Contents

2 Univariate Probability

6

	462 6	9 Dan aghiothphana.) 1 49 0a883(88 40 (y6 3 4(6)700(s634(.92 482. 483).76 498.8 632.6-143 9 (ÿ)	164D-00630H3(6)(4)90
	4.1	Introduction	46
4	Par	ameter Estimation	46
	3.7	Exercises	43
	3.6	The central limit theorem	42
	3.5	Multivariate normal distributions	41
		3.4.1 The multinomial distribution	40
	3.4	The binomial distribution	40
		3.3.5 Variance of the sum of random variables	39
		3.3.4 Correlation	39
		3.3.3 Covariance and scaling random variables	39
		3.3.2 Covariance	38
		3.3.1 Linearity of the expectation	38

8	B Hierarchical Models				
	8.1	Introduction	163		
	8.2	Parameter estimation in hierarchical models	165		
		8.2.1 Point estimation based on maximum likelihood	167		

	$\Lambda.7$ Combinatorics $\binom{n}{r}$	220
	A.8 Basic matrix algebra	220
	A.9 Miscellaneous notation	
В	More probability distributions and related mathematical constructs 3.1 The gamma and beta functions	223 223

Chapter 2

Univariate Probability

This chapter briefly introduces the fundamentals of univariate probability theory, density

trials. For a frequentist, to say that $P(Heads) = \frac{1}{2}$

2.4 Conditional Probability, Bayes' rule, and Independence

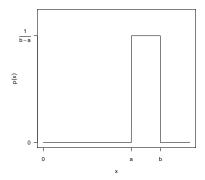
The

2.4.1 Bayes' rule

Bayes' rule

We have been given the value of the two terms in the numerator, but let us leave the

one can apply either to avoid this computation or to drastically simplify it; you will see several examples of these tricks later in the book.



(a) Probability density function

 $1000\ B.C.E.$ to $500\ B.C.E.$ (Beyer, 1986). With only this information, a crude estimate of the dis

2.8 Normalized and unnormalized probability distributions

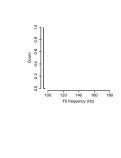
$$F(SOV) = 123$$

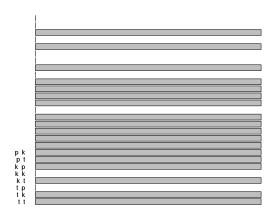
Variance of Bernoulli and uniform distributions

The variance of a Bernoulli-distributed random variable needs to be calculated explicitly, by using the definition in Equation (2.19) and summing over the possible outcomes as in



	SB, DO	SB, IO	DO, SB	DO, IO	IO, SB	IO, DO	Total
Count	478	59	1	3	20	9	570





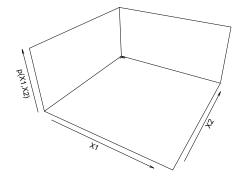
$$P(y|\hat{P}) =$$

3.3 Linearity of expectation, covariance, correlation, and variance sums of random variables

Since the covariance between conditionally independent random variables is zero, it follows that the variance of the sum of pairwise independent random variables is the sum of their variances.

3.4 The binomial distribution

a sequence of r random variables X_1, \ldots, X_r whose joint distribution is characterized by r parameters: a size parameter n denoting the number of trials, and r-1 parameters $x_1, \ldots, x_{r-1}, x_r$ where x_1, \ldots, x_r



$$P\left(X\,=\,k;r,p\right)=\begin{array}{cc} a \\ b \end{array} \left(1-p\right)^c p^d, k \quad \left\{r,r+1,\cdots\right\}$$

for some choice of a, b, c, d. Complete the specification of the distribution (i.e., say what a, b, c, d are) and justify it.

Exercise 3.8: Linearity of expectation

You put two coins in a pouch; one coin is weighted such that it lands heads ⁵

Chapter 4

Parameter Estimation

Thus far we have concerned ourselves primarily with *probability theory*: what events may occur with what probabilities, given a model family and choices for the parameters. This is useful only in the case where we know the precise model family and parameter values for the

for passivization will in fact be realized as a passive.

4.2.1 Consistency

An estimator is consistent if the estimate ^ it constructs is guaranteed to converge to the true parameter value as the quantity of data to which it is applied increases. Figure 4.1 demonstrates that Estimator 1 in our example is consistent: as the sample size increases, the

only the first n/

(Tool, 1949, cited in Language Log by Benjamin Zimmer, 18 October 2007)

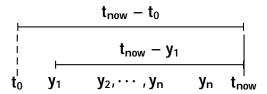
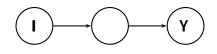
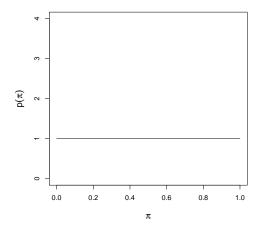


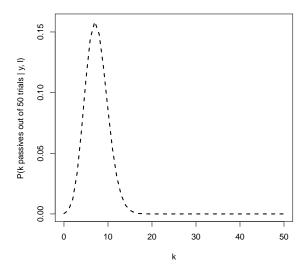
Figure 4.3: The bias of the MLE for uniform distributions





given transitive sentence will be in the passive voice. For Bayesian statistics, we must first

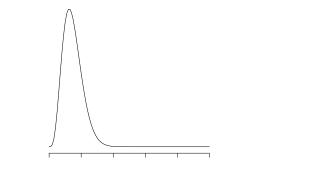
we covered a moment ago) and the *mean*. For our example, the posterior mode is $\frac{4}{}$



```
> plot(density(res[[1]][, 1]), xlab = expression(pi), ylab = expression(paste("p(", pi, ")")))
```

simply expresses th $\boldsymbol{\mu}$	nat observations y ar	e drawn from a norma	l distribution parameterized by

dous modeling flexibility. The only real limits are conceptu



testing in particular works just like any other type of Bayes

In our second hypothesis

$$\begin{split} P\;(y|H_2) &= \prod_{i} P\;(y|_{-i})\,P\;(_{-i}|H_2) \\ &= P \quad y|_{-} = \frac{1}{3} \quad P \quad = \frac{1}{3}|H_2 \quad + \quad P \quad y|_{-} = \frac{2}{3} \quad P \quad = \frac{2}{3}|H_2 \\ &= \frac{6}{4} \quad \frac{1}{3} \quad \frac{2}{3} \quad \times 0.5 + \quad \frac{6}{4} \quad \frac{2}{3} \quad \frac{1}{3} \quad \times 0.5 \\ &= 0.21 \end{split}$$

thus

We use the critical trick of recognizing this integral as a beta function (Section 4.4.2), which gives us:

$$= \frac{6}{4} B(5,$$

5.2.3	Example:	Learning	contextual	contingenci	es in seque	ences

Here's an example, where we will explain the standard error of the mean. Suppose

1. The null hypothesis is true, but we reject it (probability

Quantifying association: odds ratios

In Section 3.3 we already saw one method of quantifying the strength of association between two binary categorical variables: $\frac{1}{2} \int_{-\infty}^{\infty} \frac{1}{2} \left(\frac{1}{2} \int_{-\infty}^{\infty} \frac{1}{2}$

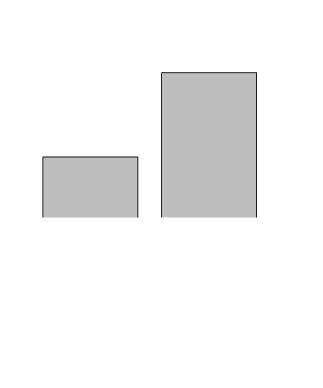
Fisher's exact test

Fisher's exact test applies to 2 \times

3. The linguist	writes up	her research	results a	and sends	them	to a	prestigious	journal.

(c) B B B B A A A A B B B B B A A A A B B B B

The file $spillover_word_rts$

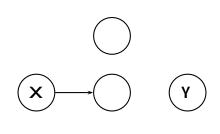


Ngarrka- ngku ka wawirri panti- rni. (Hale, 1983) man erg aux kangaroo spear nonpast

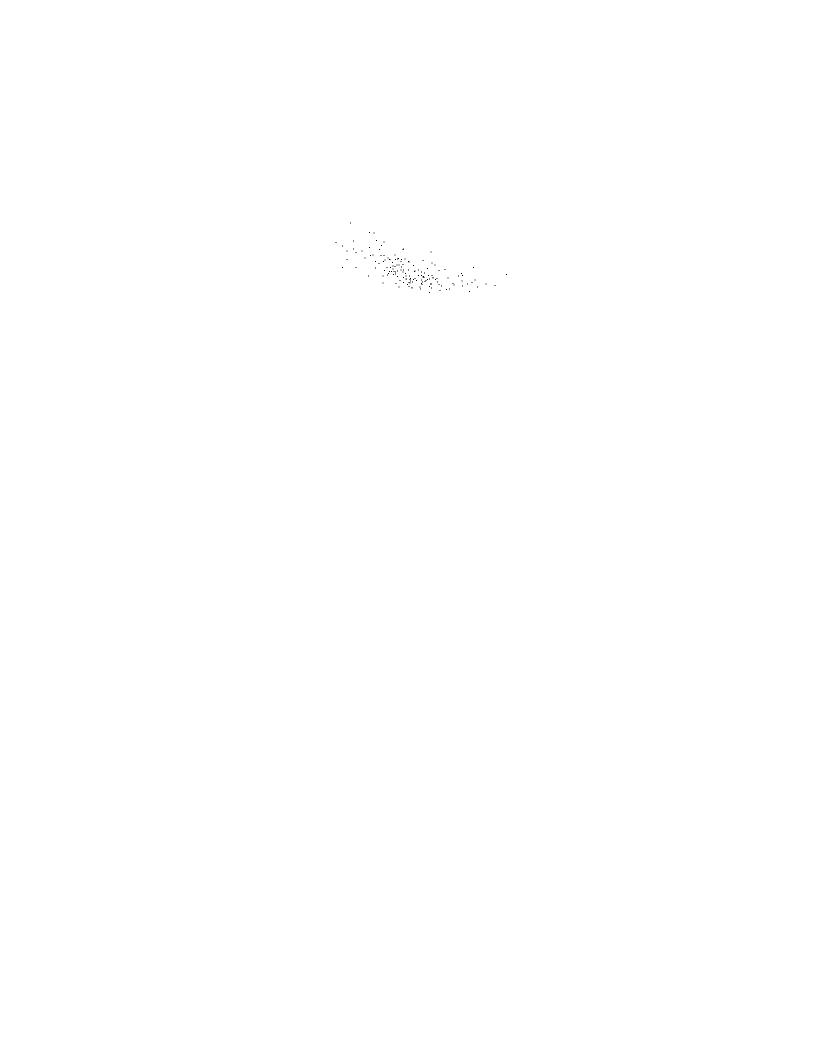
"The man is spearing the kangaroo".

In some dialects of Warlpiri, however, using the ergative case is not obligatory. Note

Chapter 6 Generalized Linear Models

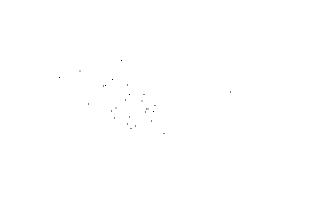


>			
	x		



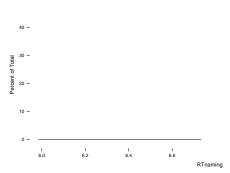


Observed data



respectively to the boxes $M_A - M_0$ and Unexplained. Thus, using the $\mathsf F$ statistic for hy-

ncmrcl70n factor could have arbitrarily di erent c



Frication	Age	X_1	X_2	X_3	X_4	X_5	X_6	X_7
burst	old	0	0	0	0	0	0	0
frication	old	1	0	0	0	0	0	0
long	old	0	1	0	0	0	0	0
short	old	0	0	1	0	0	0	0
burst	young	0	0	0	1	0	0	0
frication	young	1	0	0	1	1	0	0
long	young	0	1	0	1	0	1	0

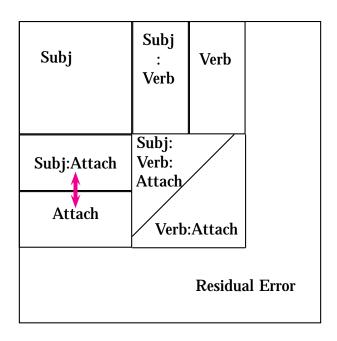
an extremely rich topic, and we take them up in Chapter 8 in full detail. There is also, however, a body of analytic techniques which uses the partit

precise reason for this. Suppose that we were to test for the presence of an interaction



```
+ result
+ }
> get. z. score <- function(response, conds. list) {
+ means <- tapply(response, conds. list, mean)</pre>
```

Subject 1 2 3 4 5 ... Verb Attachment



Error: subj:verb:attachment

So sentences with nonpronominal recipients are realized ro

is approximately distributed as a $\ ^2_{k}$ random variable, where k

6.9	Log-linear	and	multinomial	logit	models

and idiosyncratically to the probability of other words wit

ž.

 λ_1

[sr], for example, would now have the paired-segment feature ${\bf sr}$

covered in Section XXX. In this model, there is a collection of feature functions $\boldsymbol{f}_{\boldsymbol{j}}$ each of which maph

This is a new expression of the same model, but with fewer para

is McCullagh and Nelder (1989). For GLMs on categorical data, Agresti (2002) and the more introductory Agresti (2007) are highly recommended. For more information specific to

- Word frequency
- Speaker sex
- Speech rate

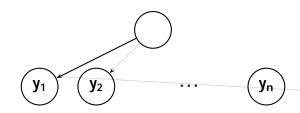
Exercise 6.5: Decomposition of variance

= 0 against an alternative-hypothesis model M_A with unconstrained $\ \ b \ A$

(Pro)noun	9192
Verb	904
Coordinator	1199
Number	237
(Pre-)Determiner	3427
Adverb	1846
Preposition or Complementizer	2418
hand	

wh-word

Chapter 8 Hierarchical Models



•		

1. Construct point estimates of the parameters of interest,	_b and	, using the principle

	lower bound	upper bound	posterior mode
μ	613.6	644.7	618.9
b	35.4	60.4	43.7
У	13.5	20.9	19.1

*		

parameter (as indicated by the

Fully Bayesian analysis

We can try a similar analysis using fully Bayesian techniques rather than the point estimate. We'll present a slightly simpler model in which the speaker-

study of language. We move from generalized linear models (G

To illustrate the approach, we construct a model with the length, animacy, discourse accessibility, pronominality, and definiteness of both the recipient and theme arguments as predictors, and with verb as a random e ect. We use log-transformed length predictors (see Section 6.7.4 for discussion).



| 0.008 - | 0.006 - | 0.004 - | 0.002 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.000 - | 0.00

F1

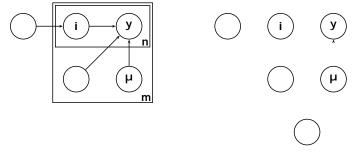
safely conclude that the e ects of these factors truly are equal and opposite? (Hint: the easiest way to construct the simpler model is to define new quanti

estimated density, and overlay

notto

Chapter 9

Dimensionality Reduction and Latent Variable Models

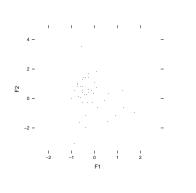


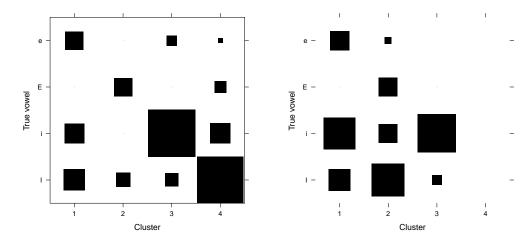
(a) Simple model

9.1	.1	Inference	for	Gaussian	mixture	models
-----	----	------------------	-----	----------	---------	--------

Cluster

1 -





(a) With learning of category frequencies

$$P(z_{ij}|w_{-ij},z-ij, ,) =$$

NP VP V Det N^{\prime} the N'PP NP N P dog Det N'near the N cat

probabilities as follows: for each rule, write a tree-search expression corresponding using a tgrep2, Tregex, or TIGERSearch

Appendix A

Mathematics notation and review

This appendix gives brief coverage of the mathematical notation and concepts that you'll encounter in this book. In the space of a few pages it is of course impossible to do justice to topics such as integration and matrix algebra. Readers interested in strengthening their fundamentals in these areas are encouraged to consult XXX [calculus] and Healy (2000).

A.1 Sets ({}, , ,)

a six, with the other five outcomes all being equally likely (i.e. 10% each). If we define a discrete random variable X representing the outcome of a roll of this die, then the clearest way of specifying the probability mass function for X is by splitting up the real numbers

it doesn't have the normalizing constant $\frac{1}{\sqrt{2}-2}.$ In order to determine the value of this

A.7 Combinatorics (

The $n \times n$ identity matrix is sometimes notated as I_n ; when the dimension is clear from context, sometimes the simpler notation I is used.

Transposition: For any matrix X of dimension $m \times n$, the transpose of X, or

Appendix C

but we can use the following conditional independencies, which can be read o the connec-



- Benor, S. B. and Levy, R. (2006). The chicken or the egg? A probabilistic analysis of English binomials. Language, 82(2):233–278.
- Beyer, K. (1986). *The Aramaic language, its distributions and subdivisions.* Vandenhoeck & Ruprecht.

Lisker, L. and Abramson, A. S. (1967). Some $e\ ects$ of context on vo

Ross, J. R. (1967). Constraints on Variables in Syntax