

# SoCal: supervised genotype calling via ellipsoidal separation for Affymetrix SNP microarray

Huwenbo Shi (603-778-363) shihuwenbo@ucla.edu

## Abstract

**Background:** Genotype calling for SNP microarray is an important step in array-based association studies. Most supervised genotype calling methods for SNP microarray fit generative models (e.g. Gaussian models) on log-transformed allele intensities of SNPs from samples having reference genotype calls.

**Methods:** In this article, I present SoCal, a supervised genotype calling method that efficiently finds ellipsoidal decision regions for each genotype via ellipsoidal separation instead of fitting a generative model. By assigning different weights to the criteria of finding separating ellipsoids, SoCal can control the effect of outliers in training data.

**Results:** In leave-one-out cross-validation with HapMap reference calls, SoCal achieved a concordance rate of 98.94% at a call rate of 100% and 99.71% at a call rate of 95%, comparable to many state-of-the-art methods. After being trained using HapMap reference calls, SoCal achieved a concordance rate of 95.10% with genotype calls made by CRLMM at a call rate of 100% on a set of validation calls excluding the training calls. SoCal also shows more robustness than RLMM, a supervised method that uses Gaussian decision regions to call genotypes, when outliers are present in training data. Overall, SoCal is a novel and promising genotype calling method.

## 1 Introduction

Accurate genotyping of SNPs is essential to discovering true causal variants in association studies [1]. Although next generation sequencing technology provides cheap whole-genome sequences for genotyping SNPs, SNP microarray is still a cost-effective genotyping technology for many specific association studies [2]. In an Affymetrix SNP microarray, oligonucleotide probes are used to match and bind DNA fragments containing biallelic SNPs. Then a fluorescence scanner scans the microarray to quantify perfect match and mismatch for these fragments. Most genotype calling procedures for Affymetrix SNP microarray consist of two steps. In the first step, raw information from microarray is summarized to obtain the intensities,  $\theta_A$  and  $\theta_B$ , of the two alleles, denoted by A and B, of each SNP. In the second step, SNPs are classified into genotype AA, AB, or BB based on the allele intensities they generate. The focus of this article is on the second step of the genotype calling procedure—genotype calling using summarized allele intensities.

For a specific SNP, if a sample has genotype AA or BB, the allele intensity,  $\theta_A$  or  $\theta_B$ , will be higher respectively. If a sample has genotype AB, the intensities,  $\theta_A$  and  $\theta_B$ , will be similar. If one plots  $\log(\theta_A)$  versus  $\log(\theta_B)$  of a SNP for a number of samples, normally 3 ellipsoidal clusters can be observed, one for each genotype, as shown in Figure 1. Many genotype calling algorithms use model-based unsupervised clustering to identify these clusters and then assign genotypes to each cluster [3, 4, 5]. Although these methods are applicable to a wide range of microarrays because they only require information from microarrays, they don't take advantage of genotype calls that are already available. Also, these methods use EM algorithm to estimate model parameters, which is sensitive to starting parameters and slow to converge [6]. To utilize reference genotype calls, Rabbee and Speed proposed the RLMM algorithm, a supervised genotype calling method that forms decision regions for each genotype by fitting bivariate Gaussian distributions on log-transformed allele intensities with reference genotype calls [7]. These Gaussian decision regions are then used to call SNPs for samples with unknown genotypes. However, fitting a Gaussian distribution is known to be non-robust to outliers [8]. And for SNP microarrays, outliers can be caused by genomic structural variations [9].

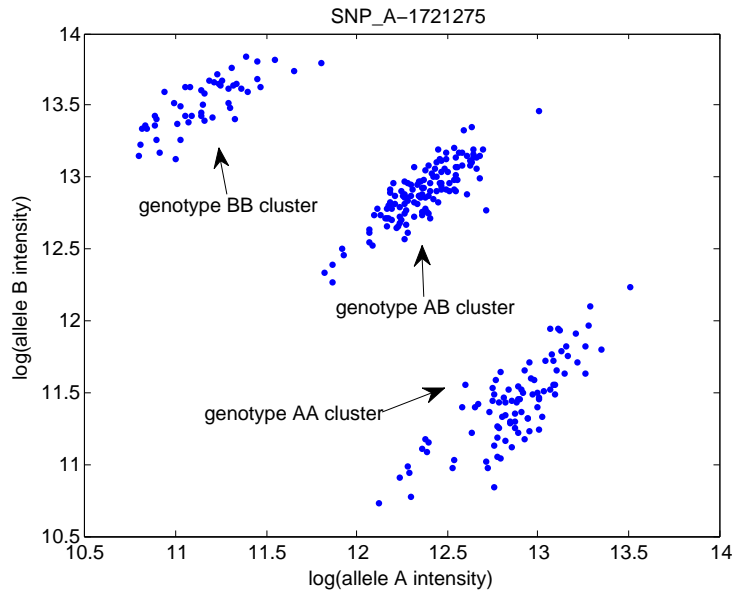


Figure 1: Genotype clusters obtained by plotting log-transformed allele intensities of the two alleles of a SNP for a number of samples. Each point in the plot represents a sample.

In this article, I present SoCal, a supervised genotype calling method for Affymetrix SNP microarray. Instead of fitting Gaussian distributions on log-transformed allele intensities with reference genotype calls, SoCal efficiently finds ellipsoidal decision regions for each genotype via ellipsoidal separation by solving a conic programming problem. SoCal can control the effect of outliers by assigning different weights to the criteria of finding separating ellipsoids—separation ratio, ellipsoid volume, and inclusion of points. After SoCal finds the ellipsoidal decision regions for each genotype, it uses them to call SNPs for samples with unknown genotypes using minimum distance classification.

Using reference genotype calls from the HapMap Project as training and validation

data, SoCal achieved a concordance rate of 98.94% at a call rate of 100% and 99.71% at a call rate of 95% in leave-one-out cross-validation. Furthermore, SoCal showed more robustness than the RLMM method when outliers were present in training data.

## 2 Methods

### 2.1 Overview of SoCal's genotype calling procedure

SNP allele intensities are first normalized and summarized from raw microarray data using the RMA method, an important preprocessing step that reduces cross-chip and cross-lab non-biological effects from raw data [10, 11]. After the preprocessing step, SoCal calls genotypes in two steps. In the first step, SoCal finds ellipsoidal decision regions for each genotype of a SNP using reference genotype calls. In the second step, SoCal uses these ellipsoidal decision regions to call SNPs for samples with unknown genotypes through minimum distance classification.

In this section, I first introduce the problem of pattern separation by ellipsoid. Then I describe how SoCal finds ellipsoidal decision regions for each genotype and then calls genotypes using these ellipsoidal decision regions.

### 2.2 Pattern separation by ellipsoid

An ellipsoid  $\mathcal{E} \subseteq \mathbb{R}^n$  can be expressed as  $\mathcal{E} = \{x \in \mathbb{R}^n | (x - c)^T E (x - c) \leq 1\}$ , where  $c$  is the center of the ellipsoid, and  $E$  a positive definite matrix denoting the shape and orientation of the ellipsoid. Let  $\{a_i\}$  be the points to be included in an ellipsoid, and  $\{b_j\}$  be the points to be excluded, the problem of ellipsoidal separation is to find  $c$  and  $E$  such that  $(a_i - c)^T E (a_i - c) \leq 1 \forall i$  and  $(b_j - c)^T E (b_j - c) > 1 \forall j$ .

### 2.3 Forming ellipsoidal decision regions for each genotype

Let  $G = \{AA, AB, BB\}$  be the set of genotypes of SNP  $n$ , and  $J_{AA}, J_{AB}, J_{BB}$  the index set of samples with the corresponding genotype. Let  $X = \{(\log(\theta_A), \log(\theta_B))_i | i = 1, \dots, |J_{AA}| + |J_{AB}| + |J_{BB}|\}$  be the set of log-transformed allele intensities of all the samples, and  $X_{AA} = \{x_j | x_j \in X, j \in J_{AA}\}$ ,  $X_{AB} = \{x_j | x_j \in X, j \in J_{AB}\}$ ,  $X_{BB} = \{x_j | x_j \in X, j \in J_{BB}\}$  the set of log-transformed allele intensities from samples having the corresponding genotype for SNP  $n$ .

To find the ellipsoid that includes  $X_{AA}$  and excludes  $X_{AB} \cup X_{BB}$ , one sets  $\{a_i\} = X_{AA}$  and  $\{b_j\} = X_{AB} \cup X_{BB}$ , and solves the following conic programming problem.

$$\begin{aligned} & \text{minimize} && -\beta_1 k + \beta_2 \text{trace}(T) + \beta_3 \|u - \mathbb{1}\|_1 \\ & \text{subject to} && (1, a_i)^T \tilde{E} (1, a_i) \leq u_i \quad \forall i \\ & && (1, b_j)^T \tilde{E} (1, b_j) \geq k \quad \forall j \\ & && \tilde{E} = \begin{bmatrix} s & v^T \\ v & F \end{bmatrix} \geq 0 \\ & && \begin{bmatrix} F & I \\ I & T \end{bmatrix} \geq 0 \end{aligned}$$

Here,  $I$  denotes the identity matrix. Derivation of the problem formulation is largely followed from [12]. For the sake of space, detailed derivation of the problem formulation is not presented here. In the problem formulation above,  $\beta_i > 0$  are the weights assigned

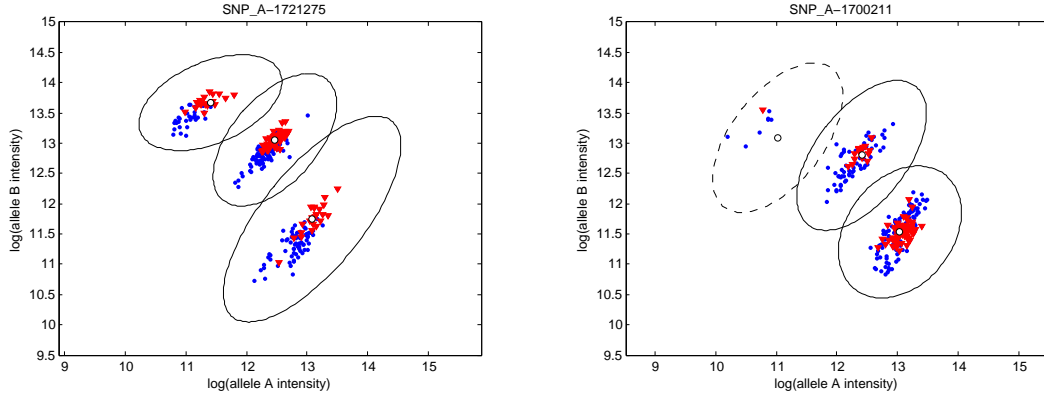
to the criteria of finding separating ellipsoid—separation ratio, ellipsoid volume, and inclusion of points. By increasing  $\beta_2$ , one finds ellipsoid with smaller volume. And by increasing  $\beta_3$ , one finds ellipsoid that tries to include more data points. In SoCal, default values for  $\beta_1$ ,  $\beta_2$ ,  $\beta_3$  are empirically set to 1,  $10^4$ , and  $10^2$  respectively.

Let  $\tilde{E}^* = \begin{bmatrix} s & v^T \\ v & F \end{bmatrix}$  be the optimal solution to the problem above. The separating ellipsoid  $\mathcal{E}^*$  is defined as  $\mathcal{E}^* = \{x \in \mathbb{R}^n | (x - c^*)^T E^* (x - c^*) \leq 2(1 + k)\}$ , where  $c^* = -F^{-1}v$ ,  $E^* = \frac{F}{(1 - s + c^{*T} F c^*)}$ .

To find the ellipsoid that includes  $X_{AB}$  and excludes  $X_{AA} \cup X_{BB}$ , one sets  $\{a_i\} = X_{AB}$  and  $\{b_j\} = X_{AA} \cup X_{BB}$ , and solves the above conic programming problem. The same procedure also applies to finding the ellipsoid that includes  $X_{BB}$  and excludes  $X_{AA} \cup X_{AB}$ .

## 2.4 Handling sparse or missing genotype clusters

If a SNP has moderate minor allele frequency (MAF), the genotype clusters of that SNP are well defined, and SoCal obtains three ellipsoidal decision regions for that SNP, one for each genotype cluster (Figure 2a). However, if a SNP has lower MAF, some genotype cluster may be sparse or missing. For these SNPs, SoCal estimates the ellipsoid for the sparse or missing genotype cluster using the ellipsoids for the other two genotypes through simple geometric transformations (Figure 2b).



(a) SNP with 3 well defined genotype clusters. SoCal obtains one ellipsoid for each genotype cluster.

(b) SNP with sparse genotype BB cluster. SoCal first obtains ellipsoids for genotype AA and AB clusters, and then estimates the ellipsoid (drawn in dashed line) for genotype BB cluster.

Figure 2: Ellipsoids obtained by SoCal for SNPs with well-defined and sparse genotype clusters. Each dot in the plots represents a sample, with samples having HapMap reference genotype calls marked as red triangles. The ellipsoids are obtained using all the reference calls.

### 2.4.1 Missing genotype AA or BB cluster

If genotype AA cluster of a SNP has less than 3 reference calls, SoCal first finds the ellipsoids for genotype AB and BB clusters, and then estimates that for genotype AA cluster through simple geometric transformations.

Let  $\mathcal{E}_{AB} = \{x \in \mathbb{R}^n | (x - c_{AB})^T E_{AB} (x - c_{AB}) \leq 1\}$  and  $\mathcal{E}_{BB} = \{x \in \mathbb{R}^n | (x - c_{BB})^T E_{BB} (x - c_{BB}) \leq 1\}$  be the ellipsoids obtained for genotype AB and BB clusters, and  $n_{AB}$ ,  $n_{BB}$  the unit vectors pointing in the direction of the major axis of the corresponding ellipsoid. SoCal estimates the center of  $\mathcal{E}_{AA}$ , the ellipsoid for genotype AA cluster, by reflecting  $c_{BB}$ , the center of  $\mathcal{E}_{BB}$ , across the major axis of  $\mathcal{E}_{AB}$ . To estimate the orientation of  $\mathcal{E}_{AA}$ , SoCal first determines the angle between  $n_{AB}$  and  $n_{BB}$ , and then applies a rotation matrix of that angle on  $E_{AB}$ .

Formally, let  $\mathcal{E}_{AA} = \{x \in \mathbb{R}^n | (x - c_{AA})^T E_{AA} (x - c_{AA}) \leq 1\}$  be the estimated ellipsoid for genotype AA cluster, and  $\alpha$  the angle between  $n_{AB}$  and  $n_{BB}$ , then  $c_{AA} = -c_{BB} + 2c_{AB} + 2n_{AB}((c_{BB} - c_{AB})^T n_{AB})$ , and  $E_{AA} = R^T E_{AB} R$ , where  $R$  is a rotation matrix of angle  $\alpha$ .

If genotype BB cluster is missing, the center and orientation of the ellipsoid for that cluster is estimated in a similar way. Formally, let  $\mathcal{E}_{BB} = \{x \in \mathbb{R}^n | (x - c_{BB})^T E_{BB} (x - c_{BB}) \leq 1\}$  be the estimated ellipsoid for genotype BB cluster, and  $\alpha$  the angle between  $n_{AB}$  and  $n_{AA}$ , then  $c_{BB} = -c_{AA} + 2c_{AB} + 2n_{AB}((c_{AA} - c_{AB})^T n_{AB})$ , and  $E_{BB} = R^T E_{AB} R$ , where  $R$  is a rotation matrix of angle  $-\alpha$ .

#### 2.4.2 Missing genotype AB cluster

Although SNPs with genotype AB cluster missing were not observed in HapMap reference genotype calls, for completeness, for these SNPs SoCal first obtains,  $\mathcal{E}_{AA}$  and  $\mathcal{E}_{BB}$ , the ellipsoids for genotype AA and BB cluster, and then estimates the center of  $\mathcal{E}_{AB}$ , the ellipsoid for the missing cluster, using the mid-point between the centers of  $\mathcal{E}_{AA}$  and  $\mathcal{E}_{BB}$ . The orientation of  $\mathcal{E}_{AB}$  is obtained by applying a rotation to the ellipsoid with the minimum volume among  $\mathcal{E}_{AA}$  and  $\mathcal{E}_{BB}$ .

Formally, let  $\mathcal{E}_{AB} = \{x \in \mathbb{R}^n | (x - c_{AB})^T E_{AB} (x - c_{AB}) \leq 1\}$  be the estimated ellipsoid for genotype AB cluster, and  $\alpha$  the angle between  $n_{AA}$  and  $n_{BB}$ , then  $c_{AB} = (c_{AA} + c_{BB})/2$ , and  $E_{AB} = R^T \hat{E} R$ , where  $\hat{E}$  is the matrix of the ellipsoid with the minimum volume among  $\mathcal{E}_{AA}$  and  $\mathcal{E}_{BB}$ , and  $R$  a rotation matrix of angle  $\pm\alpha/2$ . The sign of the angle of rotation is dependent on the choice of ellipsoid on which rotation is applied—positive for  $\mathcal{E}_{AA}$  and negative for  $\mathcal{E}_{BB}$ .

### 2.5 Genotype calling

After the ellipsoidal decision regions,  $\mathcal{E}_g = \{x \in \mathbb{R}^n | (x - c_g)^T E_g (x - c_g) \leq 1\}$ ,  $\forall g \in \{AA, AB, BB\}$  of a SNP are obtained, SoCal uses them to classify SNPs for samples with unknown genotypes using minimum distance classification.

If a sample has allele intensities  $\theta_A$  and  $\theta_B$  at SNP  $n$ , SoCal first computes  $D_g = \sqrt{(x - c_g)^T E_g (x - c_g)}$ , where  $x = (\log(\theta_A), \log(\theta_B))$ , for each  $g \in \{AA, AB, BB\}$ . SoCal then calls the genotype,  $\mathcal{G}$ , of that sample at SNP  $n$  as the genotype having minimum  $D_g$ , that is,  $\mathcal{G} = \arg \min_{g \in \{AA, AB, BB\}} D_g$ .

SoCal defines  $\lambda = 1 - D_{\mathcal{G}}/(D_{AA} + D_{AB} + D_{BB})$  to quantify the confidence of each genotype call. By increasing the threshold for  $\lambda$ , SoCal can achieve higher call accuracy at the cost of decreasing call rate.

## 3 Materials

The microarray used for evaluation in this project was the Affymetrix GeneChip Human Mapping 50K Xba Array, which contains 58,960 SNPs. Raw microarray data

for 270 samples was obtained from HapMap FTP, and reference genotype calls were obtained from HapMap using HapMart [13].

After removing strand-ambiguous SNPs and SNPs not present on HapMart from the set of SNPs on the microarray, 16,387 SNPs were left. Figure 3 shows the minor allele frequency distribution for these SNPs.

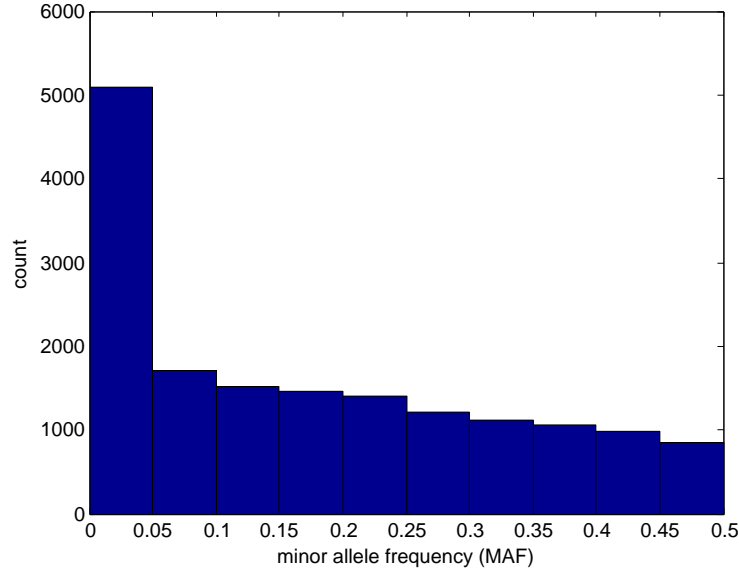


Figure 3: Minor allele frequency distribution for the 16,387 SNPs.

From the 16,387 SNPs, 4,064 SNPs with two genotype clusters having less than 3 reference genotype calls were further removed. Among these SNPs, 3,596 are monomorphic SNPs with MAF equal to 0. In total, 12,323 SNPs were left for evaluation. On average, each of these SNPs has 83 reference genotype calls.

## 4 Results

### 4.1 Cross-validation with HapMap reference calls

To evaluate the accuracy of SoCal, I compared the genotype calls made by SoCal with the reference calls from HapMap through leave-one-out cross-validation. For each SNP, I used one sample from the reference set as validation data and the rest as training data. I repeated this process until all the samples in the reference set were used as validation data exactly once. Concordance rate is defined as the ratio between the number of calls that are concordant with HapMap reference calls and the total number of calls made by SoCal.

First, I compared the accuracy of SoCal under different choices of  $\beta_i$ , the weights assigned to the criteria of finding ellipsoidal decision regions for each genotype cluster. Figure 4 shows the concordance rate of SoCal in the leave-one-out cross-validation at a wide range of call rates for different values of  $\beta_i$ . Because the weights,  $\beta_1 = 1$ ,  $\beta_2 = 10^4$ ,  $\beta_3 = 10^2$ , had the highest call rates at fixed concordance rates, they are set to be the default of SoCal. And all other experiments presented in this article used this choice of weights.

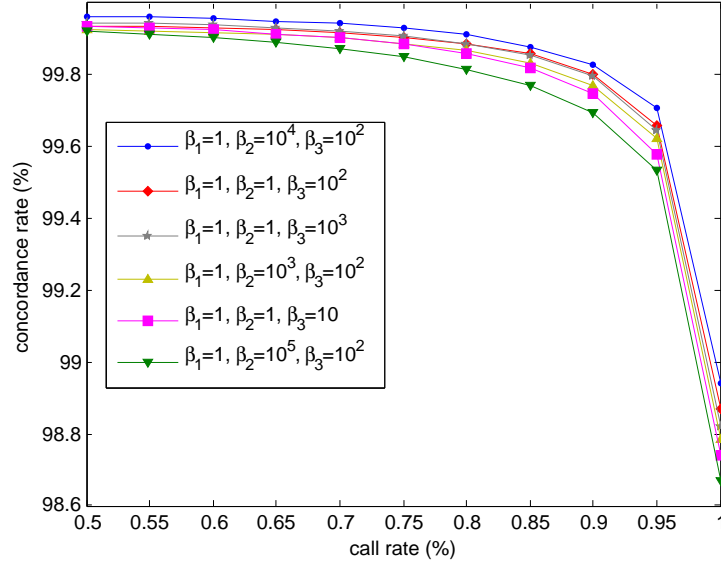


Figure 4: Concordance rate of SoCal in the leave-one-out cross-validation with HapMap reference calls as a function of call rate, for different choices of  $\beta_i$

Table 1 shows the genotype calls made by SoCal and the reference calls from HapMap in leave-one-out cross-validation. At a call rate of 100%, SoCal made 1,081,319 calls in total, out of which 1,069,857 were concordant with HapMap calls, achieving a concordance rate of 98.94%.

HapMap/SoCal	AA	AB	BB	No Call
AA	360,289	2,282	1,058	0
AB	2,667	341,012	2,257	0
BB	851	2,347	368,556	0

Table 1: At a call rate of 100%, SoCal achieved 98.94% concordance rate in leave-one-out cross-validation with HapMap reference calls.

Table 2 shows detailed comparison between SoCal and HapMap calls at a call rate of 95%. At a call rate of 95%, SoCal made 1,028,258 calls in total, out of which 1,025,242 were concordant with HapMap calls, achieving a concordance rate of 99.71%. These results are comparable to those achieved by previous methods [7, 15].

HapMap/SoCal	AA	AB	BB	No Call
AA	348,221	390	298	14,720
AB	710	319,394	775	25,057
BB	410	427	357,627	13,290

Table 2: At a call rate of 95%, SoCal achieved 99.71% concordance rate in leave-one-out cross-validation with HapMap reference calls.

## 4.2 Comparison with CRLMM calls

As another way of evaluating the accuracy of SoCal, I compared the genotype calls made by SoCal and those made by CRLMM, a state-of-the-art supervised genotype calling method for SNP microarrays that uses a two-level hierarchical model to model variations in allele intensities across SNPs and across chips [11].

I first trained SoCal using all the samples with HapMap reference calls, and then made genotype calls on the rest of the samples. When comparing SoCal with CRLMM, I excluded the training samples and compared these two methods only at samples not in the training set.

Table 3 shows detailed comparison between SoCal and CRLMM at a call rate of 100%. In total, SoCal made 2,245,891 calls, out of which 2,134,868 were concordant with those made by CRLMM, achieving a concordance rate of 95.10%. The high concordance rate between SoCal and CRLMM suggests that SoCal has the potential to become an alternative genotype caller.

CRLMM/SoCal	AA	AB	BB	No Call
AA	781,903	23,244	10,405	0
AB	22,340	564,280	22,533	0
BB	7,730	24,771	788,685	0

Table 3: At a call rate of 100%, SoCal achieved 95.10% concordance rate with the calls made by CRLMM. Training samples for SoCal were excluded during comparison.

## 4.3 Comparison with RLMM in the presence of outliers

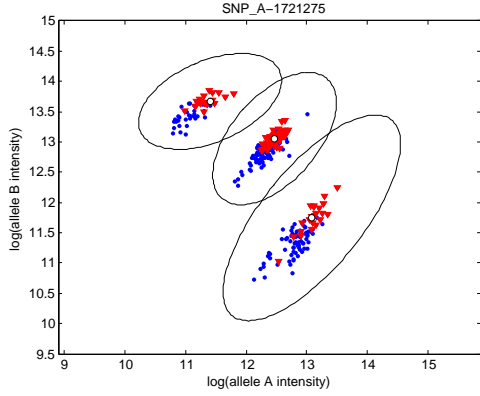
I investigated how robust SoCal is when outliers are present in training data. For comparison, I implemented the RLMM algorithm, which fits bivariate Gaussian distributions on log-transformed allele intensities of each genotype cluster and then classifies SNPs with unknown genotype into the distribution having minimum Mahalanobis distance based on the allele intensities they generate [7].

For accurate comparison, I selected a subset of 3,442 SNPs that have more than 10 reference calls for each genotype cluster from the set of 12,323 SNPs used for evaluation in previous experiments. To simulate outliers, for each SNP, I first estimated  $\mu_g$ , the mean of log-transformed allele intensities of each genotype cluster, and then drew one outlier for each genotype cluster from the Gaussian distribution  $N(\mu_g, \gamma I)$ , where  $I$  is the identity matrix and  $\gamma$  a positive constant controlling the variance of the distribution—by increasing  $\gamma$ , one increases the effect of outliers. In total, I simulated 3 outliers for each SNP, one for each genotype cluster.

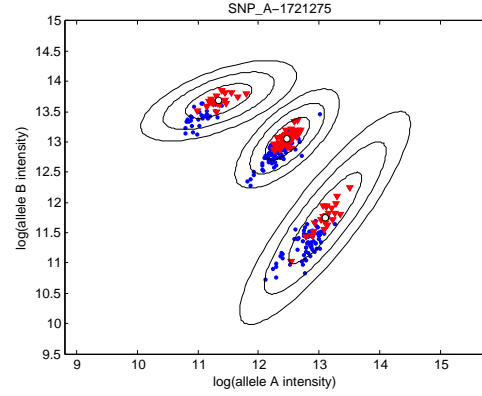
To illustrate the robustness of SoCal and RLMM, Figure 5 shows the ellipsoidal decision regions obtained by SoCal and the level curves of the Gaussian decision regions obtained by RLMM before and after an outlier is introduced to the genotype AA cluster. Before an outlier is introduced, both SoCal and RLMM can find appropriate decision regions for each genotype cluster and make accurate genotype calls. However, after an outlier is introduced into the genotype AA cluster, the estimated variance of the Gaussian decision region obtained by RLMM for the genotype AA cluster is significantly affected, making genotype calling much less accurate. On the other hand, because SoCal not only considers outliers but also jointly uses data points from other genotype clusters when forming the decision regions, the decision region for the genotype AA cluster,



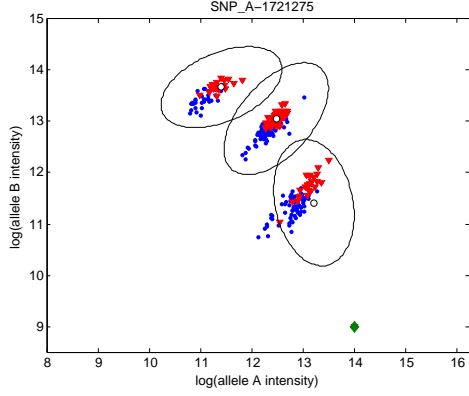
although affected, is still accurate enough to classify all samples into correct genotypes.



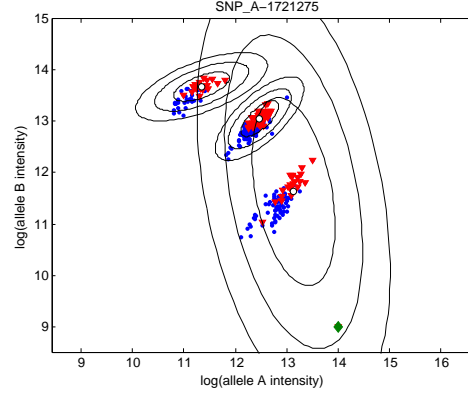
(a) Decision regions formed by SoCal when there is no outlier.



(b) Decision regions formed by RLMM when there is no outlier.



(c) Decision regions formed by SoCal when an outlier (diamond-shaped point) is introduced into the genotype AA cluster. Although affected, these decision regions can still classify all samples into correct genotypes.



(d) Decision regions formed by RLMM when an outlier (diamond-shaped point) is introduced into the genotype AA cluster. The estimated variance of the bivariate Gaussian distribution is significantly affected, and RLMM may mistakenly classify samples with genotype AB into genotype AA.

Figure 5: Decision regions formed by SoCal and RLMM before and after an outlier is introduced to the genotype AA cluster. Samples with reference genotype calls are marked as red triangles.

Figure 6 shows the decrease in concordance rate of SoCal and RLMM at call rate of 100% in leave-one-out cross-validation with HapMap reference genotype calls as the variance of simulated outliers varies from 1 to 10. Clearly, the concordance rate of SoCal decreases much more slowly than does the RLMM method. Thus, SoCal is in general more robust to outliers than RLMM.

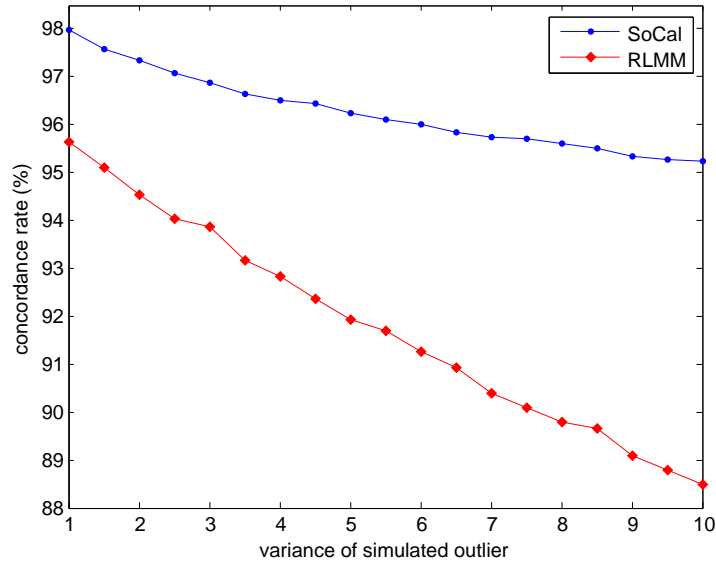


Figure 6: Concordance rate of SoCal and RLMM in the leave-one-out cross-validation with HapMap reference calls as a function of outlier variance.

#### 4.4 Software implementation

SoCal is implemented in Python. To solve the conic programming problem of finding separating ellipsoids, SoCal uses CVXOPT [14]. Source code of SoCal is available at [https://github.com/huwenboshi/wqe/tree/master/genotype\\_caller](https://github.com/huwenboshi/wqe/tree/master/genotype_caller).

## 5 Discussion

I have presented SoCal, a supervised genotype calling algorithm for Affymetrix SNP microarray. Unlike most existing supervised genotype calling algorithms that try to fit generative models (e.g. Gaussian models) on log-transformed allele intensities of SNPs from samples with reference genotype calls, SoCal uses these data to efficiently find ellipsoidal decision regions for each genotype cluster via ellipsoidal separation by solving a conic programming problem. Both cross-validation with HapMap reference calls and comparison with genotype calls made by CRLMM show that SoCal is comparable in accuracy to many of the state-of-the-art genotype calling methods. Also, when outliers are present in training data, SoCal outperforms RLMM, a genotype calling method that uses Gaussian decision regions to call genotypes, demonstrating the robustness of SoCal over existing methods. Overall, SoCal is a novel and promising genotype caller for Affymetrix SNP microarray.

Like many supervised genotype calling methods, SoCal has its limitations. First, SoCal is not directly applicable to SNPs that don't have reference genotype calls. In this case, one can first call genotypes from microarrays using unsupervised genotype calling methods. Then one can treat these calls as training data and use SoCal to form refined ellipsoidal decision regions for each genotype cluster. Because SoCal is robust to outliers in training data, the refined ellipsoidal decision regions can be directly used to accurately call genotypes for future samples. Another limitation of SoCal is that users need to tune the weights assigned to the criteria of finding the separating ellipsoids for

different microarrays. However, experiments with SoCal using different weights show that SoCal is relatively less sensitive to weight parameters if  $\beta_3$ , the weight assigned to the criterion of inclusion of points, is set to  $10^2$ .

SoCal is still in development, and can be improved and extended in many directions. First, the current approach that SoCal uses to handle SNPs with sparse or missing genotype clusters is through simple and fixed geometric transformations. This approach assumes that genotype AA cluster and genotype BB cluster are symmetric around genotype AB cluster. However, this is not true in general. An improvement to this approach is to estimate the positions and orientations of the ellipsoids for sparse or missing clusters using information from SNPs with well-defined clusters that generate similar allele intensities patterns. Second, SoCal currently only uses allele intensities data from samples having reference genotype calls. However, allele intensities data for samples having structural variations is also available from the HapMap Project [13]. A possible improvement for SoCal is to include these data in finding the ellipsoidal decision regions for each genotype cluster to further refine the decision regions of each genotype. These refined decision regions can then be used to call genotypes more accurately and to detect outliers and possible structural variations.

To summarize, SoCal presents a novel and promising method for genotype calling. It's efficient in that it finds decision regions for each genotype via ellipsoidal separation by solving a conic programming problem, which is solvable in polynomial time with guaranteed global optimum [12]. Also, SoCal is comparable in accuracy to many state-of-the-art methods. Although SoCal has the limitation that training data must be available, this limitation is also present in other supervised genotype calling methods, and has been addressed previously [11]. Finally, SoCal can also be extended and improved to be more accurate and to have more functionality.

## References

- [1] Gordon D, Finch SJ. Factors affecting statistical power in the detection of genetic association. *J Clin Invest.* 2005;115(6):1408-18.
- [2] Roh SW, Abell GC, Kim KH, Nam YD, Bae JW. Comparing microarrays and next-generation sequencing technologies for microbial ecology research. *Trends Biotechnol.* 2010;28(6):291-9.
- [3] Norlén, H., Pettersson, E., Ahmadian, A., Lundberg, J., & Sundberg, R. (2008). Classification of SNP genotypes by a Gaussian mixture model in competitive enzymatic assays. *Mathematical Statistics Stockholm University Research Report*, 3, 1-26.
- [4] Lin Y, Tseng GC, Cheong SY, Bean LJ, Sherman SL, Feingold E. Smarter clustering methods for SNP genotype calling. *Bioinformatics.* 2008;24(23):2665-71.
- [5] Fujisawa H, Eguchi S, Ushijima M, et al. Genotyping of single nucleotide polymorphism using model-based clustering. *Bioinformatics.* 2004;20(5):718-26.
- [6] Wu, C. F. (1983). On the Convergence Properties of the EM Algorithm. *The Annals of Statistics*, 11, 95-103.
- [7] Rabbee N, Speed TP. A genotype calling algorithm for affymetrix SNP arrays. *Bioinformatics.* 2006;22(1):7-12.
- [8] Huber, P.J., 1981. *Robust Statistics*. Wiley, New York.
- [9] Marioni JC, Thorne NP, Valsesia A, et al. Breaking the waves: improved detection of copy number variation from microarray-based comparative genomic hybridization.

- Genome Biol. 2007;8(10):R228.
- [10] Bolstad BM, Irizarry RA, Astrand M, Speed TP. A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics*. 2003;19(2):185-93.
- [11] Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007;8(2):485-99.
- [12] Glineur F. (1998). Pattern separation via ellipsoids and conic programming. (MS Thesis). Facult Polytechnique de Mons, Mons, Belgium.
- [13] The International HapMap Consortium. The International HapMap Project. *Nature* 426, 789-796 (2003).
- [14] M. S. Andersen, J. Dahl, and L. Vandenberghe. CVXOPT: A Python package for convex optimization, version 1.1.7. Available at [cvxopt.org](http://cvxopt.org), 2014.
- [15] Di X, Matsuzaki H, Webster TA, et al. Dynamic model based algorithms for screening and genotyping over 100 K SNPs on oligonucleotide microarrays. *Bioinformatics*. 2005;21(9):1958-63.