



COMP7/8118 M50

Data Mining

Classification Basics

Xiaofei Zhang

Slides compiled from Jiawei Han and Raymond C.W. Wong's work

THE UNIVERSITY OF
MEMPHIS

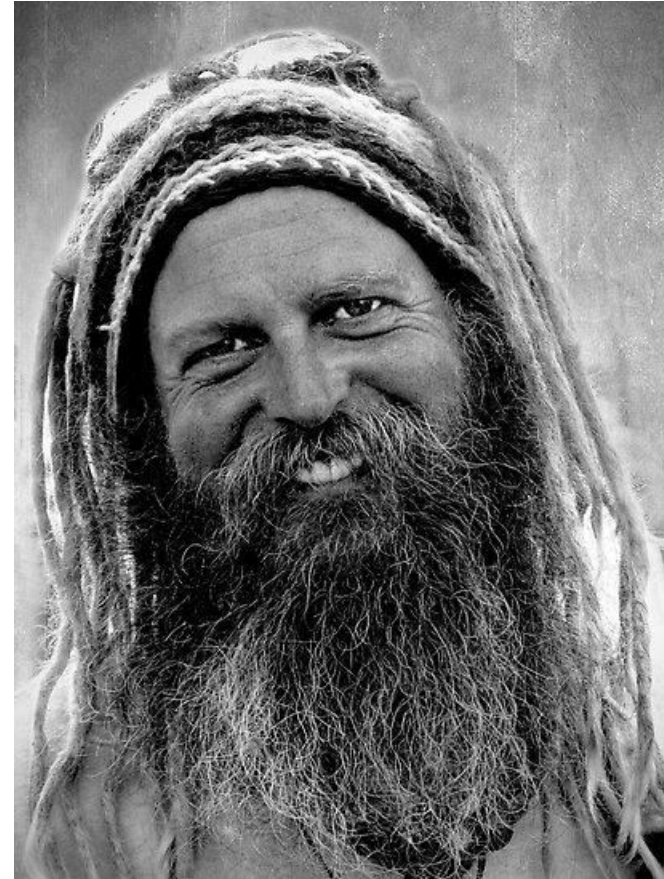
What is a hipster?

- Examples of hipster look



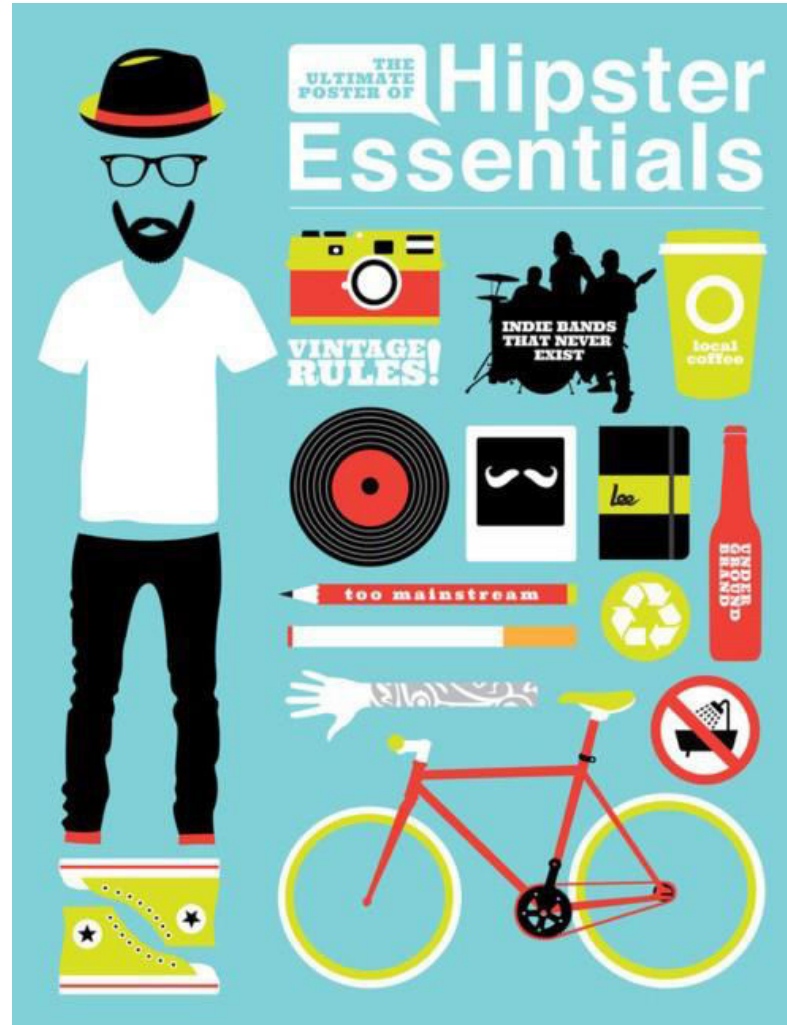
- A hipster is defined by facial hair

Hipster or Hippie?



Facial hair alone is not enough to characterize hipsters

How to be a hipster



There is a big set of **features** that defines a hipster

Classification

- The problem of discriminating between different **classes** of objects
 - In our case: Hipster vs. Non-Hipster
- Classification process:
 - Find **examples** for which you know the class (**training set**)
 - Find a set of **features** that discriminate between the examples within the class and outside the class
 - Create a **function** that given the features decides the class
 - **Apply** the function to new examples.

Catching tax-evasion

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Tax-return data for year 2011

A new tax return for 2012
Is this a cheating tax return?

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

An instance of the classification problem: learn a method for discriminating between records of different **classes** (**cheaters** vs **non-cheaters**)

What is classification?

- **Classification** is the task of *learning a target function* f that maps attribute set x to one of the predefined class labels y

Why classification?

- The target function f is known as a **classification model**
- **Descriptive modeling:** **Explanatory tool** to distinguish between objects of different classes (e.g., understand why people cheat on their taxes, or what makes a hipster)
- **Predictive modeling:** Predict a class of a **previously unseen** record

Examples of Classification Tasks

- Predicting tumor cells as **benign** or **malignant**
- Classifying credit card transactions as **legitimate** or **fraudulent**
- Categorizing news stories as **finance**, **weather**, **entertainment**, **sports**, etc
- Identifying spam email, spam web **pages**, **adult content**
- Understanding if a web query has **commercial intent** or not

Classification is **everywhere** in data science.

General approach to classification

- **Training set** consists of records with **known class labels**
- Training set is used to **build** a classification model
- A **labeled test set** of **previously unseen** data records is used to **evaluate** the quality of the model.
- The classification model is **applied** to new records with **unknown class labels**

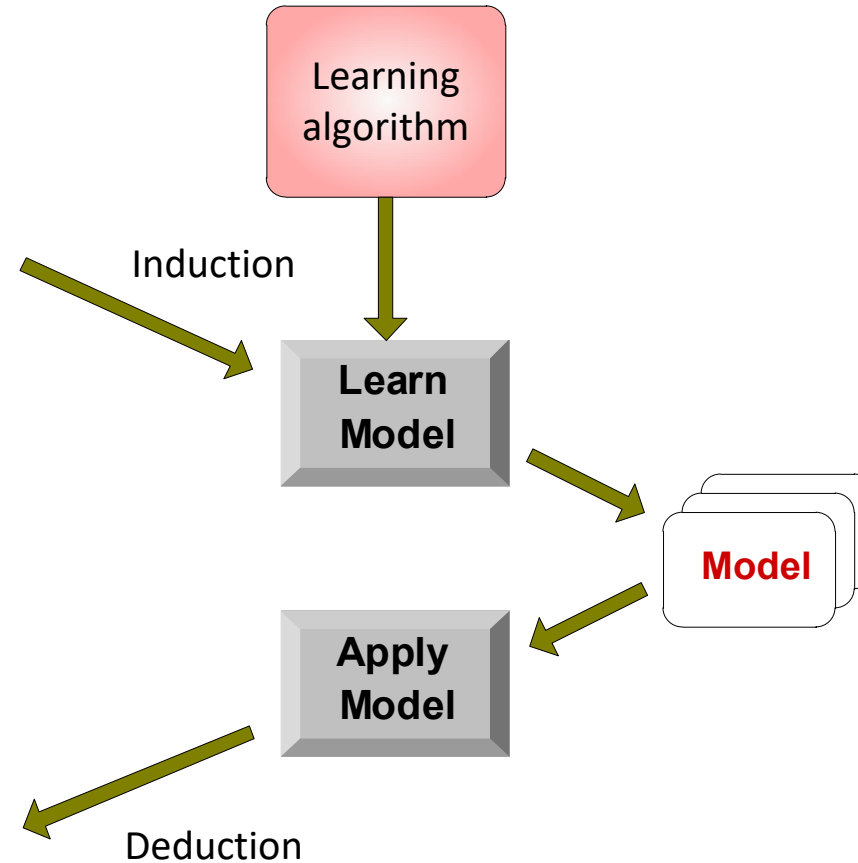
Illustrating Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Evaluation of classification models

- Counts of **test records** that are correctly (or incorrectly) predicted by the classification model
- **Confusion matrix**

Actual Class	Predicted Class	
	Class = 1	Class = 0
	Class = 1	Class = 0
Class = 1	f_{11}	f_{10}
Class = 0	f_{01}	f_{00}

$$\text{Accuracy} = \frac{\# \text{ correct prediction s}}{\text{total \# of prediction s}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

$$\text{Error rate} = \frac{\# \text{ wrong prediction s}}{\text{total \# of prediction s}} = \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

Classifiers to be covered

- Decision Tree
- Bayesian Model
- Nearest Neighbor Model
- Support Vector Machine
- Neural Network