

NONPARAMETRIC ADJUSTMENT FOR MEASUREMENT ERROR IN TIME TO EVENT DATA: APPLICATION TO RISK PREDICTION MODELS

Malka Gorfine

Tel Aviv University, Israel

Joint work with

Danielle Braun and Giovanni Parmigiani,

Department of Biostatistics, Harvard School of Public Health

Setting

Risk prediction model has first been developed based on **error-free time to event data**, and subsequently implemented in practical setting where **time to event data can be error-prone**.

Motivation I

Mendelian risk prediction models in genetic counseling:

- Calculating the probability that an individual carries a cancer causing inherited mutation based on his/her **family history**.
- Predicting the absolute risk of developing the disease over time given his/her mutation-carrier status and **family history**.

These models are in wide clinical use and web-based patient-oriented tools: breast cancer, ovarian cancer, Lynch syndrome, pancreatic cancer, melanoma etc.

These models were developed based on error-free data.

Studies about accuracy of self-reported family history show that sensitivity and specificity for reported disease status vary by degree of relative and cancer type.

Breast cancer:

65% of truly affected are reported affected

2% of truly unaffected are reported as affected

Disease status, sensitivity 65% - 95%; specificity 98% - 99%.

Age of diagnosis was misreported for 3.1% of relatives, average of 4.5 years between the true and misreported ages (Mai et al 2011; Ziogas and Anton-Culver, 2003).

Ovarian cancer:

Age of diagnosis was misreported for 4.2% of relatives, average of 4.2 years between the true and misreported ages (Ziogas and Anton-Culver, 2003).

Misreporting of family history, especially in disease status, leads to distortions in predictions (Katki, 2006).

Q: Is it possible to develop prediction models based on **error-prone data?**

A: Not in the context of Mendelian risk prediction models which relies on penetrance estimates from the literature, based on error-free data.



disease probability given carrier status

Motivation II

Survival prediction models:



Time to progression – length of time until the disease gets worse or spread

In some disease settings, such as cancer, TTP is one of the predictors for survival.

Assume a model has been developed based on **error-free TTP**.

In practice, **TTP is error-prone**:

- Tumor assessment is done using imaging, which varies by observers.
- Scans are taken at regularly intervals.

The current setting vs the common settings

The usual measurement error setting:

- ✓ Error-prone covariate observed in the main study.
- ✓ The goal is estimating the relationship between the outcome and the true covariate.

Current setting:

- ✓ The relationship between the outcome and the true covariate is known.
- ✓ The goal is to use this model for risk estimation based on an error-prone covariate.

Naïvely using the error-prone covariate will lead to biased results.

The data:

Y - outcome

T^o - the true failure time

C - the true right-censoring time

Error-free predictor: $H = (T, \delta)$ (for simplicity, one relative)

where $T = \min(T^o, C)$ $\delta = I(T^o \leq C)$

Example: $Y = 0$ or 1 and T^o the mother's age at onset.

Error-free predictor: $H = (T, \delta)$

Error-prone predictor: $H^* = (T^*, \delta^*)$

Example: the counselee doesn't know that his/her relative had the disease, or the correct age at onset.

Assumption: We have a validation study with

$$H = (T, \delta) \quad \text{and} \quad H^* = (T^*, \delta^*)$$

but no need for the outcome Y .

The risk prediction model: $\Pr(Y | H)$

Our goal is estimating $\Pr(Y | H^*)$

Our goal is estimating $\Pr(Y | H^*)$

Main idea:

$$\begin{aligned} \Pr(Y | H^*) &= \int_H \Pr(Y, H | H^*) dH \\ &= \int_H \Pr(Y | H, H^*) \Pr(H | H^*) dH \\ &= \int_H \Pr(Y | H) \Pr(H | H^*) dH \end{aligned}$$

Assumptions:

- H^* contains no information on predicting Y beyond H .
- The measurement error model $\Pr(H | H^*)$ is transportable.

surrogacy assumption

A non parametric estimator of $\Pr(H | H^*)$

The distribution is left unspecified

Assume, for simplicity, one family member

$$\rightarrow H = (T, \delta) \quad H^* = (T^*, \delta^*)$$

Main idea:

$$P(T, \delta | T^*, \delta^*)$$

$$= \lambda(T | T^*, \delta^*)^\delta S(T | T^*, \delta^*) h(T | T^*, \delta^*)^{1-\delta} G(T | T^*, \delta^*)$$

Conditional hazard and survival of true failure time

Conditional hazard and survival of true censoring time

Assumptions:

Conditional independence of event and censoring times given H^* .

$$\begin{aligned}
 & P(T, \delta | T^*, \delta^*) \\
 &= \lambda(T | T^*, \delta^*)^\delta S(T | T^*, \delta^*) h(T | T^*, \delta^*)^{1-\delta} G(T | T^*, \delta^*)
 \end{aligned}$$

These hazards and survival functions can be estimated non-parametrically by using the validation data – a large study population that does not involve the counselee.

Validation data

$$T_i = \min(T_i^o, C_i) \quad \delta_i = I(T_i^o \leq C_i) \quad i = 1, \dots, n$$

$$H_i = (T_i, \delta_i) \quad H_i^* = (T_i^*, \delta_i^*)$$



dependent individuals

Use kernel smoothed Kaplan-Meier estimator (Beran, 1981).

Kernel smoothed Kaplan-Meier estimator (Beran, 1981):

Nadaraya-Watson weight

$$W_i(t; b_{nl}, l) = I(\delta_i^* = l) K\left(\frac{t - T_i^*}{b_{nl}}\right) \quad i = 1, \dots, n \quad l = 0, 1$$

Bandwidth sequences

Known kernel function

$$\hat{S}(t | t^*, l) = \prod_{T_i \leq t, \delta_i = 1, \delta_i^* = l} \left(1 - \frac{W_i(t^*; b_{nl}, l)}{\sum_{j=1, \delta_j^* = l}^n W_j(t^*; b_{nl}, l) I(T_j \geq T_i)} \right)$$

and we get $\hat{S}(t | t^*, 0)$, $\hat{S}(t | t^*, 1)$ for $t, t^* \in (0, \tau]$.

For estimating the survival function of the censoring time, apply the above while treating the censoring times as events and the event times as censoring.

Summary

$$\hat{\mathbf{P}}(Y | H^*) = \int_H \mathbf{P}(Y | H) \hat{\mathbf{P}}(H | H^*) dH$$

provided

$$H = (H_1, \dots, H_R)$$

$$H^* = (H_1^*, \dots, H_R^*)$$

$$\mathbf{P}(H | H^*) = \prod_{j=1}^R \mathbf{P}(H_j | H_j^*)$$

$$\hat{\mathbf{P}}(H | H^*) = \prod_{j=1}^R \hat{\mathbf{P}}(H_j | H_j^*)$$

$$\hat{\mathbf{P}}(H_j | H_j^*) = \hat{\mathbf{P}}(T_j, \delta_j | T_j^*, \delta_j^*)$$

$$= \hat{\lambda}(T_j | T_j^*, \delta_j^*)^{\delta_j} \hat{S}(T_j | T_j^*, \delta_j^*) \hat{h}(T_j | T_j^*, \delta_j^*)^{1-\delta_j} \hat{G}(T_j | T_j^*, \delta_j^*)$$

$$\hat{P}(Y | H^*) = \int_H \mathbf{P}(Y | H) \hat{P}(H | H^*) dH$$

provided

In case integrating over all possible values of H is computational challenging, a Monte-Carlo estimator can be used, by sampling

$$H^{(1)}, \dots, H^{(B)}$$

From $\hat{P}(H | H^*)$ and the final proposed estimator is given by

$$\hat{P}(Y | H^*) = \frac{1}{B} \sum_{b=1}^B \hat{P}(Y | H^{(b)})$$

Application: Mendelian Risk Prediction Model

A counselee provides information on R relatives:

$$H_i^* = (T_i^*, \delta_i^*) \text{ instead of } H_i = (T_i, \delta_i) \quad i = 1, \dots, R$$

Let

$$\gamma_i = (\gamma_{i1}, \dots, \gamma_{iM}), \quad \gamma_{im} = 0 \text{ or } 1$$

$\gamma_{im} = 1$ indicates carrying the genetic variant that confer disease risk

Aims:

- Estimating $P(\gamma_0 | H_0, H_1^*, \dots, H_R^*)$
- Estimating $P(T_0^o > t | H_0, H_1^*, \dots, H_R^*)$

Carrier probability $P(Y | H) = P(\gamma_0 | H_0, H_1, \dots, H_R)$

Write

$$P(\gamma_0 | H_0, H_1, \dots, H_R) = \frac{P(\gamma_0) \sum_{\gamma_1, \dots, \gamma_R} \prod_{i=0}^R P(H_i | \gamma_i) P(\gamma_1, \dots, \gamma_R | \gamma_0)}{\sum_{\gamma_0} P(\gamma_0) \sum_{\gamma_1, \dots, \gamma_R} \prod_{i=0}^R P(H_i | \gamma_i) P(\gamma_1, \dots, \gamma_R | \gamma_0)}$$

conditional independence of family members' phenotype given their genotypes.

BRCAPRO estimate it via meta-analysis, and family history information is verified using medical records

and in practice $P(Y | H^*) = P(\gamma_0 | H_0, H_1^*, \dots, H_R^*)$

is naively being used. We propose

$$\hat{P}(\gamma_0 | H^*) = \int_H P(\gamma_0 | H) \hat{P}(H | H^*) dH$$

Survival probability $\mathbf{P}(Y | H) = \mathbf{P}(T_0^o > t | H_0, H_1, \dots, H_R, \gamma_0)$

Write

$$\begin{aligned} & \mathbf{P}(T_0^o > t | H_0, H_1, \dots, H_R, \gamma_0) \\ &= \frac{\mathbf{P}(T_0^o > t | \gamma_0) \sum_{\gamma_1, \dots, \gamma_R} \prod_{i=1}^R \mathbf{P}(H_i | \gamma_i) \mathbf{P}(\gamma_1, \dots, \gamma_R | \gamma_0)}{\sum_{\gamma_1, \dots, \gamma_R} \prod_{i=0}^R \mathbf{P}(H_i | \gamma_i) \mathbf{P}(\gamma_1, \dots, \gamma_R | \gamma_0)} \end{aligned}$$

and in practice $\mathbf{P}(Y | H^*) = \mathbf{P}(\gamma_0 | H_0, H_1^*, \dots, H_R^*)$

is naively being used. We propose

$$\hat{\mathbf{P}}(T_0^o > t | H^*) = \int_H \mathbf{P}(T_0^o > t | H) \hat{\mathbf{P}}(H | H^*) dH$$

Simulation Study

- Setting: $P(Y | H) = P(\gamma_0 | H_0, H_1, \dots, H_R)$ with single gene BRCA1

- Two datasets were generated, one to model the measurement error distribution, and the other represents the counselees.

50,000 counselees

100,000 families, each with 5 members (mother, father, 3 daughters)

- The BRCA1 carrier probability 0.006098.

- The penetrance function $P(H | \gamma)$ from BRCAPRO version 2.08.

- Normal censoring, mean 55, SD 10.

- **Measurement error in disease status:** sen=0.954, spec=0.974; sen=0.649 and spec=0.990.

- **Measurement error in age:**

$$T^* = T + \varepsilon \quad \varepsilon \sim N(0, \sigma^2) \quad \sigma = 1, 3, 5$$

$$T^* = TU \quad U \sim \text{Exp}(1)$$

Results

Table 1: Mendelian Risk Prediction Simulation Results. MSEP and O/E improve using the adjusted proposed method, ROC-AUC either improves or remains the same depending on the setting.

Counselee	Sens/Spec	Error in Age	$\sqrt{MSEP^\dagger} * 1000$			O/E			ROC-AUC		
			Error-Free	Error-Prone	Adjusted	Error-Free	Error-Prone	Adjusted	Error-Free	Error-Prone	Adjusted
Mother	0.954, 0.974	a: $N(0, 5^2)$	0.0000	19.1351	16.7405	0.9773	0.8190	0.9712	0.8160	0.8090	0.8086
		a: $N(0, 3^2)$	0.0000	18.1006	15.9490	0.9773	0.8280	0.9746	0.8160	0.8098	0.8078
		a: $N(0, 1^2)$	0.0000	17.5526	15.6430	0.9773	0.8327	0.9746	0.8160	0.8102	0.8115
		m: $exp(1)$	0.0000	43.2855	21.3037	0.9833	0.6063	0.9484	0.8145	0.7185	0.8020
	0.649, 0.990	a: $N(0, 5^2)$	0.0000	21.3122	20.8859	0.9773	1.0466	0.9783	0.8160	0.7814	0.7803
		a: $N(0, 3^2)$	0.0000	20.7947	20.5099	0.9773	1.0556	0.9792	0.8160	0.7821	0.7815
		a: $N(0, 1^2)$	0.0000	20.5213	20.1459	0.9773	1.0604	0.9737	0.8160	0.7826	0.7818
		m: $exp(1)$	0.0000	36.0314	23.8184	0.9817	0.8026	0.9565	0.8155	0.7140	0.7752
Daughter	0.954, 0.974	a: $N(0, 5^2)$	0.0000	18.8437	16.5614	0.9719	0.8166	0.9680	0.8171	0.8070	0.8082
		a: $N(0, 3^2)$	0.0000	17.7033	15.6918	0.9719	0.8256	0.9659	0.8171	0.8083	0.8086
		a: $N(0, 1^2)$	0.0000	17.1421	15.1775	0.9719	0.8301	0.9680	0.8171	0.8093	0.8083
		m: $exp(1)$	0.0000	43.4717	21.0365	0.9785	0.6028	0.9445	0.8162	0.7146	0.7976
	0.649, 0.990	a: $N(0, 5^2)$	0.0000	20.1613	20.0763	0.9719	1.0573	0.9872	0.8171	0.7895	0.7862
		a: $N(0, 3^2)$	0.0000	19.6759	19.5498	0.9719	1.0661	0.9834	0.8171	0.7913	0.7904
		a: $N(0, 1^2)$	0.0000	19.4507	19.1871	0.9719	1.0708	0.9803	0.8171	0.7928	0.7911
		m: $exp(1)$	0.0000	35.2381	23.0851	0.9760	0.8087	0.9535	0.8165	0.7229	0.7745

† MSEP: difference between (adjusted) error-prone and error-free predictions.

a: indicates a classical additive model; $T^* = T + \epsilon$ where $\epsilon \sim N(0, \sigma^2)$, m: indicates a multiplicative measurement error model, $T^* = TU$, $U \sim exp(\lambda)$.

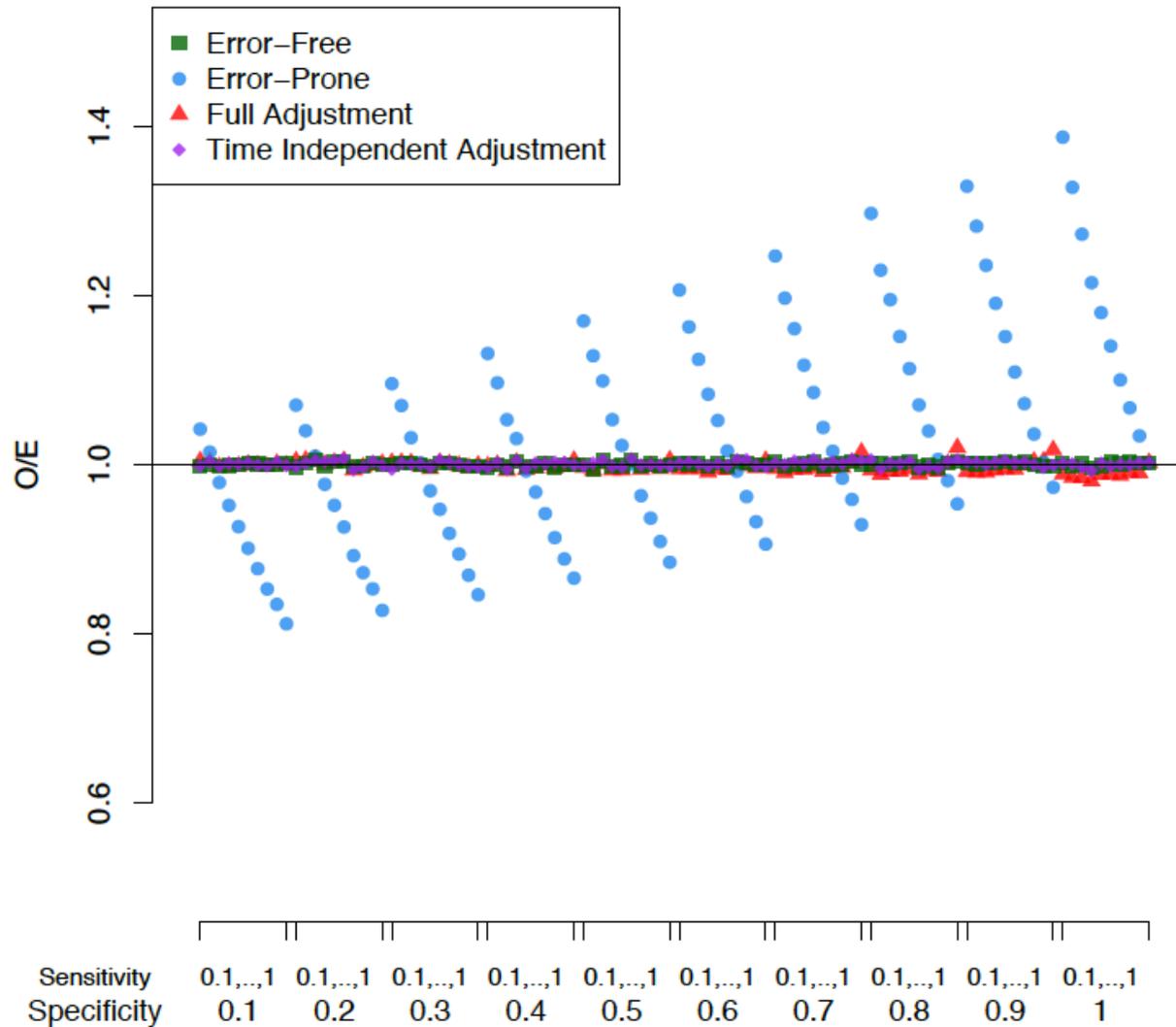
$$O/E = \sum_i I(\gamma_{0i} = 1) / \sum_i \hat{P}_i$$

$$MSEP = n^{-1} \sum_i \left\{ \hat{P}_i - \hat{P}_i(\gamma_0 | H) \right\}^2$$

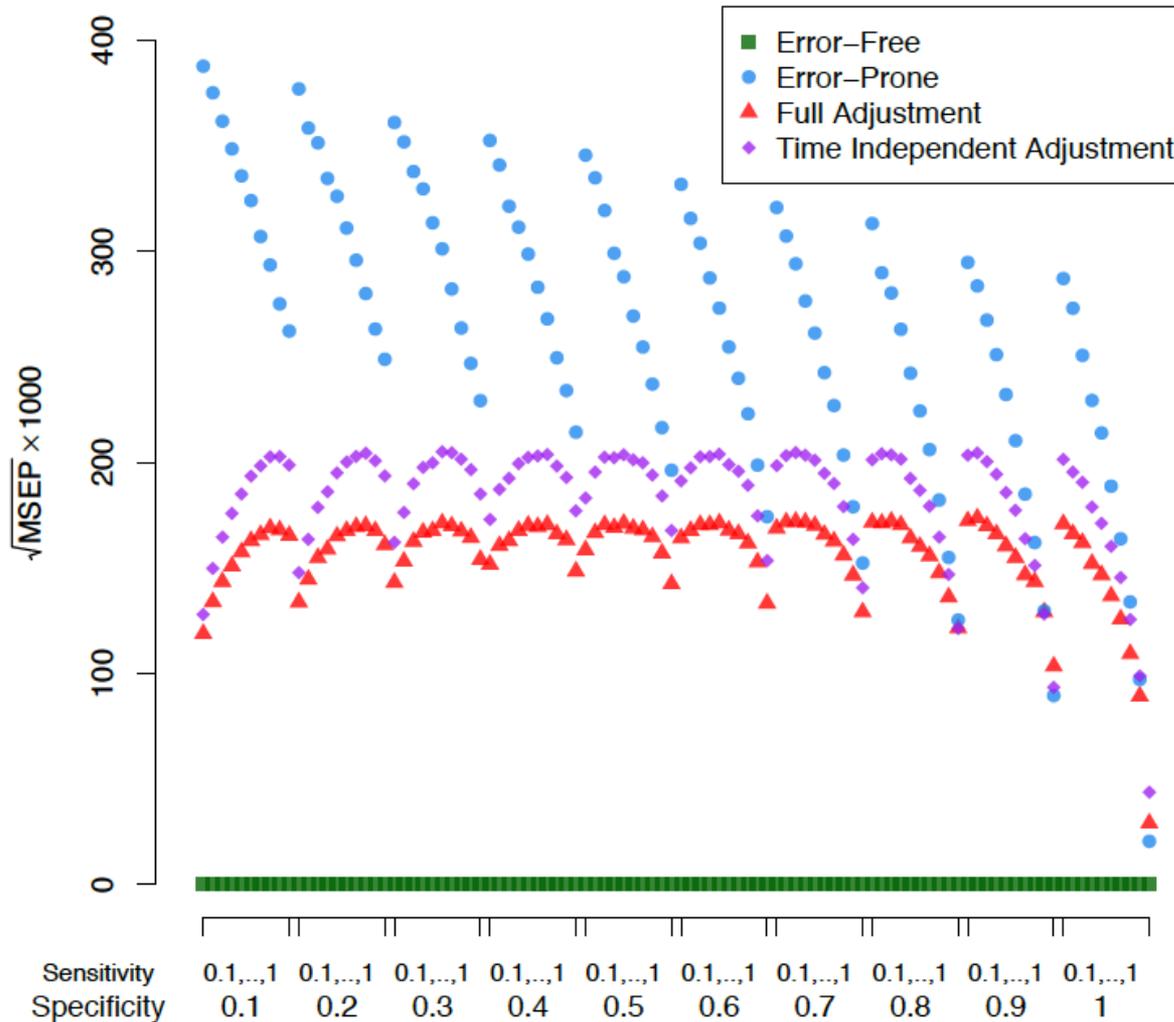
Summary of simulation results

- We are able to eliminate almost all the bias induced by ME in histories (O/E).
- We are able to improve accuracy (MSEP).
- **We are able to improve discrimination (ROC-AUC) only to some degree.**

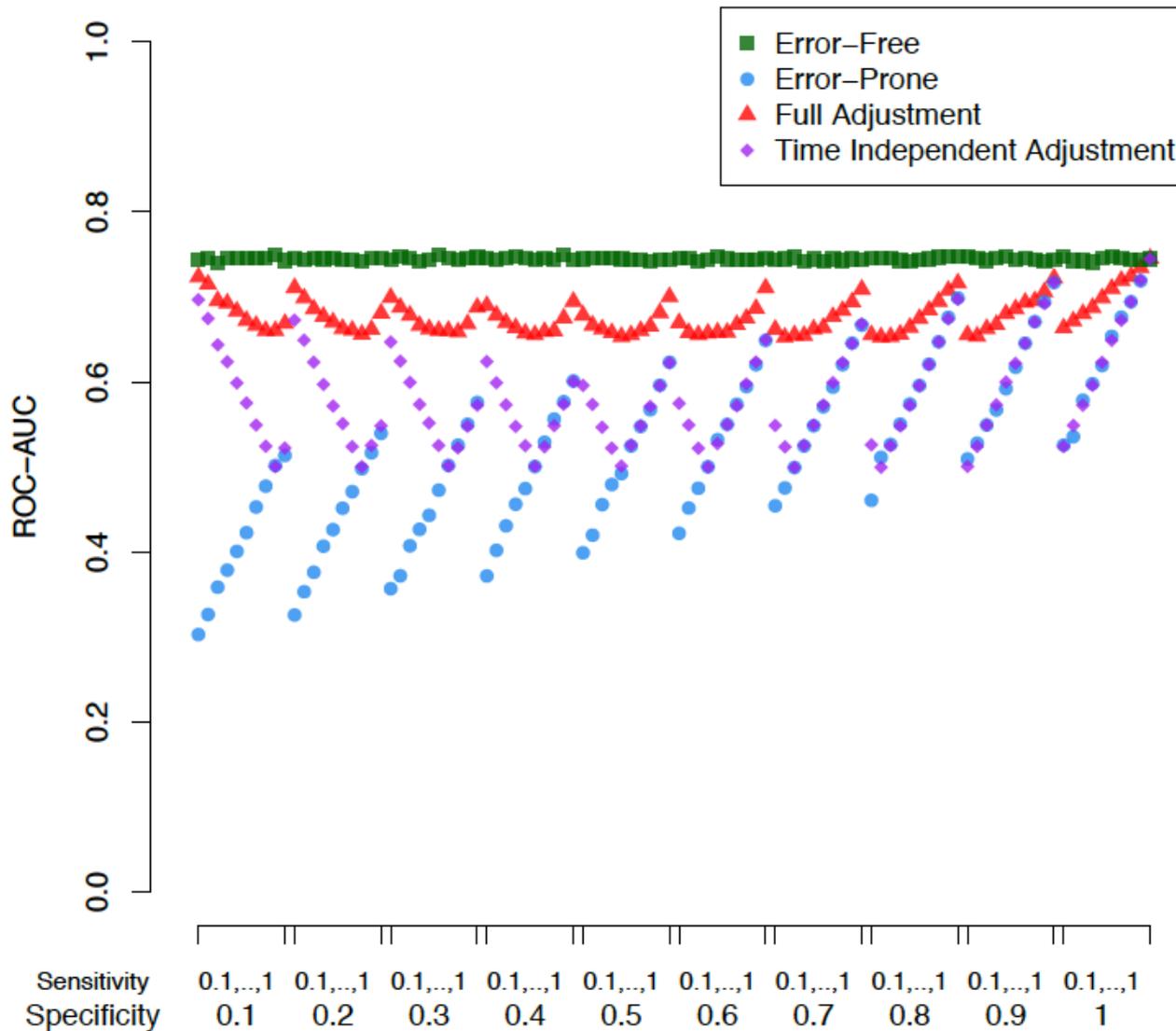
Survival prediction - summary of simulation results (multip.)



Survival prediction - summary of simulation results (multip.)



Survival prediction - summary of simulation results (multip.)



Concluding remarks

- ✓ A non-parametric adjustment is provided, for measurement error in time to event predictor.
- ✓ Ignoring the measurement error, provides miscalibrated models.
- ✓ The proposed adjustment improves calibration and total accuracy.
- ✓ The proposed method can be easily incorporated in BayesMendel R package for direct clinical use.

Model discrimination only partially improved.

END

Example – misreporting breast cancer

Counselees:

- Data from the Cancer Genetics Network (CGN) Model Evaluation Study, with known carrier status.
- 2038 families, 34310 relatives.
- 9.2% of the relatives have breast cancer.
- Only error-prone self-reported family history is available.

Validation data:

- Data from U of California at Irvine (UCI).
- 719 cancer affected counselees (breast, ovarian or colon cancer).
- 1521 female relatives, 19.3% with breast cancer.
- Error-prone and error-free family history are available.

Example – misreporting breast cancer

Log of O/E and 95% confidence intervals for being a BRCA carrier for counsees in CGN dataset, stratified by risk decile:

- Transportability?
- Small sample?

Very small improvement in Brier score and ROC-AUC.

