# Inference for quantitation parameters in q-PCR via branching processes with random effects

## Authors: Bret Hanlon and Anand Vidyashankar

Wesley Brooks

UW-Madison

April 25, 2012

# Outline

# Goal

- Compare expression of some gene between treatment conditions

# Table of Contents

# What is qPCR?

- qPCR stands for quantitative polymerase chain reaction
- PCR is the chemical reaction by which DNA replicates
  - e.g. during cell division
- The technology is designed so that PCR replicates only a target gene
- "Quantitative" because we want to count the number of gene copies in a sample
  - Absolute quantitation: count the gene copies from the original sample
  - Relative quantitation: find the ratio of copies of one gene, compared to another
  - The paper covers both; we'll only consider absolute quantitation here

# qPCR as a branching process

- PCR is controlled so that each replication cycle $k$ is discrete
- Each particle either doubles or does not during each trial
  - Probability of replication is typically high $(0.9 < p)$
- This defines a supercritical branching process that leads to exponential growth
- During early cycles ($k < 15$, say), the count is obscured by noise
- Availability of reaction chemicals attenuates the reaction after $\sim 30$ cycles
- The cycles between 15 and 30 are called the exponential phase.

# The experimental setup

A typical experimental setup is to have:

- Two treatment groups (T=treatment, C=control)
- Two genes under study
- Three replicates for each gene-treatment combination
- So an experiment typically involves twelve reactions

# The experimental procedure

Do the following for each reaction:

```
for (cycle in 1:40):
    count the gene copies
    use PCR to produce a new generation
```
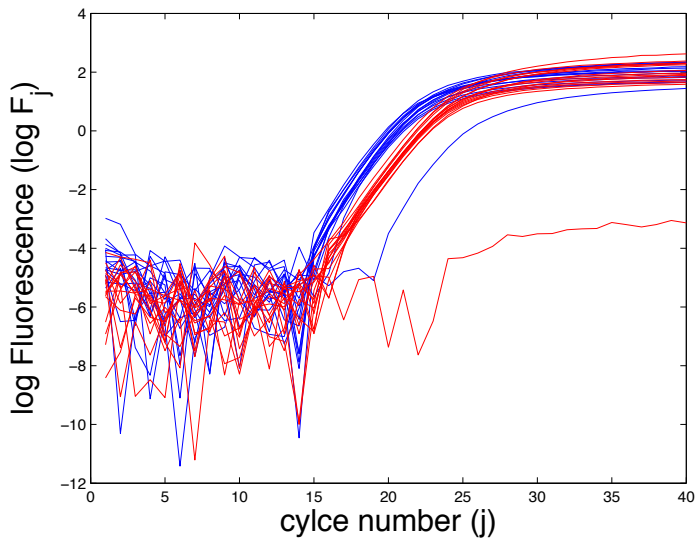
# Sample experimental data

# Table of Contents

## Sources of randomness

- $N_0$, the initial number of gene copies, is random

  - $E[N_0] = m_a$

  - $var(N_0) = \sigma_a^2$

- At each cycle of the reaction, each gene copy replicates randomly

  - $N_{n+1} = N_n + \text{Bin}(N_n, p)$

# qPCR as a branching process

Note:

$$E[N_n] = E[E(N_n|N_{n-1})] = E[(1+p)N_{n-1}]$$
$$= \cdots = (1+p)^n \times E(N_0)$$

- So $W_n = \frac{N_n}{(1+p)^n}$ is a positive martingale

- Thus, $W_n \to W$ almost surely for some $W$

- $E[W] = m_a$

# qPCR as a branching process

Consider an idealized reaction experiment:

- If we knew p and could let the number of reaction cycles $n \to \infty$:

    - $W_i = \lim_{n \to \infty} \frac{N_{i,n}}{(1+p)^n}$

    - $W_1, W_2, \ldots, W_r \overset{\text{iid}}{\sim} W$

    - So $\frac{1}{r} \Sigma_{i=1}^{r} W_i \overset{\text{a.s.}}{\to} E[W] = m_a$

- But p is unknown and we can only observe $\approx 15$ reaction cycles, so we need some other estimator.

## Estimating p

- Since $p$ is unknown, we estimate it with $\hat{p}$ via weighted least squares:

$$\begin{pmatrix} N_n \\ N_{n-1} \\ \vdots \\ N_1 \end{pmatrix} - \begin{pmatrix} N_{n-1} \\ N_{n-2} \\ \vdots \\ N_0 \end{pmatrix} = p \times \begin{pmatrix} N_{n-1} \\ N_{n-2} \\ \vdots \\ N_0 \end{pmatrix} + \epsilon$$

- Where $\epsilon_j \overset{\text{approx.}}{\sim} \text{Normal}(0, p(1-p)N_{j-1})$

- With weights $W_j = (N_{j-1})^{-1}$ the resulting estimator is:

$$\hat{p} = \frac{\Sigma_{i=1}^{n}(N_i - N_{i-1})}{\Sigma_{i=1}^{n}N_i}$$

## Making the most of a finite sample

Reminder: our idealized estimator was $W(n) = \frac{N_n}{(1+p)^n}$

- $W$ uses only the final observation ($N_n$)

- More efficient: use the sum $Y_n = \Sigma_{i=1}^n N_i$

- By the Toeplitz Lemma, $\frac{Y_n}{(1+p)^n} \overset{\text{a.s.}}{\to} \frac{1+p}{p} W \Rightarrow \frac{pY_n}{(1+p)^{n+1}} \overset{\text{a.s.}}{\to} W$

- Plug in $\hat{p}$ and the limit still holds.

# Strategy for quantitation

- Collect data on $r$ independent reactions

- For reaction $i$ ($i = 1, 2, \ldots, r$), compute the statistic $M_i = \frac{\hat{p}_i Y_{n_i}}{(1 + \hat{p}_i)^{n_i + 1}}$

- Average $M_1, M_2, \ldots, M_r$ to get $\bar{M}$

- $\sqrt{r}(\bar{M} - m_a) \xrightarrow{d} \text{Normal}(0, \sigma_L^2)$

- Where $\sigma_L^2 = \sigma_a^2 + m_a E[\frac{1-p}{1+p}]$

# Variance of the estimator

$$\sigma_L^2 = \text{var}[\frac{N_n}{(1+p)^n}] = E(\text{var}[\frac{N_n}{(1+p)^n}|p]) + \text{var}(E[\frac{N_n}{(1+p)^n}|p])$$

$$= E(\text{var}[\frac{N_n}{(1+p)^n}|p]) + \text{var}(m_a)$$

$$= E(\text{var}[\frac{N_n}{(1+p)^n}|p])$$

## Variance of the estimator

$$\text{var}[\frac{N_n}{(1+p)^n}|p] = \frac{1}{(1+p)^{2n}}\text{var}[N_n|p]$$

$$= \frac{1}{(1+p)^{2n}}(E(\text{var}[N_n|N_{n-1}, p]|p) + \text{var}(E[N_n|N_{n-1}, p]|p))$$

$$= \frac{1}{(1+p)^{2n}}(E[N_{n-1}p(1-p)|p] + \text{var}((1+p)N_{n-1}|p))$$

$$= \frac{1}{(1+p)^{2n}}(m_a(1+p)^{n-1}p(1-p) + (1+p)^2\text{var}[N_{n-1}|p])$$

$$= \frac{m_a p(1-p)}{(1+p)^{n+1}} + \frac{\text{var}[N_{n-1}|p]}{(1+p)^{2n-2}}$$

$$= \dots$$

$$= \frac{m_a p(1-p)}{(1+p)^{n+1}} + \frac{m_a p(1-p)}{(1+p)^n} + \dots + \frac{m_a p(1-p)}{(1+p)^2}$$

$$+ \frac{\text{var}[N_0|p]}{(1+p)^{2n-2n}}$$

# Variance of the estimator

$$\text{var}[\frac{N_n}{(1+p)^n}|p] = \frac{m_a p(1-p)}{(1+p)^2} \Sigma_{k=0}^{n-1} \frac{1}{(1+p)^k} + \sigma_a^2$$
$$\rightarrow m_a \frac{1-p}{1+p} + \sigma_a^2$$

So:

$$\text{var}[\frac{N_j}{(1+p)^n}] = E(\text{var}[\frac{N_n}{(1+p)^n}|p])$$
$$\rightarrow E(m_a \frac{1-p}{1+p} + \sigma_a^2)$$
$$= m_a E(\frac{1-p}{1+p}) + \sigma_a^2 = \sigma_L^2$$

# Table of Contents

# Experimental data - luteinizing hormone

- The goal with the experimental data was *relative* quantitation
    - Estimate ratio of gene expression between conditions C and T
- The sample was divided into two parts
- One part was diluted to one-third the original concentration
- Sixteen reactions were run under each condition (diluted, normal)

# Experimental data - luteinizing hormone

Luteinizing Hormone

|       | BP                 | $C_T$ method       | Std. Curve         | Adj $C_T$          |
|-------|--------------------|--------------------|--------------------|--------------------|
| $\hat{R}$ | 2.8221         | 3.3108             | 3.5558             | 3.9567             |
| GCI   | $[1.8624, 3.7817]$ | .                  | .                  | .                  |
| tCI   | $[1.7719, 3.8722]$ | .                  | .                  | .                  |
| BCI   | $[1.6870, 3.6013]$ | $[2.7935, 3.7477]$ | $[2.8646, 4.1355]$ | $[2.9024, 5.4157]$ |

# Experimental data - Strongylus vulgaris

- Again, the goal of the was relative quantitation
- One part diluted to one-tenth the original concentration
- Ten reactions run under each condition (diluted, normal)

# Experimental data - Strongylus vulgaris

S.vulgaris

|  | BP | $C_T$ method | Adj $C_T$ |
|---|---|---|---|
| $\hat{R}$ | 10.1793 | 6.3622 | 303.1662 |
| GCI | [2.8686, 17.4900] | . | . |
| tCI | [1.7414, 18.6172] | . | . |
| BCI | [1.0446, 15.7493] | [5.2733, 7.3176] | [29.7583, 892.6178] |