
Error Generating Reflection on Data Modeling

The slide features decorative horizontal lines: a thick teal line at the top, a thin teal line below it, and another thick teal line at the bottom. Two short, thick grey dashes are positioned horizontally, one on the left and one on the right, centered vertically between the middle thin teal line and the bottom thick teal line.

Research Experiments

- 1) Random Domain Active Error Generating on the data (Anime Type)
- 2) Typo Error Generating in the data (Anime Genre)
- 3) Gaussian Noise in the highest ranked feature (Members who voted to the Anime)

RAD - Anime Type

1) **Generating**

Shuffling the Anime Type (Tv series, Movie, NaN)

2) **Detecting**

Depending on the Episodes numbers

one Episodes -> Movie, otherwise it is Tv-series

3) **Repairing**

Using Episodes numbers

RAD - Error on 25% of the Training Dataset

	Before		Error in Anime Type		After	
	MSE on test dataset	Var on test dataset	Pct of detected error as movie	Pct detected error as tv series	MSE	Var on test dataset
Lasso Without HP	1.13	-29.05	26.5625%	23.0625%	0.85	0.12
Lasso With HP	1.12	-29.05	26.5625%	23.0625%	0.85	0.12

RAD - Error on 50% of the Training Dataset

	Before		Error in Anime Type		After	
	MSE on test dataset	Var on test dataset	Pct of detected error as movie	Pct detected error as tv series	MSE	Var on test dataset
Lasso Without HP	1.13	-27.79	26.125%	23.5%	0.89	0.16
Lasso With HP	1.12	-33.29	26.125%	23.5%	0.84	0.16

RAD - Error on 75% of the Training Dataset

	Before		Error in Anime Type		After	
	MSE on test dataset	Var on test dataset	Pct of detected error as movie	Pct detected error as tv series	MSE	Var on test dataset
Lasso Without HP	1.14	-16.99	25.25%	22%	0.86	0.14
Lasso With HP	1.14	-16.99	25.25%	22%	0.90	-0.04

Typo - Anime Genre

1) **Generating**

Create a Typo error in each tuple of data gener (eg. romantic → ronctzs)

2) **Detecting**

Using spelling checker python library to detect the misspelled words

3) **Repairing**

Correct the misspelled words, However, there are some strange words that could not correct them.

Typo - Error on 25% of the Training Dataset

	Before		Error in Anime Genre	After	
	MSE on test dataset	Var on test dataset	Pct of detected error as typo	MSE	Var on test dataset
Lasso Without HP	1.06	-20.27	60.94%	0.86	0.18
Lasso With HP	1.06	-20.27	60.94%	0.86	0.18

Typo - Error on 50% of the Training Dataset

	Before		Error in Anime Genre	After	
	MSE on test dataset	Var on test dataset	Pct of detected error as typo	MSE	Var on test dataset
Lasso Without HP	1.19	-31.46	73.5%	0.87	0.11
Lasso With HP	1.12	-34.64	73.5%	0.89	0.05

Typo - Error on 75% of the Training Dataset

	Before		Error in Anime Genre	After	
	MSE on test dataset	Var on test dataset	Pct of detected error as typo	MSE	Var on test dataset
Lasso Without HP	1.11	-23.10	85.1%	0.86	0.18
Lasso With HP	1.07	-19.55	85.1%	0.88	0.09

Gaussian Noise - Voting Members

1) Generating

Apply Gaussian Noise ($\mu=0$, $\sigma=10000$) on the members to check the effects of Gaussian signal on the distribution of the most correlated feature.

2) Detecting

Plotting the distribution shows us there is a noise in the data

3) Repairing

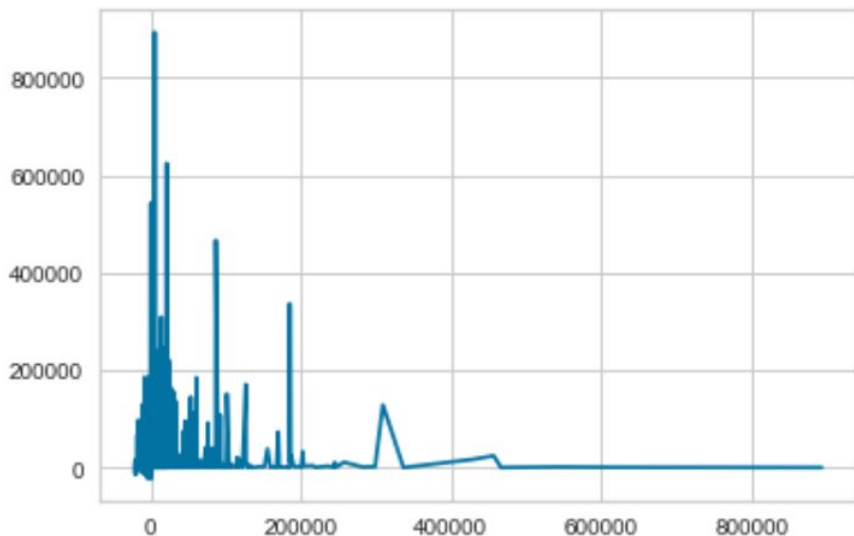
Using `savgol_filter` to filter the data from the noise

Gaussian Noise - Error on 30% of the Training Dataset

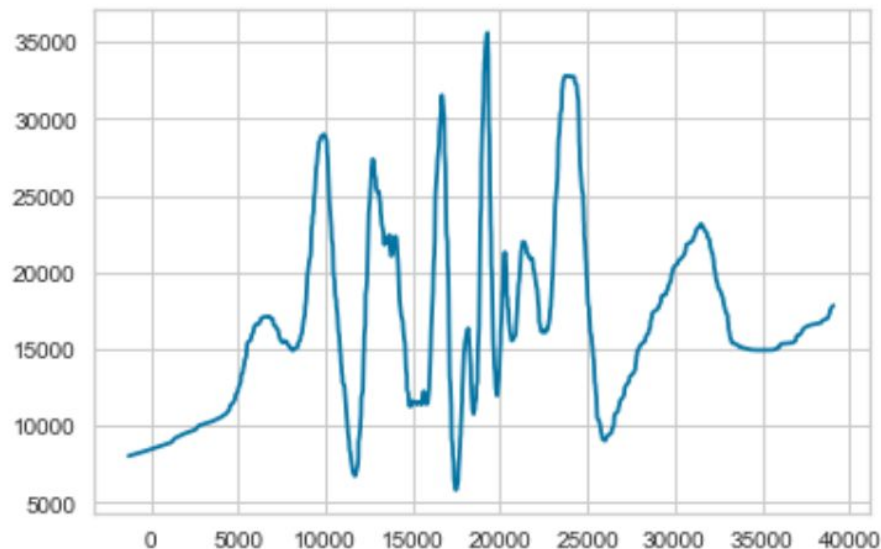
	Before		After	
	MSE on test dataset	Var on test dataset	MSE	Var on test dataset
Lasso Without HP	0.73	-8.37	0.94	-0.01
Lasso With HP	0.73	-8.37	0.94	0

Members Attribute Before and After Filtering - 30%

W-Gaussian Noise



Filter-Gaussian Noise

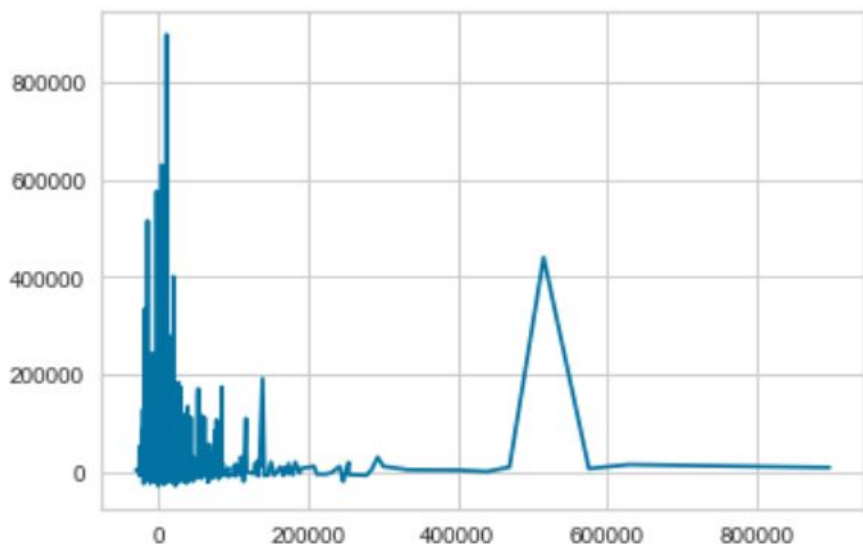


GN - Error on the whole of the Training Dataset

	Before		After	
	MSE on test dataset	Var on test dataset	MSE	Var on test dataset
Lasso Without HP	0.71	-5.57	0.96	-0.03
Lasso With HP	0.71	-5.57	0.95	0

Members Attribute Before and After Filtering - 100%

W-Gaussian Noise



Filter-Gaussian Noise

