

# Bayesian Inference

Thursday 24<sup>th</sup> November, 2016

1 Introduction - order of the day

2 Bayesian Inference

3 Conclusion

# Order of the Day

- ▶ diving a bit deeper into statistical inference
- ▶ statistics - some religious aspects
- ▶ estimation - there's more to life than maximum likelihood estimation

# "Probability Arithmetics"

## Probability - simple arithmetics

Two *fair* coins  $X$  and  $Y$  and one *fair* dice  $Z$

# "Probability Arithmetics"

## Probability - simple arithmetics

Two *fair* coins  $X$  and  $Y$  and one *fair* dice  $Z$

▶ sum  $\approx$  logical OR

- What's the probability of the dice landing on 1 or 6

$$P(Z = 1) + P(Z = 6) = 1/6 + 1/6$$

# "Probability Arithmetics"

## Probability - simple arithmetics

Two *fair* coins  $X$  and  $Y$  and one *fair* dice  $Z$

▶ sum  $\approx$  logical OR

- What's the probability of the dice landing on 1 or 6

$$P(Z = 1) + P(Z = 6) = 1/6 + 1/6$$

▶ product  $\approx$  logical AND

- What's the probability of both coins landing on heads?

$$P(X = H) \cdot P(Y = H) = 0.5 \cdot 0.5$$

# Logical OR

$$P(A \text{ or } B) \stackrel{?}{=} P(A) + P(B)$$

# Logical OR

$$P(A \text{ or } B) \stackrel{?}{=} P(A) + P(B)$$

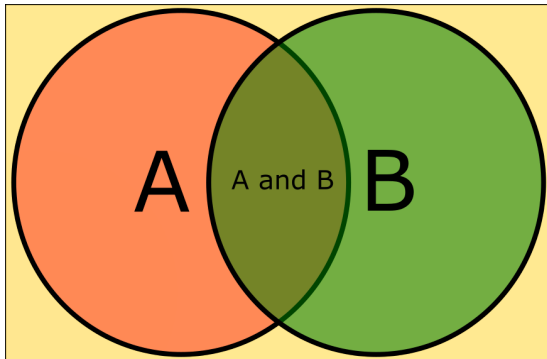
- ▶ When is this true?
  - $A$  and  $B$  are disjoint events
- ▶ And if not?
  - $X, Y$  are fair coins

$$\begin{aligned} P(\text{or } (X = 1, Y = 1)) &= 0.5 + 0.5 \\ &\stackrel{???}{=} 1 \end{aligned}$$



# Logical OR

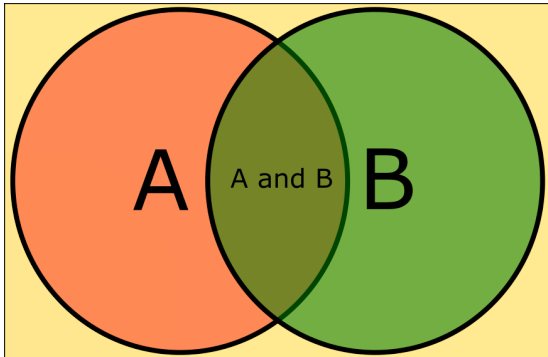
## Darts



areas

# Logical OR

## Darts



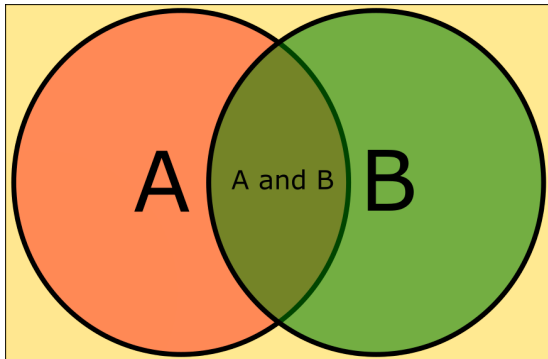
$$P(A) = P(\text{and}(A, \overline{B})) + P(\text{and}(A, B))$$

$$P(B) = P(\text{and}(\overline{A}, B)) + P(\text{and}(A, B))$$

problem?

# Logical OR

## Darts



$$P(A) = P(\text{and}(A, \overline{B})) + P(\text{and}(A, B))$$

$$P(B) = P(\text{and}(\overline{A}, B)) + P(\text{and}(A, B))$$

$$P(\text{or}(A, B)) = P(A) + P(B) - P(\text{and}(A, B))$$

# Joined Probability

$$P(X, Y)$$

- ▶  $X$  someone falling asleep in the first row
- ▶  $Y$  someone falling asleep in the last row

# Joined Probability

$$P(X, Y) \stackrel{?}{=} P(X) \cdot P(Y)$$

# Joined Probability

$$P(X, Y) \stackrel{?}{=} P(X) \cdot P(Y)$$

- ▶  $X$  a fair coin
- ▶  $Y$  the same coin

The probability  $X = Y = 1$ ?

$$P(X = 1) \cdot P(Y = 1) \stackrel{???}{=} 0.25$$

# Joined Probability

$$P(X, Y) \stackrel{?}{=} P(X) \cdot P(Y)$$

- ▶ When is this true?

# Joined Probability

$$P(X, Y) \stackrel{?}{=} P(X) \cdot P(Y)$$

- ▶ When is this true?
  - $X$  is independent of  $Y$



# Joined Probability

$$P(X, Y) \stackrel{?}{=} P(X) \cdot P(Y)$$

- ▶ When is this true?
  - $X$  is independent of  $Y$
- ▶ And if not?
  - $P(X, Y) = P(X | Y) P(Y)$

# Conditional Probability

$P(X | Y)$  - probability of  $X$  conditional on  $Y$   
returning to the previous example:

$$P(X | Y) = \begin{cases} 1, & X = Y \\ 0, & X \neq Y \end{cases}$$

therefore

$$\begin{aligned} P(X = 1, Y = 1) &= P(X = 1 | Y = 1) \cdot P(Y = 1) \\ &= 0.5 \end{aligned}$$

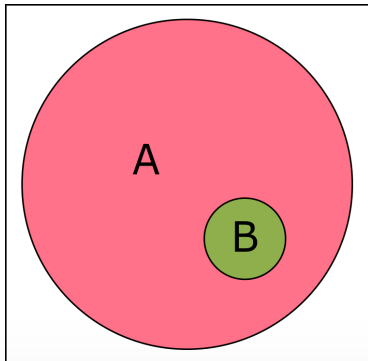
# Conditional Probability

## Another perspective

- ▶  $A$  a Fellow sitting in class
- ▶  $B$  the Fellow is listening

Calculate  $P(B = 1 \mid A = 1)$

- ▶ start with  $P(B = 1, A = 1)$



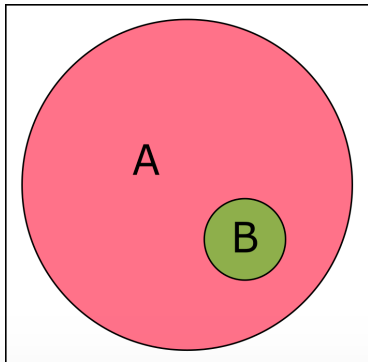
# Conditional Probability

## Another perspective

- ▶  $A$  a Fellow sitting in class
- ▶  $B$  the Fellow is listening

Calculate  $P(B = 1 \mid A = 1)$

- ▶ start with  $P(B = 1, A = 1)$
- ▶ divide by the probability of what's given  $P(A = 1)$



# The Law of Total Probability

description	probability
summery days	0.65
wintery days	0.35
rain on summery days	0.03
rain on wintery days	0.1

what's the probability of rain?

# The Law of Total Probability

description	probability
summery days	0.65
wintery days	0.35
rain on summery days	0.03
rain on wintery days	0.1

what's the probability of rain?

$$\begin{aligned}P(\text{rain}) &= P(\text{rain} \mid \text{summery}) \cdot P(\text{summery}) \\&\quad + P(\text{rain} \mid \text{wintery}) \cdot P(\text{wintery}) \\&= 0.03 \cdot 0.65 + 0.1 \cdot 0.35\end{aligned}$$

# The Law of Total Probability

$$P(X) = \sum_y P(X | Y = y) P(Y = y)$$

# Bayes Theorem

- ▶ We know that

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

- ▶ Can we compute  $P(Y | X)$  from

$$P(X), P(Y), P(X | Y)?$$



# Bayes Theorem

$$P(Y | X) = \frac{P(X, Y)}{P(X)}$$

# Bayes Theorem

$$\begin{aligned} P(Y | X) &= \frac{P(X, Y)}{P(X)} \\ &= \frac{P(X | Y) P(Y)}{P(X)} \end{aligned}$$

# Bayes Theorem

$$P(Y | X) = \frac{P(X | Y) P(Y)}{P(X)}$$

# Bayesian Vs. Frequentist Statistics



**Question:** What's the probability of rain today?

# Bayesian Vs. Frequentist Statistics

**Question:** What's the probability of rain today?

- ▶ We can estimate
- ▶ Given  $N$  days
  - we can estimate what's the proportion of rainy days
  - rain on a given day is deterministic
  - **no probability**

What's your opinion?

# Bayesian Vs. Frequentist Statistics

## **Bayesian**

Probability can model our belief or uncertainty regarding the world

## **Frequentist**

Probability is a limit of frequency

probability model

1 Introduction - order of the day

2 Bayesian Inference

3 Conclusion

# Maximum Likelihood Revisited

**Example:** estimating coin parameter  $p$

- ▶ coin tossed 10 times - 8 heads, 2 tails
- ▶ how would you estimate  $p$ ?



# Maximum Likelihood Revisited

**Example:** estimating coin parameter  $p$

- ▶ coin tossed 10 times - 8 heads, 2 tails
- ▶ how would you estimate  $p$ ?
- ▶ inside information - previous coins

$$p = \begin{cases} 0.3, & 9 \text{ out of } 10 \\ 0.8, & 1 \text{ out of } 10 \end{cases}$$

- ▶ can you use this information?

# Maximum Likelihood Revisited

**Example:** estimating coin parameter  $p$

- ▶ coin tossed 10 times - 8 heads, 2 tails
- ▶ how would you estimate  $p$ ?
- ▶ inside information - previous coins

$$p = \begin{cases} 0.3, & 9 \text{ out of } 10 \\ 0.8, & 1 \text{ out of } 10 \end{cases}$$

- ▶ can you use this information?
- ▶ MLE - no
- ▶ if  $p$  is a constant - what's the point?

# Maximum Likelihood Revisited

**Example:** estimating coin parameter  $p$

- ▶ if you are willing to go Bayesian, there is a way
- ▶ compute probability (uncertainty) for  $p$  using Bayes Theorem

in this problem	name	formula
our inside information	<b>prior distribution</b>	$P(p)$
our uncertainty regarding $p$ given the data	<b>posterior distribution</b>	$P(p   X)$

# Maximum a Posteriori Estimation

- ▶ an alternative to MLE

$$\theta_{map} = \arg \max_{\theta} P(\theta \mid data)$$

# Maximum a Posteriori Estimation

- ▶ an alternative to MLE

$$\begin{aligned}\theta_{map} &= \arg \max_{\theta} P(\theta \mid data) \\ &= \arg \max_{\theta} \frac{P(data \mid \theta) P(\theta)}{P(data)}\end{aligned}$$

# Maximum a Posteriori Estimation

- ▶ an alternative to MLE

$$\theta_{map} = \arg \max_{\theta} P(\theta \mid data)$$

- ▶ problem

$$\theta_{map} = \arg \max_{\theta} \frac{P(data \mid \theta) P(\theta)}{P(data)}$$

# Maximum a Posteriori Estimation

- ▶ an alternative to MLE

$$\theta_{map} = \arg \max_{\theta} P(\theta \mid data)$$

- ▶ problem

$$\theta_{map} = \arg \max_{\theta} \frac{P(data \mid \theta) P(\theta)}{P(data)}$$

- ▶ data is given - same for all  $\theta$

# Maximum a Posteriori Estimation

$$\theta_{map} = \arg \max_{\theta} P(data | \theta) P(\theta)$$



# MAP - examples

## What season is this?

we observe a week with two days of rain

prior	likelihood
$P(\text{summery}) = 0.65$	$P(\text{rain} \mid \text{summery}) = 0.03$
$P(\text{wintery}) = 0.35$	$P(\text{rain} \mid \text{wintery}) = 0.1$

# MAP - examples

## What season is this?

we observe a week with two days of rain

prior	likelihood
$P(\text{summery}) = 0.65$	$P(\text{rain} \mid \text{summery}) = 0.03$
$P(\text{wintery}) = 0.35$	$P(\text{rain} \mid \text{wintery}) = 0.1$

$$\theta_{\text{map}} = \arg \max_{\text{wintery}, \text{summery}} \begin{cases} 0.65 \cdot 0.03^2 \cdot 0.97^5, & \theta = \text{summery} \\ 0.35 \cdot 0.1^2 \cdot 0.9^5, & \theta = \text{wintery} \end{cases}$$

# MAP - examples

## What season is this?

we observe a week with two days of rain

prior	likelihood
$P(\text{summery}) = 0.65$	$P(\text{rain} \mid \text{summery}) = 0.03$
$P(\text{wintery}) = 0.35$	$P(\text{rain} \mid \text{wintery}) = 0.1$

$$\theta_{\text{map}} = \arg \max_{\text{wintery}, \text{summery}} \begin{cases} 0.0005, & \theta = \text{summery} \\ 0.002, & \theta = \text{wintery} \end{cases}$$

# MAP - examples

Israel<sup>tëch</sup>  
challenge

## Back to the coin

whiteboard

# MAP - examples

## Back to the coin

whiteboard

$$p_{map} = \arg \max_{0.3, 0.8} \begin{cases} 0.00067, & p = 0.8 \\ 0.000029, & p = 0.3 \end{cases}$$

# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$

# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$

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

$$= \arg \max_{\theta} \left( \frac{1}{\sqrt{2\pi\zeta^2}} e^{-\frac{(\theta - \phi)^2}{2\zeta^2}} \right) \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i - \theta)^2}{2\sigma^2}} \right)$$

# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$

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

$$= \arg \max_{\theta} \left( \frac{1}{\sqrt{2\pi\zeta^2}\sqrt{2\pi\sigma^2}} e^{-\frac{(\theta - \phi)^2}{2\zeta^2} - \sum_{i=1}^n \frac{(X_i - \theta)^2}{2\sigma^2}} \right)$$



# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$

$$\begin{aligned}\theta_{map} &= \arg \max_{\theta} \left\{ -\frac{(\theta - \phi)^2}{2\zeta^2} + \sum_{i=1}^n -\frac{(X_i - \theta)^2}{2\sigma^2} \right\} \\ &= \arg \min_{\theta} \left\{ \frac{\sigma^2}{\zeta^2} (\theta - \phi)^2 + \sum_{i=1}^n (X_i - \theta)^2 \right\}\end{aligned}$$

# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$
- ▶ if the problem is nice (this one is)
- ▶ find the unique minimizer

derivative

# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$

$$\frac{d}{d\theta} \left\{ \frac{\sigma^2}{\zeta^2} (\theta - \phi)^2 + \sum_{i=1}^n (X_i - \theta)^2 \right\} \stackrel{!}{=} 0$$

$$\frac{\sigma^2}{\zeta^2} (\theta - \phi) = -n\theta + \sum_{i=1}^n X_i$$

# MAP - examples

## normal - normal

estimating mean height in Israel from samples

$$D = X_1, \dots, X_n$$

- ▶ height -  $H \sim \mathcal{N}(\theta, \sigma^2)$
- ▶ mean height across countries -  $\theta \sim \mathcal{N}(\phi, \zeta^2)$

$$\theta_{map} = \frac{\sum_{i=1}^n X_i + \phi \frac{\sigma^2}{\zeta^2}}{n + \frac{\sigma^2}{\zeta^2}}$$

$$\zeta^2 \leq \sigma^2$$

# Bayesian Inference

## coin example

- ▶ prior -  $p = 0.8$  only 10% of the time
- ▶ intuition -  $p_{mle} = 0.8 \Rightarrow P(p = 0.8)$  will increase
- ▶ update -  $P(p \mid data) \rightarrow P(p)$

## Bayesian Updating

we all do it, all the time

- ▶ bus to class 10min
- ▶ walk to class 15min
- ▶ a bus is quicker 95% of the time
- ▶ last few days - horrible traffic - bus time 25min
- ▶ estimation of - best method of travel - will begin to change

# Is Bayesian Inference Good?

## Is the philosophy good?

- ▶ a matter of opinion
- ▶ probability is a modeling tool
- ▶ don't be an extremist

# Is Bayesian Inference Good?

## Is the practice good?

- ▶ are assumptions good?
- ▶ depends on the assumptions



# Is Bayesian Inference Good?

## Is the practice good?

- ▶ are assumptions good?
- ▶ depends on the assumptions
  - you lost something at your house
  - you think you remember where
  - start from there, expand search
    - you remember correctly - find item faster
    - you remember incorrectly - find item slower

# Is Bayesian Inference Good?

## Is the practice good?

- ▶ are assumptions good?
- ▶ depends on the assumptions
  - you see people going right or left to avoid an obstacle
  - assume best choice is normally distributed
  - choose the mean
  - bump into the obstacle

# Is Bayesian Inference Good?

## Is the practice good?

- ▶ are assumptions good?
- ▶ depends on the assumptions

discussion - assumptions

- 1 Introduction - order of the day
- 2 Bayesian Inference
- 3 Conclusion

# Conclusion

- ▶ joined and conditional probability
- ▶ Bayesian Vs. Frequentist statistics
- ▶ Bayesian Inference
- ▶ MAP estimator