

Measuring The Performance Of Medical Screening Tests: Average Precision Vs. Area Under ROC Curve

Wanhua Su^{†a}, Yan Yuan^{†b} and Mu Zhu^c



[†]These authors contributed equally to this project ^aDepartment of Mathematics and Statistics, Grant MacEwan University, Edmonton, Canada, ^bSchool of Public Health, University of Alberta, Edmonton, Canada, ^cDepartment of Statistics and Actuarial Science, University of Waterloo, Waterloo, Canada

Motivation

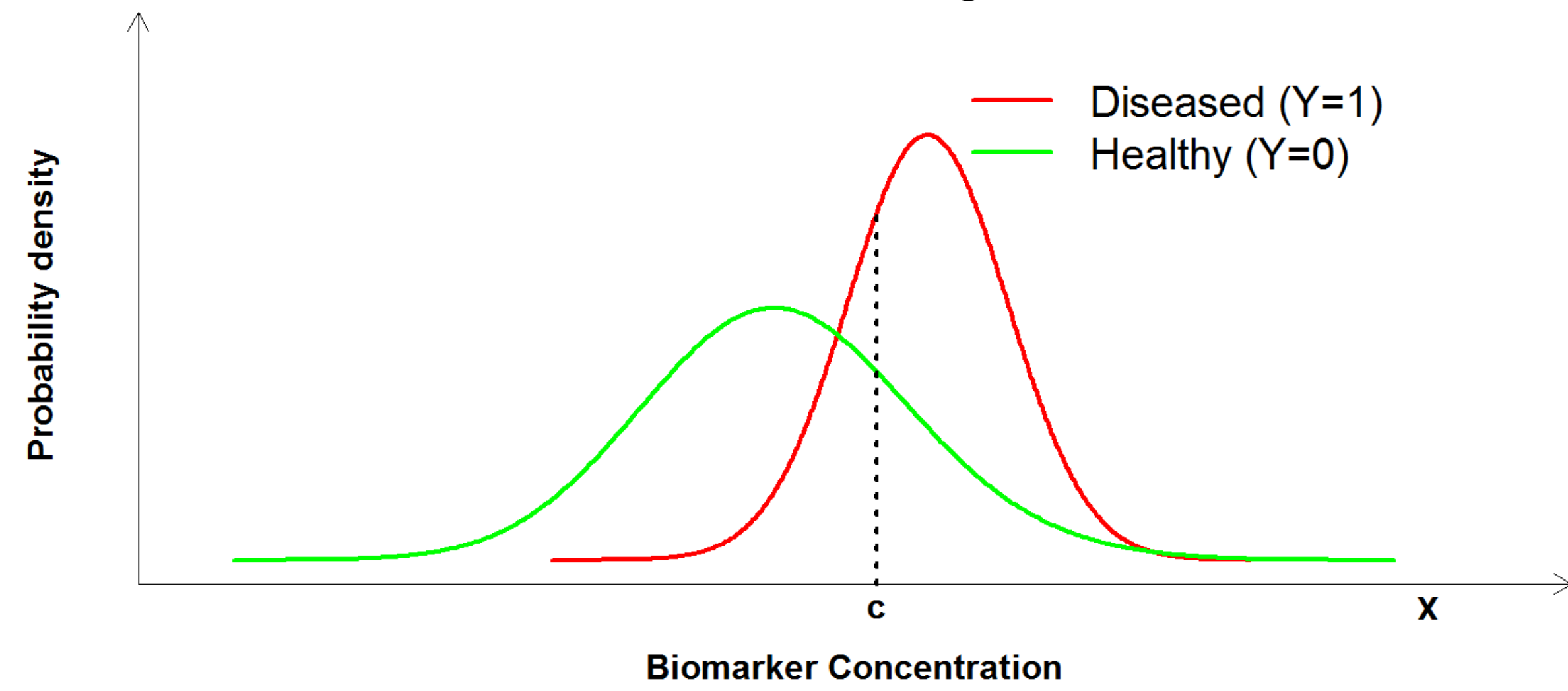
Early detection of disease in its asymptomatic stage could greatly enhance the likelihood of cure. For example, routine breast, prostate and colon cancer screening for certain age groups are recommended by medical experts in North America. The objective of medical screening tests is to detect from the underlying population the diseased asymptomatic subjects who make up a very small proportion (typically < 0.5%). Because the false positive results from screening is considered one major harm of this type of tests¹, the ability to identify diseased subjects with minimal false positive findings should be an important consideration when performance of a screening test is evaluated.

In practice, the Receiver Operating Characteristic (ROC) curve has been the primary performance measure for medical tests including the screening tests. The Area Under the ROC Curve (AUC) provides a single numerical summary of the ROC curve, thus is used as the basis of inferential statistics for comparing ROC curves. We introduce here a performance measure called Average Precision (AP) as it emphasizes more on early detection. AP is a summary measure of the so-called hit curve, which is analogous to ROC curve.

The ROC and Hit curve

Through an example, let's look at the hit curve and the ROC curve, and the definition of AUC and AP.

Suppose that in asymptomatic diseased ($Y=1$) subjects, the distribution of a biomarker concentration shifts to the right, as illustrated below



Using a threshold c , the hit curve plots $P(X>c, Y=1)$ vs. $P(X>c)$, and the ROC plots true positive fraction $P(X>c|Y=1)$ vs. false positive fraction $P(X>c|Y=0)$ for the entire set of thresholds.

Notations: $r = P(X>c)$, the hit function $h(r) = P(X>c, Y=1)$, and the proportion of diseased subjects $\pi = P(Y=1)$

Relationship between AP and AUC

Continuous framework:

If we approximate the hit curve by a quasi-concave curve and let β be the initial true positive rate of the underlying test, we can show that

$$\widehat{AP}(h) \approx \beta \times \widehat{AUC}(h),$$

where $\widehat{AP}(h)$ and $\widehat{AUC}(h)$ are re-scaled to lie between 0 and 1 for any hit (or ROC) curve h .

This approximation implies that if two medical screening tests have the same AUC, then the AP will "reward extra points" to the one with the larger β , i.e. AP places more emphasis on the initial part of the hit (ROC) curve.

Discrete framework:

A radiology screening test typically generates ordinal test results, which use categories such as highly suspicious, possibly malignant, possibly benign, etc. We consider the setup with a multinomial distribution as described in Table 1.

Table 1: A screening/diagnostic test partitions n subjects into K groups (K distinct scores). The broken bars ($|$) indicates the case where all those with scores $\geq x_k$ are declared to be test positive (class 1), while all those with scores $< x_k$ are declared to be test negative (class-0).

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

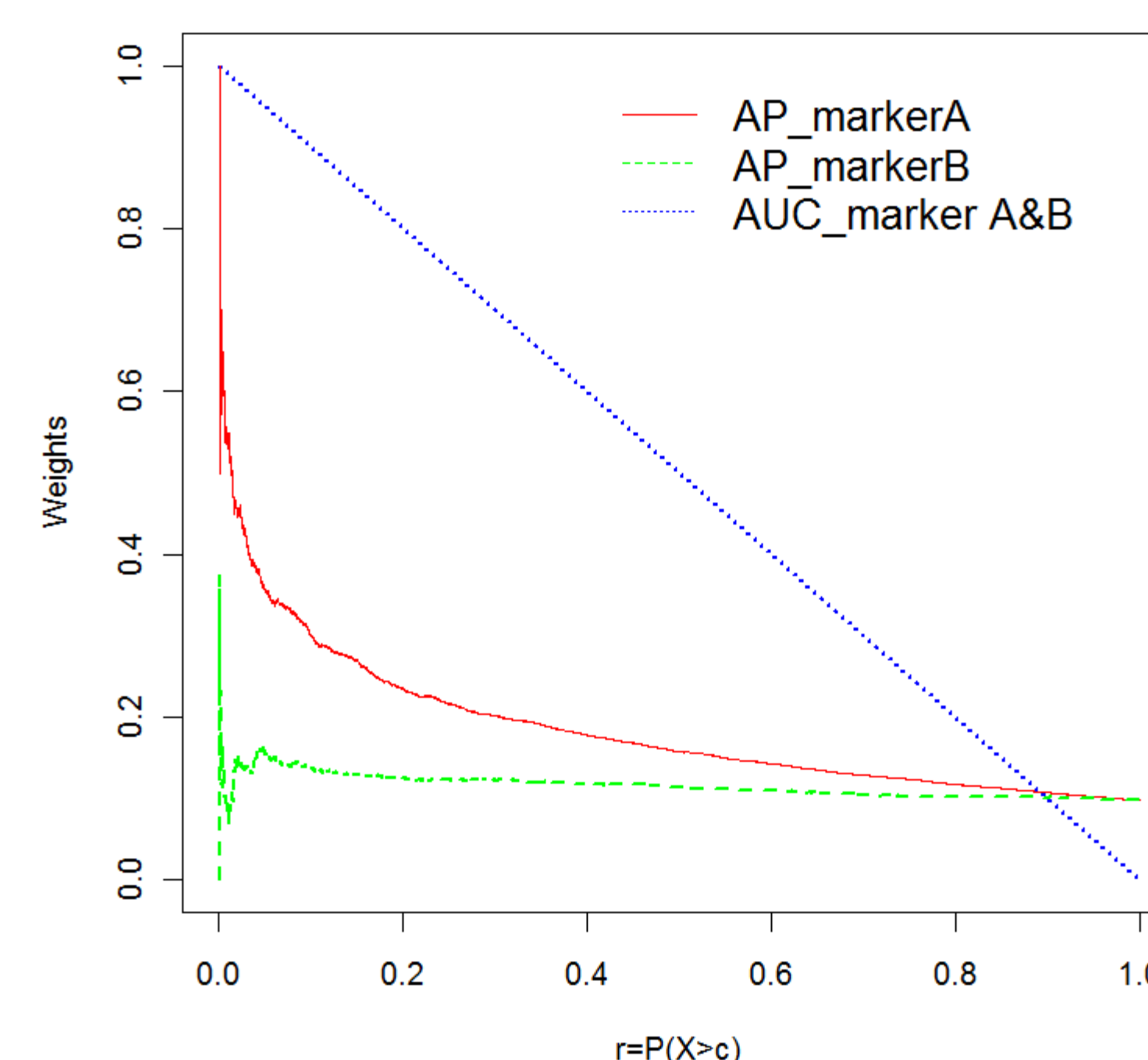
R_k = the k^{th} region in the space of test scores; S_k = total number of subjects in R_k ; Z_k = total number of class-1 subjects in R_k ; \bar{Z}_k = total number of class-0 subjects in R_k .

We can show that

$$\begin{aligned} AP &= \underbrace{\frac{Z_1}{S_1}}_{w_1} \left[\frac{Z_1}{n_1} \right] + \underbrace{\frac{Z_1 + Z_2}{S_1 + S_2}}_{w_2} \left[\frac{Z_2}{n_1} \right] + \dots + \underbrace{\frac{Z_1 + Z_2 + \dots + Z_K}{S_1 + S_2 + \dots + S_K}}_{w_K} \left[\frac{Z_K}{n_1} \right] \\ &= \sum_{k=1}^K w_k \left[\frac{Z_k}{n_1} \right] \\ AUC &= \underbrace{\frac{S_1 + S_2 + \dots + S_K}{n}}_{w'_1} \left[\frac{Z_1}{n_1} \right] + \underbrace{\frac{S_2 + \dots + S_K}{n}}_{w'_2} \left[\frac{Z_2}{n_1} \right] + \dots + \underbrace{\frac{S_K}{n}}_{w'_K} \left[\frac{Z_K}{n_1} \right] - \frac{1}{2} \left(\frac{n_1}{n_0} \right) \\ &= \sum_{k=1}^K w'_k \left[\frac{Z_k}{n_1} \right] - \frac{1}{2} \left(\frac{n_1}{n_0} \right) \end{aligned}$$

These weights, w_k and w'_k , again show that AP places an emphasis on the ability of a test to give the diseased subjects a high score.

Figure 1: Weights for AP and AUC in a simulated example. The concentration of markers A and B in healthy subjects $\sim N(0,1)$, the concentration of markers A and B in diseased subjects $\sim N(1, 1)$ and $N(0.25, 1)$, respectively. The proportion of diseased (π) subjects is 0.1.



By the simulation design, marker A is better than marker B for detecting the diseased subjects. Figure 1 shows that unlike AUC, the weights of AP favors marker A. Thus, AP favors marker A more than AUC does.

Asymptotic Variance of AP

In order to use AP as an evaluation metric, we need to know its variance. The data in each row of Table 1 follows a multinomial distribution, and the two multinomial distributions are independent given π , the proportion of diseased subject. Therefore, the log-likelihood function is

$$\ell(p_k, q_k, \pi) = \sum_{k=1}^K z_k \log p_k + \sum_{k=1}^K \bar{z}_k \log q_k + [n_1 \log \pi + (n - n_1) \log(1 - \pi)],$$

subject to the constraints

$$\sum_{k=1}^K p_k = 1 \quad \text{and} \quad \sum_{k=1}^K q_k = 1,$$

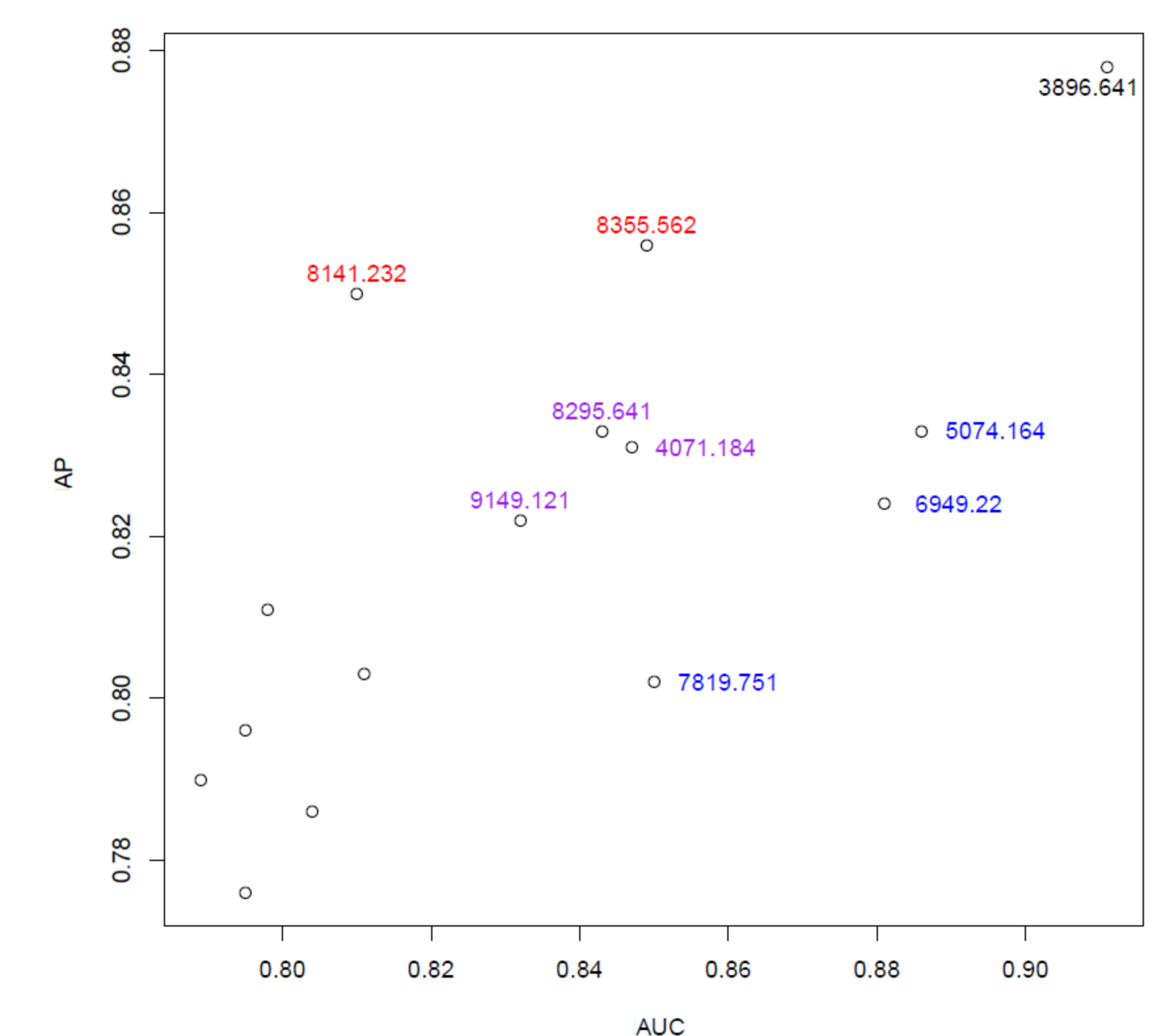
AP is a function of these parameters, i.e. $AP = f(p_k, q_k, \pi)$. We obtain the asymptotic variance of AP by using the delta method,

$$\text{var}(AP) = (\nabla f)^T J^{-1} \nabla f,$$

where J is the expected information matrix.

Examples

Continuous score for potential biomarkers: a case-control study ($\pi=0.5$)
Figure 2: AP vs. AUC of the top 15 potential biomarkers for prostate cancer.²



We see that AP and AUC rank the biomarkers differently.

Ordinal score for a screening test ($\pi=0.00784$)

Table 2: The film vs. digital mammography for breast cancer screening³

Breast cancer diagnosis (455 day follow-up)	Malignancy scores	
	\widehat{AUC} (s.e.)	\widehat{AP} (s.e.)
Film mammography	0.735 (0.012)	0.166 (0.017)
Digital mammography	0.753 (0.012)	0.144 (0.018)

Discussion:

- Among tests that have similar AUCs, the test that finds the largest proportion of diseased subjects with minimal false positives in the early part of the hit (ROC) curve will have the highest AP. It is consistent with the goal of screening. Therefore, AP is a more relevant performance measure for a screening test than the AUC.
- AUC is $Pr(X_D > X_{\bar{D}} | \text{a randomly selected pair of diseased and healthy subjects})$, where X_D and $X_{\bar{D}}$ are the test scores for the diseased and healthy subjects in the randomly selected pair, respectively. It is a conditional probability, and ignores an important parameter π , the proportion of diseased subjects in the asymptomatic population.
- The relative small numerical value of AP for a screening test (Example 2) is advantageous. Large values of AUC give clinicians and patients a false sense of accuracy of the test results, which aggravates the harm of a screening test.⁴ Thus, the relative small valued single summary measure AP offers a useful alternative to summarize the test performance.

References:

- Fletcher SW (2011) Breast Cancer Screening: A 35-year perspective. *Epidemiologic Reviews*. 33: 165-175.
- Wang Z. and Change Y. (2011) Marker selection via maximizing the partial area under the ROC curve of linear risk scores. *Biostatistics*. 12(2): 369-385.
- Pisano et al. (2005) Diagnostic Performance of Digital versus Film Mammography for Breast-Cancer Screening. *The New England Journal of Medicine*. 353(17):773-1783.
- Gigerenzer, et al. (2008) Helping Doctors and Patients Make Sense of Health Statistics. *Psychological Science in the Public Interest*. 8(2):53-90.

Acknowledgement: The authors would like to thank Maoji Li for her technical assistance in making the graphs, checking the R code for the results presented.