

# Codettes & Geekettes Code Pub Munich Data Science Workshop



```
if _operation == "MIRROR_X":  
    mirror_mod.use_x = True  
    mirror_mod.use_y = False  
    mirror_mod.use_z = False  
  
elif _operation == "MIRROR_Y":  
    mirror_mod.use_x = False  
    mirror_mod.use_y = True  
    mirror_mod.use_z = False  
  
elif _operation == "MIRROR_Z":  
    mirror_mod.use_x = False  
    mirror_mod.use_y = False  
    mirror_mod.use_z = True  
  
#selection at the end -add back the deselected mirror modifier object  
mirror_ob.select= 1  
modifier_ob.select=1  
bpy.context.scene.objects.active = modifier_ob  
print("Selected" + str(modifier_ob)) # modifier ob is the active ob  
#mirror_ob.select = 0  
base = bpy.context.selected_objects[0]  
base.select = True
```

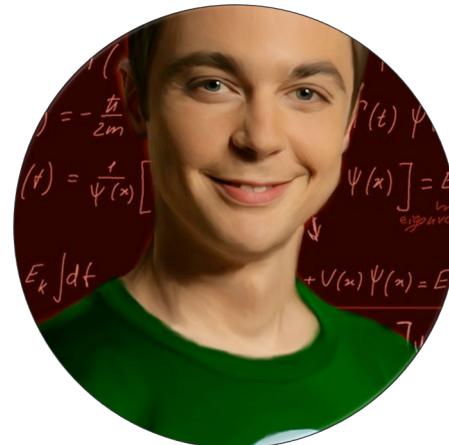


Dr. Lydia Nemec



Data Scientist @ Zeiss  
Theoretical Physicist by training

but I come in a  
wrapping similar to  
Penny



I am a bit like  
Dr. Sheldon Cooper,



I'm am nerdy



Veronika Zellner  
Mathematician  
Architect for Data & AI @Microsoft

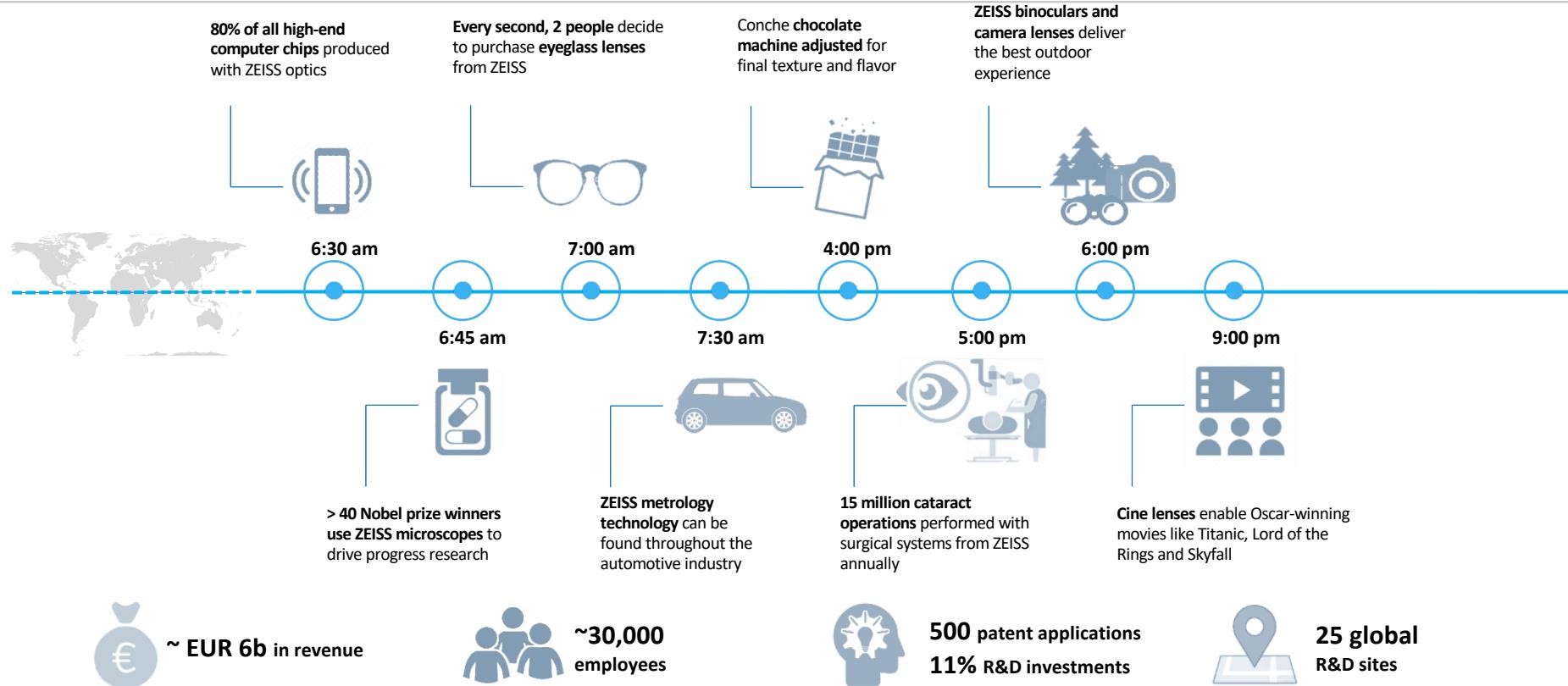


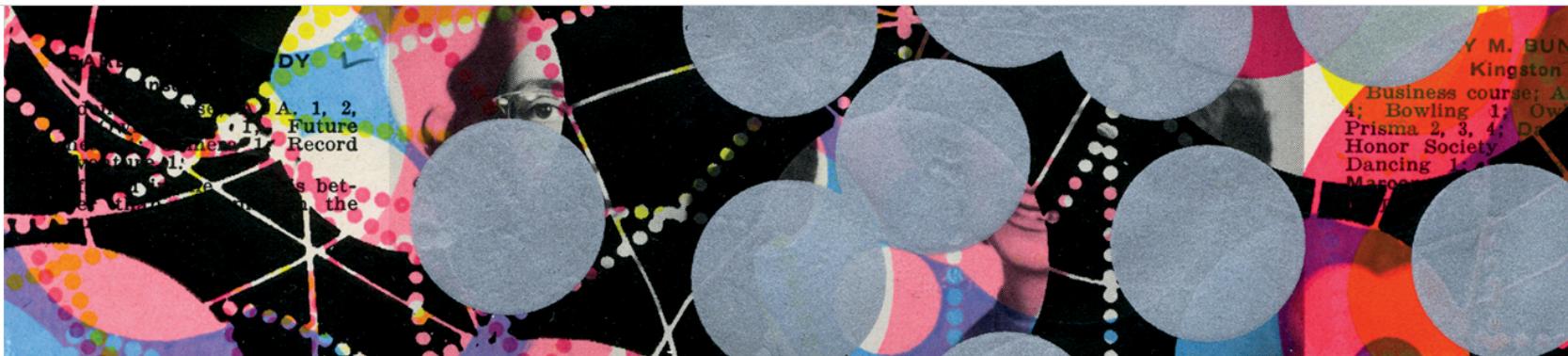
But also girly



# One day without ZEISS?

## A “user journey”





ARTWORK: TAMAR COHEN, ANDREW J BUBOLTZ, 2011, SILK SCREEN ON A PAGE FROM A HIGH SCHOOL YEARBOOK, 8.5" X 12"

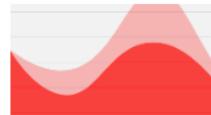
## DATA

# Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

## WHAT TO READ NEXT



Visualizations That Really Work

Data Scientists apply numerical methods like Machine Learning to extract insights from data.

## Math, Numeric & Statistic

- Machine Learning (AI)
- Statistical modelling
- Linear Algebra & Optimization

## The Scientific Mind

- Logical & independent mind
- Planning, conducting & evaluate experiments
- Excellent analytical skills
- Meticulous attention to quality and accuracy



## Computer Science & Programming

- Software development
- Programming Language (e.g. python)
- Databases (SQL/ No-SQL)
- Cloud Computing

## Communication, Soft Skills & Visualisation

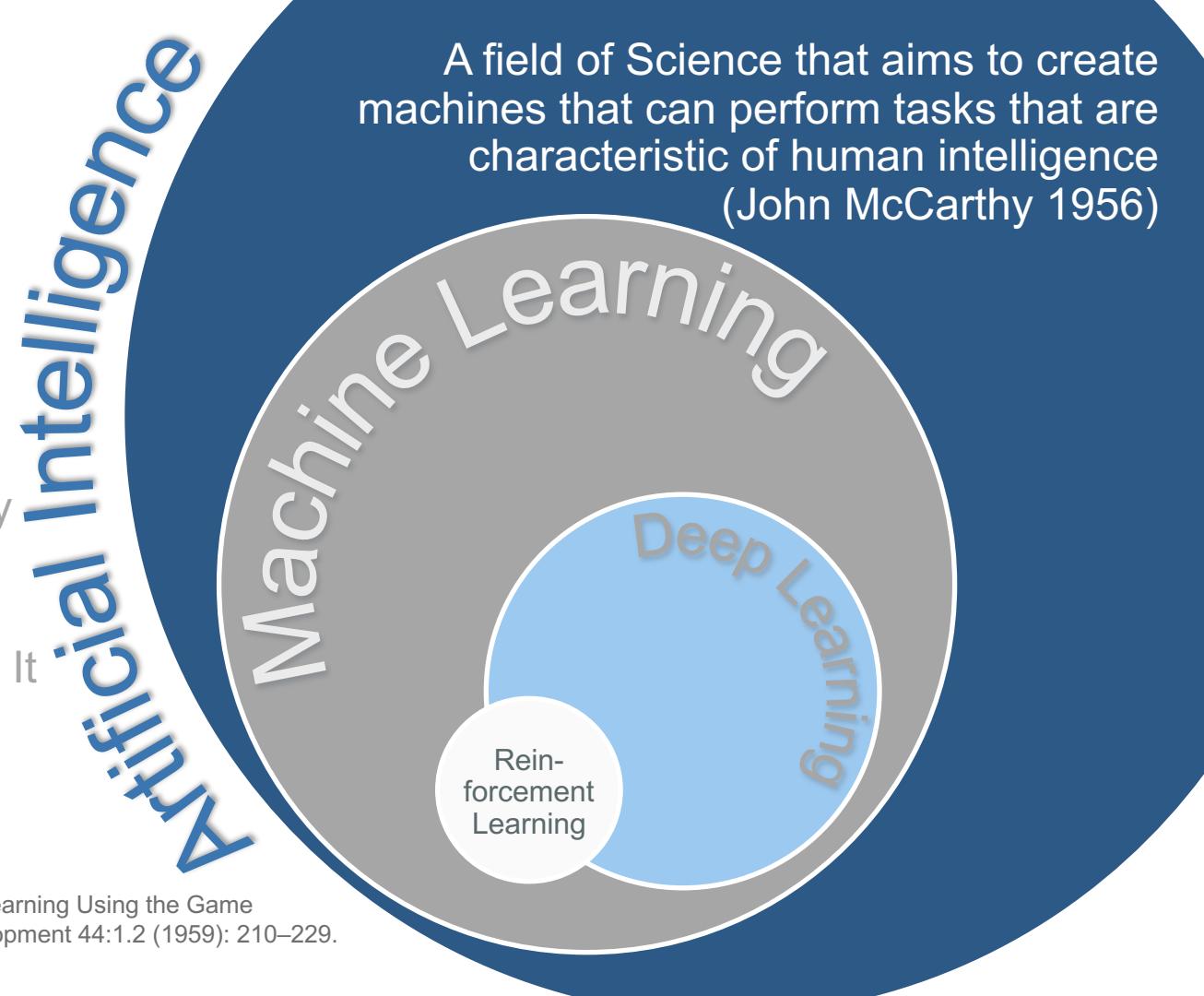
- Collaborative, strategic, proactive, creative and innovative
- Influence without authority
- Translate data-driven insights into impactful decisions and actions
- Data Visualisation

**Data Science**  
combines the complexity of  
**Software Development,**  
the challenges of applied **numerical analysis**  
with the additional **dynamic** introduced by **data!**

Machine Learning refers to a set of algorithms that allow computers to learn from data without being explicitly programmed. [1]

Deep Learning is part of a broader family of machine learning methods based on artificial neural networks. It belongs to the class of hierarchical learning algorithm.

A field of Science that aims to create machines that can perform tasks that are characteristic of human intelligence  
(John McCarthy 1956)



[1] Samuel, Arthur L. „Some Studies in Machine Learning Using the Game of Checkers,“ IBM Journal of Research and Development 44:1.2 (1959): 210–229.

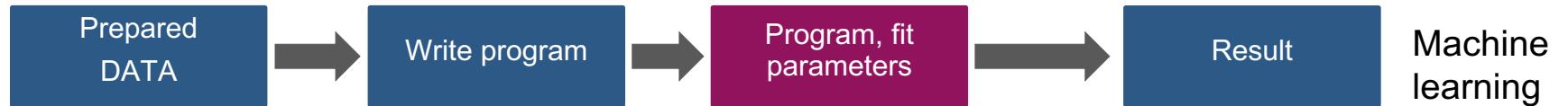
# Machine Learning: A different way of software development



## Traditional Software

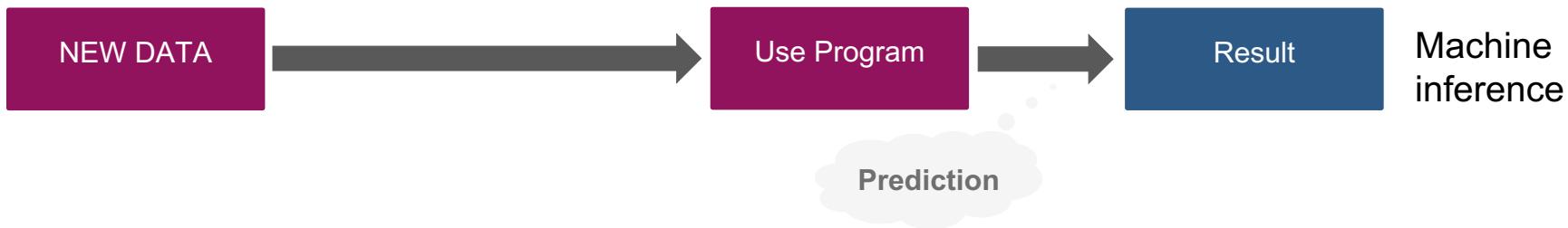


## Machine Learning Software

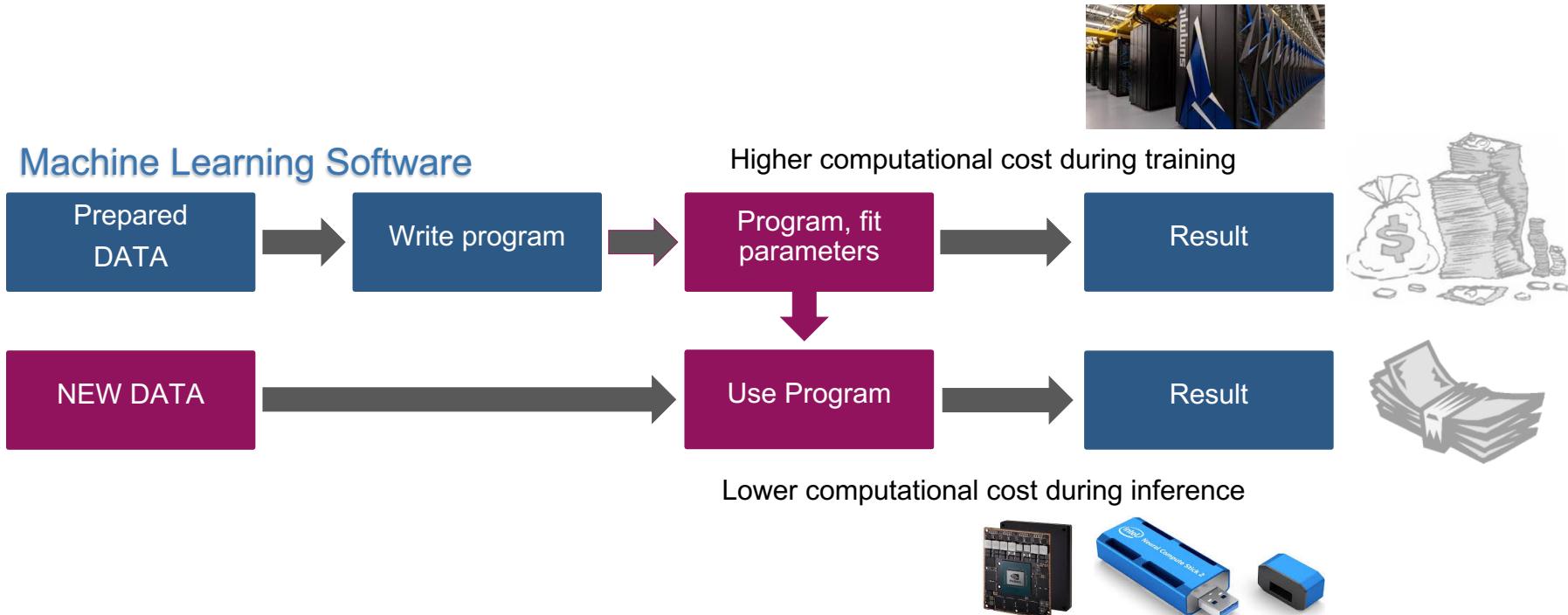


Prediction

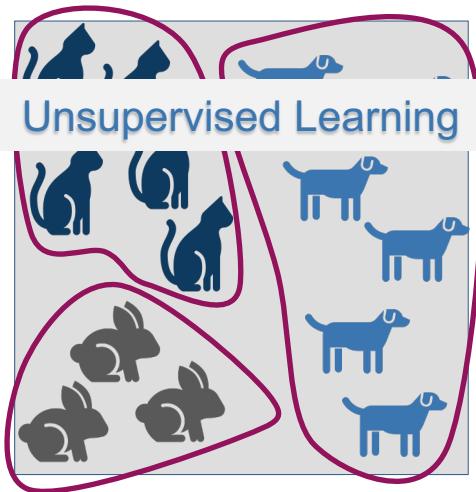
Prediction  
is the process of filling in missing information!  
Prediction takes information you have, often called "data",  
and uses it to extract information you need.



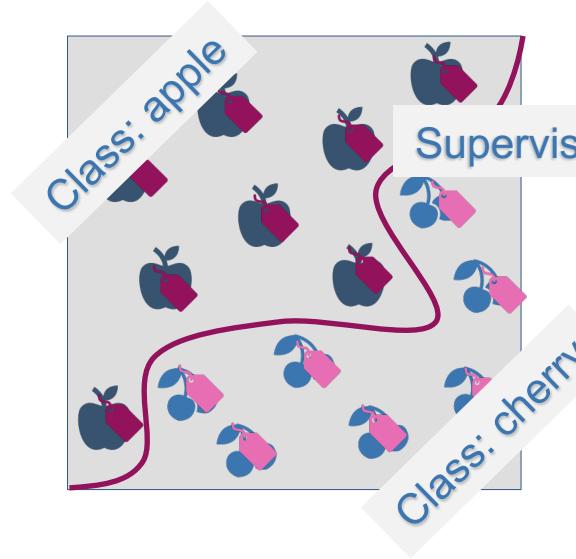
# Machine Learning cost asymmetry in training and inference



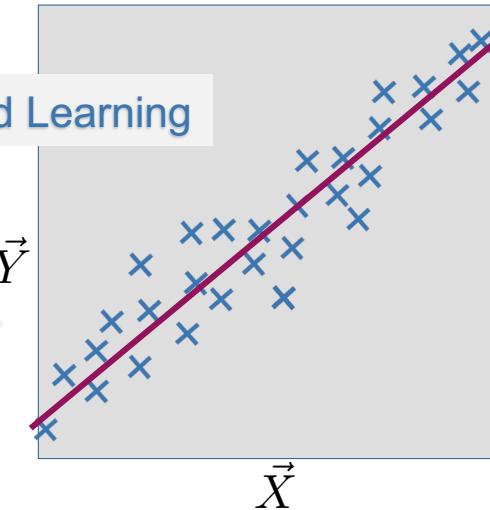
## Clustering



## Classification



## Regression



- Clustering algorithm try to find structure or pattern in uncategorized data.
- Clustering algorithm map input data to output labels.
- Regression algorithm map input data to a continuous output value.

Classification and regression model development needs labelled input data.

## Unsupervised Learning

- In unsupervised learning model, only input data is given.
- Uses unlabelled data
- Number of classes is not known.
- Less accurate and trustworthy method.
- Difficult to validate the reliability of the Machine Learning model.
- Algorithm examples:** Cluster algorithms (K-means, Hierarchical clustering)  
Signal separation (Principal & Independent component analysis)  
Neural Networks (autoencoders),  
anomaly detection (isolation forest).

## Supervised Learning

- In a supervised learning model, input and output variables are given.
- Uses labelled data
- Learning a link between the input and the outputs.
- Possible to develop highly accurate and trustworthy method.
- Difficult to get a sufficient amount of high quality labelled data.
- Algorithm examples:** Support vector machine, (deep) neural network, Linear and logistics regression, random forest, and Classification trees (forests)

## Linear Regression; the simplest example of a Machine Learning Algorithm

$$\vec{Y} \approx f(\vec{X}; \{p_1, p_2, \dots, p_N\})$$

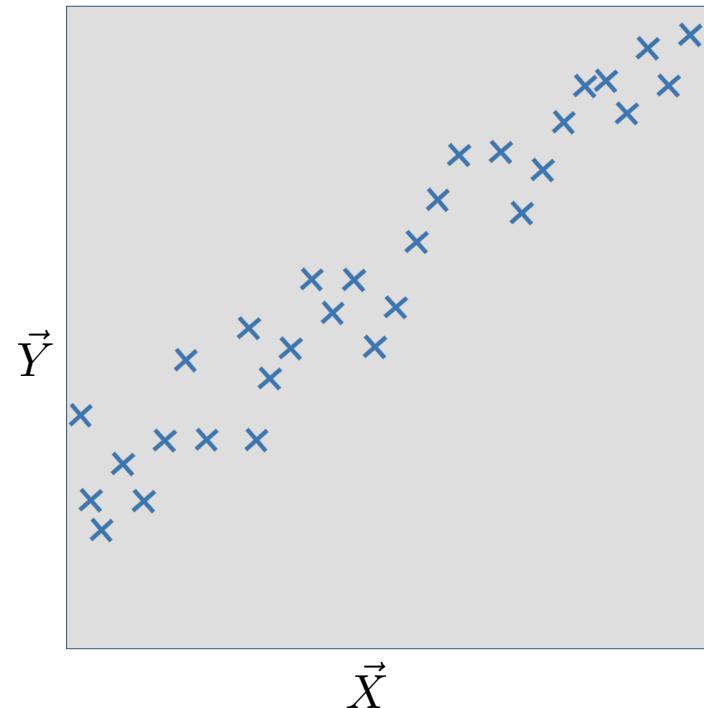
- The function  $f$  describes the relationship between  $\vec{X}$  and  $\vec{Y}$ .
- In general,  $f$  is not known.
- $\vec{X}$  and  $\vec{Y}$  can be multidimensional vectors

### Example:

x duration of a taxi drive in seconds  
y price of the Taxi drive

### Example multidimensional Input:

$\vec{X}$  ( $x_1$ ) duration and ( $x_2$ ) distance of a taxi drive  
y price of the Taxi drive



## Linear Regression; the simplest example of a Machine Learning Algorithm

$$\vec{Y} \approx f(\vec{X}; \{p_1, p_2, \dots, p_N\})$$

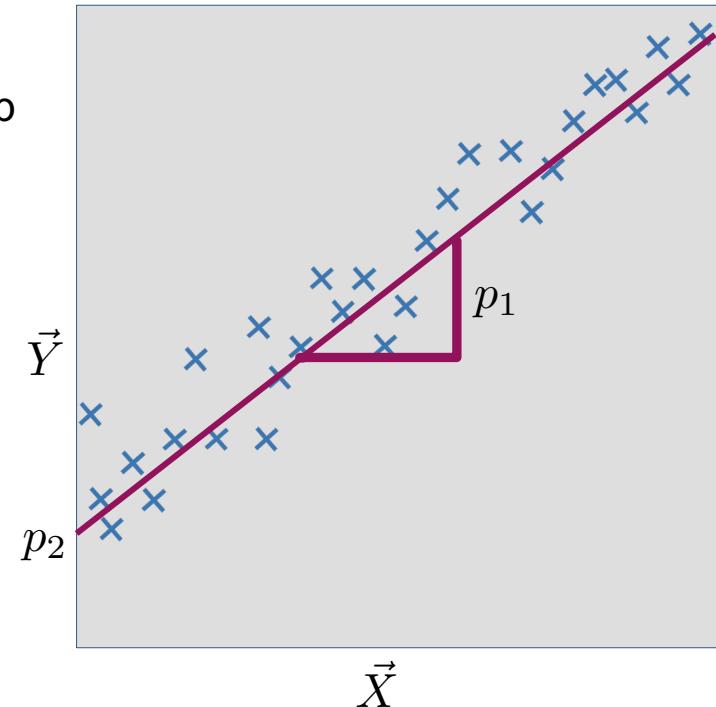
- The simplest form of  $f$  approximating the relationship between  $\vec{X}$  and  $\vec{Y}$  is a straight line.
- $\{p_1, p_2, \dots, p_N\}$  are the parameters determining the function  $f$

$$Y = f(X; p_1, p_2)$$

$$N = 2$$

$$Y = p_1 X + p_2$$

- $p_1$  slope and  $p_2$  Y-intersection of the linear equation
- The equation is fully defined by  $p_1, p_2$



## Linear Regression; the simplest example of a Machine Learning Algorithm

$$\vec{Y} \approx f(\vec{X}; \{p_1, p_2, \dots, p_N\})$$

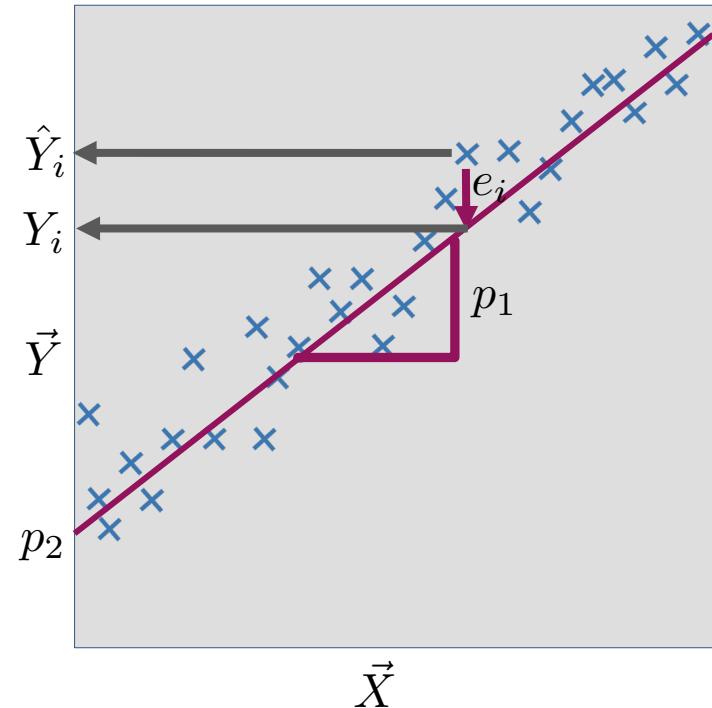
- Error  $e_i$  between the function value  $\hat{Y}_i$  and the data point  $\hat{Y}_i$

$$e_i = |\hat{Y}_i - Y_i|^2$$

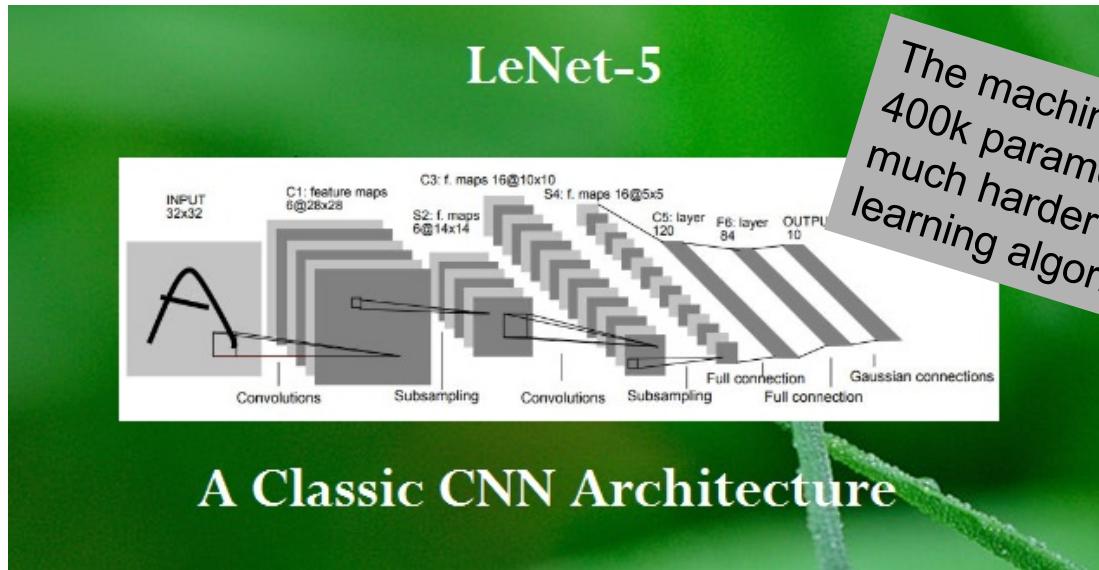
- The machine Learning part is finding  $p_1, p_2$
- $p_1, p_2$  is determined by minimising the sum of  $e_i$

$$MSE = \frac{1}{N} \sum_{i=1}^N |\hat{Y}_i - Y_i|^2$$

- Using gradient descent: updating  $p_1, p_2$  to reduce the cost function (MSE).



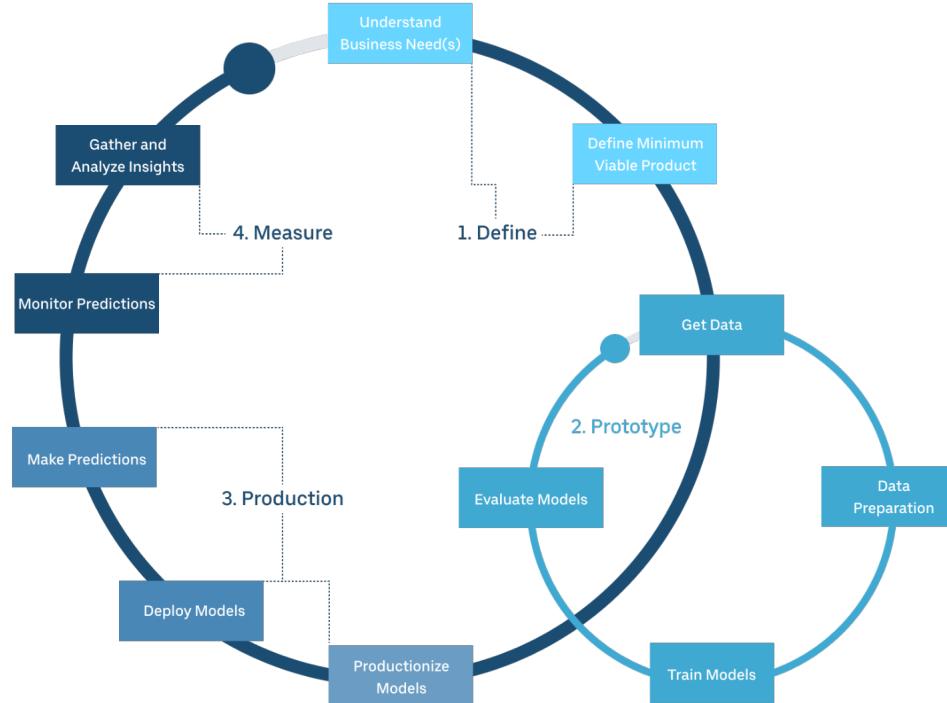
## Example of a Machine Learning Algorithm: Deep neural network for computer vision



The machine Learning part of fitting the 400k parameter  $\{p_1, p_2, \dots, p_{388272}\}$  is much harder in advanced machine learning algorithm like deep neural nets.

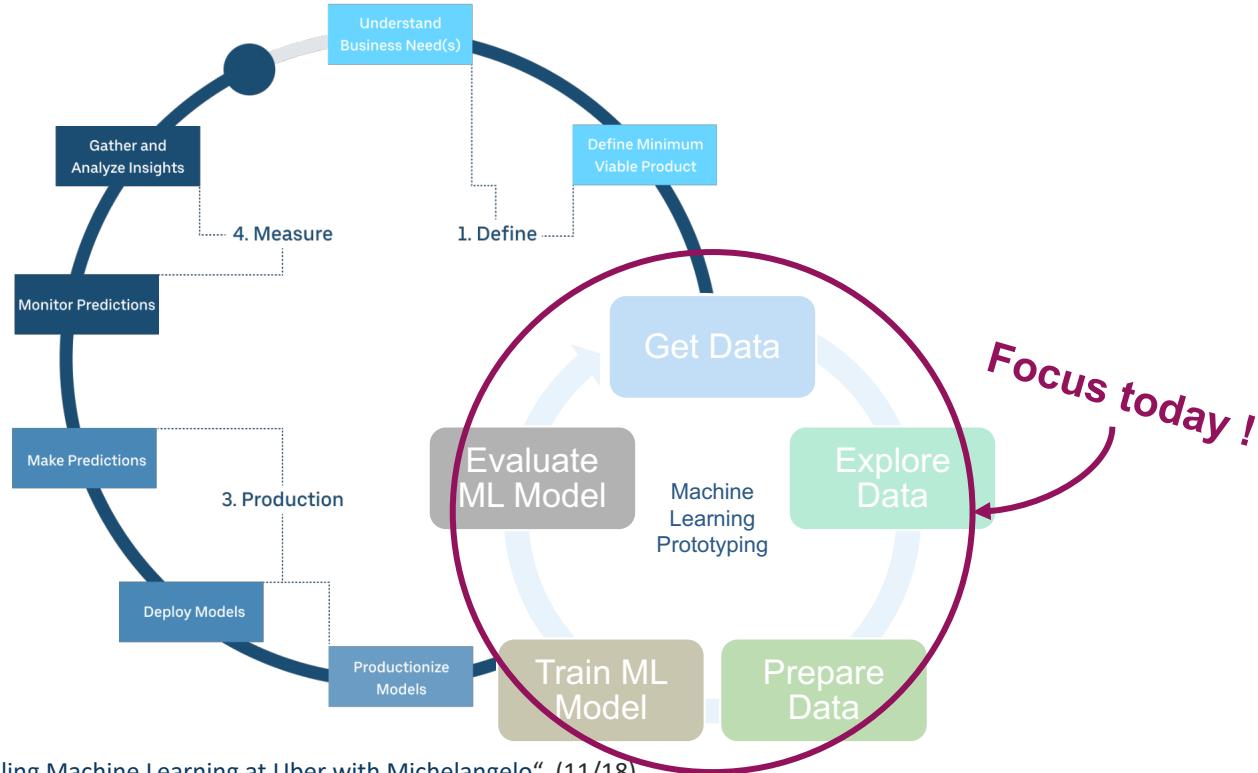
- LeNet-5 is a typical example of a deep neural network for computer vision for a  $32 \times 32$  pixel image the Net contains **~400k parameters** to optimise.

# Data Science Projects: A Highly Iterative Process



[1] Jeremy Hermann and Mike Del Balso; „[Scaling Machine Learning at Uber with Michelangelo](#)“ (11/18)

# Data Science Projects: A Highly Iterative Process



[1] Jeremy Hermann and Mike Del Balso; „[Scaling Machine Learning at Uber with Michelangelo](#)“ (11/18)

# Data Science: The surrounding Infrastructure is Vast and Complex

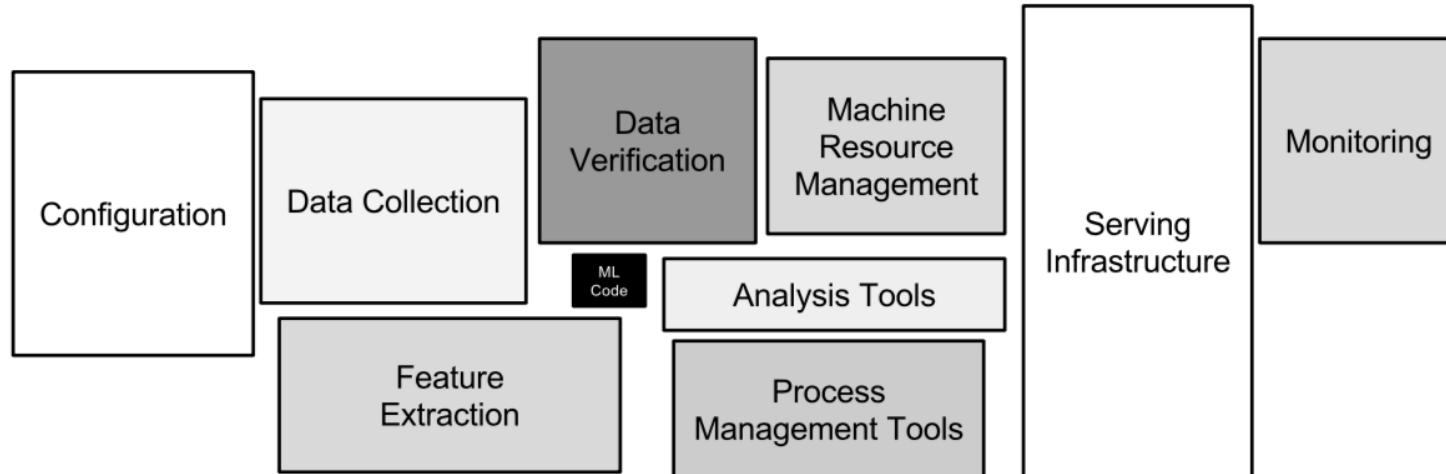


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

# The Use Case today New York City Taxi .....



# Thank you for your attention



Dr. Lydia Nemec  
Data Scientist @ Zeiss



<https://www.linkedin.com/in/lydianemec/>  
[@LydiaNemec](https://twitter.com/LydiaNemec)



Veronika Zellner  
Architect for Data & AI @Microsoft



<https://www.linkedin.com/in/veronika-zellner/>

