# Overview

Role of data in machine learning

Features and labels

The machine learning workflow
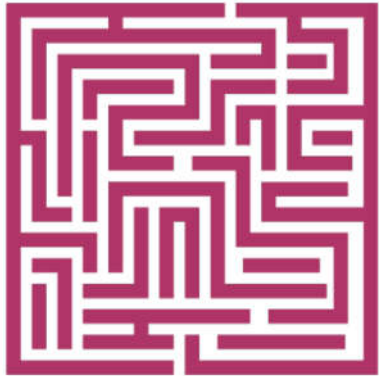
Feature engineering to convert data to features

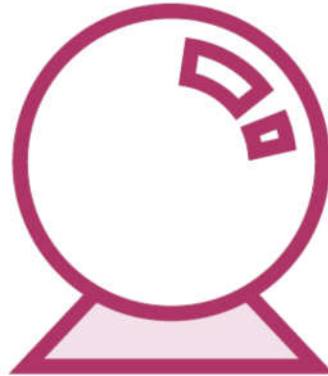Training, test, and validation data

# Features and Labels in Machine Learning

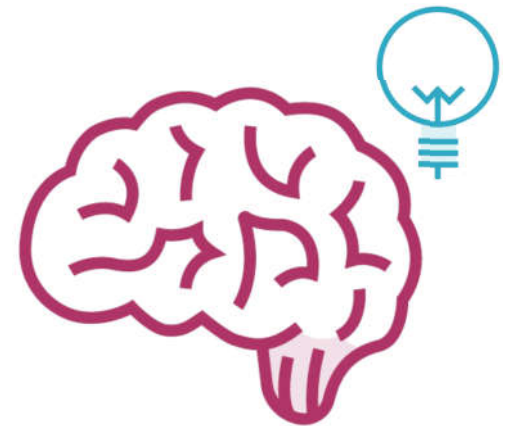A machine learning algorithm is an algorithm that is able to learn from data

# Machine Learning

Work with a huge
maze of data

Find patterns

Make intelligent
decisions

# Types of Machine Learning Problems



**Classification**          **Regression**          **Clustering**          **Dimensionality Reduction**

# Types of Machine Learning Problems

**Classification**

Regression

Clustering

Dimensionality
Reduction

# Whales: Fish or Mammals?

**Mammals**

Members of the infraorder
*Cetacea*

**Fish**

Look like fish, swim like fish,
move with fish

# ML-based Binary Classifier

# ML-based Binary Classifier



Breathes like a mammal

Gives birth like a mammal

ML-based Classifier

Mammal

Corpus

# ML-based Binary Classifier



**Corpus**

**Classification Algorithm**

**ML-based Classifier**

# ML-based Binary Classifier



Breathes like a mammal

Gives birth like a mammal

ML-based Classifier

Mammal

Corpus

# ML-based Binary Classifier

Breathes like a mammal

Gives birth like a mammal

ML-based Classifier

Mammal

Corpus

Input: Feature Vector

# ML-based Binary Classifier



Breathes like a mammal
Gives birth like a mammal

ML-based Classifier

Mammal

Output: Label

Corpus

# ML-based Binary Classifier

Moves like a fish,
Looks like a fish

**ML-based Classifier**

Fish

**Corpus**

# ML-based Binary Classifier

**Moves like a fish,**
**Looks like a fish**

ML-based Classifier

Fish

Corpus

Input: Feature Vector

# ML-based Binary Classifier

Moves like a fish,
Looks like a fish

ML-based Classifier

Corpus

Fish

Predicted Label
≠
Actual Label

# x Variables

The attributes that the ML algorithm focuses on are called features

Each data point is a list - or vector - of such features

Thus, the input into an ML algorithm is a feature vector

Feature vectors are usually called the x variables

# y Variables

The attributes that the ML algorithm tries to predict are called **labels**

**Labels are usually called the y variables**

**Types of labels**

- categorical (classification)

- continuous (regression)

# Garbage In, Garbage Out

If data fed into an ML model is of poor quality, the model will be of poor quality

# The Machine Learning Workflow

# Basic Machine Learning Workflow

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Training algorithm → Fit a model → Model → Choose a validation method → Validation method

Validation method → Examine fit and update → Satisfied?

Satisfied? — No → Update model → Model

Satisfied? — Yes → Use fitted model for predictions

New data → Use fitted model for predictions → Prediction

# What Data Do You Have to Work With?

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

Examine fit and update → Satisfied? → **No** → Update model

**Yes**

New data → Use fitted model for predictions → Prediction

# Load and Store Data

Raw data → **Prepare data** → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

Examine fit and update → Satisfied? — No → Update model

Satisfied? — Yes

New data → Use fitted model for predictions → Prediction

# Data Preprocessing



Raw data → Prepare data → **Cleaned data** → Choose an algorithm → Training algorithm → Fit a model → Model → Choose a validation method → Validation method → Examine fit and update → Satisfied? → No → Update model → Model; Yes → Use fitted model for predictions → Prediction; New data → Use fitted model for predictions

# Selecting and Extracting Features



Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm → Fit a model → Model → Choose a validation method → Validation method → Examine fit and update → Satisfied? → No → Update model; Yes → Use fitted model for predictions → Prediction; New data → Use fitted model for predictions

# Critical and Time-consuming Steps



Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm → Fit a model → Model → Choose a validation method → Validation method → Examine fit and update → Satisfied? — No → Update model; Yes → Use fitted model for predictions → Prediction; New data → Use fitted model for predictions

# Decision Trees, Support Vector Machines?



Raw data → Prepare data → Cleaned data → **Choose an algorithm** → **Training algorithm** → Fit a model → Model → Choose a validation method → Validation method → Examine fit and update → Satisfied? — No → Update model; Yes → Use fitted model for predictions → Prediction; New data → Use fitted model for predictions

# Training to Find Model Parameters

| Raw data | → | Prepare data | → | Cleaned data | → | Choose an algorithm | → | Training algorithm |
|----------|---|--------------|---|--------------|---|---------------------|---|--------------------|

| Validation method | ← | Choose a validation method | ← | **Model** | ← | **Fit a model** |
|-------------------|---|----------------------------|---|-----------|---|-----------------|

| Examine fit and update | → | Satisfied? | —No→ | Update model |
|------------------------|---|------------|------|--------------|

| New data | → | Use fitted model for predictions | → | Prediction |
|----------|---|----------------------------------|---|------------|

Yes

# Evaluate the Model

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← **Choose a validation method** ← Model ← Fit a model

Examine fit and update → Satisfied? → Update model

Satisfied? — No → Update model

Satisfied? — Yes ↓

New data → Use fitted model for predictions → Prediction

# Score the Model

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Training algorithm → Fit a model → Model → Choose a validation method → Validation method

Validation method → **Examine fit and update** → Satisfied?

Satisfied? → No → Update model

Update model → Model

Satisfied? → Yes

New data → Use fitted model for predictions → Prediction

# Different Algorithm, More Data, More Training?

# Iterate Till Model Finalized

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

Validation method → Examine fit and update → **Satisfied?** → **No** → Update model → Model

**Satisfied?** → Yes

New data → Use fitted model for predictions → Prediction

# Model Used for Predictions

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

Examine fit and update → **Satisfied?** → No → Update model

Yes

New data → **Use fitted model for predictions** → **Prediction**

# Retrained Using New Data

# Basic Machine Learning Workflow

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Training algorithm → Fit a model → Model → Choose a validation method → Validation method

Validation method → Examine fit and update → Satisfied?

Satisfied? — No → Update model → Model

Satisfied? — Yes → Use fitted model for predictions

New data → Use fitted model for predictions → Prediction

# Feature Engineering

# Basic Machine Learning Workflow

Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

↑ Training algorithm → Fit a model

Validation method → Examine fit and update → Satisfied?

Satisfied? —No→ Update model → Model (↑)

Satisfied? —Yes↓

New data → Use fitted model for predictions → Prediction

# Selecting and Extracting Features



Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

Examine fit and update → Satisfied? → No → Update model

Yes

New data → Use fitted model for predictions → Prediction

# Feature Engineering



Raw data → Prepare data → Cleaned data → Choose an algorithm → Training algorithm

Validation method ← Choose a validation method ← Model ← Fit a model

Examine fit and update → Satisfied? — No → Update model

Yes

New data → Use fitted model for predictions → Prediction

# Feature Engineering

Engineering your features so that you get the best out of your ML model.

# Feature Engineering



Block and tackle work

Bespoke - specific to:

- Problem

- Data

Not quite art, not quite science...

...More just engineering

# Scope of Feature Engineering

Feature selection

Feature learning

Feature extraction

Feature combination

Dimensionality reduction

# Scope of Feature Engineering

| Feature selection | Feature learning | Feature extraction |
|---|---|---|

| Feature combination | Dimensionality reduction |
|---|---|

# Feature Selection

Choosing the best subset from within an existing set of features (x-variables), without substantially transforming them.

# Choosing Feature Selection

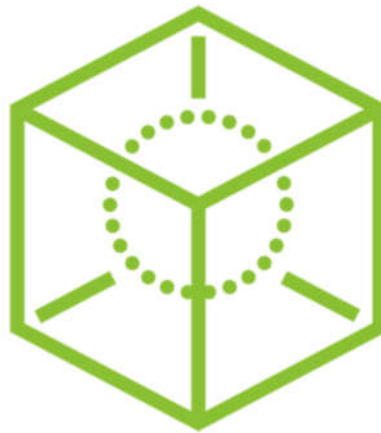| Use Case | Possible Solution |
|---|---|
| Many X-variables | |
| Most of which contain little information | Feature selection |
| Some of which are very meaningful | |
| Meaningful variables are independent of each other | |

# Feature Selection Techniques



**Filter
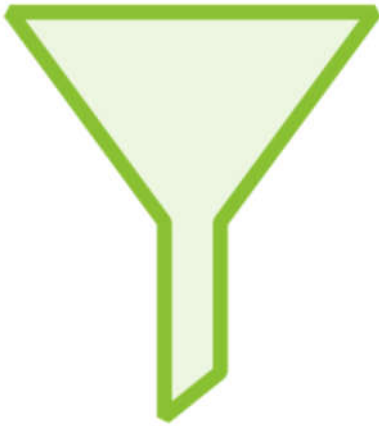methods**

**Embedded
methods**

**Wrapper
methods**

# Filter Methods

**Applying statistical techniques to select the most relevant features**

# Embedded Methods



**Relevant features selected by training a machine learning model i.e. Lasso regression, decision trees**

# Wrapper Methods

**Build candidate models by selecting feature subsets - choose the subset which gives the best model**

# Scope of Feature Engineering

| | | |
|---|---|---|
| **Feature selection** | **Feature learning** | **Feature extraction** |
| | **Feature combination** | **Dimensionality reduction** |

# Feature Learning

Rely on ML algorithms rather than human experts to "learn" the best representations of complex data such as images, videos.

(Also known as Representation Learning)

# Supervised Feature Learning

Features are learnt using labeled data

Neural networks are classic example

Greatly reduce need for expert judgment

**"Traditional"** ML-based systems rely on experts to decide what features to pay attention to

**"Representation"** ML-based systems figure out by themselves what features to pay attention to

Neural networks are examples of such systems

# Unsupervised Feature Learning

**Features need to be learned in absence of labeled corpus**

- Clustering

- Dictionary learning

- Autoencoders

# Scope of Feature Engineering

Feature selection

Feature learning

Feature extraction

Feature combination

Dimensionality reduction

# Feature Extraction

Differs from feature selection in that input features are fundamentally transformed into derived features, which are often unrecognizable and hard to interpret.

# Feature Extraction

Image descriptors for images

Principal components for matrices

Tf-Idf for documents

# Feature Extraction



Feature extraction usually also leads to dimensionality reduction

However explicit objective is to re-express feature in a "better" form

Not to reduce number of X columns

# Scope of Feature Engineering

Feature selection

Feature learning

Feature extraction

Feature combination

Dimensionality reduction

# Feature Combination

Some features naturally work better when considered together

Original feature might be raw or too granular

Improve the predictive power of features

# Feature Combination

**Feature cross in predicting traffic**

- Day-of-week

- Time-of-day

**Feature cross in predicting temperature**

- Season

- Time-of-day

# Scope of Feature Engineering

Feature selection

Feature learning

Feature extraction

Feature combination

Dimensionality reduction

# Dimensionality Reduction

Apply pre-processing algorithms to reduce complexity of raw features

Specifically aim to reduce number of input features

Excessive number of features leads to severe problems

- Curse of Dimensionality

# Dimensionality Reduction

Dimensionality reduction explicitly aim to solve Curse of Dimensionality

While also preserving as much information as possible

Form of unsupervised learning

# Dimensionality Reduction

Principle Components Analysis (PCA)

Manifold Learning

Latent Semantic Analysis

Autoencoding

# Training, Test and Validation Data

# Basic Machine Learning Workflow

```
Raw data  →  Prepare data  →  Cleaned data  →  Choose an algorithm  →  Training algorithm
                                                                              ↓
Validation method  ←  Choose a validation method  ←  Model  ←  Fit a model
        ↓
Examine fit and update  →  Satisfied?  →[No]→  Update model
                               ↓[Yes]                ↑
New data  →  Use fitted model for predictions  →  Prediction
```

# Validate and Iterate Till Model Finalized

# Data

| All data |
|:---:|

**All the data available**

# Training Data

| All data |
|---|

**Use all data to train your model**

# Training Data

Data used to train a model cannot be used to **evaluate** a model

Model may have memorized training instances

Model robustness cannot be measured on instances it has seen before

# Training Data, Test Data

| All data |
|:---:|

| Training data | Test data |
|:---:|:---:|

**Typically 80% of the data used to train the model**

# Training Data, Test Data

| All data |
|:---:|

| Training data | Test data |
|:---:|:---:|

**20% set aside to sanity-check or measure model performance**

# Training Data, Test Data

| All data |
| Training data | Test data |

**One training process to
generate one candidate model**

# Training Data, Test Data

| All data |
|:---:|

| Training data | Test data |
|:---:|:---:|

**For N candidate models, run N training and
N test processes**

# Training Data, Test Data

**Test set can be used to choose the best candidate model**

**Model evaluation on instances the model has not seen during training**

Evaluation can become biased

# Overfitting on Test Set

Choosing best candidate model on the Test Set
leads to this form of overfitting. Occurs when data
is split into just two sets: Training and test.

# Cross-validation

Carve out a separate validation set of data points; use this to evaluate different candidate models. Data now split into three sets: Training, validation and test.

# Training, Test, Validation Data

| All data |
|:---:|

| Training data | Validation data | Test data |
|:---:|:---:|:---:|

**Hold out 2 subsets of the original data, validation data and test data**

# Training, Test, Validation Data

| All data |
|----------|

| Training data | Validation data | Test data |
|---------------|-----------------|-----------|

**Training data to produce candidate models -
validation data to evaluate models**

# Training, Test, Validation Data

| All data |
|:---:|

| Training data | Validation data | Test data |
|:---:|:---:|:---:|

**Test data applied to the selected model to
provide an unbiased evaluation of the final model**

# Training, Test, Validation Data

| All data |
|:---:|

| Training data | Validation data | Test data |
|:---:|:---:|:---:|

**Now can have multiple candidate models, and
select the best one - Hyperparameter Tuning**

# Training, Test, Validation Data

| All data |
|:---:|

| Training data | Validation data | Test data |
|:---:|:---:|:---:|

**For N candidate models, run N training and N validation processes but just 1 test process**

# Singular Cross-validation

| All data |
|:---:|

| Training data | Validation data | Test data |
|:---:|:---:|:---:|

# Singular Cross-validation

| | | | |
|---|---|---|---|
| **All data** | | | |

| | | |
|---|---|---|
| Training data | Validation data | Test data |

Candidate Model 1 — Training data — Validation data

Candidate Model 2 — Training data — Validation data

Candidate Model 3 — Training data — Validation data

Candidate Model 4 — Training data — Validation data

Candidate Model 5 — Training data — Validation data

Hyperparameter Tuning

Test data

# Singular Cross-validation

| All data |
|---|

| Training data | Validation data | Test data |
|---|---|---|

**Candidate Model 1**

| Training data | Validation data |
|---|---|

# Singular Cross-validation

| All data | | |
|---|---|---|

| Training data | Validation data | Test data |
|---|---|---|

**Candidate Model 1** | Training data | Validation data |

**Candidate Model 2** | Training data | Validation data |

# Singular Cross-validation

| All data | | |
|---|---|---|
| Training data | Validation data | Test data |

| | Training data | Validation data |
|---|---|---|
| Candidate Model 1 | Training data | Validation data |
| Candidate Model 2 | Training data | Validation data |
| Candidate Model 3 | Training data | Validation data |
| Candidate Model 4 | Training data | Validation data |
| Candidate Model 5 | Training data | Validation data |

# Singular Cross-validation

| All data | | |
|---|---|---|
| Training data | Validation data | Test data |

| | Training data | Validation data | |
|---|---|---|---|
| Candidate Model 1 | Training data | Validation data | |
| Candidate Model 2 | Training data | Validation data | |
| Candidate Model 3 | Training data | Validation data | Hyperparameter Tuning |
| Candidate Model 4 | Training data | Validation data | |
| Candidate Model 5 | Training data | Validation data | |

# Singular Cross-validation

The model's performance on the validation set is incorporated into the model itself - this may introduce bias

# K-fold Cross-validation

For each candidate model, repeatedly train, and validate using different subsets of training data. Much more computationally intensive, but very robust - does not "waste" data.

# K-fold Cross-validation

For each candidate model, repeatedly train, and validate using different subsets of training data. Much more computationally intensive, but very robust - does not "waste" data.

# K-fold Cross-validation

For each candidate model, repeatedly train, and validate using different subsets of training data. Much more computationally intensive, but very robust - does not "waste" data.

# K-fold Cross-validation

| All data | |
|:---:|:---:|
| **Training data** | **Test data** |

# K-fold Cross-validation

| All data |
|---|

| Training data | Test data |
|---|---|

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|

# K-fold Cross-validation

| All data |
|:---:|

| Training data | Test data |
|:---:|:---:|

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|:---:|:---:|:---:|:---:|:---:|

Split 1

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|:---:|:---:|:---:|:---:|:---:|

# K-fold Cross-validation

| All data | | | | |
|---|---|---|---|---|

| Training data | | | | Test data |
|---|---|---|---|---|

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|

Split 1

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|

Split 2

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|

# K-fold Cross-validation

| All data | | | | |
|---|---|---|---|---|

| Training data | | | | Test data |
|---|---|---|---|---|

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|

| | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|---|
| Split 1 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 2 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 3 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

# K-fold Cross-validation

| All data | | | | |
|---|---|---|---|---|

| Training data | | | | Test data |
|---|---|---|---|---|

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|

| | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|---|
| Split 1 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 2 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 3 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 4 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 5 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

**Find the best candidate model**

# K-fold Cross-validation

| All data | | | | |
|---|---|---|---|---|

| Training data | | | | Test data |
|---|---|---|---|---|

| | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
|---|---|---|---|---|---|
| Split 1 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 2 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 3 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 4 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |
| Split 5 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

**Train each candidate model and average performance**

**Test data**