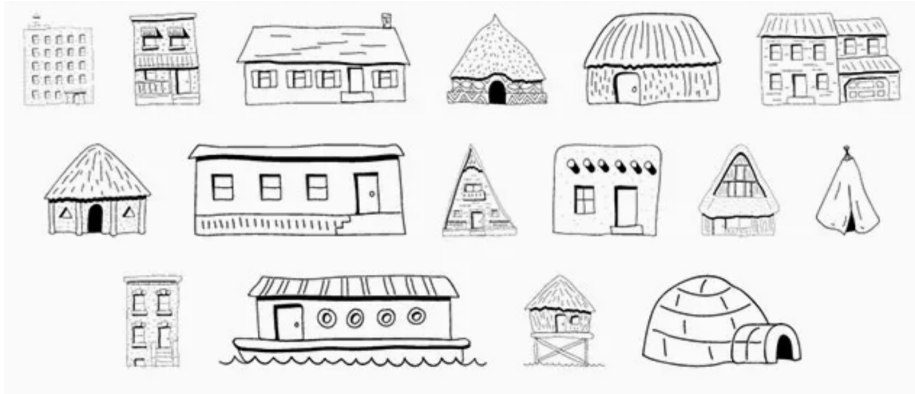


7 Machine Learning and Real-world Data (sht25)

You are part of a team of developers whose task is to classify hand-drawn images into classes defined by the objects they depict, in the framework of a new search engine that is based on shapes. The images are in black and white and are represented by 150×200 pixel bitmaps. Consider the figure below, which shows some examples of such images for the classes “house”, “tent”, “boat” and “igloo”.



- (a) You want to use a Naive Bayes Classifier with a “bag of pixels” representation.
- (i) Describe which approach you would use for classification, including which features you would use and how many there are. State the relevant formula. [5 marks]
 - (ii) What does the training material for the task look like, how was it gathered from human annotators, and how is it used to estimate the model’s parameters? [4 marks]
- (b) To improve classification, you plan to preprocess the raw pixel features in some form so that a smaller, more informative set of features results. Describe the best such manipulation you can think of, and why this should improve results. [4 marks]
- (c) Your colleague suggests that instead of pixel-based features, you should derive some higher-level features for the classifier, which capture holistic visual properties of the image instead of exact pixel location.
- (i) Describe three such features which you think will be most helpful in classification, giving reasons why this feature should work and using the above figure for illustration. Choose maximally different features. [3 marks]
 - (ii) Describe how the final classifier can be modified so that it combines pixel-based features with your features from Part (c) above. What possible problems do you foresee, and how could they be addressed? [4 marks]