



UNIVERSITY OF ALBERTA
SCHOOL OF PUBLIC HEALTH



A Summary Index of Prediction Accuracy for Binary and Censored Time to Event Data

Yan Yuan, PhD

April 23, 2019

OICR Biostatistics Seminar

Clinical Prediction: Examples of Prevention and Planning

- WHO risk charts for cardiovascular disease for most countries
- Numerous risk score systems ($n > 40$) for diabetes risk in general population
- Sepsis risk prediction (CMAJ 2019)

Risk Score as a Screening Tool

- Characteristics of typical condition that risk scores are developed for
 - seriousness (may result in mortality or significantly affect the quality of life);
 - early detection/intervention can make a difference in disease prognosis but may be expensive or invasive;
 - the event rate is low

Motivating Data – Binary outcome

Digital Mammography Imaging Screening Trial (Pisano et al. 2005 *New England Journal of Medicine*)

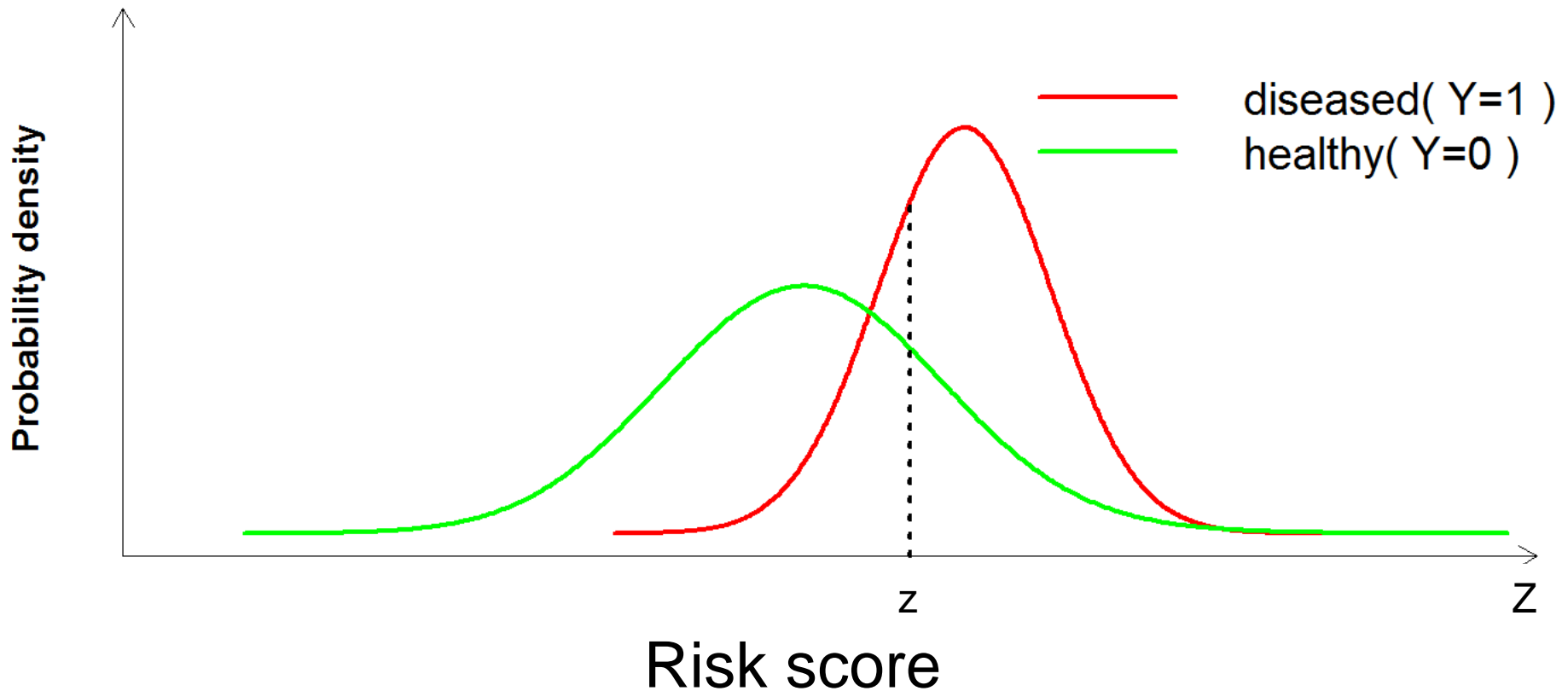
Malignancy score		7	6	5	4	3	2	1	Total
Digital M	Category	11	29	69	1061	2224	6588	32588	42570
	Total								
Film M	Cancers	10	18	25	85	49	25	122	334
	Category	17	29	70	942	2291	6910	32486	42745
	Total								
	Cancers	13	24	25	74	35	33	131	335

42,760 screening participants underwent two screening technology, 335 were diagnosed with breast cancer by the end of 15 months follow-up.

Evaluating Model Performance when Predicting Low Prevalence Events

- Threshold Dependent Measure (predictor needs to be binary)
 - ~~Misclassification rate~~
 - Sensitivity (TPF): $P(\text{test positive} \mid \text{disease present}) = P(\hat{Y} = 1 \mid Y = 1)$
 - ~~Specificity (FPF): $P(\text{test negative} \mid \text{disease absent}) = P(\hat{Y} = 0 \mid Y = 0)$~~
 - Positive Predictive value (PPV): $P(Y = 1 \mid \hat{Y} = 1)$
 - ~~Negative Predictive Value (NPV): $P(Y = 0 \mid \hat{Y} = 0)$~~

When risk score is continuous or ordinal



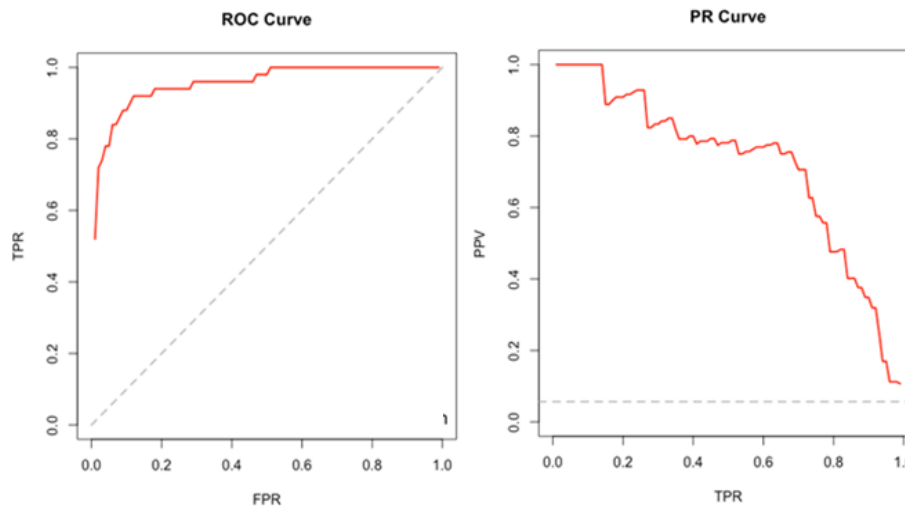
Threshold-free Summary Measure

- Area Under the ROC Curve (AUC)

$$AUC \equiv \int_R TPF(z) dFPP(z)$$

- Area under the Precision-Recall curve

$$AP \equiv \int_R PPV(z) dTPF(z)$$



MLE of AP

Score	x_1	$>$	x_2	$> \dots >$	x_k	$>$	x_{k+1}	$> \dots >$	x_K	
Partition	R_1		R_2	\dots	R_k		R_{k+1}	\dots	R_K	Total
Class-1	Z_1		Z_2	\dots	Z_k		Z_{k+1}	\dots	Z_K	n_1
Class-0	\bar{Z}_1		\bar{Z}_2	\dots	\bar{Z}_k		\bar{Z}_{k+1}	\dots	\bar{Z}_K	n_0
Total	S_1		S_2	\dots	S_k		S_{k+1}	\dots	S_K	n

Data in the above 2 X K table follow

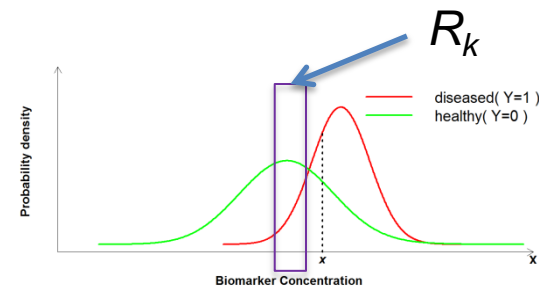
$$(Z_1, Z_2, \dots, Z_K) | n_1 \sim \text{multinomial}(n_1; p_1, p_2, \dots, p_K),$$

$$(\bar{Z}_1, \bar{Z}_2, \dots, \bar{Z}_K) | n_1 \sim \text{multinomial}(n - n_1; q_1, q_2, \dots, q_K),$$

$$n_1 \sim \text{binomial}(n, \pi),$$

For continuous risk scores

$$p_k = \int_{R_k} f_1(x) dx, \quad q_k = \int_{R_k} f_0(x) dx,$$



Asymptotic Variance of AP

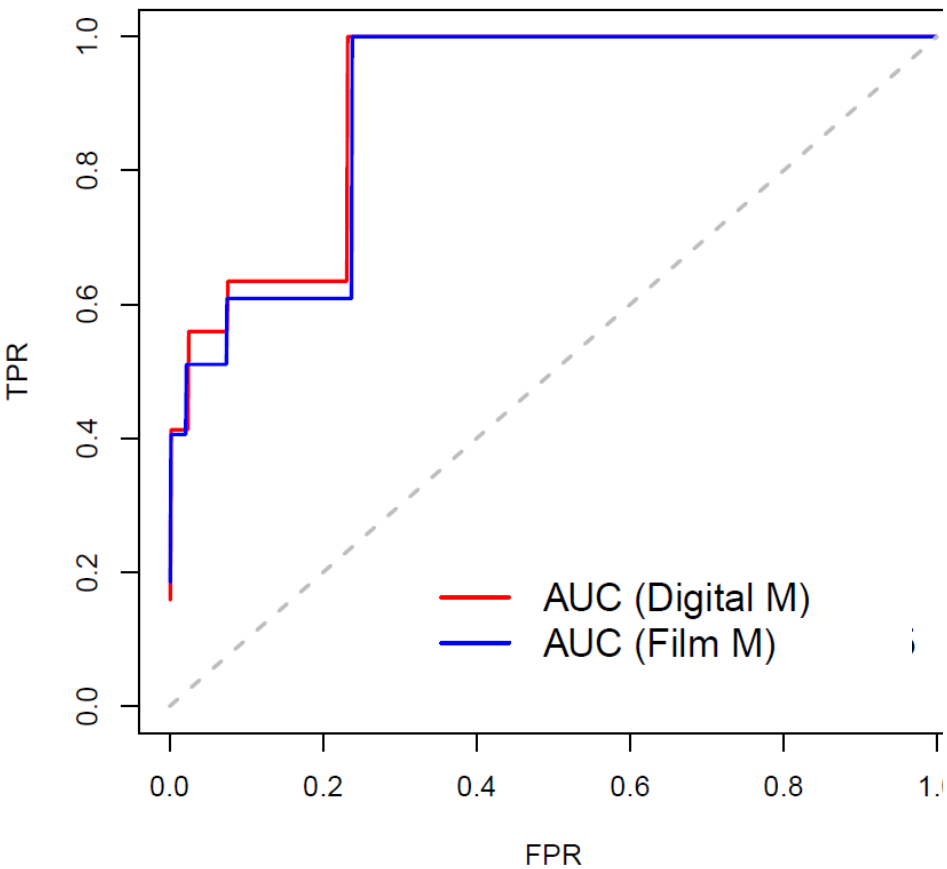
$$\widehat{AP} = g(\widehat{p}_k, \widehat{q}_k, \widehat{\pi}) = \sum_{k=1}^K \left[\widehat{p}_k \left(\frac{\widehat{\pi} \sum_{k' \leq k} \widehat{p}_{k'}}{\widehat{\pi} \sum_{k' \leq k} \widehat{p}_{k'} + (1 - \widehat{\pi}) \sum_{k' \leq k} \widehat{q}_{k'}} \right) \right]$$

Applying the Delta method, we get the variance estimator

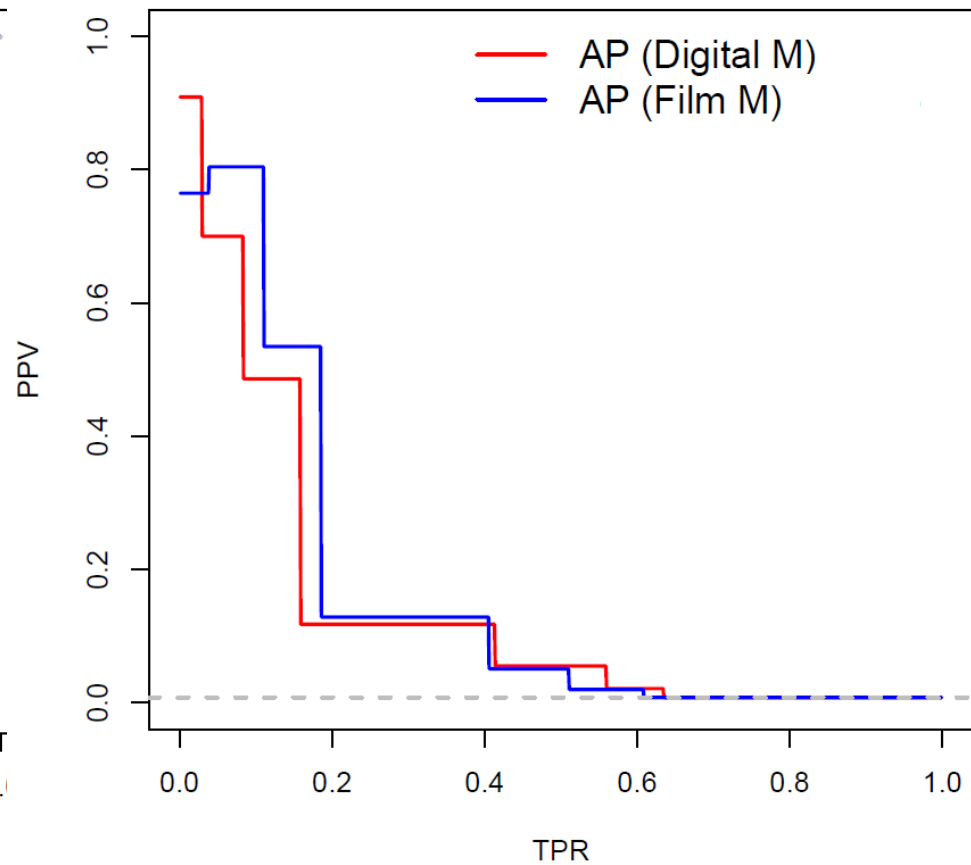
$$\widehat{var}(\widehat{AP}) \approx (\nabla g)^T \widehat{J}^{-1} (\nabla g)$$

Yuan et al. (2015)

ROC Curves



PR Curves



Malignancy score		7	6	5	4	3	2	1	Total
Digital M	Category Total	11	29	69	1061	2224	6588	32588	42570
	Cancers	10	18	25	85	49	25	122	334
Film M	Category Total	17	29	70	942	2291	6910	32486	42745
	Cancers	13	24	25	74	35	33	131	335

Given that 335 breast cancer diagnosed in 42,760 screening participants at 15 months follow-up, the prevalence π is 0.78%.

Seven-point Malignancy Scale

	\widehat{AUC} (s.e.)	\widehat{AP} (s.e.)
Film mammography	0.735 (0.012)	0.166 (0.022)
Digital mammography	0.753 (0.012)	0.144 (0.021)

Remark: Resampling method can be used for the inference of the difference in AP if we have paired data.

Yuan et al. (2015)

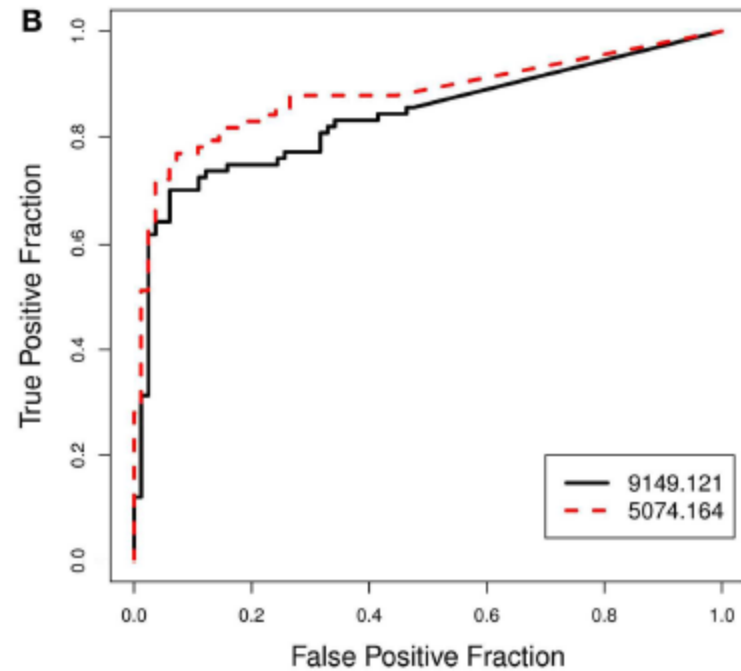
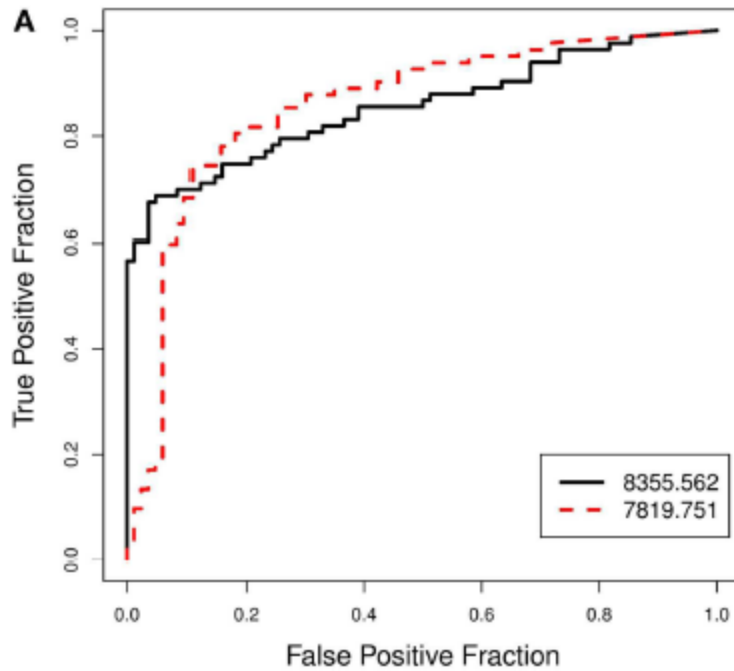


Table 1 | Prostate cancer example.

Biomarkers	AUC			AP		
	$n_0 \times 1 (\pi \approx 0.5)$	$n_0 \times 10 (\pi \approx 0.09)$	$n_0 \times 100 (\pi \approx 0.01)$	$n_0 \times 1 (\pi \approx 0.5)$	$n_0 \times 10 (\pi \approx 0.09)$	$n_0 \times 100 (\pi \approx 0.01)$
A	8355.562	0.849	0.783	0.856	0.606	0.571
	7819.751	0.850	0.857	0.802	0.370	0.062
B	5074.164	0.886	0.869	0.833	0.306	0.043
	9149.121	0.832	0.793	0.822	0.512	0.225

A simple thought experiment showing changes in the estimated AUC and AP as a result of artificially inflating the number of control subjects (n_0) to mimic real-life screening settings, where the prevalence (π) of disease is low.

Score Partition	x_1	$>$	x_2	$> \dots >$	x_k	$>$	x_{k+1}	$> \dots >$	x_K	Total
	R_1		R_2	\dots	R_k		R_{k+1}	\dots	R_K	
Class-1	Z_1		Z_2	\dots	Z_k		Z_{k+1}	\dots	Z_K	n_1
Class-0	\bar{Z}_1		\bar{Z}_2	\dots	\bar{Z}_k		\bar{Z}_{k+1}	\dots	\bar{Z}_K	n_0
Total	S_1		S_2	\dots	S_k		S_{k+1}	\dots	S_K	n

$$\widehat{AP} = \underbrace{\left[\frac{Z_1}{S_1} \right]}_{w_1} \left[\frac{Z_1}{n_1} \right] + \underbrace{\left[\frac{Z_1 + Z_2}{S_1 + S_2} \right]}_{w_2} \left[\frac{Z_2}{n_1} \right] + \dots + \underbrace{\left[\frac{Z_1 + Z_2 + \dots + Z_K}{S_1 + S_2 + \dots + S_K} \right]}_{w_K} \left[\frac{Z_K}{n_1} \right]$$

$$= \sum_{k=1}^K w_k \left[\frac{Z_k}{n_1} \right]$$

$$\widehat{AUC} = \frac{n}{n_0} \left\{ \underbrace{\left[\frac{S_1 + S_2 + \dots + S_K}{n} \right]}_{w'_1} \left[\frac{Z_1}{n_1} \right] + \underbrace{\left[\frac{S_2 + \dots + S_K}{n} \right]}_{w'_2} \left[\frac{Z_2}{n_1} \right] + \dots + \underbrace{\left[\frac{S_K}{n} \right]}_{w'_K} \left[\frac{Z_K}{n_1} \right] - \frac{1}{2} \left(\frac{n_1}{n_0} \right) \right\} - \frac{1}{2} \left(\frac{n_1}{n_0} \right)$$

$$= \frac{n}{n_0} \sum_{k=1}^K w'_k \left[\frac{Z_k}{n_1} \right] - \frac{1}{2} \left(\frac{n_1}{n_0} \right)$$

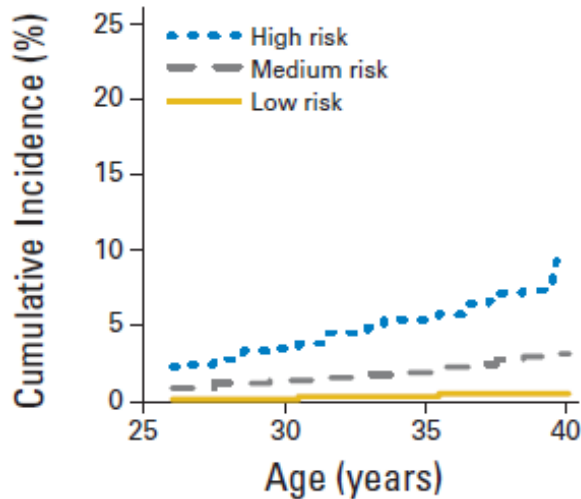
AP – AUC Relationship

- When two risk scores U_1 and U_2 are compared
 - If ROC curve of U_1 dominates that of U_2 everywhere, then PR curve of U_1 dominates that of U_2 everywhere. $AUC_1 > AUC_2$ and $AP_1 > AP_2$
 - If ROC curves of U_1 and U_2 crosses, the ranking of U_1 and U_2 based on of AUC and AP may differ.
- Both AUC and AP are semi-proper scoring rule.

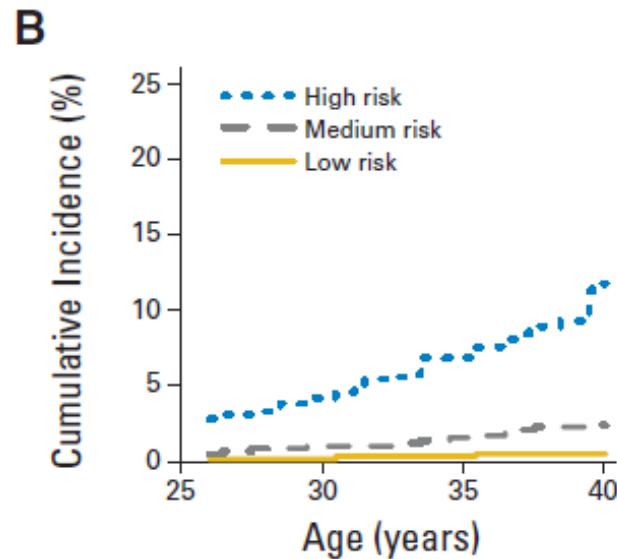
Motivating Data – Time to Event outcome

- Late effects of cancer treatments in childhood cancer survivors – e.g. Congestive heart failure (Chow et al. 2015, *Journal of Clinical Oncology*)
- Cumulative risk of CHF is ~3% by 35 years post diagnosis

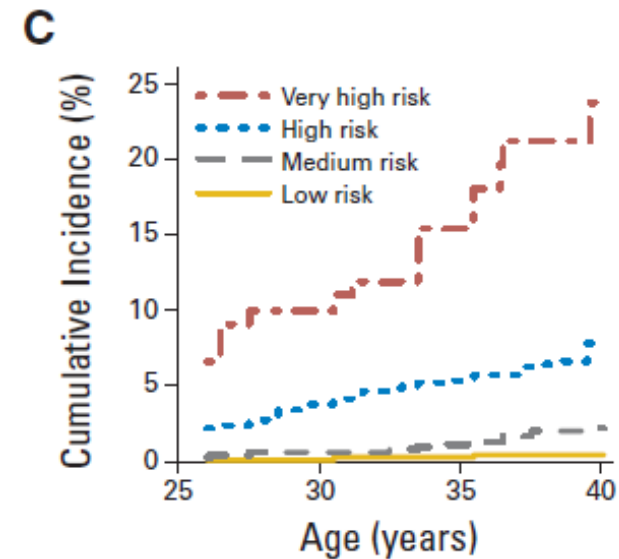
Simple Model



Standard Model



Standard + Heart Dose Model



AP_{t_0} for Time-to-Event Outcome

- Time-dependent Average Positive predictive value (AP_{t_0})

$$AP_{t_0} = \int_{\mathcal{R}} PPV_{t_0}(z) dTPF_{t_0}(z).$$

Yuan et al. 2018

Nonparametric Estimator for Event Status

Let (X, δ, Z) be the standard time to event data notation,
 X : the censored event time, δ : the censoring indicator
 Z : the risk score

$$\widehat{AP}_{t_0} = \frac{\sum_{j=1}^n I(X_j \leq t_0) \widehat{w}_{t_0,j} \widehat{PPV}_{t_0}(Z_j)}{\sum_{j=1}^n I(X_j \leq t_0) \widehat{w}_{t_0,j}}.$$

where

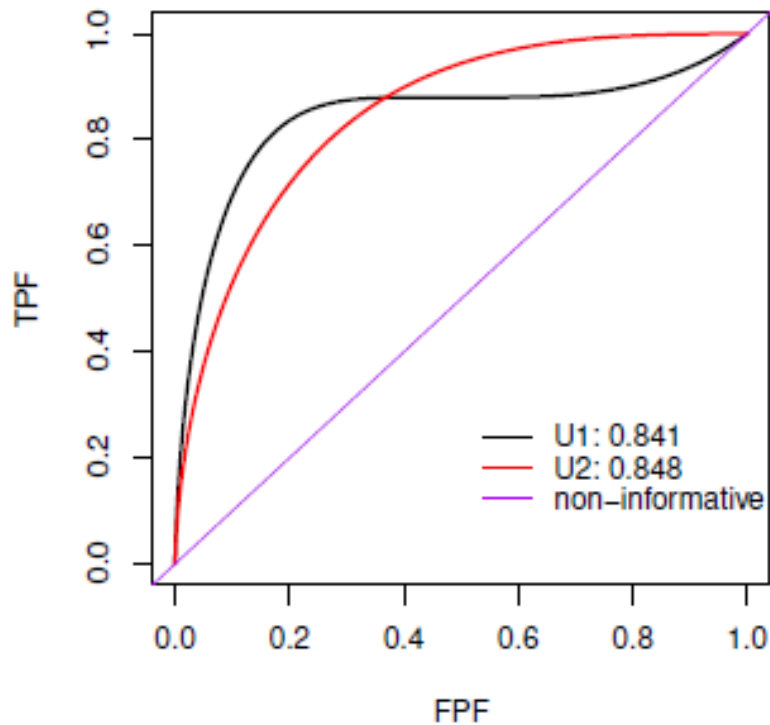
$$\widehat{w}_{t_0,i} = \frac{I(X_i < t_0) \delta_i}{\widehat{G}(X_i)} + \frac{I(X_i \geq t_0)}{\widehat{G}(t_0)}$$

$$\widehat{PPV}_{t_0}(z) = \frac{\sum_{i=1}^n \widehat{w}_{t_0,i} I(Z_i \geq z) I(X_i < t_0)}{\sum_{i=1}^n \widehat{w}_{t_0,i} I(Z_i \geq z)}.$$

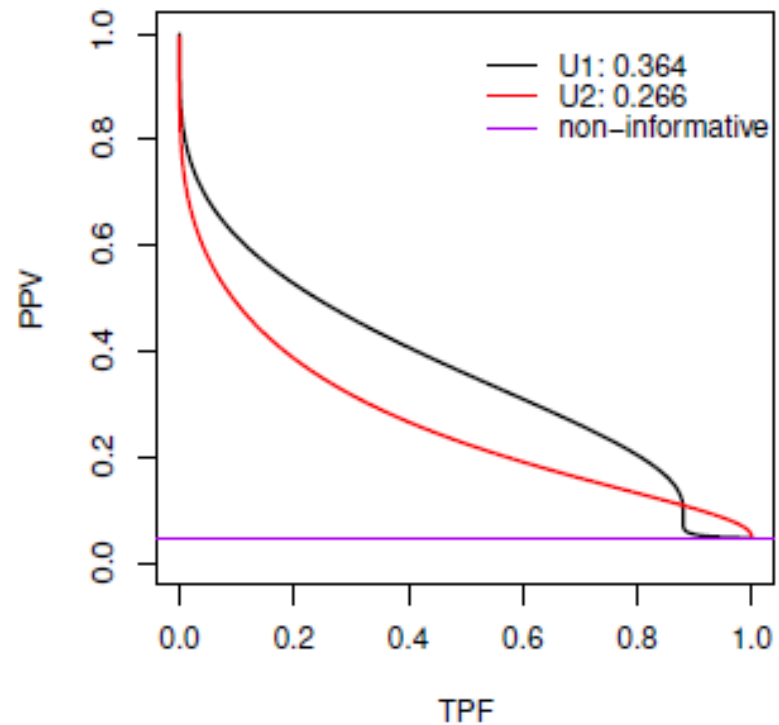
Simulation Study

$$\log(T_i) = 7.2 - 1.1U_{i1} - 2.5U_{i2} - 1.5\log(U_{i1}^2) + \epsilon_T,$$

$ROC_{t_0=8}$



$PR_{t_0=8}$



Results (n=2000)

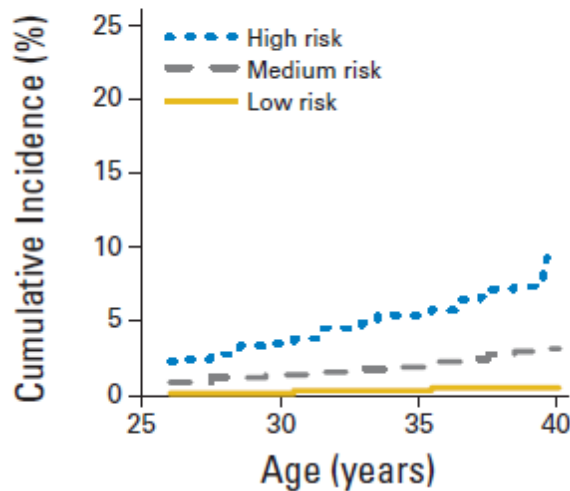
t_0	Event rate	Risk score	AP				$ECOVP^b(\%)$	AUC TRUE
			TRUE	BIAS	ESE	ASE^b		
0.5	0.0101	U_1	0.182	0.0361	0.0806	0.0794	92.2	0.920
		U_2	0.124	0.0339	0.0687	0.0679	94.1	0.904
		Δ	0.058	0.0251	0.102	0.116	96.1	0.016
		Ratio	1.47	0.4820	1.470	1.740	92.4	1.02
8	0.0495	U_1	0.364	0.0085	0.0508	0.0499	94.4	0.841
		U_2	0.266	0.0121	0.0435	0.0439	94.8	0.848
		Δ	0.098	-0.0028	0.0707	0.072	96.3	-0.007
		Ratio	1.37	0.0123	0.310	0.322	95.8	0.99
36	0.0991	U_1	0.462	0.0060	0.0416	0.0431	94.2	0.786
		U_2	0.375	0.0074	0.0387	0.0393	96.3	0.824
		Δ	0.087	-0.0045	0.0655	0.0633	95.7	-0.038
		Ratio	1.23	-0.0010	0.189	0.187	94.5	0.95

Results (n=5000)

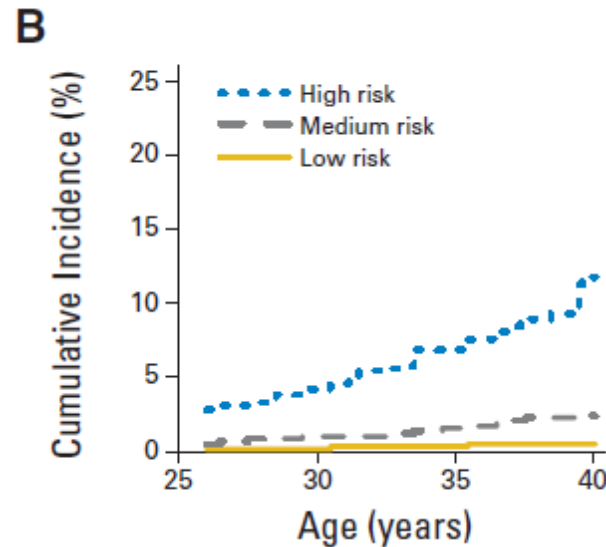
t_0	Event rate	Risk score	AP				$ECOV P^b(\%)$	AUC TRUE
			TRUE	BIAS	ESE	ASE^b		
0.5	0.0101	U_1	0.182	0.0185	0.0498	0.0503	93.6	0.920
		U_2	0.124	0.0154	0.0415	0.0415	93.6	0.904
		Δ	0.058	0.0056	0.0696	0.0712	94.2	0.016
		Ratio	1.47	0.1490	0.709	0.756	92.9	1.02
8	0.0495	U_1	0.364	0.0041	0.0327	0.0324	94.0	0.841
		U_2	0.266	0.0043	0.0285	0.0280	95.5	0.848
		Δ	0.098	-0.0005	0.0473	0.0460	96.3	-0.007
		Ratio	1.37	0.0099	0.209	0.204	94.5	0.99
36	0.0991	U_1	0.462	0.0023	0.0273	0.0275	95.0	0.786
		U_2	0.375	0.0015	0.0247	0.0251	95.5	0.824
		Δ	0.087	0.0003	0.0398	0.0402	95.1	-0.038
		Ratio	1.23	0.0058	0.117	0.120	95.0	0.95

Time to event outcome: CCSS CHF Risk

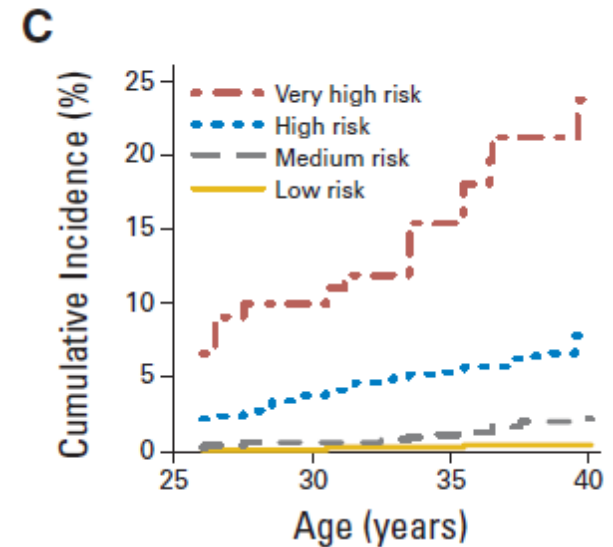
Simple Model



Standard Model



Standard + Heart Dose Model

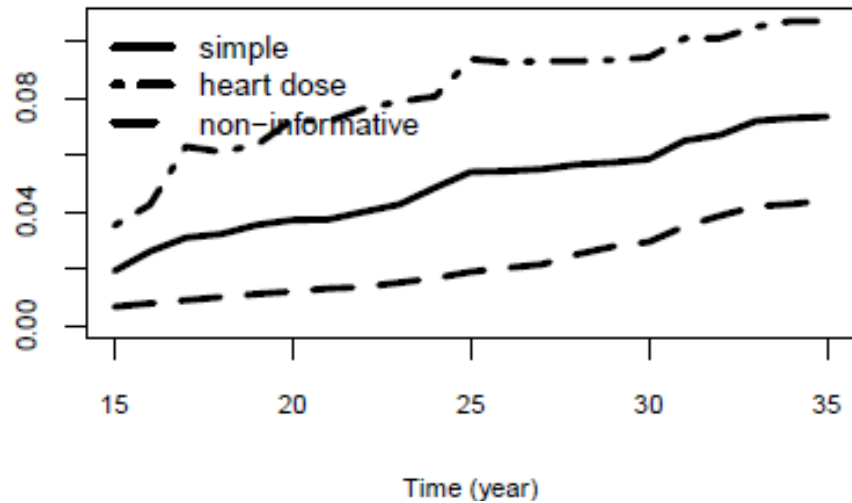


$$PPV_{t_0}^{CHF}(z) = Pr\{T < t_0, \Delta = 1 \mid Z \geq z\} \quad \text{and} \quad TPF_{t_0}^{CHF}(z) = Pr\{Z \geq z \mid T < t_0, \Delta = 1\}.$$

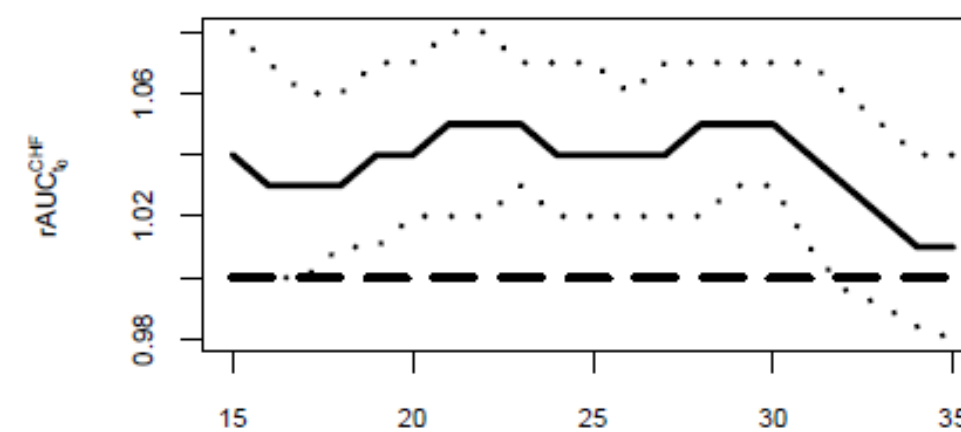
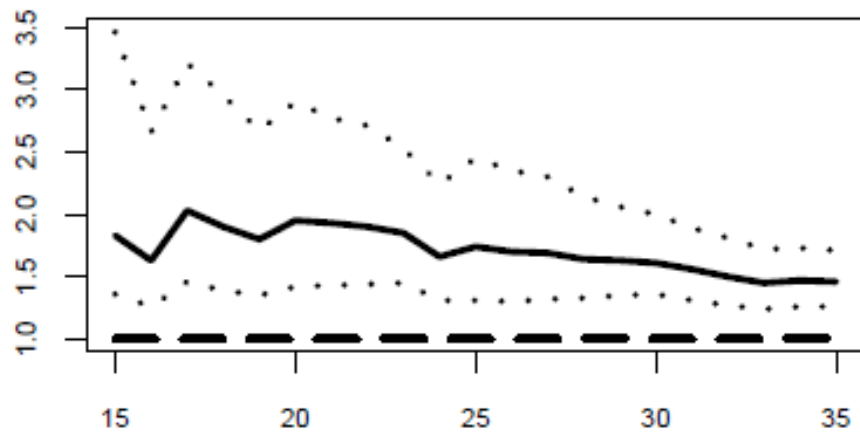
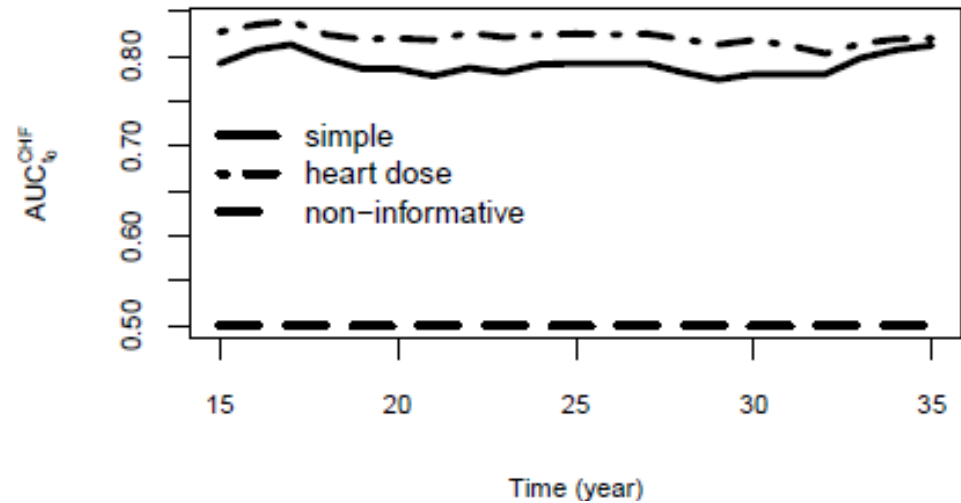
$$\widehat{PPV}_{t_0}^{CHF}(z) = \frac{\sum_{i=1}^n \hat{w}_{t_0,i} I(Z_i \geq z) I(X_i < t_0) I(\Delta_i = 1)}{\sum_{i=1}^n I(Z_i \geq z)}$$

$$\widehat{TPF}_{t_0}^{CHF}(z) = \frac{\sum_{i=1}^n \hat{w}_{t_0,i} I(Z_i \geq z) I(X_i < t_0) I(\Delta_i = 1)}{\sum_{i=1}^n \hat{w}_{t_0,i} I(X_i < t_0) I(\Delta_i = 1)}$$

AP_{t_0} vs. t_0



AUC_{t_0} vs. t_0



Incremental Value

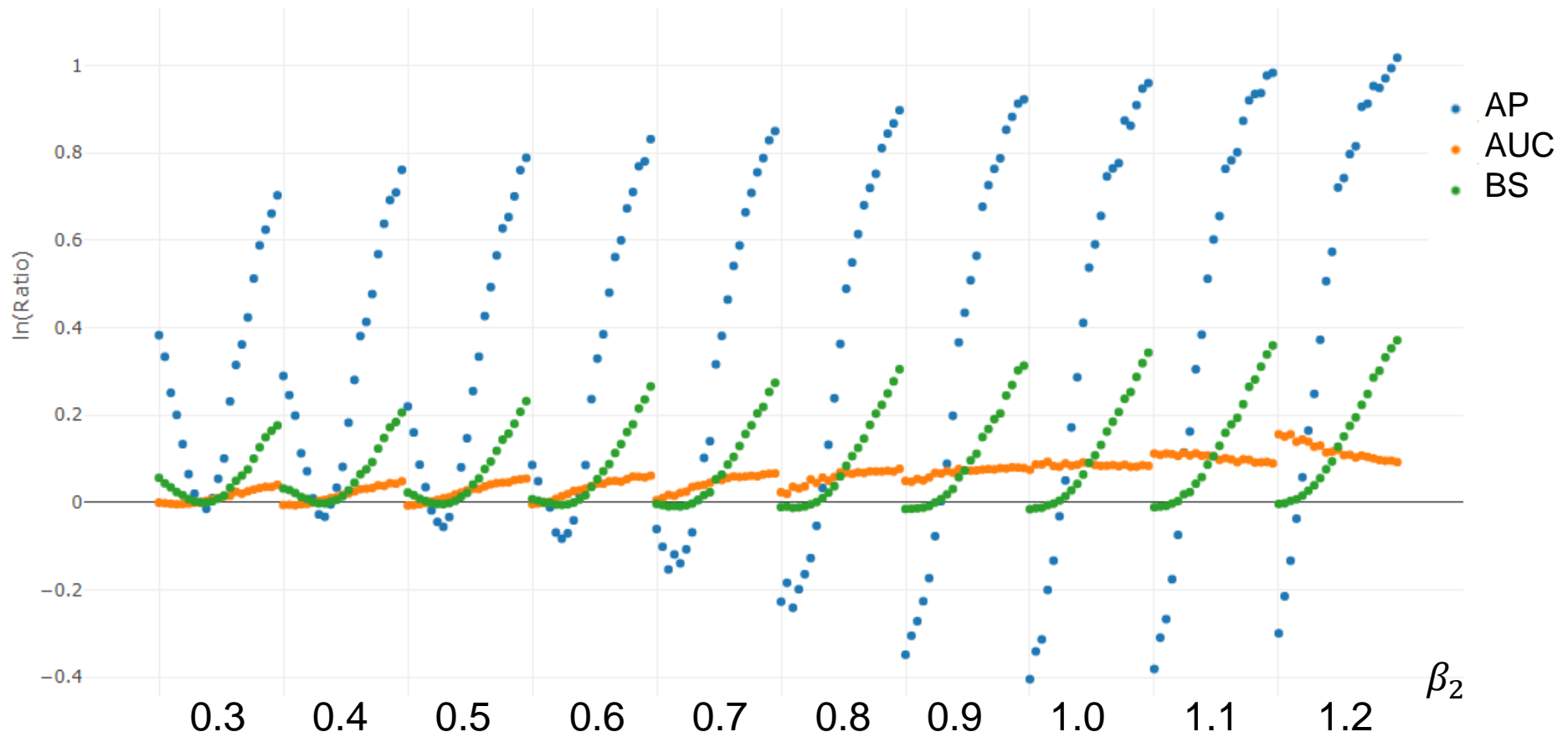
- Risk factor & outcome association vs. information/calibration gain in prediction
- Existing metrics
 - Changes in AUC and Brier scores (BS)
 - NRI (net reclassification improvement)
 - IDI (integrated discrimination improvement)

How does AP changes, in comparison to changes in AUC and BS?

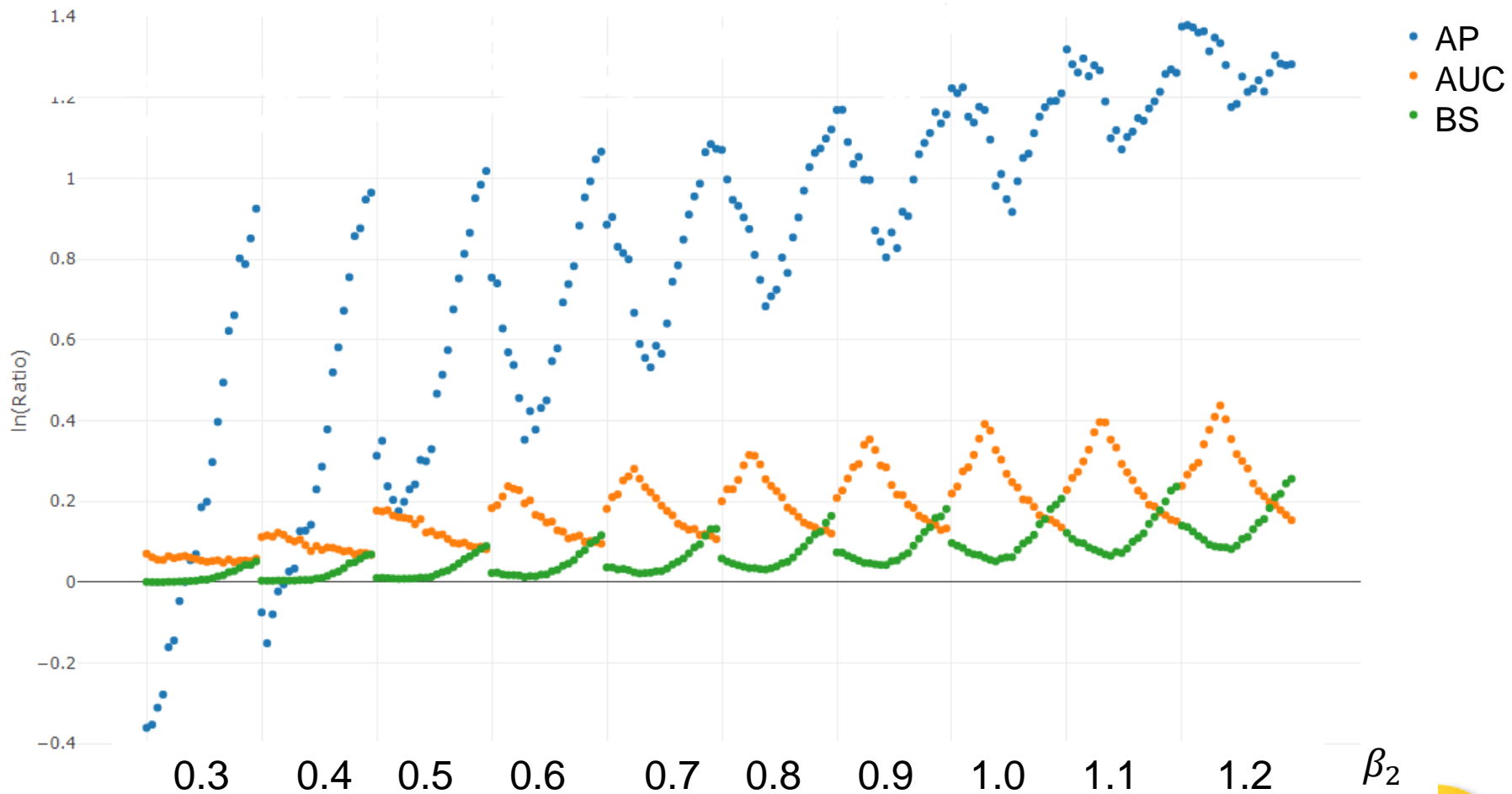
Simulation Study

- True model: $\text{logit}(\pi) = \beta_0 + \beta_1 U_1 + \beta_2 U_2 + \beta_3 U_1 U_2$,
 - β_1 and β_2 range: [0.3, 1.2]
 - β_3 range: [-1, 1]
 - Independent U_1 & $U_2 \sim \text{iid } N(0, 1)$
 - Event rate: ~5%
- Working model
 - Model 1: $\text{logit}(\pi) = \beta_0 + \beta_1 U_1$
 - Model 2: $\text{logit}(\pi) = \beta_0 + \beta_1 U_1 + \beta_2 U_2$
- Metrics
 - rAUC, rAP and rBS

$$\beta_1 = 1.2$$



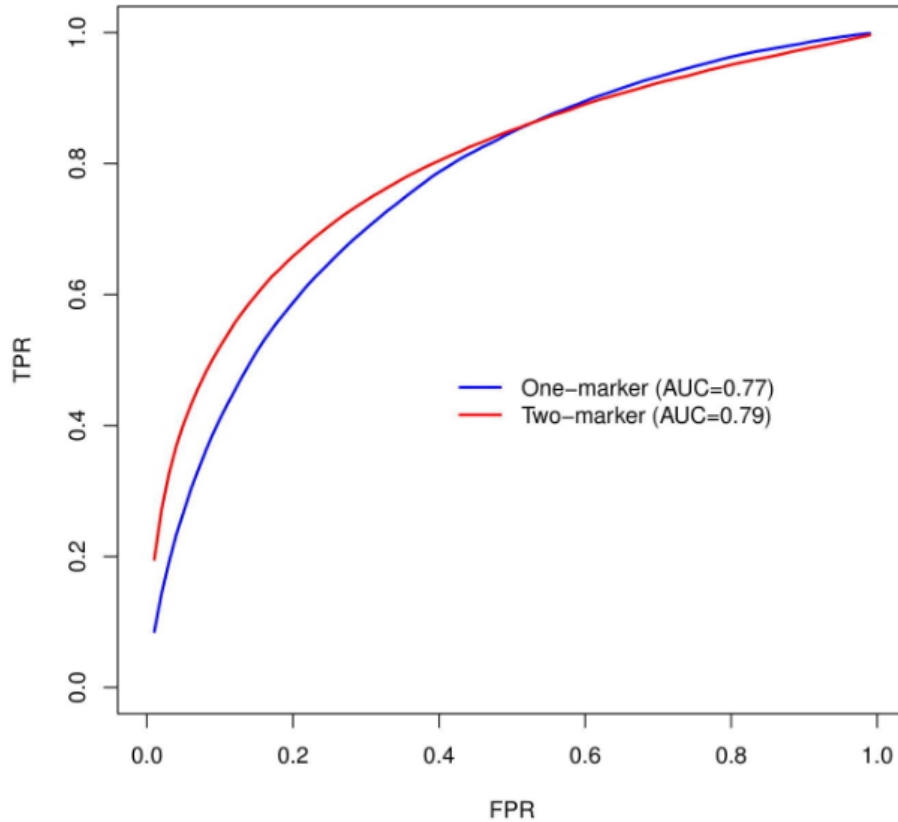
$$\beta_1 = 0.3$$



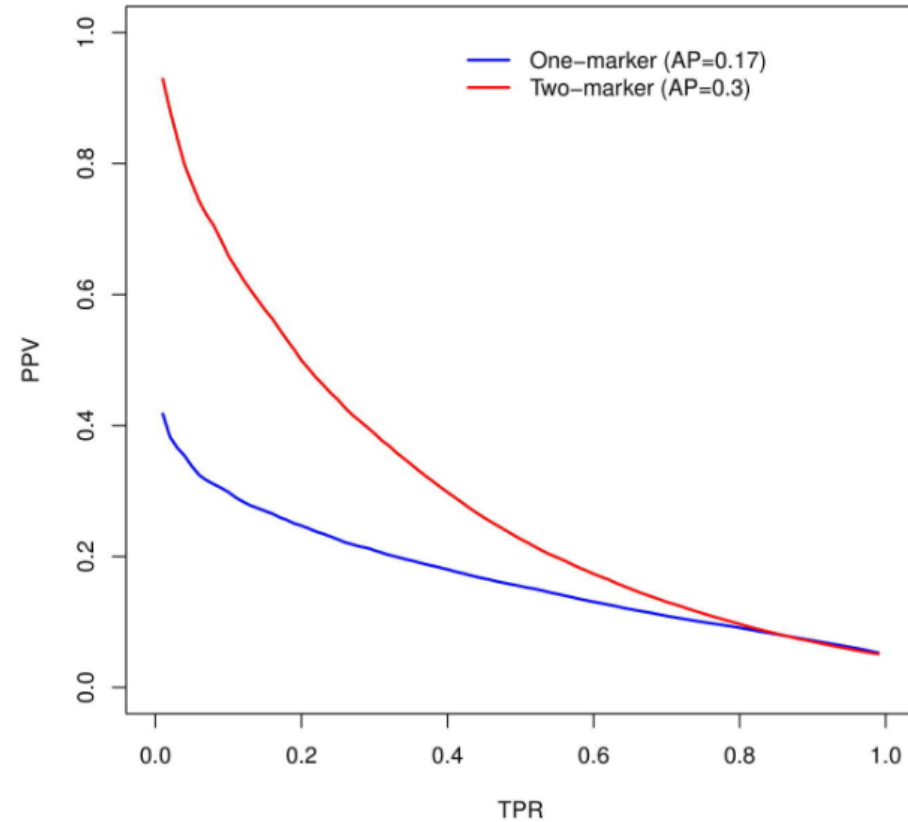
Metrics	Correlation	
	Pearson	Spearman
Log(ratio of metrics: M2/M1)		
-ln(rBS) and ln(rAUC)	0.083	0.30
-ln(rBS) and ln(rAP)	0.76	0.89
ln(rAUC) and ln(rAP)	0.48	0.51

$$\beta_1 = 0.9, \beta_2 = 0.3, \beta_3 = 0.6$$

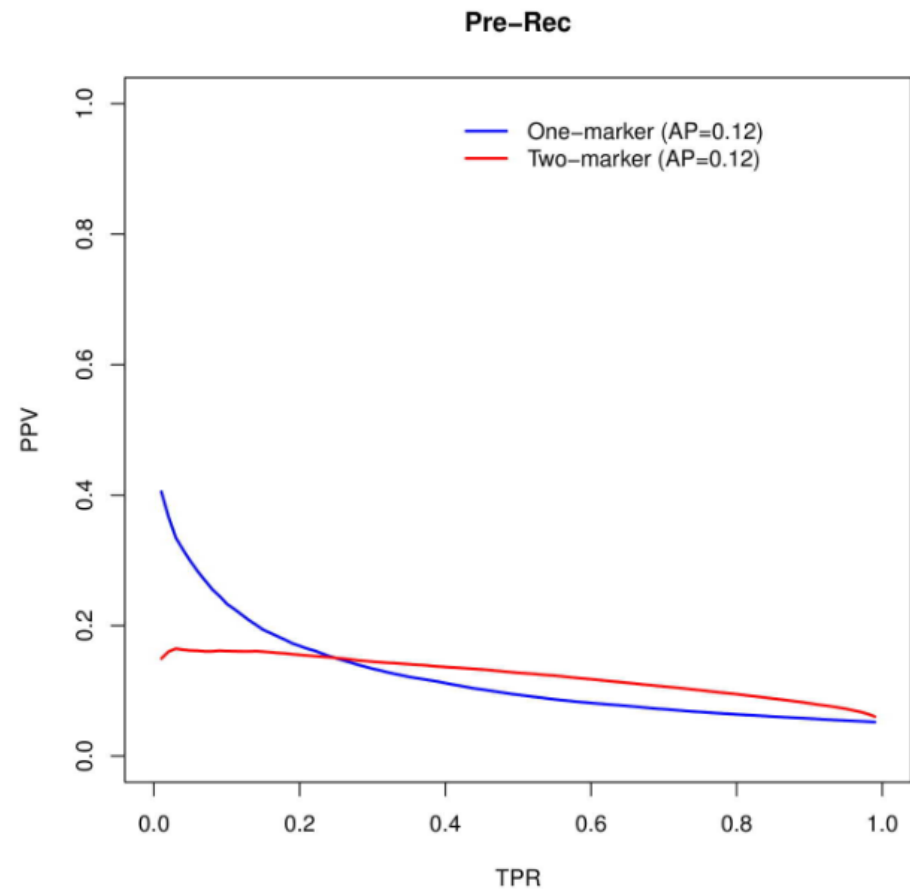
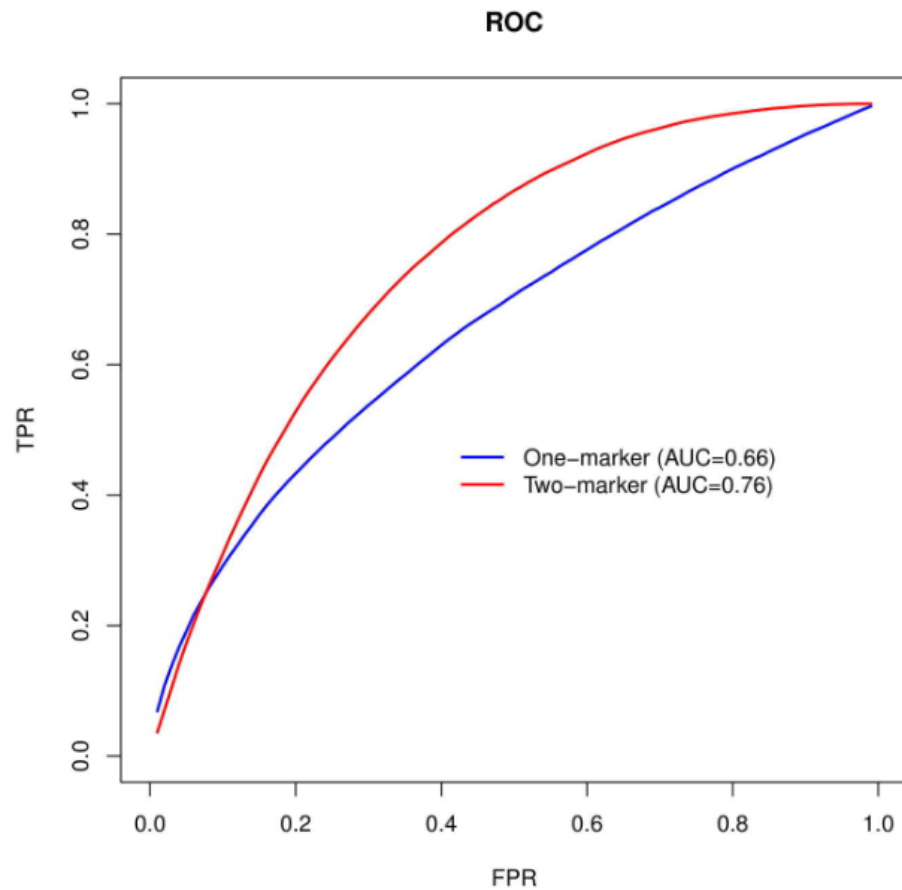
ROC



Pre-Rec



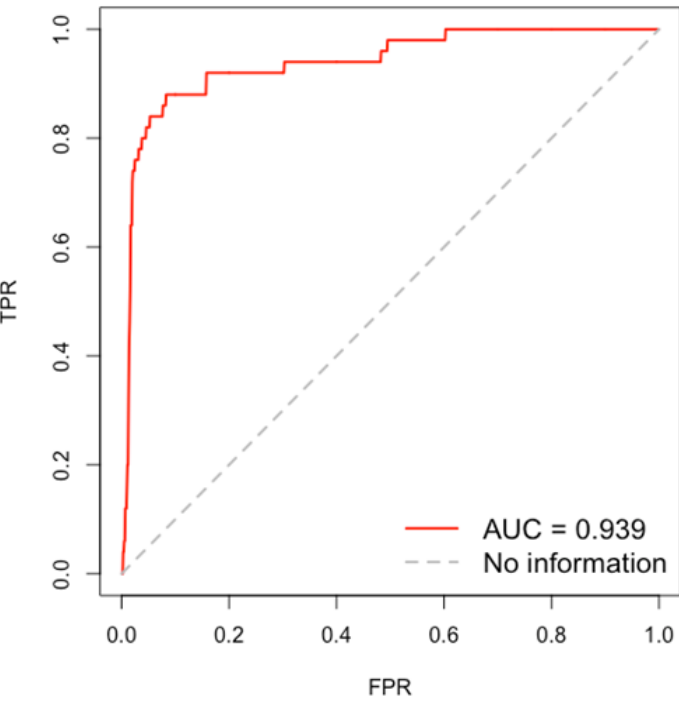
$$\beta_1 = 1, \beta_2 = 1, \beta_3 = -0.6$$



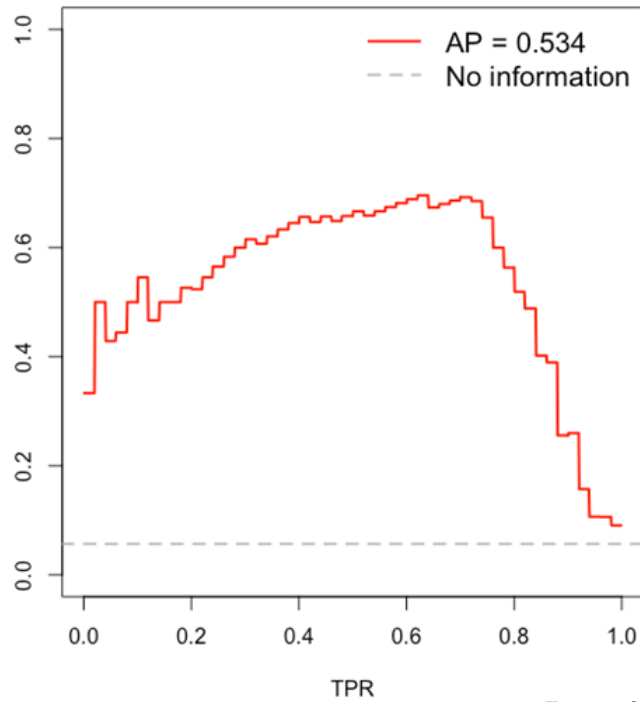
Risk Prediction for Ovarian Failure

- Goal
 - Developing risk prediction model for ovarian failure (OF) in childhood cancer survivors (CCS)
- Data
 - About 6000 female CCS (dx 1970-1999)
- Methods
 - Logistic regression; Random Forest; and Support Vector Machines
- Results
 - AUC 0.82 and AP 0.50 for Acute OF (Internal validation)

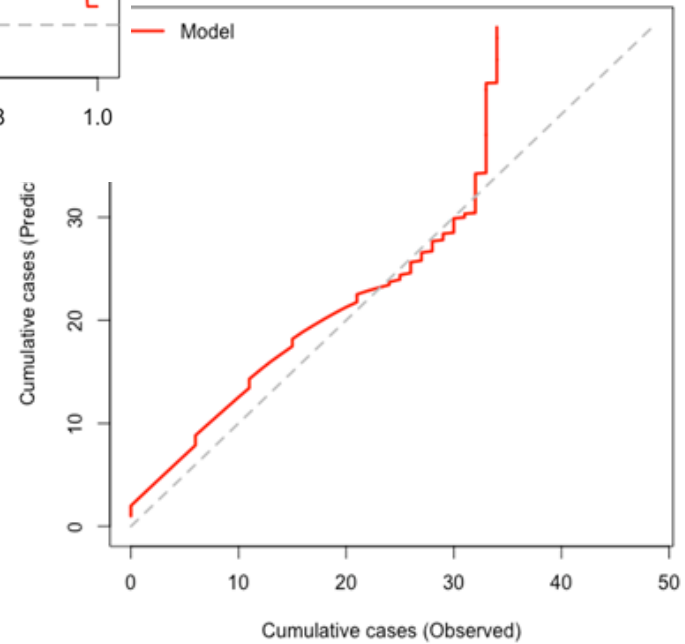
ROC Curve



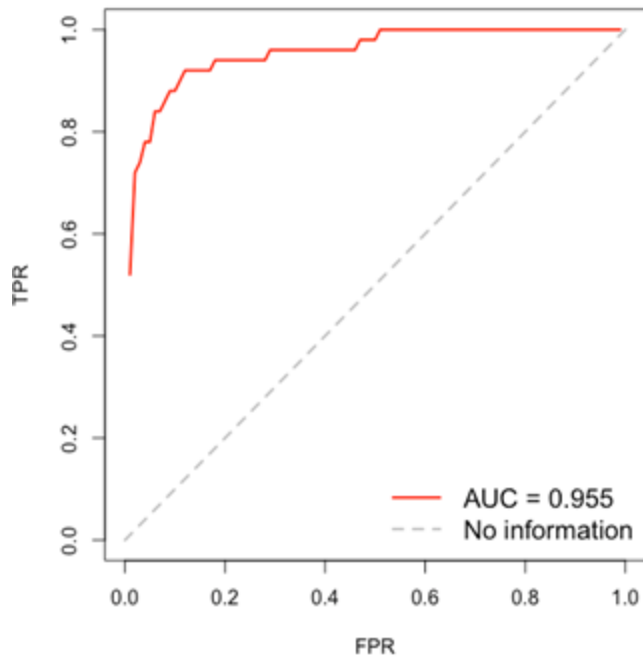
PR Curve



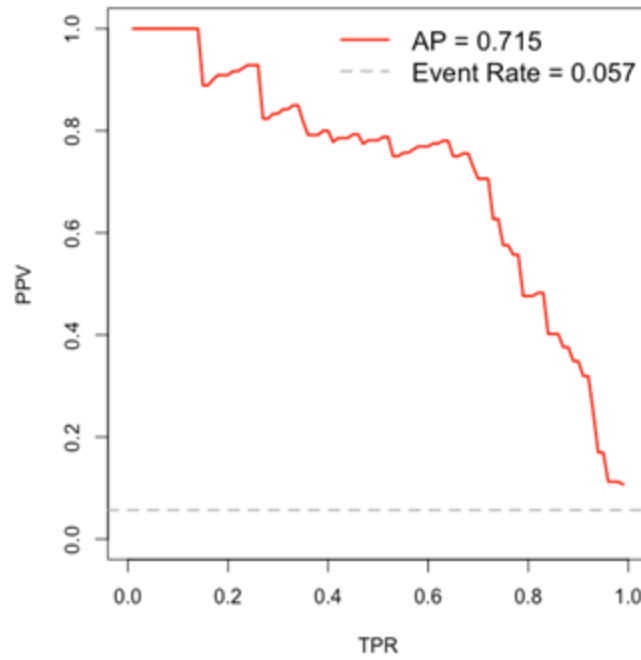
Calibration Curve



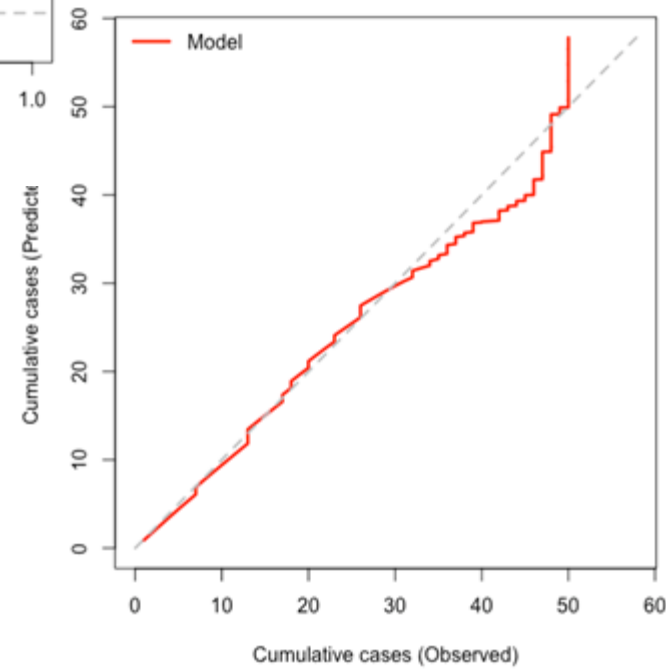
ROC Curve



PR Curve



Calibration Curve



Discussion

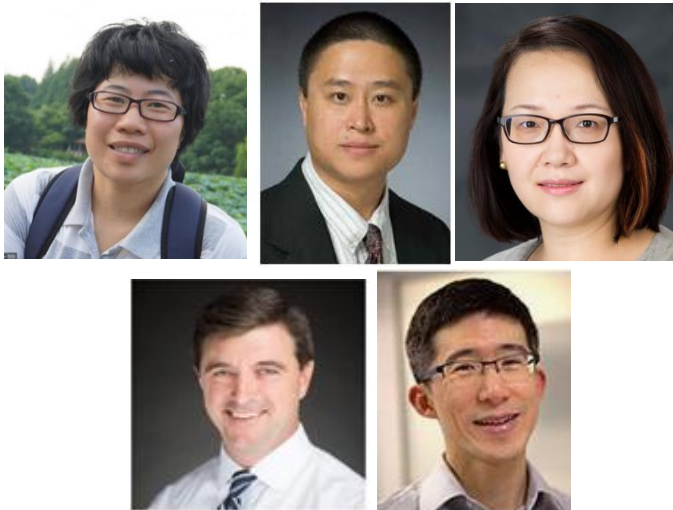
- AP is a single numerical measure, in this respect it is similar to AUC.
- A summary measure of positive predictive value, useful for evaluating and comparing prospective prediction performance of risk scores.
- More sensitive than AUC.
- Better aligned with the strict proper scoring rule Brier score than AUC (under misspecified working models)
- Event rate dependent, AP should be estimated in a prospective cohort or population-based study
- R package <APtools> and SAS macro for binary and survival time data <https://sites.ualberta.ca/~yyuan/software.html>

Acknowledgement

Students and staff

Maoji Li, Dr. Khanh Vu

Doris Li, Hengrui Cai, Zorina Han, Rebecca Clark, Michael Lu



M.S.I. Foundation



CIHR IRSC



Canadian Institutes of Health Research
Instituts de recherche en santé du Canada



UNIVERSITY OF ALBERTA
SCHOOL OF PUBLIC HEALTH





UNIVERSITY OF ALBERTA
SCHOOL OF PUBLIC HEALTH



Thank you!

Questions???