

## 8 The Gaussian probability model



The above photograph is of the front of the 10 (zehn) Deutsche Mark banknote, before Germany and the other EC countries adopted the Euro.

You can find other images by Googling ‘Gauss 10dm image.’

For biographies of and historical notes on Gauss, see specialized sites such as <http://www.mathunion.org/general/prizes/gauss/details/> or general ones such as [http://en.wikipedia.org/wiki/Carl\\_Gauss](http://en.wikipedia.org/wiki/Carl_Gauss). The story of the discovery of Ceres <http://www.keplersdiscovery.com/Asteroid.html> is interesting as it mentions the Method of Least Squares. There is a debate as to who first used the principle of *Least Squares*. It seems from this website that Gauss did, but it also seems he did not publish the method, so we cannot tell if he did in fact used it in the re-discovery of Ceres.

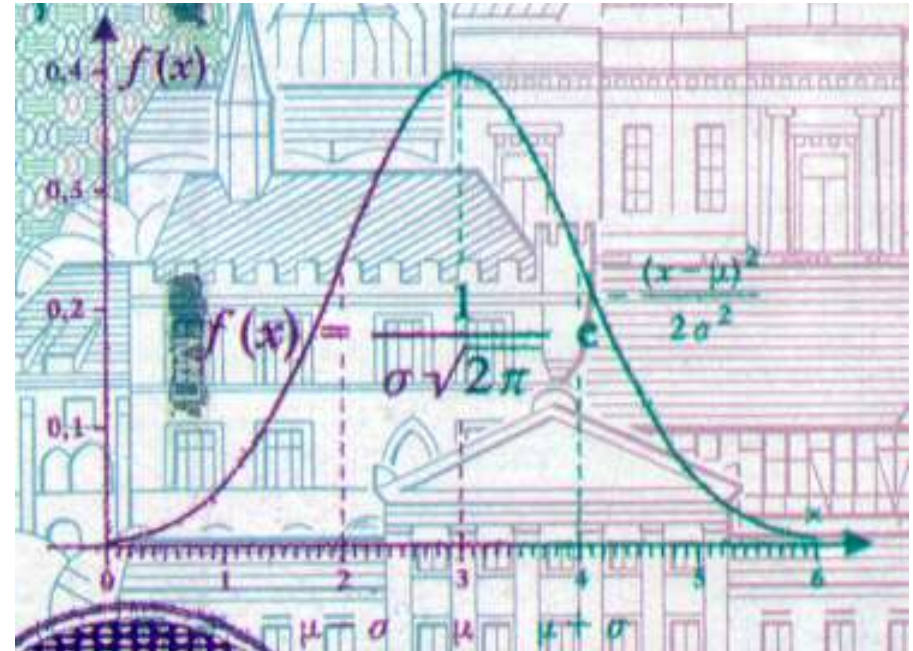
Legendre<sup>1</sup> published his method in 1806, and described the procedure very much the way it is taught today. Gauss seems to get the credit in this website [http://en.wikipedia.org/wiki/Least\\_squares](http://en.wikipedia.org/wiki/Least_squares).

There is a chapter on ‘Least Squares and the Combination of Observations’ in Stephen Stigler’s most readable and interesting book “The History of Statistics: The Measurement of Uncertainty before 1900.”

Even though we generally credit Fisher with the development of Maximum Likelihood methods, it seems that Gauss used the principle, but in reverse to

<sup>1</sup>[http://en.wikipedia.org/wiki/Adrien-Marie\\_Legendre](http://en.wikipedia.org/wiki/Adrien-Marie_Legendre)  
<http://www.nndb.com/people/891/000093612/>  
<http://www.nndb.com/people/891/000093612/>  
<http://www.britannica.com/EBchecked/topic/334063/least-squares-approximation>

derive the Normal (‘Gaussian’) Law of Errors, with density proportional to the exponentiated half of the square of the size of the error.



Note that although C&H say “in the Gaussian model, the total probability of 1 is shared among many values,” this statement applies to any model for a ‘continuous’ random variable.

‘Probability density per unit value’: this is a good description. This is how we should label the vertical axis of a pdf graph.

It is also why we can write the likelihood contribution of an ‘observed’ y value, e.g., blood pressure of 90, as  $pdf(y_{mid}, \theta) h$ , where  $y_{mid}$  is the midpoint of the interval (of width  $h$ ) that is reported as a ‘90’. Technically, we should use as  $CDF(y_{upper}, \theta) - CDF(y_{lower}, \theta)$  but usually  $h$  is sufficiently narrow that the rectangular area  $pdf(y_{mid}, \theta) \times h$  is a good approximation. Since  $h$  is constant, and does not involve  $\theta$ , it is often omitted from the likelihood contribution.

### 8.1 Standard Gaussian distribution

Because the pdf tends to be written with  $1/\sqrt{2\pi}$  in front, then, unless we evaluate this expression, we don’t get to see that it is indeed very close to 0.4.

Another point that gets overlooked is *where* (i.e., how far up the vertical axis) the ‘point of inflection’ is to be found: i.e., imagine one were to rub one’s finger upwards along the curve: where does it change from concave upwards to convex upwards?

### Supplementary Exercise 8.1.

Determine, analytically/numerically, where on the  $z$  scale the point of inflexion is located, and at what height (express its vertical location as a fraction of the max. height of  $1/\sqrt{2\pi} = 0.3989$ ). *The point of inflection helps one to draw a reasonably accurate pdf curve ‘freehand.’*

The diagram on the next page shows how one could use a ‘spinner’ or ‘roulette wheel’ to generate **random numbers that follow a Gaussian distribution**, while the one on the page thereafter shows how to do so using the reverse CDF. In fact, both use the same idea. Francis Galton<sup>2</sup> proposed a method that used 6 dice.

### Supplementary Exercise 8.2.

1. Derive the 24 ordinate values used by Galton for his *Die I*.<sup>3</sup> Do so in R (or Excel), by dividing up the  $(0, \infty)$  scale into 24 bins each containing equal 1/48ths of the total probability mass, and finding (as Galton did) the ordinate at the mid point of (the base of) each bin.
2. Out of curiosity, how many would change if he were to use the centres of gravity (mass) rather than the midpoints? Would the ‘ $\sigma$ ’ be closer to 1?
3. Then, use R (or Excel) to simulate 12 throws of *Die I*, and the requisite number of throws of *Dies II* and *III*, to produce 12 values from a

<sup>2</sup>Dice for Statistical Experiments. Nature(1890) 42 13-14 – in Resources, with related material.

<sup>3</sup>Note that whereas today we use the *Standard Deviation* (SD) as a measure of the spread of a Normal distribution [and use the ‘68-95-99.7 rule’ for 1, 2 and 3 SD’s – see [http://en.wikipedia.org/wiki/68-95-99.7\\_rule](http://en.wikipedia.org/wiki/68-95-99.7_rule)] – Galton used the smaller ‘*Probable Error*’ or ‘PE’. The PE is approximately 2/3rds (0.6745) times the SD. The *Probable Error* gets its name from the fact that a deviate from the mean is just as likely (50%, ‘as probable as not’), to be bigger than as smaller than the Probable Error: the interval  $\mu \mp 1PE$  contains the middle 50% of the  $N(\mu, PE)$  distribution, whereas the the interval  $\mu \mp 1SD$  contains the middle 68%. Notice also that by a using a smaller measure of spread, Galton’s grades or ‘degrees’ of ‘extremeness’ ran from  $-5^\circ$  to  $+5^\circ$ , where we might speak or write of Z values or Z scores from say -3 to +3.

$N(\mu = 100, \sigma = 15)$  distribution.

4. In the study carried out at the University of Canterbury [cf Resources], the investigators say that “To determine the success of this experiment, we formulate the following question as a statistical hypothesis test: Are our sampled values taken independently and identically from an appropriate discrete distribution which approximates Galton’s normal distribution?” To answer it, they used the data reported in their Appendix A. If they had approached you about how to address their question, how would you have answered them? Would you have advised them to collect results of actual throws of the dice? Would you have calculate the optimal number of trials needed for the test in the same way that they did?

## 8.2 General Gaussian Model

*Shouldn’t textbooks use  $y$  rather than  $x$  when dealing with random variables?*

This notation is important when we come to regression: most applied work involves  $y$ ’s, each of which is the realization from a conditional-on- $x$  distribution whose parameters are governed by a linear combination of a (possibly-vector-valued)  $\theta$  and a (possibly-vector-valued)  $x$ , with  $x$  treated as fixed-by-design. JH notes that the designers of the German 10DM banknote to honour Gauss also used  $x$  rather than  $y$ .

*The square of sigma is called the variance:*

In the very interesting website on Earliest Known Uses of Some of the Words of Mathematics<sup>4</sup> we find: MODULUS (in the Theory of Errors). In his first theory of least squares based on the normal distribution and presented in Gauss’s *Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientum* (1809) Gauss used a measure of precision (“mensura praecisionis observationum” (p. 245) which he denoted by  $h$ : the reciprocal of  $h$  is  $\sqrt{(2)\sigma}$ , where  $\sigma$  is the standard deviation. Both  $h$  and its reciprocal have been called the modulus: the reciprocal in G. B. Airy’s *On the Algebraical and Numerical Theory of Errors of Observation and the Combination of Observations* (1861, p. 15) and  $h$  in E. T. Whittaker & G. Robinson’s *Calculus of Observations* (1924, p. 175).

<sup>4</sup><http://jeff560.tripod.com/m.html>

In Merriman's sections 29 and 30, in his 'Discussion of the Curve  $y = ke^{-h^2x^2}$ ', "The maximum value of  $y$  is for  $x = 0$ , when  $y = k$ ;  $k$  is, hence, the probability of the error 0. [...] The constant  $h$  is a quantity of the same kind as  $1/x^2$  since the exponent  $h^2x^2$  must be an abstract number. The probability of an assigned error  $x'$  decreases as  $h$  increases; and hence, the more precise the observations, the greater is  $h$ . For this reason  $h$  may be called 'the measure of precision.' The constant  $k$  is an abstract number; and, since it is the probability of the error 0, it is larger for good observations than for poor ones. The more precise the measurements, the larger is  $k$ ."

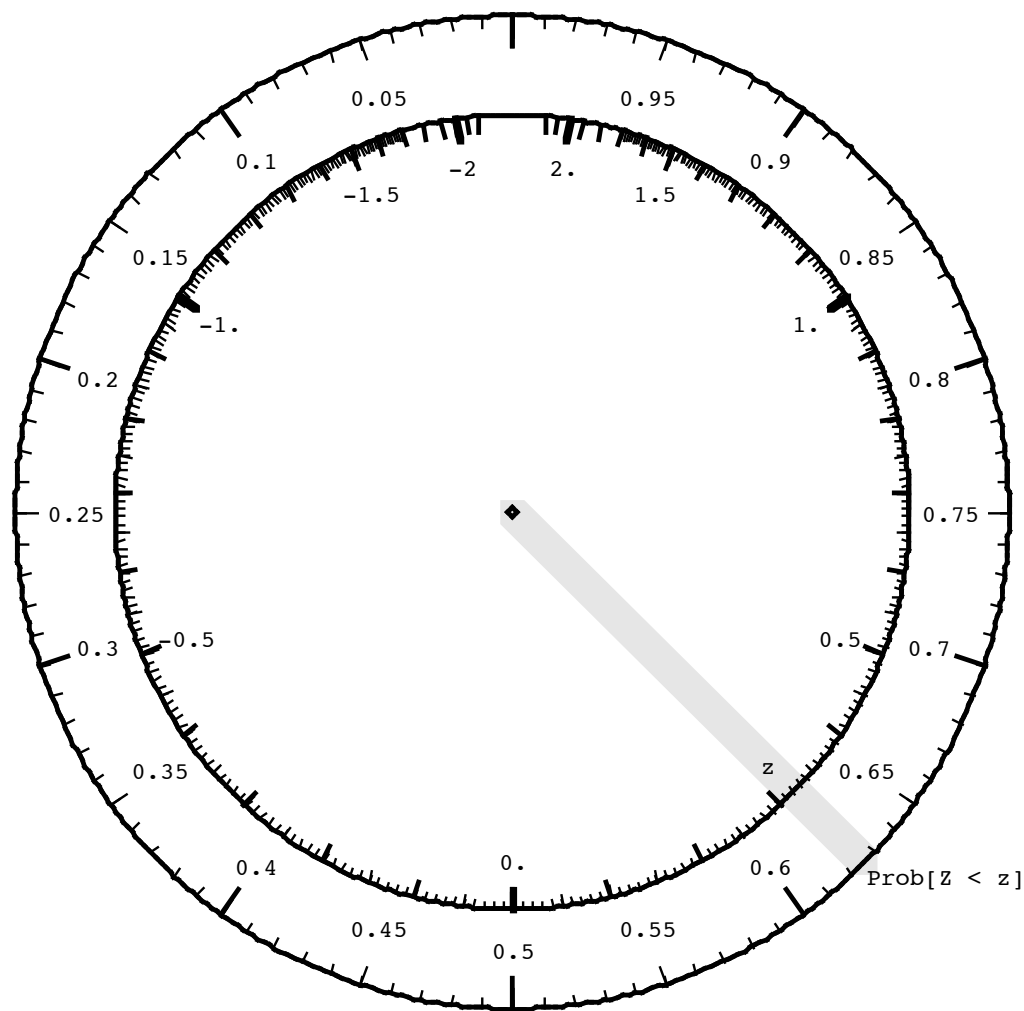
Th3 'Earliest Known Uses' website could keep one occupied for many hours, as on the same page of 'M' alone, there are entries for MARGIN OF ERROR, MARKOV CHAIN, MONTE CARLO, MARKOV CHAIN MONTE CARLO, MARTINGALE, MAXIMUM LIKELIHOOD, MEAN, MEDIAN, META-ANALYSIS, MINIMUM CHI-SQUARED, MODE, MOMENT, Moment generating function, MONTY HALL PROBLEM, MORAL EXPECTATION, MOVING AVERAGE, MULTICOLLINEARITY, MULTINOMIAL DISTRIBUTION, and MULTIVARIATE.

In the 'S' page (also full of other interesting terms) we can read that "The term STANDARD DEVIATION was introduced by Karl Pearson (1857-1936) in 1893, "although the idea was by then nearly a century old" (Abbott; Stigler, page 328). According to the DSB: The term "standard deviation" was introduced in a lecture of 31 January 1893, as a convenient substitute for the cumbersome "root mean square error" and the older expressions "error of mean square" and "mean error."

### 8.3 The Gaussian Likelihood

C&H finally explain why they adopted a cutoff on 1.353 back in chapter 3.

**Generating random numbers from a Gaussian Distribution / Connections with the Normal (Gaussian) Tables**



Imagine a disk or "Spinner" with 2 concentric circles, and a spindle through the centre. Suppose that when spun it is equally likely to come to rest at any point on the outer circumference. This is reflected in markings of 0 to 1 (or, if you prefer, % to 100%) uniformly on the circumference of the outer circle.

**Q:** How should we mark the circumference of the inner circle so that repeated spins produce values with a Gaussian  $N(0,1)$  distribution? [see "spinner" in fig 4.9 page 317 of M&M]

**A:** Use the z values corresponding to the percentiles of the Gaussian Distribution!

Then, the spinner shown will produce Z values from minus to plus infinity..

**IMPLICATIONS FOR MONTE CARLO (SIMULATION) WORK**

1 Generate numbers with a Uniform Distribution on (0,1)

e.g. in Excel use the RAND() function

i.e. generate  $P = \text{RAND}()$

2 Calculate percentile corresponding to P

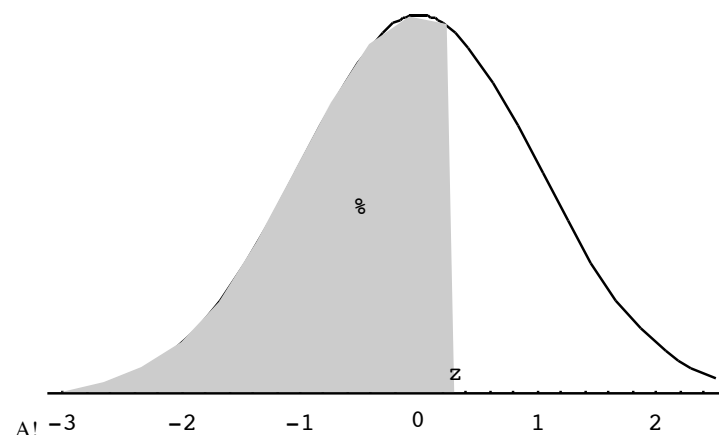
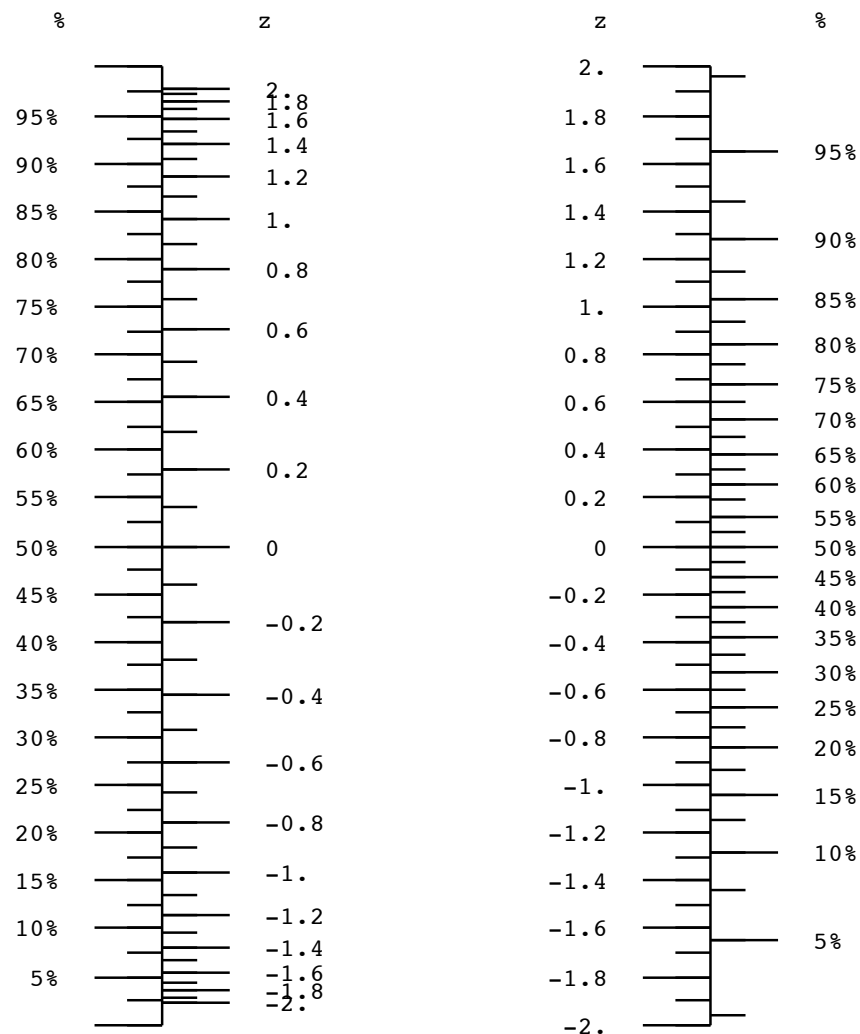
i.e.  $z = Z$  value such that  $\text{Prob}(Z < z) = P$

in Excel, use NORMINV function,

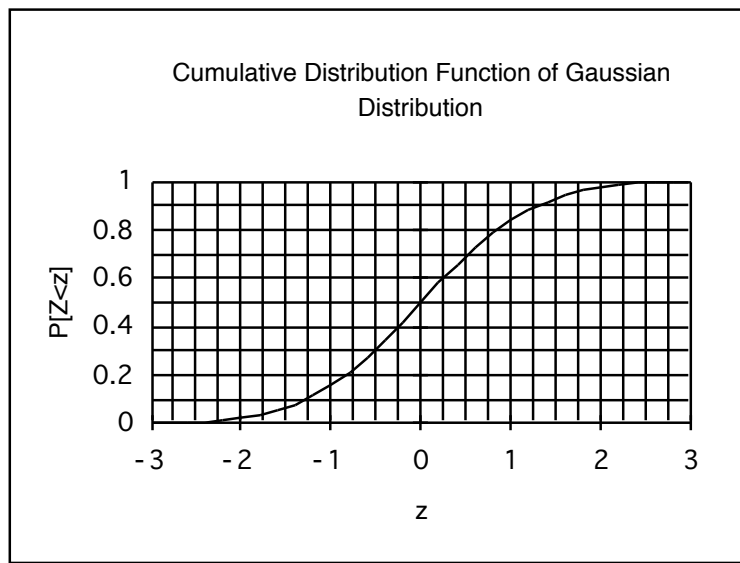
i.e.

calculate  $z = \text{NORMINV}(P, \mu=0, \sigma=1)$

### Generating random numbers from a Gaussian Distribution / Connections with the Normal (Gaussian) Tables



Another way of visualizing the Table is given below. To generate a random Z, enter randomly at the vertical axis and find corresponding Z value!



The above **nomograms** illustrate the same idea: the function links the shaded area under the Gaussian curve with the corresponding z value. It is shown, first with area or Percent or  $\Pr(Z < z)$  as a function of z, and then vice-versa (as is done in Table A of M&M). Table A tabulates  $\text{Prob}[Z < z]$  as a function of z, but one can travel in either direction.

## 8.4 The Likelihood with $N$ observations

JH prefers  $n$  **over**  $N$ . To him,  $N$  refers to the size of some universe of units, not to the size of a sample of units from it. In the social sciences, and particularly in the Publication Manual of the American Psychological Association<sup>5</sup>, its JH's understanding that if the *overall* sample is  $N = 20$  people, 9 men and 11 women, the *sub-sample* sizes of 9 and 11 would be referred to as  $n = 9$  and  $n = 11$ . I understand from a reliable source that Clayton pursued a PhD in psychology, so that might explain his notation.

Likewise, JH is unsure why C&H use the letter  $M$  where we would normally write  $\bar{x}$  [ or  $\bar{y}$  ! ]

### Supplementary Exercise 8.3.

C&H say it requires only elementary algebra to rearrange the log likelihood. Do the algebra, and verify that that is indeed true.

### Supplementary Exercise 8.4.

Chapter II of the late 19th century Least Square textbook by M Merriman, and Chapter 7 of the early 21st century one on Probability Theory by ET Jaynes, give accounts of how the Normal probability model can/does arise.

Summarize, in a paragraph each,

1. 'Hagen's demonstration' (First Deduction of the Law of Error, Merriman page 17- );
2. 'Gauss's demonstration' (Second Deduction of the Law of Error, Merriman page 22- );
3. The Herschel-Maxwell derivation (Jaynes page 200- )
4. The Gauss derivation (Jaynes page 202- )

---

<sup>5</sup><http://www.apastyle.org/manual/index.aspx>

## 9 Approximate Likelihoods

This is a central chapter; you will use the **Normal Approximation to the (sampling) distribution of ML parameter estimates throughout your career.** The key is to work in the **most appropriate parameter scale**, the one where the sampling distribution is ‘closest to Gaussian.’

Even though the title says approximate *likelihoods*, in fact the chapter is entirely about approximate *log likelihoods*. Indeed, just as R. A. Fisher did in his very first paper on this topic, exactly 100 years ago, we should always focus on the *log likelihood*. In that 1912 paper, Fisher never defined the likelihood, only its log. (In fact he did not call it the log *likelihood* ... that term came later. He just called it an absolute ‘*criterion*’).

C&H write of there being no *simple algebraic expression* (or closed form) for the supported range for the parameters of the Binomial and Poisson models. In fact, the same is true *much more broadly*: this was also the case in just about all of the ML estimation problems we have dealt with so far (e.g.  $\sigma$  in Fisher’s binned errors data, and in Tibshirani et al’s ‘accuracy of dart throws’ data;  $\mu$  and  $\sigma$  in Galton’s data on the speeds of homing pigeons; the shape and rate/scale parameters of the gamma model for tumbler longevity; the HIV infection rate parameter  $\lambda$  in the circumcision studies; the proportion ( $\pi$ ) of persons infected by West Nile virus; the parameters in the mutation rate function behind the genetic data from Iceland, etc. etc.. And this absence of a closed form will also be the norm for the parameters in many many *regression* models. Indeed, in most parameter-fitting applications, we do *not even have a closed form algebraic expression for the point estimates*, let alone their standard errors. So, it is important that we learn how to use a *more general approach*.

“The quadratic approximation becomes closer to the true log likelihood as the amount of data increases”: The authors are referring to the role of the *Central Limit Theorem*, and to the near-Gaussian sampling distribution of  $\hat{\theta}_{MLE}$ .

### 9.1 Approximating the log likelihood

“Consider a general likelihood for the parameter,  $\theta$ , of a probability model and let  $M$  be the most likely value of  $\theta$ . Since the quadratic expression

$$-\frac{1}{2} \left( \frac{M - \theta}{S} \right)^2$$

has a maximum value of zero when  $\theta = M$ , it can be used to approximate the true log likelihood ratio.”

If JH were writing this, he would have said

“**Consider a general likelihood for the parameter,  $\theta$ , of a probability model and let  $\hat{\theta}_{ML}$  be the most likely value of  $\theta$ . Since the quadratic expression**

$$-\frac{1}{2} \left( \frac{\hat{\theta}_{ML} - \theta}{S} \right)^2$$

**has a maximum value of zero when  $\theta = \hat{\theta}_{ML}$ , it can be used to approximate the true log likelihood ratio.**”

“We shall refer to  $S$  as the standard deviation of the estimate of  $\theta$ . Alternatively, it is sometimes called the standard error of the estimate”

Again, JH would have written:

“**We shall refer to  $SE[\hat{\theta}_{ML}]$  as the standard deviation of the estimate of  $\theta$ . Alternatively, it is sometimes called the standard error of the estimate**”

and (dropping the awkward *ML* subscript in the *SE*) that

“**Since the quadratic expression**

$$-\frac{1}{2} \left( \frac{\hat{\theta}_{ML} - \theta}{SE[\hat{\theta}]} \right)^2$$

**has a maximum value of zero when  $\theta = \hat{\theta}_{ML}$ , it can be used to approximate the true log likelihood ratio.**”

**Formulae for  $S$**  [for now, without justification].

THE RISK [i.e., the PROPORTION,  $\pi$  ] PARAMETER

As C&H illustrate in the next section, it is better to work in the  $\text{logit}[\pi]$  scale, unless the Normal approximation to the binomial itself is adequately accurate. Had the data been 40+ and 60–, rather than 4+ and 6–, the log-likelihood in the original, untransformed, (i.e.,  $\pi$ ) scale would have looked a lot closer to quadratic, i.e., the sampling distribution of a sample proportion would be close-enough-to-Normal for much more of the  $\pi$  range.

A common problem with using the Normal approximation in the case of limited-range parameters is that it works well enough at the less extreme limit (in this case, for values near the middle of the range) but poorly for parameter values near the more extreme limit (in this case, for values near the bottom of the 0-1 range). JH likes to say that there isn't enough room for a Gaussian distribution if we are at the lower end (0/10, 1/10, 2/10, ...) of the Binomial(10,  $\pi$ ) distribution, particularly if we entertain  $\pi$  values < 0.5 (which case we would 'expect' to have  $n \times \pi < 5$  in the sample with the characteristic/event of interest).

THE RATE [i.e., the INTENSITY,  $\lambda$  ] PARAMETER<sup>6</sup>

*“The value of  $S$  (i.e., the value of  $SE[\hat{\lambda}]$ ) which gives the best approximation to the log likelihood ratio is  $S = \sqrt{D}/PT$ ”*

The reason for this form is because  $\hat{\lambda}_{ML}$  has the form  $\hat{\lambda}_{ML} = D/PT$ , where  $D$  is the realization of a Poisson r.v., and  $PT$  is a known constant.

The lower limit of 5.3/1000 is the result of using the Normal approximation to a Poisson distribution. But a rate of 5.3/1000 means that with  $PT=500$  person-time units, we would expect  $\mu = (5.3/1000) \times 500 = 2.65$  'events', and we know that the Poisson[ $\mu = 2.65$ ] distribution is quite skewed, with a lot of probability mass on its lower tail values of 0, 1 and 2, so it is not possible to approximate it by a Normal distribution – the lower tail of a  $N[\mu = 2.65, \sigma = \sqrt{2.65}]$  distribution has an embarrassingly large amount its mass below 0!

Imagine what the lower limit would be if the observed count was  $D = 2$ : then the lower limit for the rate, based on the Normal approximation, would be  $(2 - 1.645 \times \sqrt{2})/500$ , a negative rate! Using the form  $(D \mp 1.645 \times \sqrt{D})/500$  (i.e. dealing first with the statistical uncertainty in the numerator, then dividing

<sup>6</sup>JH has changed the possibly confusing Y (for the Years of observation denominator) to PT (for amount of Population-Time in which the events occurred). He left the 'morbid' term D (for Deaths) as is, even though he prefers to use C for 'number of Cases (instances) of', a *neutral* term that covers both 'bad' and 'good' types of events/characteristics.

by the 'constant', 500) makes it much clearer where the 'weak link' is –it's the *numerator*!

Again, a parameter-transformation would help (although, with a  $\mu$  this low, it is difficult to come up with any parameter scale on which the Normal distribution would be a good approximation – there just isn't enough 'granularity').

## 9.2 Transforming the parameter

As was emphasized at the beginning, this is the key to more accurate approximations.

*“The solution to this problem is to find some function (or transformation) of the parameter which is unrestricted and to first find an approximate supported range for the transformed parameter.”*

Indeed! And it is often possible to use the range of the parameter to determine which transformation is likely to lead to a scale on which the log-likelihood is more Gaussian. For a parameter  $\theta$ , such as a *proportion*, that takes values in the (0,1) range, the 'canonical' transformation is the *logit*, i.e.,  $\log \left[ \frac{\theta}{1-\theta} \right]$ , which takes on values all the way from  $-\infty$  to  $+\infty$ .

For a parameter  $\theta$ , such as a *rate*, that takes values in the  $(0, \infty)$  range, the canonical transformation is the *log*, i.e.,  $\log[\theta]$ . Indeed, for many of the examples to be considered from now on, JH may well use  $\theta$  for the parameter measured on the full  $(-\infty, +\infty)$  scale.

THE LOG RATE PARAMETER [with some editing by JH]

“The rate parameter  $\lambda$  can take only positive values, but its logarithm is unrestricted. To calculate an approximate supported range for  $\lambda$ , it is better, therefore, to first calculate a range for  $\log[\lambda]$ , and then to convert this back to a range for  $\lambda$ . (...) To find the approximate range for  $\log[\lambda]$ , we need a new value of  $S$  that which gives the best Gaussian approximation to the log likelihood ratio curve when plotted against  $\log[\lambda]$ . When a rate  $\lambda$  is estimated from  $D$  failures over  $PY$  person-years, i.e., as

$$\hat{\lambda} = D/PY,$$

this value of  $S = \frac{SE[\hat{\lambda}]}{\{Var[\hat{\lambda}]\}^{1/2}}$  is given by

$$S = \frac{SE}{SE[\hat{\lambda}]} = \sqrt{\frac{1}{D}}.$$

(By ignoring the possibility of a Poisson count of zero, and using a Taylor series approximation often referred to as “the *Delta Method*” ), at the beginning



of the course, when dealing with Poisson variation, we derived the **square root of the variance of the log of a Poisson r.v.**,  $Y$ , finding it to be (approximately)

$$\sqrt{\frac{1}{\mu_Y}},$$

where  $\mu_Y$  is the expected value,  $E[Y]$ . And by using the observed value ( $D$  in C&H notation,  $y$  or  $c$  in JH's) as an estimate of its expectation, and plugging it into the square root of the expression for the variance, we have a good example of a standard error – the *estimated* standard deviation of a *statistic* or parameter-estimate.

**“A more convenient way of carrying out this calculation – using a (multiplicative) ‘error factor.’ ”**

JH encourages this Multiplicative Margin of Error (MME) i.e., using  $\hat{\theta} \times \div MME$ , rather than  $\exp[\log[\hat{\theta} \pm Z_{\alpha/2}SE[\log[\hat{\theta}]]]$  approach; think of it as expressing the uncertainty as X-fold (or, en français, ‘X-fois’): the upper limit is X *times larger* than the point estimate; and the lower limit is X *times smaller* than the point estimate.

THE LOG ODDS PARAMETER [editing and notation-changes by JH]

Here we use the logit, or the log-odds, i.e.,  $\theta = \log[\Omega] = \log\left[\frac{\pi}{1-\pi}\right]$ , so that  $\hat{\theta} = \log\left[\frac{\hat{\pi}}{1-\hat{\pi}}\right]$ .

“When an odds,  $\Omega$ , is estimated from  $y+$  ‘positives’ and  $y-$  ‘negatives’, i.e., as

$$\hat{\Omega} = \frac{y+}{y-}, \quad \text{or } \hat{\theta} = \widehat{\log[\Omega]} = \log\left[\frac{y+}{y-}\right],$$

the value of  $S = SE[\hat{\theta}] = \{Var[\hat{\theta}]\}^{1/2}$  is given by

$$S = SE = SE[\hat{\theta}] = \sqrt{\frac{1}{y+} + \frac{1}{y-}}.$$

(By ignoring the possibility of a Binomial count of 0 or  $n$ , and using a Taylor series approximation often referred to as “the Delta Method”), at the beginning of the course, when dealing with Binomial variation, we derived the **square root of the variance of the log of the ratio of Binomial r.v.**,  $Y+$ , to its complement  $Y-$ , finding it to be (approximately)

$$\sqrt{\frac{1}{\mu_{Y+}} + \frac{1}{\mu_{Y-}}},$$

where  $\mu_Y$  is the expected value. And by using the observed value ( $D$  in C&H notation,  $y+$  in JH's) as an estimate of the one expectation, and the observed

value of its complement ( $N - D$  in C&H notation,  $y-$  in JH's) and plugging it into the square root of the expression for the variance, we have a good example of a standard error – the *estimated* standard deviation of a *statistic* or parameter-estimate.

Epidemiologists are (overly) fond of  $2 \times 2$  tables, with  $E$  (‘Exposed’) and  $\bar{E}$  (not) and with  $D$  (‘Diseased’) and  $\bar{D}$  (not) as the rows / columns and with  $a, b, c, d$  as frequencies

	$E$	$\bar{E}$	
$D$	$a$	$b$	$n_D$
$\bar{D}$	$c$	$d$	$n_{\bar{D}}$
	$n_E$	$n_{\bar{E}}$	$n$

JH prefers this more neutral and more general notation (with X on x-axis, and going from 0 to 1):

	$X = 0$	$X = 1$
$Y = 1$	$y_+$	$y_+$
$Y = 0$	$y_-$	$y_-$
	$n_{X=0}$	$n_{X=1}$

and looking ahead to regression, with traditional X and Y axes, where the values need not be restricted to 0's and 1's, he would argue for this layout:

	$1$	$y_+$	$y_+$
$Y$	$0$	$y_-$	$y_-$
		$0$	$1$
		$X$	

In any event, at this stage, where all subjects have the same X (say  $X=0$ ), our only interest is in the left ( $X=0$ ) column of the table, and in the  $y_+$  to  $y_-$  ratio. So, for the logit, when we substitute observed for expected values, we have the biostatistician's expression

$$SE[\text{logit}[\hat{\pi}]] = SE[\widehat{\log[\Omega]}] = SE[\hat{\theta}] = \sqrt{\frac{1}{y_+} + \frac{1}{y_-}}.$$

or the epidemiologist's expression, with  $a$  and  $c$  as the focus for now (or  $a$  and  $b$  if they happened to – or were taught to – exchange the row and columns labels) :

$$SE[\text{logit}] = SE[\log[a/b]] = \sqrt{1/a + 1/b}.$$

**Supplementary Exercise 9.1.**

Suppose we measured the sex ratio in a sample of  $n = 20$ , obtaining the ratio  $\sigma : \varphi = 14:6$ . Calculate a 90% CI for the true  $\sigma : \varphi$  ratio,  $\Omega$ , by first calculating a CI for  $\log[\Omega]$ .

$$I[\lambda]_{\hat{\lambda}} = \frac{PT^2}{y}$$

and using  $\{I[\lambda]\}^{-1} = \frac{y}{PT^2}$  as  $SE[\hat{\lambda}]^2$  yields

$$SE[\hat{\lambda}] = \frac{\sqrt{y}}{PT},$$

just as we had established (directly!) at the beginning of the course, treating  $\hat{\lambda} = \frac{y}{PT} \sim \frac{Poisson(\mu)}{PT}$ .

The *point of this long exercise* is that we can check in simple cases that the two approaches lead to the same  $SE[\hat{\lambda}]$ , but that there are many instances where there is no closed form, and where the ‘curvature’ approach is the only way to calculate a SE.

THE RISK ( $\pi$ ) PARAMETER

C&H go through the same curvature exercise for this parameter, and arrive at the familiar ‘binomial’ SE of  $\sqrt{\frac{\hat{\pi} \times (1-\hat{\pi})}{n}}$ .

**9.4 Approximate likelihoods for transformed parameters**

First, JH applauds C&H for ‘transforming’ the *parameter* before transforming the random variable. See the American Statistician article “The PDF of a Function of a Random Variable: Teaching its Structure by Transforming Formalism into Intuition” by JH and D Teltch – it is under Reprints/Talks on JH’s website. Its more about a change of *scale* than the ‘change of variable’ that is usually taught. After all, the entity called ‘Montreal temperatures’ is the same (and temperatures in January are equally cold), whether we choose to measure them in °F or °C !

**Supplementary Exercise 9.2.**

C&H repeat the above SE calculations, but for the parameter  $\beta = \log[\lambda]$  rather than  $\lambda$ . Fill in the steps they don’t show. [And note their wise choice of the letter  $\beta$  – they are thinking ahead to linear predictors in multiple regression. In this simple 1-sample example, think of it as a  $\beta_0$  ! ]

“In general, derivations such as that above can be simplified considerably by using some further elementary calculus which provides a general rule for the relationship between the values of  $S$  (the SE) on the two scales. In the case of the log transformation, this rule states that multiplying the value of  $S$  on the scale of  $\lambda$  by the gradient of  $\log(\lambda)$  at  $\lambda = M$  gives the value of  $S$  on the

**9.3 Finding the best quadratic approximation**

“To do this we need some elementary ideas of calculus summarized in Appendix B. In particular, we need to be able to find the gradient (or slope) of the log-likelihood curve together with its curvature, which is defined as the rate of change of the gradient. The mathematical terms for these quantities are the first and second derivatives of the log likelihood function.”

JH would add that “the *statistical* terms for these quantities are the **Score** and (with the sign changed to have a positive value) the **Information**.”

“The value of  $\hat{\theta}$  can be found by a direct search for that value of  $\theta$  which maximizes the log likelihood, but it is often easier to find the value of  $\theta$  for which the gradient of the log likelihood is zero; this occurs when  $\theta = \hat{\theta}$ . The value of  $S$  is chosen to make the curvature of the quadratic approximation equal to that of the true log likelihood curve at  $M$ , thus ensuring that the true and approximate log likelihoods are very close to each other near  $\theta = \hat{\theta}$ .

**We therefore choose the value of  $S$  to make  $-1/(S)^2$  equal to the curvature of the true log likelihood curve at its peak.**

THE RATE [i.e.,  $\lambda$ ] PARAMETER (estimated as ‘event count’/ $PT = y/PT$ )

$$Lik[\lambda] = \exp[-\mu] \mu^y, \quad \text{where } \mu = \lambda \times PT$$

$$LogLik[\lambda] = -\mu + y \log[\mu] = -\lambda \times PT + y \log[\lambda] - y \log[PT]$$

$$Score[\lambda] = \frac{d \text{LogLik}[\lambda]}{d\lambda} = \frac{y}{\lambda} - PT$$

$$I[\lambda] = -E \left[ \frac{d^2 \text{LogLik}[\lambda]}{d\lambda^2} \right] = E \left[ \frac{y}{\lambda^2} \right] = \frac{\lambda \times PT}{\lambda^2} = \frac{PT}{\lambda}$$

Substituting  $\hat{\lambda} = y/PT$  for  $\lambda$  yields

scale of  $\log(\lambda)$ . The rules of calculus tell us that, at  $\lambda = M$ , the gradient of the graph of  $\log(\lambda)$  against  $\lambda$  is  $1/M$ . (... )”

### Supplementary Exercise 9.3

Show that the ‘simplification’ that C&H describe is none other than the Delta Method, with its use of Jacobians to go from one scale to another. The answer you obtained by applying the ‘Delta Method’ to the transformed r.v.  $\log[y/PT]$  is a check on their algebra.

“This agrees with the expression obtained by the longer method.”

By the *longer method*, they mean the use of the *curvature/information* (and its reciprocal) at  $\theta = \hat{\theta}$ . And the ‘*this*’ refers to the *Delta Method*. To JH, the difference is that the *curvature/information* approach emphasizes the *parameter* scale, while the *Delta Method* emphasizes the *random variable* scale. The ‘change-the-parameter-scale’ is more general, and avoids having to worry about realizations of random variables (e.g., a count of zero) that do not map nicely to another (e.g., the  $\log[\text{count}]$  or  $\log[y]$  scale).

## 9.5 Déjà vu: supp. exercises in Ch. 3 (Likelihood)

There was no *simple algebraic expression* (or closed form) for the supported range for the parameters of the models involving  $\sigma$  in Fisher’s binned errors data, and in Tibshirani et al’s ‘accuracy of dart throws’ data;  $\mu$  and  $\sigma$  in Galton’s data on the speeds of homing pigeons; the shape and rate/scale parameters of the gamma model for tumbler longevity; the HIV infection rate parameter  $\lambda$  in the circumcision studies; the proportion ( $\pi$ ) of persons infected by (the sero-prevalence of) West Nile virus; the parameters in the mutation rate function behind the genetic data from Iceland, etc. etc.. So, ...

### Supplementary Exercise 9.4

Revisit each these examples, and decide which, if any, parameter transformation might lead to a ‘closer-to-Gaussian’ i.e., ‘closer-to-quadratic’ log-likelihood. Then re-run the graphs on these new scales<sup>7</sup>, and see if your intuition was borne out. Can you come up with any ‘rule-of-thumb’ to decide whether the extra work involved will be worth it?

<sup>7</sup>Once you have set up the LogLik function on one scale, it is usually easy to plot it on another. Or you can change the argument and insert the transformation at the beginning of the old function.

## 9.6 (pas?) déjà vu: corresponding chapter from ‘Likelihood’ textbook by Edwards

### Supplementary Exercise 9.5

In example 5-4.1 (page 80- ) Edwards<sup>8</sup> shows where “we may sometimes take advantage of a transformation to improve the approximation of the support (log likelihood) curve by a parabola”. His example involves the parameter (C&H, and JH would call it  $\mu$ ) of a Poisson distribution.

Repeat the derivation, but with the parameter  $\pi$  of the Binomial distribution.

<sup>8</sup>The chapter from Edwards is available in the Resources for Gaussian Model webpage, shared with that for approx. livelihoods.