

How to start a machine learning project - a practical example

Goal

What?

Why?

How?

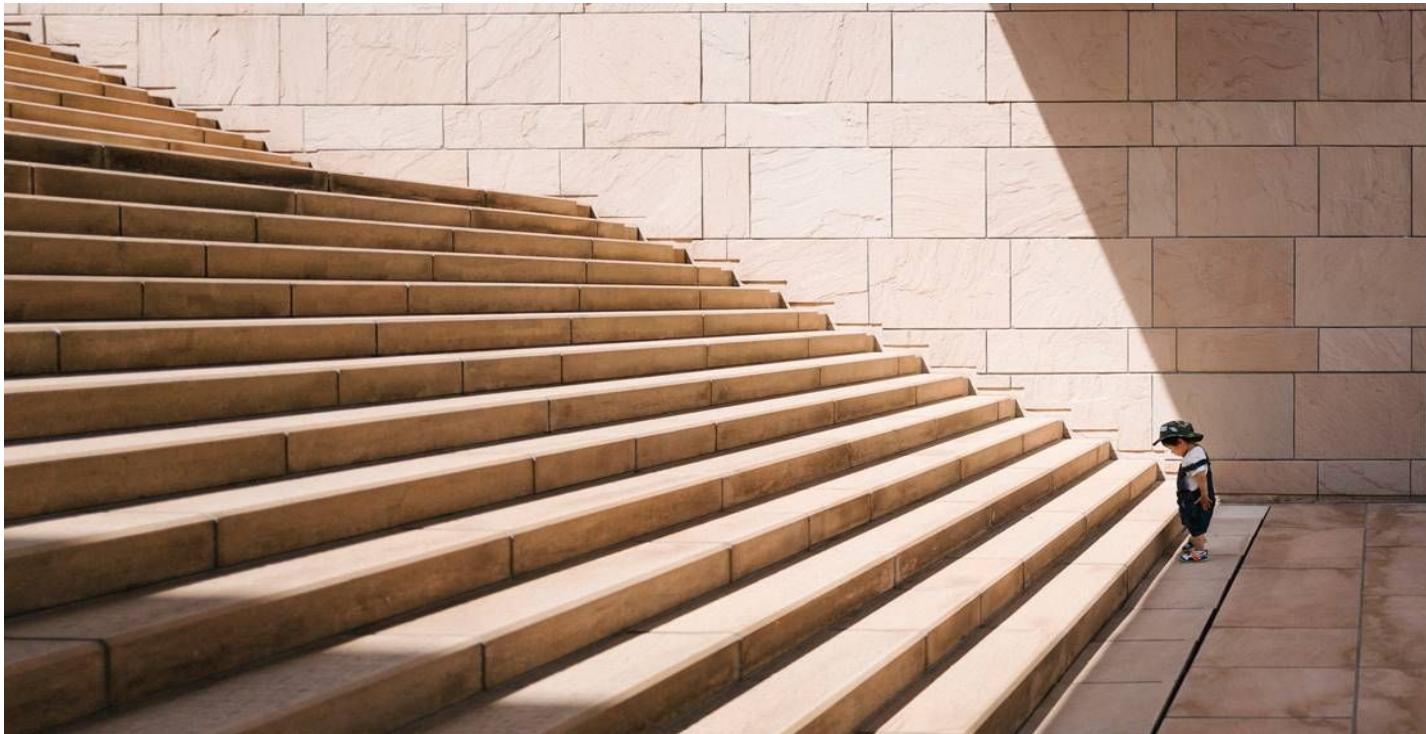
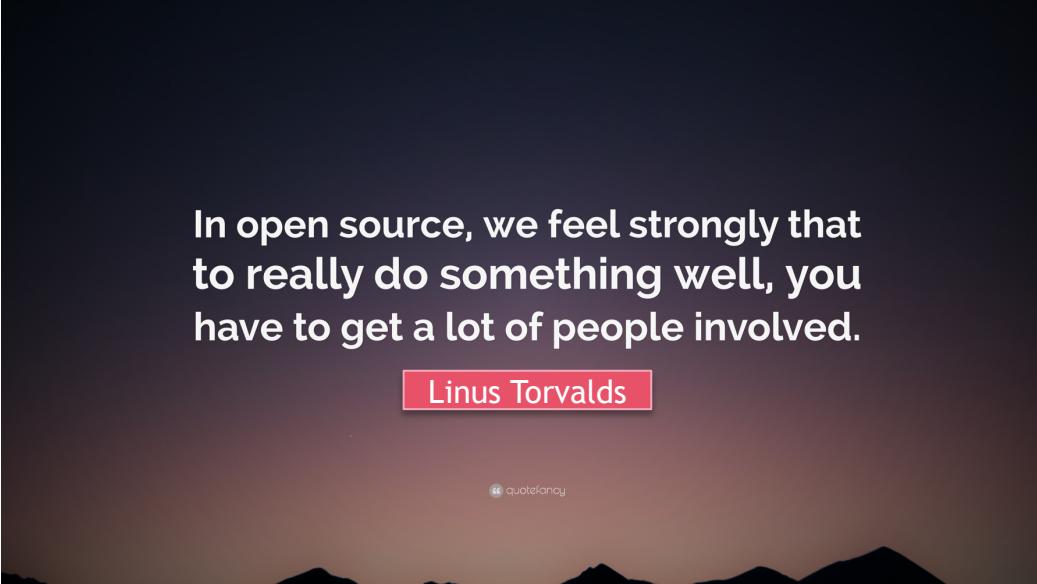


Image source: unknown

Conflicts of interest

None



In open source, we feel strongly that to really do something well, you have to get a lot of people involved.

Linus Torvalds

Outline

- End-to-end machine learning project
- Only one of many (many!) approaches...
- Snippets of code



Outline

```
# To support both python 2 and python 3
from __future__ import division, print_function, unicode_literals

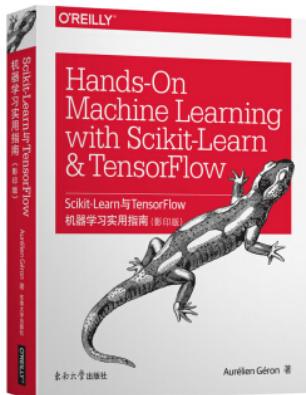
# Common imports
import numpy as np
import os
import pandas as pd

# To make this notebook's output stable across multiple runs
np.random.seed(42)

# To plot pretty figures
%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt
mpl.rc('axes', labelsize=14)
mpl.rc('xtick', labelsize=12)
mpl.rc('ytick', labelsize=12)
```

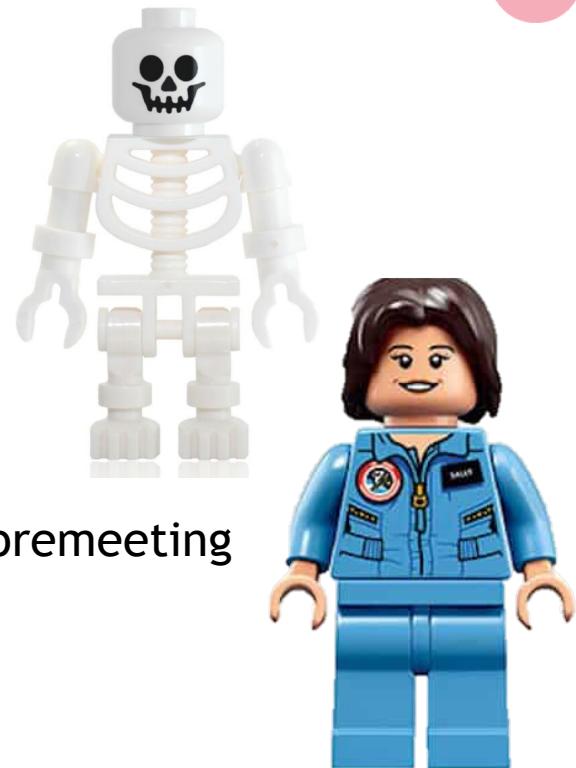
Outline

- End-to-end machine learning project
- Only one of many (many!) approaches...
- Snippets of code, available on...



github.com/JDhont/ESTRO38_premeeting

Python.org
Jupyter.org
Scikit-learn.org



An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
 2. Feature engineering
 3. Get the data - data pre-processing
 4. Creating a training, (validation) and test set
 5. Feature selection
 6. Model selection
 7. Fitting the model
 8. Model evaluation
 9. QA - implement
 10. share!
-
-

An end-to-end machine learning project

Step 1: Frame the problem

In-house experience, confirmed by literature

4DCT significantly underestimates breathing-induced
tumour motion during treatment

Margins Determined Using 4DCT Often Underestimate
Tumor Motion in Thoracic Tumors

X. Shi, S. Chen, W.D. D'Souza, N.N. Mistry
University of Maryland School of Medicine, Baltimore, MD

The long- and short-term variability of breathing induced tumor motion
in lung and liver over the course of a radiotherapy treatment 

Jennifer Dhont ^{a,b,c,d,*}, Jef Vandemeulebroucke ^{a,b}, Manuela Burghelea ^d, Kenneth Poels ^e, Tom Depuydt ^e,
Robbe Van Den Begin ^d, Cyril Jaudet ^d, Christine Collen ^d, Benedikt Engels ^d, Truus Reynders ^d,
Marlies Boussaer ^d, Thierry Gevaert ^d, Mark De Ridder ^d, Dirk Verellen ^{c,f}

Both four-dimensional computed tomography and four-dimensional cone beam computed tomography under-predict lung target motion during radiotherapy

Implications using ITV technique?

Elisabeth Steiner ^{a,*}  , Chun-Chien Shieh ^a, Vincent Cailliet ^{a,b}, Jeremy Booth ^{b,c}, Ricky O'Brien ^a, Adam Briggs ^b, Nicholas Hardcastle ^{d,e}, Dasantha Jayamanne ^b, Kathryn Szymura ^b, Thomas Eade ^{b,f}, Paul Keall ^a

An end-to-end machine learning project

Step 1: Frame the problem

HYPOTHESIS

*Local tumour control probability of thoracic tumours is related to motion amplitude
and type of motion management*

GOAL

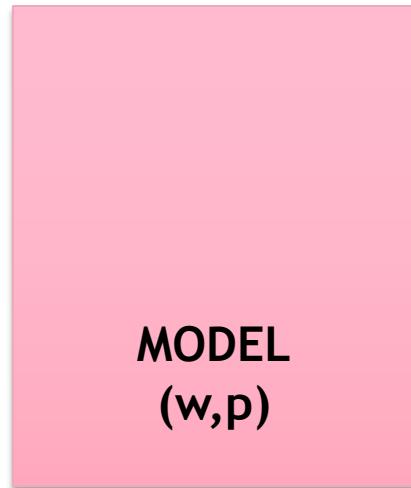
Two-fold (1) *To predict local recurrence, (2) evaluate relevance of motion
amplitude and type of motion management.*

GOAL

Two-fold (1) To predict local recurrence, (2) evaluate relevance of motion amplitude and type of motion management.



Characteristics
(a,b,c,d,e,f,...)



Local control
or
Local recurrence

An end-to-end machine learning project

Step 1: Frame the problem & look at the big picture

What type of problem is this?

1. Supervised, unsupervised, or reinforcement learning?
2. Is it a classification task or a regression task?

Supervised: fit a model on available input and output data

Classification: will this patient have a local recurrence or not?

Machine learning

Supervised learning in a classification task



Characteristics
(a,b,c,d,e,f,...)



Local control
or
Local recurrence

An end-to-end machine learning project

Step 1: Frame the problem & look at the big picture

What type of problem is this?

1. Supervised, unsupervised, or reinforcement learning?
2. Is it a classification task or a regression task?
3. Performance measure?

An end-to-end machine learning project

Step 1: Frame the problem & look at the big picture

Performance measure

Regression task

- RMSE
- MAE

Classification task

- Sensitivity - % of patients with local recurrence identified as such
 - Specificity - % of patients without local recurrence identified as such
-
-

An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
 2. Feature engineering
 3. Get the data - data pre-processing
 4. Creating a training, (validation) and test set
 5. Feature selection
 6. Model selection
 7. Fitting the model
 8. Model evaluation
 9. QA - implement
 10. share!
-
-

An end-to-end machine learning project

Step 2: Feature engineering

PREDICT TCP

Model input?

Tumor biology, tumor microenvironment, radiation dosimetry, patient characteristics, imaging features...

“Coming up with features is difficult, time-consuming and requires expert knowledge.” - Andrew Ng

An end-to-end machine learning project

Step 2: Feature engineering - MI vs. DL

Feature formulation

- Tumor biology
- Tumor microenvironment
- Radiation dosimetry
- Patient characteristics
- Imaging features...



Storing raw data



Feature selection

What is available ???

An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
 2. Feature engineering
 3. Get the data - data pre-processing
 4. Creating a training, (validation) and test set
 5. Feature selection
 6. Model selection
 7. Fitting the model
 8. Model evaluation
 9. QA - implement
 10. share!
-
-

An end-to-end machine learning project

Step 3: Get the data

CENSORED

dossier	Institution	nr_per_RT	Lesion_nr	GTV_nr	Location_exact	Location_group	Node	Technique	prescription	Metastasis % % Gy Gy Gy Gy Gy						Mode 4D (mm)				prot_v isolate n_LC	start_RT				
										GTV-ITV_50Gy	ITV_D95_50Gy	PTV_50Gy	PTV_D95_50Gy	PTV_D98_50Gy	GTV_D95_50Gy	GTV_D98_50Gy	GTV_cc	PTV	LR	AP	CC	EUC_L	comments		
A4502011E00T	1	1	2	2	RBK	1	0	ITV_Vero	50Gy 80%	57,3	99,7	65,8	51,6	67,1	60,5	63,5	0,76	555	-2,0	7,0	6,0	9,4	0	10/07/2012	
A460211WNQ01R	1	1	1	1	RMK	1	0	ITV_Vero	50Gy 80%	81,9	99,9	65,9	51,9	67	60,3	64,2	2,2	555	-2,1	1,6	3,3	4,2	0	03/07/2012	
B430615YMD01R	1	1	5	1	lever S7	2	0	tracking	50Gy 80%	78,1	99,1	65,4	51,4	65,6	60,8	63,6	19	555	2,4	0,0	6,0	6,5	0	09/08/2012	
A420713DE00L	1	1	6	1	klier retroper	4	1	ITV_Vero	50Gy 80%	48,5	99,2	64,6	50,9	65	60,4	50,5	4,8	555	-1,4	2,2	3,0	4,0	0	22/07/2012	
A420713DE00L	1	1	6	1	ROK	1	0	ITV_Vero	50Gy 80%	81,1	99,4	66,3	51	67	60,8	64,3	3,7	555	-0,7	-2,6	7,9	8,3	0	12/09/2012	
A360225HU00P	1	1	7	1	RBK	1	0	ITV_Vero	50Gy 80%	65	100	65,1	52	65,3	57,5	61,7	1,3	555	0,5	-1,1	1,4	1,8	0	03/09/2012	
A340526KG00Q	1	2	8	2	RBK	1	0	ITV_Vero	50Gy 80%	64,3	98	66,3	50	66,8	59,1	64,6	12,03	555	-2,8	3,4	3,4	4,6	0	08/09/2012	
				9	1	ROK	1	0	tracking	50Gy 80%	94,5	98,16	69	50,1	70,2	61,5	66,9	2,42	555	1,3	5,4	15,3	16,3	0	08/09/2012
				2	2	10	1	0	ITV_Vero	50Gy 80%	0,22	85,8	61,4	54,9	61,4	55,4	59,31	1,78	555	3,7	3,8	11,6	12,0	1	11/08/2013
				11	2	in vet voor lever	3	0	ITV_Vero	50Gy 80%	50,04	99,96	63,41	52,2	63,88	61,21	62,42	2,27	888	5,4	2,1	7,7	9,6	0	11/08/2013
BS708280S00A	1	1	12	1	lever S2	2	0	tracking	50Gy 80%	89,4	98,6	64,8	50,9	65,1	61,8	63,9	3,44	555	2,6	2,0	3,1	4,5	no 4D: weird 3set	0	04/10/2012
A450509YED01D	1	1	13	1	lever SB	2	0	tracking	50Gy 80%	78,4	100	66,02	52,4	66,3	60,6	63,9	17,02	555	3,1	3,6	14,2	15,0	0	28/07/2012	
B371128YED01	1	1	14	1	lever lobulus caudatus	2	0	ITV_Vero	50Gy 80%	80,87	100	65,81	52,91	66,35	61,24	64,44	4,91	555	5,8	5,1	8,1	11,2	0	27/11/2012	
BS90310TS00S	1	2	15	1	ROK	1	0	ITV_Vero	50Gy 80%	82,06	100	64,75	52,4	64,78	62,5	64,51	0,13	555	1,5	1,8	4,2	4,8	0	22/11/2012	
				16	2	LBK	1	0	ITV_Vero	50Gy 80%	92,72	100	64,54	52,3	65,47	62,94	64,29	0,2	555	1,0	0,7	3,9	2,3	0	22/11/2012
				2	1	3	1	0	ITV_Vero	50Gy 80%	90,33	100	67,77	52,57	66,28	61,71	65,41	1,9	555	0,8	1,0	7,0	7,1	0	11/03/2013
A610409VK	1	2	18	1	LBK	1	0	ITV_Vero	50Gy 80%	72,42	99,96	64,42	52,7	65,53	62,62	64,04	0,34	555	12,0	8,7	14,1	20,5	0	05/12/2012	
				19	2	klier retroper	4	1	ITV_Vero	40Gy 100%	0	0	40,86	36,03	40,63	38,04	39,99	5,06	777	0,0	0,0	0,0	0,0	1	06/12/2012
AS80419NM	1	4	20	1	ROK	1	0	ITV_Tomo	50Gy 80%	80,22	99,01	65,08	50,44	65,68	62,48	64,48	1,85	888	3,0	3,3	6,9	8,2	0	28/02/2013	
				21	2	RMK	1	0	ITV_Tomo	50Gy 80%	55,14	97,88	64,42	49,97	55,87	48,67	64,12	1,69	888	4,8	3,5	21,0	21,8	0	28/02/2013
				22	3	LBK	1	0	ITV_Tomo	50Gy 80%	70,68	98,44	64,42	50,19	65,27	62,44	64,48	0,466	888	1,1	3,2	7,0	7,8	0	28/02/2013
				23	4	LBK (long)	1	0	ITV_Tomo	50Gy 80%	64,08	98,33	64,57	50,34	65,65	61,46	64,54	0,572	888	1,0	2,5	10,5	10,8	0	28/02/2013
a451209PO	1	2	25	1	LBK	1	0	Tracking	50Gy 80%	84,78	99,84	66,28	51,87	66,69	61,13	64,45	2,53	555	3,8	4,8	4,7	7,7	0	07/03/2013	

- Tumor characteristics
- 4DCT motion amplitude
- Treatment data
- Dose parameters
- Date of FU and local recurrence
- Patient characteristics
- Imaging data
- Biological data

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

1. GTV volume (cc)
2. Tumor location
3. Motion Management (ITV vs. tracking)
4. 4DCT motion amplitude (LR, AP, CC)
5. Six dose parameters (3 PTV, 3 GTV)
6. Start of treatment date
7. Dates of FU
8. Presence of local recurrence [0,1]

Input vs. Output
Features vs. Labels

Features
Categorical or continuous

Size ??
Total: 109, 28% LR [1]

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

```
# To read the data from an Excel file and illustrate the first few rows;  
# Change the location to the Excel file location on your PC!  
data = pd.read_excel(r'/Users/jd/Documents/Work/ESTRO/ESTRO 2019/pre-meeting/ESTRO_premeeting_dataset.xlsx',target='LR.  
data.head()
```

	Location_exact	Technique	PTV_V50Gy	PTV_D2%	PTV_D98%	GTV_D2%	GTV_D98%	GTV_D50%	GTV_cc	LR	AP	CC	EUCL	LR.1
0	LBK	Tracking	100.0	64.6	51.3	65.9	57.9	62.5	1.06	1.0	3.0	5.7	6.518435	0
1	RBK	ITV	99.7	65.8	51.6	67.1	60.5	63.5	0.76	-2.0	7.0	6.0	9.433981	0
2	RMK	ITV	99.9	65.9	51.9	67.0	60.3	64.2	2.20	-2.1	1.6	3.3	4.226109	1
3	ROK	ITV	99.4	66.3	51.0	67.0	60.8	64.3	3.70	-0.7	-2.6	7.9	8.346257	0
4	RBK	ITV	100.0	65.1	52.0	65.3	57.5	61.7	1.30	0.5	-1.1	1.4	1.849324	0

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

```
# To show some info on the dataset
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 109 entries, 0 to 108
Data columns (total 14 columns):
```

Location_exact	109 non-null object
Technique	109 non-null object
PTV_V50Gy	109 non-null float64
PTV_D2%	109 non-null float64
PTV_D98%	109 non-null float64
GTV_D2%	109 non-null float64
GTV_D98%	109 non-null float64
GTV_D50%	109 non-null float64
GTV_cc	109 non-null float64
LR	109 non-null float64
AP	109 non-null float64
CC	109 non-null float64
EUCL	109 non-null float64

Categorical features

Numerical features

LR.1	109 non-null int64
------	--------------------

dtypes: float64(11), int64(1), object(2)
memory usage: 12.0+ KB

Label

An end-to-end machine learning project

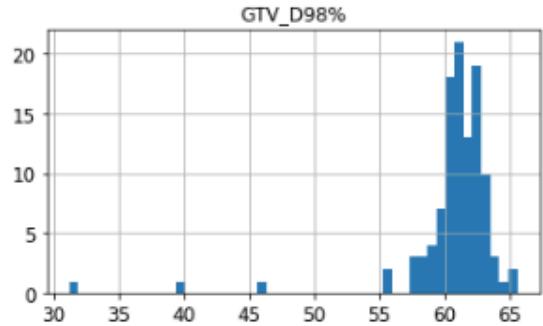
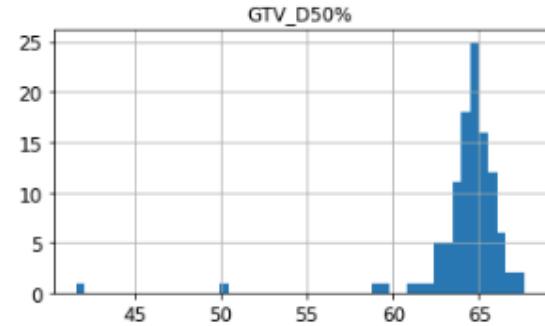
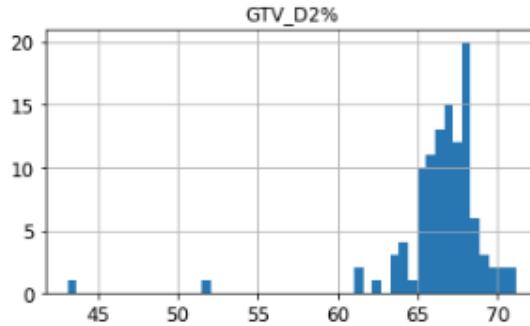
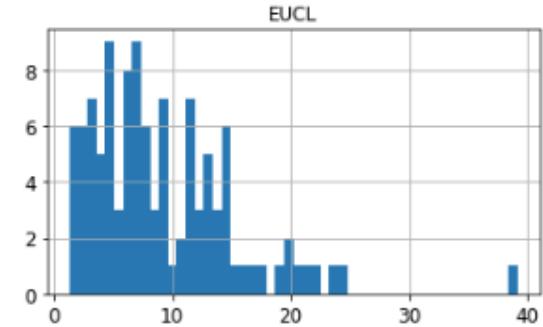
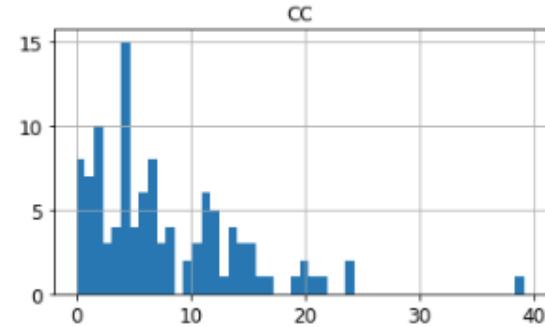
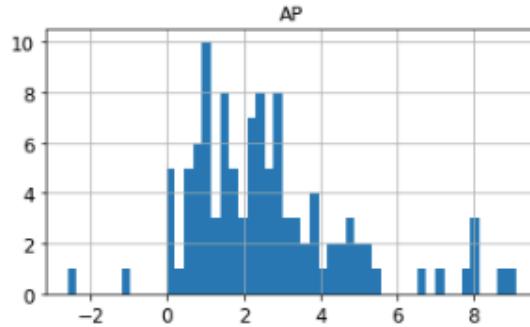
Step 3: Get the data... and have a first glimpse

```
# To give some statistics on the numerical features
data.describe()
```

	PTV_V50Gy	PTV_D2%	PTV_D98%	GTV_D2%	GTV_D98%	GTV_D50%	GTV_cc	LR	AP	CC	EUCL	LR.1
count	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000	109.000000
mean	98.774587	65.781284	51.757064	66.475321	60.640826	64.188624	4.216156	1.781284	2.610000	7.665688	9.047311	0.293578
std	4.420853	1.933504	1.934911	3.192854	4.173766	2.921886	7.219889	1.976544	2.098946	6.522449	6.134108	0.457504
min	65.760000	54.570000	43.110000	43.120000	31.210000	41.590000	0.100000	-2.800000	-2.600000	0.000000	1.277576	0.000000
25%	99.390000	64.800000	51.130000	65.700000	60.230000	63.850000	0.600000	0.800000	1.100000	3.100000	4.736032	0.000000
50%	99.870000	66.140000	51.820000	67.000000	61.210000	64.550000	1.300000	1.300000	2.300000	5.900000	7.502666	0.000000
75%	99.990000	66.980000	52.590000	67.980000	62.450000	65.240000	3.700000	2.600000	3.400000	11.600000	12.287392	1.000000
max	100.000000	69.000000	56.180000	71.110000	65.580000	67.600000	31.590000	12.000000	9.100000	39.100000	39.141538	1.000000

An end-to-end machine learning project

```
# To plot the histograms of all numerical features
%matplotlib inline
data.hist(bins=50, figsize=(20,15))
plt.show()
```



An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

```
patient_data["Location_exact"].value_counts()
```

LBK	24
RBK	17
LOK	15
ROK	14
RMK	5
LOK ant	2
LBK medial	2
Li apex	1
LBK (laag)	1
ROK, craniaal	1
post re pleura	1
Re	1
Li	1
RMK perifeer	1
LBK sup	1
LOK	1
LBK lateral	1
ROK cranial	1
LBK caudal	1
li hilus	1
ROK, anterieur	1
Lingula	1
LOK top (cave: LBK was al gereisceerd)	1
LBK inf	1

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

Preliminary pre-processing

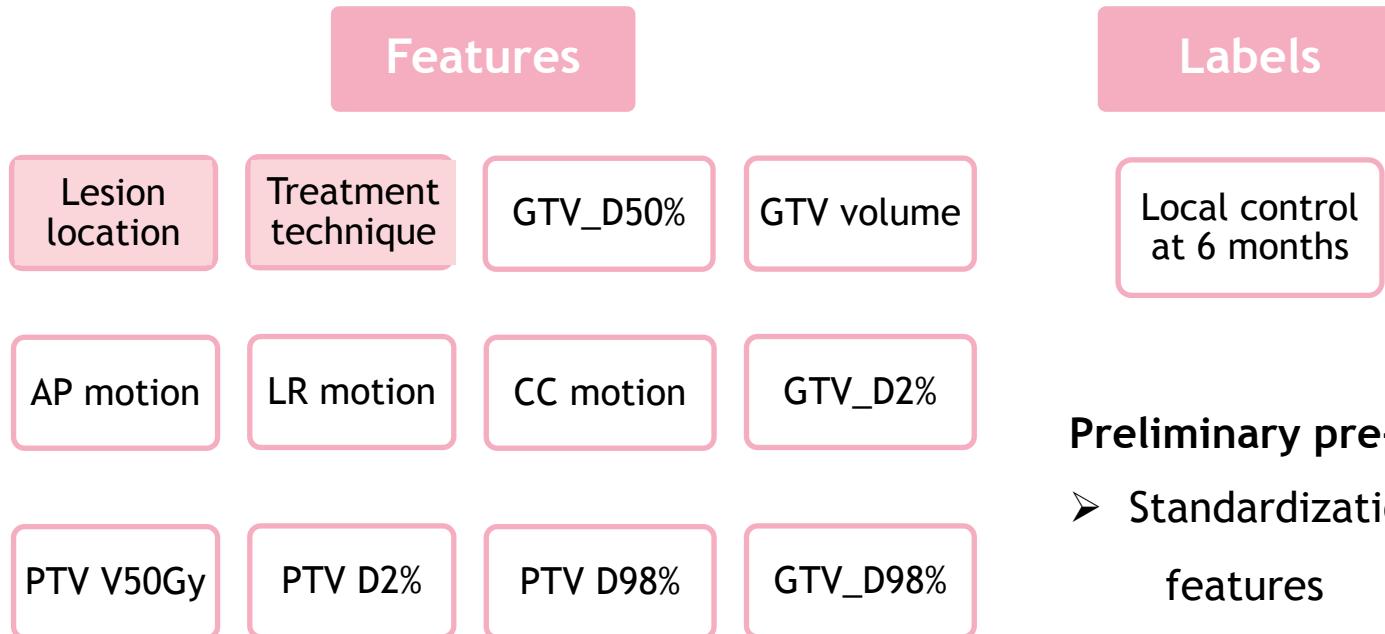
Post re pleura												Lesion location													
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V				
	Location_label	Technique	PTV_mm	PTV_xy	PTV_z	PTV_xy2	PTV_z2	ETV_xy	ETV_z	LR	AP	CC	L	start_RT	last_PV	date_update	PTV_xy3	PTV_z3	LR_V	AP_V	CC_V	LC_data	LC_meeting		
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24		
RBK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK	ROK		
RMK	Li	Re	ROK, caudaal	ROK, caudal	ant ROK	LBK medial	LBK ant	ROK (laag)	Li	Re	ROK, anterieur	ROK, cranial	ROK, craniaal	ROK, cranial	ROK, ant	ROK ant	RE LONG	RMK perifeer	ROK top (...)	LBK inf	LBK sup	LBK caudal	RMK pleura	Lingula	
ROK	LBK	LOK	ROK, caudaal	ROK caudal	ant ROK	LBK medial	LBK ant	ROK, craniaal	ROK cranial	ROK, anterieur	ROK, cranial	ROK, cranial	ROK, cranial	ROK, cranial	ROK, ant	RE LONG	LONG	RMK perifeer	ROK top (...)	LBK inf	LBK sup	LBK caudal	RMK pleura	Li apex	LOK ant
LBK	LOK	RE	ROK, caudal	ROK caudal	ant ROK	LBK medial	LBK ant	ROK (laag)	Li	Re	ROK, anterieur	ROK, cranial	ROK, craniaal	ROK, cranial	ROK, ant	RE LONG	LONG	RMK perifeer	ROK top (...)	LBK inf	LBK sup	LBK caudal	RMK pleura	Lingula	LOK ant
LOK	RE	ROK	ROK, caudaal	ROK caudal	ant ROK	LBK medial	LBK ant	ROK, craniaal	ROK cranial	ROK, anterieur	ROK, cranial	ROK, cranial	ROK, cranial	ROK, cranial	ROK, ant	RE LONG	LONG	RMK perifeer	ROK top (...)	LBK inf	LBK sup	LBK caudal	RMK pleura	Lingula	LOK ant



1 dr. dataset

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

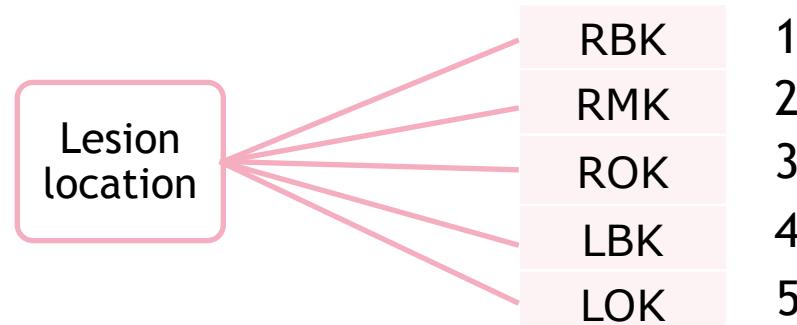


Preliminary pre-processing

- Standardization of continuous features
- Change categorical to binary

An end-to-end machine learning project

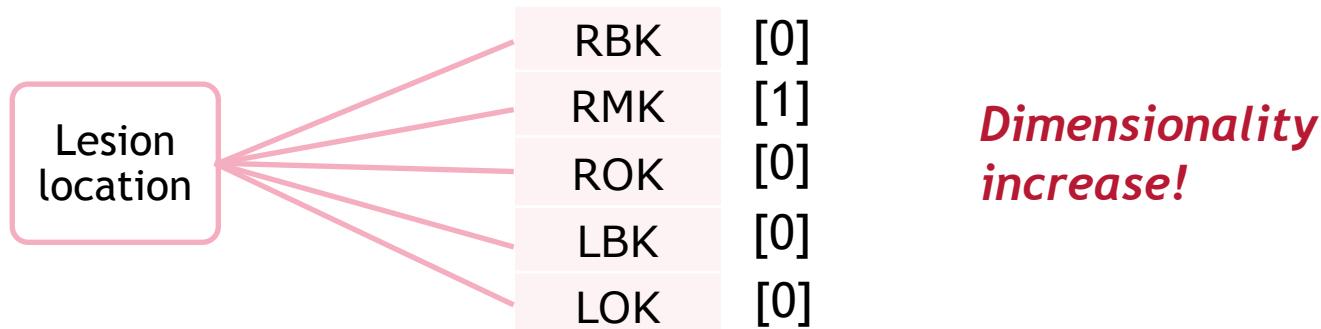
Step 3: Get the data... and have a first glimpse



Nominal vs. Ordinal features

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse



An end-to-end machine learning project

```
# To drop irrelevant features and change categorical features to binary features
from sklearn.preprocessing import OneHotEncoder
from sklearn.preprocessing import StandardScaler

onehot = OneHotEncoder(dtype=np.int, sparse=True)

nominals = pd.DataFrame(
    onehot.fit_transform(data[['Location_exact', 'Technique']])\
    .toarray(),
    columns=['LBK', 'LOK', 'RBK', 'RMK', 'ROK', 'ITV', 'tracking'])

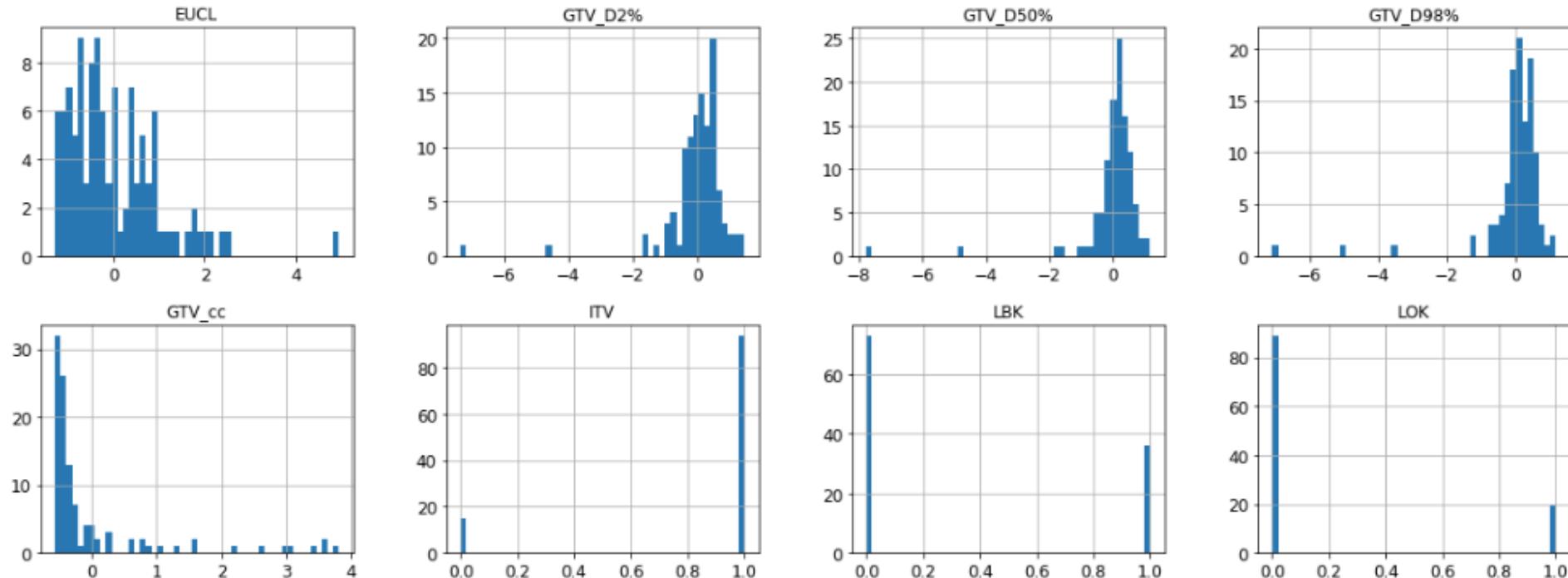
data_encoded = data.drop(["Location_exact", "Technique", "AP", "LR", "CC", "PTV_V50Gy"], axis=1)
data_encoded = nominals.join(data_encoded)

# To standardize all continuous features
scaler = StandardScaler()
data_encoded["PTV_D98%"] = scaler.fit_transform(data_encoded["PTV_D98%"].values.reshape(-1,1))
data_encoded["PTV_D2%"] = scaler.fit_transform(data_encoded["PTV_D2%"].values.reshape(-1,1))
data_encoded["GTV_D2%"] = scaler.fit_transform(data_encoded["GTV_D2%"].values.reshape(-1,1))
data_encoded["GTV_D98%"] = scaler.fit_transform(data_encoded["GTV_D98%"].values.reshape(-1,1))
data_encoded["GTV_D50%"] = scaler.fit_transform(data_encoded["GTV_D50%"].values.reshape(-1,1))
data_encoded["GTV_cc"] = scaler.fit_transform(data_encoded["GTV_cc"].values.reshape(-1,1))
data_encoded["EUCL"] = scaler.fit_transform(data_encoded["EUCL"].values.reshape(-1,1))

%matplotlib inline
data_encoded.hist(bins=50, figsize=(20,15))
plt.show()
```

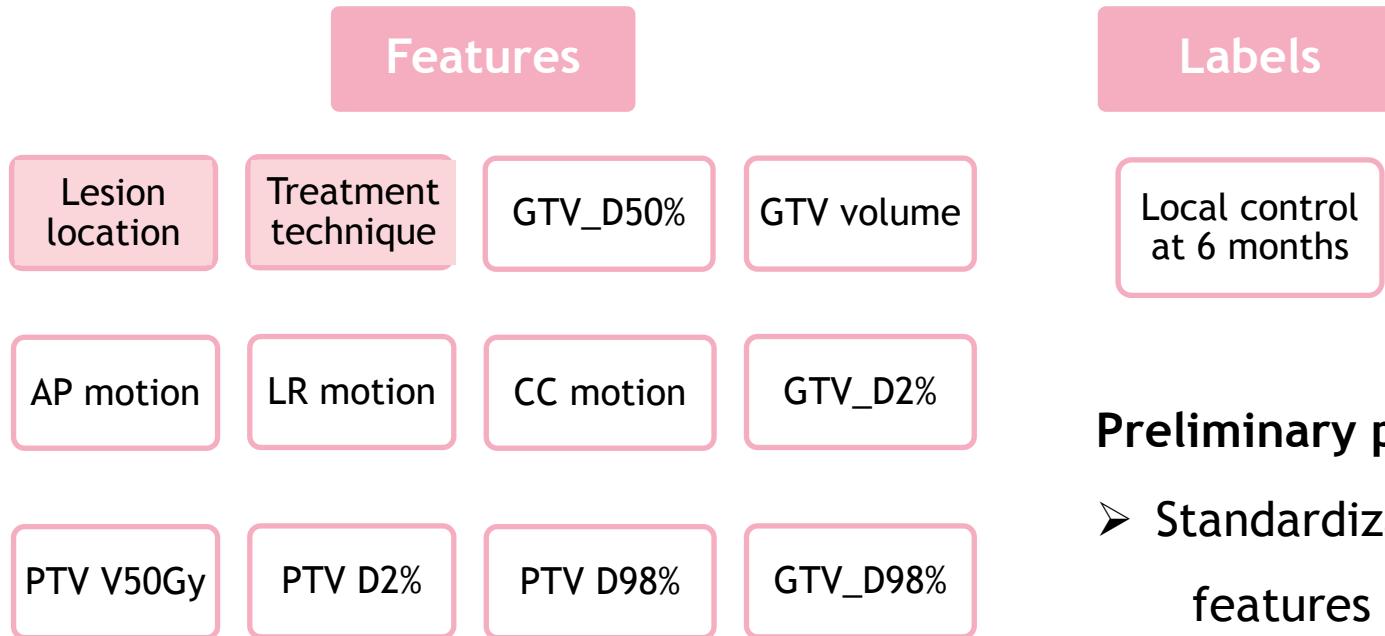
An end-to-end machine learning project

Step 3: Explore the data to get insights



An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse



Preliminary pre-processing

- Standardization of continuous features
- Change categorical to binary

An end-to-end machine learning project

Step 3: Get the data... and have a first glimpse

Features				Labels
Tracking	ITV	GTV volume	LBK	
LR motion	AP motion	CC motion	LOK	Local control at 6 months
GTV_D2%	GTV_D98%	GTV_D50%	ROK	
PTV D98%	PTV D2%	PTV V50Gy	RMK	RBK

An end-to-end machine learning project

Some take-home messages on data

Your application will only ever be as good as your data.

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016

BUSINESS NEWS OCTOBER 10, 2018 / 5:12 AM / 6 MONTHS AGO

Amazon scraps secret AI recruiting tool that showed bias against women

With machine learning, you assume nothing a priori, you learn from the data.



Data from old treatment techniques will not be accurate now.

An end-to-end machine learning project

Some take-home messages on data

A matter of quality and quantity

The unreasonable effectiveness of data¹.

“Significant benefit of training data available in the wild”

Small and medium-sized datasets are still common...

So don't abandon algorithms just yet...

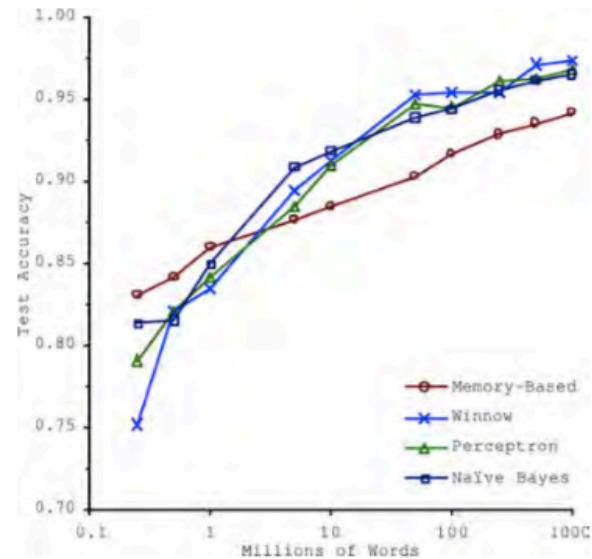


Figure 1-20. The importance of data versus algorithms⁹

[1] “The Unreasonable Effectiveness of Data,” Peter Norvig et al. (2009).

An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
 2. Feature engineering
 3. Get the data - data pre-processing
 4. Creating a training, (validation) and test set
 5. Feature selection
 6. Model selection
 7. Fitting the model
 8. Model evaluation
 9. QA - implement
 10. share!
-
-

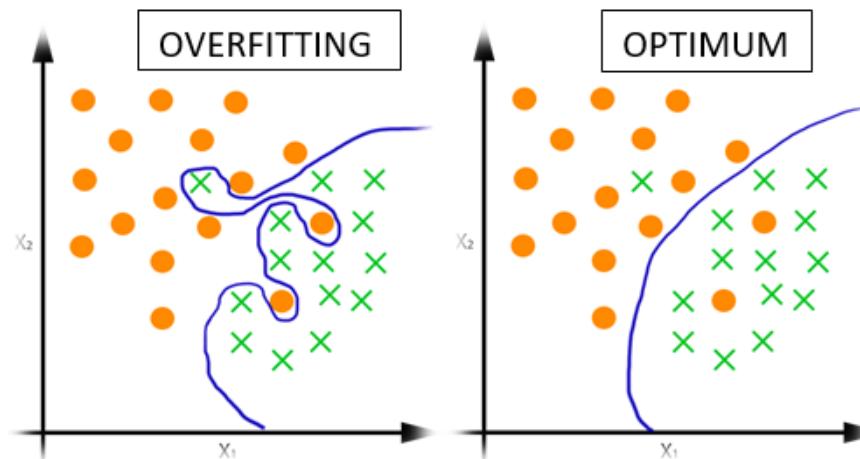
An end-to-end machine learning project

Step 4: Creation of a training, (validation) and test set



GOAL in ML, DL, ...

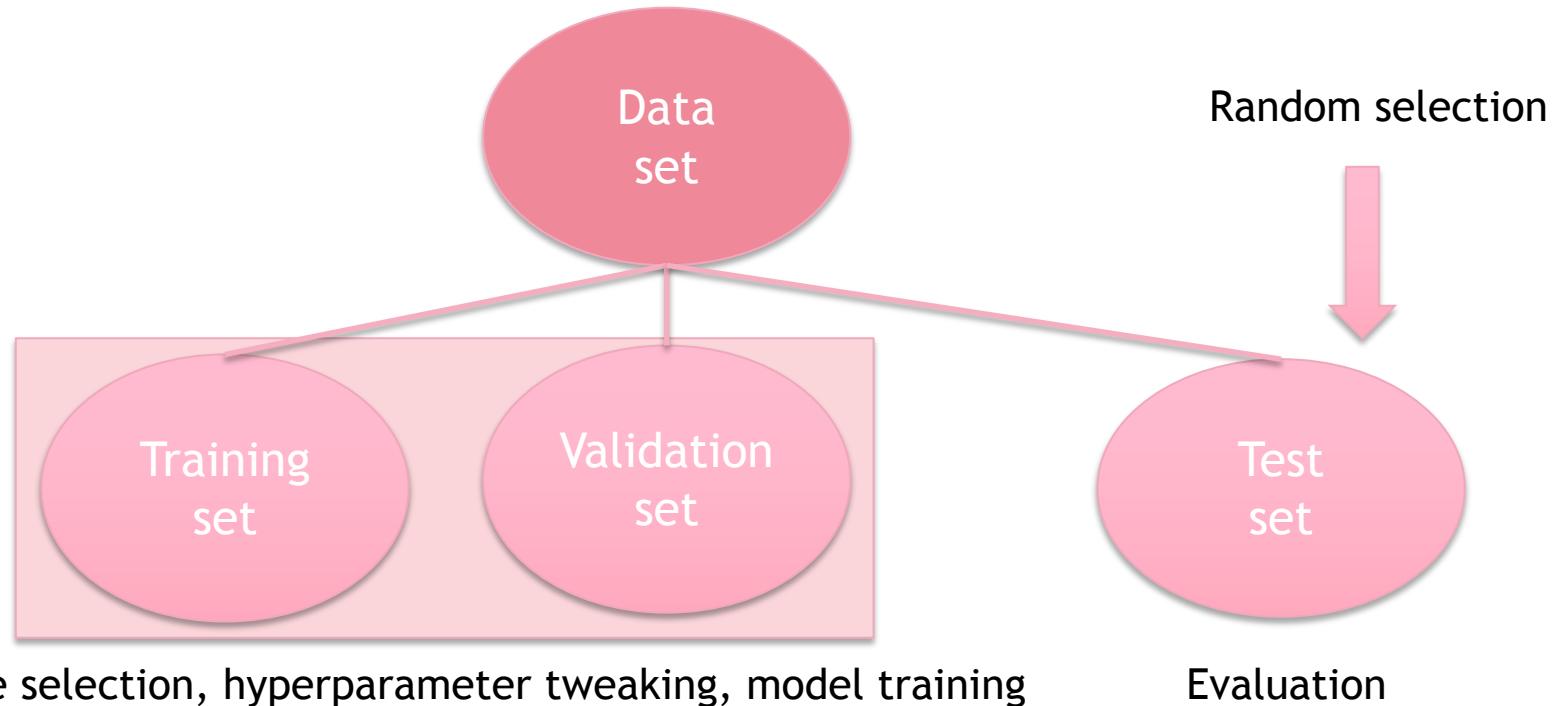
Model generalization vs. Overfitting



Source: Overfitting and human behaviour, Medium (2018)

An end-to-end machine learning project

Step 4: Creation of a training, (validation) and test set



An end-to-end machine learning project

Step 4: Creation of a training, (validation) and test set

In case of smaller datasets; n-fold cross-validation



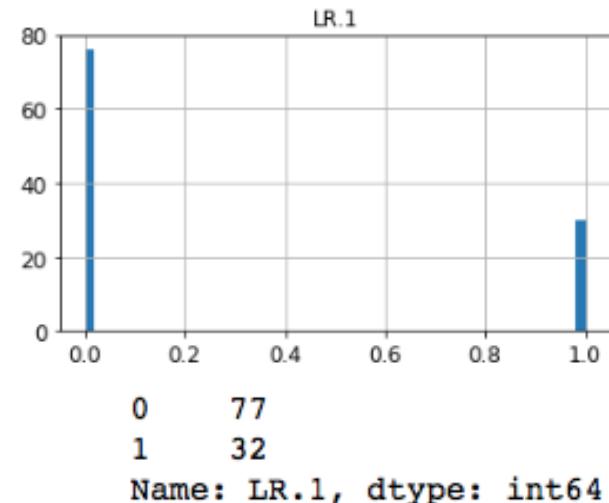
An end-to-end machine learning project

Step 4: Creation of a training, (validation) and test set

“Random sampler” - 80% training set

Large datasets: Train & test sets will represent population distributions

Data imbalance



Sensitivity - % of patients with local recurrence identified as such

Specificity - % of patients without local recurrence identified as such

An end-to-end machine learning project

Step 4: Creation of a training, (validation) and test set

- Undersampling

65 - 15 → 15 - 15

- Synthetic minority oversampling

Small datasets: Imbalanced AND not representative

```
# To show number of "Technique" categories in the dataset
data["Technique"].value_counts()
```

```
ITV        94
Tracking   15
Name: Technique, dtype: int64
```

imbalanced-learn.org

+ stratified sampling (= population representation)

An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
2. Feature engineering
3. Get the data - data pre-processing
4. Creating a training, (validation) and test set
5. Feature selection
6. Model selection
7. Fitting the model
8. Model evaluation
9. QA - implement
10. share!

An end-to-end machine learning project

Step 5: Feature selection

Feature formulation

- Tumor biology
- Tumor microenvironment
- Radiation dosimetry
- Patient characteristics
- Imaging features...

Feature extraction

- Patient records
- Pathology
- Planning software
- Radiology...

Feature selection

- Relevant feature selection
- The curse of dimensionality (data sparsity)
- Reduces chance of overfitting
- Improved model accuracy,
- Decreases training time /complexity
- Better understanding

→ Storing raw data

→ Consistency

An end-to-end machine learning project

Step 5: Feature selection

Curse of dimensionality

Multiple representations of every option

Ex. 2 binary features (gender, cancer diagnosis)

4 combinations, 40 training samples

Ex. 5 binary features (gender, cancer diagnosis,)

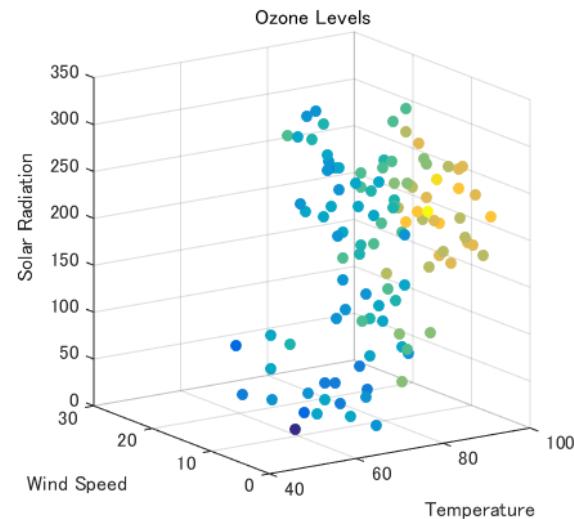
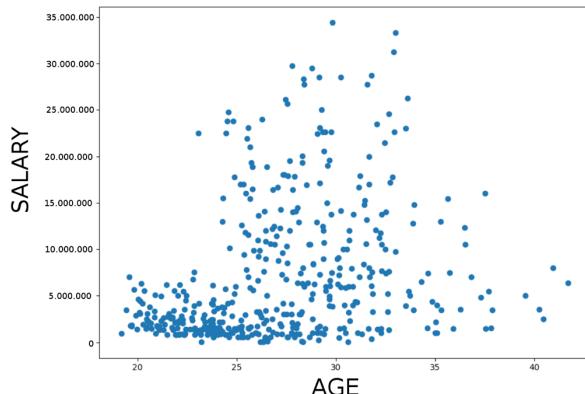
32 combinations, 320 training samples

An end-to-end machine learning project

Step 5: Feature selection

Curse of dimensionality

Multiple representations of every option



An end-to-end machine learning project

Step 5: Feature selection

Feature formulation

- Tumor biology
- Tumor microenvironment
- Radiation dosimetry
- Patient characteristics
- Imaging features...

Feature extraction

- Patient records
- Pathology
- Planning software
- Radiology...

Feature selection

- Relevant feature selection
- The curse of dimensionality
(data sparsity)
- Reduces chance of overfitting
- Improved model accuracy
- Decreases training time
/complexity
- Better understanding

→ Consistency

→ Storing raw data

An end-to-end machine learning project

Step 5: Feature selection **How?**

A matter of balance, enough but not too much

*Evaluate all possible feature combinations...
but nobody has time for that, so a suboptimal methodology has to be applied.*

1. Start with all of them.
2. *Look for redundant features*
3. *Dimensionality reduction*
4. *Automated selection methods*

An end-to-end machine learning project

Step 5: Feature selection

How?

Features

Tracking

ITV

GTV volume

LBK

LR motion

AP motion

CC motion

LOK

GTV_D2%

GTV_D98%

GTV_D50%

ROK

PTV D98%

PTV D2%

PTV V50Gy

RMK

RBK

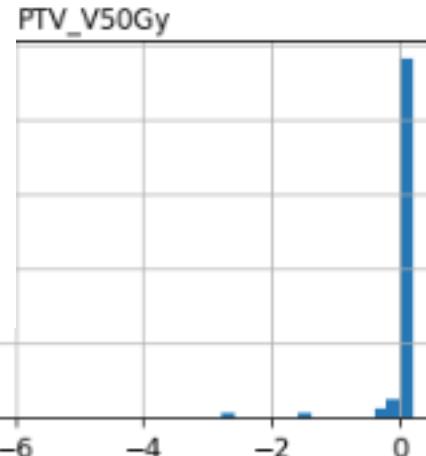
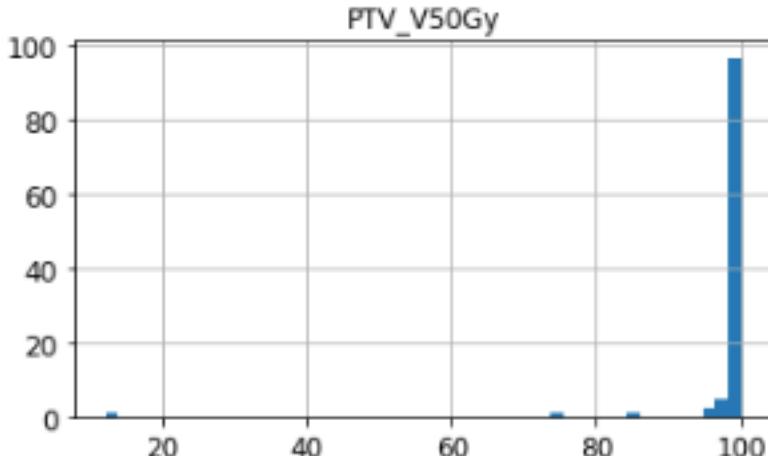
1. Start with all of them.
2. *Look for redundant features*
3. *Dimensionality reduction*
4. *Automated selection methods*

An end-to-end machine learning project

Step 5: Feature selection

How?

PTV V50Gy



1. Start

2. Look

3. Dimer

4. Auton

PTV_V50Gy

count 106.000000

mean 98.352358

std 8.915657

min 12.150000

25% 99.382500

50% 99.860000

75% 99.980000

max 100.000000

An end-to-end machine learning project

Step 5: Feature selection

How?

Features

Tracking

ITV

GTV volume

LBK

LR motion

AP motion

CC motion

LOK

GTV_D2%

GTV_D98%

GTV_D50%

ROK

PTV D98%

PTV D2%

PTV V50Gy

RMK

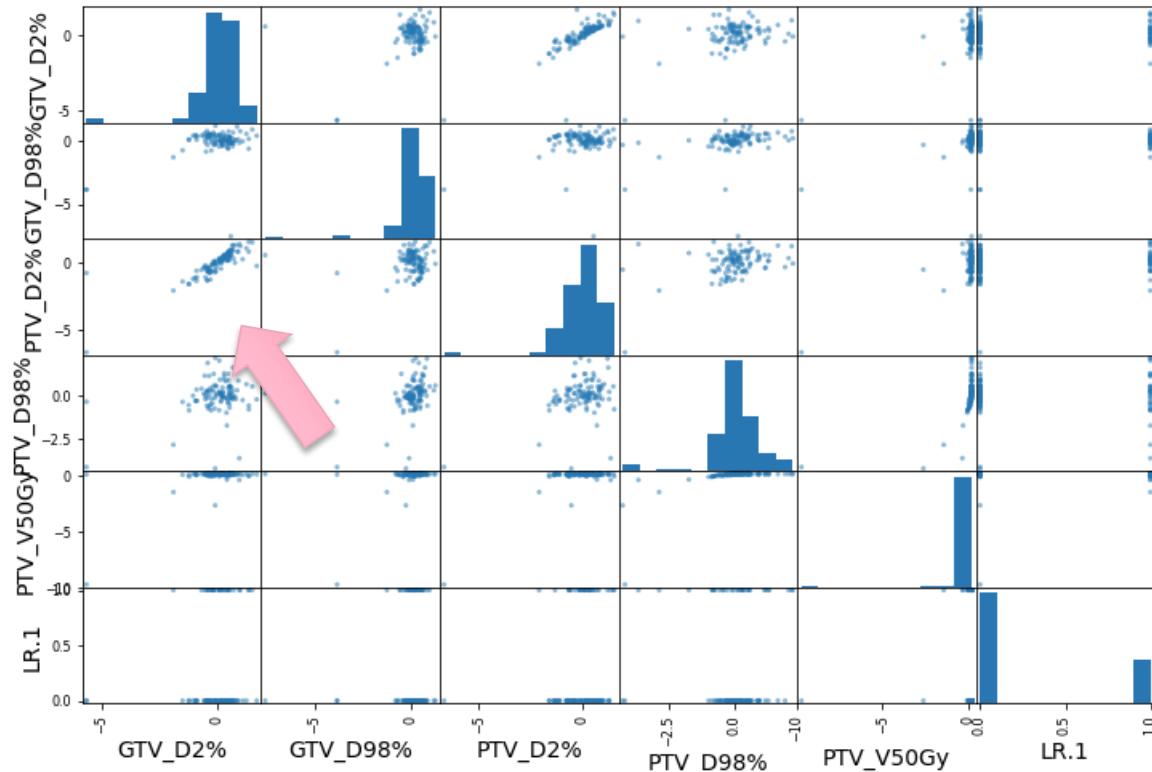
RBK

1. Start with all of them.
2. *Look for redundant features*
3. *Dimensionality reduction*
4. *Automated selection methods*

An end-to-end machine learning project

Step 5: Feature selection

How?



of them.
undant features
ty reduction
election methods

An end-to-end machine learning project

Step 5: Feature selection

How?

Features

Tracking

ITV

GTV volume

LBK

LR motion

AP motion

CC motion

LOK

GTV_D2%

GTV_D98%

GTV_D50%

ROK

PTV D98%

PTV D2%

RMK

RBK

1. Start with all of them.
2. *Look for redundant features*
3. *Dimensionality reduction*
4. *Automated selection methods*

An end-to-end machine learning project

Step 5: Feature selection

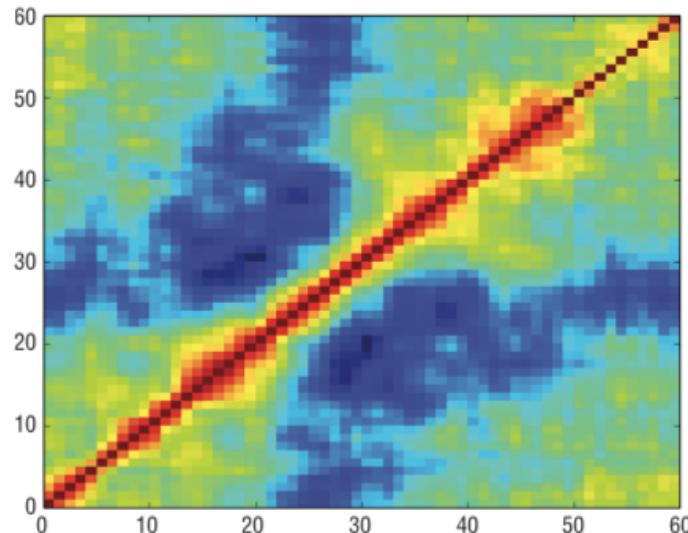


Figure 2-8: Heat map showing attribute cross-correlations

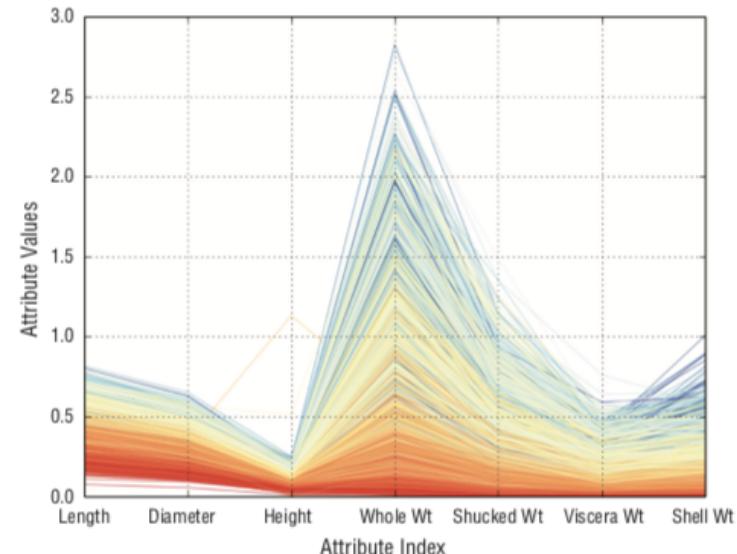


Figure 2-12: Color-coded parallel coordinate plot for abalone

An end-to-end machine learning project

Step 5: Feature selection

How?

Features

Tracking

ITV

GTV volume

LBK

LR motion

AP motion

CC motion

LOK

GTV_D2%

GTV_D98%

GTV_D50%

ROK

PTV D98%

RMK

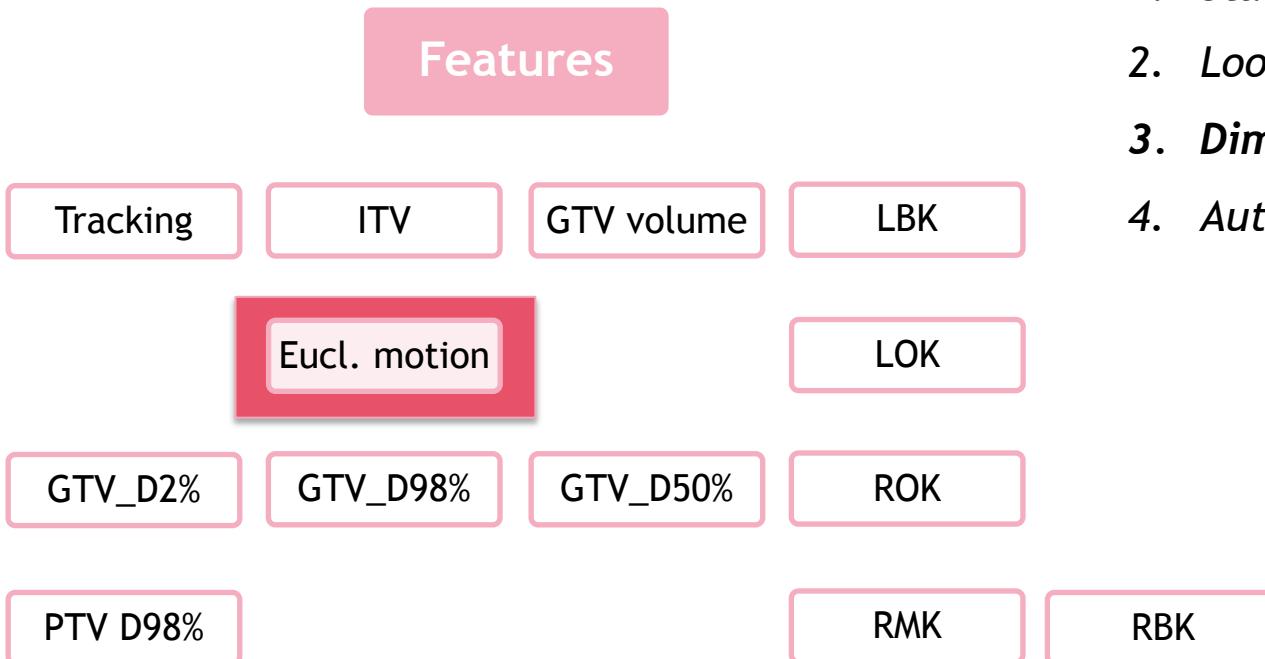
RBK

1. Start with all of them.
2. Look for redundant features
3. Dimensionality reduction
4. Automated selection methods

An end-to-end machine learning project

Step 5: Feature selection

How?



1. Start with all of them.
2. Look for redundant features
3. Dimensionality reduction
4. Automated selection methods

An end-to-end machine learning project

Step 5: Feature selection

Filter

Features are ranked according to relevance.

Pearson correlation (linear), Mutual information (MI), ...

Only highest ranked features are selected.

- univariate

4. Automated selection/elimination methods

Wrapper

Search algorithm finds subset of features with highest predictor performance.

Objective function steered maximum performance with minimum of features.

Some: backtracking take into account correlation between features.

Embedded

Embedded in the training process of the model.

- Artificial neural networks

An end-to-end machine learning project

Step 5: Feature selection

Filter

Wrapper

Cross-validation recursive feature selection on the training set

- Takes into account collinearity
- Only works on linear models

```
# To perform automatic cross-validation recursive feature selection on the training set

import matplotlib.pyplot as plt
from sklearn.svm import SVC
from sklearn.model_selection import StratifiedKFold
from sklearn.feature_selection import RFECV
from sklearn.datasets import make_classification

# Create the RFE object and compute a cross-validated score.
svc = SVC(kernel="linear")
# The "accuracy" scoring is proportional to the number of correct
# classifications
rfecv = RFECV(estimator=svc, step=1, cv=StratifiedKFold(5),
              scoring='accuracy')
rfecv.fit(train_features, train_labels)

print("Optimal number of features : %d" % rfecv.n_features_)
print(rfecv.support_)

# Plot number of features VS. cross-validation scores
plt.figure()
plt.xlabel("Number of features selected")
plt.ylabel("Cross validation score (nb of correct classifications)")
plt.plot(range(1, len(rfecv.grid_scores_) + 1), rfecv.grid_scores_)
plt.show()
```

An end-to-end machine learning project

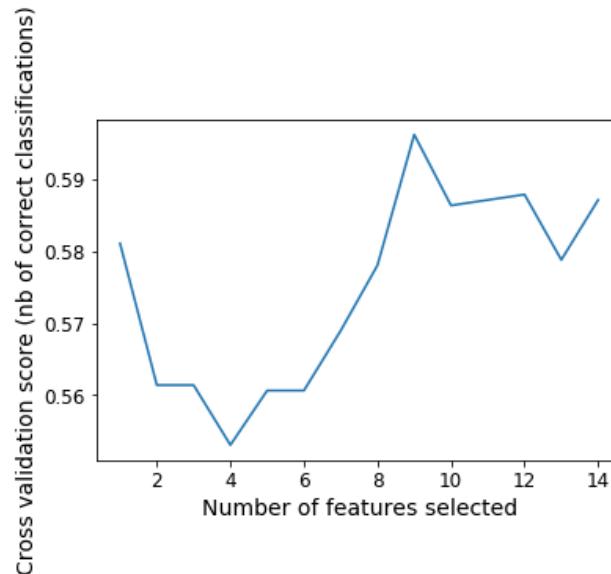
Step 5: Feature selection

Filter

Wrapper

Cross-validation recursive feature selection on the training set

```
Optimal number of features : 9
[ True  True  True  True  True False False False  True  True  True
  True False]
```



- Takes into account collinearity
- Only works on linear models

An end-to-end machine learning project

Artificial dataset

```
# To show some info on the dataset
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 450 entries, 0 to 449
Data columns (total 13 columns):
Location_exact      450 non-null object
Technique            450 non-null object
PTV_V50Gy           450 non-null float64
PTV_D2%              450 non-null float64
PTV_D98%             450 non-null float64
GTV_D2%              450 non-null float64
GTV_D98%             450 non-null float64
GTV_cc               450 non-null float64
LR                   450 non-null float64
AP                   450 non-null float64
CC                   450 non-null float64
EUCL                450 non-null float64
LR.1                 450 non-null int64
dtypes: float64(10), int64(1), object(2)
memory usage: 45.8+ KB
```

May or may not have imposed
some correlations



An end-to-end machine learning project

Step 5: Feature selection

Filter

Wrapper

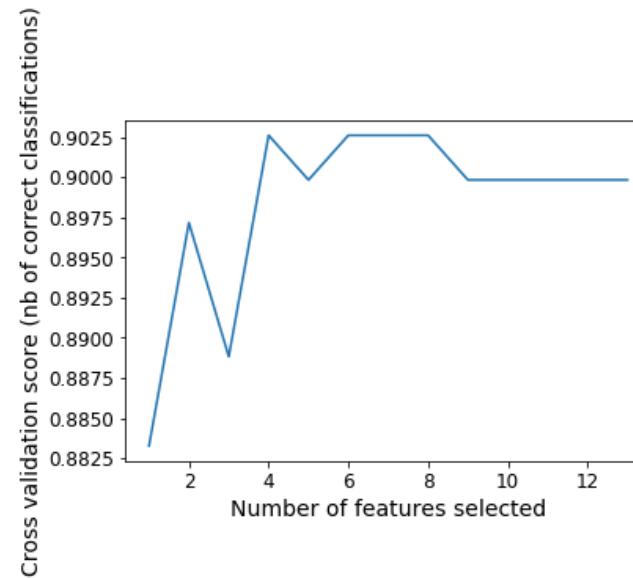
Euclidean motion

GTV volume

GTV_D2%

Cross-validation recursive feature selection on the training set

```
Optimal number of features : 4
[False False  True False False False False False  True False  True
 True]
```



An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
2. Feature engineering
3. Get the data - data pre-processing
4. Creating a training, (validation) and test set
5. Feature selection
6. Model selection
7. Fitting the model
8. Model evaluation
9. QA - implement
10. share!

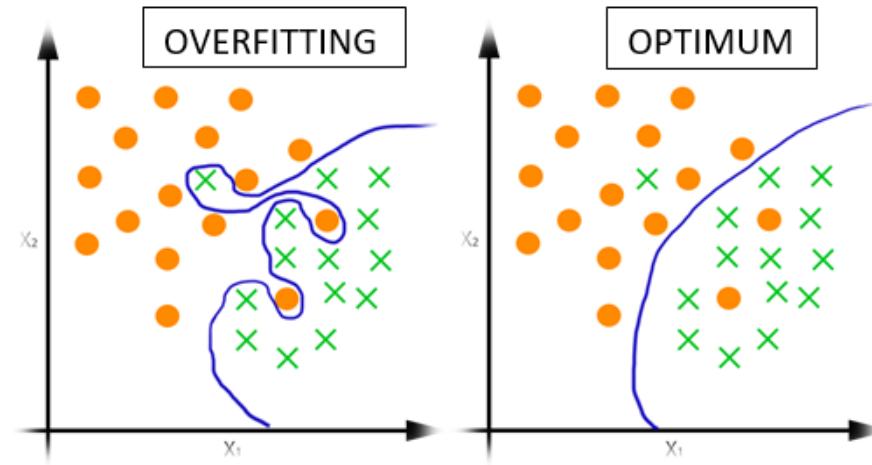
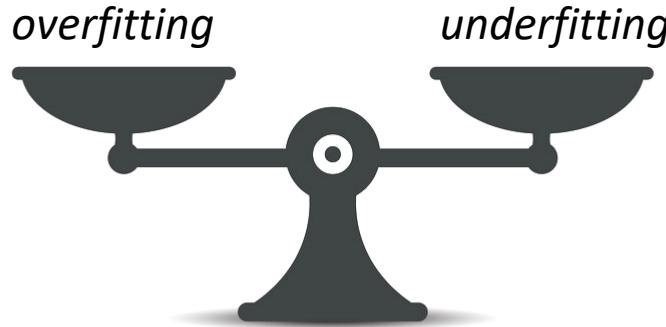
An end-to-end machine learning project

Step 6: Model selection

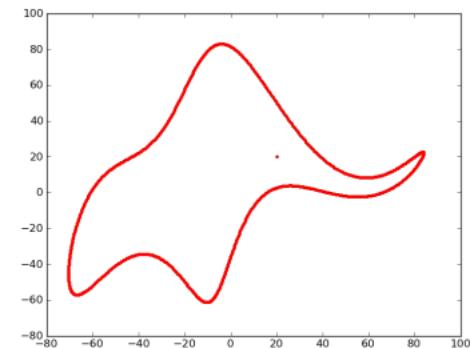
- Explore many models and short-list the best once.
 Based on problem framing.
- Have a peak at literature, you don't have to re-invent the wheel, and nobody builds a neural network from scratch! But be critical...
- Try short-list on training - validation set.

An end-to-end machine learning project

Step 6: Model selection



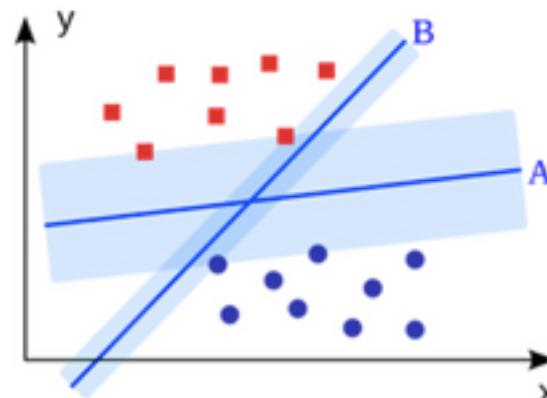
"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk." – John von Neumann



An end-to-end machine learning project

Step 6: Model selection

Support Vector Machine (SVM) - Support Vector Classifier (SVC)



Separates classes by finding plane between classes...
with maximum margin that ideally does not contain any instances

An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
2. Feature engineering
3. Get the data - data pre-processing
4. Creating a training, (validation) and test set
5. Feature selection
6. Model selection
7. Fitting the model
8. Model evaluation
9. QA - implement
10. share!



An end-to-end machine learning project

Step 7: Fitting the model

```
# To select and train the linear SVC (Support-Vector classification)

from sklearn.svm import SVC
clf = SVC(gamma='auto')
clf.fit(train_features_forfit, train_labels)

SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
2. Feature engineering
3. Get the data - data pre-processing
4. Creating a training, (validation) and test set
5. Feature selection
6. Model selection
7. Fitting the model
8. Model evaluation
9. QA - implement
10. share!

An end-to-end machine learning project

```
from sklearn.metrics import classification_report

predictions_train = clf.predict(train_features_forfit)
print(classification_report(train_labels, predictions_train))

predictions_test = clf.predict(test_features_forfit)
print(classification_report(test_labels, predictions_test))
```

	precision	recall	f1-score	support
0	0.90	0.98	0.94	258
1	0.93	0.73	0.81	102
micro avg	0.91	0.91	0.91	360
macro avg	0.91	0.85	0.87	360
weighted avg	0.91	0.91	0.90	360

Training set

	precision	recall	f1-score	support
0	0.88	0.95	0.92	64
1	0.86	0.69	0.77	26
micro avg	0.88	0.88	0.88	90
macro avg	0.87	0.82	0.84	90
weighted avg	0.88	0.88	0.87	90

Test set

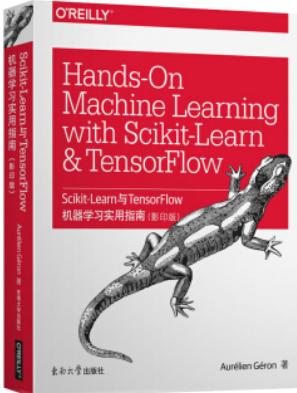
An end-to-end machine learning project

OUTLINE

1. Frame the problem, look at the big picture
2. Feature engineering
3. Get the data - data pre-processing
4. Creating a training, (validation) and test set
5. Feature selection
6. Model selection
7. Fitting the model
8. Model evaluation
9. QA - implement
10. share!

An end-to-end machine learning project

The end



github.com/JDhont/ESTRO38_premeeting

[Python.org](https://python.org)
[Jupyter.org](https://jupyter.org)
[Scikit-learn.org](https://scikit-learn.org)



YOUR PREPARATION FOR THE
REAL WORLD IS NOT IN THE
ANSWERS YOU'VE LEARNED,
BUT IN THE QUESTIONS YOU'VE
LEARNED TO ASK YOURSELF.
-BILL WATTERSON