

Supplementary Table 3. Evaluation of different machine learning approaches' prediction power on training dataset.

	Apparent error rate	
	10X cross validation on training data	Human whole cell lysate
Support Vector Machine	0.28	0.4
Random Forest	0.23	0.33
Neural network	0.23	0.27

Supplementary Table 5. Physiochemical features of peptides used in PPA.

Feature	Description	Publication
BHAR880101	Average flexibility indices	Bhaskaran-Ponnuswamy, 1988
DAWD720101	Size	Dawson, 1972
DAYM780101	Amino acid composition	Dayhoff et al., 1978a
JUKT750101	Amino acid distribution	Jukes et al., 1975
KLEP840101	Net charge	Klein et al., 1984
KYTJ820101	Hydropathy index	Kyte-Doolittle, 1982
OLSK800101	Average internal preferences	Olsen, 1980
VENT840101	Bitterness	Venanzi, 1984
ZIMJ680103	Polarity	
ZIMJ680104	Isoelectric point	
ZIMJ680105	RF rank	Zimmerman et al., 1968
FASG760101	Molecular weight	Fasman, 1976
GRAR740102	Polarity	Grantham, 1974
ARGP820101	Hydrophobicity index	Argos et al., 1982
HOPA770101	Hydration number	Hopfinger, 1971

Supplementary Table 6. Area under the ROC for the seven validation datasets of PPA (Figures 2C, 4, Supplementary Figures 3, 5). Datasets that haven't been generated in this study are cited. Results from ESP predictor (ESP, Fusaro *et al.*, 2009) were compared to PPA in three conditions. a) PPA-15 was based only on 15 physicochemical features without any protein abundance information, b) PPA-eSC used the 15 features and external protein coverage from other datasets/publications (referred to as external sequence coverage, eSC), c) PPA-iSC used 15 features and each dataset's own sequence coverage (referred to as internal sequence coverage), d) PPA-CID was trained on a CID-based dataset (Orbitrap Velos; Wisniewski, *et al.*, 2012), PPA-iBAQ used the 15 features and a normalized iBAQ value as protein abundance (referred to as intensity based absolute quantification, iBAQ; Schwanhäusser *et al.*, 2011). n.pos: Number of detected peptides per experiment. n. neg: Number of not detected peptides.

Sample	Platform	ESP	PPA-15	PPA-eSC	PPA-iSC	PPA-CID	PPA-iBAQ	n.pos	n.neg
Single digests	Q-Exactive	0.726	0.791	0.831	0.841			1445	1397
Pooled single digests	Q-Exactive	0.621	0.699	0.706	0.748			3755	8917
Human lysate	Q-Exactive	0.701	0.747	0.791	0.834	0.758	0.846	4392	14924
Human lysate (Wisniewski, <i>et al.</i> , 2012)	LTQ-Orbitrap Velos	0.765	0.787	0.815	0.835			9240	10076
Yeast lysate	Q-Exactive	0.78	0.783	0.815	0.861			3529	8581
Yeast lysate	TripleToF 5600	0.689	0.7	0.769	0.815		0.819	7402	4708
Yeast lysate (Hebert <i>et al.</i> , 2013)	Orbitrap-Fusion	0.755	0.761	0.805	0.83	0.762		21617	34361

Supplementary Table 7. Statistical analysis of the area under the ROC for the validation datasets (Figures 2C, 4). We applied a modified z-test to compute a p-value (http://www.vassarstats.net/roc_comp.html): a) PPA-15 was based only on 15 physicochemical features, b) PPA-eSC used the 15 features and external protein coverage from other datasets/publications (referred to as external sequence coverage, eSC), c) PPA-iSC used 15 features and each dataset's own sequence coverage (referred to as internal sequence coverage).

Sample	Platform	PPA-15 vs ESP		PPA-eSC vs ESP		PPA-iSC vs ESP		PPA-eSC vs PPA-15		PPA-iSC vs PPA-15		PPA-iSC vs PPA-eSC	
		z value	p-value	z value	p-value	z value	p-value	z value	p-value	z value	p-value	z value	p-value
Single digests	Q-Exactive	5.2	<1E-06	8.694	<1E-06	9.661	<1E-06	3.498	<1E-04	4.469	<1E-03	0.973	0.165
Pooled single digests	Q-Exactive	10.158	<1E-06	11.083	<1E-06	16.894	<1E-06	0.917	0.179	6.655	<1E-06	5.484	<1E-06
Human lysate	Q-Exactive	6.926	<1E-06	14.101	<1E-06	21.385	<1E-06	7.14	<1E-06	14.364	<1E-06	7.2	<1E-06
Human lysate (Wiśniewski, et al, 2012)	LTQ-Orbitrap Velos	4.528	<1E-03	10.646	<1E-06	15.307	<1E-06	6.119	<1E-06	10.784	<1E-06	4.67	<2E-06
Yeast lysate	Q-Exactive	0.426	0.335	5.114	<1E-06	12.609	<1E-06	4.687	<1E-05	12.025	<1E-06	7.336	<1E-06
Yeast lysate	TripleToF 5600	1.647	0.049	12.629	<1E-06	20.793	<1E-06	10.98	<1E-06	19.133	<1E-06	8.175	<1E-06
Yeast lysate (Hebert et al, 2013)	Orbitrap-Fusion	1.943	0.026	16.802	<1E-06	25.842	<1E-06	14.857	<2E-04	23.895	<1E-06	9.034	<1E-06

Supplementary Table 8. Neural network classification feature weights in three neural network models. Model 1: 15 physiochemical features (see Table S1).

Model 2: PPA based on protein abundance as sequence coverage. Model 3 PPA based on protein abundance as average peak intensity of the three highest intense tryptic peptides (TOP-3, Silva *et al.*, 2006). Model 4 PPA based on protein abundance as iBAQ (Schwanhäusser *et al.*, 2011)

		Model 1		Model 2		Model 3		Model 4	
		To hidden neuron 1	To hidden neuron 2	To hidden neuron 1	To hidden neuron 2	To hidden neuron 1	To hidden neuron 2	To hidden neuron 1	To hidden neuron 2
First layer input nodes	Intercept	8.281	1.992	-9.251	-1.485	-1.492	-9.483	-8.265	2.126
	BHAR880101	1.844	-1.481	-1.539	0.854	0.868	-1.522	-1.211	0.366
	DAWD720101	0.392	1.478	0.384	-0.568	-0.589	0.411	0.428	-0.116
	DAYM780101	7.566	0.092	-6.237	-0.409	-0.481	-6.41	-5.285	-0.221
	JUKT750101	-8.684	-0.99	3.896	1.349	1.382	4.038	3.626	0.832
	KLEP840101	-1.267	-0.409	-0.298	0.194	0.164	-0.328	-0.321	0.126
	KYTJ820101	2.001	0.303	-0.609	-0.433	-0.442	-0.622	-0.116	-0.244
	OLSK800101	-3.249	-1.279	-0.025	1.342	1.356	-0.026	-0.889	0.671
	VENT840101	0.86	0.573	-0.514	-0.298	-0.303	-0.527	-0.403	-0.139
	ZIMJ680103	-0.29	0.332	-0.398	-0.272	-0.313	-0.42	-0.174	-0.060
	ZIMJ680104	0.21	0.038	9.001	-0.067	0.223	9.305	6.533	-0.539
	ZIMJ680105	4.761	-1.132	-3.209	0.125	0.111	-3.274	-2.419	0.016
	FASG760101	7.447	-5.158	-10.612	3.283	3.309	-10.743	-8.246	1.443
Hidden layer nodes	GRAR740102	-2.57	8.498	1.637	-5.992	-6.31	1.432	2.611	-2.395
	ARGP820101	0.44	0.663	-0.779	-0.145	-0.164	-0.795	-0.449	0.012
	HOPA770101	-0.35	-2.584	1.687	1.469	1.585	1.75	0.779	0.462
	Abundance	NA	NA	-0.179	-0.431	-0.452	-0.185	0.399	-2.307
	Intercept	-13.854		2.998		2.908		10.883	
	Neuron 1	9.138		-173.345		-8.644		-197.752	
	Neuron 2	6.641		-9.078		-208.71		-17.944	
		Hidden layer neuron to output		Hidden layer neuron to output		Hidden layer neuron to output		Hidden layer neuron to output	