# Analysis of the Genealogy Process in Investigative Genetic Genealogy

Lawrence M. Wein

Mine Su Erturk

Graduate School of Business, Stanford University

1

---

## PREVIOUS CRIME-SOLVING RESEARCH

US-VISIT Program
- *PNAS* (2005) and Congressional testimony (2005)
- Led to switch from 2 to 10 fingers

UIDAI (Aadhaar) Program in India
- *PLOS ONE* (2014) and India Supreme Court briefing (2014)

Ballistic imaging
- *Journal of Forensic Sciences* 2014, 2017
- *AFTE Journal* 2018 (Stockton CA PD)

Sexual assault kits
- Testing the backlog: *Journal of Forensic Sciences* 2018
  *cnn.com* 2018
- How to test the backlog: *PNAS* 2020 (SFPD)

2

---

## ACKNOWLEDGMENTS

Colleen Fitzpatrick
IdentiFinders

Margaret Press
DNA Doe Project

3

## OUTLINE

Introduction: research questions and data

Main results

Parameter estimation

Modeling the two-stage genealogy process

Proposed strategy

Limitations

Conclusions

4

## STEPS IN THE IGG PROCESS

1. Obtain DNA sample from crime scene or unidentified remains and perform genotyping (SNPs)

5

## STEPS IN THE IGG PROCESS

1. Obtain DNA sample from crime scene or unidentified remains and perform genotyping (SNPs)

2. Upload SNP data to third-party service to obtain relatives of target    [GED match]   FamilyTreeDNA

6

## STEPS IN THE IGG PROCESS

1. Obtain DNA sample from crime scene or unidentified remains and perform genotyping (SNPs)

2. Upload SNP data to third-party service to obtain relatives of target  [GED match]  FamilyTree**DNA**

➡ 3. Perform genealogy research to identify target

---

## STEPS IN THE IGG PROCESS

1. Obtain DNA sample from crime scene or unidentified remains and perform genotyping (SNPs)

2. Upload SNP data to third-party service to obtain relatives of target  [GED match]  FamilyTree**DNA**

➡ 3. Perform genealogy research to identify target

4. Obtain confirmatory DNA sample from identified target or family member

---

## RESEARCH QUESTIONS

First mathematical analysis of the backend of the IGG process
(i.e., GEDmatch/FTDNA output → identify target)

## RESEARCH QUESTIONS

First mathematical analysis of the backend of the IGG process
(i.e., GEDmatch/FTDNA output → identify target)

Performance analysis:
Given the GEDmatch/FTDNA output, compute:
Probability of identifying target
Expected workload (= size of final family tree)

## RESEARCH QUESTIONS

First mathematical analysis of the backend of the IGG process
(i.e., GEDmatch/FTDNA output → identify target)

Performance analysis:
Given the GEDmatch/FTDNA output, compute:
Probability of identifying target
Expected workload (= size of final family tree)

Optimization:
How many and which matches to investigate?
When/if to descend from (possible) MRCAs
(MRCA = Most Recent Common Ancestor)

## 17 CASES FROM DNA DOE PROJECT

| Case Number $i$ | Case Solved $p^i$ | Number of Available Matches | Number of Investigated Matches | Number of Identified Matches |
|---|---|---|---|---|
| 1 | Yes | 2134 | 31 | 10 |
| 2 | No | 545 | 246 | 25 |
| 3 | No | 4257 | 50 | 41 |
| 4 | No | 509 | 86 | 22 |
| 5 | Yes | 2000 | 56 | 20 |
| 6 | No | 1373 | 232 | 171 |
| 7 | Yes | 633 | 34 | 29 |
| 8 | Yes | 795 | 72 | 68 |
| 9 | No | 928 | 56 | 37 |
| 10 | No | 5059 | 288 | 212 |
| 11 | No | 5007 | 39 | 24 |
| 12 | Yes | 313 | 107 | 44 |
| 13 | Yes | 2136 | 80 | 66 |
| 14 | Yes | 2417 | 199 | 149 |
| 15 | Yes | 308 | 74 | 54 |
| 16 | No | 221 | 31 | 29 |
| 17 | No | 610 | 72 | 43 |
| Total | 8 | 29245 | 1753 | 1044 |

13



14



15

16



AUTO CLUSTER TOOL IN GEDmatch

17

## OUTLINE

Introduction: research questions and data

➡ Main results

Parameter estimation

Modeling the two-stage genealogy process

Proposed strategy

Limitations

Conclusions

18

## BENCHMARK  STRATEGY

Search for MRCAs between two matches,
   and immediately descend from these MRCAs

19

## BENCHMARK  STRATEGY

Search for MRCAs between two matches,
   and immediately descend from these MRCAs


Investigate n matches prioritized by highest total cM

20

## BENCHMARK  STRATEGY

Search for MRCAs between two matches,
   and immediately descend from these MRCAs


Investigate n matches prioritized by highest total cM


Vary n to generate Pr(identify target) vs. E[workload] curve

21

## BENCHMARK STRATEGY



22

## BENCHMARK STRATEGY



Pr(identify target) goes to 1 as E[workload] goes to 30,000
    i.e., as number of investigated matches goes to 300

23

## BENCHMARK STRATEGY



Pr(identify target) goes to 1 as E[workload] goes to 30,000
    i.e., as number of investigated matches goes to 300
Decreasing returns after E[workload] = 13,000
    i.e., after number of investigated matches = 50

24

## PROPOSED vs. BENCHMARK STRATEGY



It solves cases much more quickly: at E[workload] = 2000,
it solves 71% of cases vs. 27% for Benchmark

## OUTLINE

Introduction: research questions and data

Main results

➡ Parameter estimation

Modeling the two-stage genealogy process

Proposed strategy

Limitations

Conclusions

## PARAMETERS

Pr(can identify a match) = 0.59

## PARAMETERS

Pr(can identify a match) = 0.59

Pr(can identify someone's spouse) = 1

Pr(can identify someone's child) = 0.98 (also considered 0.90)

Pr(can identify someone's parents) = 0.60  (by simulation)

28

## PARAMETERS

Pr(can identify a match) = 0.59

Pr(can identify someone's spouse) = 1

Pr(can identify someone's child) = 0.98 (also considered 0.90)

Pr(can identify someone's parents) = 0.60  (by simulation)

Number of children per couple

29

## NUMBER OF CHILDREN PER COUPLE



30

## OUTLINE

Introduction: research questions and data

Main results

Parameter estimation

→ Modeling the two-stage genealogy process

Proposed strategy

Limitations

Conclusions

31

---

## ASCENDING STAGE

Given: list of GEDmatch/FTDNA matches to investigate

Ascending: build family tree up (backwards in time) from matches

32

---

## ASCENDING STAGE

Given: list of GEDmatch/FTDNA matches to investigate

Ascending: build family tree up (backwards in time) from matches

Goal: Find Most Recent Common Ancestors (MRCAs) between target and each match

33

## MOST RECENT COMMON ANCESTORS



2nd cousin      Target

---

## MOST RECENT COMMON ANCESTORS



2nd cousin      Target

Cluster = ancestral couple of target
7 clusters in this example

---

## ASCENDING STAGE

Given: list of GEDmatch/FTDNA matches to investigate

Ascending: build family tree up (backwards in time) from matches

Goal: Find Most Recent Common Ancestors (MRCAs) between target and each match

State of system during ascending stage:

For each generation g and cluster $c = 1,\ldots,2^{g-1}$

$L_{g,c}$ = number of possible MRCAs identified

$P_{g,c}$ = Pr(one of the $L_{g,c}$ MRCAs is the correct MRCA)

## MOST RECENT COMMON ANCESTORS



2nd cousin      Target

List size L = 0, 1, 2, 3 or 4

37

## MOST RECENT COMMON ANCESTORS



2nd cousin      Target

List size L = 0, 1, 2, 3 or 4
P=0 if L=0, and P=1 if L=4

38

## MOST RECENT COMMON ANCESTORS



2nd cousin      Target

List size L = 0, 1, 2, 3 or 4
P=0 if L=0, and P=1 if L=4
Other P's ≠ 1/4, 1/2, 3/4

39

## MOST RECENT COMMON ANCESTORS



2nd cousin    1st cousin    Target

40

## MOST RECENT COMMON ANCESTORS



2nd cousin    1st cousin    Target

Now state changes from (L=4,P=1) to (L=1,P=1)

41

## DESCENDING  STAGE

Given: State of system during ascending stage:

For each generation g and cluster c = 1,...,$2^{g-1}$

$L_{g,c}$ = number of possible MRCAs identified

$P_{g,c}$ = Pr(one of the $L_{g,c}$ MRCAs is the correct MRCA)

Descending: build family tree down (forwards in time) from possible MRCAs between target and match

Goal: Find intersection of (i.e., marriage between) maternal and paternal family trees

42

**FIND INTERSECTION OF FAMILY TREES**

2nd cousin    Target    1st cousin

43



**FIND INTERSECTION OF FAMILY TREES**

2nd cousin    Target    1st cousin

44

## DESCENDING STAGE

Given: State of system at end of ascending stage:

For each generation g and cluster c = 1,…,$2^{g-1}$

$L_{g,c}$ = number of possible MRCAs identified

$P_{g,c}$ = Pr(one of the $L_{g,c}$ MRCAs is the correct MRCA)

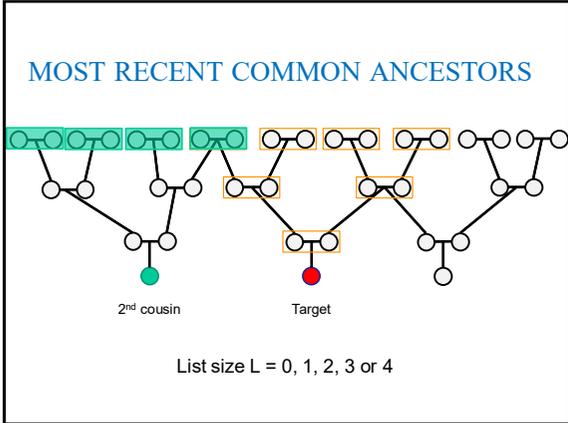Descending: build family tree down (forwards in time) from possible MRCAs between target and match

Goal: Find intersection of (i.e., marriage between) maternal and paternal family trees

Compute: 1) probability of finding intersection of family trees
          2) expected workload

45

## OUTLINE

Introduction: research questions and data

Main results

Parameter estimation

Modeling the two-stage genealogy process

➡️ Proposed strategy

Limitations

Conclusions

46

## PROPOSED vs. BENCHMARK STRATEGY



It solves cases much more quickly: at E[workload] = 2000, it solves 71% of cases vs. 27% for Benchmark

47

## STOCHASTIC DYNAMIC PROGRAMMING

Observe state, take action, observe probabilistic transition to new state, take new action,…to maximize objective

48

## STOCHASTIC DYNAMIC PROGRAMMING

Observe state, take action, observe probabilistic transition to new state, take new action,…to maximize objective

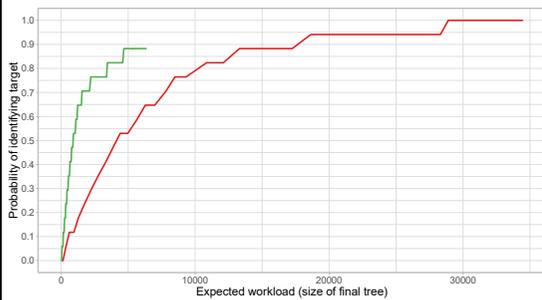State: 1) list of uninvestigated matches with cM and cluster
      2) for each cluster and generation, Pr(list contains correct MRCA of target) and size of list of MRCAs
      3) list of cluster-generation pairs for which a descending search has been performed

49

---

## STOCHASTIC DYNAMIC PROGRAMMING

Observe state, take action, observe probabilistic transition to new state, take new action,…to maximize objective

State: 1) list of uninvestigated matches with cM and cluster
      2) for each cluster and generation, Pr(list contains correct MRCA of target) and size of list of MRCAs
      3) list of cluster-generation pairs for which a descending search has been performed

Actions: 1) start an ascending search of a new match with cM and cluster
       2) start a descending search from cluster-generation
       3) end the search

50

---

## STOCHASTIC DYNAMIC PROGRAMMING

Observe state, take action, observe probabilistic transition to new state, take new action,…to maximize objective

State: 1) list of uninvestigated matches with cM and cluster
      2) for each cluster and generation, Pr(list contains correct MRCA of target) and size of list of MRCAs
      3) list of cluster-generation pairs for which a descending search has been performed

Actions: 1) start an ascending search of a new match with cM and cluster
       2) start a descending search from cluster-generation
       3) end the search

Objective: maximize Pr(target identified) – (cost x E[workload])

51

## PROPOSED STRATEGY

SDP problem was too hard to solve (huge state space)

Use a greedy (ie, myopic) approach

52

## PROPOSED STRATEGY

SDP problem was too hard to solve (huge state space)

Use a greedy (ie, myopic) approach

Given current state:
1) for a descending search, compute $\Delta P$ (increase in probability of identifying target) and $\Delta W$ (increase in E[workload]) for each undescended cluster-generation

53

## PROPOSED STRATEGY

SDP problem was too hard to solve (huge state space)

Use a greedy (ie, myopic) approach

Given current state:
1) for a descending search, compute $\Delta P$ (increase in probability of identifying target) and $\Delta W$ (increase in E[workload]) for each undescended cluster-generation

2) for an ascending search, compute $\Delta P$ and $\Delta W$ for each uninvestigated match, assuming that we descend right after ascending

54

## PROPOSED STRATEGY

SDP problem was too hard to solve (huge state space)

Use a greedy (ie, myopic) approach

Given current state:
1) for a descending search, compute $\Delta P$ (increase in probability of identifying target) and $\Delta W$ (increase in E[workload]) for each undescended cluster-generation

2) for an ascending search, compute $\Delta P$ and $\Delta W$ for each uninvestigated match, assuming that we descend right after ascending

3) compute the marginal increase in the objective (ie, $\Delta P -$ cost x $\Delta W$) from 1) and 2) and choose the action with the maximum increase (and stop after n investigations)

55

## PROPOSED STRATEGY

SDP problem was too hard to solve (huge state space)

Use a greedy (ie, myopic) approach

Given current state:
1) for a descending search, compute $\Delta P$ (increase in probability of identifying target) and $\Delta W$ (increase in E[workload]) for each undescended cluster-generation

2) for an ascending search, compute $\Delta P$ and $\Delta W$ for each uninvestigated match, assuming that we descend right after ascending

3) compute the marginal increase in the objective (ie, $\Delta P -$ cost x $\Delta W$) from 1) and 2) and choose the action with the maximum increase (and stop after n investigations)
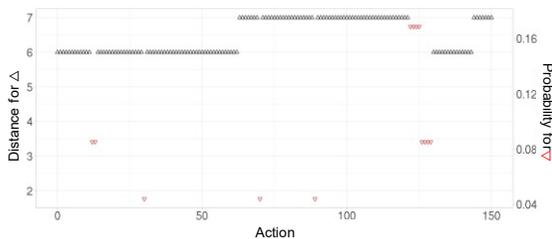
Vary n to generate Pr(target identified) vs. E[workload] curve

56

## PROPOSED STRATEGY FOR CASE #1



Closest matches in Case #1 are at distance 6
△ = ascending from distance (left vertical axis)
▽ = descending from cluster-generation with probability of true MRCA (right vertical axis)

57

## OUTLINE

Introduction: research questions and data

Main results

Parameter estimation

Modeling the two-stage genealogy process

Proposed strategy

➡ Limitations

Conclusions

58

## LIMITATIONS

Cases: small sample size and not chosen randomly

59

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

60

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

No half-relationships

61

---

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

No half-relationships

No geographical information

62

---

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

No half-relationships

No geographical information

No ethnicity information

63

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

No half-relationships

No geographical information

No ethnicity information

AutoCluster information is perfect

64

---

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

No half-relationships

No geographical information

No ethnicity information

AutoCluster information is perfect

No Y-STR data to infer surname (Gymrek, *Science* 2013)

65

---

## LIMITATIONS

Cases: small sample size and not chosen randomly

No endogamy

No half-relationships

No geographical information

No ethnicity information

AutoCluster information is perfect

No Y-STR data to infer surname (Gymrek, *Science* 2013)

Search probabilities do not depend on generation

66

## OUTLINE

Introduction: research questions and data

Main results

Parameter estimation

Modeling the two-stage genealogy process

Proposed strategy

Limitations

➡ Conclusions

67

---

## CONCLUSIONS

Pr(identify target) and E[workload] are useful only in relative terms
- But police departments and IGG companies need to
  assess solvability and workload upfront

68

---

## CONCLUSIONS

Pr(identify target) and E[workload] are useful only in relative terms
- But police departments and IGG companies need to
  assess solvability and workload upfront

Hard cases appear to be solvable but require high workload
- Tradeoff curves allow for identification of sweet spot

69

## CONCLUSIONS

Pr(identify target) and E[workload] are useful only in relative terms
- But police departments and IGG companies need to
  assess solvability and workload upfront

Hard cases appear to be solvable but require high workload
- Tradeoff curves allow for identification of sweet spot

Proposed Strategy solves cases faster by:
- looking for MRCAs between the target and each match
- aggressively descending from possible MRCAs

70

## CONCLUSIONS

Pr(identify target) and E[workload] are useful only in relative terms
- But police departments and IGG companies need to
  assess solvability and workload upfront

Hard cases appear to be solvable but require high workload
- Tradeoff curves allow for identification of sweet spot

Proposed Strategy solves cases faster by:
- looking for MRCAs between the target and each match
- aggressively descending from possible MRCAs

Proposed Strategy is meant to aid, not to replace, genealogists'
decisions

71