

# Bayesian Linear Regression

---

Nipun Batra

July 2, 2019

IIT Gandhinagar

# MLE, MAP, Bayesian

---

# Bayes Rule - 1

- $P(A|B)P(B) = P(B|A)P(A)$
- Let us consider an example from Wikipedia:
  - A particular drug is 99% sensitive and 99% specific
  - i.e. test will produce 99% true positive results for drug users and 99% true negative results for non-drug users
  - 0.5% of people are users of the drug
  - Question: What is the probability that a randomly selected individual with a positive test is a drug user?

## Bayes Rule - 2

- Test will produce 99% true positive results for drug users and 99% true negative results for non-drug users  $\implies$ 
  - $P(\text{Test} = + | \text{User} = \text{Drug}) = 0.99$ , or,  $P(+ | \text{User}) = 0.99$
  - and  $P(- | \overline{\text{User}}) = 0.99$
- 0.5% of people are users of the drug  $\implies P(\text{User}) = 0.005$
- Question: What is the probability that a randomly selected individual with a positive test is a drug user?  
 $\implies P(\text{User} | +) = ?$
- $$P(\text{User} | +) = \frac{P(+ | \text{User})P(\text{User})}{P(+)} =$$
- $$\frac{P(+ | \text{User})P(\text{User})}{P(+ | \text{User})P(\text{User}) + P(+ | \overline{\text{User}})P(\overline{\text{User}})} = \frac{0.99 \times 0.005}{0.99 \times 0.005 + 0.01 \times 0.995} \approx .332$$

## Another example on Bayes rule

# Bayes Rule for Machine Learning

- $P(A|B)P(B) = P(B|A)P(A)$
- Let us consider for a machine learning problem:
  - A = Parameters ( $\theta$ )
  - B = Data ( $\mathcal{D}$ )
- We can rewrite the Bayes rule as:
  - $P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})}$
  - Posterior:
  - Prior:
  - Likelihood
  -

- Likelihood is a function of  $\theta$
- Given a coin flip and 5 H and 1 T, what is more likely:  $P(H) = 0.5$  or  $P(H) = 1$

# Bayesian Learning is well suited for online settings

content...



# Coin flipping

- Assume we do a coin flip multiple times and we get the following observation: {H, H, H, H, H, H, T, T, T, T}: 6 Heads and 4 Tails
- What is  $P(\text{Head})$ ?
- Is your answer: 6/10. Why?

## Coin flipping: Maximum Likelihood Estimate (MLE)

- We have  $\mathcal{D} = \{D_1, D_2, \dots, D_N\}$  for  $N$  observations where each  $D_i \in \{H, T\}$
- Assume we have  $n_H$  heads and  $n_T$  tails,  $n_H + n_T = N$
- Let us have  $P(H) = \theta, P(T) = 1 - \theta$
- We have Likelihood,  $L(\theta) = P(\mathcal{D}|\theta) = P(D_1, D_2, \dots, D_N|\theta)$
- Since observations are i.i.d.,  
$$L(\theta) = P(D_1|\theta).P(D_2|\theta)...P(D_N|\theta)$$

# Coin flipping: Maximum Likelihood Estimate (MLE)

•

$$P(\mathcal{D}_i|\theta) = \begin{cases} \theta, & \text{for } \mathcal{D}_i = H \\ 1 - \theta, & \text{for } \mathcal{D}_i = T \end{cases}$$

- Thus,  $L(\theta) = \theta^{n_H} \times (1 - \theta)^{n_T}$
- Log-Likelihood,  $LL(\theta) = n_H \log \theta + (n_T)(\log(1 - \theta))$
- $\frac{\partial LL(\theta)}{\partial \theta} = \frac{n_H}{\theta} + \frac{n_T}{1-\theta}$
- For maxima, set derivative of LL to zero
- $\frac{n_H}{\theta} + \frac{n_T}{1-\theta} = 0$

$$\theta = \frac{n_H}{n_H + n_T}$$

## Maximum A Posteriori estimate (MAP)

- **MLE does not handle prior knowledge:** What if we know that our coin is biased towards head?
- **MLE can overfit:** What is the probability of heads when we have observed 6 heads and 0 tails?

# Maximum A Posteriori estimate (MAP)

Goal: Maximize the Posterior

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmin}} P(\theta|\mathcal{D})$$
$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmin}} P(\mathcal{D}|\theta)P(\theta)$$

## Prior distributions

# Beta Distribution

# Beta Distribution



## Coin toss: MAP estimate