# Hidden Markov Models

# Hidden Markov Models

# Outline

- CG-islands
- The "Fair Bet Casino"
- Hidden Markov Model
- Decoding Algorithm
- Forward-Backward Algorithm
- Profile HMMs
- HMM Parameter Estimation
- Viterbi training
- Baum-Welch algorithm

# CG-Islands

- Given 4 nucleotides: probability of occurrence is ~ 1/4.  Thus, probability of occurrence of a dinucleotide is ~ 1/16.

- However, the frequencies of dinucleotides in DNA sequences vary widely.

- In particular, *CG* is typically underrepresented (frequency of *CG* is typically < 1/16)

# Why CG-Islands?

- *CG* is the least frequent dinucleotide because *C* in *CG* is easily *methylated and* has the tendency to mutate into T afterwards

- However, the methylation is suppressed around genes in a genome.  So, *CG* appears at relatively high frequency within these *CG* islands

- So, finding the *CG* islands in a genome is an important problem

# CG Islands and the "Fair Bet Casino"

- The *CG* islands problem can be modeled after a problem named *"The Fair Bet Casino"*

- The game is to flip coins, which results in only two possible outcomes: **H**ead or **T**ail.

- The **F**air coin will give **H**eads and **T**ails with same probability ½.

- The **B**iased coin will give **H**eads with prob. ¾.

# The "Fair Bet Casino" (cont'd)

- Thus, we define the probabilities:
  - $P(H|F) = P(T|F) = \frac{1}{2}$
  - $P(H|B) = \frac{3}{4}$, $P(T|B) = \frac{1}{4}$
  - The crooked dealer changes between Fair and Biased coins with probability  10%

# The Fair Bet Casino Problem

- **Input:** A sequence $x = x_1 x_2 x_3 \ldots x_n$ of coin tosses made by two possible coins (**F** or **B**).

- **Output:** A sequence $\pi = \pi_1\ \pi_2\ \pi_3 \ldots \pi_n$, with each $\pi_i$ being either $F$ or $B$ indicating that $x_i$ is the result of tossing the Fair or Biased coin respectively.

# Problem…

*Fair Bet Casino Problem*

Any observed outcome of coin tosses could have been generated by any sequence of states!

# Problem…

**Fair Bet Casino Problem**

Any observed outcome of coin tosses could have been generated by any sequence of states!

Need to incorporate a way to grade different sequences differently.

# Problem…

*Fair Bet Casino Problem*

Any observed outcome of coin tosses could have been generated by any sequence of states!

Need to incorporate a way to grade different sequences differently.

*Decoding Problem*

# Hidden Markov Model (HMM)

- Can be viewed as an abstract machine with *k hidden* states that emits symbols from an alphabet Σ.

- Each state has its own probability distribution, and the machine switches between states according to this probability distribution.

- While in a certain state, the machine makes 2 decisions:

  - What state should I move to next?

  - What symbol - from the alphabet Σ - should I emit?

# Why "Hidden"?

- Observers can see the emitted symbols of an HMM but have *no ability to know which state the HMM is currently in*.

- Thus, the goal is to infer the most likely hidden states of an HMM based on the given sequence of emitted symbols.

# HMM Parameters

Σ: set of emission characters.

  Ex.: Σ = {H, T} for coin tossing

   Σ = {1, 2, 3, 4, 5, 6} for dice tossing

Q: set of hidden states, each emitting symbols from Σ.

   Q={F,B} for coin tossing

# HMM Parameters (cont'd)

A = ($a_{kl}$): a |Q| x |Q| matrix of probability of changing from state *k* to state *l*.

$$a_{FF} = 0.9 \qquad a_{FB} = 0.1$$

$$a_{BF} = 0.1 \qquad a_{BB} = 0.9$$

E = ($e_k(b)$): a |Q| x |Σ| matrix of probability of emitting symbol *b* while being in state *k*.

$$e_F(0) = \tfrac{1}{2} \qquad e_F(1) = \tfrac{1}{2}$$

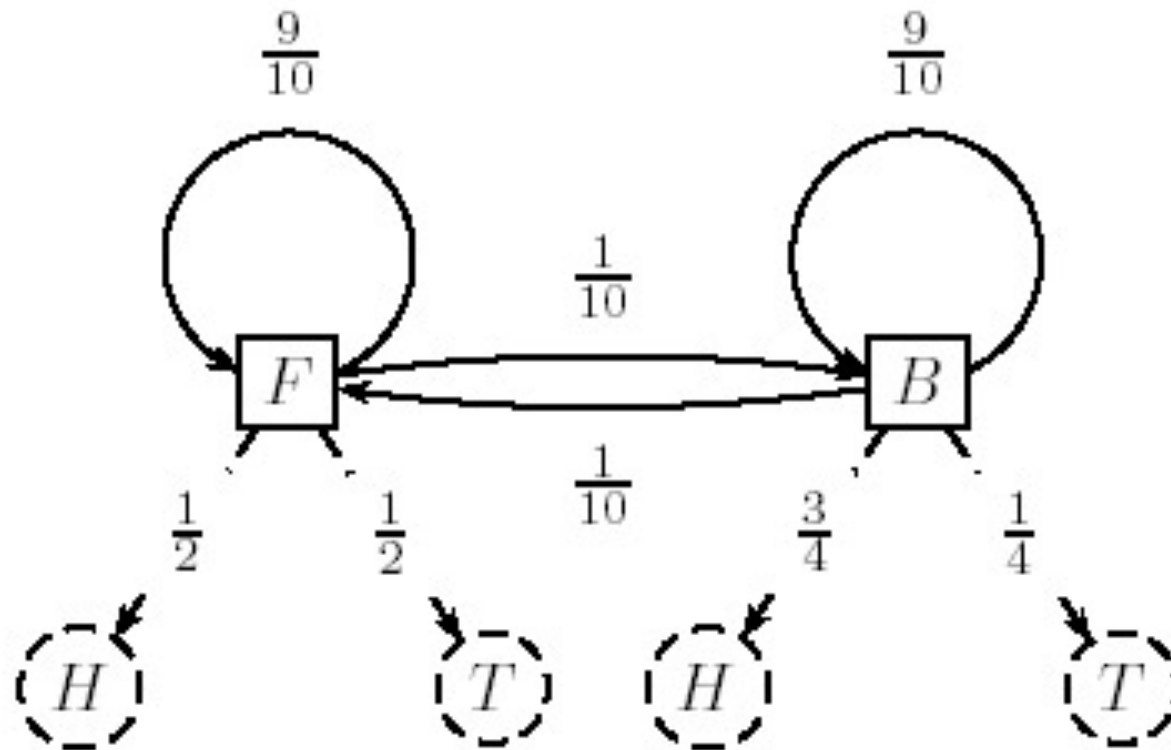$$e_B(0) = \tfrac{1}{4} \qquad e_B(1) = \tfrac{3}{4}$$

# HMM for Fair Bet Casino

- The *Fair Bet Casino* in *HMM* terms:

  Σ = {0, 1} (0 for **T**ails and 1 **H**eads)

  Q = {*F,B*} – *F* for Fair & *B* for Biased coin.

- Transition Probabilities *A* *** Emission Probabilities *E*

|  | Fair | Biased |
|---|---|---|
| Fair | $a_{FF} = 0.9$ | $a_{FB} = 0.1$ |
| Biased | $a_{BF} = 0.1$ | $a_{BB} = 0.9$ |

|  | Tails(0) | Heads(1) |
|---|---|---|
| Fair | $e_F(0) = ½$ | $e_F(1) = ½$ |
| Biased | $e_B(0) = ¼$ | $e_B(1) = ¾$ |

# HMM for Fair Bet Casino (cont'd)



**HMM model for the *Fair Bet Casino* Problem**

# Hidden Paths

- A *path* $\pi = \pi_1 \ldots \pi_n$ in the HMM is defined as a sequence of states.

- Consider path $\pi$ = FFFBBBBBFFF and sequence $x$ = 01011101001

Probability that $x_i$ was emitted from state $\pi_i$

| x | | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\pi$ | = | F | F | F | B | B | B | B | B | F | F | F |
| $P(x_i\|\pi_i)$ | | ½ | ½ | ½ | ¾ | ¾ | ¾ | ¼ | ¾ | ½ | ½ | ½ |
| $P(\pi_{i-1} \to \pi_i)$ | | ½ | $^9/_{10}$ | $^9/_{10}$ | $^1/_{10}$ | $^9/_{10}$ | $^9/_{10}$ | $^9/_{10}$ | $^9/_{10}$ | $^1/_{10}$ | $^9/_{10}$ | $^9/_{10}$ |

Transition probability from state $\pi_{i-1}$ to state $\pi_i$

# P(x|π) Calculation

- P(*x*|π): Probability that sequence *x* was generated by the path *π:*

$$P(x|\pi) = P(\pi_0 \to \pi_1) \cdot \prod_{i=1}^{n} P(x_i | \pi_i) \cdot P(\pi_i \to \pi_{i+1})$$

$$= a_{\pi_0, \pi_1} \cdot \prod e_{\pi_i}(x_i) \cdot a_{\pi_i, \pi_{i+1}}$$

# P(x|π) Calculation

- P($x$|π): Probability that sequence $x$ was generated by the path π:

$$P(x|\pi) = P(\pi_0 \to \pi_1) \cdot \prod_{i=1}^{n} P(x_i | \pi_i) \cdot P(\pi_i \to \pi_{i+1})$$

$$= a_{\pi_0, \pi_1} \cdot \prod e_{\pi_i}(x_i) \cdot a_{\pi_i, \pi_{i+1}}$$

$$= \prod e_{\pi_{i+1}}(x_{i+1}) \cdot a_{\pi_i, \pi_{i+1}}$$
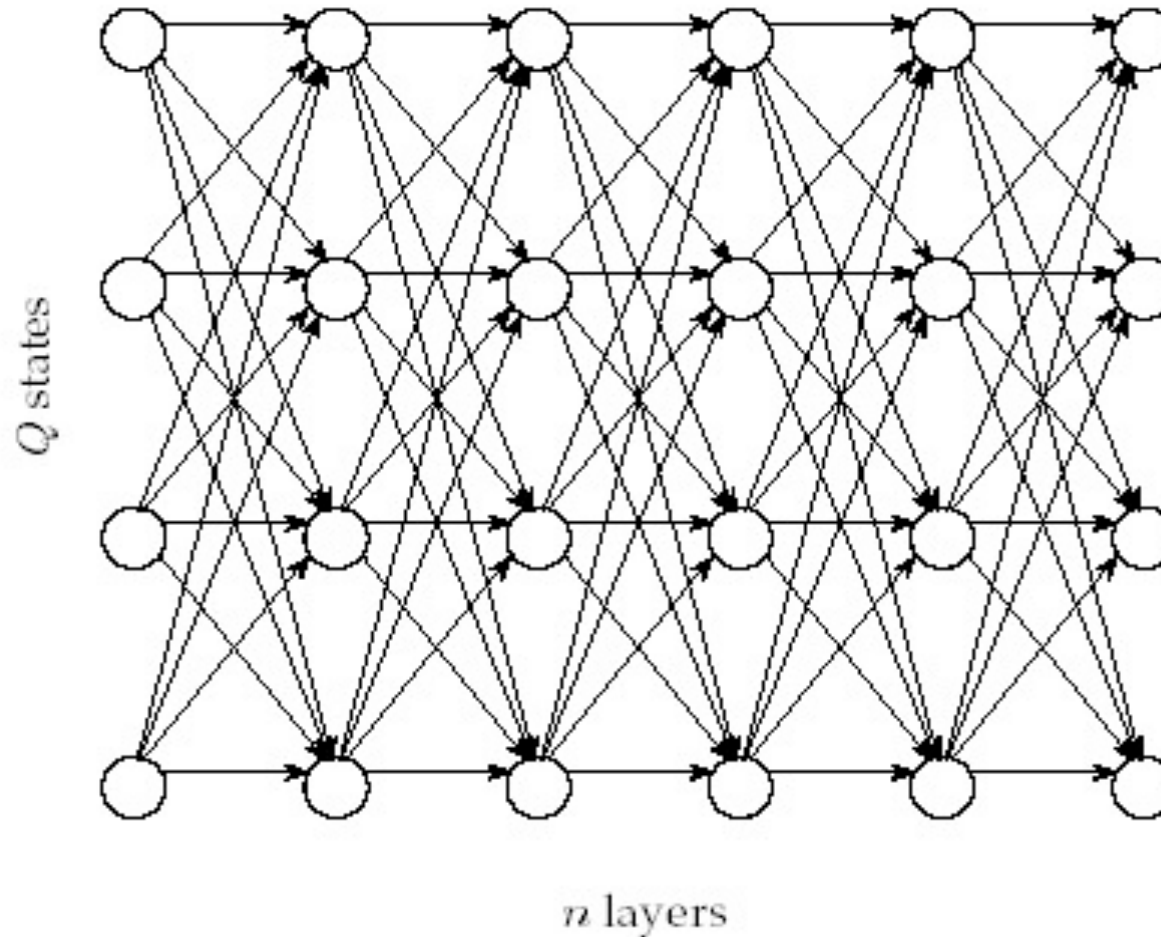
if we count from $i=0$ instead of $i=1$

# Decoding Problem

- **Goal:** Find an optimal hidden path of states given observations.

- **Input:** Sequence of observations $x = x_1 \ldots x_n$ generated by an HMM $M(\Sigma, Q, A, E)$

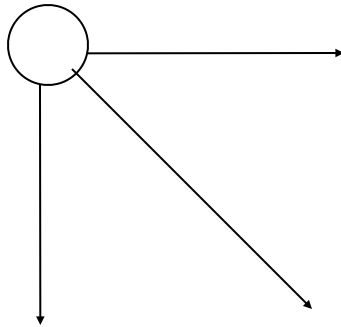- **Output:** A path that maximizes $P(x|\pi)$ over all possible paths $\pi$.

# Building Manhattan for Decoding Problem

- Andrew Viterbi used the Manhattan grid model to solve the *Decoding Problem*.

- Every choice of $\pi = \pi_1 \ldots \pi_n$ corresponds to a path in the graph.

- The only valid direction in the graph is *eastward.*
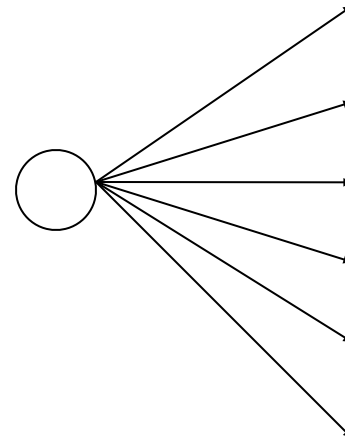
- This graph has $|Q|^2(n-1)$ edges.

# Edit Graph for Decoding Problem

# Decoding Problem vs. Alignment Problem

Valid directions in the *alignment problem.*
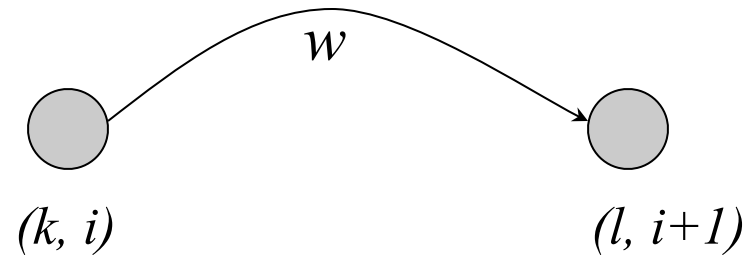
Valid directions in the *decoding problem.*

# Decoding Problem as Finding a

- The *Decoding Problem* is reduced to finding a longest path in the *directed acyclic graph (DAG)* above.


- **<u>Notes:</u>** the length of the path is defined as the *product* of its edges' weights, not the *sum.*

# Decoding Problem (cont'd)

- Every path in the graph has the probability $P(x|\pi)$.

- The Viterbi algorithm finds the path that maximizes $P(x|\pi)$ among all possible paths.

- The Viterbi algorithm runs in $O(n|Q|^2)$ time.

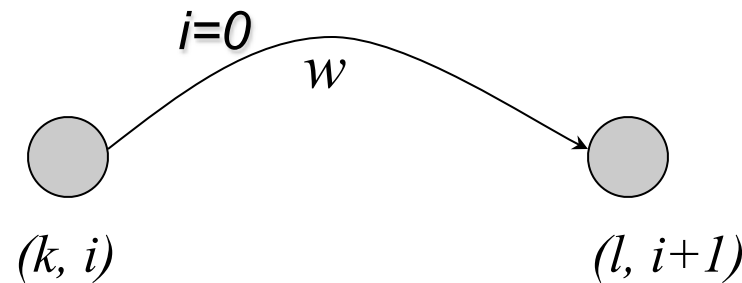# Decoding Problem: weights of edges

$w$

(k, i)                    (l, i+1)

## The weight *w* is given by:

## ???

# Decoding Problem: weights of edges

$$P(x|\pi) = \prod_{i=0}^{n} e_{\pi_{i+1}}(x_{i+1}) \cdot a_{\pi_i, \pi_{i+1}}$$
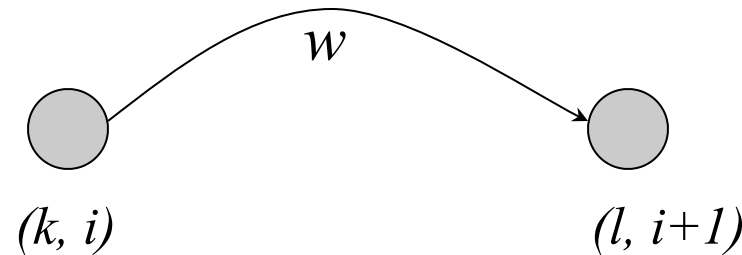
$w$

$(k, i)$    $(l, i+1)$

The weight **w** is given by:

**??**

# Decoding Problem: weights of edges

$i$-th term = $e_{\pi_{i+1}} (x_{i+1}) \cdot a_{\pi_i, \pi_{i+1}}$

$w$

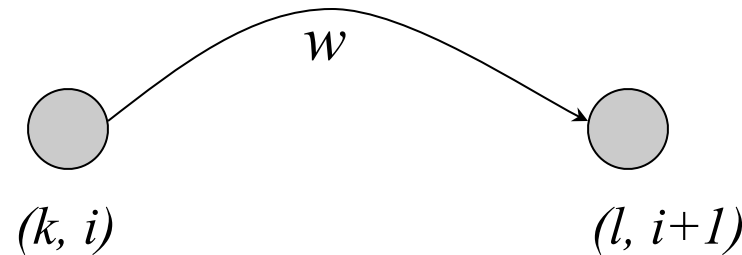$(k, i)$                    $(l, i+1)$

## The weight *w* is given by:

*?*

# Decoding Problem: weights of edges

$i$-th term $= e_{\pi_i}(x_i) \cdot a_{\pi_i, \pi_{i+1}} = e_l(x_{i+1}) \cdot a_{kl}$  for  $\pi_i = k,\ \pi_{i+1} = l$

$$w$$

$(k,\ i)$ $\qquad\qquad\qquad\qquad$ $(l,\ i{+}1)$

The weight  $w = e_l(x_{i+1}) \cdot a_{kl}$

# Decoding Problem and Dynamic Programming

$$s_{l,i+1} = \max_{k \in Q} \{s_{k,i} \cdot \text{weight of edge between } (k,i) \text{ and } (l,i+1)\} =$$

$$\max_{k \in Q} \{s_{k,i} \cdot a_{kl} \cdot e_l(x_{i+1})\} =$$

$$e_l(x_{i+1}) \cdot \max_{k \in Q} \{s_{k,i} \cdot a_{kl}\}$$