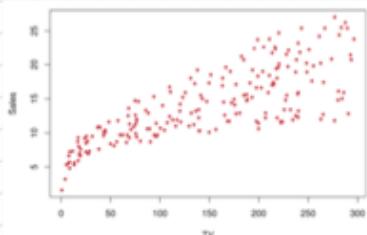


LECTURE 2 - 21/08/24



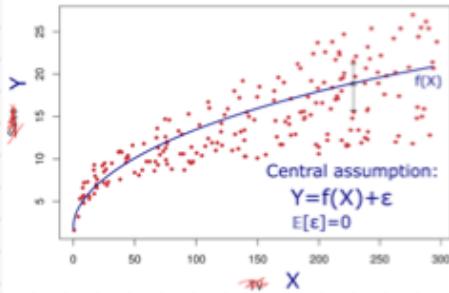
- This is a classic case of regression problem!
- The \dots are the past observations for a company
- A **question** might be, How much will I sell if I spend 50€ on TV advertising?
- ⋮

We are talking now about a linear regression problem. We see the classical input which in this case are TV ads & we wanna predict how much sales will we be able to make.

Approach:

- Basically what happens is that we wanna try to fit a **line**
- This line should be as precise as possible, which means that in all the parts of the graph it **should** be as close as possible to all the samples (\dots)

Result:



THE HARD PART OF REGRESSION RESIDES IN DECIDING WHICH FUNCTION TO CHOOSE TO COMPUTE OUR REGRESSION!

WE KNOW FROM HIGH SCHOOL THAT A LINE IS:
 $y = mx + c$

NOTE: REGRESSION NOT LINEAR R.

Regression Settings

- We assume a functional relationship between X & Y :

$$y = f(x) + \epsilon$$

- The **noise** has mean zero in this case & most of the times uncorrelated between observations, it also has constant variance.

- We can also describe this as $\rightarrow E(y|x) = f(x)$

THIS MEANS GIVEN X AS INPUT, THEN OUTPUT Y IS GIVEN BY $f(x)$

Linear Regression

- In linear regression to get back to what we said before, we choose $f(x)$ as a **linear function**:

$$y = \beta_0 + \beta_1 x + \epsilon$$

β_0 is the intercept which sometimes is called bias

This is the slope

BOTH ARE COEFFICIENTS

THE NOISE TERMS ARE ASSUMED TO BE

- INDEPENDENT
- GAUSSIAN WITH MEAN 0
- CONSTANT VARIANCE σ^2

- In matrix formulation our $f(x)$ becomes:

$$y = X\beta + \epsilon \xrightarrow{\text{Expands}} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

THE MATRIX X IS
CALLED DESIGN MATRIX

- Now what happens if $n=p=1$?

$$\begin{aligned} [y_1] &= [1 \ x_{11}] \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + [\epsilon_1] \\ y_1 &= \beta_0 + \beta_1 x_{11} + \epsilon_1 \end{aligned}$$

- Sometimes though variables are not continuous:

Has kids? YES/NO (no continuity)

Is a student? YES/NO (no continuity)

↳ This is when we THEN INTRODUCE DUMMY VARIABLES $\forall A$ OF FEATURE X

$$1_{\{X=A\}} = \begin{cases} 1, & X=A \\ 0, & X \neq A \end{cases}$$

This is called ONE-HOT ENCODING

- There are many ways of parameterising a model with a group-specific mean:

i.e.

$$y = \beta_0 1_{\{X=A\}} + \beta_1 1_{\{X=B\}} + \epsilon \quad \text{WHERE THE COEFFICIENTS ARE THE GROUP MEANS}$$

or

$$y = \beta_0 + \beta_1 1_{\{X=B\}} + \epsilon \quad \text{WHERE } \beta_0 \text{ IS THE MEAN IN GROUP A, } \beta_1 \text{ THE DIFF IN GROUP MEANS BETWEEN GROUPS B \& A}$$

- Sometimes something more complex is needed:

* not always relationships are linear in most variables x_1, x_2

* we might need a non-linear term $\rightarrow x_1 x_2$

$$\text{This means to } \rightarrow y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \epsilon$$

? now how can we RENAME this to make it look linear...

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

* we tend to call this complex terms \rightarrow INTERACTION TERMS

* the models with this term have a nice interpretation

THIS IS CALLED
TRANSFORMATION
OF FEATURES

Now consider what happens between y & x_1 for \neq values of $x_2=d$

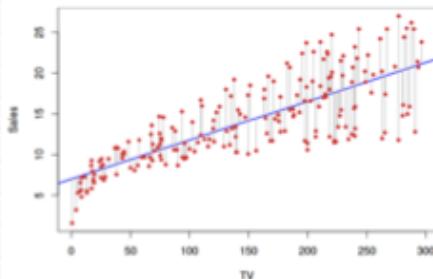
$$y = \beta_0 + \beta_1 x_1 + \beta_2 d + \beta_3 x_1 d$$

IT THEN REDUCES TO:

$$y = (\beta_0 + \beta_2 x) + \beta_1 (x_1 + \beta_3 x)$$

First Approach to model building!

- ① MAKE SOME SCATTER PLOTS AGAINST SINGLE FEATURES! (Always!)
- ② ESTIMATE THE COEFFICIENTS!



- Clearly a bad model in the picture due to the BIG DISTANCE between points & the regression line!

- * In this case we talk about ORDINARY LEAST SQUARES REGRESSION, it aims to FIND β that minimises the SUM OF SQUARED ERRORS, which is also named - the RESIDUAL SUM OF SQUARES!

$$RSS = \sum_{i=1}^n (y_i - \hat{y})^2 = \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

\hat{y} = estimated value
 y_i = the real feature
Basically, if the DIFF between the two is very HIGH, our model is BAD!
We are trying to minimize this DIFF by finding the best β 's

$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \epsilon$

ESTIMATED PARAM

\hat{y}^2 is just to REMOVE THE NEGATIVE SIGN (-)

- * In general we prefer to fit our model to MAXIMIZE the MODEL LIKELIHOOD

↳ This basically means:

- WHAT ARE THE X 'S (FEATURES) THAT WILL MORE LIKELY REPRODUCE REALITY?
- GIVEN (x_i, y_i) , the \hat{y} is as close as possible to y_i
- LIKELIHOOD IS HOW LIKELY PARAMS WILL GIVE US THE DATA WE WANT

↳ REMEMBER THAT WE ASSUME GAUSSIAN NOISE, THIS MEANS MINIMIZING RSS & USING MAXIMUM LIKELIHOOD WILL GIVE THE SAME RESULT.

↳ Maximum Likelihood starts with DEFINING A LIKELIHOOD FUNCTION:

$$p_Y(y_i | X) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{1}{2} \frac{(y_i - x_i^T \beta)^2}{\sigma^2}}$$

For a choice of β & σ^2 the likelihood of the model given all data points would be:

$$\prod_{i=1}^n p_Y(y_i | X) = \prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{1}{2} \frac{(y_i - x_i^T \beta)^2}{\sigma^2}}.$$

The goal is to find the β & σ^2 that maximize this.

Or minimize the following negative logarithm

$$\ell(\beta, \sigma^2) = -\frac{1}{2} \left(n \log \sigma^2 + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right)$$

Strategy:

① Solve for β

② Plug in $\hat{\beta}$ & maximize the profile likelihood for σ^2

* Minimizing RSS to get β

$$\text{RSS} = \sum_{i=1}^n (y_i - x_i^T \beta)^2 = (Y - X\beta)^T (Y - X\beta)$$

What we wanna do is differentiate & set equal to zero for each $\beta_r = 0 \dots p$

$$\frac{\partial \text{RSS}}{\partial \beta_r} = 2 \sum_{i=1}^n x_{ir} (y_i - x_i^T \beta)$$

Why do we differentiate?

From calculus we know that when we want to minimize or maximize a function, one way to do that is by finding that point by taking the derivative of the function with respect to the variables we are trying to optimize

Set it equal to zero to find the min which is where the function is neither up or down

Remember!!

Differentiate w/ β

③ Estimate Residual Variance

- We've seen that by fixing β at $\hat{\beta}$ & minimizing log-likelihood in σ^2 :

we get $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \hat{\beta})^2 \rightarrow$ also known as MSE

What we normally do is normalization for unbiasedness:

$$\text{RSS} / (n-p)$$

Given a large n & small p, will note see the difference

(6) Make Predictions from the Linear Model

- Given a set of features $x_i = (x_{i1}, \dots, x_{ip})$ now we can predict y_i
- We say that a natural prediction is the mean of y_i
↓
this is simply the val of the regression line at x_i
- We replace $\hat{\beta}$ with the expression we found for it using max likelihood:

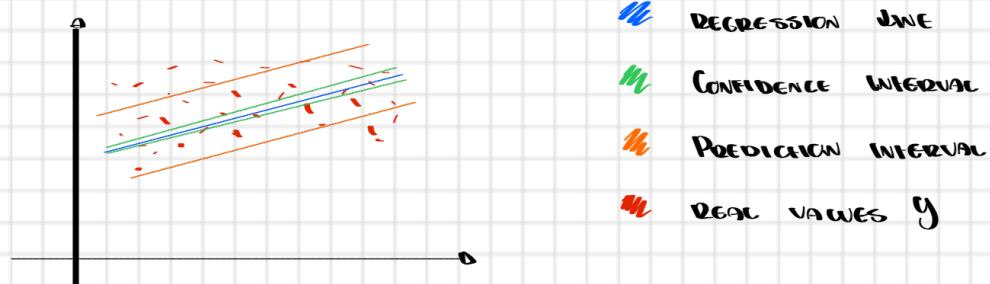
$$\hat{y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T y}_H$$

This matrix is called HAT MATRIX because "PUTS A HAT ON Y" on the transformation $y \rightarrow \hat{y}$.
It's diagonal elems h_{ii} are called LEVERAGE & are important in MODEL CHECKING

$$\begin{bmatrix} h_{11} & & \\ & h_{22} & \\ & & h_{33} \end{bmatrix}$$

(5) Confidence Intervals & Prediction Intervals

- Prediction Interval is always wider than confidence interval
 - Confidence Interval → How much the regression line will change given new data
 - Prediction Interval → Where will the predicted value \hat{y}_i fall



From the plot:

- Normally the p-value is $p=0.05$ or lower, this means:
 - We are 95% confident that for new data points, if the \hat{y}_i will change, it will stay in \parallel range.
 - We are 95% sure that the predictions \hat{y}_i of our model will fall in the range \parallel .

Balancing & Scrutinizing models! (Recap)

- Choose features to be in the model by doing Scatter plots
- Features should enter the model by:
 - TRANSFORMATIONS
 - INTERACTIONS BETWEEN FEATURES
- Fit the model.
- Do model check to see if the model fits well! (scrutinize)

SCRUTINIZE

- Scrutinize means validate / test.
- I.e. a test for $\beta_j = 0$ can be tested by comparing $\frac{\beta_j - 0}{\text{SE}(\beta_j)}$ to a t-distribution with $n-2$ degrees of freedom.
- Model output always gives the estimated SE for coefficients, $\hat{\text{SE}}(\beta_j)$.
Alternatively ... use F-stat.

This is normally used to check whether specific β 's are significantly $\neq 0$ in our linear model.

In simple terms it sees if the selected β is actually influencing a lot.

- We can also test whether several coefficients are zero

Ex.

$$M_1: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \epsilon$$

$$M_0: Y = \beta_0 + \beta_1 x_1 + \epsilon$$

We can see that this model is a special M_1 model where $\beta_2 = 0$ & $\beta_3 = 0$

* Note!

Now we have created two different models, M_0 & M_1 , they will have a different regression line which means different result, so now we will simply compare the sum of residuals (RSS) of both & we will be able to see how influential are β_2 & β_3 are.

The F-test statistic measures how much the RSS changes when we use the simpler model instead of the more complex one.

$$F = \frac{(\text{RSS}_{M_0} - \text{RSS}_{M_1})/q}{\text{RSS}_{M_1}/(n-p-1)} = \left(\frac{\text{RSS}_{M_0} - 1}{\text{RSS}_{M_1}} \right) / \frac{q}{(n-p-1)}$$

Where q is the number of parameters dropped in reducing M_1 to M_0 . Higher F if more reduction in RSS for M_1 or/and q is small compared to $n-p-1$.

- We can also have a special test where all coefficients are 0.

$$M_1: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \epsilon$$

$$M_0: Y = \beta_0 + \epsilon$$

Hypothesis Testing

- As said two nested models can be compared by F-tests.
- We might, can also compare by information criterion. \rightarrow AIC or BIC

STRATEGIES:

- FORWARD SELECTION (START BY INCLUDING NO/FEW VARS)
- BACKWARD SELECTION (START BY INCLUDING ALL/MANY VARS)
- ALTERNATING FORWARD & BACKWARD

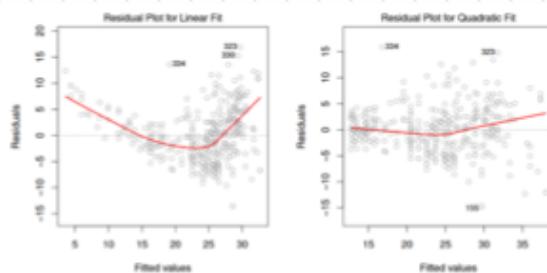
REMEMBER A

DON'T TEST TOO MUCH!
USE MODEL INSPECTION TO SEE WELL-FITTING MODEL

RESIDUALS

- WE CHECK RESIDUALS TO SEE WHETHER OUR MODEL IS WRONG & IN WHICH WAYS.
- THE (RAW) RESIDUALS ARE THE ESTIMATED ERRORS:
 $e_i = y_i - \hat{y}_i = y_i - x_i^T \beta$
- WHAT WE CAN DO IS ALSO HAVE A LOOK TO STANDARDIZED RESIDUALS THAT HAVE BEEN SCALED BY THEIR STANDARD ERROR:
$$\frac{e_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}}$$

Ex.



! WE SEE THE FITTED VALUES VS THE PREDICTED WHICH MEANS REAL VS PREDICTION

! WE SEE THE ---- AT 0 BECAUSE THE RESIDUALS ARE THE DIFF(REAL,ESTIMATED) SO THE CLOSER THIS DIFF TO 0, THE MORE PRECISE THE ESTIMATE WAS.