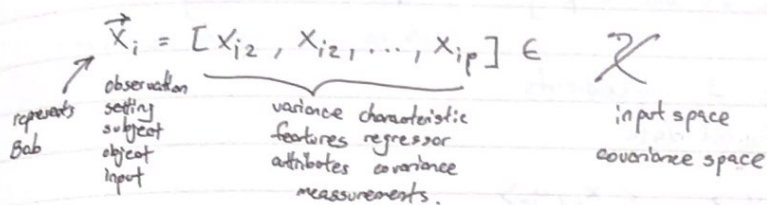


01/31/2019

Assumption:  $y: \mathcal{Z}_1, \dots, \mathcal{Z}_p$  it is not the model; it is the reality.

Next Big Thing!

Obtain  $x_1, \dots, x_p$  which hopefully... captures much information in the  $\mathcal{Z}_s$



$x_1: \text{credit score} \in \mathbb{R}$   
 continuous variables

$x_2: \text{Criminality}$   
 many metrics

$x_2 \in \{ \text{has past history, does not have} \}$   
 indicator variable  
 binary variable  
 dummy...

1                      0

$x_2 \in \{ \text{none, infraction, misdemeanor, felony} \}$   
 factor variable, categorical variable  
 with  $L=4$  levels

Two Strategies to use factor var's in math models.

(a) Ordinal encoding

$x_2 \in \{ 0, 1, 2, 3 \}$   
 ordinal factor variable.

misdemeanor felony

major downside: encoding is arbitrary

(b) Nominal encoding

$x_{2a} \in \{ 0, 1 \}$  infraction or not?

$x_{2b} \in \{ 0, 1 \}$  misdemeanor or not?

$x_{2c} \in \{ 0, 1 \}$  felony or not?

$x_{2a} = x_{2b} = x_{2c} = 0 \Rightarrow \text{"none"}$

downside: too many inputs

$p=3 \Rightarrow p=5 \quad L-1=3.$

Mathematical Model.

Can you say  $y = f(x_1, \dots, x_p)$ ?

No  $y = f(x_1, \dots, x_p) + \delta$   $\rightarrow$  error due to ignorance  
 $\delta = \hat{f}(\vec{z}) - f(\vec{z})$

in two cases even this is identical but there still thing that you don't know.

How to minimize  $\delta$ :

increase number of relevant variable.

Income?  
government employee?  
happy or not?

base on data, measurement  
observation

Find  $f$ . The approach we use is called "learning from data" an "empirical approach".  
the type of learning from data we will employ is "supervised learning"

historical data oversee the learning

Supervised Learning needs 3 ingredients

① "training data", "historical data"

one goes in something similar comes out.

$$\mathcal{D} = \{ \langle \vec{x}_1, y_1 \rangle, \langle \vec{x}_2, y_2 \rangle, \dots, \langle \vec{x}_n, y_n \rangle \}$$

$n$ : # of historical exmples (sample size).

$\vec{x}_1$  is Bob measurements  $y_1 = 1$  (he paid)

$\vec{x}_2$  " Jill "  $y_2 = 1$  (he paid)

$\vec{x}_3$  " Bill "  $y_3 = 0$  (he did not pay).

$$X = \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_n \end{bmatrix}, \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

historical data

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

Arbitrary  
know nothing  
eg.  $f: \mathbb{R}^p \rightarrow \mathbb{R}$

set of all linear function of  $x$ 's

②  $\mathcal{H}$ : a set of candidate functions  $h$  this can approximate  $f$ .  
Assumption!

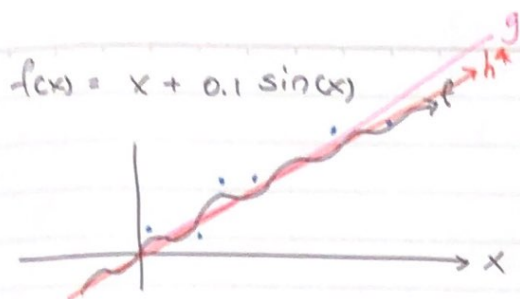
③  $A$ : an algorithm that takes  $\mathcal{H}$  and  $\mathcal{D}$  and provides  $g \in \mathcal{H}$  as the best approximation of  $f$ , which is  $h^*$ .

$$g = A(\mathcal{D}, \mathcal{H})$$

Is  $f \in \mathcal{H}$ ? Generally speaking  $\dots$  NO.

However  $\exists h^* \in \mathcal{H}$  this is the best approximation. of.

$$y = h^*(x_1, \dots, x_p) + \underbrace{f(\vec{x}) - h^*(\vec{x})}_{\text{misspecification error}} + \underbrace{f(\vec{x}) - f(\vec{x})}_{\text{error due to ignorance}}$$



the more complex is the closer is to the line.

$$\mathcal{H} = \{ \text{all linear functions of } x \}$$

$$= \{ \beta_0 + \beta_1 x : \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R} \}$$

$$\underline{h^*(x)} = x$$

want  $g_{\text{us}} = b_0 + b_1 x$

Example:

there are six data point

$$y = \underbrace{g(\vec{x})}_{\text{Model}} + \underbrace{h^*(\vec{x}) - g(\vec{x})}_{\text{estimation error}} + \underbrace{f(\vec{x}) - h^*(\vec{x})}_{\text{misspecification error}} + \underbrace{\epsilon(\vec{x}) - f(\vec{x})}_{\text{error due to ignorance}}$$

$e$  (residual).

How to predict?

For a new object  $\vec{x}_*$ , how to predict  $y$ ?

$$\hat{y} = g(\vec{x}_*)$$

How to minimize misspecification error?

Make  $\mathcal{H}$  richer, make  $A$  better?

How to minimize estimation error?

Increase  $n$  (sample size)