

More Bikes: Experiments in Univariate Regression

Tim Lawson

December 21, 2023

1 Task description

The assignment is to predict the number of available bikes at 75 rental stations in three hours' time for a period of three months, beginning in November 2014, i.e., a supervised univariate regression problem. It is divided into three sub-tasks, which differ in the information that is available:

- **Sub-task 1.** The number of available bikes at each of the 75 stations for the month of October 2014. This sub-task may be approached by building a separate model for each station or a single model for all 75 stations.
- **Sub-task 2.** A set of linear models that were trained on the number of available bikes at each of a separate set of 200 stations for a year. For the first ten stations, this data is available for analysis but not training.
- **Sub-task 3.** Both of the above.

Sub-tasks 2 and, optionally, 3 require the use of ensemble methods. The predictions are evaluated by the mean absolute error (MAE) between the predicted and true numbers of available bikes over the period of three months, beginning in November 2014. The evaluation data is not available to participants but the score achieved on a held-out test set is reported on the task leaderboard. This report begins with a preliminary analysis of the data and then describes the approach taken to each sub-task and the cross-validation results obtained.

2 Data analysis

The data is recorded at hourly intervals. A summary of the features is given in fig. 1. Notably, the meteorological features are constant for all stations at a given timestamp. The 'profile' features, i.e., the features derived from the numbers of available bikes at preceding times, are not defined for the first week of instances at each station, except for `bikes_3h`. Naturally, the number of bikes available at a given station is bounded by zero and the number of docks at that station. The variances and pairwise correlations of the features are described in section 3.1.

3 Methods

Throughout this report, I used `scikit-learn` to conduct experiments (Pedregosa et al. 2011). In each case, preprocessing and feature selection were performed by *estimators* that implemented the *transformer* interface; prediction was performed by estimators that implemented the *predictor* interface; and estimators were composed into `Pipeline` objects over which hyperparameter search was performed (Buitinck et al. 2013, pp. 4–9).

Category	Feature	Data type	Kind
Station	station	int	ordinal
	latitude	float	
	longitude	float	
	docks	int	
Temporal	timestamp	int	ordinal
	year	int	
	month	int	
	day	int	
	hour	int	
	weekday	str	
	weekhour	int	
	is_holiday	bool	categorical
Meteorological	wind_speed_max	float	
	wind_speed_avg	float	
	wind_direction	float	
	temperature	float	
	humidity	float	
	pressure	float	
	precipitation	float	
Bikes	bikes_3h	int	
	bikes_3h_diff_avg_full	float	
	bikes_avg_full	float	
	bikes_3h_diff_avg_short	float	
	bikes_avg_short	float	
	bikes	int	

Figure 1: A summary of the features of the data. Except where indicated, the features are quantitative.

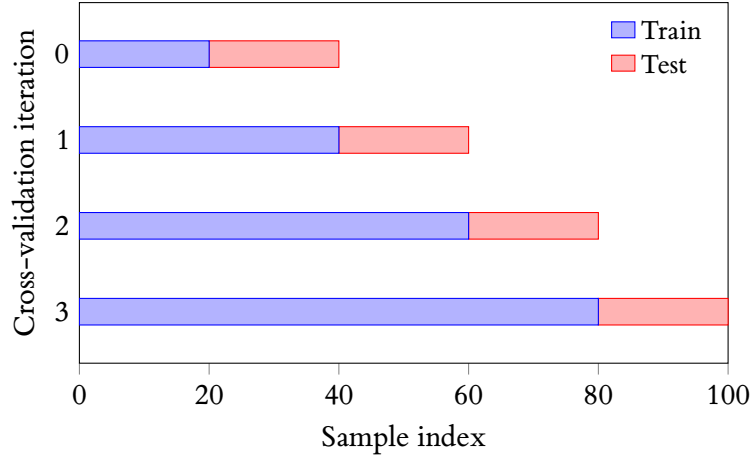


Figure 2: A visualisation of the nested time-series cross-validation behaviour, after the visualisation of `sklearn.model_selection.TimeSeriesSplit`.

Generally, standard k -fold cross-validation is disfavoured for time-series data due to the inherent correlation between successive folds (Bergmeir et al. 2018). Instead, nested time-series cross-validation¹ with ten folds was performed to evaluate the models. This behaviour is illustrated in fig. 2.

- Hyperparameter search.
- Evaluation metric.
- Statistical significance tests (paired t -tests).

3.1 Feature selection

In each case, zero-variance features were removed automatically from the data². The ‘station’ features (fig. 1) are constant for all instances at a given station; hence, for the first case of sub-task 1, the variances of these features are zero. Additionally, the precipitation feature is zero for all instances.

- Which features are most informative?
- Describe the distributions of the features and their correlations.
- Describe the distributions of the fraction of available bikes in terms of temporal features and the possible derived features.

References

- Bergmeir, Christoph, Rob J. Hyndman, and Bonsoo Koo (2018). “A Note on the Validity of Cross-Validation for Evaluating Autoregressive Time Series Prediction”. In: *Computational Statistics & Data Analysis* 120, pp. 70–83.
- Buitinck, Lars, Gilles Louppe, and Mathieu Blondel (2013). “API Design for Machine Learning Software: Experiences from the Scikit-Learn Project”. In: *ECML/PKDD 2013 Workshop: Languages for Data Mining and Machine Learning*.
- Pedregosa, Fabian et al. (2011). “Scikit-Learn: Machine Learning in Python”. In: *The Journal of Machine Learning Research* 12, pp. 2825–2830.

¹`sklearn.model_selection.TimeSeriesSplit`

²`sklearn.feature_selection.VarianceThreshold`