



Institute for
Health Metrics
and Evaluation

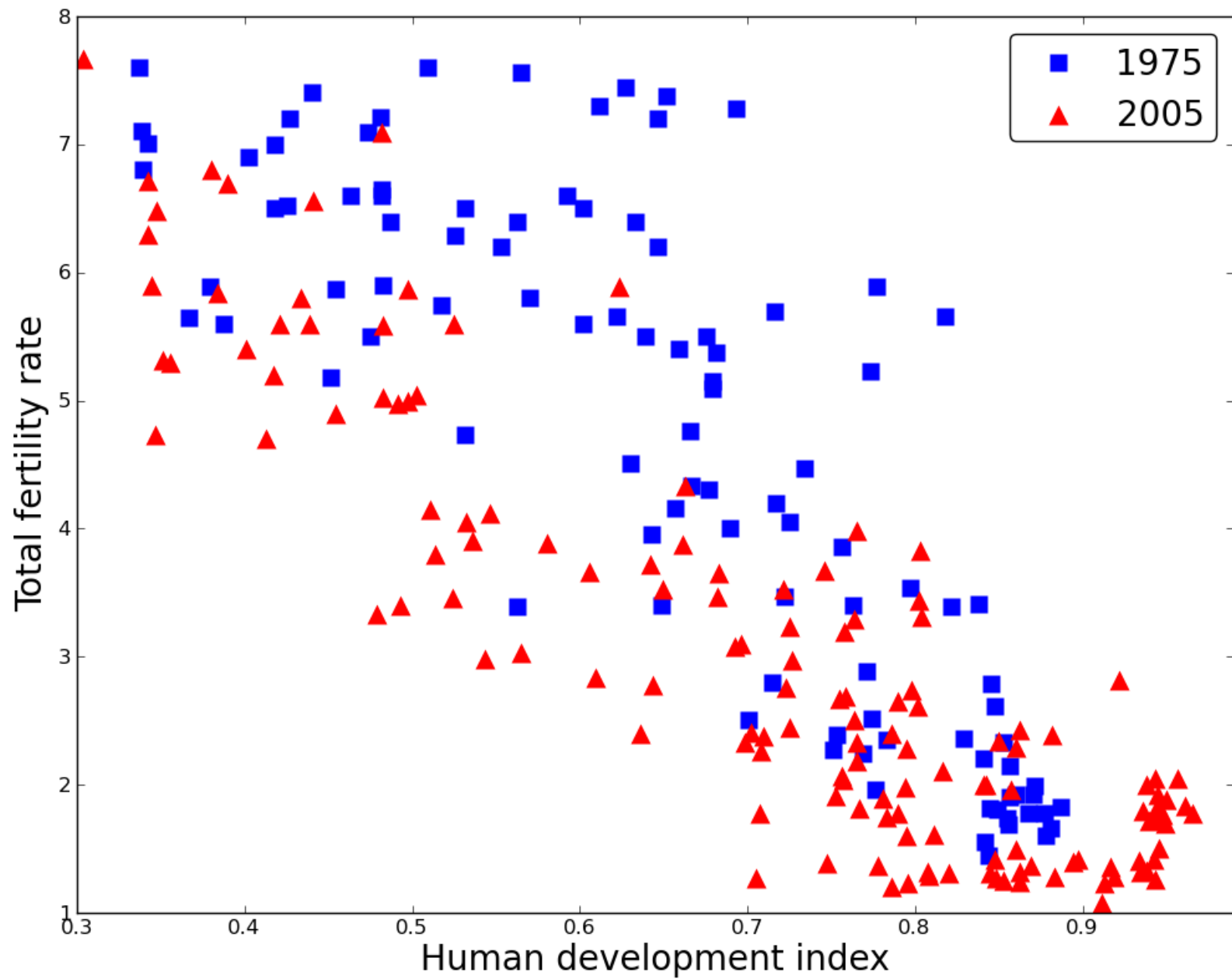
PyMC by Example: Human Fertility

SciPy 2011

Abraham D Flaxman

Institute for Health Metrics and Evaluation

UNIVERSITY OF WASHINGTON



Access

To read this story in full you will need to login or make a payment (see right).

[nature.com](#) > [Journal home](#) > [Table of Contents](#)

Letter

Nature **460**, 741–743 (6 August 2009) | doi:10.1038/nature08230; Received 1 April 2009; Accepted 17 June 2009

Advances in development reverse fertility declines

Mikko Myrskylä¹, Hans-Peter Kohler¹ & Francesco C. Billari²

1. Population Studies Center, University of Pennsylvania, 3718 Locust Walk, Philadelphia, Pennsylvania 19104, USA

2. DONDENA "Carlo F. Dondena" Centre for Research on Social Dynamics, Department of Decision Sciences and IGIER, Università Bocconi, via Röntgen 1, 20136 Milan, Italy

Correspondence to: Hans-Peter Kohler¹ Correspondence and requests for materials should be addressed to H.-P.K.

(Email: hpkohler@pop.upenn.edu).

During the twentieth century, the global population has gone through unprecedented increases in economic and social development that coincided with substantial declines in human fertility and population growth rates^{1, 2}. The negative

▲ Top

ARTICLE LINKS

- ▶ [Figures and tables](#)
- ▶ [Supplementary info](#)

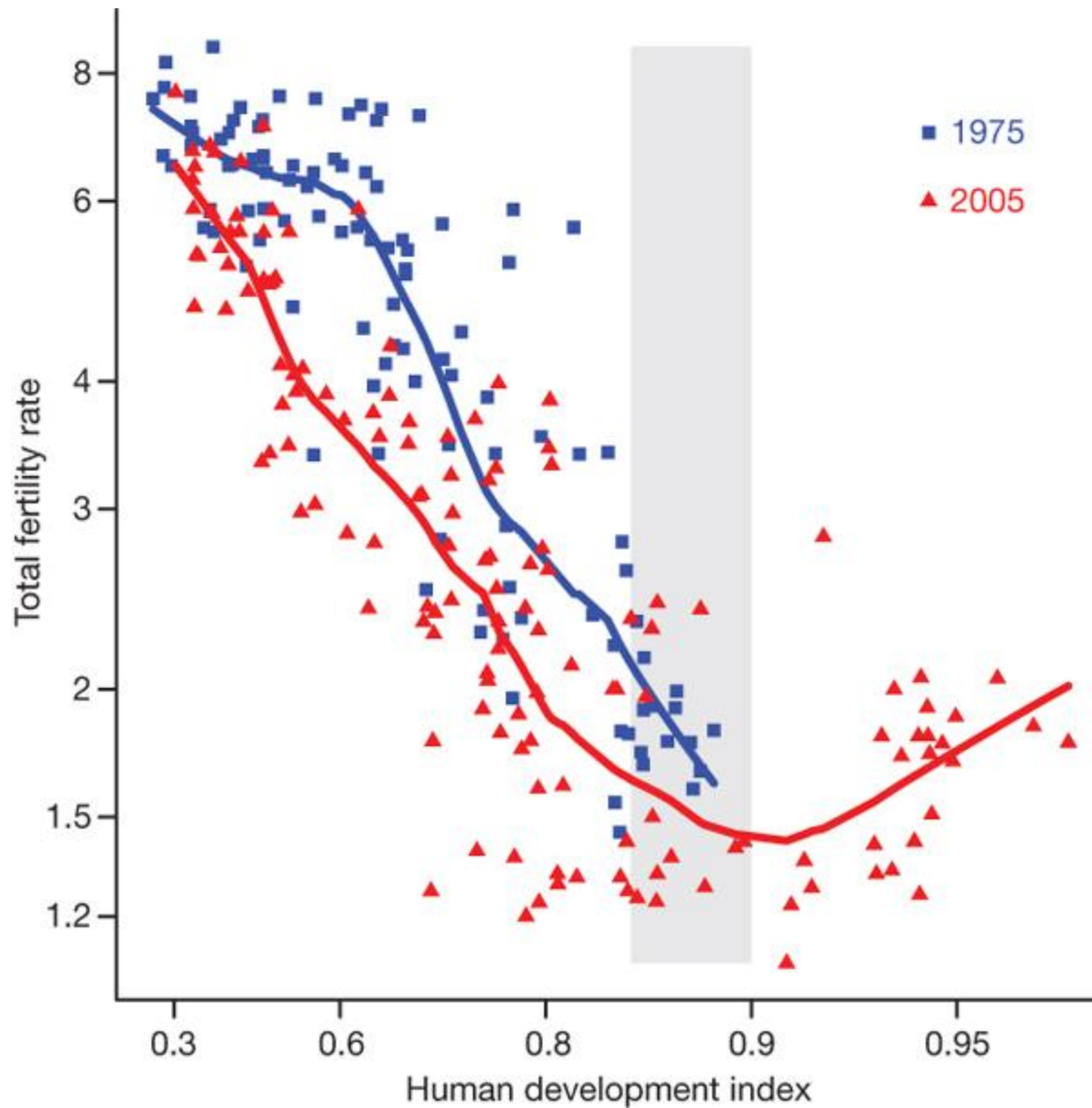
SEE ALSO

- ▶ [News and Views by Tuljapurkar](#)
- ▶ [Editor's Summary](#)

ARTICLE TOOLS

-  [Send to a friend](#)
-  [Export citation](#)
-  [Export references](#)
-  [Rights and permissions](#)
-  [Order commercial reprints](#)
-  [Bookmark in Connotea](#)

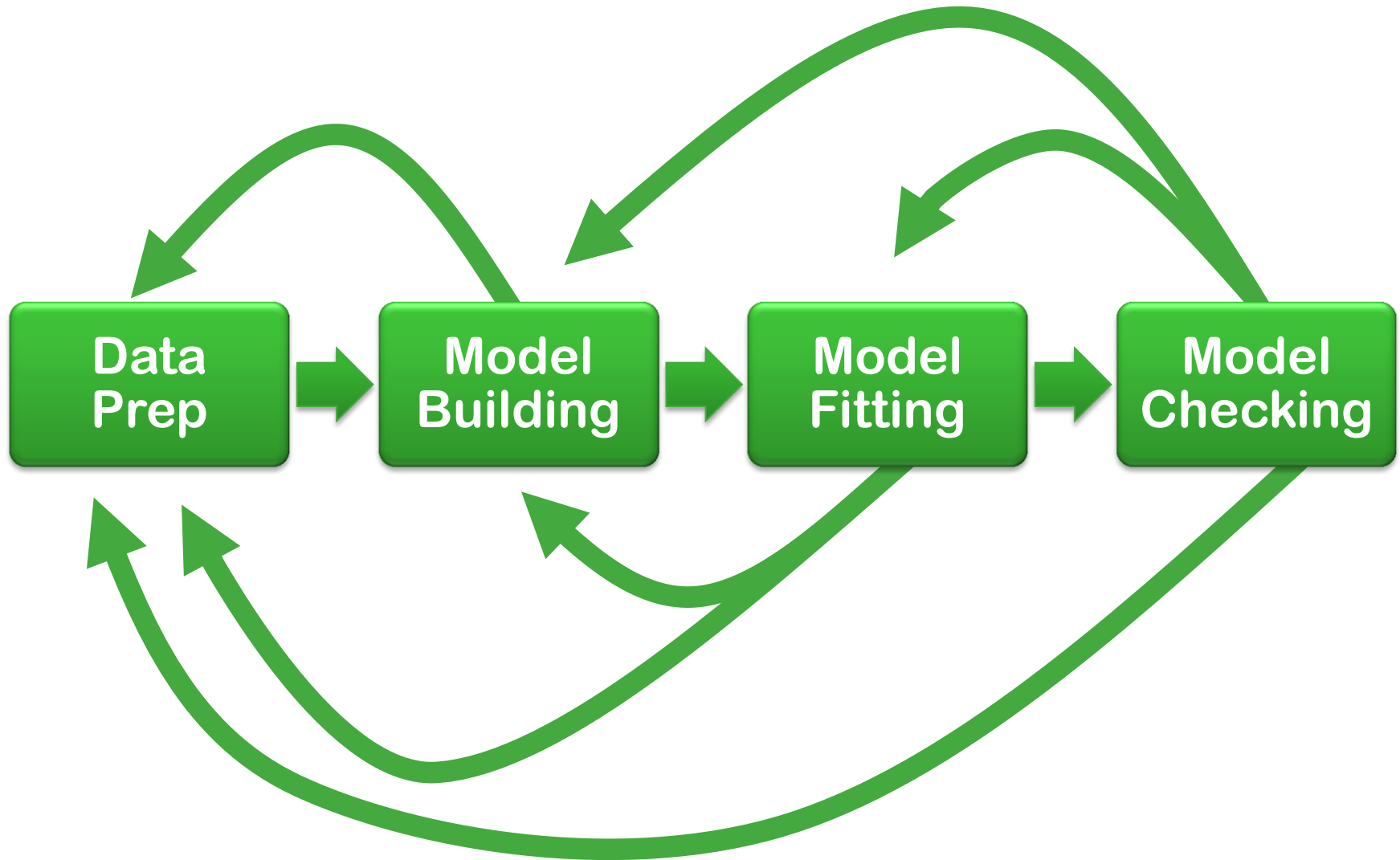
Cross-sectional relationship between TFR and HDI in 1975 and 2005.



Learning Objective

Know enough PyMC to check this out


A Generic Template for Bayesian data analysis



pymc-project-template / src

name	age	mess
..		
 data.py	May 24, 2011	init
 graphics.py	May 24, 2011	init
 models.py	May 24, 2011	init
 tests.py	May 24, 2011	init

The Data

 100755 | 145 lines (144 sloc) | 94.991 kb

```
1 country,HDI.1975,HDI.1976,HDI.1977,HDI.1978,HDI.1979,HDI.1980,
2 Albania,,,,,,,,0.731273056,0.734932246,0.737145554,0.738451574,0
3 Algeria,0.565066871,0.57096611,0.575103358,0.581994626,0.58797
4 Angola,0.427032369,0.427252059,0.4274633,0.428155305,0.4288394
5 Argentina,0.796895676,0.796294282,0.800532019,0.798749194,0.80
6 Armenia,0.700787712,0.702149028,0.703451552,0.705567048,0.7096
7 Australia,0.856418407,0.857742049,0.859507393,0.86295765,0.865
8 Austria,0.844760239,0.849874426,0.854348125,0.854713814,0.8594
9 Azerbaijan,,,,,,,,,,,,,0.770722523,0.76666317,0.745701042,0
10 Bahrain,0.773655346,0.774279264,0.774879186,0.779222144,0.7835
11 Bangladesh,0.387413871,0.389184115,0.389373185,0.395334908,0.4
12 Belarus,,,,,,,,,,,,,0.7830714,0.780204413,0.773003586,0.763
13 Belgium,0.854899715,0.857868123,0.858153017,0.861384937,0.8643
14 Belize,0.646159734,0.644969023,0.647105514,0.651537112,0.65573
15 Benin,0.338790582,0.337774288,0.338946694,0.339810674,0.343432
16 Bolivia,0.570203508,0.571454562,0.572813625,0.576806088,0.5790
```



```
import pylab as pl
```

src/data.py

```
orig = pl.csv2rec('nature08230-s2.csv')
```

```
country = []; year = []; hdi = []; tfr = []
for row in orig:
    for y in range(1975, 2006):
        if pl.isnan(row['hdi%d'%y]) \
            or pl.isnan(row['tfr%d'%y]):
            continue
        country.append(row['country'])
        year.append(y)
        hdi.append(row['hdi%d'%y])
        tfr.append(row['tfr%d'%y])
```

```
all = pl.np.core.rec.fromarrays([country, year, hdi, tfr],
                                names=['country', 'year', 'hdi', 'tfr'])
```

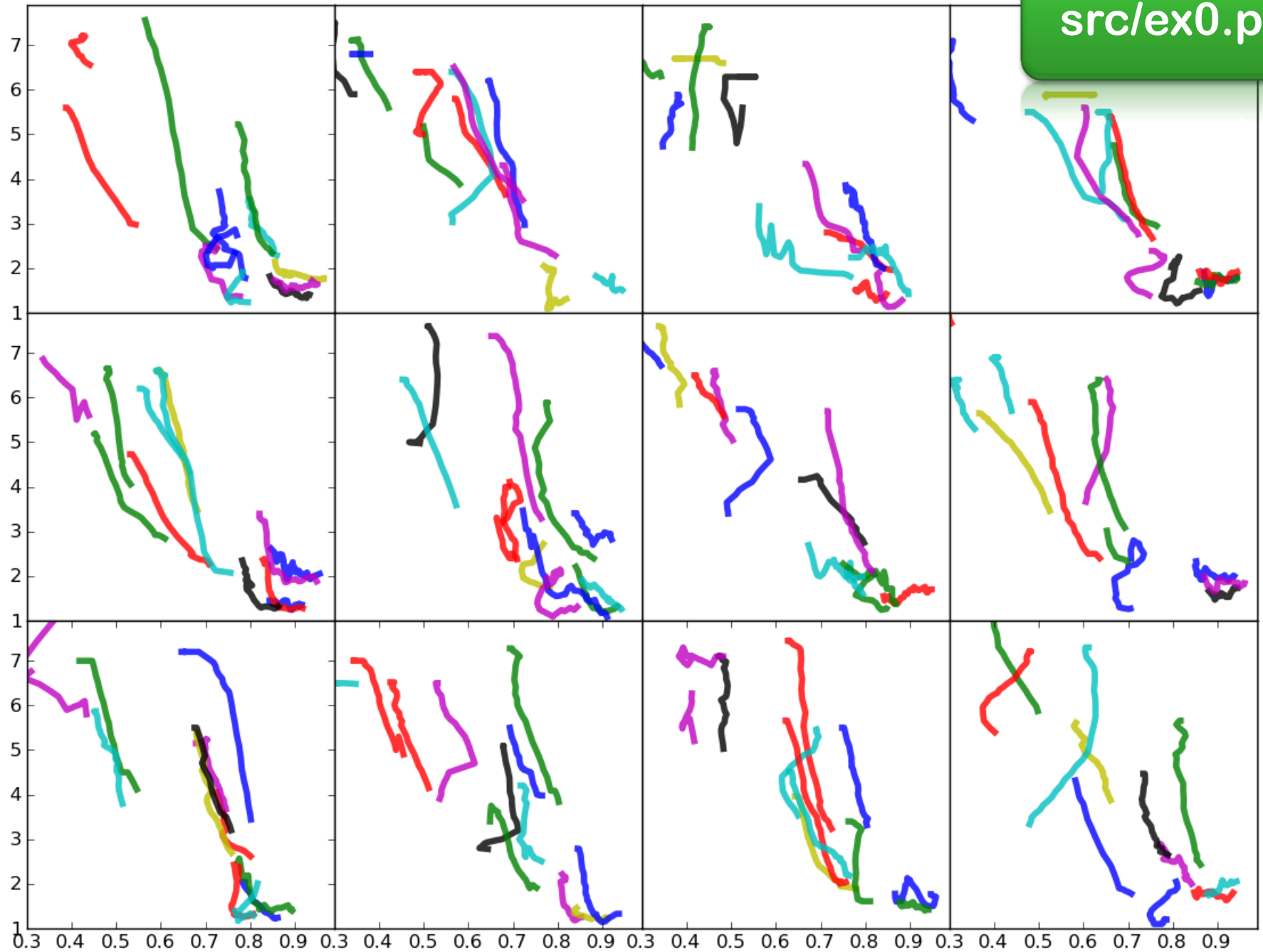
```
hdi2005 = all.hdi[all.year==2005]
tfr2005 = all.tfr[all.year==2005]
```

```
def plot_each_country(axis_bounds):
    years = range(1975, 2006)
    for i, c in enumerate(pl.unique(data.all.country)):
        pl.subplot(3, 4, i/12+1)
        pl.plot(data.all.hdi[data.all.country==c],
                data.all.tfr[data.all.country==c],
                linewidth=4, alpha=.8)
        pl.axis(axis_bounds)

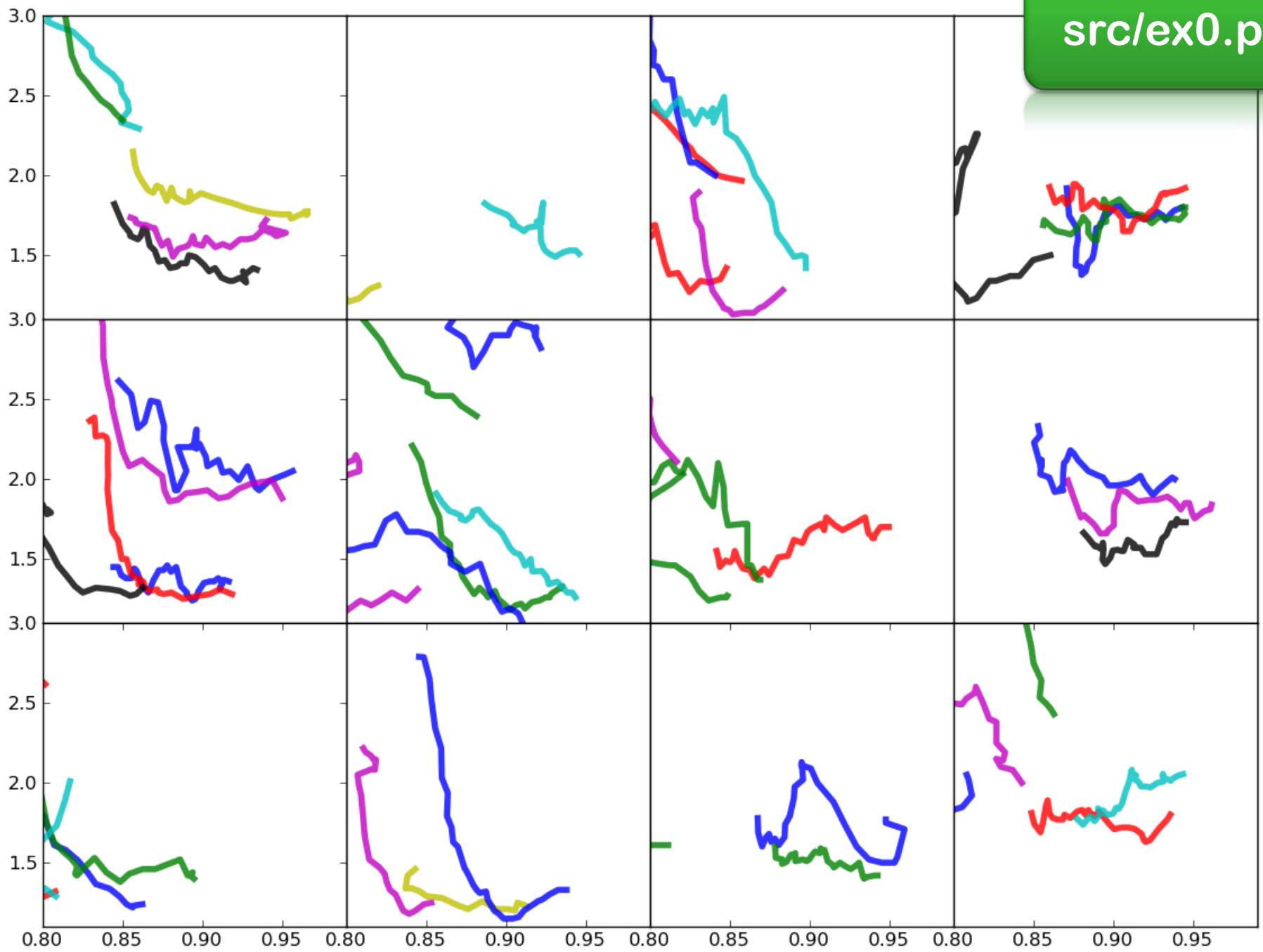
    for r in range(3):
        for c in range(4):
            subplot(3, 4, r*4+c+1)
            if r != 2:
                pl.xticks([])
            if c != 0:
                pl.yticks([])

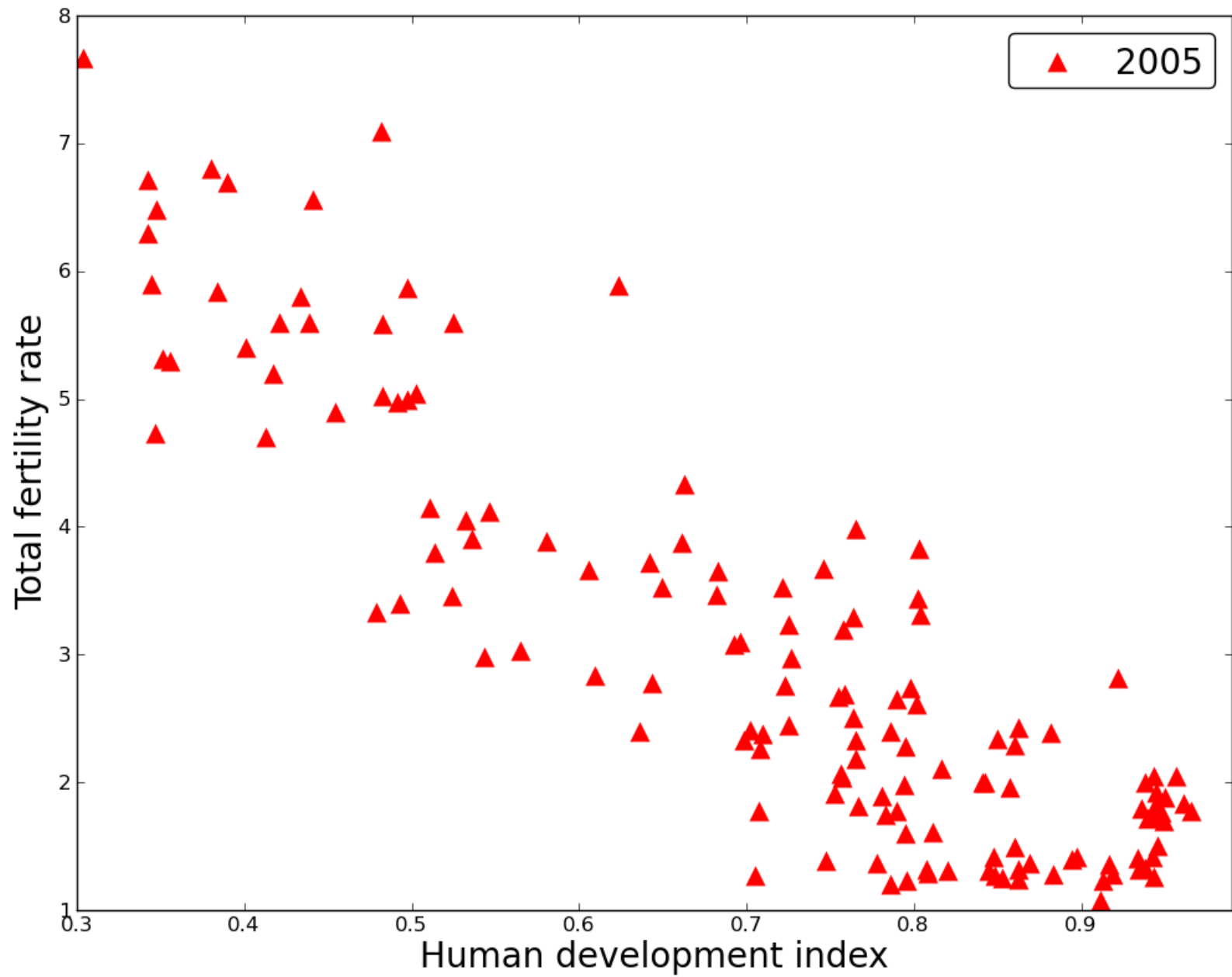
    pl.subplots_adjust(.05, .05, .95, .95, 0, 0)
```

src/ex0.py



src/ex0.py





Simple model

$$\text{TFR}_i = \beta_0 + \beta_1 \text{HDI}_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

```
1 def linear():
2     beta = mc.Uninformative('beta', value=[0., 0.])
3     sigma = mc.Uniform('sigma', lower=0, upper=100, value=1.)
4
5     @mc.deterministic
6     def y_mean(beta=beta, X=data.hdi2005):
7         return beta[0] + beta[1]*X
8
9     y_obs = mc.Normal('y_obs', value=data.tfr2005,
10                       mu=y_mean, tau=sigma**-2,
11                       observed=True)
12
13     return vars()
```

```
In [1]: import models
```

```
In [2]: vars = models.linear()
```

```
In [3]: vars
```

```
Out[3]:
```

```
{'beta': <pymc.distributions.Uninformative 'beta' at 0xae9b60f0>,  
 'sigma': <pymc.distributions.Uniform 'sigma' at 0xae9b60f0>,  
 'y_mean': <pymc.PyMCObjects.Deterministic 'y_mean' at 0xae9b60f0>,  
 'y_obs': <pymc.distributions.Normal 'y_obs' at 0xae9b60f0>}
```



```
In [10]: vars['beta'].value
```

```
Out[10]: array([ 0.,  0.])
```

```
In [11]: vars['y_obs'].logp
```

```
Out[11]: -947.69862197605573
```

```
1 def fit_linear():  
2     vars = linear()  
3  
4     mc.MAP(vars).fit(method='fmin_powell')  
5  
6     m = mc.MCMC(vars)  
7     m.sample(iter=10000, burn=5000, thin=5)  
8     return m
```

```
In [1]: import models
```

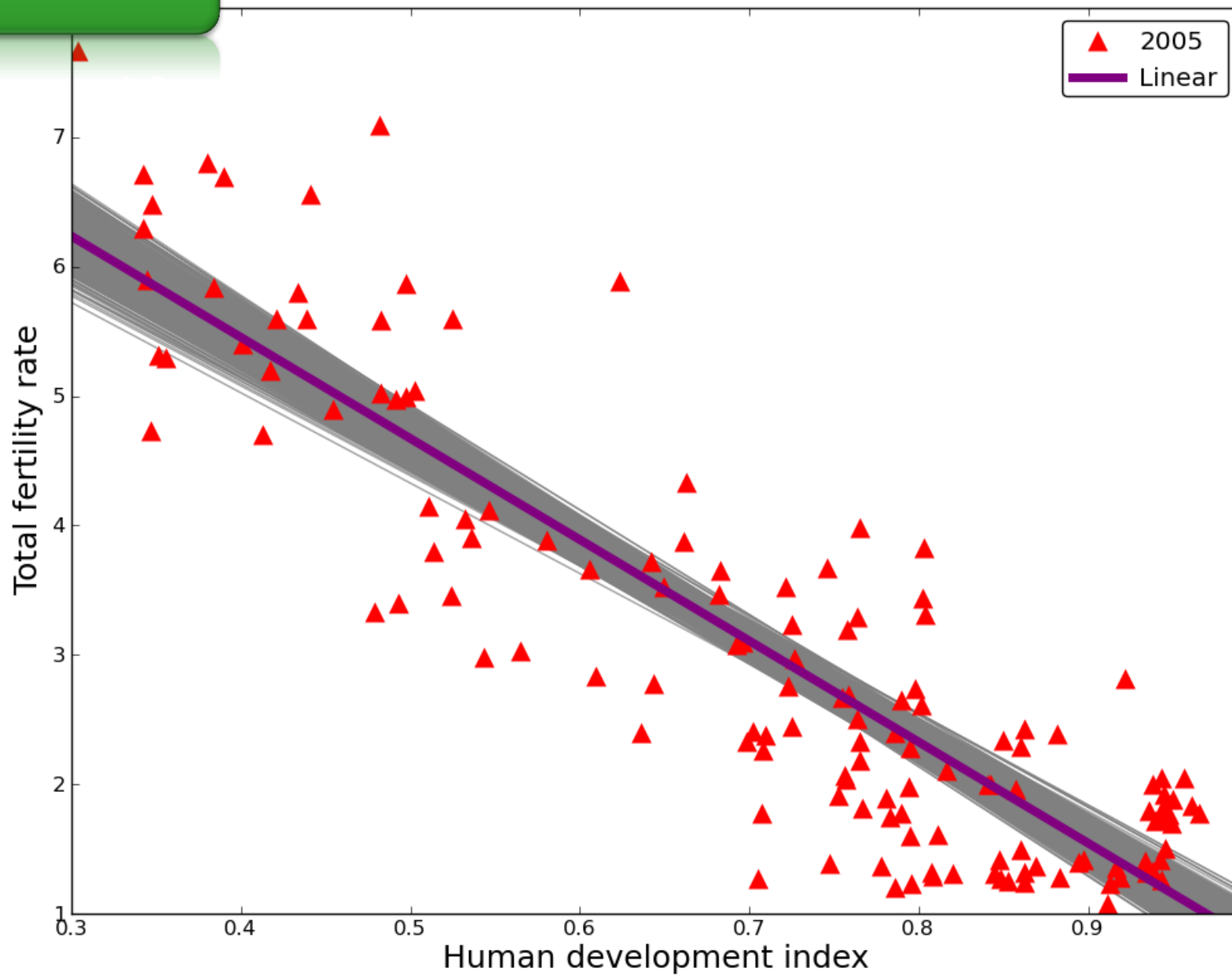
```
In [2]: m = models.fit_linear()  
Sampling: 100% [000] Iterations: 10000
```

```
In [3]: import graphics
```

```
In [4]: graphics.plot_all_data()
```

```
In [5]: graphics.plot_linear_model(m)
```

src/ex1.py



Less simple model

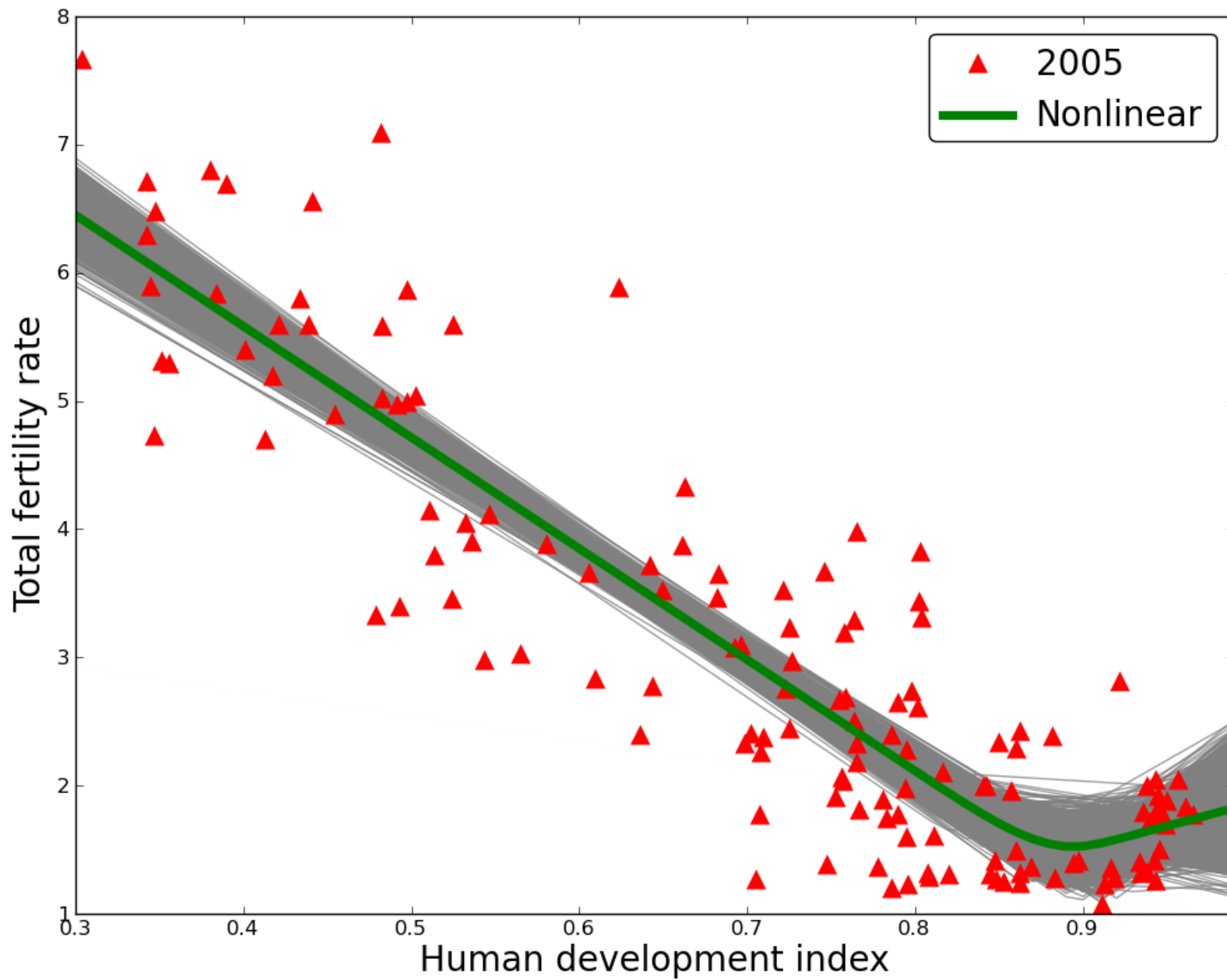
$$\text{TFR}_i = \beta_0 + \beta_1 \text{HDI}_i + \beta_2 (\text{HDI}_i - \gamma)^+ + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

```

1 def nonlinear():
2     beta = mc.Uninformative('beta', value=[0., 0., 0.])
3     gamma = mc.Normal('gamma', mu=.9, tau=.05** -2)
4     sigma = mc.Uniform('sigma', lower=0, upper=100, value=1.)
5
6     @mc.deterministic
7     def y_mean(beta=beta, gamma=gamma, X=data.hdi2005):
8         return beta[0] + beta[1]*X \
9             + beta[2]*pl.maximum(0., X-gamma)
10
11     y_obs = mc.Normal('y_obs', value=data.tfr2005,
12                       mu=y_mean, tau=sigma** -2,
13                       observed=True)
14
15     return vars()

```



Your turn to do an extension

Suggestions:

- **Alternative Data**
 - $\text{logit}(\text{HDI})$ vs $\log(\text{TFR})$
 - Different time period than 2005
 - Just one country over time
- **Alternative Priors**
 - More informative betas
 - Less informative gamma
- **Alternative Models**
 - Quadratic
 - “Country-level Random Effects”
 - Others?



Institute for
Health Metrics
and Evaluation

PyMC by Example: Human Fertility II

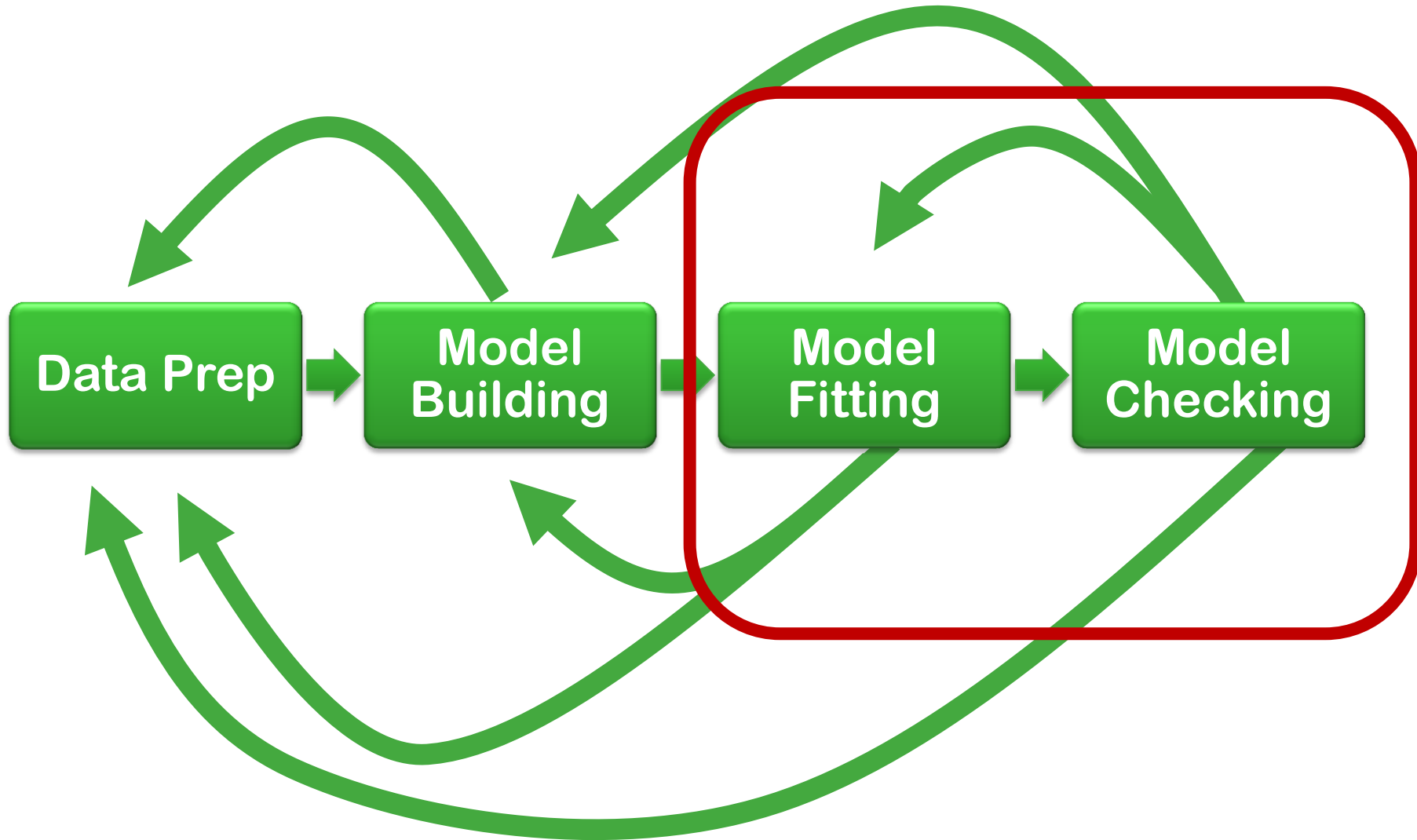
SciPy 2011

Abraham D Flaxman

Institute for Health Metrics and Evaluation

UNIVERSITY OF WASHINGTON

A Generic Template for Bayesian data analysis



$$\text{TFR}_i = \beta_0 + \beta_1 \text{HDI}_i +$$

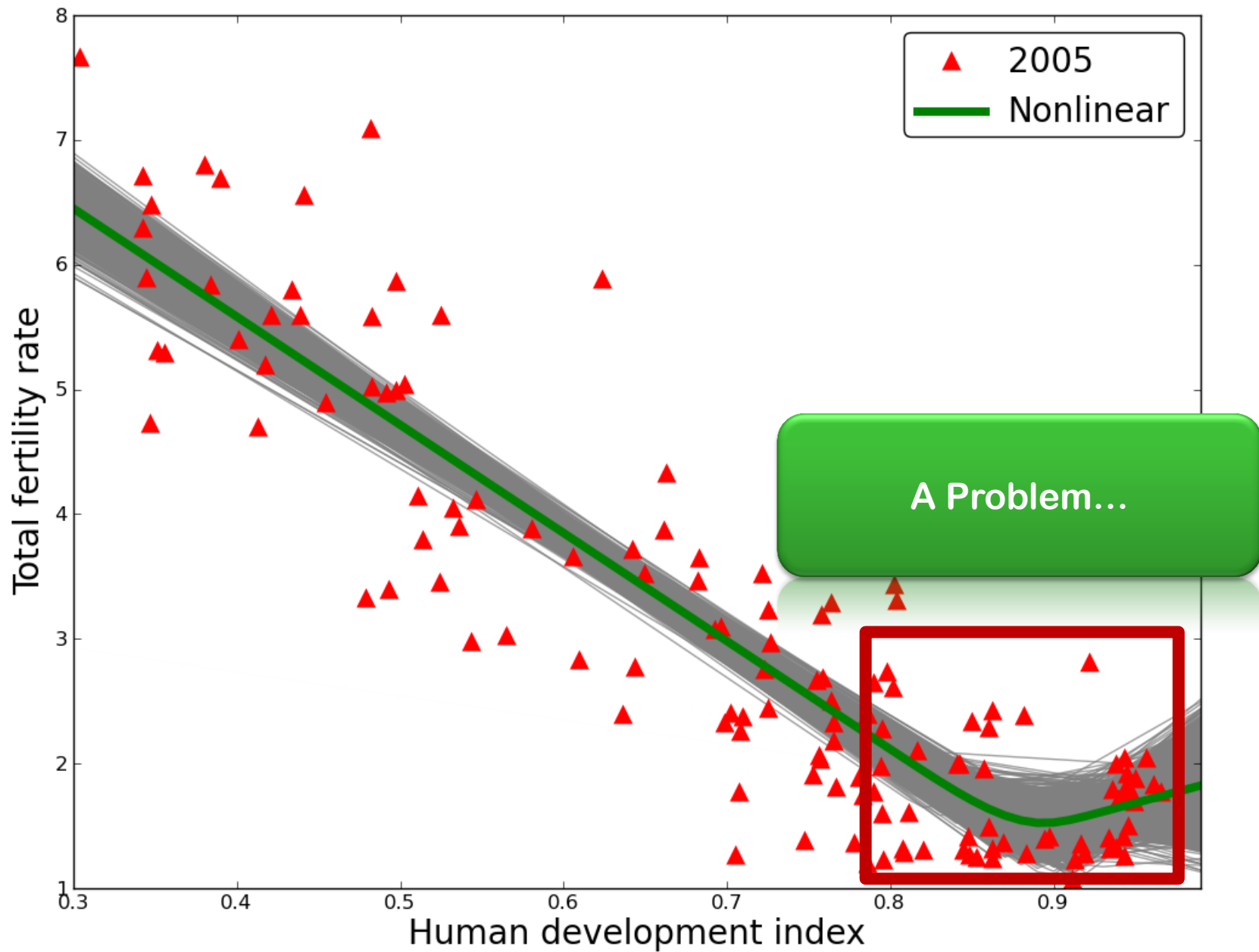
$$\beta_2(\text{HDI}_i - \gamma)^+ + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

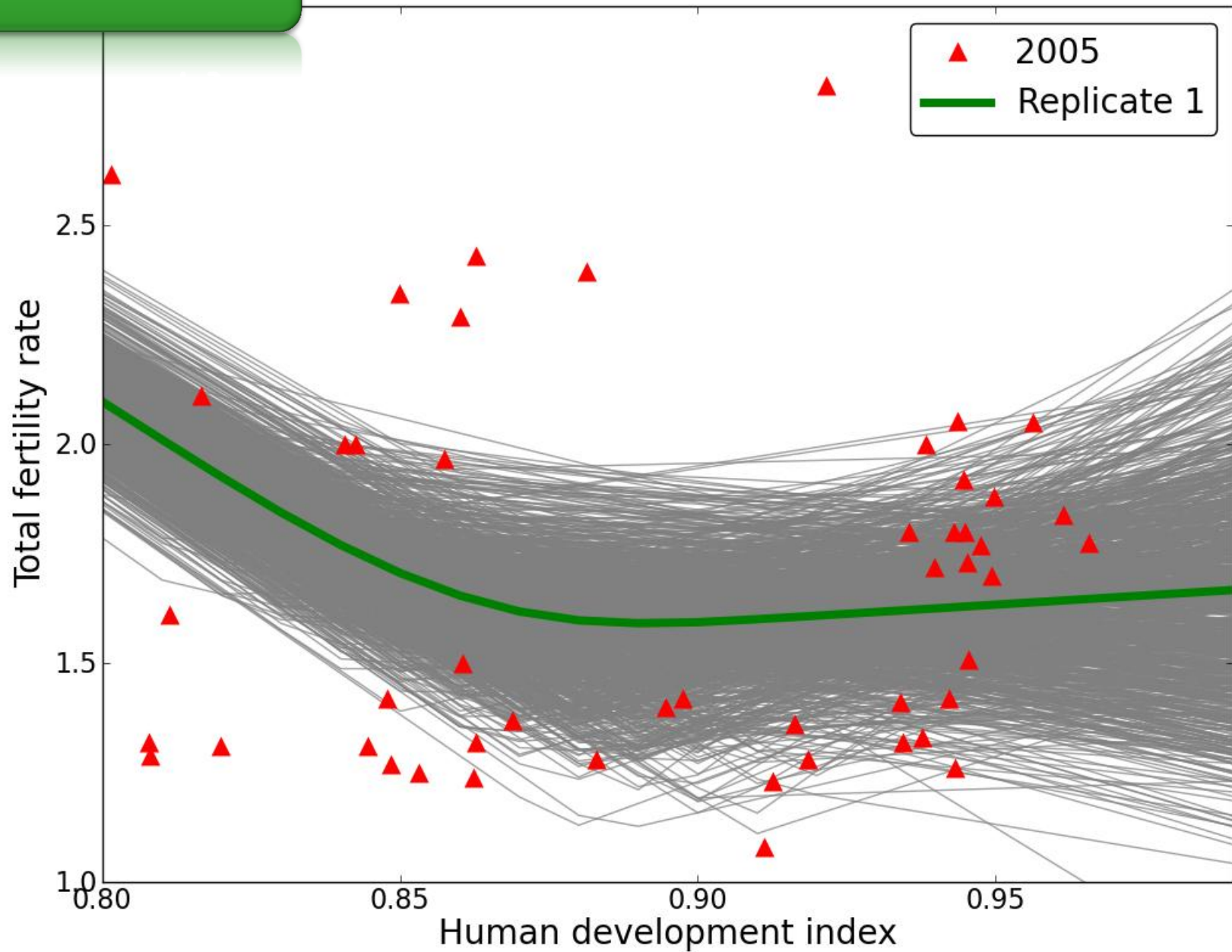


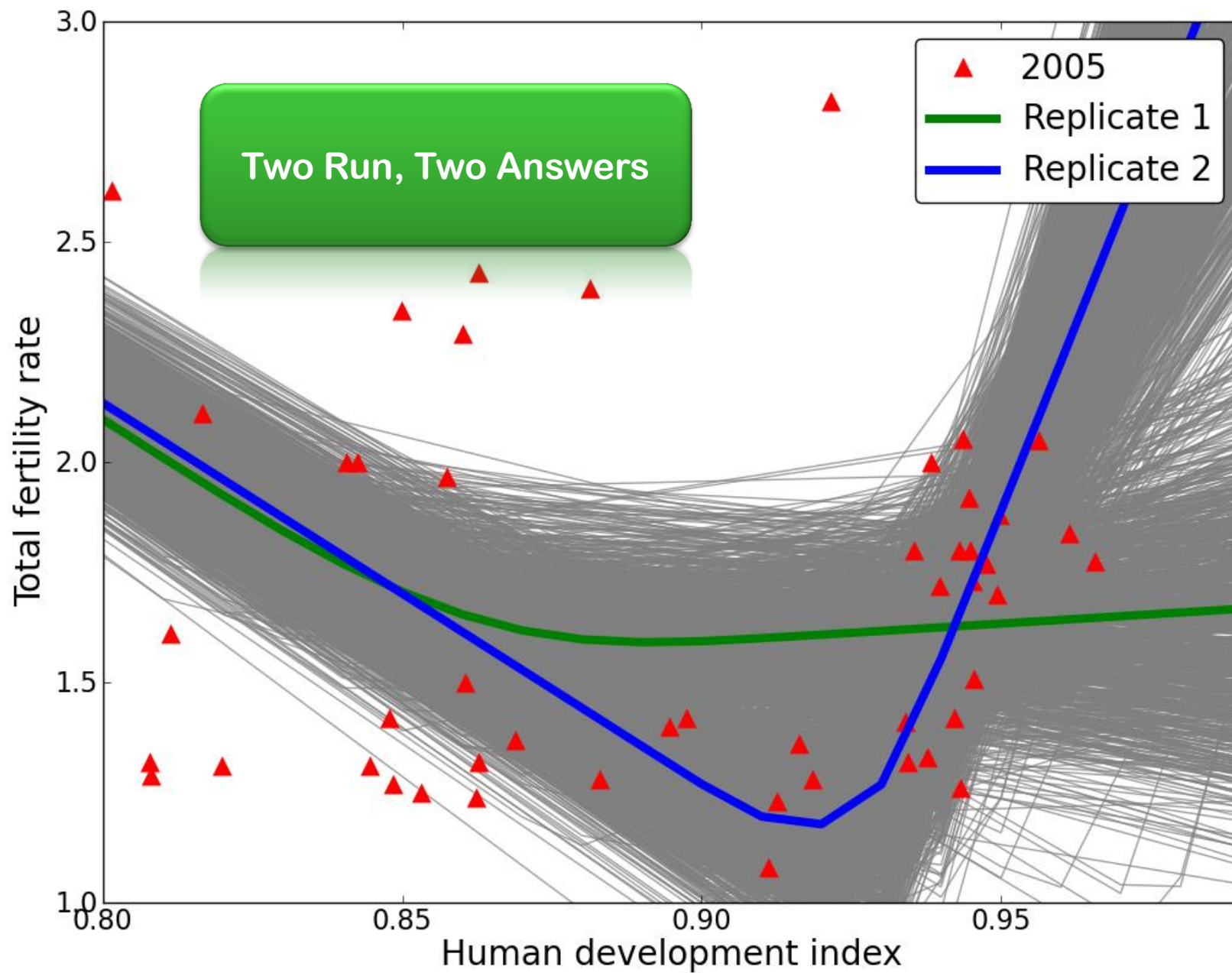
```
@mc.deterministic
def y_mean(beta=beta, gamma=gamma, X=data.hdi):
    return beta[0] + beta[1]*X \
           + beta[2]*pl.maximum(0., X-gamma)

y_obs = mc.Normal('y_obs', value=data.tfr,
                  mu=y_mean, tau=sigma**-2,
                  observed=True)
```



src/ex3.py



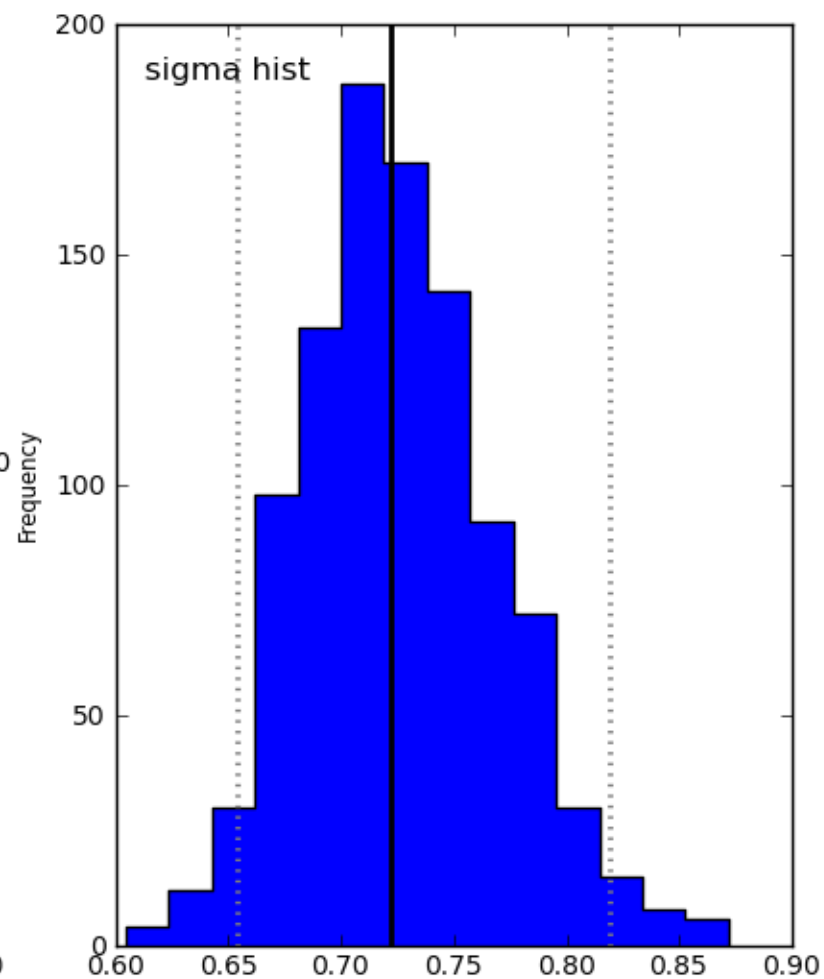
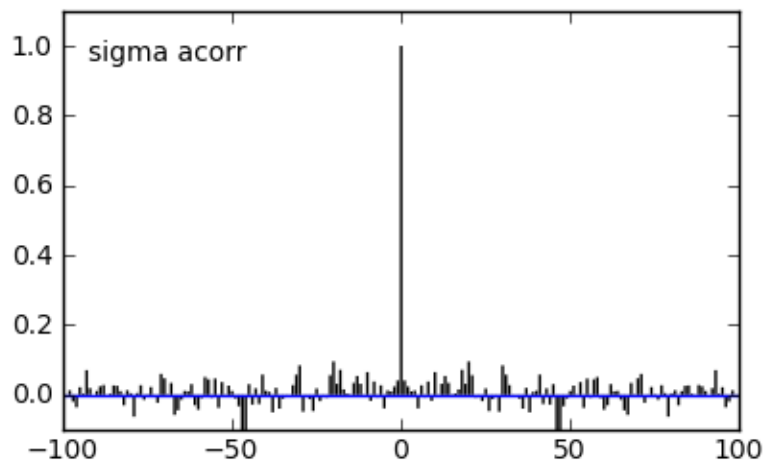
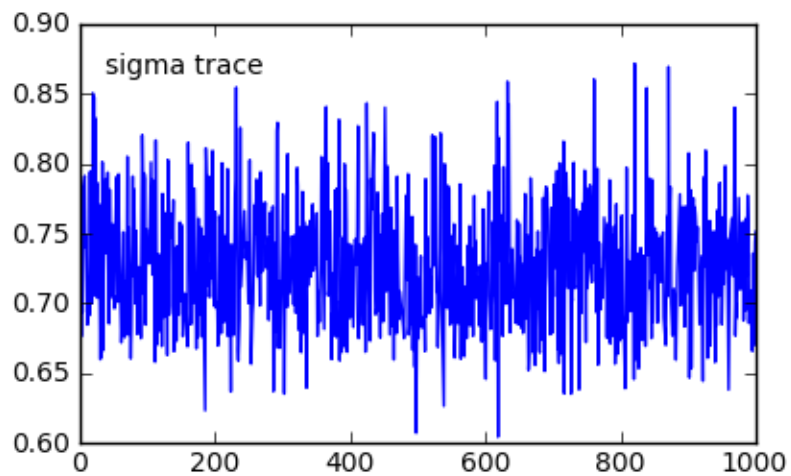


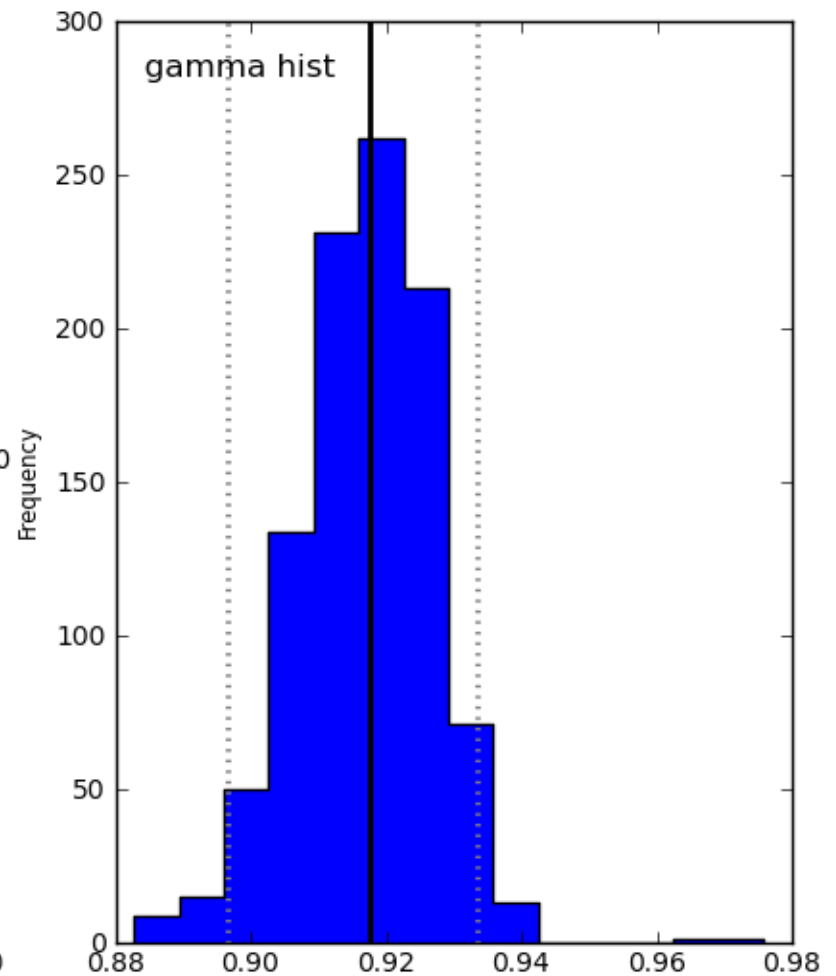
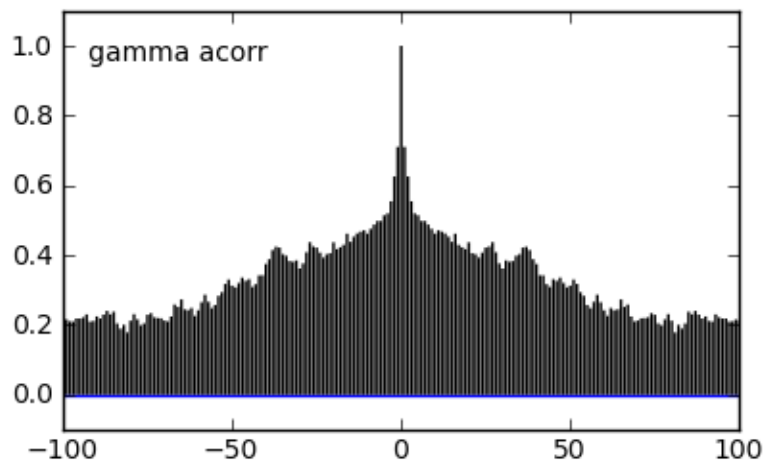
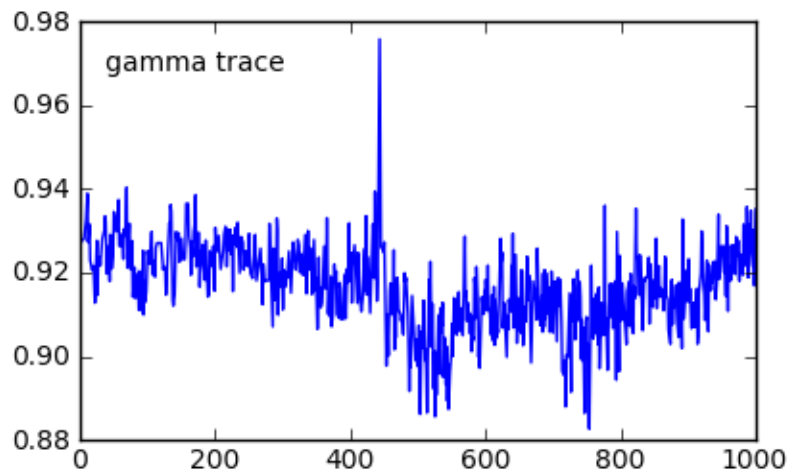
Model Fitting in PyMC

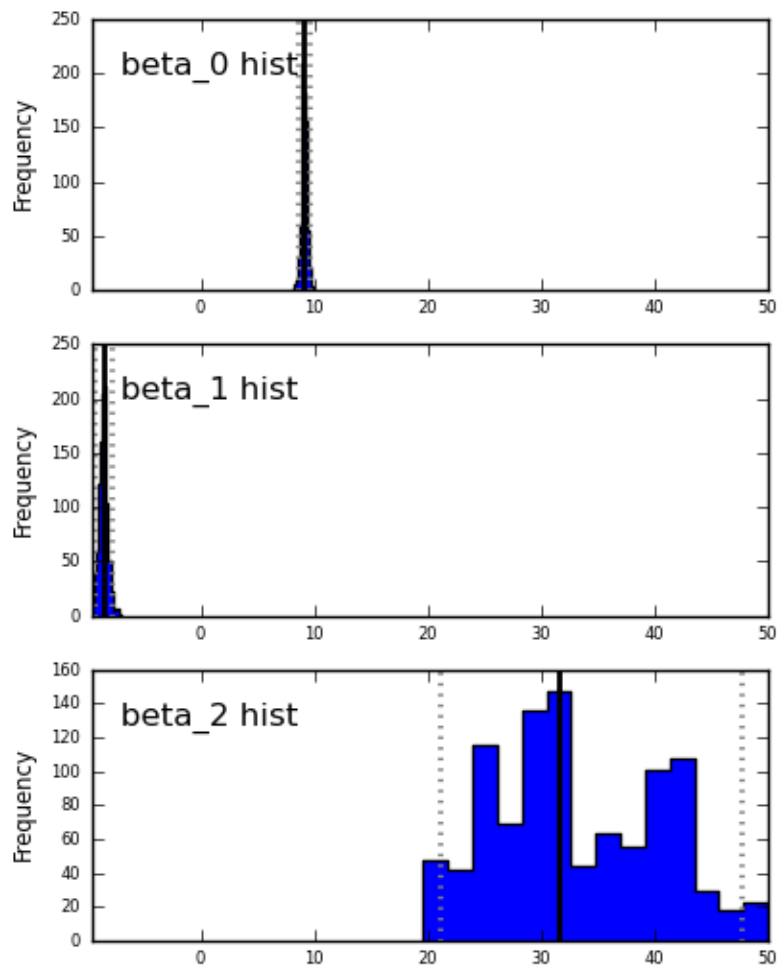
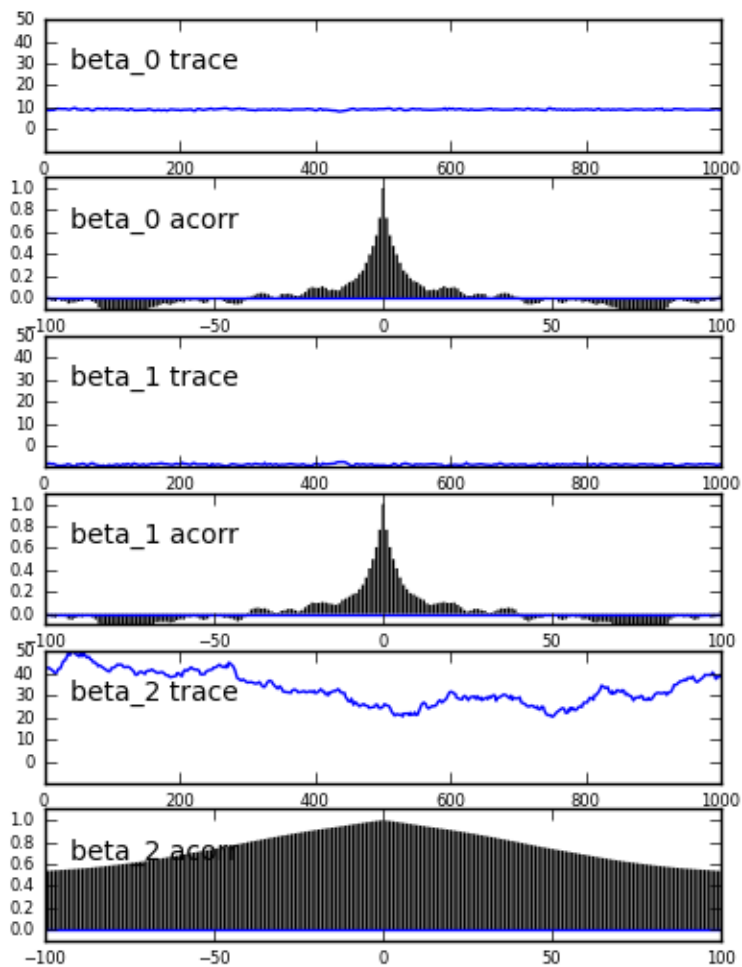
```
map = mc.MAP(vars)  
map.fit(method='fmin_powell')
```

```
mcmc = mc.MCMC(vars)  
mcmc.sample(iter=10000, burn=5000, thin=5)
```

```
# see results of mcmc with this  
mc.Matplot.plot(mcmc)
```







How to cope when MCMC don't converge

```
# run chain longer
```

```
m.sample(50000, 25000, 50)
```

```
# change initial values
```

```
mc.MAP(vars).fit(method='fmin_powell')
```

```
# change step methods
```

```
m = mc.MCMC(vars)
```

```
m.use_step_method(mc.AdaptiveMetropolis,  
                  m.beta)
```

Step Method Movies

Convergence diagnostics in PyMC

Graphics:

`mc.Matplotlib.plot(m)`

Other:

start from a bunch of different initial values and confirm that they all yield the same results.

Stats:

`mc.diagnostics`

`.discrepancy`

`.gelman_rubin`

`.geweke`

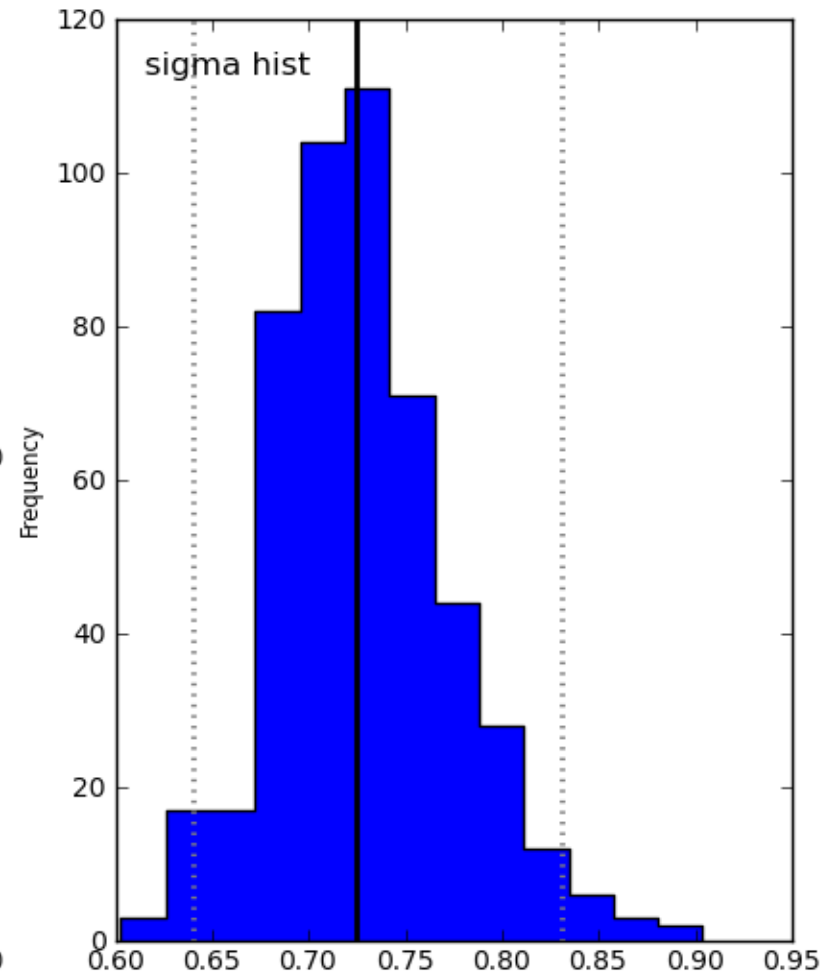
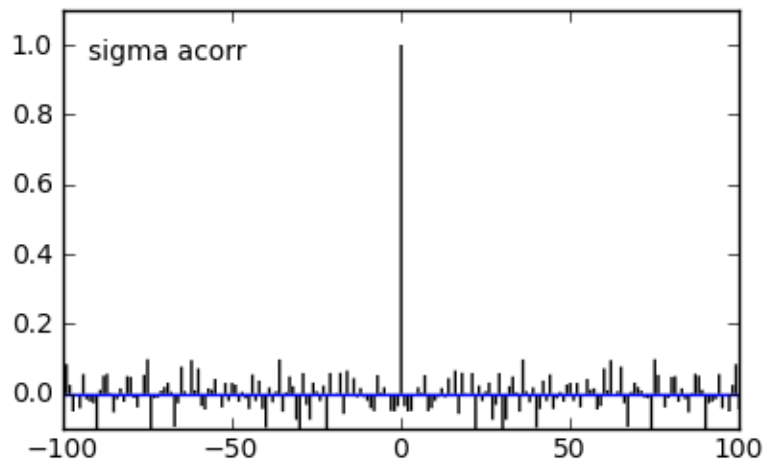
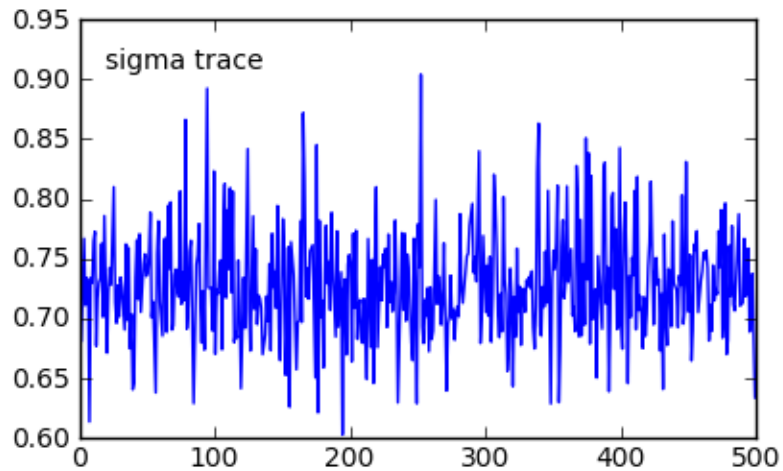
`.iat`

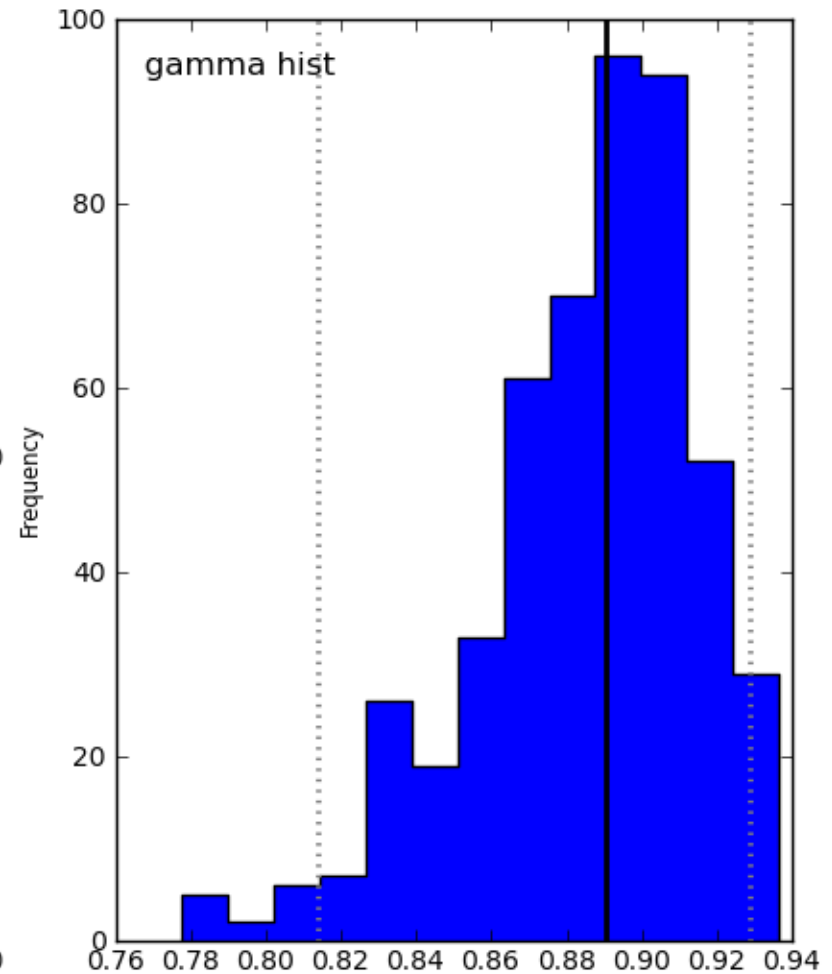
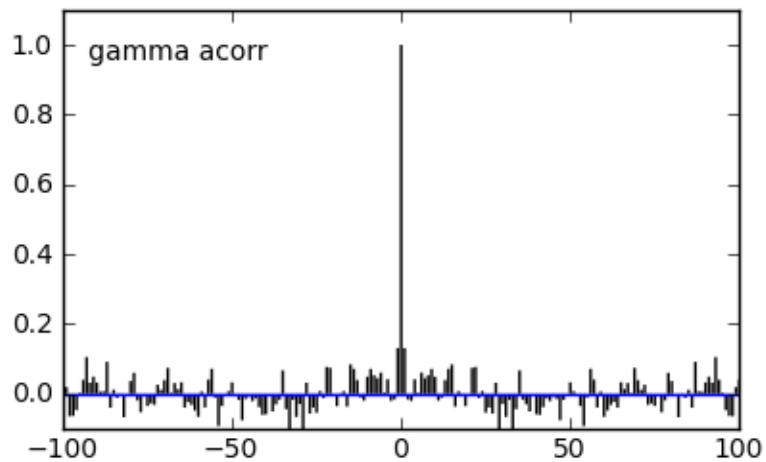
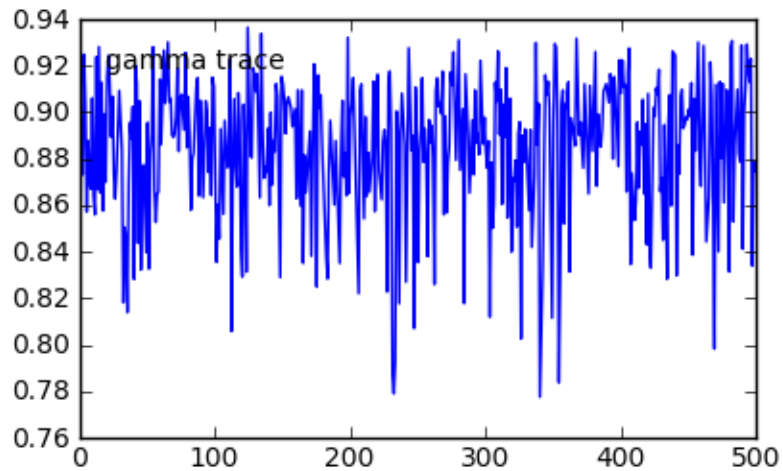
`.raftery_lewis`

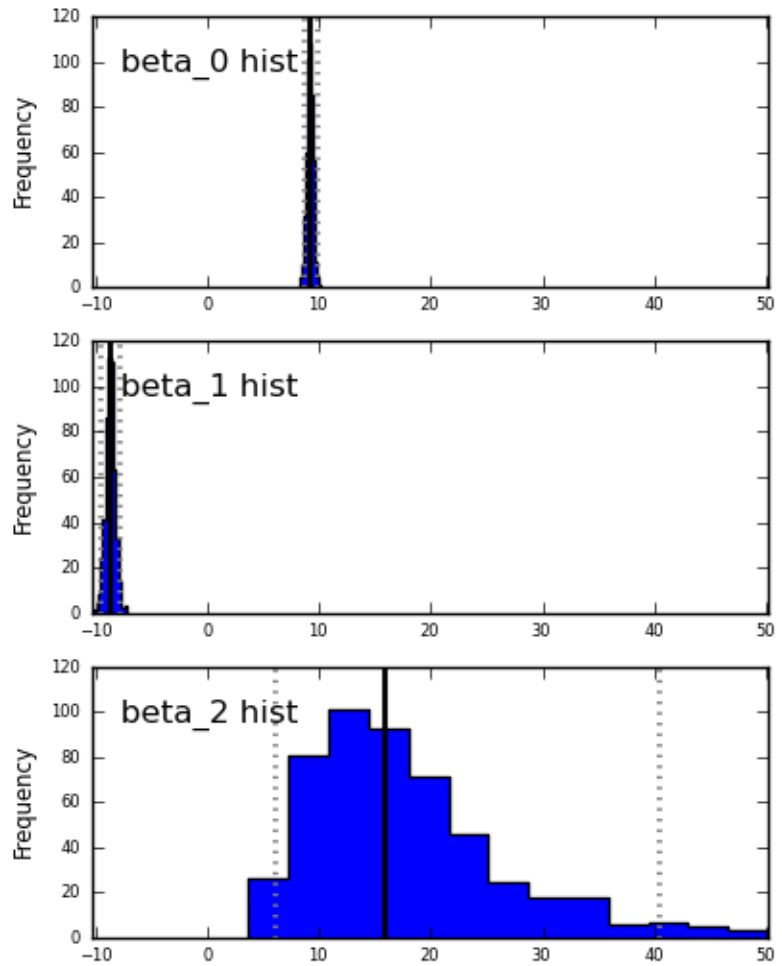
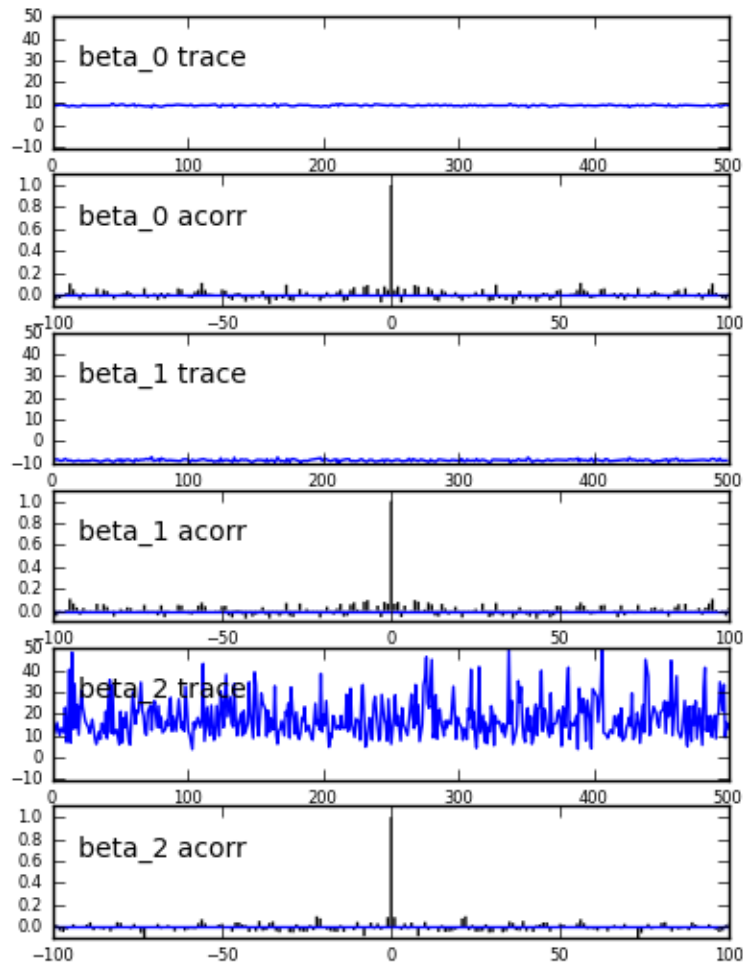
In our case, Adaptive Metropolis gets the job done

```
vars = models.nonlinear()  
  
# fit with MCMC, but not with  
# default step methods  
m = mc.MCMC(vars)  
m.use_step_method(mc.AdaptiveMetropolis,  
                  [m.beta, m.gamma])  
m.sample(50000, 25000, 50)  
  
# explore model convergence  
mc.Matplot.plot(m)
```

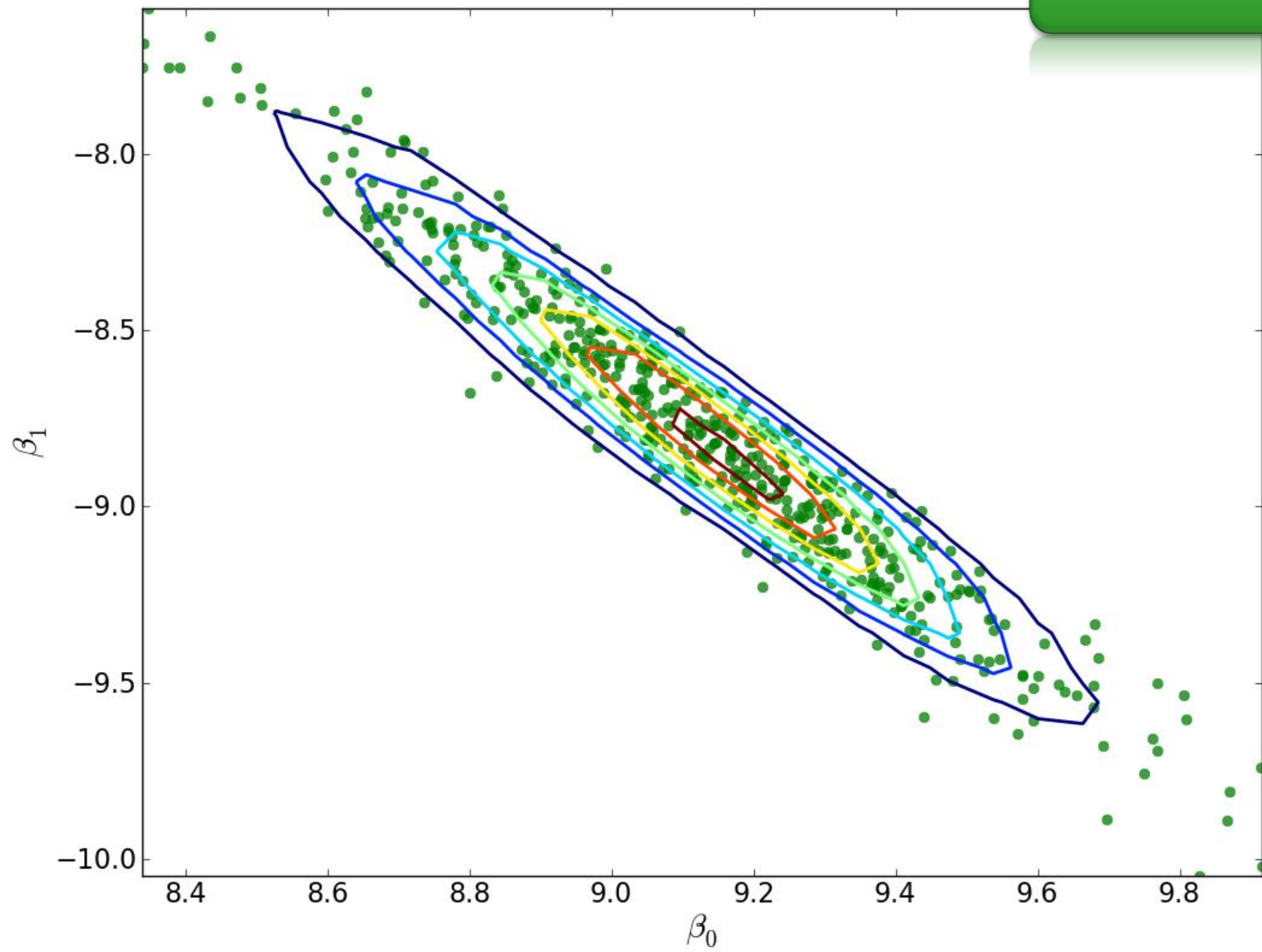
src/ex4.py



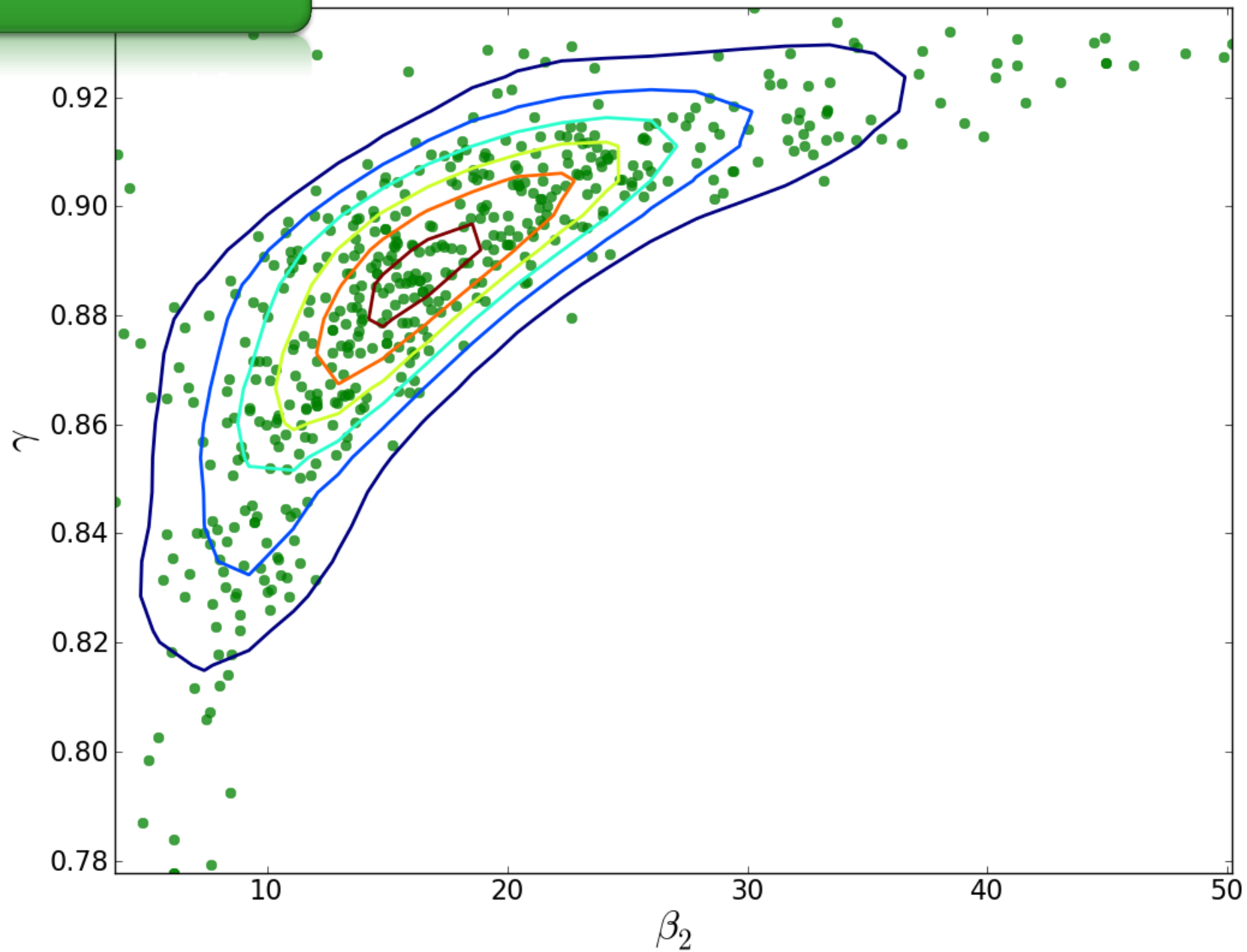




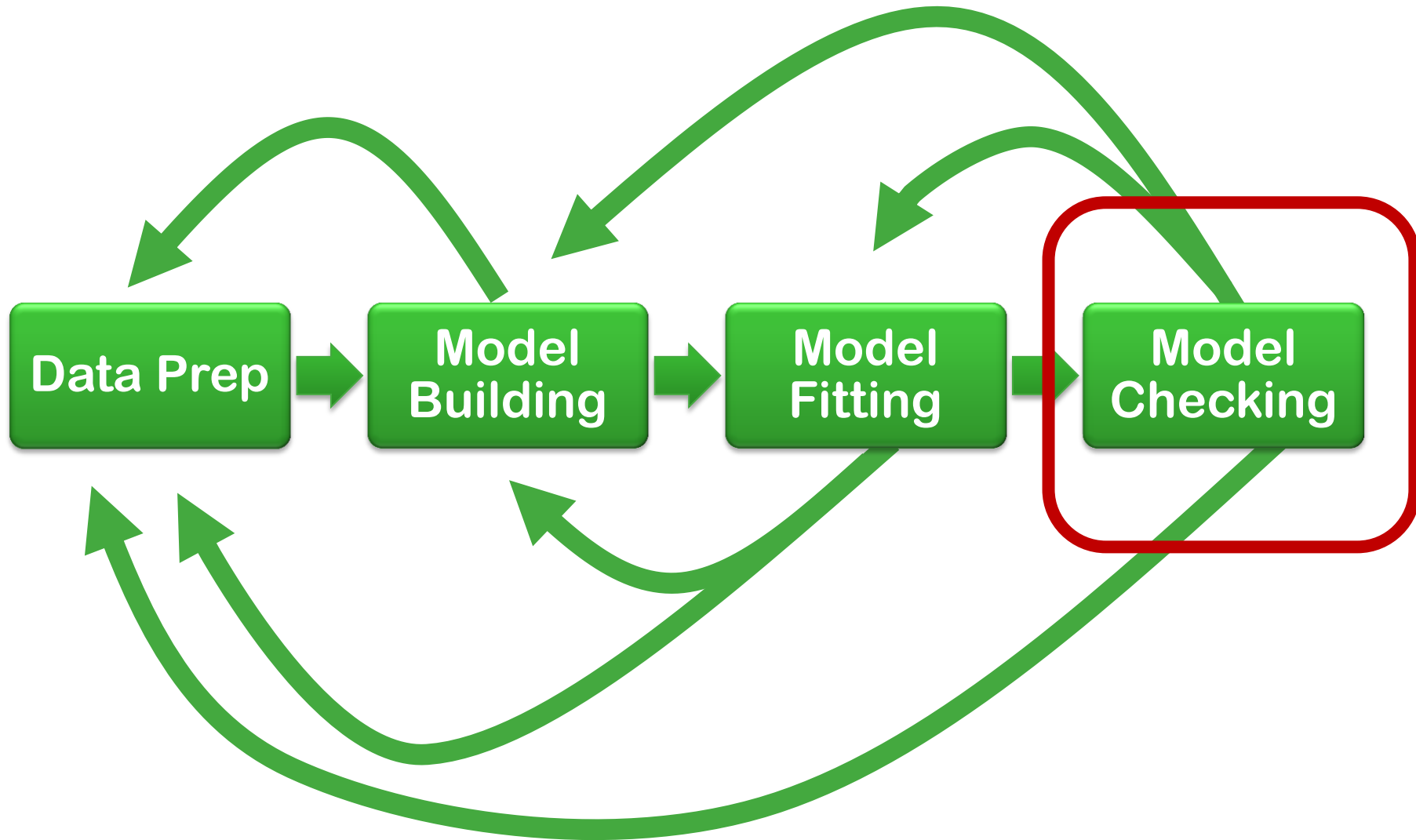
**In hindsight,
why was that model tough to fit?**



src/ex4.py



A Generic Template for Bayesian data analysis



Information Criterion

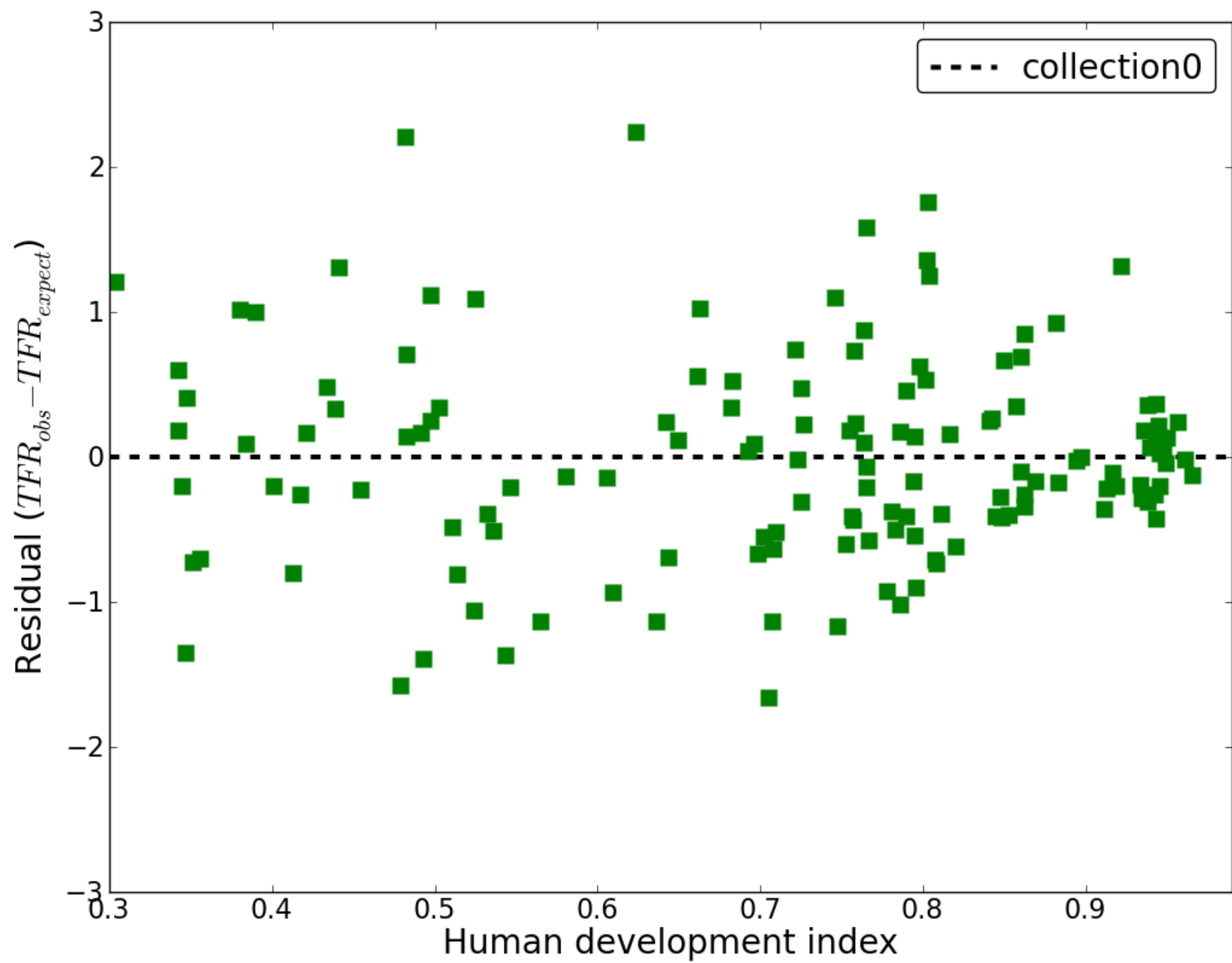
```
map = mc.MAP(vars)  
map.fit(method='fmin_powell')  
map.AIC, map.BIC
```

```
mcmc = mc.MCMC(vars)  
mcmc.sample(iter=10000, burn=5000, thin=5)  
mcmc.dic
```

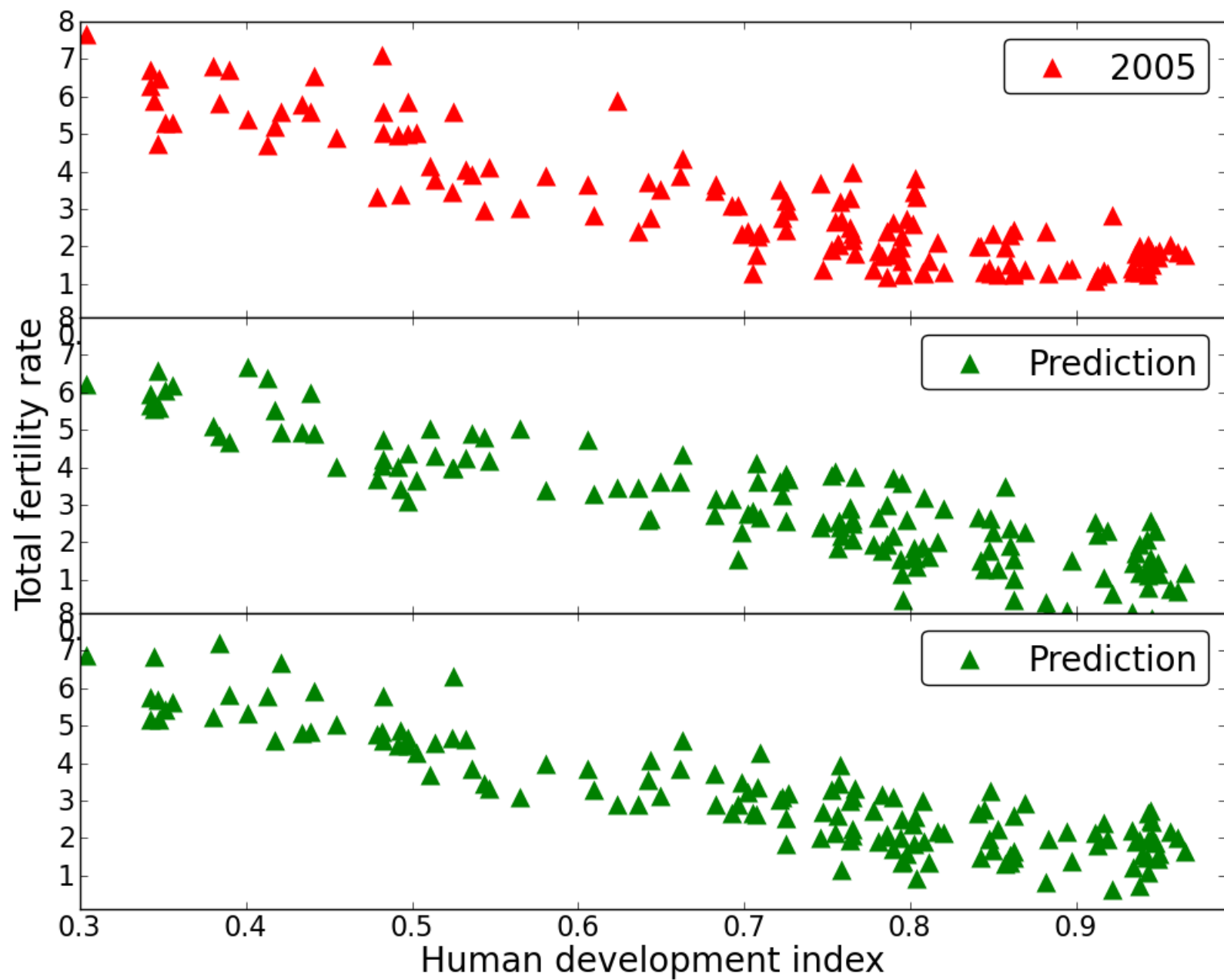
```
# or plot residuals
```

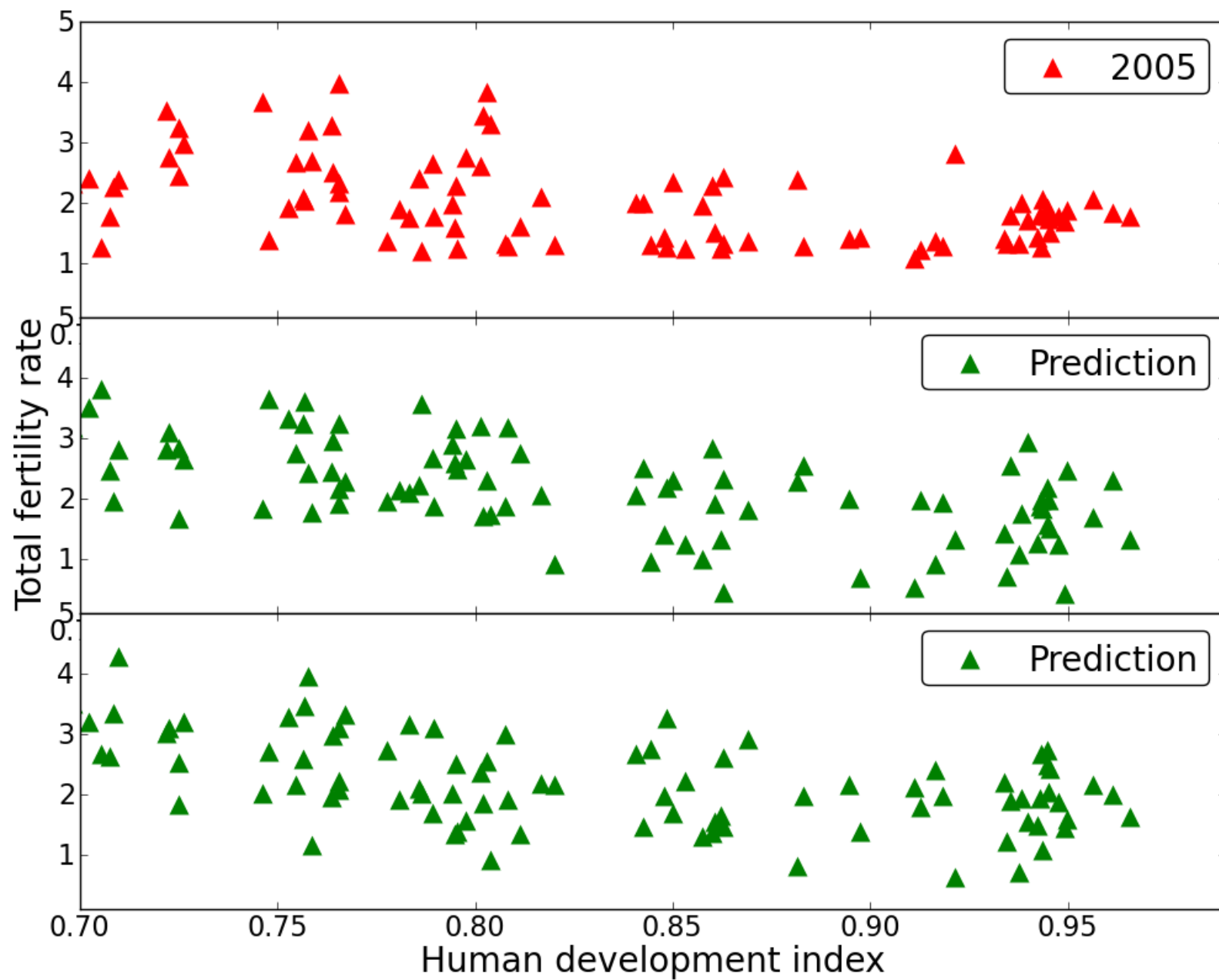
```
# or posterior predictive checks
```

```
y_err = mcmc.y_obs.value - mcmc.y_mean.stats()['mean']  
pl.hlines([0], 0, 1, linewidth=3, linestyle='dashed')  
pl.plot(data.hdi2005, y_err, 'gs', mew=0, ms=10)
```




```
# add a data posterior prediction deterministic
@mc.deterministic
def y_pred(mu=vars['y_mean'], sigma=vars['sigma']):
    return mc.rnormal(mu, sigma**-2)
vars['y_pred'] = y_pred
```





Saving Output

```
import pymc as mc
import models
vars = models.nonlinear()

# fit with MCMC, save results
# databases: ram, pickle, txt, sqlite, mysql, hdf5
mcmc = mc.MCMC(vars, db='txt',
               dbname='nonlinear', dbmode='w')
mcmc.sample(5000, 2500, 5)

# load the database from disk
db = mc.database.txt.load('nonlinear')
print db.beta.stats()
```

Brainstorm – what do you want to do with this?

- More complicated models
- Alternative step methods
- Alternative graphics
- Go through a similar exercise with a different application
- Do something with MCMC that has nothing to do with statistics, e.g. making mazes or jigsaw puzzles, or randomly coloring graphs
- Something else?