

Pareto Charts Plus

Pareto Charts,
improved Pareto Charts, and
intercomparison of percentages

Giles Crane, MPH, ASA, NJPHA

New Jersey R User Group,
Sponsored by Mango Solutions,
also Central Jersey R User Group
October 22, 2013



References

- *Vilfredo Pareto (1896-7) .Cours d'économie politique professé à l'université de Lausanne, 3 volumes.*
- Juran, J.M. (1951). "The Economics of Quality" in *Quality Control Handbook*. ed. J.M. Juran, New York, McGraw-Hill.
- Leland Wilkinson (2006). Revising the Pareto Chart. *The American Statistician*, November, 60:4:332-334.

The "7 jewels in the Crown" (George Box), applied simply, well, and consistently, provide realistic quality control and improvement.

- ❑ Pareto Charts
- ❑ Cause & Effect Diagrams
- ❑ Control charts
- ❑ Check lists
- ❑ Scatter diagrams
- ❑ Histograms
- ❑ Graphs

Example 1: Patient falls in U. of Texas Health Science Center, Houston

71 reports of patient falls at the hospital

15 Attributed causes with frequencies

9 9 9 7 6 5 4 4 3 3 3 3 3 2 1

Example 2: Illinois Asylum circa 1860

Probable Causes of Insanity (play: Mrs. Packard).

<input type="checkbox"/> Miasmatic fevers	25	<input type="checkbox"/> Loss of wife	2
<input type="checkbox"/> Indigestion	14	<input type="checkbox"/> Excessive blood loss	2
<input type="checkbox"/> Religious anxiety	11	<input type="checkbox"/> Sudden cessation of bleeding from the nose	1
<input type="checkbox"/> Loss of property	12		1
<input type="checkbox"/> Epilepsy	12	<input type="checkbox"/> Paralysis	1
<input type="checkbox"/> Disappointed love	10	<input type="checkbox"/> Concussion from steamboat explosion	1
<input type="checkbox"/> Puerperal fever	8	<input type="checkbox"/> Jealousy	1
<input type="checkbox"/> Masturbation	8	<input type="checkbox"/> Shipwreck	1
<input type="checkbox"/> Domestic unhappiness	4	<input type="checkbox"/> Fear of mob	1
<input type="checkbox"/> Spiritual rappings	4	<input type="checkbox"/> Loss of husband	1
<input type="checkbox"/> Intemperance	4	<input type="checkbox"/> Jaundice	1
<input type="checkbox"/> Intense study	4	<input type="checkbox"/> Scurvy	1
<input type="checkbox"/> Unkind relatives	4	<input type="checkbox"/> Slander	1
<input type="checkbox"/> Loss of children	3		5
<input type="checkbox"/> Blows on the head	2	<input type="checkbox"/> Unknown	56

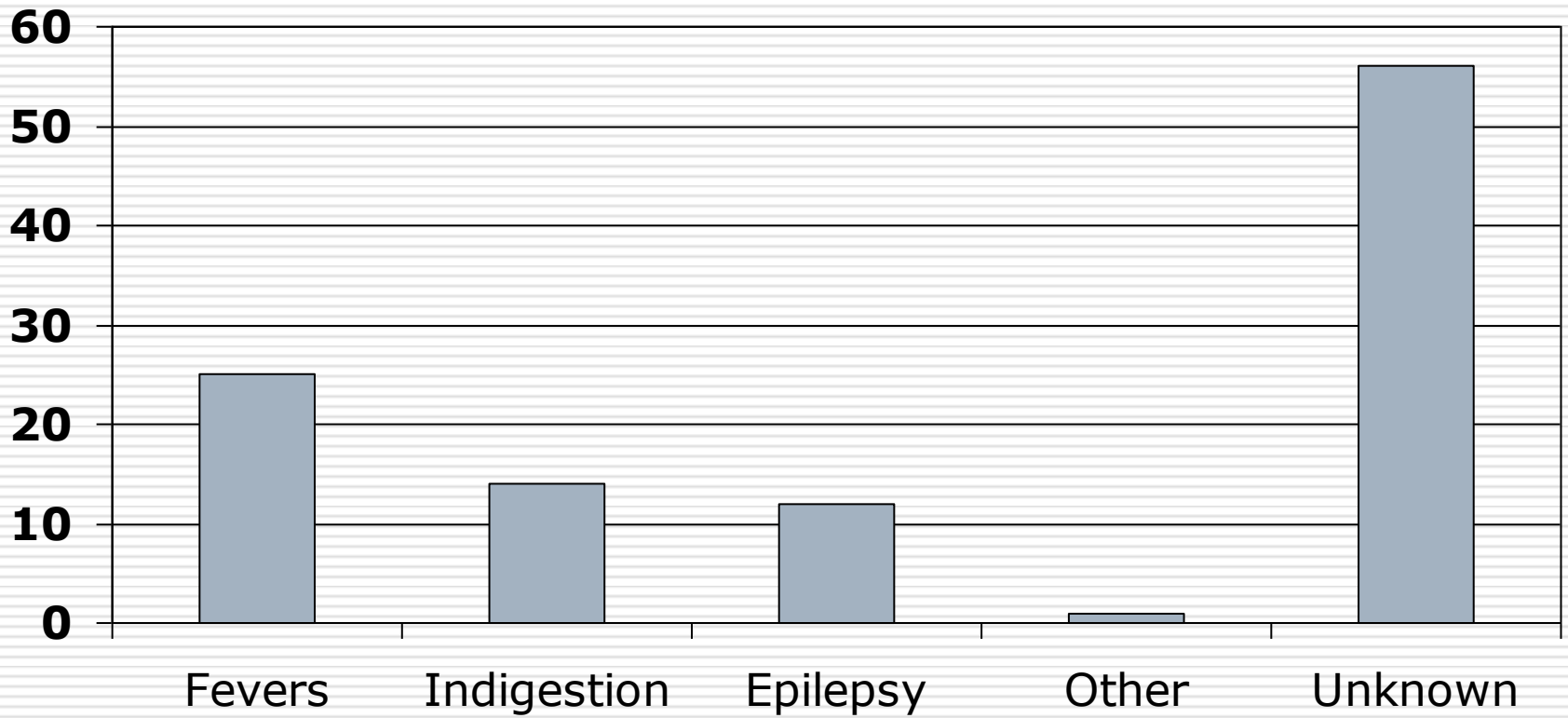
Jacksonville State Hospital Illinois



EXAMPLE 3: Survey of HIV positive women: reasons for not seeking appropriate gynecological Care (NJ).

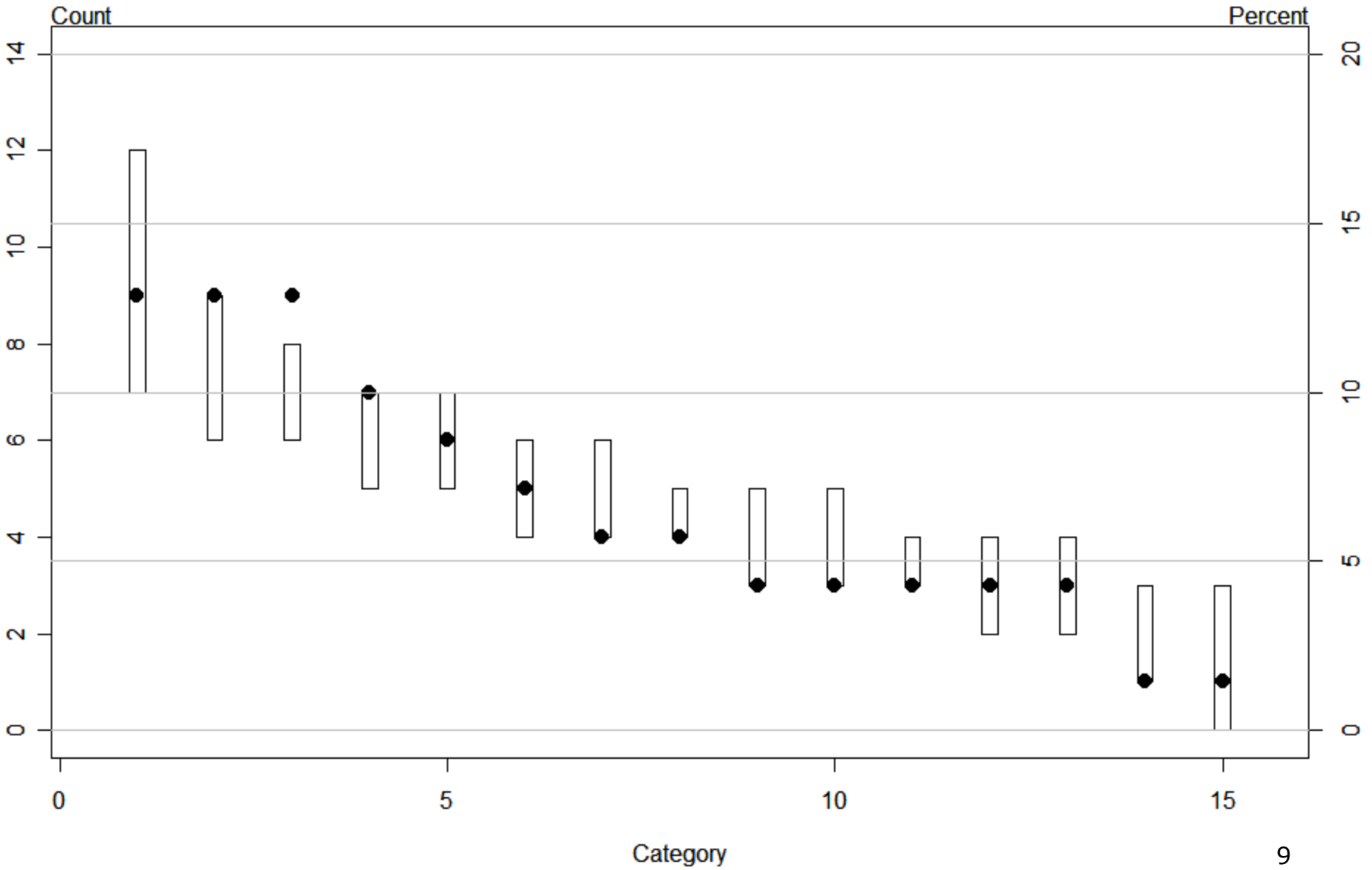
<input type="checkbox"/> Fear of results	42
<input type="checkbox"/> Procedural Discomfort	39
<input type="checkbox"/> Cost	18
<input type="checkbox"/> Competing Priorities	12
<input type="checkbox"/> Appointment Time	10
<input type="checkbox"/> Dislike Provider	6
<input type="checkbox"/> Disclosure	2

A Pareto Diagram is basically a bar chart of counts, arranged in decreasing order, but see the statistical limits in the next slides.

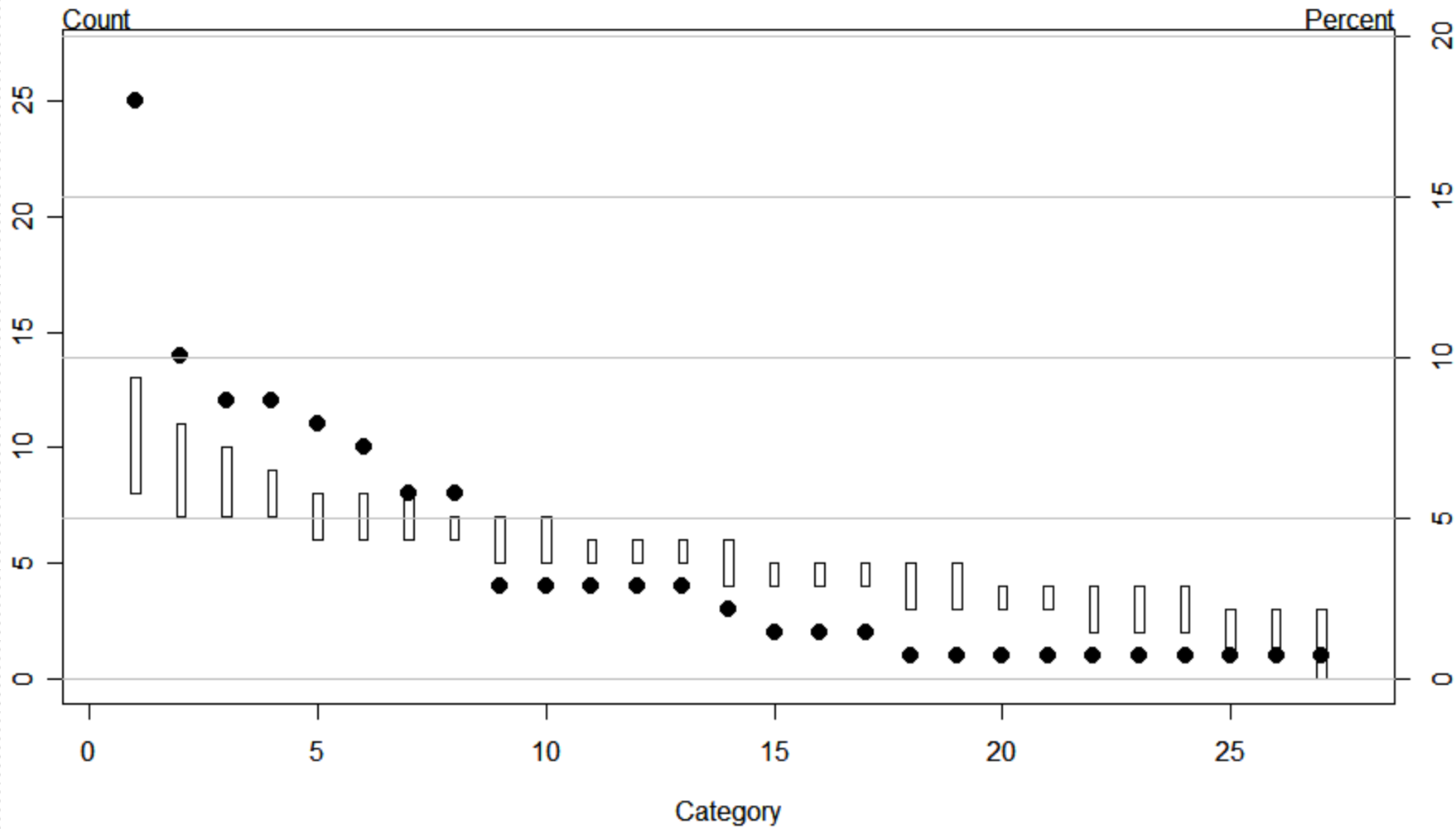


(Example 2) Causes of Insanity, abbreviated table.

15 Causes of Falls in a Hospital

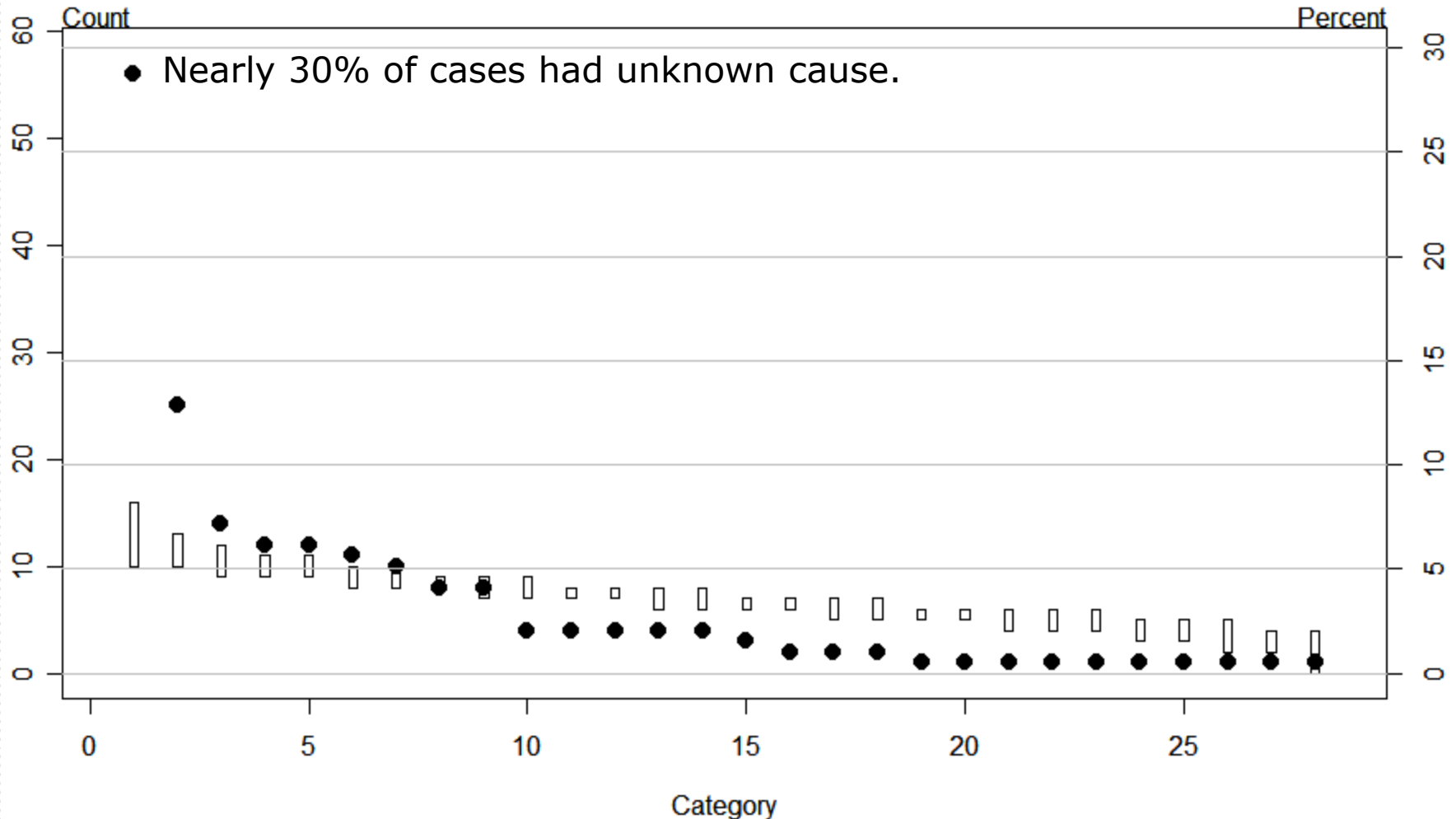


Probable Causes of Insanity

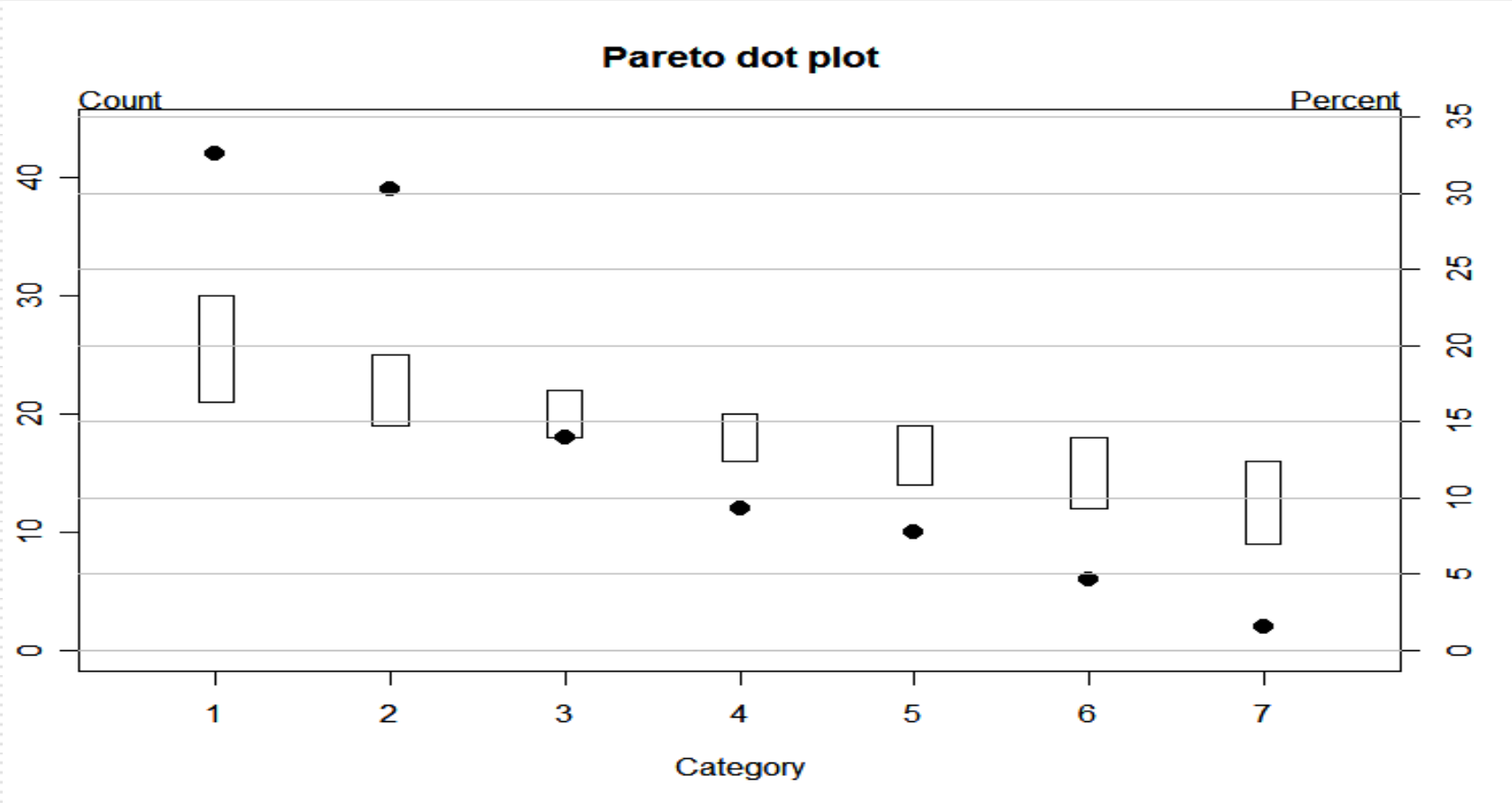


Including the unknown cause category, there are 28 categories.

Probable Causes of Insanity



With R, we can calculate limits indicating that counts within them may well be due to chance, not a real difference.



(Example 1) Reasons of HIV Positive women for not seeking gyn. care.

Statistical limits stimulate realistic considerations about the data.

- Do we have enough data?
- Difficulty of remedies for each cause?
- Cost for remedies?
- Side effects of remedies?
- Can certain causes be grouped?

In situations involving 3 or more counts, the Multinomial Distribution provides a theoretical background.

□ $P(n_1, n_2, \dots, n_k \mid \text{given } p_1, p_2, \dots, p_k) = \frac{N!}{(n_1! n_2! \dots n_k!)} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$
where $N = \text{sum of } n_i$

□ Mean & variance of each n_i are

$$\text{mean} = N * p_i$$

$$\text{var} = N * p_i * (1 - p_i)$$

But covariance n_i and $n_j = -N p_i p_j$





Acceptance Limits on multinomial counts in Pareto charts can be calculated by a computer simulation.

- Sample from multinomial distribution using number of categories, total count & equal probabilities, then sort counts, and save as a column.
- Do this several thousand times.
- For each row (thousands of columns), take the 2.5% count and the 97.5% count (quantiles), then plot as the lower and upper assurance limits.

Computing acceptance limits can be done by simulation in R. file `pardot.r`

Obtain lower and upper acceptance limits on the ordered categories.

```
parlim <- function(vn, vp=c(0.025,0.975), nsample=2000) {  
  
  ncat = length(vn)  
  nsize = sum(vn)  
  
  m = rmultinom(nsample, size = nsize, prob=rep(1,ncat) )  
  
  dsort <-function(v) sort(v, decreasing=TRUE)  
  m = apply(m, 2, dsort) # sort each column (Re-use m name)
```

Continue by taking quantiles and appending them to a data frame xlu. file **pardot.r**

Obtain lower and upper assurance limits on the ordered categories.

```
xlu = data.frame()      # Create buckets
```

```
for (i in 1:ncat) {  
  limits = quantile(m[i,],vp)  
  xlu = rbind(xlu, as.vector(limits))  
}
```

```
row.names(xlu) = NULL  # Renumber rows 1,2,..#categories  
xlu = data.frame(vn,xlu) # Put original counts as first column.  
return(xlu)  
} #End of parlim function
```

We can now make a similar plot to that in Wilkinson's paper. File: **pardot.r**

```
pardot <- function(xlu, catlabs,
                   heading="Pareto Dot Chart") {
  ncat = nrow(xlu)
  nsize = sum(xlu[,1])

  opar = par(mar=c(5,2,4,3) + 0.1 )
  plot(xlu[,1], main = heading,
       xlab = "Category",
       ylab = "",
       xlim = c(0.5, ncat + 0.5),
       ylim = c(0, max(xlu[1,3],xlu[1,1]) + 2),

  pch = 19, cex=1.5)
```

Boxes are drawn around the points to dramatize the limits. File: **pardot.r**

```
# Assurance limits boxes
hw = 0.10 # Specify Half width for boxes
rect(1:ncat-hw, xlu[,2],1:ncat+hw,xlu[,3],
     border="blue")

ticks = axis(4,at=5*0:20*nsiz/100,
             labels=as.character(5*0:20))
abline(h=5*0:20*nsiz/100,col="gray")

mtext("Count", side=3,line=0, adj=0)
mtext("Percent", side=3,line=0, adj=1)

par(opar)
}
```

Here is our NJ data example:

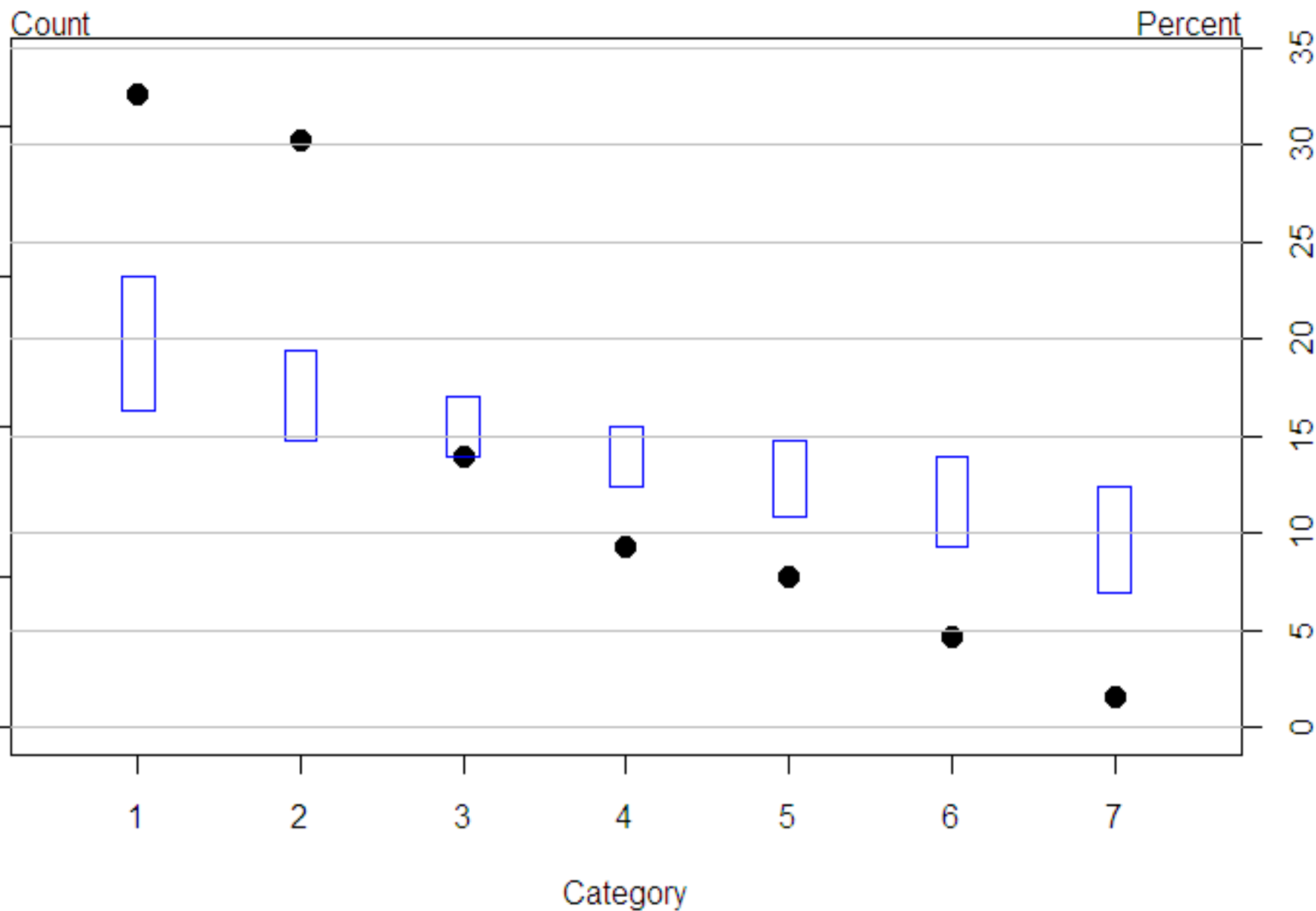
```
head1 = "Reasons of HIV Positive  
women for not seeking gyn. care."
```

```
vn=c(42,39,18,12,10,6,2)
```

```
xlu = parlim(vn)
```

```
pardot(xlu, ,head1)
```

Reasons of HIV Positive women for not seeking gyn. care.



The classical chi-square test also gives an overall test of equality of proportions:
Are any p_i different from any other?

- Chi-square Test: Consider Probability of outcome n_1, n_2, \dots, n_k , given $p_1 = p_2 = \dots = p_k$, i.e. equal probabilities.

For all k categories, Expected = $n = N/k$ with equal probability.

Then

$X = \text{Sum of } (n_i - \text{Expected})^2 / \text{Expected}.$

is distributed as chi-square [sum of k $N(0,1)$ variables]

If in tail of chi-square distribution, this is a rare event when all p_i are equal.

Hence our sample may be from events of unequal p_i .

Thus there are several, statistical approaches to answer: Are the percentages really different?

Omnibus Chi-square statistical test.

Alan Agresti's exact methods. (counts more extreme sets of count)

Simulation, including Pareto dot chart.

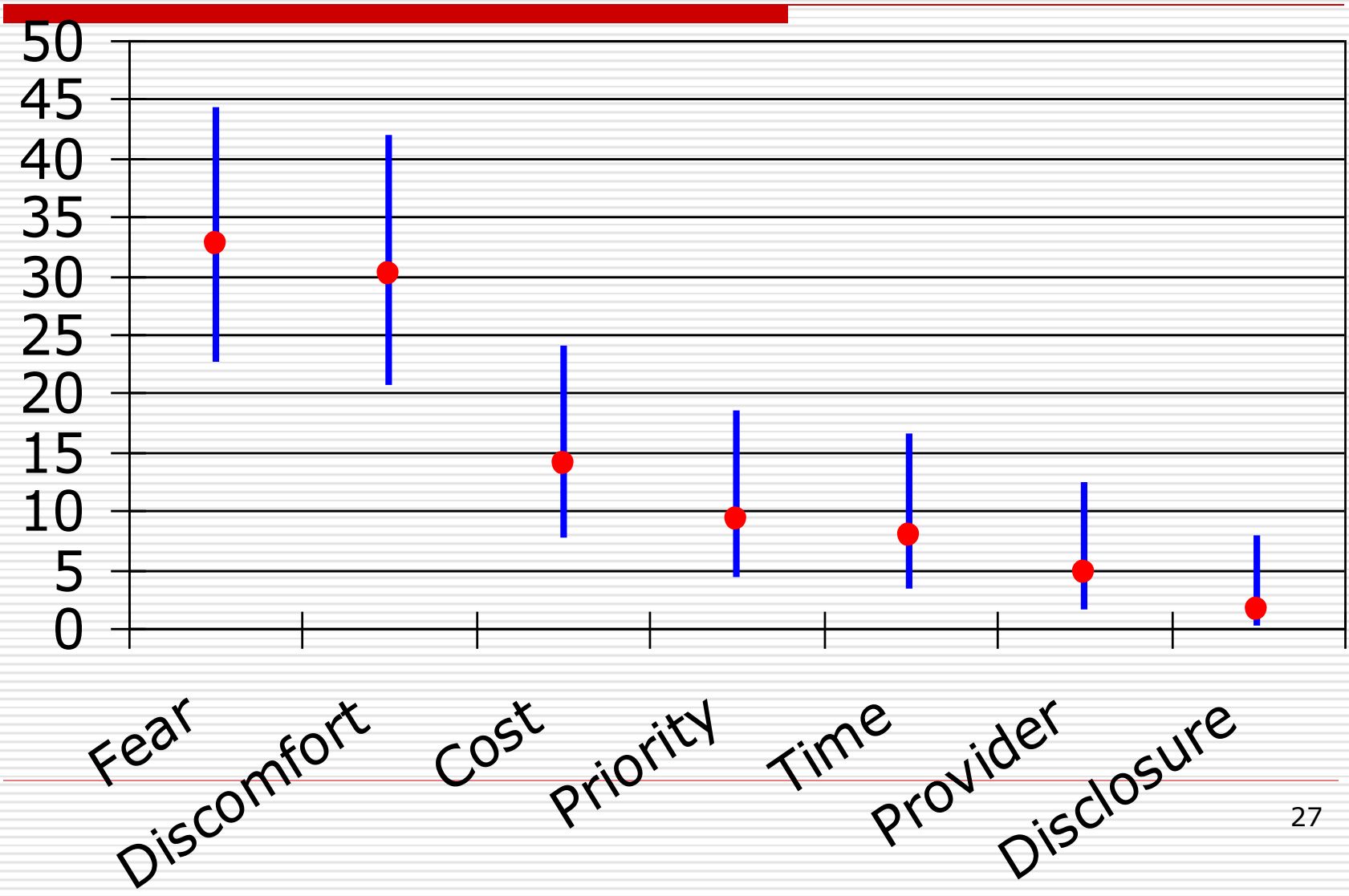
Inter-comparison of percentages.

Here are simultaneous confidence limits on the patient percentages for each reason HIV pos.women avoid proper care.

Estimate	Description	LCL95	UCL95
p1: 0.32558	Fear of results	0.226451	0.44324
p2: 0.30233	Procedural Discomfort	0.206457	0.41919
p3: 0.13953	Cost	0.076559	0.24081
p4: 0.09302	Competing Priorities	0.044293	0.18499
p5: 0.07752	Appointment Time	0.034370	0.16555
p6: 0.04651	Dislike Provider	0.016416	0.12478
p7: 0.01550	Disclosure	0.002859	0.07962

Fractions rather than percentages are shown in this table.

Statistical evidence for real differences can be seen in a plot of confidence intervals.



We computed simultaneous confidence limits in **R**. (multici.r)

Source of multici.r() Program: Mayo Foundation for Medical Education and Research (Ross Dierkhising, author)

Easily converted code from Splus to **R**

For one multinomial population or several multinomial populations.

Computes confidence limits on percentages, pairwise differences and contrasts.