

ES 2.1 - CHAPTER 2, EXERCISE 7 PART A, B AND C

⑦ THE TABLE BELOW PROVIDES A TRAINING DATA SET CONTAINING SIX OBSERVATIONS, THREE PREDICTORS, AND ONE QUALITATIVE RESPONSE VARIABLE.

OBS.	$X_1$	$X_2$	$X_3$	$Y$
1	0	3	0	RED
2	2	0	0	RED
3	0	1	3	RED
4	0	1	2	GREEN
5	-1	0	1	GREEN
6	1	1	1	RED

SUPPOSE WE WISH TO USE THIS DATA SET TO

MAKE A PREDICTION FOR  $Y$  WHEN

$X_1 = X_2 = X_3 = 0$  USING K-NEAREST NEIGHBORS

a) COMPUTE THE EUCLIDEAN DISTANCE BETWEEN

EACH OBSERVATION AND THIS TEST POINT,

$$X_1 = X_2 = X_3 = 0.$$

b) WHAT IS OUR PREDICTION WITH  $k=1$ ? WHY?

c) WHAT IS OUR PREDICTION WITH  $k=3$ ? WHY?

a) DISTANCES

$$\text{OBS 1: } d_1 = \sqrt{(0-0)^2 + (3-0)^2 + (0-0)^2} = \sqrt{0+9+0} = 3$$

$$\text{OBS 2: } d_2 = \sqrt{(2-0)^2 + (0-0)^2 + (0-0)^2} = \sqrt{4+0+0} = 2$$

$$\text{OBS 3: } d_3 = \sqrt{(0-0)^2 + (1-0)^2 + (3-0)^2} = \sqrt{0+1+9} \approx 3,16$$

$$\text{OBS 4: } d_4 = \sqrt{(0-0)^2 + (1-0)^2 + (2-0)^2} = \sqrt{0+1+4} \approx 2,24$$

$$\text{OBS 5: } d_5 = \sqrt{(-1-0)^2 + (0-0)^2 + (1-0)^2} = \sqrt{1+0+1} \approx 1,41$$

$$\text{OBS 6: } d_6 = \sqrt{(1-0)^2 + (1-0)^2 + (1-0)^2} = \sqrt{1+1+1} \approx 1,73$$

b)  $K=1$

SINCE  $K=1$ , WE SELECT ONLY THE OBSERVATION WITH THE SMALLEST DISTANCE TO THE RECOMMENDED POINT.

FROM THE PREVIOUS ANSWER (a) WE KNOW THAT THE CLOSEST IS THE OBSERVATION NUMBER 5 (GREEN) = 1,41. SINCE IT IS THE CLOSEST NEIGHBOR, THE  $Y$  OF  $X_1=X_2=X_3=0$  IS GREEN.

c)  $K=3$

IN THIS CASE WE SELECT THE 3 OBSERVATIONS WITH THE SMALLEST DISTANCE TO THE RECOMMENDED POINT.

FROM THE PREVIOUS ANSWER (a) WE KNOW THAT THIS TOP3 IS COMPOSED BY: OBSERVATION 5 (GREEN), OBSERVATION 6 (RED), OBSERVATION 2 (RED).

SINCE THE MAJORITY OF THE VALUE  $Y$  IS RED, THE PREDICTED  $Y$  OF  $X_1=X_2=X_3=0$  IS RED

OBS	d	Y	
OBS 5	1,41	GREEN	{ b }
OBS 6	1,73	RED	{ c }
OBS 2	2	RED	
OBS 4	2,24	SNOW	
OBS 1	3	RED	
OBS 3	3,16	RED	

ES 2.2 - CHAPTER 3, EXERCISE 3

Suppose we have a data set with 5 predictors,  $X_1 = \text{GPA}$ ,  $X_2 = \text{IQ}$ ,  $X_3 = \text{LEVEL}$  (1 for College and 0 for High School),  $X_4 = \text{INTERACTION BETWEEN GPA AND IQ}$ ,  $X_5 = \text{INTERACTION BETWEEN GPA AND LEVEL}$ . This response is starting a salary after promotion (in thousands of dollars). Suppose we use least squares to fit the model, we get  $\hat{\beta}_0 = 50$ ,  $\hat{\beta}_1 = 20$ ,  $\hat{\beta}_2 = 0.07$ ,  $\hat{\beta}_3 = 35$ ,  $\hat{\beta}_4 = 0.01$  and  $\hat{\beta}_5 = -10$ .

2) Which answer is correct and why?

- i) For a fixed value of IQ and GPA, High School salaries from more, on avg, than College ones.
- ii) " " " " " , College salaries from more, on avg, than High School ones.
- iii) " " " " " , High School salaries from more, on avg, than college ones provided that the GPA is high enough.
- iv) " " " " " , College salaries from more, on avg, than High School ones provided that the GPA is high enough.

The salary prediction model is:  $Y = 50 + 20X_1 + 0.07X_2 + 35X_3 + 0.01X_4 - 10X_5$

$$\text{In promotion, for HS salaries: } Y_{HS} = 50 + 20X_1 + 0.07X_2 + 35 \cancel{0} + 0.01X_4 - 10 \cancel{0} = \\ = 50 + 20X_1 + 0.07X_2 + 0.01X_4$$

$$\text{And for College salaries: } Y_{COL} = 50 + 20X_1 + 0.07X_2 + 35 \cancel{1} + 0.01X_4 - 10X_5 = \\ = 85 + 10X_1 + 0.07X_2 + 0.01X_4$$

So, the salary difference between HS and Col salaries is given by this expression:

$$Y_{COL} - Y_{HS} = (85 + 10X_1 + 0.07X_2 + 0.01X_4) - (50 + 20X_1 + 0.07X_2 + 0.01X_4) = \\ = 35 - 10X_1 \rightarrow X_1 = \frac{35}{10} = 3.5 \text{ is the threshold}$$

2 options:

- 1) If  $X_1 = \text{GPA}$  is low,  $35 - 10X_1$  is positive, so Col salaries from more
- 2) If  $X_1 = \text{GPA}$  is high,  $35 - 10X_1$  is negative, so HS salaries from more

So, since HS salaries from more than Col ones only if GPA is high enough (more than 3.5), the correct answer is the number iii.

b) PREDICT THE SALARY OF A COLLEGE STUDENT WITH IQ OF 110 AND SPA OF 4,0

In this case I use the previous year model knowing that:

$$X_1 = 4,0 \quad X_2 = 110 \quad \text{and} \quad X_3 = 1$$

so,

$$X_4 = X_1 \cdot X_2 = 4 \cdot 110 = 440 \quad \text{and} \quad X_5 = X_2 \cdot X_3 = 4 \cdot 1 = 4$$

thus,

$$\begin{aligned} Y_{\text{COL}} &= 50 + 20 \cdot 4 + 0,07 \cdot 110 + 35 \cdot 1 + 0,01 \cdot 440 - 10 \cdot 4 = \\ &= 50 + 80 + 7,7 + 35 + 4,4 - 40 = \\ &= 137,1 \end{aligned}$$

which means that this college student will earn a salary equal to 137'100 dollars

c) True or False: since the coefficient for the SPA-IQ interaction term is very small, there is very little business or no interaction effect. Justify your answer.

This is False, because even if the  $\beta_4 = 0,01$  is very small, the  $X_4$  can have high values, and as we saw in the answer (b) the end result of this interaction is still relevant. To understand better if this interaction has an effect on the model we should study its significance with the p-value ( $p\text{-value} < 0,05 = \text{significant}$ ).

### ES 2.3 - CHAPTER 4, EXERCISE 6

Suppose we collect data from a group of students in a statistics class, with variables:

$X_1$  = Hours Studied,  $X_2$  = average SPA, and  $Y$  = grade in A.

We fit a logistic regression and predict coefficients constant,  $\beta_0 = -6$ ,  $\beta_1 = 0,05$ ,  $\beta_2 = 1$ .

- a) Estimate the probability that a student who studies for 40 h and has an average SPA = 3,5 gets an A in the class.

Firstly, the Logistic Regression Model, given the constants, is:

$$P(Y=1) = \frac{e^{-6 + 0,05X_1 + 1X_2}}{1 + e^{-6 + 0,05X_1 + 1X_2}}$$

and knowing that  $X_1 = 40$  and  $X_2 = 3,5$  it becomes:

$$\begin{aligned} P(Y=1) &= \frac{e^{-6 + 0,05 \cdot 40 + 1 \cdot 3,5}}{1 + e^{-6 + 0,05 \cdot 40 + 1 \cdot 3,5}} = \frac{e^{-6+2+3,5}}{1 + e^{-6+2+3,5}} = \frac{e^{-0,5}}{1 + e^{-0,5}} \approx \\ &\approx \frac{0,606}{1 + 0,606} = 0,377 \rightarrow |P = 37,7\%| \end{aligned}$$

- b) How many hours would the student in point (a) need to study to have a 50% chance of getting an A in that class?

In this case we need to find  $X_1$  s.t.  $P(Y=1) = 0,5$ .

$$\text{so, } \frac{e^{-6 + 0,05X_1 + 1 \cdot 3,5}}{1 + e^{-6 + 0,05X_1 + 1 \cdot 3,5}} = 0,5 \rightarrow e^{(\dots)} = 0,5 (1 + e^{(\dots)}) \rightarrow$$

$$\rightarrow e^{-6 + 0,05X_1 + 3,5} - 0,5 \cdot e^{-6 + 0,05X_1 + 3,5} = 0,5$$

$$\rightarrow 0,5 \cdot e^{-6 + 0,05X_1 + 3,5} = 0,5^2 \rightarrow \ln(e^{(\dots)}) = \ln(1) = 0$$

$$\rightarrow -6 + 0,05 \cdot X_1 + 3,5 = 0 \rightarrow 0,05X_1 = 2,5 \rightarrow X_1 = \frac{2,5}{0,05} = 50$$

The student will need to study for 50 hours.

## EJ 2.4 - Topic DM

A RETAILER WANTS TO MANUFACTURE PERSONAL DISTINGUISH BRIEFS CUSTODIAN YOUNGER THAN 35 AND OLDER THAN 35. THE FOLLOWING TABLE SUMMARIZES THE DATA SET IN THE DATA BASE OF THE RETAILER IN AN ABSTRACT FORM.

		# OF INSTANCES	
A	B	Y	O
Q <sub>1</sub>	b <sub>1</sub>	4	10
Q <sub>2</sub>	b <sub>1</sub>	6	2
Q <sub>3</sub>	b <sub>1</sub>	8	6
Q <sub>1</sub>	b <sub>2</sub>	2	8
Q <sub>2</sub>	b <sub>2</sub>	6	2
Tot: 26		Tot: 28	

THE RELEVANT ATTRIBUTES, DETERMINED BY DOMAIN KNOWLEDGE, ARE FOR CONVENIENCE DENOTED BY A

AND B. THE VALUE FOR A ARE Q<sub>1</sub>, Q<sub>2</sub> AND Q<sub>3</sub>.

THE VALUE FOR B ARE b<sub>1</sub> AND b<sub>2</sub>.

THE RETAILER WANTS TO USE THE DM TECHNIQUE

TO CLASSIFY THE CUSTOMER IN THE CLASS "YOUNG",  
DENOTED BY "Y", AND "OLD" DENOTED BY "O".

- a) Assume a new customer enters the web-store and the retailer has no information at all about this customer. How will this new customer be classified based on the above data and explain why.

In this case, since no information about this is given, the customer is given, the retailer should try to predict if his class is "Y" or "O" using the data in this table in this way:

$$\text{Total } Y = 4 + 6 + 8 + 2 + 6 = 26$$

$$\text{Total } O = 10 + 2 + 6 + 8 + 2 = 28$$

Since "O" has more instances, the probable prediction for the new customer is "O".

- b) Now assume a new customer comes in for which the retailer knows the value for attribute A = Q<sub>3</sub> and B = b<sub>2</sub>. Is it possible to apply the standard (van Natta) Bayes to classify this new customer? Explain what the problem is.

In this case, to apply the standard Bayes we need to calculate the probability of the customer belonging to both "Y" and "O" knowing that A = Q<sub>3</sub> AND B = b<sub>2</sub>.

RULE OF BAYES:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

So:

$$P(Y|a_3, b_2) = \frac{P(a_3, b_2 | Y) \cdot P(Y)}{P(a_3, b_2)} ; \quad P(O|a_3, b_2) = \frac{P(a_3, b_2 | O) \cdot P(O)}{P(a_3, b_2)}$$

With:  $P(a_3, b_2 | Y)$ ,  $P(a_3, b_2 | O)$  and  $P(a_3, b_2)$  unknown.

UNFORTUNATELY IN THIS CASE WE CANNOT ANSWER THE SUMMARY QUESTIONS IN THE TABLE

TABLE IT THERE IS NO DATA WITH A COMBINATION OF  $a_3$  AND  $b_2$  IN THE SAME ROW.  
IN FACT  $P(A=a_3, B=b_2) = 0$ .

c) Hence the relevant numbers in USA are new  $B_{13}$  years (new birth rate) for the combination of this new customer. How will this customer now be classified based on their values of  $A$  and  $B$ ? Explain your answer.

Using the same probability arguments of the answer (b), and knowing that the denominator  $P(a_3, b_2)$  is the same for both  $P(Y|a_3, b_2)$  and  $P(O|a_3, b_2)$  we can ignore them (but following the "minimum likelihood" idea, we computed it).

So, solving now, we can say that  $a_3$  and  $b_2$  are reasonable improvements which manner:

$$P(a_3, b_2 | Y) = \underbrace{P(a_3 | Y)}_{\text{known from TABLE}} \cdot \underbrace{P(b_2 | Y)}_{\text{known from TABLE}} = \frac{8}{26} \cdot \frac{2+6}{26} = \frac{8}{26} \cdot \frac{8}{26} = \frac{16}{169} = 0,0946$$

$$P(a_3, b_2 | O) = \underbrace{P(a_3 | O)}_{\text{known from TABLE}} \cdot \underbrace{P(b_2 | O)}_{\text{known from TABLE}} = \frac{6}{28} \cdot \frac{8+2}{28} = \frac{6}{28} \cdot \frac{10}{28} = \frac{15}{196} = 0,0765$$

Also

$$\left| \begin{array}{l} P(a_3, b_2) = P(a_3, b_2 | Y) \cdot P(Y) + P(a_3, b_2 | O) \cdot P(O) = \\ = P(a_3 | Y) \cdot P(b_2 | Y) \cdot P(Y) + P(a_3 | O) \cdot P(b_2 | O) \cdot P(O) = \\ = 0,0946 \cdot \frac{26}{54} + 0,0765 \cdot \frac{28}{54} = 0,0455 + 0,0396 = 0,0851 \end{array} \right. \quad \text{DEMONSTRATION}$$

So,

$$P(Y|a_3, b_2) = \frac{8}{26} \cdot \frac{8}{26} \cdot \frac{26}{54} = \frac{16}{351}$$

$$\text{SINCE } \frac{16}{351} > \frac{15}{378}$$

$$P(O|a_3, b_2) = \frac{6}{28} \cdot \frac{5}{28} \cdot \frac{28}{54} = \frac{15}{378}$$

THE PROBABILITY OF "Y" IS HIGHER, SO  
THE NEW CUSTOMER IS CLASSIFIED AS "Y"

ES 2.5 - TOPIC UM

Consider the same dataset as in the previous question. Now this is a classification (discrete) units to use decision trees to classify new customers.

a) What is the classification error rate for attribute A?

A	Y	O	LABEL	# OF ERRORS
a <sub>1</sub>	6	18	O	6
a <sub>2</sub>	12	4	Y	4
a <sub>3</sub>	8	6	Y	6

$$\text{Class Error Rate (A)} = \frac{6+4+6}{54} = \frac{16}{54} = 0,296$$

b) What is the classification error rate for attribute B?

B	Y	O	LABEL	# of ERRORS
b <sub>1</sub>	18	18		18
b <sub>2</sub>	8	10		8

$$\text{Class Error Rate (B)} = \frac{18+8}{54} = \frac{26}{54} = 0,481$$

c) What will be the splitting attribute in the top (root) of the decision tree if we use the classification error rate?

Since the choice is based on the attribute with the lowest error rate,

the Root of the decision tree will be A ( $0,296 < 0,481$ ).

d) What is the Gini index for attribute A?

$$\text{Gini index for attribute A: } G_I(\text{node m}) = \sum_{n=1}^k \hat{p}_{mn} (1 - \hat{p}_{mn}), \quad \text{where}$$

$\hat{p}_{mn}$  is the probability of class n observations in node m.

$$G_I(a_1) = \frac{6}{6+18} \cdot \frac{18}{6+18} + \frac{18}{6+18} \cdot \frac{6}{6+18} = 2 \cdot \frac{6}{34} \cdot \frac{18}{34} = \frac{27}{72} = 0,375$$

$$G_I(a_2) = 2 \cdot \frac{3}{34} \cdot \frac{31}{34} = \frac{3}{8} = 0,375$$

$$G_I(a_3) = 2 \cdot \frac{4}{34} \cdot \frac{3}{34} = \frac{24}{49} = 0,490$$

$$\text{Overall } G_I(A) = \frac{12}{54} \cdot 0,375 + \frac{27}{54} \cdot 0,375 + \frac{7}{54} \cdot 0,490 = 0,4048$$

e) What is the SINI index for attribute B?

$$S1(b_1) = 2 \cdot \frac{18}{36} \cdot \frac{18}{36} = 0,5$$

$$S1(b_2) = 2 \cdot \frac{8}{18} \cdot \frac{10}{18} = \frac{40}{81} = 0,4938$$

$$\text{Overall } S1(B) = \frac{2}{9} \cdot 0,5 + \frac{1}{3} \cdot 0,4938 = \frac{2}{3} \cdot 0,5 + \frac{1}{3} \cdot 0,4938 = 0,498$$

f) What will be the splitting attribute in the top (Root) of the decision tree if we use the SINI index?

Since the chart is based on the attribute with the lowest Error rate, the

Root of the decision tree will be A ( $0,4048 < 0,498$ ).

g) Construct the full decision tree, using that Error rate or accuracy, and what is the number changing error rate on this about dataset?

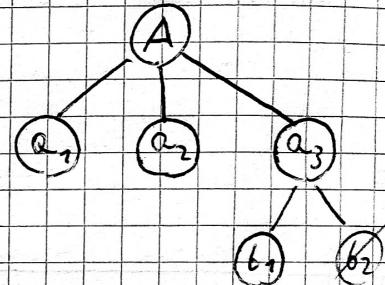
From answer (a) and (b) the Root is A

Then we split A into its 3 branches:  $a_1, a_2, a_3$

$$\text{Err Rate } (a_1) = \frac{6}{24} = 0,25$$

$$\text{Err rate } (a_2) = \frac{4}{16} = 0,25$$

$$\text{Err rate } (a_3) = \frac{6}{14} = 0,43$$



Since this  $a_3 = 0,43$  is slightly higher than the previous two, we decide that  $a_1 \rightarrow 0$ ,  $a_2 \rightarrow Y$  and  $a_3 \rightarrow$  split into b<sub>1</sub> and b<sub>2</sub> (but we have only b<sub>1</sub>).

$$\text{Err rate } (a_3, b_1) = \frac{6}{14} = 0,43 \quad \text{so we change this from } Y \text{ to } B \text{ (B > 6).}$$

$$\text{Overall changing Err rate} = \frac{6+4}{54} = 0,203$$

h) Is this changing Err rate an optimistic or pessimistic estimate of the Err rate on unseen data? Explain your answer.

Since the accuracy is  $1 - 0,203 = 0,797$ , we can say that it is an optimistic estimator.

## \* TOPIC DEP

PAPER AND PBN EXAMS: ES 1.3.2 - ASSIGNMENT 1: FACTS AND DIMENSIONS

ES 1.3.5 - ASSIGNMENT 4: "MOBILE APP B-TUTORIAL SERVICE"

PROGRAMMING EXAMS:

ES 1.3.3 - ASSIGNMENT 2: A SIMPLE ETL PROGRAM TO WORK WITH

ES 1.3.4 - ASSIGNMENT 3: AN IT SURVEY

### ES 1.3.2 - ASSIGNMENT 1: FACTS AND DIMENSIONS

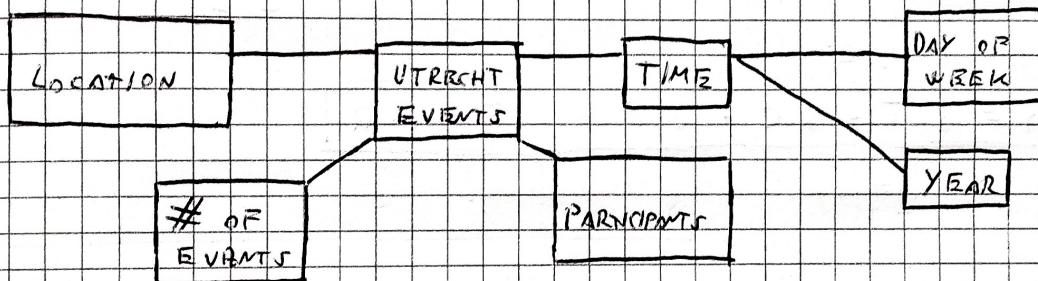
SEE FIGURE 1.1 WHICH CONTAINS DATA ON EVENTS HAPPENING IN THE CITY OF UTRrecht. WE ILLUSTRATE THE MEANING OF THE DATA BY EXPLAINING WHAT THE FIRST ROW MEANS: THERE WERE IN TOTAL 2 EVENTS ORGANIZED IN EXHIBITION ON A FRIDAY IN 2015 WHICH BOTH TOGETHER DRAW 15'000 PARTICIPANTS.

- a) WHICH ATTRIBUTES ARE THE FACTS; WHICH ATTRIBUTES ARE THE DIMENSIONS?

FACTS (NUMERIC DATA): "NO OF EVENTS" AND "TOTAL PARTICIPANTS"

DIMENSIONS (DESCRIPTIVE DATA): "DAY OF WEEK", "YEAR", AND "LOCATION".

- b) DRAW A CORECTIVE STAR SCHEMA FROM THIS CUBE



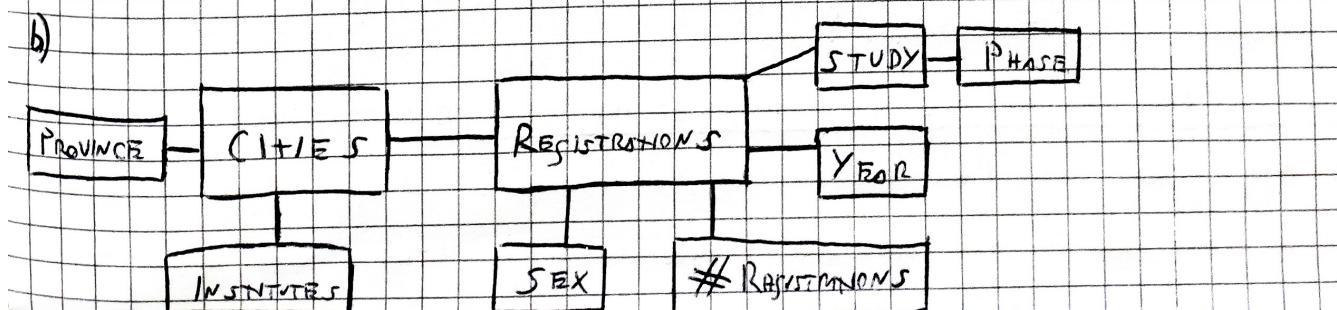
SEE FIGURE 1.2 WHICH CONTAINS DATA ON THE NUMBER OF REGISTERED STUDENTS IN THE NETHERLANDS. WE ILLUSTRATE THE MEANING OF THE DATA BY EXPLAINING WHAT THE FIRST ROW OF TABLE "REGISTRATIONS" MEANS:

IN 2015 THERE WERE 13 FEMALE STUDENTS REGISTERED FOR THE BACHELOR STUDY "TECHNICAL COMPUTER SCIENCE" AT THE UT WHICH IS LOCATED IN THE CITY OF EINDHOVEN IN THE PROVINCE OF OVERIJssel.

- 2) FACTS (NUMERIC DATA): "# OF REGISTRATIONS"

DIMENSIONS (DESCRIPTIVE DATA): "INSTID", "STUDY", "PHASE", "YEAR" AND "SEX" (AND "CITYID").

- b)



## ES 1.3.5 - Assignment 4: Multinational Marketing. Case "Mobile App Beta Tester service"

a) For Beta testing, one needs to determine the right factor(s) for the Star Schema. An important question is: "What is good", i.e., "What makes a participant seem a good beta tester for a particular game or game?". More concretely, how can you determine a result that summarizes two "Scorecard" or "suitability", which obviously refers to a calculation from the available data. Propose a formula for "Scorecard scale" and explain why a high value is an indication for a good beta tester.

- To answer the question "What makes a participant seem a good beta tester for a particular game or game?" we thought about these characteristics / requirements:
- 1) Plays the same for a given amount of time ( $T$ ) [Time spent]  $\rightarrow$  the higher, the better
  - 2) Plays consistently during the testing period ( $C$ ) [Consistency]  $\rightarrow$  " " , " "
  - 3) Is he/she good at playing it? ( $P$ ) [Pro - level]  $\rightarrow$  " " , " "
  - 4) Provides meaningful feedbacks? ( $F$ ) [Feedback]  $\rightarrow$  " " , " "
  - 5) ~~No~~ of matches/duels/games quitted in the testing period. ( $Q$ ) [Quitting rate]

So, a possible formula could be:

$$| \xi = a \cdot T + b \cdot C + c \cdot P + d \cdot F - e \cdot Q |$$

with  $a, b, c, d, e$  as coefficients of the requirements.

b) What is the business question, on what are the business questions in this case?

Formulate them as accurately as possible.

1) How can Peer find the perfect beta tester for a specific game?

$\rightarrow$  Should Peer choose them between Streamers/Youtubers, or also between common people who enjoy that game's service? Only 1 tester or more? From which country/countries? How can Peer know that the tester's feedbacks are trustworthy?

2) How can Peer measure the effectiveness of its service?

$\rightarrow$  Which data can Peer collect during this process? Since Peer should all the data to the developer? How to deal with the privacy of the tester? Which KPIs should Peer use?

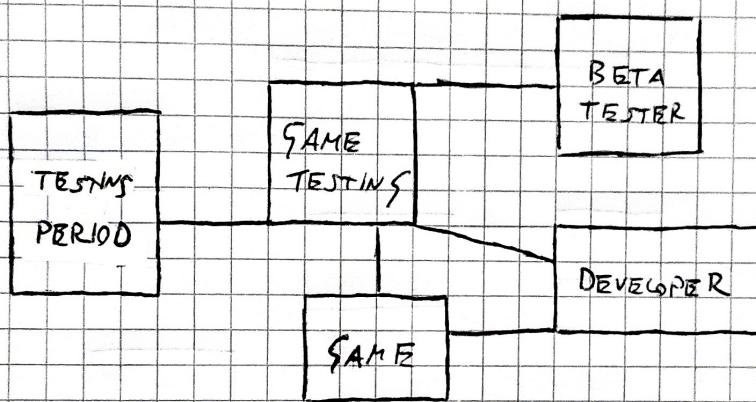
### 3) How can Plan Monitorize its Service?

→ WHICH IS THE BEST PRICING STRATEGY? HOW MUCH SHOULD THIS DEVELOPER PAY FOR THE SUBSCRIPTION, AND HOW MUCH FOR HAVING THE TESTER?  
 MORE TESTERS = MORE EXPENSES? OR MAYBE A "TESTER STAFF" DISCOUNTS PRICE?  
 WILL LOYAL CUSTOMERS (DEVELOPERS) RETURN DISCOUNTS IN SUBSEQUENT YEARS?

### 4) How can Plan Optimize/Improve its Service?

→ WHICH DATA CAN BE USED TO FIND BETTER TESTER-DEVELOPER MATCHES?  
 SHOULD PLAN TAKE CARE OF DEVELOPERS FEEDBACKS BEFORE/AFTER THE GAME?  
 SHOULD PLAN USE TEXTUAL COMMENTS IN AUGMENT THE SCORE LEVEL, OR AN INTERNAL ITEM IS ENOUGH?

### c) Give a Star Schema. Explain your design by detailing the most important design choices and considerations



THE CENTRE OF THE STAR SCHEMA IS THE "GAME TESTING" TABLE WHICH CONTAINS THE DATA IN THE INFO OF A TESTING SESSION (TEST ID = PRIMARY KEY). IT IS CONNECTED TO 4 OTHER TABLES THROUGH 4 FOREIGN KEYS: BTESTERID, DEVELOPERID, SAMEID, PERIODID.

THE "BETA TESTER" TABLE CONTAINS ALL THE INFO ABOUT THE B TESTER (e.g.: ID, NAME, AGE, GENDER, FEEDBACKS RATE...).

THE "DEVELOPER" TABLE CONTAINS ALL THE INFO ABOUT THE DEVELOPER (e.g.: ID, COMPANY NAME, # OF GAMES RELEASED...).

THE "SAME" TABLE CONTAINS ALL THE INFO ABOUT THE SAME (e.g.: ID, GENDER, # OF LEVELS, MULTIPLAYER, DEVELOPERID [FK]...).

THE "TESTING PERIOD" TABLE CONTAINS ALL THE INFO ABOUT THE TIME (e.g.: ID, STARTING DATE, ENDING DATE, YEAR...).