

Additional File

Additional file Tables

Code (Newman)	Definition (Newman dataset)	Category	Number of rows in original file	Curated and mapped unique structures
N	unaltered natural product, 1997	NP	71	62
NB	botanical drug (defined mixture), 2012	NP	14	3
ND	natural product derivative, 1997	NP	356	333
S	synthetic drug, 1997	Synthetic	463	452
S*	synthetic drug (NP pharmacophore), 1997	Synthetic	65	61
/NM	mimic of natural product, 2003	Synthetic	207 (S*/NM) + 217 (S/NM)	205 (S*/NM) + 217 (S/NM)

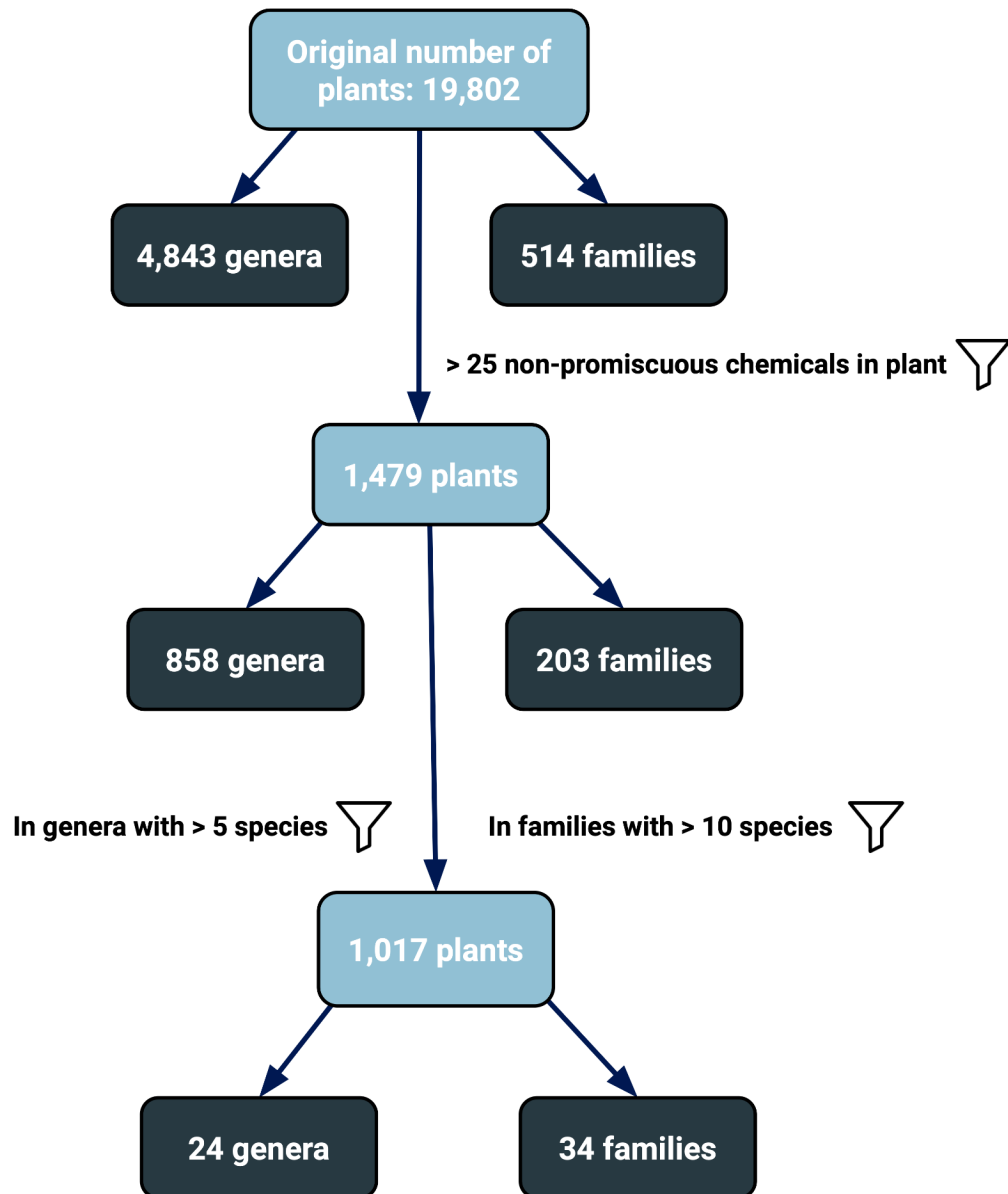
Additional file 1: Table S1. Comparison of the original Newman dataset and the dataset used in this work after normalization and filtering.

NCBITaxonomy ID	Name	Medicinal plant? (use in traditional medicine)
ncbitaxon:46220	<i>Taxus brevifolia</i>	Yes
ncbitaxon:48119	<i>Glehnia littoralis</i>	Yes
ncbitaxon:405945	<i>Leonurus sibiricus</i>	Yes
ncbitaxon:431156	<i>Croton stellatopilosus</i>	Yes
ncbitaxon:555479	<i>Nigella sativa</i>	Yes
ncbitaxon:4682	<i>Allium sativum</i>	Yes
ncbitaxon:108594	<i>Campsis grandiflora</i>	Yes
ncbitaxon:296036	<i>Phyllanthus emblica</i>	Yes
ncbitaxon:147273	<i>Taxus wallichiana</i>	Yes
ncbitaxon:4058	<i>Catharanthus roseus</i>	Yes
ncbitaxon:191701	<i>Cephalotaxus hainanensis</i>	Yes
ncbitaxon:99806	<i>Taxus cuspidata</i>	Yes
ncbitaxon:66169	<i>Cephalotaxus fortunei</i>	Yes

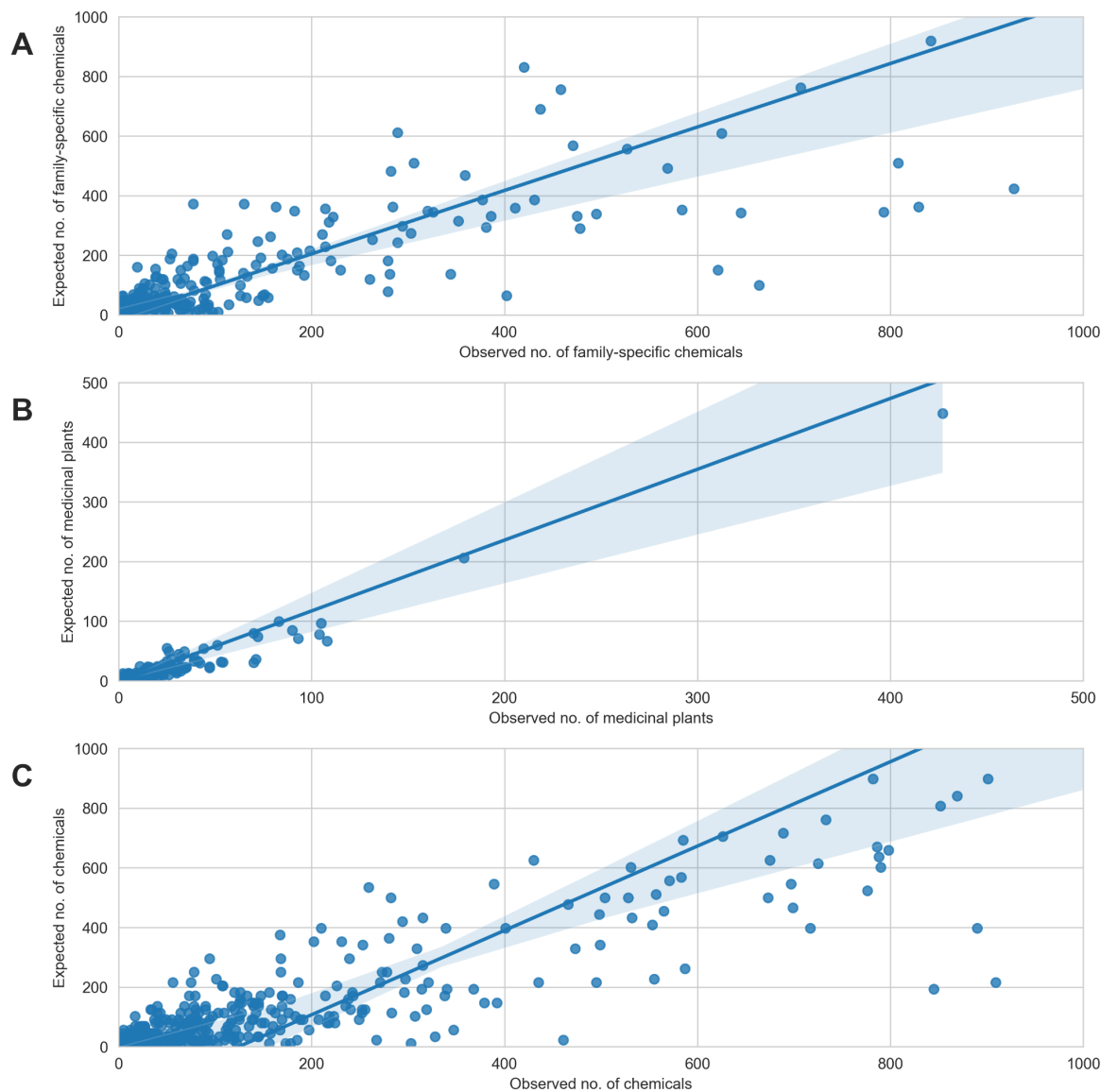
ncbitaxon:65561	<i>Hypericum perforatum</i>	Yes
ncbitaxon:25629	<i>Taxus baccata</i>	Yes
ncbitaxon:109792	<i>Citrus natsudaidai</i>	Yes
ncbitaxon:376254	<i>Ranunculus ternatus</i>	Yes
ncbitaxon:2711	<i>Citrus sinensis</i>	Yes
ncbitaxon:43166	<i>Citrus aurantium</i>	Yes
ncbitaxon:154990	<i>Euphorbia helioscopia</i>	Yes
ncbitaxon:58029	<i>Cephalotaxus harringtonia</i>	Yes
ncbitaxon:3483	<i>Cannabis sativa</i>	Yes
ncbitaxon:330167	<i>Dioscorea villosa</i>	Yes
ncbitaxon:93608	<i>Sinopodophyllum hexandrum</i>	Yes
ncbitaxon:4227	<i>Flaveria trinervia</i>	Yes
ncbitaxon:4182	<i>Sesamum indicum</i>	Yes
ncbitaxon:329759	<i>Solanum aculeastrum</i>	Yes
ncbitaxon:35933	<i>Podophyllum peltatum</i>	Yes
ncbitaxon:29780	<i>Mangifera indica</i>	Yes
ncbitaxon:246360	<i>Commiphora wightii</i>	Yes
ncbitaxon:126910	<i>Withania somnifera</i>	Yes
ncbitaxon:4072	<i>Capsicum annuum</i>	Yes
ncbitaxon:137221	<i>Rheum palmatum</i>	Yes
ncbitaxon:37690	<i>Citrus trifoliata</i>	Yes
ncbitaxon:3993	<i>Euphorbia esula</i>	Yes
ncbitaxon:39354	<i>Salvia abrotanoides</i>	Yes
ncbitaxon:224740	<i>Juniperus sabina</i>	Yes
ncbitaxon:212925	<i>Euphorbia lathyris</i>	Yes
ncbitaxon:107238	<i>Croton sublyratus</i>	Yes
ncbitaxon:2067815	<i>Commiphora mukul</i>	Yes
ncbitaxon:88032	<i>Taxus canadensis</i>	Yes
ncbitaxon:85957	<i>Taxus x media</i>	No
ncbitaxon:2708766	<i>Cola ballayi</i>	No
ncbitaxon:69918	<i>Micranthemum umbrosum</i>	No
ncbitaxon:89484	<i>Cephalotaxus sinensis</i>	No
ncbitaxon:417013	<i>Livistoneae incertae sedis</i>	No
ncbitaxon:1721085	<i>Pentzia eonii</i>	No

Additional fil 1: Table S2. List of plants containing phytochemicals that are now approved-drugs. The last column indicates whether the plant is considered a medicinal plant (e.g., has been traditionally used to treat indications).

Additional file Figures

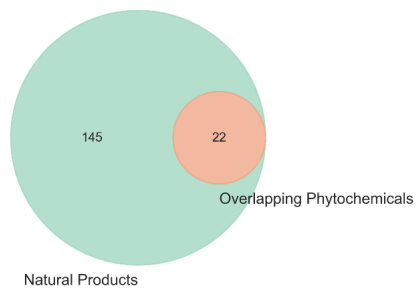


Additional file 1: Figure S1. Number of species (plants) and their corresponding genera and families after each filtering step for the showcase of the chemotaxonomy.

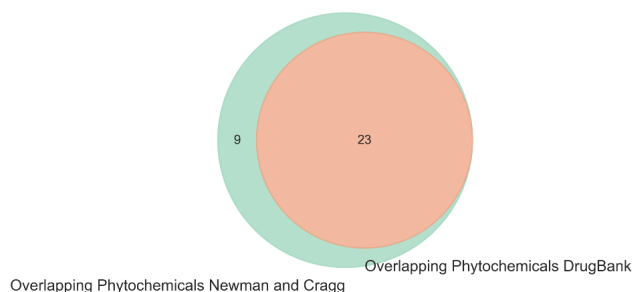


Additional file 1: Figure S2. **A)** Correlation between the observed chemicals specific to the family against the expected chemicals specific to the same family, corrected by the total number of species in the family. **B)** Correlation between the observed number of medicinal plants found in the family against the expected number of medicinal plants belonging to the same family, corrected by the total number of species in the family. **C)** Correlation between the observed number of chemicals found in the family against the expected number of chemicals belonging to the same family, corrected by the total number of species in the family. The y-axis range is set to 1,000, although a few families have over 2,000 chemicals.

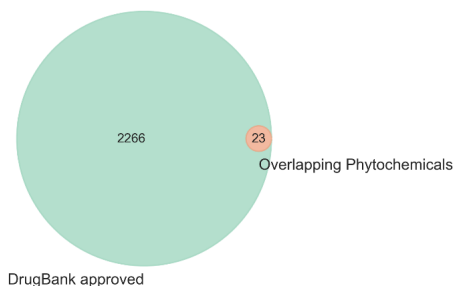
A) Phytochemicals overlap (Newman and Cragg dataset)



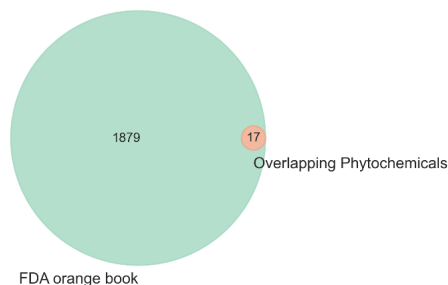
B) Phytochemicals overlap between Wishart (DrugBank) and Newman and Cragg datasets



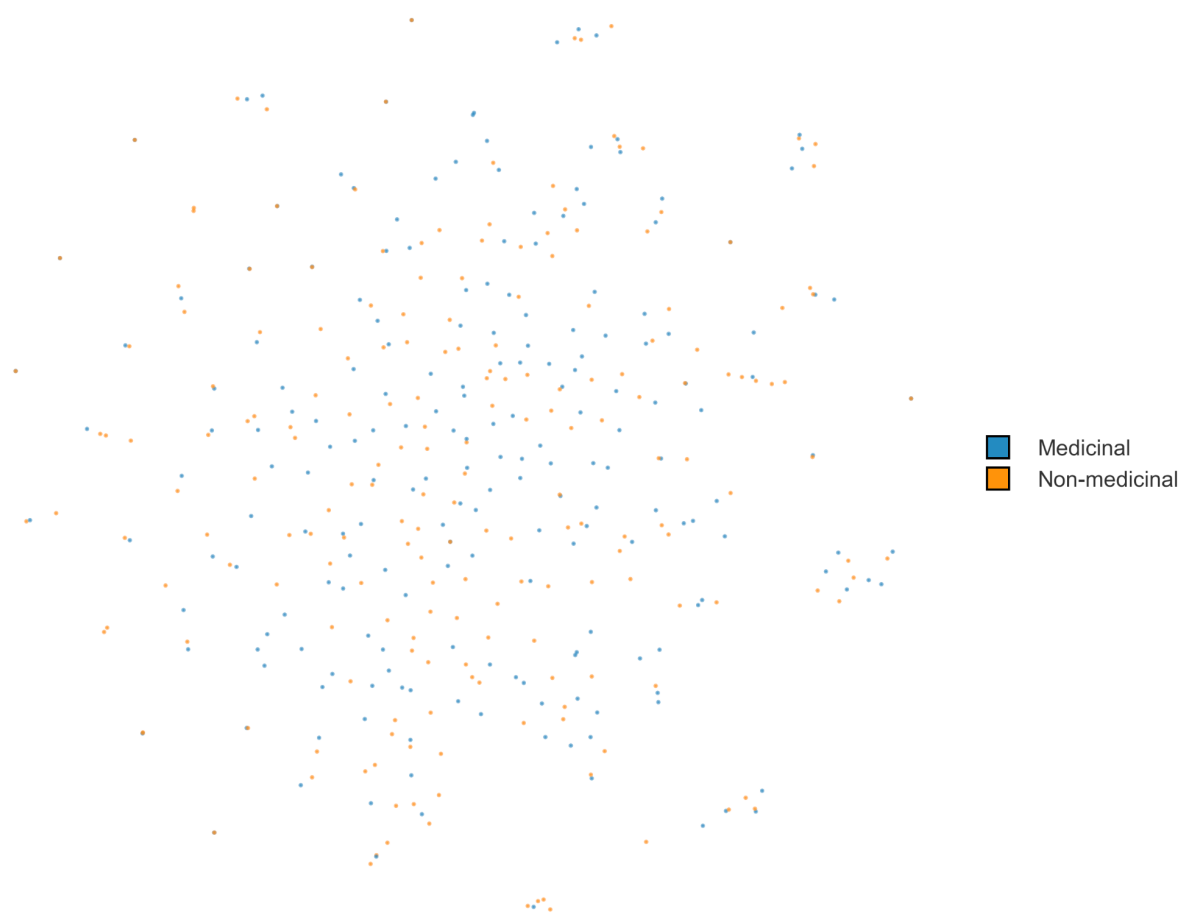
C) Phytochemicals in Wishart dataset (DrugBank)



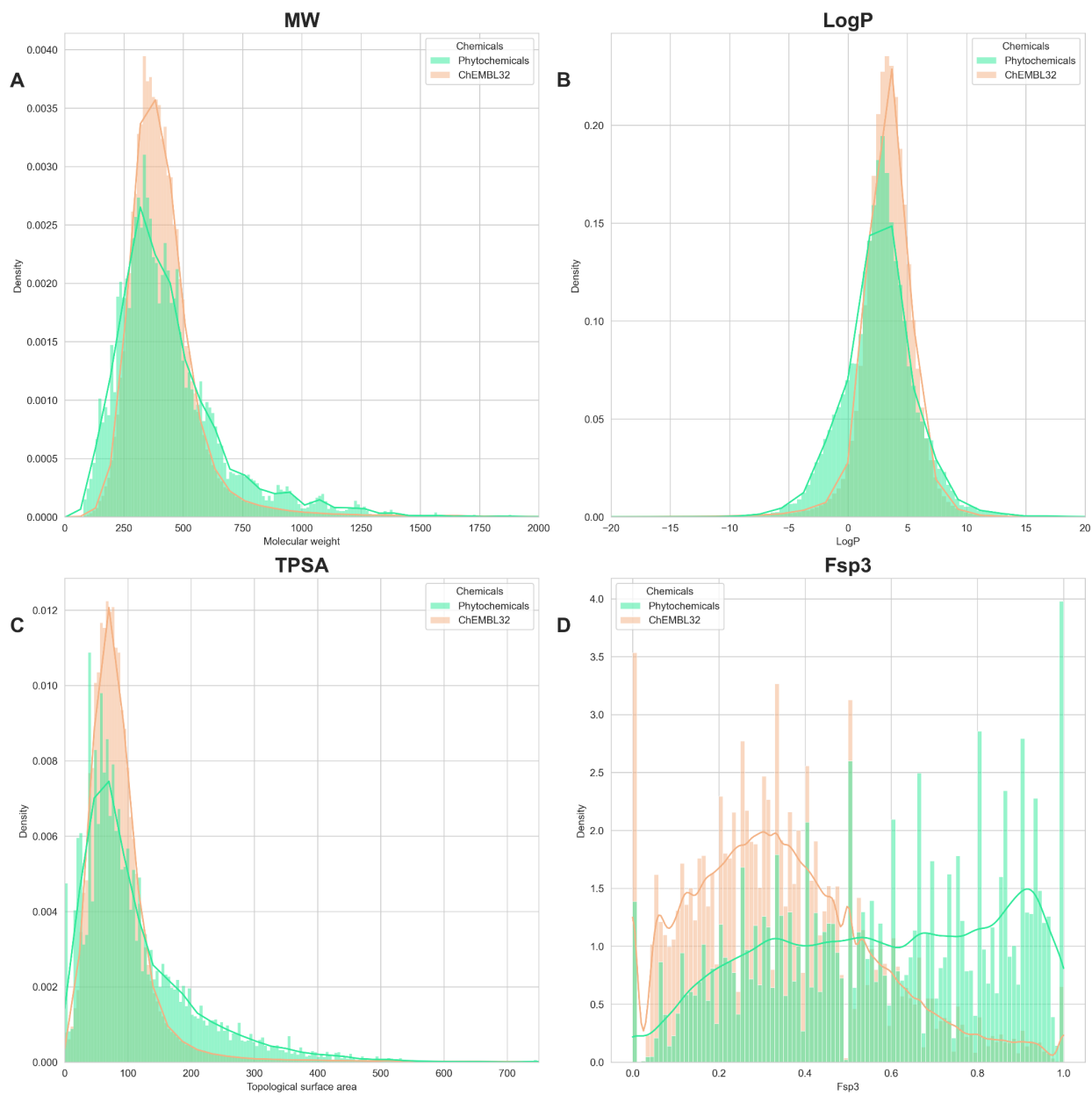
D) Phytochemicals in FDA orange book



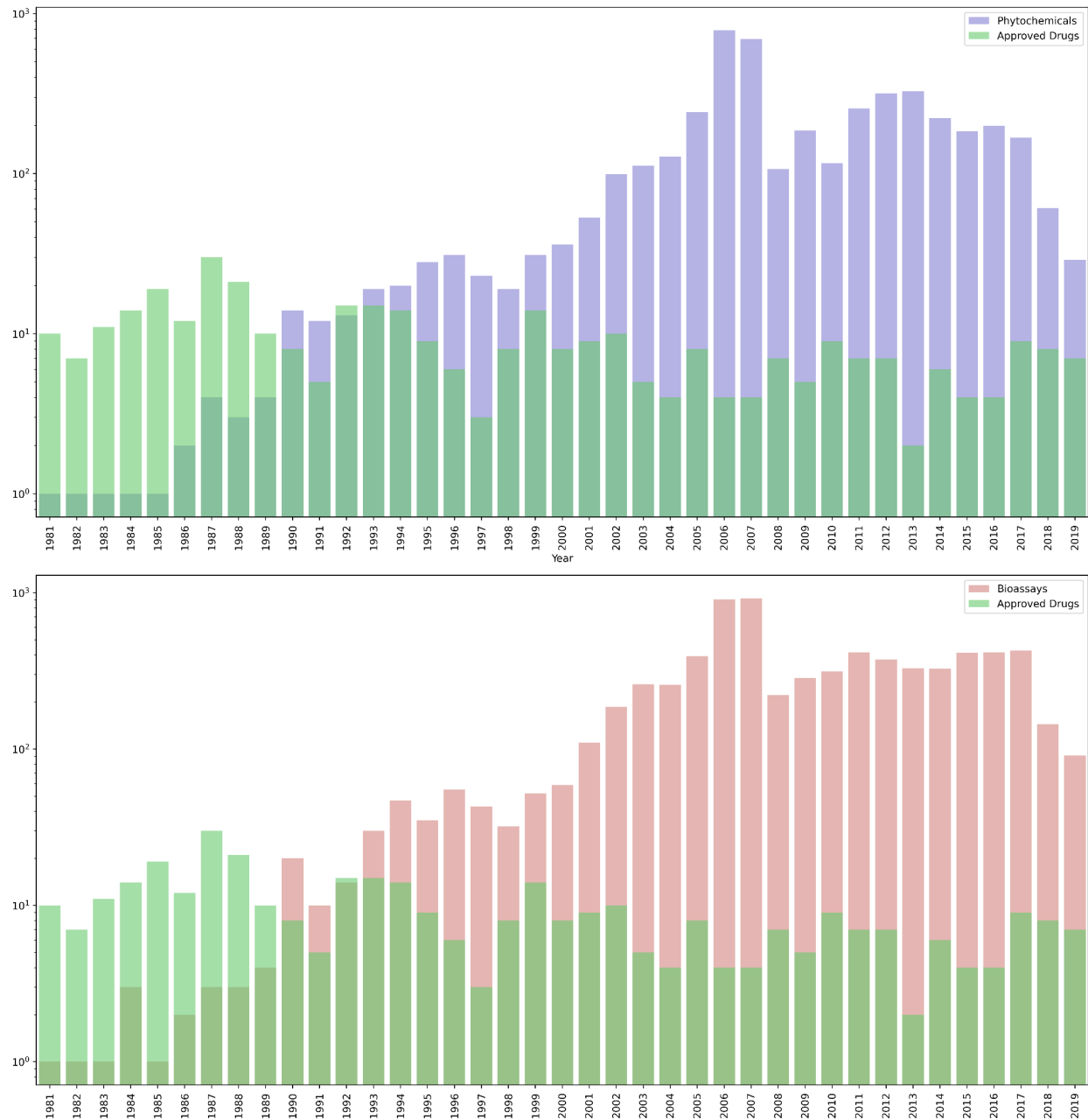
Additional file 1: Figure S3. **A)** Proportion of plant-specific compounds present in the NP approved-drugs curated by Newman and Cragg (2020). **B)** Overlap of the matching phytochemicals of the two datasets: Newman and Cragg (2020) and Wishart *et al.* (2018). **C)** Number of plant-specific compounds present in the dataset curated by Wishart *et al.* (2018). **D)** Number of plant-specific compounds present in the FDA orange book.



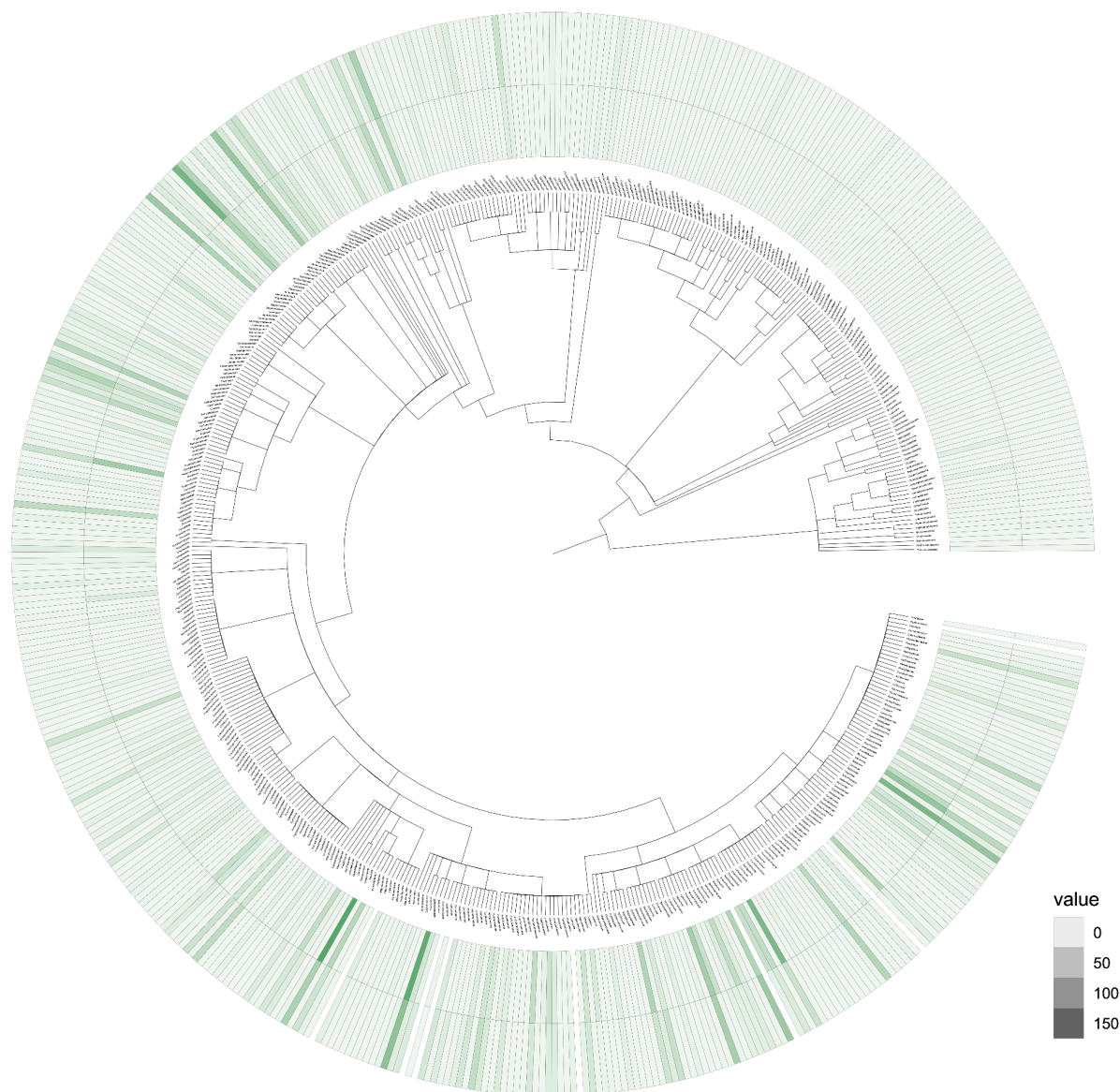
Additional file 1: Figure S4. t-SNE of the relative abundance of the chemical classes from NP-classifier.



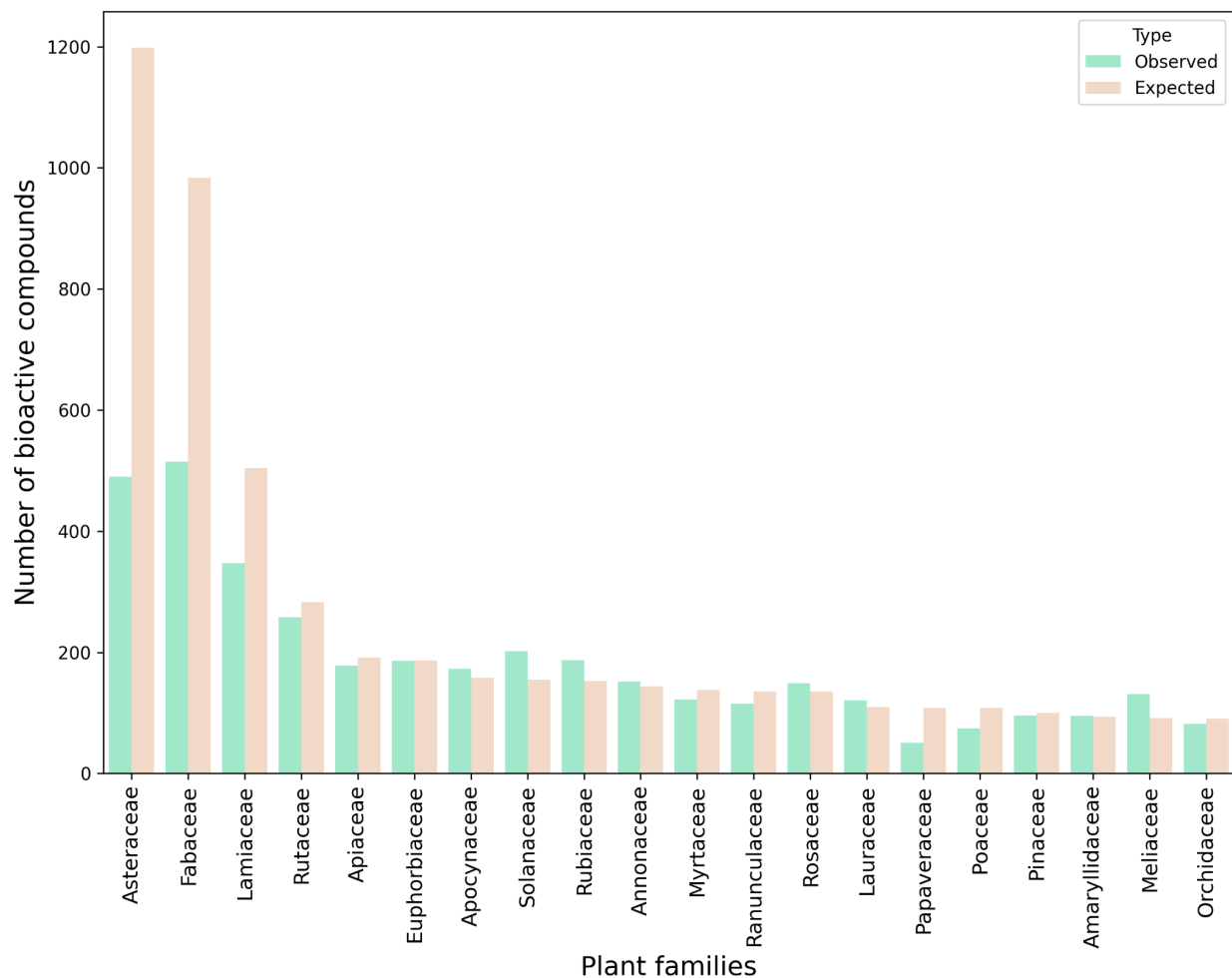
Additional file 1: Figure S5. Distribution of four chemical properties between ChEMBL compounds (version 32) and all phytochemicals in medicinal and non-medicinal plants used in our work. A) Distribution of the molecular weights (MW), B) Distribution of the LogP. C) Distribution of the topological polar surface area (TPSA) D) Distribution of the fraction of sp³ hybridized carbon atoms (Fsp₃).



Additional file 1: Figure S6. Comparison of the log distributions of NP approved drugs versus phytochemicals and their corresponding bioassays over time. The plots highlight a decline in NP approved drugs around the late 1980s. At the same time, the plots indicate an increase in the number of phytochemicals being tested for bioactivity, especially in the last 15 years.



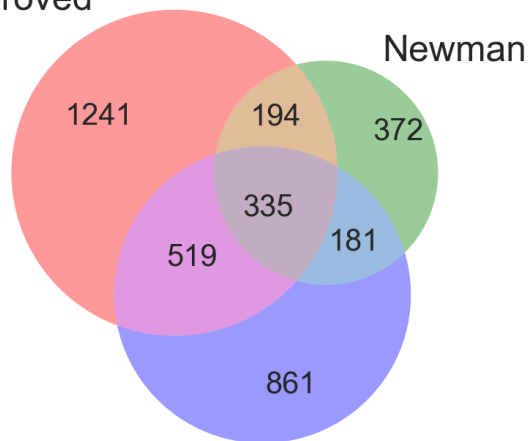
Additional file 1: Figure S7. Distribution of the number of bioactive compounds in medicinal (inner circle) and non-medicinal (outer circle) plants across plant families. Due to a few families having a disproportionately large number of bioactive compounds compared to the rest, we set their values to white any family with more than 150 bioactive compounds in both medicinal and non-medicinal plants to able to spot easier the differences between the two groups for the rest with a smaller range in the color palette (intensity).



Additional file 1: Figure S8. Distribution of the number of observed and expected bioactive compounds across the plant families with the highest expected values. The expected number of bioactive compounds is calculated by multiplying the average number of bioactive compounds per plant by the number of species in a family. The plot shows that a very low percentage of bioactive compounds have been identified in plant families such as Asteraceae and Fabaceae, unlike Pinaceae where the number of bioactive compounds identified is relatively similar to the number expected for this plant family.

Overlap between DrugBank, FDA orange book, and Newman datasets

DrugBank approved



FDA orange book

Additional file 1: Figure S9. Overlap between the DrugBank, FDA orange book and Newman and Cragg datasets.