**Slide 1**

Spring 2024

# BIOS 477/877

*Bioinformatics and Molecular Evolution*

## Lecture 18

BIOS477/877 L18 - 1

**1**

---

**Slide 2**

## TODAY'S TOPICS

➢ **PSI-BLAST**

➢ **Hidden Markov model and profile HMM**

➢ **Protein family/domain databases**
(InterPro, HMMER, etc.)

BIOS477/877 L18 - 2

**2**

---

**Slide 3**

## PSI-BLAST

1st iteration = a regular BLASTP

BIOS477/877 L18 - 3

**3**

---

**Slide 4**

## PSI-BLAST

➢ Two different E-value thresholds

Threshold E-value for reporting matches

PSI-BLAST threshold:
E-value threshold for sequences to be included to create the PSSM for the next iteration

BIOS477/877 L18 - 4

**4**

---

**Slide 5**

## PSI-BLAST

PSI-BLAST threshold (0.005)

Sequences with E-value BETTER than threshold

Used to construct PSSM and used for the next search

E-value < 0.005 (PSI-BLAST threshold)

Sequences with E-value WORSE than threshold

0.05 > E-value > 0.005

E-value < 0.05 (reporting threshold)

BIOS477/877 L18 - 5

**5**

---

**Slide 6**

## PSI-BLAST

PSSM is constructed from the chosen sequences

Run PSI-BLAST iteration 2 with max number of sequences — Start the 2nd iteration of search

Sequences with E-value WORSE than threshold

More sequences can be included, but be very careful!

• If unrelated sequences are included, generated PSSM loses the specificity.
• Errors can be amplified quickly with iterations.
➔ Profile corruption problem

BIOS477/877 L18 - 6

**6**

---

1

## Slide 7

**PSI-BLAST**

Sequences producing significant alignments    Download ⌄    Select columns ⌄    Show [500 ▼] ❓

500 sequences selected    ☐ sequences newly added this iteration ❓    GenPept    Graphics    Distance tree of results    Multiple alignment    MSA Viewer

Sequences with E-value BETTER than threshold

☑ select all  500 sequences selected    Skip to the first new sequence    PSI-BLAST iteration 2

| Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession | Select for PSI blast | Used to build PSSM | Newly added |
|---|---|---|---|---|---|---|---|---|---|---|---|
| uncharacterized protein LOC6587669 [Drosophila persimilis] | Drosophila pe... | 422 | 422 | 100% | 5e-144 | 75.00% | 420 | XP_026850866.1 | ☑ | ✓ | |
| uncharacterized protein LOC108154648 [Drosophila miranda] | Drosophila mir... | 417 | 417 | 100% | 6e-144 | 75.00% | 304 | XP_017140468.1 | ☑ | ✓ | |
| uncharacterized protein LOC4801909 [Drosophila pseudoobscura] | Drosophila ps... | 421 | 421 | 100% | 9e-144 | 74.68% | 420 | XP_001358928.5 | ☑ | ✓ | |
| uncharacterized protein LOC6726927 [Drosophila simulans] | Drosophila si... | 408 | 408 | 94% | 4e-140 | 97.26% | 308 | XP_002102329.1 | ☑ | ✓ | |
| uncharacterized protein LOC122621875 [Drosophila teissieri] | Drosophila tei... | 407 | 407 | 94% | 9e-140 | 95.21% | 308 | XP_043655819.1 | ☑ | ✓ | |
| osiris 1 [Drosophila melanogaster] | Drosophila me... | 406 | 406 | 94% | 3e-139 | 100.00% | 308 | NP_649618.1 | ☑ | ✓ | |

• • •

| hypothetical protein JYU34_019236 [Plutella xylostella] | Plutella xylost... | 147 | 147 | 84% | 1e-38 | 23.13% | 266 | KAG7297281.1 | ☑ | | ⊘ |
| hypothetical protein DOY81_001298 [Sarcophaga bullata] | Sarcophaga b... | 141 | 141 | 40% | 2e-37 | 67.19% | 148 | TMW53623.1 | ☑ | ✓ | |
| uncharacterized protein LOC124358906 isoform X2 [Homalodisca vitripennis] | Homalodisca... | 144 | 144 | 66% | 3e-37 | 22.42% | 282 | XP_046667114.1 | ☑ | ✓ | |
| uncharacterized protein LOC6573135 isoform X2 [Drosophila mojavensis] | Drosophila mo... | 139 | 139 | 45% | 5e-37 | 66.90% | 136 | XP_032589606.1 | ☑ | ✓ | |
| hypothetical protein B566_EDAN009200 [Ephemera danica] | Ephemera da... | 142 | 142 | 84% | 2e-36 | 25.09% | 279 | KAF4528605.1 | ☑ | | ⊘ |
| hypothetical protein LSTR_LSTR004646 [Laodelphax striatellus] | Laodelphax st... | 139 | 139 | 77% | 4e-35 | 21.43% | 278 | RZF36958.1 | ☑ | ✓ | |
| uncharacterized protein LOC124358906 isoform X1 [Homalodisca vitripennis] | Homalodisca... | 137 | 137 | 80% | 2e-34 | 22.99% | 280 | XP_046667112.1 | ☑ | ✓ | |
| hypothetical protein RR48_04877 [Papilio machaon] | Papilio machaon | 132 | 132 | 64% | 1e-33 | 27.59% | 196 | KPJ20279.1 | ☑ | | |

BIOS477/877 L18 - 7

**7**

---

## Slide 8

**PSI-BLAST**

**First iteration (BLASTP search)**    [0.001]

| | | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession | |
|---|---|---|---|---|---|---|---|---|---|
| ☑ uncharacterized protein LOC123680163 [Harmonia axyridis] | Harmonia axyridis | 52.4 | 52.4 | 49% | 0.001 | 25.47% | 276 | XP_045473886.1 | ☑ |
| ☑ uncharacterized protein LOC105385699 isoform X1 [Plutella xylostella] | Plutella xylostella | 52.0 | 52.0 | 77% | 0.001 | 24.60% | 265 | XP_037967634.1 | ☑ |
| ☑ uncharacterized protein LOC113498740 isoform X1 [Trichoplusia ni] | Trichoplusia ni | 52.4 | 0.001 | | 0.002 | 25.65% | 299 | XP_026734681.1 | ☑ |
| ☑ uncharacterized protein LOC123677560 [Maniola jurtina] | Maniola jurtina | | | | 0.002 | 25.15% | 268 | XP_045780333.1 | ☑ |
| ☑ hypothetical protein LSTR_LSTR004646 [Laodelphax striatellus] | Laodelphax striatellus | 51.6 | 51.6 | 76% | 0.002 | 22.64% | 278 | RZF36958.1 | ☑ |
| ☑ uncharacterized protein LOC113498740 isoform X2 [Trichoplusia ni] | Trichoplusia ni | 51.6 | 51.6 | 81% | 0.003 | 24.83% | 196 | XP_026734698.1 | ☑ |
| ☑ uncharacterized protein LOC105385699 isoform X1 [Plutella xylostella] | Plutella xylostella | 50.4 | 50.4 | 77% | 0.004 | 24.90% | 269 | XP_037967632.1 | ☑ |

[0.004]

**Second iteration**    [6e-40]

| | | | | | E value | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ uncharacterized protein LOC105385699 isoform X2 [Plutella xylostella] | Plutella xylost... | 151 | 151 | 84% | 6e-40 | 23.22% | 265 | XP_037967634.1 | | ✓ | |
| ☑ uncharacterized protein LOC105385699 isoform X1 [Plutella xylostella] | Plutella xylost... | 149 | 149 | 84% | 4e-39 | 23.51% | 269 | XP_037967632.1 | | ✓ | |
| ☑ hypothetical protein JYU34_019236 [Plutella xylostella] | Plutella xylost... | 1 | | | 1e-38 | 23.13% | 266 | KAG7297281.1 | | | ⊘ | [4e-39] |
| ☑ hypothetical protein DOY81_001298 [Sarcophaga bullata] | Sarcophaga b... | 1 | | | 2e-37 | 67.19% | 148 | TMW53623.1 | | ✓ | |

> **- 1st iteration has the real E-values for the query.**
> **- After the 2nd iteration, E-values are for the PSSM.**

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ uncharacterized protein LOC124358906 isoform X2 [Homalodisca vitripennis] | Homalodisca... | 137 | 137 | 80% | 2e-34 | 22.99% | 280 | XP_046667112.1 | | ⊘ | |
| ☑ hypothetical protein RR48_04877 [Papilio machaon] | Papilio machaon | 132 | 132 | 64% | 1e-33 | 27.59% | 196 | KPJ20279.1 | | ✓ | |
| ☑ hypothetical protein RR46_04253 [Papilio xuthus] | Papilio xuthus | 132 | 132 | 64% | 2e-33 | 27.09% | 196 | KPJ03641.1 | | ✓ | |
| ☑ uncharacterized protein LOC123680163 [Harmonia axyridis] | Harmonia axy... | 130 | 130 | 62% | 9e-32 | 24.51% | 276 | XP_045473886.1 | | ✓ | |

[9e-32]

BIOS477/877 L18 - 8

**8**

---

## Slide 9

**PSI-BLAST**

**First iteration (BLASTP search)**

| Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession | Select for PSI blast | Used to build PSSM | Newly added |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ osiris 1 [Drosophila melanogaster] | Drosophila melanogaster | 638 | 638 | 100% | 0.0 | 100.00% | 308 | NP_649618.1 | ☑ | | |
| ☑ uncharacterized protein LOC6726927 [Drosophila simulans] | Drosophila simulans | 624 | 624 | 100% | 0.0 | 97.40% | 308 | XP_002102329.1 | ☑ | | |
| ☑ uncharacterized protein LOC122621875 [Drosophila teissieri] | Drosophila teissieri | 530 | 530 | 94% | 0.0 | | 308 | XP_043655819.1 | ☑ | | |
| ☑ uncharacterized protein LOC6552948 [Drosophila erecta] | Drosophila erecta | 523 | 523 | 94% | 0.0 | 93.40% | 308 | XP_001979036.1 | ☑ | | |
| ☑ uncharacterized protein LOC6614069 [Drosophila sechellia] | Drosophila sechellia | 520 | 520 | 94% | 0.0 | 95.21% | 308 | XP_002038523.1 | ☑ | | |
| ☑ uncharacterized protein LOC120454180 [Drosophila santomea] | Drosophila santomea | 513 | 513 | 94% | 0.0 | 95.55% | 308 | XP_039495195.1 | ☑ | | |

> **E-values can become lower or higher with more iterations depending on the sequences included to build PSSM!!**

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ uncharacterized protein LOC4801909 [Drosophila pseudoobscura] | Drosophila ps... | 421 | 421 | 100% | 9e-144 | 74.68% | 420 | XP_001358928.5 | ☑ | ✓ | |
| ☑ uncharacterized protein LOC6726927 [Drosophila simulans] | Drosophila si... | 408 | 408 | 94% | 4e-140 | 97.26% | 308 | XP_002102329.1 | ☑ | ✓ | |
| ☑ uncharacterized protein LOC122621875 [Drosophila teissieri] | Drosophila tei... | 407 | 407 | 94% | 9e-140 | 95.21% | 308 | XP_043655819.1 | ☑ | ✓ | |
| ☑ osiris 1 [Drosophila melanogaster] | Drosophila me... | 406 | 406 | 94% | 3e-139 | 100.00% | 308 | NP_649618.1 | ☑ | ✓ | |
| ☑ uncharacterized protein LOC6552948 [Drosophila erecta] | Drosophila er... | 404 | 404 | 94% | 1e-138 | 93.15% | 308 | XP_001979036.1 | ☑ | ✓ | |

[~ 1e-139]

BIOS477/877 L18 - 9

**9**

---

## Slide 10

**Profile hidden Markov models**

➢ **Probabilistic models of multiple alignments**
  (Krogh *et al.*, 1994; Krogh 1998)

  ➜ **Closely related to standard profiles (PSSMs) introduced by Gribskov**

  ➜ **Used in *e.g.*, PFAM, SMART, Superfamily, Panther (databases of multiple alignments and profile HMMs)**

> **What is a hidden Markov model?**
> **What is hidden?**

Eddy (2004) What is a hidden Markov model? ← DNA regulatory element example
*Nature Biotechnology* 22: 1315-1316.

BIOS477/877 L18 - 10

**10**

---

## Slide 11

**Markov models**

➢ **Markov chain**: a stochastic process from one state to the next that exhibits the **Markov property**
  (*e.g.*, States for DNA: A, T, G, C)

$X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4 \rightarrow X_5$

• **Markov property**: The next state of the system depends only on the present state of the system **[No memory]**

BIOS477/877 L18 - 11

**11**

---

## Slide 12

**Markov models**

➢ **Markov chain**: a stochastic process from one state to the next that exhibits the **Markov property**
  (*e.g.*, States for DNA: A, T, G, C)

[P: Transition probability]

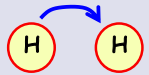$X_1 \xrightarrow{P} X_2 \xrightarrow{P} X_3 \xrightarrow{P} X_4 \xrightarrow{P} X_5$

• **Markov property**: The next state of the system depends only on the present state of the system **[No memory]**

• Transition probabilities are independent of time
  **[Time homogeneous]**
  - Markov chain if the state spaces is discrete
  - Markov process if the state space is continuous

BIOS477/877 L18 - 12

**12**

2

## Markov models

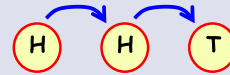> **Markov chain**
  Coin tossing example [2 states: Head and Tail]

  (H) → (H)

  - P(Head) = 0.5, P(Tail) = 0.5
  - P(HH) = P(H) x P(H)

**13**

---

## Markov models

> **Markov chain**
  Coin tossing example [2 states: Head and Tail]

  (H) → (H) → (T)

  - P(Head) = 0.5, P(Tail) = 0.5
  - P(HH) = P(H) x P(H)
  - P(HHT) = P(HH) x P(T)

**14**

---

## Markov models

> **Markov chain**
  Coin tossing example [2 states: Head and Tail]

  (H) → (H) → (T) → (H)

  - P(Head) = 0.5, P(Tail) = 0.5
  - P(HH) = P(H) x P(H)
  - P(HHT) = P(HH) x P(T)
  - P(HHTH) = P(HHT) x P(H)
  - P(HHTHT) = P(HHTH) x P(T)
  ...

**15**

---

## Markov models

> **Markov chain**
  Coin tossing example [2 states: Head and Tail]

  (H) → (H) → (T) → (H) → (T)

  - P(Head) = 0.5, P(Tail) = 0.5
  - P(HH) = P(H) x P(H)
  - P(HHT) = P(HH) x P(T) = P(H) x P(H) x P(T)
  - P(HHTH) = P(HHT) x P(H) = P(H) x P(H) x P(T) x P(H)
  - P(HHTHT) = P(HHTH) x P(T)
    = P(H) x P(H) x P(T) x P(H) x P(T)

**16**

---

## Markov models

> **Markov chain**
  Coin tossing example

  [H→T]
  0.5
  [H→H] 0.5  (H) ⇄ (T)  0.5 [T→T]
  0.5
  [T→H]

  - **States:** Head or Tail
  - **Transition probability** ⇒

  At time t+1
  |   | H | T |
  |---|---|---|
  | H | 0.5 | 0.5 |
  | T | 0.5 | 0.5 |

  Row sum = 1
  At time t

**17**

---

## Markov models

> **Markov chain**
  Coin tossing example: loaded coin (more Tail than Head)

  a=0.6
  1-a  (H) ⇄ (T)  1-b
  b=0.4

  - **States:** Head or Tail
  - **Transition probability** ⇒

  At time t+1
  |   | H | T |
  |---|---|---|
  | H | 0.4 | 0.6 |
  | T | 0.4 | 0.6 |

  Row sum = 1
  At time t

**18**

3

**19**

# Markov models

➤ **Markov chain**
**Coin tossing example**

$[H \to T]$ $a$

$[H \to H]$ $1-a$    H    T    $1-b$ $[T \to T]$

$b$
$[T \to H]$

- **States:** Head or Tail
- **Transition probability**

At time t+1

$$\begin{array}{c} & H & T \\ H & 1-a & a \\ T & b & 1-b \end{array}$$     Row sum = 1

At time t

BIOS477/877 L18 - 19

---

**20**

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**

0.6
0.4  H  T  0.6
0.4

$$P_L = \begin{array}{c} & H & T \\ H & 0.4 & 0.6 \\ T & 0.4 & 0.6 \end{array}$$

Transition probability matrix

Start

H    0.4 → H    0.4x0.4=0.16 Prob(HH) ⟸
     0.6 → T    0.4x0.6=0.24 Prob(HT)
0.4
0.6
T    0.6 → H    0.6x0.4=0.24 Prob(TH) ⟸
     0.6 → T    0.6x0.6=0.36 Prob(TT)

[1st Toss] [2nd Toss]

Prob( H @2nd Toss ) = 0.16+0.24 = 0.4

BIOS477/877 L18 - 20

---

**21**

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**

0.6
0.4  H  T  0.6
0.4

$$P_L = \begin{array}{c} & H & T \\ H & 0.4 & 0.6 \\ T & 0.4 & 0.6 \end{array}$$

Transition probability matrix

Start

H    0.4 → H    0.4x0.4=0.16 Prob(HH)
     0.6 → T    0.4x0.6=0.24 Prob(HT) ⟸
0.4
0.6
T    0.4 → H    0.6x0.4=0.24 Prob(TH)
     0.6 → T    0.6x0.6=0.36 Prob(TT) ⟸

[1st Toss] [2nd Toss]

Prob( T @2nd Toss ) = 0.24+0.36 = 0.6

BIOS477/877 L18 - 21

---

**22**

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**

0.6
0.4  H  T  0.6
0.4

$$P_L = \begin{array}{c} & H & T \\ H & 0.4 & 0.6 \\ T & 0.4 & 0.6 \end{array}$$

Transition probability matrix

Start

H    0.4 → H    0.4x0.4=0.16 Prob(HH) ⟸
     0.6 → T    0.4x0.6=0.24 Prob(HT) ⟸
0.4
0.6
T    0.4 → H    0.6x0.4=0.24 Prob(TH) ⟸
     0.6 → T    0.6x0.6=0.36 Prob(TT) ⟸

[1st Toss] [2nd Toss]

Prob( H @2nd Toss ) = 0.16+0.24 = 0.4

$$\begin{array}{cc} H_1 & T_1 \\ [0.4 & 0.6] \end{array} \times \begin{array}{c} H_1 \\ T_1 \end{array}\begin{bmatrix} H_2 & T_2 \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} = \begin{array}{cc} H_2 & T_2 \\ [0.4 & 0.6] \end{array}$$

$P_1$         $P_L$         $P_2$

BIOS477/877 L18 - 22

---

**23**

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**

0.6
0.4  H  T  0.6
0.4

$$P_L = \begin{array}{c} & H & T \\ H & 0.4 & 0.6 \\ T & 0.4 & 0.6 \end{array}$$

Transition probability matrix

Start

H    0.4 → H    0.4x0.4=0.16 Prob(HH) ⟸
     0.6 → T    0.4x0.6=0.24 Prob(HT) ⟸
0.4
0.6
T    0.4 → H    0.6x0.4=0.24 Prob(TH) ⟸
     0.6 → T    0.6x0.6=0.36 Prob(TT) ⟸

[1st Toss] [2nd Toss]

Prob( T @2nd Toss ) = 0.24+0.36 = 0.6

$$\begin{array}{cc} H_1 & T_1 \\ [0.4 & 0.6] \end{array} \times \begin{array}{c} H_1 \\ T_1 \end{array}\begin{bmatrix} H_2 & T_2 \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} = \begin{array}{cc} H_2 & T_2 \\ [0.4 & 0.6] \end{array}$$

$P_1$         $P_L$         $P_2$

BIOS477/877 L18 - 23

---

**24**

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**

0.6
0.4  H  T  0.6
0.4

$$P_L = \begin{array}{c} & H & T \\ H & 0.4 & 0.6 \\ T & 0.4 & 0.6 \end{array}$$

Transition probability matrix

Start

H    0.4 → H    0.4x0.4=0.16 Prob(HH) ⟸
     0.6 → T    0.4x0.6=0.24 Prob(HT) ⟸
0.4
0.6
T    0.4 → H    0.6x0.4=0.24 Prob(TH) ⟸
     0.6 → T    0.6x0.6=0.36 Prob(TT) ⟸

[0th Toss][1st Toss]

$P_1 = P_0 \times P_L$

$$\begin{array}{cc} H_1 & T_1 \\ [0.4 & 0.6] \end{array} \times \begin{array}{c} H_1 \\ T_1 \end{array}\begin{bmatrix} H_2 & T_2 \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} = \begin{array}{cc} H_2 & T_2 \\ [0.4 & 0.6] \end{array}$$

$P_0$         $P_L$         $P_1$

First probability: $P_0$

BIOS477/877 L18 - 24

4

## Slide 25

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**



Prob(HHH)=.064
Prob(HHT)=.096
Prob(HTH)=.096
Prob(HTT)=.144
Prob(THH)=.096
Prob(THT)=.144
Prob(TTH)=.144
Prob(TTT)=.216

[0th Toss] [1st Toss] [2nd Toss]

$$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$$

Transition probability matrix

$$P_2 = P_0 \times (P_L)^2$$

$$\begin{bmatrix} H_0 & T_0 \\ 0.4 & 0.6 \end{bmatrix} \times \begin{array}{c} H_0 \\ T_0 \end{array}\begin{bmatrix} H_1 & T_1 \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} \times \begin{array}{c} H_1 \\ T_1 \end{array}\begin{bmatrix} H_2 & T_2 \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} = \begin{bmatrix} H_2 & T_2 \\ 0.4 & 0.6 \end{bmatrix}$$

BIOS477/877 L18 - 25

**25**

## Slide 26

# Markov models

➤ **Markov chain**
**Loaded coin [states: H and T]**



[1st Toss] [2nd Toss] [n-1th Toss]

Initial Probability
$$P_0 = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

Transition Probability
$$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$$

$$P_n = P_0 \times (P_L)^n$$

$P_n$: Probability of the state H or T after n times of tossing

BIOS477/877 L18 - 26

**26**

## Slide 27

# What is hidden Markov model?

➤ **Markov chain**
**Fair coin *vs.* loaded coin [states: H and T]**



$$P_{F0} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$$
$$P_F = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_{L0} = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$
$$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$$

BIOS477/877 L18 - 27

**27**

## Slide 28

# What is hidden Markov model?

➤ **Markov chain**
**Fair coin *vs.* loaded coin [states: H and T]**



$$P_{F0} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$$
$$P_F = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

Which coin do you have?

$$P_{L0} = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$
$$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$$

**Try tossing the coin. Observe what happens!**

BIOS477/877 L18 - 28

**28**

## Slide 29

# What is hidden Markov model?

➤ **Markov chain**
**Fair coin *vs.* loaded coin [states: H and T]**



$$P_{F0} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$$
$$P_F = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

Which coin do you have?

$$P_{L0} = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$
$$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$$

**Observation: H, H, H**

$P_F$(HHH| Fair coin) = [?]      $P_L$(HHH| Loaded coin) = [?]

BIOS477/877 L18 - 29

**29**

## Slide 32

# What is hidden Markov model?

➤ **Markov chain**
**Fair coin *vs.* loaded coin [states: H and T]**



$$P_{F0} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$$
$$P_F = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

Which coin do you have?

$$P_{L0} = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$
$$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$$

**Observation: H, H, H**
$P_F$(HHH| Fair coin) $>$ $P_L$(HHH| Loaded coin)
Based on the observed sequence: HHH → Fair coin is more likely

BIOS477/877 L18 - 32

**32**

## Slide 33

# What is hidden Markov model?

➢ **Markov chain**
If fair and loaded coins are mixed!

$P_{FO} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$

$P_F = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$

If the coin can be changed…

$P_{LO} = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$

$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$

Observation: H, H, H

**33**

## Slide 34

# What is hidden Markov model?

➢ **Markov chain**
If fair and loaded coins are mixed!

$P_{FO} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$

$P_F = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$

Now we have to consider these probabilities, tool

$P_{LO} = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$

$P_L = \begin{array}{c} H \\ T \end{array}\begin{bmatrix} H & T \\ 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix}$

Observation: H, H, H

$P(H_F H_F H_F)$ or $P(H_L H_L H_L)$ or $P(H_F H_F H_L)$ or $P(H_F H_L H_L)$ or … ?

**34**

## Slide 35

# What is hidden Markov model?

➢ **Markov chain**
If fair and loaded coins are mixed!

start **(Initial probability)**

H: 0.5
T: 0.5

Fair

Loaded

H: 0.4
T: 0.6

$\begin{array}{cc} & F \quad L \\ F & \\ L & \end{array}\begin{bmatrix} F & L \\ 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$

- **States:** Fair and Loaded
- **Transition probability matrix →**
- Each state emits symbols H and T with certain **emission probabilities**

**35**

## Slide 36

# What is hidden Markov model?

➢ **State sequence is unknown (hidden): Fair or Loaded**

Observation:  H  T  H  T  T  T

State:  F or L  F or L  F or L **?** F or L  F or L  F or L

**Can we guess the hidden state sequence?**

What we know:

**Transition probabilities:** $\begin{array}{cc} & F \quad L \\ F & \\ L & \end{array}\begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$

**Initial probabilities:** (0.5, 0.5)

**Emission probabilities:** Fair: $\begin{array}{c} H \\ T \end{array}\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$  Loaded: $\begin{array}{c} H \\ T \end{array}\begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix}$

**36**

## Slide 37

# What is hidden Markov model?

➢ **State sequence is unknown (hidden): Fair or Loaded**

Observation:  H  T  H  T  T  T
(emission)  0.5  0.5  0.5  0.5  0.5  0.5

Fair:  $F_H$  $F_T$  $F_H$  $F_T$  $F_T$  $F_T$

0.5  0.9  0.9  0.9  0.9  0.9

start

**Probability of this path?**

Loaded:  $L_H$  $L_T$  $L_H$  $L_T$  $L_T$  $L_T$
(emission)

Prob=0.5×0.5×0.9×0.5×0.9×0.5×0.9×0.5×0.9×0.5×0.9×0.5=$(0.5)^7$×$(0.9)^5$=0.0046

[Prob(path) = emission prob. × transition prob]

**37**

## Slide 38

# What is hidden Markov model?

➢ **State sequence is unknown (hidden): Fair or Loaded**

Observation:  H  T  H  T  T  T
(emission)

Fair:  $F_H$  $F_T$  $F_H$  $F_T$  $F_T$  $F_T$

Prob=0.0046

start

0.5  0.9  0.9  0.9  0.9  0.9

Loaded:  $L_H$  $L_T$  $L_H$  $L_T$  $L_T$  $L_T$
(emission)  0.4  0.6  0.6  0.6  0.6  0.6

Prob=0.5×0.4×0.9×0.6×0.9×0.4×0.9×0.6×0.9×0.6×0.9×0.6
=0.5×$(0.4)^2$×$(0.6)^4$×$(0.9)^5$
= 0.0061

**38**

## What is hidden Markov model?

➢ State sequence is unknown (hidden): **Fair** or **Loaded**

Observation: (emission) H T H T T T

State

Fair: (emission) $F_H$ $F_T$ $F_H$ $F_T$ $F_T$ $F_T$

start

Loaded: (emission) $L_H$ $L_T$ $L_H$ $L_T$ $L_T$ $L_T$

Prob=0.5x0.5x0.9x0.5x0.1x0.4x0.9x0.6x 0.1x0.5x0.1x0.6

H: 0.5 T: 0.5 Fair 0.5 0.1 0.1 Loaded H: 0.4 T: 0.6 start 0.9 0.5 0.9

**39**

---

## What is hidden Markov model?

➢ State sequence is unknown (hidden): **Fair** or **Loaded**

Observation: (emission) H T H T T T

State

Fair: (emission) $F_H$ $F_T$ $F_H$ $F_T$ $F_T$ $F_T$

start

Loaded: (emission) $L_H$ $L_T$ $L_H$ $L_T$ $L_T$ $L_T$

Prob=?

**Prob(path) = all emission probabilities x all transition probabilities**

**Search the most likely path (maximum probability)**

**40**

---

## What is hidden Markov model?

➢ State sequence is unknown (hidden): **Fair** or **Loaded**

Observation: (emission) H T H T T T

State

Fair: (emission) $F_H$ $F_T$ $F_H$ $F_T$ $F_T$ $F_T$

start F F L L F L

Loaded: (emission) $L_H$ $L_T$ $L_H$ $L_T$ $L_T$ $L_T$

**Hidden state sequence!** **If this is the path with Max(Prob)**

**Path with the maximum probability shows the most likely path of the hidden state (F or L)**

**41**

---

## Profile hidden Markov models

➢ **Probabilistic models of multiple alignments**
   (Krogh *et al.*, 1994; Krogh 1998; Eddy 2004)

   ➔ **Closely related to standard profiles (PSSMs) introduced by Gribskov (used in PROSITE)**

   ➔ **Used in *e.g.*, PFAM, SMART, Superfamily, Panther (databases of multiple alignments and profile HMMs for domains and protein families)**

   • **States: Insertion, Deletion, and Match**
   • **Transition probabilities: between states**
   • **Emission probabilities: for 20 amino acids or 4 nucleotides**

**42**

---

## Basic architecture of a profile HMM

**Deletion states (d1, d2,...): model deletion (a gap) in the alignment**

Start 0.3 d1 d2 d3 End

0.06

i0 i1 i2 i3

**Insertion states (i0, i1,...): model insertions of random residues between two alignment positions**

0.015

m1 C

m2 C

m3 Y

A C 0.01 D E F G H I · · · Y

A C 0.5 D E F G H I · · · Y

A C D E F G H I · · · Y 0.01

**Match states (m1, m2,...): model the distribution of residues in the corresponding column of an alignment**

**Transition probabilities between states**

**Emission probabilities for the state**

**43**

---

## Basic architecture of a profile HMM

**Deletion states**

Start 0.3 0.48 0.48 0.37 End

0.4 0.3 0.06 0.46 0.06 0.46 0.73

0.49 0.48 0.49 0.48 1

0.46 0.015 0.015 0.015 0.7

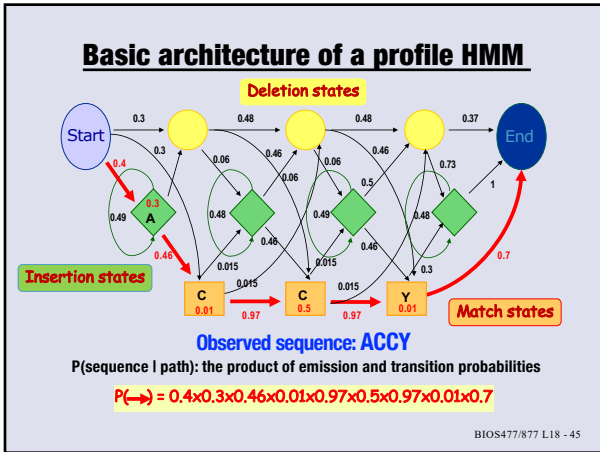**Insertion states**

0.015 0.015 0.3

0.97 0.97

**Match states**

**Observed sequence: ACCY**

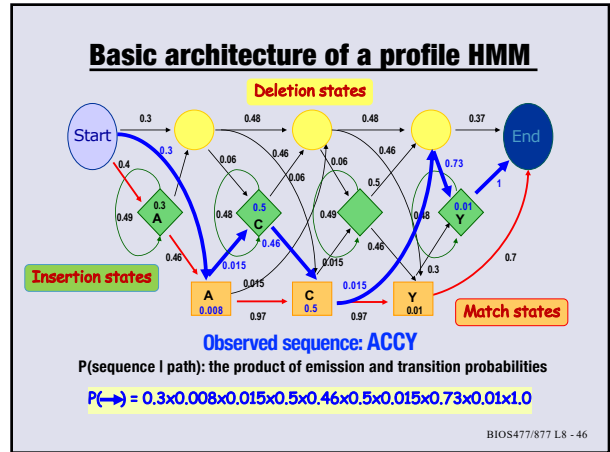• **Any sequence can be represented by a path through the model**
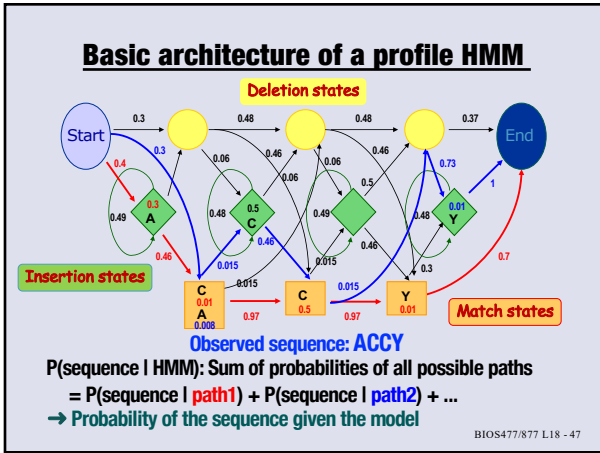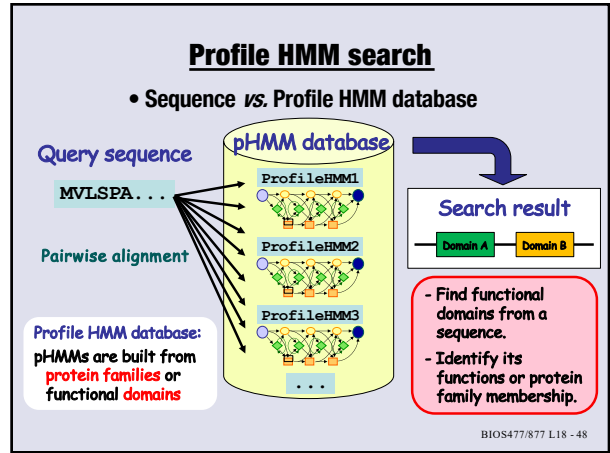• **Multiple paths are possible to model a sequence**

**44**

7

## Slide 45

**Basic architecture of a profile HMM**

Deletion states

Start → End

0.3 0.48 0.48 0.37
0.3 0.06 0.46 0.06 0.46 0.73
0.4 1
0.3 A
0.49 0.48 0.49 0.48 0.7
0.46 0.015 0.46 0.015 0.46 0.3
0.015 0.015

Insertion states

C 0.01 → C 0.5 → Y 0.01
0.97 0.97

Match states

**Observed sequence: ACCY**

P(sequence | path): the product of emission and transition probabilities

P(→) = 0.4x0.3x0.46x0.01x0.97x0.5x0.97x0.01x0.7

BIOS477/877 L18 - 45

**45**

## Slide 46

**Basic architecture of a profile HMM**

Deletion states

Start → End

0.3 0.48 0.48 0.37
0.3 0.06 0.46 0.06 0.46 0.73
0.4 1
0.3 A 0.5 C
0.49 0.48 0.46 0.49 0.48 0.01 Y
0.46 0.015 0.46 0.015 0.3

Insertion states

A 0.008 → C 0.5 → Y 0.01
0.015 0.97 0.015 0.97

Match states

**Observed sequence: ACCY**

P(sequence | path): the product of emission and transition probabilities

P(→) = 0.3x0.008x0.015x0.5x0.46x0.5x0.015x0.73x0.01x1.0

BIOS477/877 L8 - 46

**46**

## Slide 47

**Basic architecture of a profile HMM**

Deletion states

Start → End

0.3 0.48 0.48 0.37
0.3 0.06 0.46 0.06 0.46 0.73
0.4 1
0.3 A 0.5 C
0.49 0.48 0.46 0.49 0.48 0.01 Y
0.46 0.015 0.46 0.015 0.3
0.015

Insertion states

C 0.01 / A 0.008 → C 0.5 → Y 0.01
0.97 0.97

Match states

**Observed sequence: ACCY**

P(sequence | HMM): Sum of probabilities of all possible paths

= P(sequence | path1) + P(sequence | path2) + ...

➜ Probability of the sequence given the model

BIOS477/877 L18 - 47

**47**

## Slide 48

**Profile HMM search**

• Sequence *vs.* Profile HMM database

Query sequence

MVLSPA...

Pairwise alignment

pHMM database

ProfileHMM1
ProfileHMM2
ProfileHMM3
. . .

Profile HMM database: pHMMs are built from **protein families** or functional **domains**

Search result

Domain A | Domain B

- Find functional domains from a sequence.
- Identify its functions or protein family membership.

BIOS477/877 L18 - 48

**48**

## Slide 49

**Profile HMM search**

• Profile HMM *vs.* Sequence database

Query profile HMM

| HMM | A | C | D |
|-----|---|---|---|
| | m->m | m->i | |
| COMPO | 2.74297 | 3.81961 | |
| | 2.68618 | 4.42225 | |
| | 0.00206 | 6.58306 | 7. |
| 1 | 2.78640 | 5.36491 | 5. |
| | 2.68618 | 4.42225 | 2.7 |
| | 0.00206 | 6.58306 | 7.30 |

Pairwise alignment

Database (protein sequences)

Sequence1
Sequence2
Sequence3
Sequence4
Sequence5
Sequence6
Sequence7
. . .

Find sequences belonging to a protein family sharing similar functions

Search result

high similarity

Sequence28
Sequence5
Sequence11
Sequence1
Sequence73
Sequence65
Sequence33
. . .

low similarity

BIOS477/877 L18 - 49

**49**

## Slide 50

**Profile HMM databases for domain/protein families**

➢ **InterPro**: Classification of protein families
https://www.ebi.ac.uk/interpro/
➜ includes **Pfam**: a database of multiple alignments and HMMs covering many protein domains

Include external sources (e.g., Pfam, SMART, etc.)

➢ **CDD**: Conserved domain database
https://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml

➢ **SMART**: Simple Modular Architecture Research Tool
http://smart.embl-heidelberg.de/

➢ **Superfamily 2**: HMM library and genome assignments server
https://supfam.org/

➢ **Panther**: Protein ANalysis THrough Evolutionary Relationships
http://www.pantherdb.org/

• **HMMER**: profile HMM search http://www.ebi.ac.uk/Tools/hmmer/
[phmmer, hmmscan, hmmsearch, jackhmmer, ]

• **HHblits**: Homology detection by iterative HMM-HMM comparison
(MPI Bioinformatics Toolkit https://toolkit.tuebingen.mpg.de)

BIOS477/877 L18 - 50

**50**