

A brief description of the phoneme data

This dataset is a part of the original one which can be found at <http://www-stat.stanford.edu/ElemStatLearn>. We observe $n = 2000$ pairs $(\mathbf{x}_i, y_i)_{i=1,\dots,n}$ where the \mathbf{x}_i 's correspond to the discretized log-periodograms ($\mathbf{x}_i = (\chi(f_1), \chi(f_2), \dots, \chi(f_{150}))$ is the i th discretized functional data) whereas the y_i 's give the class membership (five phonemes):

$$y_i \in \{1, 2, 3, 4, 5\} \text{ with } \begin{cases} 1 \longleftrightarrow "sh" \\ 2 \longleftrightarrow "iy" \\ 3 \longleftrightarrow "dcl" \\ 4 \longleftrightarrow "aa" \\ 5 \longleftrightarrow "ao" \end{cases}$$

The phoneme dataset “npfda-phoneme.dat” contains the pairs $(\mathbf{x}_i, y_i)_{i=1,\dots,2000}$ and is organized as follows:

	Col 1	\cdots	Col j	\cdots	Col 150	Col 151
Row 1	$\chi_1(f_1)$	\cdots	$\chi_1(f_j)$	\cdots	$\chi_1(f_{150})$	y_1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
Row i	$\chi_i(f_1)$	\cdots	$\chi_i(f_j)$	\cdots	$\chi_i(f_{150})$	y_i
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
Row 2000	$\chi_{2000}(f_1)$	\cdots	$\chi_{2000}(f_j)$	\cdots	$\chi_{2000}(f_{150})$	y_{2000}

The first 150 columns correspond to the 150 frequencies whereas the last column contains the categorical responses (class number). Note that the size of each class is the same (400).