Pattern Recognition

Problem Representation

Experimental Results

A Pattern Recognition Based Model for Characterizing and Predicting Glucose-Binding Sites

Houssam Nassif

June 7, 2006



Biochemical	Background

Problem Representation 000000

Experimental Results

- Biochemical Background
 - Glucose
 - Proteins
 - Biochemical Interactions

Biochemical	Background

Problem Representation 000000

Experimental Results

- 1 Biochemical Background
 - Glucose
 - Proteins
 - Biochemical Interactions
- 2 Pattern Recognition
 - *k*-Nearest-Neighbors
 - Support Vector Machines
 - Random Forest

Biochemical	Background

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- 1 Biochemical Background
 - Glucose
 - Proteins
 - Biochemical Interactions
- 2 Pattern Recognition
 - k-Nearest-Neighbors
 - Support Vector Machines
 - Random Forest
- 3 Problem Representation
 - Data
 - Solution Approach

Biochemical	Background

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- Biochemical Background
 - Glucose
 - Proteins
 - Biochemical Interactions
- 2 Pattern Recognition
 - *k*-Nearest-Neighbors
 - Support Vector Machines
 - Random Forest
- 3 Problem Representation
 - Data
 - Solution Approach
- 4 Experimental Results
 - Learning Phase
 - Feature Selection
 - Testing Phase

Problem Desc	ription		
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results

• Glucose is a 6-carbon sugar molecule that plays a key role in many different biochemical pathways.

Problem Descri	ption		
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results

- Glucose is a 6-carbon sugar molecule that plays a key role in many different biochemical pathways.
- Glucose binds to protein molecules at specific binding sites.

Problem Desc	rintion		
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results 00000000000000000

- Glucose is a 6-carbon sugar molecule that plays a key role in many different biochemical pathways.
- Glucose binds to protein molecules at specific binding sites.
- Binding sites are specific to their respective ligands.

Definition

Ligand: The specific molecule that binds to the binding site.

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

Problem Desc	rintion		
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results 00000000000000000

- Glucose is a 6-carbon sugar molecule that plays a key role in many different biochemical pathways.
- Glucose binds to protein molecules at specific binding sites.
- Binding sites are specific to their respective ligands.
- Identifying the ligand docking at a certain binding site remains an open research problem.

Definition

Ligand: The specific molecule that binds to the binding site.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Problem Desc	rintion		
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results 00000000000000000

- Glucose is a 6-carbon sugar molecule that plays a key role in many different biochemical pathways.
- Glucose binds to protein molecules at specific binding sites.
- Binding sites are specific to their respective ligands.
- Identifying the ligand docking at a certain binding site remains an open research problem.

Definition

Ligand: The specific molecule that binds to the binding site.

Goal

Identify glucose-binding sites.

Pattern Recognition

Problem Representation 000000

Experimental Results

Biochemical Background

A glucose-binding site involves:

Pattern Recognition

Problem Representation 000000

Experimental Results

Biochemical Background

- A glucose-binding site involves:
 - A glucose molecule.

Pattern Recognition

Problem Representation 000000

Experimental Results

Biochemical Background

- A glucose-binding site involves:
 - A glucose molecule.
 - 2 A protein.

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Biochemical Background

- A glucose-binding site involves:
 - A glucose molecule.
 - A protein.
 - O Their interaction.

Biochemical	Background
0000000	

Problem Representation 000000

Experimental Results

Glucose Structure



Figure: Glucose

• Glucose is a 6-carbon sugar.

Biochemical Background ●00000000	Pattern Recognition	Problem Representation	Experimental Results
Glucose Struct	TITA		



Figure: Glucose

- Glucose is a 6-carbon sugar.
- It contains two functional groups.

Biochemical Background ●00000000	Pattern Recognition	Problem Representation	Experimental 000000000
Chucose Struct			



- Glucose is a 6-carbon sugar.
- It contains two functional groups.

Results



Figure: Glucose

(a) Carbonyl

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Glucose Struct	ure		



Figure: Glucose

- Glucose is a 6-carbon sugar.
- It contains two functional groups.



(a) Carbonyl

(b) Hydroxyl

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Biochemical	Background
0000000	

Problem Representation 000000

Experimental Results

Glucose Structure



- Glucose is a 6-carbon sugar.
- It contains two functional groups.
- Both groups can interact together.



Figure: Glucose

(a) Carbonyl (b) Hydroxyl

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Biochemical Background 0●0000000	Pattern Recognition	Problem Representation	Experimental Results
Glucose Cvcliz	ation		

• The molecule folds on itself and forms a *pyranose* ring.



Biochemical Background 0●0000000	Pattern Recognition	Problem Representation	Experimental Results
Glucose Cycliz	ation		

• The molecule folds on itself and forms a *pyranose* ring.



Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
0000000			
Glucose Cycliz	ation		

- The molecule folds on itself and forms a *pyranose* ring.
- In two different ways. Watch the star!



▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Pattern Recognitio

Problem Representation 000000

Experimental Results

Conformation Shift

• Glucose readily shifts from one conformation to another.

Pattern Recognition

Problem Representation

Experimental Results

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)
 (日)

 (日)
 (日)

 (日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

Conformation Shift

- Glucose readily shifts from one conformation to another.
- In physiological solutions: almost exclusively in the pyranose ring form.

Pattern Recognition

Problem Representation

Experimental Results

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)
 (日)

 (日)
 (日)

 (日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

Conformation Shift

- Glucose readily shifts from one conformation to another.
- In physiological solutions: almost exclusively in the pyranose ring form.
 - 36% α -pyranose

Pattern Recognition

Problem Representation

Experimental Results

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)
 (日)

 (日)
 (日)

 (日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

Conformation Shift

- Glucose readily shifts from one conformation to another.
- In physiological solutions: almost exclusively in the pyranose ring form.
 - 36% α -pyranose
 - 64% β -pyranose

Pattern Recognition

Problem Representation 000000

Experimental Results

Amino Acid Structure

• An amino acid consists of:

Pattern Recognition

Problem Representation 000000

Experimental Results

Amino Acid Structure

- An amino acid consists of:
 - a central carbon atom C, bonded to:





Pattern Recognition

Problem Representation

Experimental Results

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)
 (日)

 (日)
 (日)

 (日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

Amino Acid Structure

- An amino acid consists of:
 - a central carbon atom C, bonded to:
 - an amino group NH₂



Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Amino Acid Structure



- a central carbon atom C, bonded to:
- an amino group NH₂
- a carboxyl group COOH



Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Amino Acid Structure



• An amino acid consists of:

- a central carbon atom C, bonded to:
- an amino group NH₂
- a carboxyl group COOH
- a hydrogen atom H

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Amino Acid Structure



• An amino acid consists of:

- a central carbon atom C, bonded to:
- an amino group NH₂
- a carboxyl group COOH
- a hydrogen atom H
- and a side chains R.

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Amino Acid Structure



• An amino acid consists of:

- a central carbon atom C, bonded to:
- an amino group NH₂
- a carboxyl group COOH
- a hydrogen atom H
- and a side chains R.
- Amino acids differ by their side chain R.

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Amino Acid Structure



Figure: Amino acid

- An amino acid consists of:
 - a central carbon atom C, bonded to:
 - an amino group NH₂
 - a carboxyl group COOH
 - a hydrogen atom H
 - and a side chains R.
- Amino acids differ by their side chain R.
- The side chain confers to each amino acid its distinctive properties.

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Amino Acid Structure



Figure: Amino acid

- An amino acid consists of:
 - a central carbon atom C, bonded to:
 - an amino group NH₂
 - a carboxyl group COOH
 - a hydrogen atom H
 - and a side chains R.
- Amino acids differ by their side chain R.
- The side chain confers to each amino acid its distinctive properties.
- There are 20 different amino acids.

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

Amino Acid Structure



Figure: Amino acid

- An amino acid consists of:
 - a central carbon atom C, bonded to:
 - an amino group NH₂
 - a carboxyl group COOH
 - a hydrogen atom H
 - and a side chains R.
- Amino acids differ by their side chain R.
- The side chain confers to each amino acid its distinctive properties.
- There are 20 different amino acids.

Amino acid properties

Similar R \Leftrightarrow similar properties. Different R \Leftrightarrow different properties.
Biochemical Background ○○○○●○○○○	Pattern Recognition	Problem Representation	Experimental Results
Protein Structur	е		

◆□ ▶ < @ ▶ < E ▶ < E ▶ E ■ 9 Q @</p>

• Protein: a long chain of amino acids linked together.

Biochemical Background ○○○○●○○○○	Pattern Recognition	Problem Representation	Experimental Results
Protein Structure	9		

• Protein: a long chain of amino acids linked together.



Biochemical Background ○○○○●○○○○	Pattern Recognition	Problem Representation	Experimental Results
Protein Structur	е		

- Protein: a long chain of amino acids linked together.
- Amino acids in a protein are also called residues.



Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Protein Structu	re		

- Protein: a long chain of amino acids linked together.
- Amino acids in a protein are also called residues.
- Residues sequence and properties determine the protein shape and function.



Figure: Linked residues

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Protein Struct	lire		

- Protein: a long chain of amino acids linked together.
- Amino acids in a protein are also called residues.
- Residues sequence and properties determine the protein shape and function.



Important Fact

Similar residues can be easily interchanged in a protein.

Pattern Recognition

Problem Representation 000000

Experimental Results

Atomic Interactions

Atoms interact together through:

Pattern Recognition

Problem Representation 000000

Experimental Results

◆□> <□> <=> <=> <=> <=> <=> <=> <=>

Atomic Interactions

Atoms interact together through:

Pattern Recognition

Problem Representation 000000

Experimental Results

Atomic Interactions

Atoms interact together through:

Chemical bonds

Porces and Interactions

Pattern Recognition

Problem Representation

Experimental Results

Atomic Interactions

Atoms interact together through:

Chemical bonds

Covalent Bonds

Porces and Interactions

Pattern Recognition

Problem Representation

Experimental Results

Atomic Interactions

Atoms interact together through:

- Covalent Bonds
- Ø Hydrogen Bonds
- Porces and Interactions

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Atomic Interactions

Atoms interact together through:

- Covalent Bonds
- Ø Hydrogen Bonds
- Porces and Interactions
 - Van der Waals Forces

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Atomic Interactions

Atoms interact together through:

- Covalent Bonds
- Ø Hydrogen Bonds
- Porces and Interactions
 - Van der Waals Forces
 - Ø Hydrophobic Interactions

Biochemical Background ○○○○○○●○○	Pattern Recognition	Problem Representation	Experimental Results
Covalent Bonds			

• Close and strong interaction.

◆□ > < 個 > < E > < E > E = 9000

Figure: Covalent bond

Ο

-H

Biochemical Background ○○○○○●○○	Pattern Recognition	Problem Representation	Experimental Results
Covalent Bonds			



• Close and strong interaction.

◆□▶ ◆□▶ ◆目▶ ◆目■ のへで

• Form a molecule.

Figure: Covalent bond

Biochemical Background	Pattern Recognition	Problem Representation	Experimen 0000000
Covalent Ronds			



Figure: Covalent bond

• Close and strong interaction.

Results

◆□▶ ◆□▶ ◆目▶ ◆目■ のへで

- Form a molecule.
- Atoms share electrons.

Biochemical	Background
00000000	0

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Covalent Bonds



- Close and strong interaction.
- Form a molecule.
- Atoms share electrons.
- Electronegativity.

Figure: Covalent bond

Definition

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Covalent Bonds



Figure: Covalent bond

Definition

- Close and strong interaction.
- Form a molecule.
- Atoms share electrons.
- Electronegativity.
 - Equal \Rightarrow nonpolar

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Covalent Bonds



Figure: Covalent bond

- Close and strong interaction.
- Form a molecule.
- Atoms share electrons.
- Electronegativity.
 - Equal \Rightarrow nonpolar
 - Different \Rightarrow polar

Definition

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Covalent Bonds



Figure: Covalent polar bond

- Close and strong interaction.
- Form a molecule.
- Atoms share electrons.
- Electronegativity.
 - Equal ⇒ nonpolar
 - Different \Rightarrow polar
- Partial charges.

Definition

Pattern Recognition

Problem Representation

Experimental Results

Hydrogen Bonds

• The most relevant in protein interaction.

Pattern Recognition

Problem Representation

Experimental Results

Hydrogen Bonds



Figure: Hydrogen bond

- The most relevant in protein interaction.
- Attraction between a positively charged H and a negatively charged atoms.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Pattern Recognition

Problem Representation

Experimental Results

Hydrogen Bonds



Figure: Hydrogen bond

- The most relevant in protein interaction.
- Attraction between a positively charged H and a negatively charged atoms.
- Glucose attaches to the protein using hydrogen bonds.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Pattern Recognition

Problem Representation 000000

Experimental Results

Forces and Interactions

• Van der Waals Forces

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆○▼

Pattern Recognition

Problem Representation 000000

Experimental Results

Forces and Interactions

O Van der Waals Forces

Characteristics

Weak electrostatic attraction and repulsion forces.

Biochemical Background ○○○○○○○● Pattern Recognition

Problem Representation 000000

Experimental Results

Forces and Interactions

- Van der Waals Forces
- Ø Hydrophobic Interactions

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Forces and Interactions

- Van der Waals Forces
- Ø Hydrophobic Interactions

Definition

Hydrophobic: water hating. Hydrophilic: water loving.

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Forces and Interactions

- Van der Waals Forces
- Ø Hydrophobic Interactions

Definition

Hydrophobic: water hating. Hydrophilic: water loving.

Characteristics

Hydrophobic atoms tend to gather together.

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Forces and Interactions

- Van der Waals Forces
- O Hydrophobic Interactions

Definition

Hydrophobic: water hating. Hydrophilic: water loving.

Characteristics

- Hydrophobic atoms tend to gather together.
- Hydrophilic atoms tend to gather together.

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Forces and Interactions

- Van der Waals Forces
- O Hydrophobic Interactions
- The pyranose ring is hydrophobic

Definition

Hydrophobic: water hating. Hydrophilic: water loving.

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Forces and Interactions

- Van der Waals Forces
- Ø Hydrophobic Interactions
- The pyranose ring is hydrophobic

Definition

Hydrophobic: water hating. Hydrophilic: water loving.

Characteristics

 \Rightarrow attracted by hydrophobic residues

Biochemical	Background

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Pattern Recognition

Pattern recognition

Definition

Pattern recognition: The computer learns how to correctly classify an object.

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results 0000000000000000
Pattern Recog	nition		

 \circ

Classifier

・ロト < 団ト < 三ト < 三ト < 三ト < 〇への

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results 0000000000000000
Pattern Recog	nition		

Classifier

Classifier phases

• Training phase: Learns how to partition the feature space.

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Pattern Recog	nition		

Classifier

Classifier phases

• Training phase: Learns how to partition the feature space.

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

• Testing phase: Predicts the class of an unknown data.

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Pattern Recognit	ion		

- Classifier
 - k-Nearest-Neighbor (kNN)

Classifier phases

• Training phase: Learns how to partition the feature space.

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

• Testing phase: Predicts the class of an unknown data.

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Pattern Recog	nition		

ъ

- Classifier
 - k-Nearest-Neighbor (kNN)
 - Support Vector Machines (SVM)

Classifier phases

• Training phase: Learns how to partition the feature space.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

• Testing phase: Predicts the class of an unknown data.
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Pattern Recogn	ition		

Pattern recognition includes:

- Classifier
 - k-Nearest-Neighbor (kNN)
 - Support Vector Machines (SVM)
- Peature selection

Classifier phases

• Training phase: Learns how to partition the feature space.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

• Testing phase: Predicts the class of an unknown data.

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Pattern Recogn	ition		

Pattern recognition includes:

- Classifier
 - k-Nearest-Neighbor (kNN)
 - Support Vector Machines (SVM)
- Peature selection
 - Random Forest (RF)

Classifier phases

• Training phase: Learns how to partition the feature space.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

• Testing phase: Predicts the class of an unknown data.

Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Simple *k*NN technique



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Simple kNN technique limitations



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Simple kNN technique limitations

• Cannot handle even k's.



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Simple kNN technique limitations

- Cannot handle even k's.
- Ignores distance!



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Solution: Weighted kNN



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Solution: Weighted kNN



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Solution: Weighted kNN



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Solution: Weighted kNN



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Solution: Weighted kNN



Pattern Recognition

Problem Representation 000000

Experimental Results

k-Nearest-Neighbors (kNN)

Solution: Weighted kNN

- Puts a higher weight on lesser distances
- Can use any k



Pattern Recognition

Problem Representation 000000

Experimental Results

Support Vector Machines (SVM)



Pattern Recognition

Problem Representation 000000

Experimental Results

Support Vector Machines (SVM)

SVM technique

• Construct the optimal separating hyperplane.



Pattern Recognition

Problem Representation

Experimental Results

Support Vector Machines (SVM)

SVM technique

- Construct the optimal separating hyperplane.
- Maximize the margins.

Definition

Margin: Minimal distance away from the hyperplane.



Experimental Results

Support Vector Machines (SVM)

- Construct the optimal separating hyperplane.
- Maximize the margins.
- A small set of samples specifies the hyperplane.



Experimental Results

Support Vector Machines (SVM)

- Construct the optimal separating hyperplane.
- Maximize the margins.
- A small set of samples specifies the hyperplane.



Problem Representation

Experimental Results

Support Vector Machines (SVM)

- Construct the optimal separating hyperplane.
- Maximize the margins.
- A small set of samples specifies the hyperplane.
- Called Support Vectors (SV).



Biochemical Background	Pattern Recognition ○○●○	Problem Representation	Experimental Results
SVM Characte	ristics		

• Based on a kernel function.



Biochemical Background	Pattern Recognition ○O●○	Problem Representation	Experimental Results 000000000000000000
SVM Characte	ristics		

• Based on a kernel function.

Definition

Kernel function: A function that maps input data to a higher dimensional feature space and computes their scalar product.

$$K(\overrightarrow{a}, \overrightarrow{b}) = \phi(\overrightarrow{a})^t \cdot \phi(\overrightarrow{b}) \quad \text{for} \quad \overrightarrow{a}, \overrightarrow{b} \in X$$

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Biochemical Background	Pattern Recognition ○○●○	Problem Representation	Experimental Results
SVM Characte	eristics		

- Based on a kernel function.
- Performs scalar computation in input space \Rightarrow efficient.

Kernel function: A function that maps input data to a higher dimensional feature space and computes their scalar product.

$$K(\overrightarrow{a}, \overrightarrow{b}) = \phi(\overrightarrow{a})^t \cdot \phi(\overrightarrow{b}) \quad \text{for} \quad \overrightarrow{a}, \overrightarrow{b} \in X$$

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

Biochemical Background	Pattern Recognition ○0●○	Problem Representation	Experimental Results
SVM Characteri	stics		

- Based on a kernel function.
- Performs scalar computation in input space \Rightarrow efficient.
- Maps input to higher feature space \Rightarrow generalization problem.

Kernel function: A function that maps input data to a higher dimensional feature space and computes their scalar product.

$$K(\overrightarrow{a}, \overrightarrow{b}) = \phi(\overrightarrow{a})^t \cdot \phi(\overrightarrow{b}) \quad \text{for} \quad \overrightarrow{a}, \overrightarrow{b} \in X$$

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Biochemical Background	Pattern Recognition ○○●○	Problem Representation	Experimental Results 000000000000000000
SVM Characteri	stics		

- Based on a kernel function.
- Performs scalar computation in input space \Rightarrow efficient.
- Maps input to higher feature space \Rightarrow generalization problem.

Kernel function: A function that maps input data to a higher dimensional feature space and computes their scalar product.

$$K(\overrightarrow{a},\overrightarrow{b}) = \phi(\overrightarrow{a})^t \cdot \phi(\overrightarrow{b}) \quad \text{for} \quad \overrightarrow{a}, \overrightarrow{b} \in X$$

Definition

Generalization: Ability to correctly classify new data.

Biochemical Background	Pattern Recognition ○0●○	Problem Representation	Experimental Results
SVM Characteri	stics		

- Based on a kernel function.
- Performs scalar computation in input space \Rightarrow efficient.
- Maps input to higher feature space \Rightarrow generalization problem.
- Small number of support vectors ⇔ good generalization.

Kernel function: A function that maps input data to a higher dimensional feature space and computes their scalar product.

$$K(\overrightarrow{a}, \overrightarrow{b}) = \phi(\overrightarrow{a})^t \cdot \phi(\overrightarrow{b}) \quad \text{for} \quad \overrightarrow{a}, \overrightarrow{b} \in X$$

Definition

Generalization: Ability to correctly classify new data.



◆□> <□> <=> <=> <=> <=> <=> <=> <=>

• High number of features \Rightarrow *curse of dimensionality*.



- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙



- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Random Forest feature selection tool:



- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Random Forest feature selection tool:

• Based on multiple classification trees



- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Random Forest feature selection tool:

- Based on multiple classification trees
- Provides direct feature importance measure



- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Random Forest feature selection tool:

- Based on multiple classification trees
- Provides direct feature importance measure
- $\bullet\,$ Can be used when feature number \gg samples


- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Random Forest feature selection tool:

- Based on multiple classification trees
- Provides direct feature importance measure
- Can be used when feature number \gg samples
- Robust to noise



- High number of features \Rightarrow curse of dimensionality.
- Feature selection: select the best subset of the input features.

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Random Forest feature selection tool:

- Based on multiple classification trees
- Provides direct feature importance measure
- Can be used when feature number \gg samples
- Robust to noise
- Low bias and low variance

Pattern Recognition

Problem Representation

Experimental Results

Problem Representation

A good classifier must be trained on:

Representative data

Pattern Recognition

Problem Representation

Experimental Results

Problem Representation

A good classifier must be trained on:

- Representative data
- O The right features

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Problem Representation

A good classifier must be trained on:

- Representative data
- O The right features

This sections explains:

How the data is obtained

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Problem Representation

A good classifier must be trained on:

- Representative data
- O The right features

This sections explains:

- How the data is obtained
- e How the features are represented

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results
Database Avai	lability		

The Protein Data Bank (PDB):

• Repository of 3-D structures of biological molecules

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results
Database Avail	ability		

The Protein Data Bank (PDB):

- Repository of 3-D structures of biological molecules
- Follows a strict format

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results
Database Avail	ability		

The Protein Data Bank (PDB):

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results 000000000000000000
Database Avai	lability		

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file
- Contains many redundant, obsolete, caveat or hypothetical structures

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)
 (日)

 (日)
 (日)

 (日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results 000000000000000000
Database Avai	lability		

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file
- Contains many redundant, obsolete, caveat or hypothetical structures

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• 36,837 entries on May 30, 2006.

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results 000000000000000000
Database Avai	lability		

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file
- Contains many redundant, obsolete, caveat or hypothetical structures

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• 36,837 entries on May 30, 2006.

Perl:

Biochemical Background 000000000	Pattern Recognition	Problem Representation	Experimental Results
Database Availa	bility		

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file
- Contains many redundant, obsolete, caveat or hypothetical structures
- 36,837 entries on May 30, 2006.

Perl:

• A leading language for processing structured flat files

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Biochemical Background	Pattern Recognition	Problem Representation ●00000	Experimental Results 000000000000000000
Database Avai	lability		

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file
- Contains many redundant, obsolete, caveat or hypothetical structures
- 36,837 entries on May 30, 2006.

Perl:

- A leading language for processing structured flat files
- BioPerl: bioinformatics data processing and analysis

Biochemical Background 000000000	Pattern Recognition	Problem Representation	Experimental Results
Database Availa	bility		

- Repository of 3-D structures of biological molecules
- Follows a strict format
- Stored as a sequence of records in a flat file
- Contains many redundant, obsolete, caveat or hypothetical structures
- 36,837 entries on May 30, 2006.

Perl:

- A leading language for processing structured flat files
- BioPerl: bioinformatics data processing and analysis
- Parse a PDB file, represent it as an object and query its fields.

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Data Mining			

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)
 (日)

 (日)
 (日)

 (日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

- Covalently bonded to the protein (glycoprotein)
- Floating around in the medium
- Docking to a binding-site

Biochemical Background	Pattern Recognition	Problem Representation ○●○○○○	Experimental Results
Data Mining			

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- Covalently bonded to the protein (glycoprotein)
- Floating around in the medium
- Docking to a binding-site

Biochemical Background	Pattern Recognition	Problem Representation 0●0000	Experime 000000
Data Mining			

tal Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- X Covalently bonded to the protein (glycoprotein)
- Floating around in the medium
- Docking to a binding-site

Biochemical Background	Pattern Recognition	Problem Representation ○●○○○○	Experimental Results
Data Mining			

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- X Covalently bonded to the protein (glycoprotein)
- Floating around in the medium
- Docking to a binding-site

Biochemical Background 000000000	Pattern Recognition	Problem Representation 0●0000	Experimental Results
Data Mining			

- X Covalently bonded to the protein (glycoprotein)
- X Floating around in the medium
- Docking to a binding-site

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
000000000		0●0000	000000000000000000
Data Mining			

- X Covalently bonded to the protein (glycoprotein)
- ${\sf X}\,$ Floating around in the medium
- Docking to a binding-site

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
000000000		0●0000	000000000000000000
Data Mining			

- X Covalently bonded to the protein (glycoprotein)
- ${\sf X}\,$ Floating around in the medium
- $\checkmark\,$ Docking to a binding-site

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Data Mining			

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

In a PDB file, glucose can be:

- X Covalently bonded to the protein (glycoprotein)
- ${\sf X}\,$ Floating around in the medium
- $\sqrt{}$ Docking to a binding-site

PDB contains:

• 28,353 entries in October 2004

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
000000000		0●0000	000000000000000000
Data Mining			

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

In a PDB file, glucose can be:

- X Covalently bonded to the protein (glycoprotein)
- ${\sf X}\,$ Floating around in the medium
- $\sqrt{}$ Docking to a binding-site

PDB contains:

- 28,353 entries in October 2004
- 331 entries containing glucose

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
000000000		0●0000	000000000000000000
Data Mining			

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

In a PDB file, glucose can be:

- X Covalently bonded to the protein (glycoprotein)
- X Floating around in the medium
- $\sqrt{}$ Docking to a binding-site

PDB contains:

- 28, 353 entries in October 2004
- 331 entries containing glucose
- 59 resolved protein-glucose binding sites

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
000000000		0●0000	000000000000000000
Data Mining			

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

In a PDB file, glucose can be:

- X Covalently bonded to the protein (glycoprotein)
- X Floating around in the medium
- $\sqrt{}$ Docking to a binding-site

PDB contains:

- 28, 353 entries in October 2004
- 331 entries containing glucose
- 59 resolved protein-glucose binding sites
- 29 distinct protein-glucose binding sites

Pattern Recognition

Problem Representation

Experimental Results

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Solution Approach



Figure: The classifier algorithm outline.

Biochemical Background	Pattern Recognition	Problem Representation ○○●000	Experimental Results
Solution Approa	ch		

Non-glucose sites

Binding sites that do not bind glucose



Figure: The classifier algorithm outline.

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (H)

 (H)

00000000	0000	00000	0000000000000000	
Solution Approach				

Non-glucose sites

- Binding sites that do not bind glucose
- 2 Random protein surface regions



Figure: The classifier algorithm outline.

◆□> <□> <=> <=> <=> <=> <=> <=> <=>

Pattern Recognition

Problem Representation $\circ \circ \circ \circ \circ \circ \circ$

Experimental Results

Binding Site Representation



Figure: Glucose bound to a hydrolase, PDB entry 118A.

Pattern Recognition

Problem Representation $\circ \circ \circ \circ \circ \circ \circ$

Experimental Results

Binding Site Representation



 Sphere of 10 Å centered at the ligand

Figure: Glucose bound to a hydrolase, PDB entry 118A.

Pattern Recognition

Problem Representation 000000

Experimental Results

Binding Site Representation



- Sphere of 10 Å centered at the ligand
- Pyranose ring centroid

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Figure: Glucose bound to a hydrolase, PDB entry 118A.

Pattern Recognition

Problem Representation $\circ \circ \circ \circ \circ \circ \circ$

Experimental Results

Binding Site Representation



Figure: Glucose bound to a hydrolase, PDB entry 118A.

- Sphere of 10 Å centered at the ligand
- Pyranose ring centroid
- Divided into 8 concentric layers

Pattern Recognition

Problem Representation $\circ \circ \circ \circ \circ \circ \circ$

Experimental Results

Binding Site Representation



Figure: Glucose bound to a hydrolase, PDB entry 118A.

- Sphere of 10 Å centered at the ligand
- Pyranose ring centroid
- Divided into 8 concentric layers
- Gather features for each layer

(日)
 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (日)

 (H)

 (H)

Biochemical Background	Pattern Recognition	Problem Representation ○○○○●○	Experimental Results
Atomic Features	:		

Charge

・ロト < 団ト < 三ト < 三ト < 三日 < つへの

Biochemical Background	Pattern Recognition	Problem Representation ○○○○●○	Experimental Results
Atomic Features			

Charge

O Hydrogen Bond


Biochemical Background	Pattern Recognition	Problem Representation ○○○○●○	Experimental Results
Atomic Features			

Charge

Ø Hydrogen Bond

O Hydrophobicity



Biochemical Background	Pattern Recognition	Problem Representation ○○○○●○	Experimental Results
Atomic Features			

ChargePositive

O Hydrogen Bond

O Hydrophobicity

Biochemical Background	Pattern Recognition	Problem Representation ○○○○●○	Experimental Results
Atomic Features			

Charge

Positive

Ø Neutral

O Hydrogen Bond

O Hydrophobicity

Biochemical Background	Pattern Recognitio

Problem Representation

Experimental Results

Atomic Features

Charge

- Positive
- Ø Neutral
- 8 Negative
- O Hydrogen Bond

Output State of the state of

Biochemical	Background

Problem Representation 000000

Experimental Results

Atomic Features

- Positive
- O Neutral
- 8 Negative
- e Hydrogen Bond
 - Hydrogen-bonding
- Output State of the state of

Biochemical	Background

Problem Representation 000000

Experimental Results

Atomic Features

- Positive
- Ø Neutral
- 8 Negative
- Ø Hydrogen Bond
 - Hydrogen-bonding
 - Non hydrogen-bonding
- Output State of the state of

Biochemical	Background

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Atomic Features

- Positive
- Ø Neutral
- 8 Negative
- Ø Hydrogen Bond
 - Hydrogen-bonding
 - Non hydrogen-bonding
- Output State of the state of
 - Hydrophobic

Biochemical	Background

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Atomic Features

- Positive
- Ø Neutral
- 8 Negative
- Ø Hydrogen Bond
 - Hydrogen-bonding
 - Non hydrogen-bonding
- Output State of the state of
 - Hydrophobic
 - Ø Neutral

Biochemical	Background

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Atomic Features

- Positive
- Ø Neutral
- 8 Negative
- Ø Hydrogen Bond
 - Hydrogen-bonding
 - Non hydrogen-bonding
- Output State St
 - Hydrophobic
 - 2 Neutral
 - O Hydrophilic

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schem	les		

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schemes	5		

- Planar polar residues
- Form an extensive hydrogen-bond network with glucose

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schem	es		

- Planar polar residues
- Form an extensive hydrogen-bond network with glucose
- Remaining residues

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Ba	ackground	Pattern Recognition	Problem Represent	ation Experimental Results
Residu	e Schemes	;		
Simp	olified2			
٠	Planar pola	r residues		

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background 000000000	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schem	ies		

- Planar polar residues
- Aromatic residues
- Ring containing residues
- Stack against pyranose ring thru hydrophobic and van der Waals forces.

Asparagine	Glutamine	Arginine	Glutamate	
Aspartate	Tyrosine	Tryptophan	Phenylalanine	
Histidine	Serine	Threonine	Glycine	
Proline	Cysteine	Lysine	Alanine	
Valine	Leucine	Isoleucine	Methionine	

Biochemical	Background	

Problem Representation 00000

Experimental Results

Residue Schemes

- Planar polar residues
- Aromatic residues
- Ring containing residues
- Stack against pyranose ring thru hydrophobic and van der Waals forces.
- Remaining residues

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schen	nes		
Simplified3			

۲	Ρ	lanar	ро	lar	resid	ues	

Asparagine	Glutamine	Arginine	Glutamate	
Aspartate	Tyrosine	Tryptophan	Phenylalanine	
Histidine	Serine	Threonine	Glycine	
Proline	Cysteine	Lysine	Alanine	
Valine	Leucine	Isoleucine	Methionine	

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results 00000000000000000	
Residue Schemes				

- Planar polar residues
- Aromatic residues
- Histidine is reported to interact with glucose in a similar way to aromatic residues.

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results	
Residue Schemes				

- Planar polar residues
- Aromatic residues
- Histidine is reported to interact with glucose in a similar way to aromatic residues.
- Remaining residues

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background 000000000	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schen	nes		
Detailed			

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Schen	nes		
Detailed1			

|--|

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results	
Residue Schemes				

- Aromatic residues
- Neutral residues

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical Background 000000000	Pattern Recognition	Problem Representation ○○○○○●	Experimental Results
Residue Scheme	20		

- Aromatic residues
- Neutral residues
- Carboxylate residues (negatively charged)

Asparagine	Glutamine	Arginine	Glutamate	
Aspartate	Tyrosine	Tryptophan	Phenylalanine	
Histidine	Serine	Threonine	Glycine	
Proline	Cysteine	Lysine	Alanine	
Valine	Leucine	Isoleucine	Methionine	

Biochemical	Background

Problem Representation 00000

Experimental Results

Residue Schemes

- Aromatic residues
- Neutral residues
- Carboxylate residues (negatively charged)
- Positively charged residues

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine Alanine	
Valine	Leucine	Isoleucine Methionir	

Biochemical	Background

Problem Representation

Experimental Results

Residue Schemes

- Aromatic residues
- Neutral residues
- Carboxylate residues (negatively charged)
- Positively charged residues
- Aliphatic residues

Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical	Background

Problem Representation 00000

Experimental Results

Residue Schemes

- Aromatic residues
- Neutral residues
- Carboxylate residues (negatively charged)
- Positively charged residues
- Aliphatic residues
- Histidine

• • • • • • • • • • • • • • • • • • • •			
Asparagine	Glutamine	Arginine	Glutamate
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Biochemical	Background

Problem Representation 00000

Experimental Results

Residue Schemes

- Aromatic residues
- Neutral residues
- Carboxylate residues (negatively charged)
- Positively charged residues
- Aliphatic residues

Asparagine	Glutamine	Arginine Glutama	
Aspartate	Tyrosine	Tryptophan	Phenylalanine
Histidine	Serine	Threonine	Glycine
Proline	Cysteine	Lysine	Alanine
Valine	Leucine	Isoleucine	Methionine

Pattern Recognition

Problem Representation 000000

Experimental Results

Experimental Results

What determines a binding-site:



Pattern Recognition

Problem Representation 000000

Experimental Results

Experimental Results

What determines a binding-site:

Amino acids biochemical properties

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Experimental Results

What determines a binding-site:

- Amino acids biochemical properties
- 2 Amino acids spatial arrangement

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Experimental Results

What determines a binding-site:

- Amino acids biochemical properties
- Amino acids spatial arrangement

The approach consists of:

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Experimental Results

What determines a binding-site:

- Amino acids biochemical properties
- Amino acids spatial arrangement

The approach consists of:

Learning phase

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Experimental Results

What determines a binding-site:

- Amino acids biochemical properties
- Amino acids spatial arrangement

The approach consists of:

- Learning phase
- Peature selection

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

Experimental Results

What determines a binding-site:

- Amino acids biochemical properties
- 2 Amino acids spatial arrangement

The approach consists of:

- Learning phase
- Peature selection
- Testing phase

Pattern Recognition

Problem Representation

Experimental Results

Atomic Combinations Results

Table: Atomic properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Charge	14.55%	14.55%	78.18%
H-Bond	21.82%	16.36%	92.73%
Hydrophobicity	21.82%	20.00%	92.73%
$Charge + H\operatorname{-}Bond$	14.55%	14.55%	89.09%
Charge + Hydro	12.73%	14.55%	47.27%
H-Bond + Hydro	21.82%	18.18%	100%
Charge + H-Bond + Hydro	16.36%	16.36%	60.00%

Pattern Recognition

Problem Representation

Experimental Results

Atomic Combinations Results

• Charge outperforms hydrophobicity and hydrogen bond

Table: Atomic properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Charge	14.55%	14.55%	78.18%
H-Bond	21.82%	16.36%	92.73%
Hydrophobicity	21.82%	20.00%	92.73%
$Charge + H\operatorname{-}Bond$	14.55%	14.55%	89.09%
Charge + Hydro	12.73%	14.55%	47.27%
H-Bond + Hydro	21.82%	18.18%	100%
Charge + H-Bond + Hydro	16.36%	16.36%	60.00%

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Pattern Recognition

Problem Representation

Experimental Results

Atomic Combinations Results

- Charge outperforms hydrophobicity and hydrogen bond
- Charge, linked with another property, yields similar or slightly better results

Table: Atomic properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Charge	14.55%	14.55%	78.18%
H-Bond	21.82%	16.36%	92.73%
Hydrophobicity	21.82%	20.00%	92.73%
$Charge + H\operatorname{-Bond}$	14.55%	14.55%	89.09%
Charge + Hydro	12.73%	14.55%	47.27%
H-Bond + Hydro	21.82%	18.18%	100%
Charge + H-Bond + Hydro	16.36%	16.36%	60.00%

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □
Pattern Recognition

Problem Representation

Experimental Results

Atomic Combinations Results

 $\bullet~$ "Charge +~ Hydro" is the best atomic combination

Table: Atomic properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Charge	14.55%	14.55%	78.18%
H-Bond	21.82%	16.36%	92.73%
Hydrophobicity	21.82%	20.00%	92.73%
$Charge + H\operatorname{-}Bond$	14.55%	14.55%	89.09%
Charge + Hydro	12.73%	14.55%	47.27%
H-Bond $+$ Hydro	21.82%	18.18%	100%
Charge + H-Bond + Hydro	16.36%	16.36%	60.00%

Pattern Recognition

Problem Representation 000000

Experimental Results

Atomic Combinations Results

• kNN and SVM give similar results

Table: Atomic properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Charge	14.55%	14.55%	78.18%
H-Bond	21.82%	16.36%	92.73%
Hydrophobicity	21.82%	20.00%	92.73%
$Charge + H\operatorname{-}Bond$	14.55%	14.55%	89.09%
Charge + Hydro	12.73%	14.55%	47.27%
H-Bond $+$ Hydro	21.82%	18.18%	100%
Charge + H-Bond + Hydro	16.36%	16.36%	60.00%

Biochemical	Background

Problem Representation 000000

Experimental Results

Hydrogen Bond Property

• Hydrogen bond is a primary factor in glucose binding

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Hydrogen Bonc	l Property		

- Hydrogen bond is a primary factor in glucose binding
- Hydrogen bond preforms poorly

Biochemical	Background

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Hydrogen Bond Property

- Hydrogen bond is a primary factor in glucose binding
- Hydrogen bond preforms poorly

Hypothesis

Negative set composed mainly of non-glucose binding sites

Biochemical	Background

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Hydrogen Bond Property

- Hydrogen bond is a primary factor in glucose binding
- Hydrogen bond preforms poorly

Hypothesis

- Negative set composed mainly of non-glucose binding sites
- Ø Most binding sites rely on hydrogen bonds

Biochemical	Background

Problem Representation

Experimental Results

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Hydrogen Bond Property

- Hydrogen bond is a primary factor in glucose binding
- Hydrogen bond preforms poorly

Hypothesis

- Negative set composed mainly of non-glucose binding sites
- Ø Most binding sites rely on hydrogen bonds
- O Therefore hydrogen bond is a bad discriminating factor

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Hydrogen Bond Property

Hypothesis

- Negative set composed mainly of non-glucose binding sites
- Ø Most binding sites rely on hydrogen bonds
- O Therefore hydrogen bond is a bad discriminating factor

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Hydrogen Bond Property

Hypothesis

- Negative set composed mainly of non-glucose binding sites
- Ø Most binding sites rely on hydrogen bonds
- O Therefore hydrogen bond is a bad discriminating factor

To validate this hypothesis:

Pattern Recognition

Problem Representation

Experimental Results

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

Hydrogen Bond Property

Hypothesis

- Negative set composed mainly of non-glucose binding sites
- Ø Most binding sites rely on hydrogen bonds
- On the provide the second s

To validate this hypothesis:

Negative set composed of random protein surface regions

Pattern Recognition

Problem Representation

Experimental Results

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Hydrogen Bond Property

Hypothesis

- Negative set composed mainly of non-glucose binding sites
- Ø Most binding sites rely on hydrogen bonds
- On the provide the second s

To validate this hypothesis:

- Negative set composed of random protein surface regions
- 2 Hydrogen bond should outperform charge and hydrophobicity

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Hydrogen Bond	Property		

To validate this hypothesis:

- Negative set composed of random protein surface regions
- Ø Hydrogen bond should outperform charge and hydrophobicity

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

Biochemical Background 000000000	Pattern Recognition	Problem Representation	Experimental Results
Hvdrogen Bon	d Property		

To validate this hypothesis:

- Negative set composed of random protein surface regions
- Ø Hydrogen bond should outperform charge and hydrophobicity

The results confirm the hypothesis:

Table: Classifier training using an exclusively non-binding sites negative set.

Properties	<i>k</i> NN	SVM	SV
Charge	05.26%	05.26%	73.68%
H-Bond	03.51%	03.51%	61.40%
Hydrophobicity	05.26%	05.26%	68.42%

Pattern Recognition

Problem Representation

Experimental Results

Residue Schemes Comparison

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Residue Schemes	Comparison		

• Detailed schemes outperform simplified schemes

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical Background 000000000	Pattern Recognition	Problem Representation	Experimental Results
Residue Schemes	Comparison		

- Detailed schemes outperform simplified schemes
- Simplified schemes have better generalization

Table: Comparison of the different residue schemes.

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Residue Schemes	Comparison		

• Detailed schemes yield similar results to the atomic properties

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Residue Schemes	Comparison		

- Detailed schemes yield similar results to the atomic properties
- Incorporate in themselves some of the atomic knowledge?

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical Background 000000000	Pattern Recognition	Problem Representation	Experimental Results
Residue Schemes	Comparison		

• "Simplified2" scheme gives similar results to "simplified1" while adding the aromatic subgroup

Table: Comparison of the different residue schemes.

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

$\frac{1}{2}$	Dociduo Schon	nos Composison		
REACHARDER REACARDITION REACARDITION REACARDITION	000000000	0000	000000	OOOOOOOOOOOOOOOOO

- "Simplified2" scheme gives similar results to "simplified1" while adding the aromatic subgroup
- Aromatic residues play an important role in glucose binding

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Pasidua Sahamaa	Comparison		
Biochemical Background 00000000	Pattern Recognition 0000	Problem Representation	Experimental Results

- Residue Schemes Comparison
 - "Simplified2" scheme gives similar results to "simplified1" while adding the aromatic subgroup
 - Aromatic residues play an important role in glucose binding
 - Aromatic residues do not seem to contribute to the classification accuracy

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical Background 000000000	Pattern Recognition	Problem Representation	Experimental Results
Residue Schemes	Comparison		

• "Simplified3" improves the "simplified2" classification rate by one hit

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

	<u> </u>		
Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results

Residue Schemes Comparison

- "Simplified3" improves the "simplified2" classification rate by one hit
- Adding histidine to the aromatic subgroup gives slightly better results

Residue scheme	<i>k</i> NN	SVM	SV
Simplified1	16.36%	18.18%	76.36%
Simplified2	18.18%	18.18%	70.91%
Simplified3	16.36%	16.36%	70.91%
Detailed1	16.36%	14.55%	81.82%
Detailed2	21.82%	14.55%	81.82%

Biochemical	Background

Problem Representation

Experimental Results

Water and Ion inclusion

• Ordered water molecules and ions play a role in glucose binding

Table: Testing the importance of water and ions to glucose specificity

Mode	Properties	<i>k</i> NN	SVM	SV
Water	Charge + Hydro	12.73%	14.55%	47.27%
+ ions	Residue + Charge + Hydro	20.00%	09.09%	100%
lons	Charge + Hydro	14.55%	16.36%	83.64%
	Residue $+$ Charge $+$ Hydro	20.00%	10.91%	100%
Water	Charge + Hydro	16.36%	18.18%	83.64%
	Residue + Charge + Hydro	20.00%	10.91%	100%
	Charge + Hydro	14.55%	16.36%	52.73%
	Residue + Charge + Hydro	20.00%	12.73%	100%

Biochemical	Background

Problem Representation

Experimental Results

Water and Ion inclusion

- Ordered water molecules and ions play a role in glucose binding
- We test and validate this hypothesis

Table: Testing the importance of water and ions to glucose specificity

Mode	Properties	<i>k</i> NN	SVM	SV
Water	Charge + Hydro	12.73%	14.55%	47.27%
+ ions	Residue + Charge + Hydro	20.00%	09.09%	100%
lons	Charge + Hydro	14.55%	16.36%	83.64%
	Residue $+$ Charge $+$ Hydro	20.00%	10.91%	100%
Water	Charge + Hydro	16.36%	18.18%	83.64%
	Residue + Charge + Hydro	20.00%	10.91%	100%
	Charge + Hydro	14.55%	16.36%	52.73%
	Residue + Charge + Hydro	20.00%	12.73%	100%

Pattern Recognition

Problem Representation

Experimental Results

Atomic Properties Feature Selection

Table: Classifiers performance on atomic data using feature selection.

Property	RF	Feat.	<i>k</i> NN	SVM	Sensi-	Speci-	SV
		Nb.	Error	Error	tivity	ficity	
Charge	false	24	14.55%	14.55%	96.55%	73.08%	78.18%
	true	6	09.09%	05.45%	93.10%	96.15%	41.82%
H-Bond	false	16	21.82%	16.36%	86.21%	80.77%	92.73%
	true	5	09.09%	07.27%	96.55%	88.46%	16.36%
Hydro	false	24	21.82%	20.00%	79.31%	80.77%	92.73%
	true	5	09.09%	10.91%	96.55%	84.62%	34.55%

Biochemical	Background

Problem Representation

Experimental Results

Atomic Properties Feature Selection

Definition

Sensitivity: Ability to detect true positives (TP/P)Specificity: Ability to reject true negatives (TN/N)

Table: Classifiers performance on atomic data using feature selection.

Property	RF	Feat.	<i>k</i> NN	SVM	Sensi-	Speci-	SV
		Nb.	Error	Error	tivity	ficity	
Charge	false	24	14.55%	14.55%	96.55%	73.08%	78.18%
	true	6	09.09%	05.45%	93.10%	96.15%	41.82%
H-Bond	false	16	21.82%	16.36%	86.21%	80.77%	92.73%
	true	5	09.09%	07.27%	96.55%	88.46%	16.36%
Hydro	false	24	21.82%	20.00%	79.31%	80.77%	92.73%
	true	5	09.09%	10.91%	96.55%	84.62%	34.55%

Biochemical Background			Pattern Recognition			Problem Representation	Experimental Results
Δ.		D	 _		C I	1. A.	

Atomic Properties Feature Selection

• Classification error decreases

Table: Classifiers performance on atomic data using feature selection.

Property	RF	Feat.	<i>k</i> NN	SVM	Sensi-	Speci-	SV
		Nb.	Error	Error	tivity	ficity	
Charge	false	24	14.55%	14.55%	96.55%	73.08%	78.18%
	true	6	09.09%	05.45%	93.10%	96.15%	41.82%
H-Bond	false	16	21.82%	16.36%	86.21%	80.77%	92.73%
	true	5	09.09%	07.27%	96.55%	88.46%	16.36%
Hydro	false	24	21.82%	20.00%	79.31%	80.77%	92.73%
	true	5	09.09%	10.91%	96.55%	84.62%	34.55%

Biochemical Background				Pattern Recognition			Problem Representation	Experimental Results
		5	-			<u> </u>		

Atomic Properties Feature Selection

- Classification error decreases
- SV percentage decreases (generalization capacity increases)

Tables	Class!!!						f	
Table:	Classifiers	performance	on	atomic	data	using	teature	selection.

Property	RF	Feat.	<i>k</i> NN	SVM	Sensi-	Speci-	SV
		Nb.	Error	Error	tivity	ficity	
Charge	false	24	14.55%	14.55%	96.55%	73.08%	78.18%
	true	6	09.09%	05.45%	93.10%	96.15%	41.82%
H-Bond	false	16	21.82%	16.36%	86.21%	80.77%	92.73%
	true	5	09.09%	07.27%	96.55%	88.46%	16.36%
Hydro	false	24	21.82%	20.00%	79.31%	80.77%	92.73%
	true	5	09.09%	10.91%	96.55%	84.62%	34.55%

Pattern Recognition

Problem Representation 000000

Experimental Results

Charge Feature Selection



 Remains the best discriminating factor

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Figure: Importance of charge features

Pattern Recognition

Problem Representation 000000

Experimental Results

Charge Feature Selection



- Remains the best discriminating factor
- Negatively charged glucose binding site

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

Figure: Importance of charge features

Pattern Recognition

Problem Representation 000000

Experimental Results

Charge Feature Selection



- Remains the best discriminating factor
- Negatively charged glucose binding site

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

• Layers 1 and 2: overlap with glucose

Figure: Importance of charge features

Pattern Recognition

Problem Representation

Experimental Results

Hydrogen Bond Feature Selection



Figure: Importance of hbond features

Pattern Recognition

Problem Representation 000000

Experimental Results

Hydrogen Bond Feature Selection



- A record low SV percentage
- Both positive and negative sets rely on H-Bond

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

Figure: Importance of hbond features

Pattern Recognition

Problem Representation 000000

Experimental Results

Hydrogen Bond Feature Selection



- A record low SV percentage
- Both positive and negative sets rely on H-Bond

▲ロト ▲冊ト ▲ヨト ▲ヨト 三回 のへの

• Differ by their distribution

Figure: Importance of hbond features

Pattern Recognition

Problem Representation 000000

Experimental Results

Hydrophobicity Feature Selection



 Remains the worst discriminating factor

Figure: Importance of hydro features
Pattern Recognition

Problem Representation 000000

Experimental Results

Hydrophobicity Feature Selection



- Remains the worst discriminating factor
- Positive set: high concentration of hydrophilic atoms in layer 3

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Figure: Importance of hydro features

Pattern Recognition

Problem Representation 000000

Experimental Results

Hydrophobicity Feature Selection



- Remains the worst discriminating factor
- Positive set: high concentration of hydrophilic atoms in layer 3
- Low importance score of hydrophobic features

Figure: Importance of hydro features

Pattern Recognition

Problem Representation

Experimental Results

Simplified Residue Schemes Feature Selection



 Similar results for the 3 simplified schemes

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Figure: Importance of simplified features

Pattern Recognition

Problem Representation

Experimental Results

Simplified Residue Schemes Feature Selection



- Similar results for the 3 simplified schemes
- Importance of planar polar residues

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Figure: Importance of simplified features

Pattern Recognition

Problem Representation

Experimental Results

Simplified Residue Schemes Feature Selection



- Similar results for the 3 simplified schemes
- Importance of planar polar residues
- Low importance score of aromatic residues

Figure: Importance of simplified features

Pattern Recognition

Problem Representation

Experimental Results

Detailed Residue Schemes Feature Selection



 Similar results for both detailed schemes

Figure: Importance of detailed features

Pattern Recognition

Problem Representation

Experimental Results

Detailed Residue Schemes Feature Selection



- Similar results for both detailed schemes
- Still outperform simplified schemes.

Figure: Importance of detailed features

Pattern Recognition

Problem Representation

Experimental Results

Detailed Residue Schemes Feature Selection



- Similar results for both detailed schemes
- Still outperform simplified schemes.
- High density of the negatively charged carboxylate residues

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Figure: Importance of detailed features

Pattern Recognition

Problem Representation

Experimental Results

Combined Atomic and Residue Feature Selection

• 3.64% error for both classifiers.

Table: Classifiers performance using feature selection.

Property	RF	Feat.	<i>k</i> NN	SVM	Sensi-	Speci-	SV
		Nb.	Error	Error	tivity	ficity	
Detailed2	false	104	14.06%	07.81%	93.10%	91.43%	53.13%
+ Charge							
+ H-Bond	true	15	03.64%	03.64%	96.55%	96.15%	69.09%
+ Hydro							

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Pattern Recognition

Problem Representation 000000

Experimental Results

Selected Feature Subset

Table: The selected feature subset.

Property	Features	L1	L2	L3	L4	L5	L6	L7	L8
Charge	Negative			Х					
	Neutral	X	Х						
H-Bond	H-Bonding			X					
Hydro	Hydrophilic	Х	Х	Х					
	Neutral		Х						
	Hydrophobic							Х	X
Residues	Neutral						Х		
	Carboxylate					Х	Х		X
	Aliphatic			Х					

Biochemical	Background

Problem Representation 000000

Experimental Results

◆□> <□> <=> <=> <=> <=> <=> <=> <=>

Feature Selection Findings

• The difference in spatial configuration

Biochemical	Background

Problem Representation

Experimental Results

- The difference in spatial configuration
- Charge is the best discriminant property

Biochemical	Background

Problem Representation

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- The difference in spatial configuration
- Charge is the best discriminant property
- Discriminating atomic properties on layer 3

Biochemical	Background

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

- The difference in spatial configuration
- Charge is the best discriminant property
- Discriminating atomic properties on layer 3
- Negatively charged layer 3

Biochemical	Background

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

- The difference in spatial configuration
- Charge is the best discriminant property
- Discriminating atomic properties on layer 3
- Negatively charged layer 3
- High density of negatively charged carboxylate residues

Biochemical	Background

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

- The difference in spatial configuration
- Charge is the best discriminant property
- Discriminating atomic properties on layer 3
- Negatively charged layer 3
- High density of negatively charged carboxylate residues
- The relevance of planar polar residues

Problem Representation

Experimental Results

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- The difference in spatial configuration
- Charge is the best discriminant property
- Discriminating atomic properties on layer 3
- Negatively charged layer 3
- High density of negatively charged carboxylate residues
- The relevance of planar polar residues
- Low discrimination capacity of hydrophobic interactions and aromatic residues

Problem Representation

Experimental Results

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

- The difference in spatial configuration
- Charge is the best discriminant property
- Discriminating atomic properties on layer 3
- Negatively charged layer 3
- High density of negatively charged carboxylate residues
- The relevance of planar polar residues
- Low discrimination capacity of hydrophobic interactions and aromatic residues
- Possibly due to their one-correctly-placed relationship

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Testing Phase			

・ロト < 団ト < 三ト < 三ト < 三日 < つへの

• Parse PDB for entries newer than October 2004

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Testing Phase			

◆□> <□> <=> <=> <=> <=> <=> <=> <=>

- Parse PDB for entries newer than October 2004
- 7 distinct glucose binding sites

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Testing Phase			

- Parse PDB for entries newer than October 2004
- 7 distinct glucose binding sites

Table: Testing phase results.

Classifier	TΡ	FP	TN	FN	Error	Sensi-	Speci-	Preci-
						tivity	ficity	sion
SVM	6	0	15	1	04.55%	85.71%	100%	100%
<i>k</i> NN	6	2	13	1	13.64%	85.71%	86.67%	75.00%
COTRAN	94	27	633	12	05.09%	88.68%	95.91%	77.69%

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Testing Phase			

- Parse PDB for entries newer than October 2004
- 7 distinct glucose binding sites
- Compare with COTRAN, a galactose-binding site classifier

Table: Testing phase results.

Classifier	TΡ	FP	TN	FN	Error	Sensi-	Speci-	Preci-
						tivity	ficity	sion
SVM	6	0	15	1	04.55%	85.71%	100%	100%
<i>k</i> NN	6	2	13	1	13.64%	85.71%	86.67%	75.00%
COTRAN	94	27	633	12	05.09%	88.68%	95.91%	77.69%

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Testing Phase			

- Parse PDB for entries newer than October 2004
- 7 distinct glucose binding sites
- Compare with COTRAN, a galactose-binding site classifier
- Precision: ability to reject false positives.

Classifier	TΡ	FP	ΤN	FN	Error	Sensi-	Speci-	Preci-
						tivity	ficity	sion
SVM	6	0	15	1	04.55%	85.71%	100%	100%
<i>k</i> NN	6	2	13	1	13.64%	85.71%	86.67%	75.00%
COTRAN	94	27	633	12	05.09%	88.68%	95.91%	77.69%

Table: Testing phase results.

Pattern Recognition

Problem Representation 000000

Experimental Results

Future Work

• Add more experiments

Pattern Recognition

Problem Representation

Experimental Results

- Add more experiments
- Tune fixed parameters: sphere radius, layers number ...

Pattern Recognition

Problem Representation

Experimental Results

- Add more experiments
- Tune fixed parameters: sphere radius, layers number ...
- Wet lab experiments

Pattern Recognition

Problem Representation

Experimental Results

- Add more experiments
- Tune fixed parameters: sphere radius, layers number ...
- Wet lab experiments
- Web based interface

Pattern Recognition

Problem Representation 000000

Experimental Results

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- Add more experiments
- Tune fixed parameters: sphere radius, layers number ...
- Wet lab experiments
- Web based interface
- Extend to a glucose binding-site finder

▲ロト ▲帰ト ▲ヨト ▲ヨト 三回日 のの⊙

- Add more experiments
- Tune fixed parameters: sphere radius, layers number ...
- Wet lab experiments
- Web based interface
- Extend to a glucose binding-site finder
- Extend this study to other sugars

000000000	0000	
Acknowledgme	ents	

• Dr. Khuri for initiating and constantly supervising this work from abroad

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Acknowledgme	ents		

O

- Dr. Khuri for initiating and constantly supervising this work from abroad
- Mr. Al-Ali for providing all the data and closely assisting me

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Acknowledgmen ⁻	ts		

- Dr. Khuri for initiating and constantly supervising this work from abroad
- Mr. Al-Ali for providing all the data and closely assisting me

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

• Dr. Keyrouz for supervising my thesis and supporting me

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Acknowledgmer	nts		

- Dr. Khuri for initiating and constantly supervising this work from abroad
- Mr. Al-Ali for providing all the data and closely assisting me

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

- Dr. Keyrouz for supervising my thesis and supporting me
- Dr. Kachfe for his many discussions

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Acknowledgmer	nts		

- Dr. Khuri for initiating and constantly supervising this work from abroad
- Mr. Al-Ali for providing all the data and closely assisting me

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

- Dr. Keyrouz for supervising my thesis and supporting me
- Dr. Kachfe for his many discussions
- The committee members

Biochemical Background	Pattern Recognition	Problem Representation	Experimental Results
Acknowledgmen	ts		

- Dr. Khuri for initiating and constantly supervising this work from abroad
- Mr. Al-Ali for providing all the data and closely assisting me

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
(日)

(日)
(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)

(日)
</p

- Dr. Keyrouz for supervising my thesis and supporting me
- Dr. Kachfe for his many discussions
- The committee members
- The audience

Amino Acids Table

Table: Standard amino acids names, three- and one-letter abbreviations. They are sorted according to biological convention based on the size and properties of the side chain.

Name	3-L	1-L	Name	3-L	1-L
Glycine	Gly	G	Cysteine	Cys	С
Alanine	Ala	A	Serine	Ser	S
Valine	Val	V	Threonine	Thr	Т
Leucine	Leu	L	Aspartate	Asp	D
Isoleucine	lle		Glutamate	Glu	E
Proline	Pro	Р	Histidine	His	Н
Phenylalanine	Phe	F	Lysine	Lys	K
Tyrosine	Tyr	Y	Arginine	Arg	R
Tryptophan	Trp	W	Asparagine	Asn	Ν
Methionine	Met	М	Glutamine	Gln	Q
Atomic Features

Atom Type	Functional Group	Location	Residue	PDB Atom Symbol	Chrg	Hydrophob	H Bond
Oxygen	Amide peptide linkage	Backbone	All	0	0	-1	H Bond
Oxygen	Carboxyl – C terminus	Backbone	All	OXT	-ve	-1	H Bond
Oxygen	Carboxyl	Side Chain	GLU	OE1	-ve	-1	H Bond
Oxygen	Carboxyl	Side Chain	GLU	OE2	-ve	-1	H Bond
Oxygen	Carboxyl	Side Chain	ASP	OD1	-ve	-1	H Bond
Oxygen	Carboxyl	Side Chain	ASP	OD2	-ve	-1	H Bond
Oxygen	Amide	Side Chain	GLN	OE1	0	-1	H Bond
Oxygen	Amide	Side Chain	ASN	OD1	0	-1	H Bond
Oxygen	Hydroxyl	Side Chain	SER	OG	0	-1	H Bond
Oxygen	Hydroxyl	Side Chain	THR	OG1	0	-1	H Bond
Oxygen	Hydroxyl - Phenolic	Side Chain	TYR	он	0	-1	H Bond
Nitrogen	Amide peptide linkage	Backbone	All except PRO	N	0	-1	H Bond
Nitrogen	Amide peptide linkage	Backbone	PRO	N	0	-1	
Nitrogen	Amide	Side Chain	GLN	NE2	0	-1	H Bond
Nitrogen	Amide	Side Chain	ASN	ND2	0	-1	H Bond
Nitrogen	Amine	Side Chain	LYS	NZ	+ve	-1	H Bond
Nitrogen	Guanidino	Side Chain	ARG	NE	+ve	-1	
Nitrogen	Guanidino	Side Chain	ARG	NH1	+ve	-1	H Bond
Nitrogen	Guanidino	Side Chain	ARG	NH2	+ve	-1	H Bond
Nitrogen	Imidazole	Side Chain	HIS	ND1	0	-1	
Nitrogen	Imidazole	Side Chain	HIS	NE2	0	-1	H Bond
Nitrogen	Indole	Side Chain	TRP	NE1	0	0	
Carbon	Amide peptide linkage	Backbone	All	с	0	0	
Carbon	C-alpha	Backbone	All	CA	0	0	
Carbon	Aliphatic – neutral	Side Chain	Set A (See below)	CB, CG, CD, CE	0	0	
Carbon	Aliphatic – hydrophobic	Side Chain	LEU, VAL, ILE, MET	CB, CG, CD, CE	0	1	
Carbon	Aliphatic – Branch	Side Chain	LEU, VAL, ILE	CG1, CG2, CD1, CD2, CD1	0	1	
Carbon	Phenyl - aromatic	Side Chain	PHE, TYR	CG,CD1, CD2, CE1, CE2, CZ	0	1	
Carbon	Imidazole	Side Chain	HIS	CG, CD2, CE1	0	1	
Carbon	Aromatic	Side Chain	TRP	CG,CD1, CD2,	0	1	
Carbon	Aromatic	Side Chain	TRP	CE2, CE3, CZ2, CZ3, CH2	0	1	
Sulfur	Sulfhydril	Side Chain	CYS	SG	0	-1	H Bond
Sulfur	Thioether	Side Chain	MET	SD	0	0	
Oxygen	Sulfate	HET Group	SO4	01, 02, 03, 04	-ve	-1	H Bond
Oxygen	Phosphate	HET Group	2HP	01, 02, 03, 04	-ve	-1	H Bond
Oxygen	Water	HET Group	нон	0	0	-1	H Bond
Calcium	Ion	HET Group	CA	CA	+ve	-1	H Bond
Magnesium	Ion	HET Group	MG	MG	+ve	-1	H Bond
Zinc	lon	HET Group	ZN	ZN	+1/0	-1	H Bond

Set A = ALA, SER, THR, CYS, ASP, ASN, GLU, GLN, ARG, LYS, PRO

Table: Atomic and residue properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Residue (detailed1 scheme)	16.36%	14.55%	81.82%
Residue + Charge	18.18%	10.91%	100%
Residue + H-Bond	16.36%	10.91%	94.55%
Residue + Hydro	16.36%	10.91%	90.91%
Residue + Charge + H-Bond	16.36%	10.91%	100%
Residue + Charge + Hydro	20.00%	09.09%	100%
Residue + H-Bond + Hydro	18.18%	10.91%	98.18%
Residue + Charge + H-Bond + Hydro	18.18%	09.09%	100%

• SVM performs better

Table: Atomic and residue properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Residue (detailed1 scheme)	16.36%	14.55%	81.82%
Residue + Charge	18.18%	10.91%	100%
Residue + H-Bond	16.36%	10.91%	94.55%
Residue + Hydro	16.36%	10.91%	90.91%
$Residue + Charge + H\operatorname{-Bond}$	16.36%	10.91%	100%
Residue + Charge + Hydro	20.00%	09.09%	100%
Residue + H-Bond + Hydro	18.18%	10.91%	98.18%
${\sf Residue} + {\sf Charge} + {\sf H}{\sf -}{\sf Bond} + {\sf Hydro}$	18.18%	09.09%	100%

- SVM performs better
- *k*NN deteriorates

Table: Atomic and residue properties classification error rates.

Properties	<i>k</i> NN	SVM	SV
Residue (detailed1 scheme)	16.36%	14.55%	81.82%
Residue + Charge	18.18%	10.91%	100%
Residue + H-Bond	16.36%	10.91%	94.55%
Residue + Hydro	16.36%	10.91%	90.91%
Residue + Charge + H-Bond	16.36%	10.91%	100%
Residue + Charge + Hydro	20.00%	09.09%	100%
$Residue + H\operatorname{-Bond} + Hydro$	18.18%	10.91%	98.18%
Residue + Charge + H-Bond + Hydro	18.18%	09.09%	100%

• Residue property constrains the differences between the various combinations of atomic features

Table:	Atomic and	residue	properties	classification	error rates.
--------	------------	---------	------------	----------------	--------------

Properties	<i>k</i> NN	SVM	SV
Residue (detailed1 scheme)	16.36%	14.55%	81.82%
Residue + Charge	18.18%	10.91%	100%
Residue + H-Bond	16.36%	10.91%	94.55%
Residue + Hydro	16.36%	10.91%	90.91%
Residue + Charge + H-Bond	16.36%	10.91%	100%
Residue + Charge + Hydro	20.00%	09.09%	100%
$Residue + H\operatorname{-}Bond + Hydro$	18.18%	10.91%	98.18%
Residue + Charge + H-Bond + Hydro	18.18%	09.09%	100%

• "Charge + Hydro + Residue " is the best combination

Table:	Atomic ar	nd residue	properties	classification	error rates.
--------	-----------	------------	------------	----------------	--------------

Properties	<i>k</i> NN	SVM	SV
Residue (detailed1 scheme)	16.36%	14.55%	81.82%
Residue + Charge	18.18%	10.91%	100%
Residue + H-Bond	16.36%	10.91%	94.55%
Residue + Hydro	16.36%	10.91%	90.91%
Residue + Charge + H-Bond	16.36%	10.91%	100%
Residue + Charge + Hydro	20.00%	09.09%	100%
$Residue + H\operatorname{-Bond} + Hydro$	18.18%	10.91%	98.18%
Residue + Charge + H-Bond + Hydro	18.18%	09.09%	100%

Appendix C: Residue Experiments 0000

Hydrophobicity and Charge Properties

• Charge: the best discriminating factor



Appendix C: Residue Experiments 0000

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- Charge: the best discriminating factor
- Glucose exhibits a dual hydrophobic-hydrophilic nature

- Charge: the best discriminating factor
- Glucose exhibits a dual hydrophobic-hydrophilic nature
- Both antagonist properties are involved in glucose binding

- Charge: the best discriminating factor
- Glucose exhibits a dual hydrophobic-hydrophilic nature
- Both antagonist properties are involved in glucose binding
 - Hydrophobic: Aromatic residues stack against pyranose ring

- Charge: the best discriminating factor
- Glucose exhibits a dual hydrophobic-hydrophilic nature
- Both antagonist properties are involved in glucose binding
 - Hydrophobic: Aromatic residues stack against pyranose ring
 - **2** Hydrophilic: Planar-polar residues form hydrogen-bonds

- Charge: the best discriminating factor
- Glucose exhibits a dual hydrophobic-hydrophilic nature
- Both antagonist properties are involved in glucose binding
 - Use Hydrophobic: Aromatic residues stack against pyranose ring
 - **2** Hydrophilic: Planar-polar residues form hydrogen-bonds
- May explain why hydrophobicity is the worst discriminant

- Charge: the best discriminating factor
- Glucose exhibits a dual hydrophobic-hydrophilic nature
- Both antagonist properties are involved in glucose binding
 - Hydrophobic: Aromatic residues stack against pyranose ring
 - **2** Hydrophilic: Planar-polar residues form hydrogen-bonds
- May explain why hydrophobicity is the worst discriminant
- Hydrophobicity: secondary discriminating factor

▲ロト ▲冊 ▶ ▲ヨト ▲ヨト 三回 のへの

Detailed Schemes Comparison



Figure: Comparison of detailed1 and detailed2 schemes

Residue Schemes Feature Selection

Table: Classifiers performance on residue data using feature selection.

Property	RF	Feat.	<i>k</i> NN	SVM	Sensi-	Speci-	SV
		Nb.	Error	Error	tivity	ficity	
Detailed1	false	48	16.36%	14.55%	89.66%	80.77%	81.82%
	true	19	09.09%	12.73%	93.10%	88.46%	56.36%
Detailed2	false	40	21.82%	14.55%	96.55%	73.08%	81.82%
	true	3	10.91%	12.73%	89.66%	88.46%	56.36%
Simplified1	false	16	16.36%	18.18%	96.55%	69.23%	76.36%
	true	6	12.73%	12.73%	93.10%	80.77%	67.27%
Simplified2	false	24	18.18%	18.18%	86.21%	76.92%	70.91%
	true	4	10.91%	10.91%	96.55%	80.77%	54.55%
Simplified3	false	24	16.36%	16.36%	96.55%	69.23%	70.91%
	true	6	12.73%	12.73%	96.55%	76.92%	74.55%

Appendix C: Residue Experiments

◆□> ◆□> ◆三> ◆三> 三三 のへの

The Selected Features Importance



Figure: Importance of the selected features subset