

PROF. FRANCESCO A.N. PALTIERI
UNIV. DELLA CAMPANIA NOV. 2023 (14)

APPENDIX DISCRETE VARIABLES REPRESENTATIONS

The manipulating classifier data, as we will see in the following, its convenient to use continuous mathematics because we can compute gradients and define smooth functional transformations.

One of the problem that we have to face to train a parametric model $f(x, \theta)$ is that our training pairs may include a desired output y that belongs to a discrete, non necessarily numerical set.

More specifically, if the objective of $f(x, \theta)$ is to produce a decision, or a posterior distribution, what should be the representation of S in our use of a loss function? *Loss function.*

For example, when the classifier alphabet is $\mathcal{X} = \{ \text{HOME, GRASS, SKY, HORSE} \}$, how can we define a loss function to be used for training in comparing $f(x, \theta)$ to each element of \mathcal{X} ?

Furthermore a classification, or a regression problem, may be based on discrete inputs. We have so far always assumed that $x \in \mathcal{X} \subset \mathbb{R}^n$, but there are many problems in which x may be discrete.

To translate discrete variables into continuous representations there are various strategies:

Assume that each symbol $s \in \mathcal{X} = \{a_1, \dots, a_M\}$ A2
is translated into a real vector \underline{v}
belonging to a finite set

$$\underline{v} \in \mathcal{V} = \{\underline{v}_1, \underline{v}_2, \dots, \underline{v}_M\}$$

$$\underline{s} \in \mathcal{X} \longrightarrow \underline{v} \in \mathbb{R}^N$$

$$a_1 \longrightarrow \underline{v}_1$$

$$a_2 \longrightarrow \underline{v}_2$$

⋮

$$a_M \longrightarrow \underline{v}_M$$

MAGNITUDE ENCODING

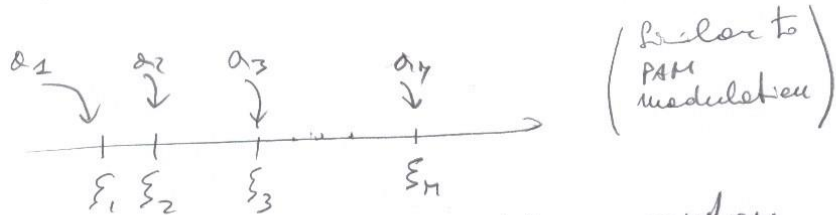
The most direct way of encoding the symbols
may be to use scalar values ($N=1$)



DIRECT BINARY CODING

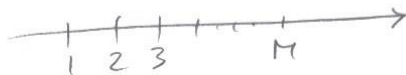
(A3)

The symbols in $\mathcal{X} = \{a_1, a_2, \dots, a_M\}$ are translated into a finite discrete set of values $\{\xi_1, \dots, \xi_M\}$ on the real line.



If the original set \mathcal{X} is not ordered the association may be arbitrary.

For example $\xi_1, \xi_2, \xi_3, \dots, \xi_M$ may be equally spaced and/or be the integers $1, \dots, M$.



This representation in neural networks, is rarely used because it lacks robustness with respect to misclassification.

TEMPERATURE CODING

$s \in \mathcal{X} = \{a_1, \dots, a_M\} \rightarrow \underline{v}$ M dimensional vector of 0 and 1's

Ex $M=5$

$$a_1 \rightarrow \underline{v}_1 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}; a_2 \rightarrow \underline{v}_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}; a_3 \rightarrow \underline{v}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}; a_4 \rightarrow \underline{v}_4 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}; a_5 \rightarrow \underline{v}_5 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

ONE-HOT ENCODING

$s \in \mathcal{X} = \{a_1, \dots, a_M\} \rightarrow \underline{v}$ M dimensional vector with all zeros except one equal to one.

Ex $M=5$

$$a_1 \rightarrow \underline{v}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}; a_2 \rightarrow \underline{v}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}; a_3 \rightarrow \underline{v}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}; a_4 \rightarrow \underline{v}_4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}; a_5 \rightarrow \underline{v}_5 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

ONE-HOT ϵ CODING

ϵ

$$S \in \mathcal{X} = \{a_1, \dots, a_M\} \rightarrow \underline{v} \quad (0 \leq \epsilon < 1 \text{ small})$$

$M=5$

$$a_1 \rightarrow \underline{v}_1 = \begin{bmatrix} 1-\epsilon \\ \epsilon/4 \\ \epsilon/4 \\ \epsilon/4 \\ \epsilon/4 \end{bmatrix} \quad a_2 \rightarrow \underline{v}_2 = \begin{bmatrix} \epsilon/4 \\ 1-\epsilon \\ \epsilon/4 \\ \epsilon/4 \\ \epsilon/4 \end{bmatrix} \quad a_3 \rightarrow \underline{v}_3 = \begin{bmatrix} \epsilon/4 \\ \epsilon/4 \\ 1-\epsilon \\ \epsilon/4 \\ \epsilon/4 \end{bmatrix} \quad a_4 \rightarrow \underline{v}_4 = \begin{bmatrix} \epsilon/4 \\ \epsilon/4 \\ \epsilon/4 \\ 1-\epsilon \\ \epsilon/4 \end{bmatrix} \quad a_5 \rightarrow \underline{v}_5 = \begin{bmatrix} \epsilon/4 \\ \epsilon/4 \\ \epsilon/4 \\ \epsilon/4 \\ 1-\epsilon \end{bmatrix}$$

In both ONE-HOT and ONE-HOT ϵ REPRESENTATIONS \underline{v} is a distribution.
 ONE-HOT ϵ representations used to avoid the presence of zeros that may create numerical instabilities during learning.

BINARY UNIPOLAR CODING (FULL)

$$M=2^q \quad S \in \mathcal{X} = \{a_1, \dots, a_M\} \rightarrow \underline{v} \in \{0, 1\}^q$$

Ex $M=4$

$$a_1 \rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}^{v_1} \quad a_2 \rightarrow \begin{bmatrix} 0 \\ 1 \end{bmatrix}^{v_2} \quad a_3 \rightarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix}^{v_3} \quad a_4 \rightarrow \begin{bmatrix} 1 \\ 1 \end{bmatrix}^{v_4}$$

BINARY BIPOLAR CODING (FULL)

$$M=2^q \quad S \in \mathcal{X} = \{a_1, \dots, a_M\} \rightarrow \underline{v} \in \{-1, 1\}^q$$

Ex $M=4$

$$a_1 \rightarrow \begin{bmatrix} -1 \\ -1 \end{bmatrix}^{v_1} \quad a_2 \rightarrow \begin{bmatrix} -1 \\ 1 \end{bmatrix}^{v_2} \quad a_3 \rightarrow \begin{bmatrix} 1 \\ -1 \end{bmatrix}^{v_3} \quad a_4 \rightarrow \begin{bmatrix} 1 \\ 1 \end{bmatrix}^{v_4}$$

BINARY UNIPOLAR CODING (SPARSE)

$$M \ll 2^q \quad S \in \mathcal{X} \in \{a_1, \dots, a_M\} \rightarrow \underline{v} \in \{0, 1\}^q$$

Ex $M=5, q=4$

$$a_1 \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix}^{v_1} \quad a_2 \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}^{v_2} \quad a_3 \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}^{v_3} \quad a_4 \rightarrow \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}^{v_4} \quad a_5 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}^{v_5}$$

BINARY BIPOLAR CODING (SPARSE)

(A5)

$$a_1 \rightarrow \begin{bmatrix} -1 \\ 1 \\ -1 \\ 1 \end{bmatrix}^{v_1} \quad a_2 \rightarrow \begin{bmatrix} -1 \\ 1 \\ -1 \\ -1 \end{bmatrix}^{v_2} \quad a_3 \rightarrow \begin{bmatrix} -1 \\ -1 \\ -1 \\ 1 \end{bmatrix}^{v_3} \quad a_4 \rightarrow \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \end{bmatrix}^{v_4} \quad a_5 \rightarrow \begin{bmatrix} -1 \\ 1 \\ -1 \\ 1 \end{bmatrix}^{v_5}$$

(The transformation from binary $b \in \{0, 1\}$ to bipolar is)

$$\begin{aligned} b=0 & \Rightarrow x = 2b - 1 \Rightarrow -1 \\ b=1 & \Rightarrow x = 2b - 1 \Rightarrow 1 \end{aligned}$$

CONTINUOUS EMBEDDING

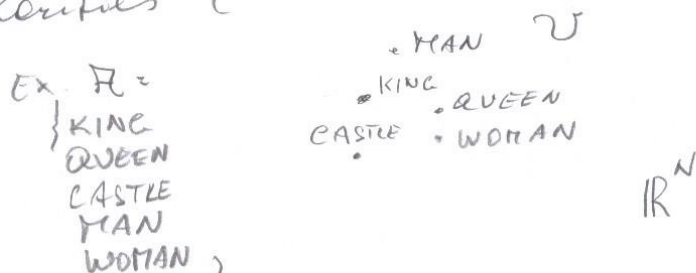
The most general case is where

$$\begin{aligned} a_1 & \rightarrow v_1 \\ a_2 & \rightarrow v_2 \\ & \vdots \\ a_n & \rightarrow v_n \end{aligned}$$

Each symbol a_i is mapped into an N -dimensional vector $v_i \in \mathbb{R}^N$.

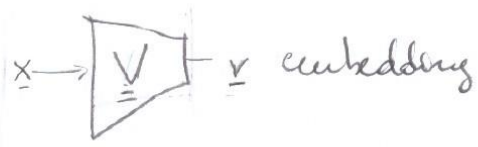
Usually the vectors are chosen in a way that reflects some kind of relation among the symbols.

For examples, in processing words, we use embeddings that tend to preserve semantic similarities (WORD2VEC)



The embeddings are usually represented as the following scheme.

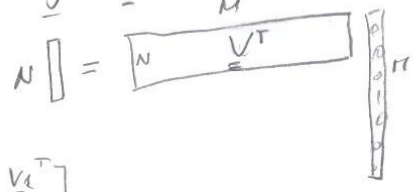
WORDS
(ONE
HOT ENCODING)



M (size of the
vocabulary)

$$y = V^T x$$

← vocabulary



$$V = \begin{bmatrix} v_1^T \\ v_2^T \\ v_3^T \\ \vdots \\ v_M^T \end{bmatrix}$$

← each row of V is
the vector for a word.

The embedding can be inserted in various architectures and be pre-learned, or learned with all the other parameters of the network.
 [see the literature for schemes to learn word embeddings]

This topic deserves a separate discussion that will be repeated elsewhere.