

The Molecular Mass and Isoelectric Point of Plant Proteomes

Tapan Kumar Mohanta^{*1}, Abdullatif Khan¹, Abeer Hashem², Elsayed Fathi Abd_Allah³, Ahmed Al-Harrasi¹

¹Natural and Medical Science Research Centre, University of Nizwa, 616, Oman

²Botany and Microbiology Department, King Saud University, Riyadh, 11451, Saudi Arabia

³Plant Production Department, King Saud University, Riyadh, 11451, Saudi Arabia

Corresponding author: Tapan Kumar Mohanta, E-mail: nostoc.tapan@gmail.com

Abstract

A proteomic analysis of proteomes from 145 plant species revealed a *pI* range of 1.99 (epsin) to 13.96 (hypothetical protein). The molecular mass of the plant proteins ranged from 0.54 to 2236.8 kDa. A putative Type-I polyketide synthase (22244 amino acids) in *Volvox carteri* was found to be the largest protein in the plant kingdom and was not found in higher plant species. Titin (806.46 kDa) and misin/midasin (730.02 kDa) were the largest proteins identified in higher plant species. The *pI* and molecular weight of the plant proteome exhibited a trimodal distribution. An acidic *pI* (56.44% of proteins) was found to be predominant over a basic *pI* (43.34% of proteins) and the abundance of acidic *pI* proteins was higher in unicellular algae species relative to multicellular higher plants. In contrast, the seaweed, *Porphyra umbilicalis*, possesses a higher proportion of basic *pI* proteins (70.09%). Plant proteomes were also found to contain the amino acid, selenocysteine (Sec), which is the first report of the presence of this amino acid in plants. Additionally, plant proteomes also possess ambiguous amino acids Xaa (unknown), Asx (asparagine or aspartic acid), Glx (glutamine or glutamic acid), and Xle (leucine or isoleucine) as well.

Key words: Proteome, amino acids, Isoelectric point, Molecular weight, Selenocysteine, Pyrrolysine

Introduction

The isoelectric or isoionic point of a protein is the pH at which a protein carries no net electrical charge and hence is considered neutral¹⁻⁴. The zwitterion form of a protein becomes dominant at neutral pH. The *pI* of polypeptides is largely dependent on the dissociation constant of the ionisable groups⁵. The major ionisable groups present in the amino acids are arginine, aspartate, cysteine, histidine, glutamate, lysine, and glutamine, where they play a major role in determining the *pI* of a protein⁶⁻⁸. Co-translational and post-translational modifications of a protein, however, can also play a significant role in determining the *pI* of a protein^{9,10}. The exposure of charged residues to the solvents, hydrogen bonds (dipole interactions) and dehydration also impact the *pI* of a protein^{11,12}. The inherent *pI* of protein, however, is primarily based on its native protein sequence. The *pI* of a protein is crucial to understanding its biochemical function and thus determining *pI* is an essential aspect of proteomic studies. During electrophoresis, the direction of movement of a protein in a gel or other matrix depends on its *pI*, hence numerous proteins can be separated based on their *pI*¹³⁻¹⁶. Given the impact of post-translational modifications and other biochemical alterations (phosphorylation, methylation, alkylation), however, the predicted *pI* of a protein will certainly be different than the predicted *pI*; the latter of which is based on the composition of amino acids in a protein^{9,17,18}. Nonetheless, an estimated isoelectric point is highly important and a commonly identified parameter.

Several studies have been conducted to understand the *pI* of proteins/polypeptides^{3,19-21}. These studies have been mainly based on animal, bacteria, and virus models and databases containing the *pI* of experimentally verified proteins. None of these databases, however, contain more than ten thousand proteins sequences which is very few relative to the availability of proteomic data. Therefore, an analysis was conducted of the *pI* and molecular weight of proteins from 144 plant

species which included 5.87 million protein sequences. This analysis provides an in-depth analysis of the *pI* and molecular mass of the proteins in the plant kingdom.

Results and Discussion

Plant proteins range from 0.54 kDa to 2236.8 kDa

A proteome-based analysis of plant proteins of 144 plant species that included more than 5.86 million protein sequences was conducted to determine the molecular mass, *pI*, and amino acid composition of proteins that exist in plant proteomes (Table 1). The analysis indicated that *Hordeum vulgare* possessed the highest number (248180) of protein sequences, while *Helicosporidium* sp. had the lowest number (6033). On average, plant proteomes possess 40988.66 protein sequences per species. The analysis also revealed that the molecular mass of plant proteomes ranged from 0.54 kDa to 2236.8 kDa. *Volvox carteri* was found to possess the largest plant protein (XP_002951836.1) of 2236.8 kDa, containing 22244 amino acids (*pI* 5.94), while *Citrus unshiu* possessed the smallest protein of 0.54 kDa, containing only four amino acids (*pI* 5.98) (id: GAY42954.1). This is the first analysis to document the largest (2236.8 kDa) and smallest (0.54 kDa) protein in the plant kingdom. These two proteins have not been functionally annotated and BLASTP analysis in the NCBI database did not identify suitable similarity with any other proteins. A few domains present in the largest protein, however, were found to be conserved with Type-I polyketide synthase. The molecular mass of some other high molecular mass proteins were: 2056.44 kDa (id: XP_001698501.1, type-1 polyketide synthase, *pI*: 6.00, aa: 21004, *Chlamydomonas reinhardtii*); 1994.71 kDa (id: XP_001416378.1, polyketide synthase, *pI*: 7.38, aa: 18193, *Ostreococcus lucimarinus*); 1932.21 kDa (id: Cz02g22160.t1, unknown protein, *pI*: 5.7, aa: 18533, *Chromochloris zofingiensis*); 1814.1 kDa (id: XP_007509537.1, unknown protein, *pI*: 4.46, aa: 16310, *Bathycoccus prasinos*); 1649.26 kDa (id: XP_011401890.1, polyketide synthase, *pI*: 5.53, aa: 16440, *Auxenochlorella protothecoides*); 1632.35 kDa (id: XP_005650993.1, ketoacyl-synt-

domain-containing protein, *pI*: 5.86, aa: 15797, *Coccomyxa subellipsoidea*); 1532.91 kDa (id: XP_002507643.1, polyketide synthase, *pI*: 7.07, aa: 14149, *Micromonas commoda*); 1370.23 kDa (id: GAX78753.1, hypothetical protein CEUSTIGMA, *pI*: 5.97, aa: 13200, *Chlamydomonas eustigma*); 1300.83 kDa (id: XP_022026115.1, unknown protein/filaggrin-like, *pI*: 11.75, aa: 12581, *Helianthus annuus*); 1269.42 kDa (id: XP_009350379.1, unknown protein, *pI*: 5.37, aa: 11880, *Pyrus bretschneideri*); 1237.34 kDa (id: XP_022840687.1, polyketide synthase, *pI*: 7.30, aa: 11265, *Ostreococcus tauri*); 1159.35 kDa (id: XP_005847912.1, polyketide synthase, *pI*: 5.91, aa: 11464, *Chlorella variabilis*); 1150.02 kDa (id: PKI66547.1, unknown protein, *pI*: 3.87, aa: 11234, *Punica granatum*); 1027.64 kDa (id: Sphfalx0133s0012.1, unknown protein, *pI*: 4.05, aa: 9126, *Sphagnum fallax*); 909.93 kDa (id: XP_002985373.1, unknown/titin-like protein, *pI*: 4.02, aa: 8462, *Selaginella moellendorffii*); 881.59 kDa (id: KXZ46216.1, hypothetical protein, *pI*: 5.80, aa: 8881, *Gonium pectorale*); 848.29 kDa (id: XP_003056330.1, *pI*: 6.12, aa: 7926, *Micromonas pusilla*); 813.31 kDa (id: GAQ82263.1, unknown protein, *pI*: 4.60, aa: 7617, *Klebsormidium nitens*), 806.46 kDa (id: XP_017639830.1, titin-like, *pI*: 4.21, aa: 7209, *Gossypium arboreum*); 806.12 kDa (id: OAE35580.1, *pI*: 4.83, hypothetical protein, aa: 7651, *Marchantia polymorpha*); and 802.74 kDa, (id: XP_012444755.1, titin-like, *pI*: 4.19, aa: 7181, *Gossypium raimondii*) (Table 1).

On average, approximately 7.38 % of the analysed proteins were found to contain ≥ 100 kDa proteins. The algal species, *V. carteri*, was found to encode largest plant protein (putative polyketide synthase); while other unicellular algae, and multi-cellular lower eukaryotic plants, including bryophytes and pteridophytes, were also found to encode some of the larger proteins (e.g. ketoacyl synthase) in the plant kingdom. The higher eukaryotic plants, including gymnosperms and angiosperms, were not found to encode a high molecular mass polyketide synthase protein. They did, however, possess the high molecular mass proteins; titin (806.46 kDa), misin/midasin (730.02 kDa), futsch (622.14 kDa), filaggrin (644.4 kDa), auxin transport

protein BIG (568.4 kDa), and von Willebrand factor (624.74 kDa) (Table 1). Titin is an extremely large protein that is greater than 1 μ M in length and found in human striated muscle^{22,23}. The largest titin protein found in plants, however, was only 806.46 kDa (*Gossypium arboreum*). The predicted formula of the 806.46 kDa titin protein was C₃₃₈₆₃H₅₄₆₁₀N₉₂₃₂O₁₃₀₆₁S₂₀₀ and its estimated half-life was 10-30 hours; whereas the predicted formula of the 2236.8 kDa protein of *V. carteri* was C₉₇₇₈₃H₁₅₇₄₀₁N₂₈₄₈₉O₃₀₂₆₅S₆₃₇. Almost all of the higher eukaryotic plants were found to possess titin, misin/midasin, and auxin transport protein BIG proteins. Species of unicellular algae were not found to possess titin or misin/midasin proteins. This suggests that titin and misin/midasin proteins originated and evolved in more complex, multicellular organisms rather than unicellular organisms. Thus, the evolution of titin, misin/midasin proteins may also be associated with the evolution of terrestrial plants from aquatic plants.

The presence of the smallest molecular mass protein, other than the tripeptide glutathione (Cys, Gly, and Glu), was also determined. A 0.54 kDa molecular mass protein, containing only four amino acids (MIMF) and starting with methionine and ending with phenylalanine, was identified in *Citrus unshiu* (id: GAY42954.1) (Table 1). Other low molecular mass plant proteins were 0.57 kDa (NP_001336532.1/ AT5G23115, *Arabidopsis thaliana*) and 0.63 kDa (AH003201-RA, *Amaranthus hypochondriacus*). Small proteins found in *A. thaliana* was MNPKS and that found in *A. hypochondriacus* was MLPYN, contained only five amino acids. These low molecular mass proteins were not present in all of the studied species and their cellular and molecular functions have not been reported yet. One of the universal small molecular weight plant proteins, however, was identified as cytochrome b6/f complex subunit VIII (chloroplast) (MDIVSLAWAALMVVFTFSLSLVWGRSGL) that contains only 29 amino acids. Cytochrome b6/f is actively involved in the electron transfer system of

photosystem II and regulates photosynthesis^{24–28}. It is commonly known that glutathione is the smallest functional polypeptide and that it plays diverse roles in cell signaling^{29–31}. The tetra and penta peptides identified in the present analysis, however, were quite different from glutathione and none of them contained Cys, Gly, or Glu amino acids, as found in glutathione. Polypeptides with less than 100 amino acids are considered small proteins and studies indicate that many small proteins are involved in cell metabolism, cell signaling, cell growth, and DNA damage^{32–35}. In the era of next-generation sequencing, small protein-coding genes are completely overlooked during genome annotation and get buried amongst an enormous number of open reading frames³⁶. Therefore, it is difficult to identify more numbers of small proteins in plants.

A previously conducted comparative study revealed that plant proteins are comparatively smaller than animal proteins, as the former are encoded by fewer exons³⁷. Longer proteins harbour more conserved domains and hence display a greater number of biological functions than short proteins. The average protein length of the studied plant species was 424.34 amino acids. A previous study reported the average length of eukaryotic proteins to be 472 amino acids and that the average length of plant proteins is approximately 81 amino acids shorter than animal proteins³⁷. Our analysis indicates, however, that plant proteins are approximately 47.66 amino acid shorter than animal proteins. In addition, studies have also indicated that eukaryotic proteins are longer than bacterial proteins and that eukaryote genomes contain approximately 7 fold more proteins (48% larger) than bacterial genomes³⁸. Although the average size of plant proteins was found to be 424.34 amino acids, the average protein size of lower, eukaryotic unicellular aquatic plant species; including *Chlamydomonas eustigma*, *Volvox carteri*, *Klebsordium nitens*, *Bathycoccus prasinos*, and *Durio zibethinus*, was found to be 576.56, 568.22, 538.73, 521.05, and 504.36 amino acids, respectively. This indicates that

unicellular plant species have an average protein size that is larger than terrestrial multicellular complex plant species, suggesting that the evolution of plant proteins involved a loss of protein size and hence gene size. The cause of the variability in protein length in the phylogenetic lineage of eukaryotic plants has yet to be elucidated. A multitude of evolutionary factors, including deletion (loss of exons) or fusion of multiple domains of proteins, may have played critical roles in shaping the size of higher plant proteins. Transposon insertion and splitting of genes increases the number of proteins but reduces the average size of the proteins^{39–42}. Higher plants contain a very large number of transposable elements and therefore these elements are the most responsible factor to expect to have played a major role in increasing protein numbers and reducing the protein size in higher plants. The percentage of transposable elements in a genome is directly proportional to the genome size of the organism and varies from approximately 3% in small genomes to approximately 85% in large genomes⁴¹. Kirag et al (2007) reported a significant correlation between protein length and the *pI* of a protein¹⁹. In our analysis, however, no correlation was found between protein length and the *pI* of a protein. For example, titin and misin are two of the larger proteins in plants and they fall in the acidic *pI* range, but not the alkaline *pI* range.

Plant encode a higher number of proteins than animals and fungi

Our analysis identified an average of 40469.47 proteins per genome (Table 1). Previously the number of proteins in plant species was reported as 36795 per genome³⁷. On average, animals and fungi encode 25189 and 9113 proteins per genome, respectively³⁷. An average of 40469.47 proteins per plant genome is 62.24% higher than in animals and 444.08% higher than in fungi. Although, plant species encode a higher number of proteins, their size is smaller than the average size of animal proteins. Notably, green algae contains a smaller number of proteins than higher plants but their average protein size is 1.27 times larger. The average protein size (low to high) in the species of green algae ranged from 273.08 (*Helicosporidium* sp.) to 576.56

(*Chlamydomonas eustigma*) amino acids, dicots ranged from 253.34 (*Trifolium pratense*) to 498.49 (*Vitis vinifera*), and monocots ranged from 111.54 (*Hordeum vulgare*) to 473.35 (*Brachypodium distachyon*) amino acids. The average protein size of monocot proteins (431.07 amino acids), however, is slightly larger than dicots (424.3 amino acids). In addition to transposons, previous studies have reported that endosymbiosis may have also played an important role in the reduction of protein size in plant genomes^{37,43,44}. This would have been due to the post endosymbiosis acquisition of thousands of genes from the chloroplast, since cyanobacterial proteins are smaller than eukaryotic proteins and cyanobacteria are the ancestors of plastids^{37,45}. In this hypothesis, the intermediate size of plant proteins would be the result of the migration of proteins from cyanobacteria (chloroplast) to the plant nucleus, thereby reducing the overall average size of the protein by a dilution effect^{46,47}.

The pI of plant proteins ranges from 1.99 to 13.96

Results indicated that the *pI* of analysed plant proteins ranged from 1.99 (id: PHT45033.1, *Capsicum baccatum*) to 13.96 (id: PKI59361.1, *Punica granatum*). The protein with the lowest *pI* (1.99) was epsin and the protein with the highest *pI* (13.96) was a hypothetical protein. This is the first study to report on the plant proteins with the lowest and highest *pI*. The *C. baccatum* protein with *pI* 1.99 contains 271 amino acids, whereas the *P. granatum* protein with *pI* 13.96 contains 986 amino acids. The epsin protein (*pI* 1.99) is composed of 16 amino acid repeats (GWIDGWIDGWIDGW), while the hypothetical protein (*pI* 13.06) is composed of 64 QKLKSGLT and 31 TRRGLTAV repeats. From among the 20 essential amino acids, the epsin protein only contained five amino acids, namely Asp (68), Gly (68), Ile (65), Met (3), and Trp (67). The amino acids were arranged in a repeating manner within the full-length epsin protein. This study is the first to report a full-length protein composed of such a minimum number of essential amino acids. Similarly, the hypothetical protein with the highest *pI* (13.96) was composed of only nine amino acids, namely Ala (62), Gly (132), Lys (127), Leu (197), Met

(M), Pro (4), Gln (64), Arg (132), and Ser (66). Intriguingly, cysteine, which is one of the most important amino acids as it is responsible for the formation of disulphide bonds, was not found in either the smallest or largest protein. Disulphide bonds maintain the conformation and stability of a protein and are typically found in extracellular proteins and only rarely in intracellular proteins⁴⁸. The absence of Cys amino acids in these proteins suggests that these proteins are localized to the intracellular compartments within the cell.

The plant proteome is primarily composed of acidic *pI* proteins rather than basic *pI* proteins (Table 1). Approximately, 56.44% of the analysed proteins had a *pI* within the acidic *pI* range. The average percentage of acidic *pI* proteins was comparatively higher in the lower eukaryotic plants, algae, and bryophytes, than in the higher land plants. A total of 64.18% of proteins in *Chlamydomonas eustigma* were found in the acidic *pI* region, followed by *Ostreococcus lucimarinus* (64.17%), *Micromonas commoda* (63.30%), *Helicosporium* sp. (62.97%), *Gonium pectoral* (62.76%), *Chromochloris zofingiensis* (62.41%), *Coccomyxa subellipsoidea* (62.12%), and *Sphagnum fallax* (61.83%). The algal species, *Porphyra umbilicalis*, had the lowest percentage (29.80%) of acidic *pI* proteins. The dicot plant, *Punica granatum*, and the algal species, *Botryococcus braunii*, had a significantly lower percentage of acidic *pI* proteins (45.72% and 47.18%, respectively) relative to other plants. Principal component analysis (PCA) of acidic *pI* protein content revealed that the acidic proteins of bryophytes and monocots cluster closely to each other compared to algae and eudicot plants (Figure 1). Similarly, in the case of basic *pI* proteins, a great variation was observed for algae, eudicot and monocot plants (Figure 2). The basic *pI* proteins of bryophytes, however, were found to be consistent. A previous study reported that protein *pI* values are correlated with the sub-cellular localization of the proteins, and that the *pI* of cytosolic proteins fall below 7²¹. Among cytosolic proteins are those involved in 26S proteasome degradation, oxidative pentose phosphate pathway, actin/tubulin, mevalonate pathway, sugar and nucleotide biosynthesis, glycolysis, RNA

processing, and several other cellular process. Our analysis indicated that the *pI* of all cytosolic proteins does not fall in the acidic *pI* range. Ribosomal proteins, pre-mRNA splicing factors, transcription factors, auxin induced protein, extensin, senescence associated protein, cyclin dependent protein kinase and other cytoplasmic proteins had a *pI* greater than 7.

In contrast to acidic *pI* proteins, plants possess a comparatively low number of basic *pI* proteins. On average, 43.34% of the analysed plant proteins possessed a *pI* in the basic range. The highest percentage of basic *pI* proteins was found in *Porphyra umbilicalis*, where 70.09% of the proteins had a basic *pI* (Table 1). *Punica granatum* also had a high percentage (54.11%) of basic *pI* proteins (Table 1). The lowest percentage of basic *pI* proteins was found in the algal species, *Chlamydomonas eustigma* (35.56%), followed by *Ostreococcus lucimarinus* (35.65%), *Micromonas commoda* (36.52%), *Helicosporidium* sp. (36.89%), and *Gonium pectorale* (37.04%). It is difficult to establish the reason that algal species contain more acidic *pI* and less basic *pI* proteins. *Porphyra umbilicalis* is a cold-water seaweed within the family, Bangiophyceae, and it is the most domesticated marine algae. The 87.7 Mbp haploid genome of *P. umbilicalis* has a 65.8% GC content and an evolutionary study reported that the genome of *Porphyra umbilicalis* had undergone a reduction in size⁴⁹. Since this species is found in the intertidal region of the ocean, it has developed the ability to cope with mid-to-high levels of tidal stress. *Porphyra* is also tolerant to UV-A and UV-B radiation⁴⁹⁻⁵¹. The high GC content in *Porphyra umbilicalis* is directly proportional to the high percentage of basic proteins. The GC content of algal species is higher relative to other plant species and algal species possess a lower percentage of basic *pI* proteins. This suggests that, in algae, percentage GC content is inversely proportional to percentage of proteins with a basic *pI*. However, this is not true in the case of higher plants.

The pI of plant proteomes exhibits a trimodal distribution

The *pI* of the analysed plant proteins ranged from 1.99 to 13.96 and exhibited a trimodal distribution (Figure 3). Schwartz et al., previously reported a trimodal distribution of the *pI* of eukaryotic proteins²¹, however, they did not provide information on the number of sequences/species considered in their study. Proteins are typically soluble near their isoelectric point and the cytoplasm possesses a pH that is close to neutral. This may be the reason for the trimodal distribution of *pI*. Although the *pI* values of proteins estimated *in silico* or experimentally might be different *in vivo*, they are typically in close agreement⁵². Kiraga et al., (2006) reported a bimodal distribution of the *pI* of proteins from all organisms, citing acidic and basic *pI* as the basis of the modality¹⁹, where modality is defined as the set of data values that appears most often. They reported that taxonomy, ecological niche, proteome size, and sub-cellular localization are correlated with acidic and basic proteins. However, no correlation was observed in the current study between either acidic or basic *pI* of proteins with regard to taxonomy, ecological niche, or proteome size. For example, *Hordeum vulgare* and *Brassica napus* possess the largest proteomes among the studied plant species, possessing 248180 and 123465 proteins, respectively. In *H. vulgare*, 53.28% of the proteins fall in the acidic and 46.50% fall in the basic *pI* ranges; while in *B. napus*, 55.28% of the proteins have an acidic *pI* and 44.48% have a basic *pI*. Other species with smaller proteomes, however, possess a higher percentage of acidic or basic proteins (Table 1). Therefore, no correlation exists between the percentage of either acidic or basic proteins and proteome size, taxonomy, or the ecological niche of an organism. Knight et al. also reported a negative correlation between the *pI* of a protein with phylogeny of the organism⁵³. The existence of a trimodal distribution of the *pI* of the plant proteome can be considered as a virtual 3D-gel of a plant's proteins where the *pI* of

the protein is plotted against the molecular weight of the protein. On average, 0.21% of the analysed proteins were found to have a neutral *pI* (*pI* 7), while only 0.09% of the proteins in *O. lucimarinus* fall in neutral *pI*.

Leu is a high- and Trp is a low-abundant amino acid in the plant proteome

The plant-kingdom-wide proteome analysis revealed that Leu was the most (9.62%) while Trp was the least (1.28%) abundant amino acid (Figure 4, Supplementary File 1). Leu is a nonpolar amino acid, whereas Trp contains an aromatic ring. The distribution of amino acids indicates that the synthesis of nonpolar amino acids is more favoured in the plant proteomes than the polar amino acids or those containing an aromatic ring. The average abundance of other nonpolar amino acids Ala, Gly, Ile, Met, and Val was 6.68%, 6.80%, 4.94%, 2.40%, and 6.55%, respectively (Table 2, Supplementary File 1). Trp and Tyr amino acid contain an aromatic ring and the abundance of these two proteins is relatively low in the plant proteome compared to other amino acids. Results of the conducted analysis indicated that the abundance of Ala (17.58%), Gly (11.76%), Pro (9.2%), and Arg (9.81%) were the highest; whereas, Tyr (1.33%), Gln (2.04%), Asn (1.53%), Met (1.45%), Lys (7.07%), Lys (2.08%), Ile (1.77%), Phe (2.01%), and Glu (3.52%) were the lowest in *Porphyra umbilicalis*. In a few algae and seaweeds Ala, Asp, Glu, Gly, Pro, Gln, Arg, Thr, and Val were found in high percentage while Asp, Glu, Phe, His, Ile, Lys, Leu, Met, Asn, Gln, and Ser were found in low percentage (Table 2). These observations indicate that the composition of amino acids in unicellular algae, seaweeds, and gymnosperms are more dynamic and variable than in angiosperms and other terrestrial land plants. Principal component analysis revealed that the low abundant amino acids, Trp, Tyr, His, Met, Cys, and Xaa (unknown), cluster in one group while the high abundant amino acids, Leu, Glu, Ile, Lys, and Ser, cluster in another group (Figure 5). None of the terrestrial land plants were located in the high- and low-abundant amino acid clusters. This suggests that the

proteome and amino acid composition of the land plants are more conserved and stable relative to the algae and seaweeds. PCA analysis further revealed that the *pI* of algae, eudicots, and monocots are lineage specific. The *pI* of algae, monocots, and eudicots were strongly correlated and clustered together (Figure 5). The question arises, however, as to why the plant proteome contains the highest percentage of Leu and of the lowest percentage of Trp amino acids. Do the energy requirements of the different biosynthetic pathways play a pivotal role in deciding the abundance of amino acids in a proteome? To address this question, an attempt was made to understand the role of amino acid biosynthetic pathways in determining the abundance of specific amino acids in the proteome.

Various amino acids are produced in different biosynthetic pathways^{54–58} (Figure 6). In some cases, a few amino acids act as the substrate for the biosynthesis of other amino acids; whereas in other cases, allosteric inhibition of the biosynthesis of amino acids occurs^{59–61}. In all of these biosynthetic pathways, ATP or NADH/NADPH are used as a source of energy, along with substrate that play a vital role in the biosynthesis of amino acids. Overall, the biosynthesis of 20 essential amino acid families are grouped by metabolic precursors⁶² (Table 3); namely α -ketoglutarate (Arg, Gln, Glu, Pro), pyruvate (Ala, Ile, Leu, Val), 3-phosphoglycerate (Cys, Gly, Ser), phosphoenolpyruvate and erythrose 4-phosphate (Phe, Trp, Tyr), oxaloacetate (Asn, Asp, Lys, Met, Thr), and ribose 5-phosphate (His) (Table 3)⁶². Ala, Ile, Leu, and Val are synthesized from pyruvate; Arg, Glu, Gln, and Pro are synthesized from α -ketoglutarate and Gly and Ser are synthesized from 3-phosphoglycerate⁶²; all of which have a higher abundance in the plant proteome relative to the other amino acids (Figure 6, Table 3). 3-phosphoglycerate and pyruvate are intermediates of glycolysis and the amino acids synthesized from these intermediates maintain a high abundance in the plant proteome. The intermediate, 3-phosphoglycerate, is formed in an early step of glycolysis⁶². The amino acids Gly and Ser are

synthesized from 3-phosphoglycerate and are also found abundantly in the plant proteome (Figure 6, Table 3). The amino acid Cys, which is also synthesized from 3-phosphoglycerate⁶², however, is present in low abundance (1.85%) in the plant proteome. The low abundance of Cys may be due to the allosteric inhibition. Phe (3.97), Trp (1.28%), and Try (2.67%) contain an aromatic ring and are synthesized via phosphoenolpyruvate and erythrose 4-phosphate. The aromatic amino acids are in low abundance in the plant proteome (Table). Since Phe also plays a role in the biosynthesis of Tyr, the abundance of Phe is relatively higher than Trp and Tyr. Glucose 6-phosphate gives rise to ribose 5-phosphate in a complex reaction of four steps⁶² and His gets subsequently synthesized from ribose 5-phosphate. It is possible that the complexity of the biosynthetic pathways of amino acids containing ring compounds might be the reason for their low abundance in the plant proteome.

Plants possess selenocysteine (Sec) and other novel amino acids

A few of the plant proteomes that were analysed had proteins containing the amino acid, selenocysteine (Sec). *C. reinhardtii*, *M. pusilla*, and *V. carteri* contained 9, 16, and 11 Sec amino acids in their proteome, respectively. Selenium containing selenoproteins are commonly found in animals but have been reported to be present in plant species. Novoselov et al., (2002) reported the presence of a selenoprotein in *C. reinhardtii*⁶³. In our analysis, nine selenoproteins (selenoprotein H, selenoprotein K1, selenoprotein M2, selenoprotein T, selenoprotein U, selenoprotein W1, selenoprotein W2, NADPH-dependent thioredoxin reductase 1, and glutathione peroxidase) were identified in *C. reinhardtii*. In addition, *M. pusilla* was found to possess 14 Sec-containing proteins [DSBA oxidoreductase (2 no.), selenoprotein T, glutathione peroxidase (4 no.), selenoprotein W, selenoprotein, selenoprotein M, selenoprotein H, selenoprotein O, methyltransferase, and peroxiredoxin). In addition, *V. carteri* was found to possess 10 Sec containing proteins (selenoprotein T, selenoprotein K1, selenoprotein H,

selenoprotein W1, selenoprotein M2, selenoprotein U, glutathione peroxidase, membrane selenoprotein, NADPH-dependent thioredoxin reductase, and peptide methionine-S-sulfoxide reductase). To the best of our knowledge, this is the first report of selenoproteins in *M. pusilla* and *V. carteri* and the first to report of H, K1, T, U, M, M2, O, W1, and W2 selenoprotein families collectively in *C. reinhardtii*, *M. pusilla*, and *V. carteri*. The I, N, P, R, S, and V selenoprotein family members are commonly found in the animal kingdom⁶⁴ but are absent in *C. reinhardtii*, *M. pusilla*, and *V. carteri*. This is also the first report of the Sec-containing proteins, DSBA oxidoreductase, methyltransferase, peroxiredoxin, peptide methionine-S-sulfoxide reductase, and membrane selenoprotein in plants lineage (algae). Notably, the selenoproteins DSBA oxidoreductase, methyltransferase, peroxiredoxin, peptide methionine-S-sulfoxide reductase, and membrane selenoprotein have not been reported in animal species. Outside of algal species, no other plant species; including bryophytes, pteridophytes gymnosperms and angiosperms, were found to possess a selenoprotein.

Some plant proteomes were also found to possess a few unknown or unspecified amino acids, commonly designated as Xaa (X). Among the analysed plant species, *Aegilops tauschii*, *Amaranthus hypochondriacus*, and *Amborella trichocarpa* encoded 149377, 55412, and 25843 X amino acids, respectively. *Solanum lycopersicum* was found to contain only one X amino acid, while at least 18 species (*Solanum pennellii*, *Solanum tuberosum*, *Sorghum bicolor*, *Sphagnum fallax*, *Spinacia oleracea*, *Spirodela polyrhiza*, *Tarenaya hassleriana*, *Theobroma cacao*, *Trifolium pratense*, *Trifolium subterraneum*, *Triticum aestivum*, *Triticum urartu*, *Vigna angularis*, *Vigna radiata*, *Vigna anguiculata*, *Vitis vinifera*, *Volvox carteri*, and *Zostera marina*) were found to lack any Xaa amino acids in their proteome. Xaa amino acids are known as non-protein amino acids as they have not been associated with any specific codons. Among the studied plant species, ten were found to contain amino acid B (Asx) that codes for the

ambiguous amino acid Asn or Asp that is translated as Asp. Species that were found to possess an Asx amino acid included *Arachis duranensis* (1), *Brachypodium stacei* (40), *Dichanthelium oligosanthes* (20), *Dunaliella salina* (31), *Glycine max* (1), *Malus domestica* (4080), *Momordica charantia* (98), *Nelumbo nucifera* (64), *Prunus persica* (1), and *Trifolium pratense* (76). At least six species were found to possess a J (Xle) amino acid. Xle amino acid can encode either Leu or Ile but during translation produces Leu. Species that were found to possess Xle amino acids included *Arabidopsis thaliana* (10), *Dichanthelium oligosanthes* (11), *Malus domestica* (2175), *Momordica charantia* (39), *Nelumbo nucifera* (29), and *Trifolium pratense* (39). At least seven species were found to possess a Z (Glx) amino acid that codes for either Glu or Gln, which is subsequently translated as Glu. Species that were found to encode a Glx amino acid included *Brachypodium stacei* (20), *Dichanthelium oligosanthes* (16), *Dunaliella salina* (7), *Malus domestica* (1552), *Momordica charantia* (28), *Nelumbo nucifera* (14), and *Trifolium pratense* (25). Among the studied species, *Malus domestica* was found to contain highest number of ambiguous amino acids (Asx, Xle, and Glx). Bodley and Davie (1966) reported the incorporation of ambiguous amino acids in a peptide chain ⁶⁵. The presence of ethanol or streptomycin or a high magnesium ion concentration induces ambiguous coding in the peptide chain ⁶⁵. They reported that poly-U (uridylic acid) in the presence of a high concentration of magnesium ions or ethanol or streptomycin induces the incorporation of Leu/Ile amino acids in a peptide chain ⁶⁵. This explains how the specificity of the protein translation process can be altered by the presence of environmental factors. A high concentration of magnesium ions, organic solvents, antibiotics, pH, and low temperature have the ability to modify the coding specificity of a peptide chain ⁶⁵. Under some conditions poly-U triggers the incorporation of Leu and Ile or Phe ⁶⁵. *Malus domestica* is rich in magnesium ions (1%) and this might explain the presence of such a high number of ambiguous amino acids in its proteome.

Conclusion

A proteomic analysis of the plant kingdom identified proteins with a great range of molecular mass and isoelectric points. Isoelectric points ranged from 1.99 to 13.96, covering almost the entire pH range. It is quite intriguing to think about the functions of protein at *pI* 1.99 or 13.96. Proteins with an acidic *pI* predominate over the proteins with an alkaline *pI*, and the presence of proteins with a *pI* that is near neutral is very negligible. The percentage of proteins with acidic or basic *pIs* is not related to the host cell, alkalinity or acidity of the environment. Additionally, the GC content of a genome or size distribution of the overall proteome are not directly proportional to the distribution (percent basic vs. percent acidic) of the *pI* in the proteome. The presence of Sec-containing proteins in some plant species needs to be further investigated to determine their functional role. Similarly, the presence of ambiguous amino acids in plant species should be further evaluated individually to understand whether these ambiguous amino acids are encoded by any specific codon in mRNA. A major question arises from this study is, whether the incorporation of ambiguous amino acids in the peptide chain of the protein brings any impact at the gene/genome level and incorporate the respective codon for the ambiguous amino acid and regulated genomic rearrangement through reverse central dogma approach. The presence of a pyrrolysine amino acid in the plant kingdom was not observed in the present study.

Methods

Protein sequences of the entire proteome of the analysed plant species were downloaded from the National Center for Biotechnology Information (NCBI) and Phytozome, DOE Joint Genome Institute (<https://phytozome.jgi.doe.gov/pz/portal.html>). All of the studied sequences were annotated nuclear-encoded proteins. The isoelectric point of each protein of each of the analysed plant species was calculated individually using the Python-based command line

“Protein isoelectric point calculator” (IPC Python) in a Linux platform². The source code used was as written by Kozlowski (2016).

Once the molecular mass and isoelectric point of the proteins in each species was determined, they were separated into acidic and basic *pI* categories. Subsequently, the average *pI* and percentage of proteins in each category was calculated using a Microsoft excel worksheet. A graph comparison of isoelectric point versus molecular mass was prepared using a python-based platform. Pearson-correlation ($r=0.19$, $p = 0$) was used for the association analysis of molecular mass and isoelectric point. The X-axis data statistics were as follows: mean, 4.717365e+01; std, 3.662983e+01; min, 8.909000e-02; 25%, 2.279452e+01; 50%, 3.874486e+01; 75%, 5.999628e+01, and max, 2.236803e+03. The Y-axis data statistics were: mean, 6.840657e+00; std, 1.594912e+00; min, 1.990000e+00; 25%, 5.537000e+00; 50%, 6.605000e+00; 75%, 8.053000e+00, and max, 1.396300e+01.

Principal component analysis

Principal component analysis of the plant proteome parameters was carried out using a portable Unscrambler software version 9.7 using the excel file format. For acidic and basic *pI*, the plant proteome data were grouped according to the plant lineage algae, bryophyte, monocot, and eudicot plants. The average of acidic and basic *pI* was used to construct the PCA plot. Similarly, amino acid abundance was also analysed in relation to algae, bryophyte, eudicot, and monocot lineage.

References

1. Kirkwood, J., Hargreaves, D., O'Keefe, S. & Wilson, J. Using isoelectric point to determine the pH for initial protein crystallization trials. *Bioinformatics* **31**, 1444–1451 (2015).
2. Kozlowski, L. P. IPC – Isoelectric Point Calculator. *Biol. Direct* **11**, 55 (2016).
3. Kozlowski, L. P. Proteome-pI: proteome isoelectric point database. *Nucleic Acids Res.* **45**, D1112–D1116 (2017).
4. Stekhoven, F., Gorissen, M. & Flik, G. The isoelectric point, a key to understanding a variety of biochemical problems: a minireview. *Fish Physiol. Biochem.* **34**, 1–8 (2008).
5. Kenneth, W., Kenneth, G. & Raymond, E. *General Chemistry*. (Saunders College Publishing, 1992).
6. Pace, C. N., Grimsley, G. R. & Scholtz, J. M. Protein ionizable groups: pK values and their contribution to protein stability and solubility. *J. Biol. Chem.* **284**, 13285–13289 (2009).
7. Grimsley, G. R., Scholtz, J. M. & Pace, C. N. A summary of the measured pK values of the ionizable groups in folded proteins. *Protein Sci.* **18**, 247–251 (2009).
8. Masunov, A. & Lazaridis, T. Potentials of Mean Force between Ionizable Amino Acid Side Chains in Water. *J. Am. Chem. Soc.* **125**, 1722–1730 (2003).
9. Zhu, K., Zhao, J., Lubman, D. M., Miller, F. R. & Barder, T. J. Protein pI Shifts due to Posttranslational Modifications in the Separation and Characterization of Proteins. *Anal. Chem.* **77**, 2745–2755 (2005).
10. Locke, D., Koreen, I. V & Harris, A. L. Isoelectric points and post-translational modifications of connexin26 and connexin32. *FASEB J.* **20**, 1221–1223 (2006).
11. Youden, W. J. & Denny, F. E. Factors Influencing the pH Equilibrium Known as the Isoelectric Point of Plant Tissue. *Am. J. Bot.* **13**, 743–753 (1926).
12. Shaw, K. L., Grimsley, G. R., Yakovlev, G. I., Makarov, A. A. & Pace, C. N. The effect of net charge on the solubility, activity, and stability of ribonuclease Sa. *Protein Sci.* **10**, 1206–1215 (2001).
13. Saraswathy, N. & Ramalingam, P. Introduction to proteomics. in *Woodhead Publishing Series in Biomedicine* (eds. Saraswathy, N. & Ramalingam, P. B. T.-C. and T. in G. and P.) 147–158 (Woodhead Publishing, 2011).
doi:<https://doi.org/10.1533/9781908818058.147>
14. Gygi, S. P., Corthals, G. L., Zhang, Y., Rochon, Y. & Aebersold, R. Evaluation of two-dimensional gel electrophoresis-based proteome analysis technology. *Proc. Natl. Acad. Sci.* **97**, 9390 LP-9395 (2000).
15. Rabilloud, T. & Lelong, C. Two-dimensional gel electrophoresis in proteomics: A tutorial. *J. Proteomics* **74**, 1829–1841 (2011).
16. Kumar, M. *et al.* An Improved 2-Dimensional Gel Electrophoresis Method for Resolving Human Erythrocyte Membrane Proteins. *Proteomics Insights* **8**, 1178641817700880 (2017).
17. Anderson, J. C. & Peck, S. C. A simple and rapid technique for detecting protein phosphorylation using one-dimensional isoelectric focusing gels and immunoblot analysis. *Plant J.* **55**, 881–885 (2008).
18. Jones, L. R., Simmerman, H. K., Wilson, W. W., Gurd, F. R. & Wegener, A. D. Purification and characterization of phospholamban from canine cardiac sarcoplasmic reticulum. *J. Biol. Chem.* **260**, 7721–7730 (1985).
19. Kiraga, J. *et al.* The relationships between the isoelectric point and: length of proteins, taxonomy and ecology of organisms. *BMC Genomics* **8**, 163 (2007).
20. Carugo, O. Isoelectric points of multi-domain proteins. *Bioinformation* **2**, 101–104 (2007).

21. Schwartz, R., Ting, C. S. & King, J. Whole proteome pI values correlate with subcellular localizations of proteins for organisms within the three domains of life. *Genome Res.* **11**, 703–709 (2001).
22. Labeit, S. *et al.* A regular pattern of two types of 100-residue motif in the sequence of titin. *Nature* **345**, 273 (1990).
23. Kurzban, G. P. & Wang, K. Giant polypeptides of skeletal muscle titin: Sedimentation equilibrium in guanidine hydrochloride. *Biochem. Biophys. Res. Commun.* **150**, 1155–1161 (1988).
24. Kurisu, G., Zhang, H., Smith, J. L. & Cramer, W. A. Structure of the Cytochrome b6f Complex of Oxygenic Photosynthesis: Tuning the Cavity. *Science* (80-.). **302**, 1009 LP-1014 (2003).
25. Leister, D. & Schneider, A. B. T.-. From Genes to Photosynthesis in *Arabidopsis thaliana*. *Int. Rev. Cytol.* **228**, 31–83 (2003).
26. Barbagallo, R. P., Finazzi, G. & Forti, G. Effects of Inhibitors on the Activity of the Cytochrome b6f Complex: Evidence for the Existence of Two Binding Pockets in the Lumenal Site. *Biochemistry* **38**, 12814–12821 (1999).
27. Cramer, W. A. *et al.* Some New Structural Aspects And Old Controversies Concerning The Cytochrome B6f Complex Of Oxygenic Photosynthesis. *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **47**, 477–508 (1996).
28. Kallas, T. The Cytochrome b6f Complex. in *The Molecular Biology of Cyanobacteria* (ed. Bryant, D. A.) 259–317 (Springer Netherlands, 1994). doi:10.1007/978-94-011-0227-8_9
29. Filomeni, G., Rotilio, G. & Ciriolo, M. R. Cell signalling and the glutathione redox system. *Biochem. Pharmacol.* **64**, 1057–1064 (2002).
30. Zhang, H. & Forman, H. J. Glutathione synthesis and its role in redox signaling. *Semin. Cell Dev. Biol.* **23**, 722–728 (2012).
31. Aquilano, K., Baldelli, S. & Ciriolo, M. R. Glutathione: new roles in redox signaling for an old antioxidant. *Front. Pharmacol.* **5**, 196 (2014).
32. Su, M., Ling, Y., Yu, J., Wu, J. & Xiao, J. Small proteins: untapped area of potential biological importance. *Front. Genet.* **4**, 286 (2013).
33. Setlow, P. I will survive: DNA protection in bacterial spores. *Trends Microbiol.* **15**, 172–180 (2007).
34. Schalk, C. *et al.* Small RNA-mediated repair of UV-induced DNA lesions by the DNA DAMAGE-BINDING PROTEIN 2 and ARGONAUTE 1. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E2965–E2974 (2017).
35. Xue, Y. *et al.* Tff3, as a Novel Peptide, Regulates Hepatic Glucose Metabolism. *PLoS One* **8**, e75240 (2013).
36. Basrai, M. a, Hieter, P. & Boeke, J. D. Small Open Reading Frames : Beautiful Needles in the Haystack Small Open Reading Frames : Beautiful Needles in the Haystack. *Genome Res.* **7**, 768–771 (1997).
37. Ramírez-Sánchez, O., Pérez-Rodríguez, P., Delaye, L. & Tiessen, A. Plant Proteins Are Smaller Because They Are Encoded by Fewer Exons than Animal Proteins. *Genomics. Proteomics Bioinformatics* **14**, 357–370 (2016).
38. Tiessen, A., Pérez-Rodríguez, P. & Delaye-Arredondo, L. J. Mathematical modeling and comparison of protein size distribution in different plant, animal, fungal and microbial species reveals a negative correlation between protein size and protein number, thus providing insight into the evolution of proteomes. *BMC Res. Notes* **5**, 85 (2012).
39. Purugganan, M. & Wessler, S. The splicing of transposable elements and its role in intron evolution. *Genetica* **86**, 295–303 (1992).

40. Huff, J. T., Zilberman, D. & Roy, S. W. Mechanism for DNA transposons to generate introns on genomic scales. *Nature* **538**, 533–536 (2016).
41. Lee, S.-I. & Kim, N.-S. Transposable elements and genome size variations in plants. *Genomics Inform.* **12**, 87–97 (2014).
42. Plasterk, R. H. A. *Transposable Elements. Encyclopedia of Genetics* (Academic Press, 2001). doi:<https://doi.org/10.1006/rwgn.2001.1316>
43. Martin, W. *et al.* Evolutionary analysis of Arabidopsis, cyanobacterial, and chloroplast genomes reveals plastid phylogeny and thousands of cyanobacterial genes in the nucleus. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 12246–12251 (2002).
44. Sloan, D. B. & Moran, N. A. Genome Reduction and Co-evolution between the Primary and Secondary Bacterial Symbionts of Psyllids. *Mol. Biol. Evol.* **29**, 3781–3792 (2012).
45. Mohanta, T. K., Pudake, R. N. & Bae, H. Genome-wide identification of major protein families of cyanobacteria and genomic insight into the circadian rhythm. *Eur. J. Phycol.* **52**, (2017).
46. Reyes-Prieto, A., Weber, A. P. M. & Bhattacharya, D. The Origin and Establishment of the Plastid in Algae and Plants. *Annu. Rev. Genet.* **41**, 147–168 (2007).
47. Martin, W. *et al.* Gene transfer to the nucleus and the evolution of chloroplasts. *Nature* **393**, 162 (1998).
48. Mallick, P., Boutz, D. R., Eisenberg, D. & Yeates, T. O. Genomic evidence that the intracellular proteins of archaeal microbes contain disulfide bonds. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 9679–9684 (2002).
49. Brawley, S. H. *et al.* Insights into the red algae and eukaryotic evolution from the genome of *Porphyra umbilicalis* (Bangioophyceae, Rhodophyta). *Proc. Natl. Acad. Sci.* **114**, E6361 LP-E6370 (2017).
50. Dring, M. J., Wagner, A., Boeskov, J. & Lüning, K. Sensitivity of intertidal and subtidal red algae to UVA and UVB radiation, as monitored by Chlorophyll fluorescence measurements: Influence of collection depth and season, and length of irradiation. *Eur. J. Phycol.* **31**, 293–302 (1996).
51. Gröniger, A., Hallier, C. & Häder, D. P. Influence of UV radiation and visible light on *Porphyra umbilicalis*: Photoinhibition and MAA concentration. *J. Appl. Phycol.* **11**, 437 (1999).
52. Sillero, A. & Ribeiro, J. M. Isoelectric points of proteins: Theoretical determination. *Anal. Biochem.* **179**, 319–325 (1989).
53. Knight, C. G., Kassen, R., Hebestreit, H. & Rainey, P. B. Global analysis of predicted proteomes: functional adaptation of physical properties. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 8390–8395 (2004).
54. Gibson, F. & Pittard, J. Pathways of biosynthesis of aromatic amino acids and vitamins and their control in microorganisms. *Bacteriol. Rev.* **32**, 465–492 (1968).
55. Tzin, V. & Galili, G. New Insights into the Shikimate and Aromatic Amino Acids Biosynthesis Pathways in Plants. *Mol. Plant* **3**, 956–972 (2010).
56. Szabados, L. & Savouré, A. Proline: a multifunctional amino acid. *Trends Plant Sci.* **15**, 89–97 (2010).
57. Galili, G. & Höfgen, R. Metabolic Engineering of Amino Acids and Storage Proteins in Plants. *Metab. Eng.* **4**, 3–11 (2002).
58. Yoshida, S. Biosynthesis and Conversion of Aromatic Amino Acids in Plants. *Annu. Rev. Plant Physiol.* **20**, 41–62 (1969).
59. Light, S. H. & Anderson, W. F. The diversity of allosteric controls at the gateway to aromatic amino acid biosynthesis. *Protein Sci.* **22**, 395–404 (2013).
60. Berg, J., Tymoczko, J. & Stryer, L. Amino Acid Biosynthesis Is Regulated by

- Feedback Inhibition. in *Biochemistry* (WH Freeman, 2002).
61. Nelson, D. & Cox, M. Amino acid biosynthesis is under allosteric regulation. in *Lehlingers Principle of Biochemistry* (ed. Lehninger, A. L.) 833–880 (2005).
62. Nelson, D. & Cox, M. Biosynthesis of amino acids, nucleotides, and related molecules. in *Lehlinger Principles of Biochemistry* 842 (W.H. Freeman and Company, 2005).
63. Novoselov, S. V *et al.* Selenoproteins and selenocysteine insertion system in the model plant cell system, *Chlamydomonas reinhardtii*. *EMBO J.* **21**, 3681–3693 (2002).
64. Reeves, M. A. & Hoffmann, P. R. The human selenoproteome: recent insights into functions and regulation. *Cell. Mol. Life Sci.* **66**, 2457–2478 (2009).
65. Bodley, J. W. & Davie, E. W. A study of the mechanism of ambiguous amino acid coding by poly U: the nature of the products. *J. Mol. Biol.* **18**, 344–355 (1966).

656 **Table 1**
657 Details of plant proteome. Table shows acidic *pI* of proteins predominates the basic *pI*. However, in sea weed *Porphyra umbilicalis*, basic *pI*
658 predominates over the acidic *pI*. Putative polyketide synthase type I found in lower eukaryote *Volvox carteri* was found to be the largest protein in
659 the plant lineage. However, titin was found to be the largest protein in the higher eukaryotic land plants. Asterisks represents no specific data
660 available for the said item.

Name of the species	Total No. Of Protein Sequences Studied	Highest Mol. Wt. (kDa) of Protein	Name of the Protein with Highest Mol. Wt.	Lowest Mol. Wt. (kDa) of Protein	Name of the Protein with lowest Mol. Wt.	Highest <i>pI</i> of Protein	Name of the Protein with highest <i>pI</i>	Lowest <i>pI</i> of Protein	Name of the Protein with lowest <i>pI</i>	No. Of Proteins in Acidic <i>pI</i>	No. Of Proteins in Basic <i>pI</i>	No. Of Proteins in Neutral <i>pI</i>	Average of Acidic <i>pI</i>	Average of basic <i>pI</i>
<i>Aegilops tauschii</i>	55713	605.202	Misin	3.169	Cyt b6/f VIII	13.159	SARMP 2-like	2.409	Nucleolin-like	30488	25116	109	5.64	8.75
<i>Amaranthus hypochondriacus</i>	23879	596.11	Unknown	0.63	Unknown	12.67	Unknown	2.498	Unknown	13009	10833	37	5.58	8.36
<i>Amborella trichopoda</i>	27313	554.72	Unknown	4.32	Unknown	12.93	Unknown	2.587	Unknown	14583	12677	53	5.53	8.56
<i>Anacardium occidentale</i>	82170	457.50	Unknown	3.10	Unknown	12.61	Unknown	2.58	Unknown	46898	35100	172	5.65	8.27
<i>Ananas comosus</i>	35775	605.85	Midasin	3.16	Cyt b6/f VIII	12.89	Unknown	3.14	PPCD VHS3-like	20268	15399	108	5.68	8.35
<i>Aquilegia coerulea</i>	41063	620.28	Unknown	3.82	Unknown	12.74	Unknown	2.85	Unknown	23818	17169	76	5.64	8.26
<i>Arabidopsis halleri</i>	26911	609.56	Unknown	3.18	Unknown	12.50	Unknown	3.10	Unknown	14727	12123	61	5.59	8.32
<i>Arabidopsis lyrata</i>	39161	611.10	Midasin	2.51	Unknown	12.74	60S RP L41	2.85	RNA Pol. II Med 17	22213	16854	94	5.62	8.26
<i>Arabidopsis thaliana</i>	48350	611.88	Misin-like	0.57	Hypothetical	12.74	60S RP L41	2.75	Glycine-rich protein	27305	20926	119	5.61	8.31
<i>Arabis alpina</i>	23286	565.05	Unknown	***	****	12.79	Unknown	2.14	Unknown	13427	9810	49	5.52	8.37
<i>Arachis duranensis</i>	52826	617.77	Misin	4.11	DDHGT 4A	12.72	PERK2	2.96	Unknown	29514	23184	128	5.67	8.25
<i>Arachis ipaensis</i>	57621	617.63	Misin	4.08	DDHGT 4A	12.36	Unknown	3.13	Small acidic protein	31471	26007	143	5.67	8.26
<i>Asparagus officinalis</i>	36763	608.71	Misin	4.28	DDHGT 4A	12.50	protein TPRXL	2.89	FS CAYBR BP	20748	15934	81	5.64	8.25
<i>Auxenochlorella protothecoides</i>	7014	1649.26	PKS	***	***	13.21	Mucin-I	2.35	Hypothetical protein	4131	2874	9	5.52	8.91
<i>Bathycoccus prasinos</i>	7900	1814.10	Unnamed	3.32	Ycf12	12.74	Unknown	3.42	Unknown	4799	3093	8	5.47	8.38
<i>Beta vulgaris</i>	32874	617.35	Misin	3.56	Unknown	12.55	Proline-rich P	3.16	Shematin-like 2	18770	14017	87	5.68	8.19
<i>Botryococcus braunii</i>	23685	522.853	Unknown	3.068	Unknown	12.82	Unknown	2.21	Unknown	11176	12439	70	5.61	8.66

<i>Brachipodium stacei</i>	36357	605.43	Unknown	3.023	Unknown	12.88	Unknown	3.05	Unknown	19318	16954	85	5.67	8.59
<i>Brachipodium sylvaticum</i>	50263	608.40	Unknown	3.14	Unknown	12.88	Unknown	3.00	Unknown	26022	24130	111	5.66	8.82
<i>Brachipodium distachyon</i>	33944	605.24	Misin	3.16	Cyt b6/f VIII	12.39	Unknown	3.03	Prothymosin α -B-like	19814	14052	78	5.67	8.35
<i>Brachypodium hybridum</i>	80980	605.62	Unknown	3.02	Unknown	13.14	Unknown	3.05	Unknown	40744	40092	144	5.66	8.79
<i>Brassica napus</i>	123465	606.90	Misin-like	3.15	petN	12.57	IQ domain 31	2.75	Shematin-like 2	68255	54929	281	5.62	8.29
<i>Brassica oleracea</i>	56687	606.74	Misin	3.48	Unknown	12.54	IQ domain 31	2.75	Shematin-like 2	31109	25432	146	5.62	8.29
<i>Brassica rapa</i>	52553	607.93	Misin	3.33	Unknown	12.57	IQ domain 31	2.48	Dentin-sialophospho Protein-like	29429	23004	120	5.62	8.29
<i>Cajanus cajan</i>	38965	619.70	Misin	3.42	Cyt b6/f VIII	12.44	CLAVATA 3/ESR	2.82	RPB1-like	22092	16793	80	5.71	8.22
<i>Camelina sativa</i>	107481	610.42	Misin-like	2.61	Peptide POLARIS	12.72	SARMP 2-like	2.13	TsetseEP-like	60623	46584	274	5.61	8.26
<i>Capsella grandiflora</i>	26561	593.86	Unknown	3.41	Unknown	12.44	Unknown	2.80	Unknown	14831	11661	61	5.62	8.30
<i>Capsella rubella</i>	34126	610.40	Misin	3.57	Unknown	12.44	Lifeguard 1	3.13	Unknown	19477	14578	71	5.63	8.25
<i>Capsicum annuum</i>	45410	617.56	Misin	3.16	Cyt b6/f VIII	12.79	GRCW protein	2.75	GRCW protein 1	25339	19970	101	5.67	8.24
<i>Capsicum baccatum</i>	35853	553.75	Auxin TP BIG	3.16	Cyt b6/f VIII	12.64	Hypothetical al	1.99	Hypothetical	20479	15300	74	5.55	8.42
<i>Capsicum chinense</i>	34973	550.87	Auxin TP BIG	3.16	Cyt b6/f VIII	13.21	Hypothetical al	2.24	Hypothetical	20497	14398	79	5.56	8.25
<i>Carica papaya</i>	26103	446.01	Unknown	3.16	Cyt b6/f VIII	12.29	50S RP L34	3.10	Small acidic protein	14512	11526	65	5.67	8.25
<i>Cephalotus follicularis</i>	36667	611.81	AAA_5	***	****	13.04	Hypothetical al	2.9	Hypothetical	18643	17959	65	5.66	8.33
<i>Chenopodium quinoa</i>	63173	621.14	Misin-like	3.16	petN	12.22	SR45-like	3.02	Unknown	36037	26973	163	5.66	8.18
<i>Chlamydomonas eustigma</i>	14161	1370.23	Unknown	2.28	Unknown	13.14	Unknown	2.18	Unknown	9089	5036	36	5.64	8.23
<i>Chlamydomonas reinhardtii</i>	14488	2056.44	PKS	1.46	Unknown	13.34	Unknown	2.40	Unknown	7427	7029	32	5.64	8.7
<i>Chlorella variabilis</i>	9780	1159.35	Unknown	4.91	Unknown	12.73	Unknown	2.88	Unknown	5883	3876	21	5.55	8.55
<i>Chromochloris zofingiensis</i>	15369	1932.21	Unknown	2.34	Unknown	12.5	Unknown	2.54	Unknown	9592	5740	37	5.61	8.35
<i>Cicer arietinum</i>	33107	616.07	Unknown	3.02	Unknown	12.25	Unknown	3.14	Unknown	19043	13979	85	5.69	8.20
<i>Citrus clementina</i>	34557	588.15	Unknown	****	****	12.44	Unknown	2.95	Unknown	19332	15151	74	5.67	8.28

<i>Citrus sinensis</i>	35648	617.00	Misin	3.16	Cyt b6/f VIII	12.25	Unknown	2.93	Circumsporozoite protein-like	20887	14680	81	5.69	8.18
<i>Citrus unshiu</i>	37970	684.122	Unknown	0.54	Unknown	12.83	Unknown	2.74	Unknown	20954	16926	90	5.67	8.34
<i>Coccomyxa subellipsoidea</i>	9839	1632.35	Ketoacyl-synt	5.06	Unknown	12.99	Unknown	3.05	Unknown	6112	3708	19	5.53	8.55
<i>Corchorus capsularis</i>	29356	606.73	Unknown	1.6	Unknown	12.74	Unknown	2.72	IMP	15540	13783	33	5.49	8.67
<i>Corchorus olitorius</i>	35704	498.84	ZRF	2.79	Unknown	12.79	Unknown	2.56	Unknown	19280	16359	65	5.49	8.62
<i>Cucumis melo</i>	29796	616.61	Misin	2.87	Unknown	13.23	Unknown	2.56	Loricin-like	17151	12562	83	5.7	8.24
<i>Cucumis sativus</i>	29796	617.61	Misin	2.87	Unknown	13.23	Unknown	2.56	Loricin-like	17151	12562	83	5.7	8.24
<i>Cucurbita maxima</i>	42777	615.95	Misin	3.58	Unknown	12.22	60S RP L39	2.54	CWP gp 1-like	24870	17817	90	5.68	8.25
<i>Cucurbita moschata</i>	43715	615.32	Misin	4.1	DDGT 4A	12.85	EPR1	2.31	PKDP	25399	18233	83	5.67	8.25
<i>Daucus carota</i>	44655	619.12	Misin	3.16	Cyt b6/f VIII	12.32	Ribo Protein L32	2.79	Loricin	26135	18423	97	5.68	8.16
<i>Dendrobium officinale</i>	34527	616.94	Misin	3.16	Cyt b6/f VIII	12.45	60S RP L39	3.23	Unknown	19029	15425	73	5.69	8.26
<i>Dichanthelium oligosanthes</i>	26468	538.41	Auxin TP BIG	1.09	Unknown	13.18	Unknown	3.00	Unknown	14146	12261	61	5.6	8.63
<i>Dorcoceras hygrometricum</i>	47778	563.43	Midasin	4.75	Unknown	12.86	Unknown	2.68	Unknown	23461	24237	80	5.44	8.92
<i>Dunaliella salina</i>	18801	603.54	Unknown	2.92	Unknown	12.69	Unknown	2.38	Unknown	9927	8833	41	5.66	8.53
<i>Durio zibethinus</i>	63007	620.96	Misin	3.57	Unknown	12.35	SR45-like	3.12	Acidic protein	37032	25800	175	5.68	8.21
<i>Elaeis guineensis</i>	41887	614.29	Misin	3.19	Cyt b6/f VIII	12.35	50S RP L34	3.28	Calsequestrin 1-like	24529	17266	92	5.69	8.29
<i>Erythranthe guttata</i>	31861	611.80	Misin	3.77	Unknown	12.14	SR45	2.6	CSF subunit 2	18284	13500	77	5.63	8.22
<i>Eucalyptus grandis</i>	52554	644.40	Futsch	3.4	Cyt b6/f VIII	12.83	Unknown	3.07	Fimbrin 1-like	31377	21034	143	5.69	8.24
<i>Eutrema salsugineum</i>	29485	609.35	Unknown	***	****	12.32	Unknown	3.05	Unknown	16105	13300	80	5.63	8.32
<i>Fragaria vesca</i>	31387	609.82	Misin	3.19	Cyt b6/f VIII	12.34	50S RP L34	3.13	Prostatic spermine BP	18862	12429	96	5.65	8.20
<i>Genlisea aurea</i>	17685	559.24	Unknown	4.94	Unknown	12.89	Unknown	2.93	Unknown	9842	7814	29	5.54	8.51
<i>Glycine max</i>	71523	619.18	Misin-like	1.34	AAPT	12.28	Unknown	2.67	HC1-like	41545	29785	193	5.68	8.22
<i>Glycine soja</i>	50399	603.79	Misin	1.59	Hypothetical	12.42	Dynein	2.21	RBP 12B	28254	22054	91	5.62	8.34
<i>Gonium pectorale</i>	16290	881.59	Unknown	5.02	Unknown	12.88	Unknown	2.48	Unknown	10225	6034	31	5.52	8.54

<i>Gossypium arboreum</i>	47568	806.46	Titin-like	3.16	Cyt b6/f VIII	12.23	SR45-like	2.91	Loricin-like	26844	20612	103	5.69	8.23
<i>Gossypium hirsutum</i>	90927	750.00	Titin-like	3.16	Cyt b6/f VIII	12.26	SARMP 2-like	2.67	ER TF TINY-like	51343	39336	248	5.67	8.25
<i>Gossypium raimondii</i>	59057	802.74	Titin-like	3.16	Cyt b6/f VIII	12.25	SR45	2.98	Arabinogalactan 11	33377	25561	119	5.67	8.25
<i>Handroanthus impetiginosus</i>	30271	475.38	Unknown	2.91	Unknown	12.66	Unknown	2.95	Unknown	16255	13959	57	5.64	8.36
<i>Helianthus annuus</i>	73839	1300.83	Unknown	2.76	Unknown	12.88	C1E8.05-like	2.62	Prostatic spermine BP	40355	33310	174	5.63	8.25
<i>Helicosporidium sp.</i>	6033	245.59	Unknown	***	***	12.79	Unknown	2.82	Unknown	3799	2226	8	5.32	8.92
<i>Herrania umbratica</i>	27748	620.76	Misin	4.23	DDGT 4A	12.66	Unknown	3.12	Small acidic protein 1	16209	11492	47	5.7	8.21
<i>Hevea brasiliensis</i>	58062	620.10	Misin	3.16	Cyt b6/f VIII	12.16	50S RP L35	3.12	Small acidic protein 1	33994	23952	116	5.69	8.20
<i>Hordeum vulgare</i>	248180	607.08	Unknown	1.75	Unknown	13.11	****	2.52	Unknown	132250	115421	509	5.60	8.62
<i>Ipomoea nil</i>	51054	636.34	Filaggrin-like	3.16	Cyt b6/f VIII	12.38	Formin 2-like	2.82	Nucleolin-like	29183	21727	144	5.69	8.18
<i>Jatropha curcas</i>	32547	628.22	Titin	3.16	petN	12.14	60S RP L39-3	3.13	Glycin-rich protein	18939	13531	77	5.69	8.18
<i>Juglans regia</i>	55627	624.15	Midasin	4.14	40S RP S29-like	12.98	Formin-like 3	2.46	Glycine-rich protein	32225	23292	110	5.69	8.22
<i>Kalanchoe fedtschenkoi</i>	45190	563.28	Unknown	3.16	Unknown	12.98	Unknown	2.88	Unknown	25024	20062	104	5.64	8.37
<i>Kalanchoe laxiflora</i>	69177	619.54	Unknown	3.01	Unknown	12.47	Unknown	3.09	Unknown	39148	29888	141	5.66	8.32
<i>Klebsormidium nitens</i>	16282	813.31	Unknown	***	***	12.83	Ser/Thr Prot Kin	3.09	Unknown	9751	6500	31	5.57	8.39
<i>Lactuca sativa</i>	45242	609.98	Misin	3.16	Cyt b6/f VIII	12.44	Glh-2-like	2.56	Ctenidin-3-like	25604	19492	146	5.67	8.20
<i>Linum usitatissimum</i>	43484	544.30	Unknown	4.95	Unknown	12.41	Unknown	2.49	Unknown	24926	18459	99	5.58	8.34
<i>Lupinus angustifolius</i>	52821	619.31	Misin	5.59	Arabinogalactan peptide 23	13.07	Collagen α -2(V) chain-like	2.76	Glutamic acid rich protein	31045	21650	126	5.67	8.21
<i>Macleaya cordata</i>	21911	624.74	Von Willbrand factor	3.87	Unknown	12.32	Unknown	2.94	Unknown	12657	9206	48	5.61	8.29
<i>Malus domestica</i>	60544	551.44	Auxin TP BIG	3.45	Unknown	12.91	CDPK	2.79	IFF6-like	34853	25548	143	5.65	8.25
<i>Manihot acuminata</i>	36528	516.85	Unknown	1.01	Unknown	12.88	Unknown	2.88	Unknown	18516	17936	76	5.64	8.58
<i>Manihot esculenta</i>	43286	621.75	Midasin	3.16	Cyt b6/f VIII	12.19	SR45-like	2.65	ASF1-like	25698	17491	97	5.68	8.18

<i>Marchantia polymorpha</i>	17956	806.12	Unknown	7.02	Unknown	12.10	Unknown	3.10	Unknown	10142	7793	21	5.54	8.48
<i>Medicago truncatula</i>	57661	611.94	Misin	1.90	NCR peptide	12.74	Unknown	2.57	LEA	30526	27027	108	5.61	8.40
<i>Micromonas commoda</i>	10137	1532.91	PKS	2.78	Antisense noncoding	12.61	Unknown	2.95	Unknown	6417	3703	17	5.4	8.61
<i>Micromonas pusilla</i>	10242	848.29	Unknown	5.07	Unknown	13.37	Unknown	2.80	Unknown	5985	4242	15	5.37	8.97
<i>Miscanthus sinensis</i>	89486	615.13	Unknown	2.90	Unknown	13.27	Unknown	2.88	Unknown	45710	43546	230	5.65	8.76
<i>Momordica charantia</i>	28666	616.95	Misin	4.21	DDG 4A	12.22	60S RP L39	3.16	Small acidic protein 1	16621	11997	48	5.69	8.20
<i>Monoraphidium neglectum</i>	16755	730.02	Misin	4.12	Unknown	13.01	Unknown	2.79	Unknown	8940	7783	32	5.44	8.91
<i>Morus notabilis</i>	26965	566.71	Auxin TP BIG	5.11	Unknown	12.32	Unknown	3.10	Unknown	13932	12984	49	5.60	8.55
<i>Musa acuminata</i>	47707	616.23	Misin	3.79	DGG 4A	12.44	Unknown	3.09	TUB8	27400	20184	123	5.69	8.28
<i>Nelumbo nucifera</i>	38191	797.88	Titin	3.16	Cyt b6/f VIII	12.22	60S RP L39	2.95	Prostatic spermine BP	22308	15782	101	5.70	8.21
<i>Nicotiana attenuate</i>	44491	616.08	Misin	4.15	DGG 4A	12.19	SR45	2.99	Unknown	23898	20492	101	5.67	8.24
<i>Nicotiana glauca</i>	48160	564.57	Auxin TP BIG	1.30	Unknown	12.44	Unknown	2.95	Unknown	26496	21565	99	5.66	8.25
<i>Nicotiana tabacum</i>	84255	539.85	Auxin TP BIG	3.16	Cyt b6/f VIII	12.44	Unknown	2.86	mRNA decay protein	46302	37768	185	5.65	8.24
<i>Nicotiana tomentosiformis</i>	48962	564.92	Aux TP BIG	3.18	Cyt b6/f VIII	12.25	Cell wall protein	3.09	Arabinogalactan peptide 14-like	27278	21567	117	5.67	8.22
<i>Olea europaea</i>	58334	567.49	Misin	3.16	Cyt b6/f VIII	12.36	Unknown	3.18	Acidic protein 2-like	32519	25690	125	5.65	8.23
<i>Oryza brachyantha</i>	26803	597.14	Misin	5.73	Unknown	12.44	NFD6	3.03	Prostatic spermine BP	16083	10659	61	5.64	8.30
<i>Oryza sativa</i>	37358	567.80	Unknown	3.22	Unknown	12.64	Unknown	2.65	Unknown	20560	16732	66	5.56	8.69
<i>Ostreococcus lucimarinus</i>	7603	1994.71	PKS	4.12	Cysteine protein	12.36	Unknown	2.85	Unknown	4879	2711	13	5.40	8.62
<i>Ostreococcus tauri</i>	7766	1237.34	Unknown	3.29	Ycf12	12.32	L39e	3.28	SVC	4732	3018	16	5.47	8.67
<i>Panicum hallii</i>	49825	608.88	Unknown	1.87	Unknown	13.07	Unknown	3.05	Unknown	24843	24878	104	5.68	8.79
<i>Phalaenopsis equestris</i>	29894	568.40	Auxin TP BIG	3.20	Cyt b6/f VIII	12.5	50S L34	2.75	Unknown	16701	13099	94	5.69	8.27
<i>Phaseolus vulgaris</i>	32720	617.41	Unknown	***	***	12.74	Unknown	3.09	Unknown	18577	14073	70	5.68	8.29
<i>Phoenix dactylifera</i>	38570	617.67	Misin	3.19	Cyt b6/f	12.45	60S RP L39	3.31	PPCD VHS3-like	22015	16450	105	5.67	8.31

<i>Physcomitrella patens</i>	35934	553.40	Unknown	3.21	Unknown	12.72	Unknown	2.77	Unknown	19655	16205	74	5.60	8.47
<i>Populus deltoides</i>	57249	567.25	Unknown	3.14	Unknown	12.54	Unknown	2.85	Unknown	31040	26084	125	5.64	8.37
<i>Populus euphratica</i>	49760	619.99	Misin	3.16	Cyt b6/f	12.22	60S RP L39	3.02	Calsequestrin -1-like	29358	20278	124	5.68	8.19
<i>Populust richocarpa</i>	45942	603.68	Misin	***	***	12.74	Unknown	2.54	Unknown	25431	20413	98	5.62	8.35
<i>Porphyra umbilicalis</i>	13360	480.60	Unknown	3.18	Cyt b6/f	13.33	Unknown	2.6	Unknown	3982	9365	13	5.47	10.40
<i>Prunus avium</i>	35009	607.05	Misin	3.84	Unknown	12.22	SR45	2.86	Unknown	20962	13972	75	5.68	8.18
<i>Prunus mume</i>	29705	678.34	Unknown	3.19	Cyt b6/f VIII	12.22	60S RP L39	2.54	Cell wall protein gp1	17588	12054	63	5.68	8.21
<i>Prunus persica</i>	32595	607.17	Misin	3.19	Cyt b6/f VIII	12.22	SR45	2.98	SCP SP60-like	19315	13201	79	5.70	8.18
<i>Punica granatum</i>	50476	1150.02	Unknown	1.21	Unknown	13.96	Unknown	2.07	Unknown	23078	27314	84	5.52	8.93
<i>Pyrus bretschneideri</i>	47086	1269.42	Unknown	4.13	Unknown	12.23	50S RP L34	3.10	Unknown	27610	19365	111	5.67	8.25
<i>Raphanus sativus</i>	61216	607.96	Misin	2.89	Unknown	12.58	IQ domain 31	2.79	Shematrin-like 2	34204	26871	141	5.61	8.30
<i>Ricinus communis</i>	27998	619.84	Misin	3.16	Cyt b6/f VIII	12.47	Unknown	3.12	Small acidic protein	16138	11789	71	5.68	8.21
<i>Salix purpurea</i>	61520	621.82	Unknown	3.26	Unknown	12.61	Unknown	2.77	Unknown	35045	26358	117	5.64	8.29
<i>Selaginella moellendorffii</i>	34746	909.93	Unknown	5.18	Unknown	12.45	Unknown	3.28	Unknown	20404	14243	99	5.66	8.27
<i>Sesamum indicum</i>	35410	614.62	Misin	3.16	Cyt b6/f VIII	12.64	L32	2.88	BCP1-like	20353	14974	83	5.68	8.21
<i>Setaria italica</i>	35844	608.65	Misin	3.16	Cyt b6/f VIII	12.44	NFD6	3.13	Prostatic spermine BP	20727	15027	90	5.7	8.36
<i>Solanum lycopersicum</i>	36008	620.33	Misin	3.16	Cyt b6/f VIII	12.29	SR45	2.46	Cell wall protein	21001	14925	82	5.67	8.20
<i>Solanum pennellii</i>	35068	620.19	Misin	3.66	Unknown	12.16	SR45	2.20	Myb-like	19944	15043	81	5.67	8.20
<i>Solanum tuberosum</i>	37960	618.76	Misin	3.16	Cyt b6/f VIII	13.24	Extension-like	2.95	Tripartite motif 44	22261	15614	85	5.67	8.19
<i>Sorghum bicolor</i>	39248	615.07	Misin	3.16	Cyt b6/f VIII	12.57	Unknown	3.00	Unknown	21947	17200	101	5.68	8.43
<i>Sphagnum fallax</i>	32298	1027.64	Unknown	3.19	Unknown	12.19	Unknown	2.62	Unknown	19972	12256	70	5.61	8.32
<i>Spinacia oleracea</i>	32794	615.44	Misin	2.50	SpolCp151	12.42	EPR1	2.56	Unknown	18350	14357	87	5.66	8.20
<i>Spirodela polyrhiza</i>	19623	596.16	Unknown	3.16	Cyt b6/f VIII	12.69	Unknown	2.74	Unknown	10680	8903	40	5.60	8.52
<i>Tarenaya hassleriana</i>	41094	614.85	Midasin	4.03	Unknown	12.41	IQ-domain 14	3.04	Unknown	23230	17762	102	5.65	8.28

<i>Theobroma cacao</i>	30854	621.12	Midasin	3.16	Cyt b6/f VIII	12.21	SR45	3.12	Small acidic protein	17897	12869	88	5.70	8.20
<i>Trifolium pratense</i>	63799	566.60	Auxin TP BIG	***	***	12.88	Unknown	2.34	Unknown	37487	26211	101	5.36	8.45
<i>Trifolium subterraneum</i>	42059	571.83	Unknown	***	***	12.33	Unknown	2.6	Unknown	24271	17701	87	5.28	8.28
<i>Triticum aestivum</i>	250	230.08	Unknown	6.54	Unknown	12.15	Unknown	3.75	Unknown	147	103	0	5.59	8.22
<i>Triticum urartu</i>	24169	559.02	UBR4	4.5	Unknown	13.27	Unknown	2.96	Unknown	13783	10339	47	5.56	8.51
<i>Vigna angularis</i>	37769	621.26	Midasin	2.96	Unknown	12.22	50S RP L34	2.81	Clumping factor A	21862	15798	109	5.7	8.21
<i>Vigna radiata</i>	42284	624.61	Midasin	3.16	Cyt b6/f VIII	12.41	Formin-like 6	2.89	ATF7IP	24545	17654	83	5.69	8.21
<i>Vigna unguiculata</i>	42287	616.90	Unknown	3.17	Cyt b6/f VIII	12.91	Unknown	2.95	Unknown	24082	18103	102	5.69	8.35
<i>Vitis vinifera</i>	41208	622.14	Midasin	3.16	Cyt b6/f VIII	12.19	Orf19	3.23	Circumsporozite	25295	15837	76	5.70	8.19
<i>Volvox carteri</i>	14436	2236.80	Putative PKS I	4.89	Unknown	13.79	Unknown	2.42	Unknown	7560	6849	27	5.64	8.53
<i>Zostera marina</i>	20450	606.86	*****	1.68	*****	12.75	*****	2.42	*****	10988	9417	45	5.61	8.36
Average	40469.47	707.23				12.62				22820	17794.26	91.36	5.62	8.37

661

662

663 **Abbreviations:** PPCD VHS3-like: phosphopantothenoylcysteine decarboxylase subunit VHS3-like isoform X2, 60S RP L41: 60S ribosomal protein
664 subunit L41, DDHGT 4A: Dolichyl-diphosphooligosaccharideprotein glycosyltransferase subunit 4A, PERK2: proline-rich receptor-like protein
665 kinase PERK2, RNA Pol. II Med 17: mediator of RNA polymerase II subunit 17, FS CAYBR BP: fibrous sheath CABYR-binding protein-like,
666 PKS: polyketide synthase, PS-I RC N: photosystem I reaction center subunit N, chloroplastic, partial; RPB1: DNA-directed RNA polymerase II
667 subunit RPB1-like isoform X3, GRCW protein: glycine-rich cell wall structural protein, SR45: serine/arginine-rich splicing factor SR45-like,
668 IMP:inosine-5'-monophosphate cyclohydrolase, ZRF: zinc ring finger-type, CWP:cell wall protein, DDGT: dolichyl-
669 diphosphooligosaccharideprotein glycosyltransferase subunit 4A, PKDP: polycystic kidney disease protein 1-like 3, CSF: cleavage stimulation
670 factor subunit 2 tau variant-like, AAPT: aminoalcoholphosphotransferase, RBP: RNA binding protein, SARMP: serine/arginine repetitive matrix
671 protein 2-like, ER TF: ethylene responsive transcription factor, CDPK: cyclin-dependent serine/threonine-protein kinase, LEA: late embryogenesis
672 associated protein, DDG: dolichyl-diphosphooligosaccharideprotein glycosyltransferase subunit 4A, NFD: nuclear fusion defective, SVC: satellite
673 virus coat protein, PPCD: phosphopantothenoylcysteine decarboxylase, SCP: spore coat protein, ATF7IP: activating transcription factor 7-
674 interacting protein 1.

Table 2

Abundance of various amino acids in different species. The second column represents the average abundance whereas the third and fourth column represent variation (high and low) in amino acid composition in different species from the average.

Amino Acids	Average abundance (%)	High Abundance (%)	Low Abundance (%)
Ala	7.68	<i>Porphyra umbilicalis</i> (17.58), <i>Monoraphidium neglectum</i> (16.58), <i>Gonium pectorale</i> (15.45), <i>Chlorella variabilis</i> (15.15), <i>Chlamydomonas reinhardtii</i> (14.57), <i>Micromonas pusilla</i> (14.20), <i>Auxenochlorella protothecoides</i> (13.69), <i>Volvox carteri</i> (13.29), <i>Helicosporidium</i> sp. (12.75), <i>Micromonas commoda</i> (12.61), <i>Coccomyxa subellipsoidea</i> (12.29), <i>Ostreococcus lucimarinus</i> (11.89), <i>Chromochloris zofingiensis</i> (11.57), <i>Dunaliella salina</i> (11.53), <i>Ostreococcus tauri</i> (11.47), <i>Klebsormidium nitens</i> (10.62), <i>Botryococcus braunii</i> (9.86), <i>Dichanthelium oligosanthos</i> (9.71), <i>Chlamydomonas eustigma</i> (9.53)	<i>Picea glauca</i> (5.12), <i>Medicago truncatula</i> (5.8), <i>Lactuca sativa</i> (5.81), <i>Zostera marina</i> (5.86)
Asp	5.32	<i>Micromonas pusilla</i> (6.68), <i>Micromonas commode</i> (6.46), <i>Ostreococcus lucimarinus</i> (6.38), <i>Ostreococcus tauri</i> (6.33), <i>Bathycoccus prasinos</i> (6.05)	<i>Picea glauca</i> (3.46), <i>Dunaliella salina</i> (4.28), <i>Chlorella variabilis</i> (4.41), <i>Chlamydomonas reinhardtii</i> (4.61), <i>Porphyra umbilicalis</i> (4.63), <i>Volvox carteri</i> (4.71), <i>Monoraphidium neglectum</i> (4.71), <i>Botryococcus braunii</i> (4.78), <i>Gonium pectorale</i> (4.78)
Glu	6.43	<i>Bathycoccus prasinos</i> (8.44), <i>Klebsormidium nitens</i> (7.2),	<i>Porphyra umbilicalis</i> (3.52), <i>Picea glauca</i>

		<i>Ostreococcus tauri</i> (7.18), <i>Ostreococcus lucimarinus</i> (7.03)	(4.02), <i>Chromochloris</i> <i>zofingiensis</i> (4.66), <i>Chlamydomonas</i> <i>reinhardtii</i> (5.15), <i>Volvox carteri</i> (5.21), <i>Monoraphidium</i> <i>neglectum</i> (5.35)
Phe	3.97	<i>Arachis duranensis</i> (6.18)	<i>Porphyra umbilicalis</i> (2.01), <i>Gonium</i> <i>pectorale</i> (2.4), <i>Monoraphidium</i> <i>neglectum</i> (2.44), <i>Volvox carteri</i> (2.49), <i>Dunaliella salina</i> (2.56), <i>Chromochloris</i> <i>zofingiensis</i> (2.58), <i>Chlamydomonas</i> <i>reinhardtii</i> (2.59), <i>Chlorella variabilis</i> (2.68), <i>Chlamydomonas</i> <i>eustigma</i> (2.89), <i>Auxenochlorella</i> <i>protothecoides</i> (2.92)
Gly	6.80	<i>Porphyra umbilicalis</i> (11.76), <i>Monoraphidium neglectum</i> (10.51), <i>Gonium pectorale</i> (10.33), <i>Chlamydomonas reinhardtii</i> (9.58), <i>Volvox carteri</i> (9.54), <i>Chlorella</i> <i>variabilis</i> (9.23), <i>Auxenochlorella</i> <i>protothecoides</i> (8.82), <i>Micromonas</i> <i>commoda</i> (8.81), <i>Micromonas</i> <i>pusilla</i> (8.59), <i>Klebsormidium</i> <i>nitens</i> (8.44), <i>Helicosporidium</i> sp. (8.29), <i>Dunaliella salina</i> (8.11), <i>Botryococcus braunii</i> (8.03)	<i>Arachis duranensis</i> (4.19)
His	2.4	<i>Dunaliella salina</i> (3.13)	<i>Bathycoccus prasinus</i> (1.87), <i>Ostreococcus</i> <i>tauri</i> (1.93), <i>Monoraphidium</i> <i>neglectum</i> (1.93), <i>Micromonas pusilla</i> (1.94), <i>Ostreococcus</i> <i>lucimarinus</i> (1.96)

Ile	4.94	<i>Zostera marina</i> (6.09)	<i>Porphyra umbilicalis</i> (1.77), <i>Monoraphidium neglectum</i> (2.34), <i>Gonium pectorale</i> (2.37), <i>Chlorella variabilis</i> (2.52), <i>Chlamydomonas reinhardtii</i> (2.60), <i>Volvox carteri</i> (2.78), <i>Helicosporidium</i> sp. (2.80), <i>Dunaliella salina</i> (2.82), <i>Auxenochlorella protothecoides</i> (2.96), <i>Micromonas pusilla</i> (2.98)
Lys	5.73	<i>Bathycoccus prasinos</i> (7.33)	<i>Porphyra umbilicalis</i> (2.08), <i>Gonium pectorale</i> (2.79), <i>Monoraphidium neglectum</i> (2.86), <i>Volvox carteri</i> (3.03), <i>Auxenochlorella protothecoides</i> (3.06), <i>Chlorella variabilis</i> (3.06), <i>Chlamydomonas reinhardtii</i> (3.07), <i>Helicosporidium</i> sp. (3.17), <i>Dunaliella salina</i> (3.39)
Leu	9.62	<i>Picea glauca</i> (13.12)	<i>Porphyra umbilicalis</i> (7.07), <i>Micromonas pusilla</i> (7.70), <i>Bathycoccus prasinos</i> (7.85), <i>Micromonas commoda</i> (7.98), <i>Ostreococcus tauri</i> (8.18), <i>Ostreococcus lucimarinus</i> (8.41)
Met	2.40	<i>Picea glauca</i> (3.75)	<i>Porphyra umbilicalis</i> (1.45), <i>Monoraphidium neglectum</i> (1.78),

			<i>Klebsormidium nitens</i> (1.84), <i>Auxenochlorella protothecoides</i> (1.95), <i>Gonium pectorale</i> (1.97), <i>Chlorella variabilis</i> (1.98), <i>Micromonas pusilla</i> (1.98), <i>Helicosporidium</i> sp. (1.99)
Asn	4.13	<i>Medicago truncatula</i> (5.09)	<i>Porphyra umbilicalis</i> (1.53), <i>Monoraphidium neglectum</i> (1.88), <i>Chlorella variabilis</i> (1.93), <i>Auxenochlorella protothecoides</i> (2.02), <i>Gonium pectorale</i> (2.11), <i>Helicosporidium</i> sp. (2.15), <i>Chlamydomonas reinhardtii</i> (2.38), <i>Micromonas pusilla</i> (2.49), <i>Volvox carteri</i> (2.49), <i>Dunaliella salina</i> (2.61), <i>Micromonas commoda</i> (2.71), <i>Coccomyxa subellipsoidea</i> (2.75), <i>Klebsormidium nitens</i> (2.82), <i>Botryococcus braunii</i> (2.89), <i>Ostreococcus tauri</i> (2.93), <i>Ostreococcus lucimarinus</i> (2.99)
Pro	5.10	<i>Porphyra umbilicalis</i> (9.2), <i>Botryococcus braunii</i> (7.10), <i>Dunaliella salina</i> (7.08), <i>Volvox carteri</i> (7.00), <i>Gonium pectorale</i> (6.85), <i>Chlamydomonas reinhardtii</i> (6.49), <i>Picea glauca</i> (6.45), <i>Chlorella variabilis</i> (6.41),	<i>Bathycoccus prasinos</i> (3.81)

		<i>Monoraphidium neglectum</i> (6.40), <i>Auxenochlorella protothecoides</i> (6.25), <i>Klebsormidium nitens</i> (6.11)	
Gln	3.74	<i>Dunaliella salina</i> (7.00), <i>Chromochloris zofingiensis</i> (6.20), <i>Monoraphidium neglectum</i> (5.60), <i>Chlorella variabilis</i> (5.56), <i>Chlamydomonas eustigma</i> (4.83), <i>Sphagnum fallax</i> (4.67), <i>Coccomyxa</i> <i>subellipsoidea</i> (4.59), <i>Volvox</i> <i>carteri</i> (4.46), <i>Botryococcus braunii</i> (4.35), <i>Chlamydomonas reinhardtii</i> (4.28), <i>Physcomitrella patens</i> (4.08), <i>Doroceras hygrometricum</i> (4.04), <i>Klebsormidium nitens</i> (4.03)	<i>Porphyra umbilicalis</i> (2.04), <i>Micromonas</i> <i>pusilla</i> (2.06), <i>Micromonas commoda</i> (2.58), <i>Ostreococcus</i> <i>tauri</i> (2.58), <i>Ostreococcus</i> <i>lucimarinus</i> (2.64)
Arg	5.68	<i>Porphyra umbilicalis</i> (9.81), <i>Micromonas pusilla</i> (8.12), <i>Ostreococcus tauri</i> (7.96), <i>Helicosporidium</i> sp. (7.74), <i>Micromonas commoda</i> (7.65), <i>Ostreococcus lucimarinus</i> (7.46), <i>Auxenochlorella protothecoides</i> (7.13)	<i>Medicago truncatula</i> (4.84), <i>Trifolium</i> <i>pratense</i> (4.95), <i>Cicer</i> <i>arietinum</i> (4.97), <i>Lactuca sativa</i> (4.99)
Ser	8.71	<i>Chlamydomonas eustigma</i> (9.72)	<i>Monoraphidium</i> <i>neglectum</i> (6.33), <i>Chlorella variabilis</i> (6.44), <i>Auxenochlorella</i> <i>protothecoides</i> (6.49), <i>Porphyra umbilicalis</i> (6.52), <i>Micromonas</i> <i>commoda</i> (6.57), <i>Ostreococcus</i> <i>lucimarinus</i> (6.72), <i>Micromonas pusilla</i> (6.74), <i>Gonium</i> <i>pectorale</i> (7.06), <i>Chlamydomonas</i> <i>reinhardtii</i> (7.11)
Thr	4.90	<i>Chromochloris zofingiensis</i> (5.80), <i>Ostreococcus tauri</i> (5.77), <i>Bathycoccus prasinos</i> (5.72), <i>Ostreococcus lucimarinus</i> (5.70), <i>Chlamydomonas eustigma</i> (5.44),	

<hr/>			
<i>Porphyra umbilicalis</i> (5.39), <i>Micromonas pusilla</i> (5.33), <i>Volvox carteri</i> (5.30)			
<hr/>			
Val	6.55	<i>Ostreococcus tauri</i> (7.42), <i>Ostreococcus lucimarinus</i> (7.35), <i>Porphyra umbilicalis</i> (7.32), <i>Micromonas commoda</i> (7.22), <i>Helicosporidium</i> sp. (7.20), <i>Micromonas pusilla</i> (7.11)	<i>Picea glauca</i> (5.26), <i>Dunaliella salina</i> (5.78)
<hr/>			

681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714

Table 3

Average abundance of different amino acids in plant proteome. Leu was high abundant whereas Trp was the low abundant amino acid in the plant kingdom. The average amino acid composition includes 5.8 million protein sequences from 145 plant species.

Biosynthetic pathways/Substrate	Amino acids	Average abundance (%) in Proteome
α -Ketoglutarate	Arginine	6.68
	Glutamate	6.43
	Glutamine	3.74
	Proline	5.10
Pyruvate	Alanine	7.68
	Isoleucine	4.94
	Leucine	9.62
	Valine	6.55
3-Phosphoglycerate	Glycine	6.80
	Cysteine	1.85
	Serine	8.71
Oxaloacetate	Asparagine	4.13
	Aspartate	5.32
	Lysine	5.73
	Methionine	2.40
	Threonine	4.90
Phosphoenolpyruvate & Erythrose 4-phosphate	Phenylalanine	3.97
	Tryptophan	1.28
	Tyrosine	2.67
Ribose 5-phosphate	Histidine	2.4

Figure legends

Figure 1

Principal component analysis (PCA) of acidic *pI* proteins. The PCA plot illustrates the relationship between the acidic *pI* of bryophytes and monocot plants which exhibit a linear correlation relative to algae and eudicots.

Figure 2

Principal component analysis (PCA) of basic *pI* proteins. The PCA plot illustrates that the basic *pI* of algae, bryophytes, eudicots, and monocot plants cluster distinctly from each other and that there is no lineage-specific correlation with basic *pI* proteins.

Figure 3

Trimodal distribution of isoelectric points (*pI*) and the molecular mass (kDa) of plant proteins. The *pI* of plant proteins ranged from 1.99 (epsin) to 13.96 (hypothetical protein), while the molecular mass ranged from 0.54 (unknown) to 2236.8 (type I polyketide synthase) kDa. The X-axis represents the *pI* and the Y-axis represents the molecular mass of the proteins.

Figure 4

Average amino acid composition of proteins in the plant kingdom. Leu is the most abundant while Trp is the least abundant. The amino acid, Sec, was only found in a few species of algae and was absent from all other species. Ambiguous amino acids were found in a few species as well.

Figure 5

Principal component analysis (PCA) of amino acid abundance in plant proteomes. The PCA plot shows that Tyr, Trp, Cys, His, Met, and Xaa (unknown) amino acids are low-abundance and cluster together. The abundance of Leu, Ser, Ile, Lys, and Gln was higher and grouped together. The plot shows that the abundance of amino acids is lineage specific. Algae, eudicots, and monocot plants exhibit a lineage specific correlation.

Figure 6

Schematic illustration of the biosynthetic pathway of amino acids. The abundance of aromatic ring containing amino acids is lower relative to other amino acids. The average abundance of the aromatic ring containing amino acid, Trp, is the lowest amongst others that are biosynthesized via phosphoenolpyruvate and erythrose 4-phosphate. Similarly, the abundance of Cys is relatively low compared to other amino acids. Ser is biosynthesized from 3-phosphoglycerate and Ser is subsequently used to produce Gly and Cys amino acids. The abundance of Cys is lower relative to Gly, suggesting the existence of allosteric feed-back inhibition of the biosynthesis of Cys by Ser.

Supplementary Data

Supplementary File 1

Supplementary file containing average amino acid composition of plant proteomes. The file is present in excel sheet and individual sheet represents details of different amino acid.

Figure 1

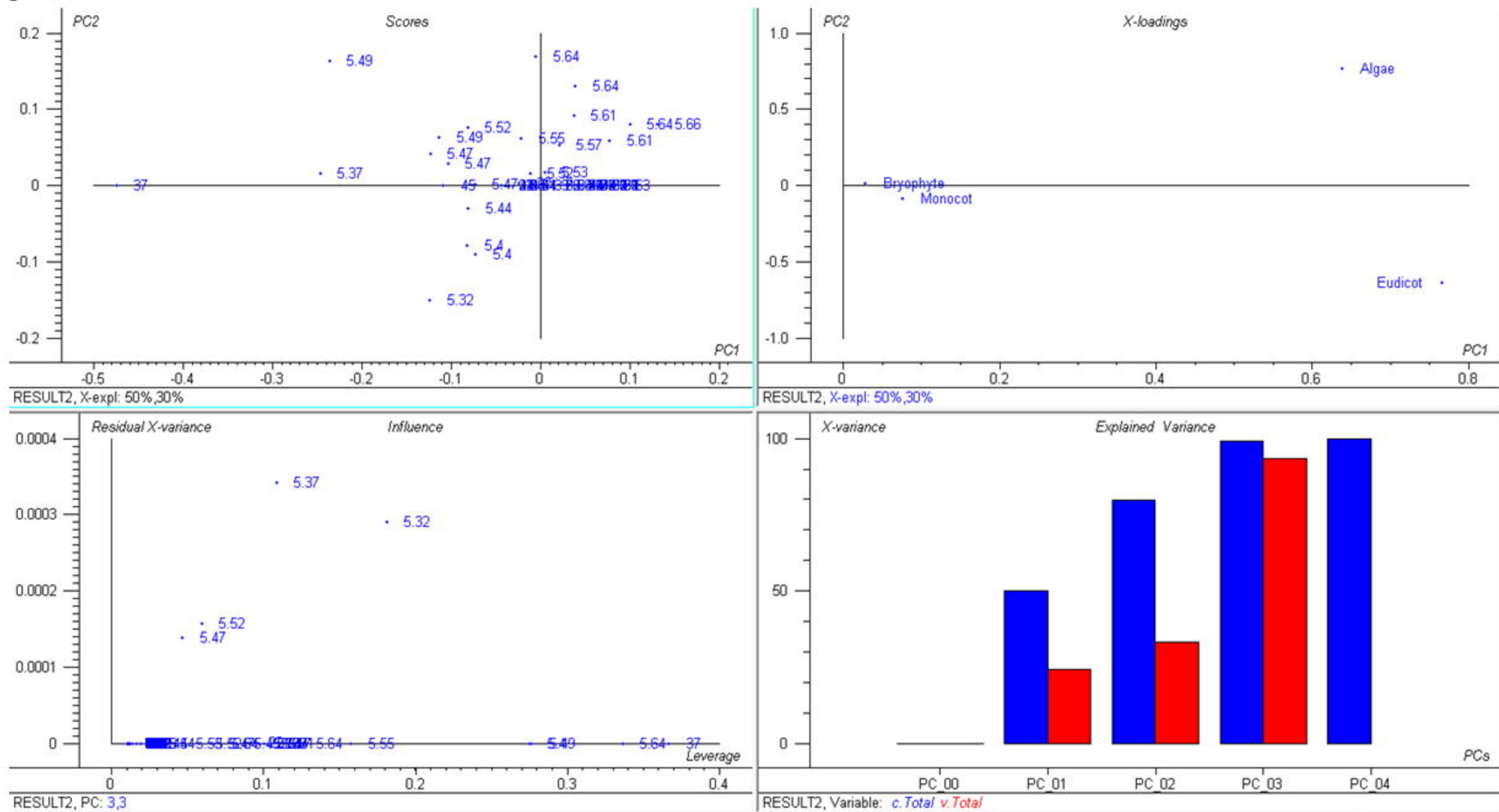


Figure 2

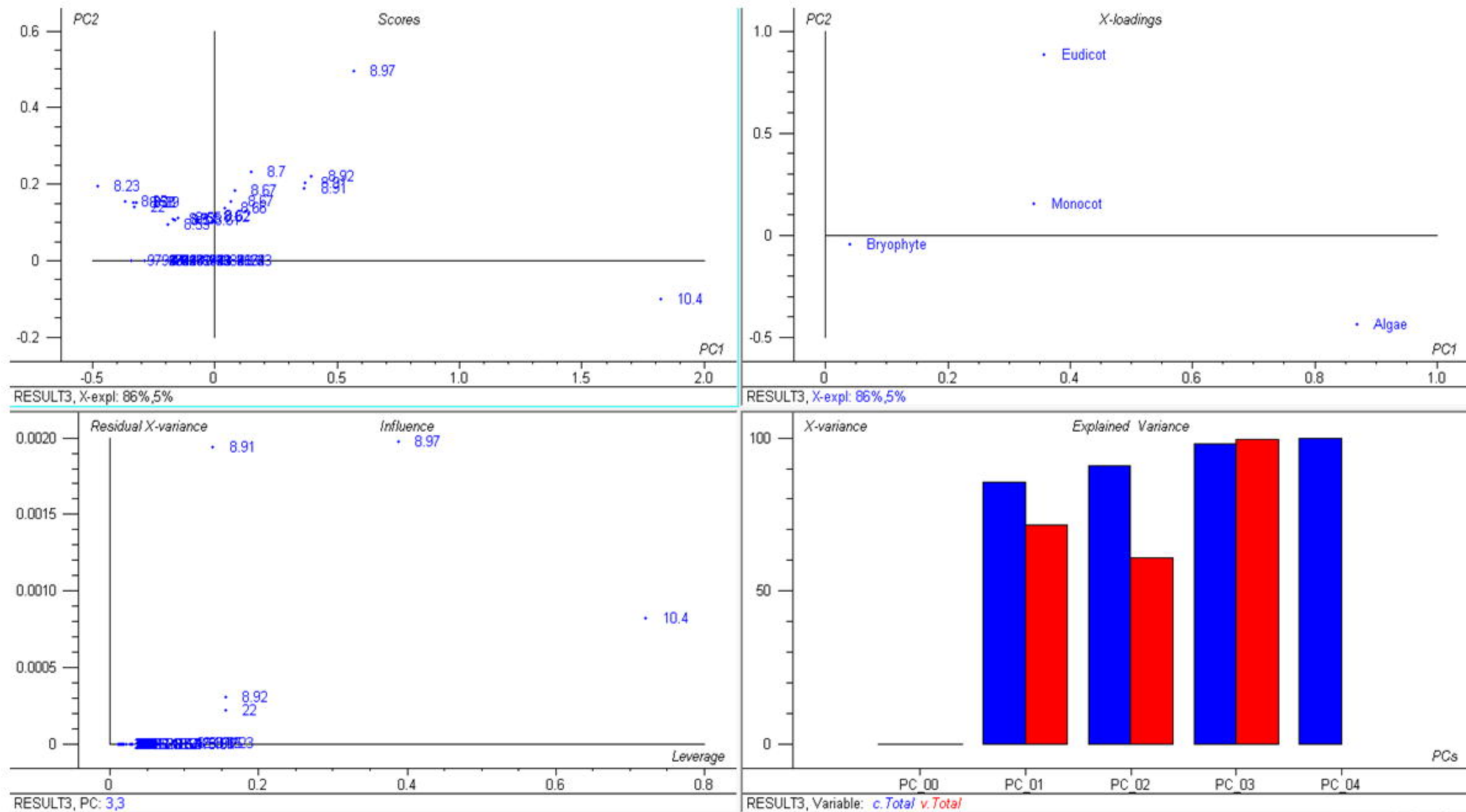


Figure 3

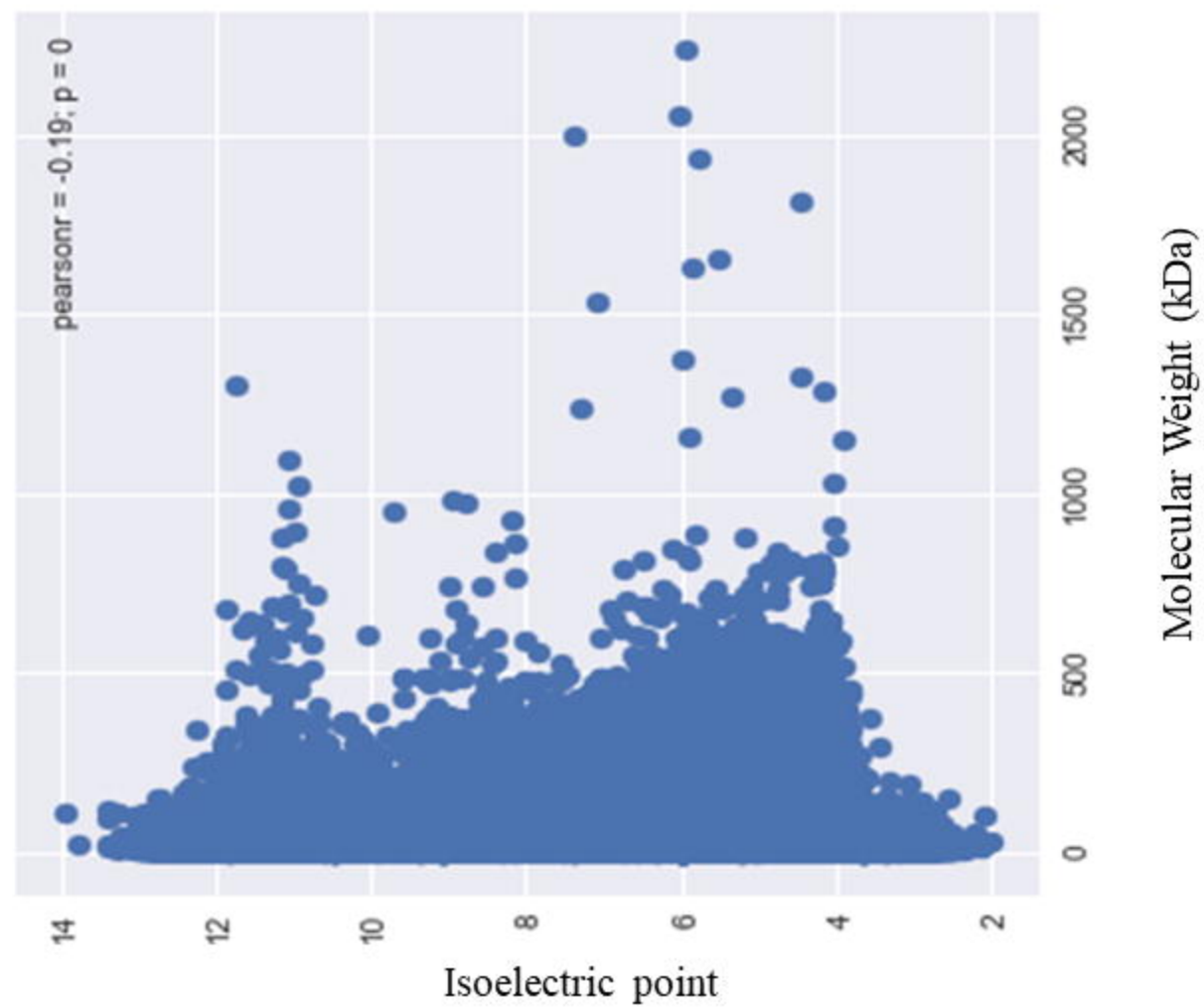


Figure 4

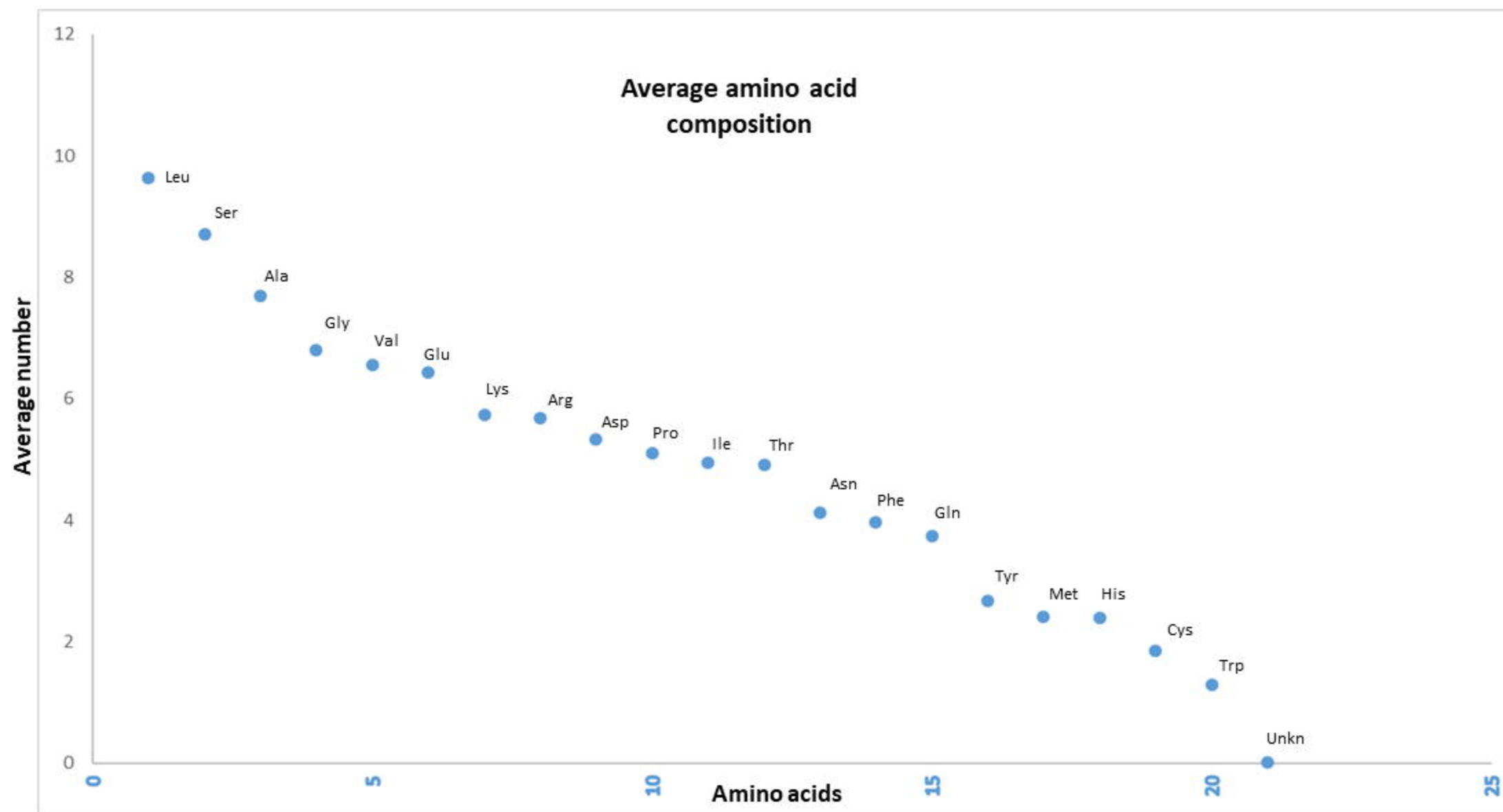


Figure 5

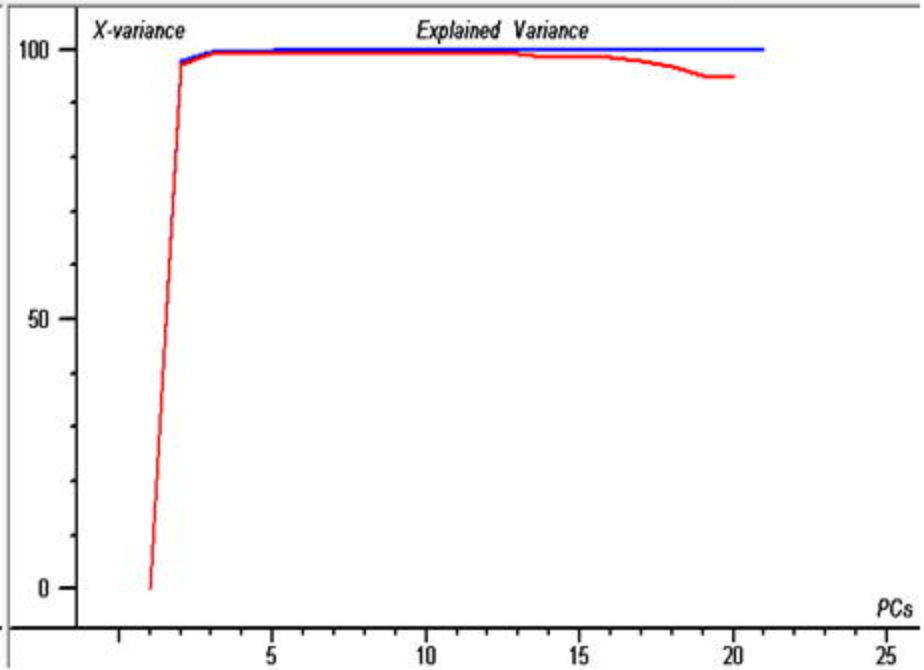
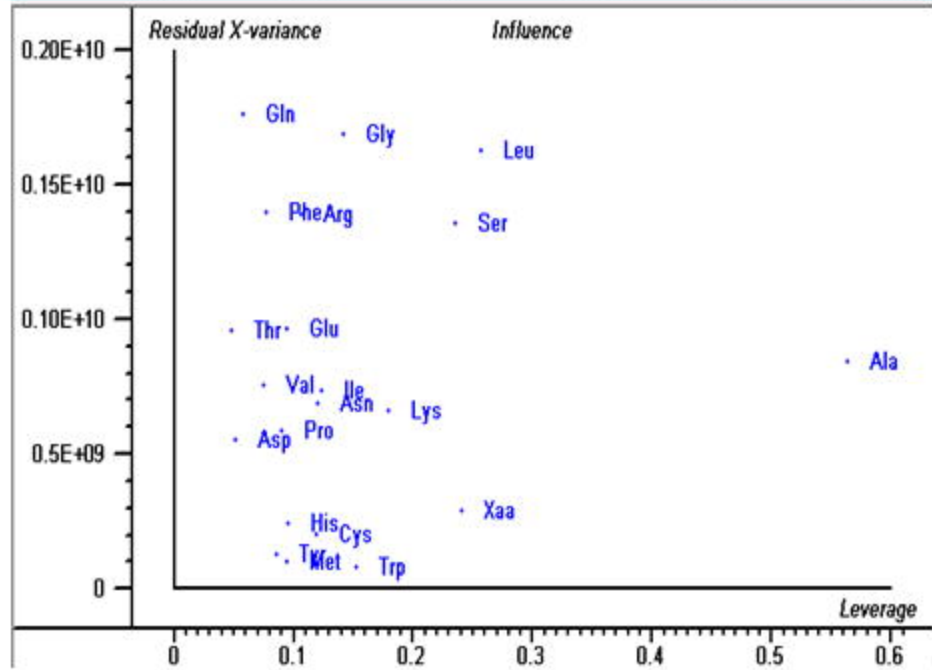
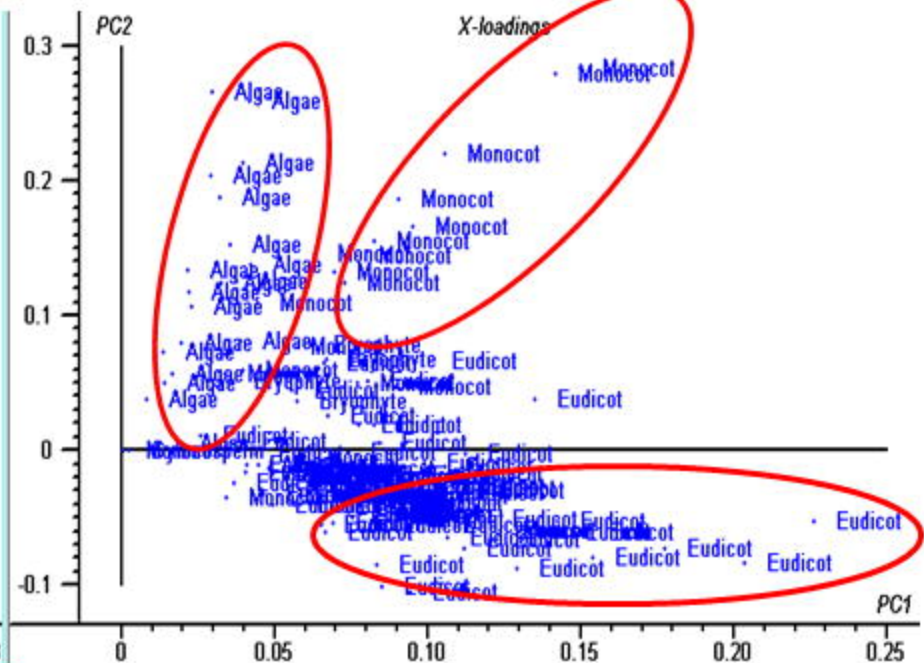
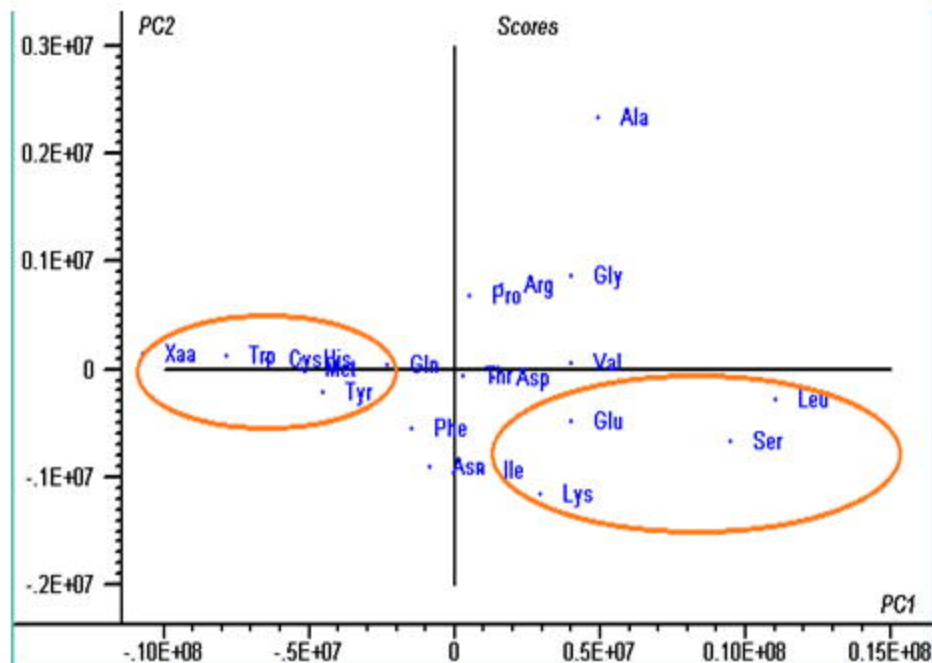


Figure 6

