

LipiDisease: associate lipids to diseases using literature mining

Piyush More^{1,2,*}, Laura Bindila³, Philipp Wild⁴, Miguel Andrade-Navarro², Jean-Fred Fontaine²

¹Department of Pharmacology, University Medical Center, 55131 Mainz, Germany

² Faculty of Biology, Johannes Gutenberg University of Mainz, 55128 Mainz, Germany

³Clinical Lipidomics Unit, Institute of Physiological Chemistry, University Medical Center, 55131 Mainz, Germany

⁴Center for Thrombosis and Hemostasis (CTH), University Medical Center, 55131 Mainz, Germany

* Author to whom correspondence should be addressed

Supplementary Information

Computational methods

The data retrieved from PubChem and PubMed was filtered to limit PubMed records corresponding to at least one disease and one lipid. Among these, the significant associations of diseases and lipids were identified by performing one-tailed Fisher's exact test as described before (Fontaine and Andrade-Navarro, 2016). In short, articles associated with diseases as well as articles associated with lipids were represented in a 2-by-2 contingency matrix to identify the over-representation of the articles associated with both diseases and lipids. All computations were limited to the resulting 709,038 PubMed records associated with 4,488 diseases and 2,771 lipids.

To evaluate the enrichment of a set of lipids to a particular disease, one-tailed Fisher's exact test was performed by creating the contingency matrix depicted in Table S1.

Table S1. The contingency matrix used for performing a one-tailed Fisher's exact test. S represents a lipid set under investigation and D represents a disease. The number of lipids in the groups a , b , c , and d were used to identify the over-representation.

| Lipids | From the S | Not from the S |
|-----------------------------|--------------|------------------|
| Associated with the D | a | b |
| Not associated with the D | c | d |

The resulting p-value was corrected for multiple testing by calculating the false discovery rate (FDR) by Benjamini and Hochberg method (Benjamini and Hochberg, 1995) using the R statistical environment (R Core Team, 2020).

Lipid-level statistics were also considered for disease enrichment. A ranked list of lipids was used for disease enrichment using the fast gene-set enrichment algorithm. This was implemented using the fgsea package in R (Korotkevich *et al.*, 2019).

Database contents

The current version of LipiDisease consists of a total of 709,038 PubMed records associated with 4,488 diseases and 2,771 lipids. Sterol lipids, polyketides, and fatty acyls are the most represented categories by those articles (Fig. S1). Among diseases, 23% of the PubMed articles discuss neoplasms with lipids. Figure S1 shows the top 10 diseases associated with any lipids as described in the literature.

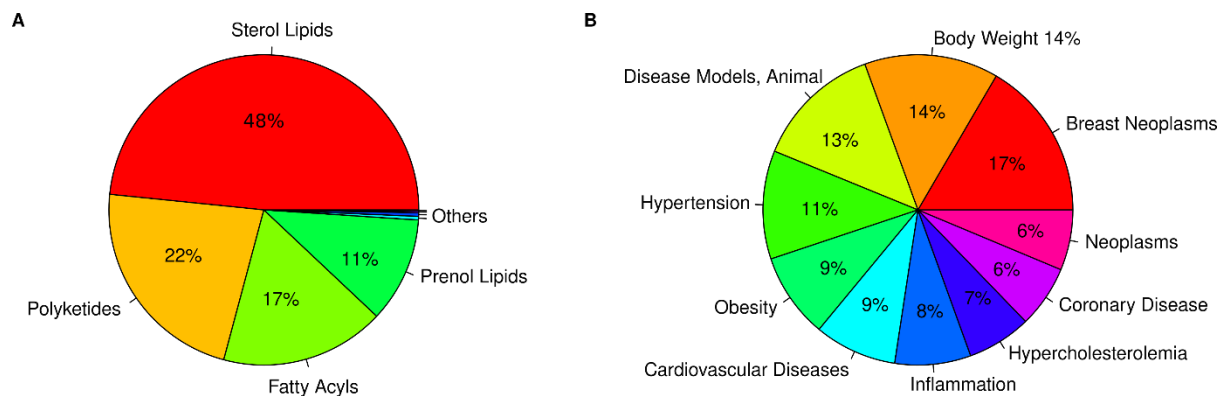


Fig. S1. Percentage of PubMed articles annotated with (A) different lipid classes and (B) diseases.

Comparison with LipidPedia

LipidPedia is one of the few tools specialized in associating lipids with biomedical information (Kuo and Tseng, 2018). It performs full-text mining to extract lipid-relevant information including the biomedical associations. However, it lacks statistical filtering and only provides diseases associated with individual lipids. Hence, we obtained lipids associated with coronary disease, type II diabetes mellitus, dyslipidemias, heart failure, and hyperlipidemias from LipidPedia and compared with the results of “Diseases to Lipids” analysis from LipiDisease performed at the default settings (minimum 5 citations associating a disease with a lipid, and 0.05 FDR). In general, LipidPedia reported high number of lipids associated with diseases compared to LipiDisease (Table S2). This could be because LipidPedia performs full-text mining to find these associations, while LipiDisease only considers MesH-curated citations. Furthermore, LipidPedia lacks statistical filtering making LipiDisease results more stringent. The output of our tool LipiDisease facilitates manual lookup of the PubMed articles from which associations are derived (Figure S2). This is not possible with LipidPedia.

Table S2. Comparison between LipiDisease and LipidPedia. The number of lipids associated with selected diseases obtained from both resources and overlap between them.

| Disease | LipiDisease | LipidPedia | Common |
|----------------------------|-------------|------------|--------|
| Coronary disease | 35 | 70 | 12 |
| Diabetes Mellitus, type II | 60 | 87 | 25 |
| Dyslipidemias | 13 | 7 | 2 |
| Heart failure | 21 | 36 | 4 |
| Hyperlipidemias | 26 | 73 | 10 |

A

| Lipid | Lipid_pubchem_cid | Disease | Citations | p-value | FDR |
|----------------------|-------------------------|----------------------------------|----------------------|----------|----------|
| Cholesterol | 5997 | Coronary_Disease | 8683 | 0.00 | 0.00 |
| Pravastatin | 54687 | Coronary_Disease | 534 | 0.00 | 0.00 |
| Epicholesterol | 5283629 | Coronary_Disease | 5466 | 0.00 | 0.00 |
| PGI2 | 5282411 | Coronary_Disease | 322 | 4.32e-73 | 1.24e-72 |
| TXA2 | 5280497 | Coronary_Disease | 194 | 1.19e-68 | 3.36e-68 |
| TXB2 | 5283137 | Coronary_Disease | 195 | 1.02e-42 | 2.53e-42 |
| 6-keto-PGF1alpha | 5280888 | Coronary_Disease | 118 | 2.56e-25 | 5.57e-25 |
| alpha-Linolenic acid | 5280934 | Coronary_Disease | 63 | 5.68e-20 | 1.17e-19 |
| EPA | 446284 | Coronary_Disease | 114 | 2.06e-16 | 4.11e-16 |
| Linolealidic acid | 5282457 | Coronary_Disease | 74 | 1.29e-11 | 2.41e-11 |

Showing 1 to 10 of 35 entries

Previous [1](#) [2](#) [3](#) [4](#) Next

B

Effects of disulfiram and pyridoxine on serum cholesterol.
1 Major LF, Goyer PF.
Cite Ann Intern Med. 1978 Jan;88(1):53-6. doi: 10.7326/0003-4819-88-1-53.
PMID: 619758
Share

Cholesterol management in patients hospitalized for coronary heart disease.
2 Boekeloo BO, Becker DM, LeBailly A, Pearson TA.
Cite Am J Prev Med. 1988 May-Jun;4(3):128-32.
PMID: 3395498
Share

[Impairment of cholesterol-acceptor function of high density lipoproteins in patients with ischemic heart disease].
3 Parfenova NS, Petrova-Maslakova LG, Kuznetsov AS, Ioffe DV, Alkhnis EG.
Cite Vopr Med Khim. 1988 Mar-Apr;34(2):42-6.
PMID: 3400190 Russian.
Share

Elevated cholesterol: fact or fancy?
4 Tamura PY.
Cite Hawaii Med J. 1988 Jun;47(6):264-6, 269.
PMID: 3403248 No abstract available.
Share

Fig. S2. Output of the “Diseases to Lipids” analysis from LipiDisease. **(A)** top 10 lipids associated with Coronary Disease and **(B)** articles from PubMed associating the lipid Cholesterol with Coronary Disease.

References

- Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B Methodol.*, **57**, 289–300.
- Fontaine, J. and Andrade-Navarro, M. (2016) Gene Set to Diseases (GS2D): disease enrichment analysis on human gene sets with literature data. *Genomics Comput. Biol.*, **2**, e33.
- R Core Team (2020) R: A Language and Environment for Statistical Computing R Foundation for Statistical Computing, Vienna, Austria.
- Korotkevich, G. *et al.* (2019) Fast gene set enrichment analysis. *BioRxiv*.
- Kuo, T.-C. and Tseng, Y.J. (2018) LipidPedia: a comprehensive lipid knowledgebase. *Bioinformatics*, **34**, 2982–2987.