

# Machine Learning SoSe23

**How much?** 2 lectures and 1 tutorial per week, 9 credits

**When?** Lectures on Mon. 14.15 h & Thu. 12.15h (in person + Zoom + Recordings)

Tutorials to pick between: Wednesday, 12 pm (in person)

Thursday, 16 pm (in person)

Fridays at 10 am (in person)

## What? Theoretical & practical foundations of Machine Learning

**More information and registration (by May 31th):** <https://cms.sic.saarland/ml24/>

# ML Team



Isabel Valera  
(Lecturer)



Ayan  
Majumdar (TA)



Jonas  
Klesen (TA)

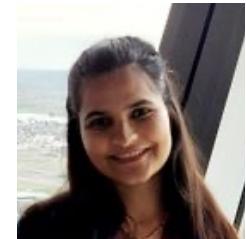


Deborah  
Kanubala (TA)



Gabriele  
Merlin (TA)

Michael Schott  
(Tutor)



Shilpa Sharma  
(Tutor)

Sofia Molotkova  
(Tutor)



Aiman Al-Azazi  
(Tutor)



Ghada Nait Said  
(Tutor)



Fatemeh Moghaddam  
(Tutor)



# ML Timeline

	Week	Date	Lecture Nr.	Title
Block I - Foundations	1	18-Apr	1	Introduction
	2	22-Apr	2	Bayesian decision Theory
	2	25-Apr	3	Empirical Risk Minimization
	3	29-Apr	4	Linear regression I
	3	2-May	5	Linear regression II
	4	6-May	6	Linear classification
	4	9-May	Holiday	No Class
	5	13-May	7	Performance measures I
	5	16-May	8	Performance measures II
	6	20-May	Holiday	No Class
Block II - Linear Models	6	23-May	9	Q&A session
	7	27-May	10	Project Description (recording)
	7	30-May	Holiday	No Class
	8	3-Jun	No Holiday	No Class
	8	6-Jun	11	Convex Optimization I
	9	10-Jun	12	Convex Optimization II
	9	13-Jun	13	Linear SVM
	10	17-Jun	14	Intro to Kernels
	10	20-Jun	15	Learning with Kernels
	11	24-Jun	16	Societal impact I
Block IV - Societal Impact	11	27-Jun	17	Societal impact II
	12	1-Jul	18	Clustering I - Kmeans
	12	4-Jul	19	Clustering II - GMMs
	13	8-Jul	20	PCA + Kernel PCA
	13	11-Jul	21	Deep learning I - MLPs + BP
Block V - Unsupervised Learning	14	15-Jul	22	Deep learning II - CNNs
	14	18-Jul	23	Deep learning III - Transformers
	15	22-Jul	24	Beyond supervised DL
	15	25-Jul	25	Q&A
Block VI - Intro to Deep Learning				

## Pre-requisites:

- linear algebra
- good basis on statistics

## Tutorials:

- Help tutorials
- Solution tutorials
- Programming tutorials

## Evaluation:

- Final written exam
- Project Report
- (Bonus) Project Challenge

# ML Tutorials

If you would like to attend the ML tutorials, please **register to (ONLY!) one of the tutorial slots by Tuesday 23rd via CMS!!**

<b>Id</b>	<b>Name</b>	<b># Students</b>	<b>Start</b>	<b>End</b>
1	<b>Tutorial group I -- Wednesdays at 14 hours</b>	0/150	18.04.24	23.04.24
2	<b>Tutorial group II -- Thursdays 16 hours</b>	0/99	18.04.24	23.04.24
3	<b>Tutorial group III -- Fridays at 10 hours</b>	0/99	18.04.24	23.04.24

# ML Evaluation

- Grade = Max( 0.75\* Final written exam + 0.25\* Project report, Final written exam )
- Grade Bonus for top 5% teams in Project Challenge
- For exam registration, no admission pre-requisites (beyond University registration)
- Project voluntary but strongly recommended!  
(more details on May 27th)

# ML Q&A

Please ask your questions using one of the following options:

- Ask during or after (if private request) class
- Ask during tutorials to our tutors
- Use the cms forum
- Do not send me emails, please!

Courses Main Page Information ▾ Registration Personal Status **Forum** Internal ▾ Configuration ▾  ivalera ▾



**Machine Learning**  
Isabel Valera

Registration for this course is open until **Friday, 31.05.2024 23:59**.

## News

# ML Materials



## Machine Learning

Isabel Valera

### Materials

+ New Category

Edit access rights

Background Material (We will upload here some material to support you refreshing your previous knowledge necessary for the course. )



ProbabilityTheorySummary

(540 KB, rev 1)



Set of slides summarizing the main definitions and concepts from Probability Theory necessary for the ML course.

Recap on linear algebra

(717 KB, rev 1)



Professor Wolf kindly agreed to share the script of her lecture so that those that did not participate in the course can revisit easily their knowledge on Probability Theory. The content of this script is part of the per-requisites.

You can also revisit your knowledge by watching the videos from the Statistics Lab (please do not distribute):

<https://www.youtube.com/playlist?list=PLSM40hbPqaQbstBPmyA4f9dV-p4ybLus>

<https://www.youtube.com/playlist?list=PLSM40hbPqaQb8WWqZxIf7xKzTcUUj8s>

<https://www.youtube.com/playlist?list=PLSM40hbPqaQbYSLQvUDGf9wuOfjsKkSm4>

<https://www.youtube.com/playlist?list=PLSM40hbPqaQy0UsDbwaCMaakoUFKQ9Y5B>

[https://www.youtube.com/playlist?list=PLSM40hbPqaQzOAYqFwsHNsmf\\_a5xelYq](https://www.youtube.com/playlist?list=PLSM40hbPqaQzOAYqFwsHNsmf_a5xelYq)

<https://www.youtube.com/playlist?list=PLSM40hbPqaQzCakhmScxCS5LsACY0ouSc>

[https://www.youtube.com/playlist?list=PLSM40hbPqaQbjvMNjWu\\_lfj0nOGhAQv5](https://www.youtube.com/playlist?list=PLSM40hbPqaQbjvMNjWu_lfj0nOGhAQv5)

<https://www.youtube.com/playlist?list=PLSM40hbPqaQybz0N2NY-FUBAscCSGo0Zn>

[https://www.youtube.com/playlist?list=PLSM40hbPqaQaozdzHAyUe2iKgrRQwaa\\_hb](https://www.youtube.com/playlist?list=PLSM40hbPqaQaozdzHAyUe2iKgrRQwaa_hb)

<https://www.youtube.com/playlist?list=PLSM40hbPqaQz4Uoq1h1K3QoiPPPF0oKf4>

<https://www.youtube.com/playlist?list=PLSM40hbPqaQz57MIE6p4RDTNWT0hiewr>

<https://www.youtube.com/playlist?list=PLSM40hbPqaQyR-YtnFa2WjIS6T6FP83n9>

[https://www.youtube.com/playlist?list=PLSM40hbPqaQZ0lymlcd\\_AipM8JVmFAMNo](https://www.youtube.com/playlist?list=PLSM40hbPqaQZ0lymlcd_AipM8JVmFAMNo)

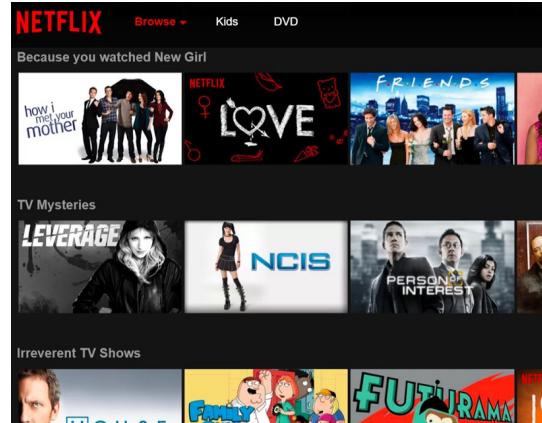


# ML is everywhere

## Frequently Bought Together



[Amazon]



[AppleNews]

## The New York Times *Banking Start-Ups Adopt New Tools for Lending*

“Big-data lending, though, relies on software algorithms...”



## AI and Personalization in Insurance: Customizing Policies and Customer Experiences

Monday, 10/04/2023 | 15:19 GMT by Finance Magnates Staff

The AI revolution is here

FORBES > LEADERSHIP

## HR Recruiting Software In The Era Of AI

DeepL Translator ▾ DeepL Pro For Business ▾ Why DeepL? API Plans and pricing Apps

Translate text 31 languages

Translate files .pdf, .docx, .pptx

Detect language ▾ German ▾

[AppleNews]

[DeepL.com]

 Maxime Legardez Coquin Forbes Councils Member  
Forbes Human Resources Council  
COUNCIL POST | Membership (Fee-based)

Apr 11, 2023, 07:30am EDT

CEO of *Maki People*, a skills assessment tool to hire better and faster.

# ML Definitions



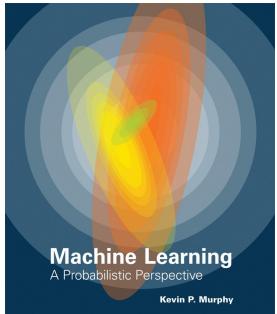
(Picture from Wikipedia)

Arthur Samuel (1959): *“Machine Learning is the field of study that gives the computer the ability to learn without being explicitly programmed.”*



(Image from Tom Mitchell's homepage)

Tom Mitchell (1998): *“a computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”*



Kevin Murphy's book (2012): *“we define machine learning as a set of methods that can automatically detect patterns in data, and then use the uncovered patterns to predict future data, or to perform other kinds of decision making under uncertainty (such as planning how to collect more data!)”*

# ML Definition by ChatGPT



what is machine learning?



Machine learning is a subfield of artificial intelligence that involves developing algorithms and models that can automatically learn from data, without being explicitly programmed. It involves building mathematical models and algorithms that can identify patterns and make predictions or decisions based on the data.

The process of machine learning typically involves several steps, including collecting and preparing data, selecting a model or algorithm, training the model on the data, testing the model to evaluate its performance, and then deploying the model to make predictions or decisions on new data.

There are several different types of machine learning, including supervised learning, unsupervised learning, and reinforcement learning. In supervised learning, the model is trained on labeled data, where the correct output is already known, and the goal is to learn to generalize to new, unseen data. In unsupervised learning, the model is trained on unlabeled data, where there is no predetermined correct output, and the goal is to discover underlying patterns or structure in the data. Reinforcement learning involves training a model to interact with an environment and learn through trial and error to maximize a reward signal.



why is it important to learn the theoretical foundations of machine learning?



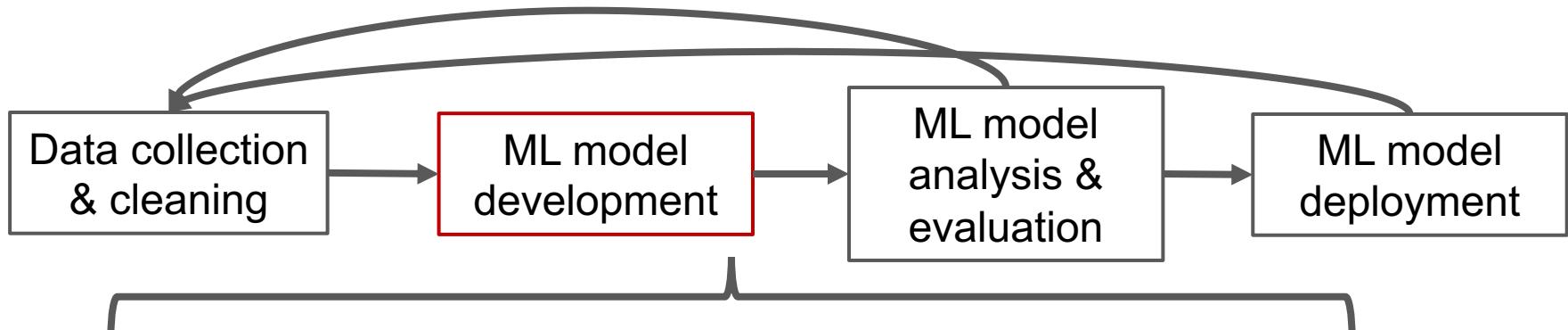
It is important to learn the theoretical foundations of machine learning because it provides a deep understanding of the underlying concepts and principles that drive the development and application of machine learning models. Understanding the theory behind machine learning can help practitioners develop better models, make more informed decisions, and avoid common pitfalls.

Some of the key benefits of learning the theoretical foundations of machine learning include:

1. Better model selection: Understanding the theoretical foundations of machine learning can help practitioners select the right model for a given problem. This can help improve the accuracy and efficiency of the model.
2. Robustness: Understanding the theory behind machine learning can help practitioners develop models that are more robust to changes in the data or the environment. This can help ensure that the model continues to perform well even when faced with new or unexpected situations.
3. Interpretability: Understanding the theory behind machine learning can help practitioners develop models that are more interpretable, meaning that they can provide insights into how the model is making its predictions or decisions.
4. Generalization: Understanding the theoretical foundations of machine learning can help practitioners develop models that generalize well to new, unseen data. This is critical for ensuring that the model can be applied to real-world problems.
5. Innovation: Understanding the theoretical foundations of machine learning can help practitioners develop new algorithms and techniques that can push the field forward and lead to new breakthroughs.

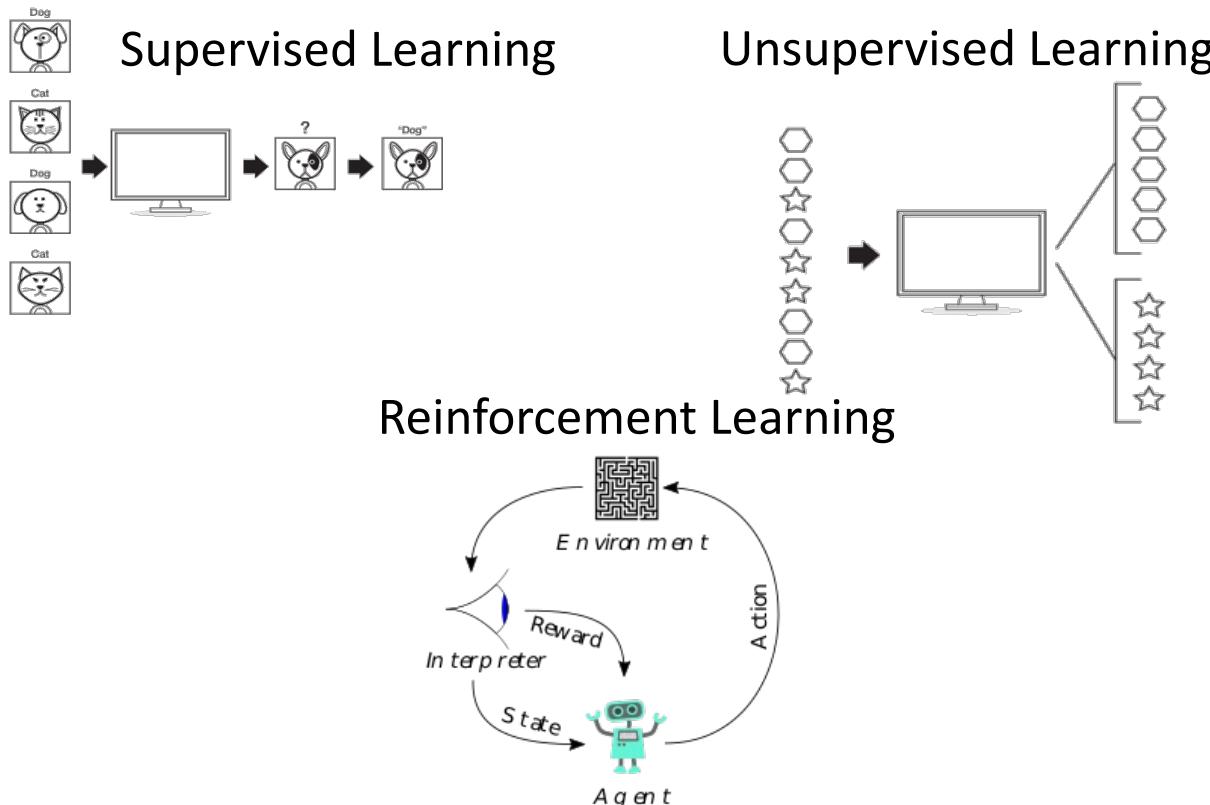
In short, learning the theoretical foundations of machine learning is critical for anyone who wants to develop and apply machine learning models effectively and efficiently.

# ML Pipeline

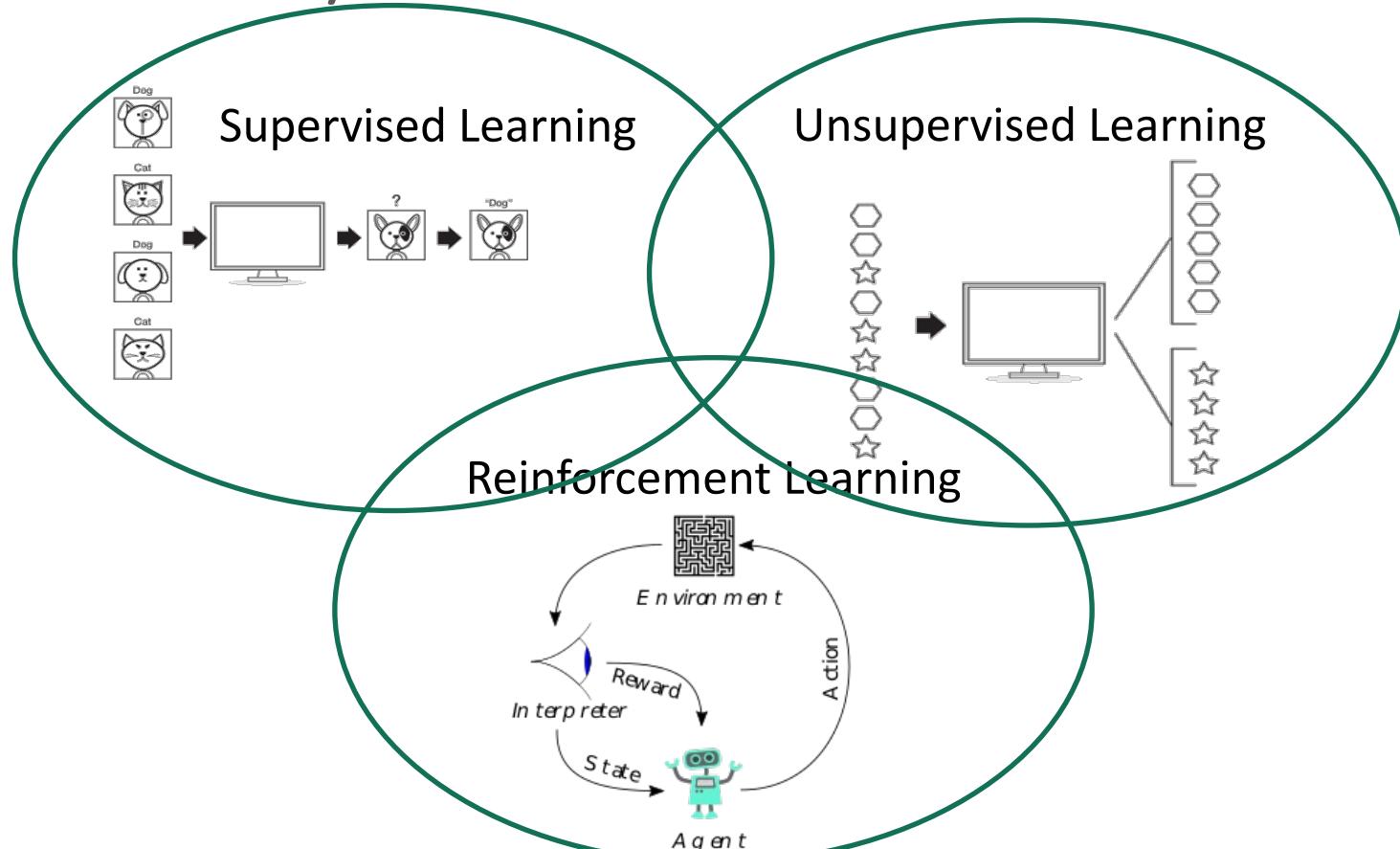


- Data pre-processing (e.g., data normalization, feature selection, etc.)
- **Model family choice (based on task, data properties, etc.)**
- **Model training & validation**

# ML Taxonomy



# ML Taxonomy



# ML Taxonomy

## What do you want the machine learning system to do?

I want to see if there are natural clusters or dimensions in the data I have about different situations.

I want to learn what actions to take in different situations.

Do you want the ML system to be active or passive?

**ACTIVE**

The system's own actions will affect the situations it sees in the future.

**PASSIVE**

The system will learn from data I give it.

Do you have access to data that describes a lot of examples of situations and appropriate actions for each situation?

Will the system be able to gather a lot of data by trying sequences of actions in many different situations and seeing the results?

Could there be patterns in these situations that humans haven't recognized before?

Could a knowledgeable human decide what actions to take based on the data you have about the situation?

No  
Yes

No

Yes

**UNSUPERVISED LEARNING MAY BE APPROPRIATE**

clustering  
anomaly detection

**SUPERVISED LEARNING MAY BE APPROPRIATE**

neural nets  
support vector machines  
regression  
recommender systems

**MACHINE LEARNING IS NOT USEFUL**

**REINFORCEMENT LEARNING MAY BE APPROPRIATE**

Credit: Thomas Malone, MIT Sloan | Design: Laura Wentzel

# ML, the basics (Block I)

- Probability theory

Bishop's book: *"A key concept in the field of pattern recognition is that of uncertainty. It arises both through noise on measurements, as well as through the finite size of data sets. Probability theory provides a consistent framework for the quantification and manipulation of uncertainty and forms one of the central foundations for pattern recognition."*

- Bayesian decision Theory

Bishop's book: *"decision theory that, when combined with probability theory, allows us to make optimal decisions in situations involving uncertainty such as those encountered in pattern recognition."*

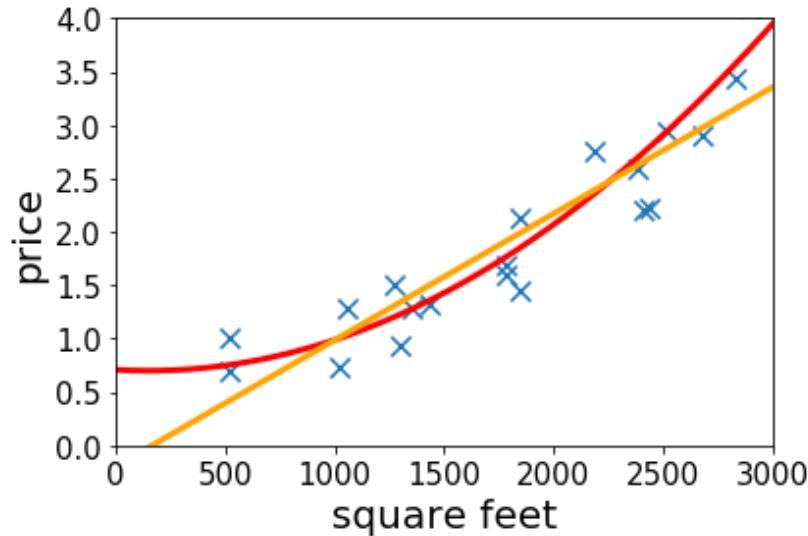
- Empirical Risk Minimization

Vapnik 1991: *"Learning is posed as a problem of function estimation, for which two principles of solution are considered: empirical risk minimization and structural risk minimization."*

# Supervised learning (Blocks II & III)

- Given a dataset that contains samples (or features/outcomes pairs):  
 $(x_n, y_n)_{n=1 \dots N}$
- Task: Predict outcome ( $y$ ) from features ( $x$ ).
- Supervision comes from  $y$ !

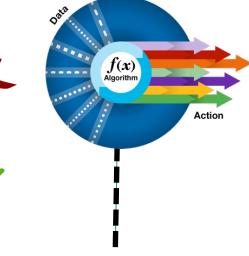
**Example of regression problem:** We observe the square feet of a house and its price. We aim to predict the price of unseen houses based on their size.



# Supervised learning (Blocks II & III)

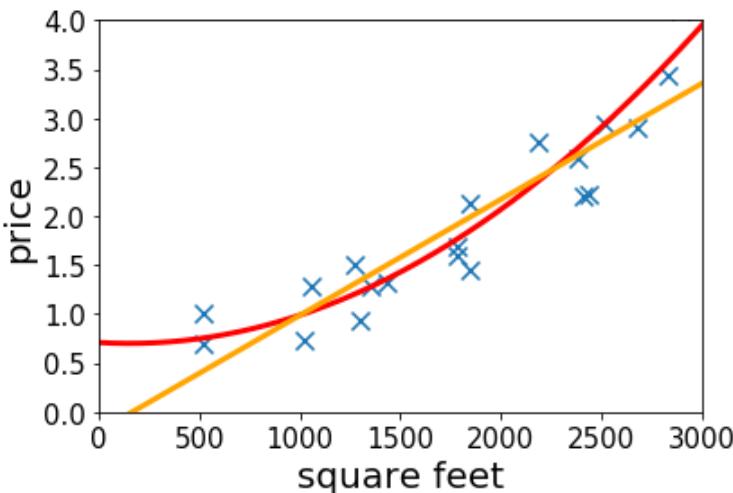
**Example of classification problem:** We observe a set of features characterizing every loan/credit applicant together with information on whether they repaid or defaulted to repay the credit in the past. The algorithm aims to predict for a new applicant if he/she will repay based on the set of features.

	x (features)			y (label)
	Income	Credit card	Previous loans	Repaid?
Applicant 1	1.5k			
Applicant 2	2.7k			
Applicant 3	2.2k			



# Supervised learning – Classification & Regression

- Regression: if  $y$  is a continuous variable
  - e.g., price prediction
- Classification: the label is a discrete (nominal) variable
  - e.g., the task of predicting loan repayment



	x (features)			y (label)
	Income	Credit card	Previous loans	Repaid?
Applicant 1	1.5k	VISA	✓	PAID IN FULL
Applicant 2	2.7k	VISA	✗	default?
Applicant 3	2.2k	VISA	✓	PAID IN FULL

# Supervised learning – Societal Impact (Block IV)

TIME

## Google Has a Striking History of Bias Against Black Girls

MIT News  
ON CAMPUS AND AROUND THE WORLD

The privacy risks of compiling mobility data  
Merging different types of location-stamped data can make it easier to discern users' identities, even when the data is anonymized.



INDEPENDENT

London terror attack: Uber slammed for being slow to turn off 'surge pricing' after rampage

Bloomberg Businessweek

Artificial Intelligence Has Some Explaining to Do

Healthcare IT News

AI in healthcare - not so fast?  
Study outlines challenges, dangers for machine learning

BBC

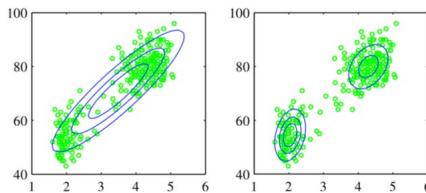
AAAS: Machine learning 'causing science crisis'

# Unsupervised learning (Block V)

- Dataset contains **no labels**.
  - Goal (vaguely-posed): to find interesting structures/patterns in the data
  - Examples:

## Generate new samples

## Density estimation



## Data exploration



In the ML course:

- Clustering: Kmeans & GMMs
  - (Kernel) PCA

# Supervised learning – Deep Learning (Block VI)

## ▪ Image Classification

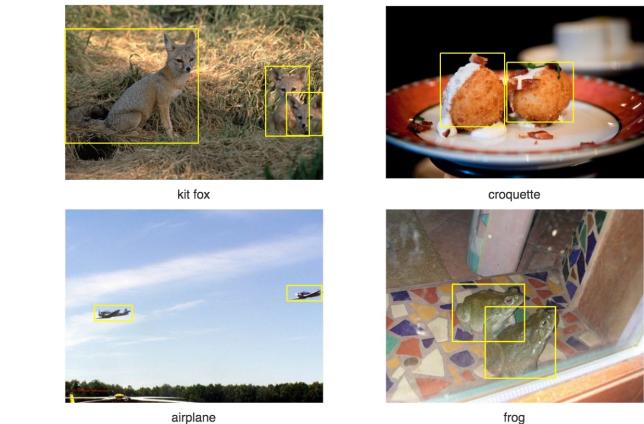
- raw pixels of the image,

ILSVRC



## ▪ Object location and detection

- raw pixels of the image
- bounding box



# Beyond supervised learning – Overview (Last lecture)

