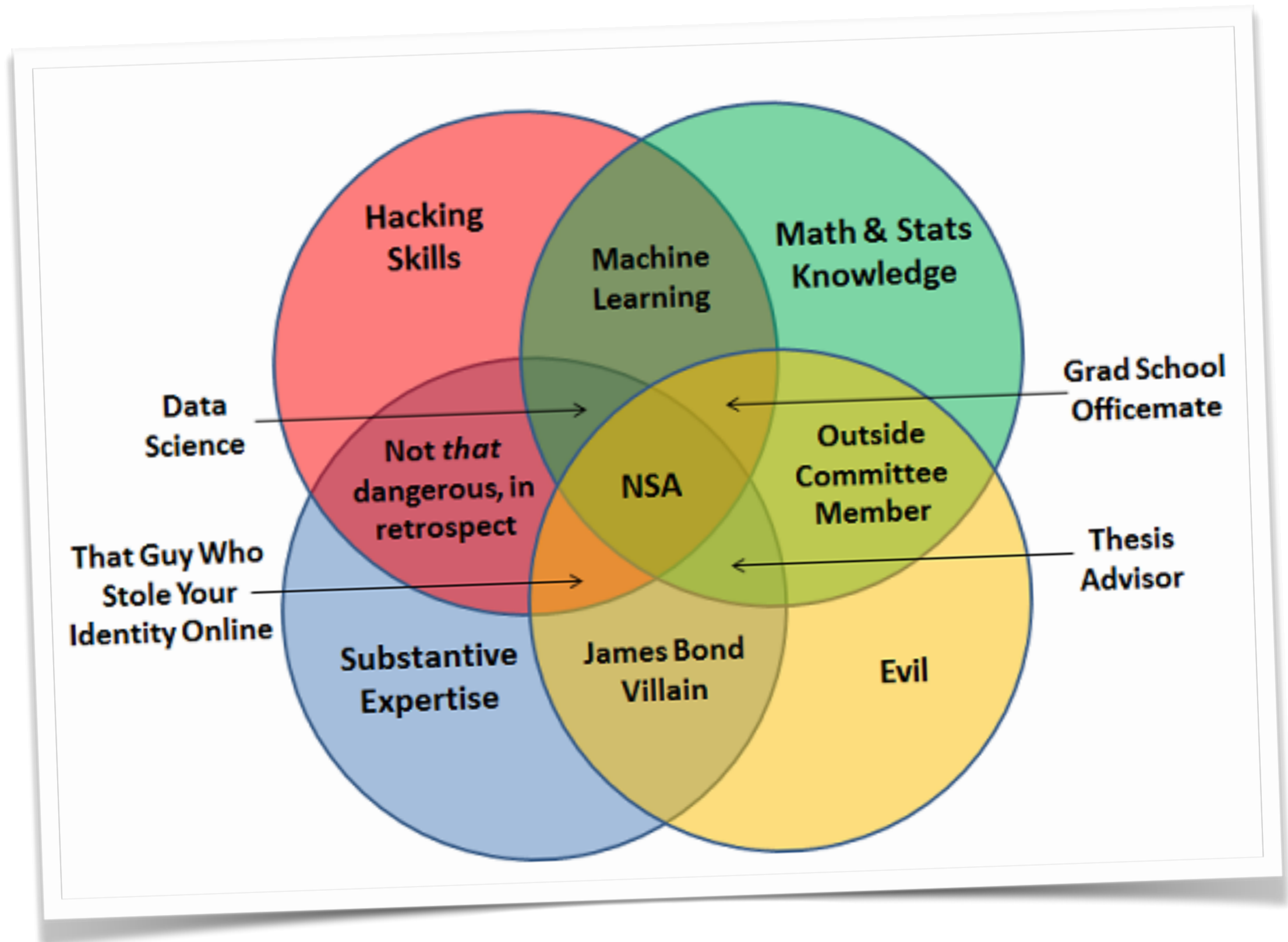
Ciência de dados...



Ciência de dados...



Ciência de dados...



Cientista de dados, The Unicorn

Data Scientist

Data Science allows front offices to better predict what will allow consumers are likely to buy. The ability to write algorithms that find relationships in datasets is valuable to generate actionable insight.

The Challenge

- Data Mining
- Analysis
- Communication

Industry Niche Titles

- Marketing Analytics Specialist
- Sales/CRM Channel Expert
- Product Services/Consumer Behaviour Analyst
- Data Analytics Expert

Urgent Need

Data Scientists - those with the technical savvy and analytical chops to derive meaning from all the information - are in high demand.

Skills by the Numbers

The skills and talents that make a fantastic Data Scientist

40%

Google's Eric Schmidt claims that every two days now we create as much information as we did from the dawn of civilization up until 2003.

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- Machine learning
- Statistical modeling
- Experiment design
- Bayesian inference
- Supervised learning: decision trees, random forests, logistic regression
- Unsupervised learning: clustering, dimensionality reduction
- Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- Computer science fundamentals
- Scripting language e.g. Python
- Statistical computing packages e.g. R
- Databases: SQL and NoSQL
- Relational algebra
- Parallel databases and parallel query processing
- MapReduce concepts
- Hadoop and Hive/Pig
- Custom reducers
- Experience with SAS-like AQS

COMMUNICATION & VISUALIZATION

- Ability to engage with senior management
- Story telling skills
- Translate data driven insights into decisions and actions
- Visual art design
- R packages like ggplot or lattice
- Knowledge of any of visualization tools e.g. Tableau, Qlik, Tableau

DOMAIN KNOWLEDGE & SOFT SKILLS

- Passionate about the business
- Curious about data
- Influence without authority
- Hacker mindset
- Problem solver
- Strategic, proactive, creative, innovative and collaborative

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- Machine learning
- Statistical modeling
- Experiment design
- Bayesian inference
- Supervised learning: decision trees, random forests, logistic regression
- Unsupervised learning: clustering, dimensionality reduction
- Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- Computer science fundamentals
- Scripting language e.g. Python
- Statistical computing packages e.g. R
- Databases: SQL and NoSQL
- Relational algebra
- Parallel databases and parallel query processing
- MapReduce concepts
- Hadoop and Hive/Pig
- Custom reducers
- Experience with SAS-like AQS

COMMUNICATION & VISUALIZATION

- Ability to engage with senior management
- Story telling skills
- Translate data driven insights into decisions and actions
- Visual art design
- R packages like ggplot or lattice
- Knowledge of any of visualization tools e.g. Tableau, Qlik, Tableau

DOMAIN KNOWLEDGE & SOFT SKILLS

- Passionate about the business
- Curious about data
- Influence without authority
- Hacker mindset
- Problem solver
- Strategic, proactive, creative, innovative and collaborative



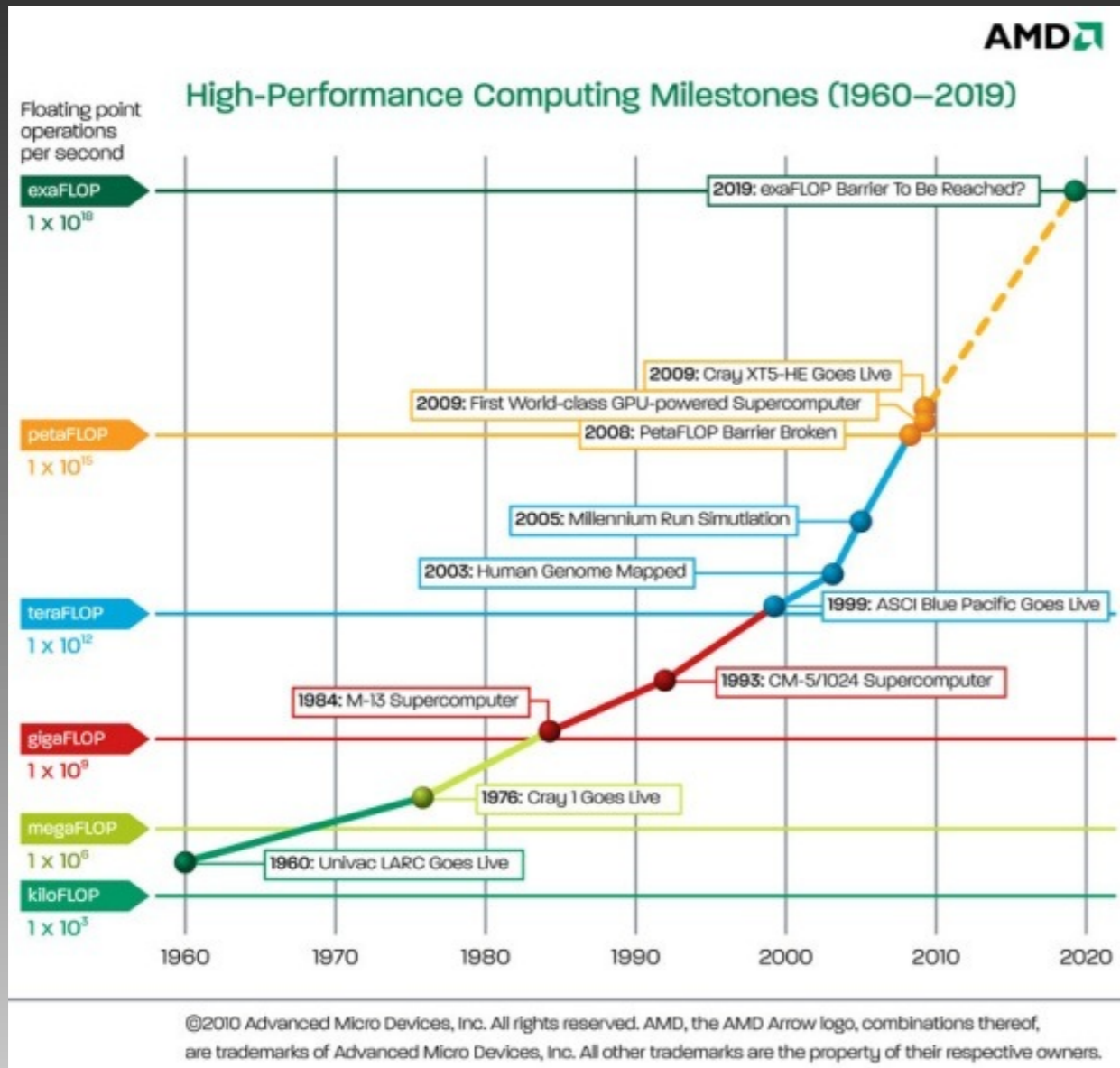
Ciência de dados...



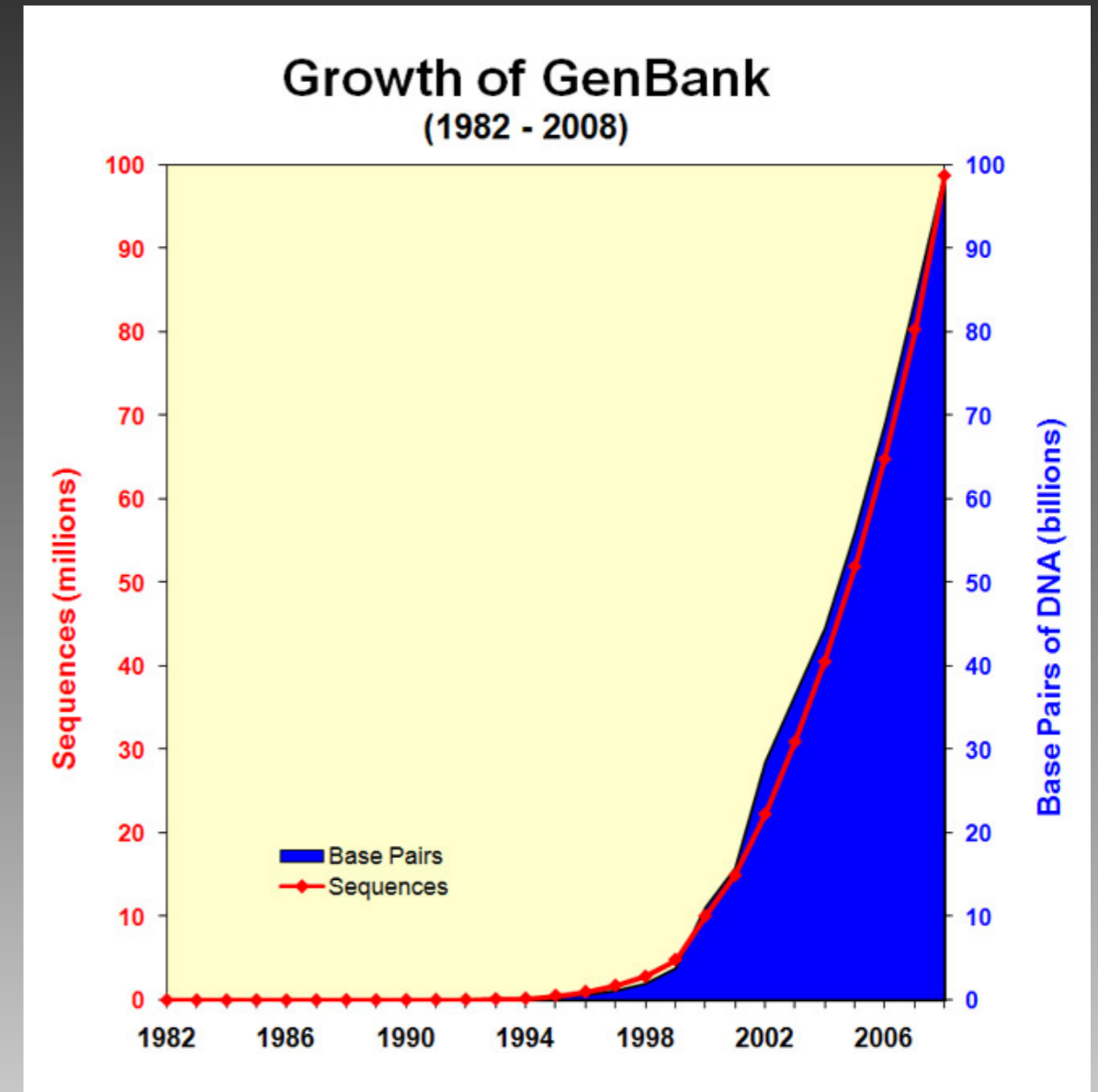
Mas, o que eu faço com isso?



Mas, o que eu faço com isso?

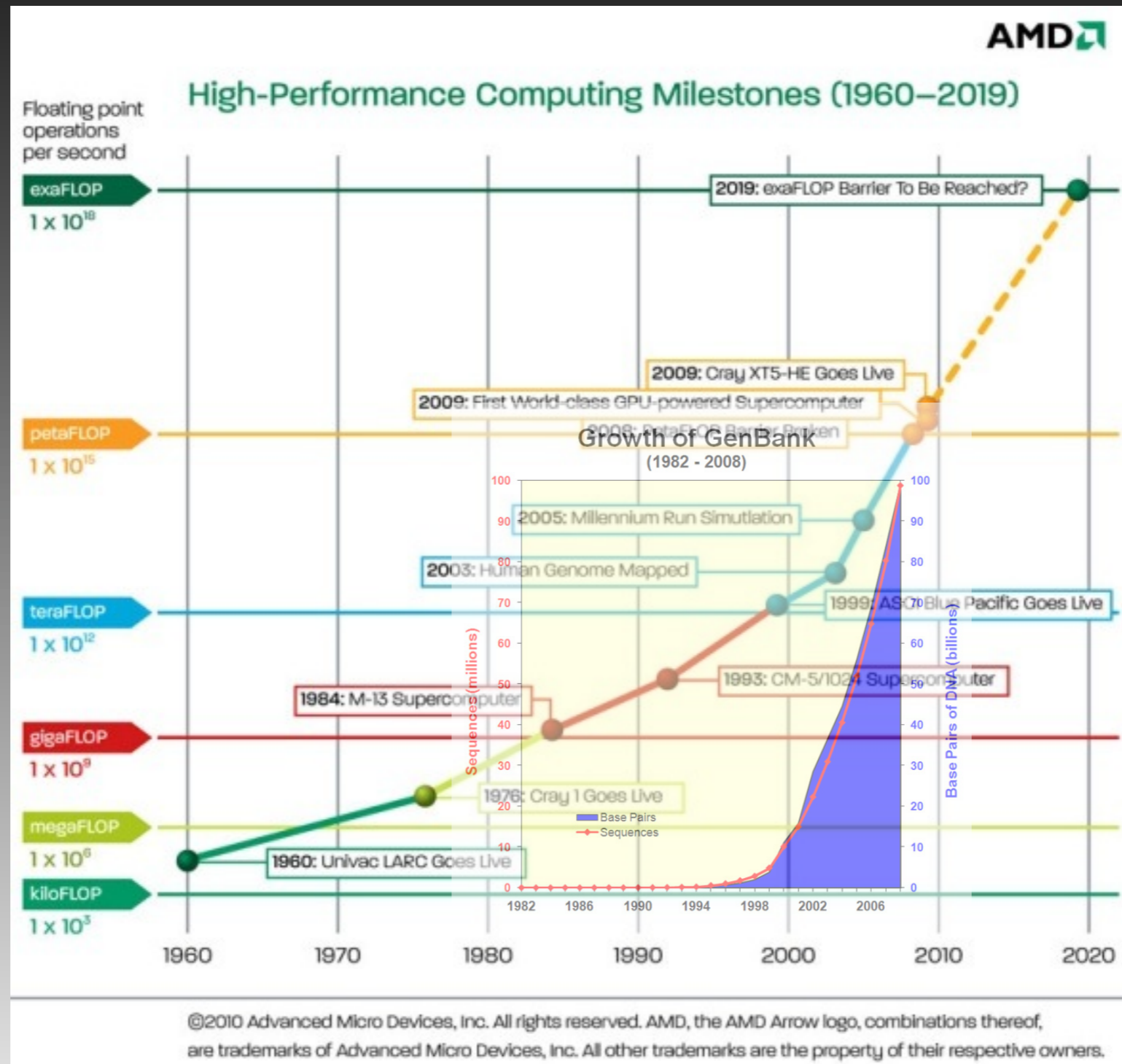


Fonte: AMD (2010)

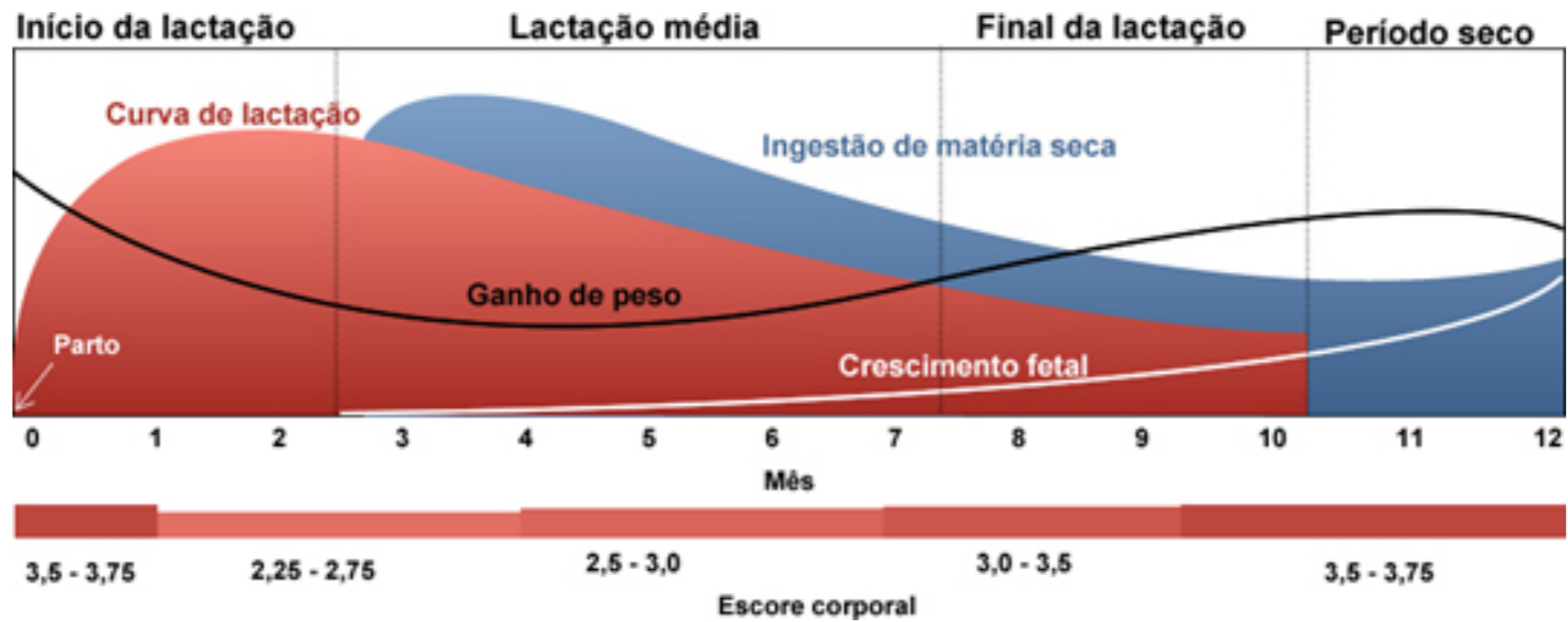


Fonte: NCBI (2008)

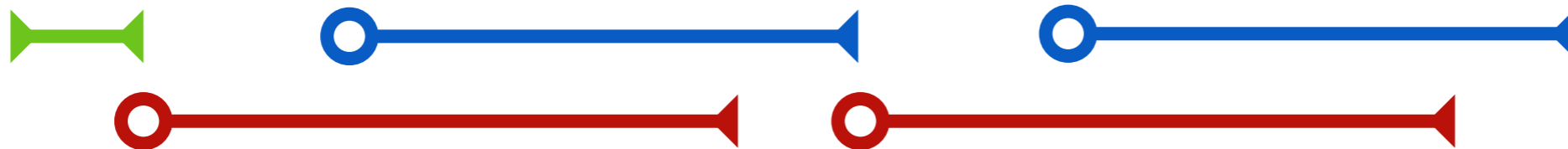
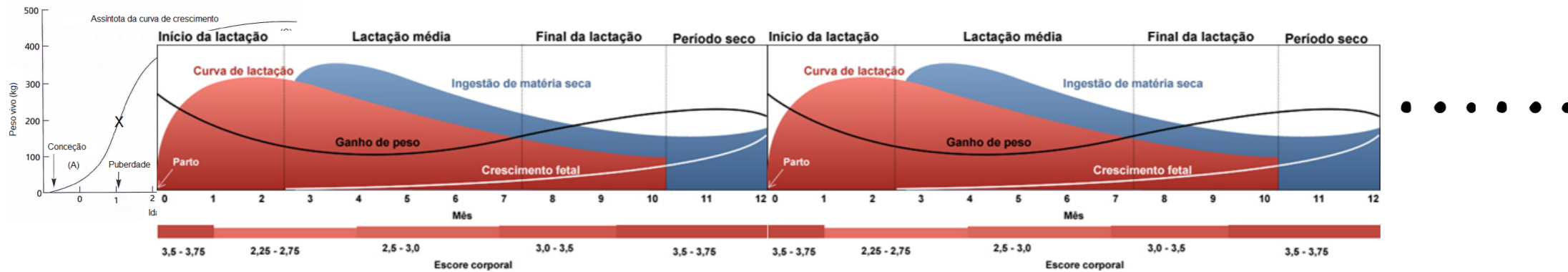
Mas, o que eu faço com isso?








Aplicações

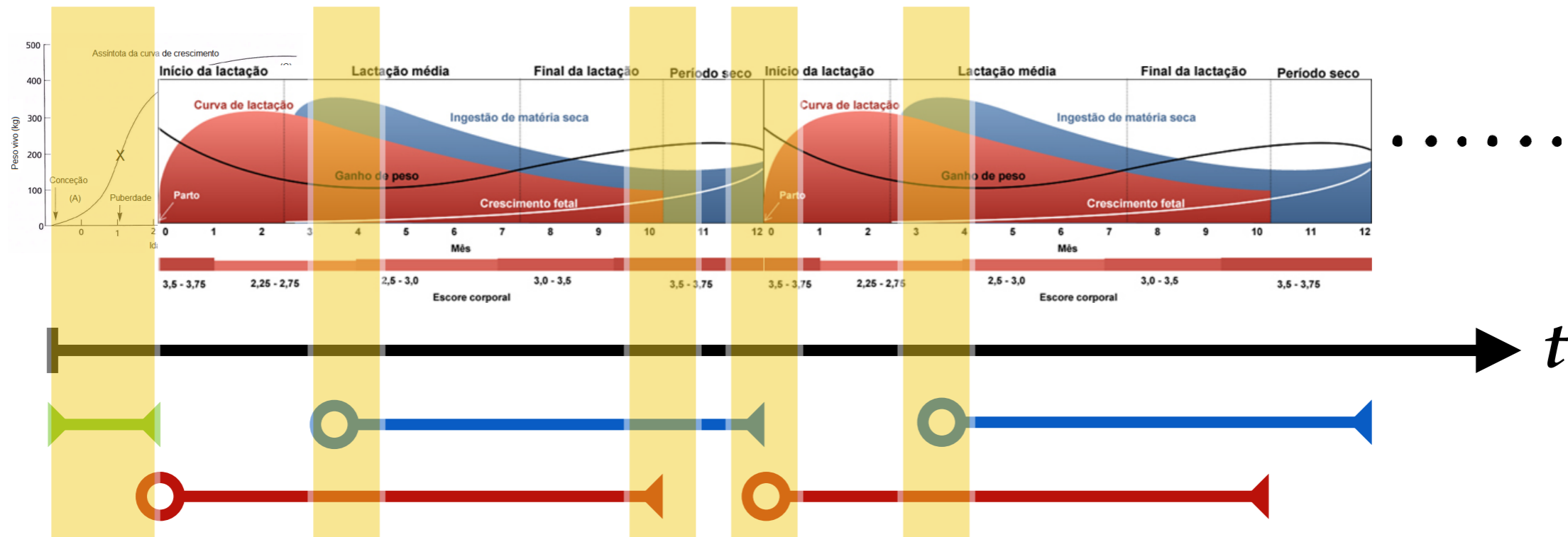







Aplicações



-  **Concepção**
-  **Início da lactação**
-  **Nascimento, crescimento, puberdade...**
-  **Parto**
-  **Secagem**

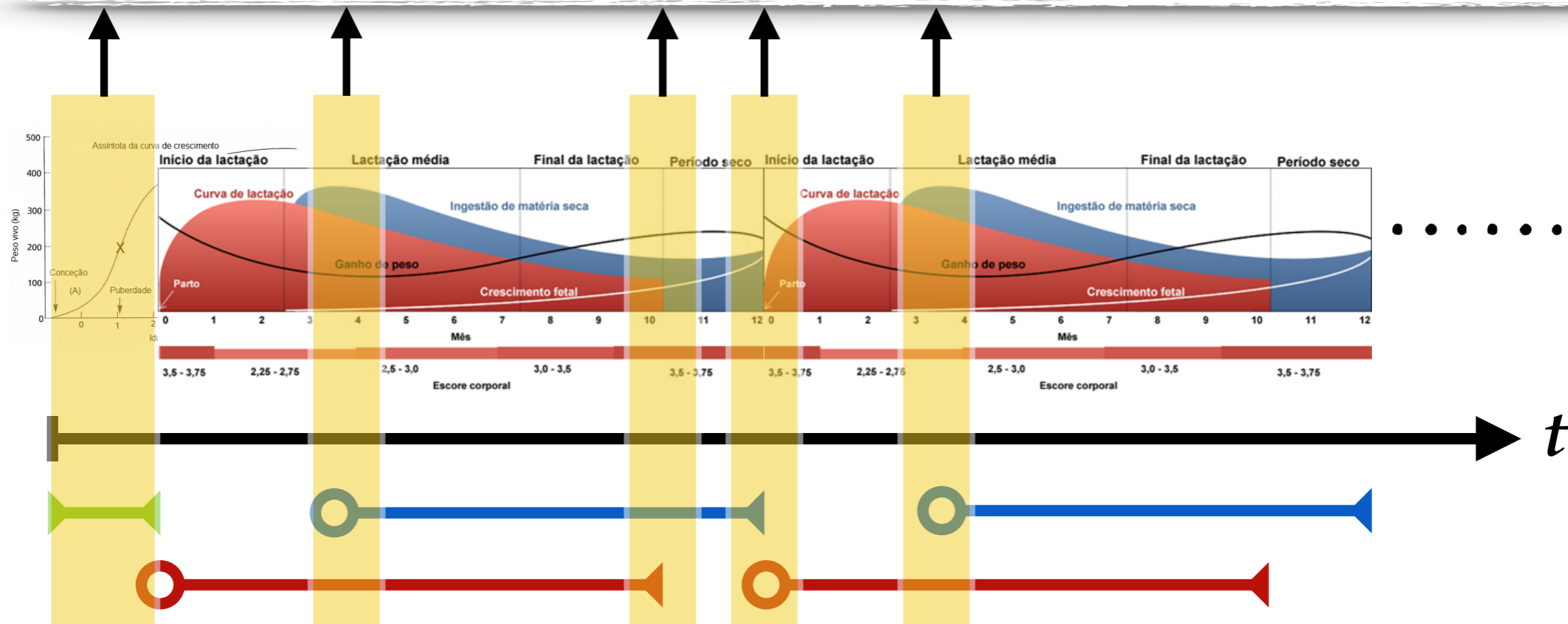
Aplicações








-  **Concepção**
-  **Início da lactação**
-  **Nascimento, crescimento, puberdade...**
-  **Parto**
-  **Secagem**

Aplicações

IoTtoMilk

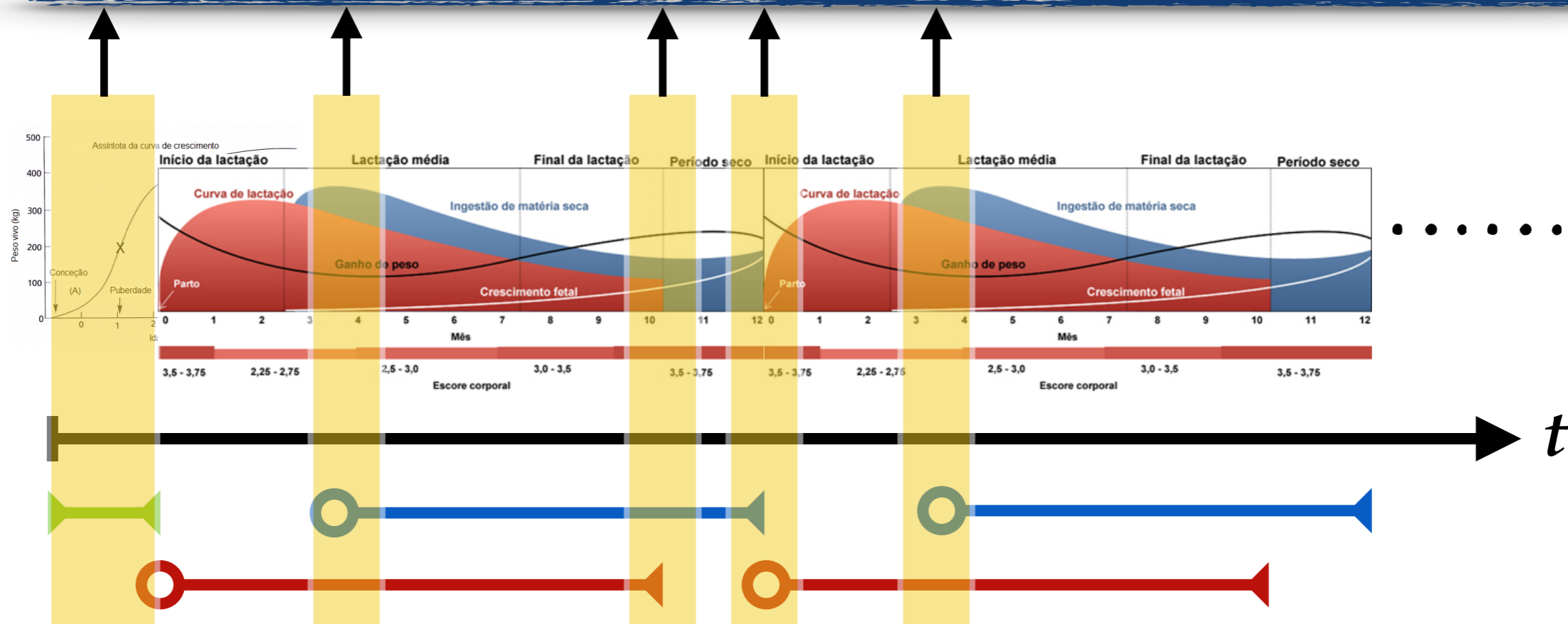


-  **Concepção**
-  **Início da lactação**
-  **Nascimento, crescimento, puberdade...**
-  **Parto**
-  **Secagem**

Aplicações

automação

IoTtoMilk



- Conceção
- Início da lactação
- ◀ Nascimento, crescimento, puberdade...
- ▶ Parto
- ◀ Secagem

Infraestrutura para ciência de dados

72

Int'l Conf. Internet Computing and Big Data | ICOMP'14 |

A Low-cost Infrastructure for Massive Storage of Phenotype Data for Dairy Cattle Genetic Improvement Programs

Wagner Arbex^{*,†}, Caio dos Santos Borsato de Carvalho^{*,‡}, Katia Santos^{*},
Vinicius Campista Brum[§] and Marcos Vinícius Barbosa da Silva^{*}
^{*}Brazilian Agricultural Research Corporation (Embrapa)
Juiz de Fora, MG, Brazil

[†]Correspondent author – wagner.arbex@embrapa.br

[‡]National Council for Scientific and Technological Development (CNPq) grant

[§]Federal University of Juiz de Fora (UFJF)

Abstract—The activities and progress of genetic breeding were always related to computing – or rather, the availability of appropriate computational resources to the behavior of genetic improvement, because, either by cross breeding or selection, the data sets involved in genetic improvement research and development activities care both in quantity and in quality. The demand for computing resources in these programs is extensive and intense, because genetic genomic evaluation

Add up to these aspects the need to interpret and understand the phenotype databases from a logical and effective structure for data storage and information retrieval.

Initiatives for data storage in genetic breeding programs are not unknown. Among others, the National Dairy Cattle Research Center (Embrapa Dairy Cattle) of the Brazilian Agricultural Research Corporation (Embrapa)

Banco de dados não convencionais

```
pc10122:Genotypes_50K arbex$ head -30 genotipo-H0L-56Kx577.txt

[Header]
BSGT Version      3.2.23
Processing Date   8/14/2008 11:44 AM
Content           USDA_Bovine_58K_271441_A.bpm
Num SNPs          56947
Total SNPs        56947
Num Samples       577
Total Samples     1630
[Data]
SNP Name          Sample ID      Allele1 - AB  Allele2 - AB
ARS-BFGL-BAC-10172 53998679      B            B
ARS-BFGL-BAC-1020  53998679      A            B
ARS-BFGL-BAC-10245 53998679      B            B
ARS-BFGL-BAC-10345 53998679      A            B
ARS-BFGL-BAC-10365 53998679      -            -
ARS-BFGL-BAC-10375 53998679      B            B
ARS-BFGL-BAC-10591 53998679      A            B
ARS-BFGL-BAC-10793 53998679      B            B
ARS-BFGL-BAC-10867 53998679      A            B
ARS-BFGL-BAC-10919 53998679      B            B
ARS-BFGL-BAC-10951 53998679      -            -
ARS-BFGL-BAC-10952 53998679      A            A
ARS-BFGL-BAC-10960 53998679      B            B
AR
AF
AF
AI
P
```

X-meeting • November 2015 • USP Genomics

Storage and recovery of dairy cattle genotype data from the data science approach

Rennan Silva, Fernanda Almeida, Wagner Arbex

Juiz de Fora, Embrapa Dairy Cattle

Abstract

s of the identification of molecular markers, which may vary naturally, the records that identify these markers are stored in a format chosen for the job. Usually, the data gathered from a marker. There are different patterns to present the genotype information about each individual, like the animal number and ID (JP) and allow this information to be accessed from a middle

Avaliação do desempenho relativo de SGBDs NoSQL para arquivos de genótipos

Vinicius Junqueira Schettino¹, Arthur Lorenzi Almeida¹,
Leojayme Rodrigues Manso Silva¹, Wagner Arbex^{1,2,*}

¹Universidade Federal de Juiz de Fora – UFJF
Dep. de Ciência da Computação – DCC
Campus Universitário, 36.036-900, Juiz de Fora, MG, Brasil

²Empresa Brasileira da Pesquisa Agropecuária – Embrapa
R. Eugênio do Nascimento, 610, 36.038-330, Juiz de Fora, MG, Brasil

schettino.vinicius@ufjf.edu.br

Seleção de variáveis e tomada de decisão...

As ações de pesquisa em genética e genômica, que antes se encontravam limitadas às "bancadas", já não existem mais sem que sejam complementadas com procedimentos computacionais, pois os dados obtidos nos laboratórios devem ser identificados, organizados, armazenados e interpretados para que possam ser recuperados, apresentados e, ainda, utilizados como um novo atributo de informação.

Se antes parecia estranho a junção de áreas com conceitos tão distantes, tais como, computação e genômica, atualmente, pode não ser totalmente compreendido, mas pela proximidade e "mistura" das aplicações desses conceitos.

O livro, **TACG - Talking About Computing and Genomics - Vol. I: Modelos e Métodos Computacionais em Bioinformática**, apresenta, em seus dois primeiros capítulos, discussões e aplicações de conceitos que tratam da mistura da computação com a genômica e, nos três capítulos finais, traz aplicações específicas em problemas técnico-científicos que estão em estudo por grupos de pesquisa em diversas partes do mundo.

Os autores e editores desta obra entregam para comunidade científica resultados de pesquisas que deixam claro como a genômica "clássica" pode beneficiar-se de métodos computacionais e matemáticos para avançar em direção à fronteira do conhecimento e, da mesma forma, apontam para onde a computação deve voltar suas pesquisas no sentido da investigação científica.

Talking About Computing and Genomics - Vol. I

Talking About Computing and Genomics TACG

Vol. I

Modelos e Métodos Computacionais em Bioinformática

Editores Técnicos
Wagner Arbex
Natália Florêncio Martins
Marta Fonseca Martins

Ministério da
Agricultura, Pecuária
e Abastecimento

GOVERNO FEDERAL
BRASIL
PAIS RICO É PAIS SEM POBREZA

Embrapa

Embrapa

Seleção de variáveis e tomada de decisão...

4

Metodologia para seleção de marcadores com máquina de vetores de suporte com regressão

Fabrízio Condé de Oliveira
Fernanda Nascimento Almeida
Fabyano Fonseca e Silva
Marcos Vinícius Gualberto Barbosa da Silva
Carlos Cristiano Hasenclever Borges
Wagner Arbex

Seleção de variáveis e tomada de decisão...

Int'l Conf. Artificial Intelligence | ICAI'14 |

497

NeuroSNP: Tool to Filter SNPs in Whole Genomic DNA

B. Zonovelli¹, C. C. H. Borges¹, and W. A. Arbex²

¹Federal University of Juiz de Fora, Juiz de Fora, MG, Brazil

²Brazilian Agricultural Research Corporation, Juiz de Fora, MG, Brazil

Abstract—The classification task the punctual differences in data found when comparing the genomes of individuals is complex. In data with noise higher than a final version, this can be in this noise difficult the process of differences how being single not. The SNPs are specific of aligned sequences of genetic variation. To the technique of the new whole genome the differ

Keywords: Bioinformatics, Machine Learning, Computation

Decision Support in Attribute Selection with Machine Learning Approach

Wagner Arbex

Brazilian Agricultural Research Corporation — Embrapa
Juiz de Fora, MG, Brazil
wagner.arbex@embrapa.br

Fabyano Fonseca e Silva

Federal University of Viçosa — UFV
Viçosa, MG, Brazil

Fabrízio Condé de Oliveira

Federal University of Juiz de Fora — UFJF
Juiz de Fora, MG, Brazil

Luis Varona

University of Zaragoza — UNIZAR
Zaragoza, Spain

Rui da Silva Verneque

Brazilian Agricultural Research Corporation — Embrapa

Guilherme Barbosa da Silva

Seleção de variáveis e tomada de decisão...

de Oliveira et al. *BMC Genomics* 2014, **15**(Suppl 7):S4
<http://www.biomedcentral.com/1471-2164/15/S7/S4>



RESEARCH

Open Access

SNPs selection using support vector regression and genetic algorithms in GWAS

Fabrízio Condé de Oliveira¹, Carlos Cristiano Hasenclever Borges¹, Fernanda Nascimento Almeida^{2,4}, Fabyano Fonseca e Silva³, Rui da Silva Verneque⁴, Marcos Vinicius GB da Silva⁴, Wagner Arbex^{1,4*}

From 9th International Conference of the Brazilian Association for Bioinformatics and Computational Biology (X-Meeting 2013)
Recife, Brazil. 3-6 November 2013

Abstract

Introduction: This paper proposes a new methodology to simultaneously select the most relevant SNPs markers for the characterization of any measurable phenotype described by a continuous variable using Support Vector Regression with Pearson Universal kernel as fitness function of a binary genetic algorithm. The proposed methodology is multi-attribute towards considering several markers simultaneously to explain the phenotype and is based jointly on statistical tools, machine learning and computational intelligence.

Results: The suggested method has shown potential in the simulated database 1, with additive effects only, and real database. In this simulated database, with a total of 1,000 markers, and 7 with major effect on the phenotype and the other 993 SNPs representing the noise, the method identified 21 markers. Of this total, 5 are relevant SNPs between the 7 but 16 are false positives. In real database, initially with 50,752 SNPs, we have reduced to 3,073

Identificação e avaliação de padrões

X-meeting • November 2015 • USP Genomics

Computational methods applied to identification of the Dairy Gir breed families

Gisele Silva, Tales Silva, Míria Bobó, Fernanda Almeida, Victor Moraes, Wagner Arbex, João Cláudio Panetto, Wagner Arbex

Federal University of Juiz de Fora (UFJF)
(Embrapa)

Abstract

The Embrapa Dairy Cattle and the Brazilian Association Brasileira dos Criadores de Gir Leiteiro (ABCGL) Improvement National Program of the Dairy Gir Leiteiro/PNMGL). Released in May 2015, the most published the results about 300 males, using molecular or contemporary thereof and more over 32,000 records and contemporary. Among several sets of data, the two databases (i) for store the certificate (registry)

X-meeting • November 2015 • USP Genomics

Applications of graph theory for verify the family relationship for genetic evaluations through the Animal Model

Pedro Bittencourt, Fernanda Almeida, Wagner Arbex
Federal University of Juiz de Fora (UFJF), Brazilian Agricultural Research Corporation (Embrapa)

Abstract

Animal genetic improvement is an usual method for increasing the productivity of the herd. Therefore, selecting animals that are potentially better than their peers and/or contemporary is important because the descendants of these animals will have enhanced characteristics. Genetic evaluations are used in animal genetic improvement programs to predict the potential genetic value of the animals and their PTAs (predicted transmitting abilities); however, to obtain good results, one needs as much information as possible about the individuals, their relatives, and their pedigree data. The Animal Model is a computational implementation of information about the kinship of animals

ciência de dados
+ computação inteligente
==
ciência da computação
+ dados inteligentes

Wagner Arbex
wagner.arbex@embrapa.br
wagner.arbex@ufjf.edu.br

