1. (a) $\int_0^1 P(S=1;\theta)d\theta + \int_0^1 P(S=0;\theta)d\theta = 1$

$$LHS = c\left(\int_0^1 \theta^1(1-\theta)^{1-1}d\theta + \int_0^1 \theta^0(1-\theta)^{1-0}d\theta\right)$$

$$= c\left(\int_0^1 \theta\, d\theta + \int_0^1 (1-\theta)d\theta\right) = 1$$

So $c \times 1 = 1, \quad c = 1$

(b) $E(S) = 1 \times P(S=1;\theta) + 0 \times P(S=0;\theta) = 1 \times \theta + 0 \times (1-\theta) = \theta$

(c) $E\left[(S-\bar{S})^2\right] = (1-\theta)^2 \times P(S=1;\theta) + (0-\theta)^2 \times P(S=0;\theta)$

$$= \theta(1-\theta)$$

2. (a) Since $\int_0^1 P_0(\theta)d\theta = \int_0^1 c'\theta^{\alpha-1}(1-\theta)^{\beta-1}d\theta = c' B(\alpha,\beta) = 1$

So $c' = \dfrac{1}{B(\alpha,\beta)}$

(b)
$$P(S=s) = \int_0^1 P(S=s|\theta)\cdot P_0(\theta)d\theta = \int_0^1 \theta^s(1-\theta)^{1-s} c'\theta^{\alpha-1}(1-\theta)^{\beta-1}d\theta$$

$$= \frac{1}{B(\alpha,\beta)}\int_0^1 \theta^{\alpha+s-1}(1-\theta)^{1-s+\beta-1}d\theta$$

$$= \frac{1}{B(\alpha,\beta)}\int_0^1 \theta^{\alpha+s-1}(1-\theta)^{\beta-s}d\theta = \frac{B(\alpha+s,\ \beta-s+1)}{B(\alpha,\beta)}$$

(c) $E(S) = 1 \times P(S=1) + 0 \times P(S=0) = P(S=1) = \dfrac{B(\alpha+1,\ \beta-1+1)}{B(\alpha,\beta)}$

$$= \frac{B(\alpha+1,\ \beta)}{B(\alpha,\beta)} = \frac{\Gamma(\alpha+1)\Gamma(\beta)/\Gamma(\alpha+\beta+1)}{\Gamma(\alpha)\Gamma(\beta)/\Gamma(\alpha+\beta)} = \frac{\alpha}{\alpha+\beta}$$

(d) $E(S^2) = 1^2 \times P(S=1) + 0^2 \times P(S=0) = \dfrac{\alpha}{\alpha+\beta}$ $\quad \therefore \operatorname{Var}(S) = \dfrac{\alpha}{\alpha+\beta} - \left(\dfrac{\alpha}{\alpha+\beta}\right)\left(\dfrac{\alpha}{\alpha+\beta}\right) = \dfrac{\alpha\beta}{(\alpha+\beta)^2}$

3. (a) Since $S \sim Ber(\theta)$, so $P(S|\theta) = \theta^S (1-\theta)^{1-S}$

$$L(\{s_i\}; \theta) = \prod_{i=1}^{N} \theta^{S_i} (1-\theta)^{1-S_i} \quad \text{and} \quad \hat{\ell}(\{s_i\}; \theta) = \frac{1}{N}\left( \log \left( \prod_{i=1}^{N} \theta^{S_i} (1-\theta)^{1-S_i} \right) \right)$$

$$\therefore \hat{\ell}(\{s_i\}; \theta) = \frac{1}{N}\left( \sum_{i=1}^{N} S_i \log \theta + \sum_{i=1}^{N} (1-S_i) \log(1-\theta) \right)$$

$$= \frac{1}{N}\left( \hat{N_1} \log \theta + \hat{N_0} \log(1-\theta) \right)$$

$$\therefore \frac{\partial \hat{\ell}}{\partial \theta} = \frac{\hat{N_1}}{N} \cdot \frac{1}{\theta} + \frac{\hat{N_0}}{N} \cdot \frac{-1}{1-\theta} = 0 \Rightarrow \hat{\theta} = \frac{\hat{N_1}}{\hat{N_1} + \hat{N_0}}$$

(b)
$$\begin{cases} P(S=1 | \hat{N_0}, \hat{N_1}) = \hat{\theta} = \frac{\hat{N_1}}{\hat{N_0} + \hat{N_1}} \\ \\ P(S=0 | \hat{N_0}, \hat{N_1}) = 1-\hat{\theta} = \frac{\hat{N_0}}{\hat{N_0} + \hat{N_1}} \end{cases}$$

(c) From the results of question (a) and (b), for next experiment,

$$\begin{cases} P(S=1 | \hat{N_0}=0, \hat{N_1}=1) = \frac{1}{1+0} = 1 \\ P(S=0 | \hat{N_0}=0, \hat{N_1}=1) = \frac{0}{1+0} = 0 \end{cases}$$

It seems ridiculous, but also make sense because we have too few samples and observations. If we have done more experiments, the estimator $\hat{\theta}$ would converge to the true value.

**4.** Since $S \sim Ber(\theta)$, $\theta \sim Beta(\alpha, \beta)$

**ⓐ** So $P(\theta | \{s_i\}) = \dfrac{\prod_{i=1}^{N} P(s_i|\theta) \, P_0(\theta)}{P(\{s_i\})} = \dfrac{\prod_{i=1}^{N} \left[\theta^{s_i}(1-\theta)^{1-s_i}\right] Beta(\alpha, \beta)}{P(\{s_i\})} = \dfrac{\theta^{\hat{M}}(1-\theta)^{\hat{N_0}} Beta(\alpha, \beta)}{P(\{s_i\})} = \dfrac{\theta^{\hat{M}}}{} (1-\theta)^{\hat{N_0}}$

and $\displaystyle\int_0^1 d\theta \, P(\theta | \{s_i\}) = 1$

$$= \dfrac{\theta^{\hat{M}}(1-\theta)^{\hat{N_0}} \cdot \theta^{\alpha-1}(1-\theta)^{\beta-1}}{P(\{s_i\}) \, B(\alpha, \beta)}$$

$$= C \cdot \theta^{(\hat{M}+\alpha-1)}(1-\theta)^{(\hat{N_0}+\beta-1)}$$

So $\displaystyle\int_0^1 d\theta \cdot C \cdot \theta^{\hat{M}+\alpha-1}(1-\theta)^{\hat{N_0}+\beta-1} = 1$

$\therefore P(\theta|\{s_i\}) = \dfrac{1}{B(\hat{M}+\alpha, \hat{N_0}+\beta)} \theta^{\hat{M}+\alpha-1}(1-\theta)^{\hat{N_0}+\beta-1}$

$\therefore C \cdot B(\hat{M}+\alpha, \hat{N_0}+\beta) = 1 \qquad \therefore C = \dfrac{1}{B(\hat{M}+\alpha, \hat{N_0}+\beta)}$

**ⓑ** $\therefore P(s=1 | \{s_i\}) = \displaystyle\int_0^1 \theta \cdot P(\theta|\{s_i\}) \, d\theta = \int_0^1 C \cdot \theta^{(\hat{M}+\alpha-1+1)}(1-\theta)^{(\hat{N_0}+\beta-1)} \, d\theta = \dfrac{B(\hat{M}+\alpha+1, \hat{N_0}+\beta)}{B(\hat{M}+\alpha, \hat{N_0}+\beta)}$

$$= \dfrac{\Gamma(\hat{M}+\alpha+1)\,\Gamma(\hat{N_0}+\beta)}{\Gamma(\hat{M}+\alpha+1+\hat{N_0}+\beta)} \Big/ \dfrac{\Gamma(\hat{M}+\alpha)\,\Gamma(\hat{N_0}+\beta)}{\Gamma(\hat{M}+\alpha+\hat{N_0}+\beta)} = \dfrac{\hat{M}+\alpha}{\hat{M}+\hat{N_0}+\alpha+\beta}$$

$$P(s=0|\{s_i\}) = 1 - \dfrac{\hat{M}+\alpha}{\hat{N_1}+\hat{M}+\alpha+\beta} = \dfrac{\hat{N_0}+\beta}{\hat{M}+\hat{N_0}+\alpha+\beta}$$

**ⓒ** $E(s | \hat{N_1}, \hat{M} ; \alpha, \beta) = 1 \times P(s=1 | \hat{N_0}, \hat{M} ; \alpha, \beta) + 0 \times P(s=0 | \hat{N_0}, \hat{M} ; \alpha, \beta)$

$$= \dfrac{\hat{N_1}+\alpha}{\hat{N_0}+\hat{M}+\alpha+\beta}$$

**ⓓ** $\begin{cases} P(s=1 | \{s_i\}) = \dfrac{1+\alpha}{1+\alpha+\beta} \\[2mm] P(s=0 | \{s_i\}) = \dfrac{\beta}{1+\alpha+\beta} \end{cases}$

At this scenario, the prior distribution parameters $(\alpha, \beta)$ counts, and would fix some problem by giving a prior to the event' belief that at which stage the event would happen.

And as the experiments go on, we gather more and more data, we would get more precise predictions.

**5.**

$$y = \sum_{d=1}^{D} x_d \theta_d + \epsilon$$

(a) For 1 observation, $E(y_i) = E\left(\sum_{d=1}^{D} x_{id}\theta_d\right) + E(\epsilon)$

$$= \sum_{d=1}^{D} x_{id}\theta_d$$

Since $\epsilon \sim N(0, \sigma^2)$ so $y_i \sim N(E(y_i), \sigma^2)$, $\quad f(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \sum_{d=1}^{D} x_{id}\theta_d)^2}{2\sigma^2}\right)$

$$P(y_i \leq y \mid x_{i,d}; \theta) = \int_0^y \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y' - \sum_{d=1}^{D} x_{id}\theta_d)^2}{2\sigma^2}\right) dy'$$

(b) Since for each $i,j$, $\epsilon_i$ is independent of $\epsilon_j$, so $y_i$ is also independent of $y_j$

so $E(\vec{y}) = \begin{pmatrix} \sum_{d=1}^{D} x_{1,d}\theta_d \\ \vdots \\ \sum_{d=1}^{D} x_{N,d}\theta_d \end{pmatrix} \triangleq \vec{\mu}$

and $\text{Cov}(y_i, y_j) = \begin{cases} 0 & \text{if } i \neq j \\ \sigma^2 & \text{if } i = j \end{cases}$ so $\Sigma = \sigma^2 I_N$.

So the join density $P\left(\{y_i\} \mid \{x_{i,d}\}; \theta\right)$ of all the $N$ observations $\{y_i, x_{i,d}\}_{i=1}^{N}$ should be:

$$f(\vec{y}) = \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\vec{y} - \vec{\mu})^T \Sigma^{-1} (\vec{y} - \vec{\mu})\right)$$

where $\vec{\mu}$ and $\Sigma$ are defined before.

(c) $\hat{\ell}(\theta; \{y_i, x_{i,d}\}) = \frac{1}{N} \sum_{i=1}^{N} \log f(y_i)$ Since $\log f(y_i) = \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \log\left(\exp\left(-\frac{(y_i - \sum_{d=1}^{D} x_{id}\theta_d)^2}{2\sigma^2}\right)\right)$

So $\hat{\ell}(\theta; \{y_i, x_{i,d}\}_{i=1}^{N}) = \frac{-1}{2}\log(2\pi\sigma^2) + \left(-\frac{1}{2\sigma^2}\right) \times \frac{1}{N} \times \left(\sum_{i=1}^{N}(y_i - \sum_{d=1}^{D} x_{id}\theta_d)^2\right)$

(d) $E(\theta; \{y_i, x_{i,d}\}_{i=1}^{N}) = -\hat{\ell}(\theta; \{y_i, x_{i,d}\}_{i=1}^{N})$

$$= \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2\sigma^2 N}\left(\sum_{i=1}^{N}(y_i - \sum_{d=1}^{D} x_{id}\theta_d)^2\right)$$

(g) The OLS method is to some extent equivalent to the ~~method~~ MLE method given that the error terms are normally distributed.

(e) $\frac{\partial E}{\partial \theta_d} = \frac{1}{\sigma^2 N} \sum_i \left(\sum_{d'=1}^{D} x_{i,d'}\theta_{d'} - y_i\right) x_{i,d}$     (f) $(X^T X)\theta = X^T y$

6. Similar to question 5, we have:

(a) $P(y_i|x_{i};\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(- \frac{(y - \sum_{d=1}^{D} h_d(x)\theta_d)^2}{2\sigma^2}\right)$

(b) $\hat{\ell}(\theta; \{y_i, x_{id}\}_{i=1}^{N}) = \frac{1}{-2N\sigma^2}\left(\sum_{i=1}^{N}(y_i - \sum_{d=1}^{D} h_d(x)\theta_d)^2\right) - \frac{1}{2}\log(2\pi\sigma^2)$

$E(\theta; \{y_i, x_{id}\}_{i=1}^{N}) = -\hat{\ell}(\theta; \{y_i, x_{id}\}_{i=1}^{N})$

$= \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2\sigma^2 N}\left(\sum_{i=1}^{N}(y_i - \sum_{d=1}^{D} h_d(x)\theta_d)^2\right)$

$= \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2\sigma^2 N}(Y - H\theta)^T(Y - H\theta)$

(c) Similarly, since for linear $X$, $(X^T X)\theta = X^T Y$

how we can consider $H$ as new $X$,

So $(H^T H)\theta = H^T Y$

(d) It is the same.

Once we get the matrix $H$ computed, the dimensionality of $x$ does not matter.

(e) The function can be expressed, so the matrix $H$ is not full-ranked.

Some rows/columns of the matrix is redundant.

We need to do Q-R decomposition of $H$.

Say $H = QR$, where $Q^T Q = I$, and $R$ is upper-triangle matrix

7.  (a) $P(\theta) = N(\theta; 0, \frac{\sigma^2}{\lambda}1)$

$$E = 2\log 2\pi\sigma^2 + \frac{1}{2\sigma^2 N}(Y - X\theta)^T(Y - X\theta) + \frac{\lambda}{2\sigma^2 N}\theta^T\theta$$

$$\frac{\partial \bar{E}}{\partial \theta_d} = \frac{1}{\sigma^2 N}\left\{\sum_{i=1}^{N} X_{id}\left(\sum_d (X_{id}\theta_d - y_i)\right) + \lambda\theta_d\right\}$$

or $\frac{\partial E}{\partial \theta} = \frac{1}{\sigma^2 N}\left\{X^T(X\theta - Y) + \lambda\theta\right\}$

(b) $P_o(\theta) = Ce^{-\lambda \sum_d |\theta_d|}$

$$E = 2\log 2\pi\sigma^2 + \frac{1}{2\sigma^2 N}\sum_i\left(y_i - \sum_d X_{id}\theta_d\right)^2 + \frac{\lambda}{N}\sum_d|\theta_d|$$

$$\frac{\partial E}{\partial \theta_d} = \frac{1}{\sigma^2 N}\left\{\sum_i X_{id}\left(\sum_{d'} X_{id'}\theta_{d'} - y_i\right)\right\} + \frac{\lambda}{N}\,\text{sgn}\,\theta_d$$

(c) $P_o(\theta) = Ce^{-\lambda \sum_d |\theta_d|} N(\theta; 0, \frac{\sigma^2}{\lambda}1)$

$$E = 2\log 2\pi\sigma^2 + \frac{1}{2\sigma^2 N}\sum_i\left(y_i - \sum_d X_{id}\theta_d\right)^2 + \frac{\lambda_1}{N}\sum_d|\theta_d| + \frac{\lambda_2}{\sigma^2 N}\sum_d\theta_d^2$$

$$\frac{\partial E}{\partial \theta_d} = \frac{1}{\sigma^2 N}\left\{\sum_i X_{id}\left(\sum_{d'} X_{id'}\theta_{d'} - y_i\right)\right\} + \frac{\lambda_1}{N}\,\text{sgn}\,\theta_d + \frac{\lambda_2}{\sigma^2 N}\theta_d$$