

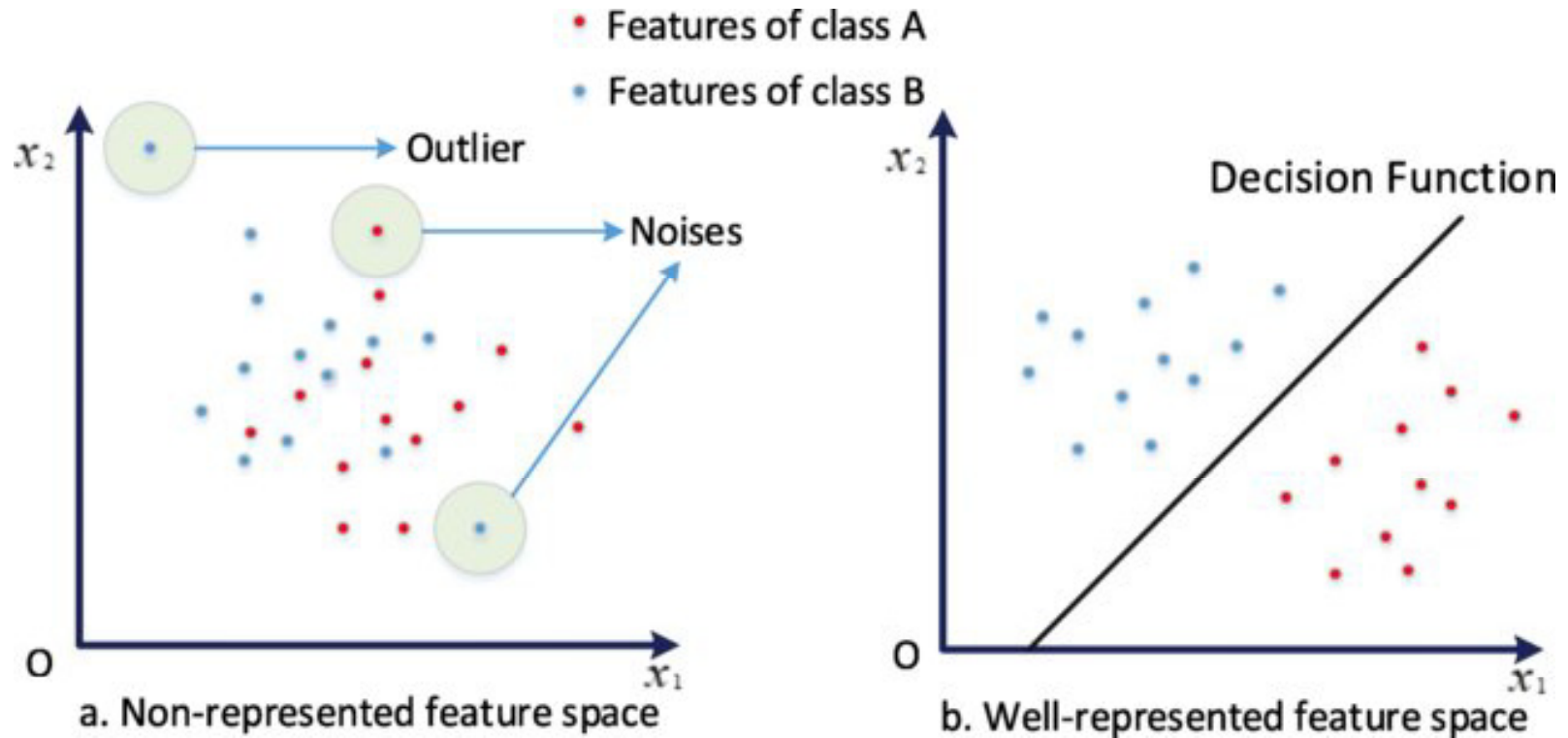
تحلیل داده‌های طرح پیشگیری از بیماری عروق کرونر زودرس توسط روش‌های یادگیری ماشین

دکتر حمید سعادت فر

تابستان ۱۴۰۳

Data Preparation

□ Feature Space: Choosing Informative Features



Some Examples

□ Feature Space: Choosing Informative Features

آزمایشات کلی					
	AZ8	SGOT (AST)		AZ1	HDL
	AZ9	SGPT (ALT)		AZ2	LDL
	AZ10	APO A		AZ3	Chol
	AZ11	APO B		AZ4	TG
	AZ12	HbA1C		AZ5	hs-CRP
	AZ12a	FBS		AZ6	IL1
				AZ7	IL6

به عنوان مثال، میانگین کلسترول برای افراد بیمار از افراد سالم کمتر می باشد.

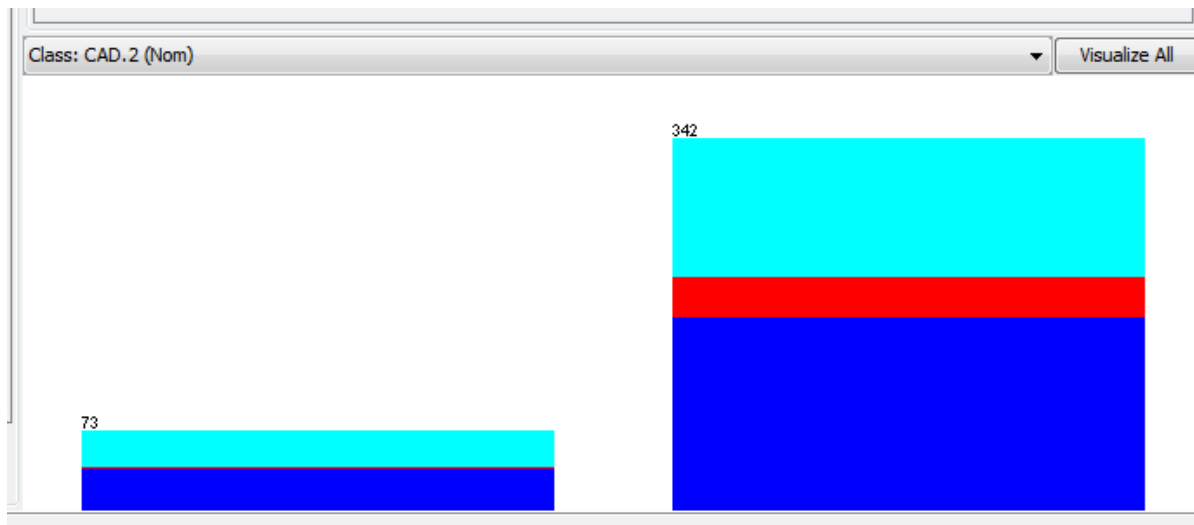
Some Examples

□ Feature Space: Choosing Informative Features

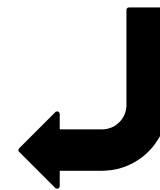
حال که به سوالات فوق پاسخ دادید، با دقت فکر کنید و پاسخ سوالات زیر را به یاد بیاورید:

P17

آیا در طول زندگی خود، غالباً سطح فعالیت فیزیکی شما در همین حدود (که در سوالات بالا پاسخ دادید) بوده است؟
(۱) بله (۲) خیر



۸۲.۴٪ از افراد به سوال فوق جواب خیر داده‌اند!!!



Some Examples

□ Feature Space: Choosing Informative Features

Ethnicity		با توجه به قومیت های موجود، یکی از اعداد را بر اساس گفته بیمار، انتخاب نمایید؟			
		(۱) نمی دانم (۲) فارس (۳) ترک آذربایجانی (۴) گیلک (۵) کرد (۶) عرب (۷) لر (به جز بختیاری) (۸) بلوچ (۹) ترکمن (۱۰) ترک قشقایی (۱۱) بختیاری (۱۲) تالشی (۱۳) ارمنی (۱۴) گرجی (۱۵) آشوری (۱۶) سایر (به دقت ثبت گردد) * در صورتی که محل تولد را نمی دانید، کد ۱ را وارد نمایید.			
	ET8	محل تولد فرد		ET1	قومیت فرد
	ET9	محل تولد پدر		ET2	قومیت پدر
	ET10	محل تولد مادر		ET3	قومیت مادر
	ET11	محل تولد پدربزرگ پدری		ET4	قومیت پدربزرگ پدری
	ET12	محل تولد مادربزرگ پدری		ET5	قومیت مادربزرگ پدری
	ET13	محل تولد پدربزرگ مادری		ET6	قومیت پدربزرگ مادری
	ET14	محل تولد مادربزرگ مادری		ET7	قومیت مادربزرگ مادری

قومیت غالب فارس بوده و بررسی
تاثیر قومیت امکان پذیر نمی باشد

Features with low variance: These features often carry little information and may even introduce noise into the model's predictions.

Some Examples

□ Feature Space: Choosing Informative Features

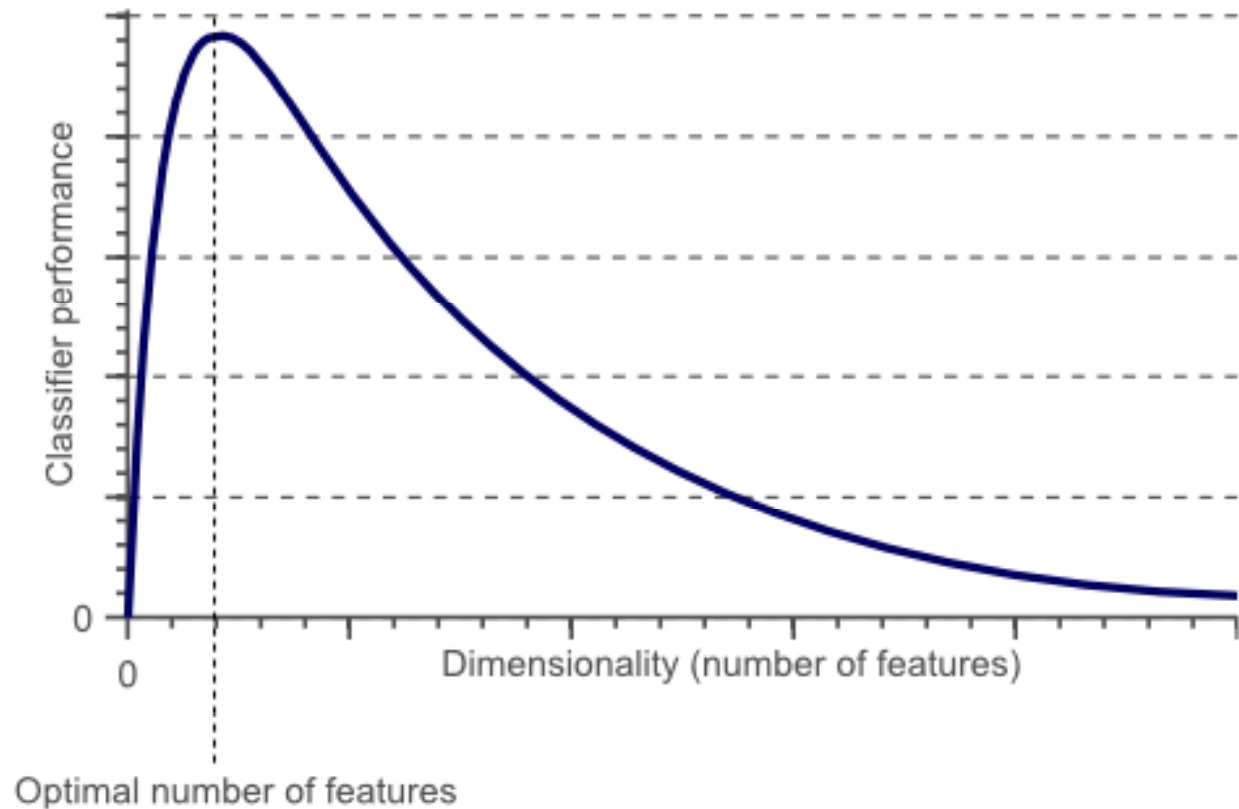
۹۶.۱۴٪ از افراد سابقه آسم،
 ۹۵.۶۶٪ حساسیت غذایی و
 ۹۳٪ حساسیت دارویی
 نداشته‌اند.

		در طول زندگی خود حدوداً چند بار علایم زیر را حس کردید؟				گزینه
گزینه	کد	۴	۳	۲	۱	
	Ar1	بیش از ده بار	شش تا ده بار	یک تا پنج بار	هرگز	حساسیت فصلی
	Ar3	بیش از ده بار	شش تا ده بار	یک تا پنج بار	هرگز	حساسیت دارویی
	Ar5	بیش از ده بار	شش تا ده بار	یک تا پنج بار	هرگز	حساسیت غذایی
	Ar7	بیش از ده بار	شش تا ده بار	یک تا پنج بار	هرگز	کهیر پوستی
	Ar9	(۱) بلی (۲) خیر				آیا سابقه ابتلا به آسم را دارا می‌باشید؟
	Ar11					اگر بلی، مدت ابتلا به آن چند سال بوده است؟
	Ar12	(۱) بلی (۲) خیر (۳) نمی‌دانم				آیا سابقه ابتلا به اگزما را دارا می‌باشید؟
	Ar14					اگر بلی، مدت ابتلا به آن چند سال بوده است؟

Data Preparation

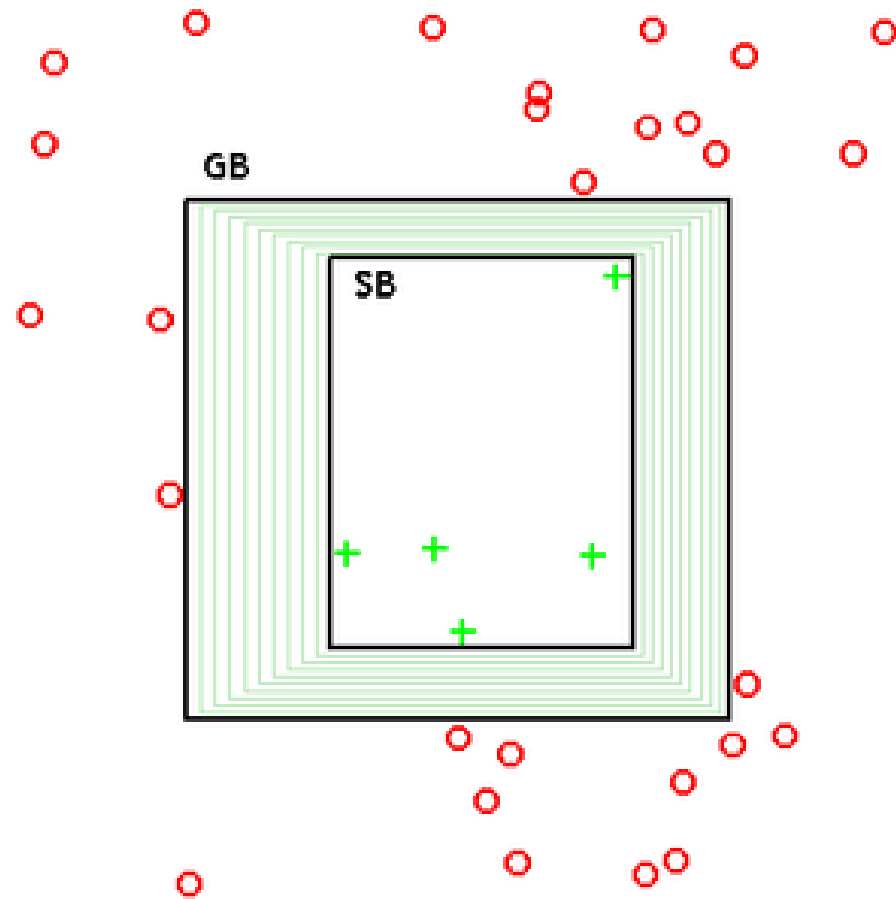
□ Feature Space: Curse of Dimensionality

بالغ بر ۲۰۰ فیلد
اطلاعاتی و تنها
۴۱۵ نمونه !!!



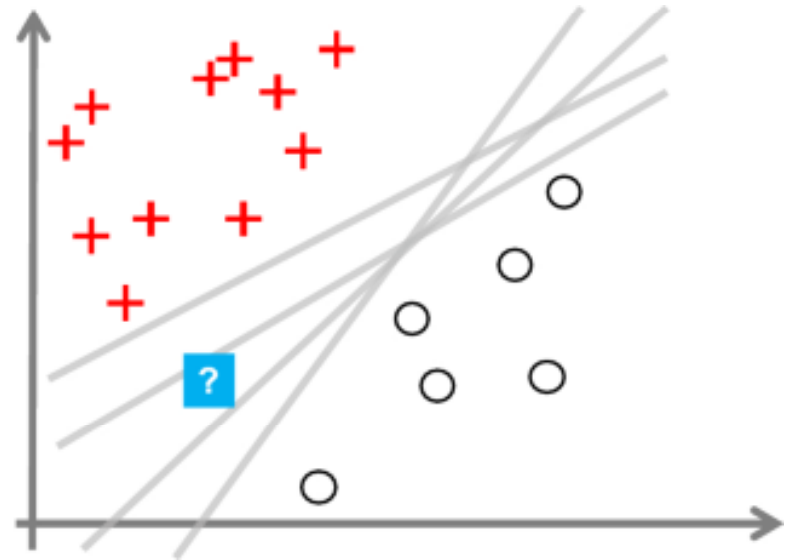
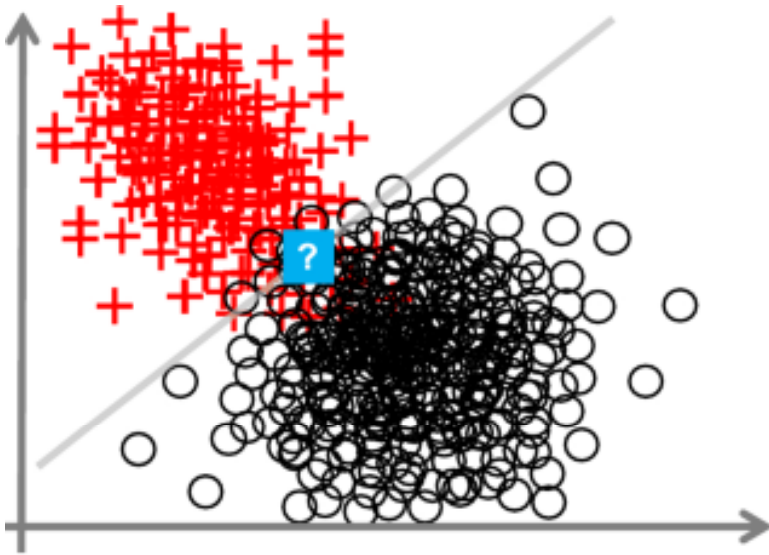
Data Preparation

□ Feature Space: Curse of Dimensionality



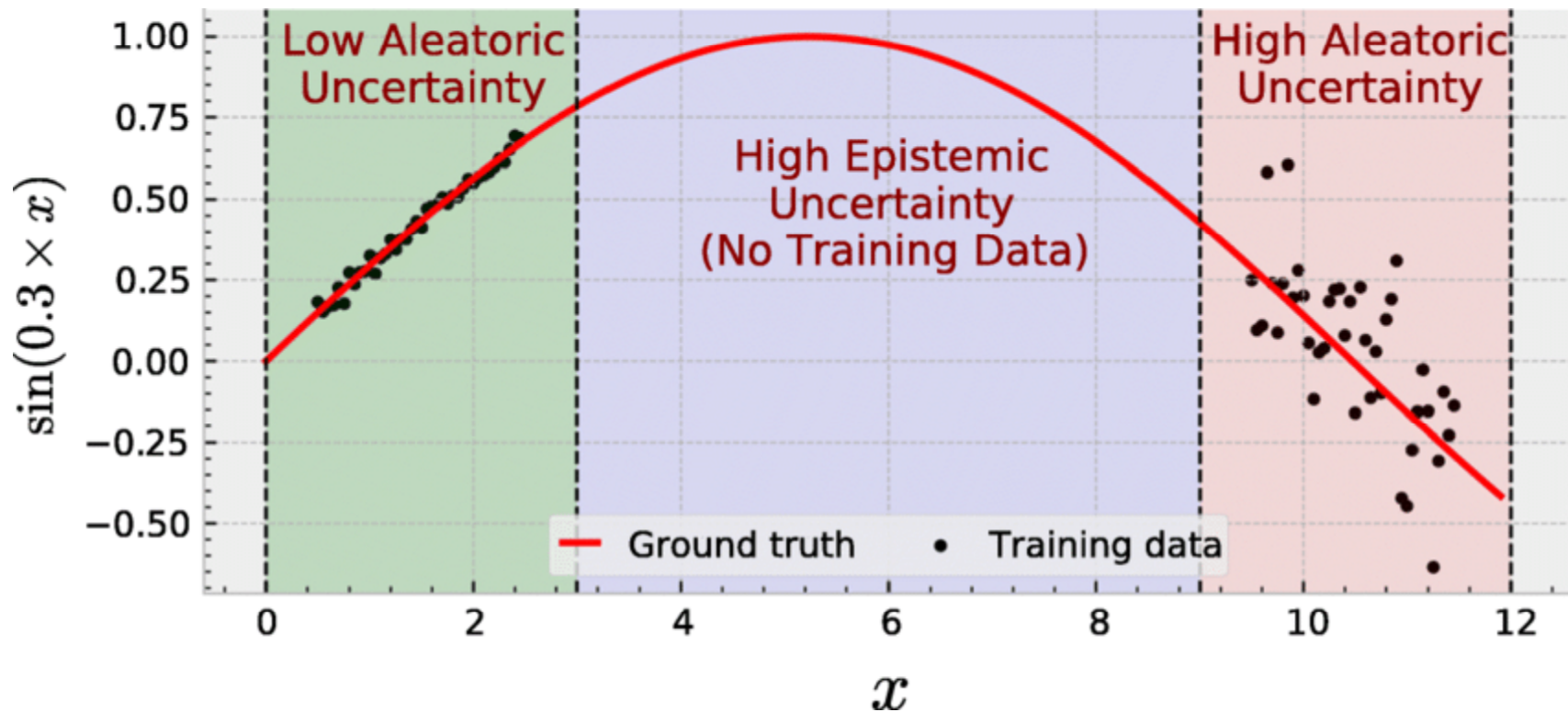
Data Preparation

□ Feature Space: Curse of Dimensionality



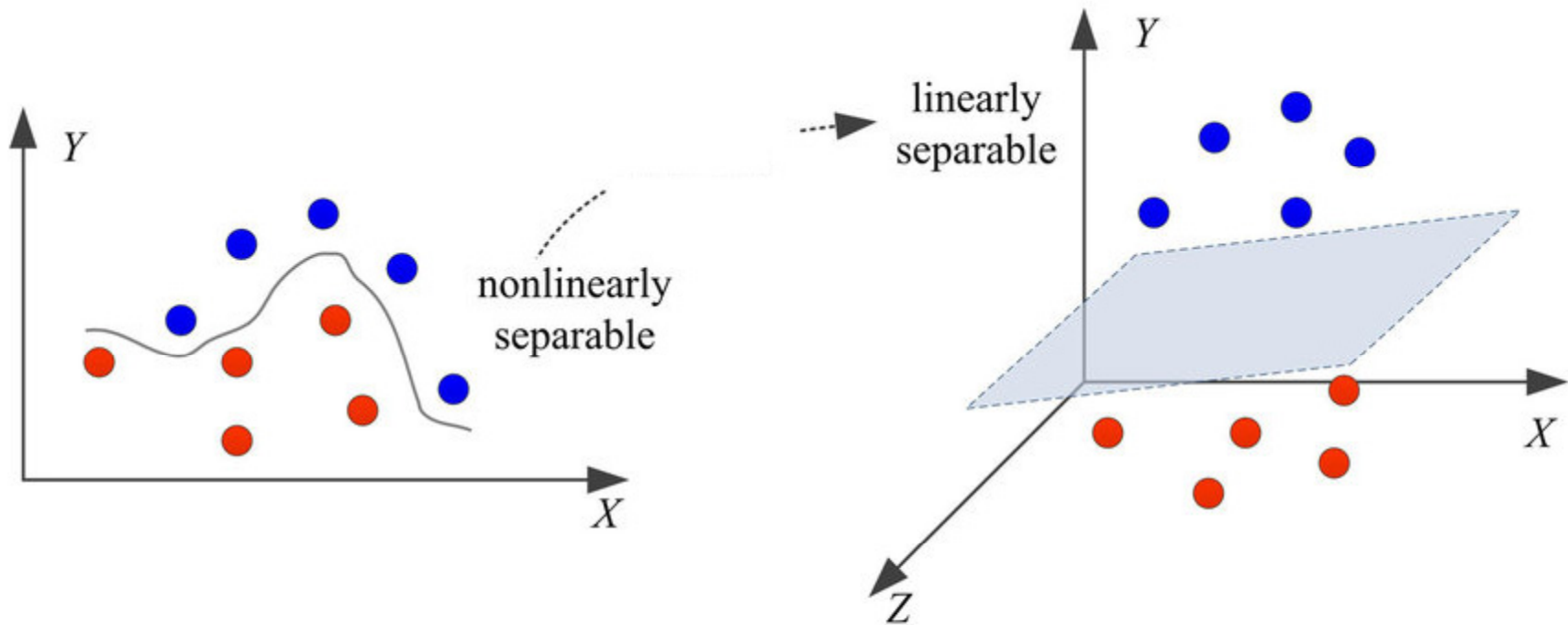
Data Preparation

□ Feature Space: Curse of Dimensionality



Data Preparation

□ **Feature Space:** Considering all important features as possible



Data Preparation

□ **Feature Space: Highly correlated features provide redundant information**

چربی خون		
(بسیار مهم: دقت شود که سوالات زیر، تا قبل از شروع علائم بیماری قلبی نظیر درد سینه یا احساس فشار روی قفسه سینه را مد نظر دارد و این نکته بایستی به درستی به بیمار تفهیم شود).		
HL1	(۱) بلی (۲) خیر	سابقه افزایش چربی خون:
HL1m	ماه:	مدت ابتلا به چربی خون:
HL1y	سال:	
		سابقه مصرف دارو جهت درمان چربی خون:
HL2	(۱) بلی (۲) خیر (۳) علیرغم تجویز پزشک دارو مصرف نمی کردم	
HL2a		اگر بلی، مدت مصرف دارو (سال):
HL3	(۱) بلی (۲) خیر	آیا به طور کلی دارو را مرتب مصرف می کردید؟

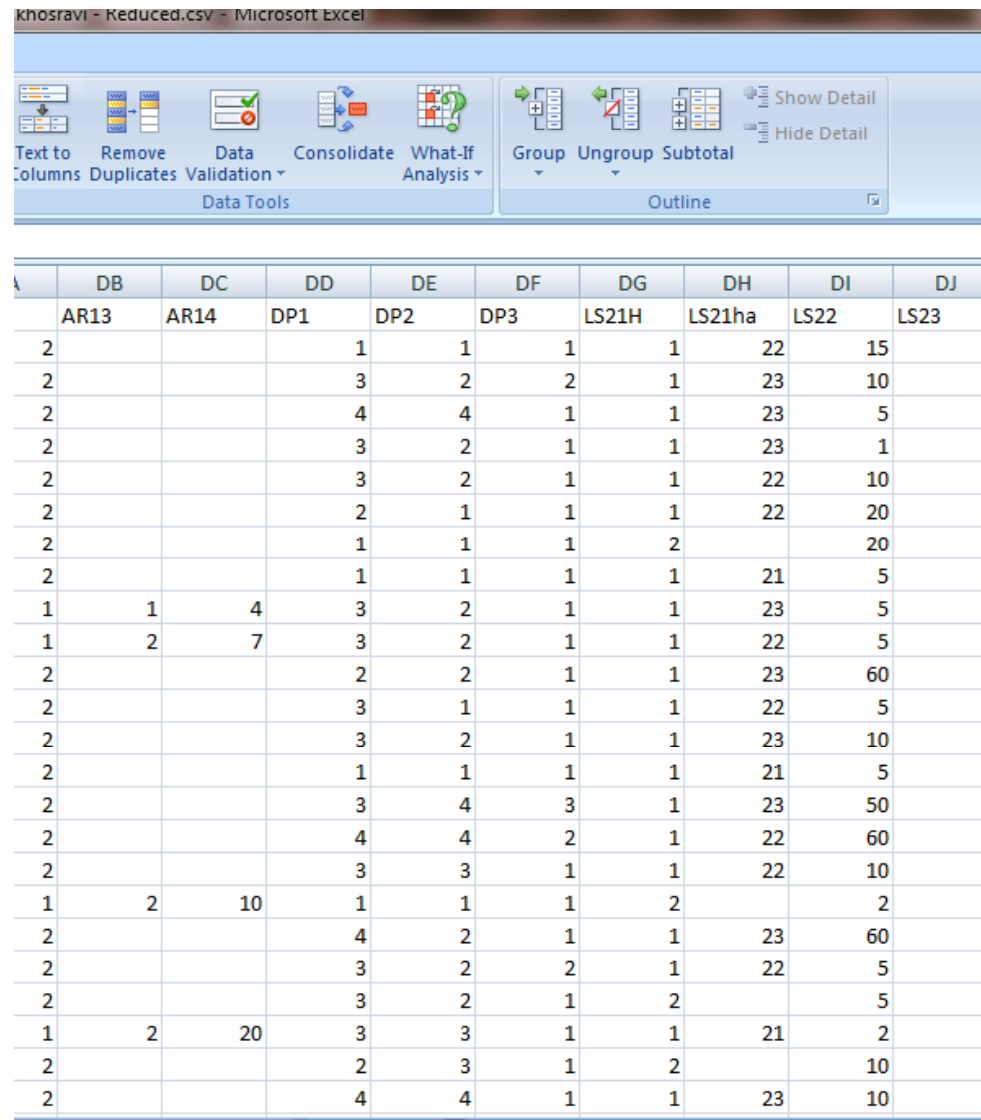


همبستگی بالای ۰.۹ که در موارد متعدد می توان شاهد بود

Data Preparation

□ Feature Space: Handling Missing Values

در مواردی نرم افزارها یا توابع
یادگیری ماشین مقادیر گم شده
را (اگر مدیریت نشوند) صفر در
نظر می گیرند که می تواند به
تحلیل ها آسیب بزند



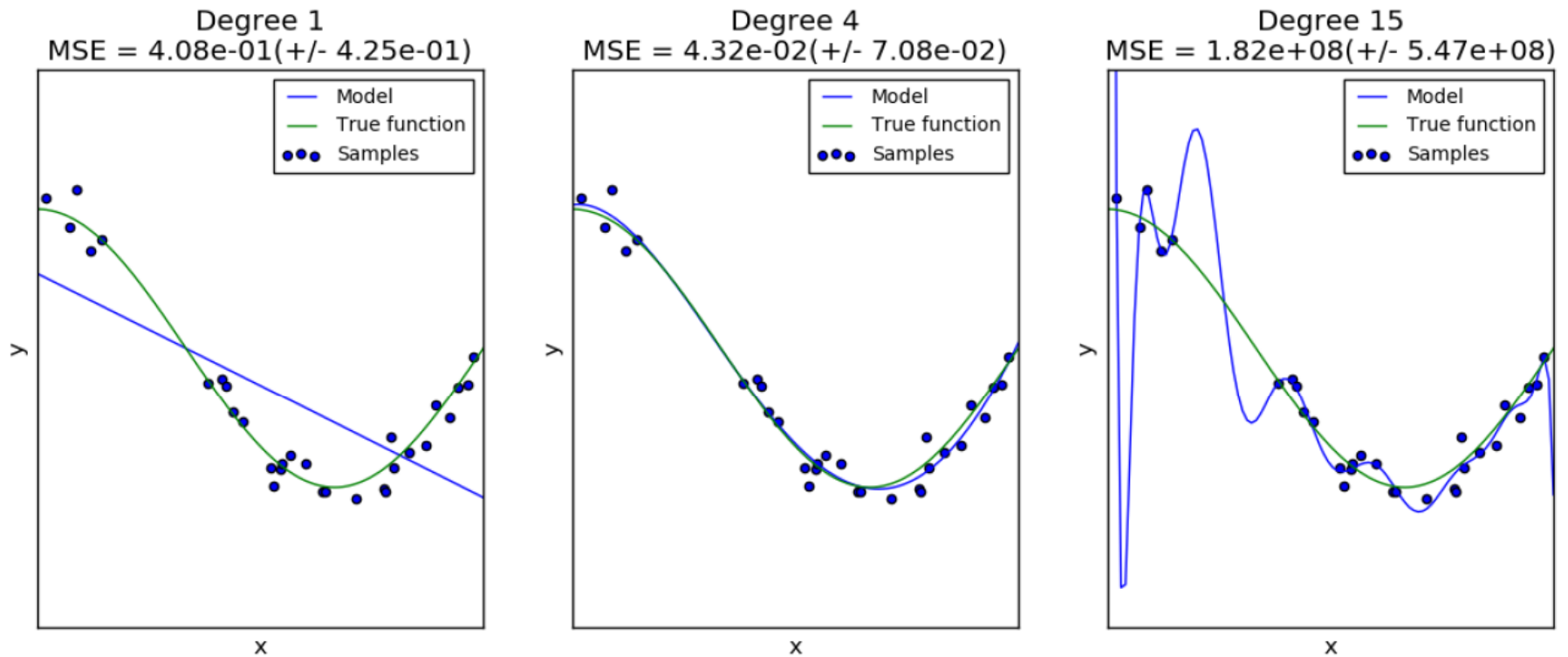
The screenshot shows the Microsoft Excel interface with the 'Data Tools' ribbon selected. The ribbon includes options like 'Text to Columns', 'Remove Duplicates', 'Data Validation', 'Consolidate', 'What-If Analysis', 'Group', 'Ungroup', 'Subtotal', 'Outline', 'Show Detail', and 'Hide Detail'. Below the ribbon is a data table with columns labeled DB, DC, DD, DE, DF, DG, DH, DI, and DJ. The rows contain numerical data, with some cells containing values like 1, 2, 3, 4, 10, 20, 60, and 15.

	DB	DC	DD	DE	DF	DG	DH	DI	DJ
	AR13	AR14	DP1	DP2	DP3	LS21H	LS21ha	LS22	LS23
2			1	1	1	1	22	15	
2			3	2	2	1	23	10	
2			4	4	1	1	23	5	
2			3	2	1	1	23	1	
2			3	2	1	1	22	10	
2			2	1	1	1	22	20	
2			1	1	1	2		20	
2			1	1	1	1	21	5	
1	1	4	3	2	1	1	23	5	
1	2	7	3	2	1	1	22	5	
2			2	2	1	1	23	60	
2			3	1	1	1	22	5	
2			3	2	1	1	23	10	
2			1	1	1	1	21	5	
2			3	4	3	1	23	50	
2			4	4	2	1	22	60	
2			3	3	1	1	22	10	
1	2	10	1	1	1	2		2	
2			4	2	1	1	23	60	
2			3	2	2	1	22	5	
2			3	2	1	2		5	
1	2	20	3	3	1	1	21	2	
2			2	3	1	2		10	
2			4	4	1	1	23	10	

High Generalization Ability

□ Important Point:

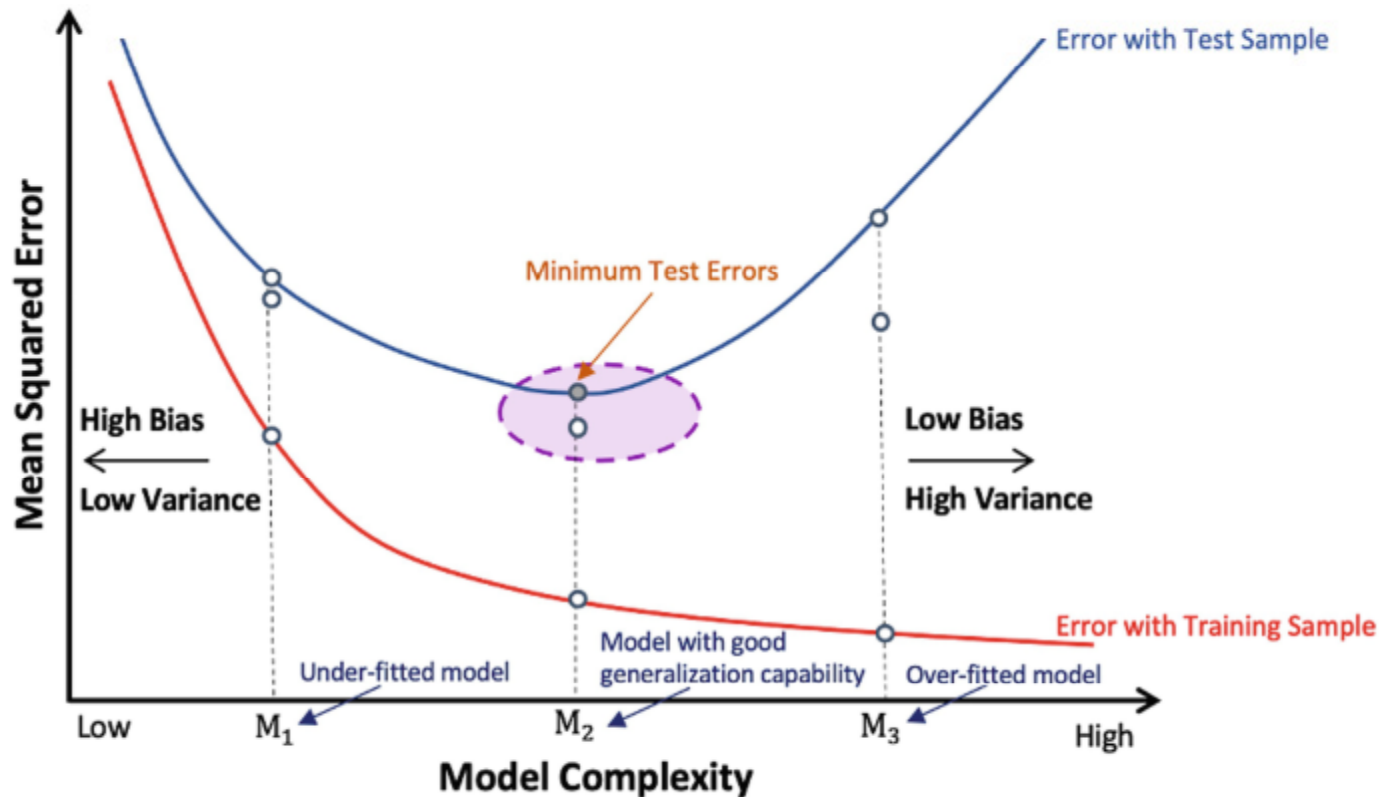
high accuracy on training data does not guarantee effectiveness in the real world



High Generalization Ability

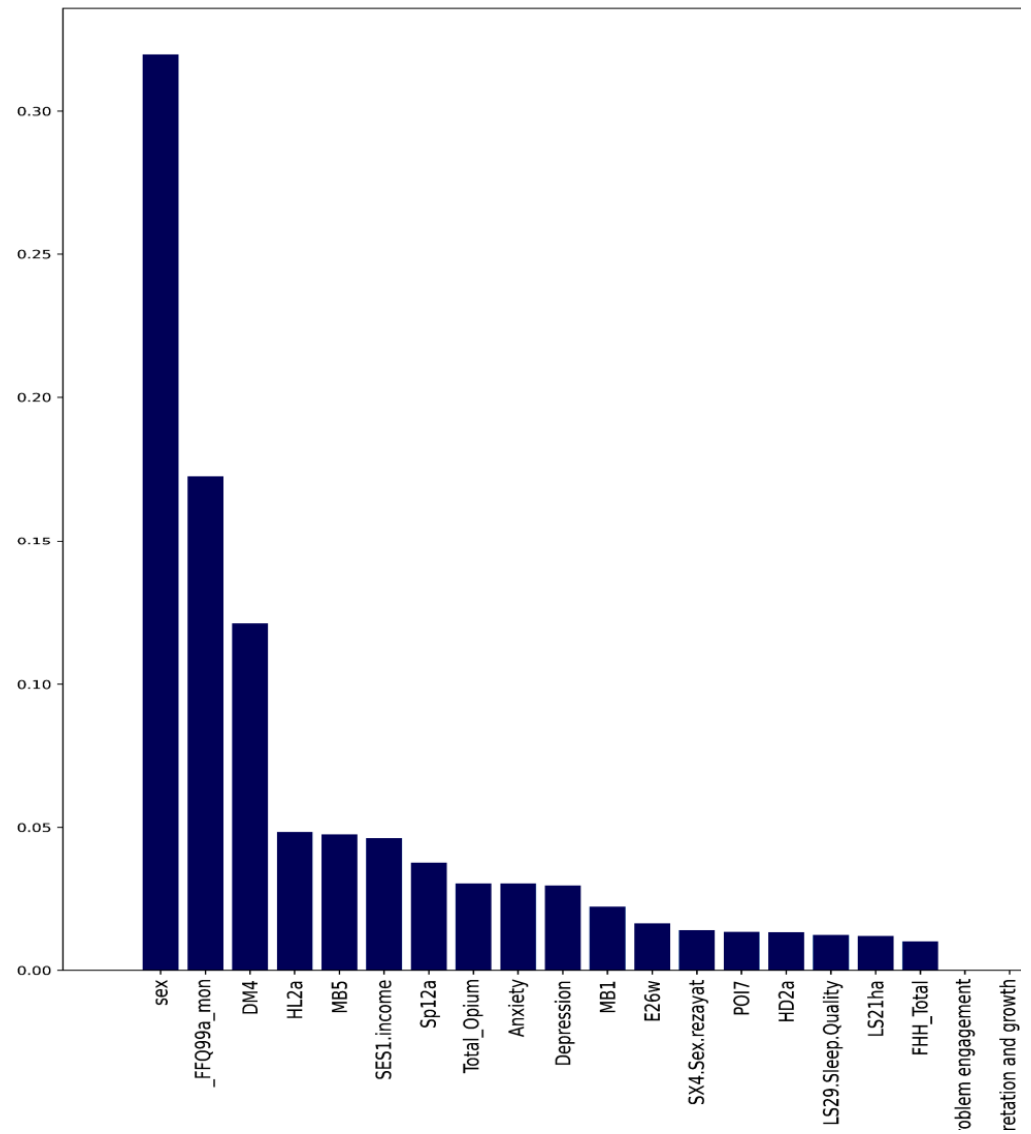
□ Important Point:

high accuracy on training data does not guarantee effectiveness in the real world



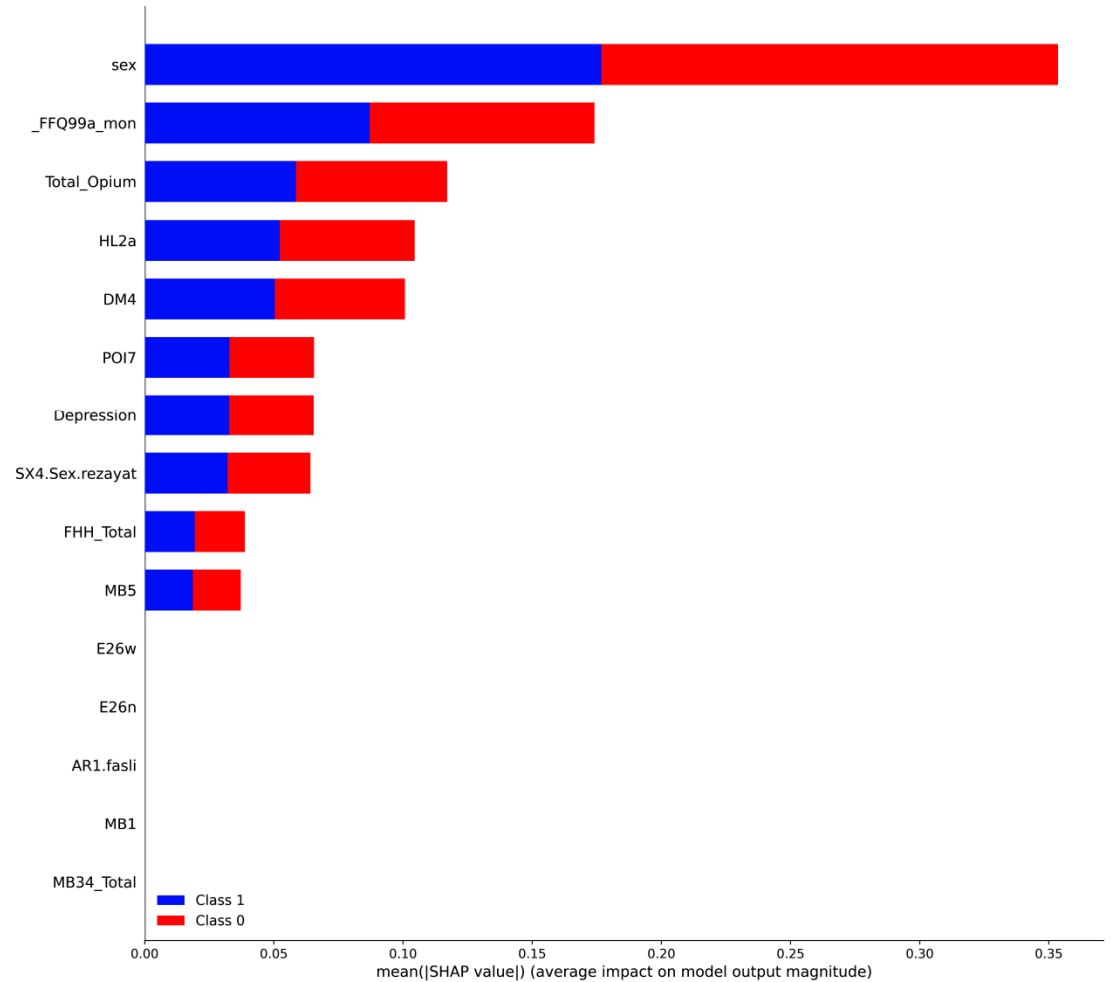
Importance of Features

□ Decision Tree



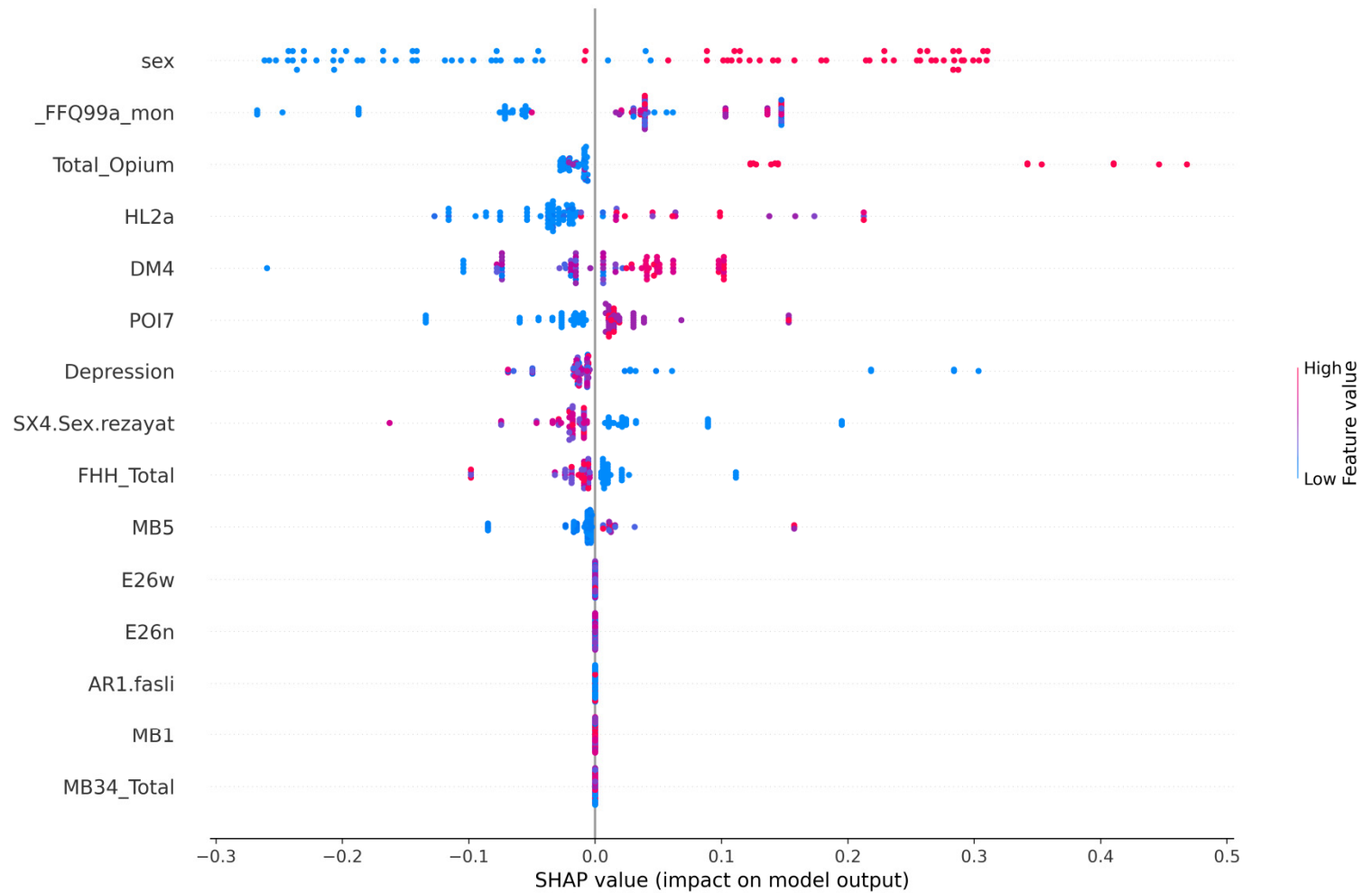
Importance of Features

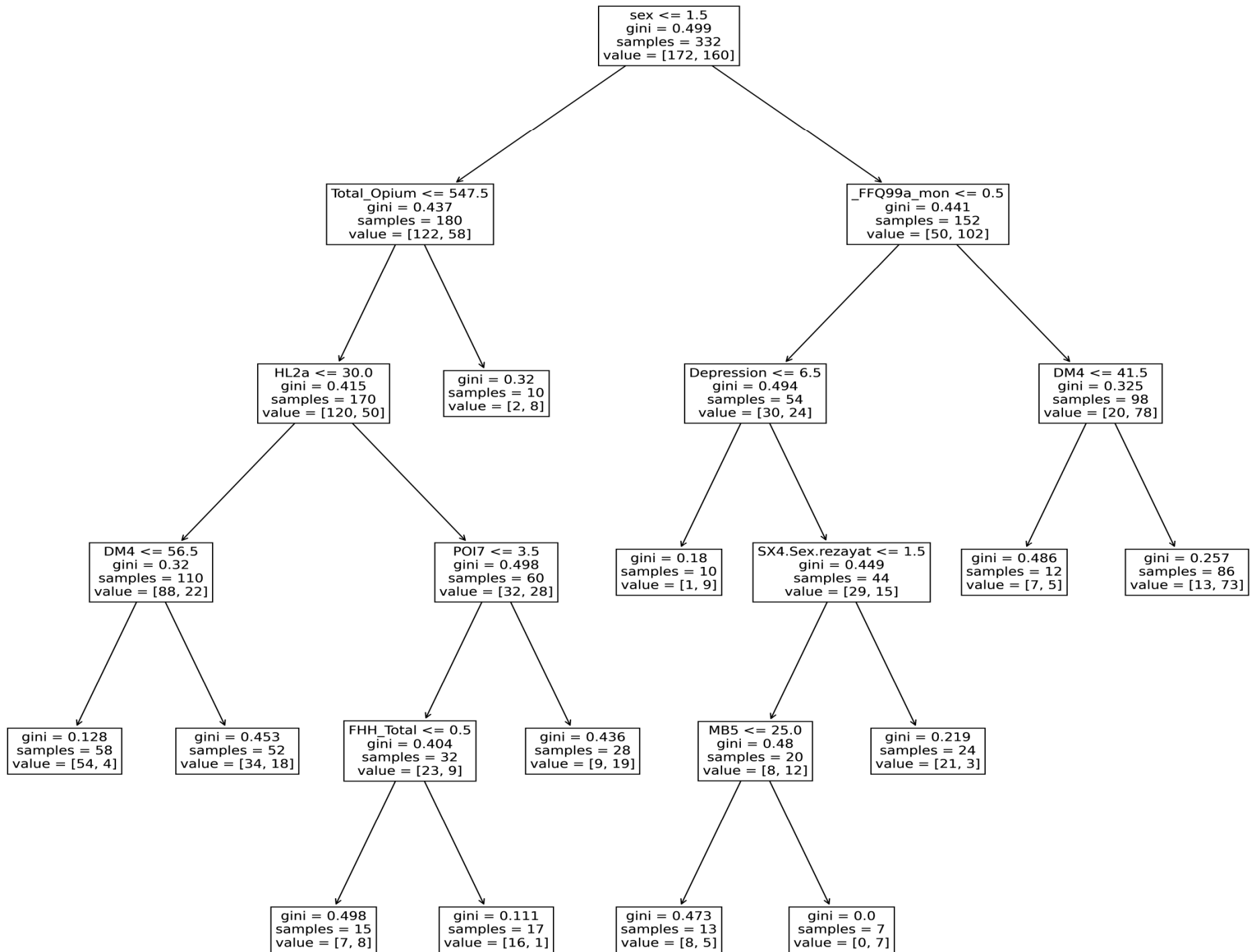
SHAP Analysis

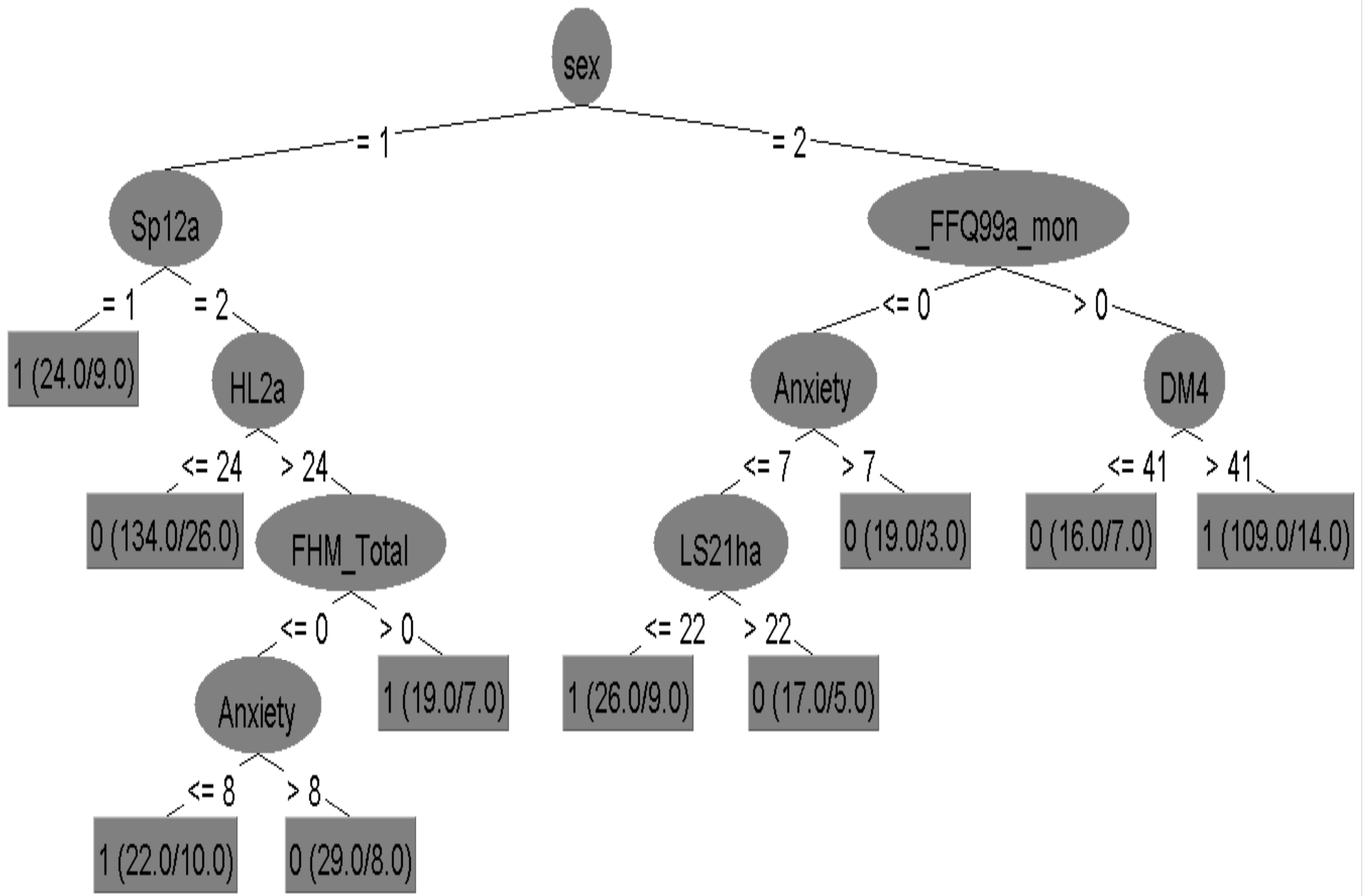


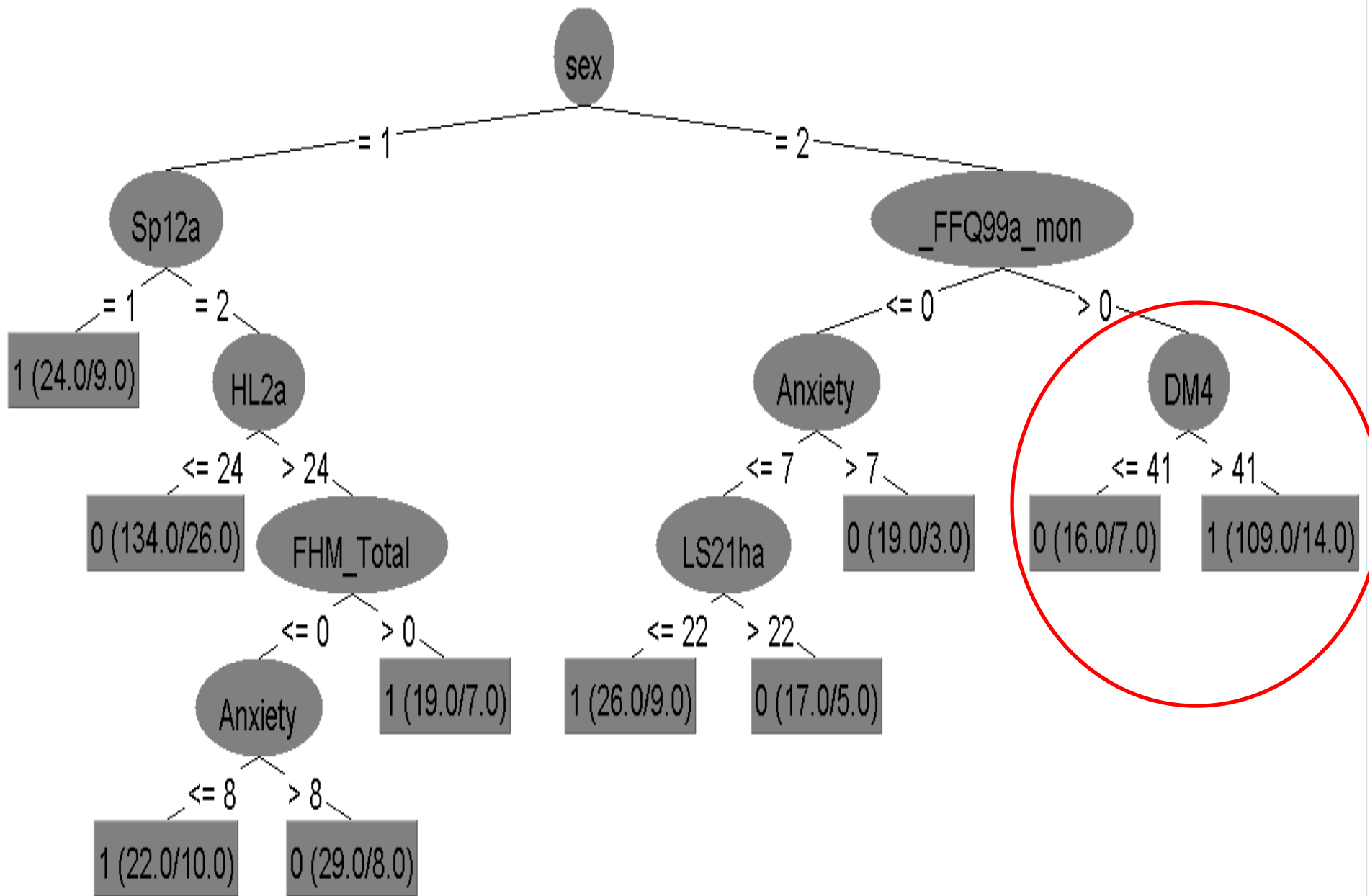
Importance of Features

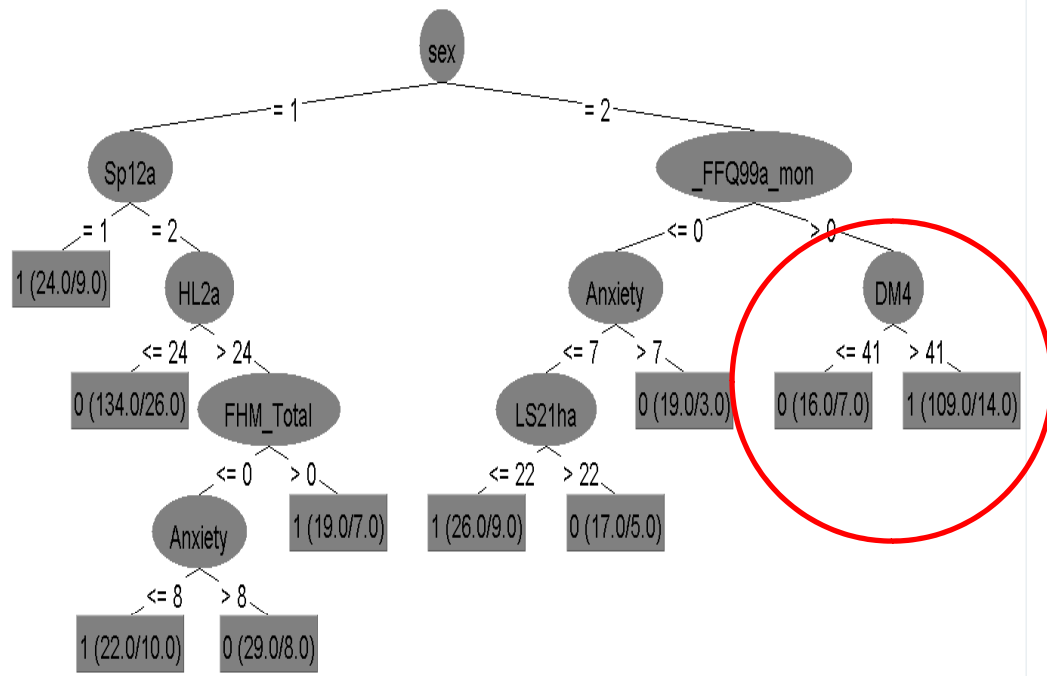
SHAP Analysis



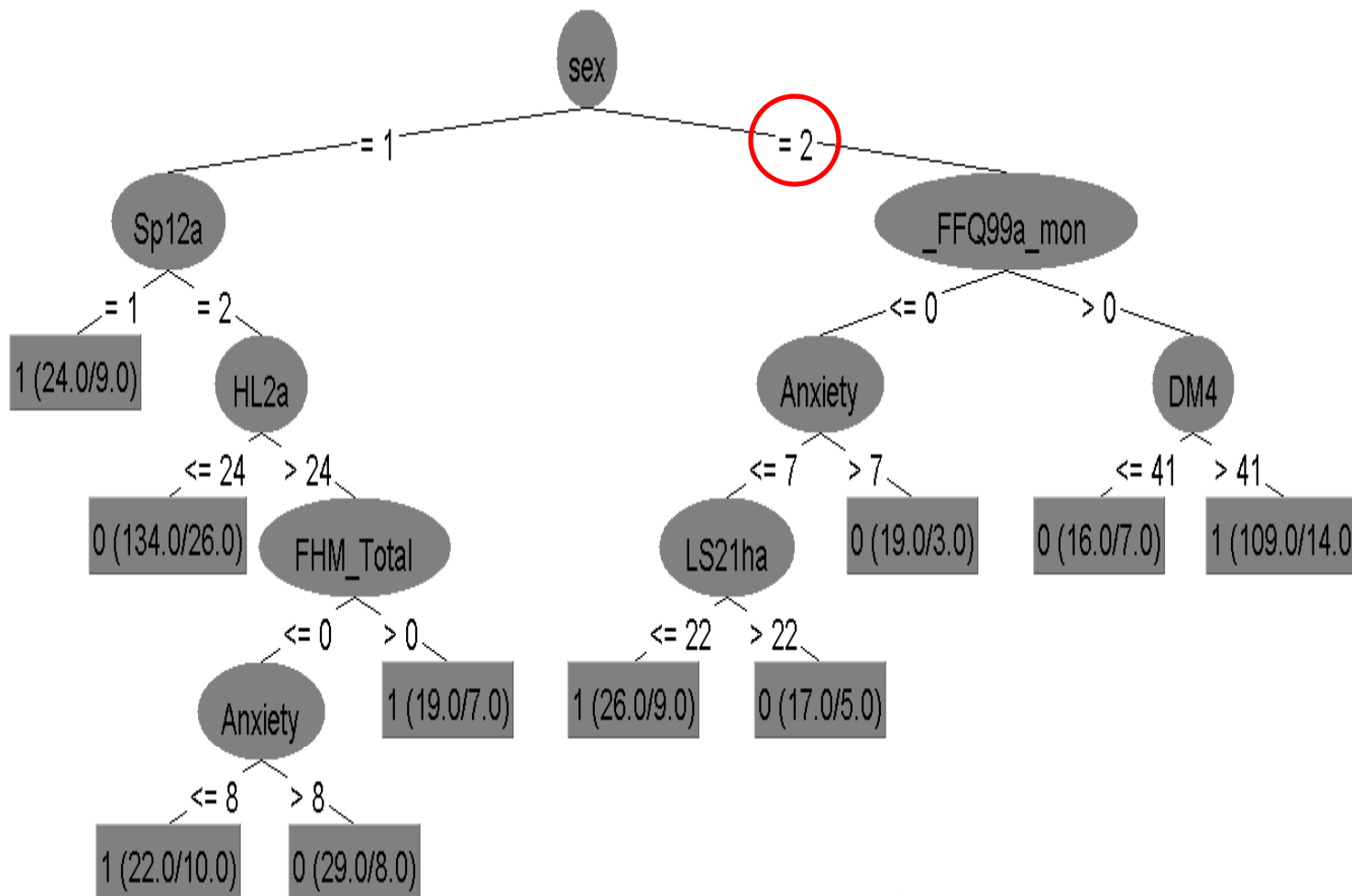




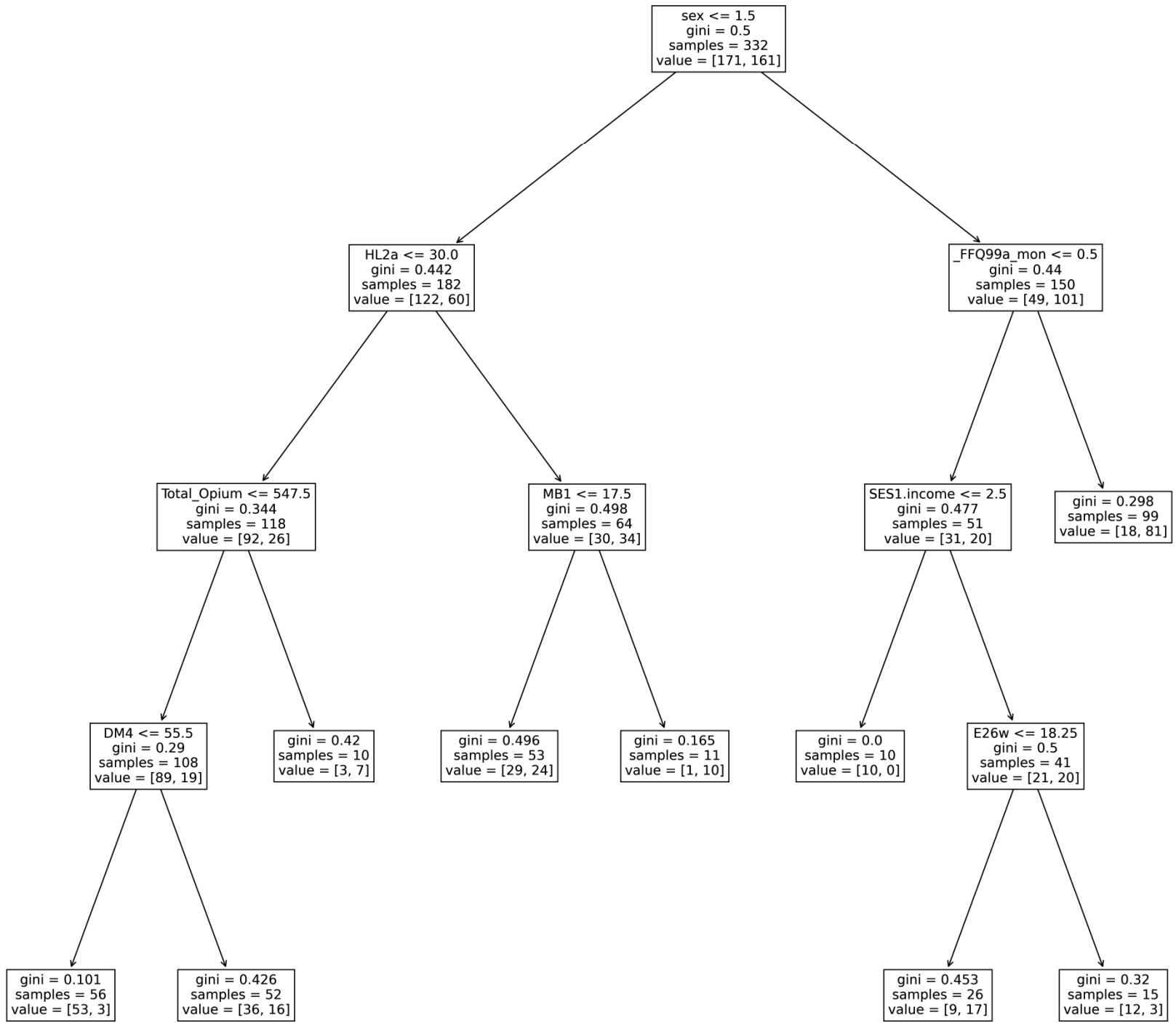




	ttest	p_value
DM4	-1.327923	1.849371e-01
Problem engagement	-1.097542	2.730439e-01
ognitive re-interpretation and growth	-1.880270	6.077483e-02
Anxiety	3.455493	6.061906e-04
Depression	2.135237	3.332944e-02
E26w	-4.201707	3.248494e-05
E26n	-5.297011	1.918409e-07
MB1	-2.023201	4.369613e-02
MB5	0.709169	4.786196e-01
MB34_Total	-4.802218	2.198318e-06
SES221	-0.904293	3.663675e-01
LS21ha	-1.684921	9.275917e-02
AR1.fasli	2.957579	3.278753e-03
LS29.Sleep.Quality	-0.633098	5.270202e-01
Total_Opium	-4.091137	5.162854e-05
Sp18a	-4.127003	4.447430e-05
HL2a	-0.870769	3.843864e-01
HD2a	-1.712770	8.750539e-02
HH2a	0.156492	8.757217e-01
FHL_Total	-0.001056	9.991579e-01
FHD_Total	0.536827	5.916764e-01
FHH_Total	0.668681	5.040729e-01
FHS_Total	0.898206	3.695988e-01
FHM_Total	-2.998210	2.879817e-03
_FFQ99a_mon	-2.023890	4.362480e-02
_FFQ100_mon	-1.704206	8.909479e-02
@_FFQ103	-2.089698	3.725624e-02
@_FFQ104	-2.119288	3.466229e-02
@_FFQ105	-1.464332	1.438640e-01
BMI	1.173332	2.413388e-01



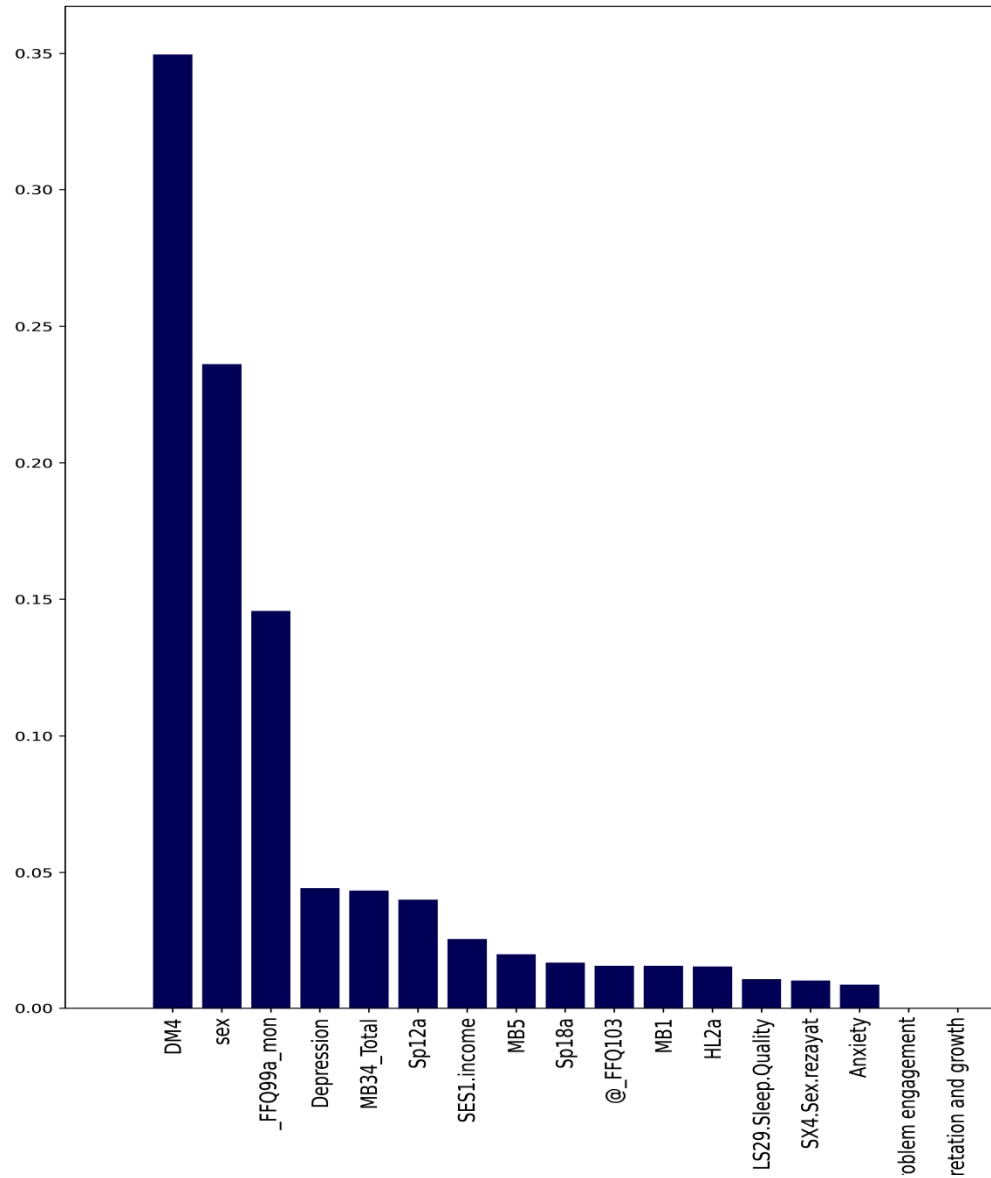
	ttest	p_value
DM4	-3.624673	0.000422
MB5	4.317094	0.000032
_FFQ99a_mon	3.177573	0.001878



Importance of Features

□ Decision Tree

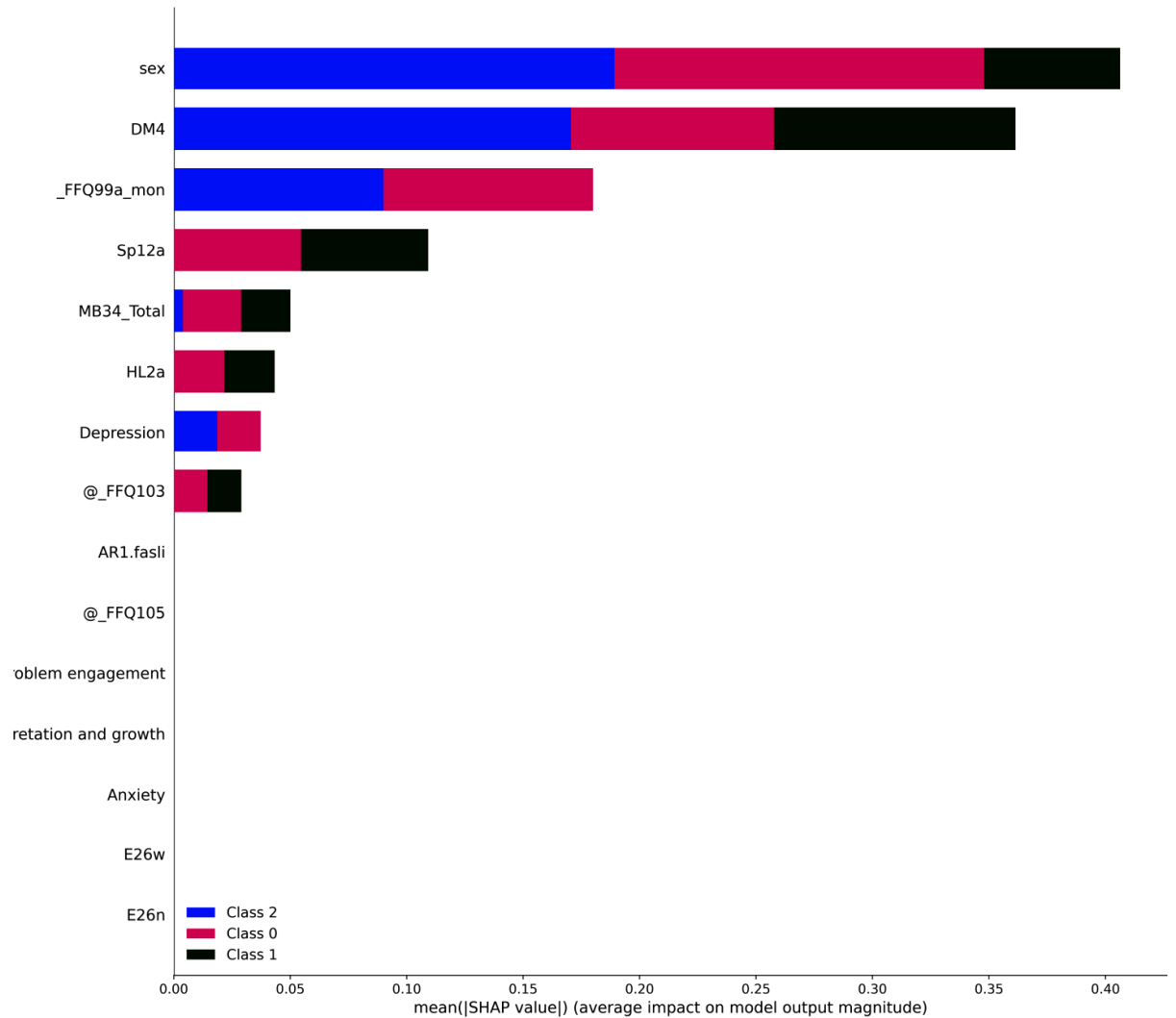
□ 3 Classes



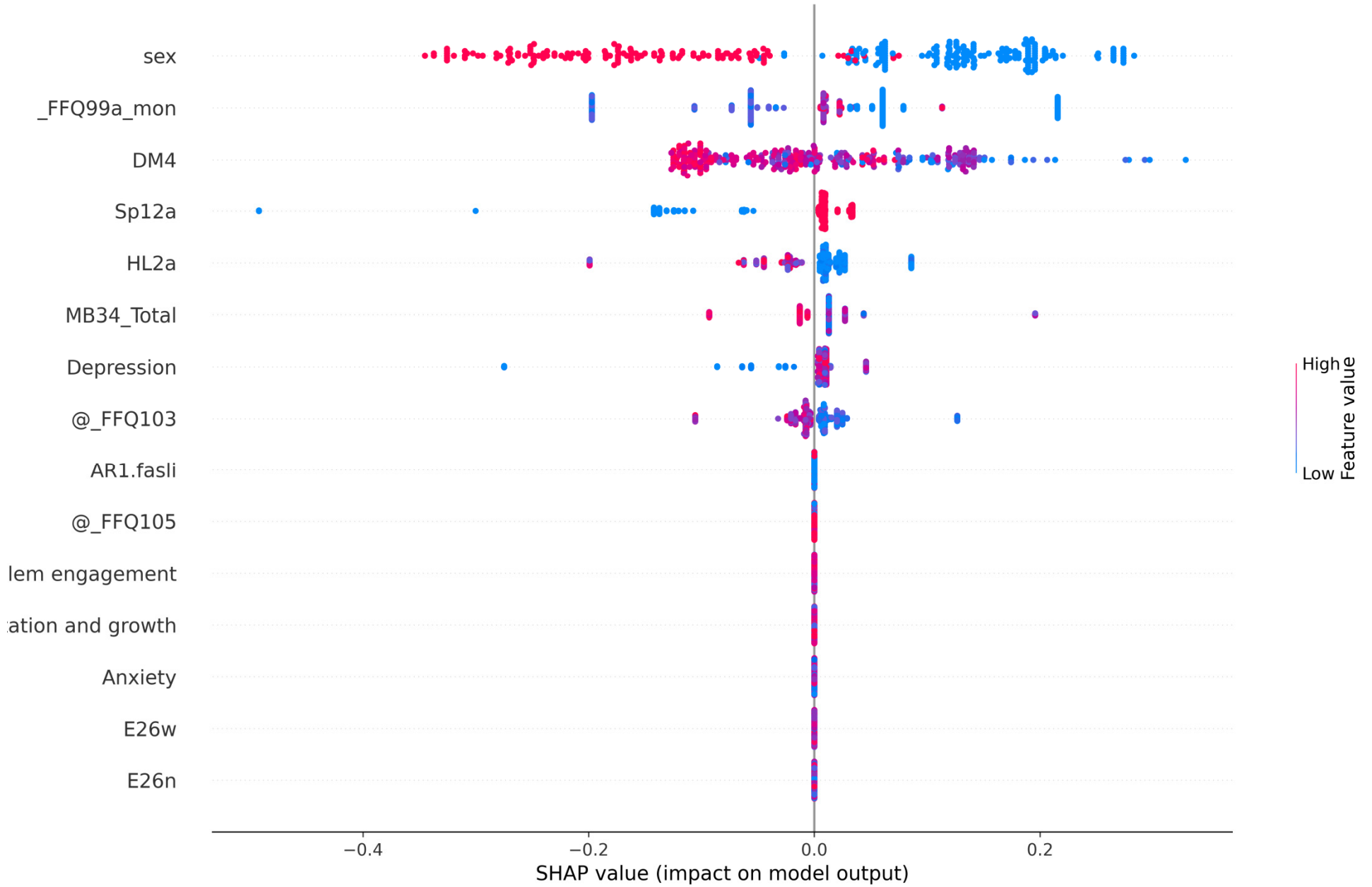
Importance of Features

□ SHAP Analysis

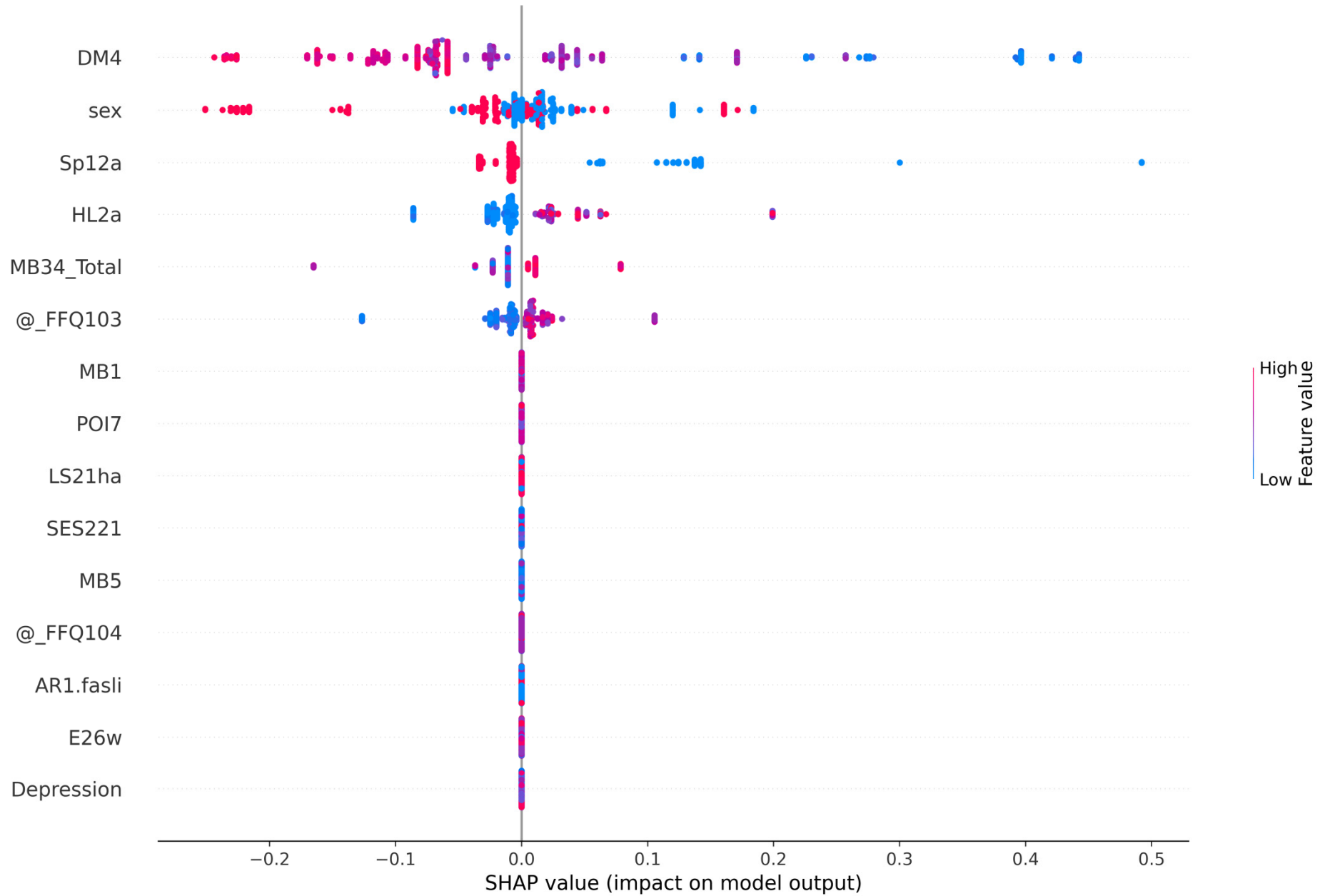
□ 3 Classes



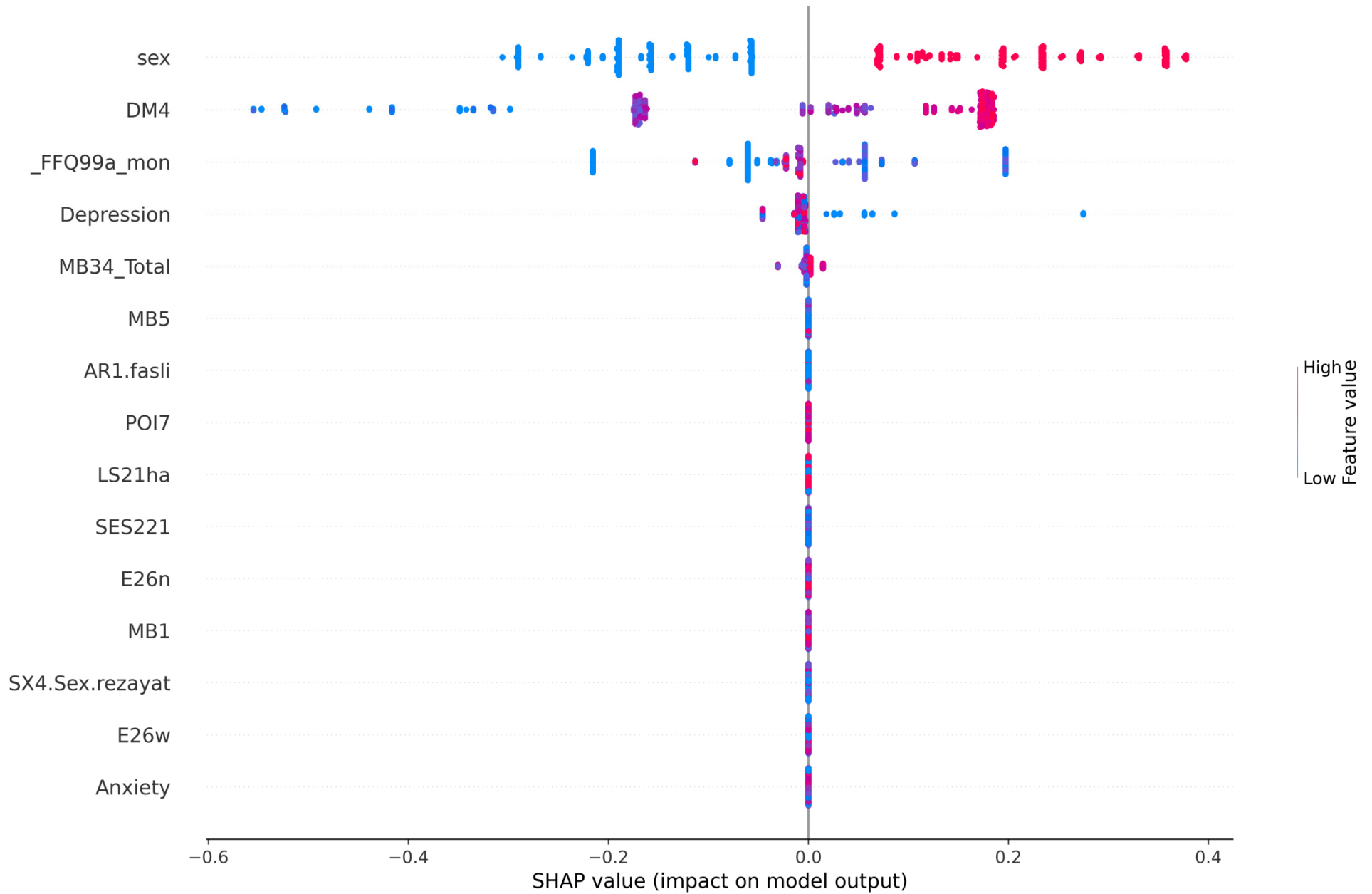
Class 0

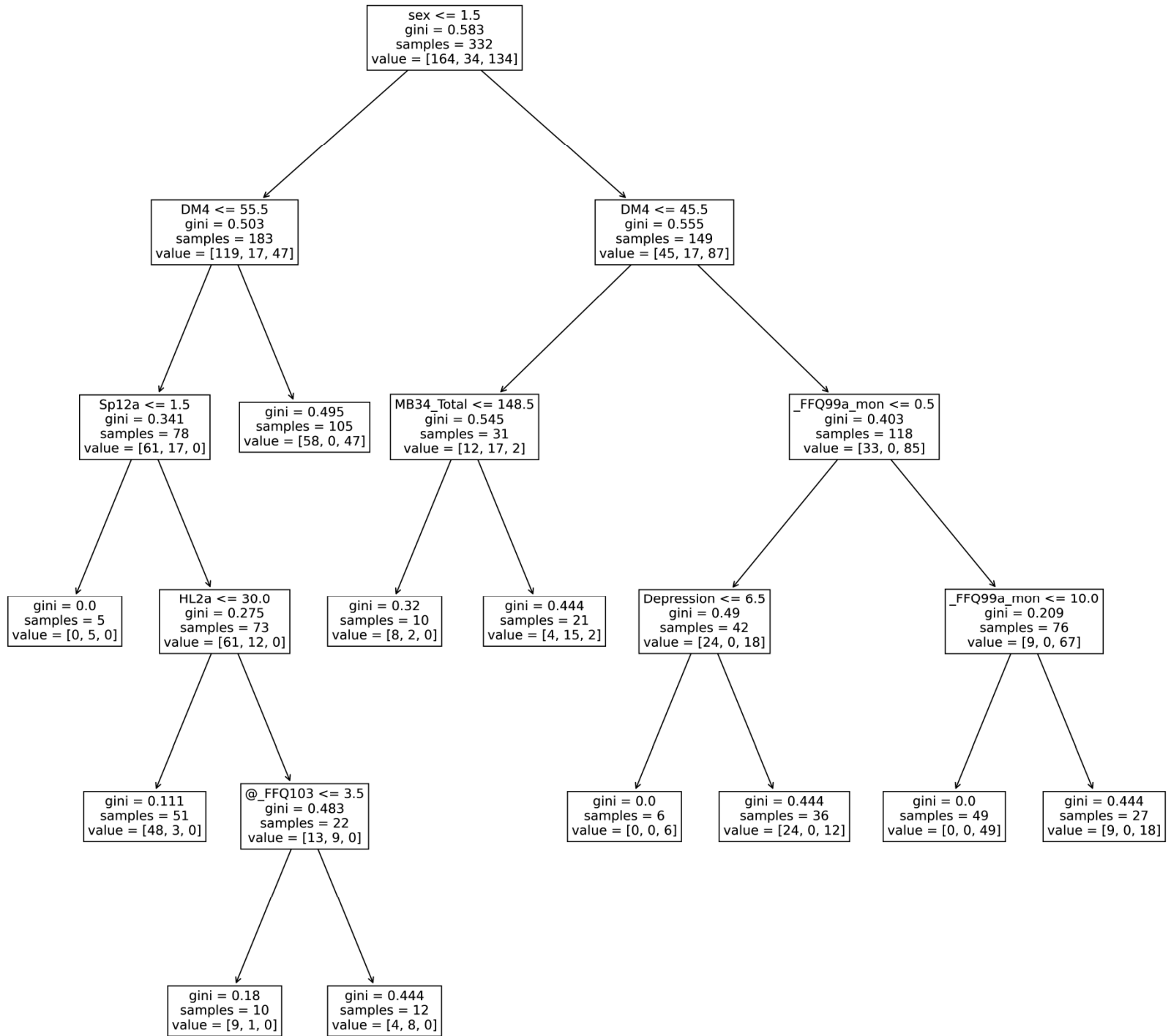


□ Class 1



Class 2



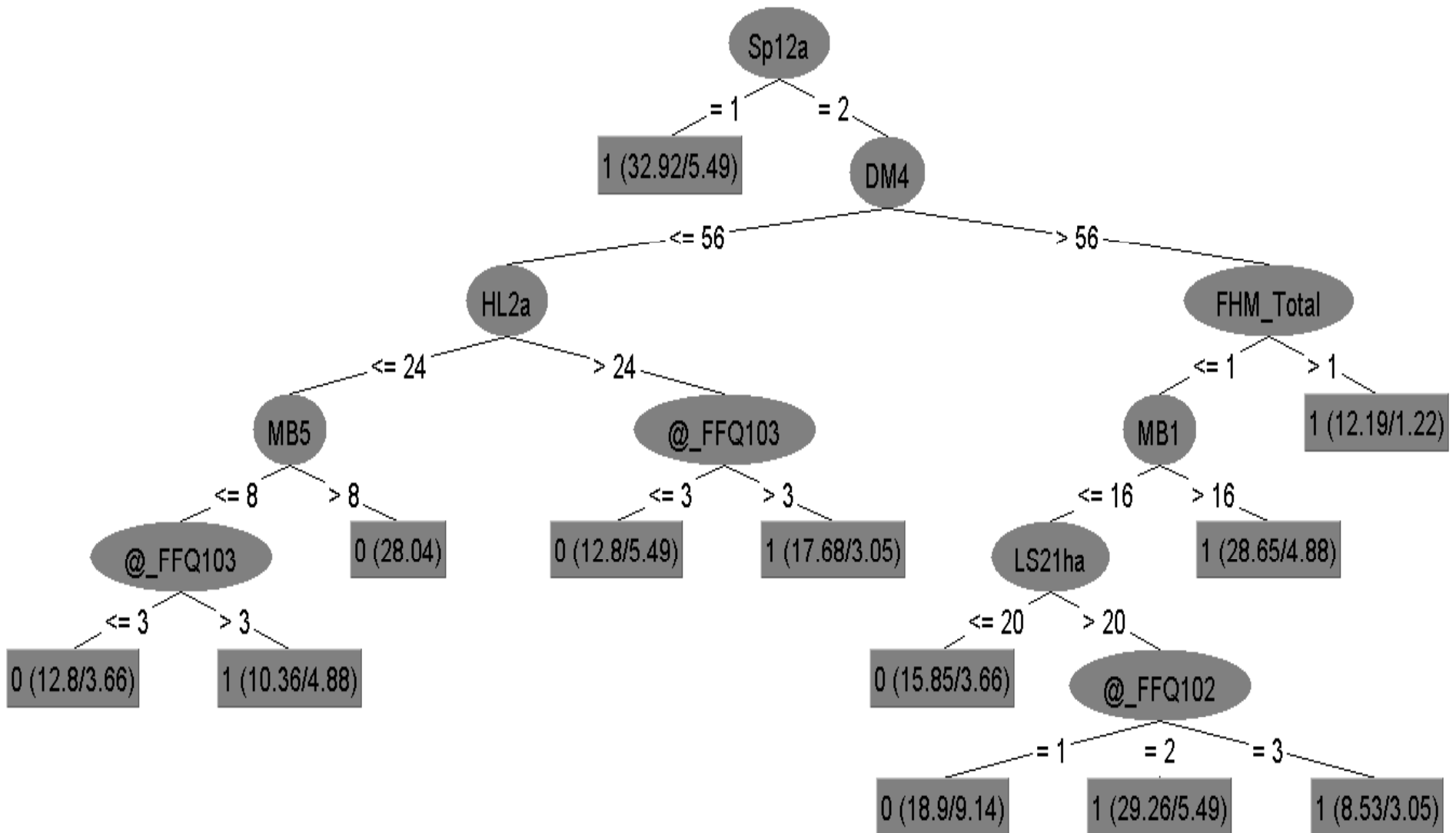


Features' Average Comparison

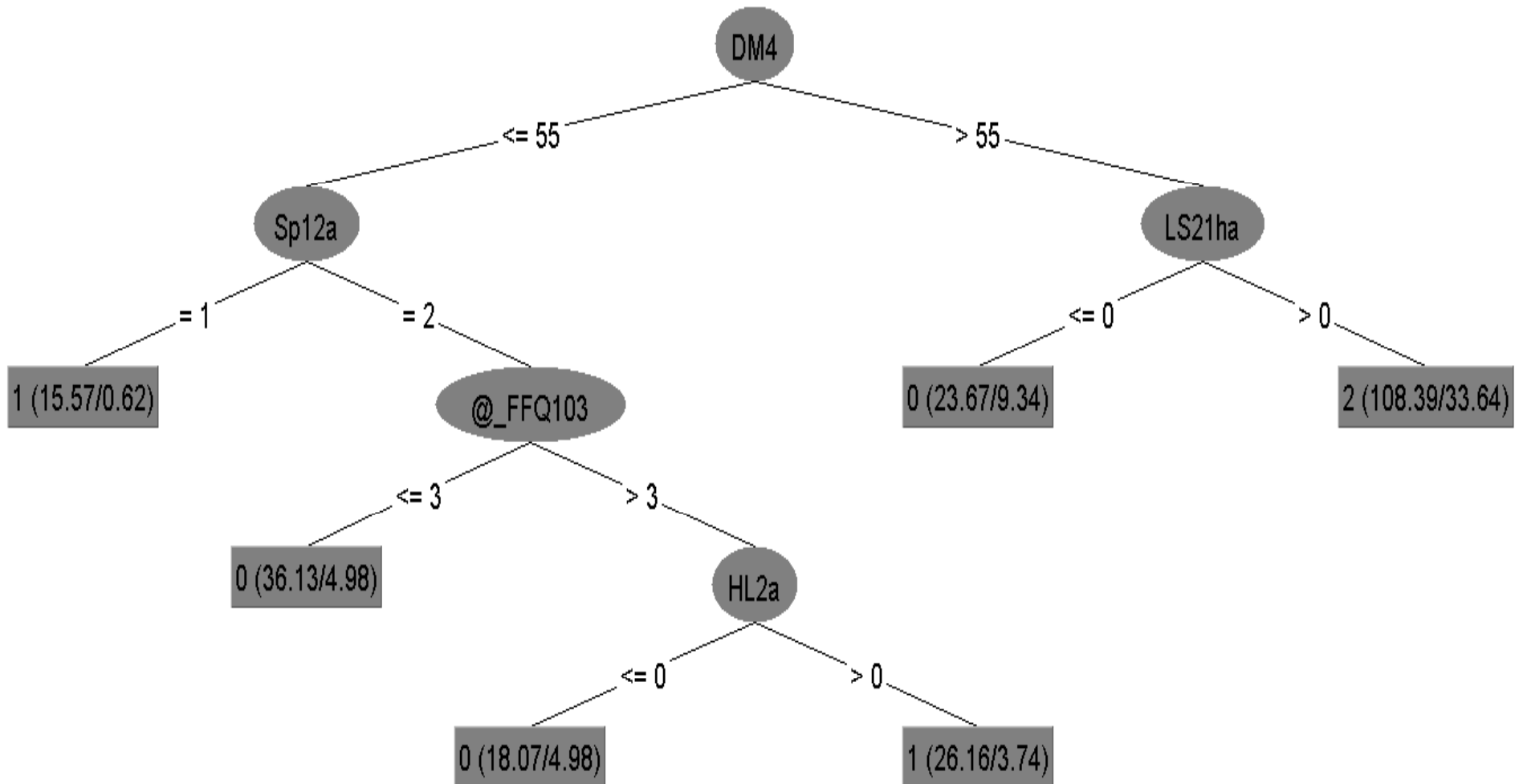
Feature Name	Women				Men		
	CAD 0 (60)	CAD 1 (20)	CAD 2 (107)		CAD 0 (155)	CAD 1 (19)	CAD 2 (54)
DM4	49.1	41.5	54.19		55.36	49.89	61.98
Problem engagement	9.17	9.6	9.44		8.94	9.42	8.63
Positive re-interpretation and growth	6.73	6.45	6.95		6.49	6.68	6.70
Anxiety	7.08	7.4	6.24		8.97	9.32	8.20
Depression	10.68	10.15	10.06		11.89	12.53	11.76
E26w	18.18	17.83	17.99		16.54	16.74	16.74
E26n	37.91	38.68	38.64		34.60	34.82	34.71
MB1	12.27	14.7	13.53		12.30	11.05	12.93
MB5	44.17	70.75	24.40		23.94	17.16	14.76
MB34_Total	143.73	178.5	158.83		75.98	72.68	67.37
AR1.fasli	0.42	0.25	0.26		0.64	0.21	0.30

Sp12a	1.67	1.65	1.56		1.94	1.68	1.83
Total Opium	219	190.3	252.96		34.33	211.37	109.11
Sp18a	5.5	6.2	8.04		1.37	3.05	2.56
HL2a	12.6	12	19.57		29.19	57.47	42.89
HD2a	7	3.6	24.11		25.39	58.11	38.44
HH2a	12	6	22.99		38.79	36	52
FHL Total	0.67	0.75	0.82		0.76	0.68	0.57
FHD Total	0.42	0.7	0.59		0.66	0.63	0.35
FHH Total	0.85	0.85	0.96		1.12	1	1
FHS Total	0.17	0.3	0.23		0.28	0.11	0.15
FHM Total	0.43	0.9	0.64		0.37	0.47	0.52
_FFQ99a_mon	6.17	8.4	7.90		5.81	5.37	7.19
_FFQ100_mon	1.28	3.25	1.78		0.55	0.11	0.39
@_FFQ102	1.9	2	1.75		1.70	1.63	1.78
@_FFQ103	4.67	5	4.66		3.19	4.21	2.85
@_FFQ104	2.87	3	2.98		2.83	2.84	2.98
@_FFQ105	5.45	4.65	6.21		5.79	6.11	6.13
SES1.income	3.73	4.55	4.23		3.24	2.74	2.83
BMI	25.88	25.89	26.79		28.20	27.65	27.38

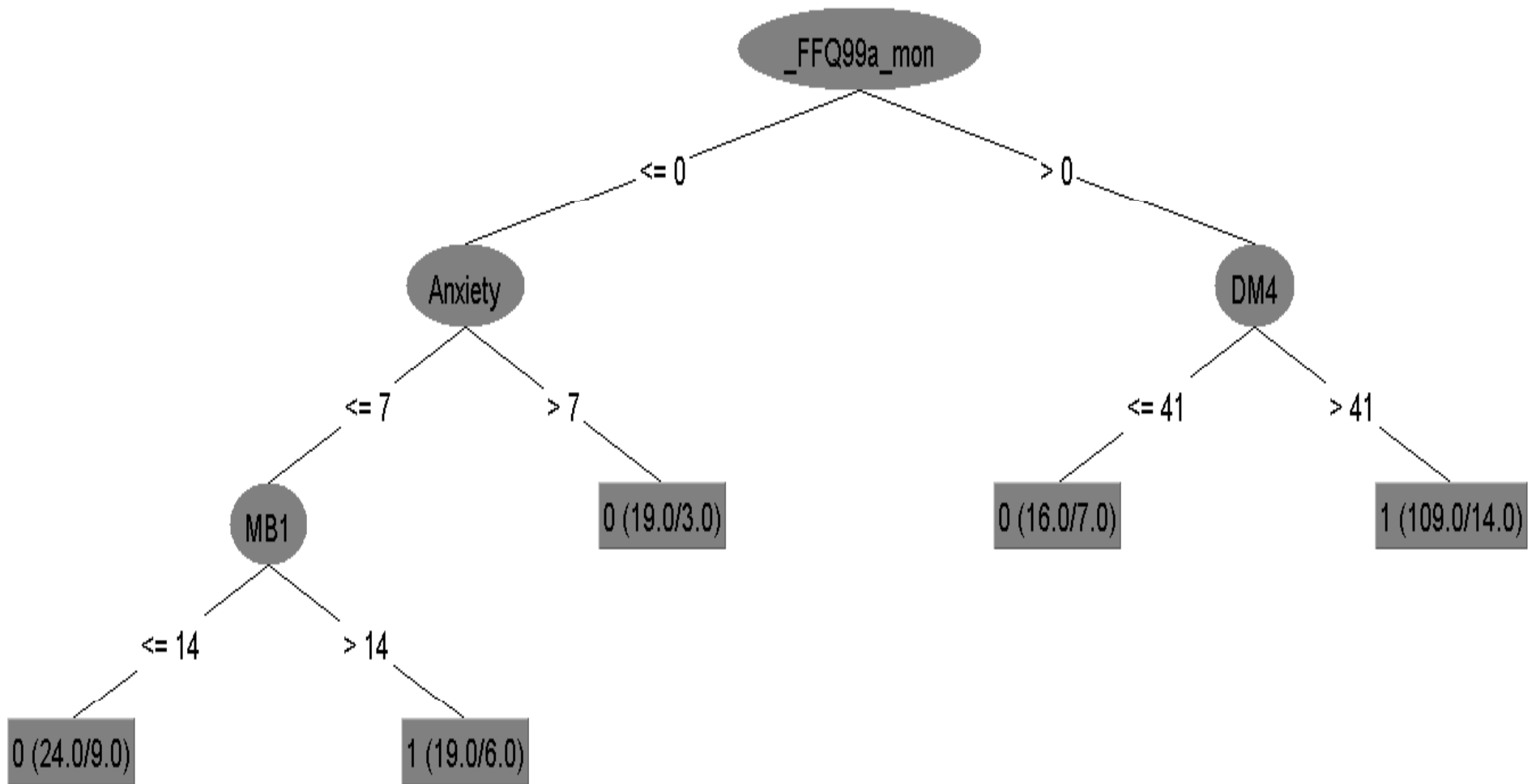
Men's Classification (2 C)



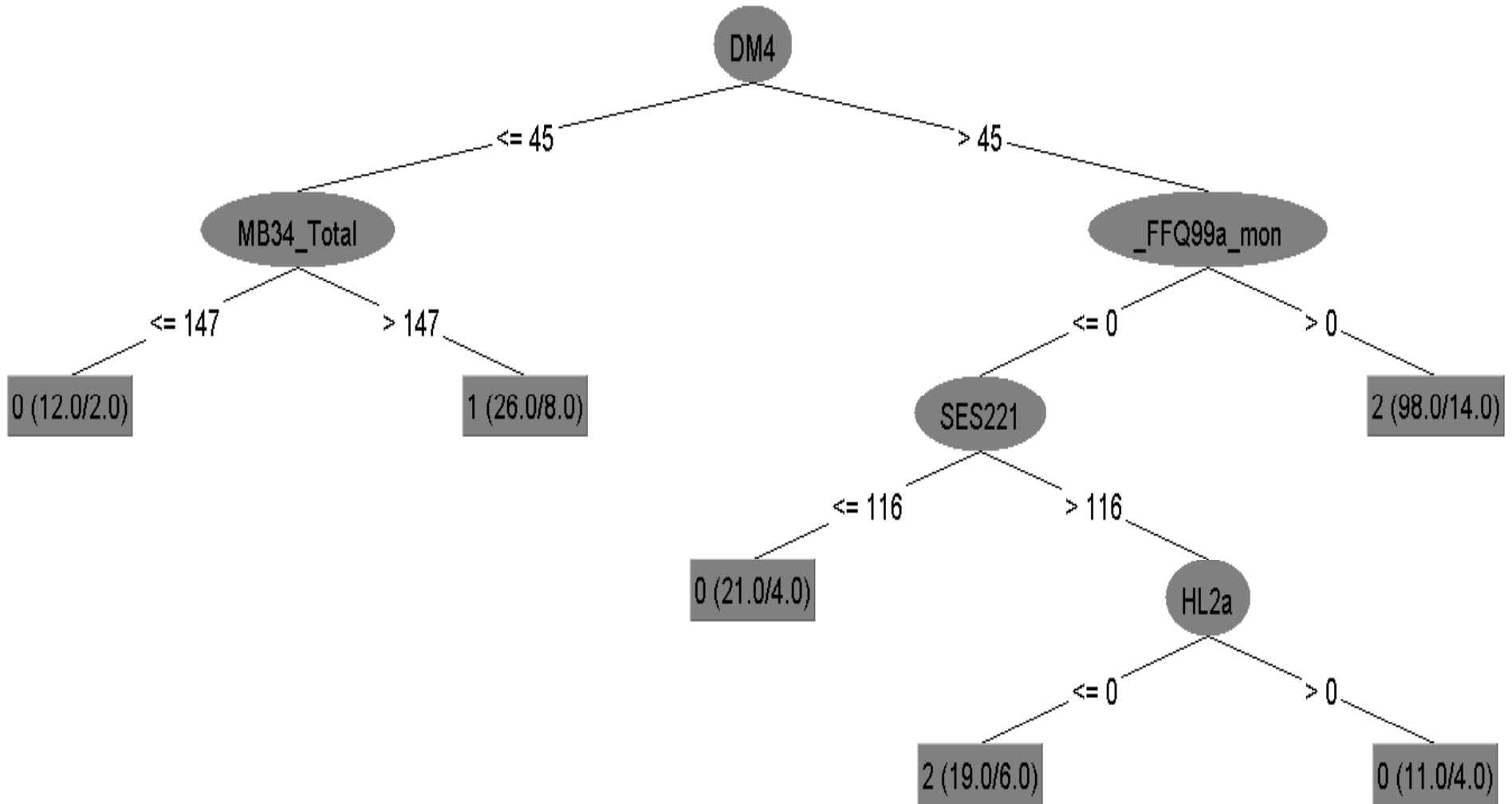
Men's Classification (3 C)



Women's Classification (2 C)



Women's Classification (3 C)





Any Questions?