

# Model Evaluation for Time to Event Studies

Daniel Saggau

11/15/2020

## Introduction

Time to event studies are unique due to the data structure of survival data. A common characteristic of survival data is the right censored nature of observations. Right censored data means that not every subject in the study has an event e.g. a time where the illness occurs or when the patient dies. The reasons for that can be various but one frequent reason for this is that the study ends prior to the event occurring. Another reason could be that we have subjects dropping out of our study, occurring with considerable frequency in clinical or epidemiological studies.

Generally speaking, there are various model evaluation metrics. This paper will introduce a number of methods, focusing on popular extensions of the mean squared error and the area under the curve, adjusted for survival studies namely the brier score and the c-index, two methods that have gained prominence within the realms of scholarship and among clinicians. The c-index does enjoy considerably prominence among clinicians due to interpretability and the ability to make comprehensive conclusions for the individual patients. Irrespective, the c-index only considers discrimination and does not account for calibration as an evaluation criterion is a pivotal shortcoming when assessing prognostic studies. Inevitably, this paper suggests that the integrated brier score is the more holistic model evaluation metric from the stance of a statistician.

The paper is structured as follows: Firstly, i will introduce the integrated brier score, providing further information on the origin, theoretical underpinning and the respective characteristics. Secondly, i will introduce the c-index, following the same procedure as with the integrated brier score. Following the analysis of these two cornerstones of model evaluation for time to event studies, i will illustrate further methods and current research within the field. Subsequently, the fourth section will provide a brief example of how to implement given methods in R. Lastly, the conclusion will summarize core findings of this brief comparative study.

# Model Evaluation metrics

When attempting to evaluate a model performance, there are various approaches. The optimal model evaluation metric inevitably will depend not only on the target but also the target audience at hand. Firstly, one needs to differentiate between the type of study at hand, namely whether we are dealing with ‘prognostic studies’ or ‘diagnostic studies’.

## Diagnostic and Prognostic studies

Prior to talking about different elements of model evaluation, we will briefly talk about the different type of studies relevant for survival data. In a clinical setting, studies can be separated into diagnostic and prognostic studies. Diagnostic studies are concerned with the problem of how to classify a patient at this point in time- Prognostic studies on the other hand are dealing with

Now, we can disentangle the different components of model evaluation. At the core, model evaluation differentiates between discrimination and calibration.

## Discrimination

Discrimination is the ability of a model to handle patients that do not have outcomes accordingly. Inevitably, when controlling for discrimination, we are controlling for how well our model is handling subjects with outcomes as compared to subjects without outcomes. Therefore, as the name suggest, we are testing how strong our model discriminates between subjects that incur an event versus subjects that don’t. As a frequent property of survival data sets is the fact that have right censored data, data that entails subjects without the outcome/event taking place. Henceforth, discrimination is an important pillar for model evaluation in survival analysis. Various measures have emerged, that deal with discrimination such as the c-index or the net reclassification index. Perfect discrimination would imply that all our subjects with the event (e.g. a disease) have higher scores than subjects that do not have an event within their time period. One should note that when using a model that only controls for discrimination, our predictive accuracy could be horrible but as long as this condition holds, inaccurate models could be evaluated falsely as superior.

## Calibration

One should mention that especially when actually applying these methods, in a clinical setting one either deals with diagnostic or prognostic tasks. Diagnostic is the analysis of a given subject at that point in time. For binary classification tasks during diagnostic studies, where we need separate between e.g. patients with and without disease, discrimination is a very important concern and potentially of greater importance. Prognostic on the other hand deals with predictive modeling, predicting e.g. in our survival analysis setting the survival of a patient.

When dealing with an prognostic analysis, calibration can become an important concern. Calibration captures the accuracy of our predictions of our model. The underlying goal is

to ensure that the predictions are as accurate as possible. More general ways to measure calibration are for instance the Hosmer Lemeshow test, the “goodness of fit” test. (Gerds and Schumacher, 2006)

For this very reason, the research community has highlighted the added value of using e.g. the integrated brier score, a score that controls for both discrimination and calibration. To further understand these methods, the following section will briefly introduce the origins of these methods, namely the ROC curve and the brier score.

## **Origin ROC / Concordance -statistics**

The Receiver Operating Characteristic Curve (ROC) is an important model evaluation tool, gaining substantial prominence in various fields of statistics. The ROC was introduced by electrical engineers during world war 2. This method is the foundation of the c-index which is one of the most prominent tools within the field of model evaluation for survival analysis. In a nutshell, the ROC takes into account two factors namely sensitivity and specificity. Firstly, Sensitivity deals with the likelihood of positive test results, specifically it deals with values above the threshold among the subject group which do endure an event e.g. the subjects with diseases (Cook, N. 2007). Sensitivity becomes more volatile when for instance dealing with milder, nuanced cases of a disease. Another name frequently used for Sensitivity is the true positive rate. On the other hand, specificity deals with false negatives, patients with a disease we classify as not having any diseases. However, specificity is especially subject to the influence of the characteristics of a subject without disease. Another name for specificity is the true negative rate. Examples of such characteristics are age or gender. The ROC takes these two factors and plots sensitivity against 1- specificity.

- The area under the curve or the c statistic ranges from 0.5 (no discrimination) to max of 1 (perfect discrimination)
- Essentially, the c statistic is equivalent to the probability that the measure or predicted risk is higher for a case than for a non-case.
- Further, c-statistic describes how well models rank case and noncase; but not a function of actual predicted probabilities

## **Time dependent ROC**

### **C-index**

The c- index is the generalization of the ROC for survival data (Cook, N., 2007). The c statistic is a rank correlation measure, focusing on Kendall’s tau (Uno et al., 2011). An alternative is working with loss functions. Because c statistics is based on ranks it is less sensitive than e.g. measures based on likelihood. Another shortcoming of the usage of rank correlation is ordering of survival times when we dont have a complete data set. Survival data frequently is censored. Therefore modifications of the traditional ROC is needed.

## Example modification

Antolini et al. (2005) propose a time dependent c-index, where discrimination is summarized over time. Their model considered the presence of a population feature rather than a shortcoming of the sample. Irrespective, they also agree with the common consent and propose the use the c-index in symbiosis with a tool to measure calibration for model evaluation.

- time to event version
- integrated AUC
- IDI, NRI
- time dependent measures
- cumulative case dynamic control

For the individual at time t:

$$L(S, t|t^*) = [(S(t^*)^2)I(t \leq t^*, \delta = 1)(\frac{1}{G(t)})] + [((1 - S(t^*))^2)I(t > t^*)(\frac{1}{G(t^*)})]$$

For the population mean:

$$L(S, t|t^*) = \frac{1}{N} \sum_{i=1}^N L(S_i, t_i|t^*)$$

## Integrated Brier Score

Integrated population mean version:

$$L(S, t|t^*) = \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^T L(S_i, t_i|t^*)$$

Uno et al.(2011) propose a modified c-statistic which is consistent for population concordance measures.

## Advantages

The c-index has gained popularity because so interpretable (Kattan and Gerds, 2018).

## Disadvantages

As mentioned above, for a more nuanced prevalence of a disease, the sensitivity is affected and henceforth problematic.(Cook,N., 2007) Specificity is also dependent on the data structure, but as suggested by Cook (2007), specificity is for instance affected by age, gender and the prevalence of concomitant risk factors.

Problem: Studies ignoring calibration (Risk prediction models in cardiovascular literature, use c-statistic, despite working with large prospective cohort studies. Nancy Coook 2007) While for diagnostic studies, discrimination is the most important feature for a model evaluation metric, the same is not true for prognostic studies. Kattan and Gerds (2018) argue that model evaluation metrics also need to differentiate between useless and harmful models. The c-index also does not account for clinical consequences and the subjective importance of false positives relative to false negatives. This problem holds for both the c-index and the brier score but generally speaking, clinical cost are different than specified in these method. The relative ability of a model to perform might not be the only question of importance in model evaluation for clinicians. For instance, other important questions could be whether introducing either model in the first place, or rephrased whether any of these models actually cause a net benefit. Further Kattan and Gerds (2018) argue that model evaluation needs to account for the time horizon.

## Introduction and origin story

The score brier was initially used for weather forecasting. The general version of the brier score is also called prediction error or mean squared error (Schoop et al.,2011; Gerds & Schumacher, 2006). Henceforth, some of the applications e.g. in the ‘pec’ package use other terminology for the brier score.

## Explanation of Method

The mean squared error in a nutshell is the incurred quadratic loss, studying the predicted and the true event status (Schoop et al.,2011). With uni-dimensional predictions the brier score is the same as the mean squared error. Graf et al. (1999) state that the “...expected brier score may be interpreted as a mean squared error of prediction when the estimated prob, which take values in interval  $[0,1]$  are viewed as prediction of event status at  $t, I(T>t)$  in  $\{0,1\}$ .” The brier score is dependent on the evaluation time. By introducing a reweighing scheme, one derives quantities that are independent on the censoring distribution and hence suitable for censored data (Graf et al.,1999) To get a comprehensive understanding of model performance, multiple time points have to be studied.

## Modifications

Wu and Lee (2014) advocate for the usage of the sBrier score. In a nutshell, the sBier score is the mean squared error for the current model divided by the mean squared error for the null model.

## Advantages

The integrated brier score is a measure accounting for both discrimination and calibration. The measure is an attempt to obtain an holistic understanding of the model, rather than just looking at e.g. discrimination. Graf et al. (1999) argue that the method is more sophisticated than the c-index because it deals with probabilities for prediction rather than classifications.

Moreover, Graf et al. (1999) suggest that the concordance statistics is merely a misclassification rate. While the c-index would look at the different labels, true/false positives/negatives, here we are working with likelihoods and can actually compare accuracy of the scores. Another feature of the integrated brier score is the ability to differentiate between useless and harmful models. Harmful models are models that make incorrect predictions while useless models always predict prevalence. With e.g. Harell’s c index, one is unable to differentiate between the two (Kattan and Gerds, 2018).

## Disadvantages

The benchmark of the different models are also dependent on the overall prevalence of the event in our data set. Henceforth, when working with data where the event rarely takes place, the benchmark becomes convoluted (Kattan and Gerds, 2018) Kattan and Gerds (2018) suggest that the evaluation is somewhat problematic with respect to numerous aspects. One pivotal shortcoming of the method is the inability to compare results independent from other models. Kattan and Gerds (2018) argue that one is only able to compare a model compared to other models and henceforth one is always at best only able to see that the one method is superior to the other models at hand. Especially in a practical setting, when undertaking a diagnosis for a patient this is rather impractical as patients don’t come in pairs. Furthermore, this implies that we are unable to see whether the implementation of the model is advisable in the first place. Steyerberg et al. (2010) further extend this argument, by arguing that one is unable to detect whether the implementation will cause more harm than benefit. Therefore, some scholars have advocated for complementary tests, to assess the overall profitability of implementing such measures, accounting for clinical consequences, or in other words more realistic weights given the preferences of clinicians. One should note that as perhaps not obvious, clinicians don’t value the different components of model accuracy equally. A false negative can be more detrimental than a false positive. For example, sending a sick patient home, in the believe of being healthy can do more damage than testing further on a healthy patient, that has been tested sick.

## Reclassification

Steyerberg, E. W., Vickers, A. J., Cook, N. R., Gerds, T., Gonen, M., Obuchowski, N., ... & Kattan, M. W. (2010)

- Change is risk stratification.
- Use observed incidence of events of the reclassification table to predicted probabilities of the orgn. model.
- Cook proposing variant og hosmer lemeshow statistic within the reclassified categories, leading to chi-squared statistic.

**Net Reclassification improvement** Cook (2008) argues that Net reclassification improvement (NRI) and calibration tests for cross classified categories can be used to study the clinical usefulness. While NRI is only a measure to study discrimination, it allows to account for the formation of categories based on clinical risk estimates. Therefore, this measure is

also focusing on the clinical application, rather than holistic model evaluation. Henceforth, reclassification might just complement existing clinicians in practical applications as opposed to providing a dominant model evaluation tools.

**Cook, N. R., & Ridker, P. M. (2009). Advances in measuring the effect of individual predictors of cardiovascular risk:**

- Integrated discrimination improvement
- IDI is equivalent to testing whether the regression coefficient in a model is equal to zero (similar to  $R^2$  or the proportion of variance explained)
- The NRI and the IDI both condition on the case-control or later disease status
- don't provide information on calibration of the estimated risk
- A limitation of NRI and other reclassification measures is that they depend on the particular categories used
- The calibration test seems to depend somewhat less on the number of categories since the degree of freedom adjust for the number of categories
- Suggest that reclassification calibration statistic and NRI may be useful in demonstrating the ability of new models and markers to change risk strata and alter treatment decision

Pepe et al. (2015) suggestion caution when evaluating models solely based on

## **Decision Analysis Curve**

One fundamental problem of the methods that we have introduced is that it does not really accommodate the interests of clinicians. From the perspective of a clinician, giving false positive and false negatives the same weight does not make any sense. A false negative entails severe repercussions relative to the false positive. For instance let's say we result in a false negative for a cancer patient, the patient is harmed to a detrimental extent and deprived of the opportunity to undertake earlier action. Moreover, these methods do not really tell us whether introducing the new model creates added value. Further, preferences may differ from a clinical standpoint. E.g. sensitivity and specificity are frequently unequal in importance to a clinician. Henceforth, scholars have proposed a new complementary framework to assess the net benefit of a model, providing a tool to assess whether implementing a new model is worth it in the first place (Vickers, A., Elkin, E., 2006). Decision analysis curve enables the use of weights, allowing optimal decision making based on subjective preferences, embodied in a net benefit equation. Vickers et al. (2016) interpret those net benefits as "clinical consequences". Further Vickers et al. (2016) illustrate that harm is transformed, using an exchange rate to put harm and benefit on one scale. This exchange rate can be obtained by asking clinicians questions based on their subjective preferences such as how many patients they would have undergo a biopsy prior to finding a cancer or weighing the benefits of getting early findings as opposed to the cost of harmful further testing. Together these elements build the net-benefit equation. Plotting different exchange rates with the net benefit equation, gives us the decision analysis curve. The curves enable the practitioner the identification of the range of threshold probabilities for when a model would be of value, providing information on the necessary benefits needed for a model to be useful and which of many models is optimal (Vickers, A., Elkin, E., 2006).

One important consideration is that decision analysis curve is a complement, not a substitute to existing models (Vickers, A., Elkin, E., 2006).

**Vickers, A. J., & Elkin, E. B. (2006). Decision curve analysis: a novel method for evaluating prediction models.**

- Compare their method to AUC method claiming that:
- AUC metric focuses solely on predictive accuracy of model
- Cannot tell us whether a model is worth using at all or which of two models is preferable
- AUC does not provide insight into usefulness aka does not account in their model that clinician may have other interests
- Two general problems: require data such as on cost or quality adjusted life years, not found in the validation data set. Cannot be evaluated without further information
- Secondly: decision analysis typically requires test or prediction model evaluated give a binary result s.t. the true and false positive and negative results can be estimated (but prediction is frequently continuously expressed)
- Interpretation requires understanding of the liking of the patient \*
- The proposed method does not require obtaining information regarding treatment preferences but need theoretical relation for the threshold probability of disease and the relative value of false positive and false negative results

## Other proposals

Wu, Y. C., & Lee, W. C. (2014)

### Lorenz Curve

#### Gini

Mean separation for the current model divided by the mean separation from an error free model

#### Pietra

Mean gain for current model divided by the mean gain for an error free model

## Implementation

Notable packages:

- ‘pec’
- ‘survival’ concordance() function
- ‘survIDINRI’



## Discrimination Plot

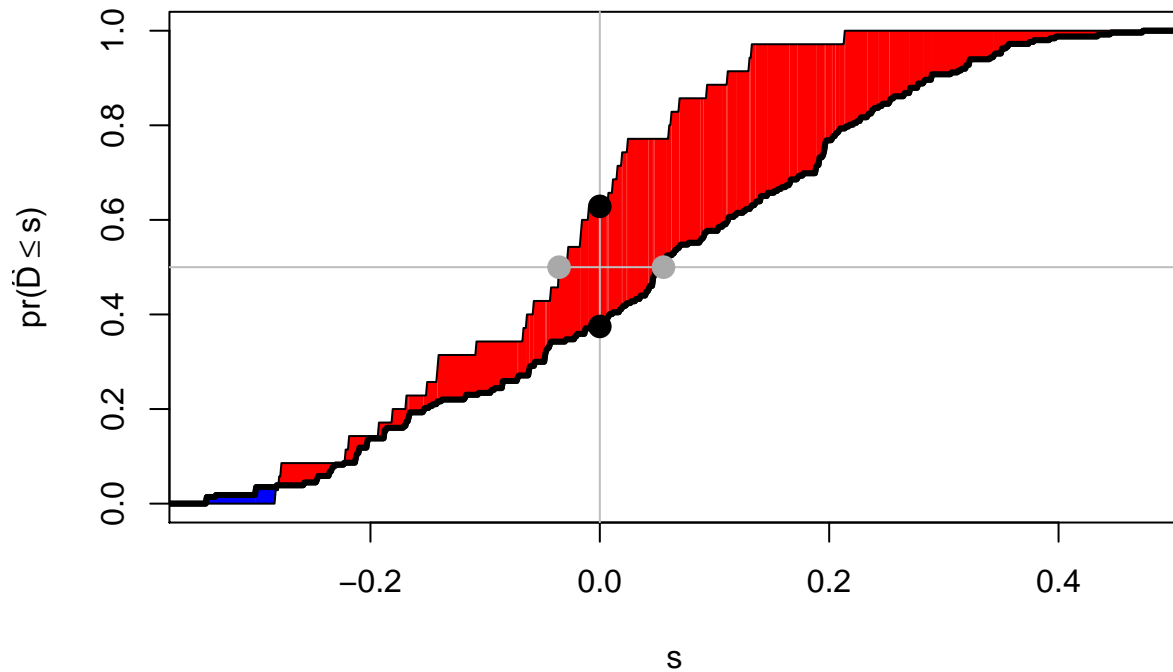
A c-index above the threshold of 0,8 can be considered good (Zhang et al.,2018).

Researchers have combined the usage of reclassification tools with discrimination measures.

### Net reclassification and integrated discrimination improvement implementation

One example package dealing

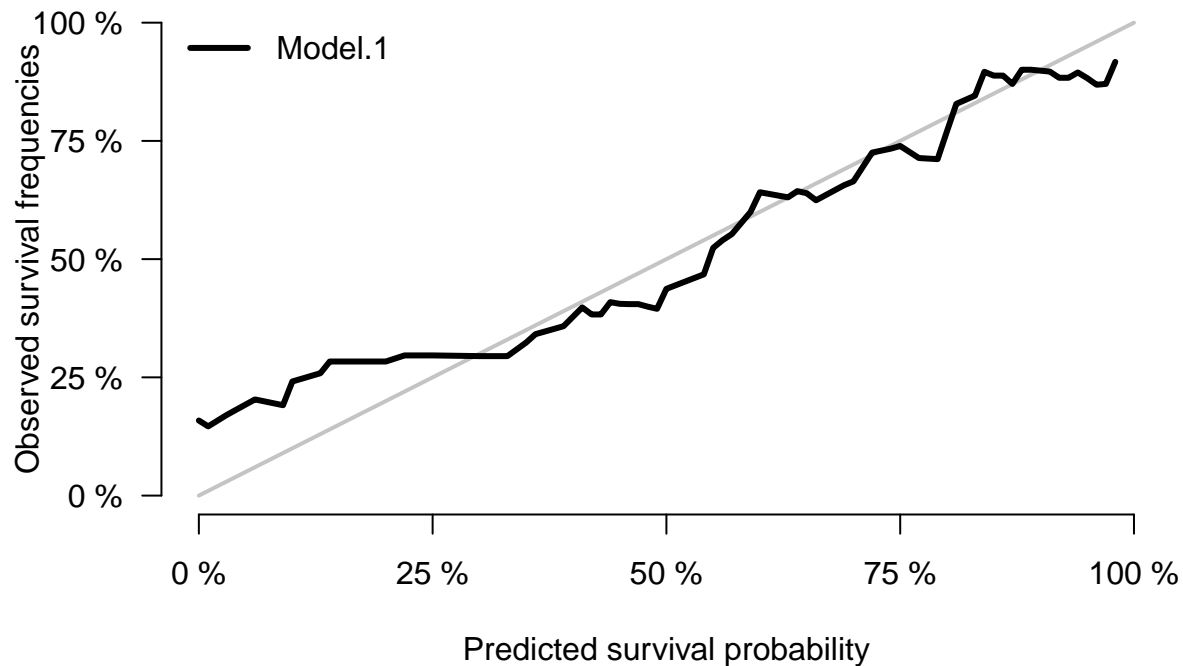
```
IDI.INF.GRAPH(res.IDI.INF)
```



## Calibration Plot

We can use calibration plots to visualize the calibration of our model. The ‘pec’ packages provides the ‘calPlot’ function.

```
calPlot(pmodel)
```



## Prediction Error plot

### mlr 3 implementation

Benchmark

- mlr3 proba:

Loss Functions/Calibration:

- Integrated Graf Score
- Integrated Log Loss
- Log Loss

Discrimination: \* Uno's AUC \* Song and Zhou's AUC

miscellaneous:

- van Houwelingen's Alpha Calibration
- van Houwelingen's Beta Calibration

## Conclusion

Time to event studies require adjusted model evaluation tools for censored survival data. At the core, studies separate between models that evaluate overall performance, discrimination and calibration. New methods such as reclassification and clinical usefulness have gained prominence among scholarship within recent research, but did not achieve the same level of recognition among clinicians and in the applied research community.

## References

- Assel, M., Sjoberg, D. D., & Vickers, A. J. (2017). The Brier score does not evaluate the clinical utility of diagnostic tests or prediction models. *Diagnostic and prognostic research*, 1(1), 1-7.
- Antolini, L., Boracchi, P., & Biganzoli, E. (2005). A time-dependent discrimination index for survival data. *Statistics in medicine*, 24(24), 3927-3944.
- Bender, A., Rügamer, D., Scheipl, F., & Bischl, B. (2020). A General Machine Learning Framework for Survival Analysis. *arXiv preprint arXiv:2006.15442*.
- Cook, N. R. (2007). Use and misuse of the receiver operating characteristic curve in risk prediction. *Circulation*, 115(7), 928-935.
- Cook, N. R. (2008). Statistical evaluation of prognostic versus diagnostic models: beyond the ROC curve. *Clinical chemistry*, 54(1), 17-23.
- Cook, N. R., & Ridker, P. M. (2009). Advances in measuring the effect of individual predictors of cardiovascular risk: the role of reclassification measures. *Annals of internal medicine*, 150(11), 795-802.
- Gerds, T. A., & Schumacher, M. (2006). Consistent estimation of the expected Brier score in general survival models with right-censored event times. *Biometrical Journal*, 48(6), 1029-1040.
- Gerds, T. A., Cai, T., & Schumacher, M. (2008). The performance of risk prediction models. *Biometrical Journal: Journal of Mathematical Methods in Biosciences*, 50(4), 457-479.
- Gerds, T. A. (2019). Package ‘pec’.
- Graf, E., Schmoor, C., Sauerbrei, W., & Schumacher, M. (1999). Assessment and comparison of prognostic classification schemes for survival data. *Statistics in medicine*, 18(17-18), 2529-2545.
- Harrell Jr, F. E., Lee, K. L., & Mark, D. B. (1996). Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Statistics in medicine*, 15(4), 361-387. Chicago
- Heagerty, P. J., Lumley, T., & Pepe, M. S. (2000). Time-dependent ROC curves for censored survival data and a diagnostic marker. *Biometrics*, 56(2), 337-344.
- Kattan, M. W., & Gerds, T. A. (2018). The index of prediction accuracy: an intuitive measure useful for evaluating risk prediction models. *Diagnostic and prognostic research*, 2(1), 7.
- Pencina, M. J., D’Agostino Sr, R. B., D’Agostino Jr, R. B., & Vasan, R. S. (2008). Evaluating the added predictive ability of a new marker: from area under the ROC curve to reclassification and beyond. *Statistics in medicine*, 27(2), 157-172
- Potapov, S., Adler, W., & Schmid, M. (2012). *survAUC: Estimators of prediction accuracy for time-to-event data*. R package version, 1-0.

- Schoop, R., Beyersmann, J., Schumacher, M., & Binder, H. (2011). Quantifying the predictive accuracy of time-to-event models in the presence of competing risks. *Biometrical Journal*, 53(1), 88-112.
- Steyerberg, E. W., & Vickers, A. J. (2008). Decision curve analysis: a discussion. *Medical Decision Making*, 28(1), 146-149.
- Steyerberg, E. W., Vickers, A. J., Cook, N. R., Gerds, T., Gonen, M., Obuchowski, N., . . . & Kattan, M. W. (2010). Assessing the performance of prediction models: a framework for some traditional and novel measures. *Epidemiology (Cambridge, Mass.)*, 21(1), 128.
- Steyerberg, E. W. (2019). *Clinical prediction models*. Springer International Publishing.
- Therneau, T. M., & Lumley, T. (2014). Package ‘survival’. *Survival analysis* Published on CRAN, 2, 3.
- Uno, H., Cai, T., Pencina, M. J., D’Agostino, R. B., & Wei, L. J. (2011). On the C-statistics for evaluating overall adequacy of risk prediction procedures with censored survival data. *Statistics in medicine*, 30(10), 1105-1117.
- van Houwelingen, H. C. (2000). Validation, calibration, revision and combination of prognostic survival models. *Statistics in medicine*, 19(24), 3401-3415.
- Vickers, A. J., & Elkin, E. B. (2006). Decision curve analysis: a novel method for evaluating prediction models. *Medical Decision Making*, 26(6), 565-574.
- Vickers, A. J., Van Calster, B., & Steyerberg, E. W. (2016). Net benefit approaches to the evaluation of prediction models, molecular markers, and diagnostic tests. *bmj*, 352, i6.
- Wang, P., Li, Y., & Reddy, C. K. (2019). Machine learning for survival analysis: A survey. *ACM Computing Surveys (CSUR)*, 51(6), 1-36. Chicago
- Wu, Y. C., & Lee, W. C. (2014). Alternative performance measures for prediction models. *PLoS One*, 9(3), e91249.