

Advanced Business Data Mining

MSIS 522 – Lesson 1

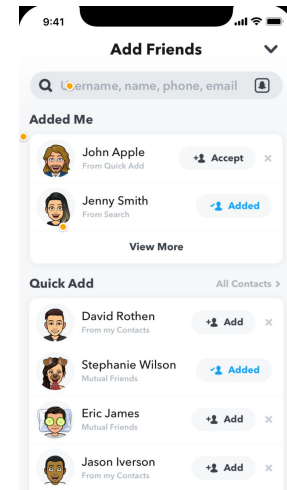
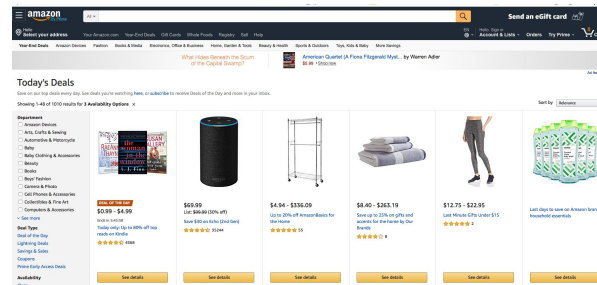
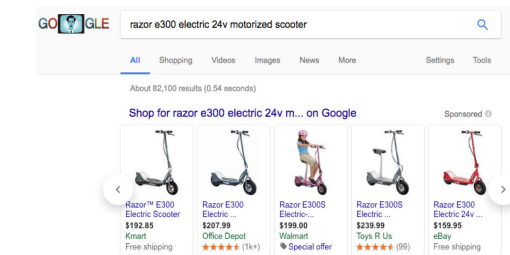


Basic Info

Time	Mon Night, Wed Night & Sat
Location	PACCAR 292 (Gold) & PACCAR 294 (Purple)
Instructor	Jun Yu (jnyu@uw.edu)
TA	Xiyuan Ge (xiyuange@uw.edu)
Textbook	<ul style="list-style-type: none">• <u>An Introduction to Statistical Learning</u> (free online)• A Course in Machine Learning (free online)• Mining of Massive Datasets (free online)• <u>Python Data Science Handbook</u> (free online)
Software	Python on <u>Google Colab</u>
Grading	<ul style="list-style-type: none">• 3 Individual assignments (90%) [20% per 24 hours applied to late submissions]• Participation (10%)

About Me

- Ph.D. in CS @ Oregon State University
 - Species Distribution Modeling for the eBird project
- Applied Researcher @ eBay
 - Paid Internet Marketing
 - Seller Risk Management
- Senior Applied Scientist @ Amazon
 - Deals Forecasting and Scheduling
 - Seller Recommendation System
- Engineering Manager @ Snapchat
 - Friend Recommendation



Jeff Bezos: The Golden Age of Artificial Intelligence



“Machine Learning and AI is a horizontal enabling layer. It will empower and improve every business, every government organization, every philanthropy, basically there is no institution in the world that cannot be improved with Machine Learning.” – Jeff Bezos

House & Vehicle Value Estimates

REDFIN City, Address, School, Agent, Zip

1-877-873-2048 Buy & Sell Real Estate Agents Tools Sign In Join

2718 Elliott Ave Seattle, WA
Status: Active

\$1,295,000 Price
\$6,361 Refund
4 Beds
3.5 Baths
3,730 Sq Ft
\$347 / Sq Ft
Redfin Estimate: \$1,294,799 On Redfin: 4 days

Overview Property Details Tour Insights Redfin Estimate Property History Public Facts Schools Neighborhood Similar Homes



Listing provided courtesy of Windermere RE Greenwood

Hot Home: This home has an offer review date in effect – go tour it now.

Go Tour This Home

MONDAY 10 A.M. TUESDAY 11 A.M. WEDNESDAY 12 A.M.

Schedule Tour

It's free, with no obligation – cancel anytime

Ask a Question

Redfin Refund: \$6,361 ©
Savings when you buy with a Redfin Agent
Start an Offer

Kelley Blue Book The Trusted Resource Home Car Values Cars for Sale Car Reviews Awards & Top 10s Research Tools

Home > What's My Car Worth > Options & Condition > Sport Utility 4D

2010 Acura MDX Sport Utility 4D near Bellevue, WA 98004

Mileage: 50,000 Condition: Very Good Save

Overall Consumer Rating 4.6 / 5
★★★★☆ 303 Ratings Write a review Check Specs

We found 1 recall. See details.

1 Compare Your Values See Overview of Values

Use these values to help make a confident decision on whether to sell, trade or donate your car.

Instant Cash Offer Trade-in Value Private Party Value Donate Your Car

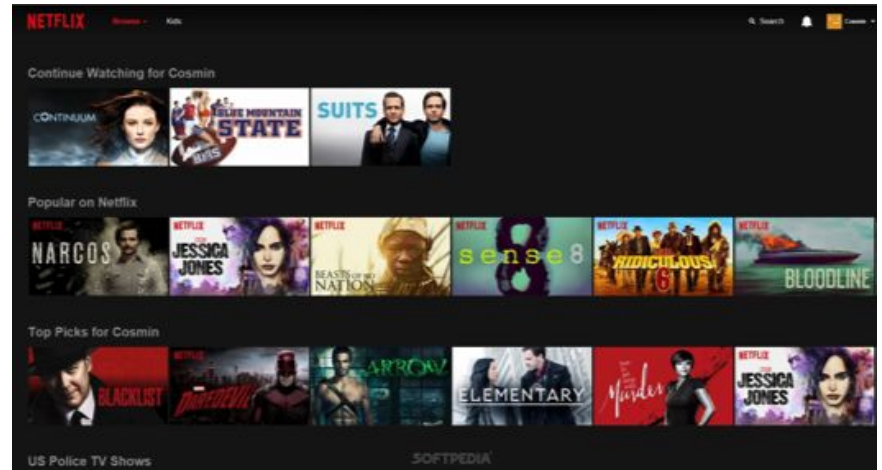
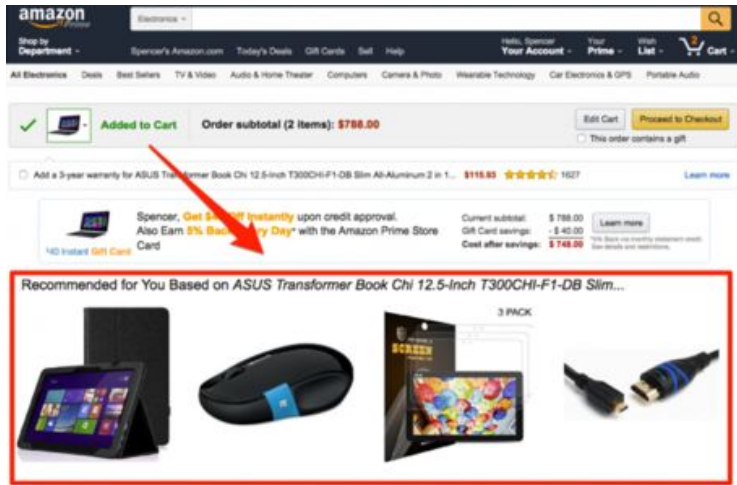
Trade-in Value

Trade-in Range: **\$13,624 - \$15,705**
Trade-in Value: **\$14,665**

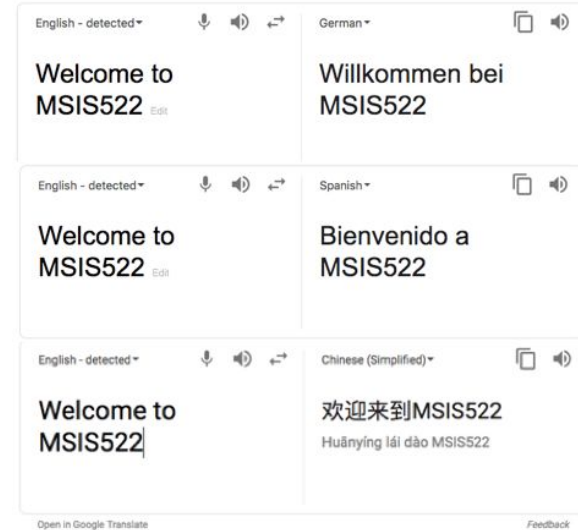
TRADE-IN VALUE
This estimated value helps you confidently negotiate with dealers.
Average Time to Trade-In: **1-7 DAYS**
Level of Effort: **Medium**

Important info & definitions Track this car's values See Overview of Values

Product & Movie Recommendation

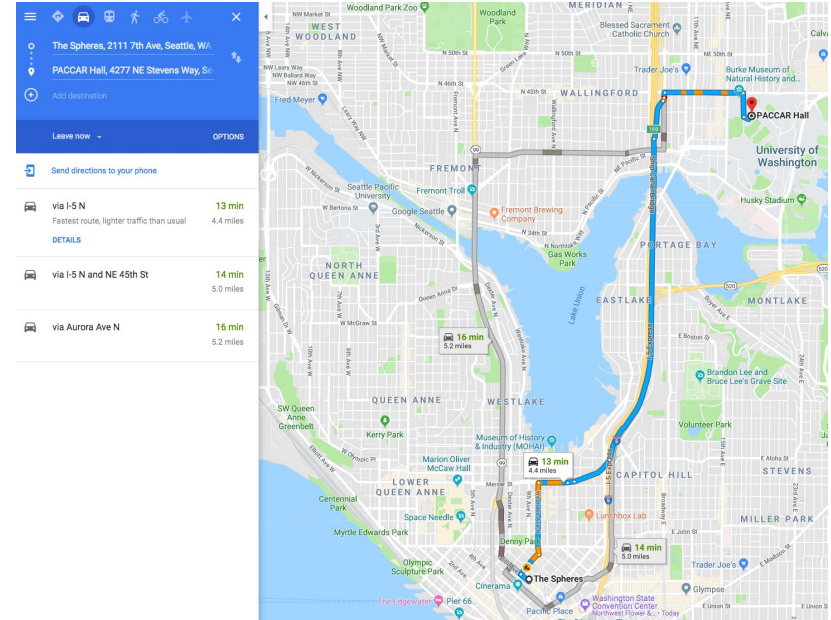
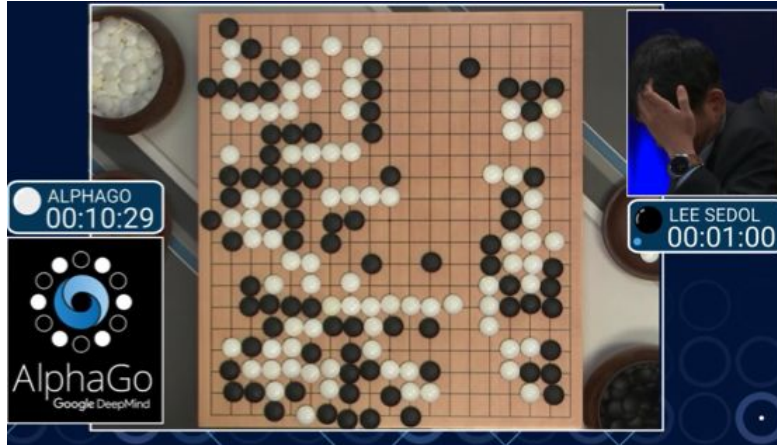


Voice Recognition & Machine Translation



Google Translate

Play Game of GO & Route Planning



Machine Learning Everywhere!



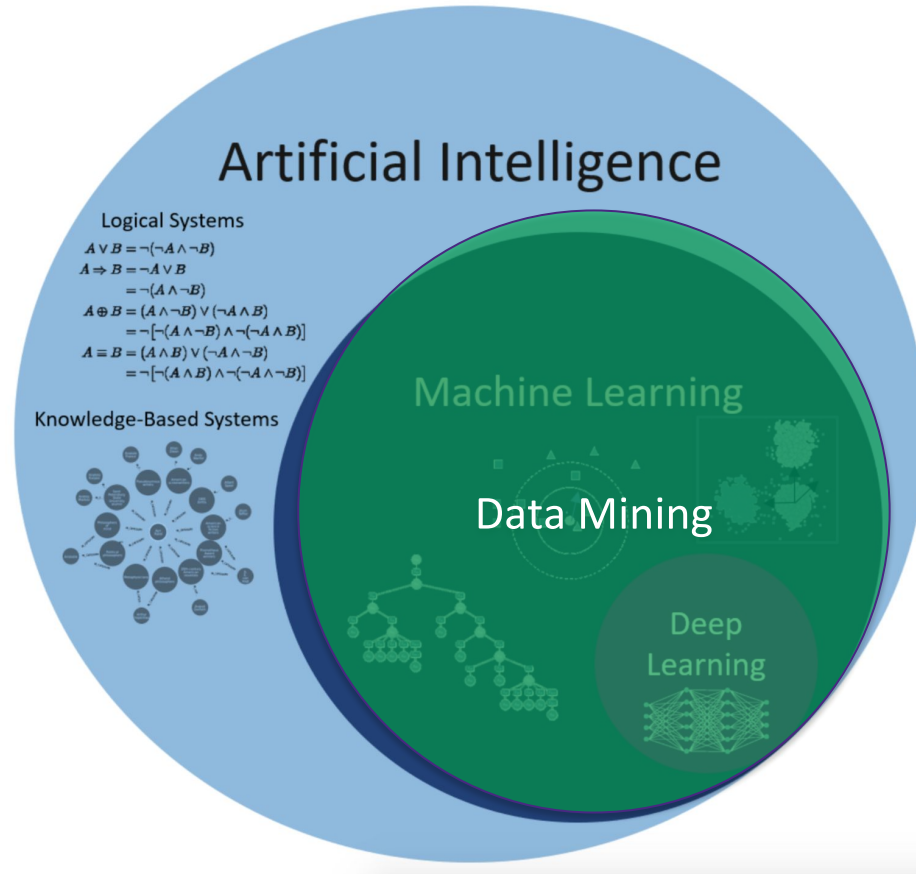
A real life example of Machine Translation



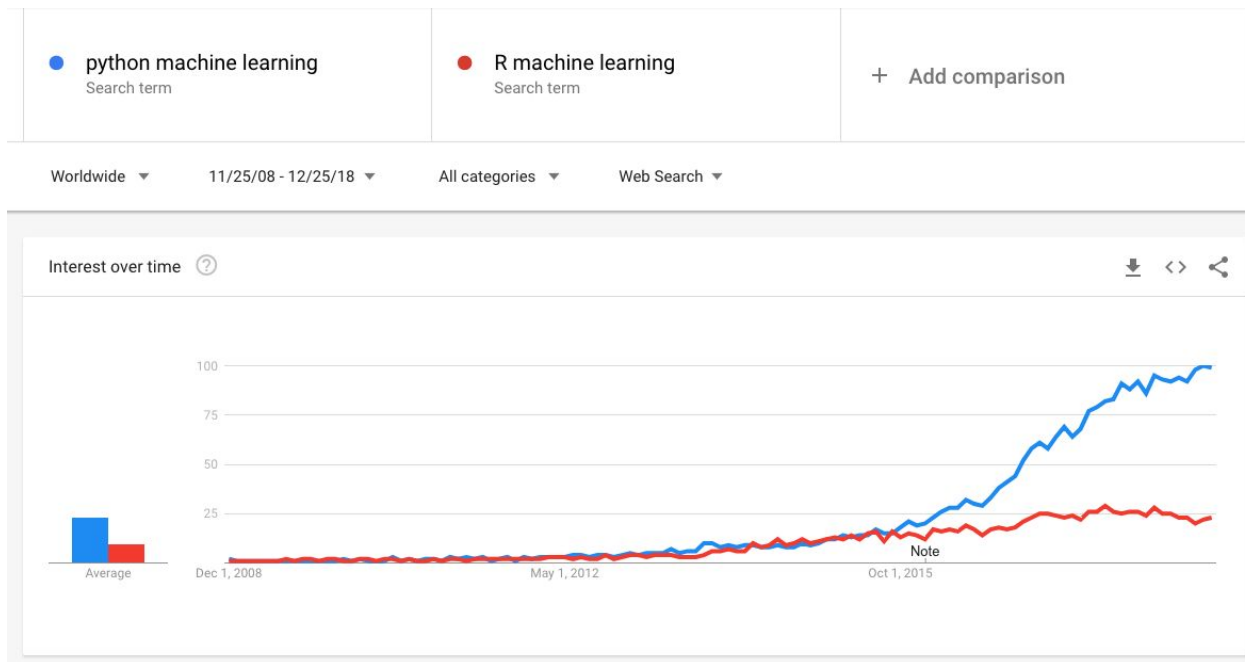
Quiz

Can you think of any real world applications driven by ML?

AI vs. ML vs. DL



Why Python (vs. R) for ML?



Data Science Wars: R vs. Python

Why Python (vs. rest)?

PYPL Popularity of Programming Language

Worldwide, Jan 2019 compared to a year ago:

Rank	Change	Language	Share	Trend
1	↑	Python	25.95 %	+5.2 %
2	↓	Java	21.42 %	-1.3 %
3	↑	Javascript	8.26 %	-0.2 %
4	↑	C#	7.62 %	-0.4 %
5	↓↓	PHP	7.37 %	-1.3 %
6		C/C++	6.31 %	-0.3 %
7		R	4.04 %	-0.2 %



Life is short, use **Python**

Theory vs. Practice

THEORY is when you know everything but nothing works.

PRACTICE is when everything works but no one knows why.

In this lab, THEORY and PRACTICE are combined: nothing works and no one knows why.

- 50% Theory + 50% Practice (lab).
- Work hard on the homework.
- Explore and participate in online ML competitions.

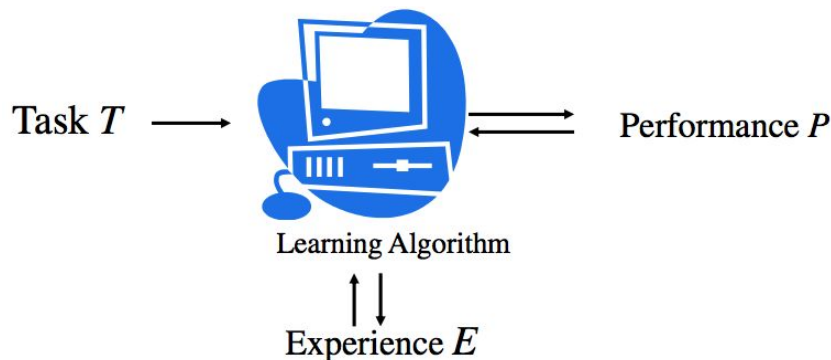
Syllabus

- Lecture 1 - Fundamentals of Machine Learning
- Lecture 2 - Decision Tree
- Lecture 3 - Ensemble Learning
- Lecture 4 - Clustering
- Lecture 5 - Recommendation System

Fundamentals of Machine Learning

What is Machine Learning?

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 . -- Tom Mitchell



Improving *performance* P with *experience* E at some *task* T .

Examples

Machine Learning: Improving *performance P* with *experience E* at some *task T*.

Task

Predict house price

Recommend products

Voice recognition

Play Go

Experience

Historical house sales

Customer's transactions

Human voice clips

Practice games of Go

Performance

Error in price estimate

Click through rate

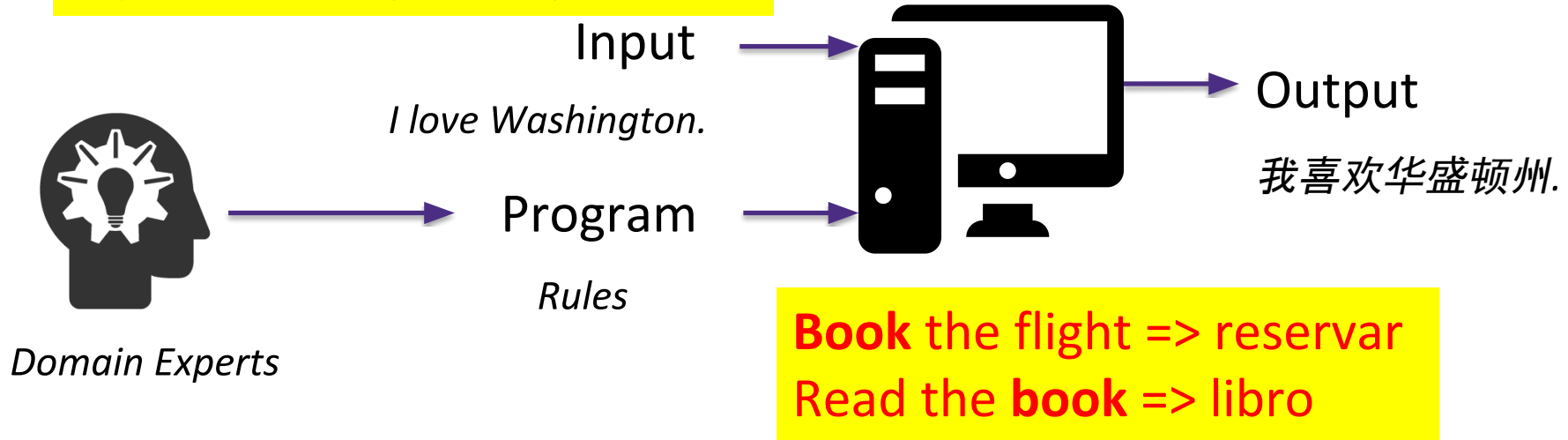
Accuracy of transcription

% of games won

Machine Translation: Traditional Programming

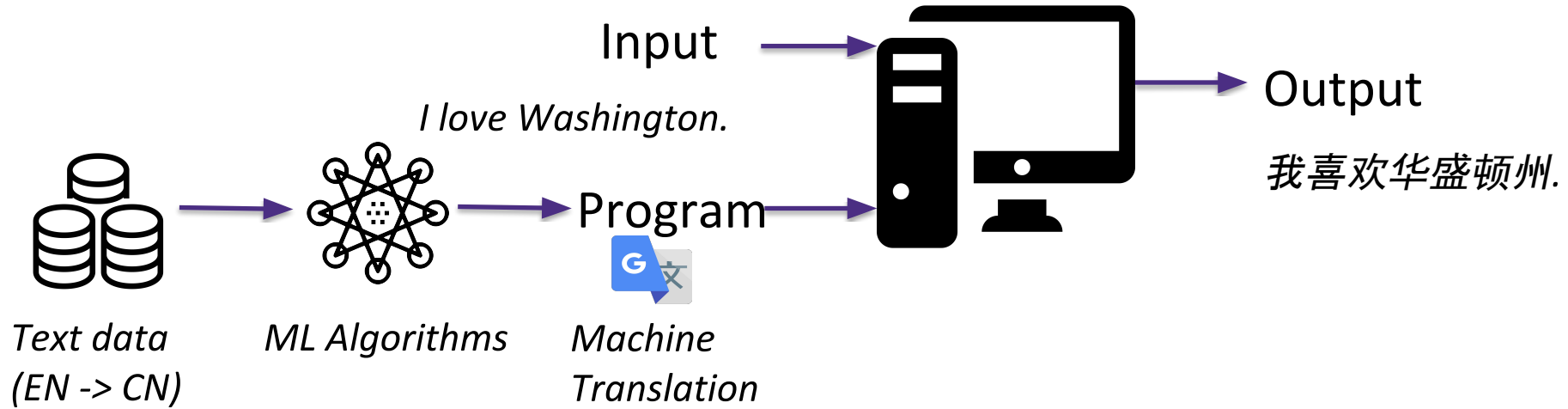
English: subject-verb-object

Japanese: subject-object-verb



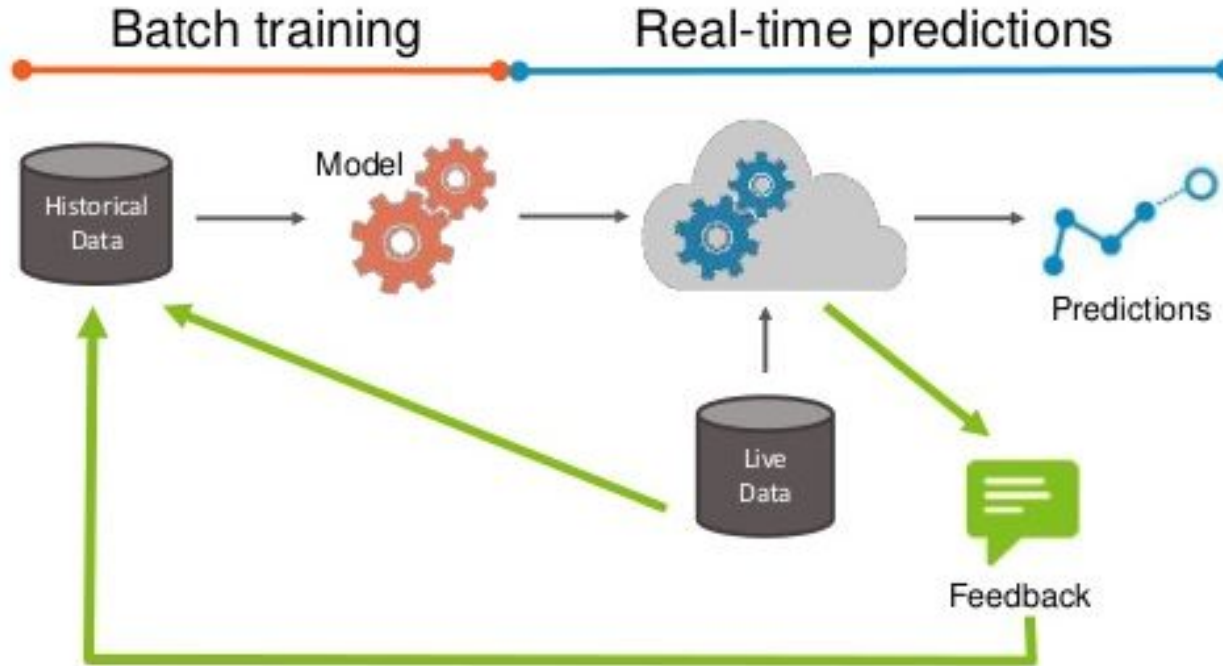
Rule-based machine translation. (1950 – 2000)

Machine Translation: Machine Learning



Google Machine Translation. (2013 – now)

Feedback Loop on Continuous Learning



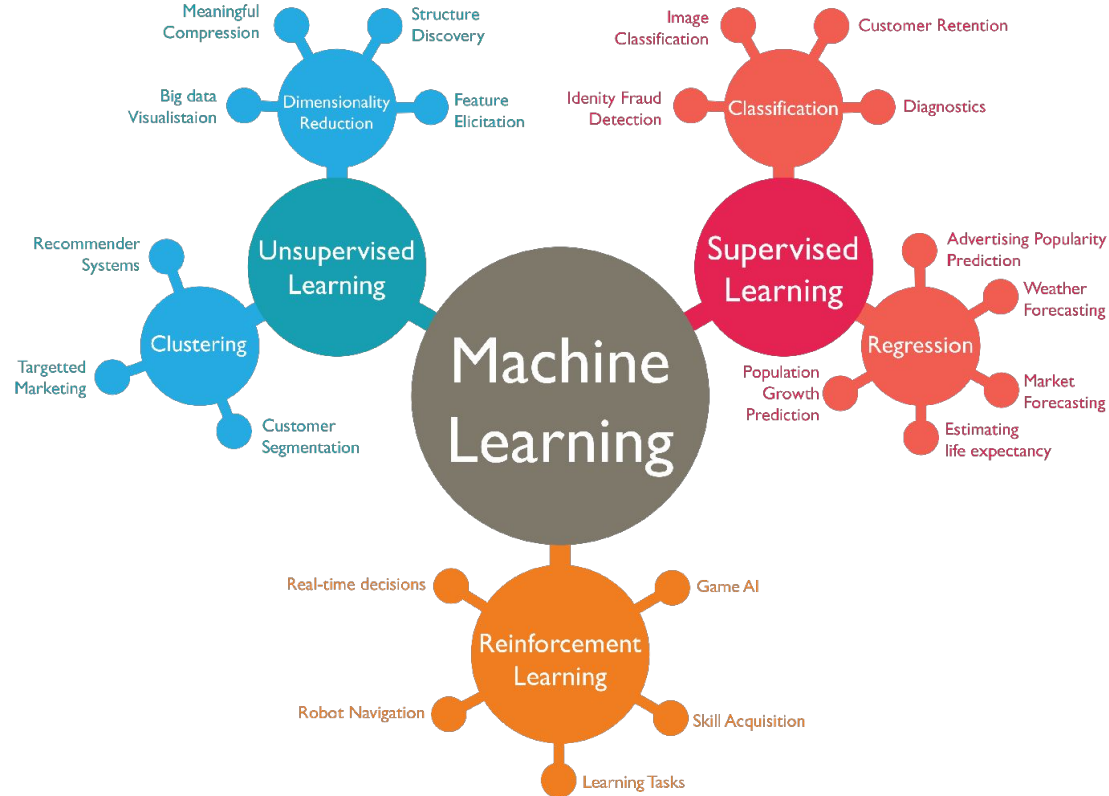
Quiz

Can you think of any scenarios when ML is preferred to human?

When do we need Machine Learning?

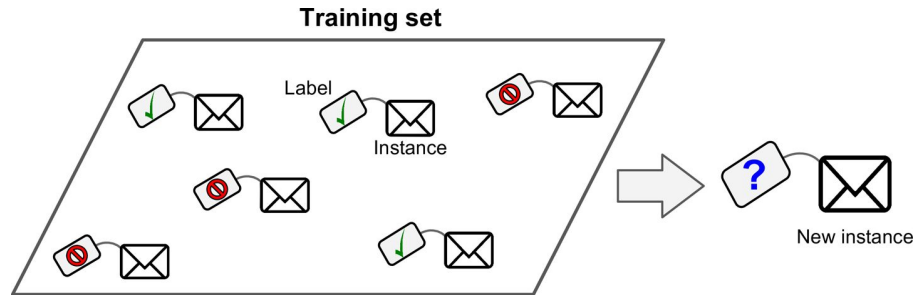
- Humans cannot do it or can do it but cannot describe how they do it (thus cannot be programmed).
 - Effectiveness of a new compound treating some disease (cannot do)
 - Object recognition (cannot describe)
- Machines can do better than humans on certain tasks.
 - Search page ranking (more data)
 - Play the game of Go (more computing power)
- Machines are cheaper than humans.
 - Handwritten digits recognition
 - Voice recognition

3 Types of Machine Learning Algorithms



Supervised Learning

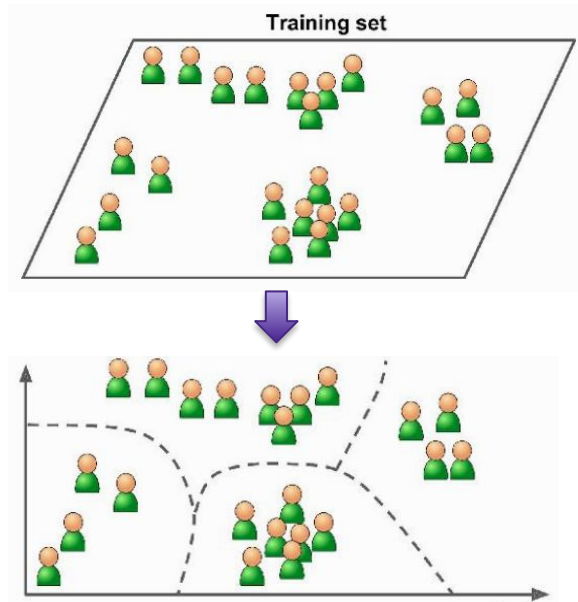
Learn to predict output from input.



- Linear/Logistic Regression
- K Nearest Neighbor
- Decision Tree
- Random Forest
- Support Vector Machine
- (Deep) Neural Network
- ...

Unsupervised Learning

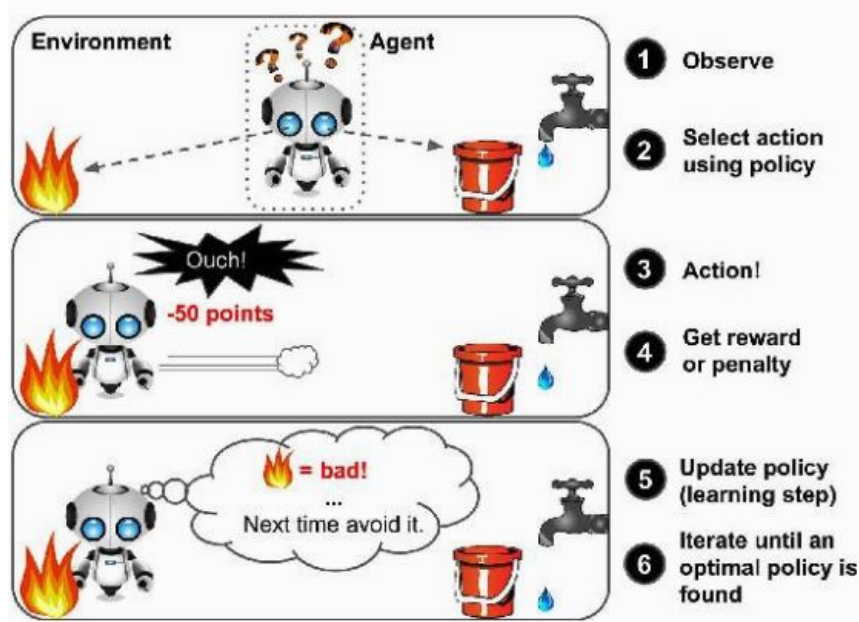
Learn to find patterns and structure in an **unlabeled** data set.



- K means
- Hierarchical clustering
- Gaussian Mixture Model
- DBSCAN
- Spectral clustering
- Mean Shift clustering
- ...

Reinforcement Learning

Learn to **take actions** in an environment to maximize expected future rewards.



- Model-based RL
- Model-free RL

Quiz

Task	Learning Type
Home value estimates on Redfin	Supervised
Product recommendations on Amazon	Unsupervised
Voice recognition in Alexa	Supervised
Play the game of Go	Reinforcement
Predict stock price	Supervised
An advertising platform segments the U.S. population into smaller groups with similar demographics and purchasing habits so that advertisers can reach their target market with relevant ads.	Unsupervised
Learn to play StarCraft	Reinforcement
Airbnb groups its housing listings into neighborhoods so that users can navigate listings more easily.	Unsupervised
File tax return	Not a learning problem

Supervised Learning

- **Regression:** A regression model predicts **continuous values**. For example, regression models make predictions that answer questions like the following:
 - What is the value of a house in California?
 - What is the demand of a product on Amazon next month?
- **Classification:** A classification model predicts **discrete values**. For example, classification models make predictions that answer questions like the following:
 - Is a given email message spam or not spam?
 - Is this an image of a dog, a cat, or a hamster?

Supervised Learning Terminology

Target Variable	The thing we're predicting, aka label. - e.g. Future stock price, cat image, spam email
Feature	An input variable, which is used to predict the label. - e.g. Historical stock price, pixels in the image, words in the email body

	Age	Gender	Weight	Height
Observation #1	12	M	80 lbs	55 in
Observation #2	11	M	85 lbs	58 in
Observation #3	12	F	73 lbs	52 in
Observation #4	10	F	71 lbs	49 in
.				
.				
.				
Observation #150				

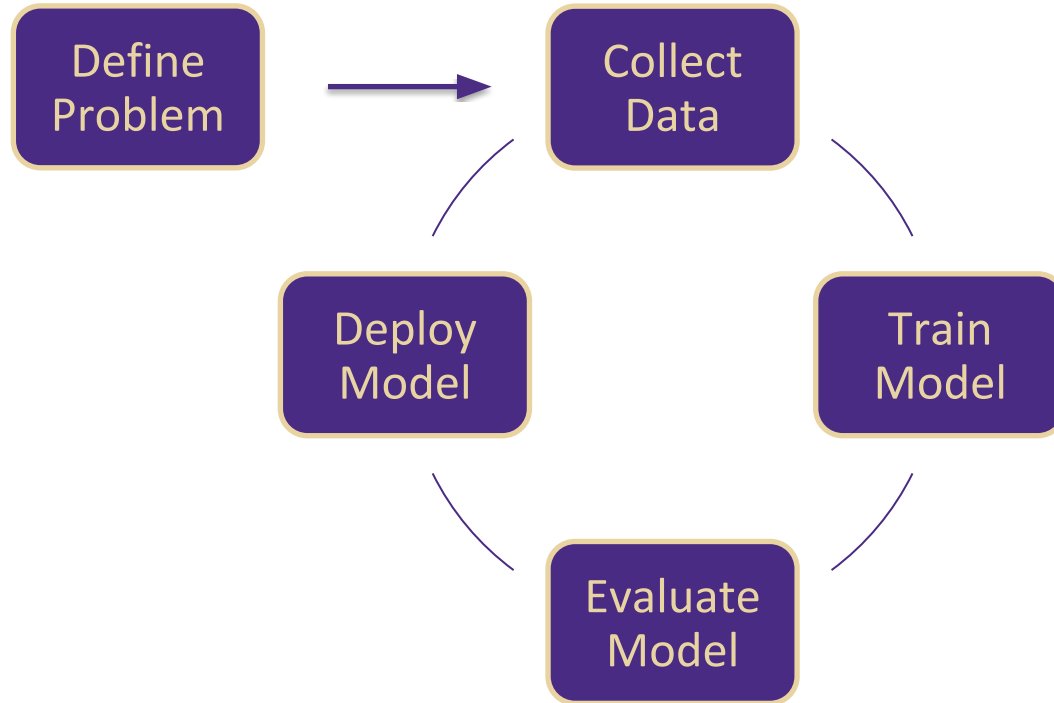
Features

Target Variable

Supervised Learning Terminology

Model Training	Build a ML model by learning the relationships between features and label. - e.g. linear regression, decision tree, neural network
Model Scoring	Apply a ML model on unlabeled examples to make predictions. - e.g. linear regression, decision tree, neural network
Model Evaluation	Measure the model's predictive accuracy. - Mean Squared Error, Accuracy, Precision and Recall and AUC.

Supervised Learning Workflow



Confusion Matrix

Misclassification Error: classifying a record as belonging to one class when it belongs to another class.

		Actual Class	
		C1 (Positive)	C0 (Negative)
Predicted Class	C1 (Positive)	Number of records with actual class 1 and predicted class 1	Number of records with actual class 0 and predicted class 1
	C0 (Negative)	Number of records with actual class 1 and predicted class 0	Number of records with actual class 0 and predicted class 0

ERRORS
!!!

Classification Performance Measures

Most classification accuracy measures are derived from the **confusion matrix**.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

		Actual Class	
		C1 (Positive)	C0 (Negative)
Predicted Class	C1 (Positive)	True Positive (TP)	False Positive (FP)
	C0 (Negative)	False Negative (FN)	True Negative (TN)

Quiz

Fill in the following confusion matrix.

		Actual Class	
		C1 (Positive)	C0 (Negative)
Predicted Class	C1 (Positive)	1	1
	C0 (Negative)	1	2

Actual class	Predicted class
C1	C0
C1	C1
C0	C0
C0	C0
C0	C1

Accuracy = ? 3 / 5

Precision = ? 1 / 2

Recall = ? 1 / 2

Regression Performance Measures

Regression Error: the difference between the predicted value (Z_i) and the actual value (Y_i).

$$\text{Root Mean-Squared Error} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - z_i)^2}$$

$$\text{Mean Absolute Error} = \frac{1}{N} \sum_{i=1}^N |y_i - z_i|$$

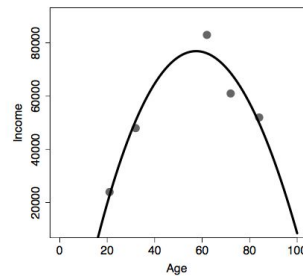
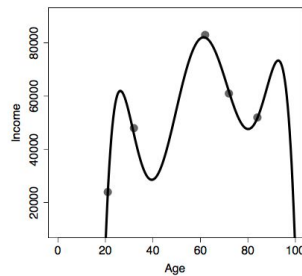
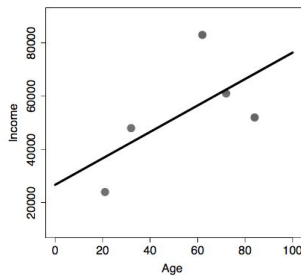
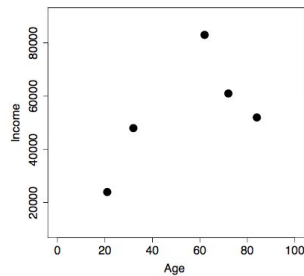
Generalization

Machine Learning is all about **generalization** to future unseen data points.

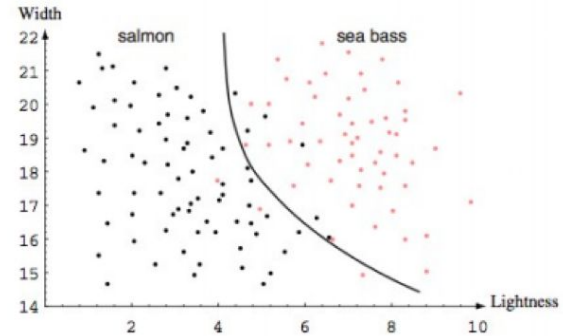
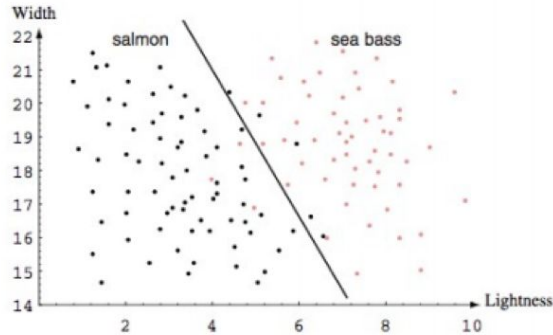
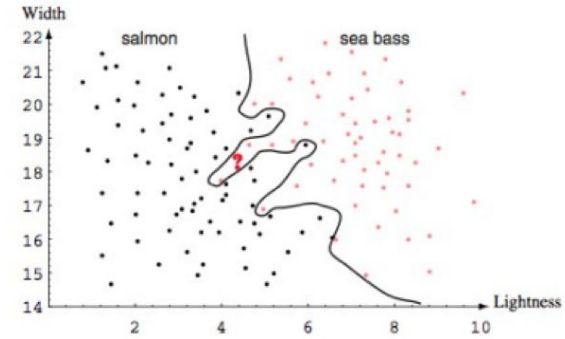
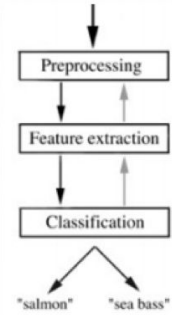
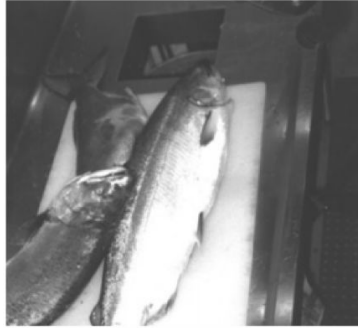
- **Underfitting** – a model is too simple and can not capture the underlying patterns within the data, thus does not perform well on new data.
- **Overfitting** – a model tries to fit the training data so closely that it does not generalize well to new data.

Table: The age-income dataset.

ID	AGE	INCOME
1	21	24,000
2	32	48,000
3	62	83,000
4	72	61,000
5	84	52,000



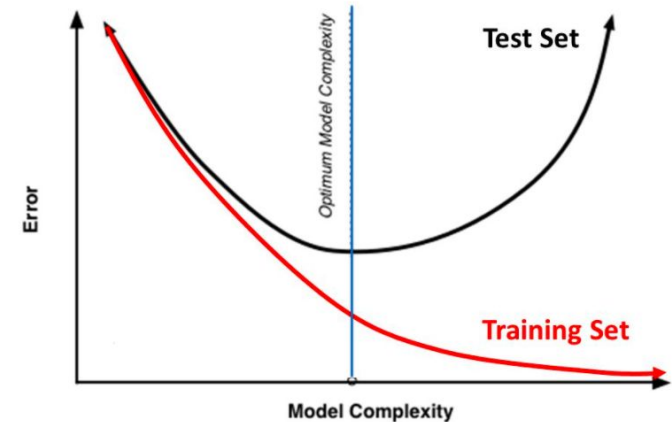
Quiz



Model Selection

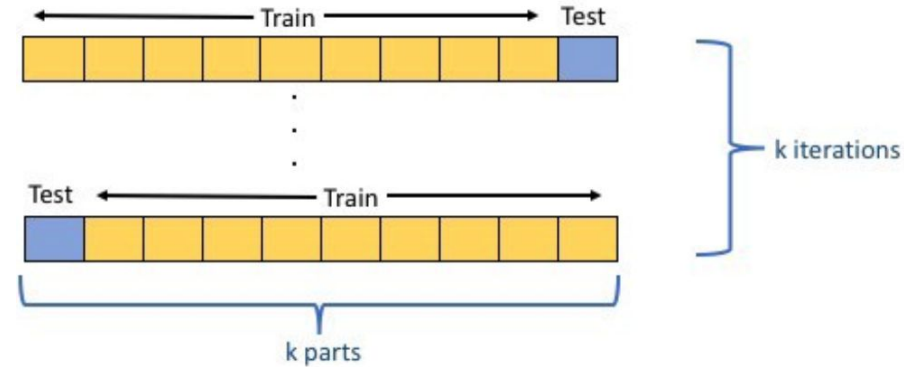
- Split the data into two sets:
 - **Training Set:** a subset of the data to train a model.
 - **Test Set:** a subset of the data to evaluate the performance of the trained model.
- Plot the errors of models with increasing complexity on both train and test sets.
- Choose the model which optimize the error on the test set.

Training Vs. Test Set Error



K-fold Cross-validation

1. Divide the training data into K parts.
2. Use $K-1$ of the parts for training and 1 for testing.
3. Repeat the procedure K times, rotating the test set.
4. Determine the performance based on the results across all K iterations.



Leave-one-out Cross-validation is the extreme case of K-fold Cross-Validation where we keep only one data point in the test set.

Quiz

What are the leave-one-out cross-validation errors for the following dataset, using 1 nearest neighbor (1-NN) and 3 nearest neighbor (3-NN)?



Answers: 1-NN: 5/10. 3-NN: 1/10.

Lab

Python Machine Learning Libraries

- ***Pandas*** provides a DataFrame object along with a powerful set of methods to manipulate, filter, group, and transform data.
- ***Matplotlib*** provides a useful interface for creation of high-quality plots and figures.
- ***Seaborn*** is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
- ***Scikit-Learn*** provides a uniform toolkit for applying common machine learning algorithms to data.