**Naive Bayes**                                                          *(find and work with a partner)*

Say we have two tests for a specific disease. Each test (features $f_1$, $f_2$) can come back either positive
"pos" or negative "neg", and the true underlying condition of the patient is represented by $y$ ($y = 1$ is
"healthy" and $y = 2$ is "disease"). We observe this training data where $n = 7$ and $p = 2$:

| $x$ | $f_1$ | $f_2$ | $y$ |
|-----|-----|-----|-----|
| $x_1$ | pos | neg | 1 |
| $x_2$ | pos | pos | 2 |
| $x_3$ | pos | neg | 2 |
| $x_4$ | neg | neg | 1 |
| $x_5$ | pos | neg | 2 |
| $x_6$ | neg | neg | 1 |
| $x_7$ | neg | pos | 2 |

1. To estimate the probability $p(y = k)$, for $k = 1, 2, \cdots, K$, we will use the formula:

$$\theta_k = \frac{N_k + 1}{n + K}$$

where $N_k$ is the count ("Number") of data points where $y = k$. Compute $\theta_1$ and $\theta_2$. What would
$\theta_1$ and $\theta_2$ be if we in fact had *no* training data?

2. To estimate the probabilities $p(x_j = v | y = k)$ for all features $j$, values $v$, and class label $k$, we will
use the formula:

$$\theta_{k,j,v} = \frac{N_{k,j,v} + 1}{N_k + |f_j|}$$

where $N_{k,j,v}$ is the count of data points where $y = k$ and $x_j = v$, and $|f_j|$ is the number of possible
values that $f_j$ (feature $j$) can take on. Fill in the following tables with these $\theta$ values.

| $y = 1$ | pos | neg |
|-----|-----|-----|
| $f_1$ |  |  |
| $f_2$ |  |  |

| $y = 2$ | pos | neg |
|-----|-----|-----|
| $f_1$ |  |  |
| $f_2$ |  |  |