

# UCI Datasets

CSE 4311 – Machine Learning

Vassilis Athitsos

Computer Science and Engineering Department

University of Texas at Arlington

# Programming Assignment - Datasets

- The assignment introduces three dataset.
- Each dataset has a training set and a test set.
  - These are the first four lines of the training set of the Yeast dataset:

0.5000	0.4600	0.6400	0.3600	0.5000	0	0.4900	0.2200	1
0.5300	0.5600	0.4900	0.4600	0.5000	0	0.5200	0.2200	1
0.5200	0.5300	0.5800	0.6900	0.5000	0	0.5000	0.2200	1
0.6700	0.6200	0.5400	0.4300	0.5000	0	0.5300	0.2200	1

# Programming Assignment - Datasets

- Each row is a training example.
- All columns **except for the last column** represent the training input, which is a vector.
- The last column is the class label, which can be a number or a string.

Training input 1	Training label 1
0.5000 0.4600 0.6400 0.3600 0.5000 0 0.4900 0.2200	1
0.5300 0.5600 0.4900 0.4600 0.5000 0 0.5200 0.2200	1
0.5200 0.5300 0.5800 0.6900 0.5000 0 0.5000 0.2200	1
0.6700 0.6200 0.5400 0.4300 0.5000 0 0.5300 0.2200	1

# Programming Assignment - Datasets

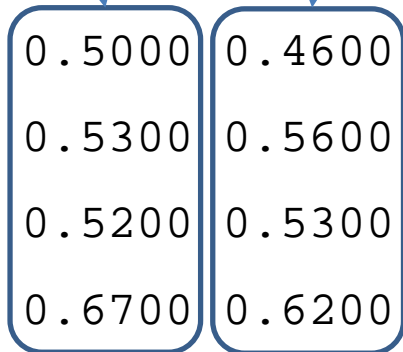
- Each row is a training example.
- All columns **except for the last column** represent the training input, which is a vector.
- The last column is the class label, which can be a number or a string.

0.5000	0.4600	0.6400	0.3600	0.5000	0	0.4900	0.2200	1
0.5300	0.5600	0.4900	0.4600	0.5000	0	0.5200	0.2200	1
0.5200	0.5300	0.5800	0.6900	0.5000	0	0.5000	0.2200	1
0.6700	0.6200	0.5400	0.4300	0.5000	0	0.5300	0.2200	1

# Programming Assignment - Datasets

- Each column (except for the last column) is called a **dimension**, or an **attribute**, or a **feature**.
- In the Yeast dataset, there are 8 attributes/dimensions/features.
- Different datasets have different numbers of attributes.

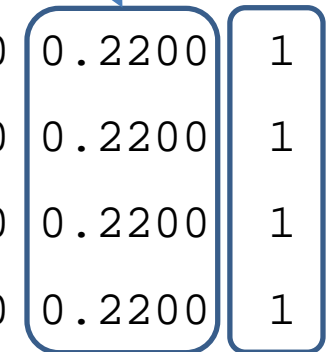
Attribute 1      Attribute 2



0.5000	0.4600	0.6400	0.3600	0.5000	0
0.5300	0.5600	0.4900	0.4600	0.5000	0
0.5200	0.5300	0.5800	0.6900	0.5000	0
0.6700	0.6200	0.5400	0.4300	0.5000	0

Attribute 8

Training  
labels



0.4900	0.2200	1
0.5200	0.2200	1
0.5000	0.2200	1
0.5300	0.2200	1

# Programming Assignment - Datasets

- So, for example, what is the value for attribute 4 of the third training input?

0.5000	0.4600	0.6400	0.3600	0.5000	0	0.4900	0.2200	1
0.5300	0.5600	0.4900	0.4600	0.5000	0	0.5200	0.2200	1
0.5200	0.5300	0.5800	0.6900	0.5000	0	0.5000	0.2200	1
0.6700	0.6200	0.5400	0.4300	0.5000	0	0.5300	0.2200	1

# Programming Assignment - Datasets

- So, for example, what is the value for attribute 4 of the third training input?
- It is 0.6900.
- We use matrix notation, indices start from 1, not 0.

0.5000	0.4600	0.6400	0.3600	0.5000	0	0.4900	0.2200	1
0.5300	0.5600	0.4900	0.4600	0.5000	0	0.5200	0.2200	1
0.5200	0.5300	0.5800	0.6900	0.5000	0	0.5000	0.2200	1
0.6700	0.6200	0.5400	0.4300	0.5000	0	0.5300	0.2200	1

# Class Labels Can Be Strings

- Your code needs to convert class labels to one-hot vectors.
- I suggest a two step approach, as shown on the slides discussing multiclass problems:
  - Step 1: map class labels to consecutive integers that start from 1.
  - Step 2: map integer class labels to one-hot vectors.
  - For easy reference, the next few slides are a copy of the slides we have seen before, that describe these two steps.



# Converting to One-Hot-Vectors

- Suppose we have this training set:
  - $\mathbf{x}_1 = (0.5, 2.4, 8.3, 1.2, 4.5)^T$ ,  $q_1 = \text{dog}$ ,
  - $\mathbf{x}_2 = (3.4, 0.6, 4.4, 6.2, 1.0)^T$ ,  $q_2 = \text{dog}$ ,
  - $\mathbf{x}_3 = (4.7, 1.9, 6.7, 1.2, 3.9)^T$ ,  $q_3 = \text{cat}$ ,
  - $\mathbf{x}_4 = (2.6, 1.3, 9.4, 0.7, 5.1)^T$ ,  $q_4 = \text{fox}$ ,
  - $\mathbf{x}_5 = (8.5, 4.6, 3.6, 2.0, 6.2)^T$ ,  $q_5 = \text{cat}$ ,
  - $\mathbf{x}_6 = (5.2, 8.1, 7.3, 4.2, 1.6)^T$ ,  $q_6 = \text{fox}$ ,
- Step 1:
  - Generate new class labels  $s_n$ , where classes are numbered sequentially starting from 1. Thus, in our example, the class labels become 1, 2, 3.

# Converting to One-Hot-Vectors

- Suppose we have this training set:
  - $\mathbf{x}_1 = (0.5, 2.4, 8.3, 1.2, 4.5)^T$ ,  $q_1 = \text{dog}$ ,  $s_1 = 1$
  - $\mathbf{x}_2 = (3.4, 0.6, 4.4, 6.2, 1.0)^T$ ,  $q_2 = \text{dog}$ ,  $s_2 = 1$
  - $\mathbf{x}_3 = (4.7, 1.9, 6.7, 1.2, 3.9)^T$ ,  $q_3 = \text{cat}$ ,  $s_3 = 2$
  - $\mathbf{x}_4 = (2.6, 1.3, 9.4, 0.7, 5.1)^T$ ,  $q_4 = \text{fox}$ ,  $s_4 = 3$
  - $\mathbf{x}_5 = (8.5, 4.6, 3.6, 2.0, 6.2)^T$ ,  $q_5 = \text{cat}$ ,  $s_5 = 2$
  - $\mathbf{x}_6 = (5.2, 8.1, 7.3, 4.2, 1.6)^T$ ,  $q_6 = \text{fox}$ ,  $s_6 = 3$
- Step 1:
  - Generate new class labels  $s_n$ , where classes are numbered sequentially starting from 1. Thus, in our example, the class labels become 1, 2, 3.

# Converting to One-Hot-Vectors

- Training set:

– $\mathbf{x}_1 = (0.5, 2.4, 8.3, 1.2, 4.5)^T$ , $s_1 = 1$	$\mathbf{t}_1 = (?, ?, ?)^T$
– $\mathbf{x}_2 = (3.4, 0.6, 4.4, 6.2, 1.0)^T$ , $s_2 = 1$	$\mathbf{t}_2 = (?, ?, ?)^T$
– $\mathbf{x}_3 = (4.7, 1.9, 6.7, 1.2, 3.9)^T$ , $s_3 = 2$	$\mathbf{t}_3 = (?, ?, ?)^T$
– $\mathbf{x}_4 = (2.6, 1.3, 9.4, 0.7, 5.1)^T$ , $s_4 = 3$	$\mathbf{t}_4 = (?, ?, ?)^T$
– $\mathbf{x}_5 = (8.5, 4.6, 3.6, 2.0, 6.2)^T$ , $s_5 = 2$	$\mathbf{t}_5 = (?, ?, ?)^T$
– $\mathbf{x}_6 = (5.2, 8.1, 7.3, 4.2, 1.6)^T$ , $s_6 = 3$	$\mathbf{t}_6 = (?, ?, ?)^T$

- Step 2: Convert each label  $s_n$  to a **one-hot vector**  $\mathbf{t}_n$ .
  - Vector  $\mathbf{t}_n$  has as many dimensions as the number of classes.
    - In our example we have three classes, so each  $\mathbf{t}_n$  is 3-dimensional.
  - If  $s_n = i$ , then set the  $i$ -th dimension of  $\mathbf{t}_n$  to 1.
  - Otherwise, set the  $i$ -th dimension of  $\mathbf{t}_n$  to 0.

# Converting to One-Hot-Vectors

- Training set:

– $\mathbf{x}_1 = (0.5, 2.4, 8.3, 1.2, 4.5)^T$ , $s_1 = 1$	$\mathbf{t}_1 = (1, 0, 0)^T$
– $\mathbf{x}_2 = (3.4, 0.6, 4.4, 6.2, 1.0)^T$ , $s_2 = 1$	$\mathbf{t}_2 = (1, 0, 0)^T$
– $\mathbf{x}_3 = (4.7, 1.9, 6.7, 1.2, 3.9)^T$ , $s_3 = 2$	$\mathbf{t}_3 = (0, 1, 0)^T$
– $\mathbf{x}_4 = (2.6, 1.3, 9.4, 0.7, 5.1)^T$ , $s_4 = 3$	$\mathbf{t}_4 = (0, 0, 1)^T$
– $\mathbf{x}_5 = (8.5, 4.6, 3.6, 2.0, 6.2)^T$ , $s_5 = 2$	$\mathbf{t}_5 = (0, 1, 0)^T$
– $\mathbf{x}_6 = (5.2, 8.1, 7.3, 4.2, 1.6)^T$ , $s_6 = 3$	$\mathbf{t}_6 = (0, 0, 1)^T$

- Step 2: Convert each label  $s_n$  to a **one-hot vector**  $\mathbf{t}_n$ .
  - Vector  $\mathbf{t}_n$  has as many dimensions as the number of classes.
    - In our example we have three classes, so each  $\mathbf{t}_n$  is 3-dimensional.
  - If  $s_n = i$ , then set the  $i$ -th dimension of  $\mathbf{t}_n$  to 1.
  - Otherwise, set the  $i$ -th dimension of  $\mathbf{t}_n$  to 0.