

# **A comparison of Bayesian estimators for unsupervised Hidden Markov Model POS taggers**

Conference on Empirical Methods in NLP, 2008

**Jianfeng Gao**

Microsoft Research

**Mark Johnson**

Brown University

Presenter: Manish Gupta

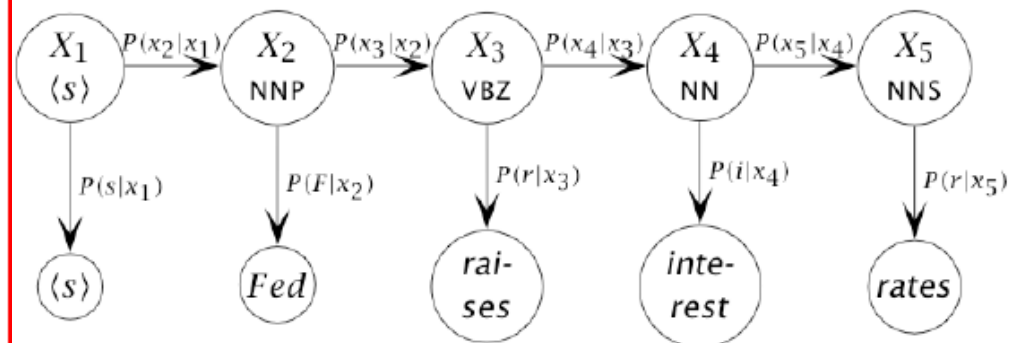
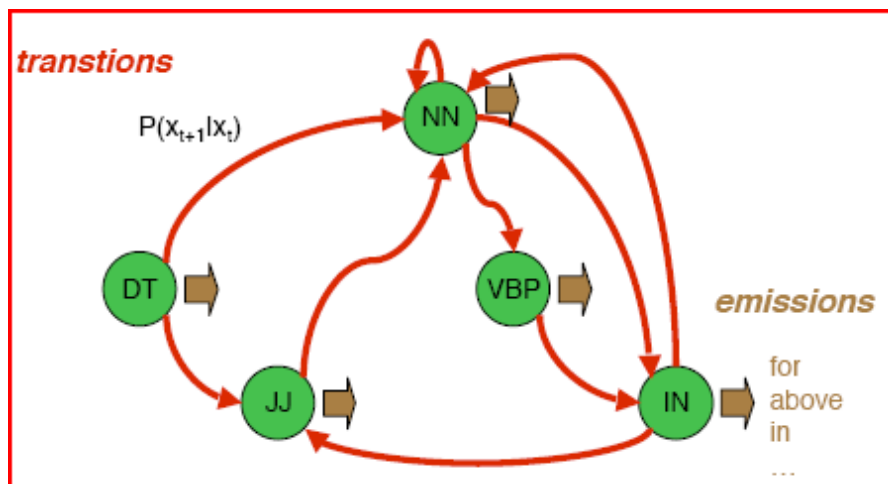
Instructor: Dr. Julia Hockenmaier

CS598

24<sup>th</sup> Feb 2010

# Basics

- Bayesian estimator: Estimator that minimizes posterior expected value of a loss function.
- Consider an unknown parameter  $\theta$  with prior distribution  $\pi$ . Let  $\delta(x)$  be an estimator where  $x$ =data. Then Bayes risk= $E_{\pi}(L(\delta, \theta))$ .  $\delta$  is Bayesian estimator that minimizes Bayes risk.
- Unsupervised: no labels/tags
- Hidden Markov Model (HMM)



# HMM and POS

- Problem: Identify label sequence given word sequence
- Observed: word sequence ( $\mathbf{w}$ ).  $|\mathbf{w}|=n$
- Hidden: POS sequence ( $\mathbf{t}$ ). #states= $m$
- Parameters:
  - Transition probabilities ( $\theta_t$ ) – Multinomial
  - Emission probabilities ( $\phi_t$ ) – Multinomial
  - Initial state distribution ( $\pi$ )
  - $\lambda = (\theta, \phi, \pi)$

# Inference for HMMs

- Parameters:
  - Transition probabilities ( $\theta_t$ ) – Multinomial
  - Emission probabilities ( $\phi_t$ ) – Multinomial

$$\begin{array}{l|l} \theta_t & \alpha \sim \text{Dir}(\alpha) \\ \phi_t & \alpha' \sim \text{Dir}(\alpha') \end{array}$$

- For experiments, they use uniform  $\alpha$  and uniform  $\alpha'$ .
- $\alpha$  controls sparsity of transition probabilities and  $\alpha'$  controls sparsity of emission probabilities.
- $\alpha' \rightarrow 0$ 
  - prior prefers models where each state emits as few words as possible
  - Situation: most words belong to a single POS

# Bayesian estimation

- As against MLE/MAP, Bayesian estimation uses multiple values of parameters.
- Posterior does not have a closed form.

$$\begin{array}{|c|} \hline P(\theta|D) \\ \hline \text{Posterior} \\ \hline \end{array} = \frac{\begin{array}{|c|c|} \hline \text{Prior} & \text{Likelihood} \\ \hline P(\theta) & P(D|\theta) \\ \hline \end{array}}{\begin{array}{|c|} \hline \int P(\theta)P(D|\theta)d\theta \\ \hline \text{Marginal Likelihood (=}P(D)\text{)} \\ \hline \end{array}}$$

- Inference methods: EM, Variational Bayes (VB) estimation (approx), 4 types of Gibbs sampler (converge to true posterior)

# Baum Welch (Forward-Backward/EM) Algorithm

- Compute forward and backward probabilities.

$$\alpha_k(t) = P(w_1, \dots, w_k \mid t_k = t, \lambda) \quad \beta_k(t) = P(w_{k+1}, w_{k+2}, \dots, w_n \mid t_k = t, \lambda)$$

- $\alpha_k(t)$  is the probability of observing a partial sequence of observables  $w_1, \dots, w_k$  given state  $t_k = t$  at time  $k$ , and  $\lambda$
- $\beta_k(t)$  is the probability of observing a partial sequence of observables  $w_{k+1}, \dots, w_n$  given state  $t_k = t$  at time  $k$ , and  $\lambda$
- Use dynamic programming to compute  $\alpha$  and  $\beta$

# E Step

- Compute counts using forward and backward probabilities
- Let  $n_{t',t}$  be the probability of being in state  $t$  at time  $k$  and at state  $t'$  at time  $k+1$ , given  $\lambda$  and  $\mathbf{w}$  sequence

$$n_{t',t}(k) = \frac{\alpha_k(t)\theta_{tt'}\phi_{t'}(w_{k+1})\beta_{k+1}(t')}{P(\mathbf{w} | \lambda)} = \frac{\alpha_k(t)\theta_{tt'}\phi_{t'}(w_{k+1})\beta_{k+1}(t')}{\sum_{t=1}^m \sum_{t'=1}^m \alpha_k(t)\theta_{tt'}\phi_{t'}(w_{k+1})\beta_{k+1}(t')}$$

$$E[n_{t',t}] = \sum_{k=1}^{n-1} n_{t',t}(k)$$

- Let  $n_t(k)$  be the probability of being in state  $t$  at time  $k$ , given  $\mathbf{w}$

$$n_t(k) = \sum_{t'=1}^m n_{t',t}(k) \quad E[n_t] = \sum_{k=1}^{n-1} n_t(k) \quad E[n'_{w,t}] = \sum_{k=1, w_k=w}^{n-1} n_t(k)$$

# M step

- Use these counts to compute updated parameters.
- Iteratively re-estimates parameters.
- Converges to local maximum

$$\theta_{t'|t}^{(\ell+1)} = \frac{E[n_{t',t}]}{E[n_t]}$$

$$\phi_{w|t}^{(\ell+1)} = \frac{E[n'_{w,t}]}{E[n_t]}$$

- $n'_{w,t}$  is #times word  $w$  occurs with state  $t$
- $n_{t',t}$  is #times state  $t'$  follows  $t$
- $n_t$  is #occurrences of state  $t$
- $O(nm^2)$  time



# Variational Bayes

- Aim: Find  $(\theta, \phi, t)$  that minimizes  $-\log P(\mathbf{w})$

$$-\log P(\mathbf{w})$$

$$= -\log \int \int \int Q(t, \theta, \phi) \cdot \frac{P(\mathbf{w}, t, \theta, \phi)}{Q(t, \theta, \phi)} dt d\theta d\phi$$

Jensen's  
inequality

$$\leq -\int \int \int Q(t, \theta, \phi) \cdot \log \frac{P(\mathbf{w}, t, \theta, \phi)}{Q(t, \theta, \phi)} dt d\theta d\phi$$

$$-\log P(\mathbf{w})$$

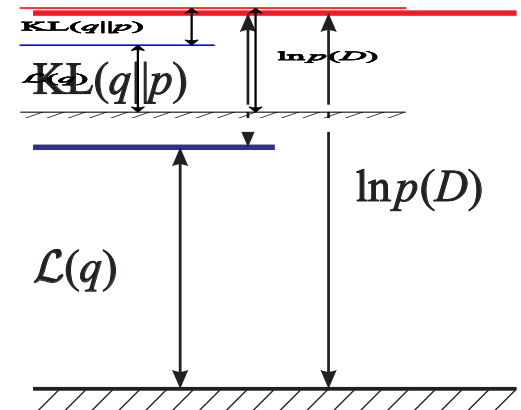
$$= -\int \int \int Q(t, \theta, \phi) \cdot \log \frac{P(\mathbf{w}, t, \theta, \phi)}{Q(t, \theta, \phi)} dt d\theta d\phi - \int \int \int Q(t, \theta, \phi) \cdot \log \frac{Q(t, \theta, \phi)}{P(t, \theta, \phi | \mathbf{w})} dt d\theta d\phi$$

$$= -\int \int \int Q(t, \theta, \phi) \cdot \log \frac{P(\mathbf{w}, t, \theta, \phi)}{Q(t, \theta, \phi)} dt d\theta d\phi - KL(q | p)$$

Variational  
free energy

# Variational Bayes

- Find a  $Q(t, \theta, \phi)$  that minimizes an upper bound to the negative log likelihood.
- Mean field assumption: local densities can be used to denote effects of global densities.
- Factorized model:  $Q(t, \theta, \phi) = Q_1(t) \times Q_2(\theta, \phi)$
- Minimize the KL divergence between desired posterior distribution and factorized approximation.
- $O(nm^2)$



# Variational Bayes

- If likelihood and prior belong to exponential family, VB is similar to Forward Backward Algorithm.

Smoothed counts

- E step is the same

- M step:  
$$\tilde{\theta}_{t'|t}^{(\ell+1)} = f(\mathbf{E}[n_{t',t}] + \alpha) / f(\mathbf{E}[n_t] + m\alpha)$$
$$\tilde{\phi}_{w|t}^{(\ell+1)} = f(\mathbf{E}[n'_{w,t}] + \alpha') / f(\mathbf{E}[n_t] + m'\alpha')$$
$$f(v) = \exp(\Psi(v))$$

Digamma is first derivative of log gamma

- $m$  and  $m'$  are #word types and states.

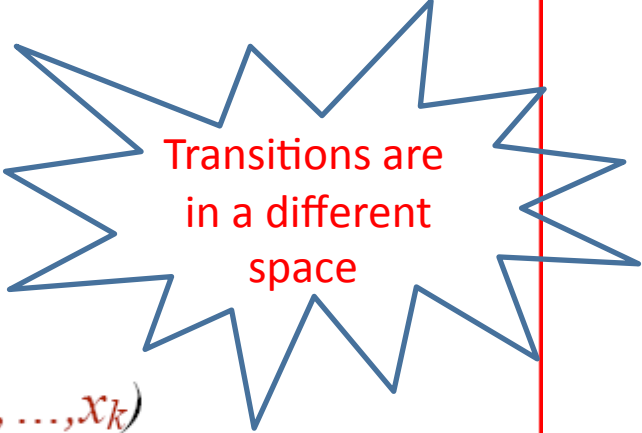
# Gibbs sampling

We will visit states according to transition probabilities  $P(\mathbf{y}|\mathbf{x})$

That is, we will go from state  $\mathbf{x} = (x_1, \dots, x_k)$   
to state  $\mathbf{y} = (y_1, \dots, y_k)$

For  $i = 1 \dots k$ :

pick  $y_i$  by sampling from  $P(Y_i | y_1, \dots, y_{i-1}, x_{i+1}, \dots, x_k)$



Transitions are  
in a different  
space

- We need all exact conditional distributions to estimate the joint probability distribution

# MCMC sampling algorithms

- Produce a stream of samples from posterior distribution  $P(\mathbf{t} | \mathbf{w}, \boldsymbol{\alpha})$
- 4 different Gibbs samplers:
  - Pointwise or blocked
  - Explicit or Collapsed
- Pointwise: Resamples a single state  $t_i$  (labeling a single word  $w_i$ ) at each step.  $O(nm)$  per iteration.
- Blocked: Resamples labels for all of the words in a sentence at a single step.  $O(nm^2)$  per iteration.
- Explicit: Samples  $\boldsymbol{\theta}$  and  $\boldsymbol{\phi}$  along with states  $\mathbf{t}$
- Collapsed:  $\boldsymbol{\theta}$  and  $\boldsymbol{\phi}$  are integrated out. Only  $\mathbf{t}$  are sampled.

# Pointwise explicit Gibbs sampler

- Resample  $\theta$  and  $\phi$  given state-to-state transition counts  $\mathbf{n}$  and state-to-word emission counts  $\mathbf{n}'$
- Resample each state  $t_i$  given word  $w_i$  and neighboring states  $t_{i-1}$  and  $t_{i+1}$

$$\theta_t \mid \mathbf{n}_t, \alpha \sim \text{Dir}(\mathbf{n}_t + \alpha)$$

$$\phi_t \mid \mathbf{n}'_t, \alpha' \sim \text{Dir}(\mathbf{n}'_t + \alpha')$$

$$P(t_i \mid w_i, \mathbf{t}_{-i}, \theta, \phi) \propto \theta_{t_i|t_{i-1}} \phi_{w_i|t_i} \theta_{t_{i+1}|t_i}$$

# Collapsed blocked Gibbs sampler

$$\theta_{t'|t}^* = \frac{n_{t',t} + \alpha}{n_t + m\alpha}$$
$$\phi_{w|t}^* = \frac{n'_{w,t} + \alpha'}{n_t + m'\alpha}$$

- Resample states for each sentence given  $\mathbf{n}$  and  $\mathbf{n}'$  for other sentences in the corpus.
- Following Metropolis-Hastings accept reject step, decide whether current state sequence be updated to  $t^*$  or whether to keep current state sequence.
- High acceptance rates: 99%

# Evaluation metrics

- Variation of information (VI): (lower the better)
  - $VI = H(C) + H(C') - 2I(C, C')$  where  $I(C, C') = H(C) - H(C|C')$
  - The variation of information (VI) between two clusterings  $C$  (the gold standard) and  $C'$  (the found clustering) of a set of data points is a sum of the amount of information lost in moving from  $C$  to  $C'$ , and the amount that must be gained.
  - Problem: Tagger that assigns all words the same POS has good VI
- Cross validation accuracy (higher the better)
  - Map each HMM state to the part-of-speech tag it co-occurs with most frequently (using train set), and use this mapping to map each HMM state sequence  $t$  to a sequence of part-of-speech tags (using validation set).
- Greedy 1-to-1 accuracy (higher the better)
  - At most 1 HMM state can be mapped to any POS tag.



# Experiments

- 8 different combinations of hyper-parameters  $\alpha$  and  $\alpha'$  (0.0001 to 1)
- Data sets of different sizes (24K – 120K – 1174K words)
- Tag sets of different sizes (Noah Smith's 17 tag set, Penn Treebank tag set)
- Run each setting 10 times with at least 1000 iterations.

	All – 50	All – 17	120K – 50	120K – 17	24K – 50	24K – 17
EM	0.40527	0.43101	0.29303	0.35202	0.18618	0.28165
VB	0.46123	<b>0.51379</b>	0.34679	0.36010	0.23823	0.36599
GS <sub>e,p</sub>	0.47826	0.43424	0.36984	0.44125	0.29953	0.36811
GS <sub>e,b</sub>	0.49371	0.46568	0.38888	<b>0.44341</b>	0.34404	0.37032
GS <sub>c,p</sub>	<b>0.49910*</b>	0.45028	<b>0.42785</b>	0.43652	<b>0.39182</b>	<b>0.39164</b>
GS <sub>c,b</sub>	0.49486*	0.46193	0.41162	0.42278	0.38497	0.36793

Figure 2: Average greedy 1-to-1 accuracy of state sequences produced by HMMs estimated by the various estimators. The column heading indicates the size of the corpus and the number of HMM states. In the Gibbs sampler (GS) results the subscript “e” indicates that the parameters  $\theta$  and  $\phi$  were explicitly sampled while the subscript “c” indicates that they were integrated out, and the subscript “p” indicates pointwise sampling, while “b” indicates sentence-blocked sampling. Entries tagged with a star indicate that the estimator had not converged after weeks of run-time, but was still slowly improving.

	All – 50	All – 17	120K – 50	120K – 17	24K – 50	24K – 17
EM	0.62115	0.64651	0.44135	0.56215	0.28576	0.46669
VB	0.60484	0.63652	0.48427	0.36458	0.35946	0.36926
GS <sub>e,p</sub>	0.64190	0.63057	0.53571	0.46986	0.41620	0.37165
GS <sub>e,b</sub>	<b>0.65953</b>	0.65606	0.57918	0.48975	0.47228	0.37311
GS <sub>c,p</sub>	0.61391*	<b>0.67414</b>	<b>0.65285</b>	<b>0.65012</b>	<b>0.58153</b>	<b>0.62254</b>
GS <sub>c,b</sub>	0.60551*	0.65516	0.62167	0.58271	0.55006	0.58728

Figure 3: Average cross-validation accuracy of state sequences produced by HMMs estimated by the various estimators. The table headings follow those used in Figure 2.

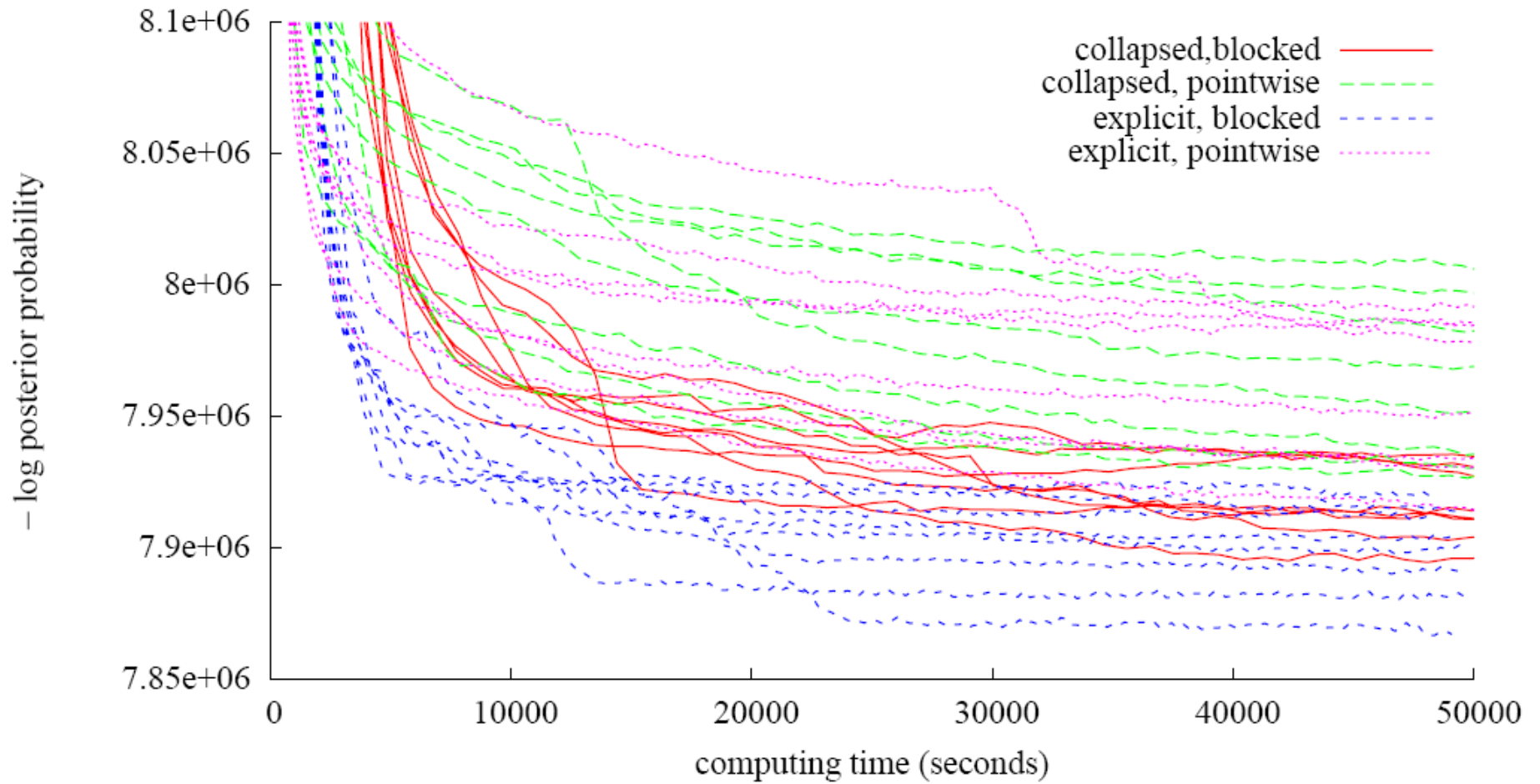
	All – 50	All – 17	120K – 50	120K – 17	24K – 50	24K – 17
EM	4.47555	3.86326	6.16499	4.55681	7.72465	5.42815
VB	4.27911	3.44029	5.00509	3.19670	4.80778	3.14557
GS <sub>e,p</sub>	4.24919	3.53024	4.30457	3.23082	<b>4.24368</b>	3.17076
GS <sub>e,b</sub>	4.04123	<b>3.46179</b>	4.22590	3.20276	4.29474	<b>3.10609</b>
GS <sub>c,p</sub>	<b>4.03886*</b>	3.52185	<b>4.21259</b>	<b>3.17586</b>	4.30928	3.18273
GS <sub>c,b</sub>	4.11272*	3.61516	4.36595	3.23630	4.32096	3.17780

Figure 4: Average Variation of Information between the state sequences produced by HMMs estimated by the various estimators and the gold tags (smaller is better). The table headings follow those used in Figure 2.

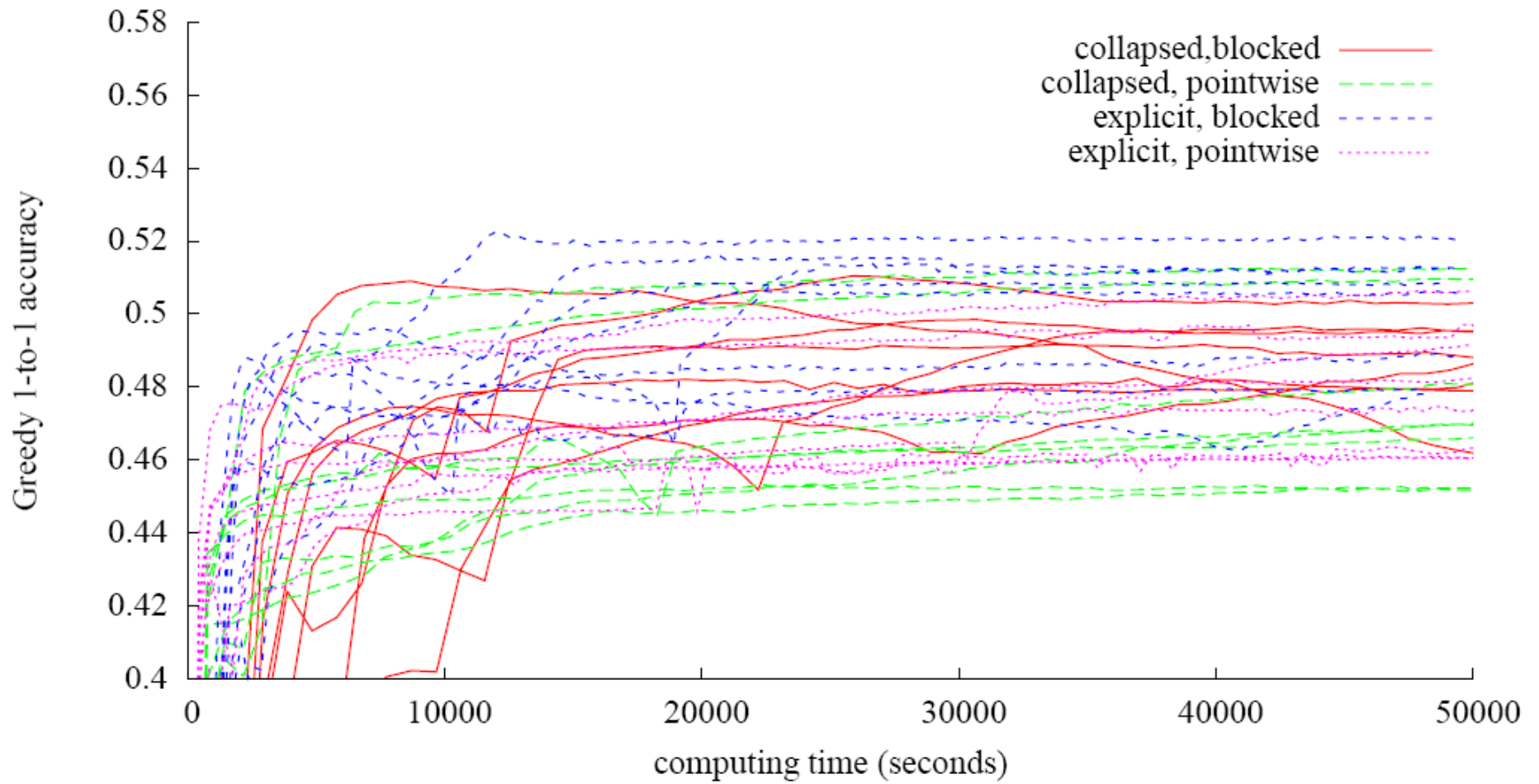
	All – 50	All – 17	120K – 50	120K – 17	24K – 50	24K – 17
EM	558	346	648	351	<b>142</b>	125
VB	<b>473</b>	<b>123</b>	<b>337</b>	<b>24</b>	183	<b>20</b>
GS <sub>e,p</sub>	2863	382	3709	63	2500	177
GS <sub>e,b</sub>	3846	286	5169	154	4856	139
GS <sub>c,p</sub>	*	34325	44864	40088	45285	43208
GS <sub>c,b</sub>	*	6948	7502	7782	7342	7985

Figure 5: Average number of iterations until the negative logarithm of the posterior probability (or likelihood) changes by less than 0.5% (smaller is better) per at least 2,000 iterations. No annealing was used.

All data, 50 states,  $\alpha = \alpha' = 0.1$



All data, 50 states,  $\alpha = \alpha' = 0.1$



# Findings

- Point-wise samplers need  $O(m)$  steps per sample. EM, VB and sentence-blocked Gibbs sampler need  $O(m^2)$  steps.
- On small datasets, all Bayesian estimators outperform EM (and to a lesser extent, VB).
  - Reasoning: Priors are imp when data is less. Also, approximation by VB would be inaccurate on small data.
- On large datasets, EM does well to cross validation accuracy.
- VB converges faster. Larger  $\alpha$  and  $\alpha'$  cause faster convergence.
  - Reasoning:  $\alpha$  and  $\alpha'$  specify how likely the samplers are to consider novel tags and so influence sampler's mobility
- Blocked samplers converge faster than pointwise samplers. Explicit samplers are faster than collapsed ones.
- Pointwise samplers initially converge faster than blocked ones → Hybrid strategy could be better.

Thanks!