

Data Mining (Week 1)

dm25s1

Topic 06 : Data Modelling

Motivating Example

Part 02 : Column Encoding

Preparation

Data Handling

Exploring Data 1

Exploring Data 2

Building Models

Dr Bernard Butler

Department of Computing and Mathematics, WIT.
(bernard.butler@setu.ie)

Autumn Semester, 2025

Prediction

Outline

- Encoding categorical features, using pandas
- Looking ahead to feature preparation
- Overview of ML and what we have achieved

Wrap up

Data Mining (Week 6)

Introduction

Motivating Example

Preparation

Data Handling

Exploring Data 1

Exploring Data 2

Building Models

Prediction

Clustering

Regression
1

Classification
1

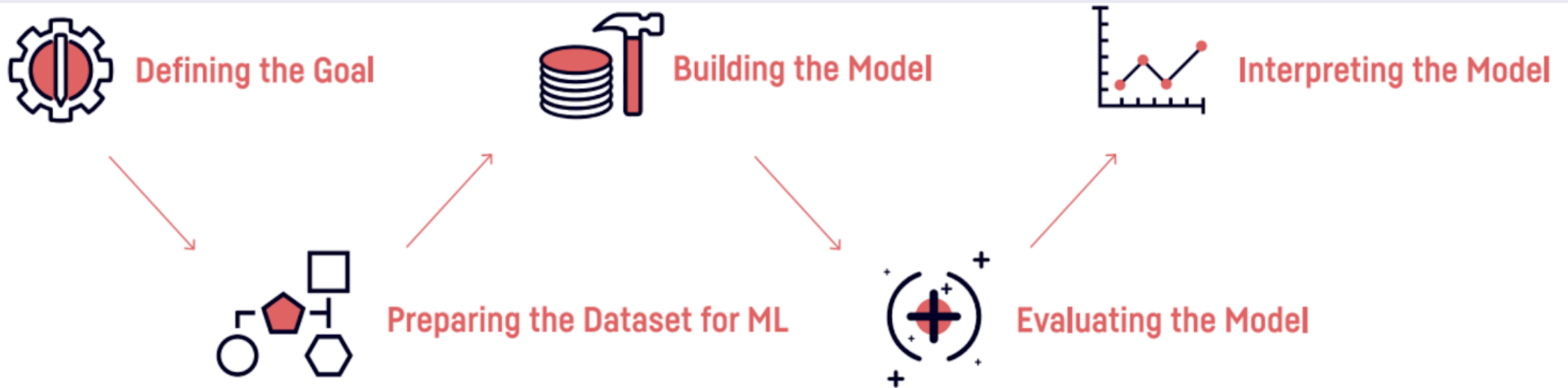
Regression
2

Classification
2

Wrap up

The Pipeline Metaphor

Model Building Pipeline

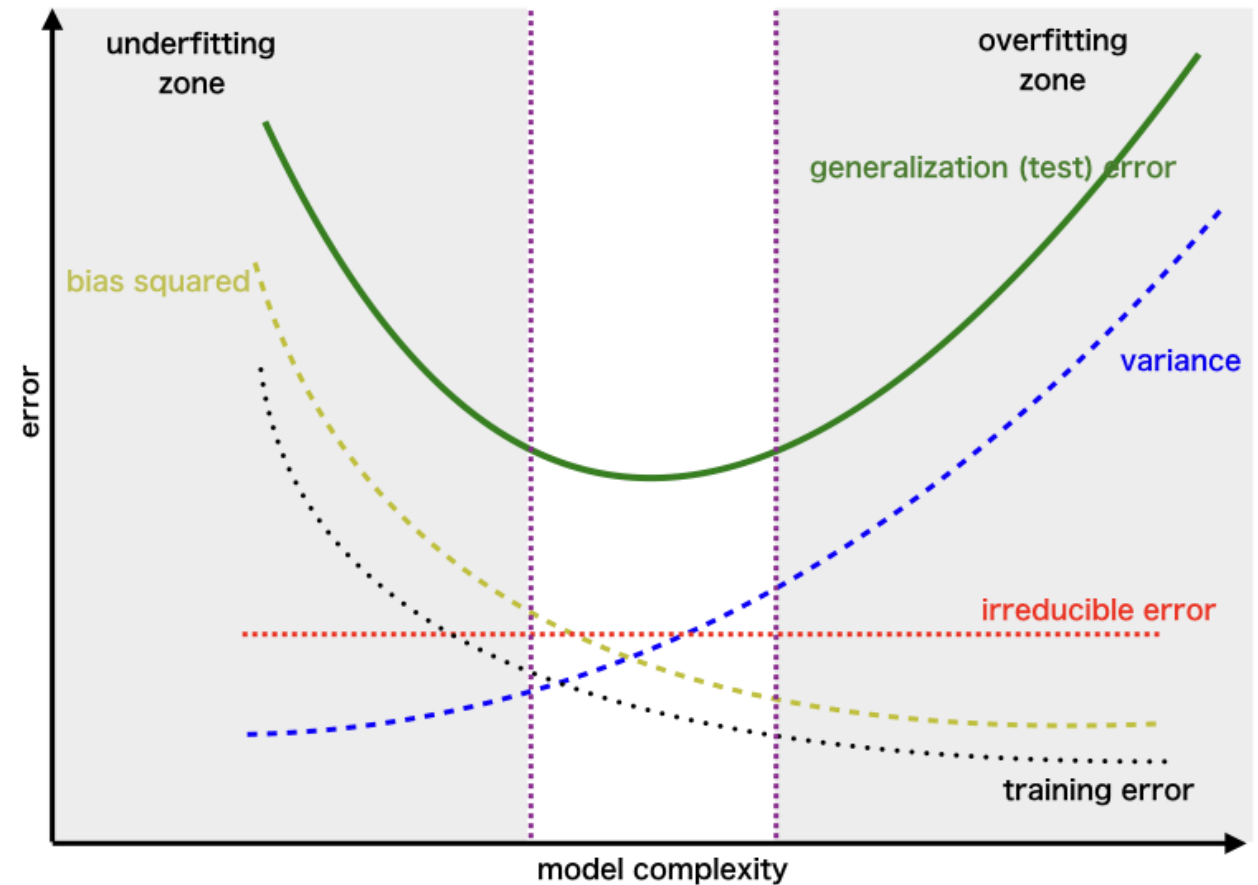
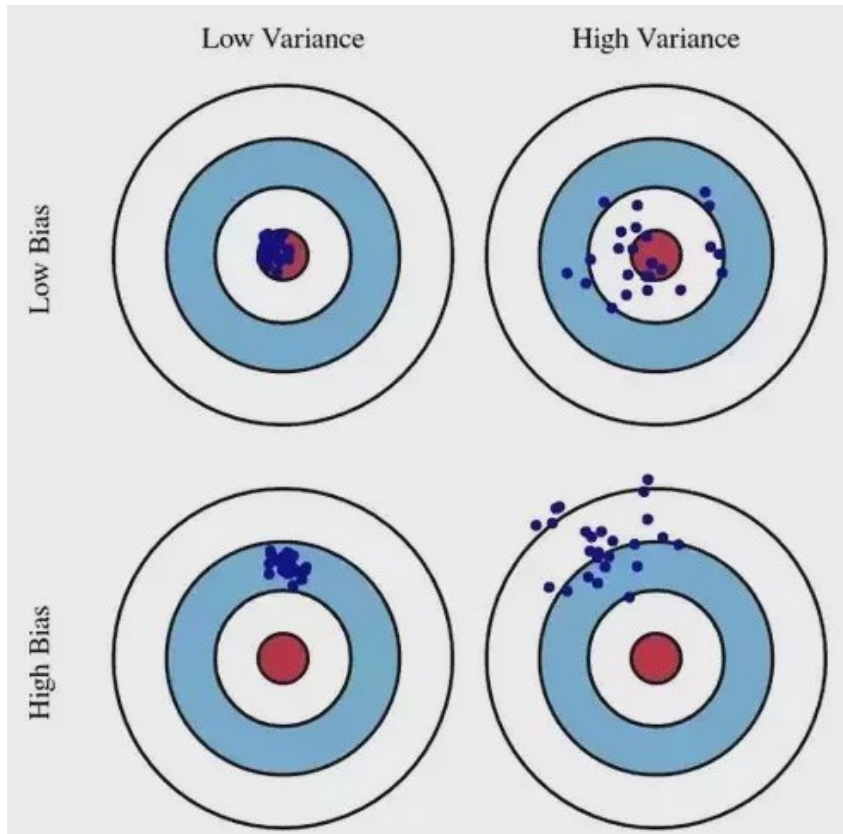


Source: Dataiku

Comments

- We saw the first two stages in previous weeks
- This week we look at the remaining stages
- Of course this pipeline is a simplification. In reality it is iterative.

Bias-Variance and Total Error



Look for parameters \mathbf{a} that minimise the generalization error (estimated using the test set that was not used during training)

Example Model Types

Model	Applications	Concerns
Logistic Regression	X-ray classification	Regression with transformed variable
Fully connected networks	Classification	Classical ANN: choose encoding and size
Convolutional Neural Networks	Image processing	deep learning - choose segmentation
Recurrent Neural Networks	Voice recognition	ANN with feedback - how much?
Random Forest	Fraud Detection	Ensemble method - how many?
Reinforcement Learning	Learning by trial and error	Choose goal and penalties
Generative Models	Text, Image creation	Choose parameters
K-means	Segmentation	Choose distance function and k
k-Nearest Neighbors	Recommendation systems	Choose distance function and k
Bayesian Classifiers	Spam and noise filtering	Deal with imbalances

Using Categorical Features in (Logistic) Regression

How can Categorical-valued features participate in linear models?

Given the following fragment of a dataset, where the goal is to predict the salary of employees in a large organisation:

```
df = pd.read_csv('data/team.csv',\
                 index_col="Name")
df
```

Role Skilled Salary			
Name			
Alice	Designer	Yes	40000
Bob	Programmer	No	25000
Carol	Tester	No	30000

How can this data be represented by a linear model, where all quantities must take numeric values?

Using pandas .getdummies() on a binary-valued column

```
dfSkilledDummies = pd.get_dummies(df['Skilled'],\n    prefix='Skilled',\n    dtype=int)\ndfSkilledDummies
```

	Skilled_No	Skilled_Yes
Name		
Alice	0	1
Bob	1	0
Carol	1	0

Note that a binary-valued column becomes 2 dummy columns

Reducing redundancy (by 1) in 2 dummy columns

```
dfSkilledIndicators = pd.get_dummies(df['Skilled'],\
    prefix='Skilled',\
    drop_first=True,\
    dtype=int)\
    .rename(columns={"Skilled_Yes": "IsSkilled"})
dfSkilledIndicators
```

IsSkilled	
Name	
Alice	1
Bob	0
Carol	0

➤ A single indicator column can replace a group of 2 dummy columns

Using pandas .getdummies() on a multi-valued column

```
dfRoleDummies = pd.get_dummies(df['Role'],\n                                prefix='Role',\n                                dtype=int)\ndfRoleDummies
```

	Role_Designer	Role_Programmer	Role_Tester
Name			
Alice	1	0	0
Bob	0	1	0
Carol	0	0	1

Note that an n -valued column becomes n dummy columns

Reducing redundancy (by 1) in n dummy columns

```
dfRoleIndicators = pd.get_dummies(df['Role'],\
    prefix='Role',\
    drop_first=True,\
    dtype=int)\
    .rename(columns={\
        "Role_Programmer": "IsProgrammer",\
        "Role_Tester": "IsTester"})
```

dfRoleIndicators

	IsProgrammer	IsTester
Name		
Alice	0	0
Bob	1	0
Carol	0	1

$n - 1$ indicator columns can replace a group of n dummy columns

Deriving and using dummy/indicator features

- Identify potential categorical features in EDA Pass 1
- Identify whether each feature is (potentially) *usable* in EDA Pass 2
- Identify whether each feature is (potentially) *useful* in EDA Pass 3
- Add all potentially usable and useful features (regardless of type) to a list F
- For each categorical feature f_j in F having n levels
 - Derive $n - 1$ indicator features \tilde{f}_j^k , where $k = 1, \dots, n - 1$
 - Replace the original categorical feature f_j in F with the derived indicator features \tilde{f}_j^k .
- Build the model using the features in F .

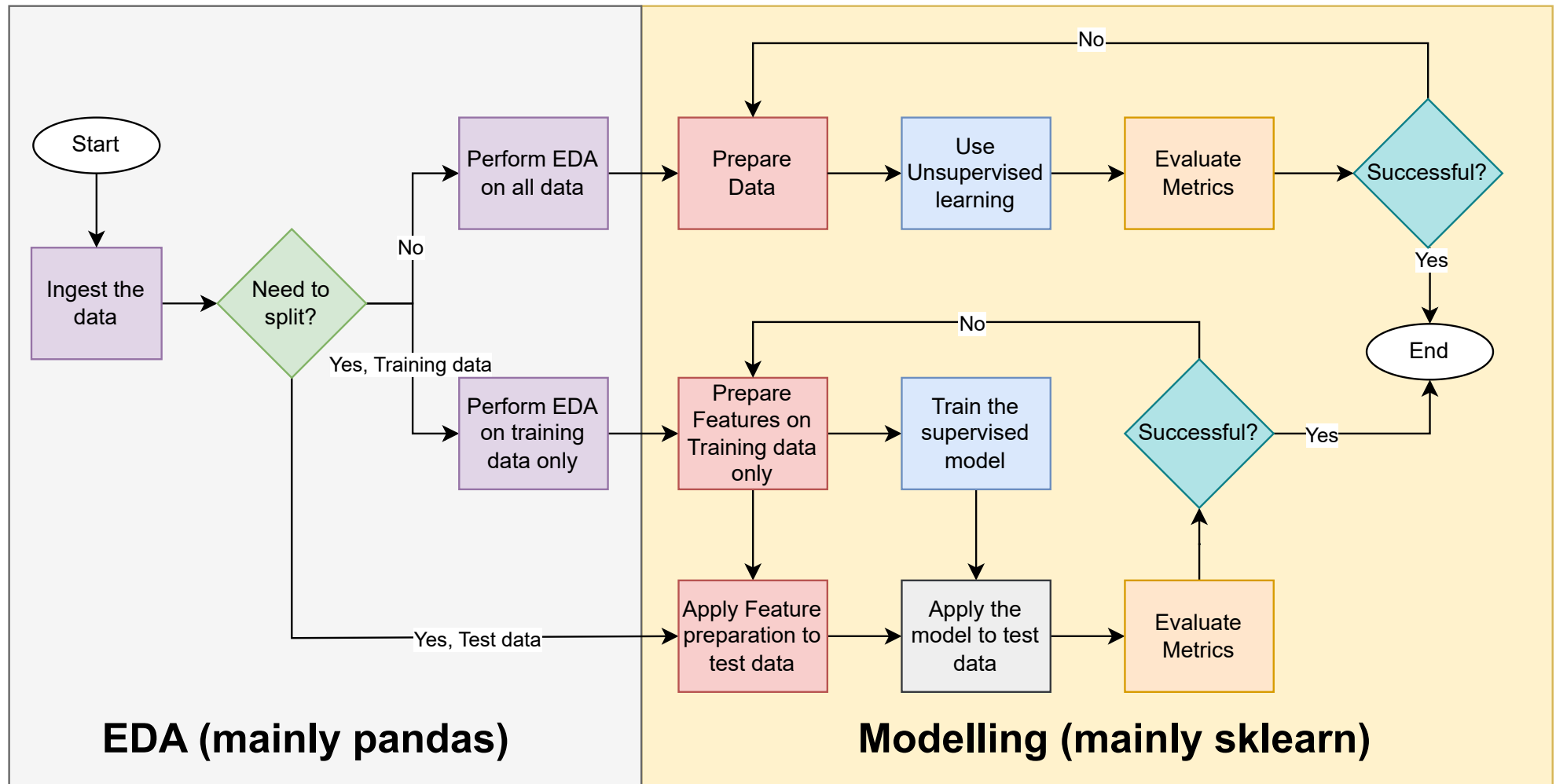
Analysis of pandas dummies

Pandas enables us to convert unordered categorical features into sets of (0,1)-valued numeric features

But what about...

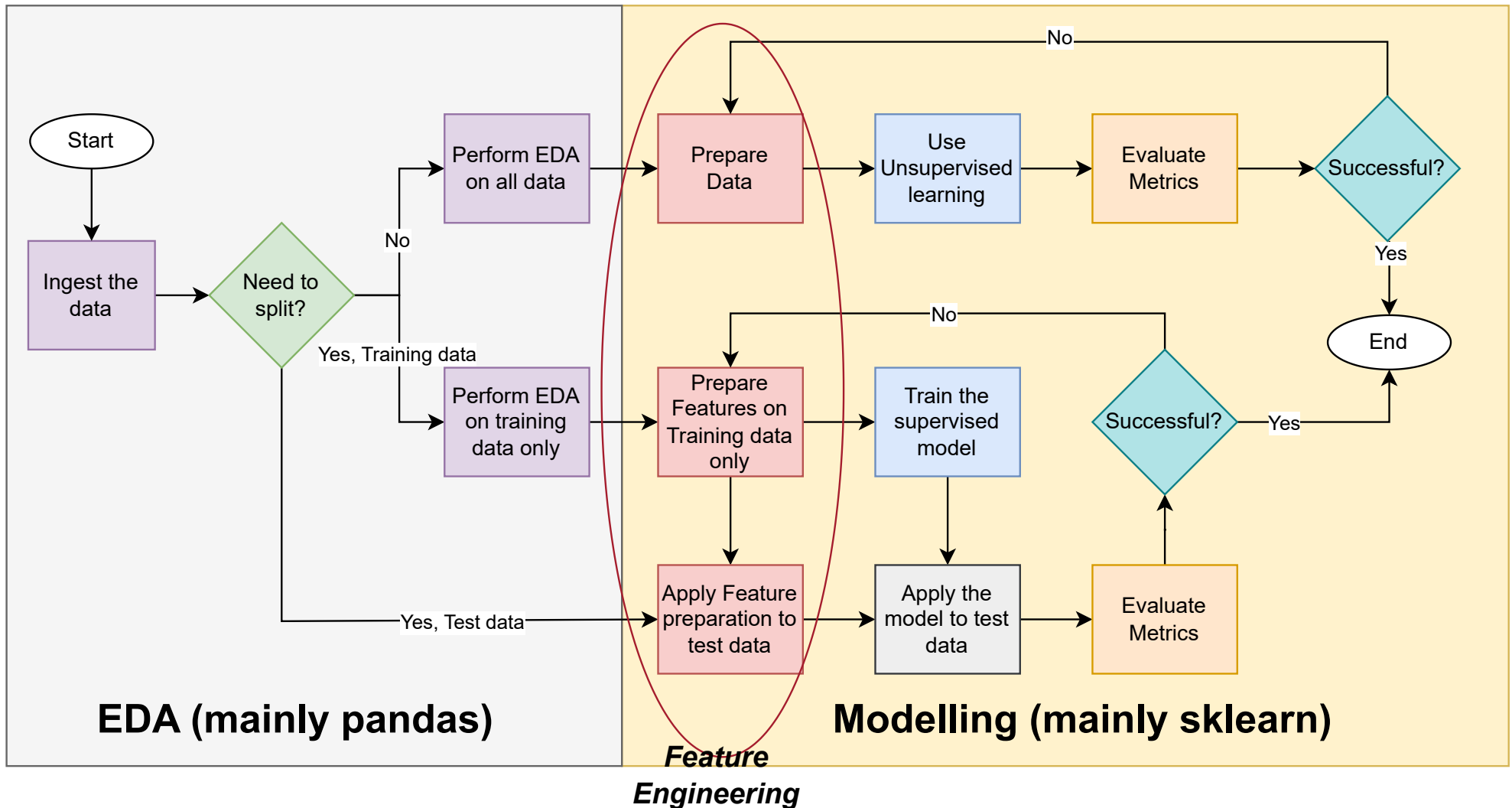
- 1 Ordered categorical features - can we do better than treating them as unordered (and losing information)?
- 2 Categorical targets (whether ordered or not) - how should we handle these?
- 3 How do we handle data (that includes categorical columns) that is split into training and test?
- 4 How can we reverse the operation (i.e., return from (0,1)-valued columns to categorical columns)?

An overview of Machine Learning



On first glance, this might seem overwhelming, but note that the boxes are colour-coded, so that related operations share the same colour.

Where feature preparation fits...



NB: Feature preparation is informed by EDA *but is not part of EDA*. So it is not part of your EDA assignment!

Feature Engineering more generally

- Use of pandas `getdummies()` requires care and is not always suitable
- Scikit-learn, imported as the `sklearn` package, supports a variety of column transformations
- Categorical columns can be features or targets, ordered or unordered
- Can also rescale numerical columns, or encode more exotic columns (other datatypes, computed columns, ...)
- As seen in the schematic, if an ML procedure is unsuccessful, more feature engineering should be considered - it can help.

Summary

- We have reviewed different types of models and considered their general form.
- We looked at the goals of modelling: minimise predictive error.
- We considered how feature engineering can help.
 - Scaling numerical features, so that variation is treated fairly between features.
 - Choosing a subset of features (more to come in future weeks...), looking for the sweet spot between under- and over-fitting.
 - Encoding categorical features as numerical dummy features (more to come in future weeks...), so they can participate in linear models
- In subsequent weeks we will put this theory into practice.

Resources

- **A Summary of the Basic Machine Learning Models**

towardsdatascience.com/a-summary-of-the-basic-machine-learning-models-e0a65627ecbe

- **Train-Test Split for Evaluating Machine Learning Algorithms**

[https://machinelearningmastery.com/](https://machinelearningmastery.com/train-test-split-for-evaluating-machine-learning-algorithms)

[train-test-split-for-evaluating-machine-learning-algorithms](https://machinelearningmastery.com/train-test-split-for-evaluating-machine-learning-algorithms)

This week I have focused on the theory rather than its (python) implementation. This is a nice article that covers the implementation side of things.

- **Cross-Validation: Estimator Evaluator**

medium.com/swlh/cross-validation-estimator-evaluator-897d28afb4ff

Nice article that covers cross-validation in a lot more detail — we will be using many of these variants in later weeks, especially k-fold stratified.