

Machine Learning for Clinical Predictive Analytics

Workshop

Wei-Hung Weng

BIG DATA FOR HEALTH WORKSHOPS AND CONFERENCE
Jul 10, 2018



Massachusetts
Institute of
Technology

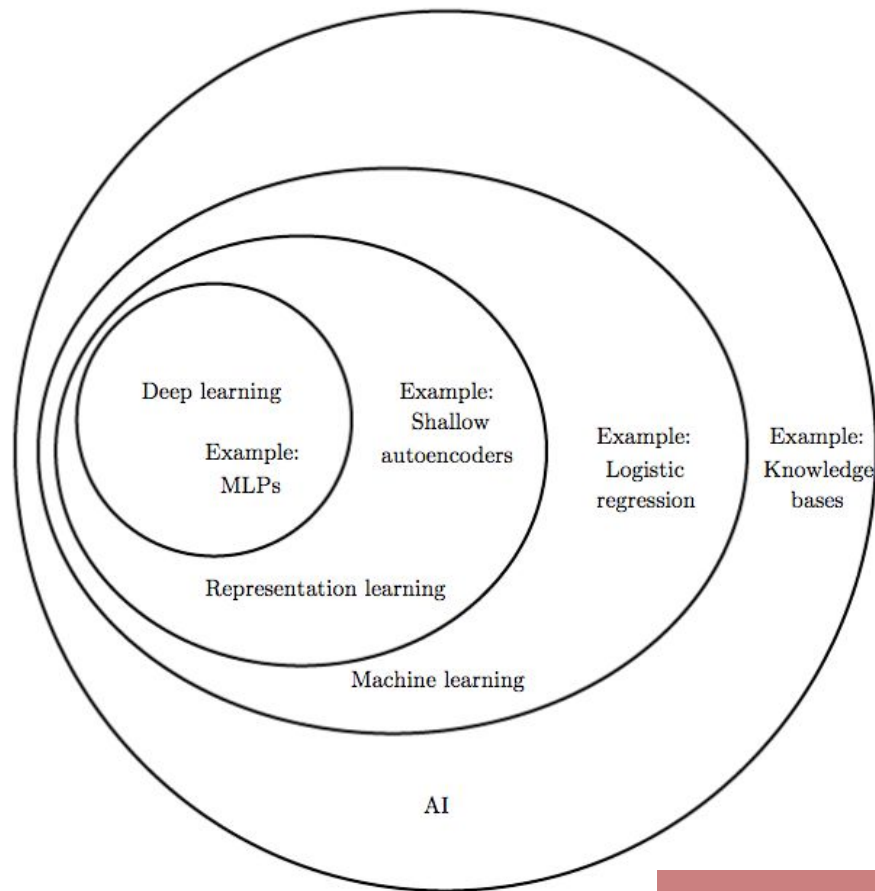


Artificial intelligence (AI) is defined by the American cognitive scientist Marvin Minsky as **the science of making machines do things that would require intelligence if done by man**

AI?

Human intelligence as a **goal**

Algorithms / Data science as an **approach**



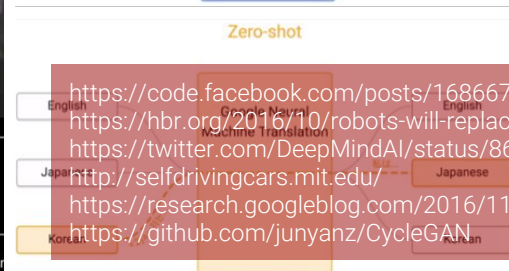
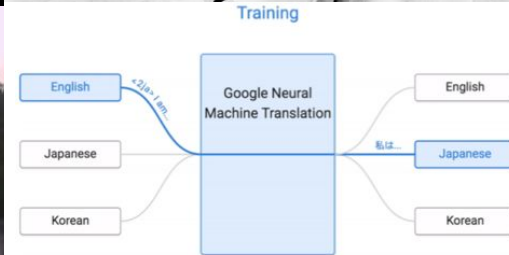


Technology Will Replace Many Doctors, Lawyers, and Other Professionals

by Leonard Susskind and Daniel Susskind

11, 2016

SHARE COMMENT TEXT SIZE PRINT \$8.95 BUY COPIES



<https://code.facebook.com/posts/1686672014972296/deal-or-no-deal-training-ai-bots-to-negotiate?pnref=story>

<https://hbr.org/2016/10/robots-will-replace-doctors-lawyers-and-other-professionals>

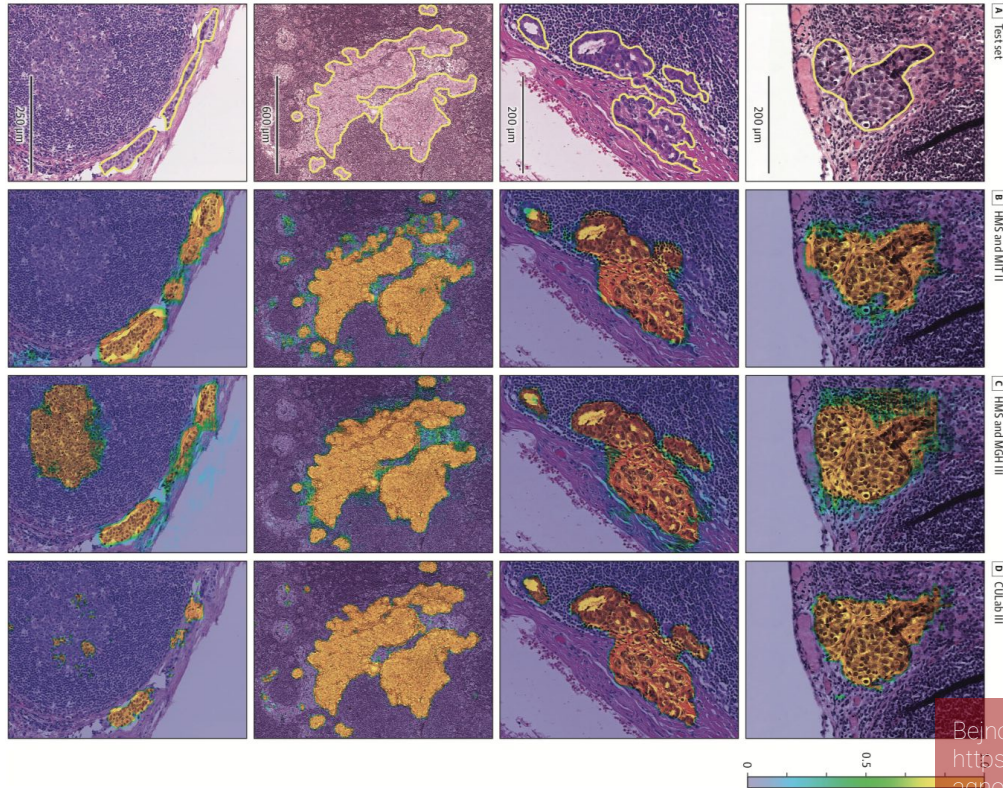
<https://twitter.com/DeepMindAI/status/867996695778410497/photo/1>

<http://selfdrivingcars.mit.edu/>

<https://research.googleblog.com/2016/11/zero-shot-translation-with-googles.html>

<https://github.com/junyanz/CycleGAN>

Biomedical Imaging Informatics



FDA approves first AI-powered diagnostic that doesn't need a doctor's help

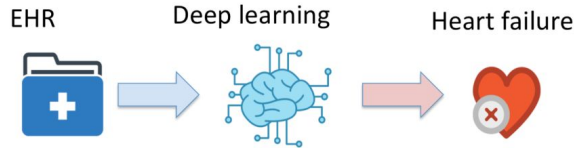
Marking a new era of "diagnosis by software," the US Food and Drug Administration on Wednesday gave permission to a company called IDx to market the first AI-powered diagnostic device.

What it does: The software is designed to detect greater than a mild level of diabetic retinopathy, which causes vision loss and affects 30 million people in the US. It occurs when high blood sugar damages blood vessels in the retina.

Bejnordi et al. JAMA 2017

<https://www.technologyreview.com/the-download/610853/fda-approves-first-ai-powered-diagnostic-that-doesnt-need-a-doctors-help/>

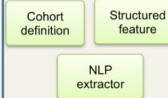
Clinical / Public Health Informatics



Data support



Data preparation



Aim 1: Prediction

1.1 Model Temporality

- Use RNN to longitudinal EHR
- HF prediction

1.2 Reduce dimensionality

- 2-level rep learning (visit, code)

Aim 2: Interpretation

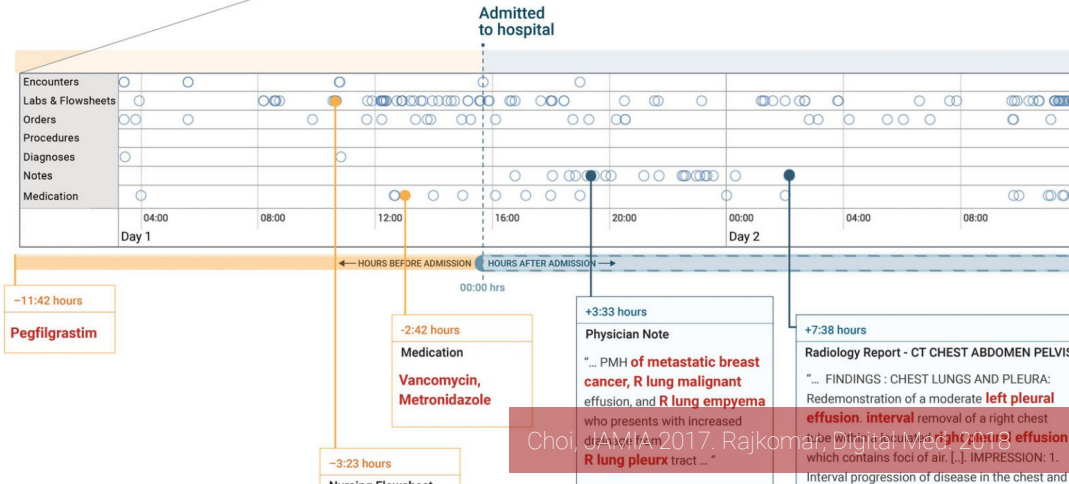
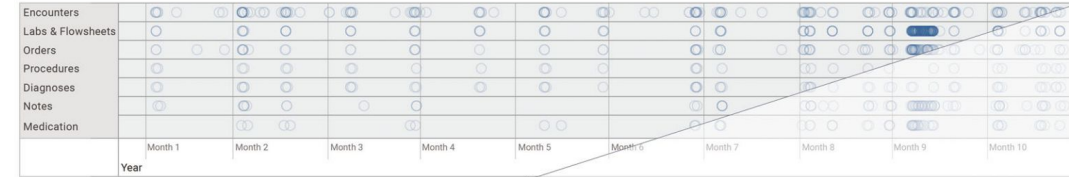
2.1 Reverse time attention

- Find
 - Key risk factors
 - Their timing

2.2 Temporal visualization

- Interactive viz for temporal models

Patient Timeline



Sciences Behind

Data organization / structuring

- Database
- Knowledge representation / Ontology
- Visualization

Learning from data (Algorithms)

- Statistics, probability, information theory, ...
- **Machine learning**
 - Computer vision, natural language processing, signal processing, ...

Machine Learning

Optimize a performance criterion using example data or past experience

- Given data X , we want to learn a function mapping $f(X)$ for certain purpose
 - Patient age, gender, vitals, labs, ... \rightarrow mortality Y/N
- ML - Given objective and evaluation metrics, and get high quality $f(X)$

Three steps

1. Modeling
 - Choose functions, features, (hyper)parameters
2. Evaluating the function
 - Define loss function, optimizer
3. Choosing the best one
 - Decide bias/variance trade-off

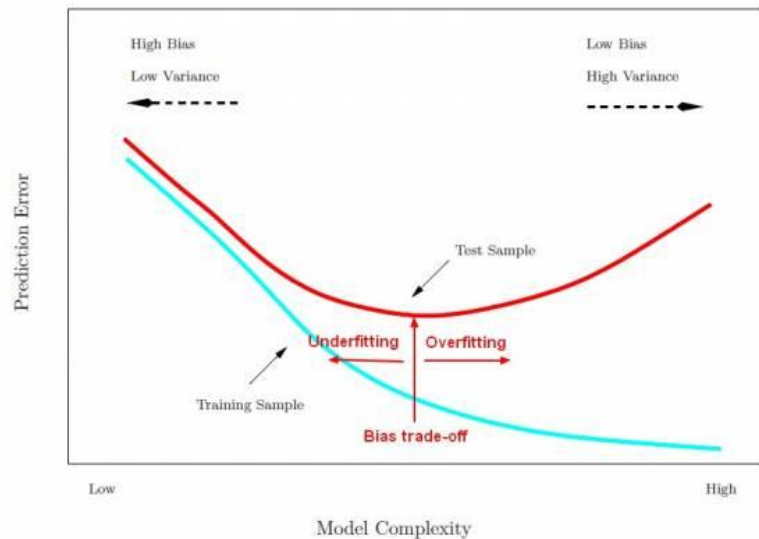
Bias vs. Variance

Bias (underfitting)

- High training/validation error
- Train more, increase model complexity, decrease regularization, add features

Variance (overfitting)

- Low training error but high validation/testing error
- More data, reduce model complexity, add regularization, reduce features

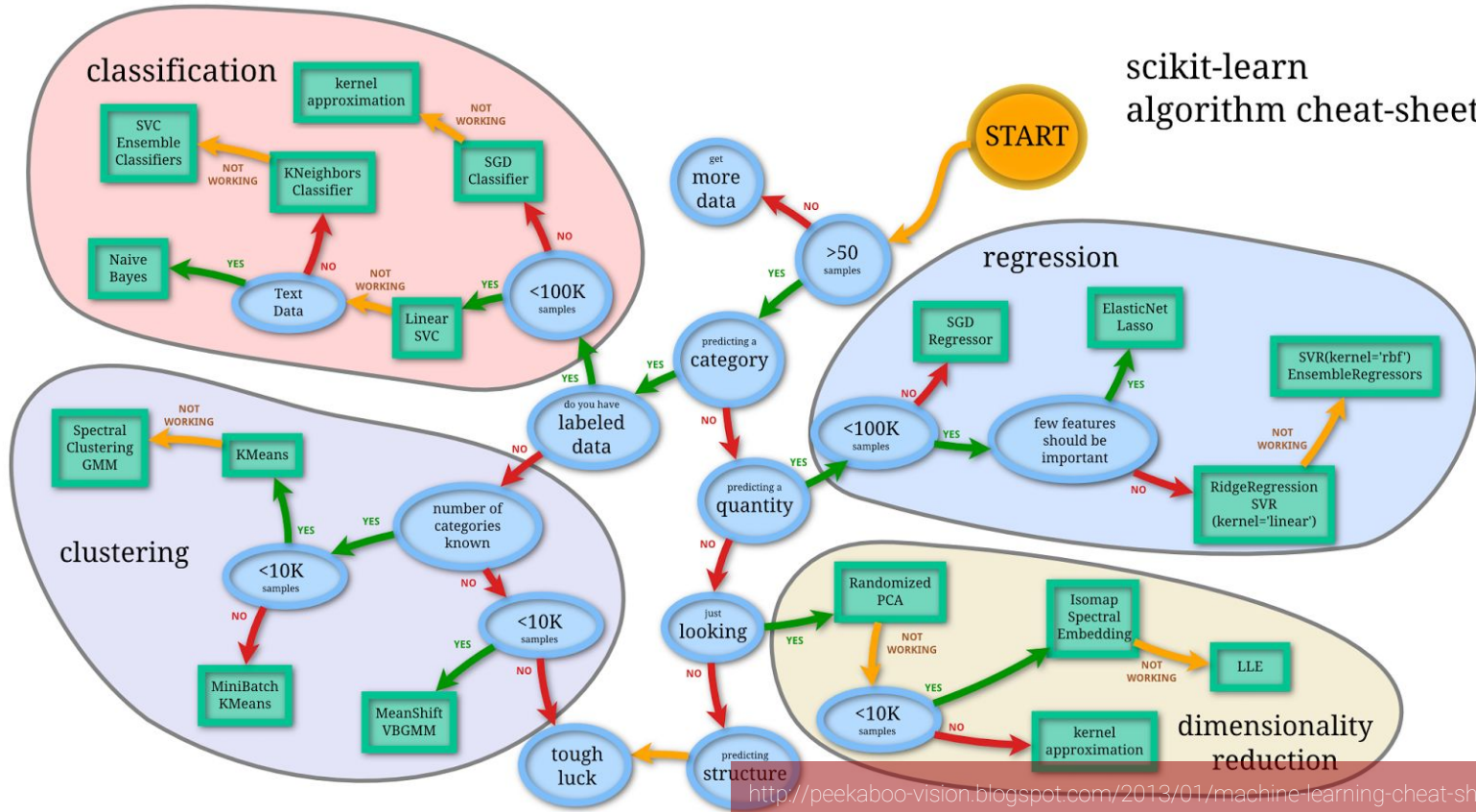


Scenario

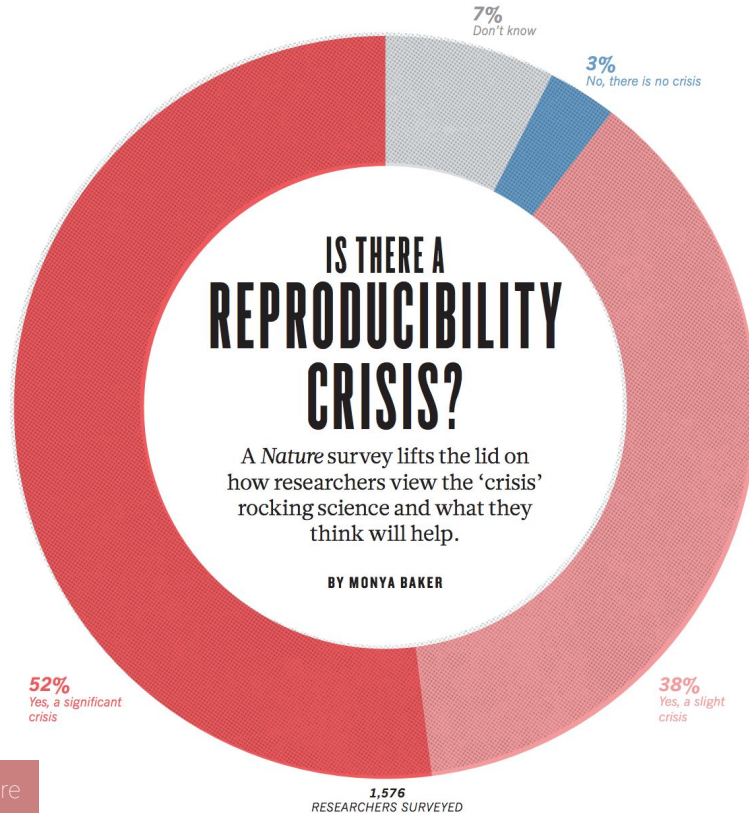
- Supervised learning
 - Regression
 - **Classification**
 - Linear
 - Non-linear (e.g. **deep learning**, SVM, decision tree, ...)
 - Structured learning
- Unsupervised learning
 - **Clustering**
 - **Dimensionality reduction**
- Transfer learning
- Reinforcement learning

Algorithms

scikit-learn
algorithm cheat-sheet



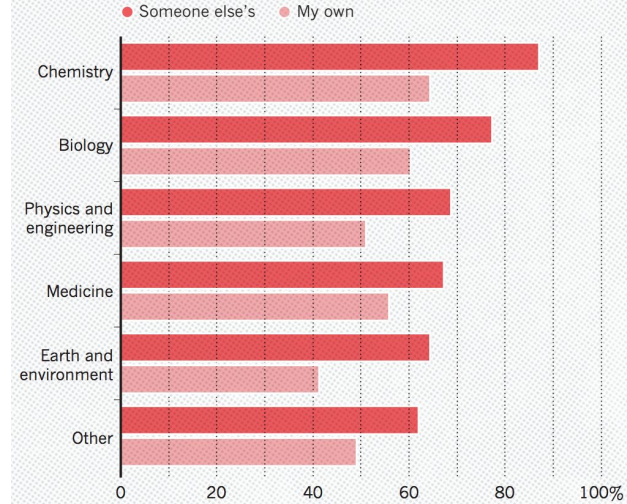
Reproducibility



Baker, 2015 *Nature*

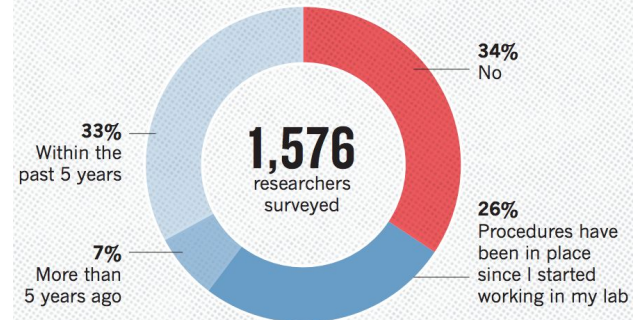
HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



HAVE YOU ESTABLISHED PROCEDURES FOR REPRODUCIBILITY?

Among the most popular strategies was having different lab members redo experiments.



ICLR 2018 Reproducibility Challenge

Background:

One of the challenges in machine learning research is to ensure that published results are reliable and reproducible. In support of this, the goal of this challenge is to investigate reproducibility of empirical results submitted to the [2018 International Conference on Learning Representations](#).

We are choosing ICLR for this challenge because the timing is right for course-based participants (see below), and because papers submitted to the conference are automatically made available publicly on [Open Review](#).

The Challenge is inspired by discussions at the ICML 2017 [Workshop on Reproducibility in Machine Learning](#).

Task Description

You should select a paper from the 2018 ICLR submissions, and aim to replicate the experiments described in the paper. The goal is to assess if the experiments are reproducible, and to determine if the conclusions of the paper are supported by your findings. Your results can be either positive (i.e. confirm reproducibility), or negative (i.e. explain what you were unable to reproduce, and potentially explain why).

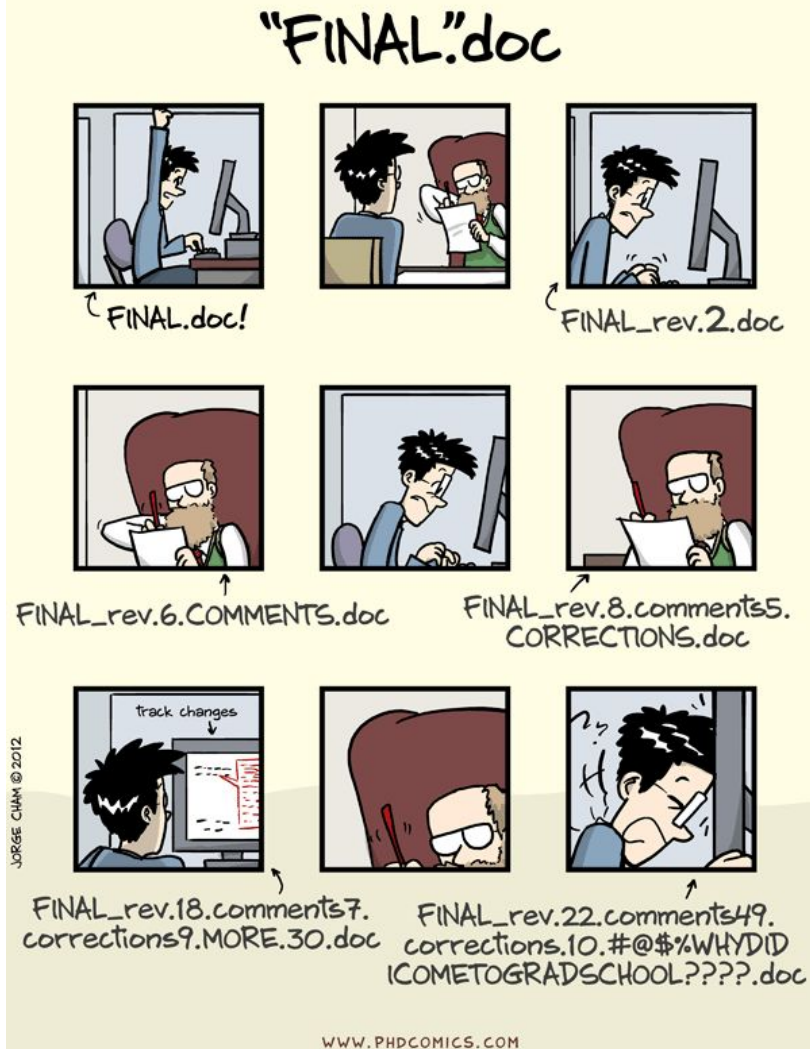
Essentially, think of your role as an inspector verifying the validity of the experimental results and conclusions of the paper. In some instances, your role will also extend to helping the authors improve the quality of their work and paper.

Reproducibility

What may be useful...

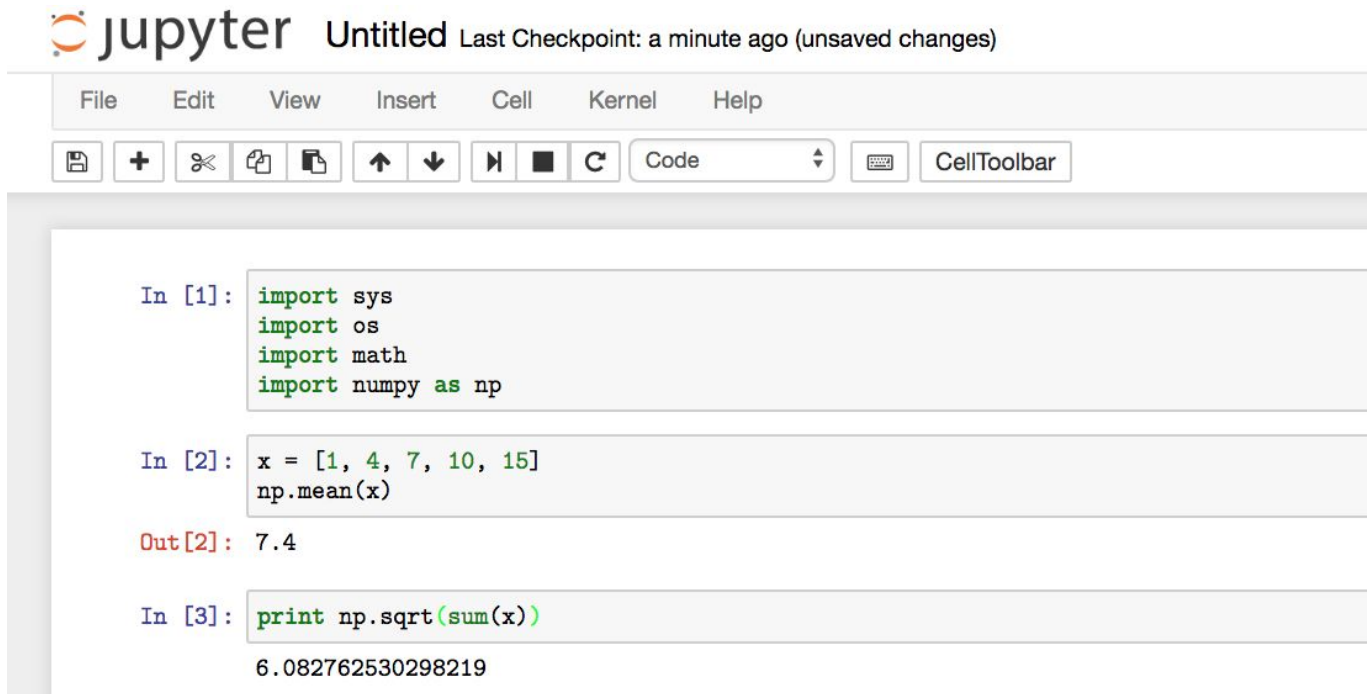
- Executable notebooks - `.ipynb`
- Code publishing platform - GitHub
- Version control - Git

Google colab



Jupyter Notebook

Execute (Shift + Enter) code cells and get your output underneath the cells



The screenshot displays the Jupyter Notebook interface. At the top, the title bar shows the Jupyter logo, the word "jupyter", and the document name "Untitled". To the right of the name, it says "Last Checkpoint: a minute ago (unsaved changes)". Below the title bar is a menu bar with options: File, Edit, View, Insert, Cell, Kernel, and Help. Underneath the menu bar is a toolbar containing icons for saving, adding a new cell, deleting a cell, duplicating a cell, moving a cell up/down, running the cell, and other functions. A dropdown menu is set to "Code", and a "CellToolbar" button is visible on the right. The main workspace contains three code cells. The first cell, labeled "In [1]:", contains four import statements: `import sys`, `import os`, `import math`, and `import numpy as np`. The second cell, labeled "In [2]:", contains two lines of code: `x = [1, 4, 7, 10, 15]` and `np.mean(x)`. Below this cell, the output is displayed as "Out[2]: 7.4". The third cell, labeled "In [3]:", contains the code `print np.sqrt(sum(x))`, and its output, "6.082762530298219", is shown below it.

jupyter Untitled Last Checkpoint: a minute ago (unsaved changes)

File Edit View Insert Cell Kernel Help

Code CellToolbar

```
In [1]: import sys
import os
import math
import numpy as np
```

```
In [2]: x = [1, 4, 7, 10, 15]
np.mean(x)
```

Out[2]: 7.4

```
In [3]: print np.sqrt(sum(x))
```

6.082762530298219

Code Publishing Platform

ckbjimmy / 2018_mlw

Unwatch

1

★ Star

0

Fork

0

<> Code

🔔 Issues 0

🔗 Pull requests 0

📁 Projects 0

📖 Wiki

📊 Insights

⚙️ Settings

2018 workshop materials

Edit

Add topics

📦 9 commits

🌿 1 branch

📦 0 releases

👤 1 contributor

📄 MIT

Branch: master

New pull request

Create new file

Upload files

Find file

Clone or download

ckbjimmy Created using Colaboratory

Latest commit ea47e02 10 days ago

data	data	25 days ago
.gitignore	data	25 days ago
LICENSE	Initial commit	25 days ago
README.md	Initial commit	25 days ago
nb1_classification.ipynb	Created using Colaboratory	11 days ago
nb2_clustering.ipynb	Created using Colaboratory	11 days ago
nb3_nn.ipynb	Created using Colaboratory	10 days ago

README.md

Version Control

- `git add .`
- `git status`
- `git commit -m 'first commit'`
- `git push`
- `git reset '.DS_Store'`

CoLab (CoLaboratory)

<http://g.co/colab>

Write code just as you would on a Jupyter Notebook

Use one of google's virtual machines to carry out your tasks

Free

Can write shell commands preceded with a '!'

- `!pip install gensim`
- `!ls`

Tutorial

- https://github.com/ckbjimmy/2018_mlw
 - Open → **"File" -> "Save a copy in Drive..."**
- Use python scikit-learn / keras
- Two datasets
- Part 1 - Supervised learning (classification)
 - When you have some labeled data
 - Given features, predict malignancy
 - ML general approaches, missing data imputation, normalization, important feature identification, ...
- Part 2 - Unsupervised learning (clustering / dimensionality reduction)
 - When you don't have labeled data
 - Grouping the similar cases
 - K-means / PCA

Tutorial

- Part 3 - Neural network
 - Deep feedforward neural network
 - ICU structured data
 - Breast cancer prediction data
 - Convolutional neural network (CNN) for image (MNIST)
 - Recurrent neural network (RNN) for text (IMDB reviews)

More ICU Problems

Classification

- Fit paCO_2 to pH / fit paCO_2 and HCO_3 to pH
- Classify ICU mortality using HCO_3 min/median & paCO_2 max/median
- Decision tree / random forest to determine mortality based on Hct, PLT min or WBC max (or based on vital signs, etc.)

Clustering

- Identify clusters of patients with $\text{HCO}_3/\text{paCO}_2$ or vital signs who did die or did not die during ICU stay
- Clustering vital signs and risk of being initiated on mechanical ventilation

Further Readings

Theory and mathematics

- Coursera Machine Learning (Andrew Ng)
- Deep learning book
- Stanford CS224n: Natural Language Processing with Deep Learning
- Stanford CS231n: Convolutional Neural Networks for Visual Recognition

Practical

- TensorFlow, PyTorch, Keras, Scikit-learn documents, guide, tutorials
- Google Machine Learning Crash Course
- Coursera Deep Learning Specialization
- Coursera Machine Learning with TensorFlow on Google Cloud Platform