

# Improvements to the Percolator Algorithm for Peptide Identification from Shotgun Proteomics Data Sets

Marina Spivak,<sup>†,‡</sup> Jason Weston,<sup>†</sup> Léon Bottou,<sup>†</sup> Lukas Käll,<sup>‡,§</sup> and William Stafford Noble<sup>\*,‡,||</sup>

NEC Labs America, Princeton, New Jersey 08540, Computer Science Department, New York University, New York, New York 10003, Department of Genome Sciences, University of Washington, Seattle, Washington 98195, Center for Biomembrane Research, Department of Biochemistry and Biophysics, Stockholm University, Sweden, and Department of Computer Science and Engineering, University of Washington, Seattle, Washington 98195

Received December 23, 2008

**Abstract:** Shotgun proteomics coupled with database search software allows the identification of a large number of peptides in a single experiment. However, some existing search algorithms, such as SEQUEST, use score functions that are designed primarily to identify the best peptide for a given spectrum. Consequently, when comparing identifications across spectra, the SEQUEST score function Xcorr fails to discriminate accurately between correct and incorrect peptide identifications. Several machine learning methods have been proposed to address the resulting classification task of distinguishing between correct and incorrect peptide-spectrum matches (PSMs). A recent example is Percolator, which uses semisupervised learning and a decoy database search strategy to learn to distinguish between correct and incorrect PSMs identified by a database search algorithm. The current work describes three improvements to Percolator. (1) Percolator's heuristic optimization is replaced with a clear objective function, with intuitive reasons behind its choice. (2) Tractable nonlinear models are used instead of linear models, leading to improved accuracy over the original Percolator. (3) A method, Q-ranker, for directly optimizing the number of identified spectra at a specified  $q$  value is proposed, which achieves further gains.

**Keywords:** shotgun proteomics • tandem mass spectrometry • machine learning • peptide identification

## 1. Introduction

A shotgun proteomics mass spectrometry experiment produces, for a given biological sample, a collection of spectra, each of which may be mapped back to its generating peptide using either de novo or database search techniques (reviewed

in refs 25 and 26). Critical to any database search procedure is the score function that evaluates the quality of the match between an observed spectrum and a candidate peptide. This function plays two complementary roles. First, the function ranks candidate peptides relative to a single spectrum, producing a single, top-scoring peptide-spectrum match (PSM) for each spectrum. Second, the function ranks the PSMs from different spectra with respect to one another. This latter absolute ranking task is intrinsically more difficult than the relative ranking task. A perfect absolute ranking function is by definition also a perfect relative ranking function, but the converse is not true because PSM scores may not be well-calibrated from one spectrum to the next.

A variety of approaches have been developed to learn PSM scoring functions from real data. Typically, the input to these PSM postprocessing methods is the relative score, as well as properties of the spectrum, the peptide, and features that represent the quality of the PSM. PeptideProphet,<sup>19</sup> for example, uses four statistics computed by the SEQUEST database search algorithm as input to a linear discriminant analysis classifier. The system is trained from labeled correct and incorrect PSMs derived from a purified sample of known proteins. Other approaches use alternative feature representations or classification algorithms, such as support vector machines (SVMs)<sup>1</sup> or decision trees.<sup>11</sup>

One drawback to these machine learning approaches is that they often do not generalize well across different machine platforms, chromatography conditions, etc. Consequently, when the experimental conditions change, a new training set must be acquired, and this acquisition and training can be expensive.

To combat this problem, several methods have been described that adjust the parameters of the model with respect to each new data set. PeptideProphet, for example, uses a fixed linear discriminant function but couples it with a postprocessor that maps the resulting unitless discriminant score to an estimated probability. In the original version of PeptideProphet,<sup>19</sup> this mapping function was learned from each data set in an unsupervised fashion (i.e., without knowing which PSMs are correct and which are incorrect) using the expectation-maximization (EM) algorithm.<sup>9</sup>

Subsequently, several algorithms have been described that use *semisupervised learning* to adjust model parameters with

\* To whom correspondence should be addressed. E-mail: noble@gs.washington.edu.

<sup>†</sup> NEC Research.

<sup>‡</sup> New York University.

<sup>§</sup> Department of Genome Sciences, University of Washington.

<sup>||</sup> Stockholm University.

<sup>||</sup> Department of Computer Science and Engineering, University of Washington.

respect to each new data set. In contrast to supervised learning, in which the given training set is fully labeled, a semisupervised learner is provided with a partially labeled training set. In the context of PSM scoring, these labels are created using a decoy database.<sup>24</sup> Each spectrum is searched once against the real (“target”) protein database and once against a decoy database composed of reversed,<sup>24</sup> shuffled,<sup>20</sup> or Markov-chain generated proteins.<sup>6</sup> Matches to the target database are unlabeled—they may or may not be correct (we expect 50–90% are false positives), but matches to the decoy database can be confidently labeled “incorrect”.

The semisupervised version of PeptideProphet<sup>5</sup> uses decoy PSMs to improve the mapping from discriminant scores to probabilities. During the EM step, PeptideProphet includes decoy PSMs, forcing them to be labeled “incorrect”. The resulting probabilities are significantly more accurate than probabilities estimated in an unsupervised fashion.

The Percolator algorithm<sup>17</sup> takes the semisupervised approach one step further. Rather than using a fixed discriminant function and employing semisupervised learning as a postprocessor, Percolator solves the entire problem in a semisupervised fashion, learning a function that consistently ranks the decoy PSMs below a subset of high-confidence target PSMs. Percolator uses an iterative, SVM-based algorithm, initially identifying a small set of high-scoring target PSMs, and then learning to separate these from the decoy PSMs. The learned classifier is applied to the entire set, and if new high-confidence PSMs are identified, then the procedure is repeated. Critical to the success of the algorithm is a statistical scoring procedure, based on estimated false discovery rates,<sup>2</sup> that prevents explosion of the high-confidence set of PSMs.

A subsequent version of PeptideProphet<sup>10</sup> extends that algorithm in a similar fashion. Like Percolator, the newest version of PeptideProphet adjusts the parameters of the discriminant function to reflect specific features of the data set and allows the algorithm to use more than one PSM for the identification of the best scoring peptide. In addition, the algorithm uses a measure of spectrum quality in its model.

Despite the good performance of Percolator, the algorithm itself is somewhat heuristic; indeed, it is unclear what exactly Percolator optimizes and whether the algorithm’s iterative optimization process provably converges. The current work proposes a novel, well-founded approach to this problem. Although only some of the matches to the target database are positive examples, we opt to treat this problem as a fully supervised classification problem with noisy labels; that is, we label all the target PSMs “correct” (but some of these are mislabeled) and all the decoy PSMs “incorrect”. However, we define a loss function that does not severely penalize examples that are far from the decision boundary. In this way, incorrect target PSMs do not strongly affect the learning procedure. We show how this choice of loss is superior to more classical choices of loss function and in the linear case how this yields results similar to the original *semisupervised* Percolator algorithm. An important benefit of using a fully supervised approach is that, in contrast to Percolator, the new approach defines a clear, intuitive objective function whose minimization is known to converge. Furthermore, the resulting classifier can be trained with tractable nonlinear models which then significantly improve the results of Percolator. Subsequently, we propose a modification of our algorithm that directly optimizes the number of PSMs relative to a user-specific statistical confidence threshold. This ability to specify the desired con-

fidence threshold a priori is useful in practice and leads to further improvement in the results. The new algorithm, called Q-ranker, is implemented in Crux version 2.0, which is available with source code at <http://noble.gs.washington/proj/crux>.

## 2. Materials and Methods

**2.1. Data Sets.** We used four previously described data sets to test our algorithms.<sup>17</sup> The first is a yeast data set containing 69 705 target PSMs and twice that number of decoy PSMs. These data were acquired from a tryptic digest of an unfractionated yeast lysate and analyzed using a 4 h reverse-phase separation. Throughout this work, peptides were assigned to spectra by using SEQUEST with no enzyme specificity and with no amino acid modifications enabled. The next two data sets were derived from the same yeast lysate but treated by different proteolytic enzymes: elastase and chymotrypsin. These data sets, respectively, contain 57 860 and 60 217 target PSMs and twice that number of decoy PSMs. The final data set was derived from a *Caenorhabditis elegans* lysate proteolytically digested by trypsin and processed analogously to the yeast data sets.

Each PSM was represented using the 17 features listed in Table 1. Note that, originally, Percolator used 20 features. In this work, we removed three features that exploit protein-level information because of the difficulty of accurately validating, via decoy database search, methods that use this type of information. We also defined 20 additional features for each peptide, also defined in Table 1, corresponding to the counts of amino acids in the given peptide. Using these additional features yields a feature vector of length 37.

**2.2. Statistical Confidence Estimates.** Throughout this work, we use the  $q$  value<sup>28</sup> as a statistical confidence measure assigned to each PSM. If we specify a score threshold  $t$  and refer to PSMs with scores better than  $t$  as *accepted* PSMs, then the *false discovery rate* (FDR) is defined as the percentage of accepted PSMs that are incorrect (i.e., the peptide was not present in the mass spectrometer when the spectrum was produced). The  $q$  value is defined as the minimal FDR threshold at which a given PSM is accepted. Note that the  $q$  value is a general statistical confidence metric that is unrelated to the Qscore method for evaluating SEQUEST results.<sup>24</sup>

We calculate  $q$  values by using decoy PSMs,<sup>18</sup> derived by searching each spectrum against a database of shuffled protein sequences. Denote the scores of target PSMs  $f_1, f_2, \dots, f_{m_f}$  and the scores of decoy PSMs  $d_1, d_2, \dots, d_{m_d}$ . For a given score threshold  $t$ , the number of accepted PSMs (positives) is  $P(t) = |\{f_i > t; i = 1, \dots, m_f\}|$ . The estimated number of false positives among the positives is given by  $E(FP(t)) = \pi_0(m_f)/(m_d)|\{d_i > t; i = 1, \dots, m_d\}|$ , where  $\pi_0$  is the estimated proportion of target PSMs that are incorrect. In this work, as previously,<sup>17</sup> we use a fixed  $\pi_0 = 0.9$ . We can then estimate the FDR at a given threshold  $t$  as

$$E\{FDR(t)\} = \frac{\pi_0 \frac{m_f}{m_d} |\{d_i > t; i = 1, \dots, m_d\}|}{|\{f_i > t; i = 1, \dots, m_f\}|}$$

The  $q$  value assigned to score  $f_i$  is then

$$q(f_i) = \min_{f_j \leq f_i} E\{FDR(f_j)\}$$

**Table 1.** Features Used to Represent PSMs<sup>a</sup>

1	XCorr	cross correlation between calculated and observed spectra
2	$\Delta C_n$	fractional difference between current and second best XCorr
3	$\Delta C_n^5$	fractional difference between current and fifth best XCorr
4	Sp	preliminary score for peptide versus predicted fragment ion values
5	ln(rSp)	the natural logarithm of the rank of the match based on the Sp score
8	mass	the observed mass $[M + H]^+$
6	$\Delta M$	the difference in calculated and observed mass
7	abs( $\Delta M$ )	the absolute value of the difference in calculated and observed mass
9	ionFrac	the fraction of matched b and y ions
10	ln(NumSp)	the natural logarithm of the number of database peptides within the specified $m/z$ range
11	enzN	Boolean: is the peptide preceded by an enzymatic (tryptic) site?
12	enzC	Boolean: does the peptide have an enzymatic (tryptic) C-terminus?
13	enzInt	number of missed internal enzymatic (tryptic) sites
14	pepLen	the length of the matched peptide, in residues
15–17	charge 1–3	three Boolean features indicating the charge state
18–37	A, ..., Y	counts of each of the 20 amino acids

<sup>a</sup> The first 10 features are computed by SEQUEST. Features 18–37 are used in section 3.6.

### 3. Results

**3.1. Supervised Algorithm for Target-Decoy Discrimination.** Given a set of examples (PSMs) ( $\mathbf{x}_1, \dots, \mathbf{x}_n$ ) (where the bold face denotes a vector) and corresponding labels ( $y_1, \dots, y_n$ ), the goal is to choose a discriminant function  $f(\mathbf{x})$ , such that

$$f(\mathbf{x}_i) > 0 \quad \text{if } y_i = 1$$

$$f(\mathbf{x}_i) < 0 \quad \text{if } y_i = -1$$

To find  $f(\mathbf{x})$ , we first choose a parametrized family of functions and then search for the function in the family that best fits the empirical data. The quality of the fit is measured using a loss function  $L(f(\mathbf{x}), y)$  which quantifies the discrepancy between the values of  $f(\mathbf{x})$  and the true labels  $y$ .

Initially, we consider the family of functions that are implemented by a linear model:

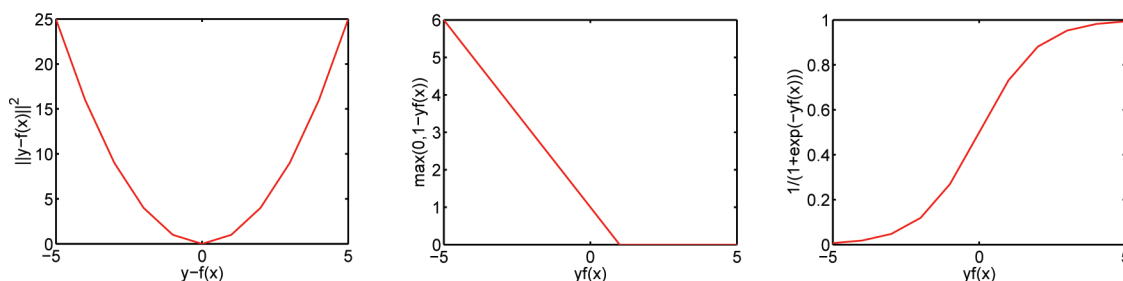
$$f(\mathbf{x}) = \sum_i w_i x_i + b$$

The possible choices of weights define the members of the family of functions.

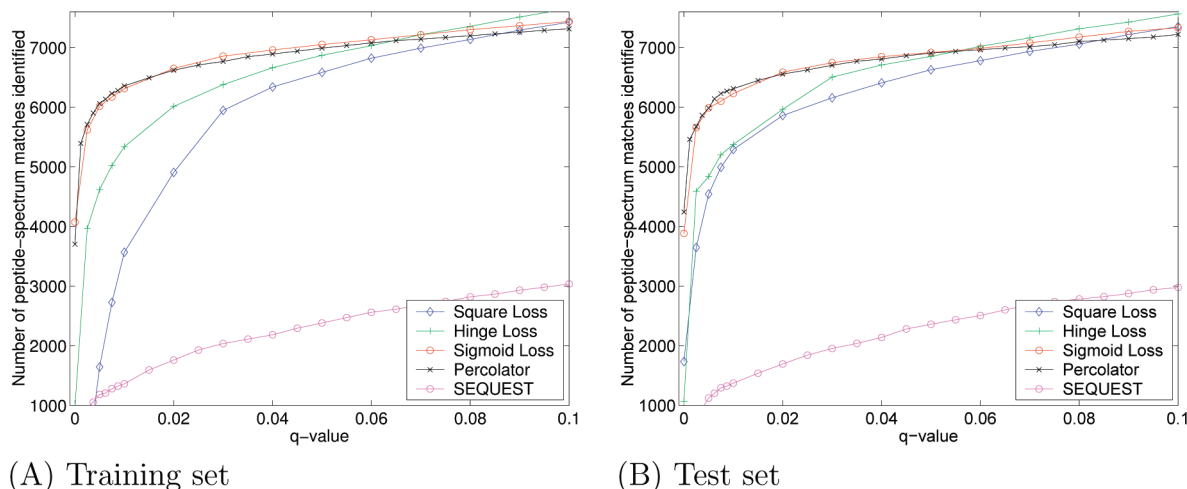
To find the function that best minimizes the loss, we choose to use gradient descent, so the loss function itself must be

differentiable. This requirement prevents us from simply counting the number of mistakes (misclassified examples), which is called the zero-one loss. Typical differentiable loss functions include the squared loss, often used in neural networks,<sup>22</sup> the hinge loss, which is used in support vector machines,<sup>8</sup> and the sigmoid loss. These loss functions are illustrated in Figure 1.

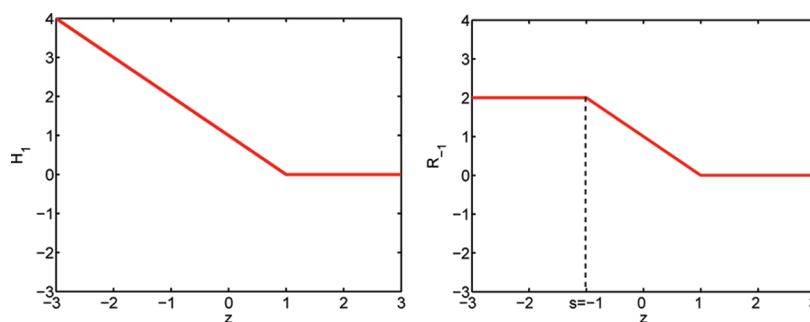
In general, choosing an appropriate loss function is critical to achieving good performance. Insight into choosing the loss function comes from the problem domain. In the current setting, we can safely assume that a significant proportion of the PSMs produced by a given search algorithm are incorrect, either because the score function used to identify PSMs failed to accurately identify the correct peptide or because the spectrum corresponds to a peptide not in the given database, to a peptide with post-translational modifications, to a heterogeneous population of peptides, or to nonpeptide contaminants. Therefore, in this scenario, a desirable loss function will be robust with respect to the multiple false positives in the data. In other words, a desirable loss function will not strongly penalize misclassified examples if they are too far away from the separating hyperplane. Considering the loss functions in Figure 1, the sigmoid loss is the only function with the desired property: when  $y/f(\mathbf{x}) < -5$ , the gradient is close to zero. The squared loss, on the other hand, has a larger gradient for misclassified examples far from the boundary than for examples close to the boundary, whereas the hinge loss penalizes examples linearly (it has a constant gradient if an example is



**Figure 1.** Three types of loss function. Each panel plots the loss as a function of the difference between the true and predicted label. The squared loss  $L(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2$  is often used in regression problems and also in classification.<sup>22</sup> The hinge loss  $L(f(\mathbf{x}), y) = \max(0, 1 - yf(\mathbf{x}))$  is used as a convex approximation to the zero-one loss in support vector machines.<sup>8</sup> The sigmoid loss  $L(f(\mathbf{x}), y) = 1/\exp(1 + f(\mathbf{x}))$  is perhaps less commonly used but is discussed in, for example, refs 23 and 27.



**Figure 2.** Comparison of loss functions. Each panel plots the number of accepted PSMs for the yeast (A) training set and (B) test set as a function of the  $q$  value threshold. Each series corresponds to one of the three loss functions shown in Figure 1, with series for Percolator and SEQUEST included for comparison.



**Figure 3.** “Cutting” the hinge loss makes a sigmoid-like loss called the *ramp loss*. Making the hinge loss have zero gradient when  $z = y_i f(\mathbf{x}_i) < s$  for some chosen value  $s$  effectively makes a piece-wise linear version of a sigmoid function.

incorrectly classified). We therefore conjecture that the sigmoid loss function should work much better than the alternatives.

**3.2. Supervised Learning Yields Performance Comparable to Percolator.** We test this conjecture by measuring the performance of the learned scoring function using a target-decoy search strategy. For this experiment, we use a collection of spectra derived via microcapillary liquid chromatography MS/MS of a yeast whole cell lysate. These spectra were searched using SEQUEST<sup>13</sup> against one target database and two independently shuffled decoy databases, producing a collection of PSMs. For a given ranking of target PSMs, we use the corresponding collection of decoy PSMs to estimate  $q$  values (section 2.2). Our goal is to correctly identify as many PSMs as possible for a given  $q$  value. Therefore, in Figure 2, we plot the number of identified PSMs as a function of  $q$  value threshold.

To ensure a valid experiment, we split the target and decoy PSMs into two equal parts. We train on the data set composed of the first half of positives and negatives, and we use the second half of the data as a testing set. The  $q$  value estimates are derived from the test set, not the training set. This approach is more rigorous than the methodology employed in ref 17, in which the positive examples were used both for training and for testing. However, the similarity between Figure 2A and B indicates that overfitting is not occurring. Nonetheless, in subsequent experiments, we retain a full separation of the train and test sets.

Figure 2 compares the performance of ranking by XCorr, Percolator, and a linear model trained using three different loss

functions. The figure shows that, for example, the Percolator algorithm identifies 5917 PSMs at a  $q$  value threshold of 0.01. As expected, the sigmoid loss dominates the other two loss functions that we considered, square loss and hinge loss.

In fact, the linear model with the sigmoid loss achieves almost identical results to the Percolator algorithm. This concordance can be explained in the following way. Percolator also uses a linear classifier (a linear SVM) with a hinge loss function. However, on each iteration, *only a subset of the positive examples is used as labeled training data* according to the position of the hyperplane. The rest of the positive examples that have a small value of  $y_i f(\mathbf{x}_i)$  are ignored during training. Consequently, one can say that their gradient is zero; hence, the hinge loss function is “cut” at a certain point so that it no longer linearly penalizes mistakes at any distance, as shown in Figure 3. A cut hinge loss is effectively a piece-wise linear version of a sigmoid function. Indeed, such a cut hinge loss has been used before and is referred to as a *ramp loss*.<sup>7</sup> By using a sigmoid loss function, we have thus developed a method that explains the heuristic choices of the Percolator algorithm but instead implements a direct, intuitive objective function. Hereafter, we refer to this method as “direct classification”.

**3.3. Nonlinear Families of Discriminant Functions Yield Improved Performance.** Having established that direct classification using a linear model performs as well as Percolator on this data set, we next consider a nonlinear family of functions by considering two-layer neural networks



$$f(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x}) + b$$

where  $h_k(\mathbf{x})$  is defined as  $\tanh((w^k)^T \mathbf{x} + b_k)$ , and  $w^k$  and  $b_k$  index the weight vector and threshold for the  $k$ th hidden unit.

We can choose the capacity of our nonlinear family of discriminant functions by increasing or decreasing the number of hidden units of our neural network. On the basis of preliminary experiments with the yeast training data set, we chose the first layer to have five linear hidden units. An experimental comparison in Figure 4 shows that a nonlinear classifier outperforms the linear model on the same data set as before. For every  $q$  value in the plot, the nonlinear model (the solid blue line with the label “direct classification (linear)”) produces as many or more PSMs than its linear counterpart (solid black line labeled “direct classification (nonlinear)”).

**3.4. Q-ranker Algorithm for Optimizing Relative to a Specified  $q$  Value.** We have established that framing our problem as a supervised classification task, utilizing nonlinear models, yields slightly improved results compared with Percolator’s semisupervised approach. We now show that reformulating the problem as a ranking task, rather than as a classification task, leads to even better performance.

Generally speaking, the goal of many shotgun proteomics experiments is to identify as many proteins as possible at a given  $q$  value threshold. For the peptide identification problem, this task corresponds to finding a ranking of PSMs that maximizes the number of accepted PSMs for a specified  $q$  value threshold. To solve this ranking problem directly, we therefore assume that the user specifies a particular desired  $q$  value threshold a priori. We then search for a ranking that is optimal with respect to the given  $q$  value. A standard formulation for solving the ranking problem is the ranking SVM,<sup>15,16</sup> which can be stated as follows:

$$\min \|w\|^2 \quad (1)$$

subject to

$$w^T \mathbf{x}_i \geq w^T \mathbf{x}_j + 1 \quad \text{if } y_i = 1 \text{ and } y_j = -1 \quad (2)$$

This algorithm reorders the examples so that larger values of  $w^T \mathbf{x}$  correspond to positive examples. Note that, compared to the classification problem posed before, this formulation no longer has a threshold  $b$  because a class label is no longer predicted, only an ordering. The ranking formulation is equivalent to optimizing the area under the receiver operating characteristic (ROC) curve<sup>14</sup> and hence would optimize all  $q$  values at once. The optimization tries to satisfy every pairwise ordering constraint. Again, as in the classification problem, because we expect 50–90% of the positive examples are false positives, the objective function will pay too much attention to these examples.

However, if optimization of only a certain  $q$  value is desired, then reordering of examples far beyond the  $q$  value threshold point on either side of the boundary will not have an effect on the  $q$  value of interest. Therefore, we instead focus on a subset of examples in the vicinity of the  $q$  value cutoff and seek to reorder the examples specifically in this region.

The proposed algorithm is thus as follows. We first find a general discriminant  $f(\mathbf{x})$  using the direct classification algo-

rithm described in the previous section. We then specify a  $q$  value to be optimized and focus sequentially on several intervals in the data set chosen in the vicinity of the specified  $q$  value. The selection of intervals is heuristic and in our case involves defining a set  $\hat{Q}$  of  $q$  value thresholds 0 to 0.1 with a step size of 0.01 and iterating over these steps. The interval  $\epsilon$  is set to equal twice the number of peptides up to the threshold point. In the course of training, we record the best result for the specified  $q$  value after each epoch. A pseudocode description of the direct ranking algorithm for specified  $q$  values (Q-ranker) is given in Algorithm 1.

Q-ranker can be extended trivially to search for optimal solutions to several  $q$  values at once by recording the best network for each of the specified  $q$  values after each epoch. In all the experimental runs presented below, the set  $\hat{Q}$  of threshold  $q$  values also served as a set of specified  $q$  values.

**Algorithm 1 The Q-ranker algorithm.** The input variables are the training set  $X$  of PSM feature vectors, the corresponding binary labels  $Y$ , indicating which PSMs are targets and which are decoys, the set  $Q$  of specified  $q$  values, the set  $\hat{Q}$  of threshold  $q$  values and the number  $n$  of training iterations. The **chooseRandom** subroutine selects a random positive or negative (depending on the first, Boolean parameter) example  $x$  that satisfies  $|f(\mathbf{x})| < \epsilon$ . The **gradientStep** subroutine makes a gradient step to satisfy the constraint  $f(\mathbf{x}^+) > f(\mathbf{x}^-) + 1$ . The algorithm returns the learned weight vector  $w$ .

```

1: procedure Q-RANKER( $X, Y, Q, \hat{Q}, n$ )
2:    $w \leftarrow$  initialize using direct classification ▷ Solve the direct classification problem.
3:   for  $q_t \in Q$  do
4:     for  $q \in \hat{Q}$  do
5:        $t \leftarrow$  computeThreshold( $X, Y, w, q$ ) ▷ Calculate the threshold corresponding to  $q$ .
6:        $\epsilon \leftarrow 2 * |\{x \in X | f(\mathbf{x}) > t\}|$ 
7:       for  $i \leftarrow 1 \dots n$  do
8:          $x^+ \leftarrow$  chooseRandom(TRUE,  $X, Y, w, \epsilon$ ) ▷ Randomly select a pair of examples.
9:          $x^- \leftarrow$  chooseRandom(FALSE,  $X, Y, w, \epsilon$ )
10:         $w \leftarrow$  gradientStep( $w, f(\mathbf{x}^+), f(\mathbf{x}^-)$ ) ▷ Update the weights.
11:      end for
12:    end for
13:    Record best result on  $q_t$ 
14:  end for
15:  return ( $w$ )
16: end procedure

```

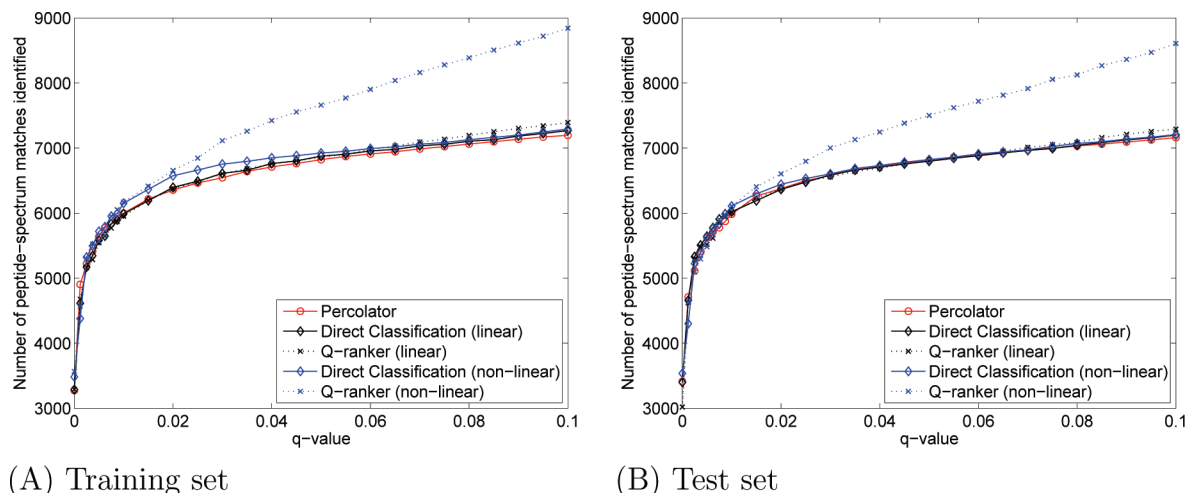
In practice, because Q-ranker focuses on a subset of the training set, we found that use of regularization techniques to control for the model complexity improves our results. In this work, we use the standard weight decay procedure, which optimizes the error function:

$$E' = E + \mu \frac{1}{2} \sum_i w_i^2$$

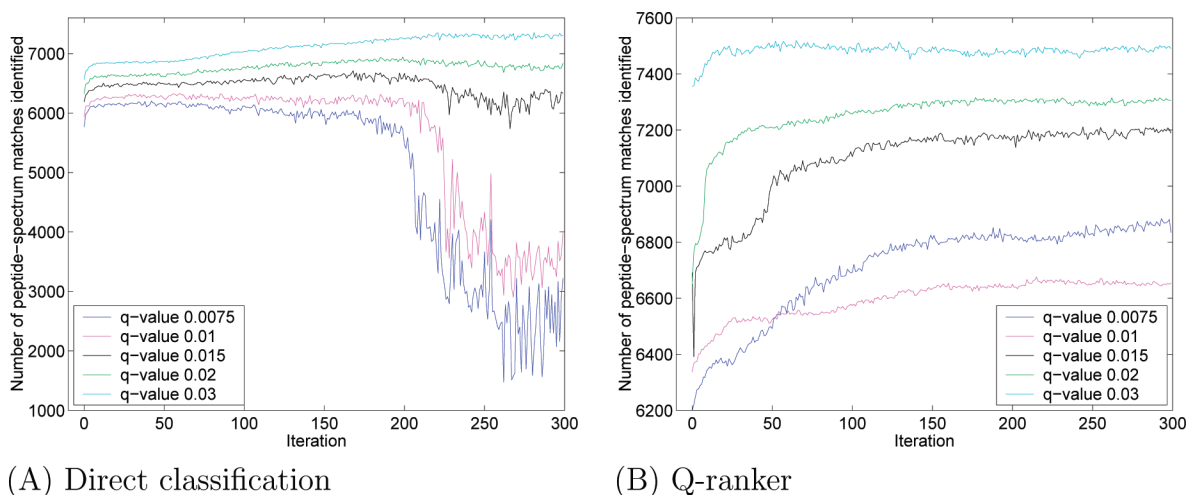
where  $w_i$  are all the weights of the discriminant function  $f(\mathbf{x})$  that we are attempting to learn,  $\mu$  is a weight decay parameter, and  $E$  is the original error function. Before training the network, we perform a three-fold cross-validation procedure to choose the learning rate and  $\mu$ .

Q-ranker generalizes the ranking SVM formulation in two ways: (i) this formulation is nonlinear (but does not use kernels); and (ii) if  $\epsilon$  is very large, then the algorithms are equivalent, but as  $\epsilon$  is reduced, our algorithm begins to focus on given  $q$  values.

Interestingly, choosing examples from a certain region of the data set is also roughly equivalent to placing the region of the sigmoid with high gradient over the region of interest about the threshold  $q$  value. Because examples further than  $\epsilon$  are not picked, this approach is equivalent to making a loss function which has gradient zero in those regions. This means that we are able to replace the sigmoid loss function used for training the general neural net with an even more intuitive choice of loss. In particular, here we use a linear



**Figure 4.** Comparison of Percolator, direct classification, and Q-ranker. The figure plots the number of accepted PSMs as a function of  $q$  value threshold for the yeast data set. Each series corresponds to a different ranking algorithm, including Percolator, as well as linear and nonlinear versions of the direct classification algorithm and Q-ranker. The nonlinear methods use five hidden units.



**Figure 5.** Comparison of training optimization methods (iteration vs error rate). The Q-ranker optimization starts from the best result of direct optimization achieved during the course of training and continues for a further 300 iterations. These results are on the training set. Note that for each  $q$  value choice, Q-ranker improves the training error over the best result from the classification algorithm.

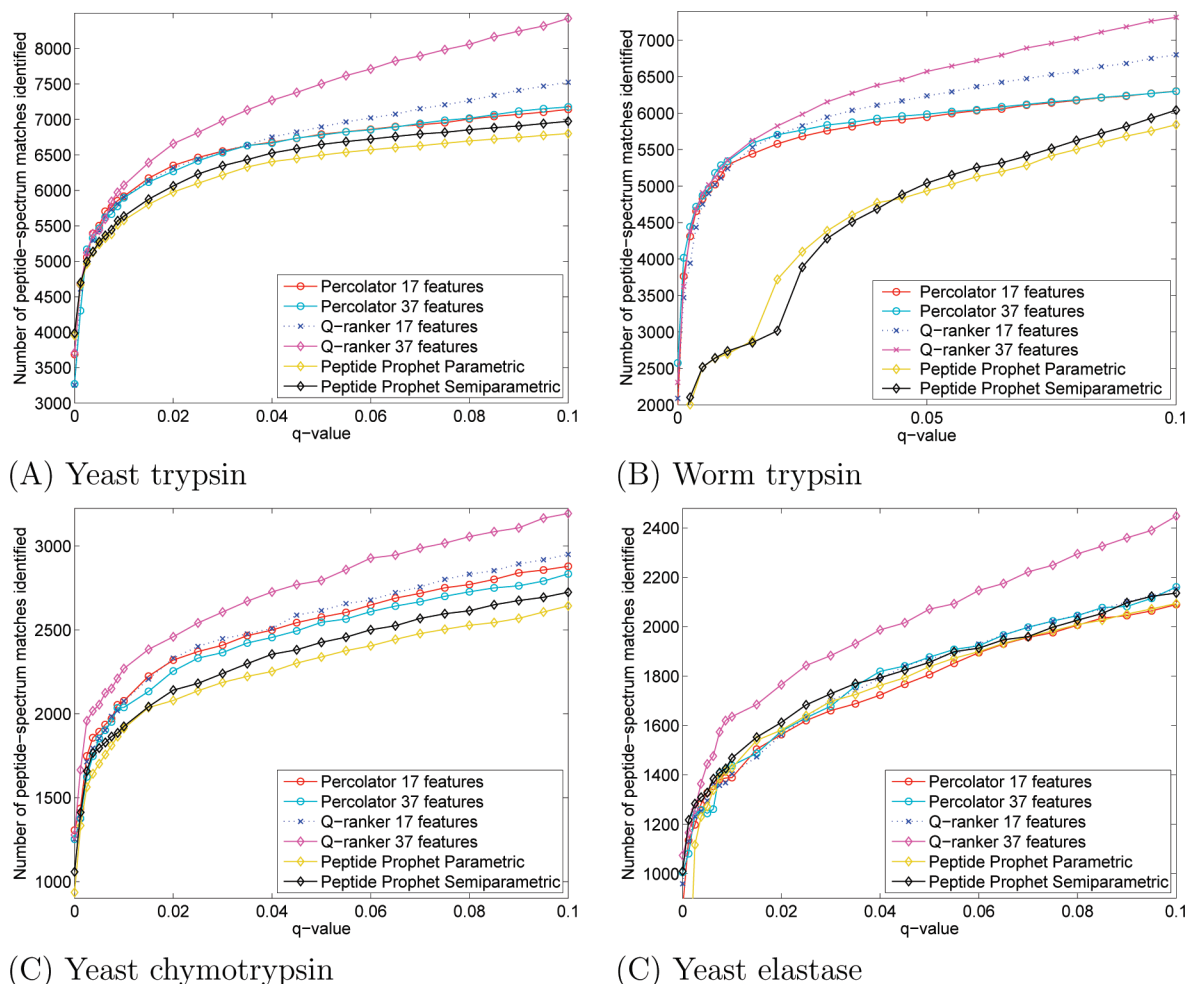
loss  $L(f(\mathbf{x}), y) = |f(\mathbf{x}) - y|$  which effectively becomes a “ramp loss” (cf. Figure 3) centered around the  $q$  value threshold with flat parts at  $\pm\epsilon$ . Because we are solving a ranking problem in the nonlinear case, we now choose a network with the following architecture:

$$f(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x}), \quad \text{where } h_k(\mathbf{x}) = \tanh((w^k)^T \mathbf{x} + b_k)$$

That is, we no longer have a final bias output.

**3.5. Q-ranker Yields Even Better Performance.** We tested our direct classification and Q-ranker algorithms on the typically digested yeast data set in Figure 4. It is clear from this figure that, although the linear Q-ranker algorithm does not improve over the direct classification algorithm, using a nonlinear architecture leads to a large improvement, especially for larger  $q$  values. Other choices of nonlinear architectures (number of hidden units) are given in Supporting Information Figure 1, each leading to improved performance relative to Percolator.

Compared to the direct classification approach described in section 3.1, Q-ranker also yields more consistent training behavior when observed for any given  $q$  value. To illustrate this phenomenon, we fix the interval  $\epsilon$  for the Q-ranker algorithm to be defined by the single threshold corresponding to the specified  $q$  value. Figure 5A shows how the results for different specified  $q$  values change during the course of training the direct classification model. The number of PSMs over lower  $q$  value thresholds (for example, 0.0075, 0.01) reach their peak early during training and then become suboptimal, while the best results for higher  $q$  value thresholds take longer to achieve. This means that, during the course of training, different  $q$  value thresholds are being optimized depending on the number of iterations. In contrast, as shown in Figure 5B, the Q-ranker algorithm learns the best decision boundary for a specified  $q$  value threshold and does not substantially diverge from the best result during further training. This behavior indicates that the algorithm in fact optimizes the desired quantity. In the following experiments, we therefore adopt Q-ranker as our



**Figure 6.** Comparison of PeptideProphet, Percolator, and Q-ranker on four data sets. Each panel plots the number of accepted target PSMs as a function of  $q$  value. The series corresponds to the three different algorithms, including two variants of Q-ranker that use 17 features and 37 features.

algorithm of choice, and we compare it further to Percolator and PeptideProphet.

**3.6. Comparison of Algorithms across Multiple Data Sets.** For our final round of experiments, we compare the performance of Q-ranker, Percolator, and two versions of PeptideProphet—the original parametric version,<sup>19</sup> which assumes that the decoy scores are distributed according to a  $\gamma$  distribution and the target scores according to a Gaussian distribution, and a newer, semiparametric approach,<sup>4</sup> which uses a mixture model of kernel functions to model the two distributions. For both sets of PeptideProphet results, we use the semisupervised version of the algorithm.<sup>5</sup> The same set of decoy PSMs is provided to Percolator, Q-ranker, and PeptideProphet. For Percolator and Q-ranker, we use 50% of the PSMs for training and 50% for testing, as before. PeptideProphet does not provide the ability to learn model parameters on one set of data and apply the learned model to the second; therefore, PeptideProphet results are generated by applying the algorithm to the entire data set. This difference gives an advantage to PeptideProphet because that algorithm learns its model from twice as much data and is not penalized for overfitting.

We report results using either 17 or 37 features, as described in Table 1, for both Percolator and Q-Ranker. Figure 6 shows the results of this experiment, conducted using the four data sets described in section 2.1. Across the four data sets, Q-ranker

consistently outperforms PeptideProphet across all  $q$  value thresholds. The left half of Table 2 shows a detailed comparison of Percolator and Q-ranker on all four data sets using 17 features as input. At  $q$  values of 0.05 or 0.10, Q-ranker yields more accepted target PSMs than either Percolator or PeptideProphet, whereas Percolator performs slightly better for  $q < 0.01$ .

Theoretically, a nonlinear network could yield a larger benefit than a linear model when the input feature space is increased, as long as the model does not overfit. We therefore experimented with extending the PSM feature vectors, adding 20 new features corresponding to the counts of amino acids in the peptide. The results of running Q-ranker with these extended vectors are shown in Figure 6, labeled “Q-ranker 37”. Increasing the number of features gives a larger boost to the performance of the nonlinear version of Q-ranker. The effect is particularly evident on data sets derived from yeast lysate digested with chymotrypsin and elastase. After this extension, Q-ranker identifies more spectra than either of the other algorithms, even at  $q < 0.01$  (right half of Table 2).

Finally, we further investigated the behavior of Q-ranker by measuring the performance of networks trained for a specified  $q$  value on other  $q$  values. We focused on specified  $q$  values 0.01, 0.05, and 0.1. Table 3 shows that, when all 37 features are employed, a network trained for a specified  $q$  value is

**Table 2.** Comparison of Percolator and Q-ranker on 17 and 37 Feature Data Sets<sup>a</sup>

data set	<i>q</i> value	17 features		37 features	
		Percolator	Q-ranker	Percolator	Q-ranker
yeast trypsin	0.01	<b>5917</b>	5885	5983	<b>6072</b>
	0.05	6793	<b>6940</b>	6813	<b>7501</b>
	0.1	7168	<b>7610</b>	7200	<b>8430</b>
yeast elastase	0.01	<b>1389</b>	1380	1491	<b>1615</b>
	0.05	1806	<b>1851</b>	1958	<b>2140</b>
	0.1	2103	<b>2196</b>	2301	<b>2561</b>
yeast chymotrypsin	0.01	2077	<b>2086</b>	2158	<b>2312</b>
	0.05	2576	<b>2620</b>	2680	<b>2844</b>
	0.1	2914	<b>2961</b>	3057	<b>3214</b>
worm trypsin	0.01	<b>5116</b>	5031	5192	<b>5238</b>
	0.05	5864	<b>6119</b>	5830	<b>6419</b>
	0.1	6169	<b>6730</b>	6146	<b>7128</b>

<sup>a</sup> Each entry in the table indicates the number of accepted PSMs for the given algorithm (column) on the given data set at the given specified *q* value (row). Entries in boldface indicate that this algorithm performed better than the other algorithm for this data set and *q* value threshold.

**Table 3.** Q-ranker successfully optimizes the specified *q* value<sup>a</sup>

specified	yeast trypsin			worm trypsin			yeast elastase			yeast chymotrypsin		
	0.01	0.05	0.10	0.01	0.05	0.10	0.01	0.05	0.10	0.01	0.05	0.10
0.01	<b>6072</b>	7453	8360	<b>5238</b>	6412	7098	<b>1615</b>	2054	2395	<b>2312</b>	2843	3199
0.05	6032	<b>7501</b>	8426	<b>5238</b>	<b>6419</b>	7047	<b>1615</b>	<b>2140</b>	<b>2561</b>	2302	<b>2844</b>	3198
0.10	6030	7500	<b>8430</b>	5213	6418	<b>7128</b>	<b>1615</b>	<b>2140</b>	<b>2561</b>	2300	2830	<b>3214</b>

<sup>a</sup> Each entry in the table lists the number of accepted PSMs at a given *q* value (column) obtained by Q-ranker with 37 features when optimizing a specified *q* value (row). Entries in boldface indicate the maximum value within each column. Note that, for each data set, all diagonal entries are in boldface.

consistently better or equal to the performance on this *q* value, compared with networks trained for other specified *q* values.

4. Discussion

In this work, we have performed all of our analyses using a combination of SEQUEST and Percolator. However, the conclusions that we draw here have implications for researchers who do not employ these particular software systems. First, the conclusions likely generalize across search engines. For example, Percolator has previously been demonstrated to work well with the Inspect<sup>17</sup> and MASCOT search engines,<sup>3</sup> so it seems likely that Q-ranker will also generalize to these search engines. Second, we have demonstrated the utility of shifting from a semisupervised framework to a supervised framework with a modified loss function, both in terms of improved understanding of the objective function being maximized and improved discriminative power. A similar shift should be straightforward to apply, for example, to the semisupervised version of PeptideProphet and may result in similar benefits.

Throughout our evaluations, we have focused on maximizing the number of spectra that are correctly assigned a peptide (i.e., the number of accepted PSMs). It is conceivable that a given algorithm might be biased in the types of peptides it can identify. In this case, the relative performance of two peptide identifications could depend on whether we count the number of accepted PSMs or the number of distinct peptides that are identified from a set of spectra. Supporting Information Figure 2 demonstrates that this bias is not occurring in our results: the relative performance of the algorithms that we considered does not change significantly when we count the number of distinct peptides identified.

One surprising result from our experiments is the relatively large benefit provided by amino acid composition features. We hypothesize that this information allows the classifier to learn to expect certain characteristics of a spectrum. For example, the presence of a proline implies a pair of high-intensity peaks

corresponding to the cleavage N-terminal to the proline; the presence of many basic residues leads to more +2 ions, and the presence of many hydrophobic residues leads to more singly charged +1 ions.<sup>21</sup> However, previous experiments with Percolator using amino acid composition features did not yield significant performance improvements. The difference, in the current setting, is that we have switched from a semisupervised to a fully supervised setting. This switch allows us to use a more complex, nonlinear model. In general, a complex model has more opportunity to improve over a simpler model if the feature space is rich. Thus, although a simple linear model such as the one in Percolator cannot fully exploit the richer, 37-dimensional feature space, the nonlinear model can. This conclusion is supported by the observation that adding compositional features also improves the performance of the direct classification method (results not shown).

An alternative, possible explanation for the added discriminative power provided by the amino acid composition feature is that they provide the algorithm with a way to “cheat”. In our experiments, we did not guarantee that the training set and the test set contain disjoint sets of peptides. Hence, an algorithm might overfit on the amino acid composition features and successfully identify the recurrence of a peptide in the train and test sets. To eliminate this alternative explanation, we performed a follow-up experiment in which we prevented the same peptide from occurring in the training and test set. The results, shown in Supporting Information Figure 4 show that the improved performance of Q-ranker over Percolator still holds.

A drawback to using a nonlinear discriminative classifier is the difficulty in interpreting the learned model. In this work, we have focused on optimizing error rate, not interpretability; sometimes it is hard to have both. Indeed, as shown in Supporting Information Figure 5, simply switching to a linear SVM in the direct classification setting yields markedly decreased performance. However, even with a nonlinear model, it is still possible to gain some insight into the relative contributions of the various features by “knocking out” each



feature individually and measuring the performance of the resulting classifier. Supporting Information Table 1 shows the percent reduction in the number of identified PSMs at  $q < 0.01$  when we knock out each feature of Q-Ranker with 17 features. Not surprisingly, the enzymatic features are most significant, followed by the score features (XCorr and  $\Delta C_n$ ). The relatively small percentage decrease for many features suggests that many provide redundant information. A more detailed interpretation of the model could be derived via further knockout experiments aimed at groups of related features, as was done in ref 17.

It is worth noting that the relative performance of the methods that we considered does not change when we use an alternative  $q$  value estimation scheme. Elias et al.<sup>12</sup> advocate estimating the FDR using target-decoy competition (i.e., searching each spectrum against a concatenated database of targets and decoys and only retaining the single top-scoring peptide) and estimating FDRs with respect to the combined collection of target and decoy PSMs. To show that our results do not depend upon our  $q$  value estimation procedure, we report in Supporting Information Figure 3 results analogous to those given in Figure 6, but using FDRs estimated by following the protocol of Elias et al. Even in this case, the Q-ranker algorithm outperforms Percolator and both versions of PeptideProphet.

In general, using a large feature space generally requires a concomitantly large number of training examples. For smaller collections of spectra, or for lower quality spectra in which the effective number of positive examples is small, we would expect a larger feature space to lead to overfitting. In the current version of the software, the user must check for overfitting explicitly and select the regularization parameter explicitly. One focus of our future work will be the implementation and validation of robust methods for avoiding such overfitting, either by adjusting the regularization parameter or reducing the complexity of the model.

## 5. Conclusions

We have described a series of algorithms that improve in various ways upon the Percolator algorithm. Given unlabeled target PSMs and negatively labeled decoy PSMs, Percolator treats the problem as a semisupervised classification problem. In this work, we instead use a supervised approach to the same problem. This change allows us to state an explicit objective function and also allows us to generalize to more powerful, nonlinear models. Finally, if the user is willing to specify a desired confidence threshold, then the Q-ranker algorithm finds an optimal ranking with respect to the specified threshold, yielding consistently improved performance relative to either Percolator or PeptideProphet. Both the direct classification and the Q-ranker algorithms are implemented in the Crux toolkit, which is available with source code from <http://noble.gs.washington.edu/proj/crux>.

**Acknowledgment.** This work was funded by NIH award R01 EB007057.

**Supporting Information Available:** Additional figures and table. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## References

- (1) Anderson, D. C.; Li, W.; Payan, D. G.; Noble, W. S. A new algorithm for the evaluation of shotgun peptide sequencing in proteomics: support vector machine classification of peptide MS/MS spectra and sequest scores. *J. Proteome Res.* **2003**, *2*, 137–146.
- (2) Benjamini, Y.; Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Statist. Soc., Ser. B* **1995**, *57*, 289–300.
- (3) Brosch, M. Yu, L. Hubbard, T., and Choudhary, J. Accurate and sensitive peptide identification with Mascot Percolator. *J. Proteome Res.* [Online early access]. DOI: 10.1021/pr800982s. Published Online: April 1, 2009.
- (4) Choi, H.; Ghosh, D.; Nesvizhskii, A. Statistical validation of peptide identifications in large-scale proteomics using target-decoy database search strategy and flexible mixture modeling. *J. Proteome Res.* **2008**, *7*, 286–292.
- (5) Choi, H.; Nesvizhskii, A. I. Semisupervised model-based validation of peptide identifications in mass spectrometry-based proteomics. *J. Proteome Res.* **2008**, *7*, 254–265.
- (6) Colinge, J.; Masselot, A.; Giron, M.; Dessingy, T.; Magnin, J. OLAV: Towards high-throughput tandem mass spectrometry data identification. *Proteomics* **2003**, *3*, 1454–1463.
- (7) Collobert, R.; Sinz, F.; Weston, J.; Bottou, L. Large scale transductive svms. *J. Mach. Learn. Res.* **2006**, *7*, 1687–1712.
- (8) Cortes, C.; Vapnik, V. Support vector networks. *Mach. Learn.* **1995**, *20*, 273–297.
- (9) Dempster, A. P.; Laird, N. M.; Rubin, D. B. Maximum likelihood from incomplete data via the EM algorithm. *J. R. Statist. Soc.* **1977**, *39*, 1–22.
- (10) Ding, Y.; Choi, H.; Nesvizhskii, A. Adaptive discriminant function analysis and reranking of MS/MS database search results for improved peptide identification in shotgun proteomics. *J. Proteome Res.* **2008**, *7*, 4878–4889.
- (11) Elias, J. E.; Gibbons, F. D.; King, O. D.; Roth, F. P.; Gygi, S. P. Intensity-based protein identification by machine learning from a library of tandem mass spectra. *Nat. Biotechnol.* **2004**, *22*, 214–219.
- (12) Elias, J. E.; Gygi, S. P. Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry. *Nat. Methods* **2007**, *4*, 207–214.
- (13) Eng, J. K.; McCormack, A. L.; Yates, J. R., III. An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database. *J. Am. Soc. Mass Spectrom.* **1994**, *5*, 976–989.
- (14) Hanley, J. A.; McNeil, B. J. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* **1982**, *143*, 29–36.
- (15) Herbrich, R.; Graepel, T.; Obermayer, K. Support vector learning for ordinal regression. In *Proceedings of the Ninth International Conference on Artificial Neural Networks* **1999**, 97, 102.
- (16) Joachims, T. Optimizing search engines using clickthrough data. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)* **2002**, 133–142.
- (17) Käll, L.; Canterbury, J.; Weston, J.; Noble, W. S.; MacCoss, M. J. A semi-supervised machine learning technique for peptide identification from shotgun proteomics datasets. *Nat. Methods* **2007**, *4*, 923–925.
- (18) Käll, L.; Storey, J. D.; MacCoss, M. J.; Noble, W. S. Assigning significance to peptides identified by tandem mass spectrometry using decoy databases. *J. Proteome Res.* **2008**, *7*, 29–34.
- (19) Keller, A.; Nesvizhskii, A. I.; Kolker, E.; Aebersold, R. Empirical statistical model to estimate the accuracy of peptide identification made by MS/MS and database search. *Anal. Chem.* **2002**, *74*, 5383–5392.
- (20) Klammer, A. A.; MacCoss, M. J. Effects of modified digestion schemes on the identification of proteins from complex mixtures. *J. Proteome Res.* **2006**, *5*, 695–700.
- (21) Klammer, A. A.; Reynolds, S. R.; Hoopmann, M.; MacCoss, M. J.; Bilmes, J.; Noble, W. S. Modelling peptide fragmentation with dynamic Bayesian networks yields improved tandem mass spectrometry identification. *Bioinformatics* **2008**, *24*, i348–i356.
- (22) LeCun, Y.; Bottou, L.; Orr, G. B., and Müller K.-R. Efficient backprop. In *Neural Networks: Tricks of the Trade*; Orr, G., Müller, K.-R., Eds.; Springer: Berlin, 1998; pp 9–50.
- (23) Mason, L.; Bartlett, P. L.; Baxter, J. Improved generalization through explicit optimization of margins. *Mach. Learn.* **2000**, *38*, 243–255.
- (24) Moore, R. E.; Young, M. K.; Lee, T. D. Qscore: An algorithm for evaluating sequest database search results. *J. Am. Soc. Mass Spectrom.* **2002**, *13*, 378–386.
- (25) Nesvizhskii, A. I.; Vitek, O.; Aebersold, A. R. Analysis and validation of proteomic data generated by tandem mass spectrometry. *Nat. Methods* **2007**, *4*, 787–797.
- (26) Hernandez, M. M. P.; Appel, R. D. Automated protein identification by tandem mass spectrometry: issues and strategies. *Mass Spectrom. Rev.* **2006**, *25*, 235–254.
- (27) Shen, X.; Tseng, G. C.; Zhang, X.; Wong, W. H. On (psi)-learning. *J. Am. Statist. Assoc.* **2003**, *98*, 724–734.
- (28) Storey, J. D. A direct approach to false discovery rates. *J. R. Statist. Soc.* **2002**, *64*, 479–498.

PR801109K