

①

Statistical Estimation: MLE, MAP, Full Bayesian

MLE

If India & Australia plays tomorrow, will India win or lose?

Past data: Win, Loss, Loss, Win, Win

$$\text{So, } p(\text{India Wins}) = \frac{3}{5}$$

Now we'll fit a statistical model & estimate MLE, and see whether we can obtain the same answer.

There can be two outcomes: $\{ \text{Win, Loss} \}$

View this as a random variable with 2 states:

$$Y = \{0, 1\}$$

Let the hypothesis class is:

$$Y \sim \text{Bernoulli}(\theta)$$

$$Y = \begin{cases} 1, & \text{prob } \theta \\ 0, & \text{prob. } (1-\theta) \end{cases}$$

↓
the estimation problem is to find θ .

The dataset is: $D = \{1, 0, 0, 1, 1\}$

Some trivial choices: $\hat{\theta} = 0$
 $\hat{\theta} = 1$ } bad estimators.

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \underbrace{P(D|\theta)}_{\text{Likelihood}}$$

[Likelihood is a function of the parameters θ]

- Let $D = \{1\} \Rightarrow L(\theta) = P(Y=1|\theta) = \theta$
- $D = \{1, 0\} \Rightarrow L(\theta) = P(Y=1|\theta) P(Y=0|\theta) = \theta(1-\theta)$
- $D = \{1, 1, 1\} \Rightarrow L(\theta) = \theta^3$

[I.I.D. samples]
(independency)

In general: $L(\theta) = \theta^{\alpha_H} (1-\theta)^{\alpha_T}$

$\alpha_H \rightarrow \# \text{ of wins}$
 $\alpha_T \rightarrow \# \text{ of losses}$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \log L(\theta)$$

[monotonic f_h]

$$\Rightarrow \hat{\theta}_{MLE} = \arg \max_{\theta} [\alpha_H \log \theta + \alpha_T \log (1-\theta)]$$

$$\frac{d}{d\theta} [\alpha_H \log \theta + \alpha_T \log (1-\theta)] = 0$$

$$\Rightarrow \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} = 0 \quad (\text{assume } \theta \neq 1, \neq 0)$$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

So in our dataset, $\alpha_H = 3, \alpha_T = 2$

$$\text{Hence, } \boxed{\hat{\theta}_{MLE} = \frac{3}{5}}$$

(matches with our previous estimation)

Why Max-Likelihood?

- Leads to "natural" estimators
- MLE is optimal if model-Class (hypothesis) is correct.

Let Y is a discrete random variable:

$$Y = y \in \{1, \dots, K\}$$

Consider log-Likelihood (for some model) average:

$$\frac{1}{N} LL(\theta) = \frac{1}{N} \sum_{i=1}^N \log P(y_i | \theta)$$

$$= \frac{1}{N} [\#(Y=1) \log P(Y=1|\theta) + \#(Y=2) \log P(Y=2|\theta) + \dots]$$

$$\begin{aligned} & \downarrow \\ &= P^*(1) \log P_\theta(1) + \dots \\ &= \sum_{y=1}^K P^*(y) \log P_\theta(y) = \sum_y P^*(y) \log \left[\frac{P_\theta(y)}{P^*(y)} \right] \\ &= \sum_y P^*(y) \log P^*(y) - \sum_y P^*(y) \log \left[\frac{P_\theta(y)}{P^*(y)} \right] \end{aligned}$$

-ve entropy
(doesn't depend on θ)

$$\text{So, } \boxed{\arg \max_{\theta} \frac{1}{N} LL(\theta) = \arg \min_{\theta} KL(P^* || P_\theta)}$$

So when we do MLE, we are trying to get closer to the reality, if we have infinite data.

$$\begin{aligned} & \text{With } \infty \text{ data,} \\ & \lim_{N \rightarrow \infty} \frac{\#(Y=1)}{N} = P(Y=1|\theta^*) \\ & \quad \uparrow \\ & \quad P^*(y) \quad \text{under the "true" model} \\ & \text{Let } P(Y=y|\theta) \Rightarrow P_\theta(y) \end{aligned}$$

MLE is a frequentist ~~estimation~~ estimation.

Our answer was 3/5 as there were 3 wins & 2 losses. What if #win=30, #loss=20? → We'll still have the same answer. MLE estimator can't convey/confident on how much data you have seen.

MAP

We can put a belief on anything (e.g. I know from my heart that India will win). That's Bayesian. In frequentist, θ is constant, it's not random. In Bayesian,

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

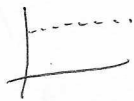
$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|D)$$

(after you observe some data, what do you believe about this)

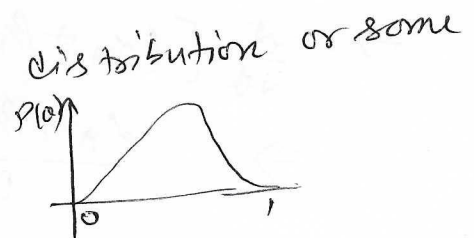
If $P(\theta)$ is constant, then MAP reduce to

MLE:

$$\hat{\theta}_{MAP} = \hat{\theta}_{MLE}$$



Prior: For $P(\theta)$, we need some probability density function between 0 & 1.



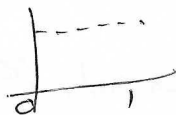
A useful prior:

$$P(\theta) = \frac{\theta^{\beta_H - 1} (1-\theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)}$$

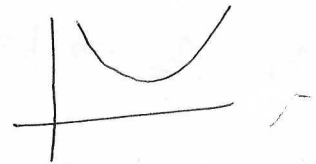
$\sim \text{Beta}(\beta_H, \beta_T)$

hyperparameters.

- $\beta_H = \beta_T = 1$



$\beta_H = 0.6$
 $\beta_T = 0.6$



- $\beta_H = 2.5$
 $\beta_T = 2.5$



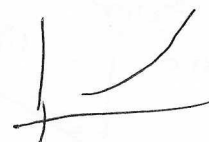
$\beta_H = 0.1$
 $\beta_T = 4.0$



- $\beta_H = 5.0$
 $\beta_T = 5.0$



$\beta_H = 4.0$
 $\beta_T = 0.6$



$$\begin{aligned}\hat{\theta}_{MAP} &= \operatorname{argmax}_{\theta} p(\theta|D) \propto p(D|\theta) p(\theta) \\ &\propto \theta^{\alpha_H} (1-\theta)^{\alpha_T} \theta^{\beta_H-1} (1-\theta)^{\beta_T-1} \\ &\propto \theta^{\alpha_H+\beta_H-1} (1-\theta)^{\alpha_T+\beta_T-1} \\ &\propto \text{Beta}(\alpha_H+\beta_H, \alpha_T+\beta_T)\end{aligned}$$

Conjugate prior

Now let's derive the MAP:

$$\begin{aligned}\hat{\theta}_{MAP} &= \operatorname{argmax}_{\theta} p(\theta|D) = \theta^{\alpha_H} (1-\theta)^{\alpha_T} \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{\text{constant}} \\ &= \operatorname{argmax}_{\theta} \log \left[\theta^{\alpha_H} (1-\theta)^{\alpha_T} \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{\text{constant}} \right]\end{aligned}$$

$$\frac{d}{d\theta} \log p(\theta|D) = \frac{d}{d\theta} [\alpha_H \log \theta + \alpha_T \log (1-\theta) + (\beta_H-1) \log \theta + (\beta_T-1) \log (1-\theta)] = 0$$

$$\Rightarrow \frac{d}{d\theta} [(\alpha_H + \beta_H - 1) \log \theta + (\alpha_T + \beta_T - 1) \log (1-\theta)] = 0$$

$$\Rightarrow \frac{\alpha_H + \beta_H - 1}{\theta} - \frac{\alpha_T + \beta_T - 1}{1-\theta} = 0$$

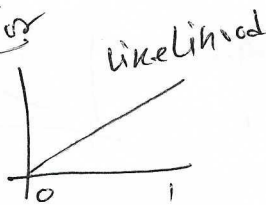
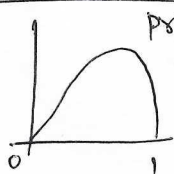
$$\Rightarrow \theta (\alpha_T + \beta_T + \alpha_H + \beta_H - 2) = \alpha_H + \beta_H - 1$$

$$\Rightarrow \hat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_T + \beta_T + \alpha_H + \beta_H - 2}$$

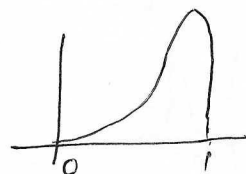
($\beta_H, \beta_T \rightarrow$ pseudo counts — reason ??)

As you see infinite data ($n \rightarrow \infty$), your prior is forgotten!
So MAP converges to MLE with ∞ data. But for small sample size, prior is important.

Effect of prior

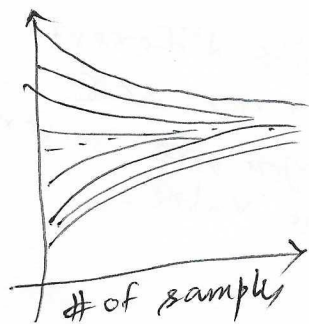


posterior



⑥

$P(x=H)$



You start with different priors, but as you see more data, you are converging to the same value.

Frequentist tool \Rightarrow MLE
Bayesian tool \Rightarrow MAP

MAP vs. Full Bayesian.

We are interested in $P(\theta|D)$

By Bayes' rule: $P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$

Now $P(D) = \int P(D|\theta) P(\theta) d\theta$ \leftarrow marginal likelihood or evidence.
 \leftarrow This integral is the root of difficulty.

Consider every possible parameter setting θ , see how well it explains the dataset D , weight it by how plausible θ was a-priori, and add everything up.

So Bayes must account for all possible models, not just the best one.

Why MAP avoids the hard part?

$$\theta_{MAP} = \arg \max_{\theta} P(D|\theta) P(\theta)$$

\Rightarrow No integral, no normalization, no averaging.
MAP doesn't care about $P(D)$ at all!!

In coin toss, θ is a number, $P(D|\theta)$ is simple, $P(\theta)$ is simple, & the integral has a closed form.

Now imagine complex ML model (like Neural Network). So the integral becomes:

$$\int_{\mathbb{R}^{10^6}} P(D|\theta) P(\theta) d\theta$$

\leftarrow Integration over million dimensional space \rightarrow no closed form soln.

MAP \rightarrow point estimate

Full Bayesian \rightarrow full posterior distribution.

MAP asks \rightarrow what single parameter value is most plausible after seeing the data?

Full Bayesian inference \rightarrow what is the entire distribution of plausible parameters?

Bayes does not ask for best θ ; Bayes asks for the contribution of all θ . MAP ignores them.

Often Bayesian involves approximating the integral.

MLE/MAP example: Gaussian

$$p(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{(x-\mu)^2}{2\sigma^2} \right]$$

Why Gaussian?

- { Central limit theorem.
- { They are easy
- { Closely related to squared loss
- { Mixtures of Gaussians are sufficient to approximate many distributions

Some properties

- Affine transformation

$$X \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$$

$$Y = aX + b$$

- Sum of independent Gaussians

$$X \sim \mathcal{N}(\mu_x, \sigma_x^2)$$

$$Y \sim \mathcal{N}(\mu_y, \sigma_y^2) \Rightarrow Z \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

$$Z = X + Y$$

Learning a Gaussian

Collect bunch of data

↳ hopefully I.I.D. samples.

Learn parameters

- mean
- variance.

⑧

MLEData: $D = \{y_1, \dots, y_N\}$ Model assumption: $y \sim \mathcal{N}(\mu, \sigma^2)$ Density fn: $p(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y-\mu)^2}{2\sigma^2}\right]$

$$(\hat{\mu}_{MLE}, \hat{\sigma}_{MLE}) = \arg \max_{\mu, \sigma} LL(\mu, \sigma)$$

$$\text{Now, } LL(\mu, \sigma) = \log p(D | \mu, \sigma) = \log \prod_{i=1}^N p(y_i | \mu, \sigma) \quad [I.I.D.]$$

$$= \sum_{i=1}^N \log p(y_i | \mu, \sigma)$$

$$= \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \mu)^2}{2\sigma^2}\right) \right]$$

$$= \sum_{i=1}^N \left[-\frac{(y_i - \mu)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]$$

$$= \sum_{i=1}^N -\frac{(y_i - \mu)^2}{2\sigma^2} - \frac{N}{2} \log(2\pi\sigma^2)$$

$$\frac{\partial LL}{\partial \mu} = -\sum_{i=1}^N \frac{2(y_i - \mu)(-1)}{2\sigma^2} = 0 \Rightarrow \sum_{i=1}^N (y_i - \mu) = 0 \quad (\text{assume } \sigma \neq 0)$$

$$\Rightarrow \boxed{\hat{\mu}_{MLE} = \frac{1}{N} \sum y_i}$$

← sample mean is the estimator that maximizes the likelihood.

$$\frac{\partial LL}{\partial \sigma} = -\sum_{i=1}^N \frac{(y_i - \mu)^2}{2} \left(-\frac{2}{\sigma^3}\right) - \frac{N}{2} \frac{1}{2\pi\sigma^2} \cdot 2\pi(2\sigma) = 0$$

$$= \frac{\sum (y_i - \mu)^2}{\sigma^3} - \frac{N}{\sigma} = 0$$

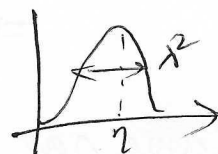
$$\Rightarrow \boxed{\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mu}_{MLE})^2}$$

MAP Assume $\sigma = (\text{unknown})$ constant.
let's put a prior distribution on μ :

$$\mu \sim \mathcal{N}(\eta, \lambda^2)$$

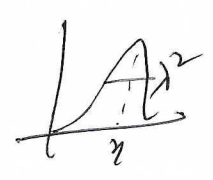
↑
hyper-parameters.

$$p(\mu | \eta, \lambda) = \frac{1}{\sqrt{2\pi}\lambda} \exp\left[-\frac{(\mu - \eta)^2}{2\lambda^2}\right]$$




$$\begin{aligned}
 \hat{\mu}_{MAP} &= \arg \max_{\mu} \log P(\mu|D) \\
 &= \arg \max_{\mu} \log \frac{P(D|\mu) P(\mu|\eta, \lambda^2)}{P(D)} \\
 &= \arg \max_{\mu} \underbrace{\log P(D|\mu)}_{\log \text{ likelihood}} + \underbrace{\log P(\mu|\eta, \lambda^2)}_{\log \text{ prior}} - \underbrace{\log P(D)}_{\text{constant w.r.t. } \mu}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \mu} (\log P(\mu|D)) &= \frac{\partial LL}{\partial \mu} + \frac{\partial}{\partial \mu} \log P(\mu|\eta, \lambda^2) \\
 &= \frac{1}{\sigma^2} \sum_{i=1}^N (y_i - \mu) + \frac{\partial}{\partial \mu} \left[-\frac{(\mu - \eta)^2}{2\lambda^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] \\
 &= \frac{1}{\sigma^2} \sum_{i=1}^N (y_i - \mu) - \frac{\mu - \eta}{\lambda^2} = 0 \\
 \Rightarrow \hat{\mu}_{MAP} &= \frac{\sum \frac{y_i}{\sigma^2} + \frac{\eta}{\lambda^2}}{\frac{N}{\sigma^2} + \frac{1}{\lambda^2}}
 \end{aligned}$$

As $\lambda \rightarrow \infty$ in $P(\mu|\eta, \lambda^2) = \frac{1}{\sqrt{2\pi}\lambda} \exp\left(-\frac{(\mu - \eta)^2}{2\lambda^2}\right)$ 

nearly uniform distribution

So, $\hat{\mu}_{MAP} = \frac{\sum \frac{y_i}{\sigma^2} + (\rightarrow 0)}{\frac{N}{\sigma^2} + (\rightarrow 0)} \Rightarrow \hat{\mu}_{MLE}$

As $\lambda \rightarrow 0$, $P(\mu|\eta, \lambda^2) \rightarrow$  (very strong prior) (delta fn)

$$\hat{\mu}_{MAP} = \frac{\frac{\sum y_i}{\lambda^2} + \eta}{\frac{N}{\lambda^2} + 1}$$

$\Rightarrow \hat{\mu}_{MAP} \rightarrow \eta$ (prior mode)

Same can be obtained from variance as well.