



Beslutsträd och andra prediktiva modeller

Mathias Lanner Sas Institute



Agenda

Introduktion till Prediktiva modeller
Beslutsträd
Pruning
Regressioner
Neurala Nätverk
Utvärdering av modeller

Predictive Modeling Applications

-  Database marketing
-  Financial risk management
-  Fraud detection
-  Process monitoring
-  Pattern detection

Predictive Modeling Training Data

Training Data

case 1:	<i>inputs</i>	<i>target</i>
case 2:	<i>inputs</i>	<i>target</i>
case 3:	<i>inputs</i>	<i>target</i>
case 4:	<i>inputs</i>	<i>target</i>
case 5:	<i>inputs</i>	<i>target</i>

————— Numeric or categorical values

Predictive Modeling Score Data

Training Data

case 1: inputs target
case 2: inputs target
case 3: inputs target
case 4: inputs target
case 5: inputs target

Score Data

case 1: inputs ?
case 2: inputs ?
case 3: inputs ?
case 4: inputs ?
case 5: inputs ?

Only input values known

5

Predictions

Training Data

case 1: inputs target
case 2: inputs target
case 3: inputs target
case 4: inputs target
case 5: inputs target

Predictions

Score Data

case 1: inputs ?
case 2: inputs ?
case 3: inputs ?
case 4: inputs ?
case 5: inputs ?

6

Predictions

Training Data

case 1: inputs target
case 2: inputs target
case 3: inputs target
case 4: inputs target
case 5: inputs target

Predictions

prediction
prediction
prediction
prediction

Score Data

case 1: inputs ?
case 2: inputs ?
case 3: inputs ?
case 4: inputs ?
case 5: inputs ?

prediction
prediction
prediction
prediction

7

Predictive Modeling Essentials



Predict new cases



Select useful inputs



Optimize complexity

8

Predictive Modeling Essentials

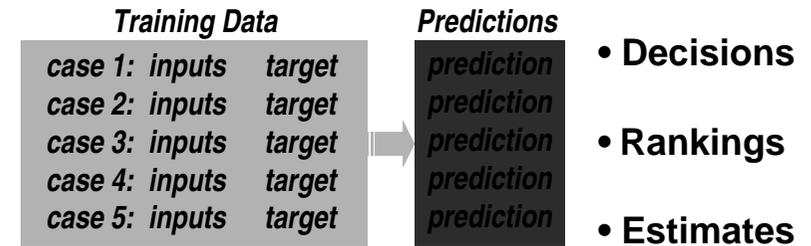
 Predict new cases

 Select useful inputs

 Optimize complexity

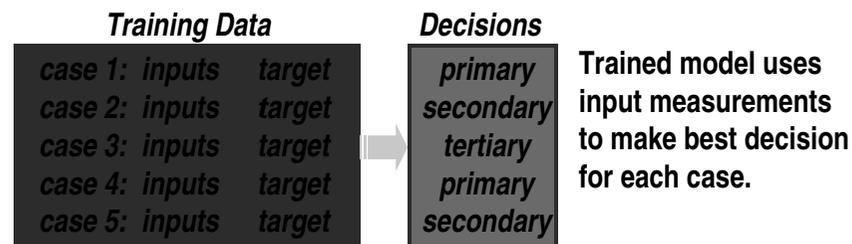
9

Three Prediction Types



10

Decision Predictions



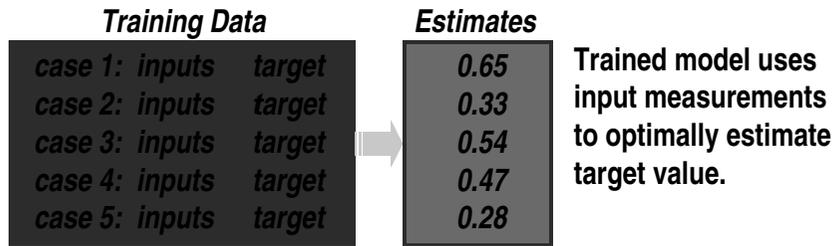
11

Ranking Predictions



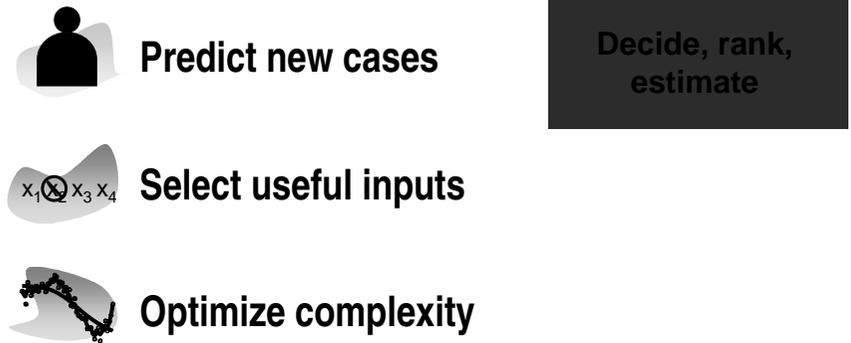
12

Estimate Predictions



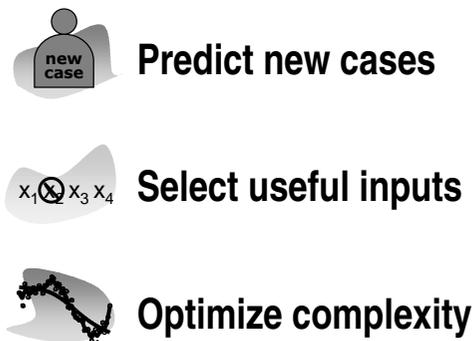
13

Model Essentials – Predict Review



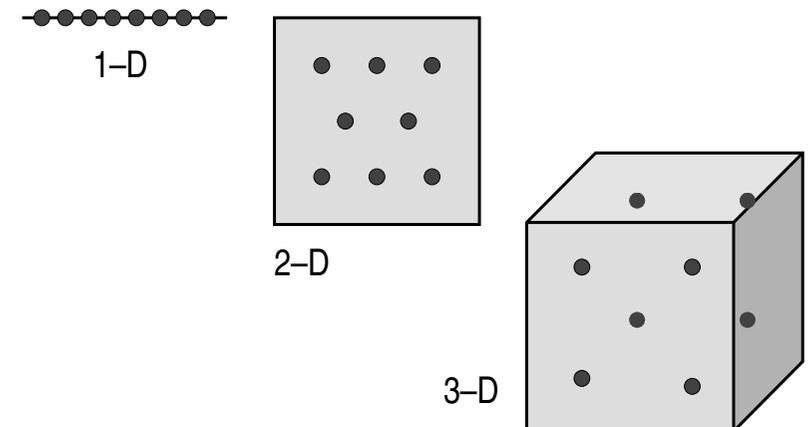
14

Model Essentials – Select Review



15

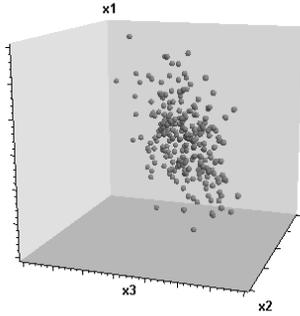
Curse of Dimensionality



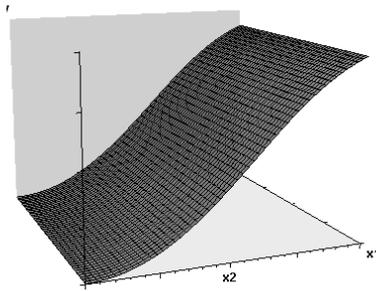
16

Input Selection

Redundancy



Irrelevancy



17

Model Essentials – Select Review



Predict new cases

Decide, rank,
estimate



Select useful inputs

Eradicate
redundancies
irrelevancies



Optimize complexity

18

Model Essentials – Optimize



Predict new cases



Select useful inputs



Optimize complexity

19

Fool's Gold

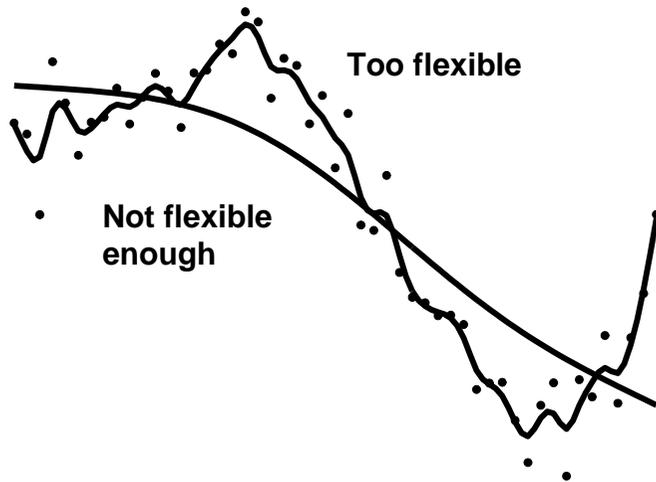
My model fits the
training data perfectly..

I've struck it rich!



20

Model Complexity



21

Data Splitting



22

Training Data Role

Training Data

case 1: inputs	target
case 2: inputs	target
case 3: inputs	target
case 4: inputs	target
case 5: inputs	target



Training data gives sequence of predictive models with increasing complexity.

Validation Data

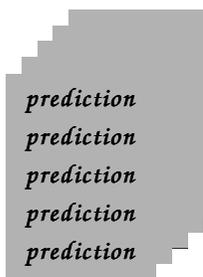
case 1: inputs	target
case 2: inputs	target
case 3: inputs	target
case 4: inputs	target
case 5: inputs	target

23

Validation Data Role

Training Data

case 1: inputs	target
case 2: inputs	target
case 3: inputs	target
case 4: inputs	target
case 5: inputs	target



Validation Data

case 1: inputs	target
case 2: inputs	target
case 3: inputs	target
case 4: inputs	target
case 5: inputs	target

Validation data helps select best model from sequence

24

Validation Data Role

Training Data

case 1: inputs target
 case 2: inputs target
 case 3: inputs target
 case 4: inputs target
 case 5: inputs target

prediction
 prediction
 prediction
 prediction
 prediction

Validation Data

case 1: inputs target
 case 2: inputs target
 case 3: inputs target
 case 4: inputs target
 case 5: inputs target

prediction
 prediction
 prediction
 prediction
 prediction

Validation data helps select best model from Sequence.

Model Essentials – Optimize



Predict new cases

Decide, rank, estimate



Select useful inputs

Eradicate redundancies irrelevancies



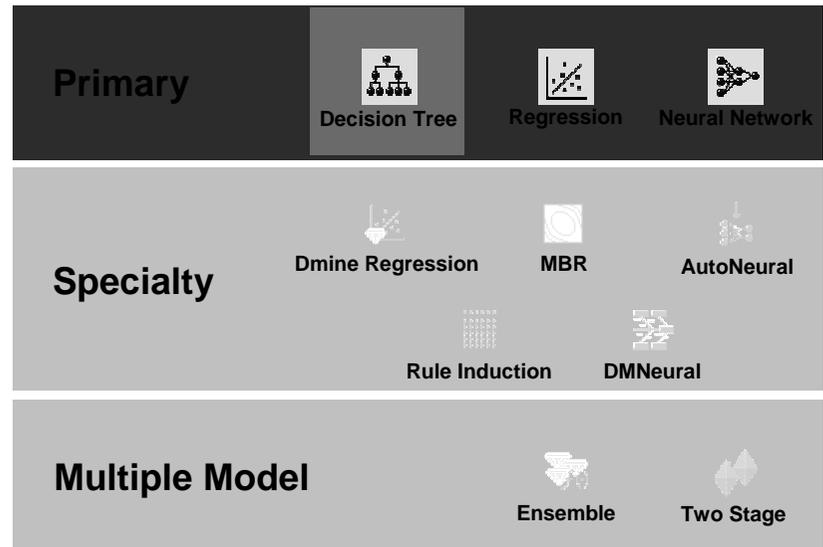
Optimize complexity

Tune models with validation data

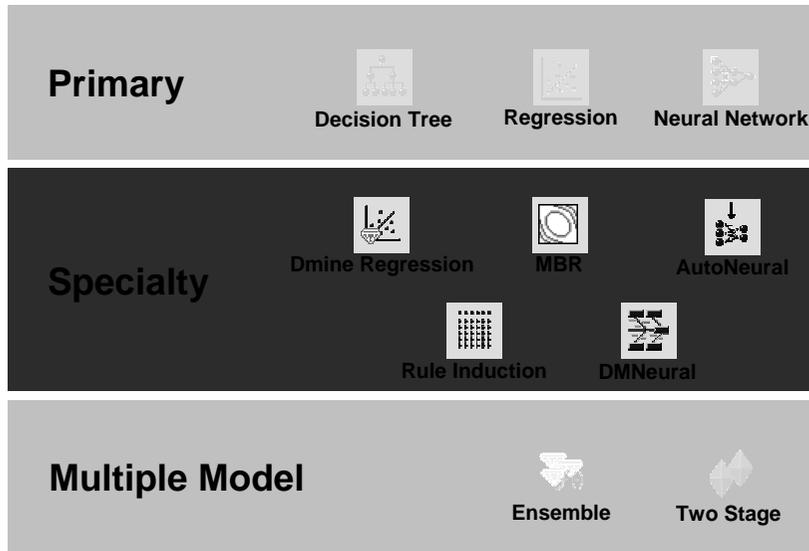
Agenda

Introduktion till Prediktiva modeller
Beslutsträd
Pruning
Regressioner
Neurala Nätverk
Utvärdering av modeller

Predictive Modeling Tools

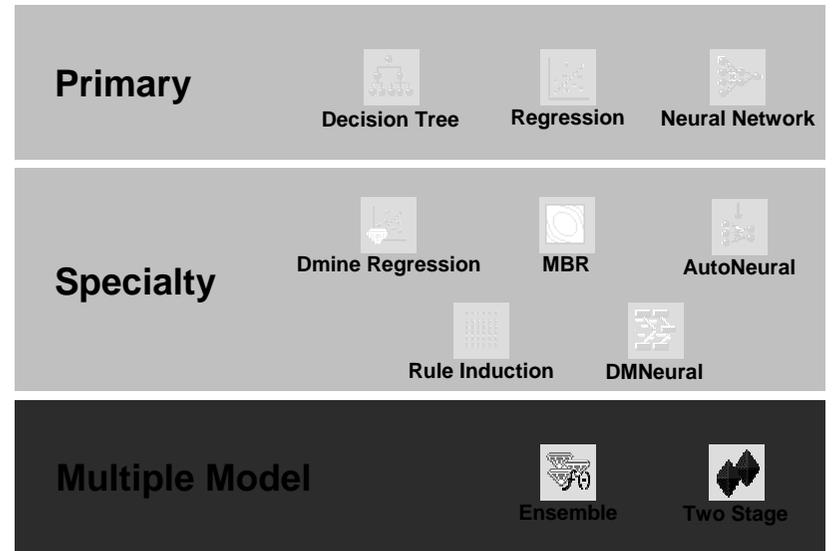


Predictive Modeling Tools



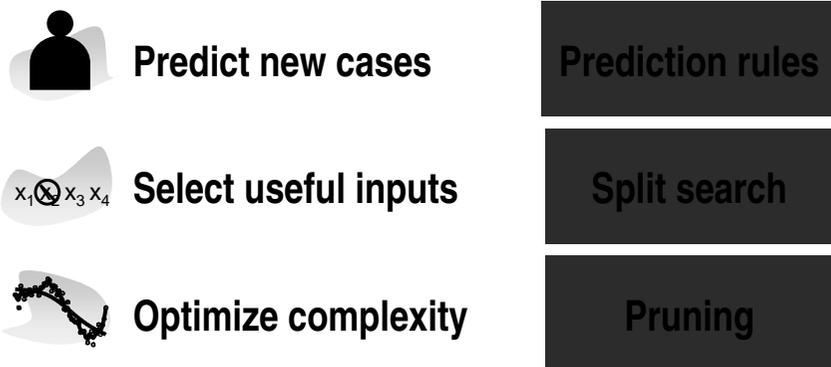
29

Predictive Modeling Tools



30

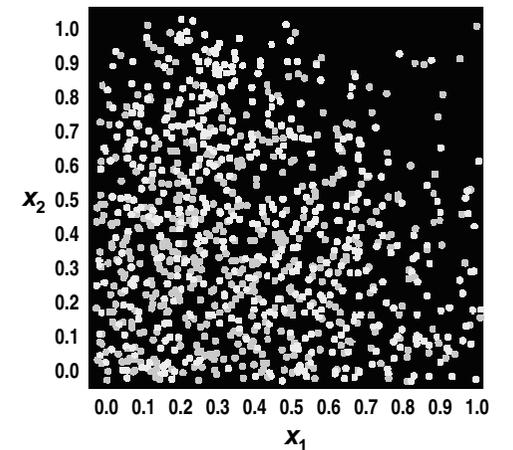
Model Essentials – Decision Trees



31

Simple Prediction Illustration

Analysis goal:
 Predict the color of a dot based on its location in a scatter plot.



32

Model Essentials – Decision Trees



Predict new cases

Prediction rules



Select useful inputs

Split search

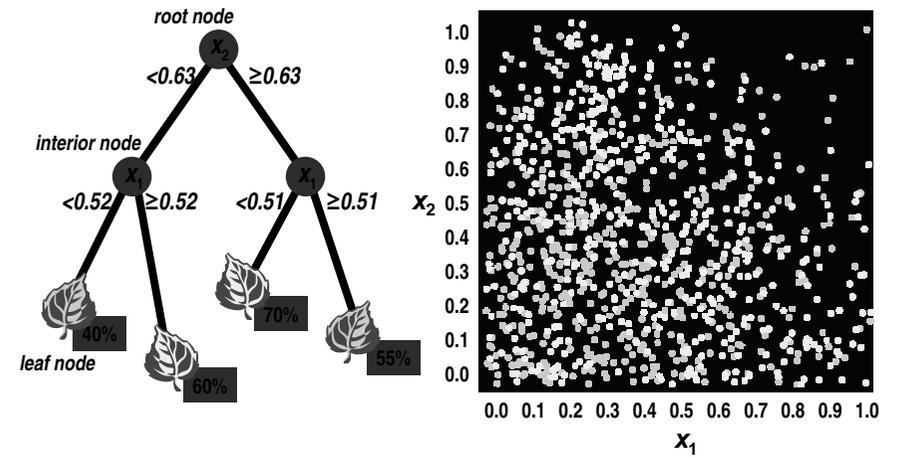


Optimize complexity

Pruning

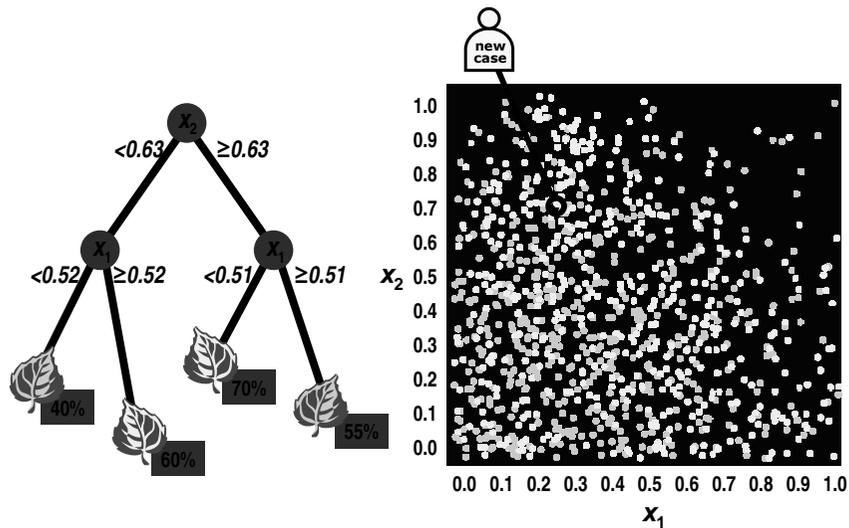
33

Decision Tree Prediction Rules



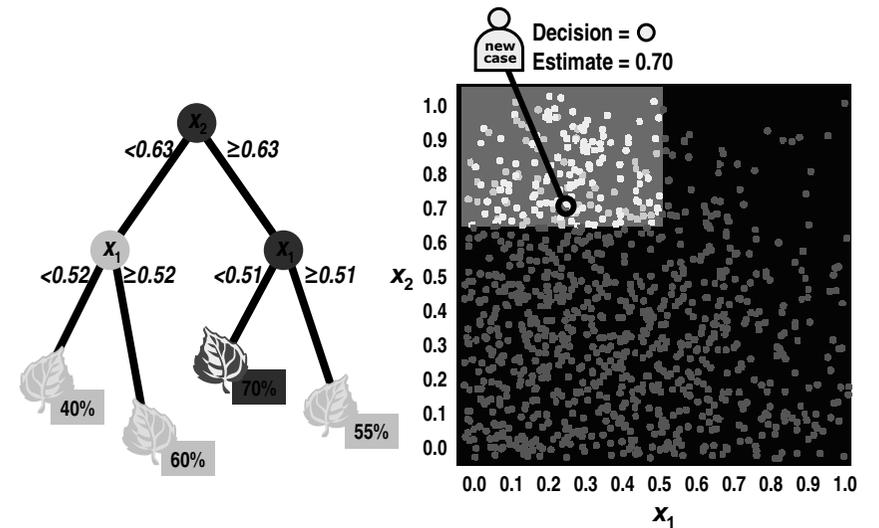
34

Decision Tree Prediction Rules



35

Decision Tree Prediction Rules



36

Model Essentials – Decision Trees



Predict new cases

Prediction rules



Select useful inputs

Split search



Optimize complexity

Pruning

37



Demo Beslutsträd

38

Agenda

Introduktion till Prediktiva modeller
Beslutsträd
Pruning
Regressioner
Neurala Nätverk
Utvärdering av modeller

39

Model Essentials – Decision Trees



Predict new cases

Prediction rules



Select useful inputs

Split search

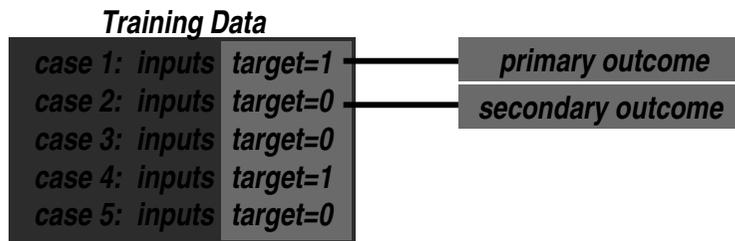


Optimize complexity

Pruning

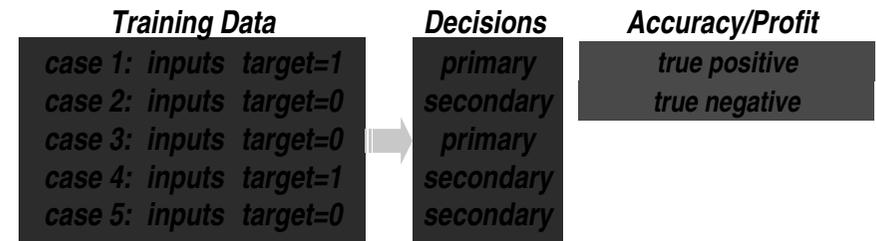
40

Binary Targets



41

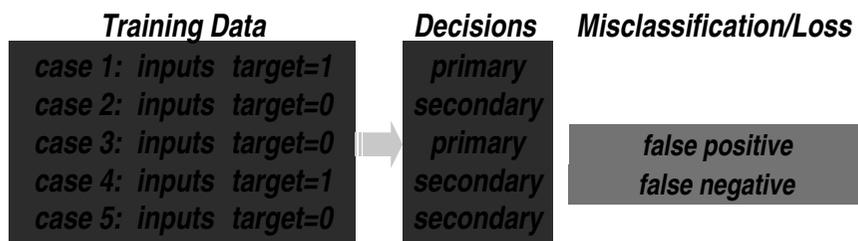
Decision Assessment



Focus on correct decisions

42

Decision Assessment (for Pessimists)



Focus on incorrect decisions

43

Ranking Assessment



Focus on correct ordering

44

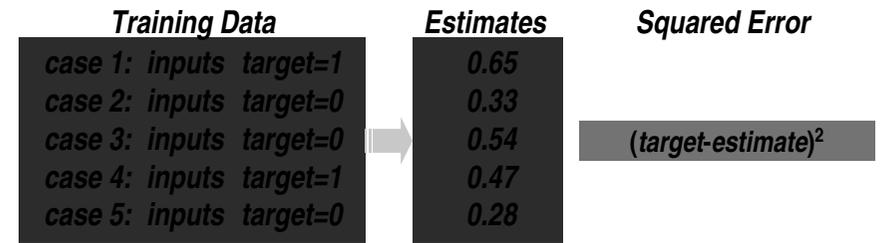
Ranking Assessment (for Pessimists)



Focus on incorrect ordering

45

Estimate Assessment (only Pessimistic!)



Focus on incorrect estimation

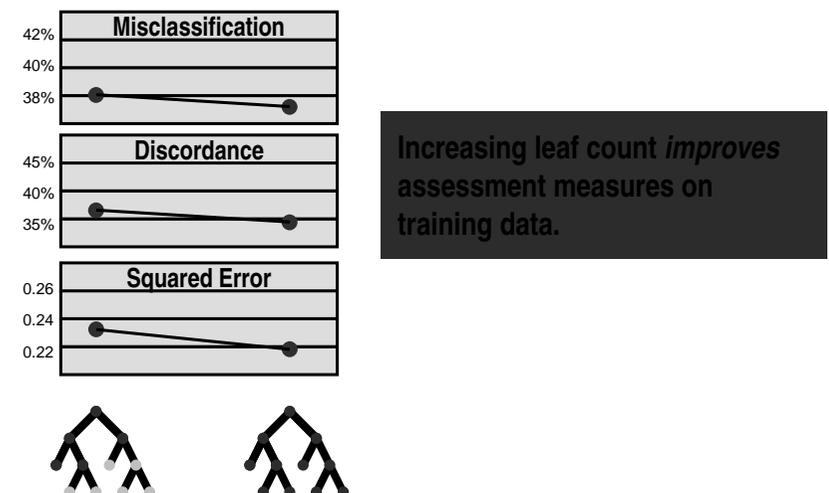
46

Predictive Modeling Assessments



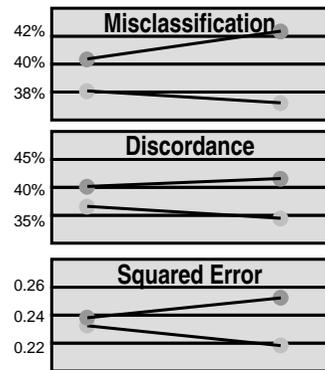
47

Optimistic Assessment, Pessimistic Stats



48

Unbiased Assessment



Increasing leaf count *might worsen* assessment measures on validation data.



Demo på Pruning

Model Essentials – Regressions

-  Predict new cases
-  Select useful inputs
-  Optimize complexity

Prediction formula

Sequential selection

Optimal sequence model

Model Essentials – Regressions

-  Predict new cases
-  Select useful inputs
-  Optimize complexity

Prediction formula

Sequential selection

Optimal sequence model

Linear Regression Prediction Formula

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$

intercept estimate
parameter estimate
prediction estimate

Choose intercept and parameter estimates to *minimize*.

squared error function

$$\sum (y_i - \hat{y}_i)^2$$

training data

53

Logistic Regression Prediction Formula

$$\log\left(\frac{\hat{p}}{1 - \hat{p}}\right) = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$

logit scores

Choose intercept and parameter estimates to *maximize*.

log-likelihood function

$$\sum \log(\hat{p}_i) + \sum \log(1 - \hat{p}_i)$$

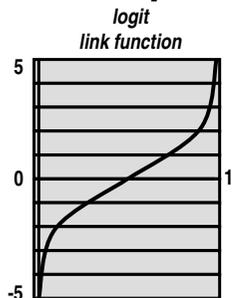
primary outcome training cases
secondary outcome training cases

54

Logit Link Function

$$\log\left(\frac{\hat{p}}{1 - \hat{p}}\right) = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$

logit scores



Δx_i consequence

<small>doubling amount</small>	$1 \Rightarrow \text{odds} \times \exp(w_i)$	<small>odds ratio</small>
	$\frac{0.69}{w_i} \Rightarrow \text{odds} \times 2$	

Model interpretation

55

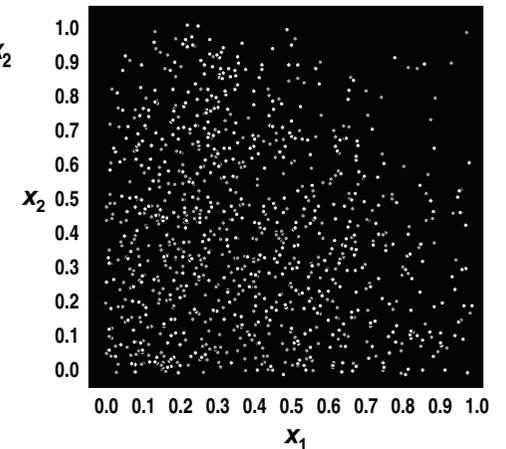
Simple Prediction Illustration – Regressions

logit equation

$$\text{logit}(\hat{p}) = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$

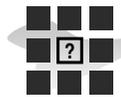
$$\hat{p} = \frac{1}{1 + e^{-\text{logit}(\hat{p})}}$$

logistic equation



56

Beyond the Prediction Formula



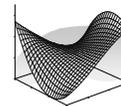
Missing values



Extreme or unusual values



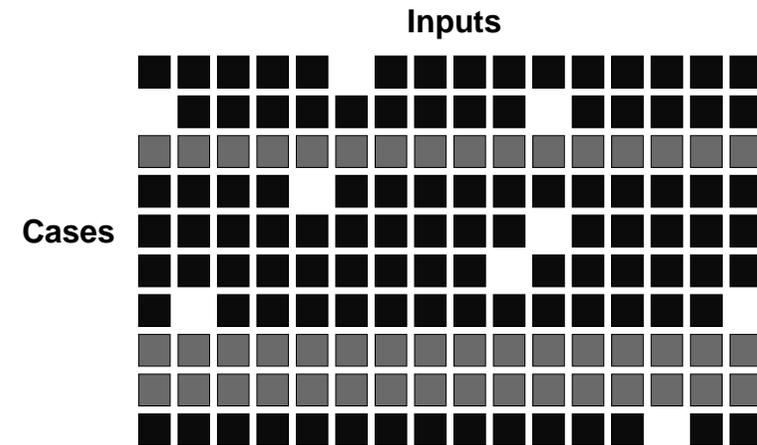
Non-numeric inputs



Nonlinearity and Non-additivity

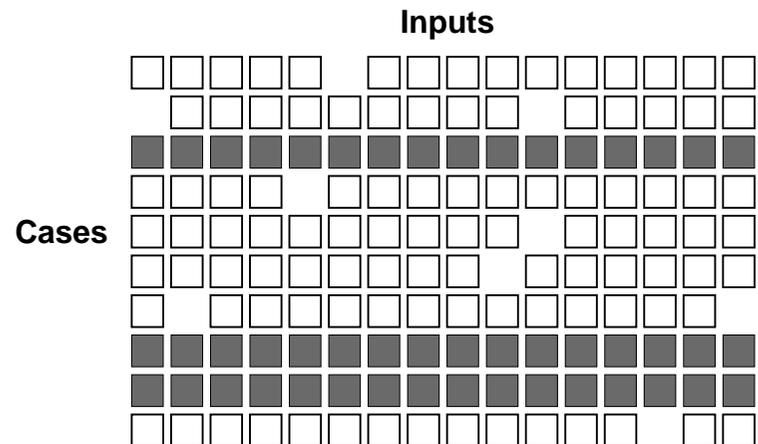
57

Missing Values and Regression Modeling



58

Missing Values and Regression Modeling



59

Missing Values and the Prediction Formula

Prediction Formula:

$$\text{logit}(p) = -2.1 + 0.072x_1 - 0.89x_2 - 1.24x_3$$

New Case:

$$(x_1, x_2, x_3) = (2, ?, -1)$$

Predicted Value:

$$\text{logit}(p) = -2.1 + 0.144x_1 - 0.8? + 1.24$$

60

Missing Value Causes

-  **N/A** Not applicable
-  **?** No match
-  **⊘** Non-disclosure

61

Missing Value Remedies

-  **N/A** Not applicable
-  **?** No match
-  **⊘** Non-disclosure

Synthetic distribution



Estimation
 $x_i = f(x_1, \dots, x_p)$

62

Model Essentials – Regressions

-  **Predict new cases**
-  **Select useful inputs**
-  **Optimize complexity**

Prediction
formula

Sequential
selection

Optimal sequence
model

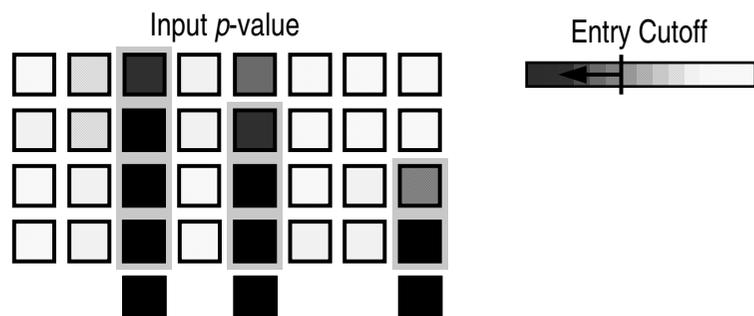
63

Sequential Selection – Forward



64

Sequential Selection – Forward



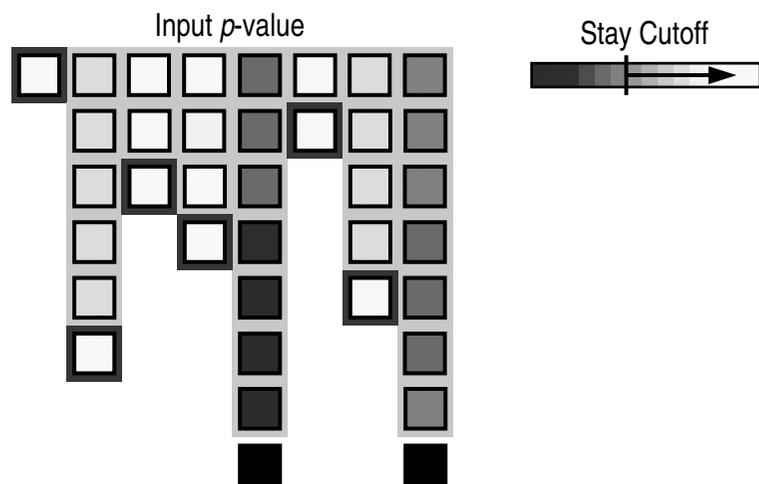
65

Sequential Selection – Backward



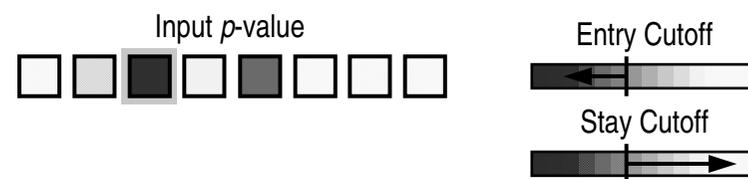
66

Sequential Selection – Backward



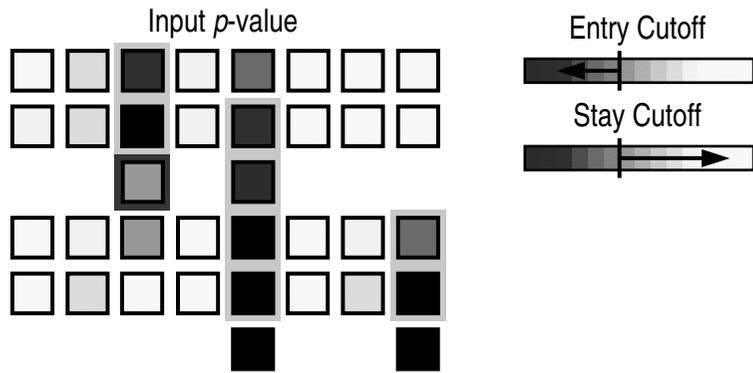
67

Sequential Selection – Stepwise



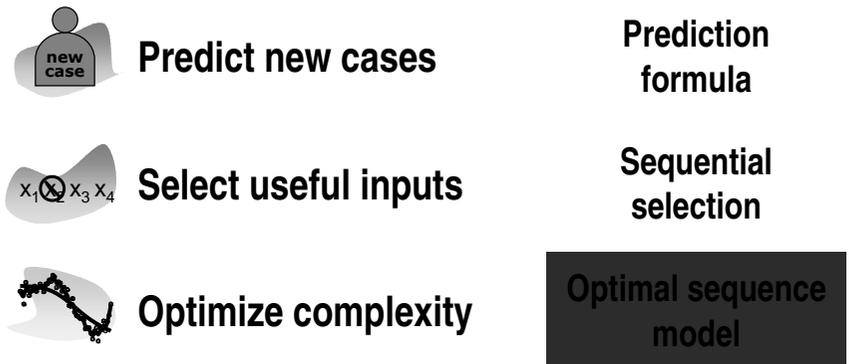
68

Sequential Selection – Stepwise



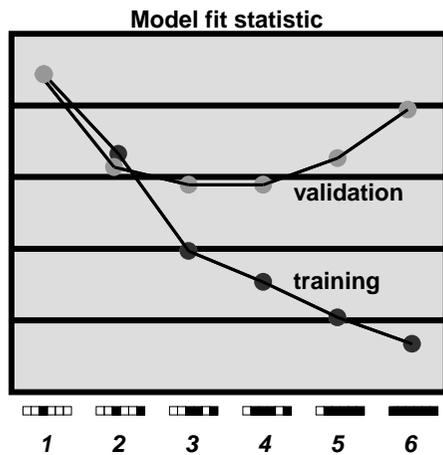
69

Model Essentials – Regressions



70

Model Fit versus Complexity

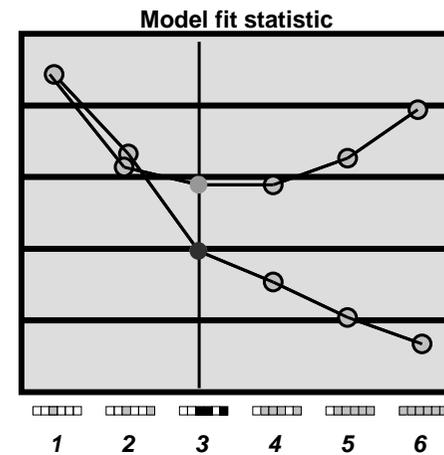


Evaluate each sequence step.

Choose simplest optimal model.

71

Select Model with Optimal Validation Fit



Evaluate each sequence step.

Choose simplest optimal model.

72

Agenda

Introduktion till Prediktiva modeller
Beslutsträd
Pruning
Regressioner
Neurala Nätverk
Utvärdering av modeller

Model Essentials – Neural Networks

 Predict new cases	Prediction formula
 Select useful inputs	None
 Optimize complexity	Stopped training

Model Essentials – Neural Networks

 Predict new cases	Prediction formula
 Select useful inputs	None
 Optimize complexity	Stopped training

Neural Network Prediction Formula

prediction estimate $\hat{y} = \hat{w}_{00} + \hat{w}_{01} H_1 + \hat{w}_{02} H_2 + \hat{w}_{03} H_3$

bias estimate
weight estimate

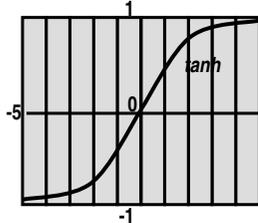
hidden unit

$H_1 = \tanh(\hat{w}_{10} + \hat{w}_{11} x_1 + \hat{w}_{12} x_2)$

$H_2 = \tanh(\hat{w}_{20} + \hat{w}_{21} x_1 + \hat{w}_{22} x_2)$

$H_3 = \tanh(\hat{w}_{30} + \hat{w}_{31} x_1 + \hat{w}_{32} x_2)$

activation function



Neural Network Binary Prediction Formula

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{w}_{00} + \hat{w}_{01} H_1 + \hat{w}_{02} H_2 + \hat{w}_{03} H_3$$

*logit
link function*

$$H_1 = \tanh(\hat{w}_{10} + \hat{w}_{11} x_1 + \hat{w}_{12} x_2)$$

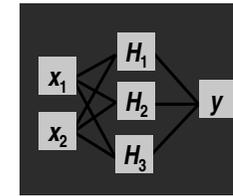
$$H_2 = \tanh(\hat{w}_{20} + \hat{w}_{21} x_1 + \hat{w}_{22} x_2)$$

$$H_3 = \tanh(\hat{w}_{30} + \hat{w}_{31} x_1 + \hat{w}_{32} x_2)$$

77

Neural Network Diagram

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{w}_{00} + \hat{w}_{01} H_1 + \hat{w}_{02} H_2 + \hat{w}_{03} H_3$$



input layer hidden layer target layer

$$H_1 = \tanh(\hat{w}_{10} + \hat{w}_{11} x_1 + \hat{w}_{12} x_2)$$

$$H_2 = \tanh(\hat{w}_{20} + \hat{w}_{21} x_1 + \hat{w}_{22} x_2)$$

$$H_3 = \tanh(\hat{w}_{30} + \hat{w}_{31} x_1 + \hat{w}_{32} x_2)$$

78

Prediction Illustration – Neural Networks

logit equation

$$\text{logit}(\hat{p}) = \hat{w}_{00} + \hat{w}_{01} H_1 + \hat{w}_{02} H_2 + \hat{w}_{03} H_3$$

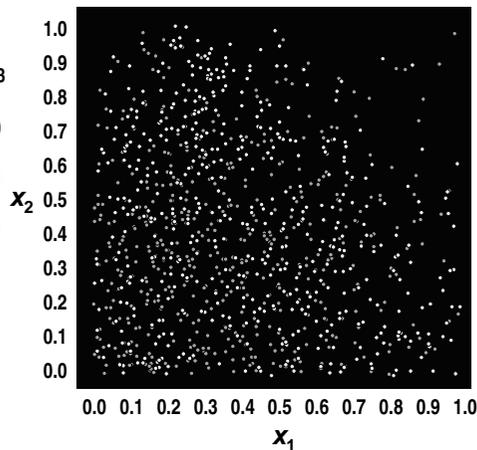
$$H_1 = \tanh(\hat{w}_{10} + \hat{w}_{11} x_1 + \hat{w}_{12} x_2)$$

$$H_2 = \tanh(\hat{w}_{20} + \hat{w}_{21} x_1 + \hat{w}_{22} x_2)$$

$$H_3 = \tanh(\hat{w}_{30} + \hat{w}_{31} x_1 + \hat{w}_{32} x_2)$$

$$\hat{p} = \frac{1}{1 + e^{-\text{logit}(\hat{p})}}$$

logistic equation



79

Agenda

Introduktion till Prediktiva modeller
Beslutsträd
Pruning
Regressionser
Neurala Nätverk
Utvärdering av modeller

80

Assessment Types

The Model Comparison tool provides

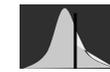
**KS^C
ASE** Summary statistics

 Statistical graphics

81

Summary Statistics Summary

Prediction Type



Decisions

Statistic

Accuracy / Misclassification
Profit / Loss
KS-statistic

1,2,3,...

Rankings

ROC Index (concordance)
Gini coefficient

$\hat{p} \approx E(Y)$

Estimates

Average squared error
SBC / Likelihood

82

Summary Statistics Summary

Prediction Type



Decisions

Statistic

Accuracy / Misclassification
Profit / Loss
KS-statistic

1,2,3,...

Rankings

ROC Index (concordance)
Gini coefficient

$\hat{p} \approx E(Y)$

Estimates

Average squared error
SBC / Likelihood

83

Summary Statistics Summary

Prediction Type



Decisions

Statistic

Accuracy / Misclassification
Profit / Loss
KS-statistic

1,2,3,...

Rankings

ROC Index (concordance)
Gini coefficient

$\hat{p} \approx E(Y)$

Estimates

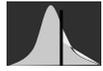
Average squared error
SBC / Likelihood

84

Statistical Graphics Summary

Prediction Type

Statistic



Decisions

1,2,3,...

Rankings

$\hat{p} \approx E(Y)$

Estimates

Sensitivity charts
Response rate charts