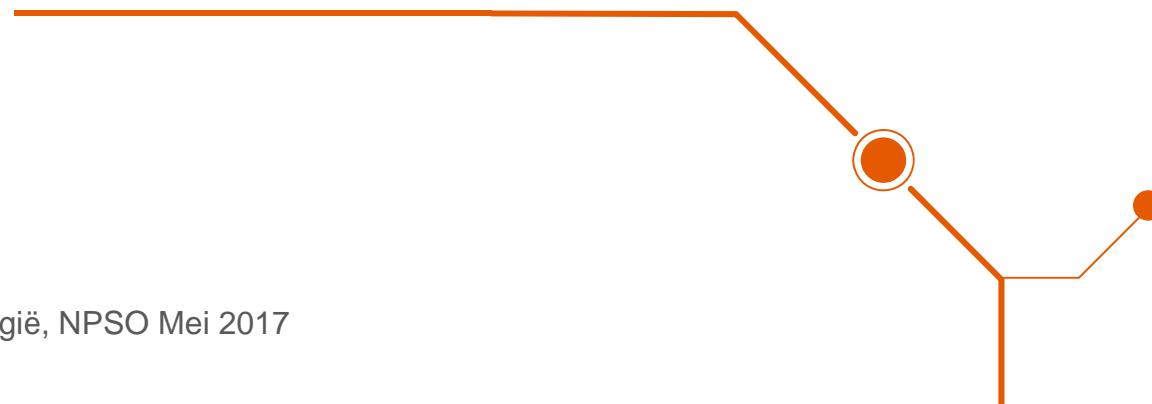




Big Data en Marktonderzoek: Opportuniteiten en gevaren



Dr. Istvan Hajnal, GfK België, NPSO Mei 2017

@IstvanHajnal



Big Data en het NPSO



Volgens de website

- ... Big data speelt in toenemende mate een rol in beleidsvorming en zeker in ook de survey wereld. Nieuwe toepassingen om gevoelens, opinies en meningen van personen te meten, zoals opinion mining, online sentiment expression, gebruik van Twitter, etc. worden in toenemende mate gebruikt.
- ... vragen rondom het nut en de bruikbaarheid van het gebruik van big data.
- **Ambigue houding, die je ook terugvindt bij marktonderzoek**



Is dit terecht? En is dit altijd zo geweest?





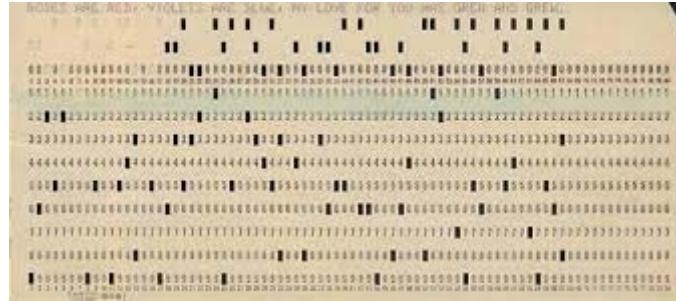




Herman Hollerith

(bron: wikipedia)

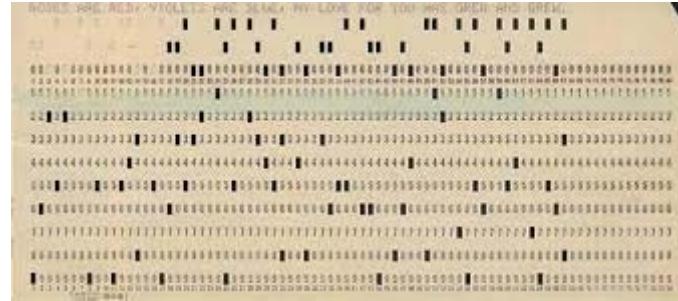
- 1860-1929
- De eerste baan van Hollerith was bij het bureau dat de volkstelling in de VS van 1880 moest uitvoeren
- Ontwikkelde een machine die het met behulp van ponskaarten mogelijk maakte grote hoeveelheden gegevens geautomatiseerd te verwerken.
- Hollerith heeft de Tabulating Machine Company opgericht
 - Later werd dit de Computer Tabulating Recording Company
 - En nog later werd de naam veranderd in International Business Machines Corporation (IBM).



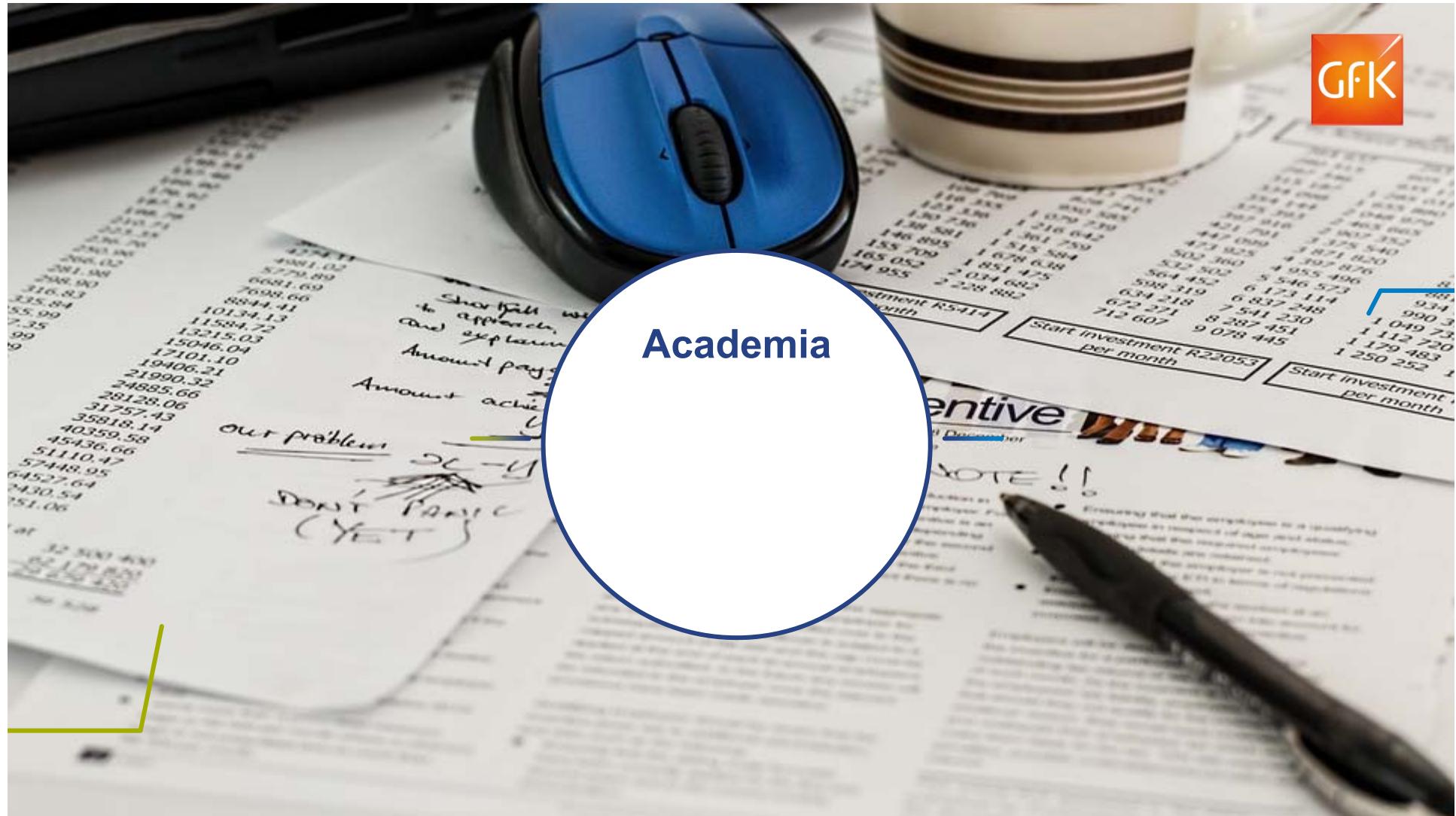
Herman Hollerith

(bron: wikipedia)

- 1860-1929
- De eerste baan van Hollerith was bij het bureau dat de volkstelling in de VS van 1880 moest uitvoeren
- Ontwikkelde een machine die het met behulp van ponskaarten mogelijk maakte grote hoeveelheden gegevens geautomatiseerd te verwerken.
- Hollerith heeft de Tabulating Machine Company opgericht
 - Later werd dit de Computer Tabulating Recording Company
 - En nog later werd de naam veranderd in International Business Machines Corporation (IBM).



Conclusie: De IBM Mainframe, de PC, Watson, e.d. komen, met wat goede wil, van een spin-off van Official Statistics



Academia

Alan Turing

- Computer pionier
- De Turingmachine
- Ontcijfering Enigma
- Princeton – Cambridge - Manchester



Academia

Geoffrey Hinton

- 1947-
- Cognitieve psycholoog en computerwetenschapper
- Universiteit van Toronto
- Expert in neurale netwerken
- Grote invloed op artificiële intelligentie en “deep learning”
- En dus op de verwerking van Big Data





Markt Onderzoek bureaus



Marktonderzoeksbureaus

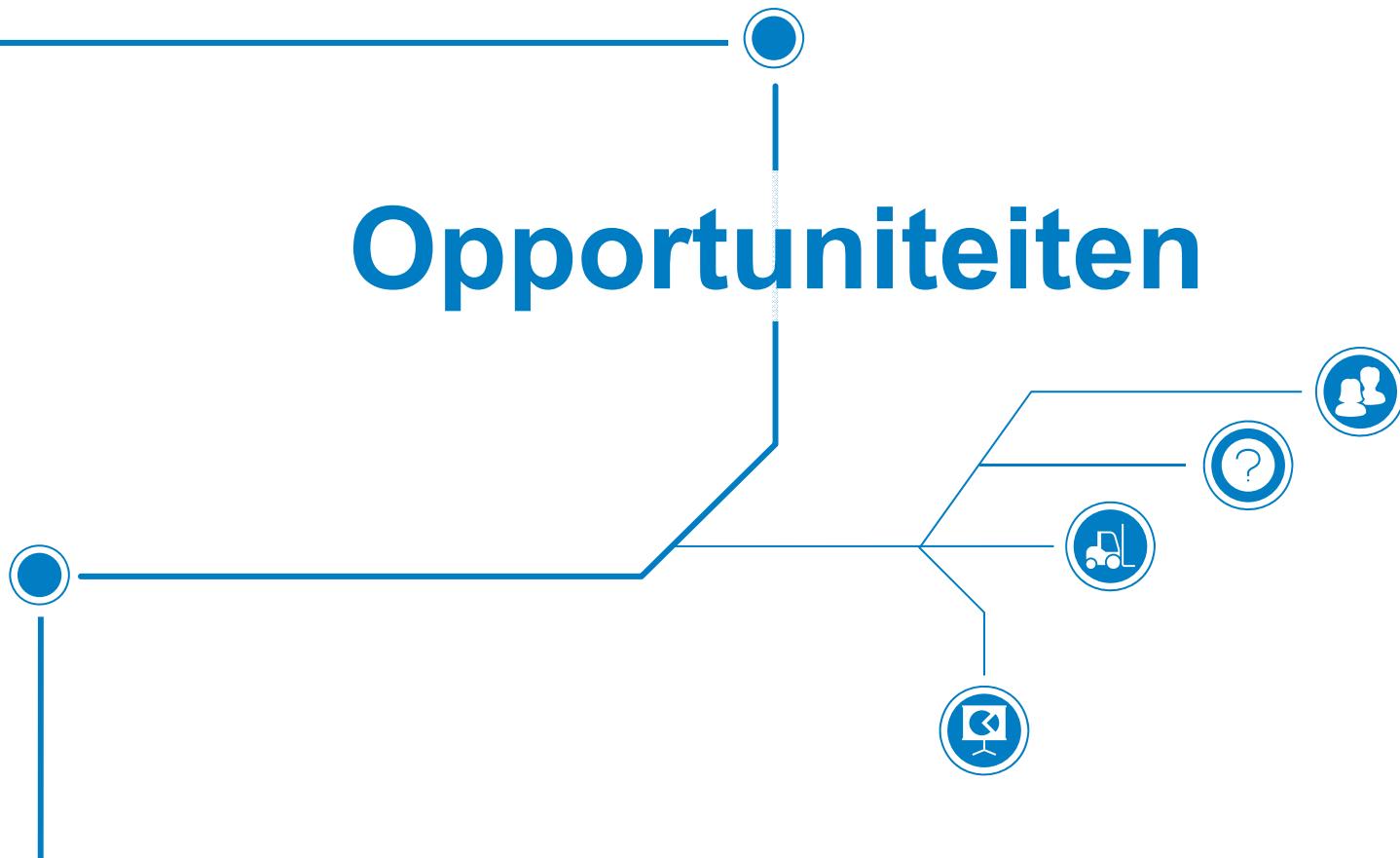
Weinig gekend feiten:

- Het eerste privé-bedrijf dat ooit een computer bestelde was een marktonderzoekbureau!
- In de jaren 60 waren marktonderzoekers pionieurs in de eerste programmeertalen





Opportunitäten





Big Data & Data Science **@GfK**

Taken from “How to Surf a Data Lake without Drowning?” a presentation by Ralph Wirth and Anna Machens on Predictive Analytics World London, 29th October 2015



As a global market research company...



...and via countless studies on diverse topics and for many industries



In an ideal world, Data Scientists could easily identify all relevant data assets for specific research questions and directly access them



Example: Potential predictors in GfK data for sales of Jacobs filter coffee



With a multi-national, heterogeneous data landscape, this is not a given

Data storage is not fully integrated, no metadata taxonomy covers all data facets

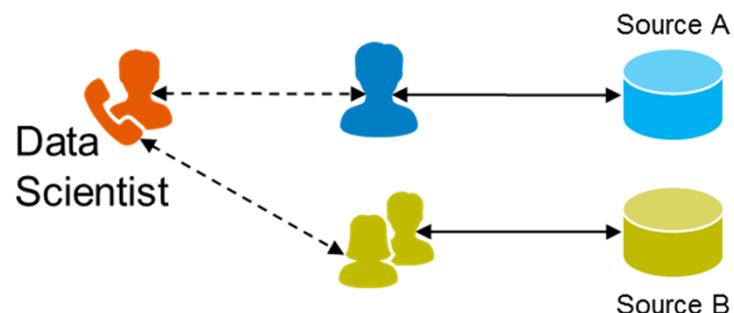
Imperfect data transparency

In some situations, a lot of communication is needed to find out what suitable data we have for a specific research problem.

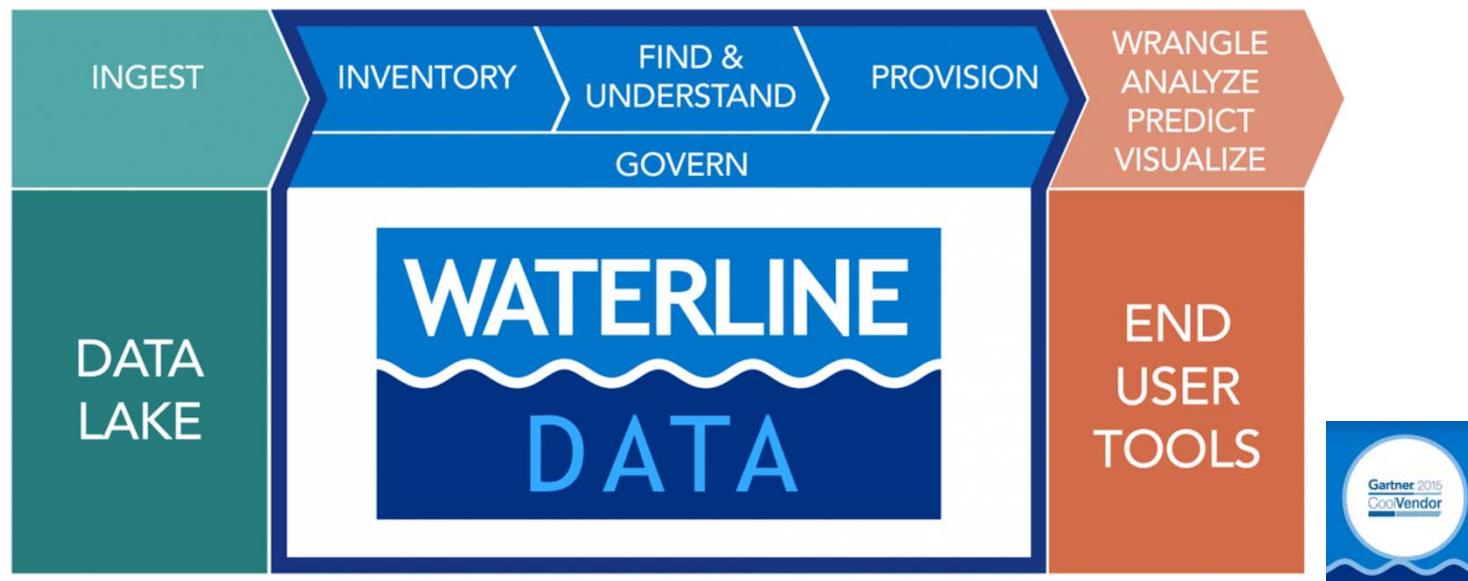


No direct access to full data landscape

Some data is stored in silos, so that access has to be handled by the respective data managers.



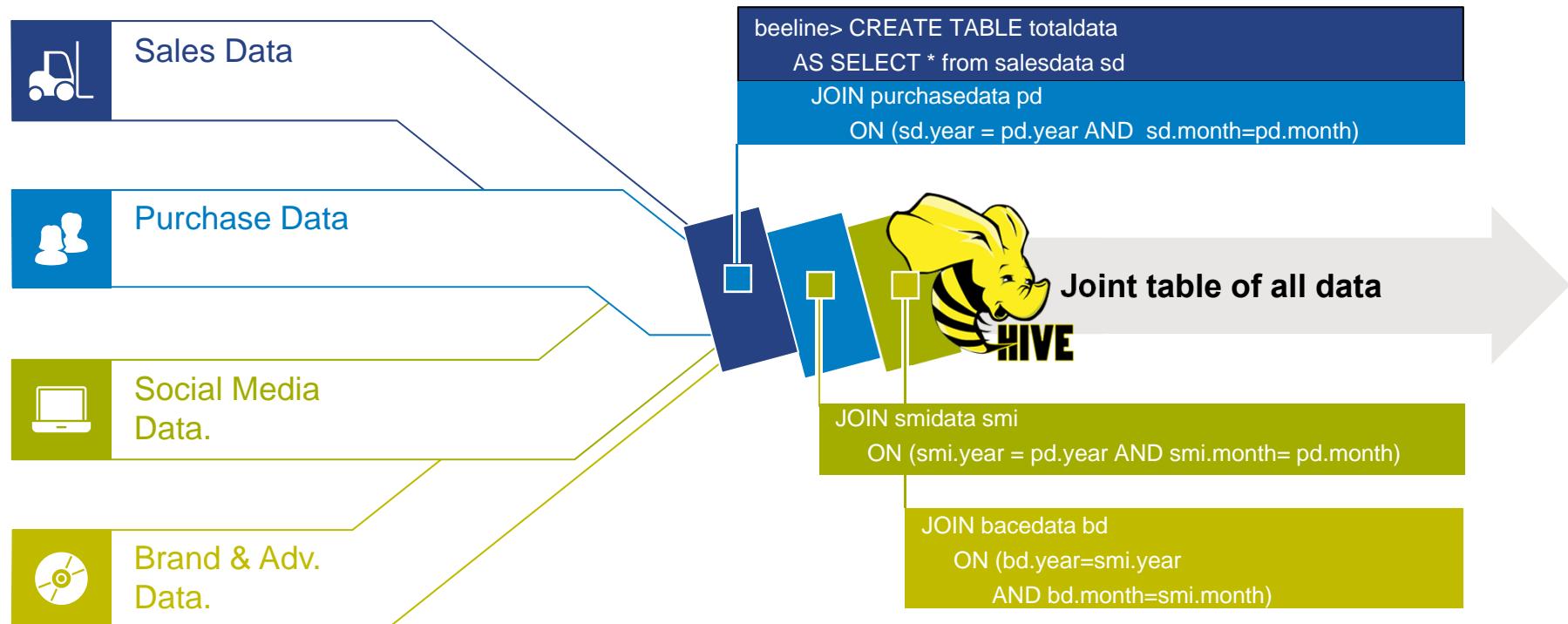
We test and use metadata management tools to find relevant data and get a first understanding of data sets



Hive is used to aggregate and integrate data from different sources with regard to the time dimension



Integration





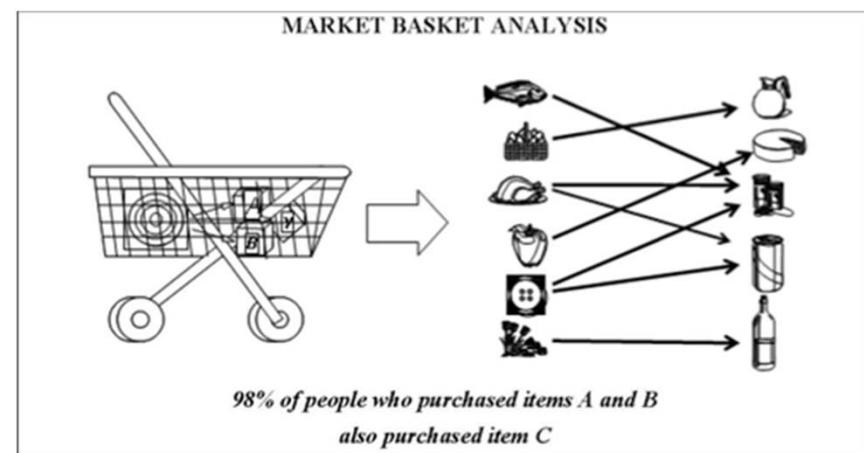
Big Data & Data Science @GfK Belgium

A sample of cases where the Marketing and Data Sciences team of GfK Belgium was involved in

Statistics and Methodology

Association Rules Mining

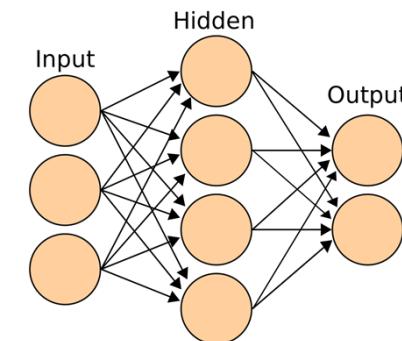
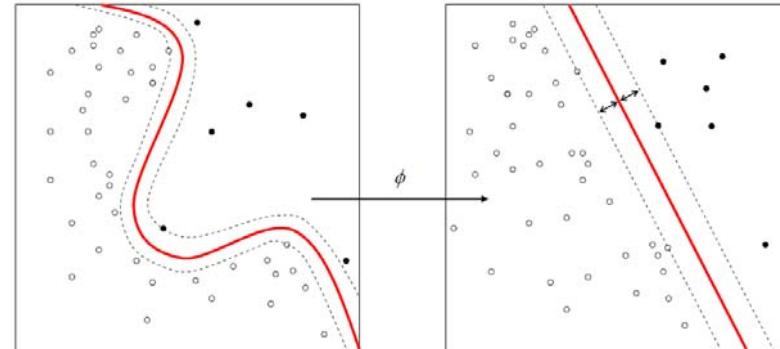
- **Association Rules Mining** is a method for discovering interesting relations between variables by identifying strong association rules in large databases.
- Association rules were introduced for discovering regularities between products in point-of-sale systems in supermarkets.
- We have applied this on Shopper data to try an identify presence of beer in a shopping basket



Statistics and Methodology

Machine Learning

- **Machine Learning** deals with algorithms that can learn from and make predictions on data. While ML is a subfield of computer science it is heavily influenced by statistics.
- Examples of such techniques are:
 - Ridge + Lasso
 - Random Forests
 - Support Vector Machines
 - AdaBoost
 - ...
- Applied on predicting presence of certain product categories in shopping baskets





Data handling

Analysis on large GfK databases

- Analyses on Media Measurement TAM data
- Ongoing proof-of-concept
- Using the Big Data infrastructure of the DataLab

GfK Data Lab - MDS Sandbox

HIVE Query Editors Metastore Manager Workflows Zoek File Browser Job Browser cobeoga

Hive Editor Query Editor My Queries Saved Queries Geschiedenis

Assist Settings

DATABASE default

Table name... airports be_port_data ext_mbrunsch_... hyperlane_de_u... page_view sample_07 sample_08 test test_sq url_categorization weather_sentiment

select * from airports where type = "small_airport";

Execute Opslaan Save as... Explain or create a New query

Recent queries Query Log Columns Results Chart

	airports.id	airports.ident	airports.type	airports.name	airports.latitude_deg	airports.longitude_deg	airports.altitude_ft
0	6524	"00AK"	"small_airport"	"Lowell Field"	59.94919967651367	-151.6959991455078	45
1	6525	"00AL"	"small_airport"	"Epps Airpark"	34.86479949951172	-86.77030181884766	82
2	6527	"00AZ"	"small_airport"	"Cordes Airport"	34.305599212646484	-112.16500091552734	36
3	6528	"00CA"	"small_airport"	"Goldstone /Gts/ Airport"	35.35049819946289	-116.88800048828125	30
4	6529	"00CO"	"small_airport"	"Cass Field"	40.62200001220703	-104.34400177001953	48
5	6531	"00FA"	"small_airport"	"Grass Patch Airport"	28.64550018310547	-82.21900177001953	53
6	6533	"00FL"	"small_airport"	"River Oak Airport"	27.230899810791016	-80.96920013427734	32
7	6534	"00GA"	"small_airport"	"Lt World Airport"	33.76750183105469	-84.06829833984375	70
8	6537	"00ID"	"small_airport"	"Delta Shores Airport"	48.145301818847666	-116.21399688720703	20
9	6539	"00IL"	"small_airport"	"Hammer Airport"	41.97840118408203	-89.5604019165039	84
10	6541	"00IS"	"small_airport"	"Hayenga's Cant Find Farms Airport"	40.02560043334961	-89.1229019165039	82

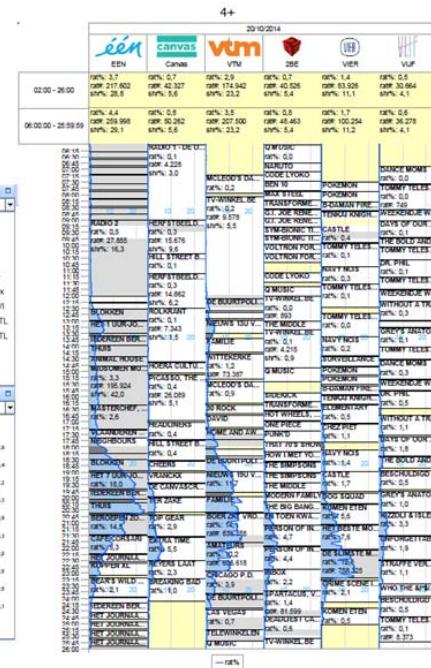
Data handling

Integration of GfK data sources

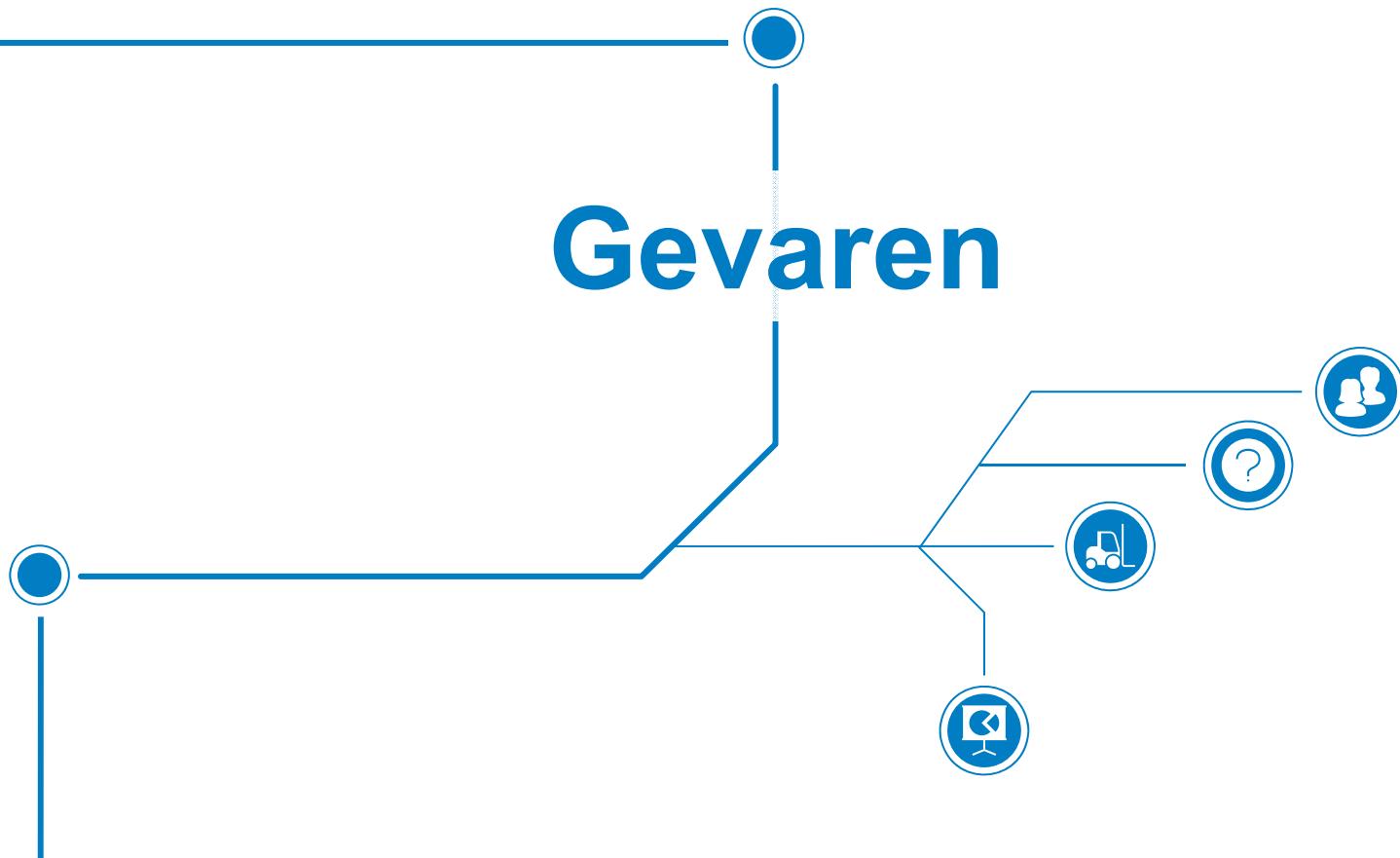
- **Data Fusion**
- In cooperation with Germany
- Example: tgo.be: Data Fusion of Media Measurement (TAM) data and Shopper (CP) data
- Not available yet

MAP software

Example output



Gevaren





Gevaren

In het kort ... Bias, Bias & Bias



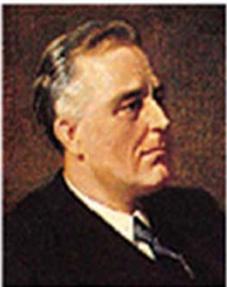


Maar eerst een opmerking vooraf

IH op IIEX 2016

I'm not buying the zillions of datapoints argument: it's not because there is a lot of water in the ocean that you need to drink it

1936 US Presidential Elections



Franklin Roosevelt
(Democrat)



Alf Landon
(Republican)

- Literary Digest
 - Sample size: 2.4 MILLION
 - Prediction: Roosevelt 43% : Landon 57%
- George Gallup
 - Sample size: 50,000
 - Prediction: Roosevelt 56% : Landon 44%

Steekproefgrootte is belangrijk

Representativiteit nog meer

AMERICA SPEAKS
THE NATIONAL WEEKLY POLL of PUBLIC OPINION

Today Election Forecast **Next Sunday** The Election in Review

Institute Forecasts the Re-election of Franklin D. Roosevelt, Gives Him 54% of Popular Vote, Minimum of 315 Electors

Major Party Percent
Is 55.7; New York in F.D.R. 'Sure' Column

By CLARENCE COLE
Director, American Institute of Public Opinion

N.Y. CITY—The Institute's latest presidential poll indicates that Roosevelt will receive approximately 55.7% of the major party vote (minor parties eliminated), or 48% for Alfred M. Landon and Frank Banks. In 1932 the President received 58.1% of the major party vote.

The reporting period, New York City, the Institute's headquarters, with a population of 7,500,000, was 100% included in the survey. It is estimated that 1,000,000 persons in the city will be eligible to vote in November, depending on how the registration lists are revised.

For the last week, one percent of the city's eligible voters were interviewed daily. The results of the first week of interviewing are as follows:

Category	Percentage
Major Party	55.7
Minor Parties	4.3
Other	40.0

The number of adults interviewed in the 1932 election was 100,000, and by general agreement it is believed that the number of adults in all the country is approximately 100,000,000. The Institute's latest presidential forecast is based on the assumption that each adult in the country has the same voting potential as each adult in the city.

Lineups of States

HIGHLIGHTING of the possibilities of victory, the Roosevelt administration has a strong lead in 37 states, according to the latest forecast of the American Institute of Public Opinion. "It is probable that he will defeat his principal opponent, Franklin D. Roosevelt, in 37 states," says the forecast. "He will win in 10 states, and will lose in 13 states." The forecast also indicates that Roosevelt will receive 315 electoral votes, while Landon will receive 87.

Election Forecast

- 1.—The American Institute of Public Opinion predicts the re-election of Franklin D. Roosevelt and John N. Garner.
- 2.—The Institute's latest presidential poll indicates that Roosevelt will receive approximately 55% of the major party vote (minor parties eliminated), or 48% for Alfred M. Landon and Frank Banks. In 1932 the President received 58.1% of the major party vote.
- 3.—With minor parties included, President Roosevelt's percentage of the total popular vote will be approximately 54%, or 49% for Landon.
- 4.—The President will receive a minimum of 315 electoral votes. The question is, will he win in 37? Should last-minute shifts in the group of states where the race is nip-and-tuck give this entire group to Roosevelt, he would receive more electoral votes than in 1932, when he polled 471.
- 5.—William Lemke, candidate of the Union Party, will poll fewer than 1,000,000 popular votes, and carry no state.
- 6.—Norman Thomas, Socialist candidate, will poll about half as many votes as in 1932, when he received 500,000.

Election Will Test Clashing Poll Methods

By CLARENCE COLE
Director, American Institute of Public Opinion

N.Y. CITY—The 1936 election will test the clashing methods of the various political organizations in their efforts to predict the outcome of the presidential election.

The American Institute of Public Opinion has been forecasting the election since 1932, and has been doing so with considerable success. The Institute's latest forecast indicates that Roosevelt will receive 55.7% of the major party vote, or 48% for Landon and Banks.

The Institute's method of forecasting the election is to interview 1,000,000 adults in the city of New York, which is 100% of the city's eligible voters. The results of the first week of interviewing are as follows:

Category	Percentage
Major Party	55.7
Minor Parties	4.3
Other	40.0

The number of adults interviewed in the 1932 election was 100,000, and by general agreement it is believed that the number of adults in all the country is approximately 100,000,000. The Institute's latest presidential forecast is based on the assumption that each adult in the country has the same voting potential as each adult in the city.

Lineups of States

HIGHLIGHTING of the possibilities of victory, the Roosevelt administration has a strong lead in 37 states, according to the latest forecast of the American Institute of Public Opinion. "It is probable that he will defeat his principal opponent, Franklin D. Roosevelt, in 37 states," says the forecast. "He will win in 10 states, and will lose in 13 states." The forecast also indicates that Roosevelt will receive 315 electoral votes, while Landon will receive 87.



Geldt zeker ook voor Big Data

Nog vaak verkeerd begrepen bij Big Data gebruikers

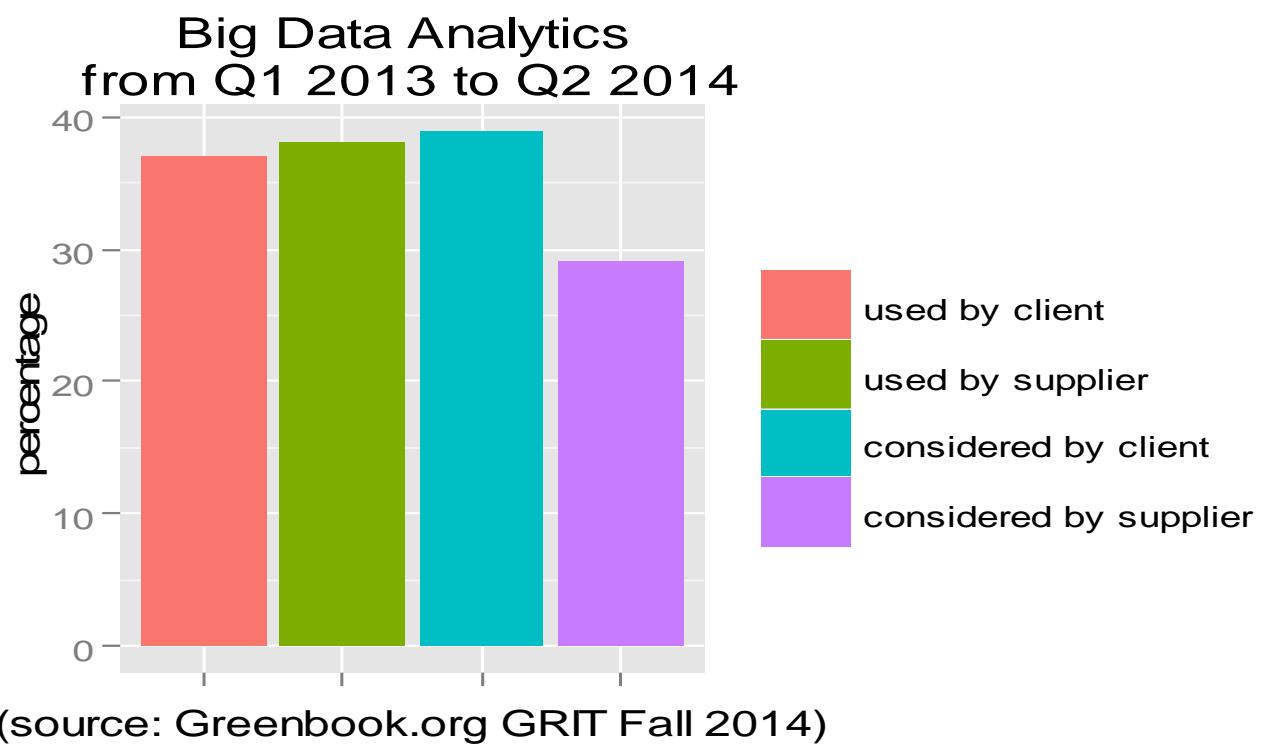
“Of course samples can and are biased as well. But there is a difference: Samples are constructed specifically with a research question in mind, and often are designed to be unbiased. Big data or other sources of data are often created for other reasons than research questions. As a consequence big data might have some disadvantages that are not offset by its bigger size”

Zie ook <http://allthingsdatascience.blogspot.be/2016/12/small-samples-versus-alternative-big.html> over een discussie rond kijkcijfers in België

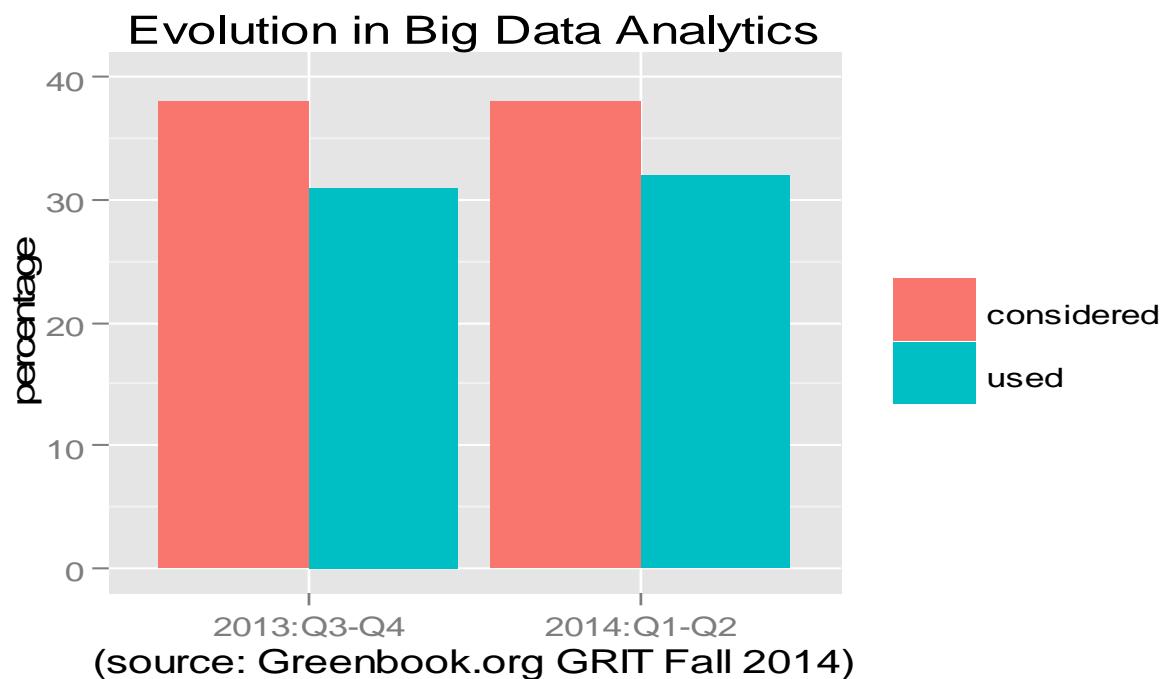
Conclusies voor marktonderzoek



Big Data: today



Big Data: evolution



Big Data fatigue?



“The Promise and Peril of Big Data”

*“Why Big Data Will Never Replace
Market Research”*

“Navigating The Big Data Hype”

(source: Greenbook)



Mijn conclusies

IIEX 2015

MRX seems to be bipolar when it comes to Big Data

My advice to suppliers of Market Research:

- Don't look at Big Data as just a fad or hype
- Don't look at Big Data as a threat to Market Research
- But embrace it as a new (business) reality
- Learn how to process large amounts of data





Dank u

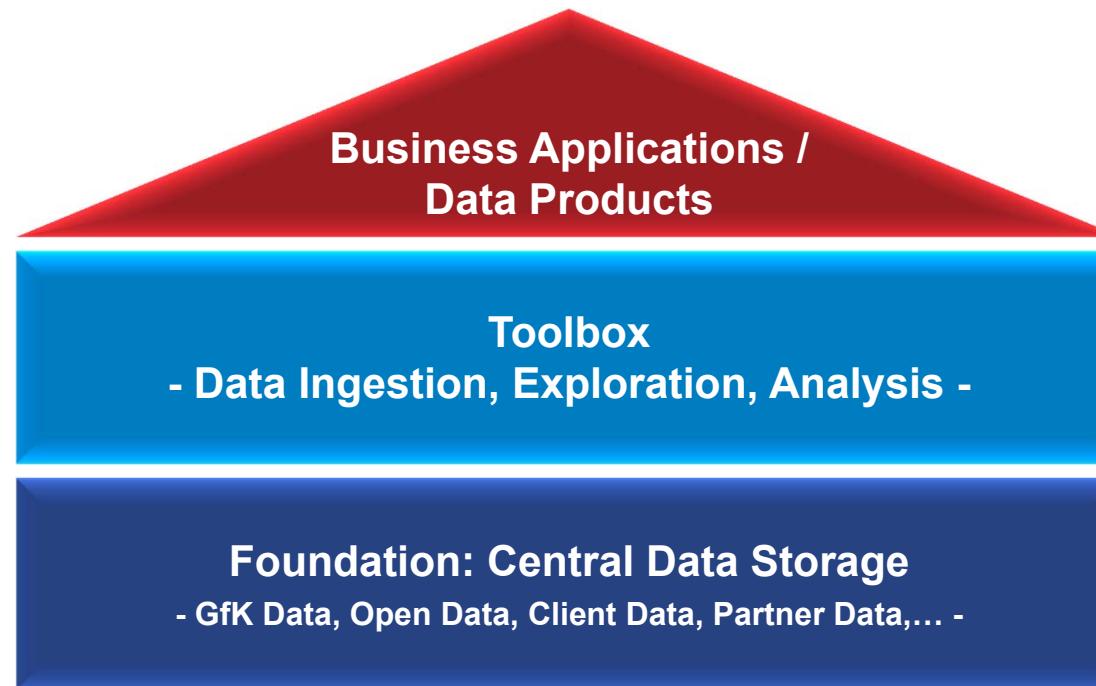


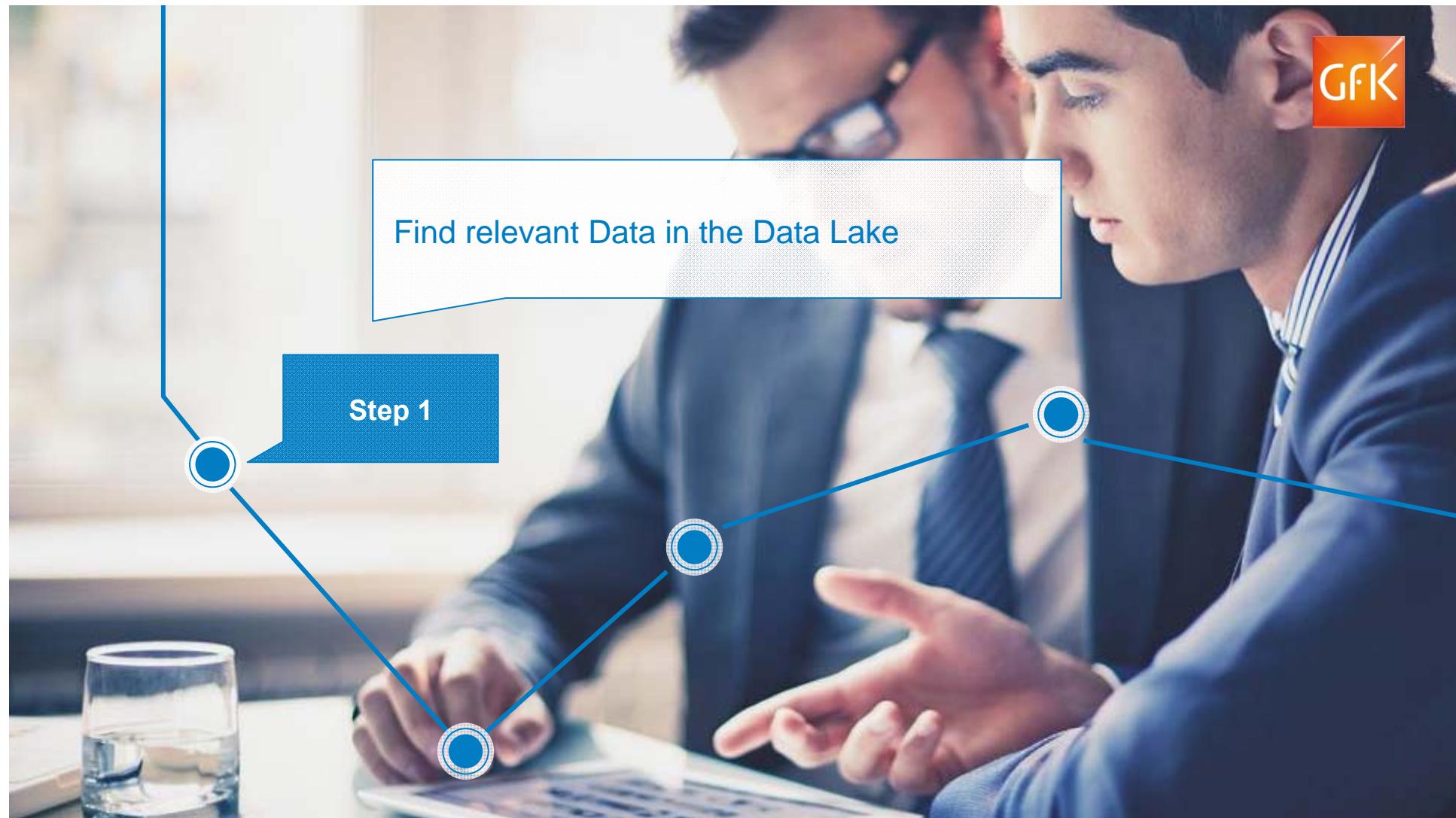
Backup Slides

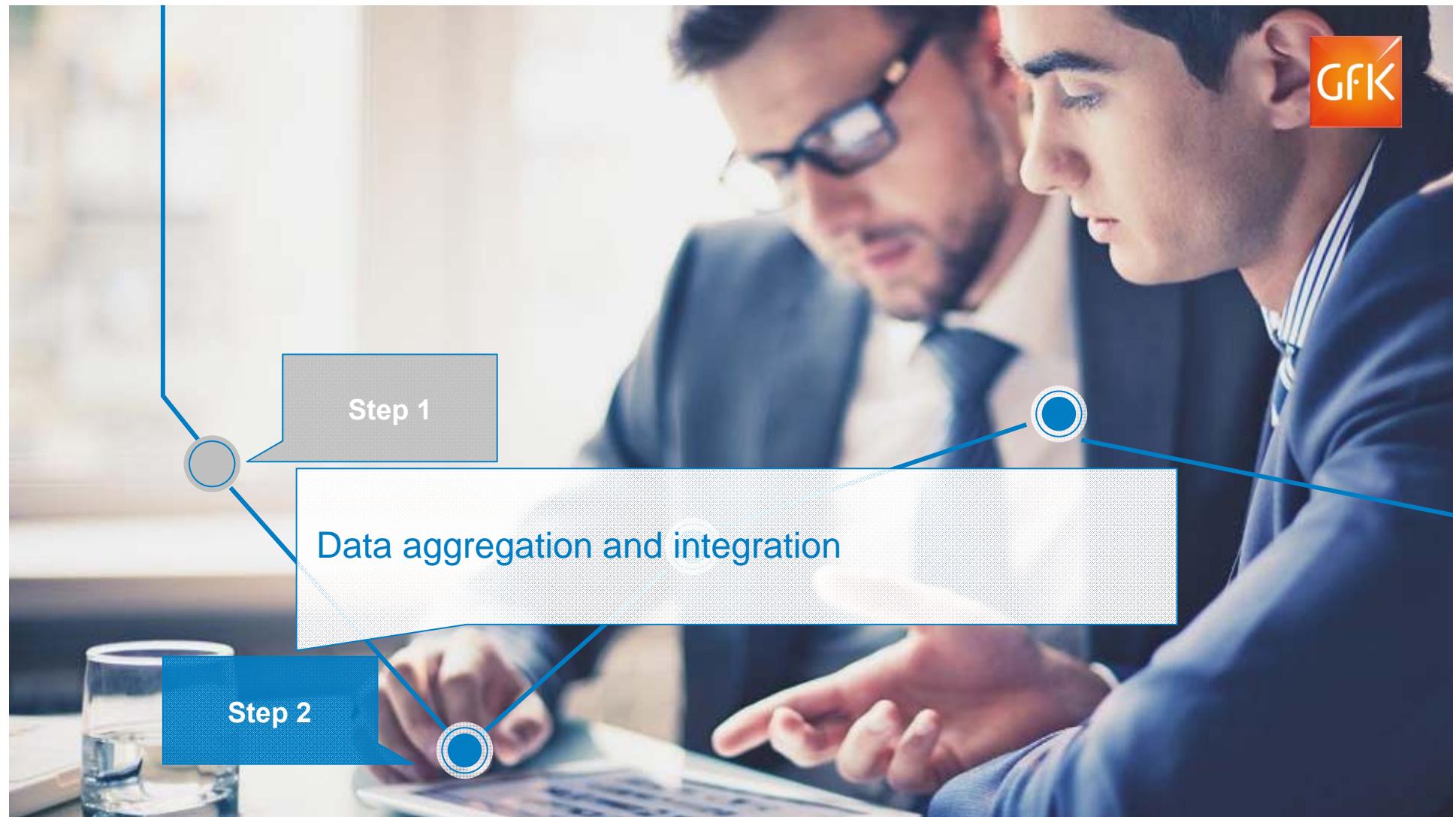
Our Data Lake combines storage for all our major data assets with tools for (Big) data ingestion, exploration, and analysis



Final objective: Valuable business applications that leverage our data landscape



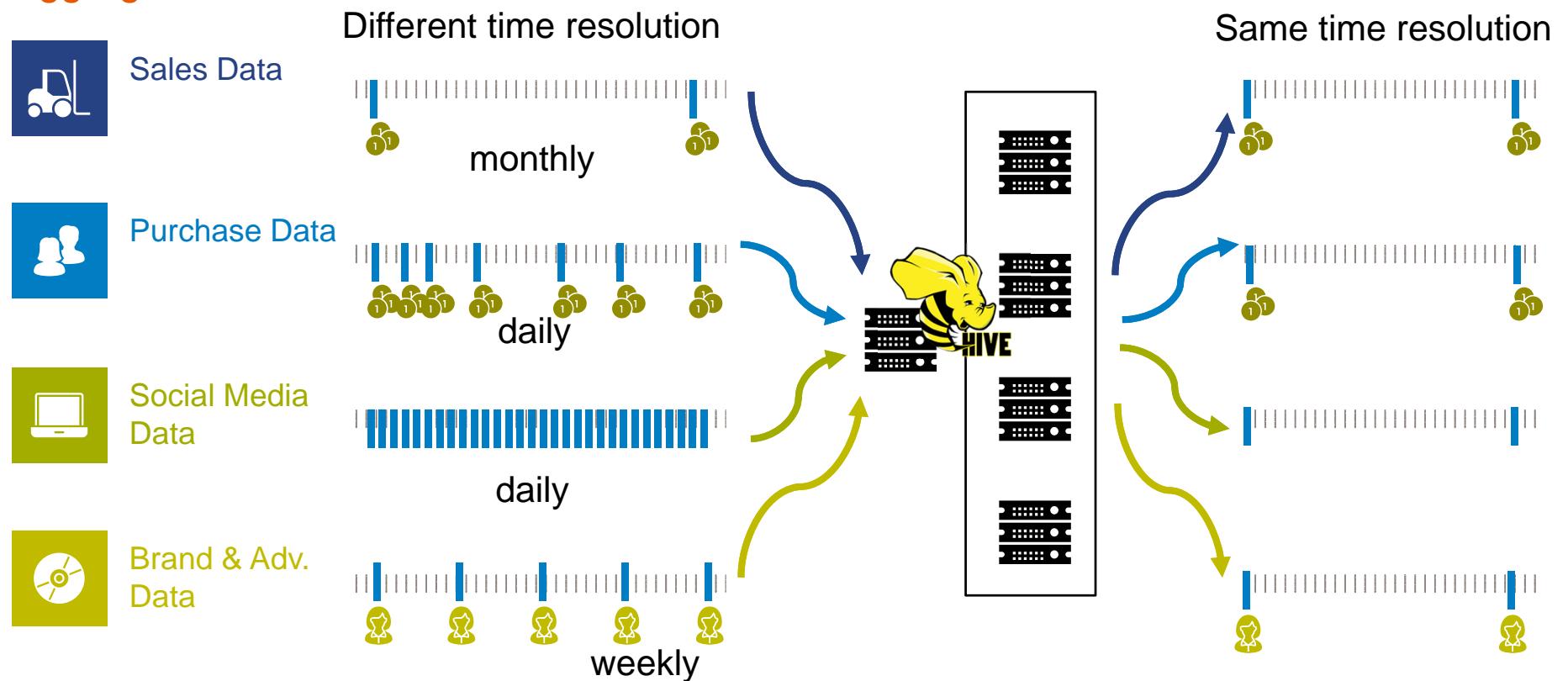




Hive is used to aggregate and integrate data from different sources with regard to the time dimension



Aggregation





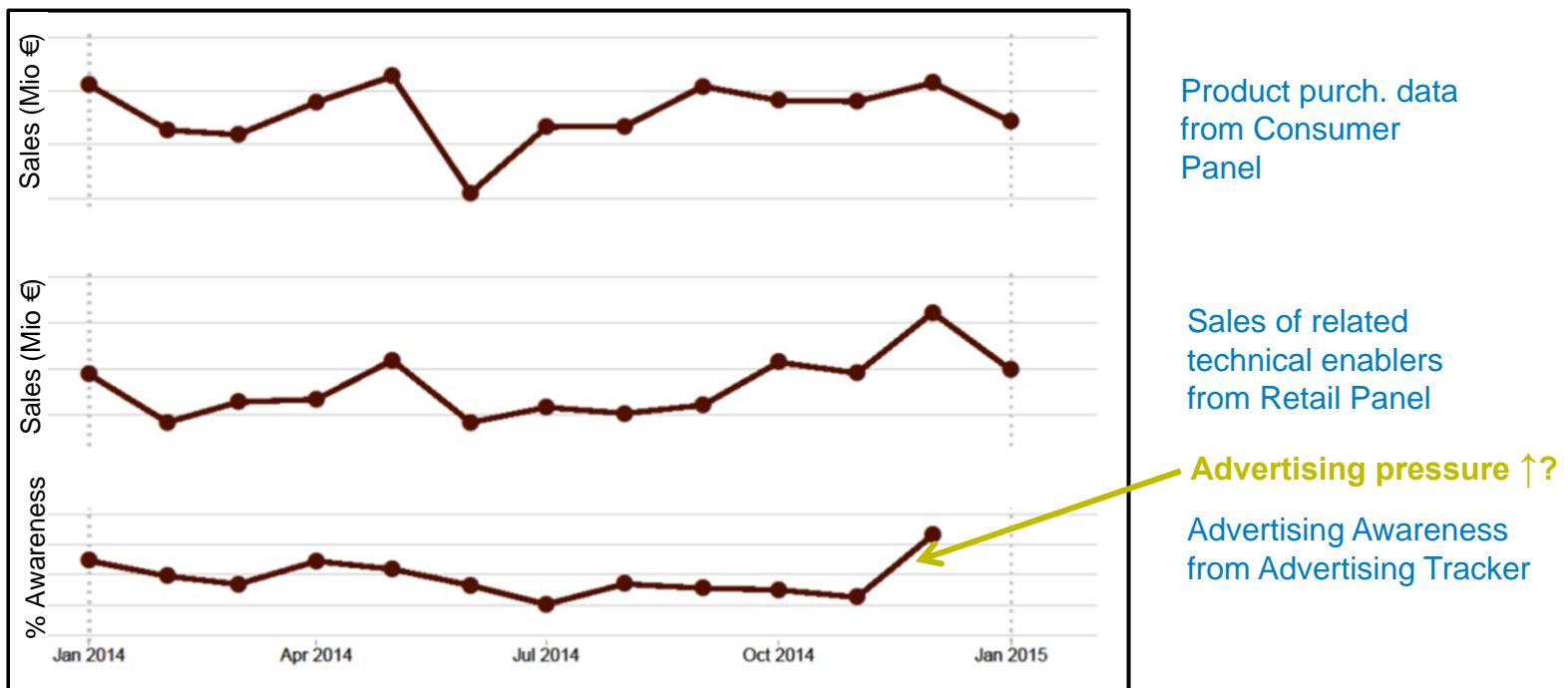
GFK

Data Scientists use preferred tools (e.g., R, Python) to explore data and understand relationships



For example, check for correlations of KPIs and insights from different data sources,

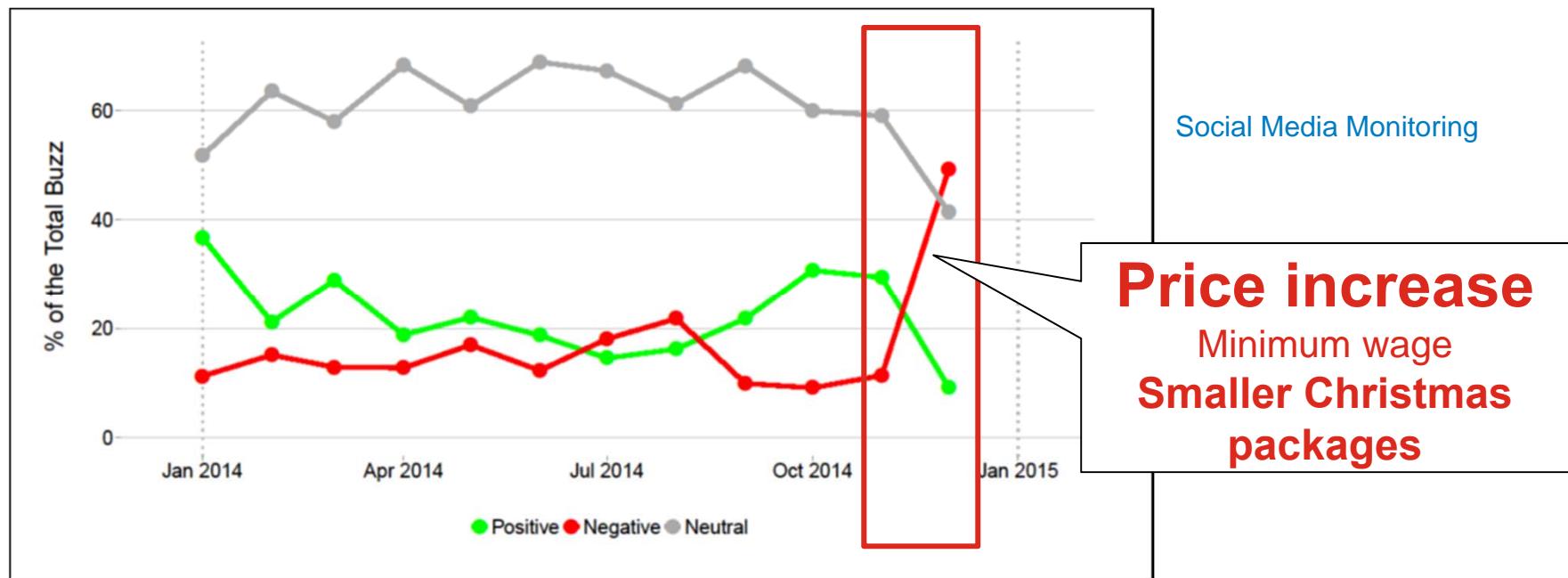
...



Data Scientists use preferred tools (e.g., R, Python) to explore data and understand relationships



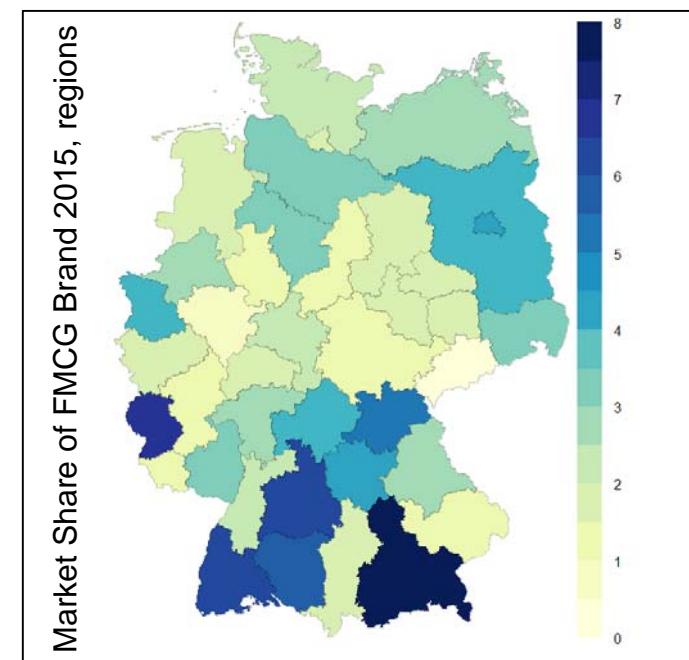
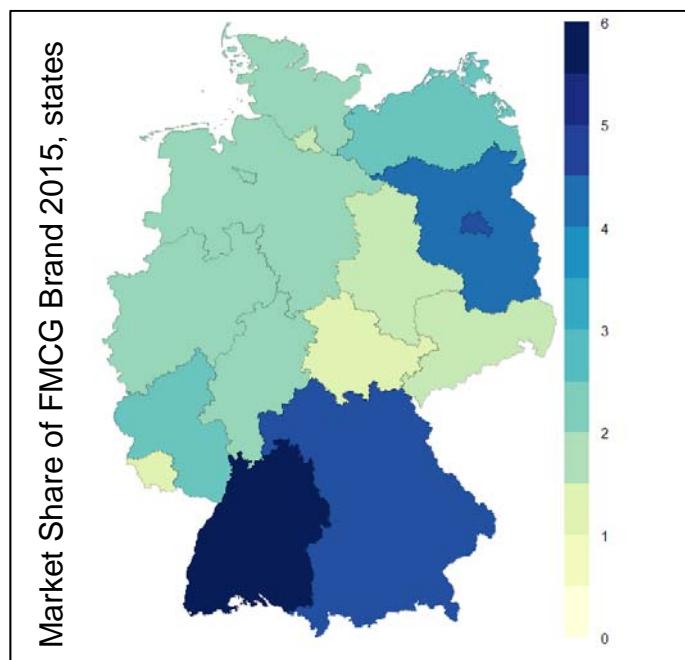
... dive into the “why” behind some data phenomena ...



Data Scientists use preferred tools (e.g., R, Python) to explore data and understand relationships



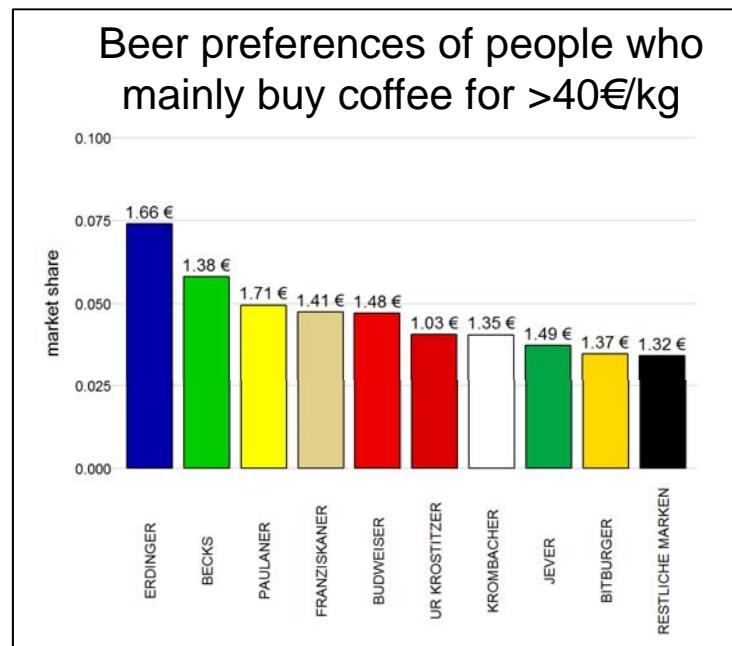
... explore regional sales patterns...



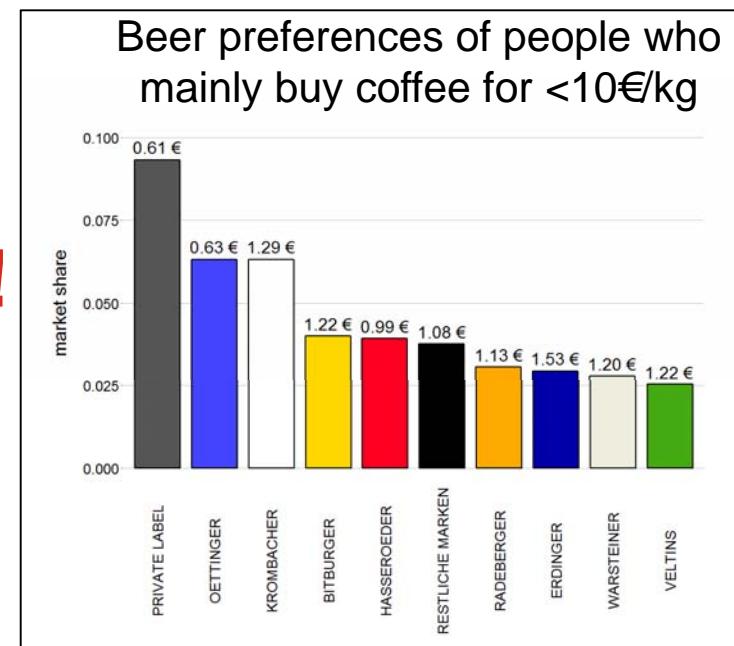
The Data Lake allows for quick and agile hypothesis-testing across data sets and product groups

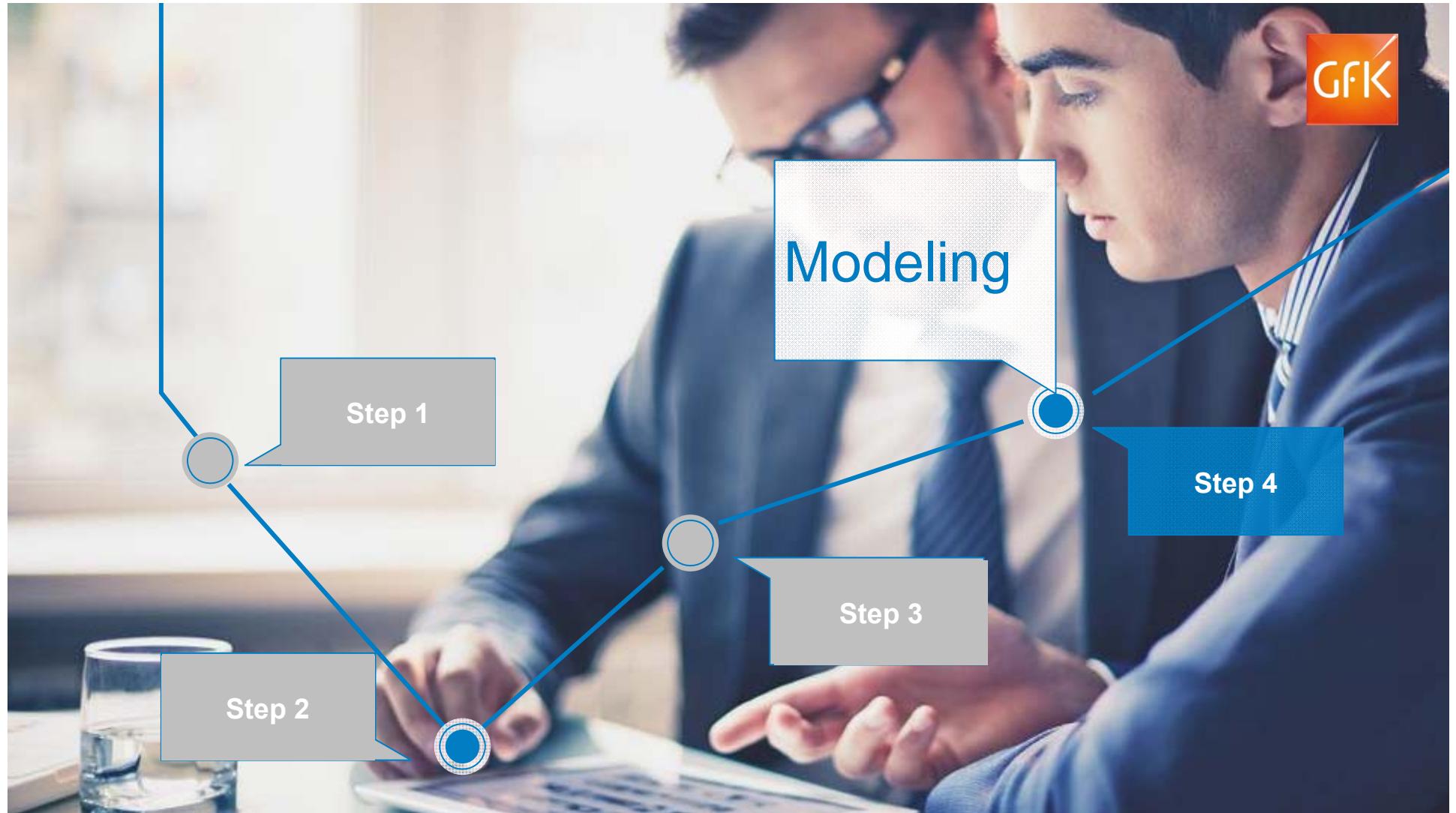


Do people who buy premium coffee also prefer more expensive beer brands?



YES!

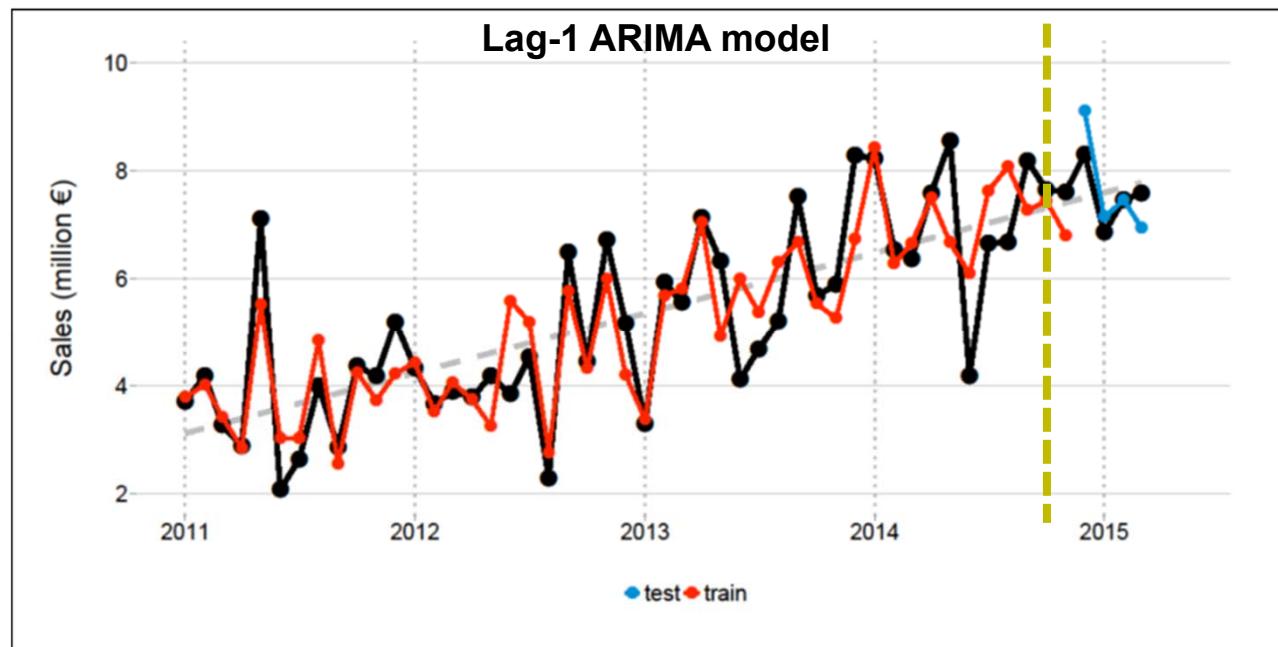




Finally, predictive modeling benefits from direct data access and freedom to explore different approaches



Prediction of Nespresso sales in 2015 – utilizing the full four years of coffee purchases...



Finally, predictive modeling benefits from direct data access and freedom to explore different approaches



...and a well-performing integrated model using only two years of information from different data sources

