

情報・システム研究機構 シンポジウム講演集
『データサイエンスの推進に向けて』

分野を超えた
データサイエンスの広がり
～自然科学から人文社会科学まで～

2017年 2月20日(月)
東京大学 伊藤謝恩ホール



目次

CONTENTS

開催趣旨 ごあいさつ	
情報・システム研究機構／機構長 北川 源四郎	1
【機構長講演】	
「情報・システム研究機構のこれまでを振り返って」	
情報・システム研究機構／機構長 北川 源四郎	2
【講演】	
「データサイエンス共同利用基盤施設の取組み」	
データサイエンス共同利用基盤施設／施設長 藤山 秋佐夫	4
【講演】	
「データサイエンス共同利用基盤施設における具体的取組みの紹介」	
・ライフサイエンス統合データベースセンター／センター長 小原 雄治	5
・社会データ構造化センター／センター長 吉野 諒三	6
・人文学オープンデータ共同利用センター／準備室長 北本 朝展	7
・ゲノムデータ解析支援センター／センター長 野口 英樹	8
・データ融合計算支援プロジェクト／中野 慎也	9
・極域環境データサイエンスセンター／準備室長 門倉 昭	10
【中継】	
南極昭和基地からの中継	
「データ発生の現場から」	11
【講演】	
「異分野融合・新分野創成を担うデータサイエンティストの育成基盤」	
情報・システム研究機構／理事・統計数理研究所／所長 樋口 知之	12
【次期機構長講演】	
「情報・システム研究機構の新時代に向けて」	
情報・システム研究機構／理事／次期機構長 藤井 良一	14
【招待講演】	
「大学におけるデータサイエンスとその教育」	
九州大学／理事・副学長 安浦 寛人 氏	16
【招待講演】	
「データは誰のものか」	
人間文化研究機構／理事 佐藤 洋一郎 氏	16
【招待講演】	
「シン・ニホン-AI×データ時代における日本の現状と人材育成課題-」	
ヤフー株式会社／チーフストラテジーオフィサー(CSO) 安宅 和人 氏	16

開催趣旨 ごあいさつ

「分野を超えたデータサイエンスの広がり」

～自然科学から人文社会科学まで～

当機構は「現代社会が直面する複雑な対象を情報とシステムの観点から捉える」という理念を掲げ、いち早くデータサイエンスの重要性を主張し実践してまいりました。

平成25年度には、大学共同利用機関の機能強化の一環として、データ中心科学リサーチコモンズ事業を開始し、機構を挙げてデータ中心科学の確立に向けた活動を進めてきました。

法人第三期開始年となる本年度からは、我が国の大学などにおけるデータ駆動型の学術研究をさらに積極的に推進するため、機構内にデータサイエンス共同利用基盤施設を設置いたしました。これにより大学などの多様な分野の研究者に対する、大規模データ共有およびデータ解析の支援と人材育成とを抜本的に強化し、科学の発展や社会のイノベーション推進に対する活動を大きく展開してまいります。



この取組みでは、これまで組織的交流が比較的少なかった異分野間や領域研究と方法論研究の間で、研究データに関連するニーズとシーズとのマッチングを実施する必要があります。そこで、本年のシンポジウムは「分野を超えたデータサイエンスの広がり～自然科学から人文社会科学まで～」をテーマとし、大学、大学共同利用機関、企業などからのご意見を伺うために、大学、他機構、および企業の著名な方々にご講演いただくとともに、新施設の具体的な活動をご紹介します。

また、当機構国立極地研究所の昭和基地（南極）における観測データ取得現場からの生中継も行います。

今回のシンポジウムが、産官学を交えて、今後のデータサイエンスの目指すべき方向性を、ともに考えていく機会となれば幸いです。

大学共同利用機関法人 情報・システム研究機構
機構長 北川 源一郎



「情報・システム研究機構のこれまでを振り返って」

大学共同利用機関法人 情報・システム研究機構／機構長 北川 源四郎



情報とシステムという観点から問題を捉える

情報・システム研究機構は、2004年に国立大学共同利用機関が法人化された際に、国立極地研究所、国立情報学研究所、統計数理研究所、国立遺伝学研究所、この4つの研究所が結集してできた大学共同利用機関法人です。その後2005年に新領域融合研究センター、2007年にライフサイエンス統合データベースセンター、2016年に新しく「データサイエンス共同利用基盤施設」を設置して、現在に至っております。

当機構は、2004年の設立以来、生命、地球、環境、社会などの現実社会における複雑な対象を、従来の「物質とエネルギー」という観点に替わって、「情報とシステム」という新しい立場から捉えるということを理念として掲げ、研究方法の開発、研究基盤の構築、およびそれらを使ったサイエンスの推進を図ってまいりました。具体的には、実験・観測・調査による大量データの取得と統合データベース化、そのデータベースに基づく知識獲得を行っており、たとえば現象の解明、予測・制御、あるいは意思決定という問題への応用を行ってまいりました。

ビッグデータがもたらす新しい歴史的瞬間へ

情報科学技術、特にセンサーの技術は飛躍的に発展し、現在アカデミックな領域でも、生命科学、金融、経済、環境、地球、物質・材料、人文学などあらゆる分野で大量・大規模データが獲得できるようになりました。

過去50年以上にわたって、ムーアの法則といわれる経験則に沿って、記憶容量、処理速度といった計算機のハードの能力はログスケールでほぼ直線的に増加してきました。しかしながら、近年のデータ生産量の伸びはログスケールでも爆発的に増大しているという圧倒的なスピードであり、これの意味するところは、従来のように計算機のハードへの依存だけでは問題が解決できない状況になっているということです。

同じような状況は人間社会の場面でも起こっています。たとえばwebやSNS等のインターネット情報、家電や自動車のセンサーからの情報、ドローンや人工知能、人工衛星からのリモートセンシングデータ、コンビニや株の取引データ、防犯カメラの画像や音声など人間社会の活動を精細かつ網羅的に捉えることによって、いわゆるビッグデータが出現しております。またこれらを活用することによって、顧客の行動を予測するマーケティング、オンラインショッピング、個別化医療、社会インフラのスマート化、自動運転などが発達しております。

人類の歴史における科学の役割を考えてみますと、専門家あるいは匠と呼ばれる人たちがこれまで経験と勘で行っていたことを、徐々にデータに基づき科学的に遂行できるようにしたということ

であろうかと思えます。たとえば占星術から天文学が発達し、錬金術からケミストリが出てきましたし、製鉄の匠の技が科学的な工業生産に変わったという事実もあります。近年では天気予報、経済予測、マネジメント、マーケティング、リスク管理なども科学的に取り扱われるようになり、いまや科学的発見（発見科学）自体やサービス、政策決定などもデータに基づく科学的方法の対象になりつつあります。

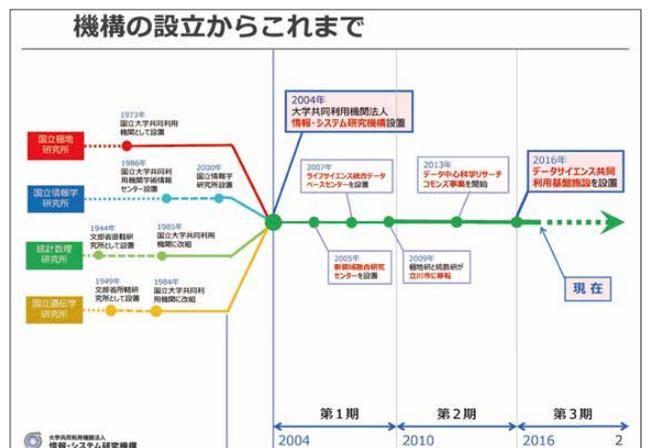
Ian Ayres (イアン・エアース) はその著書の中で、大量データ分析がもはやいろいろな分野で専門家の経験と勘を超えているという例を上げています。ビッグデータのインパクトは非常に大きく、やがて超スマート社会が実現し、数百年の歴史がある科学研究のスタイルもやがてデータサイエンス化していくと言われております。

2013年にコンピュータ将棋がプロ棋士に勝ったという出来事がありました。ついに碁の世界においても、昨年アルファ碁が世界のトップ棋士に勝ち、現在60連勝しているということがさらに衝撃的な話題となりました。

これに関連して思い浮かべるのが、1830年アメリカのボルチモア・オハイオ鉄道で行われた鉄道馬車と蒸気機関車の競争です。これはまさに産業革命を経て工業化社会が成立したことを象徴する歴史的瞬間であったわけです。いま我々は、ビッグデータ時代という次の歴史的瞬間に直面しているのではないかと思っております。

第4の科学であるデータサイエンスの必要性

ビッグデータを集めれば、今までできなかったような画期的なことができるという期待が膨らんでおりますけれども、それは楽観的過ぎるのではないかと考えます。確かにビッグデータには膨大な知識や価値が埋もれております。しかしその多くは構造化されておらず、価値密度が非常に低く、不均一であったり、スパースであったりしております。従って従来の方法をそのまま適用するだけで自動的に有用な価値が創出できるわけではありません。ビッグデータを効



果的・効率的に集約し、知識発見や価値創造を行うための革新的な方法が不可欠であり、大規模データ活用のための第4の科学的方法論の確立が必要であると考えております。

20世紀までの科学技術は、理論科学と実験科学の2つの方法をクルマの両輪として発展してきたと言えます。理論科学は演繹的方法であり、実験科学は帰納的方法あるいはデータ駆動型の方法であったわけです。20世紀の終盤、非線形系あるいは非常に複雑なシステムの理解とシミュレーションのために、第3の科学と呼ばれる計算科学が発達しました。それは京コンピュータに代表されるように著しい成果を上げております。しかしながら、ビッグデータの登場により、第4の科学としてのデータサイエンスの確立が必要です。この4つを揃えることによって、いわばクルマの4輪のように、今後の科学をドライブしていくことができると考えております。

データ中心科学リサーチコモンズ事業の役割

当機構では平成25年度、データサイエンス推進のための基盤を構築する「データ中心科学リサーチコモンズ」事業を開始いたしました。これは大きく分けて、ビッグデータの活用基盤を構築する部分と、それを使って実際のサイエンスを行う新領域融合研究センターの2つから成っております。

ビッグデータ活用基盤の形成は、データ基盤構築、モデリング・解析基盤構築、人材育成、この三位一体の基盤形成事業です。当機構の事業の特徴はこれらの分野を横糸で通す形で包括的に行っているということです。ただしデータ基盤の構築に関しては共通の方法論だけではなく、各分野固有の知識も不可欠ですから、当面は当機構が強みを持つ生命科学、地球環境科学、人間・社会科学、埋蔵分子の4つを中心に進めてまいりました。

データ基盤整備の項目で、地球環境データに関しましては、PANSYと呼ばれる南極に設置した大型レーダーのデータ取得、リアルタイム転送、アーカイブ化、さらにはそのデータの高度解析方法の開発などを行ってきました。ライフサイエンス統合データベースセンターでは、JSTとの連携によって、既に全国のライフサイエンスデータの統合化を成し遂げていますが、さらに高度化した利用を実現するために、RDF化、オントロジー整備、国際標準化に取り組んできました。人間・社会科学データに関しては、公的統計データの二次利用を可能にするオンサイト拠点の拡充、統計とweb情報の統合による高精度で高頻度のデータ作成、モバイル統計の取得などを行ってきました。データ中心ケミストリは、理論的に存在しうる化学物質を量子化学計算に基づき分子と化学反応のプロセスをデータベース化するというもので、埋蔵分子と呼ばれる新規の炭素物質の発見も行っております。

モデリング・解析基盤整備は、データから有用な知識を獲得するための方法とツールを開発する基盤整備です。データ同化・シミュレーションは、シミュレーションと観測データを統合するデータ同化手法の高度化を図るとともに、未適用分野の開発に取り組んできました。地震によって発生した微気圧変動分析からの津波の予測、細胞質流動を引き起こすせん断力分布の推定、感染症の制御をするためのシミュレーションなどを行ってきました。e-サイエンスデータ基盤技術はサイエンス3.0基盤として、リサーチマップの開発、実文書の解析手法の開発とツール公開などを行っております。その他イ

メージデータ解析、メタ知識構造解析なども行っております。

データサイエンティストの育成については、データサイエンスと融合研究の推進に必要な人材増とその育成の方法などを検討するための産官学連携の懇談会を開催し、報告書の中で提案しております。単に技術的なものではなく領域の知識、企画立案能力、コミュニケーション能力を持ったいわゆるT型、Π型の人材が必要であるということで、いくつかのレベルを想定し、それぞれの育成方法なども提案しているところでございます。

新領域融合研究センターからパラダイム創世へ

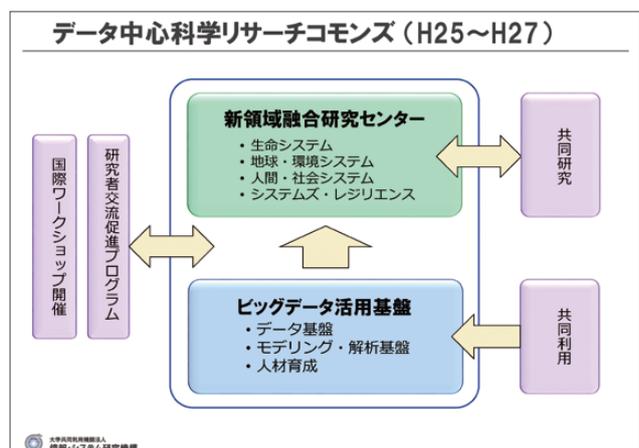
新領域融合研究センターは、ドメイン型の研究所である国立極地研究所、国立遺伝学研究所の「縦型の研究所」と、国立情報学研究所、統計数理研究所という情報基盤でありデータ解析を図る「横型の研究所」、この縦横クロスした所で、新しいパラダイム創成を目指すためのセンターです。平成25年度からは、遺伝機能(生命)システム融合研究、地球・環境システム融合研究、社会コミュニケーション融合研究、システムズ・レジリエンス融合研究の4つのプロジェクトを推進してまいりました。

地球・環境システム融合研究に関しては、南極で採取した3000mの氷柱コアで72万年前までの実データが取得できますが、そこに含まれる微量なゲノム解析を実現するための技術開発や、「コケ坊主」という特殊な生態系の全ゲノム解析なども行ってまいりました。

遺伝機能(生命)システム融合研究では、次世代シーケンサーのデータ生産とその解析方法の開発、ゲノム情報と表現型多様性データの統計手法開発、大量データに基づくゲノム機能と遺伝的ネットワークの抽出方法の開発などを行ってまいりました。一つの例としては、地球上のイネの系統をすべて解析しております。

その他、社会コミュニケーション融合研究では、自殺予防のための時空間統計データ分析などを、またシステムズ・レジリエンス融合研究では、減災を情報とシステムの観点から捉える取組みを進めてまいりました。

そして、いよいよこれらの成果を利用して、法人第三期がまさに始まったところでございます。



「データサイエンス共同利用基盤施設の取組み」

データサイエンス共同利用基盤施設／施設長
藤山 秋佐夫



データサイエンス共同利用基盤施設設置の目的

情報・システム研究機構には、国立情報学研究所があり、国立遺伝学研究所があり、国立極地研究所があり、統計数理研究所があって、世の中の通常の常識で考えると、一緒になっても何も起こらないと思われる研究所の集まりなんです。それらを集めて混ぜて、そこからさらに新しいものを作ろうとする努力が法人第一期、第二期に渡って行われ、いろいろと新しい成果も出てきたと思います。

現在の第三期では、データサイエンス推進という形で次のステップを目指している状況でございます。

じゃあ、こういう活動をして結局何ができたかということ、大学共同利用機関としての成熟度がかなり増大したということです。もともと普通においておくとなかなか混ざりにくい4つの研究所を強引に混ぜて、いろいろなことをやった結果として、共同研究体制がかなりできあがってきたということです。

特に第二期の後半あたりから、実際にわたくしが関係しているゲノム科学の世界ですと、データの生産量というのは飛躍的に増大しましたので、データ駆動型サイエンス、データサイエンスというのを生物学の世界でも実現できました。もちろん地球、物理、大気データの観測などでは、それ以前からビッグデータ時代に入れた科学が進んできたわけでございます。

大量・大規模データ共有のための仕組みを作る

データサイエンスを推進するためには、大規模なデータを生産するということが必要です。一言に大規模データといっても、ただデカければいいというわけではなく、データのどこかを見れば必ず全部がわかるという網羅性が大事です。そういう網羅的なデータは各研究所で作っていただくとして、この施設では、そのデータをどうやって大学の先生方に還元するかということを考えていきます。つまりデータ共有を進めるための仕組みを作りましょうということをまず掲げています。

それからもう一つ大事なことは、大規模データは通常のパソコンのレベルは超えてしまいますので、どうやって解析しようかということになります。そういう解析を支援するというユニットで、一つはこの機構の強みを生かしてゲノムデータのデータ解析支援ユニット、それからデータ融合計算支援ユニットの、二つのユニットを作っております。

データサイエンス共同利用基盤施設は、データ駆動型サイエンス（データサイエンス）の観点から、大学などの多様な分野の研究者に対し、大規模データの共有およびデータ解析の支援事業と人材育成を行い、我が国の大学などの機能強化に貢献する。これが本施設の基本目標でございます。事業内容は、データ共有支援、データ解析支援、T型・II型人材育成となっております。

それからもちろん大学の先生方からいろいろなアイデアをいただきたいし、内部から上がってくるいろいろなアイデアも伸ばしたいと思しますので、戦略プログラムという形でいろいろなプログラムを走らせていこうと、さまざまな計画を進めているところでございます。

データサイエンス推進に向けて果たす役割

設立からちょうど一年たったところなんです、当施設が狙っていることはこういうことになるかと思っております。

一つは大量・大規模データを活用し、発見・予測・シミュレーションといったデータサイエンス推進のための研究基盤構築と研究支援を行うことです。

それからもう一つは、データということに落とし込みますと、かなりいろいろな横串を通せるようになっておいて、そこからかなり新しい芽も出てくるのではないかと考えております。既存分野を超えた知識移転や汎化から、異分野交流・新分野創成を推進したいと考えております。

ひとまずここは大学共同利用機関なので、アカデミック・ビッグデータの活用研究拠点として、データサイエンスの支援事業を進めていくということでございます。

我々がミッションとしてやらなければいけないことは、データサイエンスを自由自在に使えるような環境を大学などの先生方へ提供し、実際に解析の支援などを行っていくことでございます。

やがて行きつく先はオープンサイエンスですが、その意味では機構そのものに国立情報学研究所もございまして、オープンサイエンスを進めるための基盤を提供する仕組みはもうできあがっております。そういったところにまで持ち込んでいって、機構全体としてデータドリブンのサイエンスを進めていきたいと思っております。

大学におけるデータ駆動型学術研究力強化のための共同利用推進事業

データサイエンス共同利用基盤施設

●データサイエンス推進のための研究基盤構築と研究支援

大量・大規模データの活用
発見・予測・シミュレーション

●既存分野を超えた知識移転や汎化→異分野交流・新分野創成の推進

多分野連携による新分野創成の推進
学術研究の発展に大きなインパクト



アカデミック・ビッグデータを活用した学術研究の推進

講演／データサイエンス共同利用基盤施設における具体的取組みの紹介 「ライフサイエンス統合データベースセンターの取組み」

ライフサイエンス統合データベースセンター／センター長
小原 雄治



データベース戦略、中核センター確立に向けて

そもそもライフサイエンスのデータベース(DB)は、たくさんの大学・病院でそれぞれの研究所が作ったので、バラバラでどこに何があるのかわからない、信頼性の高い注釈が付いていないから使いにくい、大型プロジェクトはあるがなかなかDBが公開されない、バラバラに構築・管理されていて検索・解析・応用が困難、さらにこれをまとめる戦略や中核センターがない、という数々の問題がありました。

そのためにライフサイエンス統合データベースセンターを作って、統合DBプロジェクトが始まったわけです。DBの所在情報を明らかにするために、Integbioという1543件(うち国内1105件)のDBカタログができていますし、簡単な横断検索(キーワード検索)もできています。DBを作ったのはいいが、研究終了後のメンテナンスが非常に難しいという点に関して受け入れ体制ができています。権利関係に関してはクリアなもののみを受け入れていますので、自由に再利用できるという形になっています。

こういうことを実現するために平成23年度以降、JSTのNBDC(バイオサイエンスデータベースセンター)と共同で統合DB事業を進めて参りました。DBの利用に必要な、整理する・探す・つなげるの技術を開発し、データそのものにアクセスして、意味をちゃんとわかった上で活用できるようにしましょうということです。

RDFによるデータ統合とSPARQLによる検索

データにはどうやって取ったのかという説明が必要ですし、再利用するためにはいわゆるメタデータが必要です。使っている用語に関しても取り決めや辞書が必要です。次世代のデータドリブンサイエンスのためのDB、統合DBを作るために、RDFというセマンティックウェブの技術を使うことを考えてきました。データそのものはポータルもできていますし、どこに何があるかという情報もどんどん貯まっていますが、それを統合的に利用する技術を開発しているということです。

RDFとは、Resource Description Framework の略です。データを主語＝モノのID、述語＝オントロジーで定義された属性、目的語＝別のモノのIDまたは値、この3つの組で表現するわけです。それぞれに、URIというユニークな番地を付けて管理するというものです。

たとえば遺伝子には名前がありますし、塩基配列がありますし、たとえば薬だったら作用・副作用がありますし、構造もあります。遺伝子が何かの薬剤に対して阻害されるとか、色々な関係も出てきます。そういうものを重ね合わせていくと、どこかでつながってくるわけですね。それをつなげて行きますと、2つのものの関係が見えてくる。これまでわからなかった関係が大量のデータの中から見えてくるということです。こういう情報を提供したい、これが統合DBの概念です。

それを探すときに使うのが、SPARQLという検索用の言語です。SPARQLの文を書きますと、適切な関係のものがデータとして出てくるということです。例えばLODQAというシステムでは、自然文で「アルツハイマー病に関連する遺伝子は何ですか?」と質問すると、SPARQL検索を行い、関係する遺伝子のリストが出てくるというものです。これをアカデミックなデータに全部適応できるようにしようということで頑張っております。

RDFデータの利用促進と国際的標準化

これまでJSTと一緒に、微生物ゲノム、タンパク質立体構造、遺伝子発現などさまざまな研究グループと協力をして、DBをRDF化していくことを進めてきました。私たちのセンターはそれをサポートするということです。すでに10数件のDBがRDF化できましたので、これらの利用促進を進めたいと思います。

DBの標準化は非常に重要でして、国際的なBioHackathon会議を通じて、ヨーロッパ、アメリカの方々と一緒になってRDFに関する標準化を進めてきた次第です。セマンティックウェブは平成13年に提唱され、その後だんだん広がってきています。ヨーロッパのUniProtというのが一番早かったですけれども、私たちも結構早く始めていて、阪大のPDBjも割と早く始めています。最近ではアメリカ、ヨーロッパのDBがどんどんRDF化を進めているという状況です。データのメンテナンスが非常に楽になりますので、コストも下がるというメリットもあります。こうやってBioHackathonをやって標準化を進めてまいりますと、日本、アメリカ、ヨーロッパ、たくさんのDBがつながっていくということになり、分野や国境を越えたDB統合が実現してくるのです。

次世代生命科学データベースの実現に向けて

- 次世代生命科学データベース＝データ駆動型サイエンスを実現するデータベース
- データ駆動型サイエンスにおいては、新規データ生成も必要ながら、膨大に蓄積されたデータを効率的・効果的に再利用する必要がある→データインフラの整備
- そのためには、データのセマンティクス(データの意味)を扱うことが不可欠
- また、データ処理の大幅な省力化も必要



RDFの採用

=セマンティック・ウェブ
= Linked Open Data

「人文学オープンデータ共同利用センターの取組み」

人文学オープンデータ共同利用センター／準備室長
北本 朝展



オープンデータの3つの価値と3つの関係者

オープンデータの価値にはさまざまな側面があり、「利用」が促進できる、「透明性」が確保できる、「みんなが参加」できるという基準があります。我々のセンターは、この「参加」という側面を重視したいと思っております。

「本は図書館に行って読めばいい」とも言われますが、実際そうはいかないという側面もあります。例えば、図書館が1000km離れた場所にあったらどうするのか。人文学の本は一点物が多く、1000km移動しないと読めない場合も多々あります。また、例えば米国の研究者が東アジアの本を研究したいと思ったら海を渡って来ないと読めません。そのとき、例えば日本と中国を比較して中国のデータの方がオープン化されていたら、そちらの方が研究しやすいので研究対象に選ばれやすいという状況も生じています。いろいろな人が参入できないと、研究分野自体が盛り上がらないという状況があり、そういう点からもオープンデータ化を進めていくことが重要な課題となっています。

オープン化にあたって、3つの関係者について考えています。一つは「研究者」です。ここでは、データをより深く読み、分析して新しい知見を得ていくという方向性が重要になります。

一方、大量の情報を高速に処理するためには機械の力が不可欠ですから、いかに「機械」が使いやすいデータに整えるかということも重要な課題になってきます。

また人文学データは非常に多様ですので、定型化されたデータを自動的に処理するだけでは済まない場合も結構あります。そこには「市民」の協力が必要ですし、また市民自体が学習して賢くなっていくということも、オープンサイエンスの重要な側面だと考えております。

情報学と人文学の協働で歴史的典籍活用を推進

人間文化研究機構の国文学研究資料館との共同研究についてご紹介します。これは我々の人文学オープンデータ共同利用センターと統計数理研究所、国立情報学研究所の研究者が一緒になって取り組んでいるものです。国文学研究資料館では、日本の歴史的典籍30万点をデジタル化し、国際共同研究を推進する大型プロジェクトを進めています。そこでまず、研究者のためのオープンデータとして、「日本古典籍データセット」を公開しました。ここで公開した古典籍は、くずし字で書かれていて、江戸時代のくずし字を読めればこのデータも活用できます。最終的に30万点をデジタル化する計画ですが、現在はまだ700点の公開に留まっていますので、今後どんどん増やしていくことが重要です。なお、翻刻テキストは一部の古典籍のみで、基本的には画像ファイルに書誌のメタデータなどを同梱したオープンデータとなっています。こうした画像公開によって、本を見ることは可能となりましたが、データを活用できるのはくずし字が読める研究者だけです。くずし字をスラスラ読める人は日本に数千人程度しかいないと言われており、これだけだと日本人でもほとんど活用できないことになります。

そこで、機械が使えるデータにしないといけません。では機械が使えるデータとは何か？ それはどこに何の文字が書かれているかということを機械に教えるための、いわゆる教師付き学習データというものです。「日本古典籍字形データセット」はそのためのもので、これを

使えば機械が学習して文字を認識するということが可能になります。

たとえば、「あ」という文字には、我々が読める「あ」と別の形の「あ」があります。これは変体仮名と言っていて、元となる漢字が異なる「あ」なんですね。江戸時代は「あ」が何種類もあったのですが、今は一つの「あ」しか残っていません。

機械が学習するデータとして、現在86,176文字のデータがあり、今年度末には40万文字以上になる見込みです。座標情報が入っていますので、この画像のどこに何の文字があるかという認識もできます。このようなデータを使って学習すると、文字認識のプログラムが作れます。ただ、この40万文字というのは実は全然十分ではないと考えていて、もう2桁ぐらい多い文字数がないとなかなか全部は読めないのかなとも思っています。

江戸料理レシピデータセットに未来を感じる！

江戸時代の料理の本に冷やし卵羊羹というレシピが載っていますが、これを読んで作れる人はほとんどいません。しかし、それを「江戸料理レシピデータセット」として、現代の人が料理できるような形にすれば、使えるようになるわけです。

デジタル化しただけでは十分とはいえません。くずし字を翻刻すればいわゆる古文の状態になりますが、古文も読めない。そこで現代語訳し、さらにレシピ化すると、そのまま使えるようになります。さらにこれをクックパッドで公開します。センターのウェブサイトで公開するだけでなく、クックパッドのような一般の人が使うプラットフォームに乗せることによって、データをファインダブル(可視化)することに意義があります。これには、非常に大きな反応がありまして、国立情報学研究所のTwitterアカウント史上最大の反応がありましたし、「江戸料理+クックパッド」に未来を感じるというツイートもありました。関心が非常に高いデータ公開だったと言えます。

近代の本もOCRが難しく自動解読できていません。統計数理研究所と国立国語研究所との共同プロジェクトとして、今後取り組んでいきたいと思っております。

このような課題を、我々はディープアクセス技術と呼びたいと考えております。メタデータにアクセスするというのではなく、データの中身、コンテンツの中身にまでアクセスできる技術が必要です。

市民のためのオープンデータ
http://codh.rois.ac.jp/edo-cooking/

江戸料理レシピデータセット (CODH制作) 日本古典籍データセット (国文研所蔵) を翻案

2017/2/20 情報・システム研究機構シンポジウム

講演

講演／データサイエンス共同利用基盤施設における具体的取組みの紹介 「ゲノムデータ解析支援センターの取組み」



ゲノムデータ解析支援センター／センター長
野口 英樹

次世代シーケンサーがゲノム解析を変えた

我々のセンターは、大学などいろいろな研究機関の研究者がお持ちのゲノムデータをバイオインフォマティクスの技術を用いて解析し、そのための解析技術を提供するセンターになります。

ではゲノムデータというのはどういうデータなのかということですが、動物、植物、微生物までどんな生物種であっても、細胞の中にゲノムDNAを遺伝情報として持っています。DNA以外にもそこから転写された転写産物(RNA)なども存在していて、これを計測したデータがゲノムデータということになります。

計測機器にはいろいろなものがあり、例えばDNAマイクロアレイとか、最近ですとアイリスとかいろいろ装置がありますが、現在メインで使われているのはDNAシーケンサーです。DNAの情報、ACGTの配列を直接読み取るという装置を使って出てきた塩基配列データが、当センターで主に取り扱うデータということです。

ただ問題はこれらのデータというのは基本的にはどれもこれもACGTの配列なんですけれども、生物種が違ったり、実験条件が違ったりすると、その意味することというのは大きく変わってくるわけで、当然解析方法も変わってきます。そういう意味で多様であるということと、もう一つの大きな問題はデータが大量に出てくるということなのです。

これによって、なかなか解析が困難だという状況がございます。そういう大量のデータが出てくるようになった背景に、この次世代シーケンサーと呼ばれる新しいシーケンサーの登場があります。

国立遺伝学研究所の先端ゲノミクス推進センターに置いてあるシーケンサーには、Illumina HiSeq2500とPacBioRS IIというものがあります。ゲノムのDNAというのは、例えば人の場合、3G、30億塩基対あるわけなんですけれども、このシーケンサーはその30億塩基対を端から端までずらっと連続で読み取れるというものではありません。ゲノムの非常に短い断片を読むことしかできないのです。その一個一個の断片を配列リードと呼んでいます。そのリードの長さも、Illumina HiSeq2500で最大250 bp、PacBioRS IIの場合、若干長いのですが、それでも平均15 Kbp塩基対しか読めない。非常に短い配列しか出て来ないということなのです。

ただ数の方は大量で、Illumina HiSeq2500で一連あたり、最小単位で約2.5億本くらい出てくる。つまり掛け算して約60 Gbのデータというものが一度に得られるわけです。PacBioRS IIの場合、若干数として少ないのですが、それでも前世代のシーケンサーが96本とか384本でしたので、それに比べると、100倍以上のデータが一度に得られるようになってきています。

シーケンシングに係るコストも、次世代シーケンサーの登場した2007年以降、非常に速いスピードで低下してきています。コストが下がってくると、皆さん使いたくなっていくわけです。ゲノム情報とい

うのは生物の基盤の情報になりますので、使えるのなら使いたいということで、いろいろな分野の先生方がゲノムデータを取得するようになってきました。ただ、これまで配列を扱ってこなかった先生が、いきなり1ファイル50 Gbとか60 Gbになるようなデータを数十億本とか扱わなければいけないということになっててもなかなか対応ができないわけです。ですから、解析されずに眠っているデータもたくさん出てくるという状況になっているわけです。

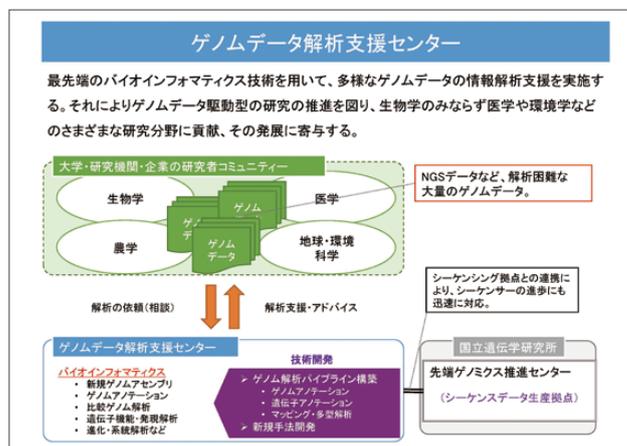
多様なゲノムデータの情報解析で研究を支援

我々のセンターでは、最先端のバイオインフォマティクスの解析技術を用いてこれらのデータを解析し、共同研究という形で研究を回したり、あるいはその研究者が求めるような結果の解析を行うためにはどういうデータを用意して、どういう解析手法をとるのかという技術的なアドバイスなどをしていくことで、ゲノムデータ駆動型の研究の推進に貢献することを目的としております。

実際にどういうデータ解析を行っているかということですが、新規にゲノムを決定する「De novoゲノムシーケンス」、まだゲノム配列が決まっていない生物種のゲノムを決定する、あるいはもうすでに決まっている別の個体とか変異系のものや変異仮想を調べる「ゲノムリシーケンス」、あるいは発現している遺伝子の解析をする「トランスクリプトーム解析」、さらに「メタゲノム解析」と幅広いデータ解析を行っています。

今年度は、11課題で支援依頼がありまして、生物種としてもかなり幅広く哺乳類、昆虫、魚類、植物、あるいは原核生物と、非常にさまざまな種類の生物種の解析依頼が来ております。

また、この解析支援を続けていく上で、効率化のためのパイプラインの構築、高速・省メモリな解析アルゴリズムの開発、解析講習会などの人材育成などに取り組んでいきたいと思っております。



講演／データサイエンス共同利用基盤施設における具体的取組みの紹介 「データ融合計算支援プロジェクトの取組み」



データ融合計算支援プロジェクト
中野 慎也

シミュレーションと観測を融合した研究手法

データ融合計算というのは一般的な用語ではないかもしれませんが、シミュレーションとデータサイエンスを融合した、この研究手法のアプローチをいろいろな分野で紹介していくことが、このプロジェクトの概要です。

すでに研究相談の受付を開始しておりまして、共同研究の実施を通して他分野に展開しており、さらに随時共同研究を受付中といったところなんです。

データ融合計算に関するノウハウや手法、ソフトウェア化したものを整備し、今後の公開に向けての準備も進めております。

また来年度以降、セミナーやハンズオンの活動を通じて、我々の研究手法をいろいろなところに紹介していこうと考えております。

具体的には、シミュレーションと観測を統合した研究手法というものに取り組みようとしております。

シミュレーションというのは、第3の科学と言われて、ひとつ前の世代の研究手法かもしれませんが、今でもスパコンなどいろいろな分野で高精度なシミュレーションが展開されています。

ただ、スパコンの性能が上がるとともに計算量がどんどん増大したことで、インプットを与えないと結果が出て来ないという問題がシミュレーションにはあります。

一方、観測データですが、最近はいろいろな種類のデータが取られていますけれども、そのデータから意味のある情報や知識をどう抽出するかが問題で、本質的な情報が見えにくくなっているという現状もあります。

そこで我々は、シミュレーションと観測を融合させ、第3の科学のアプローチと第4の科学のアプローチを融合させたような手法を提案させていただき、いま展開しているところであります。

シミュレーションというのは、入力変数と入力パラメータを与えて出力を何か返すというものです。大体出力というのは部分的には観測できるものになっています。シミュレーションを実行するときの計算量がすごく重いという状況の下で、何を入力したのか、どういうパラメータ設定で実行したのかということを知りたいという方法として「データ同化」と呼ばれる手法があります。

もう一つは、データサイエンスのアプローチを使って、シミュレーションを模倣するモデルを作るという手法があります。シミュレーションは計算時間がすごくかかるものなので、これを簡単なモデルで、ある程度ローコストで予測を行う「エミュレータ」というものを作るという手法です。この2つがシミュレーションと観測の融合の主な手法です。

システム設計と検証の作業を一体化

観測データからいろいろなデータを組み合わせで解析し、それを

シミュレーションのデータ同化と呼ばれる手法を使って、シミュレーションに食べさせてやります。さらにそのシミュレーションの振る舞いを模倣するようなモデルを作れば、どういったところのデータを取ればシステムの改善ができるかといった判断ができますので、そこからまたデータを取り、解析するというサイクルになります。

これによって、システム設計やアウトプットの定量的予測、不確実性の評価などの実現を進めていけるような方法論を提示していきたいと考えているところです。

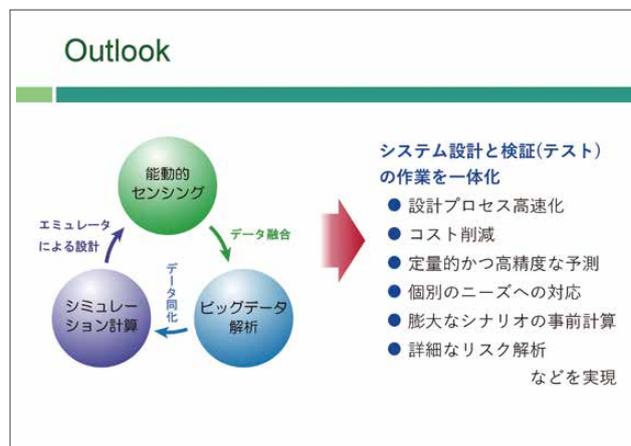
現在、我々は方法を研究しているところがございますので、具体的なテーマは外からいただきたいという立場にありますので、相談の窓口を開きましてメールなどで受付をしております。我々の持っているノウハウの提供で解決できる問題もあると思いますが、新しい方法論を開発しないといけないようなテーマの場合は、共同研究として実施させていただくといったことを考えています。

もちろん、ノウハウの提供で解決しそうな問題の場合でも、統計的な考えに詳しくない方もいらっしゃると思いますので、そういった場合には共同研究として実施させていただいて、より密な情報交換をしながら進めていくといった出口も用意しております。

データ融合計算プロジェクトの実績

これまでの実績としましては、細胞質シミュレーションのパラメータ推定(国立遺伝学研究所との共同研究)、南極氷床コア年代推定手法の開発(国立極地研究所との共同研究)、マルチモデルによる将来気候変動予測(防災科学技術研究所との共同研究)、人の動きとデータのシミュレーションから人の動きを推定するといった人流シミュレーションの状態推定(東京大学生産技術研究所との共同研究)など、いろいろな応用分野で研究を進めております。

また、民間企業との共同研究も受け付けておりますし、現在進行中のものもございます。



講演

講演／データサイエンス共同利用基盤施設における具体的取組みの紹介 「極域環境データサイエンスセンターの取組み」



極域環境データサイエンスセンター／準備室長
門倉 昭

4つのグループから多種多様なデータが集結

極域環境データサイエンスセンターの目的とするところは、国立極地研究所の所有するデータの公開と共同利用、有効利用を促進し、極域科学のデータ活動の中心を担い、地球環境研究に貢献するということです。

「宙空圏」「気水圏」「地圏」「生物圏」という、4つのグループがさまざまなデータを取っており、時間的に連続した時系列データと、大気や海のサンプルやアイスコアなどといった試料系データの2種類があります。

「宙空圏」では、IGY(1957-1958年)という昭和基地開設以来のオーロラデータがあり、昔の400ftのフィルムで記録されたデータもあります。南極、北極、多点の観測機を使ったオーロラデータも取得され、100Hzで取った最先端のデータなども蓄積されています。最近ではPANSYレーダーという最先端の大型大気レーダーが昭和基地に設置され、ここから大気の流れを観測するデータが生み出されています。また、北欧のEISCATという大型のレーダー組織の国際プロジェクトに日本の代表として国立極地研究所が参加しており、電離圏の観測データといった最先端のデータもあります。

「気水圏」では長年、温室効果ガスを観測しており、二酸化炭素の増加、地球温暖化に関するエアロゾルや雲の観測、人工衛星データからの南極域の雲や海氷の観測データもあります。ドームふじ基地では、地下3000mのアイスコアの分析から72万年前の歴史の解析が行われています。時系列データはせいぜい30年程度の蓄積ですが、試料系データは地球46億年という長いスパンの歴史を分析・解析できるデータになっております。

「地圏」では、岩石、隕石の試料が蓄積されています。時系列データで、地震、重力を観測しており、地球の重力の変化、地震波形を使った地球内部のモニタリング観測なども行われています。東日本大震災の影響は昭和基地まで及び、地震波、重力の大きな変化が観測されました。昭和基地とホバートの間の距離が年々離れているという測地データも蓄積されております。

「生物圏」では、ペンギン個体数の長期モニタリングや、最近ではペンギンあるいは大型動物の行動を解析する非常に先端的な研究がされています。陸上生物については、コケ類、種子類、さまざまな試料データが蓄積されデータベースが構築されています。

また、海洋生物から遺伝子解析をする設備やアイスコアラボラトリー、隕石ラボラトリーといった非常に大きなラボラトリーもあり、各種試料の高次処理・解析データも蓄えられ、データベースも作られています。

多彩な公開用(汎用)データベースシステム

さまざまなデータを外側に公開する仕組みとして、学術データベースがあります。これはメタデータベースで、いわばデータのカタログです。4つの研究分野のほぼすべての観測メタデータを見ることが出来るシステムも整備されてきています。

SCARという南極関係のコミュニティのデータ管理委員会と深く連携した形で開発していて、このメタデータをアメリカのGCMDに提供するために、GCMDのメタデータ形式に変換されて送られる仕組みもできています。南極観測や北極観測は常に国際的な共同研

究や連携が進みますので、データについても国際的な連携を意識した仕組みを作っております。

最近では北極域の観測が集中的に行われておりまして、GRENEあるいはArCSといった北極プロジェクトの関係のデータ・情報を扱うシステムとして、ADS(Arctic Data archive System)というものが開発されています。これは、メタデータを検索するのみではなく、実際極域で取れるいろいろな実データを検索・表示・オンライン可視化・解析もできるアプリケーションも備えた総合的なデータベースシステムになっています。

IUGONETというデータシステムもあり、大学間連携プロジェクトで作られた超高層大気関係のデータのメタデータシステムです。特に各機関が所有する超高層大気関係のデータを横断的に検索できるというものです。

ソフトウェアも開発され、アメリカの同様のシステムであるSPEDASとも連携しています。

南極GISというシステムもあり、これは主に南極域の地形図、地質図、航空写真、衛星写真などの地理情報を検索し、地図上に表示するための地理情報システムで、南極観測隊活動支援としても利用されています。

極域環境データサイエンスセンターの目標

個別のデータベースや分野限定のデータベースはあるが、極域科学全体を横断的に俯瞰出来る総合的な仕組みがないというのが現状です。データベース化や公開の進み方にもばらつきがあります。また、公開用データベースシステムは作られてきた目的や背景が異なり、それぞれに特化したシステムが並列する状態にあります。

当センターでは、極域データ全体を抱合し、検索・可視化、解析までできる総合的なシステムの設計や構築を一つの目標としております。

また、データを有効利用できるという意味では、国立極地研究所が今年の1月19日に創刊した「Polar Data Journal」というデータジャーナルを有効利用したデータ出版の積極的な促進も図っていきたいと思います。

国内・国際コミュニティとの積極的な連携やデータサイエンス・共同研究の推進を図り、極域科学のデータ活動の中心を担っていきたいと考えております。



全球的視野に立った観測とデータの共有

こちら昭和基地です。第58次南極地域観測隊の越冬隊長の岡田と申します。「データ発生の現場から」というテーマで、南極昭和基地における観測データの収集、発信の現状について紹介させていただきたいと思います。

現在の気温はマイナス5度前後、風速5m、快晴となっております。昭和基地は南極大陸から4kmほど離れた東オングル島という場所にあります。半径1kmあまりの範囲に70棟近くの建物があり、分散して基地の運営および観測活動を行っております。

基地の東地区と呼ばれる一帯には、衛星観測のデータを受信する大型アンテナのドーム、衛星受信棟、光学観測棟、情報処理棟、観測棟、環境科学棟という観測系の建物が並んでおります。衛星データの受信、オーロラの観測、気象観測、生物系の観測データを扱う建物が並ぶ一帯です。

我々第58次観測越冬隊は、2月1日に越冬交代を行い、第57次観測越冬隊から基地の運営・観測を引き継ぎました。2月15日に「観測船しらせ」が基地を離れ、最終便が飛び立ってから、第58次観測越冬隊33名による生活が始まってやっと一週間というところになります。

第58次観測越冬隊は、設営系の隊員18名のほか、観測系の隊員14名、私を合わせて33名という構成になっております。今年の隊には6名の女性隊員が参加しており、越冬隊の中では最多の人数ということになります。

大学の研究室からの隊員や、気象観測に携わる気象庁からの5名の隊員など、幅広い分野の隊員が参加して観測にあたっております。特に若手の研究者については、特任という枠になりますが、観測隊に参加して、観測や研究をするという隊員も増えてきております。

現在の昭和基地は、衛星回線を経由して常時3Mbpsの回線速度で、立川の国立極地研究所とつながっております。一日あたりにしますと、最大30GB前後のデータを電送できる回線容量ということになります。

昭和基地におけるすべての観測データを衛星回線で電送するというはまだまだ困難ですが、準リアルタイムに国内と観測データを共有する、あるいは国内・海外の共同研究者とリアルタイムで観測データを共有することができるというような体制になることは、今後の全球規模での観測の共有あるいは発信という意味では、非常に重要なインフラになってきていると考えております。

「データ発生の現場から」

南極昭和基地からの中継

第58次南極地域観測隊



現地との中継



南極昭和基地大型大気レーダー (PANSY)

PANSYを中心とした重点研究観測、定常観測、モニタリング観測、一般研究観測と4つの分類で実施しています。観測課題数としては全体で100以上の課題があり、14名の観測隊員が日々忙しく観測作業にあたっております。今日は宙空系の女性隊員2名、重点研究観測の江尻隊員と一般研究観測の鈴木隊員にも来ていただいております。

今後どの分野においても観測装置が非常に高性能化し、観測データ全体としては全球規模の観測というものが重要になってくると考えております。観測現場からという観点で申し上げますと、データの迅速な配信、それから国際的な枠組みの中での合意というのが、これから非常に重要になってくるのではないかと考えております。

重点観測課題である大型大気レーダーを使った観測、通称「PANSY」の観測を担当しております橋本隊員にも来ていただいておりますので、大型大気レーダーを使った観測について紹介してもらおうと思っております。

大型大気レーダーPANSYの果たす使命

京都大学大学院情報学研究所の橋本と申します。PANSYは地上30kmぐらまでの風、80~90kmぐらまでの風、電離圏の電子密度の観測などを行うことができる、非常に大きな大気レーダーです。PANSYのおもな観測ターゲットは、ブリザード、地表付近、成層圏、中間圏、オーロラといった南極特有の特殊な現象となっております。

これらの領域を幅広くカバーできる観測範囲がありますので、データ同化による地球大気モデル計算を組み合わせ、地球の大気循環システムの解明を行うということがPANSYの主たる使命となっております。

PANSYには全部で1045本のアンテナがあります。大気からのエコーというのは非常に弱いので、これぐらいたくさんのアンテナが必要になってきます。大体面積で東京ドーム2個分くらいあります。

観測データは、生データと呼ばれている状態で、一日に465GBぐら発生します。このデータをそのまま国内には3Mbpsの回線では送れませんので、オンラインで信号処理を行います。小さくして、スペクトラムと呼んでいるのですが、大体これで605MBぐら発生します。これぐらであれば衛星回線で送れるので、そのままリアルタイムで国内に送って、スピーディな研究ができるようにしています。生データも特異な観測の場合などはこちらでは解析できませんので、「観測船しらせ」が持ち帰ってからの研究という形になります。

「異分野融合・新分野創成を担う データサイエンティストの育成基盤」

情報・システム研究機構／理事・統計数理研究所／所長
樋口 知之



ビッグデータの利活用のための専門人材育成へ

情報・システム研究機構は、ユニークな構造をもっておりまして、分野横断的な統計数理研究所、国立情報学研究所、ドメインサイエンスである国立遺伝学研究所、国立極地研究所、この2つの縦軸と横軸でもって作られる、クロスポイントが機構全体をカバーするという仕組みを持っております。また、このクロスポイントが最近のさまざまなイノベーションの起点になるということが多々あります。

本機構では与えられたクロスポイントだけで活躍するのではなく、自らそういうクロスポイントを作っていくような人間、T型・II型人材の輩出を人材育成の目標としております。

本機構はこれまで機構本部でさまざまな事業を拡大、あるいは展開してまいりまして、現在は法人第三期にあたります。

各センターや準備室のプロジェクトで行われているさまざまな事業の中で、ポストクや若手の研究者などを雇用しておりまして、実務やサービスに携わりながらの人材育成を行ってまいりました。

一年半ほど前に、文部科学省、産業界から何名かの方にご参加いただきまして「ビッグデータの利活用に係る専門人材育成に向けた産学官懇談会」が行われました。その中で、社会全体のリテラシーやアウェアネスを向上させるために、全学的教養教育の実施や国家レベルのフラッグシップ・プロジェクトの推進などの取組みを盛り込んだ「ビッグデータの利活用のための専門人材育成について」という提言をさせていただきました。

その中では、10 大学程度で本報告書の提案に基づく人材育成をスタートすると提言したのですが、現在まだ6大学しか採択されておりませんので、引き続き増えていくことを強く期待しているところでございます。

我が国の問題の根源は、棟梁レベルの決定的不足にあります。本

機構は、この棟梁レベルという、高いレベルの人材育成に注力しているところでございます。

機構、各センターの具体的な人材育成への取組み

情報・システム研究機構では、若手研究者育成のための「ROIS/II-URIC 若手研究者クロストーク」という合宿形式の討論会を実施しております。異分野融合と交流の機会創出を目的として、最近はこの機構だけではなく、他の大学共同利用機関法人4機構の若手研究者などを巻き込んで実施されました。

今年のテーマは「分かり合えるコミュニケーション」ということで、研究を推進するためのコミュニケーション力を付けるという課題に取り組んでおります。

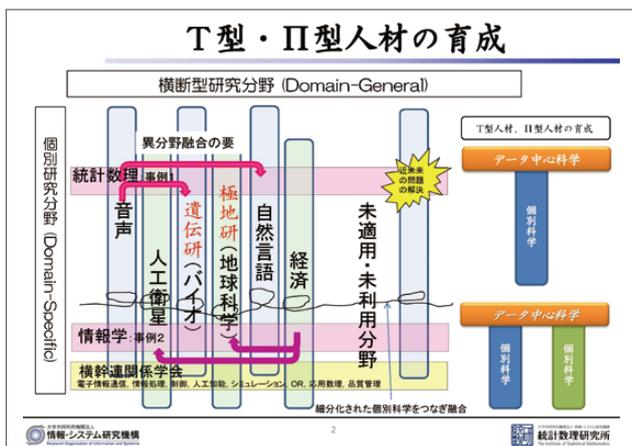
組織的に人材育成を行っております、いくつかの機関の取組みをご紹介します。ご協力いただきたいと思います。

ライフサイエンス統合データベースセンターでは3つの取組みを行っております。

まずは、国際開発者会議BioHackathonです。これは、生命科学分野の統合データベース構築のための技術基盤の確立を目的としています。これも合宿形式で行われ、密度の高い情報交換と生産性の高い技術開発が会議の場で活発に議論されています。

ライフサイエンス統合データベースセンターはこの分野で世界をリードするハブとして高く評価され、今後も継続的にリーダーシップを発揮することが期待されております。その成果は、ライフサイエンス統合データベースセンターの開発する統合データベースの基盤技術として活用されているほか、研究者コミュニティの共有資産となり論文発表も行われている事業です。

2つ目が、統合データベース講習会:AJACSです。これは少し初



情報・システム研究機構 統計数理研究所

【提言】ビッグデータの利活用のための専門人材育成について

我が国の問題の根源は、**棟梁レベルの決定的不足**にある。この解決のために**国家レベルの拠点を設置**して、年500名規模の「棟梁レベル」の人材育成をめざし、上層への成長や下層へのトリプルダウン効果も狙う。

リテラシーレベルや独り立ちレベルの**大学教育を加速させる**ために、主要10大学程度で本報告書の提案に基づく人材育成をスタートすると共に、MOOCなどのオンライン教材を整備し、全国への波及効果を狙う。

社会全体のリテラシーやアウェアネスを向上させるために、全学的教養教育の実施、国家レベルのフラッグシップ・プロジェクトの推進、コンテストの開催、映像素材の充実などの取組を行う。

大学共同利用機関法人 情報・システム研究機構
ビッグデータの利活用に係る専門人材育成に向けた産学官懇談会
「ビッグデータの利活用のための専門人材育成について」
平成27年7月30日

これらの方策を実現するにあたっては、**データサイエンスを副専攻とするダブルディグリー制人材育成の推進やスキル認定制度も有効と考えられる。**

心者の方を狙っています、実際に手を動かして、生命科学系のデータベースやツールの使い方を学ぶという、初心者向けのハンズオン講習会を全国各地で数多く開催しております。

3つ目が、SPARQLthonです。これはセマンティックウェブ技術による生命科学DBの統合をテーマとして、さまざまな課題解決を実際にやっていくという開発会議で、これまでの参加者は第一回からの累積で1275人となっております。

平成26年度からは、統合化促進プログラムから多くの方にご参加いただいております、データのRDF化に関する意識共有・技術情報共有・コラボレーションが著しく進み、また、ここで多くのRDFデータやオンコロジーが作られてきたという事業です。

統計数理研究所には、「統計思考院」という人材育成を目的とした組織があります。統計数理研究所の今期中期目標・中期計画にある「統計思考力を備えたT型人材育成による融合研究の推進」を実現するため、新しい統計学の創成を目指す研究者や学生、固有分野の研究で統計学の必要性を感じた人などさまざまな人が集い、切磋琢磨しながら「統計思考」の訓練や研究をするというものです。

こちらの活動とプログラムをいくつか紹介させていただきます。

まず、共同研究スタートアップです。2つの目的があり、一つは持ち込まれた課題を解決することです。もう一つは、課題に関わらず必ずメンターと若い研究者を組み合わせ、課題解決を通して人材育成を図るということです。

また長年に渡り、公開講座を実施しております、大体年間10～10数講、参加者は大体900人くらいです。今は2/3が民間からの参加者となっております。

さらに、もう少しハイエンドなもので、組織連携(MOUや組織長同士の合意)に基づくデータサイエンス講座企画も進めております。理化学研究所のAIPという最も高いレベルの「機械学習速習コース」を開きつつあります。これをe-learning教材にして、高いレベルの人材育成を行っております。

また、夏期大学院という連続10日間、朝から晩までぶっ通しの講義(いわば、統計数理ブートキャンプ)なども開催しております。

研究分野におけるオールジャパン一流の講師陣で、さらに外国からの著名な講師も招へいし、実践プログラミングまでのスキル向上を目標とするコースとなっております。

それから、民間の方々が自学自習できる、データサイエンス・リサーチプラザという仕組みもあります。「受託研究員制度」を利用した産学官(特に企業)への共同利用提供で、企業側のニーズに応じた人材育成や人材交流面でのフォローを行うというものです。

国立情報学研究所でもいくつかの取組みがあります。その中で2つ紹介させていただきます。まずは、トップエスイーというもので、トップレベルのIT技術者育成を目的とした教育プログラムで、座学的なものから、実際の演習を通して学んでいくというものです。毎年約40名の企業IT技術者を教育、過去11年で348名の育成実績があるというものです。

もう一つの取組みは、情報セキュリティに関するものです。今年度始まったばかりですが、学術ネットワークSOCを通じた技術職員の育成となります。大学間連携に基づく情報セキュリティ体制の基盤構築を目的とし、国立情報学研究所の方で、橋渡し人材育成を通して全国の大学の面倒を見るというものです。

また、来年度の計画として、実学として効率よく履修してもらうために、学術研究機関へのベンチマークデータの提供なども企画されています。

基盤機関としての大学院教育に関する取組み

最後は大学院教育に関する取組みです。4研究所とも、総合研究大学院大学が入っております、複合科学研究科が生命科学研究科を受け持っております。専攻の学生は本機構のクロストークなどで交流するということも行っております。大学共同利用機関が有する優れた人材と研究環境を活用して博士課程の教育を行い、一流の研究者を養成する総合研究大学院大学の基盤機関として大学院教育を実施するとともに、全国の国公立大学とも連携して学生を受け入れるなど、教育機関としても重要な役割を担っております。

情報・システム研究機構
若手研究者育成：
「ROIS/I-URIC 若手研究者クロストーク」

①目的等：異分野融合と交流の機会創出。
合宿形式。
今年度は講演・研究紹介・討議の全てを英語で実施。

②対象等：4機構の若手研究者+学生。
総研大との共同開催。
(H28実績：参加総数50名：
ROIS 32名、3機構8名、
総研大生10名)

③平成28年度テーマ：「分り合えるコミュニケーション」



情報・システム研究機構

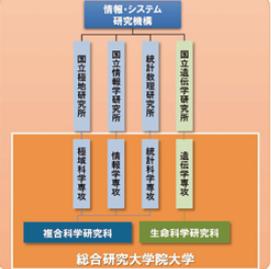
情報・システム研究機構
総合研究大学院大学
大学院教育に関する取組

総合研究大学院大学の基盤機関として

- 大学共同利用機関が有する優れた人材と研究環境を活用して博士課程の教育を行い、一流の研究者を養成する総研大の基盤機関として大学院教育を実施。
- 各機構長が総研大の経営協議会や学長・機構長連絡協議会に、また、大学共同利用機関等との連携を強化することを目的とした総研大のアドバイザリーボードに各機構長が推薦する理事が参画するなど、大学院教育の運営に積極的に関与。

全国の国公立大学との連携

- 各大学共同利用機関は、全国の国公立大学から、特別共同利用研究員等として学生を受け入れ。
- 連携大学院等を通して各大学院独自の教育方針と大学共同利用機関の優れた研究環境が融合した多様な取組が行われており、その役割は重要。



平成27年5月1日現在	情報・システム研究機構
総合研究大学院大学	170
特別共同利用研究員等	31
連携大学院等	55

情報・システム研究機構

「情報・システム研究機構の新時代に向けて」



情報・システム研究機構／理事／次期機構長
藤井 良一

研究所のミッション遂行と機構の果たす役割

情報・システム研究機構のミッションは、生命、地球、環境、社会などの複雑な問題を、物質とエネルギーの観点に替って「情報とシステム」という立場から捉えるための方法の研究、研究基盤の整備および融合研究による新分野の開拓を行なうというところでございます。

機構を構成する4研究所は、国立極地研究所と国立遺伝学研究所という分野型の研究所と、国立情報学研究所と統計数理研究所という分野によらない基盤的な科学を行う研究所から成っております。各研究所の歴史は、2004年の機構の設立より十分に早く、各々のコミュニティからの付託を受けて、先端的な学術研究、共同利用、人材育成を行うことをミッションとしてきております。

これらの研究所は各学術領域でCOEであることが求められております。これは大学共同利用機関として、最も優先度が高い大学共同利用、共同研究事業を行うための前提条件になるものと考えております。世界をリードする技術レベルを持って初めて、全国の大学などの研究機関や企業の方々が共同研究をしたいと思っただけ、かつ大学などの機能強化に貢献できると考えるからでございます。

機構は研究所による最先端研究と科学の発展、それを支える基盤を支援することが役割と考えております。研究所と機構の共同により、社会や大学などと連携して基盤的な環境を提供し、基礎的そして応用的な成果を社会に還元することが求められていると自覚しております。

2016年に法人第三期に入りまして、それまでの事業の成果を踏まえ、研究所の力を結集して、今後訪れるオープンサイエンス時代に向けて、機構本部の機能をより高めようとしております。

一つは、本部に戦略企画本部を立ち上げ、各研究所の執行部から参加をいただき、機構全体の研究戦略や将来計画、そしてガバナンス施策などの提案を行う組織を立ち上げました。

また、データサイエンス共同利用基盤施設を立ち上げまして、研究所が作り出すデータのデータベース統合を行い、オープンデータ、オープンサイエンス推進の加速を目指しております。これにより大学の共同利用の拡大とともに、統合された異分野データから新分野創成につながることを期待しております。

社会と学術のデータサイエンスへの要請

情報やそれに関連するテクノロジーは現在非常に急激に発展を見せております。急激なICT化と多種多様なビッグデータの出現、それに呼応しようとする計算能力の成長は、社会を変容させ、研究環境を大きく変化させております。機械学習は発達し、社会での応用は拡大しつつあるということです。人間の能力は、例えば2500年前の孔子の時代からあまり変わっていない、進歩していないように見

えますけれども、テクノロジーの進展は凄まじくて、素人目にはビッグデータやAIは人間が制御できる範囲を既に超えているのではなかという不安に陥るほどでございます。

国立極地研究所が共同運営しております、欧州にある超高層大気観測所の大型レーダーEISCAT(アイスキャット)では、数年後の2020年過ぎに新しいレーダーを建設して観測を開始する予定ですが、最速サンプリングでのデータ量は1.6 Tbit/秒(生データ)でございます。リアルタイム処理で1/30に減らしますが、それでも54 Gbit/秒、3.2 Tbit/分、4.7 Pbit/日、1.7 Ebit/年という巨大な量になります。もちろん連続運転は不可能でございますけれども、いかにインテリジェントに有用なデータを引き出すか、数年後には情報・システム研究機構として対応が迫られるチャレンジングな課題となっております。

今後データベースがすべての学問と産業の基本となることが予想され、それに呼応して第4の科学としての「データサイエンス」時代が到来して、データ共有を通じて研究主体が個人からグループ中心へ変化すると言えます。

また研究対象も生命、地球、環境、人間社会の複雑な現象と解決すべき諸問題が重要課題となって、さらに社会の課題に応える分野融合研究、新学術創成が求められているということでもあります。

すなわち社会的要請の変化によりまして、課題解決型研究への移行、科学技術イノベーションの牽引・推進、超スマート社会への貢献が求められ、これらの基礎となるオープンサイエンス化の推進が強く求められているということでもあります。

ビッグデータには膨大な知識や価値が埋もれておりますけれども、現在の方法・技術では必ずしも十分な有効活用はなされていないと言えます。新たな科学的な手法によりまして、知識発見や価値創造を行うことが必要だということで、それが正に第4の科学と呼ばれているデータサイエンスに他なりません。

「情報とシステム」の現況＝データサイエンス時代 情報・システム研究機構

急激なICT化と計算能力の成長 ～研究環境の変化～

- 計算能力の向上、ICTの発達、多種多様な**ビッグデータ**の出現
- AIが発達し、社会での応用が拡大

データサイエンスの時代 ～研究方法の変化～

- データベースが全ての学問と産業の基盤
- 第4の科学としての「**データサイエンス**」時代の到来
- データ共有を通じ、研究主体が個人からグループ中心へ変化

複雑化する社会に対応した研究の要請 ～研究対象の変化～

- 生命、地球、環境、人間社会の複雑な現象と解決すべき諸問題
- 社会の課題に応える分野融合研究、新学術創成

より密接に社会と関わる科学へ ～社会的要請の変化～

- 課題解決型研究への移行
- 超スマート社会への貢献
- 科学技術イノベーションの牽引・推進
- オープンサイエンス化の推進

今後、データサイエンス、オープンサイエンスは急速に進展

6

ビッグデータの利用には、レーダーの例で申しましたけれども、大量の散在するデータをリアルタイム処理する技術、大規模なデータ処理技術、それから膨大な高次元データや計算結果を人間が把握可能にするデータ可視化、そしてビッグデータから意味ある知識獲得のための方法、データ解析手法が必須となります。

分野融合・新分野創成のための機構の施策

情報・システム研究機構では、中期的対応としまして、データサイエンス、オープンサイエンスを推進して、その結果として分野融合研究や新分野創成を行うために、データサイエンス共同利用基盤施設をフラッグシップとして、データ共有支援、データ解析支援を行い、オープンサイエンスを推進していきたいと考えております。

また、学術動向や社会の要請の変化に対応するためには、従来の縦型の教員配置に加えまして、横断型の教員配置を作りまして、スピーディーに対応していくことが重要で、これは各研究所で実施されつつあります。

このデータサイエンス共同利用基盤施設は、機構の研究所の所有する、または他機構との融合で得られる、今までは個々に利用されてきた生命、遺伝子、地球環境、人間社会、人文学に関するデータベースにメタデータやRDF化などを行いまして、データ共有の支援をし、統合データベース化します。さらにデータ同化などのデータ解析支援を行うことにより、データベースを相互に利用することを可能にして、データ活用を格段に高めようとするものでございます。その結果、データサイエンスが推進されて、異分野融合が進むことで、新分野創成が可能になると期待しております。

統合されたデータは、大学などのすべての研究者が利用可能になるというものであります。また、国立情報学研究所は全国の大学のために、オープンサイエンス研究のデータ基盤の準備を進めております。さまざまなデータベースを管理・公開するための基盤で、データサイエンス共同利用基盤施設のデータベース、これも将来、この基盤の上に乗ることになると考えられます。

一方、大学や研究機関などの学術界や、国や自治体、それから企業と産業界など社会にはさまざまなデータがあります。データサイエンス共同利用基盤施設での統合データベース構築のためのデータ共有支援やデータ解析支援を横展開することによりまして、学術界、社会全体のデータベースの統合と相互活用への道が拓かれて、オープンサイエンス実現への貢献ができると考えております。

これによって、全体のデータが統合されるということと、大学の機能強化や異分野融合、新分野創成の促進にもつながると考えております。共同利用はいくつかのパターンがございますけれども、全研究者のための基盤の提供は大学共同利用機関の大きな役割であると考えております。全国に張り巡らされたSINET5、セキュリティ、クラウドそれからさまざまなCiNiiとかJAIROなどの学術情報の提供などがこの基盤にあたる考えます。

時代が求める研究者・高度技術者を社会・大学などに輩出

また、データサイエンス共同利用基盤施設も全国に提供される基盤施設の一つでございます。共同利用機関の大きな役割の一つは

人材育成であります。情報・システム研究機構の各研究所はデータや情報に関するさまざまな人材育成を行ってきております。統計数理科学人材の育成、それから情報科学人材の育成、そしてバイオインフォマティクスの育成と、大変精力的に行われてきておりまして、ほかにも多くのプログラムが計画されているとお聞きしております。機構としては、こういうもの見える化を図り、支援していきたいと思っております。

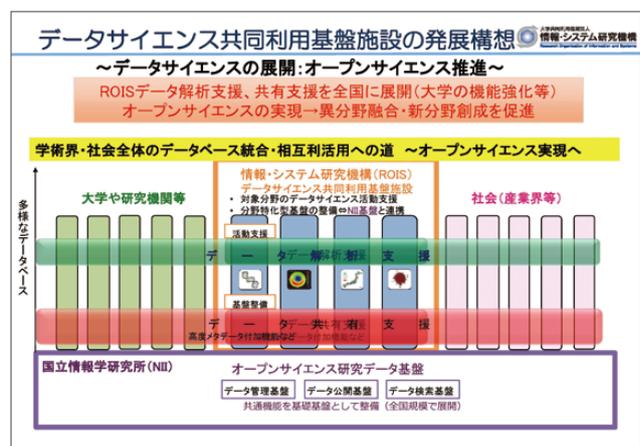
一つの課題として挙げたいのは、社会が必要としている人材数と研究所が育成できる人材数との関係です。統計数理研究所では棟梁レベルの人が毎年500人は必要と言われております。統計数理研究所はその主要な部分を育成することを期待されていると思いますが、これは大変なことだと思います。人材育成制度の整備と拡大が急務であると機構でも考えております。社会や大学などのニーズに、質だけでなく量的にも応えることができるように、機構として支援をしていきたいと考えております。

もう一つ重要なのは、人の交流、循環です。共同研究のレベルを超えて、研究者の相互の循環は両者の活力を高めて、大学の機能強化とともに、我々研究所・機構の発展、双方に貢献すると信じます。機構では、研究者交流プログラムがありますが、今後さらにクロスアポイントメント制度などを用いた双方向の人材循環を支援していきたいと考えております。

一方、大学院生の教育は、機構の最も重要なミッションの一つでございます。機構は、総合研究大学院大学の一員として、大学院生を受け入れております。研究所で最先端の研究に接して、また野外のフィールドで生の自然の神秘に接して感動するということは、学生にとって得難い経験でありまして、学術研究に進む人材だけでなく、社会において指導的な立場に立つ総合的能力を持つ人材を輩出することを可能にすると思っております。

できるだけ多くの学生たちにこの教育研究環境を経験してもらうために、総合研究大学院大学を中心に、連携大学院や特別共同利用研究員制度、インターンシップなどを十分活用できるように機構として支援していきたいと思っております。

各研究所の学理をもとに学術基盤の充実を図って、データサイエンス、オープンサイエンスを推進して、今まで以上に大学などの機能強化や社会のイノベーションに機構一丸となって貢献していきたいと思っております。



大学における データサイエンスと その教育

九州大学
理事・副学長
安浦 寛人 氏



低年次から高年次一貫型のデータサイエンス教育

九州大学では、全学生にPCを必携としたBYOD (Bring Your Own Device) を徹底し、電子教科書でしか勉強できないコースもあります。出席管理やレポート管理などの「e-Learningシステム」、授業日誌の記録「eポートフォリオシステム」、講義スライドの共有「デジタル教科書配信システム」を結び付けて教育に活かすとともに、学生たちの学び方の実態をデータ(1日約18万件のログ)で蓄積しています。例えば電子教科書の閲覧状況や引いたマーカーの数などを調べることができて、先生は学生の予習の達成度を事前に把握できます。達成度が低い場合はイントロダクションから始める、達成度が高い場合はいきなりクエスチョン&アンサーから入るなど授業内容を適応的に変更できるわけです。これは学生の学び方の改善のためだけでなく、先生にとっても使えるツールで、学生の閲覧パターンから教材の改善点も見えますし、授業の設計の良し悪しを

データは誰のものか

人間文化研究機構
理事
佐藤 洋一郎 氏



人文系データの特質とオープン化の事例

人文系データ(一次資料)の作出時期は、古いものになると1000年以上前のものもあり、そういう古典籍などは十分にオープン化されており、公開という点では意外と先行的です。ただ、変形文字や多言語などの問題で、可視化にはときに「翻訳」が必要です。

また、権利関係が非常に複雑です。権利などデータにとって付随的なものと思われるかもしれませんが、実はこれがそうでもなくて、人文系データは権利そのものなんですね。どうやって社会的に混乱が起きないようにオープン化していくかについて熟慮することは我々に迫られている課題だと申し上げておきます。

所有者がいたからこそ今に伝わる資料がある

たとえば「何々家の文書をありがたく拝受しました。については画像にして公開してよろしいでしょうか」と何うと、大概の所有者は首を

シン・ニホン -AI×データ時代における 日本の現状と人材育成問題-

ヤフー株式会社
チーフストラテジーオフィサー
(CSO)
安宅 和人 氏



ワクワク感を形にできる人が国富を生む

産業革命並みの変革期にあることは間違いないと思っています。激変期に生きているという意識を持たなければいけない時に、10年先、20年先の時間軸で考えていることが危ないということをまず訴えたいと思います。

世界の事業価値のランキングMarket capで、上の方はみんなICT企業です。Apple、Google、Amazon、Facebook、これらの企業は圧倒的に利益よりも事業価値が大きい。これは何を意味しているのかという「世の中を変えている感」が富に繋がっているということです。国富を生む方程式が変化したわけで、こういうワクワク感を形にできる人を育てないと我々の未来はないわけです。

情報産業革命の第1フェーズがほぼ終わりつつあって、日本は大敗です。そこはもうしょうがないので、やがてやってくる第2、第3フェーズにこそ勝負をかけるべきだと思います。

確認できて教え方の見直しもできます。

また、来年度から全学部の学生がサーバーセキュリティの授業を受けないと卒業できないという形にしております。データサイエンスについても今後同じような形に持っていくことを考えております。

低年次は一般的なリテラシーとしてのデータサイエンスの観点が必要になってまいりますので、基幹教育院による全学部向け科目モジュール「数理・データサイエンス実践基幹教育」を入れ、高年次では専門分野別に特化したモジュールを入れ、さらにそれは大学院でのプログラムにつながっていく階層的な構造になっております。

オープンサイエンスへ向かう時代の潮流を受け止めて

研究データを公開し公正性や成果の再利用性を高めようとするオープンサイエンスの流れが来ています。公開するのは研究者の責任、保全する環境を作るのは所属する研究機関の責任、そして流通

横に振ります。なぜか。その文書に何が書いてあるかわからない、ひょっとすると自分たちの祖先の負の記録かもしれないからです。所有者にすれば、データの中身は本質的な問題なのです。また、ある人が絵画を持っています。その所在の情報なり絵そのものを公開して、それが時価何億円の価値があったとなると、これは税務上の問題にもなり得ます。

宗教上、習慣上、センシティブな土地で調査をする際には、情報の提供先は明かしてくれるなど強く言われることもあります。

さらにゲノム情報になるともっとセンシティブになってきます。

個人情報とはもかくとして、所有権の問題はもう反故にしたらいいのではというお考えもあると思います。ただ、1000年前の文書が残ってきた背景を考えますと、誰かがお家の宝として持ち続けてきたからこそ今に伝わっているわけで、あながち所有権はいらないとは言えないわけです。

この国は何度もスクラップ&ビルドでのし上がってきた。今度も立ち上げられる。まだやれるんじゃないかと思っています。

これまでとは似て非なるData professional人材を育成

そのためには、データの力を解き放った上で、見る力・決める力・伝える力を持った人を育てる必要があります。今生まれている変化にエキサイトして変えていこうとする人であり、統計だけの専門家というのではなく、統計的素養を持った上で情報科学的な知恵と技を課題解決に使う人。また、これまでのSIer的な単なるプログラマーを超えて、課題を俯瞰し、柔軟にビッグデータ処理を実験環境から本番環境まで実現できる人。こういう、これまでとは似て非なるData professional人材が必要になってきます。

データリテラシーとアントレプレナーシップという武器を持った若者、境界・応用域を含む専門家層とリーダー層、ICTエンジニアと

をサポートするのは図書館の責任という流れに整理されつつあります。ということは、大量に生まれるビッグデータを保全するしかけを大学が持つことが非常に重要になってくるわけです。

そこで我々が頼ろうとしているのが当機構でございまして、国立情報学研究所で、データの公開・管理・検索の基盤を整備してご提供いただくというお話を聞いております。

人類が経験したことのない時代、歴史に学べない時代に我々は突入しています。データサイエンスによる新しい時代をしっかりと受け止めるためには、教育や研究に全く新しい発想が必要になってまいりますし、データに基づく新しい科学哲学、倫理、社会の指導原理が求められます。科学や技術と社会をどう調和させるかという、非常に大きな課題を提起されているという認識でおります。

人文系データとは、そういうもの、つまり人間を扱っているわけですから、データのあり方、保存の仕方、データ化の進め方、管理の仕方、安全の配慮など十分に考えておく仕掛けが同時に必要だと思います。

さまざまな種類のデータをオープン化するには、さまざまな手続きがどうしても必要で、いろいろな歯止めが必要、また生じるであろうさまざまな混乱に対するいろいろな手当が必要なことをご理解いただきたいと思います。

ミドル・マネジメント層、それぞれの育成と再生が不可欠です。

千載一遇のチャンスを活かし世界の才能もかき集めた方がいい。

また、基盤となる思考、表現の武器としての国語の刷新を第一に考えるべきです。分析的、構造的に文章や話を理解し課題を洗い出す。論理的かつ建設的に物を考え、明確かつ力強く伝えるという“コミュニケーション”というべきものに刷新していかないと、データリテラシーが入っても動いていかないんですね。

こういう研究大学での人材開発に向けて、10兆円規模の国家的なendowmentを立ち上げなければまずいんじゃないかと思っています。

以上の実現に向けて、国家全体のリソースの最適化を検討すべきではないかということが、いま国関連で自分が接するあらゆる人に私が提言していることとございます。



大学共同利用機関法人

情報・システム研究機構

Research Organization of Information and Systems

発行: 大学共同利用機関法人 情報・システム研究機構 戦略企画本部 URAステーション