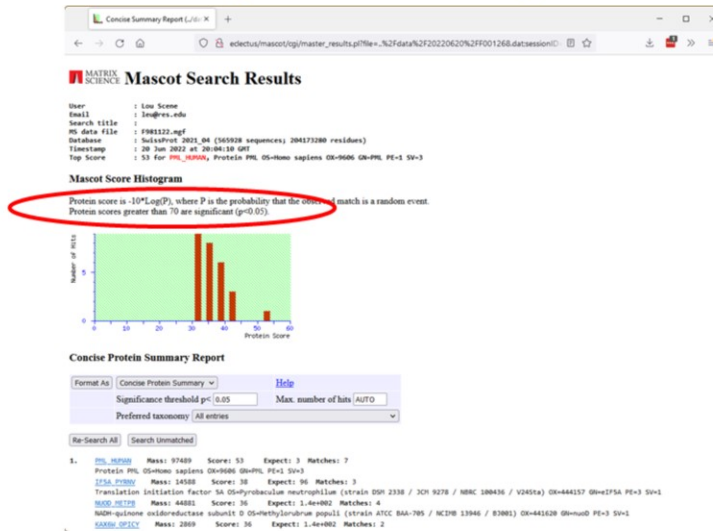


# Scoring & Statistics

**MASCOT**

 **MATRIX  
SCIENCE**

# Probability based scoring



**MASCOT**

: Scoring & Statistics

© 2007-2022 Matrix Science

**MATRIX  
SCIENCE**

This is the Mascot result report for a peptide mass fingerprint search. There is a list of proteins, each of which matches some of the experimental peptide masses, but 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%

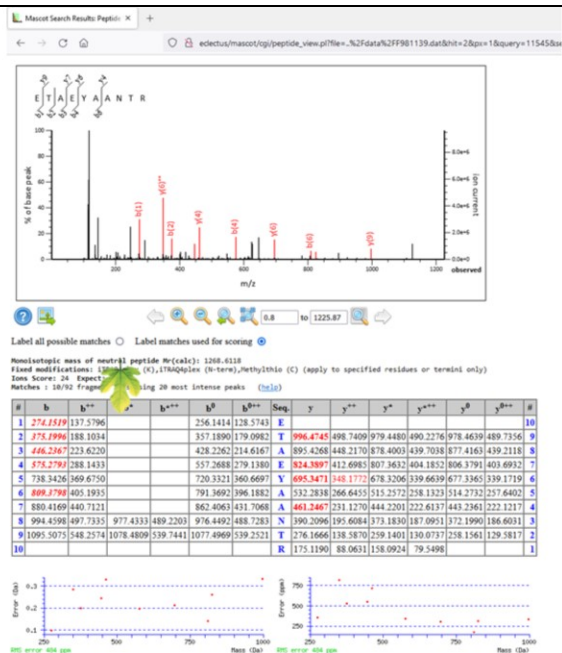
## Why is probability based scoring important?

- 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



**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



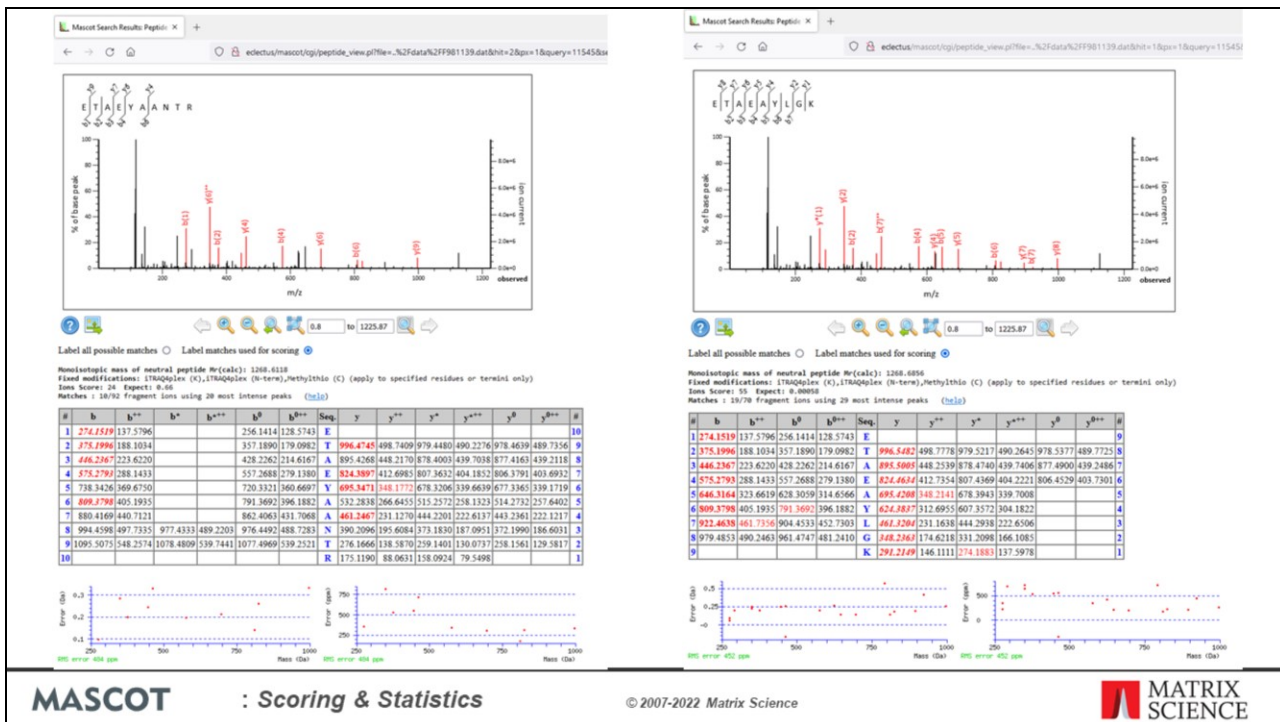
This match has a good number of matches to y and b ions, highlighted in red. Almost all the major peaks above 200 Da seem to be labelled. Could such a good match have occurred by chance?

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

## Can we calculate a probability that a match is correct?

**Yes, if it is a test sample and you know what the answer should be**

- Matches to the expected protein sequences are defined to be correct
- Matches to other sequences are defined to be wrong

**If the sample is an unknown, then you have to define “correct” very carefully**

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.



Query Dupes	Observed	Mr (expt)	Mr (calc)	ppm	M Score	Expect	Rank	U	114/113	115/113	116/113	117/113	118/113	119/113	121/113	1 2	Peptide
#175711 ▶4	566.7453	2829.6901	2828.6987	-3.04	1	32	0.0011	▶1	---	---	---	---	---	---	---	■ ■	R.YIYGRPVQVATYR.F + 2 iTRAQplex (Y)
#175901 ▶1	570.5222	2847.5746	2847.5786	-1.41	1	38	0.00029	▶1	1.033	1.029	0.595	1.227	1.358	1.050	1.059	■ ■	R.YLQRTQKSTLPVETK.D
Locus:29.1.1.3449.3 File:"iTRAQ_Splex.wiff" Score > 33 indicates identity Score > 24 indicates homology																	
#176112	952.7923	2855.3551	2855.3590	-1.38	0	80	1.2e-007	▼1	1.100	1.051	1.185	0.978	1.288	1.062	1.253	■ ■	R.TLEIPGNSDFNMIPGDFNSYVR.V + Deamidated (NQ)
				-1.38	0	60	1.4e-005	2								■ ■	R.TLEIPGNSDFNMIPGDFNSYVR.V + Deamidated (NQ)
				-1.38	0	21	0.11	3								■ ■	R.TLEIPGNSDFNMIPGDFNSYVR.V + Deamidated (NQ)
				-41.9	1	6	3.2	4									TLQELVHKAASYMDR + Deamidated (NQ); Oxidation (M)
				-1.32	1	2	7.6	5									TLETTPANDRAEPNSQLDSTHSGR + Deamidated (NQ); 2 C
				-1.32	1	2	7.6	5									TLETTPANDRAEPNSQLDSTHSGR + Deamidated (NQ); 2 C
				-30.2	1	2	8	7									NPALGNQTVAGLFLANSSEALERAVR + Deamidated (NQ);
				-30.2	1	2	8	7									NPALGNQTVAGLFLANSSEALERAVR + Deamidated (NQ);
				-30.2	1	2	8	7									NPALGNQTVAGLFLANSSEALERAVR + Deamidated (NQ);
				-33.7	0	2	8	7	10								TEIISERNQIYEKDK + 2 Deamidated (NQ); iTRAQplex
#176114 ▶8	714.8489	2855.3665	2855.3590	2.62	0	75	1.4e-007	▶1	1.092	1.044	1.151	1.048	1.342	1.127	1.153	■ ■	R.TLEIPGNSDFNMIPGDFNSYVR.V + Deamidated (NQ)
#176150 ▶1	715.8371	2859.3193	2859.3161	1.12	0	68	1.2e-006	▶1	---	---	---	---	---	---	---	■ ■	R.AACAGLQNPVKGTYQKGVV + Deamidated (NQ); iTRAQplex
#176504 ▶1	580.5195	2897.5611	2897.5609	0.992	1	48	3.2e-005	▶1	---	---	---	---	---	---	---	■ ■	R.EGALHRSSELYTEANPLDHR.G + iTRAQplex (Y)
#176508 ▶2	483.9355	2897.5693	2897.5609	2.93	1	37	0.00038	▶1	---	---	---	---	---	---	---	■ ■	R.EGALHRSSELYTEANPLDHR.G + iTRAQplex (Y)
#177105 ▶16	587.9223	2934.5751	2934.5830	-2.68	0	75	9.5e-008	▶1	---	---	---	---	---	---	---	■ ■	R.AVGSQATFSYTYMILR.G + 2 iTRAQplex (Y)

MASCOT

: Scoring & Statistics

© 2007-2022 Matrix Science



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.

NCBI Blast unnamed protein p: x

https://blast.ncbi.nlm.nih.gov/Blast.cgi?aln=tblastn&db=nr

**annexin, partial [Klebsiella pneumoniae]**  
Sequence ID: [WP\\_143442485.1](#) Length: 122 Number of Matches: 1

Range 1: 64 to 79 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Identities	Positives	Gaps
52.8 bits(117)	2e-06	16/16(100%)	16/16(100%)	0/16(0%)

Query 1 SEDFGVNEDLADSDAR 16  
Sbjct 64 SEDFGVNEDLADSDAR 79

**annexin A1 [Otoleum garnettii]**  
Sequence ID: [XP\\_003783004.1](#) Length: 346 Number of Matches: 1

Range 1: 189 to 204 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Identities	Positives	Gaps
50.3 bits(111)	1e-05	15/16(94%)	15/16(93%)	0/16(0%)

Query 1 SEDFGVNEDLADSDAR 16  
Sbjct 189 SEDFGVNEDLADSDAR 204

**annexin A1 [Equus asinus]**  
Sequence ID: [XP\\_014704224.2](#) Length: 346 Number of Matches: 1

Range 1: 189 to 204 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Identities	Positives	Gaps
50.3 bits(111)	1e-05	15/16(94%)	16/16(100%)	0/16(0%)

Query 1 SEDFGVNEDLADSDAR 16  
Sbjct 189 SEDFGVNEDLADSDAR 204

**MASCOT : Scoring & Statistics** © 2007-2022 Matrix Science **MATRIX SCIENCE**

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.



**MASCOT**

: Scoring & Statistics

© 2007-2022 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 screenshot displays a BLAST search interface with three results. Each result shows a query sequence (IISNALK) and a subject sequence (IISNALK) with a length of 7. The expectation value for all matches is 3745, and the identity is 100% (7/7). The first two results are from *Orbicella faveolata* (LOC110051226 isoform X1 and X2), and the third is from *Polynucleobacter paneuropaeus* (MBT8554431.1). The interface includes navigation links like 'Download', 'GenPept', 'Graphics', 'Next', 'Previous', and 'Descriptions'.

**Match 1:** uncharacterized protein LOC110051226 isoform X1 [*Orbicella faveolata*]  
 Sequence ID: [XP\\_020612903.1](#) Length: 4875 Number of Matches: 1  
 Range 1: 1327 to 1333  
 Score: 24.4 bits(50) Expect: 3745 Identities: 7/7(100%) Positives: 7/7(100%) Gaps: 0/7(0%)  
 Query 1 IISNALK 7  
 Sbjct 1327 IISNALK 1333

**Match 2:** uncharacterized protein LOC110051226 isoform X2 [*Orbicella faveolata*]  
 Sequence ID: [XP\\_020612904.1](#) Length: 4874 Number of Matches: 1  
 Range 1: 1326 to 1332  
 Score: 24.4 bits(50) Expect: 3745 Identities: 7/7(100%) Positives: 7/7(100%) Gaps: 0/7(0%)  
 Query 1 IISNALK 7  
 Sbjct 1326 IISNALK 1332

**Match 3:** filamentous hemagglutinin N-terminal domain-containing protein [Polynucleobacter paneuropaeus]  
 Sequence ID: [MBT8554431.1](#) Length: 3465 Number of Matches: 1  
 See 1 more title(s) See all Identical Proteins(IPG)  
 Range 1: 795 to 801  
 Score: 24.4 bits(50) Expect: 3747 Identities: 7/7(100%) Positives: 7/7(100%) Gaps: 0/7(0%)  
 Query 1 IISNALK 7  
 Sbjct 795 IISNALK 801

**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



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  $-10\log_{10}(P)$ , 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.
- 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.

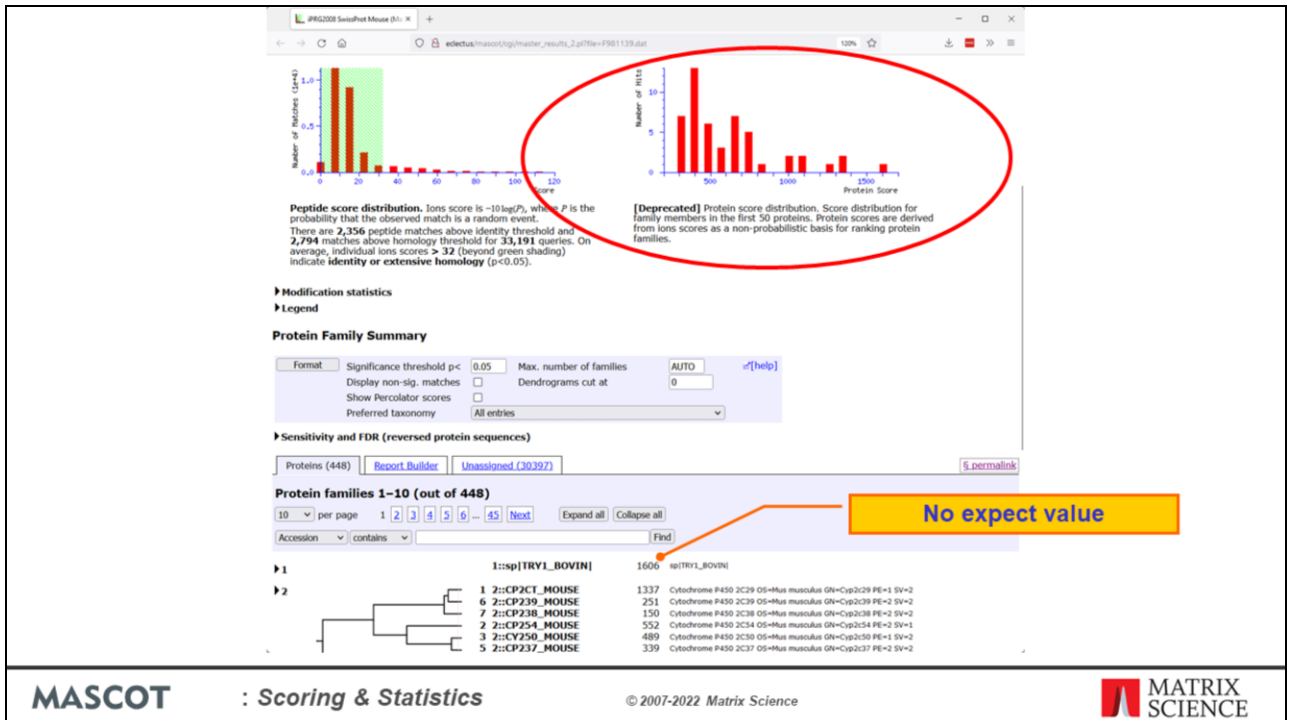
For an MS/MS search, the protein score is **not** statistically rigorous. It is just a way of ranking the protein hits

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 random.

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

$$S = -10\log P = 70$$

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.

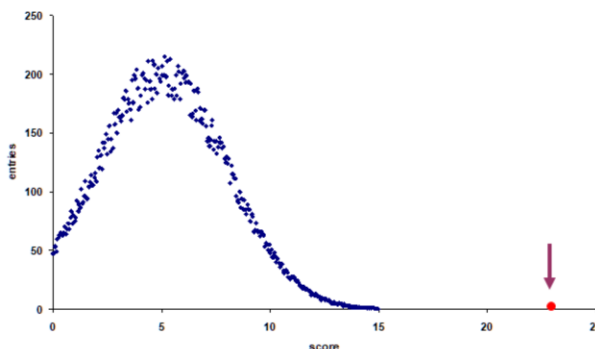
So, for example, if we are comfortable with a 1 in a 20 chance of getting a false positive match, and we are doing a PMF search of a database that contains 500,000 entries, we are looking for a probability of less than  $1 / (20 \times 500,000)$  which is a Mascot score of 70

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



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

Threshold (0): 0

	Score	Mass	Matches	Sequences	empAI
<b>B.1</b> <i>c2:RDH7_MOUSE</i>	1005	38455	45 (45)	12 (12)	<0.02
<b>B.2</b> <i>c2:H17B_MOUSE</i>	597	30949	23 (23)	7 (7)	1.37

**47 peptide matches (16 non-duplicate, 31 duplicate)**  
☒ Auto-fit to window

Query Index	Observed	Mr (exp1)	Mr (calc)	Delta M	Score	Expect	Rev	V	I	Z	Peptide
<i>c21004</i>	580.8522	1150.3079	1150.4043	-0.0927	0	0.002	<b>P</b>				R. TERNERGLUB.V
<i>c21004</i>	581.8107	1180.5149	1179.4370	0.3779	0	0.003	<b>P</b>				R. GYVINGLGLG.V + Glu-Hydro-Lys (B-term)
<i>c21004</i>	583.3202	1184.4259	1184.4049	0.0190	0	4.5e-005	<b>P</b>				R. TONUMER.V
<i>c21004</i>	587.3600	1192.7054	1192.4354	0.0430	0	0.0024	<b>P</b>				R. VYNEISOMER.V
<i>c21011</i>	602.9400	1201.8659	1201.4941	0.1719	0	0.0014	<b>P</b>				R. VESTERGLUB.V
<i>c21011</i>	605.4708	1208.9310	1208.4305	0.3005	0	0.0029	<b>P</b>				R. VYNEISOMER.V + Oxidation (M)
<i>c21024</i>	619.8346	1229.6584	1229.6900	-0.0314	0	0.0015	<b>P</b>				R. ESEYFUNK.V
<i>c21040</i>	641.5011	1280.3074	1280.7074	-0.2800	0	7.3e-006	<b>P</b>				R. VLACLSYK.O
<i>c21040</i>	644.8396	1286.8646	1286.7114	0.1532	0	0.0013	<b>P</b>				R. MESTYK.A
<i>c21040</i>	694.4422	1386.8490	1386.4301	0.1797	0	1.9e-006	<b>P</b>				R. VORAPFGLG.A
<i>c21070</i>	700.4364	1399.8587	1399.7205	0.1322	0	0.019	<b>P</b>				R. VORAPFGLG.A
<i>c21070</i>	710.0087	1418.0029	1417.8624	0.1393	0	3.2e-005	<b>P</b>				R. LATVILLYK.V
<i>c21070</i>	740.9415	1479.8605	1479.7013	0.0872	0	5.4e-006	<b>P</b>				R. LARGYFUNK.V
<i>c21070</i>	815.8047	1543.1924	1542.7912	0.4011	1	0.015	<b>P</b>				R. VORAPFGLG.A
<i>c21083</i>	841.1233	1620.3481	1619.9124	0.4354	0	0.0063	<b>P</b>				R. VORAPFGLG.A
<i>c21084</i>	811.2202	1620.4259	1619.9124	0.5133	0	7.7e-007	<b>P</b>				R. TROIVNATQYK.V

**9** **2:MGST1\_MOUSE** 863 **Proteome 5 to protease 1** OS=Mus musculus (Uniprot) PE=1 SV=1  
**10** **2:RLZA\_MOUSE** 770 **Non-reducing protein L7a** OS=Mus musculus (Uniprot) PE=2 SV=2

10 per page 1 2 3 4 5 6 7 8 9 10 Next Expand all Collapse all

Not what you expected? Try refine select summary.

**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



In Mascot 2.0, we also started displaying 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

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

If  $P_{\text{threshold}} = 0.05$  and  $S_{\text{threshold}} = 50$

score = 40 corresponds to  $E = 0.5$

score = 50 corresponds to  $E = 0.05$

score = 60 corresponds to  $E = 0.005$

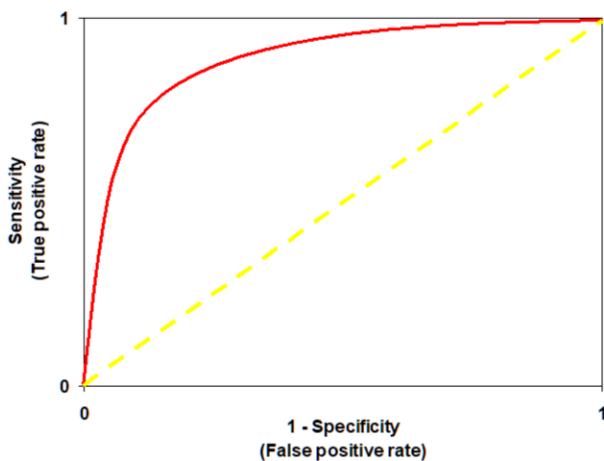
The expectation value does not contain new information. It can be derived directly from the score and the threshold. The advantage is that it tells you everything you need to know in a single number.

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.

# Sensitivity & Specificity

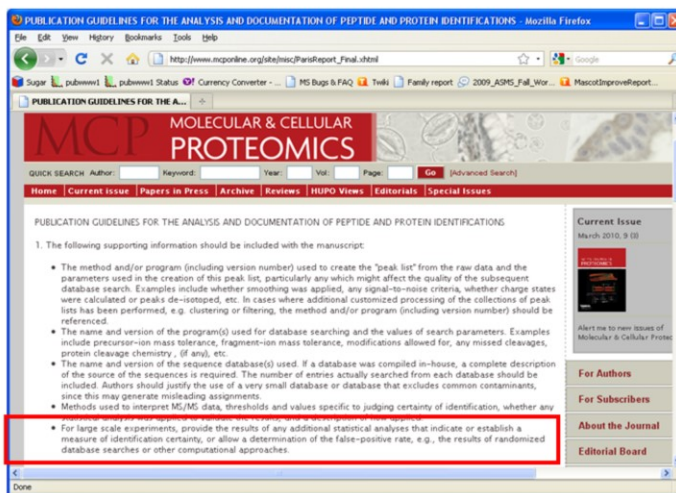


The most important attributes of a scoring scheme are sensitivity and specificity. That is, you want as many correct matches as possible, and as few incorrect matches as possible.

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



**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



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 has just recently been revised.

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)

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, 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 screenshot shows the Mascot MS/MS Ions Search web interface. The 'Decoy' checkbox is highlighted with a red circle. The interface includes fields for 'Your name' (Lou Cline), 'Email' (lou@res.edu), 'Search title' (IPRG2008 SwissProt Mouse), 'Database(s)' (cRAP (AA), SwissProt (AA)), 'Taxonomy' (Mus.), 'Enzyme' (Trypsin/P), 'Quantitation' (None), 'Fixed modifications' (ITRAQ4plex (K), ITRAQ4plex (N-term), Methylthio (C)), 'Variable modifications' (Acetyl (Protein N-term), Glu->pyro-Glu (N-term Q), Oxidation (M)), 'Peptide tol.' (0.9 Da), 'MS/MS tol.' (0.6 Da), 'Data file' (D:\IPRG2008\imgf\merged.mgf), 'Data format' (Mascot generic), 'Instrument' (ESI-TRAP), 'Precursor' (m/z), 'Error tolerant' (checkbox), 'Report top' (AUTO hits), and 'Start Search' and 'Reset Form' buttons.

**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



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. 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.

Because more and more people wish to perform decoy searches routinely, we've added this into Mascot as a built-in part of the search. 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.



The screenshot displays the Mascot search results interface. At the top, the browser address bar shows the URL: `archive-win10/mascot_2_7_00/cgi/master_results_2_`. The page title is "Protein Family Summary".

**Protein Family Summary**

Format: Significance threshold p < 0.0050! Max. number of families: AUTO [help]  
 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: 302 11 3.64%  
 PSMs: above homology 1821 18 **0.99%** Adjust to: 1%  
 Decoy results are available in the decoy report.

Proteins (302) [Report Builder] [Unassigned (31370)] [permalink]

**Protein families 1-10 (out of 280)**

10 per page 1 2 3 4 5 6 28 Next Expand all Collapse all  
 Accession contains Find Clear

**1:sp|TRY1\_BOVIN|** 1153 sp|TRY1\_BOVIN|  
 854 Endoplasmic reticulum chaperone BIP OS=Mus musculus OX=10090 GN=...  
 102 Heat shock 70 kDa protein 1-like OS=Mus musculus OX=10090 GN=Hsp...  
 202 Heat shock cognate 71 kDa protein OS=Mus musculus OX=10090 GN=H...

**2:CYB5\_MOUSE** 830 Cytochrome b5 OS=Mus musculus OX=10090 GN=Cyb5a PE=1 SV=2  
 697 Cytochrome P450 2C29 OS=Mus musculus OX=10090 GN=Cyp2c29 PE=...  
 141 Cytochrome P450 2C39 OS=Mus musculus OX=10090 GN=Cyp2c39 PE=...  
 311 Cytochrome P450 2C34 OS=Mus musculus OX=10090 GN=Cyp2c34 PE=...  
 251 Cytochrome P450 2C50 OS=Mus musculus OX=10090 GN=Cyp2c50 PE=...  
 151 Cytochrome P450 2C37 OS=Mus musculus OX=10090 GN=Cyp2c37 PE=...

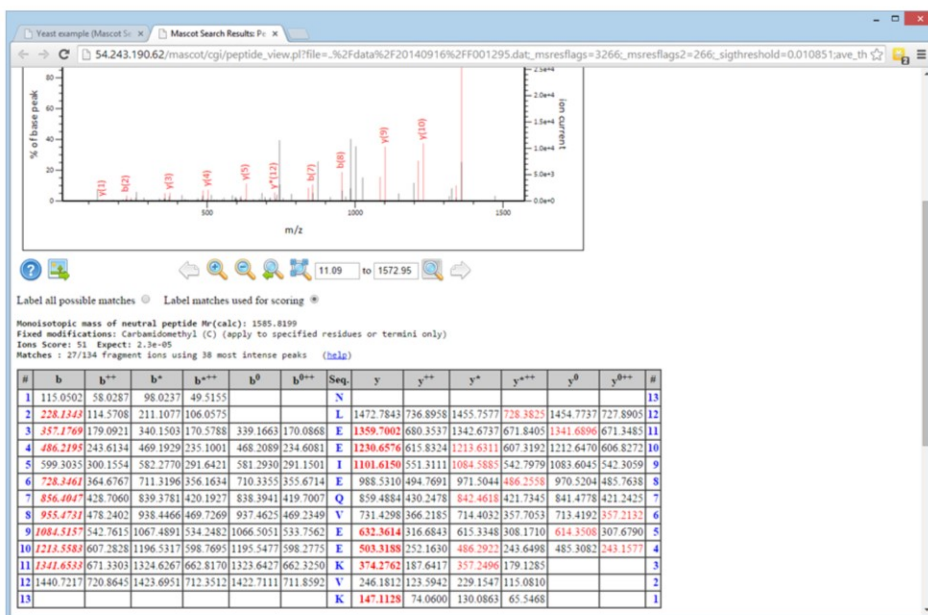
**MASCOT : Scoring & Statistics** © 2007-2022 Matrix Science **MATRIX SCIENCE**

The significance threshold has been automatically adjusted from 0.05 to 0.005054.

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 link here, then you will see the results from searching the randomised database.







**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



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	4699	219	4.66%	
<input type="text" value="Sequences"/> above <input type="text" value="homology"/>	25588	296	1.16%	<input type="button" value="Adjust to"/> <input type="text" value="1%"/>

Decoy results are available in [the decoy report](#).

New in Mascot Server 2.7 is the Protein FDR value which is automatically calculated when the decoy search option is selected.

It is displayed above the peptide/PSM decoy results.

## Based on the following assumptions and definitions:

Only significant peptide sequence matches(PSMs) are used

A protein family member may represent multiple same set proteins

Protein count used for FDR is count of family members

A protein ID is considered true positive if it contains at least one positive PSM

A protein ID is a false positive when all the PSM's are false positives

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.

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/8/11/2405><https://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  $5/500 = 1\%$**

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	<a href="#">[help]</a>
	Display non-sig. matches	<input type="checkbox"/>	Min. number of sig. unique sequences	1	
			Dendrograms cut at	0	
	Preferred taxonomy	All entries			

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

	Target	Decoy	FDR
Protein family members	4699	219	4.66%
Sequences > above homology	25588	296	1.16%

Adjust to 1%

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

Format	Significance threshold p<	0.05	Max. number of families	AUTO	<a href="#">[help]</a>
	Display non-sig. matches	<input type="checkbox"/>	Min. number of sig. unique sequences	2	
			Dendrograms cut at	0	
	Preferred taxonomy	All entries			

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

	Target	Decoy	FDR
Protein family members	3366	6	0.18%
Sequences > above homology	25588	296	1.16%

Adjust to 1%

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 ~4.5%. If we want to lower the Protein FDR we can try adjusting the peptide FDR to 1%. In this case the default peptide FDR is already close to 1% so we can 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 0.18%.

## Adjusting the Protein FDR

The screenshot shows the 'Format' tab in the Mascot Scoring & Statistics interface. The 'Significance threshold p<' is set to 0.005, which is circled in red. Other settings include 'Max. number of families' set to AUTO, 'Display non-sig. matches' unchecked, 'Min. number of sig. unique sequences' set to 1, 'Dendrograms cut at' set to 0, and 'Preferred taxonomy' set to 'All entries'. Below this, the 'Sensitivity and FDR (reversed protein sequences)' section shows a table with columns 'Target', 'Decoy', and 'FDR'. The table has two rows: 'Protein family members' with values 4198, 38, and 0.91%, and 'Sequences' with values 19483, 49, and 0.25%. The 'Sequences' row has dropdown menus for 'Sequences' (set to 'above') and 'homology' (set to 'homology'). An 'Adjust to' button and a '1%' dropdown are also present. A note at the bottom states 'Decoy results are available in [the decoy report](#)'.

	Target	Decoy	FDR
Protein family members	4198	38	0.91%
Sequences	19483	49	0.25%

**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science

**MATRIX**  
SCIENCE

Alternatively if 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 0.005 and clicked the format button. The resulting Protein FDR is approximately 1%.

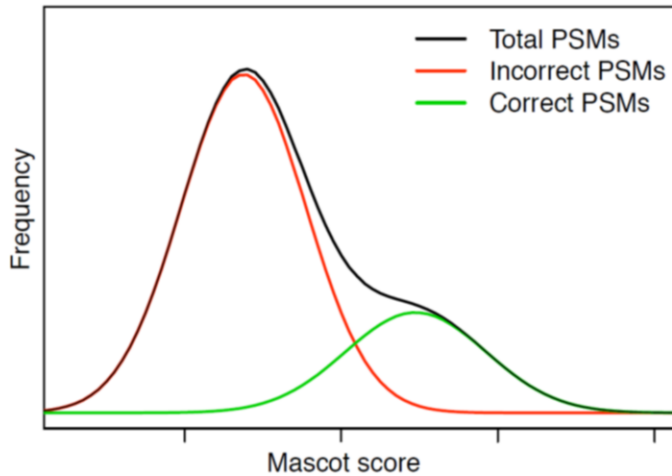
The current HUPO guidelines Interpretation Guidelines for large-scale results recommend adjusting the settings to lower than 1% protein-level global FDR so after formatting this search result would meet those guidelines.

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 optimisation



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.

## Sensitivity optimisation

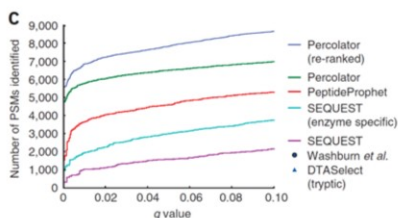


*Anal. Chem.* 2002, 74, 5383–5392

### Empirical Statistical Model To Estimate the Accuracy of Peptide Identifications Made by MS/MS and Database Search

Andrew Keller,<sup>\*,†</sup> Alexey I. Nesvizhskii,<sup>\*,‡</sup> Eugene Kolker, and Ruedi Aebersold

*Institute for Systems Biology, 1441 North 34th Street, Seattle, Washington 98103*



**NATURE METHODS** | VOL.4 NO.11 | NOVEMBER 2007 | 923

### Semi-supervised learning for peptide identification from shotgun proteomics datasets

Lukas Käll<sup>1</sup>, Jesse D Canterbury<sup>1</sup>, Jason Weston<sup>2</sup>, William Stafford Noble<sup>1,3</sup> & Michael J MacCoss<sup>1</sup>

**MASCOT**

: Scoring & Statistics

© 2007-2022 Matrix Science

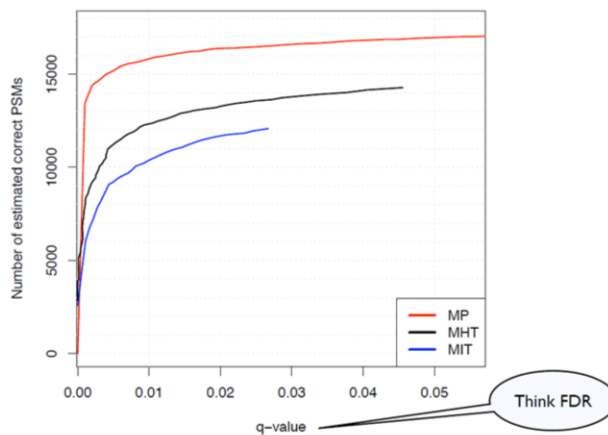
**MATRIX  
SCIENCE**

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 termini.

A more recent development has been to use the matches from a decoy database as negative examples for the classifier. 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).

This 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 screenshot shows the Mascot search engine interface. The 'Protein Family Summary' section is visible, with the 'Show Percolator scores' checkbox checked and highlighted by a red circle. Below this, the 'Sensitivity and FDR (reversed protein sequences)' section is shown, with the 'Report Builder' tab selected. The 'Protein families 1-10 (out of 3613)' section displays a list of protein families with their accession numbers and descriptions. The 'Accession' column shows '1 3:Q15149' and '2 3:PSB107'. The 'Description' column shows '3947 Protein OS=Homo sapiens OX=9606 GN=PLEC PE=1 SV=3' and '958 Egr1-like OS=Homo sapiens OX=9606 GN=EGFRC1 PE=1 SV=3'.

**MASCOT** : Scoring & Statistics

© 2007-2022 Matrix Science



The developers of Percolator have kindly agreed to allow us to distribute and install Percolator as part of Mascot 2.3 and later. 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	■	■	R.VIVETIK.G
0.2841	0	17	0.018	▶1	U	■	■	R.VIVETIK.G
0.4581	0	18	0.015	▶1	U	■	■	R.VIVETIK.G
0.0517	0	21	0.0087	▶1	U	■	■	K.VLESDLK.K
0.2227	0	25	0.0031	▶1	U	■	■	K.ITITNDQNR.L

Score > 13 indicates identity

**MASCOT** : Scoring & Statistics

© 2007-2022 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.