



THRESHOLD - A novel gene saturation analysis GUI

Authors

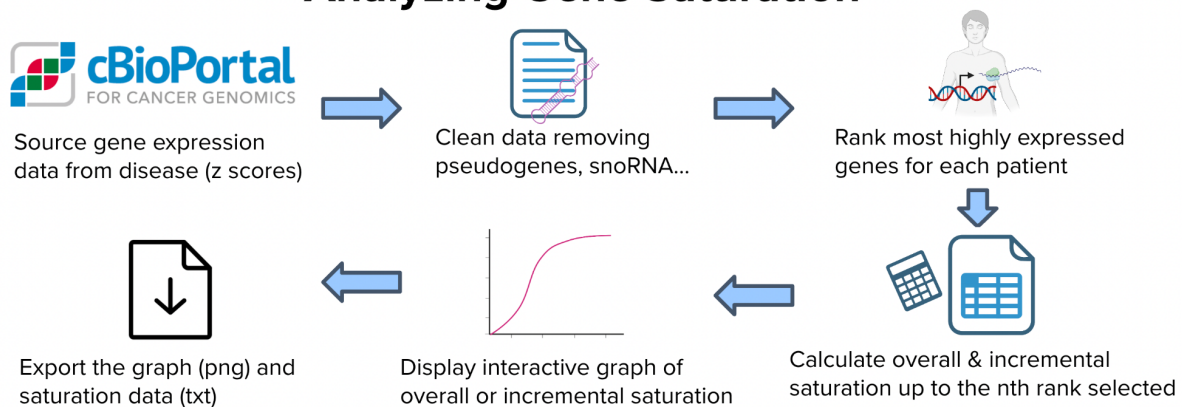
This project was developed by Finán Gammel, Jennifer Li and Dr. Alper Uzun in the Uzun Lab at Brown University.

Purpose

THRESHOLD analyzes transcriptomic data across large samples of patients to understand the cohesion of the most upregulated/downregulated genes in a given disease. This lends researchers knowledge to inform network topology analyses or to assess the relative amount of genes influencing a given disease. **THRESHOLD** offers several features to aid in analysis including user-inputted saturation type, restriction factors, and rank type. The tool outputs an interactive graph of saturation permitting the calculation of specific saturation thresholds and most saturated genes.

Overall Workflow

Analyzing Gene Saturation



- **Gene saturation:** measures the number of shared highly expressed genes in a population
 - **Overall:** total saturation up to and including nth level
 - **Incremental:** new saturation at nth level



Created with BioRender.com

**Alternatively, genes may be ranked from low to high.*

Overall Saturation - The saturation of all the genes up to and including the nth level. Quantifies the count of genes up to the nth level that meet or exceed the inputted restriction level, divided by the total count of genes up to and including the nth level.

Formula:

Overall Saturation =

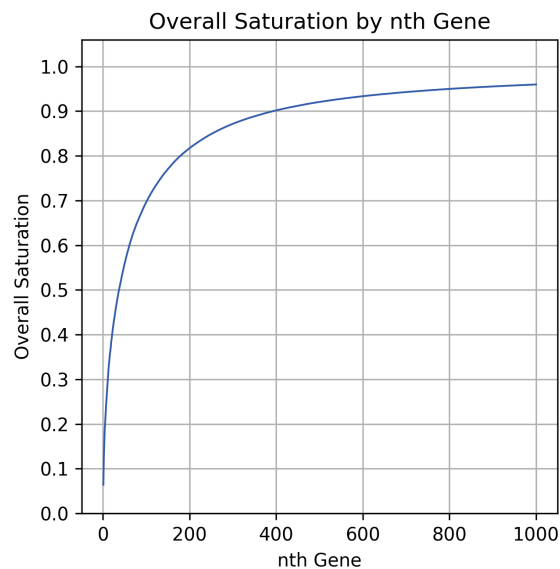
$$\frac{\sum_{e \in \{y | y \in A\}} \max(|\{y \in A | y = e\}| - (x - 1), 0)}{|A|}$$

A = the multiset of genes up to and including the nth level

x = the restriction level

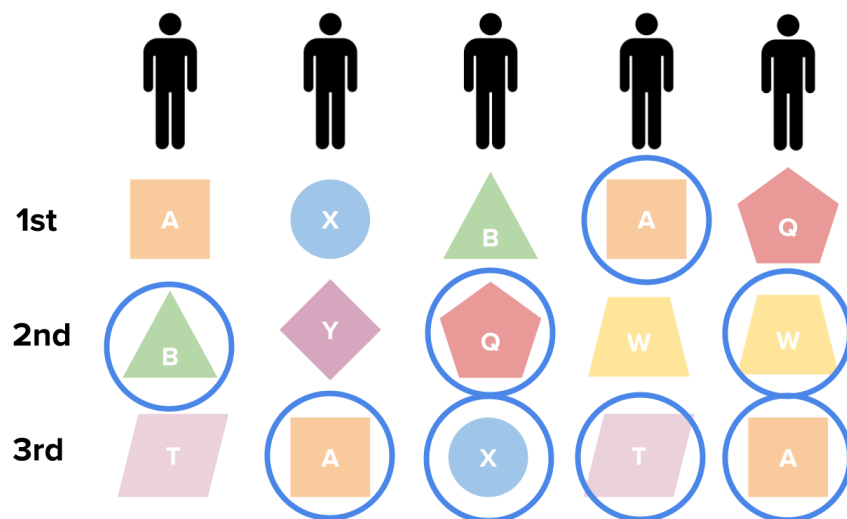
Conceptual Idea - The sum of each unique element in the multiset A's non-negative greater value between the element count in multiset A minus one less than the restriction level divided by the number of elements in the multiset A.

Graphical Output Example:



Visual Example:

Find the **overall saturation** up to the **3rd level** ($n = 3$) given a **restriction level** of **2**:



Saturated Genes; occur ≥ 2 nd time
(restriction level = 2)

8 genes saturated / 15 genes
up to and including nth level
= ~53% overall saturation

| nth level | Patient 1 | Patient 2 | Patient 3 | Patient 4 | Patient 5 |
|-----------|-----------|-----------|-----------|-----------|-----------|
| 1st | Gene A | Gene X | Gene B | Gene A | Gene Q |
| 2nd | Gene B | Gene Y | Gene Q | Gene W | Gene W |
| 3rd | Gene T | Gene A | Gene X | Gene T | Gene A |

- All genes present up to and including $n = \{A, X, B, A, Q, B, Y, Q, W, W, T, A, X, T, A\}$
 - Occurrences by element: $A = 4, X = 2, B = 2, Q = 2, Y = 1, W = 2, T = 2$
- Saturated genes (subtracted one minus the restriction level):

$A: 4 - (2-1) = 3$

$B: 2 - (2-1) = 1$

$Y: 1 - (2-1) = 0$

$T: 2 - (2-1) = 1$

$X: 2 - (2-1) = 1$

$Q: 2 - (2-1) = 1$

$W: 2 - (2-1) = 1$
- Occurrences ≥ 2 (restriction level): $3 + 1 + 1 + 1 + 0 + 1 + 1 = 8$
- 8 saturated genes / 15 genes up to and including $n = \text{~53\% overall saturation}$

Incremental Saturation - The saturation only at the incremental nth gene level.

Quantifies the count of genes at the nth level that meet or exceed the inputted restriction level up to and including that level, divided by the total count of genes at the nth level.

Formula:

Incremental Saturation =

$$\frac{\sum_{e \in \{y | y \in B\}} \max(\min(|\{y \in A | y = e\}| - (x - 1), |\{y \in B | y = e\}|), 0)}{|B|}$$

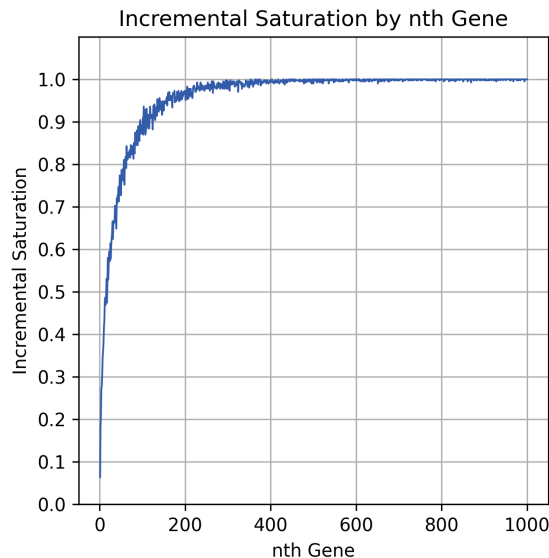
A = the multiset of genes up to and including the nth level

B = the multiset of genes at the nth level

x = restriction level

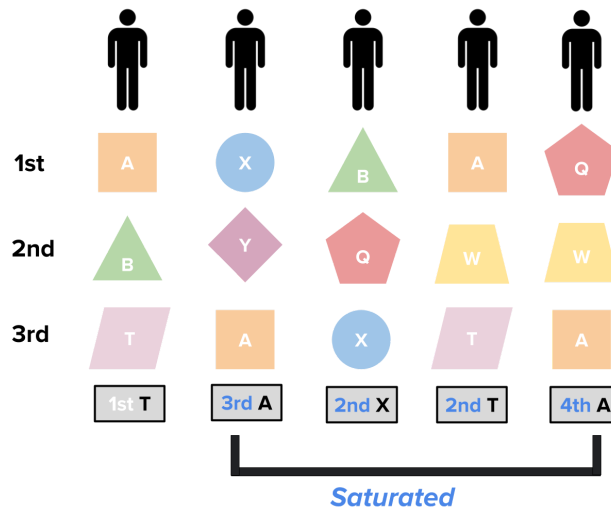
Conceptual Idea - The sum of each unique element in the multiset B's non-negative lesser value between the difference in the element count in multiset A minus one less than the restriction level and the element count in multiset B divided by the number of elements in multiset B.

Graphical Output Example:



Visual Example:

Find the **incremental saturation** of the **3rd level** ($n = 3$) given a **restriction level** of **2**:



At $n = 3$ (third level) there are 4 saturated genes out of 5 genes total;
 $4 / 5 = 80\%$ incremental saturation

| nth level | Patient 1 | Patient 2 | Patient 3 | Patient 4 | Patient 5 |
|-----------|-----------|-----------|-----------|-----------|-----------|
| 1st | Gene A | Gene X | Gene B | Gene A | Gene Q |
| 2nd | Gene B | Gene Y | Gene Q | Gene W | Gene W |
| 3rd | Gene T | Gene A | Gene X | Gene T | Gene A |

- Genes present up to and including n ; $A = \{A, X, B, A, Q, B, Y, Q, W, W, T, A, X, T, A\}$
 - Occurrences by element: $A = 4, X = 2, B = 2, Q = 2, Y = 1, W = 2, T = 2$
- Genes present at n ; $B = \{T, A, X, T, A\}$
 - Occurrences by element: $T = 2, A = 2, X = 1$
- Saturated genes (lesser between count in multiset B minus one less than the restriction level and count in multiset A ; must be non-negative)
 - $T: 2 - (2-1) = 1 \dots 2 \text{ in } B \dots 1 < 2 \rightarrow 1$
 - $X: 2 - (2-1) = 1 \dots 1 \text{ in } B \dots 1 = 1 \rightarrow 1$
 - $A: 4 - (2-1) = 3 \dots 2 \text{ in } A \dots 2 < 3 \rightarrow 2$
- Occurrences ≥ 2 (restriction level): $1 + 1 + 2 = 4$
- 4 saturated genes / 5 genes at $n = 80\%$ incremental saturation

Standardization

To compare saturation across datasets, restriction levels can be standardized.

- Increasing the occurrence level by a factor of the data set size increase and the restriction level accordingly produces the same saturation results and can be used for standardization.
- For example, if a data set of 100 patients has a restriction level of 3 genes, the occurrence level is 2 (occurrence level = restriction level - 1). To standardize a data set of 200 patients (twice the size), you would need to double the occurrence level $2 \times 2 = 4$. This corresponds to a restriction level of 5.
- This relationship can be summarized by the following equation (in genes):
$$\text{new restr. level} = \text{initial restr. level} \times \text{factor increase} - (\text{factor increase} - 1)$$

As a result, a restriction level of x% does not necessarily = the x% restriction level of another data set. You must use the above gene number conversion to achieve standardization.

Dataset Validation:

The following testing exists to verify the calculations outputted by **THRESHOLD**. From our testing, the tool has demonstrated consistent accuracy.

Saturation Validation Testing

Given the following test dataset:

| Hugo_Symbol | Entrez_Gene_Id | Pat1 | Pat2 | Pat3 | Pat4 |
|-------------|----------------|------|------|------|------|
| A | - | 1 | 5 | 1 | 5 |
| B | - | -3 | 3 | 5 | 0 |
| X | - | 5 | 0 | -1 | -2 |
| W | - | -1 | -1 | -2 | 3 |
| Q | - | -2 | -3 | 3 | 1 |
| Y | - | 3 | -2 | -3 | -3 |
| T | - | -4 | 1 | 0 | -1 |

THRESHOLD passed verified calculations all of the following tests:

- | | |
|-------------------------------------|-------------------------------------|
| - Regular rank (level = 2, n = 5) | - Reverse rank (level = 2, n = 5) |
| - Regular rank (level = 3, n = 5) | - Reverse rank (level = 3, n = 5) |
| - Regular rank (level = 4, n = 5) | - Reverse rank (level = 4, n = 5) |
| - Regular rank (level = 1% n = 5) | - Reverse rank (level = 1% n = 5) |
| - Regular rank (level = 80%, n = 5) | - Reverse rank (level = 80%, n = 5) |

Usage:

Download

To install **THRESHOLD** simply click the installation link below.

LINK

Once you have downloaded the above folder, open the entire folder in a VScode IDE. Ensure the Python and Java extensions are added. When ready, simply navigate to the threshold.gui file and press run to begin. A window should appear, prompting the input of a file.

Input File Format

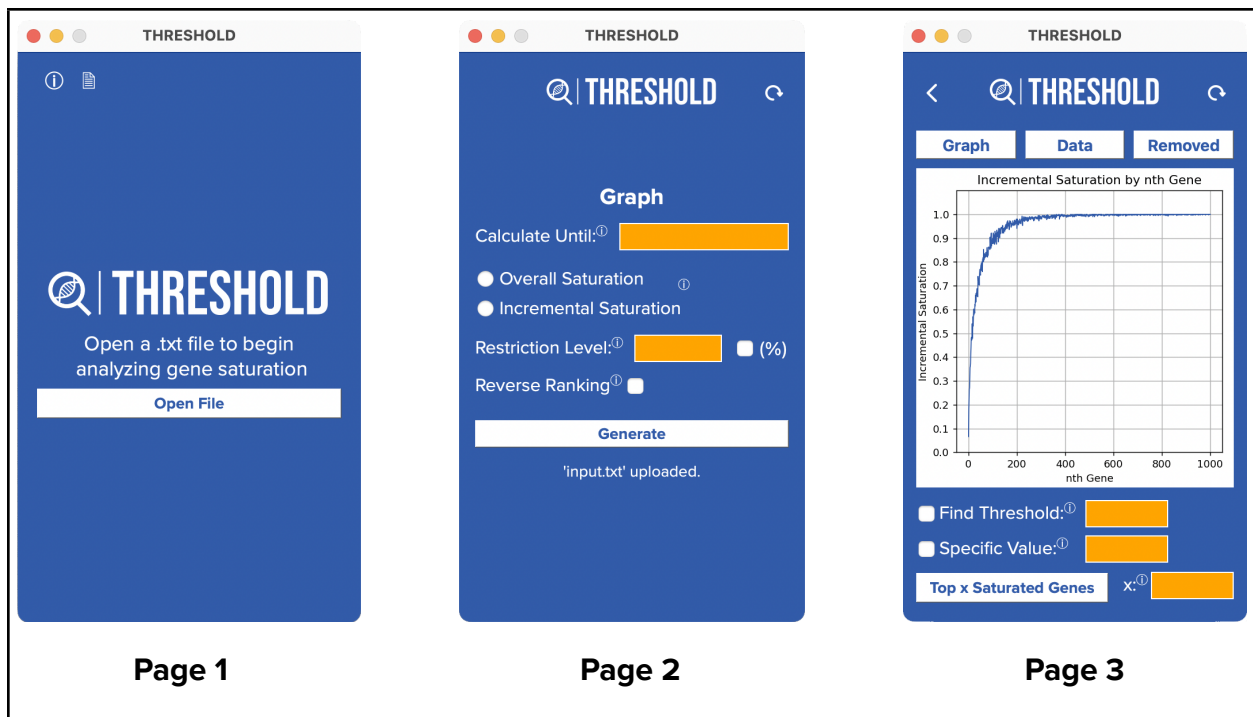
THRESHOLD requires the input of a file of patient transcriptomic data with gene expression zscores ideally compared with a control population.

The inputted .txt file must be formatted as such:

| Hugo_Symbol | Entrez_Gene_Id | {Patient1ID} | {Patient2ID}... |
|-------------|----------------|--------------|-----------------|
| ... | ... | zscores... | zscores... |
| ... | ... | zscores... | zscores... |

- The Entrez_Gene_Id column is never actually used; it is a placeholder. As long as there's a space between the Hugo_Symbol and {Patient1ID} column it will run.
- Note the first heading **MUST** be called Hugo_Symbol
- The columns must be separated by a tab '\t'

GUI Navigation:



Page 1

- To begin the analysis, simply upload a .txt file in the appropriate file format by clicking on the Open File button. Wait about 5 seconds while the pseudogenes are removed, and you will automatically move to the next page.
- To view info and documentation of the tool, press the appropriate icons in the GUI toolbar

Page 2

- Confirm your file is uploaded with the dialog below the Generate button.
- Input the parameters according to your desired analyses. Clicking on the 'i' icons will provide a detailed explanation of the parameter's domains and purpose.
- Once satisfied with your inputs, press Generate. If the inputted file is formatted correctly, and the inputs accurate, an In progress... text dialog will appear. For a

dataset of about 400 patients, you will need to wait about 2 minutes for the calculations to be made.

- If you would like to restart your analyses with a different file, press the restart button in the top right corner

Page 3

- A saturation curve will appear indicating saturation type by nth ranked gene. See documentation for methods.
- Hover over the graph to find specific saturation values. Alternatively, use the Find Threshold input to find when a certain saturation level is reached and the Specific Value input to evaluate the saturation of a specific nth gene. Clicking on the ⓘ icons will provide a detailed explanation of the parameter's domains and purpose.
- Additionally, you can find the most saturated genes with the Top x Saturated Genes button. Simply enter the number of top genes you want to find in the x: parameter and press the button.
- To export a .png file of the Graph, a .txt file of the Data or a .txt file of removed pseudogenes, press the corresponding buttons above the graph.
- To begin a new analysis with the same file, press the back button. If you would like to restart your analysis entirely, press the restart button.