# PCV Modeling Update 02/2017
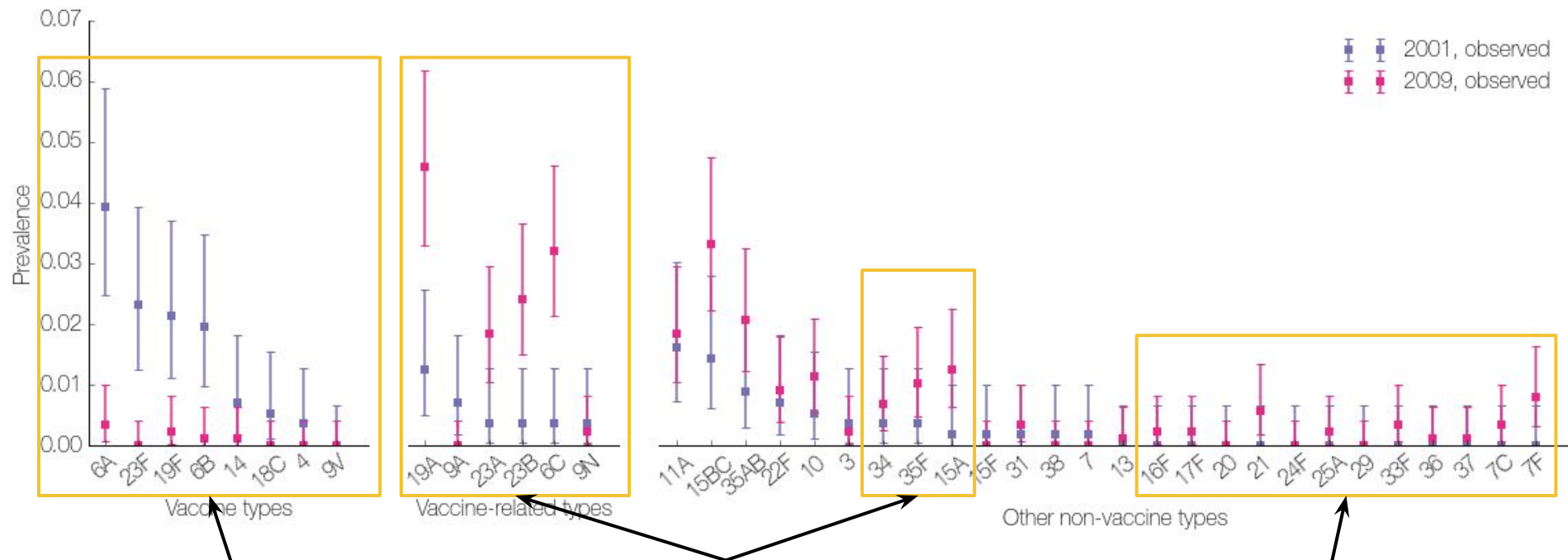
Marc Lipsitch, PI and Francisco Cai, Programmer

# Recap

- **Fitness parameters of the model are fit using peri-PCV7 (2001) data.**

- **Model reproduces peri-PCV7 serotype-specific prevalences.**

    - Expected, since number of free parameters = number of observed quantities.

- **Model had trouble reproducing PCV7-era serotype-specific prevalences.**

    - Not surprising, since there are 40 more quantities, but only 1 new parameter, vaccine efficacy.

    - 4 of 6 vaccine-related types (VRTs) were consistently underestimated (19A, 23A, 23B, 6C).

    - In general, model could not accommodate changes to the fitness ordering of serotypes.

    - Serogroup cross-immunity and shortened colonizations led to mild improvements.

    - Details in Phase II Tasks 3 and 4 Report.

# Next step

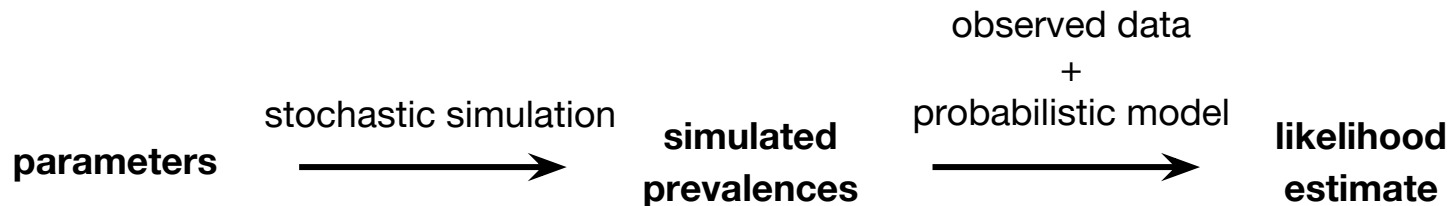- Fit model using both peri-PCV and PCV-era data.



Peri-PCV data is more informative of vaccine types

PCV7-era data is more informative for VRTs and non-vaccine types (NVTs) that expanded after vaccine introduction...

...particularly for NVTs that were not sampled in 2001.

3

# Model fitting: Before

**Goal**: Maximize expected likelihood of parameters given 2001 data.

observed data
+
stochastic simulation | probabilistic model
**parameters** → **simulated prevalences** → **likelihood estimate**

**Challenges**: The likelihood…

… is a function of many (40) parameters } **Parameter space is big.**
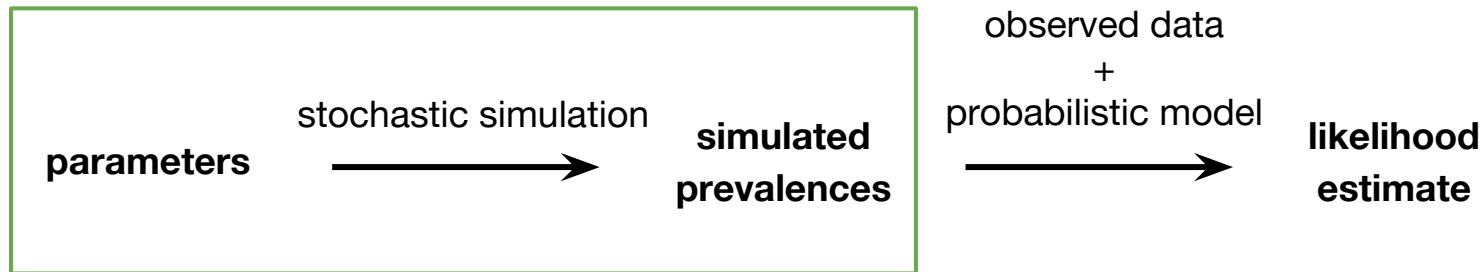
… has no closed form
→ cannot calculate gradient

**Difficult to determine which direction to explore next in parameter space.**

… cannot even be calculated directly

→ we only have a noisy estimate of it by running a stochastic simulation, which takes time (1-2 min.)

# Model fitting: Shortcuts exploited

**Goal**: Maximize expected likelihood of parameters given 2001 data.
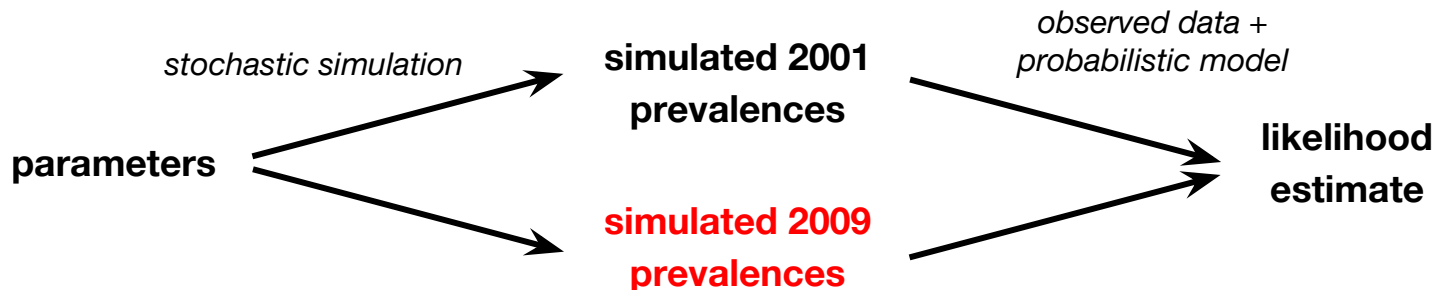


**Things we had going for us:**

1. If the simulated prevalences matched the observed prevalences, this will maximize the likelihood.
   → Focus on finding parameters that reproduce the observed prevalences (green box).

2. Monotonic relationship between fitness rank and prevalence.
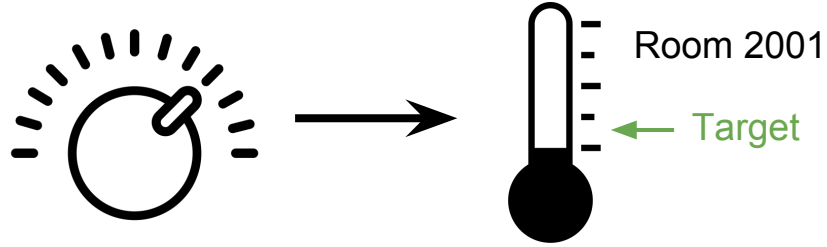   → Based on our current simulated prevalence, we know how to adjust the fitness rank.

3. Adjusting a serotype's fitness rank does not affect prevalences of other serotypes excessively.
   → Fitness parameters can be fit in parallel.

# Model Fitting: Now

**Goal**: Maximize expected likelihood of parameters given 2001 **and 2009** data.

*stochastic simulation* → **simulated 2001 prevalences**

**parameters**

**simulated 2009 prevalences**

*observed data + probabilistic model*

**likelihood estimate**

**Things we used to have going for us (now with complications):**

1.  If the simulated prevalences matched the observed prevalences, this will maximize the likelihood.
    → Focus on finding parameters that reproduce the observed prevalences.
    **The simulated prevalences now have to match the observed data at two time points.**
2.  Monotonic relationship between fitness rank and prevalence.
    → Based on our current simulated prevalence, we know how to adjust the fitness rank.
    **We now adjust one rank and hope it reproduces the observed prevalence at two time points.**
3.  Adjusting a serotype's fitness rank does not affect prevalences of other serotypes excessively.
    → Fitness parameters can be fit in parallel.
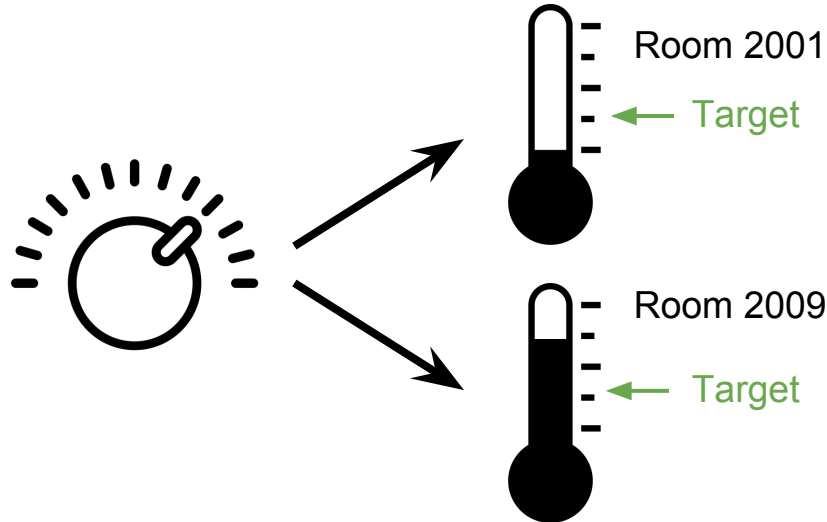    **Parameters to be fit now also include vaccine efficacy.**

# Analogy

**Before:**

Room 2001

← Target

Raise thermostat a little bit.

**Now:**

Room 2001

← Target

Room 2009

← Target

Less clear what to do.

# Initial plan

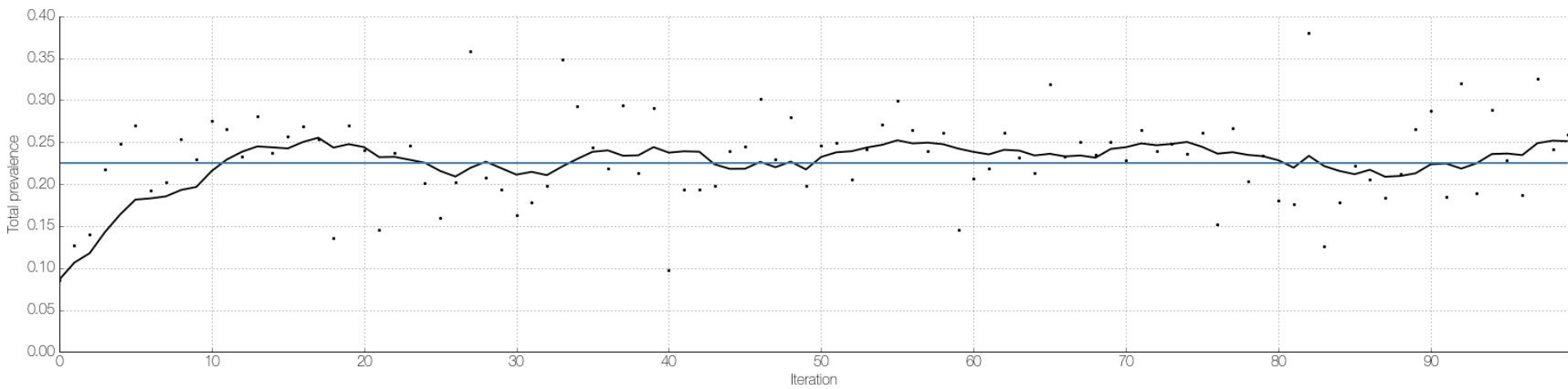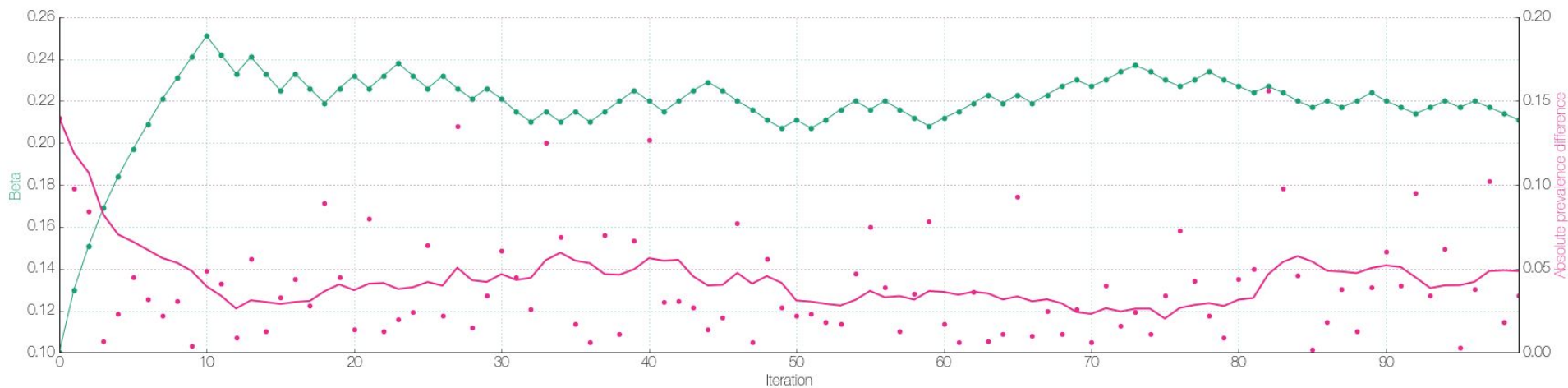**Challenges and how to address them:**

1.  Difficult to determine which direction to explore next in parameter space (no analytical gradient)

    -   **Simultaneous perturbation stochastic approximation**

        -   Simultaneously perturb all parameters to estimate all components of the gradient.

            -   ($n$ parameters requires only 2 simulations / iteration)

        -   However, simulations may be too noisy… If so, will try increasing population size or averaging results from multiple simulations.

2.  Parameter space is large.

    -   **Try fitting fitness ranks in parallel**, i.e. adjust each fitness rank as if it only affected the terms in the log-likelihood involving its serotype.

    -   Adjust the fitness ranks and the vaccine efficacy on alternate iterations.

# First milestone

- Try to reproduce previous results, i.e. fit only peri-PCV7 data using new algorithm

- As before, use absolute prevalence error, rather than log-likelihood, to quantify how well we are doing.
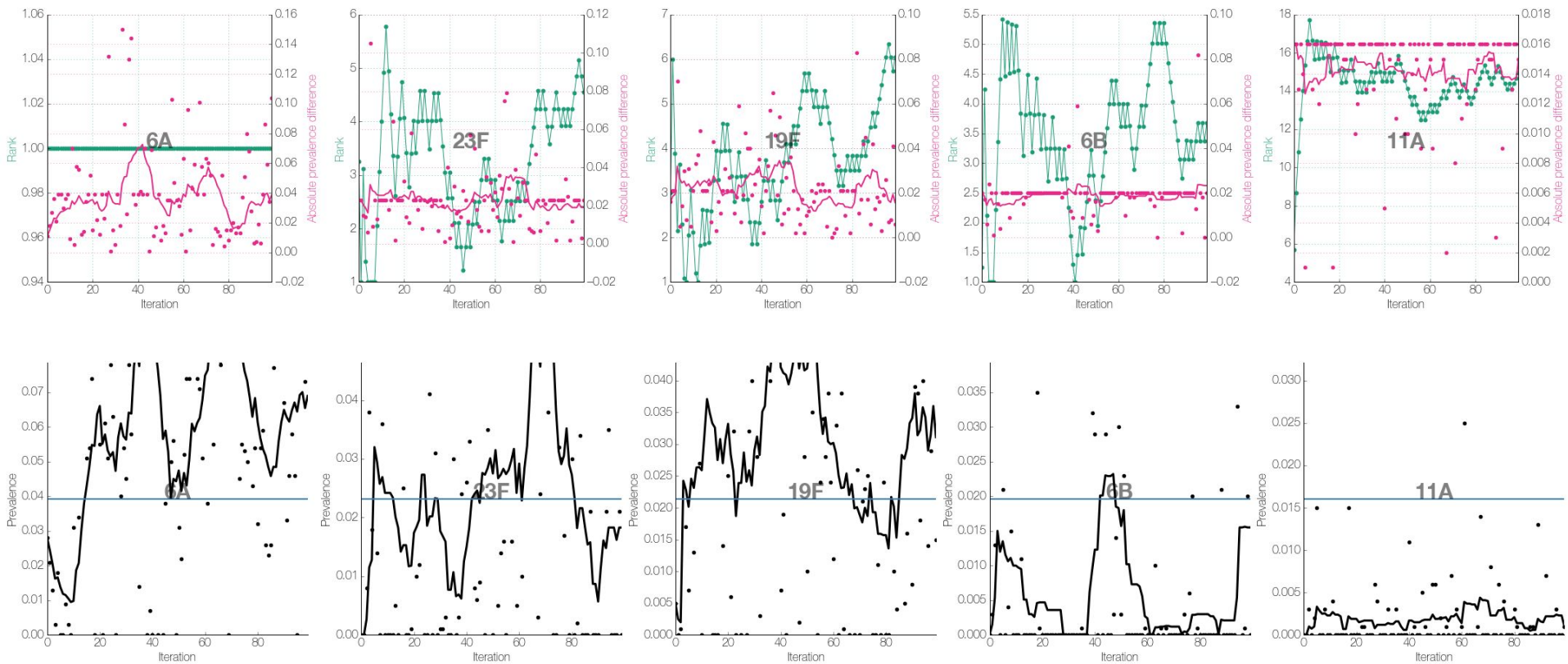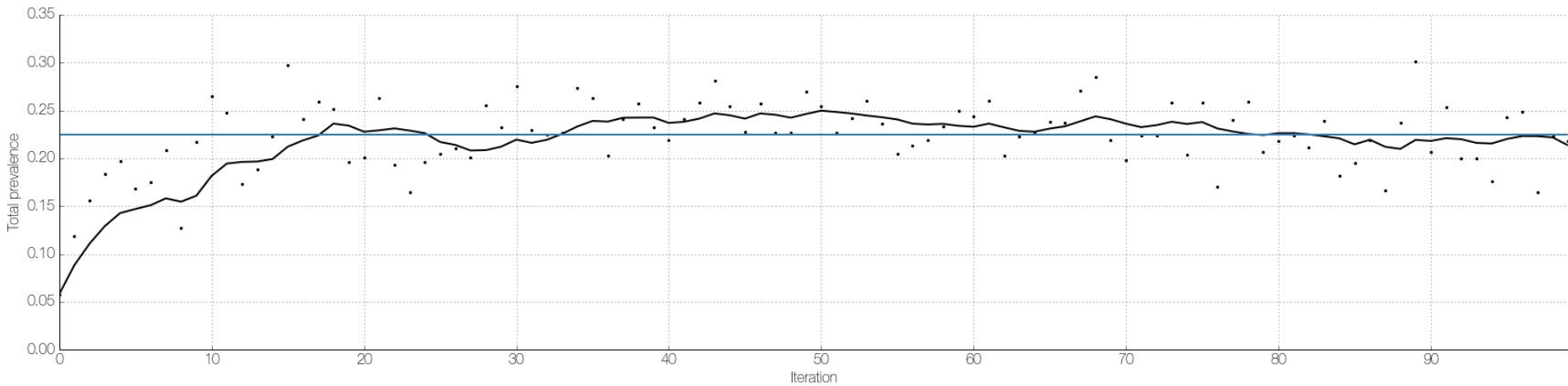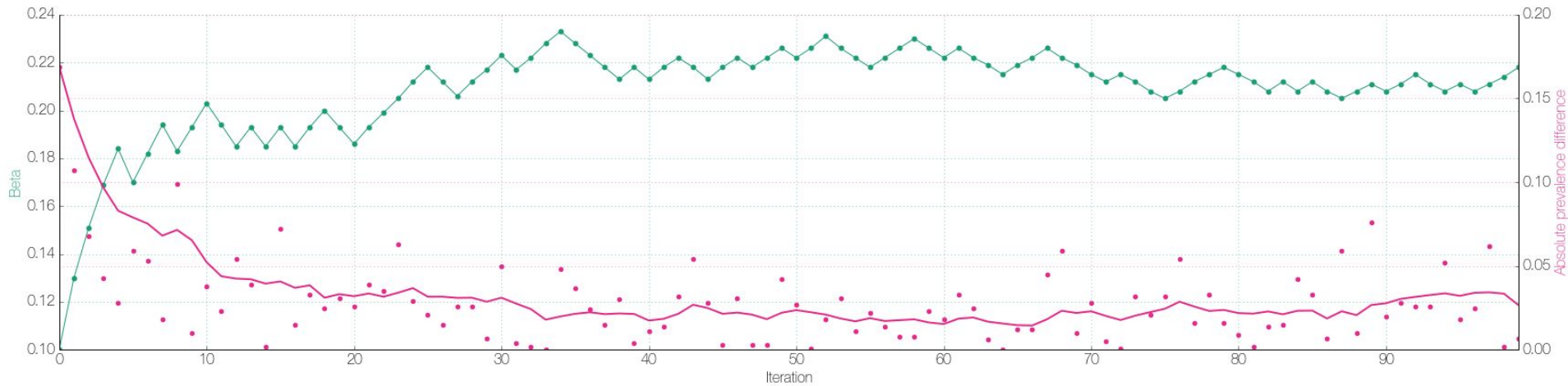
# First test

# First test

2.5x population size

# 2.5x population size

# Averaging 5 simulations

# Averaging 5 simulations

# Next steps

- Immediate challenge seems to be stochastic noise affecting the fitting of individual serotype parameters.

- How do we reduce noise, without increasing computational time too much?

- Try simple ideas:

    - Use a combination of larger population sizes and averaging more simulations.

    - Instead of averaging, consider the *distribution* of results in each simulation set.

    - Perturb one parameter at a time for gradient estimation (runtime would increase dramatically, however).

- Look for a more principled approach to statistical inference in individual-based models

    - Currently exploring this as a possible project, with Professor Pierre Jacob at Harvard (Statistics)