
IMDB reviews classification

Alejandro Salazar Mejía, Pasang Tenzin, Jan Frąckowiak

Sample Reviews:

"This was a really interesting Halloween film. I wasn't too thrilled with the whole Thorn theory but it still makes for a good film. I liked getting to see Tommy Doyle back but sadly Donald Pleasance died right after shooting. (...)"

"come undone, I love you! I could not have come to a better conclusion than you did about this movie and it's ending. My family has not seen this movie yet, but I know them too well; they will hate it.(...)"

Project Outline

Objective: Accurate classification the sentiment of reviews: positive or negative

Data: 25k reviews for training, 25k reviews for testing

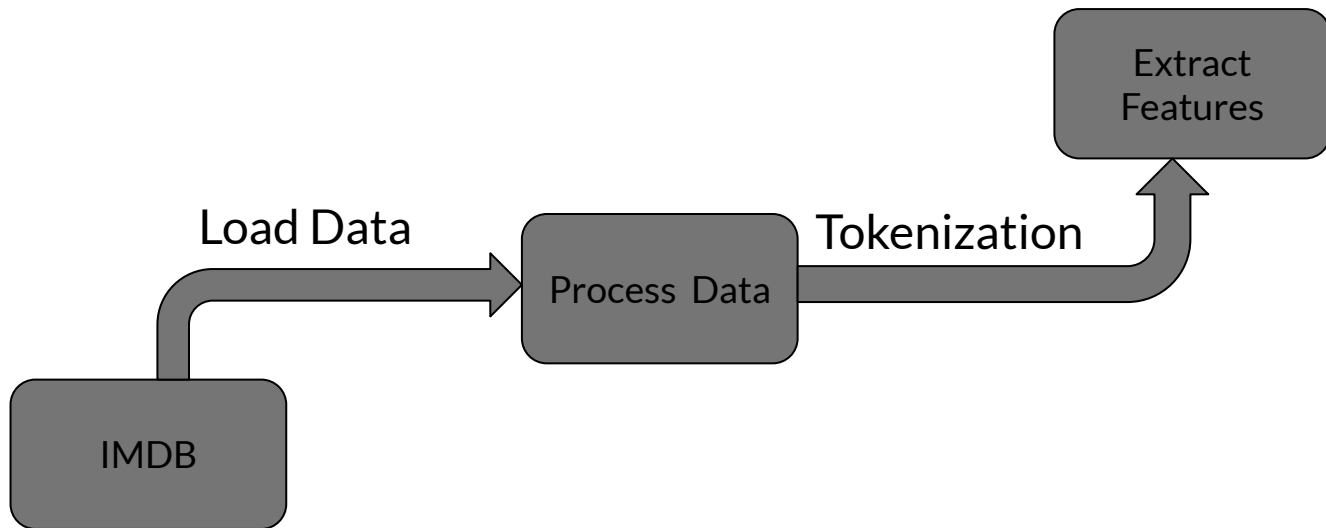
Methods: AutoAI, Manually trained models basing on embeddings

Metrics: Accuracy, Confusion Matrix, ROC-AUC

Github: <https://github.com/thually/cloudTech-Summer23-24>

Benchmark Model

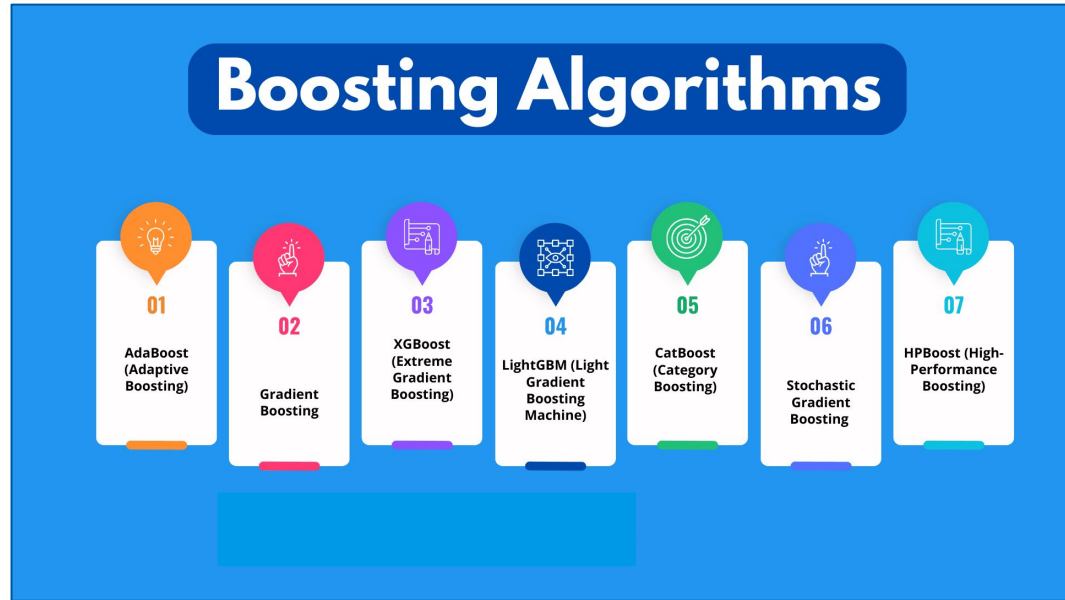
Download IMDB from IBM Cloud Object Storage



Sentiment Analysis

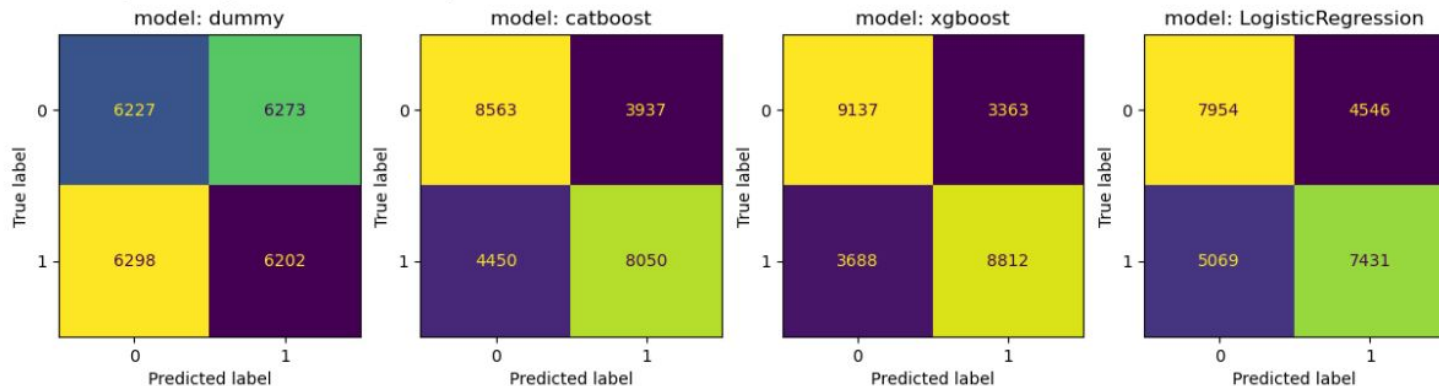
1. Load and process data.
2. Create dataframe and normalized.
3. Sampled the data and visualized.
4. Extract features.

Custom tokenization



Result with TfidfVectorizer and light normalization

```
model=dummy, accuracy=0.49716  
model=catboost, accuracy=0.66452  
model=xgboost, accuracy=0.71796  
model=LogisticRegression, accuracy=0.6154
```



XGBoost performed best with accuracy approximately 71.8%

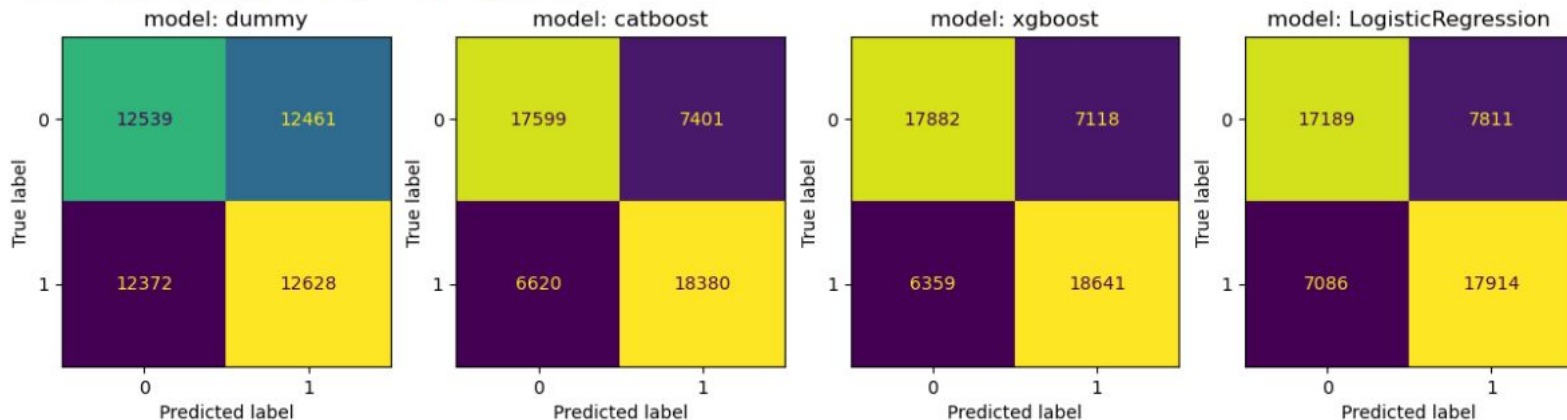
Improvement in the pre-processing of the text

1. Lowercasing all text.
2. Removing URLs.
3. Replace “
” with space.
4. Remove punctuation marks.
5. Tokenize text into words.
6. Remove non-alphabetic words.
7. Remove stop words except for “not”

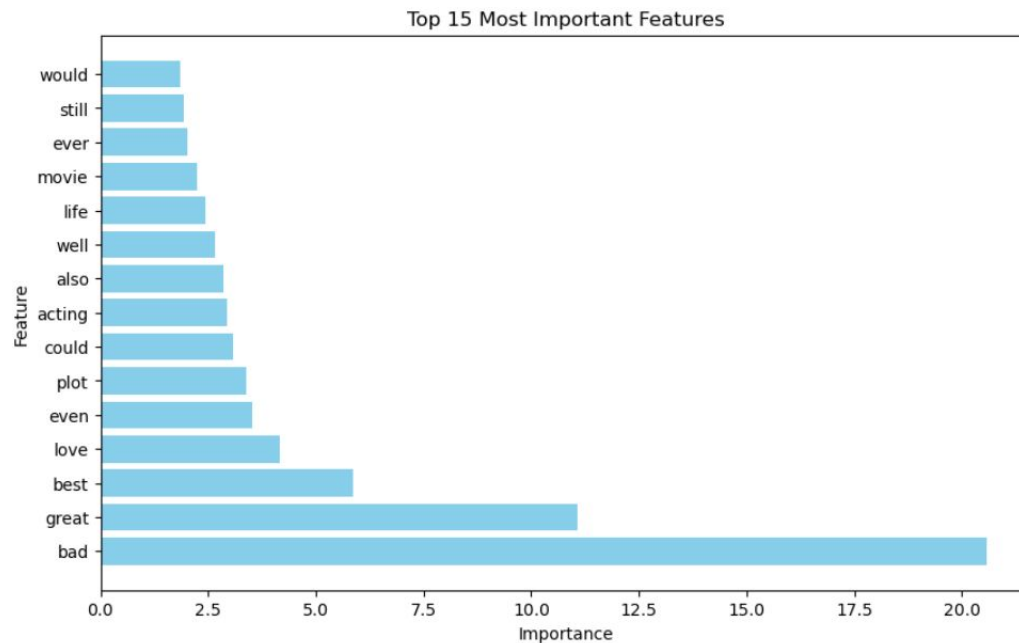
```
{ 're', 'does', 'couldn', 'hasn', 'on', "you'd", 'doesn', 'there', 'myself', 'have', 'about', 'some', 'a', 'weren', 'aren', 'and', 'if', 'wasn', 'haven', 'hers', 'o', 'll', 'that', 'don', 'whom', 'in', 'up', 'here', 'most', "won't", 'i', "hadn't", "you've", 'didn', "didn't", 'had', 'ma', 'for', 'she', 'ourselves', "wasn't", "doesn't", 'the', "hasn't", 'with', 'when', "that'll", 'is', 'now', 'an', 'of', 't', 'mightn', 'needn', 'ain', "wouldn't", 'herself', 'it', 'itself', 'he', "it's", 'will', 'out', 'once', 'my', "mustn't", 'by', 'through', 'before', 'because', 'off', 'ours', 'having', 'between', 'only', 'd', 'where', 'on', "shouldn't", 'hadn', "isn't", 'any', 'isn', 'these', 'few', 'aren't", 'then', 'themselves', "haven't", "weren't", 'our', "you'll", 'until', 'this', 'were', 's', 'yours', 'but', 'at', 'after', 'should', "shan't", 'your', 'against', 'y', 'are', 'm', 'them', 'mustn', 'was', 'being', 'nor', 'did', 'what', "couldn't", 'down', 'can', 'won', 'why', 'theirs', 'all', 'doing', 'as', 'from', 'while', 'other', 'to', 'yourselves', 'been', 'shan', 'wouldn', 'you', 'more', 'yourself', 'her', 'they', 'am', 'which', 'we', 'above', 'just', 've', 'be', "you're", "mightn't", 'how', 'both', "don't", 'into', 'no', "she's", 'its', 'same', 'during', "should've", 'him', 'such', 'very', 'those', 'again', 'who', 'shouldn', 'their', 'under', 'further', 'below', 'too', 'himself', 'over', 'his', 'each', 'do', 'so', 'own', "needn't", 'me', 'than', 'has' }
```

Results with our pre-processing

model=dummy, accuracy=0.50334
model=catboost, accuracy=0.71958
model=xgboost, accuracy=0.73046
model=LogisticRegression, accuracy=0.70206



Feature importance of CATboost



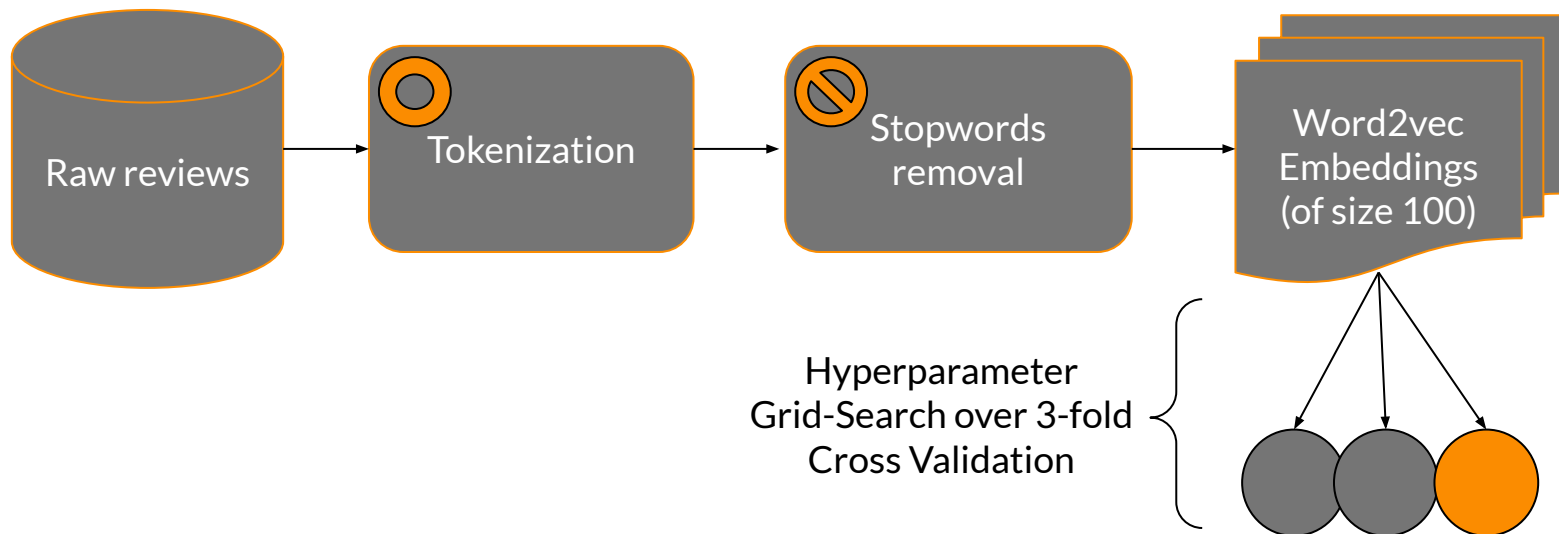
Accuracy: 0.71958
Precision: 0.7129281253636399
Recall: 0.7352
F1-score: 0.7238927945491423

Pipeline with embeddings

Embeddings Pipeline Training

Train-test split: 25k to 25k reviews (perfectly balanced)

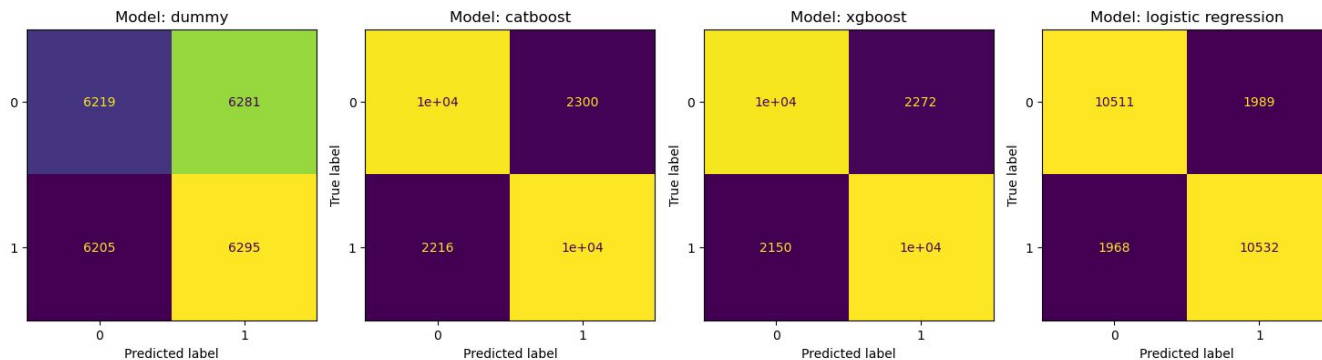
Models: Dummy, XGBoost, CatBoost, LogisticRegression



Embeddings Pipeline Results

	Dummy	Catboost	XGBoost	LogisticReg.
Tested Params:	1	32	48	5
CV Accuracy:	50.2%	91.8%	100%	84.4%
Test Accuracy:	49.6%	81.9%	82.3%	84.2%

Test C-Matrix:



AutoAI

No *Text* feature engineering

	Rank	↑	Name	Algorithm	Accuracy (Optimized) Cross Validation	Enhancements	Build time
★	1		Pipeline 8	Extra Trees Classifier	0.505	HPO-1 FE HPO-2	00:04:37
	2		Pipeline 6	Extra Trees Classifier	0.504	HPO-1	00:01:39
	3		Pipeline 4	Decision Tree Classifier	0.504		
	4		Pipeline 7	Extra Trees Classifier	0.504		
	5		Pipeline 2	Decision Tree Classifier	0.503		
	6		Pipeline 3	Decision Tree Classifier	0.502		
	7		Pipeline 1				








labels		texts
17797	1	A film that deserved theatrical release. This ...
6458	0	What a drawn out painful experience. <br ...
14478	1	Every American who thinks he or she understand...
14825	1	I think a person would be well-advised to read...
7523	0	this is horrible film. it is past dumb. first,...
		...

Text feature engineering

When enabled, columns detected as text will be transformed into vectors to better analyze semantic similarity between strings. Enabling this setting may increase run time.

☐ Use text feature engineering

With *Text feature engineering* - 1/2

	Rank	↑	Name	Algorithm	Accuracy (Optimized) <small>Cross Validation</small>	Enhancements	Build time
★	1		Pipeline 6	 XGB Classifier	0.722	TFE HPO-1	00:01:35
	2		Pipeline 8	 XGB Classifier	0.721	TFE HPO-1 FE HPO-2	00:05:47
	3		Pipeline 2	 LGBM Classifier	0.721	TFE HPO-1	00:02:34
	4		Pipeline 1	 LGBM Classifier	0.721		
	5		Pipeline 4	 LGBM Classifier	0.720		
	6		Pipeline 5	 XGB Classifier	0.719		
	7		Pipeline 7	 XGB Classifier	0.718		

Text feature engineering

When enabled, columns detected as text will be transformed into vectors to better analyze semantic similarity between strings. Enabling this setting may increase run time.

☒ Use text feature engineering

Text feature engineering columns ⓘ

☒ Automatically select text columns

☐ Manually specify text columns

Number of vectors per column ⓘ

20

Specify how many vectors to create for each column during text feature engineering. A lower number is faster and a higher number is more accurate but slower.

—

+

With *Text feature engineering* - 2/2

	Rank	↑	Name	Algorithm	Accuracy (Optimized) <small>Cross Validation</small>	Enhancements	Build time
★	1		Pipeline 4	Logistic Regression	0.792	TFE HPO-1 FE HPO-2	00:10:56
	2		Pipeline 3	Logistic Regression	0.792	TFE HPO-1 FE	00:06:10
	3		Pipeline 8	LGBM Classifier	0.791	TFE HPO-1 FE HPO-2	00:29:29
	4		Pipeline 7	LGBM Classifier	0.791		
	5		Pipeline 2	Logistic Regression	0.791		
	6		Pipeline 1	Logistic Regression	0.790		
	7		Pipeline 6	LGBM Classifier	0.789		

Text feature engineering

When enabled, columns detected as text will be transformed into vectors to better analyze semantic similarity between strings. Enabling this setting may increase run time.

☒ Use text feature engineering

Text feature engineering columns ⓘ

☒ Automatically select text columns

☐ Manually specify text columns

Number of vectors per column ⓘ

30

-

+

	Rank	↑	Name	Algorithm	Accuracy (Optimized) Cross Validation	Enhancements	Build time
★	1		Pipeline 2	🟪 XGB Classifier	0.707	HPO-1	00:01:01
	2		Pipeline 1	🟪 XGB Classifier	0.706	None	00:00:07
	3		Pipeline 6	🟢 LGBM Classifier	0.705	HPO-1	00:04:28
	4		Pipeline 5	🟢 LGBM Classifier	0.704		
	5		Pipeline 4	🟪 XGB Classifier	0.704		
	6		Pipeline 3	🟪 XGB Classifier	0.704		

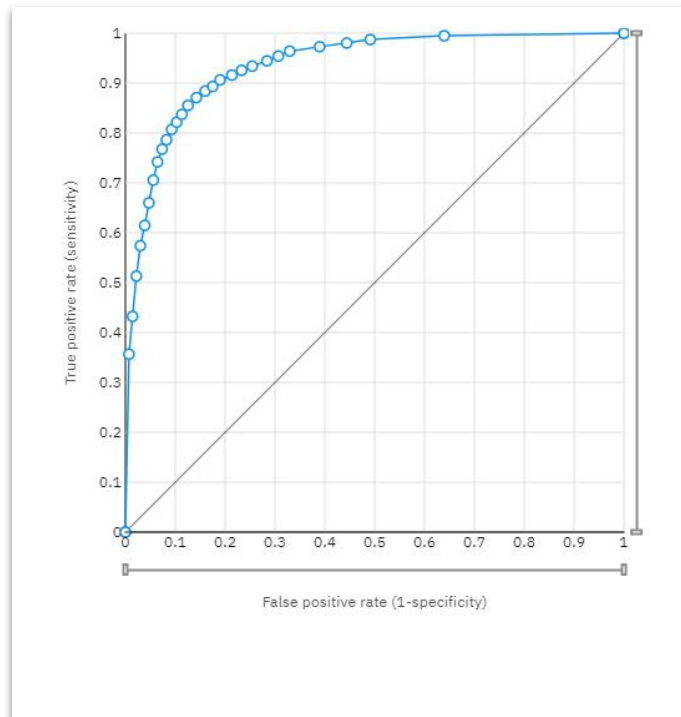
	movie	film	not	one	like	good	would
31809	0	1	1	1	1	0	0
34702	1	1	1	1	1	1	0
14222	0	1	1	1	0	0	0
26294	0	1	1	0	1	1	1

[illegible]

	Rank	↑	Name	Algorithm	Accuracy (Optimized) Cross Validation	Enhancements	Build time
★	1		Pipeline 2	🟡 Snap Logistic Regression	0.859	HPO-1	00:00:28
	2		Pipeline 1	🟡 Snap Logistic Regression	0.859	None	00:00:08
	3		Pipeline 6	🟢 LGBM Classifier	0.858	HPO-1	00:04:08
	4		Pipeline 4	🟡 Snap Logistic Regression	0.858	HPO-1 FE HPO-2	00:02:20

[illegible]

P2 Snap Logistic Regression



Observed	Predicted		
	1	0	Percent correct
1	2196	304	87.8%
0	386	2114	84.6%
Percent correct	85.1%	87.4%	86.2%

Less correct More correct

Measures	Holdout score	Cross validation score
Accuracy	0.862	0.860
Area under ROC	0.936	0.935
Precision	0.851	0.851
Recall	0.878	0.873
F1	0.864	0.862
Average precision	0.934	0.932
Log loss	0.323	0.327



Environment definition ⓘ

Large: 8 CPU and 32 GB RAM

This environment definition consumes **20 capacity units per hour** for training. For details, see [Watson Machine Learning plans](#).

Resource usage

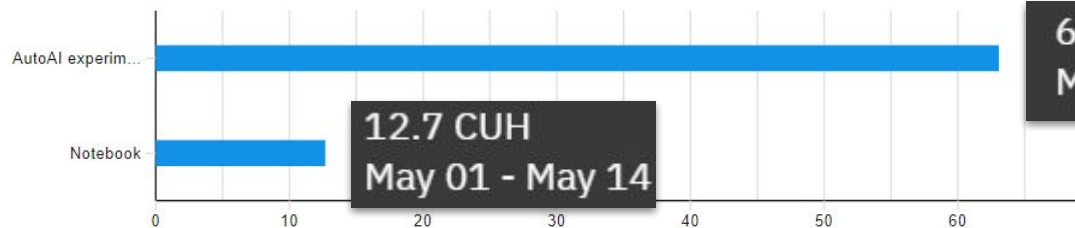
Usage summary ⓘ

For this month in this project

75.8 CUH

Usage by type ⓘ

CUH



63.1 CUH
May 01 - May 14

12.7 CUH
May 01 - May 14

Let's jump to application!
