

PNAS



1

2 **Supporting Information for**

3 **OASIS: An interpretable, finite-sample valid alternative to Pearson's X^2 for scientific discovery**

4 **Tavor Z. Baharav, David Tse, Julia Salzman**

5 **Julia Salzman.**

6 **E-mail: julia.salzman@stanford.edu**

7 **This PDF file includes:**

8 Supporting text

9 Figs. S1 to S16

10 Tables S1 to S2

11 SI References

12 Supporting Information Text

13 S.1. Data availability

14 *Mycobacterium tuberculosis* data is publicly available under Accession ID PRJEB41116 (1). SARS-CoV-2 data is publicly
 15 available under Accession ID PRJNA817806 (2). Both datasets were chosen arbitrarily from the sequence read archive (3) and
 16 were analyzed with the same method with the same default parameters. Code for SPLASH (previously called NOMAD) is
 17 publicly available at <https://github.com/salzman-lab/splash>, along with code for OASIS inference. An optimized implementation of
 18 SPLASH called SPLASH 2 has been developed (4) using the bounds and techniques described in this work, and is available here
 19 <https://github.com/refresh-bio/R-SPLASH>. *M. tuberculosis* data was generated with this improved implementation of SPLASH.
 20 Then, the same optimization procedures were run to generate the optimized p-value bounds.

21 S.2. Proofs

22 As discussed, OASIS’s test statistic S has a simple linear algebraic characterization. For $X \in \mathbb{N}^{I \times J}$ the expected matrix E
 23 can be written as the rank 1 outer product between the row and column sums (dividing by M). The centered matrix is then
 24 right normalized to account for the unequal sampling depth, dividing by the square root of the column sums to appropriately
 25 normalize, yielding the centered and right normalized matrix \tilde{X} . Multiplying from the left by \mathbf{f} yields a deviation statistic for
 26 each column, which is aggregated in a \mathbf{c} -weighted sum to yield the final test statistic S .

27 To analyze OASIS’s test statistic, we define the random variables $\{Z_{j,k}\}$ for $j = 1, \dots, J$ and $k = 1, \dots, n_j$. $Z_{j,k}$ denotes
 28 the row identity of the k -th observation from the j -th column, so $Z_{j,k} \in [I]$. The table can then be equivalently constructed
 29 from the $\{Z_{j,k}\}$ by taking $X_{i,j} = \sum_{k=1}^{n_j} \mathbb{1}\{Z_{j,k} = i\}$. Under the null hypothesis, $\{Z_{j,k}\}$ are all independently drawn from the
 30 common probability distribution \mathbf{p} over $[I]$. This ordering $k = 1, \dots, n_j$ represents a random ordering of the counts observed
 31 and is for analysis purposes only. This reformulation allows OASIS’s test statistic to be expressed as a weighted sum of
 32 independent random variables, enabling the use of classical concentration inequalities.

S.2.A. Proof of original p-value bound. The original p-value bound proposed in (5) is:

$$\mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq s) \leq 2 \exp\left(-\frac{2(1-\xi)^2 s^2}{\sum_j c_j^2}\right) + 2 \exp\left(-\frac{2\xi^2 M s^2}{\left(\sum_j c_j \sqrt{n_j}\right)^2}\right)$$

$$\text{where } \xi = \left(1 + \sqrt{\frac{M \sum_j c_j^2}{\left(\sum_j c_j \sqrt{n_j}\right)^2}}\right)^{-1}.$$

Using the quantity γ defined in this paper, this can be expressed as:

$$\mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq s) = 2 \exp\left(-\frac{2(1-\xi)^2 s^2}{\|\mathbf{c}\|^2}\right) + 2 \exp\left(-\frac{2\xi^2 s^2}{\|\mathbf{c}\|^2 \gamma}\right)$$

$$\text{where } \xi = \left(1 + \frac{1}{\sqrt{\gamma}}\right)^{-1}.$$

33 The p-value bound in (5) can be proved in the following manner. First, we estimate the expectation (unconditional on
 34 sample identity) of \mathbf{f} on the observations as $\bar{\mu}$. Then, we estimate the expectation of \mathbf{f} for each column j separately, as $\hat{\mu}_j$. We
 35 then construct an estimate for the deviation of a column j from the table average by computing $S_j = \sqrt{n_j}(\hat{\mu}_j - \bar{\mu})$, normalizing
 36 by $\sqrt{n_j}$ to ensure that each S_j will have essentially constant variance (up to the correlation between $\bar{\mu}$ and $\hat{\mu}_j$). Finally, we
 37 construct our overall test statistic S as the \mathbf{c} -weighted sum of the S_j , i.e. $S = \sum_j c_j S_j$. In summary:

$$\begin{aligned} \bar{\mu} &= \frac{1}{M} \sum_{i,j} f_i X_{i,j} = \frac{1}{M} \sum_{j,k} f_{Z_{j,k}} \\ \hat{\mu}_j &= \frac{1}{n_j} \sum_{i,j} f_i X_{i,j} = \frac{1}{n_j} \sum_{k=1}^{n_j} f_{Z_{j,k}} \\ S_j &= \sqrt{n_j}(\hat{\mu}_j - \bar{\mu}) \\ S &= \sum_j c_j S_j \end{aligned} \tag{1}$$

38 To provide a p-value bound, we use Hoeffding’s inequality for sums of independent, bounded random variables:

Lemma 1 (Hoeffding's inequality, Prop. 2.1 (6)). *Suppose that the random variables X_i , $i = 1, \dots, n$ are independent, and X_i has mean μ_i and $a_i \leq X_i \leq b_i$ almost surely. Then for all $t \geq 0$ we have:*

$$\mathbb{P} \left(\sum_{i=1}^n (X_i - \mu_i) \geq t \right) \leq \exp \left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right).$$

Taking the true (unknown) distribution of the $\{Z_{j,k}\}$ under the null to be \mathbf{p} , we define $\mu \triangleq \mathbb{E}_{Z \sim \mathbf{p}}[f_Z]$. Then, as in (5),

$$\begin{aligned} \mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq \epsilon) &\leq \mathbb{P} \left(\left| \sum_{j,k} c_j \frac{f_{Z_{j,k}} - \bar{\mu}}{\sqrt{n_j}} \right| \geq \epsilon \right) \\ &= \mathbb{P} \left(\left| \sum_{j,k} c_j \frac{f_{Z_{j,k}} - \mu}{\sqrt{n_j}} + (\mu - \bar{\mu}) \sum_j c_j \sqrt{n_j} \right| \geq \epsilon \right) \\ &\leq \min_{\xi \in (0,1)} \left[\mathbb{P} \left(\left| \sum_{j,k} c_j \frac{f_{Z_{j,k}} - \mu}{\sqrt{n_j}} \right| \geq (1 - \xi)\epsilon \right) + \mathbb{P} \left(\left| (\mu - \bar{\mu}) \sum_j c_j \sqrt{n_j} \right| \geq \xi\epsilon \right) \right] \\ &\stackrel{(a)}{=} \min_{\xi \in (0,1)} \left[\mathbb{P} \left(\left| \sum_{j,k} \frac{c_j}{\sqrt{n_j}} (f_{Z_{j,k}} - \mu) \right| \geq (1 - \xi)\epsilon \right) + \mathbb{P} \left(\left| \frac{1}{M} \sum_{j,k} (f_{Z_{j,k}} - \mu) \right| \geq \frac{\xi\epsilon}{\sum_j c_j \sqrt{n_j}} \right) \right] \\ &\stackrel{(b)}{\leq} \min_{\xi \in (0,1)} \left[2 \exp \left(-\frac{2(1 - \xi)^2 \epsilon^2}{\sum_{j,k} c_j^2 / n_j} \right) + 2 \exp \left(-\frac{\frac{2\xi^2 M^2 \epsilon^2}{(\sum_j c_j \sqrt{n_j})^2}}{M} \right) \right] \\ &= \min_{\xi \in (0,1)} \left[2 \exp \left(-\frac{2(1 - \xi)^2 \epsilon^2}{\sum_{j:n_j > 0} c_j^2} \right) + 2 \exp \left(-\frac{2\xi^2 M \epsilon^2}{(\sum_j c_j \sqrt{n_j})^2} \right) \right] \end{aligned} \quad [2]$$

39 The first inequality comes from a union bound (if neither of the two conditions are met, then $|S| < \epsilon$. (a) assumes that the
40 denominator $\sum_j c_j \sqrt{n_j} \neq 0$, as otherwise ξ can simply be taken to be 0. In (b) we utilize Hoeffding's inequality on the two
41 terms; in the first we note that the j, k -th term, $\frac{c_j}{\sqrt{n_j}}(f_{Z_{j,k}} - \mu)$ is indeed 0 mean, and is bounded as $\mu \in [0, 1]$ is simply a
42 constant offset that shifts both the min and max by the same amount (retaining the dynamic range of $f_{Z_{j,k}} \in [0, 1]$). Concretely,

$$43 \quad Y_{j,k} \triangleq \frac{c_j}{\sqrt{n_j}}(f_{Z_{j,k}} - \mu) \implies \frac{-c_j \mu}{\sqrt{n_j}} \leq Y_{j,k} \leq \frac{c_j(1 - \mu)}{\sqrt{n_j}} \quad [3]$$

44 Similarly, for the second term, each summand has a range of 1. This bound can easily be optimized over $\xi \in (0, 1)$ to within a
45 factor of 2 of optimum by equating the two terms, which is achieved when:

$$\xi = \left(1 + \sqrt{\frac{M \sum_{j:n_j > 0} c_j^2}{\sum_j c_j \sqrt{n_j}}} \right)^{-1}.$$

46 This leads to the stated p-value bound in (5) of

$$47 \quad \mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq s) \leq 2 \exp \left(-\frac{2(1 - \xi)^2 s^2}{\sum_{j:n_j > 0} c_j^2} \right) + 2 \exp \left(-\frac{2\xi^2 M s^2}{(\sum_j c_j \sqrt{n_j})^2} \right) \quad \text{with} \quad \xi = \left(1 + \sqrt{\frac{M \sum_j c_j^2}{(\sum_j c_j \sqrt{n_j})^2}} \right)^{-1}. \quad [4]$$

48 **S.2.B. Proof of improved p-values.** We restate below the finite-sample p-value bound proved in the main text.

Proposition 1 (Restatement of Proposition 1). *Under the null hypothesis, for any fixed $\mathbf{f} \in [0, 1]^I$ and $\mathbf{c} \in \mathbb{R}^J$ with $\|\mathbf{c}\|_2 \leq 1$, if $\gamma < 1$, the OASIS test statistic $S = S(\mathbf{f}, \mathbf{c})$ satisfies*

$$\mathbb{P}(|S| \geq s) \leq 2 \exp \left(-\frac{2s^2}{1 - \gamma} \right).$$

As before, we define $\bar{\mu} = \frac{1}{M} \mathbf{f}^\top X \mathbf{1}$ as the estimate of $\mu = \mathbb{E}_{Z \sim \mathbf{p}} [f_Z]$, where \mathbf{p} is the unknown common row distribution under the null. Due to the structure of the test statistic, the p-value bound in Equation (4) can be improved and simplified. Fixing \mathbf{f}, \mathbf{c} , the test statistic can be simplified as

$$\begin{aligned}
S(\mathbf{f}, \mathbf{c}) &= \sum_j c_j \sqrt{n_j} (\hat{\mu}_j - \bar{\mu}) \\
&= \left[\sum_j c_j \sqrt{n_j} \left(\frac{1}{n_j} \sum_{k=1}^{n_j} f_{Z_{j,k}} \right) \right] - \left[\sum_j c_j \sqrt{n_j} \left(\frac{1}{M} \sum_{\ell k} f_{Z_{\ell,k}} \right) \right] \\
&= \sum_{jk} \frac{c_j}{\sqrt{n_j}} f_{Z_{j,k}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \sum_{jk} f_{Z_{j,k}} \\
&= \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) f_{Z_{j,k}} \\
&= \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) (f_{Z_{j,k}} - \mu). \tag{5}
\end{aligned}$$

49 Note that, ignoring $f_{Z_{j,k}}$ in the second to last line, this summation is identically 0, to allow for the test statistic to have mean
50 0 independent of \mathbf{f} and \mathbf{p} . This allows us to subtract the mean $\mu \triangleq \mathbb{E}_{Z \sim \mathbf{p}} [f_Z]$ inside the summation in the last line, yielding a
51 sum of mean 0 terms.

52 Note that $\gamma = 1$ if and only if $\mathbf{c} = \pm \sqrt{n/M}$ by Cauchy-Schwarz, in which case each coefficient will be equal to 0, i.e.
53 $\left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) = 0$ for all j, k . Then $S(\mathbf{f}, \mathbf{c}) = 0$ with probability 1.

54 If $\gamma < 1$, then examining each term in the summation in Eq. (5) yields

$$Y_{j,k} \triangleq \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) (f_{Z_{j,k}} - \mu) \implies \left| \frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right| \left(\min_i f_i - \mu \right) \leq Y_{j,k} \leq \left| \frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right| \left(\max_i f_i - \mu \right). \tag{6}$$

55 At right we are simply claiming that $Y_{j,k}$ is almost surely bounded with the provided range. Thus, observing that each $Y_{j,k}$ has mean 0, and is bounded as above, we can apply Hoeffding's inequality to the sum of these $Y_{j,k}$. For the simplification below we consider the case where $\max_i f_i - \min_i f_i = 1$, but this is without loss of generality as dividing by $\max_i f_i - \min_i f_i$ simply rescales ϵ . Note that this condition of $\gamma < 1$ ensures that the denominator is nonzero.

$$\begin{aligned}
\mathbb{P} (|S(\mathbf{f}, \mathbf{c})| \geq \epsilon) &= \mathbb{P} \left(\left| \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) (f_{Z_{j,k}} - \mu) \right| \geq \epsilon \right) \\
&\leq 2 \exp \left(- \frac{2\epsilon^2}{\sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right)^2} \right) \\
&= 2 \exp \left(- \frac{2\epsilon^2}{\sum_j \left(c_j - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \sqrt{n_j} \right)^2} \right) \\
&= 2 \exp \left(- \frac{2\epsilon^2}{\sum_j \left[c_j^2 - 2c_j \sqrt{n_j} \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} + \frac{(\sum_{\ell} c_{\ell} \sqrt{n_{\ell}})^2}{M^2} n_j \right]} \right) \\
&= 2 \exp \left(- \frac{2\epsilon^2}{\|\mathbf{c}\|^2 - \frac{1}{M} \left(\sum_j c_j \sqrt{n_j} \right)^2} \right). \tag{7}
\end{aligned}$$

56 The only inequality used is Hoeffding's; the rest is simplification and manipulation.

This improves upon the previous result, due to the smaller denominator. It also follows much more similarly to the asymptotic normality argument, where the variance of S appears in the denominator, up to the variance of $f_{Z_{j,k}}$. Essentially,

this finite-sample valid bound takes the same form as the asymptotic bound, but upper bounds the variance of $f_{Z_{j,k}}$ to be at most $\frac{1}{4}$, as we discuss later. Thus, our final p-value bound is

$$\begin{aligned} \mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq s) &\leq 2 \exp \left(- \frac{2s^2}{(\max_i f_i - \min_i f_i)^2 \left(\|\mathbf{c}\|^2 - \frac{1}{M} \left(\sum_j c_j \sqrt{n_j} \right)^2 \right)} \right) \\ &= 2 \exp \left(- \frac{2s^2}{(\max_i f_i - \min_i f_i)^2 \|\mathbf{c}\|^2 (1 - \gamma)} \right) \end{aligned} \quad [8]$$

57 We state this general proposition below, noting that Proposition 1 follows as an immediate corollary.

Proposition 2. For any fixed $\mathbf{f} \in \mathbb{R}^I$ and $\mathbf{c} \in \mathbb{R}^J$, if $\gamma = 0$ then $S = 0$ with probability 1. If $\gamma < 1$, then

$$\mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq s) \leq 2 \exp \left(- \frac{s^2}{2\|\mathbf{f}\|_\infty^2 \|\mathbf{c}\|_2^2 (1 - \gamma)} \right),$$

58 where $\|\mathbf{f}\|_\infty$ can be tightened to $(\max_i f_i - \min_i f_i)/2$.

59 **S.2.C. Asymptotic distribution of OASIS test statistic.** We begin by restating our asymptotic normality result.

Proposition 3 (Restatement of Prop. 2). Consider any fixed $\mathbf{f} \in \mathbb{R}^I$, $\mathbf{c} \in \mathbb{R}^J$, probability distribution $\mathbf{p} \in \Delta^I$ with $\sigma_{\mathbf{f}}^2 > 0$, and any sequence of column counts $\{\mathbf{n}^{(t)}\}_{t=1}^\infty$ where each $\mathbf{n}^{(t)} \in \mathbb{N}^J$, with $\min_{j \in [J]} n_j^{(t)} \xrightarrow{t \rightarrow \infty} \infty$ and $\gamma(\mathbf{n}^{(t)}, \mathbf{c}) < 1$ for all t . Then, the random sequence of OASIS test statistics $\{S_t\}_{t=1}^\infty$, where $S_t = S(X_t, \mathbf{f}, \mathbf{c})$ and $X_t^{(j)} \sim \text{multinomial}(n_j^{(t)}, \mathbf{p})$ independently across j and t , satisfies

$$\frac{1}{\sqrt{1 - \gamma(\mathbf{n}^{(t)}, \mathbf{c})}} S_t \xrightarrow{D} \mathcal{N}(0, \sigma_{\mathbf{f}}^2 \|\mathbf{c}\|^2).$$

60 We show that OASIS's test statistic is asymptotically normally distributed using the Lyapunov Central Limit Theorem.

Theorem 1 (Theorem 27.3 (7)). Suppose that for each M the sequence X_{M1}, \dots, X_{MM} is independent and satisfies

$$\mathbb{E}[X_{Mk}] = 0, \quad \sigma_{Mk}^2 = \mathbb{E}[X_{Mk}^2], \quad s_M^2 = \sum_{k=1}^M \sigma_{Mk}^2$$

for all M, k , where the means and variances are assumed to be finite, and $s_M^2 > 0$ for large M . Defining $S_M = \sum_{k=1}^M X_{Mk}$, if

$$\lim_M \sum_{k=1}^M \frac{1}{s_M^{2+\delta}} \mathbb{E}[|X_{Mk}|^{2+\delta}] = 0$$

61 holds for some positive δ (Lyapunov's condition), then $S_M/s_M \xrightarrow{D} Z$ where $Z \sim \mathcal{N}(0, 1)$.

Proof. Continuing from Eq. (5), the test-statistic S can be expressed as

$$\begin{aligned} X_{M\tau(j,k)} &\triangleq \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_\ell c_\ell \sqrt{n_\ell}}{M} \right) (f_{Z_{j,k}} - \mu) \\ S(\mathbf{f}, \mathbf{c}) &= \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_\ell c_\ell \sqrt{n_\ell}}{M} \right) (f_{Z_{j,k}} - \mu) = \sum_{\tau=1}^M X_{M\tau}, \end{aligned}$$

62 for an appropriate reindexing τ which maps j, k indices to $[M]$. Since $f_i \in [0, 1]$ for all i , each $X_{M\tau}$ has finite variance, and has
63 mean 0 due to the centering.

Computing the variance of S_M , which is s_M^2 as the terms in S_M have mean 0 and are independent, yields

$$\begin{aligned} s_M^2 &= \sum_{jk} \mathbb{E} \left[\left(\left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_\ell c_\ell \sqrt{n_\ell}}{M} \right) (f_{Z_{j,k}} - \mu) \right)^2 \right] \\ &= \sigma_{\mathbf{f}}^2 \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_\ell c_\ell \sqrt{n_\ell}}{M} \right)^2 \\ &= \sigma_{\mathbf{f}}^2 \left(\|\mathbf{c}\|^2 - \frac{1}{M} \left(\sum_\ell c_\ell \sqrt{n_\ell} \right)^2 \right) \\ &= \sigma_{\mathbf{f}}^2 \|\mathbf{c}\|^2 (1 - \gamma), \end{aligned} \quad [9]$$

64 where we use the definition of $\sigma_{\mathbf{f}}^2 = \text{Var}_{Z \sim \mathcal{P}}(f_Z) = \sum_i p_i f_i^2 - (\sum_i p_i f_i)^2$. This is greater than 0 by assumption in the
 65 Proposition statement, where we also assumed that γ was bounded away from 1. Thus, our desired result is

$$66 \quad \frac{1}{\sqrt{1-\gamma}} S(\mathbf{f}, \mathbf{c}) \xrightarrow{D} \mathcal{N}(0, \|\mathbf{c}\|^2 \sigma_{\mathbf{f}}^2). \quad [10]$$

67 **Satisfying Lyapunov's condition:** we show that this condition holds for $\delta = 1$, which entails bounding third moments. We
 68 assume that $\max_i |f_i - \mu| \leq 1$, which is without loss of generality as \mathbf{f} is fixed and everything can simply be rescaled, and so
 69 $\mathbb{E}[|f_{Z_{j,k}} - \mu|^3] \leq \mathbb{E}[|f_{Z_{j,k}} - \mu|^2] = \sigma_{\mathbf{f}}^2$. By Cauchy-Schwarz, even the largest coefficient is vanishing:

$$70 \quad \left| \frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right| \leq \left| \frac{c_j}{\sqrt{n_j}} \right| + \frac{\|\mathbf{c}\|}{\sqrt{M}} \leq \frac{2\|\mathbf{c}\|}{\min_j \sqrt{n_j}}. \quad [11]$$

Using this, we can upper bound the sum of the third moments:

$$\begin{aligned} \sum_{\tau=1}^M \mathbb{E}[|X_{M\tau}|^3] &= \sum_{jk} \mathbb{E} \left[\left| \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) (f_{Z_{j,k}} - \mu) \right|^3 \right] \\ &= \sum_{jk} \left| \frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right|^3 \mathbb{E}[|f_{Z_{j,k}} - \mu|^3] \\ &\leq \sum_{jk} \left| \frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right|^3 \sigma_{\mathbf{f}}^2 \\ &\leq \left(\frac{2\|\mathbf{c}\|}{\min_j \sqrt{n_j}} \right) \sigma_{\mathbf{f}}^2 \left(\sum_{jk} \left| \frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right|^2 \right). \end{aligned} \quad [12]$$

Verifying Lyapunov's condition:

$$\begin{aligned} \frac{1}{s_M^3} \sum_{\tau=1}^M \mathbb{E}[|X_{M\tau}|^3] &= \frac{\sum_{jk} \mathbb{E} \left[\left| \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) f_{Z_{j,k}} \right|^3 \right]}{\left(\sum_{jk} \mathbb{E} \left[\left| \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right) f_{Z_{j,k}} \right|^2 \right] \right)^{3/2}} \\ &\leq \frac{\left(\frac{2\|\mathbf{c}\|}{\min_j \sqrt{n_j}} \right) \sigma_{\mathbf{f}}^2 \left(\sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right)^2 \right)}{\left(\sigma_{\mathbf{f}}^2 \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right)^2 \right)^{3/2}} \\ &= \frac{2\|\mathbf{c}\|}{\min_j \sqrt{n_j}} \left(\sigma_{\mathbf{f}}^2 \sum_{jk} \left(\frac{c_j}{\sqrt{n_j}} - \frac{\sum_{\ell} c_{\ell} \sqrt{n_{\ell}}}{M} \right)^2 \right)^{-1/2} \\ &= \frac{2\|\mathbf{c}\|}{\sigma_{\mathbf{f}} \min_j \sqrt{n_j}} (\|\mathbf{c}\|^2 (1-\gamma))^{-1/2} \\ &= \frac{2}{\sigma_{\mathbf{f}} \sqrt{(1-\gamma) \min_j n_j}} \end{aligned} \quad [13]$$

71 Since the original quantity is nonnegative, this last line goes to 0 as all $n_j \rightarrow \infty$ simultaneously as long as $\sigma_{\mathbf{f}} > 0$ and $\gamma < 1$,
 72 which were assumed in the Theorem statement, and thus by the Lyapunov CLT we have the desired result. The variance of
 73 $S(\mathbf{f}, \mathbf{c})$ has already been computed as s_M^2 , and so

$$74 \quad \frac{1}{\sqrt{1-\gamma}} S(\mathbf{f}, \mathbf{c}) \xrightarrow{D} \mathcal{N}(0, \|\mathbf{c}\|^2 \sigma_{\mathbf{f}}^2). \quad [14]$$

75 Note that if $\sigma_{\mathbf{f}} = 0$ then $S = 0$ w.p. 1, which matches this result under the convention that a Gaussian with variance 0 is a
 76 Dirac point mass at 0. \square

77 With this proposition, we can construct an asymptotically valid p-value.

Corollary 3.1 (Restatement of Corollary). *Under the conditions of Proposition 3,*

$$2\Phi\left(-\frac{|S|}{\hat{\sigma}_{\mathbf{f}}\|\mathbf{c}\|\sqrt{1-\gamma}}\right),$$

78 *is an asymptotically valid p-value.*

79 For comparison with our finite-sample p-value bound, we construct an asymptotically valid p-value bound using standard
80 Gaussian tail bounds as

$$81 \quad \mathbb{P}(|S(\mathbf{f}, \mathbf{c})| \geq s) = 2\Phi\left(-\frac{s}{\hat{\sigma}_{\mathbf{f}}\|\mathbf{c}\|\sqrt{1-\gamma}}\right) \leq 2\exp\left(-\frac{s^2}{2\hat{\sigma}_{\mathbf{f}}^2\|\mathbf{c}\|^2(1-\gamma)}\right), \quad [15]$$

82 where Φ is the CDF of a standard Gaussian random variable. This holds by Slutsky's theorem, as $\hat{\sigma}_{\mathbf{f}}^2$ is an asymptotically
83 consistent estimator of $\sigma_{\mathbf{f}} > 0$. Note that this form precisely matches that of the finite-sample p-value bound in Proposition 1,
84 up to bounding $\hat{\sigma}_{\mathbf{f}}^2 \leq \frac{1}{4}$ (bounding the variance of a $[0, 1]$ random variable as $1/4$). Since $f_{Z_j, k} \in [0, 1]$, the rate of convergence
85 can be quantified with the Berry-Esseen inequality (8) if desired.

86 **S.2.D. Effect Size.** The effect size measure for OASIS is motivated by a simple two group alternative, where each sample originate
87 from either group A or group B . Each group has a characteristic target distribution, \mathbf{p}_A and \mathbf{p}_B respectively, where observations
88 from a sample in group A (resp. B) are drawn i.i.d. from \mathbf{p}_A (resp. \mathbf{p}_B). Then, denoting Z as the identity of the random
89 row drawn from \mathbf{p}_A or \mathbf{p}_B , the effect size estimate $\hat{\Delta}$ is nothing but the plug-in estimate of $\Delta = |\mathbb{E}_{Z \sim \mathbf{p}_A}[f_Z] - \mathbb{E}_{Z \sim \mathbf{p}_B}[f_Z]|$
90 the difference between $\mathbb{E}[f_Z]$ where $Z \sim \mathbf{p}_A$ or $Z \sim \mathbf{p}_B$. This is an intuitive measure of discrepancy, as OASIS declaring
91 a discovery based on a specific \mathbf{f} and \mathbf{c} indicates that when partitioned into groups by \mathbf{c} , the two groups of samples have
92 significantly different row distributions as measured by \mathbf{f} . This quantity Δ is nonnegative, and is bounded above by the total
93 variation distance between the empirical row distributions of columns where $c_j > 0$ and where $c_j < 0$. This is attained by
94 taking the supremum over all $\mathbf{f} \in [0, 1]^I$, yielding a value of 0 if and only if the empirical row distributions are the same
95 between the two clusters, and 1 if and only if they are disjoint. Note that due to the rescaling of each point's contribution by
96 $\sqrt{n_j}$ (downweighting) in OASIS's test statistic, for a given \mathbf{c} constructing \mathbf{f} to minimize the p-value bound will not necessarily
97 maximize the effect size.

98 In the main text we stated that $0 \leq \hat{\Delta} \leq \delta_{\text{TV}}(\hat{\mathbf{p}}_+, \hat{\mathbf{p}}_-)$, which we prove below.

Proof of Effect Size Bound. The lower bound on $\hat{\Delta}$ follows trivially from the absolute value. The upper bound follows from
the reformulation of $\hat{\Delta}$ as:

$$\hat{\Delta} = \left| \sum_i (\hat{p}_{+,i} - \hat{p}_{-,i}) f_i \right| \leq \max_{g \in [0,1]^I} |(\hat{\mathbf{p}}_+ - \hat{\mathbf{p}}_-)^\top g| = \frac{1}{2} \|\hat{\mathbf{p}}_+ - \hat{\mathbf{p}}_-\|_1 = \delta_{\text{TV}}(\hat{\mathbf{p}}_+, \hat{\mathbf{p}}_-)$$

We can rewrite this as an ℓ_1 norm as since $\hat{\mathbf{p}}_+$ and $\hat{\mathbf{p}}_-$ are probability distributions, shifting g by a constant does not change
the objective value. Thus, we can equivalently maximize over $g \in [-1/2, 1/2]^I$. Then, for any vector x ,

$$\max_{g: \|g\|_\infty \leq 1/2} g^\top x = \frac{1}{2} \|x\|_1,$$

99 by the dual norm characterization of the ℓ_1 norm. □

Considering general \mathbf{c} beyond the simple binary case, there is no clear extension of this effect size measure. We discuss some
candidate measures and their drawbacks below. For simplicity, in this section we consider $\mathbf{f} \in [-1, 1]^I$ instead of $\mathbf{f} \in [0, 1]^I$, to
avoid centering. One candidate measure which encourages binary \mathbf{f} is:

$$\hat{\Delta} = \left| \frac{\mathbf{f}^\top X \mathbf{c}}{\mathbf{1} X |\mathbf{c}|} \right|.$$

A more gradual and natural alternative which allows for non-binarized \mathbf{f} is:

$$\hat{\Delta} = \left| \frac{\mathbf{f}^\top X \mathbf{c}}{|\mathbf{f}^\top X \mathbf{c}|} \right|.$$

100 These quantities are both bounded between 0 and 1. The first attains a value of 1 only if $f_i X_{i,j} c_j = X_{i,j} |c_j|$ (extending the
101 range of \mathbf{f} to $[-1, 1]$). The second requires that $f_i X_{i,j} c_j = |f_i X_{i,j} c_j|$. However, note that these can both be trivially achieved
102 by taking \mathbf{f} and \mathbf{c} to both be all ones vectors. This will not yield a significant p-value, however, due to the structure of \tilde{X} .

S.2.E. Statistical validity of data-splitting. The p-value bound outputted by the data-splitting optimization procedure is statistically valid, by a classical data-splitting argument (used in e.g. (9)), which we detail here for completeness. We begin by defining the random table generated by the null for probability vector \mathbf{p} with column counts $\{n_j\}_{j=1}^J$. The data generated by the null is $X \sim \text{nullTable}(\{n_j\}_{j=1}^J, \mathbf{p})$. Data-splitting entails selecting $\{n_j^{\text{Train}}\}$ as a function of the given column counts $\{n_j\}_{j=1}^J$, where $n_j^{\text{Train}} \leq n_j$ for all j . In this procedure we select uniformly at random n_j^{Train} counts uniformly at random from the n_j total counts in that column and assign them to X_{Train} , with the rest being assigned to X_{Test} . Then, $X_{\text{Train}} \sim \text{nullTable}(\{n_j\}_{j=1}^J, \mathbf{p})$, $X_{\text{Test}} \sim \text{nullTable}(\{n_j - n_j^{\text{Train}}\}_{j=1}^J, \mathbf{p})$, and the two are independent. This is due to the fact that each count in the multinomial $X^{(j)}$ is an i.i.d. draw from \mathbf{p} , and so splitting the counts randomly between $X_{\text{Test}}^{(j)}$ and $X_{\text{Train}}^{(j)}$ is simply separating these i.i.d. draws into two groups, resulting in a multinomial for each with the appropriate counts.

As shown by Proposition 1, for any fixed \mathbf{f}, \mathbf{c} , the random variable given by the p-value bound $p(X_{\text{Test}}, \mathbf{c}, \mathbf{f})$ (where only X_{Test} is random) is a statistically valid p-value bound. We show that even when $\mathbf{c} = \mathbf{c}(X_{\text{Train}})$ and $\mathbf{f} = \mathbf{f}(X_{\text{Train}})$ are functions of the training data, $p(X) = p(X_{\text{Test}}, \mathbf{c}(X_{\text{Train}}), \mathbf{f}(X_{\text{Train}}))$ is a valid p-value bound.

$$\begin{aligned}
\mathbb{P}(p(X) \leq u) &\stackrel{(a)}{=} \mathbb{P}(p(X_{\text{Test}}, \mathbf{c}(X_{\text{Train}}), \mathbf{f}(X_{\text{Train}})) \leq u) \\
&\stackrel{(b)}{=} \sum_{x_{\text{Train}}} \mathbb{P}(p(X_{\text{Test}}, \mathbf{c}(x_{\text{Train}}), \mathbf{f}(x_{\text{Train}})) \leq u \mid X_{\text{Train}} = x_{\text{Train}}) \mathbb{P}(X_{\text{Train}} = x_{\text{Train}}) \\
&\stackrel{(c)}{=} \sum_{x_{\text{Train}}} \mathbb{P}(p(X_{\text{Test}}, \mathbf{c}(x_{\text{Train}}), \mathbf{f}(x_{\text{Train}})) \leq u) \mathbb{P}(X_{\text{Train}} = x_{\text{Train}}) \\
&\stackrel{(d)}{\leq} \sum_{x_{\text{Train}}} u \mathbb{P}(X_{\text{Train}} = x_{\text{Train}}) \\
&= u
\end{aligned} \tag{16}$$

where (a) holds by definition, (b) by law of total probability, (c) by the fact that X_{Test} is independent of X_{Train} , and (d) as once we condition on X_{Train} , \mathbf{c} and \mathbf{f} are no longer random, and so the p-value bound from the theorem statement can be used. Finally, the last line follows from summing over the probabilities. The resulting p-value bound is not just random with respect to the data X (which is to be expected), but also with respect to the random splitting procedure. However, since the p-value bound holds for any fixed \mathbf{f}, \mathbf{c} , once we condition on X_{Train} , the p-value bound can be applied.

Note that more generally, there are two sources of randomness used in the algorithm; one from splitting the data into train and test sets, and one from the generation of $\mathbf{c}(X_{\text{Train}})$ and $\mathbf{f}(X_{\text{Train}})$. The first source (splitting into train and test sets) is fundamental, whereas the second is simply for computational efficiency; if computational complexity is not an issue, we could enumerate over all possible $\mathbf{f} \in \{0, 1\}^I$ and exactly solve the inner optimization problem deterministically. However, the randomness in data-splitting is necessary, and can greatly impact performance. If the training data is not representative of the test data, then even if the optimal \mathbf{f} and \mathbf{c} for the training data are identified, it need not yield a significant p-value bound on the test data. Thus, a natural approach is to generate multiple random splits of the data, perform Bonferroni correction over these multiple splits, and take the minimum as a valid p-value bound. To this end, OASIS generates independent random splits of the data by using the Poisson distribution of the counts under the null (assuming the column counts are Poisson). The Poisson nature of these counts is critically important, as under negative binomial overdispersion the counts observed are no longer independent, and so more sophisticated methods are needed, which do not yield fully independent sets of counts unless the overdispersion parameter is known (10).

S.3. Optimization procedure

Recall the p-value bound derived in Proposition 1 for the test statistic $S = \mathbf{f}^\top \tilde{X} \mathbf{c}$, restated below for convenience, where we constrain $0 \leq \mathbf{f} \leq 1$:

$$2 \exp\left(-\frac{2S^2}{\|\mathbf{c}\|^2(1-\gamma)}\right). \tag{17}$$

For simplicity, we shift and rescale \mathbf{f} to be bounded as $|\mathbf{f}| \leq 1$, which does not change the structure of the problem: an optimal solution for the original problem can be obtained by modifying the solution to the shifted and rescaled problem.

S.3.A. Reformulating the optimization problem (proof of Lemma 1). Recalling the containment lemma in the main text:

Lemma 2 (Restatement of Lemma 1). *The set of optimal solutions for the p-value bound can be expressed as*

$$\left(\underset{0 \leq \mathbf{f} \leq 1, \|\mathbf{c}\|_2 \leq 1}{\text{argmax}} \mathbf{f}^\top \tilde{X} \mathbf{c}\right) \subseteq \underset{0 \leq \mathbf{f} \leq 1, \|\mathbf{c}\|_2 \leq 1}{\text{argmin}} 2 \exp\left(-\frac{2(\mathbf{f}^\top \tilde{X} \mathbf{c})^2}{1 - \frac{1}{M} \langle \mathbf{c}, \sqrt{\mathbf{n}} \rangle^2}\right)$$

We prove this lemma by decomposing \mathbf{c} into the part parallel to $\sqrt{\mathbf{n}} \triangleq \sqrt{X^\top \mathbf{1}}$ and the part orthogonal to $\sqrt{\mathbf{n}}$, i.e. $\mathbf{c} = \mathbf{c}^{(1)} + \alpha \sqrt{\mathbf{n}}$ where $\mathbf{c}^{(1)} \perp \sqrt{\mathbf{n}}$ and $\alpha \in \mathbb{R}$. This yields an objective value of

$$\begin{aligned}
\frac{(\mathbf{f}^\top \tilde{X} \mathbf{c})^2}{\|\mathbf{c}\|^2 - \frac{1}{M} (\mathbf{c}^\top \sqrt{\mathbf{n}})^2} &= \frac{(\mathbf{f}^\top (\tilde{X} (\mathbf{c}^{(1)} + \alpha \sqrt{\mathbf{n}})))^2}{\|\mathbf{c}^{(1)} + \alpha \sqrt{\mathbf{n}}\|^2 - \frac{1}{M} (\alpha \sqrt{\mathbf{n}}^\top \sqrt{\mathbf{n}})^2} \\
&= \frac{(\mathbf{f}^\top (\tilde{X} \mathbf{c}^{(1)} + \alpha \tilde{X} \sqrt{\mathbf{n}}))^2}{\|\mathbf{c}^{(1)}\|^2 + \alpha^2 \|\sqrt{\mathbf{n}}\|^2 - \frac{1}{M} (\alpha \|\sqrt{\mathbf{n}}\|^2)^2} \\
&= \frac{(\mathbf{f}^\top \tilde{X} \mathbf{c}^{(1)})^2}{\|\mathbf{c}^{(1)}\|^2 + \alpha^2 M - \frac{1}{M} (\alpha M)^2} \\
&= \frac{(\mathbf{f}^\top \tilde{X} \mathbf{c}^{(1)})^2}{\|\mathbf{c}^{(1)}\|^2} \tag{18}
\end{aligned}$$

As

$$\tilde{X} \sqrt{\mathbf{n}} = \left(X - \frac{1}{\mathbf{1}^\top X \mathbf{1}} X \mathbf{1} \mathbf{1}^\top X \right) \text{diag}(1/\sqrt{\mathbf{n}}) \sqrt{\mathbf{n}} = \left(X - \frac{1}{\mathbf{1}^\top X \mathbf{1}} X \mathbf{1} \mathbf{1}^\top X \right) \mathbf{1} = 0$$

Thus, any component of \mathbf{c} in the direction of $\sqrt{\mathbf{n}}$ does not contribute to the numerator, and cancels out in the denominator, and so an optimal solution can be obtained by simply optimizing over \mathbf{c} orthogonal to $\sqrt{\mathbf{n}}$. Since the denominator is equal to $\|\mathbf{c}\|^2$, and the numerator and denominator both scale quadratically in $\|\mathbf{c}\|$, we can simply optimize the numerator subject to $\|\mathbf{c}\| \leq 1$ to identify a pair of optimal \mathbf{f}, \mathbf{c} , as:

$$\underset{\mathbf{c} \in \mathbb{R}^J, \mathbf{f} \in [0,1]^I}{\text{argmax}} \frac{(\mathbf{f}^\top \tilde{X} \mathbf{c})^2}{\|\mathbf{c}\|^2 - \frac{1}{M} (\mathbf{c}^\top \sqrt{\mathbf{n}})^2} \supseteq \underset{\mathbf{f} \in [0,1]^I, \|\mathbf{c}\|_2 \leq 1}{\text{argmax}} (\mathbf{f}^\top \tilde{X} \mathbf{c})^2 \tag{19}$$

$$\supseteq \underset{\mathbf{f} \in [0,1]^I, \|\mathbf{c}\|_2 \leq 1}{\text{argmax}} \mathbf{f}^\top \tilde{X} \mathbf{c}. \tag{20}$$

138 where the last line follows as $\|\mathbf{c}\| \leq 1 \implies \|\mathbf{c} - \sqrt{\mathbf{n}}\| \leq 1$. This proves Lemma 1.

139 To optimize Eq. (20), we begin by enlarging the constraint set of \mathbf{f} to be $[-1, 1]^I$ for symmetry (an optimal \mathbf{f} for the
140 original problem can be obtained from such a $\tilde{\mathbf{f}}$ by considering $\mathbf{f} = (1 + \tilde{\mathbf{f}})/2$ and $\mathbf{f} = (1 - \tilde{\mathbf{f}})/2$, optimizing $\mathbf{c} \propto \tilde{X}^\top \mathbf{f}$, and
141 comparing their objective values). For fixed \mathbf{c} , \mathbf{f} is then optimized as $\mathbf{f} = \text{sign}(\tilde{X} \mathbf{c})$, where $\text{sign}(x)$ is +1 if $x > 0$, 0 if $x = 0$,
142 and -1 if $x < 0$. \mathbf{c} is always optimized as $\mathbf{c} \propto \tilde{X}^\top \mathbf{f}$. Rolling these updates together, since the sign of an entry doesn't change
143 based on scaling we can write our iterates $\{\mathbf{f}^{(t)}\}$ as

$$\mathbf{f}^{(t+1)} = \text{sign}(\tilde{X} \tilde{X}^\top \mathbf{f}^{(t)}), \tag{21}$$

145 where the iterates $\{\mathbf{c}^{(t)}\}$ are implicitly computed as $\mathbf{c}^{(t)} \propto \tilde{X}^\top \mathbf{f}^{(t)}$.

146 **S.3.B. Iterative testing of contingency tables.** Recalling the initial intuition for the OASIS test, where there are two groups
147 of samples each with a characteristic probability distribution over the rows, the optimal \mathbf{c} vector in this case partitions the
148 samples by type (positive on one group of samples and negative on the other). However, these scalar-valued vectors \mathbf{c} are
149 limited by their one dimensionality, and so a natural question is whether OASIS can identify more than 2 clusters within a
150 table. In genomics contexts this translates to detecting subclusters or multiple cell-types.

151 Considering the matrix viewpoint of OASIS, the test statistic can be thought of as determining whether the centered and
152 right normalized contingency table has a single sufficiently large direction of deviation. We show that a simple extension of
153 OASIS allows for detection of how many ways the contingency table significantly deviates from the null, through a statistical
154 stopping condition.

155 Consider that one initial partitioning vector $\mathbf{c}^{(1)}$ is identified, for example distinguishing Delta vs Omicron samples in the
156 running SARS-CoV-2 example, and we want to determine if there are additional statistically significant ways of partitioning
157 this data not along this same direction $\mathbf{c}^{(1)}$. One way of accomplishing this is by attempting to identify new vectors \mathbf{c}, \mathbf{f} which
158 yield a significant p-value bound, requiring that this new vector \mathbf{c} be orthogonal to the previously found vector $\mathbf{c}^{(1)}$. With
159 SARS-CoV-2, this could be a vector separating Omicron BA.1 and BA.2. The constraint $\mathbf{c} \perp \mathbf{c}^{(1)}$ maintains the convexity of the
160 optimization problem. This naturally generalizes to multiple orthogonality constraints where \mathbf{c} is restricted to be orthogonal to
161 a subspace. As discussed, OASIS's initial test statistic can be thought of as identifying one direction of substantial signal.
162 Restricting this new \mathbf{c} to be orthogonal to the first one is analogous to identifying the second largest singular value and
163 its corresponding right singular vector in the case of the SVD. Thus, this iterative procedure of finding \mathbf{c} orthogonal to all
164 previously found \mathbf{c} and terminating when this no longer yields a significant p-value bound parallels analyzing the spectrum of
165 \tilde{X} until a significant drop-off occurs. In this iterative setting, the optimizing \mathbf{c} is constructed as

$$\mathbf{c}^* = \underset{\substack{\|\mathbf{c}\|_2 \leq 1, \\ \mathbf{c} \perp \mathbf{c}^{(1)}}}{\text{argmax}} \mathbf{f}^\top \tilde{X} \mathbf{c} \implies \mathbf{c}^* \propto \tilde{X}^\top \mathbf{f} - \frac{\mathbf{f}^\top \tilde{X} \mathbf{c}^{(1)}}{\|\mathbf{c}^{(1)}\|_2^2} \mathbf{c}^{(1)}. \tag{22}$$

167 \mathbf{f} is constructed as before, with iterates constructed as $\mathbf{f} = \text{sign}(\tilde{X} \mathbf{c})$.

Algorithm 1 OASIS-iter

```

1: Input: Contingency table  $X$ 
2: # Split data randomly into train and test portions
3:  $X_{\text{train}}, X_{\text{test}} \leftarrow \text{splitDataset}(X)$ 
4: for  $i = 1, 2, \dots$  do
5:   # Run OASIS-opt with constraint:  $\mathbf{c}^{(i)} \perp \mathbf{c}^{(j)} \forall j < i$ 
6:    $\mathbf{c}^{(i)}, \mathbf{f}^{(i)} \leftarrow \text{OASIS-perp}(X_{\text{train}}, \{\mathbf{c}^{(j)}\}_{j=1}^{i-1})$ 
7:   # Compute OASIS p-value bound
8:    $p^{(i)} \leftarrow \text{pvBound}(X_{\text{test}}, \mathbf{f}^{(i)}, \mathbf{c}^{(i)})$ 
9:   if  $p^{(i)} > 0.05$  then
10:     Break
11: return  $[\mathbf{c}^{(1)}, \mathbf{c}^{(2)}, \dots], [\mathbf{f}^{(1)}, \mathbf{f}^{(2)}, \dots], [p^{(1)}, p^{(2)}, \dots]$ 

```

168 **S.3.C. Asymptotic p-value.** The asymptotic p-value provided in Corollary 3.1 leads to the following optimization objective:

$$169 \quad \operatorname{argmax}_{\mathbf{f}, \mathbf{c}} \frac{(\mathbf{f}^\top \tilde{X} \mathbf{c})^2}{\|\mathbf{c}\|^2 \hat{\sigma}_{\mathbf{f}}^2}. \quad [23]$$

170 Here, $\hat{\sigma}_{\mathbf{f}}^2 = \sum_i \hat{p}_i f_i^2 - (\sum_i \hat{p}_i f_i)^2 = \sum_i \hat{p}_i (f_i - \hat{\mathbf{p}}^\top \mathbf{f})^2 = \mathbf{f}^\top (\operatorname{diag}(\hat{\mathbf{p}}) - \hat{\mathbf{p}} \hat{\mathbf{p}}^\top) \mathbf{f}$ and $\hat{\mathbf{p}} \triangleq \frac{1}{M} X \mathbf{1}$. Note that $\hat{\mathbf{p}} > 0$, as otherwise
 171 these rows would not exist in X . For a fixed \mathbf{f} , the optimal \mathbf{c} is easily computed as $\mathbf{c}^* \propto \tilde{X}^\top \mathbf{f}^*$, by Cauchy-Schwarz. This
 172 allows us to reduce the optimization problem to that of optimizing over \mathbf{f} by dual norm characterizations as

$$173 \quad \max_{\mathbf{f}, \mathbf{c}} \frac{(\mathbf{f}^\top \tilde{X} \mathbf{c})^2}{\|\mathbf{c}\|^2 \hat{\sigma}_{\mathbf{f}}^2} = \max_{\mathbf{f}} \frac{\|\tilde{X}^\top \mathbf{f}\|^2}{\hat{\sigma}_{\mathbf{f}}^2}. \quad [24]$$

174 We show that an optimal \mathbf{f} for this objective can be computed efficiently as an eigenvector problem. Defining $v_{\max}(A)$
 175 the principal eigenvector of a symmetric matrix A (selecting an arbitrary one if the maximum eigenvalue has multiplicity
 176 greater than 1), we state the following Proposition.

Proposition 4. An optimal \mathbf{f}^* for Eq. (24) can be constructed as

$$\mathbf{f}^* = \operatorname{diag}(\hat{\mathbf{p}})^{-1/2} v_{\max}(\operatorname{diag}(\hat{\mathbf{p}})^{-1/2} \tilde{X} \tilde{X}^\top \operatorname{diag}(\hat{\mathbf{p}})^{-1/2}).$$

Proof. Observe that the matrix used in the computation of $\hat{\sigma}_{\mathbf{f}}^2$ is positive semidefinite, and that restricting $\mathbf{f} \perp \hat{\mathbf{p}}$ the resulting operator is positive definite:

$$\begin{aligned} \operatorname{diag}(\hat{\mathbf{p}}) - \hat{\mathbf{p}} \hat{\mathbf{p}}^\top &\succ 0 \\ \operatorname{diag}(\hat{\mathbf{p}}) &\succ 0 \end{aligned}$$

177 To identify an optimal \mathbf{f} for Eq. (24), we show that the objective is scale and shift invariant with respect to \mathbf{f} . Formally,
 178 defining the maximization objective of Eq. (24) as

$$179 \quad h(\mathbf{f}) \triangleq \frac{\|\tilde{X}^\top \mathbf{f}\|^2}{\hat{\sigma}_{\mathbf{f}}^2}, \quad [25]$$

we claim that for any $\beta \in \mathbb{R}$, $h(\mathbf{f} + \beta \mathbf{1}) = h(\mathbf{f})$. It is clear that

$$\hat{\sigma}_{\mathbf{f} + \beta \mathbf{1}}^2 = \sum_i \hat{p}_i (f_i + \beta - \hat{\mathbf{p}}^\top (\mathbf{f} + \beta \mathbf{1}))^2 = \sum_i \hat{p}_i (f_i - \hat{\mathbf{p}}^\top \mathbf{f})^2 = \hat{\sigma}_{\mathbf{f}}^2, \quad [26]$$

and since $\mathbf{1}^\top X \mathbf{1} = M$,

$$\begin{aligned} (\mathbf{f} + \beta \mathbf{1})^\top \tilde{X} &= \mathbf{f}^\top \tilde{X} + \beta \mathbf{1}^\top \left(X - \frac{1}{M} X \mathbf{1} \mathbf{1}^\top X \right) \operatorname{diag}(1/\sqrt{X^\top \mathbf{1}}) \\ &= \mathbf{f}^\top \tilde{X} + \beta \left(\mathbf{1}^\top X - \frac{\mathbf{1}^\top X \mathbf{1}}{M} \mathbf{1}^\top X \right) \operatorname{diag}(1/\sqrt{X^\top \mathbf{1}}) \\ &= \mathbf{f}^\top \tilde{X}. \end{aligned} \quad [27]$$

180 Combining eqs. (26) and (27) together yields that

$$181 \quad h(\mathbf{f} + \beta \mathbf{1}) = h(\mathbf{f}). \quad [28]$$

Additionally, the maximization objective is scale invariant. For any scalar $\kappa \neq 0$:

$$\begin{aligned} h(\kappa \mathbf{f}) &= \max_{\mathbf{f}} \frac{\|\tilde{X}^\top(\kappa \mathbf{f})\|^2}{(\kappa \mathbf{f})^\top (\text{diag}(\hat{\mathbf{p}}) - \hat{\mathbf{p}}\hat{\mathbf{p}}^\top) (\kappa \mathbf{f})} \\ &= \max_{\mathbf{f}} \frac{\kappa^2 \|\tilde{X}^\top \mathbf{f}\|^2}{\kappa^2 \mathbf{f}^\top (\text{diag}(\hat{\mathbf{p}}) - \hat{\mathbf{p}}\hat{\mathbf{p}}^\top) \mathbf{f}} \\ &= h(\mathbf{f}) \end{aligned} \tag{29}$$

Since the numerator and denominator both scale quadratically in the magnitude of \mathbf{f} , the denominator ($\hat{\sigma}_{\mathbf{f}}^2$) can instead be treated as a constraint, i.e.

$$\max_{\mathbf{f}} \frac{\|\tilde{X}^\top \mathbf{f}\|^2}{\hat{\sigma}_{\mathbf{f}}^2} = \max_{\mathbf{f}: \hat{\sigma}_{\mathbf{f}}^2 \leq 1} \|\tilde{X}^\top \mathbf{f}\|^2 \tag{30}$$

$$= \max_{\mathbf{f}: \mathbf{f}^\top \text{diag}(\hat{\mathbf{p}}) \mathbf{f} \leq 1, \mathbf{f}^\top \hat{\mathbf{p}} = 0} \mathbf{f}^\top \tilde{X} \tilde{X}^\top \mathbf{f} \tag{31}$$

182 The right hand side of Eq. (30) is clearly at most the left hand side, and an optimal \mathbf{f} for the left hand side can be rescaled
 183 so that $\hat{\sigma}_{\mathbf{f}}^2 = 1$, yielding a feasible \mathbf{f} for the right hand side with the same objective value. Eq. (31) follows from the shift
 184 invariance in Equations (26) and (27).

We now show that this can be reformulated as an eigenvector problem. Defining $D = \text{diag}(\hat{\mathbf{p}})$ and $\mathbf{y} = D^{-1/2} \mathbf{f}$:

$$\begin{aligned} \max_{\substack{\mathbf{f}: \mathbf{f}^\top \text{diag}(\hat{\mathbf{p}}) \mathbf{f} \leq 1 \\ \mathbf{f}^\top \hat{\mathbf{p}} = 0}} \mathbf{f}^\top \tilde{X} \tilde{X}^\top \mathbf{f} &= \max_{\substack{\mathbf{y}: \|\mathbf{y}\|^2 \leq 1 \\ \mathbf{y}^\top D^{-1/2} \hat{\mathbf{p}} = 0}} \mathbf{y}^\top D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2} \mathbf{y} \\ &\leq \max_{\mathbf{y}: \|\mathbf{y}\|^2 \leq 1} \mathbf{y}^\top D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2} \mathbf{y} \\ &= \|D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}\|. \end{aligned} \tag{32}$$

The upper bound is due to removing the orthogonality constraint. The final equality is attained by taking \mathbf{y} as the principal eigenvector of the given matrix scaled to have unit norm, with objective value equal to the largest eigenvalue λ^* as below:

$$\mathbf{y}^* = v_{\max} (D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}), \tag{33}$$

$$\lambda^* = \lambda_{\max} (D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}). \tag{34}$$

This upper bound in Eq. (32) is attained with equality, as for this optimizing \mathbf{y}^* , $\mathbf{y}^{*\top} \hat{\mathbf{p}} = 0$:

$$\begin{aligned} \hat{\mathbf{p}}^\top D^{-1/2} \mathbf{y}^* &= \hat{\mathbf{p}}^\top D^{-1/2} \left(\frac{1}{\lambda^*} D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2} \mathbf{y}^* \right) \\ &= \frac{1}{\lambda^*} \hat{\mathbf{p}}^\top D^{-1} \tilde{X} \tilde{X}^\top D^{-1/2} \mathbf{y}^* \\ &= \frac{1}{\lambda^*} \mathbf{1}^\top \tilde{X} \tilde{X}^\top D^{-1/2} \mathbf{y}^* \\ &= 0. \end{aligned} \tag{35}$$

185 Thus, since the unconstrained optimum lies within the constraint set, it is also the constrained optimum. An optimizing \mathbf{f} for
 186 Eq. (30) can then be constructed as

$$\mathbf{f}^* = \text{diag}(\hat{\mathbf{p}})^{-1/2} v_{\max} (\text{diag}(\hat{\mathbf{p}})^{-1/2} \tilde{X} \tilde{X}^\top \text{diag}(\hat{\mathbf{p}})^{-1/2}). \tag{36}$$

188 Since the constructed \mathbf{y} has $\|\mathbf{y}\| = 1$, $\hat{\sigma}_{\mathbf{f}}^2 = 1$. Thus, this \mathbf{f}^* is also optimal for the original p-value objective in Eq. (23). \square

189 To construct an optimal solution for Eq. (23), we observe that by Cauchy-Schwarz for any fixed \mathbf{f} the maximum is attained
 190 when $\mathbf{c} \propto \tilde{X}^\top \mathbf{f}$. Thus, using the computed \mathbf{f}^* above and taking \mathbf{c}^* as

$$\mathbf{c}^* = \tilde{X}^\top \mathbf{f}^* \tag{37}$$

192 yields an optimal solution for Eq. (23).

193 A previous version of this work stated that the optimal \mathbf{f} could be computed as a principal right eigenvector of the matrix
 194 $D^{-1} \tilde{X} \tilde{X}^\top$. This is equivalent (assuming without loss of generality that the principal eigenvalue has multiplicity 1):

$$\begin{aligned}
\mathbf{f} &= D^{-1} \tilde{X} \tilde{X}^\top \mathbf{f} \\
\mathbf{f} &= D^{-1/2} (D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}) D^{1/2} \mathbf{f} \\
D^{1/2} \mathbf{f} &= (D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}) D^{1/2} \mathbf{f} \\
D^{1/2} \mathbf{f} &= v_{\max} (D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}) \\
\mathbf{f} &= D^{-1/2} v_{\max} (D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2})
\end{aligned}$$

195 as desired. We provide a method for modifying this vector to improve its finite sample p-value bound performance in
196 Section S.3.C.1, and provide computational details on optimizing \mathbf{f} in Section S.3.D.1.

197 **S.3.C.1. Modification for finite-sample bound.** The \mathbf{f} optimizing the asymptotic p-value, derived in the previous section, will in
198 general be suboptimal for the finite sample p-value bound. However, if the objective is to optimize the finite-sample bound, the
199 asymptotically optimal \mathbf{f}^* can be used as an initialization. One simple approach would be to binarize the vector as $f'_i = 1$
200 if $f_i^* \geq 0$, and 0 otherwise. For this \mathbf{f}' , the optimal \mathbf{c} can be identified as $\mathbf{c} \propto \tilde{X}^\top \mathbf{f}'$. Introducing randomness, \mathbf{f}^* can be
201 randomly rounded to $0, 1^I$ with probabilities proportional to the corresponding entry of \mathbf{f} (after normalizing \mathbf{f}^* vector to
202 $[0, 1]^I$), allowing for multiple trials and a better final objective value.

203 **S.3.D. More computationally intensive optimization methods.** Analyzing the reformulated finite-sample optimization problem,
204 observe that the form of the finite-sample p-value bound objective is simply a quadratic program:

$$205 \quad \max_{\mathbf{f} \in [0, 1]^I} \mathbf{f}^\top \tilde{X} \tilde{X}^\top \mathbf{f}. \quad [38]$$

206 Since $\tilde{X} \tilde{X}^\top \succeq 0$ is positive semi-definite, this is an integer constrained quadratic program. This is a known combinatorial
207 optimization problem, and specialized algorithms for optimizing it have been devised. One well-known approach is based on
208 semidefinite programming (SDP) relaxations, relaxing the vector \mathbf{f} into an $I \times I$ matrix optimization variable A and optimizing

$$209 \quad \max_{A \succeq 0, A_{ii} = 1 \forall i} \text{Tr}(\tilde{X} \tilde{X}^\top A), \quad [39]$$

210 where Tr denotes the matrix Trace operator. The objective value of Eq. (39) upper bounds that of Eq. (38). We can convert
211 the solution to Eq. (39) to an integer solution via Goemans-Williamson rounding (11), which picks a random vector \mathbf{v} on the
212 unit sphere in \mathbb{R}^I , and assigns $\mathbf{f}_i = 1$ if $A_i \mathbf{v} \geq 0$ and $\mathbf{f}_i = -1$ otherwise. This attains an expected approximation ratio of at
213 least $2/\pi$, for Eq. (38). Multiple random \mathbf{v} can be selected to slightly improve the resulting solution.

214 In practice however, solving an SDP is computationally intensive; in the biological setting of interest, contingency tables have
215 many rows due to sequencing error in observing the targets. This makes methods whose sample complexities scale superlinearly
216 in the number of rows too computationally intensive to use. Future work on computationally efficient methods with provable
217 approximation guarantees can be readily utilized in this optimization framework; the primary contribution of OASIS is in its
218 formulation of minimizing the p-value as an optimization problem, for which we propose an alternating maximization-based
219 solution which empirically yields good performance in an efficient manner.

220 **S.3.D.1. Empirical performance.** We study the performance of these different optimization methods by simulating the “strong
221 signal” setting described in Section S.6.A. The tables generated have 12 rows and 10 columns, and a clear two group structure,
222 with signal focused in the first two rows of the tables. When run on 1000 random tables, the aforementioned SDP relaxation
223 followed by rounding attained OPT 88% of the time. Alternating maximization yields similarly good performance, attaining
224 OPT over 82% of the time. Computing the asymptotically optimal \mathbf{c}, \mathbf{f} and rounding (Section S.3.C.1) attains OPT 72% of
225 the time on this synthetic dataset. Computing \mathbf{f}, \mathbf{c} from the SVD of \tilde{X} yields OPT 0% of the time, even if we postprocess \mathbf{f} by
226 rounding its entries

227 While solving the SDP relaxation yields slightly better performance than alternating maximization, it is significantly more
228 computationally intensive and requires more complicated machinery (cvxpy (12)). To test the different methods’ computational
229 complexities, we increase the counts per column from 10 to 100 to make sure that all rows have counts. We then test for various
230 numbers of rows, to show the poor scaling of SDP solvers. Results are described in Table S2. The asymptotic approach identifies
231 the principal eigenvector of $D^{-1/2} \tilde{X} \tilde{X}^\top D^{-1/2}$ where $D = \text{diag}(\mathbf{p})$, which we naively compute with an eigendecomposition.
232 However, using power iteration, only matrix vector products with $\tilde{X}, \tilde{X}^\top$ are required, which can be much more efficient.

233 **S.3.E. Comparison with the SVD.** The SVD is frequently used for matrix decomposition and interpretation tasks, but doesn’t
234 inherently minimize the p-value bound of Proposition 1 or yield desirable partitionings. Mathematically, the SVD computes

$$\arg\max_{\mathbf{c}, \mathbf{f}} \frac{\mathbf{f}^\top \tilde{X} \mathbf{c}}{\|\mathbf{f}\|_2 \|\mathbf{c}\|_2},$$

as opposed to the p-value bound objective of

$$\arg\max_{\mathbf{c}, \mathbf{f}} \frac{\mathbf{f}^\top \tilde{X} \mathbf{c}}{\|\mathbf{f}\|_\infty \|\mathbf{c}\|_2}.$$

235 An optimal \mathbf{f} exists at a corner point, and that if entries of $\tilde{X}\mathbf{c}$ are not equal to 0 then this optimizing \mathbf{f} is unique up to the
 236 absolute value (\mathbf{f} and $1 - \mathbf{f}$ yielding the same objective value). For any $\tilde{\mathbf{f}}$ with entries not in $\{0, 1\}$, at least as good of an
 237 objective value can be attained by binarizing $\tilde{\mathbf{f}}$ according to the entries of $\tilde{X}\mathbf{c}$. As discussed in the previous section, computing
 238 $\tilde{\mathbf{f}}, \mathbf{c}$ from the SVD of \tilde{X} yields OPT 0% of the time, even if we postprocess $\tilde{\mathbf{f}}$ by rounding (binarizing) its entries, as it is
 239 still not optimizing the desired objective. We take a closer look at a toy example where an SVD does not yield the desired
 240 partitioning is shown in Figure S8.

241 S.4. SARS-CoV-2 analysis details

242 The SARS-CoV-2 analysis was based on dehosted sequencing data of nasopharyngeal swabs with coinfections available under
 243 accession PRJNA817806 (2). SPLASH contingency tables were tested using OASIS-opt with a 25% train/test data split. We
 244 then utilize anchors whose BY corrected p-value bound is less than 0.05, whose effect size (when \mathbf{c}^* is binarized) is in the top
 245 10% of anchors, and have a total number of counts $M > 1000$. This yields 2495 anchors. We then iterate over these 2495
 246 anchors, and rerun the alternating maximization procedure to generate $\mathbf{c}^*, \mathbf{f}^*$ on the full table. Beyond the aforementioned
 247 thresholds, there are additional potential avenues towards filtering that we did not utilize or fine-tune, including how balanced
 248 the clustering is (ensure that the table is not simply one deviating column), how similar a \mathbf{c} vector is to others (not detect
 249 splits due to SNPs), and reducing similar anchors (noting that one single base-pair mutation may yield up to k significant
 250 anchors that are one base-pair offsets of each other).

251 The alignment was performed using Bowtie (13). We used four SARS-CoV-2 reference genomes for the strains circulating
 252 during this time period. We manually constructed “archetype” genomes for Delta, Omicron BA.1, and Omicron BA.2 by
 253 editing the Original (Wuhan) reference (NCBI assembly NC_045512.2) to contain all (and only) the defining mutations of that
 254 strain, detailed and made publicly available in (5). A bowtie index comprised of the four SARS-CoV-2 genomes was constructed
 255 (bowtie-build on a FASTA file containing four reads, one for each assembly). We then construct a FASTA file from the anchors
 256 passing the above filtering steps, and their targets which comprise $> 5\%$ of the total counts for that anchor. This FASTA file is
 257 aligned to the bowtie index using (-v 0 -a) options to ensure that we only get exact matches, and record all possible matches.

S.4.A. Embedding-aggregation. To aggregate the $\mathbf{c}^{(a)}$ vectors obtained from experiments on different anchors a (different tables), one natural objective is

$$\begin{aligned} \mathbf{c}^* &= \operatorname{argmax}_{\|\mathbf{v}\| \leq 1} \sum_a \langle \mathbf{c}^{(a)}, \mathbf{v} \rangle^2 \\ &= \operatorname{argmax}_{\|\mathbf{v}\| \leq 1} \|\mathbf{C}\mathbf{v}\|_2. \end{aligned}$$

258 Here, we stacked the $\mathbf{c}^{(a)}$ vectors to obtain a matrix \mathbf{C} . Analyzing this objective, an optimizing \mathbf{v} is attained by the principal
 259 right singular vector of \mathbf{C} . To visualize the data further, we also plot the second singular vector, which can be understood as
 260 finding a vector with maximal sum of squared inner products with $\mathbf{c}^{(a)}$ subject to being orthogonal to the first \mathbf{c} identified, and
 261 see that this separates Omicron BA.1 from BA.2.

262 We provide additional analysis in Figure S9, where we show the spectrum of the embedding-aggregation matrix, the p-values
 263 generated by OASIS-iter, and

264 **S.4.B. Counts-aggregation.** Counts-aggregation iteratively identifies orthogonal \mathbf{c} vectors on the stacked contingency tables. To
 265 merge the tables, we discard all targets that have fewer than 10 counts in a given table. Then, we normalize each table by
 266 $1/\sqrt{M}$, its total number of counts, to ensure that high count tables do not completely outweigh low count ones, and stack the
 267 resulting tables together. This yields a visually similar clustering to embedding-aggregation, but does not perfectly predict
 268 with $\mathbf{c}^{(1)} \geq 0$ whether the sample has Delta or not. There is still perfect classification accuracy on $\mathbf{c}^{(1)}$ for predicting whether
 269 a sample has Delta or not, and 91% accuracy for predicting from $\mathbf{c}^{(1)}$ and $\mathbf{c}^{(2)}$ whether a sample has Omicron BA.1 as its
 270 primary strain or not (correctly classified 94/103).

271 For predicting target strain, of the 51,557 anchor-target pairs (rows of X_{agg}), we filtered for those targets which constitute
 272 at least 5% of the total counts for that anchor. Then, we label a target as Delta if it maps to the Delta assembly and neither of
 273 the Omicron ones, and Omicron if it maps to at least one of the Omicron assemblies but not the Delta one. This yields 4836
 274 sufficiently abundant anchor-target pairs which map, constituting 2,479 of the 2,495 anchors tested.

275 S.5. OASIS extensions

276 In this section we discuss several extensions of OASIS.

277 **S.5.A. Metadata-based construction of \mathbf{f}, \mathbf{c} .** OASIS is stated generally, and treats the rows as the set $[I]$ and the columns
 278 as the set $[J]$, with no additional information regarding shared structure. However, considering the motivating biological
 279 application, rows have additional information (the associated target sequence in $\{A, T, C, G\}^{27}$ (5)). Considering this, one
 280 potential direction to construct more biologically meaningful inference is to ensure that similar sequences have similar f_i values.
 281 This can be accomplished by constructing \mathbf{f} to be Lipschitz with respect to the Levenshtein distance between the targets of the
 282 associated rows. From a different direction, to improve interpretability, regularization can be imposed on \mathbf{f} in the optimization
 283 formulations to yield sparser vectors resulting in more parsimonious descriptions of deviations between samples.

284 Similarly, in the presence of sample metadata (e.g. cell-type) in the form of a vector $\tilde{\mathbf{c}} \in [-1, 1]^J$, the user may wish to
 285 incorporate this information into constructing \mathbf{c} . In the simplest case, the metadata is fully trusted, and we want to see if this
 286 sample partitioning yields a significant p-value bound. In this case, $\tilde{\mathbf{c}}$ is utilized, along with the optimizing \mathbf{f}^* from the split
 287 data, for inference.

288 One additional approach of potential interest is if metadata $\tilde{\mathbf{c}}$ is given, but we are still free to reweight samples: not all
 289 samples are equally representative. If we require that each c_j retains the same sign as \tilde{c}_j , then we can formulate this as:

$$290 \quad \mathbf{c}^* = \underset{\substack{\mathbf{c} : \|\mathbf{c}\| \leq 1, \\ c_j \tilde{c}_j \geq 0 \forall j}}{\operatorname{argmax}} \left| \mathbf{f}^\top \tilde{X} \mathbf{c} \right|. \quad [40]$$

291 While the overall objective is non-convex, it is still biconvex. The maximization of \mathbf{f} is the same as before, and \mathbf{c} is now
 292 maximized as, defining $\tilde{S} = \mathbf{f}^\top \tilde{X}$:

$$293 \quad c_j \propto \begin{cases} S_j & \text{if } \tilde{c}_j S_j \geq 0, \\ 0 & \text{if } \tilde{c}_j S_j < 0. \end{cases} \quad \text{or} \quad c_j \propto \begin{cases} -S_j & \text{if } \tilde{c}_j S_j \leq 0, \\ 0 & \text{if } \tilde{c}_j S_j > 0. \end{cases} \quad [41]$$

294 This enables efficient approximate optimization via alternating maximization. Many other forms of metadata (e.g. category of
 295 cell-type beyond just binary) can yield different metadata-aided approaches, and can be framed as optimization problems that
 296 admit efficient alternating maximization solutions.

297 **S.5.B. OASIS as a coefficient of correlation.** OASIS can also be used to test the dependence between two random variables.
 298 Consider observations of two real valued random variables (X, Y) , where we observe draws from the joint distribution (X_i, Y_i)
 299 and want to test against the null where X and Y are independent. One recent non-parametric measure of dependence was
 300 proposed, *xicor* (14), with good theoretical properties and empirical performance. OASIS can also be used to test independence
 301 by quantizing the random variables (e.g. into deciles) and constructing a contingency table from these categorical (quantized)
 302 observations. For any quantization, OASIS can yield a finite-sample valid p-value bound against the null of independence.
 303 This empirically performs quite well, and can yield an effect size estimate comparable to *xicor*, as shown in Figure S16. The
 304 example functions used are detailed in (14) and are plotted at left, showing that OASIS has additional power against balanced,
 305 structured alternatives, such as circular and heteroscedastic.

306 **S.5.C. Bernstein-based p-value bounds.** The finite-sample valid p-value bounds provided in Proposition 1 suffer from the
 307 fact that they do not incorporate the variance of f_Z where $Z \sim \mathbf{p}$, and simply rely on the boundedness of f_Z . One natural
 308 concentration inequality that utilizes both the boundedness of a random variable as well as its variance is Bernstein’s inequality,
 309 of which an empirical variant was derived in (15). An improved p-value bound can potentially be obtained by bounding
 310 the true variance σ_f^2 with high probability by some function of the empirical variance $\hat{\sigma}_f^2$ on the training set, and utilizing
 311 Bernstein’s inequality instead of Hoeffding’s, applied to the same test statistic. Preliminary analysis shows that in low count
 312 regimes this approach yields much worse bounds, as the variance estimate concentrates much slower than the mean estimate,
 313 but improved bounds and analyses could yield better performance.

314 **S.5.D. Tensor analysis.** In computational genomics, it is of high priority to discover effects or partitionings that are reproducible.
 315 In the setting of (5) this motivates going beyond treating each table independently, and instead analyzing the anchor by sample
 316 by target tensor. OASIS’s framework enables this, where \mathbf{f} and \mathbf{c} can be generated and optimized jointly.

317 Note that more generally OASIS, with its simple linear structure, suggests a simple test for tensors. If we have a 3
 318 dimensional tensor, we can construct and analyze a test statistic similar to OASIS by centering the tensor, appropriately
 319 normalizing by sampling depth, then taking inner products along each of the 3 dimensions to reduce to a scalar.

320 S.6. Additional discussion

321 For brevity and flow we defer some additional plots from the main text to here.

322 **S.6.A. Comparison of OASIS and X^2 in planted setting.** In this section we define and provide additional simulations (Figure S1)
 323 for the example in main text Figure 1. We define a class of alternatives parameterized by their total number of counts M and
 324 a corruption parameter ϵ as follows (the number of rows and columns are fixed as 12 and 10 respectively). In our uncorrupted
 325 planted model with $\epsilon = 0$, the probability vector for the first half (5) of the samples is the first standard basis vector $\mathbf{p}^{(j)} = \mathbf{e}_1$
 326 for $j = 1, \dots, J/2$, and the latter half of the samples have $\mathbf{p}^{(j)} = \mathbf{e}_2$ for $j = J/2 + 1, \dots, J$. To vary between structured
 327 and unstructured signal, we mix each column’s probability distribution with another probability distribution $\mathbf{q}^{(j)}$, which is
 328 generated independently for each column. Concretely,

$$329 \quad \mathbf{p}^{(j)} = \begin{cases} \mathbf{e}_1 + \epsilon \mathbf{q}^{(j)} & \text{if } j \leq J/2, \\ \mathbf{e}_2 + \epsilon \mathbf{q}^{(j)} & \text{if } j > J/2. \end{cases} \quad [42]$$

330 The observed matrix is generated by drawing $X^{(j)} \sim \text{multinomial}(M/J, \mathbf{p}^{(j)})$, independently for each of the J columns. Each
 331 $\mathbf{q}^{(j)}$ was independently generated by taking each entry to be independently uniformly distributed between $[0, 1]$, and then
 332 normalized to a valid probability distribution.

333 Considering in more detail the specifics of the simulations, in subplot C) we showed two tables highlighting the shortcomings
 334 of Pearson’s X^2 test. In both examples, power is hampered by the large number of rows. If the number of rows with
 335 nonzero counts doubles, and the value of the test statistic changes minimally, then the p-value X^2 provides will decrease
 336 dramatically (square root of the original p-value, $10^{-10} \rightarrow 10^{-5}$). Additionally, X^2 normalizes by the square root of the row
 337 sums, downweighting abundant rows and upweighting rare rows (which often arise from sequencing errors).

338 This is a critical drawback, as seen in the strong signal setting: because the first two rows are abundant, X^2 downweights
 339 their deviations, and yields only a moderately significant p-value. In the weak signal setting, even though no one row constitutes
 340 a large deviation, X^2 upweights the deviation in each row, and sums, yielding a similarly significant p-value. In contrast,
 341 OASIS identifies the strong and reproducible deviation in the the first setting, yielding an extremely significant p-value bound
 342 by focusing on the difference in expression of the first two rows. In the weak signal setting, OASIS struggles to find a strong
 343 separation between the samples, and only yields a slightly significant p-value bound. A spectral analysis of X_{corr} for the two
 344 tables is shown in Figure S1c, highlighting the flat spectrum of the weak signal table and concentrated spectrum of the strong
 345 signal table.

346 **S.6.B. Robustness against simulated noise.** OASIS provides robustness against many undesirable alternatives, beyond just
 347 negative binomial overdispersion. We provide simulations showing OASIS’s robustness against corruption of each individual
 348 column’s probability distribution, and additionally show OASIS’s improved robustness to deviation of one sample in our
 349 approximate power calculation in Section S.7. In all settings, **OASIS-opt** (OASIS-opt in legends) was run with 5 random splits,
 350 each with 25% training data. **OASIS-rand** (OASIS in legends) was run with 50 \mathbf{c} vectors were drawn uniformly at random from
 351 $\{-1, +1\}^J$ and 10 \mathbf{f} vectors uniformly at random from $\{0, 1\}^I$, and all 500 pairs were tested.

352 In our simulations we need to generate distributions over targets. The simplest option is to simply take a uniform distribution
 353 over the targets. Modelling the more realistic setting where some targets are more likely than others, we can also generate a
 354 vector $\tilde{\mathbf{p}}$ where each entry is drawn independently from some distribution, and renormalize this to a probability distribution,
 355 $\mathbf{p} = \tilde{\mathbf{p}} / \|\tilde{\mathbf{p}}\|_1$. In our simulations we utilize two base distributions. One obvious parameter free choice is a uniform distribution.
 356 To model the fact that some rows are much more abundant than others we also simulate with an exponential distribution. We
 357 select a rate parameter of 5 for this, arbitrarily.

358 **S.6.B.1. Negative binomial overdispersion.** As discussed, OASIS is robust to negative binomial overdispersion, as modelled in (16).
 359 Under the statistical null, for a common row distribution \mathbf{p} , if the number of observations in each column is Poisson distributed
 360 with rate λ , then the observed counts satisfy $X_{i,j} \sim \text{Pois}(\lambda p_i)$, independently drawn for each entry. However, as shown by (16),
 361 sequencing data is overdispersed in practice, and counts should be modelled as negative binomials. For an overdispersion
 362 parameter θ (where $\theta = 0$ corresponds to the true null), we thus generate synthetic data as $X_{ij} \sim \text{Pois}(\Gamma(\lambda p_i / \theta, \theta))$. Following
 363 the trends depicted in (16), we model our overdispersion parameter as a function of the expected number of counts to be
 364 observed n , approximated as $\theta(n) = 3/n$ if $x < 150$, and $\theta(n) = 0.02$ otherwise. In these plots, the number of observations per
 365 column was varied from 10 to 1000 in log-scale through 10 equally spaced points. For the tables we show one index which
 366 shows the full dynamic range, index 6/10 with a value of $\exp(\ln(10) + \frac{5}{9}(\ln(1000) - \ln(10))) \approx 129$.

367 In both settings, with uniform or exponential target distributions (Figures S3 and S4), OASIS provides significantly better
 368 control of the discovery rate against this uninteresting alternative stemming from overdispersion.

369 **S.6.B.2. Robustness again ℓ_1 corruption.** In this setting, we study an alternative where each column is uniformly distributed, but is
 370 then mixed (with weight ϵ) with a probability distribution $\mathbf{q}^{(j)}$, which is generated independently for each column. Concretely,

$$371 \quad \mathbf{p}^{(j)} = (1 - \epsilon) \frac{1}{r} \mathbf{1} + \epsilon \mathbf{q}^{(j)}, \quad [43]$$

372 where the synthetic data is generated by taking a Poisson number of observations from each column, multinomial with that
 373 column’s probability vector $\mathbf{p}^{(j)}$, independently for each column. To avoid additional parameters, a uniform target distribution
 374 was used as the true null. $\mathbf{q}^{(j)}$ were generated independently by taking each entry to be independent and uniformly distributed
 375 between $[0, 1]$ and normalizing. This is illustrated in Figure S5.

376 In simulations, we observe that OASIS is extremely robust to these types of perturbations; whereas X^2 immediately calls
 377 that these tables deviate from the null, OASIS does not prioritize these unstructured deviations and maintains control of the
 378 discovery rate for much larger corruption magnitudes, especially for larger tables.

379 **S.6.C. Power against simulated alternatives.** OASIS has power against alternatives with unique per-sample expression (motivated
 380 by V(D)J recombination) like those in Figure S7, and can be shown to have more power against splicing-type alternatives
 381 (Figure S5) than the X^2 test in some settings.

382 **S.6.C.1. Two-group setting.** Here we show OASIS’s power against alternative splicing type alternatives across a broader range of
 383 parameters. Alternative tables were generated by taking two groups of samples of equal size, and defining the two group’s
 384 probability distributions as, for $k = 1, 2$

$$385 \quad p_i^{(k)} = \begin{cases} 1 - \epsilon + \epsilon/I & \text{if } i = k, \\ \epsilon/I & \text{otherwise.} \end{cases} \quad [44]$$

386 This indicates that $\mathbf{p}^{(1)}$ has most of its counts in the first row, and $\mathbf{p}^{(2)}$ in the second. The number of observations per sample
 387 is drawn independently from a Poisson distribution with mean 20 for each sample. Figure S6 shows that OASIS can have more
 388 power than X^2 in a variety of parameter regimes.

389 **S.6.C.2. Unique per-sample expression.** In this setting, the probability distribution of column $j \in [J]$ is $\mathbf{p}^{(j)} \in [0, 1]^I$, with $I = J$:

$$390 \quad p_i^{(j)} = \begin{cases} (1 - \epsilon) + \epsilon/r & \text{if } i = j, \\ \epsilon/r & \text{otherwise.} \end{cases} \quad [45]$$

391 This setting is illustrated in Figure S7, where a corruption of $\epsilon = 0$ corresponds to pure signal, and $\epsilon = 1$ corresponds to a
 392 uniform target distribution for all columns (no signal). We provide simulation results for this setting in Figure S7, where in
 393 most settings Pearson's X^2 has more power than OASIS, but not by a significant margin.

394 S.7. Power analysis

395 Providing exact power calculations of OASIS and X^2 under different alternatives is difficult, as the random multinomial draws
 396 affect both X and E . Even asymptotic analyses would require more advanced tools from random matrix theory. To provide
 397 intuition for the power profile of OASIS relative to X^2 , we instead perform approximate power calculations in a toy setting.
 398 Concretely, we assume that we observe a scaled version of the underlying alternative matrix, and compute OASIS and X^2
 399 p-values for this observed X . We show that OASIS indeed has substantial power in many common settings, and can be more
 400 powerful than X^2 in certain regimes. We conjecture that while OASIS may not exhibit the optimal asymptotic rate against
 401 all classes of alternatives, e.g. the unique per-sample expression example, it does have power going to 1 as the number of
 402 observations goes to infinity across a broad class of alternatives.

403 In our power computations, we additionally assume that OASIS identifies an optimal pair of \mathbf{c}, \mathbf{f} . It is reasonable to assume
 404 that this occurs in the asymptotic regime, and further since OASIS's test statistic is bilinear in \mathbf{c}, \mathbf{f} , its p-value bound is
 405 continuous in these parameters, and so OASIS only needs to identify a pair of near optimal \mathbf{c}, \mathbf{f} to yield a similar p-value
 406 bound.

407 The X^2 p-value for k degrees of freedom is asymptotically distributed as $\mathcal{N}(k, 2k)$. Since the X^2 test statistic is $X^2 =$
 408 $\left\| \text{diag}(\sqrt{M/X\mathbf{1}})\tilde{X} \right\|_F^2$, plugging in $k = (I - 1) \times (J - 1)$, and examining large values of the test-statistic,

$$409 \quad P_{X^2} \approx \exp\left(-\frac{X^2}{4IJ}\right). \quad [46]$$

S.7.A. Two-group alternative. Here we assume for simplicity that all n_j are equal, and that the clusters are of equal sizes. We
 describe the results in the context of differentially regulated alternative splicing, but the analysis holds for general 2-group
 alternatives with the below structure. For samples $j = 1, \dots, J/2$, the target (row) distribution is $\mathbf{p}^{(1)}$, and for the second half
 of columns $\mathbf{p}^{(2)}$. In this setting we have:

$$X = n \left(\mathbf{p}^{(1)} \left[\mathbf{1}_{J/2}^\top \ \mathbf{0}_{J/2}^\top \right] + \mathbf{p}^{(2)} \left[\mathbf{0}_{J/2}^\top \ \mathbf{1}_{J/2}^\top \right] \right)$$

and thus

$$E = n \left(\frac{\mathbf{p}^{(1)} + \mathbf{p}^{(2)}}{2} \right) \mathbf{1}_J^\top$$

To simplify our results, we utilize the symmetric chi-squared distance between two probability distributions $\chi^2(\mathbf{p}, \mathbf{q}) =$
 $2 \sum_i \frac{(p_i - q_i)^2}{p_i + q_i}$. The Pearson X^2 statistic is then

$$\begin{aligned} X^2 &= \sum_{i,j} \frac{(X_{i,j} - E_{i,j})^2}{E_{i,j}} \\ &= n \sum_{i,j < J/2} \frac{\left(p_i^{(1)} - \frac{p_i^{(1)} + p_i^{(2)}}{2} \right)^2}{\left(p_i^{(1)} + p_i^{(2)} \right) / 2} + n \sum_{i,j \geq J/2} \frac{\left(p_i^{(2)} - \frac{p_i^{(1)} + p_i^{(2)}}{2} \right)^2}{\left(p_i^{(1)} + p_i^{(2)} \right) / 2} \\ &= \frac{nJ}{2} \sum_i \frac{\left(p_i^{(1)} - p_i^{(2)} \right)^2}{p_i^{(1)} + p_i^{(2)}} \\ &= \frac{M}{4} \chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}) \end{aligned} \quad [47]$$

410 For disjoint $\mathbf{p}^{(1)}, \mathbf{p}^{(2)}$, the χ^2 distance is equal to 1. From the $(I - 1) \times (J - 1)$ degrees of freedom, this yield an approximate
 411 p-value of

$$P_{X^2} \approx \exp\left(-\frac{M^2 (\chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}))^2}{16I \times J}\right)$$

For OASIS, a pair of optimal \mathbf{c} and \mathbf{f} can be constructed where \mathbf{c} is +1 for the first $J/2$ components, and -1 for the latter half, and \mathbf{f} is $\text{sign}(\mathbf{p}^{(1)} - \mathbf{p}^{(2)})$. Then, the OASIS test statistic is

$$\begin{aligned} S &= \sum_{j=1}^{J/2} \sqrt{n_j} c_j (\mathbf{f}^\top \mathbf{p}^{(1)} - \bar{\mu}) + \sum_{j=J/2+1}^J c_j \sqrt{n_j} (\mathbf{f}^\top \mathbf{p}^{(2)} - \bar{\mu}) \\ &\stackrel{(a)}{=} \sum_{j=1}^{J/2} \sqrt{n} \mathbf{f}^\top (\mathbf{p}^{(1)} - \mathbf{p}^{(2)}) \\ &= \frac{\sqrt{MJ}}{2} \mathbf{f}^\top (\mathbf{p}^{(1)} - \mathbf{p}^{(2)}) \\ &\stackrel{(b)}{=} \frac{\sqrt{MJ}}{2} \delta_{\text{TV}}(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}) \end{aligned} \quad [48]$$

412 where in (a) we used the fact that all $n_j = n$ were the same and that the clusters were equally balanced, and so the opposing
413 signs of c_j cancelled the effect of $\bar{\mu}$. In (b) we used the optimal $\mathbf{f} = \text{sign}(\mathbf{p}^{(1)} - \mathbf{p}^{(2)})$.

414 **S.7.A.1. Summary.** The X^2 p-value can be approximated as

$$P_{X^2} \approx \exp\left(-\frac{M^2}{I \times J} (\chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}))^2\right). \quad [49]$$

416 OASIS has significant power against alternative splicing-type alternatives. Assuming that OASIS identifies the optimal \mathbf{c}, \mathbf{f} ,
417 which will be the case asymptotically, it will yield a p-value bound of

$$P_{\text{OASIS}} \leq 2 \exp\left(-\frac{M}{4} \delta_{\text{TV}}^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)})\right). \quad [50]$$

419 Assuming the total variation distance and χ^2 distance between the distributions are constant, OASIS has more power than
420 X^2 when $I \times J \geq M$; this is often the case we are in, where there are many rows with very low counts. Assumptions on the
421 balanced classes can be alleviated by using a class-weighted mean estimate, which computes the mean f value over samples
422 with $c_j > 0$, and averages this with the mean f value over samples with $c_j < 0$, requiring a different theoretical analysis.

S.7.B. Time series. Consider a simple time series setting, where at time $t = 0$ cells have target distribution $\mathbf{p}^{(1)}$, at time
 $t = T$ distribution $\mathbf{p}^{(2)}$, and their distributions vary smoothly in between, i.e. at time t they have the mixture distribution
 $\mathbf{p}^{(1)} + t(\mathbf{p}^{(2)} - \mathbf{p}^{(1)}) = (1 - t)\mathbf{p}^{(1)} + t\mathbf{p}^{(2)}$. This can model microbial evolution in bulk RNA-seq data, cell-cycle in single
cell RNA-seq, and many other settings. For $T + 1$ evenly spaced time intervals (T is assumed to be odd), define the vector
 $\mathbf{t} = [0, 1/T, \dots, (T - 1)/T, 1]$. We observe that, for n observations per column, our expected data matrix X and expected
matrix E are

$$\begin{aligned} X &= (\mathbf{p}^{(1)}(1 - \mathbf{t})^\top + \mathbf{p}^{(2)}\mathbf{t}^\top) n, \\ E &= \frac{n}{2} (\mathbf{p}^{(1)} + \mathbf{p}^{(2)}) \mathbf{1}^\top. \end{aligned}$$

Pearson's X^2 test statistic is:

$$\begin{aligned} X^2 &= M \left\| \text{diag}(1/\sqrt{X\mathbf{1}})(X - E) \text{diag}(1/\sqrt{X^\top \mathbf{1}}) \right\|_F^2 \\ &= nM \left\| \text{diag}(1/\sqrt{X\mathbf{1}}) (\mathbf{p}^{(1)} - \mathbf{p}^{(2)}) \left(\frac{1}{2}\mathbf{1} - \mathbf{t}\right)^\top \right\|_F^2 \\ &= nM \sum_{i=1}^I \frac{(p_i^{(1)} - p_i^{(2)})^2}{\frac{M}{2} (p_i^{(1)} + p_i^{(2)})} \left\| \frac{1}{2}\mathbf{1} - \mathbf{t} \right\|_2^2 \\ &= \frac{M}{12} \chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}) + o(T), \end{aligned} \quad [51]$$

simplifying the Frobenius norm of a rank one matrix. $\left\|\frac{1}{2}\mathbf{1} - \mathbf{t}\right\|_2^2 = \frac{(T+1)(T+2)}{12T} = \frac{T}{12} + o(T)$. This leads to an approximate p-value of

$$P_{X^2} \approx \exp\left(-\frac{M^2}{48I \times T} (\chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}))^2\right)$$

We claim that a pair of optimal (\mathbf{c}, \mathbf{f}) can be constructed as

$$\begin{aligned}\mathbf{c} &= \frac{1}{2} - \mathbf{t}, \\ \mathbf{f} &= \text{sign}(\mathbf{p}^{(1)} - \mathbf{p}^{(2)}).\end{aligned}$$

This yields

$$\begin{aligned}S &= \sqrt{n}\mathbf{f}^\top (\mathbf{p}^{(2)} - \mathbf{p}^{(1)}) \left(\mathbf{t} - \frac{1}{2}\mathbf{1}\right)^\top \mathbf{c} \\ &= 2\sqrt{n}\delta_{\text{TV}}(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}) \left\|\mathbf{t} - \frac{1}{2}\mathbf{1}\right\|^2 \\ &= 2\sqrt{n}\delta_{\text{TV}}(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}) \frac{(T+1)(T+2)}{12T}\end{aligned}$$

This gives OASIS a p-value bound of

$$\begin{aligned}P_{\text{OASIS}} &\leq 2 \exp\left(-\frac{2S^2}{\|\mathbf{c}\|^2 \|\mathbf{f}\|_\infty^2}\right) \\ &\leq 2 \exp\left(-n(T+1)\delta_{\text{TV}}^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)})/6\right) \\ &= 2 \exp\left(-M\delta_{\text{TV}}^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)})/6\right)\end{aligned}\tag{52}$$

S.7.B.1. Summary. X^2 provides an approximate asymptotic p-value of

$$\begin{aligned}P_{X^2} &\approx \exp\left(-\frac{M^2 (\chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}))^2}{I \times T}\right) \\ &\approx \exp\left(-M (\chi^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}))^2 \times \frac{n}{I}\right).\end{aligned}$$

OASIS yields a p-value bound of

$$P_{\text{OASIS}} \leq \exp\left(-M\delta_{\text{TV}}^2(\mathbf{p}^{(1)}, \mathbf{p}^{(2)})/6\right)\tag{53}$$

423 This setting displays the same behavior as in alternative splicing, where OASIS can yield improved performance when $\frac{n}{T}$ is
424 small.

S.7.C. Unique per-sample expression. In this setting, the alternative matrix is $c \times c$, where the j -th sample expresses target (row) j with probability $1 - \alpha$, and with probability α is uniform over the rest of the $c - 1$ rows. The number of observations per column are assumed to be the same and c is even. Formally, the alternative matrix A is as below:

$$\begin{aligned}A &= \left(1 - \frac{c}{c-1}\alpha\right) I_c + \frac{\alpha}{c-1} \mathbf{1}\mathbf{1}^\top \\ X &= A \text{diag}(\mathbf{n}) \\ E &= \frac{1}{c} \mathbf{1}\mathbf{1}^\top \text{diag}(\mathbf{n})\end{aligned}$$

In this setting, Pearson's X^2 statistic is

$$\begin{aligned}X^2 &= \sum_{ij} \frac{(X_{ij} - E_{ij})^2}{E_{ij}} \\ &= \sum_j \frac{n_j^2 \left((1 - \alpha) - \frac{1}{c}\right)^2}{n_j/c} + \sum_{i \neq j} \frac{n_j^2 \left(\frac{\alpha}{c-1} - \frac{1}{c}\right)^2}{n_j/c} \\ &= c \left(1 - \alpha - \frac{1}{c}\right)^2 \sum_j n_j + c \left(\frac{\alpha}{c-1} - \frac{1}{c}\right)^2 \sum_{i \neq j} n_j \\ &= cM(1 - \alpha - 1/c)^2 + \underbrace{c(c-1)M \left(\frac{\alpha}{c-1} - \frac{1}{c}\right)^2}_{\text{negligible}}.\end{aligned}\tag{54}$$

Thus, for fixed c , the X^2 p-value can be approximated as

$$P_{X^2} \approx \exp(-M^2(1-\alpha)^4)$$

Analyzing OASIS, we assume that it identifies an optimal pair of \mathbf{f} , \mathbf{c} , where since n_j are balanced

$$f_i = \begin{cases} 1 & \text{if } 1 \leq i \leq c/2 \\ 0 & \text{if } c/2 < i \leq c \end{cases} \quad \text{and} \quad c_j = \begin{cases} 1 & \text{if } 1 \leq j \leq c/2 \\ -1 & \text{if } c/2 < j \leq c \end{cases}$$

constitute an optimal pair. Since $n_j = M/c$ for all j , $\bar{\mu} = 1/2$, and

$$\hat{\mu}_j = \begin{cases} 1 - \frac{c}{c-1} \frac{\alpha}{2} & \text{if } 1 \leq i \leq c/2, \\ \frac{c}{c-1} \frac{\alpha}{2} & \text{if } c/2 < i \leq c. \end{cases}$$

This yields

$$\begin{aligned} S &= \sum_j c_j \sqrt{n_j} (\hat{\mu}_j - \bar{\mu}) \\ &= \sum_j \sqrt{M/c} \frac{1}{2} \left(1 - \frac{c}{c-1} \alpha\right) \\ &= \sqrt{cM} \frac{1}{2} \left(1 - \frac{c}{c-1} \alpha\right) \end{aligned} \tag{55}$$

and a corresponding p-value bound of

$$P_{\text{OASIS}} \leq \exp\left(-\frac{S^2}{\|\mathbf{c}\|_2^2}\right) = \exp\left(-\frac{M}{4} \left(1 - \frac{c}{c-1} \alpha\right)^2\right) \approx \exp(-M(1-\alpha)^2)$$

425 **S.7.C.1. Summary.** In this toy-setting of unique per-sample expression, X^2 will yield a p-value of

$$426 \quad P_{X^2} \approx \exp(-M^2(1-\alpha)^4). \tag{56}$$

427 Comparatively, assuming OASIS identifies an optimal pair (\mathbf{c}, \mathbf{f}) , its p-value bound can be approximated as

$$428 \quad P_{\text{OASIS}} \leq \exp\left(-\frac{M}{4} \left(1 - \frac{c}{c-1} \alpha\right)^2\right) \approx \exp(-M(1-\alpha)^2). \tag{57}$$

429 While this is worse for constant α and M tending to infinity, OASIS still has power going to 1 as M increases.

S.7.D. One deviant sample. Consider the setting where one sample is maximally deviating, i.e. the observed matrix is:

$$X = \begin{bmatrix} 0 & n & \dots & n \\ n & 0 & \dots & 0 \end{bmatrix}.$$

430 Then, the expected matrix is

$$E = \begin{bmatrix} n(1-1/J) & n(1-1/J) & \dots & n(1-1/J) \\ n/J & n/J & \dots & n/J \end{bmatrix}.$$

Expressing the summation as grouped by $X_{00}, X_{10}, X_{0,1}, X_{1,1}$:

$$\begin{aligned} X^2 &= \frac{(n(1-1/J))^2}{n(1-1/J)} + \frac{(n-n/J)^2}{n/J} + (J-1) \frac{(n-n(1-1/J))^2}{n(1-1/J)} + (J-1) \frac{(n/J)^2}{n/J} \\ &= n(1-1/J) + nJ(1-1/J)^2 + n/J + n(1-1/J) \\ &= nJ \end{aligned}$$

Thus, the Pearson X^2 test will yield an approximate p-value of

$$P_{X^2} \approx \exp\left(-\frac{(nJ)^2}{2(2J)^2}\right) = \exp(-n^2/8)$$

Comparatively, OASIS will choose $\mathbf{f} = [1 \ 0]^\top$, which yields $S_j = \sqrt{n}/J$ for $j = 2, \dots, n$ and $S_1 = -\sqrt{n}(1-1/J)$. This yields an optimizing $\mathbf{c} = [-(1-1/J) \ 1/J \ \dots \ 1/J]$ with $\|\mathbf{c}\|^2 = 1-1/J$. The test statistic is then $S = \sqrt{n}\|\mathbf{c}\|^2$. Thus,

$$P_{\text{OASIS}} \leq \exp(-S^2/\|\mathbf{c}\|^2) = \exp(-n).$$

431 **S.7.D.1. Summary.** X^2 provides a p-value of

432
$$P_{X^2} \approx \exp\left(-\frac{(nJ)^2}{2(2J)^2}\right) = \exp(-n^2/8), \quad [58]$$

433 in comparison to OASIS-opt's p-value bound of

434
$$P_{\text{OASIS}} \leq \exp(-S^2/\|c\|^2) = \exp(-n). \quad [59]$$

435 This shows how X^2 can be significantly overpowered against alternatives where only one column deviates.

436 **S.7.E. Conclusions from power analysis.** As can be seen, there are several regimes where the Pearson X^2 test has less power
437 than OASIS, primarily when the table is sparse, i.e. $\frac{M}{I \times J}$ is small. This highlights one of the difficulties with the X^2 test; its
438 p-value depends heavily on the number of rows present, even if many of them are inconsequential.

439 **S.8. Supplemental figures**

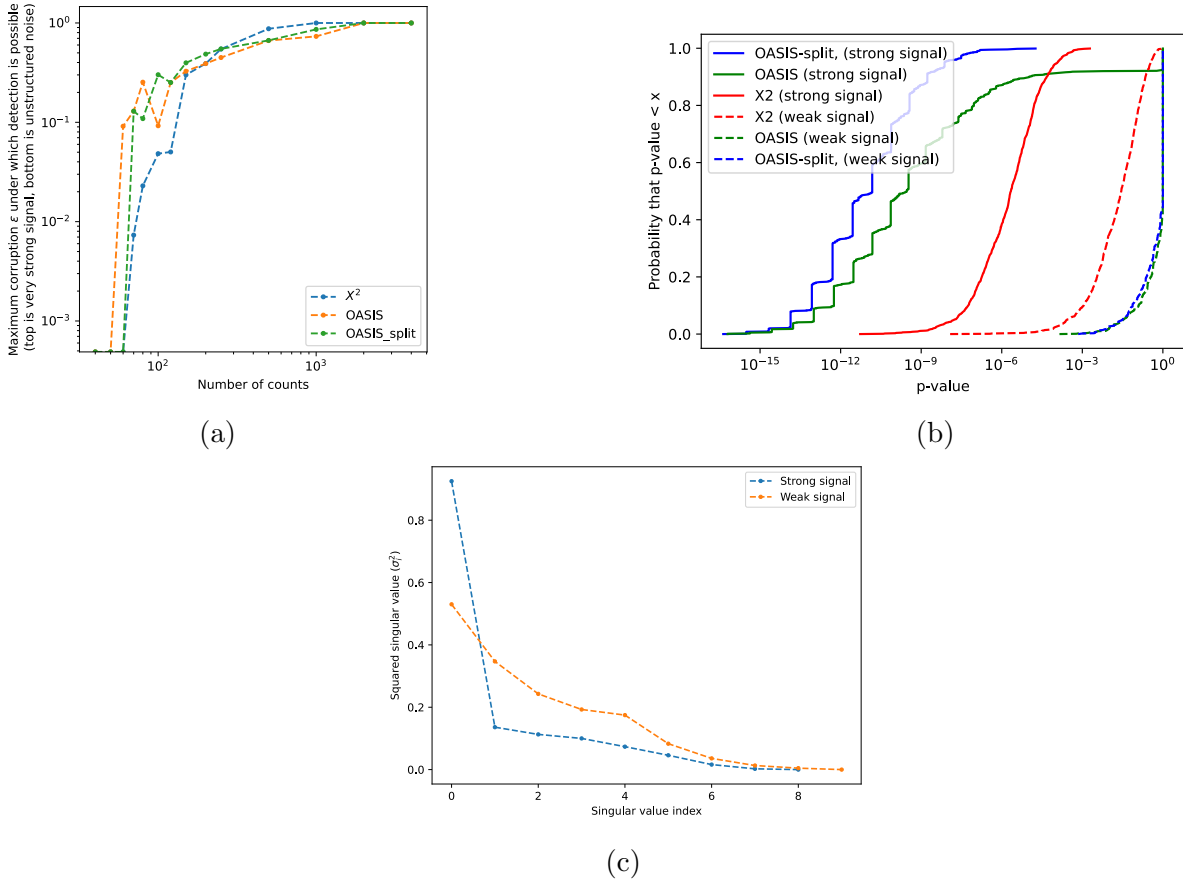
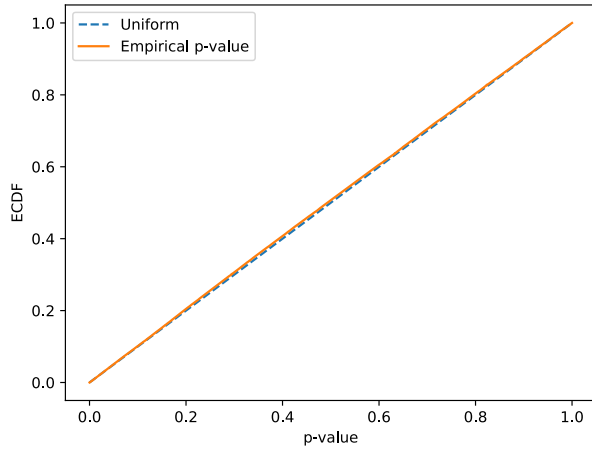
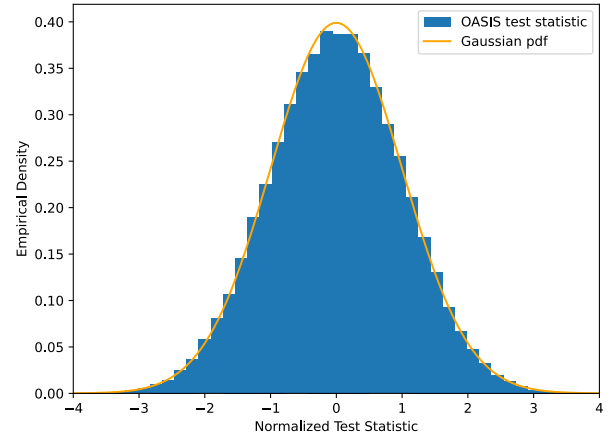


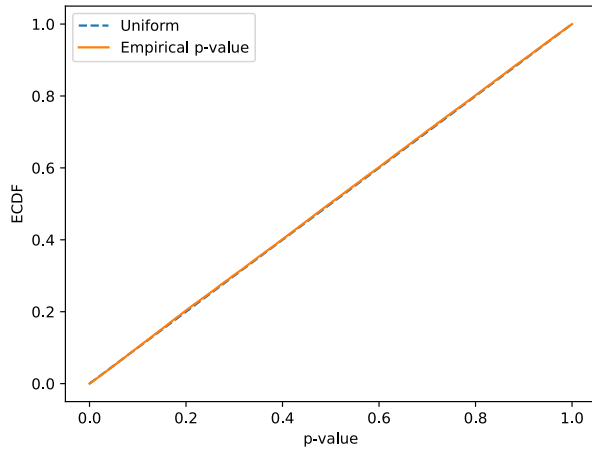
Fig. S1. Following Figure 1D, this plot shows the performance of OASIS-opt and X^2 over the class of 2-group alternatives defined in Eq. (42). (a) shows an empirical achievability curve, where the y-axis indicates how structured versus unstructured the alternative is, and the x-axis indicates the total number of counts in the table. The more up and to the left a curve is, the better. For any alternative (any $\epsilon \in [0, 1]$), a statistical test should reject the null with enough observations. For a given number of counts, we performed binary search to identify the largest ϵ for which over half of the p-values (bounds) were less than $1E - 10$. Empirically, OASIS is more likely to reject highly structured tables with lower counts, whereas Pearson's X^2 test is more likely to reject unstructured tables with moderate counts. OASIS-opt is run with 5 random splits as in Figure 1D, which we denote as OASIS-split, and we additionally show results (labelled OASIS) when only 1 random split is used. (b) plots an ECDF of p-values in the two settings for different methods. Performing 5 random splits removes the tail of the p-value bounds, generally improving performance. The flat behavior between adjacent data points for OASIS-opt is due to the data-splitting: an increase in the number of observations does not necessarily translate to an increase in the number of observations in the test data. (c) shows the spectrum of X_{corr} for the two example tables in main text Figure 1c. X^2 yields a similar p-value for the two tables, as the sum of the squared singular values is similar. OASIS prioritizes the “strong signal” table with the more concentrated spectrum.



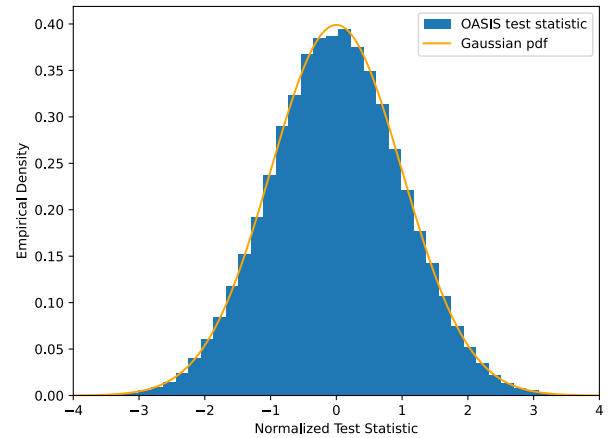
(a)



(b)



(c)



(d)

Fig. S2. Distribution of the normalized OASIS test statistic $\frac{S}{\hat{\sigma}_f \|c\| \sqrt{1-\gamma}}$ (right). As predicted by Proposition 3 this follows a Gaussian distribution, leading to uniformly distributed asymptotic p-values (left). The synthetic tables generated have few counts, but the approximation is still extremely good. Each plot shows 100k randomly generated tables with 5 rows, 8 columns, and a random target distribution (each entry i.i.d. uniform, normalized). Column counts for simulations in the first row are independently drawn from a Poisson with mean 10, while the second row has mean 30. For each trial a random p was generated, $\{n_j\}$ were drawn randomly from Poisson distributions, and X was generated from this null. Then, c and f were independently randomly generated, where the coordinates of f were i.i.d. $\text{Ber}(1/2)$, and the coordinates of c were i.i.d. $U([-1, 1])$. This process was repeated until a table with $\hat{\sigma}_f > 0$ and $\gamma < 1$ was generated to yield a valid run (as otherwise the test statistic is identically 0). KS distance (maximum distance between OASIS's asymptotic p-value's ECDF and the uniform distribution) of .0079 for $n = 10$ counts per column, .0038 for $n = 30$ counts per column.

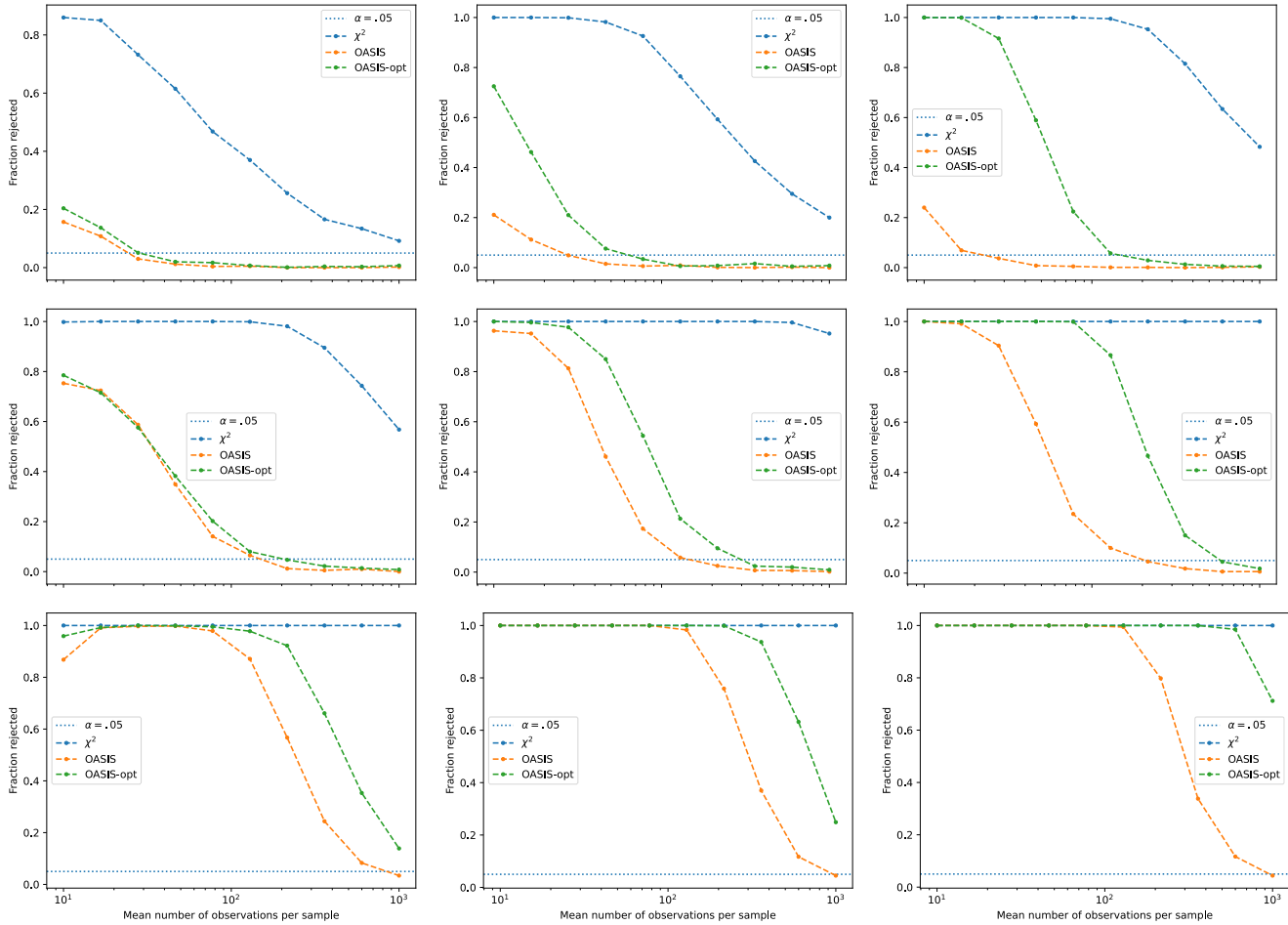


Fig. S3. Simulated data for negative binomial overdispersion. Null target distribution generated from exponential distributions, described in Section S.6.B.1. OASIS has significantly better control of the false discovery rate than the χ^2 test, and requires substantially fewer samples per column to control the FDR. The 3x3 grid of plots displays results for tables varying in number of columns (10, 50, 400 from left to right) and number of rows (5, 20, 100 from top to bottom).

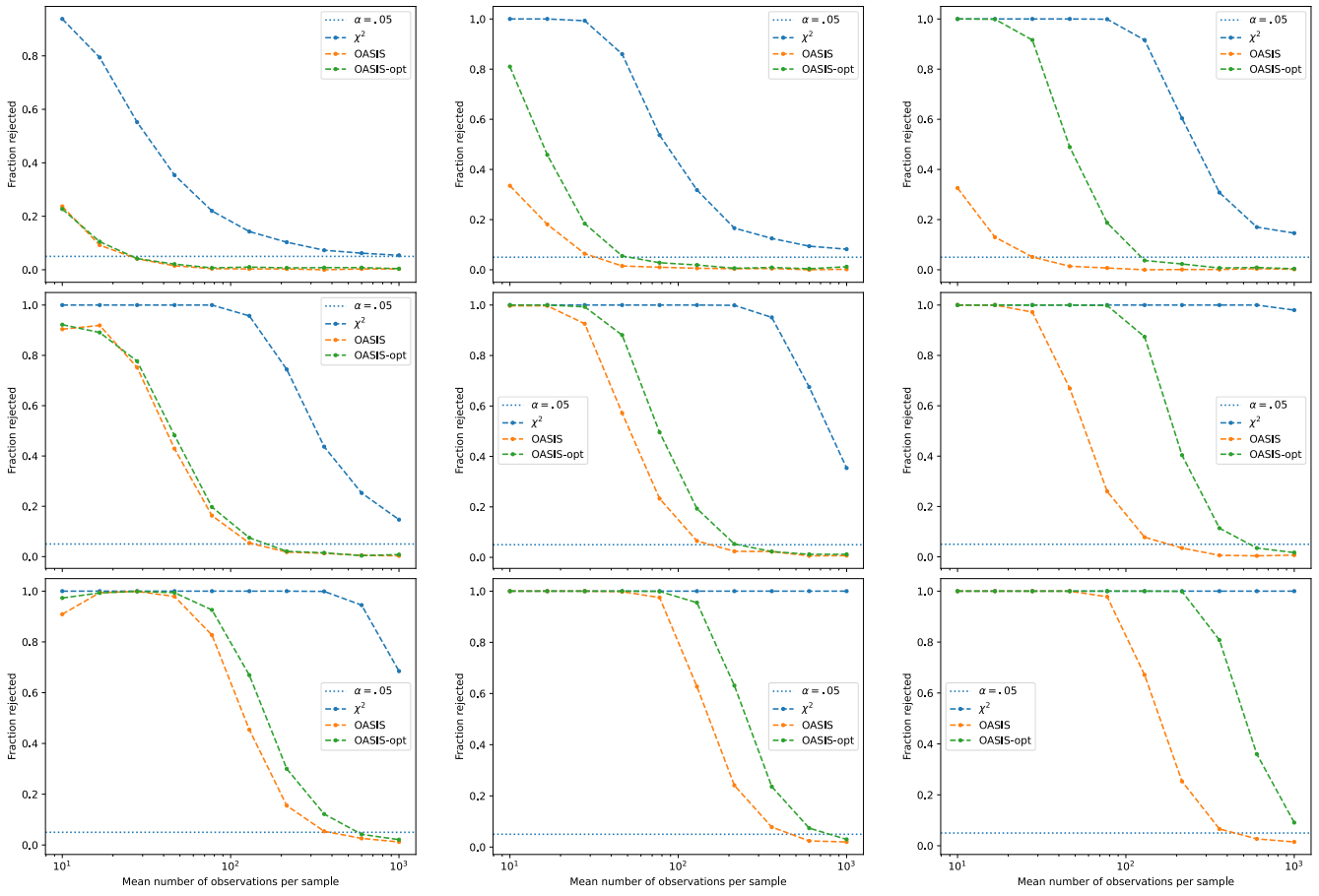


Fig. S4. Simulated data for negative binomial overdispersion. Null target distribution is uniform over rows. OASIS has significantly better control of the false positive rate than the χ^2 -test, and requires substantially fewer samples per column to control the FDR. The 3x3 grid of plots displays results for tables varying in number of columns (10, 50, 400 from left to right) and number of rows (5, 20, 100 from top to bottom).

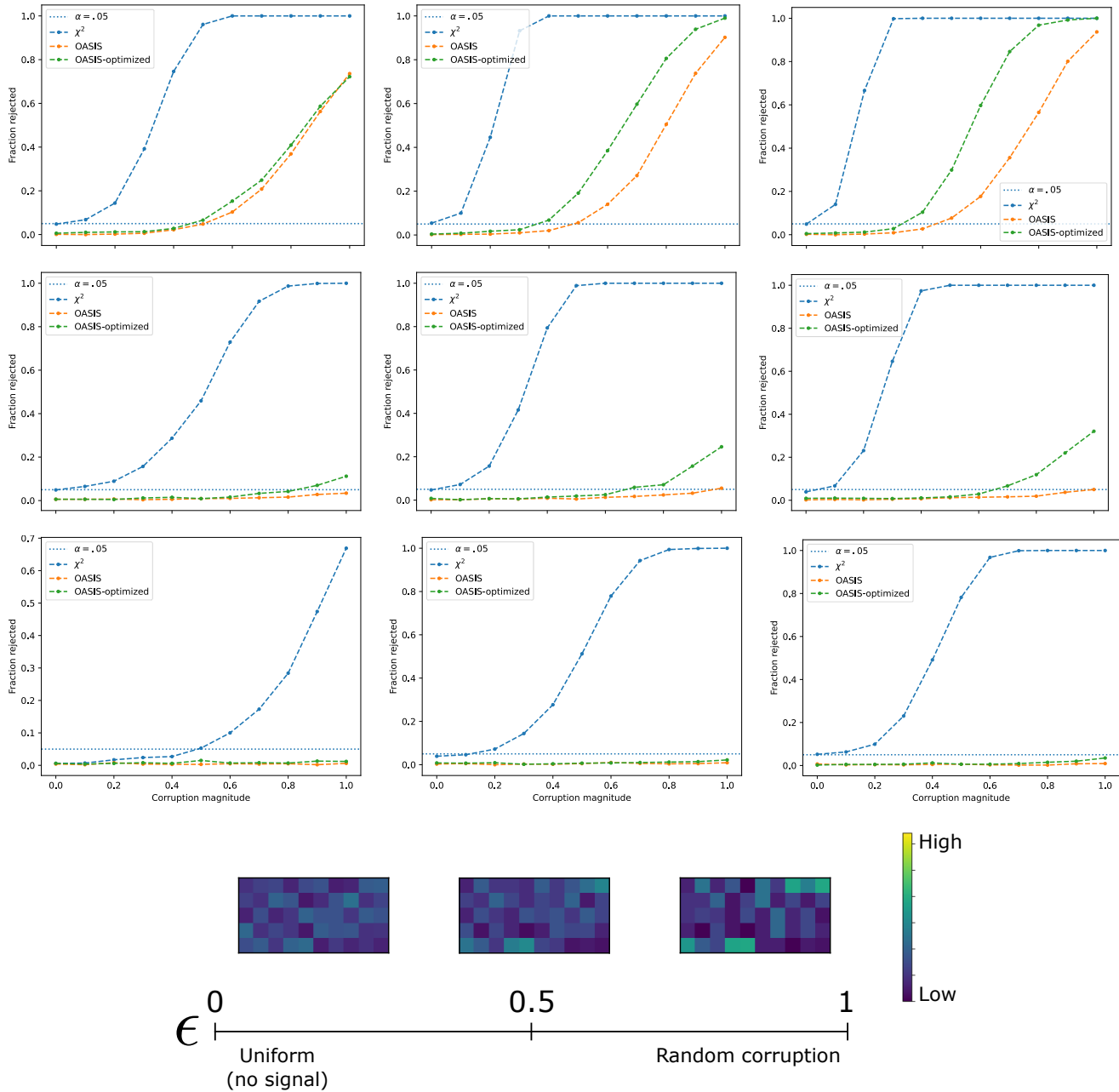


Fig. S5. Simulated data for independent ℓ_1 corruption of each column (Eq. (43)). At bottom, a diagram showing tables under different corruption magnitudes ϵ . Number of observations in each column is Poisson distributed with mean 100, independent across columns. Observations are drawn independently, and are multinomially distributed for each column. The 3x3 grid of plots displays results for tables varying in number of columns (10, 50, 100 from left to right) and number of rows (20, 10, 500 from top to bottom). χ^2 has significant power against this alternative, even though it represents contamination and a “biological null”.

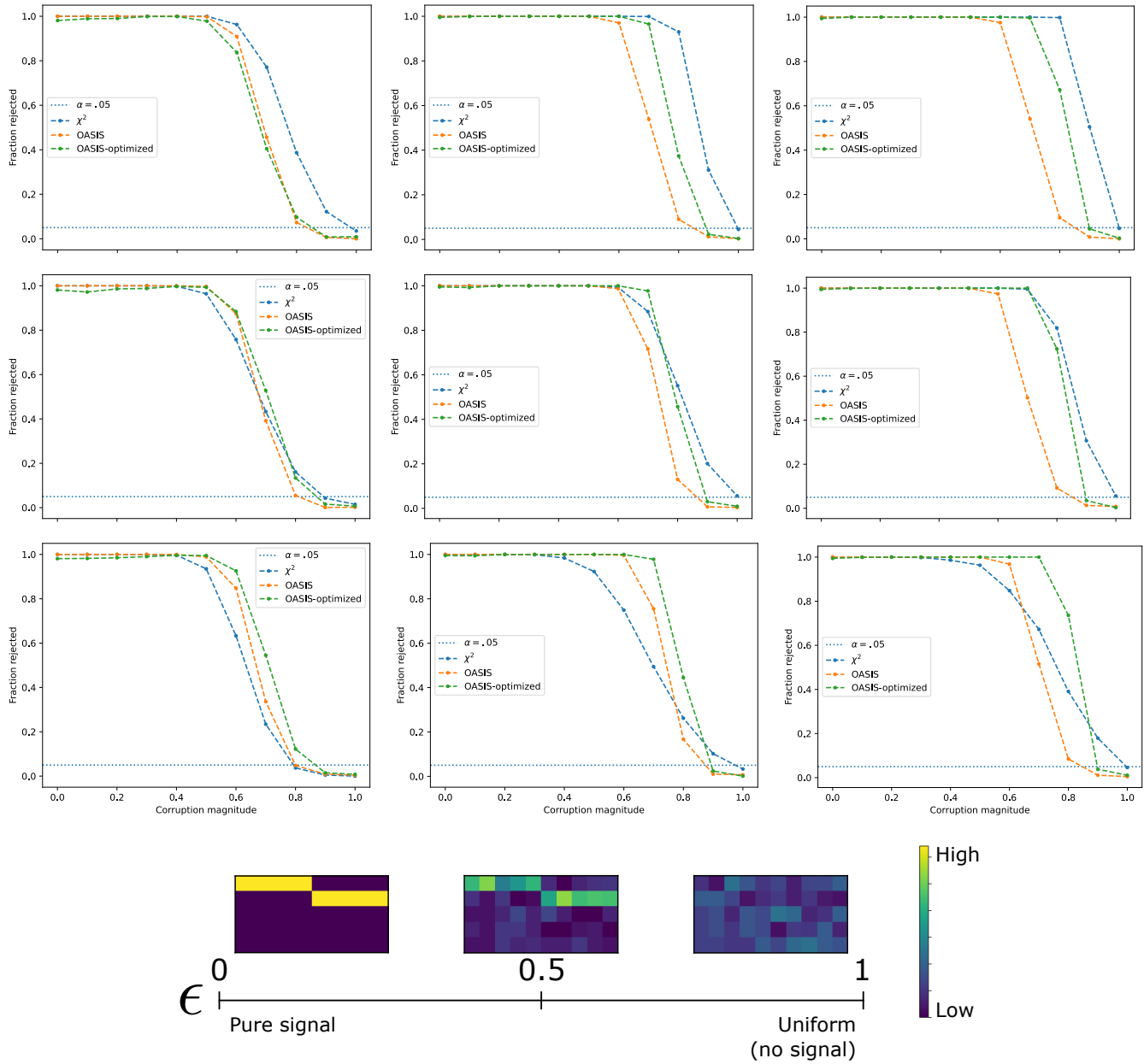


Fig. S6. Simulated data for alternative splicing-type alternative (Eq. (44)). Number of observations in each column is Poisson distributed with mean 20. The 3x3 grid of plots displays results for tables varying in number of columns (10, 50, 100 from left to right) and number of rows (20, 100, 400 from top to bottom). OASIS generally has more power than χ^2 .

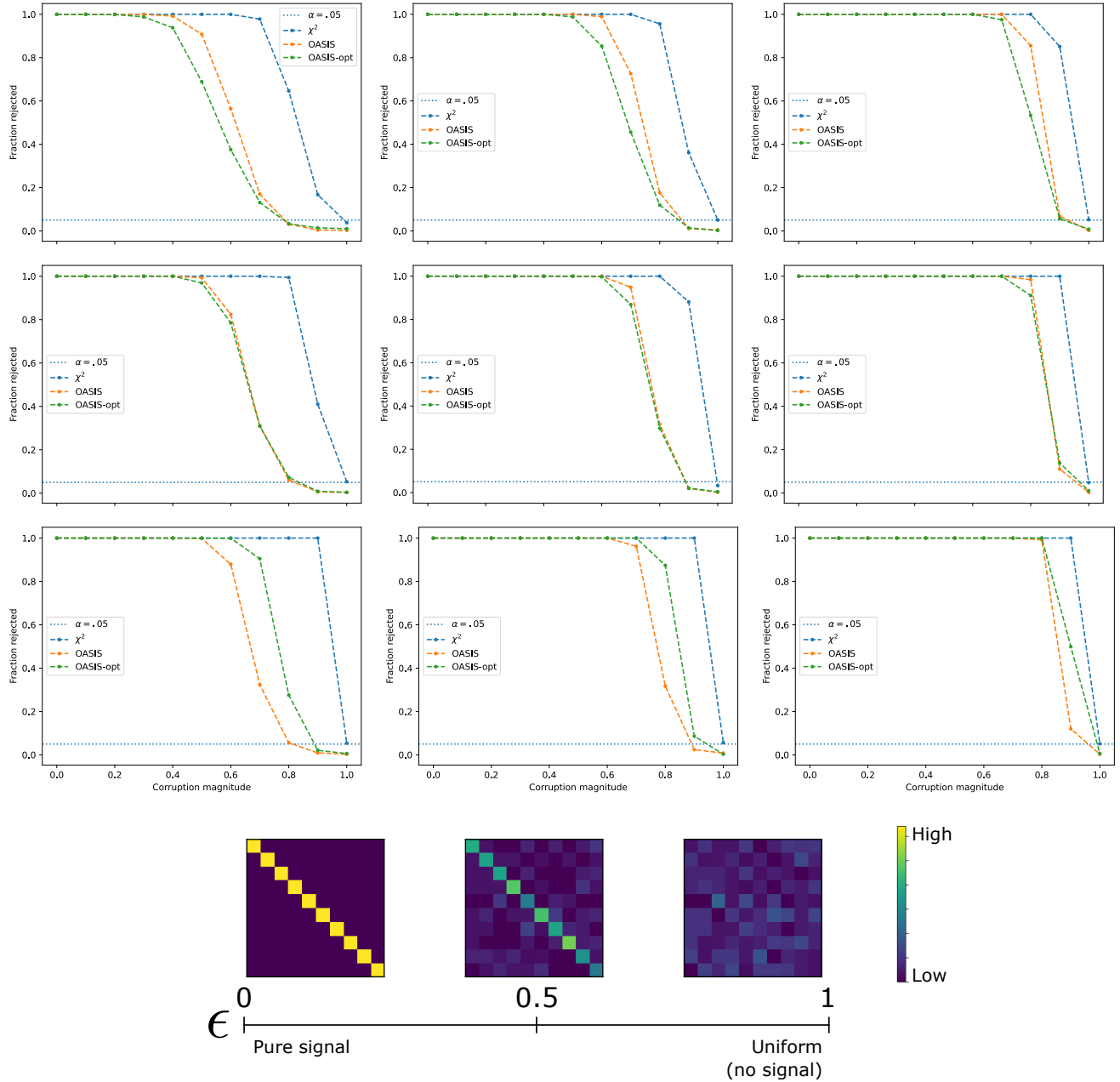


Fig. S7. Simulated data for unique per-sample expression alternative (Eq. (45)). Number of observations in each column is Poisson distributed with the given mean, independent across columns. Detailed probabilistic model in Eq. (45). The 3x3 grid of plots displays results for tables varying in number of rows (10, 25, 100 from top to bottom) and mean number of observations per column (10, 20, 50 from left to right). χ^2 has more power, but all methods perform well.

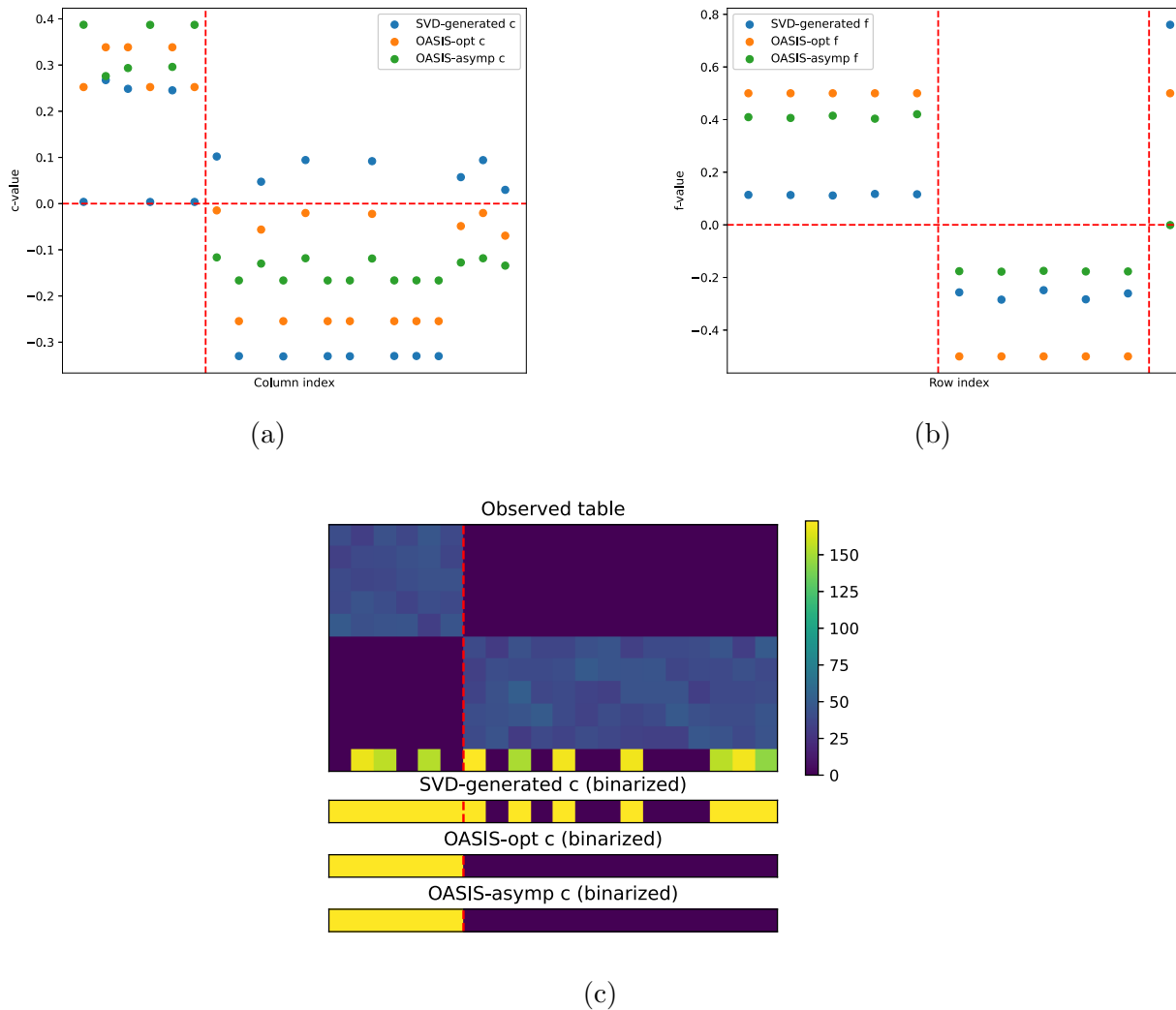


Fig. S8. SVD does not identify the visually clear planted structure in the toy example, while OASIS does. Shown in (c) is the contingency table X with two groups of columns, one which expresses rows 1-5 (the left 6 columns), and the other which expresses rows 6-10 (the right 14 columns). The values in these entries is approximately 40. However, there is a confounding final row, which randomly groups the columns into two new groups, and yields counts either 160 or 0. The last row should be considered noise, as this signal does not match any other row's expression. Multinomial sampling provides slight noise to make this more realistic. Both methods call this table as significant: an extreme example chosen for clarity of exposition. The latent structure is identified by an embedding f which takes value 1 on the first 5 rows, -1 on the next 5, and 0 (i.e., disregarding), the final row. OASIS, both from its asymptotic-optimized c (OASIS-asymp) and its finite-sample bound optimized c (OASIS-opt), perfectly identify the latent clustering. However, an SVD restricts to an ℓ_2 constraint, and since the last row has more counts, it is selected as the dominant component in the SVD-based f . (b) plots the embeddings identified by the two methods, subtracting .5 from OASIS-opt's f , so that it centers around 0. The SVD-based f does identify the row groupings, but utilizes only half the dynamic range for the first 10 rows, assigning the last row a value of 0.76, where all other components have magnitude at most 0.28. OASIS more clearly identifies the row structure, with the asymptotic objective forcing the bottom row to a value less than 9×10^{-4} in magnitude where all other components have a magnitude of at least 0.17. OASIS-opt utilizes the full dynamic range for the first 10 rows. (a) plots c , highlighting that the SVD-identified c cannot identify the latent structure. OASIS-opt identifies the clustering in a well separated manner, and the asymptotic optimized c yields even greater separation (as it fully discounts the last "noisy" row).

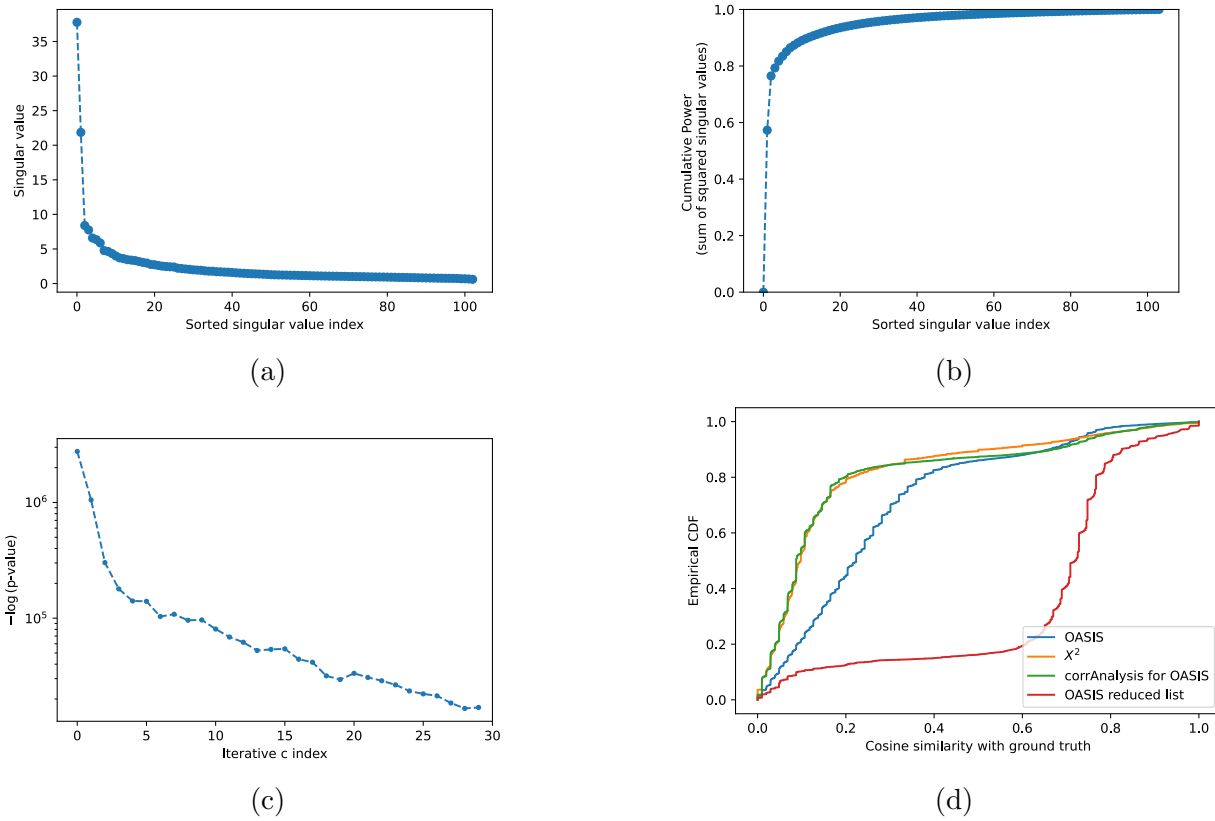


Fig. S9. Analysis of 100,914 tables generated by SPLASH (5) from SARS-CoV-2 data (2). (a,b) show the spectrum of the C matrix from embedding-aggregation for SARS-CoV-2 data, illustrating why a 2D embedding is sufficient. (c) plots the negative log of the p-value bounds for counts-aggregation for SARS-CoV-2. All are 0, up to machine precision, due to the large number of counts. Thus, while for smaller matrices significance testing provides a good cutoff, for larger tables one may still need to find an “elbow” in the curve. As can be seen, after the second index (the elbow) the remaining points follow a linear trend, indicating that only 2 components should be utilized. This threshold can be programmatically identified by looking at the ratio of successive log p-value bounds, and finding the maximum. (d) Plots the empirical CCDF of absolute cosine similarities between the binarized sample embeddings per table (c for OASIS, principal right singular vector from correspondence analysis for X^2) and ground truth vector indicating whether a patient (sample) has Delta or not. Only tables that a method declares significant after multiple hypothesis correction are used. The tables that X^2 rejects do not yield signal that correlates well with the ground truth; for example, the 0.5 and 0.9 quantiles of these empirical distributions are 0.22 and 0.66 for OASIS, as opposed to 0.10, 0.52 for X^2 . This plot additionally shows that the improved predictive power of OASIS with respect to sample metadata is primarily driven by the tables OASIS prioritizes; when performing correspondence analysis on the tables OASIS rejects, similar results are obtained, but when analyzing all calls made by Pearson’s X^2 , many have worse correspondence with sample metadata.

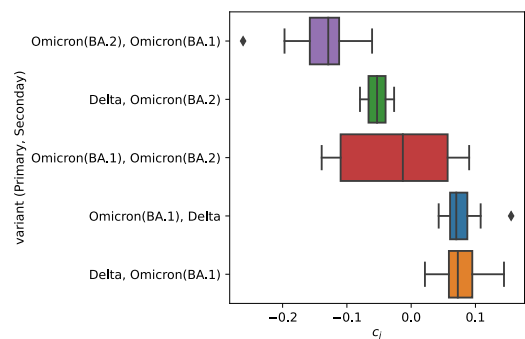
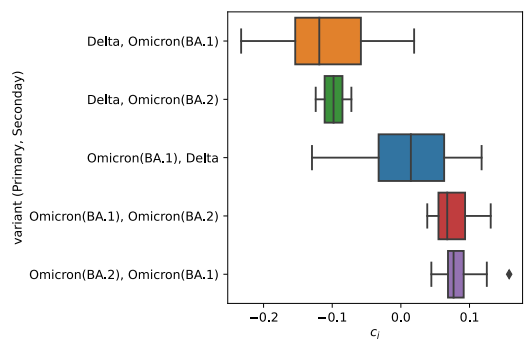


Fig. S10. Top 2 tables in terms of OASIS-opt p-value bound, not called by OASIS-rand. Both tables yield a vector c with a binarized absolute cosine similarity of 0.76 with ground truth metadata (was the patient infected with Delta or not).

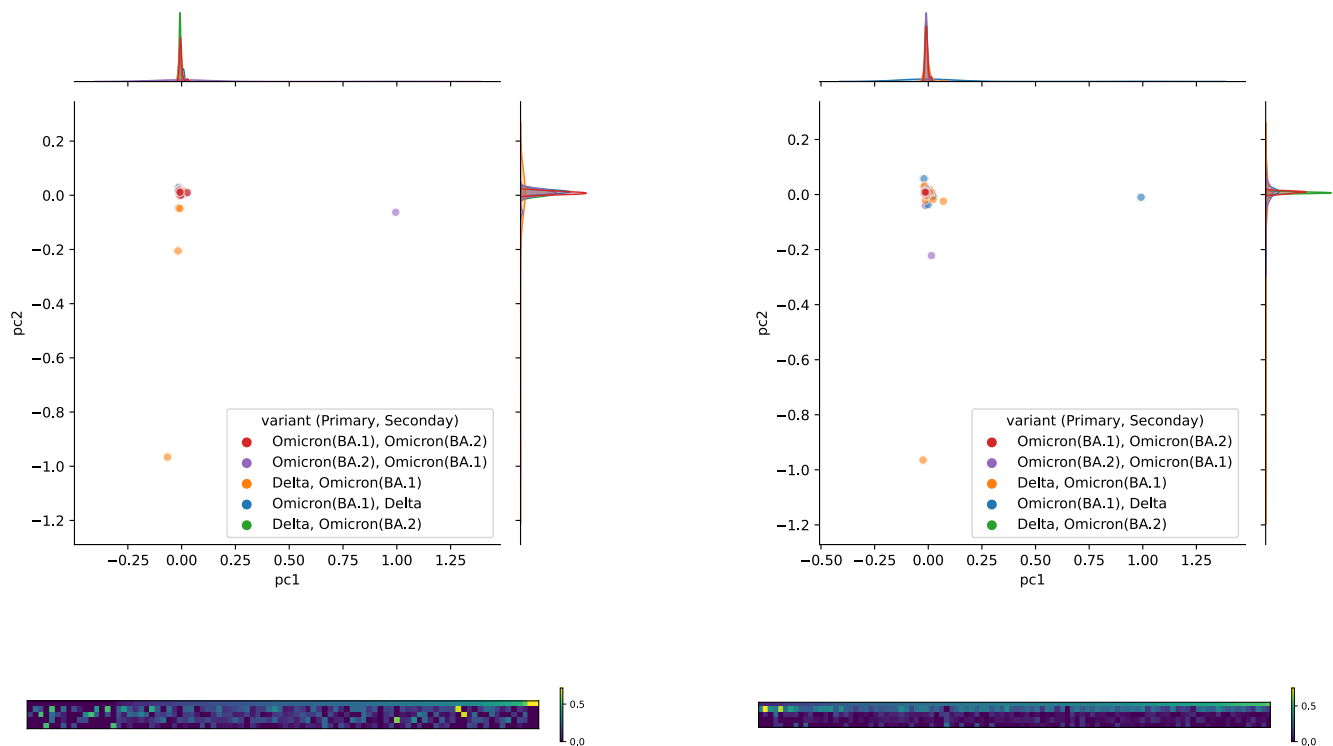


Fig. S11. 2 tables with X^2 q-values of 0, not called by OASIS-opt (q-value > 0.05). Since there were more than 2 tables with a X^2 q-value of 0, we chose the two with the smallest OASIS-opt p-value bounds. The output 1D embedding by correspondence analysis yields an absolute cosine similarity of 0.15 and 0.02 with ground truth (whether a patient has Delta or not). Even when extended to 2 dimensions, the clusters do not separate.

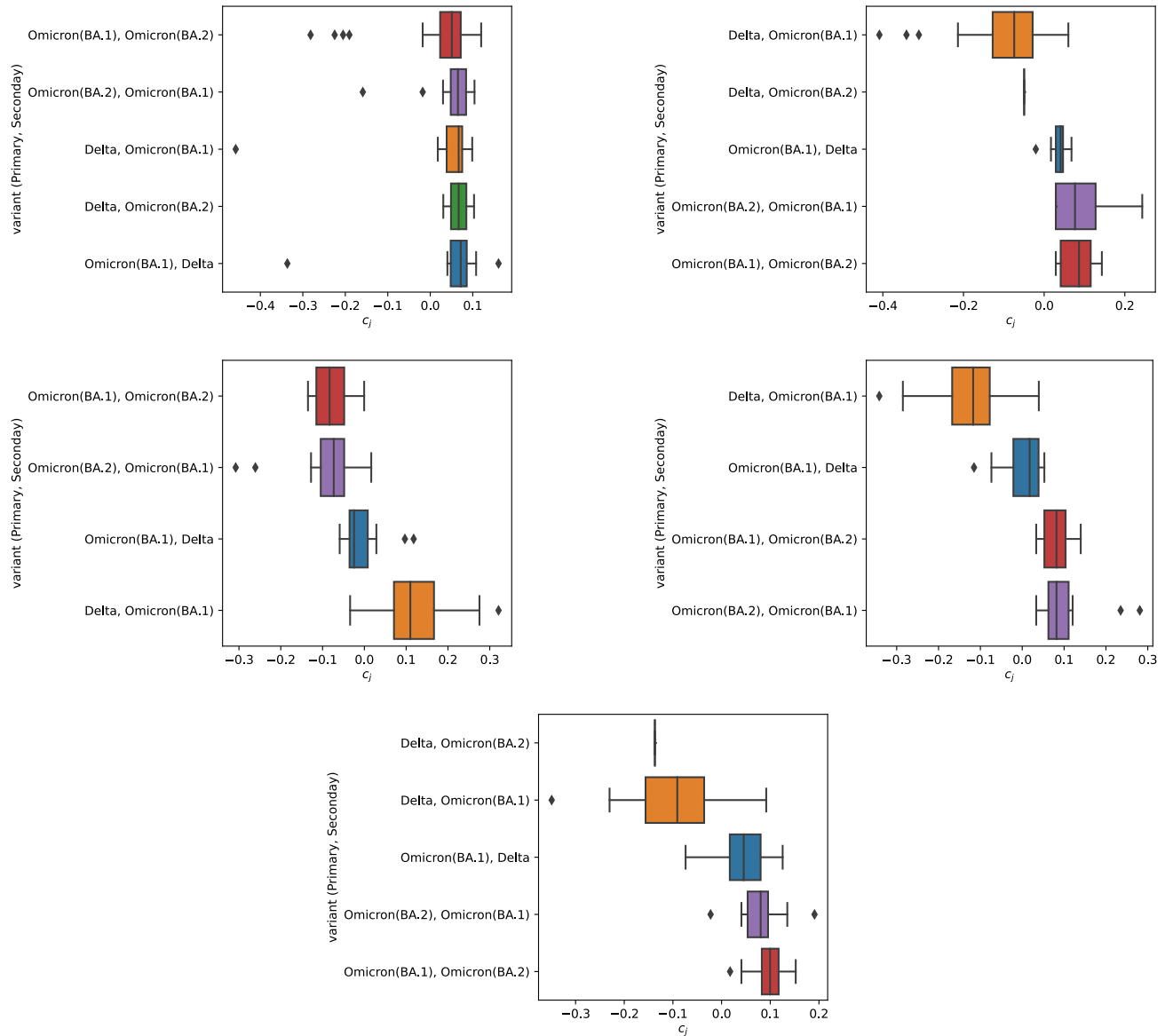


Fig. S12. All 5 tables in OASIS-opt's reduced list of anchors (effect size in the top 10%, $M > 1000$, significant post BY correction OASIS-opt p-value bound), not called by X^2 . First is 482×100 , 3688 counts, cosine similarity of 0.2. Second is 98×77 , with 2540 counts, cosine similarity of 0.61. Third is 74×71 , 1486 counts, cosine similarity of 0.69. Fourth is 67×69 , 1299 counts, cosine similarity of 0.62. Fifth is 88×91 , 1256 counts, 0.63 cosine similarity.

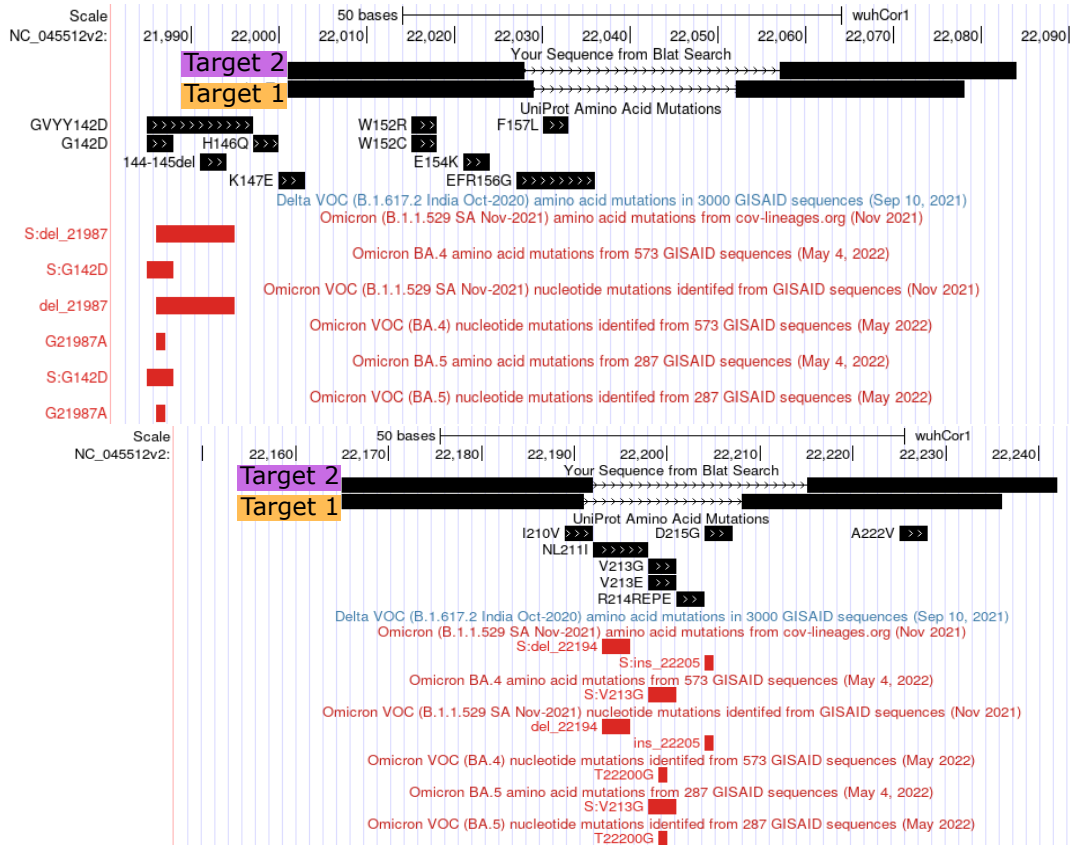
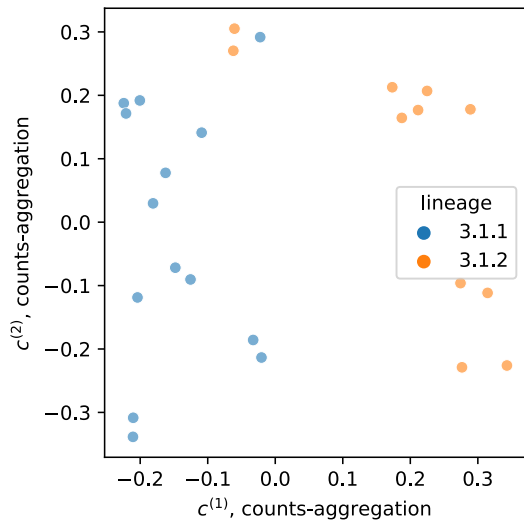
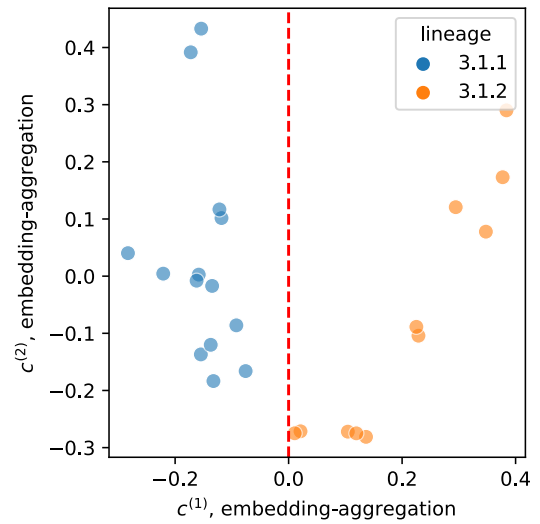


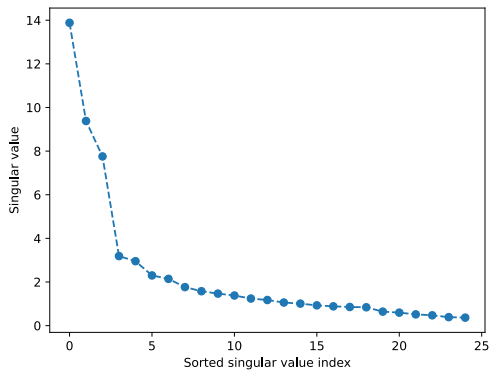
Fig. S13. BLAT (17) to SARS-CoV-2 genome via UCSC genome browser (18). For each target, we align the anchor concatenated with the target, where for each of these two anchors, only 2 targets had abundance greater than 5%. Note that these targets do not come immediately after the anchor, but are taken a fixed distance ahead (called the lookahead distance (5)). In both examples, target 1 takes a value $f_i = 1$, and target 2 takes a value $f_i = 0$. The reported 93% agreement predicts that $f_i = 1$ corresponds to Delta and $f_i = 0$ corresponds to Omicron. Thus, target 1 is predicted to be Delta, and target 2 to be Omicron. Not shown in the first example, analyzing the raw reads from these samples shows that those from Delta samples follow the genome exactly, whereas those from Omicron samples exhibit a 6 basepair deletion in this gap, leading to the resulting contingency table. This deletion is not annotated in the UCSC genome browser, but corresponds to a known Omicron deletion. In the second example, we observe the same behavior (Omicron has a larger gap between anchor and target), this time due to an annotated Variant of Concern (VOC), a deletion.



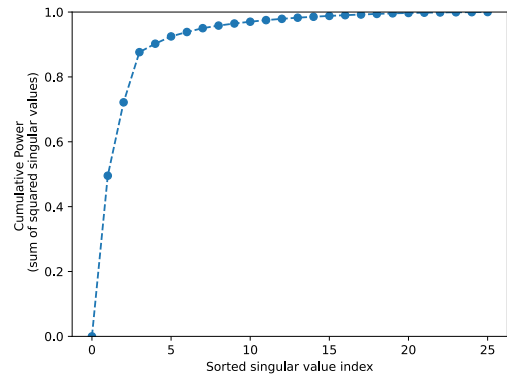
(a)



(b)

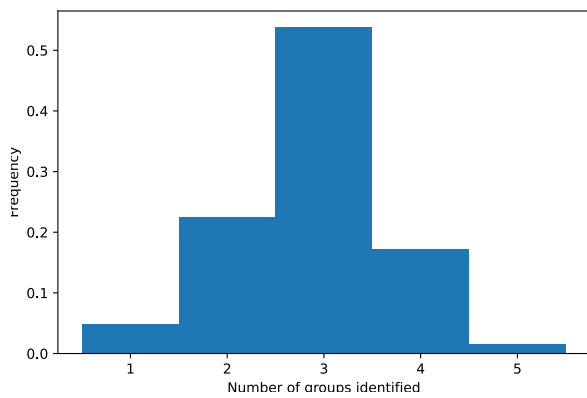


(c)

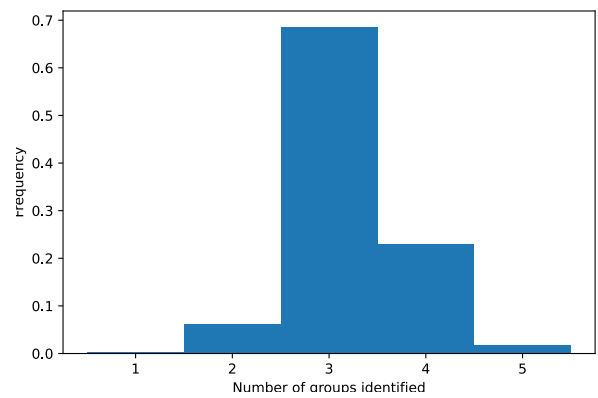


(d)

Fig. S14. Subplots (a,b) show 2D embeddings of *M. tuberculosis* data. The second dimension does not help with sub-sub-lineage classification. Plots (c,d) show the spectrum of the C matrix from embedding-aggregation.



(a)



(b)

Fig. S15. OASIS, unaware of the underlying multi-group structure, is able to give a statistically valid stopping criterion for subclustering. We show evaluation in a planted setting with 3 groups, 10 rows, and 20 columns. Counts per column are Poisson distributed: (a) shows mean 20 observations per column, (b) shows mean 30. We plot the number of identified c vectors using `OASIS-iter`; finding orthogonal $c^{(k)}$ until p-value is no longer significant. OASIS correctly identifies the number of clusters in 55% of simulations, increasing to 69% with mean 30 observations per column. OASIS is only derived as a test on contingency tables, but this simple analysis shows its utility in other applications such as subclustering.

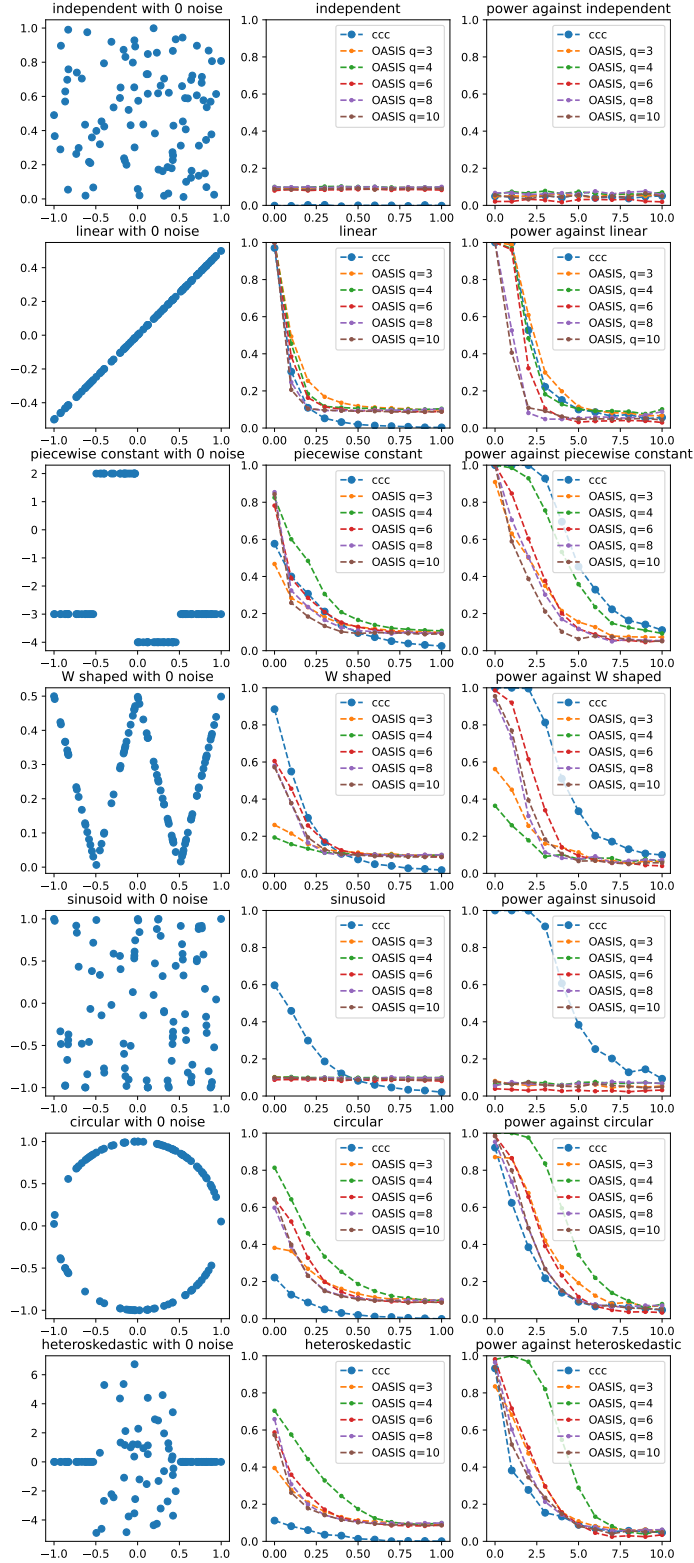


Fig. S16. Comparison of OASIS-opt and xicor (14). OASIS asymptotic p-value used, 25% train test split. Each row represents a different setting, where the left column shows the noiseless problem instance $y = f(x)$, the middle column shows the magnitude of the correlation coefficient averaged over 1000 iterations per point, and the right column shows the fraction of time the null was rejected. For OASIS, the middle column shows the effect size, and the right column shows the rejection fraction utilizing the asymptotic p-value. The quantization level q , which maps the input continuous-valued random variable to a discrete categorical one, is varied across trials. This is performed independently for the row and column random variables, and binned into q bins of equal counts ($1/q$ quantiles).

num rows	num cols	X^2	OASIS-rand	OASIS-opt
5	10	0.370	0.005	0.007
	50	0.765	0.009	0.006
	400	0.995	0.001	0.057
10	10	0.792	0.011	0.015
	50	1.000	0.007	0.046
	400	1.000	0.010	0.309
20	10	0.999	0.065	0.080
	50	1.000	0.058	0.213
	400	1.000	0.100	0.866
100	10	1.000	0.871	0.978
	50	1.000	0.983	1.000
	400	1.000	0.994	1.000

Table S1. Power at $\lambda \approx 129$, noise model as in (16). Exponentially distributed target distribution, described in text. Full plots in Figure S3.

number of rows (J)	SDP+rounding	Alternating maximization	Asymptotically-optimal (eigenvector + rounding)
12	14.5	0.86	0.20
24	194	1.2	0.25
48	-	1.26	0.59
96	-	1.31	5.9
192	-	1.48	16.7

Table S2. Computational time (in seconds) for computing \mathbf{c}, \mathbf{f} using the stated method, on 1000 matrices. As can be seen, the SDP solver quickly becomes infeasible. Forming the matrix $\tilde{X}\tilde{X}^\top$ also becomes costly as the number of rows increases. Alternating maximization stays efficient and performant. Timing results generated on a MacBook Air (M2, 2022), runs declared as a failure after 5 minutes.

441 **References**

- 442 1. V Dreyer, et al., High fluoroquinolone resistance proportions among multidrug-resistant tuberculosis driven by dominant
443 H2 mycobacterium tuberculosis clones in the mumbai metropolitan region. *Genome Medicine* **14**, 95 (2022).
- 444 2. A Bal, et al., Detection and prevalence of sars-cov-2 co-infections during the omicron variant circulation in france. *Nat.*
445 *Commun.* **13**, 1–9 (2022).
- 446 3. R Leinonen, H Sugawara, M Shumway, INSD Collaboration, The sequence read archive. *Nucleic acids research* **39**,
447 D19–D21 (2010).
- 448 4. M Kokot, R Dehghannasiri, TZ Baharav, J Salzman, S Deorowicz, Splash2 provides ultra-efficient, scalable, and
449 unsupervised discovery on raw sequencing reads. *bioRxiv* pp. 2023–03 (2023).
- 450 5. K Chaung, et al., Splash: a statistical, reference-free genomic algorithm unifies biological discovery. *Cell* **186**, 5440–5456
451 (2023).
- 452 6. MJ Wainwright, *High-dimensional statistics: A non-asymptotic viewpoint*. (Cambridge University Press) Vol. 48, (2019).
- 453 7. P Billingsley, *Probability and measure*. (John Wiley & Sons), (2008).
- 454 8. AC Berry, The accuracy of the gaussian approximation to the sum of independent variates. *Transactions american*
455 *mathematical society* **49**, 122–136 (1941).
- 456 9. F Chen, S Roch, K Rohe, S Yu, Estimating graph dimension with cross-validated eigenvalues. *arXiv preprint*
457 *arXiv:2108.03336* (2021).
- 458 10. A Neufeld, J Popp, LL Gao, A Battle, D Witten, Negative binomial count splitting for single-cell rna sequencing data.
459 *arXiv preprint arXiv:2307.12985* (2023).
- 460 11. MX Goemans, DP Williamson, Improved approximation algorithms for maximum cut and satisfiability problems using
461 semidefinite programming. *J. ACM (JACM)* **42**, 1115–1145 (1995).
- 462 12. S Diamond, S Boyd, CVXPY: A Python-embedded modeling language for convex optimization. *J. Mach. Learn. Res.* **17**,
463 1–5 (2016).
- 464 13. B Langmead, C Trapnell, M Pop, SL Salzberg, Ultrafast and memory-efficient alignment of short dna sequences to the
465 human genome. *Genome biology* **10**, 1–10 (2009).
- 466 14. S Chatterjee, A new coefficient of correlation. *J. Am. Stat. Assoc.* **116**, 2009–2022 (2021).
- 467 15. A Maurer, M Pontil, Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*
468 (2009).
- 469 16. MI Love, W Huber, S Anders, Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome*
470 *biology* **15**, 1–21 (2014).
- 471 17. WJ Kent, Blat—the blast-like alignment tool. *Genome research* **12**, 656–664 (2002).
- 472 18. WJ Kent, et al., The human genome browser at ucsc. *Genome research* **12**, 996–1006 (2002).