

Mass spectrometrists should search for all peptides, but assess only the ones they care about

Adriaan Sticker^{1, 2, 3, 4}, Lennart Martens^{2, 3, 4,*}, and Lieven Clement^{1, 4,*}

¹Department of Applied Mathematics, Computer Science & Statistics, Ghent
University, Belgium

²Medical Biotechnology Center, VIB, Ghent, Belgium

³Department of Biochemistry, Ghent University, Ghent, Belgium

⁴Bioinformatics Institute Ghent, Ghent University, Ghent, Belgium

*These authors contributed equally to this work and are joint corresponding
author, email: lennart.martens@vib-ugent.be & lieven.clement@ugent.be

August 12, 2016

Abstract

In shotgun proteomics identified mass spectra that are deemed irrelevant to the scientific hypothesis are often discarded. Noble (2015)¹ therefore urged researchers to remove irrelevant peptides from the database prior to searching to improve statistical power. We here however, argue that both the classical as well as Noble's revised method produce suboptimal peptide identifications and have problems in controlling the false discovery rate (FDR). Instead, we show that searching for all expected peptides, and removing irrelevant peptides prior to FDR calculation results in more reliable identifications at controlled FDR level than the classical strategy that discards irrelevant peptides post FDR calculation, or than Noble's strategy that discards irrelevant peptides prior to searching.

1 Introduction

Reliable peptide identification is key to every mass spectrometry-based shotgun proteomics workflow. The growing concern on reproducibility triggered leading journals to require that all peptide-to-spectrum matches (PSMs) are reported along with an estimate of their statistical confidence. The false discovery rate (FDR), i.e. the expected fraction of incorrect identifications, is a very popular statistic for this purpose. In many experiments, however, researchers want to focus on proteins of particular pathways, or few organisms in a metaproteomics sample. Hence, a large fraction of identified peptides are deemed irrelevant for their scientific hypothesis. Considering all PSMs induces an overwhelming multiple testing problem, which leads to few identifications of relevant peptides and thus to underpowered studies.² Currently, there is much debate on the optimal search strategy to boost the statistical power within this context (e.g. Noble, 2015¹ and <http://www.matrixscience.com/nl/201603/newsletter.html>).

The common approach, here referred to as the search-all-assess-all (all-all) strategy, involves (1) searching against all expected peptides in a sample, (2) calculating an FDR for each PSM and, (3) filtering irrelevant peptides from the candidate list. Recently, Noble (2015)¹ pointed out that this strategy is suboptimal because many unnecessary hypotheses are evaluated. Therefore, he

38 proposed to remove irrelevant peptides from the database prior to searching. He argues that a
39 search-subset-assess-subset (sub-sub) strategy improves the statistical power in two ways: (1) each
40 individual spectrum is tested against less candidate sequences and, (2) some spectra that originally
41 matched to proteins that were not of interest will lack a match in the subset search, decreasing
42 the number of PSMs for which an FDR estimate has to be provided.

43 However, we show that both the all-all and sub-sub strategy are suboptimal and often lead to
44 poorly controlled FDR. From a statistical perspective the all-all strategy FDR is biased as the
45 fraction of incorrect PSMs can differ substantially between the complete set and the subset leading
46 to too conservative or too liberal PSM lists. We also show that the sub-sub strategy forces many
47 good spectra derived from peptides deemed irrelevant to instead match subset peptides. This
48 because the irrelevant peptides were removed from the database prior to searching. Noble correctly
49 pointed out that this issue will not lead to statistical problems as long as a correct FDR procedure
50 is adopted. However, we illustrate that the popular target-decoy FDR procedure³ cannot avoid
51 these statistical problems when the sub-sub search strategy is adopted on small to moderate sized
52 subsets. We also argue that the Noble approach still sacrifices statistical power by testing more
53 hypotheses than necessary, i.e. PSMs that would match well to irrelevant peptides in the complete
54 search could actually be discarded because it is highly unlikely that these are subset peptides.

55 We therefore propose a search-all-assess-subset (all-sub) strategy by (1) searching the mass spectra
56 against a database with all proteins that are expected in the sample, and (2) discarding PSMs
57 matching to irrelevant peptides in the complete search prior to (3) FDR calculation, which has
58 the promise to further boost the statistical power. The filtering strategy in step (2) is independent
59 from the subsequent data analysis steps and can reduce the multiple testing problem considerably
60 without compromising the FDR calculation.⁴

61 2 Case studies

62 We first evaluate all three strategies on a *Pyrococcus furiosus* dataset with 15,365 high-resolution
63 spectra⁵ (see Supplementary Text). The complete proteome contains 2,051 proteins. We assess
64 36 subsets ranging from 17 to 381 proteins based on their Gene Ontology (GO) annotations. The
65 spectra are searched with the MS-GF+ search engine⁶ and the FDR is calculated with the target-
66 decoy approach³ (TDA, see Supplementary Text). Below, we focus on 175 proteins belonging to
67 the GO term “cytoplasm”. Results for all remaining subsets can be found in supplementary.

68 In figure 1 we illustrate that the fraction of incorrect target PSMs (π_0) in the complete search
69 is substantially different from the one in the cytoplasm subset. Based on the TDA approach we
70 estimate that 13.9% of the target PSMs are incorrect hits when adopting the all-all search, while
71 the actual fraction in the subset is probably lower, i.e. $\pi_0 = 7.2\%$ as estimated with the all-sub
72 strategy. This is also reflected in the distributions of the all-all and the all-sub MS-GF+ scores
73 in figure 1, which are bimodal. The first mode, corresponding to incorrect PSMs, is much higher
74 for the all-all strategy than for the all-sub method. Hence, the FDR cutoff using the all-all search
75 strategy is probably too conservative for the cytoplasm example. This is also reflected by the
76 increased number of subset PSMs that are returned by the all-sub method (2,578 vs 2,553 PSMs).
77 However, the FDR of the all-all method can also be too liberal, i.e. when the fraction of incorrect
78 PSMs in the subset is higher than the one in the complete search (e.g. the ATPase activity subset
79 in supplementary Fig. 10). Hence, the FDR in the all-all strategy is often not representative for
80 that of the subset leading to suboptimal PSM lists, which can be either too long or too short
81 depending on the scenario.

82 Figure 2 illustrates that the Noble sub-sub method is also suboptimal. Noble argues that his
83 strategy leads to a decrease of the number of PSMs that has to be tested. We indeed observe
84 that many target PSMs found in the all-all strategy are no longer matched or matched against a

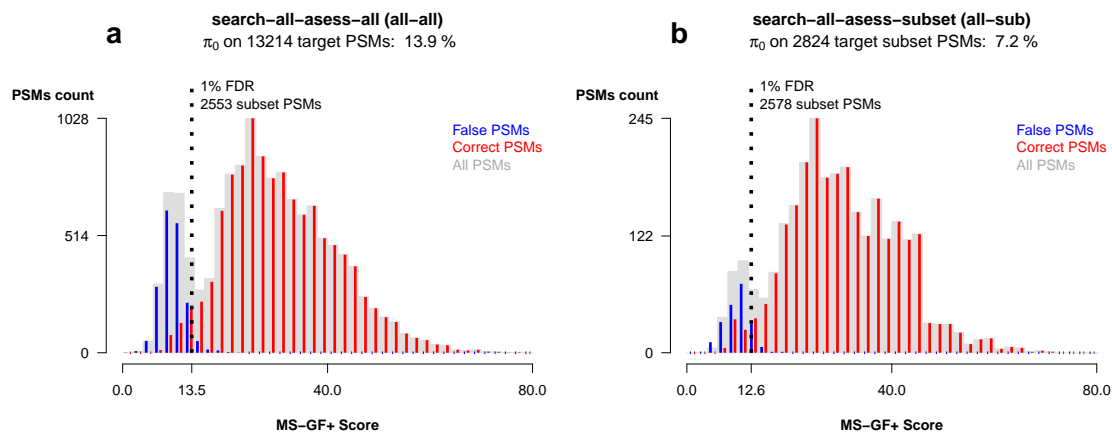


Figure 1: Histograms of MS-GF+ scores (grey) with estimated number of correct PSMs (red, $\#target - \#decoys$) and incorrect PSMs ($\#decoys$, blue), 1% FDR cutoff (dashed line). The left modes in the two distributions show that the fraction of null PSMs is higher in the all-all (a) than in the all-sub (b) search suggesting an incorrect FDR control, which is too conservative for the all-all strategy.

85 decoy sequence in the sub-sub strategy (left panel of Fig. 2), i.e. only 8,877 target PSMs remain.
 86 However, we show that the sub-sub method forces many spectra to match to incorrect sequences,
 87 which is increasing the fraction of incorrect PSMs to be tested. In the sub-sub strategy 5,831 out
 88 of 8,877 target PSMs ($\pi_0 = 65.7\%$, Fig 2. a) are expected to be incorrect compared to only 1,839
 89 out of 13,214 PSMs ($\pi_0 = 13.9\%$, Fig. 1 a) in the all-all strategy. Because of the higher fraction of
 90 incorrect PSMs in the sub-sub strategy, the score cutoff at 1% FDR increases from 13.5 to 15.0.
 91 A similar increase of the score cutoff at 1% FDR was also observed in all other GO subsets (see
 92 Supplementary Fig. 1-35). Due to this higher score cutoff, the sub-sub method will return a lower
 93 number of subset PSMs that match to the same peptide as in the all-all search strategy. Despite
 94 the more stringent cutoff, however, the sub-sub strategy returns more subset PSMs in 13 out of
 95 the 36 GO subsets. This happens because several spectra, derived from a peptide that was deemed
 96 irrelevant are now forced to match subset peptides. In the cytoplasm example (Fig. 2. b), the
 97 sub-sub PSM list contains many such PSMs that switched peptide, i.e. 53 on 2,559 PSMs at 1%
 98 FDR. This actually amounts to 2.2% of all returned PSMs, which strongly suggests an error rate
 99 of at least twice the adopted 1% FDR. Indeed, most of these switched PSMs are likely to be false
 100 positives: their MS-GF+ score is always lower for the sub-sub match than for the all-all match (see
 101 Fig. 2 b). The same behavior is also observed for many of the other subset searches, and this both
 102 at 1% and 5% FDR (Supplementary Fig. 36). Hence, the increased number of PSMs returned by
 103 the sub-sub method are highly questionable at best. We hypothesize that this is partially induced
 104 by an erratic behavior of TDA when dealing with small search libraries. This erratic behavior is
 105 illustrated in figure 2 a) where an unexpected enrichment in the number of expected correct PSMs
 106 can be observed at low MS-GF+ scores. This suggests that low quality spectra in a restricted
 107 search space tend to match incorrectly to target sequences because there is insufficient sequence
 108 variation in the decoys. This in turn leads to a poorly controlled FDR. The latter concern was
 109 also mentioned by Noble¹ and considerably reduces the reliability of the sub-sub search method
 110 for small to moderate subsets. In contrast, our proposed all-sub strategy (Fig. 1. b) also reduces
 111 the multiple testing problem considerably, but does not force good spectra to incorrectly match
 112 with subset peptide sequences because we search against all possible sequences.

113 In a second example we evaluated the methods on a malaria parasite *Plasmodium Falciparum*
 114 dataset containing 55,036 spectra.⁷ The *Plasmodium Falciparum* proteins in the dataset were
 115 captured from a mixed culture with human red blood cells. Hence, the samples are likely to be
 116 contaminated with a considerable amount of human proteins. We assessed two subsets of interest:

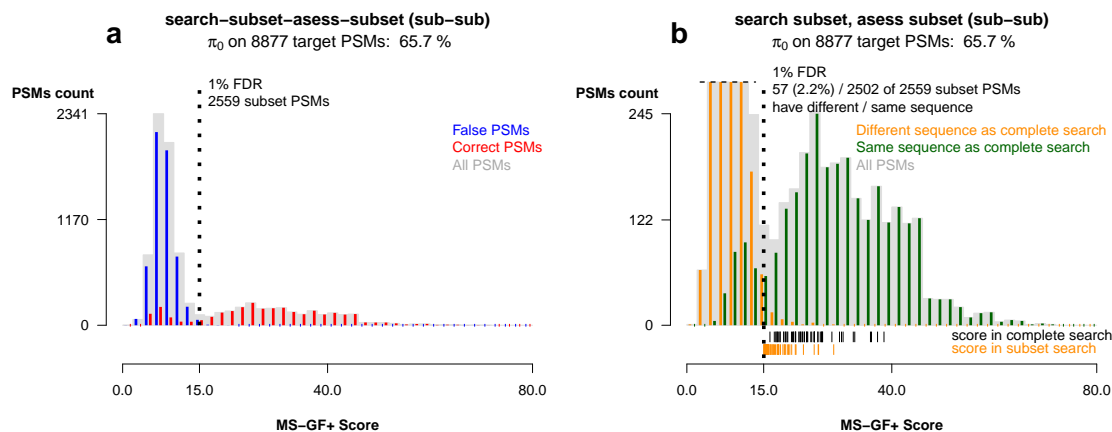


Figure 2: Histograms of MS-GF+ scores (grey) for the Noble search-subset-assess-subset strategy. a) estimated number of correct PSMs (red, $\#target - \#decoys$) and incorrect PSMs ($\#decoys$, blue). Compared to all-all strategy in figure 1 a) the score cutoff at 1% FDR (vertical dashed line) shifts to higher values, but still more subset PSMs are found. It also shows that many PSMs are forced on incorrect subset PSMs (huge first mode of the distribution). b) common PSMs (green) and PSMs that switched peptide sequences (orange) in the sub-sub strategy as compared to the all-all strategy in figure 1 a). These switched PSMs all have lower scores than in the all-all search and are therefore questionable at best (black and orange rug plot below histogram).

117 the plasmodium subset and the human subset consisting of 5,136 and 20,201 proteins, respectively.
 118 We do not expect problems with the TDA-FDR estimation in the sub-sub strategy because each
 119 subset contains more than 1,000 proteins.¹ This is confirmed in Supplementary Figures 37 and
 120 38, which do not show bimodal patterns in the distribution of the expected correct PSMs. Similar
 121 to the *Plasmodium* example in the paper of Noble,¹ the sub-sub approach retrieves more PSMs
 122 than the all-all strategy at 1% FDR (14,277 vs 1,3869 PSMs, respectively). However, the opposite
 123 is observed for the human subsample (6,108 vs 6,474 PSMs).

124 Again, the all-all FDR-control seems inadequate. π_0 in the complete set and the subset seems
 125 to differ considerably, i.e. the all-all strategy seems too conservative in the plasmodium subset
 126 ($\pi_0^{\text{all-all}} = 38.1\%$ and $\pi_0^{\text{all-sub}} = 24.5\%$, Supplementary Fig. 38), and too liberal in the human
 127 subset ($\pi_0^{\text{all-all}} = 38.1\%$ and $\pi_0^{\text{all-sub}} = 54.6\%$, Supplementary Fig. 39).

128 The Noble sub-sub strategy, on the other hand, is more conservative than our all-sub strategy at
 129 1% FDR and returned less subset PSMs for both subsets (Supplementary Fig. 38-39). Moreover, a
 130 considerable fraction of the returned PSMs still switched peptides between the all-all and sub-sub
 131 searches, i.e. 0.6% in the *Plasmodium* and 1.6% in the human subset at 1% FDR. As before, these
 132 switched PSMs are likely to contain many induced false positives. In the *Pyrococcus* example, we
 133 observed that the higher π_0 of the sub-sub strategy increased the score cutoffs at 1% FDR. Here,
 134 we observe a trade-off between the change in π_0 and in the decoy distribution (Supplementary Fig.
 135 38, panel a vs b, blue bars). The latter might be due to differences in overall protein composition
 136 in *Pyrococcus* and human subsets. In the *Plasmodium* subset the FDR cutoff still decreases from
 137 20.5 to 20.2 despite the increase in π_0 (from $\pi_0^{\text{all-all}} = 38.1\%$ to $\pi_0^{\text{sub-sub}} = 51.1\%$). This explains
 138 the increase in the number of returned subset PSMs as compared to the all-all method, i.e. an
 139 additional 324 subset PSMs with the same label are returned but they come at the expense of 84
 140 switched PSMs. Supplementary Figure 39 (panel a vs b) shows a much larger π_0 increase for the
 141 human subset ($\pi_0^{\text{sub-sub}} = 72.0\%$). This cannot be compensated by distributional changes of the
 142 null component, resulting in a higher threshold than with the strategies based on the complete
 143 search. The higher threshold leads to a lower number of returned human PSMs in common with
 144 the all-all search, and a high fraction of switched PSMs (1.6% at the 1% FDR level). Hence, even

145 for large subsets the sub-sub TDA FDR seems questionable.

146 3 Challenges and future directions

147 Both the *Pyrococcus* and *Plasmodium* example clearly indicate that mass spectrometrists will
148 benefit from searching for all peptides, but by only assessing the ones they care about. Our new
149 strategy returns more PSMs than the Noble approach and can be expected to provide a better
150 FDR control within peptide subsets than the two leading strategies, which discard peptides prior
151 to searching or post FDR calculation. Adopting the target decoy approach in our all-sub strategy
152 involving small subsets, however, often provides unstable FDR estimates due to a considerable
153 sample to sample variability in the number of subset decoys above a particular score cutoff. For
154 instance, the transmembrane transport GO subset has few PSMs and only 7 decoy hits. The
155 subset π_0 is higher than in the complete search and a more stringent MS-GF+ score cutoff seems
156 to be required. But due to the specific empirical distribution of the decoy scores the subset TDA
157 still results in a lower cutoff than in the all-all approach (Supplementary Fig. 27).

158 We observe that the location and shape of (1) the decoy distribution and (2) the estimated correct
159 distributional component of the target distribution in the all-all and all-sub strategy are very
160 similar (e.g. Supplementary Fig. 1-35, 38-39), indicating that the overall properties of the PSMs
161 in subsets and the complete set remain alike. We therefore propose to exploit the full information
162 available in the complete search for improving the estimation of the all-sub FDR.

163 Hence, we can estimate the distributional components using the complete search and only have to
164 rely on the subset decoy and target PSMs for the calculation of $\pi_0^{\text{sub-sub}}$. The latter FDR estimate
165 is more stable as it involves all subset decoys in contrast to the all-sub TDA, which only consists
166 of the relatively few subset decoys above a score cutoff. The reweighted FDR indeed results in a
167 more stringent MS-GF+ score cutoff for the transmembrane transport GO subset (Supplementary
168 Table 1).

169 We also developed a user-friendly web-based tool in R⁸ that provides (1) the all-sub FDR, (2)
170 the rescaled all-all FDR and (3) diagnostic plots for assessing the location-scale assumption.
171 (<http://iomics.ugent.be/saas/> and Supplementary Code) In our application, π_0 is estimated based
172 on the ratio of the number of subset decoys and the number of subset targets in a concatenated
173 target-decoy search. We feel that our approach can be further optimized for small subsets by using
174 the location and shape assumption explicitly when estimating π_0 .

175 References

176 ¹ William Stafford Noble. Mass spectrometrists should search only for peptides they care about.
177 *Nature Methods*, 12(7):605–608, 2015.

178 ² Alexey I Nesvizhskii. Proteogenomics: concepts, applications and computational strategies.
179 *Nature Methods*, 11(11):1114–1125, 2014.

180 ³ Joshua E Elias and Steven P Gygi. Target-decoy search strategy for increased confidence in
181 large-scale protein identifications by mass spectrometry. *Nature Methods*, 4(3):207–214, 2007.

182 ⁴ Richard Bourgon, Robert Gentleman, and Wolfgang Huber. Independent filtering increases de-
183 tection power for high-throughput experiments. *Proceedings of the National Academy of Sciences*
184 *of the United States of America*, 107(21):9546–51, may 2010.

185 ⁵ Marc Vaudel, Julia M. Burkhart, Daniela Breiter, René P. Zahedi, Albert Sickmann, and Lennart
186 Martens. A complex standard for protein identification, designed by evolution. *Journal of*
187 *Proteome Research*, 11(10):5065–5071, 2012.

188 ⁶ Sangtae Kim and Pavel a Pevzner. MS-GF+ makes progress towards a universal database search
189 tool for proteomics. *Nature communications*, 5:5277, 2014.

190 ⁷ Dingyin Tao, Ceereena Ubaida-Mohien, Derrick K Mathias, Jonas G King, Rebecca Pastrana-
191 Mena, Abhai Tripathi, Ilana Goldowitz, David R Graham, Eli Moss, Matthias Marti, and Rhoel R
192 Dinglasan. Sex-partitioning of the Plasmodium falciparum stage V gametocyte proteome provides
193 insight into falciparum-specific cell biology. *Molecular & cellular proteomics : MCP*, 13(10):2705–
194 24, 2014.

195 ⁸ R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for
196 Statistical Computing, Vienna, Austria, 2016.