

Lecture Outline

Thus, the lecture will contain:

1. Introduction
2. How to frame the business problem
3. How to transfer it to a problem which can be solved with analytics methods
4. Data identification and prioritisation, data collection and data harmonisation
5. Identification of problem solving approaches and appropriate tools (not only R even though this is important)
6. How to set up and validate models
7. **The deployment of a model**
8. Model lifecycle
9. Some words about soft skills needed by statistical and mathematical professionals

Chapter 7

The deployment of a model

7.1 Introduction

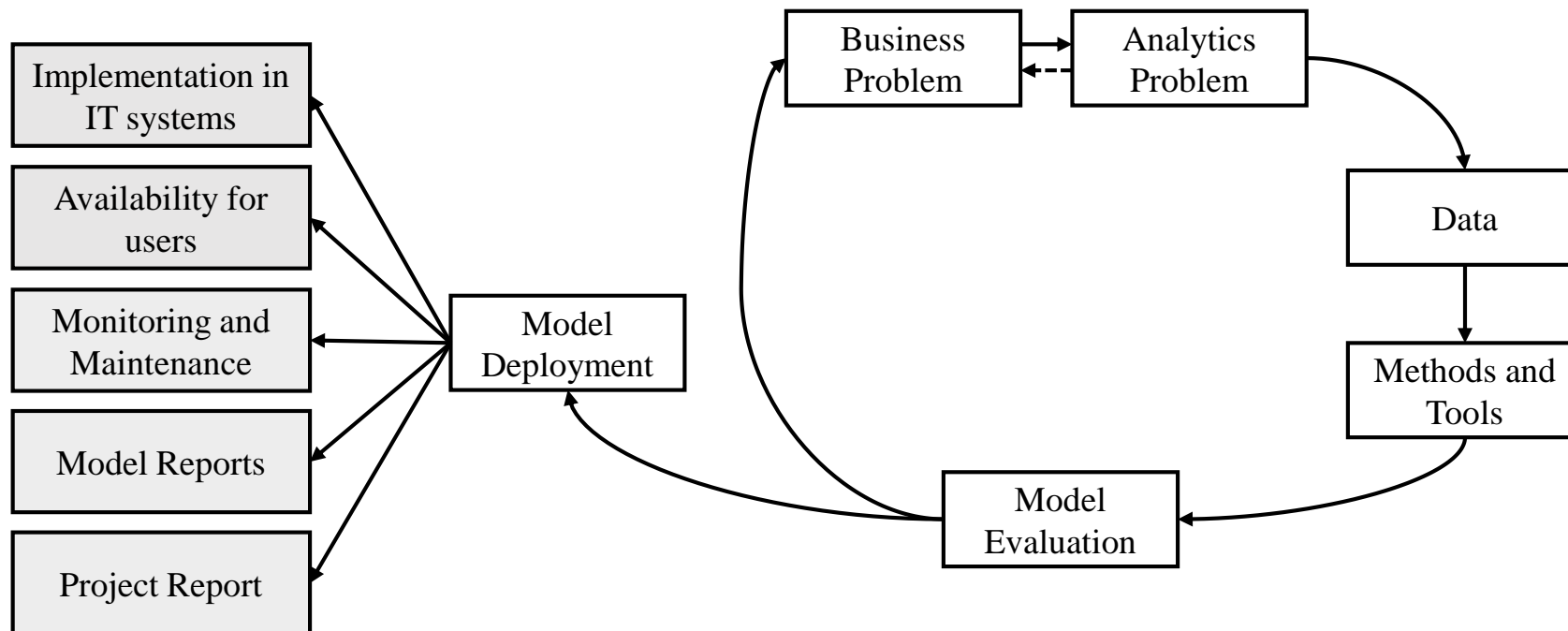
What is deployment?

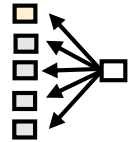
Deployment is the implementation of the data analytics model into an IT or computer system for its use on a regular basis. It contains:

- Implementation and coding into the IT or computer system
- Make the model available to the users and train the users
- Monitor and maintain the model
- Producing reports out of the model
- And last but not least the production of a final project report

7.1 Introduction

What is deployment?





7.2 Deployment Steps

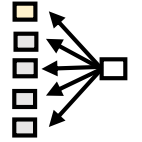


Implementation plan: A successful implementation requires:

- A comprehensive implementation plan
- Considering the most common and crucial trade offs between the technical and business aspects of the implementation
- Being aware of the most common implementation mistakes in order to avoid them

The following checklist is a help for elaborating the implementation plan:

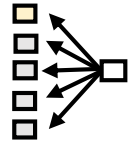
Question	Purpose / Considerations
1. What is the concept of operations (CONOP) of the current system/process?	<ul style="list-style-type: none">• The solution must be accepted by the end users• Therefore, bear in mind their feedback• Such solution should not completely disrupt their normal work process• The solution should interact with the system and user in a way that is maximally helpful and minimally disruptive



7.2 Deployment Steps

Implementation and coding into the IT or computer systems

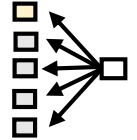
Question	Purpose / Considerations
2. What is the legacy / computing / processing / operational environment?	<ul style="list-style-type: none">• Need to familiarize yourself with the system's configuration management system, e.g. dates that the systems and application software were installed or updated• It often happens that there are leftover software, hardware and methods from earlier versions of the system e.g. data paths that don't seem to make sense, a mixture of old and new data formats
3. What are the available interface mechanisms/processes?	<ul style="list-style-type: none">• This questions aims at designing the system interface, the communication with the system and user• Typically there will be an existing Application Program Interface as part of the system specification that your solution can use to interact with the system and user• Ask for the corresponding documentation• Coordinate with the system developer in case you need to communicate with the system and/or user by direct use of the system services in order to avoid conflict



7.2 Deployment Steps

Implementation and coding into the IT or computer systems

Question	Purpose / Considerations
4. At what points in the processing stream can data be injected?	<ul style="list-style-type: none">• Typically your application will read / write / update data, therefore, you want to know at what times and places in the system architecture this is possible and appropriate• If you are processing during operational periods when other other systems activities are also taking place, then you need to synchronize with actions
5. What are the political / organizational considerations for interaction with the systems?	<ul style="list-style-type: none">• Systems are often operated following organizational policies or implicit agreements about what is allowed and when• In order to avoid violating any policy it is in your best interest to clearly understand these policies and consider them in the implementation plan



7.2 Deployment Steps

Implementation and coding into the IT or computer systems

The *most common* 4 ways of deploying models in data mining are:

1. Data mining tool (or cloud)



2. Programming language (Java, C, VB...)

```
Public Function RegressionModel(ByVal Age As Double, ByVal Sal:
                                ByVal Car_location As String) i

'Declare
Dim X1, X2, X3 As Double

'Numerical variables
If IsNumeric(Age) = False Then X1 = 25 Else X1 = Age
If IsNumeric(Salary) = False Then X2 = 31000 Else X2 = Sal:

'Categorical variables
Select Case Car_location
Case "carpark"
    X3 = 41.1
Case "street"
    X3 = 325.03
Case Else
    X3 = 100
End Select

'Output
Return 132.37 + 7.1 * X1 + 0.01 * X2 + X3

End Function
```

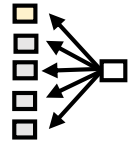
3. Database and SQL script

```
CREATE FUNCTION RegressionModel (@Age real, @Salary real, @Car_locat:
RETURNS real
AS
BEGIN
    declare @Return varchar(30)
    declare @X1 real
    declare @X2 real
    declare @X3 real
    --Age
    select @X1= case
                when @Age is null then 25
                else @Age
            end
    --Salary
    select @X2= case
                when @Salary is null then 31000
                else @Salary
            end
    --Car location
    select @X3= case @Car_location
                when 'carpark' then 41.1
                when 'street' then 325.03
                else 100
            end
    --Output
    set @return = 132.37 + 7.1 * @X1 + 0.01 * @X2 + @X3
    return @return
END
```

4. PMML: Predictive model mark up language

```
<Discretize field="Profit">
  <DiscretizeBin binValue="negative">
    <Interval closure="openOpen" rightMargin="0"/>
    <!-- left margin is -infinity by default -->
  </DiscretizeBin>
  <DiscretizeBin binValue="positive">
    <Interval closure="closedOpen" leftMargin="0"/>
    <!-- right margin is +infinity by default -->
  </DiscretizeBin>
</Discretize>

<MapValues outputColumn="longForm">
  <FieldColumnPair field="gender" column="shortForm"/>
  <InlineTable>
    <row><shortForm>m</shortForm><longForm>male</longForm>
  </row>
    <row><shortForm>f</shortForm><longForm>female</longForm>
  </row>
  </InlineTable>
</MapValues>
```

7.2 Deployment Steps

Implementation and coding into the IT or computer systems

Results of survey on the different ways of deploying models in data mining (source: KDnuggets):

1. Data mining tool (or cloud)

Approx.
45%

2. Programming language (Java, C, VB...)

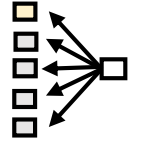
Approx.
15%

3. Database and SQL script

Approx.
25%

4. PMML: Predictive model mark up language

Approx.
15%

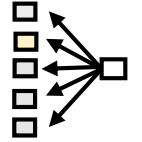


7.2 Deployment Steps

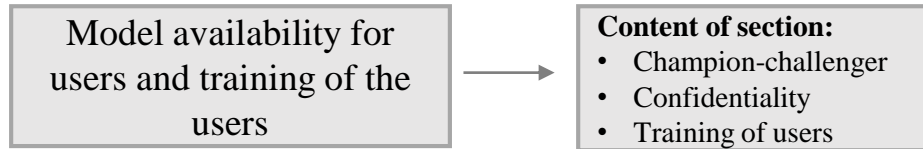
Implementation and coding into the IT or computer systems

Last comments:

- Before requesting data and using it for developing any model, make sure that the data **will continue to be captured in the future**
- Bear in mind the implementation checklist **BEFORE** developing any model
- **Failing in the implementation phase means failing in all the project**

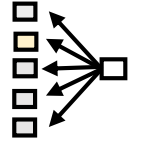


7.2 Deployment Steps



Models will only be available for users after:

- The model is tested in production for several months with volunteers
- The new model runs in parallel to the previous decisioning tool
- Champion-challenger strategies are implemented
- There is enough evidence that the new challenger strategy that uses the new model yields better results
- There is a proper take over of the tool and possibly a phase where the volunteers train other new users



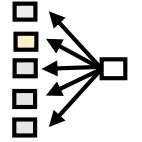
7.2 Deployment Steps

Model availability for
users and training of the
users

The specification of the model (i.e. the different weights for each possible value of the variables that are part of the model) should not be revealed to the end users

Example:

- In order to understand the reason for this secrecy, think for ex. of a sales agent at a mobile phone point of sale. His incentive is to sell as many lines and mobile phones as possible in order to obtain high sales figures and consequently a good bonus
- On the other hand, the risk department is concerned with limiting the credit risk exposure
- These 2 departments have conflicting objectives
- The risk department puts in place the risk model and the corresponding strategy/rules that determine which offer the sales agent is allowed to offer to the prospects
- If the sales agent knew exactly how the prospects are being rating (knew the weights/scores of each variable) then there would be the risk that the prospects are not evaluated properly and therefore, the model would not serve its purpose



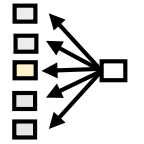
7.2 Deployment Steps

Model availability for
users and training of the
users

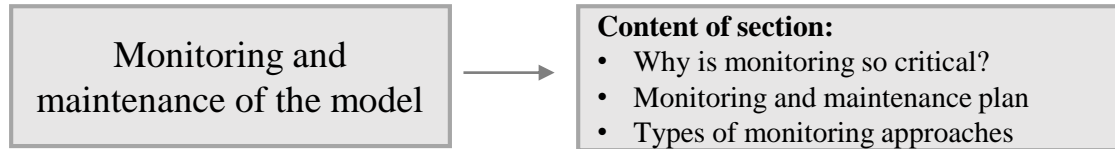
It is essential that the users are well trained and understand the background of the models and the new decisioning tool

The training should cover:

- The objective of the model
- The principles of the tool and how they work
- The limitations
- The methods for using the model
- What will the tool contribute
- What this implementation means from an operational and organizational point of view



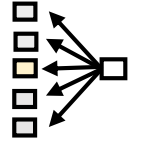
7.2 Deployment Steps



‘Putting your models in auto-pilot is dangerous’, it is key to understand:

- How is the portfolio performing
- If the models are being used in the most effective way
- When is the right time to fine-tune, recalibrate or rerun the models

Therefore, preparing a monitoring and maintenance strategy is of maximum importance. It will help avoid unnecessary long periods of incorrect usage of the data mining results. The strategy will heavily depend on the specific type of deployment of the data mining result.



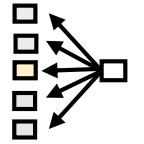
7.2 Deployment Steps

Monitoring and
maintenance of the model

The monitoring and maintenance plan specifies how the deployed results are to be maintained.

Topics to be covered:

- Overview of results deployment and indication of which results may require updating (and why)
- Description of how updating will be triggered (regular updates, trigger event, performance monitoring)
- Description of how updating will be performed
- Summary of the results updating process

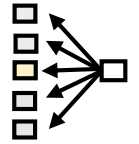


7.2 Deployment Steps

Monitoring and
maintenance of the model

Consider the following activities at the moment of preparing the monitoring and maintenance strategy plan:

- Check for dynamic aspects (i.e., what things could change in the environment?)
- Decide how accuracy will be monitored
- Determine when the data mining result or model should not be used any more. Identify criteria (validity, threshold of accuracy, new data, change in the application domain, etc.), and what should happen if the model or result could no longer be used (update model, set up new data mining project, etc.)
- Will the business objectives of the model change over time? Fully document the initial problem what was the model attempting to solve.

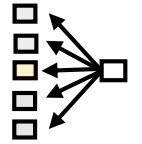


7.2 Deployment Steps

Monitoring and maintenance of the model

There are 2 types of monitoring: One-off and continuous:

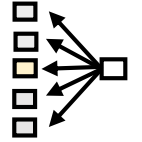
One-Off monitoring	<ul style="list-style-type: none">• Whenever a new data mining application is put into use, the results must be evaluated• Check for ex. that higher scores correspond to lower number of observed defaults or higher response rates• If there are discrepancies, then those cases should be evaluated separately and considering all the characteristics of those individuals (not just the variables that are part of the model)
Ongoing monitoring	<p>1. Population stability - Understand how the target population (e.g. new applicants for credit) changes over time due to:</p> <ul style="list-style-type: none">• New marketing strategies• Changes in product mix• Pricing updates• Changes in collection strategies• Competition• Economy



7.2 Deployment Steps

Monitoring and maintenance of the model

Ongoing monitoring	<p>2. Scorecard performance - Understand the benefits of having the right information available when making decisions for ex. In relation to the following areas:</p> <ul style="list-style-type: none">• Override management• Cutoff changes• Auto-decisioning• Pricing• Scorecard updates• Collection effectiveness• Credit line management• Authorizations management• Reissue management• Retention <p>3. Decision management - Understand how the model degradation may affect the quality of the portfolio:</p> <ul style="list-style-type: none">• Model no longer meets designed purpose• Increased delinquency and losses• Lower approval ratio• Reduced confidence in the model• Increased overrides
---------------------------	---



7.2 Deployment Steps

Monitoring and
maintenance of the model

Example: Consider a **new business scoring** model that has been deployed and now we are the stage of monitoring and maintenance of the model

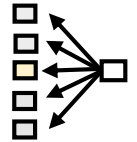
1. Population stability

In order to identify:

- The degree to which the applicant population has shifted over time
- The scorecard components driving the shift
- Acquisition trends

Create the following reports (see examples in next slides):

- Actual versus expected score distribution
- Actual versus expected characteristic distributions
- Mean applicant score
- Volumes, approval



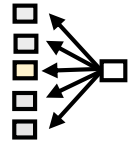
7.2 Deployment Steps

Monitoring and maintenance of the model

Actual vs. Expected Score distribution

Interval	Actual %	Expected %	Diff.
> Low	100.0	100.0	0.0
> 100	92.3	90.0	2.3
> 200	89.9	86.2	3.7
> 300	84.2	79.8	4.4
> 400	74.0	69.2	4.8
> 500	59.7	55.2	4.5
> 600	43.2	40.1	3.1
> 700	29.4	26.6	2.8
> 800	20.5	18.5	2.0
> 900	9.4	8.0	1.4

- Kolmogorov-Smirnov-Test (K-S) difference calculation
 - $= 100 \times K / (N)^{1/2}$
 - K is a constant representing a 95% level of confidence
 - N represents the number of scored accounts
 - This is a one-sample test
- When the maximum observed difference exceeds K-S test difference, the difference between actual and expected distributions is statistically significant
- As a rule of thumb, differences > 20% are of concern, as they sometimes go hand in hand with score degradation
 - **Maximum observed difference = 4.8**
 - **K-S test difference = 6.08**
 - = > No significant change has occurred**



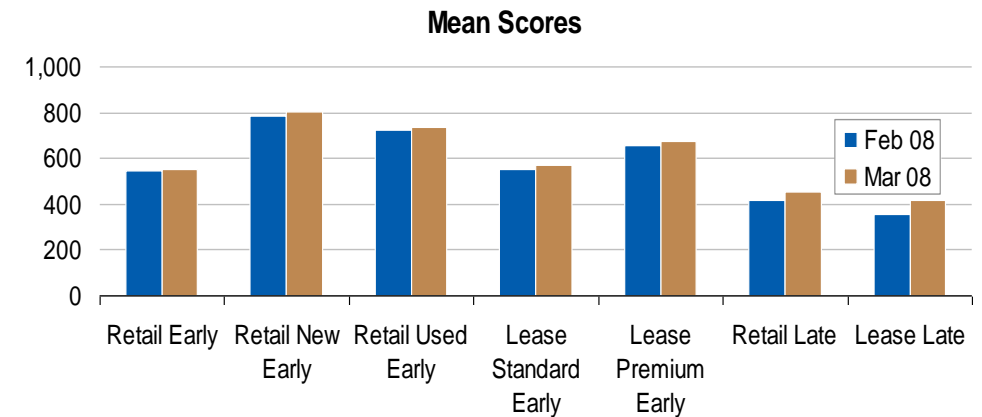
7.2 Deployment Steps

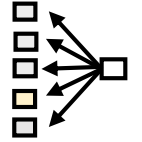
Monitoring and maintenance of the model

Actual vs. Historical Scores

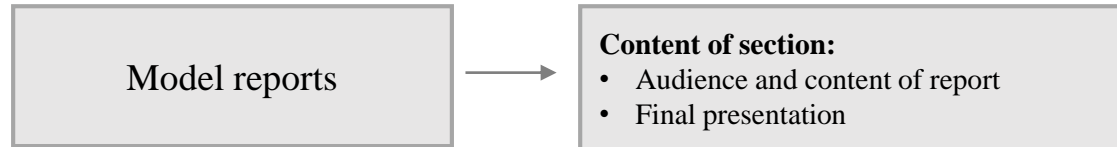
	Mar 2008 %	Feb 2008 %	Feb 2008 Diff	Jan 2008 %	Jan 2008 Diff	Expected %	Expected Diff
0-553	22.27	21.87	0.40	21.09	1.18	21.55	0.72
553-688	37.93	38.16	-0.23	37.55	0.38	38.42	-0.49
688-807	51.83	52.26	-0.43	52.89	-1.06	55.38	-3.55
807-840	69.63	70.28	-0.65	70.13	-0.50	71.24	-1.61
840-871	77.21	79.16	-1.95	78.92	-1.71	87.14	-9.93
871-886	86.74	86.96	-0.22	87.31	-0.57	87.77	-1.03
886-887	87.35	87.53	-0.18	88.24	-0.89	89.89	-2.54
887-921	93.94	93.94	0.00	94.99	-1.05	98.17	-4.23
921-936	98.41	98.44	-0.03	98.32	0.09	98.53	-0.12
936-999	100.00	100.00	0.00	100.00	0.00	100.00	0.00
Total Accts	74,905	70,957		75,981		74,871	
Mean Score:	692	693		698		693	
Std Deviation	230	227		223		224	
K-S Test:	MOD:		1.95		1.71		9.93
	Req Dif:		0.36		0.35		0.35

Mean scores by product over time

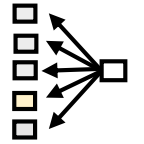




7.2 Deployment Steps



- At the end of the project, the project team writes a final model report
- Depending on the deployment plan, this report may be a summary of the project, or a final presentation of the data mining result(s)
- The actual detailed content of the report depends on the intended audience
- In general, it will describe the results obtained, the process, show which costs have been incurred, define any deviations from the original plan, describe implementation plans, and make any recommendations for future work
- When reaching this stage of the project, identify what reports are needed (slide presentation, management summary, detailed findings etc.)
 - Identify target groups for report
 - Outline structure and contents of report(s)
 - Select findings to be included in the reports

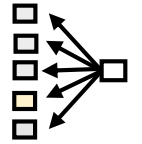


7.2 Deployment Steps

Model reports

As a reference, this type of reports include:

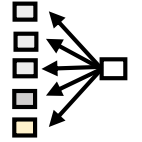
- Description of data received, who delivered the data, when, specific name files, size, number of columns, rows, scripts
- Data treatment, consistency checks, description of how data was handled, tables merged, key identifiers, assumptions etc.
- Data mining result, detailed description of variables used, in what table they appear, possible values that they can take, weights etc.
- Performance assessment including all the KPI studied in previous sections, on developing, hold-out and out-of-time sample
- Implementation and deployment plan and considerations
- Monitoring and maintenance strategy
- Reference to clarifications received from the client in relation to the data and agreed definitions used throughout the project
- Reference to intermediary reports and presentations that were created during the duration of the project



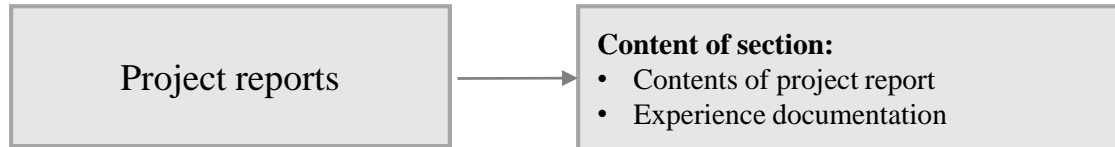
7.2 Deployment Steps

Model reports

- Apart from the final report, it is usually necessary to make a final presentation to summarize the project for example to the project / management sponsor
- The presentation normally contains a subset of the information contained in the final report, structured in a different way
- When preparing the presentation consider:
 - Which is the target group for the final presentation and determine if they will already have received the final report
 - Select which items from the final report should be included in final presentation



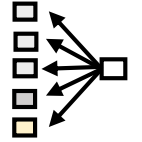
7.2 Deployment Steps



The final report is used to summarize the *entire project and its results*

Typically it contains:

- Summary of business understanding: background, objectives, and success criteria
- Summary of data mining process
- Summary of data mining results
- Summary of results evaluation
- Summary of deployment and maintenance plans
- Cost/benefit analysis
- Conclusions for the business
- Conclusions for future data mining



7.2 Deployment Steps

Project reports

- At this stage of the project, it is also recommended to perform a project review
- The main objective is to assess what went right and what went wrong, what was done well, and what needs to be improved
- As a result of this phase, an “experience documentation” is created to summarize important experience gained during the project. For example, pitfalls, misleading approaches, or tips for selecting the best-suited data mining techniques in similar situations could be part of this documentation. In ideal projects, experience documentation also covers any reports that have been written by individual project members during the project

Lecture Outline

Thus, the lecture will contain:

1. Introduction
2. How to frame the business problem
3. How to transfer it to a problem which can be solved with analytics methods
4. Data identification and prioritisation, data collection and data harmonisation
5. Identification of problem solving approaches and appropriate tools (not only R even though this is important)
6. How to set up and validate models
7. The deployment of a model
- 8. Model lifecycle**
9. Some words about soft skills needed by statistical and mathematical professionals