

DATA INVESTIGATION

APPROACHES

(Adapted from Data Science
for Business by Provost and
Fawcett, pp. 20 - 23)

Data Investigation ①

DATA = TABLE OF VALUES

Telecom Churn data set

	Senior Citizen	Tenure (#months)	Churn
Customer 1	0	1	No
Customer 2	0	34	No
.	.	.	.
.	.	.	.
Customer 7043	1	23	Yes

Data Investigation (2)

PATTERN 1: REGRESSION

	Senior Citizen	...	Features	Churn
customer 1	0	1		No
customer 2	0	34		No
.	:	:		:
:	:	:		:
Customer 7043	1	23		Yes

Predict the tenure of a customer given the values of all the other features for that customer.

Data Investigation ③

PATTERN 2: CLASSIFICATION

	Senior Citizen	...	Tenure (#months)	Churn
Customer 1	0	1	No	
Customer 2	0	34	No	
.
.
Customer 704	1	23	Yes	

Predict if a customer will churn or not given the values of all the other features for that customer.

Data Investigation (4)

PATTERN 3 : SIMILARITY

Features

	Senior Citizen	Tenure (#months)	Churn
Customer 1	0	1	No
Customer 2	0	34	No
:	:	:	:
Customer 7043	1	23	Yes

→ Similarity across attributes

Given a set of features and a value or range of values for each feature, which customers are similar?

Data Investigation ⑤

PATTERN 4: CLUSTERING

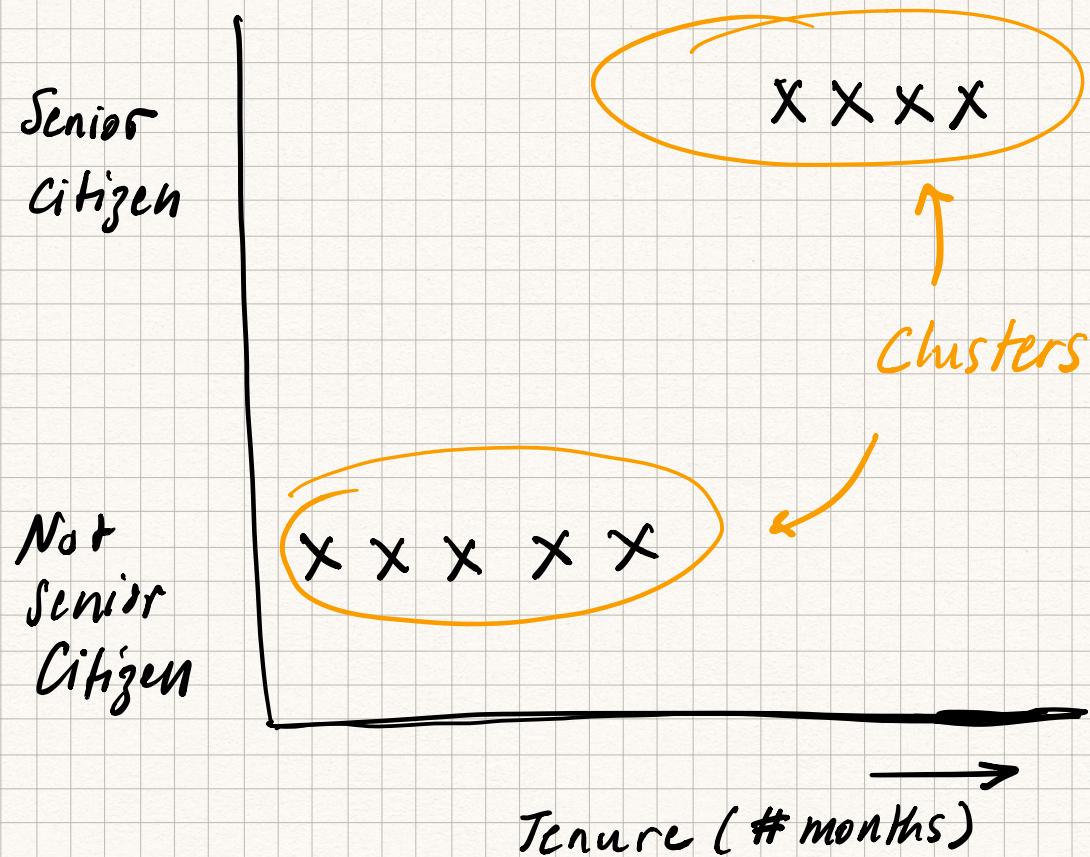
	Senior Citizen	Tenure (#months)	Churn
Customer 1	0	1	No
Customer 2	0	34	No
:	:	:	:
Customer 7043	1	23	Yes

Do customers cluster together?

(on some or all attributes)

Data Investigation ⑥

PATTERN 4: CLUSTERING EXAMPLE



Data Investigation (7)

PATTERN 5: CO-OCCURRENCE (ACROSS ATTRIBUTES)

	Senior Citizen	Tenure (#months)	Churn
Customer 1	0	1	No
Customer 2	0	34	No
:	:	:	:
Customer 7043	1	23	Yes

Features

- Which feature values always/often occur together?
- Which feature values never/rarely occur together?

Data Investigation ⑧

PATTERN 5: CO-OCCURRENCE (IN THE SAME ATTRIBUTE)

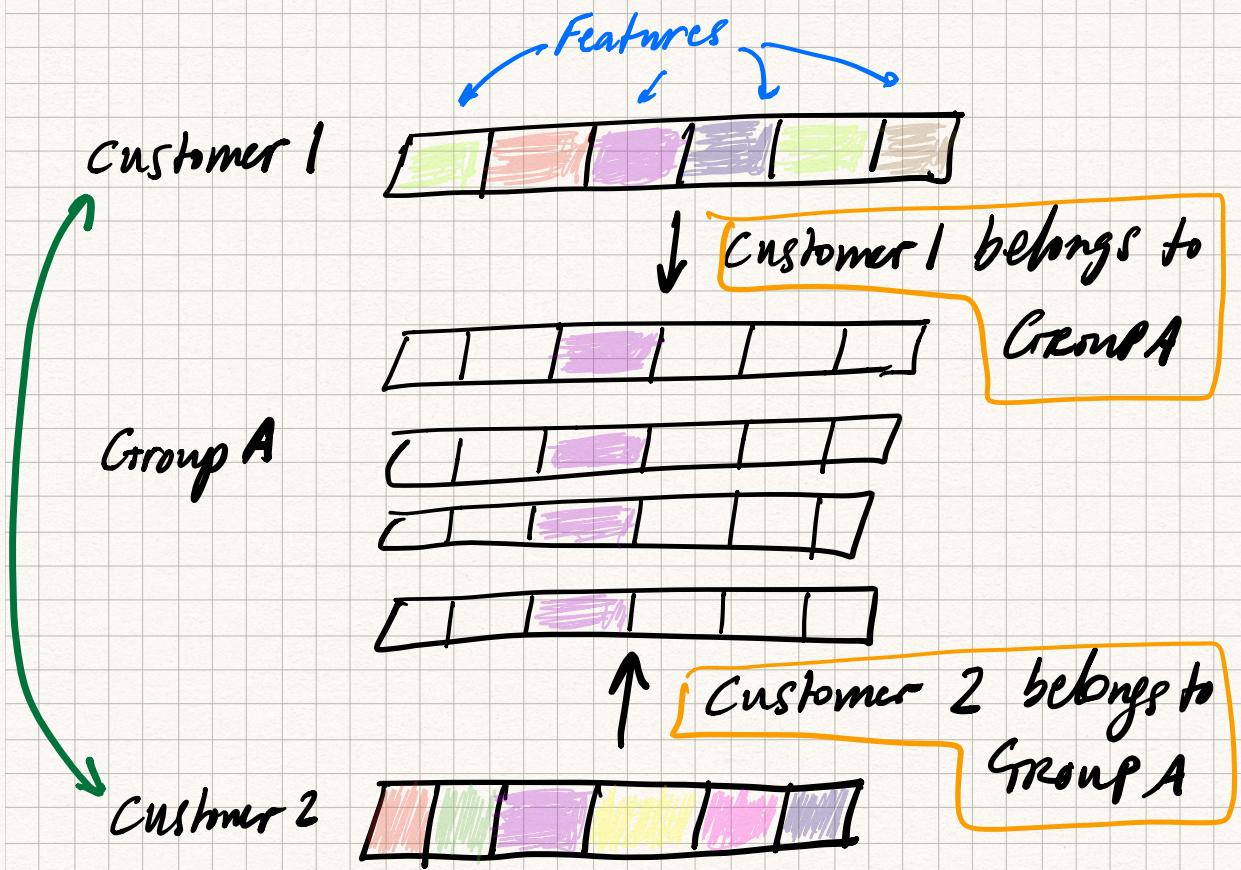
	Senior Citizen	Friends With	Tenure (#months)	Churn
Customer 1	0	16, 22, 45	1	No
Customer 2	0	16, 22, 260	34	No
.
Customer 7043	1	12, 4264	23	Yes

Customers 1 and 2 have 2 friends in common.

- which values occur always/often together in the same feature?
- which values never/rarely occur together in the same feature?

Data Investigation ⑨

PATTERN 6: LINK PREDICTION THROUGH GROUP SIMILARITY

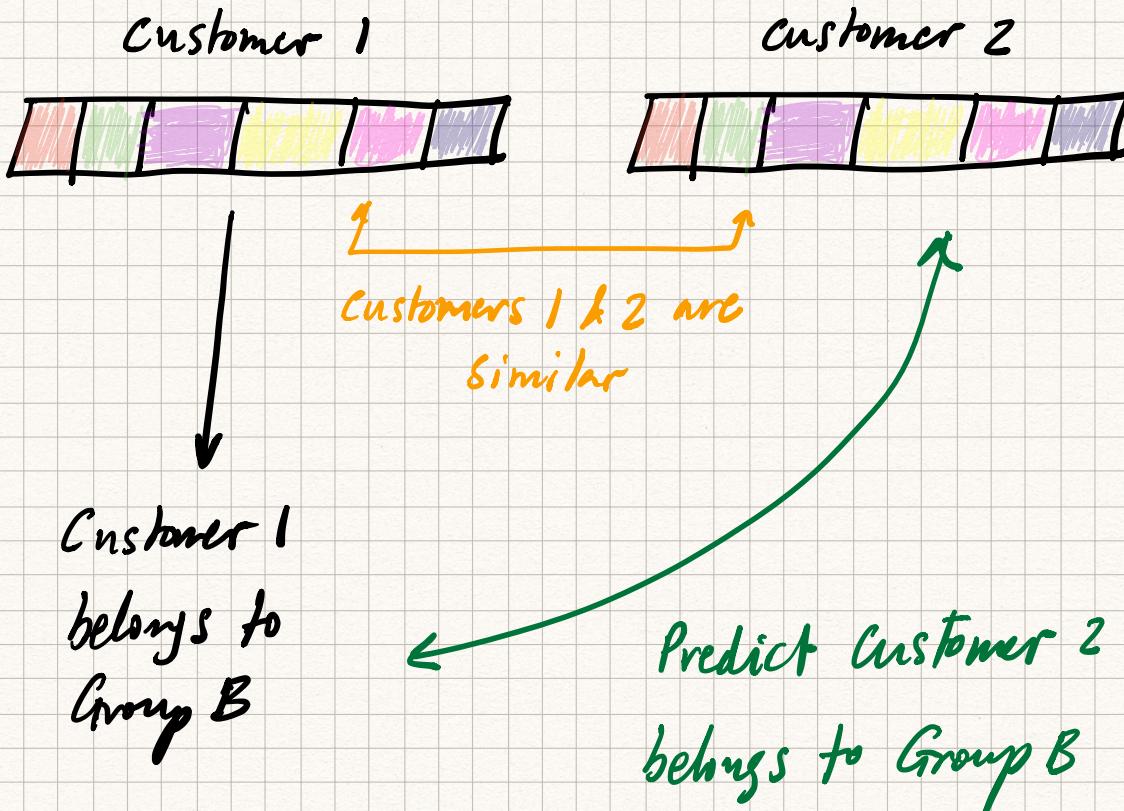


Group A: The set of customers who share the same 3 friends.

Predict that customers 1 and 2 will "respond" similarly or have the same interests - they are linked.

Data Investigation (10)

PATTERN 6: LINK PREDICTION THROUGH INDIVIDUAL SIMILARITY



Group B: The set of customers who responded to the last marketing campaign.

Data Investigation (11)

PATTERN 7: FEATURE ENGINEERING

- Which features can be ignored/deleted?
 - We'll look at some techniques later based on total variance explained
- Which features need creation?
 - Combination of existing features (e.g., monthly charge / tenure)
 - Invention of new features (e.g., # of senior citizen friends)
 - Extension of existing features (e.g., tenure², log(monthly charge))