# An Intro to Machine Learning
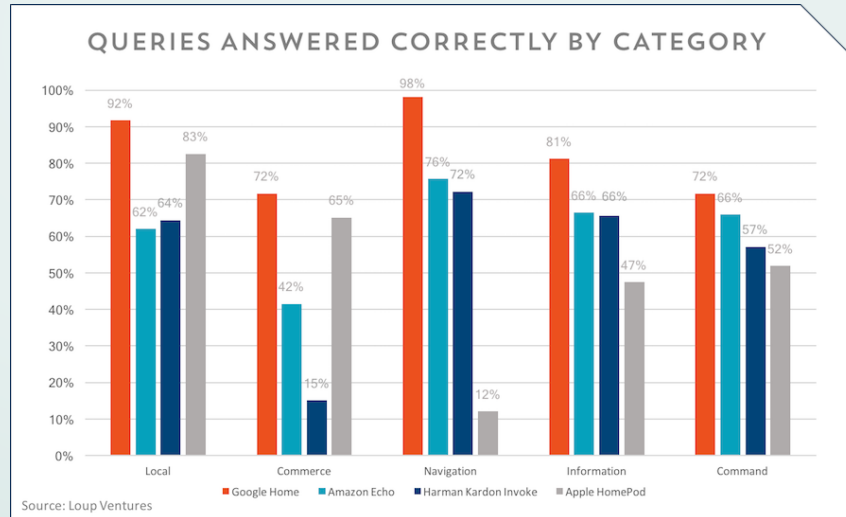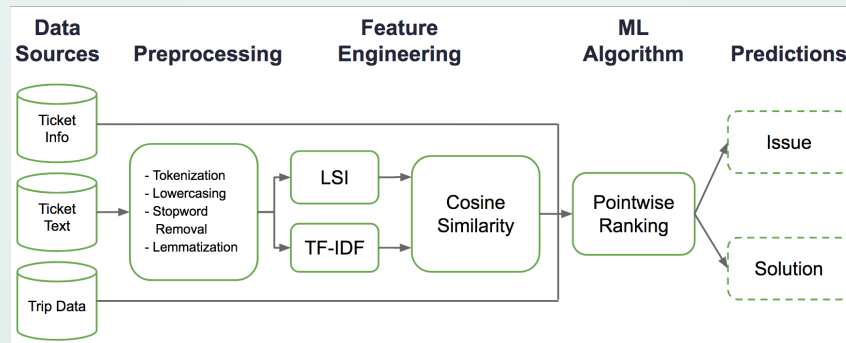
Rouhollah Abolhasani

# Speech Recognition

- **Alexa and Siri**
  - Extract vowels and consonants from the signals
  - Recognize Words, Sentences, and Sentiment
  - Act accordingly.



QUERIES ANSWERED CORRECTLY BY CATEGORY

Source: Loup Ventures

Google Home · Amazon Echo · Harman Kardon Invoke · Apple HomePod

Local: 92%, 62%, 64%, 83%
Commerce: 72%, 42%, 15%, 65%
Navigation: 98%, 76%, 72%, 12%
Information: 81%, 66%, 66%, 47%
Command: 72%, 66%, 57%, 52%

# Time-series Forecasting

- Time-series
  - Stock Market
  - Texts, Tweets, ...

- COTA: Customer Obsession Ticket Assistant
  - Handling Uber support tickets
  - Processing tickets and proposing solutions
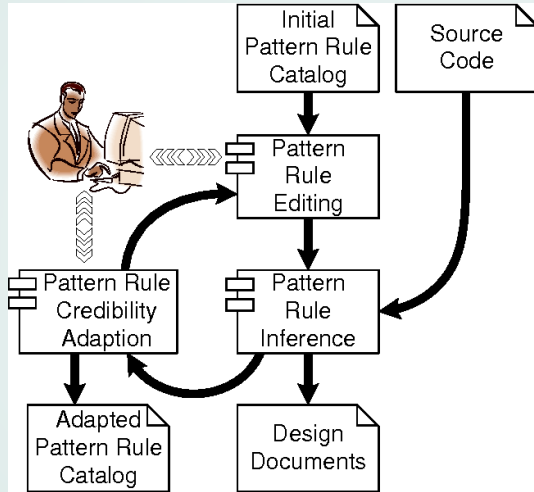


https://eng.uber.com/cota/

# Classification

- Image Classification
  - Cat vs. dog
  - Image segmentation

- Old-fashioned Spam Detector
  - What would we do?
  - Writing rules
  - More rules
  - And more rules …

- Better Spam Detector
  - Automatic feature extraction
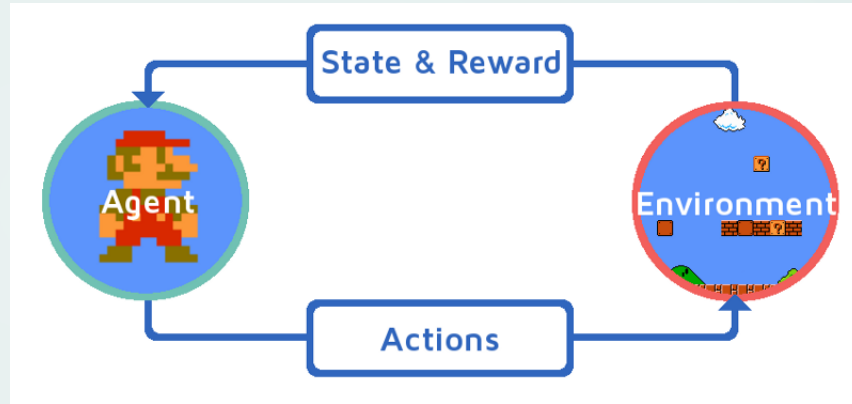  - Model training
  - Classifying new emails

# Rule-based vs. Learning



## Rule –based Programming

Write some rules. Evaluate them. Add some again, and again, …

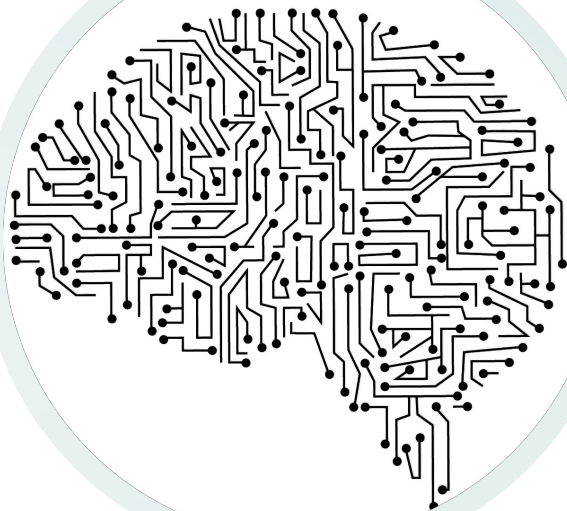## Learning

Machine programming machine with experience/data.

# Two Definitions of ML

**Arthur Samuel(1959):**
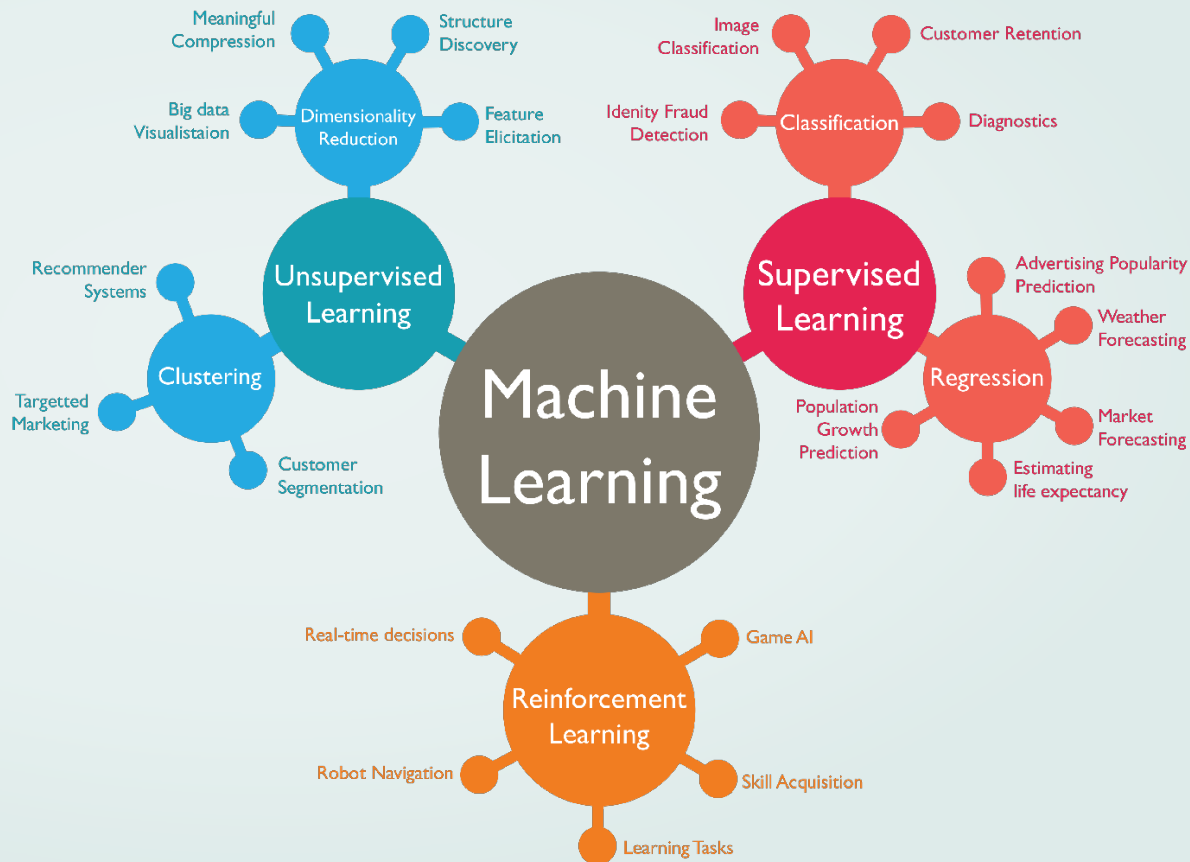- "Machine Learning is a field of study that gives computers, the ability to learn without explicitly being programmed."

**Tom Michel(1999)**
- "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."
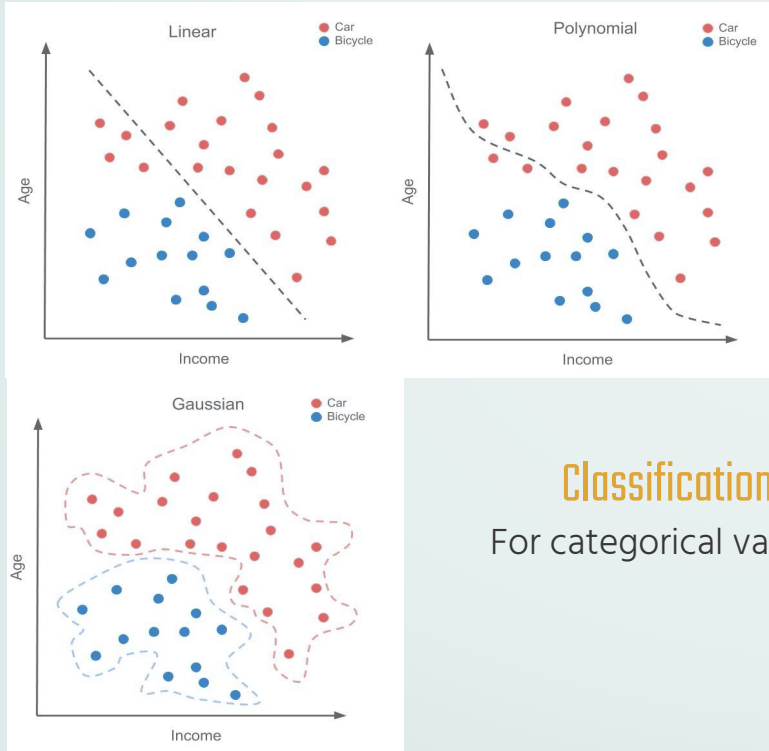
# The Big Picture

- Supervised Learning
  - Labeled data
  - Classification
  - Regression

- Unsupervised Learning
  - Unlabeled data
  - Clustering
  - Dimensionality Reduction

- Reinforcement Learning
  - Learning by living
  - Entirely different realm
  - Not covered in this course

**Dimensionality Reduction**
- Meaningful Compression
- Structure Discovery
- Big data Visualisation
- Feature Elicitation

**Unsupervised Learning**

**Clustering**
- Recommender Systems
- Targetted Marketing
- Customer Segmentation

**Classification**
- Image Classification
- Customer Retention
- Idenity Fraud Detection
- Diagnostics

**Supervised Learning**

**Regression**
- Advertising Popularity Prediction
- Weather Forecasting
- Population Growth Prediction
- Market Forecasting
- Estimating life expectancy

**Machine Learning**

**Reinforcement Learning**
- Real-time decisions
- Game AI
- Robot Navigation
- Skill Acquisition
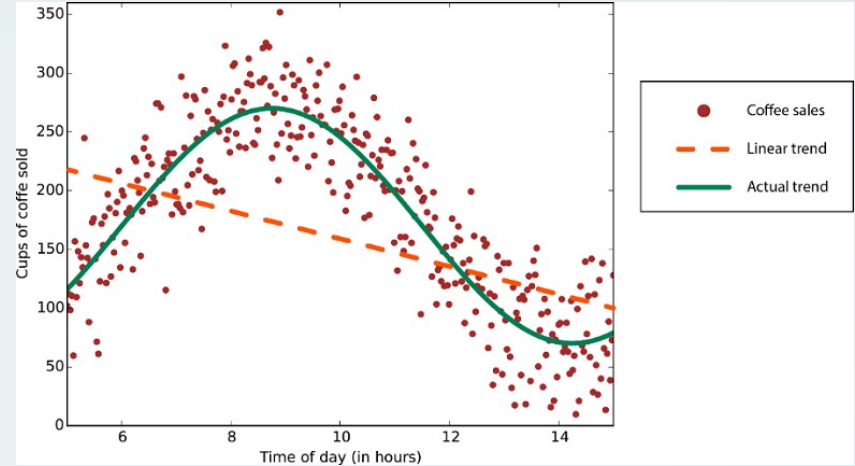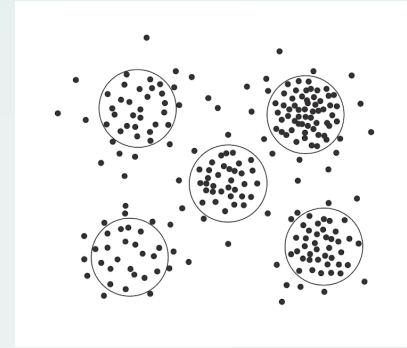- Learning Tasks

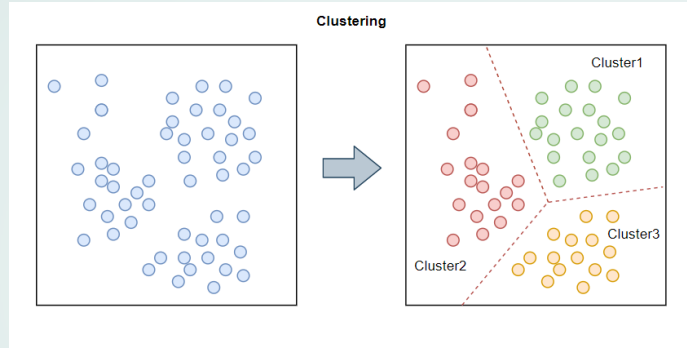# Supervised Learning



## Classification

For categorical variables
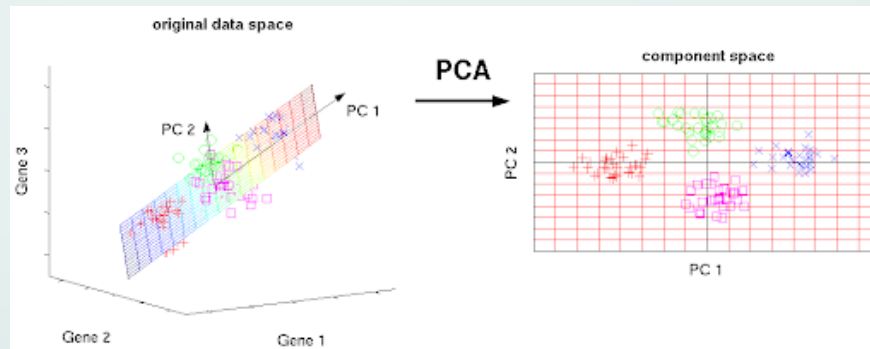
## Regression

For numeric variables
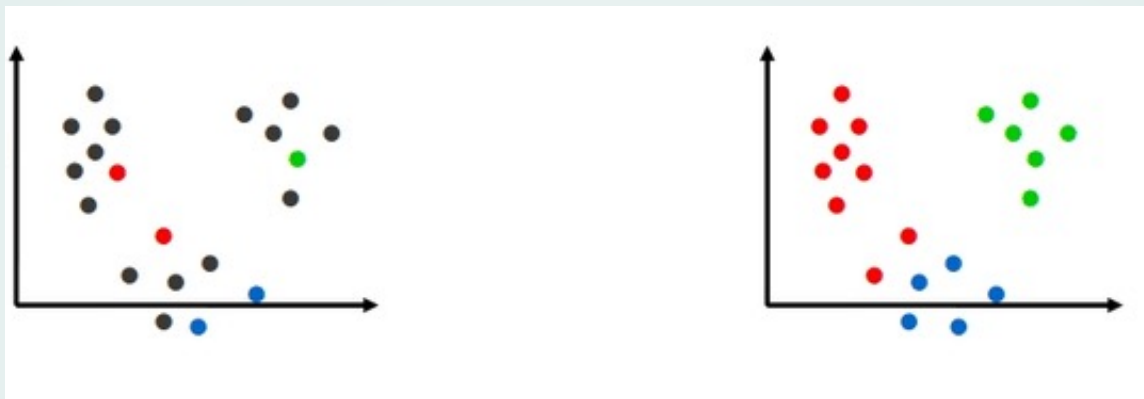
# Unsupervised Learning



Clustering



Clustering



Dimensionality Reduction

To avoid the curse of dimensionality

# Semi-supervised Learning



**Clustering and Classifying**

Deep Belief Networks for image classification, video recognition, ...

SciKit-Learn

2

# What is SciKit-Learn

- An open-source machine-learning library

- Built on Numpy, Scipy, and Matplotlib

- Contains many ML algorithms and models

- Good documentation and support

- Good for beginners

# Data in Scikit-Learn

- Data In ML typically consists of:
  - Features
  - Labels

- Features are stored in a 2D <u>matrix</u>
  - Shape=(n features, n data)

- Labels are stored in a 1D <u>vector</u>
  - Shape=(n data, )
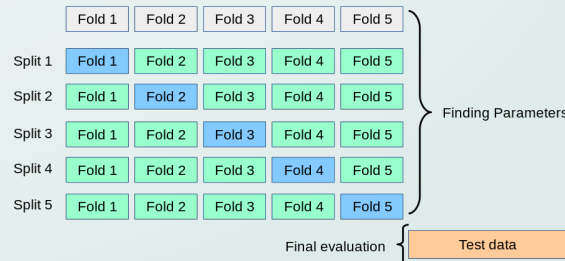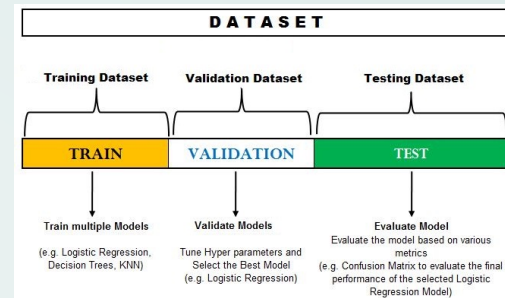
# Train-Test Data

- We always split data to:

  - Train data

  - Validation data

  - Test data: We only see in production

- Cross-validation

  - Holdout

  - K-fold

- In Scikit-Learn:
  **from sklearn.model_selection import train_test_split**
  **X_train, X_test, y_train, y_test = train_test_split(X, y, train_size)**

# Some Basic Regression Model

- Using LinearRegression model

from sklearn.linear_model import LinearRegression

- Constructing the model

reg = LinearRegression()

# Model API In Scikit-Learn

- Train the model using .fit(X, y) method

- Transform the data using .transform(X) method

- Do the combination using .fit_transform(X, [y])

- Test the model using .predict(X) method

# K-Nearest Neighbor Classifier

- For each new data:
  - Find the k-closest data
  - Take majority vote
  - Assign the label