

Probability

Meaning

Long Run Proportion
Estimate of (Un)certainity
Amount prepared to bet

Use

Describe likely behaviour of data
Communicate (un)certainity
Measure how far data are from
some hypothesized model

How Arrived At

Subjectively

Intuition, Informal calculation, consensus

Empirically

Experience (actuarial, ...)

Pure Thought

Elementary Statistical Principles

If necessary, breaking Complex
outcomes into simpler ones

Advanced Statistical Theory

calculus e.g. Gauss' Law of Errors

References

WMS5, Chapter 2 • Moore & McCabe Chapter 4

• Colton, Ch 3

• Freedman et al. Chapters 13,14,15

• Armitage and Berry, Ch 2

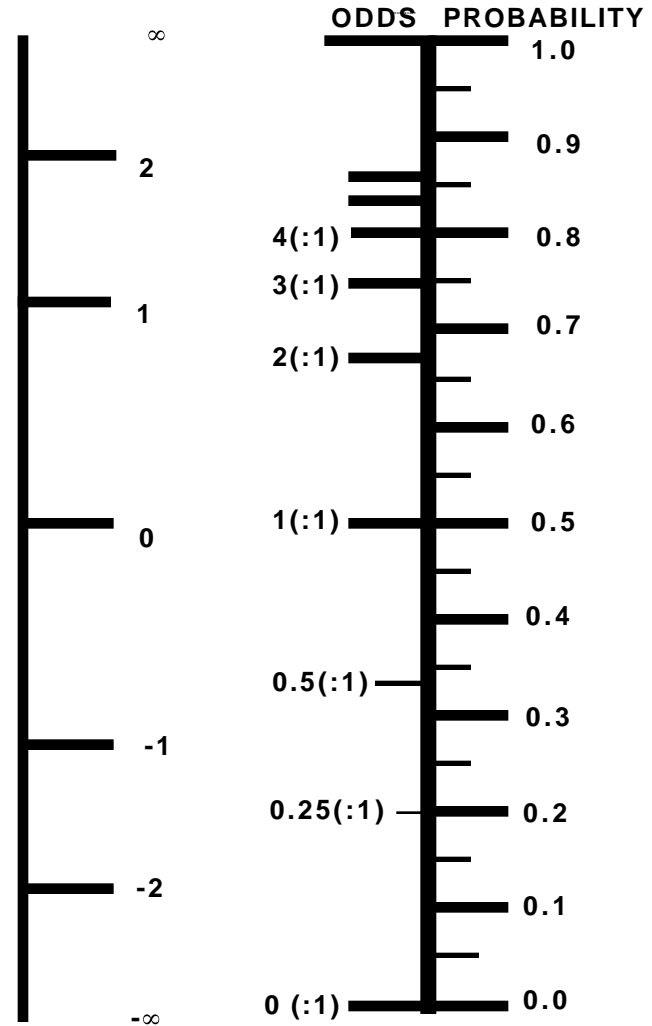
• Kong A, Barnett O, Mosteller F, and Youtz C. "How Medical Professionals Evaluate Expressions of Probability" NEJM 315: 740-744, 1986 ... *on reserve*

• Death and Taxes • Rain tomorrow • Cancer in your lifetime • Win lottery in single try • Win lottery twice • Get back 11/20 pilot questionnaires • Treat 14 patients get 0 successes • Duplicate Birthdays • Canada will use \$US before the year 2010
• OJ murdered his wife • DNA matched • OJ murdered wife | DNA matched

Probability Scales

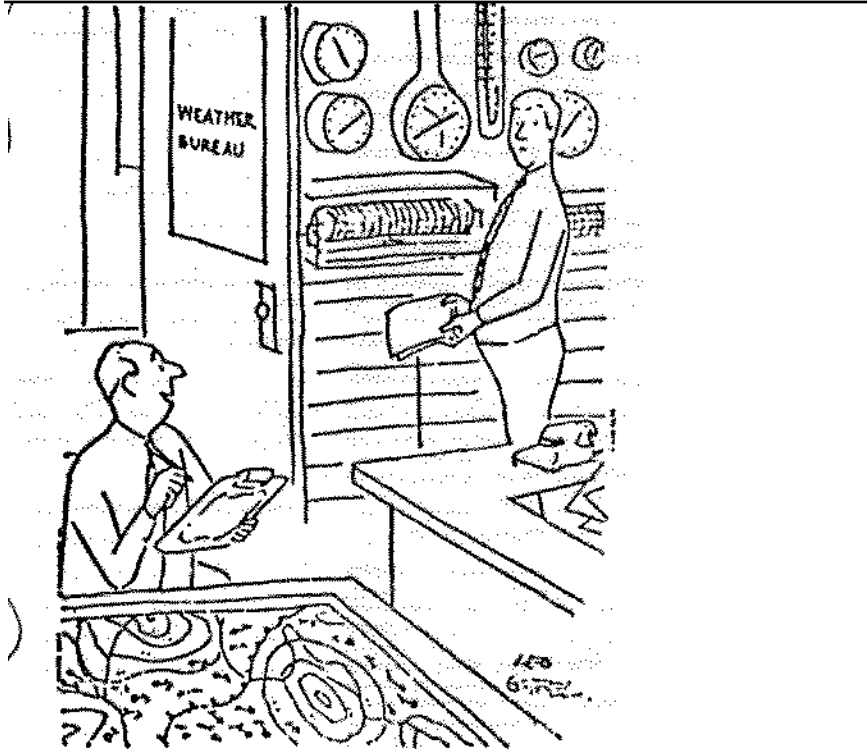
LOG-ODDS
(logit)

ODDS = PROBABILITY / (1 - PROBABILITY)
PROBABILITY = ODDS / (ODDS + 1)



• 50 year old has colon ca • 50 year old with +ve haemoccult test has colon ca • child is Group A Strep B positive • 8 year old with fever and v. inflamed nodes is Gp A Strep B positive • There is life on Mars

How to calculate probabilities

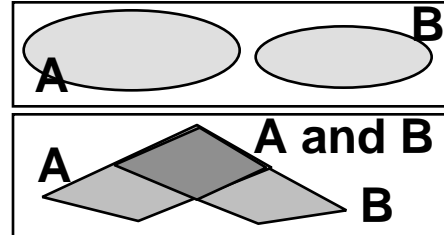


Wall Street Journal

"I figure there's a 40% chance of showers, and a 10% chance we know what we're talking about"

Probability Calculations

Basic Rules



Probabilities add to 1

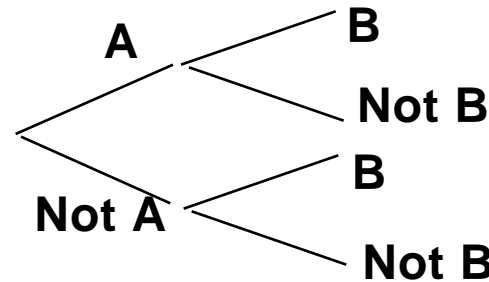
Prob(event) =
1 - Prob(complement)

ADDITION FOR "EITHER A OR B"

If mutually exclusive
 $P(A \text{ or } B) = P(A) + P(B)$

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

"PARALLEL"



MULTIPLICATION FOR "A AND B" OR "A THEN B"

If independent
 $P(A \text{ and } B) = P(A) \cdot P(B)$

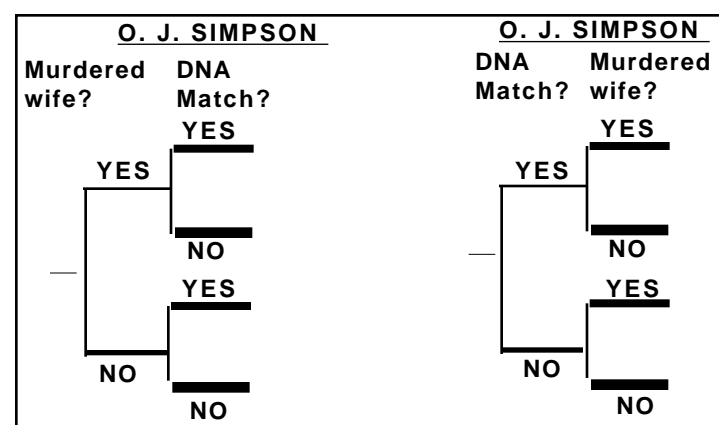
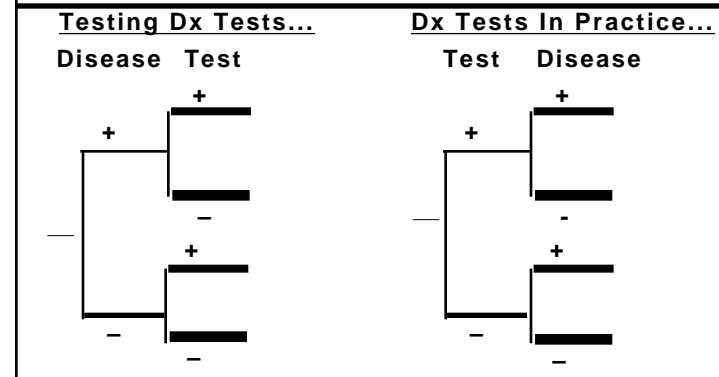
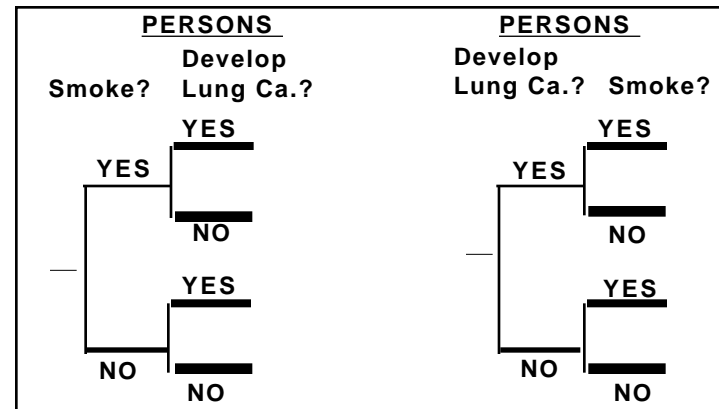
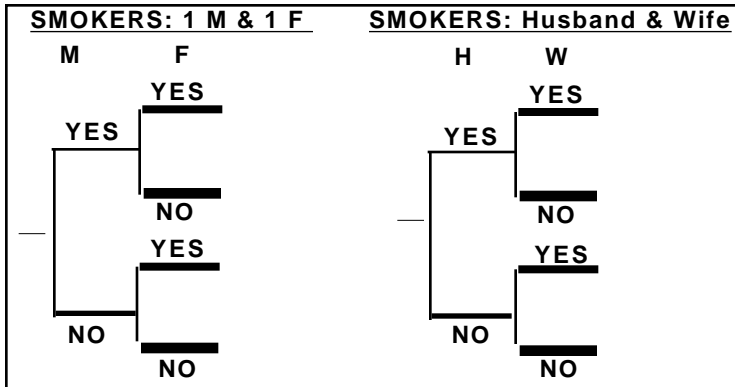
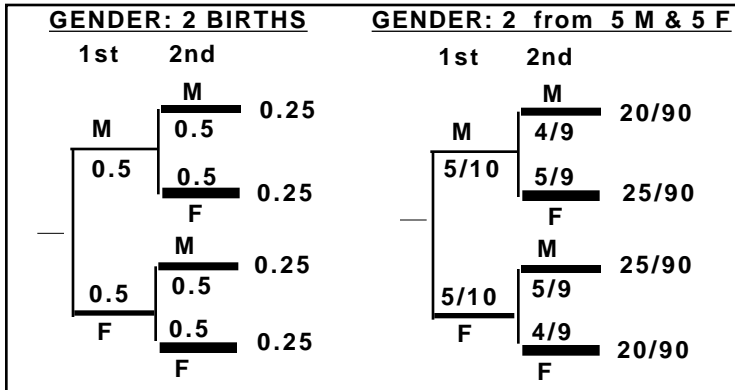
If dependent
 $P(A \text{ and } B) = P(A) \cdot P(B | A)$

"SERIAL"

Conditional Probability $P(B | A)$ = Probability of B "given A" or "conditional on A"

More Complex: • Break up into elements • Look for already worked-out calculations; • Beware of intuition, especially with "after the fact" calculations for non-standard situations

Examples of Conditional Probabilities...



US National Academy of Sciences under fire over plans for new study of DNA statistics. Confusion leads to retrial in UK.

NATURE p 101-102 Jan 13, 1994]

... He also argued that one of the prosecution's expert witnesses, as well as the judge, had **confused two different sorts of probability.**

One is the probability that DNA from an individual selected at random from the population would match that of the semen taken from the rape victim, a calculation generally based solely on the frequency of different alleles in the population.

The other is the separate probability that a match between a suspect's DNA and that taken from the scene of a crime could have arisen simply by chance -- in other words that the suspect is innocent despite the apparent match. This probability depends on the other factors that led to the suspect being identified as such in the first place.

During the trial, a forensic scientist gave the first probability in reply to a question about the second. Mansfield convinced the appeals court that the error was repeated by the judge in his summing up, and that this slip -- widely recognized as a danger in any trial requiring the explanation of statistical arguments to a lay jury -- justified a retrial.

In their judgement, the three appeal judges, headed by the Lord Chief Justice, Lord Farquharson, explicitly stated that their decision "should not be taken to indicate that DNA profiling is an unsafe source of evidence".

Nevertheless, with DNA techniques being increasingly used in court cases, some forensic scientists are worried that flaws in the presentation of their statistical significance could, as in the Deen case, undermine what might otherwise be a convincing demonstration of a suspect's guilt.

Some now argue, for example, that quantified statistical probabilities should be replaced, wherever possible, by a more descriptive presentation of the conclusions of their analysis. "The whole issue of statistics and DNA profiling has got rather out of hand," says one.

Others, however, say that the Deen case has been important in revealing the dangers inherent in the 'prosecutor's fallacy'. They argue that this suggests the need for more sophisticated calculation and careful presentation of statistical probabilities.

"The way that the prosecution's case has been presented in trials involving DNA-based identification has often been very unsatisfactory," says David Balding, lecturer in probability and statistics at Queen Mary and Westfield College in London. "Warnings about the prosecutor's fallacy should be made much more explicit." After this decision, people are going to have to be more careful."

"The prosecutor's fallacy"

Who's the DNA fingerprinting pointing at?

New Scientist, 29 Jan. 1994, 51-52. David Pringle

Pringle describes the successful appeal of a rape case where the primary evidence was DNA fingerprinting. In this case the statistician Peter Donnelly opened a new area of debate. He remarked that

forensic evidence answers the question

"What is the probability that the defendant's DNA profile matches that of the crime sample, assuming that the defendant is innocent?"

while the jury must try to answer the question

"What is the probability that the defendant is innocent, assuming that the DNA profiles of the defendant and the crime sample match?"

Apparently, Donnelly suggested to the Lord Chief Justice and his fellow judges that they imagine themselves playing a game of poker with the Archbishop of Canterbury. If the Archbishop were to deal himself a royal flush on the first hand, one might suspect him of cheating. Assuming that he is an honest card player (and shuffled eleven times) the chance of this happening is about 1 in 70,000.

But if the judges were asked whether the Archbishop were honest, given that he had just dealt a royal flush, they would be likely to place the chance a bit higher than 1 in 70,000.

The error in mixing up these two probabilities is called the "the prosecutor's fallacy", and it is suggested that newspapers regularly make this error.

Apparently, Donnelly's testimony convinced the three judges that the case before them involved an example of this and they ordered a retrial

from Vol 3.02 of Chance News

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

Random Variables & Probability Distributions

What they are:

Random Variable	Possible Outcomes (abbreviated)	Corresponding Probabilities
E.g. <u>the blood group</u> of n = 1 randomly selected person	A B AB O	P(A) P(B) P(AB) <u>P(O)</u> 1.00
<u>How many</u> of n = 20 randomly selected persons will return questionnaire in pilot study	0 1 2 ... 20	P(0) P(1) P(2) ... <u>P(20)</u> 1.00
<u>Mean cholesterol level</u> in n=30 randomly selected persons	<100 100-101 ... 249-250 >250	P(<100) P(100-101) ... P(249-250) <u>P(>250)</u> 1.00
<u>the value of the test-statistic</u> if 2 populations sampled from had the same mean	< -2.0 -2 to -1 -1 to 0 0 to 1 1 to 2 > 2.0	.028 .136 .341 .341 .136 <u>.028</u> 1.000

- we use probabilities or fractions as relative frequencies (like a histogram with an infinite number of entries)
- typically, the random quantity is obtained from an aggregate of elements e.g. a sum or a mean

References •Colton, Ch 3 •M & MCh 4.2 and 4.3

Expectation (Mean) & Variance of Random Variable

• If X takes on the DISCRETE values

x_0	with probability	p_0
x_1	with probability	p_1
...
x_k	with probability	p_k

then the expected value of X (written "E(X)") is

$$x_0 \cdot p_0 + x_1 \cdot p_1 + x_2 \cdot p_2 + \dots + x_k \cdot p_k \text{ or } \sum_{i=1}^{i=K} x_i \cdot p_i$$

Compare the formula for E(X) with that for xbar:-

- E(X) is a mean that uses expected (i.e. unobservable or theoretical or long run) relative frequencies (p's)
- xbar uses observed relative frequencies (f / n)'s.

• If X takes on the CONTINUOUS values

$$x - \frac{\Delta x}{2} \text{ to } x + \frac{\Delta x}{2} \text{ with probability } p = f(x) \cdot \Delta x,$$

$$\text{then } E(X) = \sum_{x_{\min}}^{x_{\max}} x \cdot f(x) \cdot \Delta x$$

Variance of a Random Variable

$$\text{Var}(X) = \sigma^2 = E[(x - \mu)^2] = \sum_{i=1}^{i=K} [x_i - \mu]^2 \cdot p_i$$

i.e. the Expected Squared Deviation from μ

Just as there was a computational shortcut for calculating σ^2 , we can write

$$\text{Var}(X) = \sigma^2 = E [x^2] - \mu^2$$

"ave(square) - squared ave"

References

•Colton, Ch 3 •Moore & McCabe Chapter 4.3

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

Relevance of Expectation of a Random Variable

- 1 **ACTS AS A MEAN FOR A VARIABLE THAT HAS A (CONCEPTUAL) REPETITION OR AN INFINITE N**

- 2 **THE EXPECTED VALUE OF A RANDOM VARIABLE X WILL USUALLY BE IN TERMS OF POPULATION PARAMETERS**

A STATISTIC WITH EXPECTED VALUE θ IS AN "UNBIASED ESTIMATOR" OF θ .

e.g.1 $X =$ Proportion of YES' in sample

$$E(X) = \text{PROPORTION of YES' in POP} \\ \text{THEN } \hat{X} = X$$

(X is an unbiased estimator of θ)

e.g.2 Likewise, if we use divisor of $n - 1$,

$$E(s^2) = \sigma^2, \text{ so...}$$

$\hat{s}^2 = s^2$ is an unbiased estimator of σ^2

{ \hat{s}^2 stands for "estimate of " σ^2 }

If we use divisor of n

$$E(s^2 \text{ with divisor of } n) = \frac{n-1}{n} \sigma^2 \text{ (too small on average)}$$

e.g. Expected years of life at birth (using Québec 1990 mortality data)

Length of life = age at death;

Assume for sake of illustration that deaths in decade are all at midpoint of interval (not quite true; calculations done one year rather than one decade at a time would be more exact)

decade	mid-point age	<u>Males:</u> proportion (p) dying in this decade	age • p	<u>Females:</u> proportion (p) dying in this decade	age • p
0-10	5	0.010	0.050	0.008	0.040
10-20	15	0.006	0.089	0.002	0.030
20-30	25	0.012	0.295	0.004	0.099
30-40	35	0.016	0.544	0.007	0.242
40-50	45	0.030	1.335	0.017	0.749
50-60	55	0.074	4.079	0.040	2.223
60-70	65	0.180	11.697	0.096	6.233
70-80	75	0.301	22.610	0.214	16.049
80-90	85	0.279	23.680	0.358	30.442
90-100	95	0.093	8.822	0.254	24.136
All (Σ)		1.000	73.2	1.000	80.2

Expectation of Life at Birth (average longevity)

Males: 73.2 years Females: 80.2 years

Variance[longevity] = average[square] – squared average
 Males: Ave[square] = $5^2 \cdot 0.010 + 15^2 \cdot 0.006 + \dots + 95^2 \cdot 0.093$
 = 5619.38, so

Var[longevity] = $5619.38 - 73.2^2 = 261.14$ or

SD[longevity] = $\sqrt{261.14} = 16.2$

[Think of it as the SD when the 'n' is 1 000 or 1 000 000]

However, we see from diagram in earlier section that the distribution of longevity is not Gaussian, so a standard deviation would not be very helpful in describing limits of individual variation (%-iles would be better)

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

If waiting for one of 3 unevenly spaced elevators
(all equally likely to arrive next),
where (x) do you stand? what criterion does it imply?

	0	1		5		<--elevators
						average distance
x						
0	1			5		2.00
	x					
0.5	0.5			4.5		1.83
	x					
1	0			4		1.67
	x					
1.5	0.5			3.5		1.83
	x					
2	1			3		2.00
	x					
2.5	1.5			2.5		2.17
	x					
3	2			2		2.33
	x					
3.5	2.5			1.5		2.50
	x					
4	3			1		2.67
	x					
4.5	3.5			0.5		2.83
	x					
5	4			0		3.00

The mean minimizes average squared deviation.
The median minimizes the average absolute deviation.

	0	1		5		<--elevators
						average squared distance
x						
0	1			25.00		8.67
	x					
0.25	0.25			20.25		6.92
	x					
1	0			16.00		5.67
	x					
2.25	0.25			12.25		4.92
	x					
4	1			9.00		4.67
	x					
6.25	2.25			6.25		4.92
	x					
9	4			4.00		5.67
	x					
12.25	6.25			2.25		6.92
	x					
16	9			1.00		8.67
	x					
20.25	12.25			0.25		10.92
	x					
25	16			0.00		13.67

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

e.g. Expectation & Variance of Random Digits 0 - 9

x	Prob	x ²	x•prob	x ² •prob
0	0.1	0	0.0	0.0
1	0.1	1	0.1	0.1
2	0.1	4	0.2	0.4
...
...
7	0.1	49	0.7	4.9
8	0.1	64	0.8	6.4
9	0.1	81	0.9	8.1
	1.0		4.5	28.5

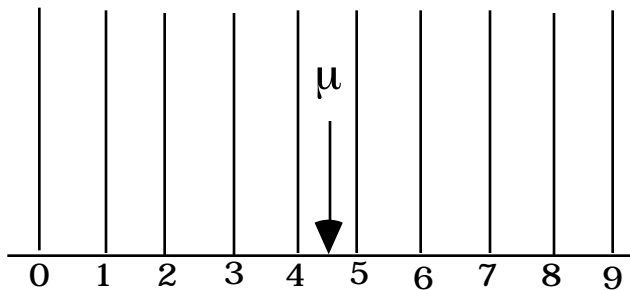
$$\text{Var}(x) = E[x^2] - \{E[x]\}^2 = 28.5 - 4.5^2 = 8.25$$

[Variance = ave. square minus squared ave.]

$$\text{SD}(x) = \sqrt{\text{Var}[x]} = 2.9$$

Relative frequency

0.1



Expectation, Variance, and SD of a Binary [0 / 1] RV

X = 0 with probability p(0) = 1-
 X = 1 with probability p(1) =

In other words...

A proportion π of the individual elements in the population are positive (X = 1); the remaining fraction or proportion $1-\pi$ are negative (X=0)

$$\begin{aligned} E(X) &= 0 \cdot p(0) + 1 \cdot p(1) \\ &= 0 \cdot (1 - \pi) + 1 \cdot \pi \\ &= \pi \end{aligned}$$

$$\begin{aligned} \text{VAR}(X) &= E(X^2) - \{E(X)\}^2 \\ &= 0^2 \cdot p(0) + 1^2 \cdot p(1) - \pi^2 \\ &= 0 + 1 \cdot \pi - \pi^2 \\ &= \pi - \pi^2 = \boxed{\pi(1 - \pi)} \end{aligned}$$

$$\text{SD}(X) = \sqrt{\text{VAR}} = \boxed{\sqrt{\pi(1 - \pi)}}$$

*This "Bernoulli" Random Variable is a **key** one in **Epidemiology** -- it is the 'kernel' or 'atom' in the molecules called Binomial Random Variables. The unit variance $\pi[1-\pi]$ and its square root show up whenever we deal with 0/1 data.*

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

Application: Meta-analyses...Reducing uncertainty by averaging estimates

[From Midterm Exam, May 94] A overview of randomized clinical trials of antiplatelet therapy as prophylaxis against deep venous thrombosis [BMJ on 22 Jan. 1994] found the following:

Category of trial	% odds reduction (SD)
general surgery	37% (8)
traumatic orthopaedic surgery	31% (13)
elective orthopaedic surgery	49% (11)

- a Is 37% a statistic or parameter?
STATISTIC, calculated from sample of data.
- b Does each SD refer to (i) variation of individuals or (ii) sampling variation associated with the estimate? Explain your reasoning.

(ii) sampling variation associated with estimate of % reduction. Basic data on individuals are 1's and 0's [thrombosis or not]. In each study, % reduction is derived from two binomial statistics.

- c Use the SD of each estimate to argue that the apparent heterogeneity in the percent reductions, i.e. the spread from 31% to 49%, could simply reflect random variation alone [differences among three estimates are more difficult than we have learned to dealt with, so for simplicity, concentrate on the difference of two estimates]

Don't know how symmetric/Gaussian the sampling variation estimate of % reduction would be, but as a first approximation, could expect, even if the reduction were the same in the two subtypes, the random difference in any two samples would be non-zero, and would be Gaussian with SD approximately equal to $\sqrt{11^2 + 13^2} = 17$. So an observed difference of $49 - 31 = 18$ would not be that unusual. One sees same thing if plots the CI's. Moreover, the SD of 17 refers to 2 random samples, not the 2 furthest apart of 3 random samples.

Some of you took SD of the 3 estimates; but SD associated with each estimate reflects sample sizes etc. Also, SD and SE interchangeable here.

Since we have neither a statistical nor a biologic basis for assuming different size effects for different types of patients, in the spirit of Occam's razor, we can construct one overall estimate from the three. One way to do this is take a simple average of the three reductions, giving each estimate a weight of 1/3 i.e. $(37+31+49)/3 = 39\%$.

- d If we create this equal-weighted average, the uncertainty {SD} associated with it should be smaller than the SD of components. Calculate SD for

$$\frac{1}{3}\text{estimate}_1 + \frac{1}{3}\text{estimate}_2 + \frac{1}{3}\text{estimate}_3$$

[§4.3 & 4.64–4.66 should help]

using rules for Var(sum) and var(constant X) and defn. of SD

$$\text{var} \left[\frac{1}{3}\text{estimate}_1 + \frac{1}{3}\text{estimate}_2 + \frac{1}{3}\text{estimate}_3 \right]$$

$$= \frac{1}{9} 64 + \frac{1}{9} 169 + \frac{1}{9} 121 = 39.33$$

$$\text{so SD} = \sqrt{39.33} = \underline{6.27}$$

Since we have three estimates with different degrees of uncertainty, it makes sense to calculate an average of them which gives more weight to the individual estimates with smaller SD's. It can be shown mathematically that the weighted average with the lowest SD is the one with weights that are inversely proportional to the individual variances. In our example here, this would lead to weights that are proportional to $\frac{1}{64}$, $\frac{1}{169}$ and $\frac{1}{121}$ respectively, or an overall estimate of

$$0.52\text{estimate}_1 + 0.20\text{estimate}_2 + 0.28\text{estimate}_3.$$

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

This gives a weighted average of just over 39%. [the fact that the two methods give almost the same answer is a coincidence in this example; it doesn't happen generally]

- e This information-weighted average has a lower uncertainty {SD} associated with it than a simple equally-weighted ave.. Calculate SD for

$$0.52\text{estimate}_1 + 0.20\text{estimate}_2 + 0.28\text{estimate}_3 ;$$

compare with SD in d.

$$0.52\text{est}_1 + 0.20\text{est}_2 + 0.28\text{est}_3 =$$

$$0.52(37\%) + \dots + 0.28(49\%) = \underline{39.16\%}$$

$$\text{var}[0.52\text{estimate}_1 + 0.20\text{estimate}_2 + 0.28\text{estimate}_3]$$

$$= 0.52^2 \text{ var}[\text{est}_1] + 0.20^2 \text{ var}[\text{est}_2] + 0.28^2 \text{ var}[\text{est}_3]$$

$$= 0.2704 (64) + 0.0400 (169) + 0.0784 (121)$$

$$= \underline{33.55}$$

so $SD = \sqrt{33.55} = 5.79$, smaller (by definition) than SD above.
NOTE THAT WHEN WEIGHTS ARE THE INVERSE OF THE VARIANCES, THE VARIANCE OF THE WEIGHTED AVERAGE EQUALS THE HARMONIC MEAN OF THESE VARIANCES DIVIDED BY THE NUMBER OF COMPONENTS IN THE WEIGHTED AVERAGE

The overview reported estimate of 42% (SD 17) for high risk medical patients.

- f Combine the single estimate for surgical patients and the 42% for medical patients. Calculate its SD. Why does the 'averaging' of the two estimates not diminish the SD very much?

$$\text{estimate}_{\text{surg}} = 39.16\% \text{ with } SD[\text{estimate}_{\text{surg}}] = 5.79,$$

$$\text{var} = 33.55;$$

$$\text{estimate}_{\text{med}} = 42.00\% \text{ with } SD[\text{estimate}_{\text{med}}] = 17.00;$$

$$\text{var} = 289.00;$$

optimal weights proportional to

$$\frac{1}{5.79^2} \text{ and } \frac{1}{17^2}$$

$$\text{or } \underline{0.90} \text{ and } \underline{0.10}$$

$$0.90\text{est}_{\text{surgical}} + 0.10\text{est}_{\text{medical}}$$

$$= 0.90(39.16\%) + 0.10(42\%)$$

$$= 39.44\%$$

Overall estimate closer to 39.16% than 42% because weighted 9:1

$$\text{var}[0.90\text{estimate}_{\text{surgical}} + 0.10\text{estimate}_{\text{medical}}]$$

$$= 0.90^2 (33.55) + 0.10^2 (289.00) = 30.07$$

$$\text{so } SD = \sqrt{30.07} = 5.49$$

SD only slightly smaller: estimate dominated by estimate_{surgical}.

Random Variables; Probability Distributions ; Expectation and Variance of a Random Variable

Law of cancellation of extremes and reduction of uncertainty (how insurance companies stay solvent)

Possible Earnings from single insurance policy and from pool of n insurance policies:

Earnings from a single policy (n=1)

X=Earnings	Prob	X • Prob	X ² • Prob
-\$19,900	0.00183	-\$36.417	724,698.3
-\$19,800	0.00186	-\$36.828	729,194.4
-\$19,700	0.00189	-\$37.233	733,490.1
-\$19,600	0.00191	-\$37.436	733,745.6
-\$19,500	0.00193	-\$37.635	733,882.5
\$500	0.99058	\$495.290	247,645.0
	1.00000	\$309.741	3,902,655.9

Expected (average) Earnings per policy
 = Earnings x Probability = \$309.74

Variance(Earnings) = ave(Earnings²)
 - (ave Earnings)²
 = 3,902,655.9 - 309.74²
 = 3,806,717 (\$²)

Std Dev(Earnings) = $\sqrt{\text{Var}(\text{Earnings})}$ = \$1,951

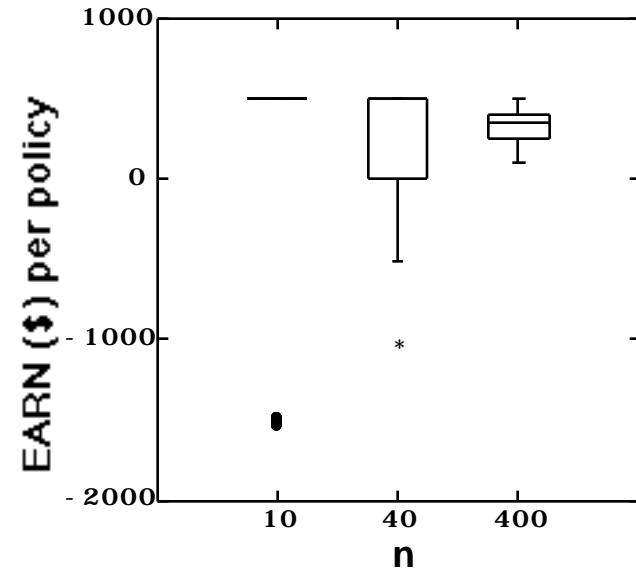
Earnings per policy from a pool of n policies

Statistics for earnings from pooled policies based on several simulations per pool size

n:	<u>1</u>	<u>10</u>	<u>40</u>	<u>400</u>
MINIMUM	-29,900	-1,540	-1,022	96
MAXIMUM	500	500	500	500
MEAN	309	268	318	320
STD DEV	1,951	645	309	92
\$1,951/ n	\$1,951	\$617	\$308	\$98

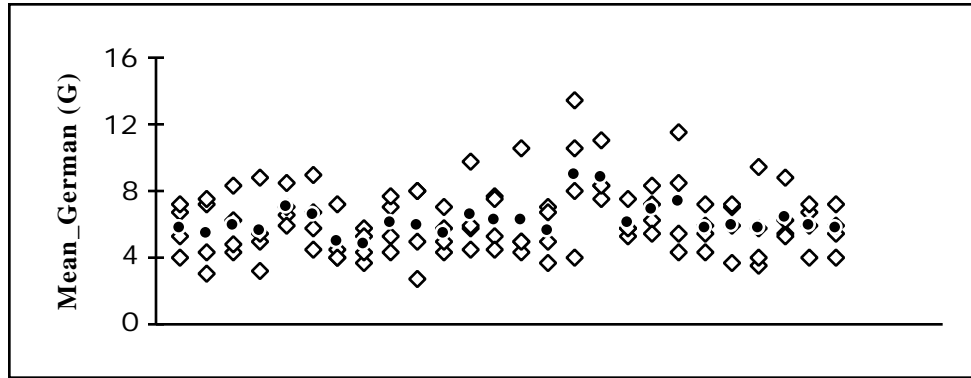
Note: This example is from Q5.22 page 358 of 1st Edition of Moore and McCabe [Q4.48 in 2nd edition has \$100,000 policy and \$250 premium per year, but principle is same]

Earnings from pool of n policies



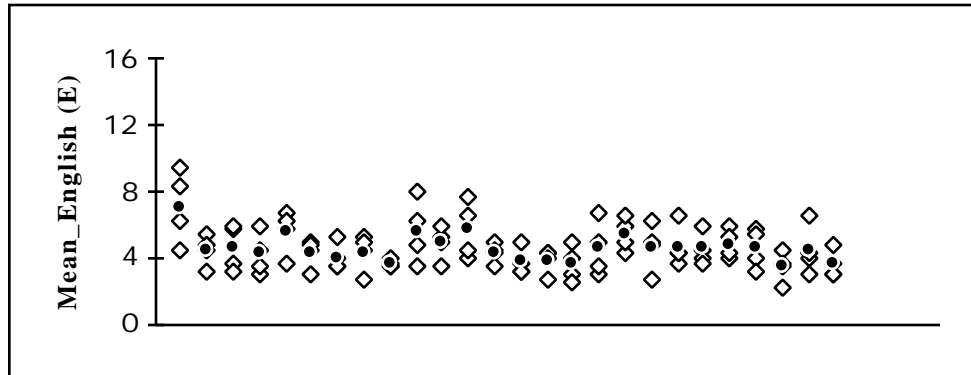
Note that the SD[mean of n policies] in simulations is quite close to that predicted theoretically, namely \$1,951/ n

Variation in the mean word length in samples of sizes $n=4(\diamond)$ and $n=16(\bullet)$, and in the differences of two means (G - E)
 [each \diamond and \bullet represents a sample from a student in course 607 in a previous year]



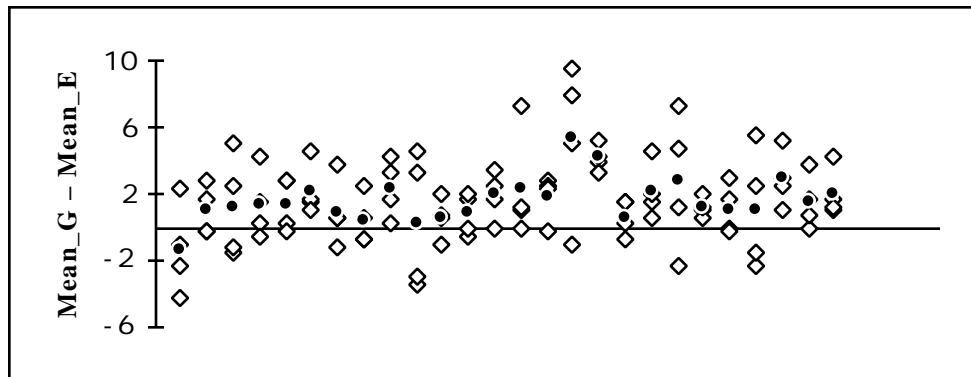
---- GERMAN ----

SD of means based on	SD of means based on
() n=4	(•) n=16
1.97	0.97



---- ENGLISH ----

SD of means based on	SD of means based on
() n=4	(•) n=16
1.33	0.81



GERMAN – ENGLISH

SD of means based on	SD of means based on
() n=4	(•) n=16
2.40	1.30