# INTRODUCTION TO BIOINFORMATICS

## PRACTICE — 9th week

Hidden Markov Models

# PROBABILITY THEORY

- CONDITIONAL PROBABILITY    $P(A \mid B)$

- JOINT PROBABILITY    $P(A, B)$

$$P(X) = \sum_{Y} P(X, Y) = \sum_{Y} P(X \mid Y)P(Y)$$
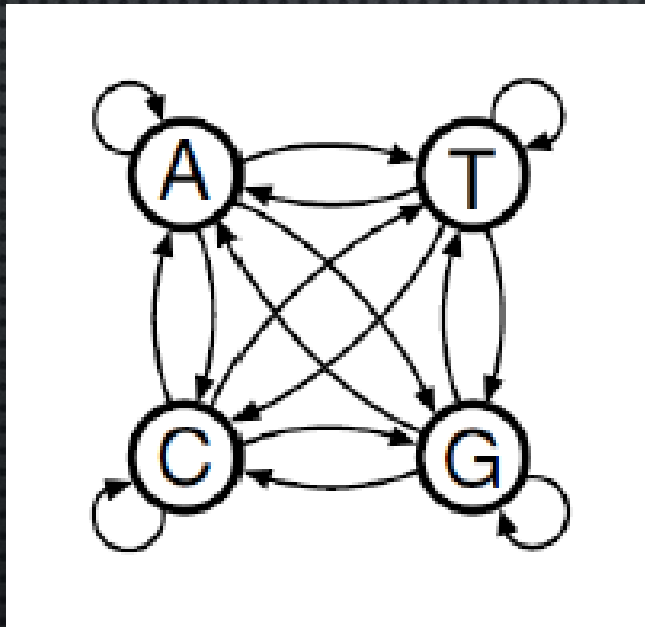
- BAYES THEOREM    $P(X \mid Y) = \dfrac{P(Y \mid X)P(X)}{P(Y)}$

# MARKOV CHAINS

Markov chain for generating random nucleotide sequences:



- Each state represents a symbol
- The sequence is the ('time') series of states

# MARKOV CHAINS

- MARKOV CHAINS CAN BE DESCRIBED BY STATE TRANSITIONS

- GIVEN A FINITE SET OF STATES: $x = \{x_1\ x_2\ x_3 \ldots x_L\}$

- AT TIME POINT *T+1* THE PROCESS STAYS IN THE SAME STATE AS IN TIME STEP *T* OR MOVES TO ANOTHER STATE

- **TRANSITION PROBABILITY**: $P_{IJ}$ (THE PROBABILITY OF MOVING FROM THE *ITH* STATE TO THE *JTH*)

  → WE CAN BUILD A **STATE TRANSITION MATRIX** FROM THEM

$$P(x) = P(x_L, x_{L-1}, \ldots, x_1) = P(x_L \mid x_{L-1}, \ldots, x_1) \cdot P(x_{L-1} \mid x_{L-2}, \ldots, x_1) \cdot \ldots \cdot P(x_1)$$
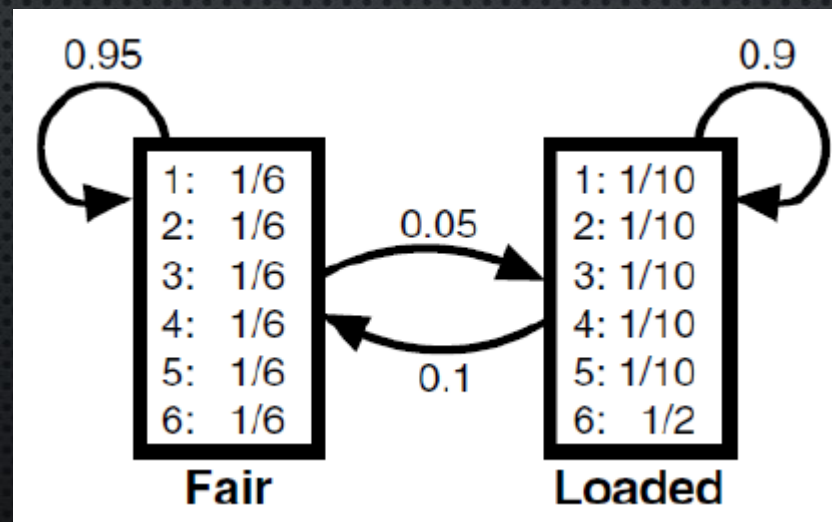
- **MARKOV PROPERTY**: NEXT STATE DEPENDS ONLY ON THE ACTUAL ONE:

$$P(x) = P(x_L \mid x_{L-1}) \cdot P(x_{L-1} \mid .x_{L-2}) \cdot \ldots \cdot P(x_2 \mid x_1) \cdot P(x_1)$$

# HMM – UNFAIR CASINO

- AN UNFAIR CASINO USES 2 KINDS OF DICES: 99% OF THE DICES IS NORMAL, 1% IS LOADED. ROLLING 6 WITH THE LOADED DICE HAS A PROBABILITY OF 0.5.

- HMM MODEL OF THE UNFAIR CASINO:

# HMM - UNFAIR CASINO

- GENERATING SEQUENCES OF SYMBOLS (I.E. NUMBERS)
  - THE SAME SEQUENCE COULD BE GENERATED BY EITHER A LOADED OR NORMAL DICE

- ONE CAN OBTAIN THE SEQUENCE ITSELF, BUT NOT THE STATES!

- HOW DO YOU FIGURE OUT WHAT WERE THE STATES (DICES), THAT ARE POSSIBLY GENERATED THE UNDERLYING SEQUENCES
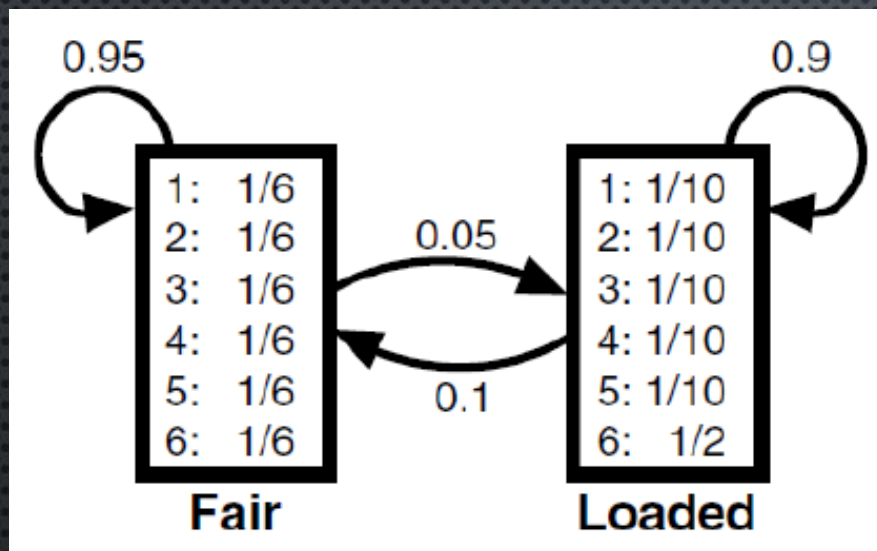
WHICH ONE WAS MORE LIKELY GENERATED BY A LOADED DICE?

SEQ1: 45**663**2**66**164**6666**411**6566**5**6**312**6**33**66**233**663**621353

SEQ2: 2514542426**25161523445655543643453256**34324**6**14

# HMM - COMPONENTS



- **HIDDEN STATE:** (LOADED/NORMAL DICE) AND (OBSERVABLE) **EMITTED SYMBOLS** ({ACTG}/{1,2,3,4,5,6})

- **EMISSION MATRIX:** PROBABILITY OF OBSERVABLE SYMBOLS GIVEN THAT THE HIDDEN MODEL IS IN A CERTAIN HIDDEN STATE

- **INITIAL DISTRIBUTION:** THE PROBABILITY OF THE MODEL BEING IN A CERTAIN HIDDEN STATE AT TIME 0

- **STATE TRANSITION MATRIX:** PROBABILITY OF MOVING FROM ONE HIDDEN STATE TO ANOTHER HIDDEN STATE

# DECODING - VITERBI ALGORITHM

QUESTION: WHAT IS THE SEQUENCE OF THE STATES THAT GENERATED THE FOLLOWING SEQUENCE?

- FIGURING OUT THAT AT EACH POSITION THE LOADED OR THE FAIR DICE HAD BEEN USED

45**66**32**66**164**6666**411**6**5**66**5**6**312**6**33**66**233**66**3**6**213532514542**426**251**6**15234456555436434532**56**34324**6**14

SOLUTION:
ASSUMING THAT WE KNOW THAT UNTIL THE POSITION 8 THE LOADED DICE WAS USED, HOW WOULD YOU DECIDE WHICH ONE IS THE NEXT STATE? SUGGESTIONS?
45**66**32**66**1

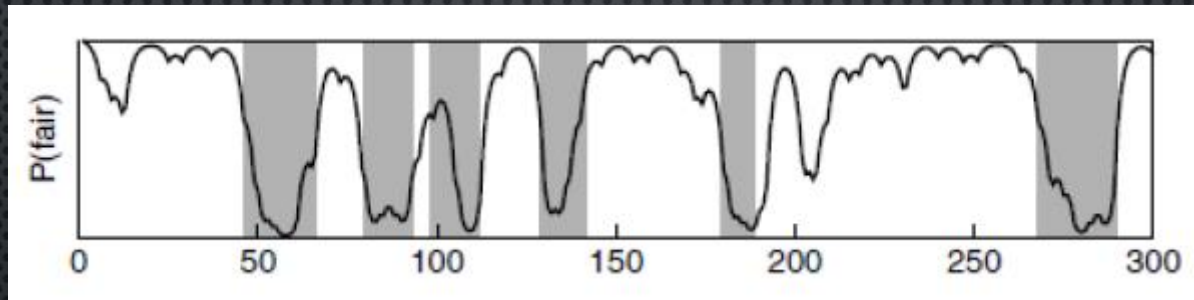~~LLLLLLLL~~?

# DECODING - VITERBI ALGORITHM

- VITERBI: FINDING THE MOST PROBABLE PATH  (PATH IS DENOTED AS $\pi$)

- MOST PROBABLE PATH:  $$\pi^* = \arg\max_{\pi} P(x, \pi)$$

- LET $v_K(I)$ BE THE PROBABILITY THAT THE MOST PROBABLE PATH ENDS AT $K$. STATE WHERE IT EMITS SYMBOL $X_I$

$$v_l(i+1) = e_l(x_{i+1}) \max_k (v_k(i) a_{kl})$$

WHERE $e_L(X_{I+1})$ IS THE PROBABILITY OF EMITTING $X_{I+1}$ AT THE $L$ STATE; $a_{KL}$ IS THE TRANSITION PROBABILITY FROM $K$ TO $L$

# Viterbi - Unfair casino

- WITH DYNAMIC PROGRAMMING IT IS POSSIBLE TO FIND THE MOST PROBABLE PATH



- FINDING GENES IN A NUCLEOTIDE SEQUENCE

# HMM - Training

- Until now we assumed that we know the model exactly. (model: transition and emission probabilities)

- In real life it is not the case

- The model parameters (emission and transition probabilities) must be estimated somehow (~training the model)

- Must have a training set, i.e. annotated genomes, annotated sequences

- Estimation:

  - Expectation-maximization algorithm

  - Baum Welch training

# PROFILE HMM

- GIVEN A PROTEIN FAMILY WE WANT TO FIND OTHER MEMBERS OF THIS FAMILY IN A GIVEN DATABASE

- WE CAN MAKE A PROFILE THAT REPRESENTS THE PROTEIN FAMILY AND WE COMPARE THE DATABASE ELEMENTS TO THIS PROFILE

- AN HMM CAN DESCRIBE SUCH A PROFILE

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **A** | .72 | .14 | 0 | 0 | .72 | .72 | 0 | 0 |
| **T** | .14 | .72 | 0 | 0 | 0 | .14 | .14 | .86 |
| **G** | .14 | .14 | .86 | .44 | 0 | .14 | 0 | 0 |
| **C** | 0 | 0 | .14 | .56 | .28 | 0 | .86 | .14 |

# PROFILE HMM

- States:

  - M: match state

  - I: insertion state

  - D: deletion state (silent state – no emission)

- n: length of sequence (number of matches)