

# Guest Lecture BAN400

05. October, 2021

**1** Introduction

**2** Why R?

**3** Project example: Valuation of Holiday Houses

**4** Code along and exercises

**5** Closing remarks

# 1

## Introduction



**Thomas Hansen**

Direktør

+47 952 60 254

[thomas.hansen@pwc.com](mailto:thomas.hansen@pwc.com)



**Line Melby**

Senior Associate

+47 948 21 449

[line.melby@pwc.com](mailto:line.melby@pwc.com)



**Nora Hansen**

Manager

+47 418 58 060

[nora.hansen@pwc.com](mailto:nora.hansen@pwc.com)



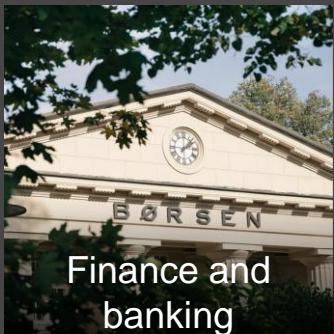
**Fredrik Angelvik**

Senior Associate

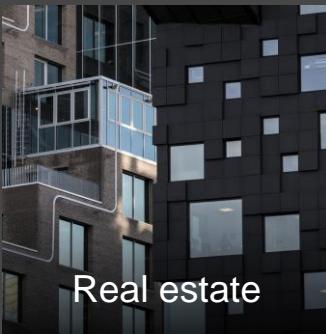
+47 480 48 017

[fredrik.angelvik@pwc.com](mailto:fredrik.angelvik@pwc.com)

# PwC are working with changing industries



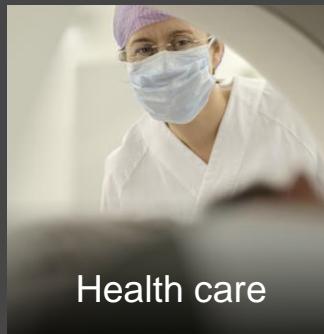
Finance and banking



Real estate



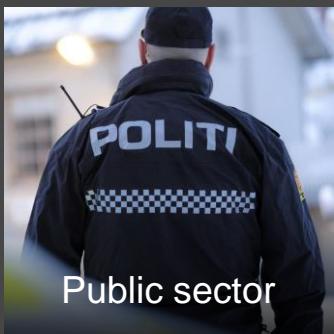
Energy



Health care



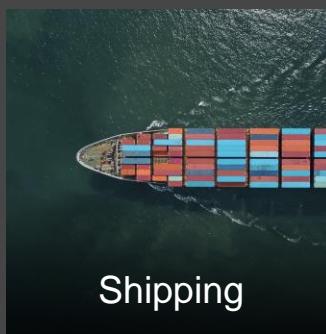
Retail



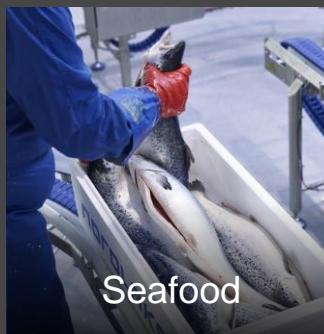
Public sector



Oil and gas



Shipping

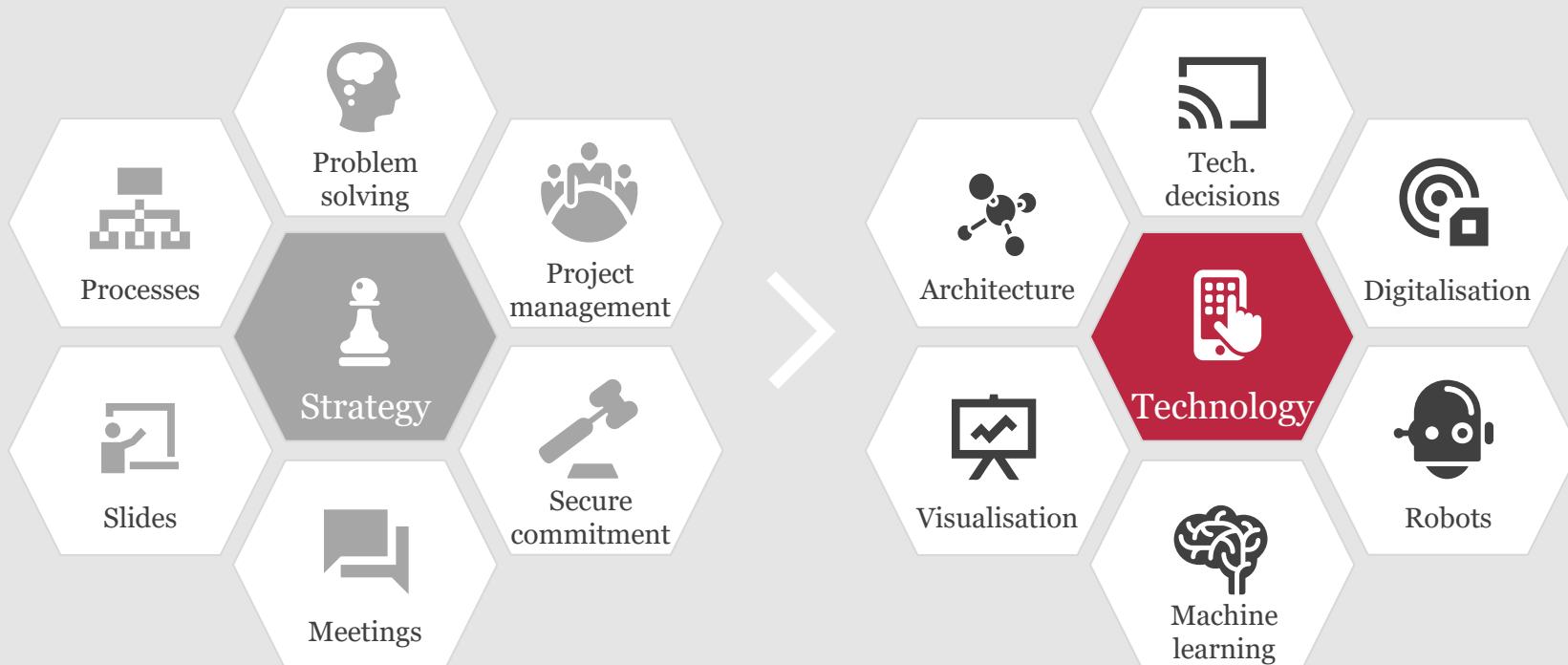


Seafood



Telecommunicatio  
n and media

Technological development creates needs and expectations. Tomorrow's consultants are **dependent on mastering the technological landscape**





Why R?

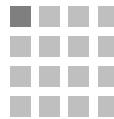


# Why we love R (We use Python as well)

# 1. A common syntax which is easy to write – but also understand



Database



Big data



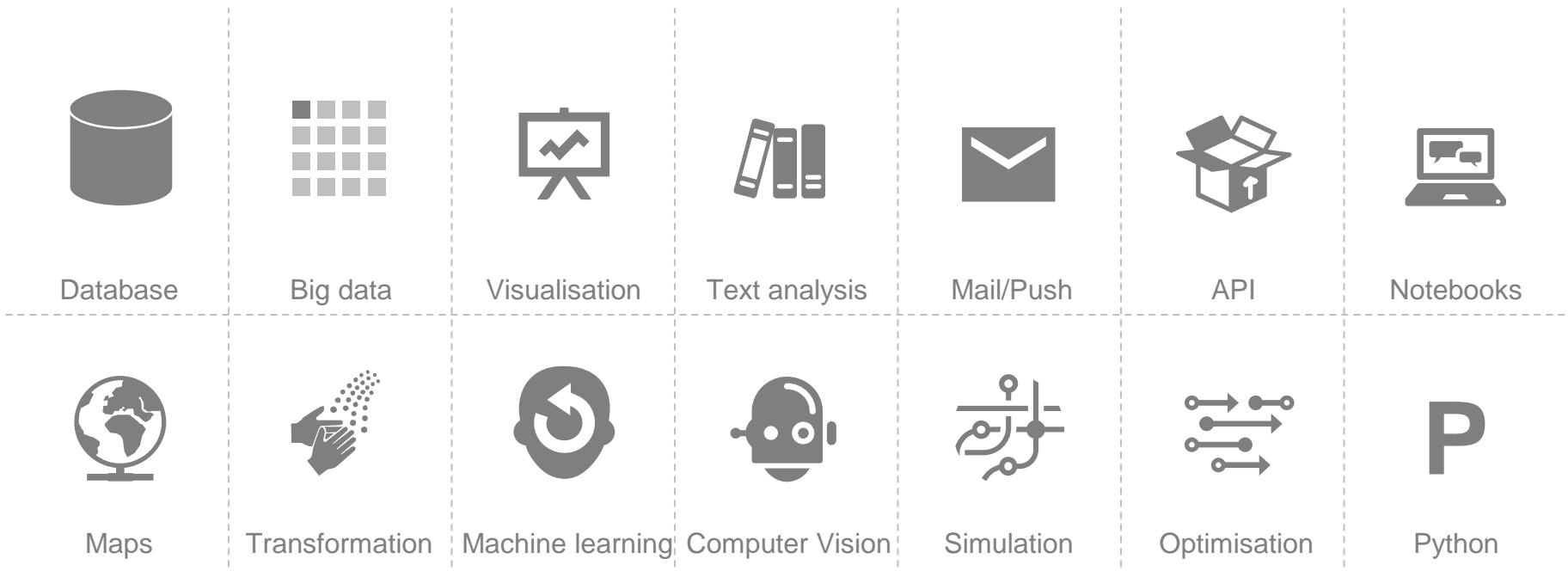
Map



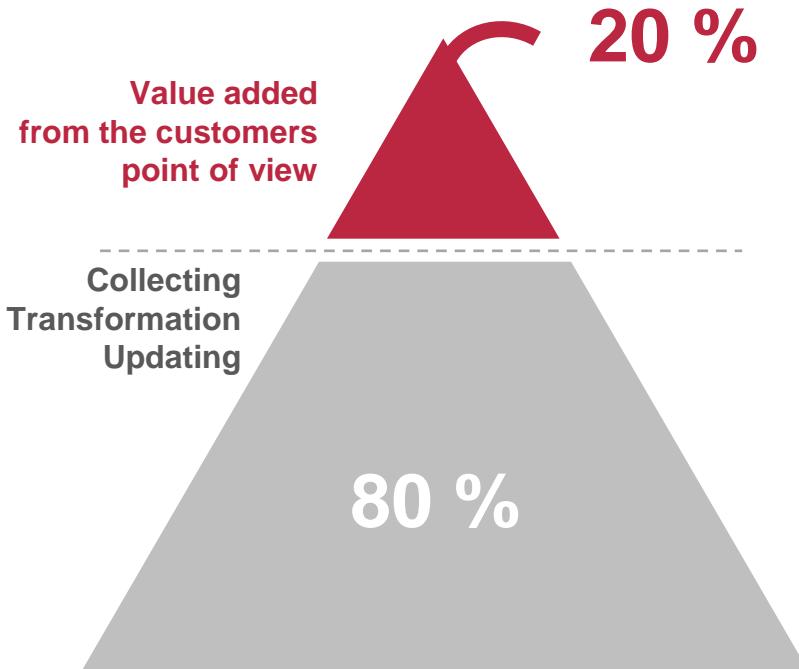
Calculate average square meter price for 3-room apartments over 60 square meters for each district in Oslo. Sort results by district with lowest square meter price.

```
data <- data_raw %>%
  filter(fylke == "Oslo", sqm > 60, bedrooms = 2) %>%
  group_by(bydel) %>%
  summarise(price_sqm = avg(price_sqm)) %>%
  arrange(price_sqm)
```

## 2. Functionality covering 99.9 % of the use cases for a Data Scientist with the same coherent syntax



### 3. It's the best tool for data wrangling and visualization (Opinionated)



↑ atomic\_explosion 20 points · 6 days ago

↓ I agree with [u/baconshoplifter](#), it's the libraries that matter not the language.

I use Python to clean and wrangle my data as libraries like Pandas and Numpy make it super easy especially with the ton of CSV files I work with. I use R to analyze and plot my data as I love ggplot and the ton of statistical packages in R.

I also use Python to build APIs, using Flask, and run them on AWS. So it really depends on what you are trying to do but I find it useful knowing both.

[Reply](#) [Share](#) [Report](#) [Save](#)

↑ foxhollow 47 points · 6 days ago

↓ I use python for many tasks and generally like it, but IMO, if you're comfortable in R, you're crazy to use pandas for data wrangling instead of dplyr. dplyr is so consistent and beautiful, while pandas is a design disaster.

[Reply](#) [Share](#) [Report](#) [Save](#)

↑ ACrispWinterDay 11 points · 6 days ago

↓ Yeah I was going to say the same thing. RStudio with dplyr and the tidyverse in general is heaven when it comes to data wrangling

[Reply](#) [Share](#) [Report](#) [Save](#)

↑ novokaoi 7 points · 6 days ago

↓ I had that exact thought the other day when I went back to using python for a project after using R for 2 years. How did I ever cope with the mess that is the Pandas API?

[Reply](#) [Share](#) [Report](#) [Save](#)



Project example:  
Valuation of Holiday Houses











Market Value

41,1  
MNOK

Kygo



Asset Value

10,3M  
NOK

Kygo

10,7M  
NOK

Ruth

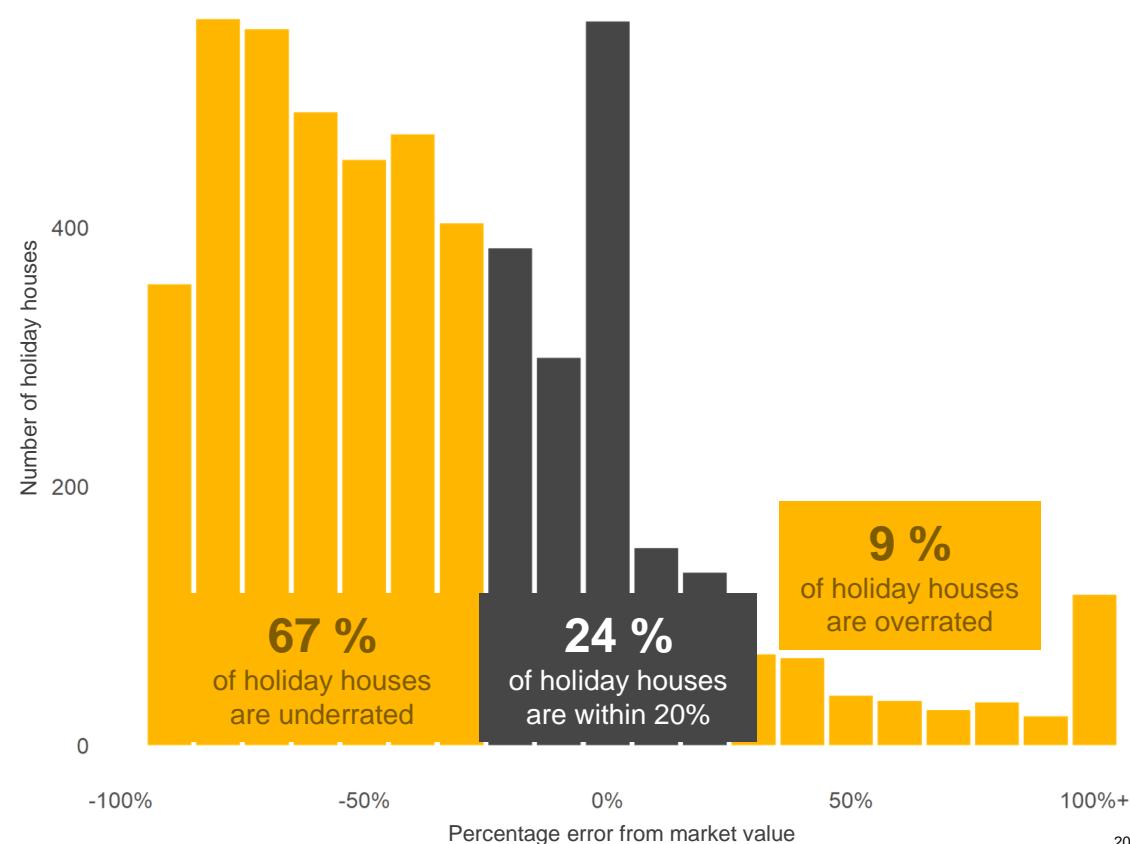
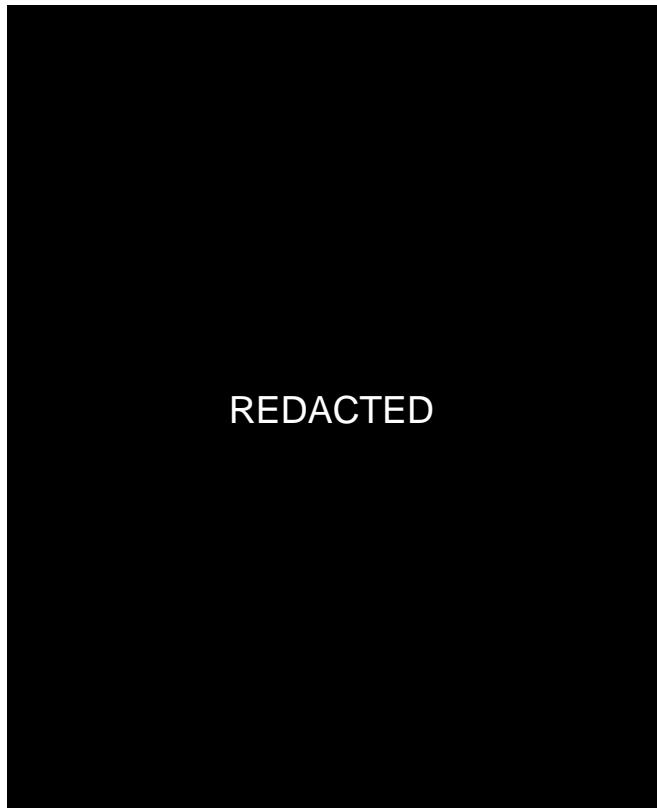
VS

The current model for valuation of primary homes **only accounts for four variables**:  
Area, age, home type and location (price zone)

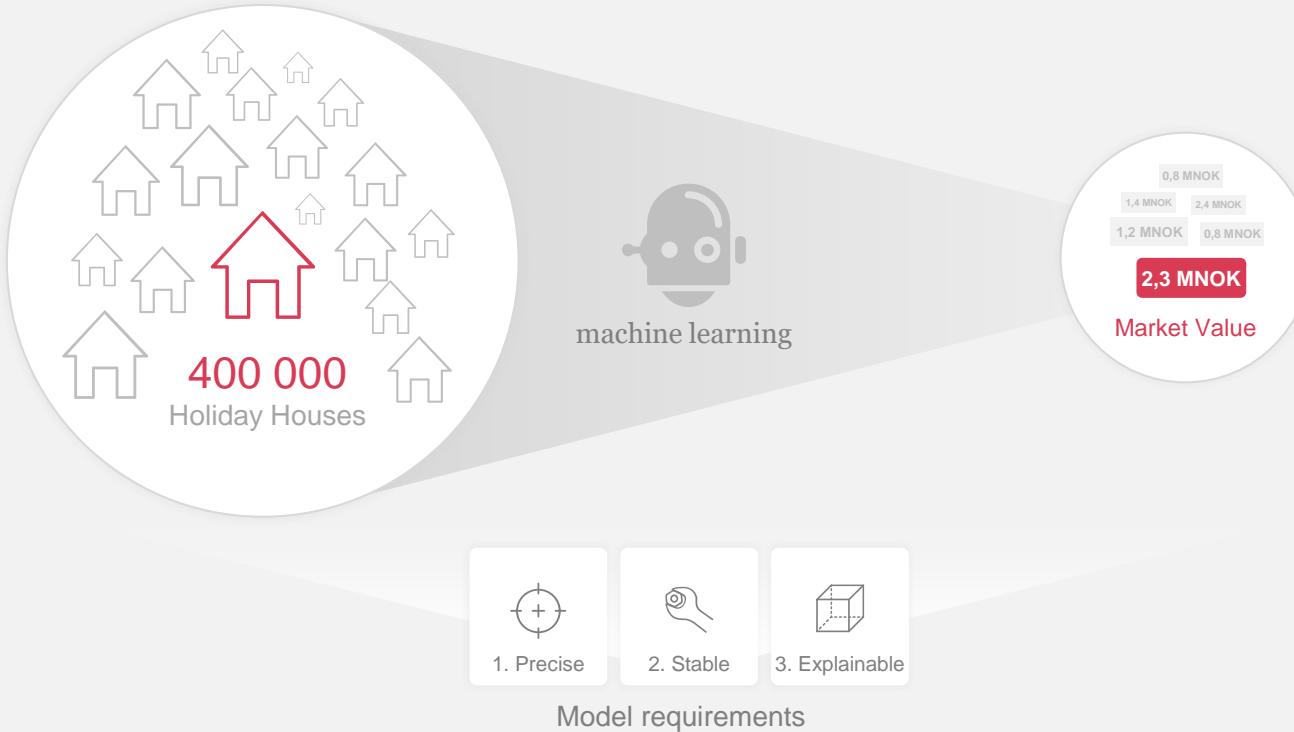




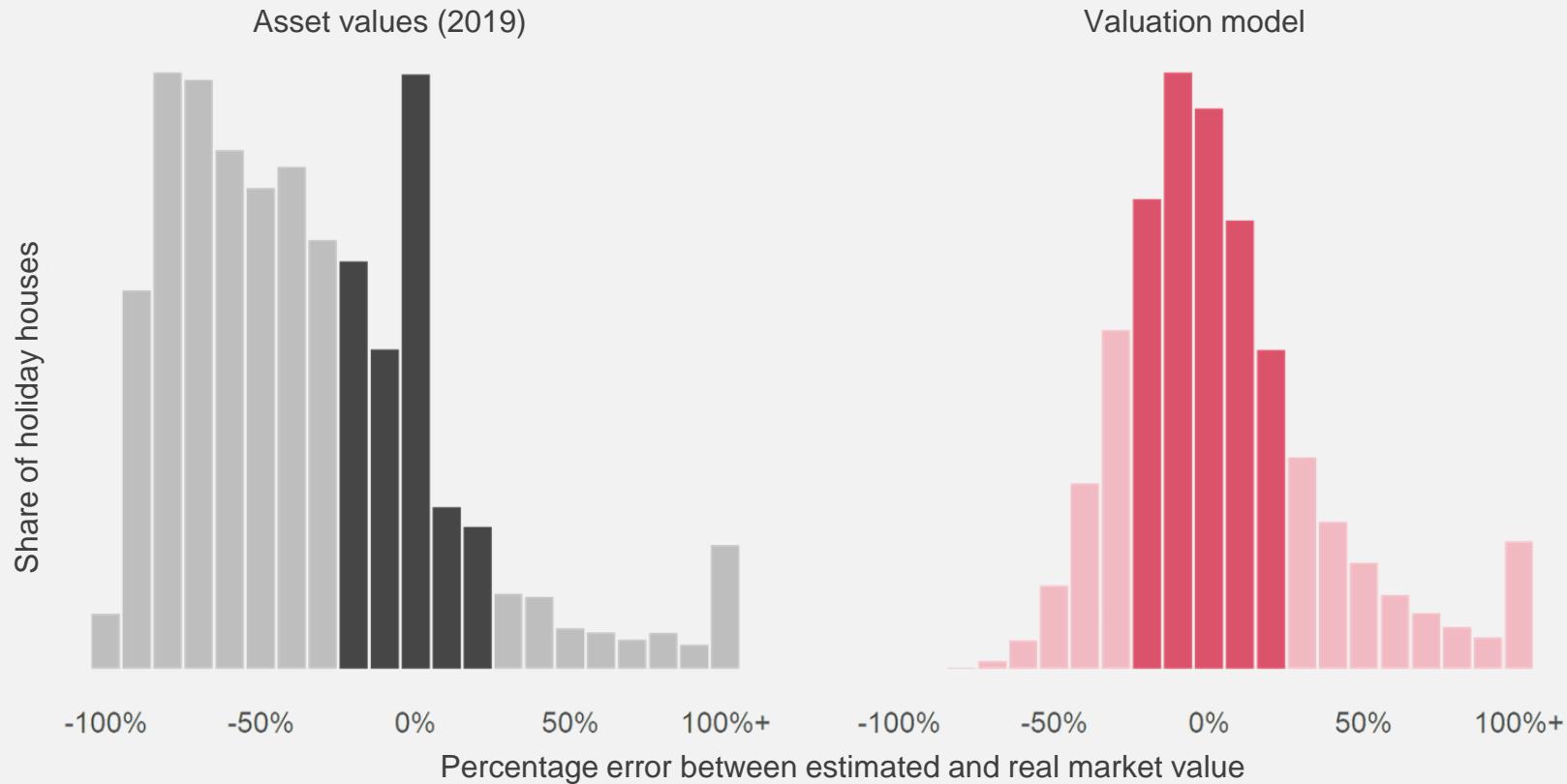
In 2019, the asset values for holiday homes sold during the year had a mean absolute percentage error of **48 percent** measured against the real market value



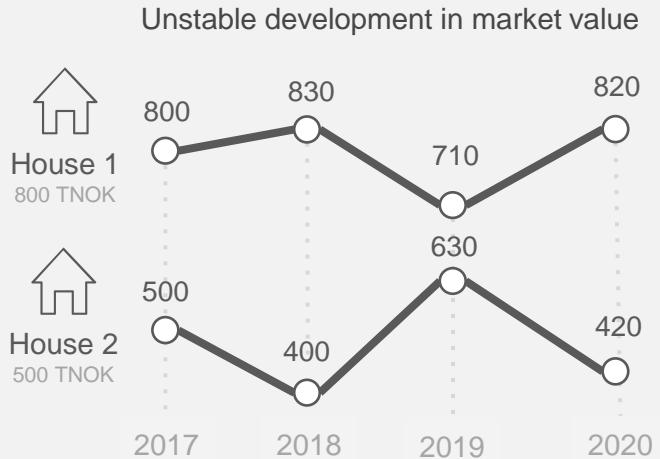
The purpose of the project is to estimate the market value of all Norwegian holiday homes in a **precise**, **stable** and **explainable** manner using current available data



## 1. Precision: The model should give more precise and fair asset values



## 2. Stable: The valuations should be stable over time



	Annual market development				COV
House 1	800	830	710	820	► 7%
House 2	500	400	630	420	► 21%
Average coefficient of variation				14%	

	Annual market development				COV
House 1	800	800	820	840	► 2%
House 2	500	510	510	540	► 4%
Average coefficient of variation				3%	

## 2. Stable: The valuations should be stable within the same area

Unstable market value development in municipalities



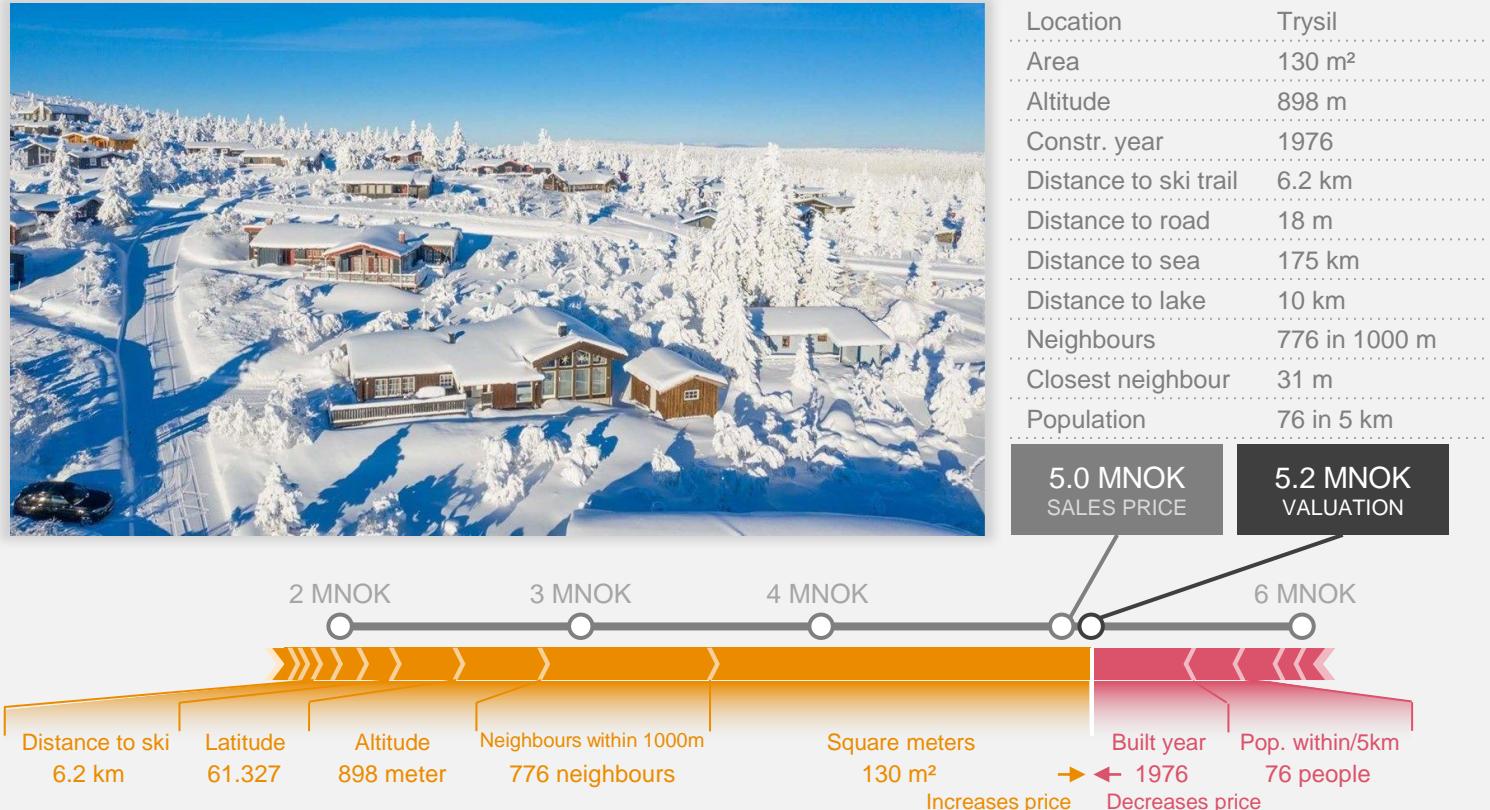
	Annual market development			Std.dev
Municipality 1	8%	-10%	2%	9%
Municipality 2	5%	-20%	-12%	13%
Average standard deviation				
				11%

Stable market value development in municipalities

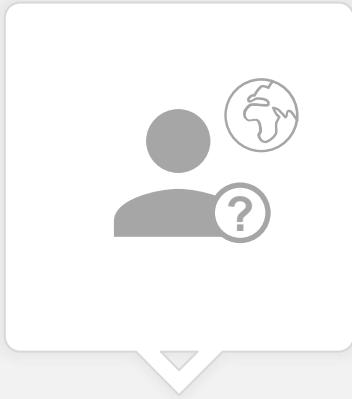


	Annual market development			Std.dev
Municipality 1	2%	-1%	0%	2%
Municipality 2	2%	0%	2%	1%
Average standard deviation				
				1,5%

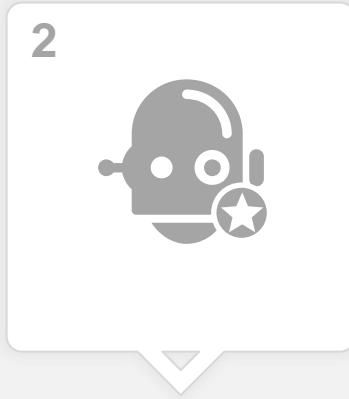
### 3. Explainable: Each valuation must be explainable in a simple manner



Today's valuation models are too crude and there is a great potential for utilizing more granular data



New data sources for valuation



Machine learning for valuation



Explainability of machine learning

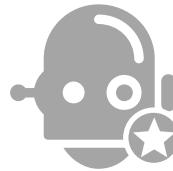
Today's valuation models are too crude and there is a great potential for utilizing more granular data

1



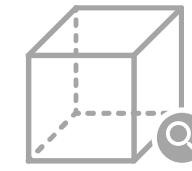
New data sources for valuation

2



Machine learning for valuation

3



Explainability of machine learning

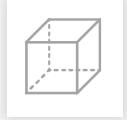
The Norwegian property register currently lacks use area and year of construction for **30%** and **60%** of the holiday homes, respectively



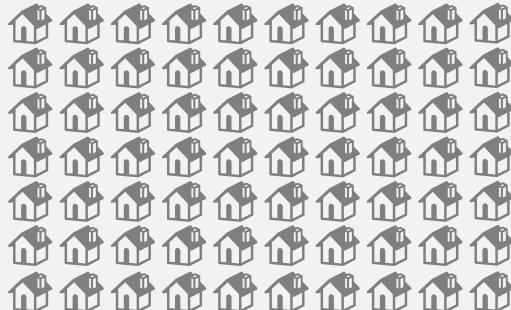
MACHINE LEARNING



EXPLAINABILITY



30 % of holiday homes lacks registered use **area**

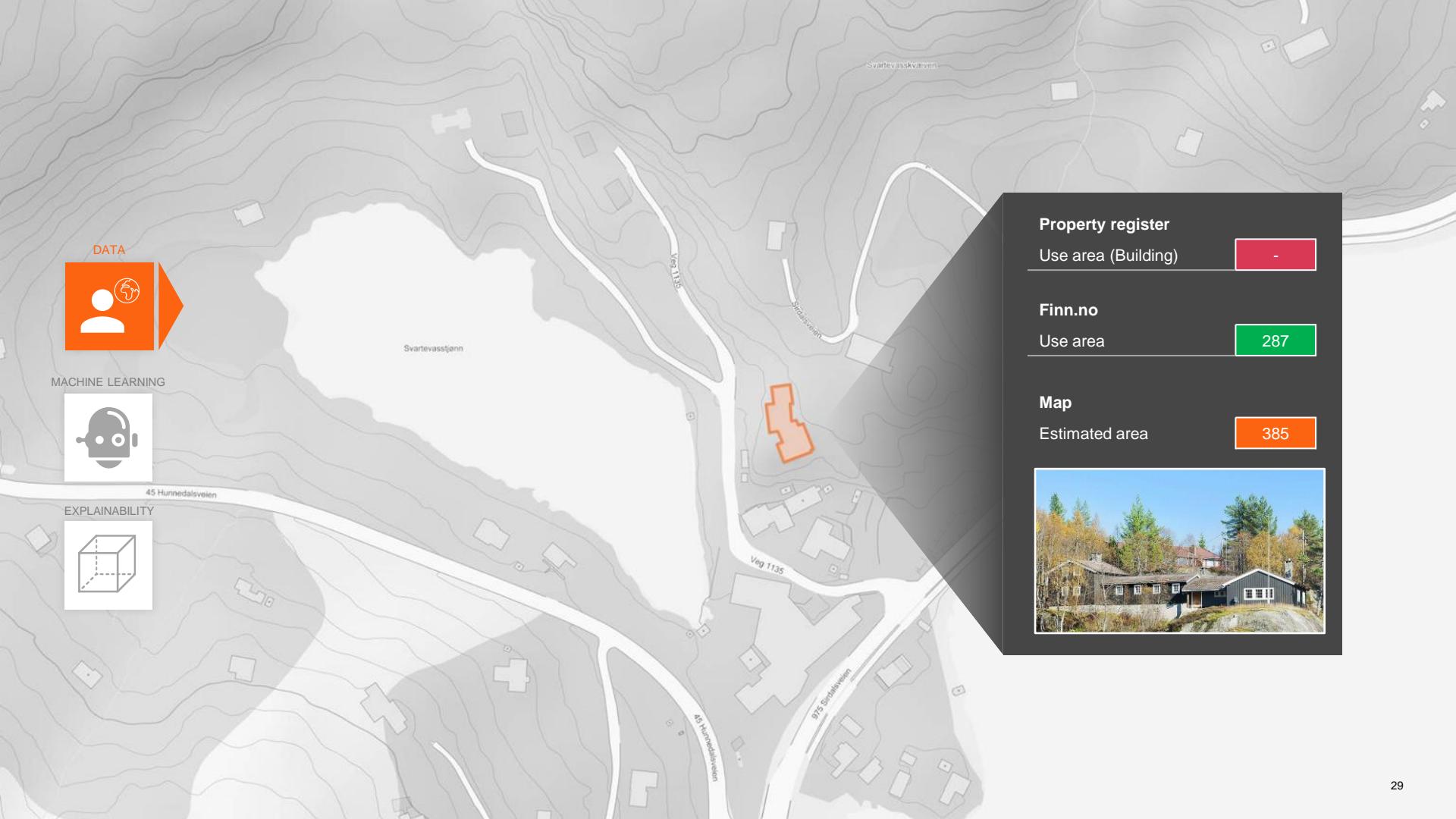


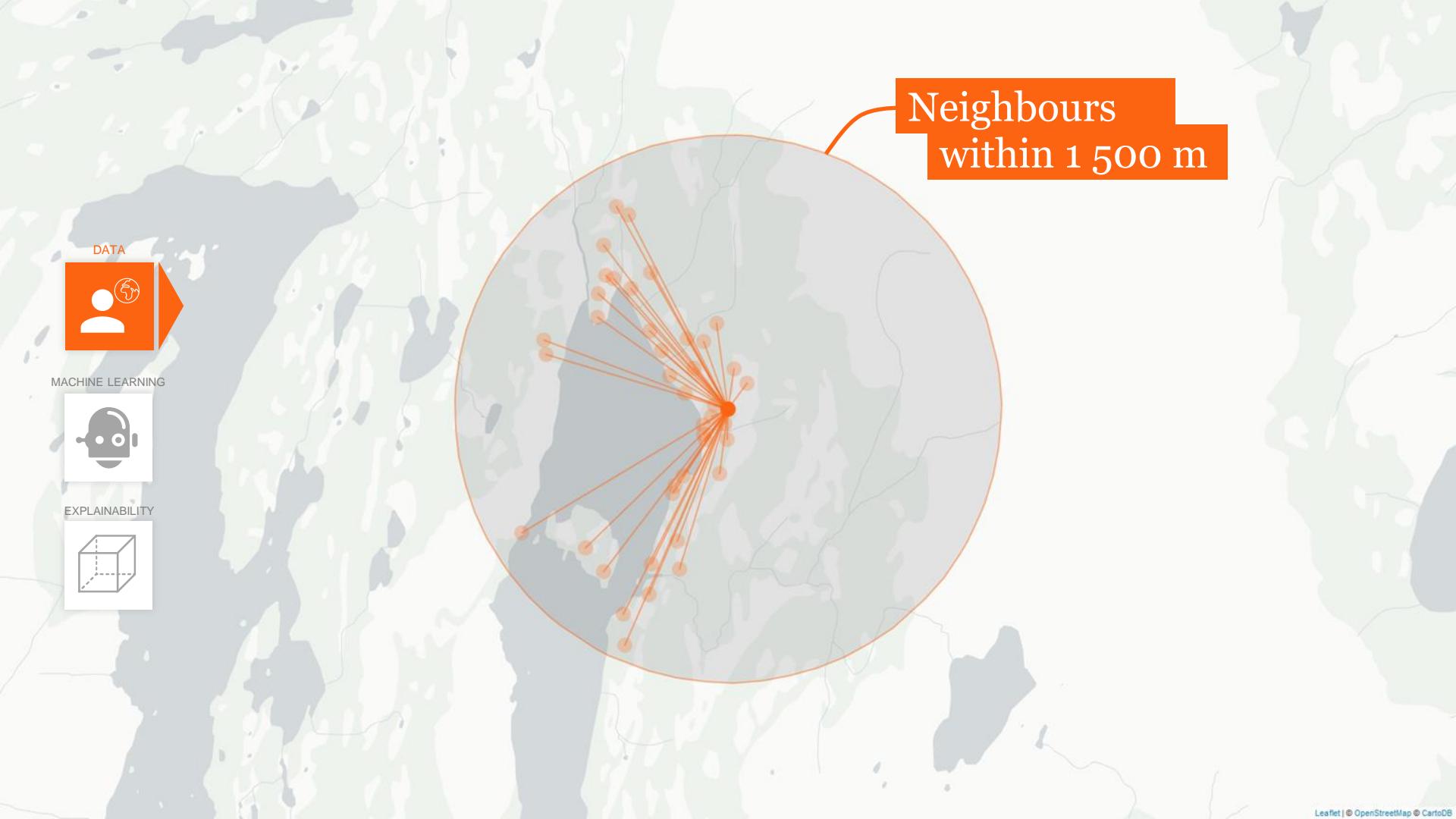
60 % of holiday homes lacks registered **year of construction**



... and further information is seldom registered

STANDARD	RUNNING WATER
STORIES	ELECTRICITY
PARKING	BEDROOMS





Neighbours  
within 1 500 m

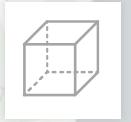
DATA



MACHINE LEARNING



EXPLAINABILITY







Distance to  
lake

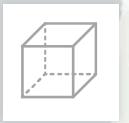
DATA



MACHINE LEARNING



EXPLAINABILITY



Distance to  
sea

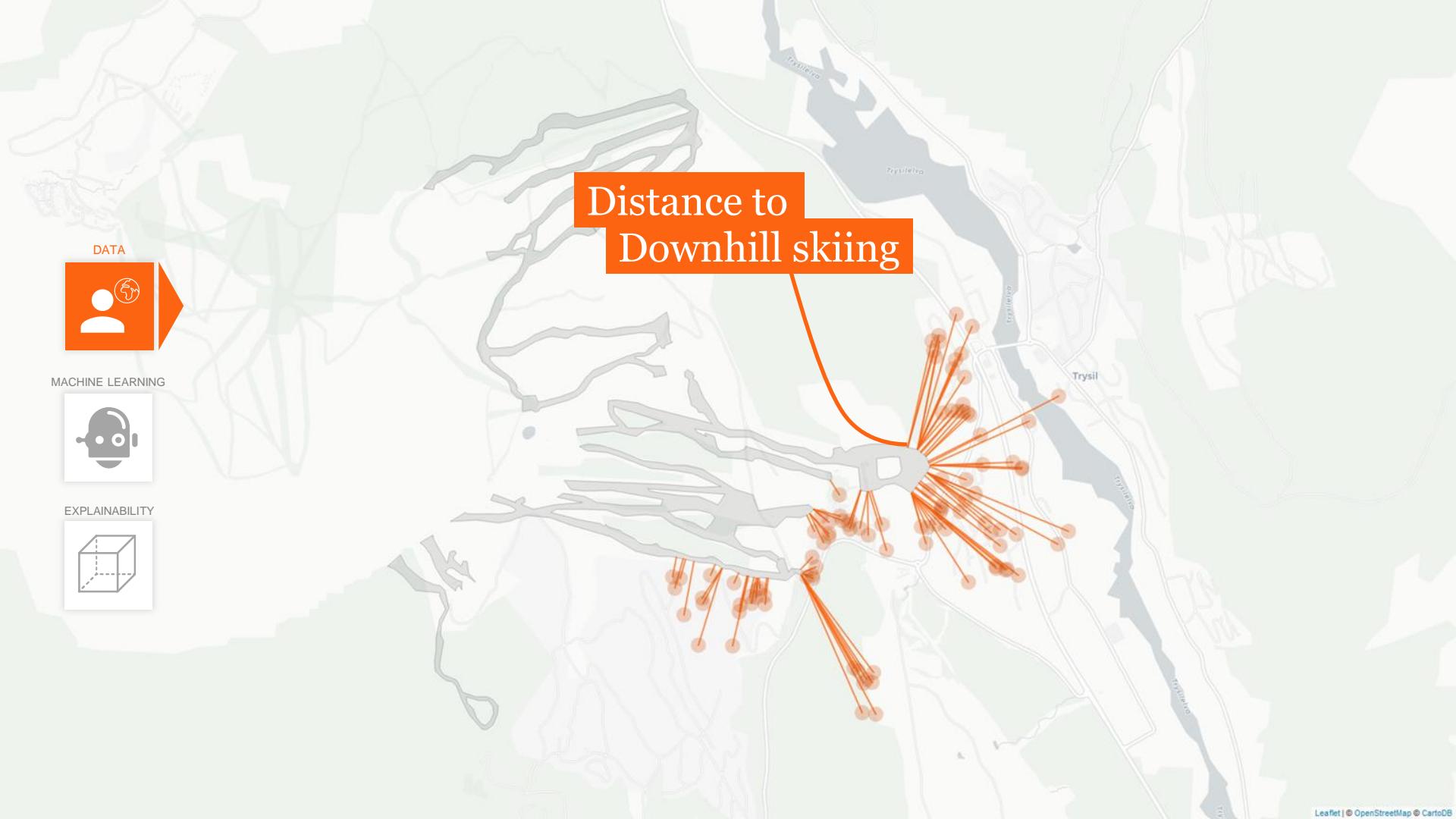


MACHINE LEARNING



EXPLAINABILITY





# Distance to Downhill skiing

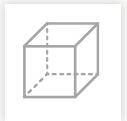
DATA

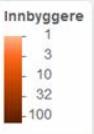


MACHINE LEARNING

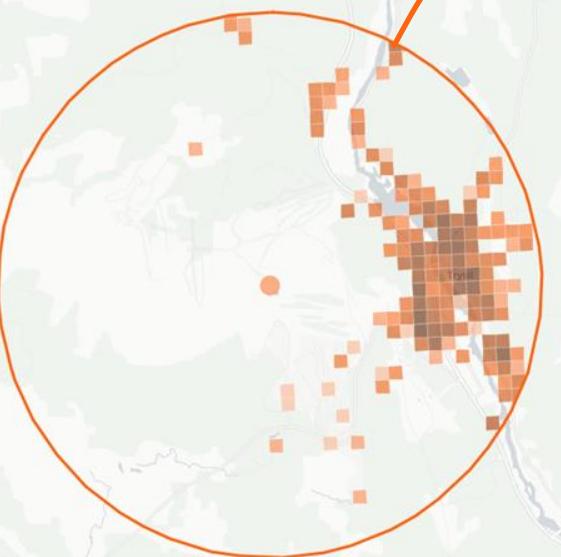


EXPLAINABILITY





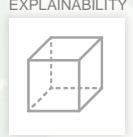
Population  
within 1 000 m



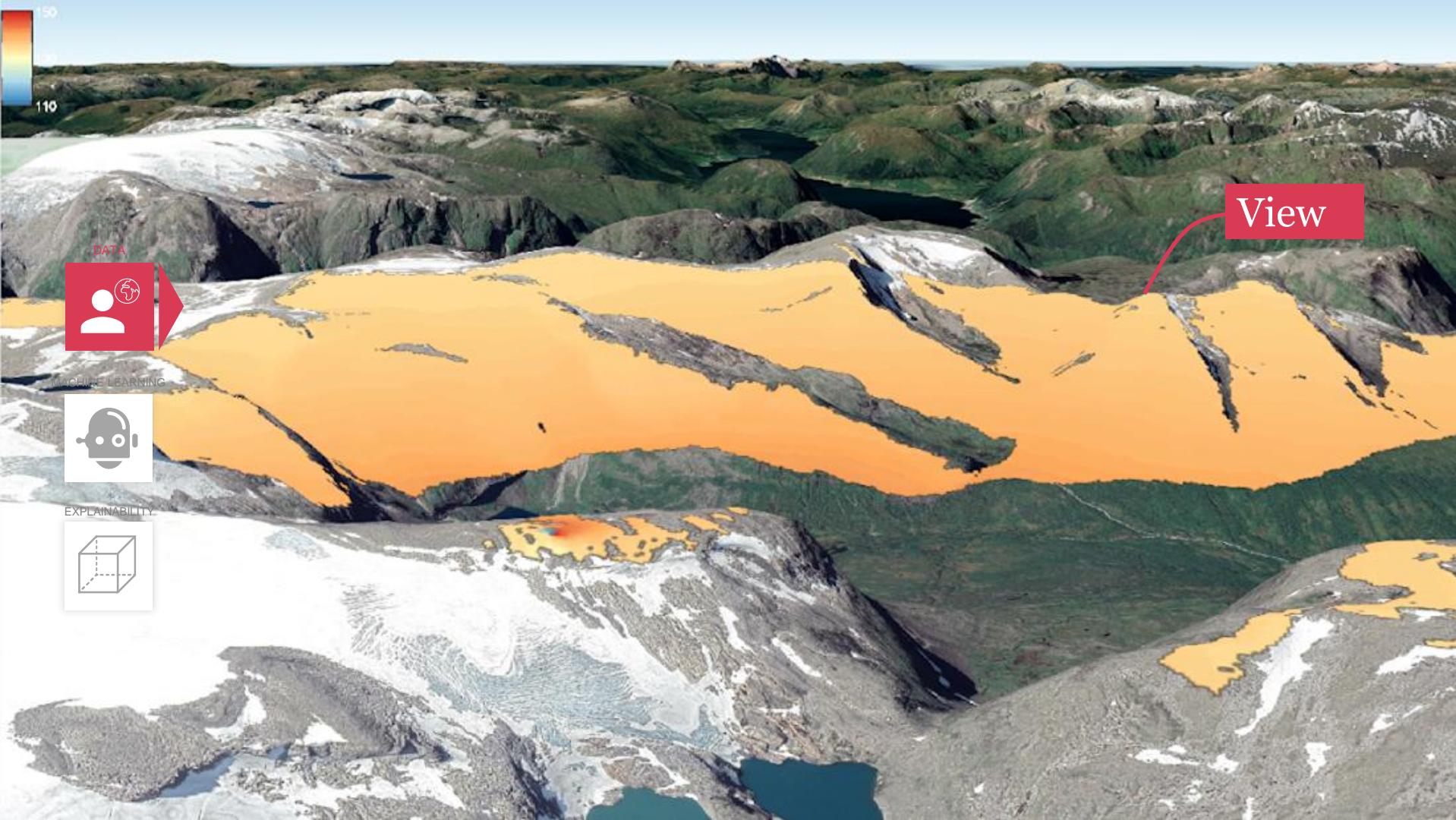
DATA

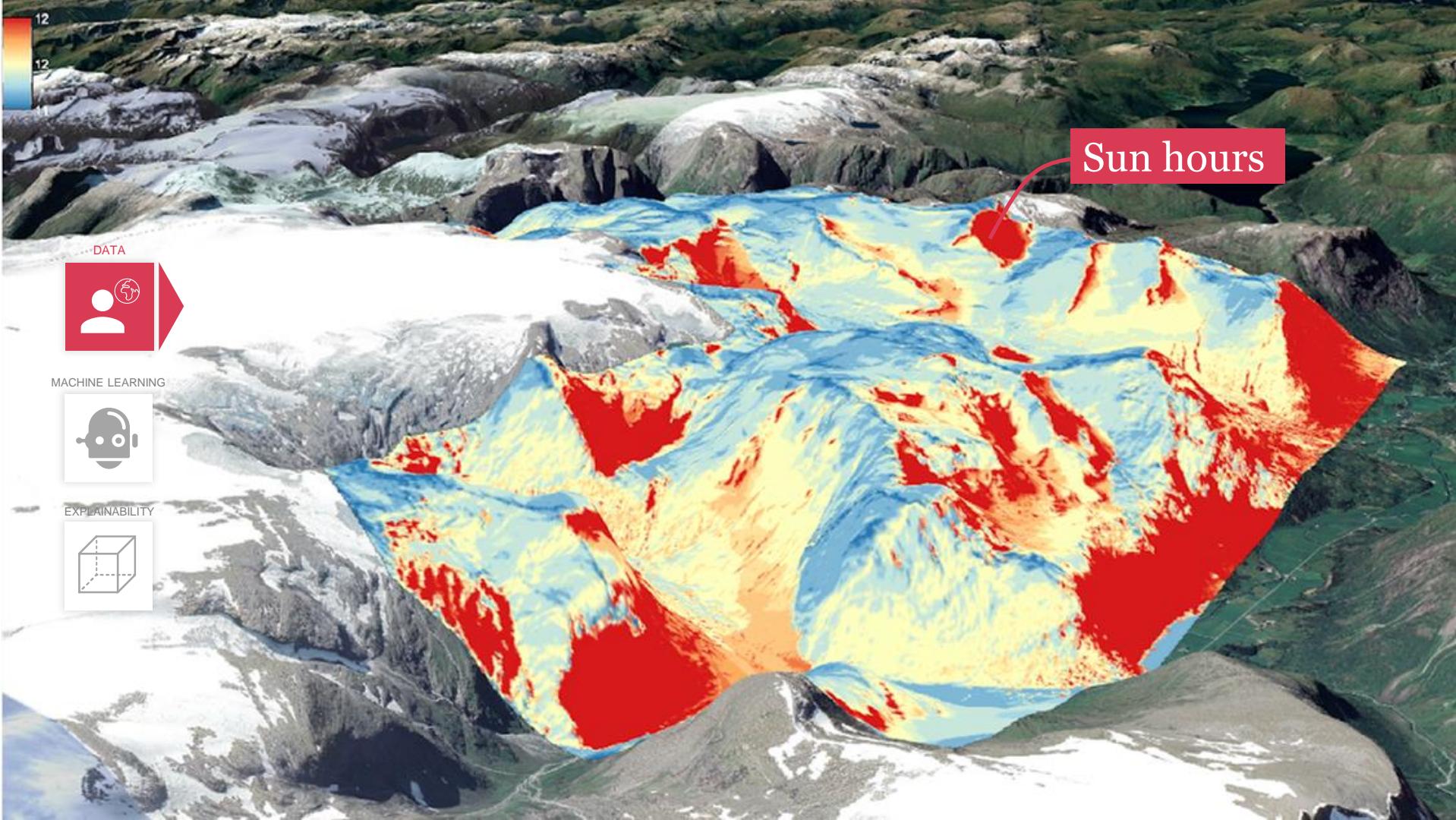


MACHINE LEARNING



EXPLAINABILITY





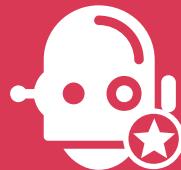
Traditional linear models does not capture the complex relationships that determine the market value of holiday homes

1



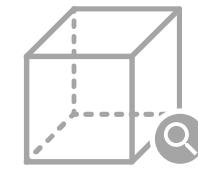
New data sources for valuation

2



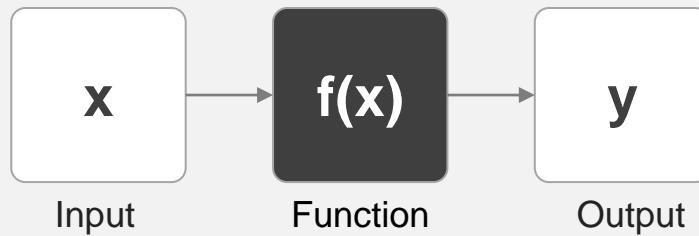
Machine learning for valuation

3

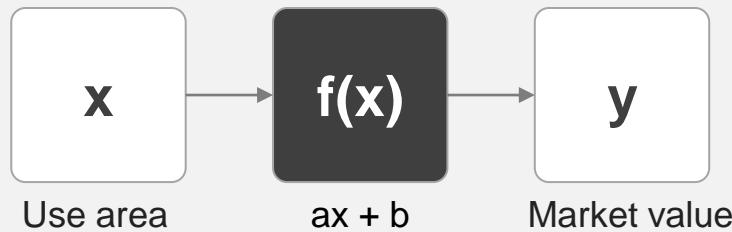


Explainability of machine learning

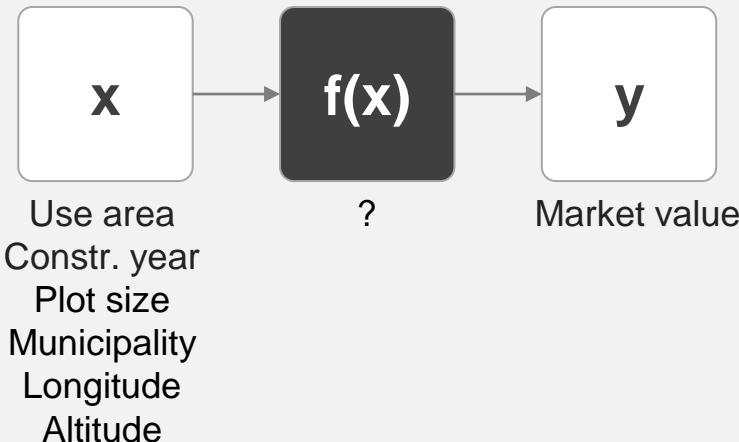
In traditional analyzes we use functions to calculate results (output) based on variables (input)



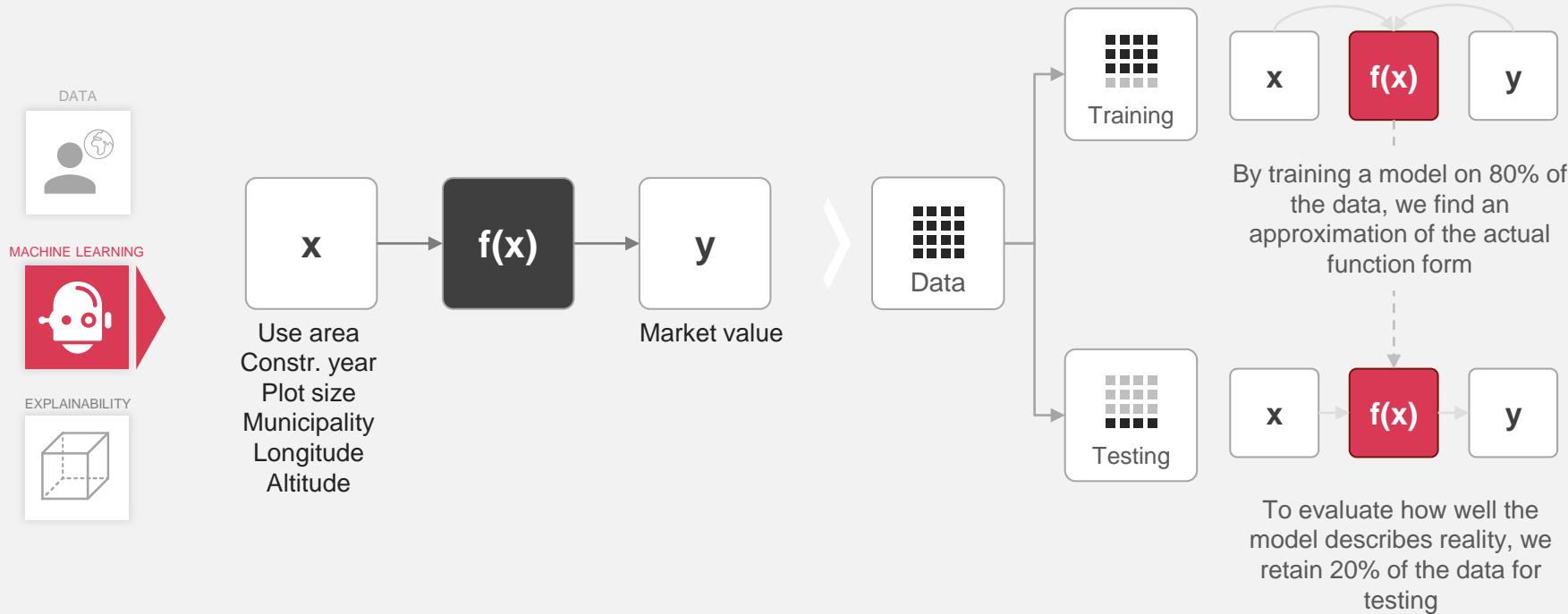
In traditional analyzes we use functions to calculate results (output) based on variables (input)



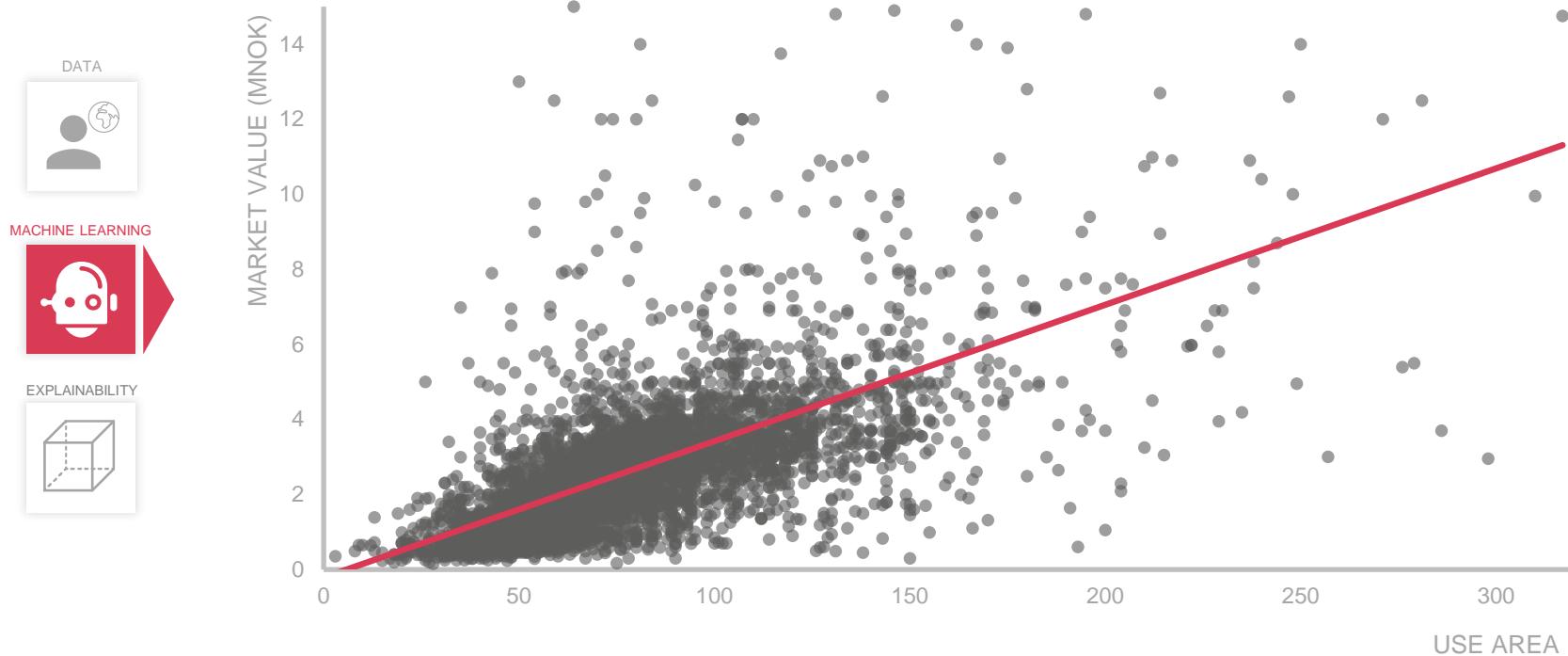
# Machine learning involves using data to find a suitable function



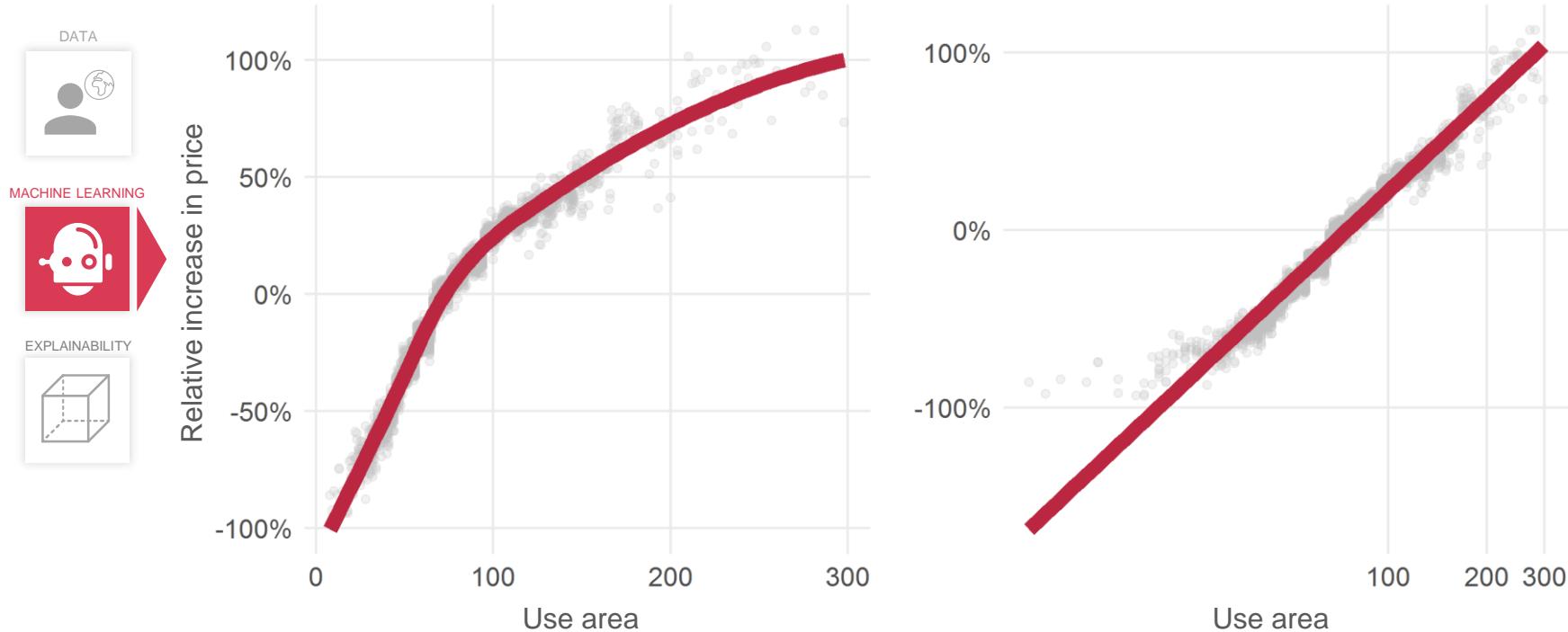
In a typical case, we divide the data in two: We use most (80%) of the data to train the model to find a suitable function. The remaining data (20%) is used for validation



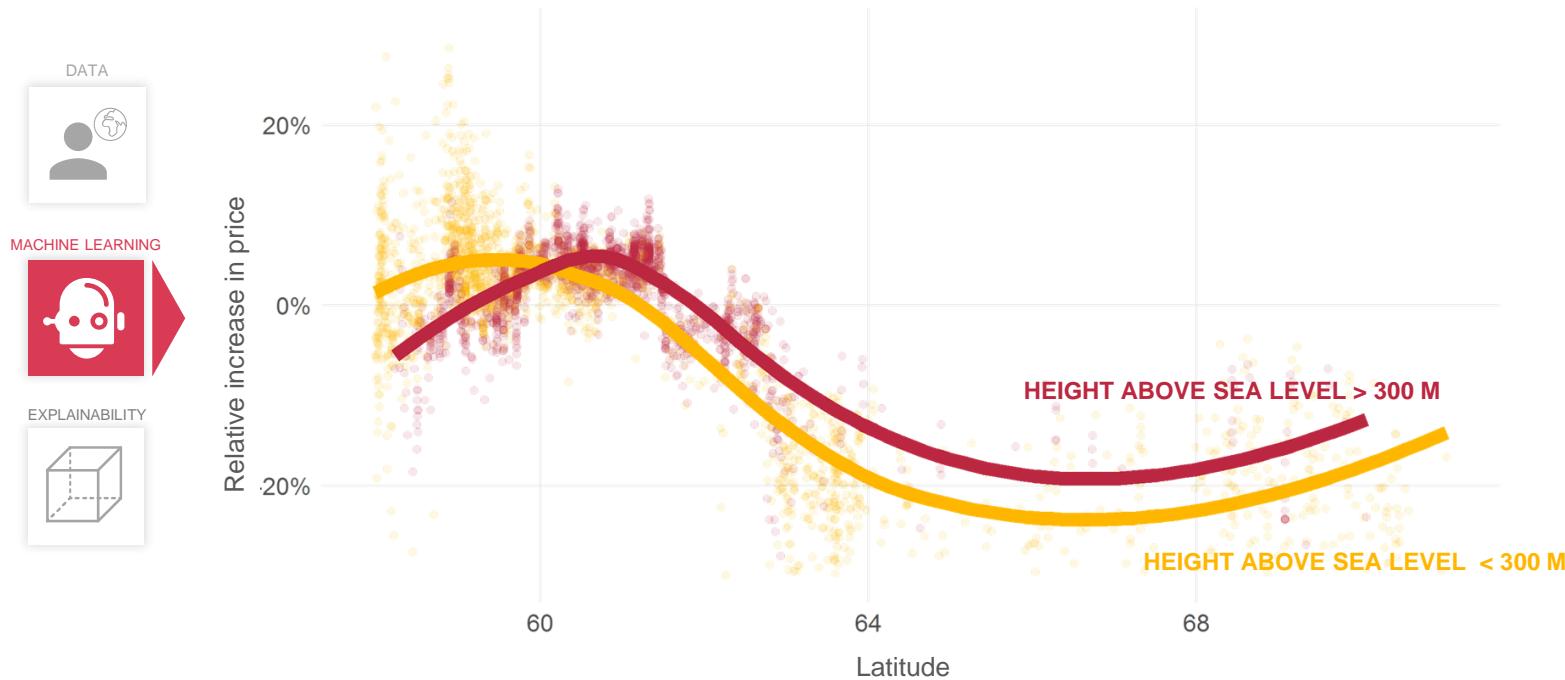
Machine learning can be used to find the linear relationships between input and output that constitute the best linear regression model



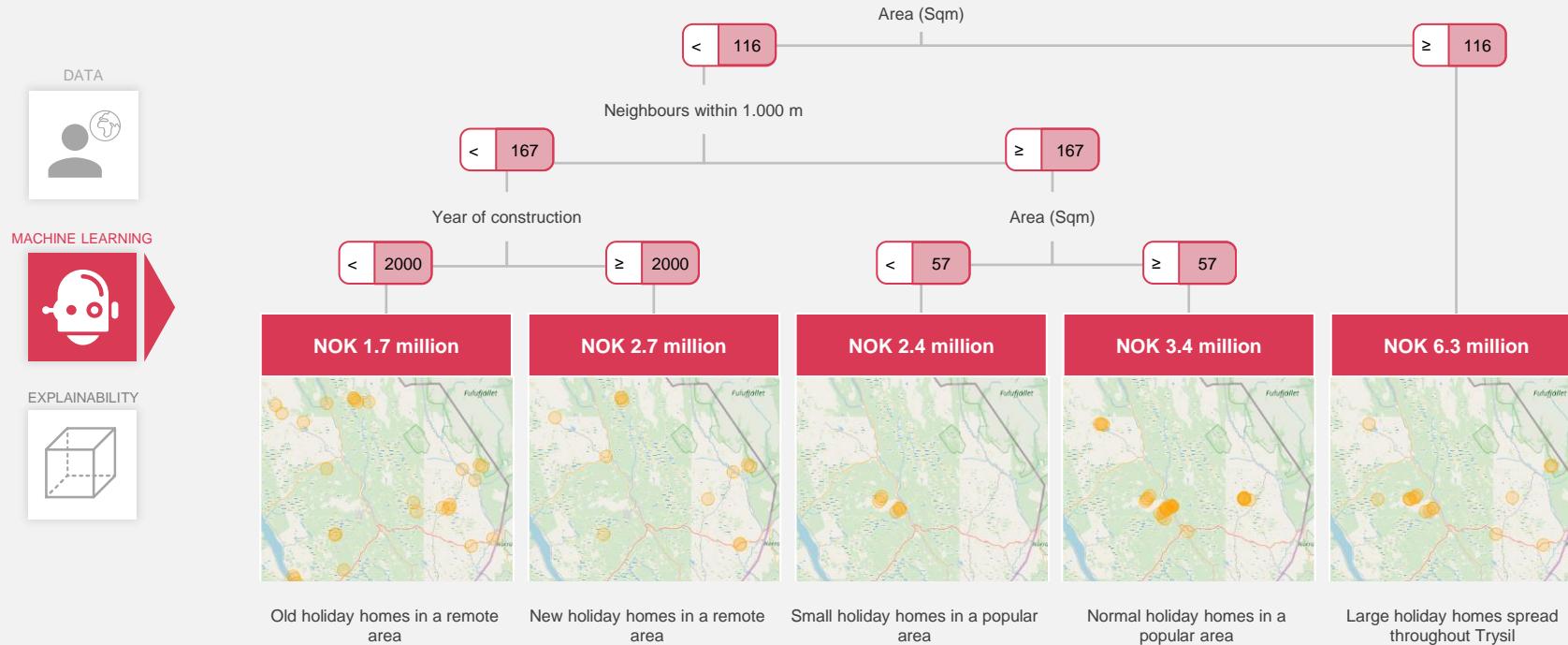
Real world relations are usually **non-linear**. To capture these using linear models, we need to transform the data in advance



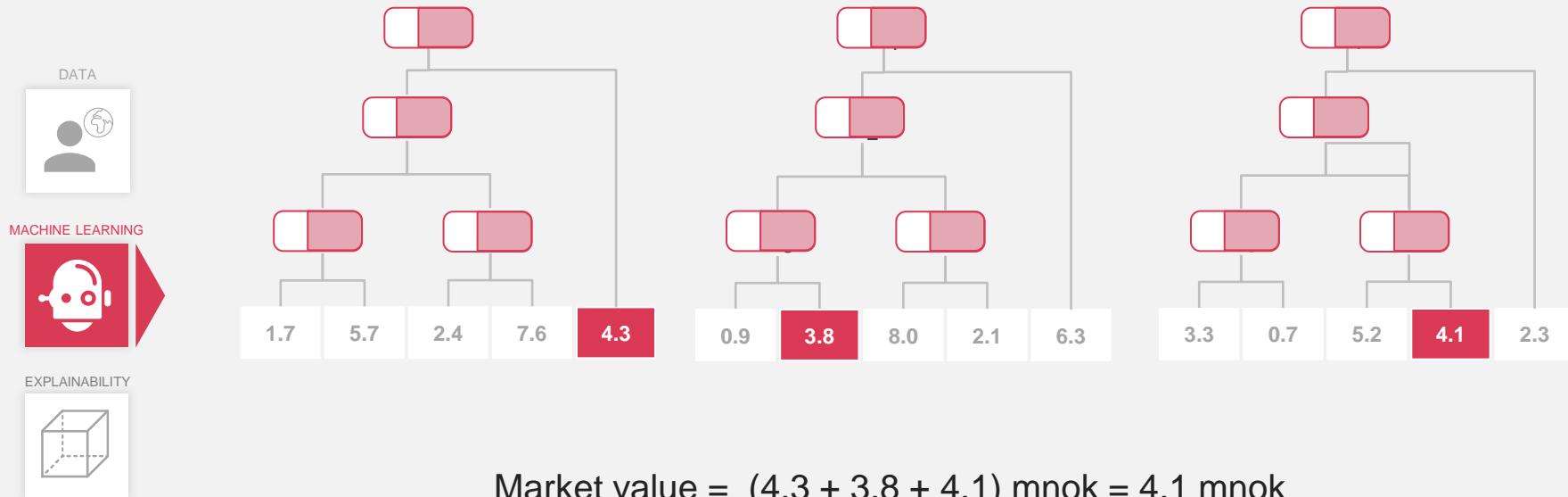
Machine learning can be used to develop models that capture complex relationships and underlying patterns in the data



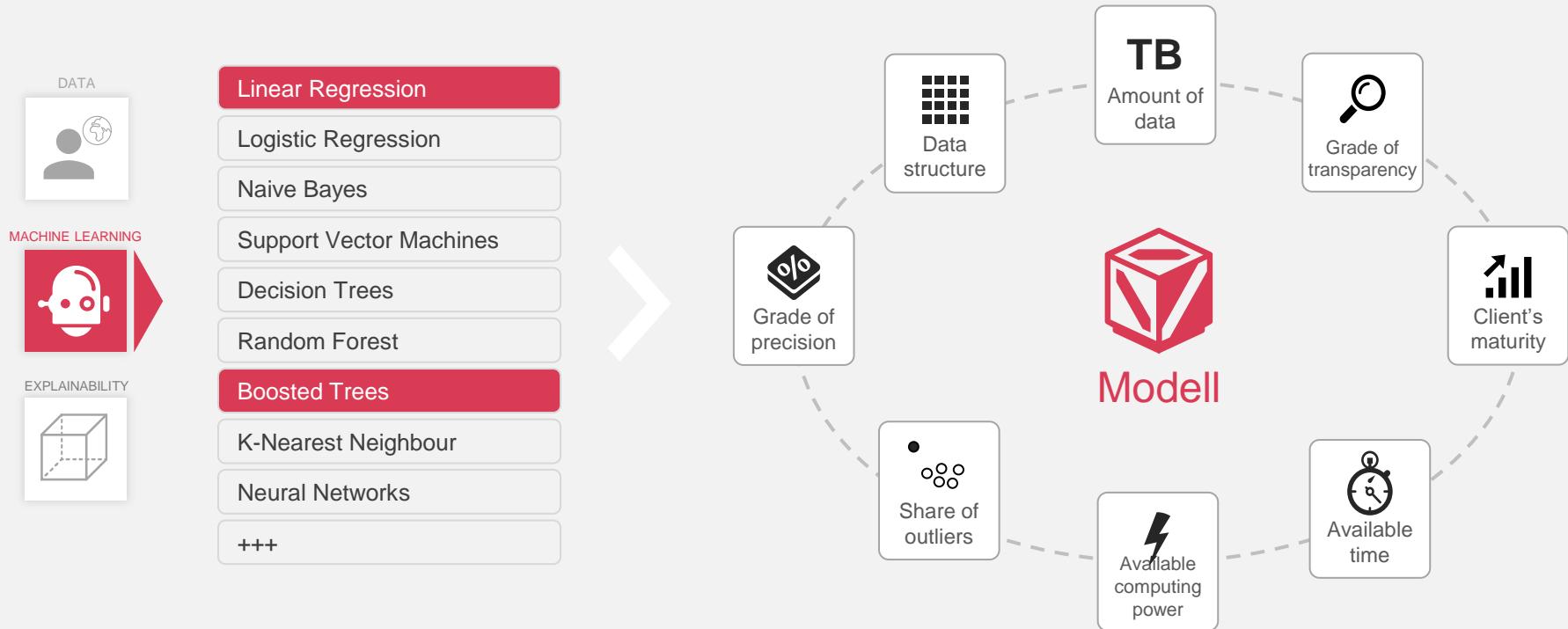
# Machine learning can be used to create decision trees, which capture complex interactions and underlying relationships better than linear models



To increase precision, we can create multiple decision trees and use the average of all predictions as market value



There are many different machine learning models. The right model choice depends on both the problem to be solved and the customer's maturity.



Using machine learning, the project have developed a model that is **more precise** than the current asset values and where the deviations are more evenly distributed



REDACTED

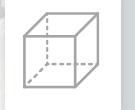
REDACTED



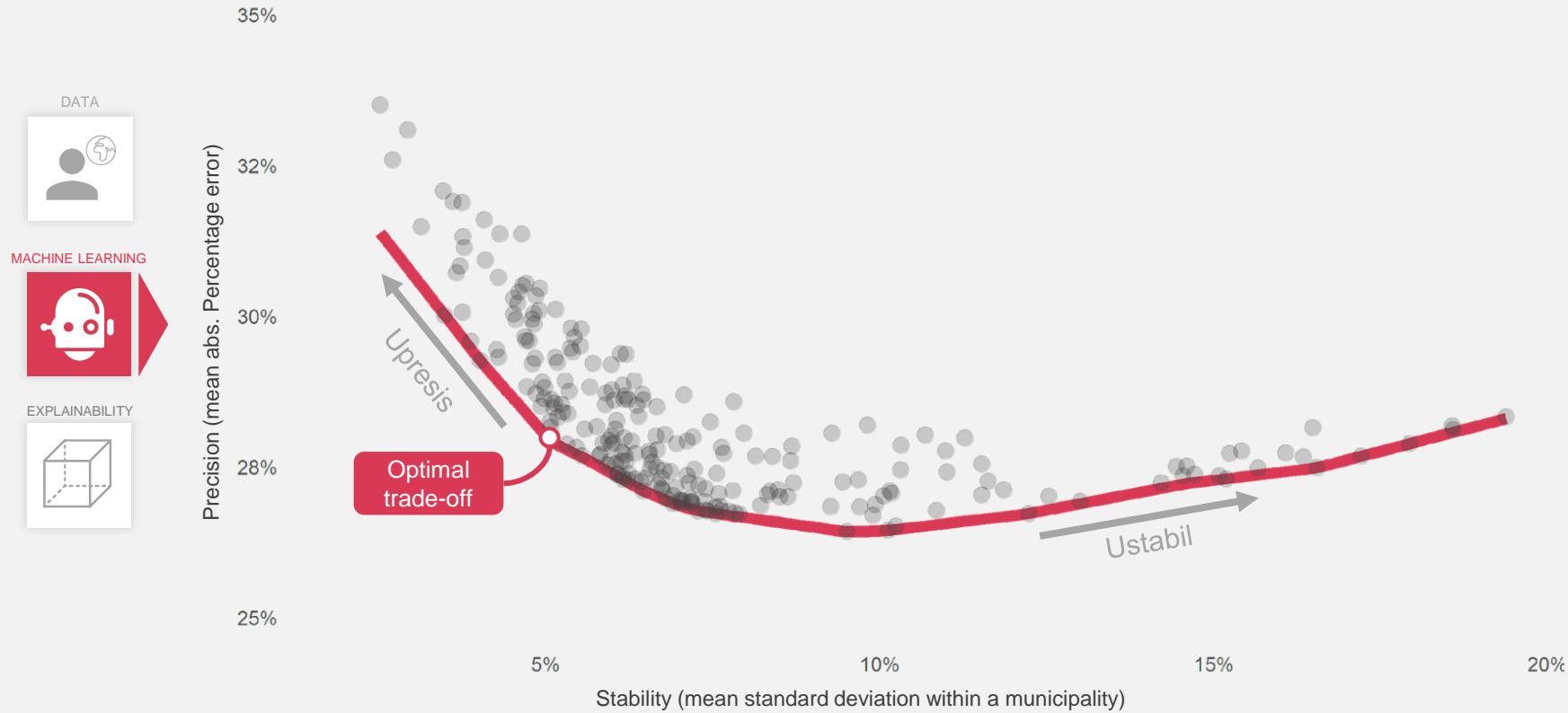
DATA



MACHINE LEARNING



We used hyperparameter tuning to find the optimal trade-off between precision and stability





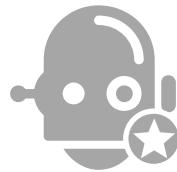
# Valuation models must be explainable to taxpayers

1



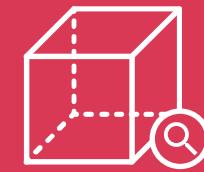
New data sources for valuation

2



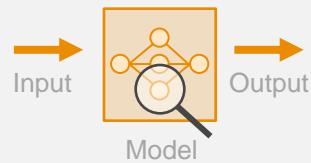
Machine learning for valuation

3



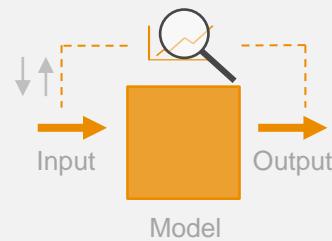
Explainability of machine learning

# Three different methods to explain models (Molnar, 2020)



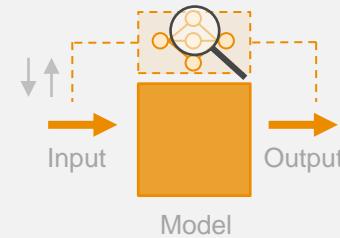
## Methods analysing model components

Linear models and decision trees are examples of naturally explainable models that can be explained by analysing components of the models, ie. model structure and parameters



## Methods analysing model sensitivity

The methods are typically model agnostic, where the model is considered a closed system. The methods analyse how the estimates are affected by changes in the variables values.



## Methods based on surrogate models

Surrogate models are naturally explainable models designed to copy the behaviour of unexplainable models. The methods are based on the explanations of the surrogate model

A photograph of two men in dark tuxedos shaking hands. The man on the left has a full white beard and grey hair, smiling warmly at the other man. The man on the right has white hair and glasses, looking down at the handshake. They are positioned in front of a blurred background of warm, autumn-like colors.

DATA



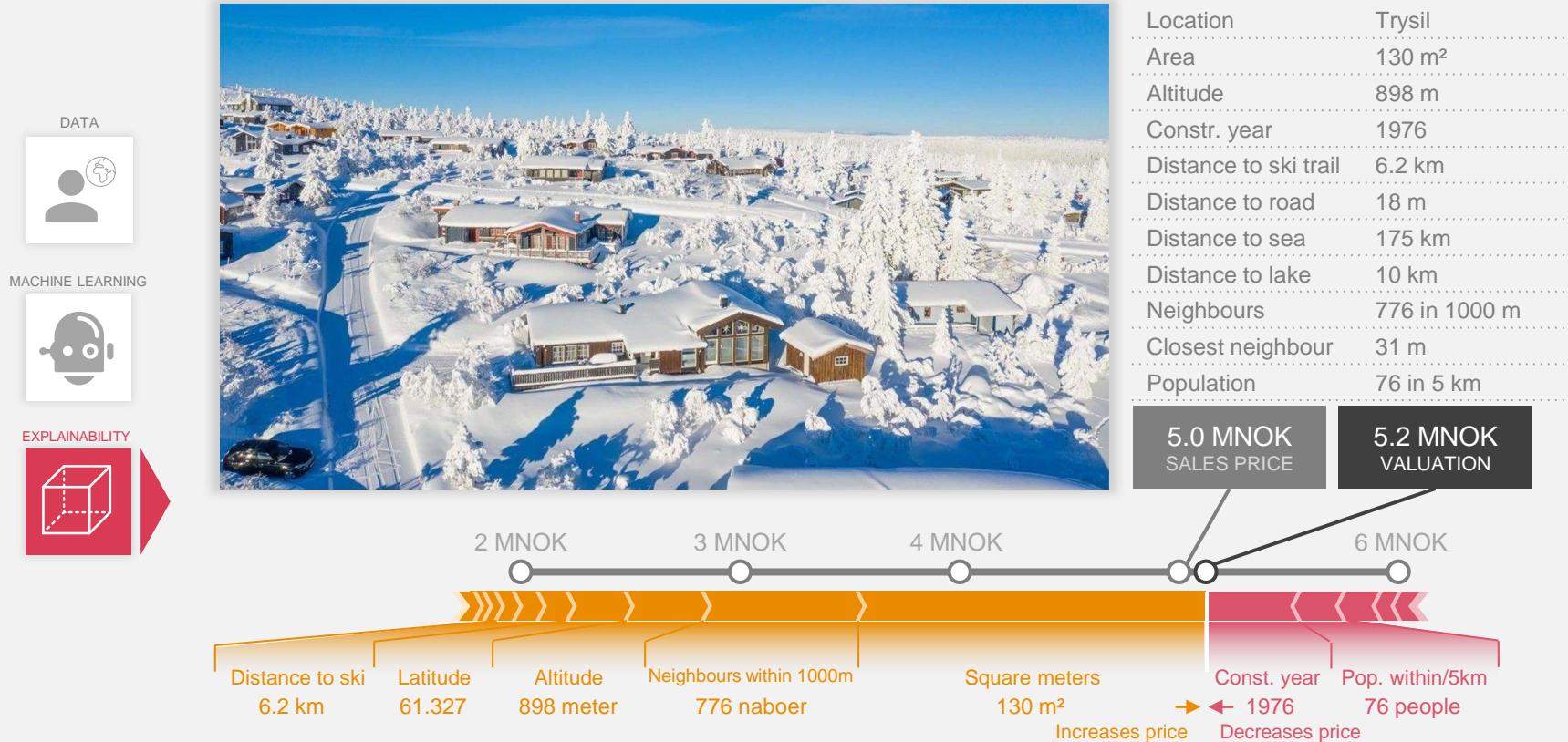
MACHINE LEARNING



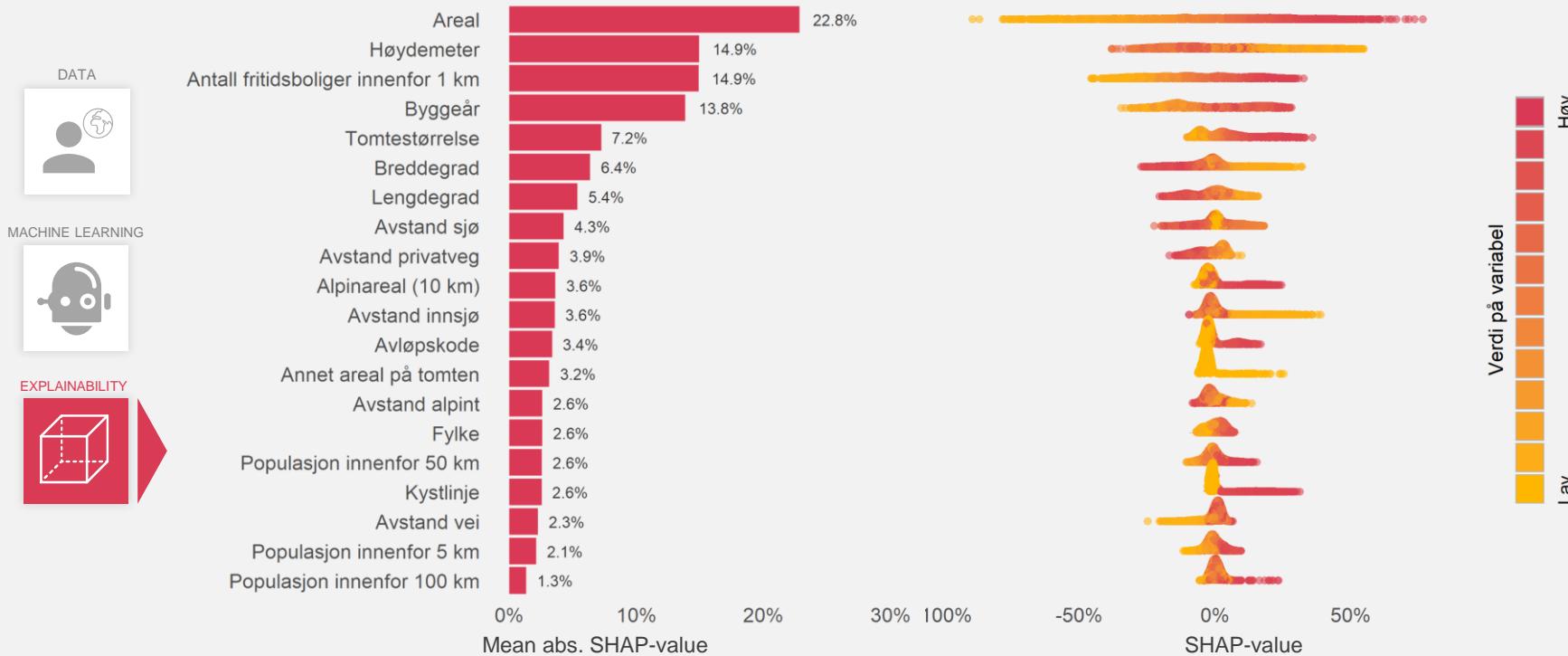
EXPLAINABILITY



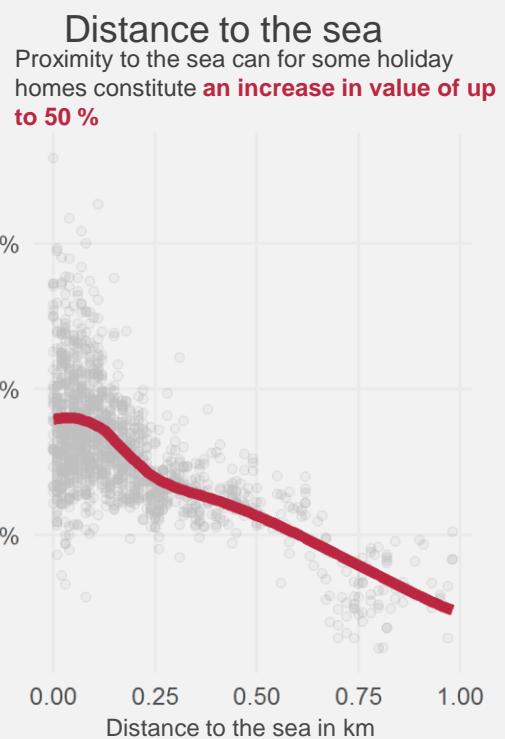
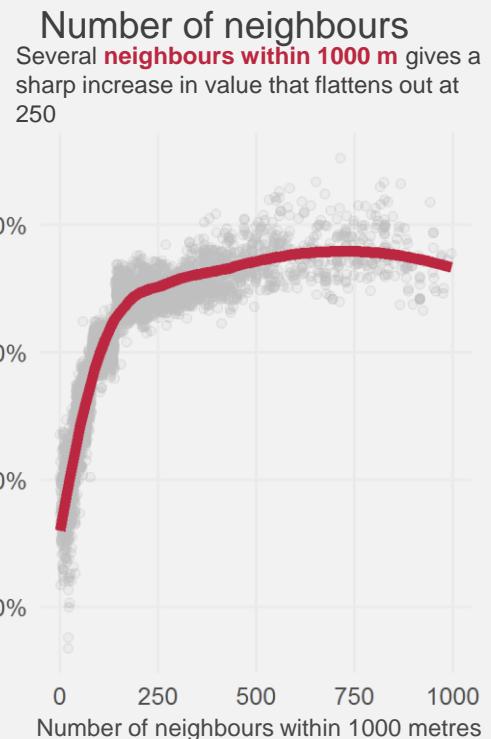
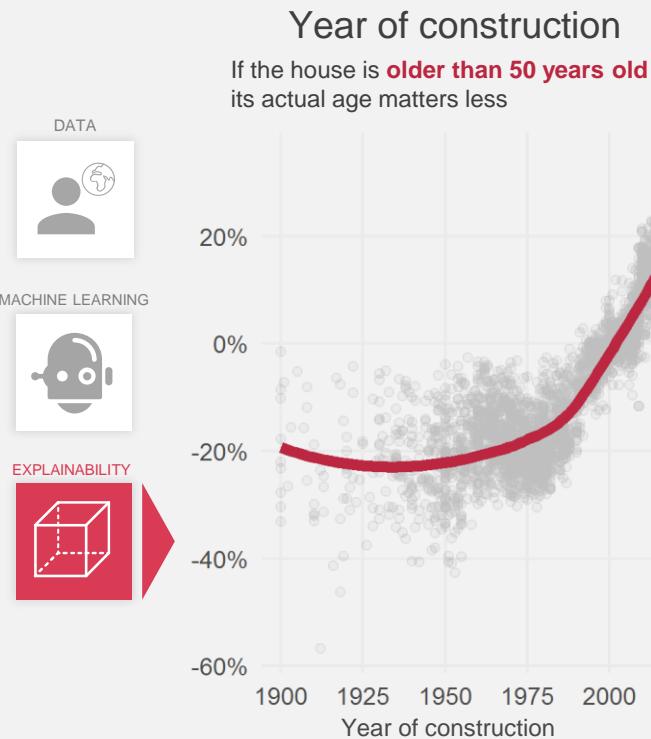
# SHAP values makes it possible to explain how the various price drivers affect the final prediction for each individual holiday home



By aggregating the SHAP values, the **most important** variables can be derived. Structural variables such as area, build year and plot size **explains less than 50 %** of the valuation on average



By analyzing how complex machine learning algorithms emphasise different variables, it is possible to gain insight on the underlying patterns in the data



# 03

## Code along and exercises

# Case: Buying a house



# Apply for relevant positions at [pwc.no/karriere](http://pwc.no/karriere)



**Full-time position**

**Relevant for** 5th grade

**When** August 2022

**Where** Bergen

**Deadline** October 10th

**Data & Analytics Internship**

**Relevant for** 4th grade

**When** January 2022

**Where** Bergen

**Deadline** TBA