

Maximum A Posteriori Estimation (MAP)

$$f_{\text{Bayes}}(\vec{x}) = E[Y | \vec{X} = \vec{x}]$$
$$= \int_{\mathcal{Y}} y P_{Y|\vec{X}}(y|\vec{x}) dy$$

If $P_{Y|\vec{X}}$ is a function of some parameter $w \in \mathcal{W}$ then we just need to estimate that parameter

$$\text{MLE: } \arg \max_{w \in \mathcal{W}} \underbrace{p(D|w)} = \prod_{i=1}^n p(z_i|w)$$

"find w that maximizes the likelihood of the data"

$$\text{MAP: } \arg \max_{w \in \mathcal{W}} \underbrace{p(w|D)} = \text{"posterior"}$$

"find w that is the most likely given the data"

$$w_{\text{MAP}} = \arg \max_{w \in \mathcal{W}} p(w|D)$$
$$= \arg \max_{w \in \mathcal{W}} \frac{p(w, D)}{P(D)}$$

prod rule

$$= \underset{w \in W}{\operatorname{argmax}} \frac{p(D|w) p(w)}{P(D)}$$

$$= \underset{w \in W}{\operatorname{argmax}} \underbrace{p(D|w)}_{\text{likelihood}} \underbrace{p(w)}_{\text{prior}}$$

MAP Basics

Ex: $Z_i \sim \mathcal{N}(\mu^*, 1)$ is the i -th person's height

$$p_z(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z-\mu^*)^2}{2}\right)$$

$$\mu \sim \mathcal{N}(160, \sigma^2), \quad p(\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mu-160)^2}{2\sigma^2}\right)$$

$$\mu_{\text{MAP}} = \underset{\mu \in \mathbb{R}}{\operatorname{argmax}} p(D|\mu) p(\mu)$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} -\log(p(D|\mu) p(\mu)) \quad \log(ab) = \log(a) + \log(b)$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} \left[-\log(p(D|\mu)) - \log(p(\mu)) \right]$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{arg\,min}} \left[\sum_{i=1}^n \frac{(z_i - \mu)^2}{2} - \log(p(\mu)) \right]$$

$$\log(p(\mu)) = \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(\mu - 160)^2}{2\sigma^2} \right) \right)$$

$$= \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left(\exp \left(-\frac{(\mu - 160)^2}{2\sigma^2} \right) \right)$$

$$= \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(\mu - 160)^2}{2\sigma^2}$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{arg\,min}} \left[\sum_{i=1}^n \frac{(z_i - \mu)^2}{2} - \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{(\mu - 160)^2}{2\sigma^2} \right]$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{arg\,min}} \left[\underbrace{\sum_{i=1}^n \frac{(z_i - \mu)^2}{2} + \frac{(\mu - 160)^2}{2\sigma^2}}_{g(\mu)} \right]$$

$$\frac{dg}{d\mu}(\mu) = -\sum_{i=1}^n (z_i - \mu) + \frac{(\mu - 160)}{\sigma^2} = 0$$

$$\Rightarrow -\sum_{i=1}^n z_i + n\mu + \frac{\mu}{\sigma^2} - \frac{160}{\sigma^2} = 0$$

$$\Rightarrow \sum_{i=1}^n z_i + \frac{160}{\sigma^2} = \left(n + \frac{1}{\sigma^2}\right)\mu$$

$$\Rightarrow \mu_{\text{MAP}} = \mu = \frac{\sum_{i=1}^n z_i + \frac{160}{\sigma^2}}{n + \frac{1}{\sigma^2}}, \quad \mu_{\text{MLE}} = \frac{\sum_{i=1}^n z_i}{n}$$

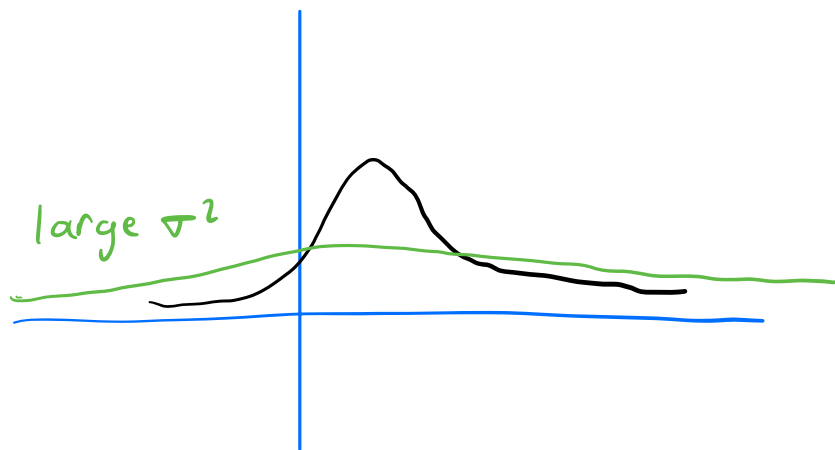
if n is large then $\mu_{\text{MAP}} \approx \frac{\sum_{i=1}^n z_i}{n} = \mu_{\text{MLE}}$

if σ^2 is small then $\mu_{\text{MAP}} \approx \frac{\frac{160}{\sigma^2}}{\frac{1}{\sigma^2}} = 160$

if σ^2 is large then $\mu_{\text{MAP}} \approx \mu_{\text{MLE}}$

if $p(\mu) \approx C \in \mathbb{R}$

then $\mu_{\text{MAP}} = \mu_{\text{MLE}}$



Estimating \vec{w} for $P_{Y|\vec{X}}$

$$D = ((\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n)) \in (\mathcal{X} \times \mathcal{Y})^n, P_D, P_D$$

(\vec{X}_i, Y_i) are i.i.d with $P_{\vec{X}, Y}$ and $P_{\vec{X}, Y}$

$$\text{Assume } Y_i | \vec{X}_i = \vec{x}_i \sim \mathcal{N}(\vec{x}_i^T \vec{w}^*, 1)$$

$$P_{Y|\vec{X}=\vec{x}}(y|\vec{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \vec{x}^T \vec{w}^*)^2}{2}\right)$$

Assume $w_j \sim \mathcal{N}(0, \frac{1}{\lambda})$ are i.i.d. for $j \in \{1, \dots, d\}$

and $w_0 \sim \mathcal{N}(0, \sigma^2)$ for very large σ

$\approx \text{Uniform}(-a, a)$ for large a

w_0 is independent of w_j for all $j \in \{1, \dots, d\}$

Calculating \vec{w}_{MAP} :

$$\vec{w}_{MAP} = \underset{\vec{w} \in \mathbb{R}^{d+1}}{\operatorname{argmax}} p(\vec{w} | D)$$

$$= \underset{\vec{w} \in \mathbb{R}^{d+1}}{\operatorname{argmax}} p(D | \vec{w}) p(\vec{w})$$

$$= \underset{\vec{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} -\log(p(D | \vec{w}) p(\vec{w}))$$

$$= \underset{\vec{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \left[-\log(p(D | \vec{w})) - \log(p(\vec{w})) \right]$$

$$= \underset{\vec{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \left[\sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2} - \log(p(\vec{w})) \right]$$

$$\log(p(\vec{w})) = \log(p(w_0, w_1, \dots, w_d))$$

$$\text{Independence} \Rightarrow \log(p(w_0) p(w_1) \dots p(w_d))$$

$$= \log(p(w_0) \prod_{j=1}^d p(w_j))$$

$$= \log(p(w_0)) + \sum_{j=1}^d \log(p(w_j))$$

$$= \log(p(w_0)) + \sum_{j=1}^d \log\left(\frac{1}{\sqrt{2\pi/\lambda}} \exp\left(-\frac{w_j^2}{2/\lambda}\right)\right)$$

$$= \log(p(w_0)) + \sum_{j=1}^d \left[\log\left(\frac{1}{\sqrt{2\pi/\lambda}}\right) - \frac{w_j^2}{2/\lambda} \right]$$

$$= \log(p(w_0)) + d \log\left(\frac{1}{\sqrt{2\pi/\lambda}}\right) - \sum_{j=1}^d \frac{\lambda w_j^2}{2}$$

$$= \operatorname{argmin}_{\vec{w} \in \mathbb{R}^{d+1}} \left[\sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2} - \log(p(w_0)) - d \log\left(\frac{1}{\sqrt{2\pi/\lambda}}\right) + \sum_{j=1}^d \frac{\lambda w_j^2}{2} \right]$$

$$= \operatorname{argmin}_{\vec{w} \in \mathbb{R}^{d+1}} \left[\underbrace{\sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2}}_{= \frac{n}{2} \hat{L}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d w_j^2}_{\text{almost regularizer}} \right]$$

$$= \operatorname{argmin}_{\vec{w} \in \mathbb{R}^{d+1}} \left[\frac{n}{2} \hat{L}_\lambda(\vec{w}) \right] \quad \frac{1}{n} \sum_{i=1}^n (y_i - \vec{x}_i^T \vec{w})^2 + \frac{\lambda}{n} \sum_{j=1}^d w_j^2$$

$$= \operatorname{argmin}_{\vec{w} \in \mathbb{R}^{d+1}} \left[\hat{L}_\lambda(\vec{w}) \right] = \vec{w}_{\text{MAP}} = \hat{\vec{w}}_\lambda$$

$$f_{\text{Bayes}}(\vec{x}) = E[Y | \vec{X} = \vec{x}]$$

$$= \int_{\mathcal{Y}} y p_{Y|\vec{X}}(y|\vec{x}) dy$$

$$\approx \int_{\mathcal{Y}} y p(y|\vec{x}, \vec{w}_{\text{MAP}}) dy$$

$$= E[Y' | \vec{X} = \vec{x}] \quad \text{where } Y' | \vec{X} = \vec{x} \sim \mathcal{N}(\vec{x}^T \vec{w}_{\text{MAP}}, 1)$$

$$= \vec{x}^T \vec{w}_{\text{MAP}}$$

$$= \vec{x}^T \hat{\vec{w}}_\lambda = \hat{f}$$

Assume $w_j \sim \text{Laplace}(0, 1/\lambda)$ are i.i.d for $j \in \{1, \dots, d\}$

and $w_0 \sim \text{Laplace}(0, b)$ for very large b

$\approx \text{Uniform}(-a, a)$ for large a

w_0 is independent of w_j for all $j \in \{1, \dots, d\}$

$$\vec{w}_{\text{MAP}} = \arg \max_{\vec{w} \in \mathbb{R}^{d+1}} p(\vec{w} | D)$$

$$= \arg \min_{\vec{w} \in \mathbb{R}^{d+1}} \left[\sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2} + \lambda \sum_{j=1}^d |w_j| \right]$$

Lasso regression