

2/3/2021

lets pretend there are three causal drivers:

Z_1 : has sufficient funds to pay back loan at the time it's due?

$Z_1 \in \{0, 1\}$

Z_2 : unforeseen emergency?

$Z_2 \in \{0, 1\}$

Z_3 : criminal intent?

$Z_3 \in \{0, 1\}$

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

$$= Z_1(1-Z_2)(1-Z_3)$$

Problem in practice?

- (1) you don't know the 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 into the approximate y , we denote the 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 OR

x_2 : missing payment previously $\in \{0, 1\}$

x_3 : criminal charge in the past $\in \{0, 1\}$

$\Rightarrow p=3$ x_i 's are called features, characteristics, attributes, variables, regressors

If the variable is "nominal categorical" meaning no inherent order, you must do it to be able to model.
ex: $x \in \{\text{red, blue, green, ...}\}$

Can we say that $y = f(x_1, x_2, \dots, x_p)$? ^{No! Only approximating at best}
 $y = f(z_1, \dots, z_t)$ or $y = f(x_1, \dots, x_p) + \delta$, s.t. $\delta = y - f(x_1, \dots, x_p)$

δ is an error in the model, it's error due to ignorance of the true causal drivers, the fact that the proxies aren't the real thing. You are missing info

How do we decrease δ , 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, "Empirical approach".
There are many flavors. We will concentrate on "supervised learning" from "historical data". This requires 3 ingredients:

(1) Training Data $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$
these are n historical examples of inputs/outputs

Alternate notation:

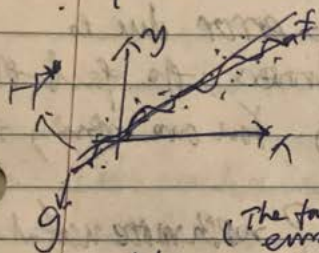
$D = \{(\vec{x}, \vec{y})\}$ where $X = \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_n \end{bmatrix}$, $\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

(2) H = a set of candidate functions with elements h that approximate f . We need this because the space of all functions is too large and ill-defined to directly find the "best" one

$$f(x_1, x_2) = \frac{1}{1 + e^{ax + bx}}$$

(3) We need A = the algorithm that takes in D, H and returns g , an approximation to f , $g = A(D, H)$.
Is it true that $f \in H$? NO, f is arbitrarily complicated and unknown and the set $\text{curly-}H$ contains mostly simple functions that can be fit with $\text{curly-}A$.
However, there is a $h^* \in H$ which is the candidate model that most closely approximates f .

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



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

$$H = \{ \text{all linear models} \} \quad \text{give more models}$$

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

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

model mispec. risk
error

$$y = h^*(\vec{x}) + \epsilon$$

$$= h^*(\vec{x}) + (f(\vec{x}) - h^*(\vec{x}))$$

$$+ (t(\vec{x}) - f(\vec{x}))$$

f(ignorance error)

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

estimation error

$$= g(\vec{x}) + [h^*(\vec{x}) - g(\vec{x})]$$

$$+ f(\vec{x}) - h^*(\vec{x})$$

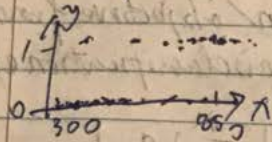
$$+ t(\vec{x}) - f(\vec{x})$$

* How do we decrease model mispec. risk?
expand the set of candidate functions H to be more complicated and thus more expressive of complex relationships

* How do we decrease estimation error?
increase sample size n (more historical examples). The more

Back to the loan example where $y \in \{0, 1\}$
 Let's say we have $P=1$ feature, the credit score:
 $x \in [300, 850]$. So your training data looks like
 $D = \{(x, y)\} = \left\{ \begin{bmatrix} 813 \\ 340 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right\}$

Let's plot the data. What is the 'null model' g_0 ?
 which is the model if you
 didn't have any x 's whatsoever.



$$g_0 = \text{Mode}[\vec{y}]$$

What is the simplest possible candidate space H ?

$$H = \{1x \geq 0: 0 \leq x\} \quad \text{e.g. } g(x) = 1x \geq 600$$

2/8/2021

null model $g_0 = \text{Mode}[\vec{y}]$ if you don't have any x 's

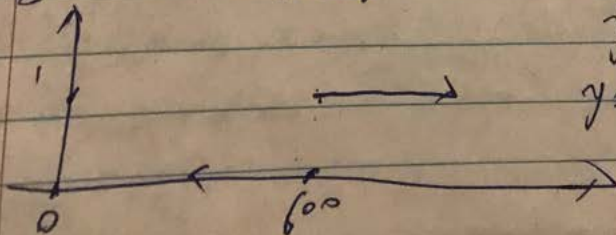
$$g_0 = \text{Mode}[\vec{y}]$$

The simplest possible candidate space H ?

$$H = \{1x \geq 0: 0 \leq x\} \quad \text{e.g. } g(x) = 1x \geq 600$$

$$H = \{1x \geq 0: \theta \in \Theta\}$$

\uparrow model parameter \rightarrow parameter space.



$$\begin{aligned} \hat{y} &= g(\vec{x}) \\ y &= g(\vec{x}) + e = \hat{y} + e \\ &= \hat{y} + (y - \hat{y}) \end{aligned}$$

Covariates

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

To learn from data, you measure the x_{ij} 's on subjects $i=1, \dots, n$.

Let $\vec{x}_i = [x_{i1}, x_{i2}, \dots, x_{ip}] \in X$, input space.

Subjects are also called: observations, settings, records, objects, input

$x_2 \in \{0, 1\}$ ^{binary variable} } types/values of variables
 $x_1 \in \mathbb{R}$ ^{continuous variable}

Let's consider measuring x_3 differently:

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

How do we make this a metric?

Method (1) Code it in order of severity spacing by 1:

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

Downside: coding is arbitrary.

Method (2) Binarize / Dummy the categorical variable.

$x_{3a} \in \{0, 1\}$ infraction or not?

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

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

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

I had 4 ^{levels} variables, but now only 3.

You can capture the last category by setting all dummy variables to zero.