

HOSEA Aim I – Report

Simon Fontaine

July 26, 2022

Calibration & threshold

Threshold	TPR	PPV	Det. prev.	Threshold	TPR	PPV	Det. prev.
10	99.51	0.14	79.20	200	47.81	0.32	16.36
15	99.19	0.14	76.68	210	45.03	0.34	14.88
20	99.05	0.15	74.69	220	41.84	0.34	13.53
25	98.74	0.15	72.82	230	39.35	0.35	12.34
30	98.46	0.15	70.91	240	36.96	0.36	11.25
35	97.93	0.16	68.92	250	34.85	0.38	10.28
40	97.26	0.16	66.88	260	32.82	0.39	9.40
45	96.77	0.17	64.81	270	31.52	0.41	8.62
50	95.82	0.17	62.74	280	29.52	0.41	7.91
55	94.66	0.17	60.67	290	27.87	0.43	7.26
60	93.86	0.18	58.59	300	26.04	0.43	6.67
65	92.80	0.18	56.51	325	22.01	0.45	5.43
70	91.54	0.19	54.44	350	18.67	0.46	4.46
75	89.93	0.19	52.41	375	16.15	0.49	3.66
80	88.77	0.20	50.40	400	14.04	0.51	3.03
85	87.40	0.20	48.42	425	12.57	0.55	2.52
90	86.03	0.21	46.47	450	10.74	0.57	2.10
95	84.35	0.21	44.57	475	9.23	0.58	1.76
100	82.91	0.22	42.70	500	8.11	0.61	1.47
105	80.84	0.22	40.87	600	4.32	0.64	0.75
110	79.26	0.23	39.08	700	2.81	0.79	0.39
115	77.89	0.23	37.35	800	1.76	0.90	0.22
120	76.06	0.24	35.67	900	1.02	0.92	0.12
125	74.34	0.24	34.03	1000	0.49	0.79	0.07
130	72.27	0.25	32.44	1100	0.28	0.74	0.04
135	69.95	0.25	30.91	1200	0.18	0.73	0.03
140	68.30	0.26	29.46	1300	0.07	0.50	0.02
145	66.34	0.26	28.06	1400	0.07	0.81	0.01
150	64.41	0.27	26.71	1500	0.07	1.18	0.01
155	62.97	0.27	25.44	1600	0.04	0.90	0.00
160	61.21	0.28	24.21	1700	0.04	1.33	0.00
165	59.39	0.29	23.05	1800	0.04	2.00	0.00
170	57.63	0.29	21.94	1900	0.04	2.94	0.00
175	55.70	0.30	20.88	2000	0.04	4.00	0.00
180	53.88	0.30	19.89				
185	52.61	0.31	18.93				
190	50.97	0.31	18.02				
195	49.49	0.32	17.17				

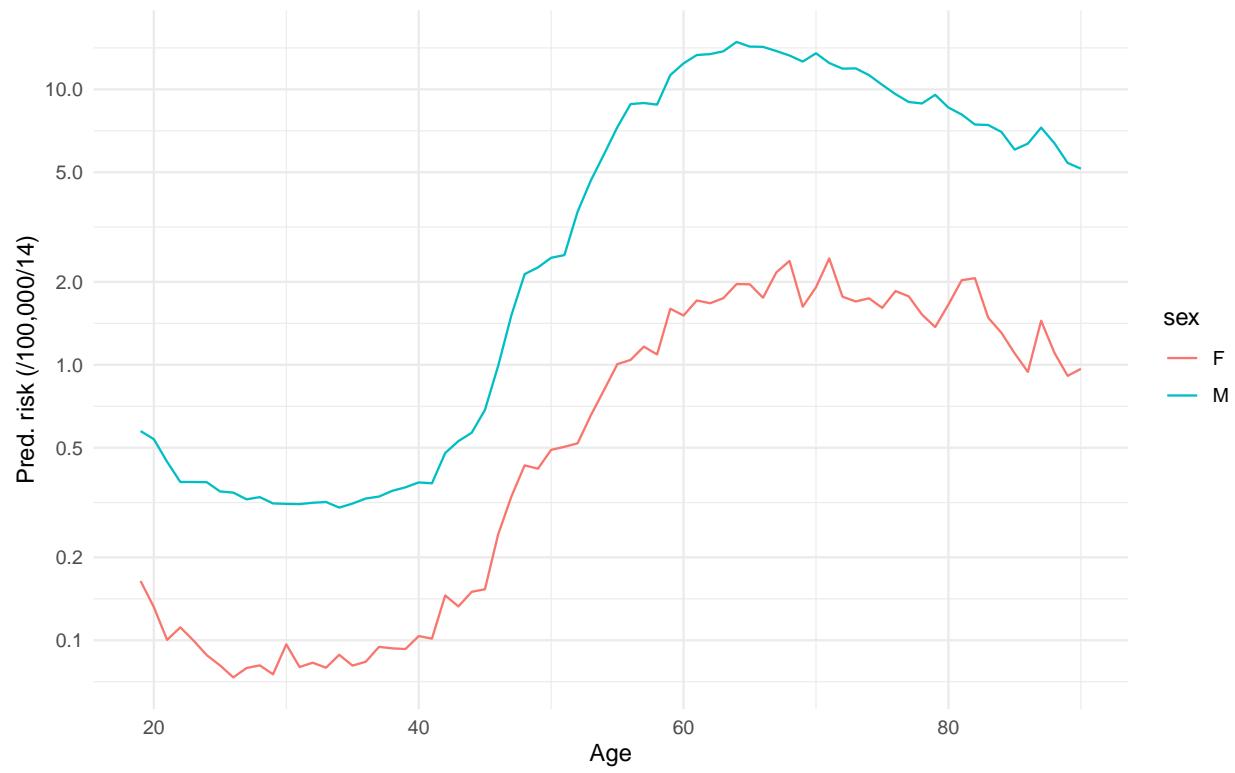
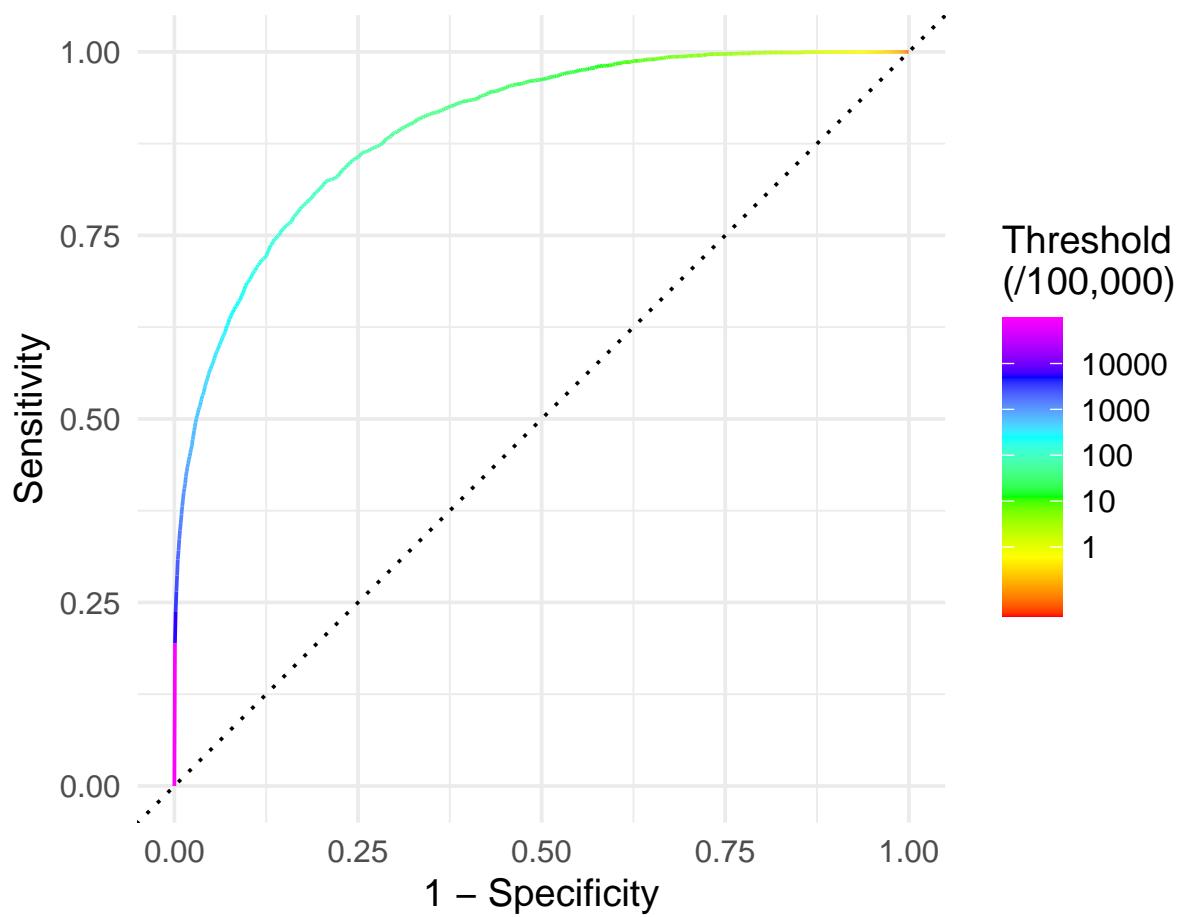
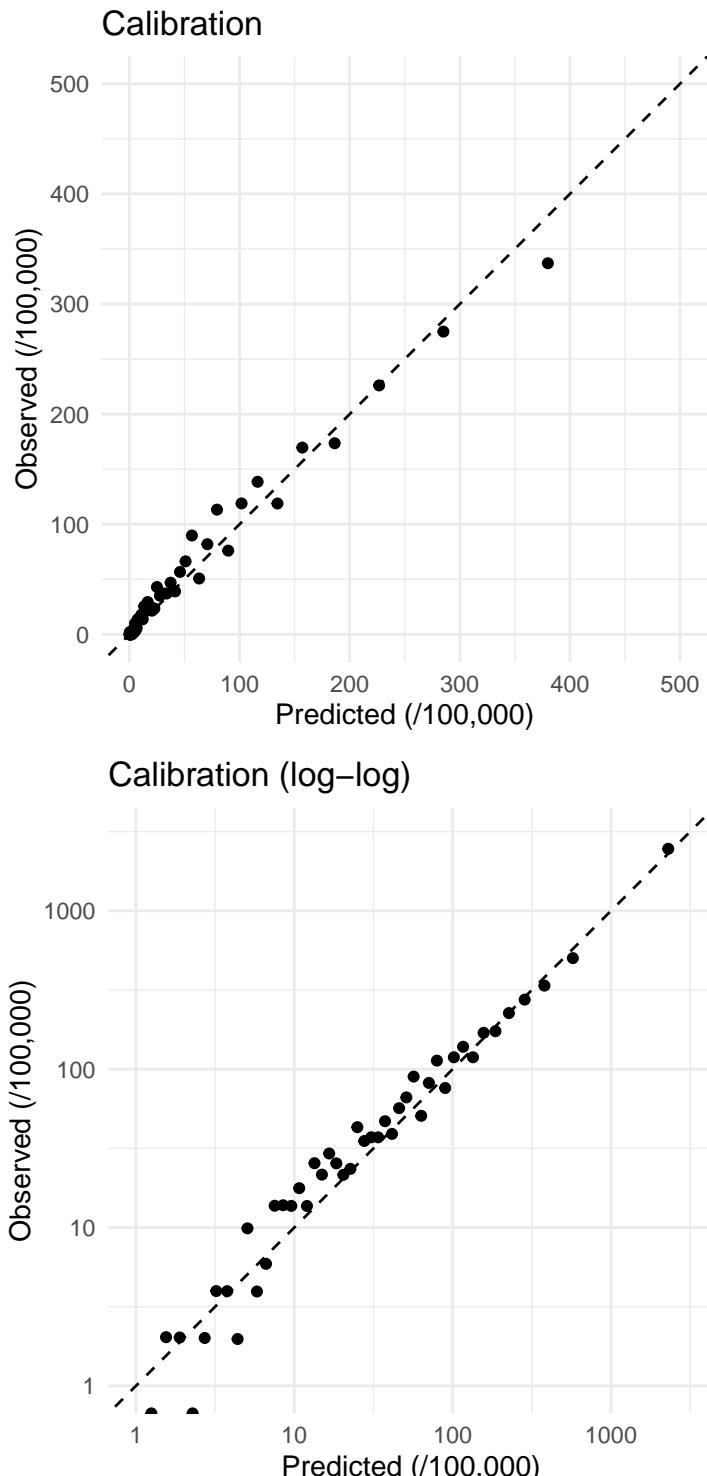


Figure 1: The average annualized risk is 7.6/100,000: for males, it is 8.2/100,000 and, for females, it is 0.7/100,000.

ROC curve (AUC: 0.898 [0.892,0.903])



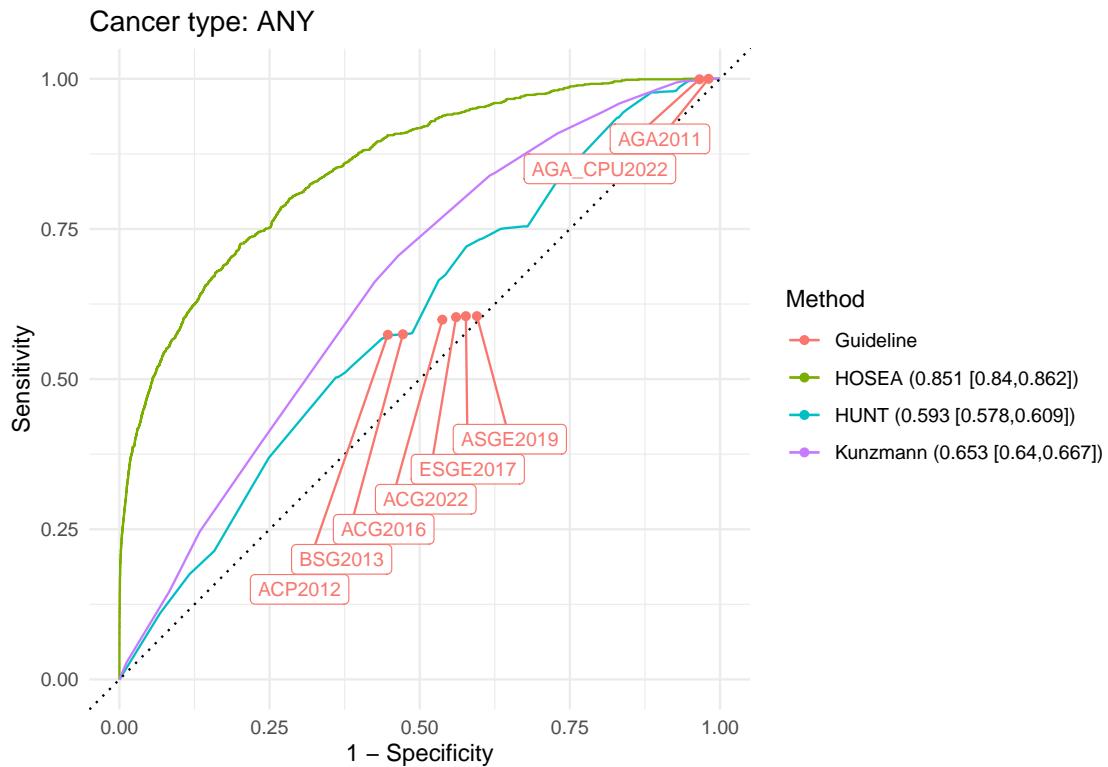


Each point represent 2% of the test data, split using predicted risk quantiles. The top plot is cropped on the right and top so we can focus on the more important region. HL test:
 "H=49.8507485754787, df=49, p=0.439294514740096"

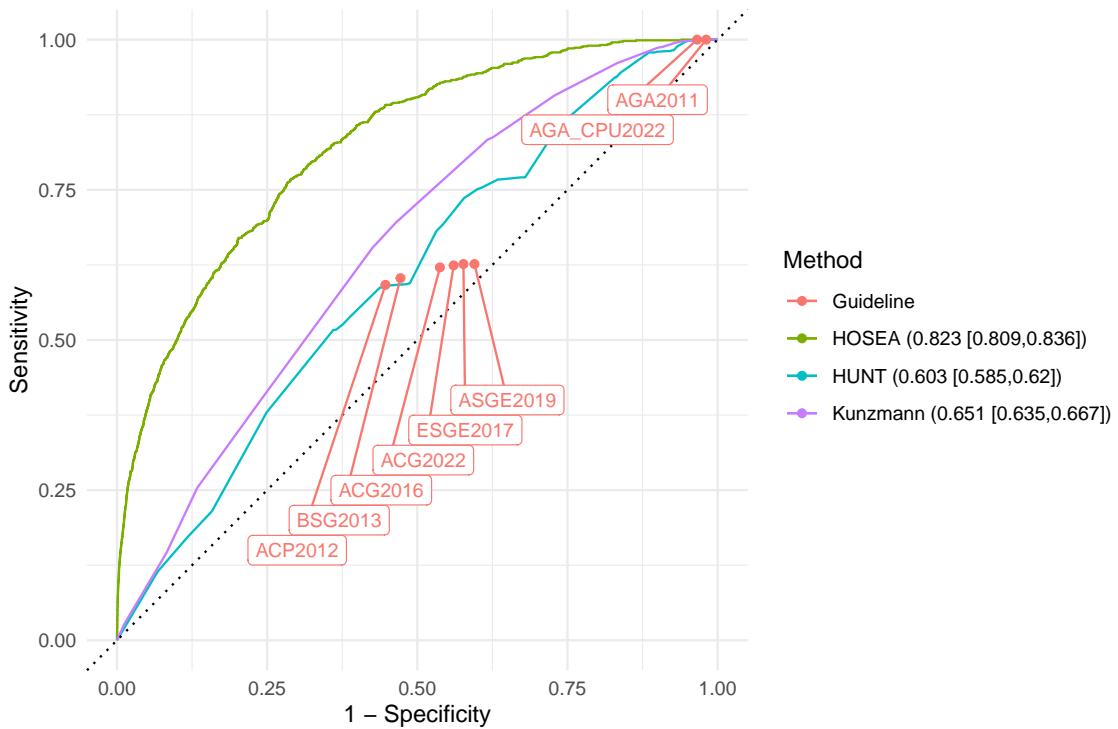
Comparison

For these comparisons:

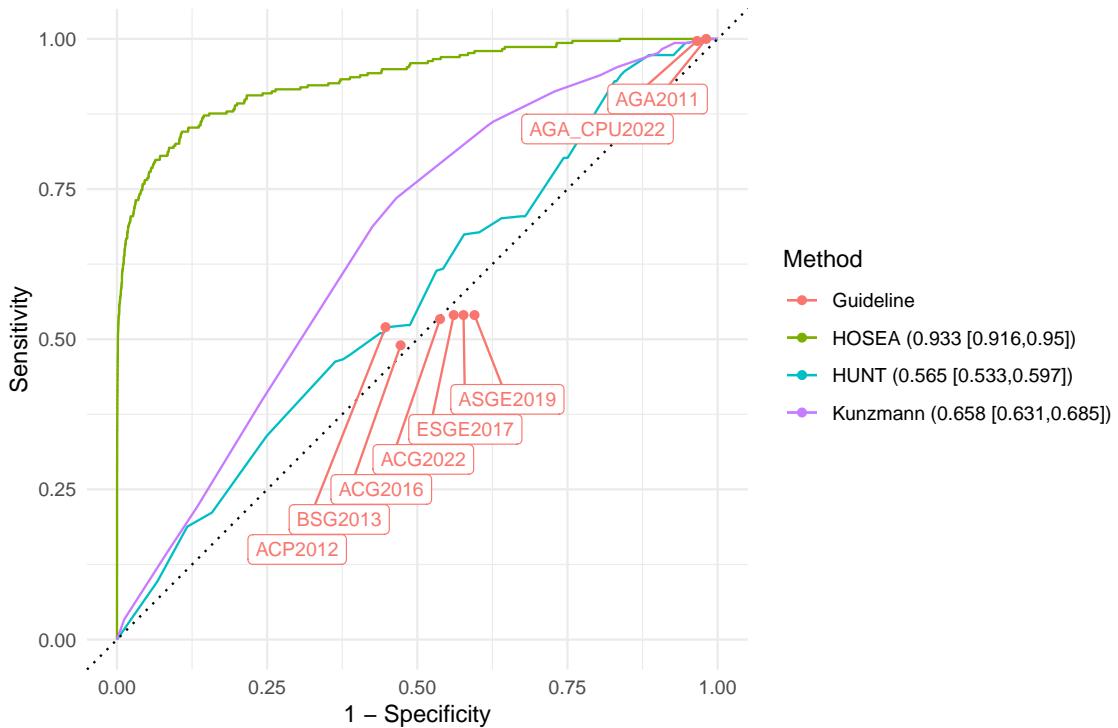
- Subset to test observations
- Filter out patients with missing information for HUNT/Kunzmann/Guidelines
- i.e., require age, BMI, race, smoking status, gerd, h2r/ppi
- 407K patients, 1192 cases (292/100,000)
- AUC + 95% CI using DeLong method



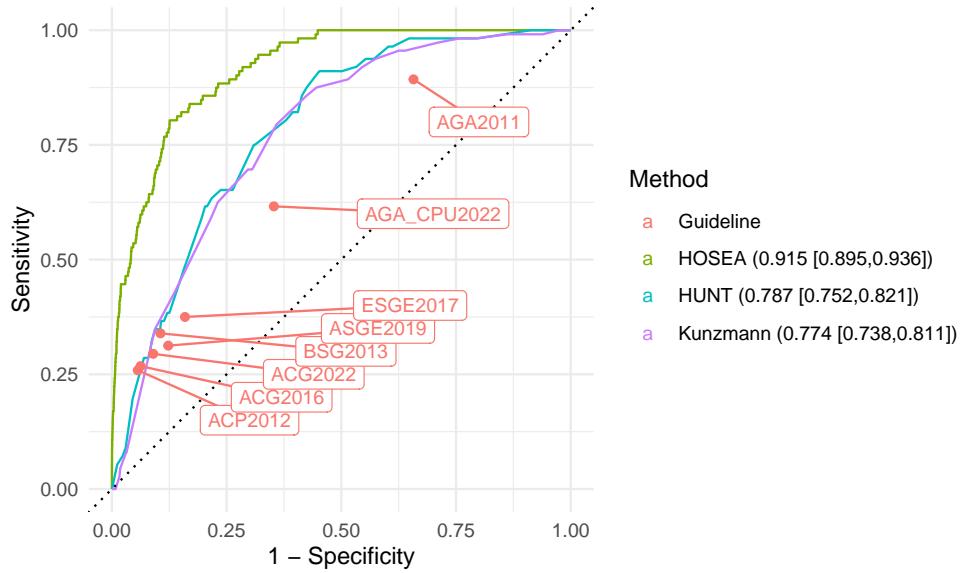
Cancer type: EAC



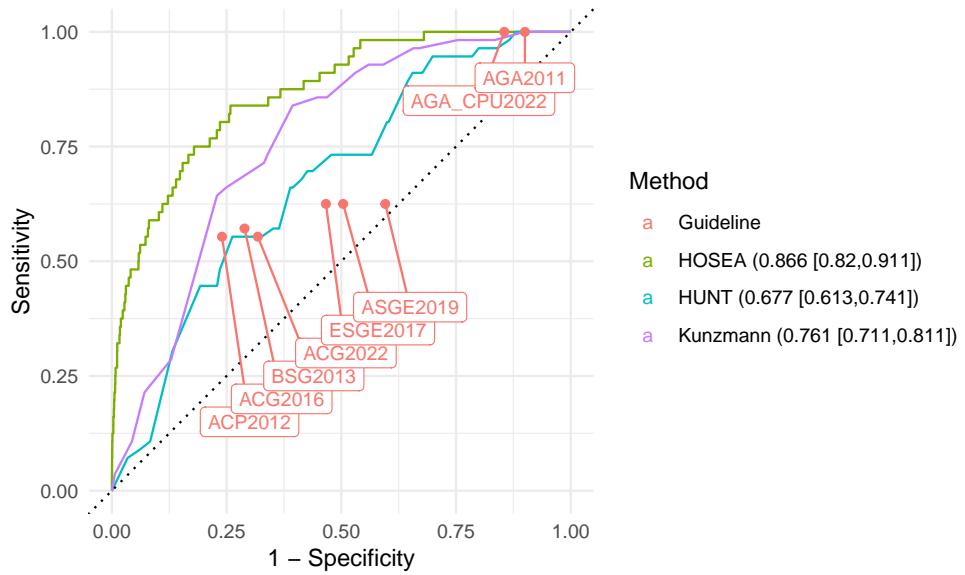
Cancer type: EGJAC



Representative sample (sex): imputed



Representative sample (sex): complete

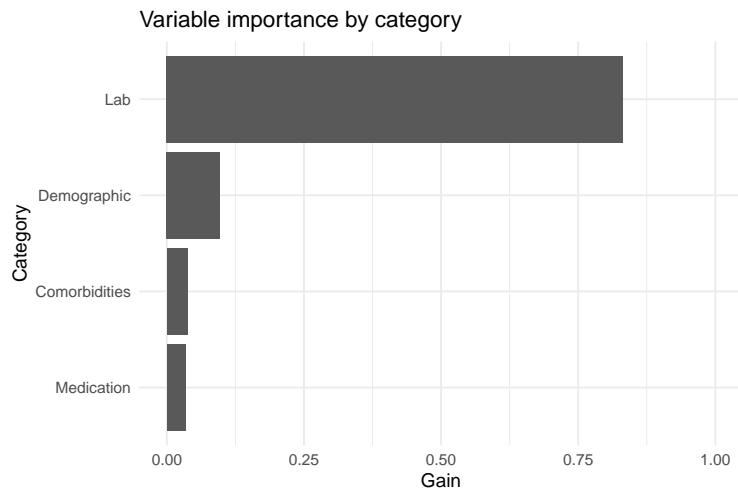


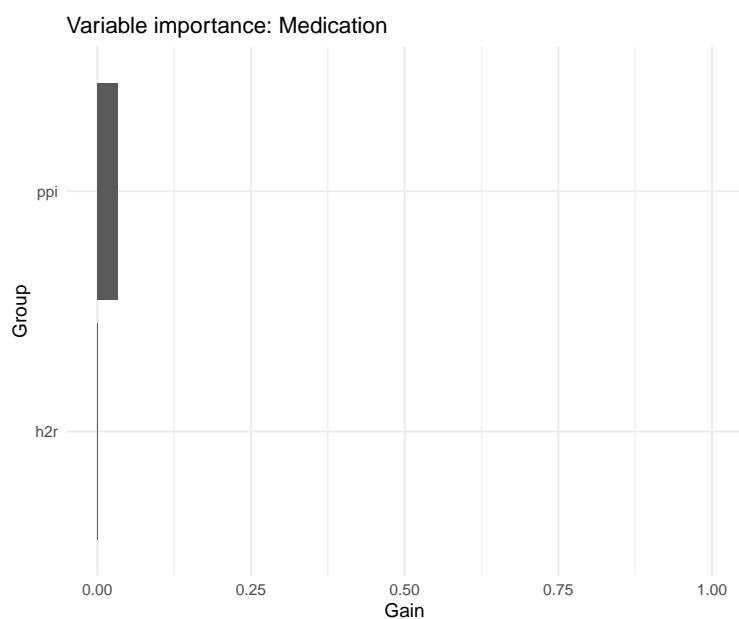
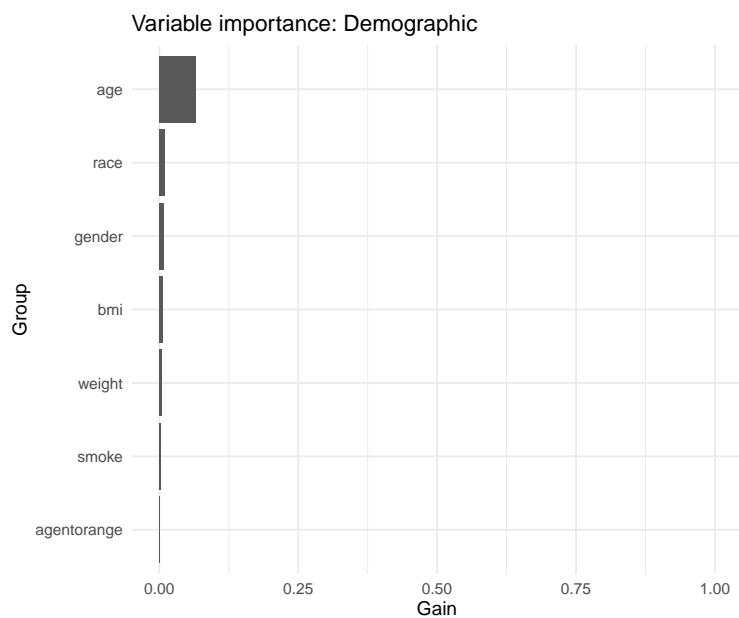
Gain Variable importance

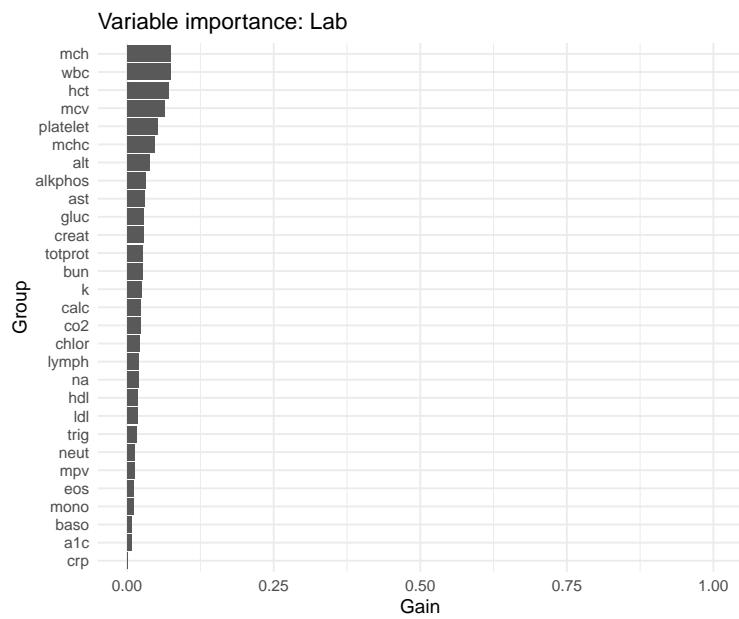
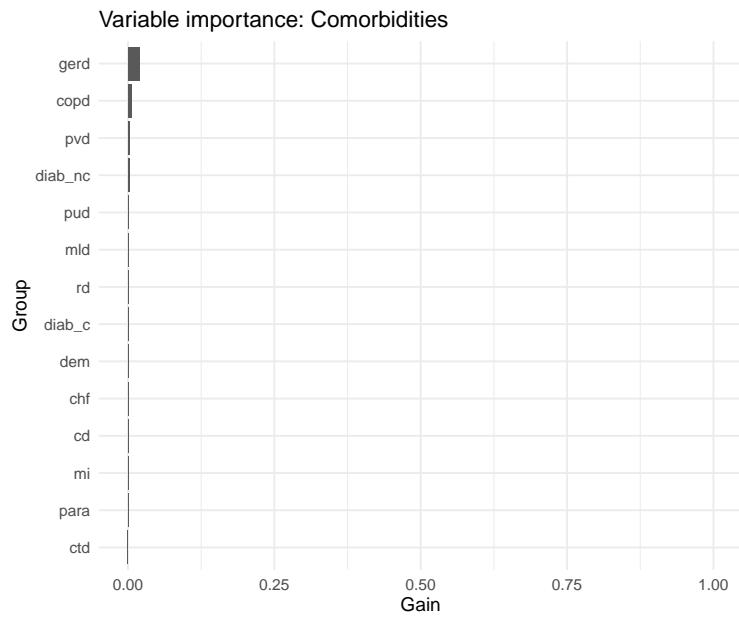
“Gain represents fractional contribution of each feature to the model based on the total gain of this feature’s splits. Higher percentage means a more important predictive feature.”

Some notes:

- These are additive, so we can compute the importance of a group of features.
- They sum up to 1
- We can only look at the relationship with the feature value (e.g., mostly positive, mostly negative) for single features







SHAP Variable Importance

- As Gain, this is additive, but does not sum to 1
- Local measure, can be aggregated using mean absolute value
- Understood as change in log-odds due to this variable (“Given the current set of feature values, the contribution of a feature value to the difference between the actual prediction and the mean prediction is the estimated Shapley value.”)
- Scale to be understood as $\beta_j x_{ij}$ in logistic regression

$$\text{logit}P[Y_i = 1] = \beta_0 + \beta_1 x_{i1} + \cdots \beta_p x_{ip}$$

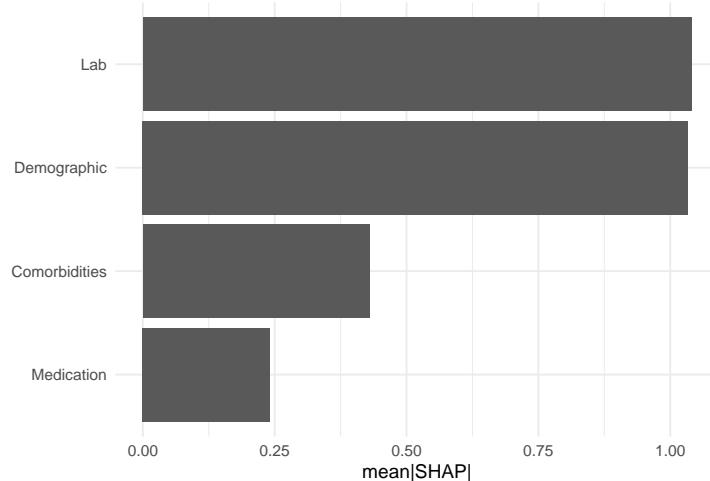
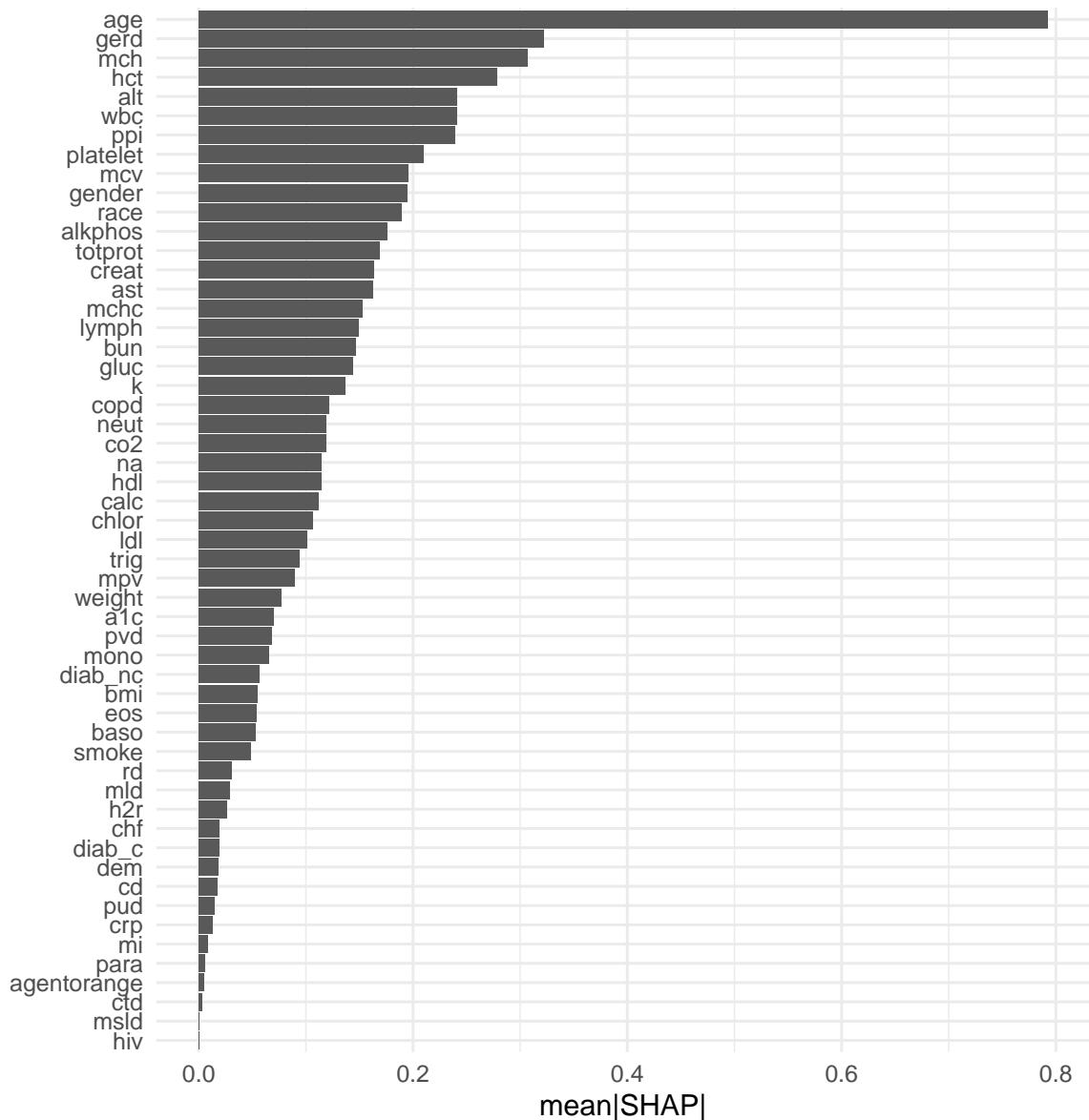


Figure 2: This has mean updated using a representative sample; previously, I was using a sample that over represented cases so age was much less important.



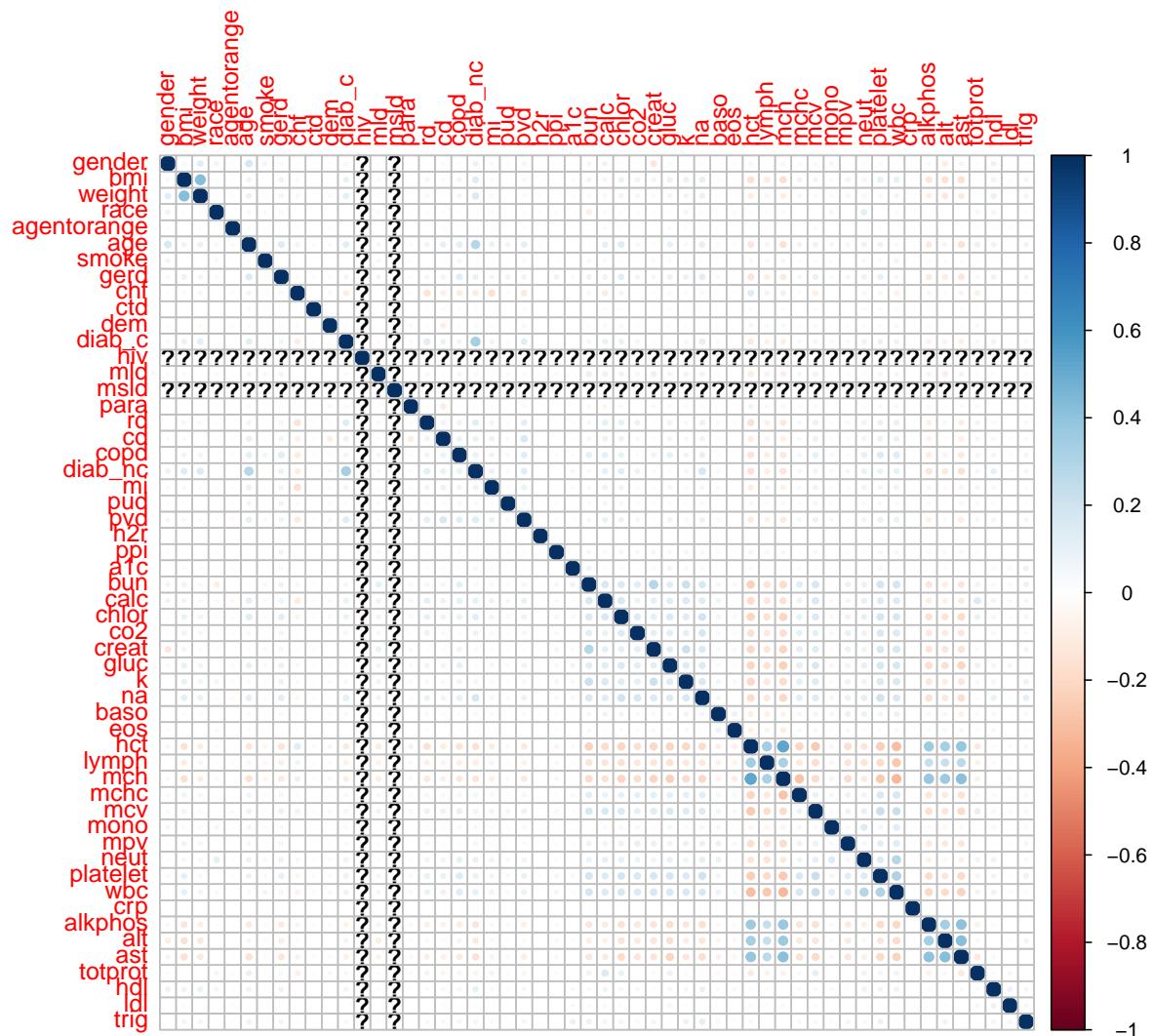
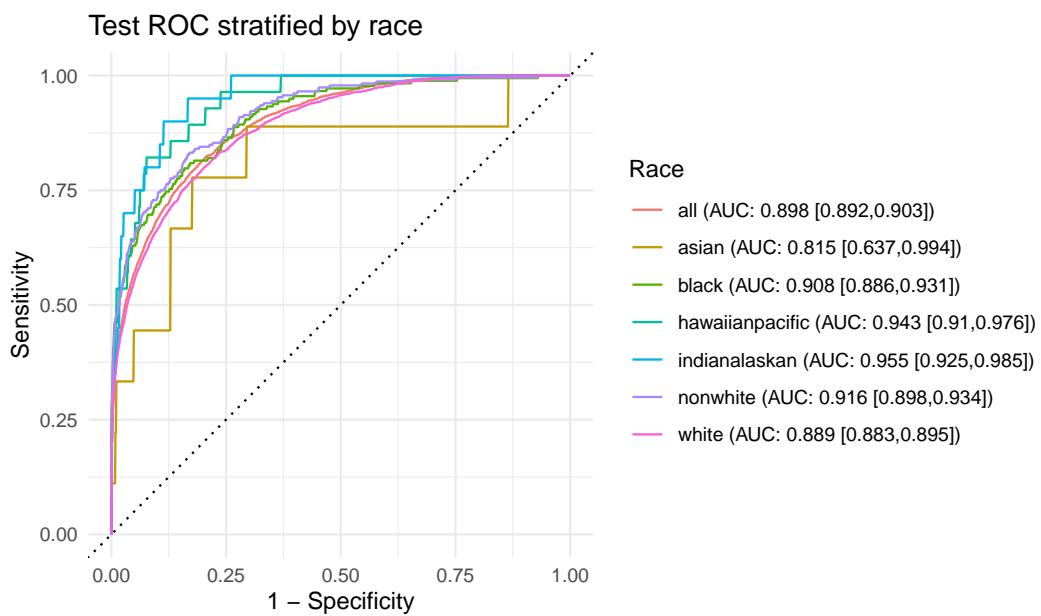
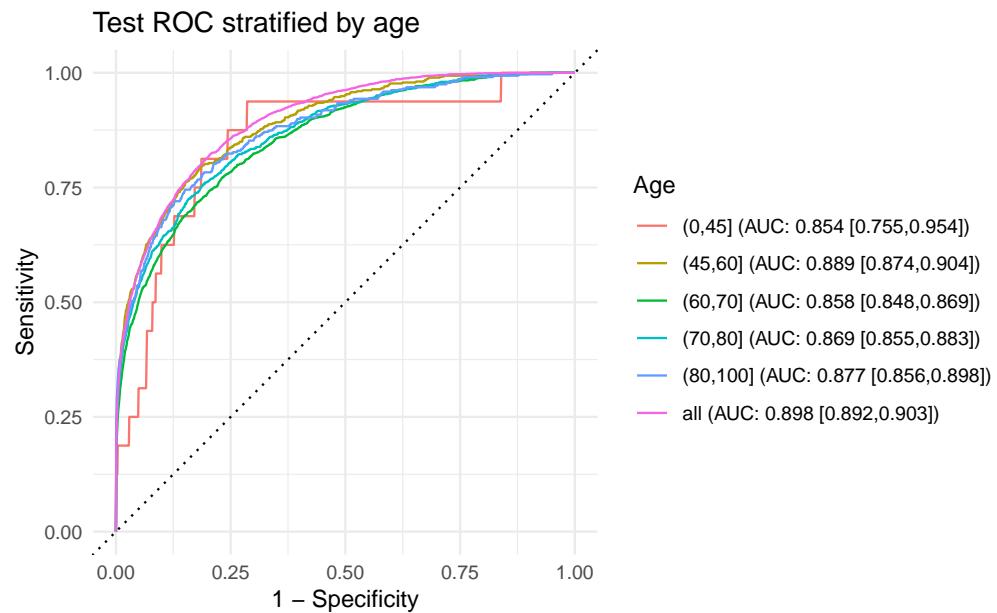


Figure 3: Correlation between group-level Shapley values. Largest absolute correlation in the following table. If two Shapley values are highly correlated, that means the underlying variables contribute essentially in the same way to the prediction and are therefore redundant.

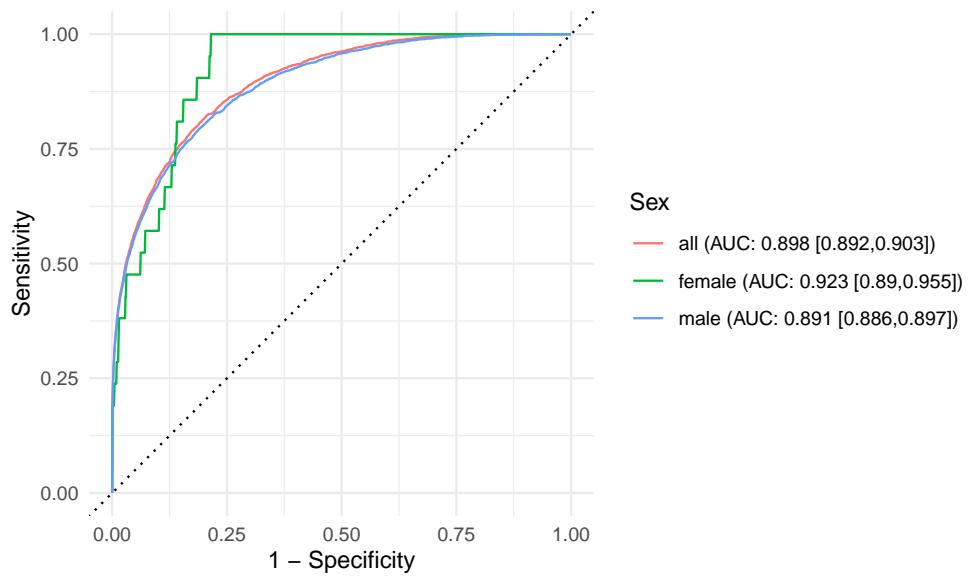
Feature pair		SHAP correlation
diab_c	diab_nc	0.31
bun	creat	0.28
bun	k	0.24
mch	mcv	0.23
bun	na	0.19
gender	creat	-0.18
diab_nc	gluc	0.18
chlor	co2	0.17
gluc	k	0.17
neut	wbc	0.17
gluc	na	0.17
creat	na	0.16
age	diab_nc	0.16
chlor	gluc	0.15
bmi_weight	diab_nc	0.15
alt	ast	0.15
k	na	0.14
gluc	mch	0.14
gender	age	0.14
cd	pvd	0.14

Table 1: Since the largest Shapley correaltion is fairly low, we can be reassured we do not have redundant variables.

Identity groups



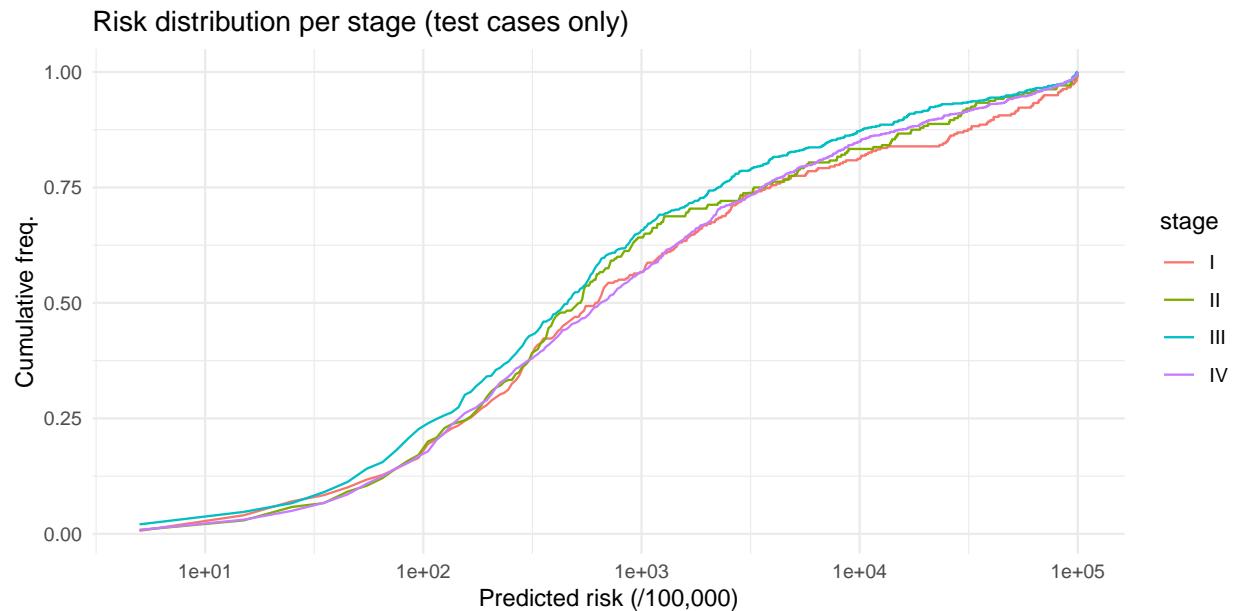
Test ROC stratified by sex



Cancer stage

Stage	Test. AUC	Nb. cases
Any	0.773 [0.765,0.78]	2810
I	0.795 [0.773,0.817]	298
II	0.798 [0.774,0.822]	240
III	0.772 [0.757,0.787]	631
IV	0.767 [0.755,0.778]	1041
I+	0.775 [0.767,0.783]	2254
II+	0.772 [0.764,0.78]	1956
III+	0.768 [0.759,0.777]	1716
IV+	0.767 [0.755,0.778]	1041

Table 2: Interestingly, it seems slightly harder to predict late stages?



Years prior

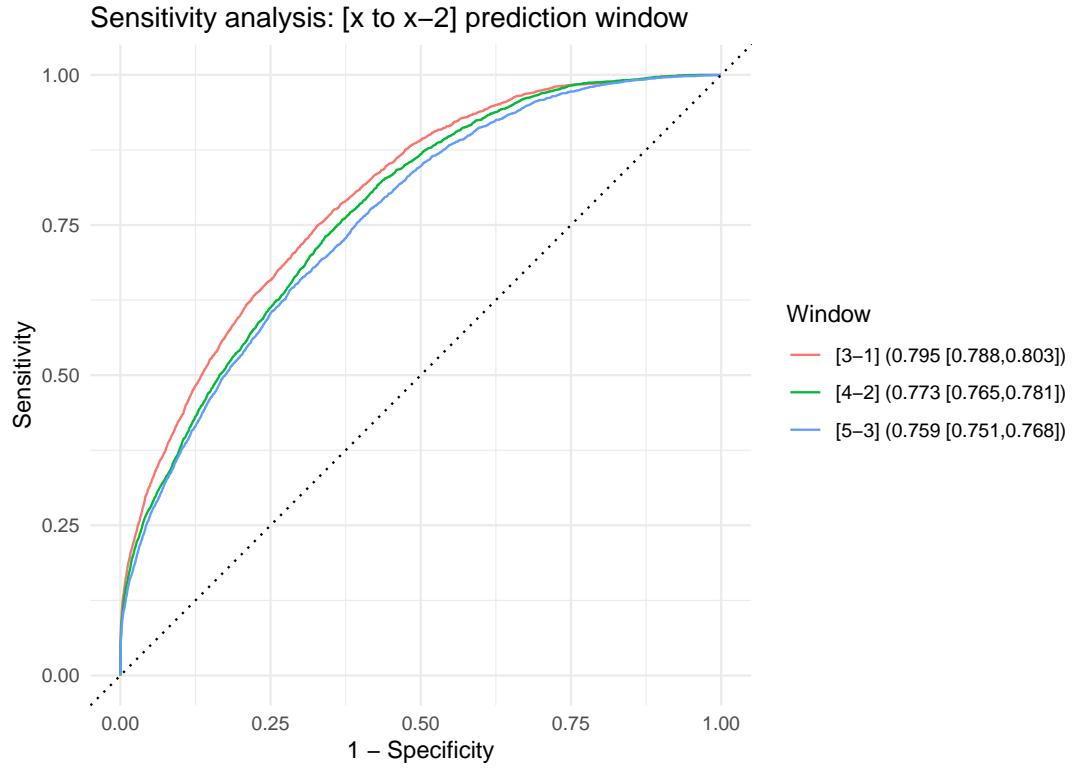


Figure 4: In this experiment, all three test set utilizes the same amount of data (2 years), but we predict farther in the future (1yr, 2yrs or 3yrs). As expected, we do worse and worse as we predict further in time, but only by a small margin. This seem to indicate are predictions are somewhat valid beyond the 1yr window we used.

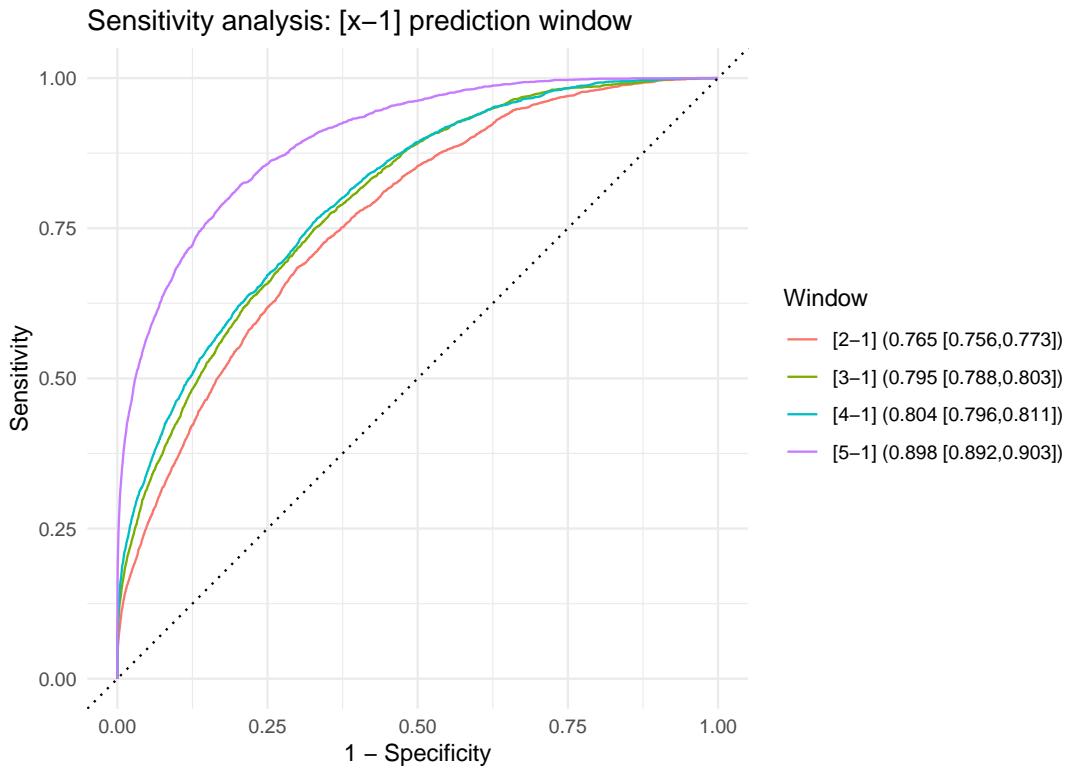


Figure 5: In this experiment, we decrease the amount of data (from 4yrs to 1yr) in the test set, but keep the prediction horizon to 1yr. As expected, performance decreases as the number of years decrease (we have less and less data). However, we notice a large jump in performance from 3yrs to 4yrs. A possible explanation is that our model is specifically trained for 4 years of data and the summary features change significantly. In particular, minimum and maximum statistics are highly sensitive to having more data. This also increases the amount of missing data.

EAC v EGJAC

- They have different predicted risk
 - EGJAC have higher predicted risk, much better separation from controls
 - Esp. for the region of interest, EAC has a lot more subject around the threshold
- Previous comparison of feature distribution was incorrect (was done after imputation)
 - see new results below
- There is a stark difference in proportion of missing values between EAC and EGJAC cases for HCT, LYMPH and MCH.
 - Almost always missing for EGJAC
- Some feature have different impact between the two models; EAC is generally very similar to ANY

	Mean (control)	Mean (EAC)	Mean (EGJAC)	pvalue.adj
black	0.169	0.041	0.089	0.001
gerd	0.232	0.357	0.269	0.005
baso_max	1.017	0.328	1.791	0.033

Table 3: Features with different means between EAC and EGJAC (at 5% level, BH). It appears that race and GERD are less important for EGJAC. We also find that baso could be different between the two.

Prop. Control	Prop. EAC	Prop. EGJAC	pvalue.adj

Table 4: Features with different proportions of non-NAs between EAC and EGJAC (at 5% level, BH). There are None!

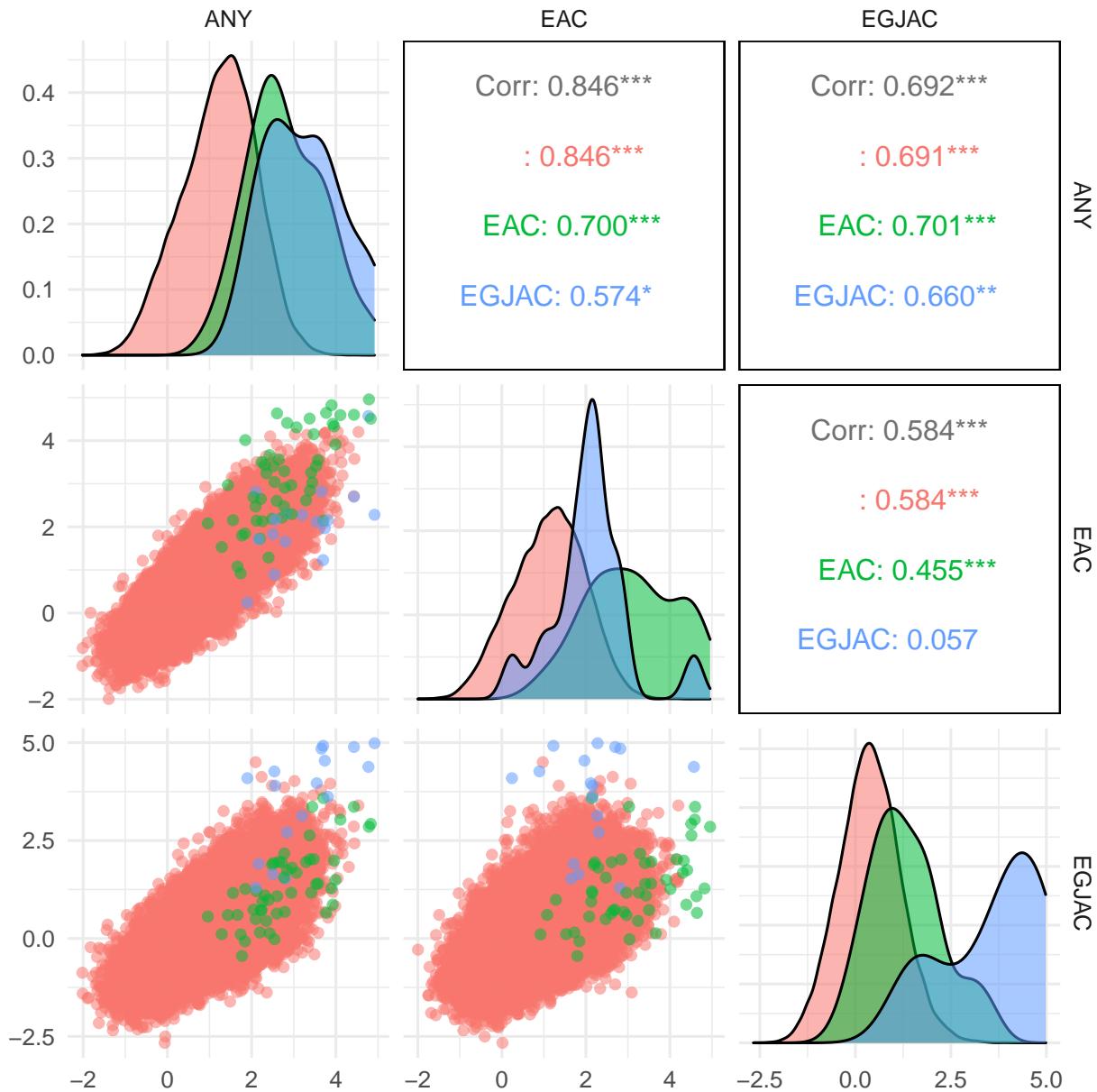


Figure 6: Scatter plots of (log) predicted risk by all three models stratified by case/control status (none, EAC, EGJAC). We have very strong correlation between EAC and ANY, but much less between ANY and EGJAC. There is no correlation between EAC and EGJAC for EGJAC patients. EGJAC assigns very high predicted risk to EGJAC patients compared to what ANY/EAC does to their respective targets.

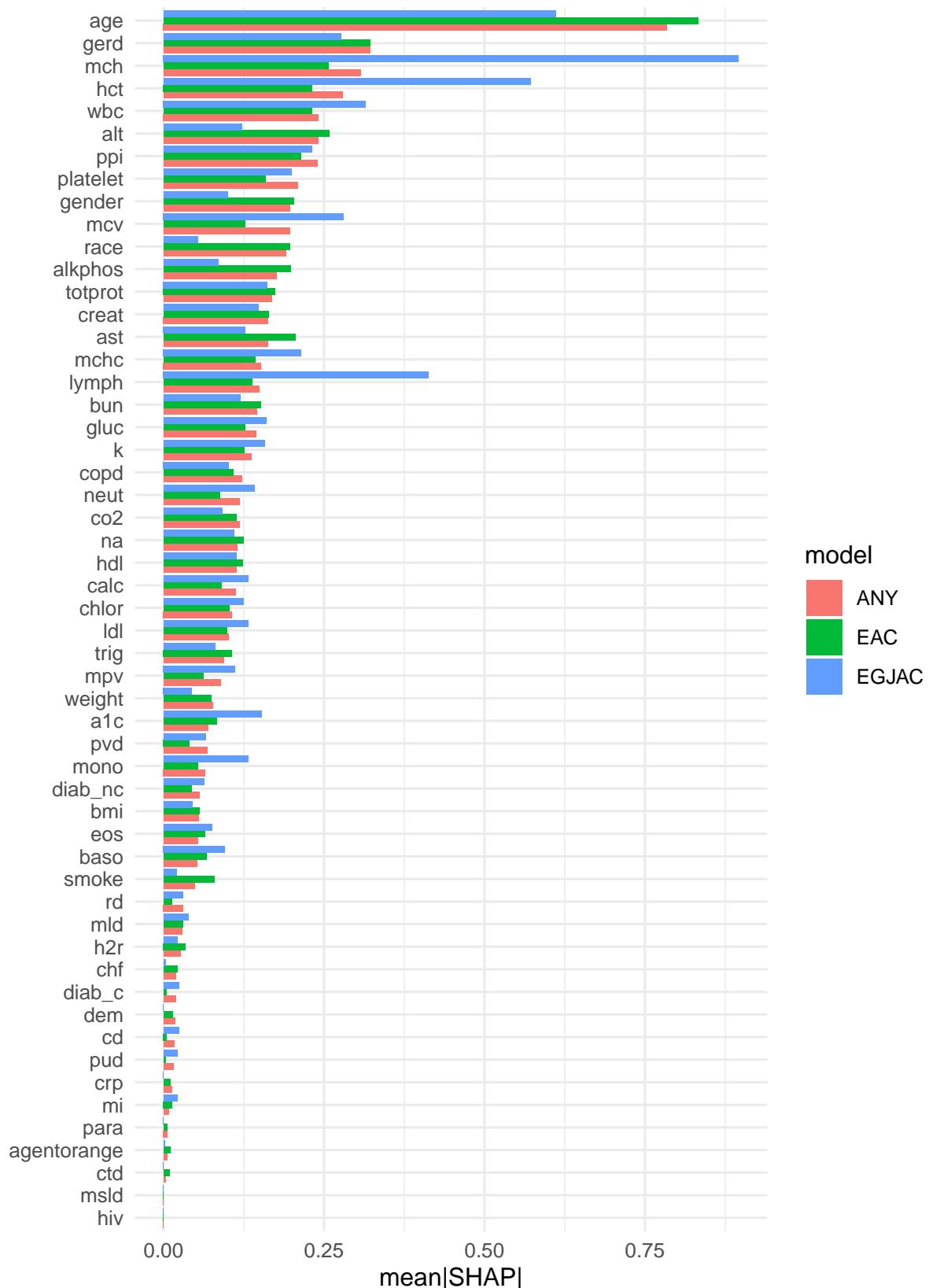


Figure 7: SHAP variable importance within all three models. Notably, we see large differences **age**, **mch**, **hct**, **mcv**, **lymph**, **race** and a few others.

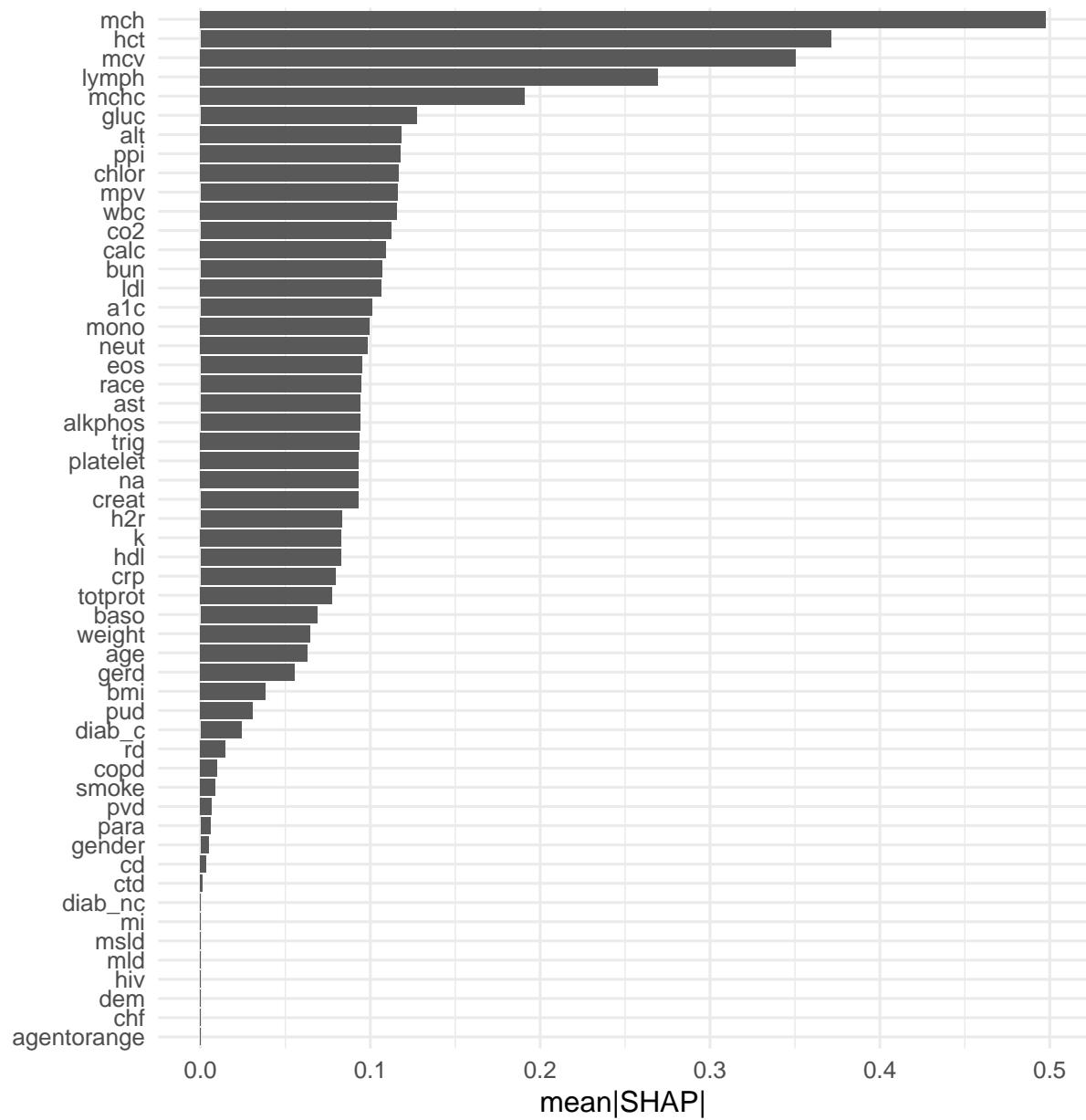


Figure 8: SHAP variable importance for a model differentiating EAC from EGJAC. We again find `mch`, `hct`, `mcv`, `lymph` and `mchc` at the top

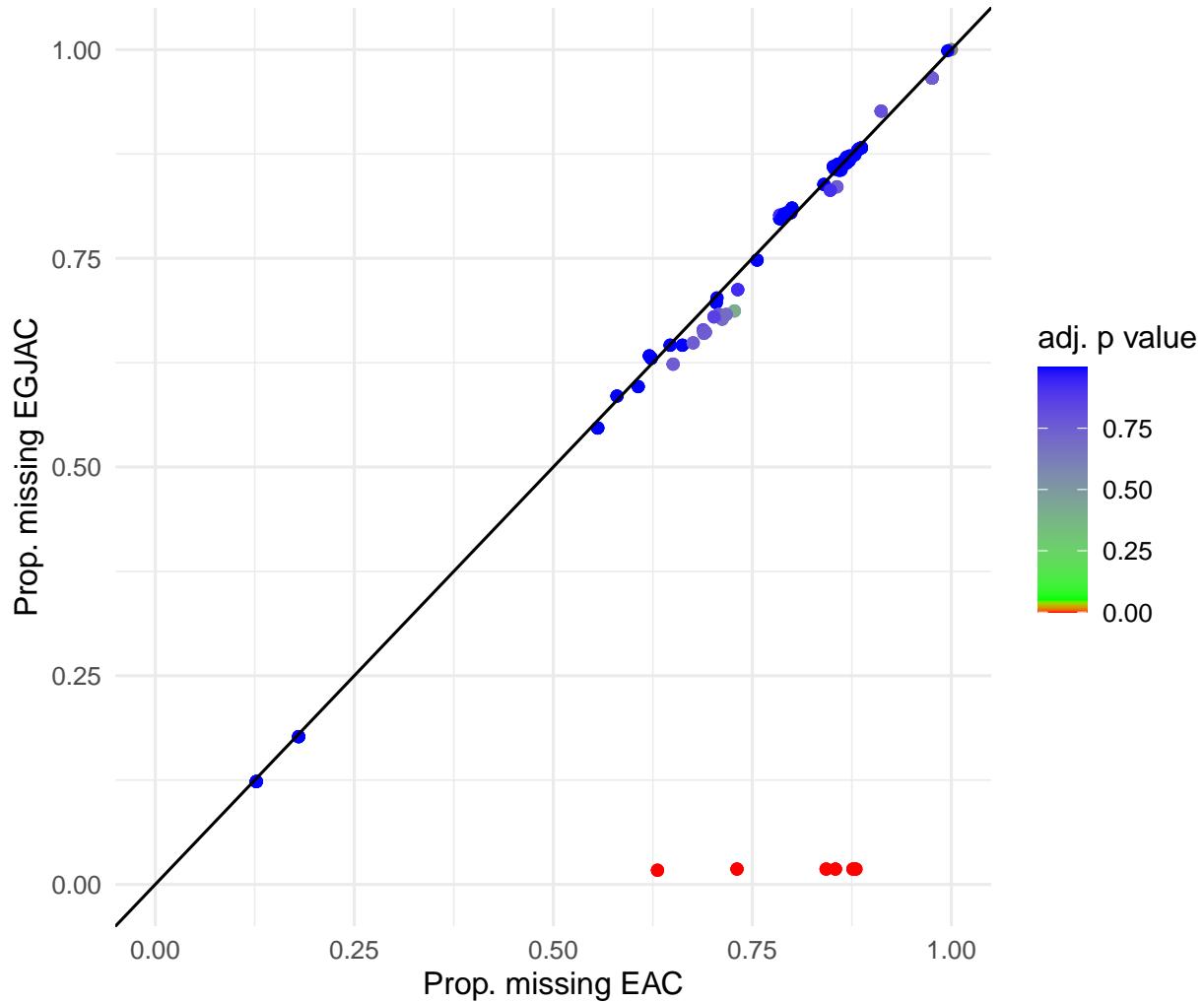
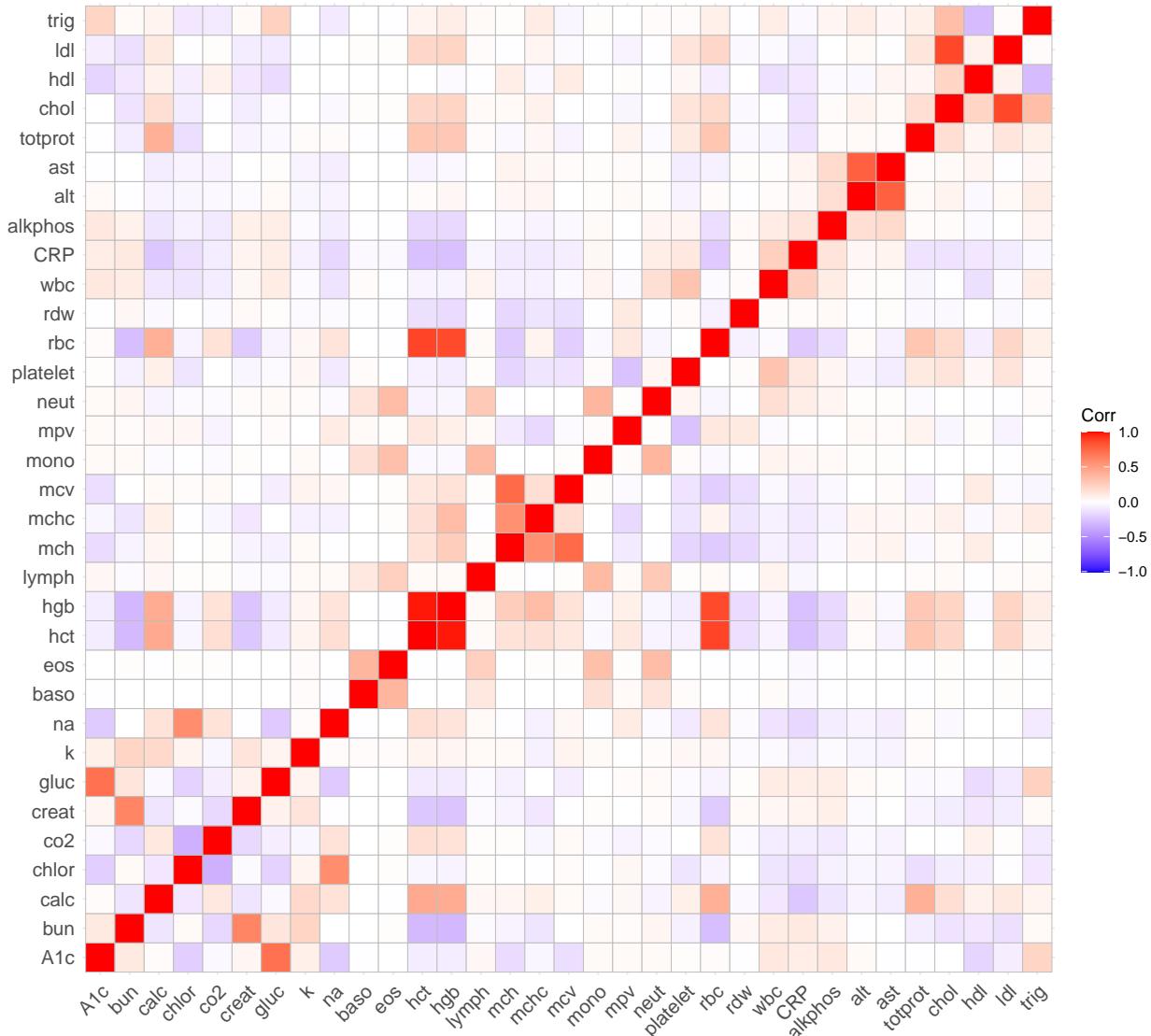


Figure 9: Labels are wrong and should read “non-missing.” This shows it really is just there three that have different missing proportions.

Old stuff to update from here

Feature selection

Correlation:



	var1	var2	correlation
2	mchc	mch	0.569
3	chlor	na	0.578
6	creat	bun	0.617
8	gluc	A1c	0.707
10	mcv	mch	0.742
11	alt	ast	0.782
13	hgb	rbc	0.857
15	chol	ldl	0.875
18	rbc	het	0.884
20	hgb	hct	0.976

Table 5: HGB and HCT are indeed very highly correlated

On chol/ldl/hdl: Missingness patterns does not help

Pattern (chol,hdl,ldl)	000	001	010	011	100	101	110	111
Prop (%)	59.5	0.7	0.5	0.6	0.2	0.0	0.2	38.3

Using performance:

- Drop various features:
 - (colonoscopy, fobt)
 - two of (hct, hgb, rct)
 - one or two of (chol, ldl, hdl)
- Refit with 1M controls & evaluate
- Proposed is dropping (colonoscopy, fobt, hbg, rct, chol)
- Conclusions:
 - Not much difference
 - Probably best to avoid dropping hct

Drop	Valid. AUC	Test AUC
None	0.832	0.827
Colonoscopy, fobt	0.830	0.826
rct, hgb	0.831	0.826
rbc, hct	0.824	0.821
hct, hgb	0.824	0.826
chol	0.834	0.831
ldl	0.831	0.829
hdl	0.832	0.825
hdl, ldl	0.831	0.830
(proposed)	0.827	0.827