

# Motif finding with Gibbs sampling example

Colin Dewey

2007.02.13

# The Data

- Hidden motif of width 7 in 4 sequences of length 10
- Each motif occurrence differs from consensus (GATTACA) in two positions

ACCATGACAG  
GAGTATAACCT  
CATGCTTACT  
CGGAATGCAT

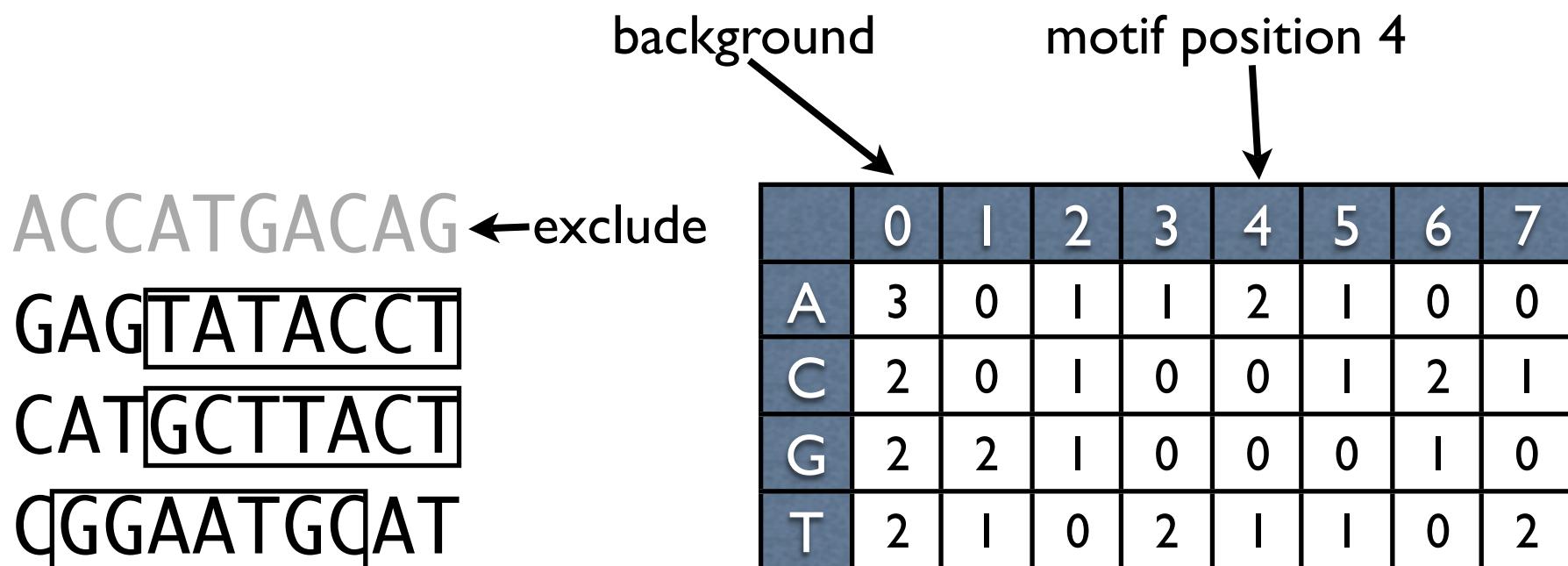
# Initialization

- Choose initial positions of motif at random

ACCATGACAG  
GAGTATAACCT  
CATGCTTACT  
CGGAATGCAT

# Predictive update step

- Update profile matrix based on motif and background frequencies and pseudocounts



# Predictive update step

- Calculate profile matrix from frequencies and pseudocounts

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| A | 3 | 0 | 1 | 1 | 2 | 1 | 0 | 0 |
| C | 2 | 0 | 1 | 0 | 0 | 1 | 2 | 1 |
| G | 2 | 2 | 1 | 0 | 0 | 0 | 1 | 0 |
| T | 2 | 1 | 0 | 2 | 1 | 1 | 0 | 2 |

$$p_{1,A} = \frac{c_{1,A} + b_A}{N - 1 + B} = \frac{0 + 0.5}{4 - 1 + 2} = 0.1$$

$$p_{1,C} = \frac{c_{1,C} + b_C}{N - 1 + B} = \frac{0 + 0.5}{4 - 1 + 2} = 0.1$$

$$p_{1,G} = \frac{c_{1,G} + b_G}{N - 1 + B} = \frac{2 + 0.5}{4 - 1 + 2} = 0.5$$

$$p_{1,T} = \frac{c_{1,T} + b_T}{N - 1 + B} = \frac{1 + 0.5}{4 - 1 + 2} = 0.3$$

$$p_{0,A} = \frac{c_{0,A} + b_A}{\sum_{i \neq 1} (\ell_i - W) + B} = \frac{3 + 0.5}{3(3) + 2} = \frac{7}{22}$$

$$p_{0,C} = \frac{c_{0,C} + b_C}{\sum_{i \neq 1} (\ell_i - W) + B} = \frac{2 + 0.5}{3(3) + 2} = \frac{5}{22}$$

$$p_{0,G} = \frac{c_{0,G} + b_G}{\sum_{i \neq 1} (\ell_i - W) + B} = \frac{2 + 0.5}{3(3) + 2} = \frac{5}{22}$$

$$p_{0,T} = \frac{c_{0,T} + b_T}{\sum_{i \neq 1} (\ell_i - W) + B} = \frac{2 + 0.5}{3(3) + 2} = \frac{5}{22}$$

|   | 0    | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
|---|------|-----|-----|-----|-----|-----|-----|-----|
| A | 0.31 | 0.1 | 0.3 | 0.3 | 0.5 | 0.3 | 0.1 | 0.1 |
| C | 0.23 | 0.1 | 0.3 | 0.1 | 0.1 | 0.3 | 0.5 | 0.3 |
| G | 0.23 | 0.5 | 0.3 | 0.1 | 0.1 | 0.1 | 0.3 | 0.1 |
| T | 0.23 | 0.3 | 0.1 | 0.5 | 0.3 | 0.3 | 0.1 | 0.5 |

# Sampling step

- For each possible motif start position, calculate ratio of likelihood of next  $W$  positions from motif vs. background

|       | 1    | 2    | 3    | 4     |
|-------|------|------|------|-------|
| $A_i$ | 0.16 | 0.13 | 0.26 | 0.017 |

$$A_1 = \frac{p_{1,A} \cdot p_{2,C} \cdot p_{3,C} \cdot p_{4,A} \cdot p_{5,T} \cdot p_{6,G} \cdot p_{7,A}}{p_{0,A} \cdot p_{0,C} \cdot p_{0,C} \cdot p_{0,A} \cdot p_{0,T} \cdot p_{0,G} \cdot p_{0,A}} \approx \frac{0.1 \cdot 0.3 \cdot 0.1 \cdot 0.5 \cdot 0.3 \cdot 0.3 \cdot 0.1}{0.31 \cdot 0.23 \cdot 0.23 \cdot 0.31 \cdot 0.23 \cdot 0.23 \cdot 0.31} \approx 0.16$$

ACCATGACAG

|   | 0    | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
|---|------|-----|-----|-----|-----|-----|-----|-----|
| A | 0.31 | 0.1 | 0.3 | 0.3 | 0.5 | 0.3 | 0.1 | 0.1 |
| C | 0.23 | 0.1 | 0.3 | 0.1 | 0.1 | 0.3 | 0.5 | 0.3 |
| G | 0.23 | 0.5 | 0.3 | 0.1 | 0.1 | 0.1 | 0.3 | 0.1 |
| T | 0.23 | 0.3 | 0.1 | 0.5 | 0.3 | 0.3 | 0.1 | 0.5 |

# Sampling step

- Sample new position  $i$  in chosen sequence based on  $A_i$

|       |      |      |      |       |
|-------|------|------|------|-------|
|       | 1    | 2    | 3    | 4     |
| $A_i$ | 0.16 | 0.13 | 0.26 | 0.017 |



normalize

|       |      |      |      |      |
|-------|------|------|------|------|
|       | 1    | 2    | 3    | 4    |
| $A_i$ | 0.28 | 0.23 | 0.46 | 0.03 |



draw random sample  
from distribution

$$a_3 = 2$$



ACCATGACAG

# Calculate likelihood

- Calculate likelihood (or some related value) after each iteration
- Iterate:
  - choose sequence
  - predictive update
  - sample new motif position in sequence
- After many iterations, choose motif positions and corresponding profile matrix