

# Supporting Information

## 1 Analytical Calculation

Under the model assumptions we can write down an expression for  $P_{n,\alpha}(f)$  which reflects a process of drawing  $N_g$  independent "true" correlations  $Z_t$  from the distribution  $q(Z_t)$ , each of which is submitted to a Gaussian noise of variance  $\sigma_n^2$ , and identifying the  $N_{TOP}(= \alpha N_g)$  top genes. Submitting the  $N_g$  true values to another realization of the noise, we obtain another list of  $N_{TOP}$  genes. For finite  $n$ , the lists are expected to be different due to noise (non-vanishing  $\sigma_n^2$ ). The probability to obtain an overlap  $f$  between two PGLs,  $P_{n,\alpha}(f)$ , is given by eq. (2) in the paper, which is specified here in more details:

$$\begin{aligned}
 P_{n,\alpha}(f) = & \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \sum_{h,l \in \{0,1\}^{N_g}} \left\{ \delta\left(\sum_{j=1}^{N_g} h_j - N_{TOP}\right) \delta\left(\sum_{j=1}^{N_g} l_j - N_{TOP}\right) \delta\left(\sum_{j=1}^{N_g} h_j l_j - f N_{TOP}\right) \right. \\
 & \prod_{j=1}^{N_g} \left[ (1-h_j) \int_{-x_1}^{x_1} dZ_{mj}^1 \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(Z_{mj}^1 - Z_j)^2}{2\sigma_n^2}} + h_j \left(1 - \int_{-x_1}^{x_1} dZ_{mj}^1 \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(Z_{mj}^1 - Z_j)^2}{2\sigma_n^2}}\right) \right] \\
 & \left. \prod_{k=1}^{N_g} \left[ (1-l_k) \int_{-x_2}^{x_2} dZ_{mk}^2 \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(Z_{mk}^2 - Z_k)^2}{2\sigma_n^2}} + l_k \left(1 - \int_{-x_2}^{x_2} dZ_{mk}^2 \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(Z_{mk}^2 - Z_k)^2}{2\sigma_n^2}}\right) \right] \right\}. \quad (1)
 \end{aligned}$$

Here  $\delta(\cdot)$  is the Kronecker delta,  $Z_j$  is the true correlation of the  $j$ th gene, and  $Z_{mj}^1, Z_{mk}^2$  are the measured correlations of the  $j$ th gene in the first and second realizations respectively.  $Nr$  is a normalization factor.  $h = (h_1, \dots, h_{N_g})$  and  $l = (l_1, \dots, l_{N_g})$  are binary vectors of size  $N_g$  whose non-zero elements correspond to the genes included in the PGLs of the first and the second realizations respectively. The integration variables  $x_1, x_2$  can be thought of as artificial 'thresholds' which separate the  $N_{TOP}$  top correlations from the rest in the two realizations. The density of  $x_1, x_2$  has no effect in the large  $N_g$  limit, and is thus omitted here. Replacing the delta functions in eq. (1) by their integral representations one obtains:

$$\begin{aligned}
 P_{n,\alpha}(f) = & \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \sum_{h,l \in \{0,1\}^{N_g}} \\
 & \left\{ \int_{-\pi}^{\pi} \frac{dy dz dw}{(2\pi)^3} e^{iy(\sum_{j=1}^{N_g} h_j - N_{TOP}) + iz(\sum_{j=1}^{N_g} l_j - N_{TOP}) + iw(\sum_{j=1}^{N_g} h_j l_j - f N_{TOP})} \right. \\
 & \left. \prod_{j=1}^{N_g} \left[ (1-h_j)P(x_1, Z_j) + h_j(1 - P(x_1, Z_j)) \right] \left[ (1-l_j)P(x_2, Z_j) + l_j(1 - P(x_2, Z_j)) \right] \right\}, \quad (2)
 \end{aligned}$$

where  $P(x, Z) \equiv P(x, Z, \sigma_n) = \int_{-x}^x dZ_m \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(Z_m - Z)^2}{2\sigma_n^2}}$ . (From now on, we shall omit the dependence on  $\sigma_n$

in  $P$ ). Simple manipulations yield

$$P_{n,\alpha}(f) = \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \int_{-\pi}^\pi \frac{dydzdw}{(2\pi)^3} \sum_{h,l \in \{0,1\}^{N_g}} \left\{ \prod_{j=1}^{N_g} \left[ (1-h_j)P(x_1, Z_j) + h_j(1-P(x_1, Z_j)) \right] \right. \\ \left. \left[ (1-l_j)P(x_2, Z_j) + l_j(1-P(x_2, Z_j)) \right] e^{(iyh_j + izl_j + iwh_j l_j)} e^{(-iyN_{TOP} - izN_{TOP} - iwfN_{TOP})} \right\}, \quad (3)$$

which results in

$$P_{n,\alpha}(f) = \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \int_{-\pi}^\pi \frac{dydzdw}{(2\pi)^3} \prod_{j=1}^{N_g} B(x_1, x_2, y, z, w, Z_j), \quad (4)$$

where

$$B(x_1, x_2, y, z, w, Z) = P(x_1, Z)P(x_2, Z) + P(x_1, Z)(1-P(x_2, Z))e^{iz} \\ + (1-P(x_1, Z))P(x_2, Z)e^{iy} + (1-P(x_1, Z))(1-P(x_2, Z))(e^{iy} + e^{iz} + e^{iw}). \quad (5)$$

For  $N_g \gg 1$ , one can approximate summation over the  $Z_j$  by integrating  $dq(Z)$ , which for symmetric  $q(Z)$  gives:

$$P_{n,\alpha}(f) = \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \int_{-\pi}^\pi \frac{dydzdw}{(2\pi)^3} \\ \exp\left(-N_g(i\alpha y + i\alpha z + i\alpha f w)\right) \exp\left(2N_g \int_0^\infty dZ q(Z) \ln(B(x_1, x_2, y, z, w, Z))\right). \quad (6)$$

Defining  $A$  as

$$A(x_1, x_2, y, z, w, Z) = P(x_1, Z)P(x_2, Z) \left( e^{-i(y+z+w)} - e^{-i(y+w)} - e^{-i(z+w)} + 1 \right) \\ + P(x_1, Z)(e^{-i(y+w)} - 1) + P(x_2, Z)(e^{-i(z+w)} - 1) + 1 \quad (7)$$

yields

$$P_{n,\alpha}(f) = \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \int_{-\pi}^\pi \frac{dydzdw}{(2\pi)^3} \\ \exp\left(-N_g(i(1-\alpha)y - i(1-\alpha)z - i(1-\alpha)f w)\right) \exp\left(2N_g \int_0^\infty dZ q(Z) \ln(A(x_1, x_2, y, z, w, Z))\right). \quad (8)$$

The above integral can be written as:

$$P_{n,\alpha}(f) = \frac{1}{Nr} \int_0^\infty dx_1 dx_2 \int_{-\pi}^\pi \frac{dydzdw}{(2\pi)^3} \exp(-N_g F), \quad (9)$$

where

$$F(x_1, x_2, y, z, w; f) = -i(1-\alpha)y - i(1-\alpha)z - i(1-\alpha)f w - 2 \int_0^\infty q(Z) dZ \ln(A(x_1, x_2, y, z, w, Z)). \quad (10)$$

The Saddle Point (SP) equations  $\nabla F = 0$  (where  $f$  is treated as a parameter) are:

$$\begin{aligned}
1 - \alpha &= 2 \int_0^\infty q(Z) dZ \frac{e^{-i(y+w)}(e^{-iz} - 1)P(x_1, Z)P(x_2, Z) + e^{-i(y+w)}P(x_1, Z)}{A(x_1, x_2, y, z, w, Z)} \\
1 - \alpha &= 2 \int_0^\infty q(Z) dZ \frac{e^{-i(z+w)}(e^{-iy} - 1)P(x_1, Z)P(x_2, Z) + e^{-i(z+w)}P(x_1, Z)}{A(x_1, x_2, y, z, w, Z)} \\
1 - \alpha f &= 2 \int_0^\infty q(Z) dZ \frac{e^{-iw}(e^{-i(y+z)} - e^{-iy} - e^{-iz})P(x_1, Z)P(x_2, Z) + e^{-i(y+w)}P(x_1, Z) + e^{-i(z+w)}P(x_2, Z)}{A(x_1, x_2, y, z, w, Z)} \\
0 &= \int_0^\infty q(Z) dZ \frac{P_{x_1}(x_1, Z) [(e^{-i(y+z+w)} - e^{-i(y+z)} - e^{-i(z+w)} + 1)P(x_2, Z) + e^{-i(y+w)} - 1]}{A(x_1, x_2, y, z, w, Z)} \\
0 &= \int_0^\infty q(Z) dZ \frac{P_{x_2}(x_2, Z) [(e^{-i(y+z+w)} - e^{-i(y+z)} - e^{-i(z+w)} + 1)P(x_1, Z) + e^{-i(z+w)} - 1]}{A(x_1, x_2, y, z, w, Z)}, \tag{11}
\end{aligned}$$

where  $P_x(x, Z)$  is the derivative of  $P(x, Z)$  with respect to  $x$ . For  $Nr$  the SP equations are those shown in (11) and an additional equation,  $w = 0$ , which is obtained from  $\frac{\partial F}{\partial f} = 0$ . Substituting  $w = 0$  in the last two SP equations in (11) one obtains  $y, z = 0$ , and the first three SP equations become:

$$\begin{aligned}
1 - \alpha &= 2 \int_0^\infty q(Z) dZ P(x_1, Z) \\
1 - \alpha &= 2 \int_0^\infty q(Z) dZ P(x_2, Z) \\
1 - \alpha f &= 2 \int_0^\infty q(Z) dZ (P(x_1, Z) + P(x_2, Z) - P(x_1, Z)P(x_2, Z)). \tag{12}
\end{aligned}$$

Note that the last equation can be written in a more meaningful way as:

$$\alpha f = 2 \int_0^\infty q(Z) dZ (1 - P(x_1, Z))(1 - P(x_2, Z)). \tag{13}$$

For very large  $N_g$ , the SP expansion gives:

$$P_{n,\alpha}(f) \sim \sqrt{\frac{N_g \det R}{\det H}} e^{-N_g(F(f) - F(f_n^*))}, \tag{14}$$

where  $R_{5 \times 5}$  and  $H_{6 \times 6}$  are the second derivative matrices of  $F$  at the saddle point with respect to  $(x_1, x_2, y, z, w)$  and  $(x_1, x_2, y, z, w, f)$  respectively,  $\det$  denotes matrix determinant and  $f_n^*$  is the value of  $f$  minimizing  $F$ . Thus,  $f_n^*$  is obtained by taking the value of  $f$  in the solution of the set of the SP eqs. (12) (which are easily solved numerically). For large  $N_g$ ,  $P_{n,\alpha}(f)$  gets a sharp maximum at  $f = f_n^*$ , and as  $N_g \rightarrow \infty$ ,  $P_{n,\alpha}(f)$  tends to a delta function at  $f = f_n^*$ . The meaning of this result is that for  $N_g \rightarrow \infty$ , and for a given finite number of samples  $n$ , the values of both  $x$  and  $f$  are independent of the specific selection of the  $n$  samples. Expanding eq. (14) into a series around  $f_n^*$  and keeping the leading term one obtains our final expression for  $P_{n,\alpha}(f)$ :

$$P_{n,\alpha}(f) = \frac{1}{\sqrt{2\pi\Sigma_n}} e^{-\frac{(f-f_n^*)^2}{2\Sigma_n^2}}, \tag{15}$$

where the variance  $\Sigma_n^2$  is given by:

$$\Sigma_n^2 = \frac{\det H}{2\pi N_g \det R}. \tag{16}$$

## 2 Simulations

### 2.1 Adjusting $\mu_g(i) - \mu_p(i)$ to fit the true distribution

As described in the paper, the two Gaussians  $G(\mu_g(i), \sigma_g(i))$  and  $G(\mu_p(i), \sigma_p(i))$  are approximating the probability distribution of the expression of gene  $i$  for  $n = N_s$ . However, we are interested in the true distributions,

namely those corresponding to infinite  $N_s$ . Therefore we have to re-scale  $\Delta\mu(i) \equiv \mu_g(i) - \mu_p(i)$  so that the distribution of the resulting correlations will fit  $q(Z_t)$ . The re-scaling can be done, for example, by keeping  $\mu_g(i)$  (and  $\sigma_g(i), \sigma_p(i)$ ), and changing  $\mu_p(i)$  such that we get:

$$\Delta\mu(i) = Z_{mi} \sqrt{\frac{V_t}{V_t + \sigma_n^2} \frac{P_L \sigma_g(i)^2 + (1 - P_L) \sigma_p(i)^2}{P_L(1 - P_L)(1 - Z_i^2)}}, \quad (17)$$

where  $P_L$  is the relative fraction of good outcome patients in the dataset.

## 2.2 Motivation for creating the Simulation model

The most straightforward way to perform simulations is the following: For each  $n$ , divide the dataset into the maximal number of non-overlapping training sets of size  $n$ ,  $K(n) \equiv \lfloor \frac{N_s}{n} \rfloor$ . (Clearly, allowing overlaps between the training sets will result in an overestimate of  $f$ ). For each training set generate a PGL, ending up with  $K(n)$  lists. Then calculate the overlaps between the  $K(n)/2$  independent pairs of lists. Repeat this procedure  $T$  times to obtain  $T \cdot K(n)/2$  overlap values whose mean and variance have to match the analytical prediction of  $f_n^*$ , and  $\Sigma_n^2$  respectively. We have found that performing the simulations in this way gives strong data-dependent fluctuations in the estimates of  $f_n^*$  (for a fixed value of  $n$ ), resulting in sometimes a non-monotonic behavior of  $f_n^*$  in  $n$ . This was observed both for the biological and simulated datasets (data not shown). We note that this instability in the estimate of  $f_n^*$  cannot be attributed to the lack of computational resources (i.e., too small number of repeats  $T$ ), and may occur even if one enumerates all possible  $\binom{N_s}{n \dots n}$  partitions for a given dataset.

To overcome this problem, we created our model, which in addition to eliminating the aforementioned phenomenon, allows to produce simulation results for unlimited  $n$  as opposed to the aforementioned procedure which is limited to  $n = N_s/2$ .

## 3 Checking the Model Assumptions on the Real Datasets

Our analytical calculation is based on four main assumptions:

*Assumption 1:* The distributions of the measured  $Z$ 's are Gaussian, centered for each gene around its  $Z_t$ .

*Assumption 2:* The variance  $\sigma_n^2$  is the same for all genes.

*Assumption 3:* The noise variables  $Z - Z_t$  are independent (i.e. uncorrelated noise for different genes).

*Assumption 4:*  $q(Z_t)$ , the distribution of the true correlations, can be approximated by a Gaussian with variance  $V_t$ . This assumption is easily generalized to represent  $q(Z_t)$  as a mixture of Gaussians.

The successful application of our method to real datasets depends on the extent to which our assumptions hold. We checked these assumptions on the six datasets analyzed in this paper. The validity of *Assumption 1* is demonstrated in Fig. 1, for five randomly peaked genes from each dataset. We have checked it also for many other genes, and the  $Z$ 's distributions of almost all genes were very well approximated by Gaussians.

Results for *Assumption 2* appear in Fig. 2. The histograms were generated by selecting 1000 random pairs of non-overlapping training sets of size  $n = N_s/2$ . The correlation of each gene with survival was calculated in each training set, and the variance of its correlation within each pair was recorded. This resulted in 1000 variances for each gene as obtained from the 1000 randomly generated pairs of training sets. The average of these 1000 values was the estimate for  $\sigma_n^2$ . The variance histograms of the six datasets are tightly centered around the mean value  $\sigma_n^2 = 1/(n - 3)$  (red vertical lines), implying that the datasets can be well described using *Assumption 2*. A relatively less centered histogram is obtained for Lung Cancer (3) which may explain the relatively high deviations between analytical prediction and simulations observed in this dataset.

Results for *Assumption 3* appear in Table 1. In most datasets the fitted  $a, b$  values satisfy  $a \approx 1$  and  $b \approx -1$ , which implies that the noise of the genes is uncorrelated (see explanation in **Materials and Methods**)

Results for *Assumption 4* are exhibited in Fig. 3. Since we used  $n = N_s$ , our sampling noise  $\sigma_n^2$  is rather small. The nice match between the fits and the real histograms of the measured  $q_n$  therefore reflects the validity of this assumption also for the true distribution  $q$ .

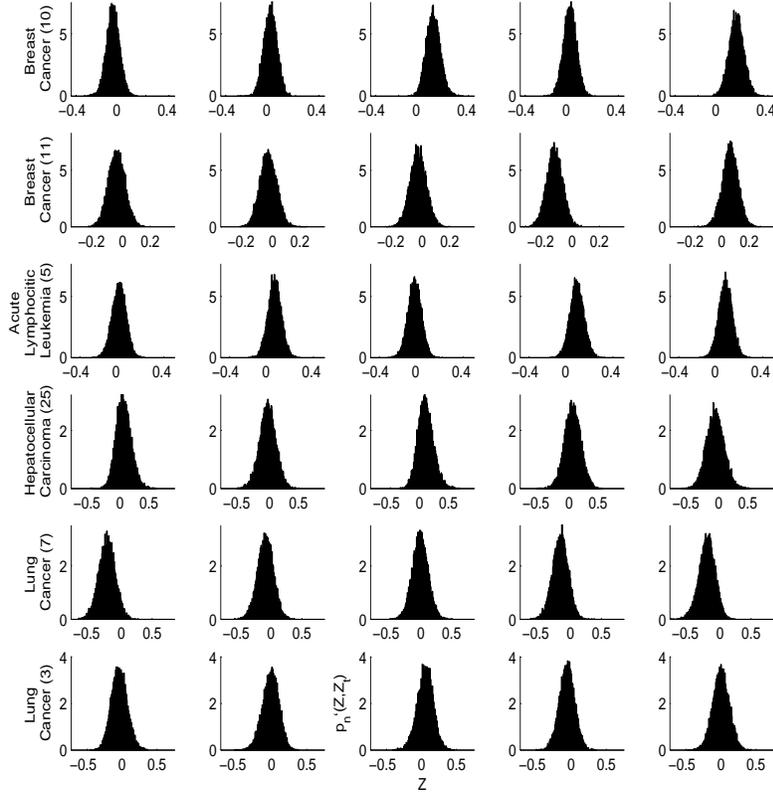


Figure 1: Histograms of the genes' Z-score for 5 randomly picked genes for each of the six datasets. The Z-scores were measured using  $\frac{N_s}{2}$  samples, which were randomly picked 10000 times out of the whole  $N_s$  samples. One can see that the measured Z-scores are to a very good approximation Normally distributed.

|                                 | $a$    | $b$     | $V_t$          |
|---------------------------------|--------|---------|----------------|
| Breast Cancer (10)              | 1.008  | -0.9913 | 0.004451       |
| Breast Cancer (11)              | 0.9806 | -0.9920 | 0.0030         |
| Acute Lymphocytic Leukaemia (5) | 0.9255 | -0.9917 | 0.0025         |
| Hepatocellular Carcinoma (25)   | 0.9645 | -0.9851 | 0.00824        |
| Lung Cancer (7)                 | 0.855  | -0.9719 | $5.618e - 005$ |
| Lung Cancer (3)                 | 1.011  | -1.003  | 0.002549       |

Table 1: Effective noise as measured in the six real datasets (see **Materials and Methods** for Details). The values  $a$ ,  $b$  and  $V_t$  are the coefficients of the fit  $V(n) = a(n - 3)^b + V_t$  where  $a(n - 3)^b$  and  $V_t$  are estimations for the effective  $\sigma_n^2$  and for the variance of the true distribution  $q(Z_t)$  respectively.

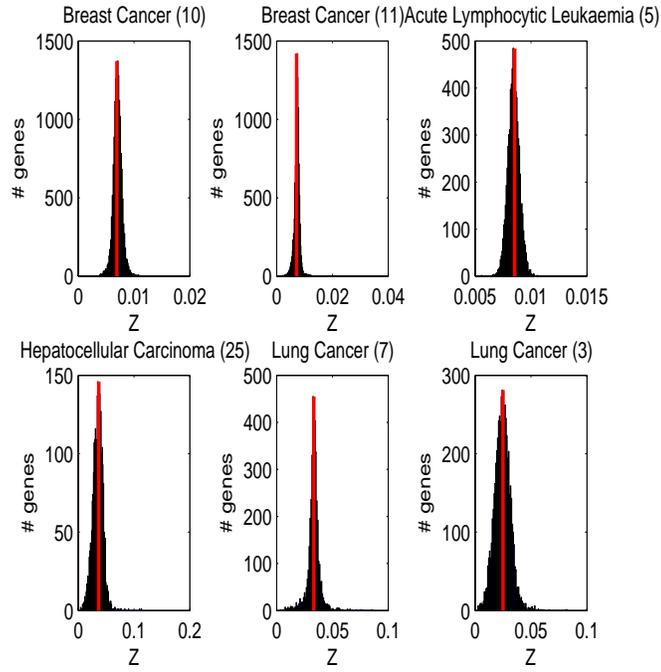


Figure 2: Histograms of  $\sigma_n^2$  in the six datasets. We used  $n = N_s/2$  which is the maximal training set size for a pair of non-overlapping training sets. The red vertical lines denote the expected value of the variance  $\sigma_n^2 = 1/(n-3)$ .

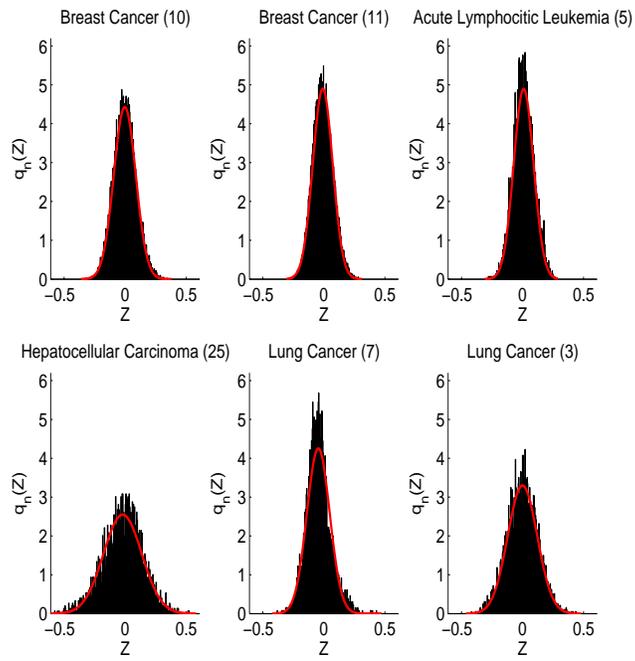


Figure 3: Histograms of the genes' Z-score for the six datasets. The Z-scores were measured using all  $N_s$  samples. The Gaussian fit of each histogram is plotted in red.