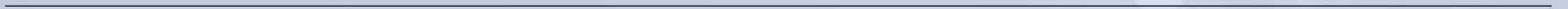
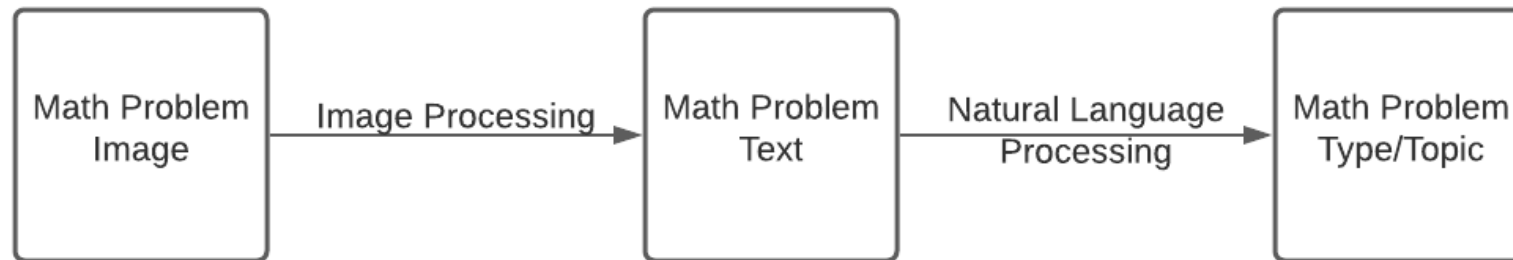


Classification for Math Problems

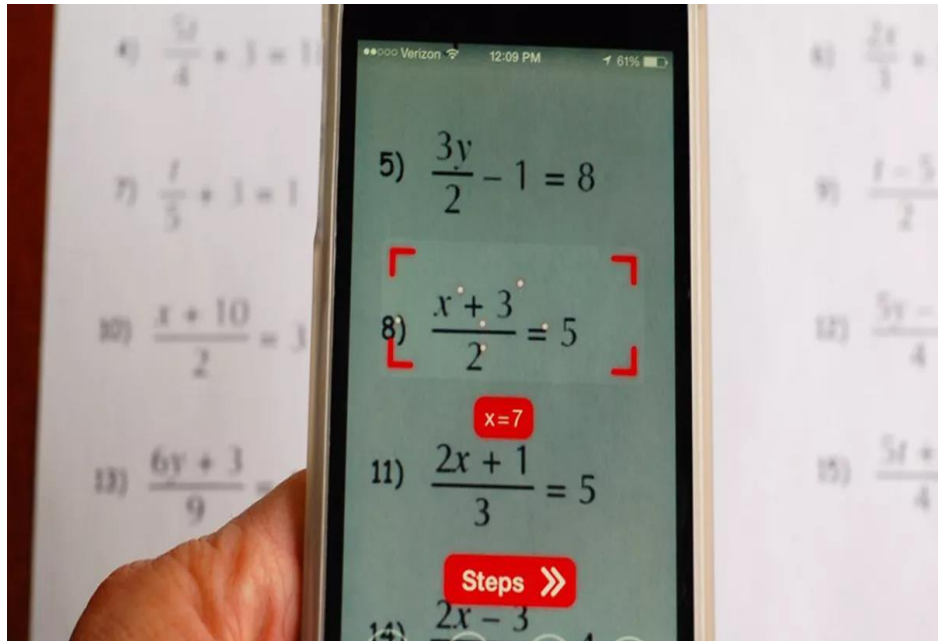
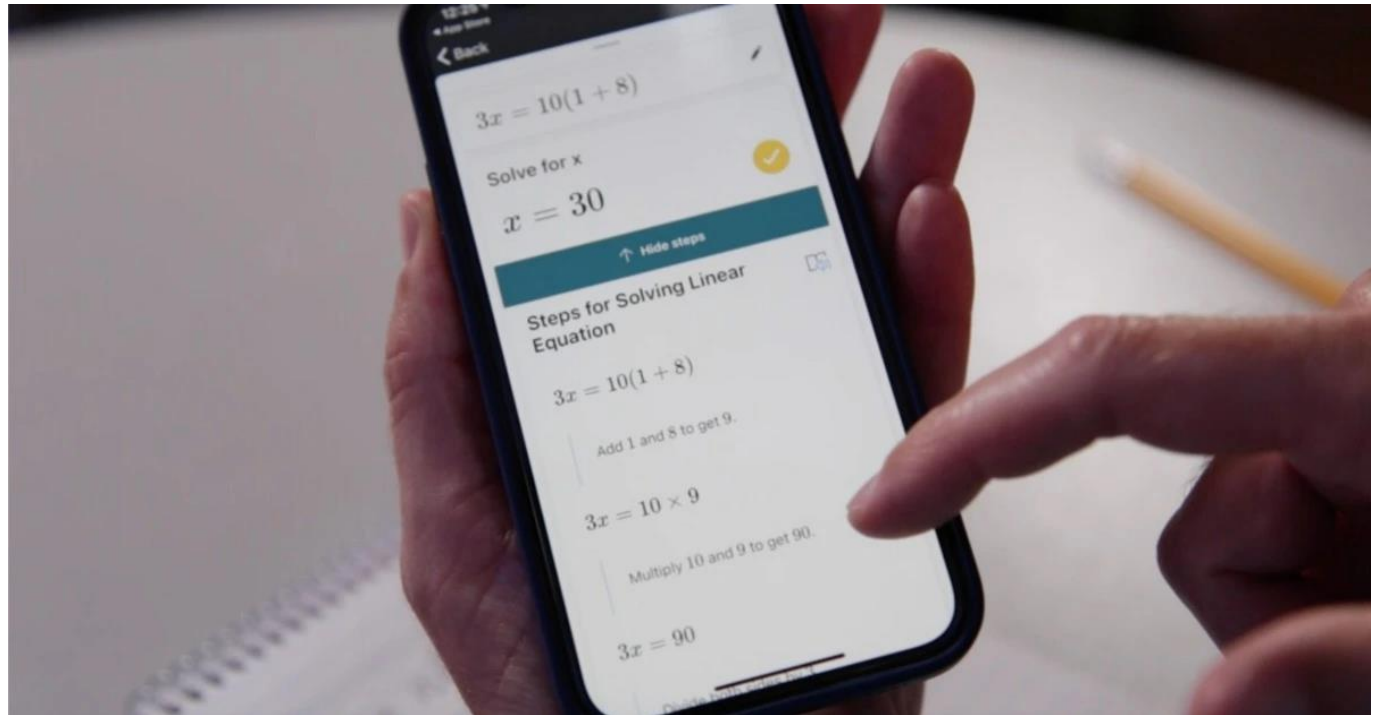


Problem Statement

- We are looking to classify math problems by problem type



Processing Flow



Use of Classification

- Recommend learning resources
- Future Considerations:
 - Provide answer check
 - Gamification?

Analytic Approach

- Multi-output/Multilabel classification problem



- We trained a Decision Tree algorithm to perform multi-output classification of text into categories of math problems.

Data Processing

- Sourced from GitHub community called Hugging Face
- Subset of larger data set
 - 56 categories of math problems
 - We took 10,000 questions from each category
- 56 categories of math problems
 - Examples: Algebra Linear 1d, Arithmetic Add or Sub, Calculus differentiate, Numbers is prime, Polynomials expand
- Problems are synthetically created

Some examples of math problems

question	
0	Find the first derivative of $2d^4 - 35d^2 - 695$ wrt d .
1	Find the third derivative of $-a^3g^3t^3 + 642a^3gt^3 + 16a^3gt^2 - 5a^2t^2 + ag^3$ wrt t .
2	What is the second derivative of $12518f^3 + 3760f$?
3	What is the third derivative of $-t^4 - 880t^3 + 152t^2$ wrt t ?
4	What is the second derivative of $2cn^2z^3 + 30cn^2 + 2cnz^2 - 2c + n^2z^2 - 3nz^3 - 2nz$ wrt n ?
5	Find the first derivative of $-1373u^3 + 81$ wrt u .
6	What is the derivative of $-2612k - 37$?
7	Find the third derivative of $w^6 - 2w^5 + 579w^4 - 5032w^2$.
8	What is the third derivative of $-70f^5 + 6f^2$?
9	Differentiate $745b^4 - 287$ with respect to b .
10	What is the second derivative of $-600l^4 - 5l^3 - 6l - 529$?

question	
0	What is $-5 - 110911$?
1	What is $-0.188 + -0.814$?
2	Sum 259 and -46.
3	Sum -10 and -52539.
4	What is the difference between -2 and 251860?
5	$-9259432 + 1$
6	What is 1141.09 less than 1?
7	What is $-5 - 72726$?
8	Total of 0.3 and 170.7.
9	Work out $29.8 + -0.18$.
10	What is 19450 minus -0.8?

question	
0	Solve $0 = 4b + b + 15$ for b .
1	Solve $-3d = -0d + 3$ for d .
2	Solve $-4h + 9 = 41$ for h .
3	Solve $2514m = 2508m - 24$ for m .
4	Solve $-7a + 6a = 4$ for a .
5	Solve $288w - 298w = -70$ for w .
6	Solve $-14h = -4h - 10$ for h .
7	Solve $5w + 3 = -2$ for w .
8	Solve $-15f + 21f - 12 = 0$ for f .
9	Solve $-22 = 6c - 4$ for c .
10	Solve $13z - 7z + 30 = 0$ for z .

Model Form and Fitting

- Decision Tree algorithm:
 - LGBM Classifier (Light Gradient Boosted Machine Classifier)
- Two methods:
 - Classifier Chain and Multi Output Classifier
- Used TF-IDF measure
 - Determines how relevant a word is to a question and a category
 - Counts appearances of words in questions (Term Frequency) and across categories (Inverse Document Frequency)
 - Higher the TF-IDF score, the more relevant the word is

Model Results and Analysis

- Used Three-fold Cross Validation
- Calculated F-scores to measure accuracy of the fitted models
- F-score for LGBM with Multi Output Classifier: .949699
- F-score for LGBM with Classifier Chain: .948426
- Chose LGBM with Multi Output Classifier as better model

Model Results and Analysis

Classification Report:

- Most Categories have 1.00 Precision, Recall, & F score
- Some categories with poor scores likely due to related categories

	precision	recall	f1-score	support
algebra__linear_1d	0.95	1.00	0.97	1938
algebra__linear_1d_composed	0.84	0.85	0.84	1973
algebra__linear_2d	1.00	1.00	1.00	1949
algebra__linear_2d_composed	0.86	0.85	0.85	2084
algebra__polynomial_roots	0.93	0.47	0.62	2099
algebra__polynomial_roots_composed	0.68	0.91	0.78	1915
algebra__sequence_next_term	1.00	1.00	1.00	1980
algebra__sequence_nth_term	1.00	1.00	1.00	2049
arithmetic__add_or_sub	1.00	0.79	0.88	2056
arithmetic__add_or_sub_in_base	1.00	1.00	1.00	1989
arithmetic__add_sub_multiple	0.85	0.99	0.91	1955
arithmetic__div	1.00	1.00	1.00	1970
arithmetic__mixed	0.95	0.67	0.78	1981
arithmetic__mul	0.94	1.00	0.97	1904
arithmetic__mul_div_multiple	0.80	0.93	0.86	1999
arithmetic__nearest_integer_root	1.00	1.00	1.00	2007
arithmetic__simplify_surd	1.00	1.00	1.00	2091
calculus__differentiate	0.91	0.45	0.60	2046
calculus__differentiate_composed	0.64	0.96	0.77	1937

comparison__closest	1.00	1.00	1.00	2097
comparison__closest_composed	1.00	1.00	1.00	1994
comparison__kth_biggest	1.00	1.00	1.00	2027
comparison__kth_biggest_composed	1.00	1.00	1.00	1996
comparison__pair	1.00	1.00	1.00	1979
comparison__pair_composed	1.00	1.00	1.00	1980
comparison__sort	1.00	1.00	1.00	2010
comparison__sort_composed	1.00	1.00	1.00	2013
measurement__conversion	1.00	1.00	1.00	1966
measurement__time	1.00	1.00	1.00	1966
numbers__base_conversion	1.00	1.00	1.00	1960
numbers__div_remainder	1.00	1.00	1.00	1973
numbers__div_remainder_composed	1.00	1.00	1.00	2002
numbers__gcd	1.00	1.00	1.00	1977
numbers__gcd_composed	1.00	1.00	1.00	2104
numbers__is_factor	1.00	1.00	1.00	1937
numbers__is_factor_composed	1.00	1.00	1.00	2055
numbers__is_prime	0.99	1.00	1.00	1996
numbers__is_prime_composed	1.00	0.98	0.99	1993
numbers__lcm	1.00	1.00	1.00	2088

numbers__lcm_composed	1.00	0.99	1.00	1970
numbers__list_prime_factors	1.00	1.00	1.00	2058
numbers__list_prime_factors_composed	1.00	1.00	1.00	2042
numbers__place_value	0.99	1.00	1.00	2007
numbers__place_value_composed	1.00	0.99	0.99	2082
numbers__round_number	1.00	1.00	1.00	1990
numbers__round_number_composed	1.00	1.00	1.00	1991
polynomials__add	0.87	0.82	0.84	2017
polynomials__coefficient_named	1.00	1.00	1.00	1932
polynomials__collect	1.00	1.00	1.00	2008
polynomials__compose	0.87	0.80	0.83	1968
polynomials__evaluate	0.99	0.92	0.96	1962
polynomials__evaluate_composed	0.91	0.92	0.92	2012
polynomials__expand	1.00	1.00	1.00	1938
polynomials__simplify_power	1.00	1.00	1.00	1934
probability__swr_p_level_set	1.00	1.00	1.00	2026
probability__swr_p_sequence	1.00	1.00	1.00	2028

Conclusion

- Using a LGBM Classifier Algorithm to classify math problem text into concept categories
- The model worked well with some limitations
- Could improve the model and use case of the model by making changes to the dataset
 - Changing the categories and how they are comprised
 - Adding more questions from other sources to introduce diversity of math problems

Citations

Hugging Face Datasets.

Github: <https://github.com/huggingface/datasets/tree/master/datasets>

Disaster Message NLP Pipeline.

Github: <https://github.com/ChristopherCochet/Disaster-Message-NLP-Pipeline/blob/master/notebooks/ML%20Pipeline%20Preparation.ipynb>

Text Classification in Python by Miguel

Fernandez Zafra. Url: <https://towardsdatascience.com/text-classification-in-python-dd95d264c802>

A gentle introduction to OCR. Url: <https://towardsdatascience.com/a-gentle-introduction-to-ocr-ee1469a201aa>.

Text Detection CTPN

by Shaohui Ruan. Github: <https://github.com/eragonruan/text-detection-ctpn>

Understanding Classification Report. Url: <https://muthu.co/understanding-the-classification-report-in-sklearn/>