

# モード2科学としてのデータサイエンスと YOKOHAMA D-STEP

坂巻 顕太郎・汪 金芳

横浜市立大学

第 17 回 統計教育の方法論ワークショップ



## 1 Data Science as Mode 2 Science

- Science Revisited
- From Data Analysis to Data Science
- Teaching Data Science
  - Teaching Mode 2 Data Science at YCU
  - Data Science Club: Learning by Practicing
  - Data Science PBL for Juniors
- Practical Data Science in Master's Program

## 2 YOKOHAMA D-STEP

## 1 Data Science as Mode 2 Science

- Science Revisited
- From Data Analysis to Data Science
- Teaching Data Science
  - Teaching Mode 2 Data Science at YCU
  - Data Science Club: Learning by Practicing
  - Data Science PBL for Juniors
- Practical Data Science in Master's Program

## 2 YOKOHAMA D-STEP



## Definition (SCIENCE)

The intellectual and practical activity encompassing the systematic study of the structure and behaviour of the physical and natural world through observation and experiment.

- ① a particular area of science.
- ② a systematically organized body of knowledge on a particular subject.
- ③ archaic knowledge of any kind.





## Definition (SCIENCE)

Science (from the Latin word *scientia*, meaning “knowledge”) is a systematic enterprise that builds and organizes knowledge in the form of testable explanations and predictions about the universe.

Modern science is typically divided into three major branches.

- ① natural sciences (e.g., biology, chemistry, and physics);
- ② social sciences (e.g., economics, psychology, and sociology)
- ③ formal sciences (e.g., logic, mathematics, and theoretical computer science).

Disciplines that use existing scientific knowledge for practical purposes, such as engineering and medicine, are described as **applied sciences**.



## Definition (Mode 1 Science)

The essential ingredients of Mode 1 Science are

- ① academic context;
- ② disciplinary;
- ③ homogeneity;
- ④ autonomy;
- ⑤ traditional quality control(peer review),

## 1 Data Science as Mode 2 Science

- Science Revisited
- From Data Analysis to Data Science
- Teaching Data Science
  - Teaching Mode 2 Data Science at YCU
  - Data Science Club: Learning by Practicing
  - Data Science PBL for Juniors
- Practical Data Science in Master's Program

## 2 YOKOHAMA D-STEP



**4. Sciences, mathematics, and the arts.** The extreme cases of science and art are clearly distinguished, but, as the case of the student who was eligible for Phi Beta Kappa because mathematics was humanistic and for Sigma Xi because it was scientific shows, the place of mathematics is often far from clear. There should be little surprise that many find the places of statistics and data analysis still less clear.

There are diverse views as to what makes a science, but three constituents will be judged essential by most, viz:

- (a1) intellectual content,
- (a2) organization into an understandable form,
- (a3) reliance upon the test of experience as the ultimate standard of validity.

# “Data Analysis” as Science

John W Tukey (1962). The future of data analysis. *AMS*, 1-67.

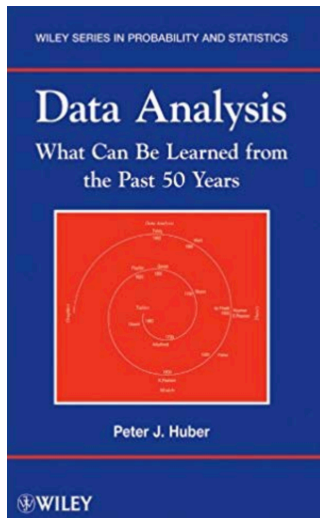


## I. GENERAL CONSIDERATIONS

**1. Introduction.** For a long time I have thought I was a statistician, interested in inferences from the particular to the general. But as I have watched mathematical statistics evolve, I have had cause to wonder and to doubt. And when I have pondered about why such techniques as the spectrum analysis of time series have proved so useful, it has become clear that their “dealing with fluctuations” aspects are, in many circumstances, of lesser importance than the aspects that would already have been required to deal effectively with the simpler case of very extensive data, where fluctuations would no longer be a problem. All in all, I have come to feel that my central interest is in *data analysis*, which I take to include, among other things: procedures for analyzing data, techniques for interpreting the results of such procedures, ways of planning the gathering of data to make its analysis easier, more precise or more accurate, and all the machinery and results of (mathematical) statistics which apply to analyzing data.

# Huber (2011) on Tukey (1962)

Peter J. Huber (2011). Data Analysis: What Can Be Learned From the Past 50 Years





Half a century ago, Tukey, in an ultimately enormously influential paper redefined our subject... [The paper] introduced the term “data analysis” as a name for what applied statisticians do, differentiating this term from formal statistical inference. But actually, as Tukey admitted, he “stretched the term beyond its philology” to such an extent that it comprised all of statistics.

# Driving Forces on Tukey's "Data Analysis"

David Donoho (2015,

<https://courses.csail.mit.edu/18.337/2015/docs/50YearsDataScience.pdf>)



Tukey identified four driving forces in the new science:

- 1 The formal theories of statistics;
- 2 Accelerating developments in computers and display devices;
- 3 The challenge, in many fields, of more and ever larger bodies of data;
- 4 The emphasis on quantification in an ever wider variety of disciplines.





## Science and data science

「統計的」、「計算的」、「人間的」の視点から  
データサイエンスを俯瞰することが出来よう。  
それぞれの視点がデータサイエンスを構成する不可欠な成分であるが、  
これらの成分の有機的結合こそがデータサイエンスのエッセンスである。

David M. Blei<sup>a,b,c,1</sup> and Padhraic Smyth<sup>d,e</sup>

Edited by Peter J. Bickel, University of California, Berkeley, CA, and approved June 16, 2017 (received for review March 15, 2017)

Data science has attracted a lot of attention, promising to turn vast amounts of data into useful predictions and insights. In this article, we ask why scientists should care about data science. To answer, we discuss data science from three perspectives: statistical, computational, and human. Although each of the three is a critical component of data science, we argue that the effective combination of all three components is the essence of what data science is about.

data science | statistics | machine learning

The term “data science” has attracted a lot of attention. Much of this attention is in business (1), in government (2), and in the academic areas of statistics (3, 4) and computer science (5, 6). Here, we discuss data science from the perspective of scientific research. What is data science? Why might scientists care about it?

Our perspective is that data science is the child of statistics and computer science. While it has inherited some of their methods and thinking, it also seeks to blend them, refocus them, and develop them to

help them more effectively navigate and understand the contours of society, finding relevant sources to their work and identifying hard to spot patterns of language that suggest new interpretations and theories. Third, modern telescopes create digital sky surveys that have transformed observational astronomy, generating hundreds of terabytes of raw image data about billions of sky objects. A catalog of these objects, if available, would give astronomers an unprecedented window into the structure of the cosmos.



**Data science is science's second chance to get causal inference right.**

**A classification of data science tasks**

Miguel A. Hernán,<sup>1,2</sup> John Hsu<sup>3,4</sup>, Brian Healy<sup>5,6</sup>

1. Departments of Epidemiology and Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA
2. Harvard-MIT Division of Health Sciences and Technology, Boston, MA
3. Mongan Institute, Massachusetts General Hospital, Boston, MA
4. Department of Health Care Policy, Harvard Medical School, Boston, MA
5. Department of Neurology, Harvard Medical School, Partners MS Center, Brigham and Women's Hospital, Boston, MA
6. Biostatistics Center, Massachusetts General Hospital, Boston, MA

Correspondence: Miguel Hernán, Department of Epidemiology, Harvard T.H. Chan School of Public Health, 677 Huntington Avenue, Boston, MA 02115; email:

[miguel\\_hernan@post.harvard.edu](mailto:miguel_hernan@post.harvard.edu)



## Abstract

「データサイエンス」という用語の導入は、  
データ解析の再定義を行う歴史的機運をもたらし、  
観察データに基づく因果推論を推し進める自然な枠組みを与えよう。

Causal inference from observational data is the goal of many data analyses in the health and social sciences. However, academic statistics has often frowned upon data analyses with a causal objective. The introduction of the term “data science” provides a historical opportunity to redefine data analysis in such a way that it naturally accommodates causal inference from observational data.

Like others before, we organize the scientific contributions of data science into three classes of tasks: description, prediction, and causal inference. An explicit classification of data science tasks is necessary to discuss the data, assumptions, and analytics required to successfully accomplish each task.

We argue that a failure to adequately describe the role of subject-matter expert knowledge in data analysis is a source of widespread misunderstandings about data science. Specifically, causal analyses typically require not only good data and algorithms, but also domain expert knowledge. We discuss the implications for the use of data science to guide decision-making in the real world and to train data scientists.

データ解析における現場知識の重要性  
に対する認識の欠如は、  
「データサイエンス」における幅広い誤解の源泉である。



## Definition (Mode 2 Science)

The essential ingredients of production of knowledge by Mode 2 Science are

- ① context of application;
- ② transdisciplinary;
- ③ heterogeneity/diversity;
- ④ reflexivity/social accountability;
- ⑤ novel quality control source.



|                                 | <b>Mode 1</b>                                | <b>Mode 2</b>  |
|---------------------------------|--|--|
| Problem                         | Academic context                             | Context of application                                       |
| Relation with other disciplines | Disciplinary                                 | transdisciplinary  |
| Institute                       | Homogeneity                                  | Heterogeneity  |
| Goal                            | Autonomy                                     | Reflexivity/social accountability<br>(social accountability) |
| Criterion                       | Traditional quality control<br>(peer review) | Novel quality control  |

## 1 Data Science as Mode 2 Science

- Science Revisited
- From Data Analysis to Data Science
- Teaching Data Science
  - Teaching Mode 2 Data Science at YCU
  - Data Science Club: Learning by Practicing
  - Data Science PBL for Juniors
- Practical Data Science in Master's Program

## 2 YOKOHAMA D-STEP



masters of the field. “Data analysts”, even if professional statisticians, will have had far less exposure to professional data analysts during their training.

Three reasons for this hold today and can at best be altered slowly:

(c1) Statistics tends to be taught as part of mathematics.

(c2) In learning statistics *per se* there has been limited attention to data analysis.

(c3) The number of years of intimate and vigorous contact with professionals is far less for statistics Ph.D.’s than for physics (or mathematics) Ph.D.’s.

# How to Teach DA: Tukey (1962)

John W Tukey (1962). The future of data analysis. *AMS*, 1-67.



Data analysis, and the parts of statistics which adhere to it, must then take on the characteristics of a science rather than those of mathematics, specifically:

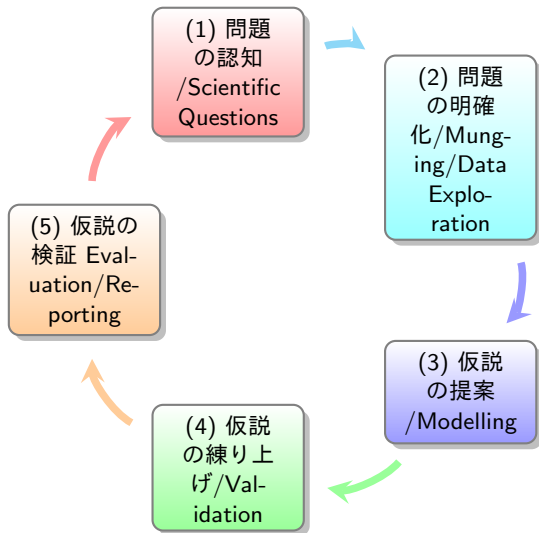
(b1) Data analysis must seek for scope and usefulness rather than security.

(b2) Data analysis must be willing to err moderately often in order that inadequate evidence shall more often *suggest* the right answer.

(b3) Data analysis must use mathematical argument and mathematical results as bases for judgment rather than as bases for proof or stamps of validity.



# Teach DS Based on John Dewey's Structured Reflective Thinking



## 1 Data Science as Mode 2 Science

- Science Revisited
- From Data Analysis to Data Science
- Teaching Data Science
  - Teaching Mode 2 Data Science at YCU
  - Data Science Club: Learning by Practicing
  - Data Science PBL for Juniors
- Practical Data Science in Master's Program

## 2 YOKOHAMA D-STEP

# Data Science Club: Learning by Practicing

## YCUデータサイエンス倶楽部

### 目次



1. YCUデータサイエンス倶楽部とは
2. 活動内容
  - 2.1. セミナー記録
3. 準備編
  - 3.1. 統計学入門
  - 3.2. 統計的推測（入門編）

## YCUデータサイエンス倶楽部とは

統計学・機械学習を初めとするデータサイエンスの基礎から、最先端までの研究やデータサイエンスがもたらす価値創造まで探求する研究会です。メンバーは横浜市立大学データサイエンス学部 of 学生から構成されます。他学部の学生はオブザーバーとして参加したい場合、必ず指導教員を通して倶楽部顧問までご連絡ください。

倶楽部長：十鳥大地（DS学部一年）

倶楽部顧問：汪金芳（DS学部教授）

## 活動内容

2018年8月から次の統計学・機械学習の分野における名著を輪読する予定です。

統計的学習の基礎 —データマイニング・推論・予測—

Trevor Hastie, Robert Tibshirani, Jerome Friedman

## サッカーの勝利要因 と予測について

横浜市立大学データサイエンス学部1年

橘川 和司 千田 晟也 羽田野 佑奈  
宇井 裕翔 長島 正悟 浦 優太 汪 金芳

第8回スポーツデータ解析コンペティション：分析部門（サッカー）

1

# Can't Open File!



ファイル全体を読み込むことができませんでした。

OK

CSV 開いています... 2017102111\_仙台vs清水\_1stHalf.csv

キャンセルするには command キーを押しながらピリオド (.) キーを押します。

リモートディスク

BOOTCAMP

2017102111\_仙台vs清水\_1stHalf.csv

2017102112\_大宮vs柏\_1stHalf.csv

2017102112\_大宮vs柏\_2ndHalf.csv

2017102113\_C大阪vs甲府\_1stHalf.csv

2017102113\_C大阪vs甲府\_2ndHalf.csv

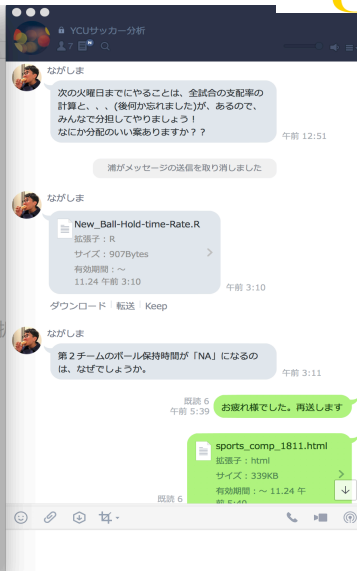
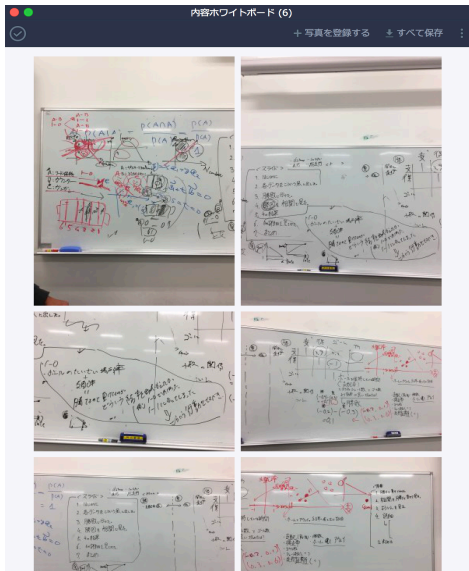
2017102114\_神戸vs鳥栖\_1stHalf.csv

2017102114\_神戸vs鳥栖\_2ndHalf.csv

2017102902\_FC東京vs清水\_2ndHalf.csv

2017102903\_甲府vs神戸\_1stHalf.csv

# Most of Time Spent on “what to do”



# 5. 付録：Rのコード

```

1 #走行距離の自動計算
2 #csvファイル名のベクトル
3 library("data.table")
4 #csvのファイルがたくさん入っているディレクトリ内で実行する
5 #ディレクトリの変更をやっておくこと！
6 Files=list.files()
7 #ボールタッチの試合情報_TBをcsvにして読み込む
8 #これも、ディレクトリ注意！
9 gameInfo=fread("GameInformation_TB.csv",header = T, sep = ',',
10               encoding = "UTF-8")
11 colNames(gameInfo)=c("gameID", "date", "half", "homeID",
12                    "homeName", "homeNA", "awayID",
13                    "awayName", "awayNA", "stadiumID",
14                    "stadiumName", "homeScore", "homeResultID",
15                    "awayScore", "awayResultID",
16                    "vlewer", "weather1", "weather2",
17                    "weather3", "temp", "humid",
18                    "wind", "grandState", "grandSurface",
19                    "startTime")
20 #チームIDだけ抜き出す
21 teamID=unique(subset(gameInfo,select = "homeID"))
22 colNames(teamID)="id"
23 #試合idとチームidの照らし合わせをする
24 idSet=subset(gameInfo,select = c("gameID", "homeID", "awayID"))
25
26 funDistanceNext=function(){
27   #前後半に分かれているので、
28   #全ファイル数の半分のずつに計測する
29   for(i in 1:(length(Files)/2)){
30     #i番からゲームIDを抜き出す
31     gameID=idSet$gameID[i]
32     #前半と後半のデータをそれぞれ読み込む
33     firstHalf=fread(Files[2*i],header = T, sep = ',',
34                   encoding = "UTF-8")
35     secondHalf=fread(Files[2*i+1],header = T, sep = ',',
36                     encoding = "UTF-8")
37     #変数名を日本語から英語にする
38     colNames(firstHalf) <- c("id", "frame", "home", "system",
39                             "number", "X", "Y", "speed")
40     colNames(secondHalf) <- c("id", "frame", "home", "system",
41                               "number", "X", "Y", "speed")
42
43     #ホームの出場選手一覧
44     homePlayer=cbind(unique(subset(firstHalf,home==1,select = c("number"))),
45                     unique(subset(secondHalf,home==1,select = c("number"))))
46     colNames(homePlayer)=c("firstNumber", "secondNumber")
47     #前半で出場する選手の人数が違うので調整する
48     if(length(unique(subset(firstHalf,home==1,select = c("number"))$number))!=
49         length(unique(subset(secondHalf,home==1,select = c("number"))$number)))){

```

```

41         "number", "X", "Y", "speed")
42
43     #ホームの出場選手一覧
44     homePlayer=cbind(unique(subset(firstHalf,home==1,select = c("number"))),
45                     unique(subset(secondHalf,home==1,select = c("number"))))
46     colNames(homePlayer)=c("firstNumber", "secondNumber")
47     #前後半で出場する選手の人数が違うので調整する
48     if(length(unique(subset(firstHalf,home==1,select = c("number"))$number))!=
49         length(unique(subset(secondHalf,home==1,select = c("number"))$number)))){
50       #NAとなる部分（人数の関係）ができた。
51       #アウェイの出場選手一覧
52       awayPlayer=cbind(unique(subset(firstHalf,home==2,select = c("number"))),
53                       unique(subset(secondHalf,home==2,select = c("number"))))
54       colNames(awayPlayer)=c("firstNumber", "secondNumber")
55     #ホームの選手と同様に調整を行う
56     if(length(unique(subset(firstHalf,home==2,select = c("number"))$number))!=
57         length(unique(subset(secondHalf,home==2,select = c("number"))$number)))){
58       #試合の途中結果と最終結果は
59       Judge=subset(GameResult,gameID==gameID)$Judge
60     #前半終了時
61
62     #選手の走行距離を入れるためのデータフレームを作る
63     homeRunDistance=data.frame("firstNumber"=homePlayer$firstNumber,
64                               "firstDistance"=rep(0, length(homePlayer$firstNumber)),
65                               "secondNumber"=homePlayer$secondNumber,
66                               "secondDistance"=rep(0, length(homePlayer$secondNumber)))
67     awayRunDistance=data.frame("firstNumber"=awayPlayer$firstNumber,
68                               "firstDistance"=rep(0, length(awayPlayer$firstNumber)),
69                               "secondNumber"=awayPlayer$secondNumber,
70                               "secondDistance"=rep(0, length(awayPlayer$secondNumber)))
71
72     #走行距離のデータフレームを完成させる
73     #ホーム前半
74     for(j in 1:length(homePlayer$firstNumber)){
75       #ホーム後半
76       for(j in 1:length(homePlayer$secondNumber)){
77         #アウェイ前半
78         for(j in 1:length(awayPlayer$firstNumber)){
79           #アウェイ後半
80           for(j in 1:length(awayPlayer$secondNumber)){
81             #完成
82
83             #節約用データフレームを作る
84             #チーム別の合計Summary（総走行距離、返り出場選手、試合数、
85             #一試合における走行距離、一試合における出場選手、選手一人における走行距離
86             if(i==1){
87               #初回走行距離
88             }

```



# Data Science PBL for Juniors



## PBL実習(インターンシップ)実施計画テンプレート

- 実習の種類  
実績追体験型    解決プロセス現在進行型    課題未着手型    課題理解型    その他

- 実習テーマ

- 実習内容

【課題】

【内容詳細】

- 期間： 月 日（曜日）～ 月 日（曜日）（日） [最短1週間～最長3週間]

- スケジュール

|     | 第1週 |    | 第2週 |    | 第3週 |    |
|-----|-----|----|-----|----|-----|----|
|     | AM  | PM | AM  | PM | AM  | PM |
| 1日目 |     |    |     |    |     |    |
| 2日目 |     |    |     |    |     |    |
| 3日目 |     |    |     |    |     |    |
| 4日目 |     |    |     |    |     |    |
| 5日目 |     |    |     |    |     |    |

- 受け入れ人数：  人（グループワーク）

- 期待スキル

- 事前学習

2019.08.02

# Practical Data Science in Master's Program

## DS専攻 博士前期課程カリキュラム

■必修科目 □選択科目 ※ () 内は単位数

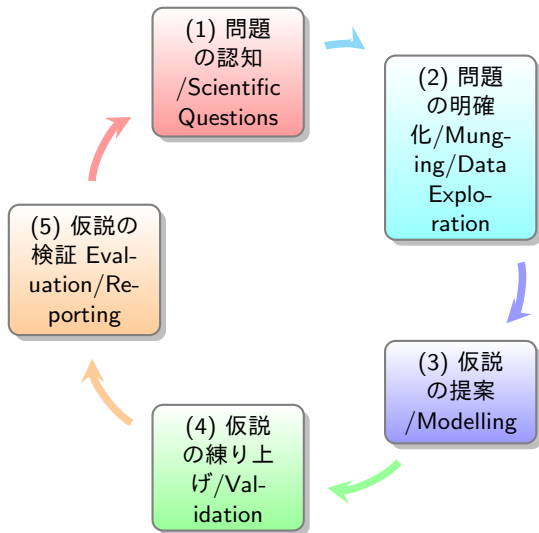
| M1前期  |                 | M1後期              |                 | M2前期        |  | M2後期 |  |
|---|-----------------|-------------------|-----------------|-------------|--|------|--|
| DS研究指導 I~IV (計8/1年・2年)・修士論文 (0)                         |                 |                   |                 |             |  |      |  |
| PDS I (2)   |                 | PDS II (2)        |                 | PDS III (2) |  |      |  |
| 統計学特論 (2)   |                 | 機械学習特論 (2)        |                 |             |  |      |  |
| デザイン思考特論 (1)  |                 | データマニング特論 (2)     |                 |             |  |      |  |
| 応用論理学 (1)   |                 |                   |                 |             |  |      |  |
| (統計科学中心型)   |                 |                   |                 |             |  |      |  |
| 多変量統計解析特論 (2)   | 最適化の基礎と応用特論 (2) | 時系列データ解析特論 (2)    | 都市環境データ解析特論 (2) |             |  |      |  |
| 実験計画と因果推論特論 (2)   | 標本調査特論 (2)      | 他専攻・他研究科科目 (2)    |                 |             |  |      |  |
| (計算機科学中心型)  |                 |                   |                 |             |  |      |  |
| クラウドコンピューティング特論 (2)                                     | 計算機統計学特論 (2)    | ビッグデータ処理基盤特論 (2)  | 非構造化データ特論 (2)   |             |  |      |  |
| プログラミング特論 (2)   |                 | 実験とシミュレーション特論 (2) | データ可視化特論 (2)    |             |  |      |  |
|   |                 | 自然言語処理特論 (2)      |                 |             |  |      |  |
| データアナリティクス特別講義・データエンジニアリング特別講義・データサインス展開特別講義(集中講義等) (2) |                 |                   |                 |             |  |      |  |

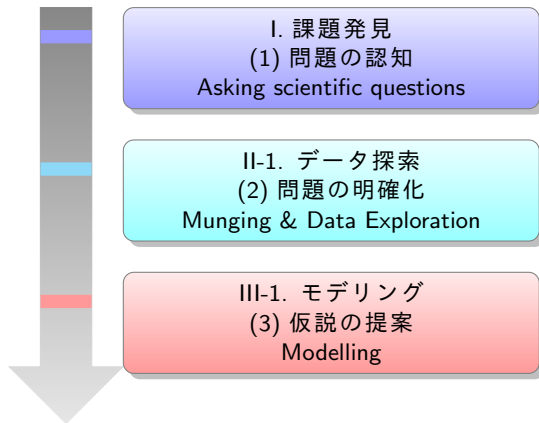
修了要件：30 単位

- 特別研究・特別演習：14単位 PDS(必修)6単位、ゼミ/修士論文(必修)8単位
- 共通科目：16単位 講義/演習(必修) 8単位、講義/演習(選択) 8単位以上

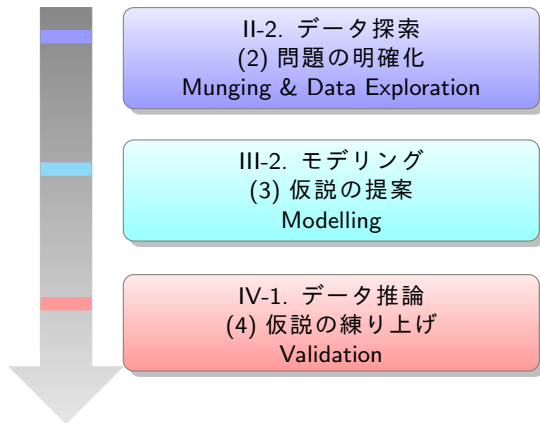
・学部4年生対象の早期履修科目は、特論科目を5科目10単位とする

# PDS Based on John Dewey's Structured Reflective Thinking



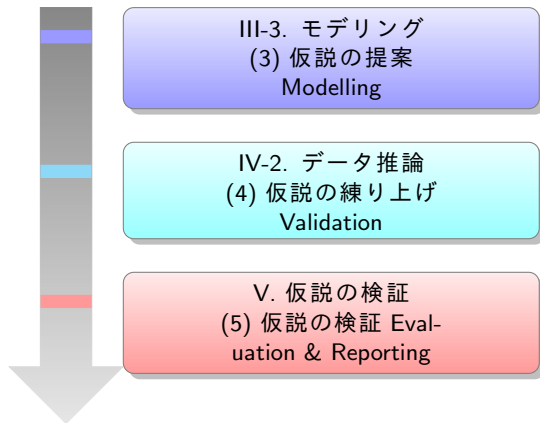


M1 前期 2 単位 (Semester 1, 2 credits)

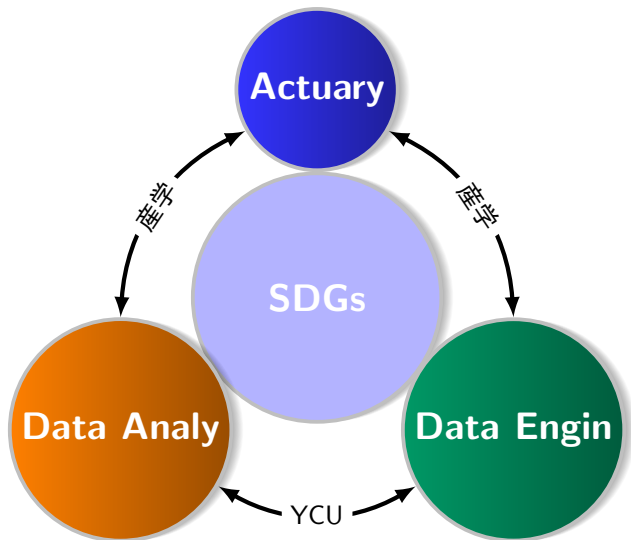


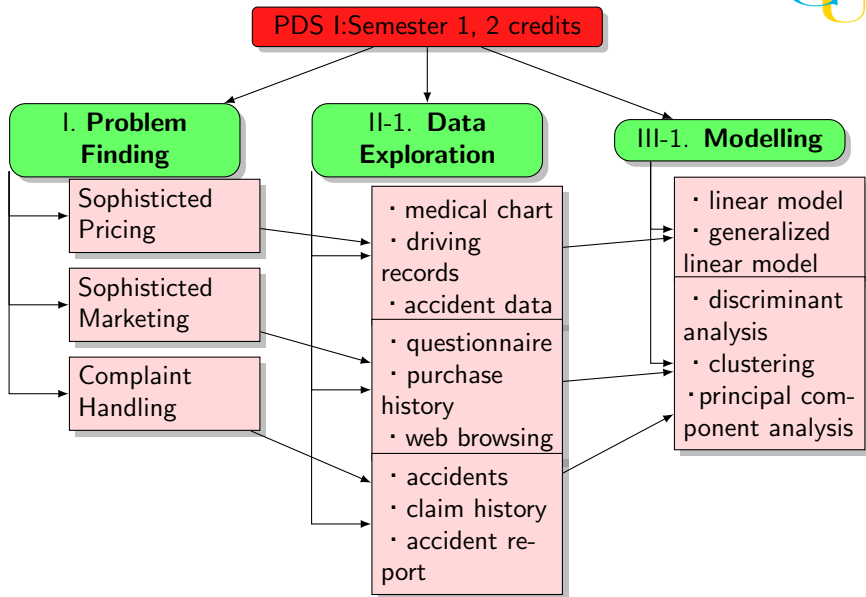
M1 後期 2 単位 (Semester 2, 2 credits)





M2 前期 2 単位 (Semester 3, 2 credits)





PDS II:Semester 2, 2 credits

## II-2. Data Exploration

descriptive statistics

correlation

visualization

application of linear models

## III-2. Modelling

- linear model
- generalized linear model
- Bayesian models

- discriminant analysis
- clustering
- principal component analysis
- Bayesian networks
- decision trees

## IV-1. Inference

least squares

maximum likelihood

bayesian estimation

model selection

LASSO

model averaging

PDS III:Semester 3, 2 credits

## III-3. Modelling

optimum prediction model

optimum bayesian model

optimum discriminant model

## IV-2. Inference

least squares

maximum likelihood

bayesian estimation

model selection

LASSO

model averaging

## V. Validation

theoretical validity

empirical validity

- Proposal of Sophisticated Pricing
- Proposal of Sophisticated Marketing
- Proposal of Complaint Handling



文理融合・実課題解決型データサイエンティスト育成

**YOKOHAMA D-STEP**

Data Scientist Educational Program

# D-STEPの取り組み

横浜市立大学 データサイエンス推進センター  
坂巻顕太郎

sakamaki@yokohama-cu.ac.jp

横浜市立大学  
YOKOHAMA CITY UNIVERSITY

## データサイエンティスト育成プログラムの募集開始！文部科学省「超スマート社会の実現に向けたデータサイエンティスト育成事業」始動

横浜市立大学は、東京理科大学、明治大学と3大学共同で「文理融合・実課題解決型データサイエンティスト育成※1」事業を実施します。これは、平成30年度文部科学省「超スマート社会の実現に向けたデータサイエンティスト育成事業※2」に採択された事業で、横浜市内のみならず我が国の社会的、経済的課題をデータサイエンスで解決できる人材育成のプラットフォームを構築し、5か年で約200人の高度データサイエンティスト、約800名のデータエキスパートを産学官連携のもとで実践的に育成することを目指すものです。特に、データサイエンティストが社会活動の中で直面する〔課題設定⇒データ収集・分析⇒新たな価値創造⇒実装〕という一連のプロセスへの対応を念頭に置いたPBL※3を行い、受講生の社会展開力を涵養することに力を置いています。

本事業では、2018年に横浜市立大学に開設されたデータサイエンス学部の知見を生かして、大学間、産業界、行政と連携して実社会から求められるデータサイエンティストを育成するために、3つの教育プログラムを開設し、本年2月19日より受講生の募集を開始します。

## プログラム

受講生の保有知識・スキルを基に、独り立ち可能なデータサイエンティストを要請。また、データサイエンティストを社会・企業に広げ、データサイエンティストの活躍の場を創るビジネスパーソン文化まで促成

# A

課題発見・解決型  
データサイエンティスト  
育成コース

データサイエンスに係る一連の  
流れを“習得したサイエンティスト”を輩出

- ・データサイエンス分野“非”専攻分野者が対象
- ・データサイエンスの基礎知識から発展までを習得
- ・通年のPBLで課題の発見・解決までの一連の流れを学習

# B

データ分析型 データ  
サイエンティスト育成  
コース

データサイエンス分野で自走できる“習熟したサイエンティスト”を輩出

- ・データサイエンス分野専攻者が対象
- ・データサイエンスに関する高度な理論・方法を習得
- ・半期のPBLで実課題解決から価値創出までを経験

# C

社会人（データエキスパート）育成コース

文化を創りサイエンティストを使いこなせる“エキスパート”を輩出

- ・短時間でデータサイエンスに関する基礎知識や活用手法を習得
- ・現場で活かせる軽微な分析業務を履成
- ・データサイエンティストの活躍の場、活用方法を理解

なお、本プログラムは、国が定める大学等における履修証明制度に該当する「履修証明プログラム」となります。「履修証明プログラム」は、学校教育法第 105 条及び学校教育法施工規則第 164 条の規定に基づき、大学が特別に社会貢献等を趣旨として、主として社会人向けに体系的な教育プログラムを開発し、その修了者（120 時間以上の履修を完了する者）に対し、学長名の履修証明を交付するものです。



# 講義概要 (A, Bコース)

Aコース: 課題発見・解決型データサイエンティスト育成コース  
 Bコース: データ分析型データサイエンティスト育成コース

| 授業科目          | A | B | 開講時期       | 単位 | 授業科目          | A | B | 開講時期 | 単位 |
|---------------|---|---|------------|----|---------------|---|---|------|----|
| 統計学基礎         | ◎ |   | e-learning | 1  | 統計・機械学習モデリング  | ◎ | ◎ | 後期   | 1  |
| 応用線形代数        | ◎ |   | e-learning | 1  | データマニング       | ○ | ○ | 後期   | 2  |
| データ分析基礎       | ◎ |   | e-learning | 1  | 最適化と計算機科学     | ○ | ○ | 後期   | 2  |
| ITセキュリティ・情報倫理 | ◎ | ◎ | e-learning | 1  | 非構造化データ特論     | ○ | ○ | 後期   | 1  |
| 多変量データ解析      | ○ | ○ | 前期         | 2  | 実験とシミュレーション特論 | ○ | ○ | 後期   | 1  |
| 標本調査法         | ○ | ○ | 前期         | 1  | プログラミング特論     | ○ | ○ | 後期   | 1  |
| 欠測データ解析       | ○ | ○ | 前期         | 2  | 計算機統計学特論      | ○ | ○ | 前期集中 | 1  |
| ノンパラメトリック法    | ○ | ○ | 前期         | 2  | 時系列データ解析特論    | ○ | ○ | 前期集中 | 1  |
| 最適化理論         | ○ | ○ | 前期         | 2  | 都市環境データ解析特論   | ○ | ○ | 前期集中 | 1  |
| 数理ファイナンス*     | ○ | ○ | 前期         | 2  | データ可視化特論      | ○ | ○ | 前期集中 | 1  |
| デザイン思考        | ◎ | ◎ | 前期集中       | 1  | 数理医学          | ○ | ○ | 集中   | 2  |

◎: 必修科目, ○: 選択科目, △: 選択必修科目

| 授業科目        | A | B | 開講時期 | 単位 |
|-------------|---|---|------|----|
| 政策課題解決PBL   | ◎ | △ | 前期   | 2  |
| ビジネス課題解決PBL | ◎ | △ | 後期   | 2  |

<https://www.yokohama-cu.ac.jp/academics/ds/syllabus.pdf>

# 政策課題解決PBL

## 講義概要

本講義では、個票データ（擬似データ）を用いて、グループワークを通じた実践的な形で問題解決の一連のプロセスを学習する。ビジネス経験、データ分析経験が豊富な講師陣が、各グループのアウトプットに対するフィードバックを提供する。問題解決における他者との協働を身に着けるだけでなく、模擬インタビューやプレゼンテーションを実施し、問題解決に必要な能力の素地を総合的に養う。

## 授業計画

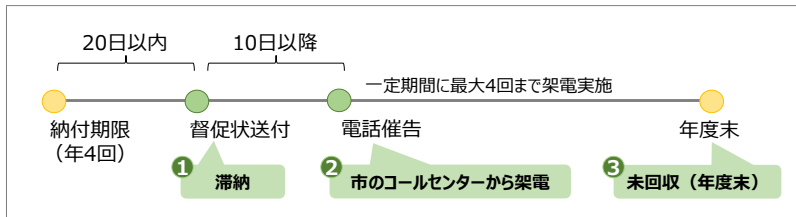
- 1 講義①（データ分析における心構え）
- 2 講義②（ロジカルシンキング）
- 3 課題説明
- 4 講義③（データ分析）
- 5 講義④（インタビューの心得・インタビューによる情報収集）
- 6 グループワーク（課題設定・データ分析）
- 7 グループワーク（データ分析・発表資料作成）
- 8 中間発表
- 9 講義⑤（解決策提言に向けたストーリー構築）
- 10 グループワーク（中間発表でのレビューを踏まえた再検討）
- 11 グループワーク（データ分析・発表資料作成）
- 12 グループワーク（データ分析・発表資料作成）
- 13 最終発表リハーサル
- 14 最終発表に向けた資料のアップデート
- 15 最終発表

# プロジェクトゴール

## ・ ゴール：今年の滞納未回収額を10%削減

- ・A市の納税課では、「1年で滞納未回収額を10%減」という目標値を設定しています。
- ・今年度中の回収額増加を狙うため、6月末までに施策をつくる必要があります。
- ・大きな社会環境変化は見込まれないため、何も施策を行わなかった場合、今年度も約50億円の未回収が発生するものと想定しています。

# 督促業務



- ① 滞納とは、納税者が納付期限までに納税しないために督促状を送達したにもかかわらず、完納されないことをいう。督促状が送達されていることが滞納処分の原則的な前提条件となる。
- ② 督促状送付から一定期間後に、市のコールセンターから電話による納税の催告を行う。A市では、電話が繋がらなかった場合、最大4回まで架電するルールとする。コールセンターでは市民税以外にも担当しており、現状の架電回数を増加することは難しい。
- ③ 未回収は、電話による督促を実施しても年度末までに納付されなかった分を指す。

# ビジネス課題解決PBL

- 現実の課題に対してデータから具体的な施策を提案
  - 2019年度はベ이스ターズ、資生堂、ハレックスとの協同
  
- 講義内容
  - 企業と課題の共有
    - 課題とデータの説明（初回、中間レビュー）
  - 課題解決のためのデザイン
    - アウトカム / 解析単位の設定と説明変数の考案
  - データ加工
    - SQLによるデータ抽出や説明変数の加工など
  - 基本的な分析手法の適用
    - Rによる分析（glm, step, glmnetなど）
  - 機械学習手法による高精度な予測
    - 各種パッケージとクロスバリデーション
  - 施策の立案とプレゼンテーション
    - 分析結果から施策を立案

# 課題とデータの説明

資生堂ジャパン



横浜DeNAベイスターズ

ハレックス



## • 中間レビュー

- 受講生によるビジネス課題解決策のプレゼン
- 協力企業から提案事項に対する有益なコメントを得て、提案をさらにブラッシュアップ

# 自治体向け データエキスパート 育成コース開講



田東 正典



山竹 竹彦



坂巻 龍太郎



西内 啓

## EBPMのためのデータ利活用方法

2016年に施行された官民データ活用推進基本法では、多様な大量の情報を活用することにより超少子高齢社会における諸課題を解決するため、地方公共団体を含む様々な主体が保有するデータの適正かつ効果的な利用を推進しています。

基本理念は、データ活用により、自立的で個性豊かな地域社会の形成などを図り、根拠(エビデンス)に基づく政策立案や評価を行うことです。この理念に基づいて官民データを十分に活用するには、データ活用の知識や技術を習得することが重要になります。

## 自治体職員のための 首都圏初回のデータサイエンス学部が 短期集中講座を開講

本セミナーでは、地方自治体の職員を対象に、  
1) エビデンスを理解するために必要な知識とデータ分析、2) すでに蓄積されているエビデンスを探索するために必要な文献検索の方法、3) 課題解決のためのデータ分析、の3つについて具体例に基づいた講義と演習を行います。これにより、エビデンスに基づいた政策立案や評価の効果的かつ効率的な推進に必要であるデータ活用に関する基本的な素養の獲得が期待されます。

2020年 1月29日(水)・30日(木)

会場:神奈川産業振興センター 第2会議室

(※この研修は、首都圏および全国で初の開催となります。研修の申し込みは、研修申込書(研修申込書)と研修参加申込書(研修参加申込書)を、郵送(1月17日(金)まで)。

主催:公立大学法人横浜国立大学

- 定 員:50名
- 対 象:自治体職員
- 受講料:無料
- 持ち物:ノートPC
- 講師:田東 正典 (横浜国立大学) 坂巻 龍太郎 (横浜国立大学) 西内 啓 (横浜国立大学株式会社データビークル) など
- 事業責任者:山竹 竹彦 (横浜国立大学)



詳細はこちらから <https://d-step.yokohama>

## Part. 1 2020年1月29日(水)9:30-12:00 エビデンスとデータ活用

どのような人が健康に関する課題があるかをデータから探索すること、スポーツジムなどに適いやすくなる政策が健康づくりに寄与するかをデータから考えることは異なりませう。この違いを理解するには、エビデンスとは何か、データ分析で何ができるか、などを知ることが重要です。本パートでは、データの読み解きに関する講義からエビデンスやデータ活用の理解を促進することを目的とします。

## Part. 2 2020年1月29日(水)13:00-17:00 既存のエビデンスを探索するための文献検索

課題や解決策に関する知見が既に存在するかどうかを事前に文献などから検討しなければ、それらの重要さやデータに基づく検討が必要かどうかの判断はできません。課題に対する知見(仮説)なしにビッグデータを解析することで、Google Flu Trends がインフルエンガの流行予測に失敗したことが顕著な例です。特定分野(業界や領域など)に関する知見の積み重ねはデータサイエンスを実践するうえでは重要になります。本パートでは、文献検索に関する講義と演習により、政策立案や評価をエビデンスに即して行うために必要となる事前の知見の収集の理解を目的とします。

## Part. 3 2020年1月30日(木)9:30-17:30 データに基づく課題の検討

データ分析の目的に応じて、用いる分析手法や分析結果の解釈は異なります。例えば、課題解決策をデータから評価する際、解決策を策定した集団と実施していない集団のデータを単純に比較するのは、年齢や性別などの特徴を考慮して評価するのは、どのようなデータを評価に用いるかによって異なります。本パートでは、データ分析の基礎的な演習を行い、データハンドリングやデータ分析に対する理解を深めることで、実際のデータ分析とエビデンスの関係を理解することを目的とします。

<注意事項>

ノート PC は各自が持参してください。Excel (または csv ファイルを開読できるソフトウェア)、RStudio がインストールされていることが必要となります。RStudio がインストールされていない場合は、事前に必要なソフトウェアのインストールについて別途ご案内いたします。RStudio とインターネットのインストールには管理者権限が必要となるため、個人が所有する PC へ、管理者アカウントで必要なソフトウェアをインストールする PC のインストールした PC をご使用ください。管理者アカウントが利用できない場合は各自でご確認ください。



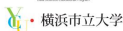
### お問い合わせ先

公立大学法人 横浜国立大学 D-STEP 事務局  
〒236-0027 神奈川県横浜市中区新港1-2-2  
☎ 045-787-8906  
E-mail: d\_step@yokohama-cu.ac.jp  
<https://d-step.yokohama>

### Yokohama City University

### Yokohama D-STEP

Yokohama City University



# D-STEPの取り組み

- A, Bコース
  - データサイエンティストを育成するための通年コース
  - 統計や機械学習, デザインに関する講義
    - on-siteだけでなくonlineでも提供
  - Project Based Learningによるグループ演習
    - 実課題解決にデータサイエンスをどう生かすかを学ぶ機会
- Cコース
  - データサイエンスに関わる非専門家が必要な基礎的な知識を提供
    - 専門家と非専門家のコミュニケーションを円滑にする
  - データサイエンティストへの入門
    - A, Bコースへの接続



# D-STEPの取り組み

- A, Bコース
  - データサイエンティストを育成するための通年コース
  - 統計や機械学習, デザインに関する講義
    - on-siteだけでなくonlineでも提供
  - Project Based Learningによるグループ演習
    - 実課題解決にデータサイエンスをどう生かすかを学ぶ機会
- Cコース
  - データサイエンスに関わる非専門家が必要な基礎的な知識を提供
    - 専門家と非専門家のコミュニケーションを円滑にする
  - データサイエンティストへの入門
    - A, Bコースへの接続