

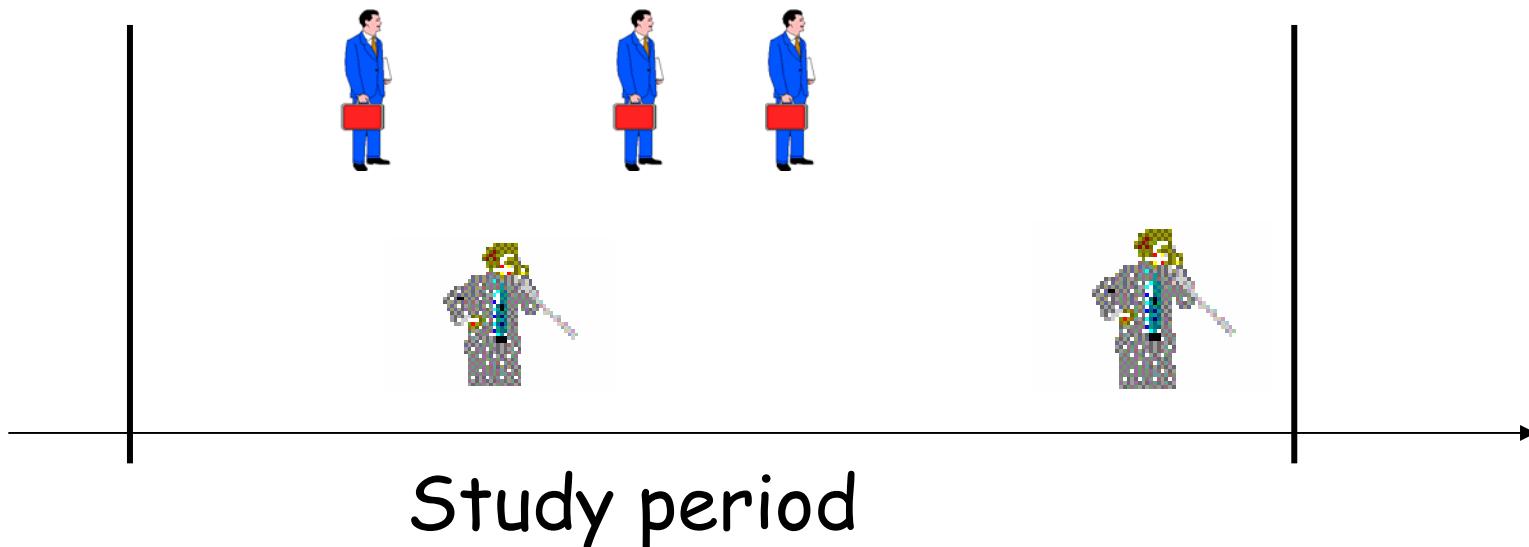
Nonparametric Mixture Maximum Likelihood Truncated Count Distributions

Dankmar Böhning
Professor and Chair in Applied Statistics
School of Biological Sciences
University of Reading, UK

*Round Table: Microbial Biodiversity Research
11th International Symposium on Microbial
Ecology
20-25 August 2006*

Counts of capture-recaptures as outcome of continuous time CR-experiment

- CR of Wildlife Populations
- CR in Public Health and Surveillance



Situation in Continuous CR Experiment

$$f_1, f_2, f_3, \dots, f_m$$

frequencies of units identified 1, 2, 3, ..., m times

f_0 is unobserved

population size: $N = f_0 + f_1 + \dots + f_m = f_0 + n$

if probability p_0 for zero-count known:

$$N = Np_0 + n \Rightarrow \hat{N} = n/(1 - p_0)$$

Illustration: Project with DEFRA on Scrapie in the UK 2002

$$f_1, f_2, f_3, \dots, f_m$$

frequencies of holdings reporting 1, 2, 3, ..., m cases:

$$f_1 = 74, f_2 = 23, f_3 = 15, \dots, f_{11} = 3, f_{12} + = 7$$

f_0 is number of hidden holdings with scrapies

adjusted size of scrapie: $N = f_0 + n = f_0 + 177$

Idea of Modelling

$$f_0, f_1, f_2, f_3, \dots, f_m$$

look at associated probabilities:

$$p_0, p_1, p_2, p_3, \dots, p_m$$

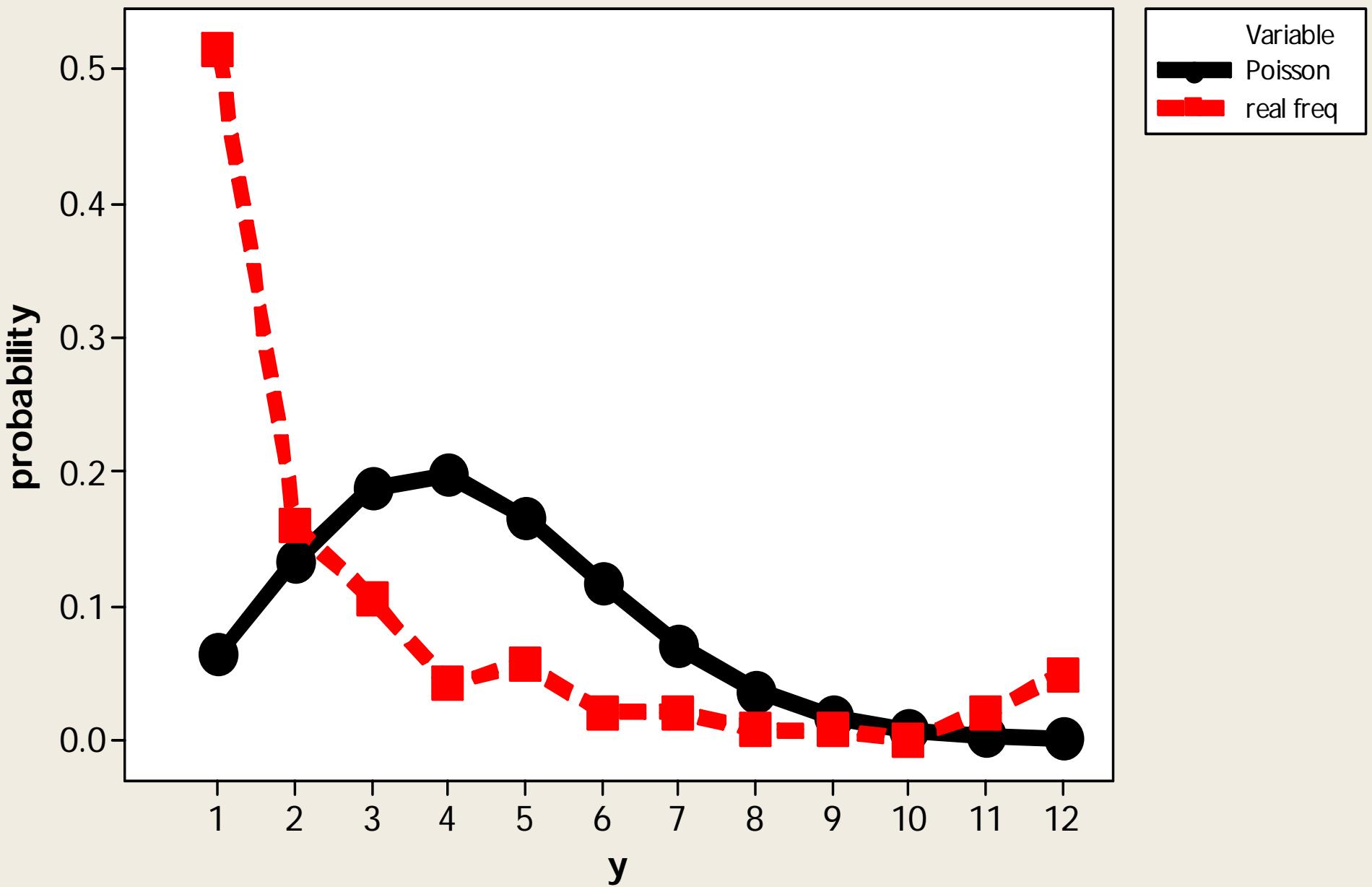
and choose a model (Poisson)

$$p_0 = e^{-\theta}, p_1 = e^{-\theta}\theta, p_2 = e^{-\theta}\theta^2/2, \dots,$$

estimate θ with $\hat{\theta}$, get $\hat{p}_0 = e^{-\hat{\theta}}$

$$\hat{N} = n / (1 - \hat{p}_0)$$

Relative Frequencies and Poisson Probabilities for Count of Cases



Idea of Mixed Modelling

instead of simple Poisson

$$p_j = e^{-\theta} \theta^j / j!$$

look at mixed Poisson:

$$p_j = \int_0^\infty e^{-\theta} \theta^j / j! f(\theta) d\theta$$

(to capture heterogeneity in θ)

Idea of Mixed Modelling

possible to estimate $f(\theta)$ in mixed Poisson:

$$\hat{p}_j = \int_0^{\infty} e^{-\theta} \theta^j / j! \hat{f}(\theta) d\theta$$

by means of nonparametric mixtures

Estimate of $f(\theta)$ for Scrapie 2002

Poisson mixture

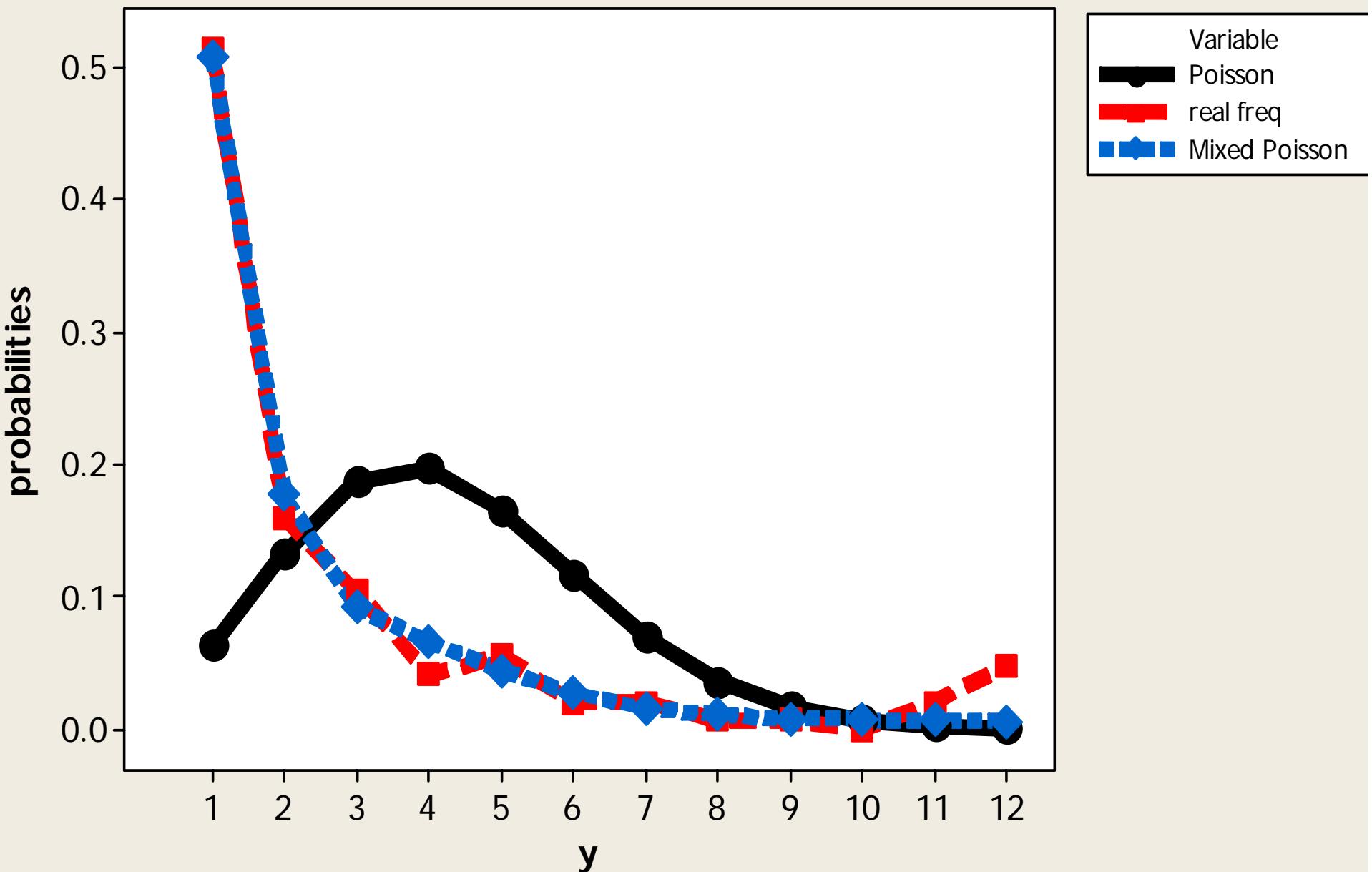
$$\hat{p}_j = \int_0^{\infty} e^{-\theta} \theta^j / j! \hat{f}(\theta) d\theta = \sum_{l=1}^k e^{-\hat{\theta}_l} \hat{\theta}_l^j / j! \hat{f}_l$$

we find

$$\hat{f}(\theta) = \begin{pmatrix} \hat{\theta}_1 & \hat{f}_1 \\ \hat{\theta}_2 & \hat{f}_2 \\ \hat{\theta}_3 & \hat{f}_3 \\ \hat{\theta}_4 & \hat{f}_4 \end{pmatrix} = \begin{pmatrix} .475 & 0.7955 \\ 3.406 & 0.1594 \\ 10.362 & 0.0277 \\ 22.367 & 0.0173 \end{pmatrix}$$

(Round Table Vienna August 2006)

Relative Frequencies, Simple Poisson, Mixed Poisson



Estimates of hidden Scrapie 2002

Simple Poisson:

$$\hat{N} = n / (1 - e^{-\hat{\theta}}) = 144 + 2$$

Mixed Poisson :

$$\hat{N} = n / \left(1 - \sum_{l=1}^4 e^{-\theta_l} \hat{f}_l\right) = 144 + 143$$

Two other nonparametric ideas

- Zelterman's robustness idea
- Chao's lower bound estimate

Idea of Zelterman

look at count probabilities:

$$p_0, p_1, p_2, p_3, \dots, p_m$$

and modelled by Poisson

$$p_0 = e^{-\theta}, p_1 = e^{-\theta}\theta, p_2 = e^{-\theta}\theta^2/2, \dots,$$

ratios of consecutive Poissons

$$p_1 / p_0 = \theta, p_2 / p_1 = \theta/2, p_3 / p_2 = \theta/3, \dots,$$

$$\Rightarrow \hat{\theta} = 2f_2 / f_1$$

$$\hat{N}_z = n / (1 - \exp(-2f_2 / f_1))$$

Round Table Vienna August 2006

Idea of Chao

look at mixed Poisson:

$$p_j = \int_0^{\infty} e^{-\theta} \theta^j / j! f(\theta) d\theta$$

Cauchy-Schwartz:

$$\left(\int_0^{\infty} e^{-\theta} \theta f(\theta) d\theta \right)^2 \leq \int_0^{\infty} e^{-\theta} f(\theta) d\theta \int_0^{\infty} e^{-\theta} \theta^2 f(\theta) d\theta$$

$$p_1^2 \leq p_0 p_2 \Rightarrow f_0 \geq f_1^2 / (2f_2)$$

Chao's lower bound estimate

Estimates of hidden Scrapie 2002

Simple Poisson:

$$\hat{N} = n / (1 - e^{-\hat{\theta}}) = 144 + 2$$

Mixed Poisson :

$$\hat{N} = n / (1 - \sum_{l=1}^4 e^{-\theta_l} \hat{f}_l) = 144 + 143$$

Zelterman:

$$\hat{N}_Z = n / (1 - e^{-2f_2/f_1}) = 144 + 167$$

Chao:

$$\hat{N}_C = n + f_1^2 / (2f_2) = 144 + 119$$

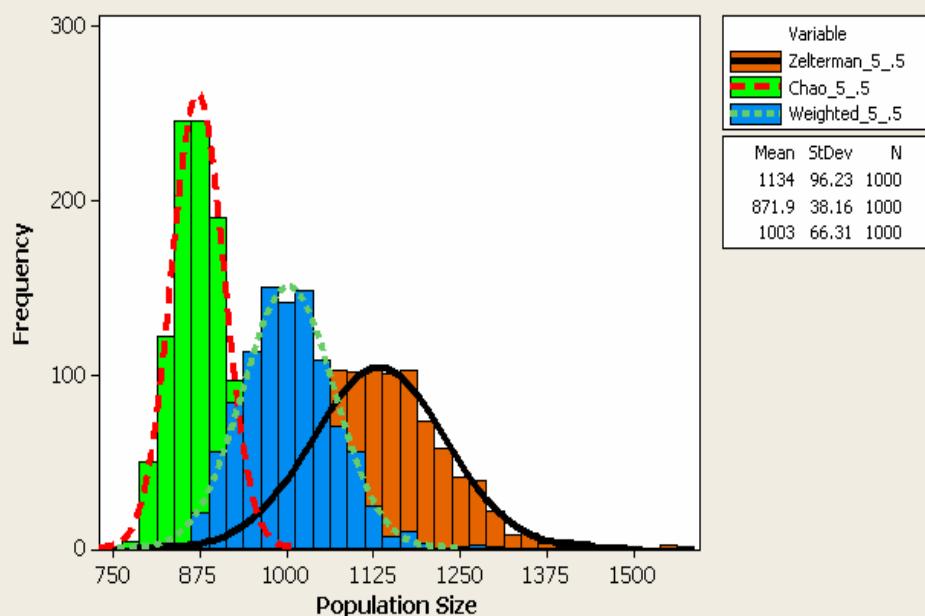
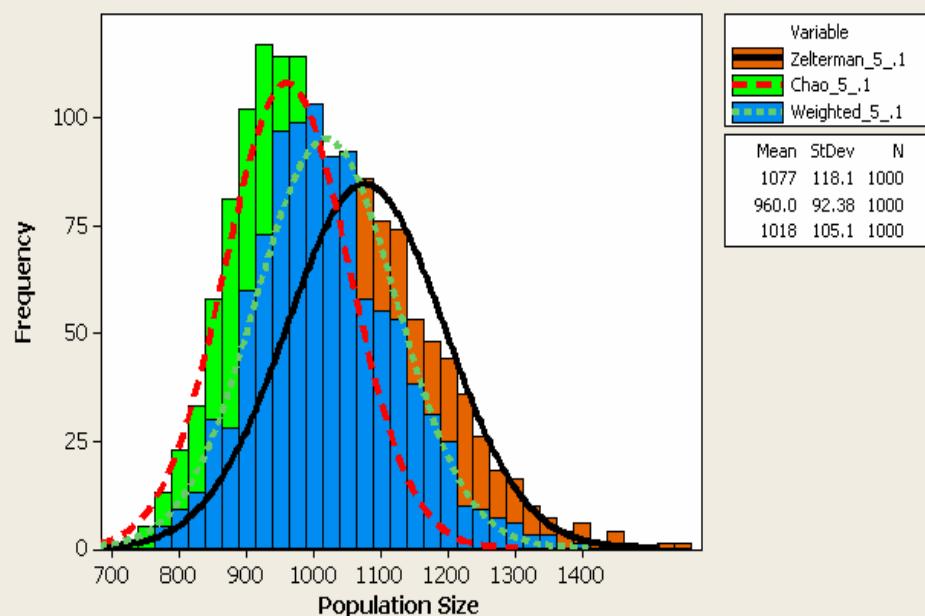
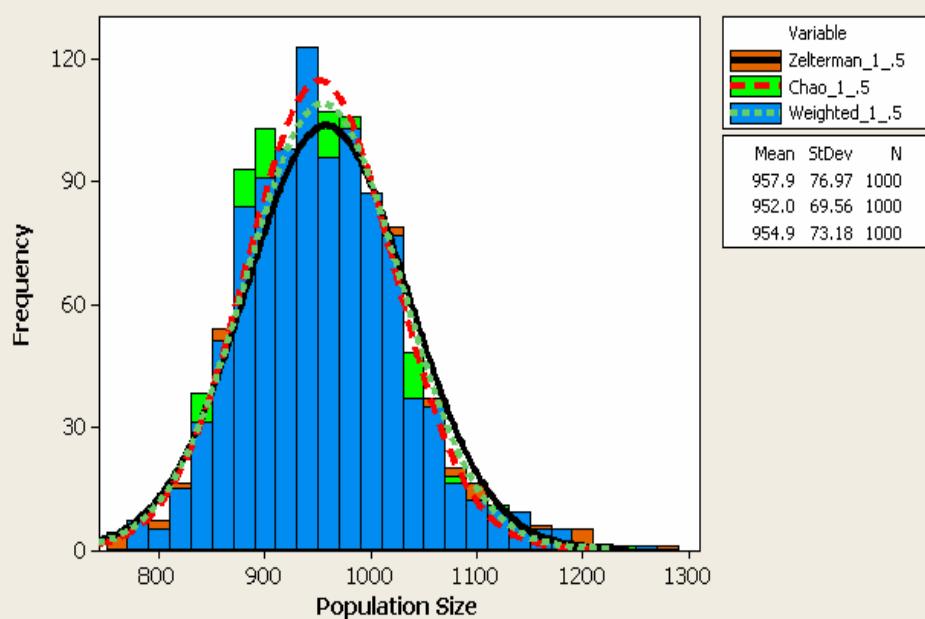
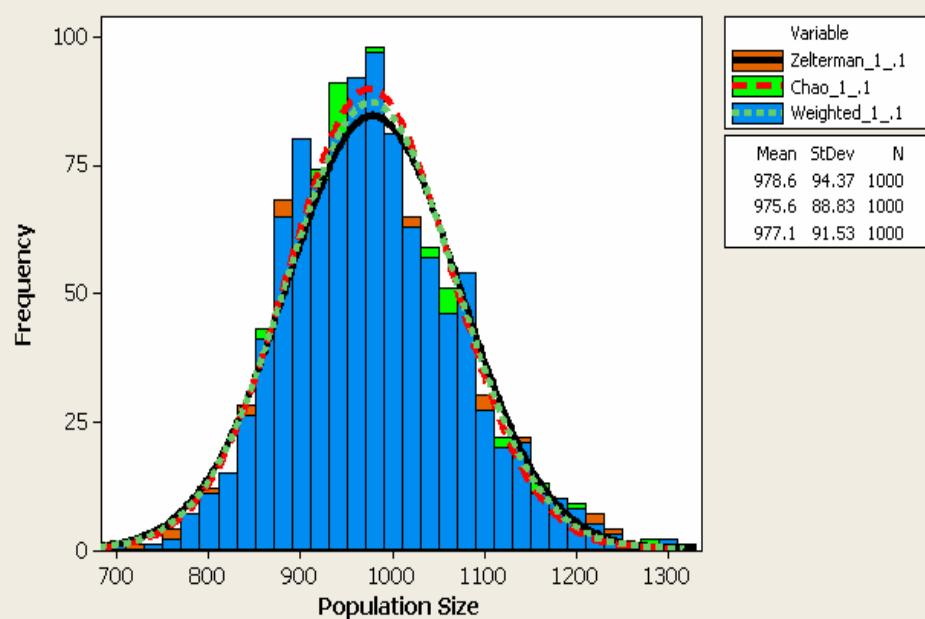
Comparing Zelterman and Chao Simulation: 4 Experiments with Poisson Model Violation

$$1: \quad f(\theta) = Po(0.5)0.5 + 0.5Po(5)$$

$$2: \quad f(\theta) = Po(0.5)0.9 + 0.1Po(5)$$

$$3: \quad f(\theta) = Po(0.5)0.5 + 0.5Po(1)$$

$$4: \quad f(\theta) = Po(0.5)0.9 + 0.1Po(1)$$

Experiment 1**Experiment 2****Experiment 3****Experiment 4**

Conclusions

- Chao's estimator keeps the lower bound
- New Weighted estimator: average between Zelterman and Chao?
- For more valid inference, mixed modelling can't be avoided

Key-References

Böhning, D. and Kuhnert, R. (2006). The Equivalence of Truncated Count Mixture Distributions and Mixtures of Truncated Count Distributions. *Biometrics* (to appear).

Böhning, D. and Schön, D. (2005). Nonparametric maximum likelihood estimation of the population size based upon the counting distribution. *Journal of the Royal Statistical Society, Series C, Applied Statistics* **54**, 721-737.

Böhning, D. and Patilea, V. (2005). Asymptotic Normality in Mixtures of Power Series Distributions. *Scandinavian Journal of Statistics* **32**, 115-132.

Papers download at (also copy of this talk):

www.reading.ac.uk/~sns05dab