

Mixing height derivation from aerosol lidar using machine learning: KABL and ADABL algorithms

Thomas Rieutord¹, Sylvain Aubert², and Tiago Machado^{1,2}

¹Centre National de Recherches Meteorologiques, Université de Toulouse, Météo-France, CNRS, Toulouse, France

²Direction des Systèmes d'Observation, Météo-France, Toulouse, France

Correspondence: Thomas RIEUTORD (thomas.rieutord@meteo.fr)

Abstract. The atmospheric boundary layer height (BLH) is a key parameter in several meteorological applications such as air quality forecasts. A common method to measure BLH is via aerosol lidar, where a strong decrease in the backscatter signal indicates the top of the boundary layer. This paper describes and compares two machine learning methods, the K-means algorithm and the AdaBoost algorithm to derive the BLH from backscatter profiles. The K-means for Atmospheric Boundary Layer (KABL) and AdaBoost for Atmospheric Boundary Layer (ADABL) algorithm codes used in this study are free and open source. Both methods are compared to the lidar manufacturer's algorithm and to reference BLHs derived from colocated radiosonde data. The radiosonde data were used as the reference for all methods. A comparison was carried for a two-year period (2017-2018) for two Météo-France operational network sites (Trappes and Brest). A large discrepancy in the results was observed between the two sites. At the Trappes site, KABL and ADABL outperformed the manufacturer's algorithm, while the performance was clearly reversed at the Brest site. We conclude that ADABL is a promising algorithm but has training issues that need to be resolved, KABL has a lower performance than ADABL but is much more versatile, and the manufacturer algorithm is performing well with little tuning but is not open-source.

1 Introduction

The atmospheric boundary layer is the lowest part of the troposphere that is directly influenced by surface forcings. As a high-impact area where most human activities take place, all pollutants emitted from the ground are diluted in this layer. The key parameter used to model this dilution is the depth of this layer, i.e., the boundary layer height (BLH). Because BLH can vary from a few tens meters to about 2 km within a single day, the amount of dilution can vary considerably and result in air quality warnings (Stull, 1988; Dupont et al., 2016). In addition, BLH is one of the largest sources of uncertainty in air quality models (Mohan et al., 2011) and there is a need to better evaluate this parameter (Arciszewska and McClatchey, 2001). In numerical weather prediction models, physical processes change within the boundary layer (Seity et al., 2011). Therefore it is important to compare BLH calculated in models with that derived from measurements.

However, measuring BLH is not straightforward. As stated in Seibert et al. (2000), there are no systems that match all the requirements to make reliable BLH estimations. The best BLH estimation can be achieved via the synergetic use of multiple instruments. However, adding instruments limits the number of sites where estimations can be made. In this paper we focus

25 on a single instrument, aerosol lidar (see Sect. 2.1.1 for more information), already widely used (Haeffelin et al., 2012). The boundary layer is detected by the decrease of the lidar signal at its top.

However this decrease can be blurred or perturbed by other strong signals (e.g., clouds, residual layers, and small scale structures) and instrumental noise. For these reasons, numerous studies exist concerning the derivation of BLH from aerosol lidar. Melfi et al. (1985) use a simple thresholding of the signal. Others methods are based on derivatives. For example, Hayden
30 et al. (1997) take the minimum of the gradient, Menut et al. (1999) use the height where the second derivative is zero (the inflection point) as well as the variance of the signal, Senff et al. (1996) use the derivative of the logarithm of the backscattered intensity along the height. One of the most used methods is the wavelet covariance transform which searches for the maximum in the convolution between the signal profile and a Haar wavelet (Gamage and Hagelberg, 1993; Cohn and Angevine, 2000; Brooks, 2003). More recent studies have been based on backscatter signal analysis such as STRAT (Morille et al., 2007) and
35 CABAM (Kotthaus and Grimmond, 2018). Graph theory has also been used to impose continuity constraints (both vertically and in time) in BLH estimations, e.g., Pathfinder (De Bruine et al., 2017). Inspired by image processing, some methods use Canny edge detection in addition to backscatter signal analysis (Morille et al., 2007; Haeffelin et al., 2012). STRAT and Pathfinder have also been merged into PathfinderTURB by Poltera et al. (2017). These studies demonstrate that deriving BLH from aerosol lidar is still an open area of research.

40 In addition, artificial intelligence (AI) has reemerged in the last decade because of the simultaneous increase in the amount of available data and computational power. Both have reached levels that enable previously impossible applications. AI is capable of tackling complex classification problems, especially in image classification (Krizhevsky et al., 2012). Such breakthroughs were made possible by deep convolutional neural networks (LeCun et al., 2015); however, AI encompasses many other techniques that also benefit from larger data and increased computational power (Besse et al., 2018). In this paper, we explore how
45 the derivation of BLH from backscatter profiles can be formulated as a classification problem and how appropriate algorithms can be applied to solve this problem. Toledo et al. (2014) have already described a method that falls into the scope of AI. They used unsupervised learning to classify whether the measurement points were within the boundary layer. This method has yielded convincing results in previous studies (Toledo et al., 2017; Caicedo et al., 2017; Rieutord et al., 2014) and is pursued here with the K-means for Atmospheric Boundary Layer (KABL) algorithm. KABL was extensively tested and is shared via an
50 open-source code. In addition, we test an alternative adaptive boosting (AdaBoost) machine learning algorithm, the AdaBoost for Atmospheric Boundary Layer (ADABL) algorithm. Both algorithms classify whether the measurement points are inside or outside of the boundary layer; however, ADABL learns the characteristics of both groups from a training set. The training set consists of atmospheric boundary layer identifications made by human experts, which is acknowledged as being more reliable than available automatic methods (Seibert et al., 2000). Such supervised algorithms make it possible to automatically reproduce
55 human expertise in boundary layer identification. To our knowledge, this is the first time that a supervised algorithm has been applied to this problem. This study is of practical interest because it includes the publication of the source code, which only uses free software.

In Sect. 2, we describe the data used in this study, i.e., the lidar data in the algorithm inputs, reference radiosonde data, and ancillary data used to sort the meteorological conditions. In Sect. 3, we describe the two machine learning algorithms (KABL

60 and ADABL) and the procedures used to evaluate them. In Sect. 4, we present the results of our study, which consists of a sensitivity analysis of the KABL algorithm, a comparison of the methods with the radiosonde data over a two-year period, and a case study. In Sect. 5, we discuss the results, limitations and prospects of our study. The final section is dedicated to the conclusions that can be drawn from our study.

2 Material

65 Our study used data from Météo-France operational network. We used colocated radiosonde and aerosol lidar data over two sites: Brest (a coastal city in an extreme Western region of France) and Trappes (a sub-urban area of Paris in an inland region of France). The dataset spanned two years: 2017 and 2018. A case study was taken conducted on August 2, 2018, for the Trappes site.

2.1 Lidar data

70 2.1.1 Lidar network

Since 2016, Météo-France has deployed a network of six automatic backscatter lidars to help the Volcanic Ash Advisory Center of Toulouse characterize volcanic ash and aerosol layers. One of the six sensor can be quickly redeployed at a more suitable geographic location depending on the transport event being tracked. The network, fully operational since April 2017, is continuously functioning and has detected aerosol events at an altitude up to 17 km. It is part of the wider automatic lidar and ceilometers network of the E-PROFILE program described in Haefele et al. (2016).

Two sampling sites in this network were selected: Brest (48.444° N, 4.412° W, 94 m above sea level) and Trappes (48.773 N, 2.0124 E, 166 m above sea level). Both sites are equipped with a Mini Micro Pulse LiDAR (MiniMPL), built by Sigma Space Corporation with an exterior casing provided by Envicontrol. A typical MiniMPL unit from Météo-France network is shown in figure 1. MiniMPL is a compact version of the standard micro pulse lidar systems deployed in the NASA global lidar network MPLNET. A comprehensive description of MiniMPL can be found in Ware et al. (2016).

2.1.2 Data processing

MiniMPL acquires profiles of atmospheric backscattering at high frequency (2500 Hz) using a low energy pulse (3.5 μ J) emitted by a Nd:YAG laser at 532 nm. The profiles are acquired in photon-counting mode and, in our present configuration, averaged over 5 min and 30 m vertical resolution bins. The instrument uses a monostatic coaxial design where the laser beam and the receiver optics share the same axis. Because of geometrical limitations, only a fraction of the signal can be recovered in the near field. Therefore, in our system, the first usable data are available at 120 m above ground level.

The instrument has polarization capabilities with the collection of photons on two different channels (for more details see Flynn et al. (2007)); the measured raw signals on the "copolarized" and "crosspolarized" channels are respectively *co* and *cr* suffixed. These raw signals are processed to obtain the quantity of interest, i.e., the range corrected signal *RCS*, also called



Figure 1. Typical MiniMPL unit from Météo-France network

90 the normalized relative backscatter. This industrial processing consists of several procedures including background, overlap, afterpulse, and dead-time corrections. A comprehensive description of the processing is given in Campbell et al. (2002). The "copolarized" and "crosspolarized" range corrected signals, RCS_{co} and RCS_{cr} , respectively, as delivered by the industrial software, are used as predictors for the machine learning algorithms described in Sect. 3.

95 The raw data type and format depends on the instrumental device used. To make the algorithms usable on other devices, we converted the files to a normalized format using the *raw211* routine and then used these files as the algorithm input. *raw211* was developed by the Site Instrumental de Recherche par Télédétection Atmosphérique and is publicly available¹

2.2 Radiosonde data

The algorithms were evaluated with respect to radiosonde (RS) estimations. Météo-France operates several RS sites for the WMO Global Observing System. Two RS sites are colocated with the lidars of Brest and Trappes. They are equipped with a
100 Meteomodem robotsonde and typically launch a Meteomodem M10 sonde at 11:15 UTC and 23:15 UTC every day.

Many methods exist to derive BLH from RS data, and several of which have been used in the literature. Some of these methods are listed below.

- Parcel method: BLH is the height at which the profile of the potential temperature θ reaches its ground value.
- Humidity gradient method: BLH is the height at which the gradient of the relative humidity is strongly negative.
- 105 – Bulk Richardson number method: BLH is the height at which the bulk Richardson number exceeds 0.25 (this threshold varies among studies).

¹<https://gitlab.in2p3.fr/ipsl/sirta/raw211>

- Surface-based inversion: BLH is the height at which the gradient temperature profile reaches zero.
- Stable layer inversion: BLH is the height at which the gradient of the potential temperature profile reaches zero.

Hennemuth and Lammert (2006) used the parcel and humidity gradient methods. Collaud Coen et al. (2014) uses all the techniques mentioned above and recommend the bulk Richardson number method for all cases. Guo et al. (2016) used the bulk Richardson number for a two-year climatology. Seidel et al. (2010) compared the parcel, the humidity gradient, the surface-based inversion methods and other methods over a period of 10 years at 505 sites worldwide. Seidel et al. (2012) compared several methods and recommended the bulk Richardson number method.

Following the recommendations in Figure 10 of Seibert et al. (2000), we chose to derive BLH using the parcel method for the 11:15 UTC sounding and bulk Richardson number for the 23:15 UTC sounding.

2.3 Ancillary data

Ancillary data were used to describe the meteorological situations at the observation sites. These data were not used by the machine learning algorithms. All the instruments were colocated with the lidar and radiosonde launchings.

- Rain gauges were used to detect rain events.
- Vaisala Ceilometer CL31 instruments were used to detect the cloud base height and distinguish cases with clouds on top or inside the boundary layer. Even though the MiniMPL is capable of detecting clouds, we relied on the cloud detection with the CL31, because MiniMPL algorithm was found to report non-existent clouds.
- Scatterometers were used to estimate the visibility and detect the occurrence of fog.

3 Machine learning methods

Machine-learning techniques are categorized into two broad families: supervised learning (mimicking a reliable reference) and unsupervised learning (learning without a reference) (Hastie et al., 2009). First, we present the supervised algorithm leading to ADABL. Then, we present the unsupervised learning leading to KABL.

3.1 Supervised learning method

Supervised methods learn from a reference. Such methods are divided into two families: classification, which aims to find the frontiers between groups, and regression, which aims to approximate a function. In this study, we treat the BLH derivation as a classification problem. We need to classify the measurement points of the lidar into two classes: ‘boundary layer’ or ‘free atmosphere’. Then, the highest point of the ‘boundary layer’ class indicates the BLH estimate. Boosting algorithms are a very powerful family of algorithms, that were developed for classification but can also be used for regression (Hastie et al., 2009). In particular, the AdaBoost algorithm is designed for binary classification (Freund and Schapire, 1997), and is therefore well suited to our problem.

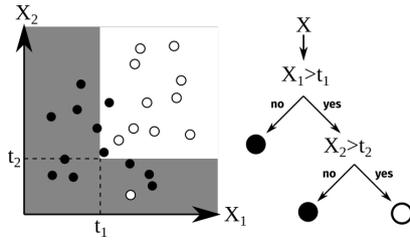


Figure 2. Illustration of binary classification with decision trees on two-dimensional artificial data.

3.1.1 AdaBoost algorithm

Let us consider the following problem. We have N vectors $x_i \in \mathbb{R}^p$ (here $p = 4$: seconds since midnight, height above ground, copolarized channel and cross-polarized channel), and for each vector, we have a binary indicator $y_i \in \{-1, 1\}$ (-1 for ‘boundary layer’, 1 for ‘free atmosphere’). From the sample $(x_i, y_i)_{i \in \llbracket 1, N \rrbracket}$, where $\llbracket 1, N \rrbracket$ is the ensemble of integers from 1 to N , we want to predict the output indicator y_{new} of any new vector x_{new} . To do so, we must find a rule based on the x_{new} coordinate values (the features) to cast it into the appropriate class. Decision tree classifiers (Breiman et al., 1984) perform this casting one feature at a time. For example, in Figure 2, there are black and white points in a two-dimensional space. The black points are mostly located where X_1 is low, hence the rule "if $X_1 < t_1$, then the point is black." However, in the other region, where $X_1 > t_1$, there are still some black points, all with low X_2 . Therefore, we add the rule "if $X_2 < t_2$, then the point is black, else it is white." Decision trees are classifiers made up of such "if" statements with various depths and thresholds. The deeper the tree, the more accurate the border but the more complex the decision and the longer it takes to train. Deep trees are strongly subject to overfitting and are less efficient than other methods. However, shallow decision trees are valuable because of their simplicity and their speed, even though their performances are quite limited (Hastie et al., 2009). They are often used as *weak learners*, that is, classifiers with poor (but better than random) performances and are very simple (Freund and Schapire, 1997). In this study, weak learners in AdaBoost are trees with a maximum depth of five (a maximum of five forks between the root and leaves).

AdaBoost is based on decision tree classifiers. It aggregates these classifiers to determine the most accurate border. The concept behind AdaBoost is illustrated in Figure 3. First, a shallow decision tree is fitted to the entire dataset using the Classification and Regression Tree (CART) algorithm (Hastie et al., 2009). All points have the same weight in this first step. Some points in the dataset are misclassified, and the error of the classifier is the weighted average of the misclassified points. Another shallow decision tree is then fitted on a resampled dataset where the previously misclassified points are over-represented. This new tree has new misclassified points, that will be over-represented in the training of the next tree, and so on, up to the specified number of trees ($M = 200$ in our case). The detailed algorithm is described in Hastie et al. (2009), algorithm 10.1, and in Schapire (2013).

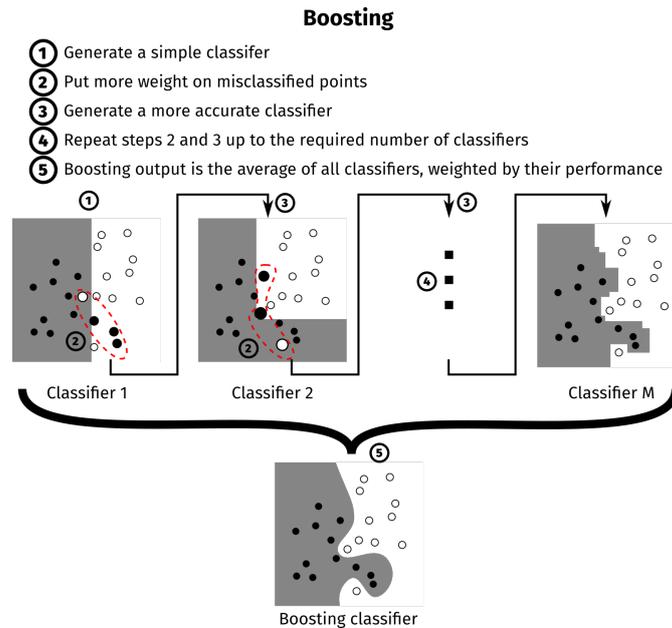


Figure 3. Illustration of boosting on two-dimensional artificial data with two classes.

160 3.1.2 Training of the algorithm

Such an algorithm needs to be trained using a trustworthy reference. On days where the boundary layer is easily visible for a human expert, the top of the boundary layer can be drawn by hand; all points below this limit are in the class ‘boundary layer’ and all points above this limit are in the class ‘free atmosphere’.

In this study, two days were labelled by hand. These two days were chosen because the boundary layers on these days were easily visible; the two labelled days were at different sites at different seasons. The first labelled day was a clear summer day in Trappes, shown in Figure 4 (left); a stable boundary layer is present near the ground during the night, topped by a residual layer and few clouds between 02:00 UTC and 04:00 UTC. A mixed layer started to develop at 09:00 UTC and remained approximately 2000 m for the rest of the day. At approximately 22:00 UTC, a new stable layer appeared to develop near the ground; however, it is not very clear where this layer started and what its extent was. The second labelled day was a clear winter day in Brest, shown in Figure 4 (right): a stable boundary layer was present near the ground during the night, topped by a residual layer, which was shallower than what was observed at the Trappes site. The mixed layer started to develop at 09:00 UTC and remained at approximately 1000 m with the height of the layer gradually decreasing throughout the day. At approximately 17:00 UTC, aerosols appeared to accumulate in a thin layer close to the ground, therefore we chose to drop BLH to that level.

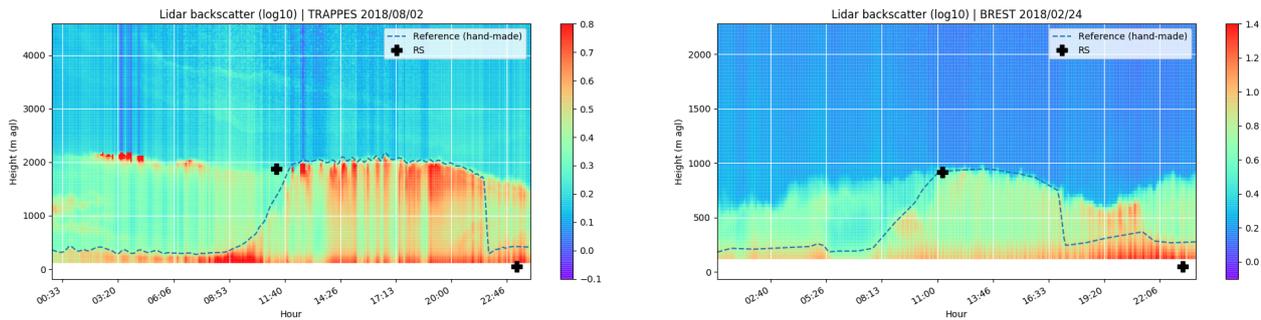


Figure 4. Hand-drawn references and radiosonde estimates overlaying the lidar range-corrected intensity signal for two days: August 2, 2018, at the Trappes site (left) and February 24, 2018, at the Brest site (right).

175 The coordinates of the points on the curves of the hand-drawn BLHs were obtained using the VGG Image Annotator software². Then, the output curves were interpolated with a cubic spline to match the lidar temporal resolution. Given the resolution of the lidar, this method of labelling the data results in $N = 86400$ individuals in total.

3.1.3 Retained configuration

180 Four predictors were used: the two lidar channels, time (number of seconds since midnight) and altitude (meters above ground level). The ADABL configuration used was

- Weak learner: decision tree of depth five;
- Number of weak learners: 200; and
- Predictors: time, altitude, RCS_{co} , and RCS_{cr} .

185 This configuration was chosen because more complex classifiers do not necessarily improve the performance. The computation time of the algorithm was still reasonable: training took 23s on the full dataset and predicting BLH for a full day took 3.7s. AdaBoost was chosen after a comparison of multiple classification algorithms, i.e., random forests, nearest neighbors, decision trees, and label spreading (study not shown here). The benchmark score was the accuracy as measured by the percentage of individuals that were well classified. The accuracy was estimated by group K-fold, where labelled data-sets are grouped into chunks of four consecutive hours, one group was used as a testing set and all the rest as a training set. This operation was
 190 repeated until each group was used as the testing set. The resulting accuracy was 96%. However, this figure overestimates the generalization ability of AdaBoost. A more correct estimation would be obtained with an independent validation set (e.g. a new labelled day). An independent validation set was not used here because the accuracy was only used to discriminate between the classification algorithms.

²Publicly available online at <https://www.robots.ox.ac.uk/~vgg/software/via/via-1.0.6.html>.

It is possible to quantify the relative importance of the predictors (Breiman et al., 1984; Hastie et al., 2009). After training, the relative importance of the time, $RC_{S_{co}}$, $RC_{S_{cr}}$, and altitude predictors was 30.3%, 28.4%, 26.5%, and 14.8%, respectively.

3.2 Unsupervised learning methods

Unsupervised methods aim to identify groups in the data. In our case, we want to identify the group ‘boundary layer’. The BLH estimate is then the upper boundary of this group. Two unsupervised learning algorithms were tested: K-means and expectation-maximisation (EM).

3.2.1 K-means algorithm

The K-means algorithm is a well proven and commonly used algorithm for data segmentation (Jain et al., 1999; Pollard et al., 1981) and consists of three steps.

1. Initialization: K centroids m_1, \dots, m_K are initialized at random locations inside the feature space.
2. Attribution: The distances from all points to all centroids $(d(x_i, m_k))_{k \in \llbracket 1, K \rrbracket, i \in \llbracket 1, N \rrbracket}$ are computed, and points are attributed to the closest centroid:

$$C(i) = \operatorname{arg\,min}_k \{d(x_i, m_k)\}.$$
3. Update: The centroids are re-defined as the average point of the cluster: $m_k = \frac{\sum_{i=1}^N x_i \mathbf{1}_{C(i)=k}}{\sum_{i=1}^N \mathbf{1}_{C(i)=k}}$

Steps 2 and 3 are repeated until the centroids stop moving. It has been shown that this algorithm converges to a local minimum of the intra-cluster variance Selim and Ismail (1984). Figure 5 (left) illustrates this algorithm.

3.2.2 EM algorithm

The EM algorithm assumes that each group $k \in \llbracket 1, K \rrbracket$ is generated by a Gaussian distribution (μ_k, Σ_k) . The algorithm iteratively estimates the parameters $\hat{\mu}_k, \hat{\Sigma}_k$ and the *responsibility* for each Gaussian $\hat{\gamma}_k^i$, where the *responsibility* is the probability of the point x^i being generated by the k -th Gaussian. Points are then attributed to the group with the highest responsibility: $C(i) = \operatorname{arg\,max}_k (\hat{\gamma}_1^i, \dots, \hat{\gamma}_K^i)$ Figure 5 (right) illustrates this algorithm.

The K-means and EM algorithms are very similar. If we assume that all Gaussians have the same fixed variance and that this variance tends to zero, the EM and K-means algorithms are the same. However, K-means does not rely on a Gaussian assumption.

3.3 KABL flowchart

The parameters of the KABL software are detailed here:

- **algo**: The applied machine learning algorithm. Possible values are
 - ‘gmm’ for the EM algorithm (Gaussian mixture model) and

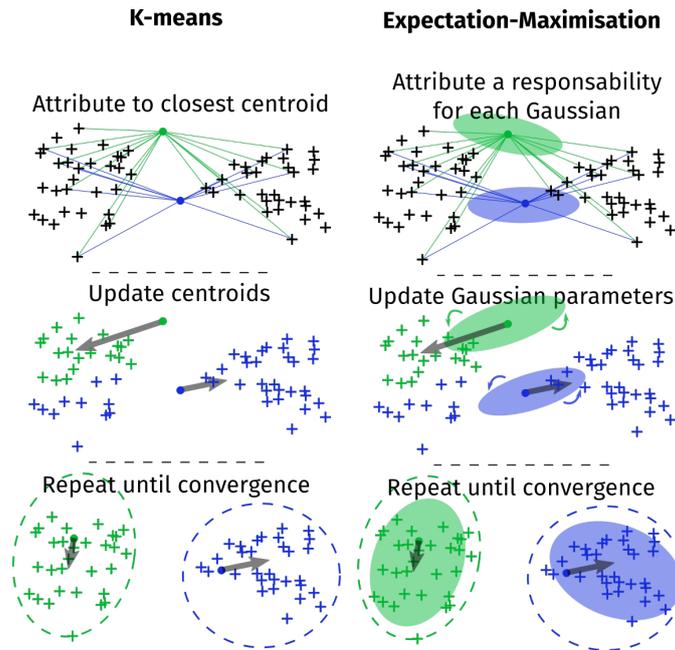


Figure 5. Illustration of the K-means and expectation-maximisation algorithm on two-dimensional artificial data and two clusters.

- ‘kmeans’ for the K-means algorithm
- **classif_score:** The internal score used to automatically choose the number of clusters (only used when `n_clusters=‘auto’`). See Sect. 3.4 for a description of the internal scores.
- 225 – **init:** initialisation strategy for both algorithms. Three choices are available:
 - ‘random’: randomly pick an individual as starting point (both Kmeans and EM);
 - ‘advanced’: use a more sophisticated initialization (kmeans++ for Kmeans (Arthur and Vassilvitskii, 2007) and the output a K-means pass for EM); and
 - ‘given’: start at explicitly passed point coordinates
- 230 – **max_height:** The height (meters above ground level) at which the profiles are cut.
- **n_clusters:** The number of clusters to be formed (between two and six). This is either explicitly given or determined automatically to optimize the score given in `classif_score`.
- **n_inits:** The number of repetitions of the algorithm. When this number is larger, the algorithm is more likely to find the global optimum but requires more time.

- 235 – **n_profiles**: The number of profiles concatenated prior to the application of the algorithm. For example, if `n_profiles=1`, only the current profile is used. If `n_profiles=3`, the current profile and the two previous profiles are concatenated and input to the algorithm.
- **predictors**: The list of variables used in the classification. These variables can be different at night and during the day. For both time periods, the variables can be chosen from
- 240 – $RC S_{co}$: the copolarized range-corrected backscatter signal; and
- $RC S_{cr}$: the cross-polarized range-corrected backscatter signal

A simplified flowchart of KABL is shown in Figure 6 and the parameters of the KABL software are highlighted in bold in the following explanation of the KABL algorithm. A netCDF file generated by the *raw2ll* software needs to be provided as input data to KABL. The data, namely, the altitude vector z (size N_z), the time vector t (size N_t), the range-corrected signals

245 $RC S_{co}$ and $RC S_{cr}$ ($N_t \times N_z$ matrices), are extracted from this file. Such data are prepared to fulfil the machine-learning algorithms requirements. For each time, the **n_profiles** last profiles are extracted. Then, the data they contain are normalized (by removing the mean and dividing by the standard deviation); this provides a matrix X ($N \times p$, where $N = \mathbf{n_profiles} \cdot N_z$ and $p = |\mathbf{predictors}|$ is the number of elements in the list). The matrix X is the usual input for a machine-learning algorithm; it has one line for each individual observation and one column for each variable (or predictor) observed. For the BLH retrieval,

250 the preparation also provides a vector Z (size N) containing the altitude of each individual observation. The algorithm (either K-means or EM, as specified by **algo**) is applied to the matrix X , with the parameters **n_clusters**, **init** and **n_inits**. This results in a vector of *labels* (size N) that contains the cluster attribution of each individual. Finally, we look for the first change in the cluster attribution, starting from the ground level. This gives us the value of BLH for this profile. These operations are repeated until reaching the end of the netCDF file.

255 3.4 Performance metrics

Two types of metrics were used.

- External scores: These metrics compare the result to a trustworthy reference. They have the advantage of providing a meaningful evaluation of the performance, but depend strongly on the quality of the reference (i.e., its accuracy and availability).
- 260 – Internal scores: These metrics rate how well the classification performs based only on the distances between points. They have the advantage of being always computable but are not linked to any physical property and therefore are not always meaningful.

None of these metrics are perfect; however, the information they provide allows a broader understanding of the algorithm performance.

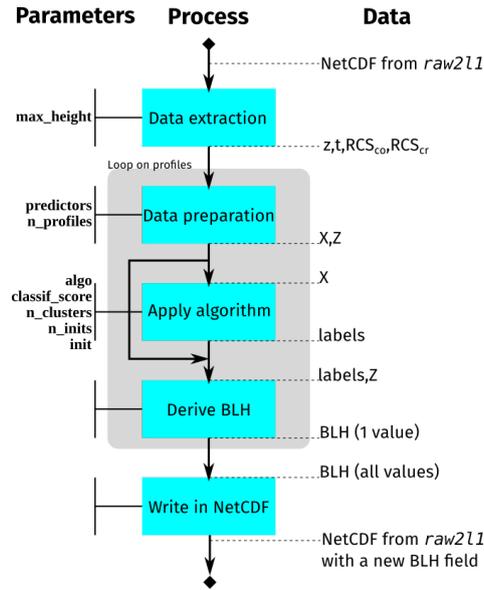


Figure 6. Simplified flowchart of the KABL algorithm.

265 3.4.1 External scores

External scores use a reference to assess the quality of the result. In our case, the reference is the RS-estimated BLH and, when available, the human expert-estimated BLH. If we denote \hat{Z} as the estimated BLH (by any of the previously introduced algorithms) and Z_{ref} as the reference, the external scores used in this study are the root mean square error (RMSE) (E_2 , equation 1) and the Pearson correlation (ρ , equation 2).

$$270 \quad E_2 = \sqrt{\mathbb{E}[(\hat{Z} - Z_{ref})^2]} \quad (1)$$

$$\rho = \frac{cov(\hat{Z}, Z_{ref})}{\sigma(\hat{Z})\sigma(Z_{ref})} \quad (2)$$

Here, \hat{Z} and Z_{ref} are random variables. When these scores are estimated, the random variables are replaced by a sample vector and the expectation and standard deviation are replaced by their usual estimators.

275 3.4.2 Internal scores

The quality of a classification can be quantified using scores that are based only on the labels and the distances between points. Such scores estimate how trustworthy an estimation is, without any external input. Many such scores exist with different formulation and different strengths and weaknesses (Desgraupes, 2013). In this study, three internal scores were used:

Table 1. Table of metrics used to measure the performance of KABL algorithm

Metric	Type	Description	Best/worst value
corr	External	Pearson correlation coefficient	1/0
errl2	External	Root mean square error (RMSE)	0/+∞
s_score	Internal	Silhouette score	1/−1
db_score	Internal	Davies-Bouldin index	0/+∞
ch_score	Internal	Calinski-Harabasz index	+∞/0
chrono	Other	Time to perform 24 h of BLH estimations	0/+∞
n_invalid	Other	Number of invalid values (NaN or Inf) in 24 h of BLH estimations	0/+∞

– the silhouette score (Rousseeuw, 1987).

280 which compares the average distance to its own group (a) to the average distance to the neighboring group (b): $S_{sil} = \frac{b-a}{\max(a,b)}$, where 1 is the best classification, 0 is neutral, and -1 is the worst classification;

– the Calinski-Harabasz index (Caliński and Harabasz, 1974),

285 which compares the between-cluster dispersion (B) to the within-cluster dispersion (W_k): $S_{ch} = \frac{(N-K)B}{(K-1)\sum_{k=1}^K W_k}$, where $+\infty$ is the best classification and 0 is the worst classification; and

– the Davies-Bouldin index (Davies and Bouldin, 1979).

which compares the average distance to its group center ($\bar{\delta}_k$) to the distance between the group centers ($d(\mu_k, \mu_{k'})$): $S_{db} = \max_{k' \neq k} \left(\frac{\bar{\delta}_k + \bar{\delta}_{k'}}{d(\mu_k, \mu_{k'})} \right)$, where 0 is the best classification and $+\infty$ is the worst classification.

290 These three scores were chosen to diversify the metrics and are all implemented in Scikit-learn (version ≥ 0.20).

3.4.3 Other metrics

In addition to the internal and external scores, the computation time and the number of invalid values (NaN or Inf) were recorded. BLH estimates of NaN or Inf can occur when all the points of the profile are assigned to the same cluster; this reflects a faulty configuration of the algorithm. Even though these metrics do not measure how well a program is performing, they are
295 useful to the user.

All the metrics used to measure the performance of KABL are summarized in Table 1.

4 Results

4.1 Sensitivity analysis of the KABL algorithm

A sensitivity analysis was performed on KABL code to identify the "best" configuration. Various KABL configuration were extensively tested on a single day: August 2, 2018, at the Trappes site, for which we have a hand-made reference (Figure 4 (top)). The most relevant configurations were retained and tested on the two-year dataset.

There are height parameters in the KABL code (see Sect. 3.3 for their descriptions). To assess the sensitivity of KABL to these parameters, the performance metrics (given in Sect. 3.4) were estimated with the hand-made BLH as Z_{ref} and with the output of KABL as \hat{Z} for different combinations of input parameters. The tested values for the input parameters are given in Table 2 and the output metrics are given in Table 1. We refer to a set of values for the KABL parameters as a *configuration*. Screening all the possible values listed in Table 2 would require 3240 different configurations.

To obtain an overview of these 3240 configurations, we started by estimating the influence of the parameters (listed in Table 2) on the different metrics (listed in Table 1). The influence of the parameters was quantified using first-order Sobol indices (Sobol, 2001; Iooss and Lemaître, 2015; Rieutord, 2017), that is, the ratio of the variance of the metric when the parameter was fixed over the total variance of the metric. If we denote Y as the metric and X as the vector of parameters, all considered as random variables, the first-order Sobol index of the i -th parameter is defined as $S_i = V(\mathbb{E}[Y|X_i]) / V(Y)$, where $V(\cdot)$ denotes the variance and $\mathbb{E}[\cdot]$ denotes the expectation. A higher Sobol index indicates a larger influence.

Figure 7 shows the Sobol indices obtained on the KABL computer code. Examining the matrix line by line, one can see that the metrics are sensitive to different parameters. For example, the silhouette score is very sensitive to **n_clusters** while the Calinski-Harabasz index is sensitive to **n_profiles** and **predictors**. Examining the matrix column by column, one can see that some parameters are more influential than others (e.g. **classif_score** is much less influential than **n_clusters**). This matrix highlights the main effects of changing a parameter and, therefore, how to set each parameter appropriately. For each parameter, we examined the metrics that it influences and determine the preferred configuration.

Critical parameters are indicated in Figure 7 by the darkest blue columns, namely, **n_clusters**, **algo**, **predictors** and **init**³. For each parameter, Figure 8 shows the distribution of the relevant output given the parameter value (violin plots are explained in Hintze and Nelson (1998)). For example, Figure 8a has the value of **algo** on the x -axis and the computing time on the y -axis. The 3240 different configurations were divided into two groups; those with **algo**='kmeans' and those with **algo**='gmm'. Figure 8a shows a smoothed histogram of the computing time for the divided populations. The other panels in Figure 8 were constructed in the same way. Each line corresponds to a critical parameter, and we represent the two most influenced outputs according to Figure 7.

The parameters values were chosen to give the most optimal values for the metrics they influence. The optimal values are indicated by a yellow star in each plot. To set **algo**, we examined the computing time (Figure 8a) and the Davies-Bouldin index (Figure 8b). These figures indicate that 'kmeans' is the best choice for both metrics (resulting in a lower computing time and a

³Even though **n_profiles** has a large Sobol index for the Calinski-Harabasz index, this influence was not explored because it is due to the increase of this index with the number of points.

Table 2. Possible values for the parameters of the KABL code. The dependencies between parameters result in 3240 different configurations.

Parameter	Possible values	Meaning
algo	'kmeans'	The K-means algorithm is used
	'gmm'	The EM algorithm is used (Gaussian mixture model)
classif_score	'silh'	The silhouette score is used
	'db'	The Davies-Bouldin index is used
	'ch'	The Calinski-Harabasz index is used
init	'random'	Starting points are chosen randomly
	'advanced'	Starting points are chosen with a smarter strategy
	'given'	Starting points are explicitly given
max_height	3500	The height (meters above ground level) at which profiles are cut
	4500	
n_clusters	2	The number of clusters to be formed is explicitly passed and is always the same
	3	
	4	
	5	
	'auto'	The number of clusters is automatically chosen to optimize classif_score
n_inits	10	The number of times the algorithm is repeated with different initializations (when init is not 'given')
	80	
n_profiles	1	Only the current profile is used
	2	The current profile and the previous profiles are used
	3	The current profile and the two previous profiles are used
	4	The current profile and the three previous profiles are used
predictors	'co'	The copolarized range-corrected signal is used at all times
	'co/co+cr'	The copolarized range-corrected signal is used during the daytime, and both polarization channels are used separately during nighttime
	'co+cr'	Both polarization channels are used separately at all times

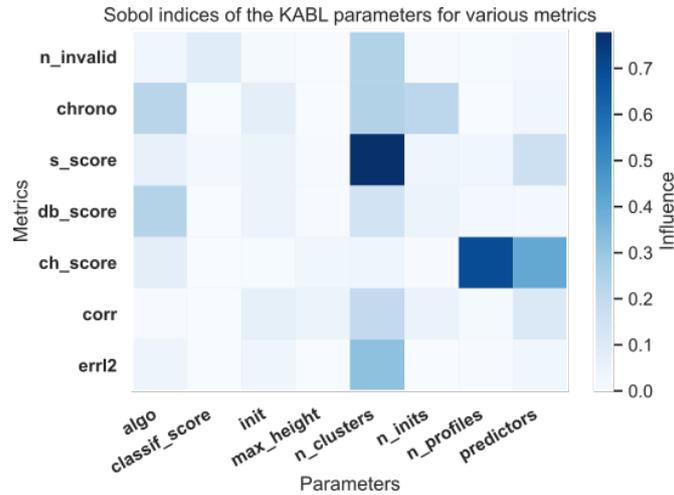


Figure 7. Relative influence of parameters on the different metrics. The x -axis indicates the parameters of the code, and the y -axis indicates the metrics. The shading represents the influence of the parameter on the metric with darker shading indicating larger influence.

Table 3. Retained values for the parameters of the KABL code after the sensitivity analysis

Parameter	Retained values
algo	'kmeans'
classif_score	'db'
init	'given'
max_height	4500
n_clusters	3
n_inits	10
n_profiles	1
predictors	'co'

lower Davies-Bouldin index). To set **init**, we examined the correlation (Figure 8-c) and the computing time (Figure 8d). In this case, 'given' appears to be the best choice. To set **n_clusters**, we examined the RMSE (8e) and the silhouette score (8f). They indicate the best number of cluster is respectively three and 'auto'. We chose to give priority to RMSE because silhouette score has also very high values for two clusters, which is suspicious given the presence of a cloud and a residual layer this day. To set **predictors**, we examined the silhouette score (8-g) and the Calinski-Harabasz index (8h); here, 'co' appears to be the best choice. Following this methodology, we can identify a few configurations worth trying. These configuration were tested on the two-year dataset. The configuration used to generate the results in Sect. 4.2.1 is given in Table 3. It was chosen to maximize correlation between KABL and RS at the Trappes site.

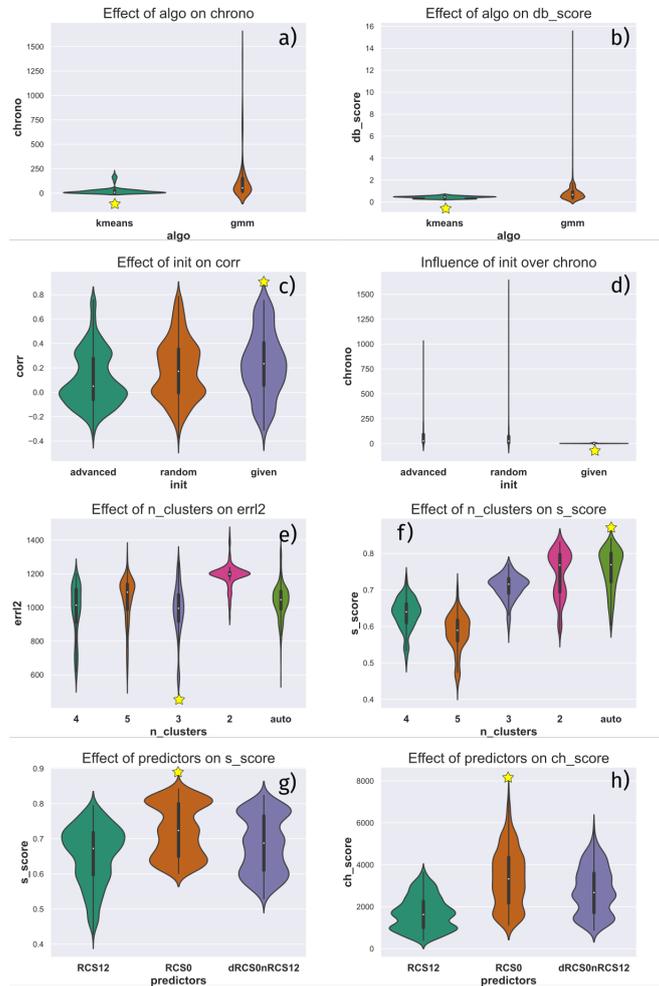


Figure 8. Distribution of the relevant outputs for the critical inputs. The effect of **algo** on (a) the computing time and (b) the Davies–Bouldin index. The effect of **init** on (c) the correlation and (d) the computing time. The effect of **n_clusters** on (e) the root mean square error (RMSE) and (f) the silhouette score. The effect of the predictors on (g) the silhouette score and (h) the Calinski–Harabasz index. For each panel, the best parameter value is highlighted by a yellow star.

4.2 Two-year comparison

The three methods (KABL, ADABL and the manufacturer’s algorithm) were compared to RS estimates over a two-year period.

4.2.1 Overall comparison

340 As explained in Sect. 3.4, two external scores, RMSE and the correlation, were used to assess the quality of the estimates. In equations 1 and 2, the reference BLH Z_{ref} was set to the RS estimate, as described in Sect. 2.2. To compute the scores, the BLH estimates from the lidar and radiosonde must be co-located. At the time of each RS estimation, the corresponding lidar estimate is the average of all the available estimates within the 10 min following the release of the radiosonde (this translates to one or two lidar estimates). The following meteorological conditions were discarded:

- 345
- rain (rain gauge measures the rainfall as >0 mm);
 - fog (scatterometer measures the visibility as <1000 m);
 - low level cloud (ceilometer measures cloud base height <3000 m);
 - RS estimation below 120m (blind zone for lidar); and
 - nighttime (RS launched at 23:15 UTC)

350 This selection rejects a large part of the dataset but ensures that only well-defined cases are retained for the comparison. In total, 178 RS measurements from Trappes and 101 RS measurements from Brest were used for the overall comparison. The meteorological conditions were measured using the ancillary instruments presented in Sect. 2.3. The results of the comparison are shown in Figure 9.

In Figure 9, we can see the results of the comparison between the KABL and RS estimates (blue bars), between the ADABL
355 and RS estimates (grey bars), and between the manufacturer’s and RS estimates (orange bar). The first column represents RMSE E_2 (lower is better), and the second column represents the correlation ρ (higher is better). The upper row shows the results for the Brest site, the lower row show the results for the Trappes site. One can see very different results depending on the site. While both KABL and ADABL outperform the manufacturer’s algorithm at the Trappes site, neither algorithm does at the Brest site. The correlation is strongly affected by the site. While the correlation for both KABL and ADABL is higher
360 than that for the manufacturer’s algorithm at the Trappes site, it collapses to close to zero for KABL at the Brest site (0.07 for ADABL). The RMSE values can be compared to the values given in Haeffelin et al. (2012). For KABL, we find 770 m at the Brest site and 798 m at the Trappes site, while for ADABL, we find 675 m at the Brest site and 552 m at the Trappes site. Our values are notably higher than those in Haeffelin et al. (2012). This is likely due to the larger extent of our dataset (178 RS at Trappes, 101 at Brest, spanning over 2 years) and the low maturity of the algorithms. Between ADABL and KABL, ADABL
365 has better correlation and RMSE values than KABL at both sites. The manufacturer’s algorithm performs well without any specific tuning on our part. It uses a wavelet covariance transform as described in Brooks (2003). This result is not surprising

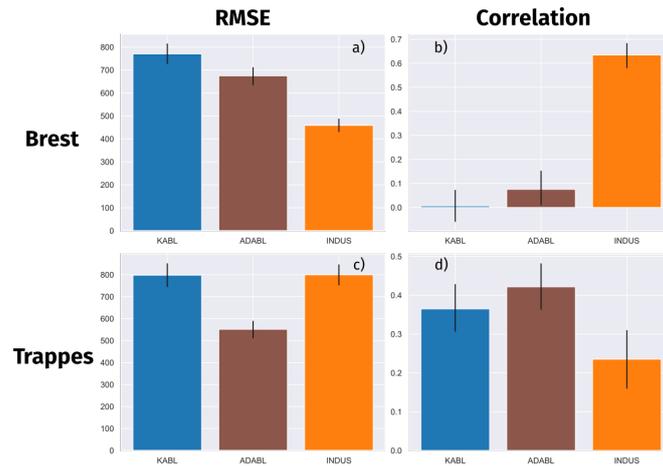


Figure 9. Results of a two-year comparison with the radiosonde (RS) estimates at both sites for two metrics: RMSE and correlation. Cases at night or with rain, fog, an RS-estimated BLH of under 120 m, or clouds under 3000 m were removed. The 95% confidence intervals were estimated using percentile bootstrapping (Davison and Hinkley, 1997)

because the wavelet method has been shown to be robust in numerous studies, especially in Caicedo et al. (2017), who included a cluster analysis method and concluded that the wavelet method should be preferred.

4.2.2 Seasonal and diurnal cycles

370 To quantify the ability of the algorithms to provide a consistent BLH estimation, Figure 10 shows the seasonal cycle (monthly average) and the diurnal cycle (six-minute average) at both sites. For each estimator, the thick line represents the average BLH estimate and the shaded area represents the inter-quartile gap. Rain, fog, and low-cloud conditions were discarded. For the monthly average, the night values were also removed. If we include only night estimates, the seasonal cycle is reversed for RS estimates, that is, the BLH estimates are lower in summer. For other estimators, we do not see such a difference between day
375 and night seasonal cycles (not shown).

At the Brest site (Figure 10a), estimates made by the manufacturer's algorithm are lower than those made by KABL and ADABL and estimates made by ADABL are usually higher than those made by KABL (except in July). The RS estimation gave BLH values that were low in summer (June–October), high in February and March (higher than the KABL estimates), and between the manufacturer's and KABL estimates during the rest of the year. Overall, the manufacturer's algorithm displays the seasonal cycle that is closest to that of the RS estimates, while KABL and ADABL both overestimate BLH. The inter-quartile distances (shaded areas) are large for all estimations, reflecting the wide range of the BLH estimates.
380

At the Trappes site (Figure 10c), KABL and ADABL also overestimate BLH in comparison to the RS estimates, while the manufacturer's estimate is close. The seasonal cycle is more visible in Trappes than in Brest, and all BLH estimates are higher in summer than in winter. The most pronounced cycle is given by KABL, while the least pronounced cycle is given by the

385 RS estimation. The inter-quartile distances are also very large, especially in summer, because the difference between the BLH estimates during the day and at night is larger.

Figures 10b and 10d show the diurnal cycle, where all values within the same six-minute period in the day were averaged. Because the radiosondes are only launched twice a day, at 11:15 UTC and 23:15 UTC, an equivalent RS-estimated diurnal cycle cannot be drawn. However, we used the average and quartile values at these times as checkpoints for the other estimates. 390 The manufacturer's and KABL estimates both have very smooth diurnal cycles, with lower BLH at night and maximum BLH around 15:00 UTC at the Trappes site and around 13:00 UTC at the Brest site. The KABL average is always higher than that calculated by the manufacturer's algorithm. The ADABL estimation has a very different diurnal cycle, similar to the conceptual image we have of the boundary layer. Indeed, ADABL was trained using hand-made BLHs that reflect this conceptual image. Therefore, it is not surprising that ADABL reproduces this image well; however, it may fail to adapt to special cases. It 395 appears that the "time" predictor (the number of seconds since midnight) has a large influence that is not balanced by the other predictors. This is likely because ADABL was trained on only two days, resulting in an unbalanced importance for sunrise and sunset on these particular days and at these locations. To balance this importance, the AdaBoost algorithm needs to be trained on more days and at more sites with a representative selection of cases.

4.3 Case study

400 The chosen case study was for April 19, 2017, at the Trappes site. The boundary layer was clearly visible and had nearly all the features of the conceptual image. The case study day must be different than the days used for the training of ADABL, otherwise the comparison would be biased in favor of ADABL.

Figure 11 represents the copolarized backscatter intensity ($RC S_{co}$) in shaded colors. The x -axis indicates the hour of the day (UTC), and the y -axis indicates the height (meters above ground level). The different BLH estimates are represented by dotted 405 lines: blue indicates KABL, orange indicates the manufacturer's algorithm, and green indicates ADABL. At the beginning of the day, there is a thick residual layer containing some plumes. Both KABL and the manufacturer's algorithm include these plumes in the boundary layer. Conversely, ADABL gives a very low estimate where there is no visible frontier. In the morning (from 08:00 UTC to 12:00 UTC), all the algorithms capture the transition reasonably well. However, KABL includes more irrelevant estimates (hitting what remains of the surface layer) than the other methods and ADABL gives an estimate that 410 is too high for no apparent reason around 12:00 UTC. During the day, ADABL sticks to the top of the boundary layer, the manufacturer's algorithm sticks to the surface layer (which is very visible), and KABL oscillates between the two. The evening transition is blurry; the surface layer slowly sends back increasing amounts of signal, finally turning the mixed layer into a residual layer. KABL locates this transition very early (around 17:00 UTC), when it stops oscillating and sticks to the surface layer. ADABL makes the transition more smoothly, from 19:00 UTC to 22:00 UTC. The manufacturer's algorithm is the last 415 to make the transition, at around 23:00 UTC, and the transition then occurs very sharply. We can conclude from this case study that none of the algorithms perfectly capture the boundary layer. Some of the limitations are physical, e.g., the evening transition is ill defined, resulting in disagreement between the algorithms. The RS estimate at 23:15 UTC is close to the lower boundary of the lidar range. This highlights the fact that BLHs below 120 m are not rare and cannot be detected with lidar

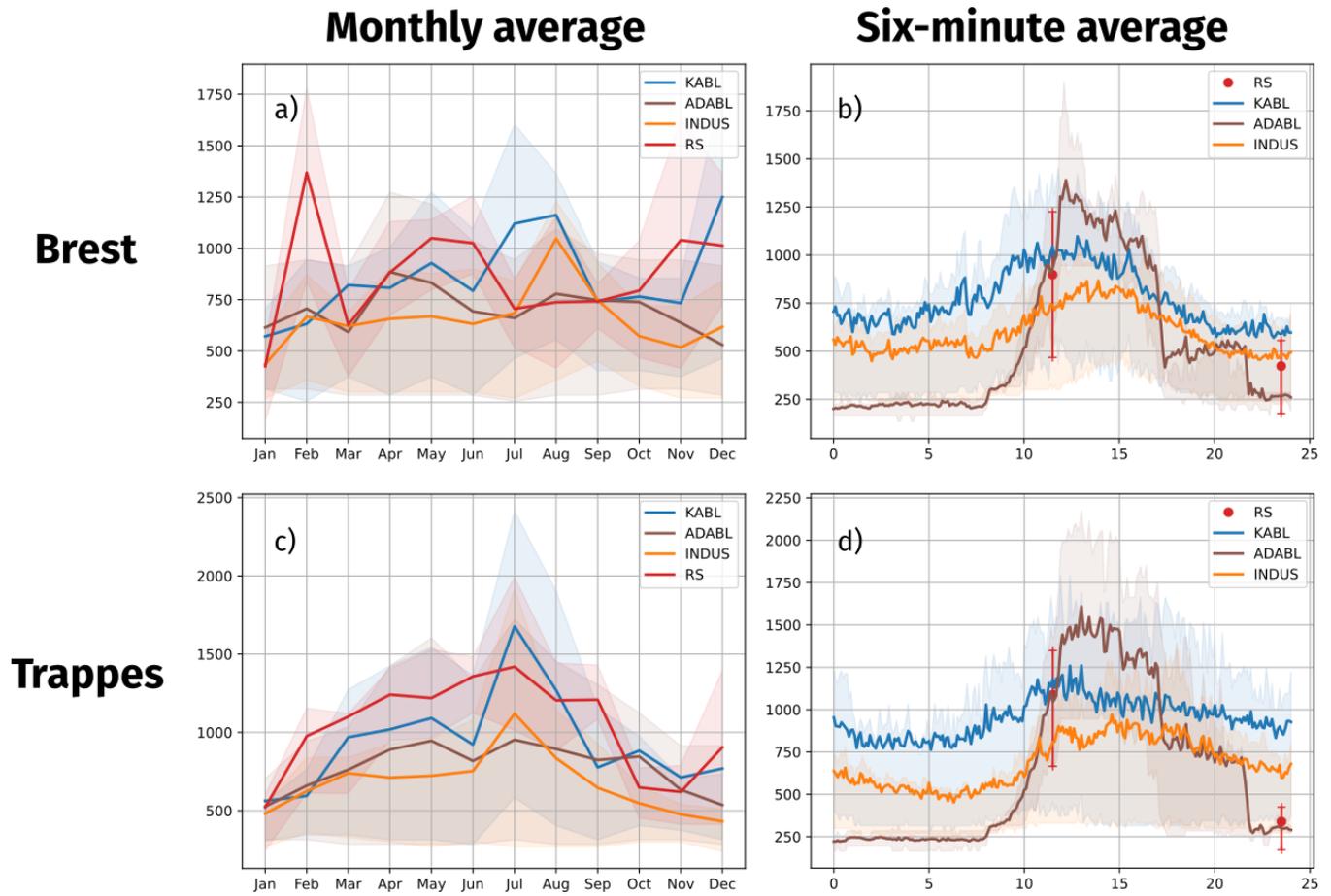


Figure 10. (a, c) Seasonal and (b, d) diurnal cycles of all BLH estimates at both sites. INDUS indicates the manufacturer’s algorithm. Thick lines represent the average, shaded area the quartiles.

alone, whatever the method. Some of the other limitations are algorithmic; KABL has an unfortunate tendency to oscillate
 420 between several candidates for the top of the boundary layer (surface layer or clouds), and ADABL too closely reproduces the
 features of the days it has been trained on (e.g., night estimates and morning transitions).

5 Discussion

This section discusses various aspects of the results and the methodology. For the sake of readability, it has been split into several short subsections.

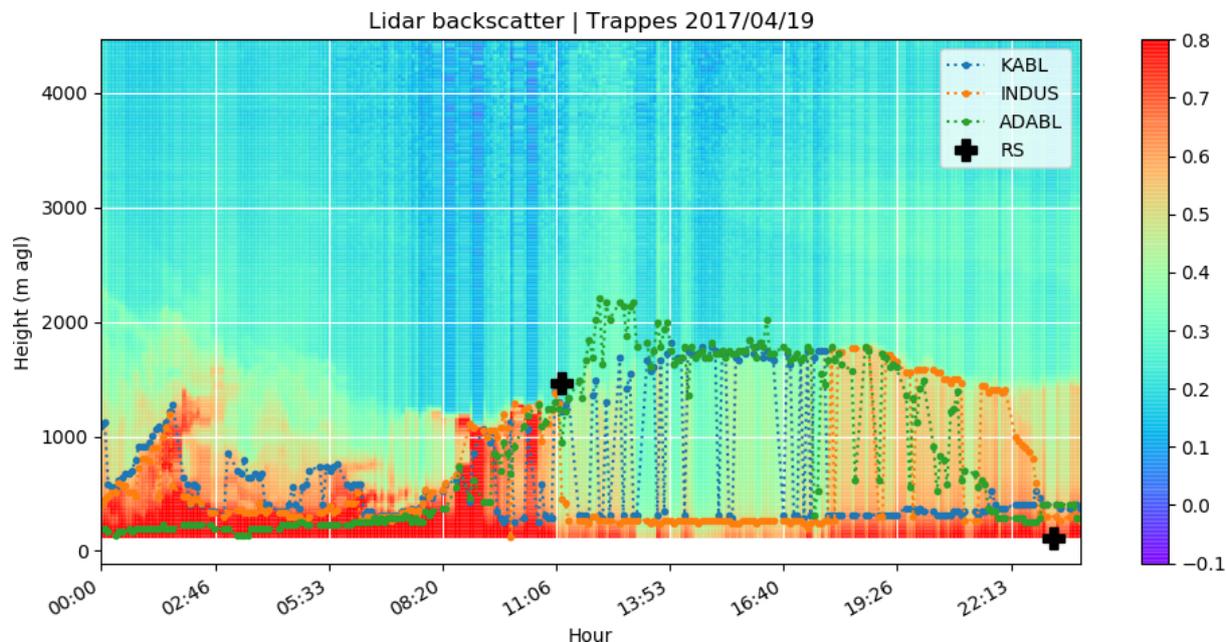


Figure 11. Case study: BLH estimates using different methods on April 19, 2017, at the Trappes site.

425 5.1 Algorithms maturity

Both of the examined algorithms are very recent. K-means algorithms have already been used to detect BLHs in previous studies (Toledo et al., 2014; Caicedo et al., 2017; Toledo et al., 2017; Rieutord et al., 2014); therefore, it is a more mature method. This is visible in this paper via the level of investigation which was much higher for KABL than for ADABL. Concerning boosting, this the first time, to our knowledge, that such an algorithm has been tested on this type of problem; therefore
 430 ADABL is a completely new algorithm. Yet, it outperforms KABL and competes favorably with the manufacturer’s algorithm despite raising training issues.

5.2 Time and altitude continuity

The oscillations observed in the figure 11 are unrealistic and need to be avoided. They occur with KABL because clusters do not always have vertical persistence (some points are identified as free atmosphere in the middle of the boundary layer).
 435 Vertical persistence needs to be enforced, for example, by adding altitude to the KABL predictors. Another way to filter out these oscillations is to increase the time continuity in the post-processing, for example, with a moving average, or by imposing a maximum BLH growth rate (Poltera et al., 2017). The distance used in K-means could also be modified to incorporate these constraints, for example, by adding penalty terms. In ADABL, time and altitude continuity are ensured because they are within the predictors. However, ADABL yields BLH estimates that are too similar to the BLHs in the training set. Removing time

440 and/or altitude from the predictors should be considered to force the algorithm to rely more on the measurements. Further, the sensitivity analysis presented here for KABL needs to be performed for ADABL.

5.3 Real time estimations

Even though it was not necessary for this study, all of the algorithms studied here can be used in real time. As soon as a backscatter profile is available, BLH estimations can be performed instantaneously⁴. However, KABL suffers from undesirable
445 oscillations from one profile to the next. A method to filter these oscillations is needed but would disable the "real-time" feature of the algorithm. In addition, the hour of the day needs to be explicitly passed in a periodic function. This has not been done here because we worked only on 24-hour time periods.

5.4 Quality of the evaluation

Even though we made an effort to sort the meteorological conditions using ancillary data, the two-year comparison still mixes
450 heterogeneous conditions. In addition, the results are clearly different at the sites studied here, emphasizing the importance of local conditions. A more precise casting of the meteorological conditions with atmospheric stability indices or large-scale insights would lead to a better understanding of the strengths and weaknesses of the algorithms. The importance of the site needs to be investigated by extending the study to a larger number of sites with different environments. A more careful examination of cloudy days also needs to be performed. Cases where the cloud bases are below to 3 km were filtered from our
455 study. However, cases where clouds reside inside the inversion should be detected by KABL as an extra cluster; further studies are required to confirm this behavior. In addition, ADABL was not specifically trained to deal with cloudy situations. Further studies to determine how ADABL behaves without training and how it could be appropriately trained would be interesting.

5.5 Quality of the reference

Radiosondes unquestionably provide the proper reference for altitude measurements. However, the derivation of BLH from
460 such measurements is contentious because several methods exist and some strongly disagree. Moreover, RS measurements cannot be used to assess the full diurnal BLH cycle. This is a clear limitation of this study because the RS measurements cannot determine if the difference between the diurnal cycle of ADABL and those of the other methods represents an improvement. Therefore, a very interesting project would be to use a dedicated field experiment with high-frequency radiosonde or other continuously running instruments as a reference. For example, microwave radiometers are good candidates because they
465 provide information that is not based on aerosols and BLH derivations from these instruments are routine (Cimini et al., 2013).

5.6 ADABL: Training

ADABL already shows good performance when trained on only two days. Most of its bad estimates result from the short length of its training period. Therefore, a short-term project would be to label more days with various meteorological conditions.

⁴Both KABL and ADABL need less than 1 s to run a single profile.

However, the dependence of ADABL on training makes it sensitive to instrumentation settings and calibrations. Even though
470 the effect of a calibration or the evolution of an instrumental device has not been studied, it is likely that training needs to
be repeated after each calibration or change in the instrumental device. Therefore, two strategies are possible for training
ADABL: remove the influence of calibration prior to training (this would require knowing the instrumental constants for all
of the devices) or train it to deal with differences (this would require including as many different devices as possible in the
training set, which would then become very large). In any case, the main limitation will be the need to label the entire dataset
475 (a priori by human experts).

5.7 KABL: Training-less

KABL appeared to perform the least well in this study; however, there are interesting prospects to improve its performance.
KABL does not require any training; therefore, it is less dependent on instrumentation settings and calibrations. Because it
is not strongly dependent on the instrumental devices, it can be used on backscatter profiles made by other instruments (e.g.,
480 ceilometers). Moreover, other profiles besides the backscatter intensity can be added as additional predictors for unsupervised
learning after normalization. Therefore, the concept of KABL can be advanced further to create synergy between multiple
remote sensing instruments. Microwave radiometers are good candidates because they have comparable time resolution to
lidar and provide independent information concerning the thermal stratification of the boundary layer. Cloud radars also have
comparable time resolution to lidar and provide additional independent information.

485 5.8 Quality flags

Currently, no quality flags for the estimation are provided. One approach would be to use the internal scores (i.e., silhouette,
Davies–Bouldin, and Calinski–Harabasz defined in Sect. 3.4) as quality flags; however, further study is required to determine
whether these metrics can serve as reliable quality flags.

6 Conclusions

490 This paper described two algorithms based on machine learning to estimate the mixing layer height from aerosol lidar mea-
surements. The first, KABL is based on the K-means algorithm. the second, ADABL, is based on the AdaBoost algorithm.
Both algorithms take the same input file: one day of data generated by the *raw211* routine and produce a similar output, a BLH
time series for the input day. KABL is a non-supervised algorithm that looks for a natural separation in the backscatter signals
between the boundary layer and the free atmosphere. ADABL is a supervised algorithm that fits a large number of decision
495 trees in a labelled dataset and aggregates them in an intelligent manner to provide a good prediction. KABL, ADABL and
the lidar manufacturer’s algorithm were tested on a two-year dataset taken from Météo-France operational lidar network. The
Trappes and Brest sites were chosen because of their different climates and the availability of regular RS measurements, which
were used as a reference.

A large discrepancy in the results was observed between the two sites. At the Trappes site, KABL and ADABL outperformed
500 the manufacturer's algorithm while the opposite occurred at the Brest site. At both sites, ADABL performed better than KABL
(higher correlation and lower error) and manufacturer's algorithm (using s wavelet covariance transform) performed well. By
analyzing the seasonal and diurnal cycles, we determined that the KABL and manufacturer's estimates have similar behavior;
however, the KABL estimates are always higher by approximately 200 m. ADABL generates the most pronounced diurnal
cycle, with a pattern that is very similar to the expected diurnal cycle; however, its results depend greatly on the days it has
505 been trained on. In particular, the sunset and sunrise times of these days over-influenced the ADABL estimation. In the case
study, we saw that both algorithms perform well overall; however, we identified several algorithmic limitations, e.g., KABL
tended to oscillate between several candidates for the top of the boundary layer (surface layers or clouds) and ADABL was
overly constrained by the days it was trained on (e.g., the night estimate and morning transition). In summary, ADABL is
promising but has training issues that need to be resolved, KABL has a lower performance but is much more versatile, and the
510 manufacturer's algorithm using a wavelet covariance transform performs well with little tuning but is not open source. A wide
range of future developments is available for ADABL and KABL, the most immediate being that the training set of ADABL
can be enhanced, time and altitude continuity can be enforced in the KABL estimation, and both can be compared to high
temporal resolution RS measurements.

Code availability. The KABL source code is available to and usable by all users, including commercial users. The code is freely available
515 under an open-source license at the following link: <https://github.com/ThomasRieutord/kabl>. It is made in Python 3.7 with regular statistics
and machine learning packages, namely, Scikit-learn 0.20 (Pedregosa et al., 2011) and SALib 1.3.7 (Herman and Usher, 2017), which are
open source and available under free licenses. The repository contains all the necessary features to run the code on *raw211* outputs. Several
days of data are also provided as examples.

Author contributions. Tiago Machado implemented an initial version of the KABL code and performed the first comparisons to the RS data.
520 Sylvain Aubert extracted and processed all the data (lidar, radiosonde, and ancillary), made one of the hand-labelled BLHs, and participated
actively in the writing of the manuscript. Thomas Rieutord implemented the current versions of KABL and ADABL, made one of the
hand-labelled BLHs, produced the figures, and actively participated in the writing of the manuscript.

Competing interests. The authors declare that they have no conflicts of interest.

Acknowledgements. The authors would like to thank Alexandre Paci, Alain Dabas, and Olivier Traullé for their helpful reading and com-
525 ments. We would like to thank Marc-Antoine Drouin, from the Site Instrumental de Recherche par Télédétection Atmosphérique, for pro-

viding us with the link to *raw211* and all the Météo-France personnel agents who install and maintain the lidar network. We thank Martha Evonuk from Evonuk Scientific Editing (<http://evonukscientificediting.com>) for editing a draft of this manuscript.

References

- Arciszewska, C. and McClatchey, J.: The importance of meteorological data for modelling air pollution using ADMS-Urban, *Meteorological Applications*, 8, 345–350, 2001.
- Arthur, D. and Vassilvitskii, S.: k-means++: The advantages of careful seeding, in: *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1027–1035, Society for Industrial and Applied Mathematics, 2007.
- Besse, P., Guillouet, B., and Laurent, B.: Wikistat 2.0: Educational Resources for Artificial Intelligence, arXiv preprint arXiv:1810.02688, 2018.
- Breiman, L., Friedman, J., Olshen, R., and Stone, C.: *Classification and Regression Trees*, Wadsworth, 1984.
- Brooks, I. M.: Finding boundary layer top: Application of a wavelet covariance transform to lidar backscatter profiles., *Journal of Atmospheric & Oceanic Technology*, 20, 2003.
- Caicedo, V., Rappenglück, B., Lefer, B., Morris, G., Toledo, D., and Delgado, R.: Comparison of aerosol lidar retrieval methods for boundary layer height detection using ceilometer aerosol backscatter data, *Atmospheric Measurement Techniques*, 10, 2017.
- Caliński, T. and Harabasz, J.: A dendrite method for cluster analysis, *Communications in Statistics-theory and Methods*, 3, 1–27, 1974.
- Campbell, J. R., Hlavka, D. L., Welton, E. J., Flynn, C. J., Turner, D. D., Spinhirne, J. D., Scott, III, V. S., and Hwang, I. H.: Full-Time, Eye-Safe Cloud and Aerosol Lidar Observation at Atmospheric Radiation Measurement Program Sites: Instruments and Data Processing, *Journal of Atmospheric and Oceanic Technology*, 19, 2002.
- Cimini, D., De Angelis, F., Dupont, J.-C., Pal, S., and Haeffelin, M.: Mixing layer height retrievals by multichannel microwave radiometer observations, 2013.
- Cohn, S. A. and Angevine, W. M.: Boundary layer height and entrainment zone thickness measured by lidars and wind-profiling radars, *Journal of Applied Meteorology*, 39, 1233–1247, 2000.
- Collaud Coen, M., Praz, C., Haeffele, A., Ruffieux, D., Kaufmann, P., and Calpini, B.: Determination and climatology of the planetary boundary layer height above the Swiss plateau by in situ and remote sensing measurements as well as by the COSMO-2 model, *Atmospheric Chemistry and Physics*, 14, 13 205–13 221, 2014.
- Davies, D. L. and Bouldin, D. W.: A cluster separation measure, *IEEE transactions on pattern analysis and machine intelligence*, pp. 224–227, 1979.
- Davison, A. C. and Hinkley, D. V.: *Bootstrap methods and their application*, 1, Cambridge university press, 1997.
- De Bruine, M., Apituley, A., Donovan, D., Klein Baltink, H., and de Haij, M.: Pathfinder: Applying graph theory for consistent tracking of daytime mixed layer height with backscatter lidar, *Atmospheric Measurement Techniques Discussions*, 10, 1893–1909, 2017.
- Desgraupes, B.: Clustering indices, *University of Paris Ouest-Lab Modal’X*, 1, 34, 2013.
- Dupont, J.-C., Haeffelin, M., Badosa, J., Elias, T., Favez, O., Petit, J., Meleux, F., Sciare, J., Crenn, V., and Bonne, J.: Role of the boundary layer dynamics effects on an extreme air pollution event in Paris, *Atmospheric environment*, 141, 571–579, 2016.
- Flynn, C. J., Mendoza, A., Zheng, Y., and Mathur, S.: Novel polarization-sensitive micropulse lidar measurement technique, *Opt. Express*, 15, 2785–2790, 2007.
- Freund, Y. and Schapire, R. E.: A decision-theoretic generalization of on-line learning and an application to boosting, *Journal of computer and system sciences*, 55, 119–139, 1997.
- Gamage, N. and Hagelberg, C.: Detection and analysis of microfronts and associated coherent events using localized transforms, *Journal of the atmospheric sciences*, 50, 750–756, 1993.

- 565 Guo, J., Miao, Y., Zhang, Y., Liu, H., Li, Z., Zhang, W., He, J., Lou, M., Yan, Y., Bian, L., et al.: The climatology of planetary boundary layer height in China derived from radiosonde and reanalysis data, *Atmospheric Chemistry and Physics*, 16, 13 309, 2016.
- Haefele, A., Hervo, M., Turp, M., Lampin J-L, Haeffelin, M., Lehmann, V., et al.: The E-PROFILE network for the operational measurement of wind and aerosol profiles over Europe, *Proceedings of WMO Technical Conference on Meteorological and Environmental Instruments and Methods of Observation (CIMO TECO 2016, Madrid)*, 2016.
- 570 Haeffelin, M., Angelini, F., Morille, Y., Martucci, G., Frey, S., Gobbi, G., Lolli, S., O'dowd, C., Sauvage, L., Xueref-Rémy, I., et al.: Evaluation of mixing-height retrievals from automatic profiling lidars and ceilometers in view of future integrated networks in Europe, *Boundary-Layer Meteorology*, 143, 49–75, 2012.
- Hastie, T., Tibshirani, R., and Friedman, J.: *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer Science & Business Media, 2009.
- 575 Hayden, K., Anlauf, K., Hoff, R., Strapp, J., Bottenheim, J., Wiebe, H., Froude, F., Martin, J., Steyn, D., and McKendry, I.: The vertical chemical and meteorological structure of the boundary layer in the Lower Fraser Valley during Pacific'93, *Atmospheric Environment*, 31, 2089–2105, 1997.
- Hennemuth, B. and Lammert, A.: Determination of the atmospheric boundary layer height from radiosonde and lidar backscatter, *Boundary-Layer Meteorology*, 120, 181–200, 2006.
- 580 Herman, J. and Usher, W.: SALib: An open-source Python library for Sensitivity Analysis, *The Journal of Open Source Software*, 2, <https://doi.org/10.21105/joss.00097>, <https://doi.org/10.21105/joss.00097>, 2017.
- Hintze, J. L. and Nelson, R. D.: Violin Plots: A Box Plot-Density Trace Synergism, *The American Statistician*, 52, 181–184, 1998.
- Iooss, B. and Lemaître, P.: A review on global sensitivity analysis methods, in: *Uncertainty management in simulation-optimization of complex systems*, pp. 101–122, Springer, 2015.
- 585 Jain, A. K., Murty, M. N., and Flynn, P. J.: Data clustering: a review, *ACM computing surveys (CSUR)*, 31, 264–323, 1999.
- Kotthaus, S. and Grimmond, C. S. B.: Atmospheric boundary-layer characteristics from ceilometer measurements. Part 1: A new method to track mixed layer height and classify clouds, *Quarterly Journal of the Royal Meteorological Society*, 144, 1525–1538, 2018.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E.: Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- 590 LeCun, Y., Bengio, Y., and Hinton, G.: Deep learning, *nature*, 521, 436–444, 2015.
- Melfi, S., Spinhirne, J., Chou, S., and Palm, S.: Lidar observations of vertically organized convection in the planetary boundary layer over the ocean, *Journal of climate and applied meteorology*, 24, 806–821, 1985.
- Menut, L., Flamant, C., Pelon, J., and Flamant, P. H.: Urban boundary-layer height determination from lidar measurements over the Paris area, *Applied Optics*, 38, 945–954, 1999.
- 595 Mohan, M., Bhati, S., Sreenivas, A., and Marrapu, P.: Performance evaluation of AERMOD and ADMS-urban for total suspended particulate matter concentrations in megacity Delhi, *Aerosol and Air Quality Research*, 11, 883–894, 2011.
- Morille, Y., Haeffelin, M., Drobinski, P., and Pelon, J.: STRAT: An automated algorithm to retrieve the vertical structure of the atmosphere from single-channel lidar data, *Journal of Atmospheric and Oceanic Technology*, 24, 761–775, 2007.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V.,
600 Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E.: Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research*, 12, 2825–2830, 2011.
- Pollard, D. et al.: Strong consistency of k -means clustering, *The Annals of Statistics*, 9, 135–140, 1981.

- Poltera, Y., Martucci, G., Collaud Coen, M., Hervo, M., Emmenegger, L., Henne, S., Brunner, D., and Haeefele, A.: PathfinderTURB: an automatic boundary layer algorithm. Development, validation and application to study the impact on in-situ measurements at the Jungfraujoch, Atmospheric Chemistry and Physics Discussions, 2017.
- 605 Rieutord, T.: Sensitivity analysis of a filtering algorithm for wind lidar measurements, Ph.D. thesis, 2017.
- Rieutord, T., Brewer, W. A., and Hardesty, R. M.: Automatic detection of boundary layer height using Doppler lidar measurements, 2014.
- Rousseeuw, P. J.: Silhouettes: a graphical aid to the interpretation and validation of cluster analysis, *Journal of computational and applied mathematics*, 20, 53–65, 1987.
- 610 Schapire, R. E.: Explaining adaboost, in: *Empirical inference*, pp. 37–52, Springer, 2013.
- Seibert, P., Beyrich, F., Gryning, S.-E., Joffre, S., Rasmussen, A., and Tercier, P.: Review and intercomparison of operational methods for the determination of the mixing height, *Atmospheric environment*, 34, 1001–1027, 2000.
- Seidel, D. J., Ao, C. O., and Li, K.: Estimating climatological planetary boundary layer heights from radiosonde observations: Comparison of methods and uncertainty analysis, *Journal of Geophysical Research: Atmospheres*, 115, 2010.
- 615 Seidel, D. J., Zhang, Y., Beljaars, A., Golaz, J.-C., Jacobson, A. R., and Medeiros, B.: Climatology of the planetary boundary layer over the continental United States and Europe, *Journal of Geophysical Research: Atmospheres*, 117, 2012.
- Seity, Y., Brousseau, P., Malardel, S., Hello, G., Bénard, P., Bouttier, F., Lac, C., and Masson, V.: The AROME-France convective-scale operational model, *Monthly Weather Review*, 139, 976–991, 2011.
- Selim, S. Z. and Ismail, M. A.: K-means-type algorithms: a generalized convergence theorem and characterization of local optimality, *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, pp. 81–87, 1984.
- 620 Senff, C., Bösenberg, J., Peters, G., and Schaberl, T.: Remote sensing of turbulent ozone fluxes and the ozone budget in the convective boundary layer with DIAL and Radar-RASS: A case study, *Contributions to atmospheric physics*, 69, 161–176, 1996.
- Sobol, I. M.: Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates, *Mathematics and computers in simulation*, 55, 271–280, 2001.
- 625 Stull, R. B.: *An introduction to boundary layer meteorology*, vol. 13, Springer, 1988.
- Toledo, D., Córdoba-Jabonero, C., and Gil-Ojeda, M.: Cluster analysis: A new approach applied to lidar measurements for atmospheric boundary layer height estimation, *Journal of Atmospheric and Oceanic Technology*, 31, 422–436, 2014.
- Toledo, D., Córdoba-Jabonero, C., Adame, J. A., De La Morena, B., and Gil-Ojeda, M.: Estimation of the atmospheric boundary layer height during different atmospheric conditions: a comparison on reliability of several methods applied to lidar measurements, *International journal of remote sensing*, 38, 3203–3218, 2017.
- 630 Ware, J., Kort, E. A., DeCola, P., and Duren, R.: Aerosol lidar observations of atmospheric mixing in Los Angeles: Climatology and implications for greenhouse gas observations, *Journal of Geophysical Research: Atmospheres*, 121, 9862–9878, 2016.