

Model assessment

23 Feb 2023

Table of contents

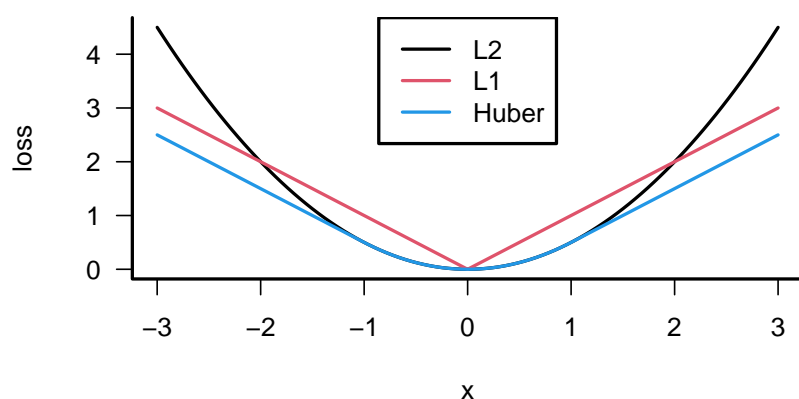
loss functions (regression/quantitative outcome)	2
loss functions (classification)	2
a short rant about categorical loss functions	2
from loss functions to model quality metrics	3
quality metrics	4
within- and out-of-sample error	4
effective number of parameters	5
errors	5
selection vs assessment	6
train-validation-test	6
calibration	6

```
## use help("image-methods", "Matrix")
## lattice graphics: ?lattice:xyplot for details on scales
ifun <- function(x, title = "", ck = FALSE, raster = TRUE) {
  image(Matrix(x),
    sub = "", xlab = "", ylab = "",
    colorkey = ck,
    aspect = "fill",
    scales = list(x = list(draw = FALSE),
                  y = list(draw = FALSE)),
    main = title,
    useRaster = raster
  )
}
```

loss functions (regression/quantitative outcome)

- continuous: L2, L1, **Huber** loss:

```
par(las = 1, bty = "l", lwd = 2)
huber <- function(x, d) ifelse(abs(x)<d, x^2/2, d*abs(x)-d/2)
curve(x^2/2, from = -3, to = 3, ylab = "loss")
curve(abs(x), add = TRUE, col = 2)
curve(huber(x, 1), add = TRUE, col = 4)
legend("top", c("L2", "L1", "Huber"), col = c(1, 2, 4), lty = 1)
```



loss functions (classification)

- 0-1
- **deviance**: $-2 \sum I(G = k) \log \hat{p}_k = -2 \log\text{-likelihood}$
- deviance generalizes to other distributions

a short rant about categorical loss functions

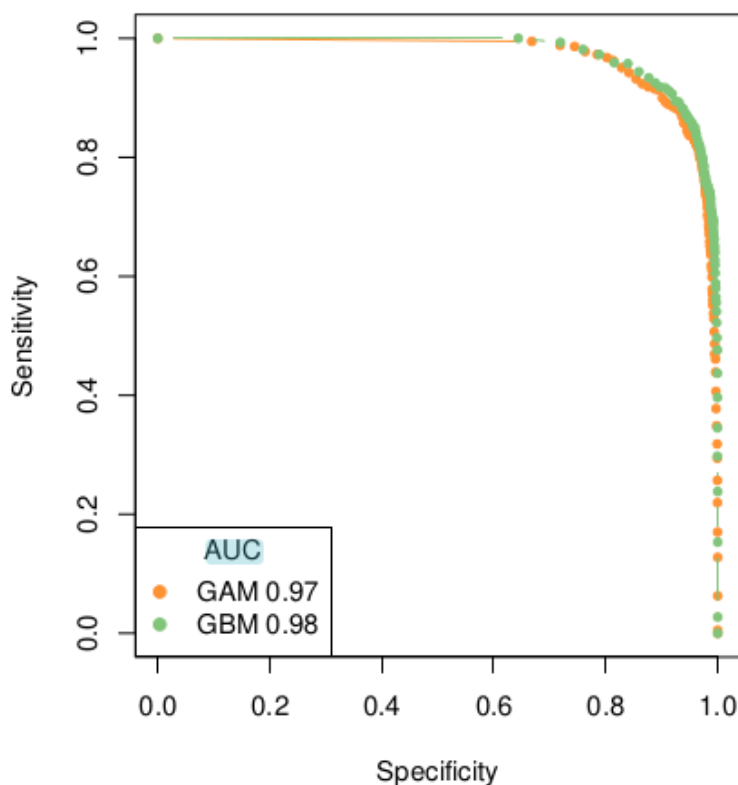
- 0-1 scoring dichotomizes prematurely
- leads to lots of confusing discussion about balancing data sets
- lots of discussion of what to do about imbalanced data sets (SMOTE etc.) (Goorbergh et al. 2022)
- when **should** we balance?

Goorbergh, Ruben van den, Maarten van Smeden, Dirk Timmerman, and Ben Van Calster. 2022. "The Harm of Class Imbalance Corrections for Risk Prediction Models: Illustration and Simulation Using Logistic Regression." *Journal of the American Medical Informatics Association*, June, ocac093. <https://doi.org/10.1093/jamia/ocac093>.

- when we have to use 0-1 scoring for some technical reason
- when we have too **much** data (downsampling, i.e., throw away majority class)
- (cf. discussion of variable selection)

from loss functions to model quality metrics

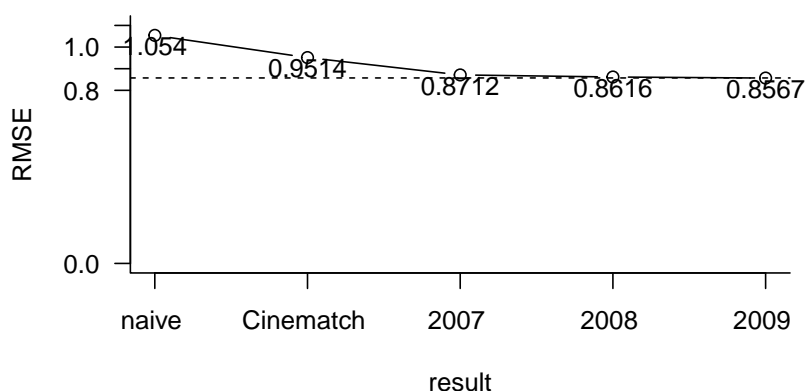
- categorical predictors:
- accuracy (total fraction correct); same problems as 0-1 classification
- AUC (area under the curve)
 - may be problematic in terms of implied misclassification costs? (Hand 2009)



Hand, David J. 2009. “Measuring Classifier Performance: A Coherent Alternative to the Area Under the ROC Curve.” *Machine Learning* 77 (1): 103–23. <https://doi.org/10.1007/s10994-009-5119-5>.

quality metrics

- some combination of loss functions per point
- scaled for **interpretability**
 - how good is good enough?
 - how much difference in model predictions matters?
 - e.g. [Netflix prize](#)
- R^2
- MSE \rightarrow RMSE \rightarrow scaled RMSE (or mean-squared log error?)
- always a **business** or **scientific** decision (*value of information*)



within- and out-of-sample error

- within-sample: R^2
- out-of-sample: adjusted R^2 (scaled by $n - p$), PRESS (predicted out-of-sample error) = LOOCV SSQ, AIC ($-2 \log L + 2p$), Mallows' C_p ($\frac{RSS + 2p\hat{\sigma}^2}{n}$). AIC and C_p equivalent for Gaussian models. AIC asymptotically \rightarrow LOOCV for linear models.
- ESL gives a weird/unusual (scaled-by- N) definition of AIC
- GCV, AUC
- BIC: $-2 \log L + (\log N)p$; higher penalty for $N > e^2$ (almost always)

- derived from a **Laplace approximation** to the **Bayes factor** (quadratic approx; \approx multivariate normal posterior) given equal priors on models

- BIC is **consistent**, AIC is **predictive** (Yang 2005)

effective number of parameters

- (generalized, penalized) linear models: $\text{trace}(\hat{H})$
- additive-error models: $\sum (\text{Cov}(\hat{y}_i, y) / \sigma_\epsilon^2)$

Yang, Yuhong. 2005. “Can the Strengths of AIC and BIC Be Shared? A Conflict Between Model Identification and Regression Estimation.” *Biometrika* 92 (4): 937–50. <https://doi.org/10.1093/biomet/92.4.937>.

errors

- test error (generalization error): prediction error over a **fixed** independent sample
- **expected** prediction error: test error averaged over test sets
- training error (within-sample): expectation

Journal of the American Medical Informatics Association, 2022, Vol. 00, No. 0

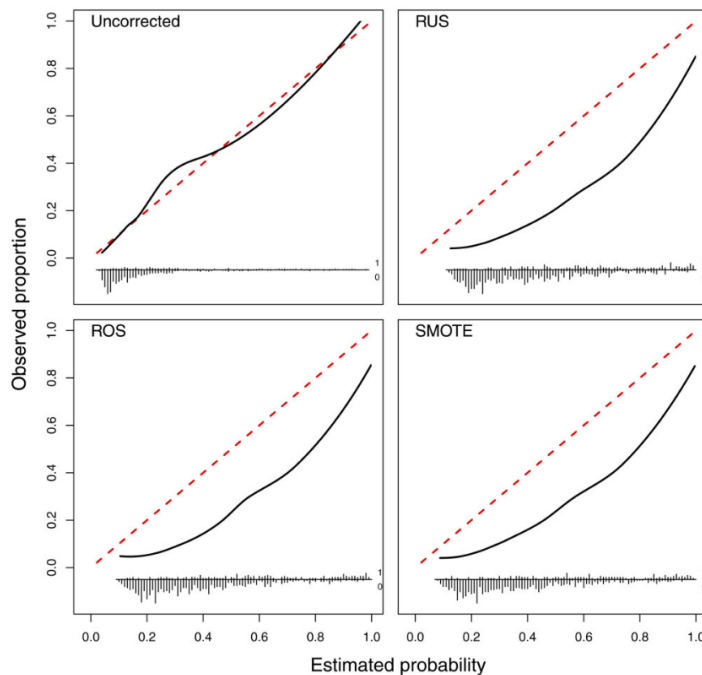


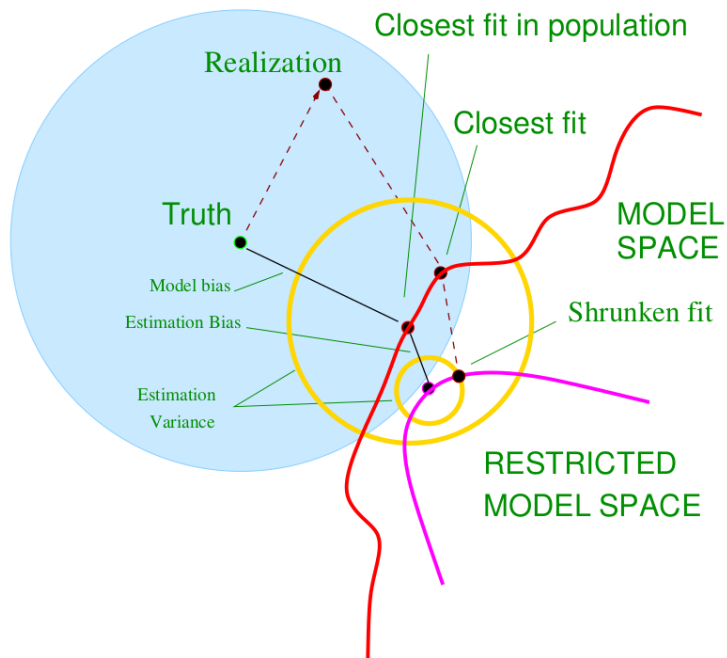
Figure 2. Flexible calibration curves on the test set for the Ridge models to diagnose ovarian cancer.

selection vs assessment

train-validation-test

$$E[f(x_0) - x_0^\top \beta^*]^2 + E[x_0^\top \beta^* - E x_0^\top \hat{\beta}_\alpha]^2$$

- estimation bias = 0 for linear regression etc., positive for ridge etc.



- in-sample error:
 - $C_p = \text{err} + 2 d/N\sigma_\epsilon^2$
- leakage:
 - non-independence
 - data-dependence of training
- jackknife, bootstrap etc.

calibration