

why study stat. learning?

→ Business: use data collected from daily ops. to
↑ efficiency, predict val., or explore relationships

Δ Developing accurate model to pred. sales from 3
media budgets → pg. 15

- When trying to model 1 var. w/ p other vars., we assume
a relationship b/w. target (Y) & features ($X = (x_1, \dots, x_p)$)

→ $Y = \underbrace{f(X)}_{\text{Arbitrary function}} + \underbrace{\epsilon}_{\text{random error}} \rightarrow \text{ind. from } X \text{ \& } \mu = 0$
→ what info X provides
abt. Y

What's the difference between prediction & inference?

PREDICTION

- Set of X but not their Y s
- Don't need highly interpretable models
- Predict Y given X :
 $\hat{Y} = \hat{F}(X)$

Δ Predict a patient's
risk of reaction from
blood characteristics

INFERENCE

PREDICTION

- Reducible error is error caused by est.

$f(x)$

- Irreducible error is simply caused by existence of rand. error
- Improve prediction by reducing reducible error:

↳ How? → select proper model
↳ tune model

How do we est. f ?

- ALWAYS assume that we have obs. set of n data pts.
↳ Data used to Train / teach method how to estimate f
= training data

↳ w/ multiple (p) features & n obs. we have data:

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \text{ where}$$

$$x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$$

Parametric vs Non-Parametric Methods?

- PARAMETRIC: involves 2-step model-based approach:

↳ ① Make assumption about SHAPE of f (relationship b/w

y & X) ⇒ model will have a functional form that you can write

△ Linear Regression

↳ ② After selecting a model, use some criteria that uses the training data to fit/train the model

↳ Δ Est. $\beta_0, \beta_1, \dots, \beta_p$ using OLS

PROS

- Simplifies model fitting process

CONS

- Chosen form is ^{almost} never going to be exactly right \rightarrow risk having poor estimates

How to address cons of parametric models?

- Choose flexible models that fit many functional forms

↳ BUT makes fitting model more complex \rightarrow \uparrow parameters

$\&$ \uparrow complex models risk overfitting

↳ Models works well on train but doesn't generalize well \rightarrow follow ϵ too closely

Parametric example: $\text{Income} = \beta_0 + \beta_1 \text{Education} + \beta_2 \text{Seniority}$

↳ Doesn't do a perf. job but may be good enough / reasonable

↳ May be best one can do w/ small data

NON-PARAMETRIC: **NO** explicit assumptions about form of F

↳ Instead use "steps" or series of procedures to get as close to data pts. w/o being too rough or wiggly

PROS

- No risk of getting functional form wrong

↳ Can be very accurate

CONS

- Unlike parametric models, there is no simplification \rightarrow

↑↑ complexity \Rightarrow requires

a lot of obs. to be accurate

Non-Parametric Example: Train plate spline for Income \rightarrow

tries to est. \hat{y} by getting as close as possible to data

while keeping hypersurface smooth

↳ Analyst has to decide on smoothness \Rightarrow if too high the model could fit perf. \Rightarrow leads to overfitting

PREDICTION ACCURACY VS MODEL INTERPRETABILITY

- Models can be plotted on a range of flexibility

↳ ↑ flexible: Can fit many shapes of data to est. f

△ Splines

↳ ↓ flexible: More restricted in shapes that it can est. for f

△ Linear Regression

Why ever choose a less flexible model?

① If interested in inference → restrictive models = ↑ interpretable

↳ LR allows for easy interpretation of coeff. is something like splines

② Even when only concerned w/ prediction, ↓ flexible models can outperform ↑ flexible models b/c they avoid overfitting

SUPERVISED VS UNSUPERVISED LEARNING

- Supervised learning: for each obs. of input(s) we have an associated output

↳ Goal = fit model that relates inputs to outputs ⇒ aim to accurately pred. future vals. of output OR better understand relationship between input & output

△ LR, GAMs, Splines, RF, NN, SVMs

- unsupervised learning: for every obs. $i=1, \dots, n$, we see a vec. of measurements x_i but **NO** response y_i

↳ Can't fit supervised models w/o y

→ What can we even do in an unsupervised learning situation?

- Look for relationships betw. vars. or obs.

- Common tool: Cluster Analysis

find, on basis of x_1, \dots, x_n , whether obs.

fall into DISTINCT groups

Δ (Pg. 25) Mkt. segmentation $\xrightarrow{\text{find}}$ groups & see if they have diff spend habits

- For 2D plots, can sometimes just use eyes \Rightarrow for $\mathbb{R}^n, n > 2$ plots, much harder to plot all obs. across all vars.

↳ Clustering becomes much more useful

- Sometimes have y only for fraction of n obs \rightarrow
semi-supervised learning

REGRESSION VS CLASSIFICATION

- Vars are either: (1) Quantitative: Have numerical vals.

Δ Age, height

(2) Qualitative: Have vals. in 1 of k diff
classes / categories

Δ Marital status, brand

- Regression Problems = Quantitative Response }
- Classification Problems = Qualitative Response } Not always this
clean-cut!

Δ Use LR for quant. & Log. R for qual.

↳ KNN & Boosting (along w/ some others) can be used
for BOTH

* Regression vs Classification depends on OUTPUT type NOT
input ⇒ can use quant. OR qual. for both types of models
(as long as qual. vars. are coded properly) *