



Q&A session & feedback Titanic

Process for predictive (and descriptive) analytics

**Project
definition**



**Data
preparation**



**Model
building**



**Model
validation**



**Model
usage**



Process for predictive (and descriptive) analytics

**Project
definition**



**Data
preparation**



**Model
building**



**Model
validation**



**Model
usage**



Make sure to also focus on these

Model validation

Metrics

- **Accuracy** – dependent on probability threshold (scikitlearn's default threshold is 0.5)
 - Specificity, true positive rate (recall), false positive rate
- **AUC** – Area-Under-ROC-Curve, threshold-independent metric

Graphical

- **Lift** – how much better is model performing than a random guess?
- **Cumulative response** – how much better is model performing than a random guess?
- **Cumulative gains** – how many members of the positive group are we reaching with selection?

Process for predictive (and descriptive) analytics

**Project
definition**



**Data
preparation**



**Model
building**



**Model
validation**



**Model
usage**



Iterative process to improve models, with improved
feature engineering and feature preprocessing

Feature selection

Only select relevant features, because:

- Overfitting
- Interpretability
- Efficiency (computational)
- Effort into industrializing more features

Potential approaches:

- Univariate feature selection
- Stepwise feature selection

Univariate variable selection

Variable 1

Variable 2

Variable 3

Variable 4

Variable 5

Variable 6

Build models with one variable

Model
1

Model
2

Model
3

Model
4

Model
5

Model
6

Compute performance metric (ex: AUC)

Performance 1

Performance 2

Performance 3

Performance 4

Performance 5

Performance 6

Select variables with performance higher than threshold

Variable 1

Variable 2

Variable 3

Variable 4

Variable 5

Variable 6

Stepwise forward feature selection

Step 1 : Model with 1 variable

Variable 1

Variable 2

Variable 3

Variable 4

Variable 5

Variable 6

Select variable with best AUC

Variable 1

Variable 2

Variable 3

Variable 4

Variable 5

Variable 6

Stepwise forward feature selection

Step 2 : Model with 2 variables

Variable 2
Variable 1

Variable 2
Variable 3

Variable 2
Variable 4

Variable 2
Variable 5

Variable 2
Variable 6

Select variables with best AUC

Variable 2
Variable 1

Variable 2
Variable 3

Variable 2
Variable 4

Variable 2
Variable 5

Variable 2
Variable 6

Stepwise forward feature selection

Step 3 : Model with 3 variables

Variable 2
Variable 4
Variable 1

Variable 2
Variable 4
Variable 3

Variable 2
Variable 4
Variable 5

Variable 2
Variable 4
Variable 6

Select variables with best AUC

Variable 2
Variable 4
Variable 1

Variable 2
Variable 4
Variable 3

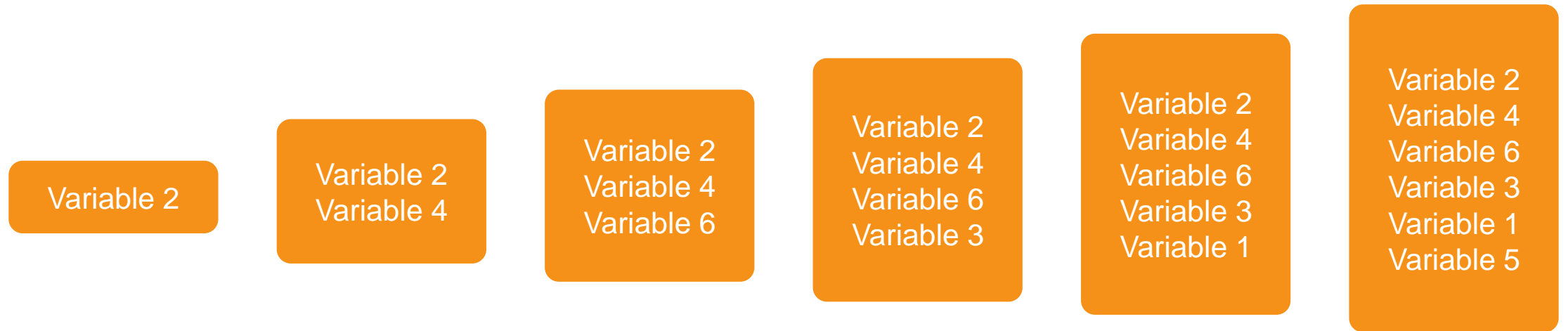
Variable 2
Variable 4
Variable 5

Variable 2
Variable 4
Variable 6

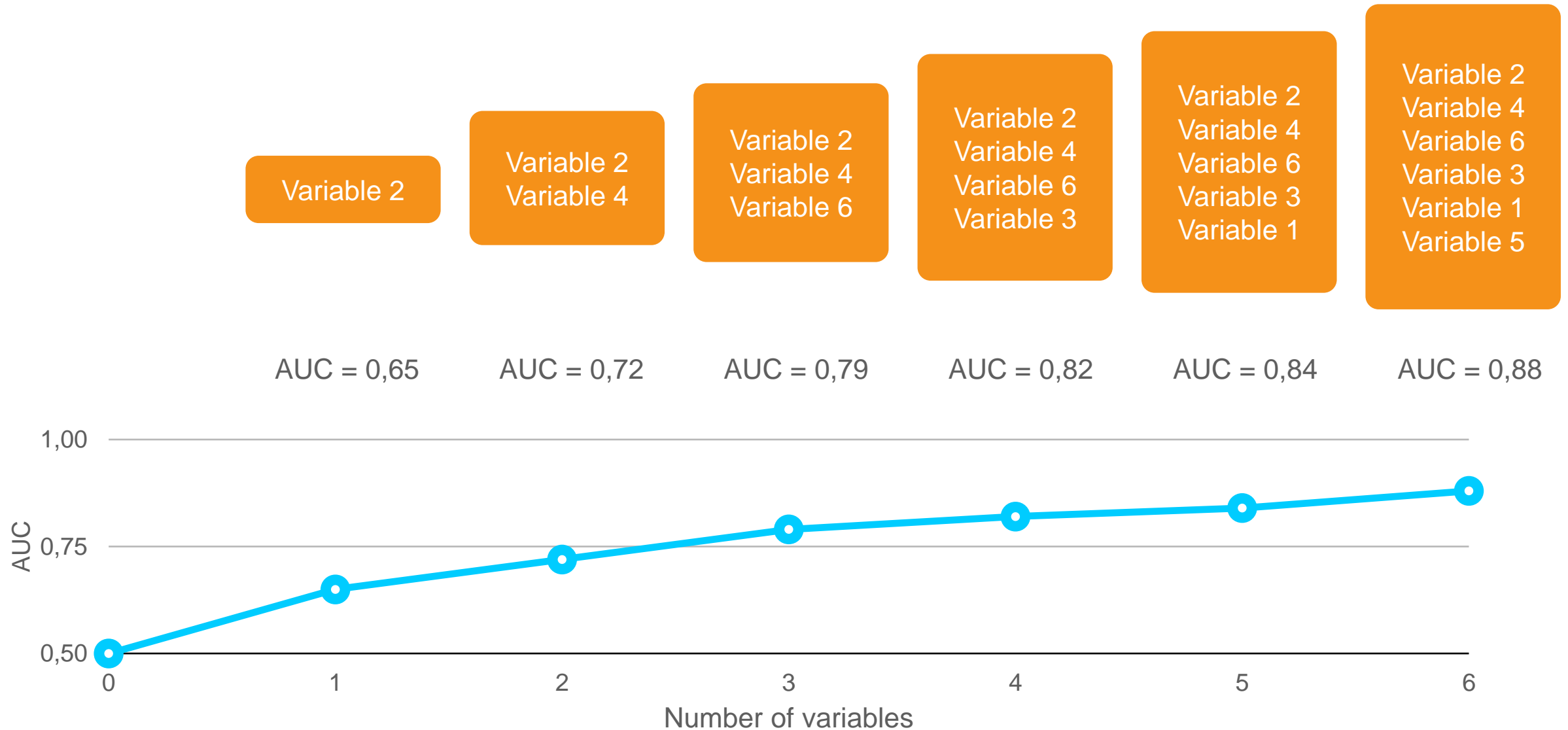
Stepwise forward feature selection

Step N : Model with all variables

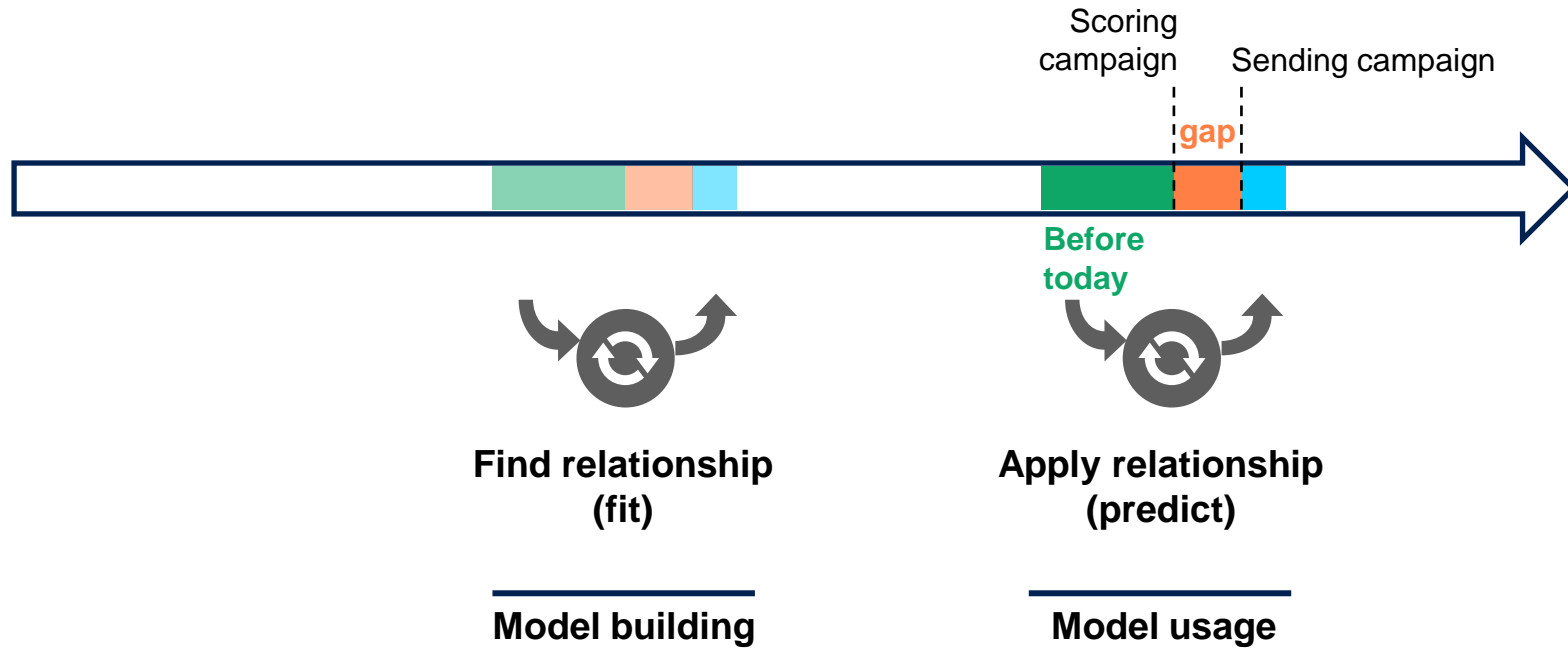
Result : N models, each with 1 additional variable



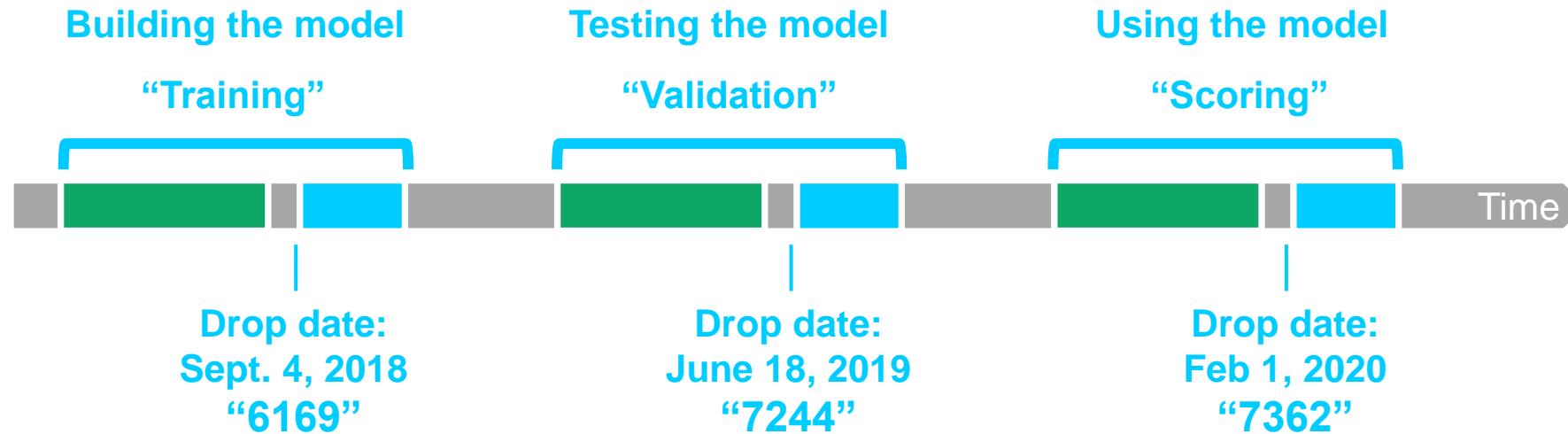
Stepwise forward feature selection



Watch your timeline



Watch your timeline



Business case assumes €0.80 per sent letter
Don't consider the costs from the old campaigns.

How can you push your model further?

1. Better features – do you capture all important things?
2. More training data
3. Double-check timeline
4. Different target (?)
5. Different model (linear vs. non-linear) & hyperparameter tuning

