# Evaluating Classifiers and Annotating Data

Kenneth (Kenny) Joseph

**University at Buffalo**
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

# Announcements

- Quiz 6 is "out"
- Midterm is **Thursday**
  - In class
  - One page handwritten notes, front and back
  - **Nothing else** (except pen/pencil)
- Two quick review things
- Questions?

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# Evaluating Classification Models

- How should we evaluate (part 1)?
- What is the best we can do?
- What is the worst we can do?
- Class Imbalances
- How should we evaluate (part 2)?
- Dealing w/ Class Imbalance through Modeling

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

UB has created a predictive algorithm to determine who should be admitted to the CSE MS program.

The algorithm takes the ~~ACT~~ GRE score and School Ranking as features, and past decisions on admissions as the outcome

The algorithm is used to **admit or reject students** starting next year

# Back to regression

- How would we evaluate this with regression, i.e. what would our evaluation metric be?

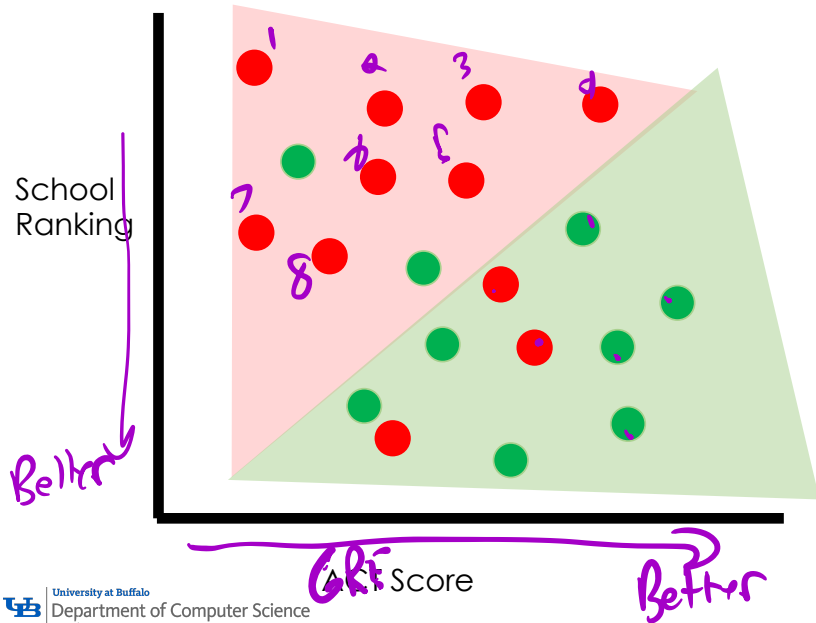$$\frac{1}{N} \sum (y - h(x))^2$$

What values can y take on?

$+1 \quad -1$

$0$

$\sum (+1 - +1)^2$

$0$

What about $h(x)$?

$+1 \quad -1$

$\sum (1 - -1)^2 = 4$

5

@_kenny_joseph

# Evaluating classification models

# The Confusion Matrix



relevant elements

false negatives | true negatives

true positives | false positives

retrieved elements

https://en.wikipedia.org/wiki/Precision_and_recall

@_kenny_joseph

# Accuracy – how many did we get correct?



School Ranking

ACT Score

Our guess:

"Truth"

|  | | |
|---|---|---|
| **8** | 3 |
| 2 | **7** |

Accuracy =
(8 + 7) / (8 + 7 + 2 + 3)
= .75

8

@_kenny_joseph

# What is the best we can do?

## The Bayes optimal classifier

$$y^* = h_{Best}(x) = \underset{x}{argmax}\ P(y|x)$$

← what if we know this

$$\begin{bmatrix} P(+1|\mathbf{x}) = .8 \\ P(-1|\mathbf{x}) = .2 \end{bmatrix}$$ ← wrong 20% of the time!

$$\epsilon_{opt} = 1 - P(h_{Best}(x)|x)$$

Discussion follows: https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote02_kNN.html

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

9

@_kenny_joseph

# What is the *worst* we can do?

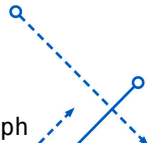Constant classifier / average label/majority

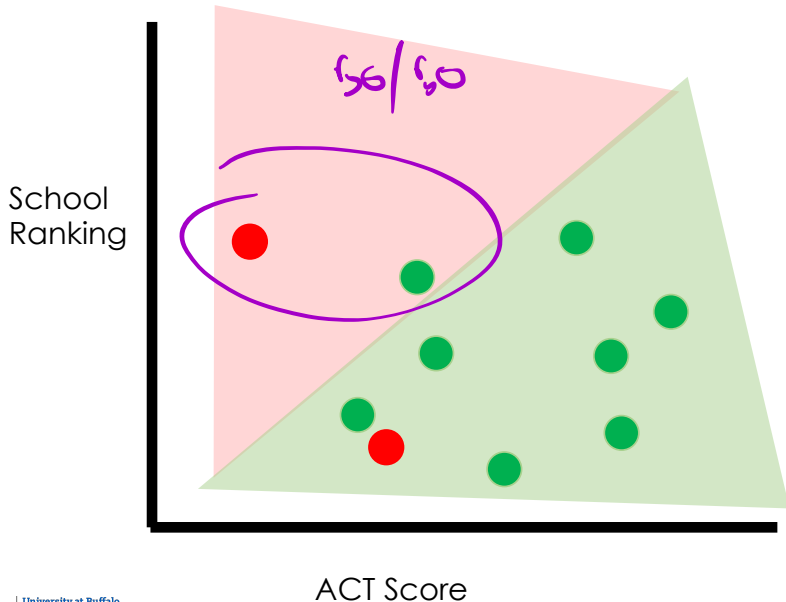Random Guessing: 50%

$80\%$ if $P(+1) = .8$

Class imbalance

**Always compare to the simple baseline for your model**

10

@_kenny_joseph

# The problem with class imbalance



School Ranking

ACT Score

50/50

**Our guess:**

"Truth"

|  |  |  |
|--|--|--|
| 🔴 | ①  | 1 |
| 🟢 | 1  | ⑦ |

**What is our accuracy?**

$$(1+7) / (1 + 7 + 1 + 1))$$

$$80\%$$

@_kenny_joseph

# The problem with class imbalance

Our guess:

|  | | |
|---|---|---|
| 🔴 | 1 | 1 |
| 🟢 | 1 | ⑦ |

"Truth"

School Ranking

ACT Score

**Is this really a good classifier?** Not really
**How does a majority classifier do?** 90%

@_kenny_joseph

# Aside – dealing with class imbalance

Decision function of LogisticRegression



As a result, the majority class does not take over the other classes during the training process.

https://imbalanced-learn.org/stable/over_sampling.html

Will cover, along with a few other things, in a "practical issues" lecture at some point after the break

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

13

@_kenny_joseph

# Precision - Of ⊕ guesses, how many actually +s?



**Our guess:**

|  | | |
|---|---|---|
| 🔴 | 8 | **3** |
| 🟢 | 2 | **7** |

"Truth"

Precision =
7 / (7 + 3) = .7

@_kenny_joseph

# Recall - Of actual +, how many do we guess?

Recall of − class:
$$8/(8+3)$$



Our guess:

|  | 8 | 3 |
|---|---|---|
| "Truth" | **2** | **7** |

Recall =
7 / (7 + 2) = .78

@_kenny_joseph

# To compute precision and recall, you have to pick a class!



School Ranking

ACT Score

Recall − : 50%

Precision − : 50%

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# To compute precision and recall, you have to pick a class!

School Ranking

ACT Score

If you were applying to UB,

Recall +

Precision +

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# Which would you prefer?



School Ranking

ACT Score

School Ranking

ACT Score

@_kenny_joseph

# There are many other metrics

Other metrics can be included in a confusion matrix, each of them having their significance and use.

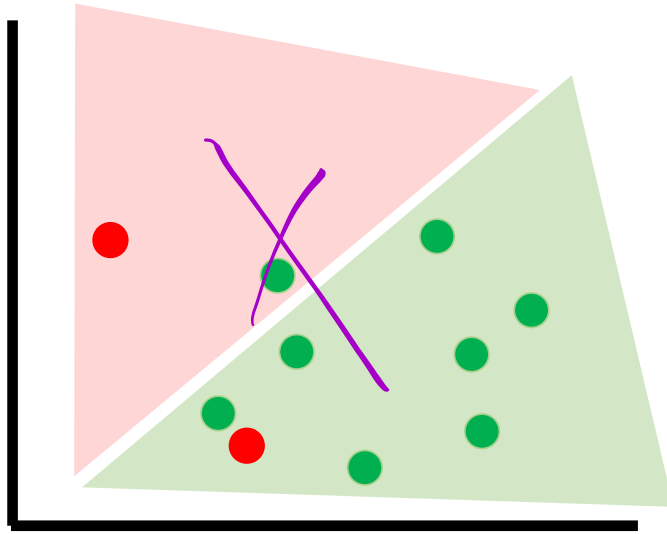| | | Predicted condition | | | |
|---|---|---|---|---|---|
| Total population = P + N | | Positive (PP) | Negative (PN) | Informedness, bookmaker informedness (BM) = TPR + TNR − 1 | Prevalence threshold (PT) = $\frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$ |
| Actual condition | Positive (P) | True positive (TP), hit | False negative (FN), type II error, miss, underestimation | True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power = $\frac{TP}{P}$ = 1 − FNR | False negative rate (FNR), miss rate = $\frac{FN}{P}$ = 1 − TPR |
| | Negative (N) | False positive (FP), type I error, false alarm, overestimation | True negative (TN), correct rejection | False positive rate (FPR), probability of false alarm, fall-out = $\frac{FP}{N}$ = 1 − TNR | True negative rate (TNR), specificity (SPC), selectivity = $\frac{TN}{N}$ = 1 − FPR |
| Prevalence = $\frac{P}{P+N}$ | | Positive predictive value (PPV), precision = $\frac{TP}{PP}$ = 1 − FDR | False omission rate (FOR) = $\frac{FN}{PN}$ = 1 − NPV | Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$ | Negative likelihood ratio (LR−) = $\frac{FNR}{TNR}$ |
| Accuracy (ACC) = $\frac{TP + TN}{P + N}$ | | False discovery rate (FDR) = $\frac{FP}{PP}$ = 1 − PPV | Negative predictive value (NPV) = $\frac{TN}{PN}$ = 1 − FOR | Markedness (MK), deltaP (Δp) = PPV + NPV − 1 | Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$ |
| Balanced accuracy (BA) = $\frac{TPR + TNR}{2}$ | | $F_1$ score = $\frac{2PPV \times TPR}{PPV + TPR}$ = $\frac{2TP}{2TP + FP + FN}$ | Fowlkes–Mallows index (FM) = $\sqrt{PPV \times TPR}$ | Matthews correlation coefficient (MCC) = $\sqrt{TPR \times TNR \times PPV \times NPV}$ − $\sqrt{FNR \times FPR \times FOR \times FDR}$ | Threat score (TS), critical success index (CSI), Jaccard index = $\frac{TP}{TP + FN + FP}$ |

Sources: [20][21][22][23][24][25][26][27] view · talk · edit

https://en.wikipedia.org/wiki/Confusion_matrix

Many different metrics … we'll dive into a few now, but not all

@_kenny_joseph

# Critical Idea: Accounting for Thresholds

Remember that, e.g., logistic regression predicts a continuous value, and then we threshold

$$score(x) < \text{threshold: } -1$$
$$\text{otherwise: } +1$$

The threshold is in some ways a *hyperparameter* … we can get different, e.g., accuracies with different thresholds.

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# Looking at Thresholds, V1: Precision/Recall Curve

$$Score(x) < threshold : -1$$

threshold: $\infty$ ; recall + : $O$

recall $-$ : 1

precision + : $\%$ = 1

Negatives as positives ("False Positives")

True positives as negatives ("False Negative")

Threshold

True Neg
$h(x)_{y=-1}$

True Pos

$-1$   $p(+1)$   $+1$

@_kenny_joseph

# Looking at Thresholds, V1: Precision/Recall Curve



Precision

precision=1 recall=0 : threshold=∞

Best model

recall

everything is +
threshold → -∞

recall=1 ; precision=0

- What does the best classifier look like?

- Which is the better classifier?

@_kenny_joseph

# Looking at Thresholds, V1: Precision/Recall Curve

- How to summarize this?

$$\boxed{F_1} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{\text{TP}}{\text{TP} + \frac{1}{2}(\text{FP} + \text{FN})}$$

⭐ A metric that unifies precision + recall

@_kenny_joseph

Best Model

True Positive Rate

$$\frac{\text{\# True Positives}}{\text{True P + FN}}$$

Rate FP < Rate TP

False Positive Rate

@_kenny_joseph

# Looking at Threshold Changes, V2: ROC

- How to summarize this?

$$AUC(f) = \frac{\sum_{t_0 \in \mathcal{D}^0} \sum_{t_1 \in \mathcal{D}^1} \mathbf{1}[f(t_0) < f(t_1)]}{|\mathcal{D}^0| \cdot |\mathcal{D}^1|},$$

University at Buffalo
Department of Computer Science
and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# Looking at Threshold Changes, V3: Precision @ k

- Final idea: State a number k of observations that you care about, look at precision there
- **Where might this be useful?** Search!

① Rank predictions $p(y|x)$ ↑ predicted probability

② Pick k := 10

③ What % of the top k are +1

precision @ k

# Which metric do we want?

- Diagnosing cancer    *Recall*
- Putting someone in jail    *Precision*

@_kenny_joseph

# Evaluation Review

- Big ideas:
  - Different metrics for different things
  - Evaluation metrics != loss function
  - Beware of class imbalances
  - Use a lot of metrics!
  - But ultimately, the right metric is tied to your application area

@_kenny_joseph

# What is missing from these evaluations?

# Adding a new feature: height

**Geoffrey Hinton**
@geoffreyhinton

Suppose you have cancer and you have to choose between a black box AI surgeon that cannot explain how it works but has a 90% cure rate and a human surgeon with an 80% cure rate. Do you want the AI surgeon to be illegal?

3:37 PM · Feb 20, 2020 · Twitter Web App

**1,125** Retweets    **615** Quote Tweets    **5,065** Likes

Income Distributions by ACT Score, 2018

**Family Income**
- About $0 to $24,000
- About $24,000 to $36,000
- About $36,000 to $50,000
- About $50,000 to $60,000
- About $60,000 to $80,000
- About $80,000 to $100,000
- About $100,000 to $120,000
- About $120,000 to $150,000
- More than $150,000

https://twitter.com/JonBoeckenstedt/status/1447584690932629511/photo/1

32

# Annotation

# Annotation Discussion - Overview

- Where do annotations come from?
- How do you know if they're any good?
  - Accuracy on downstream "expert annotated" data
  - Agreement
    - Percent agreement
    - Krippendorf
- Can we do annotation differently?
  - Aggregation models
  - Snorkel, etc.
  - Considering annotator demographics

*This Class*

*A*

3/15/22

@_kenny_joseph

# Where does data come from?

- Ultimately, most datasets come from *people*
- What might be problematic about that?

@_kenny_joseph

# Where do *annotations* come from?

"Expert" Annotators (e.g. domain experts)

slow but accurate

amazon mechanical turk™
Artificial Artificial Intelligence

CrowdFlower

…

@_kenny_joseph

# Challenges with crowd annotation

- How can you incentivize good-faith labels?
- How do you know that you're getting good faith labels?
- How do you aggregate responses across a bunch of people?

@_kenny_joseph

# Incentivizing Good-faith Labels

- Treat people with respect
  - Pay them
  - Be nice to them

@_kenny_joseph

# Ensuring Good-faith Labels

- Gold standards – have some observations you know the answer to
- Attention checks – have some questions like "are you awake"
- Redundancy – make sure multiple annotators per observation
- Really, redundancy + **agreement statistics**

@_kenny_joseph

# Agreement Statistics

- Pairwise agreement: basically, accuracy per annotator

A  B

0  1

0  1

0  0

0  0

% of observations
where they agree
Class imbalanced.

- Krippendorf's Alpha

@_kenny_joseph

# Krippendorf's Alpha

*Simpledorff* (handwritten annotation)

| | document_id | annotator_id | annotation |
|---|---|---|---|
| 0 | 1 | B | 1 |
| 6 | 3 | B | 2 |
| 7 | 3 | C | 2 |
| 9 | 4 | B | 1 |
| 10 | 4 | C | 1 |

+%/ (handwritten annotation)

| annotator_id | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 1 | 2 | 3 | 3 | 2 | 1 | 4 | 1 | 2 | nan | nan | nan |
| B | 1 | 2 | 3 | 3 | 2 | 2 | 4 | 1 | 2 | 5 | nan | 3 |
| C | nan | 3 | 3 | 3 | 2 | 3 | 4 | 2 | 2 | 5 | 1 | nan |
| D | 1 | 2 | 3 | 3 | 2 | 4 | 4 | 1 | 2 | 5 | 1 | nan |

$1 - \frac{Do}{De}$ (handwritten annotation)

Calculate (1-) the ratio between:

**Do** – observed disagreements
**De** – disagreement by chance

https://www.lighttag.io/blog/krippendorffs-alpha/

@_kenny_joseph

# Krippendorf's Alpha (cont.)

| annotator_id | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 1 | 2 | 3 | 3 | 2 | 1 | 4 | 1 | 2 | nan | nan | nan |
| B | 1 | 2 | 3 | 3 | 2 | 2 | 4 | 1 | 2 | 5 | nan | |
| C | nan | 3 | 3 | 3 | 2 | 3 | 4 | 2 | 2 | 5 | 1 | nan |
| D | 1 | 2 | 3 | 3 | 2 | 4 | 4 | 1 | 2 | 5 | 1 | nan |

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 3 | 0 | 0 | 2 | 0 |
| 2 | 0 | 3 | 0 | 0 | 4 | 1 | 0 | 1 | 4 | 0 | 0 | 0 |
| 3 | 0 | 1 | 4 | 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| 4 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |

*(handwritten annotations: "Tweet", "Pro +", "Neutral", "Anti", "#annotators")*

https://www.lighttag.io/blog/krippendorffs-alpha/

42

@_kenny_joseph

# Krippendorf's Alpha - observed

Rc

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 3 | 0 | 0 | 2 | 0 |
| 2 | 0 | 3 | 0 | 0 | 4 | 1 | 0 | 1 | 4 | 0 | 0 | 0 |
| 3 | 0 | 1 | 4 | 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| 4 | 0 | 0 | 0 | 0 | 0 | 1 | 4 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |

\# disagreements = 0

\# disagreements = 3

https://www.lighttag.io/blog/krippendorffs-alpha/

43

@_kenny_joseph

# Krippendorf's Alpha – by chance

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 3 | 0 | 0 | 2 | **0** |
| 2 | 0 | 3 | 0 | 0 | 4 | 1 | 0 | 1 | 4 | 0 | 0 | **0** |
| 3 | 0 | 1 | 4 | 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | **1** |
| 4 | 0 | 0 | 0 | 0 | 0 | 1 | 4 | 0 | 0 | 0 | 0 | **0** |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | **0** |

$$V_1 \cdot V_2 + V_1 \cdot V_3 \ \text{---}$$

$V_1$

$V_2$

$V_3$

$V_4$

$V_5$

https://www.lighttag.io/blog/krippendorffs-alpha/

@_kenny_joseph

# Krippendorf's Alpha – simple, worked through

Items judged:

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Meg:** | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| **Owen:** | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

Values:

| | 0 | 1 | |
|---|---|---|---|
| **0** | $o_{00}$ | $o_{01}$ | $n_0$ |
| **1** | $o_{10}$ | $o_{11}$ | $n_1$ |
| Number of Values: | $n_0$ | $n_1$ | $n=2N$ |

| | 0 | 1 | |
|---|---|---|---|
| **0** | 10 | 4 | 14 |
| **1** | 4 | 2 | 6 |
| | 14 | 6 | 20 |

④ Compute α-**reliability** (most simple form):

$$_{binary}\alpha = 1 - \frac{D_o}{D_e} = 1 - (n-1)\frac{o_{01}}{n_0 \cdot n_1}$$

In the example:

$$_{binary}\alpha = 1 - (20-1)\frac{4}{14 \cdot 6} = 0.095$$

https://repository.upenn.edu/cgi/viewcontent.cgi?article=1043&context=asc_papers

45

@_kenny_joseph

# Aggregation

- The most common approach is **majority vote**
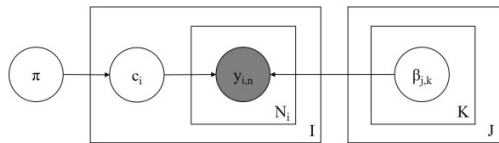- More recently, people have come up with better ways


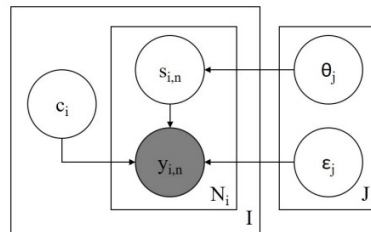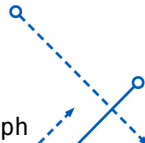
Figure 2: Plate diagram of the Dawid and Skene model.

1979



Figure 3: Plate diagram for the MACE model.

https://watermark.silverchair.com/tacl_a_00040.pdf

@_kenny_joseph

# Moving forward: Smarter Annotation...

- **Data Programming & Weak Supervision**
- **Data Augmentation**
- **Self-Supervision**
- **Data Selection**

↳ NLP

- More: https://github.com/HazyResearch/data-centric-ai

Active learning

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph