

1 BELLSOUTH TELECOMMUNICATIONS, INC.

2 DIRECT TESTIMONY OF EDWARD J. MULROW, PH.D.

3 BEFORE THE FLORIDA PUBLIC SERVICE COMMISSION

4 DOCKET NO. 000121-TP

5 MARCH 21, 2001

6

7

8 Q. PLEASE STATE YOUR NAME, AND BUSINESS NAME AND ADDRESS.

9

10 A. My name is Edward J. Mulrow. I am employed by Ernst & Young LLP as a
11 Senior Manager in the Quantitative Economics and Statistics Group. I have
12 been retained by BellSouth as a statistical advisor. My business address is
13 1225 Connecticut Ave., NW, Washington, DC 20036.

14

15 Q. ARE YOU THE SAME EDWARD J. MULROW THAT FILED DIRECT
16 TESTIMONY IN THIS DOCKET?

17

18 A. Yes. I filed direct testimony in this docket on March 1, 2001.

19

20 Q. WHAT IS THE PURPOSE OF YOUR REBUTTAL TESTIMONY?

21

22 A. The purpose of my rebuttal testimony is to respond to portions of the direct
23 testimonies of Dr. Robert M. Bell representing the ALEC Coalition, and Dr.
24 George S. Ford representing Z-Tel Communications. In responding to the
25 direct testimony of these witnesses, I address the following issues:

- 1 • The appropriate statistical methodology for making performance measure
- 2 parity comparisons.
- 3 • Dr. Bell's analysis of the impact of "delta."
- 4 • The use of a floor for the balancing critical value.
- 5

6 *1. The appropriate statistical methodology for making performance measure*
7 *parity comparisons.*

8 Q. THE ALEC COALITION, REPRESENTED BY DR. ROBERT BELL,
9 PROPOSES THAT THE FLORIDA COMMISSION ORDER THE MODIFIED
10 Z AS A COMPONENT OF THE STATISTICAL METHODOLOGY.
11 PLEASE RESPOND.

12
13 A. As I said in my direct testimony, the appropriate methodology to use in
14 situations where transaction level data is available and a BellSouth retail
15 analog exists is the Truncated Z with Error Probability Balancing. This
16 methodology is described in the Louisiana PSC "statistician's report" which is
17 attached to my direct testimony as Exhibit EJM-1. One of the more interesting
18 things about Dr. Bell's position is that it was another AT&T witness, now
19 retired, who basically created the truncated Z formulas that BellSouth is now
20 offering. Indeed, as I mentioned in my direct testimony, the methodology was
21 developed in a joint effort between AT&T's statistical expert Dr. Colin
22 Mallows (the AT&T witness who is now retired), and the Ernst & Young
23 statistical team. I find it difficult to understand how AT&T and BellSouth
24 could have expended such effort to reach a methodology that was satisfactory
25 to the experts representing each party, only to have AT&T seemingly walk

1 away from that methodology.

2

3 Q. CAN YOU EXPLAIN THE FUNDAMENTAL DIFFERENCES BETWEEN
4 WHAT DR. BELL PROPOSES AND WHAT BELL SOUTH IS
5 PROPOSING?

6

7 A. As with many things involving statistics, the explanation is a bit complex, but I
8 will try to explain in as clear a fashion as possible. You will recall that in my
9 direct testimony, I discussed how BellSouth's methodology took various
10 measures down to what I called the individual "cell" level. The purpose of
11 creating "cells" was to break each comparison down to its most basic
12 components, so that we could be relatively sure that we were comparing
13 "apples-to-apples." For instance one of the cells would be a new residential
14 provisioning order that is non-dispatched with less than 10 circuits that
15 occurred in the first part of the month in a particular wire center. We would
16 compare BellSouth's transactions that met those criteria as well as the ALEC
17 transactions that met the criteria. We would determine the mean for both
18 samples and would calculate a modified Z statistic for that cell. After doing
19 this, we would roll this cell up with other cells related to plain old telephone
20 service and would essentially aggregate all of the individual modified Z
21 statistics into a single statistic. As I explained in my direct, when we rolled
22 these individual statistics into a single statistic, we assign a value of zero to all
23 of the statistics that have a positive value, so that we do not mask the impact of
24 any negative values. This changing of positive values to zero is why we call
25 the resulting statistic a truncated Z statistic.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

On the other hand, Dr. Bell's process essentially stops with the calculation of the modified Z statistic for each of his "sub-measures" which can generally be thought of as being conceptually the same as BellSouth's cells, except that Dr. Bell has disaggregated his sub-measures in a different way, and has not disaggregated them to the levels that BellSouth has proposed. For instance, and this is taking an extreme example, Dr. Bell's proposal would essentially have BellSouth stop at the cell level discussed above, and make the final comparison about whether parity is being provided right there.

Q. SO, BASED ON THIS DISCRPTION, WHAT IS CONCEPTUALLY WRONG WITH DR. BELL'S APPROACH?

A. I can explain this most clearly by looking again at what BellSouth has done. Lets assume that there were 2000 "cells" associated with the provisioning of plain old telephone service. If we looked at the individual cells, we might find 75 that revealed an apparent discrepancy between BellSouth's performance and that provided to the ALEC. Based on these failures, relying solely on the modified Z statistic, BellSouth would be expected to pay a penalty. The problem is that we know that we are going to get some Type 1 errors in a statistical analysis like this. For instance, if we were willing to accept that 5% of the observations were going to be Type I errors, you would expect to see 100 failures. Viewed in this light, 75 failures would be well within the expected parameters. The point is, if you looked at the individual "cells" you would conclude there was a problem, but when you look at the whole picture

1 you see that there is not. Dr. Bell's approach, relying solely on the modified Z
2 statistic for individual "cells" or "sub-measures" doesn't allow this to happen.

3

4 Q. IS THE IDEA OF "AGGREGATING" MANY STATISTICAL RESULTS TO
5 MAKE AN OVERALL DETERMINATION OF PARITY SOMETHING
6 THAT WAS DEVELOPED SPECIFICALLY FOR THE LOUISIANA
7 STATISTICIAN'S REPORT?

8

9 A. No. The concept has been around for quite some time. It is sometimes
10 referred to as a "multiple testing problem."

11

12 Q. WILL YOU BRIEFLY EXPLAIN WHAT A "MULTIPLE TESTING
13 PROBLEM" IS?

14

15 A. Certainly. As I pointed out in my example above, whenever one performs
16 numerous statistical tests at one time, it needs to be recognized that some of the
17 tests will provide results indicating a problem even when there is no problem.
18 By looking at all the results in a more global way, one can determine if the
19 "failed" tests that are observed represent a true problem, or just random chance.

20

21 Q. IS BELLSOUTH THE ONLY COMPANY SUGGESTING THAT THIS
22 FORM OF AGGREGATING SHOULD BE DONE?

23

24 A. No. The four states where the FCC has granted an RBOC the right to market
25 long distance services have performance comparison plans that aggregate the

1 results of many comparisons into an overall result that determines
2 parity/disparity.

3

4 In New York, Verizon uses a weighted average of performance scores to make
5 parity judgments. In Texas, Oklahoma, and Kansas, Southwestern Bell uses
6 the "K-value" method. This "K-value" methodology was described by
7 AT&T's Dr. Mallows in an affidavit to the FCC in May 1998¹. Thus, both of
8 the methods of aggregation that AT&T's expert has suggested have been
9 adopted by former Bell Companies for use in their performance plans. AT&T
10 however, appears reluctant to accept either of these methodologies.

11

12 Q. WHY IS THE TRUNCATED Z STATISTIC A BETTER AGGREGATE
13 STATISTIC THAN THE OTHERS THAT ARE IN USE, SAY IN TEXAS?

14

15 A. As I explained in my direct testimony, the truncated Z statistic was created so
16 that it possesses five important properties.

17

18 1. It is a single, overall index on a standard scale; that is, you can use a
19 probability distribution to make judgments.

20 2. If transaction counts for BellSouth and the ALEC across comparison cells
21 (classifications) are exactly proportional, the aggregate index should be
22 very nearly the same as if we had not disaggregated. This means that

¹ Affidavit of Dr. Colin L. Mallows before the Federal Communications Commission
1998.

1 granular disaggregation I have discussed really wasn't necessary, you will
2 still get the same results.

3 3. The contribution of each cell depends on the number of transactions in the
4 cell.

5 4. As far as possible, systematic discriminatory performance in some cells is
6 not masked by good performance in other cells.

7 5. The final result does not depend critically on minor details in the data; that
8 is, small changes in transaction values only induce small changes in the
9 final result.

10

11 In addition to these important properties, the error probabilities of the truncated
12 Z test can be balanced. The other statistics used when aggregating results do
13 not meet all of the criteria that I have outlined above, but I would note that any
14 of them would be better than what Dr. Bell is proposing, which is that no
15 aggregation be done at all.

16

17 Q. ARE THERE OTHER DIFFERENCES BETWEEN WHAT DR. BELL HAS
18 PROPOSED AND THE BELL SOUTH METHODOLOGY?

19

20 A. Yes, and there is one in particular that should be considered. The formulae that
21 Dr. Bell proposes for testing proportion and rate measures do not easily lend
22 themselves to balancing Type I and II error probabilities. This creates a
23 methodological inconsistency between the test Z statistics he recommends and

1 the balancing critical value. I will discuss this in more detail later in my
2 testimony. In order to explain the root of the problem, however, I need to tell
3 you something more about statistics. What we have been discussing in the
4 examples above is a comparison of “means,” that is, we take the average of the
5 BellSouth transactions in the “cell” and compare that “average” or “mean” to
6 the comparable “mean” of the ALEC transactions. Not all observations lend
7 themselves to the calculation of “means,” however. For instance, consider
8 “missed appointments.” With “missed appointments” you are looking at the
9 percentage of the total number of scheduled appointments that were missed.
10 As a result, you end up with a proportion, such as a tenth of a percent or 5
11 percent or whatever figure is appropriate. You do not have a mean per se.
12 Another example is what we call a “rate” such as the “customer trouble report
13 rate”, where you are looking at the number of troubles BellSouth or the ALEC
14 has per the number of available lines. Unlike the “proportional” measures
15 described above, which would always have to be less than 1, the measurement
16 of a “rate” could exceed 100 percent. For instance, if you had ten access lines
17 and 12 reported troubles (that is some lines have more than a single trouble
18 during the reporting period) you can get more than a figure of 100 percent.
19 Again, these two special categories are simply different measures than the
20 “means” calculation that we have been talking about.
21
22 The root of the problem is that Dr. Bell uses the modified Z concept
23 irrespective of whether the measure is one based on “means,” “proportions,” or

1 “rates.” The difficulty from a statistical perspective is that the concept the
2 modified Z statistic is based on should not be applied across the board to all
3 measure types. Specifically, the basis for the modified Z statistic is that you
4 take the difference between the two “means” in the particular “cell” or sub-
5 measure, and divide the result by the standard deviation of BellSouth’s mean.
6 This is done to make the test sensitive to changes in the ALEC standard
7 deviation (compared to the BellSouth standard deviation) that would be
8 harmful to the ALEC. In other words, BellSouth could try to give the same
9 average service to ALEC customers as to its own customers, but do so in a way
10 that some ALEC customers receive longer completion times. For example,
11 suppose that BellSouth always services its own customers in 2 days. BellSouth
12 could service one-third of the ALEC customers in 1 day, one-third in two days,
13 and the remaining third in 3 days. On the average, the ALEC service times are
14 the same as BellSouth’s, but one-third of the ALEC customers received service
15 that was “below” average. Dividing the difference between the means by only
16 BellSouth’s standard deviation avoids masking this problem.

17
18 The same situations cannot occur for “proportion” or “rate” measures. In the
19 case of a proportion, such as “missed appointments” that is stated as a
20 percentage of total appointments scheduled, you only have one parameter to
21 consider, the proportionality. As a result, BellSouth cannot separately control
22 the proportion value and the variability about that value.

1 Q. ARE THERE DIFFERENT FORMULAS THAT DR. BELL COULD HAVE
2 USED TO ADDRESS THESE ISSUES?

3

4 A. Yes, but he did not do so, which is another reason why BellSouth's approach
5 makes more sense.

6

7 **II. Dr. Bell's analysis of the impact of 'delta.'**

8 Q. IN ORDER TO SUPPORT HIS CHOICE OF A "DELTA" VALUE OF 0.25,
9 DR. BELL PROVIDES A TABLE SHOWING THE PERCENTAGE OF
10 ALEC CUSTOMERS RECEIVING BAD SERVICE, BY BELLSOUTH
11 PERCENT AND DELTA. CAN YOU COMMENT ON THIS TABLE?

12

13 A. Well, there are a couple of interesting points I can make. First, the
14 methodology that Dr. Bell advocates in his testimony and exhibit is not the
15 method that he used to calculate the numbers in the table. Second, the table
16 does not accurately represent the way that BellSouth proposes to carry out
17 balancing for proportion measures.

18

19 Q. WOULD YOU EXPLAIN YOUR FIRST POINT MORE FULLY?

20

21 A. The overall concepts for balancing error probabilities were first developed for
22 mean performance measures. As previously described, these are measures that
23 represent the average of a measured amount, for example the average time it
24 takes to complete an order. In this case, "delta" represents the difference
25 between the ILEC and ALEC averages in terms of an ILEC standard deviation.

1 When solving for the balancing critical value, it turns out to be mathematically
2 convenient to define the alternative hypothesis this way, given the form of the
3 modified z statistic for a mean measure.

4
5 A proportion measure, on the other hand, measures the fraction of transactions
6 that possess a certain quality or attribute out of all transactions. For example,
7 percent missed installations measures the fraction of all provisioning orders in
8 a month where service was not completed on or before the assigned due date.
9 Often in statistics, methods that are worked out for means can be used on
10 proportions because a proportion can be considered as a special type of mean.
11 However, a proportion is a fraction of the whole, so it can only be a number
12 between 0 and 1 (or equivalently between 0 percent and 100 percent).

13
14 Now, if we want to describe “delta” for a proportion measure as the difference
15 between the ILEC and ALEC proportions in terms of an ILEC standard
16 deviation we have to be careful. The mathematical convenience present in the
17 mean measure case is not present with a proportion measure. Thus a different
18 method is needed.

19
20 Q. WHAT METHODS ARE AVAILABLE FOR BALANCING A
21 PROPORTION MEASURE?

22
23 A. We have identified two ways to approach balancing for proportion measures.
24 One way is to transform the proportion using the arcsine square root
25 transformation. This is what Dr. Bell used to create Table 1 on page 13 of his

1 direct testimony. The other way is to use a concept called the “odds” ratio.

2

3 Q. WHAT IS THE REASONING BEHIND USING THE ARCSINE SQUARE

4 ROOT METHOD?

5

6 A. Because “delta” for a proportion measure cannot be defined using a

7 straightforward analogy with the definition for a mean measure, a

8 transformation is used. This allows us to use the same formula to compute the

9 balancing critical value as was used in the mean measure case. However, two

10 problems arise: 1) the interpretation of “delta” is related to the transformed

11 measure, and 2) the z statistic that is used in the test should also use the

12 transformed measurement.

13

14 Q. WHAT IS THE INTERPRETATION OF “DELTA” WHEN THE ARCSINE

15 SQUARE ROOT TRANSFORMATION IS USED?

16

17 A. “Delta” becomes twice the difference of the transformed ILEC proportion with

18 the transformed ALEC proportion. It is no longer the difference between the

19 performance measures in terms of an ILEC standard deviation.

20

21 Q. YOU STATED THAT THE Z STATISTIC USED IN THE TEST SHOULD

22 ALSO USE THE TRANSFORMED MEASUREMENTS. COULD YOU

23 EXPLAIN THIS?

24

25 A. Yes. In order to arrive at the same balancing critical value fr

1 proportion measure as that of a mean measure, you must redefine the basic Z
2 statistic. When using the arcsine square root transformation the Z statistic
3 should be

$$Z = \frac{2 \left(\arcsin \left(\sqrt{\hat{p}_{ILEC}} \right) - \arcsin \left(\sqrt{\hat{p}_{ALEC}} \right) \right)}{\sqrt{\frac{1}{n_{ILEC}} + \frac{1}{n_{CLEC}}}}.$$

6
7 Q. IS THIS WHAT DR. BELL IS RECOMMENDING TO USE FOR
8 PROPORTION MEASURES?

9
10 A. No. The formula he recommends is given in Exhibit RMB-1, page 14, of his
11 direct testimony. This formula is a direct analog of the mean measure formula,
12 but as I have already explained, we need to be cautious in directly applying
13 mean measure formulae to other types of measures when we are using a
14 balancing methodology.

15
16 Q. WHAT ARE THE CONSEQUENCES OF USING A BALANCING
17 CRITICAL VALUE BASED ON THE ARCSINE SQUARE ROOT
18 TRANSFORMATION WITH A Z STATISTIC THAT IS NOT BASED ON
19 THE TRANSFORMATION?

20
21 A. There are many scenarios where the use of the wrong type of Z statistic would
22 find BellSouth to be out of parity when the use of the proper Z statistic would
23 find them in parity. Consider a simple example. Let's suppose that there are

1 1000 BellSouth provisioning orders, and that BellSouth “missed” 214 of the
2 appointments, that is, the orders where not completed on or before the due
3 date. Thus, BellSouth “missed” 21.4 percent of their orders. For the same
4 time period, suppose there were 30 comparable ALEC provisioning orders, and
5 that 8 of these were “missed.” So, BellSouth “missed” 26.7 percent of the
6 ALEC’s orders. Now the balancing critical value based on the arcsine square
7 root transformation and the “delta” of 0.25 that Dr. Bell uses is -0.675. If we
8 use the modified Z formula given in Dr. Bell’s direct testimony, we will get a
9 Z score of -0.693. Since this is less than the critical value (further from zero on
10 the negative side), we would conclude that there is a lack of parity, and
11 BellSouth would pay a penalty. On the other hand, if we use the Z formula
12 given above, which is based on the arcsine square root transformation, we get a
13 Z value of -0.666. In this case, we would say that BellSouth is compliant, and
14 there would be not a penalty assessment.

15
16 Q. SO YOU ARE SAYING THAT THE BASIC METHODOLOGY THAT IS
17 USED TO CALCULATE THE BALANCING CRITICAL VALUE NEEDS
18 TO BE MATCHED WITH THE SAME BASIC METHODOLOGY THAT IS
19 USE TO CALCULATE THE Z TEST STATISTIC.

20
21 A. Yes, and there appears to be an inconsistency in what Dr. Bell is
22 recommending for proportion measures as well as rate measures.

23
24 Q. YOU SAID THERE IS ANOTHER METHOD FOR BALANCING
25 PROPORTION MEASURES THAT IS BASED ON THE “ODDS” RATIO.

1 WHAT IS AN "ODDS" RATIO?

2
3 A. The "odds" ratio is what BellSouth has used when the information in the
4 "cells" involves proportions, which I have been discussing, rather than
5 "means." The "odds" methodology is relatively straightforward. First we need
6 to define the odds of an event such as a missed installation occurring. Odds are
7 the ratio of the probability of an event occurring to the probability that the
8 event won't occur. So, if BellSouth "missed" 21.4 percent of the installations
9 to their own customers, then the odds of a customer experiencing a "miss" is
10 found by dividing the probability of a "miss," 0.214, by the probability of an
11 "on-time" installation, 0.786 ($= 1 - 0.214$). This gives the odds of a "miss" as
12 0.276. In odds terminology, we might say that the odds of a BellSouth
13 customer experiencing a "miss" are approximately 1 to 2.6.

14
15 The odds ratio for "missed" provisioning installations is the ALEC customer's
16 odds of a "miss" divided by the BellSouth customer's odds of a "miss." When
17 this odds ratio is one or less, BellSouth is delivering parity or better service to
18 the ALEC's customers. When this odds ratio is greater than one, then
19 BellSouth is not necessarily delivering parity service. Under a balancing
20 approach, we need to determine an odds ratio greater than one to use for the
21 balancing alternative hypothesis.

22
23 Q. IS THE ODDS RATIO EASIER TO INTERPRET THAN THE ARCSINE
24 SQUARE ROOT METHOD?

1 A. Not necessarily. Many people have trouble interpreting odds, and relating the
2 value back to the probability of an event occurring. However, the
3 interpretation in terms of odds is straightforward. If the odds ratio for “missed”
4 installations is set at 3, then we know that an ALEC customer’s odds of a
5 “miss” is three times greater than that of a BellSouth customer. We would still
6 need a table, such as Dr. Bell’s Table 1, to interpret the actual difference in the
7 performance. I want to say, however, that setting the “odds” ratio at 3, which
8 is what the Louisiana Commission has done for Tier 1 measures, does not
9 necessarily mean that the probability of actually having a disparity is that great.

10

11 Q. CAN YOU PROVIDE US WITH SUCH A TABLE?

12

13 A. Certainly. Figure 1 below will help one interpret the actual difference between
14 the BellSouth proportion and the ALEC proportion for a given “odds” ratio.
15 The table shows the percentage of the time an ALEC customer will experience
16 a miss by the BellSouth percentage “missed,” for two values of the odds ratio:
17 2 and 3.

18

19

20

21

22

Figure 1
ALEC Percentage of “Missed” Installations
By BST Percentage and
The Odds Ratio of the Alternative Hypothesis

BST PERCENTAGE MISSED	Odds Ratio	
	2	3
1	2	3
5	10	14
10	18	25
20	33	43

23

1 We see from the first row of this table that for an alternative hypothesis with an
2 odds ratio of 3, the ALEC percentage of “missed” installations is about 3
3 percent when the BST percentage is 1 percent. However, the ALEC
4 percentage is about 43 percent when the BST percentage is 20 percent. So
5 when the BST percentage is close to 0, the ALEC percentage is about 3 times
6 larger at the balancing alternative hypothesis. As the BST percentage get
7 larger, the ratio of the ALEC percentage to the BST percentage gets smaller;
8 converging to 1 as the BST percentage approaches 100 percent.

9
10 Q. THIS SEEMS TO SUGGEST THAT IF BELL SOUTH HAS A MISS OF 20
11 PERCENT, THAT A MISS OF UP TO 43 PERCENT WOULD BE
12 ACCEPTABLE FOR THE ALECS. IS THIS CORRECT?

13
14 A. No, that misses the point completely and that is what is wrong, in large
15 measure with Dr. Ford’s analysis, which I will discuss in more detail below.
16 However, to put point on this, with numbers like that, with a very small sample
17 size the methodology would show BellSouth out of parity almost 60 percent of
18 the time and as the sample size approached a thousand transactions for
19 BellSouth and only fifty for the ALEC, the probability that parity will not be
20 concluded approaches 100 percent (see Table 3 below). I realize this is not
21 intuitive, and I will discuss it more below, but it would be a mistake to
22 conclude that the odds ratio balancing test allows the ALECs to experience
23 significantly worse performance than BellSouth without detecting a failure to
24 provide parity on BellSouth’s part. I would also note that the same holds true
25 for Dr. Bell’s calculations using the arcsine square root method where he

1 shows a similar disparity. Once the sample size gets to the levels that I have
2 just mentioned, the probability of finding a disparity at those levels approaches
3 100 percent.

4
5 Q. IF THE ODDS RATIO METHOD IS USED FOR DEFINING THE
6 BALANCING CRITICAL VALUE, HOW DOES THAT EFFECT THE
7 FORMULA THAT IS USED TO CALCULATE THE CRITICAL VALUE?

8
9 A. The balancing critical value for a proportion measure is based on a different
10 formula than that of a mean measure when an odds ratio approach is used. The
11 formula is more complicated than the mean measure formula, and it is given in
12 Appendix C of the Louisiana "Statistician's Report" (Exhibit EJM-1 of my
13 direct testimony).

14
15 Q. DOES THE Z STATISTIC USED TO COMPARE THE PERFORMANCE
16 MEASURES NEED TO BE MODIFIED WHEN USING THE ODDS RATIO
17 APPROACH?

18
19 A. I was able to derive the balancing critical value formula based on the odds ratio
20 because it "fit in" with the method used to calculate the cell level Z statistic.
21 This Z statistic that I refer to is given in Appendix A of the Louisiana
22 Statistician's Report. As previously alluded to, it differs from the Z statistic
23 given by Dr. Bell in his testimony.

24
25 Q. SO THE Z STATISTIC FOR PROPORTIONS PROFFERED BY DR. BELL

1 NEEDS TO BE MODIFIED, REGARDLESS OF THE BALANCING
2 APPROACH, IN ORDER TO HAVE THE BALANCING METHODOLOGY
3 CONSISTENT WITH THE BASIC Z STATISTIC METHODOLOGY.
4

5 A. Yes. I believe that we should try to be consistent with the Z statistic
6 methodology when developing the methods for balancing. That's not to say
7 that a balancing methodology cannot be worked out for the proportion Z
8 statistic in Dr. Bell's testimony, but I think it would make a complex problem
9 messier. Dr. Bell may also be able to show that a balancing critical value
10 based on a method different from the one used to create the LCUG proportion
11 Z statistic is a reasonable approximation under certain circumstances. The data
12 that we have examined so far exhibit many different characteristics, so it is
13 easy to find cases when the approximations break down. In fact all of the
14 balancing methods break down when both BellSouth and ALEC transaction
15 counts get very small. So, none of the methods we've looked at are perfect. I
16 do believe that we should do our best to avoid problems that we can identify,
17 and consistency between Z statistic methods and balancing methods helps.
18

19 Q. YOU SAID THAT YOU DID NOT THINK THAT DR. BELL'S TABLE 1
20 REPRESENTS THE WAY IN WHICH BELL SOUTH WILL CARRY OUT
21 BALANCING FOR PROPORTION MEASURES. WILL YOU EXPLAIN
22 THIS?
23

24 A. BellSouth has chosen to use the "odds" ratio approach to balancing. In fact,
25 the Louisiana Public Service Commission has ordered BellSouth to use an

1 odds ratio of 3 for Tier I testing of proportion measures, and an odds ratio of 2
2 for Tier II testing. So Figure 1 above shows the impact of the choice of an
3 “odds” ratio based on BellSouth’s proportion measure balancing position.
4

5 Q. IS THERE ANY WAY TO TRANSLATE BETWEEN THE TWO
6 METHODS?
7

8 A. Yes, the Louisiana “Statistician’s Report” provides equations that can be used
9 to translate between the two methods. Things are not that straightforward
10 however. You must have an idea of what the BellSouth proportion is in order
11 to translate between methods. For a proportion measure, we can determine
12 what the largest “delta” value will be for a fixed odds ratio over the whole
13 range of proportion values. For instance, with an odds ratio of 3, the largest
14 value of “delta” based on the arcsine square root method is about 0.54. This
15 occurs when the BellSouth percentage of “misses” is about 37 percent. For
16 percentages smaller or larger than 37 percent, the equivalent delta for an odds
17 ratio of 3 is smaller than 0.54. The equivalent delta gets very close to zero
18 when the BellSouth percentage of “misses” is close to 0 or 100 percent.
19

20 **III. The use of a floor for the balancing critical value.**

21 Q. DR. FORD STATES IN HIS TESTIMONY THAT HE BELIEVES THERE IS
22 A SERIOUS FLAW IN THE ERROR PROBABILITY BALANCING
23 METHODOLOGY, AND THAT A LIMIT ON THE BALANCING
24 CRITICAL VALUE NEEDS TO BE ESTABLISHED TO CORRECT THE
25 FLAW. PLEASE RESPOND.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

A. In reading through his arguments I sense that he is confusing hypothesis testing issues. The key issue that he is confused about is that there is a difference between the probability of a Type I error and the probability of detecting disparity. I also do not believe that Dr. Ford appreciates the problems imposed by the observational nature of a monthly performance incentive plan. I will briefly address these issues, and discuss an error in one of Dr. Ford's graphics.

When all of the statistical issues are properly understood and considered as a whole, I believe that there are no serious flaws in the balancing methodology. Therefore, there is no need for the "fix" that Dr. Ford suggests, namely, a floor on the balancing critical value.

Q. YOU SAY THAT DR. FORD IS CONFUSING THE PROBABILITY OF A TYPE I ERROR WITH THE PROBABILITY OF DETECTING DISPARITY. PLEASE EXPLAIN?

A. Dr. Ford makes several comments in his testimony that suggest that a statistical test with a small Type I error probability has very little power to detect discrimination. For instance, on page 21, lines 15 – 17, he states "At significance levels less than 0.0001 (assuming no more than 500 tests are conducted), balancing performs no function other than to make it nearly impossible to detect discrimination (i.e., reject the null hypothesis)." This is simply not true.

1 First, it needs to be understood that the significance level, i.e. the probability of
2 a Type I error, is the probability of rejecting the null hypothesis (concluding
3 that disparity exists) when, in fact, the null hypothesis is true (BellSouth is
4 providing parity service). This is not the probability that the null hypothesis
5 will be rejected when there is truly a certain amount of disparity in the system.
6 Statisticians refer to that probability as the power of a test because it allows us
7 to know how well a test can detect departures from parity. We can evaluate the
8 power of statistical test based on a balancing methodology, and we can show
9 that the power to detect discrimination beyond the materiality level defined by
10 one-half "delta" is above 50 percent.

11

12 Q. WOULD YOU DISCUSS MATERIALITY AGAIN IN THE CONTEXT
13 THAT WE ARE USING THE TERM IN THIS PROCEEDING?

14

15 A. Certainly. Recall from my direct testimony that as long as the average time
16 taken to provide the relevant service to an ALEC does not exceed the
17 BellSouth mean plus one-half "delta" times the BellSouth standard deviation,
18 then the apparent difference in mean service times would not be material. That
19 is, we would not conclude that BellSouth is providing discriminatory service.
20 To state this another way, one-half delta, the parameter that defines the
21 alternative hypothesis for balancing, is a materiality threshold for the disparity
22 in the service system when a balancing method is used for a mean measure test.

23

24 Q. WOULD YOU PROVIDE US AN EXAMPLE OF THIS?

25

1 A. Yes. Figure 2 shows the probability that a mean measure statistical test will
2 detect a difference in the mean performance of BellSouth and an ALEC when
3 the balancing alternative hypothesis uses a “delta” of 1. To calculate these we
4 assume that the true disparity is 0, 0.2, 0.45, etc. For the purpose of this
5 example I am defining the “true disparity” as the numbers indicated across the
6 top of the chart. This is not an observable figure; I am assuming the disparity
7 to exist to illustrate what I am talking about. If we have used a delta of 1, this
8 chart would tell us that any “true discrepancy” below 0.5 is immaterial and any
9 “true discrepancy” above 0.5 is material. The chart shows the probability of
10 detecting this condition. Using an example from the chart, assume a very
11 small sample size, which is always going to be problematic. In the first line,
12 even if the “true disparity” was zero, that is there was no disparity, the
13 statistical analysis is going to show that there is disparity 32 percent of the
14 time. On the other end of the scale, at 1, the analysis is only going to show a
15 material difference 68 percent of the time, when we know that the disparity
16 actually exists and is material. These are essentially examples of Type 1 and
17 Type II errors, where the Type II error at the 1 disparity level is 32 percent (the
18 complement of the probability of detection). Importantly, as the sample size
19 increases, the analysis rapidly approaches an accuracy level of 100 percent,
20 meaning that the Type I and Type II errors are essentially eliminated.
21

**Figure 2: The Probability of Detecting Disparity
Mean Measure Test with Delta = 1**

BST Sample Size	ALEC Sample Size	Balancing Critical Value	True Disparity Level						
			0	0.2	0.45	0.5	.55	0.8	1
10	1	-0.477	0.317	0.387	0.481	0.5	0.519	0.613	0.683
100	5	-1.091	0.138	0.256	0.457	0.5	0.543	0.744	0.862
1000	50	-3.45	0	0.019	0.365	0.5	0.635	0.981	1
12000	800	-13.693	0	0	0.085	0.5	0.915	1	1
100000	2500	-24.693	0	0	0.007	0.5	0.993	1	1

Q. IT SEEMS THEN THAT A MEAN MEASURE TEST BASED ON A
BALANCING METHODOLOGY DOES MAKE IT POSSIBLE TO DETECT
DISCRIMINATION AS LONG AS THE TRUE DISPARITY IS BEYOND
THE MATERIALITY THRESHOLD. IS THAT TRUE?

A. Yes, a mean measure test based on balancing and large sample sizes has a high
likelihood of detecting disparity beyond the materiality threshold, but a low
probability of detecting disparity that falls under the threshold.

Q. ISN'T IT TRUE THAT THESE CONDITIONS ARE THE SAME ONES
THAT LEAD TO BALANCING CRITICAL VALUES THAT ARE
FURTHER FROM ZERO THAN THOSE THAT ARE CONVENTIONALLY
USED?

A. Yes. Large sample sizes lead to critical values that are further from zero than
those that are used in many applications. Such critical values, in turn, lead to
small significance levels. But, as I have shown, those small significance levels
(which are the probabilities corresponding to a true disparity of 0 in Figure 2)

1 do not imply that BellSouth will get away with any amount of discrimination.
2 Those levels of disparity that are lower than the materiality threshold, which is
3 defined by the choice of delta, will not be considered discriminatory.
4 However, levels of disparity beyond the materiality threshold will be detected
5 as discriminatory with a high likelihood.

6

7 Q. IS THE SAME THING TRUE FOR PROPORTION MEASURES?

8

9 A. A similar statement can be made for a proportion measure test. When using an
10 odds ratio approach to balancing, the materiality threshold is not one-half of
11 the odds ratio used in the balancing alternative hypothesis, but the threshold is
12 at a point close to this. Figure 3 below illustrates this by showing the
13 probability that the testing procedure will determine disparity (reject the null
14 hypothesis), for a range of disparity levels and BST/ALEC sample sizes when
15 the BellSouth proportion of missed installations is 0.20 and balancing is done
16 for the alternative hypothesis with an odds ratio of 3.

17

18 Notice that for a balancing alternative with odds ratio of 3 (BST proportion of
19 0.20 and CLEC proportion of 0.43), there is a significant probability of
20 determining disparity for odds ratio levels less than 3. For example, with a
21 CLEC proportion of misses of 0.30 there is at least a 50% chance, regardless of
22 sample size, that disparity will be determined and a remedy paid. Here we
23 have an odds ratio of 1.75, much less than the balancing alternative of 3.

24

**Figure 3: The Probability Of Determining Disparity
When the BellSouth Proportion of Missed Installations is 0.20 and
the Balancing Critical Value is Determined at an Odds Ratio of 3**

Number of Transactions		Level of Disparity in Terms of Odds Ratio						
		<i>Level of Disparity in Terms of CLEC Proportion</i>						
		1*	1.25	1.75	2	2.25	2.75	3**
BST	ALEC	0.20	0.24	0.30	0.33	0.36	0.41	0.43
10	1	0.4110	0.4440	0.5000	0.5220	0.5410	0.5750	0.5890
100	5	0.2920	0.3730	0.5040	0.5570	0.6030	0.6790	0.7080
1000	50	0.0410	0.1530	0.5130	0.6750	0.7960	0.9300	0.9590
12000	800	0.0000	0.0000	0.5520	0.9640	0.9990	1.0000	1.0000
100000	2500	0.0000	0.0000	0.5930	0.9990	1.0000	1.0000	1.0000

Q PLEASE RECAP YOUR POINT REGARDING DR. FORD'S TESTIMONY
THAT YOU HAVE BEEN DISCUSSING.

A. Dr. Ford seems to believe that low significance levels means that actual and
material disparities will not be discovered, particularly with large sample sizes.
That is simply not true, as I have demonstrated above.

Q. LET'S MOVE ON TO THE SECOND ISSUE YOU BELIEVE DR. FORD IS
CONFUSED ABOUT. CAN YOU DESCRIBE YOUR POINT WITH MORE
SPECIFICITY?

A. Dr. Ford seems concerned about the large critical values that can result from
the analysis that is proposed in BellSouth's plan. He believes that some sort of
"standard analysis" would preclude the use of significance levels below one
percent. For example, on page 20 of his direct testimony, lines 11 – 13, he
states, "Recall that standard significance levels of a means-difference test are

* An odds ratio of one assumes that there is parity. Thus, the probability of determining disparity in this situation is the probability of a Type I error.

** The probability of determining disparity increases as the level of disparity goes beyond an odds ration of three.

1 5%, or in some cases as low as 1%. A 1% significance level is considered
2 quite small. Rarely are significance levels chosen below this value.”

3 Basically, he is suggesting that large critical values in and of themselves
4 suggest some sort of problem and that there ought to be a floor on critical
5 values to eliminate any such problems.

6

7 The problem with appealing to the “standard,” or “conventional” testing
8 approach that is described by most introductory statistical textbooks, and even
9 more advanced textbooks, is that there is almost always an assumption that the
10 data in a study are collected according to a designed plan and that there is more
11 than ample time to evaluate critically the data that is being used. In the
12 simplest of cases, the assumption is that a simple random sample has been
13 collected. In more complex cases, such as agricultural experiments or clinical
14 trials, the sampling plans call for collecting data in specific ways. In all these
15 cases, the sample size of the data collected is usually under the control of the
16 data collector.

17

18 Most statistical textbooks also warn users of statistics to think about the results
19 that they are observing. Just because a test results in a statistically significant
20 difference between two means or proportions, one should also make sure that
21 the observed difference makes sense from a practical point of view. This is
22 especially true when sample sizes are very large. In these cases, Z statistics
23 may have a large magnitude even when the actual difference between the
24 performance measures is quite small.

25

1 Q. WHY ARE THESE POINTS IMPORTANT IN THE CONTEXT OF THIS
2 PROCEEDING?

3

4 A. There are two reasons. First, the performance assessment plans that we are
5 dealing with involve observational studies. This is a process where the
6 subjects select themselves into one of the groups that are being compared. In
7 our case customers select the telephone company that they want. We have
8 very little control over this, and unlike the situations that textbooks usually
9 cover; we have no control over the sample sizes that will be used every month.

10

11 Q. WHAT IS THE SECOND ISSUE?

12

13 A. The analysis of this data must be completed in a short amount of time, for
14 many measures, every month. Normally, a good statistician would explore the
15 data, and try to answer many questions about the data. This is particularly true
16 when seemingly large Z values are calculated, which seems to be Dr. Ford's
17 concern. That is, normally you should try to discover why such large Z values
18 occurred. Was it due to a large discrepancy in the performance measure? Or,
19 maybe it is the case that, from a practical point of view, there is very little
20 difference in performance and the large Z value was simply caused by large
21 sample sizes.

22

23 Q. DOES THE FACT THAT THESE PLANS REQUIRE VERY SPEEDY
24 REPORTING AND PROVIDE ALMOST NO TIME FOR ANY CRITICAL
25 ANALYSIS MEAN THAT THE STATISTICAL METHOD PROPOSED IS

1 SIMPLY INAPPROPRIATE TO USE?

2

3 A. Absolutely not. Indeed, the speed with which the data is to be reported and
4 penalties paid if owed is one of the reasons why the Ernst & Young statistical
5 team felt a balancing method was valuable. Large Z values that go beyond a
6 balancing critical value are most likely caused by truly disparate treatment.
7 But Z values that don't go beyond the balancing critical value are immaterial in
8 terms of the difference in performance.

9

10 Q. IS THERE ANYTHING ELSE IN DR. FORD'S TESTIMONY THAT YOU
11 FEEL IS IMPORTANT TO DISCUSS?

12

13 A. Yes. I think it is important to discuss the opposite of the small significance
14 level issue that Dr. Ford raises. That is, the use of significance levels that are
15 much larger than what is conventionally used when sample sizes are small. I
16 would also like to discuss a graph in Dr. Ford's testimony that is very
17 misleading.

18

19 Q. WHAT HAPPENS TO THE SIGNIFICANCE LEVEL OF A BALANCED
20 STATISTICAL TEST WHEN SAMPLES ARE SMALL?

21

22 A. The significance levels can be 3, 5, or even up to almost 10 times larger than a
23 conventional value of 5 percent. For example, with a BellSouth sample of
24 1000, an ALEC sample of 30, and a "delta" value of 0.25, the balancing critical
25 value of a mean measure test is -0.675. This gives a significance level for the

1 test of about 25 percent. This means that BellSouth would be found to be out
2 of parity 25 percent of the time.

3

4 Q. DOES THIS MEAN THAT YOU HAVE AN OBJECTION TO BALANCING
5 FOR SMALL SAMPLES?

6

7 A. No. This is what balancing is supposed to do. When sample sizes are small it
8 gives the benefit of the doubt to the ALEC. On the flip side, the data must
9 show that there is a material difference, not just a conventionally significant
10 difference, in the performance measure when the sample sizes are large.

11

12 Q. ISN'T IT MORE LIKELY THAT SAMPLE SIZES WILL BE LARGE?

13

14 A. On the contrary, in the performance measure data that I have looked at sample
15 sizes tend to be small enough in such areas as UNE services and other special
16 types of services that the balancing critical value of a Tier I test tends to be
17 between 0 and -1. In the example I give above, samples of 1000 (BellSouth)
18 and 30 (ALEC) lead to a balancing critical value of -0.675 for a "delta" of
19 0.25, and -1.35 for a "delta" value of 0.5. While a sample of size 30 for the
20 ALEC is not huge, many would not consider it to be overly small.

21

22 Q. DOES DR. FORD RECOGNIZE THIS FACT ABOUT LARGE
23 SIGNIFICANCE LEVELS FOR SMALL SAMPLE SIZES?

24

1 A. I believe he does, since he suggests using a “delta function” to choose “delta”
2 based on the ALEC sample size. But I do not believe that this is the correct
3 concept. Balancing error probabilities is not about searching for critical values
4 that in some sense makes the two sides happy. When one adopts a balancing
5 approach it has to be understood that you are really trying to determine what
6 type of difference in performance truly has a material impact on an ALEC’s
7 business.

8
9 Q. YOU HAVEN’T MENTIONED HOW DR. BELL FEELS ABOUT THE
10 EFFECTS OF BALANCING OF THE SIGNIFICANCE LEVEL OF THE
11 TEST. DOES HE THINK THERE NEEDS TO BE A “FIX” FOR THE
12 METHOD?

13
14 A. In discussing large negative Z scores that do not trigger a test failure because
15 the balancing critical value is larger (further from zero than the Z score), Dr.
16 Bell states on page 14, lines 16-17, “Such an outcome would be justified only
17 if one could be certain that delta has not been set too large.” He goes on to say
18 that he feels no floor is warranted if the “delta” he advocates, 0.25, is used.

19
20 From this statement, I infer that Dr. Bell understands that a balanced test has
21 sufficient power to detect truly discriminatory performance on the part of
22 BellSouth. However, this will only be true if “delta” is chosen so that it
23 effectively defines the materiality threshold.

24
25

1 Q. DO YOU AGREE WITH DR. BELL?

2

3 A. In principal yes. I am not convinced, however, that a "delta" of 0.25 is correct.
4 The Louisiana Public Service Commission has ordered BellSouth to use a
5 "delta" of 1. The Georgia Public Service Commission has ordered that a
6 "delta" of 0.5 be used. In both situations, there will be periodic reviews of the
7 effectiveness of the methodology. I assume that if these commissions find that
8 "delta" was set too large, they will lower the value. It's also possible that a
9 review will find that the values are too low. Only time will tell.

10

11 Q. YOU STATED THAT DR. FORD HAS INCLUDED A GRAPH IN HIS
12 DIRECT TESTIMONY THAT IS MISSEADING. PLEASE EXPLAIN
13 THIS TO US.

14

15 A. Exhibit No. ___(GSF-3) of Dr. Ford's direct testimony is supposed to be a
16 graph that shows the alternative distribution with different "delta" values. Dr.
17 Ford does not identify the exact distribution he is using, but based on the bell-
18 shapes he uses, and the language in his testimony, I assume that he is using a
19 normal distribution. Given that, there is no way his graph illustrates
20 distributions that are shifted 0.25, 0.5 and 1 standard deviations from the
21 BellSouth distribution.

22

23 Q. HOW CAN YOU TELL THAT?

24

1 A. The normal distribution has certain properties about it that indicate to you the
2 size of its standard deviation based on the spread of the bell-curve. Figure 4
3 below illustrates this.

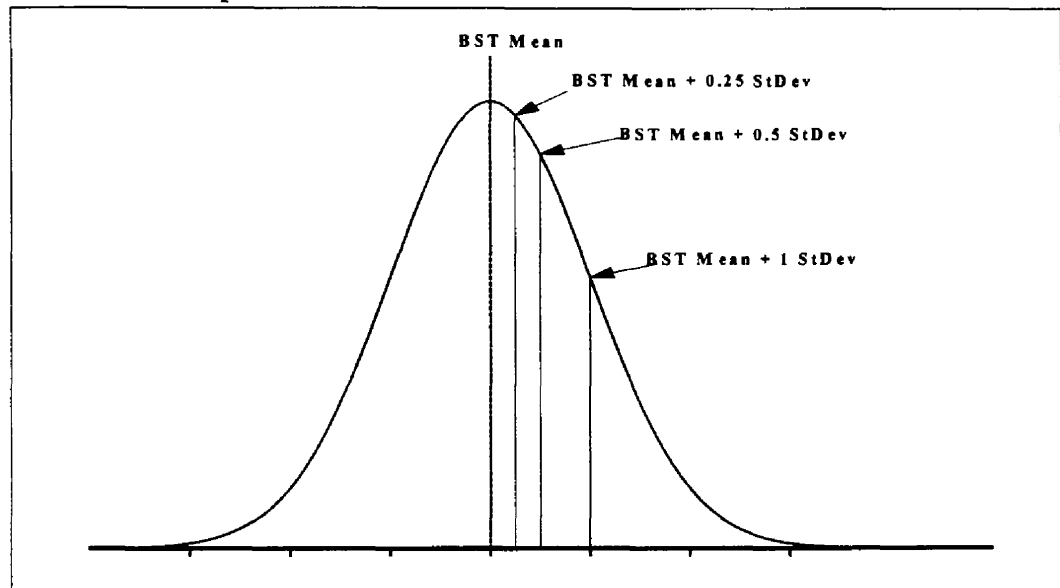
4

5

6

7

FIGURE 4: The Normal Distribution
Illustration of the Relationship Between
The Spread of the Bell-Curve and the Standard Deviation



8

9

10 Figure 4 shows the location of the points that are 0.25, 0.5 and 1 standard
11 deviations (StDev) from the mean of the distribution (BST Mean is at the
12 center of the bell-curve). We can also look up the area under a normal bell-
13 curve to the left of each of these values. These areas are approximately 60, 70,
14 and 84 percent of the total area under the curve for the points 0.25, 0.5 and 1
15 standard deviation from the mean, respectively. A visual inspection of Figure
16 4 will indicate that this graph exhibits these area features.

17

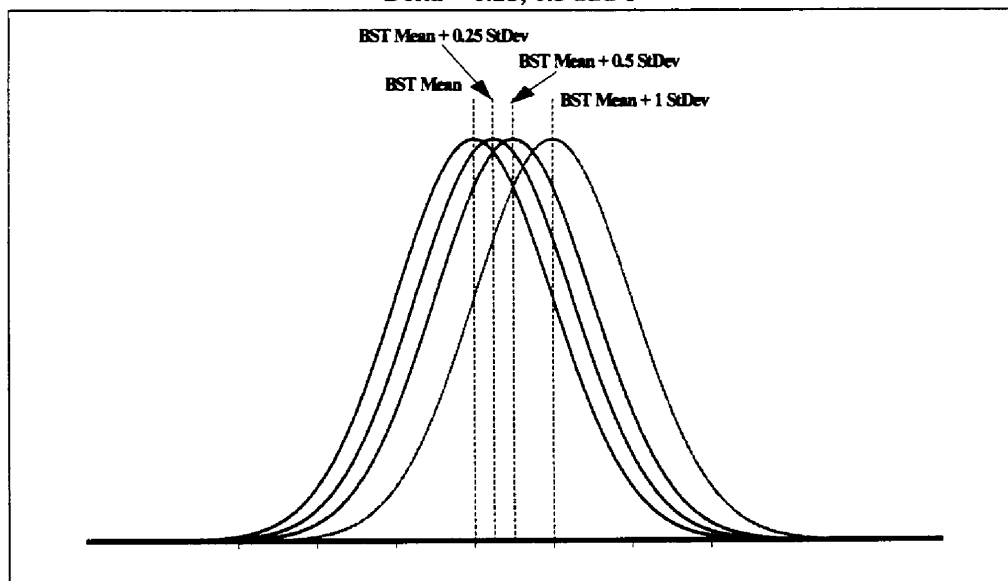
18 Looking at Dr. Ford's graph in Exhibit No.__(GSF-3), he does not place these
19 points correctly on his graph. The point where he places the mean plus 0.25

standard deviation appears to really be about 2 standard deviations from the mean. I can only guess that the point that is supposed to be 1 standard deviation from the mean is located about 8 standard deviations from the mean.

Q. WHAT SHOULD THE CONCEPT DR. FORD IS ATTEMPTING TO ILLUSTRATE REALLY LOOK LIKE?

A. Figure 5 shows 4 bell-curves. The first one on the left represents the BellSouth service time distribution. The second one represents the alternative hypothesis distribution for an ALEC that has the same standard deviation as the BellSouth distribution, but its mean is larger than BellSouth's by 0.25 standard deviations. The third and fourth bell-curve are similar, representing ALEC means that are 0.5 and 1 standard deviations larger than the BellSouth means.

**FIGURE 5: Location of Alternative Normal Distributions
With Respect to the BellSouth Distribution**
Delta = 0.25, 0.5 and 1



1 Q. THIS CERTAINLY GIVES A MUCH DIFFERENT VISUAL
2 REPRESENTATION THAN DR. FORD'S GRAPH. WHY IS DR. FORD'S
3 GRAPH SO DIFFERENT?

4
5 A. I am not sure. Perhaps he is not using a normal distribution. But his curves are
6 symmetric bell-shapes, and while there are other distributions with similar
7 shapes, the relationship between the curve and the point that is one standard
8 deviation from the mean is not that much different from where it is located
9 based on the normal distribution. I can only conclude that Dr. Ford either
10 doesn't understand what he is doing, or he is deliberately trying to be
11 misleading.

12
13 Q. DOES THIS CONCLUDE YOUR REBUTTAL TESTIMONY?

14
15 A. Yes.