# Multiple States in Cancer Cell Population and Transplantation Experiment Inference

Yue Wang

Institut des Hautes Études Scientifiques (IHÉS), France

Mar. 10, 2021

- Two mathematical biology projects: Given certain biological data, what biology can we learn, and what mathematics can we work on?
- Cancer cell population: When there are enough data, what patterns can we reveal? Explain biological phenomena with ODEs and stochastic processes.
- Tissue transplantation: When there are not enough data, can we infer unknown data? Turn experimental design into combinatorial problems.
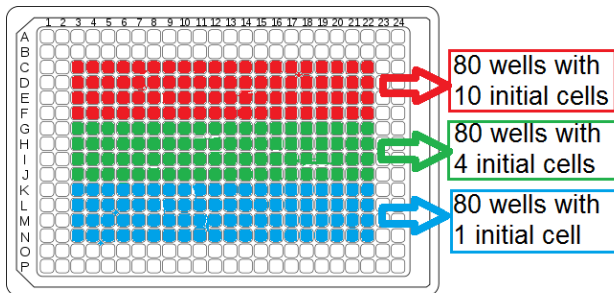
# Multiple States in Cancer Cell Population

- Cancer cell population is often thought to be homogeneous.
- Analyze experimental data to reveal the existence of multiple cell states.
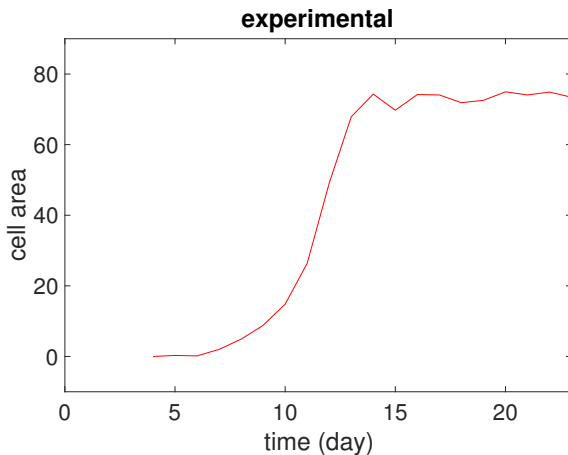- Theoretical explanations of related new phenomena.

- Cultivate HL60 leukemia cells *in vitro*.



- 80 wells with 10 initial cells
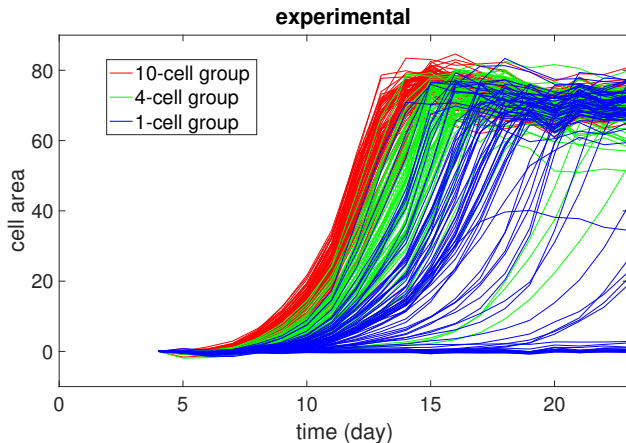- 80 wells with 4 initial cells
- 80 wells with 1 initial cell

- Initial cells are sampled randomly from a large population.
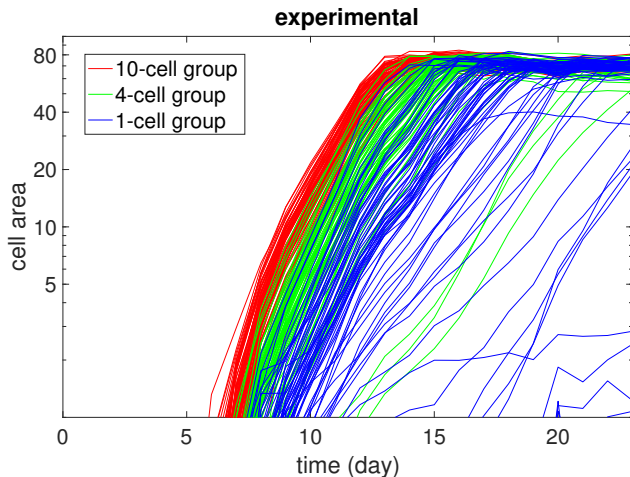- For each well, the cell area (proportional to cell number) is measured everyday.

**experimental**

Growth curve of one well, describing how the population changes along time. In general, the population grows exponentially until saturation.

Each growth curve corresponds to one well. Red: 10 initial cells; green: 4 initial cells; blue: 1 initial cell.

Growth curves with *y*-axis in log scale. Red: 10 initial cells; green: 4 initial cells; blue: 1 initial cell.
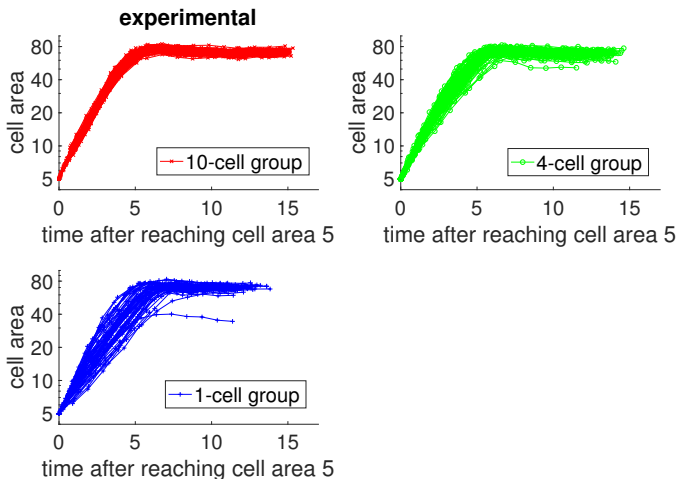
Figure: Translated population curves, starting from cell area 5. *Y*-axis is in log scale. Some 1-cell-wells never reach cell area 5, thus are not shown.

# Section I.1: Growth rates

For one well, denote the population at day $n$ as $c_n$, and the population at day $n + 1$ as $c_{n+1}$. Then the growth rate is $g_n = (c_{n+1} - c_n)/c_n$. For each well in each day, draw the growth rate $g_n$ versus the population $c_n$. The point cloud near $(75, 0)$ corresponds to saturated wells.

# Section I.1: Experimental phenomena

- After reaching the same population, all 10-cell-wells grow fast; some 1-cell-wells grow much slower.
- Some 1-cell-wells keep at low population levels for a long time.
- When a 1-cell-well grows to have 10 cells, it is different from a 10-cell-well.
- Cells cannot be homogeneous.

# Section I.1: Analysis

We assume that there are at least three cell states with different growth rates: fast, moderate, and slow.

- Build a multi-type branching process model.
- Initial cells have three possible states, determining the growth rate. Growth rate is inheritable, and decreases as total population increases.
- For each time period, each cell has a probability to divide, and a probability to die.
- In simulation, this model can reproduce most experimental phenomena within a wide range of parameters.
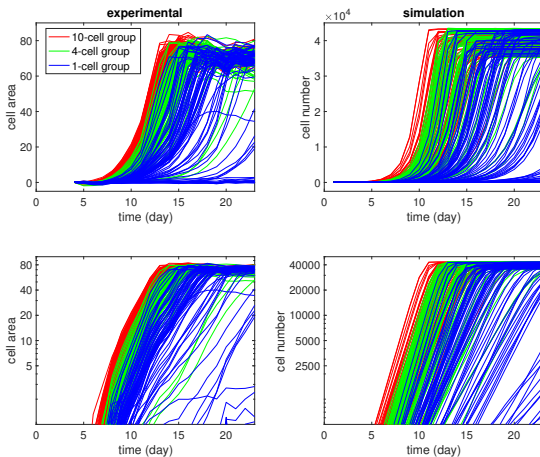
Figure: Population growth curves.

Figure: Translated population curves, starting from cell area 5 or cell number 2500. *Y*-axis is in log scale.

Figure: Growth rate versus population size.

- Experimental data reveal the existence of multiple states in cancer cell population.
- Corresponding model can reproduce experimental phenomena.

- Questions?

# Section I.2: Multiple states

- The existence of multiple states has been verified in some other cancers.
- SUM159 breast cancer cell population has three states: stem, luminal, basal (distinguished by cell-surface markers).
- Why could multiple states (possibly with different growth rates) survive simultaneously?
- There exist epigenetic transitions between different states. Such change of state is inheritable.

# Section I.2: Multiple states

- Starting from any one state, other states will emerge, and the population gradually recovers the equilibrium proportions.
- It is called the "state equilibrium phenomenon".

How to explain such state equilibrium phenomenon?

## Section I.2: Deterministic model

- Cells can divide, die or transform into other states. Assume cells do not interact, and there is no carrying capacity.
- The population vector $\vec{x}$ of different states satisfies a linear ODE system:

$$\mathrm{d}\vec{x}/\mathrm{d}t = \vec{x}\mathbf{A},$$

where $\mathbf{A} = \{a_{ij}\}$, the matrix of transition rates.

- The population proportion vector $\vec{w} = \vec{x}/||\vec{x}||_1$ satisfies a quadratic system:

$$\frac{\mathrm{d}\vec{w}}{\mathrm{d}t} = \vec{w}[\mathbf{A} - (\vec{w}\vec{b}')\mathbf{I}],$$

where $\vec{b} = \vec{1}\mathbf{A}'$.

Perron-Frobenius Theorem states that **A** has a real eigenvalue $\lambda_1$, which is larger than the real parts of any other eigenvalues. Its normalized eigenvector is denoted by $\vec{u}_1$.

### Theorem

*If $\lambda_1$ is a simple root of the characteristic polynomial (in reality, this holds in general), the system $\mathrm{d}\vec{w}/\mathrm{d}t = \vec{w}[\boldsymbol{A} - (\vec{w}\vec{b}')\boldsymbol{I}]$ has a unique stationary fixed point $\vec{u}_1$.*

Therefore the proportion vector $\vec{w}$ always converges to $\vec{u}_1$.

## Section I.2: Stochastic model

- We can describe this population with a branching process.
- One cell of state $i$, $Y_i$, can branch into a (stochastic) combination of cells with different states:
  $Y_i \overset{\alpha_i}{\to} d_{i1} Y_1 + d_{i2} Y_2 + \cdots + d_{in} Y_n$. The waiting time is exponential with rate $\alpha_i$.
- Here $d_{ij}$ are random variables. For example,
  $d_{11} = 2$, $d_{12} = 0$ means division $Y_1 \to 2 Y_1$;
  $d_{11} = d_{12} = 0$ means death $Y_1 \to \emptyset$;
  $d_{11} = 0$, $d_{12} = 1$ means transition $Y_1 \to Y_2$.
- If we take expectations for population, the branching process model returns to the ODE model.

## Section I.2: Stochastic model

- Due to stochasticity, it is possible that all cells die out, and the proportions cannot be defined.
- We focus on the stochastic trajectories that no state dies out forever (called "non-extinction").
- If $\lambda_1 > 0$, as the initial cell number increases, the probability of non-extinction tends to 1.

### Theorem

*Assume that $\lambda_1 > 0$ and $\lambda_1$ is a simple root of the characteristic polynomial. Conditioned on non-extinction, the proportion vector $\vec{w}$ converges to $\vec{u}_1$ with probability 1.*

- This is a strong law of large numbers for branching processes. It improves a result by Svante Janson in 2004.
- It provides a stochastic explanation for the state equilibrium phenomenon.

- We have explained the state equilibrium phenomenon in ODE model and branching process model.

- Questions?

# Tissue Transplantation Experiments: Inference and Experimental Design

## Section II: Outline

- Tissue transplantation experiments are important in developmental biology. However, most experimental results are unknown.
- Penalty function-based method to infer unknown experimental results.
- How to design experiments, so that the inference method can be applied most efficiently?

Donor    Host

- Tissue transplantation experiments: For an embryo, excise a piece of one tissue (donor tissue), and transplant it to another tissue (host tissue).
- E.g., the transplantation experiment with donor tissue D and host tissue H is denoted as {D,H}.
- The transplanted tissue is placed in an unnatural environment. Therefore, its development might be normal (N) or abnormal (A).

- Developmental biology: Why could a zygote (in natural environment) develop into an adult animal?
- To understand why the developmental process in natural environment works, we also need to understand why the developmental process in unnatural environment does not always work.
- Tissue transplantation experiments describe how tissues behave in unnatural environments.

|      |      | | | Donor | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
|      | | AM19 | PM19 | PM15 | UL11 | LL11 | LL15 | LL19 |
| Host | AM19 | ? | N | A | A | A | A | N |
|      | PM19 | ? | N | ? | N | N | ? | ? |
|      | PM15 | ? | ? | ? | ? | ? | ? | ? |
|      | UL11 | ? | N | ? | N | N | ? | ? |
|      | LL11 | ? | N | ? | N | N | ? | ? |
|      | LL15 | ? | ? | ? | ? | ? | ? | ? |
|      | LL19 | ? | ? | ? | ? | ? | ? | ? |

Table: Results for *Xenopus laevis*, reported by Krneta-Stankic et al. 2010

N=normal; A=abnormal; ?=unknown. AM=anterior paraxial mesoderm; PM=presomitic mesoderm; UL=upper lateral lip; LL=lower lip; Number=developmental stage.

- There are many possible tissue transplantation experiments. Only a small portion has been conducted. We need a method to infer the unknown results.
- Core idea: Similar experiments should have similar results. For similar experiments, we can use known results to infer unknown results.
- Assume we have known the similarities between experiments.

Experiment similarity chart:



$\{PM19, PM19\} = N$         $\{PM15, PM15\} = N$

$\left\{\begin{matrix} PM19 \\ PM15 \end{matrix}\right\}$

$\{AM19, PM15\} = A$         $\{AM19, PM19\} = N$

Black/red terms are experiments with known/unknown results.
Linked experiments are similar.
The result of {PM19,PM15} can be inferred by the known
results of similar experiments.

## Section II.1: Ideas

- How to determine the similarities between experiments?
- We can describe the similarities between tissues qualitatively or quantitatively. (Transcriptome information, concentration of certain molecules, distance on the developmental tree, etc.)
- With the similarities between tissues, we can establish the similarities between experiments. Experiments are similar if they have similar donor tissues and similar host tissues.
- For this project, we do not have enough data. Thus the experiment similarities are assigned subjectively and rather arbitrarily.

- We take guesses of unknown experimental results, and use a penalty function to evaluate such guesses. Then we can find the best guesses.
- There is a penalty if two similar experiments have different results.
- For the concrete form of this penalty function, we can get inspirations from the Ising model.

# Section II.1: Ising model

The Ising model describes ferromagnetism in statistical mechanics. Consider a set of lattice sites, where each site $k$ has a variable $\sigma_k$ that takes $+1$ or $-1$.

- For a configuration $\sigma$ of $\pm 1$, its energy function (no external field) is

$$\mathrm{H}(\sigma) = -\sum_{i \sim j} J_{ij}\sigma_i\sigma_j,$$

  where $i \sim j$ means site $i$ and site $j$ are neighboring, and $J_{ij} \geq 0$ is the interaction coefficient. For neighboring sites $i, j$, when $\sigma_i = \sigma_j$, the energy is lower.

- The probability of a configuration $\sigma$ is

$$\mathbb{P}_\beta(\sigma) = e^{-\beta \mathrm{H}(\sigma)}/Z_\beta,$$

  where $\beta = (k_B T)^{-1}$, $Z_\beta$ is the normalization constant.

- Configuration with lower energy (smaller penalty) has higher probability. Neighboring sites tend to have the same value.

Analogies between tissue transplantation experiments and the Ising model:

| Tissue transplantation | Ising model |
| --- | --- |
| Experiment similarity chart | Lattice |
| Experiment | Site |
| Similar experiments | Neighboring sites |
| Result: normal/abnormal | Value: +1/-1 |
| | |
| Penalty: similar experiments have different results | Penalty: neighboring sites have different values |
| Penalty function? | Energy function |

Pure analogies, not physical correspondence.

- For tissue transplantation experiments, we take guesses for unknown experimental results $\{\sigma_i\}$.
- The penalty function is

$$\mathrm{H}(\sigma) = -\sum_{i,j} J_{ij}\sigma_i\sigma_j.$$

- The probability of a configuration $\sigma$ is

$$\mathbb{P}_\beta(\sigma) = e^{-\beta \mathrm{H}(\sigma)}/Z_\beta.$$

- Regard N as $+1$, and A as $-1$.
- Use experiment similarities to determine parameter $J_{ij}$.

$$H(\sigma) = -\sum_{i,j} J_{ij}\sigma_i\sigma_j, \ \mathbb{P}_\beta(\sigma) = e^{-\beta H(\sigma)}/Z_\beta.$$

$\{PM19, PM15\} = N$:

$$H = -1 \times 1 \times 1 - 1 \times 1 \times 1 - 1 \times 1 \times (-1) - 1 \times 1 \times 1 = -2.$$

$$\mathbb{P} = e^{-0.1 \times (-2)}/Z_\beta = 0.60.$$
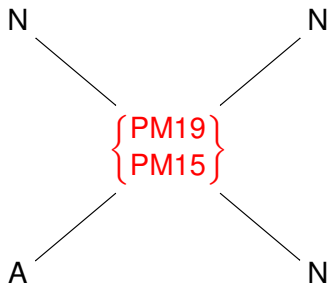
$$H(\sigma) = -\sum_{i,j} J_{ij}\sigma_i\sigma_j, \ \mathbb{P}_\beta(\sigma) = e^{-\beta H(\sigma)}/Z_\beta.$$

$\{PM19, PM15\} = A:$

$$H = -1\times(-1)\times 1 - 1\times(-1)\times 1 - 1\times(-1)\times(-1) - 1\times(-1)\times 1 = 2.$$

$$\mathbb{P} = e^{-0.1\times 2}/Z_\beta = 0.40.$$

Result=N is the most probable guess. $\mathbb{P}(N) = 0.60.$

## Section II.1: Ideas

Configuration of guesses    Penalty Probability

| | | | | | | |
|---|---|---|---|---|---|---|
| -1 | -1 | -1 | -1 | -1 | 14 | 0.0019 |
| 1 | -1 | -1 | -1 | -1 | 10 | 0.0029 |
| -1 | 1 | -1 | -1 | -1 | 14 | 0.0019 |
| 1 | 1 | -1 | -1 | -1 | 2 | 0.0064 |
| -1 | -1 | 1 | -1 | -1 | 14 | 0.0019 |
| 1 | -1 | 1 | -1 | -1 | 2 | 0.0064 |
| -1 | 1 | 1 | -1 | -1 | 6 | 0.0043 |
| 1 | 1 | 1 | -1 | -1 | -14 | 0.0319 |
| -1 | -1 | -1 | 1 | -1 | 16 | 0.0016 |
| 1 | -1 | -1 | 1 | -1 | 12 | 0.0024 |
| -1 | 1 | -1 | 1 | -1 | 8 | 0.0035 |
| 1 | 1 | -1 | 1 | -1 | -4 | 0.0117 |
| -1 | -1 | 1 | 1 | -1 | 12 | 0.0024 |
| 1 | -1 | 1 | 1 | -1 | 0 | 0.0079 |
| -1 | 1 | 1 | 1 | -1 | -4 | 0.0117 |
| 1 | 1 | 1 | 1 | -1 | -24 | 0.0868 |

For each configuration of the unknown results (guesses), we can calculate its probability. We can determine the most probable guesses:

|      |      |      | Donor |      |      |      |      |
|------|------|------|-------|------|------|------|------|
|      | AM19 | PM19 | PM15  | UL11 | LL11 | LL15 | LL19 |
| AM19 | *N*  | N    | A     | A    | A    | A    | N    |
| PM19 | *N*  | N    | <u>N</u> | N | N | <u>N</u> | <u>N</u> |
| PM15 | *A*  | <u>N</u> | *N* | <u>N</u> | <u>N</u> | <u>N</u> | <u>N</u> |
| UL11 | *A*  | N    | <u>N</u> | N | N | <u>N</u> | <u>N</u> |
| LL11 | *A*  | N    | <u>N</u> | N | N | <u>N</u> | <u>N</u> |
| LL15 | *A*  | <u>N</u> | <u>N</u> | <u>N</u> | <u>N</u> | *N* | <u>N</u> |
| LL19 | *N*  | <u>N</u> | <u>N</u> | <u>N</u> | <u>N</u> | <u>N</u> | *N* |

Host (label for rows)

## Section II.1: Ideas

Since each configuration of guesses has a probability, we can take expectations, and obtain the probability for each experimental result to be "N":

|      |      | Donor |      |      |      |      |      |      |
| ---- | ---- | ----- | ---- | ---- | ---- | ---- | ---- | ---- |
|      |      | AM19  | PM19 | PM15 | UL11 | LL11 | LL15 | LL19 |
|      | AM19 | *100%* | 100% | 0%   | 0%   | 0%   | 0%   | 100% |
|      | PM19 | *100%* | 100% | 65%  | 100% | 100% | 49%  | 56%  |
| Host | PM15 | *0%*  | 65%  | *100%* | 62% | 62%  | 53%  | 54%  |
|      | UL11 | *0%*  | 100% | 62%  | 100% | 100% | 81%  | 81%  |
|      | LL11 | *0%*  | 100% | 62%  | 100% | 100% | 90%  | 90%  |
|      | LL15 | *0%*  | 49%  | 53%  | 81%  | 90%  | *100%* | 86% |
|      | LL19 | *100%* | 56%  | 54%  | 81%  | 90%  | 86%  | *100%* |

- For now, we have designed an inference method that works for experiments with binary deterministic results.
- What if the known experimental results are not deterministic, but stochastic?

# Section II.1: Another situation

|  |  | Donor | | | | |
|------|-----------|-------|-------|-------|-------|-------|
|  |  | PLE11 | PLE12 | PLE14 | PLE16 | PLE19 |
|  | LFR\PLE14 | 61% | 58% | 82% | ? | ? |
| Host | LFR\PLE16 | ? | ? | ? | ? | ? |
|  | LFR\PLE19 | 4% | 24% | 83% | ? | 100% |

Table: Results for *Xenopus laevis*, reported by Henry et al. 1987

- Percentage is the probability of normal development (N).
- PLE11: presumptive lens ectoderm, stage 11.
- LFR\PLE14: lens-forming region without presumptive lens ectoderm, stage 14.

- Sample deterministic configurations from these stochastic results (assume different experiments are independent).
- For each deterministic configuration, apply our method to obtain the expectation of guesses.
- For example, assume we have three similar experiments: [61%N ? 58%N].

    Sample deterministic results: $\mathbb{P}([N \ ? \ N]) = 61\% \times 58\% = 35\%$.

    Apply the inference method: $\mathbb{P}(?=N \mid [N \ ? \ N]) = 98\%$.

$\mathbb{P}([N \ N \ N]) = \mathbb{P}([N \ ? \ N]) \times \mathbb{P}(?=N \mid [N \ ? \ N]) = 35\% \times 98\% = 35\%$.

## Section II.1: Another situation

Similarly, we can calculate for other deterministic configurations of [61%N ? 58%N]:

$$\mathbb{P}([N\ ?\ A]) = 61\% \times (100\% - 58\%) = 26\%.$$

$$\mathbb{P}([N\ N\ A]) = 26\% \times 50\% = 13\%.$$

$$\mathbb{P}([A\ ?\ N]) = (100\% - 61\%) \times 58\% = 23\%.$$

$$\mathbb{P}([A\ N\ N]) = 23\% \times 50\% = 11\%.$$

$$\mathbb{P}([A\ ?\ A]) = (100\% - 61\%) \times (100\% - 58\%) = 16\%.$$

$$\mathbb{P}([A\ N\ A]) = 16\% \times 2\% = 0\%.$$

Then average over these deterministic configurations:

$$\mathbb{P}([?{=}N]) = \mathbb{P}([N\ N\ N]) + \mathbb{P}([N\ N\ A]) + \mathbb{P}([A\ N\ N]) + \mathbb{P}([A\ N\ A]) = 59\%.$$

Final results: [61%N 59%N 58%N].

# Section II.1: Another situation

|      |           | Donor |       |       |       |       |
|------|-----------|-------|-------|-------|-------|-------|
|      |           | PLE11 | PLE12 | PLE14 | PLE16 | PLE19 |
|      | LFR\PLE14 | 61%   | 58%   | 82%   | ?     | ?     |
| Host | LFR\PLE16 | ?     | ?     | ?     | ?     | ?     |
|      | LFR\PLE19 | 4%    | 24%   | 83%   | ?     | 100%  |

Table: Results for *Xenopus laevis*, reported by Henry et al. 1987

|      |           | Donor |       |       |       |       |
|------|-----------|-------|-------|-------|-------|-------|
|      |           | PLE11 | PLE12 | PLE14 | PLE16 | PLE19 |
|      | LFR\PLE14 | 61%   | 58%   | 82%   | 93%   | 94%   |
| Host | LFR\PLE16 | 39%   | 53%   | 88%   | 97%   | 97%   |
|      | LFR\PLE19 | 4%    | 24%   | 83%   | 96%   | 100%  |

Table: Inferred results

## Section II.1: Yet another situation

- The methods are for experiments with binary results.
- What if the known experimental results are not binary?
- The penalty function is:

$$\mathrm{H}(\sigma) = -\sum_{i \sim j} J_{ij} \sigma_i \sigma_j.$$

  The cross term $\sigma_i \sigma_j$ measures the similarity between $\sigma_i$ and $\sigma_j$.

- Rewrite the penalty function:

$$\mathrm{H}(\sigma) = -\sum_{i \sim j} J_{ij} f(\sigma_i, \sigma_j).$$

- If $\sigma_i, \sigma_j$ are more similar, $f(\sigma_i, \sigma_j)$ is larger. Also, $f(\sigma_i, \sigma_j) = f(\sigma_j, \sigma_i)$. For binary case, $f(\sigma_i, \sigma_j) = \sigma_i \sigma_j$.

## Section II.1: Summary

- Based on the similarities between experiments, we have designed methods to infer the unknown experimental results.
- The results are not necessarily deterministic or binary.
- In the future, we hope to have more experimental data to verify the inference results and determine the parameters.
- Such methods should not be limited to tissue transplantation experiments.

- Questions?

- Assume there are many tissue transplantation experiments, and we do not have any results yet.
- To know all the results, we can choose some experiments to conduct, and apply our method to infer the unknown experimental results.
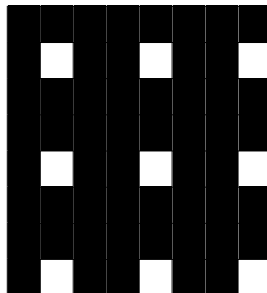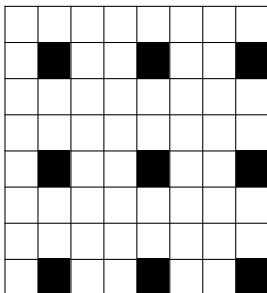- How to choose experiments to conduct?

# Section II.2: Experimental design

- Assume the experiment similarity chart is 2-D lattice. Each unit is an experiment, and neighboring units are similar experiments.
- Black units are conducted experiments, and white units are non-conducted experiments.
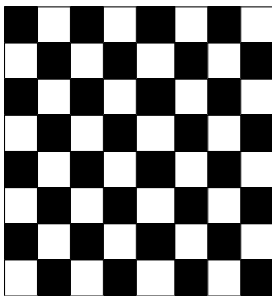- Experimental design (choosing experiments to conduct) becomes coloring the chart.

We need enough data to apply the inference method. We should minimize the experimental cost.
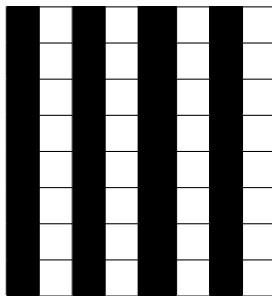
- The results of non-conducted experiments are inferred by similar conducted experiments.
- To guarantee the inference quality, one non-conducted experiment should be similar to at least $k$ conducted experiments.
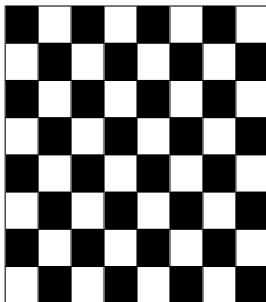


$k = 4$ $\qquad\qquad$ $k = 2$

- If one non-conducted experiment is similar to at least $k$ conducted experiments, how to minimize the number of conducted experiments?
- The most efficient design: no conducted experiments are similar, and each non-conducted experiment is similar to exactly $k$ conducted experiments.
- How to color the 2-D square lattice $\mathbb{Z}^2$, so that two black units are not neighboring, and each white unit is neighboring to exactly $k$ black units?
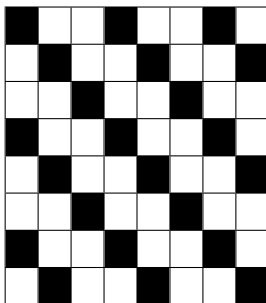
No neighboring black units, and each white unit is neighboring to $k = 4$ black units (ignore the boundary cases).



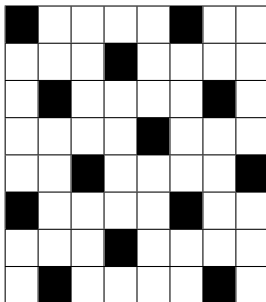We need to conduct $1/2$ experiments.

No neighboring black units, and each white unit is neighboring to $k = 2$ black units (ignore the boundary cases).



We need to conduct $1/3$ experiments.

No neighboring black units, and each white unit is neighboring to $k = 1$ black units (ignore the boundary cases).



We need to conduct $1/5$ experiments.

# Section II.2: Experimental design

- In practice, the experiment similarity chart is not 2-dimensional, but 4-dimensional. Each experiment has a coordinate $(x, y, z, w)$ that stands for donor tissue type, donor tissue developmental stage, host tissue type, host tissue developmental stage.
- For now, we assume the chart is $\mathbb{Z}^4$.
- Similar coloring problems for such 4-dimensional figures.
- We need some abstract methods.

## Section II.2: Experimental design

- In $\mathbb{Z}^4$, each unit is neighboring to 8 units.
- For $k = 8$, color a unit $(x, y, z, w)$ if

$$x + y + z + w \equiv 0 \mod 2.$$

We need to conduct $1/2$ experiments.

- For $k = 4$, color a unit $(x, y, z, w)$ if

$$x + y + z + w \equiv 0 \mod 3.$$

We need to conduct $1/3$ experiments.

- For $k = 2$, color a unit $(x, y, z, w)$ if

$$x + 2y + z + 2w \equiv 0 \mod 5.$$

We need to conduct $1/5$ experiments.

- For $k = 1$, color a unit $(x, y, z, w)$ if

$$x + 2y + 3z + 4w \equiv 0 \mod 9.$$

  We need to conduct $1/9$ experiments.

- Such methods can be generalized to $\mathbb{Z}^n$ and $k$ that $k \mid 2n$.

- Color a unit $(x_1, \ldots, x_n)$ if $a_1 x_1 + \cdots + a_n x_n \equiv 0$ mod $(2n/k + 1)$. Here if $k \mid n$, $(a_1, \ldots, a_n)$ are $k$ groups of $1, 2, \ldots, n/k$; otherwise, $(a_1, \ldots, a_n)$ are $k/2$ groups of $1, 2, \ldots, 2n/k$.

- In practice, $k = 2$ or $k = 1$ is enough to conduct satisfactory inference. Therefore we only need to conduct $1/5 - 1/3$ experiments (two-dimensional) or $1/9 - 1/5$ experiments (four-dimensional).

- We have solved a problem: how to choose experiments to conduct, so that other unknown results can be inferred properly with minimal cost.

- In mathematical biology, sometimes the right key (mathematics) is unexpected.

- Questions?

- Reveal the existence of multiple states in cancer cell population, and prove the state equilibrium phenomenon with ODEs and branching processes.
- Develop methods to infer the results of tissue transplantation experiments, and solve the experimental design problem with combinatorics.
- The same flavor: analyze biological data; build models; extract and solve meaningful mathematical problems.

# Thank you!

# Other related works

- Mathematical biology.
- Design algorithms to find "jumping genes" in gene sequences. Design algorithms to calculate the distance between developmental trees.
- Analyze the notion of "positional information" in developmental biology. Build models for embryo development.
- Ongoing: Infer the structure of gene regulatory networks. Build models for cell membrane electric potential.

- Applied probability (especially in machine learning theory).
- Causal inference: Impossibilities in quantifying causal relationships.
- Reinforcement learning: Policy evaluation in pricing processes with historical data (possibly polluted).
- Statistical physics: Entropy production of lifted Markov chains.