

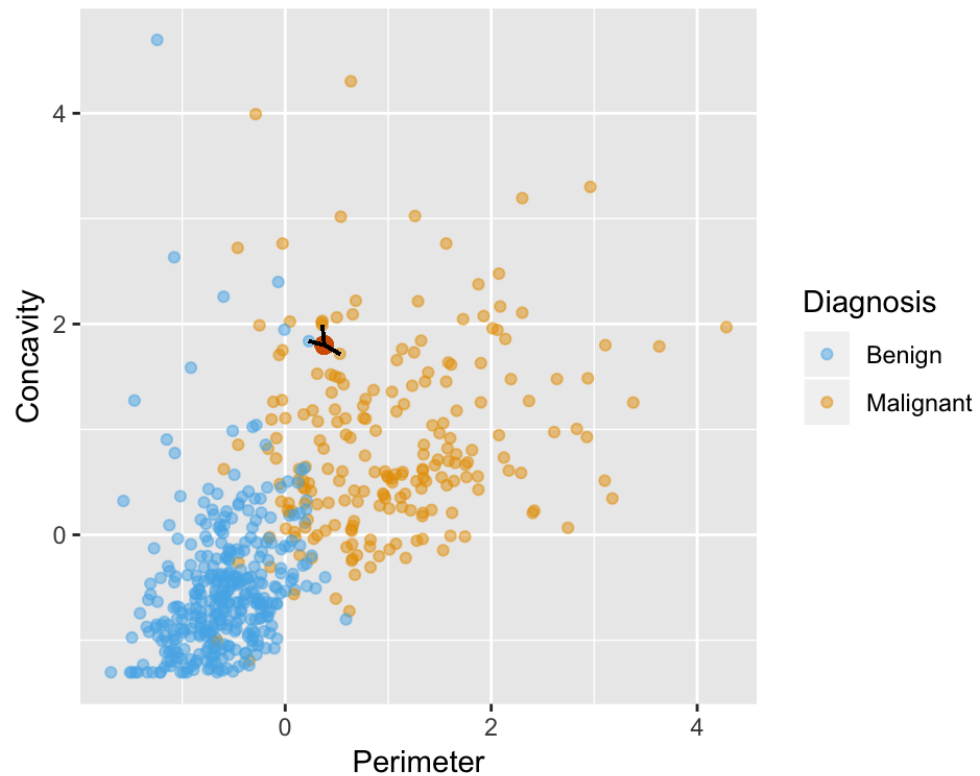
DSCI 100 - Introduction to Data Science

Lecture 7 - Classification continued

2019-02-14

Continuing with the classification problem

Can we use data we have seen in the past, to predict something about the future?

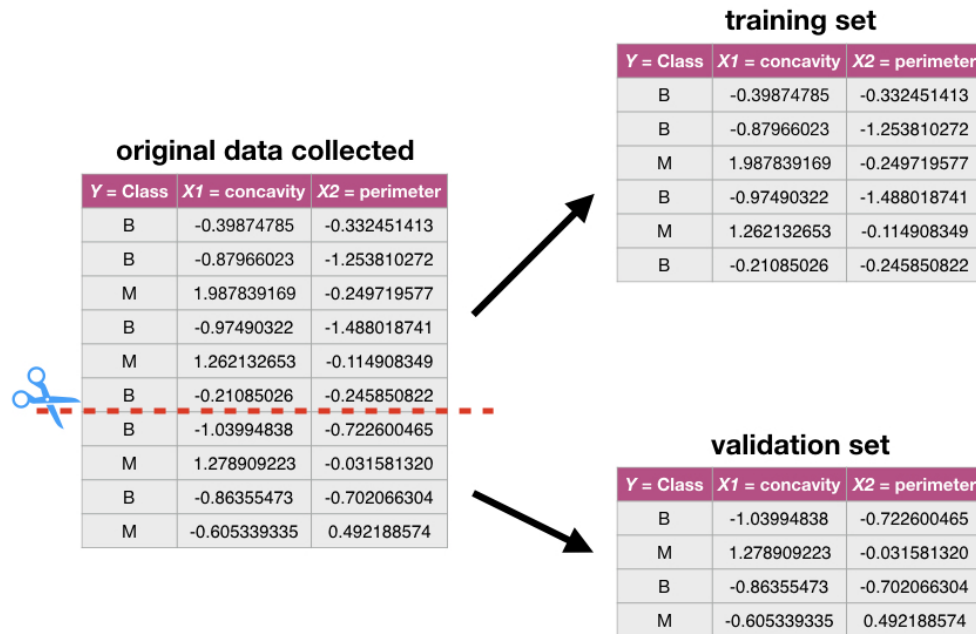


Unanswered questions from last week

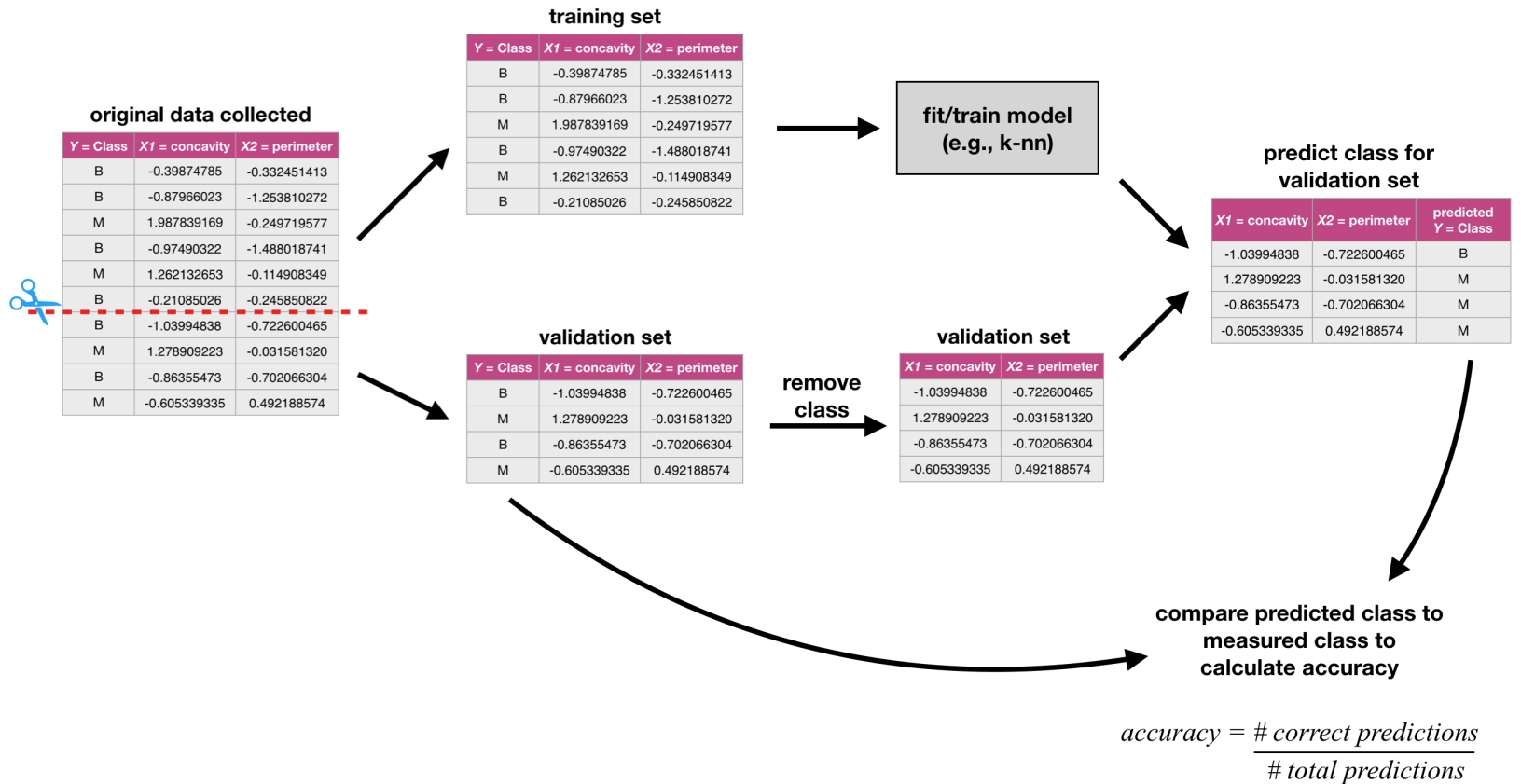
1. Is our model any good?
2. How do we choose k ?

1. Is our model any good?

Creating the training and validation sets

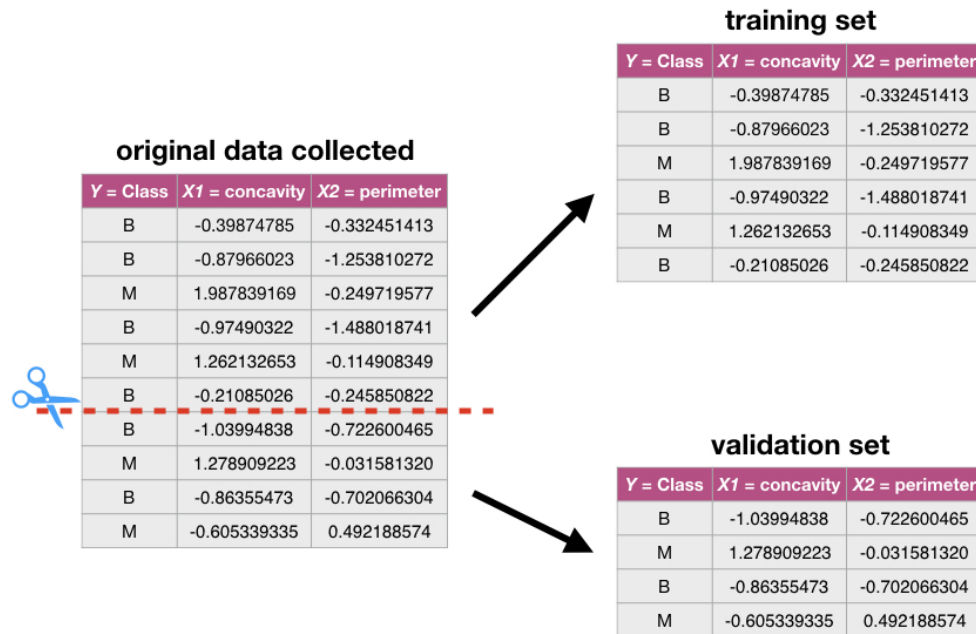


Assumption -> there is no order to the original data collected



Is one accuracy measurement good enough?

Creating the training and validation sets



Assumption -> there is no order to the original data collected

Cross-validation as an alternative approach

			fold 1	fold 2	fold 3	fold 4	fold 5
Y = Class	X1 = concavity	X2 = perimeter					
M	2.166629	2.087306	validation	training	training	training	training
M	0.451585	0.17376					
M	0.482262	-0.05241					
M	2.216001	1.289747					
B	-0.5749	-1.48802	training	validation	training	training	training
M	1.094846	-0.13276					
B	0.157276	-0.51728					
B	-0.21085	-0.24585					
B	-1.03995	-0.7226	training	training	validation	training	training
B	-0.80139	-0.33275					
B	-0.86355	-0.70207					
M	-0.0474	0.828473					
M	-0.00474	0.310655	training	training	validation	training	training
B	-0.88964	-0.75742					
M	0.554641	1.870061					
M	0.414676	0.251135					
B	-1.29584	-0.50655	training	training	training	validation	training
B	-0.74482	-1.40945					
B	-1.20076	-1.11305					
B	-0.87324	-1.34517					
M	2.013244	0.352318	training	training	training	validation	training
B	-0.63506	-1.24012					
M	0.824144	0.135073					
M	0.144813	0.2184					
B	-0.85354	-1.19727	training	training	training	training	validation
B	-0.61253	-0.8708					
B	4.696536	-1.2428					
B	-0.80666	-0.64255					
M	-0.01672	1.74507	training	training	training	training	validation
M	1.842602	1.319507					
M	0.027377	0.090433					
B	-0.48838	-0.52232					
M	0.051344	0.640987	training	training	training	training	validation
B	-0.89538	-0.47321					
M	0.510063	1.274867					

$accuracy_1$

$accuracy_2$

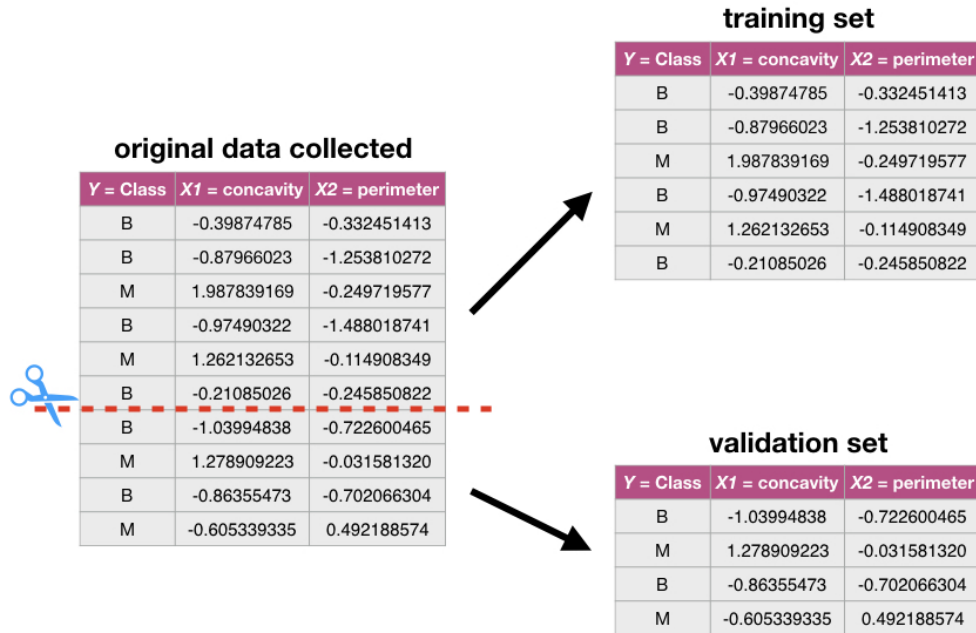
$accuracy_3$

$accuracy_4$

$accuracy_5$

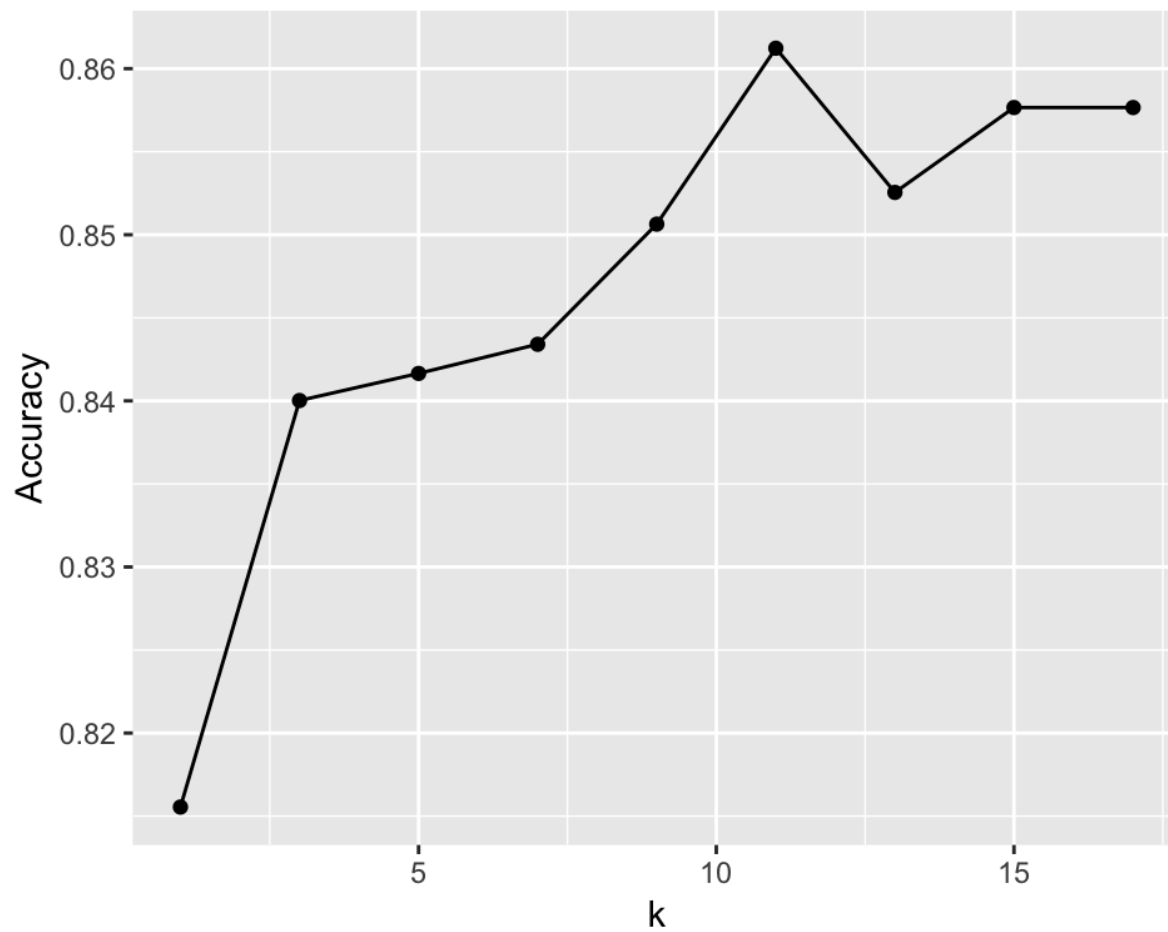
$$cross-validation\ accuracy = \frac{accuracy_1 + accuracy_2 + accuracy_3 + accuracy_4 + accuracy_5}{\# folds}$$

Creating the training and validation sets

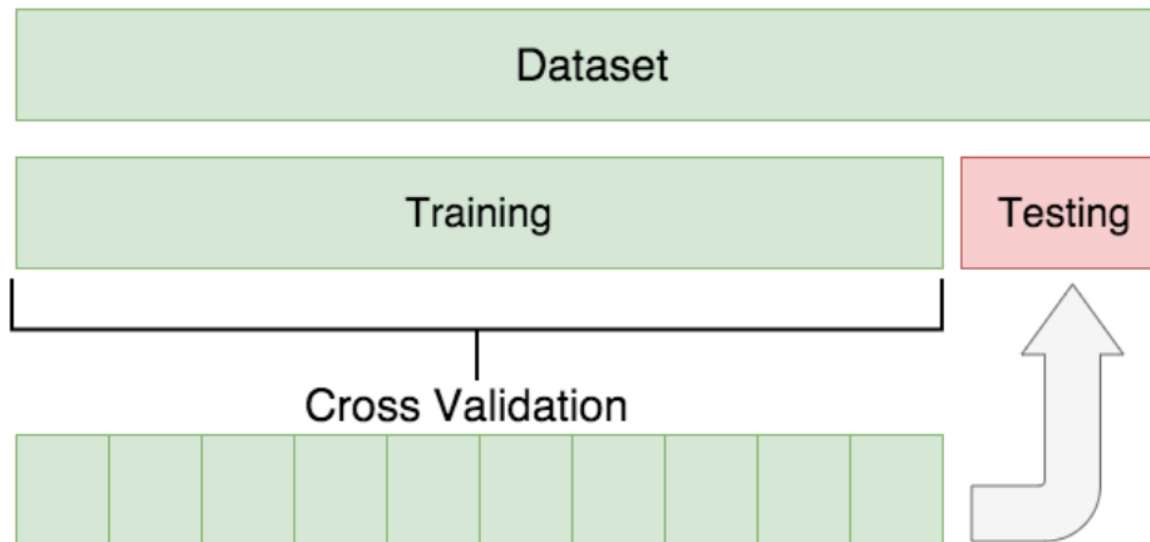


Assumption -> there is no order to the original data collected

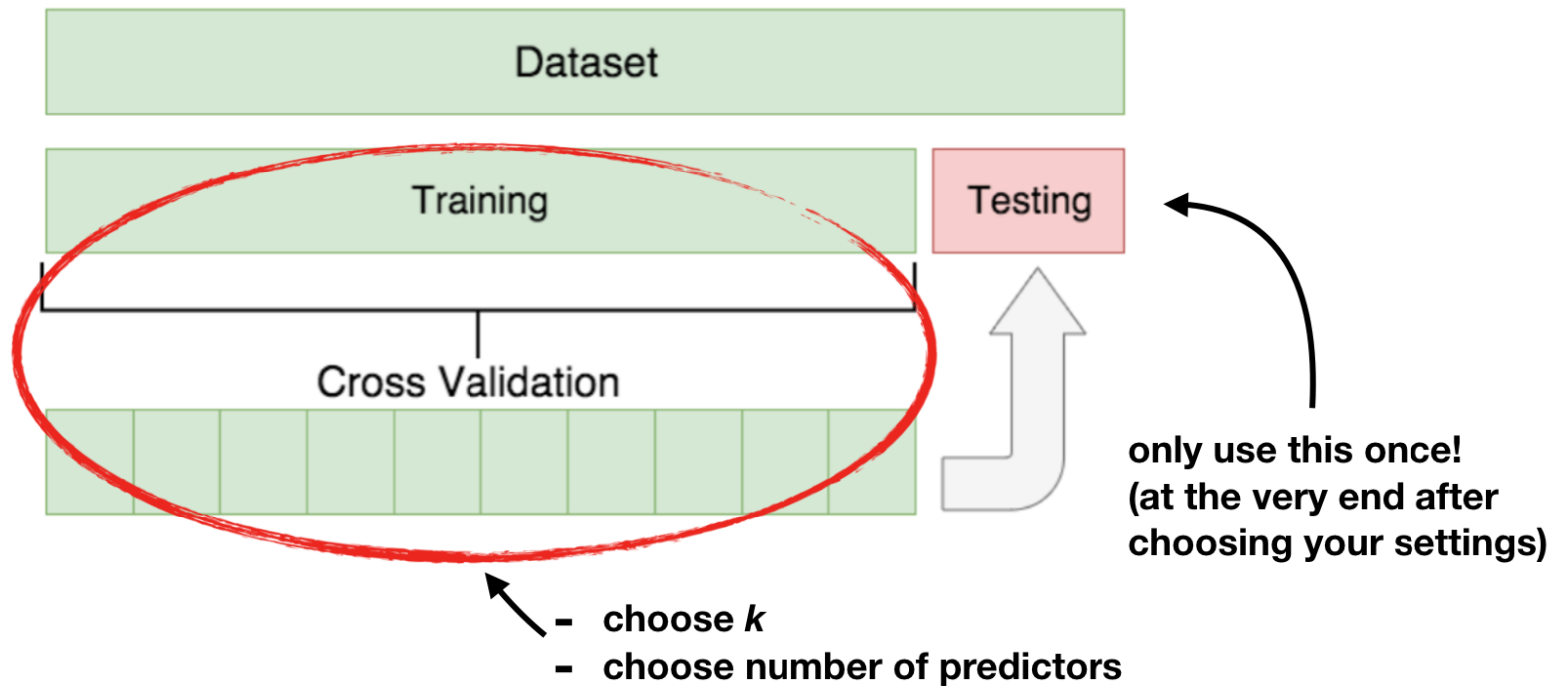
2. How do we choose k ?



The big picture



The big picture



Why are we doing all this???



Our question:

- Can we use past information to predict the class labels of new observations we don't have labels for?
- We can always do this, but we might only want to do this if we have evidence we can do this well.

Class activity 1

- In your group, discuss and explain cross-validation in your own words. Post your group's answer as a response to this post in Piazza.

Class activity 2

- In your group, discuss and explain what a test, validation and training data set are in your own words. Post your group's answer as a response to this post in Piazza.