

## Lecture 02

Let's pretend there are 3  
casual drivers

\* For a loan

Z-1: has sufficient funds to pay  
back loan at the time it's  
due?  $Z_1 \in \{0, 1\}$

Z-2: Unforeseen emergencies?  
 $Z_2 \in \{0, 1\}$

Z-3: Criminal Intent?  
 $Z_3 \in \{0, 1\}$

$$y = f(Z_1, Z_2, Z_3) = Z_1 (1 - Z_2) (1 - Z_3)$$

Problems in practice?

- (1) You don't know  $z$ 's because they are realized in the future
- (2) You may not know the function  $f$   
which can be very complicated

What is the next best thing since you have to make a decision now and you need a model that works now?

You obtain information that approximates the information in the  $z$ 's and combine this information to approximate  $y$ . We denote these proxies that do this approximation the  $x$ 's and we denote  $P$  to be the number of such proxies:

$X_1, X_2, \dots, X_P$  For example:

$X_1$ : Salary at the time of loan application  $\in \mathbb{R}$

$X_2$ : missing payments previously  $\in \{0, 1\}$

$X_3$ : criminal charge in the past  $\in \{0, 1\} \leq P=3$

$X_j$ 's are called features, characteristics, attributes, variables, independent variables, covariates, regressors.

What is normally done in the real world?  
You use the features that are available

To learn from data you measure  $x_i$ 's  
On subjects  $i = 1 \dots n$

Let  $\vec{x}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,p}] \in X$ , the input  
space

Subjects are also called observations,  
Settings, records, objects, inputs.

|                                          |                         |
|------------------------------------------|-------------------------|
| $x_2 \in \{0,1\}$ binary variable        | } types of<br>variables |
| $x_1 \in \mathbb{R}$ continuous variable |                         |
| $x_3$ is a binary variable               |                         |

Let's consider measuring  $x_3$  differently,

$x_3 \in \{\text{none, infraction, misdemeanor, felony}\}$   
this is an ordinal category variable

How do we make this a metric?

(1) code it in order of severity spacing by 1:

$x_3 \in \{0,1,2,3\}$

Downside: coding is arbitrary

(2) binarise / dummy the categorical variable:

$X_{3a} \in \{0,1\}$  interaction or not?

$X_{3b} \in \{0,1\}$  misdemeanor or not?

$X_{3c} \in \{0,1\}$  felony or not?

One variable became 3 variables  $\Rightarrow p=5$

I had 4 levels ( $L=4$ ) but now I made

$L-1=3$  variables. Why?

You can capture the last category (called the reference category) by setting all "dummies" / binary variables to zero

If the variable is "nominal categorical" meaning no inherent order, you must do #2 to be able to use it in a model e.g.

$X \in \{\text{red, blue, green, yellow, purple, brown, ...}\}$

Can we say that  $y = f(x_1, x_2, \dots, x_p)$ ?

"NO It's only approximating at best"

- Gabriel

$y = f(z_1, \dots, z_t)$  where you don't know  $t$  or the  $z$ 's

$$y = f(x_1, \dots, x_p)$$

or

$$y = f(x_1, \dots, x_p) + \delta \quad \text{s.t.} \quad \delta = y - f$$

What is delta? It's an error. It's error due to... Ignorance. Ignorance of the true causal drivers, It's errors due to the fact that the proxies aren't the real thing. You're missing information.

How do we decrease delta? Increase  $p$  with more useful variables

How do we get  $f$ ? Note that there is no "analytical solution". The approach we use is "learning from data". This is an "empirical approach". There are many flavors. We will concentrate on "supervised learning" from "historical data". This requires three ingredients:

(1) Training data

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

these are  $n$  historical examples of inputs/outputs



$$\mathcal{D} = \langle X, \vec{y} \rangle \text{ where } X = \begin{bmatrix} \leftarrow x_1 \rightarrow \\ \leftarrow x_2 \rightarrow \\ \vdots \\ \leftarrow x_n \rightarrow \end{bmatrix},$$

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

2.  $H$  := set of candidate functions with elements  $h$  that approximate  $f$ . we need this because the space of all functions is too large and too ill-defined to directly find the "best one". You need to limit the space

3. we need  $A$  := the algorithm that takes in  $\mathcal{D}, H$  and returns  $g$ , an approximation to  $f$ ,  $g = A(\mathcal{D}, H)$

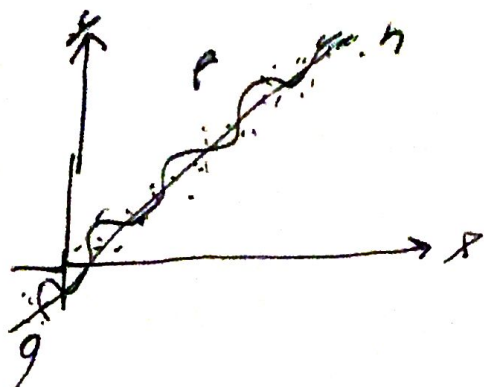
Is it true that  $f \in H$ ?

No  $f$  is arbitrarily complicated and unknown and the set curry-H contains usually simple functions that can be fit

with curly - A

However there is a  $h^* \in H$  which is the candidate model that most closely approximates  $f$ . Here is an example:

$$p=1, x \in \mathbb{R}, y \in \mathbb{R}$$



$$f(x) = x + 0.1 \sin(x)$$

$$H = \{ \text{all linear models} \}$$

$$= \{ b_0 + b_1 x : b_0 \in \mathbb{R}, b_1 \in \mathbb{R} \}$$

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

$$g \approx h^*(\vec{x}) + \epsilon$$

$$= h^*(\vec{x}) + \underbrace{[f(\vec{x}) - h^*(\vec{x})] + [h^*(\vec{x}) - f(\vec{x})]}_{\epsilon}$$

model misspecification error

\(\oint\) ignorance error

model residual

(the "full error" the difference between predicted and observed)

$$y = g(\vec{x}) + e$$

$$= g(\vec{x}) + h^*(\vec{x}) - g(\vec{x}) \quad \text{estimation error}$$

$$\left\{ \begin{array}{c} + f(\vec{x}) - h^*(\vec{x}) \\ + +(\vec{x}) - f(\vec{x}) \end{array} ; \sigma \right\} e$$

How do we decrease model misspecification error?

Expand the set of candidate functions  $H$  to be more complex and thus more expressive of complex relationships

How do we decrease estimation error?

Increase sample size  $n$  (more historical examples). The rows in  $D$

Back to the loan example where  $y$  is  $\epsilon_{0.1}$

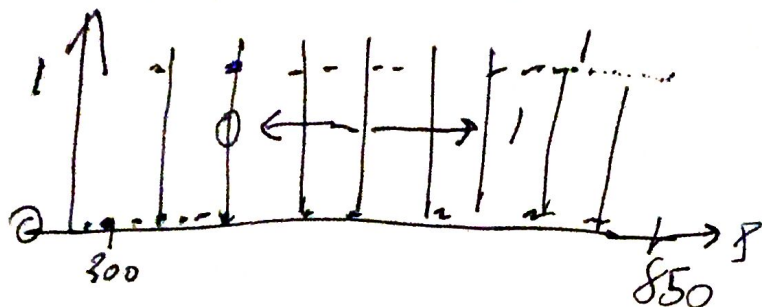
Let's say we have  $p=1$  feature, the credit score:



$X$  in  $[300, 850]$ . 50500 training data 100000s lines

$$\theta = \langle X, \vec{y} \rangle = \left( \begin{bmatrix} 810 \\ 390 \\ 750 \\ \vdots \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix} \right)$$

Let's plot the training data



What is the "null model"  $\theta_0$  which is the model if you didn't have any  $x$ 's whatsoever

$$\theta_0 = \text{mode}[\vec{y}]$$

What is the simplest candidate space  $H$ ?

$$H = \{ \mathbb{I}_{x \geq \theta} : \theta \in X \}$$

$$\text{e.g. } \theta(x) = \mathbb{I}_{x > 600}$$