Chapter 2: Statistical Learning



Credit: <u>https://www.instagram.com/sandserifcomics/</u>

statistical machine learning is more than just statistics and more than just machine learning. We choose methods based on data AND our goals.

1 What is Statistical Learning?

A scenario: We are consultants hired by a client to provide advice on how to improve sales of a product.

$\sim \times_1$	Xa	X ₃	<u> </u>
TV	radio	newspaper	sales
230.1	37.8	69.2	22.1
44.5	39.3	45.1	10.4
17.2	45.9	69.3	9.3
151.5	41.3	58.5	18.5

We have the advertising budgets for that product in 200 markets and the sales in those markets. It is not possible to increase sales directly, but the client can change how they budget for advertising. How should we advise our client?



input variables "predictors", "independent variables", "features"

```
advertising budgets

X1 - TV

X2 - radio

X3 - nouspoper

output variable "response", "dependent variable"
```

More generally - observe quartitative response & and p predictors X Xp

Assure there is some relationship between response and predictors.

$$Y = f(X) + e.$$

$$Y = f(X) + e$$

I can involve more then one variable Ce.g. TV, radio, rewspaper).

Essentially, statistical learning is a set of approaches for estimating f.

1.1 Why estimate f?

There are two main reasons we may wish to estimate f.

Prediction

and actual y

In many cases, inputs X are readily available, but the output Y cannot be readily obtained (or is expensive to obtain). In this case, we can predict Y using

prediction for
$$\rightarrow \hat{Y} = \hat{f}(X)$$
 remember error averages to 0

In this case, \hat{f} is often treated as a "black box", i.e. we don't care much about it as long as exact from not as important it yields accurate predictions for Y.

The accuracy of \hat{Y} in predicting Y depends on two quantities, *reducible* and *irreducible* error.

reducible:
$$\hat{f}$$
 is not a partial ostimute for \hat{f} but we can reduce error by using an appropriate
statistical learning method to astimute it.
irreducible: Even if \hat{f} unas estimated perfectly, we would still have some error because $\hat{f} = \hat{f}(x)$ but
 Y is a function \hat{f} e! We cannot reduce this no matter how well we estimate \hat{f} .
why? e contains unmeasured variables that might be useful in predicting?, or measurement error.
Consider an estimate \hat{f} and predictor $X(fixed)$.
expedied value \hat{f} $\implies E[(Y - \hat{f})^2] = E[(\hat{f}(x) + e - \hat{f}(x))^2]$
unique difference $\hat{f} = E[(\hat{f}(x) - \hat{f}(x))^2] + Var(e)$ unique
 $irreducible$ irreducible

We will focus on techniques to estimate f with the aim of reducing the reducible error. It is important to remember that the irreducible error will always be there and gives an upper bound on our accuracy. almost clump unknown in practice.

Inference

Sometimes we are interested in understanding the way Y is affected as X_1, \ldots, X_p change. We want to estimate f, but our goal isn't to necessarily predict Y. Instead we want to understand the relationship between X and Y.

i.e. how Y charges is a function of $X_{1,2}$ -, X_P $\Rightarrow \hat{f}$ no longer a black box! We need to know its form. We may be interested in the following questions:

- 1. Which predictors are associated ~/ the response? often only a small fractions are subtentially associated ~/ Y => identifying important poedictors can be useful
- 2. What is the relationship between persponse and predictor? possitive? negative? linear? etc.
- 3. Can the relationship between Y and each predictor be adequately simurized by a linear equation or is if more complex?

To return to our advertising data,

- How much of an, increase in sales is a growing of a given meease in TV indget?

Depending on our goals, different statistical learning methods may be more attractive.

E.g. Lincar models albert interpretable interace but may not give the most accurate particitions. highly nonlinear approaches can pointle accurate padictions but much less interpretable (informal is dually or impossible). 1.2 How do we estimate f? We have observed a different dota points is we want to estimate f w/ f Goal: capply a statistical learning method to the training data more to estimate our unknown

function f.

In other words, find a function \hat{f} such that $Y \approx \hat{f}(X)$ for any observation (X, Y). We can characterize this task as either *parametric* or *non-parametric*

Parametric

1. Make an assurption about the shape of
$$f$$
.
e.g. $f(X) = B_0 + \underline{B}_1 X_1 + \dots + \underline{B}_p X_p - \sum_{i=1}^{n} f_i$ is labour in X''
parameters.

This approach reduced the problem of estimating f down to estimating a set of *parameters*.

Why? This simplifies the problem of estimating f.

Dis a drawt age: What if the wodel we doorn is very diffect from the shape of f? Then the estimate land any poredictions) will be poor. We can try a more flexible model => more parameters and can lead to overfitting if errors in training data

eg. splines, (dr. 7).

Non-parametric

Non-parametric methods do not make explicit assumptions about the functional form of f. Instead we seek an estimate of f that is as close to the data as possible without being too wiggly.

Why?

Advontage:

- fit a wider range of possible shapes for f.
- " no restrictions on shape so can't assume wrong shape for f!

Disadvontage

- They don't reduce the problem! => need a lot of data.

1.3 Prediction Accuracy and Interpretability

Of the many methods we talk about in this class, some are less flexible – they produce a small range of shapes to estimate f.

Why would we choose a less flexible model over a more flexible one?

- If you are interested in infectice, respirative modules are more interpretable.

- Flexible methods can lead to complicated estimates of f so that it is difficult to understand how any individual predicts is related to the response.

in some setting we are about prediction => floxible model may be preferred.



Lin talk about hour to charge in Ch. S.

2 Supervised vs. Unsupervised Learning

Most statistical learning problems are either supervised or unsupervised -

Supervised for each observation i=1,-,n there is an associated response y; goal: fit modul that relates predictors Xi to response yi. Ly maybe for prediction or inference.

methods : linear regression, logistic regression, GAM, boosting, AFs, SVM etc.

Un super nied

for each observation i=1,--, n we have a vector of measurements 25: but no response zi' e.g. cancer example from Ch. 1 What's possible when we don't have a response variable?

- We can seek to understand the relatopnships between the variables, or
- We can seek to understand the relationships between the observations.
 based an observations, In-> DCn discen it hay fall into distinct groups.



Sometimes it is not so clear whether we are in a supervised or unsupervised problem. For example, we may have m < n observations with a response measurement and n - m observations with no response. Why?

Maybe its expensive to collect y but not X.

In this case, we want a method that can incorporate all the information we have.

"Semi-supervised" methods outside the scope of the class.

Within a supervised approach :

3 Regression vs. Classification

Variables can be either quantitative or categorical.

V > ore & K diffect classes. Numerical values

Examples –

Age - quantitative

Height - quantitative

Income - quantit Aire

Price of stock - quantitative

Brand of product purchased - integrical

Cancer diagnosis - Categoric

Color of cat - Categorid.

We tend to select statistical learning methods for supervised problems based on whether the response is quantitative or categorical.

However, when the predictors are quantitative or categorical is less important for this choice.