

Artificial Intelligence III: Artificial Intelligence and Deep Learning

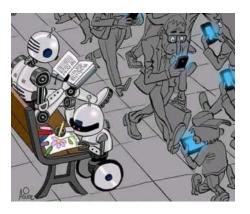
Lecture 2 **Fundamental Principles**

> Dr. Patrick Chan patrickchan@ieee.org South China University of Technology, China



Machine Learning

- Design an algorithm is not easy
- What if a machine can learn...



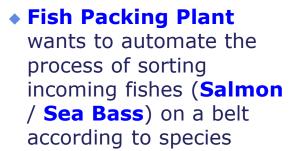
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Machine Learning: Fish Packing Plant Example

Salmon / Sea Bass





How to design a system?





Machine Learning: Fish Packing Plant Example

Salmon / Sea Bass



- Example: For a Fish
 - If Length > 10 and Fin Area > 10, Sea Bass
 - If Weight > 4.3 or Length < 4, Salmon</p>
- Factor of a rule:
 - Length, Weight, Color • Characteristic (Feature) Shape of fin / head.
 - Quantification (Threshold) Fin area > 10, Sea Bass Fin area < 10. Salmon



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Salmon / Sea Bass

- Process
- Difficult to determine a rule manually, even for an expert
 - How to define the rules (Feature and Threshold)?
 - E.g. If Weight > 4.3 or Length < 4, Salmon

Machine Learning can help

You do not know how but implement an algorithm which can learn from data



Object Sensing (camera) Preprocessing (Isolate Fish, reduce noise... Refined Image Feature Extraction (Take Measurement) Input Features Classification **↓** Output Class (Salmon / Sea Bass)

Machine Learning: Fish Packing Plant Example

Fish

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

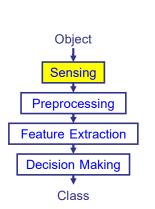


Machine Learning: Fish Packing Plant Example

Process

Sensing

- Digitize the object to the format which can be handled by machines
- Example
 - Type of Device Camera? Depth Camera? Infra-red? Ultrasound? Movement Sense? Combination?
 - Setting of Device Number? Angle? Overlap shooting range?
 - Background Lighting? Background simplicity?



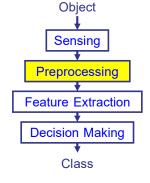


Machine Learning: Fish Packing Plant Example

Process

Preprocessing

- Refine the data
- Example
 - Lighting conditions
 - Position of fish
 - Angle of fish
 - Noise
 - Blurriness
 - Segmentation (remove object from background)







Process

Decision Making

Decision Type:

Class (Classification)

Region (Segmentation)

techniques are available

Many machine learning

Rank (Ranking)

Value (Regression, Value Prediction)

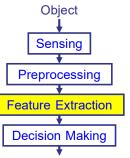
Action (Reinforcement Learning)

Machine Learning: Fish Packing Plant Example

le

Feature Extraction

- Decide which information is able to distinguish classes
- Example
 - Length, width, weight, number and shape of fins, tail shape, etc.
- Rely on technical background and common sense
 - Experts may help



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

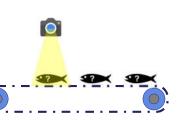


Machine Learning: Fish Packing Plant Example

Sensing & Preprocessing

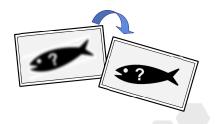


Assume a fish is put on a belt and a single camera is installed to take a photo on each fish



Preprocessing
 Remove the blueness

 and noise



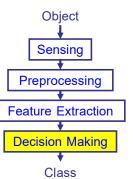
Class

Machine Learning: Fish Packing Plant Example

Feature Extraction



 Length is chosen (as a feature) as a decision criterion

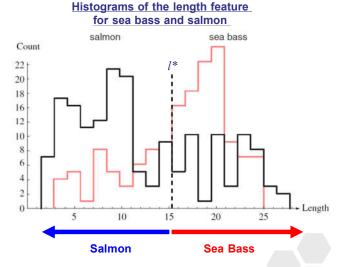




Feature Extraction



- 15 is selected as the threshold
- Although sea bass is longer in general, there are many exceptions
- The experts "may be" wrong!
- How about other features?
- E.g. lightness



Dr. Patrick Chan @ SCUT

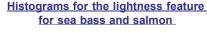
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

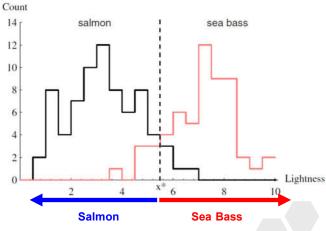
Machine Learning: Fish Packing Plant Example

Feature Extraction



- 5.5 is selected as the threshold
- "lightness" is better than "length"





Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Machine Learning: Fish Packing Plant Example

Cost Consideration



◆ Case 1: Company's view

Salmon \$ > Seabass \$

- Salmon is more expensive than sea bass. Selling Salmon with the price of sea bass will be a loss
 - If salmon is classified as sea bass: HIGH cost.
 - If sea bass is classified as salmon: LOW cost.

Case 2: Customer's view

- Customers who buy salmon will be upset if they get sea bass; Customers who buy sea bass will not be upset if they get the more expensive salmon
 - If salmon is classified as sea bass: LOW cost
 - If sea bass is classified as salmon: HIGH cost

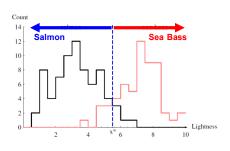
Machine Learning: Fish Packing Plant Example

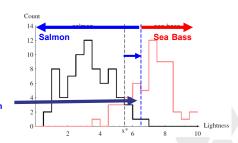
Cost Consideration

• Case 1: Company's view

- HIGH cost Salmon is classified as sea bass
- LOW cost Sea bass is classified as salmon
- Avoid classifying salmon wrongly by scarifying sea bass

More seabass Mistaken as salmon





Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

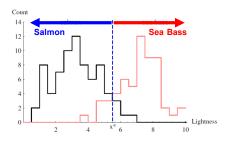


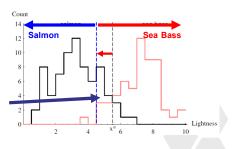
Cost Consideration



- LOW cost
 Salmon is classified as sea bass
- HIGH cost
 Sea bass is classified as salmon
- Avoid classifying sea bass wrongly by scarifying salmon

More salmon Mistaken as seabass



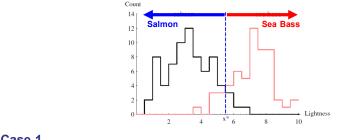


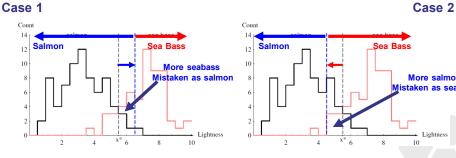
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

Machine Learning: Fish Packing Plant Example

Cost Consideration





18 Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Machine Learning: Fish Packing Plant Example

Multiple Features





More features should be considered

• Two features: Lightness (x_1) , Width (x_2)

A fish is represented by a point in a 2D
 feature space:

$$\mathbf{x} = \left[\begin{array}{c} x_1 \\ x_2 \end{array} \right]$$

The two features (lightness and width)

Width

Sea bass and salmon

sea bass

17

16

15

14

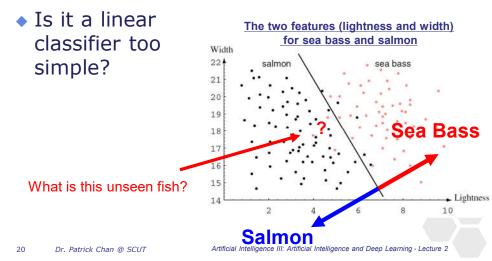
24

Artificial Intelligence Illi: Artificial Intelligence and Deep Learning - Lecture 2

Machine Learning: Fish Packing Plant Example Classifier



 A decision boundary can be drawn to divide the feature space into two regions





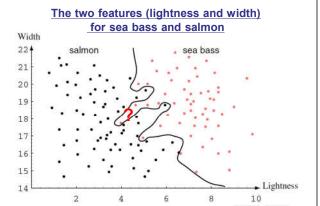
Classifier





- Will other classifiers be better?
 - More complex classifier
- Perfectly classify training samples
- Ultimate objective is to classify unseen samples correctly
- Can it be generalized to unseen sample?

Dr. Patrick Chan @ SCUT



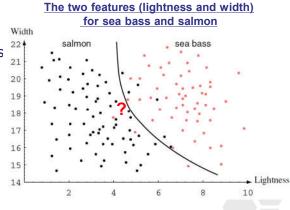
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

and complexity

Tradeoff between accuracy of training samples

Look more reasonable

- Not too complex
- Good in classifying the training samples



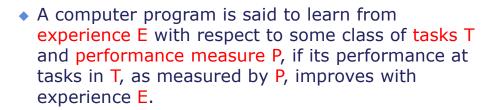
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Machine Learning: Fish Packing Plant Example

Summary



- Task T: Separate Salmon and Sea Bass
- Performance P : Accuracy on identification
- Experience E: Caught Salmon and Sea Bass



Key Concepts

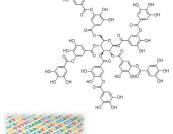
- Data Type
- Dataset (Collected Samples)
- Data Cleaning
- Model Requirement
- Evaluation
- Comparison
- Terminology



Key Concepts: Data Type

- Record-based data
 - Data matrix
 - Document data
 - Transaction data
- Graph-based data
 - World wide web
 - Molecular structure
 - Map data
- Ordered data
 - Spatial data
 - Temporal data
 - Sequential data
 - Genetic/Genomic sequence data

		Customer	Product ID Q	uantity
S00001	12/1/2012 9:00:00 AM	C0001	P025	1
S00002	12/1/2012 9:05:58 AM	C0025	P025	3
S00003	12/1/2012 9:11:33 AM	C0010	P001	2
S00004	12/1/2012 9:17:16 AM	C0017	P023	4
S00005	12/1/2012 9:23:04 AM	C0018	P016	5
S00006	12/1/2012 9:28:43 AM	C0011	P018	a 4
9007	12/1/2010 9-24-07 AM	Come	DOOS	-





Dr. Patrick Chan @ SCLIT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

Key Concepts: Data Type

Record-based Data

- Common illustration of Data
 - Excel
 - Traditional Database
- Properties / features of a specific object
- For example, human
 - Eye size (mm unit)
 - Eve color
 - Skin color
 - Height (cm unit)
 - Wear glasses or not
 - Gender
 - Age
 - Length of finger (cm unit)

Object, Instance Individual Sample

Subject

Attributes, Characteristics, Features, Variables, ..., etc

Ha	Retuna	Status	Income
1	Yes	Single	125K
2	No	Married	100K
3	No	Single	70K
4	Yes	Married	120K
5	No	Divorced	95K
6	No	Married	60K
7	Yes	Divorced	220K
8	No	Single	85K
9	No	Married	75K
10	No	Single	90K

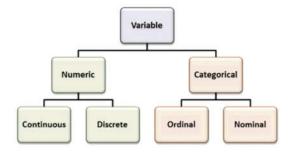
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Key Concepts: Data Type

Record-based Data: Variable



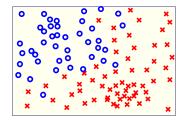
- Continuous: real number, e.g. height = 167.23cm
- Discrete: integer, e.g. age = 1
- Nominal: a natural order or rank, e.g. High Low
- Ordinal: no order, e.g. Red Blue Yellow



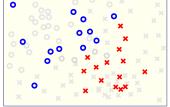
Key Concepts

Dataset (Collected Samples)

- Aim of machine learning Perform well on all samples
- What information we have? Collected samples







Population

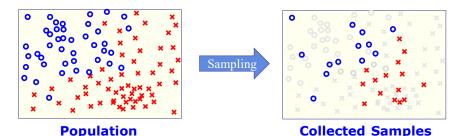
- All possible samples (usually infinite)
- Usually represented by a distribution

Collected Samples

- Limited number
- Subset of population



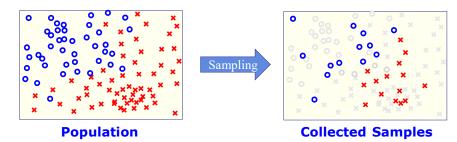
Dataset (Collected Samples)



- Any samples can be chosen?
 - Distribution of collected samples should be similar to the population's one
 - Independent and identically distributed (iid) is assumed
 - Randomly choose without any bias

Dr. Patrick Chan @ SCUT Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

Dataset (Collected Samples)



- More samples are better?
 - Usually yes when the sampling method is fair

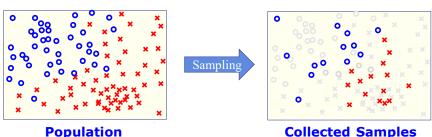
30 Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



Key Concepts

Dataset (Collected Samples)



- How many samples are enough?
 - Depend on the complexity of the problem
 - Rely on the performance of the trained model
 - General answer: more is better



Key Concepts Data Cleaning

- Once you have finished collecting samples, can we start the learning immediately?
- Garbage in Garbage out



Key Concepts

Data Cleaning

 What data scientists spend the most time doing?

Building training sets 3%

Cleaning and organizing data

60%

Collecting data set

19%

Mining data

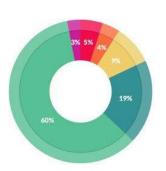
9%

Refining algorithms

4%

Other





Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

Key Concepts: Data Cleaning Bad Quality Data

- Identifying damaged or inaccurate, incomplete, incorrect, or irrelevant parts within data
- Replacing, modifying, or deleting the dirty or rough data

Quality

- Noise
- Outliers
- Missing values
- Duplicate data



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Key Concepts: Data Cleaning

Noise

- Noise can refer to any random fluctuation in a signal
- Attribute noise Errors in features
 - Value errors
 - Incorrect or erroneous values Missing values Data entries that are absent
 - Irrelevant values Does not contribute

Feature 1	Feature 2	Class
0.25	Red	+
0.25	Red	-
0.99	Yellow	-
1.02	Yellow	+
2.05	Yellow	-

Class noise

Variability or errors in the labels

- Contradictory examples Identical features but different labels.
- Mislabeled examples label is incorrect



Key Concepts: Data Cleaning

Noise

- Median/mean noise filter
 - Apply a sliding window to an image
 - Sort the pixel values within the window
 - Calculate the median/average value from the sorted values as the new value for the current pixel
 - Repeat until all the pixels are processed



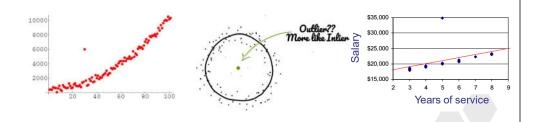
Dr. Patrick Chan @ SCUT



Input					Output						
			1	3	1	1	4	0	1	3	1
	2		2	2	3	2	1	1	1	1	3
			0	1	0	1	1	1	1	2	0
1	2	1	0	2	2	1	1	1	1	1	2
2	5	3	1	2	5	2	2	2	2	2	5
1	1	4	2	3	0	1	1	4	2	3	0



- An outlier is a data point that differs significantly from other observations
- Identification:
 - Subjective: visualization
 - Objective: statistical way, i.e. Cook's D value



Dr. Patrick Chan @ SCLIT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

Key Concepts: Data Cleaning Outlier

Deletion

If the outlier is caused by human error (i.e. typo, unrealistic response)

Replacement

- Replace observations with other values (mean, etc.)
- Handle Differently
 - Analyze outliners separately from the rest

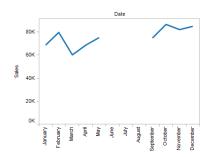
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



Key Concepts: Data Cleaning Missing Value

- No data value is stored
- Missing data are a common occurrence and can have a significant effect on the conclusions





Key Concepts: Data Cleaning Missing Value

Deletion

- Delete the whole observation with missing values
- Partial deletion (delete the part of missing values in downstream modeling)
- Replace with other values
 - mean, mode, median
 - Possibility that the model will be distorted
- Insert predicted values
 - Imputation by using statistical or machine learning methods

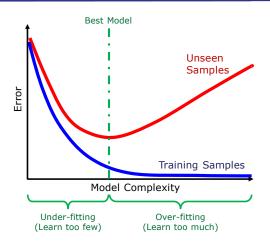


- Objects that are duplicates, or almost duplicates of one another
- Common issue when collecting data from heterogeneous sources
 - E.g. If a person has multiple email addresses
- Deletion
 - Which record should be deleted?
 - Time
 - Source Trustworthy
 - Majority

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2

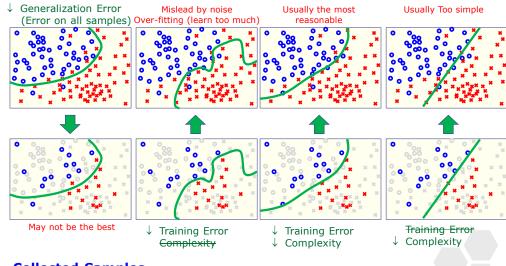
Key Concepts Model Requirements



- Be noted that this is just a common situation in machine learning, where each application varies
 - E.g. Having more training samples reduces the variance between two lines.



Population



Collected Samples

42 Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 2



Key Concepts **Evaluation**

- How to evaluate our trained model?
 - Evaluate by training samples?
 No, already see in training
 - What information we have? Collected samples
 - Some collected samples should not be used in training
 - Separate into two non-overlapping sets:
 - Training set: For training
 - Test set: For evaluation

Key Concepts Comparing Classifiers

- For a classification problem, given
 - Dataset D
 - Classifiers A and B
- How can we measure which classifier, A or B, is better for D?

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture43



Key Concepts

Comparing Classifiers

Method

- Randomly separate D into training and test sets
- Use Training Set to train A and B
- Use Test Set to measure the performances of the trained A and B
- Select the better performing classifier

Is it ok?

- The winner may just be lucky in performing better for that particular test set.
- No guarantee for different test sets

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture4



Key Concepts

Comparing Classifiers

- The bias of test set should be reduced
- Two re-sampling techniques
 - Independent Run
 - Cross-Validation



Key Concepts: Comparing Classifiers

Independent Run

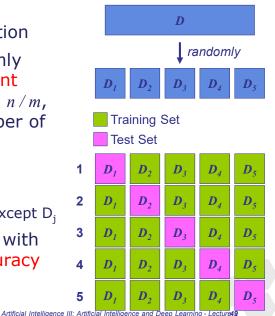
- Statistical method
- Also called Bootstrap and Jackknifing
- Repeat the experiment "n" times independently
 - Repeat n times
 - *i* is the number of running time
 - Randomly separate D into Training Set, and Test Set,
 - Use Training Set, to train A_i and B_i
 - Use Test Set, to evaluate the trained A, and B,
 - Select the classifier with higher average accuracy



Key Concepts: Comparing Classifiers

Cross-Validation

- M-fold Cross-Validation
- Dataset D is randomly divided into *m* disjoint sets D_i of equal size n/m_i where n is the number of samples in dataset
- Repeat m times
 - Trained by D_i
 - Evaluated by all D_i except D_i
- Select the classifier with higher average accuracy



Key Concepts Terminology

- Instance / Sample Observations from an application
- Feature / Attribute Property or characteristic of a sample
- Dimensionality The number of features

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



Key Concepts

Dr Patrick Chan @ SCLIT

Terminology

Training Set

A set of samples used to train a model

Test Set

A set of samples used to evaluate the performance of the trained model. Usually separate from the training set.

Unseen Samples

Any samples not in training set



Key Concepts

Terminology

Training Error

Error on training samples

Test Error

Error on test samples

Generalization Error

The ability of a model to perform well on unseen samples

In some discussion, Test Frror = Generalization Frror







Objective Function / Error Function

A mathematical function used to quantify error made by a model, closely related to the objective

Can be more than error on samples, may include any other concepts

E.g. complexity of a model

