## Bayesian Inference:

Recall MLE:

$$\hat{\theta} = \arg\max_{\theta} \left[ p(D \mid \theta) \right]$$

where $D$ is data, and $\theta$ is the unknown parameter/variable.

Here, Parameter $\theta$ is unknown but deterministic (frequentist view).

### Bayesian:

- $\theta$ is viewed as a random variable (even when it is actually deterministic).
- Use Bayes' Rule to calculate posterior distribution.

$$\underbrace{p(\theta \mid D)}_{\text{posterior distribution}} = \frac{\overbrace{p(D \mid \theta)}^{\text{likelihood}} \overbrace{p(\theta)}^{\text{prior}}}{p(D)} \propto p(D \mid \theta) p(\theta).$$

$$p(D) = \int p(D \mid \theta) p(\theta) \, d\theta.$$

### Proof of Baye's Rule:

Multiplication Rule:

$$\left. \begin{array}{l} p(\theta \cap D) = p(\theta) p(D \mid \theta) \\ p(D \cap \theta) = p(D) p(\theta \mid D) \end{array} \right\} \text{same.}$$

(the two left sides are "same")

$$\Rightarrow p(\theta) p(D \mid \theta) = p(D) p(\theta \mid D)$$

$$\Rightarrow p(\theta \mid D) = \frac{p(D \mid \theta) p(\theta)}{p(D)}.$$

Example : Did the sun just explode?

- Suppose we have a device that detects if the sun explodes with high accuracy:

$$P( X = \theta \mid \theta ) = 1 - \alpha$$
$$P( X = 1 - \theta \mid \theta ) = \alpha$$

$\alpha$ = error, known & fixed.
↑ small

where

$\theta \in \{0,1\}$ = did the sun exploded. //binary

$X \in \{0, 1\}$ = did the device alarms. //binary.

- If alarm fires $(X = 1)$, should we believe the sun exploded or not?

MLE :

$$\hat{\theta} = \underset{\theta \in \{0,1\}}{\arg\max} \left[ P( X = 1 \mid \theta ) \right]$$

$$= \begin{cases} \alpha, & \text{for } \theta = 0 \\ 1 - \alpha, & \text{for } \theta = 1 \end{cases}$$

$$= 1$$

Bayes :

Step 1 : Find the prior.

$$P(\theta) = \begin{cases} \beta \ (\text{very very small}) & \text{for } \theta = 1, \\ 1 - \beta & \text{for } \theta = 0. \end{cases}$$

Step 2 : Set up equation for posterior :

$$P(\theta \mid X = 1) = \frac{P(X = 1 \mid \theta) \, P(\theta)}{P(X = 1)} \quad \propto \quad P(X = 1 \mid \theta) \, P(\theta)$$

$$= \begin{cases} \overbrace{(1-\alpha)}^{P(x=1\mid\theta=1)} \overbrace{\beta}^{P(\theta=1)}, & \text{if } \theta = 1 \\ \underset{P(x=1\mid\theta=0)}{\alpha} \underset{P(\theta=0)}{(1-\beta)}, & \text{if } \theta = 0, \end{cases}$$

If $(1-\alpha)\beta > \alpha(1-\beta)$ : predict $\theta = 1$.

If $(1-\alpha)\beta < \alpha(1-\beta)$ : predict $\theta = 0$.

Equivalently :

predict $\theta = 1$, if $\dfrac{\beta}{1-\beta} > \dfrac{\alpha}{1-\alpha}$

predict $\theta = 0$, if $\dfrac{\beta}{1-\beta} < \dfrac{\alpha}{1-\alpha}$

$$\boxed{\text{Posterior} \propto \text{Likelihood} * \text{Prior}.}$$

Example : Predicting Commute Time :
- You moved to new apartment. Friend said commute time is $30 \pm 10$ mins.
- You drove a few times, and time = $\{25, 45, 30, 50\}$
- How should you predict commute time?

Prior : Assume $P(\theta) \sim N(\mu_0, \sigma_0^2)$, $\mu_0 = 30$, $\sigma_0 = 10$.

$$= \frac{1}{2\pi\sigma_0} \exp\left(\frac{-(\theta-\mu_0)^2}{2\sigma_0^2}\right) \propto \exp\left(\frac{-(\theta-\mu_0)^2}{2\sigma_0^2}\right)$$

Likelihood : Based on observation / data.

observe : $X_1, \cdots, X_n$. Assume noise in observation.

$$x_i = \theta + \sigma_1 \xi_i, \quad \xi_i \sim N(0,1), \quad \sigma_1 = 5.$$

$$P(x_i \mid \theta) \sim N(\theta, \sigma_1^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{(X_i-\theta)^2}{2\sigma_1^2}\right) \propto \exp\left(\frac{-(X_i-\theta)^2}{2\sigma_1^2}\right)$$

ith observation ↖  ↗ theoretical value

Posterior :

$$P(\theta \mid D) = \frac{P(D\mid\theta)P(\theta)}{P(D)} \propto P(D\mid\theta)P(\theta).$$

↑ doesn't matter much; normalization constant.

$$= \prod_{i=1}^{n} P(x_i \mid \theta)\, P(\theta)$$

$P(D\mid\theta)$ = data.

$$\propto \left[ \prod_{i=1}^{n} \exp\left( \frac{-(x_i - \theta)^2}{2\sigma_1^2} \right) \right] \exp\left( \frac{-(\theta - \mu_0)^2}{2\sigma_0^2} \right)$$

$$\overbrace{\theta^2 - 2x_i\theta + x_i^2} \qquad = \theta^2 - 2\mu_0\theta + \mu_0^2$$

$$\propto \exp\left[ -\sum_{i=1}^{n} \left( \frac{(\theta - x_i)^2}{2\sigma_1^2} \right) - \frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right]$$

$$= \exp\left[ -\sum_{i=1}^{n} \left( \frac{\theta^2 - 2x_i\theta + x_i^2}{2\sigma_1^2} \right) - \frac{\theta^2 - 2\mu_0\theta + \mu_0^2}{2\sigma_0^2} \right]$$

$$= \exp\left[ -\frac{1}{2}\left[ \sum_{i=1}^{n} \left( \frac{\theta^2 - 2x_i\theta + x_i^2}{\sigma_1^2} \right) + \frac{\theta^2 - 2\mu_0\theta + \mu_0^2}{\sigma_0^2} \right] \right]$$

$$= \exp\left[ -\frac{1}{2}\left[ \theta^2 \overbrace{\left( \sum_{i=1}^{n} \left( \frac{1}{\sigma_1^2} \right) + \frac{1}{\sigma_0^2} \right)}^{A} - 2\theta \overbrace{\left( \sum_{i=1}^{n} \left( \frac{x_i}{\sigma_1^2} \right) + \frac{\mu_0}{\sigma_0^2} \right)}^{B} \right.\right.$$

$$\left.\left. + \underbrace{\left( \sum_{n=1}^{n} \left( \frac{x_i^2}{\sigma_1^2} \right) + \frac{\mu_0^2}{\sigma_0^2} \right)}_{C} \right] \right]$$

$$= \exp\left( -\frac{1}{2}\left( A\theta^2 - 2B\theta + C \right) \right), \qquad \text{where}$$

$$A = \sum_{i=1}^{n} \left( \frac{1}{\sigma_1^2} \right) + \frac{1}{\sigma_0^2} = \frac{n}{\sigma_1^2} + \frac{1}{\sigma_0^2}$$

$$B = \sum_{i=1}^{n} \left( \frac{x_i}{\sigma_1^2} \right) + \frac{\mu_0}{\sigma_0^2}$$

$$C = \text{constant (bc } \theta \text{ not involved)}.$$

// Recall Gaussian distribution

$$\propto \exp\left[\left(1-\frac{1}{2}\right)\frac{1}{\sigma^2}(X-\mu)^2\right].$$

$$= \exp\left(-\frac{1}{2}\underset{\underset{\text{variance}}{\frac{1}{\uparrow}}}{A}\left(\theta - \underset{\underset{\text{mean}}{\uparrow}}{\frac{B}{A}}\right)^2 + \text{const.}\right) \sim N\left(\frac{B}{A}, \frac{1}{A}\right)$$

$$\mu_{\text{posterior}} = \frac{B}{A} = \frac{\left(\sum_{i=1}^{n}\frac{x_i}{\sigma_i^2} + \frac{\mu_0}{\sigma_0^2}\right)}{\left(\frac{n}{\sigma_i^2} + \frac{1}{\sigma_0^2}\right)}$$

$$\sigma^2_{\text{posterior}} = \frac{1}{A} = \left(\frac{n}{\sigma_i^2} + \frac{1}{\sigma_0^2}\right)^{-1}$$

If $n = 0$ (no data), $\mu_p = \dfrac{\mu_0/\sigma_0^2}{1/\sigma_0^2} = \mu_0$.
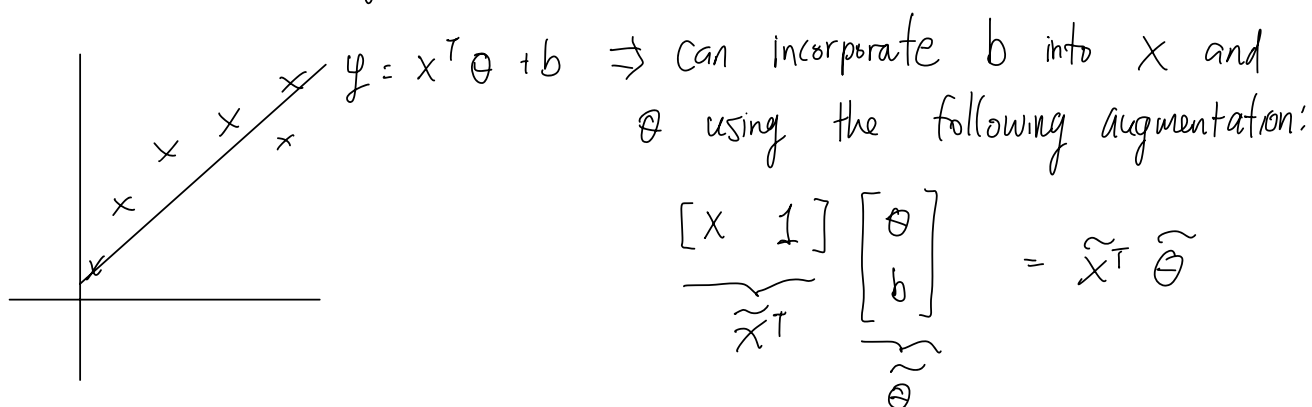
$$\sigma_p^2 = \sigma_0^2.$$

If $n = \infty$ (infinite data), $\mu_p \approx \dfrac{\sum\limits_{i=1}^{n} \dfrac{x_i}{\sigma_1^2} + 0}{\dfrac{n}{\sigma_1^2} + 0} \approx \dfrac{1}{n} \sum\limits_{i=1}^{n} x_i$

$$\sigma_p^2 \approx \left( \dfrac{n}{\sigma_1^2} \right)^{-1} = \dfrac{\sigma_1^2}{n}$$

If $n > 0$ (have data), $\mu_p$ = weighted sum of data and prior, where

$$\text{weight} = \dfrac{1}{\sigma_1^2}$$

# Bayesian Linear Regression

**Bayesian Linear Regression:** Goal is to use Bayesian Inference to solve Linear Regression.

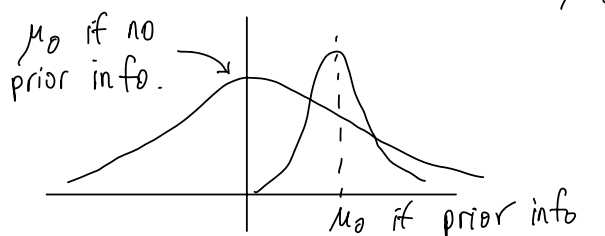- Given data points $\{x_i, y_i\}_{i=1}^{n}$, want to find $\theta$, such that $y \approx x^T \theta$.



$y = x^T \theta + b \Rightarrow$ can incorporate $b$ into $x$ and $\theta$ using the following augmentation:

$$\underbrace{[X \quad 1]}_{\widetilde{X}^T} \underbrace{\begin{bmatrix} \theta \\ b \end{bmatrix}}_{\widetilde{\theta}} = \widetilde{X}^T \widetilde{\theta}$$

## Least Square Method:

$$\hat{\theta} = \arg\min_{\theta} \sum_{i=1}^{n} (y_i - x_i^T \theta)^2$$

## Bayesian Inference: Treat $\theta$ as a random variable

**Prior:** Assume Gaussian Distribution; $P(\theta) \sim N(\mu_0, \sigma_b^2)$

The prior depends on the user. If you have information on $\theta$, then maybe you center $\mu_0$ at that, and make $\sigma$ small. If no prior information, then set $\mu_0 = 0$, $\sigma_2^2 = $ large number.



$\mu_0$ if no prior info.

$\mu_0$ if prior info.

<u>Likelihood</u> : Assume $y_i = x_i^T \theta + \sigma_1 \xi_i$ , where

$\sigma_1$ : variance,

$\xi \sim N(0, 1)$

$\Rightarrow P(\{y_i, x_i\} \mid \theta) = \underbrace{P(y_i \mid x_i, \theta)}_{\text{Gaussian.}} \underbrace{P(x_i)}_{\substack{\text{Constant} \\ \text{wrt } \theta.}}$  // multiplication rule ?

# <u>Posterior</u> :

$P(\theta \mid D) \propto P(D \mid \theta) P(\theta)$

$= \left[ \prod_{i=1}^{n} P(\{y_i, x_i\} \mid \theta) \right] P(\theta)$

$= \left[ \prod_{i=1}^{n} P(y_i \mid x_i, \theta) \underbrace{P(x_i)}_{\text{does not depend on } \theta.} \right] P(\theta)$

$\propto \left[ \prod_{i=1}^{n} P(y_i \mid x_i, \theta) \right] P(\theta)$

prior : can be known or unknown

$\propto \left[ \prod_{i=1}^{n} \exp\left( - \frac{(\overset{\text{actual}}{y_i} - \overset{\text{predicted}}{x_i^T \theta})^2}{2\sigma_1^2} \right) \right] \exp\left[ - \frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right]$

$\underbrace{\qquad\qquad\qquad\qquad}_{P(D\mid\theta) = \prod_{i=1}^{n} P(y_i \mid x_i, \theta) P(\cancel{x_i})} \qquad \underbrace{\qquad}_{P(\theta)}$

$= \exp\left[ - \sum_{i=1}^{n} \left( \frac{(y_i - x_i^T \theta)^2}{2\sigma_1^2} \right) - \frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right]$

$$= \exp\left[ -\frac{1}{2}\left( \sum_{i=1}^{n} \left( \frac{(y_i - x_i^T \theta)^2}{\sigma_i^2} \right) - \frac{(\theta - \mu_\theta)^2}{\sigma_0^2} \right) \right]$$

$$= \exp\left[ -\frac{1}{2}\left( \underbrace{\theta^T A \theta}_{\propto \theta^2} - \underbrace{2B^T\theta}_{\propto \theta} + \text{const} \right) \right] \sim N\left( \underbrace{A^{-1}B}_{\mu_p}, \underbrace{A^{-1}}_{\sigma_p^2} \right)$$

where

$$A = \sum_{i=1}^{n} \frac{x_i x_i^T}{\sigma_i^2} + \frac{I}{\sigma_0^2} \quad \text{Identity matrix.}$$

$$B = \sum_{i=1}^{n} \frac{y_i x_i}{\sigma_i^2} + \frac{\mu_\theta}{\sigma_0^2}$$

$$\theta = \begin{bmatrix} \ \\ \ \\ \ \end{bmatrix} \quad A = \begin{bmatrix} \ & \ \\ \ & \ \end{bmatrix}$$

$$d \times 1 \qquad\qquad d \times d.$$