

Assumption: $y = t(z_1, \dots, z_t)$ (not model, reality)
 Assumes world is deterministic, but we can't know z 's or t

Next best thing: Obtain x_1, \dots, x_p which hopefully capture much of the info in the z 's

then $\vec{x}_i := [x_{i1}, x_{i2}, \dots, x_{ip}] \in \mathcal{X}$ "input space"
 observation setting record object \rightarrow regressors variable covariate feature attribute measurement \rightarrow Overablespace

recall x_1 : credit score $\in \mathbb{R}$
Continuous variable

x_2 : Criminality
 many metrics

$x_2 \in \{\text{has history, has no history}\}$
 Dummy/Indicator/Binary Variable \downarrow 1 \downarrow 0

Or $x_2 \in \{\text{none, infraction, misdemeanor, felony}\}$
factor/categorical variable

$L = 4$ levels

p : number of regressors

Two strategies to factor variables into mathematical models

a) ordinal encoding

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

"ordinal factor variable"

major downside: arbitrarily scaled

b) nominal encoding

$x_{2a} \in \{0, 1\}$

$x_{2b} \in \{0, 1\}$

$x_{2c} \in \{0, 1\}$

downside:

$p=3 \rightarrow p=5$

Can you say $y = f(x_1, \dots, x_p)$?

NO, y is $f(x_1, \dots, x_p) + \mathcal{J}$

$$\mathcal{J} = t(\vec{z}) - f(\vec{y})$$

Error due to Ignorance

How to minimize \mathcal{J} :

increase # of relevant variables

Find f . The approach we use is called learning from data, an empirical approach.

Based on measured data

The type of learning from data we will employ is supervised learning

historical data oversees the learning

Supervised learning in 3 steps

1) "Training data", historical data "

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

n : # of historical examples (sample size)

each $\langle \vec{x}_i, y_i \rangle$ is an input/output pair

$$X := \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_n \end{bmatrix}, \quad \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathcal{D} = [X, \vec{y}]$$

$\vec{y} \in \mathcal{Y}^n$ (output space)

Assumption

$\dim = n \times p$

2) \mathcal{H} := a set of candidate functions h that can approximate f

3) A : an algorithm that takes \mathcal{H} & \mathcal{D} & provides $g \in \mathcal{H}$ as the best approximation of f .

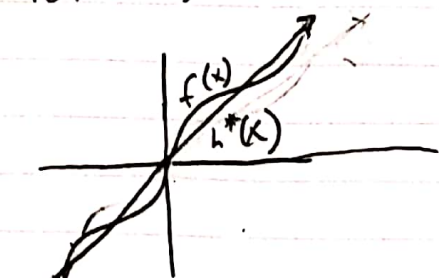
$$g = A(\mathcal{D}, \mathcal{H})$$

Is $f \in \mathcal{H}$? Generally speaking, no

However $\exists h^* \in \mathcal{H}$ that is the best approximation of f

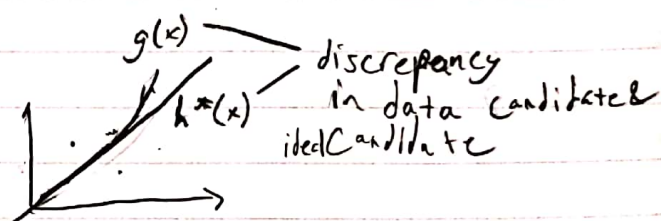
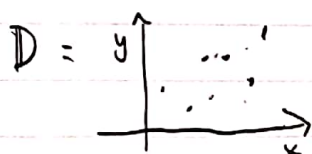
$$y = h^*(x_1, \dots, x_p) + \underbrace{f(\vec{x}) - h^*(\vec{x})}_{\text{misspecification error}} + \underbrace{t(\vec{z}) - f(\vec{x})}_{\text{ignorance error}}$$

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



machine learning minimizes \uparrow sample size minimizes

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



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

$$y = \underbrace{g(\vec{x})}_{\text{model}} + \underbrace{h^*(\vec{x}) - g(\vec{x})}_{\text{estimator error}} + \underbrace{f(\vec{x}) - h^*(\vec{x})}_{\text{misspecification error}} + \underbrace{t(\vec{z}) - f(\vec{x})}_{\text{ignorance}}$$

e (residual)

how to predict from a new \vec{x} to y ?

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

errors appear random b/c they are confounding

A tries to predict h^*

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

How to minimize misspecification error?

\mathcal{H} richer, A better (?)

How to minimize estimation error?

Increase n (sample size)