AMDIS Machine Learning Survey and Update

Michael Wang, MD Xinran (Leo) Liu, MD University of California, San Francisco slı.do

Join at slido.com #Y182

What barriers have you run into with deploying AI

Disclosures

Leo Liu:

- Clinical Advisor Advanced Clinical L.L.C
 - Based at Google Brain
- Fiancée
 - Machine Learning Scientist at Amazon

Mike Wang:

• Clinical Consultant - Commure, Inc

Survey Results

Survey Results

PICO for ML Papers

Survey Results

PICO for ML Papers

Discussing Barriers to ML

WHAT IS MACHINE LEARNING

"Machine learning is the science of getting computers to act without being explicitly programmed."

– Andrew Ng, Stanford University



Taking advantage of computational power

Survey Results

PICO for ML Papers

Discussing Barriers to ML

What best describes your general exposure to machine learning? (Select one)



How would you characterize your enthusiasm towards Artificial Intelligence (AI) and Machine Learning (ML)? (Select one)

20 responses



In an ideal world, who in your organization should be most responsible for decisions surrounding ML algorithms in the clinical setting? (Select one) ²⁰ responses



How comfortable do you feel when evaluating ML applications? 20 responses



Survey Results

PICO for ML Papers

Discussing Barriers to ML

Survey Results

PICO for ML Papers

Discussing Barriers to ML

slı.do

Join at slido.com #Y182

informatics fellow with a CS degree wants to predict readmissions in patients getting discharged from the hospital. The model he builds will most likely be a model

Supervised - You supply the "right answers" to the model and it learns how to come up these answers.

Supervised - You supply the "right answers" to the model and it learns how to come up these answers.

Unsupervised - There are no right answers and the models are trying find "patterns" in the data.

Supervised - You supply the "right answers" to the model and it learns how to come up these answers.

Unsupervised - There are no right answers and the models are trying find "patterns" in the data.



Supervised - You supply the "right answers" to the model and it learns how to come up these answers.

Unsupervised - There are no right answers and the models are trying find "patterns" in the data.

Reinforcement - Train the model by "rewarding" it for correct actions

Example: training a robot to grab a ball, training a computer to play chess

A first question

Your ambitious first year informatics fellow with a CS degree wants to predict readmissions in patients getting discharged from the hospital.

The model he builds will most likely be a _____ model.

- a) **Supervised -** Your fellow will be labeling which patients count as readmissions.
- b) Unsupervised
- c) Perfect
- d) Imaginary

P - Patient/Problem

P - Patient/Problem

I - Intervention

- P Patient/Problem
- I Intervention
- C Comparison

- P Patient/Problem
- I Intervention
- C Comparison
- 0 Outcomes

- **P Patient/Problem**
- I Intervention
- C Comparison
- 0 Outcomes

- Patient Population
 - Does this model apply to my patient population?

• Patient Population

- Does this model apply to my patient population?
 - Example: MANY models are build off of the publicly available MIMIC dataset (ICU data from a large hospital in Boston)

• Patient Population

- Does this model apply to my patient population?
 - Example: MANY models are build off of the publicly available MIMIC dataset (ICU data from a large hospital in Boston)
- Problem
 - For supervised models: How much do I trust the labelling used to build the model?

• Patient Population

- Does this model apply to my patient population?
 - Example: MANY models are build off of the publicly available MIMIC dataset (ICU data from a large hospital in Boston)

• Problem

- For supervised models: How much do I trust the labelling used to build the model?
 - Example: Some predictive models for sepsis use ICD10 sepsis codes as their "truth" label (this is only 50% sensitive)

A Real Paper!

Cardiologist-Level Arrhythmia Detection with Convolutional Neural Networks

Pranav Rajpurkar* Awni Y. Hannun* Masoumeh Haghpanahi Codie Bourn Andrew Y. Ng PRANAVSR @CS.STANFORD.EDU AWNI@CS.STANFORD.EDU MHAGHPANAHI@IRHYTHMTECH.COM CBOURN@IRHYTHMTECH.COM ANG@CS.STANFORD.EDU

A Real Paper!

Testing

We collect a test set of 336 records from 328 unique patients. For the test set, ground truth annotations for each record were obtained by a committee of three boardcertified cardiologists; there are three committees responsible for different splits of the test set. The cardiologists discussed each individual record as a group and came to a consensus labeling. For each record in the test set we also collect 6 individual annotations from cardiologists not participating in the group. This is used to assess performance of the model compared to an individual cardiologist.

How is the "true arrhythmia" being determined in the training set:

- a. A single cardiologist read in the EHR
- b. Majority from 6 different boarded cardiologists
- c. ICD10 codes
- d. Consensus committee of 3 boarded cardiologists
- e. Medical Students
- f. Magic 8 ball

slı.do

Join at slido.com #Y182

there are three committees responsible for different splits of the test set. The cardiologists discussed each individual record as a group and came to a consensus labeling. For each record in the test set we also collect 6 individual annotations

from cardiologists not
A Real Paper!

Testing

We collect a test set of 336 records from 328 unique patients. For the test set, ground truth annotations for each record were obtained by a committee of three boardcertified cardiologists; there are three committees responsible for different splits of the test set. The cardiologists discussed each individual record as a group and came to a consensus labeling. For each record in the test set we also collect 6 individual annotations from cardiologists not participating in the group. This is used to assess performance of the model compared to an individual cardiologist. How is the "true arrhythmia" being determined in the training set:

- a. A single cardiologist read in the EHR
- b. Majority from 6 different boarded cardiologists
- c. ICD10 codes
- d. Consensus committee of 3 boarded cardiologists
- e. Medical Students
- f. Magic 8 ball

PICO for Machine Learning

- **P Patient/Problem**
- I Intervention
- C Comparison
- 0 Outcomes

PICO for Machine Learning

- P Patient/Problem
- I Intervention
- C Comparison
- 0 Outcomes

- Actionable Insight
 - Will knowing a model's predictions change provider/system behavior

• Actionable Insight

- Will knowing a model's predictions change provider/system behavior
 - Example: Showing mortality predictions to your discharging doctors.

• Actionable Insight

- Will knowing a model's predictions change provider/system behavior
 - Example: Showing mortality predictions to your discharging doctors.
- Integration with workflow
 - Applications need defined ROI and incorporation into workflow.

• Actionable Insight

- Will knowing a model's predictions change provider/system behavior
 - Example: Showing mortality predictions to your discharging doctors.

• Integration with workflow

- Applications need defined ROI and incorporation into workflow.
 - Example: Showing mortality predictions in the discharge activity that then auto-populate a care transitions referral to reduce readmissions.

PICO for Machine Learning

- P Patient/Problem
- I Intervention
- **C** Comparison
- 0 Outcomes

• Static vs Dynamic

• Will the model change over time or not?

• Static vs Dynamic

- Will the model change over time or not?
 - Static Easier to implement, but often degrades over time.

• Static vs Dynamic

- Will the model change over time or not?
 - Static Easier to implement, but often degrades over time.
 - Dynamic Much harder to implement because this needs large data feeds.

• Static vs Dynamic

- Will the model change over time or not?
 - Static Easier to implement, but often degrades over time.
 - Dynamic Much harder to implement because this needs large data feeds.
 - Consider review periods: le every 2-3 years (no consensus)

PICO for Machine Learning

- P Patient/Problem
- I Intervention
- C Comparison
- **O Outcomes**

PICO for Machine Learning

- P Patient/Problem
- I Intervention
- C Comparison
- **O Outcomes**

• AUC

- AUC
- Accuracy

- AUC
- Accuracy
- Sensitivity/Specificity

• Binary Classifiers (le sepsis vs no sepsis)

- AUC
- Sensitivity/Specificity
- F1 statistic
- Accuracy

Binary Classifiers (le sepsis vs no sepsis)

- AUC
- Sensitivity/Specificity
- F1 statistic
- Accuracy

• Non-binary Classifiers (le low, mid, high risk for PE)

- Accuracy
- o Kappa

Binary Classifiers (le sepsis vs no sepsis)

- AUC
- Sensitivity/Specificity
- F1 statistic
- Accuracy

Non-binary Classifiers (le low, intermediate, high risk for PE)

- Accuracy
- Kappa
- Regressions (le what is the % chance of readmission)
 - R-squared
 - MSE

Binary Classifiers (le sepsis vs no sepsis)

- AUC
- Sensitivity/Specificity
- F1 statistic
- Accuracy

Non-binary Classifiers (le low, mid, high risk for PE)

- Accuracy
- Kappa
- Regressions (le what is the % chance of readmission)
 - R-squared
 - MSE
- Which metric(s) do you choose?

BMJ Open Multicentre validation of a sepsis prediction algorithm using only vital sign data in the emergency department, general ward and ICU

Qingqing Mao,¹ Melissa Jay,¹ Jana L Hoffman,¹ Jacob Calvert,¹

Receiver Operator Curve for severe sepsis (AUC: 0.87)



BMJ Open Multicentre validation of a sepsis prediction algorithm using only vital sign data in the emergency department, general ward and ICU

Qingqing Mao,¹ Melissa Jay,¹ Jana L Hoffman,¹ Jacob Calvert,¹

Receiver Operator Curve for severe sepsis (AUC: 0.87)

Sensitivity ~ 80%, Specificity ~80%



BMJ Open Multicentre validation of a sepsis prediction algorithm using only vital sign data in the emergency department, general ward and ICU

Qingqing Mao,¹ Melissa Jay,¹ Jana L Hoffman,¹ Jacob Calvert,¹

Receiver Operator Curve for severe sepsis (AUC: 0.87)

Sensitivity ~ 80%, Specificity ~80%

Assuming prevalence ~10%



BMJ Open Multicentre validation of a sepsis prediction algorithm using only vital sign data in the emergency department, general ward and ICU

Qingqing Mao,¹ Melissa Jay,¹ Jana L Hoffman,¹ Jacob Calvert,¹

Receiver Operator Curve for severe sepsis (AUC: 0.87)

Sensitivity ~ 80%, Specificity ~80%

Assuming prevalence ~10%

PPV: 31%



APPLE WATCH

96% sensitive, 98% specific – sounds great!

Prevalence of afib for population under 55 = 0.1% What do you think is positive predictive value?

a) 1%

b) 5%

c) 20%

d) 50%

e) 90%



slı.do

Join at slido.com #Y182

The Apple Watch is 96% sensitive and 98% specific. What do you think its positive predictive value is for the <55 population which has an Afib prevalence of 0.1%.

APPLE WATCH

96% sensitive, 98% specific – sounds great!

Prevalence of afib for population under 55 = 0.1% What do you think is positive predictive value?

1% 5% 20% 50% 90%

Overall afib prevalence (21+ yrs old): ~1% Overall PPV: 33%



• Binary Classifiers

- AUC
- Sensitivity/Specificity
- F1 statistic
- Accuracy

• Non-binary Classifiers

- Accuracy
- o Kappa
- Regressions
 - R-squared
 - MSE
- Which do you choose? It depends

The right statistic for you depends on your population, use case, and preferences

- In general:
 - Asking for more sensitivity will hurt positive predictive value and specificity.
 - Consider optimizing something like an F1 score (which penalizes false positives).

Roadmap

Survey Results

PICO for ML Papers

Discussing Barriers to ML

Roadmap

Survey Results

PICO for ML Papers

Discussing Barriers to ML

Survey Results

Have you attempted or successfully deployed any machine learning projects? (Select one)

20 responses


What barriers have you faced in implementing ML solutions (Select all that apply)

20 responses



