



Cannabis Data Science

Cananbis Data Science #144

January 25th, 2024

The phytochemical diversity of commercial *Cannabis* in the United States

Christiana J. Smith¹, Daniela Vergara², Brian Keegan³, Nick Jikomes^{4*}

1 Independent Researcher, Seattle, Washington, United States of America, **2** Department of Ecology and Evolutionary Biology, University of Colorado Boulder, Boulder, Colorado, United States of America, **3** Department of Information Science, University of Colorado Boulder, Boulder, Colorado, United States of America, **4** Department of Science and Innovation, Leafly Holdings Inc, Seattle, Washington, United States of America

* njikomes@gmail.com



Abstract

The legal status of *Cannabis* is changing, fueling an increasing diversity of *Cannabis*-derived products. Because *Cannabis* contains dozens of chemical compounds with potential psychoactive or medicinal effects, understanding this phytochemical diversity is crucial. The legal *Cannabis* industry heavily markets products to consumers based on widely used labeling systems purported to predict the effects of different “strains.” We analyzed the cannabinoid and terpene content of commercial *Cannabis* samples across six US states, finding distinct chemical phenotypes (chemotypes) which are reliably present. By comparing the observed phytochemical diversity to the commercial labels commonly attached to *Cannabis*-derived product samples, we show that commercial labels do not consistently align with the observed chemical diversity. However, certain labels do show a biased association with specific chemotypes. These results have implications for the classification of commercial *Cannabis*, design of animal and human research, and regulation of consumer marketing—areas which today are often divorced from the chemical reality of the *Cannabis*-derived material they wish to represent.

OPEN ACCESS

Citation: Smith CJ, Vergara D, Keegan B, Jikomes N (2022) The phytochemical diversity of commercial *Cannabis* in the United States. PLoS ONE 17(5): e0267498. <https://doi.org/10.1371/journal.pone.0267498>

Editor: Muhammad Qasim Shahid, South China Agricultural University, CHINA

Received: October 28, 2021

Accepted: April 8, 2022

Published: May 19, 2022

PLOS ONE

Phytochemical diversity of commercial *Cannabis*

Table 3. Sample, cultivator, and cultivar breakdown by testing lab, after data cleaning.

Testing Lab	State	# of samples with cannabinoid measurements	# of samples with cannabinoid and terpene measurements	# of unique cultivators	# of unique cultivar names
CannTest	AK	6,253	6,173	293	834
ChemHistory	OR	13,508	11,720	589	1,538
Confidence Analytics	WA	53,190	11,070	831	1,794
Modern Canna Science	FL	1,620	695	5	121
PSI Labs	MI	7,240	5,268	543	748
SC Labs	CA	8,112	7,917	1,058	1,218
Total		89,923	42,843	3,319	3,087

<https://doi.org/10.1371/journal.pone.0267498.t003>

Between-Producer Cosine Similarity

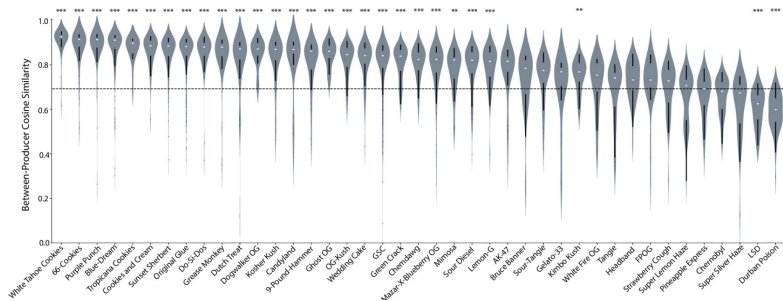


Table 3. Sample, cultivator, and cultivar breakdown by testing lab, after data cleaning.

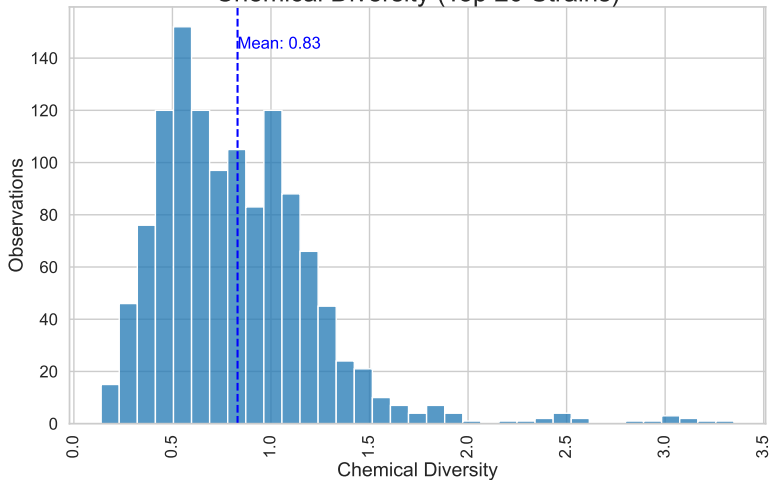
Testing Lab	State	# of samples with cannabinoid measurements	# of samples with cannabinoid and terpene measurements	# of unique cultivators	# of unique cultivar names
CannTest	AK	6,253	6,173	293	834
ChemHistory	OR	13,508	11,720	589	1,538
Confidence Analytics	WA	53,190	11,070	831	1,794
Modern Canna Science	FL	1,620	695	5	121
PSI Labs	MI	7,240	5,268	543	748
SC Labs	CA	8,112	7,917	1,058	1,218
Total		89,923	42,843	3,319	3,087

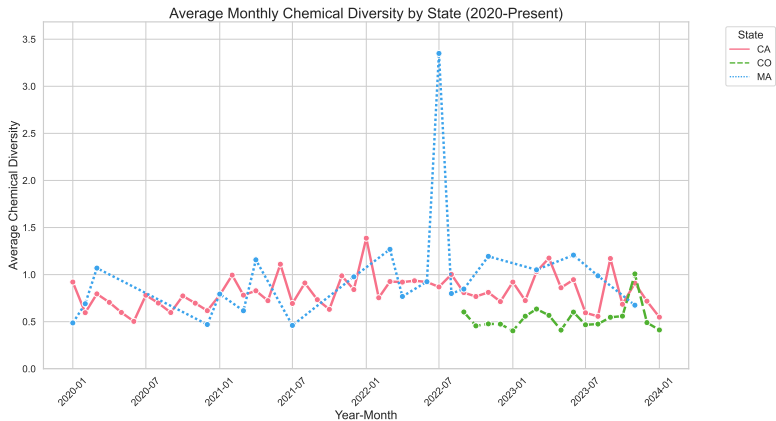
<https://doi.org/10.1371/journal.pone.0267498.t003>

Lab Results Parsed from Public COAs

State	Total Tests	Flower Tests	Terpene Tests	Producers	Strains
CA	65,082	19,850	8,342	169	6,227
CO	11,031	5,177	404	10	205
FL	7,331	4,706	-	-	-
MA	7,189	3,692	1,950	?	1,543
Total	90,633	33,425	10,696	179	7,975

Chemical Diversity (Top 20 Strains)





Hypothesis 1

Hypothesis: Cannabis chemical diversity is increasing over time.

Dep. Variable:	chemical_diversity	R-squared:	0.450
Model:	OLS	Adj. R-squared:	0.434
Method:	Least Squares	F-statistic:	28.39
Date:	Thu, 25 Jan 2024	Prob (F-statistic):	1.71e-13
Time:	15:10:22	Log-Likelihood:	-10.071
No. Observations:	108	AIC:	28.14
Df Residuals:	104	BIC:	38.87
Df Model:	3		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	0.7038	0.064	10.948	0.000	0.576	0.831
C(lab_state)[T.CO]	-0.4892	0.083	-5.895	0.000	-0.654	-0.325
C(lab_state)[T.MA]	0.2708	0.057	4.755	0.000	0.158	0.384
year_month_numeric	0.0064	0.002	3.097	0.003	0.002	0.010

Common strains in the public lab results

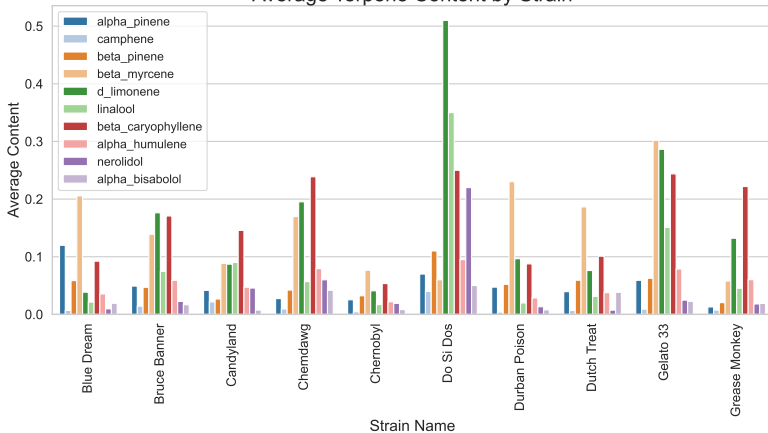
Top strains mentioned in the paper

Strain Name	Observations
Wedding Cake	277
Blue Dream	163
Sour Diesel	128
Headband	81
Mimosa	75
Tangie	71
Purple Punch	60
Green Crack	51
Kosher Kush	47
Gelato 33	41
Chemdawg	34
Grease Monkey	32
Pineapple Express	27
Super Lemon Haze	26
Lemon G	26
Durban Poison	25
Strawberry Cough	21
Sunset Sherbert	18
Sour Tangie	15
Chernobyl	15

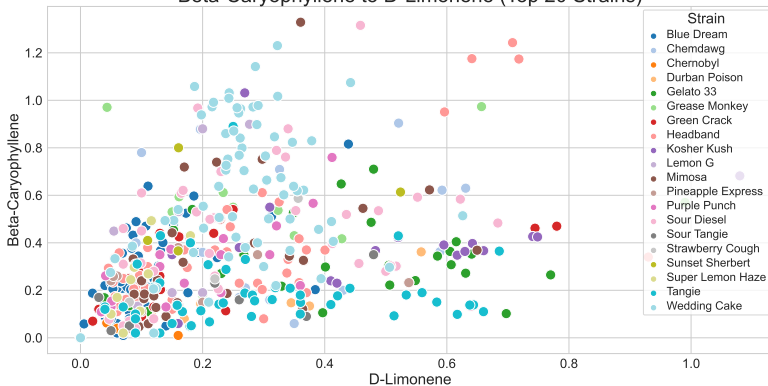
Top strains in the public lab results

Strain Name	Observations
Gelato	715
Runtz	462
Wedding Cake	278
GMO	194
Blue Dream	163
GG4	96
Zkittles	72
OG Kush	70
White Truffle	42
Dolato	27
Skywalker OG	21
Strawberry Cough	21

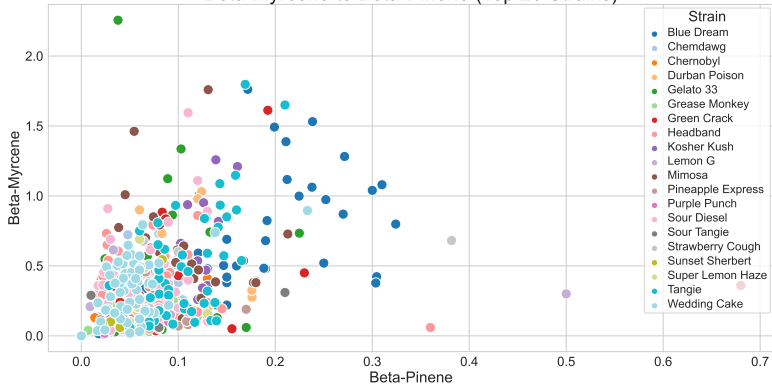
Average Terpene Content by Strain



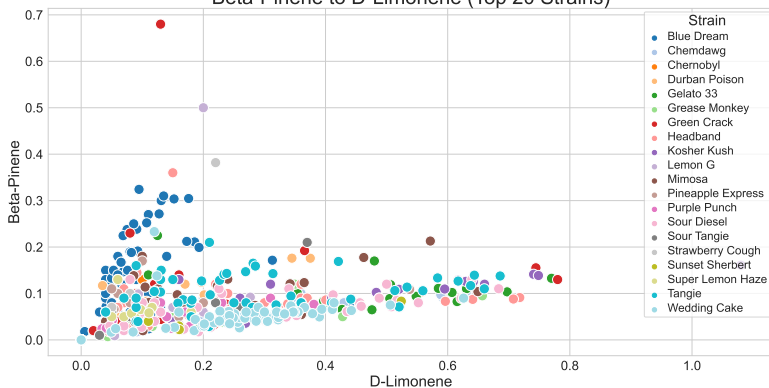
Beta-Caryophyllene to D-Limonene (Top 20 Strains)



Beta-Myrcene to Beta-Pinene (Top 20 Strains)



Beta-Pinene to D-Limonene (Top 20 Strains)

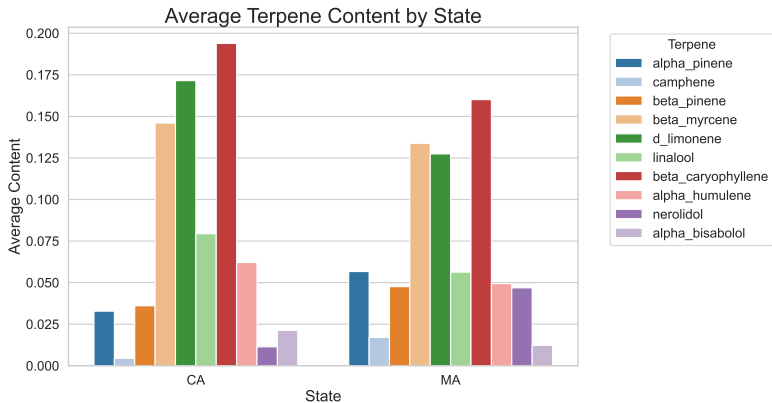


Hypothesis 2

Hypothesis: "Blue Dream" strains have higher beta-pinene than other top strains, on average.

Dep. Variable:	beta_pinene	R-squared:	0.021
Model:	OLS	Adj. R-squared:	0.021
Method:	Least Squares	F-statistic:	28.03
Date:	Thu, 25 Jan 2024	Prob (F-statistic):	1.40e-07
Time:	15:14:07	Log-Likelihood:	1828.4
No. Observations:	1288	AIC:	-3653.
Df Residuals:	1286	BIC:	-3643.
Df Model:	1		
Covariance Type:	nonrobust		

	coef	std err	t	P > t	[0.025	0.975]
const	0.0326	0.002	18.647	0.000	0.029	0.036
is_blue_dream	0.0260	0.005	5.294	0.000	0.016	0.036



Lab results by county (for the top 20 strains)

County	Observations
Denver County	121
Humboldt County	56
Mendocino County	48
Los Angeles County	42
Sacramento County	38
Sonoma County	23
Monterey County	20
Santa Barbara County	19
Alameda County	14
Lake County	13
El Paso County	13
Larimer County	11
San Francisco County	8
Trinity County	7
Riverside County	7
Santa Cruz County	6
Huerfano County	5
San Bernardino County	4
Routt County	4
Yolo County	4

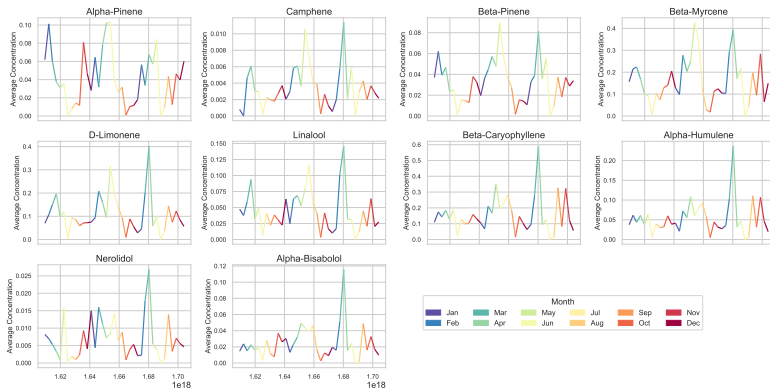
Hypothesis 3

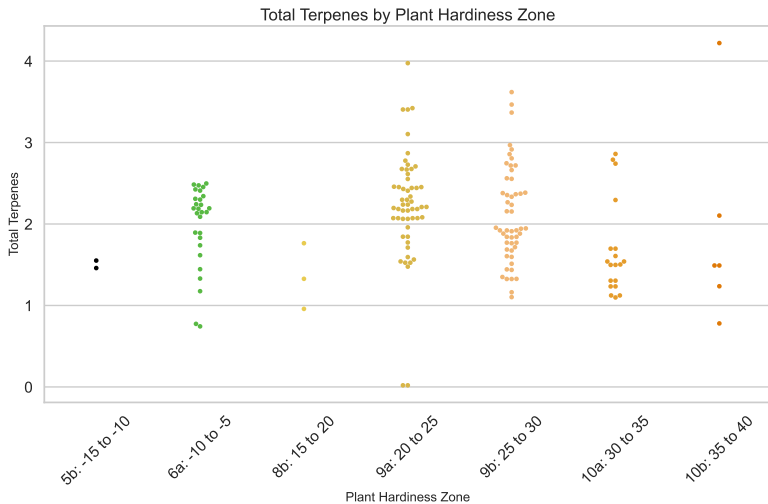
Hypothesis: "Emerald Triangle" counties, Humboldt, Mendocino, and Trinity, produce cannabis with higher terpene concentrations than cannabis produced in other counties, on average.

Dep. Variable:	total_terpenes	R-squared:	0.046
Model:	OLS	Adj. R-squared:	0.046
Method:	Least Squares	F-statistic:	1622.
Date:	Thu, 25 Jan 2024	Prob (F-statistic):	0.00
Time:	15:16:06	Log-Likelihood:	-44550.
No. Observations:	33425	AIC:	8.910e+04
Df Residuals:	33423	BIC:	8.912e+04
Df Model:	1		
Covariance Type:	nonrobust		

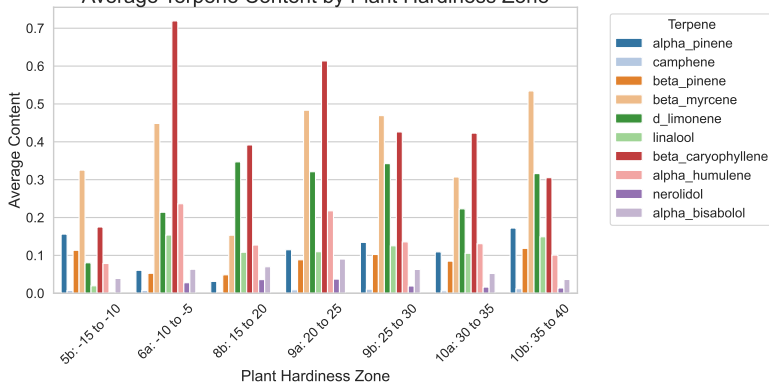
	coef	std err	t	P> t	[0.025	0.975]
Intercept	0.5267	0.005	101.567	0.000	0.517	0.537
C(is_emerald_triangle)[T.1]	0.8297	0.021	40.276	0.000	0.789	0.870

Average Terpene Concentration by Month





Average Terpene Content by Plant Hardiness Zone





Thank you for coming.

Lesson of the Day

- Always look on the bright side of life.