# Scoring & Statistics

This is the Mascot result report for a MS/MS ion search. A list of proteins is reported with a score calculated from individual peptide matches. Proteins with lots of strong peptide matches come at the top of the report. However, it is very important to understand that the protein score in an MS/MS search is not statistically rigorous. It is just a way of ranking the protein hits.

If we expand one of the protein matches, we can see a list of peptides assigned to the protein. Each peptide has a score and expect value, these are probabilistic scores. If you mouse over the rank column you can see the score cutoffs for each query. The identity and homology threshold will be described in more detail in this talk (slides 17 and 18).

This is the Mascot result report for a peptide mass fingerprint search. There is also a list of proteins each of which matches some of the experimental peptide masses. In peptide fingerprint, we use probability-based scoring, and the report tells us that these matches are not statistically significant. The score threshold for this search is 70, and the top scoring match is 53. The graph is a histogram of the scores of the top ten matches and, as you see, all of them are in the area shaded green to indicate random, meaningless matches.

## What is probability based scoring?

We compute the probability that the observed match between the experimental data and mass values calculated from a candidate protein or peptide sequence is a random event.

The 'correct' match, which is not a random event, has a very low probability.

Reject anything with a probability greater than a chosen threshold, e.g. 0.05 or 0.01

What exactly do I mean by probability based scoring?

We calculate, as accurately as possible, the probability that the observed match between the experimental data, and mass values calculated from a candidate peptide or protein sequence, is a random event.

The real match, which is not a random event, then has a very low probability.

We can then reject anything with a probability greater than a chosen threshold, e.g. 1%

## Why is probability based scoring important?

- How else would you judge whether a PMF result was meaningful?
- For MS/MS, human judgment is subjective and can be unreliable

Why is probability based scoring important?

Well, how else would you judge whether a protein hit in a peptide mass fingerprint search was meaningful?

In the case of MS/MS data, it is very difficult to judge whether a match is significant or not by looking at the spectrum. Let me illustrate this with an example.

This match has a good number of matches to y and b ions, highlighted in red. Almost all the major peaks above 200 Da seem to be labelled. Could such a good match have occurred by chance?

You cannot tell, because you can match anything to anything if you try hard enough.

If this sounds strange, here's a simple analogy. If I say that I was tossing a coin and got ten heads in a row, does that mean there was something strange about the coin, like it had two heads? You cannot tell, because you need to know how many times I tossed the coin in total. If I picked it up off the table, tossed it ten times, then put it down, yes, that would suggest this was not a fair coin. However, if I tossed it ten thousand times, I would expect to get ten heads in a row more than once.

So, it isn't just a matter of how good the match is, i.e. how many y or b ions you found, it's a case of how hard you tried to find the match. In the case of a database search, this means how large is the database, what is the mass tolerance, how many variable modifications, etc., etc. These are very difficult calculations to do in your head, but they are easy calculations for the search engine.

If we look at the expectation value for this match, it is 0.66. That is, we could expect to get this match purely by chance. It looks good, but it's a random match.

If I show you a better match, then it is easy to dismiss the previous one as inferior. We can all make that judgement very easily. This match has an expectation value of less than 1 in 5,000. It is definitely not random.

The challenge is, what if you don't have the better match to compare against? Maybe this sequence wasn't in the database. If you only had the inferior match, how would you decide by looking at it whether it was significant or not?

The other interesting question is whether this is the "correct" match. Who can say that a better match isn't possible, where we get the extra y ion or the first and last b ions?

## Why is probability based scoring important?

- How else would you judge whether a PMF result was meaningful?
- For MS/MS, human judgment is subjective and can be unreliable
- Standard, statistical tests of significance can be applied to the results.

If we use probability based scoring, we can apply standard, statistical tests of significance to the results.

If we don't do this, then the only way to know the level of false positives is a target decoy search, and this isn't always possible, e.g. when searching a small number of spectra or doing a targeted search of only a few protein sequences

Probability based scoring calculates the probability that the match is random. This is, the probability that the match is meaningless. Many people ask whether we can report the probability that the match is correct. Is this possible?

It is certainly possible if you are analysing a known protein or standard mixture of proteins. If you know what the sequences are, or think you know, then the matches to the known sequences are defined to be correct and those to any other sequence are defined to be wrong. If the sample is an unknown, then it is difficult even to define what is meant by a correct match.

This is a typical MS/MS search result, where we see a series of high scoring homologous peptides. The sequences of the top three matches are very similar, and their expectation values vary from random through to very unlikely to be random. The best match has an expectation value of 1.2E-7. However, we cannot be sure that this is an identity match to the analyte peptide. It is simply the best match we could find in the database. There is always the possibility that a better match exists, that is not in the database, so to call it the correct match would be misleading.

The important thing is that we have a mechanism to discard matches that are nothing more than random matches.

It is a similar situation in Blast, except that you have the luxury of seeing when you have a perfect identity match. Here, the identity match has an expectation value of 2E-6, which reminds us that it would be a random match if the database was 2 million times larger. The match with one different residue is not worthless, it has an expectation value of 1E-5 and is a very good match. It just isn't as good a match as the one above.

Monoisotopic mass of neutral peptide Mr(calc): 1045.6739
Fixed modifications: iTRAQ4plex (K),iTRAQ4plex (N-term),Methylthio (C) (apply to specified residues or termini only)
Ions Score: 30  Expect: 2
Matches : 7/54 fragment ions using 8 most intense peaks   (help)

| # | b | b$^{++}$ | b* | b*$^{++}$ | b$^0$ | b$^{0++}$ | Seq. | y | y$^{++}$ | y* | y*$^{++}$ | y$^0$ | y$^{0++}$ | # |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 258.1934 | 129.6003 | | | | | I | | | | | | | 7 |
| 2 | 371.2775 | 186.1424 | | | | | I | 789.4951 | 395.2512 | 772.4685 | 386.7379 | 771.4845 | 386.2459 | 6 |
| 3 | 458.3095 | 229.6584 | | | 440.2989 | 220.6531 | S | 676.4110 | 338.7091 | 659.3845 | 330.1959 | 658.4004 | 329.7039 | 5 |
| 4 | 572.3524 | 286.6798 | 555.3259 | 278.1666 | 554.3419 | 277.6746 | N | 589.3790 | 295.1931 | 572.3524 | 286.6798 | | | 4 |
| 5 | 643.3895 | 322.1984 | 626.3630 | 313.6851 | 625.3790 | 313.1931 | A | 475.3360 | 238.1717 | 458.3095 | 229.6584 | | | 3 |
| 6 | 756.4736 | 378.7404 | 739.4471 | 370.2272 | 738.4630 | 369.7352 | L | 404.2989 | 202.6531 | 387.2724 | 194.1398 | | | 2 |
| 7 | | | | | | | K | 291.2149 | 146.1111 | 274.1883 | 137.5978 | | | 1 |

**MASCOT** : *Scoring & Statistics*       © 2007-2023 Matrix Science       MATRIX SCIENCE

If we are doing probability based matching, we are not scoring the quality of the spectrum, we are scoring whether the match is random or not.

Even when the mass spectrum is of very high quality, if the peptide is so short that it could occur in the database by chance, then you will not get a very good score.

The situation in a Blast search is identical. Even though this is a perfect identity match, the expectation value is 3745. This is just a random match. Hence, the earlier tip to discard spectra from low mass precursors.

## The Mascot Score

The Mascot score is -10Log10(P), where P is the absolute probability that observed match is random event

- For a PMF, P is the probability that the set of experimental peptide molecular masses came from the enzyme digest of the protein sequence
  - assuming the sequence is unrelated to the spectrum.
- For an MS/MS search, P is the probability that the masses in the MS/MS spectrum came from the gas phase fragmentation of the peptide sequence
  - assuming the sequence is unrelated to the spectrum.

For an MS/MS search, the protein score is *not* statistically rigorous. It is just a way of ranking the protein hits

MASCOT    : *Scoring & Statistics*          © 2007-2023 *Matrix Science*          MATRIX SCIENCE

For a peptide mass fingerprint, there is just one score that matters: the protein score. This tells us whether the match is significant or not, and is determined by calculating the probability of getting the observed number of peptide mass matches if the protein sequence was random.

For an MS/MS search, we have two scores. The important one is the peptide match score or ions score. This is the probability of getting the observed number of fragment ion mass matches if the peptide sequence was an unrelated sequence.

However, most people are interested in which proteins are present, rather than which peptides have been found. So, we assign peptide matches to protein hits and provide protein scores for MS/MS searches, so that the proteins with lots of strong peptide matches come at the top of the report.

However, it is very important to understand that the protein score in an MS/MS search is not statistically rigorous. It is just a way of ranking the protein hits.

This is why there is no expect value for the protein score in an MS/MS search, and why there is a short explanation at the top of every report.

## Significance Thresholds

### The identity threshold is calculated from the number of trials

If there are 500,000 entries in the database, a 1 in a 20 chance of getting a false positive match for a peptide mass fingerprint is a probability of

$P = 1 / (20 \times 500{,}000)$

which is a score of

$S = -10\mathrm{Log}P = 70$

Because a Mascot score is a log probability, assigning a significance threshold is very simple. It is just a function of the number of trials - the number of times we test for a match. For a peptide mass fingerprint, this is the number of entries in the database. For an MS/MS search, it is the number of peptides in the database that fit to the precursor mass tolerance. For an enzyme like trypsin, and a reasonable mass tolerance, this number will be less than the number of entries in the database. For a no-enzyme search, the number of trials will often be more than the number of entries in the database.

So, for example, if we are comfortable with a 1 in a 20 chance of getting a false positive match, and we are doing a PMF search of a database that contains 500,000 entries, we are looking for a probability of less than $1 / (20 \times 500{,}000)$ which is a Mascot score of 70

If we could only tolerate a false positive rate of 1 in 200 then the threshold would be 80, 1 in 2000 90, etc.

For MS/MS searches with trypsin, and a reasonable mass tolerance, the numbers tend to be lower. The default identity threshold is typically a score of around 40.

**Significance Thresholds**

The homology threshold is an empirical measure of whether the match is an outlier

MASCOT : *Scoring & Statistics*  © 2007-2023 Matrix Science  MATRIX SCIENCE

Unfortunately, MS/MS spectra are often far from ideal, with poor signal to noise or gaps in the fragmentation. In such cases, it may not be possible to reach the identity threshold score, even though the best match in the database is a clear outlier from the distribution of random scores. To assist in identifying these outliers, we also report a second, lower threshold for MS/MS searches; the 'homology' threshold. This simply says the match is an outlier.

In practice, from measuring the actual false positive rate by searching large data sets against reversed or randomised databases, we find that the identity threshold is usually conservative, and the homology threshold can provide a useful number of additional true positive matches without exceeding the specified false positive rate.

# Expectation values



We display an expect or expectation value in addition to the score.

## Expectation values

### The number of times you could expect to get this score or better by chance

$$E = P_{threshold} * (10^{**} ((S_{threshold} - score) / 10))$$

If $P_{threshold} = 0.05$ and $S_{threshold} = 50$
  score = 40 corresponds to E = 0.5
  score = 50 corresponds to E = 0.05
  score = 60 corresponds to E = 0.005

**MASCOT**   *: Scoring & Statistics*      © 2007-2023 Matrix Science      MATRIX SCIENCE

The expectation value does not contain new information. It can be derived directly from the score and the threshold. The advantage is that it tells you everything you need to know in a single number.

It is the number of times you could expect to get this score or better by chance.

A completely random match has an expectation value of 1 or more.

The better the match, the smaller the expectation value.

The most important attributes of a scoring scheme are sensitivity and specificity. That is, you want as many correct matches as possible, and as few incorrect matches as possible.

This is often illustrated in the form of a Receiver Operating Characteristic or ROC plot. This plots the relationship between the true positive and false positive rates as the threshold is varied. The origin is a very high threshold, which lets nothing through. At the top right, we have a very low threshold, that allows everything through. Neither extreme is a useful place to be. The diagonal represents a useless scoring algorithm, that is equally likely to let through a false match as a true one. The red curve shows a useful scoring algorithm, and the more it pushes the curve up towards the top left corner, the better. Setting a threshold towards this top left corner gives a high ratio of correct matches to false matches.

A few years ago, there was a little too much focus on sensitivity and not enough consideration given to specificity, so that some of the published lists of proteins were not as accurate as the authors might have hoped.

# Validation



A growing awareness of this problem led to initiatives from various quarters. Most notably, the Editors of Molecular and Cellular Proteomics, who held a workshop in 2005 to define a set of guidelines, which have been revised multiple times, last time in 2017.

For large scale studies, there is a requirement to estimate your false discovery rate. One of the most reliable ways to do this is with a so-called decoy database.

## Validation

### Search a "decoy" database

- Decoy entries can be reversed or shuffled or randomised versions of target entries
- Decoy entries can be separate database or concatenated to target entries

### Gives a clear estimate of false discovery rate

- Elias, J. E. and Gygi, S. P., Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry, Nature Methods 4 207-214 (2007)

**MASCOT** : *Scoring & Statistics* © 2007-2023 Matrix Science MATRIX SCIENCE

---

This is very simple but very powerful. You repeat the search, using identical search parameters, against a database in which the sequences have been reversed or randomised. You do not expect to get any real matches from the decoy database. So, the number of matches that are found in the decoy database is an excellent estimate of the number of false positives in the results from the target database.

You'll read a lot of discussion in the literature about whether the decoy sequences should be reversed or randomised; whether to search a single database containing both target and decoy sequences or separate databases. I suggest the most important thing is to do a decoy search; any decoy search. What you need to know is whether your level of false positives is 1% or 10% or 100%. Its less of a concern whether its 1% or 1.1%.

Although this is an excellent validation method for large data sets. It isn't useful when you only have a small number of spectra, or a small database, because the numbers are too small to give an accurate estimate. Hence, this is not a substitute for a stable scoring scheme, but it is an excellent way of validating important results.

**Validation**

The option to run a decoy search is built-in to Mascot Server. If you choose the Decoy checkbox on the search form, then every time a protein or peptide sequence from the target database is tested, a reversed or randomised sequence of the same length is automatically generated and tested. The average amino acid composition of the random sequences is the same as the average composition of the target database. The matches and scores for the decoy sequences are recorded separately in the result file. The result is identical to searching a separate database rather than a concatenated database.

On our public web site there is a help page devoted to decoy database searches. It includes a download link to a utility program that allows you to create a randomised or reversed database in the case that you don't want to use the built-in feature. If you have an early version of Mascot, or if you want to verify the results from another search engine, you can use this utility to create a decoy database for searching.

When the search is complete, the statistics for matches to the decoy sequences are reported in the result header. If you change the significance threshold, the numbers are recalculated. There is a option to adjust the significance threshold so as to achieve a chosen FDR value. For example, if we choose 1% FDR using the homology threshold.

The significance threshold has been automatically adjusted. You can see the adjustment by expanding the Score distribution section report where it displays the change from 0.05 to 0.03661. In the report the expect values are recalculated and matches with an expect below 0.05 are significant.

Why do we get these false positives? Do they reflect some defect in the search engine? Let's have a closer look. If you click the decoy report link as shown at the bottom of the screen shot in the expanded Sensitivity and FDR section, then you will see the results from searching the randomised database.

The results from the matches to the randomised sequences are saved in a different section of the results file on the Mascot server. This means that we can view these results in exactly the same way as if we had performed a separate search against a randomised database that we had created manually. We can see matches here with scores of 37 and 51, with expect values well below 1%/0.01. If we click on the query number link to display the Peptide View of one of these matches …

This is what it looks like. A pretty decent match from a decoy sequence. Tryptic peptide, no variable modifications, good run of b and y ions, most of the larger peaks matched.

Asking whether it is correct or wrong becomes almost a philosophical question.

The fact is, when we search large numbers of spectra against large sequence databases, we can get such matches by chance. No amount of expert manual inspection will prevent this. Database matching is a statistical process and, for this search, the number and magnitude of the false positives are well within the predicted range, which is all we can ask for.

## Protein FDR

### Automatically calculated when decoy option selected
### Displayed in the decoy section:

▼*Sensitivity and FDR (reversed protein sequences)*

| | Target | Decoy | FDR |
|---|---|---|---|
| Protein family members | 3502 | 592 | 16.90% |
| PSMs ∨ above homology ∨ | 20109 | 659 | 3.28% |

Decoy results are available in ☑the decoy report.

**MASCOT** : *Scoring & Statistics*     © 2007-2023 Matrix Science     MATRIX SCIENCE

In Mascot Server 2.7, we introduced the Protein FDR value which is automatically calculated when the decoy search option is selected.

It is displayed above the peptide/PSM decoy results.

Protein FDR is based on the following assumptions and definitions:

By default, the protein family report only shows peptide sequence matches (PSMs) with significant scores and are used for the protein family assignment.

A protein family member may represent multiple same set proteins. Only members of all the protein families in the report, those that contain a unique peptide, are counted.

Protein count used for FDR is count of family members. That is, if the report contains 2 families, one with 4 members and the other with a single member, this counts as a total of 5 proteins. Same-set, sub-set and intersection proteins are not counted, because they have no independent peptide evidence.

While a protein ID is a false positive when all the PSM's are false positives. Just one true PMS would make the protein identification a true positive. This is very important.

**Based on the following assumptions and definitions:**

**Given the number of proteins and the numbers of true and false peptide sequences we use a hypergeometric model to estimate the number of proteins are truly false positive**

• Simplified approach to that used by MAYU, from the Aebersold group

Given the number of proteins and the numbers of true and false peptide sequences we use a hypergeometric model to estimate the number of proteins that are truly false positive.

The algorithm is a simplified approach to that used by MAYU, from the Aebersold group:

https://www.mcponline.org/content/8/11/2405https://www.mcponline.org/content/8/11/2405

The main differences are that we do not make a separate estimate of the FDR for one-hit wonders and we do not partition the database by protein size. We use a simpler estimate for the number of false proteins in the target database, based on the assumption that the number of decoy proteins never reaches a significant proportion of the database size.

## Example

**Target database has 1000 entries**

**Search results**

- 500 target proteins
- 10 decoy proteins
- FDR 10/500=2%?

**False PSMs distributed across 10 proteins**

- Some will also contain true PSMs
- Half proteins in target database contain true PSMs
- Estimate that only 5 target proteins contain nothing but false PSMs

**Protein FDR is 5/500 = 1%**

**MASCOT** : *Scoring & Statistics*     © 2007-2023 *Matrix Science*     MATRIX SCIENCE

Imagine the and the search results show 500 target proteins and 10 decoy proteins. Does this mean protein FDR is 10/500 = 2%? No, it does not. We can assume the false PSMs in the target are distributed across 10 proteins, but some of these will also contain true PSMs, so should not be counted as false. Since half the proteins in the target database contain true PSMs, a reasonable estimate would be that only 5 target proteins containing nothing but false PSMs, so that the protein FDR is 5/500 = 1%.

## Adjusting the Protein FDR

### Default protein FDR

| | | | | | |
|---|---|---|---|---|---|
| Format | Significance threshold p< | 0.05 | Max. number of families | AUTO | [help] |
| | Target FDR (overrides sig. threshold) | (not set) | FDR type | PSM | |
| | Display non-sig. matches | ☐ | Min. number of sig. unique sequences | 1 | |
| | Show Percolator scores | ☐ | Dendrograms cut at | 0 | |
| | Preferred taxonomy | All entries | | | |

▼Sensitivity and FDR (reversed protein sequences)

| | Target | Decoy | FDR |
|---|---|---|---|
| Protein family members | 3502 | 592 | 16.90% |
| PSMs above homology | 20109 | 659 | 3.28% |

Decoy results are available in the decoy report.

### Set Min. number sig. unique peptide sequences to 2

| | | | | | |
|---|---|---|---|---|---|
| Format | Significance threshold p< | 0.05 | Max. number of families | AUTO | [help] |
| | Target FDR (overrides sig. threshold) | (not set) | FDR type | PSM | |
| | Display non-sig. matches | ☐ | Min. number of sig. unique sequences | 2 | |
| | Show Percolator scores | ☐ | Dendrograms cut at | 0 | |
| | Preferred taxonomy | All entries | | | |

▼Sensitivity and FDR (reversed protein sequences)

| | Target | Decoy | FDR |
|---|---|---|---|
| Protein family members | 2119 | 19 | 0.90% |
| PSMs above homology | 20109 | 659 | 3.28% |

Decoy results are available in the decoy report.

**MASCOT** : *Scoring & Statistics*     © 2007-2023 Matrix Science     **MATRIX SCIENCE**

The default significance threshold for a Mascot search is usually 0.05 and this will often give a peptide FDR in the region of 5%. In this dataset the protein FDR is ~17%. If we want to lower the Protein FDR, we can try adjusting the peptide FDR to 1%. We can also try adjusting the report other ways.

We can adjust the Minimum number of significant unique sequences. This has quite a strong affect on the Protein FDR. If we change it to 2 and eliminate the "one hit wonders" the protein FDR drops to 0.9%. Note that "one hit wonders" are not necessarily false positives, just that they only have one significant peptide sequence match.

# Adjusting the Protein FDR

| Format | Significance threshold p< | (0.004) | Max. number of families | AUTO | [help] |
|--------|---------------------------|---------|-------------------------|------|--------|
| | Target FDR (overrides sig. threshold) | (not set) ∨ | FDR type | PSM ∨ | |
| | Display non-sig. matches | ☐ | Min. number of sig. unique sequences | 1 ∨ | |
| | Show Percolator scores | ☐ | Dendrograms cut at | 0 | |
| | Preferred taxonomy | All entries | | ∨ | |

**▼Sensitivity and FDR (reversed protein sequences)**

| | Target | Decoy | FDR |
|--|--------|-------|-----|
| Protein family members | 2674 | 27 | 1.01% |
| PSMs ∨ above homology ∨ | 13191 | 29 | 0.22% |

Decoy results are available in ☑the decoy report.

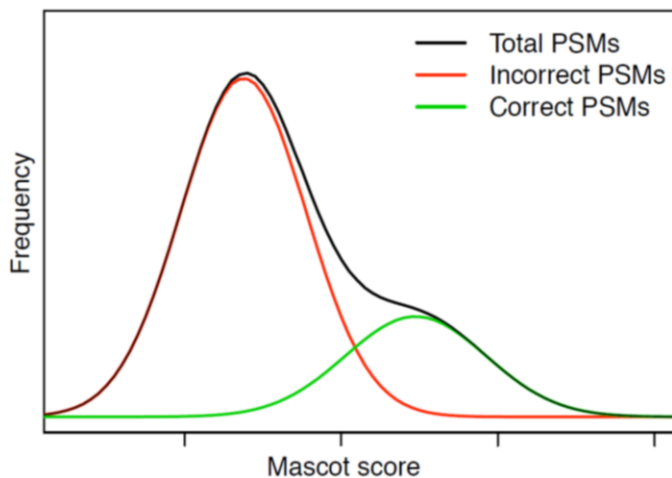MASCOT    : *Scoring & Statistics*    © 2007-2023 Matrix Science    MATRIX SCIENCE

Alternatively, we can adjust the Significance threshold for the results. I took a guess and reduced it by a factor of 10 from the default values of 0.05 to 0.004 and clicked the format button. The resulting Protein FDR is approximately 1%.

The current HUPO guidelines Interpretation Guidelines for large-scale results recommend adjusting the settings to lower than 1% protein-level global FDR so after formatting this search result would meet those guidelines.

When interpreting the results, a protein FDR of 1% only tells us that 1% of the proteins listed are wholly false. This doesn't mean the other 99% are "correct". In particular, where there are same-set proteins, we cannot say which one is "correct".

This is because database redundancy causes protein inference ambiguity and we can account for the PSM evidence using several sets of proteins. It is important to remember that a protein accession number in the summary report does not mean "this is the correct protein", it means "the correct protein is likely to be very similar to one of the set of proteins represented by this family member".

Sensitivity improvement is always a hot topic. A limitation of database matching is that even the best scoring scheme cannot fully separate the correct and incorrect matches, as shown here in a schematic way. The score distribution for the correct matches overlaps that of the incorrect matches. When we use a decoy search we are deciding where to place a threshold of some sort.

But, what if we could find ways to pull these two distributions further apart? In other words, improve the specificity of the scoring.
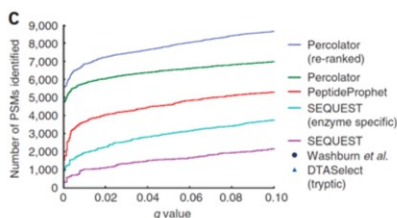
One of the first attempts to do this was Peptide Prophet from the ISB. This was and is popular for transforming Sequest scores into probabilities.

It takes information about the matches in addition to the score, and uses an algorithm called expectation maximization to learn what distinguishes correct from incorrect matches. Examples of additional information would be precursor mass error, number of missed cleavages, or the number of tryptic terminii.

We can use the matches from a decoy database as negative examples for the classifier. Another popular tool, Percolator, trains a machine learning algorithm called a support vector machine to discriminate between a sub-set of the high-scoring matches from the target database, assumed correct, and the matches from the decoy database, assumed incorrect.

Sensitivity optimisation

M. Brosch, L. Yu, T. Hubbard, J. Choudhary, *J Proteome Res* (2009).

Percolator can give very substantial improvements in sensitivity. The original Percolator was implemented mainly with Sequest in mind, but Markus Brosch at the Sanger Centre wrote a wrapper that allowed it to be used with Mascot results and published results such as this. The black trace is the sensitivity using the Mascot homology threshold and the red trace is the sensitivity after processing through Percolator. It doesn't work for every single data set. But, when it does work, the improvements can be most impressive.

# Sensitivity optimisation

The developers of Percolator have kindly agreed to allow us to distribute and install Percolator as part of Mascot Server. This option will be available for any search that has at least 100 MS/MS spectra and auto-decoy results, but it works best if there are several thousand spectra. To switch to Percolator scores, just check the box and then choose Filter. In this example we take a medium sized search result.

# Sensitivity optimisation

**Sensitivity and FDR (reversed protein sequences)**

|  | Target | Decoy | FDR |
|---|---|---|---|
| Protein family members | 3674 | 143 | 3.89% |
| PSMs ⌄ above homology ⌄ | 15520 | 155 | 1.00% |

Decoy results are available in ☑the decoy report.

**Sensitivity and FDR (reversed protein sequences)**

|  | Target | Decoy | FDR |
|---|---|---|---|
| Protein family members | 4145 | 140 | 3.38% |
| PSMs ⌄ above homology ⌄ | 20079 | 154 | 0.77% |

Decoy results are available in ☑the decoy report.

| Delta M | Score | Expect | Rank | U | 1 | 2 | 3 | Peptide |
|---|---|---|---|---|---|---|---|---|
| 0.1363 0 | 33 | 0.00049 | ▶1 | U | ■ | ■ | ■ | R.LIGDAAK.N |
| 0.3020 0 | 14 | 0.039 | ▶1 | U | ■ | | | K.VQVEYK.G |
| 0.2841 0 | 17 | 0.018 | | | ■ | | | |
| 0.4581 0 | 18 | 0.015 | ▶1 | | *Score > 13 indicates **identity*** | | | |
| 0.0517 0 | 21 | 0.0087 | ▶1 | U | ■ | | | K.VLEDSDLK.K |
| 0.2227 0 | 25 | 0.0031 | ▶1 | U | ■ | | | K.ITITHDQNR.L |

**MASCOT** : *Scoring & Statistics*     © 2007-2023 Matrix Science     ▲ MATRIX SCIENCE

Using the Mascot homology threshold for a 1% false discovery rate, there are 15520 peptide matches. Re-scoring with Percolator gives a useful increase to 20079 matches.

Note that, in general, the scores are lower after switching to Percolator. The Posterior error probability is tabulated in the expect column. A Mascot score is calculated from the expect value and the single score threshold, which we describe as the identity threshold, has a fixed value of 13 (-10 log 0.05). By keeping the score, threshold, and expect value consistent, we hope to avoid breaking any third party software that expects to find these values.