Consider a sample of n iid random variables $X_1, X_2, ..., X_n$.

Maximum What parameter θ maximizes the likelihood Likelihood of our observed data Estimator $(X_1, X_2, ..., X_n)$? (MLE)

 $L(\theta) = f(X_1, X_2, ..., X_n | \theta)$ $= \prod_{i=1}^{n} f(X_i | \theta)$

 $\theta_{MLE} = \arg\max_{\theta} f(X_1, X_2, ..., X_n | \theta)$ likelihood of data

Observations:

- MLE determines θ value that maximizes the probability of observing the sample.
- If we're estimating θ , couldn't we just maximize the probability of θ ?

Today: Bayesian estimation using the Bayesian definition of probability!

Lisa Yan, Chris Piech, Mehran Sahami, and Jerry Cain, CS109, Winter 2023

Stanford University 3

Maximum A Posteriori (MAP) Estimator

Not Review! New!

Consider a sample of n iid random variables $X_1, X_2, ..., X_n$.

Maximum Likelihood Estimator

(MLE)

What parameter θ maximizes the likelihood

of our observed data

 $(X_1, X_2, ..., X_n)$?

$$L(\theta) = f(X_1, X_2, ..., X_n | \theta)$$
$$= \prod_{i=1}^{n} f(X_i | \theta)$$

 $\theta_{MLE} = \arg\max_{\theta} f(X_1, X_2, ..., X_n | \theta)$ likeliho\psi of data)

Maximum a Posteriori

Given the sample data $(X_1, X_2, ..., X_n)$,

what is the most probable

(MAP) parameter θ ? Estimator

 $\theta_{MAP} = \underset{\theta}{\operatorname{arg max}} f(\underline{\theta} | X_1, X_2, \dots, X_n)$ posterior distribution

of θ

Maximum A Posteriori (MAP) Estimator

Consider a sample of n iid random variables $X_1, X_2, ..., X_n$.

def The Maximum a Posteriori (MAP) Estimator of θ is the value of θ that maximizes the **posterior** distribution of θ .

$$\theta_{MAP} = \arg\max_{\theta} f(\theta|X_1, X_2, \dots, X_n)$$

Intuition with Bayes' Theorem:

After seeing

data, posterior

belief of θ

posterior

 $L(\theta)$, probability of data given parameter θ

likelihood prior

 $P(\theta | \text{data}) = \frac{P(\text{data} | \theta)P(\theta)}{P(\text{data})}$

both the prior and the posterior

nitice that

Before seeing data, prior belief of θ

Lisa Yan, Chris Piech, Mehran Sahami, and Jerry Cain, CS109, Winter 2023

Stanford University 5

Solving for θ_{MAP}

- Observe data: $X_1, X_2, ..., X_n$, all iid
- Let likelihood be same as MLE: $f(X_1, X_2, ..., X_n | \theta) = \prod f(X_i | \theta)$
- Let the prior distribution of θ be $g(\theta)$.

posterior distribution

likelihood function

priori distribution

 $\theta_{MAP} = \arg\max_{\theta} f(\theta | X_1, X_2, ..., X_n) = \arg\max_{\theta} \frac{f(X_1, X_2, ..., X_n | \theta)g(\theta)}{h(X_1, X_2, ..., X_n)}$ (Bayes' Theorem)

 $= \arg\max_{\theta} \frac{g(\theta) \prod_{i=1}^{n} f(X_i | \theta)}{h(X_1, X_2, \dots, X_n)} \xrightarrow{\theta 와 관련없음} \text{ (independence)}$

 $= \arg\max_{\theta} g(\theta) \prod_{i=1}^{n} f(X_i | \theta) \qquad (1/h(X_1, X_2, ..., X_n) \text{ is a positive constant w.r.t. } \theta)$ $= \arg\max_{\theta} \left(\log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta) \right)$ Stanford University 6



Stanford University 6

θ_{MAP} : Interpretation 1

• Observe data: $X_1, X_2, ..., X_n$, all iid

• Let likelihood be same as MLE:
$$f(X_1, X_2, ..., X_n | \theta) = \prod_{i=1}^n f(X_i | \theta)$$

• Let the prior distribution of θ be $g(\theta)$.

θ_{MAP} : Interpretation 2

• Observe data: $X_1, X_2, ..., X_n$, all iid

• Let likelihood be same as MLE:
$$f(X_1, X_2, ..., X_n | \theta) = \prod_{i=1}^n f(X_i | \theta)$$

• Let the prior distribution of θ be $g(\theta)$.

$$\theta_{MAP} = \arg\max_{\theta} f(\theta|X_1, X_2, ..., X_n) = \arg\max_{\theta} f(\theta|X_1, X_2, ..., X_n) = \arg\max_{\theta} \frac{g(\theta) \prod_{i=1}^n f(X_i|\theta)}{h(X_1, X_2, ..., X_n)}$$
 (independence)
$$= \arg\max_{\theta} g(\theta) \prod_{i=1}^n f(X_i|\theta)$$
 (1/h(X_1, X_2, ..., X_n) is a positive constant w.r.t. θ)
$$= \arg\max_{\theta} \left(\log g(\theta) + \sum_{i=1}^n \log f(X_i|\theta)\right)$$
 (1/h(X_1, X_2, ..., X_n) is a positive constant w.r.t. θ)
$$\theta_{MAP} = \max_{\theta} \left(\log g(\theta) + \sum_{i=1}^n \log f(X_i|\theta)\right)$$
 (1/h(X_1, X_2, ..., X_n) is a positive constant w.r.t. θ)

Mode: A statistic of a random variable

The mode of a random variable X is defined as:

arg max p(x) $\underset{x}{\arg\max} f(x) \quad \underset{\text{PDF } f(x))}{\text{(X continuous,}}$ (X discrete, PMF p(x)

- Intuitively: The value of X that is "most likely".
- Note that some distributions may not have a unique mode (e.g., Uniform distribution, or Bernoulli(0.5))

 $\theta_{MAP} = \underset{\theta}{\operatorname{arg\,max}} f(\theta | X_1, X_2, \dots, X_n)$

 θ_{MAP} is the most likely θ given the data X_1, X_2, \dots, X_n .

Stanford University 9

Back to our happy Laplace

Consider our previous 6-sided die.

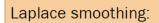
- Roll the dice n = 12 times.
- 3 ones, 2 twos, 0 threes, 3 fours, 1 fives, 3 sixes Observe:

Recall
$$\theta_{MLE}$$
: $p_1 = 3/12, p_2 = 2/12, p_3 = 0/12, p_4 = 3/12, p_5 = 1/12, p_6 = 3/12$

What are your Laplace estimates for each roll outcome?

$$p_i = \frac{X_i + 1}{n + m} \qquad \qquad \chi_3 = 0 \Rightarrow \frac{D + 1}{12 + 6} = \frac{1}{18}$$

$$p_1 = 4/18, p_2 = 3/18, p_3 = 1/18,$$
 $p_4 = 4/18, p_5 = 2/18, p_6 = 4/18$



· Easy to implement/remember

Avoids parameter estimation of 0

Conjugate distributions

MAP estimator:

$$\theta_{MAP} = \arg\max_{\theta} f(\theta|X_1, X_2, ..., X_n)$$
 The mode of the posterior distribution of θ

Distribution parameter	Conjugate distribution
Bernoulli p	Beta
Binomial p	Beta
Multinomial p_i	Dirichlet generalization of Beta
Poisson λ	Gamma
Exponential λ	Gamma
Normal μ	Normal
Normal σ^2	Inverse Gamma