

Comparison of Bayesian and Frequentist Inferences

Xuchen Yao

EMLCT
University of Groningen

11 March 2009

Outline

Introduction

Introducing Tomas Bayes
the Interpretation of Probability

Examples of Comparison

Mean

Proportion

Application



Outline

Introduction

Introducing Tomas Bayes
the Interpretation of Probability

Examples of Comparison

Mean

Proportion

Application



Thomas Bayes

(British mathematician, c. 1702 – 7 April 1761)



Figure: The correct identification of this portrait has been questioned

T. Bayes.

Figure: Signature

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

Figure: Another signature



Outline

Introduction

Introducing Tomas Bayes
the Interpretation of Probability

Examples of Comparison

Mean

Proportion

Application



Two Schools of Views

the Frequentist

- $P(x) \approx \frac{n_x}{n_t}$
- an event's probability is the limit of its relative frequency in a large number of trials.
- a long-run fraction: $P(x) = \lim_{n_t \rightarrow \infty} \frac{n_x}{n_t}$

Bayesian

- $P(H|D) = \frac{P(D|H)P(H)}{P(D)}$
- the probability is a measure of a state of knowledge.
- a degree of believability.

Outline

Introduction

Introducing Tomas Bayes
the Interpretation of Probability

Examples of Comparison

Mean

Proportion

Application



Calculation and Estimation

the average sentence length of a book

calculation

- Digital text
- Sentence segmenter

estimation

- Frequentist
- Bayesian



Calculation and Estimation

the average sentence length of a book

calculation

- Digital text
- Sentence segmenter

estimation

- Frequentist
- Bayesian



the Frequentist Approach

the Catcher in the Rye, J. D. Salinger, 1951

That's the thing about girls {5}. Every time they do something pretty, even if they're not much to look at, or even if they're sort of stupid, you fall half in love with them, and then you never know where the hell you are {38}. Girls {1}. Jesus Christ {2}. They can drive you crazy {5}. They really can {3}.

Frequentist

- 5 sentences, with length [5, 38, 1, 2, 5, 3]
- Central Limit Theorem: as n increases,

$$\bar{X}_n = S_n/n = (X_1 + \cdots + X_n)/n \sim N(\mu, \frac{\sigma^2}{n})$$
- Frequentist estimation: $\mu = 9.0, \sigma = 35.0$



the Frequentist Approach

the Catcher in the Rye, *J. D. Salinger*, 1951

That's the thing about girls {5}. Every time they do something pretty, even if they're not much to look at, or even if they're sort of stupid, you fall half in love with them, and then you never know where the hell you are {38}. Girls {1}. Jesus Christ {2}. They can drive you crazy {5}. They really can {3}.

Frequentist

- 5 sentences, with length [5, 38, 1, 2, 5, 3]
- Central Limit Theorem: as n increases,

$$\bar{X}_n = S_n/n = (X_1 + \cdots + X_n)/n \sim N(\mu, \frac{\sigma^2}{n})$$
- Frequentist estimation: $\mu = 9.0, \sigma = 35.0$



the Bayesian Approach

the prior knowledge of sentence length?

Search "sentence length distribution"

SENTENCE LENGTH DISTRIBUTION OF TOTAL INMATE POPULATION
AS OF JUNE 30, FISCAL YEARS 2004 - 2008

SENTENCE LENGTH	FY 2004		FY 2005		FY 2006		FY 2007		FY 2008	
	Number	Percent	Number	Percent	Number	Percent	Number	Percent	Number	Percent
Shock Incarceration (Court Ordered)	118	0.5%	114	0.5%	134	0.6%	129	0.5%	123	0.5%
YOA	1,519	6.3%	1,397	5.9%	1,412	6.0%	1,374	5.8%	1,442	5.8%
3 Months or Less	2	0.0%	3	0.0%	2	0.0%	2	0.0%	0	0.0%
3 Months 1 Day-1 Year	480	2.0%	562	2.4%	415	1.8%	449	1.9%	611	2.4%
1 Year	431	1.8%	421	1.8%	451	1.9%	434	1.8%	575	2.3%
1 Years 1 Day-2 Years	1,329	5.6%	1,394	5.9%	1,236	5.3%	1,204	5.0%	1,368	5.5%
2 Years 1 Day-3 Years	1,835	7.7%	1,795	7.6%	1,796	7.7%	1,898	7.9%	1,945	7.8%
3 Years 1 Day-4 Years	868	3.6%	912	3.9%	843	3.6%	882	3.7%	911	3.6%
4 Years 1 Day-5 Years	2,423	10.1%	2,215	9.4%	2,164	9.3%	2,256	9.4%	2,349	9.4%
5 Years 1 Day-6 Years	874	3.7%	799	3.4%	769	3.3%	768	3.2%	820	3.3%
6 Years 1 Day-7 Years	877	3.7%	891	3.8%	915	3.9%	954	4.0%	1,000	4.0%
7 Years 1 Day-8 Years	911	3.8%	902	3.8%	913	3.9%	931	3.9%	979	3.9%
8 Years 1 Day-9 Years	330	1.4%	322	1.4%	295	1.3%	318	1.3%	329	1.3%
9 Years 1 Day-10 Years	2,495	10.4%	2,524	10.7%	2,512	10.7%	2,468	10.3%	2,504	10.0%
10 Years 1 Day-20 Years	4,310	18.0%	4,347	18.4%	4,499	19.2%	4,734	19.8%	4,935	19.7%
20 Years 1 Day-30 Years	2,194	9.2%	2,104	8.9%	2,073	8.9%	2,071	8.7%	2,064	8.3%
Over 30 Years	883	3.7%	880	3.7%	881	3.8%	904	3.8%	932	3.7%
Life W/10 Yr. Parole Eligibility	383	1.6%	379	1.6%	357	1.5%	352	1.5%	341	1.4%
Life W/20 Yr. Parole Eligibility	859	3.6%	850	3.6%	834	3.6%	833	3.5%	829	3.3%
Life W/30 Yr. Parole Eligibility	136	0.6%	133	0.6%	131	0.6%	131	0.5%	128	0.5%
Life W/No Parole Eligibility	554	2.3%	605	2.6%	663	2.8%	705	3.0%	771	3.1%
Death	69	0.3%	71	0.3%	65	0.3%	58	0.2%	56	0.2%
Non-Judisdictional Inmates*	43	0.2%	37	0.2%	30	0.1%	32	0.1%	34	0.1%
TOTAL	23,923	100.0%	23,657	100.0%	23,390	100.0%	23,887	100.0%	25,066	100.0%
AVERAGE SENTENCE LENGTH**	11 Years 9 Months		11 Years 9 Months		11 Years 9 Months		11 Years 11 Months		11 Years 8 Months	

Figure: Sentence length distribution of the prisoners:-()



the Bayesian Approach

the prior knowledge of sentence length?

Search "sentence length distribution"

SENTENCE LENGTH DISTRIBUTION OF TOTAL INMATE POPULATION
AS OF JUNE 30, FISCAL YEARS 2004 - 2008

SENTENCE LENGTH	FY 2004		FY 2005		FY 2006		FY 2007		FY 2008	
	Number	Percent	Number	Percent	Number	Percent	Number	Percent	Number	Percent
Shock Incarceration (Court Ordered)	118	0.5%	114	0.5%	134	0.6%	129	0.5%	123	0.5%
YOA	1,519	6.3%	1,397	5.9%	1,412	6.0%	1,374	5.8%	1,442	5.8%
3 Months or Less	2	0.0%	3	0.0%	2	0.0%	2	0.0%	0	0.0%
3 Months 1 Day-1 Year	480	2.0%	562	2.4%	415	1.8%	449	1.9%	611	2.4%
1 Year	431	1.8%	421	1.8%	451	1.9%	434	1.8%	575	2.3%
1 Years 1 Day-2 Years	1,329	5.6%	1,394	5.9%	1,236	5.3%	1,204	5.0%	1,368	5.5%
2 Years 1 Day-3 Years	1,835	7.7%	1,795	7.6%	1,796	7.7%	1,898	7.9%	1,945	7.8%
3 Years 1 Day-4 Years	868	3.6%	912	3.9%	843	3.6%	882	3.7%	911	3.6%
4 Years 1 Day-5 Years	2,423	10.1%	2,215	9.4%	2,164	9.3%	2,256	9.4%	2,349	9.4%
5 Years 1 Day-6 Years	874	3.7%	799	3.4%	769	3.3%	768	3.2%	820	3.3%
6 Years 1 Day-7 Years	877	3.7%	891	3.8%	915	3.9%	954	4.0%	1,000	4.0%
7 Years 1 Day-8 Years	911	3.8%	902	3.8%	913	3.9%	931	3.9%	979	3.9%
8 Years 1 Day-9 Years	330	1.4%	322	1.4%	295	1.3%	318	1.3%	329	1.3%
9 Years 1 Day-10 Years	2,495	10.4%	2,524	10.7%	2,512	10.7%	2,468	10.3%	2,504	10.0%
10 Years 1 Day-20 Years	4,310	18.0%	4,347	18.4%	4,499	19.2%	4,734	19.8%	4,935	19.7%
20 Years 1 Day-30 Years	2,194	9.2%	2,104	8.9%	2,073	8.9%	2,071	8.7%	2,084	8.3%
Over 30 Years	883	3.7%	880	3.7%	881	3.8%	904	3.8%	932	3.7%
Life W/10 Yr. Parole Eligibility	383	1.6%	379	1.6%	357	1.5%	352	1.5%	341	1.4%
Life W/20 Yr. Parole Eligibility	859	3.6%	850	3.6%	834	3.6%	833	3.5%	829	3.3%
Life W/30 Yr. Parole Eligibility	136	0.6%	133	0.6%	131	0.6%	131	0.5%	128	0.5%
Life W/No Parole Eligibility	554	2.3%	605	2.6%	663	2.8%	705	3.0%	771	3.1%
Death	69	0.3%	71	0.3%	65	0.3%	58	0.2%	56	0.2%
Non-Judisdictional Inmates*	43	0.2%	37	0.2%	30	0.1%	32	0.1%	34	0.1%
TOTAL	23,923	100.0%	23,657	100.0%	23,390	100.0%	23,887	100.0%	25,066	100.0%
AVERAGE SENTENCE LENGTH**	11 Years 9 Months		11 Years 9 Months		11 Years 9 Months		11 Years 11 Months		11 Years 8 Months	

Figure: Sentence length distribution of the prisoners:-()



the Prior Knowledge of Sentence Length

log-normal distribution

- Ref: Contributions to the Science of Text and Language: Word Length Studies and Related Issues, By *Peter Grzybek*
- Sentence length has a right skewness. It cannot be approximated by normal distribution. Thus log-normal distribution is proposed and testified.

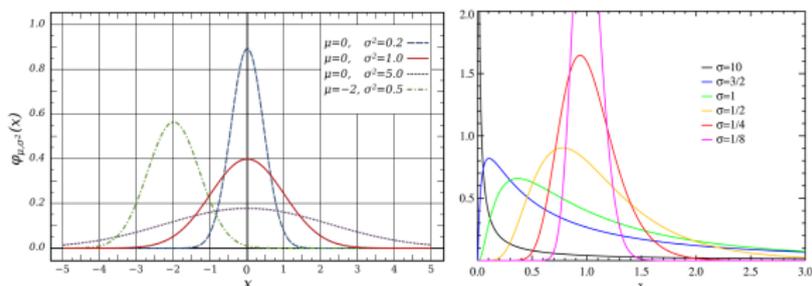


Figure: Normal and log-normal distribution

the Problem Rephrased

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

- $P(X|\mu)$: observation in normal distribution.
- $P(\mu)$: prior knowledge in log-normal distribution.

if $P(\mu) \propto \log N(\mu, \sigma^2)$, then $P(\mu_{\log}) \propto N(\mu, \sigma^2)$

- $P(X_{\log}|\mu_{\log}) \propto e^{-\frac{1}{2\sigma^2/n}(X_{\log}-\mu)^2}$
- $P(\mu_{\log}) \propto e^{-\frac{1}{2s^2}(\mu-m)^2}$
- The posterior distribution will also be normal.

the Problem Rephrased

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

- $P(X|\mu)$: observation in normal distribution.
- $P(\mu)$: prior knowledge in log-normal distribution.

if $P(\mu) \propto \log N(\mu, \sigma^2)$, then $P(\mu_{\log}) \propto N(\mu, \sigma^2)$

- $P(X_{\log}|\mu_{\log}) \propto e^{-\frac{1}{2\sigma^2/n}(X_{\log}-\mu)^2}$
- $P(\mu_{\log}) \propto e^{-\frac{1}{2s^2}(\mu-m)^2}$
- The posterior distribution will also be normal.

the Posterior Probability

$$P(\mu_{log}|X_{log}) \propto P(X_{log}|\mu_{log})P(\mu_{log}) \propto$$

$$e^{-\frac{1}{2\sigma^2}(X_{log}-\mu)^2} e^{-\frac{1}{2s^2}(\mu-m)^2} \propto e^{-\frac{1}{2\sigma^2 s^2 / (\sigma^2 + s^2)} \left(\mu - \frac{\sigma^2 m + s^2 X_{log}}{\sigma^2 + s^2}\right)^2}$$

Plug in the sample mean, we get the posterior mean and variance:

$$m_{pos} = \frac{\frac{\sigma^2}{n} m + s^2 \overline{X_{log}}}{\frac{\sigma^2}{n} + s^2} \quad s_{pos}^2 = \frac{\frac{\sigma^2}{n} s^2}{\frac{\sigma^2}{n} + s^2}$$

- n : the number of samples
- X_{log} : the natural logarithm of the sample
- $\overline{X_{log}}$: the mean of X_{log}
- (m, s^2) : the prior estimation of the mean and variance of sentence length
- σ^2 : the variance of the sentence length we already know



the Posterior Probability

$$P(\mu_{log}|X_{log}) \propto P(X_{log}|\mu_{log})P(\mu_{log}) \propto$$

$$e^{-\frac{1}{2\sigma^2}(X_{log}-\mu)^2} e^{-\frac{1}{2s^2}(\mu-m)^2} \propto e^{-\frac{1}{2\sigma^2 s^2 / (\sigma^2 + s^2)} (\mu - \frac{\sigma^2 m + s^2 X_{log}}{\sigma^2 + s^2})^2}$$

Plug in the sample mean, we get the posterior mean and variance:

$$m_{pos} = \frac{\frac{\sigma^2}{n} m + s^2 \overline{X_{log}}}{\frac{\sigma^2}{n} + s^2} \quad s_{pos}^2 = \frac{\frac{\sigma^2}{n} s^2}{\frac{\sigma^2}{n} + s^2}$$

- n : the number of samples
- X_{log} : the natural logarithm of the sample
- $\overline{X_{log}}$: the mean of X_{log}
- (m, s^2) : the prior estimation of the mean and variance of sentence length
- σ^2 : the variance of the sentence length we already know



Result

Assign values

- n: 5
- X_{log} : $\ln([5, 38, 1, 2, 5, 3])$
- m: 10, s: 10
- σ : 9.73
- Result: $m_{pos}=13.57$

Comparison

- Frequentist: 9.0
- Bayesian: 13.57
- True value: 13.64

What Others Say

by *Charles Annis*

- "probability"_{confidence interval} = long-run fraction having this characteristic.
- "probability"_{credible interval} = degree of believability.
- A frequentist is a person whose long-run ambition is to be wrong 5% of the time.
- A Bayesian is one who, vaguely expecting a horse, and catching a glimpse of a donkey, strongly believes he has seen a mule.
- $P(\text{mule}|\text{donkey}) \stackrel{?}{=} \frac{P(\text{horse})P(\text{donkey}|\text{horse})}{P(\text{donkey})}$



What Others Say

by *Charles Annis*

- "probability"_{confidence interval} = long-run fraction having this characteristic.
- "probability"_{credible interval} = degree of believability.
- A frequentist is a person whose long-run ambition is to be wrong 5% of the time.
- A Bayesian is one who, vaguely expecting a horse, and catching a glimpse of a donkey, strongly believes he has seen a mule.
- $P(\text{mule}|\text{donkey}) \stackrel{?}{=} \frac{P(\text{horse})P(\text{donkey}|\text{horse})}{P(\text{donkey})}$



What Others Say

by *Charles Annis*

- "probability"_{confidence interval} = long-run fraction having this characteristic.
- "probability"_{credible interval} = degree of believability.
- A frequentist is a person whose long-run ambition is to be wrong 5% of the time.
- A Bayesian is one who, vaguely expecting a horse, and catching a glimpse of a donkey, strongly believes he has seen a mule.
- $P(\text{mule}|\text{donkey}) \stackrel{?}{=} \frac{P(\text{horse})P(\text{donkey}|\text{horse})}{P(\text{donkey})}$



Outline

Introduction

Introducing Tomas Bayes
the Interpretation of Probability

Examples of Comparison

Mean

Proportion

Application



the Bayesian Approach

the prior: the Beta distribution

$$f(p; a, b) = \frac{1}{B(a, b)} p^{a-1} (1-p)^{b-1} \sim \text{Beta}(a, b)$$

$$E(p) = \frac{a}{a+b}$$

the posterior: the Beta distribution

$$f(p|x) \propto p^{a+x-1} (1-p)^{b+n-x-1} \sim \text{Beta}(a+x, b+n-x)$$

suppose $a = b = 1$, then

$$\hat{p}_B = \frac{1+x}{2+n}$$

$$MSE(\hat{p}_B) = \left(\frac{1-2\hat{p}_B}{n+2}\right)^2 + \left(\frac{1}{n+2}\right)^2 n \hat{p}_B (1 - \hat{p}_B)$$

still suppose $n = 16, x = 10$, then

$$MSE(\hat{p}_B) = 0.011888431641518061$$



the Bayesian Approach

the prior: the Beta distribution

$$f(p; a, b) = \frac{1}{B(a, b)} p^{a-1} (1-p)^{b-1} \sim \text{Beta}(a, b)$$

$$E(p) = \frac{a}{a+b}$$

the posterior: the Beta distribution

$$f(p|x) \propto p^{a+x-1} (1-p)^{b+n-x-1} \sim \text{Beta}(a+x, b+n-x)$$

suppose $a = b = 1$, then

$$\hat{p}_B = \frac{1+x}{2+n}$$

$$MSE(\hat{p}_B) = \left(\frac{1-2\hat{p}_B}{n+2}\right)^2 + \left(\frac{1}{n+2}\right)^2 n \hat{p}_B (1 - \hat{p}_B)$$

still suppose $n = 16, x = 10$, then

$$MSE(\hat{p}_B) = 0.011888431641518061$$



the Bayesian Approach

the prior: the Beta distribution

$$f(p; a, b) = \frac{1}{B(a, b)} p^{a-1} (1-p)^{b-1} \sim \text{Beta}(a, b)$$

$$E(p) = \frac{a}{a+b}$$

the posterior: the Beta distribution

$$f(p|x) \propto p^{a+x-1} (1-p)^{b+n-x-1} \sim \text{Beta}(a+x, b+n-x)$$

suppose $a = b = 1$, then

$$\hat{p}_B = \frac{1+x}{2+n}$$

$$MSE(\hat{p}_B) = \left(\frac{1-2\hat{p}_B}{n+2}\right)^2 + \left(\frac{1}{n+2}\right)^2 n \hat{p}_B (1 - \hat{p}_B)$$

still suppose $n = 16, x = 10$, then

$$MSE(\hat{p}_B) = 0.011888431641518061$$



Comparison of Proportion

	proportion	MSE
frequentist	0.625	0.015
Bayesian	0.611	0.012

Table: Estimation of proportion and MSE

Authorship Detection

the Federalist Papers, by *Alexander Hamilton*, *James Madison* and *John Jay*

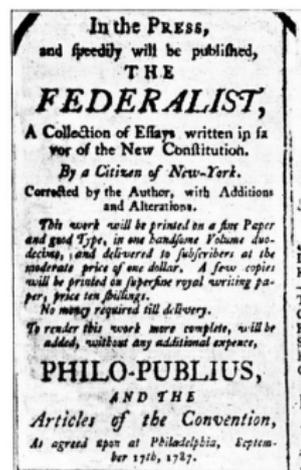


Figure: the
Federalist

- the ratification of the United States Constitution
- 85 articles: Hamilton (51), Madison (29), Jay (5)
- 12 are published under “Publius”.
- Statistical analysis based on word frequencies and writing styles.
- All 12 were written by Madison.

Authorship Detection

the Federalist Papers, by *Alexander Hamilton*, *James Madison* and *John Jay*

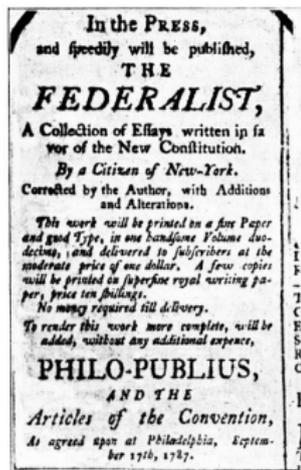


Figure: the
Federalist

- the ratification of the United States Constitution
- 85 articles: Hamilton (51), Madison (29), Jay (5)
- 12 are published under “Publius”.
- Statistical analysis based on word frequencies and writing styles.
- All 12 were written by Madison.

Bayesian POS Tagger

Combining Bayes and HMM

- $\hat{T} = \underset{T \in \tau}{\operatorname{argmax}} P(T|W) = \underset{T \in \tau}{\operatorname{argmax}} \frac{P(T)P(W|T)}{P(W)} = \underset{T \in \tau}{\operatorname{argmax}} P(T)P(W|T)$, where T: possible tags, W: word

- Incorporating the trigram model:

-

$$P(T)P(W|T) = P(t_1)P(t_2|t_1) \prod_{i=3}^n P(t_i|t_{i-2}t_{i-1}) \left[\prod_{i=1}^n P(w_i|t_i) \right]$$

- counting:

- $P(t_i|t_{i-2}t_{i-1}) = \frac{c(t_{i-2}t_{i-1}t_i)}{c(t_{i-2}t_{i-1})}$ and $P(w_i|t_i) = \frac{c(w_i t_i)}{c(t_i)}$

- smoothing...



Bayesian POS Tagger

Combining Bayes and HMM

- $\hat{T} = \underset{T \in \tau}{\operatorname{argmax}} P(T|W) = \underset{T \in \tau}{\operatorname{argmax}} \frac{P(T)P(W|T)}{P(W)} = \underset{T \in \tau}{\operatorname{argmax}} P(T)P(W|T)$, where T: possible tags, W: word

- Incorporating the trigram model:

-

$$P(T)P(W|T) = P(t_1)P(t_2|t_1) \prod_{i=3}^n P(t_i|t_{i-2}t_{i-1}) \left[\prod_{i=1}^n P(w_i|t_i) \right]$$

- counting:

- $P(t_i|t_{i-2}t_{i-1}) = \frac{c(t_{i-2}t_{i-1}t_i)}{c(t_{i-2}t_{i-1})}$ and $P(w_i|t_i) = \frac{c(w_i t_i)}{c(t_i)}$

- smoothing...



Word Sense Disambiguation

Naive Bayesian Classifier

sentence length distribution of {text, prisoners}

- *sentence.n.01*: a string of words satisfying the grammatical rules of a language
- *conviction.n.02*: a final judgment of guilty in a criminal case and the punishment that is imposed
- *prison_term.n.01*: the period of time a prisoner is imprisoned

$$\hat{S} = \underset{S \in \tau}{\operatorname{argmax}} P(\text{Sense} | \text{Context}) = \underset{S \in \tau}{\operatorname{argmax}} \frac{P(S)P(C|S)}{P(C)} =$$

$$\underset{S \in \tau}{\operatorname{argmax}} P(S)P(C|S) = \underset{S \in \tau}{\operatorname{argmax}} [\log P(S) + \log P(C|S)]$$

- Bayesian network for WordNet



Others

- ASR
- OCR
- IR
-



Summary

- Frequentist vs. Bayesian
 - comparison of mean/proportion
 - confidence interval vs. credible interval
- *a priori* knowledge, *a posteriori* probability
- Applications
 - Authorship detection, HMM & Bayes (POS tagger, ASR), WSD, IR, OCR.

