# INTRODUCTION TO DATA SCIENCE

Lectures based on:
- ➢ E. Fox and C. Guestrin, „Machine Learning and Data Analysis", Univ. of Washington
- ➢ M. Cetinkays-Rundel, „Data Analysis and Statistical Inference", Univ. of Duke

2/10/2024

WFAiS UJ, Informatyka Stosowana

I stopień studiów

# What is Data Science?

**Is mainly about extracting knowledge from data (terms "data mining" or "Knowledge Discovery in Databases" are highly related). It can be about analyzing trends, building predictive models, … etc.**

**Is an agglomerate of <span style="color:red">data collection, data modeling and analysis</span>, a decision making, and everything you need to know to accomplish your goals. Eventually, it boils down to the following fields/skills:**

- **<u>Computer science:</u>**

**Algorithms, programming (patterns, languages etc.), understanding hardware & operating systems, high-performance computing'**

- **<u>Mathematical aspects:</u>**

**Linear algebra, differential equations for optimization problems, statistics**

- **<u>Few others:</u>**

**<span style="color:red">Machine learning</span>, domain knowledge, and data visualization & communication skills**

2/10/2024

# Data Science and Machine Learning?

**Machine learning** algorithms are algorithms that learn (often predictive) models from data. I.e., instead of formulating "rules" manually, a machine learning algorithm will learn the model for you.
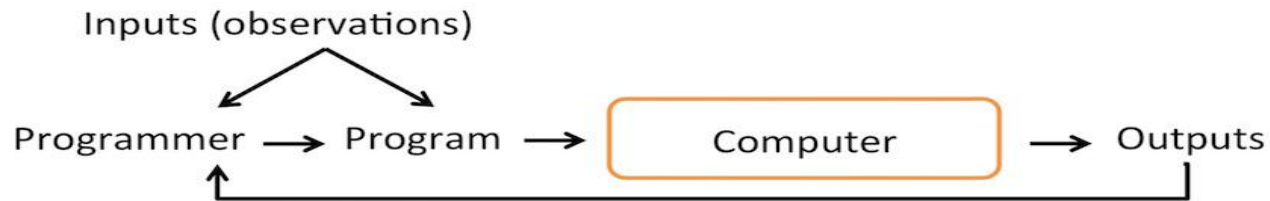
**Machine learning** - at its core - is about the use and development of these learning algorithms. **Data science** is more about the extraction of knowledge from data to answer particular question or solve particular problems.

**Machine learning is often a big part of a "data science" project**, e.g., it is often heavily used for exploratory analysis and discovery (clustering algorithms) and building predictive models (supervised learning algorithms). However, in **data science**, you often also worry about the collection, wrangling, and cleaning of your data (i.e., data engineering), and eventually, you want to draw conclusions from your data that help you solve a particular problem.

# Traditional programming paradigm and Machine Learning

**The Traditional Programming Paradigm**

Inputs (observations)

Programmer → Program → Computer → Outputs

*Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed* – Arthur Samuel (1959)
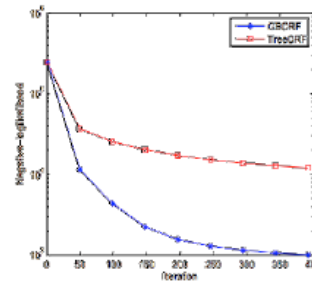
**Machine Learning**

Inputs

Outputs → Computer → Program

Sebastian Raschka, 2016

2/10/2024

# Outline of the course

- **Exploratory Data Analysis: introduction**

    **→ today**

- **Data Analysis with Machine Learning algorithms:**
    - **Regression  (October)**
    - **Classification  (November)**
    - **Retrieval & Clustering  (December)**
    - **Other ML methods, Statistical inference (January)**

# Analyse data with Machine Learning

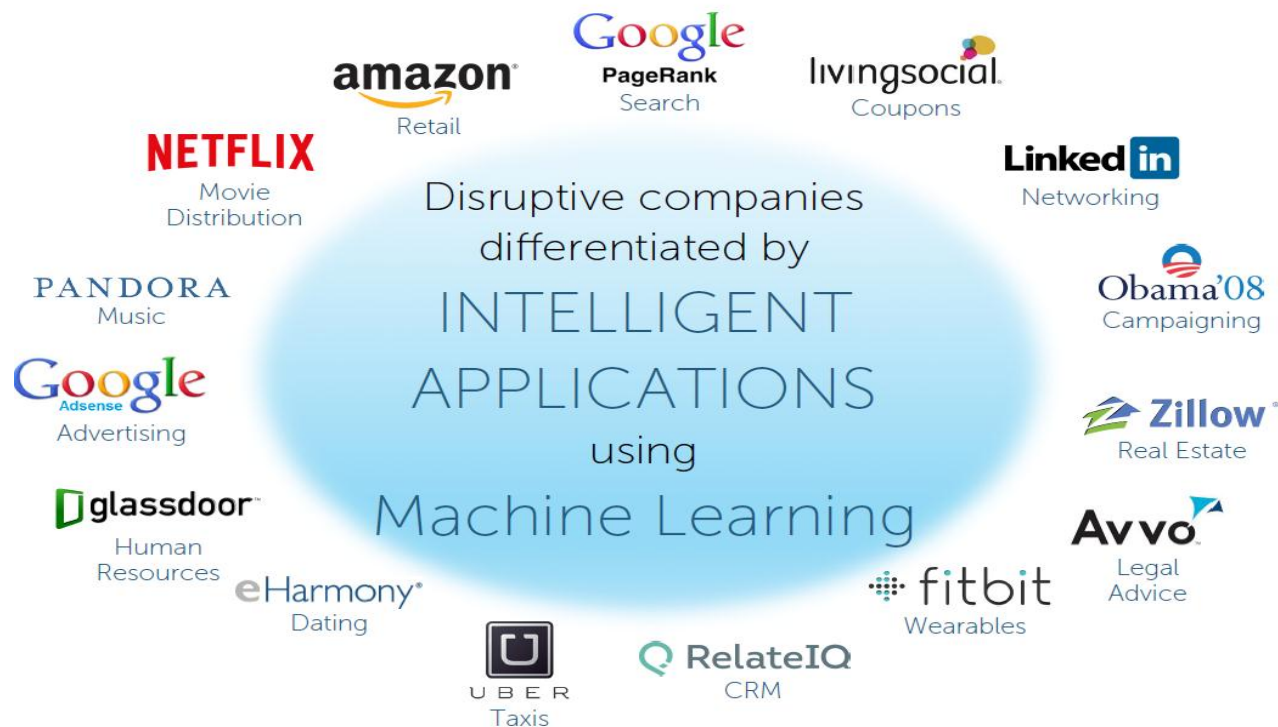- **Machine learning is changing the world.**
- **Old view of ML**



2/10/2024

# Machine learning is changing the world

- **Current view: disruptive inteligent applications are used by leading comercial companies**
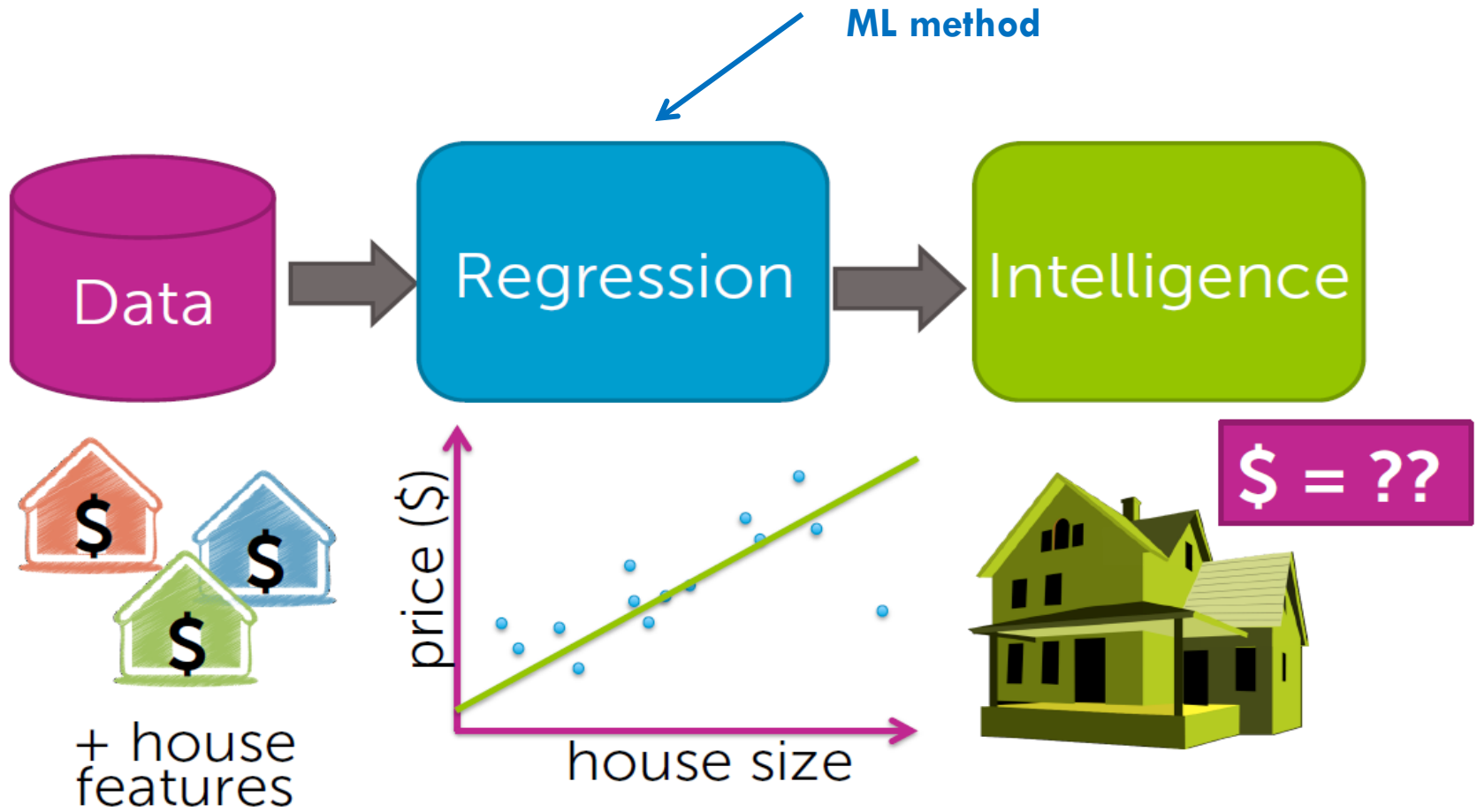


2/10/2024

# Machine learning

□ **Data → inteligence pipeline**

**New kind of analysis which brings inteligence how to solve a problem**



Data → ML Method → Intelligence

**Eg. which product to buy which film to chose connect people and taxi driver**

2/10/2024

# Case study 1: Prediction

**ML method**

Data → Regression → Intelligence

+ house features

price ($)

house size

$ = ??

# Prediction

hard work

- How much will your salary be? ($y = \$\$$)
- Depends on $x$ = performance in courses, quality of programming assignments, # of discussion responses, ...

2/10/2024

# Prediction

## Tweet popularity

- How many people will retweet your tweet?
- Depends on # followers, # of followers of followers, features of text tweeted, popularity of hashtag, # of past retweets,...

# Prediction:

| Models | • Linear regression<br>• Regularization:<br>Ridge (L2), Lasso (L1) |
|---|---|
| Algorithms | • Gradient descent<br>• Coordinate descent |
| Concepts | • Loss functions, bias-variance tradeoff, cross-validation, sparsity, overfitting, model selection |

2/10/2024

# Case study 2: Classification

# Classification

## Spam filtering



Input: **x**     Output: y

2/10/2024

# Multiclass classifier

## Output y has more than 2 categories



Input: **x**
Webpage

Education

Finance

Technology

Output: y

# Classification

## Image classification



Input: **x**
Image pixels

Top Predictions
- Labrador retriever
- golden retriever
- redbone
- bloodhound
- Rhodesian ridgeback

soft prediction

Output: y
Predicted object

2/10/2024

# Classification

# Classification:

**Models**
- Linear classifiers
  (logistic regression, SVMs, perceptron)
- Kernels
- Decision trees

**Algorithms**
- Stochastic gradient descent
- Boosting

**Concepts**
- Decision boundaries, MLE, ensemble methods, random forests, CART, online learning
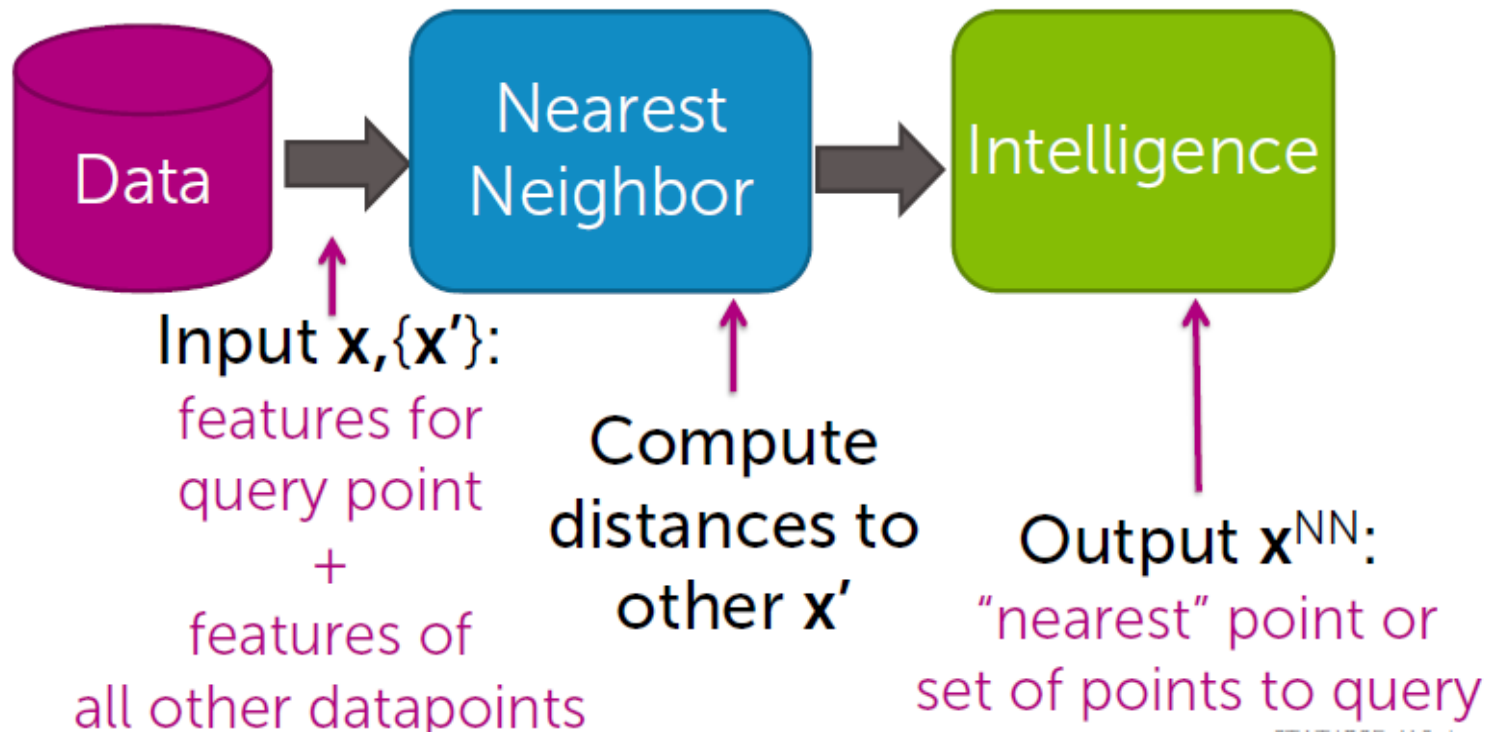
2/10/2024

# Case Study: document retrieval

# Retrieval

Search for related items



Data → Nearest Neighbor → Intelligence

Input $\mathbf{x},\{\mathbf{x'}\}$:
features for query point
+
features of all other datapoints

Compute distances to other $\mathbf{x'}$

Output $\mathbf{x}^{NN}$:
"nearest" point or set of points to query

2/10/2024

# Retrieval

## Retrieve "nearest neighbor" article

Space of all articles,
organized by similarity of text

query article

nearest neighbor

# Retrieval

## Or set of nearest neighbors

Space of all articles,
organized by similarity of text



query article

set of nearest neighbors

# Retrieval

# Case study 3++:
## Document structuring for retrieval

ML method

Data → Clustering → Intelligence

SPORTS

WORLD NEWS

ENTERTAINMENT

SCIENCE

2/10/2024

# Clustering

## Clustering images

For search, group as:

- Ocean
- Pink flower
- Dog
- Sunset
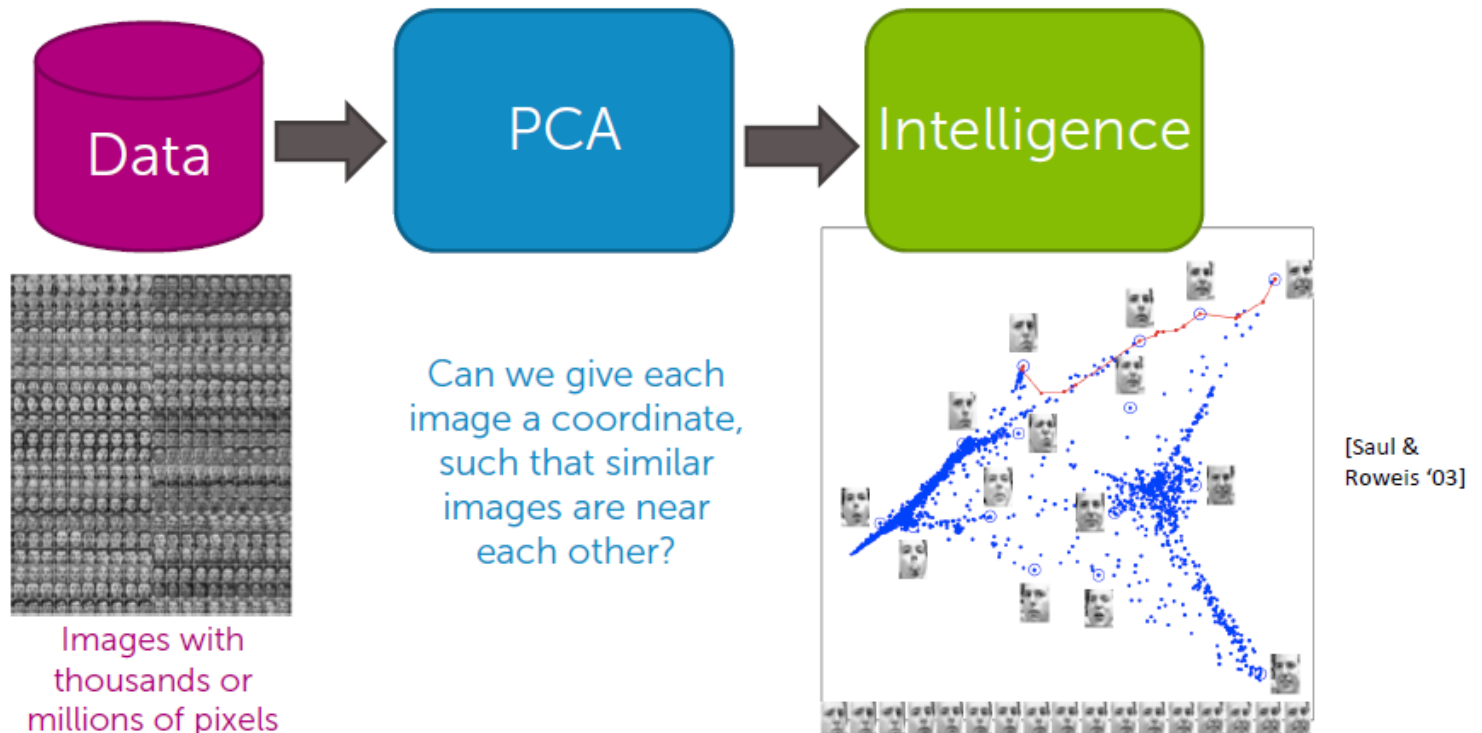- Clouds
- ...

# Clustering

## Or users on websites...

Discover groups of users for better targeting of content



2/10/2024

# Embeding

## Example: Embedding images to visualize data



Data
→ PCA →
Intelligence

Can we give each image a coordinate, such that similar images are near each other?

Images with thousands or millions of pixels

[Saul & Roweis '03]

# Clustering: Finding documents

**Models**
- Nearest neighbors
- Clustering, mixtures of Gaussians
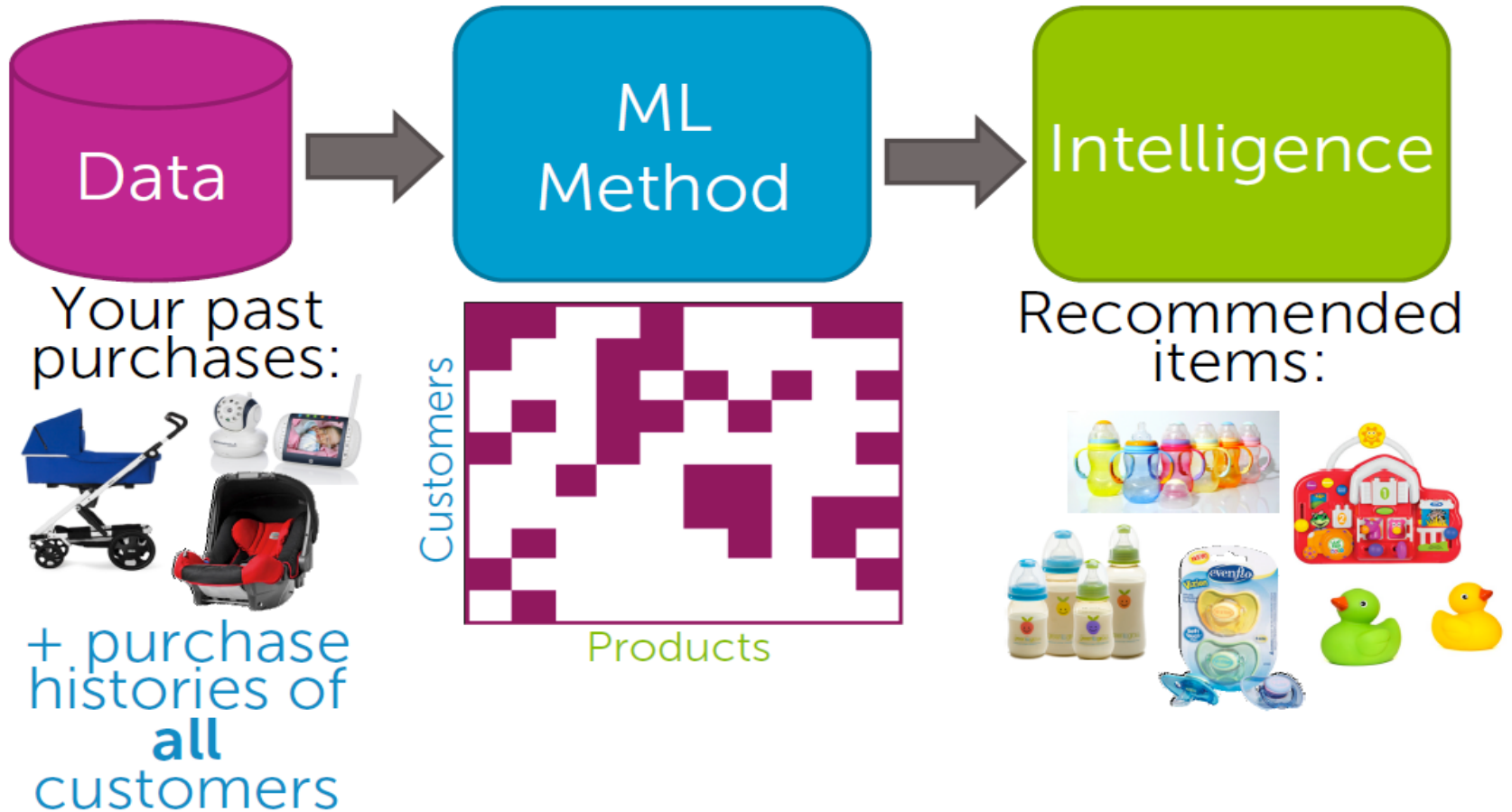- Latent Dirichlet allocation (LDA)

**Algorithms**
- KD-trees, locality-sensitive hashing (LSH)
- K-means
- Expectation-maximization (EM)

**Concepts**
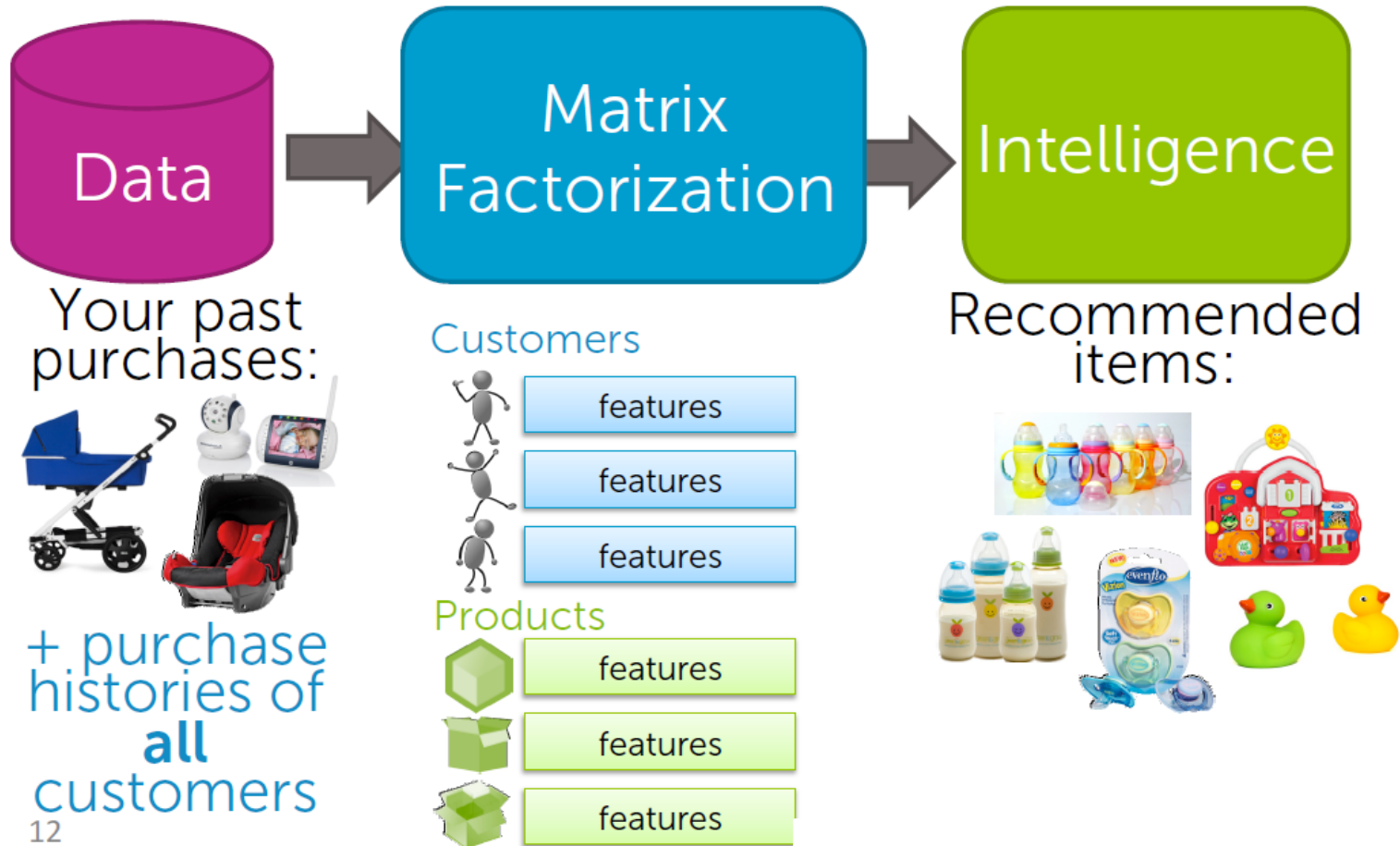- Distance metrics, approximation algorithms, hashing, sampling algorithms, scaling up with map-reduce

2/10/2024

# Case study: Product recommendation

Data → ML Method → Intelligence

Your past purchases:

+ purchase histories of **all** customers

Customers / Products

Recommended items:

# Case study: Product recommendation

# Recomender systems applications

Movies



Songs



Friends, apps, ...

2/10/2024

# Case study 5:
# Visual product recommender

# What is (supervised) deep learning?

Flexible method for performing classification or regression



Input {$x$}:
raw data or extracted features for points in dataset

Nonlinear feature representation

Output {$z$}:
label or value

2/10/2024

# Examples of deep learning success stories

- Image classification
- Image segmentation
- Image captioning
- Object detection
- Speech recognition
- Speech synthesis
- Machine translation
- Handwriting recognition
- ...

# Other ML methods

- Reinforcement learning
- Learning theory
- Active learning
- Multi-task and transfer learning
- Spectral methods
- ...

# Deploing inteligence module

**Case studied are about building, evaluating, deploying inteligence in data analysis.**



ML Method

Data

Intelligence

Task

Model & parameters

Optimization algorithm

Evaluation

Use pre-specified or develop your own

2/10/2024