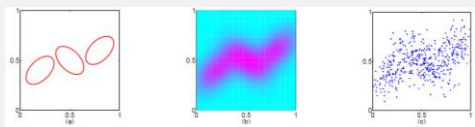


BTRY 4840/6840, CS 4775 Computational Genetics and Genomics



October 18, 2018

Final projects

- 1-page proposal due Oct 30 (a week from Tuesday)
- Find something you're excited about that is related to class topics
 - Finding overlap with other classes or your research is good
- Projects need to involve modeling, implementation, and data analysis
 - Only rare exceptions to the above; must get approval
- Graduate students expected to produce novel research, undergrads can re-implement a method
- Be ambitious but keep focused – have ~5 weeks
- In class presentations in weeks following proposal

Final projects, continued

- Final projects due December 11 by 4:30pm
 - Turn in: report, code, data (if lightweight)
- Report should be ~6-10 pages, but length is less important than the quality of the report
 - Use standard scientific paper format: Introduction, Methods, Results, Discussion
 - Figures are important: take time to make these clear and communicative, reference and explain these in the text
 - Put important excerpts from source code and/or data in Appendix
 - Keep this brief: raw data is less important than interpretation
- Remember: this is equivalent of 2½ problem sets

Today's lecture

- Learning: parameter estimation
 - Learning when some variables are latent (hidden) via Expectation-Maximization (EM)
 - Related approach: K-means clustering
 - Example: EM for Gaussian mixture model
 - Baum-Welch algorithm: EM for HMMs
 - Viterbi training

Expectation-Maximization

Recall: 2d Gaussian mixture model

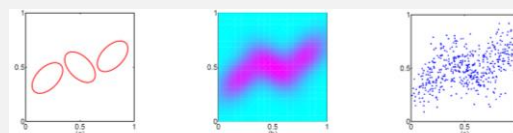


Figure 10.2: Illustration of a mixture of 3 Gaussians in a two-dimensional space showing (a) contours representing one standard deviation for each of the mixture components, (b) the marginal probability density of the mixture distribution, and (c) a sample of 500 points drawn from the marginal distribution.

- Parameters: means, covariances of 3 Gaussians
- Latent variables: which Gaussian each point from

Jordan (2003)

Recall: K-means clustering algorithm

7

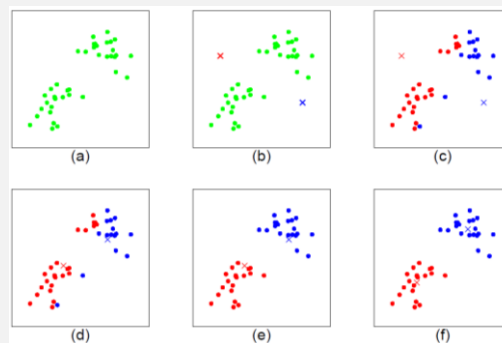
- Wish to estimate
 - Cluster means: (μ_1, \dots, μ_K)
 - Assignments of points to clusters: (z_1, \dots, z_N)
- Algorithm:
 - Make initial guess of (μ_1, \dots, μ_K)
 - Can use K random data points or some other way to initialize
 - Update assignments: [closest cluster via Euclidean distance]

$$z_n = \arg \min_j \|x_n - \mu_j\|^2$$
 - Update the cluster means:

$$\mu_i = \frac{\sum_n x_n \cdot I[z_n = i]}{\sum_n I[z_n = i]}$$
 - Iterate until convergence

Illustration of K-means

8



Jordan (2003)

Recall: More principled approach using statistics

9

- K-means assigns each point to exactly one cluster
 - So called "hard assignment"
 - Can also perform "soft assignment": probabilistic
 - Use probability that each data point is a member of each cluster
- Consider Z_n as latent variables, and x_n observed
- Seek to estimate parameters $(\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K)$
 - Can also estimate Σ_i s
- Will use posterior probability of each x_n being in i :

$$\begin{aligned} \tau_n^i &= p(Z_n = i | x_n, \theta) = \frac{p(x_n | Z_n = i, \theta_i) p(Z_n = i)}{p(x_n | \theta)} \\ &= \frac{\pi_i \cdot \mathcal{N}(x_n | \mu_i, \Sigma_i)}{\sum_j \pi_j \cdot \mathcal{N}(x_n | \mu_j, \Sigma_j)} \end{aligned}$$

Expectation-maximization for Gaussian mixtures

10

- Algorithm:
 - Make initial guess of $\theta^{(0)} = (\pi_1^{(0)}, \dots, \pi_K^{(0)}, \mu_1^{(0)}, \dots, \mu_K^{(0)})$
 - Will take Σ_i s as known, but these can also be included
 - Update posterior probability of assignments $\tau_n^{i(t)}$ "E step"
 - Here t is the iteration number
 - Note: $\tau_n^{i(t)} = E[I[Z_n = i] | x_n, \theta^{(t)}]$
 - Bernoulli RV: expectation is its probability (previous slide)
 - Equation for $\tau_n^{i(t)}$ relies on current parameters $\pi_i^{(t)}, \mu_i^{(t)}$
 - Update the parameters: "M step"

$$\mu_i^{(t+1)} = \frac{\sum_n x_n \cdot \tau_n^{i(t)}}{\sum_n \tau_n^{i(t)}} \quad \pi_i^{(t+1)} = \frac{1}{N} \sum_n \tau_n^{i(t)}$$
 - Sets parameters to their MLE given hidden variable assignments
 - Iterate until small difference in likelihood between iterations
 - Can prove that in EM, likelihood increases in each iteration

If covariance matrix unknown

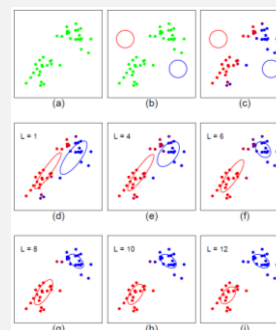
11

- Can update covariance matrices Σ_i as

$$\Sigma_i^{(t+1)} = \frac{\sum_{n=1}^N \tau_n^{i(t)} (x_n - \mu_i^{(t+1)}) (x_n - \mu_i^{(t+1)})^T}{\sum_{n=1}^N \tau_n^{i(t)}}$$

Example: $K = 2, L = \#$ iterations

12



Jordan (2003)

Recall: Desired uses of HMMs

13

Evaluation:

- Given: observed x and HMM specification

Question: what is the joint probability of x and a given z ?

Question: what is the likelihood of x based on the HMM?

Decoding:

- Given: observed x and HMM

Question: what sequence of hidden states produced x ?

- Viterbi decoding: most likely hidden state sequence

- Posterior probability of hidden states: probability of each state z_i producing each x_i

Learning:

- – Given: observed x and HMM without complete probabilities
- Question:** what emission, transition probabilities produced x ?

Hidden Markov models

How do we learn transition, emission probabilities from **labeled** data?

Supervised learning (Review)

14

Recall: Supervised learning MLEs

15

- Determining θ from labeled data: fairly simple

- Given x_1, x_2, \dots, x_N with correct z_1, z_2, \dots, z_N

- Let

$A_{kl} = \# \text{ times transition } z_i = k \rightarrow z_{i+1} = l \text{ occurs for all } i$

$E_k(b) = \# \text{ times } z_i = k \text{ emits } b$

- Can show that MLE for the parameters are

$$a_{kl} = \frac{A_{kl}}{\sum_{l'} A_{kl'}} \text{ and } e_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')}$$

- In practice often use pseudocounts (priors)
 - Avoids 0 probabilities and can aid in overfitting small datasets

Hidden Markov models

How do we learn transition, emission probabilities from **unlabeled** data?

Unsupervised learning: Baum-Welch

16

Baum-Welch algorithm overview

17

1. Guess at the values of $\theta^{(0)} = (a_{kl}^{(0)}, e_k^{(0)}(b))$
2. Estimate $A_{kl}^{(t)}, E_k^{(t)}(b)$ counts considering all paths
 - Sum the probabilities of all possible transitions / emissions
 - Not using fixed counts of labeled sequences
 - **Are expected values** and generally real valued (not integers)
 - The “E step”
3. Update $a_{kl}^{(t+1)}, e_k^{(t+1)}(b)$ using $A_{kl}, E_k(b)$
 - The “M step”
 - Uses MLE – identical to supervised learning for this
4. Repeat

Calculating $A_{kl}^{(t)}, E_k^{(t)}(b)$ by expectation

18

- Can set expected counts of transitioning using every step i as:

$$A_{kl}^{(t)} = E[A_{kl} | x, \theta^{(t)}] = \sum_i P(z_i = k, z_{i+1} = l | x, \theta^{(t)})$$

$$= \sum_i \frac{f_k^{(t)}(i) a_{kl}^{(t)} e_l^{(t)}(x_{i+1}) b_l^{(t)}(i+1)}{P(x | \theta^{(t)})}$$

- Here $f_k^{(t)}(i)$ is forward probability using $\theta^{(t)}$ as parameters
 - All other values with superscript (t) are also those using $\theta^{(t)}$

- Similarly have

$$E_k^{(t)}(b) = \frac{1}{P(x | \theta^{(t)})} \sum_i I[x_i = b] \cdot f_k^{(t)}(i) \cdot b_k^{(t)}(i)$$

Transition probability conditional on data \mathbf{x} , current parameters $\theta^{(t)}$

19

- Consider the probability of transitioning from k to l given \mathbf{x} :

$$\begin{aligned} P(z_i = k, z_{i+1} = l | \mathbf{x}, \theta^{(t)}) &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} \cdot P(z_i = k, z_{i+1} = l, \mathbf{x} | \theta^{(t)}) \\ &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} P(x_1, \dots, x_i, z_i = k, z_{i+1} = l, x_{i+1}, \dots, x_N | \theta^{(t)}) \\ &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} P(z_{i+1} = l, x_{i+1}, \dots, x_N | z_i = k, \theta^{(t)}) P(x_1, \dots, x_i, z_i = k | \theta^{(t)}) \\ &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} P(z_{i+1} = l, x_{i+1}, \dots, x_N | z_i = k, \theta^{(t)}) f_k^{(t)}(i) \\ &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} P(x_{i+2}, \dots, x_N | z_{i+1} = l, \theta^{(t)}) P(x_{i+1} | z_{i+1} = l, \theta^{(t)}) P(z_{i+1} = l | z_i = k, \theta^{(t)}) f_k^{(t)}(i) \\ &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} f_k^{(t)}(i) a_{kl}^{(t)} e_l^{(t)}(x_{i+1}) b_l^{(t)}(i+1) \\ \text{So } P(z_i = k, z_{i+1} = l | \mathbf{x}, \theta^{(t)}) &= \frac{1}{P(\mathbf{x} | \theta^{(t)})} f_k^{(t)}(i) a_{kl}^{(t)} e_l^{(t)}(x_{i+1}) b_l^{(t)}(i+1) \end{aligned}$$

- Why does this lead to different $a_{kl}^{(t+1)}$? Influenced by data \mathbf{x}

Can sum over multiple training sequences when available

20

- Can set expected counts of transitioning using every step i and all sequences \mathbf{x} as:

$$A_{kl}^{(t)} = \sum_{\mathbf{x}} \sum_i \frac{f_k^{(t)}(i) a_{kl}^{(t)} e_l^{(t)}(x_{i+1}) b_l^{(t)}(i+1)}{P(\mathbf{x} | \theta^{(t)})}$$

- Similarly have

$$E_k^{(t)}(b) = \sum_{\mathbf{x}} \frac{1}{P(\mathbf{x} | \theta^{(t)})} \sum_i I[x_i = b] \cdot f_k^{(t)}(i) \cdot b_k^{(t)}(i)$$

Baum-Welch algorithm: EM for HMMs

21

- Initialize $\theta^{(0)} = (a_{kl}^{(0)}, e_k^{(0)}(b))$ for all k, l, b
- For $t = 0$ until end:
 - Run forward algorithm: compute $f_k^{(t)}(i)$ for all k, i using $\theta^{(t)}$
 - Run backward algorithm: compute $b_k^{(t)}(i)$ for all k, i using $\theta^{(t)}$
 - Calculate $A_{kl}^{(t)}, E_k^{(t)}(b)$ using earlier equations "E step"
 - Calculate new $\theta^{(t+1)} = (a_{kl}^{(t+1)}, e_k^{(t+1)}(b))$ for all k, l, b as:

$$a_{kl}^{(t+1)} = \frac{A_{kl}^{(t)}}{\sum_{l'} A_{kl'}^{(t)}} \text{ and } e_k^{(t+1)}(b) = \frac{E_k^{(t)}(b)}{\sum_{b'} E_k^{(t)}(b')} \quad \text{"M step"}$$
 - Get updated likelihood $P(\mathbf{x} | \theta^{(t+1)}) \leftarrow$ **Guaranteed increased**
- Terminate when $\Delta P(\mathbf{x} | \theta)$ small or fixed # iterations

Notes on Baum-Welch

22

- Runtime complexity:
 - # iterations $\times O(K^2 N)$
- Will always increase $P(\mathbf{x} | \theta)$
- But may not find global maximum
 - Gives local optimum
 - Initial assignment can help guide solution to global maximum, but generally no way to guarantee this
- If large # of parameters θ , can overfit / overtrain
 - As we've talked about, overfitting \Rightarrow too heavily tuned to data

Recall: Randomized training algorithm

23

- Initialize $\theta^{(0)} = (a_{kl}^{(0)}, e_k^{(0)}(b))$ for all k, l, b
- Iterate:
 - Run forward algorithm: compute $f_k^{(t)}(i)$ for all k, i
 - Sample n hidden state vectors \mathbf{z} from $P(\mathbf{z} | \mathbf{x}, \theta^{(t)})$
 - Calculate $A_{kl}^{(t)}, E_k^{(t)}(b)$ using all paths \mathbf{z} [as in supervised]
 - Calculate new $\theta^{(t+1)} = (a_{kl}^{(t+1)}, e_k^{(t+1)}(b))$ for all k, l, b [MLE]
- Terminate when $\Delta P(\mathbf{x} | \theta)$ small or fixed # iterations
- Notes:
 - May not converge as quickly as EM: randomized
 - But stochastic so may yield higher $P(\mathbf{x} | \theta)$

Viterbi training

24

Algorithm

- Initialize $\theta^{(0)} = (a_{kl}^{(0)}, e_k^{(0)}(b))$ for all k, l, b [same]
- Iterate until convergence:
 - Run Viterbi decoding to get $\mathbf{z}^* = \arg \max_{\mathbf{z}} P(\mathbf{x}, \mathbf{z} | \theta^{(t)})$
 - Calculate $A_{kl}^{(t)}, E_k^{(t)}(b)$ based on \mathbf{z}^* [as in supervised]
 - Calculate new $\theta^{(t+1)} = (a_{kl}^{(t+1)}, e_k^{(t+1)}(b))$ for all k, l, b [MLE]
- Here convergence is guaranteed
 - Number of paths is discrete, will eventually be unchanged
- Does **not** maximize $P(\mathbf{x} | \theta)$ in general
 - Does not perform as well as Baum-Welch for most problems

Final notes

- Summary:
 - Expectation-maximization (EM)
 - Gaussian mixture example
 - Baum-Welch algorithm: EM for HMMs
 - Viterbi training