CENTERFO RBIOLOGI CALSEQU ENCEANA LYSIS CBS

Cross validation, training and evaluation of data driven prediction methods

Morten Nielsen Department of Health Technology, DTU

Data driven method training

• A prediction method contains a very large set of parameters

- A matrix for predicting binding for 9meric peptides has 9x20=180 weights
- Over fitting is a problem



years

Evaluation of predictive performance





Evaluation of predictive performance

• Train PSSM on <u>Permuted</u> (random) data

- No pseudo counts, No sequence weighting
- Fit 9*20 parameters to 9*10
 data points
- Evaluate on training data
 - -PCC = 0.97
 - -AUC = 1.0
- Close to a perfect prediction method AND
- Same performance as on the original data

	٢	- AAAMAAKLA
		AAKNLAAAA
		AKALAAAAR
n		AAAAKLATA
5		ALAKAVAAA
ב	4	IPELMRTNG
		FIMGVFTGI
		NVTKVVAWI
		LEPLNLVLK
		VAVIVSVPE
	ř	MRSGRVHAV
n		VRFNIDETE
U		ANYIGQDGI
2		AELCGDPGD
5		QTRAVADGK
_ \\		GRPVPAAHE
Ĕ		MTAQWWLDA
57		FARGVVHVI
-		LQRELTRLC
		AVAEEMTKS



Repeat on large training data (229 ligands)





FLAFFSNGV

FLAFFSNGV WLGNHGFEV TLNAWVKVV LLATSIFKL LLSKNTFYL KVGNCDETV YLNAFIPPV QLWTALVSL MLMTGTLAV QLLADFPEA

FLAFFSNGV	VLMEAQQGI	ILLLDQVLV	KMYEYVFKG
HLMRDPALL	WLVHKQWFL	ALAPSTMKI	MLLTFLTSL
FLIVSLCPT	ITWQVPFSV	RMPAVTDLV	ALYSYASAK
YFLRRLALV	FLLDYEGTL	FLITGVFDI	LLVLCVTQV
MTSELAALI	MLLHVGIPL	GLIIISIFL	IVYGRSNAI
GLYEAIEEC	SLSHYFTLV	GLYYLTTEV	AQSDFMSWV
KLFFAKCLV	VLWEGGHDL	YLLNYAGRI	RLEELLPAV
VLQAGFFLL	AIDDFCLFA	KVVSLVILA	LLVFACSAV
TLKDAMLQL	GLFQEAYPL	YQLGDYFFV	GMVIACLLV
MSDIFHALV	MVVKVNAAL	FMTALVLSL	WLSTYAVRI
GMRDVSFEL	FLGFLATAG	ILAKFLHWL	IVLGNPVFL
QLPLESDAV	SLYPPCLFK	MTPSPFYTV	LLVAPMPTA
KVGNCDETV	RIFPATHYV	IIDQVPFSV	YLNKIQNSL
ILYQVPFSV	YLMKDKLNI	AIMEKNIML	LLNNSLGSV
GLKISLCGI	ALGLGIVSL	MMCPFLFLM	FMFNELLAL
WLETELVFV	ALYWALMES	GLDPTGVAV	GMLPVCPLI
WQDGGWQSV	LLIEGIFFI	SILNTLRFL	GLSLSLCTL
VMLIGIEIL	RLNKVISEL	KVEKYLPEV	YLVAYQATV
SVMDPLIYA	IMSSFEFQV	FTLVATVSI	ILLVAVSFV
GMFGGCFAA	RLLDDTPEV	SLDSLVHLL	LVLQAGFFL
VLAGYGAGI	VILWFSFGA	VLNTLMFMV	FLQGAKWYL





--MMGMFNMLSTVLGVS----AKSSPAYPSVLGQTI-------RHLIFCHSKKKCDELAAK-

Always

How to training a method. A simple statistical method: Linear regression

Observations (training data): a set of x values (input) and y values (output).

Model: **y** = **ax** + **b** (2 parameters, which are estimated from the training data)

Prediction: Use the model to calculate a y value for a new x value

Note: the model does not fit the observations exactly. Can we do better than this?

Overfitting y = ax + b $ax^{6}+bx^{5}+cx^{4}+dx^{3}+ex^{2}+fx+q$ 2 parameter model Good description, poor fit 7 parameter model Poor description, good fit

Note: It is not interesting that a model can fit its observations (training data) exactly.

To function as a prediction method, a model must be able to generalize, i.e. produce sensible output on new data.

How to estimate parameters for prediction?

CENTERFO RBIOLOGI CALSEQU ENCEANA LYSIS CBS

A Regression Problem

Linear Regression Quadratic Regression Join-the-dots

1. Randomly choose 30% of the data to be in a test set

2. The remainder is a training set

Randomly choose
 30% of the data to be in a test set

2. The remainder is a training set

3. Perform your regression on the training set

(Linear regression example)

(Linear regression example) Mean Squared Error = 2.4 1. Randomly choose 30% of the data to be in a test set

2. The remainder is a training set

3. Perform your regression on the training set

4. Estimate your future performance with the test set

1. Randomly choose 30% of the data to be in a test set

2. The remainder is a training set

Perform your regression on the training set

(Quadratic regression example) 4. Estimate your future Mean Squared Error = 0.9

performance with the test set

The test set method

(Join the dots example) Mean Squared Error = 2.2 1. Randomly choose 30% of the data to be in a test set

2. The remainder is a training set

3. Perform your regression on the training set

4. Estimate your future performance with the test set

So quadratic function is best

CENTERFO RBIOLOGI CALSEQU ENCEANA LYSIS CBS

Cross validation

Train on 4/5 of data Test/evaluate on 1/5 =>

Produce 5 different methods each with a different prediction focus

Model over-fitting

NTEREO

CBS

What is going on?

5 fold training

Which method to choose?

5 fold training

The Wisdom of the Crowds

• The Wisdom of Crowds. Why the Many are Smarter than the Few. James Surowiecki

One day in the fall of 1906, the British scientist Fracis Galton left his home and headed for a country fair... He believed that only a very few people had the characteristics necessary to keep societies healthy. He had devoted much of his career to measuring those characteristics, in fact, in order to prove that the vast majority of people did not have them. ... Galton came across a weight-judging

competition...Eight hundred people tried their luck. They were a diverse lot, butchers, farmers, clerks and many other no-experts...The crowd had guessed ... <u>1.197</u> pounds, the ox weighted 1.198

- The highest scoring hit will often be wrong
 - Not one single prediction method is consistently best
- Many prediction methods will have the correct fold among the top 10-20 hits
- If many different prediction methods all have a common fold among the top hits, this fold is probably correct

- Use cross validation
- Evaluate on *concatenated data* and <u>not</u> as an average over each cross-validated performance

Method evaluation

Method evaluation

Model evaluation - to concatenate or not to concatenate

- Cross validation is always good!, but how many folds?
 - Few folds -> small training data sets
 - Many folds -> small test data sets
- 560 peptides for training
 - 50 fold (10 peptides per test set, few data to stop training)
 - 2 fold (280 peptides per test set, few data to train)
 - 5 fold (110 peptide per test set, 450 per training set)

- CENTERFO RBIOLOGI CALSEQU ENCEANA LYSIS CBS
- Use test set to stop training, and test set performance to evaluate training
 - Over-fitting?
- If test set is small, Yes
- If test set is large, No
- Confirm using "true" 5 fold cross validation
 - 1/5 for evaluation
 - 4/5 for 4 fold cross-validation

Conventional 5 fold cross validation

"Nested (or true)" 5 fold cross validation

- If you use the test data for model optimization (hyper-parameter optimization, early stopping ..) you should always use nested cross validation
- When data is scarce, the performance obtained used "conventional" versus "nested" cross validation can be very large
- When data is abundant the difference is in general small

Training/evaluation procedure

- Define method
- Select data
- Deal with data redundancy
 - In method (sequence weighting)
 - In data (Hobohm)
- Deal with over-fitting either
 - in method (SMM regulation term) or
 - in training (stop fitting on test set performance)
- Evaluate method using cross-validation