

# Supplementary Information for Kinpute

Mark Abney and Aisha El Sherbiny

Department of Human Genetics, University of Chicago

# Contents

1	Additional quality metrics	4
---	----------------------------	---

## List of Tables

1	Genotype error model, $\Pr(O_i G_i)$	6
2	$\Pr(G_u S_{iu} = k)$	6
3	$\Pr(G_i G_u, S_{iu} = k)$	6
4	Informative 3-way configurations. Note that because individuals 1 and 2 are unordered, we list only one of the corresponding ordered 3-way configurations (e.g. configuration (8, 9, 8) is also informative but has the same probabilities as configuration (9, 8, 8) with individuals 1 and 2 switched). Note that for configuration (8, 8, 7), (7, 7, 7) and (3, 3, 7) $\Pr(G_1, G_2 G_u, S_{u12}) = \Pr(G_1 G_u, S_{u1}) = \Pr(G_2 G_u, S_{u2})$ , though we do not seek these configurations in our algorithm.	7
5	Configuration ( $S_{u1} = 9, S_{u2} = 8, S_{12} = 8$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	7
6	Configuration ( $S_{u1} = 8, S_{u2} = 8, S_{12} = 9$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	7
7	Configuration ( $S_{u1} = 8, S_{u2} = 8, S_{12} = 8$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	8
8	Configuration ( $S_{u1} = 8, S_{u2} = 6, S_{12} = 5$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	8
9	Configuration ( $S_{u1} = 8, S_{u2} = 5, S_{12} = 6$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	8
10	Configuration ( $S_{u1} = 8, S_{u2} = 5, S_{12} = 5$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	9
11	Configuration ( $S_{u1} = 5, S_{u2} = 5, S_{12} = 2$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	9
12	Configuration ( $S_{u1} = 4, S_{u2} = 3, S_{12} = 8$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	9
13	Configuration ( $S_{u1} = 3, S_{u2} = 2, S_{12} = 5$ ), probabilities $\Pr(G_1, G_2 G_u, S_{u12})$	10
14	Heterozygote sensitivity	10
15	Heterozygote positive predictive value	11
16	Concordance rates	12

## List of Figures

- 1 The nine condensed identity states  $S_{AB}$ . A line joining two alleles indicates those alleles are IBD. . . . . 5

# 1 Additional quality metrics

In addition to  $R^2$ , we computed other measures on how much Kinpute improves the imputation quality of LD-based methods. In all measures we compare the imputed genotypes (both from Impute2 and from Kinpute) to the true genotypes on the 48 held out individuals. Genotypes were compared only for chromosome 22. The additional measures are, concordance, heterozygote sensitivity (i.e. recall), and heterozygote positive predictive value (i.e. precision). Each of these measures require hard genotype calls, which were done at a genotype probability threshold of 0.9, (i.e. genotypes classified as having “high” certainty) or by picking the genotype with the maximum probability (i.e. for genotypes classified as having “low” certainty). The concordance is defined as the frequency that the called genotype is the same as the true genotype. Heterozygote sensitivity is the fraction of true heterozygotes that are called as heterozygotes. The heterozygote positive predictive value is the fraction of called heterozygotes that are true heterozygotes. Results are shown in Tables 14 – 16.

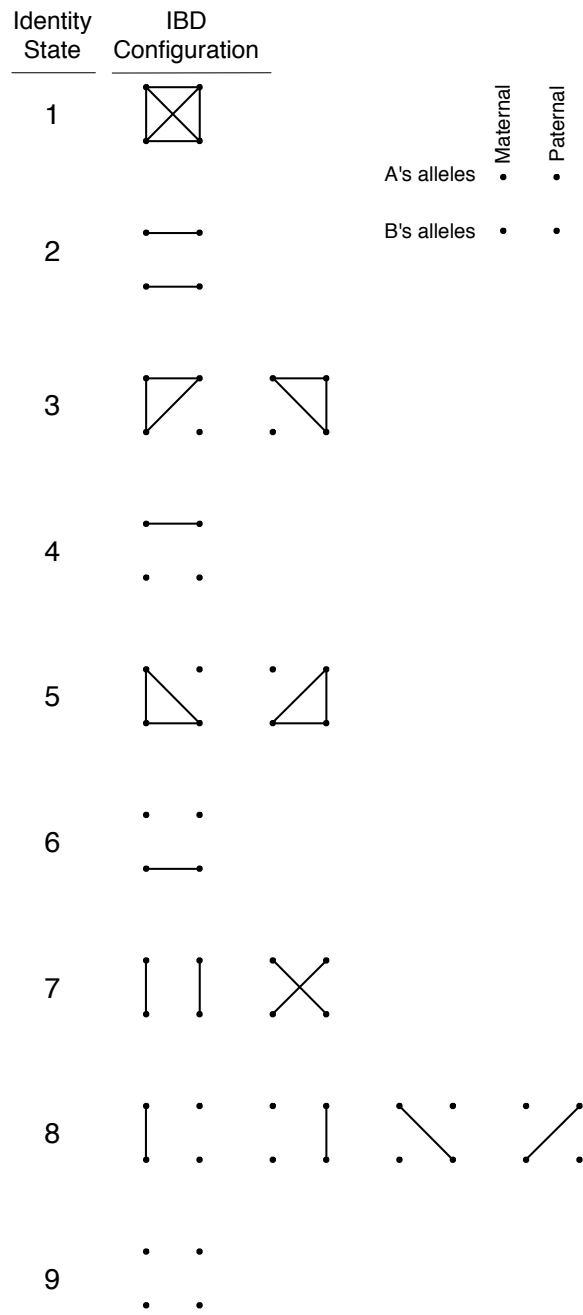


Figure 1: The nine condensed identity states  $S_{AB}$ . A line joining two alleles indicates those alleles are IBD.

$G_i$	$O_i$		
	0	1	2
0	$(1 - \epsilon)^2$	$2\epsilon(1 - \epsilon)$	$\epsilon^2$
1	$\epsilon(1 - \epsilon)$	$1 - 2\epsilon(1 - \epsilon)$	$\epsilon(1 - \epsilon)$
2	$\epsilon^2$	$2\epsilon(1 - \epsilon)$	$(1 - \epsilon)^2$

Table 1: Genotype error model,  $\Pr(O_i|G_i)$

$k$	$G_u$		
	0	1	2
1, 2	$(1 - p)$	0	$p$
3, 4	$(1 - p)^2$	$2p(1 - p)$	$p^2$
5, 6	$(1 - p)$	0	$p$
7, 8, 9	$(1 - p)^2$	$2p(1 - p)$	$p^2$

Table 2:  $\Pr(G_u|S_{iu} = k)$

$G_u$	$G_i$	0	1	2
		0	1	2
$S_{iu} = 9$	0, 1, 2	$(1 - p)^2$	$2p(1 - p)$	$p^2$
	0	$1 - p$	$p$	0
	1	$(1 - p)/2$	$1/2$	$p/2$
$S_{iu} = 8$	2	0	$1 - p$	$p$
	0	1	0	0
	1	0	1	0
$S_{iu} = 7$	2	0	0	1
	0, 2	$(1 - p)^2$	$2p(1 - p)$	$p^2$
	1	0	0	0
$S_{iu} = 6$	0	$1 - p$	$p$	0
	1	0	0	0
	2	0	$1 - p$	$p$
$S_{iu} = 5$	0, 1, 2	$1 - p$	0	$p$
	0	1	0	0
	1	$1/2$	0	$1/2$
$S_{iu} = 4$	2	0	0	1
	0, 2	$1 - p$	0	$p$
	1	0	0	0
$S_{iu} = 3$	0	1	0	0
	1	$1/2$	0	$1/2$
	2	0	0	1
$S_{iu} = 2$	0, 2	$1 - p$	0	$p$
	1	0	0	0
	0	1	0	0
$S_{iu} = 1$	1	0	1	0
	2	0	0	1
	0	1	0	0

Table 3:  $\Pr(G_i|G_u, S_{iu} = k)$

$S_{u1}$	$S_{u2}$	$S_{12}$
9	8	8
8	8	9
8	8	8
8	6	5
8	5	6
8	5	5
5	5	2
4	3	8
3	2	5

Table 4: Informative 3-way configurations. Note that because individuals 1 and 2 are unordered, we list only one of the corresponding ordered 3-way configurations (e.g. configuration (8, 9, 8) is also informative but has the same probabilities as configuration (9, 8, 8) with individuals 1 and 2 switched). Note that for configuration (8, 8, 7), (7, 7, 7) and (3, 3, 7)  $\Pr(G_1, G_2|G_u, S_{u12}) = \Pr(G_1|G_u, S_{u1}) = \Pr(G_2|G_u, S_{u2})$ , though we do not seek these configurations in our algorithm.

		$G_2$			
		$G_1$	0	1	2
$G_u = 0$	0	0	$q^2$	0	0
	1	0	$pq$	$pq$	0
	2	0	0	$p^2$	0
$G_u = 1$	0	0	$q^2/2$	$q^2/2$	0
	1	0	$qp/2$	$qp$	$qp/2$
	2	0	0	$p^2/2$	$p^2/2$
$G_u = 2$	0	0	0	$q^2$	0
	1	0	0	$pq$	$pq$
	2	0	0	0	$p^2$

Table 5: Configuration ( $S_{u1} = 9, S_{u2} = 8, S_{12} = 8$ ), probabilities  $\Pr(G_1, G_2|G_u, S_{u12})$

		$G_2$			
		$G_1$	0	1	2
$G_u = 0$	0	0	$q^2$	$pq$	0
	1	0	$pq$	$p^2$	0
	2	0	0	0	0
$G_u = 1$	0	0	0	$q^2/2$	$pq/2$
	1	0	$q^2/2$	$pq$	$p^2/2$
	2	0	$pq/2$	$p^2/2$	0
$G_u = 2$	0	0	0	0	0
	1	0	0	$q^2$	$pq$
	2	0	0	$pq$	$p^2$

Table 6: Configuration ( $S_{u1} = 8, S_{u2} = 8, S_{12} = 9$ ), probabilities  $\Pr(G_1, G_2|G_u, S_{u12})$

		$G_2$		
		0	1	2
$G_u = 0$	0	$q(q+1)/2$	$pq/2$	0
	1	$pq/2$	$p(p+1)/2$	0
	2	0	0	0
$G_u = 1$	0	$q^2/4$	$q(p+1)/4$	0
	1	$q(p+1)/4$	$(p^2+q^2)/4$	$p(q+1)/4$
	2	0	$p(q+1)/4$	$p^2/4$
$G_u = 2$	0	0	0	0
	1	0	$q(q+1)/2$	$pq/2$
	2	0	$pq/2$	$p(p+1)/2$

Table 7: Configuration ( $S_{u1} = 8, S_{u2} = 8, S_{12} = 8$ ), probabilities  $\Pr(G_1, G_2 | G_u, S_{u12})$

		$G_2$		
		0	1	2
$G_u = 0$	0	$q$	0	0
	1	0	0	$p$
	2	0	0	0
$G_u = 1$	0	$q/2$	0	0
	1	$q/2$	0	$p/2$
	2	0	0	$p/2$
$G_u = 2$	0	0	0	0
	1	$q$	0	0
	2	0	0	$p$

Table 8: Configuration ( $S_{u1} = 8, S_{u2} = 6, S_{12} = 5$ ), probabilities  $\Pr(G_1, G_2 | G_u, S_{u12})$

		$G_2$		
		0	1	2
$G_u = 0$	0	$q$	0	0
	1	$p$	0	0
	2	0	0	0
$G_u = 1$	0	0	0	$q/2$
	1	$q/2$	0	$p/2$
	2	$p/2$	0	0
$G_u = 2$	0	0	0	0
	1	0	0	$q$
	2	0	0	$p$

Table 9: Configuration ( $S_{u1} = 8, S_{u2} = 5, S_{12} = 6$ ), probabilities  $\Pr(G_1, G_2 | G_u, S_{u12})$



		$G_2$		
		0	1	2
$G_u = 0$	$G_1$	0	$q$	0
	0	$q$	0	0
	1	$p$	0	0
$G_u = 1$	$G_1$	0	$q/2$	0
	0	$q/2$	0	0
	1	$p/2$	0	$q/2$
$G_u = 2$	$G_1$	0	0	0
	0	0	0	0
	1	0	0	$q$
	2	0	0	$p$

Table 10: Configuration  $(S_{u1} = 8, S_{u2} = 5, S_{12} = 5)$ , probabilities  $\Pr(G_1, G_2|G_u, S_{u12})$

		$G_2$		
		0	1	2
$G_u = 0$	$G_1$	0	1	0
	0	1	0	0
	1	0	0	0
$G_u = 1$	$G_1$	0	0	$1/2$
	0	0	0	$1/2$
	1	0	0	0
$G_u = 2$	$G_1$	0	0	0
	0	0	0	0
	1	0	0	0
	2	0	0	1

Table 11: Configuration  $(S_{u1} = 5, S_{u2} = 5, S_{12} = 2)$ , probabilities  $\Pr(G_1, G_2|G_u, S_{u12})$

		$G_2$		
		0	1	2
$G_u = 0$	$G_1$	0	$q^2$	0
	0	$q^2$	0	0
	1	$pq$	$pq$	0
$G_u = 1$	$G_1$	0	0	0
	0	0	0	0
	1	0	0	0
$G_u = 2$	$G_1$	0	0	0
	0	0	$q^2$	0
	1	0	$pq$	$pq$
	2	0	0	$p^2$

Table 12: Configuration  $(S_{u1} = 4, S_{u2} = 3, S_{12} = 8)$ , probabilities  $\Pr(G_1, G_2|G_u, S_{u12})$

		$G_2$		
	$G_1$	0	1	2
$G_u = 0$	0	$q$	0	0
	1	0	0	$p$
	2	0	0	0
$G_u = 1$	0	0	0	0
	1	0	0	0
	2	0	0	0
$G_u = 2$	0	0	0	0
	1	$q$	0	0
	2	0	0	$p$

Table 13: Configuration ( $S_{u1} = 3, S_{u2} = 2, S_{12} = 5$ ), probabilities  $\Pr(G_1, G_2|G_u, S_{u12})$

Info	SNP type	MAF	genotype certainty	Number of genotypes (%)	Impute2	Impute2 & Kinpute
> 0.4	shared	> 0.02	High	2911525 (59.7)	0.945	<b>0.969</b>
> 0.4	shared	> 0.02	Low	469987 (9.6)	0.634	<b>0.834</b>
> 0.4	shared	$\leq 0.02$	High	387932 (8.0)	0.869	<b>0.931</b>
> 0.4	shared	$\leq 0.02$	Low	15460 (0.3)	0.640	<b>0.902</b>
> 0.4	private	> 0.02	High	201672 (4.1)	0.927	<b>0.960</b>
> 0.4	private	> 0.02	Low	31896 (0.7)	0.626	<b>0.836</b>
> 0.4	private	$\leq 0.02$	High	54182 (1.1)	0.789	<b>0.835</b>
> 0.4	private	$\leq 0.02$	Low	3178 (0.1)	0.691	<b>0.824</b>
< 0.4	shared	> 0.02	High	235481 (4.8)	0.045	<b>0.565</b>
< 0.4	shared	> 0.02	Low	106855 (2.2)	0.398	<b>0.846</b>
< 0.4	shared	$\leq 0.02$	High	284529 (5.8)	0.000362	<b>0.443</b>
< 0.4	shared	$\leq 0.02$	Low	8751 (0.2)	0.0473	<b>0.873</b>
< 0.4	private	> 0.02	High	28684 (0.6)	0.0240	<b>0.671</b>
< 0.4	private	> 0.02	Low	7988 (0.2)	0.308	<b>0.890</b>
< 0.4	private	$\leq 0.02$	High	123111 (2.5)	0.00295	<b>0.158</b>
< 0.4	private	$\leq 0.02$	Low	5913 (0.1)	0.0843	<b>0.654</b>

Table 14: Heterozygote sensitivity

Info	SNP type	MAF	genotype certainty	Number of genotypes (%)	Impute2	Impute2 & Kinpute
> 0.4	shared	> 0.02	High	2911525 (59.7)	<b>0.953</b>	0.950
> 0.4	shared	> 0.02	Low	469987 (9.6)	0.590	<b>0.766</b>
> 0.4	shared	$\leq$ 0.02	High	387932 (8.0)	0.961	<b>0.965</b>
> 0.4	shared	$\leq$ 0.02	Low	15460 (0.3)	0.758	<b>0.909</b>
> 0.4	private	> 0.02	High	201672 (4.1)	<b>0.934</b>	0.932
> 0.4	private	> 0.02	Low	31896 (0.7)	0.579	<b>0.752</b>
> 0.4	private	$\leq$ 0.02	High	54182 (1.1)	0.670	<b>0.755</b>
> 0.4	private	$\leq$ 0.02	Low	3178 (0.1)	0.330	<b>0.630</b>
< 0.4	shared	> 0.02	High	235481 (4.8)	0.531	<b>0.761</b>
< 0.4	shared	> 0.02	Low	106855 (2.2)	0.437	<b>0.741</b>
< 0.4	shared	$\leq$ 0.02	High	284529 (5.8)	0.429	<b>0.923</b>
< 0.4	shared	$\leq$ 0.02	Low	8751 (0.2)	0.502	<b>0.892</b>
< 0.4	private	> 0.02	High	28684 (0.6)	0.493	<b>0.844</b>
< 0.4	private	> 0.02	Low	7988 (0.2)	0.404	<b>0.733</b>
< 0.4	private	$\leq$ 0.02	High	123111 (2.5)	0.211	<b>0.585</b>
< 0.4	private	$\leq$ 0.02	Low	5913 (0.1)	0.109	<b>0.391</b>

Table 15: Heterozygote positive predictive value

Info	SNP type	MAF	genotype certainty	Number of genotypes (%)	Impute2	Impute2 & Kinpute
> 0.4	shared	> 0.02	High	2911525 (59.7)	0.964	<b>0.972</b>
> 0.4	shared	> 0.02	Low	469987 (9.6)	0.611	<b>0.810</b>
> 0.4	shared	$\leq$ 0.02	High	387932 (8.0)	0.987	<b>0.992</b>
> 0.4	shared	$\leq$ 0.02	Low	15460 (0.3)	0.661	<b>0.888</b>
> 0.4	private	> 0.02	High	201672 (4.1)	0.966	<b>0.973</b>
> 0.4	private	> 0.02	Low	31896 (0.7)	0.631	<b>0.810</b>
> 0.4	private	$\leq$ 0.02	High	54182 (1.1)	0.983	<b>0.987</b>
> 0.4	private	$\leq$ 0.02	Low	3178 (0.1)	0.682	<b>0.879</b>
< 0.4	shared	> 0.02	High	235481 (4.8)	0.933	<b>0.961</b>
< 0.4	shared	> 0.02	Low	106855 (2.2)	0.569	<b>0.838</b>
< 0.4	shared	$\leq$ 0.02	High	284529 (5.8)	0.970	<b>0.982</b>
< 0.4	shared	$\leq$ 0.02	Low	8751 (0.2)	0.612	<b>0.912</b>
< 0.4	private	> 0.02	High	28684 (0.6)	0.945	<b>0.975</b>
< 0.4	private	> 0.02	Low	7988 (0.2)	0.615	<b>0.862</b>
< 0.4	private	$\leq$ 0.02	High	123111 (2.5)	<b>0.989</b>	<b>0.989</b>
< 0.4	private	$\leq$ 0.02	Low	5913 (0.1)	0.901	<b>0.919</b>

Table 16: Concordance rates