

## Distinguishing Proteins From Arbitrary Amino Acid Sequences

**Authors:** Stephen S.-T. Yau<sup>1\*</sup>, Wei-Guang Mao<sup>1\*\*</sup>, Max Benson<sup>2</sup>, Rong Lucy He<sup>3</sup>

### **Affiliations:**

<sup>1</sup> Department of Mathematical Sciences, Tsinghua University, Beijing, 100084, China.

<sup>2</sup> Department of Computer Science, Seattle Pacific University, Seattle, WA 98119, USA.

<sup>3</sup> Department of Biological Sciences, Chicago State University, Chicago, IL 60628, USA.

\*To whom correspondence should be addressed. Email: [yau@uic.edu](mailto:yau@uic.edu) (SSTY).

\*\*Co-first author

### **Supplementary Information:**

#### **Databases.**

You can get the datasets we used via the links below:

Uniprot release 2013\_03\_Complete proteome\_Reviewed\_Normalized:

[http://r720.math.tsinghua.edu.cn/Data/proteome/Uniprot release 2013\\_03\\_Complete proteome\\_Reviewed\\_Normalized.fasta](http://r720.math.tsinghua.edu.cn/Data/proteome/Uniprot%20release%202013_03_Complete%20proteome_Reviewed_Normalized.fasta)

Uniprot release 2014\_03\_Complete proteome\_Reviewed\_Normalized:

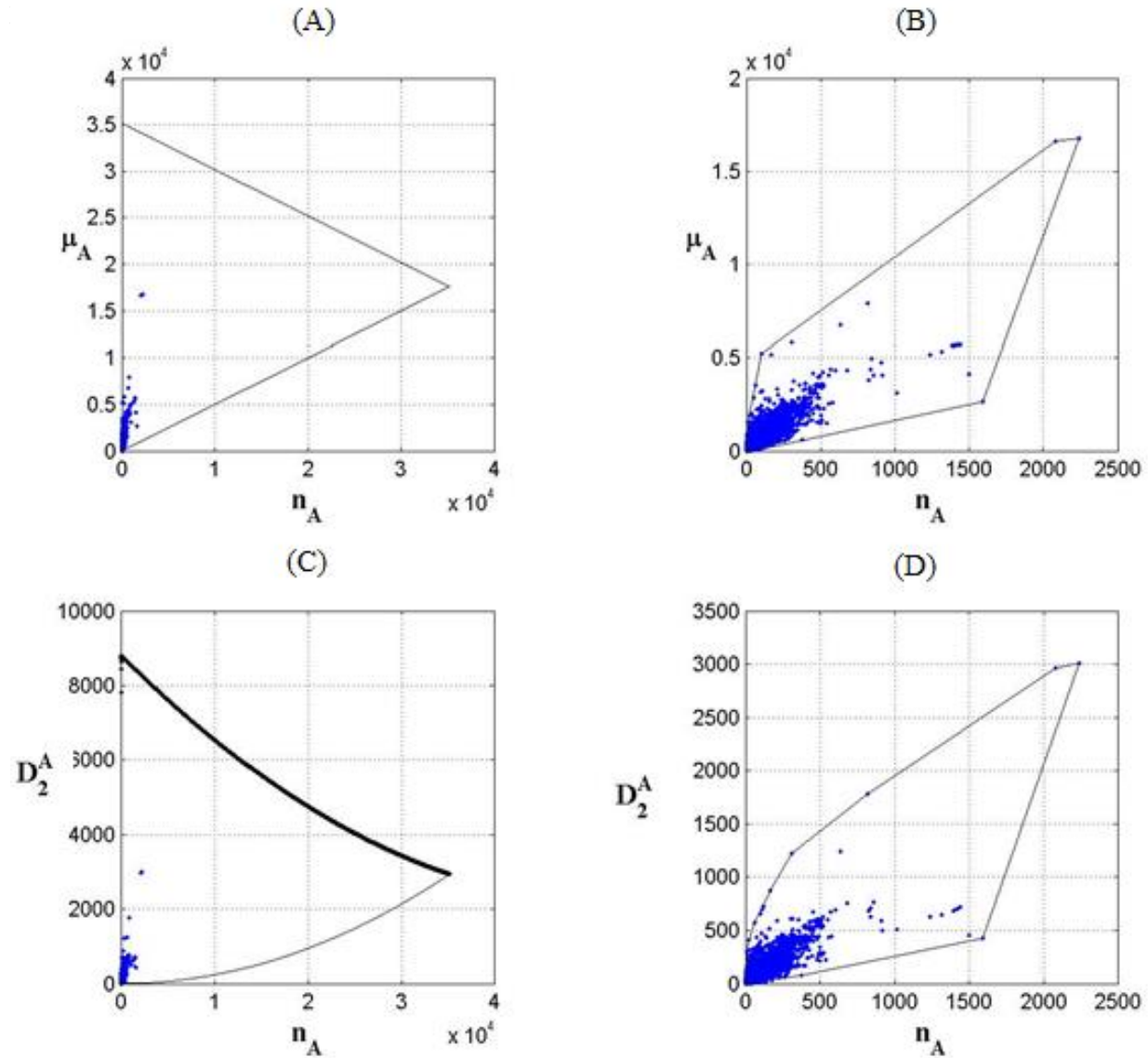
[http://r720.math.tsinghua.edu.cn/Data/proteome/Uniprot release 2014\\_03\\_Complete proteome\\_Reviewed\\_Normalized.fasta](http://r720.math.tsinghua.edu.cn/Data/proteome/Uniprot%20release%202014_03_Complete%20proteome_Reviewed_Normalized.fasta)

Uniprot release 2014\_06\_Complete proteome\_Reviewed\_Normalized.fasta

[http://r720.math.tsinghua.edu.cn/Data/proteome/Uniprot release 2014\\_06\\_Complete proteome\\_Reviewed\\_Normalized.fasta](http://r720.math.tsinghua.edu.cn/Data/proteome/Uniprot%20release%202014_06_Complete%20proteome_Reviewed_Normalized.fasta)

**Qhull Algorithm.** We used the `convhulln()` function supplied with MATLAB to compute all convex hulls. This function is based on Qhull algorithm. We used a publicly available MATLAB function called `inhull()` to test whether or not a given point was in the convex hull. The code for this function can be accessed through the link: <http://www.mathworks.com/matlabcentral/fileexchange/10226-inhull>

**Additional Figures.** The following figures show the protein areas for each of the twenty amino acids. Figure 1 below shows the same information as Figure 1 in the main text, while the other figures show the protein areas for the other nineteen amino acids. The figure were generated these figures as described in the main text, using the dataset *Uniprot 2013\_03*. Blue points in each of these four subfigures stand for natural vectors corresponding to proteins. From left to right and top to bottom, the four subfigures are named (A), (B), (C) and (D). (A) shows the picture in  $(n_k, \mu_k)$  coordinate and (C) shows the picture  $(n_k, D_2^k)$  coordinate. (B) is the enlarged view of the protein area in (A). The black lines stand for the boundaries of the convex hull for protein area. (D) is the enlarged view of the protein area in (C). The black lines stand for the boundaries of the convex hull for protein area.



**Figure 1: Protein Area Detail for Alanine**

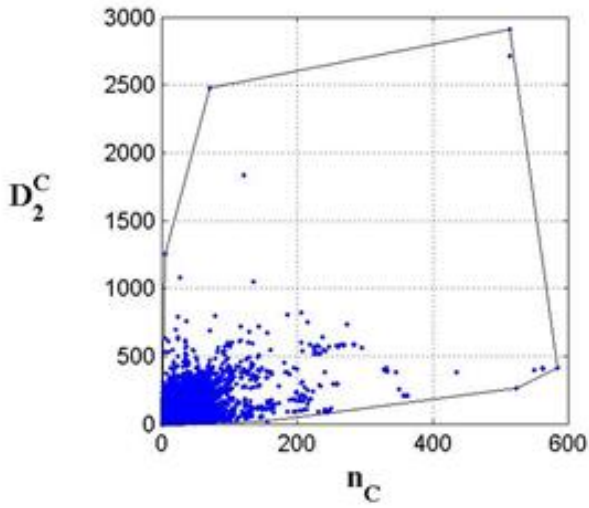
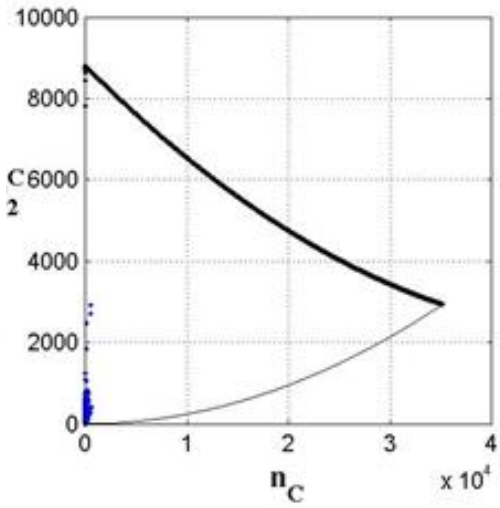
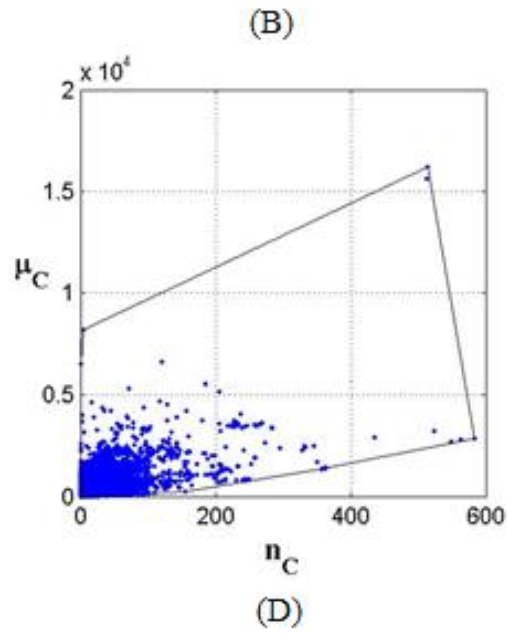
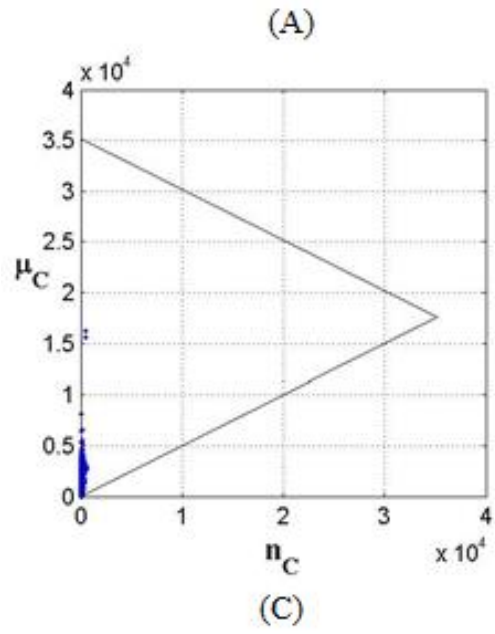


Figure 2: Protein Area Detail for Cysteine

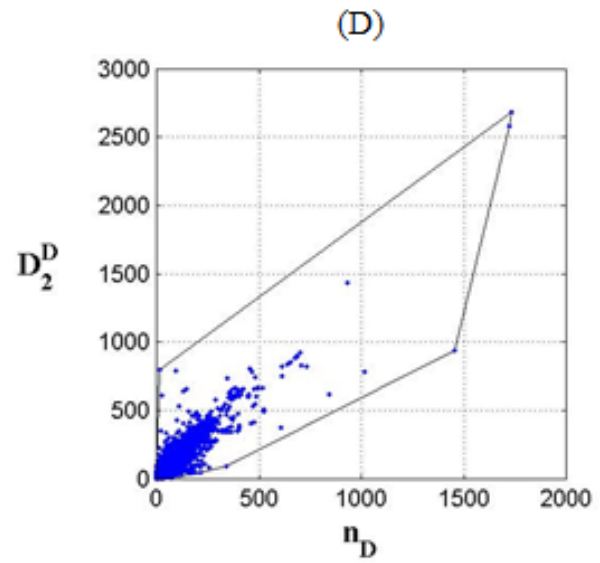
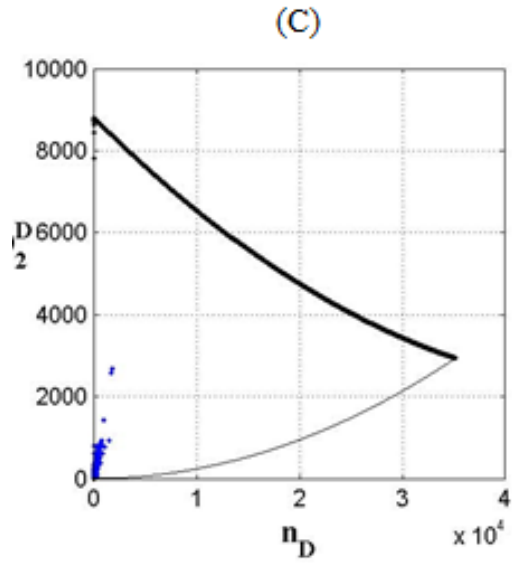
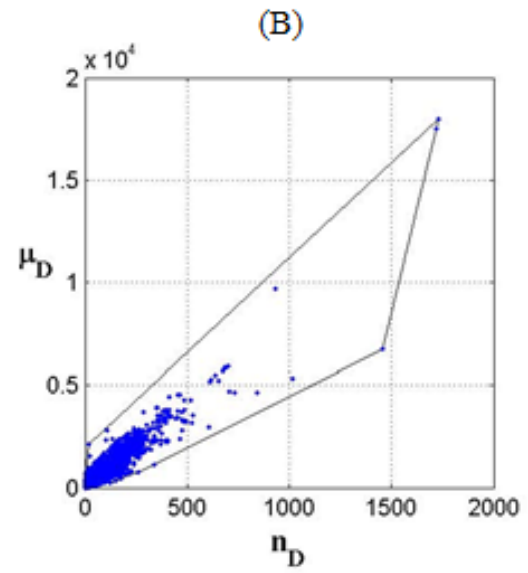
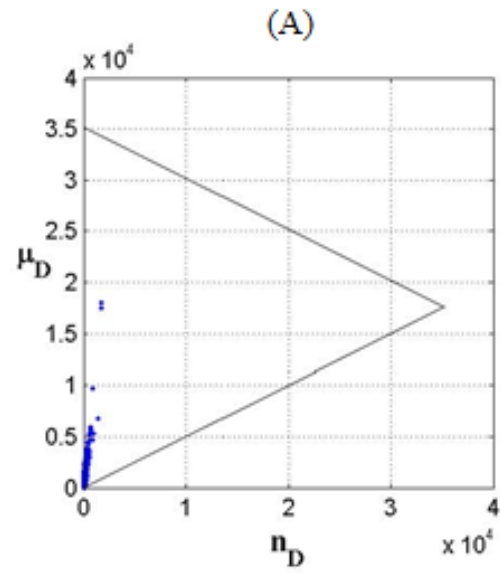


Figure 3: Protein Area Detail for Aspartic Acid

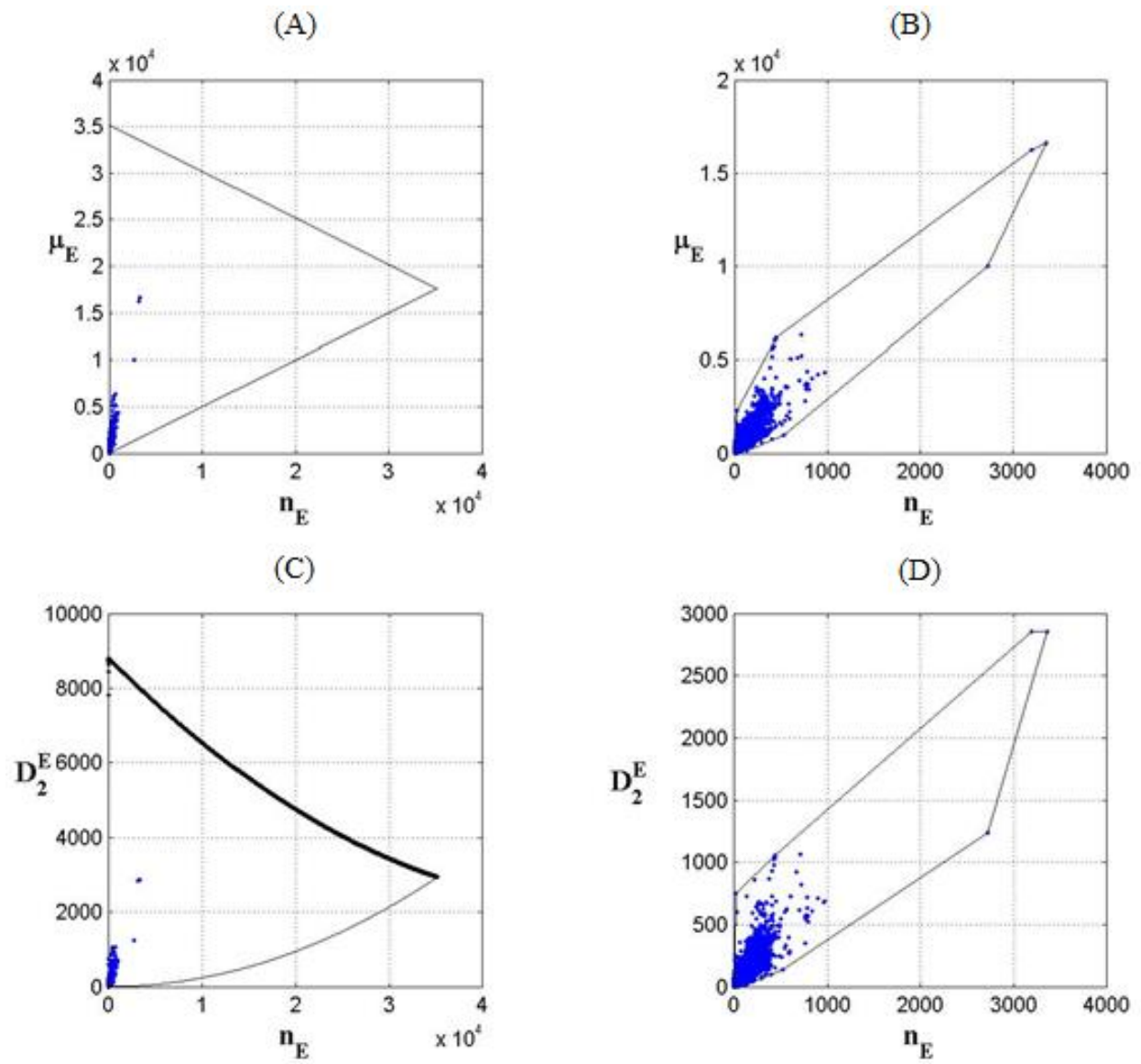


Figure 4: Protein Area Detail for Glutamic Acid

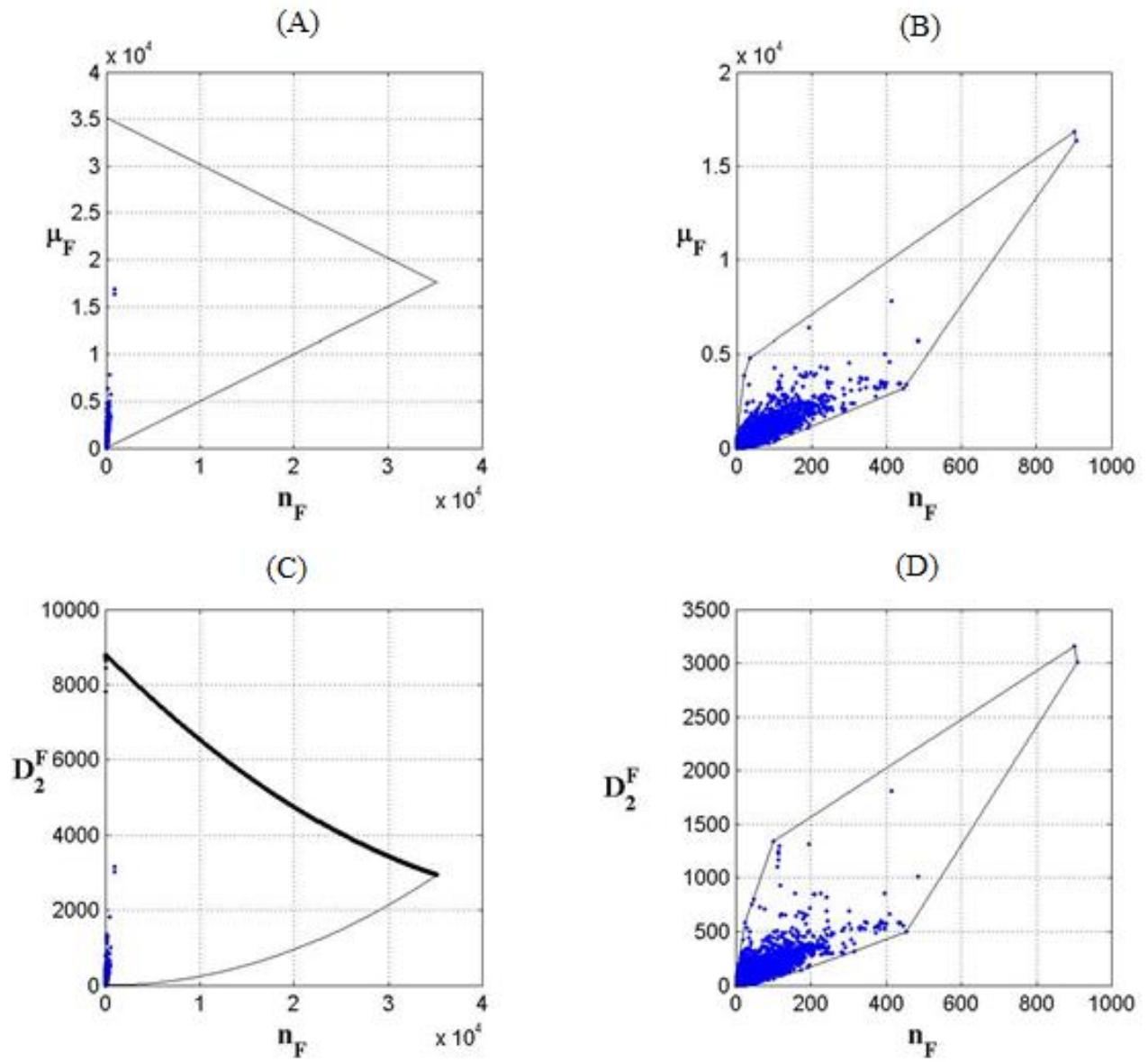


Figure 5: Protein Area Detail for Phenylalanine

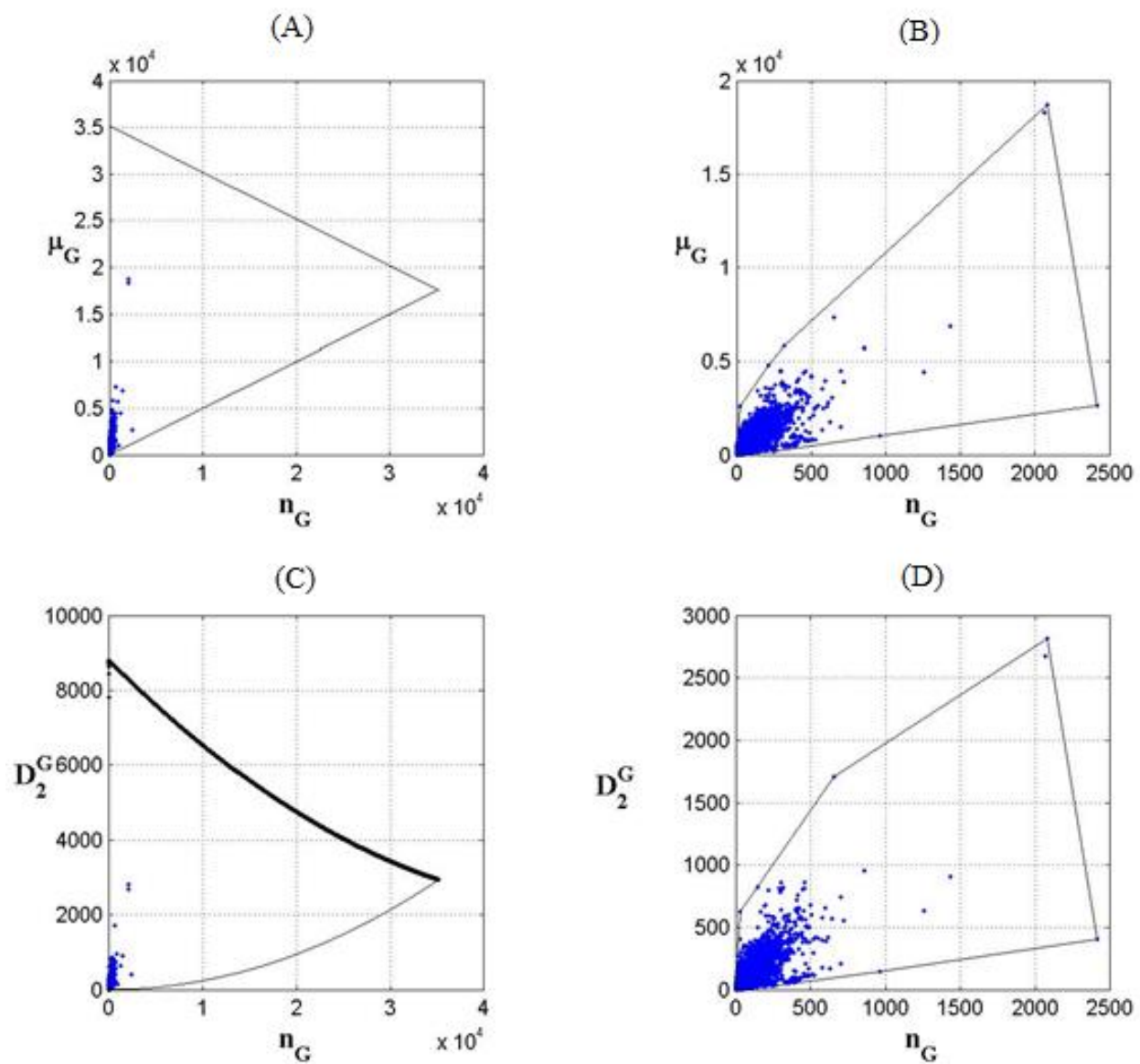
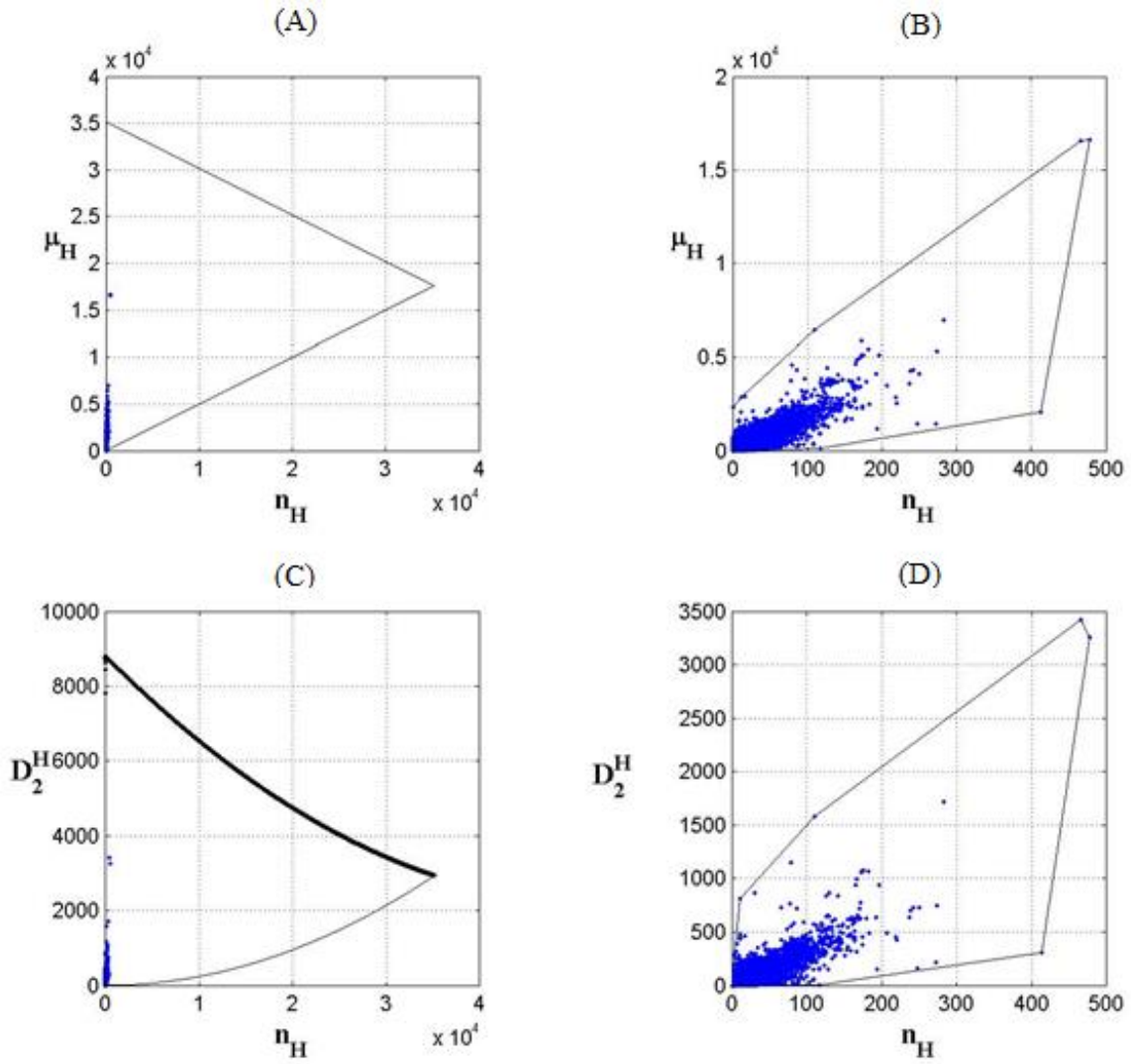
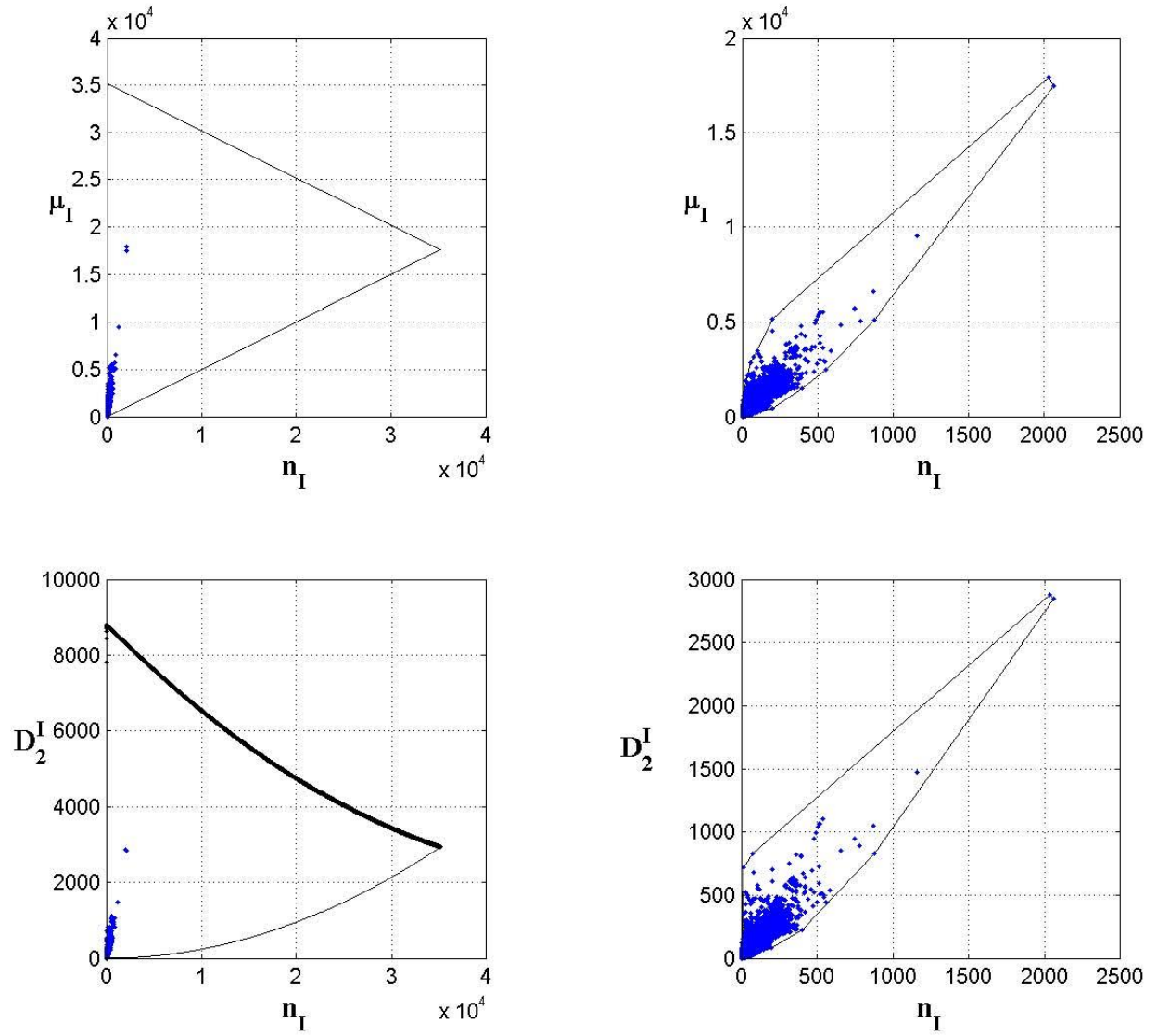


Figure 6: Protein Area Detail for Glycine



**Figure 7: Protein Area Detail for Histidine**





**Figure 8: Protein Area Detail for Isoleucine**

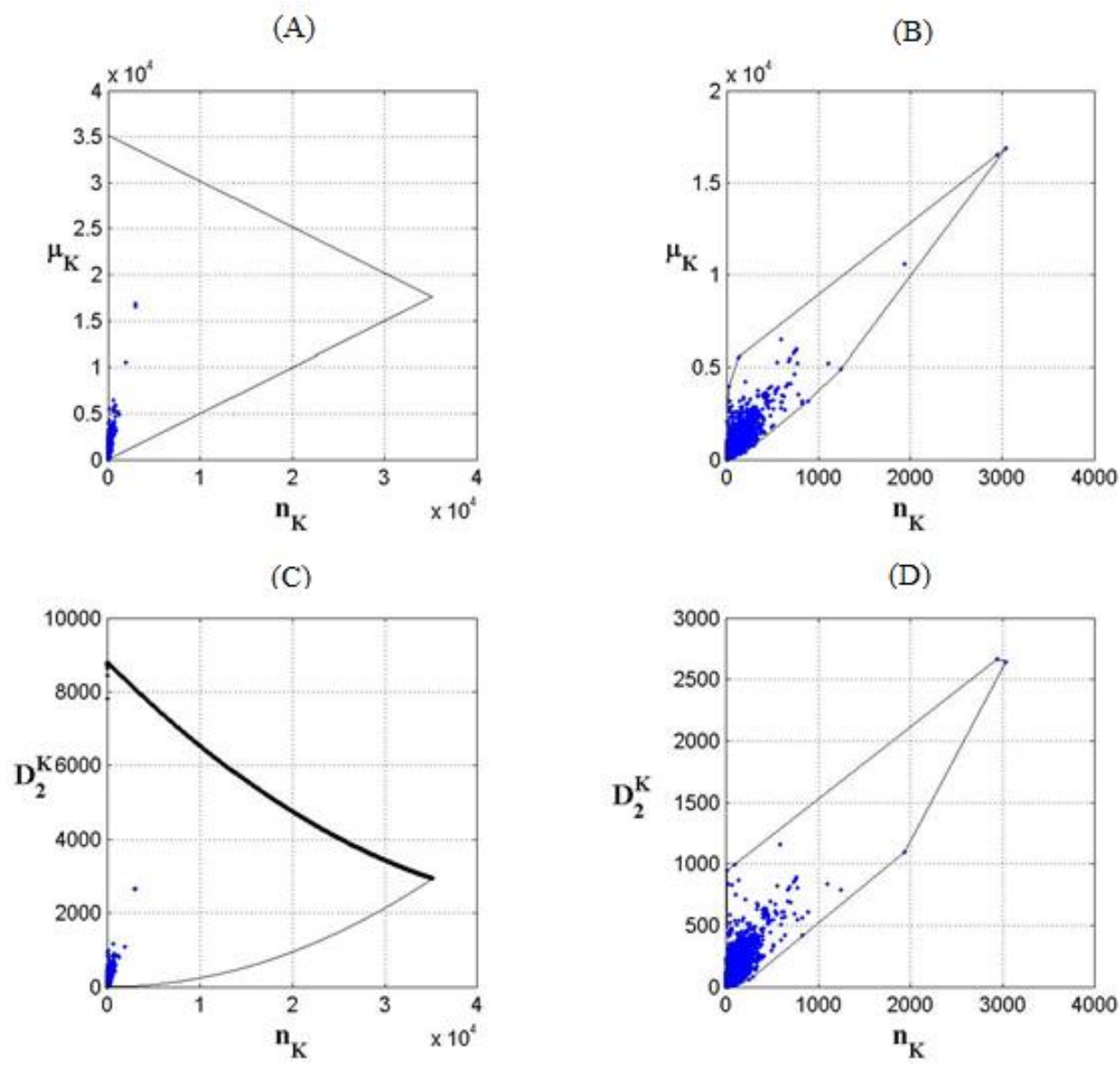


Figure 9: Protein Area Detail for Lysine

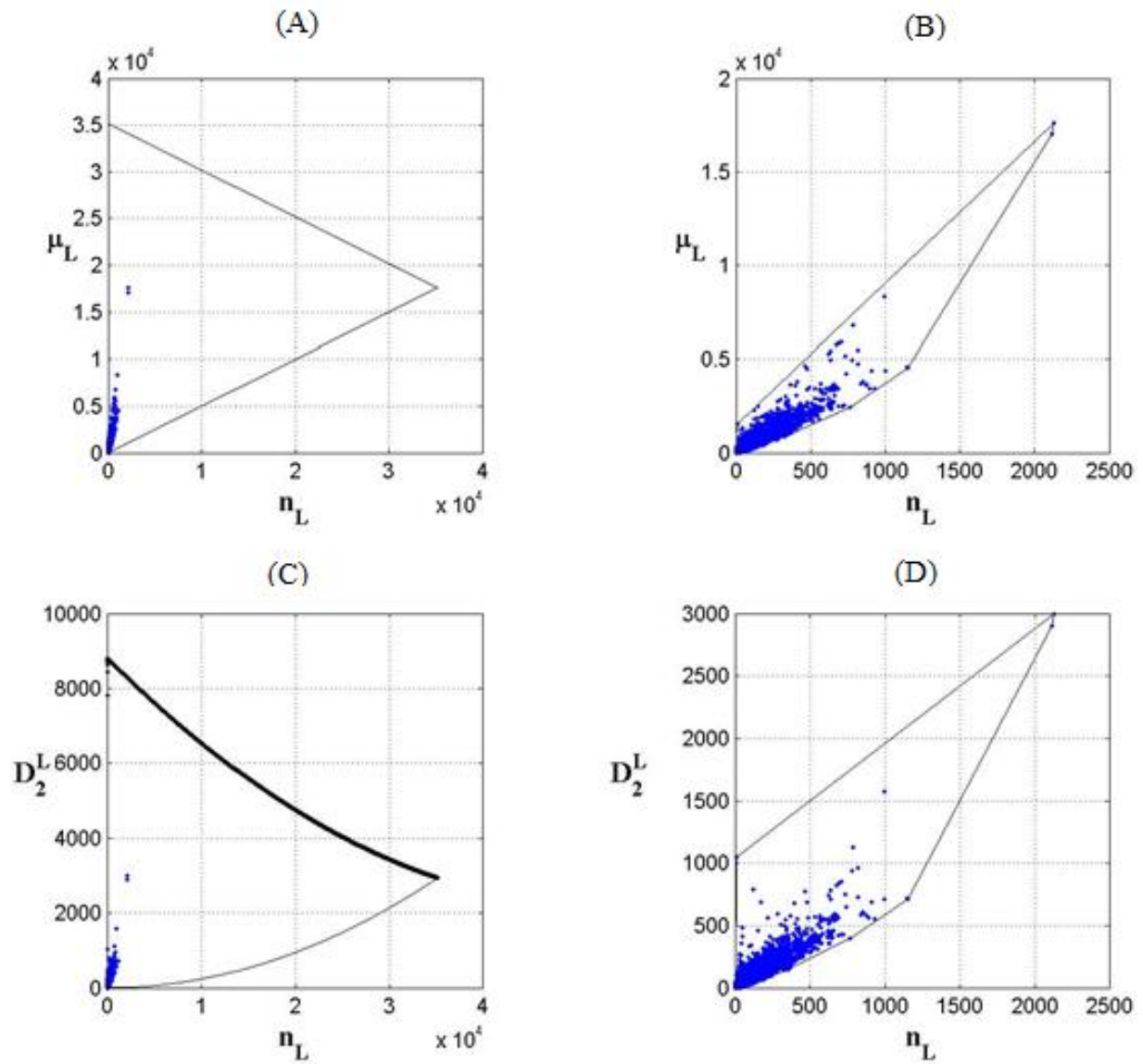


Figure 10: Protein Area Detail for Leucine

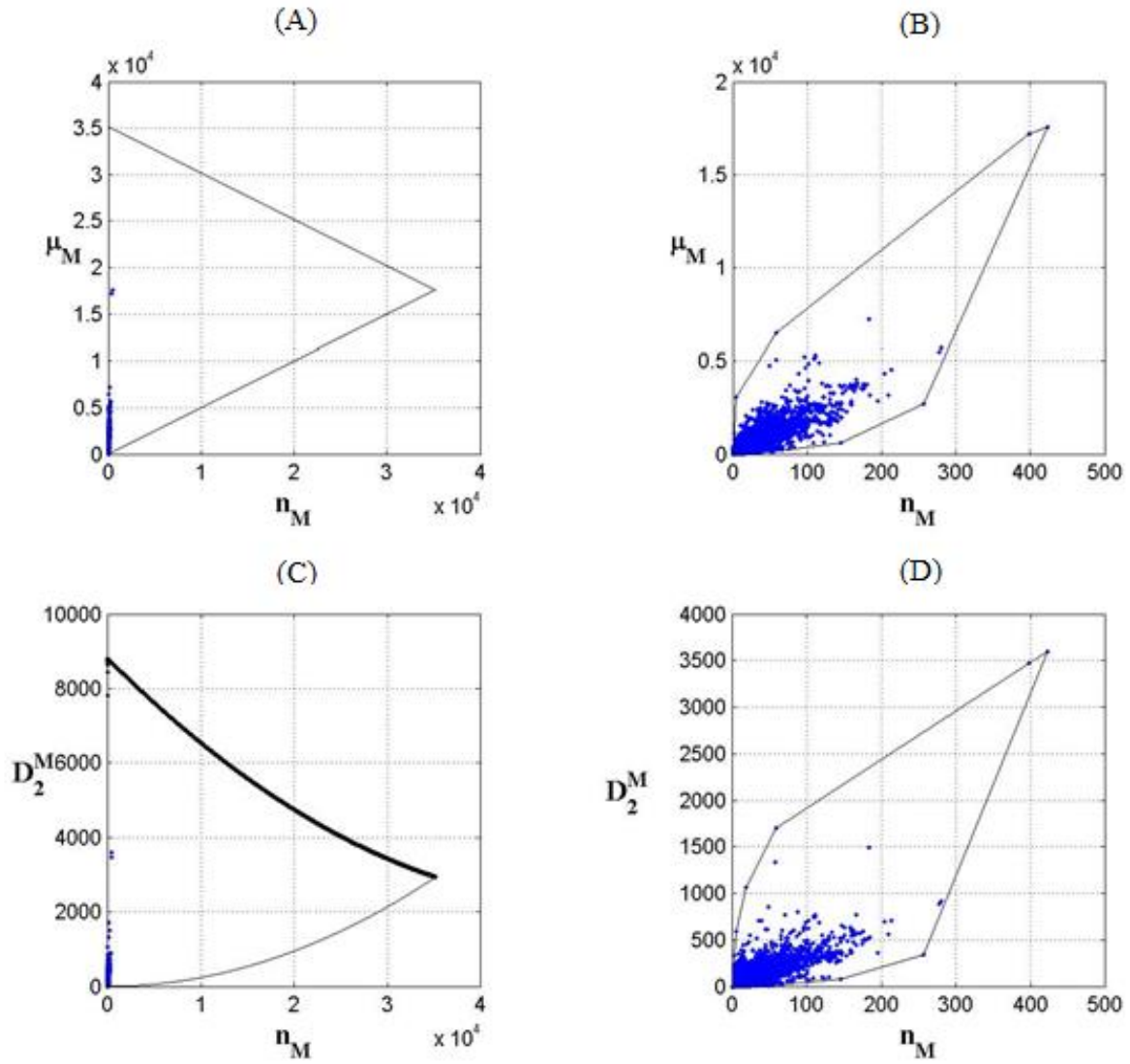


Figure 11: Protein Area Detail for Methionine

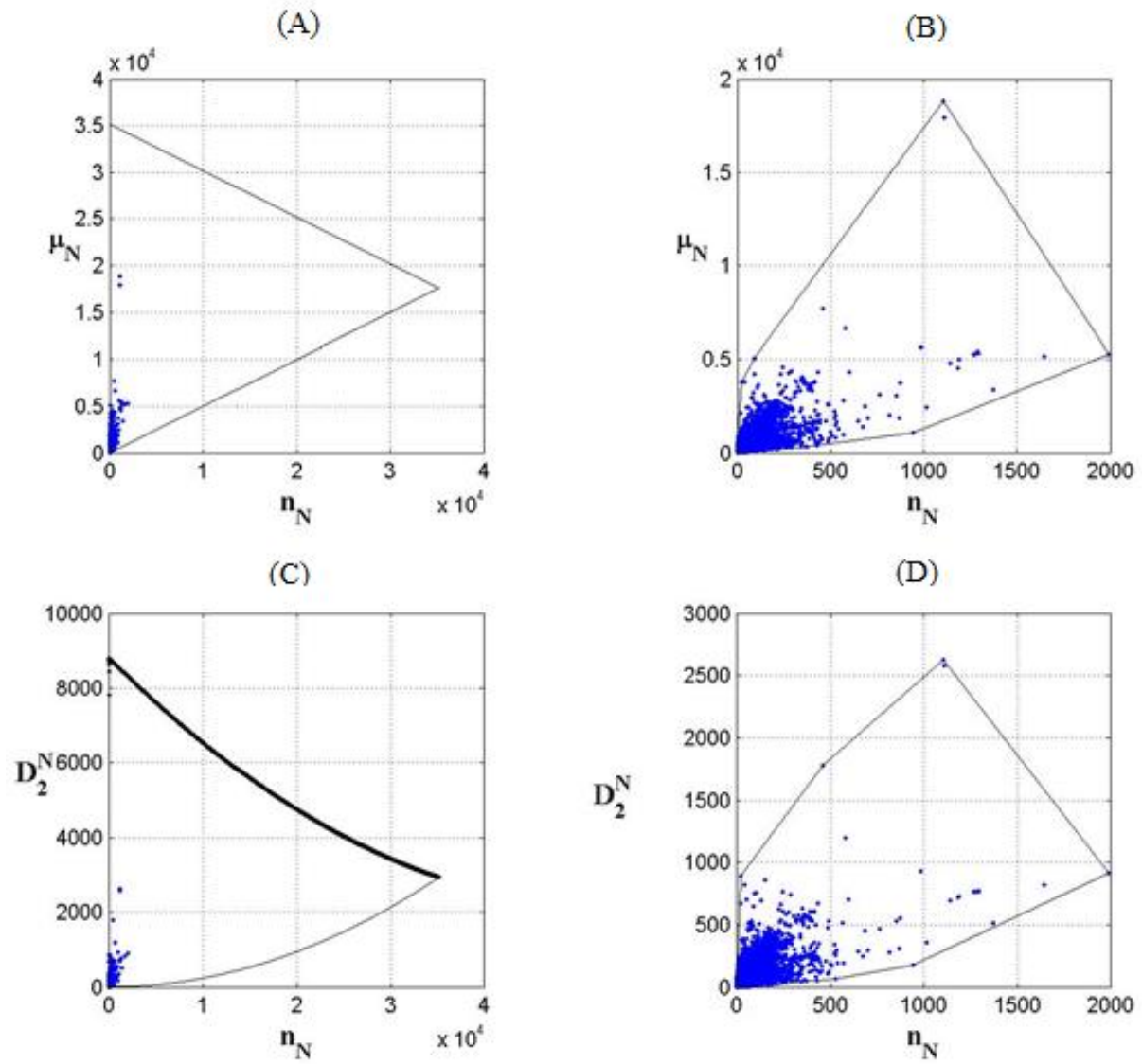


Figure 12: Protein Area Detail for Asparagine

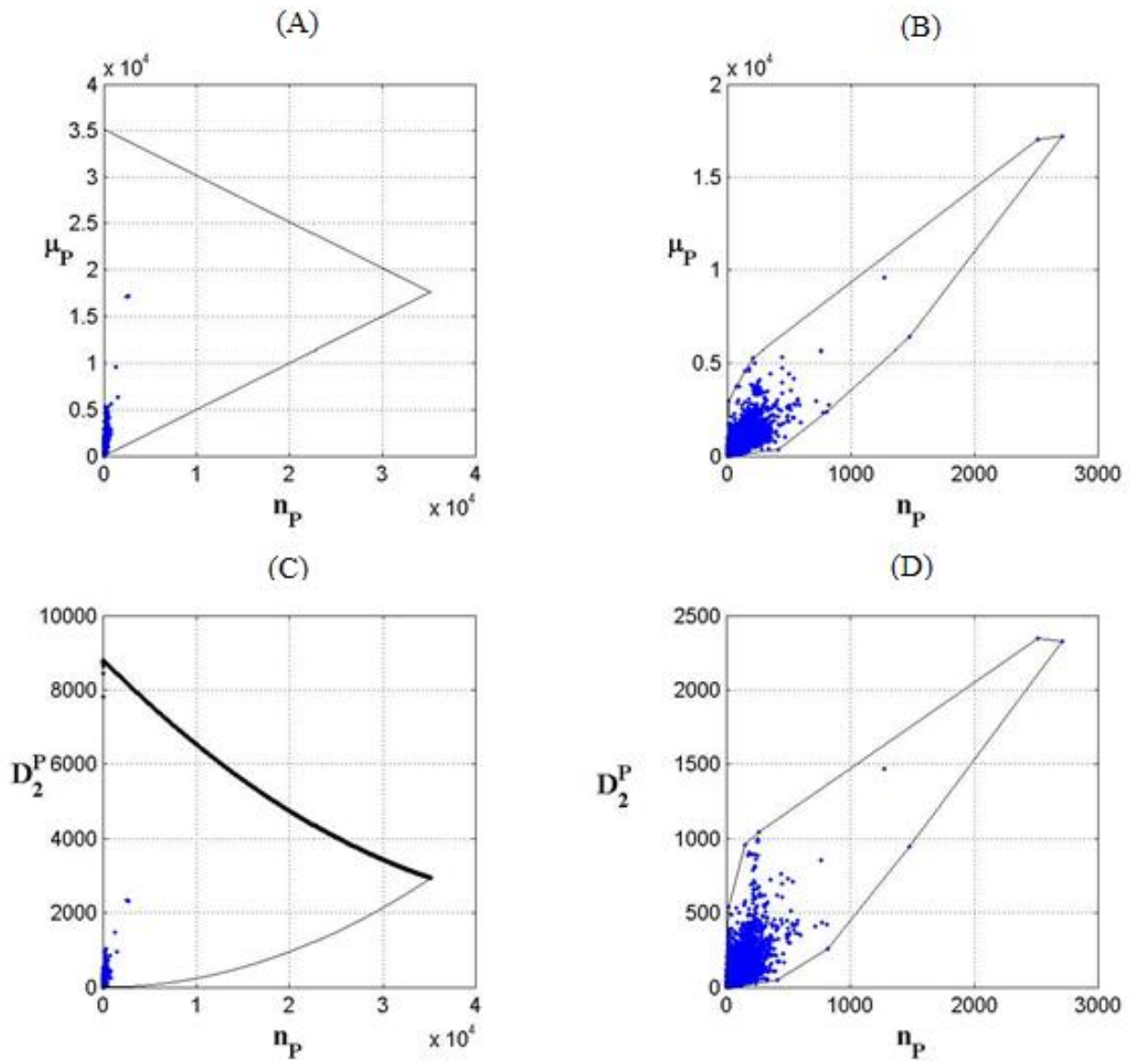


Figure 13: Protein Area Detail for Proline

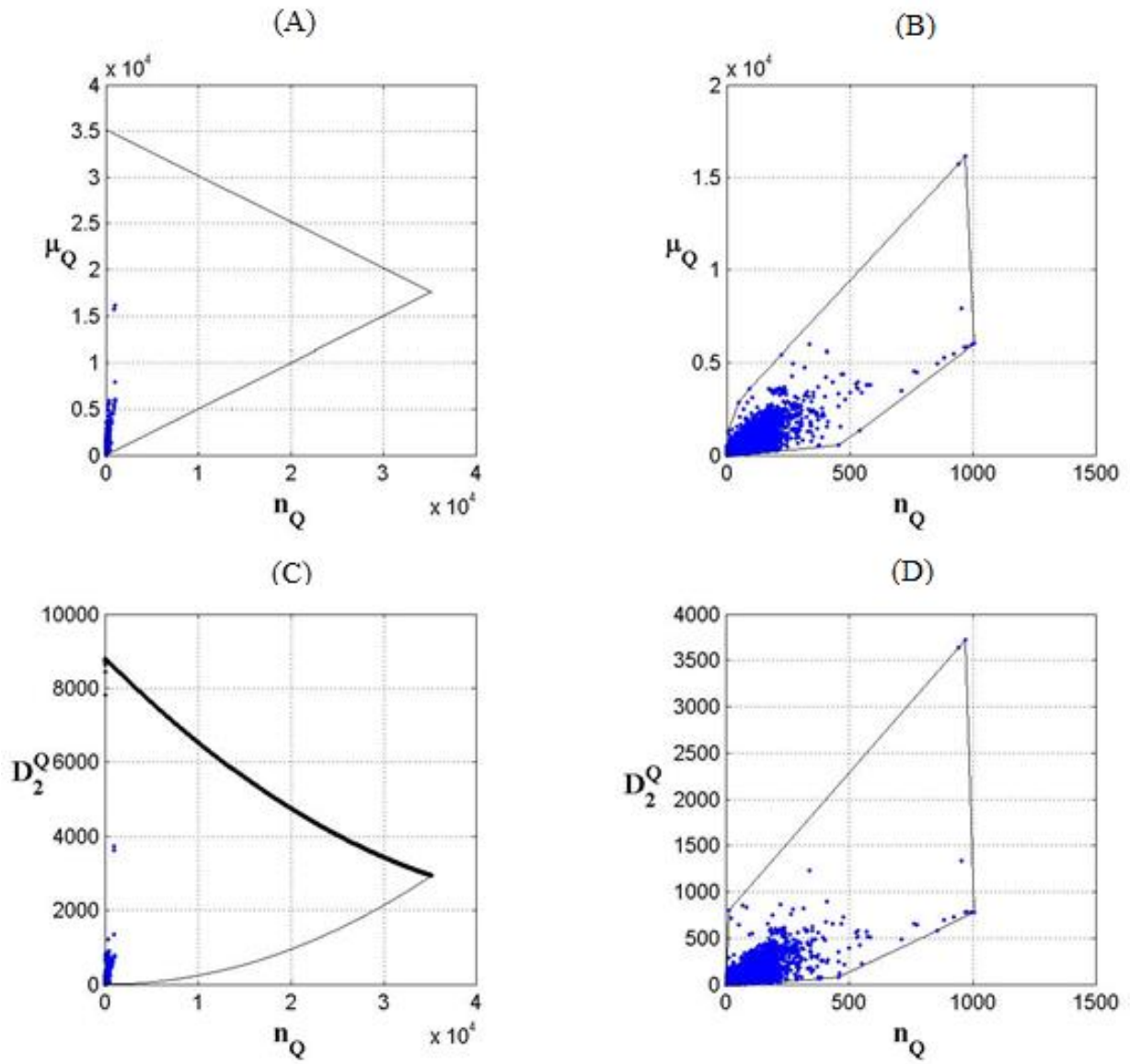


Figure 14: Protein Area Detail for Glutamine

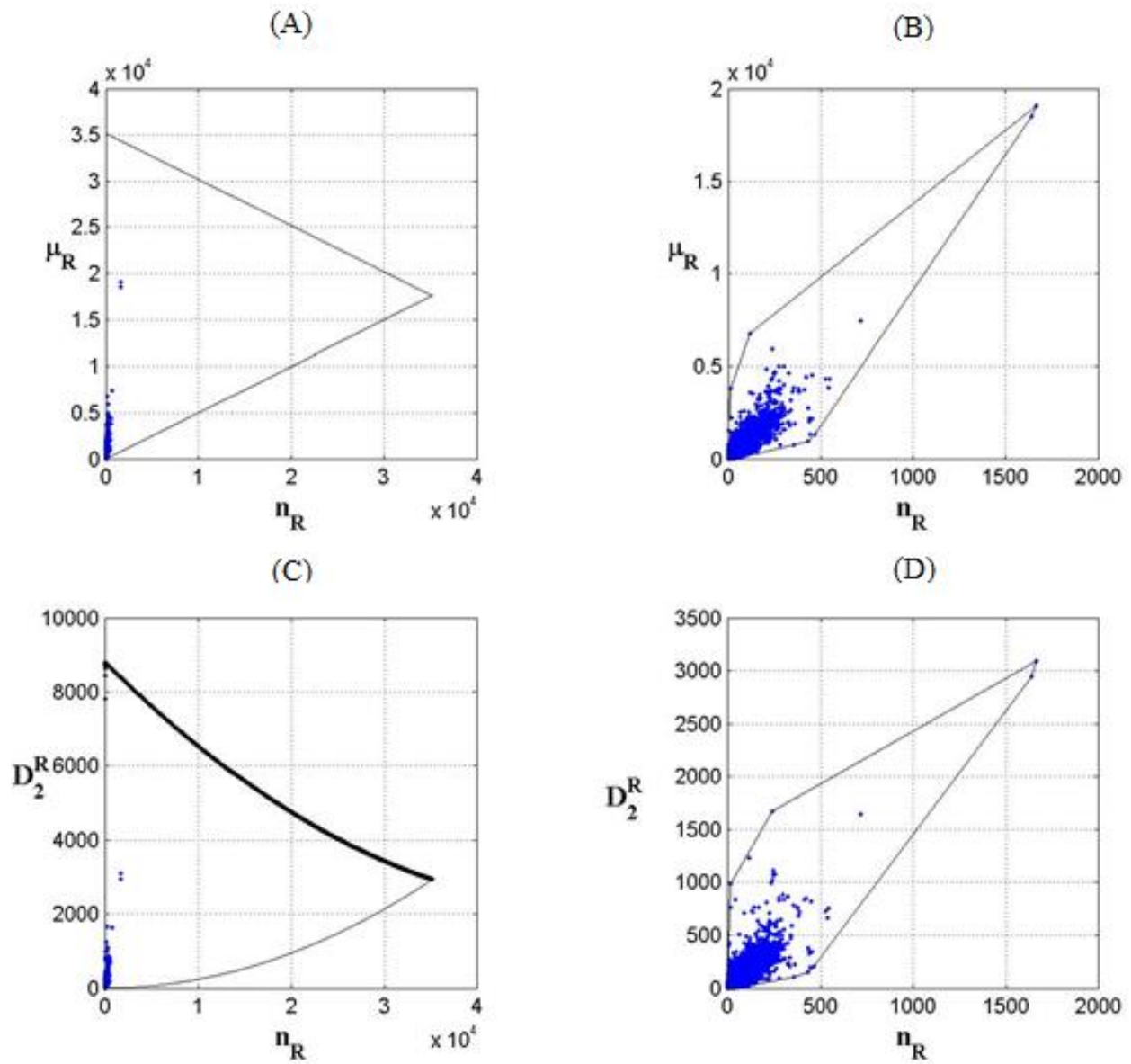


Figure 15: Protein Area Detail for Arginine



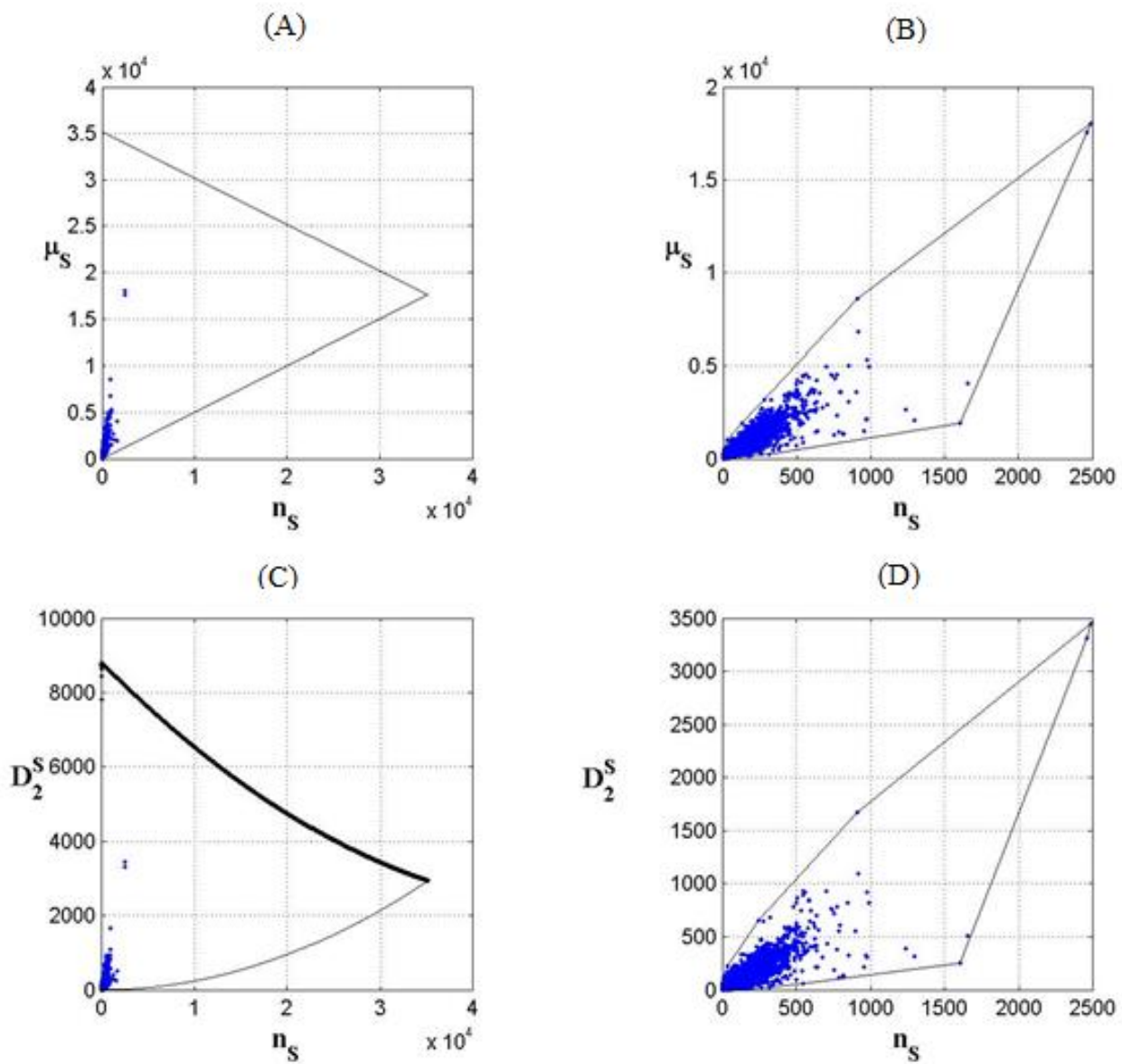


Figure 16: Protein Area Detail for Serine

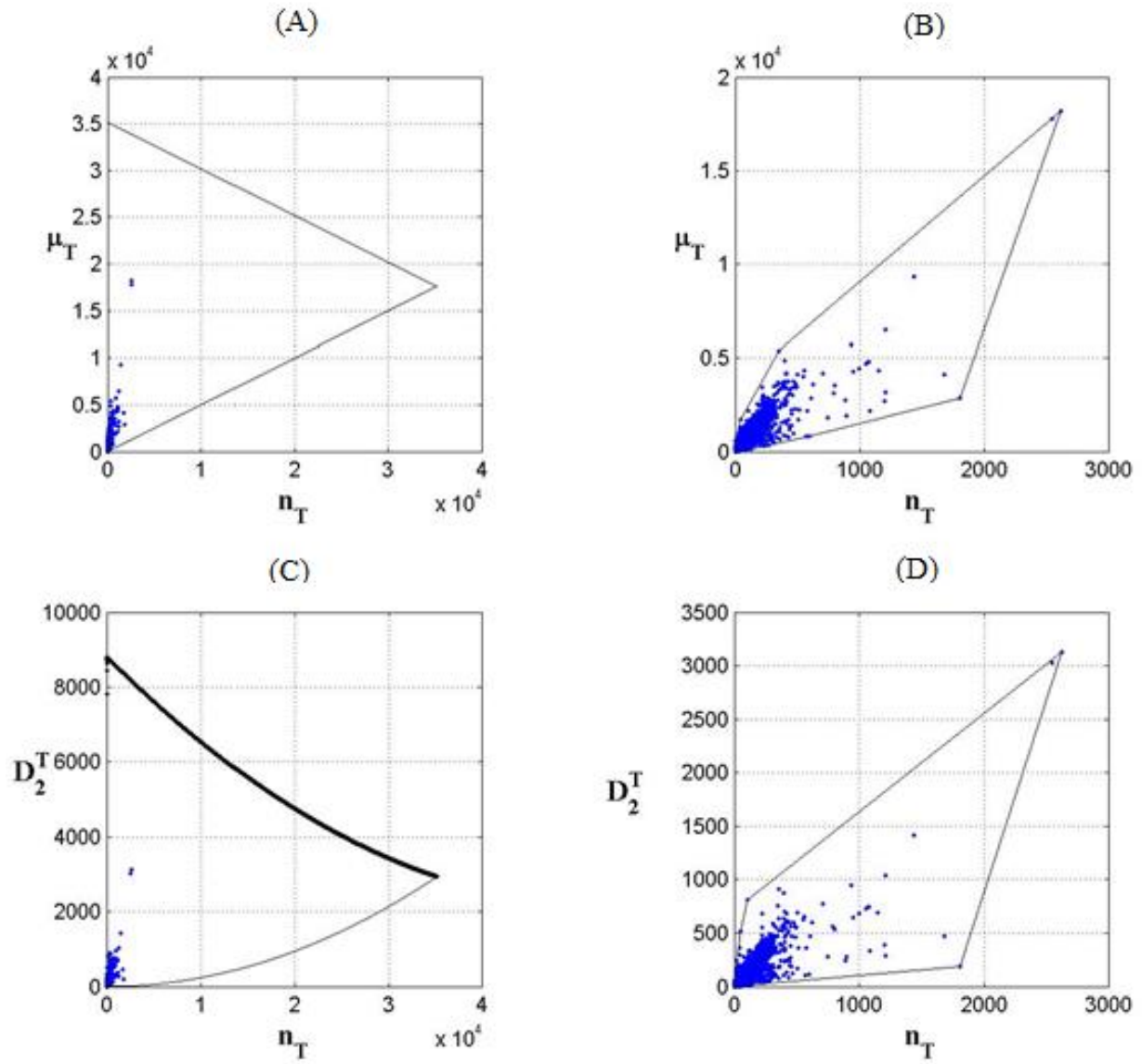


Figure 17: Protein Area Detail for Threonine

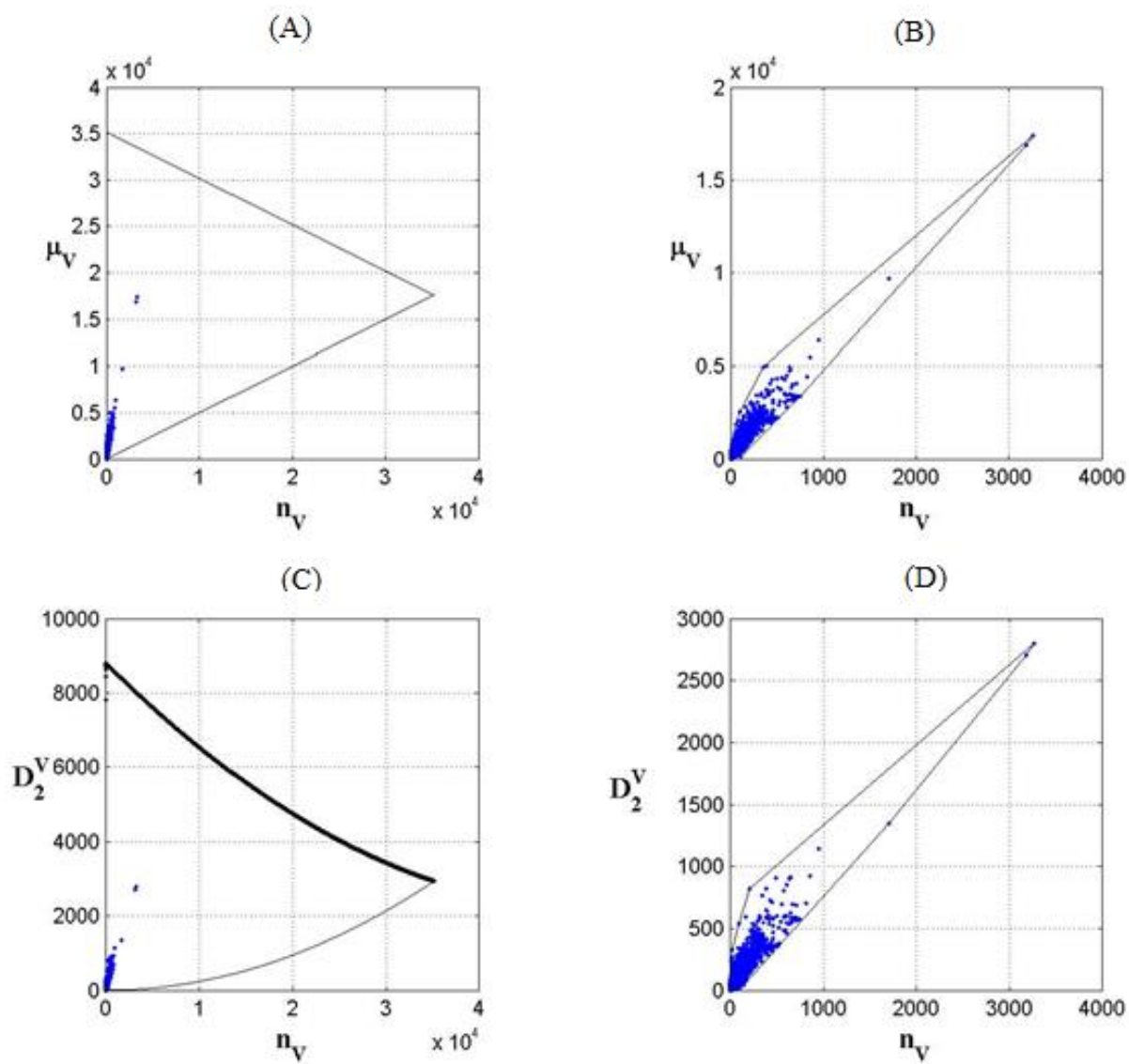


Figure 18: Protein Area Detail for Valine

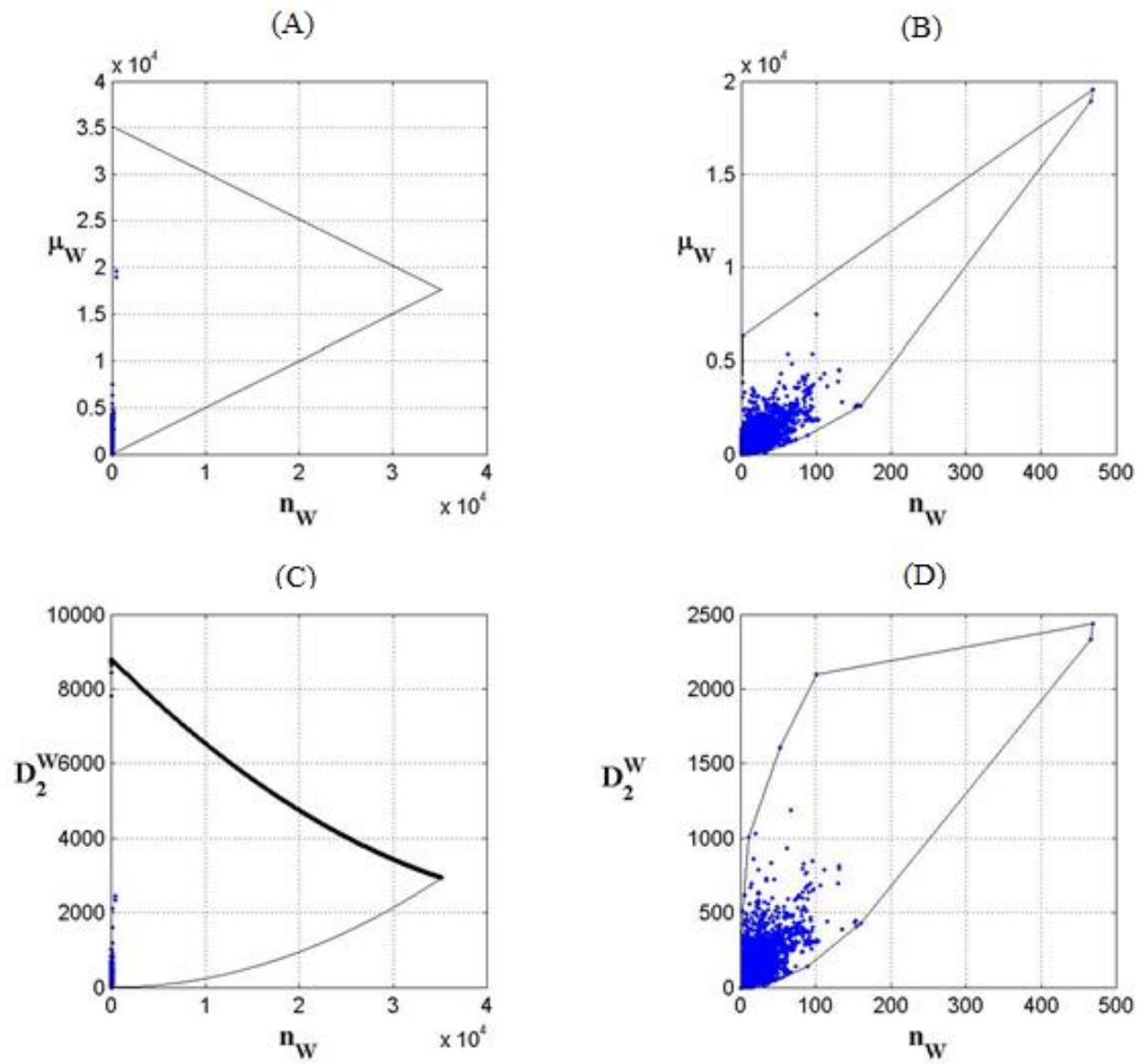


Figure 19: Protein Area Detail for Tryptophan

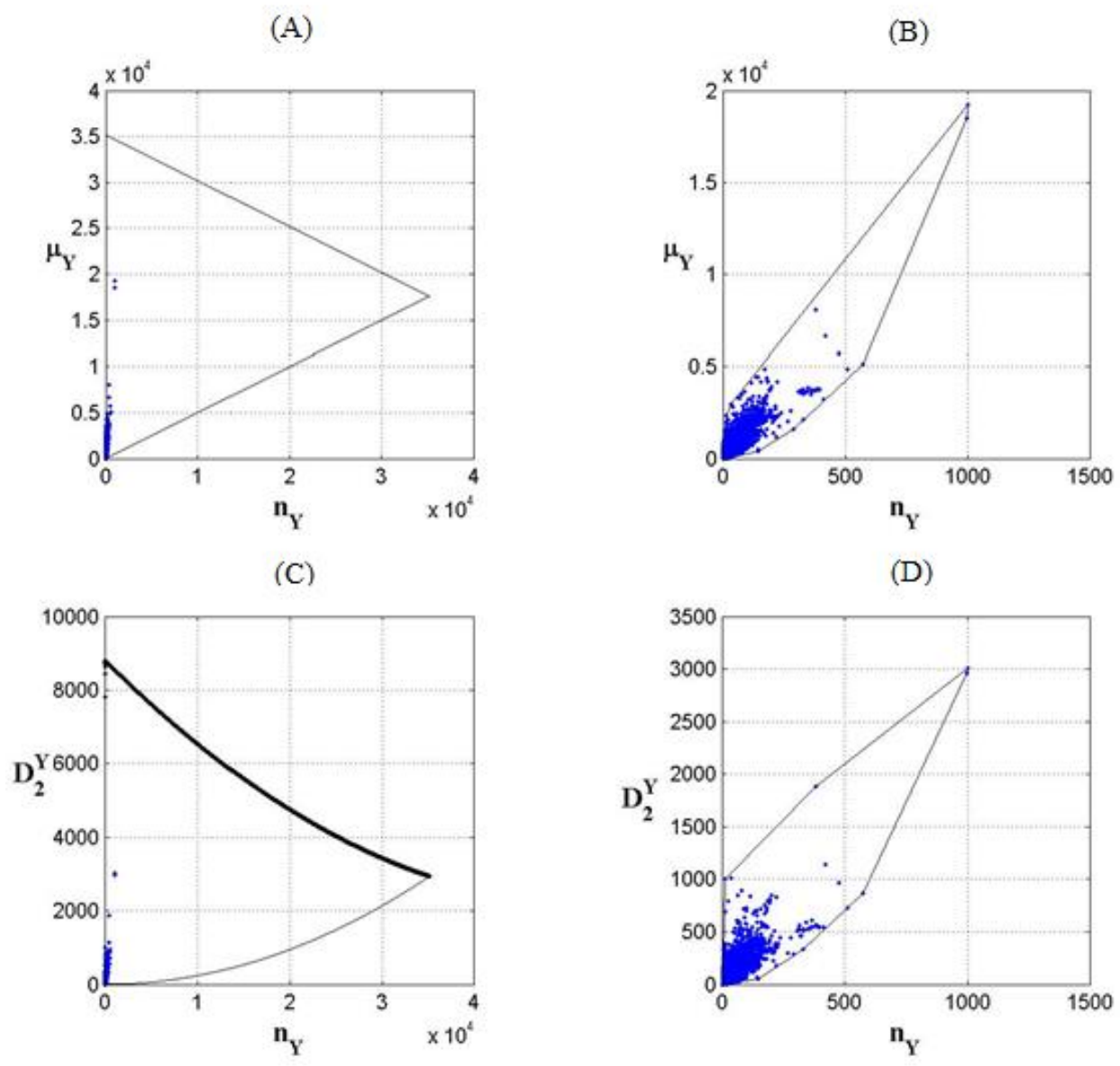


Figure 20: Protein Area Detail for Tyrosine