

機械学習 データ分析コンペについて

岩政公平
iwamasa@morphometrics.jp
九州大学理学部生物学科
数理生物学研究室

データ分析コンペティションとは

**ある課題に対して統計学や機械学習などの
手法を用いて予測精度を競い合う**

どんな予測を行うか

SIIM-FISABIO-RSNA COVID-19 Detection - Kaggle

胸部X線写真からCOVID-19の疾患箇所を予測する(画像, 物体検出)

(<https://www.kaggle.com/c/siim-covid19-detection>)

Riid Answer Correctness Prediction - Kaggle

あるユーザーが問題を正解できるか行動履歴から予測する(時系列, 分類)

(<https://www.kaggle.com/c/riid-test-answer-prediction>)

Cornell Birdcall Identification - Kaggle

鳥の鳴き声からその鳥の種を予測する(音声, 分類)

(<https://www.kaggle.com/c/birdsong-recognition>)

マイナビ × SIGNATE Student Cup 2019: 賃貸物件の家賃予測

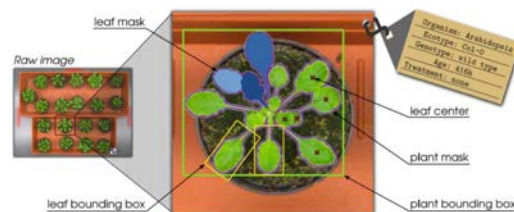
所在地や間取りから賃貸物件の価格を予測する(テーブル, 回帰)

(<https://signate.jp/competitions/182>)

など

どんな予測を行うか

INTRODUCTION: PLANT PHENOTYPING DATASETS



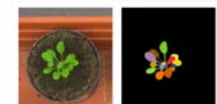
We present a collection of benchmark datasets in the context of plant phenotyping. We provide annotated imaging data and suggest suitable evaluation criteria for plant/leaf segmentation, detection, tracking as well as classification and regression problems. The figure symbolically depicts the data available together with ground truth segmentations and further annotations and metadata.

The Plant Phenotyping Datasets are intended for the development and evaluation of computer vision and machine learning algorithms such as (in parenthesis we point to general category of computer vision problems that these datasets can also be used for):

- plant detection and localization (multi-instance detection/localization)
- plant segmentation (foreground to background segmentation)
- leaf detection, localization, and counting (multi-instance detection, object counting)
- leaf segmentation (multi-instance segmentation)
- leaf tracking (multi-instance segmentation)
- boundary estimation for multi-instance segmentation (boundary detectors)
- classification and regression of mutants and treatments (general classification recognition)

The data can be used by scientists that already work in related fields but also from general computer vision scientists that work in related computer vision problems. No matter what, testing your algorithms on these data, you help us improve the state-of-the-art in phenotyping and feed the world one image at a time.

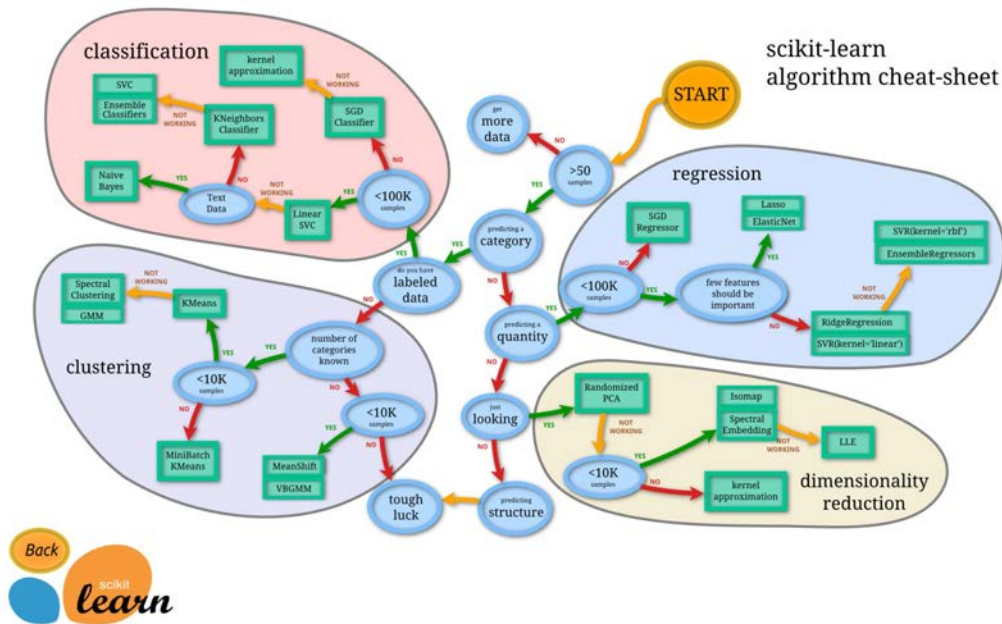
Sponsored by



<https://www.plant-phenotyping.org/datasets-home>

どんな魅力があるのか いろいろな手法(モデル)を知れる

- 機械学習のモデルは様々ある

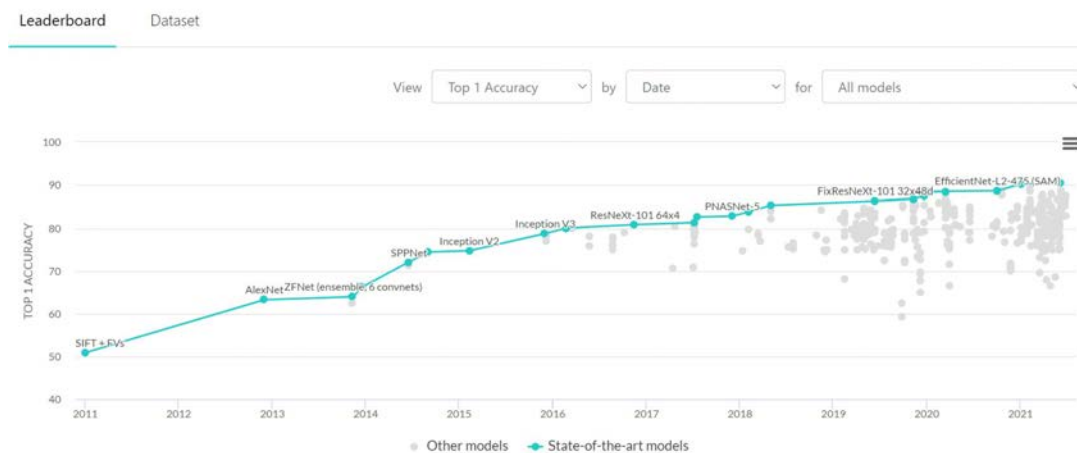


https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

どんな魅力があるのか いろいろな手法(モデル)を知れる

- より高精度なモデルが毎年提案される

Image Classification on ImageNet



<https://paperswithcode.com/sota/image-classification-on-imagenet>

どんな魅力があるのか

いろんな手法(モデル)を知れる

- コンペを通して様々なモデルを実際に試してみることができる
 - 新しく提案されたモデルが必ず精度がいいとは限らない
- 他の参加者とディスカッションすることができて互いにモデルの精度の良し悪しや実装コードを共有することができる
- 専門外の分野について知ることができる

モデル以外にも、

- データ加工の仕方や評価データの作成などの手法を知ることができる
 - 1億行あるテーブルデータや1000万枚ある画像データなど機械学習以外の部分で工夫しないとイケないコンペもある
 - データの可視化の方法なども学べる(可視化大事)
 - プログラミングスキルや英語力も上達するかも

どんな魅力があるのか

賞金・メダルがもらえる

- 上位に入賞すると賞金やメダル(称号)がもらえる
- Kaggleの場合は上位のチームにメダルが与えられる

Goldメダル	Silverメダル	Bronzeメダル
上位10チーム+α	上位5%	上位10%

- このメダルを集めると称号が与えられる(2021/07/19現在)

 226
Grandmasters

 1,591
Masters

 6,707
Experts

 60,277
Contributors

 90,549
Novices

どんな魅力があるのか

楽しい

- 競技性が高く、順位が伸びることが楽しい
 - ランキングで上位になるのは単純に嬉しい
 - コンペ期間中は常にランキングが更新されるためハラハラ感が楽しめる
- データを見てこの処理をすると精度が向上しそうと考える作業が楽しい
 - データの勘所が捉えられれば一気に精度が向上することもしばしば
- 知らない分野・モデルに触れることができる
 - 日本で行われるコンペもあるため初めてみたい人はここがいいかも
 - atmaCup: <https://www.guruguru.science/>
 - SIGNATE: <https://signate.jp/>
 - Nishika: <https://www.nishika.com/>