

Application of Decision Tree

Protein Solvent Accessibility Prediction

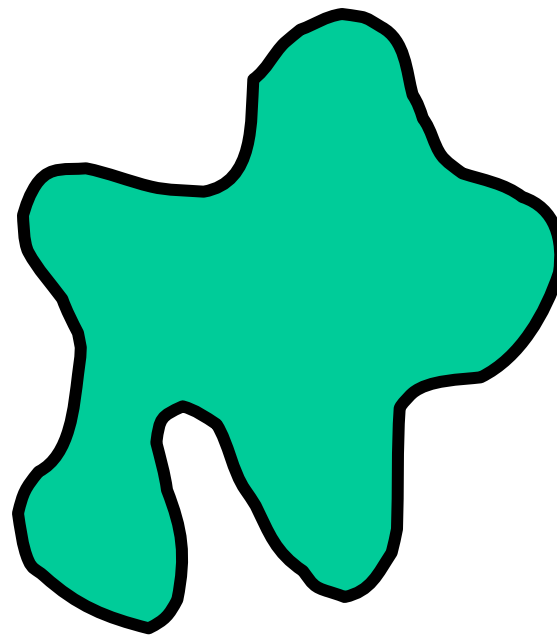
What is a molecular surface?

A molecular surface is a closed 3D "manifold".

What's a "manifold"?

Here is a 2D manifold.

A cell, for example, is a 3D manifold. It is continuous, closed, non-intersecting. It has an inside and an outside.

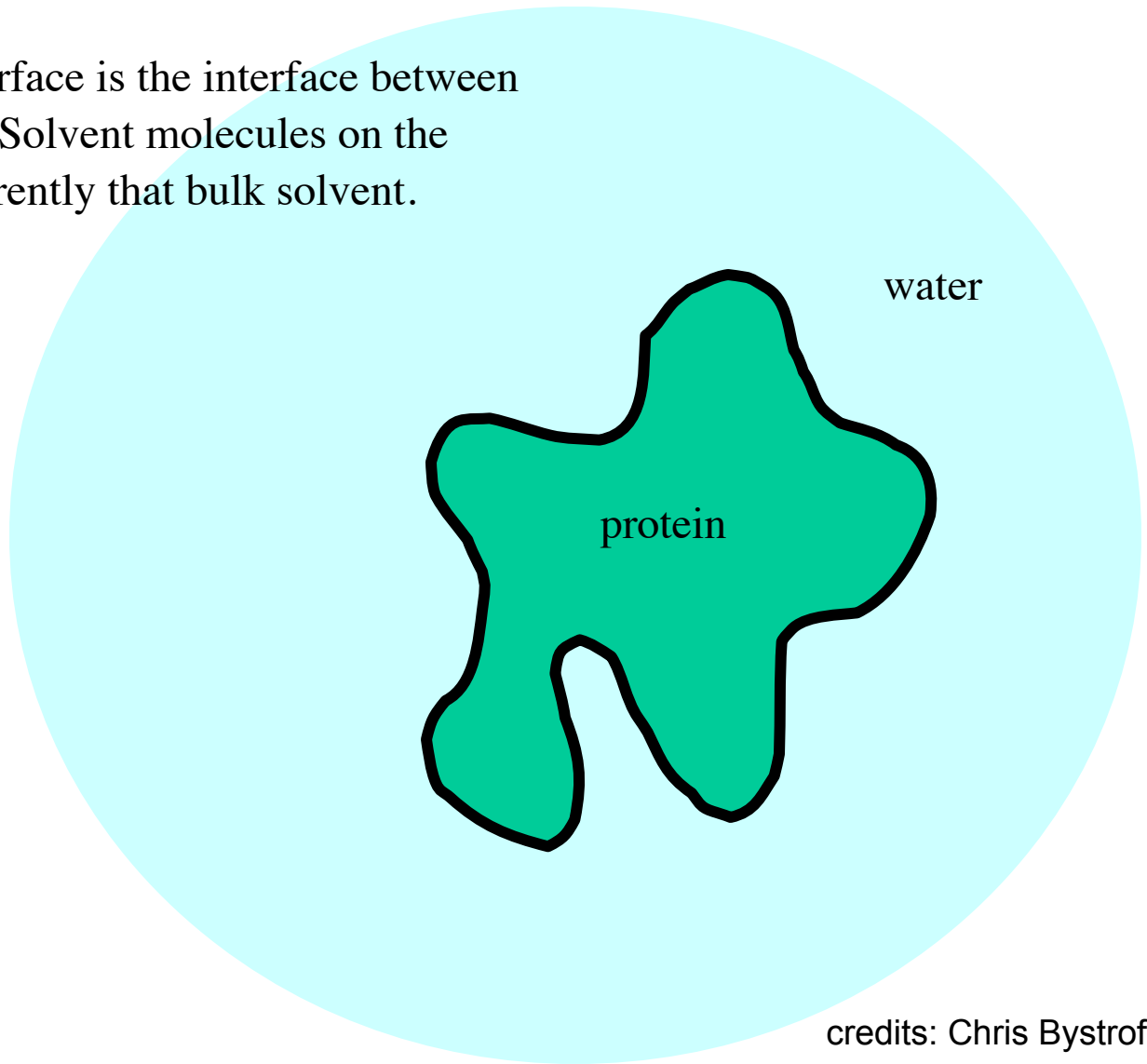


Solvent accessible surface

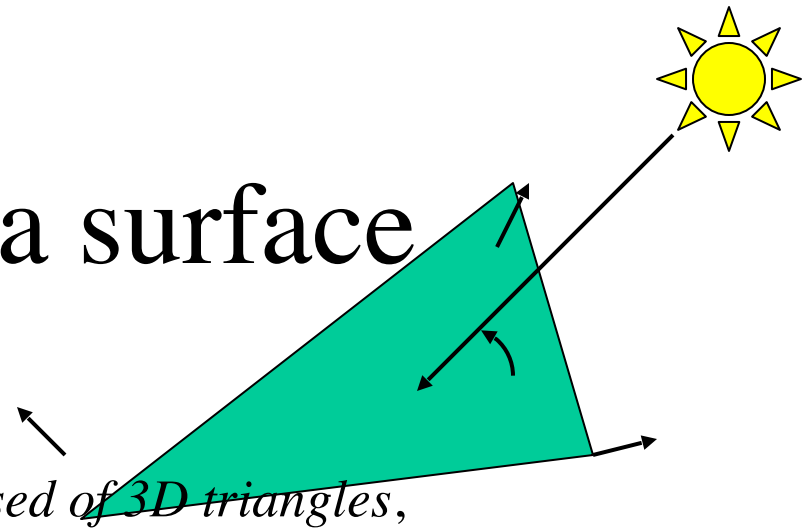
The solvent accessible surface is the interface between molecule and its solvent. Solvent molecules on the surface may behave differently than bulk solvent.

Surfaces have:

- size/area
- electrostatic properties
- shape properties

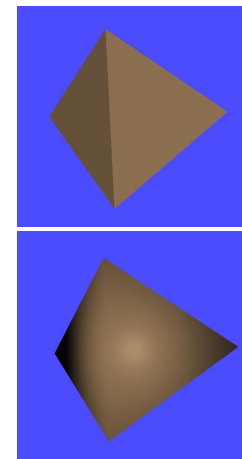


Rendering a surface



- A surface of any shape may be composed of 3D triangles, that is, 3 sets of xyz coordinates, one for each vertex.
- To display the surface on the screen, each triangle is rotated and translated according to the current **frame of reference**.
- Continuous triangles make a continuous surface.
- Then, each pixel is assigned a **brightness** according to the angle between the triangle and the light source.

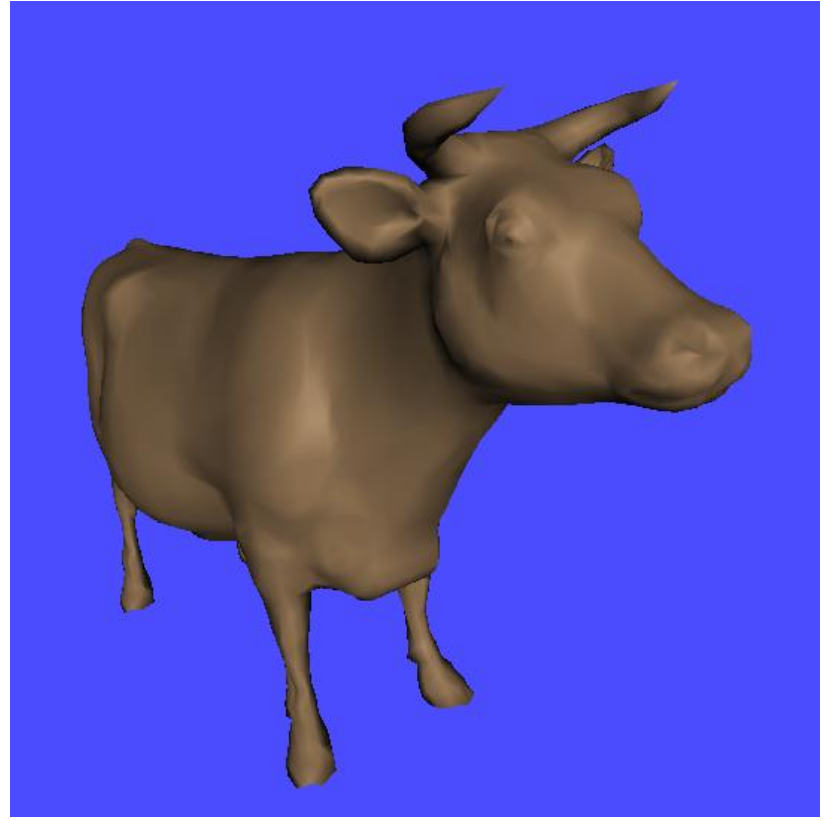
- **Phong shading** may be applied to simulate *curvature*. In this case, each pixel in the triangle has a different brightness, depending on where it is.



Surfaces maybe described as a set of connected triangles



A cow-shaped manifold made of triangles.

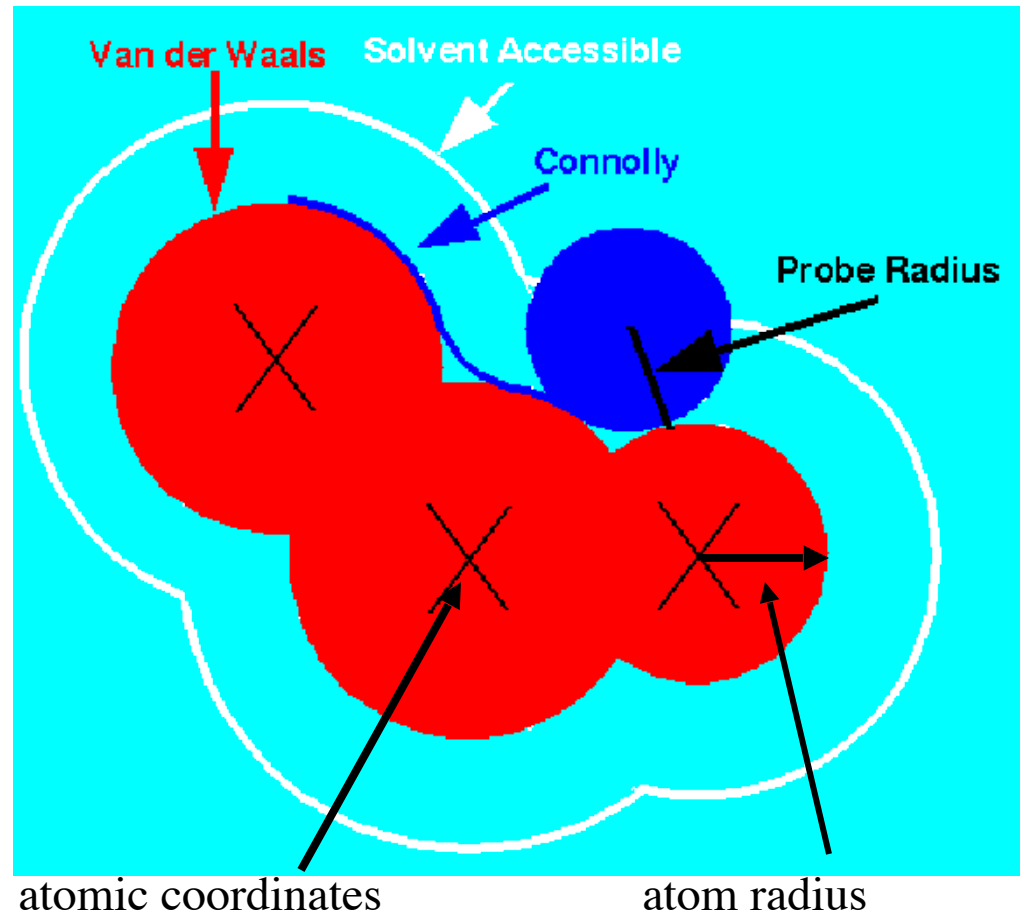


Phong-shaded cow. Shading give the illusion of higher resolution.

The Connolly surface

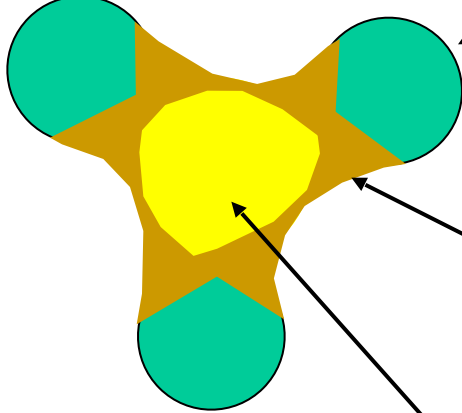
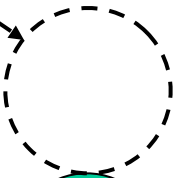
Conceptually, roll a probe sphere over the molecule...

- Everywhere the center of the sphere goes is the Solvent Accessible Surface (SAS).
- Everywhere the sphere touches (including empty space) is the Solvent Excluded (or "Connolly") Surface (SES).



Surface shapes

rolling probe sphere

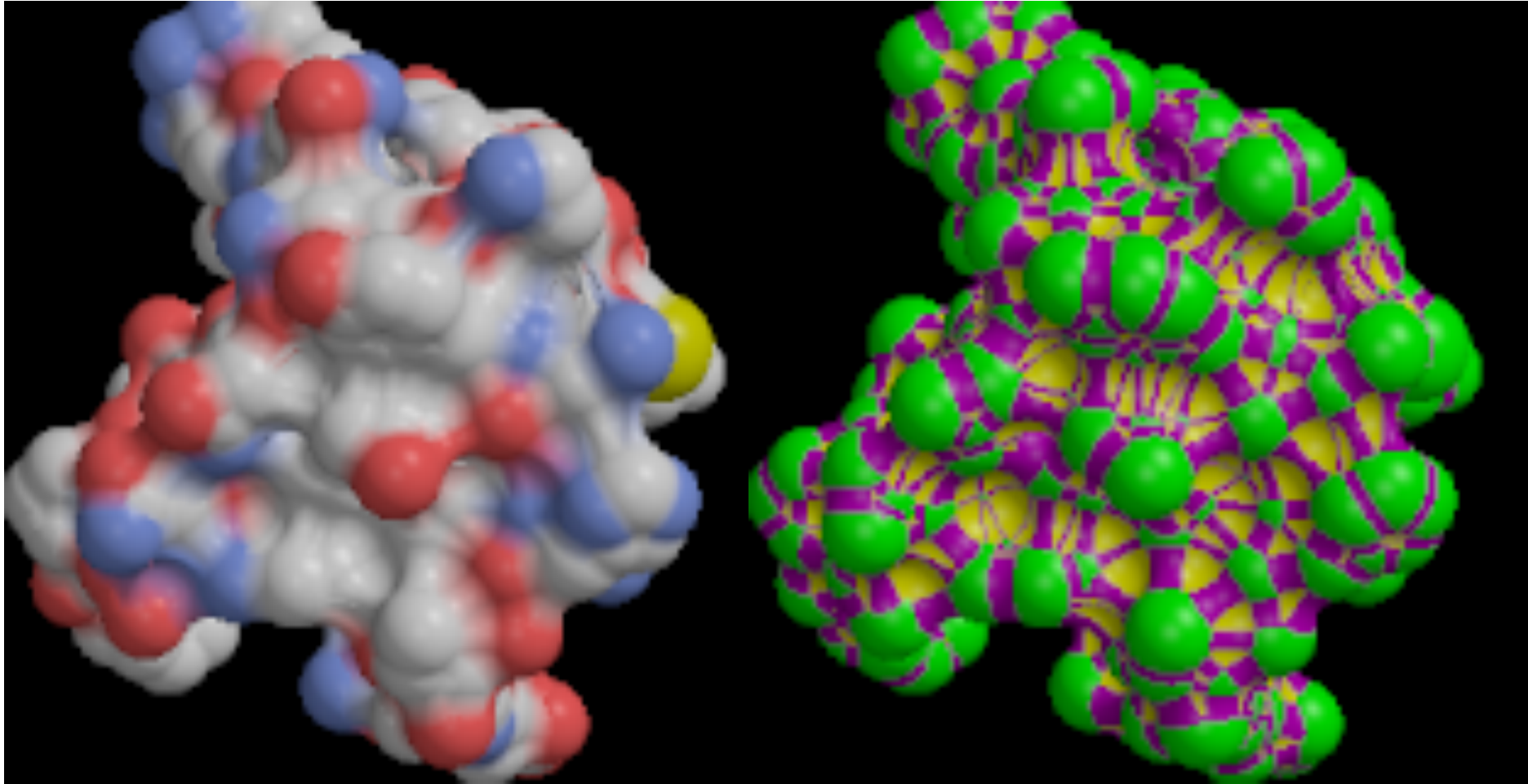


Green: convex-convex, contact surface of probe with atom.

Brown: convex-concave, toroidal surface when touches two atoms.

Yellow: concave-concave, reentrant surface when touches three atoms.

Coloring by atom, by shape



Surfaces maybe shaded by partial charge.

or by shape. Yellow parts are 'reentrant'.

credits: Chris Bystroff

Chimera demo

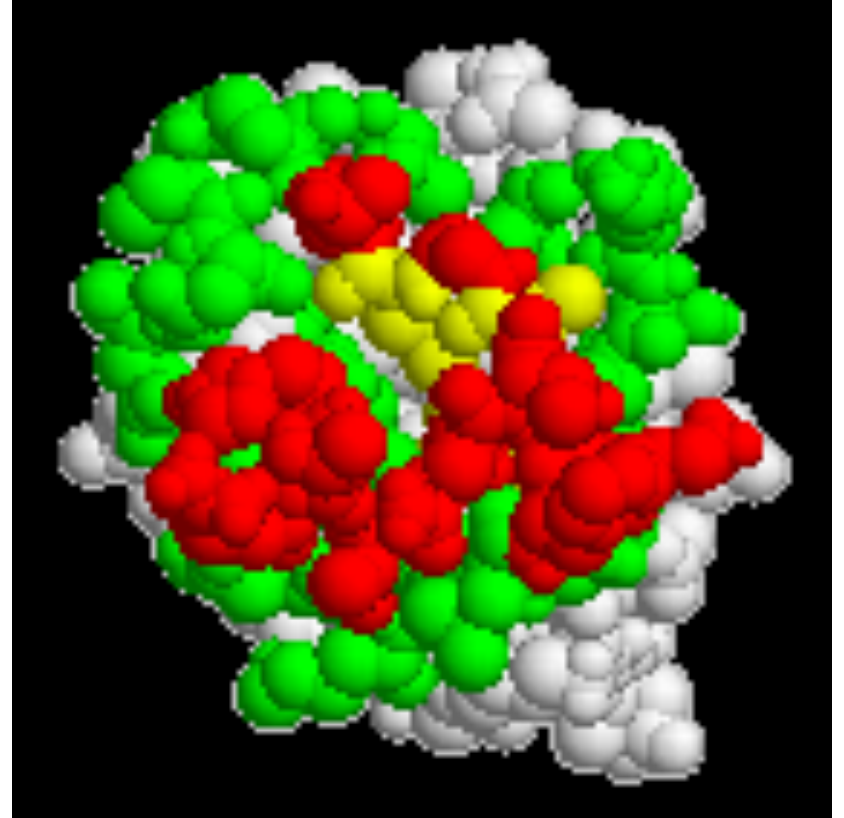
1BYI

Exposure of amino acid to solvent
is quantified by solvent accessibility

Solvent Accessibility

Size of the area of an amino acid that is exposed to solvent (water).

- Maximum solvent accessible area for each amino acid is its whole surface area.
- Hydrophobic residues like to be Buried inside (interior).
- Hydrophilic residues like to be exposed on the surface.



Structure Analysis

- If we know the 3D protein structure, we can calculate solvent accessible surface area (SASA) for each amino acids from 3D structure
- Then we can calculate the relative solvent accessibility (RSA) by calculating what percentage is exposed to solvent
- If the RSA is $< 25\%$, then buried (**b**), else exposed (**e**)
- Most widely used tool: **DSSP** (Dictionary of Protein Secondary Structure: Pattern Recognition of Hydrogen-Bonded and Geometrical Features. Kabsch and Sander, 1983)

But what this has to do
with decision tree?

Protein Relative Solvent Accessibility Prediction

- We have two classes: buried (**b**), else exposed (**e**)
 - Binary labels (Y)
- We know that amino acids has some chemical properties (X)
 - Like some are charged,
 - Instance space, X
 - Sample of labeled training data $\{ \langle x^{(i)}, y^{(i)} \rangle \}$
 - Hypothesis space, $H = \{ f: X \rightarrow Y \}$

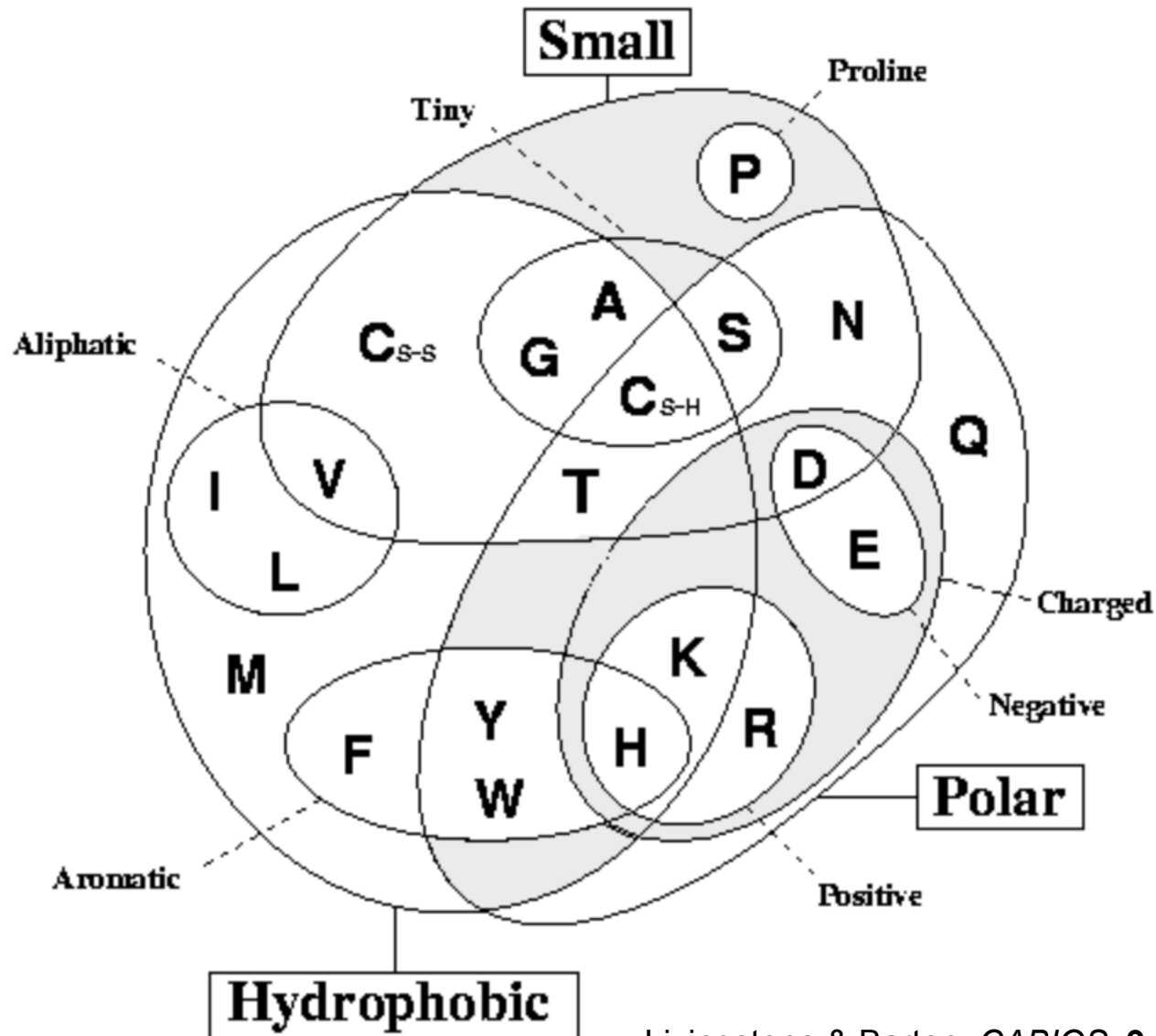
$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ (\text{CH}_2)_3 \\ \\ \text{NH} \\ \\ \text{C}=\text{NH}_2 \\ \\ \text{NH}_2 \end{array} $ <p>Arginine (Arg / R)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{C}=\text{O} \\ \\ \text{NH}_2 \end{array} $ <p>Glutamine (Gln / Q)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{C}_6\text{H}_5 \end{array} $ <p>Phenylalanine (Phe / F)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{C}_6\text{H}_4 \\ \\ \text{OH} \end{array} $ <p>Tyrosine (Tyr / Y)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{C}_8\text{H}_6\text{N} \end{array} $ <p>Tryptophan (Trp, W)</p>
$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ (\text{CH}_2)_4 \\ \\ \text{NH}_2 \end{array} $ <p>Lysine (Lys / K)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{H} \end{array} $ <p>Glycine (Gly / G)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_3 \end{array} $ <p>Alanine (Ala / A)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{C}_4\text{H}_3\text{N}_2 \end{array} $ <p>Histidine (His / H)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{OH} \end{array} $ <p>Serine (Ser / S)</p>
$ \begin{array}{c} \text{H}_2 \\ \\ \text{C} \\ / \quad \backslash \\ \text{H}_2\text{C} \quad \text{CH}_2 \\ \backslash \quad / \\ \text{H}_2\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \end{array} $ <p>Proline (Pro / P)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{COOH} \end{array} $ <p>Glutamic Acid (Glu / E)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{COOH} \end{array} $ <p>Aspartic Acid (Asp / D)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{H} - \text{C} - \text{OH} \\ \\ \text{CH}_3 \end{array} $ <p>Threonine (Thr / T)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{SH} \end{array} $ <p>Cysteine (Cys / C)</p>
$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{S} \\ \\ \text{CH}_3 \end{array} $ <p>Methionine (Met / M)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{CH} \\ / \quad \backslash \\ \text{CH}_3 \quad \text{CH}_3 \end{array} $ <p>Leucine (Leu / L)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH}_2 \\ \\ \text{C}=\text{O} \\ \\ \text{NH}_2 \end{array} $ <p>Asparagine (Asn / N)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{HC} - \text{CH}_3 \\ \\ \text{CH}_2 \\ \\ \text{CH}_3 \end{array} $ <p>Isoleucine (Ile / I)</p>	$ \begin{array}{c} \text{H} \\ \\ \text{H}_3\text{N}^+ - \text{C} - \text{C} \begin{array}{l} \nearrow \text{O}^- \\ \searrow \text{O}^- \end{array} \\ \\ \text{CH} \\ / \quad \backslash \\ \text{CH}_3 \quad \text{CH}_3 \end{array} $ <p>Valine (Val / V)</p>

20
naturally
occurring
amino acid
residues

Properties of Amino Acids

Amino acid	Abbrev.	Side chain	Hydro-phobic	Polar	Charged	Small	Tiny	Aromatic or Aliphatic	van der Waals volume	Codon	Occurrence in proteins (%)
Alanine	Ala, A	-CH ₃	X	-	-	X	X	-	67	GCU, GCC, GCA, GCG	7.8
Cysteine	Cys, C	-CH ₂ SH	X	-	-	X	-	-	86	UGU, UGC	1.9
Aspartate	Asp, D	-CH ₂ COOH	-	X	negative	X	-	-	91	GAU, GAC	5.3
Glutamate	Glu, E	-CH ₂ CH ₂ COOH	-	X	negative	-	-	-	109	GAA, GAG	6.3
Phenylalanine	Phe, F	-CH ₂ C ₆ H ₅	X	-	-	-	-	Aromatic	135	UUU, UUC	3.9
Glycine	Gly, G	-H	X	-	-	X	X	-	48	GGU, GGC, GGA, GGG	7.2
Histidine	His, H	-CH ₂ -C ₃ H ₃ N ₂	-	X	positive	-	-	Aromatic	118	CAU, CAC	2.3
Isoleucine	Ile, I	-CH(CH ₃)CH ₂ CH ₃	X	-	-	-	-	Aliphatic	124	AUU, AUC, AUA	5.3
Lysine	Lys, K	-(CH ₂) ₄ NH ₂	-	X	positive	-	-	-	135	AAA, AAG	5.9
Leucine	Leu, L	-CH ₂ CH(CH ₃) ₂	X	-	-	-	-	Aliphatic	124	UUA, UUG, CUU, CUC, CUA, CUG	9.1
Methionine	Met, M	-CH ₂ CH ₂ SCH ₃	X	-	-	-	-	-	124	AUG	2.3
Asparagine	Asn, N	-CH ₂ CONH ₂	-	X	-	X	-	-	96	AAU, AAC	4.3
Proline	Pro, P	-CH ₂ CH ₂ CH ₂ -	X	-	-	X	-	-	90	CCU, CCC, CCA, CCG	5.2
Glutamine	Gln, Q	-CH ₂ CH ₂ CONH ₂	-	X	-	-	-	-	114	CAA, CAG	4.2
Arginine	Arg, R	-(CH ₂) ₃ NH-C(NH) NH ₂	-	X	positive	-	-	-	148	CGU, CGC, CGA, CGG, AGA, AGG	5.1
Serine	Ser, S	-CH ₂ OH	-	X	-	X	X	-	73	UCU, UCC, UCA, UCG, AGU, AGC	6.8
Threonine	Thr, T	-CH(OH)CH ₃	X	X	-	X	-	-	93	ACU, ACC, ACA, ACG	5.9
Valine	Val, V	-CH(CH ₃) ₂	X	-	-	X	-	Aliphatic	105	GUU, GUC, GUA, GUG	6.6
Tryptophan	Trp, W	-CH ₂ C ₈ H ₆ N	X	-	-	-	-	Aromatic	163	UGG	1.4
Tyrosine	Tyr, Y	-CH ₂ -C ₆ H ₄ OH	X	X	-	-	-	Aromatic	141	UAU, UAC	3.2

Venn diagram of Amino Acid Properties



Attributes

Hydrophobic

Polar

Small

Proline

Tiny

Aliphatic

Aromatic

Positive

Negative

Charged

Labels

Buried

Exposed

Protein RSA Prediction using Decision Trees

- Well posed function approximation problems:
 - Instance space, X
 - Sample of labeled training data $\{ \langle x^{(i)}, y^{(i)} \rangle \}$
 - Hypothesis space, $H = \{ f: X \rightarrow Y \}$
- Learning is a search/optimization problem over H
 - Various objective functions
 - minimize training error (0-1 loss)
 - among hypotheses that minimize training error, select smallest (?)
- Decision tree learning
 - Greedy top-down learning of decision trees (ID3, C4.5, ...)
 - Overfitting and tree/rule post-pruning
 - Extensions...