

# 東京科学大学（理工学系） 応用基礎レベルの取り組み紹介

特任准教授 奥村 圭司

特任教授 宮崎 慧

東京科学大学 データサイエンス・AI全学教育機構



- データサイエンス・AI全学教育プログラムのご紹介

- 応用基礎レベル

取組概要 / カリキュラム / 履修者の声 / 修了証

- 教材提供について

- 一部抜粋としても利用可能

- 一般公開している資料 (スライドPDF版, Jupyterノートブック)

- 教員限定での提供資料 (スライドPPT版, 課題)

- プログラミング課題の自動採点について

- 学生が解答をセルフチェックする場合の例

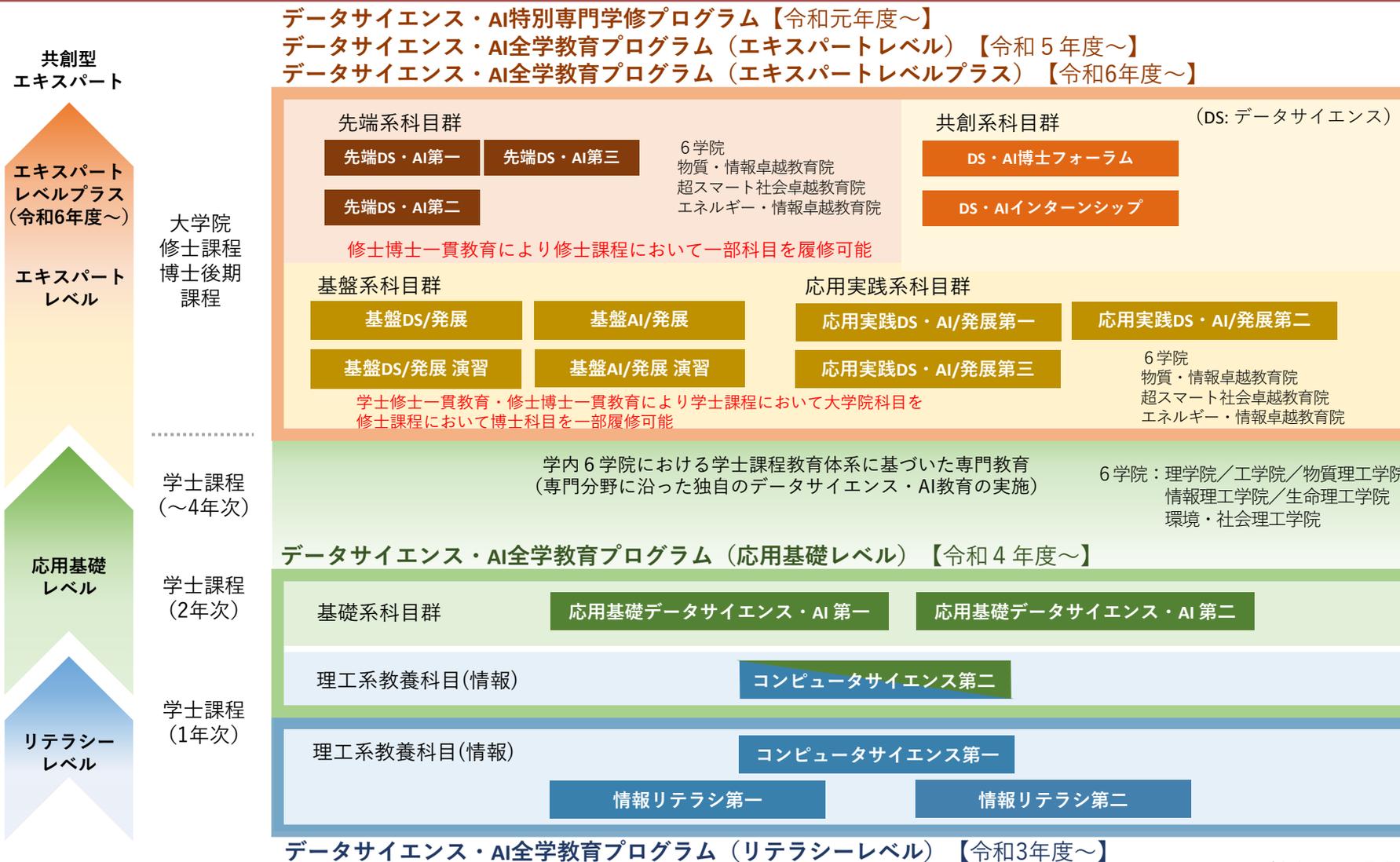
- 「応用基礎データサイエンス・AI第一／第二」 (2年次以上向け授業) について

- 科目の概要 / 講義・演習資料例 / 講義フロー・イメージ / 評価方法 / 講義実施の所感

# DS・AI全学教育プログラムのご紹介

# DS・AI全学教育プログラム(理工学系)の概要

- 理工系総合大学・学士修士一貫教育の特徴を活かし、大学院修士レベルの全学向けデータサイエンス・AI特別専門学修プログラムを創設 (令和元年度)
- リテラシーレベルからエキスパートレベルまで一貫した教育をめざし、リテラシーレベル教育プログラムを開始、博士後期課程科目を追加 (令和3年度)
- リテラシーレベルからエキスパートレベルへの橋渡しとなる応用基礎レベル教育プログラムを開始し、博士後期課程科目を充実 (令和4年度)



# 応用基礎レベル：取組概要（申請書より）

## 教育目標

### 共創型エキスパート人材育成 に向けた応用基礎レベル教育

将来的に学生が進むであろう専門分野に依らず**数理・データサイエンス(DS)・AI**を駆使して問題解決ができる**能力獲得**に向け、リテラシーレベルの学修を終えて基礎的素養を身につけた学生が、**各種手法の理論やプログラミング実践スキル**を含むより高度な学習項目を修得し、大学院課程におけるエキスパートレベルの学修への橋渡しとなる素養を獲得する。



## 実施体制

### データサイエンス・AI全学教育機構運営委員会

学部と大学院を統一した組織である6学院\*および卓越大学院プログラムに採択された3卓越教育院\*\*等から選出された委員により構成

全学教育プログラムの運営・実施、改善・進化、自己点検・評価を行う全学委員会

\* 理学院、工学院、物質理工学院、情報理工学院、生命理工学院、環境・社会理工学院  
\*\* 物質・情報卓越教育院、超スマート社会卓越教育院、エネルギー・情報卓越教育院

### 理工系教養科目（情報）実施委員会

理工系教養科目（情報）の授業を担当する教員により構成

データサイエンス・AI全学教育機構運営委員会と連携して授業の点検、評価、改善

## 全学教育プログラム（応用基礎レベル）科目構成

理工系教養科目（情報）1科目1単位 + 情報理工学院共通科目2科目2単位を取得

学士課程（1年次） 理工系教養系科目（情報）	学士課程（2年次） 情報理工学院共通科目（全学開講）
4Q コンピュータサイエンス第二 1単位	1Q 応用基礎データサイエンス・AI第一 1単位
全学教育プログラム（リテラシーレベル）科目構成 情報リテラシー第一・第二／コンピュータサイエンス第一／基礎データサイエンス・AI	2Q 応用基礎データサイエンス・AI第二 1単位

- 1Q: 前学期・第1クォーター、2Q: 前学期・第2クォーター、4Q: 後学期・第4クォーター
- 1年次4Qは複数クラス編成、講義0.5単位 + 演習0.5単位で合計1単位
- 1Q・2Qは講義1単位、演習教材を豊富に提供することにより実践スキル向上を促す
- 1年次必修の理工系教養科目（数学）（微分積分学・線形代数学）の履修を前提

## 科目の特徴

- 理工系総合大学の特徴と学士修士一貫教育システムを活かしたカリキュラム編成
- 座学一辺倒ではなく実例演習も重視し、基礎的なプログラミング実践スキルを修得
- 学院間の垣根を越えたクラス編成と系統的に深く理解しやすい授業の実施
- エキスパートレベル教育への橋渡しの役割

## 共創型エキスパート人材とは

- DS・AIを駆使できる理論的基盤を身につけ
- DS・AIで専門の境界を越えて多様な人々と交わり
- DS・AIの未来を担う若者を教えられるトップ人材

# 応用基礎レベル：カリキュラム

## ◆コンピュータサイエンス第二

データサイエンス、機械学習の基本を学ぶと共に  
より実践的なプログラミングを書けるようになる

第1回 計算とは何か (1)

チューリング機械

第2回 計算とは何か (2)

計算可能性理論と計算複雑性理論 (P・NP)

第3回 より実践的なPythonプログラム

NumPy, Matplotlib

第4回 データ分析基礎

相関、疑似相関、回帰

第5回 機械学習の基礎

教師あり学習 (サポートベクタ分類器)

第6回 コンピュータサイエンス、データサイエンス、AIに関するトピックス

教師あり学習 (汎化)

第7回 コンピュータサイエンス、データサイエンス、AIに関するトピックス

教師なし学習 (クラスタリング)

## ◆応用基礎データサイエンス・AI第一

データサイエンス・AIの基本理論、アルゴリズムを理解  
すると共に課題演習を通して実践的スキルを身に付ける

第1回 データサイエンス・AIの基礎と歴史

第2回 データエンジニアリング基礎

ビッグデータとデータエンジニアリング、データ表現、データベース

第3回 データサイエンス・AI利活用基礎1 (Python言語とライブラリ)

関数、クラス・メソッド、可視化

第4回 データサイエンス・AI利活用基礎2 (Python/pandasの利用)

オープンデータ、データクレンジング

第5回 データサイエンス・AI数学基礎

線形代数、確率、確率分布

第6回 データサイエンス基礎1

各種の離散・連続分布、チェビシェフの不等式、中心極限定理、乱数

第7回 データサイエンス基礎2

基本統計量、相関、単回帰・重回帰分析

## ◆応用基礎データサイエンス・AI第二

将来のエキスペートレベルの学修につながる知識  
や実践的スキルの修得を目指した学修内容を提供

第1回 数理統計の基礎1

標本と統計的推定、点推定、区間推定

第2回 数理統計の基礎2

仮説検定、母分散・母比率の推定と検定、適合度と独立性の検定

第3回 機械学習1

教師なし学習：クラスタリング、主成分分析

第4回 機械学習2

教師あり学習：各種の回帰、過学習と正則化、分類、性能評価

第5回 ニューラルネットワークと深層学習1

パーセプトロン、確率的勾配降下法、誤差逆伝播法

第6回 ニューラルネットワークと深層学習2

畳み込みNN、リカレントNN

第7回 機械学習・深層学習・AI応用

強化学習、生成モデル(VAE, GAN, 拡散モデル)、注意機構とトランスフォーマ

この2科目のみで  
申請書の様式1が  
求める範囲はカバー  
(教材提供も可能)

# 応用基礎レベル：履修者の声・修了証

## ●履修者の声

- 統計学は統計学、PythonはPython、データ分析はデータ分析、という別々の認識でしたが、これらを関連して理解できました。とても密度が高い授業だと思います。(環境・社会理工学院)
- Pythonのライブラリの使い方、データのクレンジングがかなり自分の思い通りにできるようになった。(情報理工学院)
- データサイエンスにおける処理のなかで数学の知見が大いに活用されていることを実感できた。(工学院)
- 資料、丁寧なご講義、オンデマンド動画があって、繰り返し学習ができますので、本当に有り難いです。一度きりの対面授業では厳しいです。オンデマンドで振り返りができるので、勉強に集中ができます。(環境・社会理工学院)

## ●修了者全員に「デジタル修了証」(オープンバッジ)を付与



# 教材提供について

# 提供資料とライセンス

授業スライド・演習問題・プログラミング  
課題(自動採点对応)などを公開・提供中

<https://www.dsai.titech.ac.jp>

Collaboration > 教材提供

## 一般公開

- スライドPDF版  
CCライセンス(CC-BY-NC-SA 4.0)

- Jupyterノートブック

CCライセンス+MITライセンス

## 教員限定提供

- スライドPPT版・課題など

教材利用規約への同意により提供



# 例) 一般公開・スライド

[www.dsai.titech.ac.jp](http://www.dsai.titech.ac.jp)

DS&AI

教材提供 (応用基礎レベル)

応用基礎データサイエンス・AI第一

[2024年度 第1クォーター版]

全学院対象・学士課程2年次向けの授業です。

Details

応用基礎データサイエンス・AI第二

## 応用基礎データサイエンス・AI第一

[2024年度 第1クォーター版]

全学院対象・学士課程2年次向けの授業です。

シラバス

教材提供

以下の各項目より資料 (主にスライドPDF版、リンク先でJupyterノートブック) を入手できます。  
本ウェブサイトから直接ダウンロードできる教材は、クリエイ

第1回 データサイエンス・AIの基礎と歴史

第2回 データエンジニアリング基礎

- 演習資料 (第2回) (Jupyterノートブック形式; Google Colabへのリンク)

【教員限定】教材 (PPT版) や課題等の提供について (Googleフォームへのリンク)

## スライド (PDF版)

### 教材の利用条件

- 本教材 (PDF版) は「クリエイティブ・コモンズ・ライセンス 表示-非営利-継承 4.0 国際 (CC BY-NC-SA 4.0)」のもと提供しています。
  - ライセンスの内容を知りたい方は以下よりご確認ください。  
<https://creativecommons.org/licenses/by-nc-sa/4.0/deed.ja>
  - クレジットの表示: 資料のいずれかのページに、提供元として下記を明記してください。
    - 東京工業大学 データサイエンス・AI全学教育機構
- 【教員限定】本教材をご利用の方へアンケートのお願い
  - ご回答はこちら: <https://forms.gle/8s7QcaWKfNZm8uBWA>
- 【教員限定】本教材 (PPT版) や課題等について
  - 日本国内の教育機関に在籍する教員に限定して提供しています。
  - お申込みはこちら: <https://forms.gle/8PAhW8taARjAsXay5>

DS&AI  
Center of Data Science  
and Artificial Intelligence

## 応用基礎データサイエンス・AI 第一

### 第2回 データエンジニアリング基礎

2024.04.17

データサイエンス・AI全学教育機構

### 講義の概要

1. ビッグデータとデータエンジニアリング
  - ・ ビッグデータとは
  - ・ ICT (情報通信技術) の進展とビッグデータ
  - ・ ビッグデータの活用事例

一部抜粋でも  
利用可能

# 例) 一般公開・ノートブック

[www.dsai.titech.ac.jp](http://www.dsai.titech.ac.jp)



応用基礎データサイエンス・AI第一

[2024年度 第1クォーター版]

全学院対象・学士課程2年次向けの授業です。

Details



応用基礎データサイエンス・AI第二

## 応用基礎データサイエンス・AI第一

[2024年度 第1クォーター版]

全学院対象・学士課程2年次向けの授業です。

シラバス



教材提供



以下の各項目より資料（主にスライドPDF版、リンク先でJupyterノートブック）を入手できます。  
本ウェブサイトから直接ダウンロードできる教材は、クリエイ

第1回 データサイエンス・AIの基礎と歴史



第2回 データエンジニアリング基礎



- 演習資料 (第2回) (Jupyterノートブック形式; Google Colabへのリンク)



【教員限定】教材 (PPT版) や課題等の提供について (Googleフォームへのリンク)



ノートブック

y24q1-BADSAI1-2\_Practice.ipynb  
ファイル 編集 表示 挿入 ランタイム ツ

+ コード + テキスト ドライブにコピー

接続

応用基礎データサイエンス・AI 第一

第2回 データエンジニアリング基礎

演習資料

東京工業大学 データサイエンス・AI全学教育機構

2024.04.17

↑ ↓ ↻ ↷ 🗑️ ⋮

▼ データ表現

Python 言語環境において

1. 数値
2. 文字 (テキスト)
3. 音声 (オーディオ信号)
4. 画像

の各データがどのように取り扱われているかを、実例を通して紹介する。

▼ 1. 数値のバイナリ表現

1.1 エンディアン (endianess)

コンピュータの主記憶 (メインメモリ) や外部記憶装置 (HDD、SSD、USBフラッシュメモリなど) のデータの格納は、通常バイト (8ビット) 単位を基本として行われる。この際、2バイト以上の複数バイトにより表現される一つの数値がメモリ上に配置されるか (バイト並び順) は、使用するコンピュータシステムによって異なる。例えば、4バイト (32ビット) で表現される整数型のデータが16進16909060、2進表記で0b0001001000110100) であるとする。

- ビッグエンディアン (big-endian)  
メモリ上に [01][02][03][04] の順番に配置される
- リトルエンディアン (little-endian)

一部抜粋でも  
利用可能

# 例) 教員限定資料の申し込み

[www.dsai.titech.ac.jp](http://www.dsai.titech.ac.jp)



応用基礎データサイエンス・AI第一

[2024年度 第1クォーター版]

全学院対象・学士課程2年次向けの授業です。

Details →



応用基礎データサイエンス・AI第二

## 応用基礎データサイエンス・AI第一

[2024年度 第1クォーター版]

全学院対象・学士課程2年次向けの授業です。

シラバス



教材提供



以下の各項目より資料（主にスライドPDF版、リンク先でJupyterノートブック）を入手できます。  
本ウェブサイトから直接ダウンロードできる教材は、クリエイ

第1回 データサイエンス・AIの基礎と歴史



第2回 データエンジニアリング基礎



- 演習資料 (第2回) (Jupyterノートブック形式; Google Colabへのリンク)



【教員限定】教材 (PPT版) や課題等の提供について  
(Google フォームへのリンク)



【教員限定】教

申込フォーム

Google にログインすると作業内容を保存できます。詳細

\* 必須の質問です

希望の教材について

リテラシーレベル

全授業回 (第1~7回) の資料を提供します。特定回のみご希望の場合は、「その他」で希望の回・教材を記入してください。

「基礎データサイエンス・AI」教材 (PPT版)

「基礎データサイエンス・AI」課題 (正誤問題等)

「基礎データサイエンス・AI」課題 (プログラミング・自動採点対応)

「教材利用規約」への同意が必要

応用基礎データサイエンス・AI第一で希望の回・教材を記入してください。

「応用基礎データサイエンス・AI第一」教材 (PPT版)

「応用基礎データサイエンス・AI第一」課題・レポート問題

「応用基礎データサイエンス・AI第二」教材 (PPT版)

「応用基礎データサイエンス・AI第二」課題・レポート問題

その他: \_\_\_\_\_

教材の利用目的について教えてください。\*

回答を入力

戻る

次へ

フォームをクリア

- 演習問題・課題等を教員限定で提供する場合の規約

[www.dsai.titech.ac.jp/terms-of-use/](http://www.dsai.titech.ac.jp/terms-of-use/)

## Terms of Use

教材利用規約

本教材利用規約は、国立大学法人東京工業大学（以下、「本学」といいます。）が、データサイエンス・AI教育の普及及び促進のために、日本国内の学校その他の教育機関に在籍する常勤・非常勤の教員（以下、「利用者」といいます。）に対して提供する教材の利用条件について定めるものです。利用者は、本教材利用規約に同意の上、本教材利用規約に定めに従って、教材を利用することができます。

### 3. 複製、配布、公開の範囲

利用者は、所属教育機関に在籍する学生及び職員に対してのみ、教材の複製、提供、配布、公開をすることができます。

### 5. 改変、翻訳、編集の許諾

利用者は、利用目的に則った利用のために、教材を改変、翻訳又は編集をすることができます。ただし、これらの改変、翻訳又は編集は、原著者の表示を維持し、著作権を尊重する形で行われなければなりません。また、改変、翻訳又は編集後の二次的著作物は、本教材利用規約の定めに従って利用しなければなりません。

### 8. 利用者の責任

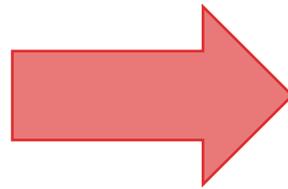
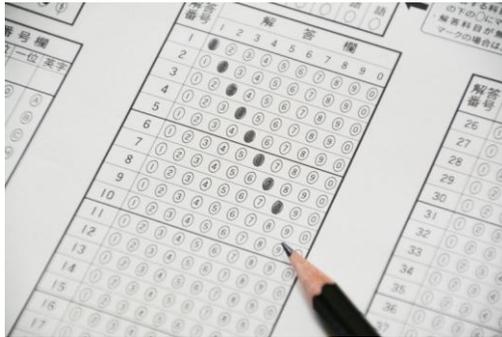
1. 利用者は、本教材利用規約に従って教材を利用する責任を負います。利用者が本教材利用規約に違反した場合には、本学の指示に従って、教材の利用停止、返還、消去等に応じなければなりません。
2. 利用者は、教材の修正、訂正、更新等に伴い、本学が、教材の利用停止、差し替え、消去等の要請を行ったときは、直ちに、かかる要請に応じなければなりません。

# プログラミング課題の自動採点について

# プログラム課題の自動採点とは？①

最も身近なイメージは、「マークシート方式」

マークシート方式  
(選択式) で回答



大量の解答を  
短時間で  
正確に採点できる 💡

番号を読み取り、膨大なデータを集計できるため、  
大学入学共通テストや資格試験など、多くの場面で使用されている。

# プログラム課題の自動採点とは？②

先の考え方の進化系として、プログラミング（ex. Pythonのコード）課題の採点に応用

Pythonのコードの穴あき部分に  
解答（コードの自由記述）を記入



実行キーひとつで

即時に

正確に採点ができる 💡

(例) `def tashizan(a,b):`

`return`  ←穴あき部分に解答を記入 (正答は、`a+b`)

- 教員の採点負担を軽減する。
- 学生の自習に活用することで、セルフチェック機能としても有効。
- 自分の書いたコードをチェックしたり、正解へのヒントを得られる。

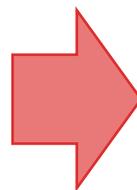
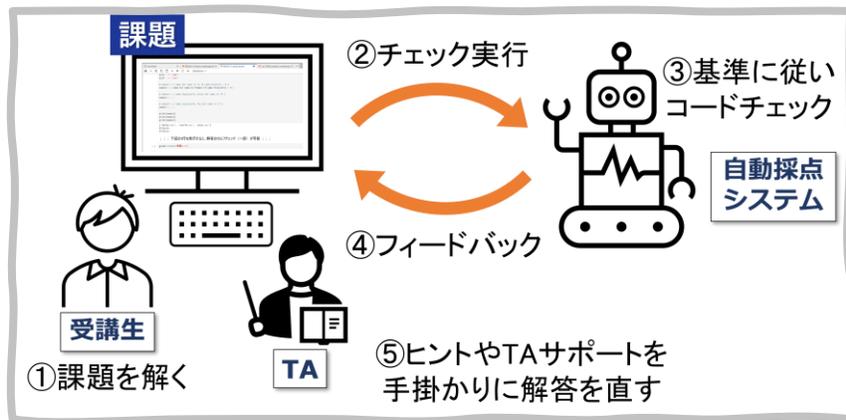


# 日本国内の大学への普及

自動採点システムを本格的に開発・運用し、  
世界でも先進的な取り組みを実施している  
カリフォルニア大学バークレー校（UC Berkeley）と連携



「基礎データサイエンス・AI」など理工学系の授業に  
自動採点システム **Otter-Grader** を導入  
Jupyterノートブックによる課題へ適用



考え方自体はシンプル

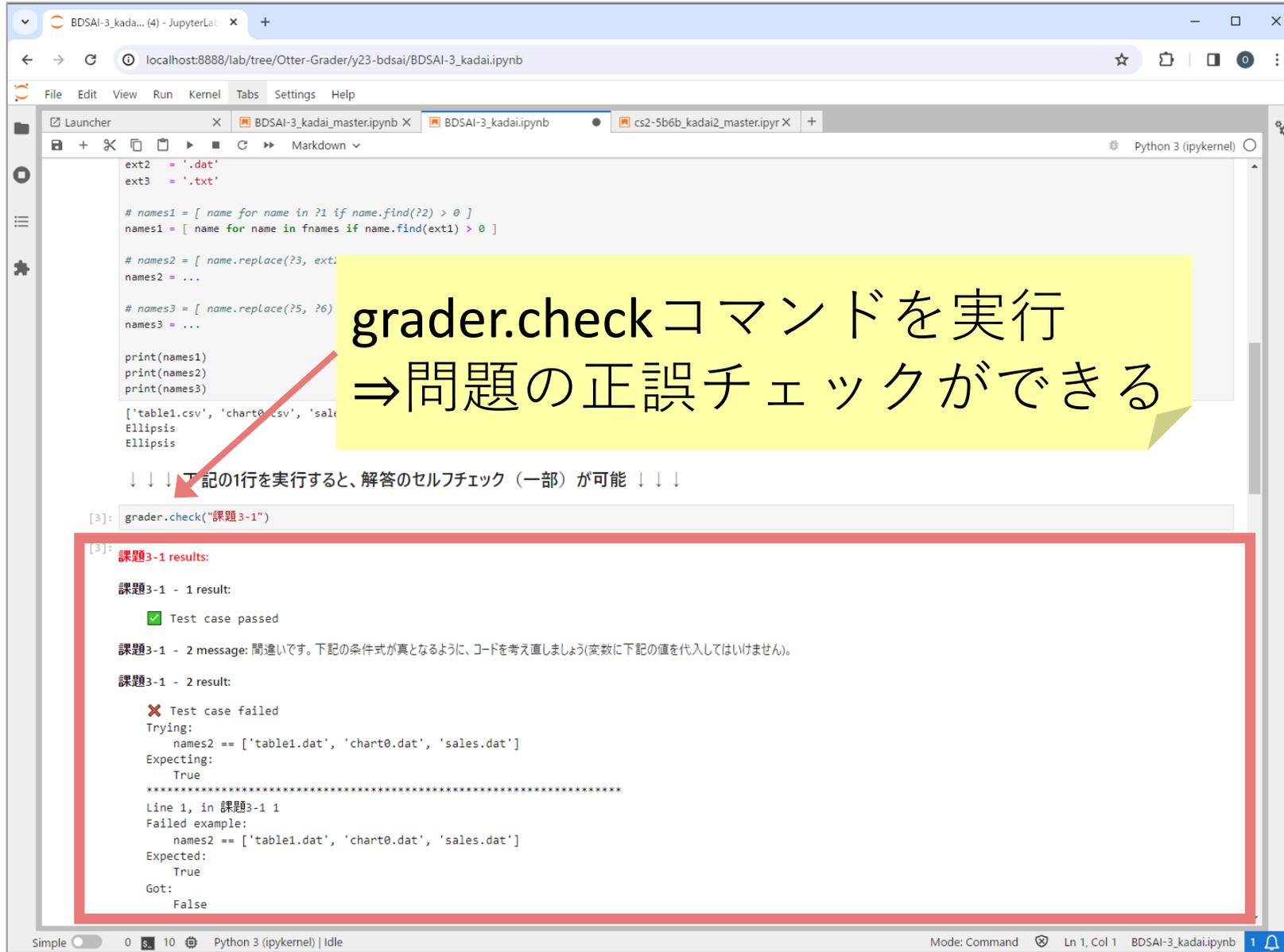
しかし、ソフトウェアは  
開発よりもメンテナンス  
が大変

## 日本国内の様々な大学への普及を推進

- 「数理・データサイエンス・AI教育強化拠点コンソーシアム」の拠点校として、他大学への普及を推進
- 東京科学大学の発足に伴い、同大学医歯学系への普及を推進

**提供教材**の中に  
導入マニュアルや対応課題あり

# 例) 学生が解答をセルフチェックする場合



The screenshot shows a JupyterLab window with a Python 3 kernel. The code in the editor includes file extensions and list comprehensions for filtering and replacing text. A yellow callout box highlights the `grader.check` command. Below the code, the execution output shows a successful test case and a failed one with a detailed error message.

```
ext2 = '.dat'
ext3 = '.txt'

# names1 = [ name for name in ?1 if name.find(?2) > 0 ]
names1 = [ name for name in fnames if name.find(ext1) > 0 ]

# names2 = [ name.replace(?3, ext:
names2 = ...

# names3 = [ name.replace(?5, ?6)
names3 = ...

print(names1)
print(names2)
print(names3)

['table1.csv', 'chart0.csv', 'sales.
Ellipsis
Ellipsis

↓↓↓ 下記の1行を実行すると、解答のセルフチェック（一部）が可能 ↓↓↓

[3]: grader.check("課題3-1")

[3]: 課題3-1 results:
課題3-1 - 1 result:
  ✓ Test case passed
課題3-1 - 2 message: 間違いです。下記の条件式が真となるように、コードを考え直しましょう(変数に下記の値を代入してはいけません)。
課題3-1 - 2 result:
  ✗ Test case failed
  Trying:
  names2 == ['table1.dat', 'chart0.dat', 'sales.dat']
  Expecting:
  True
  .....
  Line 1, in 課題3-1 1
  Failed example:
  names2 == ['table1.dat', 'chart0.dat', 'sales.dat']
  Expected:
  True
  Got:
  False
```

grader.check コマンドを実行  
⇒ 問題の正誤チェックができる



TAによるサポート  
との二本柱

モデルカリキュラム中心科目  
「応用基礎データサイエンス・AI第一／第二」  
(2年次以上向け授業)  
について

- 科目名

- 応用基礎データサイエンス・AI 第一（以下第一と略記）

- 第1Q（計7回）・各回100分・単位数：1（講義：0.5 演習：0.5）

- 応用基礎データサイエンス・AI 第二（以下第二と略記）

- 第2Q（計7回）・各回100分・単位数：1（講義：0.5 演習：0.5）

- 到達目標

- データサイエンス、データエンジニアリング、AI分野の基礎理論や手法を学び、研究に応用できる基礎力を養う
- Python環境でのプログラミングスキル習得とライブラリ活用
- アルゴリズム実装演習を行い、事例を通じて実践的な課題解決力を向上

	授業計画	課題
第1回	データサイエンス・AIの基礎と歴史	データサイエンス・AIの基礎、歴史、役割
第2回	データエンジニアリング基礎	ビッグデータとデータエンジニアリング、データ表現、データベース
第3回	データサイエンス・AI利活用基礎1 (Python言語とライブラリ)	関数、クラス・メソッド、可視化
第4回	データサイエンス・AI利活用基礎2 (Python/pandasの利用)	オープンデータ、データクレンジング
第5回	データサイエンス・AI数学基礎	線形代数、確率、確率分布
第6回	データサイエンス基礎1	各種の離散・連続分布、チェビシェフの不等式、中心極限定理、乱数
第7回	データサイエンス基礎2	基本統計量、相関、単回帰・重回帰分析

	授業計画	課題
第1回	数理統計の基礎1	標本と統計的推定、点推定、区間推定
第2回	数理統計の基礎2	仮説検定、母分散・母比率の推定と検定、適合度と独立性の検定
第3回	機械学習1	教師なし学習：クラスタリング、主成分分析
第4回	機械学習2	教師あり学習：各種の回帰、過学習と正則化、分類、性能評価
第5回	ニューラルネットワークと深層学習1	パーセプトロン、確率的勾配降下法、誤差逆伝播法
第6回	ニューラルネットワークと深層学習2	畳み込みNN、リカレントNN
第7回	機械学習・深層学習・AI応用	強化学習、生成モデル(VAE, GAN, 拡散モデル)、注意機構とトランスフォーマ

## ビッグデータとは

- 大量のデジタルデータ

- 人間の能力では処理しきれない膨大な量のデータ

### SNS

YouTube, Twitter, FaceBook, Instagram, Email, TikTok, ...

### ニュース・デジタルアーカイブ

テキスト, 画像, 音声, 動画, ...

### POS (Point of Sales) データ

コンビニやスーパーの購買行動, Eコマース取引履歴情報, ...

### 人流データ

携帯電話位置情報, Wi-Fi接続情報, 交通系ICカード情報, ...

### 観測データ

気象観測データ, 自然現象観測データ, リモートセンシングデータ,  
各種実験データ, ...

毎分当たりの流通データ量\*  
YouTube投稿：500時間分の動画  
Twitter共有：347,200ツイート  
FaceBook共有：1,700,000件  
Instagram共有：66,000枚の写真  
Email送信：231,400,000件  
日本における情報通信量\*\*  
総務省情報通信統計データベース

\* Data Never Sleeps, <https://www.domo.com/data-never-sleeps/>  
\*\* <https://www.soumu.go.jp/johotsusintokei/index.html>

## ビッグデータ

- 大量のデータ
  - 人間の行動
  - SNS
  - YouTube
  - ニュース
  - テキスト
  - POS (Point of Sale)
  - コンピュータ
  - 人流データ
  - 携帯電話
  - 観測データ
  - 気象観測
  - 各種実

## データの表現：画像（色の表現）

- RGB
  - 三原色：赤 (Red ■), 緑 (Green ■), 青 (Blue ■)

三原色の組合せで色情報を表現

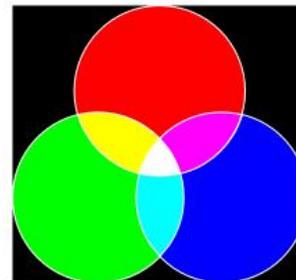
色の表現としてRGB以外にもHSV\*, CMYK\*\*などが使われている

- 色情報の量子化
  - 24ビットカラー

R/G/B 各成分の明度を256段階 (8ビット符号なし整数) で表し、それらの組合せで異なる  $256 \times 256 \times 256 = 16,777,216$  色を表現 (例)

	R:G:B=0:0:0		R:G:B=128:128:128
	R:G:B=255:128:0		R:G:B=0:255:128
	R:G:B=128:0:255		R:G:B=128:128:255
	R:G:B=192:128:64		R:G:B=255:255:255

\* 色相 (Hue)/彩度 (Saturation)/明度 (Value); \*\* Cyan/Magenta/Yellow/Key plate



カラー画像

## ビッグデータ

- 大量のデータ
  - 人間の行動
  - SNS
  - YouTube
  - ニュース
  - テキスト
  - POS (P)
  - コンピ
  - 人流データ
  - 携帯電
  - 観測データ
  - 気象観
  - 各種実

Tokyo Tech, 2024

## データの表現

- RGB
  - 三原色
  - 三原色の
- 色情報の量
  - 24ビット
  - R/G/B
  - それらの
  - 例)



Tokyo Tech, 2024

## データベースとは

- データベース (database, DB)
  - 組織化(構造化)されたデータの集まり
  - 得たい情報の検索や蓄積が容易になるように整理されたデータの集合
  - 操作：検索, 追加, 更新, 削除
  - 例) 住民基本台帳, 銀行口座情報, SNS会員登録情報, 顧客情報
  - 音楽データベース, ルート検索, 学務管理, 電子カルテ, ...
- データモデル
  - データの格納の仕方を規定
  - 階層型データベース
  - ネットワーク型データベース
  - カード型データベース
  - リレーショナルデータベース (Relational DB, RDB)
  - 関係モデル (Relational Model, RM) に基づくDB [ Edgar F. Codd, 1970]
  - 最も広く使われている

Tokyo Tech, 2024

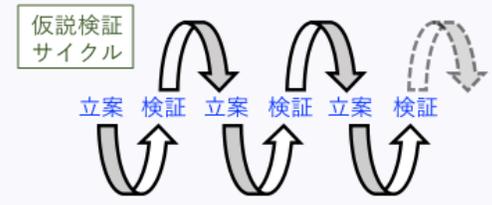
46

## データ分析と仮説検証サイクル

- データ分析の目的
  - データの利活用
    - データを分析し、有用な情報を抽出
    - 得られた発見・知見を問題解決や付加価値創出に利用
- 仮説検証サイクル
  - 仮説の立案とその検証の繰り返し
    1. 仮説の立案 (発見)
    2. 仮説の検証
- 仮説の発見
  - 調査設計とデータ収集
  - 良い仮説を発見するための手段
    - データ要約・可視化, 特徴量抽出, クラスタリング, 次元削減, etc.
- 仮説の検証
  - 仮説が正しいかどうかを判断するための手段
    - 相関分析, 回帰分析, 仮説検定, etc.

仮説の発見：例えば統計データやビッグデータの中から有用な傾向や規則性を見出す

- 仮説検証の結果、仮説が正しいと言えない場合は、新たな仮説を立て、その検証を行う
- この反復（仮説検証サイクル）を確からしい結論が得られるまで繰り返す
- 立てた仮説の検証のために、適切なデータを新たに収集する必要がある場合もある



## データ分析

- データ分析
  - データ
    - データ
    - 得られた
  - 仮説検定
    - 仮説の種類
      - 1. 仮説
      - 2. 仮説
  - 仮説の棄却
    - 調査結果
    - 良い仮説
      - データ
      - クラスタ
  - 仮説の検定
    - 仮説の種類
      - 相関

Tokyo Tech, 2024

## 量的データの可視化

- 基本的な可視化手法
  - ヒストグラム (度数分布図, histogram)
    - 度数分布から作成する棒/柱状グラフ
      - 適切な階級幅を選択することが重要
  - 円グラフ (pie chart, circle chart)
    - 全体に対する割合を扇型の中心角で表現
  - 箱ひげ図 (box plot)
    - データの分布の代表値を箱とひげ (直線) で表現
      - ヒストグラムの簡易版としての分布比較
  - 散布図 (scatter plot)
    - 量的変数同士の関係の可視化
      - 横軸・縦軸がそれぞれ異なる量的変数の値に相当
      - 各標本のそれぞれの変数の値を2次元平面の1点として表現
  - ヒートマップ (heat map)
    - 二次元データの値の大小を色により表現

Tokyo Tech, 2024

6

## データ分析

- データ分析
  - データ
    - ・ データ
    - ・ 得られた
  - 仮説検定
    - 仮説の種類
      - 1. 仮説
      - 2. 仮説
  - 仮説の検証
    - 調査設計
    - 良い仮説
      - ・ データ
      - ・ クラスタ
  - 仮説の検証
    - 仮説の種類
      - ・ 相関

Tokyo Tech, 2024

## 量的データ

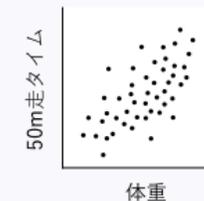
- 基本的な
  - ヒストグラム
    - 度数分布
      - ・ 適切な
  - 円グラフ
    - 全体は
  - 箱ひげ図
    - データ
      - ・ ヒストグラム
  - 散布図
    - 量的変数
      - ・ 横軸
      - ・ 各標本
  - ヒートマップ
    - 二次元

Tokyo Tech, 2024

## データ分析の基礎知識：相関関係と因果関係

- 因果関係 (因果性) (causality)
  - 「原因・要因」と「結果」の関係
  - 因果関係が成り立つ条件
    - ・ 影響の単一方向性・時間関係 → 原因が先で結果が後
    - ・ 共変関係 → 一方の変化と連動して他方も変化 = 相関関係
    - ・ 原因と結果に共通の要因がない
  - 相関関係
    - 相関関係は因果関係を包含しない
      - ・ 相関関係があっても因果関係があるとは限らない
    - 影響の方向性はわからない
      - ・ 二つの事象  $X, Y$ :  $X \rightarrow Y, X \leftarrow Y, \text{ or } X \leftrightarrow Y$
    - 疑似相関 (spurious correlation/relationship)
      - ・ 二つの事象に共通の要因  $Z$  があり因果関係があるように見える:  $Z \rightarrow X, Z \rightarrow Y \Rightarrow X \rightarrow Y$
    - 相関係数  $\neq 0$  の場合
      - ・ 線形関係にはないが無関係とは限らない

疑似相関例：小学生の計測結果 (仮想データ)



明らかに誤った因果関係の解釈：  
「体重が重いほど50m走が速くなる」

共通の要因は交絡因子 (confounding factor) と呼ばれ、その変数は交絡変数または第三の変数と呼ばれる

Tokyo Tech, 2024

17

## K-means クラスタリング：手法の概要

- 非階層的クラスタリングの代表的手法

- クラスタリングの問題設定

- データ： $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^L$
    - データ間の距離： $d(\mathbf{x}_i, \mathbf{x}_j) \stackrel{\text{def}}{=} \|\mathbf{x}_i - \mathbf{x}_j\|^2$
    - クラスタ数  $K$  を事前に決めてクラスタリング

- K-means アルゴリズム

- 初期化

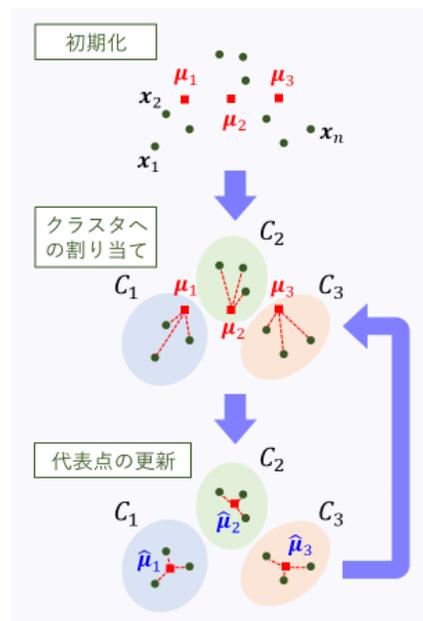
- クラスタ  $C_1, C_2, \dots, C_K$  の代表点の初期値を選ぶ
      - $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K \in \mathbb{R}^L$

- 繰り返し

- 各データを最も近い代表点を持つクラスタに割り当てる
      - $C_k = \{\mathbf{x}_i : k = \arg \min_j d(\mathbf{x}_i, \boldsymbol{\mu}_j)\}, i = 1, \dots, n, j = 1, \dots, K$
      - クラスタの代表点を更新する

$$\boldsymbol{\mu}_k = \arg \min_{\mathbf{y}} \sum_{\mathbf{x}_i \in C_k} d(\mathbf{x}_i, \mathbf{y}), \quad k = 1, \dots, K$$

- 代表点に変化しなくなったら繰り返しを終了



## K-means

- 非階層的
  - クラス
  - データ
  - データ
  - クラス
- K-mean
  - 初期化
  - クラス
  - $\mu_1, \dots, \mu_k$
  - 繰り返し
  - 各データ
  - $C_k = \{x \in X \mid r_k(x) > r_{k'}(x)\}$
  - クラス
  - $\mu_k = \frac{1}{|C_k|} \sum_{x \in C_k} x$
  - 代表

## 混合ガウスモデル：モデルパラメータ

### 混合正規分布

#### 確率密度関数

- $k$  番目の混合成分の正規分布  $f_k(\mathbf{x}|\theta_k)$

$$f_k(\mathbf{x}|\theta_k) \stackrel{\text{def}}{=} N(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \theta_k = \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$$

- 混合重み  $c_k, c_1 + c_2 + \dots + c_K = 1$

$$f(\mathbf{x}) = \sum_{k=1}^K c_k N(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \sum_{k=1}^K c_k f_k(\mathbf{x}|\theta_k)$$

#### 負担率 (responsibility)

- データ  $\mathbf{x}$  が観測されたとき、 $\mathbf{x}$  が  $k$  番目の混合成分から生成された確率 (事後確率)

$$r_k = \Pr(C_k | \mathbf{x}) = \frac{\Pr(C_k) \Pr(\mathbf{x} | C_k)}{\Pr(\mathbf{x})} = \frac{c_k f_k(\mathbf{x} | \theta_k)}{f(\mathbf{x})} = \frac{c_k f_k(\mathbf{x} | \theta_k)}{\sum_{k=1}^K c_k f_k(\mathbf{x} | \theta_k)}$$

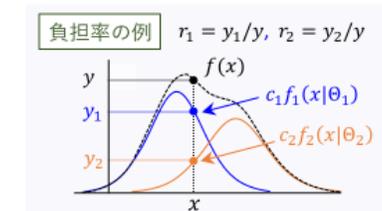
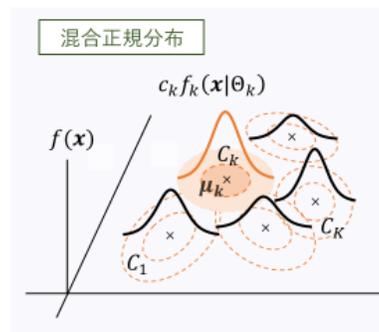
← ベイズの定理より

### GMMのモデルパラメータ

#### 混合正規分布のパラメータ

- 混合重み, 正規分布の平均と分散共分散行列の組の集合

$$\boldsymbol{\Theta} \stackrel{\text{def}}{=} \{\{c_k\}, \{\theta_k\}\} = \{c_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$$



## K-means

- 非階層的
  - クラス
  - データ
  - データ
  - クラス
- K-mean
  - 初期化
  - クラス
  - $\mu_1, \dots, \mu_k$
  - 繰り返し
  - 各データ
  - $C_k = \{x_i \mid x_i \text{ が } \mu_k \text{ に最も近い}\}$
  - クラス
  - $\mu_k = \frac{1}{|C_k|} \sum_{x_i \in C_k} x_i$
  - 代表

Tokyo Tech, 2024

## 混合ガウス

- 混合正規分布
  - 確率密度関数
  - $k$  成分
  - $f_k(x)$
  - 混合
  - $f(x)$
  - 負担率
  - データ
  - $r_k$
- GMMの
  - 混合正規分布
  - 混合
  - $\Theta = \{\mu_k, \Sigma_k\}$

Tokyo Tech, 2024

## 主成分分析の考え方

### データの射影

#### 射影元

- $L$ 次元データ:  $x_1, x_2, \dots, x_n \in \mathbb{R}^L$
- 正規直交基底:  $u_1, u_2, \dots, u_L \in \mathbb{R}^L, \langle u_k, u_l \rangle = \begin{cases} 1, & k = l \\ 0, & k \neq l \end{cases}$

#### 射影先

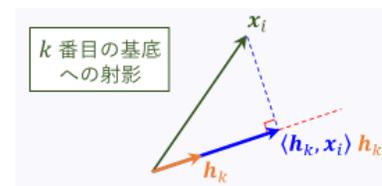
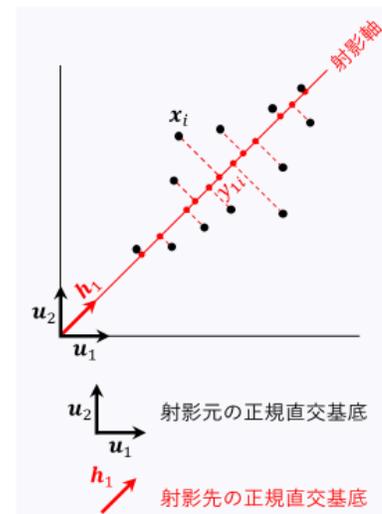
- $M (\leq L)$ 次元データ:  $y_1, y_2, \dots, y_n \in \mathbb{R}^M$
- 正規直交基底:  $h_1, h_2, \dots, h_M \in \mathbb{R}^L, \langle h_k, h_l \rangle = \begin{cases} 1, & k = l \\ 0, & k \neq l \end{cases}$

#### 射影データ

$$y_i = [h_1, h_2, \dots, h_M]^T x_i = \begin{bmatrix} h_1^T x_i \\ h_2^T x_i \\ \vdots \\ h_M^T x_i \end{bmatrix} = \begin{bmatrix} y_{1i} \\ y_{2i} \\ \vdots \\ y_{Mi} \end{bmatrix}$$

$$y_{ki} = \langle h_k, x_i \rangle = h_k^T x_i$$

Tokyo Tech, 2024



33

各手法の理論的な背景まで

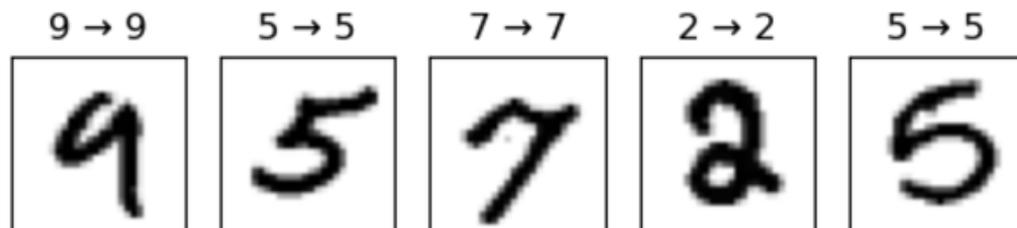
# 演習資料例：第二 ワーク・深層学習

## 第5・6回 ニューラルネット

テストデータの中からランダムに入力画像を5サンプル選び、その正解ラベルと認識結果を表示した結果を以下に示す。

```
import random

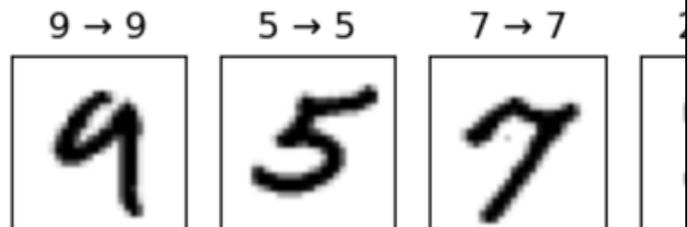
fig, axs = plt.subplots(1,5,figsize=(6,2))
idx = random.sample(range(X_test.shape[0]),5)
for i in range(5):
    image, label = X_test[idx[i]].reshape(28,28), np.argmax(y_test[idx[i]])
    result = np.argmax(z[idx[i]])
    # Test画像と認識結果表示
    axs[i].imshow(image, cmap='gray_r')
    axs[i].set_title(f'{label} → {result}')
    axs[i].get_xaxis().set_visible(False)
    axs[i].get_yaxis().set_visible(False)
plt.show()
```



テストデータの中からランダムに入力画像を5サンプル選び、その正解ラベルと認識結果を表示した結果を以下に示す。

```
import random

fig, axs = plt.subplots(1,5,figsize=(6,2))
idx = random.sample(range(X_test.shape[0]),5)
for i in range(5):
    image, label = X_test[idx[i]].reshape(28,28)
    result = np.argmax(z[idx[i]])
    # Test画像と認識結果表示
    axs[i].imshow(image, cmap='gray_r')
    axs[i].set_title(f'{label} → {result}')
    axs[i].get_xaxis().set_visible(False)
    axs[i].get_yaxis().set_visible(False)
plt.show()
```



```
import random
import numpy as np
import matplotlib.pyplot as plt

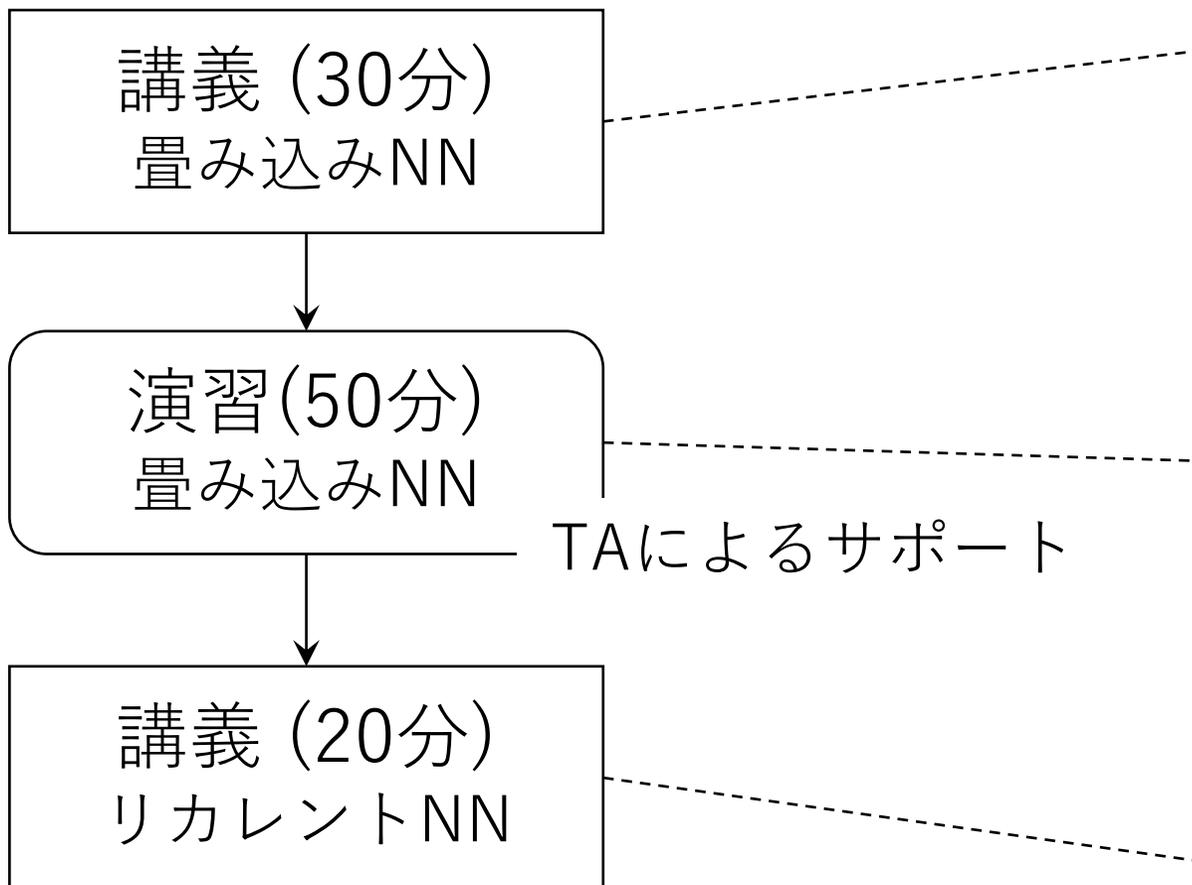
# テストデータからランダムに5個選んで画像とカテゴリを表示
classes = ('airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck')
fig, axs = plt.subplots(1,5,figsize=(6,2))
idx = random.sample(range(data_test.data.shape[0]),5)
for i in range(5):
    img, lbl = data_test[idx[i]]
    img_np = img.numpy().copy() # テンソルからndarrayへ変換
    img_np = np.transpose(img_np, (1,2,0)) # [ch, row, col] -> [row, col, ch]
    imgc = (img_np + 1) / 2 # 画素値の範囲を[0,1]に戻す
    axs[i].imshow(imgc)
    axs[i].set_title(f'{classes[lbl]}')
    axs[i].axis('off')
plt.show()
```



Python notebookを用いて講義で説明した内容の演習

# 講義フロー・イメージ

- zoom リアルタイム配信かつ演習室を用意

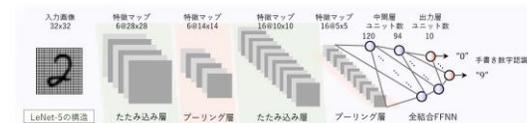


- 講義動画を後日受講生に配布

## たたみ込みNN (CNN)



- たたみ込み層とプーリング層を持つ深層NN
  - CNN: LeNetとその進化形  
LeNet-5 (1998); AlexNet (2012); VGGNet (2014); ResNet (2016)
  - CNNの基本構造
    - LeNet-5
      - たたみ込み層 (convolution layer)  
局所領域 (受容野) のフィルタリング処理に相当
      - プーリング層 (pooling layer)  
局所領域の特徴量の要約とダウンサンプリング処理に相当



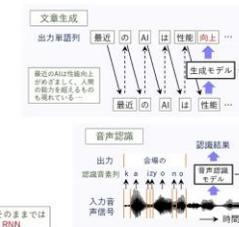
Tokyo Tech, 2024



## 系列データとニューラルネットワーク



- ニューラルネットワークにより解決したい問題の例
  - 回帰問題
    - 自然言語処理
      - 文章生成, 文章要約, 言語翻訳
    - 音声処理
      - 音声合成, 音声強調 (雑音除去, 音源分離)
  - 分類 (識別) 問題
    - 自然言語処理
      - 文章分類 (トピックス, 評定), 感情分析
    - 音声処理
      - 音声認識, 話者認識, 感情認識
  - 静的パターン (画像) との差異
    - 系列長が一定ではない  
⇒ 不定長入力への対応
    - 逐次的に特徴が変化する  
⇒ 系列データ長より短い単位 (時間) における逐次的な対応



Tokyo Tech, 2024

21

- 講義後に学内の学修管理システムを利用し、各回で多肢選択問題の小テストを実施
- または Pythonを用いた計算問題を回答しJupyter notebookファイルを提出
  - 上記2種の課題については自動採点を利用
- 期末課題：講義で配布したJupyter NotebookのPythonサンプルプログラムを参考に、独自に設定した処理を行うPythonプログラムを作成し、実行結果を提出

応用基礎データサイエンス・AI 第一

第2回 演習問題

2024.4.17

※全て半角文字・半角数字を用いて答えること

[1] 以下の記述の中で正しいものがいくつあるか、その個数を答えなさい。

- データ活用のための課題の一つにデータからの付加価値の生成が挙げられる。
- コンピュータ上のデータ表現について、最小単位はビット(bit)と呼ばれる。
- 1メガバイトは1024バイトに等しい。
- 2022年にSI接頭辞として加えられた R (ronna) とは  $10^{30}$  を意味する。
- 符号付き整数値を表す2進数表現ビット列が与えられたとき、最上位桁のビットが0であれば負数となる。

[演習1]

定義した MyCounter クラスで、カウンタのプリセットやセット時に整数ではなく実数値の3.5を与えた場合、どのような動作をするか。以下のプログラムを変更し、c.inc()、c.dec()を順に出力せよ。

```
c1 = MyCounter() # インスタンス名 c1 は変えないこと
print(c1.inc())
print(c1.dec())
```

[演習2]

カウント値として実数の整数部分を取り出し、常に整数を出力するようプログラムを修正したい。以下のコードについて、"..."と書かれている部分に適切なコードを記述しなさい。

```
class MyCounter_SI:
    def __init__(self, preset=0):
        self._cnt = preset

    def inc(self):
        self._cnt += 1
        return self._cnt
```

- 講義資料と演習資料ともに詳しい解説
- 受講生の予習や復習に適している
- 小テストと演習問題について自動採点を活用し、授業担当者の負担を軽減可能