# Genome 540 discussion

#### February 18th, 2025 Joe Min





Homework 6 overview

Poisson distributions

#### Homework 6

### Homework 6 overview

- 1. Implement the D-segment algorithm to identify regions of elevated copy number variation (CNVs)
- 2. Run your program on chromosome 16 of the CHM13 long-read-based genome assembly

Goals:

- Identify all high-scoring regions of the genome according to some scoring scheme
- In O(n) time, identify boundaries of these regions by identifying regions with large "drops" in score



#### Whole region has total score of 105 but...



Whole region has total score of 105 but...

Discontinuous high-scoring regions might have different biological significance



To find \*all\* high-scoring segments (i.e., all regions with score >= S), we could keep track of region best starts and stops to maximize region scores but this requires going backwards, which could be  $O(n^2)$  in the worst case



We can instead keep track of a region's maximum cumulative score, as well as a cumulative drop from that to identify the region's ending boundary



We can instead keep track of a region's maximum cumulative score, as well as a cumulative drop from that to identify the region's ending boundary

```
O(N) algorithm to find all maximal D-segs:
cumul = max = 0; start = 1;
for (i = 1; i \le N; i^{++})
     cumul += s[i];
     if (cumul \geq max)
          \{\max = \operatorname{cumul}; \operatorname{end} = i;\}
     if (\text{cumul} \le 0 \text{ or cumul} \le \text{max} + D \text{ or } i == N) {
          if (\max \ge S)
            {print start, end, max; }
          max = cumul = 0; start = end = i + 1; /* NO BACKTRACKING
            NEEDED! */
```

## Applying D-segment reveals CNVs



| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = 0   |
|-------------------|--------------|-----------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = 0 |
| <b>2:</b> 1.0     |              | end = $0$ |
| > <b>=3:</b> 1.25 |              | cumul = 0 |

O(N) algorithm to find all maximal D-segs: cumul = max = 0; start = 1;  $for (i = 1; i \le N; i++) \{$  cumul += s[i];  $if (cumul \ge max)$   $\{max = cumul; end = i;\}$   $if (cumul \le 0 \text{ or cumul} \le max + D \text{ or } i == N) \{$   $if (max \ge S)$   $\{print \text{ start, end, max;} \}$  max = cumul = 0; start = end = i + 1; /\* NO BACKTRACKING NEEDED! \*/  $\}$ 

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = 0          |
|-------------------|--------------|------------------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = <b>1</b> |
| <b>2:</b> 1.0     |              | end = <b>1</b>   |
| > <b>=3:</b> 1.25 |              | cumul = 0        |

O(N) algorithm to find all maximal D-segs: cumul = max = 0; start = 1;  $for (i = 1; i \le N; i++) \{$  cumul += s[i];  $if (cumul \ge max)$   $\{max = cumul; end = i;\}$   $if (cumul \le 0 \text{ or cumul} \le max + D \text{ or } i == N) \{$   $if (max \ge S)$   $\{print \text{ start, end, max;} \}$  max = cumul = 0; start = end = i + 1; /\* NO BACKTRACKING NEEDED! \*/

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = <b>0.5</b>   |
|-------------------|--------------|--------------------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = <b>2</b>   |
| <b>2:</b> 1.0     |              | end = <b>2</b>     |
| > <b>=3:</b> 1.25 |              | cumul = <b>0.5</b> |

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = <b>1.5</b>   |
|-------------------|--------------|--------------------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = 2          |
| <b>2:</b> 1.0     |              | end = <b>3</b>     |
| > <b>=3:</b> 1.25 |              | cumul = <b>1.5</b> |

O(N) algorithm to find all maximal D-segs: cumul = max = 0; start = 1;  $for (i = 1; i \le N; i++) \{$  cumul += s[i];  $if (cumul \ge max)$   $\{max = cumul; end = i;\}$   $if (cumul \le 0 \text{ or cumul} \le max + D \text{ or } i == N) \{$   $if (max \ge S)$   $\{print \text{ start, end, max; }\}$  max = cumul = 0; start = end = i + 1; /\* NO BACKTRACKING NEEDED! \*/  $\}$ 

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>D:</b> -3 | max = <b>2.75</b>           |
|--------------|-----------------------------|
| <b>S:</b> 3  | start = 2                   |
|              | end = <b>4</b>              |
|              | cumul = <b>2.75</b>         |
|              | <b>D:</b> -3<br><b>S:</b> 3 |

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = 2.75         |
|-------------------|--------------|--------------------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = 2          |
| <b>2:</b> 1.0     |              | end = <b>5</b>     |
| > <b>=3:</b> 1.25 |              | cumul = <b>2.0</b> |

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75 | <b>D:</b> -3 | max = <b>3.0</b>   |
|-----------------|--------------|--------------------|
| <b>1:</b> 0.50  | <b>S:</b> 3  | start = 2          |
| <b>2:</b> 1.0   |              | end = 6            |
| >=3: 1.25       |              | cumul = <b>3.0</b> |

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = 3.0           |
|-------------------|--------------|---------------------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = 2           |
| <b>2:</b> 1.0     |              | end = $6$           |
| > <b>=3:</b> 1.25 |              | cumul = <b>2.25</b> |

O(N) algorithm to find all maximal D-segs: cumul = max = 0; start = 1;  $for (i = 1; i \le N; i++) \{$  cumul += s[i];  $if (cumul \ge max)$   $\{max = cumul; end = i;\}$   $if (cumul \le 0 \text{ or cumul} \le max + D \text{ or } i == N) \{$   $if (max \ge S)$   $\{print \text{ start, end, max;} \}$  max = cumul = 0; start = end = i + 1; /\* NO BACKTRACKING NEEDED! \*/  $\}$ 

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75   | <b>D:</b> -3 | max = 3.0          |
|-------------------|--------------|--------------------|
| <b>1:</b> 0.50    | <b>S:</b> 3  | start = 2          |
| <b>2:</b> 1.0     |              | end = $6$          |
| > <b>=3:</b> 1.25 |              | cumul = <b>1.5</b> |

| Position          | 1     | 2    | 3   | 4    | 5     | 6   | 7     | 8     |
|-------------------|-------|------|-----|------|-------|-----|-------|-------|
| Read Start Counts | 0     | 1    | 2   | 3    | 0     | 2   | 0     | 0     |
| Score             | -0.75 | 0.50 | 1.0 | 1.25 | -0.75 | 1.0 | -0.75 | -0.75 |

#### Sample scoring scheme

| <b>0:</b> -0.75 | <b>D:</b> -3 | max = 3.0              |
|-----------------|--------------|------------------------|
| <b>1:</b> 0.50  | <b>S:</b> 3  | <mark>start = 2</mark> |
| <b>2:</b> 1.0   |              | <mark>end = 6</mark>   |
| >=3: 1.25       |              | cumul = 1.5            |

O(N) algorithm to find all maximal D-segs: cumul = max = 0; start = 1;  $for (i = 1; i \le N; i++) \{$  cumul += s[i];  $if (cumul \ge max)$   $\{max = cumul; end = i;\}$   $if (cumul \le 0 \text{ or cumul} \le max + D \text{ or } i == N) \{$   $if (max \ge S)$   $\{print \text{ start, end, max; }\}$  max = cumul = 0; start = end = i + 1; /\* NO BACKTRACKING NEEDED! \*/  $\}$ 

#### Poisson distributions

## What is the Poisson distribution?

- Used for discrete, countable data over given intervals
- Classic example: yearly deaths by horsekick in the Prussian army (from von Bortkiewicz, 1898)
- Many years had zero deaths
- Some years had one
- Fewer had two, etc.



### What is the Poisson distribution?

Because we are dealing with count data

- Support only from natural numbers; thus the distribution often shifts left toward zero
- Shape depends on the parameter (λ), which in turn relies on the likelihood of each event



### Defining the Poisson

A random variable X has a Poisson distribution if its probability function is defined by:

$$f(k;\lambda) = \Pr(X{=}k) = rac{\lambda^k e^{-\lambda}}{k!}$$

That is; the probability that the measured count X has value k is given by  $\frac{\lambda^k e^{-\lambda}}{k!}$ 

- k represents numbers of occurences (0, 1, 2, ...)
- $\lambda$  is both the mean and variance of f(k;  $\lambda$ )

### Working with the Poisson

Has nice properties:

- $\lambda$  is both the mean and variance of f(k;  $\lambda$ )
- Serves as an approximation of the binomial distribution (if the number of events is sufficiently high and the probability of each event is sufficiently low)
- Has a known conjugate prior (Gamma), giving a closed-form posterior for use in Bayesian inference (basically; it's easier to update models with new data)



Reminder:

Homework 6 is due Sunday, February 23rd at 11:59pm!