1 **Probabilistic topic modelling in food spoilage analysis: a case study with Atlantic salmon**

2 **(*Salmo salar*)**

3 L. Kuuliala[1,2*], R. Pérez-Fernández[2,3], M. Tang[2], M. Vanderroost[1], B. De Baets[2] and F.

4 Devlieghere[1]

5 [1]*Research Unit Food Microbiology and Food Preservation (FMFP), Department of Food*

6 *Technology, Safety and Health, Part of Food2Know, Faculty of Bioscience Engineering, Ghent*

7 *University, Coupure links 653, B-9000 Ghent, Belgium*

8 [2]*Research Unit Knowledge-based Systems (KERMIT), Department of Data Analysis and*

9 *Mathematical Modelling, Part of Food2Know, Faculty of Bioscience Engineering, Ghent*

10 *University, Coupure links 653, B-9000 Ghent, Belgium*

11 [3]*Department of Statistics and O.R. and Mathematics Didactics, University of Oviedo, Calle*

12 *Federico García Lorca 18, 33007 Oviedo, Spain*

13

14 [*]Corresponding author. Research Unit Food Microbiology and Food Preservation (FMFP),

15 Department of Food Technology, Safety and Health, Faculty of Bioscience Engineering, Ghent

16 University, Coupure links 653, B-9000 Ghent, Belgium. Tel.: +32-(0)9 264 92 03; *E-mail*

17 *address:* Lotta.Kuuliala@UGent.be

18    **Abstract**

19    Probabilistic topic modelling is frequently used in machine learning and statistical analysis for

20    extracting latent information from complex datasets. Despite being closely associated with

21    natural language processing and text mining, these methods possess several properties that make

22    them particularly attractive in metabolomics applications where the applicability of traditional

23    multivariate statistics tends to be limited. The aim of the study was thus to introduce probabilistic

24    topic modelling – more specifically, Latent Dirichlet Allocation (LDA) – in a novel experimental

25    context: volatilome-based (sea)food spoilage characterization. This was realized as a case study,

26    focusing on modelling the spoilage of Atlantic salmon (*Salmo salar*) at 4 °C under different

27    gaseous atmospheres (% $CO_2/O_2/N_2$): 0/0/100 (A), air (B), 60/0/40 (C) or 60/40/0 (D). First, an

28    exploratory analysis was performed to optimize the model tunings and to consequently model

29    salmon spoilage under 100 % $N_2$ (A). Based on the obtained results, a systematic spoilage

30    characterization protocol was established and used for identifying potential volatile spoilage

31    indicators under all tested storage conditions. In conclusion, LDA could be used for extracting

32    sets of underlying VOC profiles and identifying those signifying salmon spoilage, giving rise to

33    an extensive discussion regarding the key points associated with model tuning and/or spoilage

34    analysis. The identified compounds were well in accordance with a previously established

35    approach based on partial least squares regression analysis (PLS). Overall, the outcomes of the

36    study not only reflect the promising potential of LDA in spoilage characterization, but also

37    provide several new insights into the development of data-driven methods for food quality

38    analysis.

39    **Keywords**

40  Latent Dirichlet Allocation; food quality; metabolomics; potential spoilage indicator; volatile

41  organic compound

## 1. Introduction

The rapid development of metabolomics technologies has greatly improved our understanding of complex biological systems during the past few decades. In the food and nutrition sector, this global trend has had a major impact on the development of foodomics (Miguel et al., 2012), a novel interdisciplinary field where metabolomics – the study of low molecular weight (<1500 Da) molecules associated with biological samples (Castro-Puyana et al., 2017; Pinu, 2016) – has already been used for addressing various questions related to food quality and safety (Böhme et al., 2019; Klampfl, 2018; Mancano et al., 2018; Martinović et al., 2018; Xu, 2017). In particular, the latest advances in the analysis of spoilage-indicating volatile organic compounds (VOCs) have greatly benefitted both scientific knowledge (Dong et al., 2019; Odeyemi et al., 2018; Wang et al., 2016) and technology development (Ghasemi-Varnamkhasti et al., 2018; Pavase et al., 2018; Poghossian et al., 2019).

However, the complexity of the microbial metabolism poses a major challenge in food quality characterization. At any given moment during storage time, the food volatilome consists of numerous compounds that differ in terms of quantity, chemical composition, reactivity, olfactory impact and sensory acceptability. Irrespective of the applied quantification method, the extraction of information from the resulting datasets thus calls for advanced statistical analysis. Basic multivariate methods such as principal components analysis (PCA), partial least squares regression analysis (PLS) and hierarchical cluster analysis (HCA) have frequently been used for this purpose (Bermejo-Prada et al., 2015; Mansur et al., 2019; Mikš-Krajnik et al., 2016). However, the applicability of these methods in biomarker identification tends to be limited; for example, while PLS outperforms HCA and PCA as a selective tool, it still requires a linear relationship between the studied variables and has a limited capacity in distinguishing correlation from a cause-and-effect

65  relationship (Kuuliala et al., 2018). Hence, more flexible methods are needed for improving our

66  ability to identify the most useful volatile spoilage indicators.

67  Probabilistic topic modelling comprises a group of methods used in machine learning and statistical

68  analysis for extracting underlying thematic information from an unstructured collection of

69  documents (Blei, 2012). Currently, Latent Dirichlet Allocation (LDA) introduced by Blei et al.

70  (2003) represents one of the most widespread approaches. Despite the fact that these methods have

71  traditionally been closely associated with the analysis of textual data – for example,

72  consumer/customer feedback (Bastani et al., 2019; Hu et al., 2019), social media content (Curiskis

73  et al., 2019; Nolasco & Oliveira, 2019) or research interests (Xiong et al., 2019; Yang et al., 2019) –

74  LDA has also already been successfully used in different biological settings, particularly in

75  genomics (Chen et al., 2010; Perina et al., 2010; Pratanwanich & Lio, 2014; Shiraishi et al., 2015;

76  Yu et al., 2014; Zhang et al., 2012). However, to the best of the authors' knowledge, its prospects

77  within food science still remain to be elucidated.

78  The aim of the present study is to introduce LDA as an exploratory and selective statistical

79  technique for characterizing (sea)food quality on the basis of its volatilome. First, a non-technical

80  overview of the principles of LDA is given in Section 2. In the experimental part (Sections 3-4),

81  LDA is applied for modeling the quality decay of raw Atlantic salmon (*Salmo salar*) under different

82  gaseous atmospheres and for consequently identifying potential spoilage indicators. Special

83  emphasis is given on 1) optimizing the model parameters, 2) developing a systematic spoilage

84  characterization protocol, including a set of criteria for identifying VOCs that possess promising

85  potential for quality monitoring applications (referred to as "potential spoilage indicators" from here

86  on), and 3) comparing the performance of the said protocol with a previously established PLS-based

87  approach (Kuuliala et al., 2019).  The obtained results, conclusions and decisions are discussed in

88  Section 5. Finally, summarizing remarks are given in Section 6.

## 2. Latent Dirichlet Allocation

LDA is a flexible *generative probabilistic model* for discrete data (Blei et al., 2003), meaning that it can be used for determining the joint probability distribution underlying a group of known samples and consecutively generating new samples from the same distribution. This approach not only allows for exploring previously unknown (underlying or *latent*) structures in large and complex datasets, but also reducing dimensionality, detecting co-occurring variables and evaluating the similarity between individual unlabeled samples. It should be noted though that since LDA is inherently *unsupervised*, it does not involve a pre-defined output and thus cannot be directly used for analyzing the relations between independent and dependent variables. For that purpose, a complementary modelling approach and/or an extension into (semi-)supervised LDA is needed (see e.g. Fu et al., 2015; Li et al., 2018).

The key concepts of topic modelling are *word* (*term*)*, document, corpus*, *word-document matrix (WDM)* and *topic.* By definition, a word refers to a basic unit of discrete data, a document to a sequence of words, and a corpus to a collection of documents (Blei et al., 2003). A WDM (also known as a *document-term matrix* or a *bag of words*) indicates the frequency of each word in each document belonging to the corpus; in case of $n$ documents and $m$ words, an $n \times m$ WDM is obtained (Liu et al., 2016). This information can be used for extracting a set of probability distributions over all words which appear in the corpus (i.e. the *vocabulary*); these distributions – or, more specifically, the interpretations of these distributions (see Subsection 3.2.3) – can be referred to as topics. In other words, a given document is seen as a product of a generative process, consisting of 1) choosing a distribution over topics, 2) a topic from the chosen distribution, 3) a word from the chosen topic, and 4) returning to step 2 until the pre-determined document length has been reached (Blei et al., 2003).

The application of LDA in a classical text mining setting is illustrated in Example 1.

*Example 1.* *A corpus was formed from 121 abstracts (documents) of original research conducted at the Research Unit Food Microbiology and Food Preservation at Ghent University (FMFP) and published in international peer-reviewed journals between 2000 and 2018. A WDM was constructed by calculating the frequency of each individual word in each abstract (excluding digits, punctuation, symbols and English stop-words) and used for generating an LDA model with five topics. The R package **tm** (Feinerer & Hornik, 2018; Feinerer et al., 2008) was used for constructing the WDM and the package **topicmodels** (Grün & Hornik, 2011) for learning the LDA model with default parameters. The obtained results were examined and visualized in accordance with Subsection 3.2.3.*

*Figure 1 presents the 20 most common words and their distribution in the extracted topics 1-5, associatively interpreted as follows: 1) microbial aspects of food preservation technologies, 2) food quality and microbial spoilage, 3) modelling of microbial behavior in foods, with special emphasis on norovirus, 4) food packaging and shelf life, 5) microbial food safety and health. Overall, these topics could be interpreted as the central themes of research carried out at FMFP during the studied time interval. It should be noted that the number of topics affects their specificity; for example, a model with two topics could be interpreted as 1) food quality and 2) food safety (data not shown).*

*Table 1 shows the distribution of topics 1 - 5 in selected documents. For example, document 16 (“Growth of* Escherichia coli *O157:H7 and* Listeria monocytogenes *with prior resistance to intense pulsed light and lactic acid” by Rajkovic et al., 2011) could be associated with food preservation (topic 1), whereas the 64:36 relation between topics 1 and 3 in document 40 (“Multi-method approach indicates no presence of sub-lethally injured* Listeria monocytogenes *cells after mild heat treatment” by Uyttendaele et al., 2008) indicates that pathogen growth was examined by both experimental and statistical methods.*

When compared with classical multivariate methods, LDA poses several advantages that make it particularly attractive for addressing complex biological problems. It is flexible, adaptable and imposes relatively few assumptions; importantly, the number of topics ($k$) is assumed to be known, the order of words within a given document to be exchangeable and all documents to be independent of each other (Liu et al., 2016). A given document may be associated with multiple topics and a given word may belong to multiple documents (Binkley et al., 2014; Griffiths & Steyvers, 2004). Overall, LDA is compatible with any inference algorithm and can be extended or incorporated as part of a more complex approach (Blei et al., 2003). A further discussion on model tuning and associated challenges is provided in Subsection 5.1.

## 3. Materials and methods

### 3.1. Data pre-processing and notation

All statistical analyses were performed using R 3.6.1 (R Core Team, 2019) on four independent subsets (A-D) of the salmon data previously collected by Kuuliala et al. (2019). Briefly, the experimental units of the study were individually packed salmon fillet portions, stored for $\leq$ 13 days at 4 °C under specific gaseous atmospheres (% $CO_2/O_2/N_2$): 0/0/100 (A), air (B), 60/0/40 (C) or 60/40/0 (D). Each package (n=16 per atmosphere) was linked with a single value of each of the following variables: storage time (d), concentrations of 25 VOCs (C1-C25; ppb) and sensory rejection percentage ($R_\%$; %). These subsets were used for generating four WDMs, where the individual salmon packages were treated as "documents", VOCs as "words" and their concentrations rounded to the nearest whole numbers as "frequencies". Following the same semantic principles, the extracted "topics" are referred to as "VOC profiles" or "profiles" throughout the present manuscript and denoted whenever applicable as X$k$P$n$, where X indicates the subset (A-D), $k$ the number of profiles, and $n$ is a profile identifier ($n$=1,2,…,$k$)..

### 3.2. Exploratory analysis

8

The exploratory analysis aimed at optimizing model performance (Subsections 3.2.1-3.2.2) and, consecutively, establishing the principles of spoilage characterization (Subsections 3.2.3-3.2.4). All activities were performed using the subset A (100 % $N_2$) and its WDM. In case of all mentioned R functions, default parameters were used unless otherwise specified.

*3.2.1. Model tuning*

In literature, several metrics have been proposed for facilitating the selection of the number of topics (here, profiles). However, the existing methods typically apply different identification criteria; for example, the algorithm of Cao et al. (2009) returns the value of $k$ that minimizes the average cosine distance between the extracted topics, the algorithm of Arun et al. (2010) the value that minimizes the symmetric Kullback-Leibler (KL) divergence (Kullback & Leibler, 1951) between the matrices representing the word-per-topic and topic-per-document distributions, and the algorithm of Deveaud et al. (2014) the value that maximizes the sum of the divergences between topic pairs. For this reason, multiple metrics are frequently implemented for comparative purposes.

In this study, the three aforementioned metrics (denoted Cao, Arun and Deveaud) were used for selecting appropriate tunings for LDA models with 2-10 profiles, using the function FindTopicsNumber() from the package **ldatuning** (Murzintcev, 2019). Inference was estimated by the variational expectation-maximization algorithm (VEM) or Gibbs sampling with following specifications:

- VEM: *method*="VEM"
- Gibbs100A: *method*="Gibbs", *iter*=100, *burnin*=0, *thin*=1
- Gibbs100B: *method*="Gibbs", *iter*=100, *burnin*=50, *thin*=1
- Gibbs1000A: *method*="Gibbs", *iter*=1000, *burnin*=0, *thin*=1
- Gibbs1000B: *method*="Gibbs", *iter*=1000, *burnin*=500, *thin*=1

184   where *iter* refers to the number of iterations, *burnin* to the number of discarded (burned) initial

185   iterations and *thin* to the interval of retained iterations. For evaluating modelling consistency, two

186   lists consisting of ten random seeds (s1-s10) or WDM column orders (O1-O10) were created; all

187   possible seed-order combinations were tested using each of the aforementioned five methods.

188   Finally, the obtained results (optimal number of profiles; $k_{opt}$) were visualized using the function

189   heatmap.2() from the package **gplots** (Warnes et al., 2020). In conclusion (see Subsection 5.1 for a

190   discussion), the method Gibbs1000B and a single seed-order pair were selected to be used for all

191   further modelling activities.


192   *3.2.2. Cross-validation*

193   In literature, the performance of LDA is frequently assessed by means of a perplexity analysis.

194   Briefly, perplexity ($P(w)$) is a measure of the model's ability to predict an unseen dataset ($w$) and

195   can be denoted as follows:


$$P(w) = \exp\left(-\frac{\log\left(p(w)\right)}{\sum_{d=1}^{D}\sum_{j=1}^{V} n^{jd}}\right) \tag{1}$$


197   where $\log\left(p(w)\right)$ is the log-likelihood of the new data and $n^{jd}$ the number of times the $j^{\text{th}}$ term occurs

198   in the $d^{\text{th}}$ document (Grün & Hornik, 2011); an increasing perplexity signifies a decreasing model

199   performance.


200   In this study, perplexity analysis was used for specifying the selection of $k$ and for evaluating the

201   prediction ability. First, leave-one-out cross-validation was performed with functions provided in

202   the package **topicmodels** (Grün & Hornik, 2011). Each sample 1-16 was assigned once as the

203   holdout set while using the other fifteen samples as the training set; models ($k$=2,…,10) were

204   learned using the function LDA() and the selected tunings (Subsection 3.2.1). Perplexities were

205   calculated for each model with the function perplexity(), using the holdout sets as the new data.


206   *3.2.3. Profile interpretation*

207  In literature, different methods are used for interpreting the extracted LDA profiles. In classical text

208  mining applications, this task is often carried out by examining the most frequent words of each

209  extracted topic (see, e.g., Example 1). Although it may seem an intuitive process, it relies on

210  knowledge regarding the logical relations between different words: semantically speaking, the

211  concept of "topic" should thus be seen as the "label" of a word distribution rather than being

212  synonymous with the distribution itself.

213  In this study, profile interpretation was performed for identifying spoilage-associated profiles. First,

214  three exploratory models ($k$=3, 5, 9) were learned using the function LDA() from the package

215  **topicmodels** and the selected tunings (Subsection 3.2.1). Relevant distributions were extracted

216  using the function tidy() from the package **tidytext** (Silge & Robinson, 2016): the distribution of

217  the VOCs within profiles was visualized using the package **ggplot2** (Wickham, 2016) and the

218  distribution of the profiles within samples using Excel 2013 for Windows. In order to reduce the

219  risk of interpretation bias, top-8 VOCs were reported in accordance with the suggestion of Agrawal

220  et al. (2018). For assessing the similarity between the compositions of all extracted profiles,

221  hierarchical cluster analysis (HCA) was performed on the basis of the Euclidean distance and

222  average linkage, using the function pvclust() from the package **pvclust** (Suzuki et al., 2019).

223  Finally, the profiles were interpreted by examining the relations between cluster distribution,

224  storage time and $R_\%$ (Excel 2013 for Windows).

225  *3.2.4. Identification of potential spoilage indicators*

226  In this study, potential spoilage indicators were identified by evaluating the prevalence of each

227  VOC in the spoilage-associated profiles of the three exploratory models (Subsection 3.2.3). This

228  was done by establishing the relative spoilage association ($SA_\%$):

229  $$SA_\% = \frac{n_{sc}}{n_{tot}} \qquad\qquad (2)$$

where $n_{sc}$ is the number of times a given compound occurred among the top-8 VOCs of the
spoilage-associated profile(s) (n=$n_{tot}$) of a given model (k=3, 5, 9). Compounds fulfilling the
criterion $SA_{\%} \geq 0.5$ (occurrence among the top-8 VOCs in at least half of the spoilage-associated
profiles) were considered relevant. Finally, identification performance was assessed by comparing
the obtained results with the outcomes of a previously established PLS-based protocol (Kuuliala et
al., 2019; denoted here as $IP_{PLS}$).

**3.3. Selective analysis**

The selective analysis aimed at the identification of potential spoilage indicators under all tested
storage conditions. This was done by applying the spoilage characterization protocol developed
during exploratory analysis (Subsection 4.1.5) also for subsets B-D. Briefly, independent LDA
models (k=3) representing a given condition (B-D) were learned and interpreted in accordance with
Subsection 3.2.3 and potential spoilage indicators were identified in accordance with Subsection
3.2.4.

**4. Results**

*4.1. Exploratory analysis*

*4.1.1. Model tuning*

The impact of model tuning on the optimal number of profiles ($k_{opt}$) is visualized in Fig. 2. The Cao
and Devaud metrics were found to give highly similar results, the median (Md) being either four
(VEM-Cao; VEM-Deveaud; Gibbs100A-Deveaud; Gibbs100B-Deveaud), five (Gibbs100B-Cao; all
Gibbs1000-methods) or six (Gibbs100A-Cao). In contrast, the Arun metric differed from the other
two by consistently suggesting Md($k_{opt}$) = 2, irrespective of the chosen tunings. The smallest
observed variation in $k_{opt}$ and thus the best model stability was achieved with the method
Gibbs1000B; generally, increasing the iteration number (100 vs. 1000) and/or the number of burned
iterations (0 vs. 50 or 500) reduced the variation in case of the Cao metric, whereas respective

12

254 stabilization in the Deveaud metric was achieved after a parallel increase in both parameters (100A

255 vs. 1000B). Finally, no clear trends associated with WDM column order and/or seed could be

256 detected under any tested circumstances.

*4.1.2. Cross-validation*

258 The perplexities of the cross-validation models are shown in Table 2. When considering any given

259 holdout sample (ID=1,…,16), little difference could typically be observed between different levels

260 of $k \neq 2$. On the other hand, when considering any given $k$, a distinct pattern could be observed over

261 storage time. The highest values (7.60-52.30) appeared in the beginning of storage time (days 1-3),

262 followed by a decrease between days 5 and 13 (1.78-3.63). The prediction ability was thus lowest in

263 the case of samples analyzed during the early days of storage.

*4.1.3. Profile interpretation*

265 The distributions of the top-8 VOCs within the extracted profiles (A3P1-P3; A5P1-P5; A9P1-P9) of

266 the three exploratory models are shown in Fig. 3A-C and the corresponding clustering results in

267 Fig. 3D. Irrespective of the value of $k$, three main profile clusters could be observed. In cluster 1

268 (A3P1, A5P1 and A9P5), none of the top-8 VOCs (ethanol, 3-methyl-1-butanol, dimethyl sulfide,

269 carbon disulfide, acetone, ammonia, 2,3-butanedione, 3-methylbutanal, acetic acid and/or ethyl

270 acetate) accounted for over 25 % (in terms of relative abundance) of the entire volatilome, whereas

271 cluster 2 (A3P3, A5P2, A5P5, A9P3, A9P4 and A9P8) was dominated by hydrogen sulfide +

272 ethanol and cluster 3 (A3P2, A5P3, A5P4, A9P1, A9P2, A9P6, A9P7 and A9P9) by ethanol. In the

273 latter two cases, the most abundant compound(s) accounted for over 80 % of the extracted profiles.

274 The relations between profile distribution, storage time and sensory rejection are shown in Fig. 4.

275 When considering storage time (Fig. 4A-C-E), a major shift in volatilome composition could be

276 observed between days 3-5 in all three models. Cluster 1 was found most prominent during the

277 early days (1-3) and cluster 3 during the latter days (5-11), whereas cluster 2 had a minor

278  contribution and only exceeded 50 % in two late-stage samples (days 11-13). When considering

279  sensory rejection (Fig. 4B-D-F), respective observations could be made; generally, increasing $R_\%$

280  was accompanied with decreasing contribution of cluster 1 and corresponding increase in clusters 2-

281  3. Cluster 1 dominated the volatilome of samples with less than 10 % rejection, whereas clusters 2

282  and 3 had varying contributions at a certain rejection percentage. Finally, when comparing the

283  profile distributions of any given sample (Fig. 4A-C-E), increasing $k$ was seen to cause a

284  partitioning in sub-profiles while retaining the three main clusters. For example, the partitioning of

285  A3P2 ($k=3$) into A5P4 and A5P3 ($k=5$) and further to A9P1, A9P2, A9P6, A9P7 and A9P9 ($k=9$)

286  could be observed under cluster 2. In conclusion, profiles belonging to clusters 2 and 3 were

287  considered spoilage-associated (see Subsection 5.2 for discussion) and were thus used in $SA_\%$

288  calculations (Subsection 4.1.4).

289  *4.1.4. Identification of potential spoilage indicators*

290  The relative spoilage associations of all quantified VOCs are given in Table 3. A comparison

291  between the models with 3, 5 or 9 profiles showed that 11, 9 and 5 compounds had $SA_\% \geq 0.5$,

292  respectively. More specifically, three major VOC groups could be identified: 1) 13/25 compounds

293  with $SA_\% < 0.5$ in all three models (C2, C7-C10, C13, C15, C17, C18, C22-C25), 2) 5/25

294  compounds showing an initial $SA_\%=1$ and a decreasing trend along with increasing $k$ (C4, C16,

295  C19-21) and 3) 7/25 compounds fluctuating around $SA_\%=0.5$ (C1, C3, C5, C6, C11, C12, C14).

296  *4.1.5. Spoilage characterization protocol*

297  Based on the results of the exploratory analysis (Subsections 4.1.1-4.1.4), the following protocol

298  (denoted IP$_{LDA}$) was established for LDA-based salmon spoilage characterization:

299  • *Profile number criterion*: for model training, use the lowest $k$ that does not lead to a change

300      in perplexity when compared with the optimal $k$;

301      •   *Spoilage characterization criterion:* for $SA_\%$ calculation, use profiles whose contribution

302         shows a positive correlation with the metadata (here, storage time and $R_\%$);

303      •   *Spoilage association criterion*: for identifying potential spoilage indicators, use $SA_\% = 0.5$ as

304         the cut-off threshold.

305 When comparing the performance of the 3-profile IP$_{LDA}$ (denoted IP$_{LDA3}$ in Table 3) with the

306 previously established IP$_{PLS}$, a high correspondence could be observed at the full protocol level: 5/6

307 VOCs that fulfilled the three IP$_{PLS}$ selection criteria were also identified by IP$_{LDA3}$, while 4/5 VOCs

308 having $SA_\%$=1 were also identified by IP$_{PLS}$. In contrast, a lower correspondence was observed at

309 the methodological level (LDA vs. PLS): out of 15 compounds identified by at least one of the two

310 methods, three compounds (C8, C17, C23) were missed by LDA and three other compounds (C5,

311 C11, C14) by PLS. Finally, a good overall consensus was reached in recognizing irrelevant VOCs:

312 13/25 compounds were classified as irrelevant by both protocols.

313 *4.2. Selective analysis*

314 The relations between storage time, sensory rejection and profile distribution under conditions B-D

315 are visualized in Fig. 5. Overall, the differences in the evolution of profile distributions showed that

316 the applied atmosphere had a major impact on the progression of spoilage. Under air (B), 60/0/40

317 (C) and 60/40/0 (D), profiles B3P2, C3P1 and D3P2 could be associated with spoilage,

318 respectively, whereas the other extracted profiles showed little correlation with storage time and/or

319 rejection. In conclusion, only the aforementioned three profiles were considered in $SA_\%$

320 calculations.

321 All identified potential spoilage indicators are shown in Table 3. Again (see Subsection 4.1.4),

322 correspondence between the two identification approaches was found to be higher at the protocol

323 level than at the methodological level. Out of 9, 4 and 1 compounds identified by IP$_{PLS}$ under

324 conditions B-D, respectively, 8, 3 and 1 were also identified by IP$_{LDA}$. Additional compounds

325  identified by IP$_{LDA}$ were 2,3-butanediol, acetone, acetoin, methyl mercaptan and ethyl acetate (C),

326  and ethanol, acetoin, 2,3-butanedione, carbon disulfide, dimethyl sulfide, hydrogen sulfide and

327  ethyl acetate (D). Out of these twelve additional identifications, six were also achieved by PLS.


## 5. Discussion

### 5.1 Model optimization

330  In machine learning and statistical analysis, tuning refers to a process where different parameters

331  are tested in order to optimize model performance. Even though it is generally well known that

332  tuning may greatly affect the LDA output – for example, Agrawal et al. (2018) concluded that

333  neither reusing the tunings of a preceding study nor relying on "off-the-shelf" settings can be

334  recommended – relatively few efforts of evaluating and/or controlling the impact of LDA tuning on

335  model output and performance have been published so far. For this reason, special emphasis was

336  given in the present study on a systematic selection of appropriate tunings. The key points of this

337  decision-making process are elaborated in the following paragraphs.

338  *An inference algorithm* is needed for approximating the inference of the posterior distribution. The

339  two options considered in the present study – VEM and Gibbs sampling – represent two widely

340  popular and yet fundamentally different approaches. Unlike VEM, Gibbs sampling does not

341  converge to a point estimate but generates random samples from a complex distribution, meaning

342  that the true distribution will be eventually reached (Binkley et al., 2014). The fact that VEM

343  resulted in a high variation between individual models (Fig. 2) was thus not unexpected, as it was

344  likely due to converging towards different local maxima. It must be emphasized though that the

345  choice of Gibbs sampling over VEM does not guarantee finding the true distribution *per se*, as

346  ensuring the representability of the obtained results in the former case requires additional attention

347  on the *Gibbs sampling parameters* (denoted here as *iter*, *burnin* and *thin*). Briefly, a sufficiently

348  high number of burned iterations is needed for ensuring that the sampler converges to the correct

16

349      distribution, whereas thinning has traditionally been considered advantageous in addressing the

350      risks associated with autocorrelation (Binkley et al., 2014). However, with regard to the remarks of

351      Link and Eaton (Link & Eaton, 2012) on the thinning of Markov Chain Monte Carlo (MCMC)

352      chains, a 10-fold increase in the number of iterations (100 vs. 1000) was used instead of thinning

353      for reducing the risk of autocorrelation in the present study. Given the achieved level of

354      stabilization in non-burned chains and the additional beneficial impact of burning (Subsection

355      4.1.1), the method Gibbs1000B (*iter*=1000, *burnin*=500, *thin*=1) was considered appropriate for

356      further modelling activities.

357      Irrespective of the used inference algorithm, factors such as *Dirichlet hyperparameters* and *the*

358      *order of input (WDM) columns* should be taken into account. The hyperparameters $\alpha$ and $\beta$

359      represent the Dirichlet priors on the output distributions (Binkley et al., 2014); a high value of $\alpha$ (*$\beta$,*

360      respectively) indicates that a given document (topic, resp.) is likely to consist of a broad range of

361      topics (words, resp.), whereas a low value suggests that only a few topics (words, resp.) are

362      involved. The relevance of these parameters has recently been highlighted by Park et al. (2019),

363      emphasizing that the underlying assumption of unimodality may lead to biased parameter

364      estimation if the corpus consists of clusters with different topic distributions. However, since 1) all

365      extracted profiles can initially be assumed to comprise all VOCs and to be present in each sample,

366      but 2) no prior assumptions should be made on the exact composition and/or distribution of these

367      profiles, the default settings of the packaging **topicmodels** were considered appropriate in this

368      study. On the other hand, since 1) the order of VOCs within the WDM can and must be considered

369      fully exchangeable, but 2) a high level of stabilization was achieved with the method Gibbs1000B

370      (Subsection 4.1.1), examining a single random seed-order pair was considered sufficient.

371      In this study, the impact of model tuning on the *number of extracted profiles* was assessed by means

372      of three popular metrics (Subsection 3.2.1). Since no overall consensus was reached (Subsection

373      4.1.1), cross-validation was performed. The fact that the highest perplexity coincided with the

374  beginning of storage time (days 1-3) was attributed to the characteristic deterioration patterns of

375  salmon under the tested storage condition (100 % $N_2$). In the absence of both $CO_2$ and $O_2$,

376  considerable accumulation of ethanol and sulfuric compounds was observed over storage time

377  (Kuuliala et al., 2019); for this reason, the relatively small proportion of early-day samples (days 1-

378  3) in a given training set resulted in a better fit for the late-day samples (days 5-13). However, apart

379  from the ill-performing two-profile models, little difference in perplexity was observed between

380  models with different values of $k$. For this reason, three (the smallest possible $k$ according to

381  perplexity), five (Md($k$) according to the Cao and Deveaud metrics) and nine (an example of a high

382  $k$) profiles were chosen for further exploratory activities.

383  Finally, *computational power* may pose additional challenges during exploratory analysis. As the

384  order of extracted topics is exchangeable even between consecutive iterations (Binkley et al., 2014),

385  the optimization of model parameters for a specific experimental setting requires comparison

386  between multiple corresponding but independent models. While learning a single LDA model did

387  not take considerably longer than PLS, the total time needed for generating 100 models (10 orders x

388  10 seeds) with a given method ranged from 23 minutes (VEM) to over 8.5 hours (Gibbs1000B).

389  This was anyhow considered acceptable at the exploratory stage, not only because of the benefits

390  over PLS (see Section 2 and Subsection 5.2) but also because of the positive impact on the

391  analytical workflow. In other words, sufficient emphasis on model tuning at the exploratory stage

392  considerably reduced the need for computational effort during the selective stage.

393  ### 5.2. Spoilage characterization

394  Despite the increasing popularity of topic modelling in scientific discovery, information about its

395  biological applications is still rather scarce in the current literature. Since the publication of the

396  review of Liu et al. (2006), a limited number of studies have been published on topics such as

397  genetic/protein functionality (Backenroth et al., 2017; Liu et al., 2017; Liu et al., 2018), metabolic

398  sub-structures (van der Hooft et al., 2016) and mutation signatures (Matsutani et al., 2019). The

399 present study thus extends the current state-of-the-art by introducing LDA as the basis of a

400 systematic spoilage characterization protocol. The central focus points that should be considered

401 when applying LDA as an exploratory and/or selective method within this specific context are

402 discussed in the following paragraphs; for further information about the identified compounds and

403 their role in seafood spoilage, please see the preceding study of Kuuliala et al. (2019).

404 Firstly, the method of *profile interpretation* greatly determines the performance and representability

405 of LDA. This highlights the importance of metadata in spoilage analysis, particularly if no prior

406 knowledge about the expected deterioration processes and/or potential spoilage indicators exists.

407 Hence, storage time and sensory rejection were used in the present study for confirming the

408 tentative interpretations arising from the VOC distributions *per se*. For example, the correlations

409 between increasing rejection, decreasing contribution of cluster 1 and subsequent increasing

410 contributions of clusters 2-3 under 100 % $N_2$ (Subsection 4.1.3) indicated that the former cluster

411 could be associated with freshness and the latter two with spoilage. Furthermore, the emergence of

412 multiple spoilage-associated profiles under condition A (Fig. 4) and those showing stable or

413 fluctuating patterns under conditions B-D (Fig. 5) suggest that the progression of spoilage cannot

414 necessarily be comprehensively modelled with a single profile, at least in case of small datasets

415 and/or in the absence of extremely fresh/spoiled samples.

416 *The number of extracted profiles* has both mathematical and biological relevance. Even though all

417 extracted profiles and their relations provide information about salmon quality and can thus be

418 considered relevant for exploratory purposes, a systematic selection criterion was anyhow needed

419 for selective analysis. In line with the previous PLS-based protocol (Kuuliala et al., 2019), choosing

420 the lowest $k$ (condition A: 3) that does not lead to an increase in perplexity when compared with the

421 optimal $k$ (condition A: 5) was considered appropriate. Furthermore, additional perplexity analyses

422 performed for conditions B-D (results not shown) indicated that different holdout sample ID and/or

423 $k$ had little impact on perplexity; under elevated $CO_2$ and/or $O_2$ levels, the final concentrations were

considerably lower and the relations between VOCs more stable when compared to storage under

100 % $N_2$ (Kuuliala et al., 2019), meaning that the differences between all samples (and thus

between the extracted profiles) were less pronounced when compared to those under condition A

(see Subsection 5.1). Consequently, $k=3$ was consistently used for selective purposes under all

tested conditions.

When aiming at using LDA for selective analysis, relevant *definitions* should be established in the

first place. In the case of spoilage analysis, this essentially means determining what kind of VOCs

can be considered spoilage indicators. Firstly, it should be noted that all VOCs present at a given

time point do not necessarily contribute to the perceived off-odors, as this requires exceeding

certain concentration thresholds. However, it is equally important to note that these thresholds are

both compound-specific and context-dependent. For example, while the human olfactory threshold

(OT) of a given VOC indicates its minimum perceivable concentration (Devos, Patte, Rouault,

Laffort & Van Gemert, 1990), the previously reported OTs have usually been defined for pure

compounds, whereas the seafood volatilome consists of multiple compounds which may interact or

interfere with each other. Consequently, exceeding the OT does not guarantee that a VOC can be

perceived as a part of a complex volatilome or that it contributes to offensive off-odors: the

thresholds associated with these two aspects are in fact often not well known before the commence

of the intended study. Anyhow, whether this matters depends on the scientific and/or practical

context. As previously highlighted (Ioannidis et al., 2018; Kuuliala et al., 2019), a VOC that shows

a high positive correlation with microbiological and/or sensory deterioration may have great value

in quality monitoring applications even if its individual olfactory contribution remains low or

unknown. For these reasons, the concept "potential spoilage indicator" is used in the present study

to refer to all VOCs that show promising potential in monitoring quality decay, irrespective of their

individual olfactory contribution.

Along with model optimization (Subsection 5.1), the methods of *data pre-processing* have a key role in ensuring that the applied methodology is in line with the defined aims. The relevance of this remark can be seen when comparing the performance of PLS and LDA (Table 3). In the former case, the data had been standardized in order to disregard the differences in concentration magnitudes between individual VOCs; furthermore, additional data quality criteria were implemented at the identification stage to dismiss VOCs whose quantification accuracy was considered inadequate (for more information, see Kuuliala et al., 2019). In contrast, all LDA models were learned with non-standardized data and used without additional data quality screening. The high correspondence between the final outcomes of the two protocols – despite the aforementioned methodological differences – was thus attributed to the characteristics of the salmon datasets. As elaborated in the previous study (Kuuliala et al., 2019), the data quality criteria that were used for screening VOCs primarily targeted those compounds that were present in low concentrations (< 100 ppb) throughout storage time; due to the lack of standardization, this kind of compounds (for example, all quantified amines) typically had little impact on the extracted LDA profiles. However, the lack of data quality screening also increased the number of LDA-based identifications of low-range VOCs when compared to PLS, particularly under condition D where the differences in concentration magnitudes between individual VOCs were at the lowest. Overall, these results demonstrate that while LDA is less sensitive to problems arising from the presence of low-range VOCs than PLS, an initial data quality analysis can be generally considered advisable.

For finalizing the selective analysis, representative *cut-off thresholds* are needed. In case of spoilage analysis, it is of primary importance to note that the sole presence/absence of a VOC does not automatically signify a certain quality status: instead, VOC-based spoilage characterization requires considering both the evolution and relations of all quantified compounds throughout storage time. For this reason, the concept of relative spoilage association was developed in this study to select those VOCs which had the highest (top-8) relative abundance within spoilage-associated profiles. In

line with the definition of a potential spoilage indicator (see above), this concept was developed for

quality monitoring purposes and does thus not directly signify individual olfactory contribution. For

example, only a single spoilage-associated profile per condition could be identified under

conditions B-D (Subsection 4.2), meaning that a given VOC could only receive an $SA_\%$ value 0

(irrelevant) or 1 (relevant). In the case of condition A, multiple spoilage-associated profiles were

identified, meaning that a well-established numerical cut-off limit was needed. The commonly

observed increase/decrease in $SA_\%$ as a function of $k$ (Table 3) was attributed to sub-profile

partitioning (Subsection 4.1.3), which increased the overall diversity of the top-8 VOCs and thus

reduced the number of compounds with extreme $SA_\%$ values (0 or 1). In this study, setting the limit

at $SA_\%$=0.5 was experimentally found to respond to the aforementioned needs in an optimal manner

as well as to result in a high correspondence between the two identification protocols ($IP_{LDA}$ and

$IP_{PLS}$); however, it is important to note that the selection of the cut-off limit should always be done

on a case-by-case basis.

Finally, it is advisable to evaluate *the prospects* of a newly developed methodology in a broader

context. In this case, attention was thus given to the applicability of LDA in the volatilome-based

spoilage characterization of other food products. Generally, the developed methodology is expected

to be widely applicable for examining the evolution and composition of complex volatilomes,

suggesting that highly perishable products packed under gaseous atmospheres (such as vegetables,

meat and seafood) whose quality decay is manifested by the accumulation of multiple VOCs

(resulting in unacceptable off-odors) would be the most promising target group. While model

tunings and cut-off limits should be experimentally determined whenever introducing a new

product, the basic development process described in the present study could be used as a theoretical

basis for exploring corresponding research questions in new experimental settings. In the future,

emphasis should thus be given not only on testing and validating the method in these settings, but

also on product-specific planning that precedes actual modelling activities. Preferably, the

availability and quality of the input data should be considered already before setting up a storage

experiment: since each food product has its characteristic shelf-life, the experimental setup should

allow regular and representative data collection throughout storage time. Before selective analysis,

the relation between the VOCs and the dependent variable should receive particular attention, as it

is not always linear: for example, reaching 100 % rejection does not mean that the volatilome will

also be stable from there on. In general, it should be kept in mind that the prevailing assumptions

regarding the complex quality deterioration mechanisms may affect the entire spoilage

characterization process: the fact that the training of LDA models does not inherently involve

metadata could thus be considered advantageous, especially when considering new product types.

For further insights, a comparison between unsupervised and supervised LDA would be worth

examination.


## 6.  Conclusions

The outcomes of the present study show that LDA can be successfully applied for extracting

underlying information about the quality status of Atlantic salmon under different storage

conditions, suggesting that a respective approach could be well adapted for other food products

and/or quality deterioration patterns. As a selective tool, LDA was found to identify potential

volatile spoilage indicators with equal specificity and lower stringency when compared to PLS. In

particular, the flexibility and high interpretability of LDA were considered highly advantageous in

this problem setting; not only because of their beneficial impact on the experimental workflow, but

also because of the achieved insights into the development of systematic spoilage characterization

processes. Overall, the results support the state-of-the-art of data-driven food quality

characterization, an emerging field with great prospects.

**References**

Agrawal, A., Fu, W., Menzies, T., 2018. What is wrong with topic modeling? And how to fix it using search-based software engineering. Inf. Softw. Technol. 98, 74-88.

Arun, R., Suresh, V., Veni Madhavan, C.E., Narasimha Murthy, M.N., 2010. On Finding the Natural Number of Topics with Latent Dirichlet Allocation: Some Observations. In: Zaki ,M.J., Yu, J.X., Ravindran, B., Pudi, V. (Eds.). Advances in Knowledge Discovery and Data Mining. PAKDD 2010. Lecture Notes in Computer Science, vol 6118, Springer, Berlin, Heidelberg.

Backenroth, D., He, Z., Kiryluk, K., Boeva, V., Pethukova, L., Khurana, E., Christiano, A., Buxbaum, J., Ionita-Laza, I., 2018. FUN-LDA: A Latent Dirichlet Allocation Model for Predicting Tissue-specific Functional Effects of Noncoding Variation: Methods and Applications. Am. J. Hum. Genet. 102, 920-942.

Bastani, K., Namavari, H., Shaffer, J., 2019. Latent Dirichlet allocation (LDA) for topic modeling of the CFPB consumer complaints. Expert Syst. Appl. 127, 256-271.

544    Bermejo-Prada, A., Vega, E., Pérez-Mateos, M., Otero, L., 2015. Effect of hyperbaric storage at

545    room temperature on the volatile profile of strawberry juice. LWT - Food Sci. Technol. 62, 906-

546    914.

547    Binkley, D., Heinz, D., Lawrie, D., Overfelt, J., 2014. Understanding LDA in Source Code

548    Analysis. In: Proceedings of the 22nd International Conference on Program Comprehension, ACM,

549    Hyderabad, India, pp. 26–36.

550    Blei, D.M., 2012. Probabilistic Topic Models. Commun. ACM 55, 77-84.

551    Blei, D.M., Ng, A.Y., Jordan, M.I., 2003. Latent Dirichlet Allocation. J. Mach. Learn. Res. 3, 993-

552    1022.

553    Böhme, K., Calo-Mata, P., Barros-Velázquez, J., Ortea, I., 2019. Recent applications of omics-

554    based technologies to main topics in food authentication. TrAC Trend. Anal. Chem. 110, 221-232.

555    Cao, J., Xia, T., Li, J., Zhang, Y., Tang, S., 2009. A density-based method for adaptive LDA model

556    selection. Neurocomputing 72, 1775-1781.

557    Castro-Puyana, M., Pérez-Míguez, R., Montero, L., Herrero, M., 2017. Application of mass

558    spectrometry-based metabolomics approaches for food safety, quality and traceability. TrAC Trend.

559    Anal. Chem. 93, 102-118.

560    Chen, X., Hu, X., Shen, X., Rosen, G., 2010. Probabilistic topic modeling for genomic data

561    interpretation. 2010 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)

562    Hong Kong, pp. 149-152.

563    Curiskis, S.A., Drake, B., Osborn, T.R., Kennedy, P.J., 2020. An evaluation of document clustering

564    and topic modelling in two online social networks: Twitter and Reddit. Inform. Process. Manag. 57,

565    pp. 102034.

566     Deveaud, R., Sanjuan, É., Bellot, P., 2014. Accurate and Effective Latent Concept Modeling for Ad

567     Hoc Information Retrieval. Document Numérique, Lavoisier, 2014, pp.61-84.

568     Devos M., Patte F., Rouault J., Laffort P. & Van Gemert L. J. (Eds.) (1990). Standardized human

569     olfactory thresholds. New York, US: Oxford University Press.

570     Dong, D., Jiao, L., Li, C., Zhao, C., 2019. Rapid and real-time analysis of volatile compounds

571     released from food using infrared and laser spectroscopy. TrAC Trend. Anal. Chem. 110, 410-416.

572     Feinerer, I., Hornik, K., 2018. tm: Text Mining Package. R package version 0.7-5.

573     Feinerer, I., Hornik, K., Meyer, D., 2008. Text Mining Infrastructure in R. J. Stat. Softw. 25, 1-54.

574     Fu, Y., Yan, M., Zhang, X., Xu, L., Yang, D., Kymer, J.D., 2015. Automated classification of

575     software change messages by semi-supervised Latent Dirichlet Allocation. Inf. Softw. Technol. 57,

576     369-377.

577     Ghasemi-Varnamkhasti, M., Apetrei, C., Lozano, J., Anyogu, A., 2018. Potential use of electronic

578     noses, electronic tongues and biosensors as multisensor systems for spoilage examination in foods.

579     Trends Food Sci. Technol. 80, 71-92.

580     Griffiths, T.L., Steyvers, M., 2004. Finding scientific topics. Proc. Natl. Acad. Sci. 101, 5228-5235.

581     Grün, B., Hornik, K., 2011. topicmodels: An R Package for Fitting Topic Models. J. Stat. Softw.

582     40, 1-30.

583     Hu, N., Zhang, T., Gao, B., Bose, I., 2019. What do hotel customers complain about? Text analysis

584     using structural topic model. Tour. Manag. 72, 417-426.

585     Ioannidis, A., Kerckhof, F., Riahi Drif, Y., Vanderroost, M., Boon, N., Ragaert, P., De Meulenaer,

586     B., Devlieghere, F., 2018. Characterization of spoilage markers in modified atmosphere packaged

587     iceberg lettuce. Int. J. Food Microbiol. 279, 1-13.

588    Klampfl, C.W., 2018. Ambient mass spectrometry in foodomics studies. Curr. Opin. Food Sci. 22,

589    137-144.

590    Kullback, S., Leibler, R.A., 1951. On Information and Sufficiency. Ann. Math. Stat. 22, 79-86.

591    Kuuliala, L., Abatih, E., Ioannidis, A.-G., Vanderroost, M., De Meulenaer, B., Ragaert, P.,

592    Devlieghere, F., 2018. Multivariate statistical analysis for the identification of potential seafood

593    spoilage indicators. Food Control 84, 49-60.

594    Kuuliala, L., Sader, M., Solimeo, A., Pérez-Fernández, R., Vanderroost, M., De Baets, B., De

595    Meulenaer, B., Ragaert, P., Devlieghere, F., 2019. Spoilage evaluation of raw Atlantic salmon

596    (*Salmo salar*) stored under modified atmospheres by multivariate statistics and augmented ordinal

597    regression. Int. J. Food Microbiol. 303, 46-57.

598    Li, X., Ma, Z., Peng, P., Guo, X., Huang, F., Wang, X., Guo, J. , 2018. Supervised latent Dirichlet

599    allocation with a mixture of sparse softmax. Neurocomputing 312, 324-335.

600    Link, W.A., Eaton, M.J., 2012. On thinning of chains in MCMC. Methods Ecol. Evol. 3, 112-115.

601    Liu, L., Tang, L., Dong, W., Yao, S., Zhou, W., 2016. An overview of topic modeling and its

602    current applications in bioinformatics. SpringerPlus 5, 1608.

603    Liu, L., Tang, L., He, L., Yao, S., Zhou, W., 2017. Predicting protein function via multi-label

604    supervised topic model on gene ontology. Biotechnol. Biotechnol. Equip. 31, 630-638.

605    Liu, L., Tang, L., Tang, M., Zhou, W., 2018. A partially function-to-topic model for protein

606    function prediction. BMC Genomics 19, 883.

607    Mancano, G., Mora-Ortiz, M., Claus, S.P., 2018. Recent developments in nutrimetabolomics: from

608    food characterisation to disease prevention. Curr. Opin.Food Sci. 22, 145-152.

609 Mansur, A.R., Seo, D., Song, E., Song, N., Hwang, S.H., Yoo, M., Nam, T.G., 2019. Identifying

610 potential spoilage markers in beef stored in chilled air or vacuum packaging by HS-SPME-GC-

611 TOF/MS coupled with multivariate analysis. LWT Food Sci Technol. 112, 108256.

612 Martinović, T., Šrajer Gajdošik, M., Josić, D., 2018. Sample preparation in foodomic analyses.

613 Electrophoresis 39, 1527-1542.

614 Matsutani, T., Ueno, Y., Fukunaga, T., Hamada, M., 2019. Discovering novel mutation signatures

615 by latent Dirichlet allocation with variational Bayes inference. Bioinformatics 35, 4543-4552.

616 Miguel, H., Simó C. , García-Cañas V., Ibáñez E., Alejandro, C., 2012. Foodomics: MS-based

617 strategies in modern food science and nutrition. Mass Spectrom. Rev. 31, 49-69.

618 Mikš-Krajnik, M., Yoon, Y., Ukuku, D.O., Yuk, H., 2016. Volatile chemical spoilage indexes of

619 raw Atlantic salmon (*Salmo salar*) stored under aerobic condition in relation to microbiological and

620 sensory shelf lives. Food Microbiol. 53, Part B, 182-191.

621 Murzintcev, N., 2019. ldatuning: Tuning of the Latent Dirichlet Allocation Models Parameters.  R

622 package version 1.0.0.

623 Nolasco, D., Oliveira, J., 2019. Subevents detection through topic modeling in social media posts.

624 Future Gener. Comp. Sy. 93, 290-303.

625 Odeyemi, O.A., Burke, C.M., Bolch, C.C.J., Stanley, R., 2018. Seafood spoilage microbiota and

626 associated volatile organic compounds at different storage temperatures and packaging conditions.

627 Int. J. of Food Microbiol. 280, 87-99.

628 Pavase, T.R., Lin, H., Shaikh, Q., Hussain, S., Li, Z., Ahmed, I., Lv, L., Sun, L., Shah, S.B.H.,

629 Kalhoro, M.T., 2018. Recent advances of conjugated polymer (CP) nanocomposite-based chemical

sensors and their applications in food spoilage detection: A comprehensive review. Sens. Actuators

B Chem. 273, 1113-1138.

Perina, A., Lovato, P., Murino, V., Bicego, M., 2010. Biologically-aware Latent Dirichlet

Allocation (BaLDA) for the Classification of Expression Microarray. In: Dijkstra, T.M.H.,

Tsivtsivadze, E., Marchiori, E., Heskes, T. (Eds.), Pattern Recognition in Bioinformatics. PRIB

2010. Lecture Notes in Computer Science, vol 6282. Springer, Berlin, Heidelberg.

Pinu, F.R., 2016. Early detection of food pathogens and food spoilage microorganisms: Application

of metabolomics. Trends Food Sci. Technol. 54, 213-215.

Poghossian, A., Geissler, H., Schöning, M.J., 2019. Rapid methods and sensors for milk quality

monitoring and spoilage detection. Biosens. Bioelectron. 140, 111272.

Pratanwanich, N., Lio, P., 2014. Exploring the complexity of pathway–drug relationships using

latent Dirichlet allocation. Comput. Biol. Chem. 53, 144-152.

R Core Team, 2019. R: A language and environment for statistical computing. R Foundation for

Statistical Computing, Vienna, Austria. URL: https://www.R-project.org/.

Shiraishi, Y., Tremmel, G., Miyano, S., Stephens, M., 2015. A Simple Model-Based Approach to

Inferring and Visualizing Cancer Mutation Signatures. PLoS genetics 11, e1005657.

Silge, J., Robinson, D., 2016. tidytext: Text Mining and Analysis Using Tidy Data Principles in R.

J. Open Source Softw. 1, 37.

Suzuki, R., Terada Y., Shimodaira, H., 2019.pvclust: Hierarchical Clustering with P-Values via

Multiscale Bootstrap Resampling. R package version 2.2-0.

van der Hooft, J., Wandy, J., Barrett, M.P., Burgess, K.E.V., Rogers, S., 2016. Topic modeling for

untargeted substructure exploration in metabolomics. Proc. Natl. Acad. Sci. 113, 13738-13743.

652 Wang, Y., Li, Y., Yang, J., Ruan, J., Sun, C., 2016. Microbial volatile organic compounds and their

653 application in microorganism identification in foodstuff. TrAC Trend. Anal. Chem. 78, 1-16.

654 Warnes, G.R., Bolker, B., Bonebakker, L., Gentleman, R., Liaw, W.H.A., Lumley, T., Maechler,

655 M., Magnusson, A., Moeller, S., Schwartz, M., Venables, B., 2020. gplots: Various R Programming

656 Tools for Plotting Data. R package version 3.0.3.

657 Wickham, H., 2016. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag, New York,

658 Xiong, H., Cheng, Y., Zhao, W., Liu, J., 2019. Analyzing scientific research topics in

659 manufacturing field using a topic model. Comput. Ind. Eng. 135, 333-347.

660 Xu, Y., 2017. Foodomics: A novel approach for food microbiology. TrAC Trend. Anal. Chem. 96,

661 14-21.

662 Yang, M., Qu, Q., Chen, X., Tu, W., Shen, Y., Zhu, J., 2019. Discovering author interest evolution

663 in order-sensitive and Semantic-aware topic modeling. Inform. Sciences 486, 271-286.

664 Yu, K., Gong, B., Lee, M., Liu, Z., Xu, J., Perkins, R., Tong, W., 2014. Discovering Functional

665 Modules by Topic Modeling RNA-Seq Based Toxicogenomic Data. Chem. Res. Toxicol. 27, 1528-

666 1536.

667 Zhang, J., Liu, B., He, J., Ma, L., Li, J., 2012. Inferring functional miRNA–mRNA regulatory

668 modules in epithelial–mesenchymal transition with a probabilistic topic model. Comput. Biol.

669 Med.42, 428-437.

**Figure captions**

**Fig. 1.** The distribution of the top-20 words in five topics extracted from a collection of 121 original research abstracts, published by the Research Unit Food Microbiology and Food Preservation (FMFP; Ghent University, Ghent Belgium) between 2000 and 2018.

**Fig. 2.** The optimal number of profiles ($k$), extracted from the volatilome of Atlantic salmon fillet portions stored under 100 % $N_2$ (condition A) at 4 °C. For methodological specifications concerning the applied metrics (Cao, Arun, Deveaud), inference estimation methods (VEM, Gibbs100A, Gibbs100B, Gibbs1000A and Gibbs1000B), seeds (s1-s10) and input column orders (O1-O10), see Subsection 3.2.1.

**Fig. 3.** The distribution of the top-8 volatiles in the profiles extracted from the volatilome of Atlantic salmon fillet portions stored under 100 % $N_2$ (condition A) at 4 °C; A) $k = 3$, B) $k = 5$, C) $k = 9$, and D) hierarchical clustering of all profiles.

**Fig. 4.** The relations between storage time (from left to right: sample ID 1-16), sensory rejection ($R_\%$), and profile/cluster distribution in Atlantic salmon fillet portions stored under 100 % $N_2$ (condition A) at 4 °C; A-B) $k = 3$, C-D) $k = 5$, E-F) $k = 9$.

**Fig. 5.** The relations between storage time (from left to right: sample ID 1-16), sensory rejection ($R_\%$) and profile distribution ($k = 3$) in Atlantic salmon fillet portions stored under different gaseous conditions (% $CO_2/O_2/N_2$) at 4 °C; air (B), 60/0/40 (C) and 60/40/0 (D).