

1 **Comparative transcriptomics of multidrug-resistant *Acinetobacter***
2 ***baumannii* in response to antibiotic treatments**

3 Hao Qin^{1,2*}, Norman Wai-Sing Lo³, Jacky Fong-Chuen Loo¹, Xiao Lin^{1,2}, Aldrin Kay-
4 Yuen Yim^{1,2*}, Stephen Kwok-Wing Tsui⁴, Terrence Chi-Kong Lau⁵, Margaret Ip³, and
5 Ting-Fung Chan^{1,2#}

6 ¹School of Life Sciences, The Chinese University of Hong Kong, Hong Kong SAR,
7 China

8 ²Partner State Key Laboratory of Agrobiotechnology, The Chinese University of Hong
9 Kong, Hong Kong SAR, China

10 ³Department of Microbiology, The Chinese University of Hong Kong, Hong Kong SAR,
11 China

12 ⁴School of Biomedical Sciences, The Chinese University of Hong Kong, Hong Kong
13 SAR, China

14 ⁵Department of Biomedical Sciences, City University of Hong Kong, Hong Kong SAR,
15 China

16

17

18

19 #Address correspondence to Ting-Fung Chan, tf.chan@cuhk.edu.hk.

20 *Present address: Hao Qin, 3D Medicines Corporation, Shang Hai, China; Aldrin Kay-
21 Yuen Yim, Washington University School of Medicine, Saint Louis, MO, USA.

22 N.W.S.L and J.F.C.L contributed equally to this work.

23 **Supplementary Methods**

24 **Angular based linear regression.**

Variable symbol	Variable explanation
n	Total number of genes
m	Total number of samples
v	Raw expression of a certain gene in a certain sample
\vec{v}_i	Gene vector of gene i constructed by the expression of samples
\vec{r}	Reference vector constructed by the size factors of samples
L	The sum of cosine values of the angles between the gene vectors and the reference vector
R	Raw size factors of samples
s	Final size factors of samples

25

26 Angular based linear regression method assumes that in a pool of samples, there is a set
 27 of house-keeping genes, whose expressions are stable across samples. Due to biological
 28 variations, the expressions of a single gene might vary significantly revealed by the fold
 29 changes. But the overall variations are majorly contributed by the sequencing process,
 30 which are going to be normalized by this method. Finding the overall trend of
 31 variations can be approximated by least angle linear regression. Least angle linear
 32 regression calculates a line minimizing the sum of angles between the line and the gene
 33 vectors in an m -dimensional space constructed by the expressions of n genes from m
 34 samples. The angles directly indicate the variations of genes' expressions. RNA-seq

35 measured expressions tend to have larger standard deviations when the absolute
 36 magnitudes are large. Least-square regression assumes the standard deviations are the
 37 same for all the data points, which can skew the regression line in RNA-seq. Hence
 38 considering the angle between the data points and the regression line as error model is
 39 more appropriate.

40

41 Assuming the vector of gene i is \vec{v}_i , which is constructed by the expressions v of this
 42 gene in all samples. Least angle regression minimizes the sum of angles between gene
 43 vectors \vec{v}_i and reference vector \vec{r} , which is the unit vector representing the regression
 44 line. There are many ways to find the least angle regression line. Maximizing the sum
 45 of cosine values of the angles is an easy method, because cosine function is strictly
 46 monotonic when the angle is between 0° and 180° . The regression line can be
 47 approximated by this equation:

48 The sum of angles $\sim L = \sum_{i=1}^n \frac{\vec{r} \cdot \vec{v}_i}{|\vec{v}_i|}$ (1)

49 where

50
$$\vec{r} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_m \end{bmatrix}$$

51
$$\sum_{i=1}^m r_i^2 = 1$$

52 The regression line, or the regression vector, generates the maximum value of function
 53 L . To find the extreme values, the method of Lagrange multipliers is applied:

54
$$\bigwedge (r_1, r_2 \dots r_m, \lambda) = L(r_1, r_2 \dots r_m) + \lambda(r_1^2 + r_2^2 + \dots + r_m^2 - 1)$$
 (2)

55 The partial derivatives are set to zero to find the extreme values:

56 From sample 1 to sample m :

57
$$\frac{\partial \Lambda}{\partial r_i} = 0$$

58 And:

59
$$\frac{\partial \Lambda}{\partial \lambda} = 0$$

60 There is only one extreme value in function L , which yields the maximum value. The
61 vector of the regression line actually equals to the sum of all unit gene vectors, which
62 can be calculated by:

63
$$R_k = \sum_{i=1}^n \frac{v_{ik}}{|\vec{v}_i|}, \quad \text{where } k \in [1, m] \quad (3)$$

64 The size factor of a certain sample s_k can thus be calculated as:

65
$$s_k = \frac{R_k}{|\vec{R}|} \cdot \sqrt{m} \quad (4)$$

66 After getting the size factors, the normalization can be achieved by a single step:

67
$$\text{Normalized expression} = \frac{\text{Raw expression}}{\text{Size factor}}$$

68

69 **Finding the best regression model.**

Variable symbol	Variable explanation
w	Normalized expression
u	Re-scaled normalized expression
μ	Mean of all expression
σ	Standard deviation of all expressions
P	Probability of a certain gene under given model

70

71 The most critical process in angular based linear regression method is to determine the
72 set of house-keeping genes. It could be determined based on existing knowledge, which
73 is commonly not available. Hence to tackle the problem, identifying the best set of
74 genes in normalization is another direction, which is only based on the RNA-seq results.

75

76 A way to approach the best set of house-keeping genes is using a probabilistic model
77 based on Gaussian distribution. This method calculates a likelihood to indicate the
78 goodness of the model. In this model, house-keeping genes are considered following
79 the Gaussian distribution, while non-house-keeping genes are considered skewed and
80 thus becoming outliers. The best model has the largest likelihood.

81 Assuming that the normalized expression of gene i in sample k is w_{ik} , all the
82 normalized expressions should be re-scaled to u_{ik} by the following calculation to set
83 the equal weight to every gene:

$$84 \quad u_{ik} = \frac{w_{ik}}{\bar{w}_i} \quad (5)$$

85 where \bar{w}_i is the mean of the normalized expressions of gene i . Afterwards all the
86 expressions are pooled. The mean μ and standard deviation σ are calculated by:

$$87 \quad \mu = \frac{1}{m \times n} \cdot \sum_{i=1}^n \sum_{k=1}^m u_{ik}$$
$$88 \quad \sigma = \sqrt{\frac{\sum_{i=1}^n \sum_{k=1}^m (u_{ik} - \mu)^2}{m \times n - 1}}$$

89 Under the size factors calculated by the set of house-keeping genes, the likelihood can
90 be calculated by:

91
$$\text{Likelihood} = \sum_{i=1}^n \log_{10} \left(\prod_{k=1}^m P_{ik} \right) \quad (6)$$

92 Where the probability P_{ik} of gene i in sample k is:

93
$$P_{ik}$$

94
$$= \begin{cases} \Phi \left(- \left| \frac{u_{ik} - \mu}{\sigma} \right| \right) \times 2, & \text{if } i \text{ belongs to housekeeping genes.} \\ 1 - \Phi \left(- \left| \frac{u_{ik} - \mu}{\sigma} \right| \right) \times 2, & \text{if } i \text{ doesn't belong to housekeeping genes.} \end{cases}$$

95 where Φ is the cumulative function of Gaussian distribution.

96

97 The simplest algorithm to find the best set of house-keeping genes is to calculate the
 98 likelihoods of all combinations. However it is resource and time consuming, which
 99 could be impossible when the number of genes is huge. Therefore better algorithm is
 100 indispensable to reduce the combinations to be calculated. One strategy is to find the
 101 best number of house-keeping genes, and then refine the combination of the house-
 102 keeping gene set around the number. Firstly to find the best number of house-keeping
 103 genes, the initial size factors are calculated with the whole set of genes. Then in each
 104 round, the gene, which has the smallest likelihood in the house-keeping gene set of last
 105 round, is eliminated from the house-keeping gene set and the new initial size factors
 106 are calculated based on the new set of house-keeping genes, until all the genes are
 107 eliminated from the house-keeping gene set. The number of the house-keeping genes,
 108 which has the maximum likelihood in this step, is selected. Secondly, to further refine
 109 the house-keeping gene set, starting from the best gene set in the last step, the state of
 110 the gene with the smallest probability is replaced by its opposite state, which means
 111 that if the gene is in the house-keeping gene set, it is eliminated and vice versa. Then

112 the likelihood of the new set of house-keeping genes is calculated. If the new likelihood
113 is larger than the old likelihood, this step will be repeated until the new likelihood is
114 not larger anymore. After this step, the best set of house-keeping genes is determined
115 and chosen in the normalization.

116

117 **Normalizing samples in different conditions.** In this study, different samples have
118 various conditions. Between conditions different sets of house-keeping genes should be
119 assumed. Therefore the size factors should be calculated in a progressive process. This
120 process includes two steps. Firstly, the condition which has the most number of samples
121 is chosen as the base pool to set the standard scale, and the size factors s of these
122 samples are calculated. In this study, the samples in antibiotic-free medium harvested
123 at mid-log phase are normalized first. Secondly, relative to a certain sample q in the
124 base pool, the samples which are conditionally related to sample q are pooled
125 accordingly and their relative size factors in the new scale are calculated. In the third
126 step, the relative size factors can be converted to standard scale through the relative size
127 factor of sample q in the new scale:

128
$$s_k = \frac{s'_k \times s_q}{s'_q} \quad (7)$$

129 where k is the sample pooled in the new scale, which is needed to be converted to
130 standard scale, s is the size factor in the standard scale and s' is the size factor in the
131 new scale. In this study, all antibiotic-treated samples harvested at mid-log phase were
132 grouped with the normal mid-log phase samples accordingly and their size factors in
133 the standard scale were calculated. At last repeat the step 2 and step 3 to the rest of the

134 samples until all size factors in standard scales are calculated. In this study, the samples
 135 treated by antibiotics and harvested at stationary phase are grouped with the samples
 136 treated by antibiotics and harvested at mid-log phase accordingly. There is a special
 137 case. Since the mid-log phase sample of R4 treated by amikacin is missing, the
 138 stationary phase sample of R4 treated by amikacin was grouped with the stationary
 139 phase sample of R5 treated by amikacin, which is also an amikacin resistant strain.

140

141 **Prioritizing differential genes in two pooled groups.**

Variable symbol	Variable explanation
w	Normalized expression
\vec{v}	Gene vector constructed by the expressions of a certain gene in all samples
\vec{A}	The trend vector where the genes are prioritized
\vec{B}	The opposite trend vector
S	The subspace constructed by \vec{A} and \vec{B}
c	The coordinate of a certain gene in the subspace S
θ	The angle between trend vector \vec{A} and the projected gene vector in subspace S
\vec{N}	The central vector in the subspace constructed by the samples in group A
φ	The error angle between the central vector \vec{N} and the projected gene vector in subspace constructed by the

	samples in group A
Θ	The sorting parameter

142

143 It is common to sort the genes that are differentially expressed in two groups. In this
144 study, the genes that are differentially expressed between two groups, multidrug
145 resistant strains and drug sensitive strains, are targets to be further analyzed. Each
146 strains are biologically different and grouped. Two groups are labeled as **A** and **B**,
147 whose number of samples are t_A and t_B respectively. In an $m = t_A + t_B$
148 dimensional space, all genes, whose number is n , become vectors, constructed by their
149 expression w_i and annotated by \vec{v}_i for gene i :

$$150 \quad \vec{v}_i = \begin{bmatrix} w_{iA1} \\ w_{iA2} \\ \vdots \\ w_{iAt_A} \\ w_{iB1} \\ w_{iB2} \\ \vdots \\ w_{iBt_B} \end{bmatrix}, \text{ where } i \in [1, n]$$

151 If the genes need to be prioritized by the trend where the genes are only expressed in
152 samples of group **A**. The trend vector \vec{A} is defined as:

$$153 \quad \vec{A} = \begin{bmatrix} \left. \begin{matrix} 1 \\ \vdots \\ 1 \end{matrix} \right\} t_A \\ 0 \\ \left. \begin{matrix} \vdots \\ 0 \end{matrix} \right\} t_B \end{bmatrix}$$

154 Similarly, the opposite trend of \vec{A} is annotated by \vec{B} , where the genes are only
155 expressed in samples of group **B**. It is defined as:

156

$$\vec{B} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{matrix} t_A \\ t_B \end{matrix}$$

157 \vec{A} and \vec{B} construct a sub-space $S = \text{span}\{\vec{A}, \vec{B}\}$ in the m dimensional space, where158 \vec{A} and \vec{B} are a pair of orthogonal basis. Each gene in the m dimensional space can159 be projected to the sub-space S . For a gene i , the coordinates in S can be converted by:

160
$$c_A = \frac{\vec{v}_i \cdot \vec{A}}{|\vec{A}|^2}, c_B = \frac{\vec{v}_i \cdot \vec{B}}{|\vec{B}|^2} \quad (8)$$

161 To calculate the prioritizing parameter for a group, the reference line is first set to the

162 trend vector of the group. For example, the genes are going to be prioritized according

163 to trend \vec{A} . The angle θ between the projection of gene vector \vec{v}_i in subspace S and164 \vec{A} can measure the extent of differential expression of gene i between group **A** and165 group **B**. If θ is close to 0° , the gene i has higher expressions in group **A**. If θ is close166 to 90° , the gene i has higher expressions in group **B**. If θ is close to 45° , the gene i has167 similar expression levels between group **A** and group **B**. Several trigonometric168 functions can imply the magnitude of θ , which can be calculated by:

169
$$\sin \theta = \frac{c_B}{\sqrt{c_A^2 + c_B^2}} \quad (9)$$

170

171 Nevertheless, the projection process disregards the variations within groups, causing

172 ties in the ranking when two or more genes have the same number difference of samples

173 that express the gene between the two groups. The ties can be further broken by

174 integrating an error angle φ to the θ . Supposing the genes are sorted to the trend \vec{A} ,

175 the error angle is calculated in the sub-space constructed by all the expressions of

176 samples in group **A**. The error angle φ is the angle between the projected gene vector

177 $\text{proj}_{\mathbf{A}} \mathbf{v}_i$ in the sub-space of group **A** and the central vector \vec{N} .

$$178 \quad \cos \varphi = \frac{\text{proj}_{\mathbf{A}} \mathbf{v}_i \cdot \vec{N}}{|\text{proj}_{\mathbf{A}} \mathbf{v}_i| \cdot |\vec{N}|}, \text{ where } \text{proj}_{\mathbf{A}} \mathbf{v}_i = \begin{bmatrix} w_{iA1} \\ \vdots \\ w_{iAt_A} \end{bmatrix}, \vec{N} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad (10)$$

179 This equation can be further expanded to:

$$180 \quad \cos \varphi = \frac{\sum_{j=1}^{t_A} w_{iAj}}{\sqrt{\sum_{j=1}^{t_A} w_{iAj}^2} \cdot \sqrt{t_A}}$$

$$181 \quad \cos \varphi = \sqrt{\frac{\sum_{j=1}^{t_A} w_{iAj}^2 + 2(w_{iA1} \cdot w_{iA2} + \dots + w_{iAa} \cdot w_{iAb} + \dots + w_{iA(t_A-1)} \cdot w_{iAt_A})}{(\sum_{j=1}^{t_A} w_{iAj}^2) \cdot t_A}}$$

$$182 \quad \cos \varphi = \sqrt{\frac{1}{t_A} + \frac{2(w_{iA1} \cdot w_{iA2} + \dots + w_{iAa} \cdot w_{iAb} + \dots + w_{iA(t_A-1)} \cdot w_{iAt_A})}{(\sum_{j=1}^{t_A} w_{iAj}^2) \cdot t_A}} \quad (11)$$

183 where:

$$184 \quad a, b \in [1, t_A], \quad a \neq b$$

$$185 \quad w_{iAj} \geq 0, \quad j \in [1, t_A]$$

$$186 \quad \sum_{j=1}^{t_A} w_{iAj}^2 > 0$$

187 The maximum angle of φ is 90° because:

$$188 \quad \cos \varphi \geq \sqrt{\frac{1}{t_A}}$$

$$189 \quad \min \cos \varphi = \lim_{t_A \rightarrow +\infty} \sqrt{\frac{1}{t_A}} = 0$$

$$190 \quad \max \varphi = 90^\circ$$

191

192 When a certain gene i is only expressed in the samples of group **A**, even only in one

193 sample, the gene should always be considered differentially expressed in group **A**,

