

01/30/2020

$$y = f(z_1, z_2, z_3)$$

true function

↑ Phenomenon      true causal drivers  
pay back mortgage (1 or 0)

$$y \in \mathcal{Y} = \{0, 1\}$$

$z_1$ : has the money (1)

$z_2$ : Unforces emergency (1)

$z_3$ : Is criminal (1)

Next Best thing: Find features that approximate the "information" in  $z_1, z_2, z_3$ .

Here are 3 that are directly related.

$x_1$ : Salary at time of application (continuous)  $\in \mathbb{R}$

$x_2$ : Miss previous payment (binary)  $\in \{0, 1\}$

$x_3$ : (Criminal intent) Do they have a record?  
(binary)  $\in \{0, 1\}$

$x_j$ 's are called "features", "characteristics",  
"attributes", "regressors", "variable",  
"covariates", "independent variables"

Let  $p$  denote the # of features

$$\text{let } \vec{x}_i := \{x_{i1}, x_{i2}, \dots, x_{ip}\} \in \mathcal{X} \text{ input space}$$

$\vec{x}_i$  is called  $i^{\text{th}}$  "observation", "subject"

"setting"

"record"

"input"

"unit"



e.g.  $\mathbb{R}^p$

Can I measure  $x_3$  better? <sup>worst crime</sup>

$x_3 \in \{ \text{none, infraction, misdemeanor, felony} \}$

categorical or factor variable with  $L=4$  levels

Mathematical models required numerical values.  
What do we do?

Two options

① Note: this is R factor validated with a monotonic order code this variable  
variable via  $x_3 \in \{0, 1, 2, 3\}$   
 $\uparrow$  home  $\uparrow$  felony

Downside :- The coding is arbitrary

\*

e.g. why not

$x_3 \in \{0, 1, 5, 100\}$

② Create multiple features

$x_{3a} \in \{0, 1\}$  is infraction?

$x_{3b} \in \{0, 1\}$  is misdemeanor?

$x_{3c} \in \{0, 1\}$  is felony?

None is captured by  $x_{3a} = x_{3b} = x_{3c} = 0$ .  
1-1 binary variables.

Consider  $x_j \in \{ \text{Red, Green, Blue} \}$

unordered factors

$y = t(z_1, z_2, z_3) = f(x_1, \dots, x_p) + f$   
b-f is the error due to ignorance

f:- is the 'best' possible way of combining  $x_1, \dots, x_p$  to minimize f



How to get  $f$ ? ; Analytical solution? = No,

Analytical solution for  $f$ .

$$f(x) = \int x dx = x^2/2$$

The approach we will use is "learning from data"  
Use data to get an estimate of  $f$ .

This procedure is also termed "supervised learning".

There are three ingredients.

① Training Data (Data)

Notation:  $\mathbb{D}$

$$\mathbb{D} = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$$

$n$ : Historical examples. (subject with response)

↓  
sample  
size

i.e. - It happened already.

$\vec{x}_1$ : Bob's features

$y_1$ : 1 (He paid it back)

$\vec{x}_2$ : Bell's features

$y_2$ : 0 (He didn't pay back)

Standard notation:-

$$X = \begin{bmatrix} \leftarrow \vec{x}_1 \rightarrow \\ \leftarrow \vec{x}_2 \rightarrow \\ \vdots \\ \leftarrow \vec{x}_n \rightarrow \end{bmatrix}$$

$n \times p$  matrix

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

col. vector length  $n$ .

$$\mathbb{D} = [X, \vec{y}]$$



②  $\mathcal{H}$  = a set of candidate function for  $f$ .

Recall  $f: x \rightarrow y$  e.g.  $F: \mathbb{R}^p \rightarrow \mathbb{R}$   
You need to simplify the set of possible function

③  $A$ : an algorithm; such that  $g = A(D, \mathcal{H})$  away  
to select/learn a model.  $g \in \mathcal{H}$  using ②