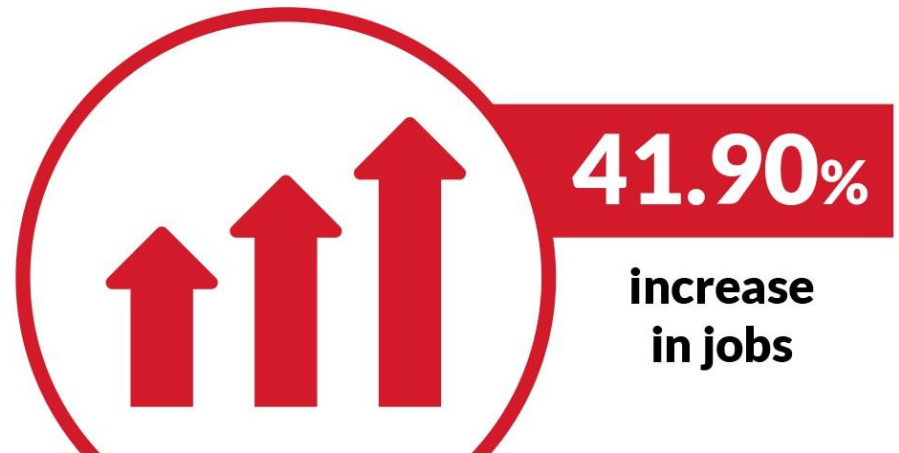


Introduction to Data Science

LECTURE 1: DS & ML

Project Engineer
LOVNISH VERMA

Demand for Data Science Positions Job Growth (2021 - 2031)

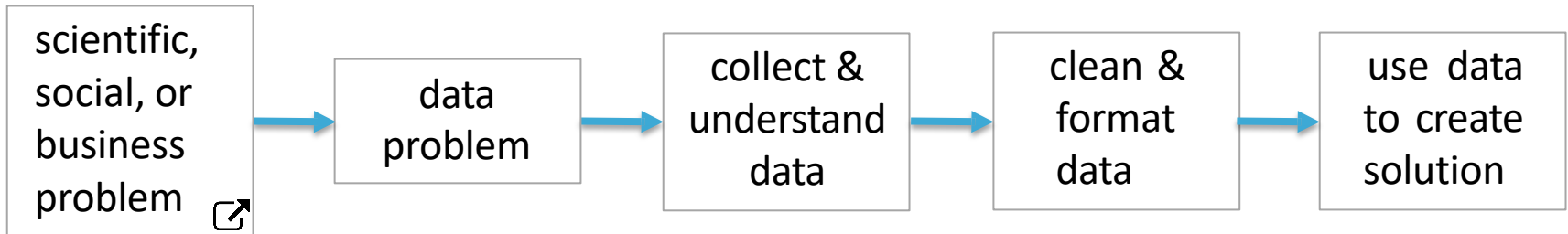


WHAT IS DATA SCIENCE?

- “Data science, also known as data-driven science, is an interdisciplinary field of scientific methods, processes, algorithms and systems to extract knowledge or insights from data in various forms, either structured or unstructured, similar to data mining.”
- “Data science intends to analyze and understand actual phenomena with ‘data’. In other words, the aim of data science is to reveal the features or the hidden structure of complicated natural, human, and social phenomena with data from a different point of view from the established or traditional theory and method.”

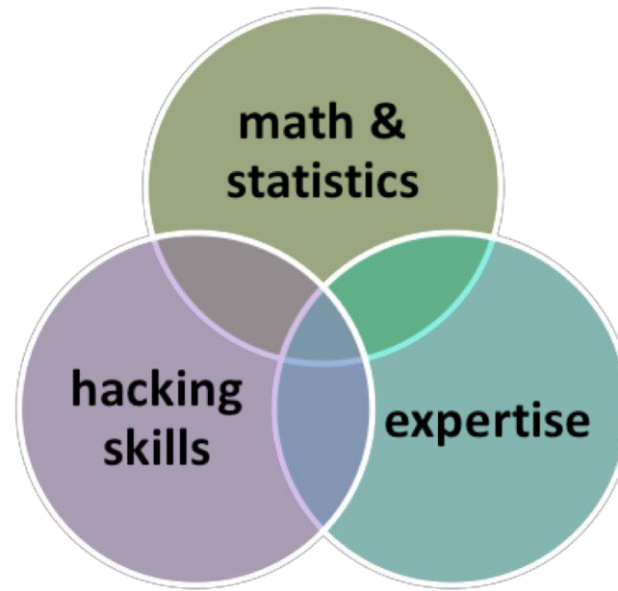
WHAT IS DATA SCIENCE?

...solving problems with data...



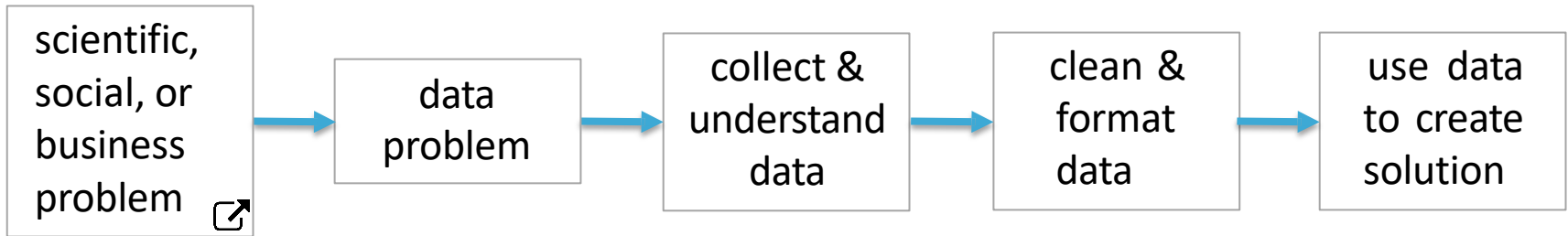
...sounds cool!

What makes a good data scientist?

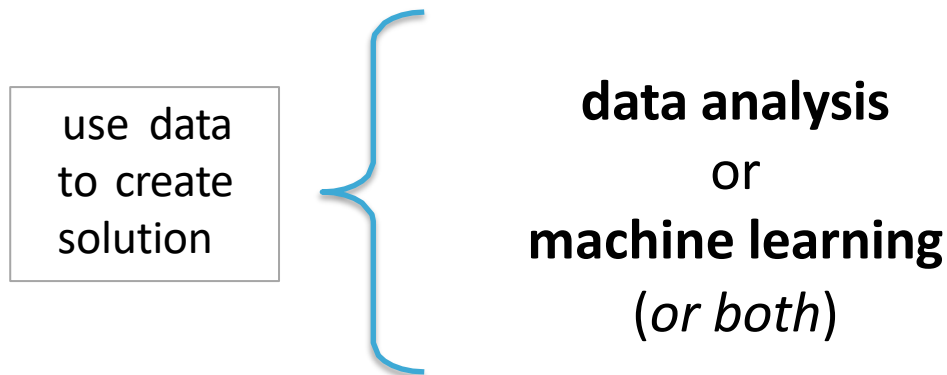


WHAT IS DATA SCIENCE?

...solving problems with data...



...which step is most challenging?



WHAT IS DATA ANALYSIS?

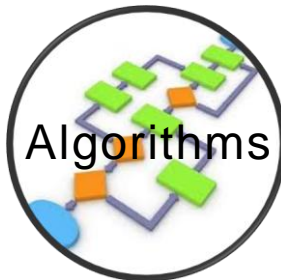
...using data to discover useful information...



- **data:** anything you can *measure* or *record*



- **statistics:** summarize (and visualize) *main characteristics* of the data



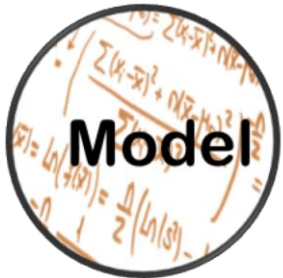
- **algorithms:** apply algorithms to find *patterns* in the data

WHAT IS MACHINE LEARNING?

...creating and using models that learn from data...



- **data:** anything you can *measure* or *record*



- **model:** specification of a (mathematical) *relationship* between different variables



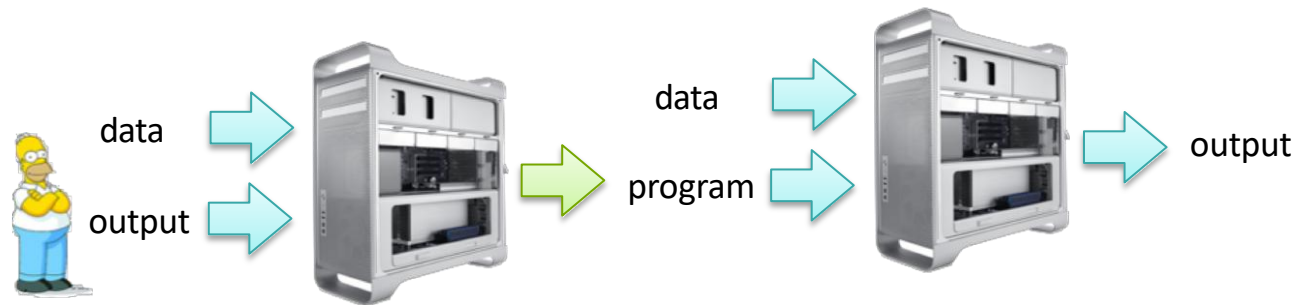
- **evaluation:** how well does the model *work*?

WHAT IS MACHINE LEARNING?

- Traditional CS



- Machine Learning

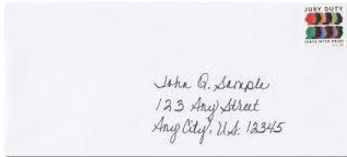


WHAT IS MACHINE LEARNING?

- ...creating and using models that learn from data...

- Examples

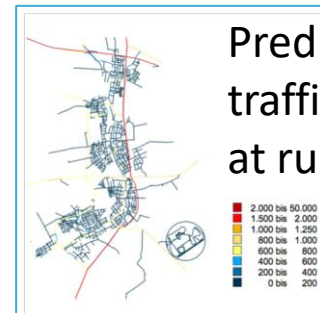
Identifying zip code
from handwritten
digits



Detecting
communities
in social
networks



Predicting the
traffic volume
at rush hour



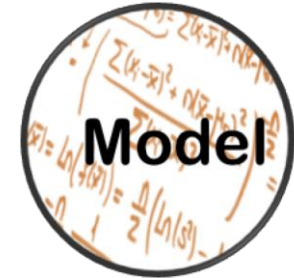
Detecting fraudulent
credit card
transactions



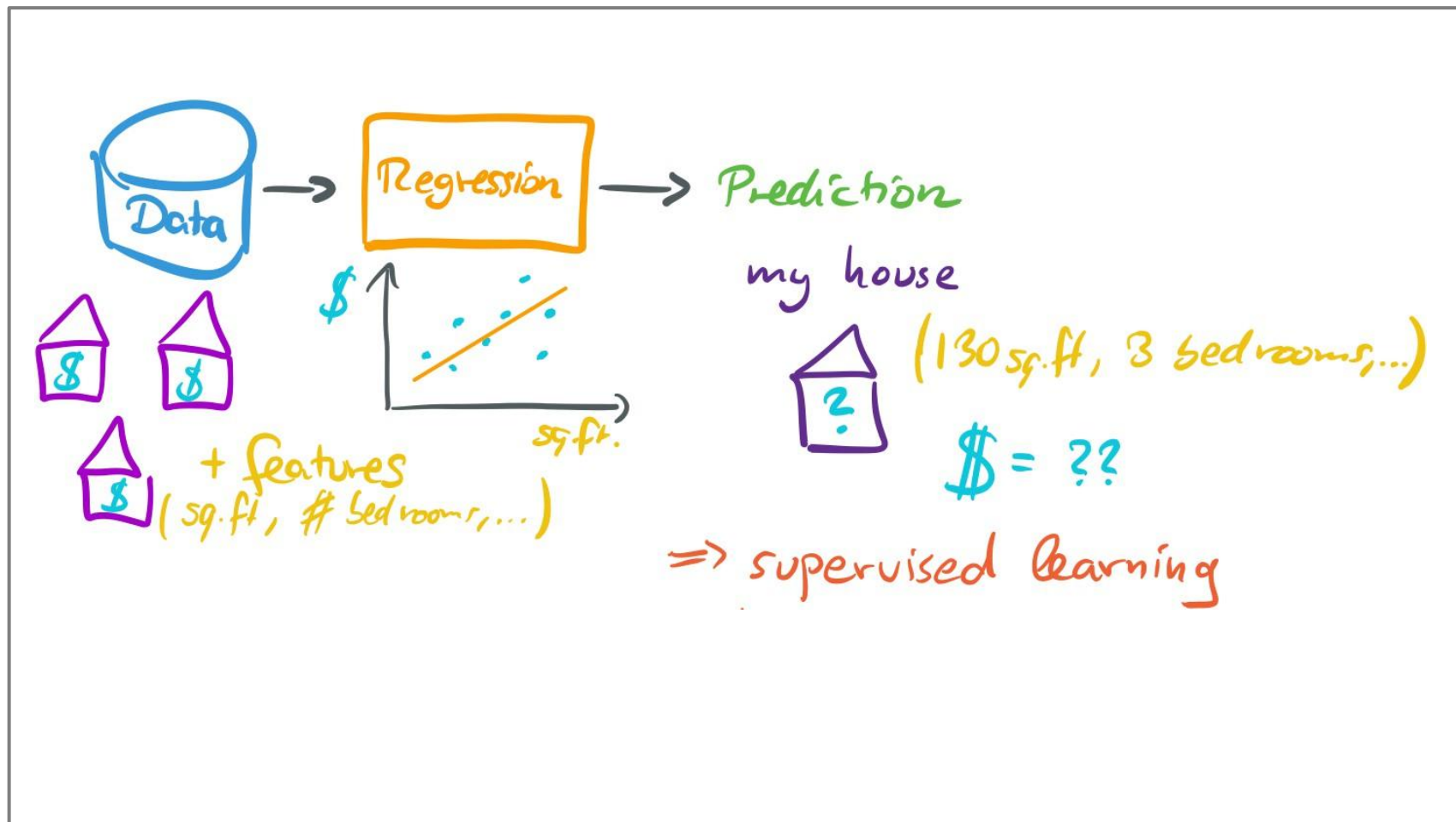
Determining the
location of distribution
centers based on
customers'
residence



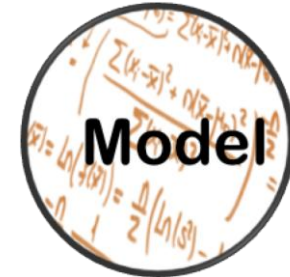
LEARNING FROM DATA



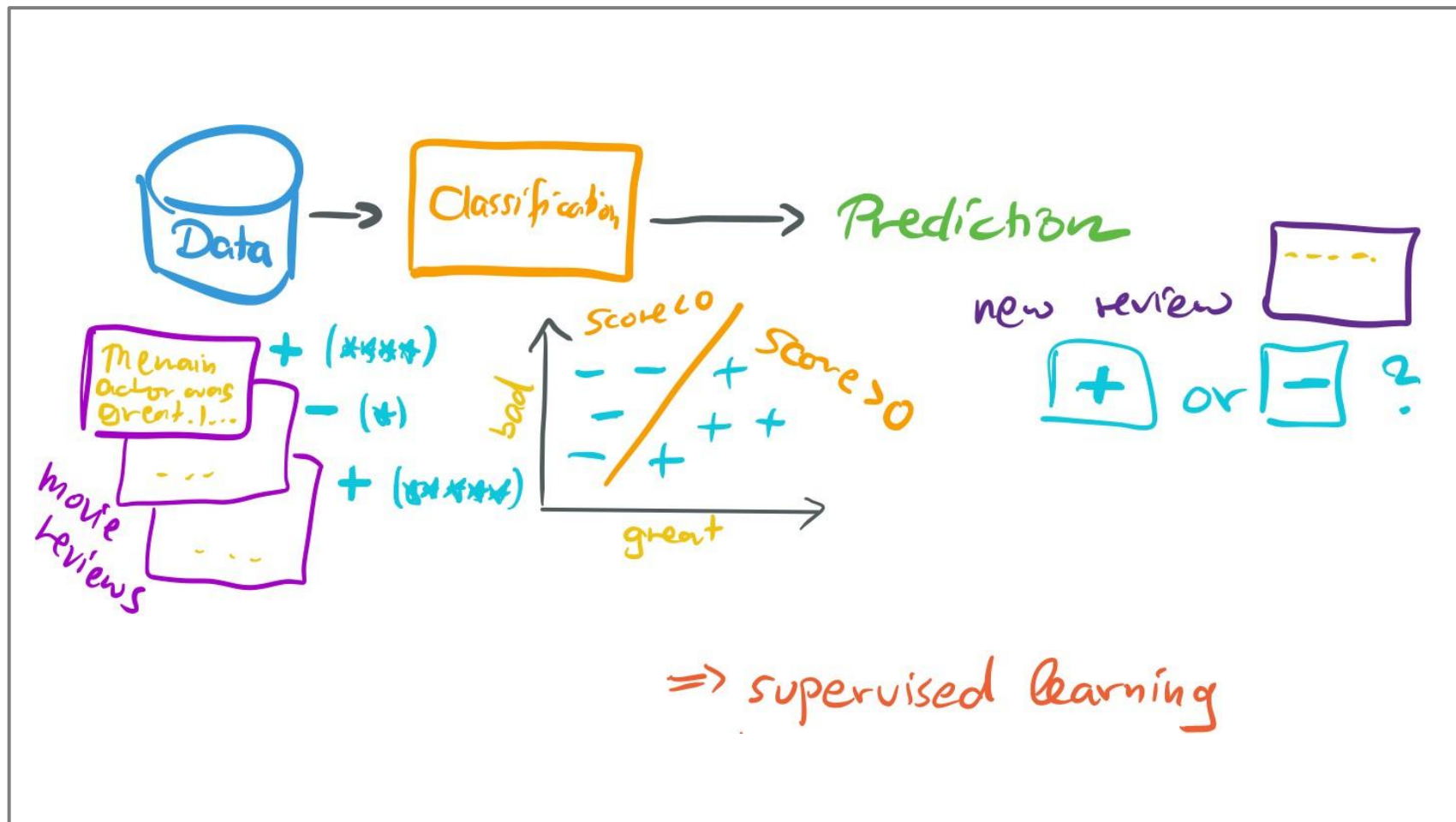
- Regression



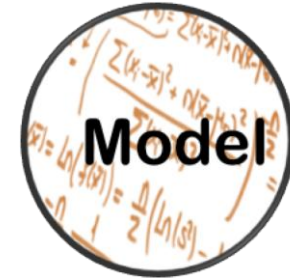
LEARNING FROM DATA



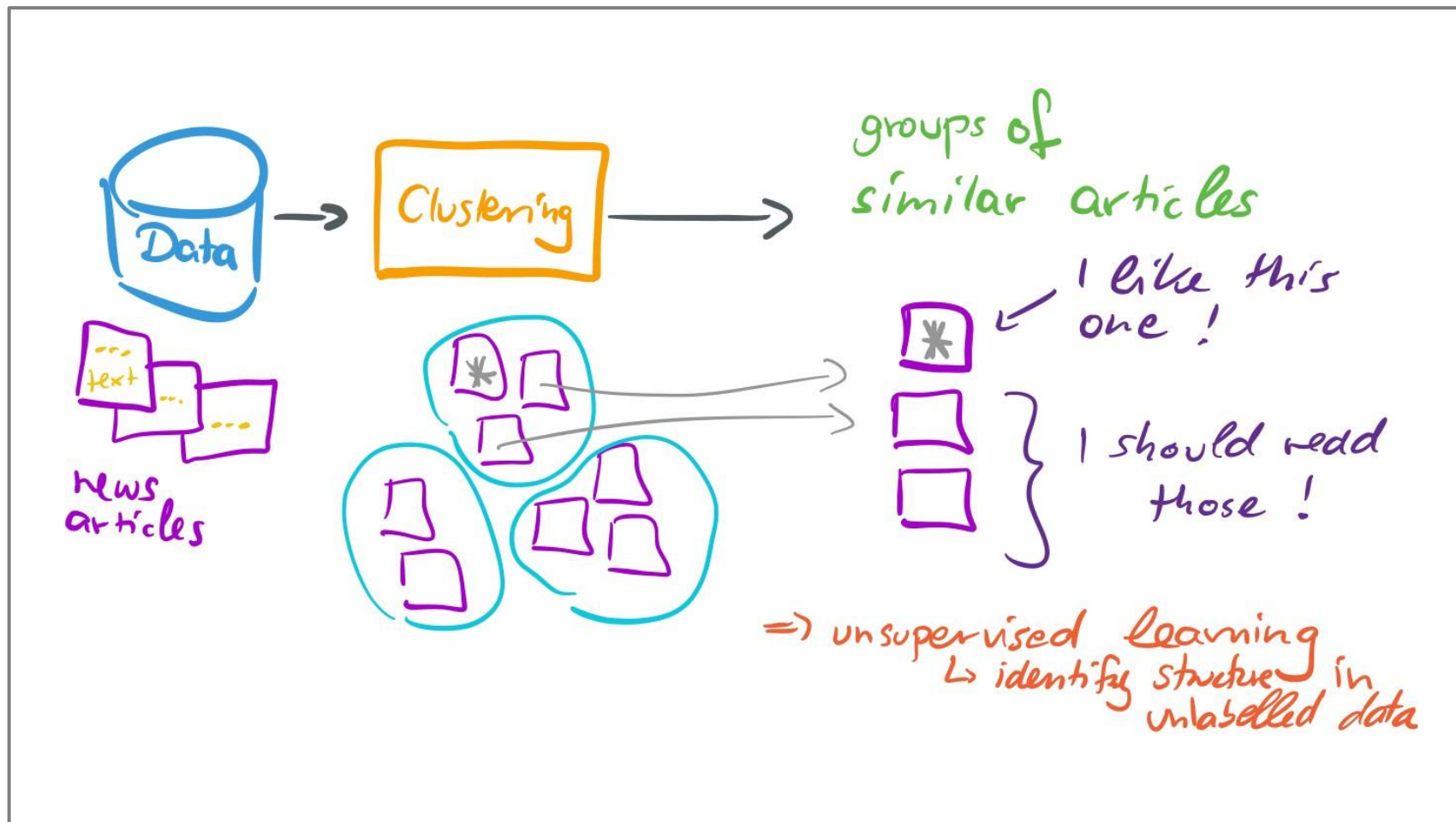
- Classification



LEARNING FROM DATA



- Clustering



WHAT IS MACHINE LEARNING?

regression, classification, clustering

...creating and using models that learn from data...

→ supervised learning/predictive modelling

- come up with predictions
- extract knowledge/insights

→ unsupervised learning/data mining

ACTIVITY 1

regression, classification, clustering

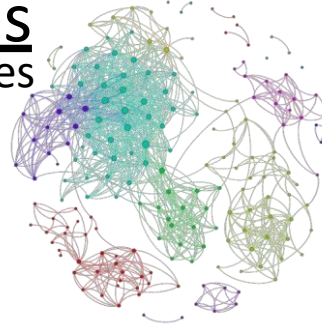
- ...creating and using models that learn from data...

- Categorize these Examples

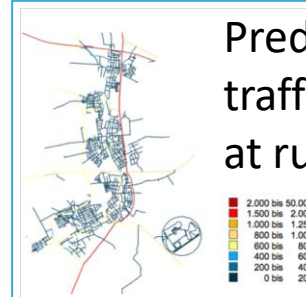
Identifying zip code
from handwritten
digits



communities
in social
networks



Predicting the
traffic volume
at rush hour



Detecting fraudulent
credit card
transactions

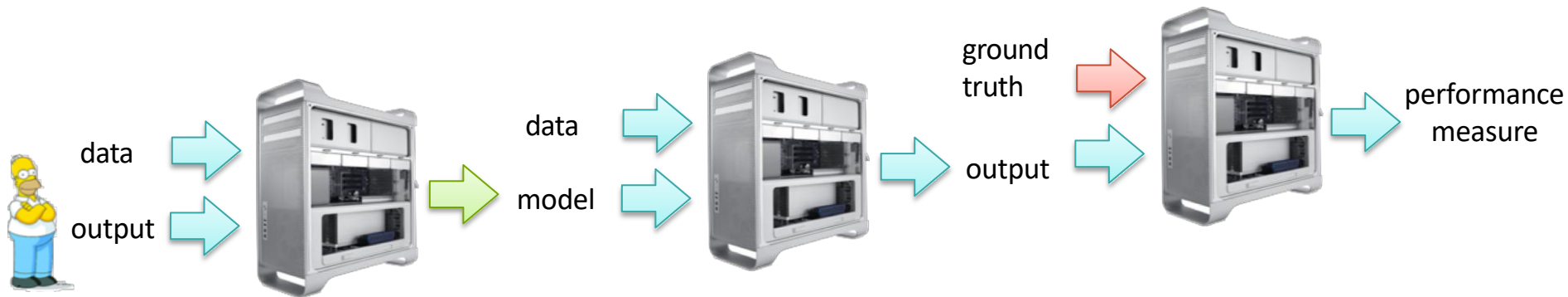


Determining the
location of distribution
centers based on
customers'
residence



MACHINE LEARNING WORKFLOW

- *training phase, test phase, evaluation phase*



→ let's have a closer look at the *data* we are using

DATA



- Notation:

- D all observed data
- X all features
- y observations
- \square_{TE} test
- \square_{TR} training
- \hat{y} predictions

Helper Notation:

n number of data points

d number of features

m number of training points

$\square_{1, \dots, i, \dots, n}$: indices for data points

$\square_{1, \dots, j, \dots, d}$: indices for features

- What data structure to use?
 - *set, list, or array?*

SUMMARY & READING

- *Data Science* is about

data, models, and evaluation

- *Data Science* can solve a wide **variety of problems** – once we have the *right* data *and* model!

