

MULTI-VARIATE ANALYSIS

DATA MINING

2/3

We are interested in extending the analysis of a dataset to the **MULTI-VARIATE** case. In particular, if previously we studied the dataset one feature (e.k.a. one column vector) at a time, now we want to consider several features jointly.

	f_1	...		
CASE SAMPLE	S_1	...		
	S_2	...		

(UNI-VARIATE VIEW)

	f_1	f_2		
CASE SAMPLE	S_{11}	S_{12}		
	S_{21}	S_{22}		

(MULTI-VARIATE VIEW)

Here the S_i are scalar values and we count the # of scalars to estimate the feature's distribution.

Here the **CASE SAMPLE** are not scalars but row-vectors.

PMF ESTIMATION

Given a dataset with discrete features p_1, p_2 , if we want to estimate their joint pmf we have to count the row vectors and compute for each unique row vector its frequency.

Thus, we have,

$$pmf_{p_1+p_2}(\underline{s}_i) := \frac{\# \text{ of rows equals to } \underline{s}_i}{\text{total } \# \text{ of rows}}$$

Q: What if a particular row never appears in the dataset? Don't we want to apply some SMOOTHING?

SCATTER PLOTS

Scatter plots allow us to visualize if there are any spheres in which all the samples of a particular class are situated.

NOISY DATA is data in which all the samples are uniformly distributed with respect to the class.

A scatter plot allows us to understand if the data is noisy or not.

TODO: Add image of scatter plots.

PMF VS PDF

Consider a dataset that contains numerical data such as the width or the height of something.

When we estimate the distribution of said feature, should we estimate a pmf or a pdf?

If we want to be rigorous, we have to treat the feature as a continuous r.v., and thus we have to estimate a pdf. However, if the initial accuracy of the dataset is low (i.e. few decimal digits), we can **DISCRETIZE** the dataset and treat the feature as a discrete r.v.

INT NUMBERS \rightarrow PMF (DISCRETE R.V.)
FLOAT NUMBERS \rightarrow PDF (CONTINUOUS R.V.)

\nearrow DISCRETIZATION

A notable method to discretize a dataset is to apply the `FLOOR()` function.