

# INTRO to DATA SCIENCE

## LECTURE 2: MACHINE LEARNING

## **LAST TIME:**

- FIRST LOOK AT DATA SCIENCE & THE DATA MINING WORKFLOW**
- DATA VISUALIZATION WITH R & GGPLOT2**
- FIRST LINEAR MODEL**

## **QUESTIONS?**

## What's big data?

The practical viewpoint:

- ①  $O(n^2)$  algorithm feasible: small data
- ② Fits on one machine: medium data
- ③ Doesn't fit on one machine: big data

**I. WHAT IS MACHINE LEARNING?**

**II. MACHINE LEARNING PROBLEMS**

**EXERCISES:**

**III. MULTIPLE REGRESSION & FEATURE EXTRACTION**

# **I. WHAT IS MACHINE LEARNING?**

---

# WHAT IS MACHINE LEARNING?

---

6

from Wikipedia:

“Machine learning, a branch of artificial intelligence, is about the construction and study of systems that *can learn from data*.”

source: [http://en.wikipedia.org/wiki/Machine\\_learning](http://en.wikipedia.org/wiki/Machine_learning)

from Wikipedia:

“Machine learning, a branch of artificial intelligence, is about the construction and study of systems that *can learn from data*.”

“The core of machine learning deals with *representation* and *generalization*...”

- *representation* – extracting structure from data
- *generalization* – making predictions from data

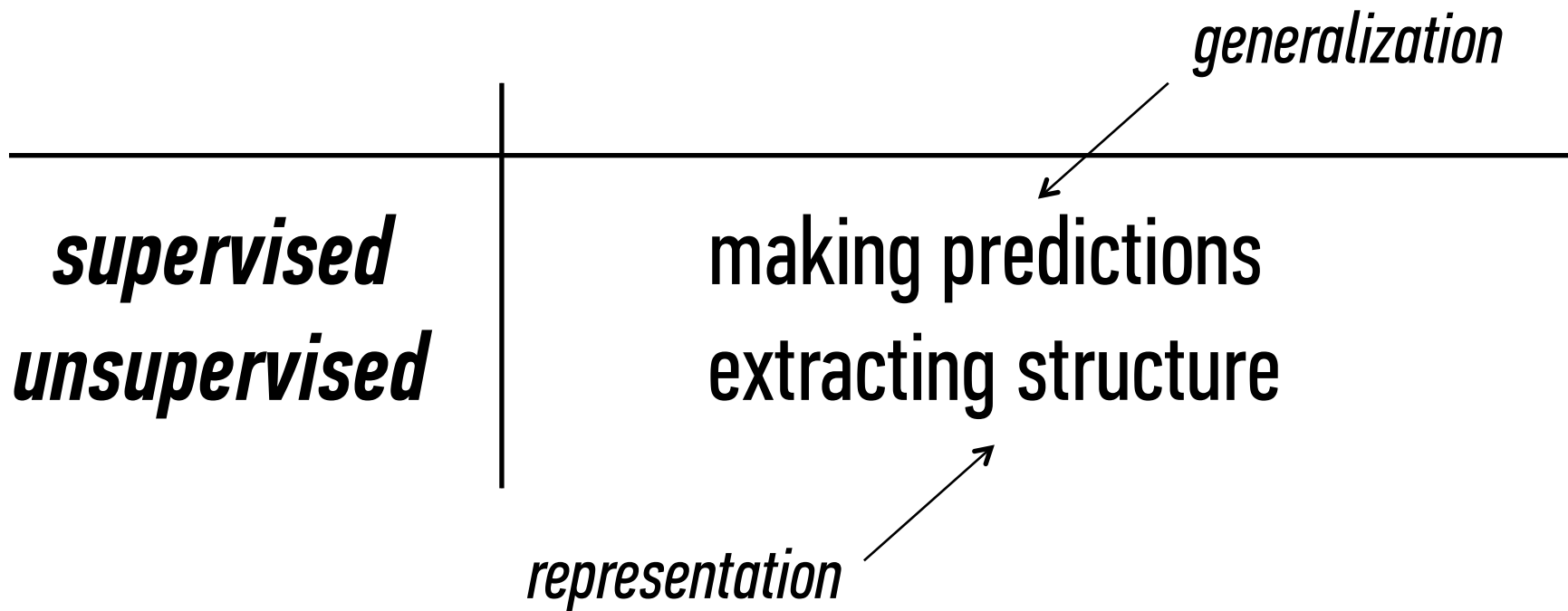
source: [http://en.wikipedia.org/wiki/Machine\\_learning](http://en.wikipedia.org/wiki/Machine_learning)

# **II. MACHINE LEARNING PROBLEMS**



---

<i><b>supervised</b></i>	making predictions
<i><b>unsupervised</b></i>	extracting structure



	<b><i>continuous</i></b>	<b><i>categorical</i></b>
	<b>quantitative</b>	<b>qualitative</b>

*continuous*

*categorical*

quantitative

qualitative

## NOTE

The space where data live is called the *feature space*.

Each point in this space is called a *record*.

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	regression	classification
<i>unsupervised</i>	dimension reduction	clustering

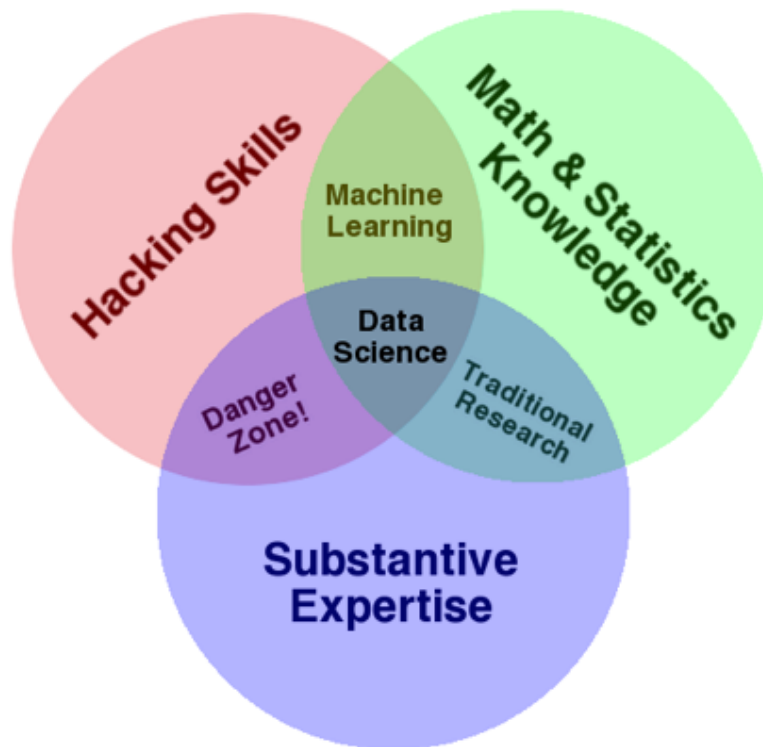
	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	regression	classification
<i>unsupervised</i>	dimension reduction	clustering

## NOTE

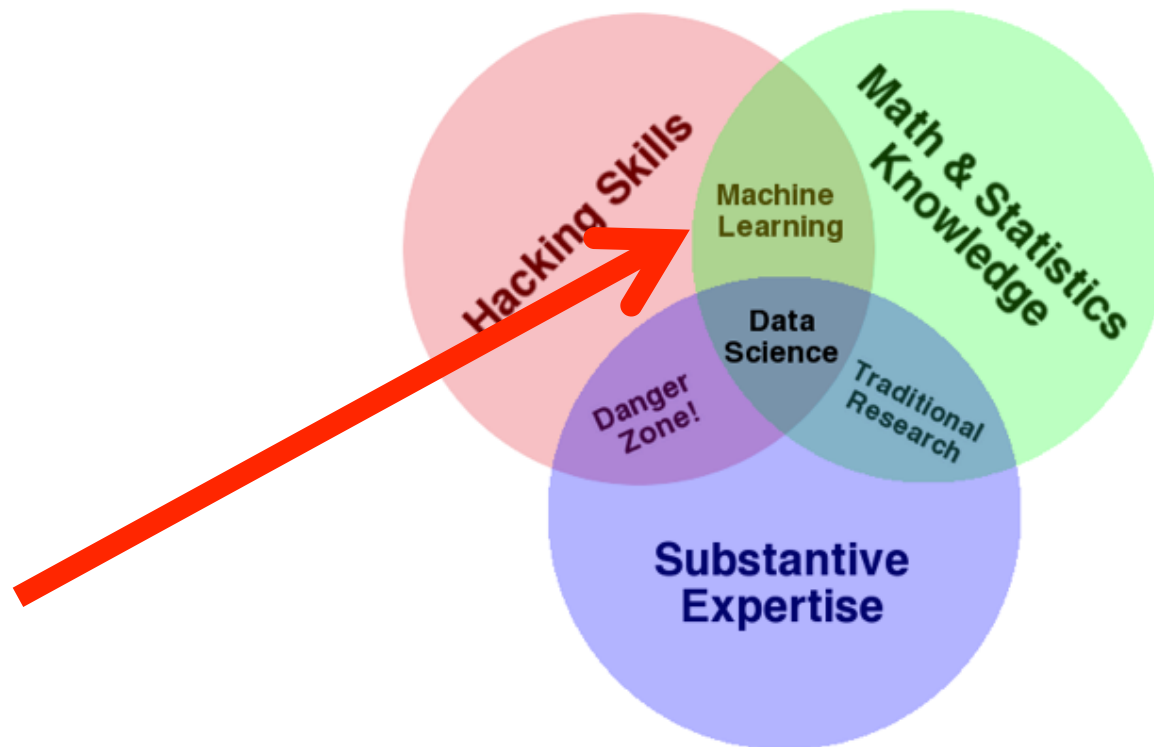
We will implement solutions using *models* and *algorithms*.

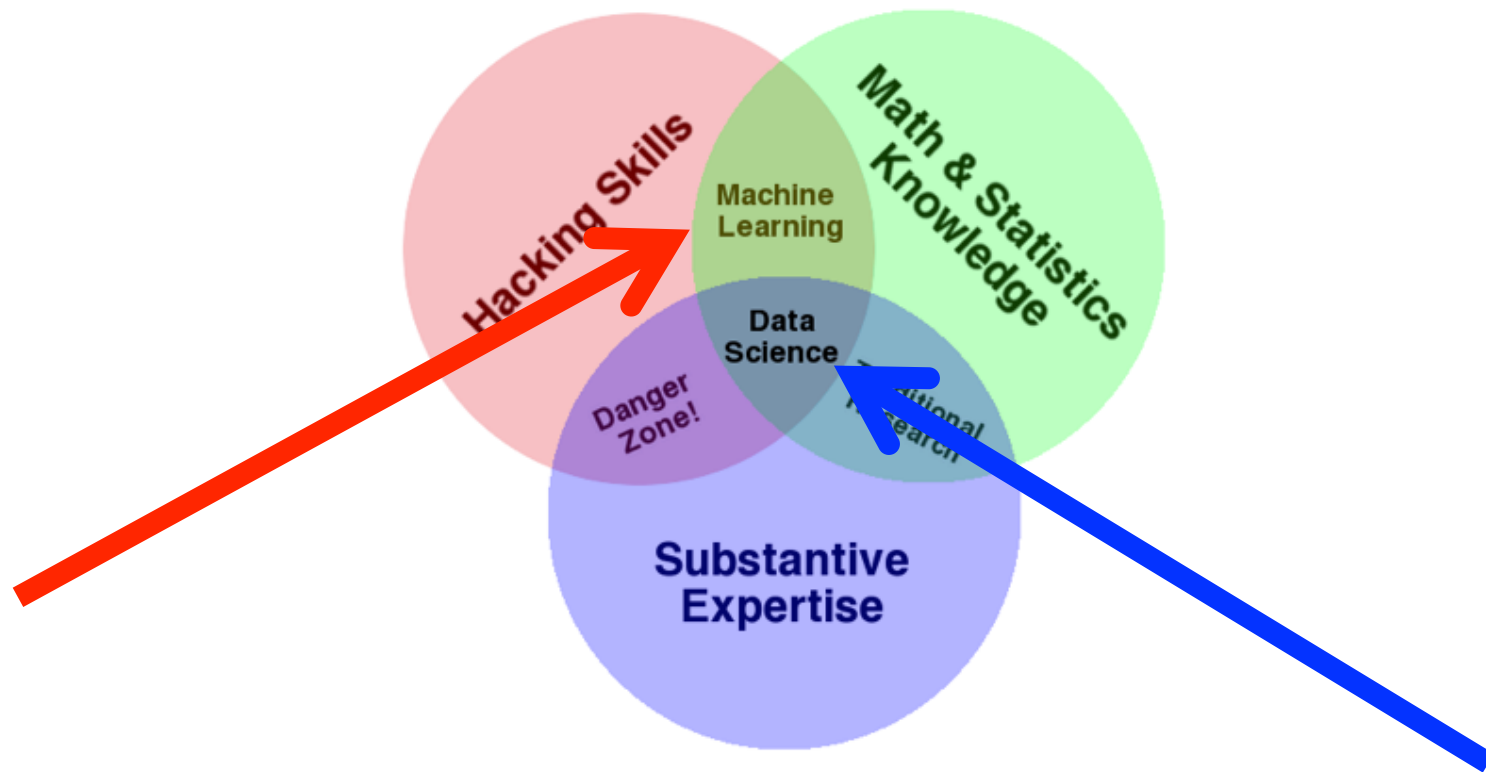
Each will fall into one of these four buckets depending on the type of problem and type of data.

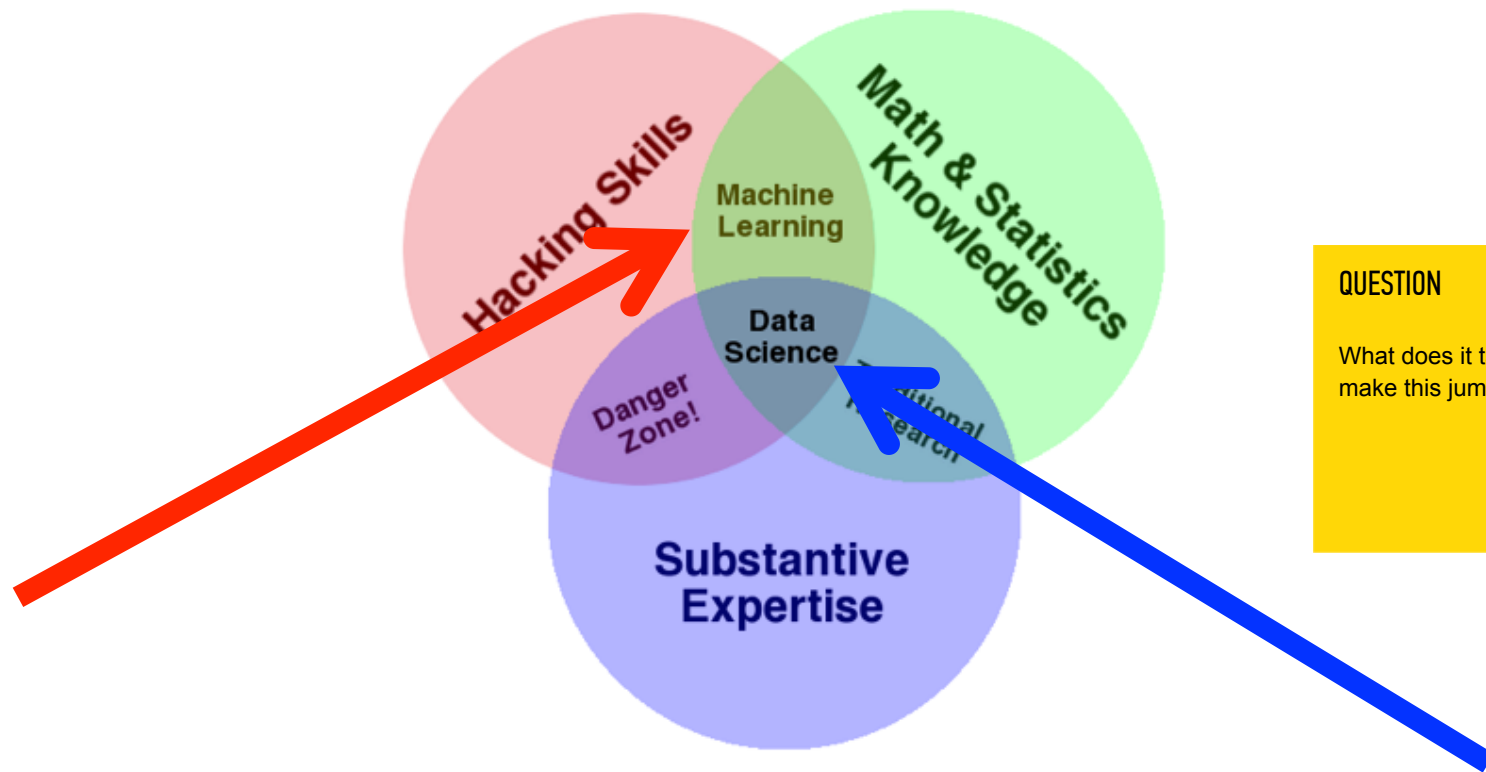
# ***DATA SCIENCE AND MACHINE LEARNING***





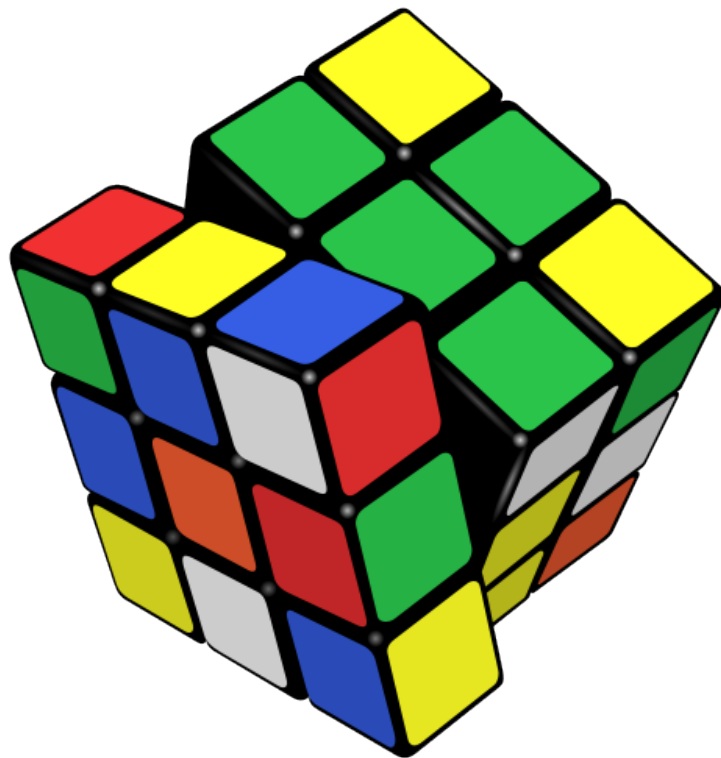






## QUESTION

What does it take to make this jump?



## NOTE

Implementing solutions to ML problems is the focus of this course!

# **III. RELATIONSHIPS AMONG SEVERAL VARIABLES**

---

## EXERCISE – MULTIPLE REGRESSION (BACKWARD ELIMINATION)

---

22

### KEY OBJECTIVES

---

- Create a regression model using several independent variables
- Extract meaningful features

### TOOLS

---

- R (plot, lm, update)

- 1) **Linearity** of the relationship between dependent and independent variables (doesn't mean the relation between  $y$  and  $x$  has to be linear since we can use transformations if  $y$  and  $x$  as well)
- 2) **Independence** of the errors
- 3) **Homoscedasticity** (constant variance of the errors)
  - 1) versus time
  - 2) Versus the predictions or any independent variables
- 4) **Normality** of the error distribution

---

**INTRO TO DATA SCIENCE**

---

**DISCUSSION**