

## # General notes on machine learning

### Training and testing data

- Training and predictions should be done on different sets of data. If training and predictions are done on same type of data, then there will be overfitting of the data
- Generalizing to new data is not possible if the training and prediction are both on the same dataset
- Save 10% of the dataset to be used as

### Unpacking NB using Bayes Rule

#### Quiz : Prior and Posterior

- Bayes rule
  - prior probability + test evidence = posterior probability **semantically Bayes incorporates some evidence from the test into prior probability to arrive at a posterior probability.**
- 1. PRIOR
  - $P(C) = 0.01 = 1\%$  **prior probability of cancer**
  - $P(\text{Pos}|C) = 0.9 = 90\%$  **prior probability of the test being positive and the person having cancer**
  - $P(\text{Neg}|\neg C) = 0.9 = 90\%$  **prior probability of the test being negative and the person not having cancer**
  - $P(\neg C) = 0.99 = 99\%$  **prior probability of not having cancer**
  - $P(\text{Pos}|\neg C) = 0.1$  **prior probability of the test being positive and the person not having cancer**
- 2. POSTERIOR
  - $P(C | \text{Pos}) = P(C) \times P(\text{Pos} | C) = 0.009 = 0.9\%$ 
    - \*  **$P(C|\text{Pos})$  is posterior of the probability of cancer given that the test says positive**
    - \*  $P(C)$  is prior probability of cancer
    - \*  **$P(\text{Pos} | C)$  is probability of positive result given that a person has cancer. This is the test sensitivity.**
  - $P(\neg C | \text{Pos}) = P(\neg C) \times P(\text{Pos} | \neg C) = 0.099 = 9.9\%$ 
    - \*  **$P(\neg C|\text{Pos})$  is posterior of the probability of not having cancer given that the test says positive**
    - \*  $P(\neg C)$  is prior probability of not having cancer
    - \*  **$P(\text{Pos} | \neg C)$  is probability of positive result given that a person does not have cancer.**

[fig1] [fig1]: nb1.png

\* But probabilities do not add upto 1