# Введение в Анализ Данных

Александр Сенов
Data Sceintist @ E-Contenta

# СУТЬ

DATA → KNOWLEDGE → ACTION

ECONTENTA

# ПРИЛОЖЕНИЯ

**Classification**    **Classification + Localization**    **Object Detection**    **Instance Segmentation**



CAT    CAT    CAT, DOG, DUCK    CAT, DOG, DUCK
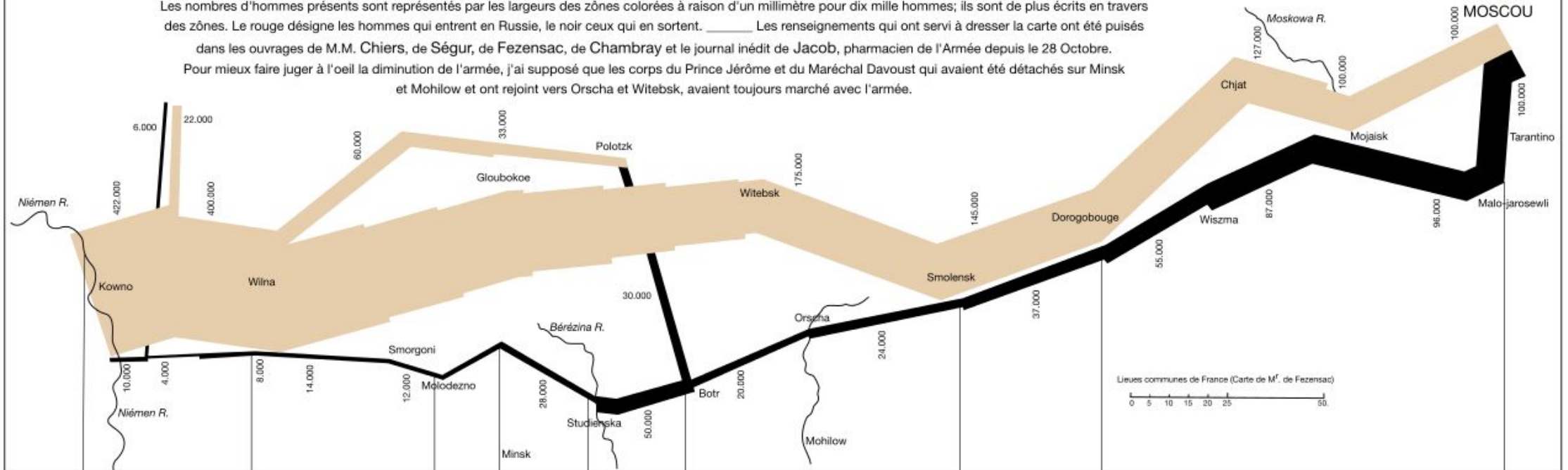
Single object    Multiple objects



Cats

Dogs

ECONTENTA

## Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.
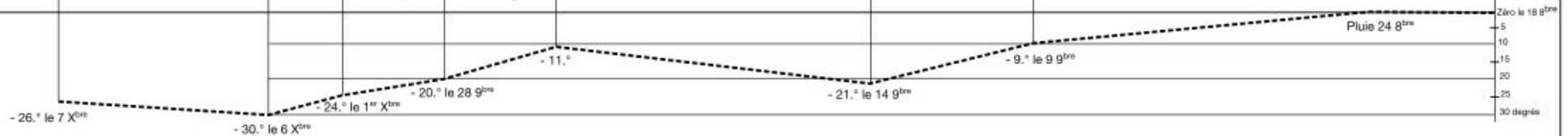
Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite. Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zônes colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zônes. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. _____ Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M.M. Chiers, de Ségur, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre. Pour mieux faire juger à l'oeil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davoust qui avaient été détachés sur Minsk et Mohilow et ont rejoint vers Orscha et Witebsk, avaient toujours marché avec l'armée.

TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Cosaques passent au galop le Niémen gelé.

- 26.° le 7 X^bre
- 30.° le 6 X^bre
- 24.° le 1^er X^bre
- 20.° le 28 9^bre
- 11.°
- 21.° le 14 9^bre
- 9.° le 9 9^bre
Pluie 24 8^bre
Zéro le 18 8^bre

ECONTENTA

HTTP://THRONESVIZ.GITHUB.IO/

ECONTENTA

ECONTENTA

# ЗАЧЕМ

**Job Trends** from Indeed.com

— Data-science



TB per Month

108% CAGR 2009-2014

3,600,000

3.6 EB
per mo

2.2 EB
per mo

1,800,000

1.2 EB
per mo

0.6 EB
per mo

0.09 EB
per mo

0.2 EB
per mo

0

2009    2010    2011    2012    2013    2014

For more details, see Appendix B: Forecast and Methodology.
Source: Cisco VNI Mobile, 2010

ECONTENTA

# THE WORLD OF DATA

**NUMBER OF EMAILS SENT EVERY SECOND**

## 2.9
MILLION

**DATA CONSUMED BY HOUSEHOLDS EACH DAY**

## 375
MEGABYTES

**VIDEO UPLOADED TO YOUTUBE EVERY MINUTE**

## 20
HOURS

**DATA PER DAY PROCESSED BY GOOGLE**

## 24
PETABYTES

**TWEETS PER DAY**

## 50
MILLION

**TOTAL MINUTES SPENT ON FACEBOOK EACH MONTH**

## 700
BILLION

**DATA SENT AND RECEIVED BY MOBILE INTERNET USERS**

## 1.3
EXABYTES

**PRODUCTS ORDERED ON AMAZON PER SECOND**

## 72.9
ITEMS

IN THE 21ST CENTURY, we live a large part of our lives online. Almost everything we do is reduced to bits and sent through cables around the world at light speed. But just how much data are we generating? This is a look at just some of the massive amounts of information that human beings create every single day.

SOURCES: Cisco; comScore; MapReduce; Radicati Group; Twitter; YouTube

A COLLABORATION BETWEEN GOOD AND OLIVER MUNDAY

IN PARTNERSHIP WITH **IBM**

ECONTENTA

**Josh Wills** @josh_wills

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

RETWEETS 1,255  LIKES 713

6:55 PM - 3 May 2012

JOB GROWTH AND DEMAND

NAMED THE
TOP JOB
IN AMERICA
FOR 2016
BY GLASSDOOR

SEXIEST JOB
OF THE
21ST CENTURY
BY HARVARD BUSINESS REVIEW

79.7%
OF DATA SCIENTISTS REPORT THERE IS A
SHORTAGE
IN THEIR FIELD

11% PROJECTED GROWTH FROM 2014 TO 2024
7% FASTER THAN GROWTH FOR ALL OCCUPATIONS

ECONTENTA

ЧТО

Data Asquisition

Data Parsing

Data Filtering & Cleaning

Exploration & Patterns Mining

Results Presentation

**Author: Swami Chandrasekaran**

**1. Fundamentals**
- Matrices & Linear Algebra Fundamentals
- Hash Functions, Binary Tree, O(n)
- Relational Algebra, DB Basics
- Inner, Outer, Cross, Theta Join
- CAP Theorem
- Tabular Data
- Data Frames & Series
- Sharding
- OLAP
- Multidimensional Data Model
- ETL
- Reporting Vs BI Vs Analytics
- JSON & XML
- NoSQL
- Regex
- Vendor Landscape
- Env Setup
- Entropy

**2. Statistics**
- Pick a Dataset (UCI Repo)
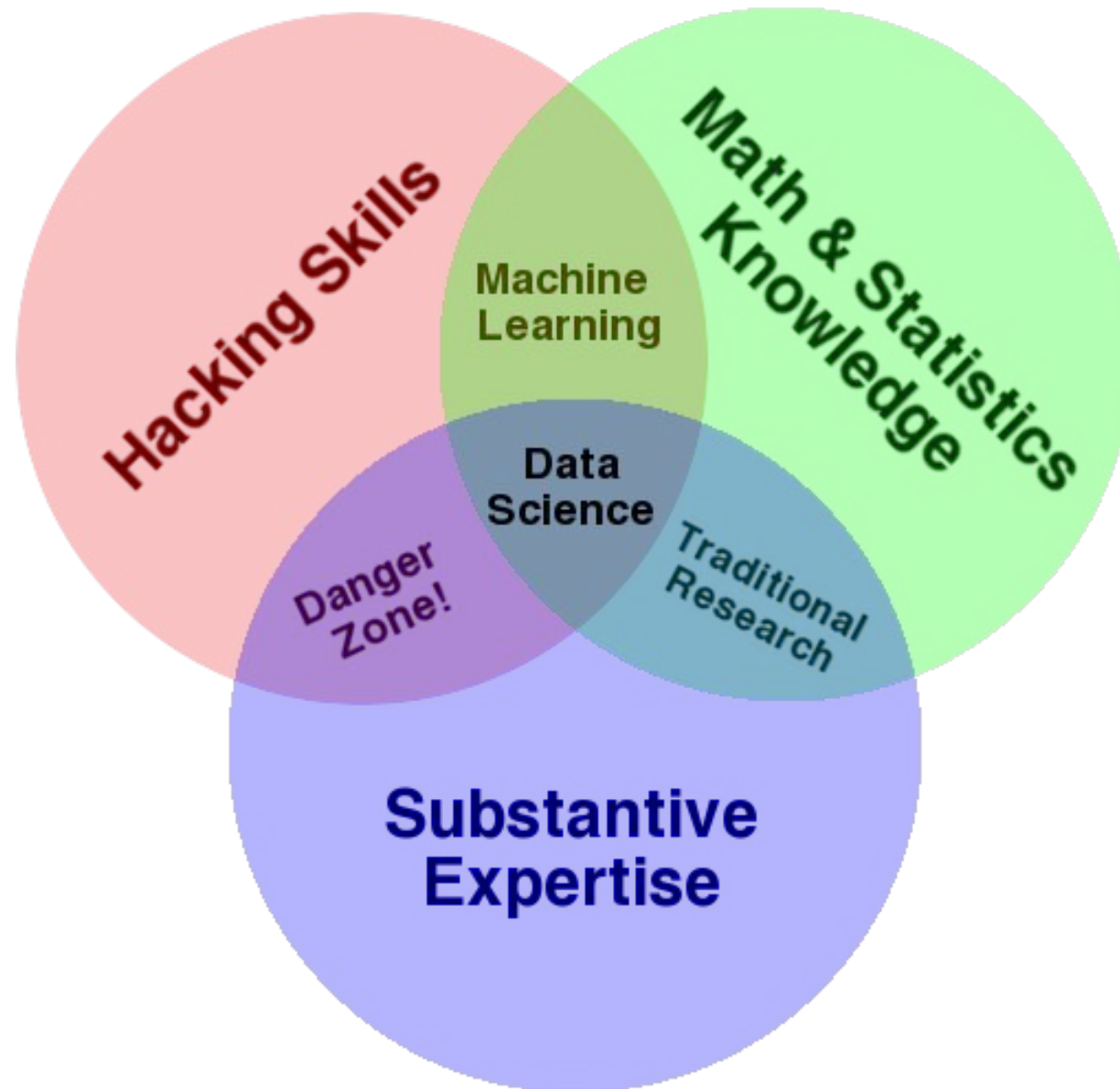- Descriptive Statistics (mean, median, range, SD, Var)
- Exploratory Data Analysis
- Histograms
- Percentiles & Outliers
- Probability Theory
- Bayes Theorem
- Random Variables
- Cumul Dist Fn (CDF)
- Continuos Distributions (Normal, Poisson, Gaussian)
- Skewness
- ANOVA
- Prob Den Fn (PDF)
- Central Limit Theorem
- Monte Carlo Method
- Hypothesis Testing
- p-Value
- Chi² Test
- Estimation
- Confid Int (CI)
- MLE
- Kernel Density Estimate
- Regression
- Covariance
- Correlation
- Pearson Coeff
- Causation
- Least² Fit
- Euclidean Distance

**3. Programming**
- Python Basics
- Working in Excel
- R Setup R Studio
- R Basics
- Expressions
- Variables
- IBM SPSS
- Rapid Miner
- Vectors
- Matrices
- Arrays
- Factors
- Lists
- Data Frames
- Reading CSV Data
- Reading Raw Data
- Subsetting Data
- Manipulate Data Frames
- Functions
- Factor Analysis
- Install Pkgs

**4. Machine Learning**
- What is ML?
- Numerical Var
- Categorical Var
- Supervised Learning
- Unsupervised Learning
- Concepts, Inputs & Attributes
- Training & Test Data
- Classifier
- Prediction
- Lift
- Overfitting
- Bias & Variance
- Trees & Classification
- Classification Rate
- Decision Trees
- Boosting
- Naïve Bayes Classifiers
- K-Nearest Neighbor
- Logistic Regression
- Ranking
- Linear Regression
- Perceptron
- Hierarchical Clustering
- K-means Clustering
- Neural Networks
- Sentiment Analysis
- Collaborative Filtering
- Tagging
- *Regression*
- *Classification*
- *Clustering*

**5. Text Mining / NLP**
- Corpus
- Named Entity Recognition
- Text Analysis
- UIMA
- Term Document Matrix
- Term Frequency & Weight
- Support Vector Machines
- Association Rules
- Market Based Analysis
- Feature Extraction
- Using Mahout
- Using Weka
- Using NLTK
- Classify Text
- Vocabulary Mapping

**6. Visualization**
- Data Exploration in R (Hist, Boxplot etc)
- Uni, Bi & Multivariate Viz
- ggplot2
- Histogram & Pie (Uni)
- Tree & Tree Map
- Scatter Plot (Bi)
- Line Charts (Bi)
- Spatial Charts
- Survey Plot
- Timeline
- Decision Tree
- D3.js
- InfoVis
- IBM ManyEyes
- Tableau

**7. Big Data**
- Map Reduce Fundamentals
- Hadoop Components
- HDFS
- Data Replication Principles
- Setup Hadoop (IBM / Cloudera / HortonWorks)
- Name & Data Nodes
- Job & Task Tracker
- M/R Programming
- Sqoop: Loading Data in HDFS
- Flume, Scribe: For Unstruct Data
- SQL- with Pig
- DWH with Hive
- Scribe, Chukwa For Weblog
- Using Mahout
- Zookeeper Avro
- Storm: Hadoop Realtime
- Rhadoop, RHIPE
- rmr
- Cassandra
- MongoDB, Neo4j

**8. Data Ingestion**
- Summary of Data Formats
- Data Discovery
- Data Sources & Acquisition
- Data Integration
- Data Fusion
- Transformation & Enrichment
- Data Survey
- Google OpenRefine
- How much Data?
- Using ETL

**9. Data Munging**
- Dimensionality & Numerosity Reduction
- Normalization
- Data Scrubbing
- Handling Missing Values
- Unbiased Estimators
- Binning Sparse Values
- Feature Extraction
- Denoising
- Sampling
- Stratified Sampling
- Principal Component Analysis

**10. Toolbox**
- MS Excel w/ Analysis ToolPak
- Java, Python
- R, R-Studio, Rattle
- Weka, Knime, RapidMiner
- Hadoop Dist of Choice
- Spark, Storm
- Flume, Scibe, Chukwa
- Nutch, Talend, Scraperwiki
- Webscraper, Flume, Sqoop
- tm, RWeka, NLTK
- RHIPE
- D3.js, ggplot2, Shiny
- IBM Languageware
- Cassandra, MongoDB

5% · 15% · 50% · 30% · 40% · 50% · 80% · 100%

# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th 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 package 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 xaaS like AWS

## DOMAIN KNOWLEDGE & SOFT SKILLS

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

## COMMUNICATION & VISUALIZATION

- ☆ Able 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. Flare, D3.js, Tableau

ECONTENTA

ЧЕМ

**SQL**

```sql
SELECT e.ID, e.LastName, e.FirstName, pn.Number
FROM Employee e
LEFT OUTER JOIN PhoneNumber pn
ON e.ID = pn.ID
```

Results | Messages

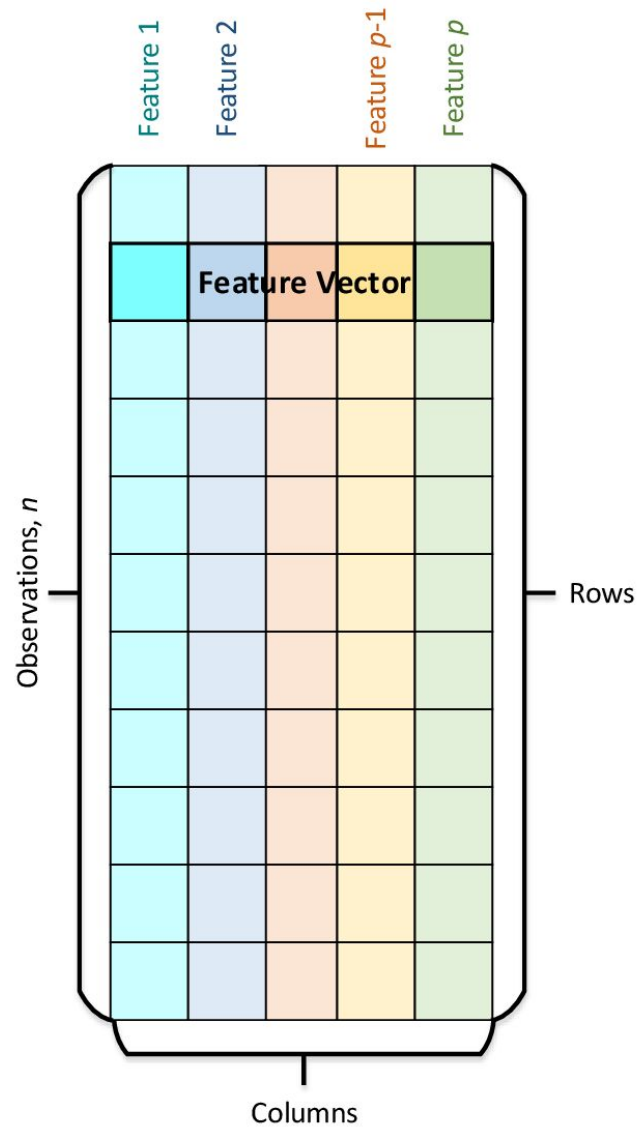| | ID | LastName | FirstName | Number |
|---|---|---|---|---|
| 1 | 1 | Johnson | Joe | 555-2323 |
| 2 | 2 | Lewis | Larry | NULL |
| 3 | 3 | Thompson | Thomas | 555-9876 |
| 4 | 4 | Patterson | Patricia | NULL |

ECONTENTA

KDnuggets Analytics/Data Science 2016 Software Poll, top 10 tools

Data Tools

ECONTENTA

# ДАННЫЕ

## Data Set

Feature 1, Feature 2, Feature p-1, Feature p

Feature Vector

Observations, n

Rows

Columns

| DOC NO | DOC DATE | DOC TYPE | QTY |
|--------|----------|----------|-----|
| REQ_01 | 01/APR/09 | REQUEST | 100 |
| REQ_02 | 02/APR/09 | REQUEST | 150 |
| ISS_01 | 03/APR/09 | ISSUE | 100 |
| ISS_02 | 04/APR/09 | ISSUE | 150 |
| ISS_03 | 04/APR/09 | ISSUE | 150 |
| | | | |
| | | | |
| | | | |
| | | | |

ECONTENTA

# Features types

- Quantitative
    - Continious
        - Real
        - Natural
        - Ratio, Percentage
        - Bounded
        - Interval
    - Discrete
        - Binary
        - Nominal
        - Ordinal

- Qualitative
    - Text
    - Images
    - Sound

ECONTENTA

Image gradients → Keypoint descriptor

(a)   (b)

|  | Document 1 | Document 2 | Document 3 | Document 4 | Document 5 | Document 6 | Document 7 | Document 8 |
|---|---|---|---|---|---|---|---|---|
| Term(s) 1 | 10 | 0 | 1 | 0 | 0 | 0 | 0 | 2 |
| Term(s) 2 | 0 | 2 | 0 | 0 | 0 | 18 | 0 | 2 |
| Term(s) 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Term(s) 4 | 6 | 0 | 0 | 4 | 6 | 0 | 0 | 0 |
| Term(s) 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Term(s) 6 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| Term(s) 7 | 0 | 1 | 8 | 0 | 0 | 0 | 0 | 0 |
| Term(s) 8 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |

← Word Vector (Passage Vector)

↑ Document Vector

ECONTENTA

# Спасибо!

## ...теперь немного практики