

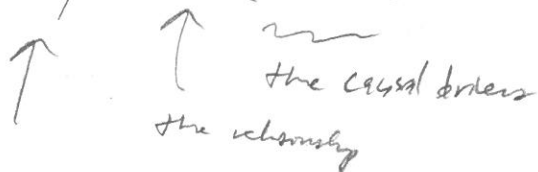
Math 390.4 11/31/19 Lec 2

4

the equals sign!

$$y = t(z_1, \dots, z_t)$$

Not a model... it's messy



the presumed phenomenon of interest

Impossible since  $z$ 's unknown.

Next best thing:  $x_1, \dots, x_p$  which hopefully captures a lot of the information in the  $z$ 's.

let  $\vec{x}_i := [x_{i1}, \dots, x_{ip}] \in \mathcal{X}$  "input space" or "common space."

observed,  
setting,  
record,  
object,  
input,  
subjects,

each variable,  
scalars,  
attributes,  
clamped variables,  
regressors,  
common

$x_1$  is called a continuous variable

$x_1$  was salary...  $\# \in \mathbb{R}$  could be negative if debt > income

yes/no...  $x_2$  reported criminal history  
 $x_2 \in \{ \text{has past crime record, does not have past crime record} \}$

$x_2$  is called a binary variable.

↓  
0

we can also define as worst crime so far committed...

$X_2 \in \{ \text{none, infraction, misdemeanor, felony} \}$   
or factor with  $L=4$  levels.

this is called a categorical variable, but statistical models need numbers. How do we make this a number? Two choices:

(a) Code it with an inherent order

$X_2 \in \{0, 1, 2, 3\}$  ordinal categorical var.  
ordinal factor var.  
ordinal var.

Major downside:

coding is arbitrary

(b) Code it without an order. How?

$X_{2a} \in \{0, 1\}$  infraction or not? nominal categorical var.  
 $X_{2b} \in \{0, 1\}$  misdemeanor ... ? nominal factor var.  
 $X_{2c} \in \{0, 1\}$  felony ... ? nominal var.

if  $X_{2a} = X_{2b} = X_{2c} = 0 \Rightarrow \text{"none"}$

price to pay  $p=3 \rightarrow p=5$  since this single feature became  $L-1=3$  binary features.

Can we say

$$y = f(x_1, x_2, x_3) \quad \text{if} \quad y = t(z_1, z_2, z_3)$$

No...  $x_1, x_2, x_3$  does not contain the perfect information of  $z_1, z_2, z_3$

Q: How can  $y$  be different if same inputs??

Instead...

$$y = f(x_1, x_2, x_3) + \delta \quad \text{or} \quad y \approx f(x_1, x_2, x_3)$$

a model

Where  $\delta = t - f$  and we call this "error due to ignorance". Why? You don't have access to all the information needed to construct  $y$  in a given setting.

How to decrease  $\delta$ ? Collect more relevant  $x$ 's!

Analytical solution? No...

How to get  $f$ ? The approach we will use is called "learning from data", an "empirical approach" <sup>from measurements</sup>.

There are many flavors. We will use "supervised learning" from historical data.

[7]

Supervised learning requires 3 ingredients

① Training data  $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$   
↳ historical samples of inputs/outputs

$\vec{x}_1$  may be Bob's information,  $y_1 = 1$  (he repaid his loan)  
 $\vec{x}_2$  ... Jill's ...  $y_2 = 1$  (she repaid her loan)  
 $\vec{x}_3$  ... Bill's ...  $y_3 = 0$  (he did not repay)  
⋮

Alternate notation  $D = \langle X, \vec{y} \rangle$  where  $X := \begin{bmatrix} \leftarrow \vec{x}_1 \rightarrow \\ \leftarrow \vec{x}_2 \rightarrow \\ \vdots \\ \leftarrow \vec{x}_n \rightarrow \end{bmatrix}$  dim  $n \times p$   
 $X \in \mathcal{X}^n, \vec{y} \in \mathcal{Y}^n$

②  $\mathcal{H} :=$  a set of candidate functions  $h$  that approx.  $f$ .  
 $f$  is an unknown and arbitrary relationship between  $\mathbb{R}^p \rightarrow \mathbb{R}$ .  
You cannot learn/estimate arbitrary functions. You need  
to simplify the candidate set.

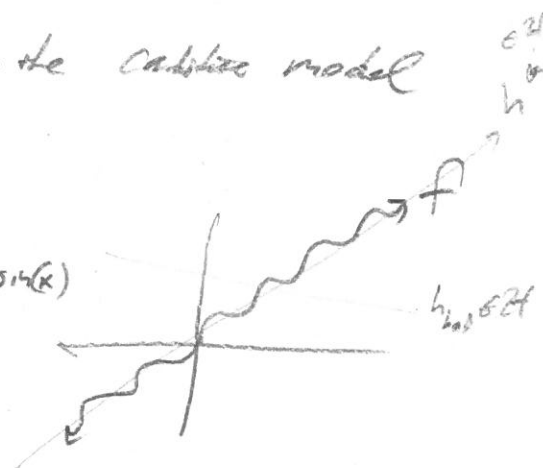
③  $A :=$  an "algorithm" that takes  $D$  and  $\mathcal{H}$  as  
input and returns  $\hat{g}$ , an approximation to  $f$ ,  
 $A = A(D, \mathcal{H})$ .

Would  $f \in \mathcal{H}$ ? Usually not!  $f$  is complex and arbitrary and  $\mathcal{H}$  contains <sup>usually</sup> simple functions that can be fit in  $\mathcal{A}$ .

However, there is  $h^* \in \mathcal{H}$  which is the candidate model the most closely approximates  $f$ .

eg.  $\mathcal{H} = \{\text{all linear models}\}$

$$f(x) = x + 0.15 \sin(x)$$



$$\Rightarrow h^*(x) = x$$

$$y = h^*(\vec{x}) + \underbrace{f(\vec{x}) - h^*(\vec{x})}_{\text{model misspecification error}} + \underbrace{\epsilon(\vec{x}) - f(\vec{x})}_{\epsilon: \text{error due to ignorance}}$$

$\epsilon$  error

How to minimize model misspecification error? Make  $\mathcal{H}$  richer, more flexible able to capture  $f$  more exactly.

Thus, because  $h^*$  exists in  $\mathcal{H}$ , doesn't mean you'll find it!  $\mathcal{A}$  won't be perfect and the  $\epsilon$ 's will appear random. So  $\mathcal{A}$  does its best to produce  $g \in \mathcal{H}$ .

Eq.  $y = f(x) + \varepsilon = x + 0.1 \sin(x) + \varepsilon$

practical  
the model we "use"



↓

$$y = g(\bar{x}) + \underbrace{h^*(\bar{x}) - g(\bar{x})}_{\text{approx. error}} + \underbrace{f(x) - h^*(x)}_{\text{misspec. error}} + \underbrace{\varepsilon(\bar{x}) - f(\bar{x})}_{\text{ignore error}}$$

$\varepsilon$   
↑  
"residual"  
(what's left over)

It has three components, the three sources of error.

How to measure error? Increase n! Better A?