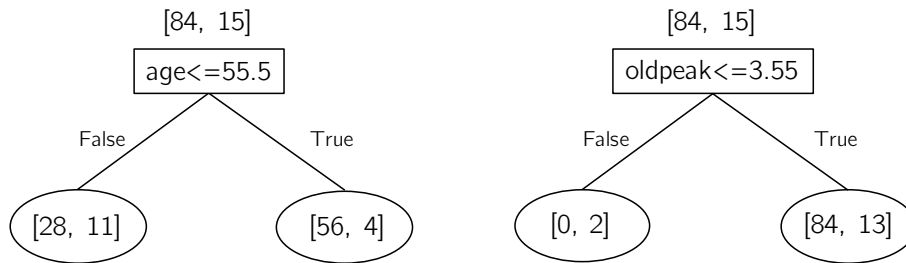
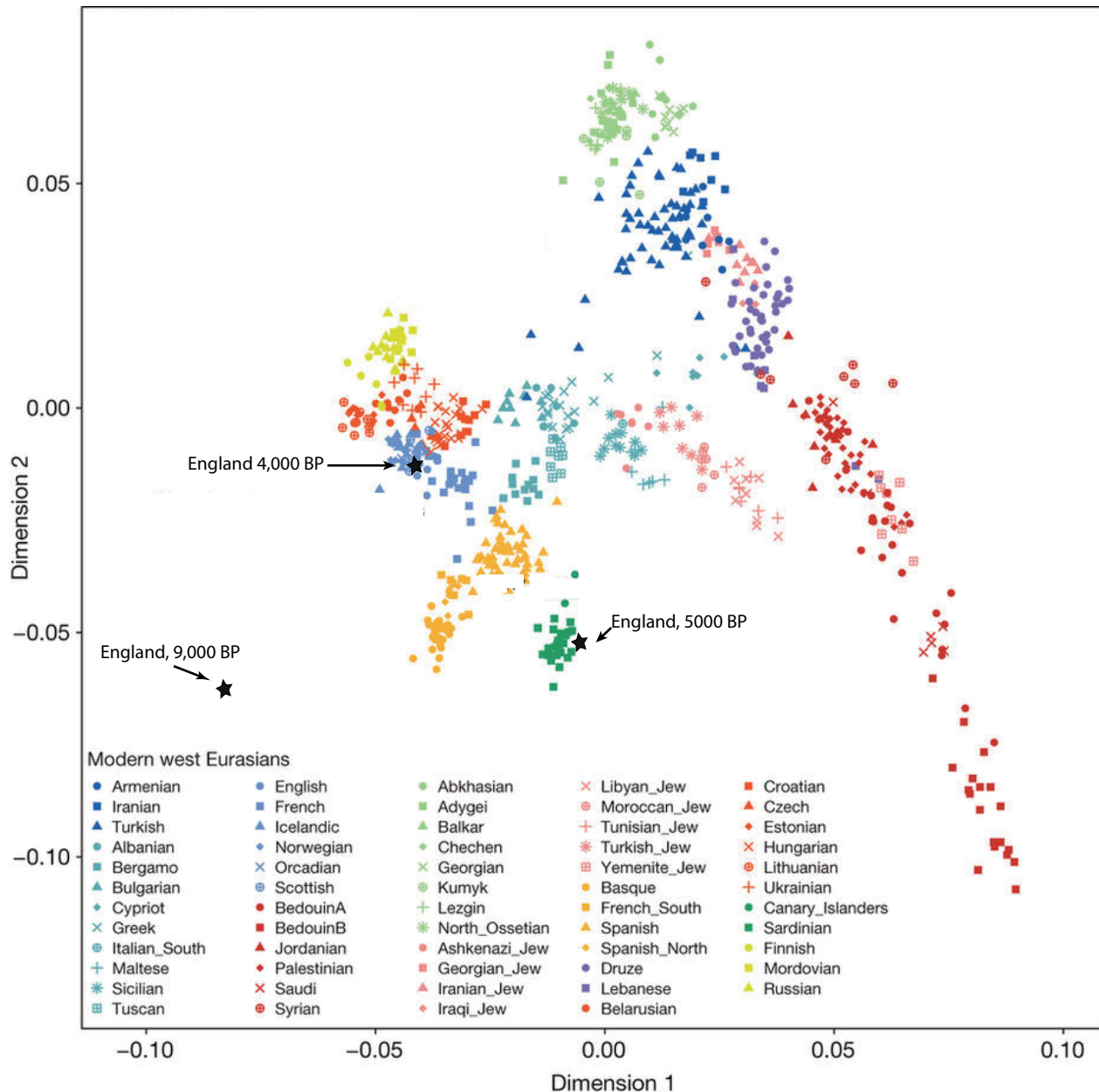


3. *Entropy*. Consider the two feature choices below (for the heart disease dataset), and their associated splits. Counts of label -1 vs. 1 are shown in brackets.



- (a) After splitting the data based on each feature, what is the *classification error* for each tree?
- (b) Before considering the feature, what is $H(Y)$, the entropy of the initial partition? (don't need to find a value, just set up the equation)
- (c) Which tree do you think produces more information gain?

4. The figure below shows the first two PCs of genome-wide data from 777 present-day people from West Eurasia, along with three ancient British people who lived 9000, 5000 and 4000 years ago (labeled stars, “BP” means “[years] Before Present”).



- (a) What can you infer about the relationship between each of the ancient people and present-day Europeans?
- (b) What does this figure suggest about the history of Britain, and the people living there, over the past ten thousand years?

Acknowledgements: Iain Mathieson