# Data Quality Services (DQS)
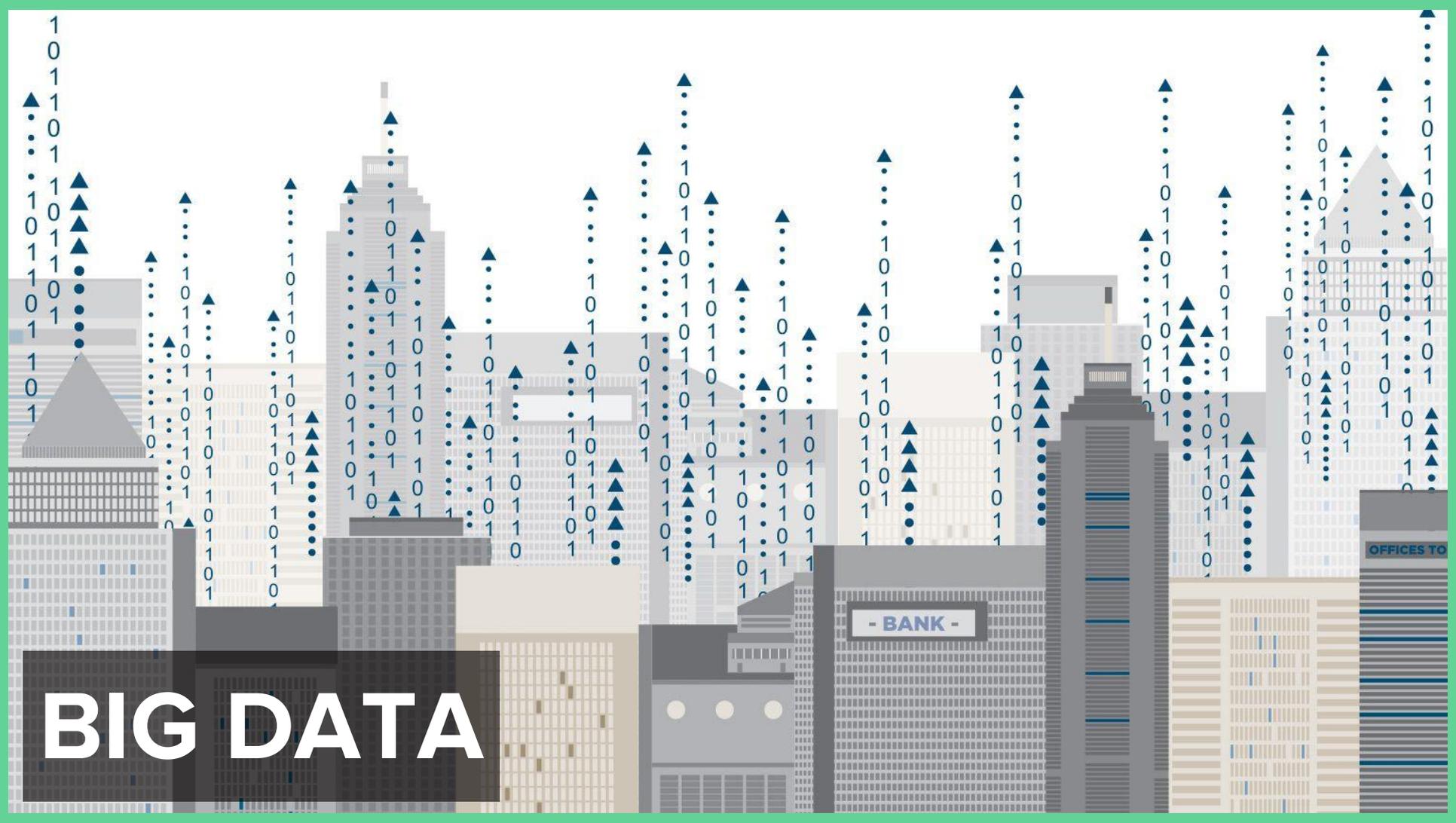
Integração de Dados e Warehousing

Pedro Henrique de Queiroz Lima

# Roteiro

- Motivação
- Qualidade de dados
  - O que é?
  - Aspectos
- Data Quality Services
  - Workflow
  - Limpeza de dados
  - Correspondência de dados
  - Profiling
  - Admnistração
- Recapitulando...
- Referências

# Motivação

BIG DATA

# 4.4ZB

1 Zetabyte = $10^{21}$ bytes

Quantidade de informação digital gerada acumulada até hoje (IDC/Forbes)

# 44ZB

1.7MB/s/pessoa

Projeção da quantidade de informação digital gerada acumulada até 2020

BUSINESS INTELLIGENCE

# US$ 3,1 tri/ano

EUA, 2015

Eram US$600 bi/ano, com dados de 2003 (IDC)

Causados por decisões de negócio baseadas em dados errados, retrabalho em TI e postagem de mala direta (US$ 600 bi nestes dois últimos).

# 15%

dos executivos confiam na qualidade geral dos dados

# 27%

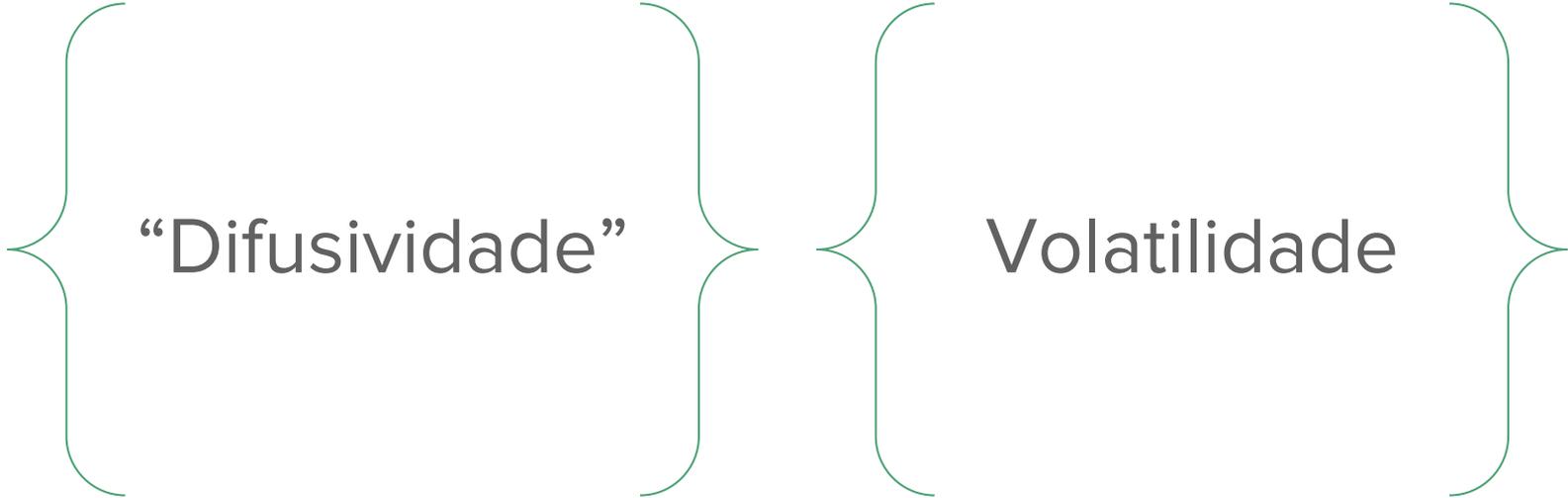dos executivos tem certeza da qualidade dos dados

# Qualidade de Dados

# O que é?

Qualidade de Dados

Adequação dos dados para as atividades de operações, decisões de negócios e planejamento da instituição.

# Dimensões de Qualidade de Dados

| Fator | Exemplo |
| --- | --- |
| Completude | 25% dos registros não contém "Sobrenome" |
| Conformidade | "Rua" e "r."; "avenida" e "Ave." ou "Av." |
| Consistência | Sexo representado por "0/1" e "M/F" |
| Precisão | Casas decimais necessárias ou dado atual |
| Validade | Altura: -15,5m |
| Unicidade | Registros devem representar entidades únicas |

# Dificuldades em Qualidade de Dados

{ "Difusividade" } { Volatilidade }

Problemas no design de software

# Validação de dados

Erros de entrada podem ser transferidos para a base de dados

# Fusão e aquisição de empresas

Possíveis duplicações, dados em formatos diferentes, incompletude de informações

# Formatação

Espaçamentos, abreviações, apelidos

| First Name | Middle Name | Last Name | Gender | Street | City | State | Country | Phone |
|------------|-------------|-----------|--------|--------|------|-------|---------|-------|
| Kimberly | B | Zimmerman | Female | 6040 S. Justine | CHICAGO | IL | USA | 123-555-0167 |
| Kim | B | Zimmerman | F | S. Justine | CHICAGO | IL | United States | (123) 555-0167 |

# Mudanças de atributos

Alterações de endereços, telefones

# Data Quality Services

# Microsoft® SQL Server® 2012

## Enterprise Information Management

Delivering Credible Consistent Data to Every Organization

# Enterprise Information Management

- Master Data Services (MDS)
- SQL Server Integration Services (SSIS)
- Data Quality Services (DQS)

Tira carga de trabalho do fluxo ETL, com interface simplificada

# Concorrentes

# Ferramentas

Monitoramento e rastreamento dos estados das atividades de qualidade e da qualidade de dados

Monitoramento

Limpeza

Correção, remoção ou enriquecimento de dados que estejam incorretos ou incompletos.

Análise dos dados de origem para prover entendimento da qualidade de dados e seus problemas.

Profiling

Correspondência

Identificação, ligação ou fusão de registros duplicados entre conjuntos de dados.

# Workflow

# SQL Server Data Quality Client

# Knowledge Base

- Amostras da base de dados
- Bases de dados de Referência (Azure Marketplace)
- Descoberta de conhecimento
- Domínios
  - Leading Values
  - Regras
- Base de conhecimento cresce com iterações e entradas de usuário

# Knowledge Discovery

# Domínios de dados

# Projeto de Qualidade de Dados

# Limpeza de Dados

# Limpeza de Dados

# Limpeza de Dados

# Limpeza de Dados

# Correspondência (Matching)

# Exemplo Matching

Dados

| Name | Address | Postal Code | City | State |
|------|---------|-------------|------|-------|
| Mag. Smith | 545 S Valley View D. # 136 | 34563 | Poughkeepsie | New York |
| Margaret smith | 545  Valley View ave unit 136 | 34563-2341 | Pughkeepsie | New-York |
| Maggie  Smith | 545 S Valley View Dr | | San Diego | CA |

Resultados do matching

| Rule | | Record Id | Cluster | Score | Name | Address | Postal Code | City | State |
|------|---|-----------|---------|-------|------|---------|-------------|------|-------|
| ⊟ | | 1000000 | 1000000 | | Mag. Smith | 545 S Valley View D. # 136 | 34563 | Poughkeepsie | New York |
| | ■ | 1000001 | 1000000 | 74% | Margaret smith | 545  Valley View ave unit 136 | 34563-2341 | Pughkeepsie | New-York |

# Exemplo Matching

Resolução de Correspondências

| Record matching drill-down | | | | |
|---|---|---|---|---|
| Fields Scores Contributions | | | | |
| Field | Weight | Matched record terms | Pivot record terms | Score contribution |
| Full Name | 0.2 | Margaret smith | Mag. Smith | 59% |
| Address Line 1 | 0.2 | 545  Valley View ave unit 136 | 545 S Valley View D. # 136 | 66% |
| Zip | 0.2 | 34563-2341 | 34563 | 63% |
| City | 0.2 | Pughkeepsie | Poughkeepsie | 81% |
| State | 0.2 | New-York | New York | 100% |

# Propriedades de correspondência

- Similaridade
- Peso
- Pré-requisito
- Pontuação mínima de correspondência

# Correspondência de dados

# Correspondência de dados

# Correspondência de dados

# Correspondência de dados

# Recapitulando...

- Remove carga do processo ETL
- Oferece ferramentas de validação e estatísticas sobre regras de domínio
- Pode limpar e remover duplicações
- Pode ser integrada ao SSIS com o DQS Transform
- Baseado em base de conhecimento
- Pode usar dados de terceiros para conhecimento

# Referências

1. Big Data: 20 Mind-Boggling Facts Everyone Must Read
2. Integrating and governing Big Data
3. Extracting Value from Chaos
4. The Cost of Poor Data Quality
5. Data Quality Services Demo