

ESTIMATING IMPORTANCE OF HIGHLY CORRELATED FEATURES USING MATRIX DECOMPOSITION

Anonymous authors

Paper under double-blind review

ABSTRACT

Hyperspectral images contain a large volume of source data that exhibits high correlations along neighboring spectral bands. This makes it necessary to select the most informative features among correlated groups of features to effectively solve various machine learning problems. A method of feature importance evaluation for hyperspectral image data is proposed. This method combines iterative training of decision tree classifiers based on spectral features with matrix decomposition to overcome sparsity. Decision trees provide intrinsic feature selection mechanism but only a small number of features are usually taken into account by the CART algorithm for training a single decision tree classifier instance. Furthermore when features are highly correlated (e.g., Pearson $\rho > 0.8$), tree-based methods like Random Forest or XGBoost arbitrarily assign importance to one feature while suppressing others, as they redundantly capture the same signal. The considered method is compared with several tree based methods for feature importance evaluation such as vanilla Gini impurity decrease and more complicated Boruta algorithm. The features are highlighted using a classification algorithm for thick cloud classification based on the marked-up satellite data. Classification accuracy testing based on significant features is performed for different types of surfaces for the set of several single images.

1 INTRODUCTION

Hyperspectral images (HSI) usually acquired from satellites capture light spectra in numerous narrow wavelength ranges. This approach boosts their informational content compared to conventional RGB images. In image processing tasks like cloud detection, monitoring, agriculture and environmental protection (Borzov & Uzilov, 2016), the large volume of HSI data and correlations across neighboring spectral bands require selection of the most informative features.

The existing methods of feature selection assess both feature properties and target variable relationships, covering forward selection and backward elimination methods, exhaustive search, and machine learning-based techniques (Guyon & Elisseeff, 2003). Feature importance evaluation can be done directly by analysing the pair relations between feature and target variable. This approach is generally called filter methods. Correlations between adjacent spectral bands gives folks not only to the idea of using dimensionality reduction algorithms for image compression but also to select relevant features (Myasnikov, 2017) from which PCA is the most popular (Zimichev et al., 2014). This approach along with direct analysis of the correlation matrix doesn't require additional labeling. Such methods as mutual information scores and ANOVA F-value are also very popular filter methods that takes into account labels but doesn't rely on training algorithms. Though these methods are computationally efficient they ignore feature interdependence which can lead to the selection of redundant features that are highly correlated with each other (Zhu et al., 2007). That is why, to select features more effectively, information from the results of HSI classification, in combination with the existing labeling of the source data, is needed. Machine learning (ML) models are widely used to evaluate the importance of features, especially in cases of feature correlations.

Model-specific methods for feature selection integrate the feature selection process directly into the model training phase. This approach offers good balance of performance and computational cost (Saito et al., 2018). There are wrapper and embedded model-specific methods. Wrapper methods identify the optimal subset of features by evaluating their different combinations using a specific

054 predictive model. Models such as the recursive feature elimination (RFE) select features iteratively,
055 maximizing the performance of the classification model. Although effective in finding optimal sub-
056 sets, these methods can be computationally demanding (Zubair et al., 2024). Embedded methods
057 integrate feature selection directly into the model training process itself, using mechanisms like reg-
058 ularization penalties to perform selection simultaneously with parameter estimation. Examples of
059 embedded methods include Lasso (Least Absolute Shrinkage and Selection Operator) and Ridge
060 regression (Fira et al., 2025), neural networks with learnable drop layer (Jiménez-Navarro et al.,
061 2024), sparse principal component analysis (Seghouane et al., 2019), tree-based methods (Tuv et al.,
062 2009).

063 Decision trees provide an intrinsic feature selection mechanism (Breiman et al., 1984), effectively
064 handling non-linear data and correlated features Kohavi & John (1997), and offering inherent ex-
065 plainability Mishra (2022). Decision tree based models like Random Forests and Gradient Boosting
066 provide feature importance scores based on how much a feature contributes to reducing impurity
067 (e.g., Gini impurity decrease) or variance at each split in the decision trees. More complicated
068 Boruta method determine feature relevance by comparing the importance of an original feature with
069 the importance of permuted counterparts known as shadow features and functions as a wrapper
070 algorithm built around the Random Forest classifier (Kursa & Rudnicki, 2010).

071 In general model specific methods give better results than filter methods but their significant draw-
072 back is that the selected set of features is inherently tied to the specific ML algorithm being used
073 (Islam et al., 2022). Stochastic nature of some ML models can lead to reproducibility issues, with
074 variations in feature importance rankings across different model training runs (Vos et al., 2024).

075 Model-agnostic methods offer greater flexibility and are widely used for interpreting complex
076 "black-box" models. Model-agnostic feature selection methods can be applied to any machine learn-
077 ing model, as they are independent of the model's internal workings (Khan et al., 2025). These meth-
078 ods assess the relevance of features based on their intrinsic properties and their relationship with the
079 target variable, without considering a specific predictive model. Examples of model-agnostic tech-
080 niques include methods based on eXplainable Artificial Intelligence (XAI) like Permutation Feature
081 Importance (PFI) (Flora et al., 2024) and SHapley Additive exPlanations (SHAP) (Lundberg & Lee,
082 2017). These techniques offer flexibility and can be used to compare the importance of features
083 across different models but still affected by feature correlations (Liang et al., 2024) especially PFI
084 (Salih, 2025).

085 The more sophisticated method of feature importance evaluation is used the more computing time
086 it requires. The more features are taken into account, the more their importance is affected by
087 multicollinearity. This paper is dedicated to the research of the possible approximation for feature
088 importance evaluated from sparse matrix that contains minimal set of values obtained by selection
089 of different feature subsets. We suggest the algorithm of feature importance recovery by means of
090 matrix factorization.

092 2 PROBLEM STATEMENT AND METHOD OF SOLUTION

095 While high-dimensional HSI data can be rich in information, it also introduces specific obstacles
096 that often reduces the effectiveness of ML models. A family of techniques designed to overcome the
097 curse of dimensionality are commonly addressed through a set of methods known as dimensionality
098 reduction. In this paper the dimensionality reduction is considered as feature selection problem.

099 The purpose of this work is to propose the algorithm for selecting a limited set of spectral bands
100 and derived features that yields the best prediction quality of dense cloud classification. A dedicated
101 set of labeled images of the HYPERION sensor with a spatial resolution of 30 m and a spectral
102 resolution of 10 m in the spectral range of 400–2500 nm was used. After converting the raw radiance
103 data into reflectance values and eliminating zero channels as well as those corresponding to strong
104 light absorption in water vapour, the selected hyperspectral images (HSI) were organised into a
105 training set. The features of this set comprised reflectance values from spectral channels 8–224,
106 combined with derived characteristics and other indices calculated on a pixel-by-pixel basis. Fig.
107 1 shows the mean spectral reflectance distribution for cloud and non-cloud pixels, suggesting that
a classification algorithm could be developed to differentiate between these pixels based on their

spectral characteristics. In this study, in addition to reflectance values, normalized indices such as NDVI (Huang et al., 2021), NDSI (Jin et al., 2022) and NDWI (Gao, 1996) were used.

We evaluate the importance of features by assess the cloud classification accuracy of a logistic regression model depending on the choice of spectral channels and derived features. The iterative training of ML models with different subset of features is used here to approximate target variable. After training of ML model the corresponding feature importance can be evaluated. By evaluating different sets of features used in the training process, we derive a sparse matrix representing the feature importances for each instance of the trained ML models. This approach yields the following observations:

1. Minimal subset of features is considered because of correlations between them;
2. Training ML models using as many features as possible is computationally expensive;
3. As a result of training ML models, the resulting importance matrix is sparse.

Logistic regression as linear model is strongly affected by features correlation. Decision trees are better but training with large feature subset is still computationally expensive. The resulting importance matrix is sparse, which gives rise to the problem of zero value reconstruction. Such reconstruction is needed to overcome possible feature importance underestimation for correlated features when using Gini impurity decrease as a metric of such evaluation. If we treat feature importance evaluation by approximating from a few values from the sparse matrix one can see that features with high importance ...

So we reformulated it by the following way. The feature importances are modeled as the following matrix decomposition:

$$\hat{T} = PQ^T, \quad (1)$$

where $\hat{T} (n_u \times n_i)$ is the predicted importances corresponding to trained ML model instance u and feature i , $P (n_u \times n_f)$ and $Q (n_i \times n_f)$ are latent factors that capture hidden preferences for features and train instances, respectively. The challenge to the matrix factorization problem is to find P and Q^T . Basically, such an algorithm is going to be used to find latent factors that represent intrinsic feature attributes in a lower dimension (the number n_f of latent factors is chosen beforehand). A learning approach is therefore developed to converge the decomposition results close to the observed importances as much as possible, while ensuring all importance values remain nonnegative. Additionally we introduce the feature bias matrix $\bar{\mu}$ as a correction of (1):

$$\hat{T} = \bar{\mu} + PQ^T. \quad (2)$$

The feature bias matrix $\bar{\mu}$ supposed to capture the tendency of features to have importance higher (or lower) than the average. So $\bar{\mu}$ measures a trained model tendency to systematically overestimate

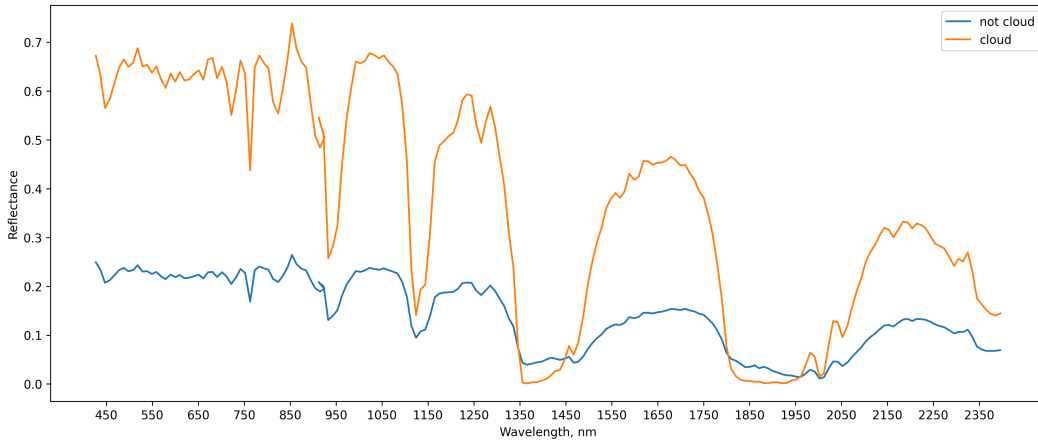


Figure 1: Mean spectral distribution of the reflectance corresponding to cloud and non-cloud pixels

or underestimate feature importances relative to the average across all trained ML model instances. Matrices P , Q and $\bar{\mu}$ can be obtained through a regularized optimization procedure:

$$\sum_u \sum_i \left((T_{u,i} - \hat{T}_{u,i})^2 + \lambda(\mu_{u,i}^2 + \|P_u\|^2 + \|Q_i\|^2) \right). \quad (3)$$

To bind classification accuracy with feature importance, the correlation between feature importances and classification error is added to the objective function:

$$\sum_u \sum_i \left((T_{u,i} - \hat{T}_{u,i})^2 + R_{u,i} + \lambda(\mu_{u,i}^2 + \|P_u\|^2 + \|Q_i\|^2) \right), \quad (4)$$

where λ is a regularization parameter, e_u is the classification error of trained ML instance and

$$R_{u,i} = \sum_u \sum_i (e_u - \bar{e}_u)(\hat{T}_{u,i} - \mu_{u,i}) = \sum_u (e_u - \bar{e}_u) \sum_i p_u q_i^T. \quad (5)$$

Here $e_u - \bar{e}_u$ is the deviation of classification error e_u from the average error of classification. The second term of equation (4) accounts for reducing the significance of features in the case of high classification error. If the classification error e_u is below average the importance of features corresponding to the trained ML instance u tends to increase; conversely, if e_u is above average, the corresponding importance values tend to decrease.

Stochastic Gradient Descent used here to solve the problem (2) of matrix factorisation is an optimization algorithm in which the model parameters (in this case, the bias $\bar{\mu}$ and the factor vectors) are repeatedly updated by adding the negative of the gradients calculated with respect to the function (4) being optimized. The algorithm essentially performs the following steps for a given number of iterations:

$$\begin{aligned} \mu_{u,i} &\leftarrow \mu_{u,i} + \gamma(\delta_{ui} - \lambda\mu_{u,i}) \\ p_u &\leftarrow p_u + \gamma((\delta_{ui} - 0.5 \cdot e_u) \cdot q_i - \lambda p_u) \\ q_i &\leftarrow q_i + \gamma((\delta_{ui} - 0.5 \cdot e_u) \cdot p_u - \lambda q_i) \end{aligned} \quad (6)$$

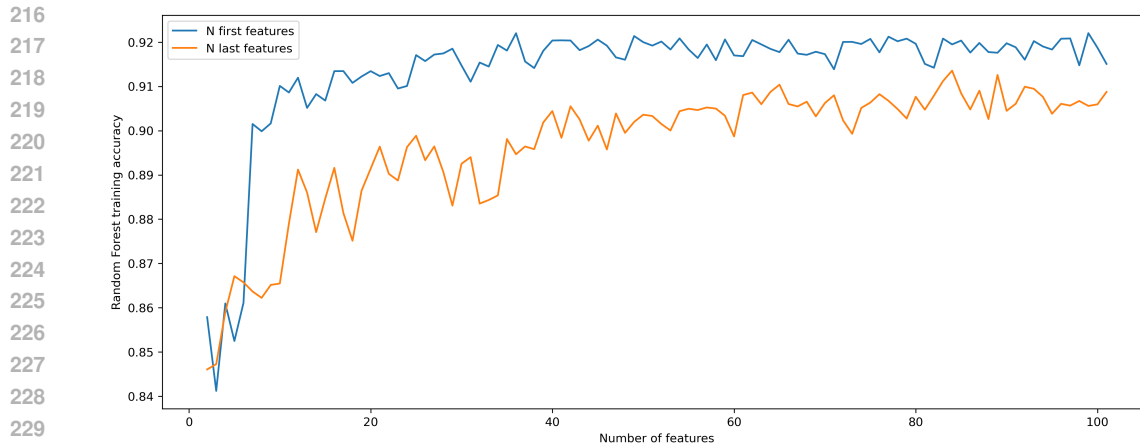
where γ is the learning rate and $\delta_{ui} = r_{ui} - \hat{r}_{u,i} = r_{u,i} - (\mu_{u,i} + p_u q_i^T)$ is the error made by the model for the pair (u, i) .

3 RESULTS AND DISCUSSION

The input data for the matrix factorization algorithm is feature importance statistics evaluated by training Decision Tree classifier for different values of hyperparameters. For Decision Tree classifier the decrease of Gini impurity is usually used to assess the feature importance. By considering the accuracy associated with different feature sets, we can calculate the initial correlation between feature importance and training accuracy and sort features according to this value of correlation. Let's refer to this as correlation importance. The most important features include NDWI index and the limited set of features from the NIR and the lower SWIR wavelength range. There is little to no similarity between the feature importance rankings obtained via the Boruta algorithm and those from the present study, with the exception of the consistently high importance assigned to the NDWI feature. This can be explained by the use of shadow features in feature analysis via the Boruta algorithm not used here. However Fig. 2 shows the graph of the random forest model's accuracy versus the selected set of features. By choosing features with higher importance scores (starting from the top of the list in descending order of importance), we obtain models with higher training accuracy.

4 CONCLUSION

The significant number of channels in HSI, combined with feature multicollinearity, leads to challenges in selecting machine learning models, reducing their accuracy and interpretability. To address this issue for thick cloud classification, decision trees are employed with a selection of a limited number of significant features based on the iterative exclusion algorithm. The proposed method applied



231 Figure 2: Random Forest training for the first and the last features by their correlation with training
232 accuracy of different Decision Tree classifiers
233

234
235 to find the set of features that yields better prediction accuracy. The most relevant features include
236 NDWI index, limited number of NIR bands and the lower part of SWIR spectrum range bands. De-
237 cision Tree model used to assess feature importance effectively handles correlated features avoiding
238 redundancy and can be combined with forward feature selection to achieve more robust statistics.
239 The final set of features selected after applying matrix factorization to overcome sparsity and tak-
240 ing into account the correlation between feature importance and model accuracy. This can be used
241 to construct a classifier for recognizing thick clouds based on their spectral characteristics. Such
242 classifier can be considered a baseline for more complex cloud classification models.

243 The ambiguity in feature importance rankings produced by traditional methods - such as permutation
244 feature importance, SHAP values, tree-based methods, and linear model coefficients - underscores
245 the need for a more empirically grounded approach. The proposed strategy of iterative feature
246 elimination with feature approximation via matrix decomposition offers an alternative solution to
247 this challenge. This approach not only reveals which features are truly critical for maintaining
248 performance but also accounts for feature interactions and helps identify their optimal subset. One
249 of the way to evaluate feature importance is to use Gini impurity decrease in combination with
250 matrix factorization to overcome sparsity. The point of the present study is to use information about
251 the classification accuracy to evaluate feature importance. This was done by the incorporation of the
252 special correlation term in the objective function for optimization.

253 ACKNOWLEDGMENTS

254 REFERENCES

- 255
256 S.M. Borzov and S.B. Uzilov. *J. Comput. Technol.*, 21(1):40 – 48, 2016.
257
258 L. Breiman, J. Friedman, R.A. Olshen, and C.J. Stone. *Classification and Regression Trees*. Chap-
259 man and Hall/CRC., 1984.
260
261 Monica Fira, Liviu Goras, and Hariton-Nicolae Costin. Evaluating sparse feature selection methods:
262 A theoretical and empirical perspective. 15(7), 2025. doi: 10.3390/app15073752.
263
264 Montgomery L. Flora, Corey K. Potvin, Amy McGovern, and Shawn Handler. A machine learning
265 explainability tutorial for atmospheric sciences. *Artificial Intelligence for the Earth Systems*, 3(1):
266 e230018, 2024. doi: 10.1175/AIES-D-23-0018.1. URL <https://journals.ametsoc.org/view/journals/aies/3/1/AIES-D-23-0018.1.xml>.
267
268 B.-C. Gao. NdwI—a normalized difference water index for remote sensing of vegetation liquid
269 water from space. *Remote Sensing of Environment*, 58(3):257–266, 1996. doi: [https://doi.org/10.1016/S0034-4257\(96\)00067-3](https://doi.org/10.1016/S0034-4257(96)00067-3).

- 270 Isabelle Guyon and André Elisseeff. An introduction of variable and feature selection. *J. Machine*
271 *Learning Research Special Issue on Variable and Feature Selection*, 3:1157 – 1182, 2003. doi:
272 10.1162/153244303322753616.
- 273 Sha Huang, Lina Tang, Joseph P. Hupy, Yang Wang, and Guofan Shao. A commentary review on
274 the use of normalized difference vegetation index (ndvi) in the era of popular remote sensing.
275 *Journal of Forestry Research*, 32(1):1–6, 2021. doi: 10.1007/s11676-020-01155-1.
- 276 Md Rashedul Islam, Aklima Akter Lima, Sujoy Chandra Das, M. F. Mridha, Akibur Rahman
277 Prodeep, and Yutaka Watanobe. A comprehensive survey on the process, methods, evaluation, and
278 challenges of feature selection. *IEEE Access*, 10:99595–99632, 2022. doi: 10.1109/ACCESS.
279 2022.3205618.
- 280 M J Jiménez-Navarro, M Martínez-Ballesteros, I S Brito, F Martínez-Álvarez, and G Asencio-
281 Cortés. Embedded feature selection for neural networks via learnable drop layer. *Logic Journal*
282 *of the IGPL*, 33(5):jzae062, 2024. doi: 10.1093/jigpal/jzae062.
- 283 Donghyun Jin, Kyeong-Sang Lee, Sungwon Choi, Noh-Hun Seong, Daeseong Jung, Suyoung Sim,
284 Jongho Woo, Uujin Jeon, Yugyeong Byeon, and Kyung-Soo Han. An improvement of snow/cloud
285 discrimination from machine learning using geostationary satellite data. *International Journal of*
286 *Digital Earth*, 15(1):2355–2375, 2022. doi: 10.1080/17538947.2022.2152886.
- 287 Adam Khan, Asad Ali, Jahangir Khan, Fasee Ullah, and Muhammad Faheem. Exploring consistent
288 feature selection for software fault prediction: An xai-based model-agnostic approach. *IEEE*
289 *Access*, 13:75493–75524, 2025. doi: 10.1109/ACCESS.2025.3558913.
- 290 Ron Kohavi and George H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97
291 (1):273–324, 1997. ISSN 0004-3702. doi: 10.1016/S0004-3702(97)00043-X. URL [https://](https://www.sciencedirect.com/science/article/pii/S000437029700043X)
292 www.sciencedirect.com/science/article/pii/S000437029700043X. Rele-
293 vance.
- 294 Miron B. Kurşa and Witold R. Rudnicki. Feature selection with the boruta package. *Journal of*
295 *Statistical Software*, 36(11):1–13, 2010. doi: 10.18637/jss.v036.i11.
- 296 Annie Liang, Thomas Jemielita, Andy Liaw, Vladimir Svetnik, Lingkang Huang, Richard Baum-
297 gartner, and Jason M. Klusowski. Challenges in variable importance ranking under correlation,
298 2024.
- 299 Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions. 12 2017. doi:
300 10.48550/arXiv.1705.07874.
- 301 P. Mishra. *Practical explainable AI using python: Artificial Intelligence Model Explanations Using*
302 *Python-Based Libraries, Extensions, and Frameworks*. New York: Apress, 2022. doi: 10.1007/
303 978-1-4842-7158-2.
- 304 E.V. Myasnikov. Hyperspectral image segmentation using dimensionality reduction and clas-
305 sical segmentation approaches. *Comput. Opt.*, 41(4):564–572, 2017. doi: 10.18287/
306 2412-6179-2017-41-4-564-572.
- 307 Shota Saito, Shinichi Shirakawa, and Youhei Akimoto. Embedded feature selection using prob-
308 abilistic model-based optimization. In *Proceedings of the Genetic and Evolutionary Computa-*
309 *tion Conference Companion*, pp. 1922–1925. Association for Computing Machinery, 2018. doi:
310 10.1145/3205651.3208227.
- 311 Ahmed M. Salih. Re-visiting explainable ai evaluation metrics to identify the most informative
312 features, 2025.
- 313 Abd-Krim Seghouane, Navid Shokouhi, and Inge Koch. Sparse principal component analysis with
314 preserved sparsity pattern. *IEEE Transactions on Image Processing*, 28(7):3274–3285, 2019. doi:
315 10.1109/TIP.2019.2895464.
- 316 Eugene Tuv, Alexander Borisov, George Runger, and Kari Torkkola. Feature selection with ensem-
317 bles, artificial variables, and redundancy elimination. *Journal of Machine Learning Research*, 10:
318 1341–1366, 07 2009. doi: 10.1145/1577069.1755828.
- 319
320
321
322
323

324 Gideon Vos, Liza van Eijk, Zoltan Sarnyai, and Mostafa Rahimi Azghadi. Stabilizing machine
325 learning for reproducible and explainable results: A novel validation approach to subject-specific
326 insights, 2024.

327 Zexuan Zhu, Yew-Soon Ong, and Manoranjan Dash. Wrapper-filter feature selection algorithm
328 using a memetic framework. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cy-*
329 *bernetics)*, 37(1):70–76, 2007. doi: 10.1109/TSMCB.2006.883267.

330
331 E.A. Zimichev, N.L. Kazanskiy, and P.G. Serafimovich. Spectral-spatial classification with
332 k-means++ partitional clustering. *Comput. Opt.*, 38(2):281–286, 2014. doi: 10.18287/
333 0134-2452-2014-38-2-281-286.

334
335 Iqbal Muhammad Zubair, Yung-Seop Lee, and Byunghoon Kim. A new permutation-based method
336 for ranking and selecting group features in multiclass classification. *Applied Sciences*, 14(8),
337 2024. doi: 10.3390/app14083156.

338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377