



# 机器学习与人工智能 Machine Learning and Artificial Intelligence

Lecture 1 Introduction/Overview

Yingjie Zhang (张颖婕)

Peking University

yingjiezhang@gsm.pku.edu.cn

2021 Fall

# Me...

Emails me: start with "[MLAI]" in the subject title



- Email: [yingjiezhang@gsm.pku.edu.cn](mailto:yingjiezhang@gsm.pku.edu.cn)
- Office hour: Wed 10-11am; or by appointment
- Website: [sites.google.com/view/yingjiezhang/home](https://sites.google.com/view/yingjiezhang/home)
- My Research:
  - Topics: Mobile and Sensor Technologies, Big Data and Smart City, User-generated Content, Sharing Economy, and Social Media.
  - Methodologies: Econometrics, Machine Learning, Text Mining, and Field Experiment.

# Teaching Assistant

## Guangxin Yang (杨广鑫)

- 2<sup>nd</sup> PhD Student in Marketing
- Email: [ygx@stu.pku.edu.cn](mailto:ygx@stu.pku.edu.cn).
- TA session: Saturday 10:00-11:00; or by appointment
- A novice at ML with you guys. Debug Python code together!

# What is MACHINE LEARNING

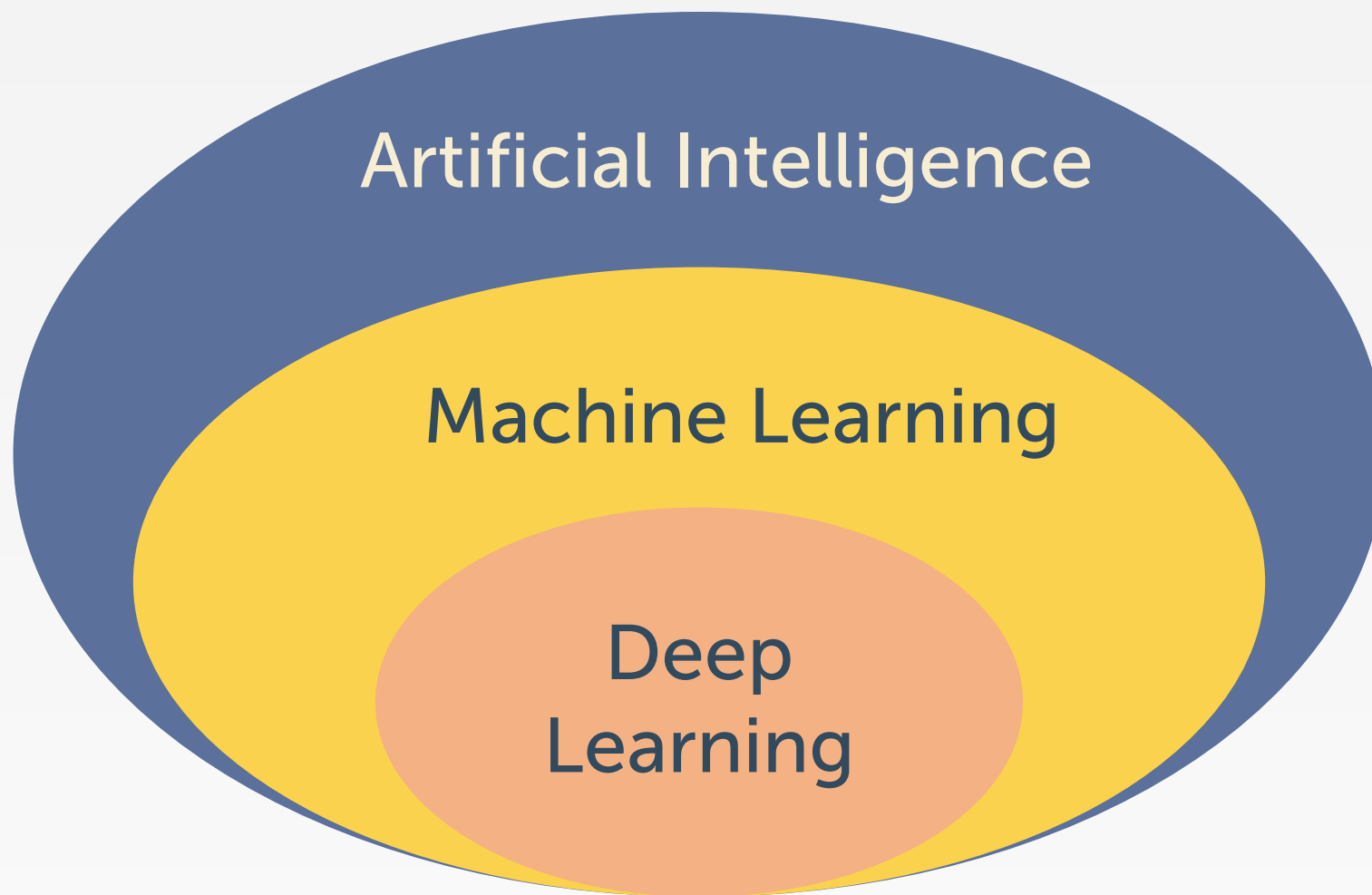
# Artificial Intelligence



Artificial Intelligence

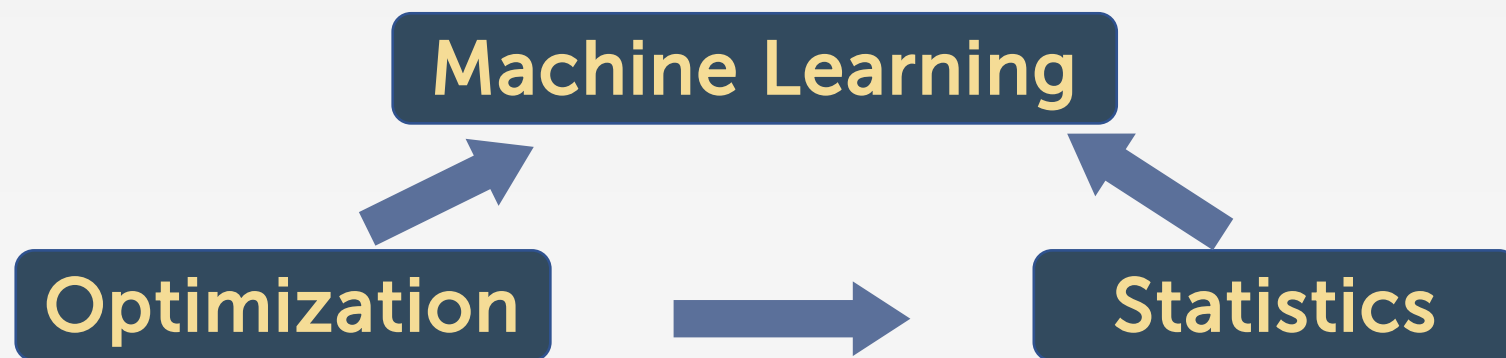
The diagram consists of two nested ellipses. The outer ellipse is blue and contains the text 'Artificial Intelligence'. The inner ellipse is yellow and contains the text 'Machine Learning'. The yellow ellipse is entirely contained within the blue ellipse, illustrating that Machine Learning is a subset of Artificial Intelligence.

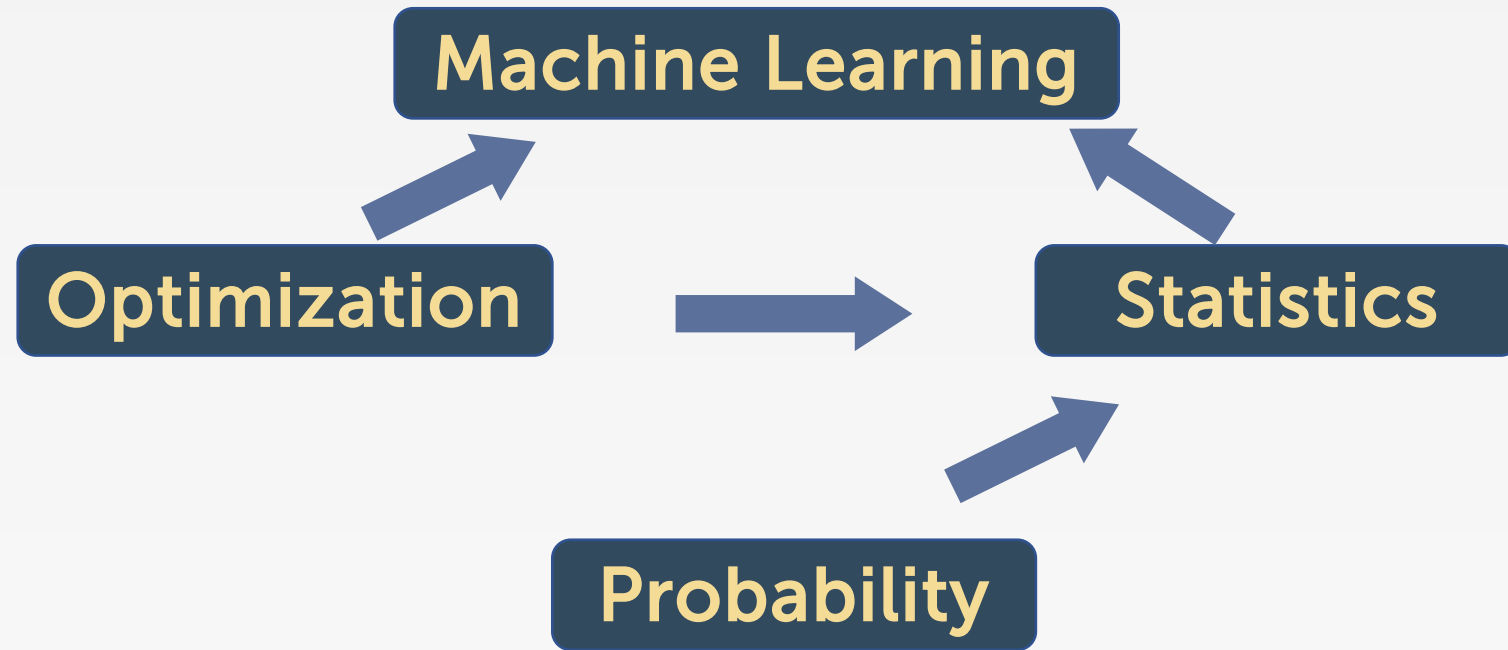
Machine Learning

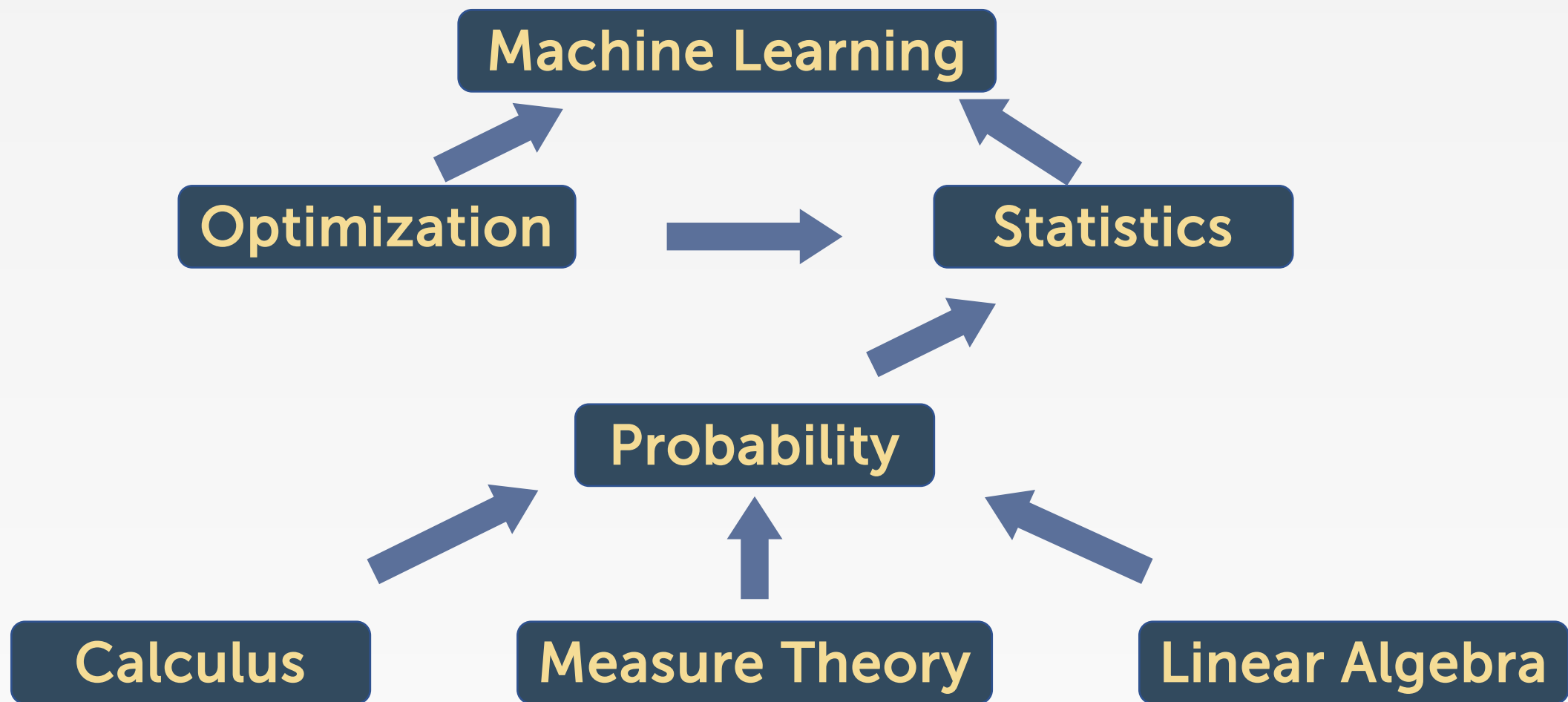


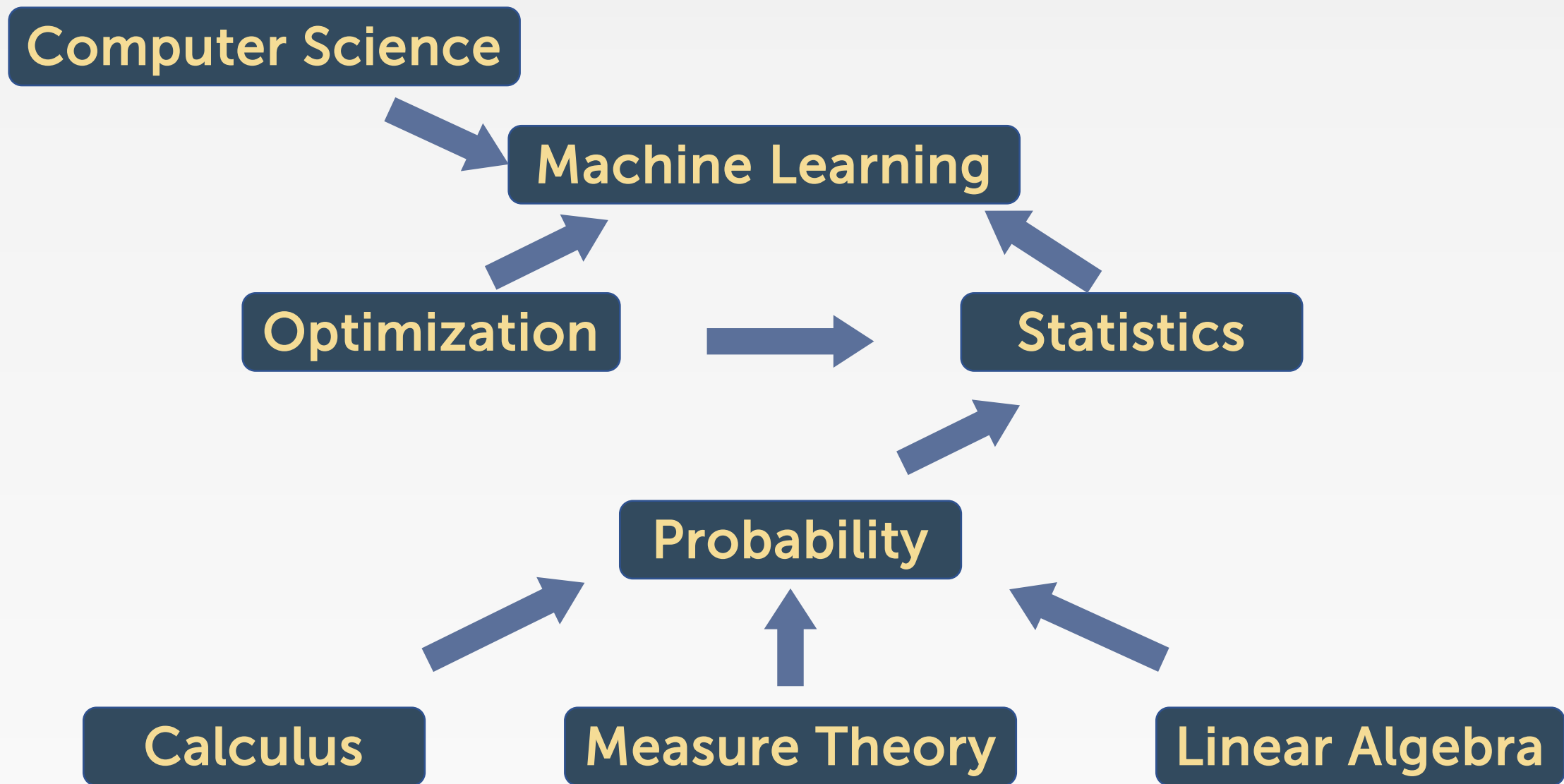
# Machine Learning

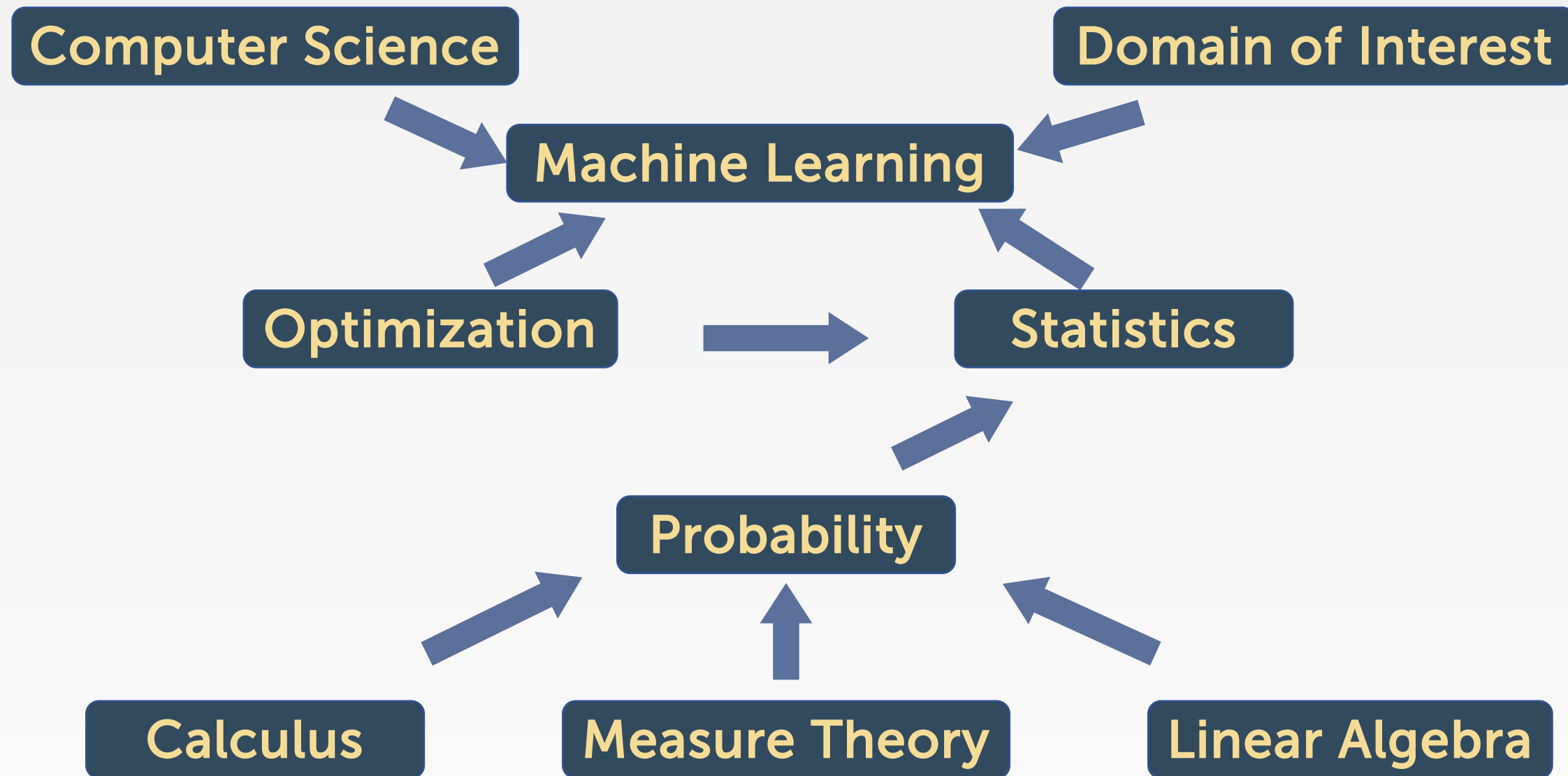






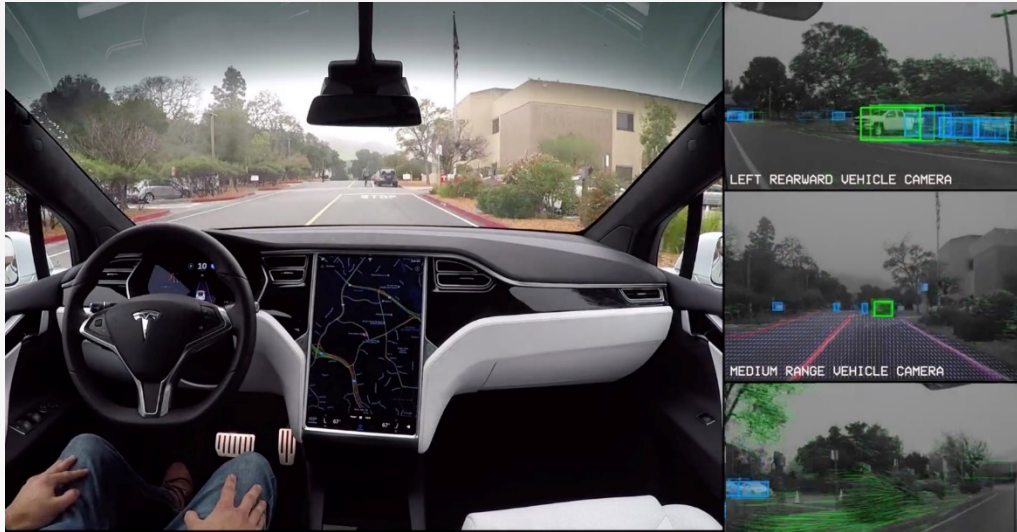






# ML is everywhere

Learning to drive an autonomous vehicle (robotics)



Tesla Self-Driving cars



Uber Self-Driving services

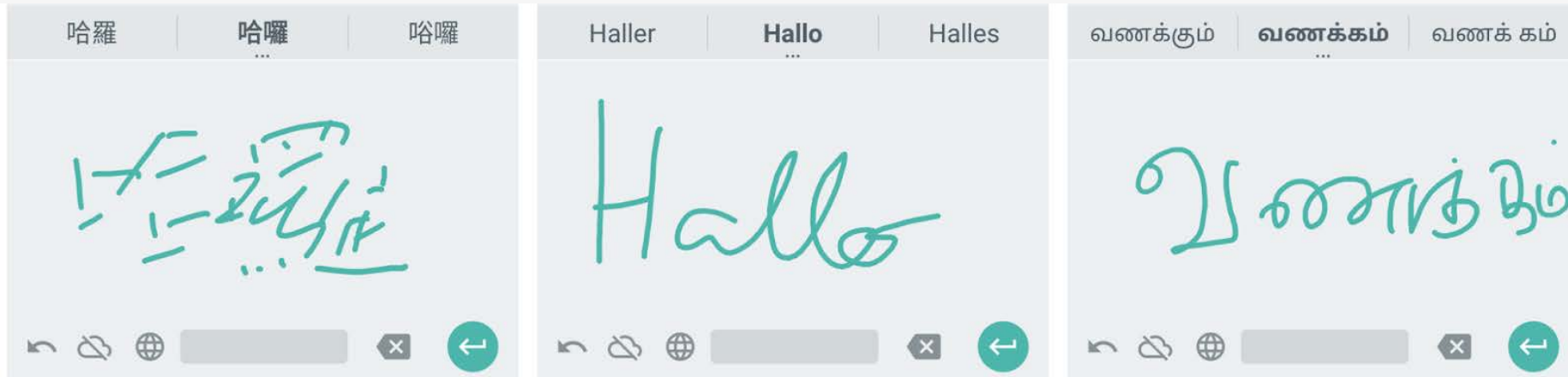
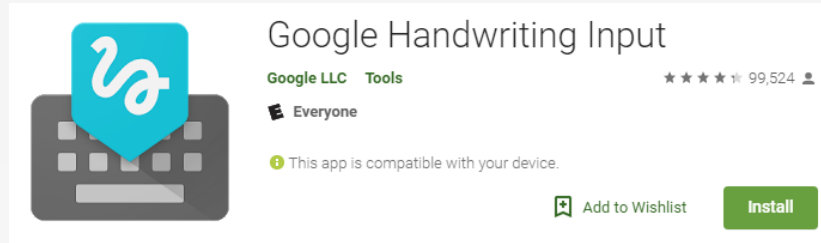
# ML is everywhere

Learning to beat the masters at go games (Games/Reasoning)



# ML is everywhere

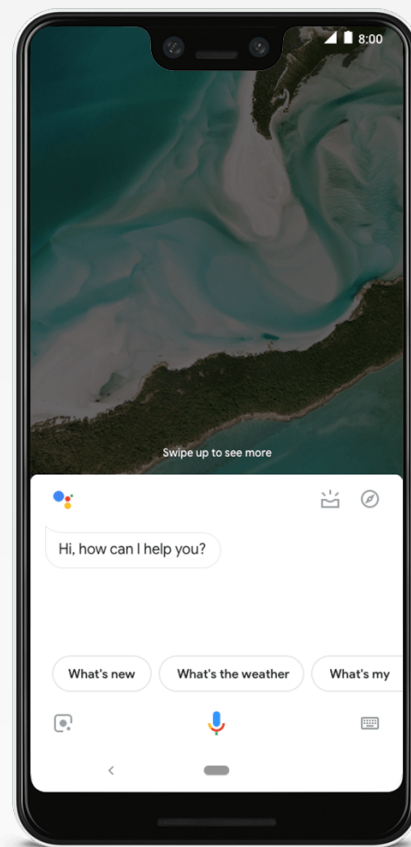
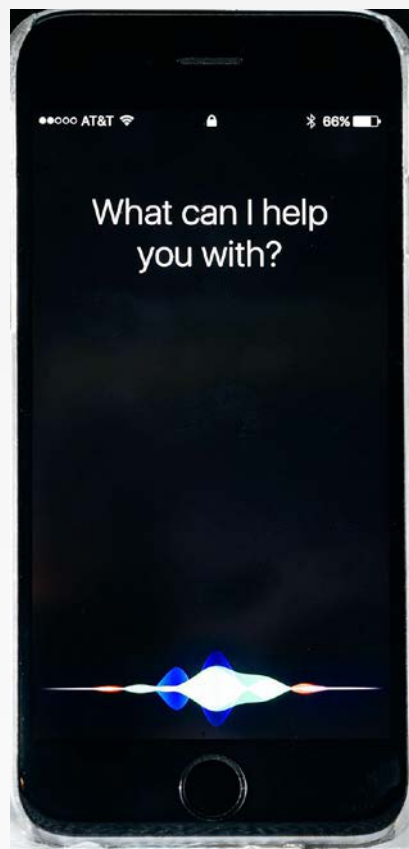
Learning to recognize handwriting (computer vision)





# ML is everywhere

Learning to recognize spoken words (speech recognition)



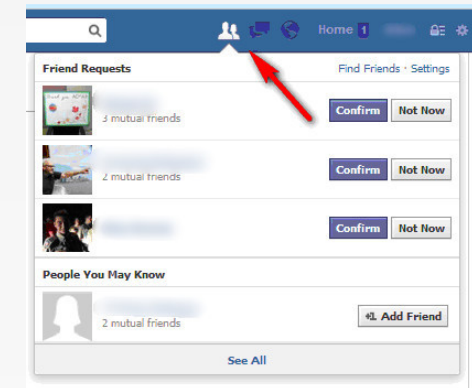
# ML is everywhere

Learning to...

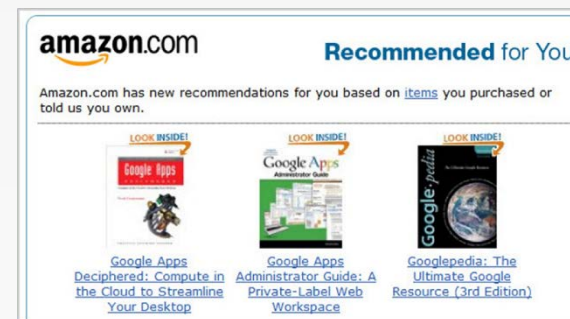
identify spam emails



suggest people you may know



recommend movies



...

# Societal/Business Impacts of ML

- Search results are optimized for ad revenues

**Baidu 百度** 出国留学 百度一下

Q 网页 资讯 视频 图片 知道 文库 贴吧 地图 采购 更多

百度为您找到相关结果约100,000,000个 搜索工具

**出国留学网【专业的留学门户网站】**

出国留学网从2005年创立至今,已经成为领先的中国教育门户网站。出国留学网下设留学、移民、签证、考试、范文、作文等频道,充分满足各类学生对不同类型资讯的需求。同时专注打造权威的工具箱及海内...

出国留学网 百度快照

**出国留学-新航道广州学校**

随着人们生活质量的提高,很多学生有出国留学的想法,而美国一直都是最受中国留学生喜欢的国家,那么如何才能去美国留学呢?以下是去美国留学的两种方法:查看详情 30 21-08 香港留学|想要申请香港...

gz.xhd.cn/cglx/ 百度快照

**第一留学网 - 出国留学中介机构\_留学费用咨询 - 全球留学...**

第一留学网专注出国留学服务,解答出国留学条件及留学途径留学费等知识,同时推荐正规合适的出国留学中介机构,并且为正在留学的用户提供高效便捷的全球留学院校排名信息查询。

www.hnrichfund.com/ 百度快照

**北京留学\_北京留学机构\_专业出国留学中介-金吉列留学官网**

金吉列北京留学服务中心,北京出国留学首选咨询服务机构,为您提供全方位北京留学,北京出国留学中介等留学权威资讯,想了解更多北京留学资讯就到金吉列留学官网。

www.jjl.cn/ 百度快照

**相关组织机构**

金吉列留学 留学中介 公费留学 中国留学服务中心

中国500家最大私企 留学一条龙服务 要求学习结束后回国 教育部直属事业单位

国际学校 中外合作办学 英国大学 多伦多大学

供母语教育的学校 高教未来发展的大趋势 拥有英国皇家特许状 加拿大一所著名的大学

爱丁堡大学 墨尔本大学 曼谷大学 平均学分绩点

英语国家古老大学 世界顶尖的研究型大学 泰国最大私立大学之一 量与质的计算

# Societal/Business Impacts of ML

- Search results are optimized for ad revenues

**Baidu 百度** 出国留学 百度一下

Q 网页 资讯 视频 图片 知道 文库 贴吧 地图 采购 更多

百度为您找到相关结果约100,000,000个 搜索工具

**出国留学网【专业的留学门户网站】**

出国留学网从2005年创立至今,已经成为领先的中国教育门户网站。出国留学网下设留学、移民、签证、考试、范文、作文等频道,充分满足各类学生对不同类型资讯的需求。同时专注打造权威的工具箱及海内...

出国留学网 百度快照

**出国留学-新航道广州学校**

网络预约享特惠 领取200抵现优惠券

随着人们生活质量的提高,很多学生有出国留学的想法,而美国一直都是最受中国留学生喜欢的国家,那么如何才能去美国留学呢?以下是去美国留学的两种方法:查看详情 30 21-08 香港留学|想要申请香港...

gz.xhd.cn/cglx/ 百度快照

**第一留学网 - 出国留学中介机构 留学费用咨询 - 全球留学...**

第一留学网专注出国留学服务,解答出国留学条件及留学途径留学费等问题,同时推荐正规合适的出国留学中介机构,并且为正在留学的用户提供高效便捷的全球留学院校排名信息查询。

www.hnrichfund.com/ 百度快照

**北京留学\_北京留学机构\_专业出国留学中介-金吉列留学官网**

金吉列北京留学服务中心,北京出国留学首选咨询服务机构,为您提供全方位北京留学,北京出国留学中介等留学权威资讯,想了解更多北京留学资讯就到金吉列留学官网。

www.jjl.cn/ 百度快照

出国留学 留学机构和课程/72条 留学经验

## 魏则西事件五周年



青年魏则西去世五周年

时间: 2016年4月12日

2016年4月12日, 21岁的魏则西因滑膜肉瘤去世, 在其生前求医过程中, 通过百度搜索到武警北京总队第二医院, 被该医院宣传的“生物免疫疗法”、“斯坦福技术”所骗, 花费不赀却未收获任何效果, 贻误合理治疗时机。魏则西去世后, 莆田系医院虚假宣传、百度搜索竞价排名、部队医院对外承包混乱等问题引发社会强烈关注。

# Societal/Business Impacts of ML

- Search results are optimized for ad revenues
- An autonomous vehicle is permitted to drive unassisted on the road



# Societal/Business Impacts of ML

- Search results are optimized for ad revenues
- An autonomous vehicle is permitted to drive unassisted on the road



# Societal/Business Impacts of ML

- Search results are optimized for ad revenues
- An autonomous vehicle is permitted to drive unassisted on the road



## 31岁企业家驾驶蔚来ES8车祸身亡，行驶数据浮出水面

2021年08月15日 19:22:36

来源：北京青年报

1595人参与 420评论



日前，一则蔚来ES8“自动驾驶”发生车祸死亡的事件引起网友关注。8月15日，疑似涉事车辆行驶数据曝光。

### 31岁创始人驾驶ES8车祸身亡

8月14日，认证名为“美一好”的个人公众号发布讣告称，2021年8月12日下午2时，上善若水投资管理公司创始人、意统天下餐饮管理公司创始人、美一好品牌管理公司创始人林文钦（昵称“萌剑客”），驾驶蔚来ES8汽车启用自动驾驶功能（NOP领航状态）后，在沈海高速涵江段发生交通事故，不幸逝世，终年31岁。



美一好 >



### 讣告 | 我们的“萌剑客”走了

2021年8月12日下午2时，上善若水投资管理公司创始人、意统天下餐饮管理公司创始人、美一好品牌管理公司创始人林文钦先生（昵称“萌剑客”），驾驶蔚来ES8汽车启用自动驾驶功能（NOP领航状态）后，在沈海高速涵江段发生交通事故，不幸逝世，终年31岁。

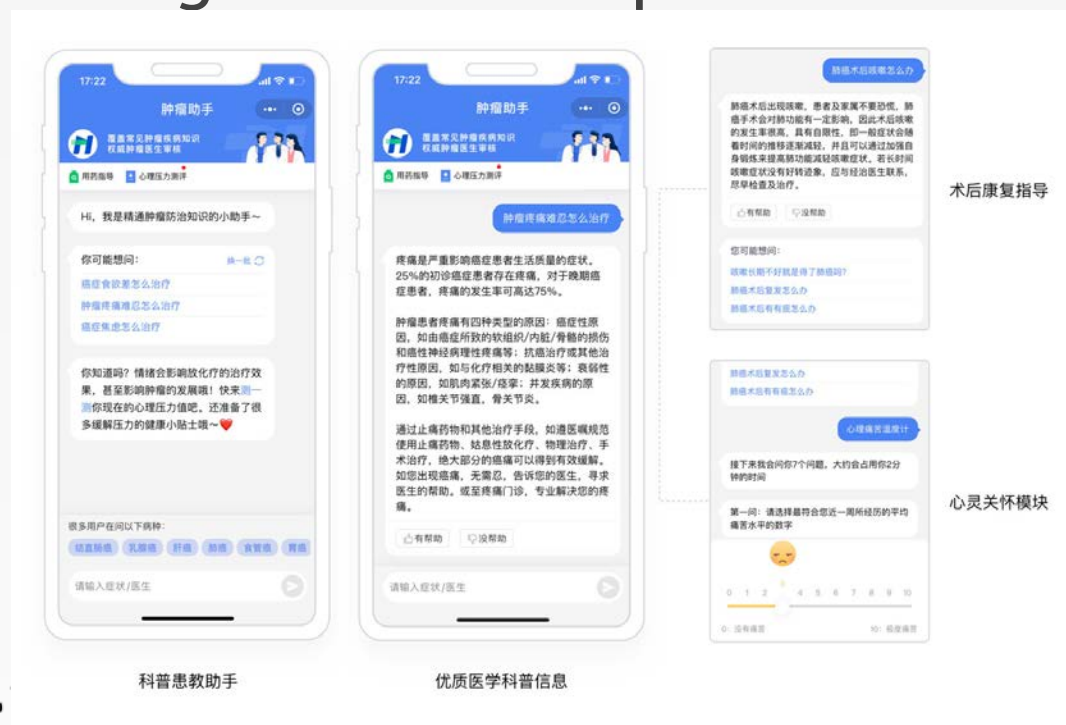
# Societal/Business Impacts of ML

- Search results are optimized for ad revenues
- An autonomous vehicle is permitted to drive unassisted on the road
- A doctor is prompted by an intelligent system with a plausible diagnosis for her patient



# Societal/Business Impacts of ML

- Search results are optimized for ad revenues
- An autonomous vehicle is permitted to drive unassisted on the road
- A doctor is prompted by an intelligent system with a plausible diagnosis for her patient

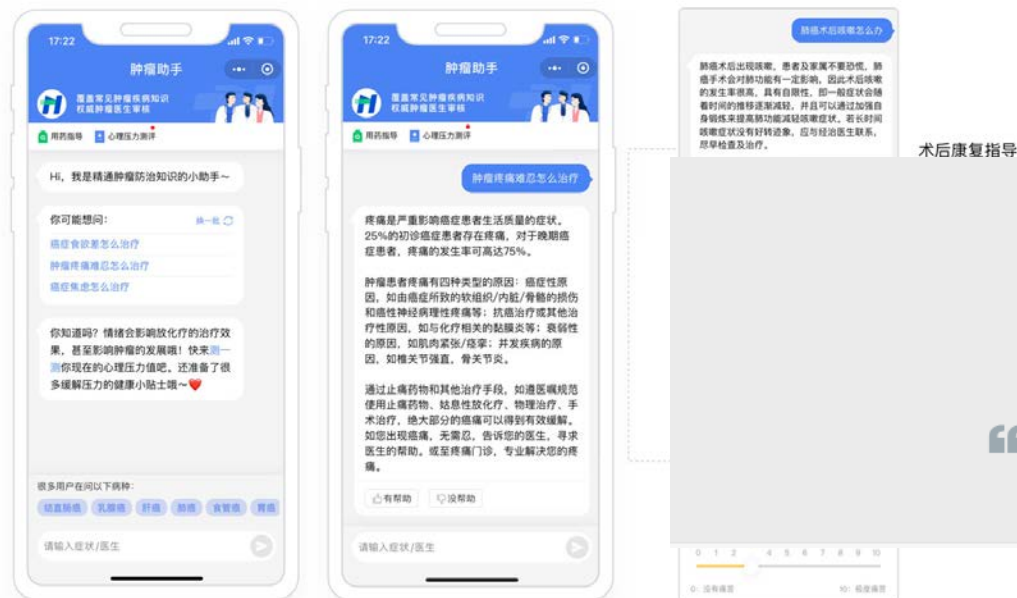


光

Guanghua School of Management

# Societal/Business Impacts of ML

- Search results are optimized for ad revenues
- An autonomous vehicle is permitted to drive unassisted on the road
- A doctor is prompted by an intelligent system with a plausible diagnosis for her patient



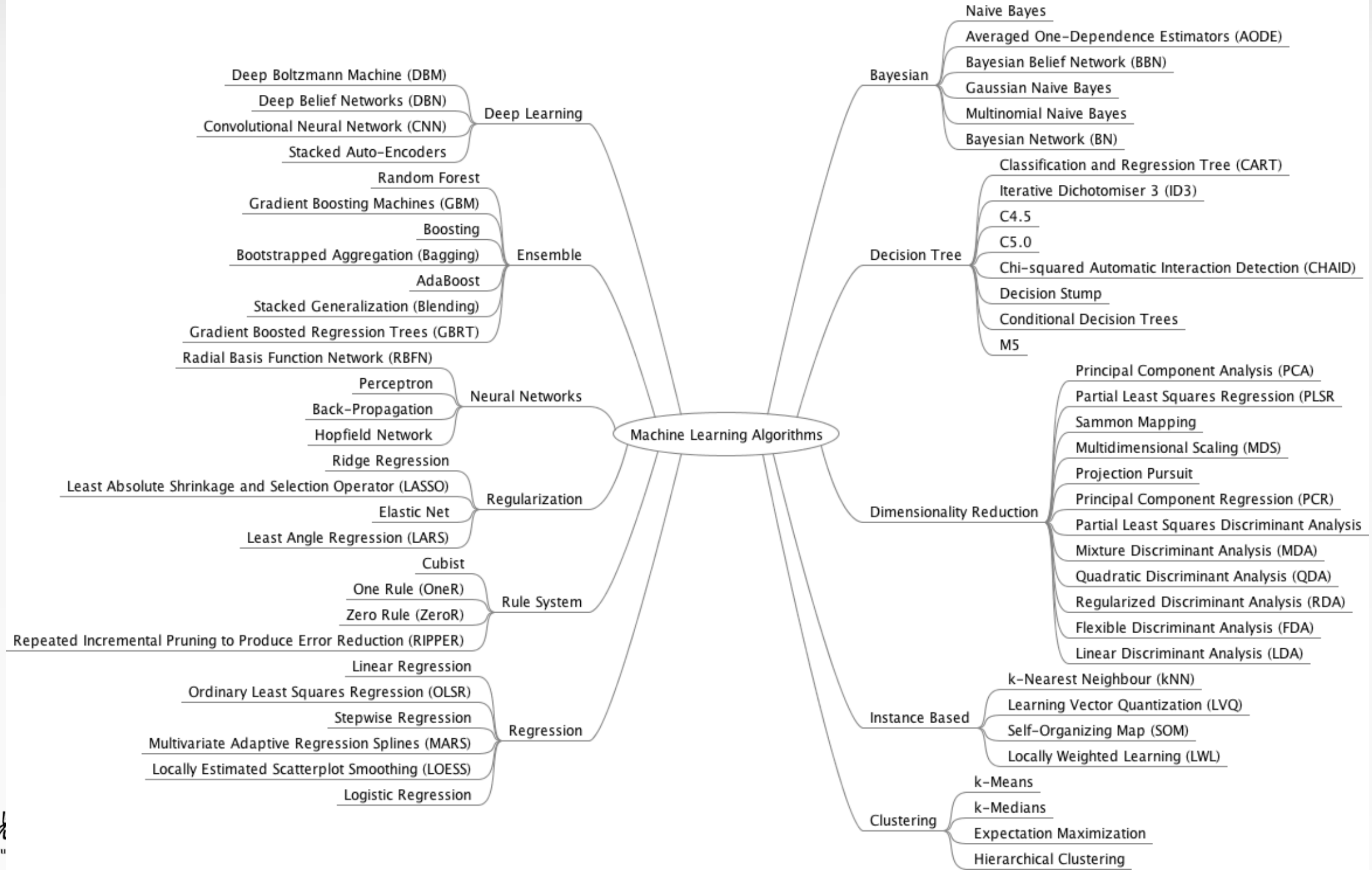
## 英国权威医学期刊 diss 医疗AI：在乳腺癌检测上取代放射科医生 是痴人说梦

本文作者：我在思考中

2021-09-06 11:04

导语：近日，《英国医学杂志》刊登了一篇研究工作。该团队工作对近年 AI 技术用于乳腺癌筛查的工作进行了检索，希望检验 AI 技术用于乳房 X 光图像识别的准确度。





# DEFINING a ML Problem

# ML brings together different areas

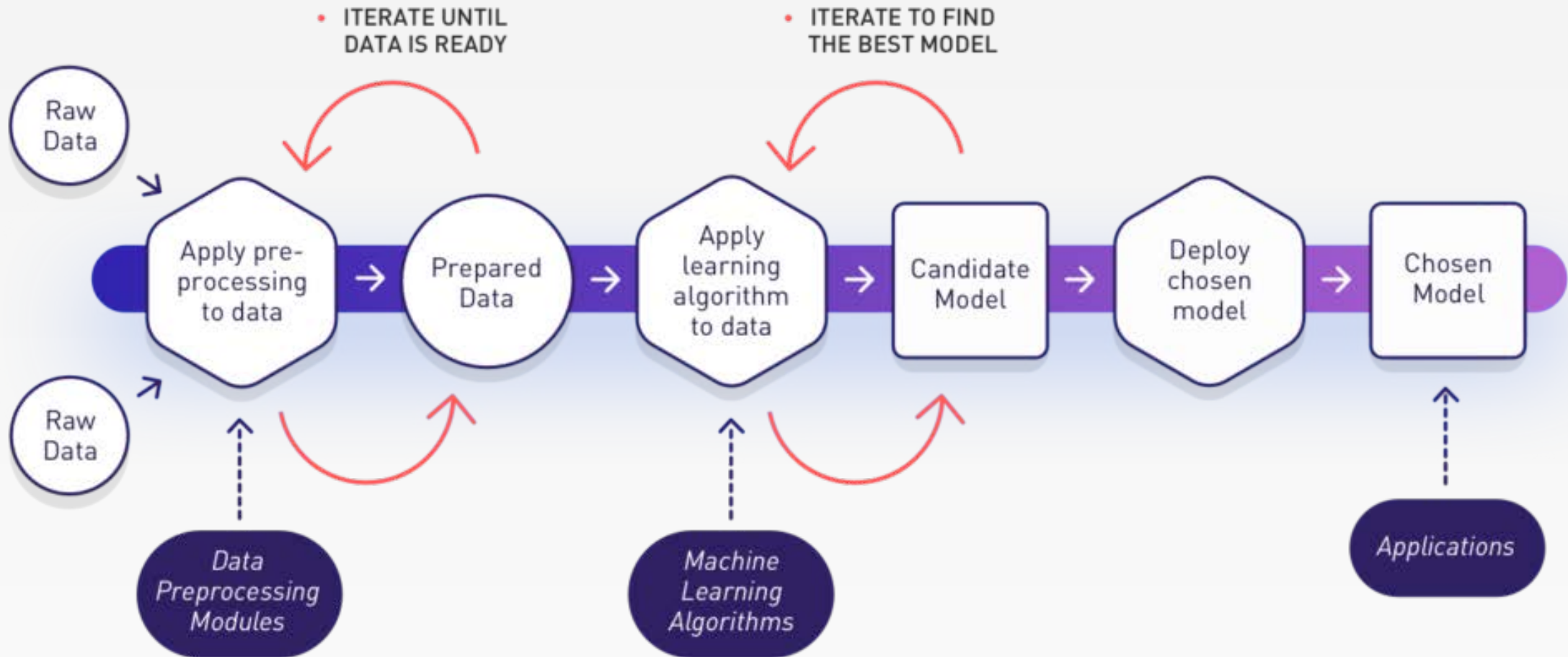
# ML brings together different areas

- Statistical methods
  - Infer conclusions from data
  - Estimate reliability of predictions
- Computer science
  - Large-scale computing architectures
  - Algorithms for capturing, manipulating, indexing, combining, retrieving and performing predictions on data
  - Software pipelines that manage the complexity of multiple subtasks

# ML brings together different areas

- Statistical methods
  - Infer conclusions from data
  - Estimate reliability of predictions
- Computer science
  - Large-scale computing architectures
  - Algorithms for capturing, manipulating, indexing, combining, retrieving and performing predictions on data
  - Software pipelines that manage the complexity of multiple subtasks
- Economics, biology, psychology
  - How can an individual or system efficiently improve their performance in a given environment?
  - What is learning and how can it be optimized?

# Machine Learning Workflow





# Define a Learning Problem

- **Definition:** A computer program learns if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

# Define a Learning Problem

- **Definition:** A computer program learns if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .
- **Three components**
  - Task,  $T$
  - Performance measure,  $P$
  - Training,  $E$

# Define a Learning Problem

- Three components
  - Task  $T$
  - Performance measure  $P$
  - Training  $E$

Example 1 : Handwriting recognition

- $T$ :
- $P$ :
- $E$ :

# Define a Learning Problem

- Three components
  - Task  $T$
  - Performance measure  $P$
  - Training  $E$

## Example 1 : Handwriting recognition

- $T$ : recognizing and classifying handwritten words with images
- $P$ : percent of words correctly classified
- $E$ : a database of handwritten words with given classifications

# Define a Learning Problem

- Three components
  - Task  $T$
  - Performance measure  $P$
  - Training  $E$

## Example II : Self-driving

- $T$ :
- $P$ :
- $E$ :

# Define a Learning Problem

- Three components

- Task  $T$
- Performance measure  $P$
- Training  $E$

## Example II : Self-driving

- $T$ : driving on public four-lane highways using vision sensors
- $P$ : average distance traveled before an error
- $E$ : a sequence of images and steering commands recorded while observing a human driver

# Define a Learning Problem

- Three components
  - Task  $T$
  - Performance measure  $P$
  - Training  $E$

**Exercise:** Siri response to voice commands

- $T$ : ?
- $P$ : ?
- $E$ : ?

# Solution #1 Expert Systems

- Over 20 years ago, we had *rule-based systems*:
  1. Put a bunch of linguists in a room
  2. Have them think about the structure of their native language and write down the rules they devise



# Solution #1 Expert Systems

- Over 20 years ago, we had *rule-based systems*:
  1. Put a bunch of linguists in a room
  2. Have them think about the structure of their native language and write down the rules they devise

Give me directions to Starbucks

If: "give me directions to X"

Then: `directions(here, nearest(X))`

# Solution #1 Expert Systems

- Over 20 years ago, we had *rule-based systems*:
  1. Put a bunch of linguists in a room
  2. Have them think about the structure of their native language and write down the rules they devise

Give me directions to Starbucks

If: "give me directions to X"

Then: directions(here, nearest(X))

How do I get to Starbucks?

If: "how do I get to X"

Then: directions(here, nearest(X))

# Solution #1 Expert Systems

- Over 20 years ago, we had *rule-based systems*:
  1. Put a bunch of linguists in a room
  2. Have them think about the structure of their native language and write down the rules they devise

Give me directions to Starbucks

If: "give me directions to X"

Then: directions(here, nearest(X))

How do I get to Starbucks?

If: "how do I get to X"

Then: directions(here, nearest(X))

Where is the nearest Starbucks?

If: "where is the nearest X"

Then: directions(here, nearest(X))

# Solution #2 Annotate Data and Learn

- Experts:
  - Very good at answering questions about specific cases
  - Not very good at telling HOW they do it
- 1990s: So why not just have them tell you what they do on SPECIFIC CASES and then let MACHINE LEARNING tell you how to come to the same decisions that they did

# Solution #2 Annotate Data and Learn

- Collect raw sentences  $\{x^{(1)}, \dots, x^{(n)}\}$
- Experts annotate their meaning  $\{y^{(1)}, \dots, y^{(n)}\}$

# Solution #2 Annotate Data and Learn

- Collect raw sentences  $\{x^{(1)}, \dots, x^{(n)}\}$
- Experts annotate their meaning  $\{y^{(1)}, \dots, y^{(n)}\}$

$x^{(1)}$  : Give me directions to Starbucks

$y^{(1)}$  : `directions(here,  
nearest(Starbucks))`

$x^{(2)}$  : Send a text to John that I'll be late

$y^{(2)}$  : `txtmsg(John, I'll be late)`

$x^{(3)}$  : Show me the closest Starbucks

$y^{(3)}$  : `map(nearest(Starbucks))`

$x^{(4)}$  : Set an alarm for seven in the morning

$y^{(4)}$  : `setalarm(7:00AM)`

# Define a Learning Problem

- Three components
  - Task  $T$
  - Performance measure  $P$
  - Training  $E$

**Exercise:** Siri response to voice commands

- $T$ :
- $P$ :
- $E$ :

# Define a Learning Problem

- Three components

- Task  $T$
- Performance measure  $P$
- Training  $E$

## Exercise: Siri response to voice commands

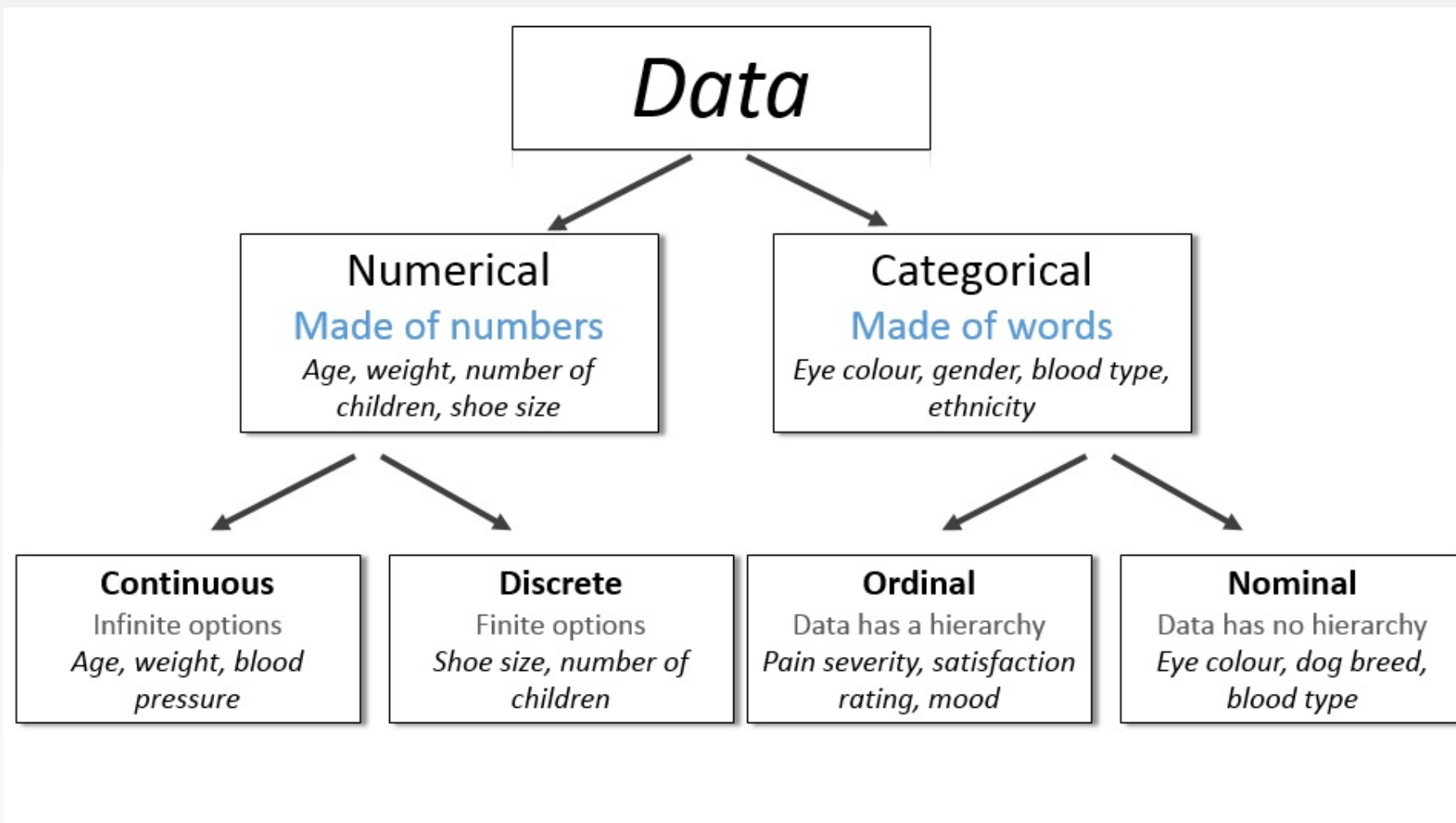
- $T$ : predicting action from speech
- $P$ : percent of correct actions taken in user pilot study
- $E$ : examples of (speech, action) pairs



# Problem Formulation

- Formulate a problem in more than one ways:
- Loan applications:
  - Credit score (regression)
  - Default probability (density estimation)
  - Loan decision (classification)

# Data Types



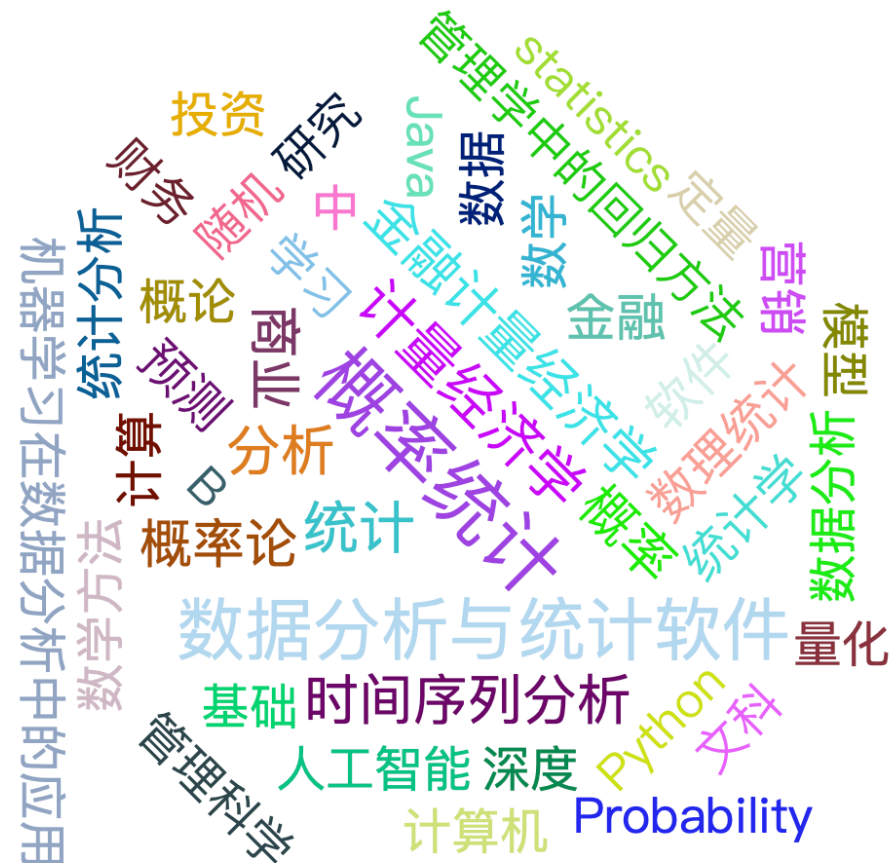
# SYLLABUS



# 专业

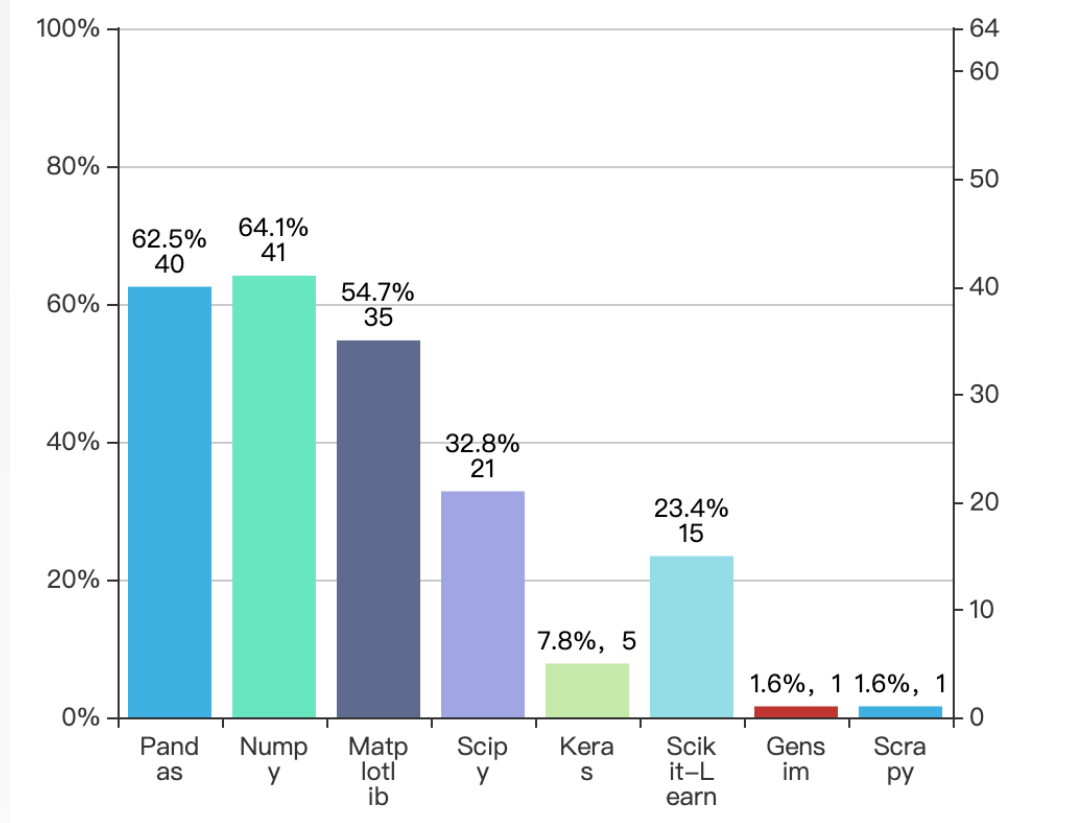
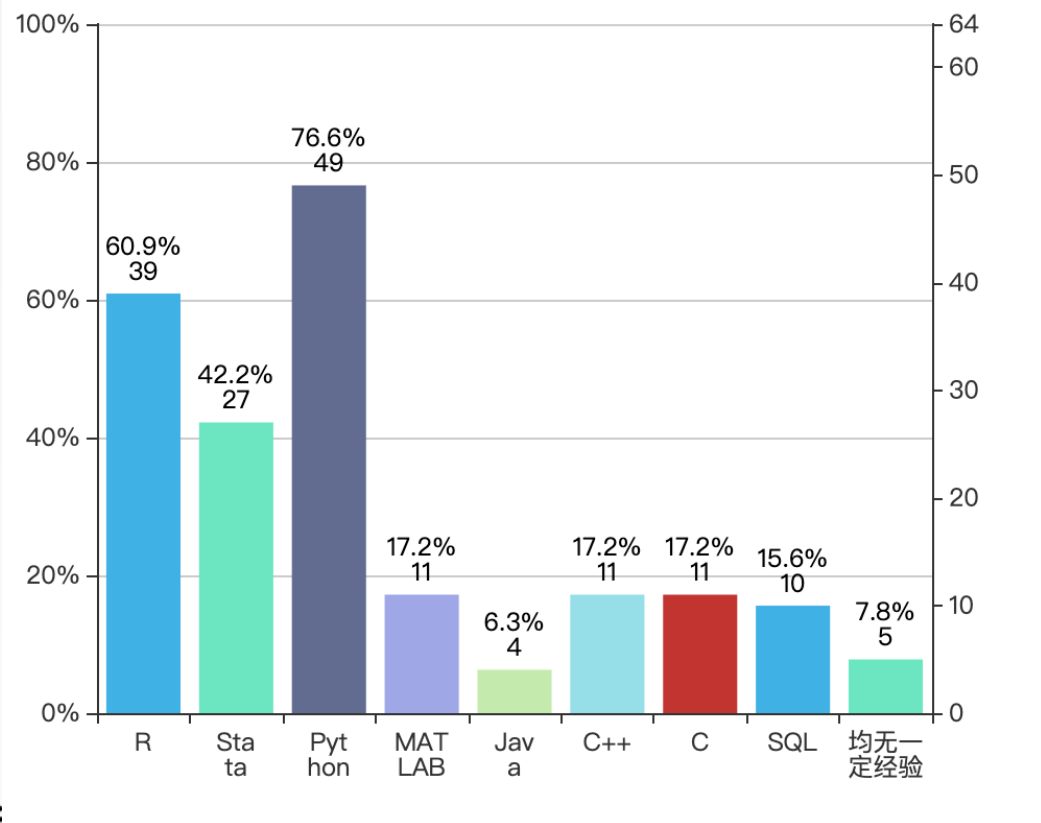


# Minor



## 先修课程

# 编程背景



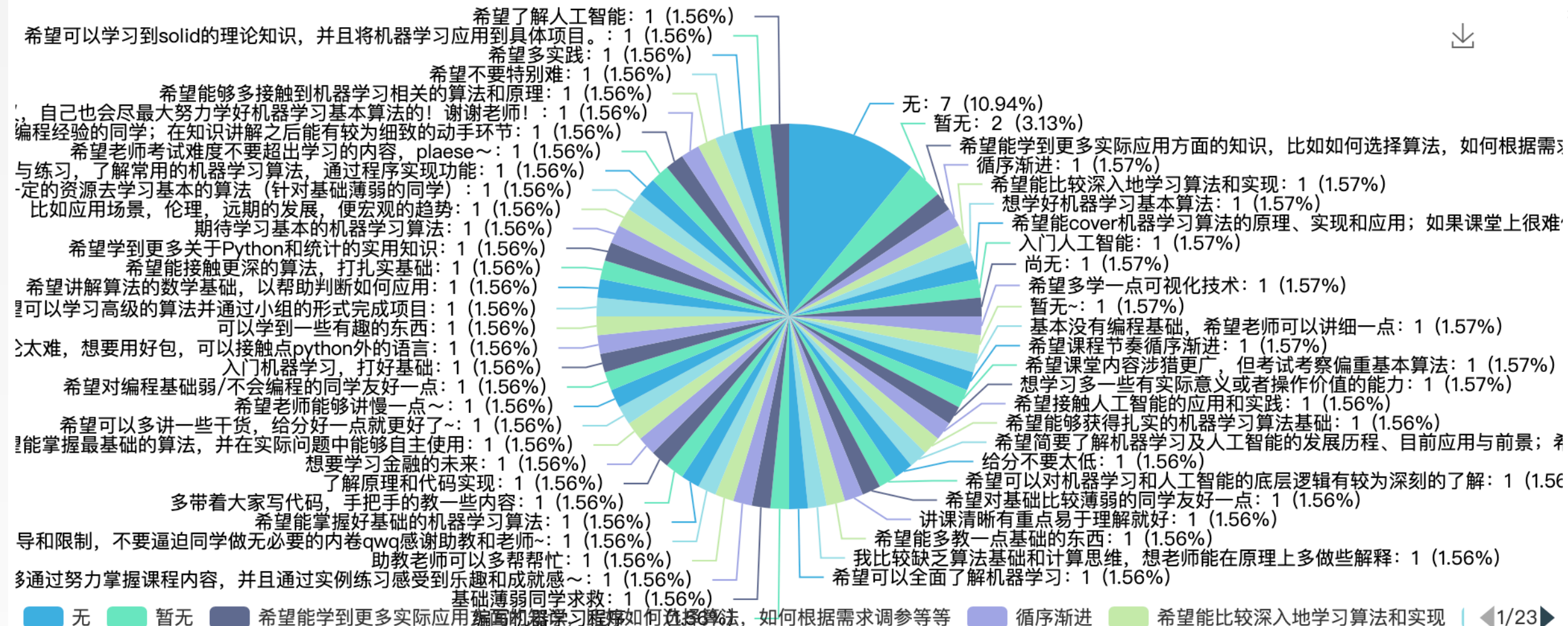
# 课程期待

42%

心有余而力不足，想学好机器学习基本算法

58%

期待，可以学习更加fancy的算法





# Logistics

# Logistics

- Homework:
  - 3 individual assignments
  - Late assignment is ***NOT*** accepted
    - Supervised Learning
    - Unsupervised Learning + Model Evaluation
    - Advanced Topics (e.g., Deep Learning, Reinforcement Learning)
  - **Re-grade**: You can appeal your grade with a one-page explanation. A re-grade may cause your grade to either go up or go down.



# Logistics

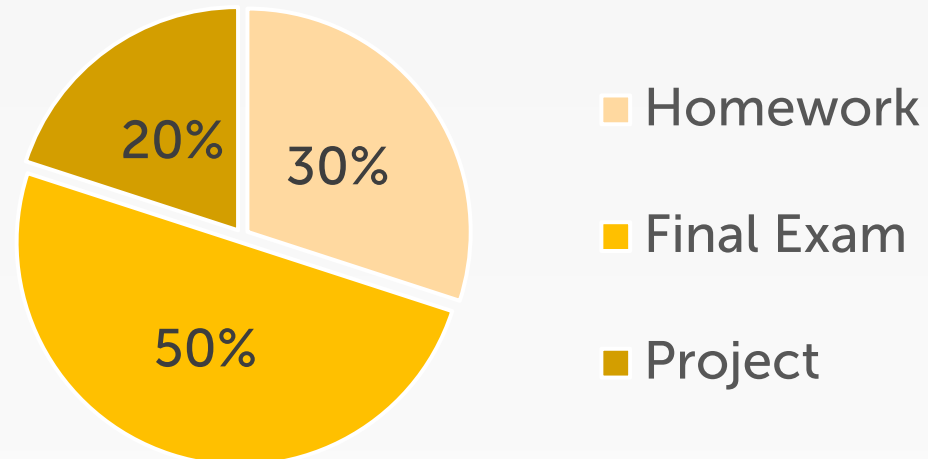
- Homework:
  - 3 individual assignments
  - Late assignment is **NOT** accepted
    - Supervised Learning
    - Unsupervised Learning + Model Evaluation
    - Advanced Topics (e.g., Deep Learning, Reinforcement Learning)
  - **Re-grade**: You can appeal your grade with a one-page explanation. A re-grade may cause your grade to either go up or go down.
- Group Project
  - At most 5 students per group (No Free-riding)
  - Proposal
  - In-class presentation
  - Final reports

# Logistics

- Homework:
  - 3 individual assignments
  - Late assignment is **NOT** accepted
    - Supervised Learning
    - Unsupervised Learning + Model Evaluation
    - Advanced Topics (e.g., Deep Learning, Reinforcement Learning)
  - **Re-grade**: You can appeal your grade with a one-page explanation. A re-grade may cause your grade to either go up or go down.
- Group Project
  - At most 5 students per group (No Free-riding)
  - Proposal
  - In-class presentation
  - Final reports
- Final Exam
  - Week 12 (in-class)

# Logistics

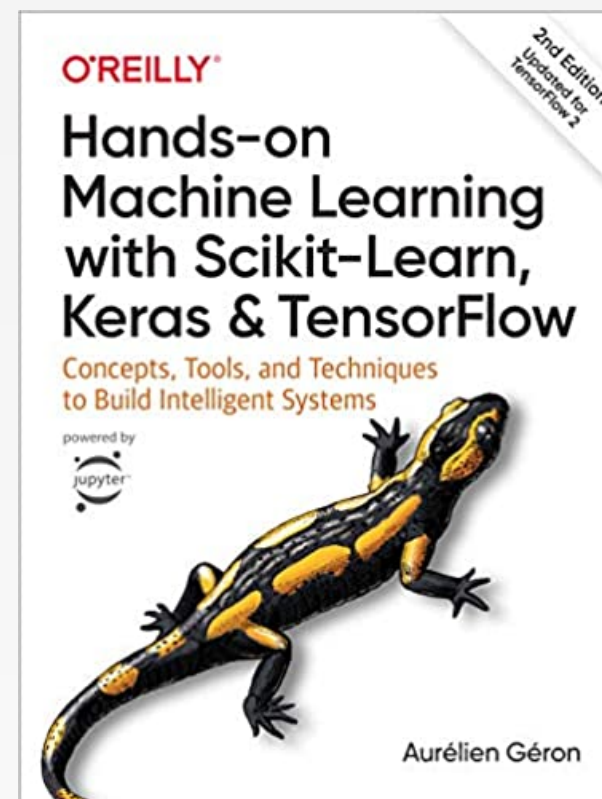
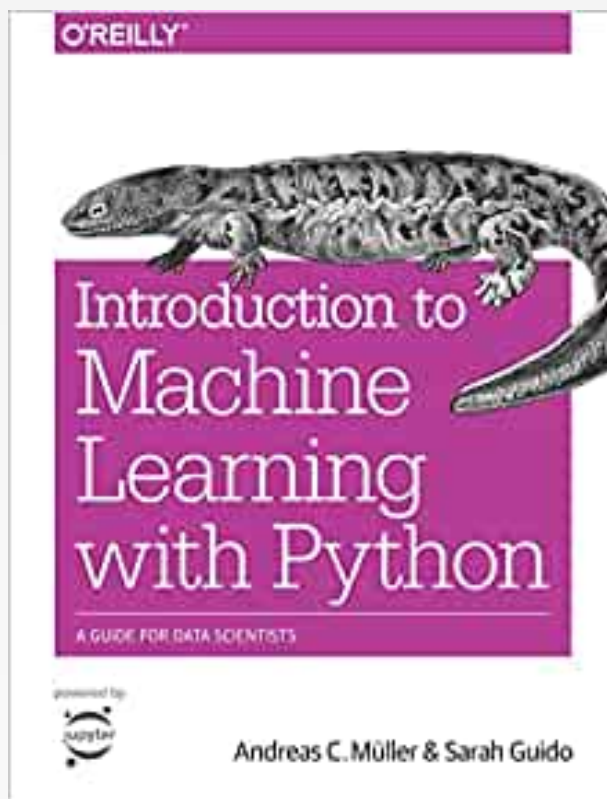
- Academic Honor Code: **No plagiarism!**
  - form study groups (with arbitrary number of people); discuss and work on homework problems in groups
  - write down the solutions independently
  - write down the names of people with whom you've discussed the homework
- Class participation



# Topics

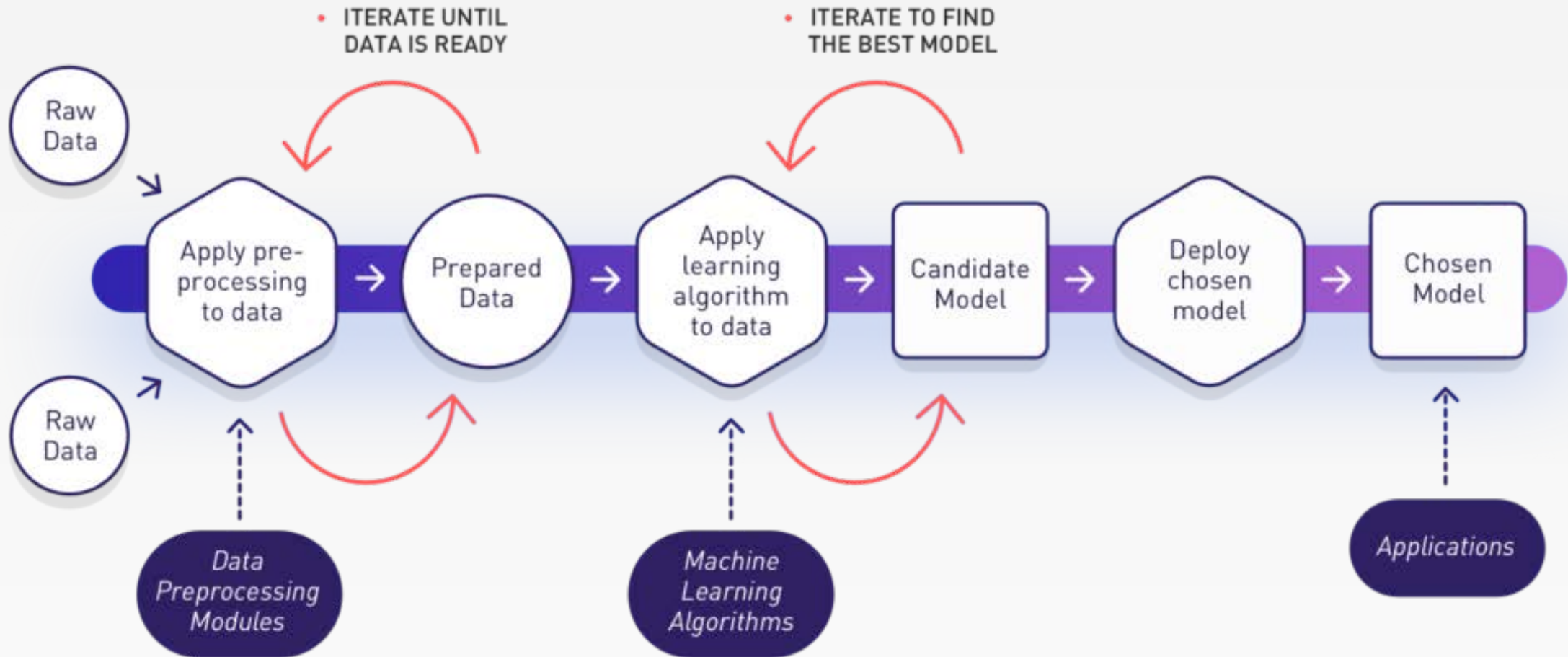
- Regression
- K-Nearest Neighbors
- Decision Trees
- Naïve Bayes
- SVM
- Ensemble Models
- Cross Validation
- Overfitting / Underfitting
- Model Evaluation and Selection
- Unsupervised Learning
  - Clustering
  - PCA
- Reinforcement Learning
- Deep learning
  - Neural network
  - Backpropagation
  - CNNs, LSTM, etc.

# Recommended Books



# Overview of ML Models

# Machine Learning Workflow



# Applied ML

- Understand basic ML concepts and workflow (require basic statistics/probability background)
- Apply properly “black-box” ML components and features
- From theory to real-world practice



# Types of ML Systems

# Types of ML Systems

## Criteria

Whether or not they are trained with human supervision

# Types of ML Systems

## Criteria

Whether or not they are trained with human supervision

### Supervised Learning

Fraud detection  
Prediction of stock markets

# Types of ML Systems

## Criteria

Whether or not they are trained with human supervision

### Supervised Learning

Fraud detection  
Prediction of stock markets

### Unsupervised Learning

Customer segmentation  
Recommendation

# Types of ML Systems

## Criteria

Whether or not they are trained with human supervision

### Supervised Learning

Fraud detection  
Prediction of stock markets

### Unsupervised Learning

Customer segmentation  
Recommendation

### Semi-supervised Learning

Photo-hosting service  
Speech analysis  
Web-content classification

# Types of ML Systems

## Criteria

Whether or not they are trained with human supervision

### Supervised Learning

Fraud detection  
Prediction of stock markets

### Unsupervised Learning

Customer segmentation  
Recommendation

### Semi-supervised Learning

Photo-hosting service  
Speech analysis  
Web-content classification

### Reinforcement Learning

Robotics  
Go games  
Self-driving cars

# Types of ML Systems

## Criteria

Whether or not they are trained with human supervision



### Supervised Learning

Fraud detection  
Prediction of stock markets



### Unsupervised Learning

Customer segmentation  
Recommendation

### Semi-supervised Learning

Photo-hosting service  
Speech analysis  
Web-content classification

### Reinforcement Learning

Robotics  
Go games  
Self-driving cars

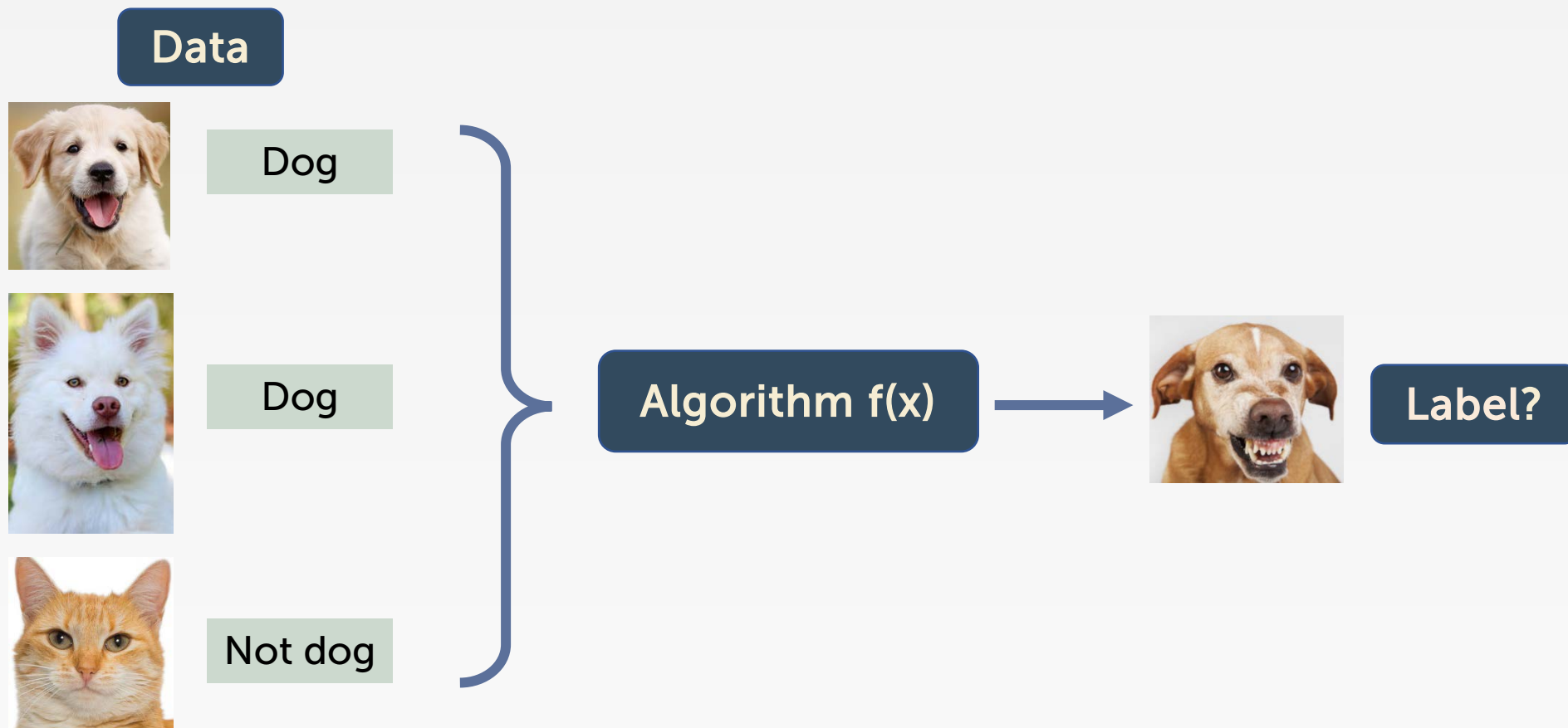
# Supervised Learning



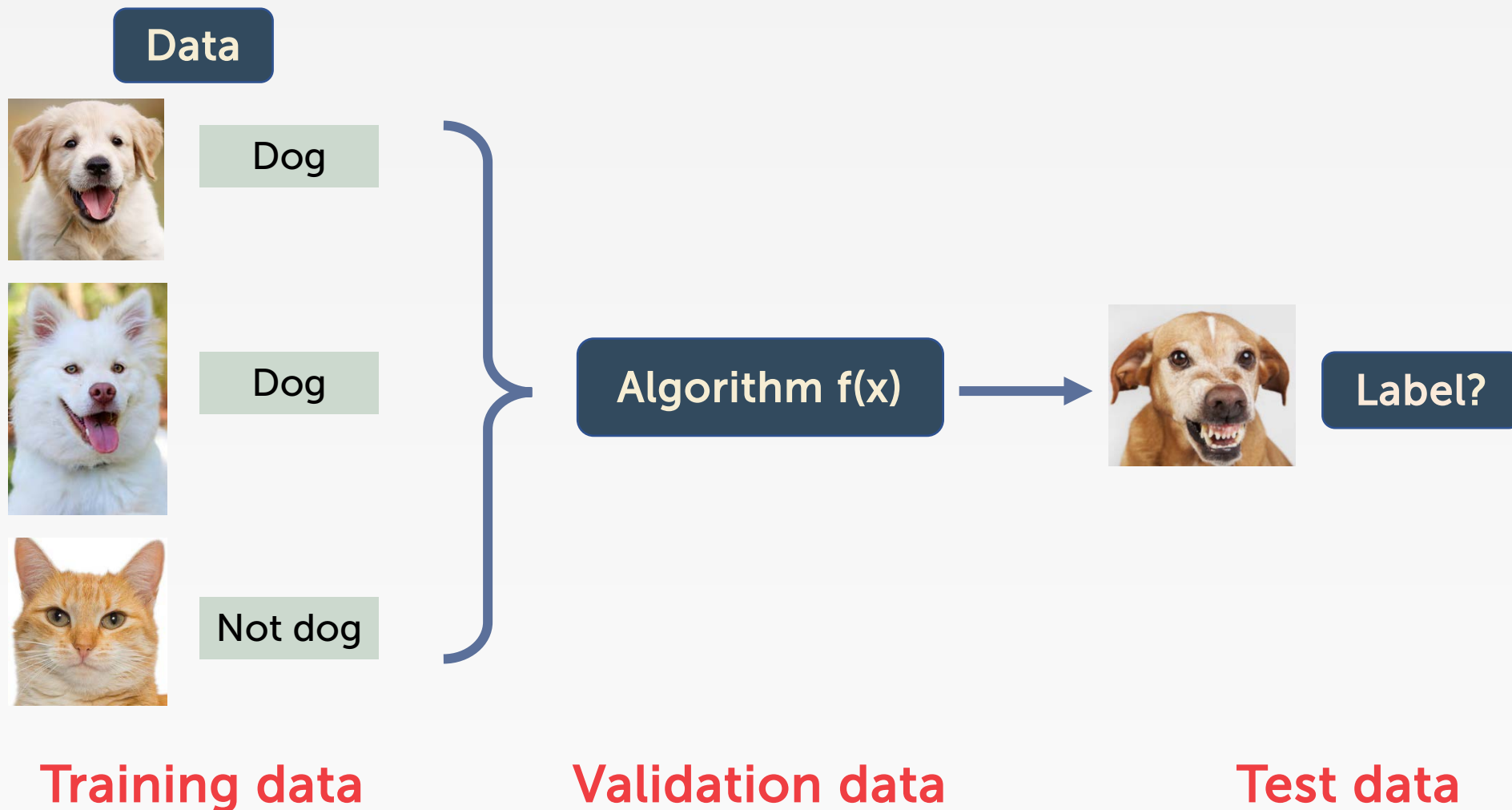
# Supervised Learning

- Regressions
- KNN
- Decision Trees
- SVM
- Naïve Bayes
- ...

# Supervised Learning



# Supervised Learning

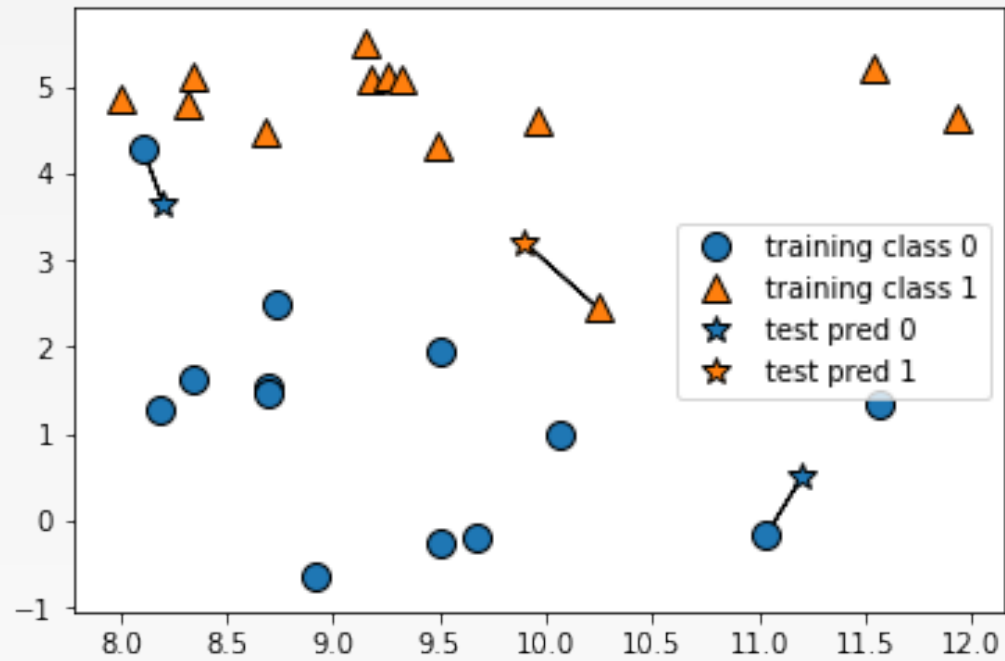


# Data Division

- **Training dataset**: the sample of data used to fit the model
- **Validation Dataset**: the sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyperparameters.
- **Test Dataset**: a set of examples used only to assess the performance (i.e. generalization) of a fully specified classifier

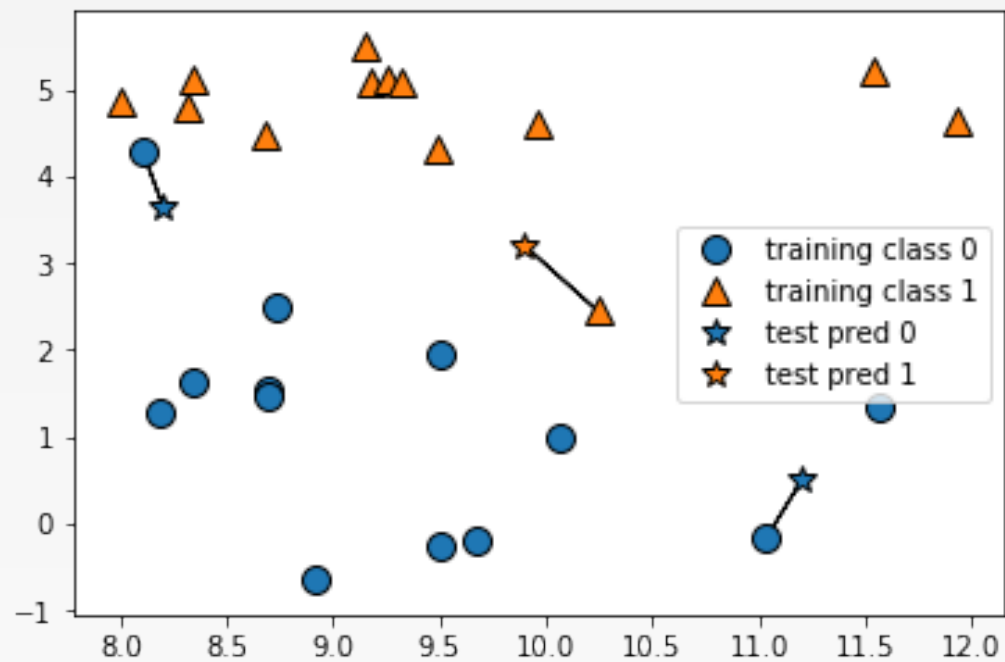
# K-Nearest Neighbors

# An Illustrative Example

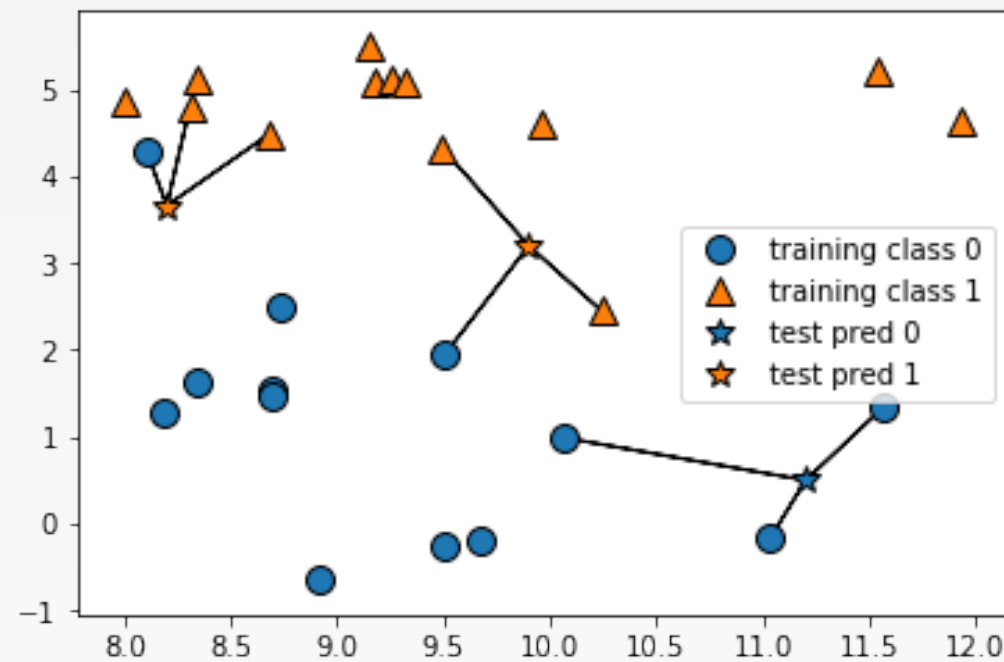


$K = 1$

# An Illustrative Example



$K = 1$



$K = 3$

# KNN: Classification Algorithm



# KNN: Classification Algorithm

- Training: Store all the examples  $(X_{train}, Y_{train})$

# KNN: Classification Algorithm

- Training: Store all the examples  $(X_{train}, Y_{train})$
- Prediction:  $X_{new}$ 
  - Let  $X_1, \dots, X_k$  be the  $k$  most similar examples to  $X_{new}$
  - Use certain method (e.g., majority vote) to determine  $Y_{new}$  based on  $(Y_1, \dots, Y_k)$

# KNN: Classification Algorithm

- Training: Store all the examples  $(X_{train}, Y_{train})$
- Prediction:  $X_{new}$ 
  - Let  $X_1, \dots, X_k$  be the  $k$  most similar examples to  $X_{new}$
  - Use certain method (e.g., majority vote) to determine  $Y_{new}$  based on  $(Y_1, \dots, Y_k)$

## Keys

1. A distance metric

Euclidean distance  $d(x_j, x_k) = \sqrt{\sum_i (x_{j,i} - x_{k,i})^2}$   
Manhattan distance  $d(x_j, x_k) = \sum_i |x_{j,i} - x_{k,i}|$

# KNN: Classification Algorithm

- Training: Store all the examples  $(X_{train}, Y_{train})$
- Prediction:  $X_{new}$ 
  - Let  $X_1, \dots, X_k$  be the  $k$  most similar examples to  $X_{new}$
  - Use certain method (e.g., majority vote) to determine  $Y_{new}$  based on  $(Y_1, \dots, Y_k)$

## Keys

1. A distance metric

Euclidean distance  $d(x_j, x_k) = \sqrt{\sum_i (x_{j,i} - x_{k,i})^2}$   
Manhattan distance  $d(x_j, x_k) = \sum_i |x_{j,i} - x_{k,i}|$

2. Value of "K"

Cross validation: larger k? smaller k?

# KNN: Classification Algorithm

- Training: Store all the examples  $(X_{train}, Y_{train})$
- Prediction:  $X_{new}$ 
  - Let  $X_1, \dots, X_k$  be the  $k$  most similar examples to  $X_{new}$
  - Use certain method (e.g., majority vote) to determine  $Y_{new}$  based on  $(Y_1, \dots, Y_k)$

## Keys

1. A distance metric

Euclidean distance  $d(x_j, x_k) = \sqrt{\sum_i (x_{j,i} - x_{k,i})^2}$   
Manhattan distance  $d(x_j, x_k) = \sum_i |x_{j,i} - x_{k,i}|$

2. Value of "K"

Cross validation: larger k? smaller k?

3. Aggregation of the classes of neighbor points

Majority vote

# KNN: Pros and Cons

# KNN: Pros and Cons

- Advantages:
  - Very simple and intuitive
  - The cost of the learning process is zero
  - No assumption about the characteristics/distributions
  - Works on both classification and regression tasks

# KNN: Pros and Cons

- Advantages:

- Very simple and intuitive
- The cost of the learning process is zero
- No assumption about the characteristics/distributions
- Works on both classification and regression tasks

- Drawbacks:

- Computationally expensive when the dataset is very large
  - Need to calculate the compare distance from new example to all other examples
- Sensitive to outliers



# Python...

# Python Quick Checks

- I can read python codes...
- I can write python functions...
- Errors and debugging...

# Coding Tips

- Comments?
- Printed messages?
- Functions?

# iPython

- Command: jupyter notebook
  - Install: Anaconda
- 
- Programming in the browser
  - Codes, instructions, and outputs are displayed “in-line”
  - Useful for writing codes that tells a story
  - Used by scientists and researchers
  - ...

# Python Packages

- Scikit-learn: Python Machine Learning Library

```
from sklearn.tree import DecisionTreeClassifier
```

- Numpy: Scientific Computing Library

- Typically, data input to scikit-learn will be in the form of a Numpy array

```
import numpy as np
```

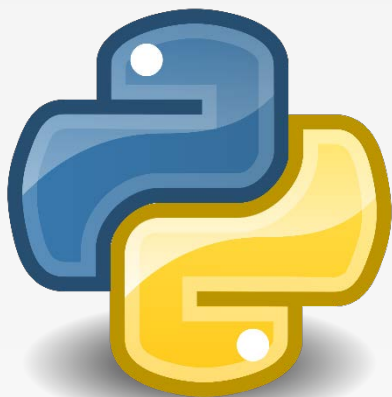
- Pandas: Data Manipulation

```
import pandas as pd
```

- Matplotlib: Plotting Library

```
import matplotlib.pyplot as plt
```

- Others: mglearn; graphviz; seaborn



# Python Practice

# Questions?

# For Next Week...

- Python Review
- Regressions
- Bring your laptop with Python (and packages) installed