Last Time: Philosophy of Modeling

Assumption:    $y = t(z_1, \ldots, z_t)$
   it is not the model, it is reality.
   But this is impossible because you don't know the z's.

Next Best Thing:
   obtain $x_1, \ldots, x_p$ which hopefully captures much
   of the information in the z's.

$$\overrightarrow{x_i} = [x_{i_1}, x_{i_2}, \ldots, x_{ip}] \in \mathcal{X}$$

BOB

observation         variables, features              "input space"
setting             attributes, characteristics
subject             Regressors                        covariate
Record              Covariates Measurements
object                                                space
Input

$X_1$: credit score $\in \mathbb{R}$
   continuous variable

$X_2$: criminality — Many Metrics
   $X_2 \in \{$ has past criminal history, does not have$\}$

indicator variable                    $\downarrow 1$              $\downarrow 0$
binary variable
dummy variable              — or —
           $X_2 \in \{$ none, infraction, misdemeanor, felony$\}$

           factor variable, categorical variable
                    with (L = 4) levels

Levels are the number of possible states
            of a factor variable.

Two strategies is use factor var's is med models :

a.) original encoding

$$X_2 \in <0, 1, 2, 3>$$

ordinal factor variable

- Major downfall : Encoding is arbitrage

b.) Nominal Encoding

$$X_{2a} \in \{0, 1\} \sim \text{infraction or not}$$
$$X_{2b} \in \{0, 1\} \sim \text{misdemeanor or not}$$
$$X_{2c} \in \{0, 1\} \sim \text{felony or not}$$
$$X_{2a} = X_{2b} = X_{2c} = 0 \Rightarrow \text{"NONE"}$$

- Downside :  $p = 3 \longrightarrow p = 5$,  $L - 1 = 3$
  ↳ More Regressors

fav. color, states, Make of a car } Things that arent nominal by design.

¿ Can you say :  $y = f(x_1, \ldots, x_p)$?   NO
if  $x$  has an approximate of  $z$  then you can't put the x's together to exactly make y.

$$y = f(x_1, \ldots, x_p) + \delta$$

error due to ignorance.
$f(z) - f(x)$

How to minimize  $\delta$ :
increase the number of relevant variables.

Find f  The approach we use is called "learning f from data"
       an empirical approach.
              based on measurements, data, observations.
       The type of data we will employ is
                   "supervised learning".
              historical data oversees the

Supervised learning needs 3 ingredients:
   ① "Training data", "Historical data"
        $D = \{< \vec{X_1}, y_1 >; < \vec{X_2}, y_2 >; \cdots < \vec{X_n}, y_n >\}$
           n: # of historical examples (sample size)

   $\vec{X_1}$ is Bob's measurements    $y_1 = 1$ (he repayed)
   $\vec{X_2}$ is Jill's measurements    $y_2 = 1$ (she repayed)
   $\vec{X_3}$ is Bill's measurements    $y_3 = 0$ (he did not repay)

$$X = \begin{bmatrix} \vec{X_1} \\ \vec{X_2} \\ \vdots \\ \vec{X_n} \end{bmatrix} \text{ by dimensions} \qquad \vec{y} \in y^n \quad \vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y \end{bmatrix}$$
           $x \in \mathcal{X}^n$   $n \times p$

$$\mathbb{D} = [X, \vec{y}]$$
   Ex: $f : \mathbb{R}^p \to \mathbb{R}$  f is an arbitrary and unknown relationship
                           between $\mathbb{R}^p \to \mathbb{R}$
   ② $\mathcal{H} :=$ a set of candidate functions h that
   Assumption: can approximate f.

   ③ $A =$ an algorithm that takes $\mathcal{H}$ and $\mathbb{D}$ and
          provides $g \in \mathcal{H}$ as the best approximation
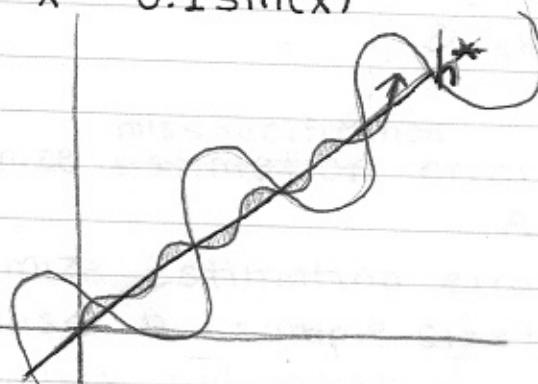          of f; which is h*.
    Is $f \in \mathcal{H}$? Generally speaking ... no.
However, $\exists$ h* $\in \mathcal{H}$ that is the best approximation of
     $y = h^*(x_1, \ldots, x_p) + f(\vec{X}) - h^*(\vec{X}) + t(\vec{Z}) - f(\vec{X})$
              Misspecification   Error due to
                 Error          ignorance

What did I misspecify?

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



your misspeeification
error just

If you make this assumption
then you make lines and are
off everywhere its not a line.

$\mathcal{H} = \{$ all linear functions of $x \}$
$= \{ \beta_0 + \beta_1 x : \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R} \}$

0 Misspecification Error → Simple
$\{ \beta_0 + \beta_1 x + \beta_2 \sin x \}$ , $\beta_0, \beta_1, \beta_2 \in \mathbb{R}$

· You have f inside h.
· your h is complex enough,
  has enough degrees of freedom.

"you want"

$$g(x) = b_0 + b_1 x$$



n = 6pts

why didnt I get $H^*$?
getting f isn't
possible cause youre
limited to lines.
getting $H^*$

$$Y = g(\vec{X}) + h^*(\vec{X}) - g(\vec{X}) + f(\vec{X}) - h^*(\vec{X}) + t(\vec{Z}) - f(\vec{X})$$

model    Estimation Error    Misspecification
                                Error

e
Residual

How to predict? For a new object $\vec{X}_*$,
how to predict y?   $\hat{y} = g(\vec{X}_*)$

$$Y = h^*(\vec{X}) + \underbrace{f(\vec{X}) - h^*(\vec{X}) + t(\vec{Z}) - F(\vec{X})}_{\varepsilon \text{ choose}}$$

How to minimize ~~estimation~~ misspecification error? Make $\mathcal{H}$ richer.

A better?

How to minimize estimation error?

Increase $n$ (sample size)

$\rightarrow$ X matrix and go like $\begin{bmatrix} \downarrow & \downarrow \\ & \end{bmatrix}$ not $\begin{bmatrix} \rightarrow \\ \rightarrow \end{bmatrix}$

get more observations.

· all of the errors appear random so you don't know where the error is coming from.