

## Journal of the Hellenic Veterinary Medical Society

Vol 72, No 4 (2021)



### A 30 year topic analysis of Veterinary Medicine literature

I FYTILAKOS, V ALEXOPOULOS

doi: [10.12681/jhvms.29430](https://doi.org/10.12681/jhvms.29430)

### To cite this article:

FYTILAKOS, I., & ALEXOPOULOS, V. (2022). A 30 year topic analysis of Veterinary Medicine literature. *Journal of the Hellenic Veterinary Medical Society*, 72(4), 3401–3414. <https://doi.org/10.12681/jhvms.29430>

## A 30 year topic analysis of Veterinary Medicine literature

I. Fytilakos<sup>1</sup>, V. Alexopoulos<sup>2</sup>

<sup>1</sup>*Department of Zoology-Marine Biology, National and Kapodistrian University of Athens*

<sup>2</sup>*School of Veterinary Medicine, University of Thessaly, Karditsa, Greece*

**ABSTRACT:** In the present study Latent Dirichlet allocation (LDA) was used as a generative probabilistic model to extract major topics in interdecadal research for the Veterinary Medicine scientific literature. A total of 22 topics were extracted during the 1991-2000 period, 23 topics during 2001-2010 and 60 topics during 2011-2020. Three different algorithms were used to validate the model: perplexity, silhouette clustering and gradient boosted trees. All three validation metrics showed that LDA performed well in extracting topics. Each decade was characterized by unique topics as well as common topics which existed throughout periods. The most frequent topics were identified and trends were quantified with the use of indexes. A list of the 30 most frequent and most associated with the term Veterinary Medicine words is provided. A shift in scientific thinking probably occurred during the 30-year-period in the process of incorporating the fields related to Veterinary students, antimicrobial resistance and animals' behavior.

**Keywords:** Veterinary Medicine, latent Dirichlet allocation, topic modeling, model validation, literature trends

*Corresponding Author:*

Ioannis Fytilakos, Department of Zoology-Marine Biology, National and Kapodistrian University of Athens, Panepistimiopolis, 15701, Ilisia, Athens, Greece  
E-mail address: ifytilakos@gmail.com

*Date of initial submission: 28-09-2020*  
*Date of acceptance: 12-11-2020*

## INTRODUCTION

The more classical approach of collecting information and exploratory analysis in scientific literature includes qualitative research methods, which offer high flexibility and focus well on understanding a problem. However, there are faster and less subjective methods to study the literature. Quantitative research methods offer solutions as they include easy data collection and analysis procedures and are not affected by the subjectivity of the researcher (Queirós et al., 2017).

The Latent Dirichlet allocation (LDA) is a generative probabilistic Bayesian model for collecting of discrete data such as text corpora (Blei et al., 2003). In LDA, text documents are used as a collection of words to identify underlying topics. Unsupervised machine learning techniques such as LDA, require little prior work from the researcher and are able to categorize big amount of data. However, LDA can also generate ambiguous topics which are hard to interpret and to classify as discussed previously in Web analysis research (Nanni, 2017). LDA has been successfully applied in the past in various fields such as in biology, biodiversity, climate change and animal communities (Zhang et al., 2019).

Studies in the literature of Veterinary Medicine with the use of LDA are scarce. However other advanced statistical techniques of machine learning have been applied on necropsy reports for detecting emerging diseases (Bollig et al., 2020) and in a literature review of urothelial cancer (Lin et al., 2020). Machine-learning-based literature mining may analyze large collections of documents, identify patterns in a dataset using statistical and computational methods or make predictions based on the discovered patterns (Lin et al., 2020). It is useful in summarizing key research themes and trends (Lin et al., 2020).

Previous works on text information extraction of literature mainly used text mining processes to study several subfields of Veterinary Medicine such as poor animal welfare (Contiero et al., 2019), epidemiology (Van der Waal et al., 2017), studies in livestock animals (Sahadevan et al., 2012), parasitology (Ellis et al., 2020), studies in antimicrobial prescribing practices (Welsh et al., 2017) or in geographic trends of science (Christopher and Marusic, 2013). Furthermore, another aspect of extracting information from text is the construction of automated electronic surveillance systems in order to predict emergency situations regarding disease outbreak (Lustgarten et al., 2020; Dórea et al., 2015), companion animal syn-

dromes (Anholt et al., 2014), temporal and spatial features of diseases (Bollig et al., 2020) or in decision support frameworks (Jones-Diette et al., 2019).

Although many information collection techniques have been applied in Veterinary Medicine in the past for the purposes of surveillance systems, historical studies on an extended temporal scale has not been conducted. This is the first attempt that aims at clarifying the major scientific terms of literature during a 30 year period, from 1991 to present. Topic extraction using LDA modeling is the main purpose of the present study to reveal the diachronically major topics in the field of Veterinary Medicine and to ascertain possible interdecadal differences in trends. Simultaneously, a lack of validation processes in related literature has been observed, thus a three-way validation approach was followed in this study to calculate the accuracy of the LDA model in predicting topics.

## MATERIALS AND METHODS

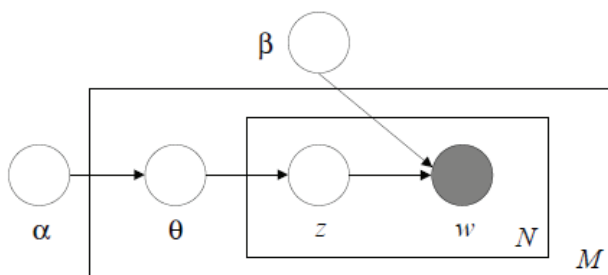
Abstracts of publications related to Veterinary Science were extracted from the Web of Science database. For this purpose the term “Veterinary Science” in quotes was entered in the search engine to extract only the Abstracts belonging to the category of Veterinary Sciences. Research articles, Reviews and Conference Proceeding papers were selected as they are the source of Abstracts and the study was planned at the level of decade; thus three decades 1991-2000, 2001-2010 and 2011-2020 were used as a filter to extract 857, 1.732 and 3.256 abstracts respectively. Abstracts were stored in three separate files representing decades.

At first, a pre-processing stage transformed Abstract texts into words: special characters, symbols, numbers and articles were excluded from the analysis with the “stop words” procedure which was common for all three decades. KH coder (Higuchi, 2016) is able to analyze English data by grouping derivatives based on a built-in dictionary (it is called lemmatization and uses the Stanford POS Tagger toolkit) or by cutting the last letters and grouping words by their stem (it is called stemming and uses the Snowball Stemmer toolkit). For instance, the term “veterinary” during the stemming process becomes “veterinari”, terms “veterinarian/ veterinarians” become “veterinarian” and terms “tumor/ tumors/ tumorous/ tumoral” become “tumor”. Both toolkits were tested for their ability to group derivatives and to extract a large number of words. After the data preparation process, a pre-processing command was used to segment the word file

into words. This is a necessary internal process to organize the results in a SQL database form, to carry out searching and tabulating (Higuchi, 2016).

Quantitative analysis of text data followed. Words in the documents were counted to obtain the number of appearances. A word frequency table was constructed with the 30 most frequent words in each decade. A comparison of term frequencies was done to study differences between decades. For this purpose word frequencies were normalized as documents contained a different number of abstracts. A Word Association Table was constructed for each decade using the Jaccard coefficient to determine the associations between words and the term “Veterinary Science”. The Jaccard coefficient emphasizes whether or not specific words co-occur, and is suitable for analyzing sparse data and is also calculated irrespectively of the term frequency (Higuchi, 2016). The values of Jaccard coefficient vary between 0 and 1. In KH Coder, words that appear frequently in the same sentence/paragraph are considered to be closely associated, and in that case the Jaccard coefficient reaches 1.

The Latent Dirichlet allocation (LDA) was used as a generative probabilistic model to extract topics from documents. LDA is based on the idea that documents represent a distribution of words which surround a topic. The model assumes that we predefine the number of topics into a document ( $k$ ), one parameter for the distribution of topics into a document ( $\alpha$ ) and one parameter for the distribution of words into topics ( $\beta$ ). These two parameters were set as  $\alpha = 50/N$  and  $\beta = 0.1$  in all runs of the algorithm. A plate notation from (Blei et al., 2003) of LDA variables is presented below, where:



$M$  denotes the number of Abstracts

$N$  is number of words in a given Abstract (Abstract  $i$  has  $N_i$  words)

$\alpha$  is the parameter of the Dirichlet prior on the per-Abstract topic distributions

$\beta$  is the parameter of the Dirichlet prior on the per-topic word distribution

$\theta_i$  is the expected topic proportion of Abstract  $M$ , which is generated by a Dirichlet distribution parameterized by parameter  $\alpha$

$z$  is the topic for the  $n$ th word in Abstract  $M$  and

$w$  is the word in the  $n$ th position word of Abstract  $M$ .

Three measurements were used to evaluate the effectiveness of the LDA modeling those of perplexity, gradient boosted trees (GBT) and silhouette clustering. Validating LDA is a hard procedure as unsupervised machine learning uses data without pre-existing labels. For this purpose perplexity was used to define the ideal number of topics into each document. GBT were used to evaluate the performance of LDA in identifying topics and silhouette clustering with a manual labeling procedure was used in order to validate topics extracted with LDA.

Perplexity is a measurement of how well a probability model predicts a sample. In language modeling, perplexity decreases in the likelihood of the test data. A lower perplexity score indicates a better generalization performance (Blei et al., 2003). Due to the large volume of words collected during the 2011-2020 period, perplexity tends to get minimized at over 150 topics which makes it difficult to interpret such a large number of topics. Considering that word frequency of the endmost topics was very low only the 60 most important topics were finally presented concerning this decade in terms of word frequency.

$$\text{perplexity (Document)} = \exp \left\{ - \frac{\sum_{d=1}^M \log p(w_d)}{\sum_{d=1}^M N_d} \right\}$$

$M$  denotes the number of Abstracts

$N$  is number of words in a given Abstract (Abstract  $i$  has  $N_i$  words)

$p$  is the probability distribution of the model

$w$  is a word-level variable.

Gradient boosted trees is an ensemble consisting of a set of alternative models using multiple learning algorithms to produce a more accurate classifier than that of the standard classifier (Opitz and Maclin, 1999). In this case the whole dataset (5.845 Abstracts) was divided in two random subsets for a total of 10 random times to compare the model performance. These subsets were used in combinations of two to produce 45 different trees to evaluate the performance of LDA in identifying topics. Each subset was used to extract 40 topics, thus 80 topics were compared each time (40 topics from the training set and 40 topics from the prediction set) which were manually labeled. Manual labeling is a time-consuming process and rules have to be followed in order to obtain meaningful topics for a neutral individual. Three rules were applied to label topics: at least two words are necessary in order to label a topic meaningful; if there were more than two common words between topics, the label was given according to a third

word etc., and words in a similar context were considered as belonging to the same topic. For instance, the words tumor, tumour, malign, lymphoma, neoplasm, cancer were all identified under the topic of cancer and the words echocardiography, cardiac, pulmonary, valve, doppler, heart, pressure, arrhythmia, myocardium were all identified under the topic of heart-related problems. The weight of each word was used as an advisory index but not as a criterion in the identification of common topics between training and prediction sets. An example of labeling is given in Table 1.

**Table 1.** Example of the manual process followed to give labels to topics according to their meaning. Topics with the same context were given same labels (Topic 1 was the label for tumor related context) while topics with different context were given different labels (Topic 2 and 3)

Training set		Prediction set	
Word	Weight	Word	Weight
Topic1		Topic1	
tumor	394.0	tumor	284.0
cell	316.0	cell	256.0
histopatholog	150.0	tumour	153.0
tumour	144.0	lymphoma	104.0
histolog	115.0	cytolog	99.0
Topic2		Topic3	
effect	87.0	vaccin	198.0
dog	72.0	infect	182.0
agent	70.0	antibodi	101.0
chemotherapi	57.0	immun	81.0
treatment	56.0	virus	63.0

An indirect use of silhouette clustering was used to validate the results of LDA. Silhouette clustering is a method of interpretation and validation of consistency and cohesion within clusters of data. The average silhouette width can be used to evaluate cluster validity (Rousseeuw, 1987). As LDA extracted topics without any labels applied, the same manual procedure described above was used to give labels to each topic. An agglomerative hierarchical clustering of the two sets combined (each set consisted of 40 topics \* five most frequent words of each topic = 200 words \* two sets = 400 words) was used to identify groups of topics at the level of sets (training and prediction). Silhouette clustering was then applied to measure the consistency and the cohesion of 45 different combinations of clusters. A silhouette value (score) of one data point can be calculated with the formula:

$$s(i) = \frac{b(i)-a(i)}{\max\{a(i),b(i)\}}, -1 \leq s(i) \leq 1$$

where  $a(i)$  is the dissimilarity 'within' a cluster and  $b(i)$  is the dissimilarity 'between' clusters.

The range of silhouette scores varies between +1 for objects that classified well in the predefined clusters (those of training and prediction sets) and -1 for objects that have been misclassified. In our case an accurate model would lead to small 'between' and 'within' distance dissimilarities between the training and the prediction sets, thus  $s(i)$  score would tend to zero. The overall average silhouette width for the entire plot (the average of the  $s(i)$  for all objects  $i$  belonging to the whole dataset)(Rousseeuw, 1987), was calculated with the Euclidean distance as a metric to indirectly calculate the accuracy of the model.

Three measures of topic diagnostic information were used to highlight the most frequent topics of each decade those of document entropy, document burstiness and corpus distance. Entropy is the degree of probability of a topic to occur in different documents (Abstracts in our case). The concept of entropy in information systems has been introduced by Shannon (1948) and has been widely used in topic extraction analysis. A topic of higher entropy is possible to have been extracted by a high number of several documents. Burstiness in natural language documents is the property of the most common words to represent a large number of topics (Boyd-Graber et al., 2014). This relationship of a few words representing a majority of topics and vice versa seems to follow a Zipfian distribution, an empirical law previously observed in social and physical sciences (Boyd-Graber et al., 2014). The combination of entropy and burstiness gives us a measure of the most common words and topics of literature. Finally the distance of topics from the corpus of each decade was measured with the Kullback-Leibler divergence distribution (Kullback and Leibler, 1951). A lesser distance from a corpus indicates that a topic is closely related to Veterinary medicine. Closely related topics consist of the most frequent words of the corpus. On the other hand, more distant topics are more distinct from the main corpus. A three-decade comparison was carried out to ascertain trends in corpus-related topics. A list of the most frequent words appearing before and after selected terms is presented in the Supplementary material, aiming not only to help but also to furtherly promote the interpretation of the topics and the comprehension of their position into the text.

Three open source free software were used for the procedures of document preprocessing, topic extraction and silhouette clustering: KH coder v. 3.Beta.01a (Higuchi, 2016), Orange v. 3.26.0 (Demsar et al., 2013) and RapidMiner v. 9.7.2+ (Mierswa et al., 2006).



## RESULTS

Word stemming extracted almost the same number of words compared with the lemmatization procedure and was preferred for its better grouping ability (Table 2). From the total number of terms approximately half of them were excluded from the analysis (Table 2). The majority of words were common between decades with a different order of appearance (Table 3). Word association analysis revealed 14 strongly associated unique words in the between decade comparison (Table 3). Four words (dure, test, drug, system) were unique during 1991-2000, four words (provide, patient, student, medic) during 2001-2010 and six words (associ, dog, perform, compare, product, group) during 2011-2020. Normalized word frequency comparison revealed an increase in the use

of common words such as dog, cat, pig, antibiotics and tumor and a decrease in words such as substance, vaccine, market, public and epidemiology during the decade 2011-2020. A detailed list of comparisons is given in Table 4.

**Table 2.** Comparison of two different algorithms for their ability to extract words from documents. The Tagger algorithm uses a built-in dictionary while the Stemmer algorithm groups words by their stem and then cuts their last letters

	1991-2000	2001-2010	2011-2020
<b>Stemmer</b>			
<b>Tokens in total</b>	122627	314517	762825
<b>Tokens in use</b>	67094	175453	431440
<b>Tagger</b>			
<b>Tokens in total</b>	126876	326773	797047
<b>Tokens in use</b>	67264	175970	433491

**Table 3.** List of the 30 most frequent and most associated words with the term Veterinary Medicine for each decade. Numbers correspond to frequencies and to Jaccard coefficient (JC). Words that appear frequently in the same abstract are closely associated thus JC reaches 1

<b>Term Frequency</b>						<b>Associated Words</b>					
<b>1991-2000</b>			<b>2001-2010</b>			<b>2011-2020</b>			<b>1991-2000</b>		
veterinari	1527		veterinari	3779		veterinari	6528		medicin	0.8672	
medicin	1164		medicin	2538		medicin	4561		anim	0.2721	
anim	813		anim	1833		anim	3809		clinic	0.2053	
dog	483		dog	1406		studi	3704		result	0.1787	
clinic	396		clinic	1192		dog	3521		studi	0.1759	
diseas	378		studi	1121		clinic	2540		develop	0.1726	
studi	358		case	962		treatment	2104		human	0.1725	
result	309		human	948		result	2020		diseas	0.1713	
human	285		diseas	914		group	1989		import	0.1521	
effect	281		student	723		diseas	1936		year	0.1448	
develop	271		result	692		effect	1846		effect	0.1367	
treatment	269		develop	666		human	1846		increas	0.1287	
test	268		health	648		signific	1760		method	0.1265	
drug	267		treatment	634		differ	1660		includ	0.1259	
method	246		resist	611		increas	1656		practic	0.1251	
case	243		report	600		case	1583		univers	0.1240	
differ	226		increas	599		evalu	1583		treatment	0.1169	
year	219		effect	573		cell	1501		differ	0.1167	
practic	206		year	569		student	1451		examin	0.1088	
system	205		differ	564		perform	1377		review	0.1074	
cat	202		cat	552		report	1372		case	0.1066	
increas	202		practic	552		associ	1330		time	0.1044	
import	201		evalu	547		cat	1321		dure	0.1028	
group	200		univers	533		resist	1307		test	0.1022	
signific	200		examin	512		product	1216		signific	0.1005	
veterinarian	193		includ	508		test	1189		determin	0.0995	
infect	188		veterinarian	508		compar	1182		evalu	0.0993	
evalu	186		patient	491		includ	1177		articl	0.0978	
examin	186		hors	485		control	1146		drug	0.0968	
includ	185		infect	481		sampl	1146		system	0.0955	

**Table 4.** Normalized word frequency comparisons of selected words presented in subcategories. Arrows indicate a frequency value over [ $\uparrow$ ], below [ $\downarrow$ ] or equal [ $\square$ ] to the average of the 1991-2020 study period. Red gradient was used to highlight higher word frequency values of each subcategory

Animals	dog	cat	hors	pig	bird	sheep	domest	mammari	cow	bovin	breed
	$\downarrow$ 14,61	$\downarrow$ 5,95	$\downarrow$ 5,03	$\downarrow$ 1,27	$\downarrow$ 1,18	$\downarrow$ 1,08	$\downarrow$ 0,99	$\downarrow$ 0,62	$\downarrow$ 2,04	$\downarrow$ 1,24	$\downarrow$ 2,50
	$\downarrow$ 21,20	$\downarrow$ 8,19	$\downarrow$ 7,17	$\downarrow$ 1,68	$\downarrow$ 0,55	$\downarrow$ 0,87	$\downarrow$ 1,36	$\downarrow$ 0,90	$\downarrow$ 1,56	$\downarrow$ 1,03	$\downarrow$ 2,72
Diet	feed	milk	dairi	nutrit	Other	infecti	care	chronic	protocol	routin	altern
	$\downarrow$ 1,58	$\downarrow$ 1,42	$\downarrow$ 0,65	$\downarrow$ 1,48		$\downarrow$ 1,14	$\downarrow$ 1,21	$\downarrow$ 0,87	$\downarrow$ 0,81	$\downarrow$ 0,81	$\downarrow$ 1,33
	$\downarrow$ 0,97	$\downarrow$ 0,98	$\downarrow$ 1,16	$\downarrow$ 0,94		$\downarrow$ 1,26	$\downarrow$ 2,72	$\downarrow$ 1,32	$\downarrow$ 0,81	$\downarrow$ 0,97	$\downarrow$ 1,09
Immune system	antibiot	antibodi	vaccin	antimicrobi	antigen	serum	plasma	immun	virus	metabol	dna
	$\downarrow$ 3,73	$\downarrow$ 1,48	$\downarrow$ 2,47	$\downarrow$ 1,67	$\downarrow$ 0,53	$\downarrow$ 2,72	$\downarrow$ 1,45	$\downarrow$ 1,24	$\downarrow$ 1,14	$\downarrow$ 0,90	$\downarrow$ 0,56
	$\downarrow$ 2,55	$\downarrow$ 0,69	$\downarrow$ 3,31	$\downarrow$ 5,29	$\downarrow$ 0,69	$\downarrow$ 2,86	$\downarrow$ 1,82	$\downarrow$ 1,47	$\downarrow$ 1,83	$\downarrow$ 0,80	$\downarrow$ 0,58
Professions	human	student	administ	econom	safeti	market	owner	materi	public		
	$\downarrow$ 8,51	$\downarrow$ 2,84	$\downarrow$ 0,93	$\downarrow$ 0,93	$\downarrow$ 0,87	$\downarrow$ 0,81	$\downarrow$ 0,81	$\downarrow$ 0,62	$\downarrow$ 2,62		
	$\downarrow$ 14,22	$\downarrow$ 10,79	$\downarrow$ 1,18	$\downarrow$ 0,84	$\downarrow$ 1,18	$\downarrow$ 0,51	$\downarrow$ 1,99	$\downarrow$ 1,64	$\downarrow$ 5,70		
Pathology / Diseases	surgeri	gastrointestin	dose	tumor	disord	therapeut	joint	pressur	pharmacokinet	salmonella	respiratori
	$\downarrow$ 2,87	$\downarrow$ 0,87	$\downarrow$ 2,32	$\downarrow$ 1,91	$\downarrow$ 1,85	$\downarrow$ 1,76	$\downarrow$ 1,73	$\downarrow$ 1,54	$\downarrow$ 1,48	$\downarrow$ 0,56	$\downarrow$ 1,33
	$\downarrow$ 2,58	$\downarrow$ 0,45	$\downarrow$ 2,69	$\downarrow$ 3,69	$\downarrow$ 1,62	$\downarrow$ 2,02	$\downarrow$ 0,60	$\downarrow$ 1,29	$\downarrow$ 1,03	$\downarrow$ 0,69	$\downarrow$ 1,41
	renal	patholog	surgic	bone	pathogen	heart	cardiac	diagnos	mastiti	biochem	diabet
	$\downarrow$ 1,51	$\downarrow$ 1,45	$\downarrow$ 1,42	$\downarrow$ 1,36	$\downarrow$ 1,36	$\downarrow$ 0,68	$\downarrow$ 0,65	$\downarrow$ 1,30	$\downarrow$ 0,59	$\downarrow$ 0,68	$\downarrow$ 0,68
	$\downarrow$ 1,13	$\downarrow$ 3,22	$\downarrow$ 2,75	$\downarrow$ 2,35	$\downarrow$ 2,43	$\downarrow$ 1,09	$\downarrow$ 1,04	$\downarrow$ 2,70	$\downarrow$ 0,43	$\downarrow$ 0,60	$\downarrow$ 0,78
	epidemiolog	physiolog	strain	genet	mechan	oral	arteri	enrofloxacin	pharmacolog	urinari	anaesthet
	$\downarrow$ 1,30	$\downarrow$ 1,27	$\downarrow$ 1,27	$\downarrow$ 1,24	$\downarrow$ 0,96	$\downarrow$ 1,11	$\downarrow$ 1,05	$\downarrow$ 1,02	$\downarrow$ 0,77	$\downarrow$ 0,71	$\downarrow$ 0,68
	$\downarrow$ 1,15	$\downarrow$ 1,18	$\downarrow$ 2,93	$\downarrow$ 1,21	$\downarrow$ 1,61	$\downarrow$ 2,35	$\downarrow$ 0,92	$\downarrow$ 0,80	$\downarrow$ 1,10	$\downarrow$ 0,81	$\downarrow$ 0,39
	1,11	1,14	2,80	1,56	2,27	2,14	1,60	0,63	0,50	1,16	0,00

LDA successfully extracted 22 topics during 1991-2000, 23 topics during 2001-2010 and 60 topics during 2011-2020 (Figure 1-3) as indicated by the perplexity score (Figure 4). Topics including the term “veterinari” differed between decades. Six topics with this term were extracted during 1991-2000, seven topics during 2001-2010 and six topics during 2011-2020 in a relevant different volume of literature Abstracts. During the first decade the terms study, animal-veterinarian, student-education, school, surgery-techniques and drugs-use-effects co-appeared with the term veterinary. During the second decade the terms practice, student-education-program, clinic-examination, human-treatment, study-group, faculty-school and animal-health co-appeared with the term veterinary. Finally the terms clinic-case, university-research, animal-human-health, disease-human-patient, clinic-studies and student co-appeared with the term veterinary during the third decade. Apart from the topics directly related to the overly generalized terms “Veterinari” and “medicine” which were common for all three decades, several other topics were extracted many of which were very specific. For instance, the topics “ultrasound-arteri-pregnanc-fetal-Doppler”, “cell-tumor-tumour-dog-carcinoma” and “resist-antimicrobi-isol-strain-antibiot” refer to artery exam-

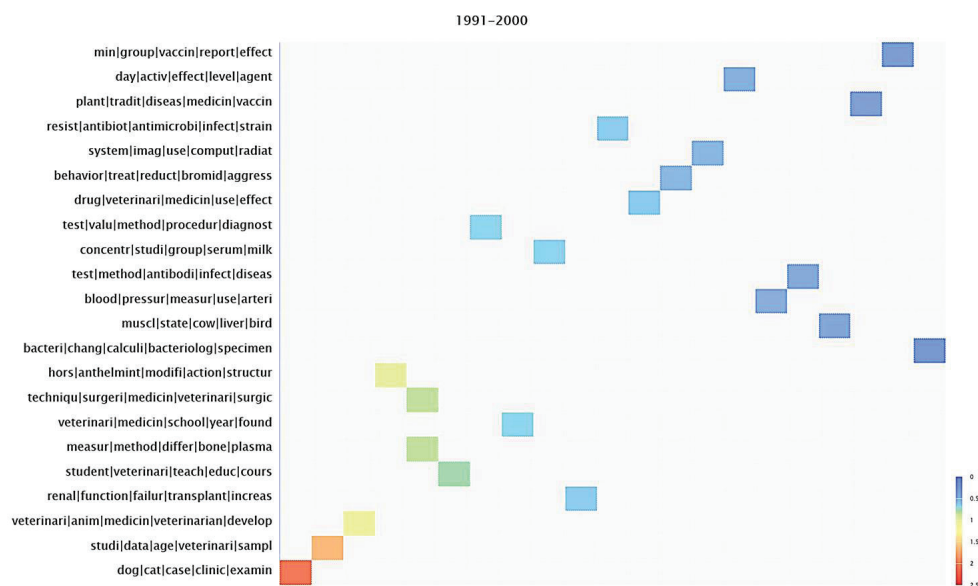
ination, cancer in dogs and antimicrobial resistance respectively.

GBT modeling indicated that LDA is able to extract topics with a root-mean-square error (RMSE) of  $2.759 \pm 0.053$ . This means that proportionally with the 40 topics extracted from each subset during the validation process, a mean of 2.759 topics were not predicted correctly. Cluster analysis of the 45 different combinations of sets grouped the 80 topics into two clusters at the level of training and prediction sets indicating the degree of correctly predicted (or respectively incorrectly predicted) topics. Silhouette scores of each cluster analysis indicated the number of topics from the training set that was predicted incorrectly (Figure 5). An overall average of  $s(i) = 0.079864 (\pm 2\sigma \text{ of } 0.02429)$  means that there are small ‘between’ and ‘within’ dissimilarities between training and prediction sets, thus the model is accurate.

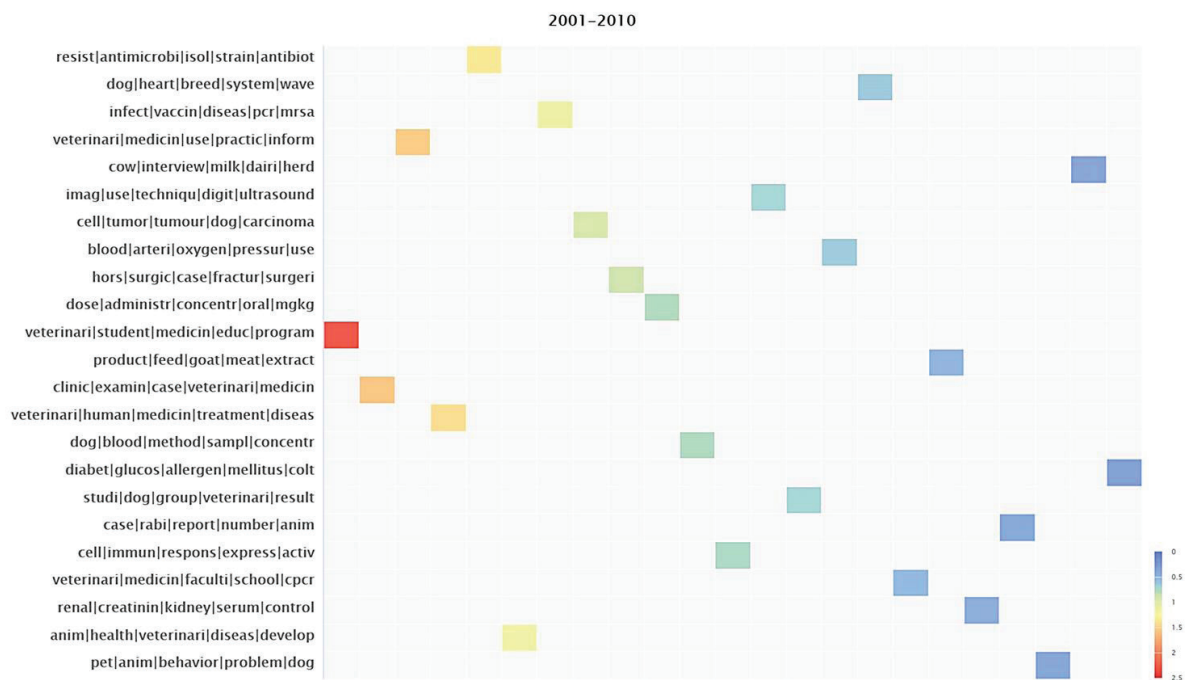
Each decade consisted of different topics of high entropy and high burstiness (Figure 1-3). Topics “dog-cat-case-clinic-examin”, “studi-data-age-veterinari-sampl” and “veterinari-anim-medicin-veterinarian-develop” were most frequent during the first decade. Topics “veterinari-student-medicin-educ-pro-

gram”, “veterinari-medicin-use-practic-inform”, “clinic-examin-case-veterinari-medicin” and “veterinari-human-medicin-treatment-diseas” were most frequent during the second decade. Finally topics “student-veterinari-medicin-studi-survey”, “anim-veterinari-medicin-health-human”, “diet-group-feed-supplement-fed”, “veterinari-medicin-diseas-human-patient” and “case-dog-veterinari-clinic-report”

were most frequent during the third decade. A distance from the corpus revealed that “antimicrobial resistance” and “veterinary students” were more distinct topics from the corpus of Veterinary medicine during 1991-2000 while “animal behavior problems” and “surgeries” were more distinct during 2011-2020 (Figure 6).

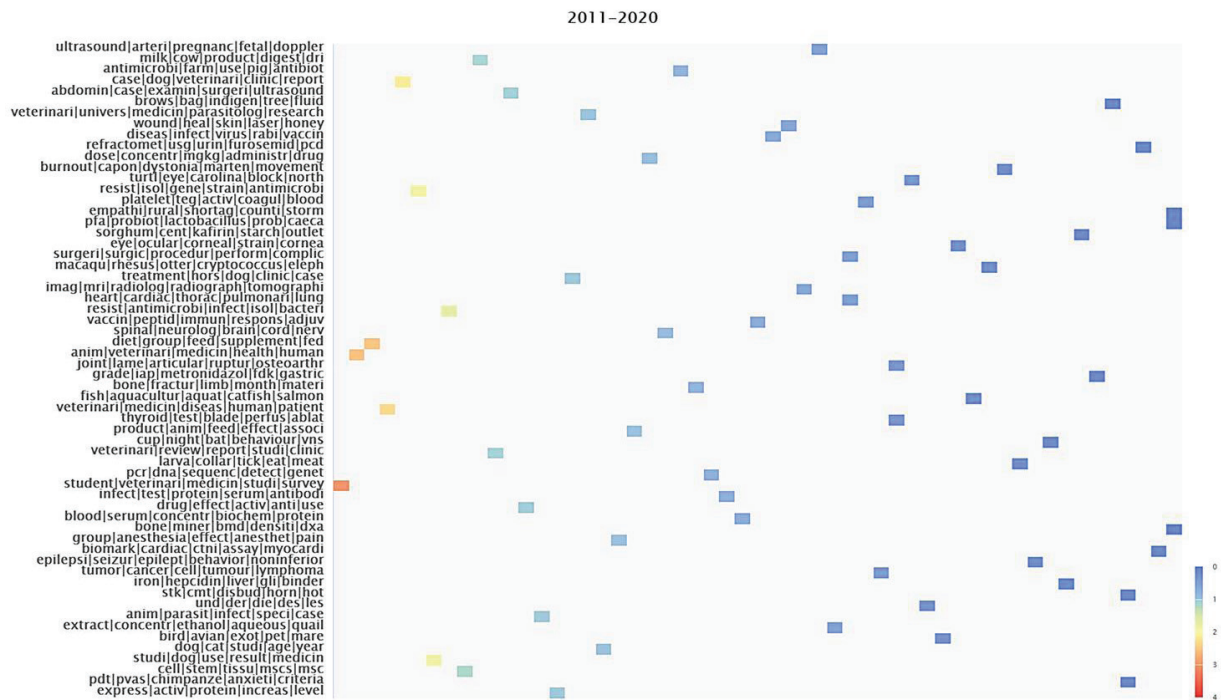


**Figure 1.** Topics extracted from the period of 1991-2000 with the LDA analysis. The combination of entropy and burstiness measures was used as an index to quantify the most frequent topics of literature (red color)

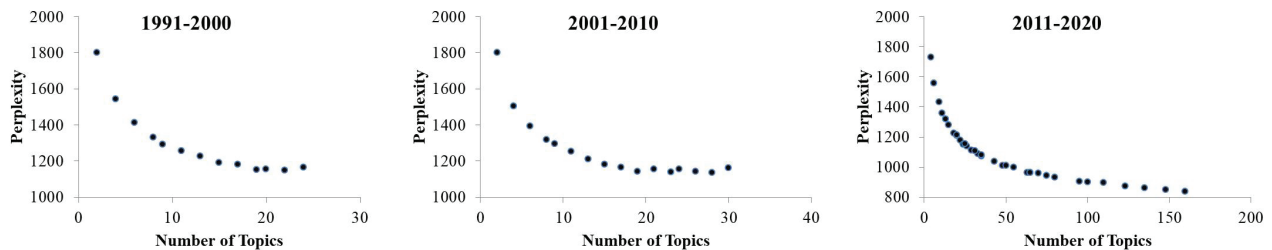


**Figure 2.** Topics extracted from the period of 2001-2010 with the LDA analysis. The combination of entropy and burstiness measures was used as an index to quantify the most frequent topics of literature (red color)

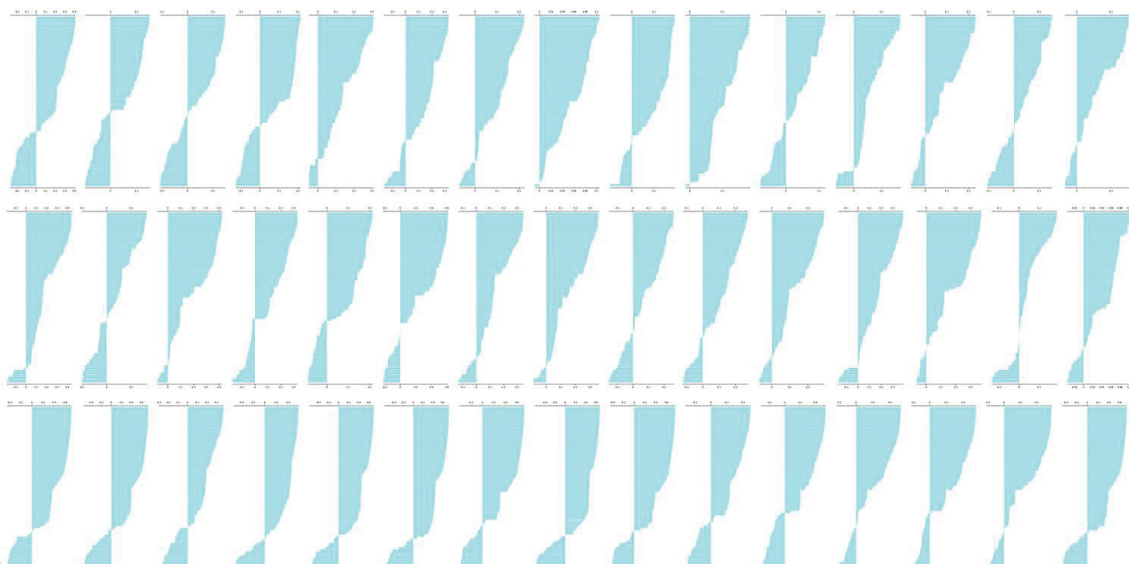




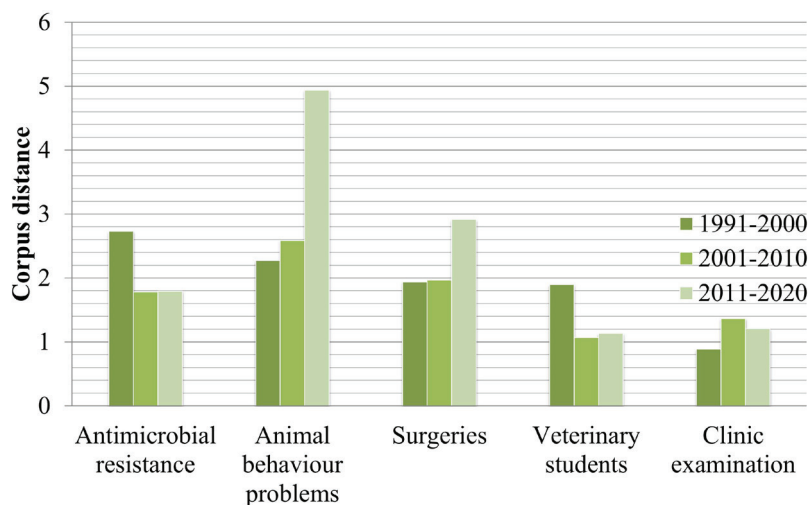
**Figure 3.** Topics extracted from the period of 2011-2020 with the LDA analysis. The combination of entropy and burstiness measures was used as an index to quantify the most frequent topics of literature (red color)



**Figure 4.** Perplexity is a measure of how well the LDA model predicted a sample. A lower perplexity score indicates better generalization performance of the model



**Figure 5.** Silhouette graphs showing the consistency and the cohesion of 45 different combinations of training and prediction sets. Small ‘between’ and ‘within’ dissimilarities between training and prediction sets, give a small silhouette score  $[s(i)]$  close to zero and would imply that the model is accurate. The overall average width for each plot is calculated with the Euclidean distance to give the average  $s(i)$



**Figure 6.** Comparison of selected topics described in terms of their context. Distance of topics from the corpus of Veterinary medicine was measured with the Kullback-Leibler divergence distribution. Lower distance from a corpus indicates that a topic is closely related with Veterinary medicine

## DISCUSSION

The interpretation of topics is not an easy mission due to the generality of some topics or because many of their words do not match in meaning. The extraction of ambiguous topics has been reported in the past (Nanni, 2017). However, the LDA model managed to uncover a satisfying number of topics, many of which are well interpretable. The term “veterinari” was included in a number of topics during each decade, showing that these topics were not able to acquire a special meaning. Apart from these strongly linked and sometimes overgeneralized topics, other more specific and more meaningful topics existed. The basic analysis in extracting the meaning of a document by giving us a number of topics is the LDA modeling. The interpretation process is subjective and shall focus on the more meaningful topics. To avoid this subjectivity the indexes of topic diagnostic information, document entropy, document burstiness and corpus distance, will constitute a more quantitative approach which helps us objectively interpret results in contrast with the direct subjective human interpretation. They can provide a useful automated summary of topic quality (Boyd-Graber et al., 2014).

LDA analysis constitutes the main technique applied in the present study to extract the major scientific topics of Veterinary literature. At first, the three measures of topic diagnostic information are discussed combined to analyze the relation between persistently present topics during the three decades and the main corpus of literature. Secondly, an effort is made to identify major shifts occurred in the overall topic rank. For this purpose, a simplification is used to shortly

describe one topic (Supplementary material). For instance the topics “student|veterinari|teach|educ|cours” and “veterinari|student|medicin|educ|program” are both described as “students” and the topics “resist|antibiot|antimicrobi|infect|strain” and “resist|antimicrobi|isol|strain|antibiot” are both described as “antibiotic resistance” because they refer to students and antibiotics resistance respectively. Thirdly, trends of selected words representing topics are discussed with the use of word frequencies (collocation statistics - Supplementary material) to further facilitate the interpretation of topics under a more analytical aspect.

The use of document entropy and document burstiness helped to identify and rank the most frequent topics. A focus on two animals (dog and cat) during the first time period, a focus on students’ educational programs during the second period and a focus on students surveys during the third period indicate different hot topics of each period. Simultaneously, student-related topics were less distinct from the corpus during the last two periods. This means that veterinary student topics became more frequent and were progressively incorporated in the main body of scientific literature. It is possible that a shift in scientific thinking had occurred during the 30-year-period. Students’ attitude, learning, motivation, competence in science, learning in practice are some of the aspects of the use of this term in literature (Mich et al., 2010; Jones et al., 2019). A relevant incorporation of a topic into the main corpus of Veterinary Medicine occurs with the topic of antibiotic resistance. The latter is less frequent during the first period and becomes a hot topic less distanced from the corpus during the past two pe-

riods. It is possible that scientists working on antibiotics take time to find novel solutions to this problem or they are aware of microorganisms' resistance to drugs (Toutain et al., 2016) thus the reporting of problems is progressively accumulated in literature.

The topic of animal behavior is extracted in all the three periods from the LDA model but it does not belong to the most frequent ones. This probably shows that veterinarians have already incorporated a perception that behavior problems are equally important with others topics such as clinical examination even before the nineties. The large distance from the corpus during 2011-2020 indicates that during the past few years this topic is linked to several words other than the most common ones. These different words probably reflect a variety of new subjects introduced in larger quantities into scientific community which were not strongly present in the past. Indeed, abstracts collected during the present study, contained the term behavior in the context of behavior of aggressive dogs (Csoltova et al., 2017), changes of behavior under a specific therapy (Packer et al., 2016), behavior studies, behavior alterations or behavior abnormalities (Tynes and Sinn, 2014).

A diachronically increasing interest of scientific literature for the topics of students and antimicrobial resistance and a decreasing interest for the topics of surgeries and animal behavior have been observed (Fig. 1-3 and Supplementary Material). On the contrary, the topic of clinical examination constantly seems to be of high interest among publications as it shows up in the five most frequent topics during the 30-years-period. During the 1991 - 2000 period the most frequent topics included clinical examination, were followed by topics including surgeries, students, antimicrobial resistance and animal behavior (Supplementary material). During the 2001 - 2010 period a shift was noticed in topic frequencies rank: the most frequent topic included students while clinical examination, antimicrobial resistance, surgeries and animal behavior followed. Lastly, during the 2011 - 2020 period the most frequent topics included students followed by topics including clinical examination, antimicrobial resistance, surgeries and behavior. The frequencies of these topics show their relative position in the overall ranking of each decade and probably reflect similar shifts in the interests of each decade. Furthermore, during the last period topics including clinical examination, antimicrobial resistance and surgeries were extracted in two different versions. It is possible

that a broadening occurred regarding their discussion in the scientific literature. Other topics emerge during specific time periods such as tumors during 2001 - 2020, vaccinations during 2001 - 2010, dosage (dose, mg/kg, concentration) during 2001 - 2020, dairy (milk, cow) during 1991 - 2000 and 2011 - 2020 and blood pressure during 1991 - 2010. It is difficult to tell whether there is a specific incident that provoked these shifts in topic ranking (for instance a pandemic) or if advances in technology have promoted the interest of scientists (for instance new tumor confrontation techniques). It is possible that a reporting increase in the national veterinary registration systems of each country acts as a signal to activate further scientific research of specific diseases.

Collocation statistics retrieved with KH coder (Supplementary material) contributed to the objective interpretation of selected words representing topics. Four of these words, those of "behavior", "resistance", "student" and "surgery", were selected as representatives of the topics that were firmly extracted during the three-decade-period and displayed changes into the topic ranking of each decade (animal behavior, antimicrobial resistance, students and surgeries respectively). Words appearing before or after these selected terms are directly depended on them (they display greater weight). It is more probable to see the word "veterinary" and "female" one word left from "student" and word "educ" one word right from "student" during 1991 - 2000. All these words combined thus indicate that they are connected under this point of view: veterinary student's education or female student's education. During 2001 - 2010 other words co-appear with the word student such as medicine, learning, graduate, interest, evaluation, training, participation or experience. All these words combined thus indicate that they are connected under the following point of view: veterinary student's learning, veterinary student's training etc. During 2011 - 2020 new words appear close to the term student: performance, perception, studies, assessment, training, attitude, experience or improvement. All these words combined indicate that they are connected under this point of view: veterinary student's performance or veterinary student's perception etc. It is possible that a shift in the scientific thinking occurred during 1991 - 2020 from the simple perspective of veterinary student's education (Heath et al., 1996) to a more competitive aspect including performance, assessment of their experience etc (Zenner et al., 2005).

It is more probable to encounter the words “veterinary”, “animal”, “invas”, “cardiovascular” one word left from “surgeri” and the words “procedur” and “depart” one word right from “surgeri” during 1991 - 2000. All these words combined thus suggest that they are connected under this point of view: veterinary surgeries or invasive surgeries or surgeries procedure. During 2001 - 2010 new words appear close to the word surgeries such as “dure”, “clinic”, “laparoscop”, “abdomin”, “convent”, “hors”, “colic” etc. All these words combined suggest that they are connected under this point of view: veterinary surgeries or laparoscopic surgeries or conventional surgeries etc. During 2011 - 2020 new words appear close to the word surgeries such as “dog”, “perform”, “course”, “spinal”, “open”. All these words combined suggest that they are connected under this point of view: surgeries duration or veterinary surgeries or surgeries performance or dog surgeries etc. It is possible that the scientific field of surgeries is constantly evolving and dealing with ever-changing topics.

It is more probable to see the words “problem”, “animal”, “therapy”, “pet”, “pharmacotherapy” close to the word “behavior” during 1991 - 2000. New words appear close to “behavior” during 2001 - 2010 those of “medicin”, “veterinary” and “cours”. During 2011 - 2020 words that appear close to “behavior” include “problem”, “change”, “relat”, “dog” and “intervent”. All these words combined suggest that they are connected under this point of view: animal behavior problem during 1991 - 2020, but with some differentiations between decades from behavior therapy/pharmacotherapy to behavior course and then to behavior change. It is possible that scientists make efforts to intervene into the animal behavior problems through therapies. However during 2001 - 2020 an effort is made to attribute biological interpretations into the aggressive behavior of animals. Some types of agitated behavior indeed have a strong genetic basis (Grandin and Deesing, 2014).

It is more probable to see the words “antibiot”, “bacteri” and “pathogen” during 1991 - 2000, the words “antimicrobi”, “isol” and “methicilin” during 2001 - 2010 and the words “antimicrobi”, “antibiot” and “multidrug” during 2011 - 2020 close to the term “resist”. All these words combined suggest that they are connected under this point of view: antibiotic resistance or antimicrobial resistance or isolates of bacteria to study the antibiotic resistance etc. Many of the terms closely present with “resist” are common

between the three decades. However it is possible that each decade is characterized by different priorities regarding research on antibiotic resistance as new terms stand out during the 2001 - 2020 period, those of multidrug, *Staphylococcus aureus* and gene nevertheless this is not so evident and cannot be distinctly supported by the specific results.

Previous text mining works have revealed a common number of terms also described in the present study. The subject of public awareness regarding the way in which farm animals are kept, the use of antimicrobials to increase animal performance and animal welfare have been reported before (Contiero et al., 2019). The subject of infectious diseases transmitted from animals to humans (or the opposite) and antimicrobial resistance due to prescriptions and results in human health have also been reported in the past (Lustgarten et al., 2020).

The validation of the LDA model carried out with several ways to verify that the accuracy in extracting the same topics was good at a satisfactory level. Each time the LDA model is applied a number of less frequent and hard to interpret topics is extracted (Nanni, 2017). We have to take into consideration that a part of the topics that LDA failed to correctly extract is these nonessential topics. On the other hand the LDA model is suitable in correctly identifying the most important of them over a large period of time. More specialized queries in the future may contribute to revealing possible trends of less studied topics of the scientific literature.

## CONCLUSIONS

LDA managed to reveal the most frequent topics of three continuous decades. The number of topics extracted during each period increases proportional to the volume of scientific literature. Differences throughout decades occur and may reflect perceptions of researchers. Topics related with veterinary students and antibiotic resistance are probably incorporated into the main corpus of literature during the 2001 - 2020 period while topics related with animal behavior were probably enriched with a variety of new subfields not recorded in the past. Quantitative literature research is an appropriate tool in identifying trends in topics.

## ACKNOWLEDGEMENTS

The authors would like to thank the two anonymous reviewers for their useful comments.



## SUPPLEMENTARY MATERIAL

Collocation statistics retrieved with KH coder. Word frequencies appearing before (L) and after (R) a node word are presented. The list shows the words that directly depend on the node word. Words that appear closer to the node words display greater weight,

thus a higher Score. Four examples of strongly connected words with “Veterinary”, persistently present in all the 1991 - 2020 period are shown here those of “behavior”, “resistance”, “student” and “surgery”. [A = 1991 - 2000, B = 2001 - 2010, C = 2011 - 2020]

Collocation Stats

Node Word

Word: behavior

POS:

Conj.:

Hits: 49

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	problem	ALL	17	4	13	1	0	1	0	2	12	0	0	0	1	14.733
2	anim	ALL	12	12	0	1	0	1	0	10	0	0	0	0	0	10.533
3	therapi	ALL	9	1	8	0	0	0	1	0	8	0	0	0	0	8.500
4	pet	ALL	6	5	1	0	0	0	0	5	0	0	0	0	1	5.200
5	pharmacotherapi	ALL	4	0	4	0	0	0	0	4	0	0	0	0	0	4.000
6	medicin	ALL	6	3	3	2	1	0	0	3	0	0	0	0	0	3.650
7	aggress	ALL	4	2	2	0	0	0	0	2	1	0	0	1	0	3.250
8	hors	ALL	3	2	1	0	0	0	0	1	0	0	0	0	0	2.200
9	confront	ALL	2	2	0	0	0	0	0	2	0	0	0	0	0	2.000
10	modif	ALL	2	2	0	0	0	0	0	2	0	0	0	0	0	2.000
11	treat	ALL	2	2	0	0	0	0	0	2	0	0	0	0	0	2.000
12	relet	ALL	4	3	1	1	0	1	0	1	0	0	1	0	0	1.867
13	veterinari	ALL	4	3	1	1	0	1	0	0	0	0	1	0	0	1.733
14	clinic	ALL	4	3	1	1	0	2	0	0	1	0	0	0	0	1.700
15	elmin	ALL	3	2	1	0	0	1	0	0	0	0	0	0	1	1.533
16	field	ALL	3	1	2	0	0	0	0	1	2	0	0	0	0	1.500
17	practic	ALL	2	2	0	0	0	0	0	1	0	0	1	0	0	1.333
18	drug	ALL	5	2	3	0	0	2	0	0	0	0	0	1	2	1.317
19	cat	ALL	3	1	2	0	0	0	0	1	0	0	1	0	1	1.250
20	common	ALL	2	2	0	1	0	0	0	1	0	0	0	0	0	1.200

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: behavior

POS:

Conj.:

Hits: 108

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	anim	ALL	22	18	4	2	1	0	1	14	0	0	0	2	2	16.050
2	problem	ALL	12	2	10	1	1	0	0	9	1	0	0	0	0	9.950
3	medicin	ALL	12	1	11	0	1	0	0	9	0	0	0	0	2	9.650
4	veterinari	ALL	15	12	3	2	4	0	0	6	0	0	0	1	2	8.050
5	course	ALL	6	1	5	0	1	0	0	5	0	0	0	0	0	5.250
6	clinic	ALL	8	2	6	0	0	0	0	2	2	0	1	2	1	5.033
7	play	ALL	6	5	1	0	1	0	1	3	0	0	1	0	0	4.083
8	sexual	ALL	4	4	0	0	0	0	0	4	0	0	0	0	0	4.000
9	abnorm	ALL	6	4	2	0	0	0	3	1	0	1	0	0	0	3.833
10	chang	ALL	4	1	3	1	0	0	0	3	0	0	0	0	0	3.200
11	biolog	ALL	5	4	1	0	1	0	2	0	1	0	0	0	0	3.083
12	behavior	ALL	12	6	6	2	3	1	0	0	0	1	3	2	0	2.967
13	aggress	ALL	6	3	3	1	0	0	0	2	0	0	2	1	0	2.900
14	reproduct	ALL	5	4	1	0	1	0	2	1	0	0	1	0	0	2.583
15	profession	ALL	4	3	1	1	0	1	0	1	0	0	0	0	0	2.533
16	identifi	ALL	5	4	1	0	2	0	1	1	0	0	0	0	0	2.500
17	canin	ALL	6	2	4	1	0	0	0	1	0	1	0	1	2	2.350
18	felin	ALL	3	3	0	0	0	1	0	2	0	0	0	0	0	2.333
19	medic	ALL	4	2	2	0	1	0	1	1	0	0	0	0	0	2.250
20	servic	ALL	3	1	2	0	1	0	0	2	0	0	0	0	0	2.250

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: behavior

POS:

Conj.:

Hits: 254

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	problem	ALL	47	7	40	4	1	2	0	39	0	0	0	0	1	40.917
2	chang	ALL	32	8	24	1	0	3	2	22	0	1	0	0	1	26.783
3	relet	ALL	16	9	7	0	1	1	3	4	2	1	0	3	0	9.517
4	dog	ALL	26	13	13	6	2	0	2	3	0	0	0	5	3	9.117
5	intervent	ALL	10	2	8	1	0	0	0	1	6	1	0	1	0	7.950
6	treatment	ALL	19	10	9	0	3	4	3	0	2	0	2	4	1	7.450
7	anim	ALL	15	9	6	4	1	0	1	3	1	1	1	1	2	7.033
8	cat	ALL	18	6	12	1	1	1	2	0	1	5	4	2	0	6.850
9	clinic	ALL	10	5	5	1	2	0	1	4	1	0	0	0	0	6.700
10	felin	ALL	6	6	0	0	0	0	0	6	0	0	0	0	0	6.000
11	effect	ALL	11	6	5	1	0	2	2	1	2	0	3	0	0	5.867
12	studi	ALL	17	9	8	4	2	2	0	1	0	1	0	2	3	5.733
13	veterinari	ALL	17	4	13	2	0	1	1	0	0	3	3	4	3	5.333
14	medic	ALL	14	9	5	2	2	3	0	1	0	2	1	0	1	5.183
15	alter	ALL	5	2	3	0	0	0	0	2	3	0	0	0	0	5.000
16	biolog	ALL	5	5	0	0	0	0	0	5	0	0	0	0	0	5.000
17	depress	ALL	6	3	3	0	0	0	0	2	1	3	0	0	0	5.000
18	undesir	ALL	5	5	0	0	0	0	0	5	0	0	0	0	0	5.000
19	pain	ALL	13	10	3	0	3	1	6	0	0	0	0	2	1	4.950
20	report	ALL	10	5	5	1	1	0	2	0	3	0	2	0	0	4.783

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: resist

POS:

Conj.:

Hits: 147

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	antibiot	ALL	30	23	7	0	3	0	1	19	0	2	2	2	2	22.567
2	bacteri	ALL	12	10	2	0	0	0	1	9	0	1	0	0	0	10.250
3	pathogen	ALL	18	3	15	1	1	0	0	1	3	2	5	3	0	7.967
4	bacteria	ALL	11	3	8	1	0	0	1	4	3	1	0	0	0	7.533
5	determin	ALL	10	4	6	1	0	1	2	4	0	0	1	1	0	5.983
6	develop	ALL	11	8	3	0	1	2	5	0	2	0	0	0	1	5.617
7	monitor	ALL	6	1	5	0	0	0	1	0	0	0	0	0	0	5.500
8	test	ALL	8	4	4	0	0	2	0	4	0	0	0	0	0	5.500
9	drug	ALL	6	5	1	0	0	0	0	5	0	0	1	0	0	5.333
10	increas	ALL	13	11	2	3	0	5	2	1	0	1	0	1	0	5.017
11	antimicrobi	ALL	6	5	1	0	1	0	4	0	0	1	0	0	0	4.667
12	veterinari	ALL	14	4	10	1	1	2	0	0	3	4	1	2	0	4.600
13	select	ALL	7	2	5	0	1	0	1	0	3	0	0	0	0	4.250
14	anim	ALL	12	3	9	0	1	1	0	0	3	3	1	2	0	4.233
15	strain	ALL	6	2	4	0	1	0	0	1	2	1	1	0	0	4.083
16	resist	ALL	12	6	6	3	0	2	1	0	1	2	0	3	0	3.533
17	pattern	ALL	4	1	3	0	0	1	0	3	0	0	0	0	0	3.333
18	percentag	ALL	5	5	0	0	0	0	4	1	0	0	0	0	0	3.000
19	medicin	ALL	8	3	5	1	0	0	0	0	3	2	0	0	0	2.700
20	mechan	ALL	5	4	1	0	1	1	0	1	0	0	0	0	0	2.583

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: resist

POS:

Conj.:

Hits: 611

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	antimicrobi	ALL	152	114	38	4	5	2	10	93	3	9	9	6	11	114.917
2	isol	ALL	92	38	54	3	9	3	19	4	17	4	11	12	10	45.017
3	antibiot	ALL	50	35	15	1	2	1	1	30	0	5	4	2	4	36.667
4	bacteria	ALL	56	7	49	1	2	2	0	24	15	4	2	4	0	36.500
5	methicillin	ALL	33	27	6	0	0	0	0	27	0	4	1	1	0	29.583
6	gene	ALL	37	7	30	2	3	1	1	0	26	0	0	2	2	28.883
7	drug	ALL	33	28	5	0	1	0	0	27	0	1	1	0	3	28.667
8	staphylococcus	ALL	29	2	27	1	0	1	0	24	1	0	2	0	0	25.533
9	strain	ALL	35	7	28	2	1	2	0	26	2	2	6	2	0	21.883
10	quionol	ALL	25	19	6	0	1	1	0	17	1	5	0	0	0	21.083
11	bacteri	ALL	24	13	11	0	0	0	2	11	3	2	3	2	1	17.000
12	resist	ALL	42	21	21	5	4	8	1	3	3	1	8	4	5	16.333
13	multidrug	ALL	18	17	1	0	0	0	1	15	0	0	0	0	0	16.000
14	level	ALL	26	18	8	0	3	0	2	6	0	0	0	2	0	15.000
15	eureus	ALL	29	2	27	1	0	1	0	0	24	1	0	0	2	13.267
16	fluoroquinol	ALL	33	16	13	2	0	0	0	11	0	1	0	1	0	12.433
17	coli	ALL	33	16	17	2	3	3	8	0	6	6	3	2	0	12.300
18	determin	ALL	23	9	14	2	2	1	4	0	7	1	1	4	1	12.267
19	medicin	ALL	20	7	13	2	2	0	0	1	1	0	2	0	0	12.050
20	develop	ALL	25	21	4	1	3	7	3	1	0	0	0	1	2	11.433

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: resist

POS:

Conj.:

Hits: 1307

Result

Collocation Stats

Node Word

Word: student POS: Conj.: Hits: 101

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	veterinari	ALL	36	24	1	2	2	3	15	0	4	3	4	1	22.100	
2	educ	ALL	9	4	5	1	2	0	0	4	0	1	0	0	5.450	
3	fearn	ALL	4	4	0	0	0	0	0	4	0	0	0	0	4.000	
4	learn	ALL	8	3	5	0	1	1	0	1	2	1	1	1	3.700	
5	medic	ALL	6	5	1	0	1	1	2	0	0	0	0	0	3.283	
6	enter	ALL	4	0	4	0	0	0	0	3	0	0	1	0	3.250	
7	facult	ALL	6	1	5	0	0	0	0	1	4	0	1	0	3.250	
8	clinic	ALL	6	1	5	1	0	0	0	2	1	1	0	1	3.233	
9	interest	ALL	4	1	3	0	0	0	0	1	1	1	1	0	2.833	
10	minor	ALL	4	3	1	0	0	1	0	2	0	0	0	0	2.667	
11	receiv	ALL	3	0	3	0	0	0	0	2	1	0	0	0	2.500	
12	course	ALL	8	4	4	0	1	2	1	0	0	1	2	1	2.450	
13	year	ALL	7	5	2	1	1	0	3	0	0	0	0	0	2.450	
14	colleg	ALL	7	2	5	1	0	0	0	0	1	2	1	2	2.233	
15	evalu	ALL	3	0	3	0	0	0	0	2	0	0	0	1	2.200	
16	underpres	ALL	3	3	0	1	0	0	0	2	0	0	0	0	2.200	
17	commen	ALL	2	0	2	0	0	0	0	2	0	0	0	0	2.000	
18	postgradu	ALL	2	2	0	0	0	0	0	2	0	0	0	0	2.000	
19	practic	ALL	6	2	4	0	0	2	0	0	1	1	1	1	1.950	
20	medicin	ALL	7	2	5	1	1	0	0	0	3	1	1	1	1.900	

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: student POS: Conj.: Hits: 723

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	veterinari	ALL	222	175	47	12	10	14	26	113	2	7	8	4	16	152.933
2	year	ALL	80	54	26	8	6	6	14	20	5	2	4	15	39.267	
3	medicin	ALL	51	32	19	7	4	6	11	0	0	13	3	3	23.083	
4	learn	ALL	46	13	33	3	4	4	1	1	10	7	3	8	5	21.933
5	medic	ALL	20	19	1	0	1	0	0	18	0	0	0	1	0	18.500
6	facult	ALL	42	19	23	3	2	4	9	1	1	1	1	2	4	17.400
7	graduat	ALL	30	16	14	0	3	4	0	9	6	3	3	2	1	16.233
8	center	ALL	18	14	4	2	1	0	1	0	14	0	0	0	0	15.150
9	student	ALL	44	22	22	7	7	4	2	2	2	2	4	7	7	14.967
10	evalu	ALL	30	19	11	2	2	6	7	4	2	2	2	1	1	14.767
11	interest	ALL	23	7	16	1	1	2	0	3	8	0	5	1	2	14.483
12	educ	ALL	37	25	12	3	8	6	6	2	1	4	1	4	2	14.333
13	teach	ALL	30	19	11	4	6	2	2	5	3	1	4	1	1	14.000
14	anim	ALL	42	18	25	4	4	1	0	5	6	8	6	3	3	13.983
15	provid	ALL	25	18	7	3	4	2	1	8	1	2	2	0	2	13.833
16	train	ALL	30	21	9	2	2	10	5	2	3	1	1	2	2	13.467
17	particip	ALL	19	7	12	0	1	1	1	4	4	1	3	0	0	12.167
18	exper	ALL	31	9	22	1	3	5	0	0	3	4	4	4	4	12.000
19	program	ALL	28	19	9	4	5	3	3	2	0	4	1	1	3	10.400
20	profession	ALL	16	13	3	0	0	4	3	6	0	1	2	0	0	10.000

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: student POS: Conj.: Hits: 1451

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	veterinari	ALL	534	435	99	21	28	23	99	264	4	34	14	18	29	368.333
2	year	ALL	167	126	41	17	10	13	14	44	2	3	9	19	10	85.483
3	medicin	ALL	120	80	40	11	6	11	16	36	0	1	25	4	10	63.200
4	medic	ALL	63	57	6	1	2	1	4	18	1	2	1	1	1	53.650
5	perform	ALL	77	29	48	3	7	8	6	5	25	12	3	3	3	47.033
6	student	ALL	116	58	58	15	12	15	3	3	3	5	12	13	15	41.500
7	learn	ALL	69	21	48	8	4	4	1	4	23	4	5	12	4	38.900
8	percept	ALL	48	7	41	0	4	1	1	1	27	2	5	2	1	34.150
9	studi	ALL	82	49	33	13	12	15	8	1	4	8	5	0	0	28.967
10	asses	ALL	52	31	21	8	4	6	4	9	5	2	9	1	4	25.650
11	course	ALL	65	39	26	5	13	7	8	6	0	2	7	9	8	23.767
12	train	ALL	53	30	23	3	5	6	5	4	5	2	9	2	6	23.217
13	attitud	ALL	37	14	23	2	4	6	2	0	16	1	5	1	0	22.767
14	anim	ALL	72	31	41	4	6	7	10	0	1	1	13	17	9	22.417
15	univers	ALL	58	26	32	5	4	2	7	0	1	8	15	5	12	21.483
16	educ	ALL	50	30	20	7	4	11	4	4	4	3	5	3	5	20.983
17	exper	ALL	49	14	35	4	4	1	3	2	7	4	4	10	0	20.467
18	group	ALL	53	30	23	3	7	9	2	1	5	6	6	5	1	19.850
19	improv	ALL	31	18	13	2	1	1	3	1	1	1	3	1	3	18.417
20	result	ALL	48	29	19	5	8	5	3	3	1	5	6	4	3	18.267

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: burger POS: Conj.: Hits: 102

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	veterinari	ALL	23	15	8	2	3	3	7	0	5	2	0	1	12.517	
2	medicin	ALL	20	11	9	1	3	3	1	0	6	2	1	7	7.150	
3	anim	ALL	8	8	0	2	0	1	5	0	1	0	0	0	6.833	
4	perform	ALL	8	4	4	0	2	0	1	1	2	1	0	0	4.333	
5	invas	ALL	4	4	0	0	0	0	0	4	0	0	0	0	4.000	
6	procedur	ALL	7	4	3	2	1	1	0	0	2	1	0	0	3.233	
7	diseas	ALL	4	2	2	0	0	0	2	1	0	0	1	0	3.200	
8	depart	ALL	5	3	2	0	0	1	0	2	0	0	0	0	3.167	
9	cardiovascular	ALL	3	3	0	0	0	0	0	0	3	0	0	0	3.000	
10	clean	ALL	4	4	0	0	0	0	2	2	0	0	0	0	3.000	
11	dure	ALL	3	2	1	0	0	0	0	2	1	0	0	0	3.000	
12	minor	ALL	3	3	0	0	0	0	0	3	0	0	0	0	3.000	
13	facult	ALL	4	0	4	0	0	0	0	2	0	1	0	1	2.583	
14	oncolog	ALL	3	3	0	0	0	0	1	2	0	0	0	0	2.500	
15	laser	ALL	6	4	2	0	2	0	1	1	0	0	1	1	2.450	
16	year	ALL	5	3	2	0	0	0	2	1	0	0	0	2	2.400	
17	reconstruct	ALL	3	3	0	0	0	0	1	2	0	0	0	0	2.333	
18	german	ALL	3	2	1	0	0	0	0	2	0	0	1	0	2.250	
19	spare	ALL	3	2	1	0	0	0	0	2	0	0	1	0	2.250	
20	includ	ALL	4	2	2	0	0	0	2	0	1	0	0	1	2.200	

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: burger POS: Conj.: Hits: 184

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	veterinari	ALL	33	16	17	1	2	4	3	6	1	3	4	3	6	15.317
2	medicin	ALL	27	17	10	2	0	2	11	2	0	0	3	4	3	11.167
3	dure	ALL	16	14	2	2	2	2	2	6	0	0	0	1	1	9.017
4	clinic	ALL	17	10	7	0	1	5	4	0	3	0	2	1	1	8.033
5	anim	ALL	17	10	7	1	1	4	1	3	0	0	3	2	2	7.183
6	laparoscop	ALL	7	7	0	0	0	0	0	7	0	0	0	0	0	7.000
7	invas	ALL	7	7	0	0	0	1	0	6	0	0	0	0	0	6.333
8	depart	ALL	10	7	3	1	1	0	5	0	2	0	1	0	1	5.200
9	abdomin	ALL	8	5	3	0	0	0	0	4	0	0	1	0	2	5.100
10	convent	ALL	5	5	0	0	0	1	0	4	0	0	0	0	0	4.333
11	hors	ALL	13	5	8	0	3	1	0	1	0	1	0	6	4	4.283
12	equin	ALL	9	8	1	1	0	1	0	6	0	1	0	0	0	4.033
13	synthetich	ALL	6	4	2	0	0	0	3	0	1	0	1	0	0	4.033
14	colic	ALL	4	4	0	0	0	0	0	4	0	0	0	0	0	4.000
15	facult	ALL	7	0	7	0	0	0	0	2	2	3	0	0	0	4.000
16	perform	ALL	6	2	4	0	0	0	1	2	0	0	1	1	1	3.950
17	patient	ALL	11	5	6	1	3	1	0	0	0	3	3	0	0	3.783
18	includ	ALL	7	3	4	1	1	0	0	1	1	1	0	1	1	3.400
19	minim	ALL	7	7	0	0	1	0	6	0	0	0	0	0	0	3.250
20	lectat	ALL	4	2	2	0	0	2	0	2	0	0	0	0	0	3.200

Copy Filter Sort: The Score Window span: LS RS

Collocation Stats

Node Word

Word: burger POS: Conj.: Hits: 432

Result

N	Word	POS	Total	LT	RT	L5	L4	L3	L2	L1	R1	R2	R3	R4	R5	The Score
1	dure	ALL	35	32	2	6	1	9	15	2	0	0	0	0	0	23.733
2	veterinari	ALL	43	23	20	4	7	5	1	6	1	5	7	4	3	18.150
3	anim	ALL	33	20	13											

List of the topics extracted during 1991 - 2020 with the use of the LDA model. Topics are presented in descending order with the use of entropy and burstiness as indexes to measure the most common words and topics of literature. Abbreviations of selected topics were used to help track their order during each time period.

1991 - 2000			2001 - 2010			2011 - 2020				
1	dog/cat/case/clinic/examin	CE	1	veterinari/student/medicin/educ/program	ST	1	student/veterinari/medicin/studij/survey	ST	31	imag/radiolog/radiograph/tomographi
2	studi/data/age/veterinari/sampl		2	clinic/veterinari/case/veterinari/medicin	CE	2	anim/veterinari/medicin/health/human	CE	32	ultrasound/arteri/pregnanc/fetal/doppler
3	veterinari/anim/medicin/veterinari/develop		3	veterinari/medicin/use/practic/inform	CE	3	diet/group/feed/supplement/feed	CE	33	extract/concentri/ethanol/aqueous/quail
4	hors/anthelmint/modifi/analysis/structur		4	veterinari/human/medicin/practic/diseas	CE	4	veterinari/medicin/diseas/human/patient	CE	34	surgic/surgic/procedure/perfor/complic
5	techniq/surgic/medicin/veterinari/surgic	SU	5	resist(antimicrobijs)/col/strain/antibiot	AR	5	case/dog/veterinari/clinic/report	CE	35	heart/cardiac/thoracic/pulmonari/lung
6	measur/method/differ/bone/plasma		6	anim/health/veterinari/diseas/develop	CE	6	resist(iso)/gene/strain/antimicrobi	AR	36	platelet/tg(activ)/coagul/blood
7	student/veterinari/teach/educ/cours	ST	7	infect/vaccin/diseas/pr/mrsa	VA	7	studi(dog)/use/result/medicin	CE	37	tumor/cancer/cell/tumour/lmphoma
8	test/valu/method/procedur/diagnost		8	cell/tumor/tumour/dog/carcinoma	TU	8	resist(antimicrobi)/infect(iso)/bacteri	AR	38	joint/lame/articular/ruptr/osteoarthr
9	veterinari/medicin/school/year/found		9	hors/surgic/case/fractur/surgic	SU	9	cell/stem/tissu/mescl/mes	CE	39	thyroid/test/blade/perfus/ablat
10	concentr(studj/group)/serum/milk	DA	10	dose/administr/concentr/oral/mg/kg	DO	10	milk/cow/product/digest/dri	DA	40	thyruley/carolina/block/north
11	renal/function/failur/transplant/increas		11	dog/blood/method/sampl/concentr	CE	11	veterinari/review/report/studi/clinic	CE	41	und(er)/de/de/sles
12	resist(antibiot/antimicrobi)/infect/strain	AR	12	cell/immun/respons/express/activ	CE	12	abdomin/case/examin/surgic/ultrasound	SU	42	bird/avian/cell/tumour/lmphomare
13	drug/veterinari/medicin/use/effect		13	imag/use/techniq/digit/ultrasound	CE	13	drug/effect/activ/antifuse	CE	43	eye/ocular/corneal/strain/cornea
14	behavior/treat/reduct(bromid)/aggress	AB	14	studi(dog)/group/veterinari/result	CE	14	anim/parasit/infect/speci/case	DA	44	fish/aquacultur/aquat/catfish/salmon
15	system/imag/use/comput/radiat		15	blood/arteri/oxy/genin/pressur/use	BP	15	express/activ/protein/increas/level	CE	45	macaqu/thesus/otter/cryptococcus/eleph
16	day/activ/effect/level/agent		16	dog/heart/breed/system/wave	CE	16	treatment/hors/dog/clinic/case	CE	46	burnout/capon/dystonia/martien/movement
17	blood/pressur/measur/use/arteri	BP	17	veterinari/medicin/facultj/school/cper	CE	17	veterinari/univers/medicin/parasitolog/research	CE	47	larva/collar/tick/ear/teat
18	test/method/antibodi/infect/diseas		18	product/feed/goat/meat/extract	CE	18	dog/cat/studi/age/year	CE	48	epileps/seizur/epilept/behavior/noninferior
19	muscl/state/cow/liver/bird		19	renal/creatinin/kidney/serum/control	CE	19	group/anesthesia/effect/anesthet/pain	CE	49	cup/night/bat/behaviour/vns
20	plant/tradit/diseas/medicin/vaccin	VA	20	case/rabireport/number/anim	CE	20	product/anim/feed/effect/associ	CE	50	iron/hepatic/liver/lg/binder
21	mini/group/vaccin/report/effect	VA	21	pet(anim)/behavior/problem/dog	AB	21	dose/concentr/mg/kg/administr/drug	DO	51	sorghum/cenit/kafirini/starch/outlet
22	bacteri/change/calculi/bacteriolog/specimen		22	cow/interview/milk/dairi/herd	DA	22	spinal/neurolog/brain/cord/nerve	CE	52	grade/inpl/metronidazol/fdm/gastric
			23	diabet/glucos/allergen/mellitus/colt	CE	23	antimicrobi/fam/use/pig/antibiot	CE	53	brows/bag/indigen/tree/fluid
			24	bone/fractur/limb/month/inateri	CE	24	bone/fractur/limb/month/inateri	CE	54	pdp/ivas/chimpanzee/anxiety/criteria
			25	pcr(da)/sequenc/detect/genet	CE	25	pcr(da)/sequenc/detect/genet	CE	55	stk(cnt)/disbud/horn/hot
			26	infect/test/protein/serum/antibodi	CE	26	infect/test/protein/serum/antibodi	CE	56	refractomet/usg/urini/furosemid/ped
			27	blood/serum/concentr/biochem/protein	CE	27	blood/serum/concentr/biochem/protein	CE	57	biomark/cardiac/ctn/assay/myocard
			28	vaccin/peptid/immun/respons/adjuv	VA	28	vaccin/peptid/immun/respons/adjuv	VA	58	empathi/rural/shortage/constru/stom
			29	diseas/infect/virus/rabi/vaccin	VA	29	diseas/infect/virus/rabi/vaccin	VA	59	pfa/probiot/lactobacillus/probi/caeca
			30	wound/heal/skin/laser/honey	BP	30	wound/heal/skin/laser/honey	BP	60	bone/miner/bmd/densit/dxa
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25			25				
			26			26				
			27			27				
			28			28				
			29			29				
			30			30				
			25		</					



## REFERENCES

- Anholt RM, Berezowski J, Jamal I, Ribble C, Stephen C (2014) Mining free-text medical records for companion animal enteric syndrome surveillance. *Prev Vet Med* 113(4): 417-422.
- Blei DM, Ng AY, Jordan MI (2003) Latent dirichlet allocation. *J Mach Learn Res* 3(Jan): 993-1022.
- Bollig N, Clarke L, Elsmo E, Craven M (2020). Machine learning for syndromic surveillance using veterinary necropsy reports. *PloS one* 15(2): p.e0228105.
- Boyd-Graber J, Mimno D, Newman D (2014) Care and feeding of topic models: Problems, diagnostics, and improvements. In: *Handbook of mixed membership models and their applications* Chapman and Hall/CRC: pp 225-255.
- Christopher MM, Marusic A (2013) Geographic trends in research output and citations in veterinary medicine: insight into global research capacity, species specialization, and interdisciplinary relationships. *BMC Vet Res* 9(1): 115.
- Contiero B, Cozzi G, Karpf L, Gottardo F (2019) Pain in Pig Production: Text Mining Analysis of the Scientific Literature. *J Agric Environ Ethics* 32: 401-412.
- Csoltova E, Martineau M, Boissy A, Gilbert C (2017) Behavioral and physiological reactions in dogs to a veterinary examination: Owner-dog interactions improve canine well-being. *Physiol Behav* 177: 270-281.
- Demsar J, Curk T, Erjavec A, Gorup C, Hocevar T, Milutinovic M, Možina M, Polajnar M, Toplak M, Staric A, Stajdohar M, Umek L, Zagar L, Zbontar J, Zitnik M, Zupan B (2013) Orange: Data Mining Toolbox in Python. *J Mach Learn Res* 14(Aug): 2349-2353.
- Dórea FC, Lindberg A, Elvander M (2015) Veterinary syndromic surveillance in practice: costs and benefits for governmental organizations. *Infect Ecol Epidemiol* 5(1): 29973.
- Ellis J, Ellis B, Velez-Estevéz A, Reichel M, Cobo M (2020) 30 years of parasitology research analysed by text mining. *Parasitology* 1: 1-58.
- Grandin T, Deesing MJ (2014) Genetics and behavior during handling, restraint, and herding. In: *Genetics and the behavior of domestic animals* 2nd ed, Academic Press, New York: pp 115-158.
- Heath TJ, Lanyon A, Lynch-Blosse M (1996) A longitudinal study of veterinary students and recent graduates: 3. Perceptions of veterinary education. *Aust Vet J* 74(4): 301-304.
- Higuchi K (2016) KH Coder 3 reference manual. Kyoto (Japan): Ritsumeikan University.
- Jones BD, Byrnes MK, Jones MW (2019) Validation of the MUSIC Model of Academic Motivation Inventory: Evidence for use with veterinary medicine students. *Front Vet Sci* 6: 11.
- Jones-Diette JS, Dean RS, Cobb M, Brennan ML (2019) Validation of text-mining and content analysis techniques using data collected from veterinary practice management software systems in the UK. *Prev Veterinary Med* 167: 61-67.
- Kullback S, Leibler RA (1951) On information and sufficiency. *Ann Math Stat* 22(1): 79-86.
- Lin H, Sheu PC, Tsai JJP, Charles CN, Wang, CYC (2020) Text mining in a literature review of urothelial cancer using topic model. *BMC Cancer* 20: 462.
- Lustgarten JL, Zehnder A, Shipman W, Gancher E, Webb TL (2020) Veterinary informatics: forging the future between veterinary medicine, human medicine, and One Health initiatives-a joint paper by the Association of Veterinary Informatics (AVI) and the CTSA One Health Alliance (COHA). *JAMIA Open* 3(2): 306-317.
- Mich PM, Hellyer PW, Kogan L, Schoenfeld-Tacher R (2010) Effects of a pilot training program on veterinary students' pain knowledge, attitude, and assessment skills. *J Vet Med Educ* 37(4): 358-368.
- Mierswa I, Wurst M, Klinkenberg R, Scholz M, Euler T (2006). Yale: Rapid prototyping for complex data mining tasks. In: *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*: pp. 935-940.
- Nanni F (2017) The Web as a Historical Corpus: Collecting, Analysing and Selecting Sources on the Recent Past of Academic Institutions. (Doctoral dissertation, alma). *Universita di Bologna* 242p.
- Opitz D, Maclin R (1999) Popular ensemble methods: An empirical study. *J Artif Intell Res* 11: 169-198.
- Packer RM, Law TH, Davies E, Zanghi B, Pan Y, Volk HA (2016) Effects of a ketogenic diet on ADHD-like behavior in dogs with idiopathic epilepsy. *Epilepsy Behav* 55: 62-68.
- Queirós A, Faria D, Almeida F (2017) Strengths and limitations of qualitative and quantitative research methods. *Eur J Educ Stud* 9: 369-387.
- Rousseeuw PJ (1987) Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *J ComputAppl Math* 20: 53-65.
- Sahadevan S, Hofmann-Apitius M, Schellander K, Tesfaye D, Fluck J, Friedrich CM (2012) Text mining in livestock animal science: introducing the potential of text mining to animal sciences. *J Anim Sci* 90(10): 3666-3676.
- Shannon CE (1948) A mathematical theory of communication. *The Bell system technical journal* 27(3): 379-423.
- Toutain PL, Ferran AA, Bousquet-Melou A, Pelligand L, Lees P (2016) Veterinary medicine needs new green antimicrobial drugs. *Front Microbiol* 7: 1196.
- Tynes VV, Sinn L (2014) Abnormal repetitive behaviors in dogs and cats: a guide for practitioners. *Vet Clin North Am Small Anim Pract* 44(3): 543-564.
- Van der Waal K, Morrison RB, Neuhauser C, Vilalta C, Perez AM (2017) Translating big data into smart data for veterinary epidemiology. *Front Vet Sci* 4: 110.
- Welsh CE, Parkin TDH, Marshall JF (2017) Use of large-scale veterinary data for the investigation of antimicrobial prescribing practices in equine medicine. *Equine Vet J* 49(4): 425-432.
- Zenner D, Burns GA, Ruby KL, De Bowes RM, Stoll SK (2005) Veterinary students as elite performers: preliminary insights. *J Vet Med Educ* 32(2), 242-248.
- Zhang Y, Tao J, Wang J, Ding L, Ding C, Li Y, Qichao Zhou, Dunhai Li, Zhang H (2019) Trends in diatom research since 1991 based on topic modeling. *Microorganisms* 7(8): 213.