

## 医療介護とデータサイエンスの新たな視点

慶應義塾大学	山内 恒人 君
スイス再保険	藤澤 陽介 君
慶應義塾大学	小林 凌雅 氏
慶應義塾大学	小暮 厚之 氏

### 山内

それでは、今日はこれから「医療介護とデータサイエンスの新たな視点」ということで、1時間半にわたるプレゼンテーションを行いたいと考えております。データサイエンスに関しましては、この場で何もいろいろと申し上げるまでもなく、今後のアクチュアリーを支える、一つの大きな柱になるだろうと思います。その中で、やはりコーディングスキルや、あるいは、その統計に関する最近の話題というものにキャッチアップするというのは、なかなか機会がないと難しいということもございまして、最先端の話だけではなくて、もちろん今日はかなりその断片も見えるわけですが、われわれが学んできたことも含めて、どのようにデータサイエンスという形で新たな進化が得られているのかということについて、3名の方からお話を伺うという予定にしております。詳細はそれぞれのお話の中で語られますので、この場では申し上げることはいたしません。

早速でございますけれども、その3名の方々をこれからご紹介したいと思います。最初はスイス再保険会社日本支店シニアヘルスソリューションズマネージャー、バイスプレジデントでいらっしゃいます藤澤陽介さんです。よろしくお願いいたします。お次でございますが、慶應義塾大学政策メディア研究科後期博士課程に在学中でいらっしゃいます、小林凌雅さんです。よろしくお願いいたします。最後でございますけれども、慶應義塾大学総合政策学部の教授でいらっしゃいます、小暮厚之先生です。それでは、早速、藤澤さんから、よろしくお願いいたします。

**藤澤** ご紹介ありがとうございます。スイス再保険会社で働いています、藤澤と申します。本日はよろしくお願いいたします。

## 医療介護とデータサイエンスの 新たな視点

スイス再保険  
藤澤陽介

私のパートは導入部分です。これまでも何度か小暮先生や小林さんと意見交換させていただいたのですが、私も一人のオーディエンスの気持ちで、早く先生方のプレゼンを聞きたいと思っています。私の役目はこのパートをなるべく早く終わらせて、次のプレゼンを、オーディエンスになった気持ちで楽しみたいと思っています。

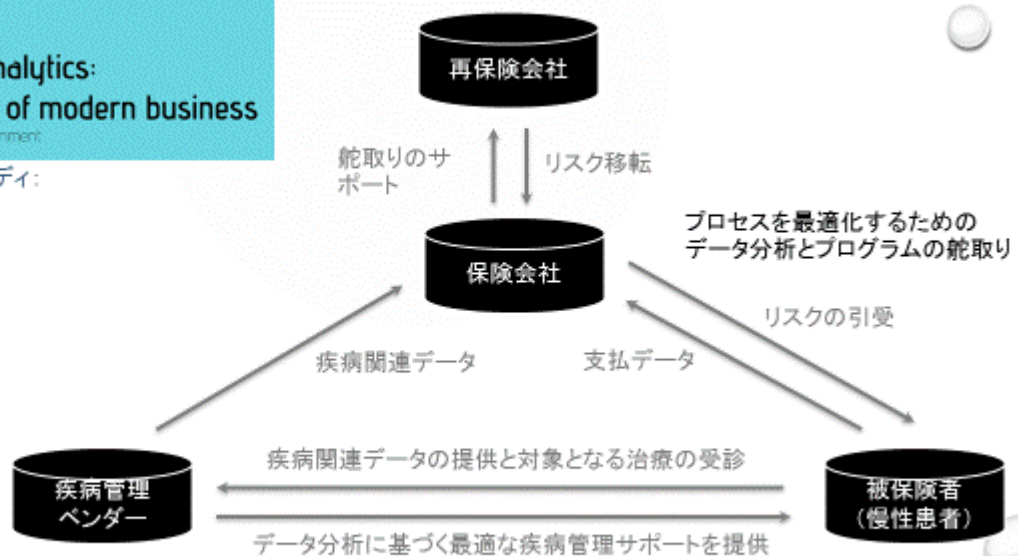
CRO FORUM

### Big Data & Analytics: the algorithm of modern business

CROs in a changing environment

2つのケーススタディ:

- 医療保険
- テレマティクス

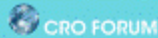


※CROフォーラム Big Data & Analyticsをもとに作成

最初は導入ということで、実務的な視点の話をしたと思って作ったスライドです。CROフォーラムというところがございます。ここが幾つか、外部環境の変化を踏まえて、今後どのようなリスク管理が必要になるのかということレポートにまとめています。その中の一つがこのビッグデータとアナリティクスというテーマでございます。このペーパーの中に二つのケーススタディが載ってまして、損保で言うとテレマティクスが有名だと思うのですが、生保、もしくは医療の世界だと、医療保険のところで

ちょっと実務的なその動向が変わりつつあるということが紹介されています。

一般的に保険会社は被保険者、リスクを引き受けて、そのデータを入手するという形のビジネスをしているケースが多いと思います。ここに疾病管理のベンダーが入ってきて、例えば被保険者の中でも典型的なものは糖尿病の患者さんですが、糖尿病の患者さんが合併症を予防するためにどのようなことができるのかという形で疾病管理のベンダーが入ってきて、どちらかと言うとパーソナライズした疾病管理のサポートをする。そこでまた新たなデータが集まってくる。これを保険会社が、ここでそのままつないで良いのかという法的な問題もあろうかと思いますが、これが例えば遺伝子になると、倫理的に本当につないで良いのかという問題も当然出てくるかと思いますが、そのようなデータを仮につなぐことができた場合に、これまでいわゆる保険会社でやっていた経験値分析に新たな変数が入ってくるという状況になり、使える統計的なスキルは少しずつ変わってくるのではないかという流れになっています。それを再保険会社がバックでサポートしているという事例が、幾つか海外の方で出てきています。



## Big Data & Analytics: the algorithm of modern business

LRDs in a changing environment

リスク管理には、これから多くの専門分野にわたるスキルが求められる。

- アクチュアリー能力
- データサイエンスのスキル

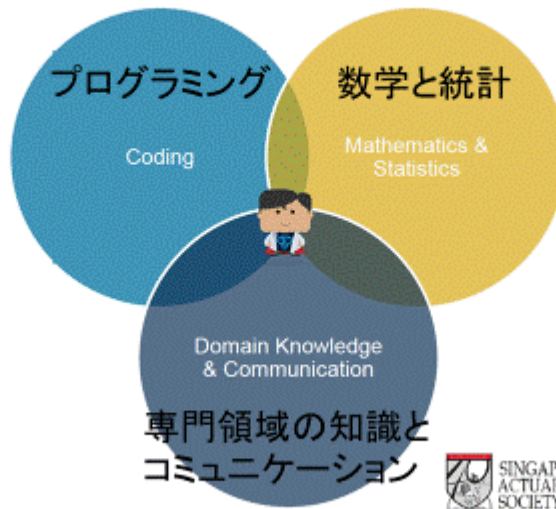
IFoA  
(英国)

■ 2017/9月に、主要組織に声を掛け、データ・サイエンス・サミットを開催 (IAJからも参加)

1. テクニカル・セッション (英国、カナダ、アイルランド)
2. 教育セッション (SOA、フランス、ドイツ、英国、CAS、シンガポール)
3. ビジネス開発セッション (オーストラリア、南アフリカ)
4. 外部環境 (コンサル、英国王立統計学会、プライバシー保護)
5. データ・サイエンスに関する国際協力とネクストステップ

このペーパーの中に書いているのですけれども、リスク管理という観点で、アクチュアリー能力のスキルはもちろん必要だと思いますが、データサイエンスのスキルを、今後磨いていかないといけないということも書かれています。このようなことを背景に、これは9月に開催されたデータサイエンスの集中セミナーでも少しお話ししたのですが、イギリスのアクチュアリー会が9月にデータサイエンスのサミットを開催しています。その導入ということで、一部分、ここは教育のセッションの、シンガポールの方の話を紹介させていただこうと思っています。

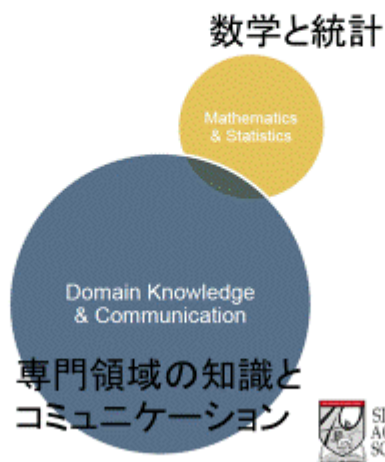
## What is a Data Scientist?



(出典) IFoA Global Data Science Summit, Colin Priest (Data Robot)

これは、データサイエンスはどのようなスキルが必要になるのか。これもよく出てくる図なのですが、プログラミングと、あと数学と統計、そして専門領域の知識とコミュニケーションが必要になる。これはシンガポール・アクチュアリー会の方のプレゼンを抜粋させていただいたものです。一方、アクチュアリーはどのような領域の教育を受けてきたのかという図が、こちらです。これも多分国によって違って、その場でもディスカッションがあったのですが、「やはり国によって、もっと数学と統計が強い国もあるのではないか」というような意見なども、サミットの中では出ていました。

## What is an Actuary Taught?



(出典) IFoA Global Data Science Summit, Colin Priest (Data Robot)

例えば、日本で言うと、多分数学的なスキルは海外のアクチュアリーに比べると、もうちょっと高いと思っています。このバブルがもう少し広いのかもしれませんが、統計の場合、なかなか日本で統計教育を受けるような機会はなかったり、現代的な統計技術が増える中で、そのような技術にキャッチアップして

いく教育がちょっと足りないのではないかという話や、あと、プログラミングを学ぶ機会がやはりなくて、これは各国のアクチュアリー会として、そのような教材であったり、勉強の機会を提供していく必要があるという話がありました。

## プログラミング

- 不足しているスキル

- HOW TO JOIN DATA TABLES
- SQL FOR DATA MANIPULATION
- R OR PYTHON SCRIPTS FOR MODELLING
- OBJECT ORIENTED DESIGN

データベース言語

RやPythonでモデリング  
(オブジェクト指向)

- ACTUARIES USE EXCEL FOR EVERYTHING, NO MATTER HOW SUITABLE IT IS, OR ISN'T, FOR A PARTICULAR TASK

つまり、アクチュアリーはエクセルが大好き

※IFoA Global Data Science Summit, Colin Priest (Data Robot)をもとに作成

例えば、プログラミングという観点で、不足しているスキルとして、ちょっと英語のままで恐縮なのですけれども、例えばSQLやデータベースの取り扱いなどです。RやPythonでのプログラミングのスキルというのは、やはり絶対的に不足しているという話がありました。下にも書いているのが、日本語に訳すとインパクトがありすぎると思って、あえて訳していないのですけれども、やはりアクチュアリーはエクセルで何でも実装してしまう。マクロを書いて、VBAを書いて、シミュレーションをするようなところが、特徴としてあるのではないかという話がございました。

# 数学と統計

- 狭い領域にフォーカス

- COMPOUND INTEREST AND FINANCIAL ECONOMICS
- LIFE CONTINGENCY TABLE MANIPULATION
- GENERALIZED LINER MODELS

所謂、生保数理(特に計算基数)

一般化線形モデル(GLM)

- 不足しているスキル

- 最適化
- 欠測値のインピュテーション
- 機械学習のアルゴリズム
- 数値解析による推定
- オペレーションズ・リサーチ

※IFoA Global Data Science Summit, Colin Priest (Data Robot)をもとに作成

次に、数学と統計の部分ですけれども、これもその場で議論になっていたのは、やはり生保数理の計算基数は、狭い領域にフォーカスしすぎているのではないかという話や、一般化線形モデルについても、これも9月の集中セミナーの時にも少し話があったのですが、やはり真のモデルというものが分からない中で、このGLMは説明がとてもしやすいということもあって、過度に使いすぎているのではないかというような意見が出ていました。

不足しているスキルとしましては、最適化や、欠損値があったときにどのように補完するのか、インピュテーションするのか。機械学習、数値解析、オペレーションズ・リサーチなど、このような分野を学ぶ機会はなかなかないのですけれども、データサイエンスの世界では、このようなスキルも必要になってくるといような話がありました。

## 専門領域の知識とコミュニケーション

### ・ 強み

- HOW INSURANCE WORKS
- LEGISLATION AND REGULATION
- ENTERPRISE RISK MANAGEMENT

所謂、日本のアクチュアリー試験とCERAの範囲

### ・ プロフェッショナル・ジャッジメントへの過度な依存

- NOT MUCH DIFFERENCE TO GUESSING
- UNABLE TO EXPLAIN ALGORITHMS TO NON-TECHNICAL CLIENTS
- POOR UNDERSTANDING OF UNCERTAINTY AND MODEL RISK

ここでの"guess"は「(根拠のない)推測」のニュアンス

アクチュアリー以外の  
人への説明力

不確実性とモデルリスク

※IFoA Global Data Science Summit, Colin Priest (Data Robot)をもとに作成

あと、専門領域の知識とコミュニケーションのところですけども、ここは強い部分ということで、保険の仕組みや法律や規制などはアクチュアリー試験にも入っていますし、その下のエンタープライズ・リスク・マネジメントも、CERAの範囲に含まれていたり、アクチュアリーの人でもリスク管理をやっている人が増えてきているので、ここは一つの強みなのではないかと言われていました。

一方で、ここも国によって違うと思いますが、プロフェッショナル・ジャッジメントに過度に依存しているのではないかと。日本の場合は、個人的にはそこまでジャッジせずに、なるべくデータを探してきて、プライシングするというカルチャーが多分あると思っていますが、海外のアクチュアリーのプライシングを見ていると、結構ジャッジに頼っているところもあります。guess はいろいろな意味がありますが、根拠のない推測に近いようなニュアンスで、ジャッジをしている人がいるのではないかという意見も出ていました。

あとは、説明力のコミュニケーションのところ、これもよく言われるところですけども。あとは不確実性やモデルリスクなど、そのような部分に対する理解が少し不足しているかもしれないというような話でした。

## 以下の分野の知識拡充



- プログラミング
  - RとPYTHON
- 数学と統計
  - 罰則付回帰とスパースモデリング(LASSO)
  - ベイズ統計
- 専門領域の知識とコミュニケーション
  - モデルリスクとパラメータリスク

そのようなことを背景に、今回小暮先生と小林さんをお願いして、主にそのような話をさせていただきたいと、私自身をお願いしたのですけれども、その以下の部分の知識を拡充するような、一種の継続教育のような場を提供したいということで、このセッションを企画した次第です。

ここに書いているような技術の中には、もしかしたら、実務で使えるような、面白い手法があるかもしれないと思いますので、そのような部分を今後もっと深掘りして、実務で使ったり、論文を書いたり、そのような方向につなげていただけたらいいなと思っている次第です。

## INTRODUCTION - 線形重回帰モデルとスパースモデリング -

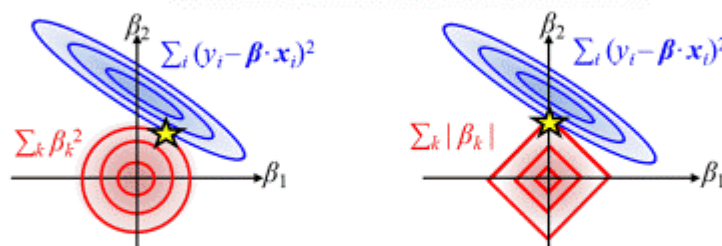
- $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$
- どうやって  $\beta_j$  を決めるのか?
- P値? それとも AIC?
- 最小二乗法  $RSS = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2$  を最小化
- RIDGE回帰  $RSS + \Lambda \sum_{j=1}^p \beta_j^2$  を最小化  $\Leftrightarrow$  RSSを最小化 SUBJECT TO  $\sum_{j=1}^p \beta_j^2 \leq s$
- LASSO回帰  $RSS + \Lambda \sum_{j=1}^p |\beta_j|$  を最小化  $\Leftrightarrow$  RSSを最小化 SUBJECT TO  $\sum_{j=1}^p |\beta_j| \leq s$ 
  - L1ノルムの罰則項を加えることで、パラメータ推定と同時に、変数選択を自動的に行う手法

導入部分ということで、一般化線形モデルや線形重回帰モデルですね。これはアクチュアリー試験にも出てきますので、皆さんご存じだと思います。どのようにパラメーターを選ぶのかというときに、例えばP値や、AICというのも一つの手法にはなるかと思っています。



一般的にこれまで教えられていた統計学に基づくと、最小二乗法を使って誤差が最小になるようにパラメーターを選ぶというのが常套手段だと思うのですが、例えばそこに罰則化項を入れて、このリッジ回帰や Lasso の部分ですけれども、この  $\lambda$  で、 $\lambda$  がついている罰則化項を入れて、その上で最小化する。これを Lagrange を使って表現すると、右側の表現になって、制約条件付きの最小化問題になります。このような手法を用いて、パラメータの推定と変数選択を自動的にやるという概念が、Lasso です。もう少し詳しい概念は小林さんの方から、この後説明があると思います。

## INTRODUCTION - 線形重回帰モデルとスパースモデリング -



【出典】第2回データサイエンス集中セミナー「統計数理研究所の歩みとこれから」  
野村俊一

- 機械学習ワークショップ@シンガポール・アクチュアリー会(2015年)
  - UCI MACHINE LEARNING WEBSITEの米国での糖尿病の入院患者データ
  - 10年間、100,000レコード、50個の説明変数(中にはカテゴリ変数も存在)
  - ミッション: 糖尿病の治療で再入院する患者の予測モデルの構築

9月のデータサイエンスの集中セミナーで、統数研の野村先生がお話しされていた中にも、この罰則付き回帰の話がございました。これはイメージなのですが、左側がリッジ回帰、右側が Lasso です。その前のスライドで、制約条件付きの最適化問題の所ですけれども、RSSを最小化するとき、制約条件が二乗の場合がリッジで、Lasso は絶対値でということですので、その制約条件が二乗の場合はこの円のような形で決まって、その中でRSSが最小化できるようにパラメーターを決めると、その  $\beta_1$  と  $\beta_2$  というパラメーター空間の中では、星で書いた所が最適化の解になるということです。

一方で、右側の Lasso の所で見ると、この  $\beta_1$  と  $\beta_2$  の所で、絶対値ですので、とがっているような形になっています。この例で言うと、その  $\beta_1$  が0になります。このような形で変数を自動的に選択して、 $\beta_1$  は0になって、 $\beta_2$  を使えば良いというように、自動的に変数選択してくれるような特徴があるので、スパースモデルと言っています。

Lasso に私が初めて出会ったのは、2015年にシンガポール・アクチュアリー会です。マシニングのワークショップに参加した時に、シンガポールのアクチュアリー会の人々が Lasso を用いていました。それが糖尿病のデータで、実際に匿名化した生データを使って、糖尿病の患者が、次に再入院する確率を予測するモデルを作るということを、ワークショップでパソコンを持ち込んでやりました。その中でも、説明変数が50を超えるような中で、どれが効くのかというところを、AICやP値でやってももちろん良いのですが、Lasso を使ってやると、数行のコードで一気にその変数が選択できるという説明を受けて、結構、私自身は衝撃的に思ったのが、その Lasso と初めて会った時の印象です。

# INTRODUCTION

## - 要介護者の死亡率 -

- 実務的な課題: 推定に用いるデータが不足
- “長寿化を考慮した要介護状態別死亡率の予測—混合リー・カーター・モデルによるアプローチ”(小暮厚之、神谷信一、伏見考弘、2016年度JARIP 研究発表大会)
  - この論文では、長寿化を考慮して要介護状態における死亡率を予測するための統計的手法について考察する。このような手法を開発する際に困難となる点は、要介護状態の人口集団に対する死亡データが十分に整備されていない点である。本稿では、全人口集団が健康状態別の部分人口集団に区分されるような状況で、各健康状態別の死亡数データは欠如しているが、対応する人口データは利用可能という状況下で、健康状態別死亡率を予測する新たな手法を提案する。この手法を我が国の介護年金制度のデータに適用し、要介護状態を反映した死亡率の予測を議論する。

あと、介護の観点で、これは介護のプライシングをされた方は皆さん、共通の苦勞をされていると思うのですが、やはりデータがとても不足していて、どのように死亡率や要介護の発生率を推定するのかというところが、かなり難しいのです。このような部分を補完するとき、小暮先生からお話いただくベイズ的なアプローチを使うと、データがない所をうまく事前分布を置きながら、要介護者の死亡率の推定が可能となります。

## 本セッションの目的

- 一種の継続教育
- データ・サイエンスに関する集中セミナー(第3回: 岩沢宏和)
  - 以下のたぐいの(もっとずっと詳細な)整理を早急に行う必要がある。
  - 機械学習の手法や道具立てと称されるもので、アクチュアリー分野での予測に直接的ないし間接的に役立つものうち、
    - I) 統計学の手法や道具立てとして理解できるもの(主成分分析, 罰則付回帰, EMアルゴリズムなど) → 統計学の枠組みの中で取り入れ
    - II) I)以外で汎用的な手法 → 統計解析の基本手法として取り入れるべきもの(クラスタリングの諸手法など)と(将来に備えて)その他の機械学習の基本手法として習得しておくべきもの(NN, SVMなど)とに分類して取り入れ



これは最後のスライドです。データサイエンスの集中セミナーの時に、岩沢先生がおっしゃっていたように、データサイエンスと言うと、何かふわっとして、定義がなく、どこまでがデータサイエンスなのか、よく分からない印象があるかと思うのですが、そのような中でも、粒々のその手法の中に、アクチュアリーが使える手法は多分幾つかあるはずで、私が個人的に思っているのは、この Lasso やベイ

ズ的アプローチは、比較的アクチュアリーに近いデータサイエンスの範ちゅうの手法なのではないのかと思っています。ですので、これは繰り返しになるのですがけれども、今日この後続く話を聞いていただいた上で、面白いなと思えば、それを深掘りして実務に使ってみたり、論文を書いてみるなど、そのような方向に向かうといいなと思っています。私の発表は以上です。

**山内** 藤澤さん、どうもありがとうございました。それでは、引き続きまして、小林さんの方からプレゼンテーションをお願いしたいと思います。

**小林** ご紹介ありがとうございます。

## アクチュアリー年次大会

小林 凌雅  
慶應大学 後期博士課程  
ryoga@sfc.keio.ac.jp

慶應大学の小林です。今日はどのような話をするかという、まず統計学と機械学習についていろいろ思うことを書いてみたので、皆さん、聞いてくれたら嬉しいです。その後に、先ほど藤澤さんからありましたスパース推定について、話していこうと思っています。

# Big Data

## Gartner の 3V<sup>[1]</sup>

- Volume
  - Variety (高次元データ)
  - Velocity (高頻度データ)
- による定義されることが有名です。



しかしながら、今回は Big Data を  
**計画されずに集められたデータ**  
として話をすすめていきます。

[1] Douglas, Laney, 2001, "3D Data Management: Controlling Data Volume, Velocity and Variety," Gartner.

3

まずビッグデータについてです。これは発表する前に藤澤さんや小暮先生、山内先生と話していて、他にもアクチュアリーの方はいたのですけれども、そのアクチュアリーの中の1人が、「私はビッグデータという言葉はなるべく使わないようにしている」と言っていました。そこで私はひねくれているので、あえてビッグデータという言葉を使ってみました。ここでビッグデータという言葉を使うために「ビッグデータとは何か」と改めて考えてみました。どのようなものとして捉えるのがいいかと思ったのですが、一般的にはこのガートナー社の3V、Volume、Variety、Velocityのような話によって定義されることが多いと思います。ただ、今までの自分の周りの現状などを見て、ビッグデータをどう定義すればいいかと考えたときに、計画されずに集められたデータのことをビッグデータとして定義するのが、とても自分の中で直感的であると思いました。なぜかと言うと、昔のコンピュータ、その左の図のような、紙に穴を開けてプログラミングするというコンピュータがあったらしいのです。自分は知らないのですが、そうですね、小暮先生。

小暮 そうなのです。

小林 ただ、今のコンピュータは右側のコンピュータのように、例えば回帰分析、逆行列の演算は大きな処理が必要です。逆行列を計算することも、あの左側のパソコンだったらどれくらいかかるのかと考えたら、到底その何も考えられずに集められたデータなどは分析できるような時代ではなかったと思われます。ただ、その逆行列の演算などもぱっと計算できるので、やはり計画されずに集められたデータを分析するような余裕が出てきたと、最近感じています。

これは、あるコンサル業界の友人から聞いた話なのですけれども、コンサル業界に集まってくるデータも、やはり計画されずに、「とりあえず何とかしてくれ、ドーン」というように渡されるようなデータが多いらしいのです。これは業界用語で「ドーン案件」と言うらしいのです。やはり計画されずに集められたデータというのが重要なのかなと、このスライドを作った時はとても思っていました。ただ、その後にあ

る企業さんに行かせていただいて、画像のデータがあるのですが、この画像のデータをパートの人が一つ一つ「どこが手で、どこが頭だ」とプロットしているような会社がありまして、この定義もちよつと古くなりつつあるのかなと、今、この場で話していると思っています。

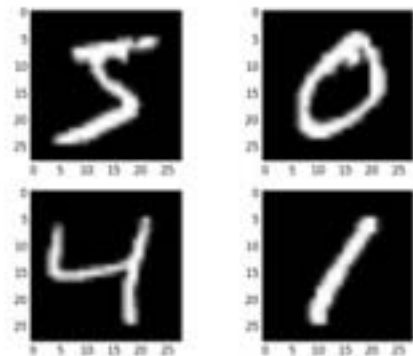
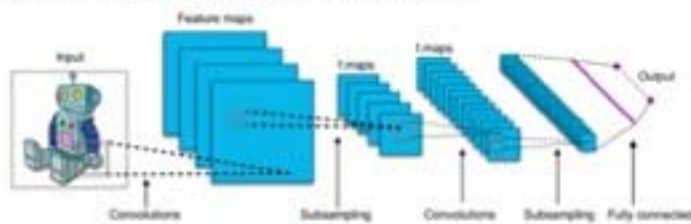
例えばですけれども、画像解析などは結構、最たる例だと思います。画像解析の分野では機械学習の方法を用いて、盛んに研究がされてきています。ここで扱われる画像データは、典型的には何か明確な目的があって、このような分析をしようなどと考えずに集められたデータではなく、ただウェブ上にある取得可能なデータを適当に引っ張ってきて、それから何ができるのかと考えていることが、とても多いように感じます。このように扱われるデータを、ここではビッグデータとして話していきます。

## Big Data

### 例 MNIST のデータ

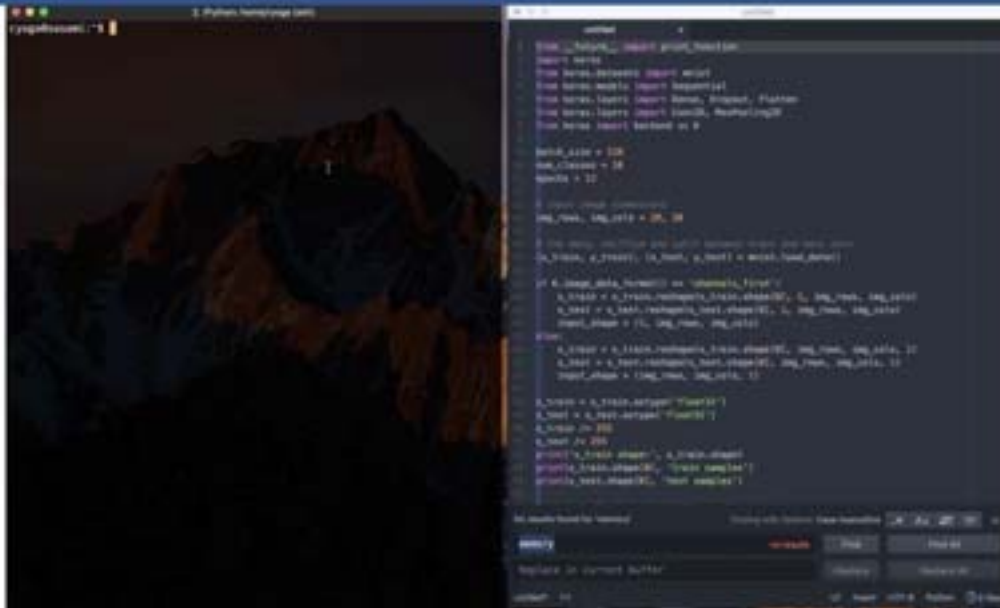
MNIST は Modified National Institute of Standards and Technology の略です。手書きの数字のデータ・セットです。

それを畳み込みニューラルネットワークによって数字を識別することを考えます。



一番有名なのは、この MNIST のデータなどなのですが、手書きの数字のデータを、はやっているのも、ニューラルネットワークのような方法を使って分析していました。

# Big Data



話ばかり聞いてしまうと、多分飽きてしまうので、実装をお見せしようと思います。小さくなってしまっているのですが、このような感じで、さっきの手書きの画像のデータを分析しています。ただの手書きの画像、この文字が幾つなのかと分類のモデルを回しています。これはPythonですね。あえてiPythonで、ちょっと派手にやっています。これがこう分析していくと、このように計算ができて……。ここですね。たしか10分強あるのですけれども、これがたった10分、この手元のパソコンで回すだけで。そろそろですね。「正答率が99%も出るって、結構すごい時代になったな」と、個人的に思います。先ほどのこのパソコンでは考えられなかったような時代になってきたのかなと思っています。

## R と Python

最近、次のように質問をされることがあります。

「R と Python はどちらを勉強したら良いでしょうか？」



これも最近よく聞かれるので、あえて入れたのですが、R と Python、どちらを勉強した方がよいで

すか」という質問を結構されます。結構、RとPythonの違いについて、真剣に考えてみました。いろいろ考えてみたのですが、このRとPythonの違いで最も大きいことは、やはりコミュニティの違いかなと思います。

## R と Python

R と Python の違いで最も大きな点は、**コミュニティの違い**です。

実現できることに関しては、あまり大きな差はないように感じます。

しかしながら、  
**R** は**統計学**を研究しているコミュニティの人が利用するのに対して、  
**Python** は**機械学習**を研究している人が利用する傾向があります。

利用している人の多さによって、過去のナレッジの蓄積が異なります。

言語でできることは、問題がきちんと定義できて、関数さえ書ければ、全部同じような処理ができるので、多分あまり違いはないと思います。ただ、扱っている人のコミュニティは大きく違うなど感じています。Rは統計学を研究しているコミュニティの人がとても利用していると思います。統計学会などへ行くと、ほとんどRでコードを書き、結構自分はびっくりしました。逆に言うと、このPython、では、機械学習のコミュニティの研究発表に行くと、ほとんどの人がPythonで書いています。利用している人の多さによって、過去のナレッジの蓄積が異なるので、それは「自分がやりたいことが、どっちが多いのかな」ということをきちんと考えないと、判断できないものかなと思いました。あまり答えになっていないですがこれがどちらを勉強したらよいかの答えです。

## 統計学 と 機械学習

次に統計学と機械学習の違いについて考えていきます。

どちらもデータから情報を得る方法という点については同じですが、**統計学は解釈を重要視し、機械学習は予測を重要視するもの**とされることがよくあります<sup>\*1</sup>。

解釈/予測をするためには何をすればよいのでしょうか？

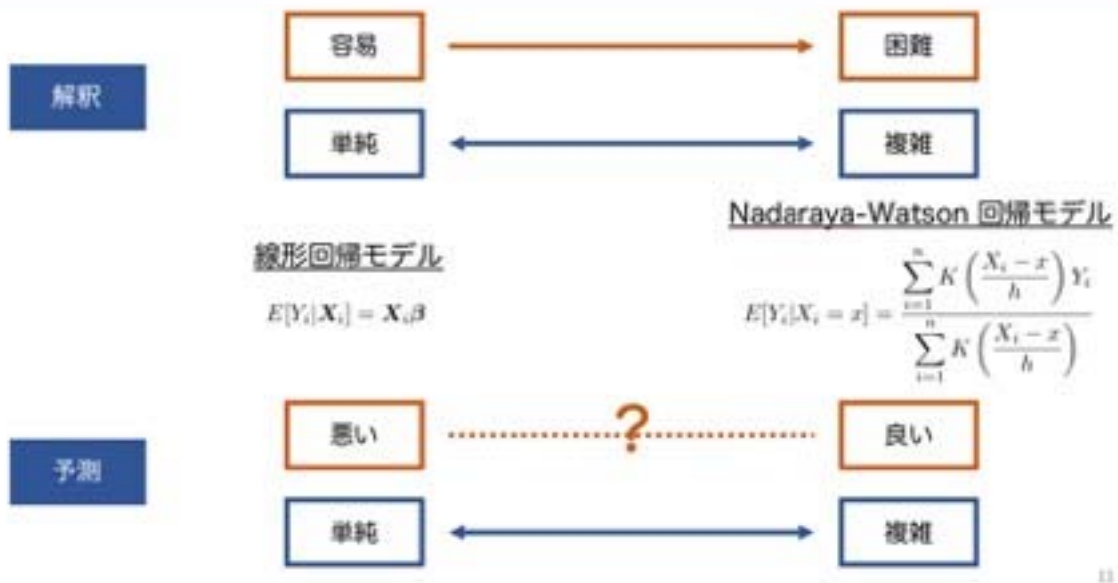
\*1 一般的に統計学として扱われるものにも解釈が難しいものや予測に重きをおいたもの、同様に機械学習として扱われるものにも解釈が簡単なものなどがあります。

次に統計学と機械学習の違いについて、考えていきたいと思います。統計学と機械学習は何が違うのかということも改めて考えてみたのです。これは人の回答を借りたのですが、どちらもデータから情報を得る方法という点については同じです。ただし、統計学は解釈を重要視していて、機械学習は予測を重要視しているものとされることがよくあります。ただ一般的に統計学として扱われるようなものにも、解釈が難しいものや、予測に重きを置いたものもあります。同様に機械学習でも解釈を簡単にするためのものもありますが、こちらは、とりあえずこれで話を進めていきたいと思います。

その解釈、予測をするためには何が必要なのでしょう。恐らくですけども、どちらもモデル、モデル化によって実現することが典型的な方法だろうと考えました。ただし、異なる点として、解釈することは、より単純なモデルを使った方が実現できるのに対して、予測することはモデルの複雑さはあまり関係ないです。



# 統計学 と 機械学習

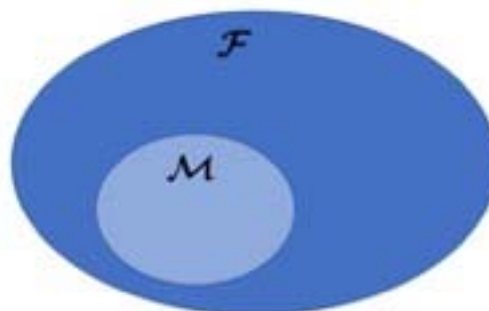


そもそも複雑さというものは何なのかと考えていくと、結構複雑さというのは人間が定義するものなので、解釈できるかどうかによって定義されると考える方が自然なのかもしれません。もちろん何か数学的にアプローチする方法もあるのですが、この解釈に関しては、直感的に容易なものや困難なものがあります。ただ予測に関しては、「単純なモデルだから悪い」「良い」というのは、正直よく分かりません。単純なモデルの方が、予測がよくなる場合もありますし、悪くない場合もあります。

## モデル

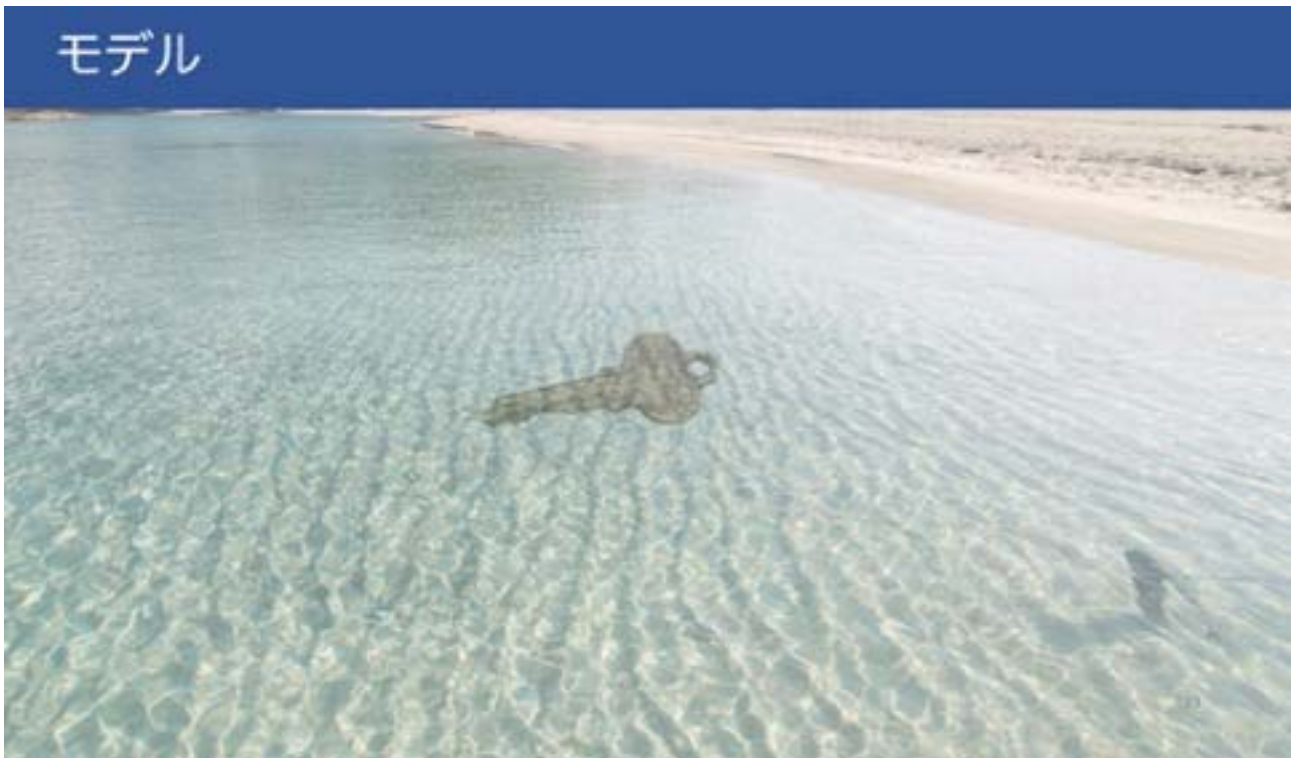
モデル  $\mathcal{M}$  とは、すべての考えられる集合を  $\mathcal{F}$  とすると、その  $\mathcal{F}$  の部分集合であると定義します。

⇒ モデル  $\mathcal{M}$  とは、**仮定を付加することによって**  $\mathcal{F}$  を絞っていくこと。



そもそも、そのモデル化によって実現できるという話をしましたが、そもそもモデルとか何かという話もしておきたいなと思いました。モデルとは何かと言うと、すべての考えられることの中に、幾つか仮定

を置くことだと。仮定を置いて、集合を小さくしていくことをモデルと言うこととします。



ちょっと分かりやすい例かどうか分からないですけども、友人が車に乗っていて海まで行ったのですけれど、鍵をなくしてしまいました。海で鍵をなくしてしまって、鍵を探そうとなったのですけれど、皆さんは海でかぎをなくしたら、どのような所を探しますか。いわゆる機械学習の人は、海全部を掘れば鍵は見つかるだろうと考える人達が先ほどの定義の機械学習だと思います。ただ、統計学の場合は、この辺りにかぎがあるだろうということを、モデルを置いて探していきます。これはとても一長一短があって、自分がいた場所から誰かが持って行ってしまったと。そうなってしまうと、そのかぎは一生見つけることはできないです、その自分のいた所だけを探しては。ただし、この全部を掘り返してしまえば、いずれかは見つかると思います。

前半のまとめです。これまで統計学と機械学習の違いを比較してきました。しかしながら、最近ではやはり統計学と機械学習は別に反発し合うのではなくて、昔の確率論と統計学のように対立し合うのではなくて、統計学と機械学習は互いにいいところを認めて、取り込んでいく方向に動いているのではないのかと個人的には思っています。また、統計学と機械学習の間のような、統計的機械学習と言われる分野も、非常に近年注目を集めたりしています。

## スパース推定

次の分析例を通して考えていきます。

下記のようなデータ・セットが得られたとします<sup>[2]</sup>。

**被説明変数:** 中性脂肪 (diabetes patients)

**説明変数:** 年齢 (age), 性別 (sex), BMI (body mass index), 血圧 (average blood pressure) など残り 6 変数。

データ数は 442 個で、1 年前をベースラインとします。

[2] Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R., 2004, "LEAST ANGLE REGRESSION," *The Annals of Statistics*, Vol. 32, No. 2, 407-499.

18

ここからは、スパース推定を話していきたいと思います。今回は特に次の小暮先生の発表に合わせて、バイズを意識して、このスパース推定について説明しています。

まず例を通して考えていきましょう。次のようなデータセットが得られたとします。例えば、説明するものとしては、中性脂肪を説明したいとします。ただ被説明変数としては年齢、性別、BMI、血圧などを説明変数となるデータを考えましょう。これは結構有名な論文なのですが、Efron などが書いている、Lars という最適化方法の具体例で扱われているデータをそのまま持ってきました。

## スパース推定

Patient	AGE	SEX	BMI	BP	Serum measurements					Response	
	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$y$
1	59	2	32.1	101	157	93.2	38	4	4.9	87	151
2	48	1	21.6	87	183	103.2	70	3	3.9	69	75
3	72	2	30.5	93	156	93.6	41	4	4.7	85	141
4	24	1	25.3	84	198	131.4	40	5	4.9	89	206
5	50	1	23.0	101	192	125.4	52	4	4.3	80	135
6	23	1	22.6	89	139	64.8	61	2	4.2	68	97
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
441	36	1	30.0	95	201	125.2	42	5	5.1	85	220
442	36	1	19.6	71	250	133.2	97	3	4.6	92	57

Efron et al. (2004) から引用。

17

このようなデータセットですね。

## スパース推定

ここでは、**Data Scientist** とは Big Data を分析する人とします。

Data Scientist が先のデータ・セットを分析することになったら  
どうすればよいでしょうか？

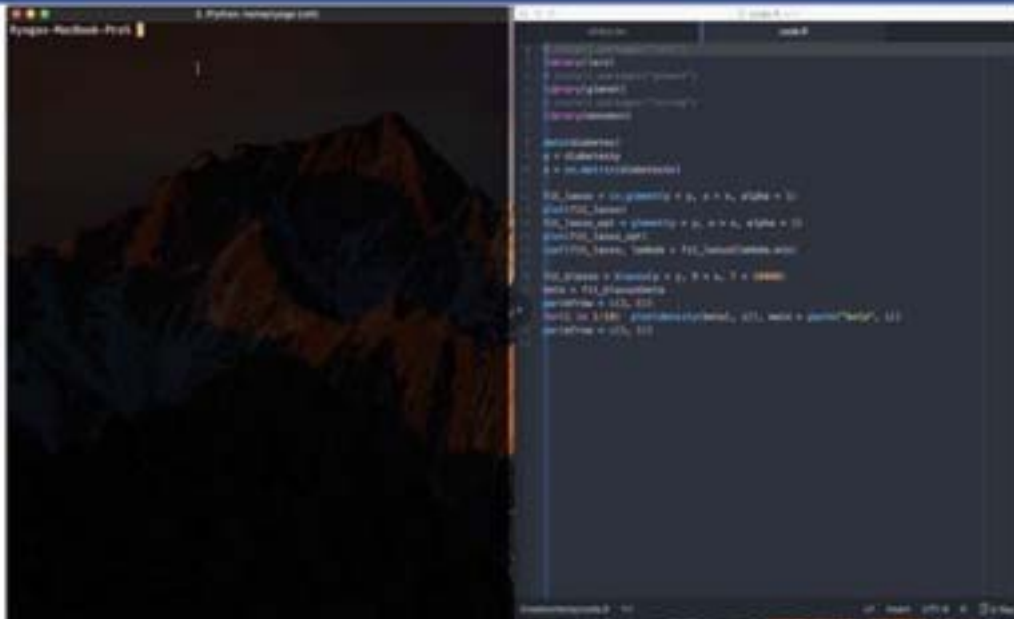
Data Scientist は、まだ新しい職種なので専門の人というよりは、  
従来の仕事をしつつ、Data Scientist の仕事もしている人が多い  
ように感じます。

⇒ 時間がない人が多い

ここで、データ・サイエンティストという言葉が最近よく聞きます。そこで、データ・サイエンティストという言葉も自分なりに定義してみようと思いました。このデータ・サイエンティストとは、ビッグデータを分析する人です。つまり、この計画されずに集められたデータを分析する人と定義します。

データ・サイエンティストが、さっきのデータセットを分析することになったら、どうすればよいでしょうか。データ・サイエンティストは、これも自分の周りの話を聞いた状況を考えて、まだとても新しい職種なので、専門の人がデータ・サイエンティストになるというよりは、従来の仕事をしつつ、データ・サイエンティストの仕事もしている人も多いように感じました。つまり、何が言いたいかというと、多分時間がない人が多いと、皆さん、言っていました。そのようなときに、スパース推定が使えるのではないのかなと思っています。

## スパース推定量



これですね。このように、今度はRで書いているのですが、もし興味があれば、参考としてコードもつけましたので、ぜひ見てみてください。

**山内** 小林さん、これは具体的に何をしようとしているものなのですか。

**小林** それは、これから説明していきます。すみません。多分、先に見せた方が後に聞くモチベーションになるかなと思ったので。

**山内** 失礼しました。

**小林** ここで推定していますね、何かしら。MCMCをして、かっこいいなと思って、このような感じで説明できたらなと思っています。

## LASSO 推定

そのような分析をするために  
LASSO 推定量が提案されています<sup>[3]</sup>.

LASSO 推定量を用いることにより  
簡単に重要そうな変数を抽出する  
ことができます.

⇒ 重要な変数を選択する時間を  
節約することができます.

[3] [Tibshirani, R.](#), 1996, "Regression shrinkage and selection via the lasso," *Journal of Royal Statistical Society, Series B.*, Vol. 58, No. 1, 267-288.

(Intercept)	152.13348
age	.
sex	-13.59236
bmi	506.66309
map	199.02672
tc	.
ldl	.
hdl	-124.16053
tch	.
ltg	441.69986
glu	.

20

今まで何をしていたかということを説明しますと、Lasso 推定量というものをやっていました。Lasso 推定量とはどのようなものかというところ、Lasso 推定量を用いることによって、簡単に重要そうな変数を抽出することができます。重要な変数を選択する時間を節約することができます。

そこに出ているように、先ほどのデータセットの線形モデルを考えたのですが、こちらのデータセットですね。そうすると、例えばこのデータセットにおいては、年齢などは関係ありませんよ。その代わりに、性別やBMIなどがとても重要な要素となっていますよ。直感的にも、中性脂肪ならBMIなどは何となく合っているのかなという感じがします。

## LASSO 推定

従来の統計学の枠組みでは、それぞれの推定量をみることによって、  
変数選択を行っていました。

それをあたかも自動的に選択してくれているように見えます。

この LASSO 推定量は、機械学習の知識を用いた統計学というよりは、  
Ridge 推定量<sup>[4]</sup>の正則化の拡張やAIC<sup>[5]</sup>などの変数選択の拡張と捉える方が  
良いという人もいるかもしれません。

[4] [Hoerl, A. E. and Kennard, R. W.](#), 1970, "Ridge Regression: Biased Estimation for Nonorthogonal Problems," *Technometrics*, Vol. 12, 55-67.

[5] [Akaike, H.](#), 1973, "Information theory and an extension of the maximum likelihood principle," *Proceedings of the 2nd International Symposium on Information Theory*, 267-281.

21

従来の統計学の枠組みでは、それぞれの推定量を見ることによって、先ほど藤澤さんもおっしゃっていましたが、P値やt値、もしくはAICを見ることによって、変数選択を行ってきました。それをあたかも自動的に選択してくれているように見えます。

このLasso推定量は、機械学習の延長として書いているのですが、本当はリッジ推定量の正則化の拡張や、AICなどの変数選択の拡張として捉える方がよいという人もいるかもしれません。

## LASSO 推定

最小二乗推定量は、

$$\hat{\beta}_{ols} = \arg \min_{\beta \in \mathbb{R}^d} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - X_i' \beta)^2 \right]$$

で定義されるのに対して、LASSO推定量は、

$$\hat{\beta}_{lasso} = \arg \min_{\beta \in \mathbb{R}^d} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - X_i' \beta)^2 + \lambda \sum_{j=1}^d |\beta_j| \right]$$

で定義されます。ただし、 $\lambda > 0$  とします。

もうちょっと具体的に話をしますと、いわゆる最小二乗推定量とはどのようなものだったかという、通常の最小二乗法と同じように、このように定義ができます。ただ、一番後ろにその推定量の絶対値の和が入倍されたものを加えたもので定義しています。

## LASSO 推定

ラグランジュ的に捉えると、

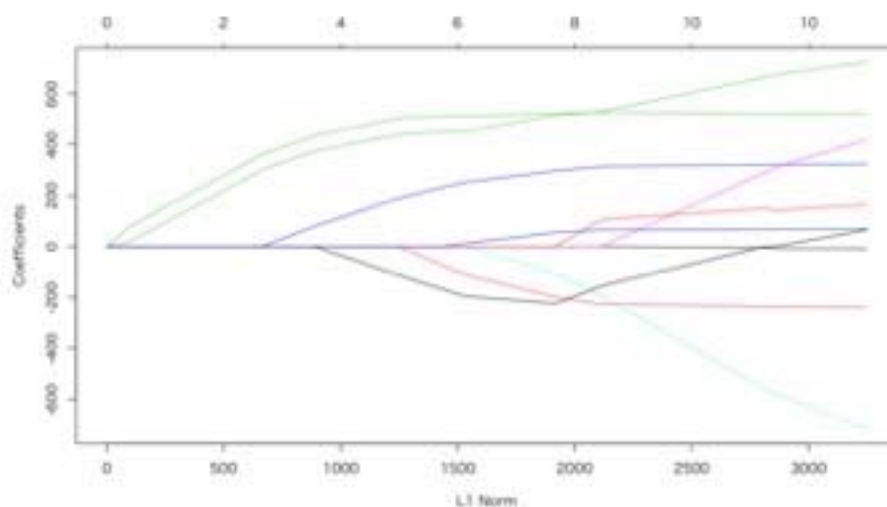
$$\begin{aligned} \min_{\beta \in \mathbb{R}^d} & \frac{1}{n} \sum_{i=1}^n (Y_i - X_i' \beta)^2 \\ \text{s.t.} & \sum_{j=1}^d |\beta_j| < \lambda \end{aligned}$$

と書き換えられます。これは、推定量の絶対値和が“予算” $\lambda$ の中で選択されると考えることができます。

これは次のような問題として見ることもできます。そもそも絶対値の和が $\lambda$ より小さくなるということはどういうことなのかというのを直感的に説明すると、普通の最小二乗推定量が $\beta$ 、推定量の回帰係数の和などは、マイナス無限から無限まで取れるような状況を想定しています。これは絶対値が $\lambda$ に収まるように、つまり「予算が $\lambda$ しかない中で、この推定量を振り分けてください」と考えることができます。この提案していた人、Hastie さんの本などにも書かれています。

## LASSO 推定

$\lambda$ の大きさによって変数の数が変わることが分かります。



そうすると、この $\lambda$ を大きくすることによって、右側に行くほど大きいのですが、左に行くほど $\lambda$ が小さくなっていて、上に書いてある小さな数字なのですが、それは選ばれる変数の数となりますね。すみません、左が大きいですね。 $\lambda$ が小さくなると、選ばれる変数の数が大きくなっていて、 $\lambda$ を小さくしていくと、すべて0になってしまう。つまり予算が0なので、全部0になってしまうということが分かります。



## LASSO 推定

LASSO 推定量は関係ない変数が含まれているかもしれないというモデルを考えています。

統計学では、推定量は**一貫性**と**漸近正規性**の2つの性質をもつ推定量を良い推定量としてきました。

しかしながら、LASSO 推定量は、**バイアス**を持ちます。また、変数選択ができることと引き換えに**一貫性を失っています**。

SCAD<sup>[6]</sup> 推定量のような（ある条件下では）一貫性をもつLASSO のような推定量も提案されています。

[6] Fan, J., and Li, R., 2001, "Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties," *Journal of the American Statistical Association*, Vol. 96, No. 456, 1348-1360. 36

Lasso 推定量は関係ない変数が含まれているかもしれないというモデルを考えています。統計学では、推定量は一貫性と漸近性の二つの性質を持つ推定量を、良い推定量としてきました。しかしながら、Lasso 推定量はバイアスを持ちます。また変数選択ができることと引き換えに、一貫性を失っています。この SCAD 推定量のように、一貫性を持つようなら、Lasso のような推定量も提案されています。他にも Lasso 推定量以外に、予算の決め方というのは他にもあるのではないかとということで、その他のスパース推定量について議論しています。

## その他のスパース推定

LASSO 推定量以外にもさまざまなスパース推定量が提案されています。

スパース推定量を次の関数最初化する推定量とする。

$$Q(\beta) = \frac{1}{n} \sum_{i=1}^n (Y_i - X_i' \beta)^2 + p(\beta)$$

ここで、説明のために、罰則項を  $p(\beta)$  と書くこととする。

とりあえず分かりやすさのために、一般に、予算のところを  $p(\beta)$ 、罰則項として書いています。

## その他のスパース推定

### Elastic Net

$$p(\boldsymbol{\beta}) = \lambda \left[ \alpha \sum_{j=1}^d |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^d \beta_j^2 \right]$$

ただし、 $\lambda > 0, 0 < \alpha < 1$  とする。

$d > n$  の場合、すべての変数を取り込むことを可能にする。

[7] Zou, H., Hastie T., 2005, "Regularization and variable selection via the Elastic Net," *Journal of the Royal Statistical Society, Series B*, Vol. 67, 301-320.

28

例えば Elastic Net という方法があります。Elastic Net とはどのようなときに非常に使えるかと言うと、例えばデータ数がデータの次元数より小さい場合です。そのような場合にも Elastic Net は推定することが可能です。このような場合は、従来の統計学ではどうしても扱えなかった方法なのですけれども、さらにこの二乗のペナルティを入れることによって、推定することが可能になります。

## その他のスパース推定

### SCAD (Smoothly Clipped Absolute Deviation)

$$p(\beta_j) = \begin{cases} \lambda |\beta_j| & \text{if } |\beta_j| \geq \lambda \\ \frac{\gamma \lambda |\beta_j| - (\beta_j^2 + \lambda^2)/2}{\gamma - 1} & \text{if } \lambda < |\beta_j| \leq \gamma \lambda \\ \frac{\lambda^2(\gamma^2 - 1)}{2(\gamma - 1)} & \text{if } \gamma \lambda < |\beta_j| \end{cases}$$

LASSO 推定量は大きな回帰係数について過剰な罰則を付加点を緩和した。

[6] Fan, J., and Li, R., 2001, "Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties," *Journal of the American Statistical Association*, Vol. 96, No. 456, 1348-1360. 29

他にも SCAD という方法が提案されています。これはどのような方法かと言うと、このいわゆる Lasso 推定量というのは、絶対値をずらす推定量なのですけれども、SCAD の方は、横軸が真のパラメータなのですけれども、真のパラメータが大きくなればなるほど、きちんと同じようにズレを修正していくようなペナ

ルティとして考えられています。このように、一番上が Lasso 推定量に対応する所なのですけれども、さらに戻っていくようにすることによって、ある特定の状況においては一致性のような性質を持つような方法が提案されています。

## その他のスパース推定

### MC+ (Minimax Concave)

$$p(\beta) = \lambda \sum_{j=1}^d \int_0^{|\beta_j|} \left\{ 1 - \frac{x}{\nu\lambda} \right\} dx$$

LASSO 推定量は大きな回帰係数について過剰な罰則を付加点を緩和した。

[8] Zhang, C. H., 2010, "Nearly unbiased variable selection under minimax concave penalty" *Annals of Statistics*, Vol. 38, 894-942.

他にもこのような方法が提案されていて、いろいろな方法が提案されています。いろいろ紹介されたは良いものの、どの方法がいいのかという話になってくると思うのですが、先ほどのデータセットを使って、Lasso、Elastic Net、SCAD、あと、MC+というような方法で比較してみました。

## その他のスパース推定

これらの推定量の比較は次の通りです。

	ols	lasso	elastic net	sacd	mc+
(Intercept)	152.13348	152.13348	152.133484	152.13348	152.1335
age	-10.01220	.	.	.	.
sex	-239.81909	-13.59236	-29.564942	-228.52846	-241.9679
bmi	519.83979	506.66309	479.262681	533.89500	518.9142
map	324.39043	199.02672	210.963526	323.96597	321.5075
tc	-792.18416	.	.	-139.93804	-618.5891
ldl	476.74584	.	.	.	351.7738
hdl	101.04457	-124.16053	-144.833396	-238.79698	.
tch	177.06418	.	.	.	128.0975
ltg	751.27932	441.69986	420.982368	550.59918	690.9516
glu	67.62539	.	9.111982	23.82546	67.1174

そうすると、どれも、例えば年齢のような変数は除いた方がいいと。ただ、この辺りやこの辺りは、あ

る方法によっては入れた方がいい、ある方法によっては入れない方がいいとなっています。なぜ、このようなことが起きるのかなということ考えたのですけれども、これはデータ数が無限大、かなり大きなデータ数があるような状況では、すべて統一された結果が得られるはずですが、しかしながら、山をどの方向から登っていくのかということで、この違いが出てきているのかなと思います。

ちなみに、これはA I Cでもやってみたのですけれども、実はA I Cだとうまく働かなかったのですね。それぞれの変数についてその変数を入れるパターン、入れないパターンの2通りあります。例えば変数が10個あるとすると、2の10乗、1,024通りあります。従来の方ですと1,024通りのA I Cをやらなければいけないとなります。A I Cは基本的には、多い場合はステップ関数のようなものを使ってやりますが、それはアルゴリズム的にうまくいくことが保証されていません。実際に今回の場合はうまくいかなかったです。スパース推定は、このように最適化の部分がしっかり議論されているので、うまくいっているのかなと思いました。

## Bayesian LASSO 推定量

この LASSO 推定量はベイズ推定量 (MAP推定量) として解釈できます。

事前分布をラプラス分布:

$$f(\beta_j) = \frac{1}{2\tau} \exp \left\{ -\frac{|\beta_j|}{\tau} \right\}$$

とすることと同様の結果が得られます。

ここで、 $\tau = 1/\lambda$  としています。

今度はベイズ的な解釈をしてみようと思います。この Lasso 推定量は、実はベイズのMAP推定量として解釈することができます。MAP推定量に関してはちょっと後で言及します。どのようなものか。事前分布をラプラス分布とすると、同様な結果を得られます。先ほどのデータに対して、ラプラス分布を事前分布として、MCMCによって事後分布を求めています。その事後分布における最大値と Lasso 推定量の関係性が見てとれると。

## Bayesian LASSO 推定量

先程のデータに対して、ラプラス分布を事前分布として、MCMCによって事後分布を求めてみます<sup>[9]</sup>。

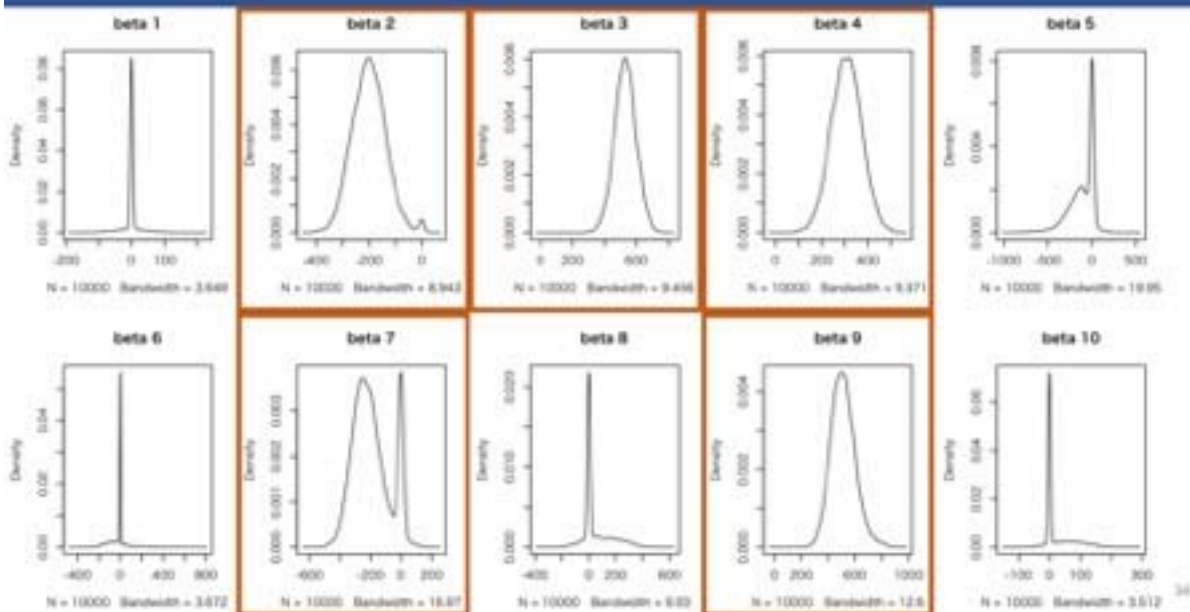
その事後分布における最大値と、LASSO 推定量の関係性がみてとれます。

[9] Park, T., Casella, G., 2008, "The Bayesian Lasso," *Journal of the American Statistical Association*, Vol. 103, No. 482, 681-686.

11

ちなみにリッジ推定量の場合は、平均0の正規分布を事前分布としたときの、事後分布のMAP推定量に対応しています。

## Bayesian LASSO 推定量



そうすると、こうなっていて、例えばこの7番目の変数などはとてもとがっていて、8番目の変数はとてもとがっていて、0の部分でとてもスパイクな分布が得られています。

MAP推定量というのは、この推定量の中の一番、尤度ではないですけど、尤度のような事後分布が大きくなる所を選択しているんで、0の部分でピークが来ている分布が幾つかあって、実はこれが Lasso 推定量と対応しているというのは、その Lasso 推定量が提案されたときから言及されています。

## まとめ（全体）

- ▶ はじめに統計学と機械学習の違いについて導入しました。
- ▶ 最近では互いの良いところを取り入れつつどちらも発展しています。
- ▶ そのひとつの例として、LASSO 推定を紹介しました。
- ▶ LASSO 推定以外のスパース推定について説明しました。
- ▶ このようなベイズ推定量とも関係が見えました。
- ▶ 機械学習と同様にベイズ推定量も Big Data とともに相性が良いのかもしれませんが。

まとめとして、初めに統計学と機械学習の違いについて、導入しました。最近では互いのいい所を取り入れつつ、どちらも発展しています。その一つの例として、Lasso 推定量を紹介しました。Lasso 推定以外のスパース推定についても、議論しました。先ほどのようなベイズ推定量とも、実はとても関係があるのだということも見えてきました。機械学習と同様に、ビッグデータなど大きなデータを扱うような場合は、ベイズ推定というものも、実は Lasso 推定と同様に相性がいいのかもしれませんが。最後に、次の小暮先生の発表の前置きをして、今回は終わらせていただきます。ありがとうございました。

**山内** 小林さん、どうもありがとうございました。でも、多分聞きたいことは山ほどあると思うのですが、ちょっとここは心に留めていただいて、最後にQ&Aの時間を設けますので、よろしく願います。それでは、小林さん、ありがとうございました。

それでは小暮先生、よろしくお願いします。

**小暮** 慶應大学の小暮と申します。どうぞよろしくお願いします。

セッションB 15:15～16:45  
「医療介護とデータサイエンスの新たな視点」

## ベイズ法の視点

小暮 厚之

慶應義塾大学 総合政策学部

2017年11月10日

1

先ほど小林君からうまく言っていて、私はベイズという観点から、最近少し研究をしているので、さっきの機械学習もそうなのですけれども、このような新しい統計の流れというのがあって、その中で機械学習自体が私は、とても大きく言ってしまうと、ベイズの一つの形なのではないかなと思っています。それくらいベイズというのは、ある意味基本的な、データ分析の骨組みを与えるものではないかなと思っています。というわけで、今日はベイズについてお話をしたいと思います。

たくさん話したいことがあるのですが、少しお話を絞って話したいと思います。最初はベイズ法の典型的な応用ということで、信頼性理論というものがありますけれども、これは保険の世界でとても有名な理論です。ベイズが使われている典型的な例として、とても分かりやすいですし、皆様はもちろん専門家ですので、ベイズについてはよくご存じかと思いますが、少しメモリーをリフレッシュしてもらおうという意味で、最初に信頼性理論の中でベイズの話をしていただきたいと思います。

その後で、ベイズ法の最近の応用ということで、これは私が最近やっていることなのですが、要介護度別死亡率の推計というものを、ベイズを使ってやってみたいと思っています。その後、もしも時間があったら、モデルリスクというような問題についても、お話しできればと思っています。以上のような感じで、進めていきたいと思っています。

## はじめに: 本日の内容

- **ベイズ法の典型的な応用：自動車保険（信頼性理論）**  
事前分布と事後分布，予測分布，信頼性理論，ベイズ法とは何か。
- **ベイズ法の最近の応用：要介護度別死亡率推計**  
介護保険データ，死亡データの欠如，混合 Lee-Carter モデル
- **死亡率予測におけるモデル・リスク**  
様々な確率的死亡率モデル，状態空間モデル，周辺尤度

2

最初なのですけれども、ベイズ法とは何かということで、自動車保険を考えていただきたいと思います。これはベンジャミン&クルーグマンという、あの有名なクルーグマンですが、彼が SOA で発表していた資料があったのですけれども、それを借りて、ここで話したいと思います。

タロウさんとハナコさんという2人の人がいるとします。これは2人ともA社の自動車保険に入っています。タロウさんの方は過去5年間に2回の自動車事故を起こしているのですね。そうすると、5年で2回ですから、0.4ですね。年間事故率のベスト・エスティメートは0.4です。一方、ハナコさんというのはとても優良なドライバーでして、全く事故を起こしていない。このような状況があるとします。タロウさんとハナコさんの保険料をどうしたらいいのかという問題です。ただし、A社は自分の会社で10万人の契約者のデータを持っています。その年間の事故率は0.05です。このような問題です。

このときに、100%の信頼性というのは、タロウさんとハナコさんの個人的な経験だけに信頼を置くと、どうなるかと言うと、タロウさんというのは0.4で、全体の平均は0.05でしたよね。だから、それに従ってしまうと、タロウさんの翌年の保険料は8倍になってしまうのですね。そうすると、タロウさんは、「8倍にもなっちゃうんなら」と言って、もうこのような保険に入らないですよ。

一方、ハナコさんの場合、ハナコさんは過去5年間に1回も事故を起こしていないですから、ハナコさんの経験だけに頼ってしまうと、翌年の保険料は0になってしまうではないですか。そうしたら、保険会社は困ってしまいますよね。だから、100%の信頼性というのは、ちょっと無理だろうと。

では、逆に0%の信頼性も、その個人の経験は使わないで、A社が持っている過去の10万人のデータを使ったらどうなるかと言うと、2人とも平均的な保険料を支払います。いいように見えますけれども、何とかフェアな感じはしますけれども、でもハナコさんは優良なのだから、そのような人に対して、「いやいや、うちはもっと安い料金で保険を提供できますよ」と言ったら、ハナコさんは抜けてしまいますよね。残るのはタロウさんです。タロウさんのような人たちだけが残ってしまう。だから、0%の信頼性というものありえないわけです。



## 例：太郎と花子の自動社保険

- 太郎と花子は二人とも A 社の自動車保険に入っている。
  - － 太郎は、過去 5 年間に 2 回の事故を経験している。この個人的な経験から、年間事故率の「ベスト・エスティメイト」は  $2/5=0.4$
  - － 花子は、過去 5 年間に全く事故を経験していない。この個人的な経験から、年間事故率の「ベスト・エスティメイト」はゼロ
- A 社の保有する 10 万人の契約者の年間事故率は 0.05

(Benjamin and Klugman, 2016)

4

どうしたらいいかと言うと、ちょうど中間を取ったらいいですよね。例えば、25%の信頼性を個人の経験に置きます。そうすると、タロウさんの見積値は 0.1375、ハナコさんは 0.0375、どちらもまあまあ受け入れるような額ですよね。保険会社もこれだったら、ペイしますよね。

## 太郎と花子の保険料

各ドライバーの個人的履歴に 25%の信頼を置いたら？

- 太郎の見積値は

$$0.25 \times 0.4 + 0.75 \times 0.05 = 0.1375.$$

となる。おそらく、太郎はこの割り増しを受け入れ、より注意深い運転に気を付けるようになるかもしれない。

- 花子の見積値は

$$0.25 \times 0 + 0.75 \times 0.05 = 0.0375.$$

となる。花子は、この保険料の割引に満足し A 社にとどまるであろう。

- A 社は、平均的には適正な保険料を集めることとなる。

6

これを信頼性理論というようにいいます。これをちょっと数式で書くと、このような保険料率の決め方というのは、 $X_{\text{バー } j}$  というのが契約者  $j$  の過去の事故率です。 $X_{\text{バー } j}$  と書いてあるのは、過去 5 年分の平均。 $j$  というのは  $j$  番目の契約者という意味です。

それから、 $\mu$  というのが保険契約者全体の事故率です。これを加重平均します。Z、さっき 0.25 を使いましたから、0.25 の  $X_{\text{バー } j}$  + 0.75 の  $\mu$  というようになります。この Z のことを信頼度といいます。

これが信頼性理論による料率の算出です。何かよさそうですね。でも、よさそうだからといって、使ってしまうというわけにもいかないですね。大体この算出法はよさそうですね、どのように正当化できるのか、それから、信頼性のZはどのように決めたらいいのかということを考えなければいけないではないですか。どう考えたらいいでしょうか。

## 信頼性理論

- このような保険料率の決め方は、

$$z\bar{X}_j + (1 - z)\mu$$

と表される。ここで、

- $\bar{X}_j$  は契約者  $j$  の過去の事故率
- $\mu$  は保険契約者全体の事故率
- $z$  は個別契約者の経験の「信頼度」

- これは、「信頼性理論」による料率算出

このような算出法はいかに正当化できるか。また、信頼度  $z$  をいかに決めるべきか。

いろいろな考え方があると思いますが、信頼性理論の目的というのは、経験に基づいて翌年度の保険料率を算出します。いいですね。これを統計学の言葉で言ったらどうなるかというと、過去のデータ、 $X_1$ 、 $X_2$ 、 $X_n$ 、さっきは  $n$  が 5 でしたけれども、これに基づいて翌年度の事故率、 $X_{n+1}$ 、これを予測するという、予測の問題なのです。

だとしたら、最適な予測を使えばいいですね。予測値を  $m$  ティルド  $X$  とします。  $X$  は過去のデータです。過去のデータの関数として、予測値を作ります。一番良いのは何かというと、ターゲットの  $X_{n+1}$  との、例えば二乗リスクを最小化するようなものです。この二乗リスクを最小化するものは、条件付き期待値です。だから、最適な予測は条件付き期待値です。これを計算すればいいわけです。どのように計算しますか。

## 信頼性理論と予測

- 信頼性理論の目的は「経験に基づいた翌年度の保険料率の算出」
- 統計学の用語でいえば、過去のデータ  $\mathbf{X} = (X_1, X_2, \dots, X_n)$  に基づく次年度の事故率  $X_{n+1}$  の予測。
- 予測値を  $\tilde{m}(\mathbf{X})$  とするとき、2乗リスク

$$E[(X_{n+1} - \tilde{m}(\mathbf{X}))^2]$$

を最小化する最適な予測は

$$\tilde{m}(\mathbf{X}) = E[X_{n+1} | \mathbf{X}]$$

8

このときに、よく考えてみたら、保険契約者の中にはタロウさんのタイプとハナコさんのタイプと、両方いるわけではないですか。異質性があるわけです。そのような人たちを相手にしているのです。だったら、その異質性を表したいですね。それを確率変数  $\theta$  で表すことにしましょう。

異質性  $\theta$  が与えられたときに、将来のデータと、それから過去のデータが独立だとすると、これをずっと計算していくと、最適な予測、 $m$  ティルド  $\mathbf{X}$  というのは、一番下の表現、 $\mu(\theta)$  というのが一番下にありますがけれども、 $\theta$  が与えられたときの次の期の、来年の事故率の予測値を、 $\theta$  に関して平均を取ったものです。それが  $m$  ティルド  $\mathbf{X}$  になります。

この  $m$  ティルド  $\mathbf{X}$  が一番下になりますけれども、ちょっと一番下は分かりにくいので、ちょっと書き直してみます。異質性  $\theta$  をしようとしたときに、過去のデータ  $\mathbf{X}$  というのが条件付き密度  $f(\mathbf{x} | \theta)$ 、いわゆるデータのモデルに従うとします。それから  $\theta$  の密度を  $p(\theta)$  とします。そうすると、さっきの  $M$  ティルド  $\mathbf{X}$  という、一番いい最適な予測値というのは、真ん中の辺りに積分記号で書いてあるものになります。ここで  $p(\theta | \mathbf{x})$  というのは、 $\mathbf{X}$  が与えたときの  $\theta$  の条件付き分布です。これをベイズでは事後分布といいます。事後分布と別に言ってもいいですし、例えば条件付き分布なのですね。ですので、そこにあるような表現です。

## 信頼性理論と異質性

- 個別契約者のリスク特性の「異質性」（太郎のタイプか、花子のタイプか）が（観測できない）確率変数  $\theta$  で表されるとする。
- 異質性  $\theta$  が与えられると、 $X_{n+1}$  と  $\mathbf{X}$  が独立であるならば、最適な予測は

$$\begin{aligned}\tilde{m}(\mathbf{X}) &= E[X_{n+1}|\mathbf{X}] \\ &= E[E[X_{n+1}|\theta, \mathbf{X}]] \\ &= E[E[X_{n+1}|\theta]] \\ &= E[\mu(\theta)]\end{aligned}$$

ここで、

$$\begin{aligned}\mu(\theta) &= E[X_{n+1}|\theta] \\ &= \text{ある特定の異質性 } \theta \text{ に対する翌年の事故率の予測値}\end{aligned}$$

$p(\theta)$  のことを事前分布という言い方をします。だから、 $m$  ティルド  $\mathbf{X}$  という最適な予測値というのは、結局、各異質性  $\theta$  に対する翌年の事故率の予測値を、 $\theta$  の事後分布で平均したものなのです。

これはベイズの世界です。だから、これを計算してあげればいいのですよ。これを計算してあげればいいのですけれども、最初にこのようなものを考え出した人たちというのは、これを直接計算しようとはしていません。どうしてかと言うと、これは式で書けば簡単ですけれども、でも積分が入っているではないですか。その積分は一般的には難しいのです。だから、ここでどうするかと言うと、2通りのアプローチがあって、積分が簡単にできるような世界で考える。このようなものは、昔、ベイズでよくやっていました。共役分布というものを使います。でも、それはあまりよくないですよ。だって数学的に簡単だから、そういつてしまうのは、よくないではないですか。

## ベイズ!

- 異質性  $\theta$  を所与とするときの過去のデータ  $\mathbf{X} = (X_1, \dots, X_n)$  は、条件付き密度  $f(\mathbf{x}|\theta)$  に従うとする。また、 $\theta$  の密度を  $p(\theta)$  とする。
- このとき

$$\begin{aligned}\tilde{m}(\mathbf{X}) &= E[\mu(\theta)|\mathbf{X}] \\ &= \int \mu(\theta)p(\theta|\mathbf{X})d\theta\end{aligned}$$

であり、

$$p(\theta|\mathbf{X}) = \frac{f(\mathbf{x}|\theta)p(\theta)}{\int f(\mathbf{x}|\theta)p(\theta)d\theta}$$

- ベイズ統計学では、異質性の密度  $p(\theta)$  を**事前密度**、データが観測された後の異質性の条件付き密度を  $p(\theta|\mathbf{X})$  を**事後密度**という。
- $\tilde{m}(\mathbf{X})$  は、各異質性  $\theta$  に対する翌年の事故率の予測値  $\mu(\theta)$  を  $\theta$  の事後分布によって平均したもの

10

特に保険の世界ではデータは重要ですから、この  $f(\mathbf{x}|\theta)$  というのがどのような損失分布に従うかというところを、重要ではないですか。そこを大事にしたいとしたら、もう一つアプローチがあって、それはどうするかと言うと、 $\bar{X}$  という最適な予測値が線形モデルというように考えてしまいましょうと。簡単な形でしまいましょうと。  $a + b\bar{X}$ 。ちょうど過去のデータから将来を予測するときに、われわれはよく分からないときは回帰モデルを考えますよね。あれは線形モデルです。あれと同じです。線形モデルで考えましょうと。

## 線形ベイズ推定

- 予測量のクラスを「線形モデル」 **過去のデータの平均**

$$\tilde{m}(\mathbf{X}) = a + b\bar{X}$$

に限定する

- このとき、2乗リスク  $E[(X_{n+1} - \tilde{m}(\mathbf{X}))^2]$  を最小にする予測量は

$$\tilde{m}(\mathbf{X}) = z\bar{X} + (1-z)\mu$$

という形をとる。ここで

$$\mu = E[\mu(\theta)], \quad z = \frac{\text{Var}(\mu(\theta))}{\text{Var}(\mu(\theta)) + E[\text{Var}(\bar{X}|\theta)]}$$

11

そうすると、二乗リスクを最小にするような予測量というのは、 $z\bar{X} + (1-z)\mu$  という形となるわけです。この  $\mu$  や  $z$  というのは、そこにある形で与えられます。だから、線形モデルにすると、信頼性理論で言っているような形が出てくるわけです。特に、もうちょっと条件を加えると、ビュールマンモ

デルというものが出てきます。

ビールマンモデルはパラメータのパラメータのような、ハイパーパラメータというものを持っていて、そこをデータから推定するという形を取ります。これを経験ベイズとって、非常にプラクティカルな方法として、よく使われているということになります。これがベイズの非常に成功している、しかも分かりやすい例ではないかと思えます。

## 経験ベイズ: Bühlmannモデル

- 特に、 $\theta$  の条件付きで  $X_1, X_2, \dots, X_n$  は共通の平均  $E[X_j|\theta] = \mu(\theta)$  及び分散  $\text{Var}(X_j|\theta) = \sigma^2(\theta)$  の分布に互いに独立に従うとする

⇒ 線形ベイズ推定量は

$$\tilde{m}(\mathbf{X}) = \frac{\lambda n}{\lambda n + \phi} \bar{X} + \frac{\phi}{\lambda n + \phi} \mu$$

ただし  $\lambda = \text{Var}(\mu(\theta)), \phi = E[\sigma^2(\theta)]$

- 未知な「ハイパーパラメータ」 $\mu, \lambda, \phi$  を対応する観察データの点推定値で置き換えると Bühlmann の推定量 (経験ベイズ)。

12

この辺で、ベイズの話をお話しますが、ベイズというのは何かということ、少しまとめておきたいですね。どうしてかと言うと、私が若いころは、もうずっと前ですけども、このベイズ対頻度法の哲学的な論争がずっとあって、あたかも神学論争のようなことが行われていたのです。

そもそも、元々の問題は同じではないですか。統計的推測の問題は一体何かと言うと、データ  $X_n$  があって、これがある密度関数  $g$  に従うわけですね。データからこの  $g$  をどのように推定するかということが、統計の問題なわけですね。

頻度論はどうするかと言うと、この  $g$  というのを直接推定するノンパラメトリックな方法もありますけれども、普通は  $g$  というのは、ある統計モデルに従う、ある統計モデルに属しているというように仮定します。

この中でデータ、データというのは経験分布ですね。これに一番近いような統計モデルを選択します。このときに近いというのが問題になりますよね。何が近いのか。特にカルバック・ライブラー・ダイバージェンス、カルバック・ライブラーの距離というのがありますよね。あれを使うと、 $\theta$  ハットが最推定値になります。頻度論というのは、統計モデルを考えて、そのデータに一番近いようなものを使いましょうという話です。

では、ベイズ。何ですか、ベイズというのは。ベイズは頻度論と全く違うものでしょうか。ベイズというのは、解こうとしている問題は同じなわけですね。ベイズはどうやるかと言うと、同じように統計モデルを使います。  $f(x|\theta)$  を使います。ただし、これに事前分布という、 $p(\theta)$  のモデルを追加します。だから、ベイズというのは、 $f(x|\theta)$  と  $p(\theta)$  がペアになっている、そのようなモデルを使います。

では、その  $g$  をどのように推定するかと言うと、 $g$  から、 $g$  の新しい観測値、 $X_{n+1}$  とします。それ

の予測分布というものを作ります。予測分布というのは、ただ単に  $X_n$  が与えられたときの、 $X$  の条件付き分布です。それは計算すると、その一番下にあるような式になります。

ベイズというのは、だから、確かに事前分布を使いますが、事前分布のモデルを使っているだけで、そのような意味では頻度論を拡張したモデルを使っているということが、私が言いたいことなのですね。

## ベイズ法 vs 頻度論

### 統計的推測の問題

データ  $X_n = (X_1, X_2, \dots, X_n)$  はある密度関数  $g(x)$  に従う。  
データから  $g$  をいかに推定するか？

- 頻度論

$g$  はある統計モデル  $\{f(x|\theta), \theta \in \Theta\}$  に属すると仮定し、データ（経験分布）に最も「近い」 $f(x|\hat{\theta})$  を用いる。

- ベイズ

$f(x|\theta)$  に、事前分布  $p(\theta)$  を追加した  $\{f(x|\theta), p(\theta), \theta \in \Theta\}$  を統計モデルとし、 $g$  の新たな観測値  $X_{n+1}$  の予測分布  $f(x|\mathbf{x}_n)$  を用いる：

$$f(x|\mathbf{x}_n) = \int f(x|\theta, \mathbf{x}_n)p(\theta|\mathbf{x}_n)d\theta$$

13

とはいえ、「じゃ、 $p(\theta)$  って何だよ」という話になりますよね。私からすると、「 $f(x|\theta)$  だって何だよ」ということなのですよ。  $f(x|\theta)$  だって主観的に選んでいるわけでしょう。  $p(\theta)$  だって主観的に選んで何が悪いのか、という感じはするのですけれども、そうは言っても、 $f(x|\theta)$  の方は何とかデータのモデルなどは見えるので、少し違いますよね、 $p(\theta)$  とは。

この  $p(\theta)$ 、事前分布をどのようにして選ぶかに対しては、二つのアプローチがあります。主観的ベイズというものと、客観的ベイズ。客観的ベイズはちょっと何か変な感じがするかもしれませんが。主観的ベイズというのは、まさにこの主観確率、事前分布は主観確率ですよというような立場です。

でも、事前分布というのは、データに関するプライベートな、必ずしも公的ではないという意味で、主観的な情報を表しています。だから、個人の意思決定やビジネスの意思決定では、主観的ベイズというのはいいかもかもしれません。さっきの自動車保険の例で言えば、自動車保険には過去のデータがあるわけですから、それを事前分布としてまとめて利用するというのは、当然やるべきことですよ。

もう一つ客観的ベイズというものがあります。このとき事前分布は、ちょっと今日は説明しませんが、情報があまりないようにします。これがある意味、公平だという意味なのでも、客観的ベイズというのが、最近、比較的良好に使われるようになってきています。

### 二つのアプローチ

- 主観的ベイズ  
事前分布はデータに関する私的（主観的）な情報を表す。  
個人やビジネスの意思決定
- 客観的ベイズ  
事前分布は、「無情報的」であるように選択されるべき。  
公的な分析

14

もう一つ、事前分布の選択について、今と昔を比べます。さっきは主観と客観ですけれども。昔というのは、「じゃ昔っていつですか」という話ですけれども、MCMC以前です。ベイズというのは、MCMCというシミュレーション法が使われるようになった前と後で、もう革命的に変わりました。

昔はベイズというのは、「まあいいけど、結局計算できないでしょう」という話だったのですね。ですので、ベイズを実際使うときは、自然共役事前分布という、事前分布と事後分布が同じ分布属に属するようなものに、限って使っていました。それは積分の計算の必要がありません。ある意味、リーズナブルですよ。

今はどうかと言うと、今はMCMC以降です。MCMCというのは、マルコフ連鎖モンテカルロ法という方法なのですけれども。これは元々、実験物理の世界で使われたような、シミュレーションで使われた技術なのですけれども。これを使うことによって、次元が非常に高い、さっき Lasso の話や、あるいは機械学習の話などが出ていましたけれども、機械学習の一つの特徴は、次元が非常に高いものを使うということです。変数の数が大きいということですね。そのような場合でも、積分の計算というか、乱数が発生して、積分の計算ができるようになったというのが、MCMC以降です。

MCMC以降は何がいいかと言うと、問題や状況に応じて、自由に事前分布を選択できるようになったわけです。昔は技術的な理由で「これしか使わないよ」と。例えば、さっきの自動車保険の例なら、モデルの形を線形に制約します。今はそれは必ずしもする必要はないです。もちろん実際問題としては、そうした方が便利かもしれませんが、事前分布は自由に選択できます。それが大きな違いなのです。

なぜベイズが今これだけはやっているかと言うと、まさにMCMCのおかげと言っても過言ではないのですね。時代でいうと、多分2000年、21世紀になってからは、ベイズの世界というものが、今の統計学の世界に、ぱっと出てきたという感じです。もちろんMCMCがすべてを解決するわけではないのですけれども、非常に大きなインパクトを与えているというのが、現在の状況ではないかなと思います。



- 昔 (MCMC 以前)  
自然共役事前分布: 各尤度関数に対して, 事前分布と事後分布が同じ分布族に属する  
⇒ 積分計算の必要なし.
- 現在 (MCMC 以降)  
問題や状況に応じて自由に事前分布を選択

15

次の話題なのですが、これは私自身がやっていることを、少しご紹介させていただきたいと思います。要介護状態別死亡率の予測ということです。この辺も皆さんはご専門家の方なので、あえて説明することもないかとは思いますが、長寿、平均余命がどんどん伸びていて、長寿リスクというのが非常に問題になっています。長寿であるということは、同時に要介護状態になるというリスクがあるわけで、だから、ロングタームケア・リスクというのが、大きな注目の的になるわけです。だから、死亡率というものを考えるときも、要介護状態に応じた死亡率というものを計算したいと、将来の死亡率を知りたいと。

だけれど、実際に見てみると、国民レベルの個票のデータでは、ある地域に限って言えば、そのようなデータが利用可能なケースもだんだん増えてきているようですけれども、国民レベルの健康状態別の死亡データというのは公開されていないという中で、どうやって要介護度別の死亡率を推定しようかということです。

## 要介護度別死亡率

- 過去の予想を上回る平均余命の伸長が世界的に進行中
  - 平均余命の伸長は、長寿リスクとともに、要介護状態となる可能性—長期介護リスク (LTC)—の増大を意味する。
- ⇒ 将来の死亡率がいかに要介護状態と関係しているかの把握が重要
- しかし、(国民レベルの) 健康状態別の死亡データは公開されていない。

17

だから、目的としては、要介護状態別の部分集団に対する死亡率モデルを開発したい。だけれど、死亡率のデータ、死亡データはないです。あとは平均余命がどんどん伸びているというの也要考えなければいけません、ということが課題です。

## 課題

### 目的

要介護状態別の部分集合に対する死亡率モデルの開発。  
ただし

1. 要介護状態別の死亡率データは存在しない
2. 平均余命の伸長を考慮

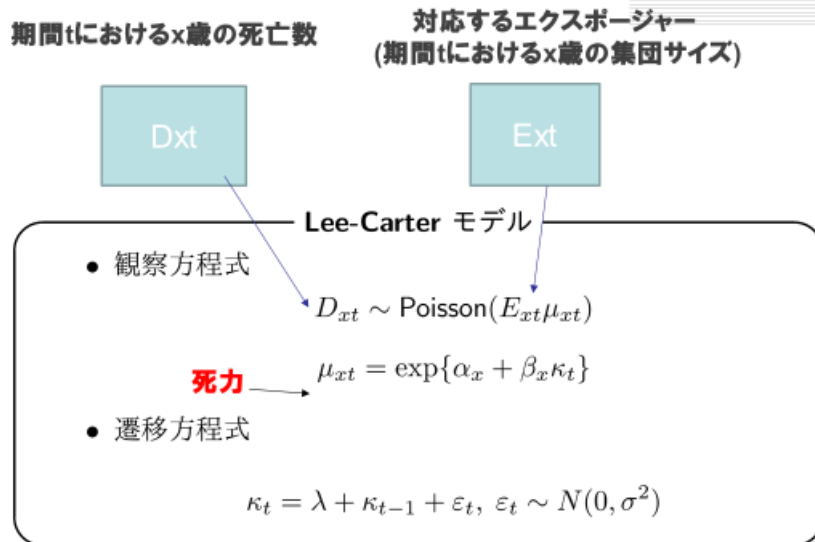
18

これに対して、死亡率の予測と言うと、一番よく使われる代表的なモデルは、リー・カーター・モデルと呼ばれるものですよね。リー・カーター・モデルは、ここに書いてあるとおり、少しだけ復習をしますと、データがあって、 $D_{x,t}$ というのは死亡のデータです。 $E_{x,t}$ がエクスポージャーです。

リー・カーターは二つの方程式があって、一つはその $D_{x,t}$ のデータのモデルで、これは観察方程式といいます。 $D_{x,t}$ というのはポアソンで平均が $E_{x,t}$ 、エクスポージャー×死力。この死力の形がいわゆるリー・カーター・モデルの形です。 $x$ が年齢を表わしていて、 $t$ が期間です。だから、 $\mu_{x,t}$ は基本的

には  $\alpha_x$  というパラメータが  $\mu_{xt}$  を決めるのだけでも、それからずれている部分というのが、 $\beta_x \kappa_t$  で表されていて、 $\kappa_t$  は一般的にどんどん死亡率が下がっているということを表わすパラメータで、ですので、この  $\kappa_t$  について、下にあるような状態方程式で書いてありますけれども、ランダムウォークを仮定しています。ただし、 $\kappa_t$  も  $\alpha_x$  も  $\beta_x$  も観測されないという、このようなモデルになっています。

## 全人口集団の死亡率予測: Lee-Carterモデル



Lee and Carter, 1992

19

これは  $D_{xt}$  と  $E_{xt}$  があれば、推定できますよね。だから、要介護状態による部分集団に分けて、同じことをやれば、できるはずですよ。けれども、問題は状態0というのが要介護を必要としない状態で、1~J と、J が一番重度が高いような、このように部分集団に分けたとして、各部分集団について、今のリー・カーター・モデルを当てはめてやればいいではないですか。

でも、 $D_{xt0}$ 、 $D_{xt1}$ 、 $D_{xtj}$  というこの右側の部分がないわけですよ。その中でどのように推定したらいいのかという問題を考えてみました。どうやるかと言うと、モデリング。モデリングに頼りましょう。死力をモデリングします。

一番上に書いてあるのは、やはりどのような部分集団についても、リー・カーター的な構造があるでしょう。けれども、状態によって違うので、その違いを表わすようなパラメータを入れて、それを  $\eta_j$  としています。ただし、 $\eta_j$  というのはだんだん  $j$  が大きくなるにつれて、大きくなるような制約を置いています。このような制約を置いた上で、推定をしましょう。さらに、一個一個の死力の和が、加重平均が全集団の死力になるというように仮定を置きましょう。このときの加重は人口の総対比率です。

## 要介護状態(健康状態)による部分集団



- 各部分集団について、 $D_{xtj}$  と  $E_{xtj}$  が利用可能ならば、各状態  $j$  について死亡率予測モデルをフィットすればよい。
- しかし、しばしば、**部分集団の死亡データ  $D_{xtj}$  は利用不可能**。

20

これを仮にミクスチュア・リー・カーター・モデルと呼んで、これをさっきと同じような方式で推定するというを試みています。これは、これだけを見ると簡単そうですが、さっきの制約を表わすような事前分布を入れたりしなければいけないので、結構それなりに大変ということになります。

## 仮定: 死力のモデリング

- 各部分集団の死力は

$$\mu_{xtj} = \exp \{ \gamma_x + \eta_j + \beta_x \kappa_t \}$$

通常の Lee-Carter に  $\eta_j$  を追加  
(Li and Lee, 2005)

ここで、 $\eta_j$  は (不) 健康要因であり、

$$0 = \eta_0 < \eta_1 < \dots < \eta_J$$

- 全集団の死力は

$$\mu_{xt} = \sum_{j=0}^J w_{xtj} \mu_{xtj} = \sum_{j=0}^J w_{xtj} \exp \{ \gamma_x + \eta_j + \beta_x \kappa_t \}$$

ここで:

$$w_{xtj} = \frac{E_{xtj}}{E_{xt}} = \text{期間 } t \text{ における } j \text{ 集団の } x \text{ 歳人口の相対比率}$$

21

ということで、最尤法でやっても可能なので、実際に試みましたが、結構複雑で、なかなかいい解が出てきませんでした。それもあって、ベイズ法を使って、推定しました。ベイズ法のいい所はパラメータ不確実性という、その予測をするときに、モデルを推定したということの誤差まで含めてリスクを把握しなければいけないという部分も含めて、予測できるということで、ベイズ法がいいと一般に言われていて、ケアンズやブレイクなどはそのような論文を書いていますし、ベイズ法で死亡率の推計をしようという試みは、いろいろあるというわけです。というわけで、そのような方向に沿って、ここでもベイズ法で推定したということです。

### Mixture Lee-Carter model

- 観察方程式

$$D_{xt} \sim \text{Poisson}(E_{xt}\mu_{xt})$$

$$\text{with } \mu_{xt} = \sum_{j=0}^J w_{xtj} \exp\{\gamma_x + \eta_j + \beta_x \kappa_t\}.$$

- 状態方程式

$$\kappa_t = \lambda + \kappa_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N(0, \sigma^2)$$

$\eta_j$  がすべて等しい :

$$\eta_0 = \eta_1 = \eta_2 = \dots = \eta_J = \eta, \text{ say,}$$

とき,  $\alpha_x = \gamma_x + \eta$  とすれば, 通常の Lee-Carter モデル.

22

## 死亡率予測モデルとパラメーター不確実性

- Lee-Carter モデルの通常の推定 (最尤法) では, 推定パラメータの誤差 (パラメータ・リスク) を考慮できない.
- ⇒ 将来の死亡率のリスクを過小評価  
Cairns et al. (2006) and Blake et al. (2008)
- ⇒ ベイズ法による考察  
Czado et al. (2005), Pedrosa (2006), Kogure, Kitsukawa and Kurachi(2008), Kogure and Kurachi (2010), etc.
- 本報告では, ベイズ法によって提案した混合 Lee-Carter モデルの推定を行う.

23

実際のデータなのですけれども、日本の介護保険のデータを応用しました。介護保険というのは、スタートした時点は要支援と、それから要介護1～要介護5の6段階に分かれていたのですが、途中で、2006年に要支援の所が1と2に分かれました。その時に、要介護1とのかねあいで、要支援1・2を設定したということなので、ここでは2001年～2005年までは状態1を、そこの表に書いてあるとおり、要支援+要介護1として、2006年以降は要支援1～2と要介護1を一つにまとめて状態1として、それ以外については、状態2～状態5というのは、それぞれ要介護2～要介護5に対応するというカテゴリズをして、状態0～状態5のときの死亡率の予測を行いました。

## 我が国介護保険データへの応用

- 介護保険制度は 2000 年 4 月に開始。開始時は要介護の状態に応じて受給者は 6 段階（要支援，要介護 1 ～ 要介護 5）に区分。
- 2006 年度から，要支援の区分は（要介護 1 との調整も行い）要支援 1 と要支援 2 に細分化。このとき経過的要介護が設けられた。

本研究では，全期間のデータの整合性のために，以下のように状態を設定：

期間	状態 0	状態 1	状態 2	状態 3	状態 4	状態 5
2001～2005	非要介護	要支援+要介護 1	要介護 2	要介護 3	要介護 4	要介護 5
2006～2008	非要介護	要支援 1～2+要介護 1 +経過的要介護	要介護 2	要介護 3	要介護 4	要介護 5
2009～2014	非要介護	要支援 1～2+要介護 1	要介護 2	要介護 3	要介護 4	要介護 5

以下では，状態 0～状態 5 の 6 状態の死亡率予測を行う。

24

## 要介護者受給者数

- 要介護（要支援）状態区分・性・年齢階級別の受給者数（『介護保険事業状況報告』）を要介護部分集団の人口サイズとした。
- 65 歳～89 歳までの，受給者数は 5 歳区間の年齢区分で与えられる。

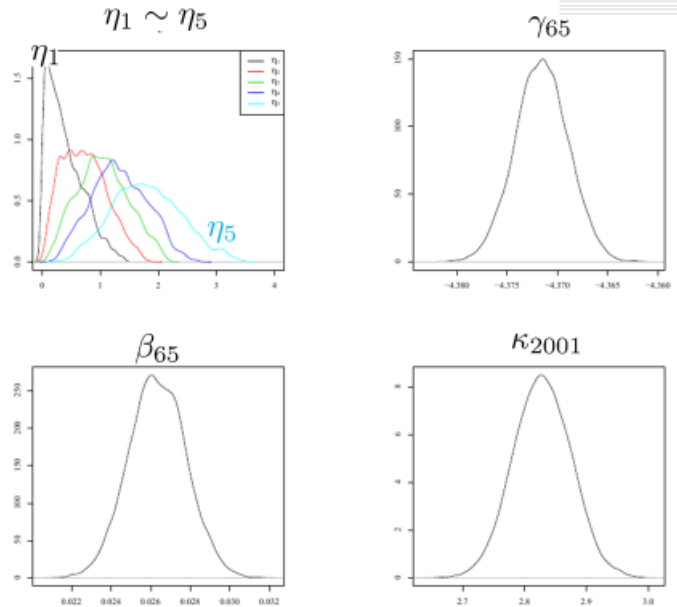
65～69, 70～74, 75～79, 80～84, 85～89,

- これらの年齢区分の受給者数の線形補間を行うことにより，各年齢における受給者数を求め， $\{E_{xtj}\}$  を作成

26

ベイズを使って、MCMCを使って、予測したのですけれども。これはスキップしますね。これもちょっとスキップします。これが結果です。 $\eta 1$ 、 $\eta 5$ というのは、さっき言った要介護状態を表わすパラメータです。 $\eta 1$ が一番軽度で、 $\eta 5$ が一番重度です。そうすると、ここにあるように、これは $\eta 1 \sim \eta 5$ の事後分布なのですけれども、だんだん重度が上がるごとにシフトしているということが分かると思います。他のパラメータについても、とてもきれいに推定されているということになります。

事後分布:  $\eta$ ,  $\gamma_{65}$ ,  $\beta_{65}$ ,  $\kappa_{2001}$ : 男性

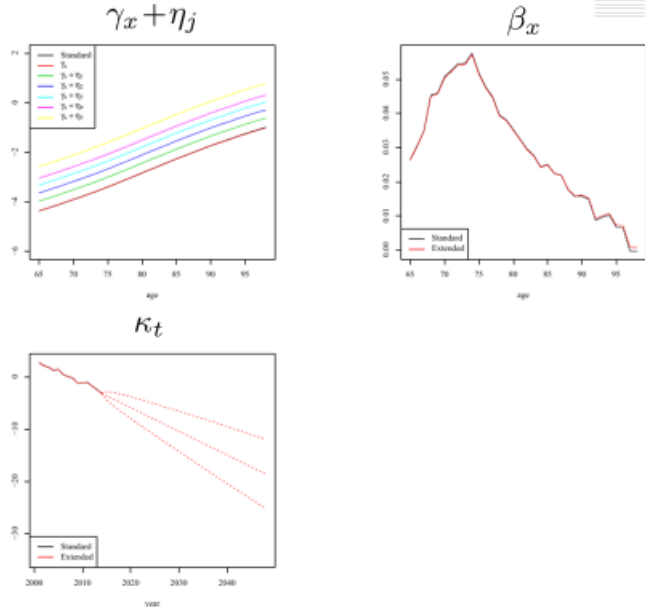


27

これが実際の推定結果の要約統計量なのですが、これを見ても  $\eta_1 \sim \eta_5$ 、これは平均値ですけれども、だんだん大きくなっているということで、比較的うまく推定できているのではないかと考えています。

これは各パラメータの時間的な推移、あるいは年齢に対する推移を表わしたものです。

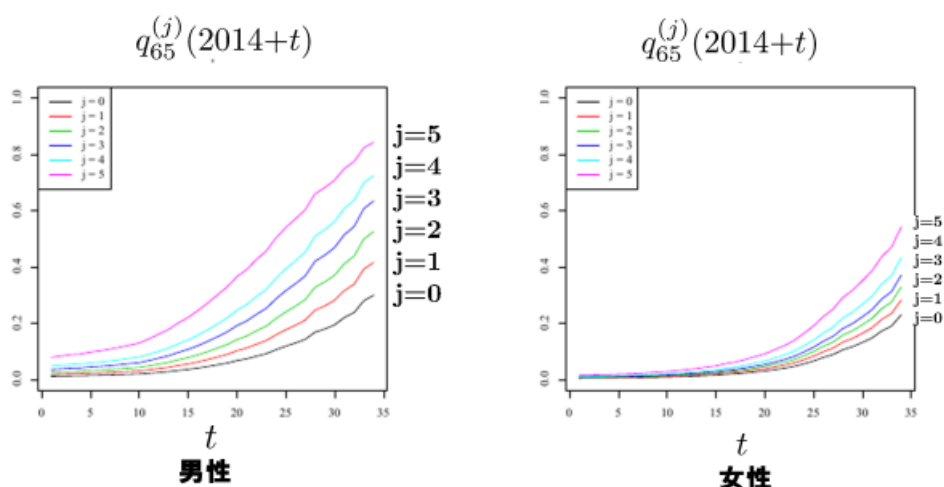
パラメータ:  $\gamma_x$ ,  $\beta_x$ ,  $\kappa_t$ の事後平均の推移: 男性



29

でも、これ自体はパラメータなので、本当に関心があるのは将来の死亡率ですよね。これを予測した結果、将来の死亡率を予測するためには、将来どうなるかという部分も予測しなければいけないのですが、そこも予測分布というものを使って、この  $\kappa$  のところですね。ランダムウォークのところを予測して、死亡率の推定をしてみましたというのが、この結果です。

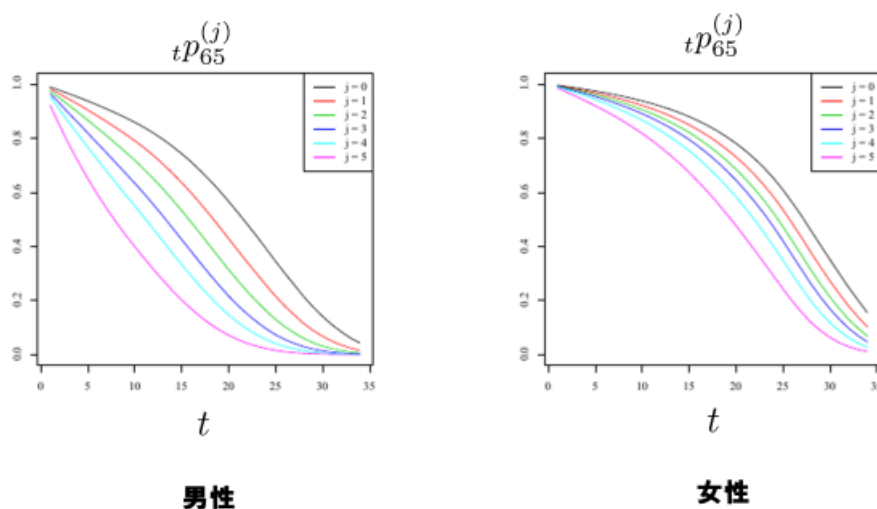
## 35年後までの将来死亡率の予測



31

これまでの結果というのは男性だけだったのですけれども、ここでは女性の結果も書いています。2014年に65歳であった人の将来の死亡率というのを、 $j=0$ というのは支援を必要としない人です。  $j=5$ が一番重度が高い人。その死亡率はこのようになります。このモデルが正しければ、このようになるということです。同じように、将来の生存率の予測、生存関数の予測もしています。これも65歳の人の将来の生存率がどうなるかというのを、各要介護状態別に推定したというものであります。

## 35年後までの将来生存率の予測



33

ただし、ちょっと戻りますけれども、これはちょっとまだプリミティブというか、プレリミナリーな研究なので、状態が全然変わらないまま死ぬという、そのように非現実的な仮定で推定しているので、実際にはその間で状態を遷移するということも、もちろん推定して、考慮しなければいけないのですけれども。そこのところまではまだやっていないので、そのような若干非現実的な想定で計算されている部分も



ありますという、一定の考慮はあります。

このような形で、たとえデータがなくとも、うまくモデリングをして、ベイズを使うことで、比較的、何となくもっともらしい、この絵を見てもっともらしいと言っても「え？」という感じかもしれませんが、でも、何かもっともらしくないですか。死亡率、生存率。何か、もっともらしくないですか。もっとぐちゃぐちゃと出てくるかと思ったのですけれども、うまく出てくるのですね。そのように制約していると言われれば、そうなのかもしれませんけれども。でも、モデリングとはそのようなものですよ。ちょっと開き直りですかね。このような形で推定できますと。

時間がちょっとなくなってきてしまったのですけれども。最後はちょっと駆け足になってしまいますけれども。ちょっと最後に言いたいといえますか、お話ししたい。2、3分いいですか。大丈夫ですか。

モデルリスクという、皆さん、言葉を聞いたことがありますかね。これはちょっと英語が小さくて読めないかもしれませんが、ボックスという有名な統計学者がいて、ボックスは有名なフィッシャーの娘婿なのですけれどもね。フィッシャーには娘さんが2人いて。どうでもいいかもしれませんが。長女と結婚したのがこのボックスで、次女と結婚したのが、私がアメリカでお世話になったアンスコム先生です。

ボックスがそのような、まさに統計学者の本流です。彼が言っている有名な言葉に、このような言葉があって、「すべてのモデルは間違ってる」。すごいですよね。「すべてのモデルは間違ってる。だけど、中には役に立つものもある」というように、ボックスは言っています。

モデルの否定のように見えますけれども、モデルの本質をととても言い当てていますよね。だって、モデルというぐらいだから、現実ではないです。間違っているに決まっています。でも、統計学で何かを語る時、われわれはモデルが正しいかのように考えて、語ってきたことがとても多いではないですか。それがよくないというのが、例えば赤池先生のAICが出てきて、気づかされるわけですから。やはり、そこはとても重要だなというのは、私自身は思っています。

そのリスクを計量化するときに、ケアンズというスコットランドにいる有名なアクチュアリーの研究者がいますけれども、彼は三つのリスクを考えなければいけないと言っています。プロセスリスクは本質的な、そもそものリスクですね。それから、パラメータリスク。モデルを使うわけですから、そのモデルを推定したとして、その推定のリスク、パラメータの不確実性。さらにモデル。どのようなモデルを使ったかという、モデルリスクもあるのだということなのです。

## モデル・リスク

リスクの計量化における3種類のリスク (Cairns, 2000)

- **プロセス・リスク**  
用いるモデルの確率的特性に基づく不確実性
- **パラメータ・リスク (パラメータ不確実性)**  
用いるモデルに含まれるパラメータの値の不確実性
- **モデル・リスク (モデル不確実性)**  
用いるモデル自体の不確実性 (「真の」確率分布がモデルに含まれていないかもしれない)

様々な死亡率モデルの開発に伴い、モデル・リスクへの関心が増大 (Haberman et al., 2014; Yang, Li and Balasooriya, 2015)

35

AICを使ってモデルを選択する。それでいいような気がしますけれども、AICでモデルを選択すること自体が、そもそも不確実性をはらんでいるわけです。どのようなモデルを使ったかというモデルリスクも、リスクの計量化の中に入れなければいけないよということを、ケアンズは言っています。

特に、私はこの死亡率モデルということに興味を持って、このところやっているのですけれども、最近、非常に多くの死亡率モデルというものが開発されているのですね。その中でモデルリスクというものの関心がとても増えています。死亡率のモデルを作って、それをもとに長寿リスクのデリバティブズのようなものを考えるわけですが、そのときに、モデルリスクというものも考慮に入れて、長寿デリバティブのプライシングをやらなければいけないというようなことを、最近イギリスのアクチュアリーの人たちの中には、そのようなことを言っている人がいるらしいです。

このような死亡率モデルというのは、状態空間モデルというもので記述できます。さっきのリー・カーター・モデルも状態空間モデルなのです。この  $y$  の所の  $f$  という、それから  $\phi$  の所のモデリングをいろいろ考えることで、例えば観測方程式として、リー・カーター・モデル以外にAPC、Age-Period-Cohortや、CBD、それからリー・アンド・リーというモデルなど、いろいろモデルは考えています。

それから遷移方程式の方も、さっきのモデルのランダムウォークは、確率的トレンドモデルと言われますけれども、定常モデルや、確定的トレンドモデルや、あるいはもっと複雑な時系列モデルなど、いろいろあるわけです。

## 死亡率モデルと状態空間モデル

- 確率的死亡率モデルは状態空間モデルとして記述できる (Fung, Peters and Shevchenko, 2016).

$$\begin{cases} \text{観測方程式} & \mathbf{y}_t \sim f(\mathbf{y}_t | \phi_t, \boldsymbol{\theta}) \\ \text{遷移方程式} & \phi_t \sim f(\phi_t | \phi_{t-1}, \boldsymbol{\theta}) \end{cases}$$

ここで、 $\mathbf{y}_t$  は時系列データであり、 $\phi_t$  は潜在変数

Lee-Carter モデルでは、 $\mathbf{y}_t = (\log m_{x_{\min}t}, \dots, \log m_{x_{\max}t})'$ ,  $\phi_t = \kappa_t$

- 主要な代替モデル
  - 観測方程式  
APC モデル, CBD モデル, Li and Lee モデル
  - 遷移方程式  
確率的トレンドモデル, 定常モデル, 確定的トレンドモデル

36

このようなものをどうやって予測するかというときに、AICを使うのも確かにいいかもしれないけれども、モデル選択はAICでやりましょうというような形で出てきますけれども、ベイズだとモデル選択自体が、ベイズの枠組みの中に入っているということ、これは言いたかったことなのです。周辺尤度を使うということになるのですね、モデル選択をする場合には、モデルの事後確率というのを求めて、それは周辺尤度に比例しているのです。

## モデル選択：周辺尤度

代替モデルを  $\mathcal{M}_1, \dots, \mathcal{M}_J$  とする。

- 各モデルはデータ  $\mathcal{D}$  のモデルとパラメータ  $\boldsymbol{\theta}$  の事前分布の組

$$\mathcal{M}_j = \{(f_j(\mathcal{D} | \boldsymbol{\theta}), f_j(\boldsymbol{\theta})), \boldsymbol{\theta} \in \Theta_j\}$$

- 各モデルの事前確率は等しいと仮定：

$$\Pr(\mathcal{M}_j) = \frac{1}{J} \quad (J \text{ 個のモデルのうち一つは真})$$

- モデル  $\mathcal{M}_j$  の事後確率は

$$\begin{aligned} \Pr(\mathcal{M}_j | \mathcal{D}) &= \frac{f(\mathcal{D} | \mathcal{M}_j) \Pr(\mathcal{M}_j)}{f(\mathcal{D})} \\ &\propto f(\mathcal{D} | \mathcal{M}_j) = \int_{\Theta_j} f_j(\mathcal{D} | \boldsymbol{\theta}) f_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \text{周辺尤度} \end{aligned}$$

⇒ 周辺尤度を最大にするモデルが望ましい。

37

その周辺尤度というのは、実は展開していくと、オッカム・ファクターと呼ばれますけれども、その複雑なモデルをちょっと罰則するような項が自動的に含まれています。このオッカム・ファクターの所をさらに展開すると、いわゆるBICというものになっていきます。

## 周辺尤度とオッカムの剃刀

- 事後分布が正規分布で近似できるとすると：

$$\begin{aligned} \text{周辺尤度} = f(\mathcal{D} | \mathcal{M}_j) &= \int_{\Theta_j} f_j(\mathcal{D} | \theta) f_j(\theta) d\theta \\ &\approx \underbrace{f_j(\mathcal{D} | \hat{\theta}_j)}_{\text{データに適合させた尤度}} \underbrace{f_j(\hat{\theta}_j) |\Sigma_{\hat{\theta}_j} / 2\pi|^{1/2}}_{\text{オッカム・ファクター}} \end{aligned}$$

ここで、 $\hat{\theta}_j$  は最大事後確率推定値であり、

$$\Sigma_{\hat{\theta}_j} = \left( -\frac{\partial^2 \log [f_j(\mathcal{D} | \theta) f_j(\theta)]}{\partial \theta \partial \theta'} \right)^{-1} \Bigg|_{\theta = \hat{\theta}_j}$$

- 「オッカム・ファクター」はモデルの複雑さに対する罰則項。

38

というわけで、ベイズのモデリングというのは、モデルリスクにも対応できるのではないかということ  
を最後に付け加えて、私の発表としたいと思います。どうもありがとうございました。

## まとめ／結論

- ベイズ法の応用として「信頼性理論」と「要介護度別死亡率の予測」を紹介した。
  - 前者は、各個人の限定されたデータと関連する外部のデータを組み合わせるために利用。事前分布は各個人の異質性を表す。**主観的ベイズ**。
  - 後者は、パラメータ・リスクやモデル・リスクを取り込むためにベイズ法を利用。事前分布はモデルの一部。**客観的ベイズ**。
- 「計算ベイズ技術」と「客観的ベイズ」の普及によって、ベイズ統計学はデータサイエンスの「実践」の重要な柱
- リスク管理の実践においても、ベイズ統計学のより一層の活用を期待したい。

39

山内 どうもありがとうございます。多分皆さん、お聞きになりたいことは山ほどあるかと思うのですが、あと5分程度しかございませんので、お二人程度お受けしたいと思います。

岩沢 皆さん、貴重な講演、ありがとうございました。私としては、本当はディスカッションが良いのですが、ちょっと私としてはメッセージを言わないといけないところがあって。

小林さんの発表自体は別に問題はなかったと思うのですが、ちょっと僕は誤解を与えるところが

あると思って。できれば、小林さん、すぐ出るなら出してほしいのですけれども、出ないですか、小林さんのスライドは。

ですので、口頭で言うと思うのですけれども、Lasso の話があって。私も講演の中で機械学習から学べる手法の一つと言ったのですけれども、原論文を見れば、彼は本当は分かっているはずだと思うのですけれども、原論文は統計学の論文で出たものなのです。逆輸入なのですね。

Lasso という方式は、統計学の人を考え出したのだけれども、スパースモデリングで使えると、後で分かったのですよ。今日の発表だと、「スパースモデリングに使うのは Lasso がいいという提案が出た」という説明になっていたのですけれども、違うのです。あれは統計学のものなのです。それを「スパースモデリングで使えるよ」と分かったものだから、機械学習の人が飛びついたのです。今は学ぶのは機械学習の教科書から学ぶのです。今、学ぼうと思うと。でも、あれは統計学の手法なのですね。

ですので、実際にはそのような逆輸入のものが幾つもあり、実はあつたりするので、そのように今日の話も捉えてもらえると良いと思います。口頭で、ちらりと最後に言ったので、もちろん彼はもう知っているわけですが、原論文にそのベイズ的な解釈が入っているのですね。Lasso の。そうですね。そのように口頭ではおっしゃっていたのだけれども。このスライドだけを見ると、それが後から、ベイズ的にも解釈できるといように誤解されてしまうと思うのですけれども。

そのように、本当に統計学の手法が機械学習の人が先に使ってしまったという。だから、われわれはそれを取り戻して、使うべき時期に来ているのではないかなということ、ちょっと補足として、お願いします。

**山内** では、お一人、いかがでございますか。では、重原さん。

**重原** すみません。紙に穴を開けてコンピュータを使った最後の世代ぐらいの重原。第一生命経済研究所の重原と申します。

今後のアクチュアリーにコーディングが必要という、あるいはコンピュータ技術が必要というときに、特に古い人だと、新しい言語を覚えれば良いのかなぐらいに思う危険性が非常にあると思うので、そのところは恐らく違うのだということ。そのことについて、補足でご説明していただければと。簡単に言えば、R と Python をという、新しい言語を学べば良いという話ではないですよ。だから、キーワードになるのは、恐らくモジュールや何かをうまく使う技術などというような話ではないかなとは、思っているのですけれども、その辺の話をお聞かせいただければ。

**山内** 例えば、データを投入して、ただ答えが出れば良いという話。

**重原** そういったプログラムが書ければ良いのは、ちょっと理屈が違うよねという話です。

**山内** なるほど。多分、お立場によって回答が変わると思うのですけれども、では藤澤さんの方から一つお願いします。どう思われますか。

**山内** いや、これはご質問というより、「どう思われますか」だと思うのですよね。

**藤澤** 私が使ったことがあるのはRですが、Rを使うといろいろなパッケージがあつて、データを入れると結果が出てくるというところはあつて、それほどプログラムの知識がなくても、統計的なモデルが実装できるというメリットはあると思います。

ただ、それだけで本当にいいのかと言うと、やはり別の話もあつて。ちょっとお答えになっているか分からないですけれども、やはり理論的なところも、どのようにして結果を解釈するのかというところは、多分とても重要で、そこを間違ってしまうと、間違ったそのモデルや、モデリングの結果を実務に使ってしまうという危険性があると思うので、そこは多分、両方必要になってくるのではないのかなと思っています。

特に言語にこだわっているというわけではなくて、Rでも Python でも Java でも、何でもいいと思うのですけれども、出てきた結果を正しく解釈するというところも、必要になってくるのではないのかなというのが、実務的な観点での意見です。

**山内** では、続けてお願いします。

**小林** 先ほどの発表でもちらっと述べましたが、何がしたいかによって、コーディングというのは、しょせん手段なので、何がしたいかによって変わってくるのかなと思いますが、どのようなことに特に興味があるのでしょうか。

もし仮に、コーディングが得意で、何も制約がないのであれば、個人的にはCやC++などを書けると、とてもいいのかなと思いますが、そこまで勉強するのは大変だと思うのであれば、やはりRやPythonがとても導入しやすいのかなと思います。

**重原** そのような意味では、どちらかと言うと、それ自体に興味があるというより、結果に興味がある人が、恐らくこの中にいる人はほとんどだと思うので、そのような意味では、RやPythonを使うということになるのでしょうかけれども、あれですよ。先ほどコミュニティと言っていましたけれども、いわゆる、コミュニティに蓄積されたルーチンや何かを使っていくということが重要になってくるのではないかなと、質問の趣旨はそのようなところだったのですけれども。

**小林** それもあまりいい回答ができないですけれども、自分は比較的実務とかけ離れた所にいるので、新しいものを作るには自分で全部コーディングしなければいけないという状況がかなり多いので、うまく回答も出せていないですね。ディープラーニングなどをやるのであれば、ライブラリだと PyTorch や Keras など、そのような大学が開発しているライブラリのようなものがあるので、そのようなものをうまく利用していった方がいいというのは、そのとおりでと思います。

**山内** では小暮先生、ぜひお願いします。

**小暮** 私も最近脳活にいいかと思い Python を始めたのですけれども、なかなか進んでいなくて、あれですけれども。

私はコーディングが本質的ではないというのは、確かにそうだと思うのですけれども、一方でコーディングのパワーというものに、とても実は最近思うところがあつて。やはり実装できる能力ということだと

思うのですね、コーディングできるというのは。それはやはり侮れないし、これからのデータ・サイエンティストと呼ばれる人にとっては、結構本質的な部分にかかわるのではないかなと思っています。ただし、自分はもうできません。

**山内** どうもありがとうございました。もうお時間ということになりました。それでは最後に、藤澤さん、それから小林さん、小暮先生、お三人に盛大な拍手をお願いいたします。