

Bayesian approach to modelling plant pest introduction

Jonathan Yuen

Swedish University of Agricultural Sciences, SE 75007 Uppsala,
Sweden

Bayesian analysis



Thomas Bayes (1702-1761)

Tomato Disease Example

- An MSc student was sent to study other diseases on virus-resistant tomatoes in Nicaragua
- Was told by the scientists in Nicaragua that the diseases were early blight and fusarium wilt
- Decision : What agar media do you send with the student
- Depends on the person making the decision and what information they might have about tomato diseases in the tropics

Non-Bayesian approach (Frequentist)

- Basic principle:
 - parameters as fixed but unknown quantities
 - probability as long-run frequency
- The true mean of the population exists and it is a number
- If you sample again and again you will find something that is based on it!
- CI_{95} contains true value 95 out of 100 times

Non-Bayesian approach

Assumptions and estimation for ANOVA

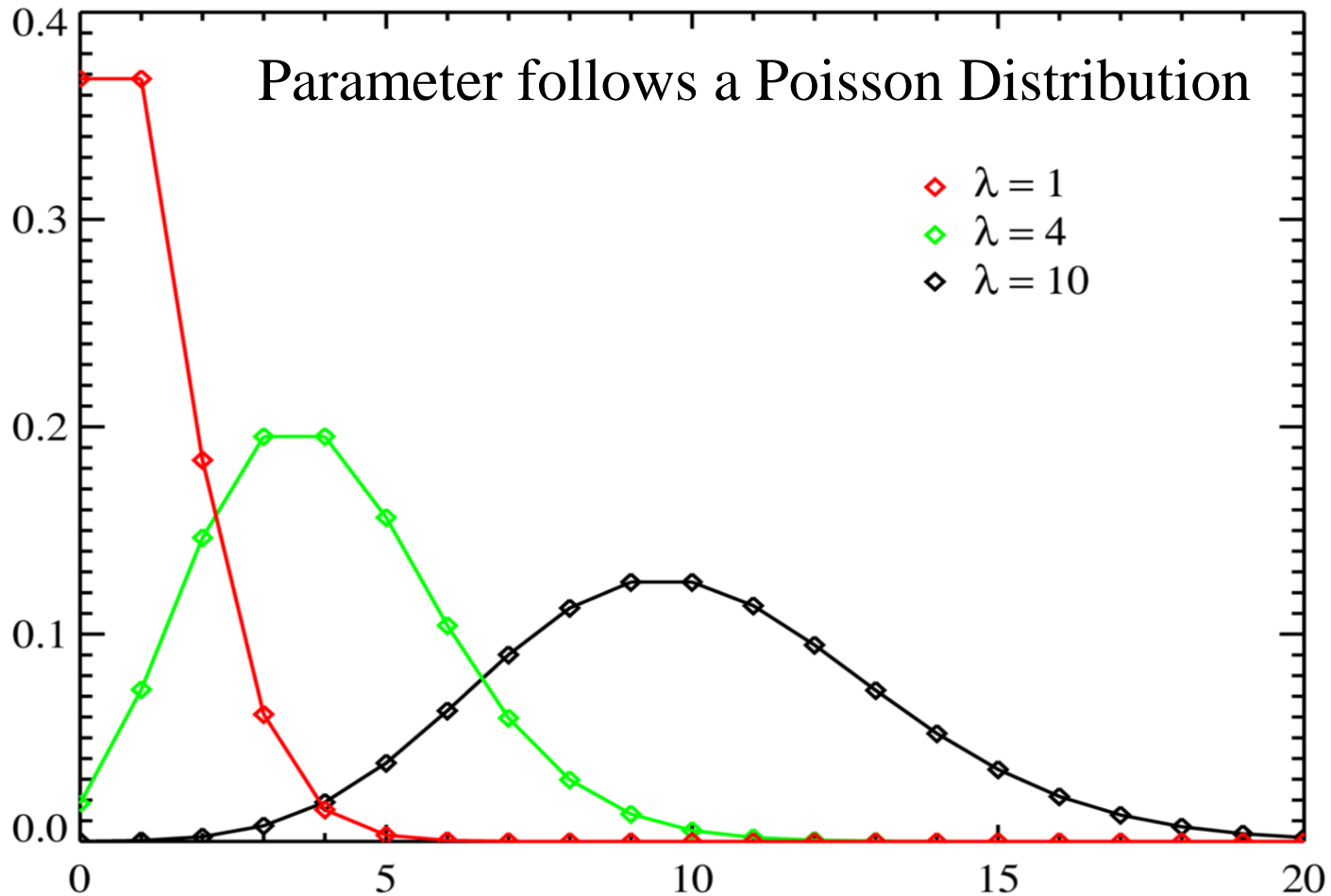
- Unknown but fixed effects contribute to the underlying value for each of the cells
- Normal distribution of deviations from these unknown underlying values and the observed values
- ML estimate is the same as a least-squares estimate

Bayesian inference: Key ideas

- parameters as random variables
- probability used to describe uncertainty about unknowns
- combine information (prior distributions and current data) to draw conclusions

Bayesian methods

Describe uncertainty with probability distributions



Some important terms (or how to talk like a Bayesian)

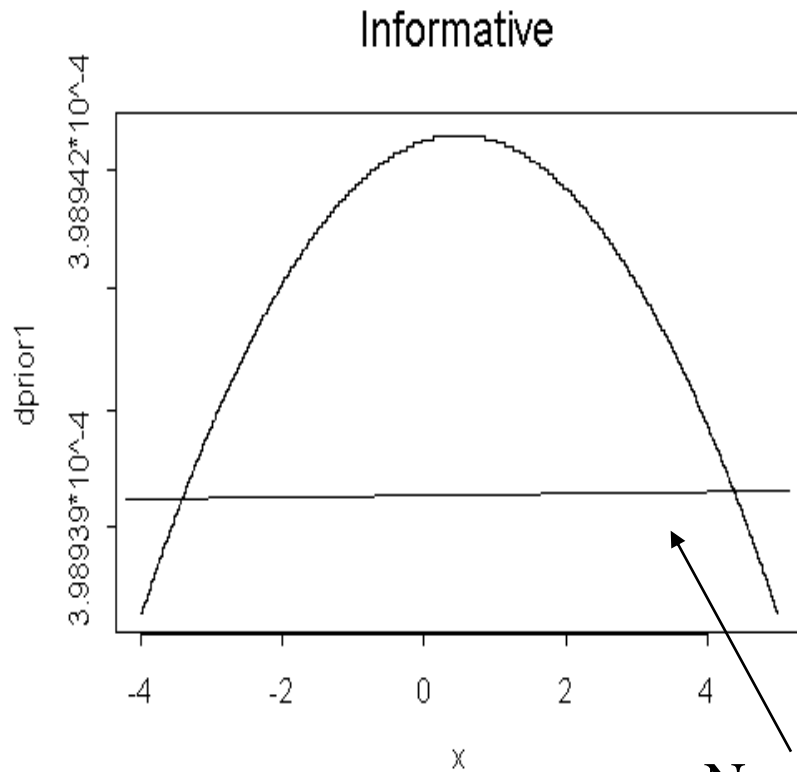
- Prior – idea before collection new information
- New Data – likelihood
- Updated knowledge – Also called posterior
- Credibility interval– two limits derived from the posterior distribution, and calculated such that a pre-determined area under the distribution (95% for example) is between the two limits. Sometimes called a Bayesian confidence interval, but should not be confused with a frequentist confidence interval.

Tomato Example

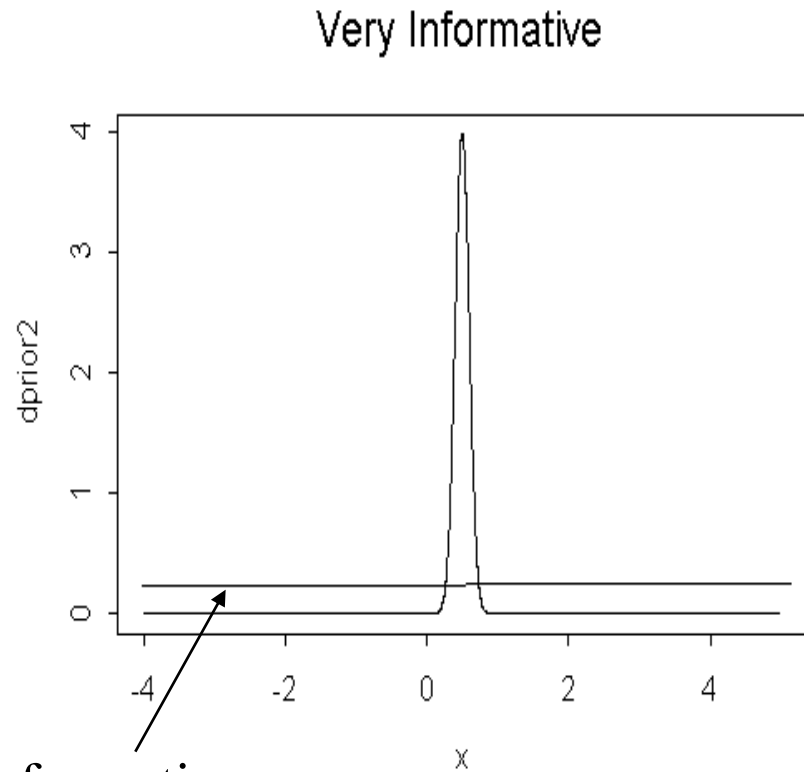
- Prior – varies from person to person.
 - What are the common diseases of tomatoes in the tropics?
- Current data – being told that the diseases were fusarium wilt and early blight
- Updated knowledge from combining data and prior knowledge
- Student's prior was relatively uninformed and PDA was sufficient
- Other priors lead to both PDA and TZC
 - My prior is that bacterial wilt (caused by *Ralstonia solanacearum*) is common in the tropics. I sent a bottle of TZC with the student (and she used it to confirm the presence of *Ralstonia solanacearum*).

Bayesian methods

Distributions to describe prior knowledge/belief



Non-informative

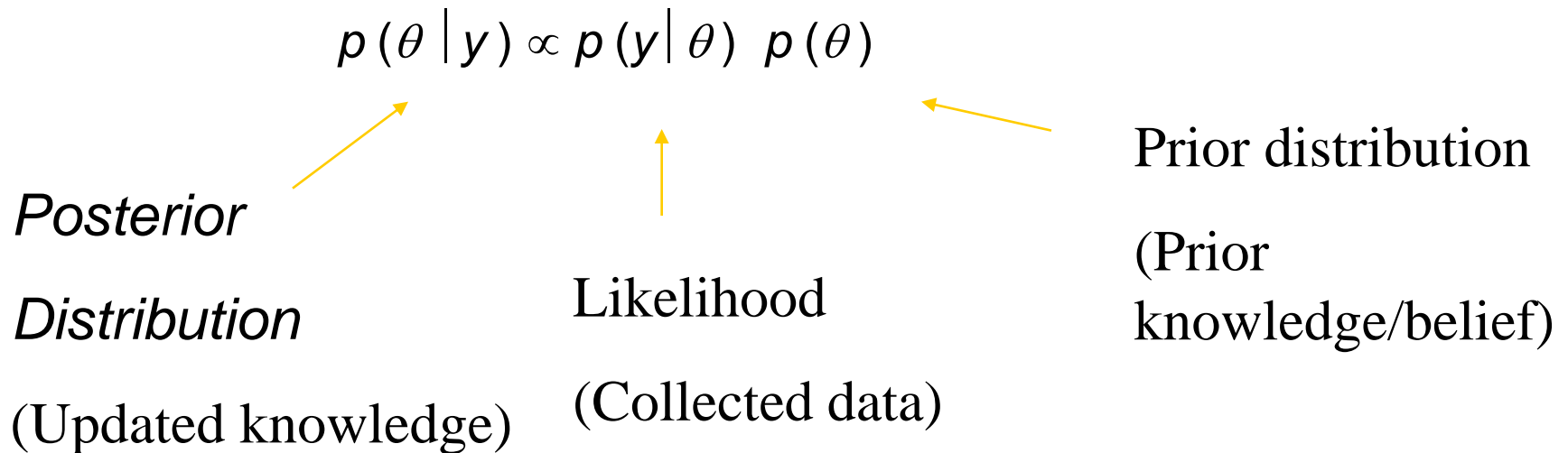


Bayes's Theorem

$$\Pr(A|B) = \frac{\Pr(B|A)\Pr(A)}{\Pr(B|A)\Pr(A) + \Pr(B|\bar{A})\Pr(\bar{A})}$$

Bayesian inference

- Posterior inference
 - Bayes' theorem to derive posterior distribution

$$p(\theta | y) \propto p(y | \theta) p(\theta)$$


Posterior Distribution
(Updated knowledge)

Likelihood
(Collected data)

Prior distribution
(Prior knowledge/belief)

ELISA and HIV

a diagnostic example

- 10,000 HIV positive individuals (tested with a Western blot, the gold standard) were tested with a new ELISA test. 9990 were positive with the ELISA test, 10 were negative
- Sensitivity is $9990/10,000$ or 99.9%
- 10,000 nuns who denied risk factors for HIV were also tested with the ELISA test. 9990 were negative with the ELISA test, 10 were positive
- Specificity is $9990/10,000$ or 99.9%
- $LR+ = 99.9/.1 = 999$
- $LR- = 0.1/99.9 = 0.001$

Likelihood ratios

- $LR+ = 99.9/.1 = 999$
- $LR- = 0.1/99.9 = 0.001$
- Prior odds are $1/99 = 0.010101$ (general population)
- $1/999 = 0.001001$ (blood donors)
- $1/9 = 0.111111$ (drug abusers)

Prior and posterior odds using likelihood ratios

Prior Odds		LR+	LR-
		(999)	(0.001001)
1 in 99	0.010101	10.1	1 x E-5
1 in 999	0.001001	1	1 x E-6
1 in 9	0.111111	111	1 x E-4

**Most problems cannot be solved
in a simple manner and require
other methods**

Numerical Methods (Software)

- Gibbs Sampling
 - BUGS (=Bayesian Using Gibbs Sampling) software WinBugs
 - Openbugs
- Approximate Bayesian Computation (ABC)
 - Useful when calculating the probability of the data, given the parameters, is difficult or impossible.
 - See Makowski et al. (2011) for an application of ABC.
- Integrated Nested Laplace Approximations (INLA)
 - See Martínez-Minaya (2015) for an example

Gibbs Sampling

- **Gibbs sampling** is an algorithm to generate a sequence of samples from the joint probability distribution of two or more random variables.
- Gibbs sampling is an example of a Markov chain Monte Carlo algorithm.

Gibbs Sampler

Posterior distributions typically are hard to obtain analytically

Gibbs sampler: iterative, numerical approach that produces a Markov chain

....that eventually converges and can give information about a posterior distribution

Gibbs Sampling

- **Purpose:** *Draw from a Joint Distribution*

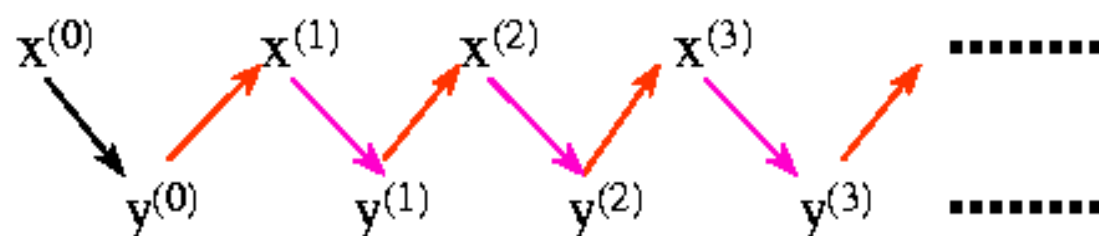
$$x = (x_1, \dots, x_n); \text{ target } \pi(x)$$

- **Method:** *Iterative Conditional Sampling*

$$\forall i, \text{ Draw } x_i \sim \pi(x_i | x_{[-i]})$$

Why does it work?

- It is a *Markov Chain*!!



If $X^{(0)} = x_0$, then distribution of $X^{(t)}$ is $N(\rho^{2t} x_0, 1 - \rho^{4t})$ which “converges” to $N(0, 1)$ as $t \rightarrow \infty$.

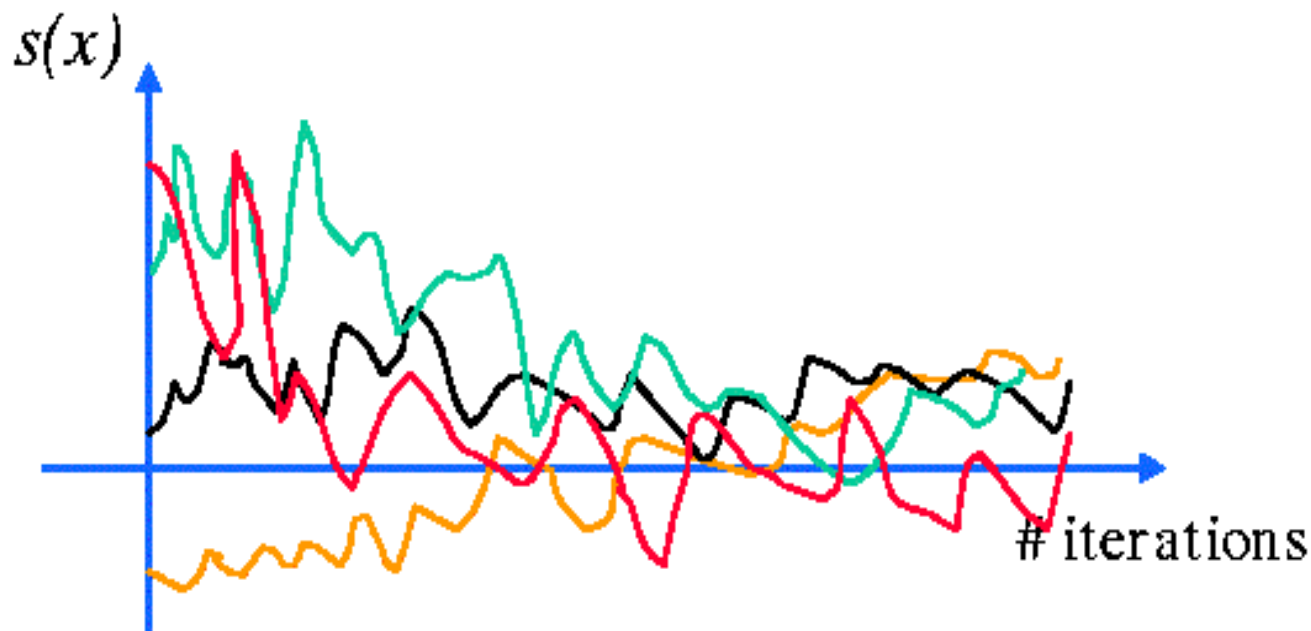
Joint distribution of $(X^{(t)}, Y^{(t)})$?

$$N\left(\begin{pmatrix} \rho^{2t} x_0 \\ \rho^{2t+1} x_0 \end{pmatrix}, \begin{pmatrix} 1 - \rho^{4t} & \rho(1 - \rho^{4t}) \\ \rho(1 - \rho^{4t}) & 1 - \rho^{4t+2} \end{pmatrix}\right) \xrightarrow{t \rightarrow \infty} N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

Convergence (con't)

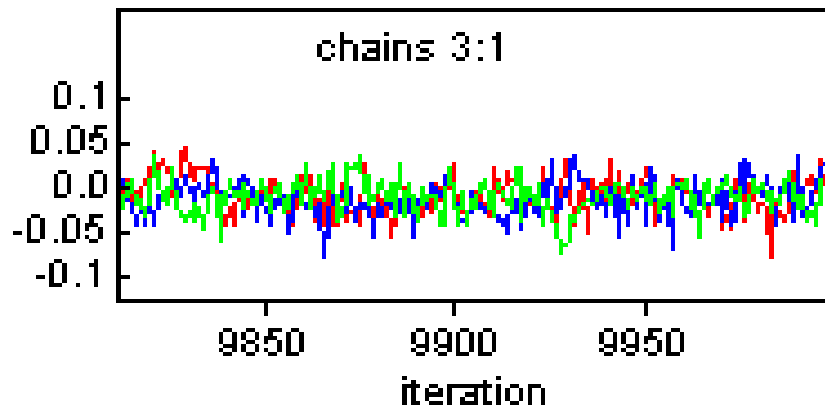
- Gelman and Rubin (1992, *Statist. Sci.*)

Multiple chains approach: Compare within-chain variation and between-chain variation.

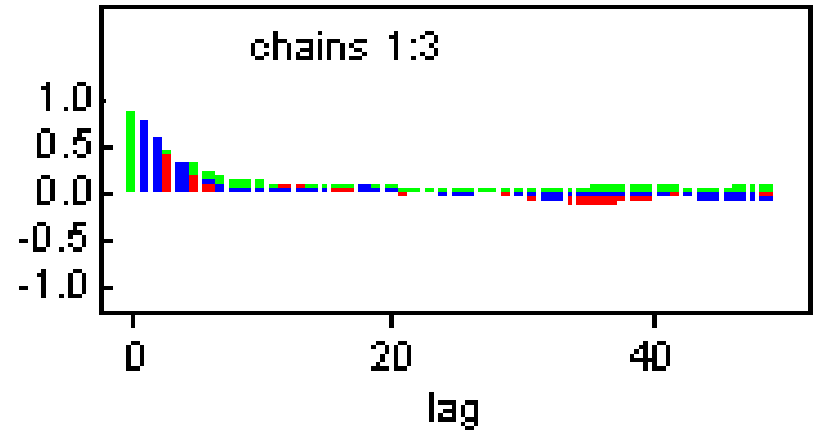


Convergence

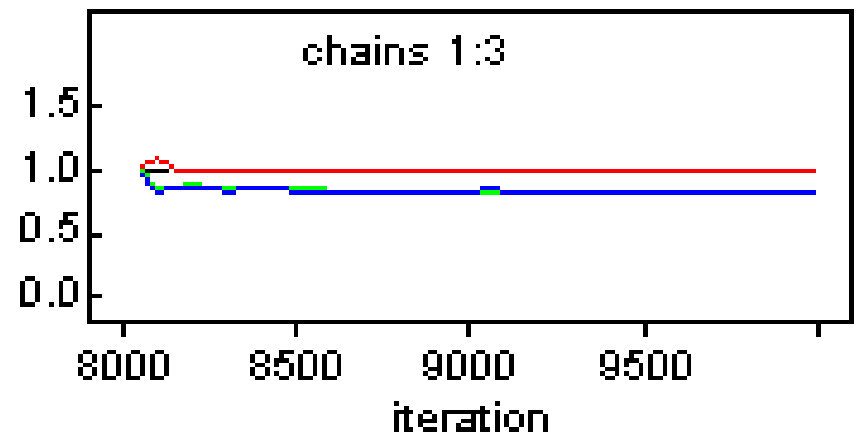
A. Trace plots



B. Autocorrelations

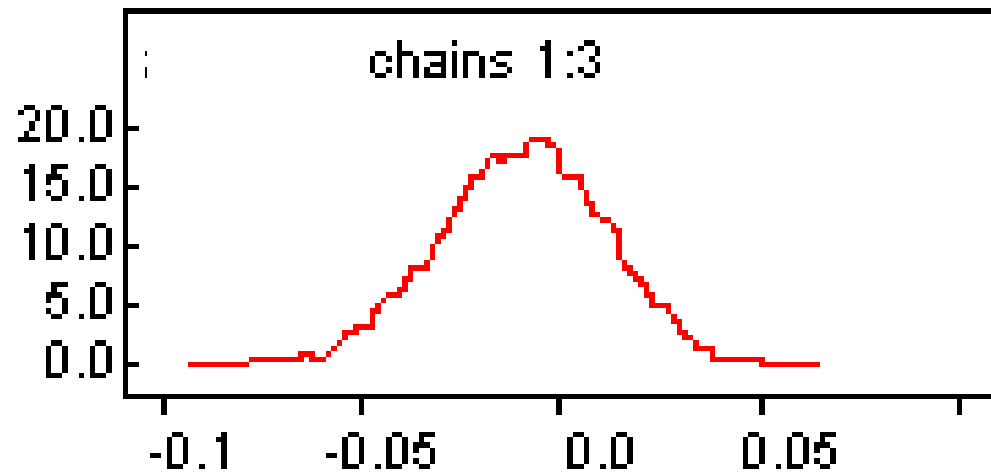
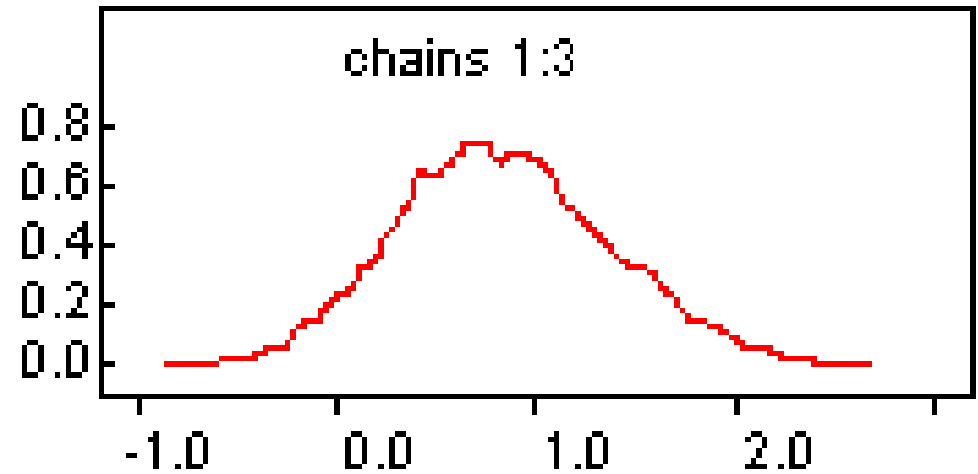


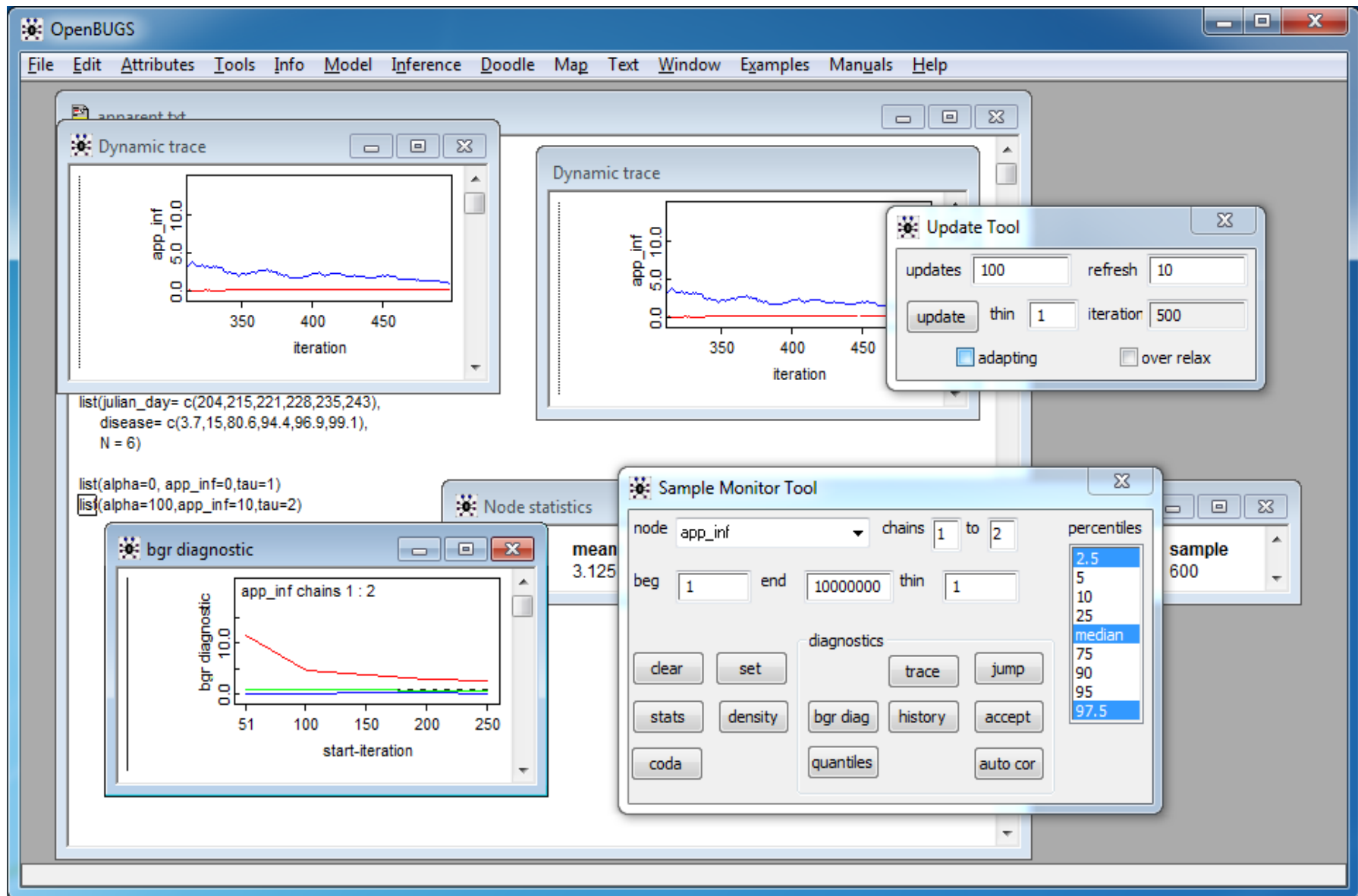
C. Gelman-Rubin (GR)



Convergence

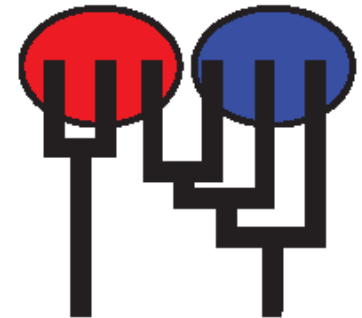
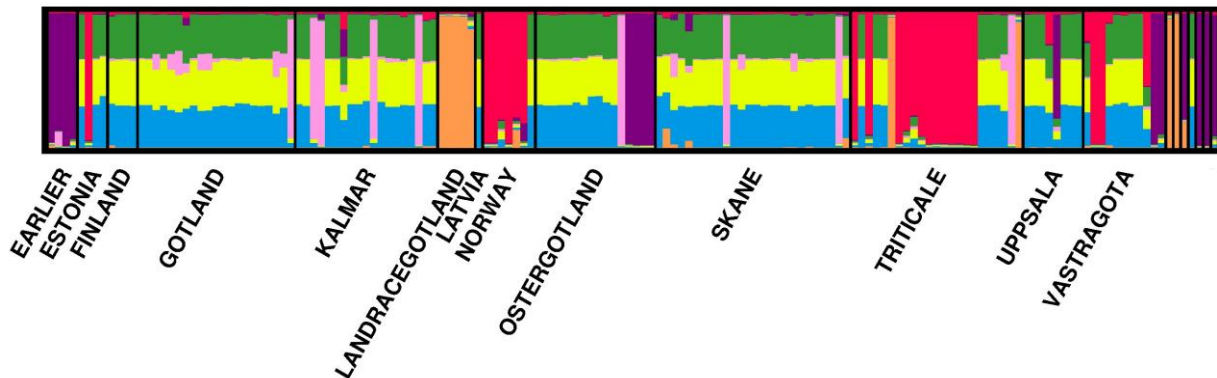
D. Kernel density





Bayesian methods are also a part of other programs

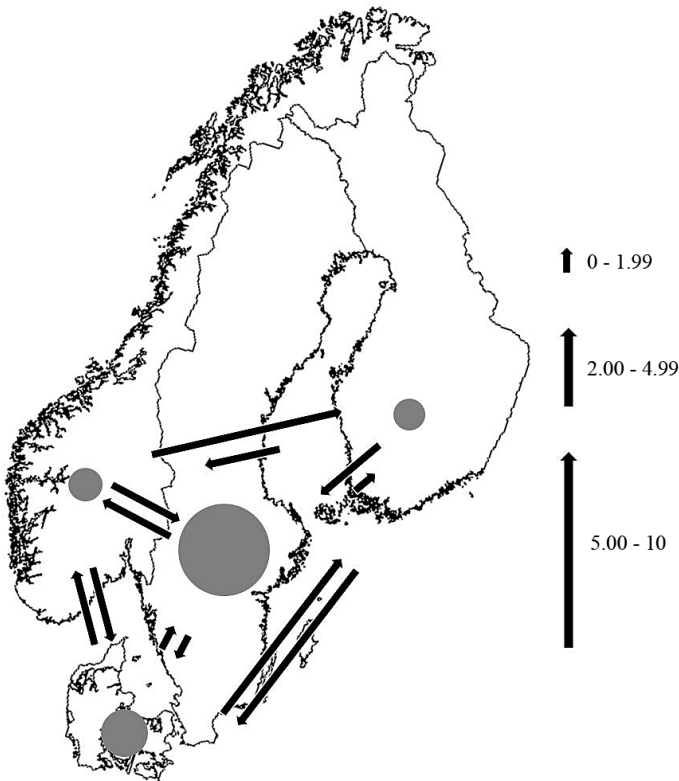
- Analysis of population groupings/origins with 'structure'
- Analysis of populations sizes and migrations patterns with 'migrate'



Nordic migration of *P. infestans*

Table 4. Marginal likelihood, Bayes factor, and model probability comparing one, two, three and four subpopulation.

	Bezier approximation of marginal likelihood	Ln of (bayes factor)	Model probability
4 populations	-2333.54	0	1
3 populations	-2386.73	-53.19	0
2 populations	-2736.30	-402.76	0
1 population	-4060.35	-1726.81	0



**Genotypic and phenotypic variation
of *Phytophthora infestans* on potato
in the two Swedish regions Bjäre and
Östergötland in 2015**

Ida Petersson



Department of Forest Mycology and Plant Pathology
Independent Project in Biology • Master's thesis • 30 HEC • Uppsala • 2015

A comparison between samples from two different regions in Sweden

- Samples taken from Bjäre and Östergötland
- Two fields in each region
- Two foci in each field

Comparison with Bayes's Factor

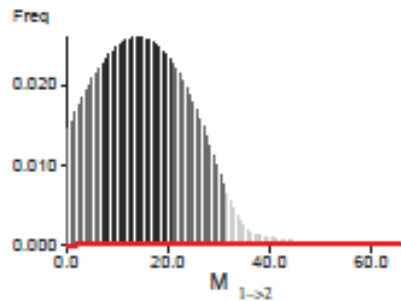
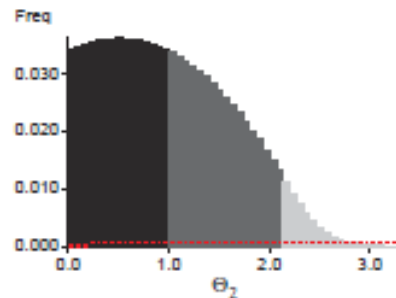
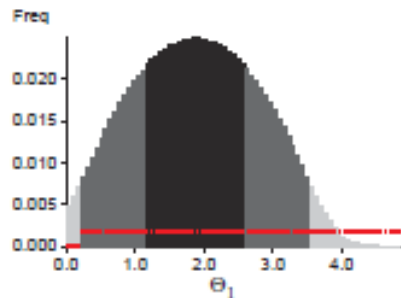
	Marginal Likelihoods				
Model	Run 1	Run 2	Run 3	Run 4	Average
1 single population	-1341.36	-1291.8	-1333.39	-1312.27	-1319.705
2 populations, 2-way migration	-1316.6	-1407.32	-1352.26	-1321.58	-1349.44
2 populations, Bjäre to Östergötland only	-1178.38	-1188.34	-1179.57	-1179.64	-1181.4825
2 populations, Östergötland to Bjäre only	-1169.57	-1209.05	-1201.77	-1199.04	-1194.8575

Larger value Marginal Likelihood → better model. Best model here is 2 populations with migration only from Bjäre to Östergötland

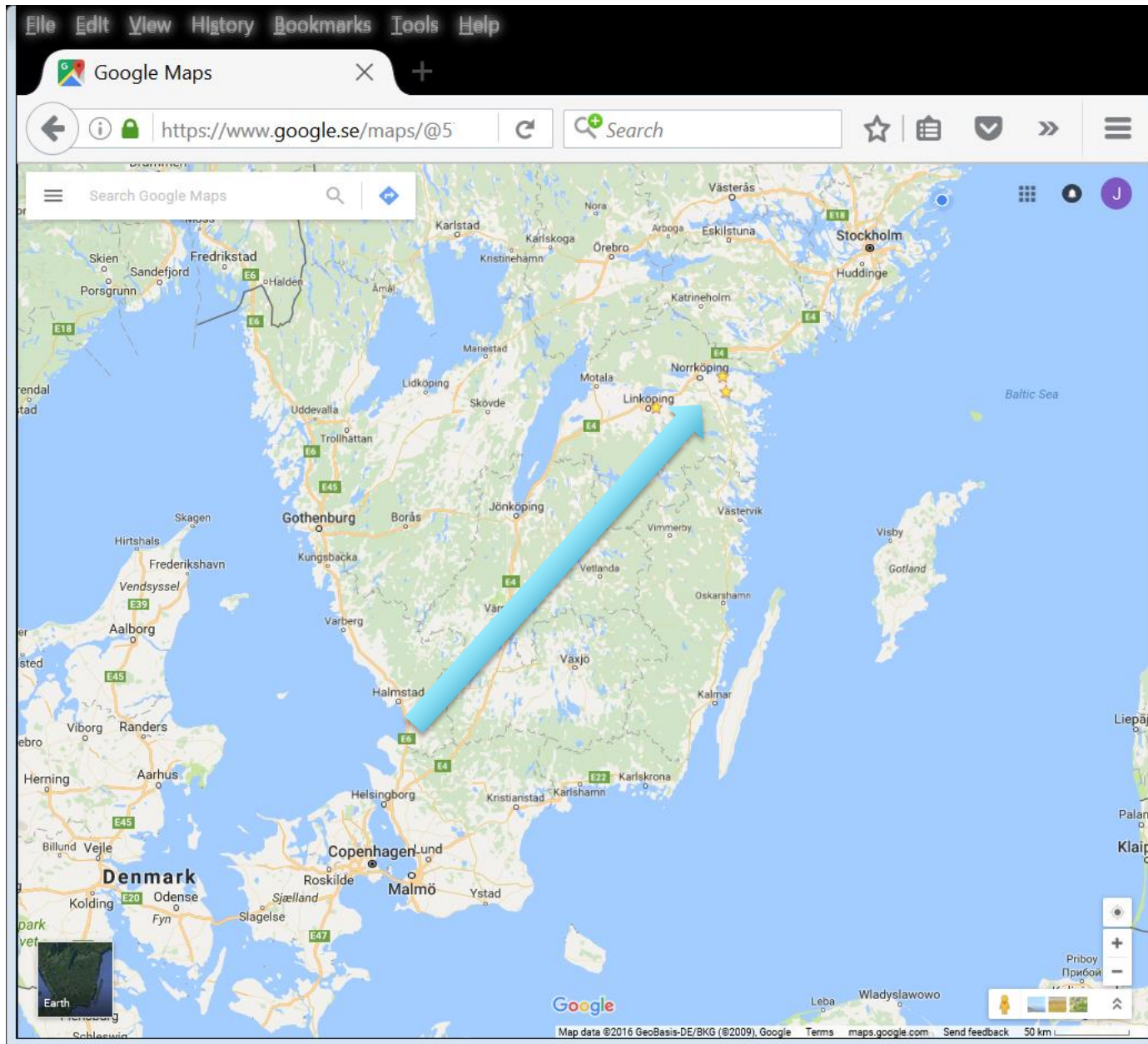
Estimates of population sizes and migration

name of the data set -- 15

Bayesian Analysis: Posterior distribution over all loci



Bjäre population
is larger than the
one in
Östergötland



Where do we stand?

- Bayesian methods are already being used in fields directly related to PRA (parameter estimation)
- Bayesian networks can be used for risk assessment either through commercial programs or open source solutions
- Bayesian methods often are a built-in part of programs used in analysis of population structure of different organisms

How can we optimize acquiring new methodology?

- We should think about how our work is related to neighboring disciplines can be a good start.
- What are our neighboring disciplines?
- Population biology?
- Epidemiology?
- What else?