Credit worthiness Model

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

we had a true system (not a model but sometimes called the "the model")

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

$z_1$: has sufficient funds

$z_2$: Unforeseen emergency

$z_3$: Criminal intentions

Problem: $\{z_1, z_2, z_3\}$ unobserved ie. impossible to obtain. What to do?

Here $\{z_1, z_2, z_3\}$

Next best thing: try to collect information "related" to $\{z_1, z_2, z_3\}$
define & ↑A

active process: you need to think about this information, how to collect it, and then actually collect it - may be expensive!

take time!

Can't always do this!!

Business as usual: use what you have (or what is easily available)

Lesson

Let's pretend we've got the resources to "define & collect".

$x_1$: Salary. How to measure: historical avg. salary

$x_2$: previous loan repayments . " " : did they ever miss a previous loan payment?
historical

$x_3$: Criminal record : " " : previous crime type?

Do $\{X_1, X_2, X_3\}$ contain the same info as $\{Z_1, Z_2, Z_3\}$? No.

Bec... $X_1, X_2, X_3$ are _possible_ to observe and $Z_1, Z_2, Z_3$ are _impossible_ to observe

$$\{X_1, X_2, X_3\} \overset{\text{has}}{\underset{\text{the same info as}}{\approx}} \{Z_1, Z_2, Z_3\}$$

Note:
row vector

Let $\vec{X} = [X_1, X_2, X_3] \in \mathcal{X}$  the "input space"
the "covariate space"

$\dim[\vec{X}] = "p"$  standard notation or "d"

an "observation", features,
an "record", attributes,
an "object", characteristics,
an "input", regressors, inputs, information, variables, dep. variable
"indep. var" covariates

What does the set $\mathcal{X}$ look like?   valid in a mathematical model? Yes!

$X_1 \in \mathbb{R}$  Why regone? Could have delta.   $X_1$ is called a "continuous variable"

$X_2 \in \{$ missed a payment, did not miss a payment $\}$

valid in a mathematical model? No!

$\Rightarrow X_2 \in \{0, 1\}$   $X_2$ is called a "binary variable" or "dummy variable"

did not miss payment ↑  did miss payment

$X_3 \in \{$ none, infraction, misdemeanor, felony $\}$

valid in a mathematical model? No!   $X_3$ is called a categorical variable.

What to do?

Two approaches for categorical variables

a) Code it with an interval order. This is known as an "ordinal factor/categorical variable"

$$X_3 \in \{ \underset{\text{none}}{\underset{\uparrow}{0}}, \underset{\text{infraction}}{\underset{\uparrow}{1}}, \underset{\text{misdemeanor}}{\underset{\uparrow}{2}}, \underset{\text{felony}}{3} \}$$

Downsides: (a) who picks the numerical ordering? Arbitrary!
(b) Is it truly ordinal?

b) Code it without an interval order. This is known as a "nominal factor/categorical variable"

How??

$X_{3a} \in \{0,1\}$  not none / none

$X_{3b} \in \{0,1\}$  not infraction / infraction

$X_{3c} \in \{0,1\}$  not misdemeanor / misdemeanor

$X_{3d} \in \{0,1\}$  not felony / felony

$X_3$ is now 4 different dummy variables.

$$p = 3 \rightarrow 6$$

This may be good or bad... we will see why later.

---

Once again... we are trying to find a model for $y$, creditworthiness for Bob.

The true model is

$y = t(z_1, z_2, z_3)$ but we cannot observe $z$'s for Bob

But we do have

$X_1, X_2, X_3$ that we observe for Bob.

Since $\{X_1, X_2, X_3\}$ has a lot of de intimeren controveil is $\{z_1, z_2, z_3\}$

Can we say $y = t(X_1, X_2, X_3)$? No... not de sure!

Can we say

$$y = f(X_1, X_2, X_3)?$$

No... we cannot use imperfect intimeren that does not exactly capture de phenomenon to explain de phenomenon precisely

Instead...

$$y \approx f(X_1, X_2, X_3)$$ ✓ some difference between approx & truth.

$$\Rightarrow y = f(X_1, X_2, X_3) + \delta$$

You cannot model $y$ exactly!

What is $f$? $f$ is de "best" funtional relationship we have.

How do we ges $f$?

Can we solve it analytically? e.g. $h(x) = x^2$ find $\min\{h(x)\}$.

the analytical sol is to take deriv. set $= 0$.

There is no analytical sol since there is no governing theory we can use to logically deduce the answer.

We can use an "empirical solution" using "historical data".

This is called "learning from data". Many flavors...

the first (most common) is "supervised learning".

"Supervised Learning" requires 3 ingredients:

① "Training data"

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

these are historical input-output examples

$\vec{x}_1$ may be Bill's characteristics where $y_1 = 1$ (ie. he paid his loan)

$\vec{x}_2$ may be Jill's " " " " $y_1 = 1$ " " "

$\vec{x}_3$ " " Tony's " " " " $y_1 = 0$ " " didn't " " " ,

⋮

there are n examples. Sometimes $D$ is denoted using vector and matrix notation:

$$D = \langle X, \vec{y} \rangle \quad \text{where} \quad X \in \mathcal{X}^n \text{ of dim } n \times p$$
$$\vec{y} \in \mathcal{y}^n = \{0,1\}^n$$

② $\mathcal{H} = \{$ all candidate functions $h$ that approximate $f \}$

Why needed? $f$ may be a very complicated function you will never be able to **learn**. So pick a large set of candidate functions that can approximate $f$.

③ $A$: an algorithm that takes in $D, \mathcal{H}$ and selects one best candidate function, $g$ $\Rightarrow g = A(D, \mathcal{H})$