# Practical Machine Learning

## A Use Case for Data Cleansing and Object Identification in Market Research

Prof. Dr. Thomas Ruf (Kynetec Germany GmbH)

Practical Machine Learning

# A Use Case for Data Cleansing and Object Identification in Market Research

Prof. Dr. Thomas Ruf (Kynetec Germany GmbH)

# Outline

- Kynetec business overview

- Panel-based market research

- Current data production workflow

  o Preparation of incoming data for matching

  o Master data coding

  o Transactional data matching

- A Use Case for Machine Learning?!

# Outline

- **Kynetec business overview**

- Panel-based market research

- Current data production workflow

  ○ Preparation of incoming data for matching

  ○ Master data coding

  ○ Transactional data matching

- A Use Case for Machine Learning?!

# kynetec

# Global market research in animal health and agriculture

_____

As **global leaders in market research** for **animal health** and **agriculture**, Kynetec helps companies around the world understand the dynamics of their marketplace, turning research into business opportunities and enabling clients to create winning strategies. We conduct **tracking studies** for monitoring **market trends**, **customized research** for answering unique business challenges and provide **market forecasts** to support your long-term vision, bringing you closer to your customers.

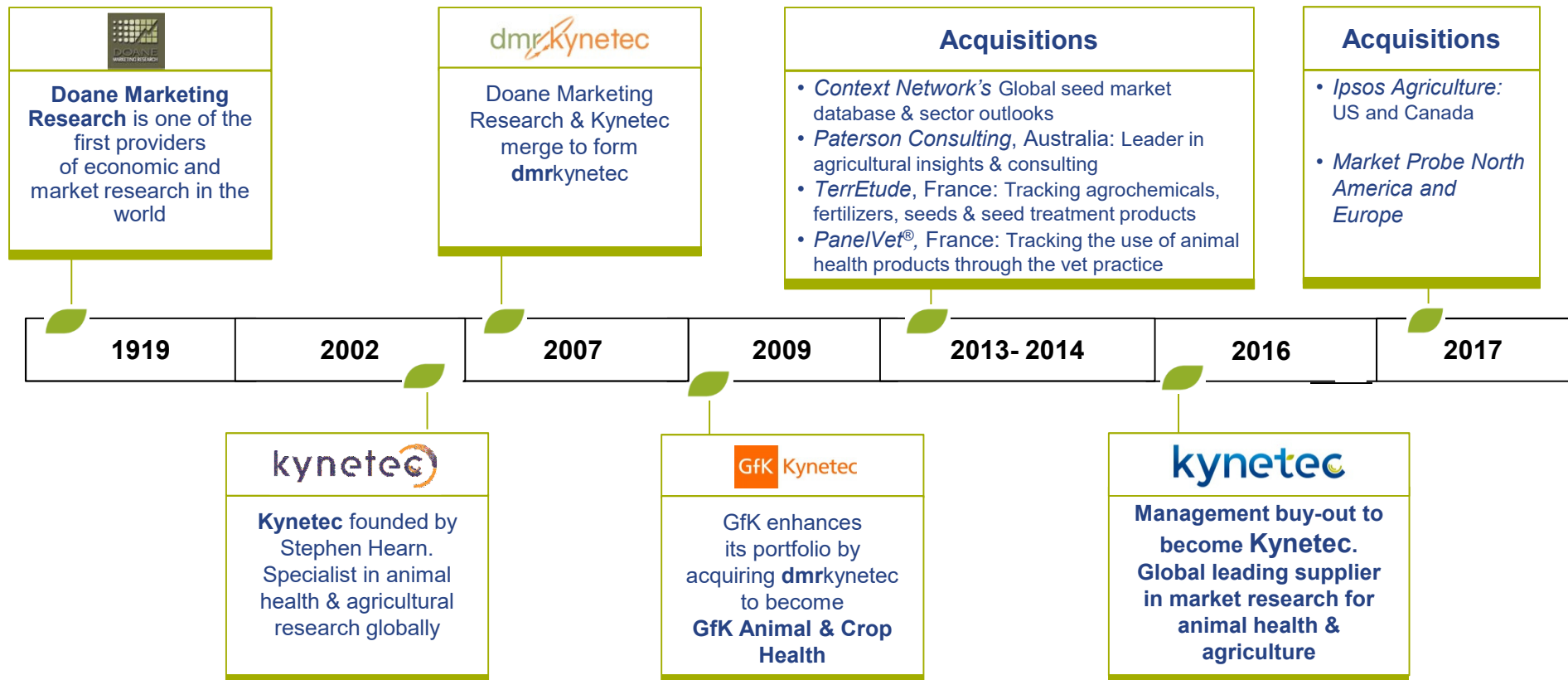| | | | |
|---|---|---|---|
| Conducting research in more than **80** countries | Connecting with **100 000s** vets, farmers and pet owners | Processing more than **1bn** data points annually | Global team of more than **650** research professionals |

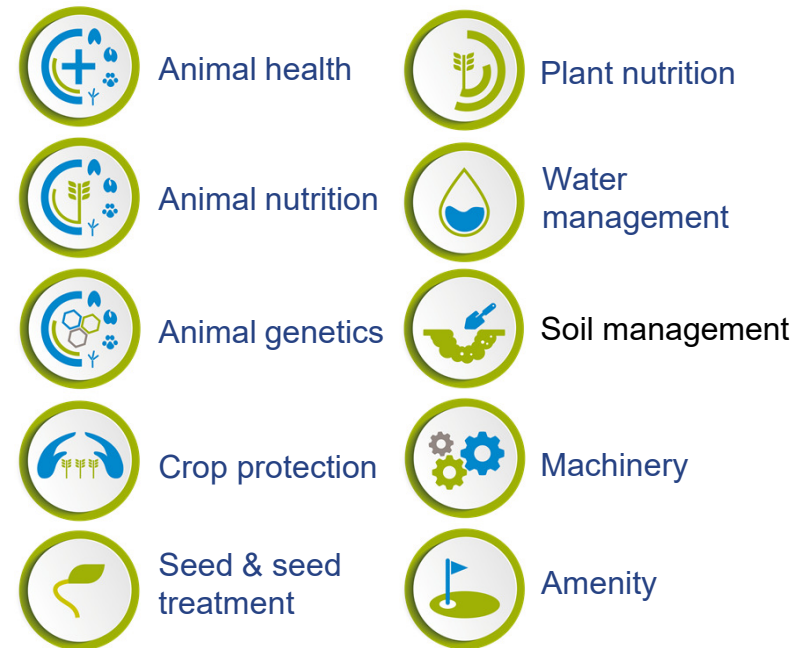# Long history in understanding animal health & agriculture

**Doane Marketing Research** is one of the first providers of economic and market research in the world

Doane Marketing Research & Kynetec merge to form **dmr**kynetec

### Acquisitions

- *Context Network's* Global seed market database & sector outlooks
- *Paterson Consulting*, Australia: Leader in agricultural insights & consulting
- *TerrEtude*, France: Tracking agrochemicals, fertilizers, seeds & seed treatment products
- *PanelVet®,* France: Tracking the use of animal health products through the vet practice

### Acquisitions

- *Ipsos Agriculture:* US and Canada

- *Market Probe North America and Europe*

| 1919 | 2002 | 2007 | 2009 | 2013- 2014 | 2016 | 2017 |
|------|------|------|------|------------|------|------|

**Kynetec** founded by Stephen Hearn. Specialist in animal health & agricultural research globally

GfK enhances its portfolio by acquiring **dmr**kynetec to become **GfK Animal & Crop Health**

**Management buy-out to become Kynetec.** **Global leading supplier in market research for animal health & agriculture**

# Delivering exceptional quality research in all sectors of our industry, covering all information needs

## Client challenges ➜ Our thinking

| | |
|---|---|
| Market opportunity & innovation | Message management: Resonance, reach, recall |
| Brand & customer experience | Forecasting |
| Digital strategy & intelligence | Launch readiness & evaluation |
| Value & pricing optimization | Segmentation & positioning |
| User experience | Concept testing |
| Retail sales tracking | Market dynamics & competitive landscaping |
| Product design optimization | Packaging design optimization |

## Sector expertise

- Animal health
- Animal nutrition
- Animal genetics
- Crop protection
- Seed & seed treatment
- Plant nutrition
- Water management
- Soil management
- Machinery
- Amenity

# Connecting with clients & markets globally

## Present in all major regions, and in most major markets



### Global Reach & local presence

- Conducting research in more than 80 countries
- Employees present in 22 major agriculture & animal health countries

### More than 650 talented professionals

- 300 research professionals
- 300 skilled interviewers
- >80 research partners

Australia    Canada    France    India    Italy    Malaysia    Poland    Spain    Thailand    UK    USA

Brazil    China    Germany    Ireland    Japan    Philippines    Russia    Switzerland    Turkey    Ukraine    Belgium

# Outline

- Kynetec business overview

→ - Panel-based market research

- Current data production workflow

  ○ Preparation of incoming data for matching

  ○ Master data coding

  ○ Transactional data matching

- A Use Case for Machine Learning?!

# Panel market research methodology

# In many regions, we provide market insight reports that cover key categories segmented by defined regions.

# Broad selection of syndicated studies to meet client needs


VetTrak™


VetTrak Insight™


VetTrak™
BDI/CDI Score


PetInsight™


SpecialtyTrak™


ParaTrak™


Consumer Perceived Recommendations


Pet Care Retail Trackers


VetTrak Forecaster™


VetTrak EU™


Sigma AH™


PressTrak™ Vet


Global Digital Trackers Studies/CEE


VetTrak™ Geo

# VetTrak™

## Overview

► Tracks product consumption (dispensed doses and revenue) and inventory levels for key parasiticides and vaccines in the veterinary channel

## Content

► Flea/tick, heartworm, GI worms, vaccines, NSAIDs, & Anti-infectives
► Metrics: Current Period, MAT, and YTD
► Includes up to10 years' worth of historical data

## Methodology

● US: survey-based clinic inventory data is blended with data extracted from a panel of practice management systems
● Mixed mode survey (n=~800) / PMS (n=~8,000)
● Projected to the universe of independent companion animal clinics (N=~25,000)

## Deliverables

● Monthly using i-map3™ platform
● Powerpoint summary report
● Excel based pivot table reports
● Available 20 working days after month end



## Benefits

- Provides ongoing tracking of key metrics needed to manage your brand in the vet channel
- Tracks sales out of vet clinic – doses and dollars

# Outline

- **Kynetec business overview**

- **Panel-based market research**

- **Current data production workflow**

  ○ Preparation of incoming data for matching

  ○ Master data coding

  ○ Transactional data matching
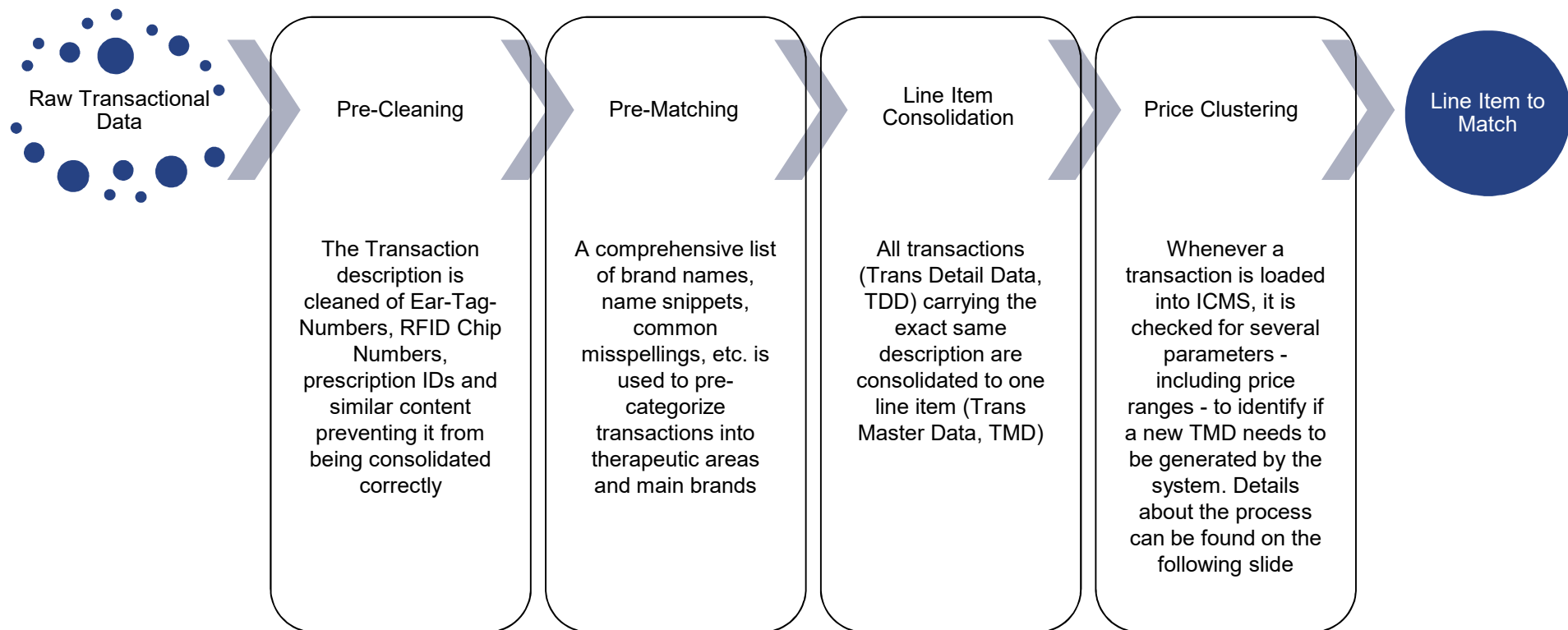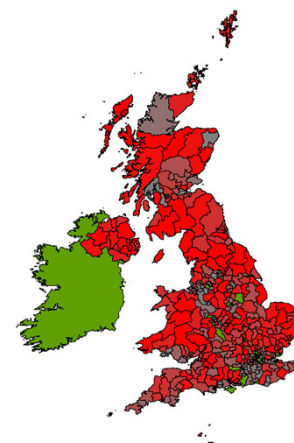
- **A Use Case for Machine Learning?!**

# Panel-based market research

## Simplified Data Production Workflow



Master data

Coding & Matching

Transaction data

# Data Production for Panel Research

## Key challenge: distinguish noise from data

### Nature of transactional data

Transactional data are sourced from incontrollable sources and may contain arbitrary references to products and services applied in a Vet practice or clinic

### Lining transactional data to master data

To be able to derive globally meaningful insights from transactional data, single transactions need to be linked to harmonizid master data

### Huge data volumes = huge problem?

In the work with big amounts of data, it is essential to minimize the need for manual work. For this reason, a specialized coding and matching solution is needed, supporting both master data management and the actual matching process. ICMS provided by Infacta is our tool of choice.

# Current data production system

## Integrated Coding and Matching System (ICMS)



- **ICMS developed by Infacta GmbH, Germany**

- **Originally created for GfK hospital panel**

- **Main competitive advantages of Infacta:**
  - 30+ years of in-depth panel production knowledge
  - Agile, pragmatic software development approach
  - Strategic technology partnerships
  - Emphasis on managing the production workflow
  - Specific component developments possible

- **No licence fees, compensation based on managed system usage**

## Outline

- Kynetec business overview

- Panel-based market research

- Current data production workflow

  o Preparation of incoming data for matching

  o Master data coding

  o Transactional data matching

- A Use Case for Machine Learning?!

## Status Quo

Transactional data of several data sources are the basis of the work in ICMS

| Source | # of Transactions |
|---|---|
| Provider 1 | 125.462.534 |
| Provider 2 | 56.847.183 |
| Provider 3 | 2.203.730 |
| Provider 4 | 490.590.998 |
| **Total** | **675.104.446** |

Matching close to 700mio transactions is not feasible in an efficient way, so we need to apply methods of data concentration

# Simplified Process of Cleaning pre ICMS

**Raw Transactional Data**

**Pre-Cleaning**

The Transaction description is cleaned of Ear-Tag-Numbers, RFID Chip Numbers, prescription IDs and similar content preventing it from being consolidated correctly

**Pre-Matching**

A comprehensive list of brand names, name snippets, common misspellings, etc. is used to pre-categorize transactions into therapeutic areas and main brands

**Line Item Consolidation**

All transactions (Trans Detail Data, TDD) carrying the exact same description are consolidated to one line item (Trans Master Data, TMD)

**Price Clustering**

Whenever a transaction is loaded into ICMS, it is checked for several parameters - including price ranges - to identify if a new TMD needs to be generated by the system. Details about the process can be found on the following slide

**Line Item to Match**

# Outline

- Kynetec business overview

- Panel-based market research

- Current data production workflow

  o  Preparation of incoming data for matching

  o  Master data coding

  o  Transactional data matching

- A Use Case for Machine Learning?!

# Example: Frontline combo spot-on, large dogs (6 x 2.68ml)

| | | |
|---|---|---|
| CEESA1 | H INSECTICIDES/ECTOPARASITICIDES (EARS EXCL. ANTIFUNGAL INCL.) | 23 categories |
| CEESA2 | H01 INSECTICIDES/ECTOPARASITICIDES (EARS EXCL. ANTIFUNGAL INCL.) PETS | 101 |
| CEESA3 | H01A INSECTICIDES/ECTOPARASITICIDES (EARS EXCL. ANTIFUNGAL INCL.) PETS Flea and ticks | 239 |
| CEESA4 | H01A03 INSECTICIDES/ECTOPARASITICIDES (EARS EXCL. ANTIFUNGAL INCL.) PETS Flea and ticks Drops | 285 |
| CEESA5 | H01A03A INSECTICIDES/ECTOPARASITICIDES (EARS EXCL. ANTIFUNGAL INCL.) PETS Flea and ticks Drops Dogs | 152 |
| ATC2 | C1 SML ANIMAL ECTOPARASITIC | 80 |
| ATC3 | C1A SML ANIMAL ECTOPARASITIC | 123 |
| ATC4 | C1A1 SML ANIMAL ECTOS-TOPICAL | 203 |
| Product group | FRONTLINE COMBO DOG | c. 6,000 product groups |
| Product | FRONTLINE COMBO SPOT ON LDOG 6X2.68ML | c. 24,000 SKUs |
| Manufacturer | MERIAL AH | c. 700 manufacturers |
| Form | SPOT ON | 128 formulations |
| Mode of admin. | TOPICAL | 21 modes of administration |
| Legal category | POM-V | 4 legal categories |



FRONTLINE COMBO SPOT ON LDOG 6X2.68ML
MAT Current Year GBP

# Attribute Management

## General Process (outlined in detail on following slides)



Definition of Product World

Set up and description of attributes per product world

Definition of inheritance model and association of attributes to layers

Accepted values for attributes to be inserted or reverse engineered from source data.

# Outline

- Kynetec business overview

- Panel-based market research

- Current data production workflow

  o Preparation of incoming data for matching

  o Master data coding

  o Transactional data matching

- A Use Case for Machine Learning?!

# Matching transactional data

## In many cases, transactional data doesn't come in a format ready to use.

For this reason, additional matching to meaningful master data is required to deliver insightful information from unharmonized data.



Selection of key functionality

Transactional data

Indicators for the quality of the proposed master data item

Master data items

# Matching transactional data

## Transactions window



| Description | Species | Value | Price | Units | GTIN | Procode1 | Procode2 | Created | Updated | Updated by | Comments |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Click here to define a filter | | | |
| $ALIM | CHIEN | 2.323 € | 24,71 € | 94 | | ME\|$ALIM | - | 21.07.2017 16:33:13 | | | |
| $ALIM | CHAT | 945 € | 20,10 € | 47 | | ME\|$ALIM | - | 21.07.2017 16:33:13 | | | |
| $ALIM | LAPIN | 38 € | 12,53 € | 3 | | ME\|$ALIM | - | 21.07.2017 16:33:14 | | | |
| $ALIM | | 48 € | 48,42 € | 1 | | ME\|$ALIM | - | 21.07.2017 16:33:19 | | | |
| $ALIM CRITAL CARE | LAPIN | 19 € | 18,67 € | 1 | | ME\|$ALIM | - | 21.07.2017 16:33:14 | | | |
| $MED | CHAT | 695 € | 19,30 € | 36 | | ME\|$MED | - | 21.07.2017 16:33:12 | | | |
| $MED | CHIEN | 3.391 € | 17,94 € | 189 | | ME\|$MED | - | 21.07.2017 16:33:13 | | | |
| $MED | CHIEN | 20 € | 20,00 € | 1 | | ME\|$ALIM | - | 21.07.2017 16:33:14 | | | |
| $SOINS | CHIEN | 105 € | 15,01 € | 7 | | AC\|$SOINS | - | 21.07.2017 16:32:59 | | | |

All relevant information of the transactional data can be displayed to enhance coding quality through more reference information

In addition to this, all columns can make use of filters and customized searches

# Matching transactional data

## Master Data window

| Item Master | Price correctness | 17 % | Score rating | ★ ★ ★ ★ ★ |
|---|---|---|---|---|

**Auto Search** | Custom Search

| GTIN | Marque | Animaux | SKU | | bricant | Ref-Preis | Score |
|---|---|---|---|---|---|---|---|
| 08713184055804 | NOBIVAC | CHIEN | NOBIVAC CHP SUSP INJ FL+SOLV. 50X1DOS | | MSD | 229,04 € | 11,75 |
| 08713184056801 | NOBIVAC | CHIEN | NOBIVAC CHP SUSP INJ FL+SOLV. 10X1DOS | | MSD | | 11,75 |
| 08714015018593 | DURAMUNE | CHIEN | DURAMUNE CHP LYO/SUSP FL 25X1 DOS | | ZOETIS | | 11,75 |
| 03660144008880 | CHP | BOVIN,CAPRIN,LAPIN,OVIN,PORCIN,VOLAILLE | CHP PDR ORALE SEAU 1X2.5KG | | COOPHAVET | | 11,59 |

*Matching transactional data*

The auto search is triggered whenever a transactional data entry is selected. It displays the best results of the search using google full text search algorithms.

Columns of additional reference data can be selected based on the full list of master data attributes. Additionally data based on recent translations of transactions to master data items such like reference price or date of latest link are available

An indicator of the correctness of the price of transaction versus the master data item is available to be able to further qualify translations. In this case the quality it too low to establish a direct link.

Also a quality indicator of the strength of the overlap of the full text search is provided. Here first experiences need to show if an automated matching from a score threshold could be implemented to minimize manual work.

# Matching transactional data

## Data Source Investigation



A click on Raw Data shows all transactions with the same transaction text. Systems that are based on invoice data can directly show the scan of the invoice with a click on show invoice.

# Outline

- Kynetec business overview

- Panel-based market research

- Current data production workflow

  o Preparation of incoming data for matching

  o Master data coding

  o Transactional data matching

- A Use Case for Machine Learning?!

# HoloClean
## Repair as an Inference Problem

- Blogpost and VLDB 2017 paper
- Open source  www.holoclean.io

HoloClean:
Weakly Supervised Data Repairing

**HoloClean: Holistic Data Repairs with Probabilistic Inference**

Theodoros Rekatsinas[*], Xu Chu[†], Ihab F. Ilyas[†], Christopher Ré [*]
{thodrek, chrismre}@cs.stanford.edu, {x4chu, ilyas}@uwaterloo.ca
[*] Stanford University and [†] University of Waterloo

epairing [30].
6, 49] on dif-
t (i) their av-
on and recall)
hese methods
did not perform any correct repairs. This is because these methods
sources, with quantitative data repairing methods, which leverage      limit themselves to only one of the aforementioned signals, and ig-
statistical properties of the input data. Given an inconsistent dataset   nore additional information that is useful for data repairing. We
as input, HoloClean automatically generates a probabilistic pro-
gram that performs data repairing. Inspired by recent theoretical

UNIVERSITY OF WATERLOO

# HoloClean [VLDB'17]

## Repair as an Inference Problem

### Input

**Dataset to be cleaned**

| | DBAName | Address | City | State | Zip |
|---|---|---|---|---|---|
| t1 | John Veliotis Sr. | 3465 S Morgan ST | Chicago | IL | 60608 |
| t2 | John Veliotis Sr. | 3465 S Morgan ST | Chicago | IL | 60609 |
| t3 | John Veliotis Sr. | 3465 S Morgan ST | Chicago | IL | 60609 |
| t4 | Johnnyo's | 3465 S Morgan ST | Cicago | IL | 60608 |

**Denial Constraints**

c1: DBAName $\rightarrow$ Zip
c2: Zip $\rightarrow$ City, State
c3: City, State, Address $\rightarrow$ Zip

**Matching Dependencies**

m1: Zip = Ext_Zip $\rightarrow$ City = Ext_City
m2: Zip = Ext_Zip $\rightarrow$ State = Ext_State
m3: City = Ext_City $\wedge$ State = Ext_State $\wedge$
$\wedge$ Address = Ext_Address $\rightarrow$ Zip = Ext_Zip

**External Information**

| Ext_Address | Ext_City | Ext_State | Ext_Zip |
|---|---|---|---|
| 3465 S Morgan ST | Chicago | IL | 60608 |
| 1208 N Wells ST | Chicago | IL | 60610 |
| 259 E Erie ST | Chicago | IL | 60611 |
| 2806 W Cermak Rd | Chicago | IL | 60623 |

### The HoloClean Framework

**1. Error Detection Module**

- Use integrity constraints
- Leverage external data
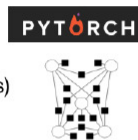- Detect outliers
- Identify possible repairs

**2. Compilation Module**

- Automatic Featurization
- Statistical analysis and candidate repair generation
- Compilation to factors/tensors

**3. Repair Module**

PYTORCH

- Ground probalistic model
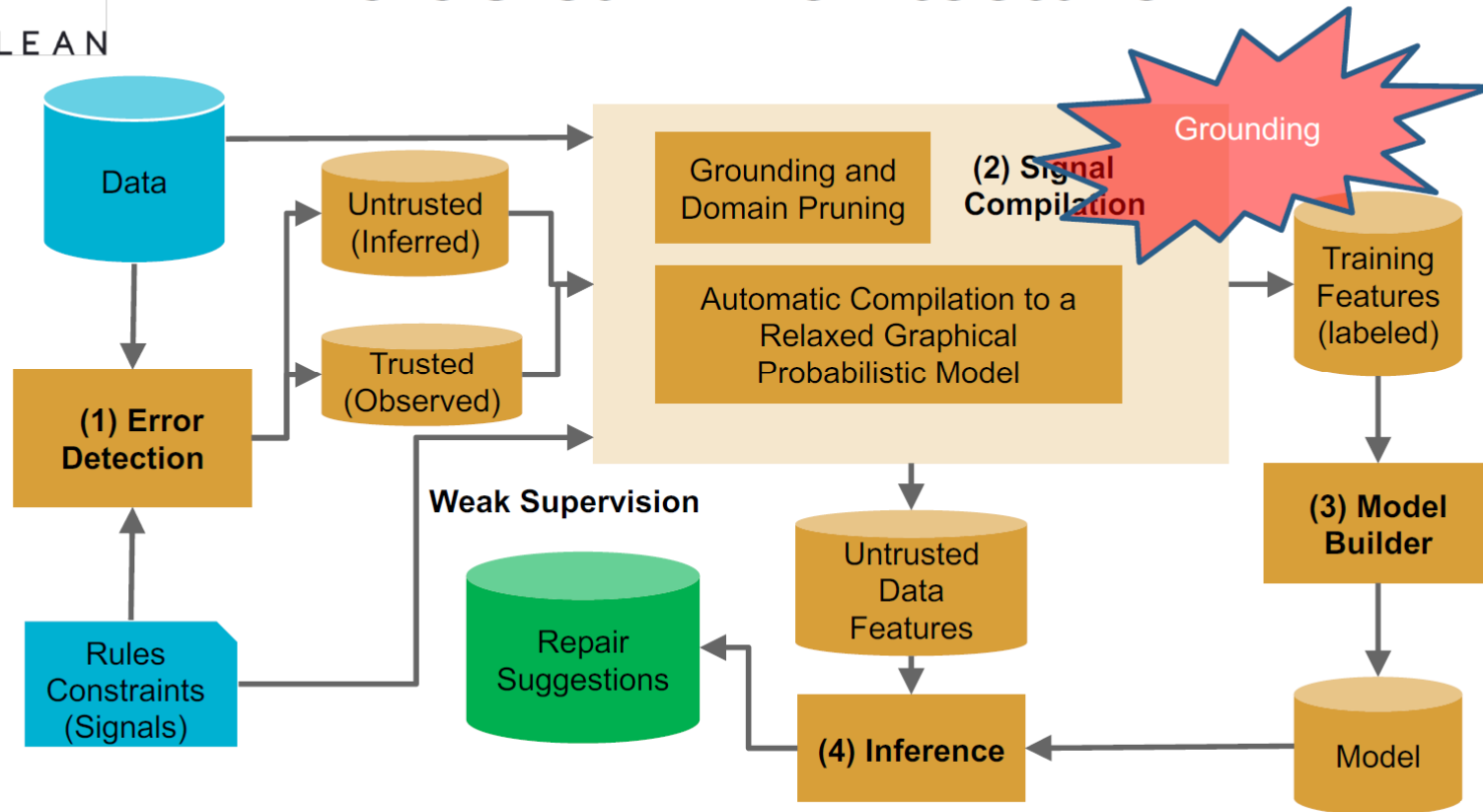- Statistical learning (weights)
- Probabilistic inference

### Output

**Proposed Cleaned Dataset**

| | DBAName | Address | City | State | Zip |
|---|---|---|---|---|---|
| t1 | John Veliotis Sr. | 3465 S Morgan ST | Chicago | IL | 60608 |
| t2 | John Veliotis Sr. | 3465 S Morgan ST | Chicago | IL | **60608** |
| t3 | John Veliotis Sr. | 3465 S Morgan ST | Chicago | IL | **60608** |
| t4 | **John Veliotis Sr.** | 3465 S Morgan ST | **Chicago** | IL | 60608 |

**Marginal Distribution of Cell Assignments**

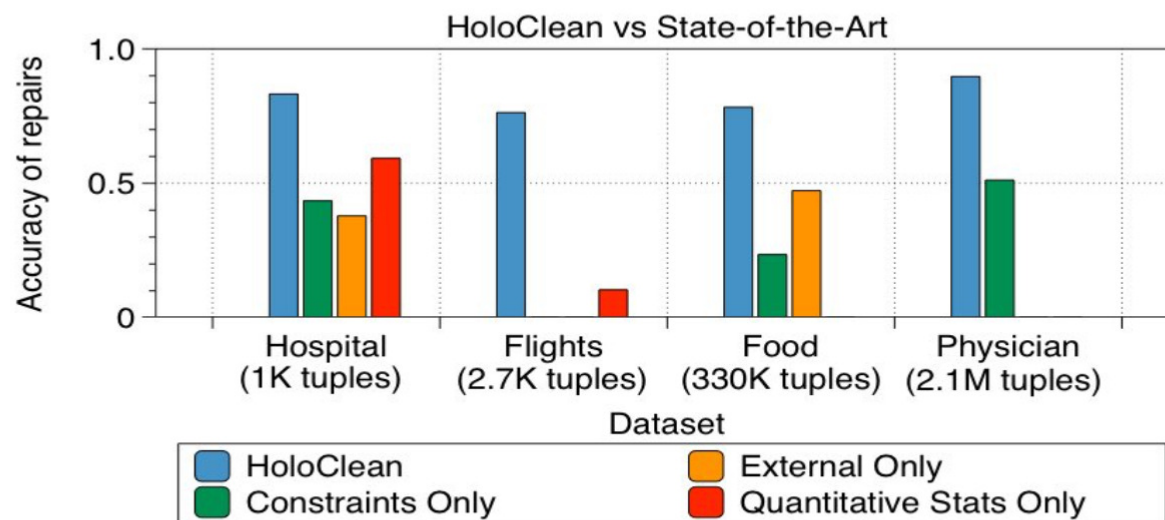| Cell | Possible Values | Probability |
|---|---|---|
| t2.Zip | 60608 | 0.84 |
| | 60609 | 0.16 |
| t4.City | Chicago | 0.95 |
| | Cicago | 0.05 |
| t4.DBAName | John Veliotis Sr. | 0.99 |
| | Johnnyo's | 0.01 |

# HoloClean Architecture

# HoloClean Key Ideas

- ***Domain Pruning***: Limit the active domain for random variable
  - Values that appear in the table
  - Use co-occurrence to prune even more aggressively
  - Borrow ideas from missing value imputation

- ***Tying Weights:*** *Learn one weight per constraint not per factor (Templated MLNs)*

- ***Constraint Relaxation***: Relax constraints over sets of random variables to features over independent random variables.

# HoloClean Results



HoloClean vs State-of-the-Art

**HoloClean:** our approach combining all signals and using inference
**Holistic[Chu,2013]:** state-of-the-art for constraints & minimality
**KATARA[Chu,2015]:** state-of-the-art for external data
**SCARE[Yakout,2013]:** state-of-the-art ML & qualitative statistics

# Results of the application of HoloClean to Kynetec data matching