Combinatorial Algorithms for Tumor Phylogenetics

Mohammed El-Kebir







Clonal Theory of Cancer [Nowell, 1976]

Mutation



Clonal Theory of Cancer [Nowell, 1976]





Clonal Theory of Cancer [Nowell, 1976]





Clonal Theory of Cancer [Nowell, 1976]



Heterogeneous Tumor



Cancer Evolution: Cell Division, Mutation & Migration



Understanding Tumor Life History Has Clinical Applications



Outline – Reconstructing Cancer Evolution



Precise mathematical models are needed to describe the evolutionary process in cancer



















Tumor Phylogeny Estimation: Given frequencies F, find (1) phylogeny T and (2) proportions U



Tumor Phylogeny Estimation

Metastatic Colorectal Cancer (Patient CRC2)

[Kim et al., Clin Cancer Res 21(19), 2015]:

- 5 primary samples (P1-P5)
- 2 metastases (M1-M2)
- 412 single-nucleotide variants (SNVs)



CRC-2

Tumor Phylogeny Estimation

Metastatic Colorectal Cancer (Patient CRC2)

[Kim et al., Clin Cancer Res 21(19), 2015]:

- 5 primary samples (P1-P5)
- 2 metastases (M1-M2)
- 412 single-nucleotide variants (SNVs)
- 41 mutate more than once (homoplasy)



CRC-2

Heuristic for Tumor Phylogeny Estimation

Metastatic Colorectal Cancer (Patient CRC2)

[Kim et al., Clin Cancer Res 21(19), 2015]:

- 5 primary samples (P1-P5)
- 2 metastases (M1-M2)
- 412 single-nucleotide variants (SNVs)
- 41 mutate more than once (homoplasy)

division/mutation history or the migration history





Assumptions:

- Infinite sites assumption: a character changes state once
- Error-free data



Tumor Phylogeny Estimation: Given frequencies *F*, find (1) phylogeny *T* and (2) proportions *U*



Assumptions:

- Infinite sites assumption: a character changes state once
- Error-free data

Tumor Phylogeny Estimation: Given frequencies *F*, find (1) phylogeny *T* and (2) proportions *U*



VAF Factorization Problem (VAFFP) [El-Kebir*, Oesper* et al., 2015] Given *F*, find *U* and *B* such that *F* = *U B*

Assumptions:

- Infinite sites assumption: a character changes state once
- Error-free data

Variant of VAFFP:

TrAp [Strino *et al.,* 2013], PhyloSub [Jiao *et al.,* 2014] CITUP [Malikic *et al.,* 2015], BitPhylogeny [Yuan *et al.,* 2015] LICHEE [Popic *et al.,* 2015], ...



VAF Factorization Problem (VAFFP) [El-Kebir*, Oesper* et al., 2015] Given *F*, find *U* and *B* such that *F* = *U B*



VAF Factorization Problem (VAFFP) [El-Kebir*, Oesper* et al., 2015] Given *F*, find *U* and *B* such that *F* = *U B*

Given F and T (or B), is there a usage matrix U?

VAFFP: Given *F*, find *U* and *B* such that *F* = *U B*



Given F and T (or B), is there a usage matrix U?

VAFFP: Given *F*, find *U* and *B* such that *F* = *U B*



Given F and T (or B), is there a usage matrix U?

VAFFP: Given *F*, find *U* and *B* such that *F* = *U B*







Lemma (Ancestry Condition): Given **F** and **T**, for all samples *p* and mutations *k* child of *j*, $f_{pj} \ge f_{pk}$

necessal



Lemma (Ancestry Condition): Given **F** and **T**, for all samples *p* and mutations *k* child of *j*, $f_{pj} \ge f_{pk}$

Ancestry graph G = (V, A); given F

- Vertex for every mutation
- Edge $(j,k) \in A$ iff $f_{pj} \geq f_{pk}$

for all samples p



Lemma (Ancestry Condition): Given F and T, for all samples p and mutations *k* child of *j*, $f_{pj} \ge f_{pk}$

Ancestry graph G = (V, A); given F

- Vertex for every mutation
- Edge $(j,k) \in A$ iff $f_{pj} \geq f_{pk}$

for all samples p

Theorem 1:

T is a solution to the VAFFP if and only if **T** is a spanning tree of **G** satisfying the Sum Condition

Solving the VAFFP: ILP formulation

max $\sum x_{jk}$ **Find the largest** set of edges in G $(v_j, v_k) \in A'$ s.t. $\sum x_{rj} = 1$ **Exactly one root node** $v_i \in \delta^+(v_r)$ $x_{kl} \leq \sum x_{jk}$ $v_j \in \delta^-(v_k)$ $\sum \quad x_{jk} \le 1$ $v_j \in \delta^-(v_k)$ $v_j \in \delta^-(v_k)$ $v_l \in \delta^+(v_k)$



 $\forall (v_k, v_l) \in A$ Connectivity

 $\forall v_k \in V$ Tree

 $\sum f_{pk}x_{jk} \ge \sum f_{pl}x_{kl} \quad \forall p \in [m], v_k \in V$ Sum condition $x_{jk} \in \{0, 1\}$ $\forall (v_i, v_k) \in A'$



Biological Problem	Assumptions	Computational Problem
Mixed tumor samples require specialized phylogeny algorithms	 Infinite sites assumption Error-free data 	VAF Factorization Problem (VAFFP) Given F, find U and B such that F = U B

Approach

- Combinatorial characterization as constrained spanning trees in a directed acyclic graph
- Integer linear programming

Probabilistic Model for Noisy Measurements

Probabilistic Model for Noisy Measurements

AncesTree

AncesTree:

 Build approximate ancestry graph *G*

2. Find the largest tree **T** in the approximate ancestry graph **G** and matrix **F** that satisfy the sum condition for **T** and **F** (mixed integer linear programming)

3. Compute the usage matrix **U**

40

Solving the I-VAFFP: MILP formulation

41

Biological Problem	Assumptions	Computational Problem
Mixed tumor samples require specialized phylogeny algorithms	Infinite sites assumption	Interval VAF Factorization Problem (I-VAFFP) Given F^- and F^+ , find F , U and B such that $F = U B$

Approach

- Combinatorial characterization as constrained spanning trees in a directed acyclic graph
- Probabilistic model for noisy variant allele frequencies
- AncesTree : mixed integer linear programming

Results on Simulated Data

Generated:

- 3,4,5 node trees
- 20 Mutations
- 5 Samples
- Average coverage: 200x

Ran:

- AncesTree (Combinatorial)
- PhyloSub (Probabilistic) [Jiao et al., 2014]
- Canopy (Probabilistic) [Jiang et al., 2016]

4 nodes 5 nodes 3 nodes AncesTree PhyloSub Canopy 0.40.00.20.40.81.00.0 0.20.40.60.81.0).0 0.20.61.00.60.8Proportion of Ancestral Relationships Correctly Inferred Proportion of Ancestral Relationships Correctly Inferred Proportion of Ancestral Relationships Correctly Inferred

Runtime on Simulated Data

······ 1 minute

----- 1 hour

- 8 hours

Biological Problem	Assumptions	Computational Problem
Mixed tumor samples require specialized phylogeny algorithms	• Infinite sites assumption	Interval VAF Factorization Problem (I-VAFFP) Given F^- and F^+ , find F , U and B such that $F = U B$

Approach

- Combinatorial characterization as constrained spanning trees in a directed acyclic graph
- Probabilistic model for noisy variant allele frequencies
- AncesTree : mixed integer linear programming

Cell Division/Mutation and Migration are Separate Processes

Cell Division/Mutation and Migration are Separate Processes

Leaf-labeled Clone Tree T

 $\mu(T, \ell) = 8$ migrations

Given **T** and **e**, **migrations** are bichromatic edges in **T**, or edges in **G**

Cell Division/Mutation and Migration are Separate Processes

Minimum Migration Analysis in Ovarian Cancer

McPherson et al. (2016). Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nature Genetics*.

- Instance of the maximum parsimony small phylogeny problem
- Can be solved in polynomial time [Fitch, 1971; Sankoff, 1975]

Biological Problem	Assumptions	Computational Problem
Reconstructing migration history of a metastatic cancer	Migrations are rareMigrations are independent	Parsimonious Migration History: Given <i>T</i> , find vertex labeling <i>e</i> with minimum number of migrations

Approach

• Dynamic programming (Sankoff/Fitch algorithm)

Are there multiple vertex labelings with $\mu^* = 13$ migrations?

McPherson et al. (2016). Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nature Genetics*.

Minimum Migration History is Not Unique

• Enumerate all minimum-migration vertex labelings in the backtrace step

54

Comigrations: Simultaneous Migrations of Multiple Clones

- Multiple tumor cells migrate simultaneously through the blood stream [Cheung et al., 2016]
- Second objective: number γ of comigrations is the number of multi-edges in migration graph G^+

ROv

SBwl

Om

B2 SBwl B1 Om Right Ovary Small Bowel

Omentum

+ Not necessarily true in the case of directed cycles

Comigrations: Simultaneous Migrations of Multiple Clones

- Multiple tumor cells migrate simultaneously through the blood stream [Cheung et al., 2016]
- Second objective: number γ of comigrations is the number of multi-edges in migration graph G^+

Tradeoff between Migrations and Comigrations

• Minimum number γ^* of comigrations is m-1 (where m is #anatomical sites)

Tradeoff between Migrations, Comigrations and Migration Pattern

Parsimonious Migration History Problem

Parsimonious Migration History Problem: Given a clone tree T and a set \mathcal{P} of allowed migration patterns, find a vertex labeling ℓ with the minimum migration number $\mu^*(T)$ and subsequently the smallest comigration number $\hat{\gamma}(T)$.

Polytomy Resolution in Ovarian Cancer 7

McPherson et al. (2016). Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nature Genetics*.

Resolving Clone Tree Ambiguities

Applying MACHINA to Metastatic Breast Cancer

brain

seeding (pM)

kidney

brain

kidney

liver

rib

seeding (mS)

MACHINA accurately infers clone trees and migration histories on simulated data

Precise mathematical models are needed to describe the evolutionary process in cancer:

- Do not try to solve everything at once; it is OK to simplify and gradually add complexity
- Understanding combinatorial structure leads to a better understanding of the problem at hand
- This leads to better and efficient algorithms