

# Probability estimation with sparse structure

Chanwoo Lee, October 1, 2020

## 1 Low rank approximation

To calculate the conditional probability  $\mathbb{P}(y = 1|\mathbf{X})$  where  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$  is not necessarily a point in training dataset, we need to calculate  $\mathbf{B}$  from the output  $\mathbf{P}_{\text{row}}$ ,  $\mathbf{P}_{\text{col}}$ , and  $\alpha$ . I calculate  $\mathbf{B}$  as

$$\mathbf{B} = \sum_{i=1}^n \alpha_i y_i (\mathbf{P}_{\text{row}} \mathbf{P}_{\text{row}}^T \mathbf{X}_i + \mathbf{X}_i \mathbf{P}_{\text{col}} \mathbf{P}_{\text{col}}^T). \quad (1)$$

I have checked  $\mathbf{B}$  calculated from the formula (1) satisfies

$$\text{fitted}_i = \langle \mathbf{B}, \mathbf{X}_i \rangle + \text{intercept},$$

where  $\text{fitted}_i$  and  $\text{intercept}$  are output of `SMMK.con`. This implies that the estimated predictor is exactly  $\langle \mathbf{B}, \mathbf{X} \rangle + \text{intercept}$ . Notice that this  $\mathbf{B}$  has rank at most  $2r$ . The following figure shows the probability estimation based on the optimal coefficient  $\mathbf{B}$ .

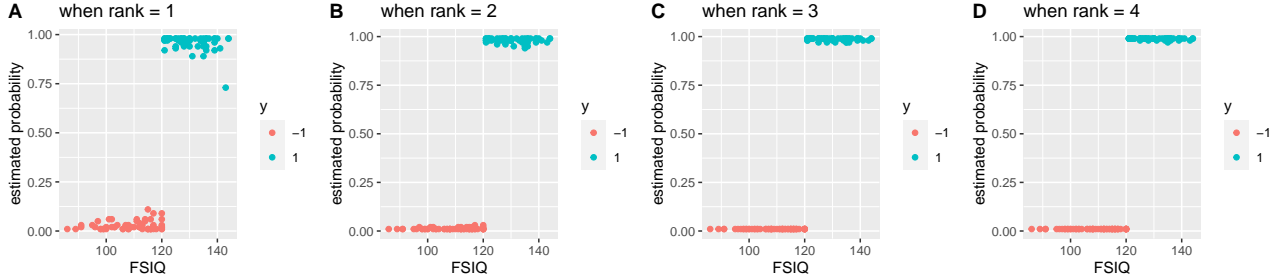


Figure 1: Probability estimation results according to different ranks based on the coefficient  $\mathbf{B}$ .

Since  $\mathbf{X}_i$  are symmetric, we make  $\mathbf{B}$  symmetric as  $\mathbf{B}' = (\mathbf{B} + \mathbf{B}^T)/2$ . I verified everything remains when we use  $\mathbf{B}'$  instead of  $\mathbf{B}$ . Since we have perfect separation even when  $\text{rank} = 1$ , we can check that the estimated probability has clear separation far from 0.5, from which we can not have information about true IQ values.

To make the coefficient matrix  $\mathbf{B}$  has rank  $r$ , we can use low rank approximation on  $\mathbf{B}$  from (1) as

$$\tilde{\mathbf{B}} = \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T, \quad (2)$$

where columns of  $\mathbf{U}_r$  and  $\mathbf{V}_r$  are the first  $r$  right and left singular vectors of  $\mathbf{B}$ , and  $\mathbf{\Sigma}_r$  is a diagonal matrix whose entries are top  $r$  singular values of  $\mathbf{B}$ . Notice that there is no guarantee that  $\tilde{\mathbf{B}}$  is the best  $r$  rank coefficient that minimizes the loss function. I have checked that the approximated  $\tilde{\mathbf{B}}$  cannot perfectly separate the training dataset when the rank is less than 4, which implies  $\tilde{\mathbf{B}}$  is not a global optimum among rank  $r$  matrices.

Figure 2 shows the probability estimation results based on  $\tilde{\mathbf{B}}$ . We can see that the estimated probability has linear relationship to the true IQ values though it looks quite noisy.

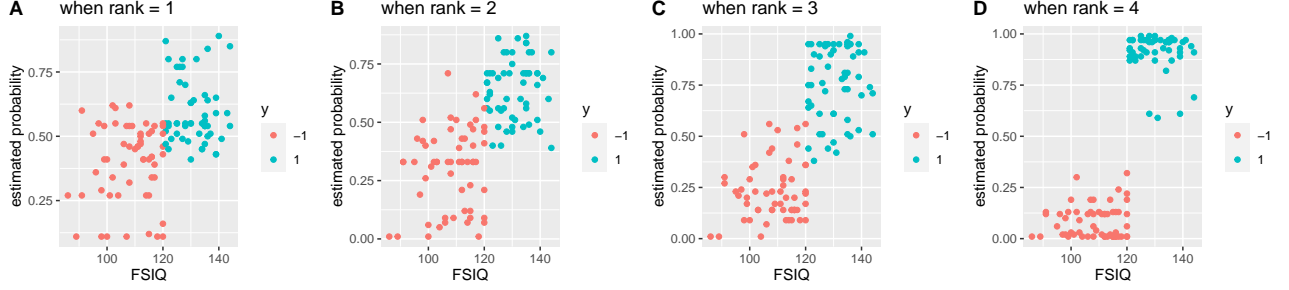


Figure 2: Probability estimation results according to different ranks based on the coefficient  $\tilde{\mathbf{B}}$ .

## 2 Sparse structure

I first, set the rank and sparsity as (1,66) and estimate the probability based on two updated algorithm ADMM and SMMK.con. Figure 3 compares 4 different ways to estimate the conditional probability  $\mathbb{P}(y = 1|\mathbf{X})$ . Figure A is the estimation based on  $\mathbf{B}$  in (1), Figure B is based on low rank approximation  $\tilde{\mathbf{B}}$  in (2), Figure C is based on the sparse structure with the algorithm ADMM, and Figure D is based on the sparse structure with the algorithm SMMK.con. We can check that methods based on the sparse structure estimate the probability values centered at 0.5. My guess for the reason is that the sparse structure does not perform well on classification. So I compare the error rates of classification on each method with the training dataset  $\{(\mathbf{X}_1, y_1)\}_{i=1}^{114}$ . The following table shows the error rate on the training set when rank = 1. Considering the response variable is

Method	SMMK	SMMK+ low rank	ADMM	SMMK+low rank+sparse
Error rate	0	0.289	0.403	0.350

Table 1: Error rates on the training dataset according to different methods when rank = 1.

binary, ADMM method has quite huge error rate.

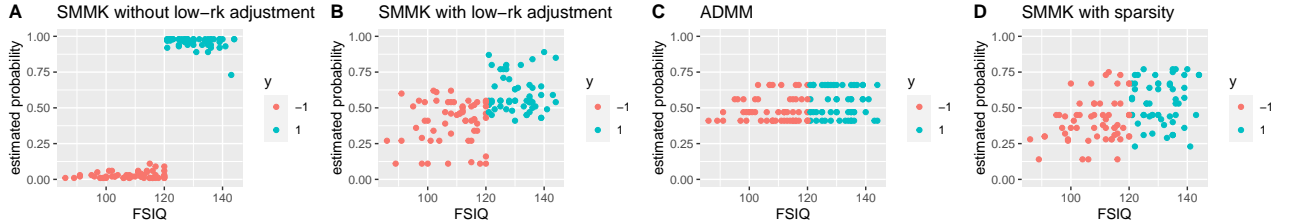


Figure 3: Comparison of the probability estimation when rank = 1.

Therefore, I find another combination of rank and sparsity that has moderate error rates. To this end, I set the rank and sparsity as (2,60) which has small cross validation error in the figure “hinge\_test\_corner.pdf”. The following table shows the comparison of the error rates of each method when the rank = 2.

Figure 4 shows the probability estimation results on each method. The performance of methods based on the sparse structure seems to improve compared to when the rank is one. If you make

Method	SMMK	SMMK+ low rank	ADMM	SMMK+low rank+sparse
Error rate	0	0.157	0.157	0.245

Table 2: Error rates on the training dataset according to different methods when rank = 2.

the classification rule as

$$y = \begin{cases} 1 & \text{if } \mathbb{P}(y = 1|\mathbf{X}) > 0.5, \\ -1 & \text{otherwise.} \end{cases}$$

We can check that there is less chance to assign  $y = 1$  when IQ value is small or  $y = -1$  when IQ value is high when we use SMMK with the low rank approximation.

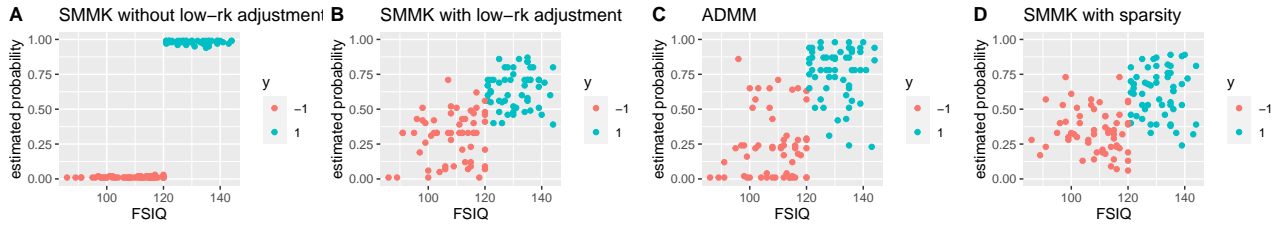


Figure 4: Comparison of the probability estimation when rank = 2.

### 3 Discussion on setting hyper parameters

We set hyper parameters ( $rank, cost, sparsity$ ) based on cross validation result when the weight  $p = 0.5$ , and use the combination of parameters to estimate the conditional probability  $\mathbb{P}(y = 1|\mathbf{X})$ . There are two questions:

1. why do we fix ( $rank, cost, sparsity$ ) to estimate the probability estimation.
2. why are we based on the best combinations ( $rank, cost, sparsity$ ) of equal weight  $\pi = 0.5$  case.

First, we pre-specify the considered function class

$$\mathcal{F} = \{f(\cdot) = \langle \mathbf{B}, \cdot \rangle + b : \text{rank}(\mathbf{B}) \leq r, \text{sparsity}(\mathbf{B}) \leq s\},$$

which means the rank and the sparsity are the global parameters that we have to stick to regardless of  $\pi$ . In addition, our consistency theorems are based on the fixed rank, sparsity, and the cost values. One example is that we have the condition  $\lambda^{-1} \geq 4\delta_n^{-2} J_{\pi}^*$ , in the theorem for the consistency probability estimation. We also used the same cost  $\lambda$  in the proof as well as the rank. I have checked that the reference papers [1, 2] keep using the same cost value across different weights.

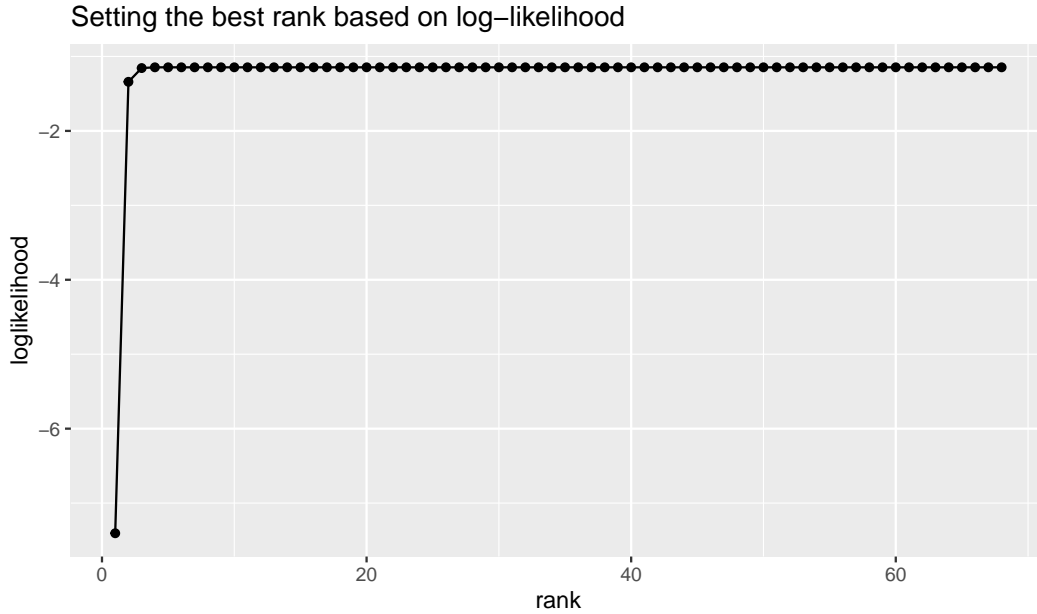
Second, the second question is a fair question. There is no specific reason set the rank based on the cross validation results at equal weight. One possible answer to this is since there is no information about which misclassification is more important (misclassify as  $y = 1$  when  $y = -1$  vs as misclassify  $y = -1$  when  $y = 1$ ) we assume the two types of misclassification are equally important setting the

case  $\pi = 0.5$  as criterion. To avoid this issue, we can use log-likelihood

$$L(cost, rank) = \sum_{i=1}^n y_i \log \left( \hat{\mathbb{P}}(y_i = 1 | \mathbf{X}_i, cost, rank) \right) + (1 - y_i) \log \left( 1 - \hat{\mathbb{P}}(y_i = 1 | \mathbf{X}_i, cost, rank) \right).$$

and choose the combination that maximizes the  $L(cost, rank)$  [2]. Another possible approach is to use estimated generalized Kullback-Leibler loss [1]. One issue for alternatives way is that it is time consuming procedure. For example, IQ data set needs at least  $68 \times 99$  computations to compare the log likelihood for each setting. If we add sparsity, we might have to calculate  $68 \times 99 \times sparsity$  combinations to pick the best one.

The following figure is the loglikelihood result



## References

- [1] Junhui Wang, Xiaotong Shen, and Yufeng Liu. Probability estimation for large-margin classifiers. *Biometrika*, 95(1):149–167, March 2008.
- [2] Yichao Wu, Hao Helen Zhang, and Yufeng Liu. Robust model-free multiclass probability estimation. *Journal of the American Statistical Association*, 105(489):424–436, 2010.

