Parsimonious Migration History Problem: Complexity and Algorithms

Mohammed El-Kebir, University of Illinois at Urbana-Champaign WABI 2018



Tumorigenesis: (i) Cell Division, (ii) Mutation & (iii) Migration



Tumorigenesis: (i) Cell Division, (ii) Mutation & (iii) Migration



Tumorigenesis: (i) Cell Division, (ii) Mutation & (iii) Migration



Goal: Given phylogenetic tree *T*, find *parsimonious* vertex labeling *e* with fewest migrations

Slatkin, M. and Maddison, W. P. (1989). A cladistic measure of gene flow inferred from the phylogenies of alleles. *Genetics*, 123(3), 603–613.

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 [Fitch, 1971; Sankoff, 1975]



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 [Fitch, 1971; Sankoff, 1975]



Minimum Migration History is Not Unique

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



6

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^+





RFTA

ROv

SBwl

Om

B2 SBwl B1 Om **Right Fallopian Tube**

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^+



Constrained Multi-objective Optimization Problem

Parsimonious Migration History (PMH): Given a phylogenetic tree T and a set $\mathcal{P} \subseteq \{S, M, R\}$ of allowed migration patterns, find vertex labeling ℓ with minimum migration number $\mu^*(T)$ and smallest comigration number $\hat{\gamma}(T)$.



El-Kebir, M., Satas, G., & Raphael, B. J. (2018). Inferring parsimonious migration histories for metastatic cancers. *Nature Genetics*, 50(5), 718–726.

Results

Parsimonious Migration History (PMH): Given a phylogenetic tree T and a set $\mathcal{P} \subseteq \{S, M, R\}$ of allowed migration patterns, find vertex labeling ℓ with minimum migration number $\mu^*(T)$ and smallest comigration number $\hat{\gamma}(T)$.



PMH is NP-hard when $\mathcal{P} = \{S\}$

3-SAT: Given $\varphi = \bigwedge_{i=1}^{k} (y_{i,1} \lor y_{i,2} \lor y_{i,3})$ with variables $\{x_1, \dots, x_n\}$ and k clauses, find $\varphi : [n] \rightarrow \{0,1\}$ satisfying φ



PMH is NP-hard when $\mathcal{P} = \{S\}$

3-SAT: Given $\varphi = \bigwedge_{i=1}^{k} (y_{i,1} \lor y_{i,2} \lor y_{i,3})$ with variables $\{x_1, \dots, x_n\}$ and k clauses, find $\varphi : [n] \rightarrow \{0,1\}$ satisfying φ

Three ideas:

- 1. Ensure that $(x, \neg x) \in E(G)$ or $(\neg x, x) \in E(G)$
- 2. Ensure that $\ell^*(r(T)) = \bot$
- 3. Ensure that ϕ is satisfiable if and only if ℓ^* encodes a satisfying truth assignment





PMH is NP-hard when $\mathcal{P} = \{S\}$

3-SAT: Given $\varphi = \bigwedge_{i=1}^{k} (y_{i,1} \lor y_{i,2} \lor y_{i,3})$ with variables $\{x_1, \dots, x_n\}$ and k clauses, find $\varphi : [n] \rightarrow \{0,1\}$ satisfying φ

Three ideas:

- 1. Ensure that $(x, \neg x) \in E(G)$ or $(\neg x, x) \in E(G)$
- 2. Ensure that $\ell^*(r(T)) = \bot$
- 3. Ensure that ϕ is satisfiable if and only if ℓ^* encodes a satisfying truth assignment





Lemma: Let B > 10k + 1 and A > 2Bn + 27k. Then, φ is satisfiable if and only if $\mu^*(T) = (B + 1)n + 25k$



Lemma: Let B > 10k + 1 and A > 2Bn + 27k. Then, φ is satisfiable if and only if $\mu^*(T) = (B + 1)n + 25k$

PMH is FPT in number m of locations when $\mathcal{P} = \{S\}$



$$d_T(u,v) \ge d_{\hat{G}}(\operatorname{lca}_{\hat{G}}(u),\ell(v)) \qquad \forall u,v \in V(T) \text{ such that } u \preceq_T v.$$

$$\ell^*(v) = \begin{cases} \operatorname{LCA}_{\hat{G}}(r(T)), & \text{if } v = r(T), \\ \sigma(\ell^*(\pi(v)), \operatorname{LCA}_{\hat{G}}(v)), & \text{if } v \neq r(T), \end{cases}$$

where $\sigma(s,t) = s$ if s = t and otherwise $\sigma(s,t)$ is the unique child of s that lies on the path from s to t in \hat{G} .

Lemma: If (1) holds then ℓ^* is a minimum migration labeling consistent with \widehat{G} .



where $\sigma(s,t) = s$ if s = t and otherwise $\sigma(s,t)$ is the unique child of s that lies on the path from s to t in \hat{G} .

Lemma: If (1) holds then ℓ^* is a minimum migration labeling consistent with \widehat{G} .

17

Simulations



Available on: https://github.com/elkebir-group/PMH-S

Conclusions

- PMH is NP-hard when restricted to a multi-tree
- PMH is FPT in *m* when restricted to a multi-tree
- FPT algorithm is practical for modest number *m* of locations

Discussion:

- Applicable to other domains, i.e. comigration of bacterial/viral strains through a single transmission event
- Polytomy resolution version is hard as well
- Open question: Hardness when *G* is unrestricted or restricted to a DAG?





