

Runs Test Paper

Lindesay Scott-Hayward

2016-07-19

Contents

Simulated Data	1
Uncorrelated data	1
Correlated data	5
What does all this mean for power analysis?	6
50% decline, $\phi = 50$	6
25% decline, $\phi = 50$	8
0% decline, $\phi = 50$	8
50% decline, $\phi = 50$, multiruns	8
25% decline, $\phi = 50$, multiruns	9
Real Data: Falls of Warness	9

Simulated Data

Make some simulated toy data using the following equation:

$$y = \beta_0 + \beta_1 X_1$$

where $\beta_0 = 1$ and $\beta_1 = 0.3$

```
dat<-makeToyData(200, length.panels=5)
head(dat)
```

	x	evph	mu	panelid
1	1.000000	0	3.669297	1
2	1.090909	0	3.770746	1
3	1.181818	0	3.874999	1
4	1.272727	0	3.982135	1
5	1.363636	0	4.092234	1
6	1.454545	0	4.205376	2

Uncorrelated data

1000 sets of noisy data are simulated from this truth using a Poisson distribution.

```
newdat<-generateNoise(nsim, dat$mu, family='poisson', d=1)
```

To test, fit a glm to one of the sets of data.

```
init_glm<-glm(newdat[,1] ~ x, data=dat, family='quasipoisson')
summary(init_glm)
```

```
Call:
glm(formula = newdat[, 1] ~ x, family = "quasipoisson", data = dat)
```

```
Deviance Residuals:
```

Min	1Q	Median	3Q	Max
-2.69510	-0.74512	0.00829	0.50767	2.69483

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.974638	0.053994	18.05	<2e-16 ***
x	0.306862	0.007009	43.78	<2e-16 ***

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(Dispersion parameter for quasipoisson family taken to be 0.9277439)
```

```
Null deviance: 2301.86  on 199  degrees of freedom
Residual deviance: 186.32  on 198  degrees of freedom
AIC: NA
```

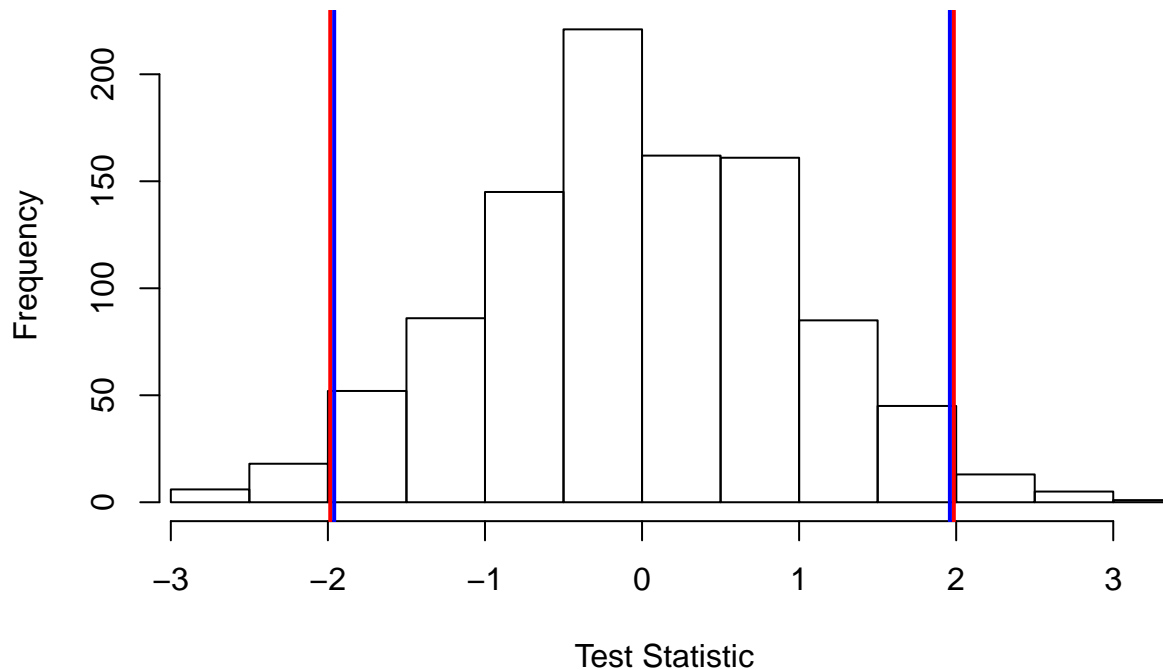
```
Number of Fisher Scoring iterations: 4
```

Models were fitted to all datasets generated above and the runs test was evaluated for each one. The distribution of the test statistics is returned. A plot is also produced, which shows the lower 2.5% and upper 97.5% critical values of the empirical distribution (red) and from the Normal ($N(0, 1)$) distribution.

```
empdistribution<-getRunsCritVals(n.sim = nsim, simData=newdat,
                                model = init_glm, data = dat, plot=TRUE,
                                returnDist = TRUE)
```

```
.....
.....
```

Empirical Distribution: Runs Test Statistic



Use the data generated to assess the 5% error rate for this test when using the Normal distribution or the empirical distribution generated above.

```
newdat<-generateNoise(nsim, dat$mu, family='poisson', d=1)
ps<-matrix(NA, nrow=nsim, ncol=2)

for(i in 1:nsim){
  sim_glm<-glm(newdat[,i] ~ x , data=dat, family='quasipoisson')
  # find both the empirical and Normal p-values
  ps[i,1]<-runs.test(residuals(sim_glm, type='pearson'), critvals = empdistribution)$p.value
  ps[i,2]<-runs.test(residuals(sim_glm, type='pearson'))$p.value
}
```

```
(length(which(ps[,1]<0.05))/nsim)*100
```

```
[1] 6.9
```

```
(length(which(ps[,2]<0.05))/nsim)*100
```

```
[1] 6.9
```

How does this change if we vary the dispersion parameter for the data?

```
nphi=seq(1,51, by=5)
errrate<-matrix(NA, nrow=length(nphi), ncol=2)
counter=1
nsim=2000
for(p in nphi){
  newdat<-generateNoise(nsim, dat$mu, family='poisson', d=p)
```

```

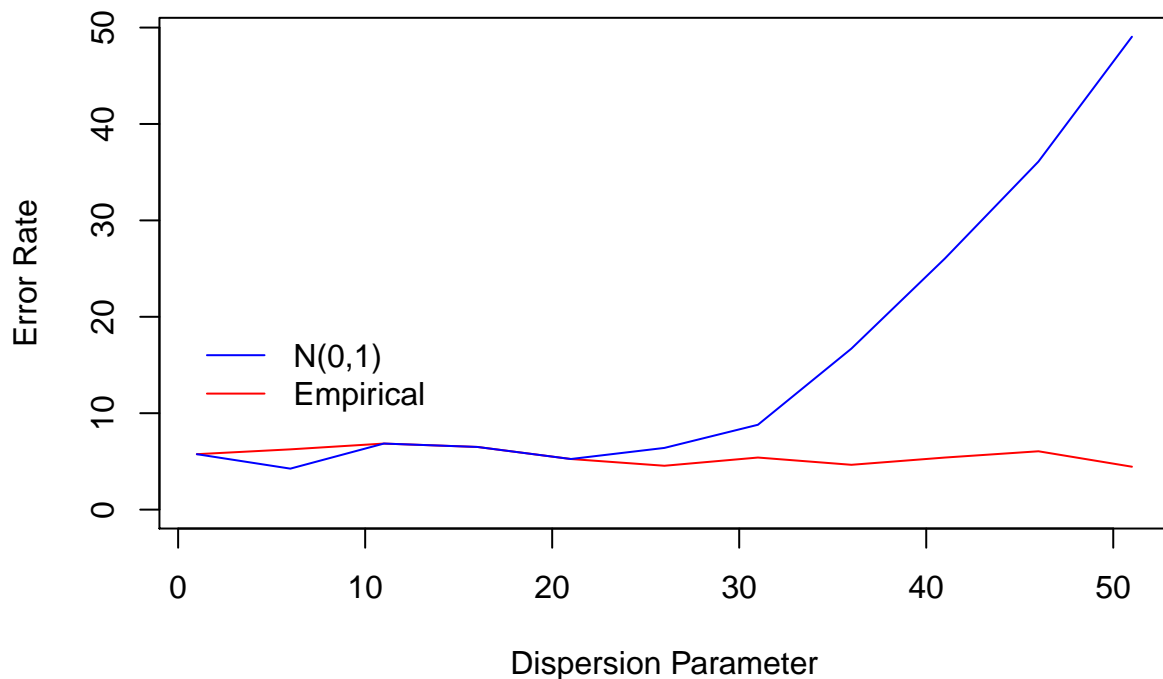
# for each change in phi, update the empirical distribution
empdistribution<-getRunsCritVals(n.sim = nsim, simData=newdat,
                                model = init_glm, data = dat, plot=FALSE,
                                returnDist = TRUE, dots=FALSE)

ps<-matrix(NA, nrow=nsim, ncol=2)
for(i in 1:nsim){
  sim_glm<-glm(newdat[,i] ~ x + as.factor(evph), data=dat, family='quasipoisson')
  # find both the empirical and Normal p-values
  ps[i,1]<-runs.test(residuals(sim_glm, type='pearson'), critvals = empdistribution)$p.value
  ps[i,2]<-runs.test(residuals(sim_glm, type='pearson'))$p.value
}

errrate[counter,1]<-(length(which(ps[,1]<0.05))/nsim)*100
errrate[counter,2]<-(length(which(ps[,2]<0.05))/nsim)*100
counter=counter+1
}

plot(nphi, errrate[,1], type='l', col='red', ylim=c(0, range(errrate)[2]), ylab='Error Rate', xlab='Disp',
lines(nphi, errrate[,2], col='blue')
legend(0, 20, legend = c('N(0,1)', 'Empirical'), col = c('blue', 'red'), lty = c(1,1), bty = 'n')

```



Correlated data

```
dat<-makeToyData(200, length.panels=5)
head(dat)

      x evph      mu panelid
1 1.000000    0 3.669297      1
2 1.090909    0 3.770746      1
3 1.181818    0 3.874999      1
4 1.272727    0 3.982135      1
5 1.363636    0 4.092234      1
6 1.454545    0 4.205376      2

nsim=2000
newdat<-generateNoise(nsim, dat$mu, family='poisson', d=1)
rho=seq(0.1, 0.9, by=0.1)
errrate<-matrix(NA, nrow=length(rho), ncol=2)
counter=1

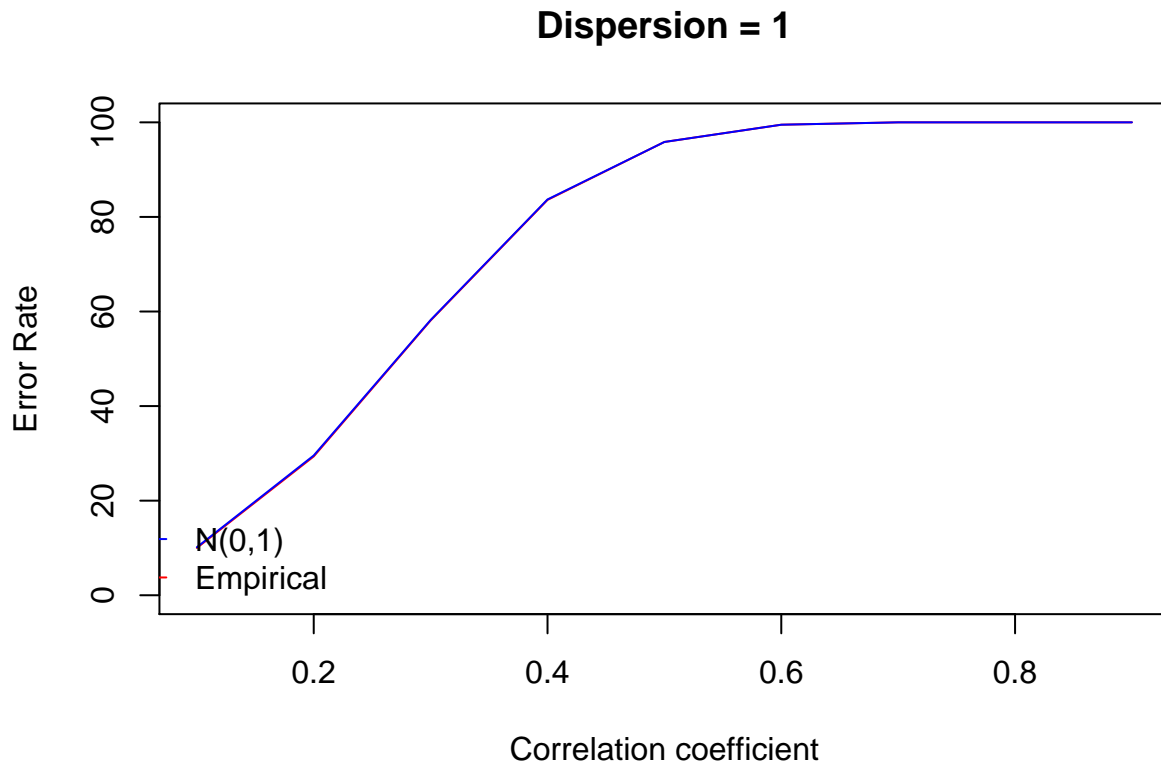
for(r in rho){
  newdatcorr<-generateIC.toy(dat, c(1, r, r^2, r^3, r^4), 'panelid', newdat, ncol(newdat), dots=FALSE)

  # for each change in phi, update the empirical distribution
  empdistribution<-getRunsCritVals(n.sim = nsim, simData=newdat,
                                   model = init_glm, data = dat, plot=FALSE,
                                   returnDist = TRUE, dots=FALSE)

  ps<-matrix(NA, nrow=nsim, ncol=2)
  for(i in 1:nsim){
    sim_glm<-glm(newdatcorr[,i] ~ x, data=dat, family='quasipoisson')
    # find both the empirical and Normal p-values
    ps[i,1]<-runs.test(residuals(sim_glm, type='pearson'), critvals = empdistribution)$p.value
    ps[i,2]<-runs.test(residuals(sim_glm, type='pearson'))$p.value
  }

  errrate[counter,1]<-(length(which(ps[,1]<0.05))/nsim)*100
  errrate[counter,2]<-(length(which(ps[,2]<0.05))/nsim)*100
  counter=counter+1
}

plot(rho, errrate[,1], type='l', col='red', ylim=c(0, range(errrate)[2]), ylab='Error Rate', xlab='Corr
lines(rho, errrate[,2], col='blue')
legend(0, 20, legend = c('N(0,1)', 'Empirical'), col = c('blue', 'red'), lty = c(1,1), bty = 'n')
```



What does all this mean for power analysis?

If data are assumed to be correlated, then a correlation structure is chosen and GEE - type model fitted. If the independent working correlation structure is specified, then robust standard errors are used for inference. Does this machinery still work when there is actually no correlation present in the data?

Question:

1. Are robust standard errors the same as raw standard errors when no correlation is present?
2. What happens when a blocking structure is chosen that is not necessary. i.e. does the answer to question 1 vary with differing block structures?
3. Is the power to detect change affected by using robust standard errors when not necessary?

Make data but add an impact. Here we repeat the data but have a 50% decline in the response for the second set.

50% decline, $\phi = 50$

```
datimp<-makeToyData(200, length.panels=5, changecoeff.link = log(0.5))
head(datimp)
```

	x	evph	mu	panelid
1	1.000000	0	3.669297	1
2	1.090909	0	3.770746	1
3	1.181818	0	3.874999	1

```
4 1.272727    0 3.982135      1
5 1.363636    0 4.092234      1
6 1.454545    0 4.205376      2
```

Generate independent noisy data with a high dispersion, e.g. $\phi = 50$

```
nsim=2000
newdat<-generateNoise(nsim, datimp$mu, family='poisson', d=50)

init_impglm<-glm(newdat[,1] ~ x + as.factor(evph), data=datimp, family='quasipoisson')
summary(init_impglm)
```

Call:

```
glm(formula = newdat[, 1] ~ x + as.factor(evph), family = "quasipoisson",
    data = datimp)
```

Deviance Residuals:

```
      Min       1Q   Median       3Q      Max
-6.9581  -3.3820  -2.2934  -0.4455  16.0946
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    0.62427    0.43070   1.449 0.148811
x              0.32972    0.05441   6.060 6.81e-09 ***
as.factor(evph)1 -0.87350    0.25678  -3.402 0.000811 ***
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for quasipoisson family taken to be 31.53828)

```
Null deviance: 5603.6  on 199  degrees of freedom
Residual deviance: 3792.4  on 197  degrees of freedom
AIC: NA
```

Number of Fisher Scoring iterations: 7

Get the empirical distribution critical values:

```
empdistribution<-getRunsCritVals(n.sim = nsim, simData=newdat,
                                model = init_impglm, data = datimp, plot=FALSE,
                                returnDist = TRUE, dots=FALSE)
```

```
runspowsim_out<-runsPowerSim(newdat, init_impglm, empdistribution, nsim, powercoefid = 3)
```

```
apply(runspowsim_out$rawrob, 2, sum)/nsim
```

```
      Emp      Norm
0.0555 0.8510
```

```
(length(which(runspowsim_out$imppvals[,1]<=0.05))/nsim)*100
```

```
[1] 72.05
```

```
(length(which(runspowsim_out$imppvals[,2]<=0.05))/nsim)*100
```

```
[1] 73.4
```

25% decline, $\phi = 50$

```
Emp  Norm
0.066 0.609
[1] 23.45
[1] 26
```

0% decline, $\phi = 50$

```
Emp  Norm
0.0510 0.4485
[1] 5.15
[1] 5.85
```

It seems that when the Normal distribution is used to test for correlation, consistently the power to detect change is higher than when the empirical distribution is used to choose which standard errors to use.

Test: Repeat for 50% decline to get a range of power values. Check to see if the discrepancy is consistent.

50% decline, $\phi = 50$, multiruns

```
wilcox.test(x = runspower[,1], runspower[,2], paired = T)
```

Wilcoxon signed rank test with continuity correction

```
data: runspower[, 1] and runspower[, 2]
V = 0, p-value < 2.2e-16
alternative hypothesis: true location shift is not equal to 0
wilcox.test(x = power[,1], power[,2], paired = T)
```

Wilcoxon signed rank test with continuity correction

```
data: power[, 1] and power[, 2]
V = 12, p-value < 2.2e-16
alternative hypothesis: true location shift is not equal to 0
apply(power, 2, mean)
```

```
[1] 72.9140 74.0305
```

```
apply(runspower, 2, mean)
```

```
[1] 0.059565 0.834885
```

Statistically we have a significant difference between the power calculated from raw/robust s.e. but the real world significance in this toy example is minimal. The mean difference is one percentage point (with raw giving a higher power).

25% decline, $\phi = 50$, multiruns

Wilcoxon signed rank test with continuity correction

```
data: runspower[, 1] and runspower[, 2]
V = 0, p-value < 2.2e-16
alternative hypothesis: true location shift is not equal to 0
```

Wilcoxon signed rank test with continuity correction

```
data: power[, 1] and power[, 2]
V = 0, p-value < 2.2e-16
alternative hypothesis: true location shift is not equal to 0
[1] 23.6045 25.3685
[1] 0.057965 0.620265
```

Real Data: Falls of Warness

```
init_glm<-glm(response ~ TideState + WindStrength + SeaState + SimpPrecipitation + CloudCover, data=dat,
nsim=500
newdat<-generateNoise(nsim, fitted(init_glm), family='poisson', d=1)
```

500 sets of noisy data are simulated from this truth using a Poisson distribution.

To test, fit a glm to one of the sets of data.

```
fowsim_glm<-glm(newdat[,1] ~ TideState + WindStrength + SeaState + SimpPrecipitation + CloudCover, data=newdat,
summary(fowsim_glm)
```

Call:

```
glm(formula = newdat[, 1] ~ TideState + WindStrength + SeaState +
    SimpPrecipitation + CloudCover, family = quasipoisson, data = dat)
```

Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-2.6354	-1.2238	-0.1367	0.5286	4.2756

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.755768	0.056922	-13.277	< 2e-16 ***
TideState	-0.120080	0.008578	-13.998	< 2e-16 ***
WindStrength	-0.363828	0.012430	-29.269	< 2e-16 ***
SeaState	0.397512	0.012605	31.537	< 2e-16 ***
SimpPrecipitationHEAVY	-0.109604	0.095631	-1.146	0.252
SimpPrecipitationLIGHT	0.955581	0.054166	17.642	< 2e-16 ***
SimpPrecipitationNONE	0.358279	0.049992	7.167	7.88e-13 ***
SimpPrecipitationSHOWERS	0.877451	0.050528	17.366	< 2e-16 ***
CloudCover	0.097307	0.002949	33.000	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for quasipoisson family taken to be 1.001435)

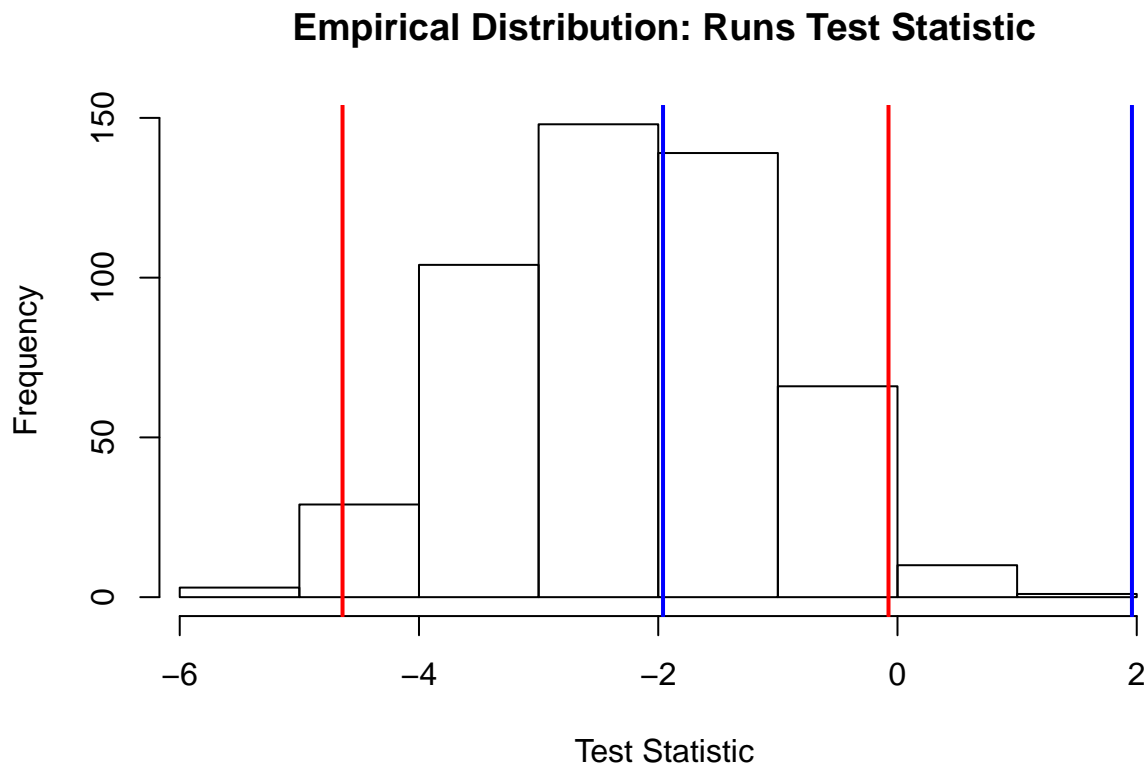
Null deviance: 34638 on 26334 degrees of freedom
Residual deviance: 29493 on 26326 degrees of freedom
AIC: NA

Number of Fisher Scoring iterations: 5

Models were fitted to all datasets generated above and the runs test was evaluated for each one. The distribution of the test statistics is returned. A plot is also produced, which shows the lower 2.5% and upper 97.5% critical values of the empirical distribution (red) and from the Normal ($N(0, 1)$) distribution.

```
empdistribution<-getRunsCritVals(n.sim = nsim, simData=newdat,  
                                model = fowsim_glm, data = dat, plot=TRUE,  
                                returnDist = TRUE)
```

.....



Use the data generated to assess the 5% error rate for this test when using the Normal distribution or the empirical distribution generated above.

```
newdat<-generateNoise(nsim, fitted(init_glm), family='poisson', d=1)  
ps<-matrix(NA, nrow=nsim, ncol=2)  
  
for(i in 1:nsim){  
  sim_glm<-update(fowsim_glm, newdat[,i]~.)  
}
```

```

# find both the empirical and Normal p-values
ps[i,1]<-runs.test(residuals(sim_glm, type='pearson'), critvals = empdistribution)$p.value
ps[i,2]<-runs.test(residuals(sim_glm, type='pearson'))$p.value
}

```

```

(length(which(ps[,1]<0.05))/nsim)*100

```

```

[1] 5.6

```

```

(length(which(ps[,2]<0.05))/nsim)*100

```

```

[1] 62.6

```

How does this change if we vary the dispersion parameter for the data?

```

nphi=seq(1,51, by=5)
errrate<-matrix(NA, nrow=length(nphi), ncol=2)
counter=1
nsim=500
for(p in nphi){
  newdat<-generateNoise(nsim, fitted(init_glm), family='poisson', d=p)
  # for each change in phi, update the empirical distribution
  empdistribution<-getRunsCritVals(n.sim = nsim, simData=newdat,
                                   model = init_glm, data = dat, plot=FALSE,
                                   returnDist = TRUE, dots=FALSE)

  ps<-matrix(NA, nrow=nsim, ncol=2)
  for(i in 1:nsim){
    sim_glm<-update(fowsim_glm, newdat[,i]~.)
    # find both the empirical and Normal p-values
    ps[i,1]<-runs.test(residuals(sim_glm, type='pearson'), critvals = empdistribution)$p.value
    ps[i,2]<-runs.test(residuals(sim_glm, type='pearson'))$p.value
  }

  errrate[counter,1]<-(length(which(ps[,1]<0.05))/nsim)*100
  errrate[counter,2]<-(length(which(ps[,2]<0.05))/nsim)*100
  counter=counter+1
}

```

```

plot(nphi, errrate[,1], type='l', col='red', ylim=c(0, range(errrate)[2]), ylab='Error Rate', xlab='Dispersion',
lines(nphi, errrate[,2], col='blue')
legend(0, 30, legend = c('N(0,1)', 'Empirical'), col = c('blue', 'red'), lty = c(1,1), bty = 'n')

```

