

10/31/ Assume

$$y = f(z_1, \dots, z_t)$$

it is not the model, it is the reality.

Next best thing

obtain  $x_1, \dots, x_p$  which hopefully  
capture much of the information in the  $z$ 's.

let  $\vec{x}_i := [x_{i1}, x_{i2}, \dots, x_{ip}] \in X$

observation, object  
setting,  
record,  
objects,

variable, features,  
regressor, measurements.

"input space"

$x_1$ : credit score  $\in \mathbb{R}$

Continuous variable

$x_2$ : criminality

many metrics.

$x_2 \in \{ \text{has past history, does not have} \}$

Binary Variable

Indicator Variable

Dummy Variable

$x_2 \in \{ \text{none, infraction, misdemeanor, felony} \}$

factor variable or categorical variable

w/  $L = 4$  levels.

now

Two strategies to use factor variables in mathematical models.

(a) ordinary encoding

$x_2 \in \{0, 1, 2, 3\}$

Ordinal factor variable

major downside: encoding is arbitrary  $\Rightarrow$  you can assign many ~~1000~~ 1, 500 etc.

## ⑥ nominal encoding

You create  $\rightarrow$   $x_{za} \in \{0,1\}$  infraction or not?  
 $x_{zb} \in \{0,1\}$  misdemeanor or not?  
 $x_{zc} \in \{0,1\}$  felony or not?

$x_{za}, x_{zb}, x_{zc} = \text{"none"}$

downside:  $p=3 \Rightarrow p=5$  1-1  
more

there are things that we can not nominalize like favorite color "your choice"

Can we say:  $y = f(x_1, \dots, x_p)$ ? <sup>not reality</sup> No

$y = f(x_1, \dots, x_p) + \underbrace{\delta}_{t(\bar{z}) - f(x)} \rightarrow$  error to ignorance  
error to thing we don't know even if  $f(x, \dots, x)$  are identical

How to minimize  $\delta$ :

increase # of relevant variables.

$\rightarrow$  after finding  $x_s$  we need

$\hookrightarrow$  Find  $f$ . The approach we use is called "learning from data" a "empirical ~~data~~ approximation". The type of ~~data~~ learning from data we will employ is "supervised learning".  
(base on measurements, data)

Historical data overseas  $\uparrow$

We learn from measurement



buy  $\rightarrow$  same amount

Supervised learning needs 3 ingredients.

① "training data", "historical data"

Examples  
of things  
come in  
and goes out  
similar

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

$n$ : # of historical examples (sample size)

$\vec{x}_1$  is Bob's measurements  $y_1 = 1$  (he paid)  
 $\vec{x}_2$  is Jill "  $y_2 = 1$  (she paid)  
 $\vec{x}_3$  is Bill "  $y_3 = 0$  (he did not paid)

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

$$D = [X; \vec{y}]^{n \times p}$$

$n \times p \Rightarrow$  has  $p$  measurements not "1"

We know nothing about  $f$

eg.  $f: \mathbb{R}^p \rightarrow \mathbb{R}$  so we need to shrink the space down

assumption, because we not know  $f$

②  $H$  = a set of candidate functions  $h$  that can approximate  $f$ .

③  $A$  = an algorithm that takes  $H$  and  $D$  and provides  $g \in H$  as the best approximation of  $f$ .  $\rightarrow$  historical data.

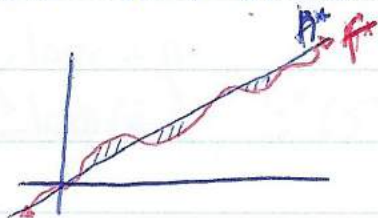
$$g = A(D, H)$$

Is  $f \in H$ ? Generally speaking ... no.

However <sup>Exists</sup>  $\exists h^* \in H$  that is the best approximation of  $f$ .

$$y = h^*(x_1, \dots, x_p) + \underbrace{f(\vec{x}) - h^*(\vec{x})}_{\text{misspecification error}} + \underbrace{f(\vec{z}) - f(\vec{x})}_{\text{error we know we don't know to ignorance.}}$$

Let  $f(x) = x + \overset{\text{scale}}{0.1} \sin(x)$



error we know we don't know to ignorance.

this error is captured in the graph we can maximize get ride

we can ~~maximize~~ get ride of this error

you did not include enough functions in  $h^*$

$$\begin{aligned} H &= \{ \text{all the linear functions of } x \} \\ &= \{ \beta_0 + \beta_1 x : \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R} \} \\ h^*(x) &= x \end{aligned}$$

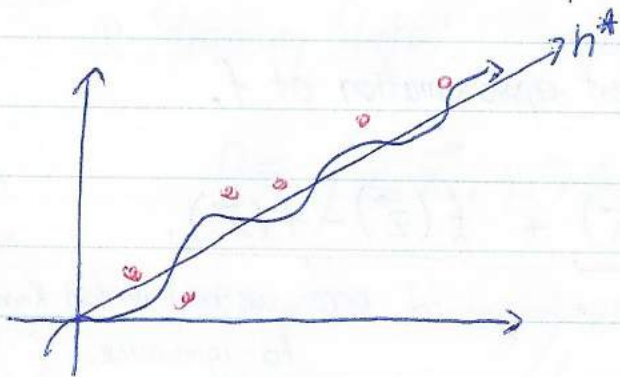
now do not have misspecification

will be good

$$\begin{aligned} f(x) &= x + 0.1 \sin(x) \\ H &= \{ \text{all linear function of } x \} \\ &= \{ \beta_0 + \beta_1 x + \beta_2 \sin x : \beta_0, \beta_1, \beta_2 \in \mathbb{R} \} \end{aligned}$$



$D =$  is 6 points in a piece paper



I want  $g(x) = \underbrace{b_0 + b_1 x}_{\text{element of } H}$

$(h^*)$  our candidate

$g \rightarrow \text{ideal}$   
 $f \rightarrow \text{ideal}$  } candidate  
 not know

$$y = g(\vec{x}) + \underbrace{h^*(\vec{x}) - g(\vec{x})}_{\text{estimation error}} + f(\vec{x}) - h^*(\vec{x})$$

$g \rightarrow$  is the only one we know.

model

$\epsilon$  (residual)

How to predict?

$\hookrightarrow$  for a new object  $\vec{x}_x$ , how to predict  $y$ ?

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

$$y = h^*(\vec{x}) + \underbrace{f(\vec{x}) - h^*(\vec{x}) + (f(\vec{z}) - f(\vec{x}))}_{\epsilon \text{ error (noise)}}$$

Residual

How to minimize misspecification Error?

Make  $H$  richer,  $A$  as well.

| ? options.

ladder notation

How to minimize estimation error?

Increase  $n$  (the sample size)

→ we can not know from where the error is coming from.

02/05

$$y = g(\bar{x}) + h^*(\bar{x}) - g(\bar{x}) +$$

$\epsilon$  (residual)

For a new observation  $x^*$ ,  $\hat{y} = g(x^*)$

$g$  comes from In supervised learning  $\xrightarrow{\text{historical data}}$   
 $g = A(D, H) \xrightarrow{\text{model space}}$   
 $\downarrow$   
algorithm

Loan Model

$y = \{0, 1\}$   $\xrightarrow{\text{pay back loan (credit)}}$

$\xrightarrow{\text{no pay back loan}}$

Model is called binary classification model if  $\{0, 1\}$ .

Null Model: You have no features and you were to create the best model  $g$ .

$H = y$   $g = \text{stochastic Model}(y)$