# CS 533C project update

## Visualization
## of the evolutions of a virus genetic sequence

Mahshid ZeinalyBaraghoush (mzeinaly@sfu.ca)

# Overview

- The target domain + my goals

- Original dataset/tasks

- Translate dataset :

Data abstraction $\longrightarrow$ Data visual encoding

- Translate task:

Task abstraction $\longrightarrow$ Visual encoding/interaction

- Next steps

# The target domain

- The IEEE Vast Challenge 2010
- Evaluate IMAS "Interactive Multi-genomic Analysis System"

Goal:

Enhance IMAS in an iterative process using infovis design guidelines

# Domain vocabulary

- Gene sequence

    ATGCACCGCCCTGCGCAGTTCATAG

- Virus replication ⟶ Virus strains

- Substitutions ⟶ Mutations of virus strains

    ATGCACCGCCCTGCGCAGTTCATAG

    ATGCACCGCCGTGCGCAGTTCATAG

- Disease characteristics

# Domain: dataset

## Tabular

### Virus strain sequences table:

100 virus strains

1000 nucleotides long

| Sequence ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 118 | A | T | C | G | A | C | T | C |
| 200 | | | | | | | | |
| 19 | A | T | C | G | C | G | T | A |
| | A | T | G | G | C | C | T | C |

### Disease characteristics table:
100 virus strains

10 disease characteristics

| Sequence ID | Symptoms | Drug Resistance |
|---|---|---|
| 118 | Mild | High |
| 200 | Severe | Mild |
| 19 | Moderate | High |

# Domain: task

- Identify substitutions that lead to an increase in different disease characteristic's severity.

- Output for a particular characteristics:

Severity-driven substitutions in the order of importance:

1. G $\rightarrow$ A, 513  +  T $\rightarrow$ A, 907

2. A $\rightarrow$ G, 39

# Data abstraction: Original dataset

## Tabular

**Row : strain sequence**
**Column: position in aligned sequences**
**Cell attribute: nominal data**

**Row: strain characteristics info**
**Column: a characteristics**
**Cell attribute: ordinal variable**

| Sequence ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------|---|---|---|---|---|---|---|---|
| 118 | A | T | C | G | A | C | T | C |
| 200 | | | | | | | | |
| 19 | A | T | C | G | C | G | T | A |
| | A | T | G | G | C | C | T | C |

| Sequence ID | Symptoms | Drug Resistance |
|-------------|----------|-----------------|
| 118 | Mild | High |
| 200 | Severe | Mild |
| 19 | Moderate | High |

# Data abstraction: derived dataset

## Tabular

### Row : strain sequence
### Column: position in aligned sequences
Derived cell variable: whether there is a substitution from the first row

| Sequence ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 118 | A | T | C | G | A | C | T | C |
| 200 |   |   |   |   |   |   |   |   |
| 19 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
|   | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |

### Row: strain characteristics
### Column: a Characteristics
### Cell: ordinal variable
Derived characteristics from the originals

| Sequence ID | Symptoms | Drug Resistance |
|---|---|---|
| 118 | Mild | High |
| 200 | Severe | Mild |
| 19 | Moderate | High |

# Data: visual encoding

| Labels | 1 | 2 | 3 | 4 | 5 | 6 |
|--------|---|---|---|---|---|---|
|        | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | . | . | . | . |
| 197 | . | . | . | . | . | . |
| 248 | . | . | . | . | C | . |
| 253 | . | . | . | . | . | . |
| 418 | . | . | C | . | . | . |
| 583 | . | . | . | . | C | . |

Note:
Color justifications

# IMAS multi-alignment view
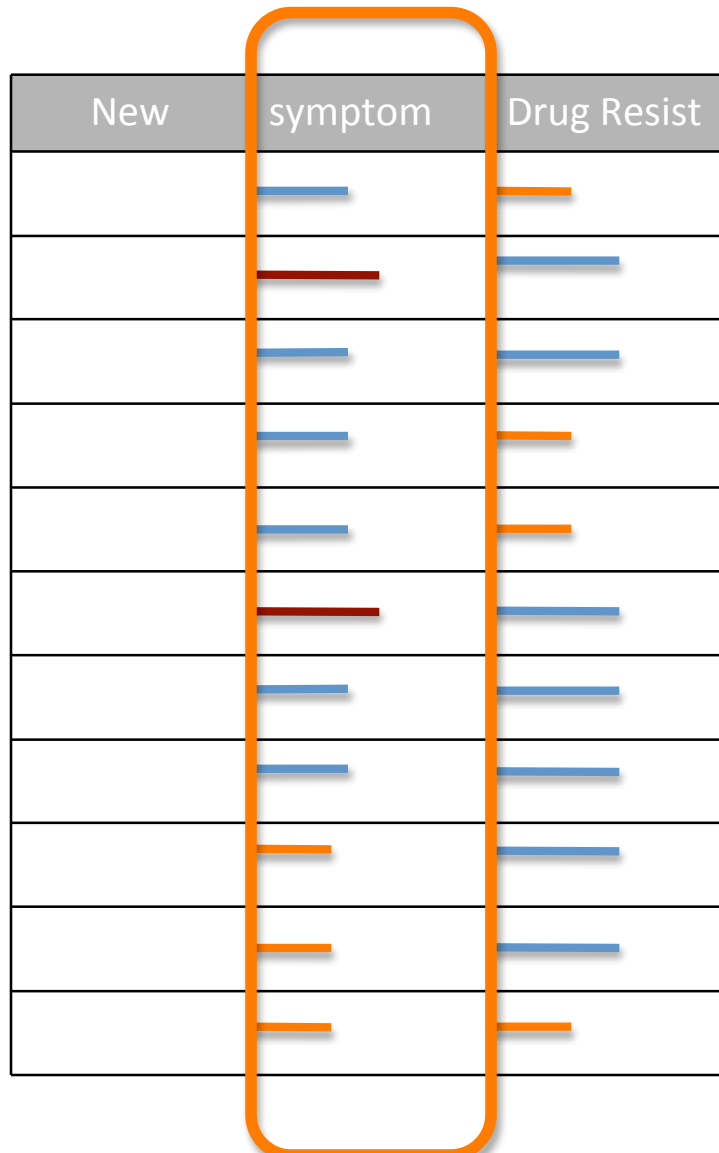
# Task abstraction

- Find columns that their substitution lead to the most severe rows. Report them individually or in a combination with other columns:

1. Sort the rows

2. Filter the columns until you find the answer:

    2.1 Find interesting columns

    2.2 Find related columns

       ( and filter the rest )

# Task abstraction

- Find columns that their substitution lead to the most severe rows. Report them individually or in a combination with other columns:

1. Sort the rows

2. Filter the columns until you find the answer:

    2.1 Find interesting columns

    2.2 Find related columns

       ( and filter the rest )

# Dataset/Task visual encoding

| New | symptom | Drug Resist |
|-----|---------|-------------|
| | ▬ (blue) | ▬ (orange) |
| | ▬ (dark red) | ▬ (blue) |
| | ▬ (blue) | ▬ (blue) |
| | ▬ (blue) | ▬ (orange) |
| | ▬ (blue) | ▬ (orange) |
| | ▬ (dark red) | ▬ (blue) |
| | ▬ (blue) | ▬ (blue) |
| | ▬ (blue) | ▬ (blue) |
| | ▬ (orange) | ▬ (blue) |
| | ▬ (orange) | ▬ (blue) |
| | ▬ (orange) | ▬ (orange) |

| | 0 | 1 | 2 | 3 | 4 | 5 |
|--------|---|---|---|---|---|---|
| Labels | A | A | T | A | T | A |
| 705 | G | . | C | . | . | . |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | G | . | C | . | . | . |
| 501 | . | . | . | . | C | T |
| 118 | . | C | . | . | C | T |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | . | . | C | G | . | . |
| 583 | . | . | . | . | C | T |

# "Table Lens" view



| Labels | 0 | 1 | 2 | 3 | 4 | 5 |
|--------|---|---|---|---|---|---|
|        | A | A | T | A | T | A |
| 705 | G | . | C | . | . | . |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | G | . | C | . | . | . |
| 501 | . | . | . | . | C | T |
| 118 | . | C | . | . | C | T |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | . | . | C | G | . | . |
| 583 | . | . | . | . | C | T |

# Table Lens view

| New | symptom | Drug Resist |
|---|---|---|



| Labels | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

# Table Lens view

| | New | symptom | Drug Resist |
|---|---|---|---|

| Labels | 0 A | 1 A | 2 T | 3 A | 4 T | 5 A |
|---|---|---|---|---|---|---|
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

# Table Lens view

| | New | symptom | Drug Resist |
|---|---|---|---|
| | | ▬ (dark red) | ▬ (blue) |
| | | ▬ (dark red) | ▬ (blue) |
| | | ▬ (blue) | ▬ (blue) |
| | | ▬ (blue) | ▬ (orange) |
| | | ▬ (blue) | ▬ (orange) |
| | | ▬ (blue) | ▬ (orange) |
| | | ▬ (blue) | ▬ (blue) |
| | | ▬ (blue) | ▬ (blue) |
| | | ▬ (orange) | ▬ (blue) |
| | | ▬ (orange) | ▬ (blue) |
| | | ▬ (orange) | ▬ (orange) |

Overall =
Symptom +
Drug Resistance

Create

| Labels | 0 A | 1 A | 2 T | 3 A | 4 T | 5 A |
|---|---|---|---|---|---|---|
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

# Dataset/Task visual encoding

# Task abstraction

- Find columns that their substitution lead to the most severe rows. Report them individually or in a combination with other columns:

1. Sort the rows ✔

1. Filter the columns until you find the answer:

    2.1  Find interesting columns and/or

    2.2  Find related columns

       ( and filter the rest )

# Task abstraction

- Find columns that their substitution lead to the most severe rows. Report them individually or in a combination with other columns:

1. Sort the rows ✔

2. Filter the columns until you find the answer:

    2.1  Find interesting columns and/or

    2.2  Find related columns

        ( and filter the rest )

# Comparison tasks challenges

- Navigation + Comparison: increases the memory load
- The order of columns should be remained

# Interaction design: filter columns

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Labels** | A | A | T | A | T | A | T | A | C | G | G | G | T |
| 118 | . | C | . | . | C | . | . | . | . | . | . | A | . |
| 123 | | | | | | | | | | | | | |
| 202 | . | C | . | . | C | . | . | . | . | . | . | . | . |
| 211 | G | . | C | . | . | . | . | . | . | . | . | A | . |
| 501 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 705 | . | . | . | . | C | T | . | . | . | . | T | . | . |
| 197 | G | . | C | . | . | . | . | . | . | . | . | . | . |
| 248 | . | . | . | . | C | T | . | . | . | . | . | . | C |
| 253 | | | | | | | | | | | | | |
| 418 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 583 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| | G | . | C | . | . | . | . | . | . | . | . | A | . |

# Filter columns: selecting the points



| Labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|--------|---|---|---|---|---|---|---|---|---|---|----|----|----|
|        | A | A | T | A | T | A | T | A | C | G | G  | G  | T  |
| 118    | . | C | . | . | C | . | . | . | . | . | .  | A  | .  |
| 123    |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 202    | . | C | . | . | C | . | . | . | . | . | .  | .  | .  |
| 211    | G | . | C | . | . | . | . | . | . | . | .  | A  | .  |
| 501    | . | . | C | G | . | . | . | . | . | . | .  | .  | .  |
| 705    | . | . | . | . | C | T | . | . | . | . | T  | .  | .  |
| 197    | G | . | C | . | . | . | . | . | . | . | .  | .  | .  |
| 248    | . | . | . | . | C | T | . | . | . | . | .  | .  | C  |
| 253    |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 418    | . | . | . | . | C | . | . | . | . | . | T  | .  | .  |
| 583    | . | . | C | G | . | . | . | . | . | . | .  | .  | .  |
|        | G | . | C | . | . | . | . | . | . | . | .  | A  | .  |

# Filter columns: selecting the points

| Labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | A | T | A | T | A | T | A | C | G | G | G | T |
| 118 | . | C | . | . | C | . | . | . | . | . | . | A | . |
| 123 | . | C | . | . | C | . | . | . | . | . | . | . | . |
| 202 | G | . | C | . | . | . | . | . | . | . | . | A | . |
| 211 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 501 | . | . | . | . | C | T | . | . | . | . | T | . | . |
| 705 | G | . | C | . | . | . | . | . | . | . | . | . | . |
| 197 | . | . | . | . | C | T | . | . | . | . | . | . | C |
| 248 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 253 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 418 | G | . | C | . | . | . | . | . | . | . | . | A | . |

# Hide/Unhide columns

| | 0 | 1 | 11 | 12 |
|---|---|---|---|---|
| Labels | A | A | G | T |
| 118 | . | C | A | . |
| 123 | . | C | . | . |
| 202 | G | . | A | . |
| 211 | . | . | . | . |
| 501 | . | . | . | . |
| 705 | . | . | . | . |
| 197 | G | . | . | . |
| 248 | . | . | . | C |
| 253 | . | . | . | . |
| 418 | . | . | . | . |
| 583 | G | . | A | . |
| | | | | |

Note:
The transition needs to be smooth

# Redundant columns

| Labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | A | T | A | T | A | T | A | C | G | G | G | T |
| 118 | . | C | . | . | C | . | . | . | . | . | . | A | . |
| 123 | | | | | | | | | | | | | |
| 202 | . | C | . | . | C | . | . | . | . | . | . | . | . |
| 211 | G | . | C | . | . | . | . | . | . | . | . | A | . |
| 501 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 705 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 197 | G | . | C | . | . | . | . | . | . | . | . | . | . |
| 248 | . | . | . | . | C | . | . | . | . | . | . | . | C |
| 253 | | | | | | | | | | | | | |
| 418 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 583 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| | G | . | C | . | . | . | . | . | . | . | . | A | . |

# Redundant columns

| Labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|--------|---|---|---|---|---|---|---|---|---|---|----|----|----|
|        | A | A | T | A | T | A | T | A | C | G | G | G | T |
| 118 | . | C | . | . | C | . | . | . | . | . | . | A | . |
| 123 |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 202 | . | C | . | . | C | . | . | . | . | . | . | . | . |
| 211 | G | . | C | . | . | . | . | . | . | . | . | A | . |
| 501 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 705 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 197 | G | . | C | . | . | . | . | . | . | . | . | . | . |
| 248 | . | . | . | . | C | . | . | . | . | . | . | . | C |
| 253 |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 418 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 583 | . | . | C | G | . | . | . | . | . | . | . | . | . |
|     | G | . | C | . | . | . | . | . | . | . | . | A | . |

# Filtering columns

Filter columns

0   1   2

| Labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | A | T | A | T | A | T | A | C | G | G | G | T |
| 118 | . | C | . | . | C | . | . | . | . | . | . | A | . |
| 123 | . | C | . | . | C | . | . | . | . | . | . | . | . |
| 202 | G | . | C | . | . | . | . | . | . | . | . | A | . |
| 211 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 501 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 705 | . | . | . | . | . | . | . | . | . | . | . | . | . |
| 197 | G | . | C | . | . | . | . | . | . | . | . | . | . |
| 248 | . | . | . | . | C | . | . | . | . | . | . | . | C |
| 253 | . | . | . | . | C | . | . | . | . | . | T | . | . |
| 418 | . | . | C | G | . | . | . | . | . | . | . | . | . |
| 583 | G | . | C | . | . | . | . | . | . | . | . | A | . |
| | . | . | . | . | C | T | . | . | . | . | . | . | . |

# Filtering columns

| Labels | 0 | 1 | 2 | 3 | 4 | 5 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|
|  | A | A | T | A | T | A | G | G | T |
| 118 | . | C | . | . | C | . | . | A | . |
| 123 |  |  |  |  |  |  |  |  |  |
| 202 | . | C | . | . | C | . | . | . | . |
| 211 | G | . | C | . | . | . | . | A | . |
| 501 | . | . | C | G | . | . | . | . | . |
| 705 | . | . | . | . | C | . | T | . | . |
| 197 |  |  |  |  |  |  |  |  |  |
| 248 | G | . | C | . | . | . | . | . | . |
| 253 | . | . | . | . | C | . | . | . | C |
| 418 | . | . | . | . | C | . | T | . | . |
| 583 | . | . | C | G | . | . | . | . | . |
|  | G | . | C | . | . | . | . | A | . |
|  | . | . | . | . | C | T | . | . | . |

# Filtering columns

Filter columns

0 1 2

| Labels | 0 | 1 | 2 | 3 | 4 | 10 | 11 |
|---|---|---|---|---|---|---|---|
|  | A | A | T | A | T | G | G |
| 118 | . | C | . | . | C | . | A |
| 123 |  |  |  |  |  |  |  |
| 202 | . | C | . | . | C | . | . |
| 211 | G | . | C | . | . | . | A |
| 501 | . | . | C | G | . | . | . |
| 705 | . | . | . | . | C | T | . |
| 197 |  |  |  |  |  |  |  |
| 248 | G | . | C | . | . | . | . |
| 253 | . | . | . | . | C | . | . |
| 418 | . | . | . | . | C | T | . |
| 583 | . | . | C | G | . | . | . |
|  | G | . | C | . | . | . | A |
|  | . | . | . | . | C | . | . |

# Task abstraction

- Find columns that their substitution lead to the most severe rows. Report them individually or in a combination with other columns:

1. Sort the rows ✔

2. Filter the columns: Basic ✔

   2.1 Find interesting columns

   2.2 Find related columns

      ( and filter the rest )
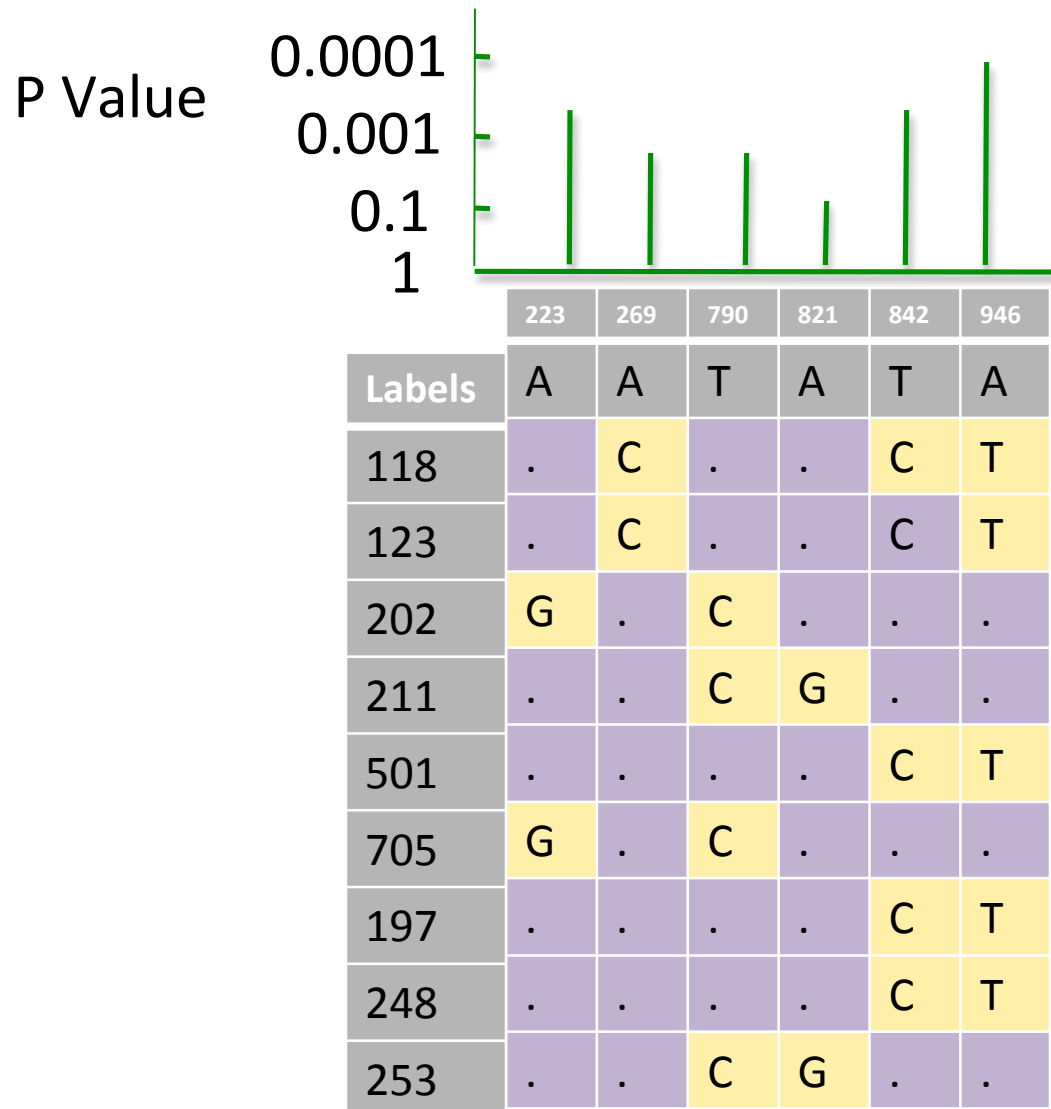
# A pattern within a column

|        | 12 |
|--------|-----|
| **Labels** | A |
| 118 | C |
| 123 | C |
| 202 | . |
| 211 | . |
| 501 | C |
| 705 | . |
| 197 | . |
| 248 | . |
| 253 | . |
| 418 | . |
| 583 | . |

Severe

Moderate
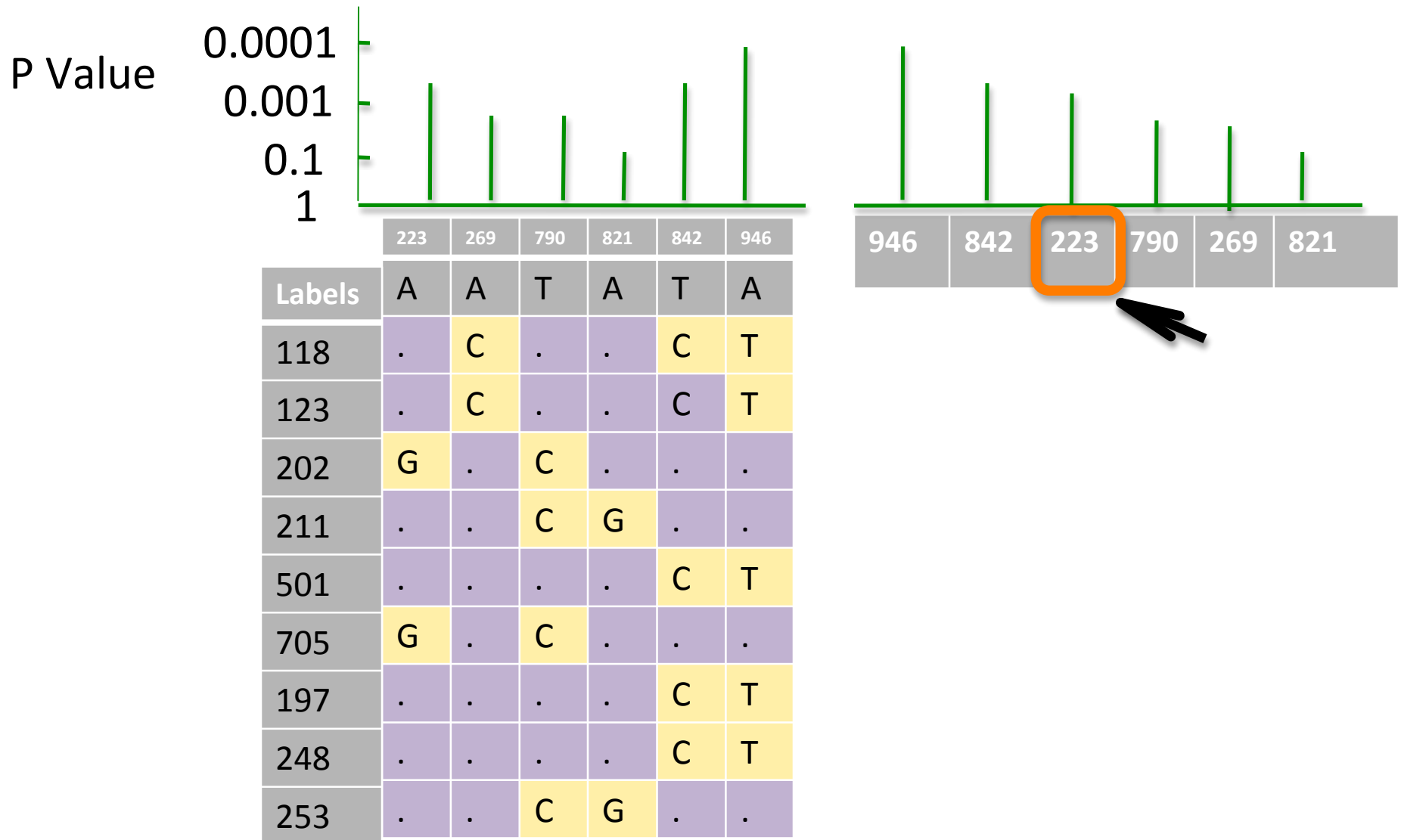
Mild

# Detecting the pattern:

- Challenging for humans
- Derived data for each column:

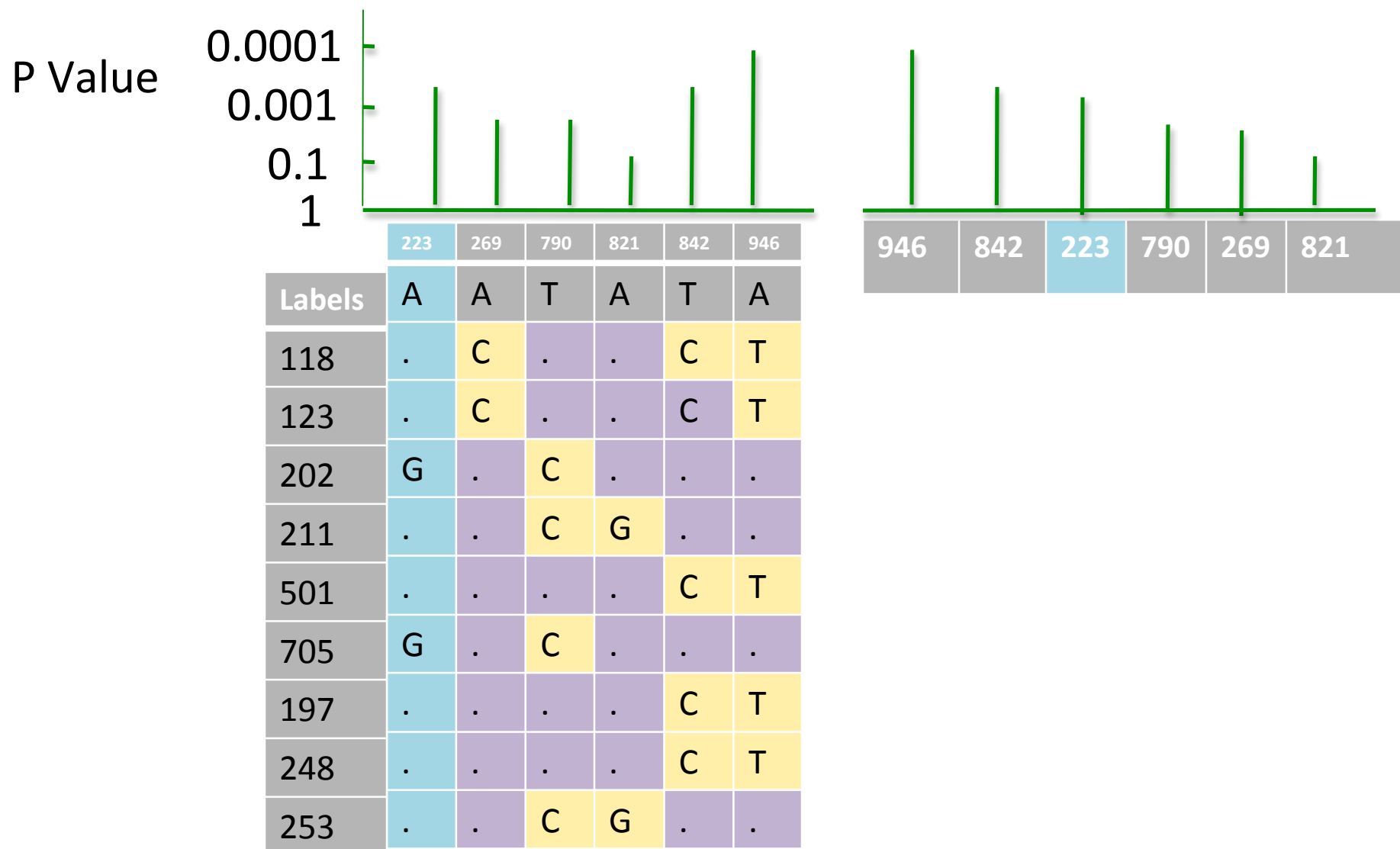Statistical significance score in the Mann-Whitney U test (the P-value)
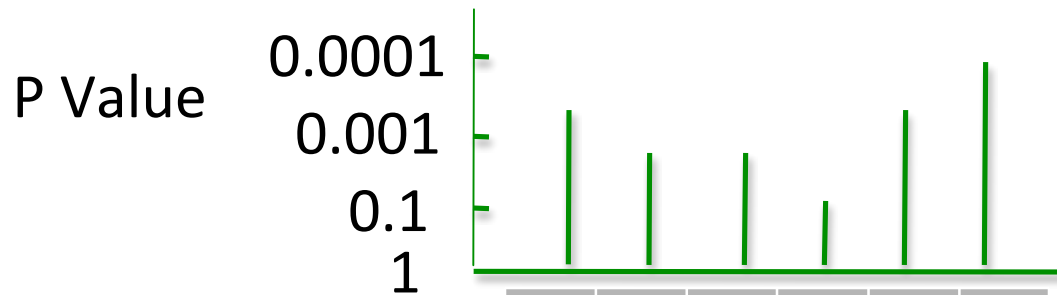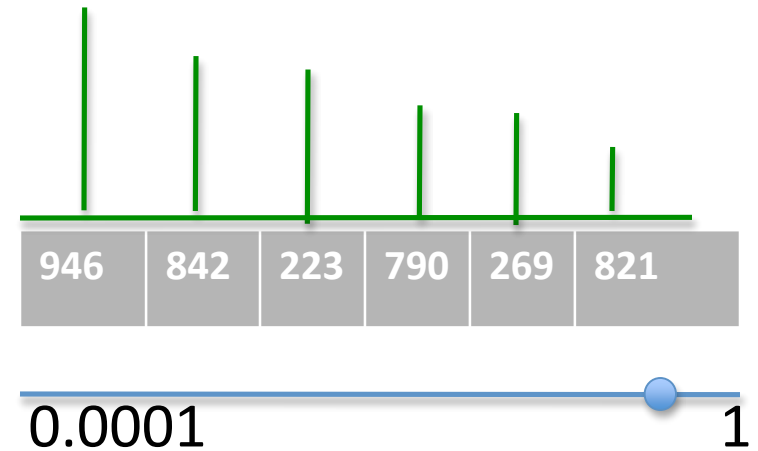
# Visual encoding of the P-value



P Value

| P Value | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|
| **Labels** | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |

# Sorting/Filtering columns by the P-value

P Value



| Labels | 223 | 269 | 790 | 821 | 842 | 946 |
|--------|-----|-----|-----|-----|-----|-----|
|        | A   | A   | T   | A   | T   | A   |
| 118    | .   | C   | .   | .   | C   | T   |
| 123    | .   | C   | .   | .   | C   | T   |
| 202    | G   | .   | C   | .   | .   | .   |
| 211    | .   | .   | C   | G   | .   | .   |
| 501    | .   | .   | .   | .   | C   | T   |
| 705    | G   | .   | C   | .   | .   | .   |
| 197    | .   | .   | .   | .   | C   | T   |
| 248    | .   | .   | .   | .   | C   | T   |
| 253    | .   | .   | C   | G   | .   | .   |

| 946 | 842 | 223 | 790 | 269 | 821 |
|-----|-----|-----|-----|-----|-----|

# Sorting/Filtering columns by the P-value

P Value



| Labels | 223 | 269 | 790 | 821 | 842 | 946 |
|--------|-----|-----|-----|-----|-----|-----|
| Labels | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |

| 946 | 842 | 223 | 790 | 269 | 821 |
|-----|-----|-----|-----|-----|-----|

# Sorting/Filtering columns by the P-value



P Value

| | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|
| Labels | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |

| 946 | 842 | 223 | 790 | 269 | 821 |
|---|---|---|---|---|---|

0.0001           1

# Sorting/Filtering columns by the P-value

P Value

| | 223 | 269 | 790 | 842 | 946 |
|---|---|---|---|---|---|
| **Labels** | A | A | T | T | A |
| 118 | . | C | . | C | T |
| 123 | . | C | . | C | T |
| 202 | G | . | C | . | . |
| 211 | . | . | C | . | . |
| 501 | . | . | . | C | T |
| 705 | G | . | C | . | . |
| 197 | . | . | . | C | T |
| 248 | . | . | . | C | T |
| 253 | . | . | C | . | . |

| 946 | 842 | 223 | 790 | 269 |
|---|---|---|---|---|

0.0001                  1

0.0001
0.001
0.1
1

# Task abstraction

- Find columns that their substitution lead to the most severe rows. Report them individually or in a combination with other columns:

1. Sort the rows ✔

2. Filter the columns  Basic ✔

    2.1  Find interesting columns ✔

    2.2  Find related columns

       ( and filter the rest )

# Correlation pattern between the columns

| Labels | 4<br>T | 5<br>A |
|---|---|---|
| 118 | C | T |
| 123 | C | T |
| 202 | . | . |
| 211 | . | . |
| 501 | C | T |
| 705 | . | . |
| 197 | C | T |
| 248 | C | T |
| 253 | . | . |
| 418 | . | . |
| 583 | C | T |

# Complementary pattern between the columns

| Labels | 2<br>T | 5<br>A |
|--------|--------|--------|
| 118 | . | T |
| 123 | . | T |
| 202 | C | . |
| 211 | C | . |
| 501 | . | T |
| 705 | C | . |
| 197 | . | T |
| 248 | . | T |
| 253 | C | . |
| 418 | C | . |
| 583 | . | T |

# Derived attribute between any pair of columns

- Correlation: shows how much two columns are dependant

$$\rho_{X,Y} = \mathrm{corr}(X,Y) = \frac{\mathrm{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y},$$

- A relationship between any two columns
- +1: highly correlated
- -1:  highly complement

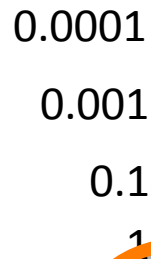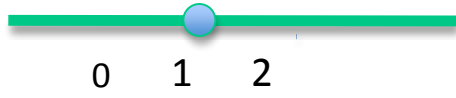# Visual encoding design of the correlation

# Filter Columns

0  1  2

| | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|
| 0.0001 | | | | | | |
| 0.001 | | | | | | |
| 0.1 | | | | | | |
| 1 | | | | | | |

| | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|

| **Labels** | A | A | T | A | T | A |
|---|---|---|---|---|---|---|
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

| New | symptom | Drug Res. |
|---|---|---|

946
842
269
821
790
223

Filter Columns

0  1  2

0.0001
0.001
0.1
1

| | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|
| **Labels** | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

| 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|

| New | symptom | Drug Res. |
|---|---|---|

946
842
269
821
790
223

Filter Columns

0  1  2

| | New | symptom | Drug Res. |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

0.0001
0.001
0.1
1

| Labels | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|
| | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

| 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|

946  842  269  821  790  223

## Filter Columns

0  1  2

0.0001
0.001
0.1
1

| 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|

| 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|

| New | symptom | Drug Res. | | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|---|---|---|
| | | | Labels | A | A | T | A | T | A |
| | | | 118 | . | C | . | . | C | T |
| | | | 123 | . | C | . | . | C | T |
| | | | 202 | G | . | C | . | . | . |
| | | | 211 | . | . | C | G | . | . |
| | | | 501 | . | . | . | . | C | T |
| | | | 705 | G | . | C | . | . | . |
| | | | 197 | . | . | . | . | C | T |
| | | | 248 | . | . | . | . | C | T |
| | | | 253 | . | . | C | G | . | . |
| | | | 418 | G | . | C | . | . | . |
| | | | 583 | . | . | . | . | C | T |

946    842    269    821    790    223

Filter Columns

0    1    2

| New | symptom | Drug Res. |
|---|---|---|

0.0001
0.001
0.1
1

| | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|

| | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|

| Labels | A | A | T | A | T | A |
|---|---|---|---|---|---|---|
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

946    842    269    821    790    223

# Filter Columns

0  1  2



| Labels | 223 | 269 | 790 | 821 | 842 | 946 |
|---|---|---|---|---|---|---|
| | A | A | T | A | T | A |
| 118 | . | C | . | . | C | T |
| 123 | . | C | . | . | C | T |
| 202 | G | . | C | . | . | . |
| 211 | . | . | C | G | . | . |
| 501 | . | . | . | . | C | T |
| 705 | G | . | C | . | . | . |
| 197 | . | . | . | . | C | T |
| 248 | . | . | . | . | C | T |
| 253 | . | . | C | G | . | . |
| 418 | G | . | C | . | . | . |
| 583 | . | . | . | . | C | T |

| New | symptom | Drug Res. |
|---|---|---|

0.0001
0.001
0.1
1

946  842  269  821  790  223

# The next step

- Finalizing the ideas
- Implementation with Processing or Protovis
- Qualitatively testing on the real domain experts

# Question and Feedbacks