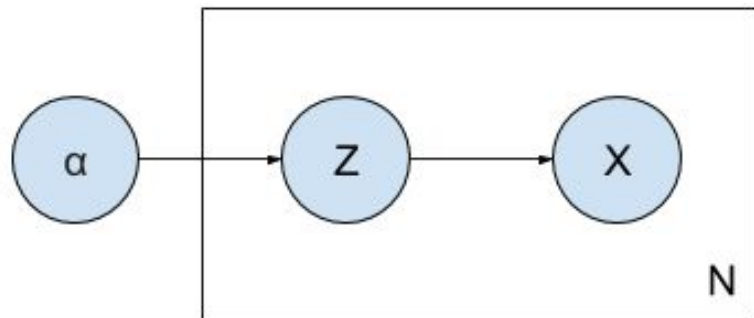
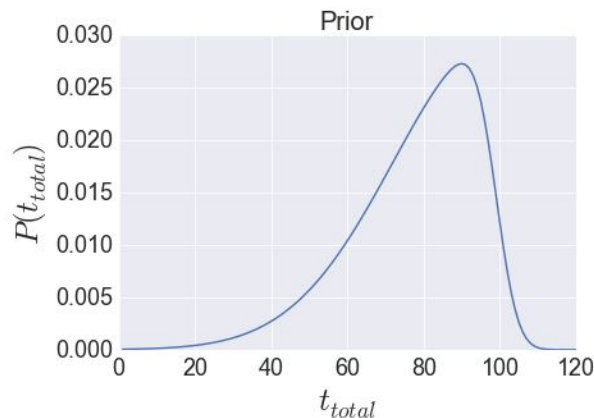


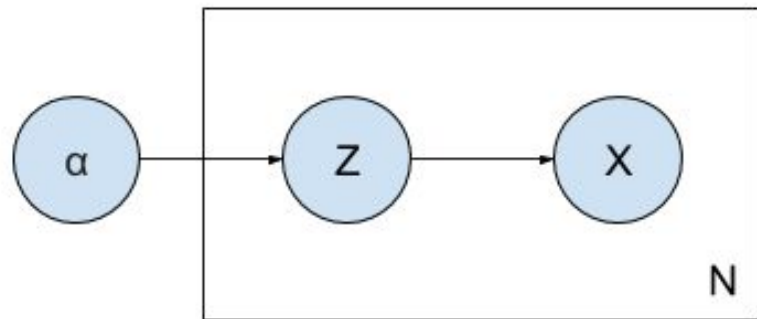
# Learning Priors

So far..

- $p(t_{\text{total}}|t) \propto p(t|t_{\text{total}})p(t_{\text{total}})$
- $p(t|t_{\text{total}}) = 1/t_{\text{total}}$  for  $t_{\text{total}} \geq t$  and 0 otherwise
- $Z \Rightarrow t_{\text{total}}$
- $X \Rightarrow t$
- $Z = \{z_1, z_2, \dots, z_N\}$
- $X = \{x_1, x_2, \dots, x_N\}$



- If  $z_1 = 80$ , then  $x_1 \leq 80$
- $X$  is superset of total lifespans
- You might encounter or hear about an alive or dead person
- Implication: even people who are alive inform our prior of total lifespan



# Expectation Maximization

$$\begin{aligned} Q(\alpha|\alpha^{(t)}) &= E_{Z|X, \alpha^{(t)}} [\log L(\alpha; X, Z)] \\ &= \sum_{i=1}^N E_{z_i|x_i, \alpha^{(t)}} [\log L(\alpha; x_i, z_i)] \\ &= \sum_{i=1}^N \sum_{z_i} p(z_i|x_i, \alpha^{(t)}) \log p(x_i, z_i|\alpha) \\ &= \sum_{i=1}^N \sum_{z_i} T(x_i, z_i) \log (p(x_i|z_i)p(z_i|\alpha)) \\ &= \sum_{i=1}^N \sum_{z_i} T(x_i, z_i) \log (p(z_i|\alpha)/z_i) \end{aligned}$$

where:

- $T(x_i, z_i) := p(z_i|x_i, \alpha^{(t)})$  is a fixed function with respect to  $\alpha^{(t)}$
- $p(x_i, z_i|\alpha) = p(x_i|z_i)p(z_i|\alpha)$  is the same formula from Tenenbaum's paper, with  $\alpha^{(t)}$  included, since  $p(x_i|z_i, \alpha) = p(x_i|z_i)$  by assumption
- Substituting  $p(x_i|z_i) = 1/z_i$  for all  $z_i$  is okay since  $T(x_i, z_i) = 0$  for  $z_i \neq x_i$  by definition.

$$T(x_i, z_i) = p(z_i|x_i, \alpha^{(t)}) = \frac{p(z_i, x_i|\alpha^{(t)})}{\sum_{z_i} p(z_i, x_i|\alpha^{(t)})}$$

- Normally - maximize expectation wrt  $\alpha$
- Substitute  $\alpha$  and repeat until convergence
- Analogous to batch gradient descent

$$\sum_{i=1}^N \sum_{z_i} T(x_i, z_i) \log (p(z_i|\alpha)/z_i)$$

# Stepwise Expectation Maximization

- Update  $\alpha$  for every incoming sample  $x_i$
- Analogous to stochastic gradient descent
- Since single sample is a bad approximation we interpolate between  $s'_i$  and  $\mu$

$$\sum_{z_i} T(x_i, z_i) \log (p(z_i|\alpha)/z_i)$$

$\mu \leftarrow$  initialization;  $k = 0$

for each iteration  $t = 1, \dots, T$ :

for each example  $i = 1, \dots, n$  in random order:

$$s'_i \leftarrow \sum_{\mathbf{z}} T(x_i, z_i) \log (p(z_i|\alpha)/z_i) \quad [\text{inference}]$$

$$\mu \leftarrow (1-\eta_k)\mu + \eta_k s'_i; k \leftarrow k+1 \quad [\text{towards new}]$$

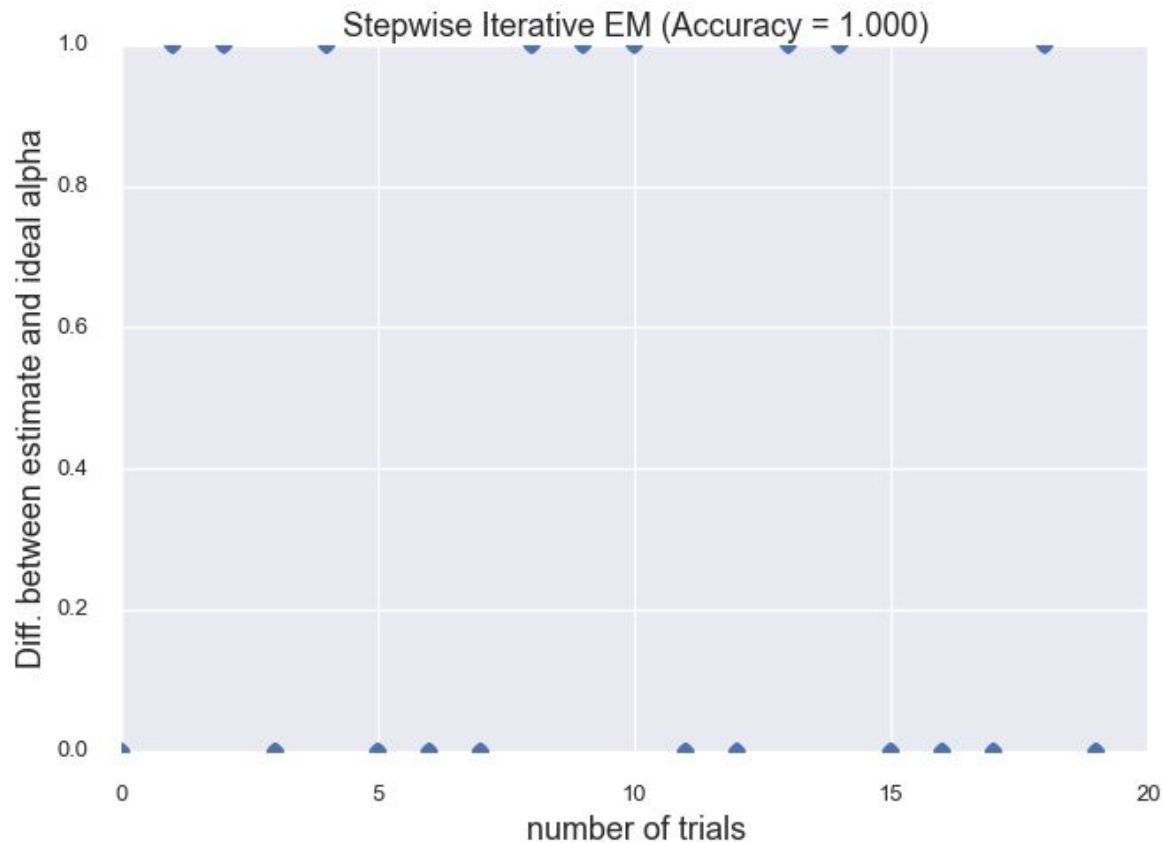
# Properties of stepwise EM

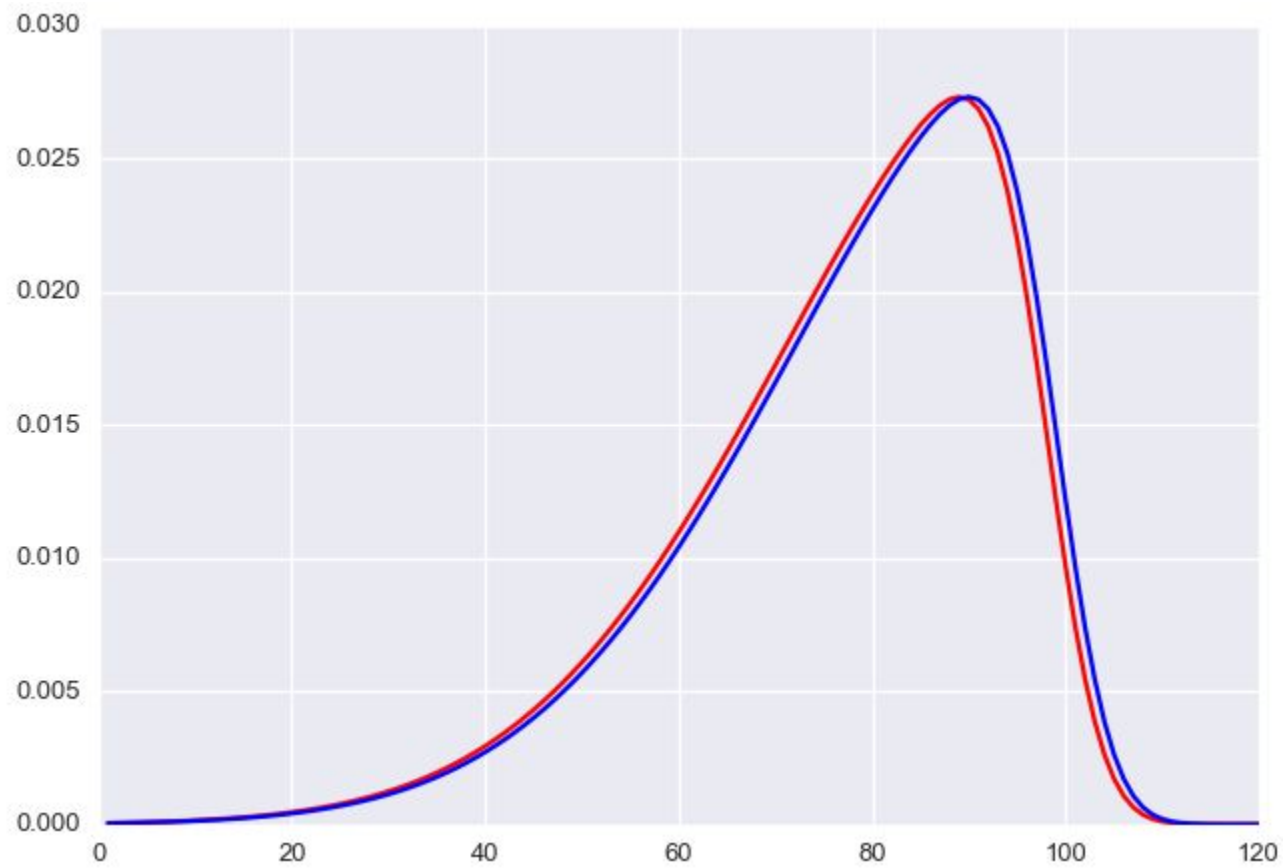
- Standard results from stochastic approximation literature
- If we take  $\eta_k = (k + 2)^{-\alpha}$ , then any  $0.5 < \alpha \leq 1$  is valid
- The smaller the alpha, the larger the updates and more quickly we forget (decay) the old statistics.
- Swift progress but generates instability in the trajectory
- One of the principal motivations: Speed
- Liang, P., & Klein, D. (2009, May)

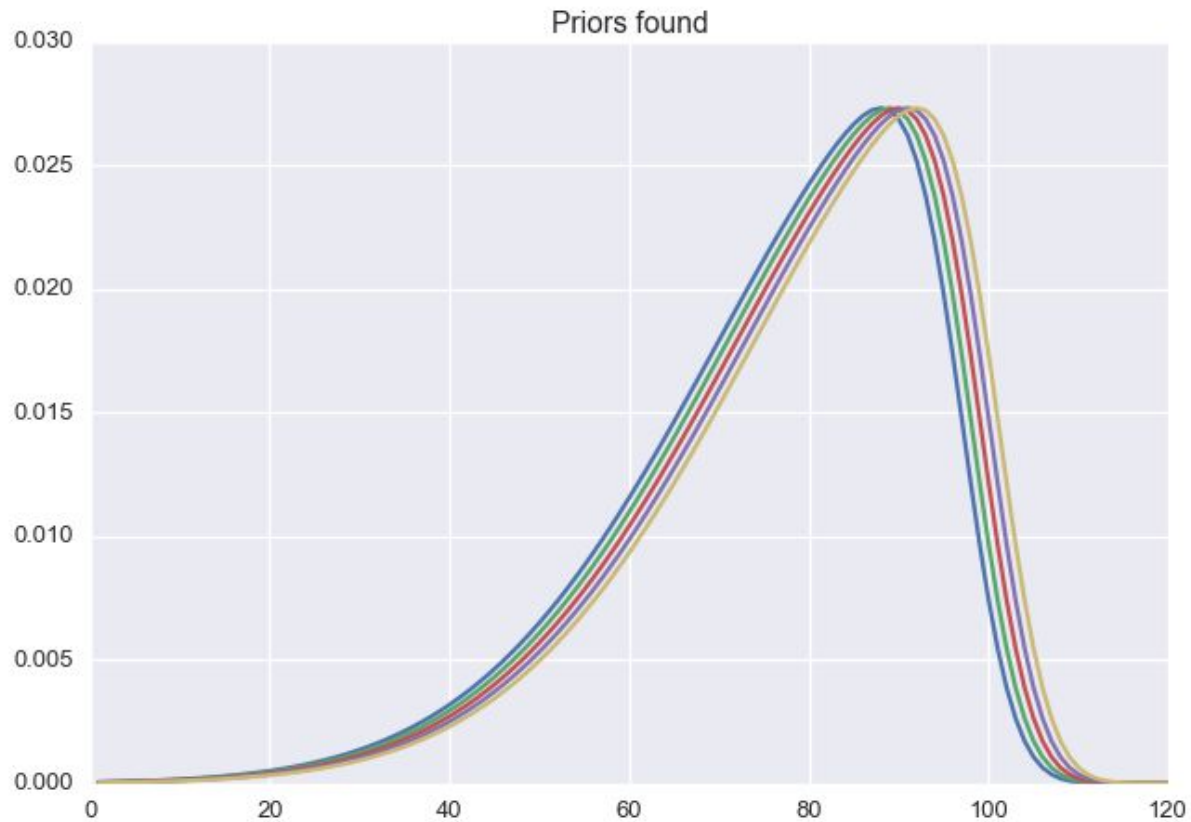
# Preliminary Results



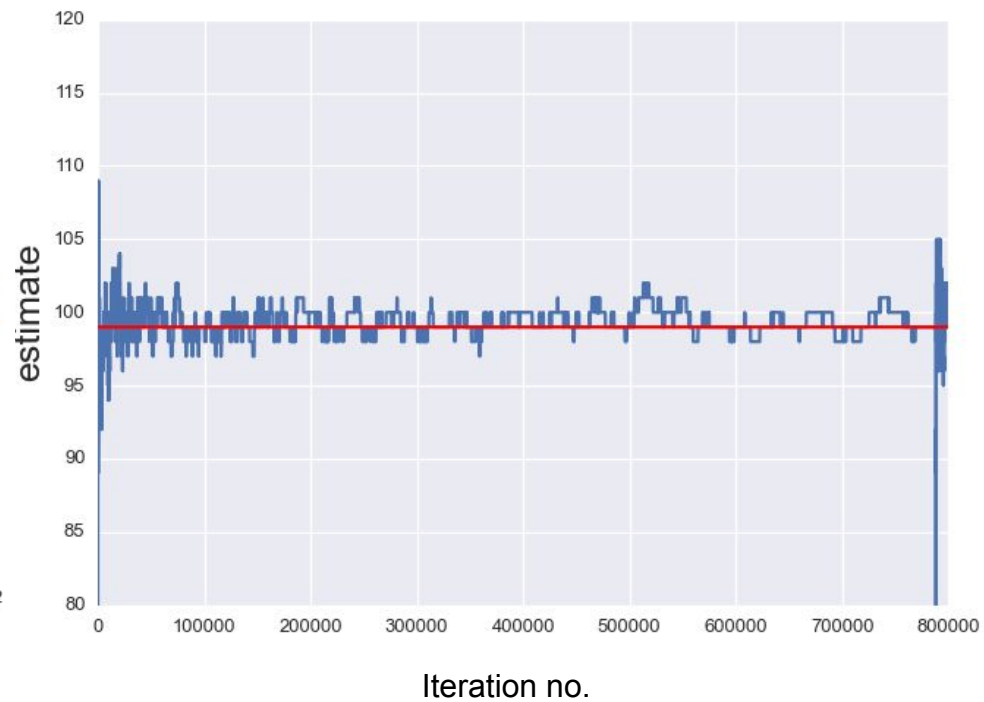
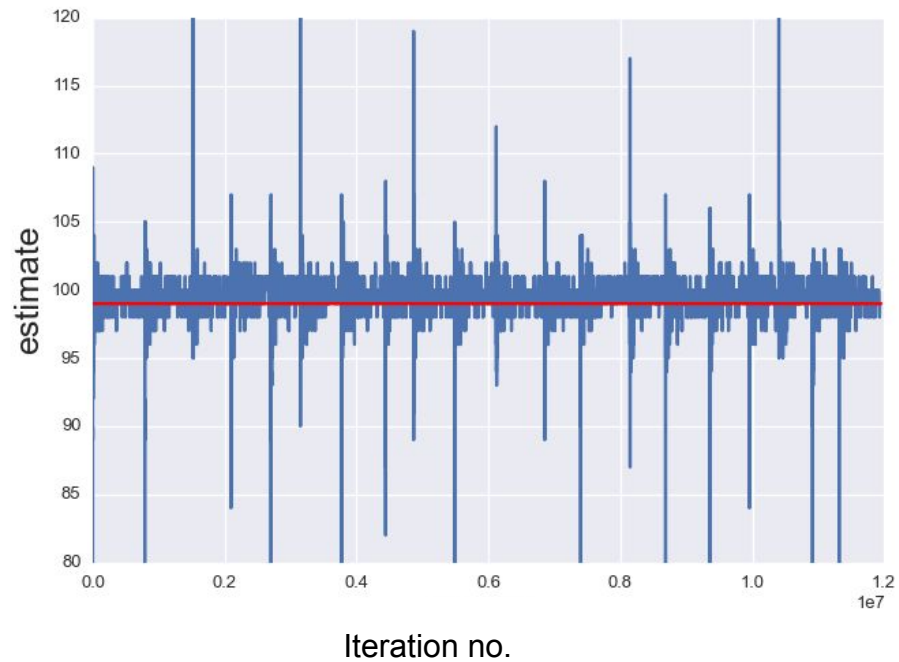
- $\alpha = 0.5$
- Iter per sample = 5
- Batch iter = 5
- Trials = 20
- Samples/trial = 50,000

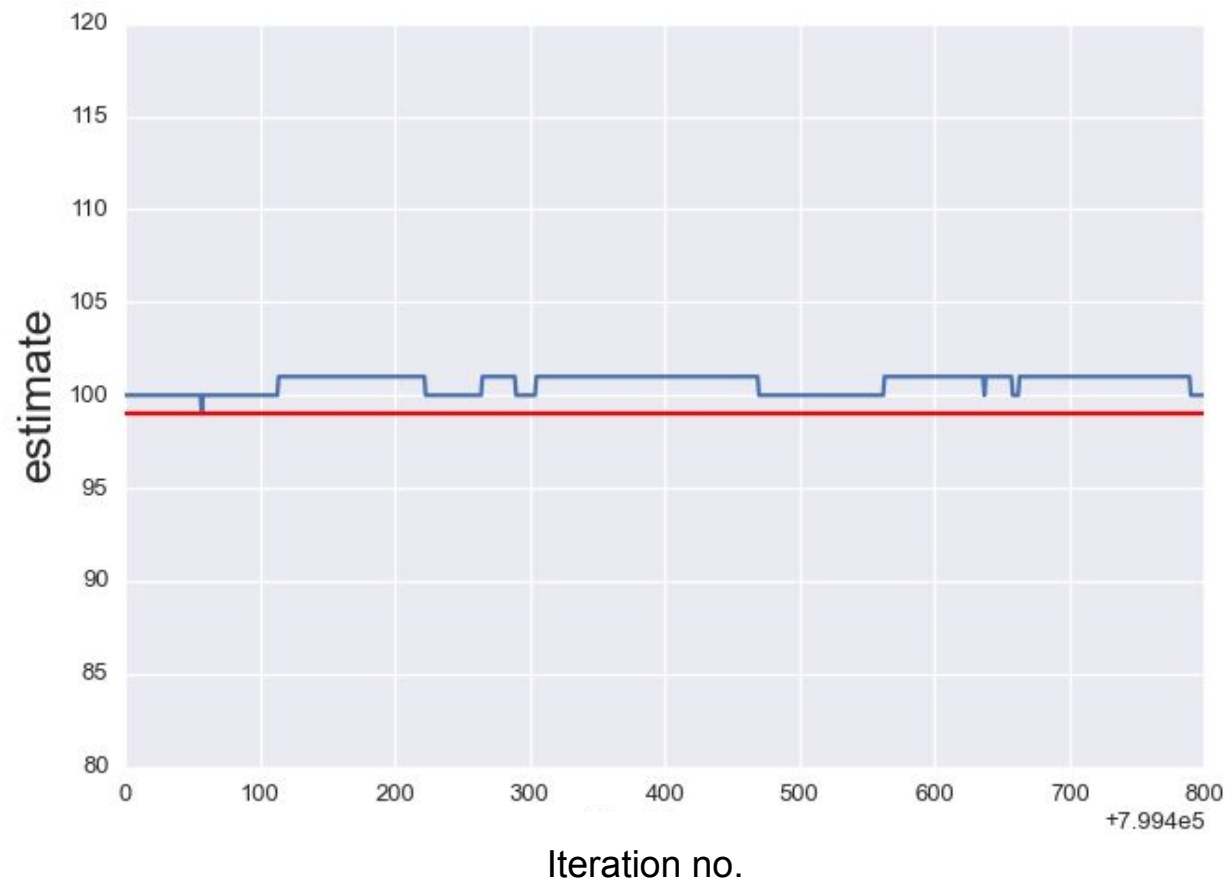




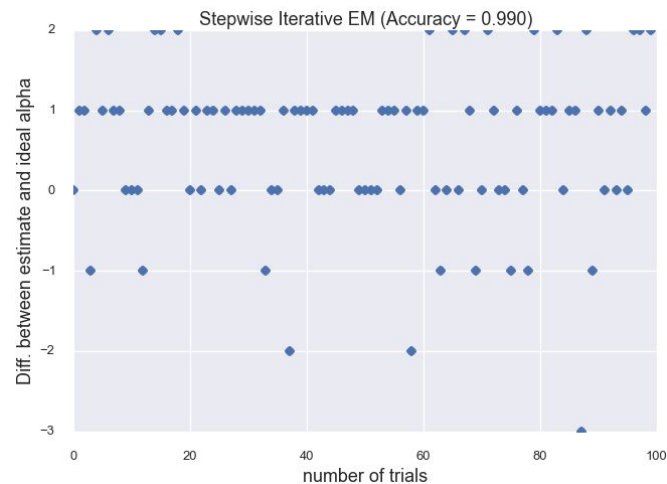
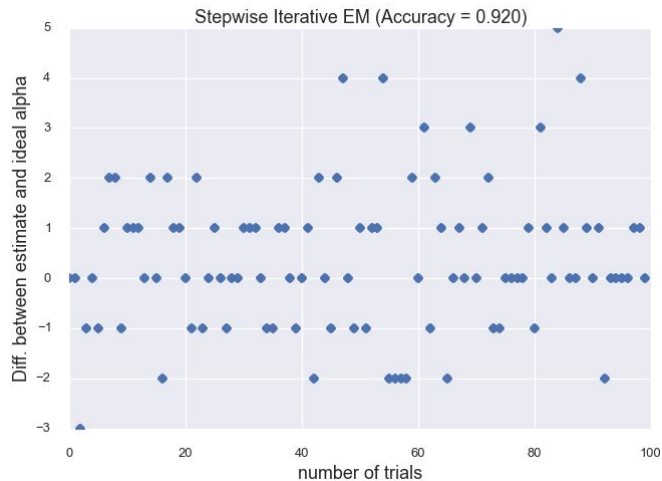
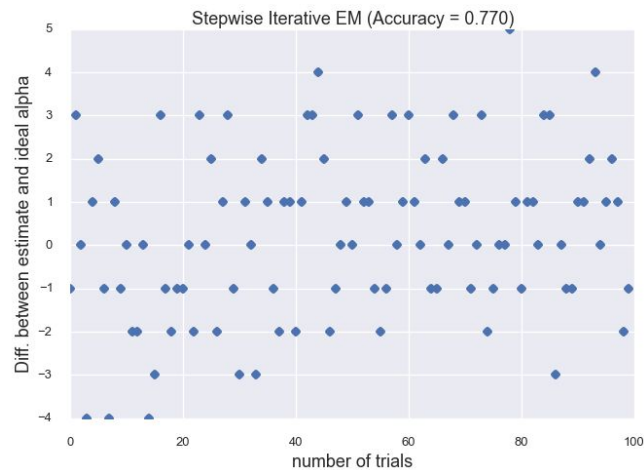


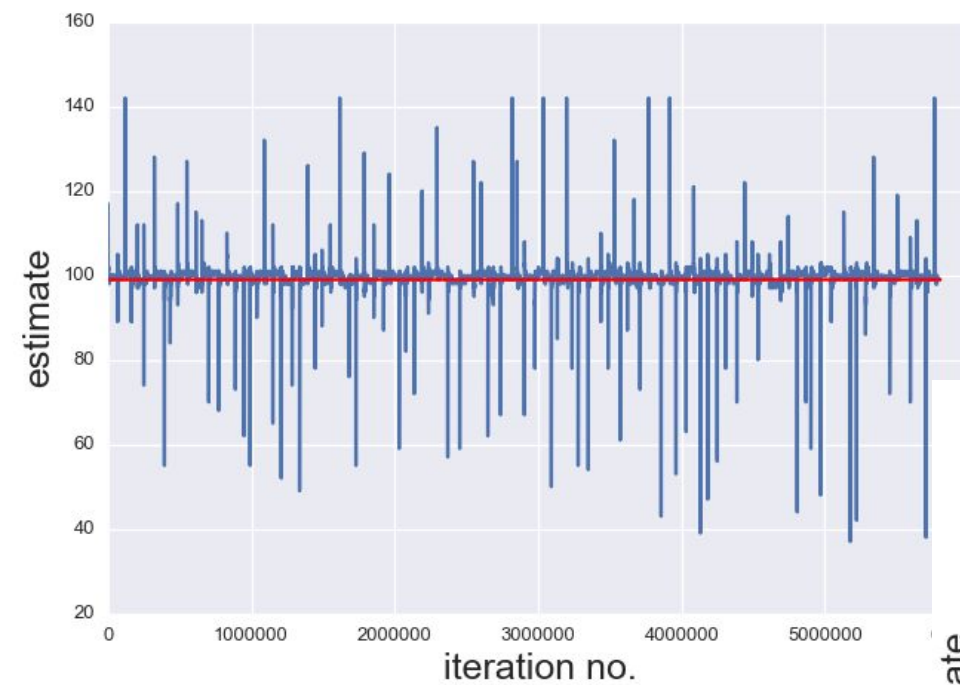
Accuracy measure = Should be within  $\pm 2$  of the optimal estimate



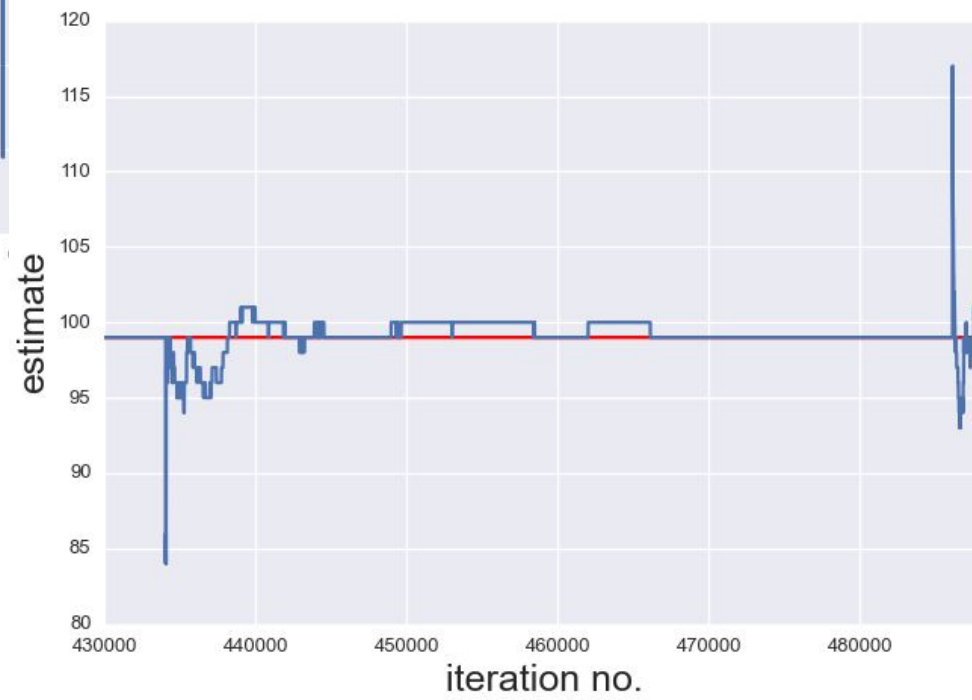


Samples/trials :  
500, 1000,  
3000, and 5000

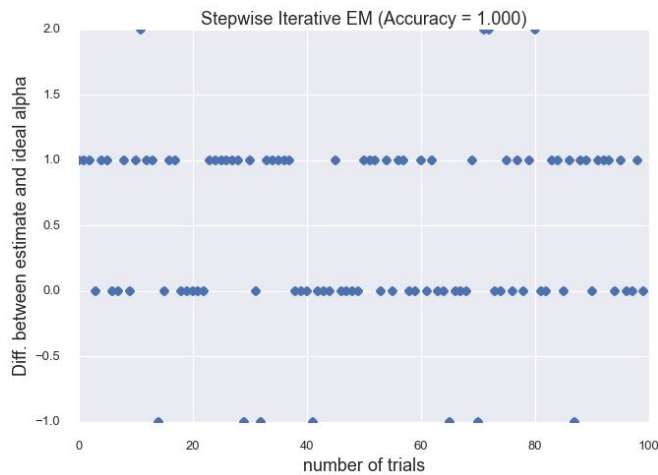
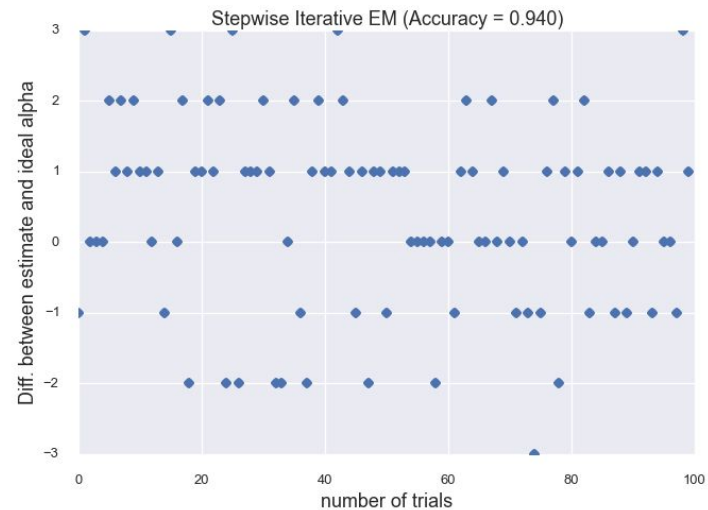
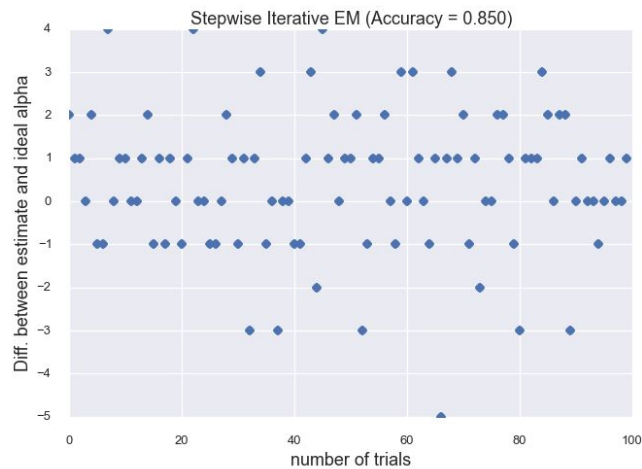




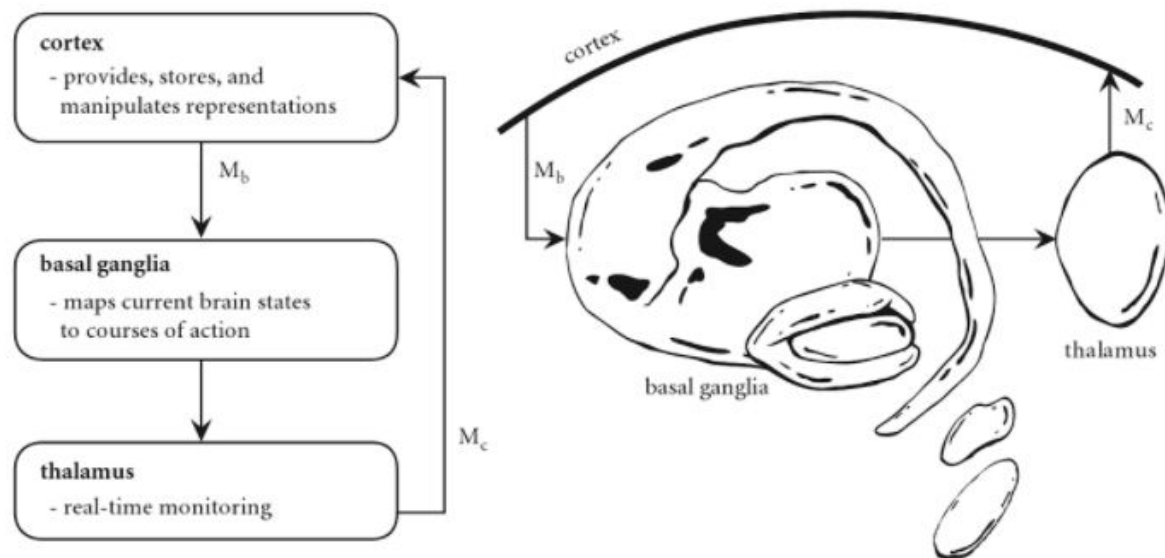
$\alpha = 0.65$   
Increase in stability



Samples/trials :  
500, 1000,  
3000, and 5000







**FIGURE 5.5** The cortex-basal ganglia-thalamus loop. Arrows indicate connections between the three areas. At a functional level, brain states from the cortex are mapped through the  $M_b$  matrix to the basal ganglia. Each row in such a matrix specifies a known context for which the basal ganglia will choose an appropriate action. The product of the current cortical state and  $M_b$  provides a measure of how similar the current state is to each of the known contexts. The output of the basal ganglia disinhibits the appropriate areas of thalamus. Thalamus, in turn, is mapped through the matrix  $M_c$  back to the cortex. Each column of this matrix specifies an appropriate cortical state that is the consequence of the selected action. The relevant anatomical structures are pictured on the right based on a simplified version of Figure 5.1.

# Nengo Implementation

- Step1: Sample  $x_i$  comes in (cortical input)
- Step2: Compute the likelihood function  $p(x/z)$  - 120 dim vector (cortical state)
 
$$p(t|t_{\text{total}}) = 1/t_{\text{total}} \text{ for } t_{\text{total}} \geq t \text{ and } 0 \text{ otherwise}$$
- Step3:  $p(x, z) = p(x/z) p(z/\alpha_0)$  where  $\alpha_0$  is our initial guess - 120 dim vector (cortical state)
- Step4:  $s'_i = (p(x, z) / \text{np.sum}(p(x, z))) \cdot \text{dot}(\log p(Z/A).T)$  (matrix  $M_b$ )
  - (120) dim . (120\*M) dim  $\Rightarrow$  M dim where M is the size of search space for  $\alpha$
- Step5:
  - $\eta_k = (k + 2)^{-\alpha}$ , Increment  $k$
  - $\mu \leftarrow (1 - \eta_k)\mu + \eta_k s'_i; k \leftarrow k + 1$  (Bg input)
- Step6:  $\bar{\alpha} = \text{argmax}(\mu)$  (automatically happens in Bg)
- Step7: Find  $p(z/\bar{\alpha})$  (Bg output)

(update cortical state in step 3, omit step 1 and 2 until a new sample comes in.)