# Bioinformatika
# Hidden Markov Models

Michael Anděl
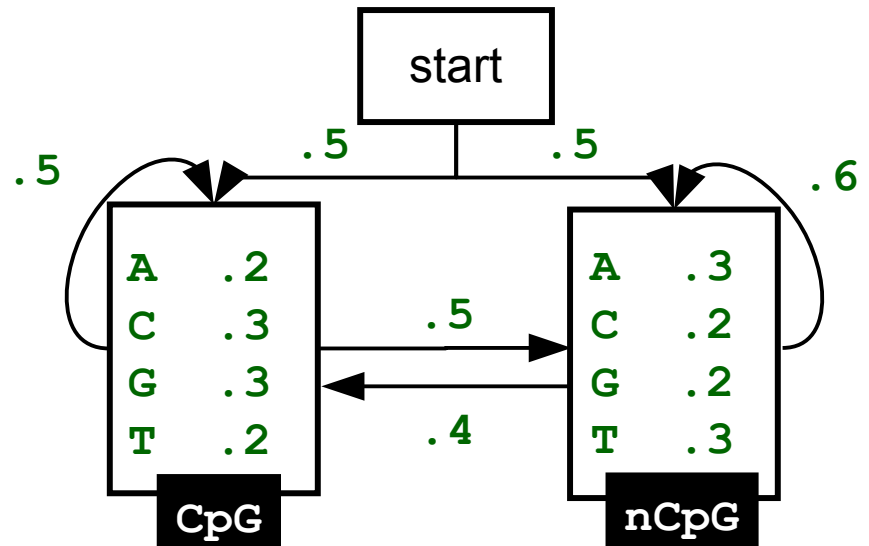(some slides are courtesy of Mark Craven, U. of Wisconsin)

# Motivation

- Sequence categorization into family of sequences (Forward alg.)

- Sequence anotation: CpG detection, gene finding (Viterbi alg.)

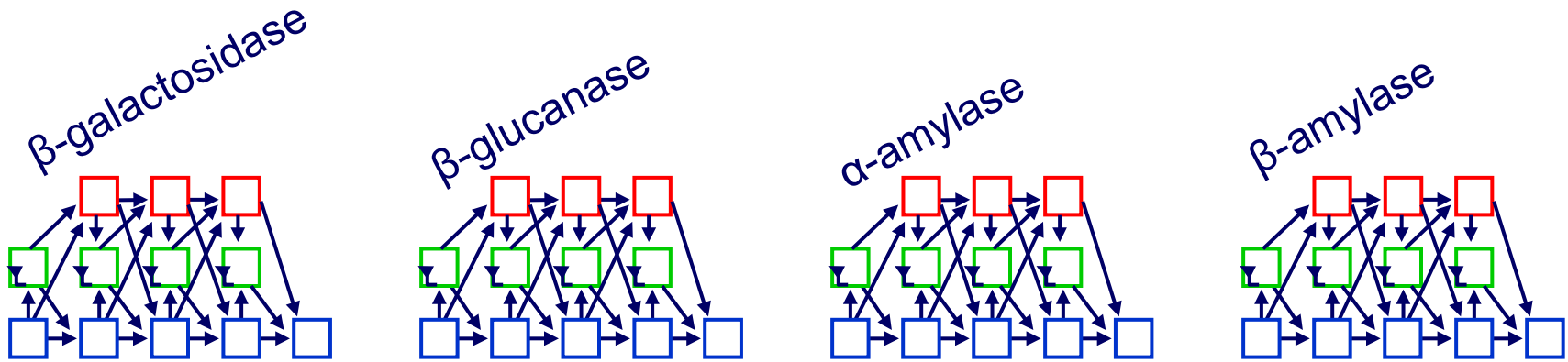- Learning hidden parameters (Baum-Welsh alg.)

# Notation

- Ex: Naïve model of CpG detection

$$p(x_i...x_N; s_i...s_N) = p(s_0)\, p(x_0|s_0) \prod_{i=1}^{N} p(x_i|s_i)\, p(s_i|s_{i-1})$$

# Forward algorithm



β-galactosidase  β-glucanase  α-amylase  β-amylase

☐ Given **K** models of **K** sequence families.

☐ Categorize a new sequence **x**.

$$p\left(\alpha-amyl.\right)p\left(x_0...x_N|\alpha-amyl.\right) \; < \; p\left(\beta-amyl.\right)p\left(x_0...x_N|\beta-amyl.\right)$$

$$p\left(x_0...x_N|\alpha-amyl.\right) \; = \; \sum_{s_0...s_N \in S^N} p\left(x_0...x_N ; s_0...s_N|\alpha-amyl.\right)$$

# Forward algorithm

$$p(x_0 x_1 \ldots x_{N-1} x_N) = \sum_{s_0 s_1 \ldots s_{N-1} s_N \in S^N} p(x_0 x_1 \ldots x_{N-1} x_N ; s_0 s_1 \ldots s_{N-1} s_N)$$

$$p(x_0 x_1 \ldots x_{N-1} x_N) = \sum_{s_0 s_1 \ldots s_{N-1} \in S^{N-1}} p(x_0 x_1 \ldots x_{N-1} ; s_0 s_1 \ldots s_{N-1})$$

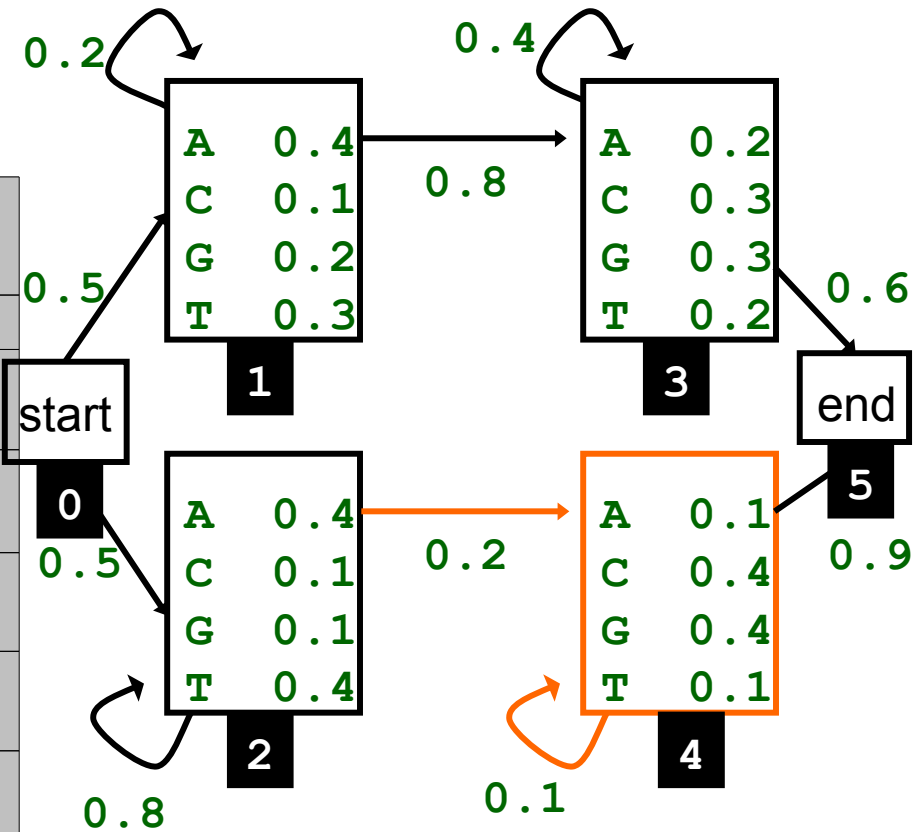$$\ldots \sum_{s_N \in S} p(s_N | s_{N-1}) p(x_N | s_N)$$

$$\boxed{p(x_0 \ldots x_i) = \sum_{s_i \in S} p(x_i | s_i) p(x_0 \ldots x_{i-1}) p(s_i | s_{i-1})}$$

# Forward algorithm (ex.)

$$\sum_{s_0 \in S} p(\epsilon|s) = 1 \ \text{ if } \ s_0 = \text{START else: } \sim 0$$

| | ε | T | A | G | A | ε |
|---|---|---|---|---|---|---|
| **0** 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| **1** | 0 | | | | | |
| **2** | 0 | | | | | |
| **3** | | | | | | |
| **4** | | | | | | |
| 5 | | | | | | |

0.2

0.4

```
A   0.4        A   0.2
C   0.1    0.8 C   0.3
G   0.2        G   0.3
0.5 T   0.3        T   0.2    0.6
        1              3
```

start

```
        0
0.5
A   0.4        A   0.1
C   0.1    0.2 C   0.4
G   0.1        G   0.4
T   0.4        T   0.1    0.9
    2              4
```

end

5

0.8

0.1

# Forward algorithm (ex.)

$$p(x_0...x_i) = \sum_{s_i \in S} p(x_i|s_i)\, p(x_0...x_{i-1})\, p(s_i|s_{i-1})$$

|   | ε | T | A | G | A | ε |
|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 0 | 0 | 0 | 0 |
| **1** | 0 | .3 x **1** x .5<br>.3 x **0** x .2<br>.15 |   |   |   |   |
| **2** | 0 | .4 x **1** x .5<br>.4 x **0** x .8<br>.2 |   |   |   |   |
| **3** |   | 0 |   |   |   |   |
| **4** |   | 0 |   |   |   |   |
| **5** |   | 0 |   |   |   |   |

0.2

0.4

**1**
```
A   0.4
C   0.1
G   0.2
T   0.3
```

**3**
```
A   0.2
C   0.3
G   0.3
T   0.2
```

0.8

0.5

0.6

start

**0**

0.5

end

**5**

0.2

0.9

**2**
```
A   0.4
C   0.1
G   0.1
T   0.4
```

**4**
```
A   0.1
C   0.4
G   0.4
T   0.1
```

0.8

0.1

# Forward algorithm (ex.)

$$p(x_0 \ldots x_i) = \sum_{s_i \in S} p(x_i | s_i)\, p(x_0 \ldots x_{i-1})\, p(s_i | s_{i-1})$$

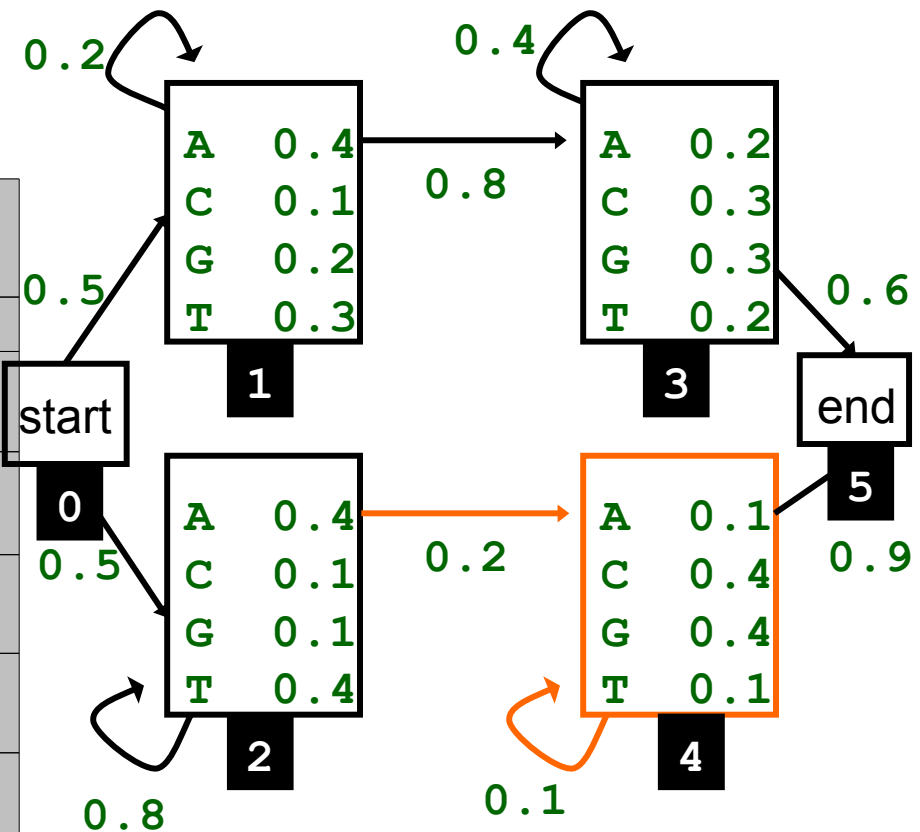|   | ε | T | A | G | A | ε |
|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 0 | 0 | 0 | 0 |
| **1** | 0 | .3 x **1** x .5<br>.3 x 0 x .2<br>.15 | .4 x **0** x .5<br>.4 x .15 x .2<br>.012 | | | |
| **2** | 0 | .4 x **1** x .5<br>.4 x 0 x .8<br>.2 | .4 x **0** x .5<br>.4 x .2 x .8<br>.064 | | | |
| **3** | | 0 | .2 x .15 x .8<br>.2 x 0 x .4<br>.024 | | | |
| **4** | | 0 | .1 x .2 x .2<br>.1 x 0 x .1<br>.004 | | | |
| **5** | | 0 | 0 | | | |

0.2

0.4

A 0.4
C 0.1
G 0.2
T 0.3

1

0.8

A 0.2
C 0.3
G 0.3
T 0.2

3

0.5

0.6

start

end

0.5

0

A 0.4
C 0.1
G 0.1
T 0.4

2

0.2

A 0.1
C 0.4
G 0.4
T 0.1

4

5

0.9

0.8

0.1

# Forward algorithm (ex.)

$$p(x_0 \dots x_i) = \sum_{s_i \in S} p(x_i|s_i)\, p(x_0 \dots x_{i-1})\, p(s_i|s_{i-1})$$

| | ε | T | A | G | A | ε |
|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 0 | 0 | 0 | |
| **1** | 0 | .3 x **1** x .5<br>.3 x 0 x .2<br>.15 | .4 x **0** x .5<br>.4 x .15 x .2<br>.012 | .2 x **0** x .5<br>.2 x .012 x .2<br>5e-4 | .4 x **0** x .5<br>.4 x 5e-4 x .2<br>4e-5 | |
| **2** | 0 | .4 x **1** x .5<br>.4 x 0 x .8<br>.2 | .4 x **0** x .5<br>.4 x .2 x .8<br>.064 | .1 x **0** x .5<br>.1 x .064 x .8<br>.00512 | .4 x **0** x .5<br>.4 x 5e-3 x .8<br>.0016 | |
| **3** | | 0 | .2 x .15 x .8<br>.2 x 0 x .4<br>.024 | .3 x .012 x .8<br>.3 x .024 x .4<br>.00576 | .2 x 5e-4 x .8<br>.2 x 6e-3 x .4<br>6e-4 | |
| **4** | | 0 | .1 x .2 x .2<br>.1 x 0 x .1<br>.004 | .4 x .064 x .2<br>.4 x .004 x .1<br>.00528 | .1 x .005 x .2<br>.1 x .005 x .1<br>1.5e-4 | |
| **5** | | 0 | 0 | 0 | 0 | |



0.2

0.4

A 0.4
C 0.1
G 0.2
T 0.3

1

A 0.2
C 0.3
G 0.3
T 0.2

3

0.8

0.5

0.6

start

0

end

5

0.5

A 0.4
C 0.1
G 0.1
T 0.4

2

0.2

A 0.1
C 0.4
G 0.4
T 0.1

4

0.9

0.8

0.1

# Forward algorithm (ex.)

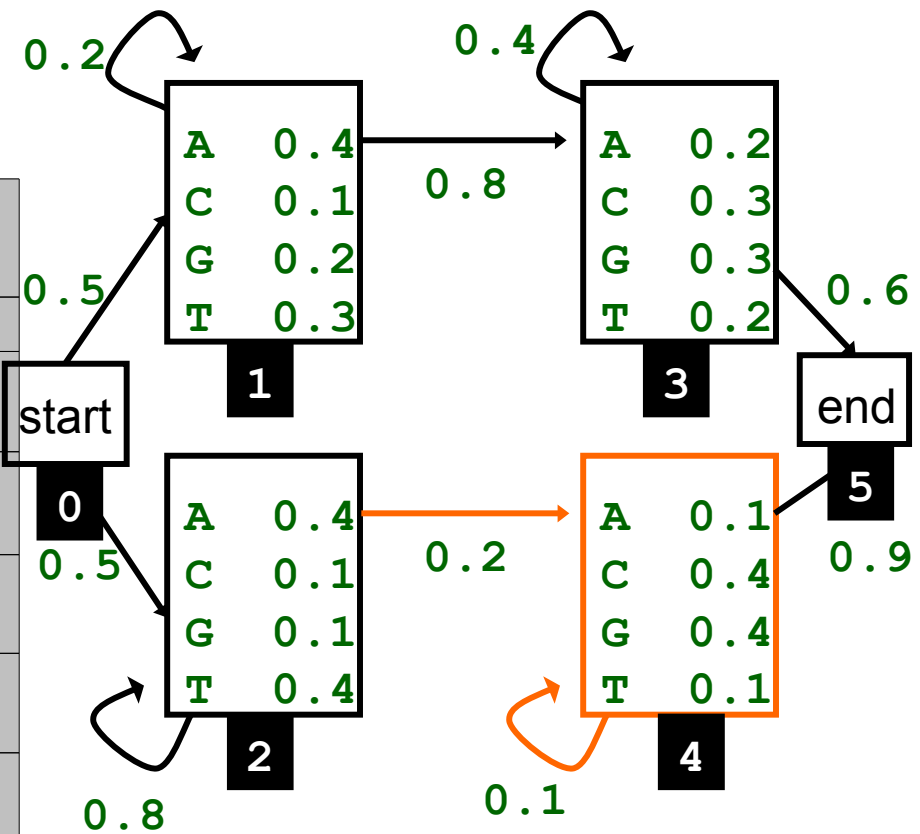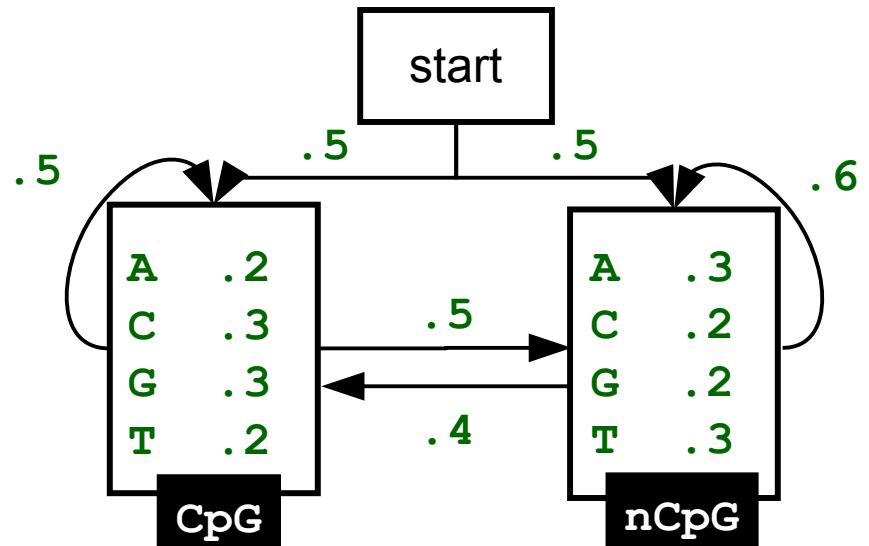$$p(x_0 \ldots x_i) = \sum_{s_i \in S} p(x_i | s_i)\, p(x_0 \ldots x_{i-1})\, p(s_i | s_{i-1})$$

|   | ε | T | A | G | A | ε |
|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 0 | 0 | 0 | |
| **1** | 0 | .3 x **1** x .5<br>.3 x 0 x .2<br>.15 | .4 x **0** x .5<br>.4 x .15 x .2<br>.012 | .2 x **0** x .5<br>.2 x .012 x .2<br>5e-4 | .4 x **0** x .5<br>.4 x 5e-4 x .2<br>4e-5 | 0 |
| **2** | 0 | .4 x **1** x .5<br>.4 x 0 x .8<br>.2 | .4 x **0** x .5<br>.4 x .2 x .8<br>.064 | .1 x **0** x .5<br>.1 x .064 x .8<br>.00512 | .4 x **0** x .5<br>.4 x 5e-3 x .8<br>.0016 | 0 |
| **3** | | 0 | .2 x .15 x .8<br>.2 x 0 x .4<br>.024 | .3 x .012 x .8<br>.3 x .024 x .4<br>.00576 | .2 x 5e-4 x .8<br>.2 x 6e-3 x .4<br>6e-4 | 0 |
| **4** | | 0 | .1 x .2 x .2<br>.1 x 0 x .1<br>.004 | .4 x .064 x .2<br>.4 x .004 x .1<br>.00528 | .1 x .005 x .2<br>.1 x .005 x .1<br>1.5e-4 | 0 |
| 5 | | 0 | 0 | 0 | 0 | 1e-4 x .6<br>6e-4 x .9<br>**4.6e-4** |

0.2

0.4

| 1 | | 3 | |
|---|---|---|---|
| A | 0.4 | A | 0.2 |
| C | 0.1 | C | 0.3 |
| G | 0.2 | G | 0.3 |
| T | 0.3 | T | 0.2 |

0.8

0.5

0.6

start

0

0.5

end

5

| 2 | | 4 | |
|---|---|---|---|
| A | 0.4 | A | 0.1 |
| C | 0.1 | C | 0.4 |
| G | 0.1 | G | 0.4 |
| T | 0.4 | T | 0.1 |

0.2

0.9

0.8

0.1

# Viterbi algorithm

- Given an observed sequence **x.**

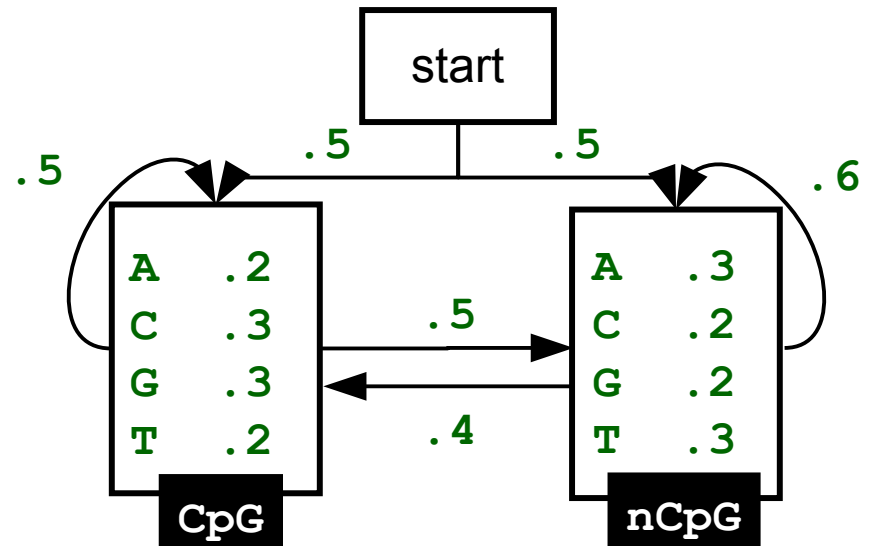- What is the most likely path **s** through the model, i.e. sequence anotation?

$$s^* = \arg\max_{s_0...s_N \in S^N} p\left(x_0...x_N ; s_0...s_N\right)$$

# Viterbi algorithm (ex.)

| | ε | A | T | G | G | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|---|
| **START** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **CpG** | 0 | | | | | | | | | |
| **nCpG** | 0 | | | | | | | | | |

$$\max p\left(\epsilon | s_0\right) = 1 \text{ if } s_0 = \text{START else: } \sim 0$$

# Viterbi algorithm (ex.)

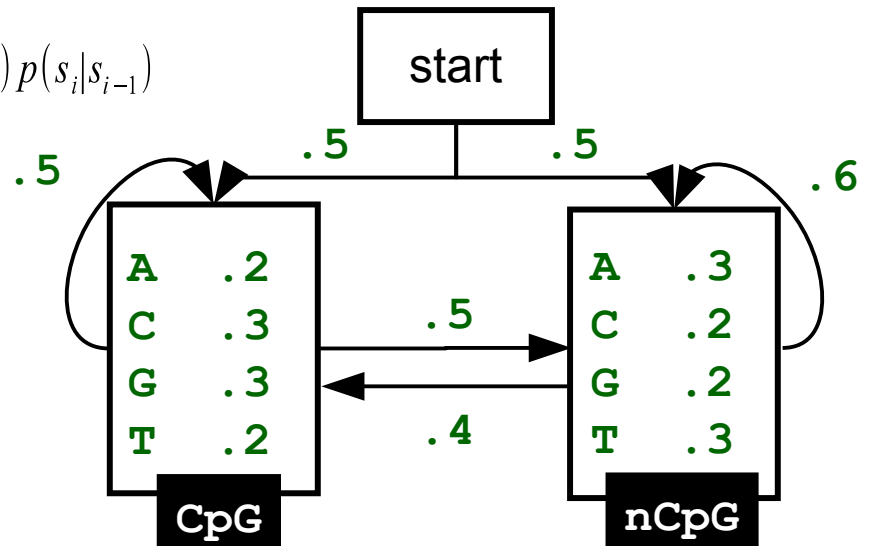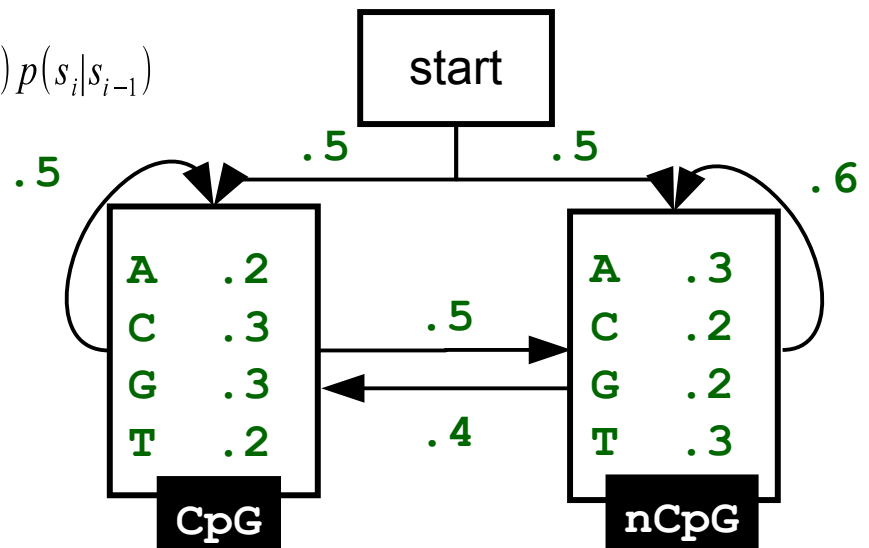| | ε | A | T | G | G | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|---|
| **START** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **CpG** | 0 | .2 x **1** x .5<br>.2 x 0 x .5<br>.2 x 0 x .4<br>.1 | | | | | | | | |
| **nCpG** | 0 | .3 x **1** x .5<br>.3 x 0 x .5<br>.3 x 0 x .6<br>.15 | | | | | | | | |

$$\max_{s_i \in S} p(x_0 \ldots x_i | s_i) = \max_{s_i \in S} p(x_i | s_i) \max_{s_{i-1} \in S} p(x_0 \ldots x_{i-1} | s_{i-1}) p(s_i | s_{i-1})$$

# Viterbi algorithm (ex.)

| | ε | A | T | G | G | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|---|
| **START** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **CpG** | 0 | .2 x **1** x .5<br>.2 x 0 x .5<br>.2 x 0 x .4<br>.1 | .2 x **0** x .5<br>.2 x .1 x .5<br>.2 x .15 x .4<br>.012 | | | | | | | |
| **nCpG** | 0 | .3 x **1** x .5<br>.3 x 0 x .5<br>.3 x 0 x .6<br>.15 | .3 x **0** x .5<br>.3 x .1 x .5<br>.3 x .15 x .6<br>.027 | | | | | | | |

$$\max_{s_i \in S} p\left(x_0 \ldots x_i | s_i\right) = \max_{s_i \in S} p\left(x_i | s_i\right) \max_{s_{i-1} \in S} p\left(x_0 \ldots x_{i-1} | s_{i-1}\right) p\left(s_i | s_{i-1}\right)$$



start

.5   .5

.5    .6

| CpG | |
|---|---|
| A | .2 |
| C | .3 |
| G | .3 |
| T | .2 |

.5

.4

| nCpG | |
|---|---|
| A | .3 |
| C | .2 |
| G | .2 |
| T | .3 |

# Viterbi algorithm (ex.)

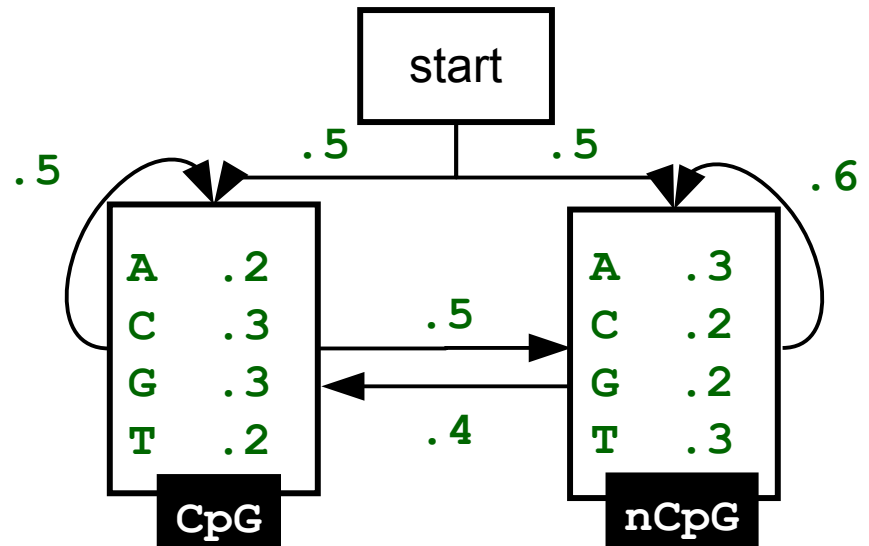| | ε | A | T | G | G | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|---|
| **START** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **CpG** | 0 | .2 x **1** x .5<br>.2 x 0 x .5<br>.2 x 0 x .4<br>.1 | .2 x **0** x .5<br>.2 x .1 x .5<br>.2 x .15 x .4<br>.012 | 0<br>.3 x .012 x .5<br>.3 x .027 x .4<br>.0032 | 0<br>.3 x .0032 x .5<br>.3 x .0032 x .4<br>5e-4 | 0<br>.3 x .012 x .5<br>.3 x .027 x .4<br>5e-5 | | | | |
| **nCpG** | 0 | .3 x **1** x .5<br>.3 x 0 x .5<br>.3 x 0 x .6<br>.15 | .3 x **0** x .5<br>.3 x .1 x .5<br>.3 x .15 x .6<br>.027 | 0<br>.2 x .012 x .5<br>.2 x .027 x .6<br>.0032 | 0<br>.2 x .0032 x .5<br>.2 x .0032 x .6<br>4e-4 | 0<br>.2 x .012 x .5<br>.2 x .027 x .6<br>4e-5 | | | | |

$$\max_{s_i \in S} p(x_0 \ldots x_i | s_i) \ = \ \max_{s_i \in S} p(x_i | s_i) \ \max_{s_{i-1} \in S} p(x_0 \ldots x_{i-1} | s_{i-1}) \, p(s_i | s_{i-1})$$

# Viterbi algorithm (ex.)

| | ε | A | T | G | G | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|---|
| **START** | 0 | -inf | -inf | -inf | -inf | -inf | -inf | -inf | -inf | -inf |
| **CpG** | -inf | ln.2+**0**+ln.5<br>ln.2+-inf+ln.5<br>ln.2+-inf+ln.4<br>-2.30 | ln.2+-**inf**+ln.5<br>ln.2+0+ln.5<br>ln.2+ln.15+ln.4<br>-2.3 | | | | | | | |
| **nCpG** | -inf | ln.3+**0**+ln.5<br>ln.3+-inf+ln.5<br>ln.3+-inf+ln.6<br>-1.9 | **-inf**<br>ln.3+ln.1+ln.5<br>ln.3+ln.15+ln.6<br>-1.9 | | | | | | | |

$$\arg\max_{s_i \in S} p(x_0 \ldots x_i | s_i) \;=\; \arg\max_{s_i \in S} \; \log p(x_0 \ldots x_i | s_i)$$

start

.5 .5

.5 .6

CpG
```
A   .2
C   .3
G   .3
T   .2
```

nCpG
```
A   .3
C   .2
G   .2
T   .3
```

.5

.4

# Assignement

- Gene finding

http://www.biostat.wisc.edu/~craven/776/hw3.html

- Need not implement viterb algorithm. You can use an arbitrary solver.

- Alternative to gene finding assignement:

  - Viterbi and Forward implementation (or gene expression)

  - Sequence assembly