# Inference on Gene Regulation and a Stochastic Model of Molecule Selection

Yue Wang

Department of Computational Medicine,
University of California, Los Angeles

Jan. 06, 2023

- Two projects: inferring gene regulatory relationships with different types of data; optimizing the protocol of a molecule selection process.
- Yue Wang, and Zikun Wang. (2022). "Inference on the structure of gene regulatory networks." Journal of Theoretical Biology, 539, 111055.
- Yue Wang, Bhaven A. Mistry, and Tom Chou. (2022). "Discrete stochastic models of SELEX: aptamer capture probabilities and protocol optimization." Journal of Chemical Physics, 156(24), 244103.

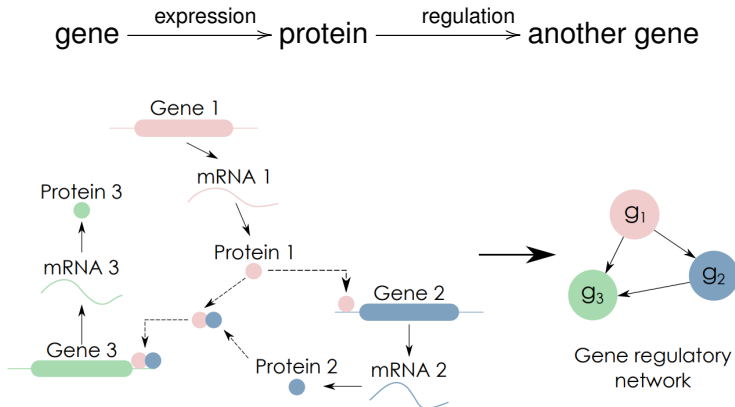# Section I: Inference on Gene Regulatory Relationships

- Introduction to gene regulatory networks (GRN).
- Types of data that can be used to infer GRN structures.
- A framework for data type classification.
- Mathematical inference methods for GRN structures.

## Section I.1: Introduction

- Gene expression: genes are transcribed to mRNAs and then translated to proteins.
- Various molecular regulators affect gene expression (change levels of mRNAs and/or proteins).
- Two types of regulation: activation and inhibition.
- Some regulators are small molecules, such as oxygen, sugars and vitamins.
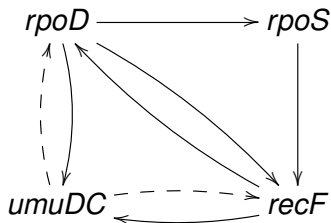
- Some regulators are proteins. We focus on regulations between genes.

$$\text{gene} \xrightarrow{\text{expression}} \text{protein} \xrightarrow{\text{regulation}} \text{another gene}$$



Gene regulatory network

- Genes and their regulatory relations form a gene regulatory network (GRN).

$rpoD \longrightarrow rpoS$

$umuDC \dashrightarrow recF$

- An example of GRN in *E. coli*. Each vertex is a gene. Two types of regulations: solid arrow means activation, and dashed arrow means inhibition.
- We aim at determining the GRN structure.
- For two genes $G_i, G_j$, does the expression of $G_i$ activate or inhibit the expression of $G_j$?

## Section I.1: Introduction

- Genes (DNAs), mRNAs and proteins are generally confined within living cells.
- It is extremely difficult or even impossible to directly determine whether one gene regulates another gene with biochemical methods.
- We have accumulated a large amount of data, e.g., bulk level gene expression data and single-cell level phenotype data.

- Certain types of data can be used to infer the GRN structure.
- They can be classified in different dimensions.
- Setup: consider a set of genes $G_1, \ldots, G_n$ that possibly regulate each other.

- Dimension 1: Gene expression vs. Phenotype.
- We can measure the expression levels of genes $G_1, \ldots, G_n$.
- We can also measure the level of a phenotype $G_0$ (e.g., growth rate, drug resistance) which is affected by these genes.
- For inferring gene regulatory relationships, gene expression data are more informative than phenotype data.

- Dimension 2: Single-cell vs. Bulk.

- The gene expression of a single cell is stochastic. We can measure the levels of $G_1, \ldots, G_n$ for a single cell and repeat many times, so as to obtain a group of random variables $X_1, \ldots, X_n$ that represent the random levels of $G_1, \ldots, G_n$.

- We can also measure these quantities over a large population of cells (bulk level), so that the randomness is averaged out. Then we obtain deterministic results $x_1, \ldots, x_n$.

- After taking expectation, single-cell data become bulk data. Thus single-cell data are more informative than bulk data.

## Section I.2: Data types

- Dimension 3: Interventional vs. Non-interventional.
- We can intervene with certain genes (siRNA, CRISPR, etc.), so that the expression levels of these genes are changed. Then other related genes are also affected.
- We can measure expression levels $x_1', \ldots, x_n'$ after interfering with certain genes, and compare with corresponding quantities before intervention $x_1, \ldots, x_n$.
- We can also observe without any intervention.
- Interventional data are more informative than non-interventional data.

# Section I.2: Data types

- Dimension 4: One-time vs. Time series.
- We can measure at a single time point, $X_i(0)$.
- We can also measure at multiple time points as a time series, $X_i(0), X_i(1), X_i(2), \ldots$.
- With time series data, we can study the dynamics of gene expression.
- Time series data are more informative than one-time data.

# Section I.2: Data types

- When we measure at single-cell level at multiple time points, we obtain a sequence of random variables $X_i(0), X_i(1), X_i(2), \ldots$.
- Extra dimension: Joint distribution vs. Marginal distribution.
- If the same cell can be measured multiple times, we obtain the joint distribution for multiple time points, $\mathbb{P}[X_i(0) = c_0, X_i(1) = c_1, X_i(2) = c_2]$.
- Most measurements are destructive, meaning that one cell can be measured only once. If so, we can only obtain the marginal distribution for each time point, $\mathbb{P}[X_i(0) = c_0], \mathbb{P}[X_i(1) = c_1], \mathbb{P}[X_i(2) = c_2]$.
- With the joint distribution, we can obtain more information, such as correlation coefficients.
- Joint distribution data are more informative than marginal distribution data.

# Section I.2: Data types

- We have four major dimensions: (1) Gene expression or Phenotype; (2) Single-cell or Bulk; (3) Non-interventional or Interventional; (4) One-time or Time series.
- According to these four dimensions, we have $2^4 = 16$ different data types (scenarios).
- In four scenarios (Single-cell + Time series), there is an extra dimension of Joint distribution or Marginal distribution, meaning a total of 20 scenarios.

# Section I.2: Data types

| | | One-Time | | Time Series | |
|---|---|---|---|---|---|
| | | Non-Intervention | Intervention | Non-Intervention | Intervention |
| Gene Expression | Single-Cell | Scenario 1 | Scenario 2 | Scenario 3a Joint<br><br>Scenario 3b Marginal | Scenario 4a Joint<br><br>Scenario 4b Marginal |
| | Bulk | Scenario 5 | Scenario 6 | Scenario 7 | Scenario 8 |
| Phenotype | Single-Cell | Scenario 9 | Scenario 10 | Scenario 11a Joint<br><br>Scenario 11b Marginal | Scenario 12a Joint<br><br>Scenario 12b Marginal |
| | Bulk | Scenario 13 | Scenario 14 | Scenario 15 | Scenario 16 |

All 20 scenarios, classified by data types.

Questions?

## Section I.2: Data types

- For each dimension, one choice is more informative than the other, such as Time series $>$ One-time. The most informative data type is Scenario 4a: Gene expression + Single-cell + Interventional + Time series + Joint distribution.
- Nevertheless, for more informative data types, generally the experiments are more difficult, more expensive, and less accurate.
- Less informative scenarios are also worth studying.

## Section I.3: Inference methods

- Different scenarios (data types) require different mathematical inference methods.
- In reality, the regulation of gene expression is very complicated, with many unknown mechanisms.
- We have to make some assumptions about GRN and data, so that the gene expression follows certain mathematical models.
- One necessary assumption is that we can observe all related factors (no hidden variable).

# Section I.3: Inference methods

- Four common assumptions:
- Path Blocking (PB): the intervention on one gene has no effect on another gene (or a phenotype), if and only if other intervened genes have already blocked all paths.
- Directed Acyclic Graph (DAG): the GRN can be described by a directed graph without directed cycles.
- Markov and Faithful (MF): the distribution of gene expression properly reflects the underlying DAG through conditional independence relations.
- Linear System (LS): the gene expression (and possibly phenotype) time series data satisfy a linear ODE system.
- Those assumptions might not hold in reality. Thus the inference results are not ground truths.

- For each scenario, we discuss what structures can be inferred, and what assumptions are required.
- Scenarios 1/3/4/8 have been extensively studied. For other scenarios, we invent new mathematical methods, or prove that the GRN structure cannot be inferred.

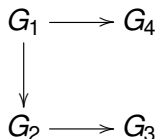| | | One-Time | | Time Series | |
|---|---|---|---|---|---|
| | | Non-Intervention | Intervention | Non-Intervention | Intervention |
| Gene Expression | Single-Cell | Scenario 1: MF+DAG: partial. | Scenario 2: PB: full. DAG: partial. MF+DAG: full. | Scenario 3 a/b: 3a Joint: UC: full. 3b Marginal: MF+DAG: partial. | Scenario 4 a/b: 4a Joint: UC: full. 4b Marginal: LS: full. PB: full. DAG: partial. MF+DAG: full. |
| | Bulk | Scenario 5: No. | Scenario 6: PB: full. DAG: partial. | Scenario 7: No. | Scenario 8: LS: full. PB: full. DAG: partial. |

Inference results for different scenarios (part I). MF, DAG, PB, LS: mathematical assumptions required by corresponding inference methods. UC: no assumption required. Full/partial/no means all/some/no GRN structures can be inferred. Blue methods are known; red methods are novel.

# Section I.3: Inference methods

| | | One-Time | | Time Series | |
|---|---|---|---|---|---|
| | | Non-Intervention | Intervention | Non-Intervention | Intervention |
| Phenotype | Single-Cell | Scenario 9: No. | Scenario 10: PB: partial. | Scenario 11 a/b: No. | Scenario 12 a/b: PB: partial. LS+DAG: partial*. PB+LS+DAG: partial*. |
| | Bulk | Scenario 13: No | Scenario 14: PB: partial. | Scenario 15: No. | Scenario 16: PB: partial. LS+DAG: partial*. PB+LS+DAG: partial*. |

Inference results for different scenarios (part II). DAG, PB, LS: mathematical assumptions required by corresponding inference methods. Partial/no means some/no GRN structures can be inferred. Asterisk means activation/inhibition cannot be determined. Red methods are novel.

- Scenario 4 (gene expression, single-cell, interventional, time series) is the most informative case, and there are many methods to fully determine the GRN structure.
- Phenotype data and non-interventional data are much less informative, and the inference results are generally limited.

- Questions?

## Section I.4: Example

- In Scenario 6 (gene expression, bulk, interventional, one-time), we can partially infer the GRN structure under the DAG assumption.
- DAG: directed acyclic graph, meaning that the GRN has no directed cycle.
- GRN is represented by a DAG. Each vertex is a gene, and each directed edge is a regulatory relation.

$$G_1 \longrightarrow G_4$$
$$\downarrow$$
$$G_2 \longrightarrow G_3$$

- In a DAG, if there is a directed path from $G_i$ to $G_j$, then $G_i$ is an ancestor of $G_j$, and $G_j$ is a descendant of $G_i$.

- $G_1 \longrightarrow G_4$

  $\downarrow$

  $G_2 \longrightarrow G_3$

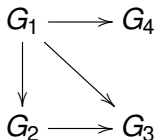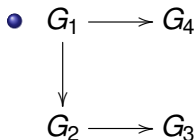  $G_1$ has descendants $G_2, G_3, G_4$; $G_2$ has descendant $G_3$; $G_3$ and $G_4$ have no descendant.

## Section I.4: Example

- After adding intervention on gene $G_i$, gene $G_j$ is also affected, if and only if $G_j$ is a descendant of $G_i$ in the DAG.
- With such intervention experiments, we can determine the ancestor-descendant relations between genes.
- Now we have a mathematical problem: given the ancestor-descendant relations of a DAG, how to infer its structure?
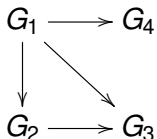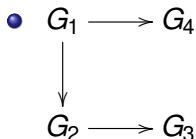
- $G_1 \longrightarrow G_4$

  $\downarrow$

  $G_2 \longrightarrow G_3$

  $G_1$ has descendants $G_2, G_3, G_4$; $G_2$ has descendant $G_3$; $G_3$ and $G_4$ have no descendant.

- $G_1 \longrightarrow G_4$      $G_1 \longrightarrow G_4$

  $G_2 \longrightarrow G_3$      $G_2 \longrightarrow G_3$
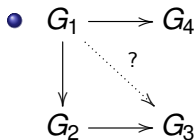
- $G_1$ has descendants $G_2, G_3, G_4$; $G_2$ has descendant $G_3$; $G_3$ and $G_4$ have no descendant.
- The same ancestor-descendant relations might correspond to multiple DAGs, and they are called "AD equivalent".
- All DAGs that are AD equivalent form an equivalent class. The above two DAGs form an equivalent class.
- From the ancestor-descendant relations, we can only determine the equivalent class that contains the true DAG.

## Section I.4: Example

- $G_1 \longrightarrow G_4 \qquad\qquad G_1 \longrightarrow G_4$

  $\qquad\downarrow \qquad\qquad\qquad\qquad \downarrow\qquad\nearrow$

  $G_2 \longrightarrow G_3 \qquad\qquad G_2 \longrightarrow G_3$

- Edges $G_1 \to G_2$, $G_2 \to G_3$, $G_1 \to G_4$ appear in all DAGs in this equivalent class. Thus the true DAG must have these edges.

- Edge $G_1 \to G_3$ appears in some but not all DAGs in this equivalent class. Thus the true DAG might have this edge.

- Other edges appear in no DAGs in this equivalent class. Thus the true DAG does not have those edges.

- $G_1 \longrightarrow G_4$

  $\qquad\downarrow \qquad ?$

  $G_2 \longrightarrow G_3$

  We can partially infer the DAG structure.

## Section I.4: Example

- 1. From the intervention experiments, obtain the ancestor-descendant relations among genes.
- 2. Find all DAGs that follow these ancestor-descendant relations.
- 3. Identify edges that appear in all these AD equivalent DAGs, and edges appear in some but not all these AD equivalent DAGs.
- 4. Produce a partially determined DAG.
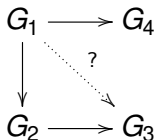- Steps 2 and 3 are cumbersome. We need to search in exponentially many DAGs.

# Section I.4: Example

A much faster approach:

### Theorem

*Given ancestor-descendant relations:*

*(1) If $G_j$ is not a descendant of $G_i$, then we can determine that the edge $G_i \rightarrow G_j$ does not exist in the true DAG.*

*(2) If $G_j$ is a descendant of $G_i$, and $G_i$ has another descendant $G_k$, which is an ancestor of $G_j$, then we cannot determine the existence of the edge $G_i \rightarrow G_j$ in the true DAG.*

*(3) If $G_j$ is a descendant of $G_i$, and $G_i$ does not have another descendant $G_k$, which is an ancestor of $G_j$, then we can determine that the edge $G_i \rightarrow G_j$ exists in the true DAG.*

$$G_1 \longrightarrow G_4$$

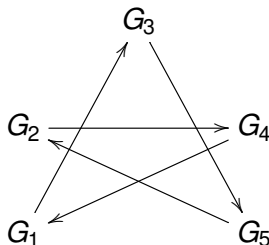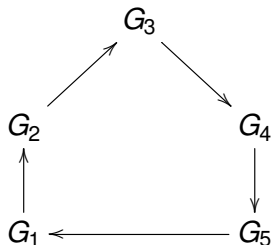$$G_2 \longrightarrow G_3$$

# Section I.4: Example

Although not all edges can be inferred, we have a lower bound for the number of edges that can be inferred.

### Theorem

*If the GRN is a connected DAG with n vertices, then we can use ancestor-descendant relations to identify at least n − 1 edges.*

- If the GRN has directed cycles, we might infer no edge.
- 

- These two GRNs share the same ancestor-descendant relations, but they have no common edges. Thus we cannot determine the existence of any edges.

- Questions?

## Section I: Discussion

- What if we have multiple types of data?
- With new technologies, there might be new types of data that do not fit in our framework.
- More informative data types are more expensive and less accurate. How can we design experiments to infer GRN structures, while the cost is minimized?
- Can we use such data to infer the existence of autoregulation (a gene that regulates its own expression)?

## Section I: Summary

- Introduce the GRN structure inference problem.
- Classify the inference problem into 20 scenarios.
- Previous studies are unified under a few scenarios. Invent mathematical methods for scenarios that have not been extensively studied.
- This work provides a unified framework to discuss the GRN structure inference problem.

- Questions?

# Section II: Protocol Optimization of a Molecule Selection Process

- Introduce SELEX: a process to select aptamers.
- Review the traditional deterministic model of SELEX.
- Build a stochastic model for SELEX and analyze its properties.
- Search for the optimal protocol of SELEX.

- Aptamers are short, single-stranded DNA or RNA molecules that bind to a specific target.
- Targets can be heavy metal ions, proteins, or even whole cells.
- Certain aptamers (linked with fluorescent tracers) can bind selectively to biomarkers on cancer cells, but not to healthy cells. This test can identify cancer cells in a tissue sample.
- Besides testing, aptamers can also be used in treatments. Therefore, aptamers are also called chemical antibodies.

- It is difficult to design and synthesize the best aptamer for a target directly.
- In general, we start with enough targets and a large library of randomly generated aptamers, and aptamers have different affinities to the targets.
- How to select the best aptamers (with the highest affinities to the targets) in an easy way?

# Section II.1: Introduction

- Systematic Evolution of Ligands by EXponential enrichment (SELEX): a convenient method to select the best aptamers.
- If we mix aptamers and targets, aptamers with higher affinities to the targets are more likely to bind to the targets. We can use the targets to pick out such aptamers.
- It is similar to a population evolution process.
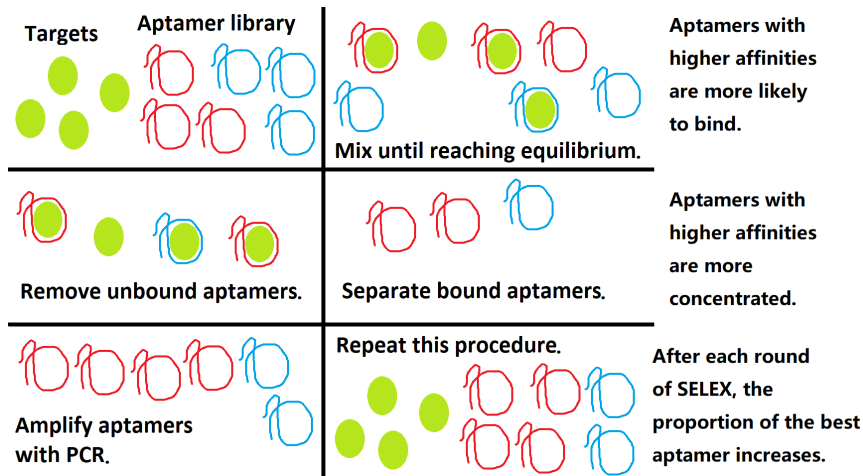
Aptamers and targets can bind and unbind reversibly.



Figure: Protocol of SELEX

- We have enough targets, and the aptamers can be amplified by PCR. When starting one round of SELEX, we can control the quantity of targets and the quantity of aptamers, but the proportions of different aptamer types cannot be controlled.

- We obtain an optimization problem: maximize the proportion of the best aptamer (with the highest affinity) after this round of SELEX.

- Questions?

- We need a mathematical model to study the optimization of SELEX protocol.
- To simplify the discussion, we combine aptamers with different affinities into two types: strong type $A_1$, weak type $A_2$. The association constants (affinities) satisfy $K_1 > K_2$.

# Section II.2: Deterministic model

- A traditional deterministic approach uses the law of mass action, which is valid when the number of molecules is sufficiently large.
- Notations: $[T]$: total concentration of targets $\mathtt{T}$; $[A_i]$: total concentration of aptamer type $\mathtt{A}_i$; $[a_i]$: concentration of aptamers $\mathtt{A}_i$ that are bound to targets at equilibrium.
- At equilibrium, for each $i = 1, 2$ and the reaction $\mathtt{T} + \mathtt{A}_i \rightleftharpoons \mathtt{TA}_i$, we have:

$$([T] - [a_1] - [a_2])\,([A_i] - [a_i])\,K_i = [a_i].$$
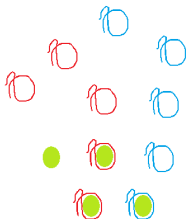
| unbound | unbound | bound |
| target | aptamer | aptamer |

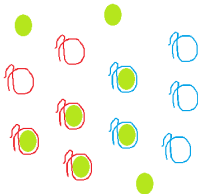- Given $[T], [A_1], [A_2], K_1, K_2$, we can solve $[a_1], [a_2]$.

- For $A_1$, the stronger aptamer, the goal is to maximize its proportion in bound aptamers: $[a_1]/([a_1] + [a_2])$.
- We can set different values of target concentration $[T]$ and aptamer concentration $[A_1]$, $[A_2]$, but the ratio $[A_1]/[A_2]$ is fixed.
- In this deterministic model, $[a_1]/([a_1] + [a_2])$ increases with $[A_1]$ (and $[A_2]$), and decreases with $[T]$.
- The optimal policy in the deterministic model: add as many aptamers as possible, and as few targets as possible.
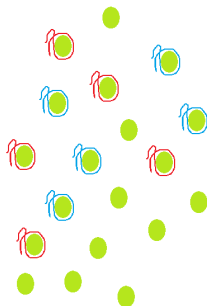
Optimal policy: $[A_i] \gg [T]$.



**Very few targets** | **Intermediate** | **Too many targets**

$$\frac{[a_1]}{[a_1] + [a_2]} = \frac{[A_1]K_1}{[A_1]K_1 + [A_2]K_2}$$

upper bound

$$\frac{[a_1]}{[a_1] + [a_2]} = \frac{[A_1]}{[A_1] + [A_2]}$$

lower bound

## Section II.2: Deterministic model

- The optimal policy in the deterministic model requires very large aptamer concentration $[A_1]$ (and $[A_2]$) and very small target concentration $[T]$.
- When $[T]$ is too small, randomness is inevitable, and the law of mass action does not hold.
- We need a stochastic model.
- We will show that something is different in this stochastic model.

- Questions?

## Section II.3: Stochastic model

- Notations: $T$: total number of targets; $A_i$: total number of $\mathbb{A}_i$ type aptamers; $a_i$: number of $\mathbb{A}_i$ aptamers that are bound to targets. $\bar{K}_i = K_i/V$: reaction coefficient, where $V$ is the system volume.

- Consider a continuous-time Markov chain on 2-dimensional lattice $\mathbb{Z}^2$, where the states are the bound aptamer counts $(a_1, a_2)$.

- The transition rates satisfy

$$\frac{r[(a_1, a_2) \to (a_1 + 1, a_2)]}{r[(a_1 + 1, a_2) \to (a_1, a_2)]} = \frac{(T - a_1 - a_2)(A_1 - a_1)}{a_1 + 1} \bar{K}_1.$$

$$\frac{r[(a_1, a_2) \to (a_1, a_2 + 1)]}{r[(a_1, a_2 + 1) \to (a_1, a_2)]} = \frac{(T - a_1 - a_2)(A_2 - a_2)}{a_2 + 1} \bar{K}_2.$$

- The stationary probability distribution satisfies

$$\mathbb{P}(a_1, a_2) = \mathbb{P}(0, 0) \times \binom{T}{T - a_1 - a_2, a_1, a_2}$$
$$\times \left[ \binom{A_1}{a_1} \binom{A_2}{a_2} \right] \times [a_1! a_2!] \times \left[ \bar{K}_1^{a_1} \bar{K}_2^{a_2} \right]$$

- In the stochastic model, when the total aptamer numbers $A_1, A_2$ and the total target number $T$ are very small, it is possible that no aptamer is bound to a target. This means $a_1 = a_2 = 0$, and $a_1/(a_1 + a_2)$ is not defined.

- In practice, we only want $\mathrm{A}_1$ aptamers. When $a_1 = a_2 = 0$, we can stipulate that $a_1/(a_1 + a_2) = 0$.

- Now we can consider the expected $\mathrm{A}_1$ proportion, $\mathbb{E}[a_1/(a_1 + a_2)]$.
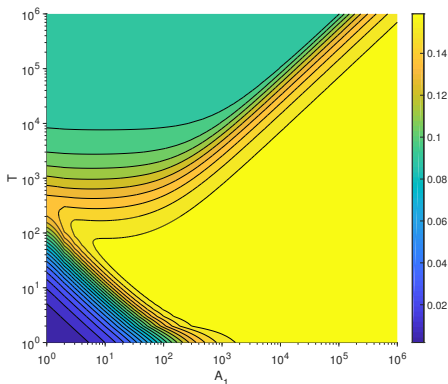
# Section II.3: Stochastic model

- In the stochastic model, we still have the same upper bound for $A_1$ proportion.

### Theorem

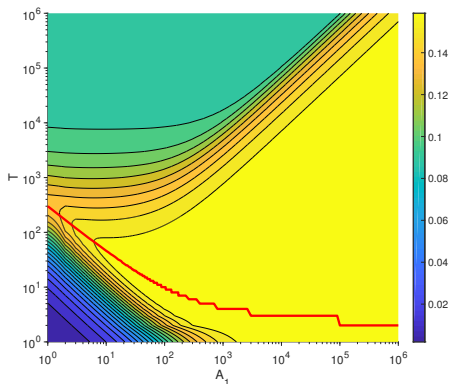$\mathbb{E}[a_1/(a_1 + a_2)] \leq A_1 \bar{K}_1/(A_1 \bar{K}_1 + A_2 \bar{K}_2).$

- When $A_1, A_2, T$ are very small, it is very likely that $a_1 = a_2 = 0$, so that $\mathbb{E}[a_1/(a_1 + a_2)] \approx 0$. Thus we do not have the same lower bound.

Contour plot of $\mathbb{E}[a_1/(a_1 + a_2)]$. Unlike the deterministic model, in the stochastic model, $\mathbb{E}[a_1/(a_1 + a_2)]$ does not always increase with $A_1$ (fix $A_1/A_2$), and does not always decrease with $T$.

The red curve indicates the optimal target number $T$ that maximizes $\mathbb{E}[a_1/(a_1 + a_2)]$ for each aptamer number $A_1$ (fix $A_1/A_2$). When the aptamer number $A_1$ (fix $A_1/A_2$) increases, the optimal target number $T$ decreases.
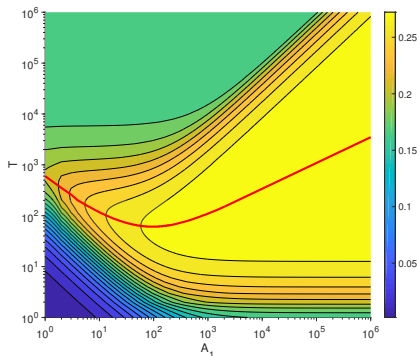
- Optimal policy in the stochastic model:
- When the aptamer number $A_1$ (and $A_2$) is not large, the target number $T$ should not too small. Otherwise, $\mathbb{P}(a_1 = 0, a_2 = 0)$ might be large.
- What if we make $A_1$ (and $A_2$) sufficiently large, so that $\mathbb{P}(a_1 = 0, a_2 = 0) \approx 0$? Can we set $T = 1$ now?

## Section II.4: Optimal policy

- For one round of SELEX, $T = 1$ and very large $A_1, A_2$ can reach the upper bound:
  $\mathbb{E}[a_1/(a_1 + a_2)] \approx A_1 \bar{K}_1/(A_1 \bar{K}_1 + A_2 \bar{K}_2)$.
- However, since there is only one target molecule, after one round of SELEX, only one aptamer type is left.
- After further rounds of SELEX, the expected $\mathtt{A}_1$ proportion does not increase.

- Contour plot of the $A_1$ proportion $\mathbb{E}[a_1/(a_1 + a_2)]$ after two rounds of SELEX:



- For the first round, a policy with large $A_1$ and very small $T$ does not perform well. The optimal target number $T$ (red curve) first decreases and then increases with $A_1$.

### Theorem

*The optimal policy for multiple rounds of SELEX in the stochastic model is $A_1, A_2 \gg T$ and $T \gg 1$.*

- After $N$ rounds of SELEX, this policy has $\mathbb{E}[a_1/(a_1 + a_2)] \approx A_1 \bar{K}_1^N/(A_1 \bar{K}_1^N + A_2 \bar{K}_2^N)$.
- Thus $1 - \mathbb{E}[a_1/(a_1 + a_2)]$ converges to 0 exponentially fast with the rate $\approx \bar{K}_2/\bar{K}_1$. This is the most important factor for the efficiency of multi-round SELEX.

- Questions?

## Section II: Discussion

- We only consider one type of target. What if we want to select aptamers with higher affinities to target $T_1$ but lower affinities to target $T_2$?

- In the current protocol, after mixing aptamers and targets, we wait until equilibrium. What if we stop mixing before reaching equilibrium?

- Specifically, if binding is much faster than unbinding ($\bar{K}_1, \bar{K}_2 \gg 1$), we can stop the mixing when binding is all done, while unbinding has not started.

- This theoretical analysis can be applied to other scenarios, such as selecting drug-resistant cells.

## Section II: Summary

- We discuss SELEX, a process to select the best aptamer for binding a target.
- In the traditional deterministic model, the optimal policy (for any rounds of SELEX) is to have very large aptamer numbers $A_1, A_2$ and a very small target number $T$.
- We develop a stochastic model, in which the optimal policy for multiple rounds of SELEX is to have very large aptamer numbers $A_1, A_2$ but a moderate target number $T$.

- Questions?

# Thank you!