**6.047/6.878/HST.507**
**Computational Biology: Genomes, Networks, Evolution**

# Lecture 07

# Hidden Markov Models
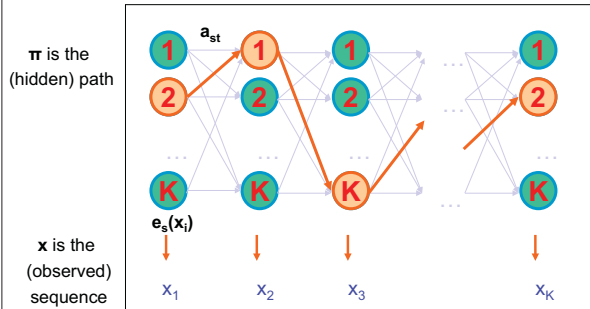# Part II



---

## Module II: Modeling genes and gene expression



- Computational Foundations
  - Hidden Markov Models (HMMs): Central tool in CS
  - Decoding, evaluation, parsing, likelihood, scoring
  - Unsupervised Learning: Expectation Maximization
  - Supervised learning: generative/discriminative models
- Biological frontiers:
  - PS2: Modeling conservation, GC content, CpG islands
  - L6/L7: Genome annotation and parsing
  - L8: Gene expression analysis: cluster genes/conditions
  - L9: Regulatory motif discovery: EM, gibbs sampling, info

---

## Goals for today: HMMs, part II

- Review: Three algorithms from last time
  - Markov Chains and Hidden Markov Models
  - Increasing the 'state' space / adding memory
  - Calculating likelihoods $P(x,\pi)$
  - Viterbi algorithm: Find $\pi^* = \text{argmax}_\pi P(x,\pi)$
- Counting over all paths
  - Forward algorithm: Find $P(x)$, over all paths
  - Model comparison: ex: "CpGs" vs. "Gs and Cs"
- Posterior decoding: Another way of 'parsing'
  - Find most likely state $k_i$, overall all possible paths
- Learning (ML training, Baum-Welch, Viterbi training)
  - Supervised: Find $e_i(.)$ and $a_{ij}$ given labeled sequence
  - Unsupervised: given only x → annotation + params

---

## Markov Chains & Hidden Markov Models



- Markov Chain
  - Q: states
  - p: initial state probabilities
  - A: transition probabilities

- HMM
  - Q: states
  - V: observations
  - p: initial state probabilities
  - A: transition probabilities
  - E: emission probabilities

---

## HMM nomenclature for this course



Transitions: $a_{kl}=P(\pi_i=l|\pi_{i-1}=k)$
Transition probability from state $k$ to state $l$

Emissions: $e_k(x_i)=P(x_i|\pi_i=k)$
Emission probability of symbol $x_i$ from state $k$

- Vector $x$ = Sequence of observations
- Vector $\pi$ = Hidden path (sequence of hidden states)
- Transition matrix $A=a_{kl}=$ probability of $k \rightarrow l$ state transition
- Emission vector $E=e_k(x_i)$ = prob. of observing $x_i$ from state k
- Bayes's rule: Use $P(x_i|\pi_i=k)$ to estimate $P(\pi_i=k|x_i)$

## One path / All paths

| | One path | All paths |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$<br><br>Prob of a path, emissions | 2. Scoring x, all paths<br><br>$P(x) = \sum_{\pi} P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = \text{argmax}_{\pi}\, P(x,\pi)$<br><br>Most likely path | 4. Posterior decoding<br><br>$\pi^{\wedge} = \{\pi_i \mid \pi_i = \text{argmax}_k\, \sum_{\pi} P(\pi_i = k \mid x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given $\pi$<br>$\Lambda^* = \text{argmax}_{\Lambda}\, P(x,\pi \mid \Lambda)$<br>6. Unsupervised learning.<br>$\Lambda^* = \text{argmax}_{\Lambda}\, \text{max}_{\pi} P(x,\pi \mid \Lambda)$<br>Viterbi training, best path | 6. Unsupervised learning<br><br>$\Lambda^* = \text{argmax}_{\Lambda}\, \sum_{\pi} P(x,\pi \mid \Lambda)$<br><br>Baum-Welch training, over all paths |

---

**1. Scoring probability of a path + sequence**

Multiply emissions, transitions

---

## Probability of given path p, emissions x



$\pi$ is the (hidden) path

$x$ is the (observed) sequence

- $P(x,\pi) = a_{0\pi_1} * \prod_i \left( e_{\pi_i}(x_i) \right) \times a_{\pi_i \pi_{i+1}}$

  start    emission    transition

---

## Example: One particular P vs. B assignment



L: (path with P and B states, transitions 0.75, 0.75, 0.15, 0.25, 0.85, 0.85, 0.85; emissions 0.25, 0.25, 0.25, 0.42, 0.42, 0.30, 0.25, 0.25)

S:   G    C    A    A    A    T    G    C

$P = P(G \mid B)P(B_1 \mid B_0)P(C \mid B)P(B_2 \mid B_1)P(A \mid B)P(P_3 \mid B_2)...P(C \mid B_7)$

$= (0.85)^3 \times (0.25)^6 \times (0.75)^2 \times (0.42)^2 \times 0.30 \times 0.15$

$= 6.7 \times 10^{-7}$

---

---

**3. Decoding: find the most likely path**

Viterbi algorithm

## Finding the most likely path



- Find path $\pi^*$ that maximizes total joint probability $P[x, \pi]$

- $P(x,\pi) = a_{0\pi_1} * \prod_i e_{\pi_i}(x_i) \times a_{\pi_i\pi_{i+1}}$

  start    emission    transition

## Calculate maximum P(x,π) recursively

**Viterbi algortithm**

Define $V_k(i)$ = Probability of the most likely path through state $\pi_i = k$

Compute $V_k(i+1)$ recursively, as a function of $\max_{k'} \{ V_{k'}(i) \}$



- Assume we know $V_j$ for the previous time step (i-1)

- Calculate $V_k(i) = e_k(x_i) * \max_j ( V_j(i-1) \times a_{jk} )$

  current max    this emission    max ending in state j at step i    Transition from state j

  all possible previous states j

## The Viterbi Algorithm

State 1
2

$V_k(i)$

K

$x_1 \ x_2 \ x_3 \ \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots x_N$

Input: $x = x_1 \ldots\ldots x_N$

**Initialization:**
$V_0(0) = 1, V_k(0) = 0$, for all $k > 0$

**Iteration:**
$V_k(i) = e_k(x_i) \times \max_j a_{jk} V_j(i-1)$

**Termination:**
$P(x, \pi^*) = \max_k V_k(N)$

**Traceback:**
Follow max pointers back

**In practice:**
Use log scores for computation

**Running time and space:**
Time: $O(K^2N)$
Space: $O(KN)$

| | One path | All paths |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$<br><br>Prob of a path, emissions | 2. Scoring x, all paths<br><br>$P(x) = \sum_\pi P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = \text{argmax}_\pi P(x,\pi)$<br><br>Most likely path | 4. Posterior decoding<br><br>$\pi^\wedge = \{\pi_i \mid \pi_i = \text{argmax}_k \sum_\pi P(\pi_i = k\mid x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given $\pi$<br>$\Lambda^* = \text{argmax}_\Lambda P(x,\pi\mid\Lambda)$<br>6. Unsupervised learning.<br>$\Lambda^* = \text{argmax}_\Lambda \max_\pi P(x,\pi\mid\Lambda)$<br>Viterbi training, best path | 6. Unsupervised learning<br><br>$\Lambda^* = \text{argmax}_\Lambda \sum_\pi P(x,\pi\mid\Lambda)$<br><br>Baum-Welch training, over all paths |

## 2. Model evaluation:
## Total P(x|M), summed over all paths

Forward algorithm

## Simple: Given the model, generate some sequence x



Given a HMM, we can generate a sequence of length n as follows:
1. Start at state $\pi_1$ according to prob $a_{0\pi_1}$
2. Emit letter $x_1$ according to prob $e_{\pi_1}(x_1)$
3. Go to state $\pi_2$ according to prob $a_{\pi_1\pi_2}$
4. … until emitting $x_n$

We have some sequence x that can be emitted by p. Can calculate its likelihood.
However, in general, many different paths may emit this same sequence x.
How do we find the <u>total probability</u> of generating a given x, over any path?

## Complex: Given x, was it generated by the model?



Given a sequence x,
What is the probability that x was generated by the model (using any path)?

– $P(x) = \sum_{\pi} P(x,\pi)$

- Challenge: exponential number of paths

- (cheap) alternative:
  – Calculate probability over maximum (Viterbi) path $\pi^*$
- (real) solution
  – Calculate sum iteratively using <u>principles</u> of dynamic programming

---

## The Forward Algorithm – derivation

Define the forward probability:

$$f_l(i) = P(x_1 \ldots x_i, \pi_i = l)$$

$$= \sum_{\pi 1 \ldots \pi i-1} P(x_1 \ldots x_{i-1}, \pi_1, \ldots, \pi_{i-2}, \pi_{i-1}, \pi_i = l)\, e_l(x_i)$$

$$= \sum_k \boxed{\sum_{\pi 1 \ldots \pi i-2} P(x_1 \ldots x_{i-1}, \pi_1, \ldots, \pi_{i-2}, \pi_{i-1}=k)}\, a_{kl}\, e_l(x_i)$$

$$= \sum_k \boxed{f_k(i-1)}\, a_{kl}\, e_l(x_i)$$

$$= e_l(x_i) \sum_k \boxed{f_k(i-1)}\, a_{kl}$$

---

## Calculate total probability $\Sigma_\pi P(x,\pi)$ recursively



- Assume we know $f_j$ for the previous time step (i-1)

- Calculate $f_k(i) = e_k(x_i) * \text{sum}_j ( f_j(i-1) \times a_{jk} )$

  current max · this emission · sum ending in state j at step i · transition from state j

  every possible previous state j

---

## The Forward Algorithm



$x_1 \; x_2 \; x_3 \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots x_N$

Input: x = x1……xN

**Initialization:**
  $f_0(0)=1$, $f_k(0) = 0$, for all k > 0

**Iteration:**
  $f_k(i) = e_k(x_i) \times \text{sum}_j \, a_{jk}\, f_j(i-1)$

**Termination:**
  $P(x, \pi^*) = \text{sum}_k \, f_k(N)$

**In practice:**
  Sum of log scores is difficult
  → approximate exp(1+p+q)
  → scaling of probabilities

**Running time and space:**
  Time:   $O(K^2 N)$
  Space:  $O(K)$

---

## Application:
## Distinguishing between two models

HMM1:  Promoters = only Cs and Gs matter
HMM2: Promoters = it's actually CpGs that matter
("C"-phosphate-"G", i.e. on the same strand!)

**(increasing the state space)**

---

## In the human genome, CpG islands matter!



- Regions of regulatory importance in promoters of many genes
  – Defined by their methylation state (epigenetic information)
  – CpGs more important than simply the abundance of Cs and Gs
  – Provide evidence of methylation state!
- Methylation process in the human genome (form of silencing):
  – Methylation signature: high chance of methyl-C mutating to T in CpG
    → CpG dinucleotides are rare, throughout the genome
  – BUT methylation is suppressed for active promoters
    → CpG dinucleotides are much more frequent than elsewhere
    • Such regions are called **CpG islands**
    • A few hundred to a few thousand bases long
- Problems:
  – Given a short sequence, does it come from a CpG island or not?
  – How to find the CpG islands in a long sequence
- How do we encode this in an hidden Markov model?

## Increasing the state of the system (looking back)

- Markov Models are memory-less
  - In other words, all memory is encoded in the states
  - To remember additional information, augment state
- Our first HMM had minimal memory
  - State, emissions, only depend on **current** state
  - Current state only encoded **one** previous nucleotide
- How do you count **di**-nucleotide frequencies?
  - CpG islands: di-nucleotides
  - Codon triplets: tri-nucleotides
  - Di-codon frequencies: six nucleotides
- ➔ Expanding the number of states



| | + | - |
|---|---|---|
| | A: .1 | A: 1/4 |
| | C: .3 | C: 1/4 |
| | G: .4 | G: 1/4 |
| | T: .2 | T: 1/4 |

---

## Modeling CpG islands: incorporating memory



$e_B(1) = \frac{1}{4}$

- Markov Chain
  - Q: states
  - p: initial state probabilities
  - A: transition probabilities
- HMM
  - Q: states
  - V: observations
  - p: initial state probabilities
  - A: transition probabilities
  - E: emission probabilities

---

## Example 2: CpG islands: incorporating memory



| | A: .1 | A: 1/4 |
|---|---|---|
| | C: .3 | C: 1/4 |
| | G: .4 | G: 1/4 |
| | T: .2 | T: 1/4 |

---

## HMM for CpG islands



- Build a single model that combines two such Markov chains:
  - **'+'** states: A+, C+, G+, T+
    - Emit symbols: A, C, G, T in CpG islands
  - **'-'** states: A-, C-, G-, T-
    - Emit symbols: A, C, G, T in non-islands
- Emission probabilities distinct for the '+' and the '-' states
  - Infer most likely set of states, giving rise to observed emissions
  - ➔ 'Paint' the sequence with + and - states

**Why we need so many states…**
**In our simple GC-content example, we only had 2 states (+|-)**
**Why do we need 8 states here: 4 CpG+ / 4 CpG- ?**
**➔ Encode 'memory' of previous state: nucleotide transitions**

---

## Training emission parameters for CpG+/CpG- states



- Count di-nucleotide frequencies:
  - 16 possible di-nucleotides. 16 transition parameters.
  - Alternative: 16 states, each emitting di-nucleotide
- Derive two Markov chain models:
  - **'+' model**: from the CpG islands
  - **'-' model**: from the remainder of sequence
- Transition probabilities for each model:
  - Encode differences in di-nucleotide frequencies

| + | A | C | G | T |
|---|---|---|---|---|
| A | .180 | .274 | .426 | .120 |
| C | .171 | .368 | .274 | .188 |
| G | .161 | .339 | .375 | .125 |
| T | .079 | .355 | .384 | .182 |

| - | A | C | G | T |
|---|---|---|---|---|
| A | .300 | .205 | .285 | .210 |
| C | .322 | .298 | .078 | .302 |
| G | .248 | .246 | .298 | .208 |
| T | .177 | .239 | .292 | .292 |

---

| | **One path** | **All paths** |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$ ✓<br><br>Prob of a path, emissions | 2. Scoring x, all paths<br><br>✓ $P(x) = \sum_\pi P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = \text{argmax}_\pi P(x,\pi)$ ✓<br><br>Most likely path | 4. Posterior decoding<br><br>$\pi^\wedge = \{\pi_i \mid \pi_i = \text{argmax}_k \sum_\pi P(\pi_i = k \mid x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given π<br>$\Lambda^* = \text{argmax}_\Lambda P(x,\pi \mid \Lambda)$<br>6. Unsupervised learning.<br>$\Lambda^* = \text{argmax}_\Lambda \max_\pi P(x,\pi \mid \Lambda)$<br>Viterbi training, best path | 6. Unsupervised learning<br><br>$\Lambda^* = \text{argmax}_\Lambda \sum_\pi P(x,\pi \mid \Lambda)$<br><br>Baum-Welch training, over all paths |

## 4. Decoding, all paths

Find the likelihood an emission $x_i$ is generated by a state

---

## Calculate most probable label at a single position



Sum over all paths

π:

x:  G   C   A   A   A   T   G   C

**P(Label$_i$=B|x)**

- Calculate most probable label, $L^*_i$, at each position i
- Do this for all N positions gives us $\{L^*_1, L^*_2, L^*_3 \ldots L^*_N\}$
- How much information have we observed? Three settings:
  - **Observed nothing: Use prior information**
  - **Observed only character at position i: Prior + emission probability**
  - **Observed entire sequence: Posterior decoding**

---

## Calculate P(π$_7$= CpG+ | x$_7$=G)

- With no knowledge (no characters)
  - Simply time spent in markov chain states
  - P( $\pi_i$=k ) = most likely state (**prior**)

- With very little knowledge (just that character)
  - Time spent, adjusted for different emission probs.
  - Use Bayes rule to change inference directionality
  - P( $\pi_i$=k | $x_i$=G ) = P($\pi_i$=κ) * P($x_i$=G|$\pi_i$=k) / P($x_i$=G)

- With knowledge of entire sequence (all characters)
  - P( $\pi_i$=k | x=AGCGCG…GATTATCGTCGTA)
  - Sum over all paths that emit 'G' at position 7
  - ➔ **Posterior** decoding

---

## Motivation for the Backward Algorithm

We want to compute

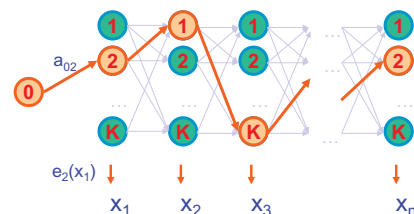P($\pi_i$ = k | x), the probability distribution on the i$^{th}$ position, given x

We start by computing

P($\pi_i$ = k, x) = P($x_1$…$x_i$, $\pi_i$ = k, $x_{i+1}$…$x_N$)

= P($x_1$…$x_i$, $\pi_i$ = k) P($x_{i+1}$…$x_N$ | $x_1$…$x_i$, $\pi_i$ = k)

= P($x_1$…$x_i$, $\pi_i$ = k) P($x_{i+1}$…$x_N$ | $\pi_i$ = k)

**Forward, f$_k$(i)**    **Backward, b$_k$(i)**

---

## The Backward Algorithm – derivation

Define the backward probability:

$b_k(i)$ = P($x_{i+1}$…$x_N$ | $\pi_i$ = k)

= $\Sigma_{\pi i+1…\pi N}$ P($x_{i+1}$,$x_{i+2}$, …, $x_N$, $\pi_{i+1}$, …, $\pi_N$ | $\pi_i$ = k)

= $\Sigma_l \Sigma_{\pi i+1…\pi N}$ P($x_{i+1}$,$x_{i+2}$, …, $x_N$, $\pi_{i+1}$ = l, $\pi_{i+2}$, …, $\pi_N$ | $\pi_i$ = k)

= $\Sigma_l e_l(x_{i+1}) a_{kl} \Sigma_{\pi i+1…\pi N}$ P($x_{i+2}$, …, $x_N$, $\pi_{i+2}$, …, $\pi_N$ | $\pi_{i+1}$ = l)

= $\Sigma_l e_l(x_{i+1}) a_{kl} b_l(i+1)$

---

## Calculate total end probability recursively



hidden states

observations   $x_i$   $x_{i+1}$

- Assume we know $b_l$ for the next time step (i+1)

- Calculate $b_k(i)$ = sum$_l$ ( $e_l(x_{i+1})$ × $a_{kl}$ × $b_l(i+1)$ )

current max         next emission     transition to next state     prob sum from state l to end

sum over all possible next states

## The Backward Algorithm

State 1
2

$b_k(i)$

K

$x_1 \; x_2 \; x_3 \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots x_N$

Input: $x = x1\ldots\ldots xN$

**Initialization:**
$b_k(N) = a_{k0}$, for all k

**Iteration:**
$b_k(i) = \Sigma_l \, e_l(x_{i+1}) \, a_{kl} \, b_l(i+1)$

**Termination:**
$P(x) = \Sigma_l \, a_{0l} \, e_l(x_1) \, b_l(1)$

**In practice:**
Sum of log scores is difficult
→ approximate exp(1+p+q)
→ scaling of probabilities

**Running time and space:**
Time:  $O(K^2N)$
Space:  $O(K)$

---

## Putting it all together:  Posterior decoding

State 1
2

P(k)

K

$x_1 \; x_2 \; x_3 \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots x_N$

- $P(k) = P(\; \pi_i = k \mid x \;) = f_k(i)*b_k(i) \, / \, P(x)$
  - Probability that $i^{th}$ state is k, given all emissions x
- Posterior decoding
  - Define most likely state for every of sequence x
  - $\pi^{\wedge}_i = argmax_k \, P(\pi_i = k \mid x)$
- Posterior decoding 'path' $\pi^{\wedge}_i$
  - For classification, more informative than Viterbi path $\pi^*$
    - More refined measure of "which hidden states" generated x
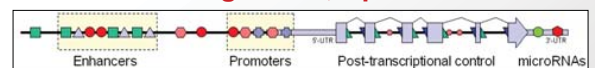  - However, it may give an invalid sequence of states
    - Not all j→k transitions may be possible

---

## Summary this far

- Generative model.  Hidden states, observed emissions.
  - Generate a random sequence
    - Choose random transition, choose random emission (#0)
- Scoring the likelihood of a sequence
  - Calculate likelihood of annotated path and sequence
    - Multiply emission and transition probabilities (#1)
  - Without specifying a path, total probability of generating x
    - Sum probabilities over all paths
    - Forward algorithm (#3)
- Decoding:  Finding the most likely path, given a sequence
  - What is the most likely path generating entire sequence?
    - Viterbi algorithm (#2)
  - What is the most probable state at each time step?
    - Forward + backward algorithms, posterior decoding (#4)
- Next:  Learning (#5 and #6)

---

|  | **One path** | **All paths** |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$  ✓<br><br>Prob of a path, emissions | 2. Scoring x, all paths<br><br>✓  $P(x) = \Sigma_\pi \, P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = argmax_\pi \, P(x,\pi)$  ✓<br><br>Most likely path | 4. Posterior decoding<br><br>✓  $\pi^\wedge = \{\pi_i \mid \pi_i = argmax_k \, \Sigma_\pi P(\pi_i = k\mid x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given $\pi$<br>$\Lambda^* = argmax_\Lambda \, P(x,\pi\mid\Lambda)$ ⭐<br>6. Unsupervised learning.<br>$\Lambda^* = argmax_\Lambda \, max_\pi P(x,\pi\mid\Lambda)$ ⭐<br>Viterbi training, best path | 6. Unsupervised learning<br>⭐<br>$\Lambda^* = argmax_\Lambda \, \Sigma_\pi P(x,\pi\mid\Lambda)$<br><br>Baum-Welch training, over all paths |

---

# Learning: How to train an HMM

**Transition probabilities**
e.g. $P(P_{i+1}\mid B_i)$ – the probability of entering a pathogenicity island from background DNA

**Emission probabilities**
i.e. the nucleotide frequencies for background DNA and pathogenicity islands

$P(L_{i+1}\mid L_i)$

B          P

P(S|B)      P(S|P)

---

## Two learning scenarios

Case 1. Estimation when the "right answer" is known

**Examples:**
GIVEN: a genomic region $x = x_1 \ldots x_{1,000,000}$ where we have good (experimental) annotations of the CpG islands

Case 2. Estimation when the "right answer" is unknown

**Examples:**
GIVEN: the porcupine genome; we don't know how frequent are the CpG islands there, neither do we know their composition

QUESTION: Update the parameters θ of the model to maximize $P(x\mid\theta)$

## Two types of learning: Supervised / Unsupervised

### 5. Supervised learning

infer model parameters given **labeled** training data

- GIVEN:
  - a HMM M, with unspecified transition/emission probs.
  - labeled sequence x,
- FIND:
  - parameters $\theta = (E_i, A_{ij})$ that maximize $P[x \mid \theta]$
- ➔ Simply count frequency of each emission and transition, as observed in the training data

### 6. Unsupervised learning

infer model parameters given **unlabelled** training data

- GIVEN:
  - a HMM M, with unspecified transition/emission probs.
  - unlabeled sequence x,
- FIND:
  - parameters $\theta = (E_i, A_{ij})$ that maximize $P[x \mid \theta]$
- ➔ Viterbi training:
  guess parameters, find optimal Viterbi path (#2), update parameters (#5), iterate
- ➔ Baum-Welch training:
  guess parameters, sum over all paths (#4), update parameters (#5), iterate

---

## 5: Supervised learning

Estimate model parameters
based on **labeled** training data

---

## Case 1.    When the right answer is known

Given $x = x_1 \ldots x_N$
for which the true $\pi = \pi_1 \ldots \pi_N$ is known,

**Define:**

$A_{kl}$       = # times $k \to l$ transition occurs in $\pi$
$E_k(b)$     = # times state k in $\pi$ emits b in x

We can show that the maximum likelihood parameters $\theta$ are:

$$a_{kl} = \frac{A_{kl}}{\Sigma_i \; A_{ki}} \qquad e_k(b) = \frac{E_k(b)}{\Sigma_c \; E_k(c)}$$

---

## Learning From Labelled Data
## ➔ Maximum Likelihood Estimation

**If we have a sequence that has islands marked, we can simply count**



**L:**

S:    G    C    A    A    A    T    G    C

| $P(L_{i+1}\mid L_i)$ | $B_{i+1}$ | $P_{i+1}$ |
|---|---|---|
| $B_i$ | | |
| $P_i$ | | |

P(S|B)
A:
T:
G:
C:
!

P(S|P)
A:
T:
G:
C:
ETC..

---

## Case 1.    When the right answer is known

**Intuition:** When we know the underlying states,
Best estimate is the average frequency of
transitions & emissions that occur in the training data

**Drawback:**
Given little data, there may be **overfitting**:
$P(x|\theta)$ is maximized, but $\theta$ is unreasonable
**0 probabilities – VERY BAD**

**Example:**
Given 10 nucleotides, we observe
     x = C, A, G, G, T, C, C, A, T, C
     $\pi$ = P, P, P, p, p, P, P, P, P, P
Then:
     $a_{PP}$ = 1;    $a_{PB}$ = 0
     $e_P(A)$ = .2;
     $e_P(C)$ = .4;
     $e_P(G)$ = .2;
     $e_P(T)$ = .2

---

## Pseudocounts

Solution for small training sets:

Add pseudocounts

$A_{kl}$       = # times $k \to l$ transition occurs in $\pi$  + $r_{kl}$
$E_k(b)$     = # times state k in $\pi$ emits b in x   + $r_k(b)$

$r_{kl}$, $r_k(b)$ are pseudocounts representing our prior belief

Larger pseudocounts $\Rightarrow$ Strong priof belief

Small pseudocounts ($\epsilon < 1$): just to avoid 0 probabilities

## Slide 1: Example: Training Markov Chains for CpG islands



- Training Set:
  - set of DNA sequences w/ known CpG islands
- Derive two Markov chain models:
  - '+' model: from the CpG islands
  - '-' model: from the remainder of sequence
- Transition probabilities for each model:

| + | A | C | G | T |
|---|---|---|---|---|
| A | .180 | .274 | .426 | .120 |
| C | .171 | .368 | .274 | .188 |
| G | .161 | .339 | .375 | .125 |
| T | .079 | .355 | .384 | .182 |

| - | A | C | G | T |
|---|---|---|---|---|
| A | .300 | .205 | .285 | .210 |
| C | .322 | .298 | .078 | .302 |
| G | .248 | .246 | .298 | .208 |
| T | .177 | .239 | .292 | .292 |

$$a_{st}^+ = \frac{c_{st}^+}{\sum_{t'} c_{st'}^+}$$

$c_{st}^+$ is the number of times letter $t$ followed letter $s$ <u>inside</u> the CpG islands

$$a_{st}^- = \frac{c_{st}^-}{\sum_{t'} c_{st'}^-}$$

$c_{st}^-$ is the number of times letter $t$ followed letter $s$ <u>outside</u> the CpG islands

## Slide 2: 6: Unsupervised learning

Estimate model parameters based on **unlabeled** training data

## Slide 3: Unlabelled Data

**How do we know how to count?**



L:

S:  G  C  A  A  A  T  G  C

?

$P(L_{i+1}|L_i)$

|  | $B_{i+1}$ | $P_{i+1}$ | End |
|---|---|---|---|
| $B_i$ | | | |
| $P_i$ | | | |
| Start | | | |

P(S|B)

A:
T:
G:
C:

P(S|P)

A:
T:
G:
C:

## Slide 4: Unlabeled Data



L:

S:  G  C  A  A  A  T  G  C

An idea:
1. Imagine we start with some parameters
2. We *could* calculate the most likely path, P*, given those parameters and S
3. We *could* then use P* to update our parameters by maximum likelihood
4. And iterate (to convergence)

$P(L_{i+1}|L_i)^0$   $P(S|B)^0$   $P(S|P)^0$

$P(L_{i+1}|L_i)^1$   $P(S|B)^1$   $P(S|P)^1$

$P(L_{i+1}|L_i)^2$   $P(S|B)^2$   $P(S|P)^2$

…

$P(L_{i+1}|L_i)^K$   $P(S|B)^K$   $P(S|P)^K$

## Slide 5: Learning case 2.    When the right answer is unknown

We don't know the true $A_{kl}$, $E_k(b)$

**Idea:**
- We estimate our "best guess" on what $A_{kl}$, $E_k(b)$ are (M step, maximum-likelihood estimation)
- We update the probabilistic parse of our sequence, based on these parameters (E step, expected probability of being in each state given parameters)
- We repeat

**Two settings:**
- Simple: Viterbi training (best guest = best path)
- Correct: Expectation maximization (all paths, weighted)

## Slide 6

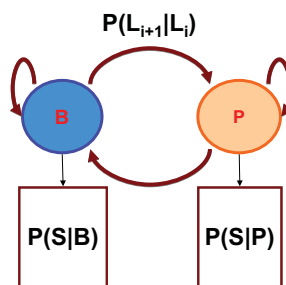|  | **One path** | **All paths** |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$<br><br>Prob of a path, emissions ✓ | ✓ 2. Scoring x, all paths<br><br>$P(x) = \sum_\pi P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = \text{argmax}_\pi P(x,\pi)$<br><br>Most likely path ✓ | ✓ 4. Posterior decoding<br><br>$\pi^\wedge = \{\pi_i \mid \pi_i = \text{argmax}_k \sum_\pi P(\pi_i = k|x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given π<br>$\Lambda^* = \text{argmax}_\Lambda P(x,\pi|\Lambda)$ ✓<br><br>6. Unsupervised learning.<br>$\Lambda^* = \text{argmax}_\Lambda \max_\pi P(x,\pi|\Lambda)$<br>Viterbi training, best path | 6. Unsupervised learning<br><br>$\Lambda^* = \text{argmax}_\Lambda \sum_\pi P(x,\pi|\Lambda)$<br><br>Baum-Welch training, over all paths |

## Simple casae: Viterbi Training

**Initialization:**

Pick the best-guess for model parameters
  (or arbitrary)

**Iteration:**
1. Perform Viterbi, to find $\pi^*$
2. Calculate $A_{kl}$, $E_k(b)$ according to $\pi^*$ + pseudocounts
3. Calculate the new parameters $a_{kl}$, $e_k(b)$

Until convergence

**Notes:**
- Convergence to local maximum guaranteed. Why?
- Does not maximize $P(x \mid \theta)$
- In general, worse performance than Baum-Welch

---

| | **One path** | **All paths** |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$  ✓<br><br>Prob of a path, emissions | 2. Scoring x, all paths<br><br>✓  $P(x) = \sum_\pi P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = \text{argmax}_\pi P(x,\pi)$  ✓<br><br>Most likely path | 4. Posterior decoding<br><br>✓  $\pi^\wedge = \{\pi_i \mid \pi_i = \text{argmax}_k \sum_\pi P(\pi_i = k \mid x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given $\pi$<br>$\Lambda^* = \text{argmax}_\Lambda P(x,\pi \mid \Lambda)$  ✓<br>6. Unsupervised learning.<br>$\Lambda^* = \text{argmax}_\Lambda \max_\pi P(x,\pi \mid \Lambda)$  ✓<br>Viterbi training, best path | 6. Unsupervised learning<br><br>$\Lambda^* = \text{argmax}_\Lambda \sum_\pi P(x,\pi \mid \Lambda)$<br><br>Baum-Welch training, over all paths |

---

# Expectation Maximization (EM)

***The basic idea is the same:***

**1.Use model to estimate missing data (E step)**
**2.Use estimate to update model (M step)**
**3.Repeat until convergence**

**EM is a general approach for learning models (ML estimation) when there is "missing data" Widely used in computational biology**

EM pervasive in computational biology
➡ Rec 3 (SiPhy), Lec 8 (Kmeans), Lec 9 (motifs)

---

# Expectation Maximization (EM)

**1. Initialize parameters randomly**

**2. E Step** Estimate <u>expected probability</u> **of hidden labels**, Q, given current (latest) parameters and observed (unchanging) sequence

$$Q = P(Labels \mid S, params^{t-1})$$

**3. M Step** Choose new <u>maximum likelihood</u> parameters over probability distribution Q, given current probabilistic label assignments

$$params^t = \arg\max_{params} E_Q\left[\log P(S, labels \mid params^{t-1})\right]$$

**4. Iterate**

**P(S|Model) *guaranteed* to increase each iteration**

---

## Case 2.    When the right answer is unknown

Starting with our best guess of a model M, parameters $\theta$:

Given $x = x_1 \ldots x_N$
  for which the true $\pi = \pi_1 \ldots \pi_N$ is unknown,

We can get to a provably more likely parameter set $\theta$

Principle: EXPECTATION MAXIMIZATION

1. Estimate probabilistic parse based on parameters (E step)
2. Update parameters $A_{kl}$, $E_k$ based on probabilistic parse (M step)
3. Repeat 1 & 2, until convergence

---

## Estimating probabilistic parse given params (E step)

To estimate $A_{kl}$:



At each position i:

Find probability transition $k \rightarrow l$ is used:

$P(\pi_i = k, \pi_{i+1} = l \mid x) = [1/P(x)] \times P(\pi_i = k, \pi_{i+1} = l, x_1 \ldots x_N) = Q/P(x)$

where $Q = P(x_1 \ldots x_i, \pi_i = k, \pi_{i+1} = l, x_{i+1} \ldots x_N) =$
  $= P(\pi_{i+1} = l, x_{i+1} \ldots x_N \mid \pi_i = k) P(x_1 \ldots x_i, \pi_i = k) =$
  $= P(\pi_{i+1} = l, x_{i+1} x_{i+2} \ldots x_N \mid \pi_i = k) f_k(i) =$
  $= P(x_{i+2} \ldots x_N \mid \pi_{i+1} = l) P(x_{i+1} \mid \pi_{i+1} = l) P(\pi_{i+1} = l \mid \pi_i = k) f_k(i) =$
  $= b_l(i+1) e_l(x_{i+1}) a_{kl} f_k(i)$

So:    $P(\pi_i = k, \pi_{i+1} = l \mid x, \theta) = \dfrac{f_k(i)\, a_{kl}\, e_l(x_{i+1})\, b_l(i+1)}{P(x \mid \theta)}$

(For one such transition, at time step $i \rightarrow i+1$)

## New parameters given probabilistic parse (M step)

(Sum over all k→l transitions, at any time step i)

So,

$$A_{kl} = \sum_i P(\pi_i = k, \pi_{i+1} = l \mid x, \theta) = \sum_i \frac{\boxed{f_k(i)}\,\boxed{a_{kl}}\,\boxed{e_l(x_{i+1})}\,\boxed{b_l(i+1)}}{P(x \mid \theta)}$$

Similarly,

$$E_k(\boxed{b}) = [1/P(x)]\sum_{\{i \mid x_i = \boxed{b}\}} \boxed{f_k(i)}\,\boxed{b_k(i)}$$

## Dealing with multiple training sequences

(Sum over all training seqs, all k→l transitions, all time steps i)

If we have several training sequences, $x^1, \ldots, x^M$, each of length N,

$$A_{kl} = \sum_x \sum_i P(\pi_i = k, \pi_{i+1} = l \mid x, \theta) = \sum_x \sum_i \frac{\boxed{f_k(i)}\,\boxed{a_{kl}}\,\boxed{e_l(x_{i+1})}\,\boxed{b_l(i+1)}}{P(x \mid \theta)}$$

Similarly,

$$E_k(\boxed{b}) = \sum_x (1/P(x))\sum_{\{i \mid x_i = \boxed{b}\}} \boxed{f_k(i)}\,\boxed{b_k(i)}$$

## The Baum-Welch Algorithm

**Initialization:**
   Pick the best-guess for model parameters
      (or arbitrary)

**Iteration:**
1. Forward
2. Backward
3. → Calculate new log-likelihood $P(x \mid \theta)$   (E step)
4. Calculate $A_{kl}$, $E_k(b)$
5. → Calculate new model parameters $a_{kl}$, $e_k(b)$  (M step)

**GUARANTEED TO BE HIGHER BY EXPECTATION-MAXIMIZATION**

Until $P(x \mid \theta)$ does not change much

## The Baum-Welch Algorithm – comments

Time Complexity:

   # iterations $\times O(K^2N)$

- Guaranteed to increase the log likelihood of the model

   $P(\theta \mid x) = P(x, \theta) / P(x) = P(x \mid \theta) / ( P(x) P(\theta) )$

- Not guaranteed to find <u>globally</u> best parameters

   Converges to local optimum, depending on initial conditions

- Too many parameters / too large model:      Overtraining

|  | **One path** | **All paths** |
|---|---|---|
| **Scoring** | 1. Scoring x, one path<br><br>$P(x,\pi)$<br><br>Prob of a path, emissions | 2. Scoring x, all paths<br><br>$P(x) = \sum_\pi P(x,\pi)$<br><br>Prob of emissions, over all paths |
| **Decoding** | 3. Viterbi decoding<br><br>$\pi^* = \text{argmax}_\pi P(x,\pi)$<br><br>Most likely path | 4. Posterior decoding<br><br>$\pi^\wedge = \{\pi_i \mid \pi_i = \text{argmax}_k \sum_\pi P(\pi_i=k\mid x)\}$<br><br>Path containing the most likely state at any time point. |
| **Learning** | 5. Supervised learning, given $\pi$<br>$\Lambda^* = \text{argmax}_\Lambda P(x,\pi\mid\Lambda)$<br>6. Unsupervised learning.<br>$\Lambda^* = \text{argmax}_\Lambda \max_\pi P(x,\pi\mid\Lambda)$<br>Viterbi training, best path | 6. Unsupervised learning<br><br>$\Lambda^* = \text{argmax}_\Lambda \sum_\pi P(x,\pi\mid\Lambda)$<br><br>Baum-Welch training, over all paths |

## What have we learned ?

- Generative model.  Hidden states, observed emissions.
  - Generate a random sequence
    - Choose random transition, choose random emission (#0)
- Scoring:  Finding the likelihood of a given sequence
  - Calculate likelihood of annotated path and sequence
    - Multiply emission and transition probabilities (#1)
  - Without specifying a path, total probability of generating x
    - Sum probabilities over all paths
    - Forward algorithm (#3)
- Decoding:  Finding the most likely path, given a sequence
  - What is the most likely path generating entire sequence?
    - Viterbi algorithm (#2)
  - What is the most probable state at each time step?
    - Forward + backward algorithms, posterior decoding (#4)
- Learning:  Estimating HMM parameters from training data
  - When state sequence is known
    - Simply compute maximum likelihood A and E (#5a)
  - When state sequence is not known
    - Viterbi training:  Iterative estimation of best path / frequencies (#5b)
    - Baum-Welch:  Iterative estimation over all paths / frequencies (#6)

# The main questions on HMMs

**1. Scoring x, one path** = Joint probability of a sequence and a path, given the model
- – GIVEN a HMM M, a path $\pi$, and a sequence x,
- – FIND Prob[ x, $\pi$ | M ]
- ➔ "Running the model", simply multiply emission and transition probabilities
- ➔ Application: "all promoter" vs. "all backgorund" comparisons

**2. Scoring x, all paths** = total probability of a sequence, summed across all paths
- – GIVEN a HMM M, a sequence x
- – FIND the total probability P[x | M] summed across all paths
- ➔ Forward algorithm, sum score over all paths (same result as backward)

**3. Viterbi decoding** = parsing a sequence into the optimal series of hidden states
- – GIVEN a HMM M, and a sequence x,
- – FIND the sequence $\pi^*$ of states that maximizes P[ x, $\pi$ | M ]
- ➔ Viterbi algorithm, dynamic programming, max score over all paths, trace pointers find path

**4. Posterior decoding** = total prob that emission $x_i$ came from state k, across all paths
- – GIVEN a HMM M, a sequence x
- – FIND the total probability P[$\pi_i$ = k | x, M)
- ➔ Posterior decoding: run forward & backward algorithms to & from state $\pi_i$ =k

**5. Supervised learning** = optimize parameters of a model given training data
- – GIVEN a HMM M, with unspecified transition/emission probs., labeled sequence x,
- – FIND parameters $\theta = (e_i, a_{ij})$ that maximize P[ x | $\theta$ ]
- ➔ Simply count frequency of each emission and transition observed in the training data

**6. Unsupervised learning** = optimize parameters of a model given training data
- – GIVEN a HMM M, with unspecified transition/emission probs., unlabeled sequence x,
- – FIND parameters $\theta = (e_i, a_{ij})$ that maximize P[ x | $\theta$ ]
- ➔ Viterbi training: guess parameters, find optimal Viterbi path (#2), update parameters (#5), iterate
- ➔ Baum-Welch training: guess, sum over all emissions/transitions (#4), update (#5), iterate