

# How to Conduct a Machine Learning Project for Clinical Medicine

Wei-Hung Weng, MD, MMSc (MIT CSAIL)  
Joe Byers (BIDMC)



Khesar Gyalpo University of Medical Sciences of Bhutan  
Oct 27-29, 2019



Massachusetts  
Institute of  
Technology

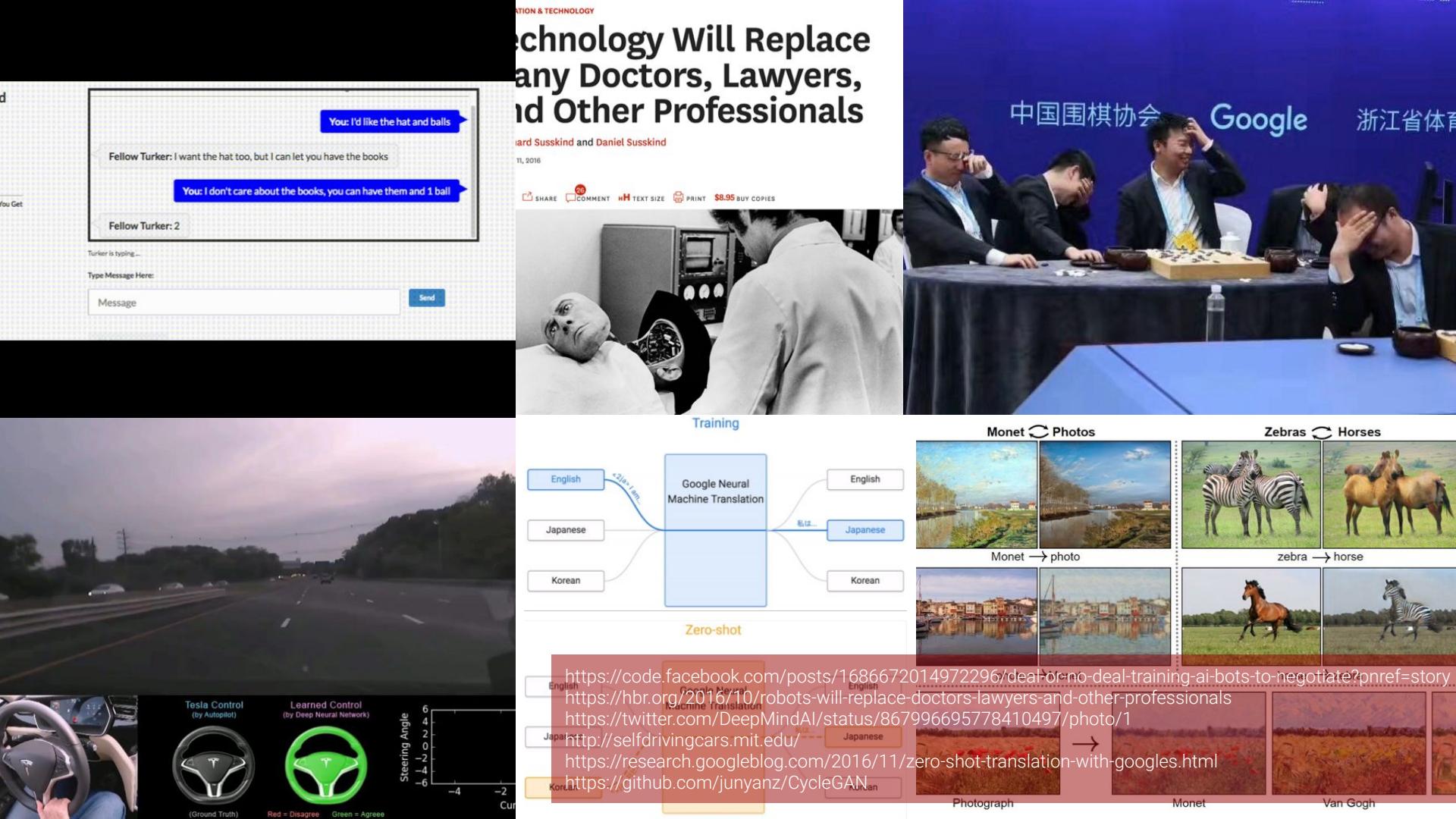


# Outline

- Hour 1
  - Why ML in clinical medicine?
  - Why reproducibility? How to make this happen?
    - Tools
    - Data and community (PhysioNet / eICU / MIMIC)
  - Caveats of ML for healthcare problems
- Hour 2
  - Some examples of clinical problems → ML tasks
  - Hands-on coding exercise / demo
- Quick survey
  - Programming / python?
  - ML?

# Why ML?

- Failure (?) of classical symbolic AI
  - A lot of knowledge is intuitive, difficult to put in rules and facts, not consciously accessible
- Principles giving rise to intelligence via learning
- ***Get knowledge directly from data and experience***
- To facilitate learning higher levels of abstraction [Bengio 2013, LeCun 2015]
  - Deep learning
- Can be described compactly since our intelligence is not just the result of a huge bag of tricks and pieces of knowledge, but of general mechanisms to acquire knowledge [Bengio 2013]



# Why ML/DL for EHR?

- Demographics, diagnoses, laboratory test results, medication prescriptions, clinical notes, medical images, ...
- Challenging
  - data quality (noisy, biased, ...)
  - data and annotation availability
  - heterogeneity of data types
- ***Traditional modeling → feature engineering***
  - labor intensive efforts
  - expert-defined phenotyping
  - ad-hoc feature engineering
  - limited generalizability across datasets or institutions
- ***Deep learning → learning hidden representations***
  - expert-driven feature engineering to data-driven feature construction



Demographics



Medications



Clinical Notes  
and Reports



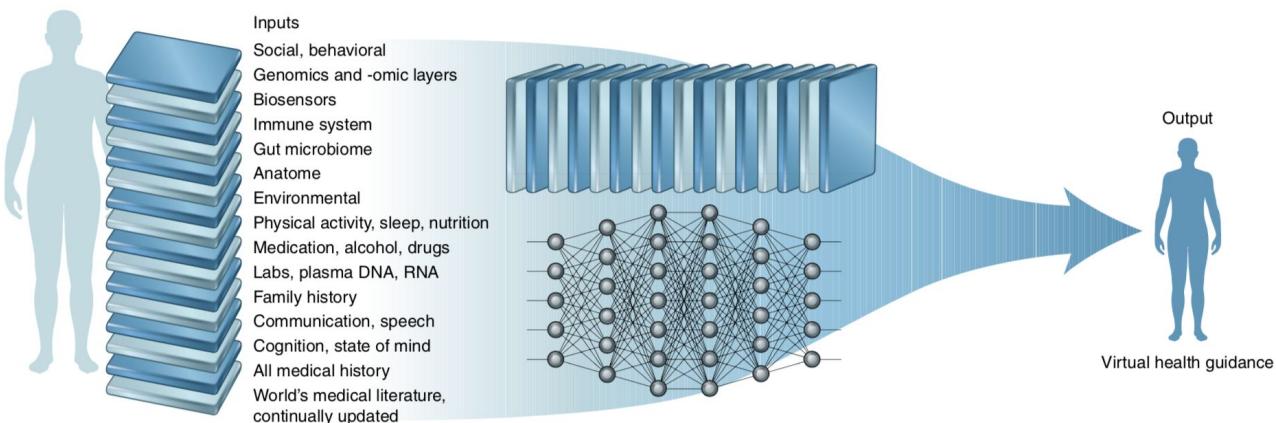
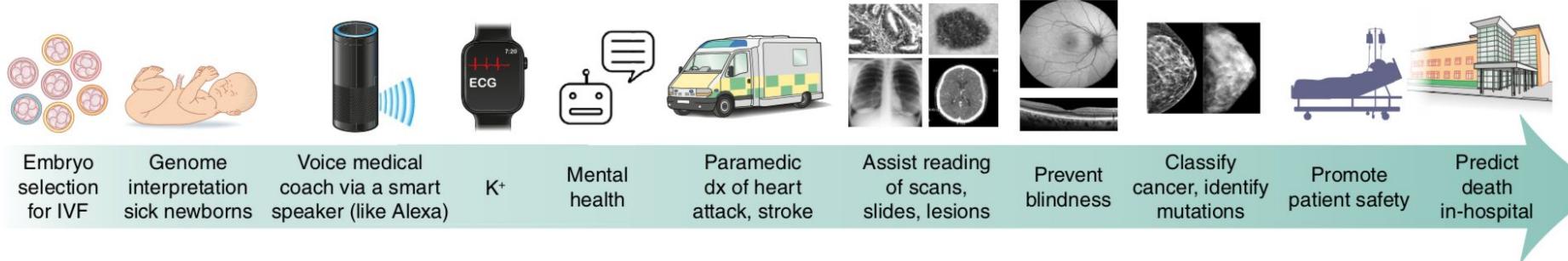
Continuous  
Monitoring Data



Multi-typed  
Medical Codes



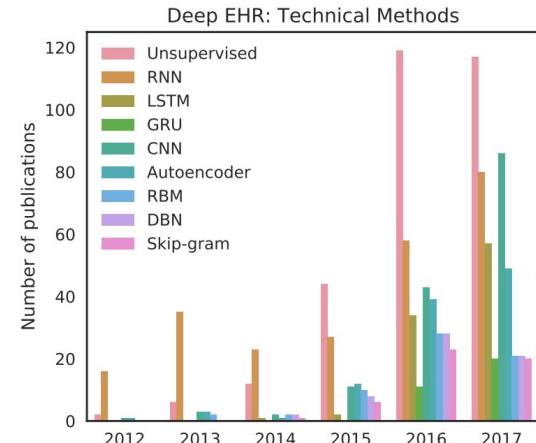
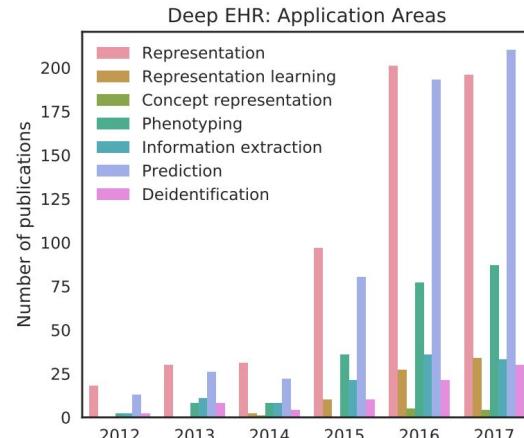
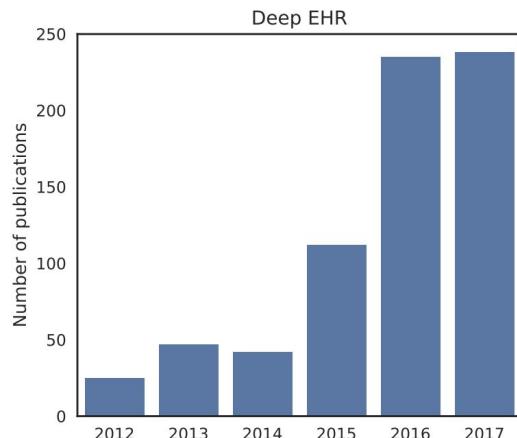
Medical  
Images



# Precision Medicine

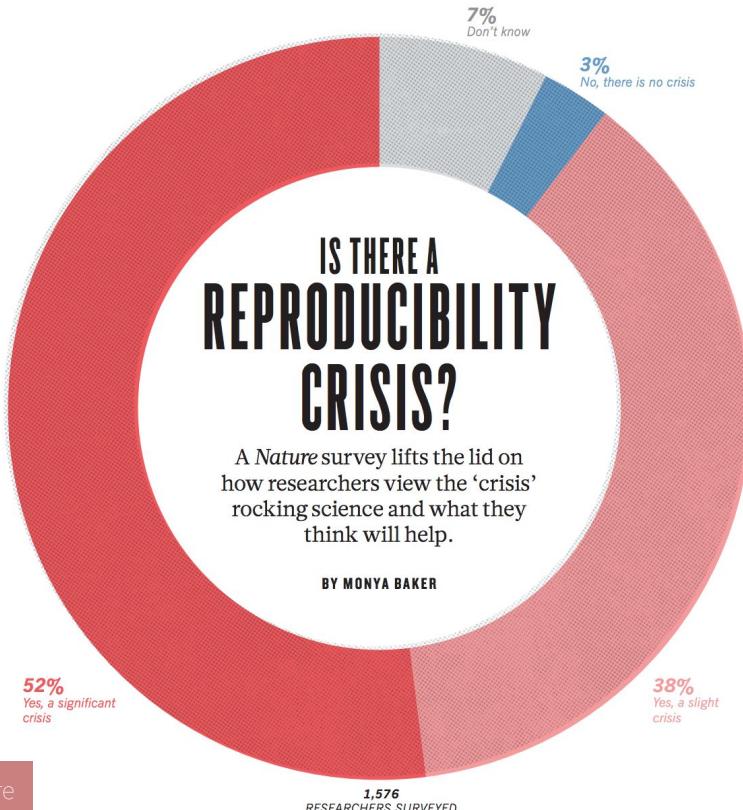
	<i>Traditional Approach</i>	<i>Precision Medicine Approach</i>		
<b>Population of Individuals</b>				
<b>Classify by Risk</b>				
<b>Surveillance for Preclinical Disease</b>				
<b>Signs or Symptoms</b>				
<b>Treat with</b>			 	
<b>Strategy</b>	<b>"One Size Fits All" Leads to Overall Mixed Results</b>		<b>Focus Existing</b>	<b>Repurpose FDA Approval</b>
	  			 
<b>Outcome</b>	 <b>Benefit</b>	 <b>No Effect</b>	 <b>Adverse</b>	 <b>Benefit</b>
			 <b>Benefit</b>	 <b>Benefit</b>

# Deep! (2017/06)

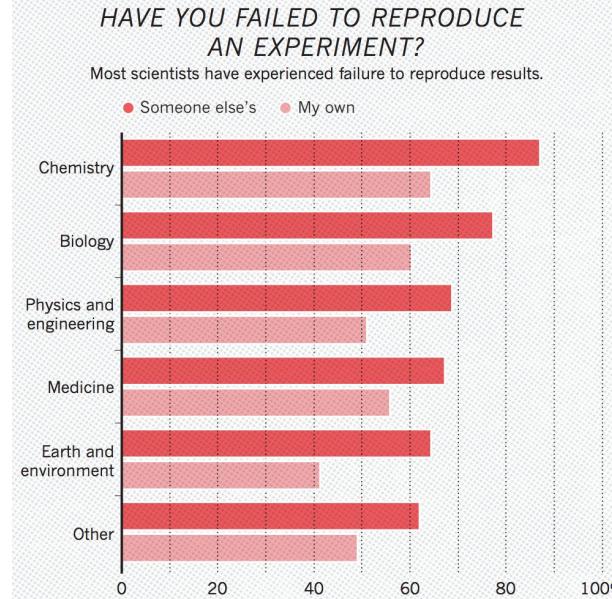


# Reproducibility

# Reproducibility

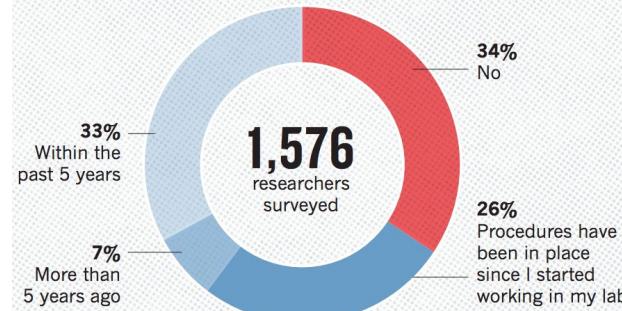


Baker, 2015 *Nature*



## 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

- Data + code (+ documentation!)
- For...
  - Validation
  - Learning
  - Collaboration



**Katerina Borodina** @kathyra\_ · Jul 7

my first coding job was with a small company, replacing their only developer who had recently quit. the code base was massive, in a language I had never used before. there was only one comment, at the end of an 1100 line function, and it said:

**//that'll do pig. that'll do.**

97 721 5.9K

# Tools

- Version control
  - Git
- Executable notebooks
  - Jupyter notebook (via anaconda)
  - Google colab
- Code publishing platform
  - GitHub
  - Gitlab

# Version Control

- git add .
- git status
- git commit -m 'first commit'
- git push
- git reset '.DS\_Store'
- ...

"FINAL".doc



JORGE CHAM © 2012

# Jupyter Notebook

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

The screenshot shows a Jupyter Notebook interface with the following components:

- Title Bar:** Shows the logo, "Untitled", and a note about the last checkpoint being a minute ago (unsaved changes).
- Toolbar:** Includes buttons for File, Edit, View, Insert, Cell, Kernel, Help, and various cell management icons.
- Cell 1:** An In [1] cell containing Python code:

```
import sys
import os
import math
import numpy as np
```
- Cell 2:** An In [2] cell containing Python code:

```
x = [1, 4, 7, 10, 15]
np.mean(x)
```
- Output 2:** An Out [2] cell showing the result of the mean calculation: 7.4.
- Cell 3:** An In [3] cell containing Python code:

```
print np.sqrt(sum(x))
```
- Output 3:** The result of the square root calculation: 6.082762530298219.

# CoLab (CoLaboratory)

<http://g.co/colab>

Write code just as you would on a Jupyter Notebook

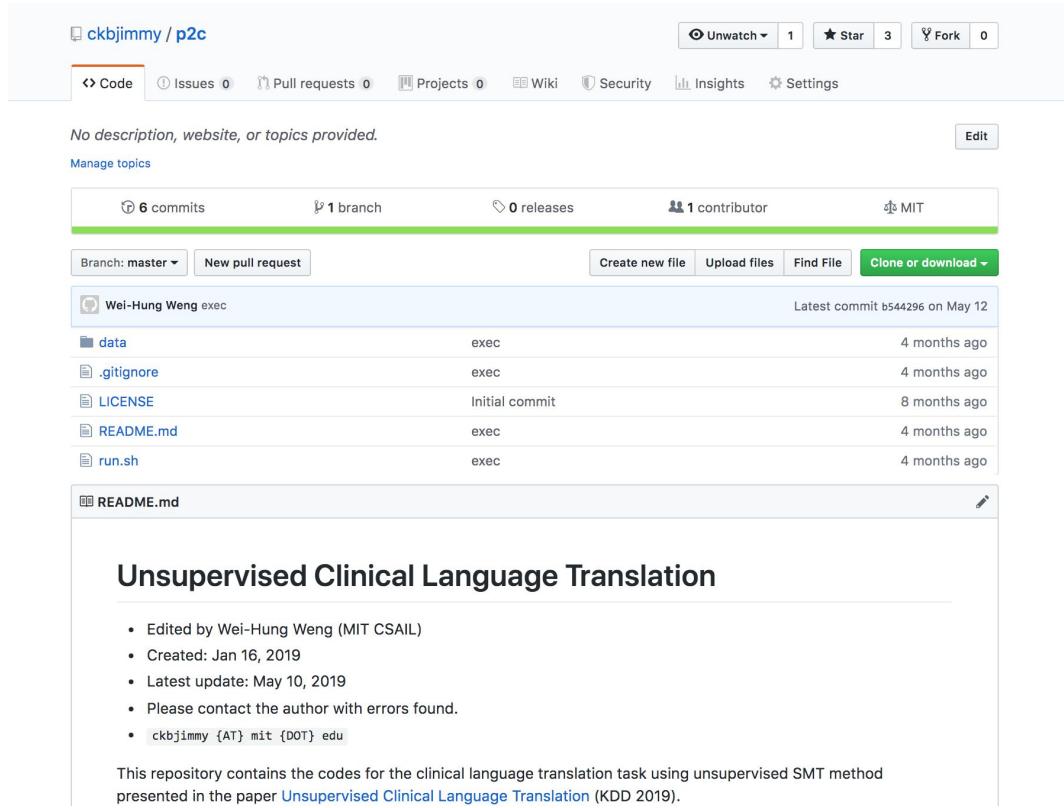
Use one of google's virtual machines to carry out your tasks

Free

Can write shell commands preceded with a ‘!’

- !pip install gensim
- !ls

# Code Publishing (scary but worth it)

 ckbjimmy / p2c

Unwatch 1 Star 3 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Security Insights Settings

No description, website, or topics provided. Edit

Manage topics

6 commits 1 branch 0 releases 1 contributor MIT

Branch: master New pull request Create new file Upload files Find File Clone or download

Wei-Hung Weng exec Latest commit b544296 on May 12

data	exec	4 months ago
.gitignore	exec	4 months ago
LICENSE	Initial commit	8 months ago
README.md	exec	4 months ago
run.sh	exec	4 months ago

README.md

## Unsupervised Clinical Language Translation

- Edited by Wei-Hung Weng (MIT CSAIL)
- Created: Jan 16, 2019
- Latest update: May 10, 2019
- Please contact the author with errors found.
- ckbjimmy {AT} mit {DOT} edu

This repository contains the codes for the clinical language translation task using unsupervised SMT method presented in the paper [Unsupervised Clinical Language Translation \(KDD 2019\)](#).

**PhysioNet / eICU / MIMIC-III**

# physionet.org

The PhysioNet website features a dark header bar at the top. On the left is the "PhysioNet" logo. To its right are links for "Find", "Share", "About", and "News". Further to the right is an "Account" dropdown menu, a search input field with a magnifying glass icon, and a small "Search" button.

# PhysioNet

The Research Resource for Complex Physiologic Signals

[Data](#)   [Software](#)   [Challenges](#)   [Tutorials](#)

 [Database](#)  [Credentialed Access](#)

## eICU Collaborative Research Database

Tom Pollard , Alistair Johnson , Jesse Raffa , Leo Anthony Celi , Omar Badawi , Roger Mark

Published: April 15, 2019. Version: 2.0

### When using this resource, please cite:

Pollard, T., Johnson, A., Raffa, J., Celi, L. A., Badawi, O., Mark, R. (2019). eICU Collaborative Research Database. PhysioNet. doi:10.13026/C2WM1R

### Additionally, please cite the original publication:

[The eICU Collaborative Research Database, a freely available multi-center database for critical care research. Pollard TJ, Johnson AEW, Raffa JD, Celi LA, Mark RG and Badawi O. Scientific Data \(2018\), DOI: <http://dx.doi.org/10.1038/sdata.2018.178>.](#)

### Please include the standard citation for PhysioNet:

Goldberger AL, Amaral LAN, Glass L, Hausdorff JM, Ivanov PCh, Mark RG, Mietus JE, Moody GB, Peng C-K, Stanley HE. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals (2003). Circulation. 101(23):e215-e220.

### Contents [▼](#)

#### Share



#### Access

##### Access Policy:

Only PhysioNet credentialed users who sign the specified DUA can access the files.

##### License (for files):

[PhysioNet Credentialed Health Data License 1.5.0](#)

## Abstract

The eICU Collaborative Research Database is a multi-center database comprising deidentified health data associated with over 200,000 admissions to ICUs across the United States between 2014–2015. The database includes vital sign measurements, care plan documentation, severity

[Database](#) [Credentialed Access](#)

## MIMIC-III Clinical Database

Alistair Johnson , Tom Pollard , Roger Mark

Published: Sept. 4, 2016. Version: 1.4

**When using this resource, please cite:**

Johnson, A., Pollard, T., Mark, R. (2016). MIMIC-III Clinical Database. PhysioNet.  
doi:10.13026/C2XW26

**Additionally, please cite the original publication:**

Johnson, A. E. W., Pollard, T. J., Shen, L., Lehman, L. H., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Celi, L. A., & Mark, R. G. (2016). MIMIC-III, a freely accessible critical care database. *Scientific Data*, 3, 160035.

**Please include the standard citation for PhysioNet:**

Goldberger AL, Amaral LAN, Glass L, Hausdorff JM, Ivanov PCh, Mark RG, Mietus JE, Moody GB, Peng C-K, Stanley HE. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals (2003). *Circulation*. 101(23):e215-e220.

## Abstract

MIMIC-III is a large, freely-available database comprising deidentified health-related data associated with over forty thousand patients who stayed in critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012.

[Contents](#) ▾**Share****Access****Access Policy:**

Only PhysioNet credentialed users who sign the specified DUA can access the files.

**License (for files):**

[PhysioNet Credentialed Health Data License 1.5.0](#)

# MIMIC-CXR

PhysioNet Find Share About News Account ▾ Search 

 Database  Credentialed Access

## MIMIC-CXR Database

Alistair Johnson , Tom Pollard , Roger Mark , Seth Berkowitz , Steven Horng 

Published: Sept. 19, 2019. Version: 2.0.0

### When using this resource, please cite:

Johnson, A., Pollard, T., Mark, R., Berkowitz, S., Horng, S. (2019). MIMIC-CXR Database. PhysioNet. doi:10.13026/C2JT1Q

### Please include the standard citation for PhysioNet:

Goldberger AL, Amaral LAN, Glass L, Hausdorff JM, Ivanov PCh, Mark RG, Mietus JE, Moody GB, Peng C-K, Stanley HE. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals (2003). Circulation. 101(23):e215-e220.

### Contents ▾

#### Share



#### Access

##### Access Policy:

Only PhysioNet credentialed users who sign the specified DUA can access the files.

##### License (for files):

[PhysioNet Credentialed Health Data License 1.5.0](#)

## Abstract

The MIMIC Chest X-ray (MIMIC-CXR) Database v2.0.0 is a large publicly available dataset of chest radiographs in DICOM format with free-text radiology reports. The dataset contains 377,110 images corresponding to 227,835 radiographic studies performed at the Beth Israel Deaconess Medical Center in Boston, MA. The dataset is de-identified to satisfy the US Health Insurance Portability and Accountability Act of 1996 (HIPAA) Safe Harbor requirements. Protected health information (PHI) has been removed. The dataset is intended to support a wide body of research in medicine including image understanding, natural language processing, and decision support.

# Warning

- Not in a clean, tidy spreadsheet
- Not generated for us to do research



reddit



r/MachineLearning



Find a community, post, or user

LOG IN

↑ captkrob 2 points · 1 year ago

↓ If you want an idea of a "real world" disease dataset, try looking at the MIMIC-III database (<https://mimic.physionet.org/>).

It's probably significantly more work to request the data and set it up in an easy format for ML than you had ever planned on, but this is the kind of thing actual data scientists and informaticians working on this kind of problem are used to dealing with. The dataset is tremendously rich, but also incredibly noisy. In other words, much like medical practice.

Share Save

# Getting Access

- Sign data use agreement (DUA)
- Take online courses (citi)

## Sign Data Use Agreement - eICU Collaborative Research Database v2.0

Sign the following data use agreement to access the files in [eICU Collaborative Research Database v2.0](#).

### PhysioNet Credentialled Health Data Use Agreement 1.5.0

If I am granted access to the database:

1. I will not attempt to identify any individual or institution referenced in PhysioNet restricted data.
2. I will exercise all reasonable and prudent care to avoid disclosure of the identity of any individual or institution referenced in PhysioNet restricted data in any publication or other communication.
3. I will not share access to PhysioNet restricted data with anyone else.
4. I will exercise all reasonable and prudent care to maintain the physical and electronic security of PhysioNet restricted data.
5. If I find information within PhysioNet restricted data that I believe might permit identification of any individual or institution, I will report the location of this information promptly by email to [PHI-report@physionet.org](mailto:PHI-report@physionet.org), citing the location of the specific information in question.
6. I have requested access to PhysioNet restricted data for the sole purpose of lawful use in scientific research, and I will use my privilege of access, if it is granted, for this purpose and no other.
7. I have completed a training program in human research subject protections and HIPAA regulations, and I am submitting proof of having done so.
8. I will indicate the general purpose for which I intend to use the database in my application.
9. If I openly disseminate my results, I will also contribute the code used to produce those results to a repository that is open to the research community.
10. This agreement may be terminated by either party at any time, but my obligations with respect to PhysioNet data shall continue after termination.

The screenshot shows the CITI Program website. At the top right, there are links for '+1 888.529.5929', 'English', 'Register', and 'Log In'. The main navigation menu includes 'Subscriptions', 'Courses', 'CE/CMEs', 'Tools', and 'Support'. A search icon is also present. The central content area features a large blue banner with the text 'Human Subjects Research (HSR)' in yellow. Below the banner, a sub-section titled 'HSR provides foundational training in human subjects research and includes the historical development of human subject protections, ethical issues, and current regulatory and guidance information.' is visible. On the left, a sidebar lists various training categories: 'View All', 'CE Certified Courses', 'Animal Care and Use (ACU)', 'Bioethics', 'Biomedical PI', 'Biosafety and Biosecurity (BSS)', 'Clinical Research Coordinator (CRC)', and 'Clinical Trial Billing Compliance (CTBC)'. At the bottom right, there are 'ORGANIZATIONS' and 'LEARNERS' sections with 'LEARN MORE' and 'BUY NOW' buttons, along with 'Questions?' and 'Contact Us' links.

# Data Sharing

- More people seeing the data → more knowledge
- Education
- Accelerating research
- You can also share the data through PhysioNet!

# MIMIC-III

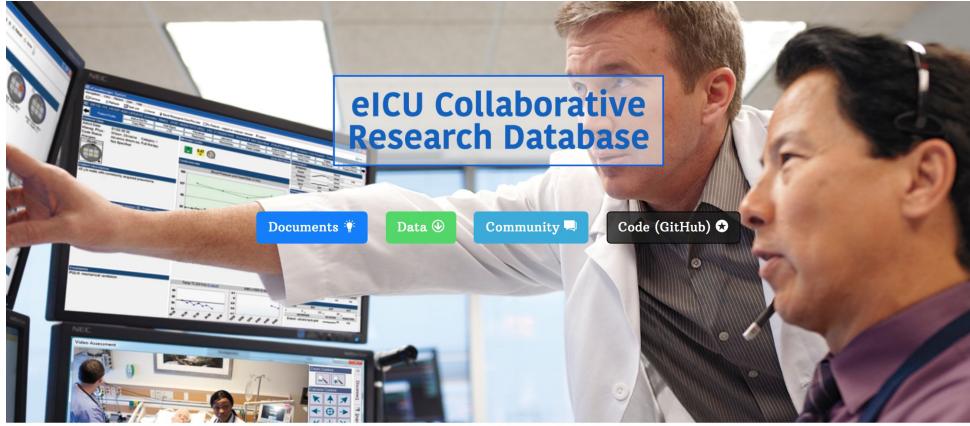
- Single center (>60K ICU stay, >40K patients)
- ICU data in details (almost everything incl. notes and waveforms)
- Limited out-of-ICU data (medication, DOD, ...)



*MIMIC-III, a freely accessible critical care database.* Johnson AEW,  
Pollard TJ, Shen L, Lehman L, Feng M, Ghassemi M, Moody B, Szolovits P,  
Celi LA, and Mark RG. *Scientific Data* (2016). DOI: 10.1038/sdata.2016.35.  
Available from: <http://www.nature.com/articles/sdata201635>

# eICU

- Multicenter (>200 hospitals, >200K ICU stay, 2 years)
- ICU data (no free text notes)
- Heterogeneous data (every hospital has their own input format)



If you use the eICU Collaborative Research Database in your work, please cite the following publication:

*The eICU Collaborative Research Database, a freely available multi-center database for critical care research. Pollard TJ, Johnson AEW, Raffa JD, Celi LA, Mark RG and Badawi O. Scientific Data (2018). DOI: <http://dx.doi.org/10.1038/sdata.2018.178>. Available from: <https://www.nature.com/articles/sdata2018178>*

# eICU vs. MIMIC-III

- Strengths of eICU
  - Latest data (2014-2015)
  - Larger dataset, multicenter
  - Severity score for all patients (APACHE)
  - Treatment plan in `carePlanGeneral` table
  - Diagnosis available during stay
- Strengths of MIMIC
  - Documentation, codebase and community support
  - More literature
  - Consistent within hospital
  - With clinical notes

# Navigating eICU [eicu-crd.mit.edu]

The screenshot shows a navigation sidebar on the left with links like About, Getting started, and various tables. The main content area displays the apachePredVar page, which includes a purpose section and a detailed text about APACHE predictions. A sidebar on the right contains sections for Important considerations, Table columns, and Further reading.

**eICU Collaborative Research Database** eICU Collaborative Research Database

About >

Getting started >

Tables in eICU-CRD ▾

- admissiondrug
- admissiondx
- allergy
- apacheApsVar
- apachePatientResult
- apachePredVar**
- carePlanCareProvider
- carePlanEOL
- carePlanGeneral
- carePlanGoal
- carePlanInfectiousDisease

## apachePredVar

**Purpose:** Provides variables underlying the APACHE predictions. Acute Physiology Age Chronic Health Evaluation (APACHE) consists of a groups of equations used for predicting outcomes in critically ill patients. APACHE II is based on the APS or acute physiology score (which uses 12 physiologic values), age, and chronic health status within one of 56 disease groups. APACHE II is no longer considered valid due to inadequate case mix index adjustments and over estimates mortality because it is based on models from the 1970s-1980s. APACHE III, introduced in 1991, improved the equation by changing the number and weights of the APS and revising the measurement of chronic health status. APACHE IVa further improved the equations and has been described as having the highest discrimination of any other adult risk adjustment model (SAPS 3, SOFA, MPM III).

**apachePredVar**

**Important considerations**

**Table columns** **Further reading**

# eICU Patient Tracking

- `patientunitstayid` → ICU stay
- `patienthealthsystemstayid` → hospital stay
- `uniquepid` → patient
- `hospitalid` → hospital
- ...offset

# Tables

- **apacheapsvar**
  - First day aggregated data for APACHE
- **apachepredvar**
- **diagnosis** → offset, string, ICD code
- **infusiondrug** → fluid, insulin, vasopressor, sedative
- **intakeoutput** → urine output
- **patient** → demographics
- **lab, medication, pasthistory, treatment** → offset, string
- **vitalperiodic, vitalaperiodic**
  
- [github.com/ckbjimmy/hst953\\_iomed](https://github.com/ckbjimmy/hst953_iomed)

# Codebase [github.com/MIT-LCP/eicu-code]

MIT-LCP / eicu-code

Unwatch 35    Unstar 94    Fork 76

Code Issues 29 Pull requests 3 Projects 0 Wiki Security Insights

Code and website related to the eICU Collaborative Research Database <https://eicu-crd.mit.edu>

database eicu-crd mimic physionet ehr healthcare

187 commits 3 branches 1 release 8 contributors MIT

Branch: master New pull request Create new file Upload files Find File Clone or download

File	Commit Message	Date
.gitignore	add ds_store	2 years ago
LICENSE	clean up	3 years ago
README.md	add doi badge	last year
styleguide.md	remove external links	last year
README.md		

eICU Collaborative Research Database Code Repository

DOI 10.5281/zenodo.1249016

# Concepts

- In BigQuery
- Derived tables
  - [/eicu-code/concepts/pivoted/](#)

Branch: master ▾ [eicu-code](#) / [concepts](#) / [pivoted](#) /

 <a href="#">alistairewj</a>	add chloride		Latest commit <a href="#">c66c0fe</a> on May 12
..			
<a href="#">pivotedsqll</a>	add pivot tables		last year
<a href="#">pivotedsqll</a>	add pivot tables		last year
<a href="#">pivotedsqll</a>	add chloride		5 months ago
<a href="#">pivotedsqll</a>	remove superfluous column		last year
<a href="#">pivotedsqll</a>	add pivoted views		last year
<a href="#">pivotedsqll</a>	add pivoted views		last year
<a href="#">pivotedsqll</a>	add pivot tables		last year
<a href="#">pivotedsqll</a>	add pivot tables		last year
<a href="#">pivotedsqll</a>	add pivoted views		last year

```
, max(case
    when drugname in
        (
            'EPI (mcg/min)'
            , 'Epinepherine (mcg/min)'
            , 'Epinephrine'
            , 'Epinephrine ()'
            , 'EPINEPHrine(Adrenalin)MAX 30 mg Sodium Chloride 0.9% 250 ml (mcg/min)'
            , 'EPINEPHrine(Adrenalin)STD 4 mg Sodium Chloride 0.9% 250 ml (mcg/min)'
            , 'EPINEPHrine(Adrenalin)STD 4 mg Sodium Chloride 0.9% 500 ml (mcg/min)'
            , 'EPINEPHrine(Adrenalin)STD 7 mg Sodium Chloride 0.9% 250 ml (mcg/min)'
            , 'Epinephrine (mcg/hr)'
            , 'Epinephrine (mcg/kg/min)'
            , 'Epinephrine (mcg/min)'
            , 'Epinephrine (mg/hr)'
            , 'Epinephrine (mg/kg/min)'
            , 'Epinephrine (ml/hr)'
        ) then 1 else 0 end)
    as epinephrine
```

# Are You the First Person to Ask → GitHub Issues

MIT-LCP / [eicu-code](#)

[Unwatch](#) 35   [Unstar](#) 94   [Fork](#) 76

[Code](#)   [Issues 29](#)   [Pull requests 3](#)   [Projects 0](#)   [Wiki](#)   [Security](#)   [Insights](#)

[Filters](#)  [Labels 7](#) [Milestones 0](#) [New issue](#)

Author	Labels	Projects	Milestones	Assignee	Sort
ragi24					
acanakoglu					
remifol					
RyuheiSo					
eruca					
shong95					

**Dialysis**  
#83 opened 6 days ago by ragi24

**Mechanical ventilation in eICU**  
#82 opened 14 days ago by acanakoglu 1 comment

**APACHE score calculation for missing data**  
#81 opened 23 days ago by remifol 2 comments

**How to differentiate cardiac deaths from non-cardiac deaths?**  
#80 opened 26 days ago by RyuheiSo

**how to define patients first access to the icu stay?**  
#79 opened on Aug 9 by eruca

**Stuck in the middle of procedure of eICU setting up**  
#78 opened on Jul 23 by shong95 10 comments

How?

# Process

- How to develop machine learning models for healthcare [Chen 2019]
- Machine Learning for Clinical Predictive Analytics [Weng 2019]
- ***Problem definition***
- ***Data curation***
- ***ML model development***
- ***Validation***
- ***Assessment of clinical impact***
- ***(Deployment and monitoring)***

comment

## How to develop machine learning models for healthcare

Rapid progress in machine learning is enabling opportunities for improved clinical decision support. Importantly, however, developing, validating and implementing machine learning models for healthcare entail some particular considerations to increase the chances of eventually improving patient care.

Po-Hsuan Cameron Chen, Yun Liu and Lily Peng

---

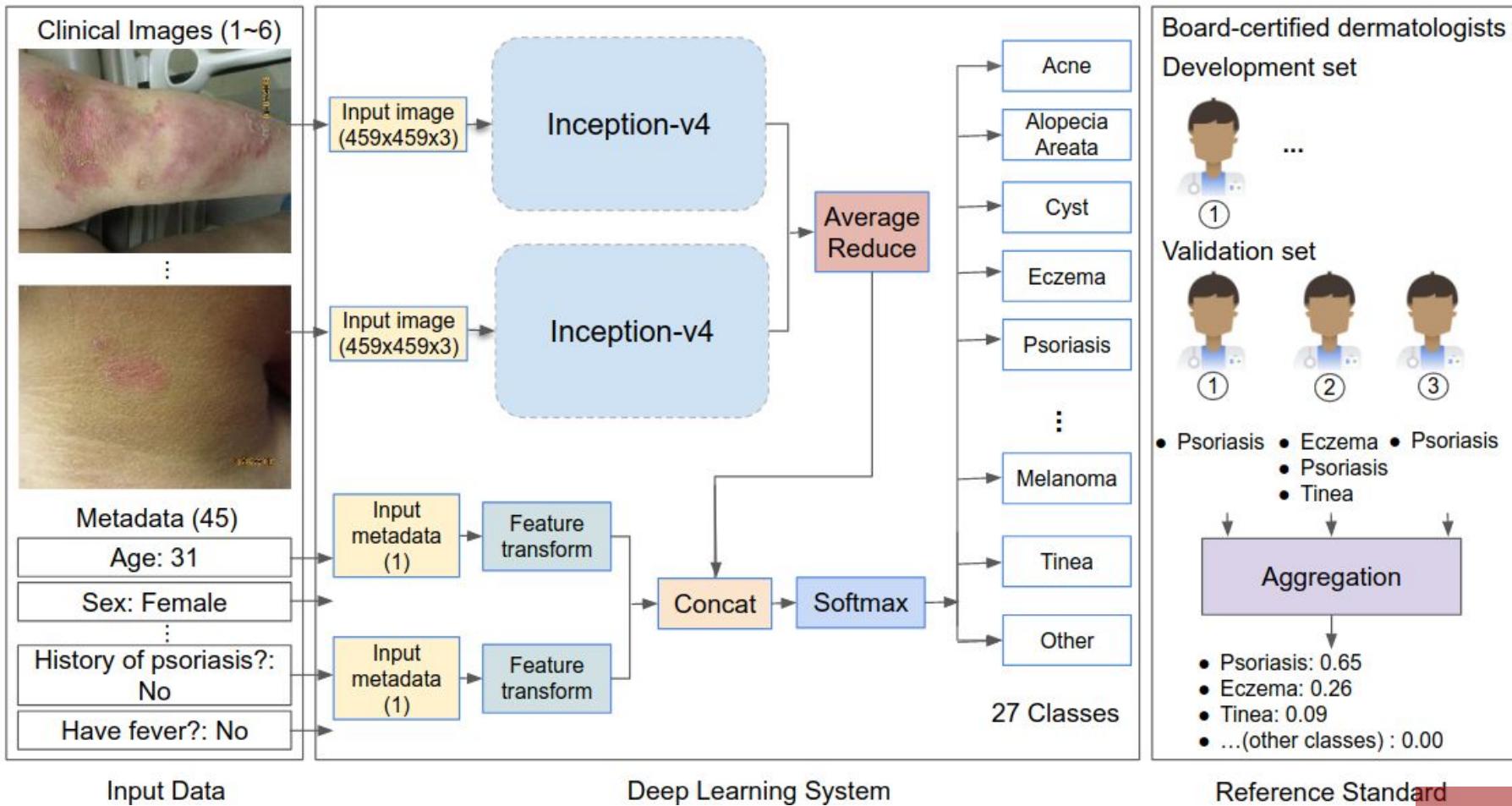
## Machine Learning for Clinical Predictive Analytics

---

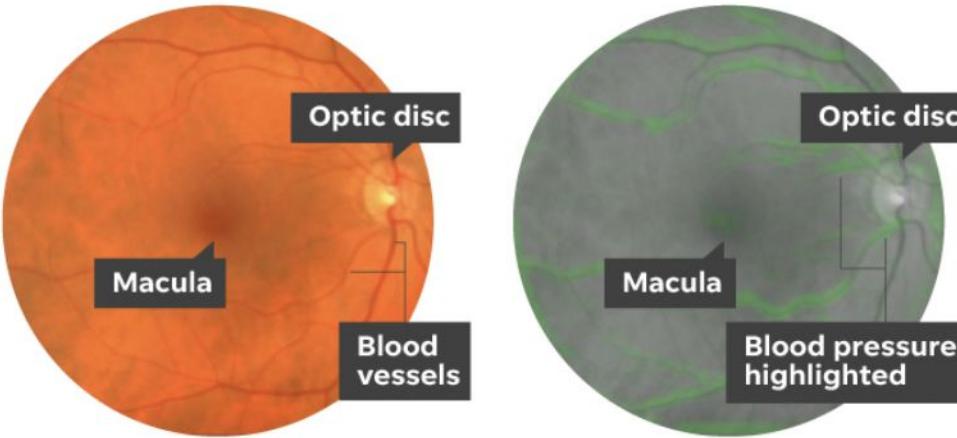
Wei-Hung Weng<sup>1</sup>

# [1] Appropriate Problem Definition

- Problem selection
  - Make a meaningful impact in patient care by ***providing actionable insights***
  - ML model should only leverage input data that are available at the time when the proposed clinical decision is being made
- Defining prediction task
  - ***Learning from humans***
    - Removing human-factor bottleneck (availability and fatigue of human raters)
    - ***E.g. fundus imaging → DR grading [Gulshan 2016], skin image → disease clf [Liu 2019]***
  - ***Enabling extraction of previously unknown insights***
    - Detection of novel signals has the potential to improve diagnosis or prognosis using cheaper and scalable modalities
    - ***E.g. use fundus imaging to predict CV risk [Poplin 2018]***
    - Must be taken to ensure that ‘novel signals’ found are not the result of confounding factors or random chance



# Google AI can predict heart problems by taking pictures of your eye



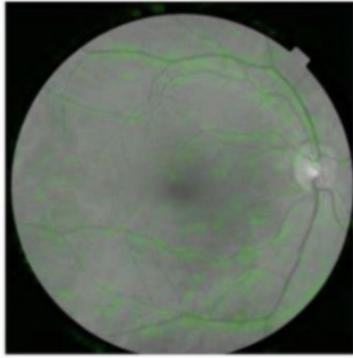
- ▶ Image of the back of the eye showing the macula (dark spot in the middle), optic disc (bright spot at the right), and blood vessels.
- ▶ Retinal image in gray, researchers can focus on blood vessels to determine the health risks associates with a patient's blood pressure.
- ▶ Images showed that each cardiovascular risk factor prediction uses a distinct pattern, such as blood vessels for blood pressure and optic disc for other predictions.

Source: Baig, Edward C. "Google Hopes AI Can Predict Heart Disease by Looking at Retinas." *USA Today*, Gannett Satellite Information Network, 19 Feb. 2018, [www.usatoday.com/story/tech/2018/02/19/google-ai-can-predict-heart-disease-looking-pictures-retina/344547002/](http://www.usatoday.com/story/tech/2018/02/19/google-ai-can-predict-heart-disease-looking-pictures-retina/344547002/)

Original



Age



Gender



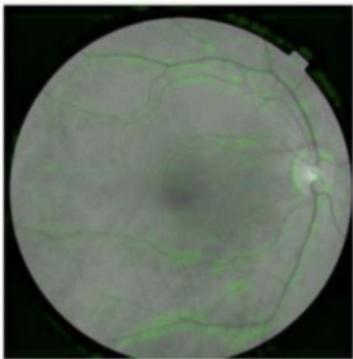
Actual: 57.6 years  
Predicted: 59.1 years

Actual: female  
Predicted: female

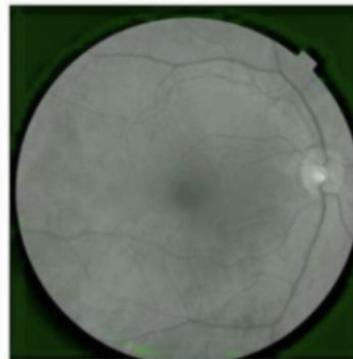
Smoking



HbA1c



BMI



Actual: non-smoker  
Predicted: non-smoker

Actual: non-diabetic  
Predicted: 6.7%

Actual:  $26.3 \text{ kg m}^{-2}$   
Predicted:  $24.1 \text{ kg m}^{-2}$

# [1] Clinical Perspective

- Cost / Risk assessment and adjustment
  - Insurance
  - Resource redistribution
- Precision / Personalized medicine
  - For oncology / rare diseases / mental disorders / ...
  - Applications
    - Clinical decision support
    - Drug discovery
    - Outcome prediction
      - Lifespan prediction / Disease progression
    - Chronic disease management
      - Early prediction of blood glucose for self-management

# [1] ML Perspective

- Risk stratification
- Causal inference
- Bias
- Time-series
- Modeling unstructured data
- Interpretability and explainability
- Disease progression modeling
- Reasoning and decision making (current ML probably not enough)

# [1] Framing into the ML Scenario

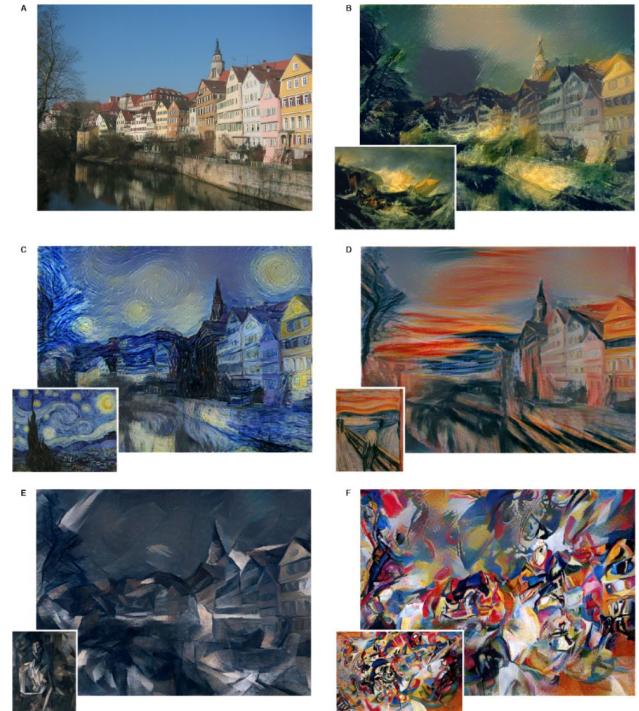
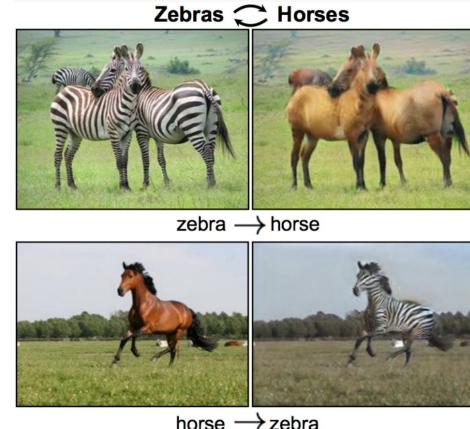
- ***Supervised learning***
  - Regression
  - Classification
    - Linear
    - Non-linear (e.g. deep learning, SVM, decision tree, ...)
  - Structured learning
- Unsupervised learning
  - Clustering
  - Dimensionality reduction
- Transfer learning
- Reinforcement learning
- ...

# Data Modality

- Structured / Tabular data (e.g. claims data, vitals, labs, demographics, ...)
- Medical imaging
- Natural language processing
- Waveform
- ...
- Multimodality

# Computer Vision / Medical Imaging

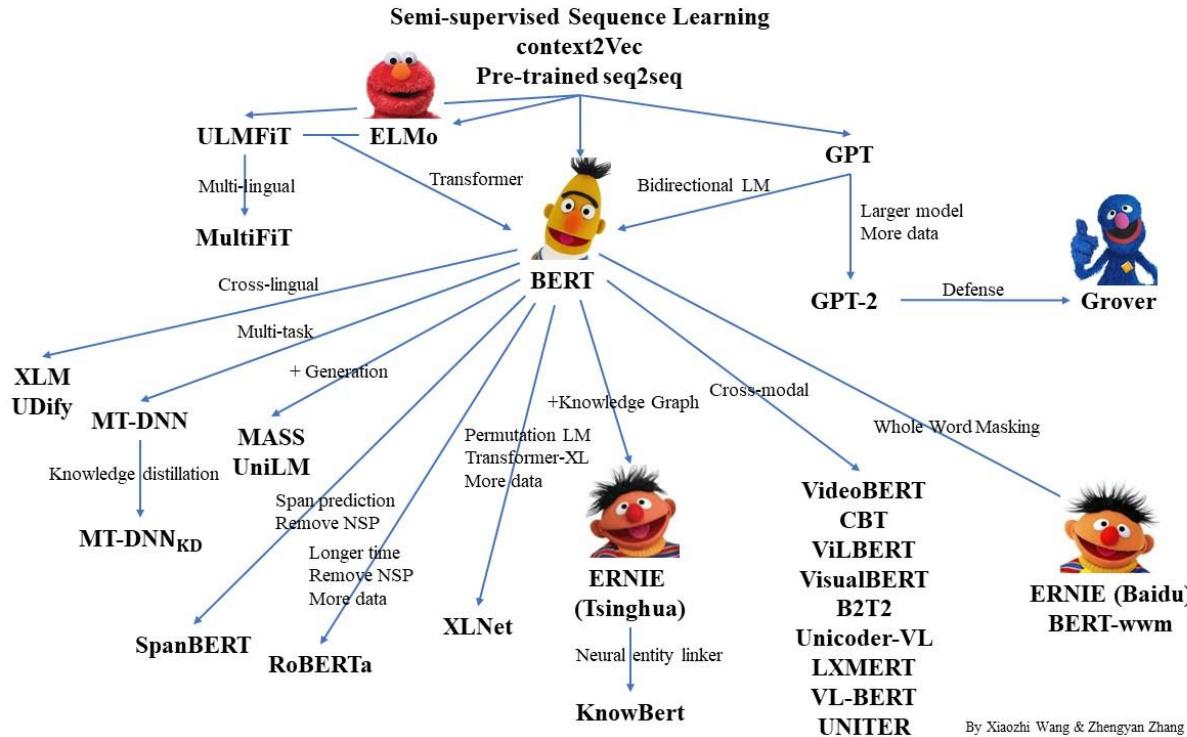
- Mainly DL-driven
- Main applications
  - Classification
  - Classification With Localization
  - Object Detection
  - Object Segmentation
  - Style Transfer
  - Colorization
  - Reconstruction
  - Super-Resolution
  - Synthesis
  - Others



# Natural Language Processing

- Embedding, encoder-decoder, attention, pretrained models / transfer learning
- [ruder.io](http://ruder.io)
- [nlpprogress.com](http://nlpprogress.com)
  - Automatic speech recognition / CCG / Common sense / Constituency parsing / Coreference resolution / Dependency parsing / Dialogue / Domain adaptation / Entity linking / Grammatical error correction / Information extraction / Language modeling / Lexical normalization / Machine translation / Missing elements / **Multi-task learning** / **Multi-modal** / Named entity recognition / Natural language inference / Part-of-speech tagging / Question answering / Relation prediction / Relationship extraction / Semantic textual similarity / Semantic parsing / Semantic role labeling / Sentiment analysis / Shallow syntax / Simplification / Stance detection / Summarization / Taxonomy learning / Temporal processing / **Text classification** / Word sense disambiguation

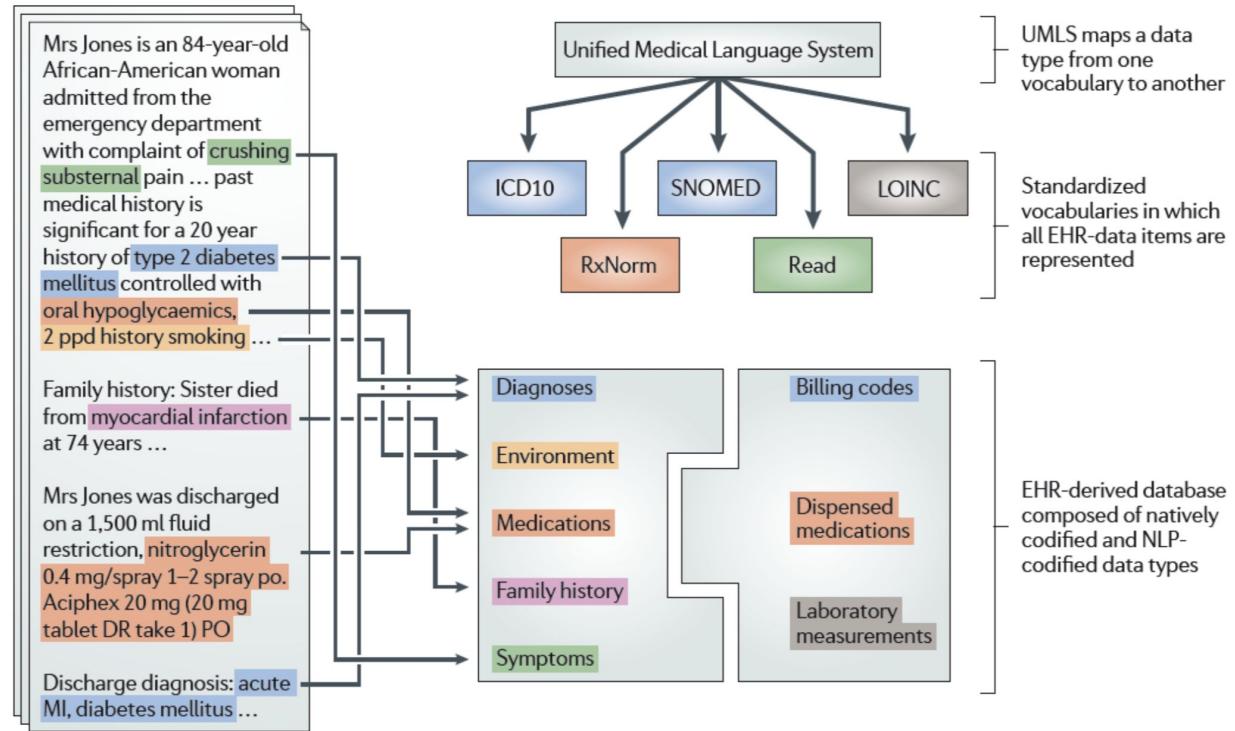
# Pretrained LM



By Xiaozhi Wang & Zhengyan Zhang @THUNLP

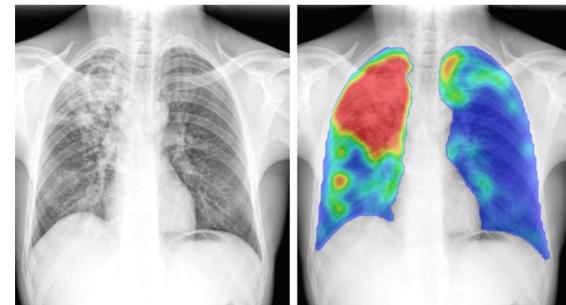
# Clinical NLP

- Smaller corpus
- Domain-specific
  - Concepts
  - Typos
  - Misspellings
  - Grammars
  - ...
- Pretrained
  - Alsentzer 2019
  - ClinicalBERT



# [1] Appropriate Problem

- **Verification / Evaluation method**
  - **Independent dataset**
  - Saliency ‘heatmaps’ for **interpretability**
  - **Ask clinicians** (who were blinded to the prediction task) for quantitative, unbiased evaluation
- Data availability
  - Lack of digitization (pathology slides)
  - Inaccessible because of patient privacy or commercial concerns
  - Lacking because the disease of interest is too rare



# [2] Curating Datasets

- **Choose a correct cohort**
- **Data split**
  - **Should be 'clean' with respect to patients**
  - Merging from multiple sources / patient-level overlap → image similarity to detect duplication
- The size of the validation set, information from **clinical trials** may be helpful
  - **Power calculations** helps determine the sample size required to confidently evaluate the model performance
  - All primary / secondary analyses should be pre-specified, avoiding 'post-hoc' analysis
  - Only perform exploratory analyses on the training set, and validate the hypotheses on the validation set
- Class imbalance
  - Data augmentation
  - Additional steps to ensure proper model calibration or adjustments in the evaluation metric



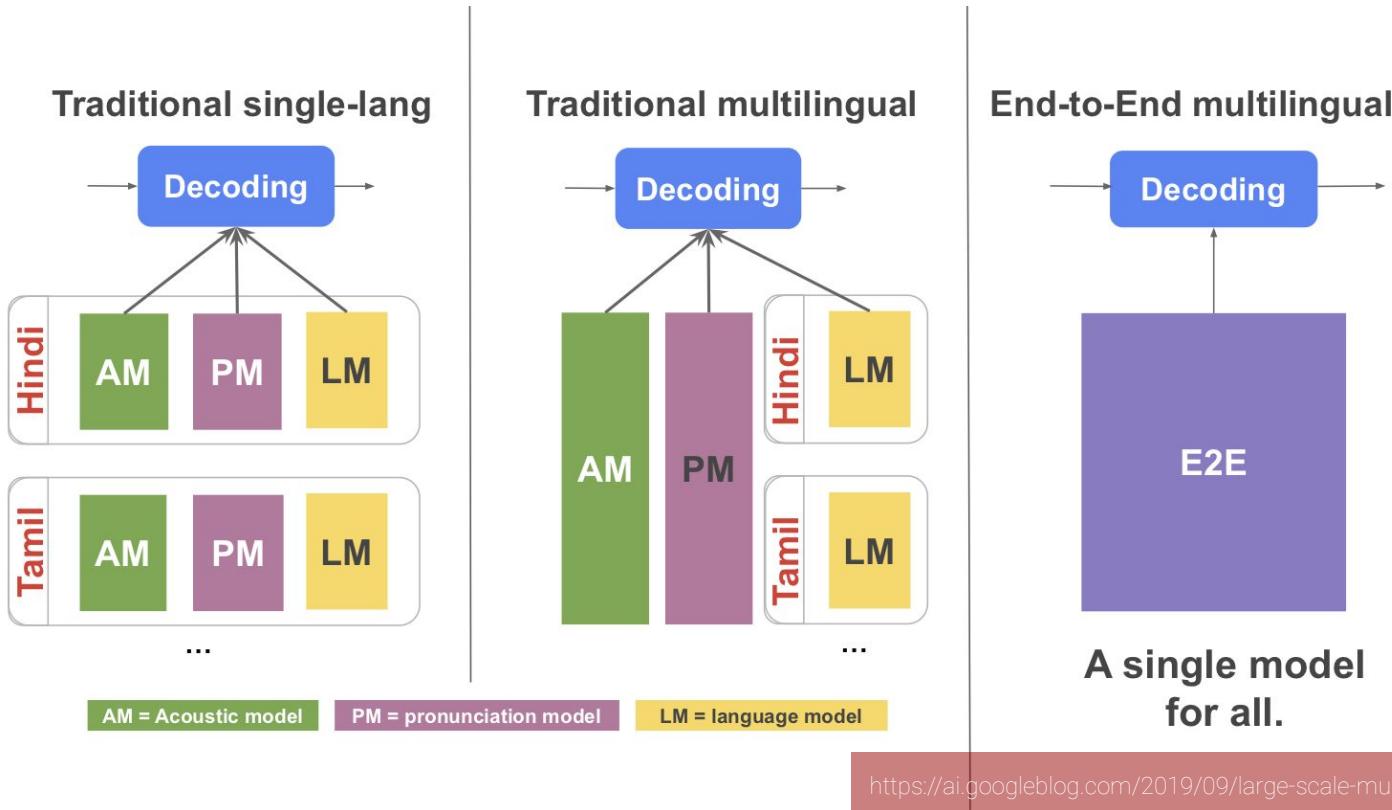
# [2] Curating Datasets

- Data Quality
  - The notion of 'good quality' may be disease-specific, for example the same fundus image may be of sufficient quality for assessing glaucomatous nerve head features but not for diabetic retinopathy
- ***Reference / Ground truth***
  - Determination of the reference truth often involves subjective judgement, introducing systematic errors, random errors or both
  - **Adjudication** by a panel of experts may be helpful but can be slow and expensive
  - Adjudication of only a subset of the data
    - Testing dataset for final evaluation
    - Validation dataset for hyperparameter optimization during the modeling

# [3] Modeling

- Several considerations for model architecture design
  - Data modality and volume, model interpretability, model inference time, balancing model overfitting and underfitting, ...
- End-to-end? ***Data size & the value of intermediate outputs***
  - Better when large datasets are available and if the final performance is the primary metric of interest (not the case in healthcare)
- ***Decomposing the model***
  - Intermediate output may be useful (e.g. interpretability of the prediction)
  - Healthcare data can differ substantially across data sources → easier generalization
  - Models with a large number of parameters usually have better predictive performance
    - May require minutes to hours for inference, which may hinder intraoperative usage where time is of the essence
    - ***Real-time need***

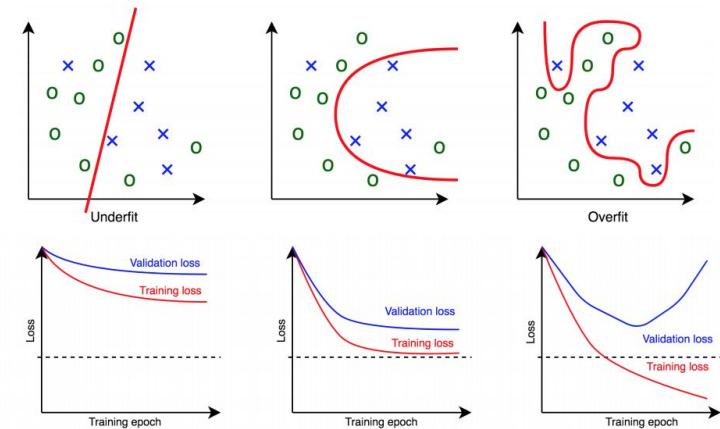
# An Example of Decomposition vs. End-to-end



# [3] Model Generalizability

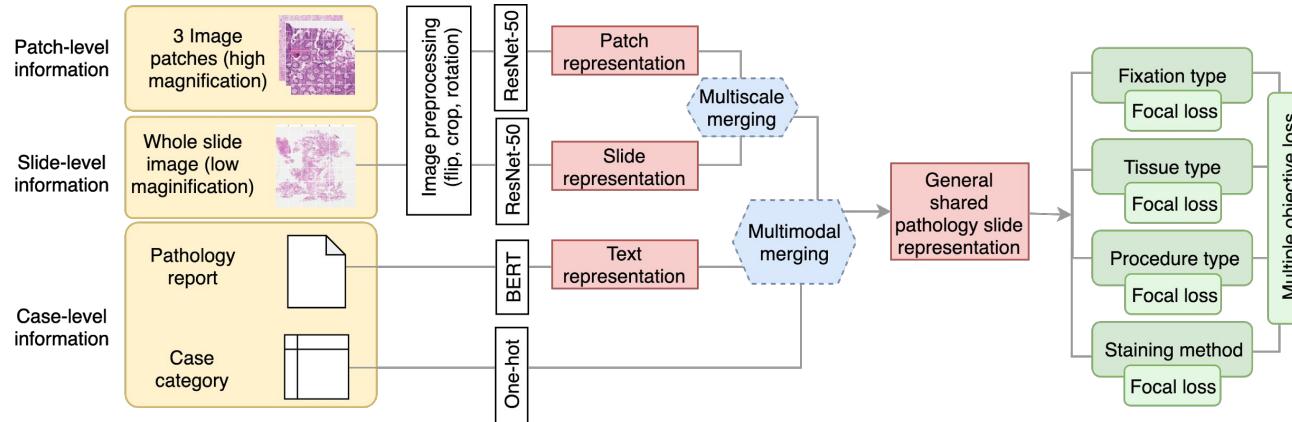
- Bias-variance trade-off
  - Balancing model underfitting/overfitting
- **Regularization**
  - Data augmentation: particularly useful in the data-limited healthcare regime
  - Inject prior knowledge
  - Perturbations that may not correspond to real-world changes may still be helpful in learning
- **Train-validation-test split must be carefully preserved**
  - Any violation of train-validation-test hygiene can result in ungeneralizable performance

	Training error	Validation error	Approach
High bias	High	Low	Increase complexity
High variance	Low	High	Decrease complexity Add more data



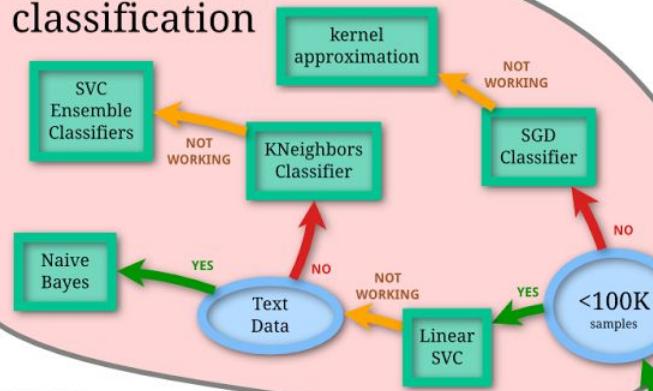
# [3] An Example

- Learning general pathology data representation [Weng 2019]
  - Limited pathology data, extremely class imbalance
  - Multimodal (image + text + structured data), multitask, data augmentation, neural network regularization, prior knowledge from pretraining networks, loss function design, ...

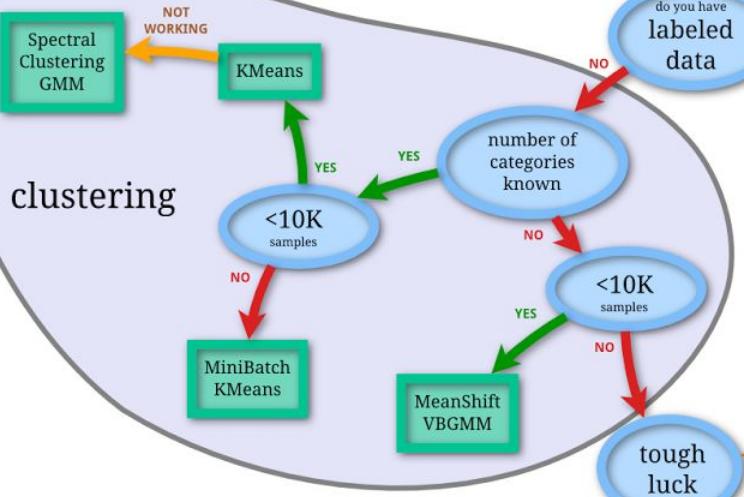


# scikit-learn algorithm cheat-sheet

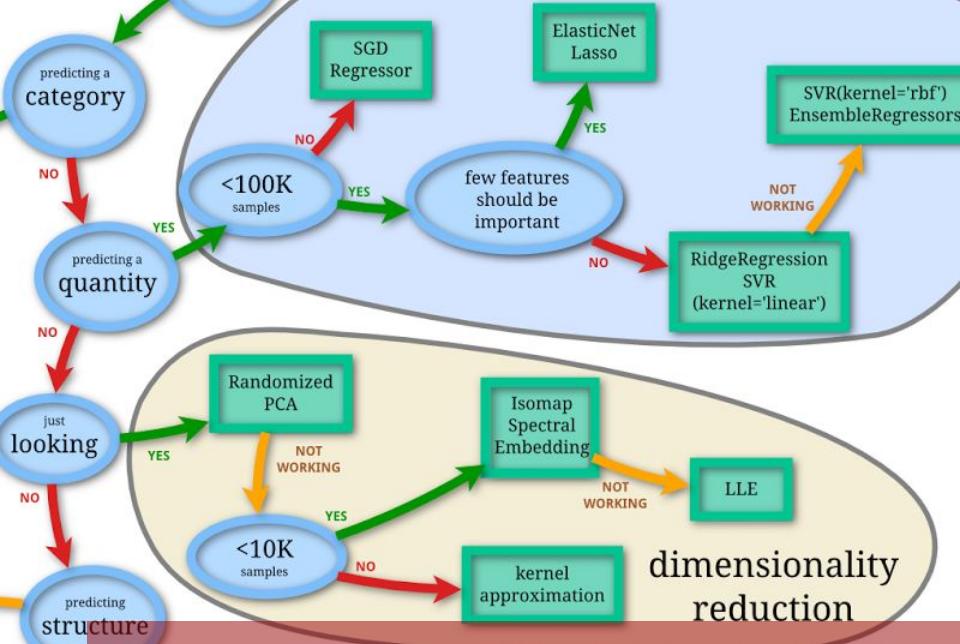
## classification



## clustering



## regression

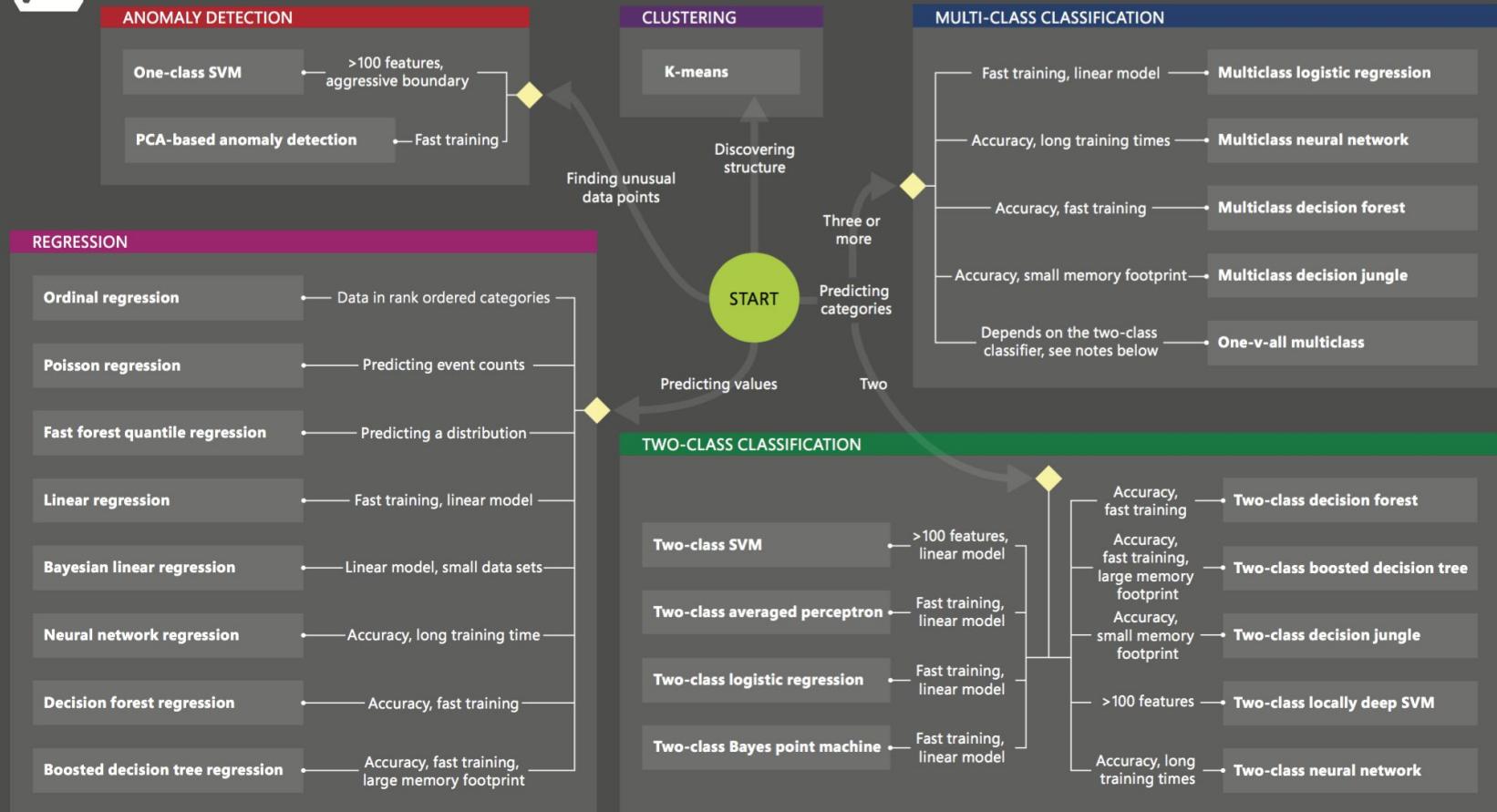


## dimensionality reduction



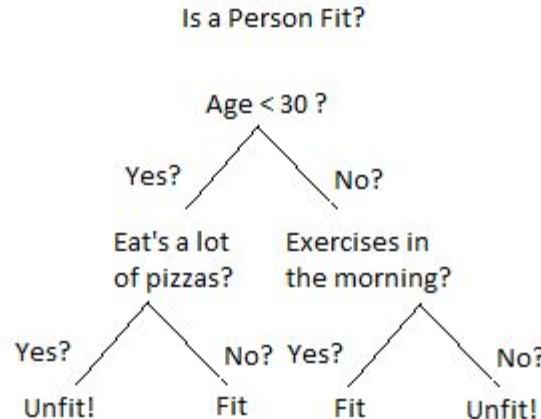
# Microsoft Azure Machine Learning: Algorithm Cheat Sheet

This cheat sheet helps you choose the best Azure Machine Learning Studio algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.



# [3] Modeling

- **Always start from simpler algorithms / models**
  - e.g. LR, then random forest
- Deep learning is not always a good approach
  - Very limited data (transfer learning might help)
  - Interpretation is not easy



```
Call:  
lm(formula = fin_loss ~ num_people + num_records + per_sensitive +  
+ dys_impact + dys_detect + cost_controls, data = hw2)  
  
Residuals:  
    Min      1Q  Median      3Q     Max  
-33.261 -4.734 -0.108  5.326 23.490  
  
Coefficients:  
              Estimate Std. Error t value Pr(>|t|)  
(Intercept) 21.1883422 3.2315604 6.557 1.39e-10 ***  
num_people   -0.0031807 0.0994055 -0.032 0.9745  
num_records   0.4291799 0.1010031 4.249 2.57e-05 ***  
per_sensitive 0.0216817 0.0133377 1.626 0.1047  
dys_impact    0.0003192 0.0007069 0.452 0.6518  
dys_detect    0.0019461 0.0011590 1.679 0.0938 .  
cost_controls 0.1096839 0.0587405 1.867 0.0625 .  
---  
Signif. codes:  
0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
  
Residual standard error: 8.175 on 493 degrees of freedom  
Multiple R-squared: 0.2264, Adjusted R-squared: 0.2169  
F-statistic: 24.04 on 6 and 493 DF, p-value: < 2.2e-16
```

# [4] Evaluating Model Performance

- Evaluation metrics should be consistent with the ones used in the community!
- ***Discrimination metric***
  - **Threshold** selection tends to play a critical role in healthcare because clinical applications commonly involve binary decisions
  - High sensitivity for screening / high specificity for diagnosis
  - Resource constraints (time, labor effort, cost limitations for screening)
- ***Calibration metric***
  - Evaluate how well the predicted probabilities match the actual probabilities
  - Although under-reported, calibration metrics (e.g., the Hosmer-Lemeshow statistic) are crucial for real-world use because these probabilities are used for expected cost-benefit analysis
- Validation should be done using large, heterogeneous datasets to ensure generalization to diverse patient populations

# [4] Evaluating Model Performance

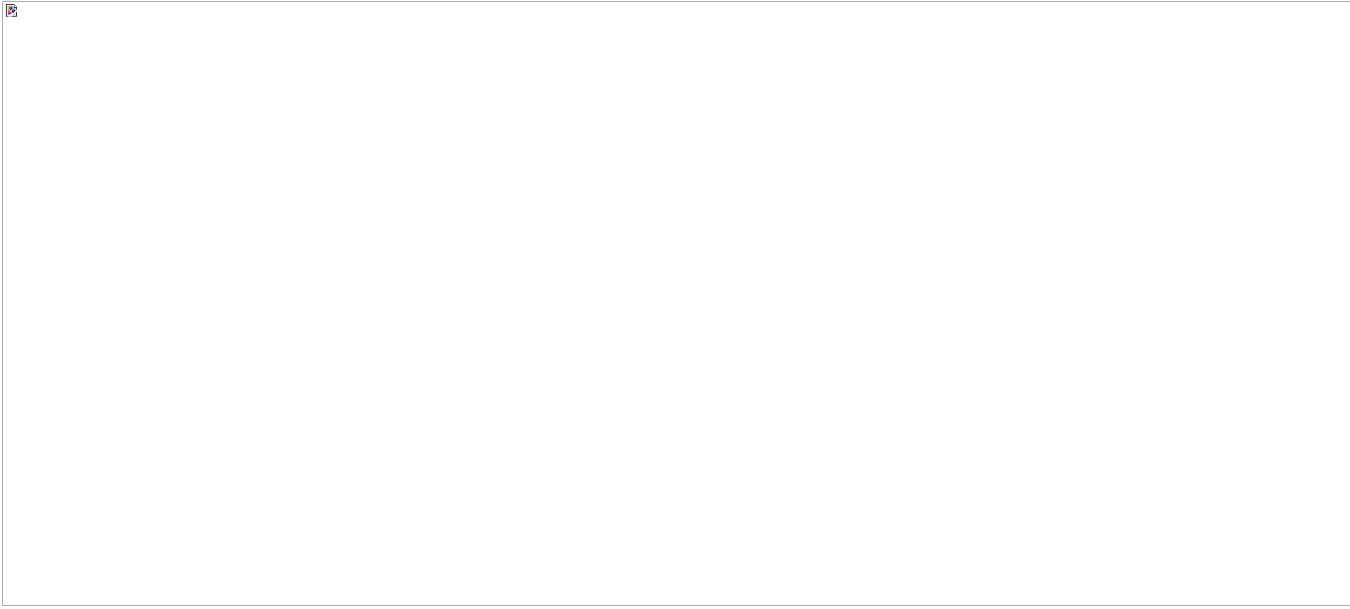
- Subgroup analysis /cluster analysis
- Sensitivity analysis
- Data augmentation for class imbalance
  - Validation set collected may have a different distribution of disease subtypes relative to real-world populations → evaluation should be adjusted according to realistic prevalence distributions
- ***Baseline comparison***
  - Comparison with a 'human baseline'
  - Comparison with a baseline (e.g. LR) ***based on variables that are readily available in the clinic (e.g. APACHE, KDIGO, CHADS2, ...)*** may be useful to evaluate the added value of the proposed novel association

# [3+4] Keys

- ***Loss (Objective) function***
- ***Metrics***
- ***Quantitative***
- ***Qualitative***
- ***Baseline***
- ***Error analysis***
- ***Ablation analysis***

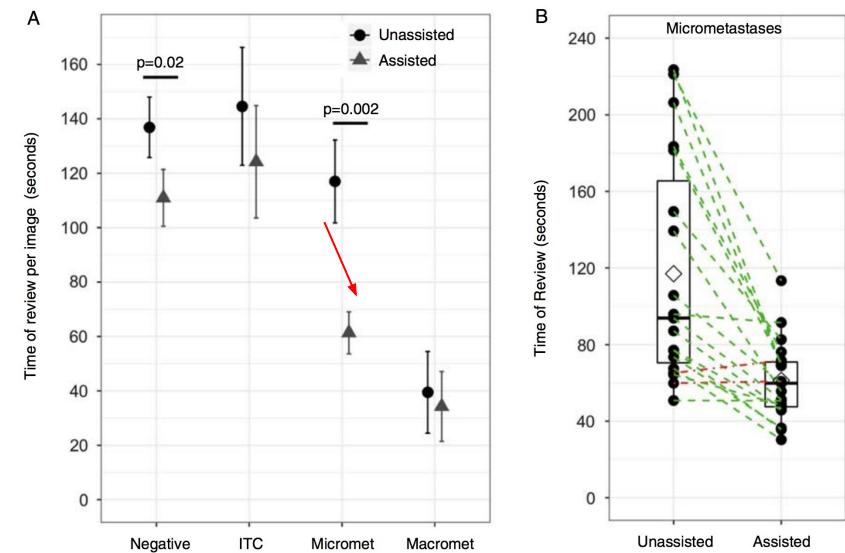
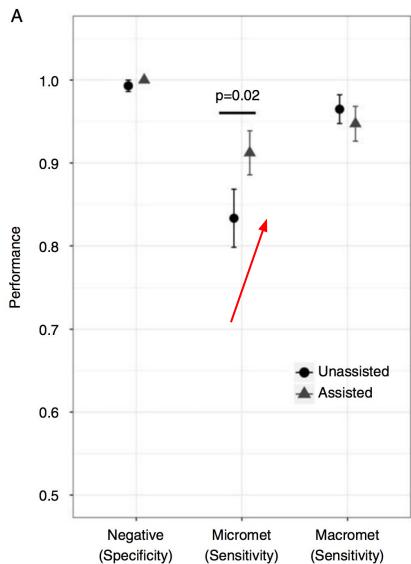
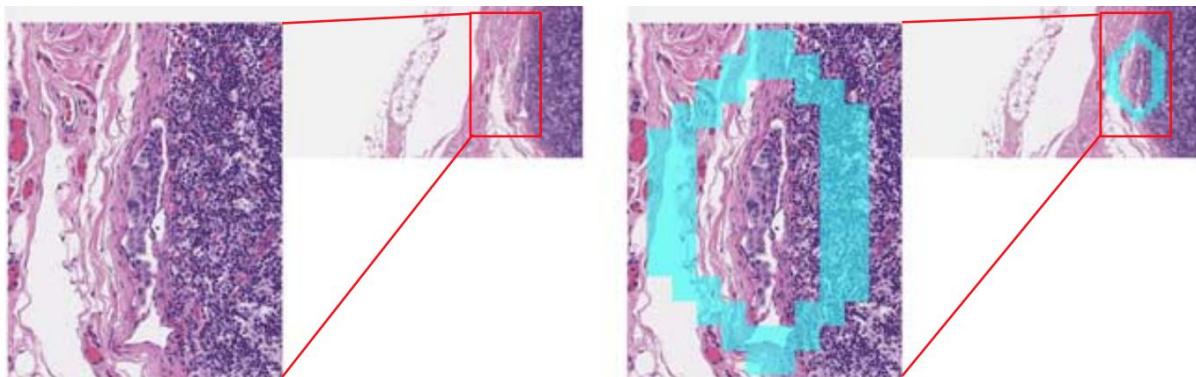
Task	Error type	Loss function	Note
Regression	Mean-squared error	$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$	Easy to learn but sensitive to outliers (MSE, L2 loss)
	Mean absolute error	$\frac{1}{n} \sum_{i=1}^n  y_i - \hat{y}_i $	Robust to outliers but not differentiable (MAE, L1 loss)
Classification	Cross entropy = Log loss	$-\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$ $= -\frac{1}{n} \sum_{i=1}^n p_i \log q_i$	Quantify the difference between two probability distributions
	Hinge loss	$\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i \hat{y}_i)$	For support vector machine
	KL divergence	$D_{KL}(p  q) = \sum_i p_i (\log \frac{p_i}{q_i})$	Quantify the difference between two probability distributions

		Predicted		
		True	False	
Actual	True	True positive (TP) Type II error	False negative (FN) Type II error	Recall = Sensitivity = $\frac{TP}{TP+FN}$
	False	False positive (FP) Type I error	True negative (TN)	Specificity = $\frac{TN}{TN+FP}$
		Precision = $\frac{TP}{TP+FP}$		Accuracy = $\frac{TP+TN}{TP+TN+FP+FN}$ F1 = $\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$



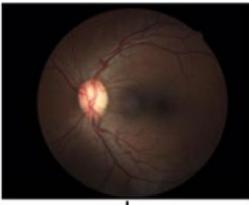
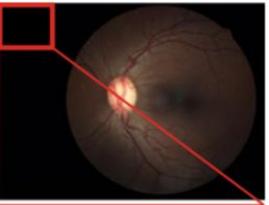
# [5] Clinical Impact

- Steiner 2018



# [5] Clinical Impact

- A performant model alone is insufficient to create clinical impact
- The system has to be designed to be useful even in cases of failure such as false positives
- ‘Pre-diagnosis’ to pre-screen images and draw attention to images of high or uncertain risk, ‘peri-diagnosis’ to highlight lesions on the image during the regular image grading process, or ‘post-diagnosis’ to resurface images that may have been graded wrongly due to inexperience or fatigue
- ***User trust → human-in-the-loop development***
  - The degree of reliance can be measured by the user-model agreement rate, comparing that with the ML prediction accuracy
- ***‘Alarm fatigue’***

Problem selection	Data collection	ML development	Validation	Assessment of impact	Deployment and monitoring
  <p>Referable diabetic retinopathy?</p> <p>Gulshan et al. (2016) Ting et al. (2017)</p>	<p>Development dataset</p> <p>128,175 images</p> <p>Validation dataset</p> <p>11,711 images</p> <p>Gulshan et al. (2016)</p>		<p>Retrospective</p> <p>Sensitivity: 0.90 Specificity: 0.98</p> <p>Gulshan et al. (2016)</p> <p>Prospective</p> <p>Sensitivity: 0.87 Specificity: 0.91</p> <p>Abramoff et al. (2018)</p>	 <p>Model predictions</p> <ul style="list-style-type: none"> <li>None</li> <li>Mild</li> <li>Moderate</li> <li>Severe</li> <li>Proliferative</li> </ul> <p>40% reduction in false negatives</p> <p>Sayres et al. (2018)</p>	<p>Future work</p>

# tl;dr

- How to develop machine learning models for healthcare [Chen 2019]
- Machine Learning for Clinical Predictive Analytics [Weng 2019]

comment

## How to develop machine learning models for healthcare

Rapid progress in machine learning is enabling opportunities for improved clinical decision support. Importantly, however, developing, validating and implementing machine learning models for healthcare entail some particular considerations to increase the chances of eventually improving patient care.

Po-Hsuan Cameron Chen, Yun Liu and Lily Peng

---

## Machine Learning for Clinical Predictive Analytics

---

Wei-Hung Weng<sup>1</sup>

# MIT CRITICAL DATA



# Secondary Analysis of Electronic Health Records

EXTRAS ONLINE



Springer Open

# Resources

- Practical
  - TensorFlow, PyTorch, Keras, Scikit-learn documents, guide, tutorials
  - **Google Machine Learning Crash Course**
  - ***Do more projects!***
- Theory and mathematics
  - Coursera Machine Learning (Andrew Ng)
  - Deep learning book
  - PMRL
  - Stanford CS224n: Natural Language Processing with Deep Learning
  - Stanford CS231n: Convolutional Neural Networks for Visual Recognition

# Seeking Papers

- Clinical-driven (journals)
  - JAMA
  - Nature Medicine
  - npj Digital Medicine
  - JAMIA
  - AMIA (conference)
- ML/AI-driven (conferences / workshops)
  - MLHC
  - NeurIPS ML4H Workshop
  - NAACL Clinical NLP Workshop
  - MICCAI
  - NeurIPS, ICML, ICLR, KDD, AAAI, ACL, NAACL, CVPR, ...

# tl;dr

- Caveats and tips for model development [Chen 2019, Weng 2019]
  - Problem definition →
  - Data curation →
  - Model development →
  - Validation →
  - Assessment of clinical impact →
  - Deployment and monitoring
- Wei-Hung Weng [[ckbjimmy@mit.edu](mailto:ckbjimmy@mit.edu)]

# Framing Your Clinical Questions into ML Tasks

- Clinical problems?
- Healthcare workflow?
- Socioeconomic problems?
- ?

# What's Special in Bhutan?

- Clinical problems?
- Healthcare workflow?
- Socioeconomic problems?
- ?

# Clinical Problem → ML Task

[Problem] ICU mortality prediction [Ghassemi KDD 2014]

[ML task] Prediction (supervised)

[Data] MIMIC-III structured data + notes

[ML model] Unupservised LDA + structured SVM

[Baseline] ***Admission baseline*** / retrospective topic / time-varying model

[Evaluation] AUC, Topic ***post-hoc analysis***

# Clinical Problem → ML Task

[Problem] Want to know better subtyping in the cohort, e.g. autism?

[ML task] Cluster analysis (unsupervised learning)

[Data] Data with features

[ML model] ***K-means***, PCA, autoencoder

[Evaluation and validation] ***Visualization, post-hoc explanation***, downstream tasks  
(then you need annotation and a model for downstream task)

# Clinical Problem → ML Task

[Problem] What may happen in the next clinical visit? [Choi MLHC 2016]

[ML task] Disease progression modeling, sequential prediction

[Data] Clinical visit claims data (input codes, output ICDs)

[ML model] RNN

[Baseline] **LR**, most frequent, last visit

[Evaluation] Recall@30

# Clinical Problem → ML Task

[Problem] How to know if the patient may readmit? [Caruana 2015, Xiao 2018]

[ML task] Prediction, supervised learning

[Data] Labs, notes, history of hospitalized patients

[ML model] Generalized additive models [Caruana 2015], topic modeling + RNN  
[Xiao 2018]

[Baseline] / **word2vec+LR**, GRU, ... [Xiao 2018]

[Evaluation and validation] AUC, **risk score** [Caruana 2015] / Acc, AUC, PR-AUC,  
**clustering** [Xiao 2018]

# Clinical Problem → ML Task

[Problem] How to find an optimal treatment for my patient's condition?

[Komorowski Nat Med 2018]

[ML task] Optimal sequential decision making (Reinforcement learning)

[Data] MIMIC-III, eICU (time-series structured data, vitals, labs, medical orders, ...)

[ML model] Policy iteration

[Baseline] Doctors' action in EHR

[Evaluation and validation] Value (this is challenging)

# Clinical Problem → ML Task

[Problem] How to detect the severe (referrable) diabetic retinopathy? [Gulshan JAMA 2016]

[ML task] Image classification (supervised learning)

[Data] DR images with annotations

[ML model] Inception V3

[Baseline] ***Ophthalmologists' manual grading***

[Evaluation and validation] ***Sensitivity and specificity***

# Clinical Problem → ML Task

[Problem] Decide the discharge ICD code based on the clinical notes? [Mullenbach NAACL 2018]

[ML task] NLP + multilabel classification

[Data] MIMIC-III

[ML model] CNN + attention + classifier

[Evaluation and validation] Precision@K, AUC, ***physician evaluation***

# Clinical Problem → ML Task

[Problem] How to translate clinical jargons to layman understandable language?

[Weng KDD 2019]

[ML task] Machine translation / text style transfer (unsupervised or supervised if we have pairs)

[Data] MIMIC-III (some professional-level, some layman-level)

[ML model] Cross-aligned representation + statistical machine translation

[Baseline] ***Dictionary-based translation***

[Evaluation and validation] Precision@K, ***mean opinion score***

# More Problems from ICU

## Classification

- Fit paCO<sub>2</sub> to pH / fit paCO<sub>2</sub> and HCO<sub>3</sub> to pH
- Classify ICU mortality using HCO<sub>3</sub> min/median & paCO<sub>2</sub> 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 HCO<sub>3</sub>/paCO<sub>2</sub> or vital signs who did die or did not die during ICU stay
- Clustering vital signs and risk of being initiated on mechanical ventilation

# More Problems from Respiratory Therapy

# Hands-on Exercise

- [https://github.com/ckbjimmy/2019\\_bt](https://github.com/ckbjimmy/2019_bt)
  - Open → "File" -> "Save a copy in Drive..."
- Requirement
  - Google account
  - Python + some ML knowledge
- scikit-learn / tensorflow-keras
- Two datasets
  - ICU structured data / Breast cancer prediction data
- Supervised learning (classification)
  - When you have some labeled data
  - Given features, predict malignancy
  - ML general approaches and modeling, missing data imputation, normalization, important feature identification, ...