

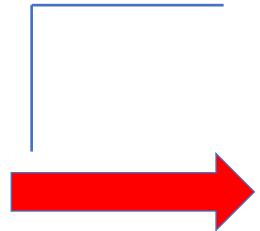
UVA CS 6316: Machine Learning

Lecture 1: Introduction

Dr. Yanjun Qi

University of Virginia
Department of Computer Science

Today

- 
- Course Logistics
 - Machine Learning Basics
 - A Rough Plan of Course Content
 - Machine Learning History

Welcome

- Course Website:
 - <https://qianjunqi.github.io/2019f-UVA-CS6316-MachineLearning/>

We focus on learning fundamental principles, mathematical formulation, algorithm design and learning theory.

Objective

- To help students able to build simple machine learning tools
 - (not just a tool user!!!)
- Key Results:
 - Able to build a few simple machine learning methods from scratch
 - We focus on developing skills, more than feeding knowledge
 - paper and understand the pipeline

Course Staff

- Instructor: Prof. Yanjun Qi
 - QI: /ch ee/
 - You can call me “professor”, “professor Qi”;
 - I have been teaching Graduate-level and Under-Level Machine Learning course for years!
 - My research is about machine learning
- TA and Office Hour information @ CourseWeb

Course Material

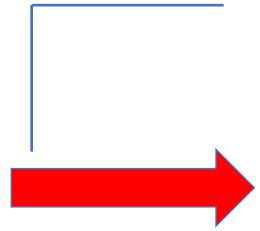
- Text books for this class is:
 - NONE
- My slides – **if it is not mentioned in my slides, it is not an official topic of the course**
- Your UVA Collab for Assignments and Project
- Google Forms for Quizzes

Course Background Needed

- **Background Needed**

- Calculus, Basic linear algebra, Basic probability and Basic Algorithm
- Statistics is recommended.
- Students should already have good programming skills, i.e. **python** is required for all programming assignments

Today

- 
- Course Logistics
 - Machine Learning Basics
 - Machine Learning History
 - Rough Plan of Course Content

OUR DATA-RICH WORLD

- Biomedicine
 - Patient records, brain imaging, MRI & CT scans, ...
 - Genomic sequences, bio-structure, drug effect info, ...
- Science
 - Historical documents, scanned books, databases from astronomy, environmental data, climate records, ...
- Social media
 - Social interactions data, twitter, facebook records, online reviews, ...
- Business
 - Stock market transactions, corporate sales, airline traffic, ...
- Entertainment
 - Internet images, Hollywood movies, music audio files, ...



What can we do with the data wealth?

→ REAL-WORLD IMPACT

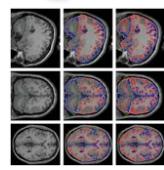
Transportation
Data



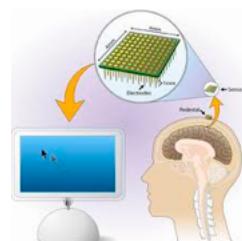
Genomic Data



Medical Images



Brain computer
interaction (BCI)



0
1
0010
1010
011

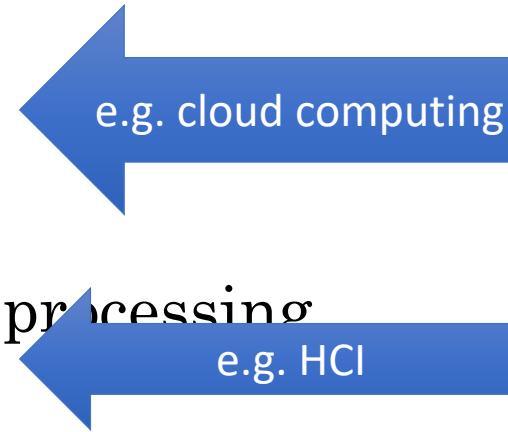
Device sensor data
9/2/19



- Business efficiencies
- Scientific breakthroughs
- Improve quality-of-life:
 - healthcare,
 - energy saving / generation,
 - environmental disasters,
 - nursing home,
 - transportation,
 - ...

BIG DATA CHALLENGES

- Data capturing (sensor, smart devices, medical instruments, et al.)
- Data transmission
- Data storage
- Data management
- High performance data processing
- Data visualization
- Data security & privacy (e.g. multiple individuals)
-



this
course

- Data analytics
 - How can we analyze this big data wealth ?
 - E.g. Machine learning and data mining

MACHINE LEARNING IS CHANGING THE WORLD

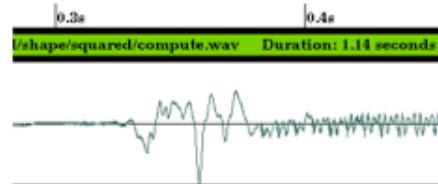
Data:

```
PatientID1 time=1 --> PatientID2 time=2 --> PatientID3 time=3
Age: 23   Age: 23   Age: 23
FirstPregnancy: no   FirstPregnancy: no   FirstPregnancy: no
Anemia: no   Anemia: no   Anemia: no
Diabetes: no   Diabetes: YES   Diabetes: no
PreviousPrematureBirth: no   PreviousPrematureBirth: yes   PreviousPrematureBirth: no
Ultrasound: ?   Ultrasound: abnormal   Ultrasound: abnormal
Elective C-Section: ?   Elective C-Section: no   Elective C-Section: no
Emergency C-Section: ?   Emergency C-Section: ?   Emergency C-Section: Yes
```

One of 18 learned rules:

```
If No previous vaginal delivery, and
Abnormal 2nd Trimester Ultrasound, and
Malpresentation at admission
Then Probability of Emergency C-Section is 0.6
```

Over training data: 26/41 = .63,
Over test data: 12/20 = .60



Speech Recognition



Control learning



Object recognition

Mining Databases

Text analysis

Peter H. van Oppen, Chairman of the Board & Chief Executive Officer
Mr. van Oppen has served as chairman of the board and chief executive officer of ADIC since its acquisition by Interpoint in 1994 and a director of ADIC since 1986. Until its acquisition by Crane Co. in October 1996, Mr. van Oppen served as chairman of the board of directors, president and chief executive officer of Interpoint. Prior to 1985, Mr. van Oppen worked as a consulting manager at Price Waterhouse LLP and at Bain & Company in Boston and London. He has additional experience in medical electronics and venture capital. Mr. van Oppen also serves as a director of Seattle FilmWorks Inc. and Spacelabs Medical, Inc.. He holds a B.A. from Whitman College and an M.B.A. from Harvard Business School, where he was a Baker Scholar.

BASICS OF MACHINE LEARNING

- “The goal of machine learning is to build computer systems that can **learn and adapt from their experience.**” – Tom Dietterich
- “**Experience**” in the form of available **data examples** (also called as instances, samples)
- Available examples are described with properties (**data points in feature space X**)

e.g. SUPERVISED LEARNING

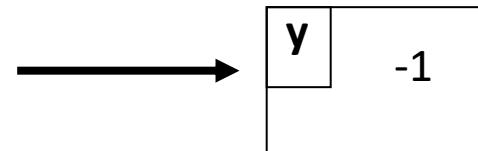
- Find function to map **input** space X to **output** space Y

$$f : X \longrightarrow Y$$

- So that the **difference** between y and $f(x)$ of each example x is small.

e.g.

x	I believe that this book is not at all helpful since it does not explain thoroughly the material . it just provides the reader with tables and calculations that sometimes are not easily understood ...
----------	--



Output Y: {1 / Yes , -1 / No }
e.g. Is this a positive product review ?

Input X : e.g. a piece of English text

SUPERVISED Linear Binary Classifier

- Now let us check out a **VERY SIMPLE** case of

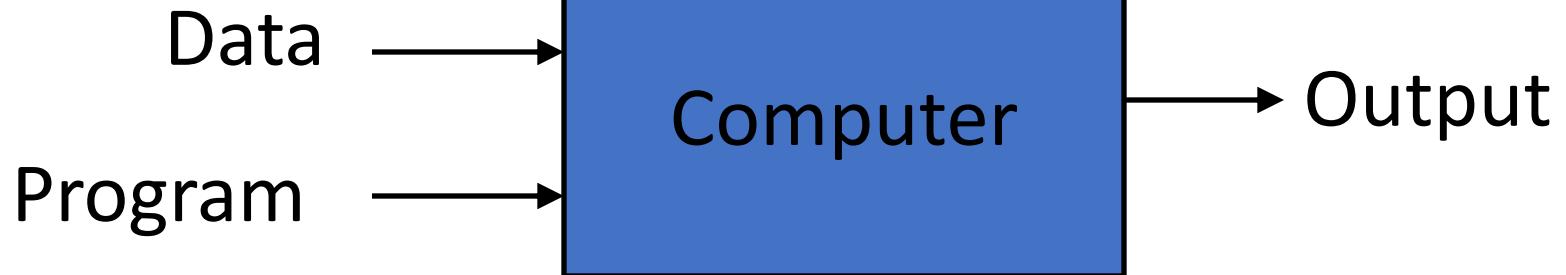


e.g.: Binary y / Linear f / X as \mathbb{R}^2

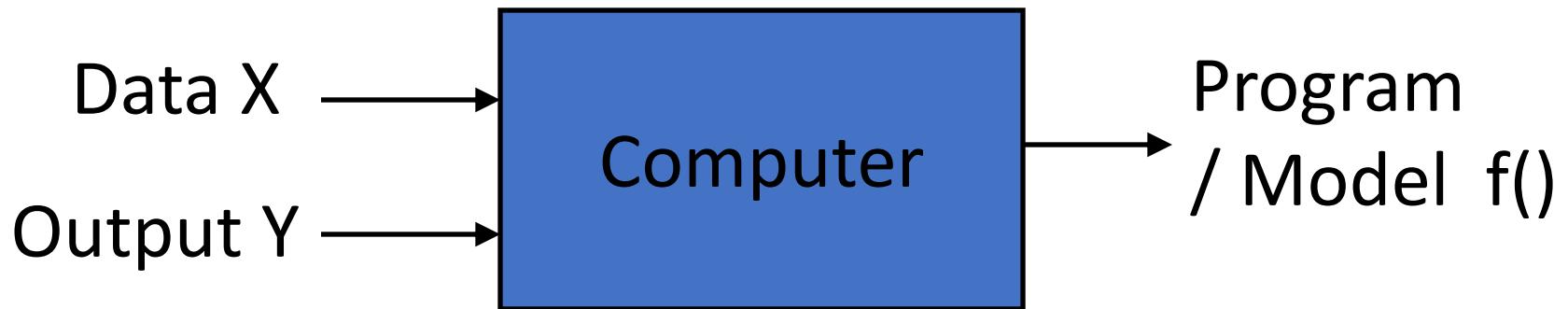
$$f(x, w, b) = \text{sign}(w^T x + b)$$

$$x = (x_1, x_2)$$

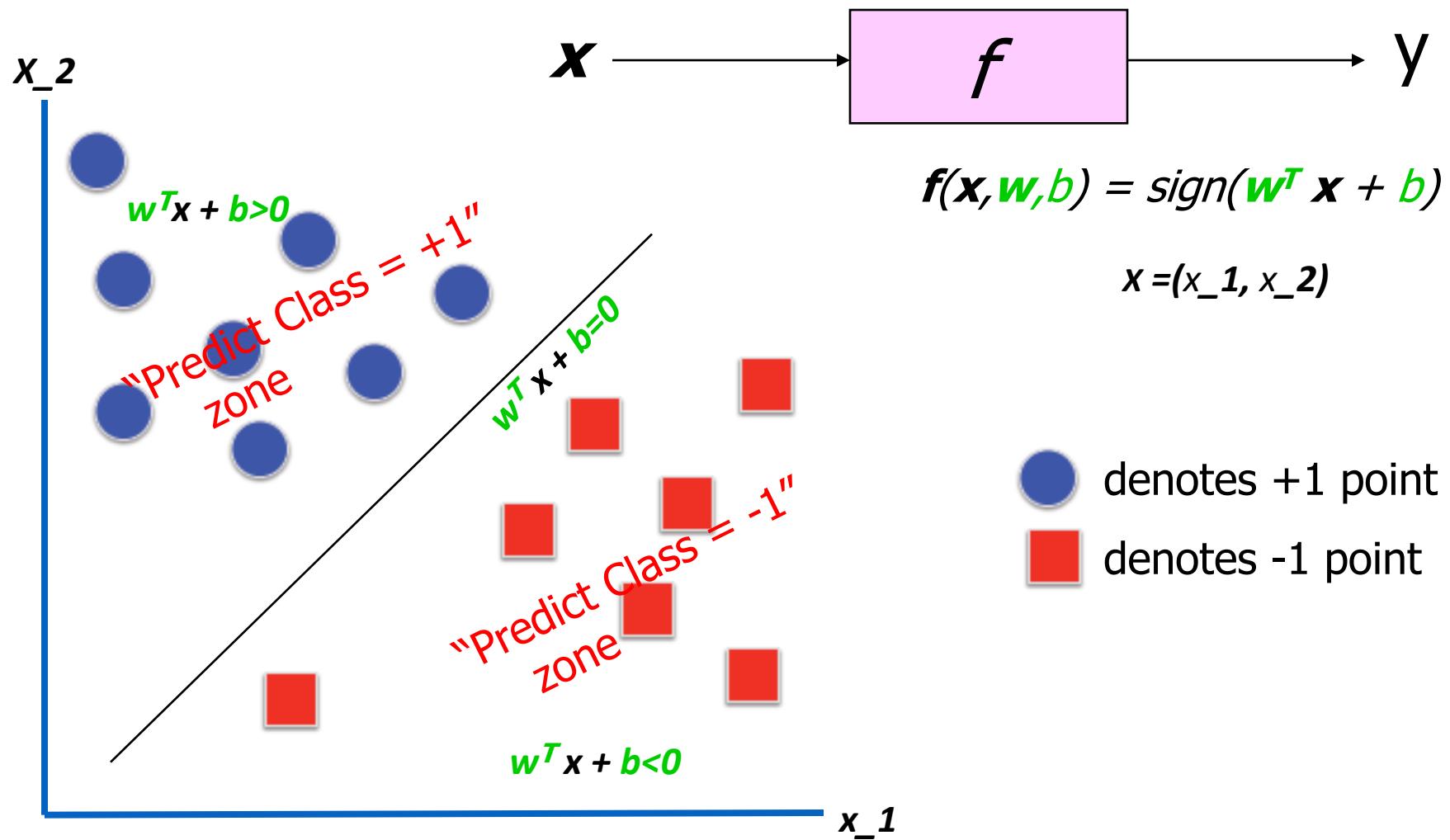
Traditional Programming



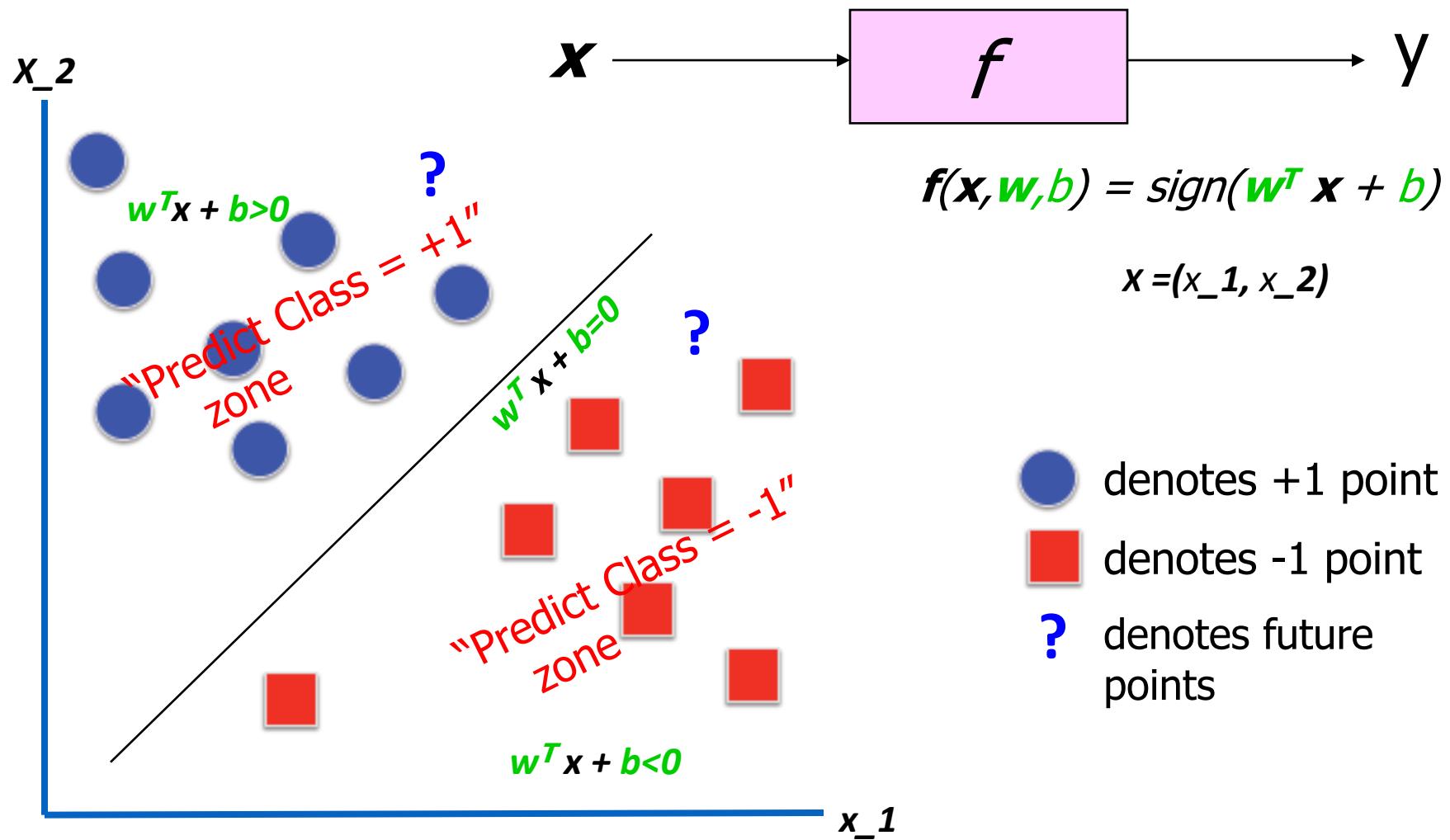
Machine Learning (training phase)



SUPERVISED Linear Binary Classifier



SUPERVISED Linear Binary Classifier



Basic Concepts

- Training (i.e. learning parameters \mathbf{w}, b)
 - Training set includes
 - available examples $\mathbf{x}_1, \dots, \mathbf{x}_L$
 - available corresponding labels y_1, \dots, y_L
 - Find (\mathbf{w}, b) by minimizing loss
 - (i.e. difference between y and $f(\mathbf{x})$ on available examples *in training set*)

$$(\mathbf{w}, b) = \underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{i=1}^L \ell(f(x_i), y_i)$$

Basic Concepts

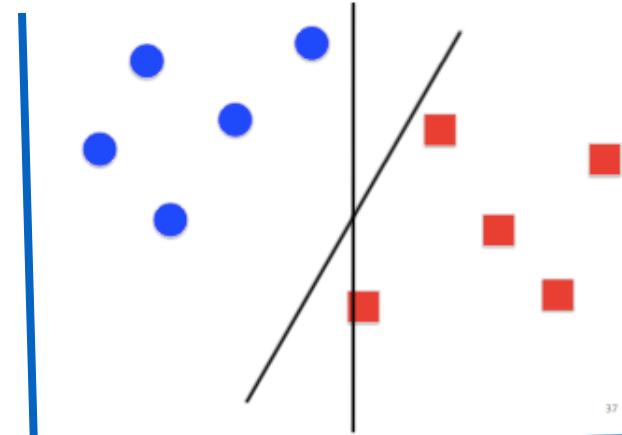
- Testing (i.e. evaluating performance on “future” points)
 - Difference between true $y_?$ and the predicted $f(\mathbf{x}_?)$ on a set of testing examples (i.e. *testing set*)
 - Key: example $\mathbf{x}_?$ not in the training set
- Generalisation: learn function / hypothesis from past data in order to “explain”, “predict”, “model” or “control” new data examples

Basic Concepts

- Loss function

- e.g. hinge loss for binary classification

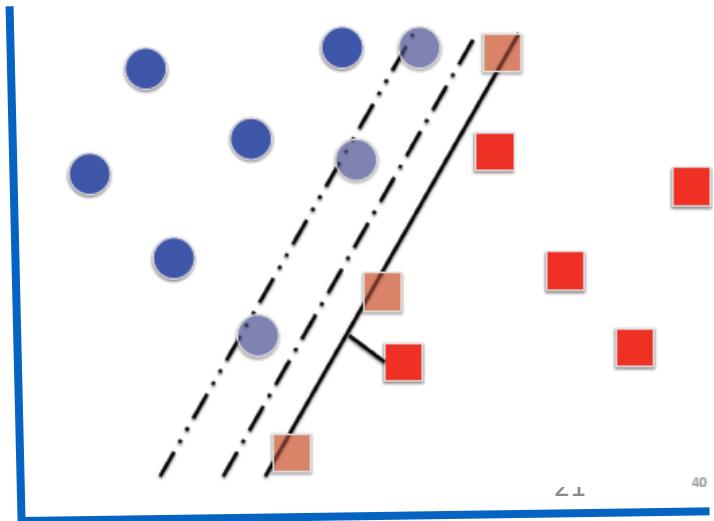
$$\sum_{i=1}^L \ell(f(x_i), y_i) = \sum_{i=1}^L \max(0, 1 - y_i f(x_i)).$$



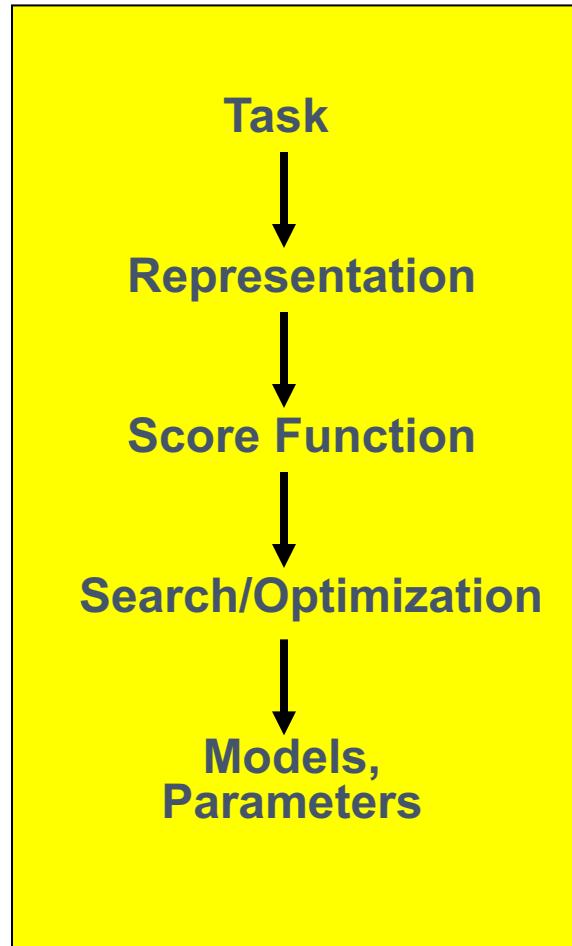
- e.g. pairwise ranking loss for ranking task (i.e. ordering examples by preference)

- Regularization

- E.g. $C \sum_{i=1}^L \ell(f(x_i), y_i) + \frac{1}{2} \|w\|^2$, an added term on loss function to control f



Machine Learning in a Nutshell



ML grew out of work in AI

Optimize a performance criterion using example data or past experience,

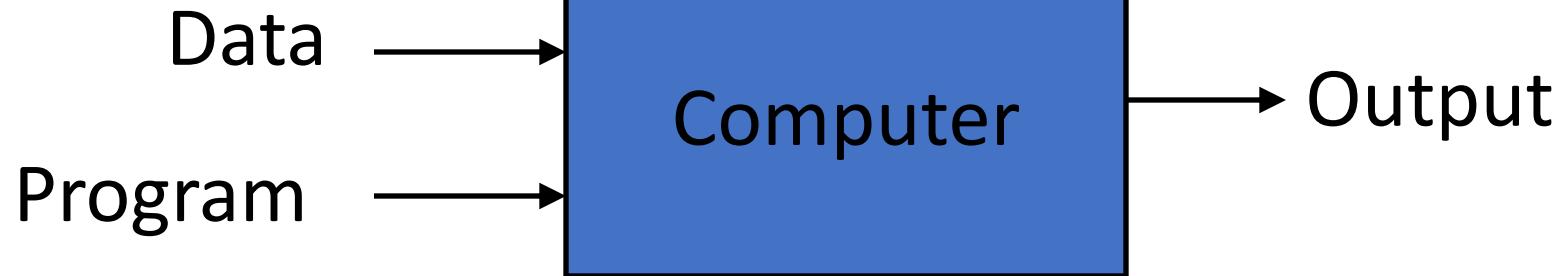
Aiming to generalize to unseen data

Today

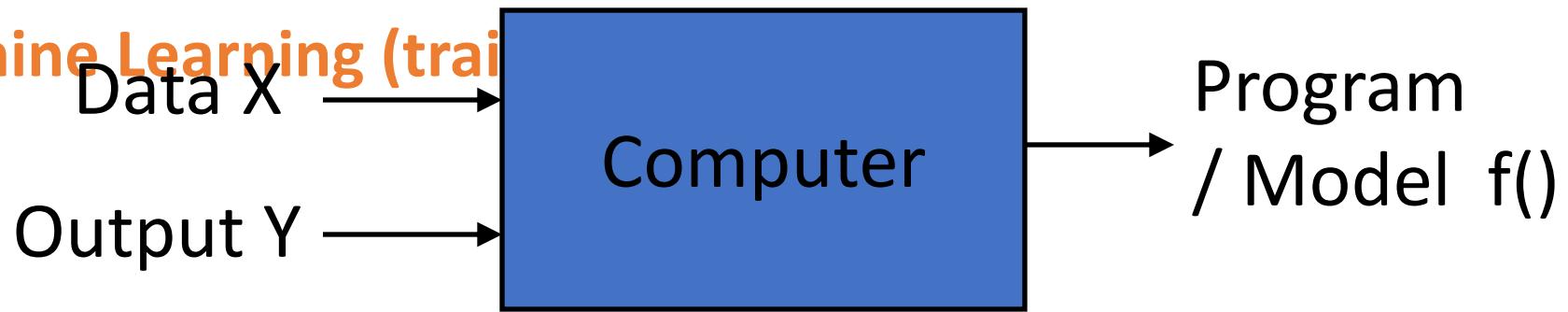
- Course Logistics
- Machine Learning Basics
- A Rough Plan of Course Content
- Machine Learning History



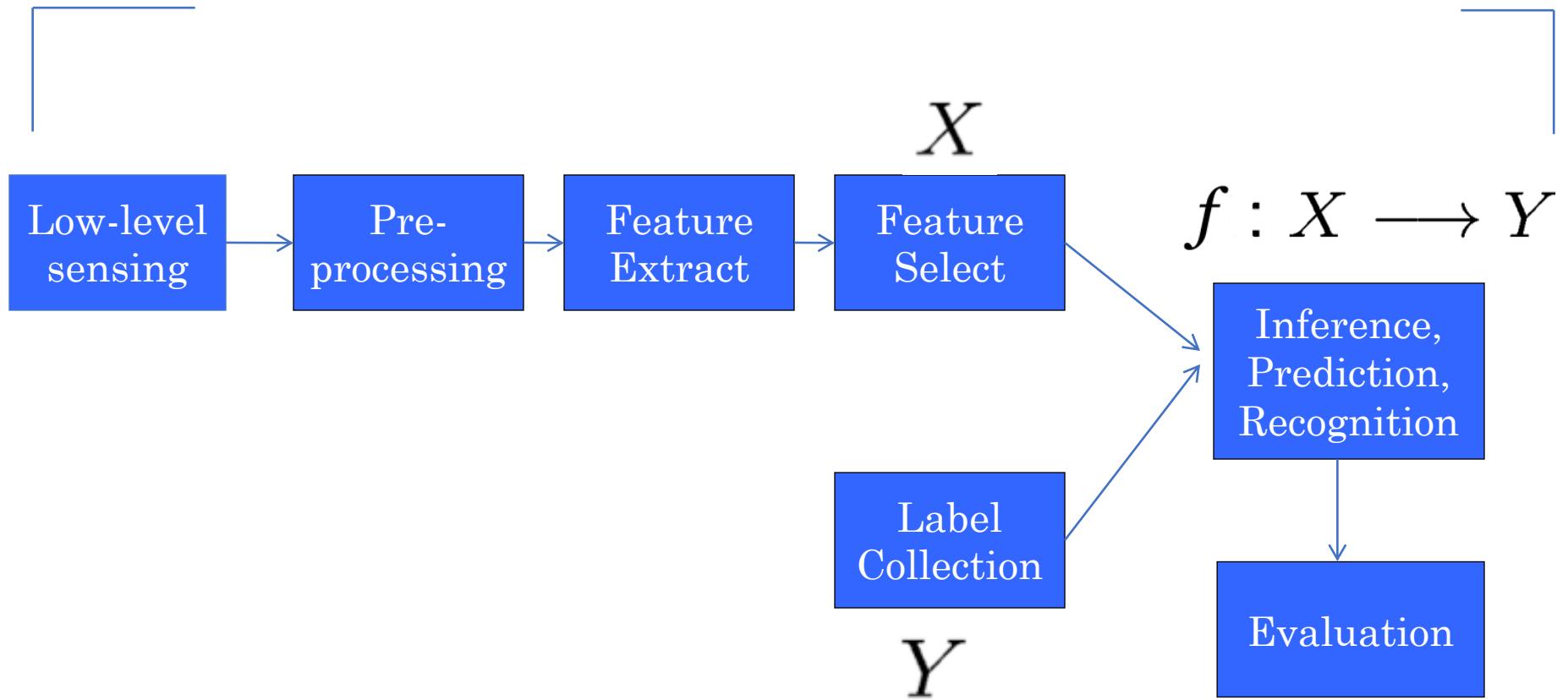
Traditional Programming



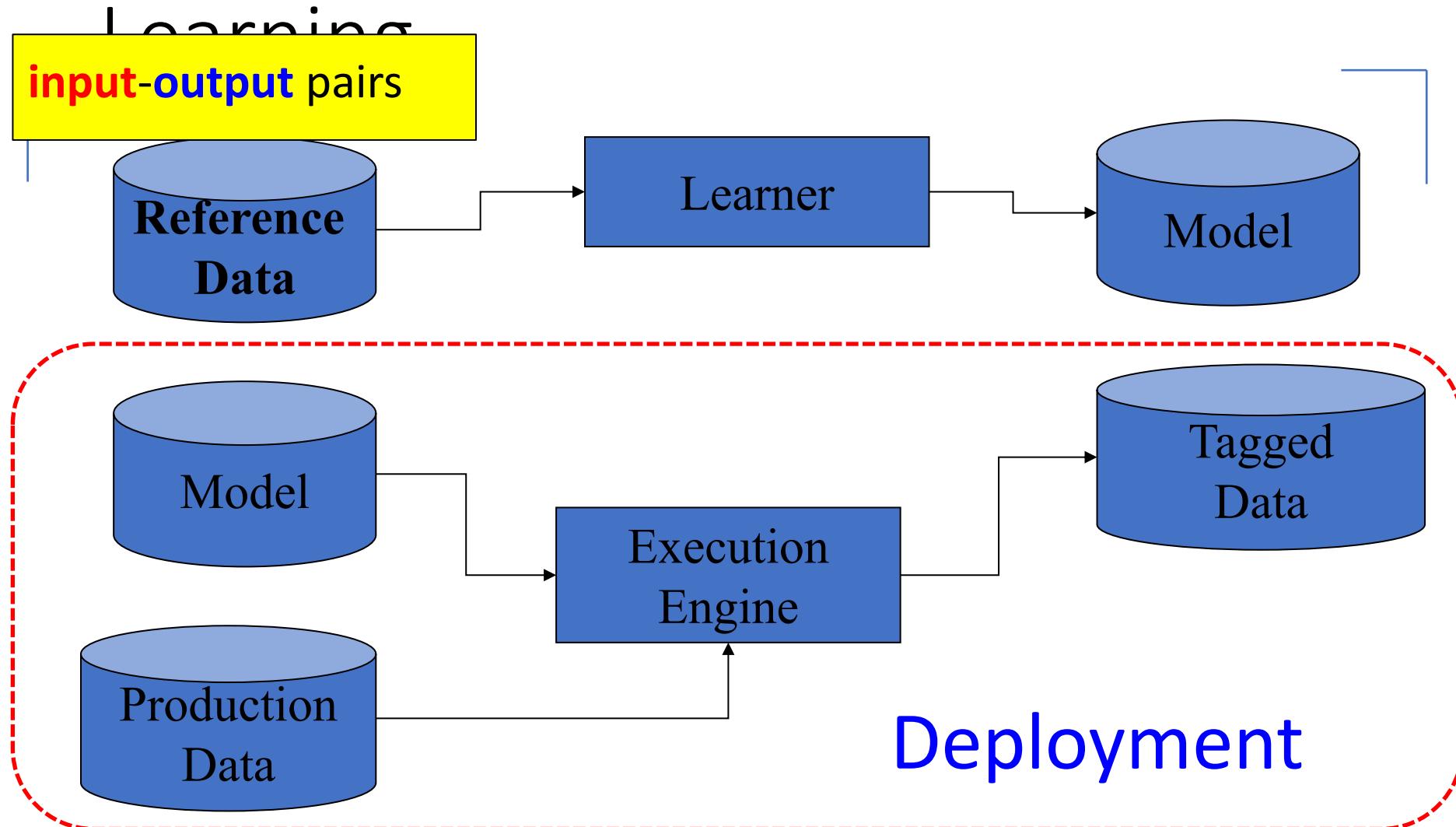
Machine Learning (train)



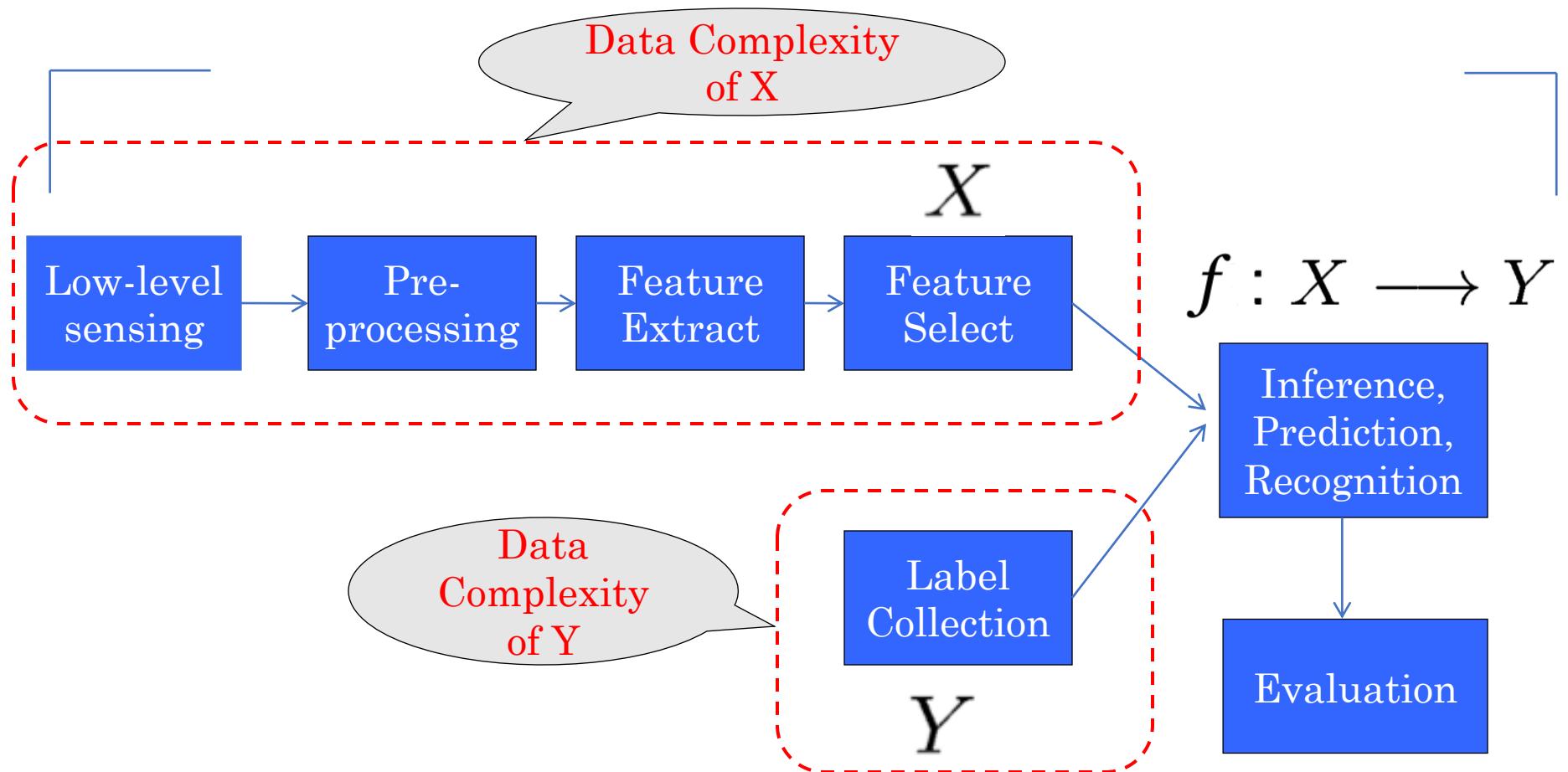
TYPICAL MACHINE LEARNING SYSTEM (Training Mode)



An Operational Model of Machine Learning

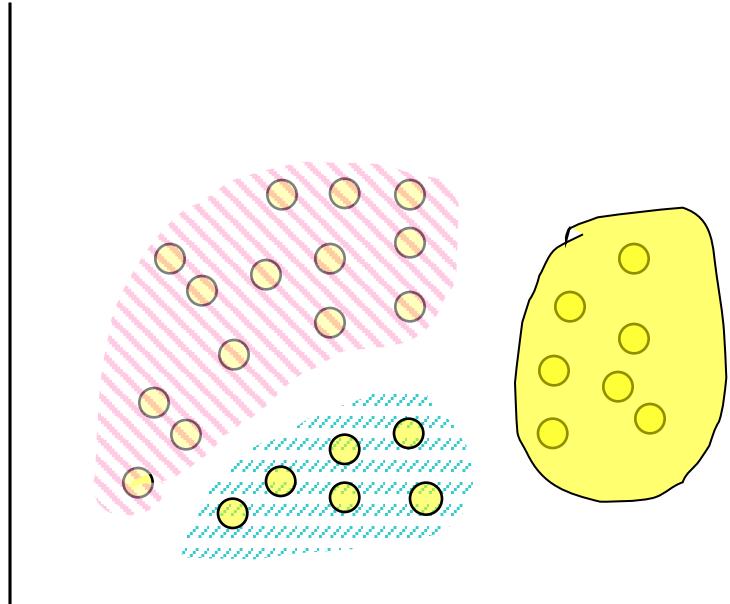


TYPICAL MACHINE LEARNING SYSTEM



UNSUPERVISED LEARNING : COMPLEXITY in Y]

- No labels are provided (e.g. No Y provided)
- Find patterns from unlabeled data, e.g. clustering



e.g. clustering => to find
“natural” grouping of
instances given un-labeled
data

STRUCTURAL OUTPUT LEARNING : [COMPLEXITY OF Y]

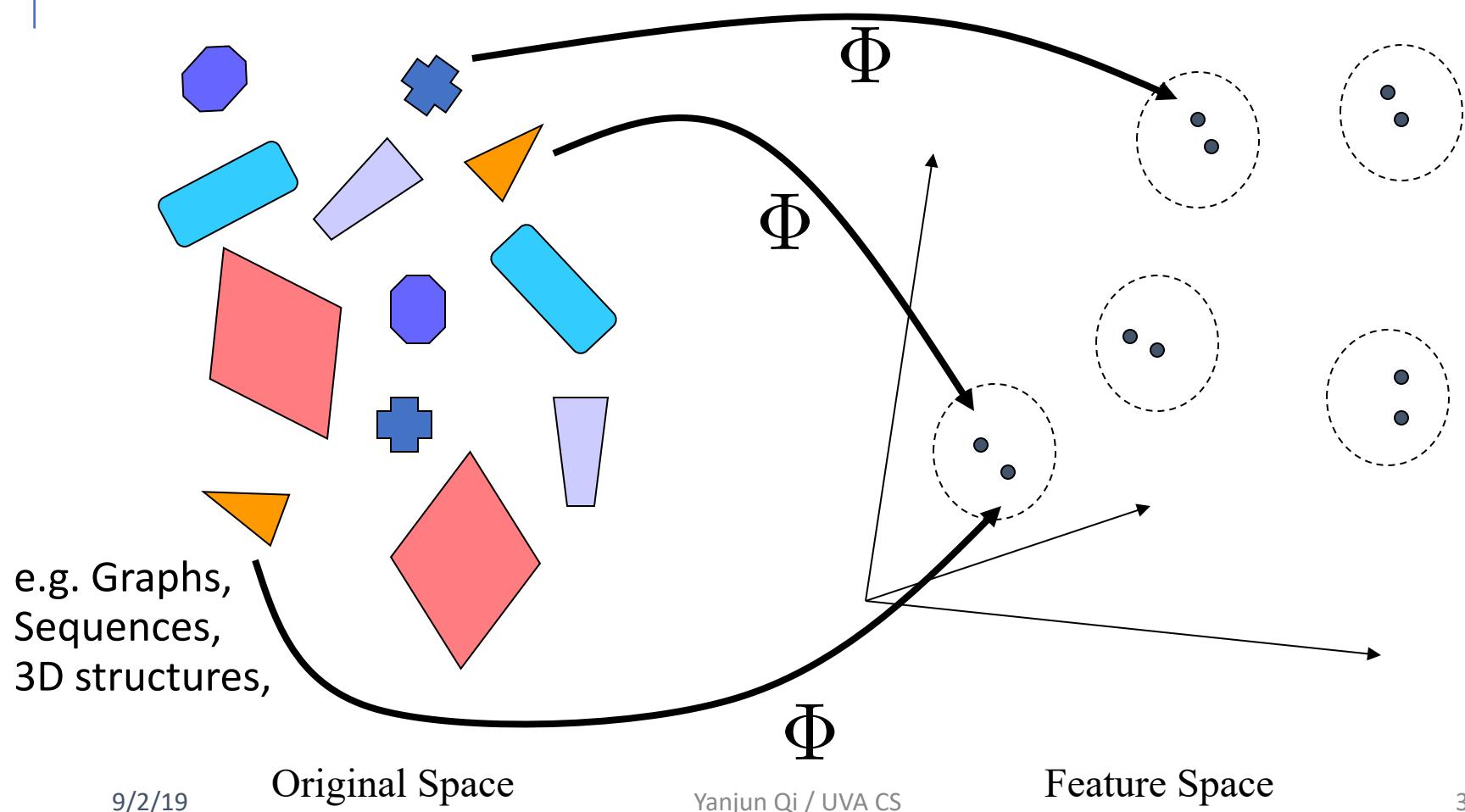
- Many prediction tasks involve **output labels having structured correlations or constraints among instances**

Structured Dependency between Examples' Y	Sequence	Tree	Grid
Input X	APAFSVSPASGACCGPECA...	The dog chased the cat	
Output Y	 CCEEEEEECCCCCCCCHHHHCCC...	<pre>graph TD; S --> NP1[NP]; S --> VP[VP]; NP1 --> Det1[Det]; NP1 --> N1[N]; VP --> V[V]; VP --> NP2[NP]; NP2 --> Det2[Det]; NP2 --> N2[N];</pre>	

Many more possible structures between y_i , e.g. **spatial** , **temporal**, **relational** ...

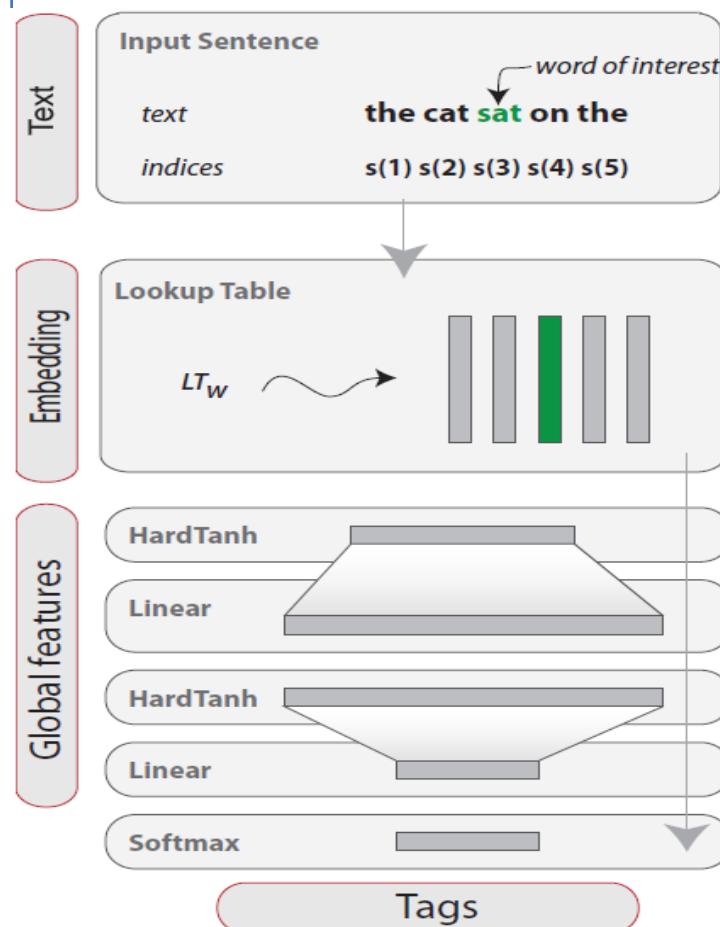
STRUCTURAL INPUT : Kernel Methods

[COMPLEXITY OF X]

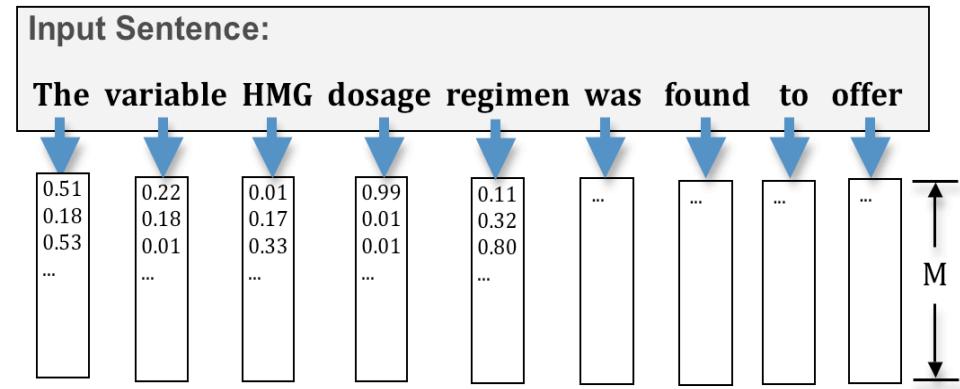


MORE RECENT: FEATURE LEARNING [COMPLEXITY OF X]

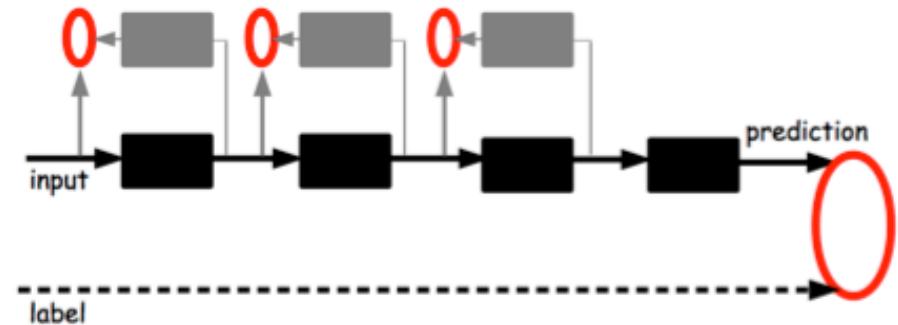
Deep Learning



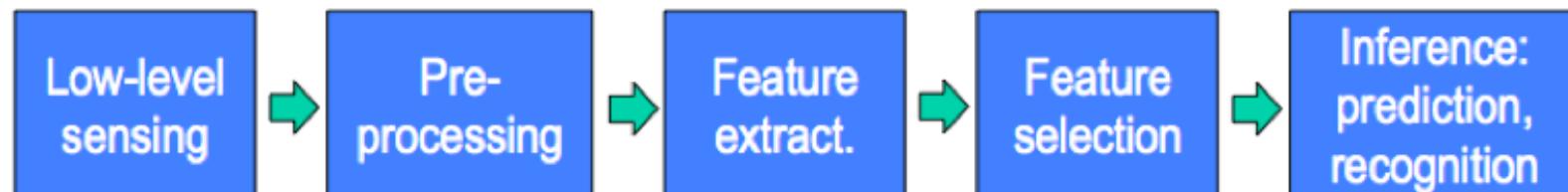
Supervised Embedding



Layer-wise Pretraining



DEEP LEARNING / FEATURE LEARNING : [COMPLEXITY OF X]



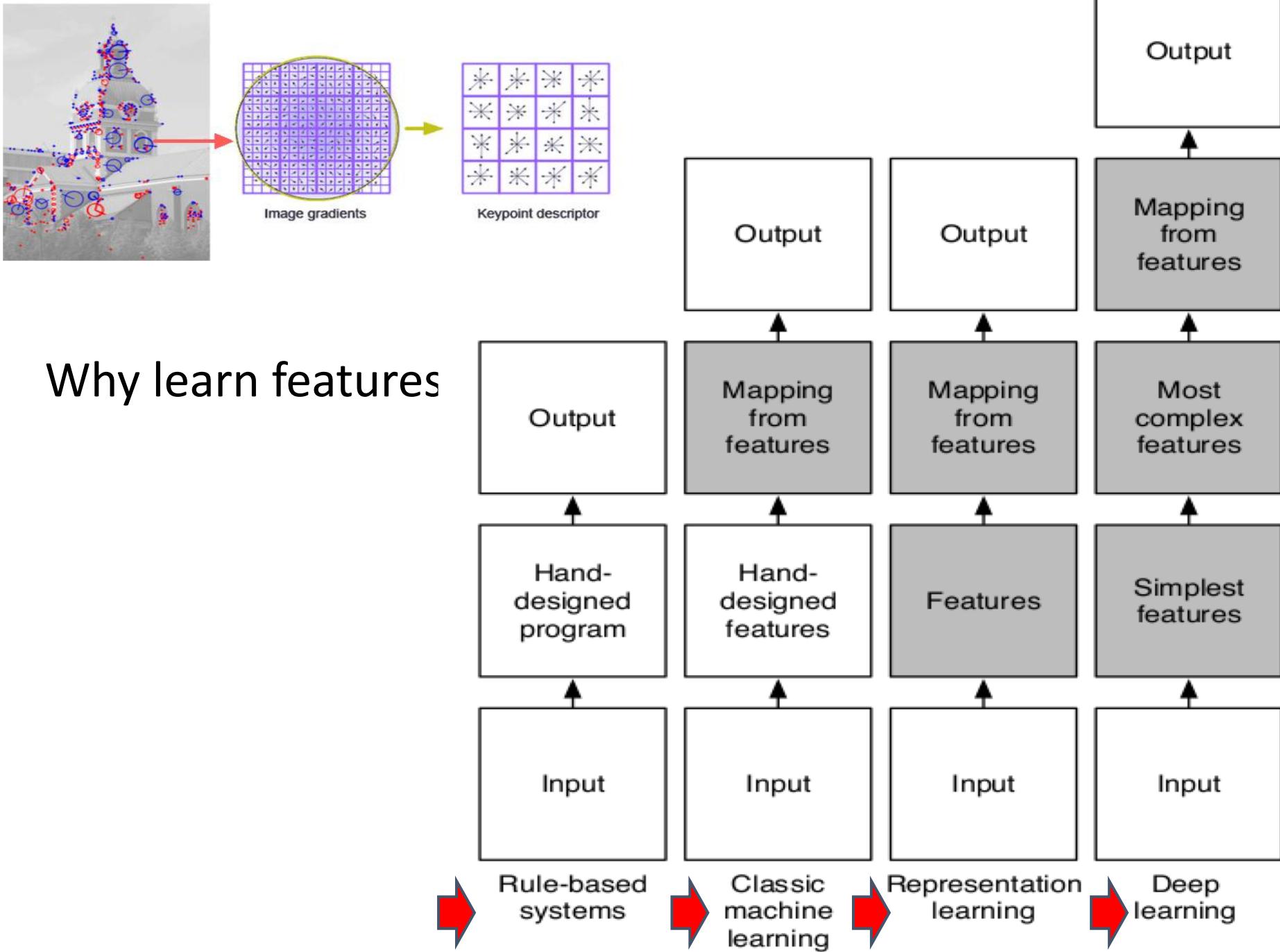
Feature Engineering

- ✓ Most critical for accuracy
- ✓ Account for most of the computation for testing
- ✓ Most time-consuming in development cycle
- ✓ Often hand-craft and task dependent in practice



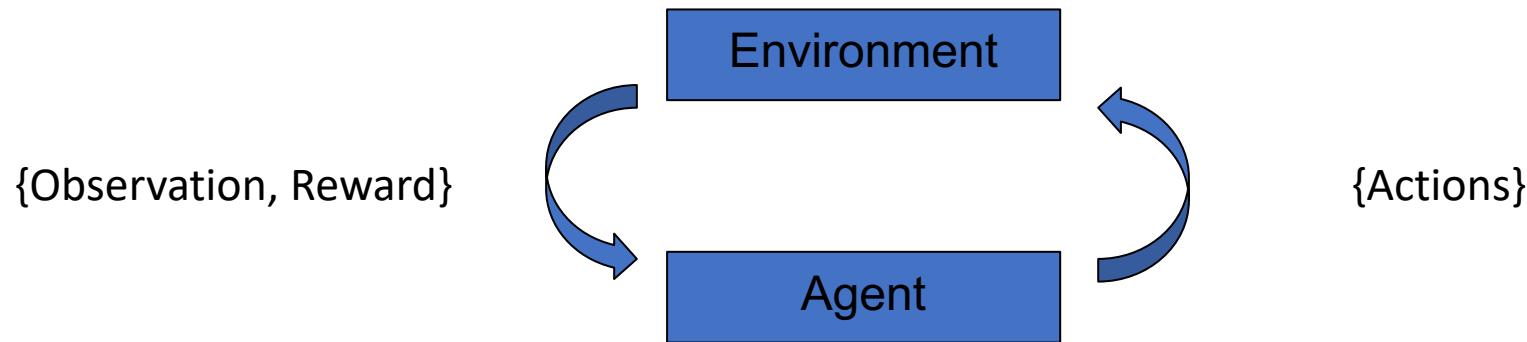
Feature Learning

- ✓ Easily adaptable to new similar tasks
- ✓ Layerwise representation
- ✓ Layer-by-layer unsupervised training
- ✓ Layer-by-layer supervised training



Reinforcement Learning (RL)

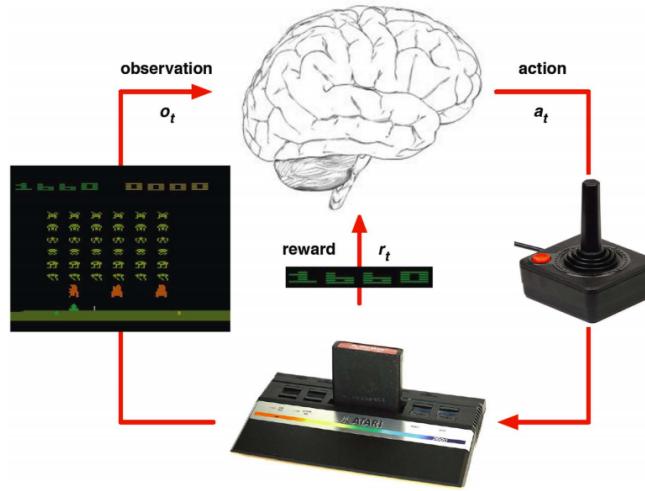
- What's Reinforcement Learning?



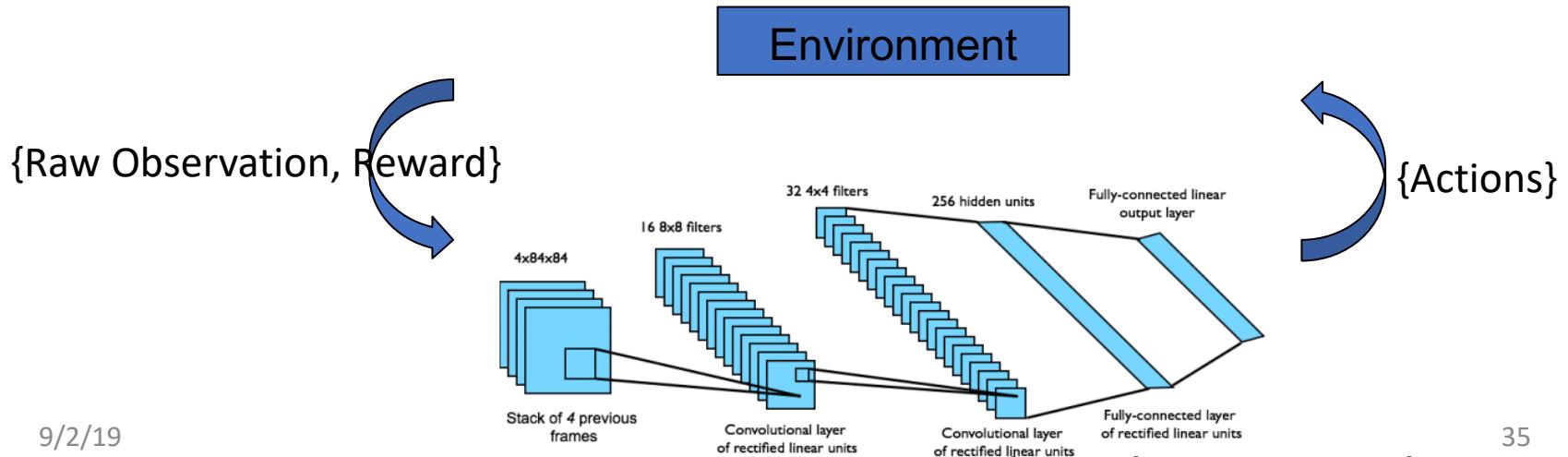
- Agent interacts with an environment and learns by maximizing a scalar reward signal
- No labels or any other supervision signal.
- Previously suffering from hand-craft states or representation.

Deep Reinforcement Learning

- Human



- So what's **DEEP RL**?



When to use Machine Learning (Adapt to / learn from data) ?

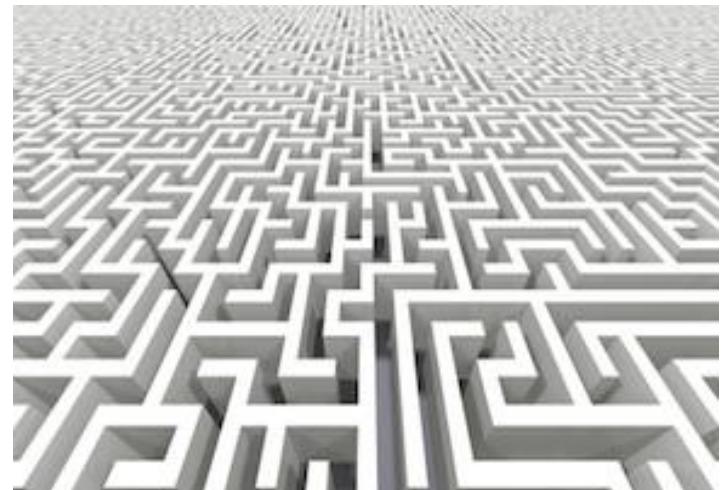
- 1. Extract knowledge from data
 - Relationships and correlations can be hidden within large amounts of data
 - The amount of knowledge available about certain tasks is simply too large for explicit encoding (e.g. rules) by humans
- 2. Learn tasks that are difficult to formalise
 - Hard to be defined well, except by examples, e.g., face recognition
- 3. Create software that improves over time
 - New knowledge is constantly being discovered.
 - Rule or human encoding-based system is difficult to continuously re-design “by hand”.

“Big Data” Challenges for Machine Learning

LARGE-SCALE



HIGH-COMPLEXITY



- || ✓ Large size of samples
- || ✓ High dimensional features

Not the focus,
being covered in
my advanced-
level course

Large-Scale Machine Learning: SIZE MATTERS

LARGE-SCALE



- One thousand data instances
- One million data instances
- One billion data instances
- One trillion data instances

Those are not different numbers,
those **are different mindsets !!!**

BIG DATA CHALLENGES FOR MACHINE LEARNING

LARGE-SCALE



Highly Complex



Most of
this
course

The variations of both **X**
(feature, representation)
and **Y** (labels) are complex
!

- ✓ Complexity of X
- ✓ Complexity of Y

Course Content Plan →

Six major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
-
- Graphical models
- Reinforcement Learning

Y is a continuous

Y is a discrete

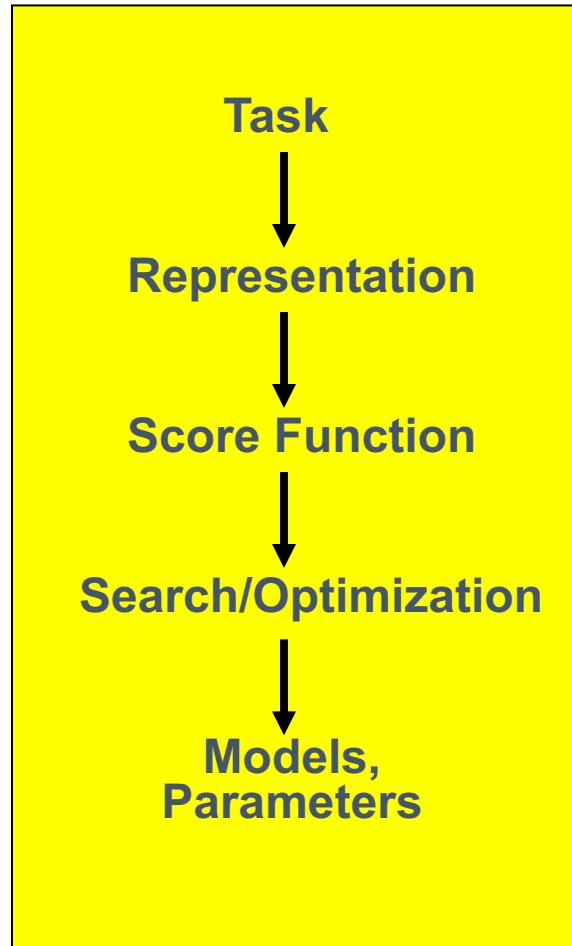
NO Y

About $f()$

About interactions among X_1, \dots, X_p

About interactions among X_1, \dots, X_p

Machine Learning in a Nutshell



ML grew out of work in AI

Optimize a performance criterion using example data or past experience,

Aiming to generalize to unseen data

What we have covered

Task	
Representation	
Score Function	
Search/Optimization	
Models, Parameters	

What we will cover

Task	Regression, classification, clustering, dimen-reduction
Representation	Linear func, nonlinear function (e.g. polynomial expansion), local linear, logistic function (e.g. $p(c x)$), tree, multi-layer, prob-density family (e.g. Bernoulli, multinomial, Gaussian, mixture of Gaussians), local func smoothness, kernel matrix, local smoothness, partition of feature space,
Score Function	MSE, Margin, log-likelihood, EPE (e.g. L2 loss for KNN, 0-1 loss for Bayes classifier), cross-entropy, cluster points distance to centers, variance, conditional log-likelihood, complete data-likelihood, regularized loss func (e.g. L1, L2) , goodness of inter-cluster similar
Search/ Optimization	Normal equation, gradient descent, stochastic GD, Newton, Linear programming, Quadratic programming (quadratic objective with linear constraints), greedy, EM, asyn-SGD, eigenDecomp, backprop
Models, Parameters	Linear weight vector, basis weight vector, local weight vector, dual weights, training samples, tree-dendrogram, multi-layer weights, principle components, member (soft/hard) assignment, cluster centroid, cluster covariance (shape), ...

	X_1	X_2	X_3	Y
s_1				
s_2				
s_3				
s_4				
s_5				
s_6				

$$f : \boxed{X} \longrightarrow \boxed{Y}$$

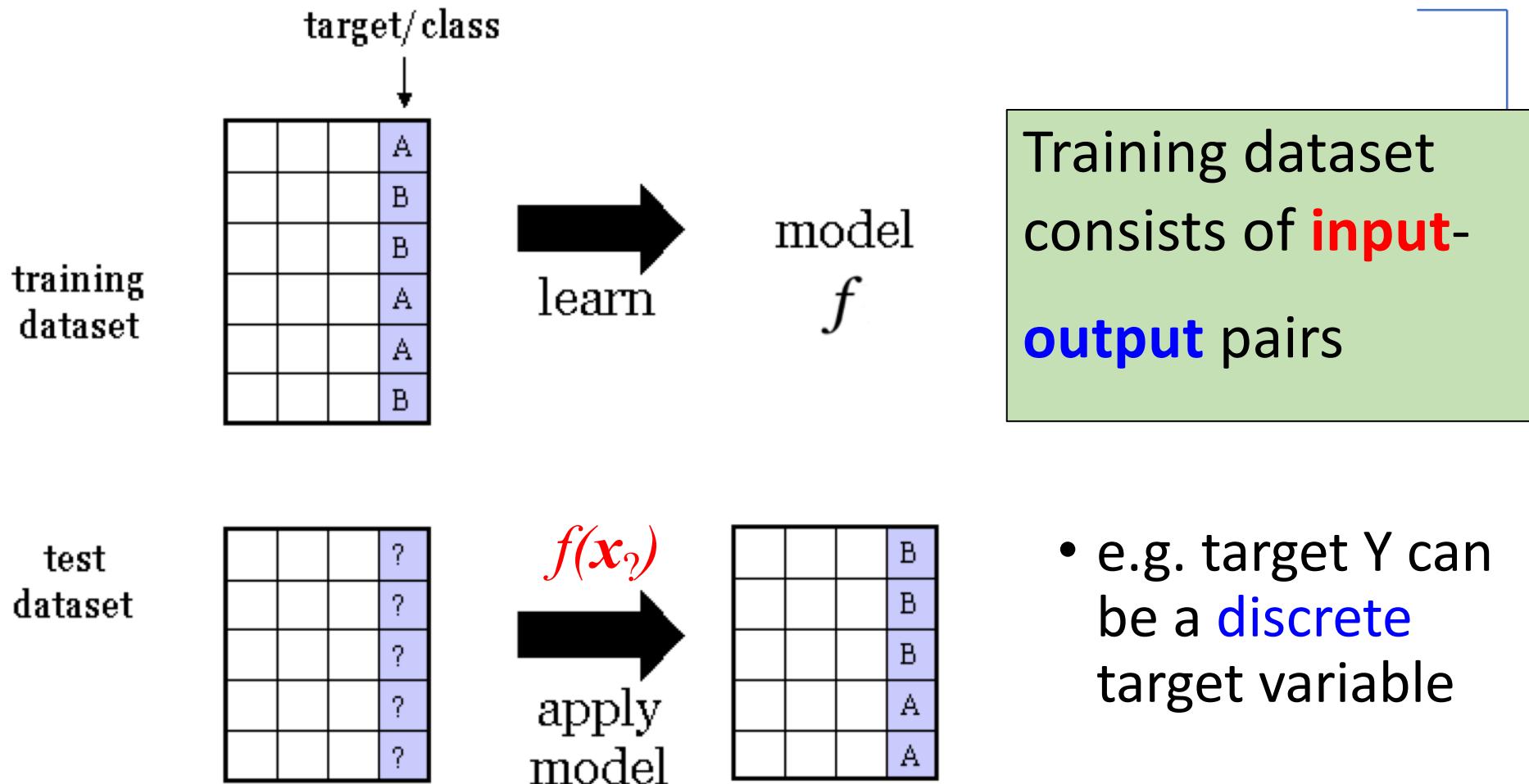
- **Data/points/instances/examples/samples/records:** [**rows**]
- **Features/attributes/dimensions/independent variables/covariates/predictors/regressors:** [**columns, except the last**]
- **Target/outcome/response/label/dependent variable:** special column to be predicted [**last column**]

Main Types of Columns

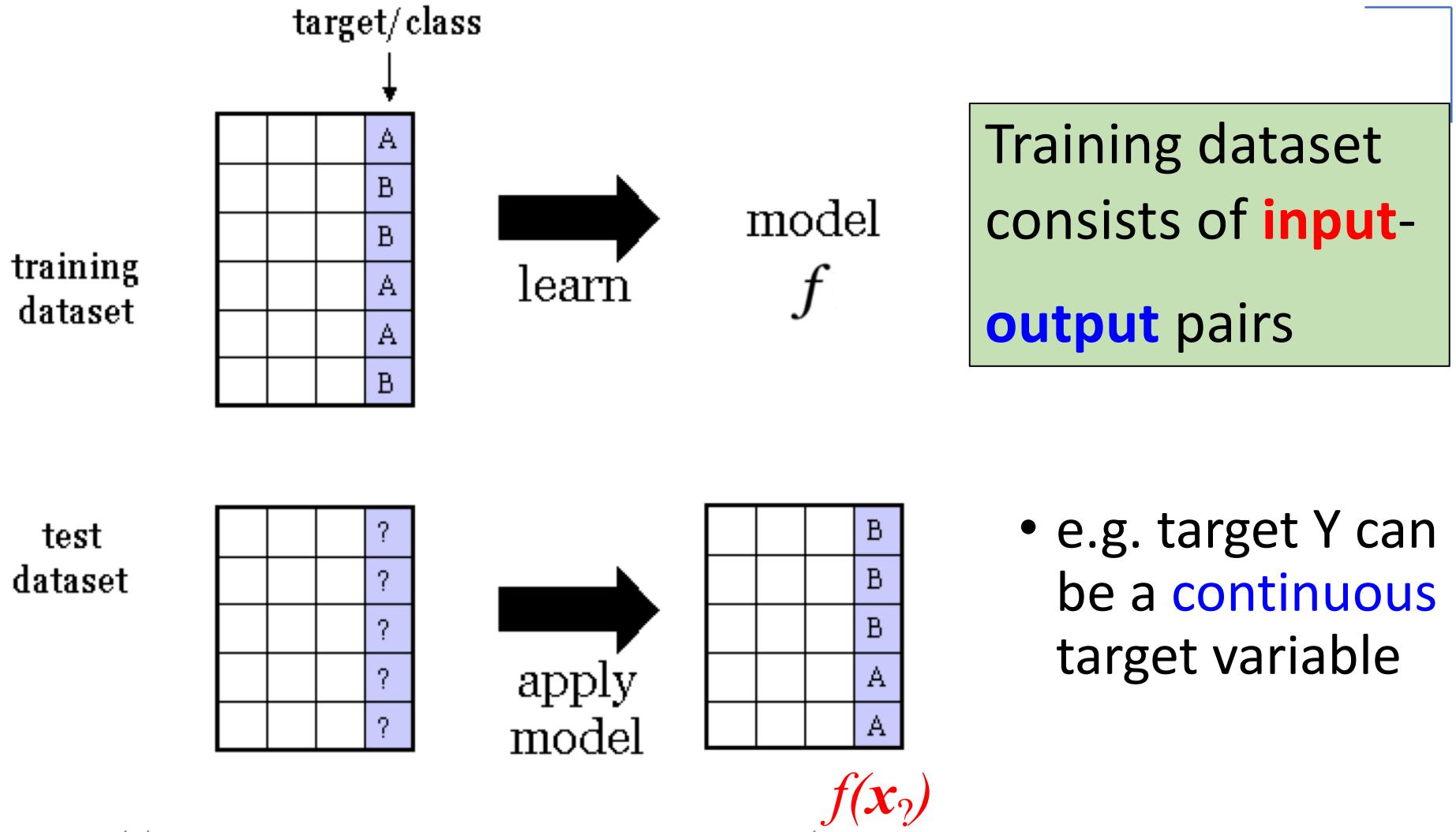
	X_1	X_2	X_3	Y
S_1				
S_2				
S_3				
S_4				
S_5				
S_6				

- *Continuous*: a real number, for example, weight
- *Discrete*: a symbol, like “Good” or “Bad”

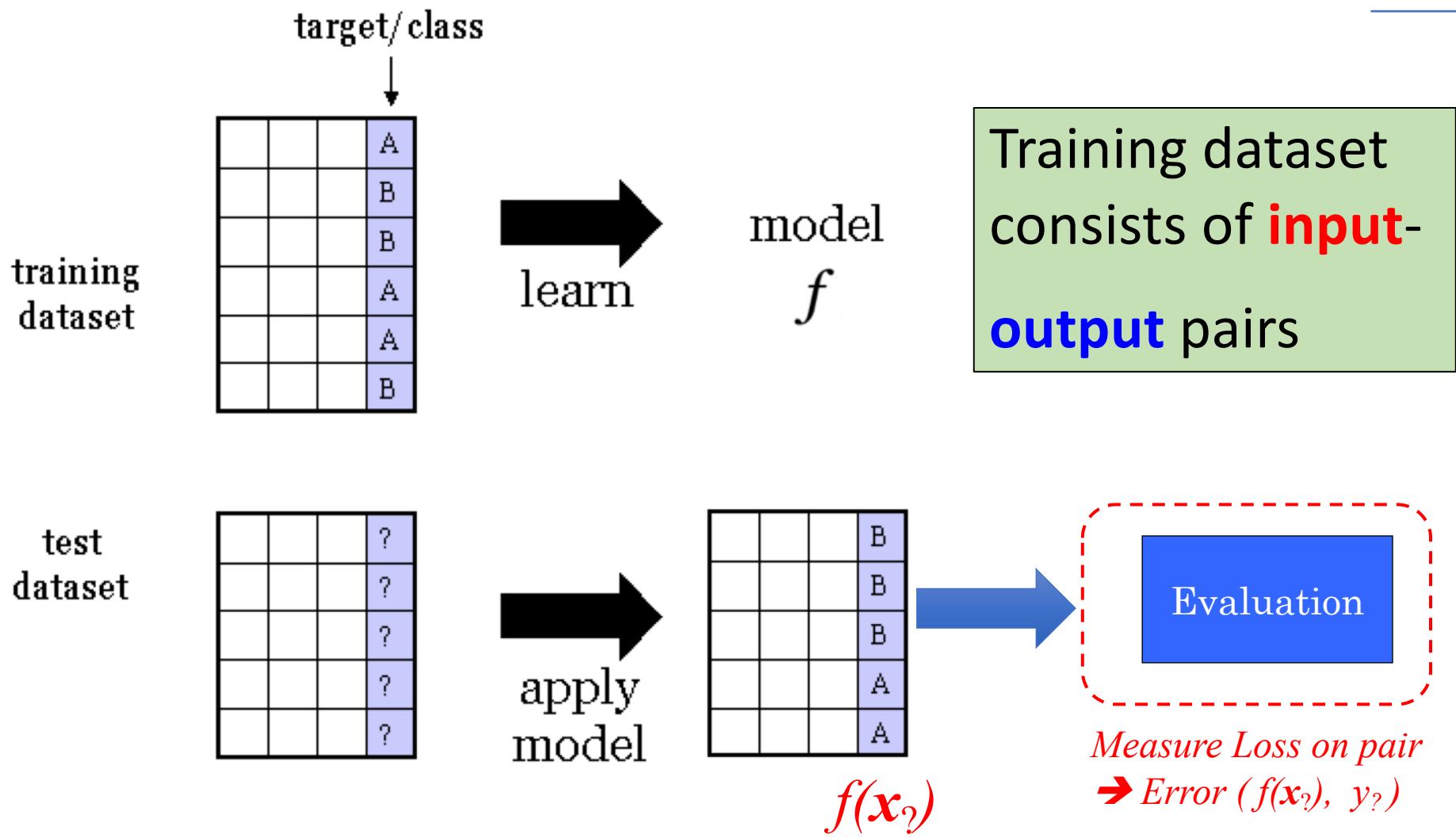
e.g. SUPERVISED Classification

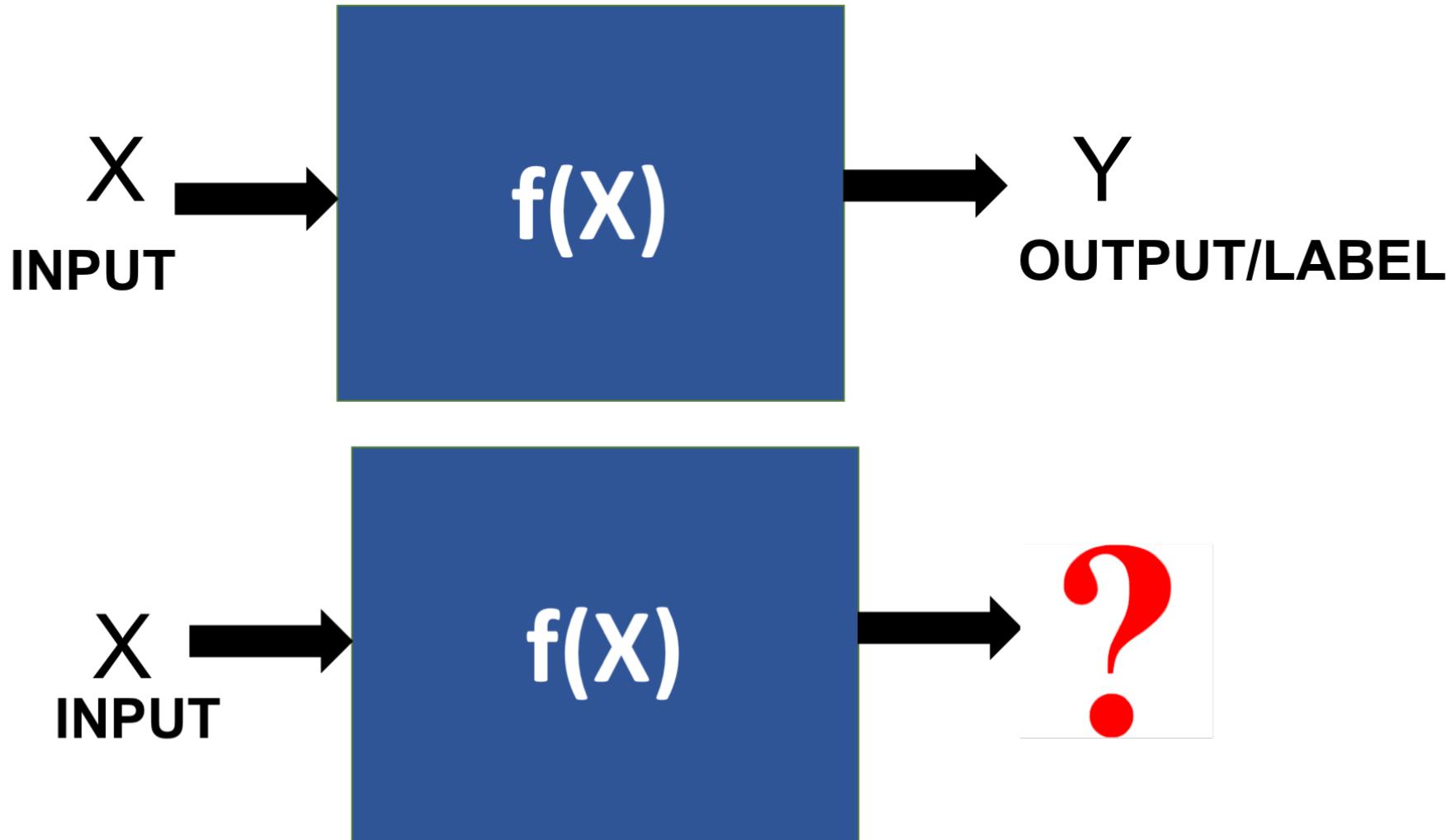


e.g. SUPERVISED Regression



How to know the program works well?





Testing

training dataset

$$\mathbf{X}_{train} = \begin{bmatrix} \cdots & \mathbf{x}_1^T & \cdots \\ \cdots & \mathbf{x}_2^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_n^T & \cdots \end{bmatrix} \quad \vec{\mathbf{y}}_{train} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

test dataset

$$\mathbf{X}_{test} = \begin{bmatrix} \cdots & \mathbf{x}_{n+1}^T & \cdots \\ \cdots & \mathbf{x}_{n+2}^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_{n+m}^T & \cdots \end{bmatrix} \quad \vec{\mathbf{y}}_{test} = \begin{bmatrix} y_{n+1} \\ y_{n+2} \\ \vdots \\ y_{n+m} \end{bmatrix}$$

Notation

- Inputs
 - X, X_j (jth element of vector X) : random variables written in capital letter
 - p #input variables, n #observations
 - \mathbf{X} : matrix written in bold capital
 - Vectors are assumed to be column vectors
 - Discrete inputs often described by characteristic vector (dummy variables)

- Outputs
 - quantitative Y
 - qualitative C (for categorical)
- Observed variables written in lower case
 - The i-th observed value of X is x_i and can be a scalar or a vector



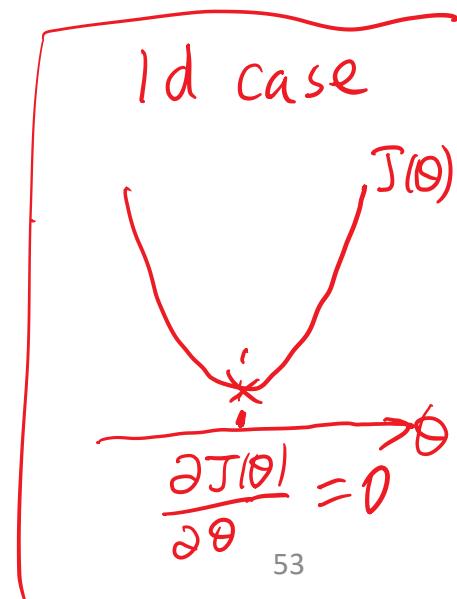
a few sample slides

$$\begin{aligned}
 J(\theta) &= (\underline{x}\theta - y)^T (\underline{x}\theta - y) \frac{1}{2} \\
 &= ((x\theta)^T - y^T)(\underline{x}\theta - y) \frac{1}{2} \\
 &= (\theta^T \underline{x}^T - y^T)(\underline{x}\theta - y) \frac{1}{2} \\
 &= (\underbrace{\theta^T \underline{x}^T \underline{x}\theta - \theta^T \underline{x}^T y - y^T \underline{x}\theta + y^T y}_{\text{since } \theta^T \underline{x}^T y = y^T \underline{x}\theta}) \frac{1}{2}.
 \end{aligned}$$

since $\theta^T \underline{x}^T y = y^T \underline{x}\theta$
 $\langle x\theta, y \rangle \quad \langle y, \underline{x}\theta \rangle$

$$= (\underbrace{\theta^T \underline{x}^T \underline{x}\theta}_{\text{1d case}} - \underbrace{2\theta^T \underline{x}^T y}_{\text{J}(\theta)} + \underbrace{y^T y}_{\frac{\partial J(\theta)}{\partial \theta} = 0}) \frac{1}{2}$$

$\Rightarrow J(\theta)$ quadratic func of θ ;



Deriving the Maximum Likelihood Estimate for Bernoulli

maximize

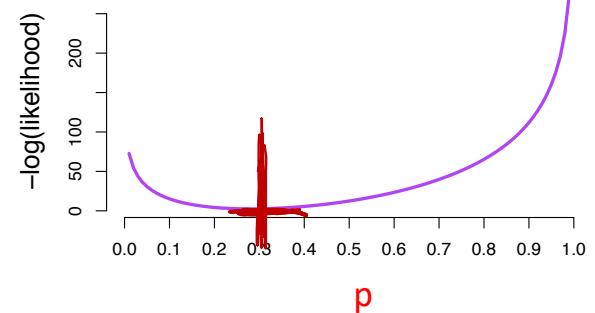
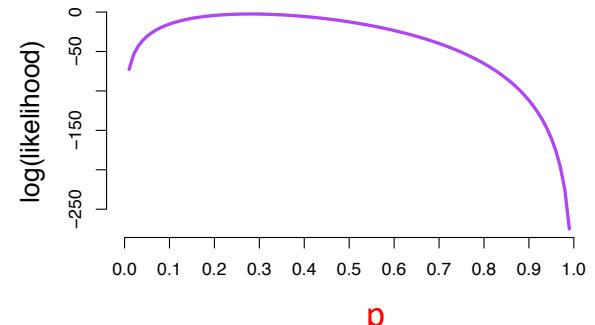
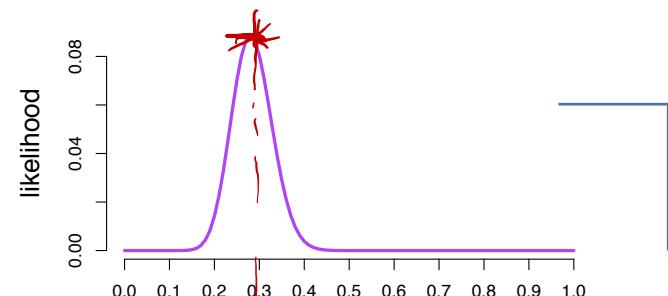
$$L(p) = p^x (1-p)^{n-x}$$

maximize

$$\log(L(p)) = \log[p^x (1-p)^{n-x}]$$

Minimize the negative log-likelihood

$$-l(p) = -\log[p^x (1-p)^{n-x}]$$

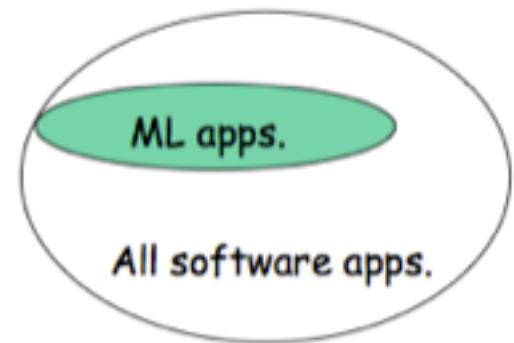


Today

- Course Logistics
- Machine Learning Basics
- A Rough Plan of Course Content
-  Machine Learning History

MACHINE LEARNING IN COMPUTER SCIENCE

- Machine learning is already the preferred approach for
 - Speech recognition, natural language processing
 - Computer vision
 - Medical outcome analysis
 - Robot control ...
- Why growing ?
 - Improved machine learning algorithms
 - Improved CPU / GPU powers
 - Increased data capture, new sensors, networking
 - Systems/Software too complex to control manually
 - Demand to self-customization for user, environment,



HISTORY OF MACHINE LEARNING

- 1950s
 - Samuel's checker player
 - Selfridge's Pandemonium
- 1960s:
 - Neural networks: Perceptron
 - Pattern recognition
 - Learning in the limit theory
 - Minsky and Papert prove limitations of Perceptron
- 1970s:
 - Symbolic concept induction
 - Winston's arch learner
 - Expert systems and the knowledge acquisition bottleneck
 - Quinlan's DT ID3
 - Michalski's AQ and soybean diagnosis
 - Scientific discovery with BACON
 - Mathematical discovery with AM

HISTORY OF MACHINE LEARNING (CONT.)

- 1980s:
 - Advanced decision tree and rule learning
 - Explanation-based Learning (EBL)
 - Learning and planning and problem solving
 - Utility problem
 - Analogy
 - Cognitive architectures
 - Resurgence of neural networks (connectionism, backpropagation)
 - Valiant's **PAC Learning** Theory
 - Focus on experimental methodology
- 1990s
 - **Data mining**
 - Adaptive software agents and web applications
 - Text learning
 - Reinforcement learning (RL)
 - Inductive Logic Programming (ILP)
 - Ensembles: Bagging, Boosting, and Stacking
 - Bayes Net learning

HISTORY OF MACHINE LEARNING (CONT.)

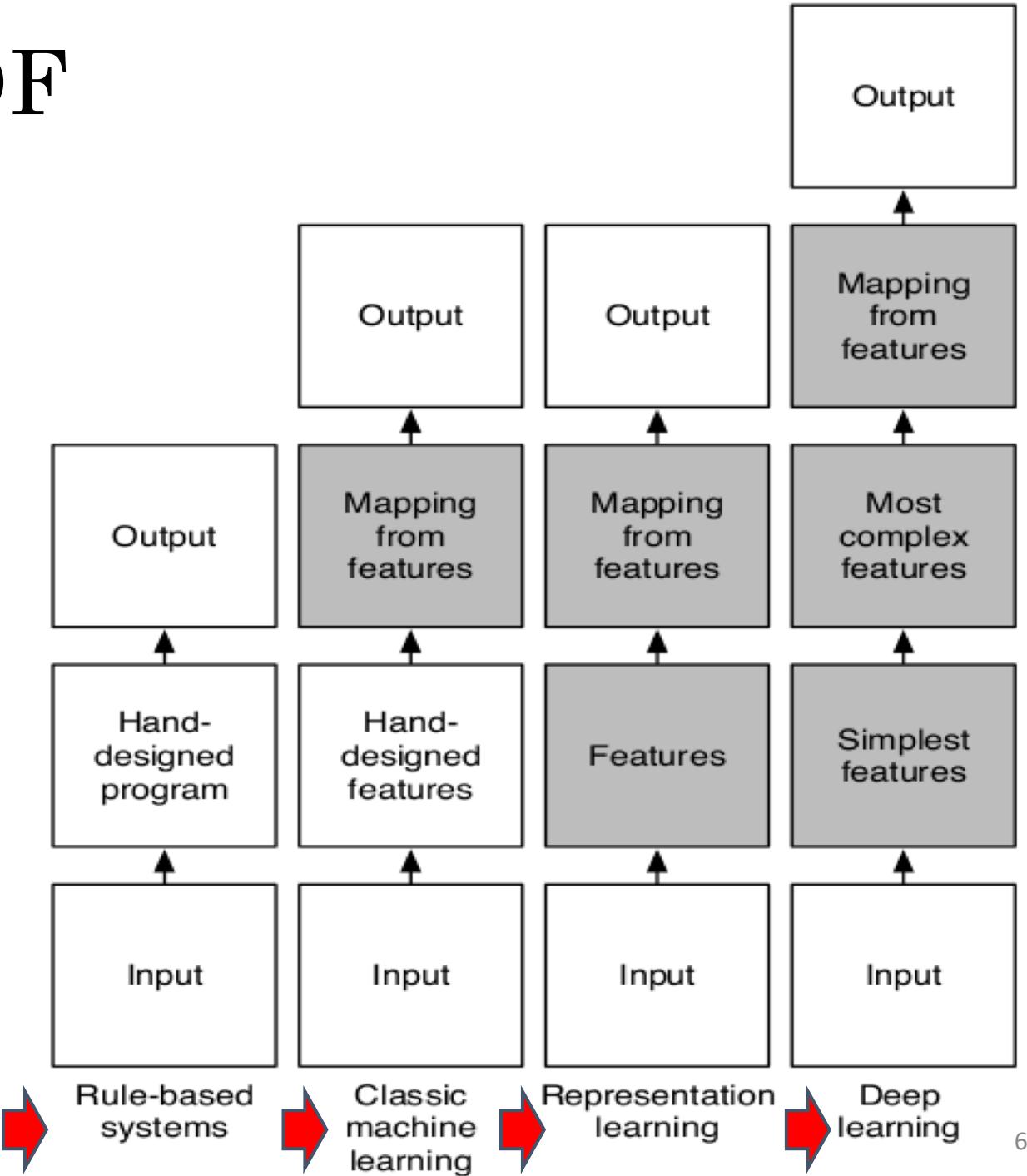
- 2000s

- Support vector machines
- Kernel methods
- Graphical models
- Statistical relational learning
- Transfer learning
- Sequence labeling
- Collective classification and structured outputs
- Computer Systems Applications
 - Compilers
 - Debugging
 - Graphics
 - Security (intrusion, virus, and worm detection)
- Email management
- Personalized assistants that learn
- Learning in robotics and vision

HISTORY OF MACHINE LEARNING (CONT.)

- 2010s
 - Speech translation, voice recognition (e.g. SIRI)
 - Google search engine uses numerous machine learning techniques (e.g. grouping news, spelling corrector, improving search ranking, image retrieval,)
 - 23 and me (scan sample of person genome, predict likelihood of genetic disease, ...)
 - DeepMind, Google Brain, ...
 - IBM Watson QA system
 - Machine Learning as a service (e.g. google prediction API, bigml.com, cloud autoML .)
 - IBM healthcare analytics
 -

HISTORY OF MACHINE LEARNING (CONT.)



RELATED DISCIPLINES

- Artificial Intelligence
- Data Mining
- Probability and Statistics
- Information theory
- Numerical optimization
- Computational complexity theory
- Control theory (adaptive)
- Psychology (developmental, cognitive)
- Neurobiology
- Linguistics
- Philosophy

What are the goals of AI research?

Artifacts that THINK
like HUMANS

Artifacts that THINK
RATIONALLY

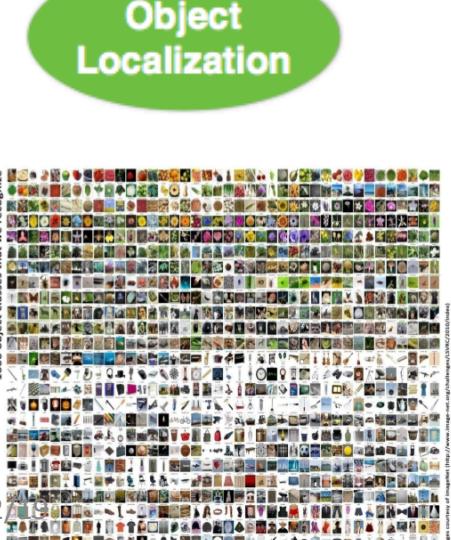
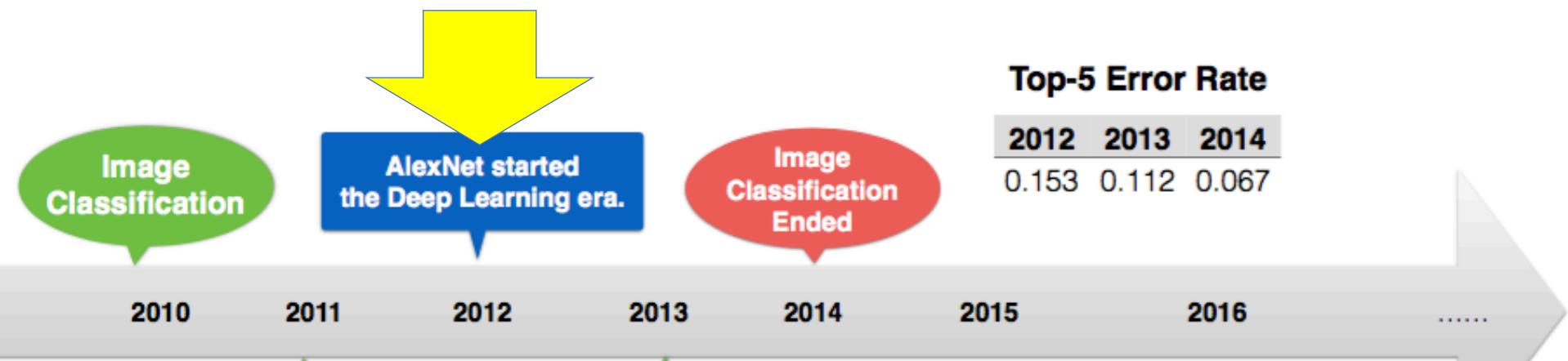
Artifacts that ACT
like HUMANS

Artifacts that ACT
RATIONALLY

How can we build more intelligent computer / machine ?

- Able to
 - **perceive the world**
 - **understand the world**
 - **react to the world**
- This needs
 - Basic speech capabilities
 - Basic vision capabilities
 - Language/semantic understanding
 - User behavior / emotion understanding
 - **Able to act**
 - **Able to think ??**

How can we build more intelligent computer / machine ? : Milestones in Recent Vision/AI Fields



72%, 2010

74%, 2011

85%, 2012

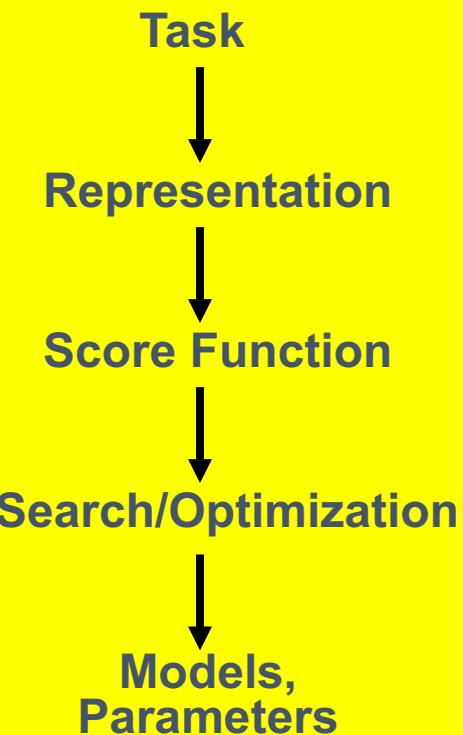
ImageNet Competition:

[Training on 1.2 million images [X] vs. 1000 different word labels [Y]]

Detour: three planned programming assignments about AI tasks

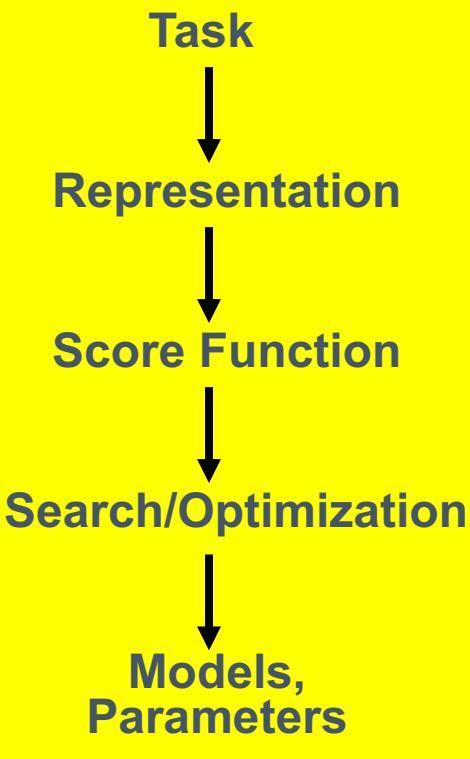
- HW: Semantic **language understanding** (sentiment classification on movie review text)
- HW: **Visual object recognition** (labeling images about handwritten digits)
- HW: **Audio speech recognition** (unsupervised learning based speech recognition task)

HW1



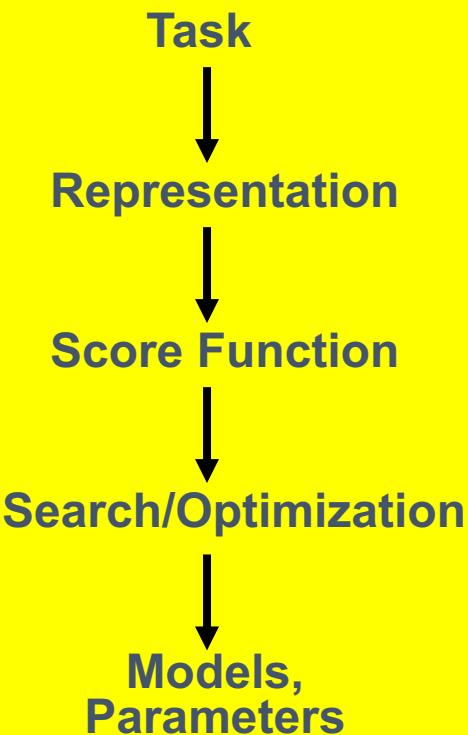
- Q1: Linear algebra review
- Q2: Linear regression + LOOCV
 - Regression
 - Evaluation pipeline
- Q3: Machine learning pipeline practice
 - Basic pipeline
 - GUI Toolbox
 - Evaluation

HW2



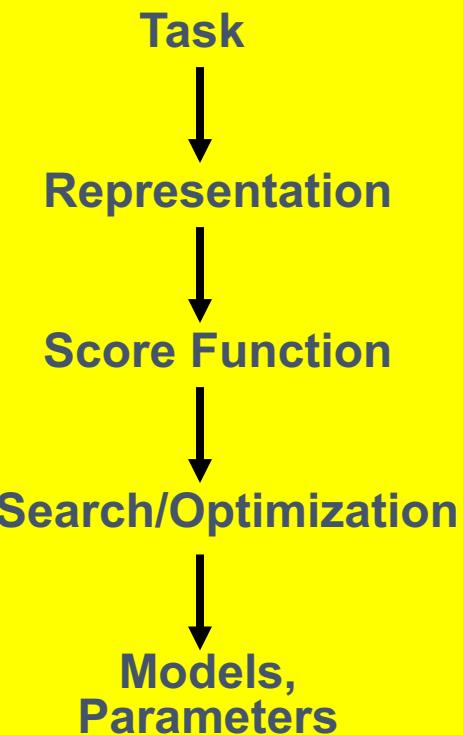
- Q1: Linear regression model fitting
 - Data loading
 - Basic linear regression
 - Three ways to train : Normal equation / SGD / Batch GD
 - Polynomial regression
- Q2: Ridge regression
 - Math derivation of ridge
 - Understand why/how Ridge
 - Model selection of Ridge with K-CV

HW5



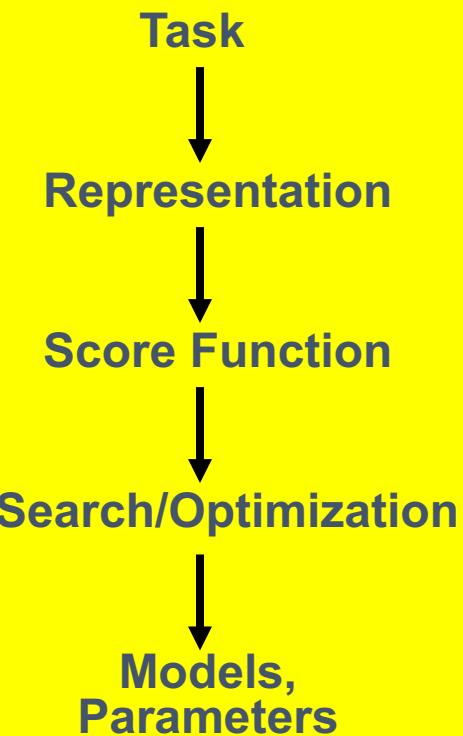
- Q1: Support Vector Machines with Scikit-Learn
 - Data preprocessing
 - How to use SVM package
 - Model selection for SVM
 - Model selection pipeline with train-vali, or train-CV; then test

HW5



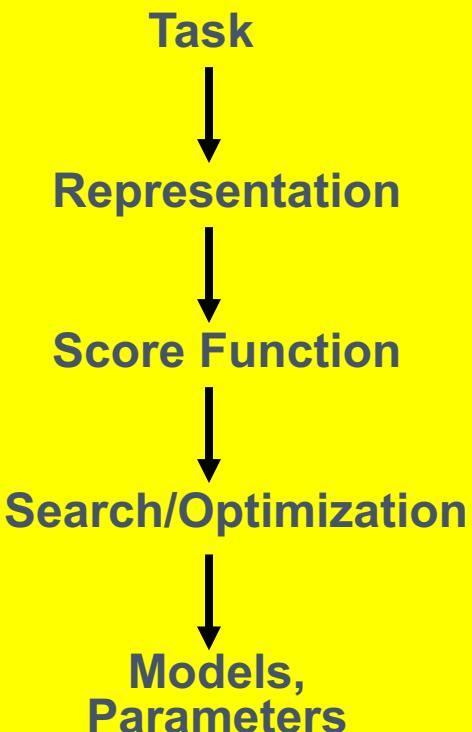
- Q1: Naive Bayes Classifier for Text-base Movie Review Classification
 - Preprocessing of text samples
 - BOW Document Representation
 - Multinomial Naive Bayes Classifier
 - BOW way
 - Language model way
 - Multivariate Bernoulli Naive Bayes Classifier

HW6



- Q1: Neural Network Tensorflow Playground
 - Interactive learning of MLP
 - Feature engineering vs.
 - Feature learning
- Q2: Image Classification
 - Tool using
 - DT / KNN / NN
 - PCA effect for image classification

HW6



- Q3: Unsupervised Clustering of audio data and consensus data
 - Data loading
 - K-mean clustering
 - GMM clustering
 - How to find K: knee-finding plot
 - How to measure clustering: purityMetric

References

- Prof. Andrew Moore's tutorials
- Prof. Raymond J. Mooney's slides
- Prof. Alexander Gray's slides
- Prof. Eric Xing's slides
- <http://scikit-learn.org/>
- Hastie, Trevor, et al. The elements of statistical learning. Vol. 2. No. 1. New York: Springer, 2009.
- Prof. M.A. Papalaskar's slides