

サブグループ同定について

水田 正弘 大学統計教員育成センター 特任教授

1. サブグループ解析とサブグループ同定

最近、「本医薬品は、50歳以上で過度な飲酒をする人に対する効果が低い」など、「特定の属性」を有する人の薬効などに関する報告を耳にする。このような主張を統計的に証明する問題は比較的明確に設定できる。これは、サブグループ解析として、個別化医療の研究の一環で注目されている。しかし、「特定の属性」を見つけ出すのは、別の課題であり、サブグループ同定法(Subgroup Identification)と呼ばれている。サブグループ同定法として、複数の手法が提案されているが、想定するサブグループの族の中で、効果(Treatment effect)が最大となるサブグループが導出される保証は全くない。そこで、想定する全てのサブグループを列挙し、それぞれの効果を計算することで、「厳密な意味で最適なサブグループ」を決定する方法を研究している。

2. 「良い」サブグループとは？

ある属性を有するサブグループにおいて、治療を受けた人の、受けない人に対する効果(Treatment effect)の差が、その属性を有しないサブグループにおける効果の差と比べて高い場合、「良い」サブグループと言える。このこと自体は、差の検定やGLM、LASSOなどで調整した効果の差として、統計的な検定やベイズ的アプローチにより評価できる。また、Pmtree (model4you.Compute model-based tree from model)、glmtree (partykit. Logistic Regression Trees)、QUINT (Qualitative INteraction Trees)などが「良い」サブグループを同定する手法(ソフトウェア)として発表されている。

3. 最適なサブグループの同定法

サブグループ同定の問題を、想定したサブグループの属の中から、治療等の効果が最も高いものを導き出すと設定することができる。先に紹介した既存の手法は、残念ながら近似最適解の検出に過ぎない。そこで、想定する全てのサブグループについて、効果を計算することで、(近似ではない)最適な解を同定することを試みた。2020年の統計関連学会連合大会では、 $p=5, n=100$ の人工データにおいて、43,199,128,758個(432億個)のサブグループが列挙されたことを報告した。しかし、実際のデータ解析を想定した場合、サンプルサイズ n 、変数数 p ともに少なすぎる。そこで、データ構造の特徴抽出を実施するなどにより、適用可能な状況の拡大および計算時間の短縮方法を検討している。

4. まとめ

病気になったとき、その治療法の選択は人生において大きな岐路となることが多い。すなわち、1つの疾病に対して、複数の治療法があるとき、どれを選ぶかは重大な問題である。特定の治療法が全ての人に有効であるとは限らないので、個別化医療が注目されている。これの基本として、特定の治療方法の有効性を集団として検討するサブグループ解析がある。ある治療方法が有効であるサブグループと、そうではないサブグループが存在することは自然である。適切な「特定の属性」を設定し、治療の対象者を限定すれば、良好な治療成績が期待できる。例えば、COVID-19関連で、一般化した表現に、「特定の年齢の人に対するワクチン接種の推奨」がある。すなわち、あるグループに対してワクチンは損失よりも利得が大きいという意味である。同様に、ある療法が有効なグループの特定、ある政策に対する賛成が多いグループの特定などは重要な事項であり、データサイエンスにおいて、普遍的な課題である。サブグループを与えたとき、療法が有効や賛成が多いなどの数値的な評価は、統計的または社会的に定義できることが多い。そこで、最適なグループの同定法、すなわちサブグループ同定法が研究されてきた。これにより、近似的に最適なグループを見つける方法は、いくつか開発されているが、最適解が得られる保証はない。厳密な意味での最適解の導出法があれば、それ自体、有効な手法であるとともに、これまでに提案されてきたサブグループ同定法の手法を評価することができる。

参考文献

- 水田正弘: Subgroup identificationのための全列挙について, 2020年度統計関連学会連合大会講演報告集
- 水田正弘: Subgroup Identificationにおける条件設定と最適解について, 2021年度統計関連学会連合大会講演報告集

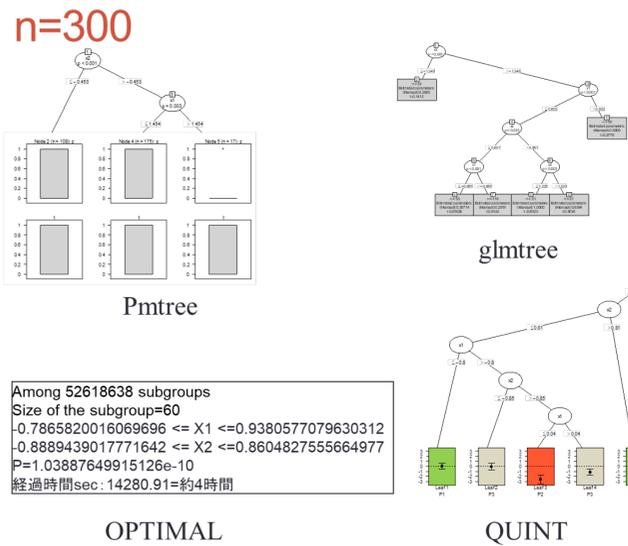


図1. 人工データ($P=2, n=300$)に対するPmtree, glmtree, QUINT, OPTIMALの結果

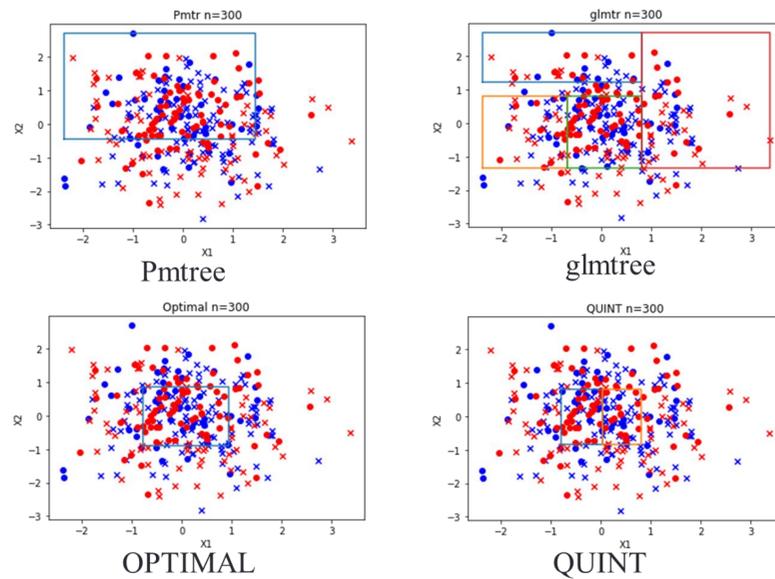


図2. 人工データ($P=2, n=300$)に対するPmtree, glmtree, QUINT, OPTIMALの結果の図示

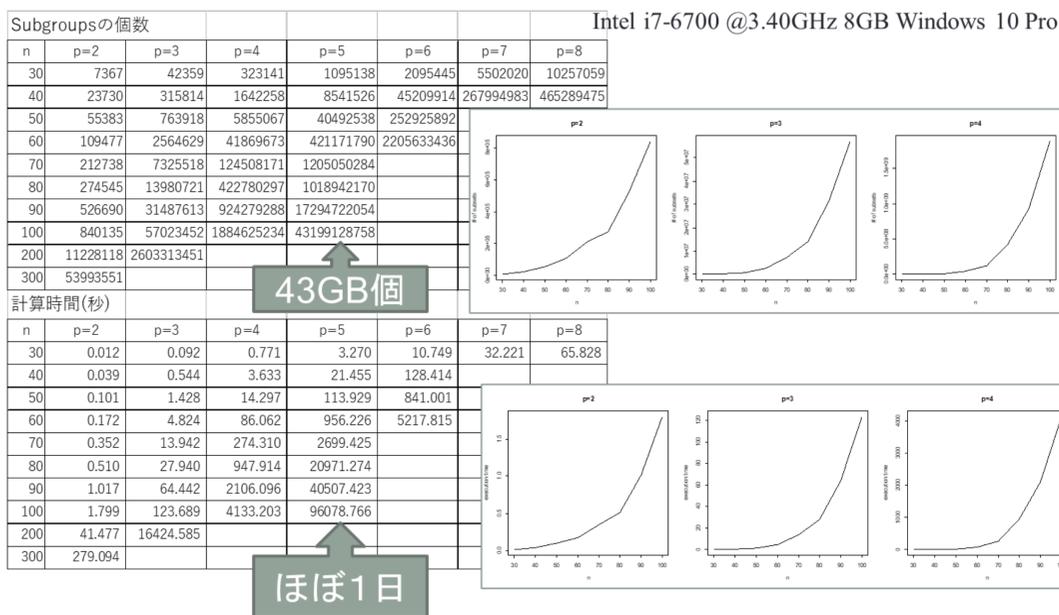


図3. 列挙の個数と計算時間