# Topic

Prediction of the burial status of transmembrane residues of helical membrane proteins
(with support vectors)

by Thorsten Will

# Overview

## Introduction

      - helical membrane proteins
      - needed definitions: burial status ,(r)SASA
      - a two step architecture: the TMX method

## Statistical methods used

      - prediction:
            - basic principles of a SVM
            - Support Vector Regression in detail

      - assessment:
            - measuring regression performance
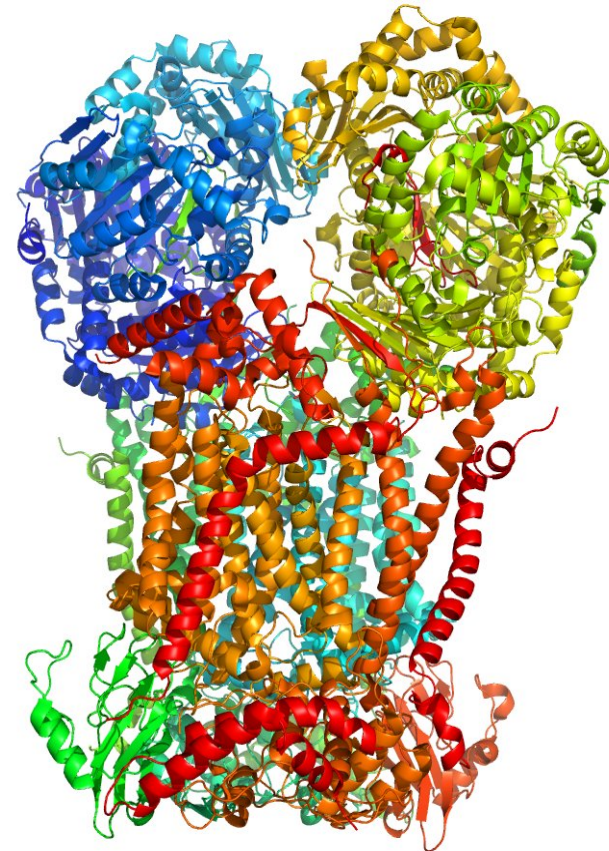            - cross-validation

## The practical part

      - the dataset and the problem
      - workflow
      - results and discussion

# Why helical membrane proteins are important

**Facts:**

- crucial role in fundamental
 cellular processes

- account for 20-30% of the
ORFs of sequenced genomes



Cytochrome bc1 complex
(respiratory chain) (from PDB: 1PP9)

# Why helical membrane proteins are interesting

**Facts:**

- crucial role in fundamental cellular processes → structure determination desirable

<span style="color:red">so far very difficult with current experimental techniques</span>

- account for 20-30% of the ORFs of sequenced genomes

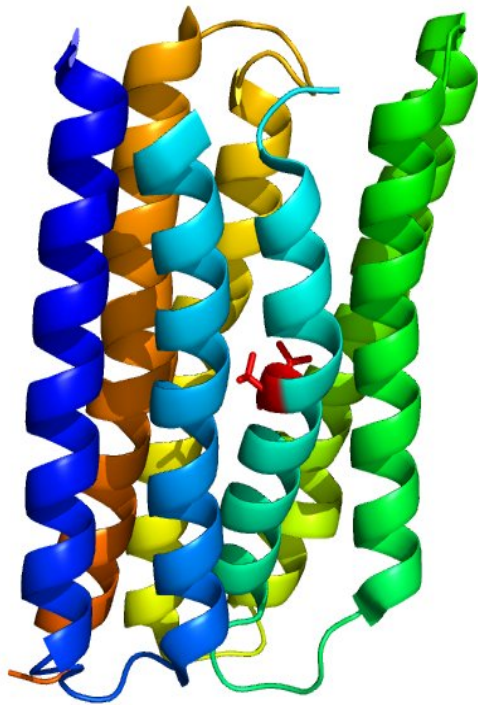Less than 1% of the proteins with known structure are HMPs !
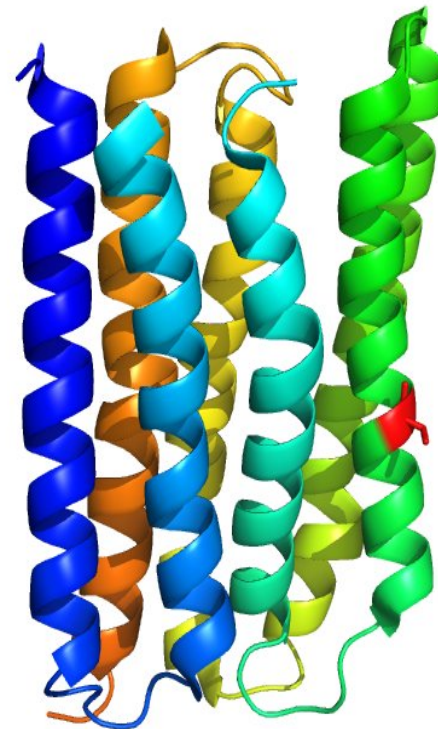
sequence-based predictors invaluable!

# What is the burial status and why is it useful

In the case of membrane proteins:
**Buried** in the protein core vs **exposed** to the membrane
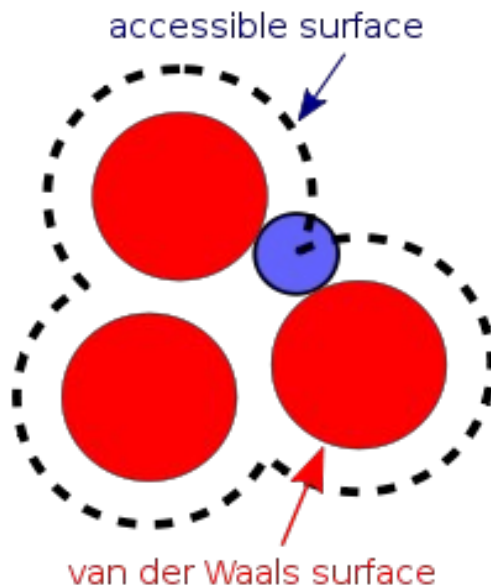


two burried threonine



an exposed serine

Example:
a bacterial rhodopsin from the dataset (from PDB: 1xio)

# Definition of the rSASA

**(S)ASA**: (Solvent) Accessible Surface Area
- the accessible area of the surface for a solvent of a specific size



accessible surface

van der Waals surface

**rSASA**: relative SASA
- normalized measurement for the SASA:
  $\rightarrow$ dividation by reference values

- rSASA = 0.00 $\rightarrow$ residue defined as burried

from Wikipedia:
http://en.wikipedia.org/wiki/Accessible_surface_area

6

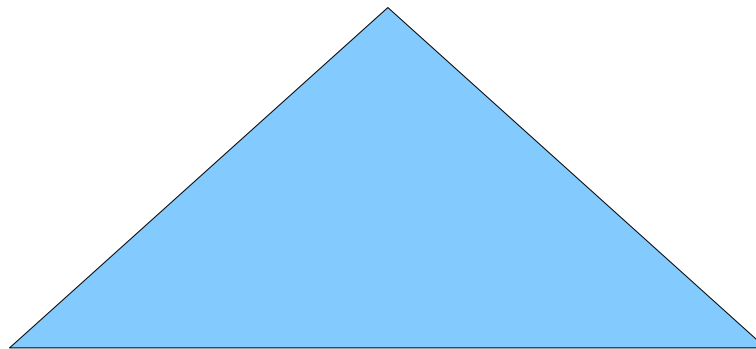few words about...

# **TMX**: TransMembrane eXposure

## A two step approach

1. positional score
2. classification

# SVM in general

**SVM**: Support Vector
Machines

**SVC**: Support
Vector Classifier

**SVR**: Support
Vector Regression

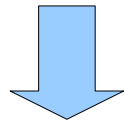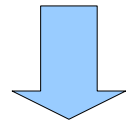# SVM basics: the formal problem simplified

SVC: classification

SVR: regression

$$\{(x_1, y_1), \ldots, (x_n, y_n)\} \subset \mathcal{X} \times \{-1, +1\}$$

Input / training data

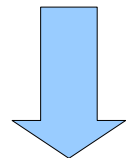$$\{(x_1, y_1), \ldots, (x_n, y_n)\} \subset \mathcal{X} \times \mathbb{R}$$

# SVM basics: the formal problem simplified

SVC: classification

SVR: regression

$$\{(x_1, y_1), \ldots, (x_n, y_n)\} \subset X \times \{-1, +1\}$$  Input / training data  $$\{(x_1, y_1), \ldots, (x_n, y_n)\} \subset X \times \mathbb{R}$$

Search for a function  $f(x_i) \approx y_i$  for "many" i. $\rightarrow$ construct specific **hyperplane H**

$$H : \langle w, x \rangle + b$$

H does best possible
class separation

$$w \in X, b \in \mathbb{R}$$

H never deviates larger than ε,
so called ε-regression

# SVM basics: the formal problem simplified

SVC: classification

SVR: regression

$$\{(x_1, y_1), ..., (x_n, y_n)\} \subset X \times \{-1, +1\}$$

Input / training data

$$\{(x_1, y_1), ..., (x_n, y_n)\} \subset X \times \mathbb{R}$$

Search for a function $f(x_i) \approx y_i$ for "many" i. → construct specific **hyperplane H**

$$H : \langle w, x \rangle + b$$

H does best possible class separation

$$w \in X, b \in \mathbb{R}$$

H never deviates larger than ε, so called ε-regression
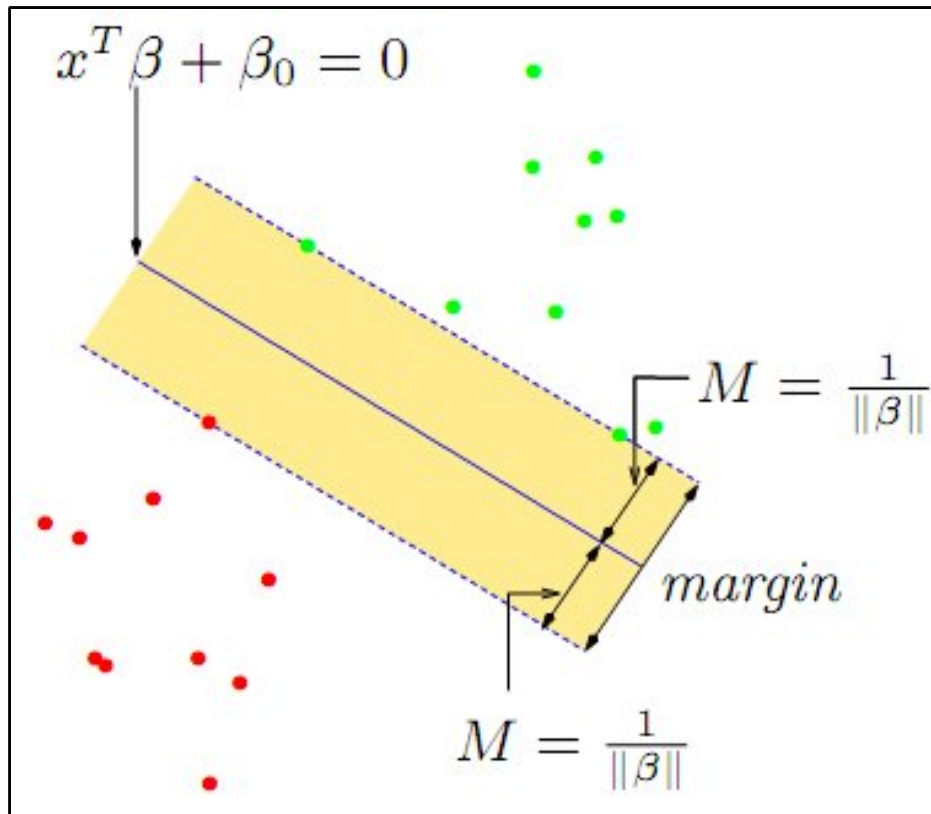
$$f(x) = sgn(\langle w, x \rangle + b)$$

Yielding predictions for unknown $x \in X$

$$f(x) = \langle w, x \rangle + b$$

11

# SVM basics: optimal hyperplanes

How to choose the **best** hyperplane



$$x^T \beta + \beta_0 = 0$$

$$M = \frac{1}{\|\beta\|}$$

*margin*

$$M = \frac{1}{\|\beta\|}$$

from "The elements of statistical learning"

Choose the hyperplane that **maximizes the margin**

$\rightarrow$ minimize $\|\beta\|$

# SVM basics: introducing slack variables

Getting around with **non-feasible** problems



from "The elements of statistical learning"

**Allow violations** of the margin constraint but minimize the extent of the violations

# SVM basics: the kernel-trick

Getting **non-linear**

Every dot product is replaced by a **non-linear kernel function**
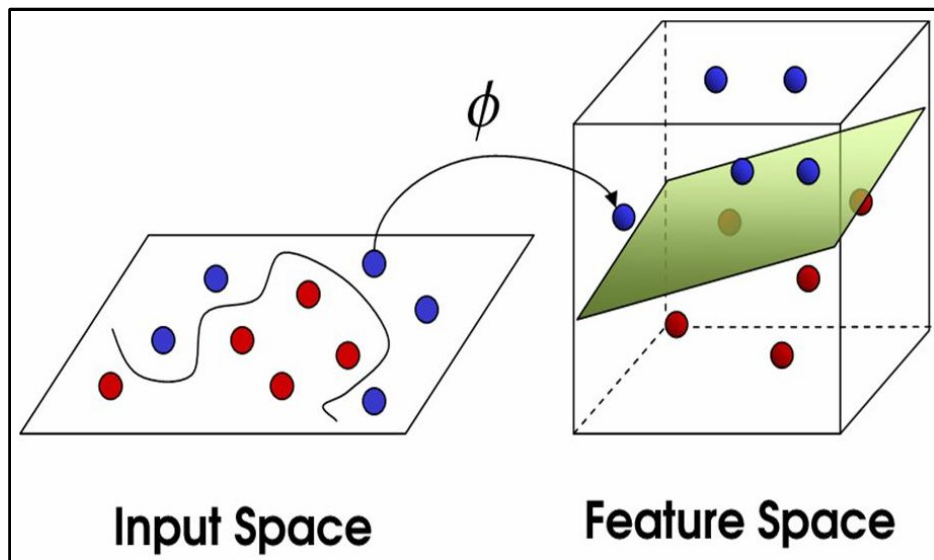$\rightarrow$ input space transported into a high-dimensional feature space



from www.imtech.res.in: rice blast prediction
http://www.imtech.res.in/raghava/rbpred/svm.jpg

# SVM basics: the kernel-trick

Getting **non-linear**

Every dot product is replaced by a **non-linear kernel function**
→ input space transported into a high-dimensional feature space



**Input Space**  **Feature Space**

Some kernel functions:

polynomial (homogenous):
$$k(x, x') = \langle x, x' \rangle^d$$

polynomial (inhomogeneous):
$$k(x, x') = (\langle x, x' \rangle + 1)^d$$

radial basis function:
$$k(x, x') = e^{-\gamma * \langle x - x', x - x' \rangle} \ for \ \gamma > 0$$

sigmoid:
$$k(x, x') = \tanh(\kappa \langle x, x' \rangle + c)$$
$$for \ some \ \kappa > 0 \wedge c < 0$$

15

# SVR in particular: from abstract to application

simplified:

$$\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathcal{X} \times \mathbb{R}$$
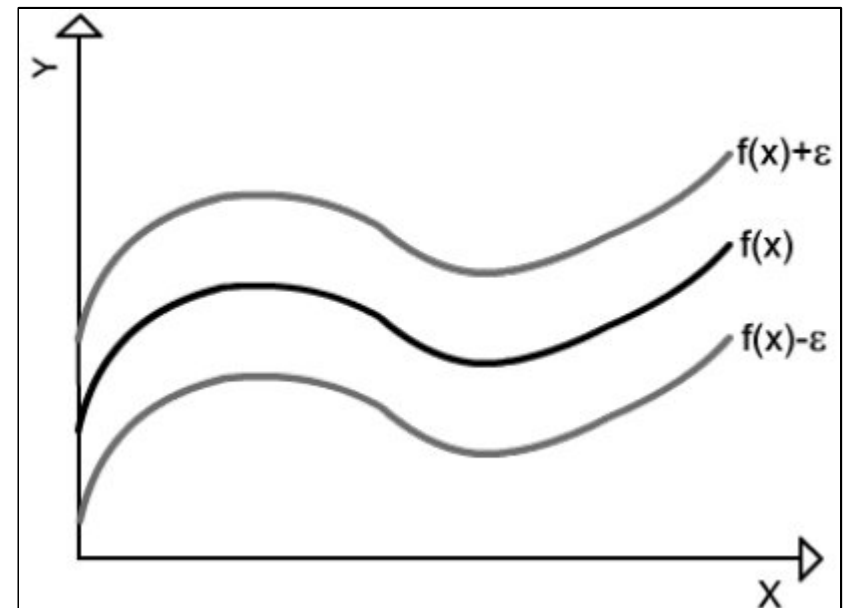
$$f(x_i) \approx y_i$$

$$f(x) = \langle w, x \rangle + b$$

Seeking for best w,b :

minimize $\dfrac{1}{2}\|w\|^2$

subject to $\begin{cases} y_i - \langle w, x_i \rangle - b \leq \epsilon \\ \langle w, x_i \rangle + b - y_i \leq \epsilon \end{cases}$
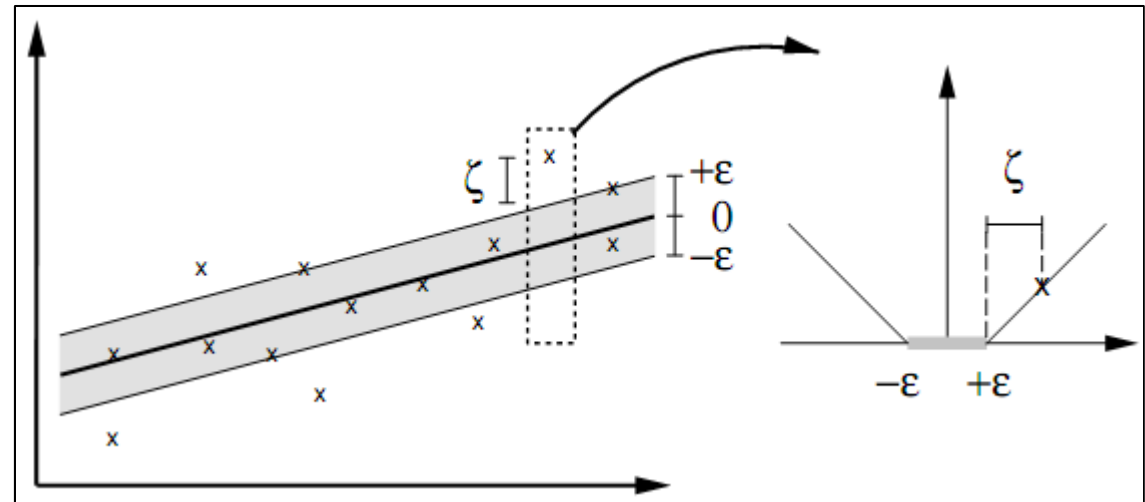


from "Online SVR", Parrella
http://onlinesvr.altervista.org/Theory/Images/03-17.png

16

Adding the „soft-margin" loss
function to obtain the ε-regression:

$$\text{minimize } \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{n}(\xi_i + \xi_i^*)$$

$$\text{subject to } \begin{cases} y_i - \langle w, x_i \rangle - b \leq \epsilon + \xi_i \\ \langle w, x_i \rangle + b - y_i \leq \epsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases}$$



from "A Tutorial on SVR", Smola / Schölkopf

Quadratic Programming
Problem

Behaviour:  ε-insensitive loss-function:

$$|\xi|_\epsilon := \begin{cases} 0 & \text{for } |\xi| \leq \epsilon \\ |\xi| - \epsilon & \text{otherwise} \end{cases}$$
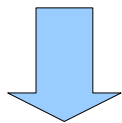
17

# SVR in particular: building the Lagrangian

$$L := \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} (\xi_i + \xi_i^*)$$

objective / primal function

$$-\sum_{i=1}^{n} (\eta_i \xi_i + \eta_i^* \xi_i^*)$$

$$-\sum_{i=1}^{n} \alpha_i (\epsilon + \xi_i - y_i + \langle w, x_i \rangle + b)$$

$\alpha_i, \alpha_i^*, \eta_i, \eta_i^* \geq 0$  are Langrangian multipliers

$$-\sum_{i=1}^{n} \alpha_i^* (\epsilon + \xi_i^* + y_i - \langle w, x_i \rangle - b)$$

$w, b, \xi_i, \xi_i^*$  are the primal variables

$$\partial_b L = \sum_{i=1}^{n} (\alpha_i^* - \alpha_i) = 0$$

$$\partial_w L = w - \sum_{i=1}^{n} (\alpha_i - \alpha_i^*) x_i = 0$$

$$\partial_{\xi_i^{(*)}} L = C - \alpha_i^{(*)} - \eta_i^{(*)} = 0$$

$$\partial_w L = w - \sum_{i=1}^{n} (\alpha_i - \alpha_i^*) x_i = 0$$

can be rewritten as $w = \sum_{i=1}^{n} (\alpha_i - \alpha_i^*) x_i$

$$f(x) = \sum_{i=1}^{n} (\alpha_i - \alpha_i^*) \langle x_i, x \rangle + b$$

w can be **completely** described by a linear combination of the training patterns!

# SVR in particular: optimization problem

Substitution yields the dual optimization problem:

$$
\text{maximize} \begin{cases} -\dfrac{1}{2} \displaystyle\sum_{i,j=1}^{n} (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*)\langle x_i, x_j \rangle \\[2ex] -\epsilon \displaystyle\sum_{i=1}^{n} (\alpha_i + \alpha_i^*) + \sum_{i=1}^{n} y_i(\alpha_i - \alpha_i^*) \end{cases}
$$

$$
\text{subject to} \quad \sum_{i=1}^{n} (\alpha_i - \alpha_i^*) = 0 \ \wedge \ \alpha_i, \alpha_i^* \in [0, C]
$$

easier Quadratic Programming Problem
(solvable by several optimization algorithms)

# Model assessment

A performance measure for regression:

## Pearson's Correlation Coefficient

$$Corr(X,Y) = \frac{Cov(X,Y)}{\sqrt{Var(X)} \cdot \sqrt{Var(Y)}} = \frac{Cov(X,Y)}{\sigma_X \cdot \sigma_Y} \in [-1,1]$$

# Model assessment

A performance measure for regression:

**Mean Squared Error**

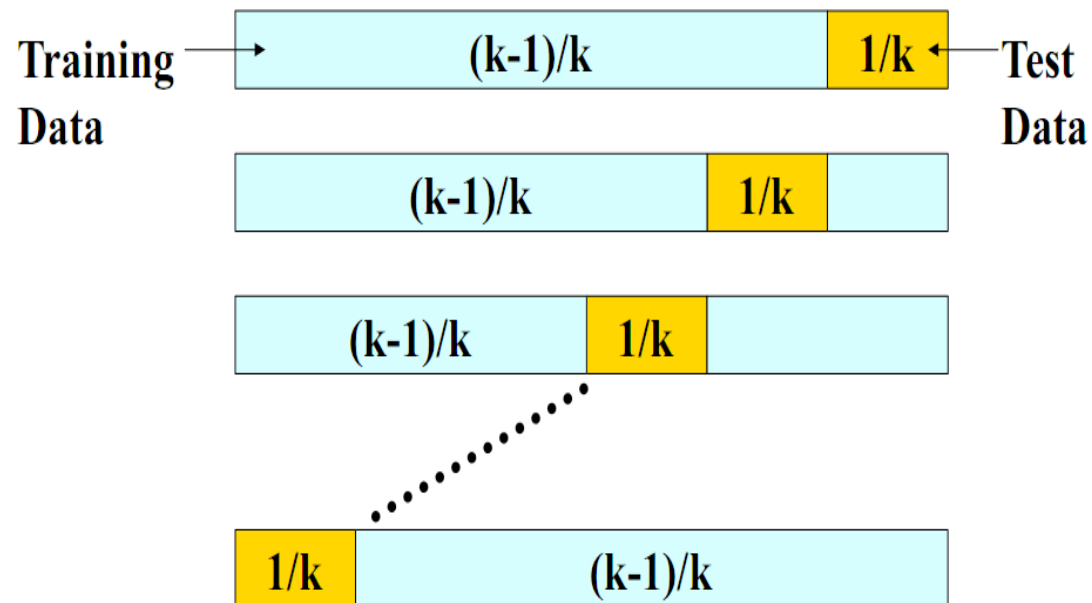$$MSE(f(X)) = E[(X - f(X))^2] = Var(f(X)) + Bias^2(f(X), X)$$

$$f(X): predictor\ for\ X$$

# Model assessment

Estimating the prediction quality with limited data:

## k-fold cross-validation



From Lecture Bioinformatik I by H.-P. Lenhof

# The dataset and the problem



- 2595 residues with computed rSASA
  - of the transmembrane regions of 28 different proteins

- 41 features for each residue
  - frequencies per aa
  - PSIBLAST-score per aa
  - conservation-score

```
> summary(model)

Call:
svm.default(x = trainset[6:ncol(dataset)], y = trainset[5])


Parameters:
   SVM-Type:  eps-regression
 SVM-Kernel:  radial
       cost:  1
      gamma:  0.02439024
    epsilon:  0.1


Number of Support Vectors:  1914


>
```

Finetuning of SVM

# Results

after 2595-fold cross-validation on training data:

Total Mean Squared Error: 0.02627267
Squared Correlation Coefficient: 0.4464628

Correlation of **0.668**

# Accuracy per amino acid

| Amino acid | Mean squared error [$10^{-2}$] |
|---|---|
| R | 0.055 |
| H | 0.058 |
| Q | 0.100 |
| D | 0.243 |
| N | 0.530 |
| E | 0.585 |
| K | 0.879 |
| S | 1.035 |
| T | 1.307 |
| G | 1.789 |

| | |
|---|---|
| Y | 1.902 |
| M | 2.160 |
| P | 2.571 |
| W | 2.603 |
| C | 2.872 |
| A | 3.003 |
| V | 3.306 |
| I | 3.327 |
| L | 3.328 |
| F | 3.704 |

# The practical part
# Accuracy per amino acid

| Amino acid | Mean squared error [$10^{-2}$] | Burried [%] |
|:---:|:---:|:---:|
| R | 0.055 | 100.0 |
| H | 0.058 | 90.9 |
| Q | 0.100 | 100.0 |
| V | 3.306 | 39.9 |
| I | 3.327 | 30.0 |
| L | 3.328 | 27.7 |
| F | 3.704 | 28.4 |

hydrophilic

hydrophobic

# Accuracy per amino acid

| Amino acid | Mean squared error [$10^{-2}$] | Burried [%] | Mean conservation score | |
|---|---|---|---|---|
| R | 0.055 | 100.0 | 1.673 | |
| H | 0.058 | 90.9 | 1.413 | hydrophilic |
| Q | 0.100 | 100.0 | 1.040 | |
| V | 3.306 | 39.9 | -0.230 | |
| I | 3.327 | 30.0 | -0.255 | hydrophobic |
| L | 3.328 | 27.7 | -0.204 | |
| F | 3.704 | 28.4 | 0.008 | |

# Contribution of individual features

| feature(s) | Correlation with rSASA |
|---|---|
| aa-freq. | 0.0806 |
| PSSM-score | 0.1481 |
| conservation | 0.4393 |

Other probably interesting features:

| | |
|---|---|
| hydrophobicity | 0.2862 |
| vdW - volumes | 0.2044 |

# The practical part
# Adding new features to the model

| Features used for prediction | Correlation of prediction and rSASA |
|---|---|
| standard | 0.668 |
| - conservation | 0.644 |
| + hydrophobicity | 0.671 |
| + hydrophobicity + VdW - volumes | 0.669 |

# Summary

**Membrane proteins** are interesting for bioinformatics

Support Vector Machines are **useful tools** in machine learning

Features that are **most likely useful** in
sequence based methods:
> $\rightarrow$ neighborhood (windowing)
> $\rightarrow$ profiles / PSSM
> $\rightarrow$ conservation

But, like always:
> it is **not** about the more the better!

# Used literature

Hastie, Tibshirani, Friedman:
  The Elements of Statistical Learning (sec. edition, 2008)


Park, Hayat, Helms:
  Prediction of the burial status of transmembrane residues of helical
  membrane proteins (paper, BMC Bioinformatics, 2007)


Smola, Schölkopf:
  A Tutorial on Support Vector Regression (paper, 2003)

# Thank you for your attention!