

Supplementary Materials

March 20, 2024

A Derivation of constraint set B

Let r_C, r_{J_1} and r_{I_1} be the proportion of the top scores with missing second score that come from C, J_1 and I_1 , respectively, where, for $X \in \mathbf{X} = \{C, J_1, I_1\}$, $r_X > 0$ and $\sum_X r_X = 1$. The proportion of top scores coming from X and having a second score is $w_X - r_X v_\Phi$. The proportion of observed second scores coming from Y ($Y \in \mathbf{Y} = \{C, J_1, I_1, J_2, J_1\}$) among all observed and missing second scores is $(1 - v_\Phi)v_Y$. Updating the constraint set A by replacing all occurrences of w_X and v_Y by $w_X - r_X v_\Phi$ and $(1 - v_\Phi)v_Y$ we get an updated constraint set based only on the spectrum where both the scores are available. For example, $w_C \leq v_{J_1} + v_{I_1}$ is replaced by $w_C - r_C v_\Phi \leq (1 - v_\Phi)(v_{J_1} + v_{I_1})$ and $v_C \leq w_{J_1} + w_{I_1}$ is replaced by $(1 - v_\Phi)v_C \leq w_{J_1} + w_{I_1} - (r_{J_1} + r_{I_1})v_\Phi$. Since r_X is an unknown quantity, we further modify the new constraints by using $r_X = 1$ ($r_X = 0$) when it appears on the left (right) hand side, since it minimizes (maximizes) $w_X - r_X v_\Phi$, making the constraints looser and valid without any assumptions on the distribution of top scores with missing second score. The two constraints in the example above get transformed to $w_C - v_\Phi \leq (1 - v_\Phi)(v_{J_1} + v_{I_1})$ and $(1 - v_\Phi)v_C \leq w_{J_1} + w_{I_1}$. All the constraints modified in this manner are given as constraint set B with $v'_Y = (1 - v_\Phi)v_Y$.

B Monotonicity of the ECM algorithm and the binary search

The ECM algorithm based component parameter updates are derived by optimizing the $\mathcal{Q}(\zeta|\bar{\zeta})$ one parameter at a time. When optimizing $\mathcal{Q}(\zeta|\bar{\zeta})$ w.r.t. component Y parameters, $\theta_Y = \{\mu_Y, \Delta_Y, \Gamma_Y\}$, the other components' parameters and the weight parameters in ζ can be ignored as they appear in additive terms, constant w.r.t. θ_Y . However, note that all parameters in the current parameter set ζ still play a role in component Y parameter updates. In effect optimizing the \mathcal{Q} -function w.r.t. component Y parameters reduces to the optimization of $\mathcal{Q}(\mu_Y, \Delta_Y, \Gamma_Y|\bar{\zeta})$ as defined below. For $Y \in \{C, J_1, I_1\}$,

$$\begin{aligned} \mathcal{Q}(\mu_Y, \Delta_Y, \Gamma_Y|\bar{\zeta}) &= -\frac{1}{2} \sum_{s_1 \in \mathbb{S}_1} \bar{w}_X(s_1) \left(\log \Gamma_Y + \frac{(s_1 - \mu_Y)^2 - 2\xi_1(s_1, \bar{\zeta})(s_1 - \mu_Y)\Delta_Y + \xi_2(s_1, \bar{\zeta})(\Delta_Y^2 + \Gamma_Y)}{\Gamma_Y} \right) \\ &\quad - \frac{1}{2} \sum_{s_2 \in \mathbb{S}_2} \bar{v}_Y(s_2) \left(\log \Gamma_Y + \frac{(s_2 - \mu_Y)^2 - 2\xi_1(s_2, \bar{\zeta})(s_2 - \mu_Y)\Delta_Y + \xi_2(s_2, \bar{\zeta})(\Delta_Y^2 + \Gamma_Y)}{\Gamma_Y} \right). \end{aligned}$$

And for $Y \in \{J_2, I_2\}$,

$$\mathcal{Q}(\mu_Y, \Delta_Y, \Gamma_Y|\bar{\zeta}) = -\frac{1}{2} \sum_{s_2 \in \mathbb{S}_2} \bar{v}_Y(s_2) \left(\log \Gamma_Y + \frac{(s_2 - \mu_Y)^2 - 2\xi_1(s_2, \bar{\zeta})(s_2 - \mu_Y)\Delta_Y + \xi_2(s_2, \bar{\zeta})(\Delta_Y^2 + \Gamma_Y)}{\Gamma_Y} \right).$$

The new parameters $\ddot{\mu}_Y, \ddot{\Delta}_Y$ and $\ddot{\Gamma}_Y$ from the ECM parameter update equations are obtained as the unique stationary points (where the first derivative is 0) of $\mathcal{Q}(\mu_Y, \Delta_Y, \Gamma_Y|\bar{\zeta})$, $\mathcal{Q}(\ddot{\mu}_Y, \Delta_Y, \Gamma_Y|\bar{\zeta})$ and $\mathcal{Q}(\mu_Y, \ddot{\Delta}_Y, \Gamma_Y|\bar{\zeta})$ as a single variable function of μ_Y, Δ_Y and Γ_Y , respectively.

$\mathcal{Q}(\mu_Y, \ddot{\Delta}_Y, \Gamma_Y|\bar{\zeta})$ and $\mathcal{Q}(\ddot{\mu}_Y, \Delta_Y, \Gamma_Y|\bar{\zeta})$ as functions μ_Y and Δ_Y , respectively, are concave, since their second derivatives are not positive everywhere. Consequently, $\ddot{\mu}_Y$ and $\ddot{\Delta}_Y$ are the global maximizers.

Now $\mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \Gamma_Y | \bar{\zeta})$ is not a concave function of Γ_Y , however, it can be expressed as a difference of two convex functions as follows.

$$\begin{aligned} \mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \Gamma_Y | \bar{\zeta}) &= -\frac{1}{2} \sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_X(s_1) \left(\log \Gamma_Y + \frac{(s_1 - \ddot{\mu}_Y)^2 - 2\xi_1(s_1, \bar{\zeta})(s_1 - \ddot{\mu}_Y)\ddot{\Delta}_Y + \xi_2(s_1, \bar{\zeta})(\ddot{\Delta}_Y^2 + \Gamma_Y)}{\Gamma_Y} \right) \\ &\quad - \frac{1}{2} \sum_{s_2 \in \mathbb{S}_2} \bar{\nu}_Y(s_2) \left(\log \Gamma_Y + \frac{(s_2 - \ddot{\mu}_Y)^2 - 2\xi_1(s_2, \bar{\zeta})(s_2 - \ddot{\mu}_Y)\ddot{\Delta}_Y + \xi_2(s_2, \bar{\zeta})(\ddot{\Delta}_Y^2 + \Gamma_Y)}{\Gamma_Y} \right) \\ &= \frac{1}{2} \left(\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_X(s_1) + \sum_{s_2 \in \mathbb{S}_2} \bar{\nu}_Y(s_2) \right) (\log 1/\Gamma_Y - \ddot{\Gamma}_Y/\Gamma_Y) + \text{constant}. \end{aligned}$$

In the equation above $g(\Gamma_Y) = \ddot{\Gamma}_Y/\Gamma_Y$ and $h(\Gamma_Y) = \log 1/\Gamma_Y$ are convex functions, since their second derivatives are non-negative everywhere. Now based on Definition 2.2.2 in [1], the ϵ -subdifferential of g and h at $\ddot{\Gamma}_Y$ is given by

$$\begin{aligned} \partial_\epsilon g(\ddot{\Gamma}_Y) &= \begin{cases} \left(-\infty, -\frac{1}{\ddot{\Gamma}_Y} + \frac{\epsilon}{(\Gamma_Y - \ddot{\Gamma}_Y)} \right] & \text{when } \Gamma_Y \geq \ddot{\Gamma}_Y \\ \left[-\frac{1}{\ddot{\Gamma}_Y} + \frac{\epsilon}{(\Gamma_Y - \ddot{\Gamma}_Y)}, \infty \right) & \text{when } \Gamma_Y < \ddot{\Gamma}_Y \end{cases} \\ \partial_\epsilon h(\ddot{\Gamma}_Y) &= \begin{cases} \left(-\infty, \frac{\log(\ddot{\Gamma}_Y/\Gamma_Y)}{(\Gamma_Y - \ddot{\Gamma}_Y)} + \frac{\epsilon}{(\Gamma_Y - \ddot{\Gamma}_Y)} \right] & \text{when } \Gamma_Y \geq \ddot{\Gamma}_Y \\ \left[\frac{\log(\ddot{\Gamma}_Y/\Gamma_Y)}{(\Gamma_Y - \ddot{\Gamma}_Y)} + \frac{\epsilon}{(\Gamma_Y - \ddot{\Gamma}_Y)}, \infty \right) & \text{when } \Gamma_Y < \ddot{\Gamma}_Y \end{cases} \end{aligned}$$

Now since $\log(\ddot{\Gamma}_Y/\Gamma_Y) \leq \frac{\ddot{\Gamma}_Y}{\Gamma_Y} - 1$,

$$\begin{aligned} \frac{\log(\ddot{\Gamma}_Y/\Gamma_Y)}{\Gamma_Y - \ddot{\Gamma}_Y} &\leq -\frac{1}{\Gamma_Y} \quad \text{when } \Gamma_Y \geq \ddot{\Gamma}_Y \\ \frac{\log(\ddot{\Gamma}_Y/\Gamma_Y)}{\Gamma_Y - \ddot{\Gamma}_Y} &\geq -\frac{1}{\Gamma_Y} \quad \text{when } \Gamma_Y < \ddot{\Gamma}_Y. \end{aligned}$$

Thus $\partial_\epsilon h(\ddot{\Gamma}_Y) \subseteq \partial_\epsilon g(\ddot{\Gamma}_Y)$ and, consequently from Theorem 2.3.1 in [1], $g(\Gamma_Y) - h(\Gamma_Y)$ has a global minimizer at $\ddot{\Gamma}_Y$. It follows that $h(\Gamma_Y) - g(\Gamma_Y)$ and consequently $\mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \Gamma_Y | \bar{\zeta})$, has a global maximizer at $\ddot{\Gamma}_Y$.

Since $\ddot{\mu}_Y, \ddot{\Delta}_Y$ and $\ddot{\Gamma}_Y$ are global maximizers of $\mathcal{Q}(\mu_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta})$, $\mathcal{Q}(\ddot{\mu}_Y, \Delta_Y, \bar{\Gamma}_Y | \bar{\zeta})$ and $\mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \Gamma_Y | \bar{\zeta})$, respectively,

$$\begin{aligned} \mathcal{Q}(\bar{\mu}_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) &\leq \mathcal{Q}(\ddot{\mu}_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) \\ &\leq \mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) \\ &\leq \mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \ddot{\Gamma}_Y | \bar{\zeta}). \end{aligned}$$

Thus the new component parameters from the ECM updates increase (precisely, do not decrease) the value of the \mathcal{Q} -function at each iteration.

The monotonicity of the \mathcal{Q} -function can be similarly established for the parameter updates from the binary search. $\mathcal{Q}(\mu_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta})$, $\mathcal{Q}(\ddot{\mu}_Y, \Delta_Y, \bar{\Gamma}_Y | \bar{\zeta})$ and $\mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \Gamma_Y | \bar{\zeta})$ each have a unique stationary point (where the first derivative is 0), $\ddot{\mu}_Y, \ddot{\Delta}_Y$ and $\ddot{\Gamma}_Y$, respectively, which is the global maximizer. Furthermore, they do not have any singular points (where the first derivative is undefined), except for $\mathcal{Q}(\ddot{\mu}_Y, \ddot{\Delta}_Y, \Gamma_Y | \bar{\zeta})$ at $\Gamma_Y = 0$, which is an endpoint. Thus none of them have a local minimum and an absolute minimum only exists at the endpoints. It follows that the three functions are unimodal with $\ddot{\mu}_Y, \ddot{\Delta}_Y$ and $\ddot{\Gamma}_Y$ giving the modes. Thus any point on the line connecting $\bar{\mu}_Y$ to $\ddot{\mu}_Y$, $\bar{\Delta}_Y$ to $\ddot{\Delta}_Y$ and $\bar{\Gamma}_Y$ to $\ddot{\Gamma}_Y$ would not decrease the value of the corresponding \mathcal{Q} -function. Precisely, in case of μ_Y , binary search finds a feasible $\hat{\mu}_Y$ on the line connecting $\bar{\mu}_Y$ and $\ddot{\mu}_Y = \arg\max_{\mu_Y} \mathcal{Q}(\mu_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta})$. Thus $\mathcal{Q}(\bar{\mu}_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) \leq \mathcal{Q}(\hat{\mu}_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta})$. Now to update Δ_Y , binary search finds a feasible $\hat{\Delta}_Y$ on the line connecting $\bar{\Delta}_Y$ and $\ddot{\Delta}_Y = \arg\max_{\Delta_Y} \mathcal{Q}(\hat{\mu}_Y, \Delta_Y, \bar{\Gamma}_Y | \bar{\zeta})$. Thus

$\mathcal{Q}(\hat{\mu}_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) \leq \mathcal{Q}(\hat{\mu}_Y, \hat{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta})$. Finally, to update Γ_Y , binary search finds a feasible $\hat{\Gamma}_Y$ on the line connecting $\bar{\Gamma}_Y$ and $\check{\Gamma}_Y = \operatorname{argmax}_{\Gamma_Y} \mathcal{Q}(\hat{\mu}_Y, \hat{\Delta}_Y, \Gamma_Y | \bar{\zeta})$. Thus $\mathcal{Q}(\hat{\mu}_Y, \hat{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) \leq \mathcal{Q}(\hat{\mu}_Y, \hat{\Delta}_Y, \hat{\Gamma}_Y | \bar{\zeta})$. Thus, in summary, $\mathcal{Q}(\bar{\mu}_Y, \bar{\Delta}_Y, \bar{\Gamma}_Y | \bar{\zeta}) \leq \mathcal{Q}(\hat{\mu}_Y, \hat{\Delta}_Y, \hat{\Gamma}_Y | \bar{\zeta})$.

C Mode of SN distribution

The mode of $\text{SN}(\mu, \sigma, \lambda)$ is given by $\mu + \sigma m_0(\lambda)$, where $m_0(\lambda) = \sqrt{2/\pi} \delta - (1 - \pi/4) \frac{(\sqrt{2/\pi} \delta)^3}{1 - (2/\pi) \delta^2} - \frac{\operatorname{sign}(\lambda)}{2} \exp(2\pi/|\lambda|)$

D Evaluating a pairwise density constraint

To evaluate a pairwise density constraint, say $f_{\text{SN}}(\cdot; \theta_1) \succ f_{\text{SN}}(\cdot; \theta_2)$, first we compare their modes. If $\operatorname{mode}(f_{\text{SN}}(\cdot; \theta_1)) \leq \operatorname{mode}(f_{\text{SN}}(\cdot; \theta_2))$, then the constraint is not satisfied. If $\operatorname{mode}(f_{\text{SN}}(\cdot; \theta_1)) > \operatorname{mode}(f_{\text{SN}}(\cdot; \theta_2))$, then we test the other two conditions. For a fast comparison, we generate a set, \mathbb{S}_{grid} , of 200 equally spaced points between the maximum and minimum of the top scores to limit the number of points for density comparison. Then we test $f_{\text{SN}}(s; \theta_1) > f_{\text{SN}}(s; \theta_2)$ ($f_{\text{SN}}(s; \theta_1) < f_{\text{SN}}(s; \theta_2)$) at all points $s \in \mathbb{S}_{\text{grid}}$ that are above (below) $\operatorname{mode}(f_{\text{SN}}(\cdot; \theta_1))$ ($\operatorname{mode}(f_{\text{SN}}(\cdot; \theta_2))$). If it is true, we consider the constraint as satisfied, otherwise it is not satisfied.

E Binary search for density constraints

In Algorithm S1, the `densityConstraintsSatisfied`($Y, \alpha_Y, \bar{\zeta}$) function is used to check if all pairwise constraints for component Y would still be satisfied if one of its three parameters in $\bar{\zeta}$ is updated to a new value α_Y . The pairwise density constraints are evaluated as described in Section D.

Algorithm S1 Binary search for density constraints.

Require: 1) $Y \in \mathbf{Y}$: the component whose parameters are being updated, 2) $\check{\alpha}_Y \in \{\check{\mu}_Y, \check{\Delta}_Y, \check{\Gamma}_Y\}$: the new, possibly infeasible, estimate of one of the three SN parameters of Y from the parameter update equations in Section 4.3.2, and 3) $\bar{\zeta}$: the current feasible estimate of all parameters. $\bar{\alpha}_Y$, the current estimate of the same parameter as $\check{\alpha}_Y$, is contained in $\bar{\zeta}$.

Ensure: $\hat{\alpha}_Y$: a feasible point on the line segment connecting $\bar{\alpha}_Y$ to $\check{\alpha}_Y$, as close to $\check{\alpha}_Y$ as possible.

```

1: if densityConstraintsSatisfied( $Y, \check{\alpha}_Y, \bar{\zeta}$ ) then
2:   ▷ All density constraints for component Y are satisfied when  $\check{\alpha}_Y$  replaces the corresponding parameter
   in  $\bar{\zeta}$ .
3:    $\hat{\alpha}_Y \leftarrow \check{\alpha}_Y$ 
4: else
5:   while  $\frac{|\check{\alpha}_Y - \bar{\alpha}_Y|}{\bar{\alpha}_Y} > 10^{-4}$  do
6:      $\hat{\alpha}_Y \leftarrow \frac{\bar{\alpha}_Y + \check{\alpha}_Y}{2}$ 
7:     if densityConstraintsSatisfied( $Y, \hat{\alpha}_Y, \bar{\zeta}$ ) then
8:        $\bar{\alpha}_Y \leftarrow \hat{\alpha}_Y$ 
9:     else
10:       $\check{\alpha}_Y \leftarrow \hat{\alpha}_Y$ 
11:    end if
12:  end while
13:   $\hat{\alpha}_Y \leftarrow \bar{\alpha}_Y$ 
14: end if
15: return  $\hat{\alpha}_Y$ 

```

Note that if the current parameter, $\bar{\alpha}$, is feasible, the binary search is guaranteed to give a feasible solution as it can find a value arbitrarily close to $\bar{\alpha}$. Consequently, if the the first set of parameters are feasible, the approach is guaranteed to give a feasible solution at each iteration.

F One-sample Model Update Rules

Below are the update rules for the parameters of the component skew normal distribution.

$$\begin{aligned}
 w_Y &= \frac{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1)}{|\mathbb{S}_1|} \\
 \ddot{\mu}_Y &= \frac{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1) \bar{m}_Y(s_1, \bar{\Delta}_Y)}{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1)}, \\
 \ddot{\Delta}_Y &= \frac{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1) \bar{d}_Y(s_1, \ddot{\mu}_Y)}{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1) \xi_2(s_1, \bar{\theta}_Y)}, \\
 \ddot{\Gamma}_Y &= \frac{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1) \bar{g}_Y(s_1, \ddot{\mu}_Y, \ddot{\Delta}_Y)}{\sum_{s_1 \in \mathbb{S}_1} \bar{\omega}_Y(s_1)},
 \end{aligned}$$

G Code Availability

<https://github.com/shawn-peng/xlms>

| | Search time w/ decoys | Search time w/o decoys | Model run time | Target database size |
|-----------|-----------------------|------------------------|----------------|----------------------|
| ALott | 35 h 13 m 10.0 s | 20 h 22 m 12.0 s | 6 m 40.7 s | 4 |
| Alinden | 192 h 10 m 52.0 s | 70 h 6 m 5.0 s | 8 m 17.6 s | 79 |
| CPSF | 2 h 0 m 4.0 s | 37 m 5.3 s | 6 m 5.2 s | 28 |
| D1810 | 18 m 34.4 s | 10 m 9.3 s | 3 m 49.8 s | 10 |
| MS2000225 | 21 m 1.4 s | 7 m 19.7 s | 4 m 6.6 s | 5 |
| QE | 8 m 15.2 s | 5 m 14.5 s | 4 m 18.1 s | 2 |
| RPA | 1 h 37 m 57.3 s | 25 m 26.7 s | 3 m 30.3 s | 480 |
| Alban | 25 m 20.3 s | 7 m 48.8 s | 5 m 43.3 s | 20 |
| Ecoli | 37 m 26.5 s | 12 m 16.1 s | 11 m 0.1 s | 56 |
| Peplib | 3 m 41.9 s | 2 m 15.3 s | 4 m 3.4 s | 9 |

Table S1: Comparison of running time for the database search and the 2SMix model. The 2SMix run time is reported as an average per restart. The experiments were executed on AMD EPYC 7452 32-Core Processor @ 2.345 GHz with 251GB RAM for single thread time.

| | self-links | cross-links | loop-links | mono-links | total | self-links ratio |
|-----------|------------|-------------|------------|------------|-------|------------------|
| ALott | 7998 | 20500 | 659 | 3758 | 24917 | 0.32 |
| Alinden | 1277 | 2406 | 423 | 2647 | 5476 | 0.23 |
| CPSF | 414 | 633 | 126 | 804 | 1563 | 0.26 |
| D1810 | 7 | 8 | 1 | 61 | 70 | 0.10 |
| MS2000225 | 248 | 410 | 20 | 407 | 837 | 0.30 |
| QE | 0 | 0 | 0 | 1 | 1 | 0.00 |
| RPA | 0 | 16 | 9 | 57 | 82 | 0.00 |
| Alban | 84 | 197 | 21 | 142 | 360 | 0.23 |
| Ecoli | 112 | 168 | 28 | 617 | 813 | 0.14 |
| Peplib | 1662 | 1663 | 0 | 858 | 2521 | 0.66 |

Table S2: Counts of PSMs of each cross-linking type. Self-links are pairs of peptides that come from the same protein. Cross-links are pairs of peptides that come from the same protein or from different proteins. The total is the sum of cross-links, loop-links and mono-links. The self-links ratio is self-links/total.

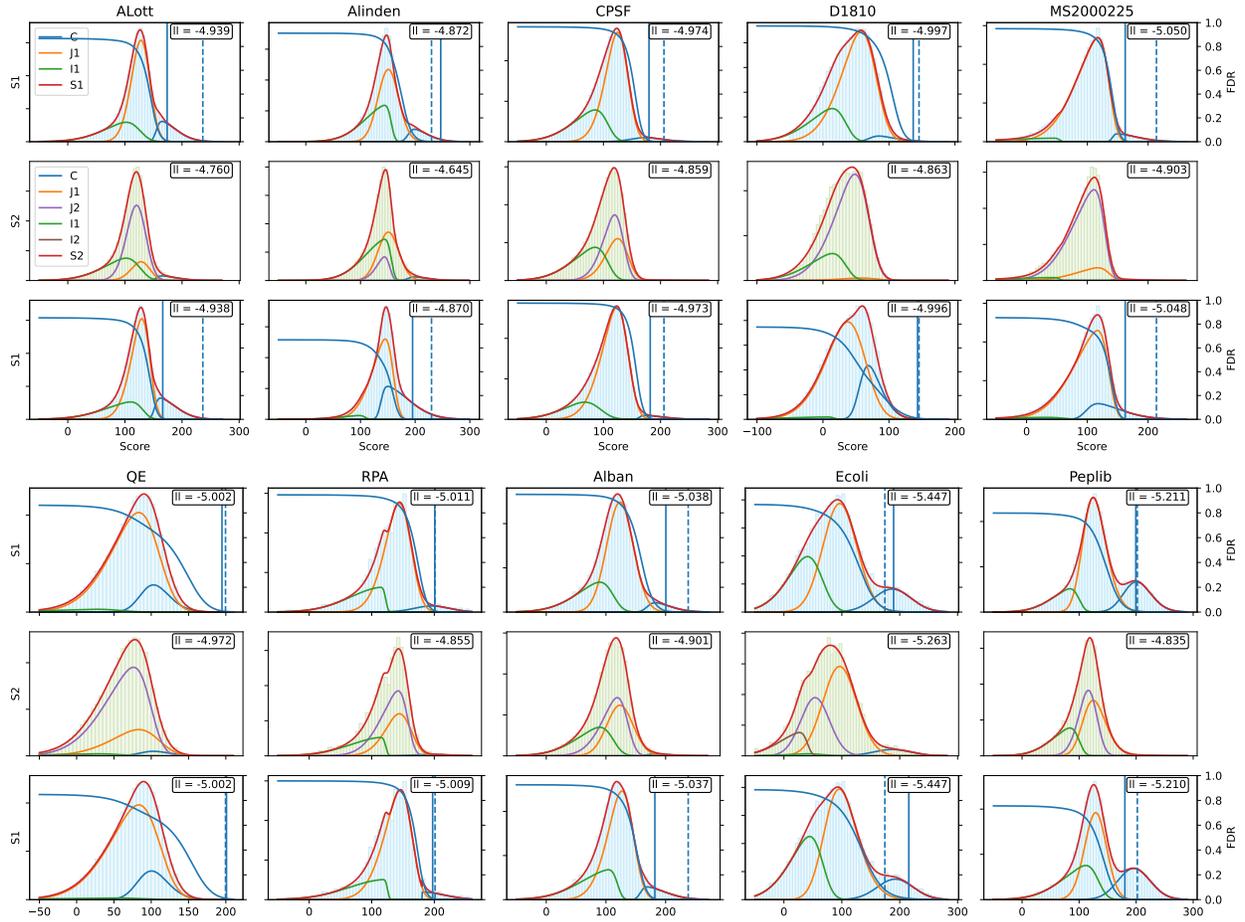


Figure S1: 1SMix and 2SMix model fits. 1SMix and 2SMix are the 1-sample and 2-sample mixture model respectively. The top two rows show the 2SMix fit on the top scores and the second scores. The third row shows 1SMix fit on the top scores. The empirical distributions of the top scores and the second scores are shown as histograms. The component densities estimated by the models are displayed. Each component density is weighted by its estimated mixture weight. The vertical solid line shows the 1% FDR threshold given by the mixture model, while the vertical dashed line shows the 1% FDR threshold by TDA for comparison. The average log-likelihood of the fitted model is shown in the upper right corner. For 2SMix the log-likelihood is computed separately for the top scores and the second scores.

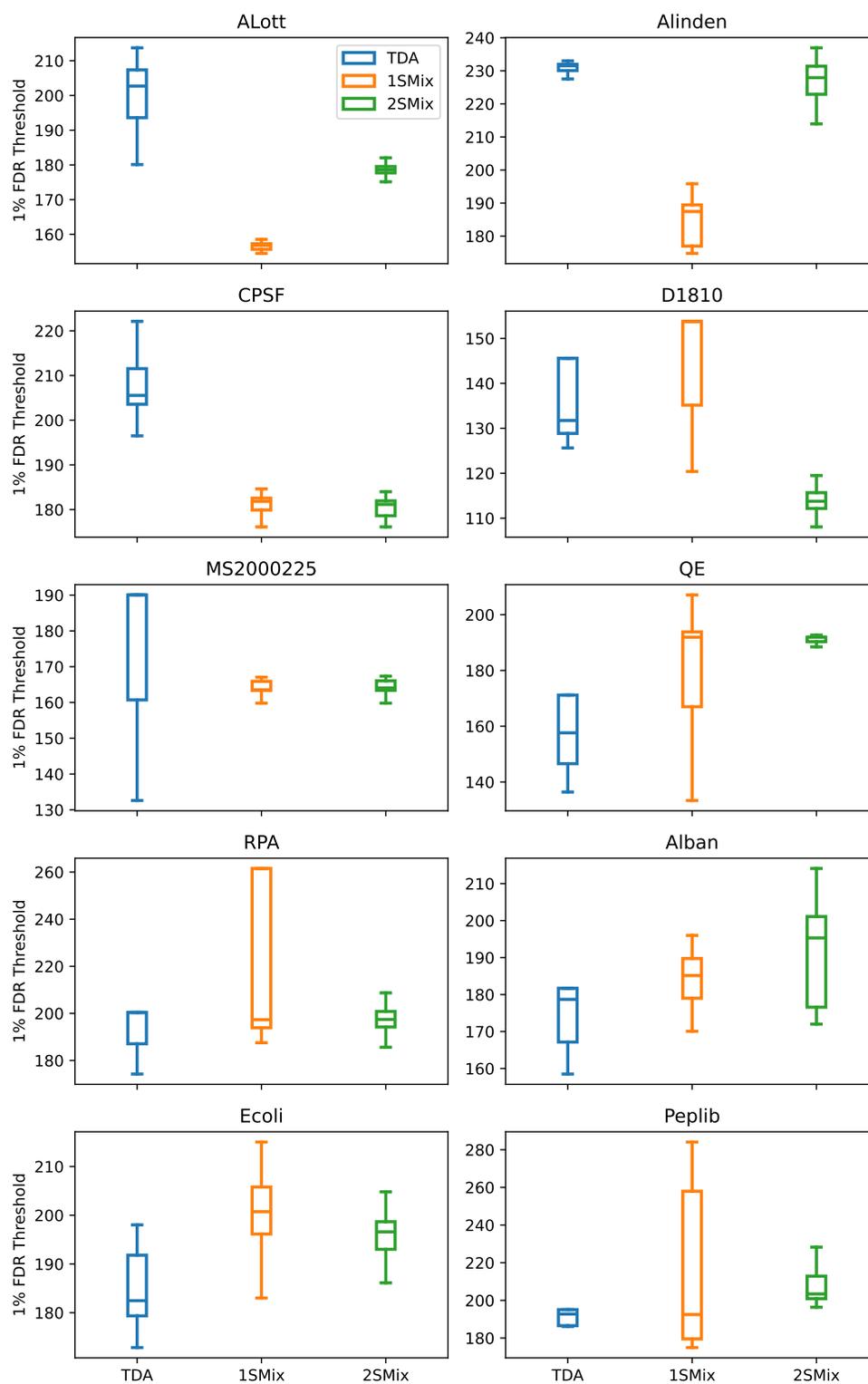


Figure S2: Stability of the 1% FDR threshold estimation. The stability of the methods was compared using 50 bootstrap samples, on which the 1% FDR thresholds were estimated. The larger spread of a boxplot indicates lower stability. Bootstrapping was performed on the set of original spectra.

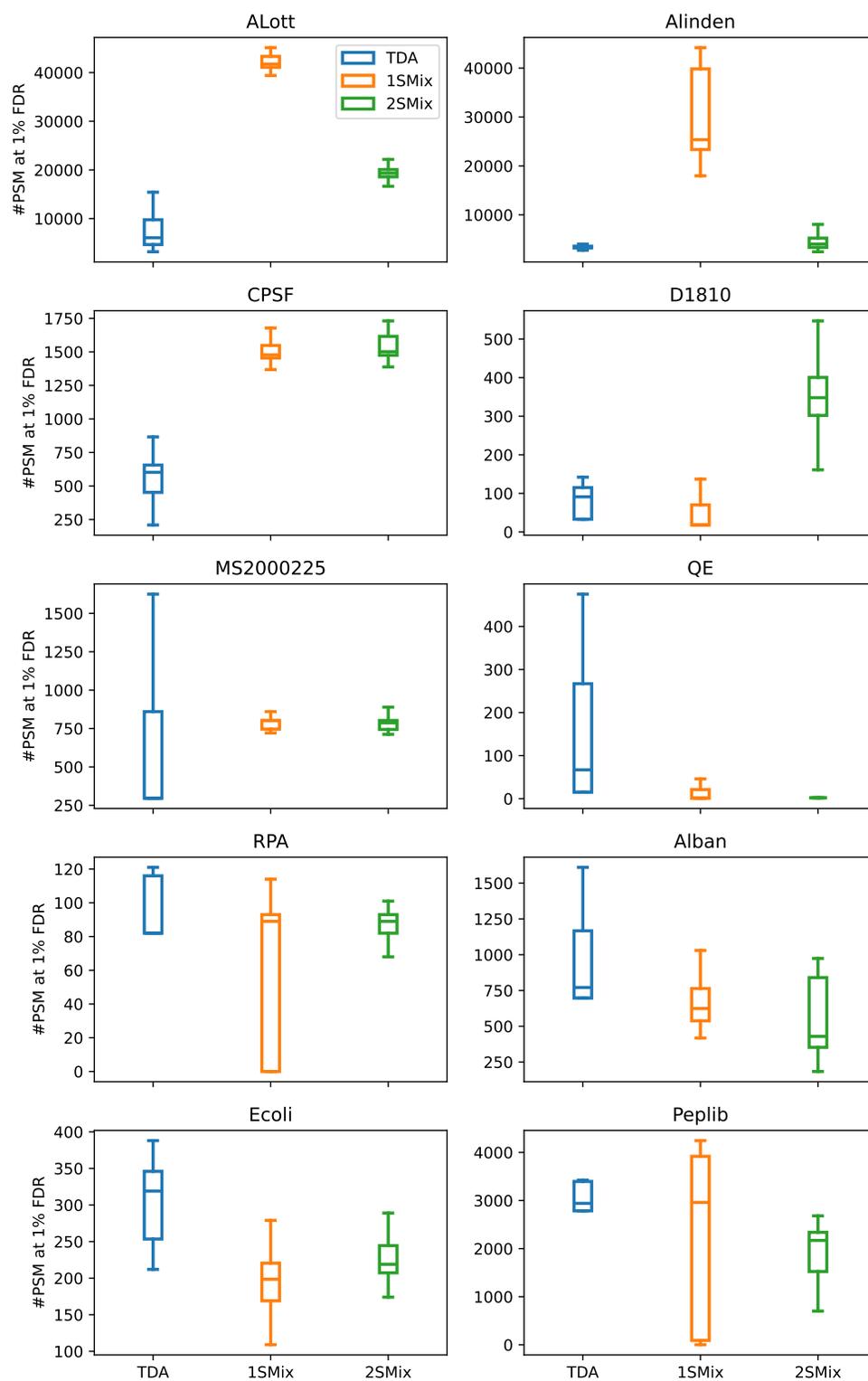


Figure S3: Stability of PSM identification. The stability of the methods was compared using 50 bootstrap samples, on which the number of PSMs scoring above 1% FDR threshold were counted. The larger spread of a boxplot indicates lower stability. Bootstrapping was performed on the set of original spectra.

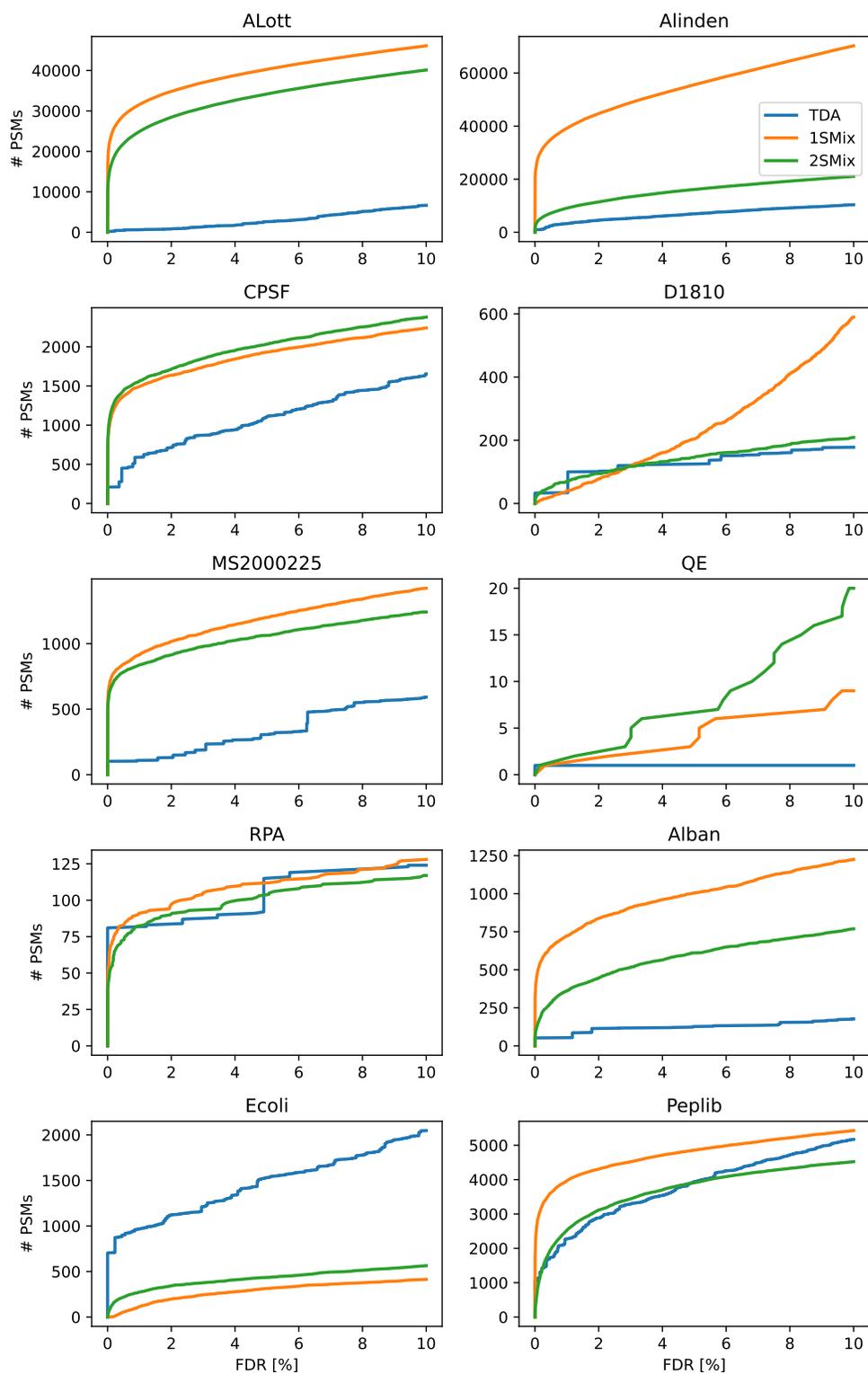


Figure S4: The number of identified PSMs as a function of FDR by 2SMix, 1SMix and TDA.



Figure S5: Comparison of PSMs identified by 2SMix and TDA at 1% FDR. PSMs identified by both 2SMix and TDA are in brown. Those identified by 2SMix, but not by TDA, are in green. Those identified by TDA, but not by 2SMix, are in pink.

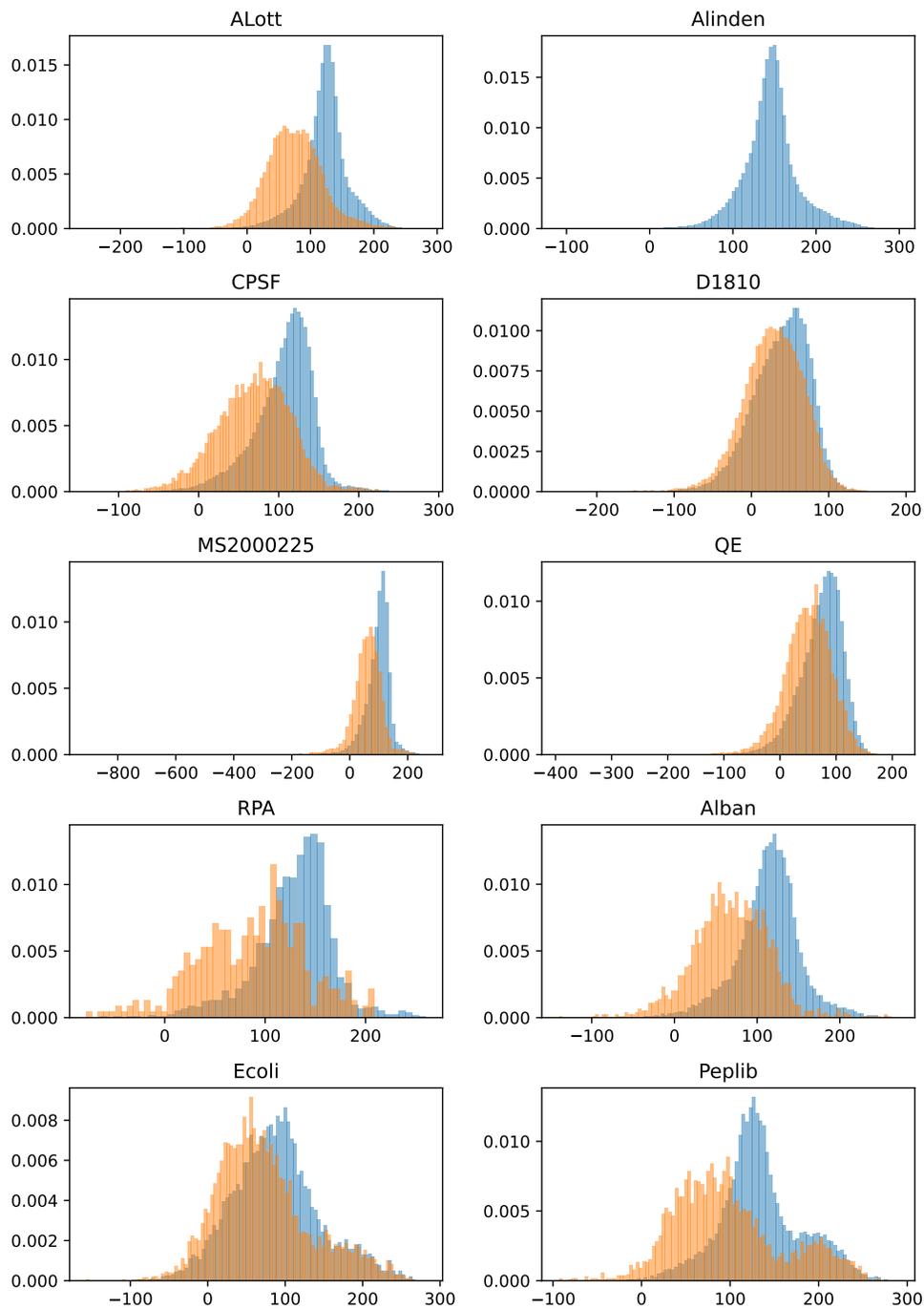


Figure S6: Second scores are not missing at random. The blue histogram is for the top scores from all spectra and the orange histogram is for the top scores whose spectra do not have a second hit. The second scores are not missing at random, since the top scores have a larger mass to the left.

References

- [1] Mirjam Dür. Duality in global optimization: optimality conditions and algorithmical aspects. 1999.