

ABSCHNITT 5

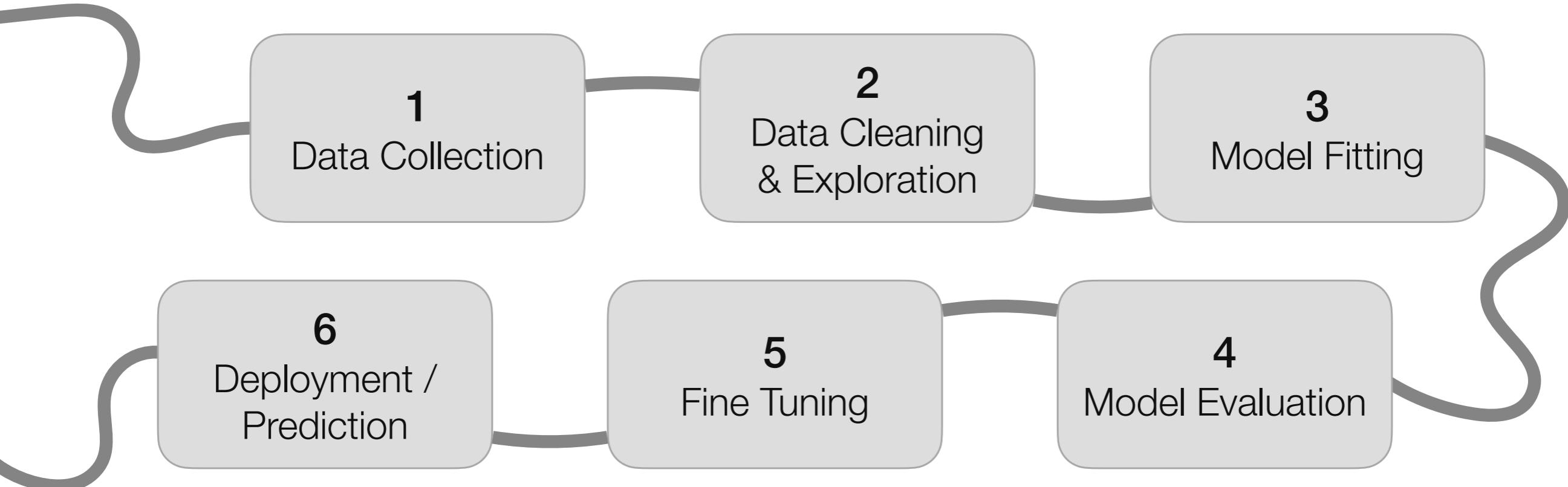
DIE DATA SCIENCE PIPELINE

EIN REZEPT FÜR MACHINE LEARNING PROJEKTE

WAS IST EINE DATA SCIENCE PIPELINE?

Eine Data Science Pipeline umfasst die prototypischen **Schritte**, die es braucht, um ein **Machine Learning** Modell erfolgreich zu **implementieren**.

EINE DATA SCIENCE PIPELINE UMFASTT (UNGEFÄHR) SECHS SCHRITTE:



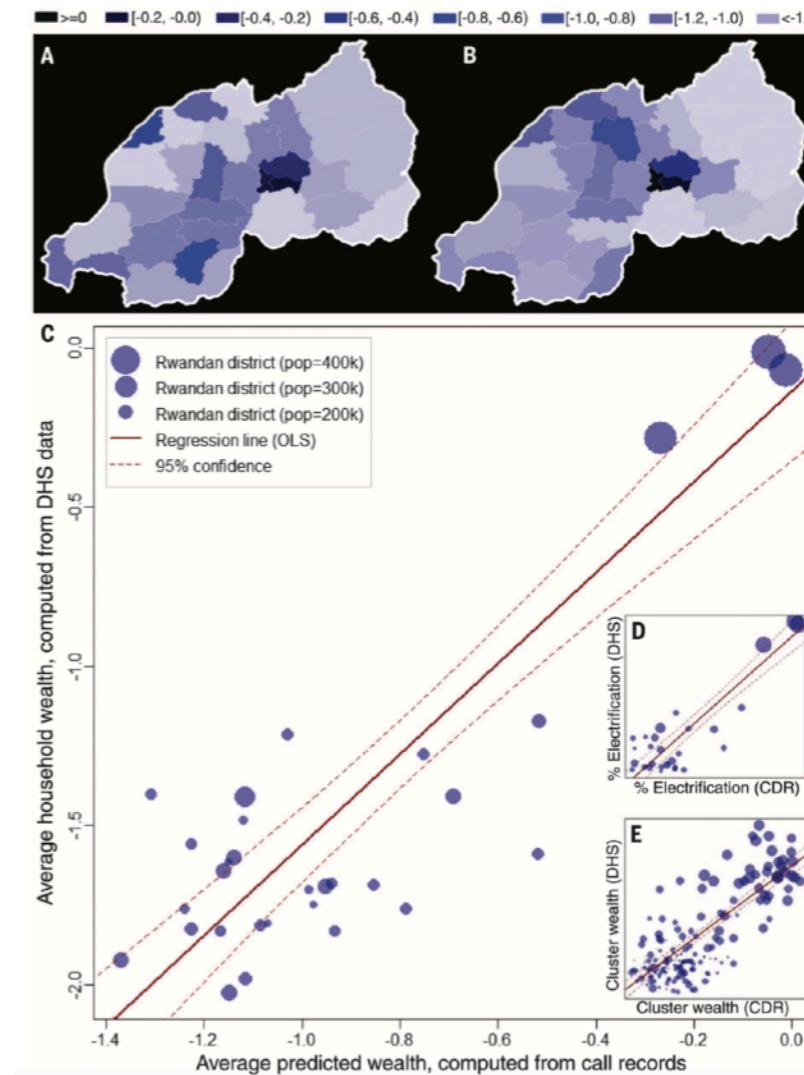
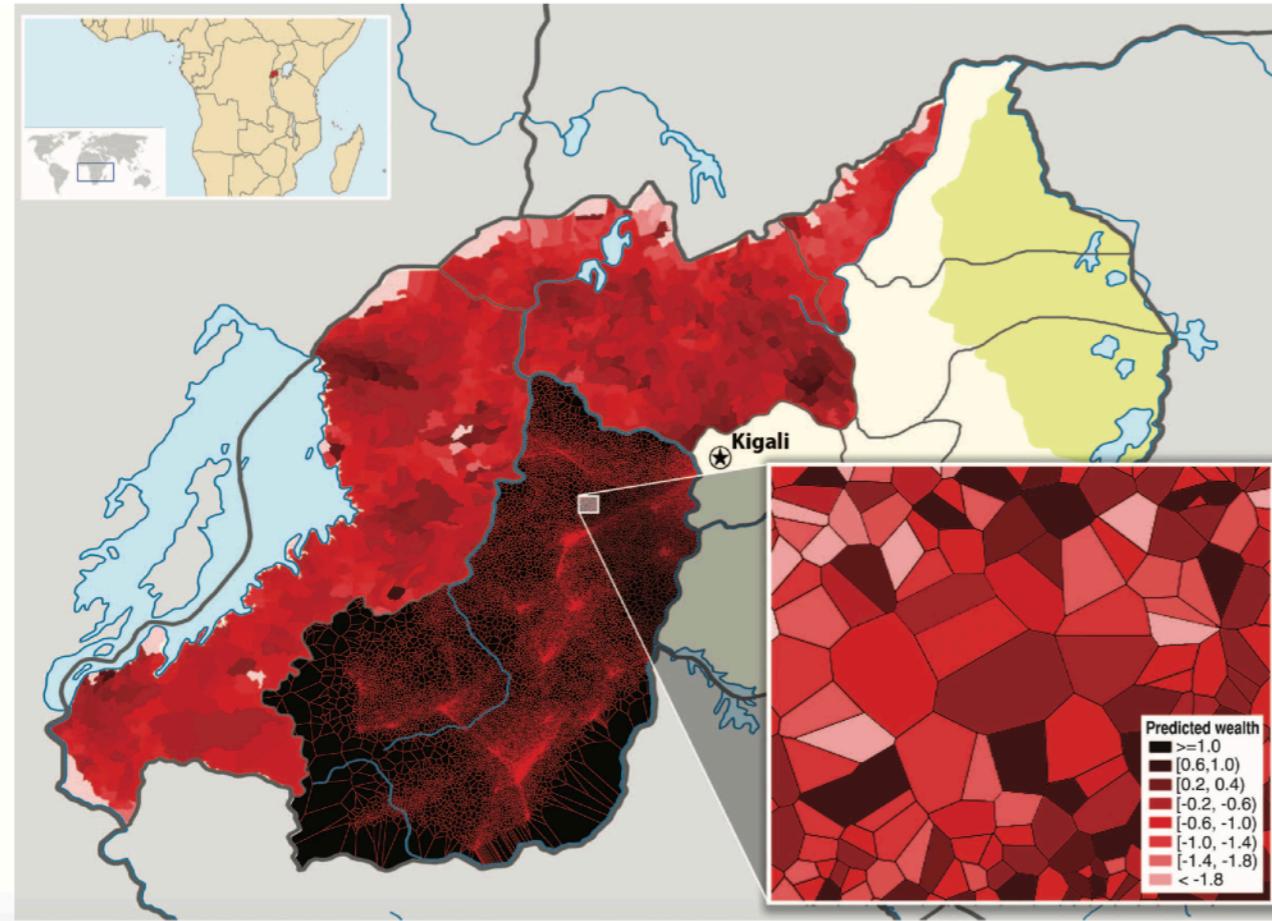
1 // DATA COLLECTION

custom-made

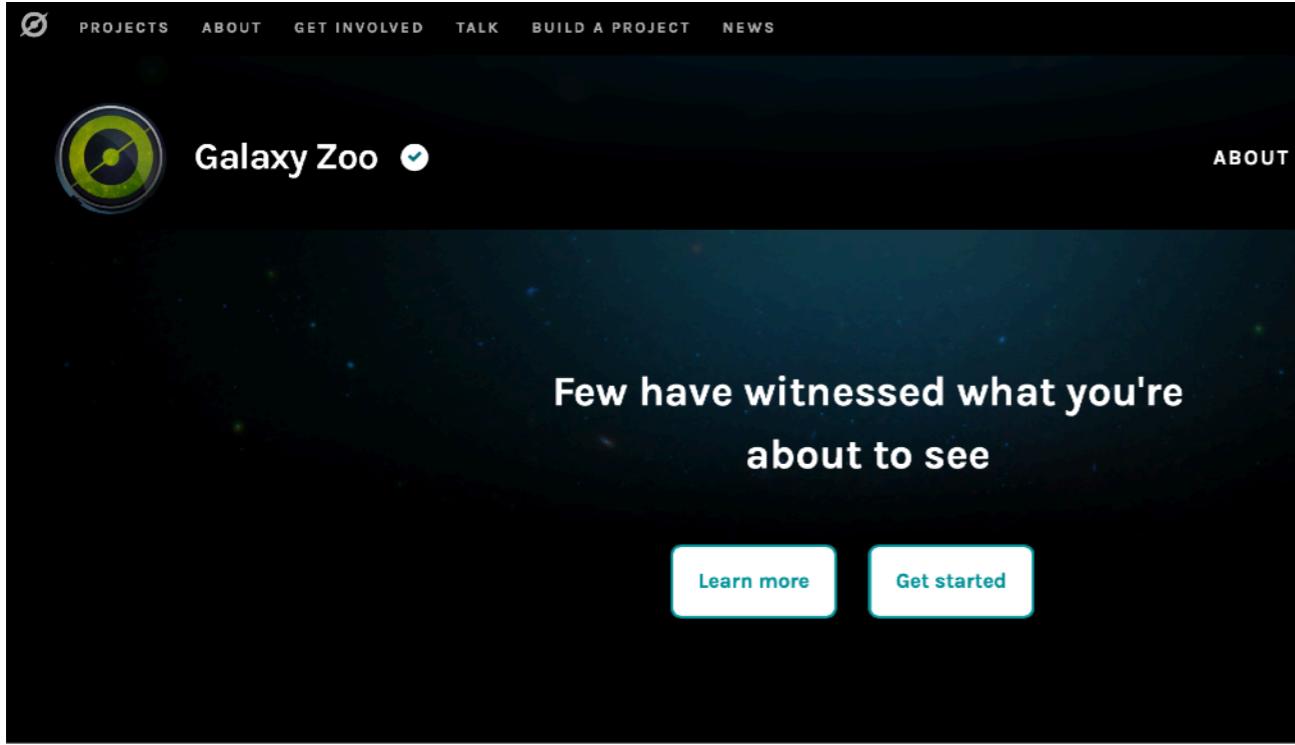


ready-made

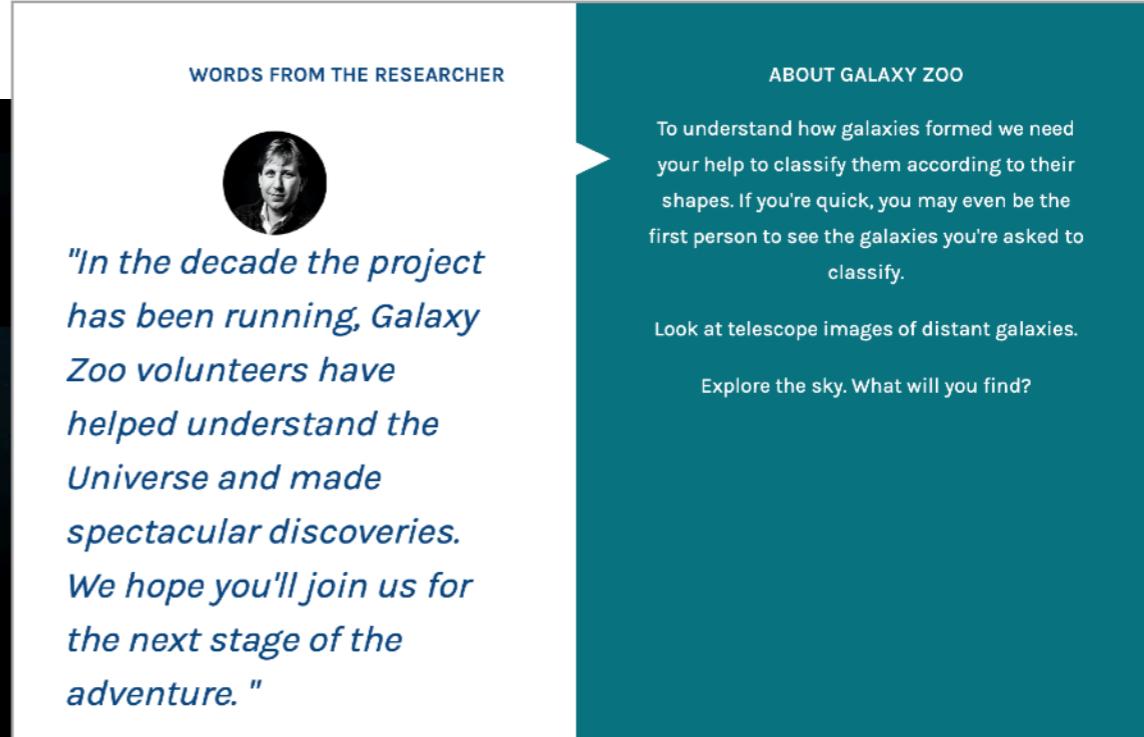
1 // EXKURS: READY UND CUSTOM MADE DATEN KOMBINIEREN.



1 // EXKURS: CUSTOM MADE DATEN IM GROSSEN STIL.



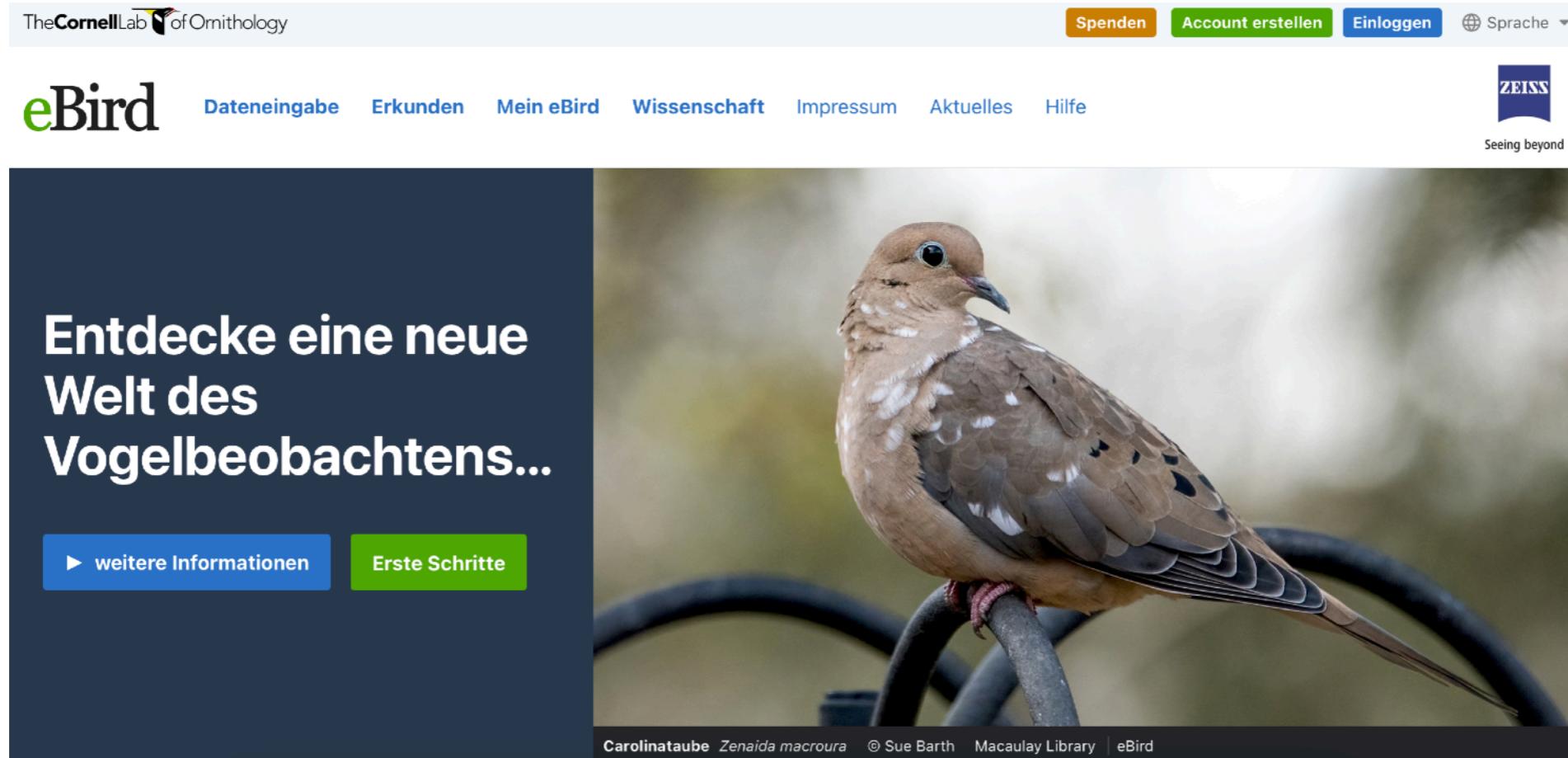
The screenshot shows the Galaxy Zoo project landing page. At the top, there's a navigation bar with links: PROJECTS, ABOUT, GET INVOLVED, TALK, BUILD A PROJECT, and NEWS. Below the navigation is the Galaxy Zoo logo, which features a stylized telescope icon. To the right of the logo is the word "ABOUT". In the center, a large text block reads: "Few have witnessed what you're about to see". Below this text are two buttons: "Learn more" and "Get started".



This screenshot shows a section of the Galaxy Zoo website. On the left, under "WORDS FROM THE RESEARCHER", there is a portrait of a man and a quote: "*In the decade the project has been running, Galaxy Zoo volunteers have helped understand the Universe and made spectacular discoveries. We hope you'll join us for the next stage of the adventure.*" On the right, under "ABOUT GALAXY ZOO", it says: "To understand how galaxies formed we need your help to classify them according to their shapes. If you're quick, you may even be the first person to see the galaxies you're asked to classify." Below this, there are two smaller text blocks: "Look at telescope images of distant galaxies." and "Explore the sky. What will you find?"



1 // EXKURS: READY MADE DATA NUTZBAR MACHEN.



The Cornell Lab of Ornithology

Spenden Account erstellen Einloggen Sprache ▾

eBird Dateneingabe Erkunden Mein eBird Wissenschaft Impressum Aktuelles Hilfe

ZEISS Seeing beyond

Entdecke eine neue Welt des Vogelbeobachtens...

► weitere Informationen Erste Schritte

Carolinataube *Zenaida macroura* © Sue Barth Macaulay Library | eBird

1 // EXKURS: WIE VIELE DATEN BRAUCHE ICH EIGENTLICH?

Es kommt drauf an.

Auf das Modell. Auf das zu lösende Problem. Auf die Qualität der Daten.

Mehr Parameter, mehr Daten.

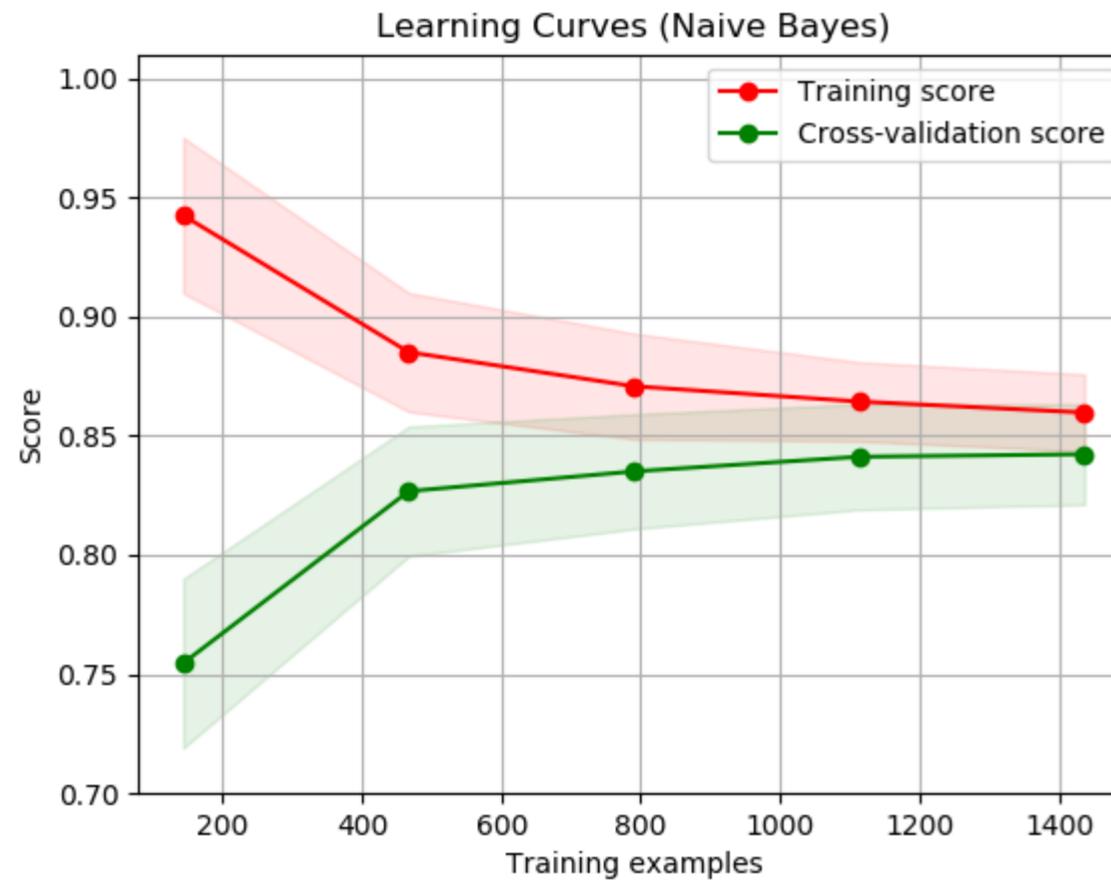
Mehr Klassen vorherzusagen, mehr Daten.

Saubere Daten, weniger Daten. Ungenaue Daten, mehr Daten.

signal-to-noise-ratio

Die richtigen Daten, weniger Daten. Die falschen Daten, mehr Daten.

1 // EXKURS: WIE VIELE DATEN BRAUCHE ICH EIGENTLICH?



Sich an ähnlichen Problemen in der Literatur orientieren und die Lernkurve (**Learning Curve**) analysieren.

2.A // DATA CLEANING

Randomisieren.

Deduplizieren.

Normalisieren.

Imputieren.

Data Augmentation.

Dimensionsreduktion.

Feature Engineering.

2.B // EXPLORATORY DATA ANALYSIS

Datenvisualisierung.

Deskriptive Statistiken.

Plausibilitätsprüfung.

Muster erkunden.

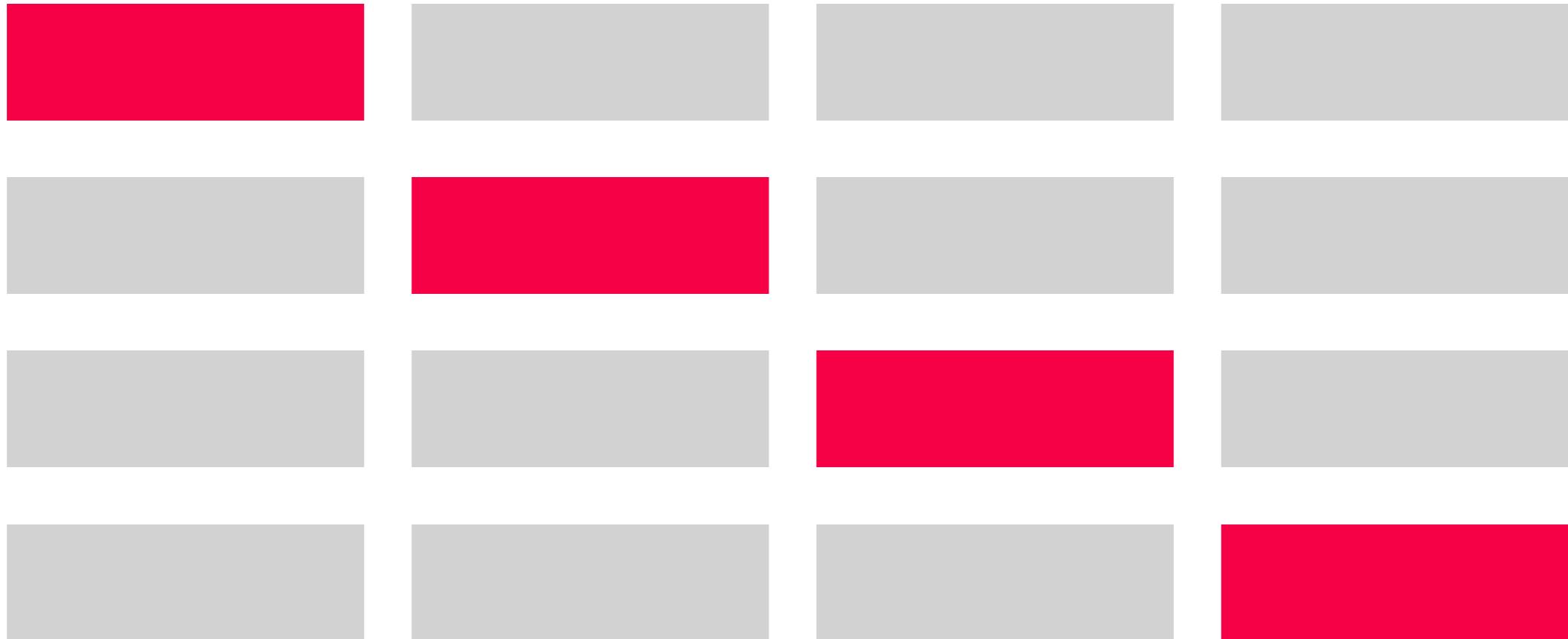
VOR DER MODELLANPASSUNG: DIE CROSS-VALIDIERUNG VORBEREITEN

Training Set: wird verwendet, um die Modellparameter zu schätzen

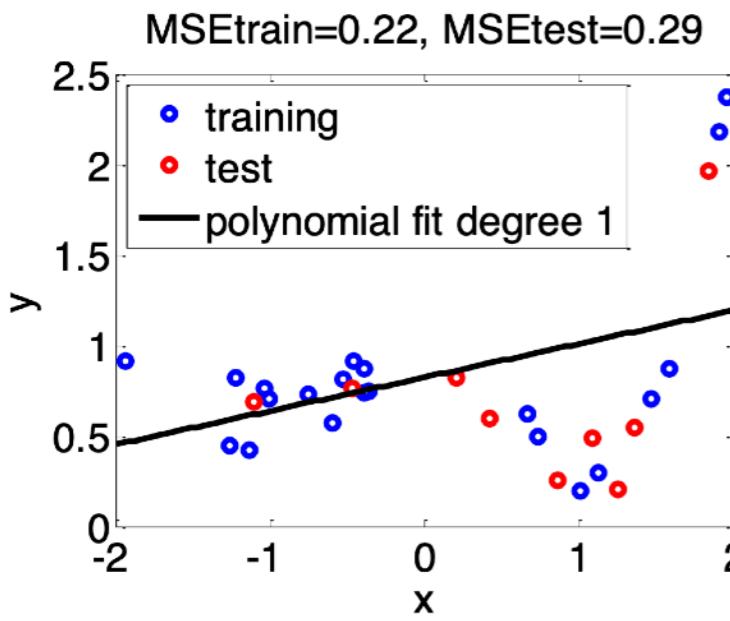
Validation Set: wird verwendet, um das Modell bereits in der Trainingsphase auf dessen Güte zu testen; das Modell mit den besten Evaluierungsergebnissen im Validation Set wird ausgewählt.

Test Set: wird verwendet, um die tatsächliche Güte des Modells auf ungesehenen Daten zu testen

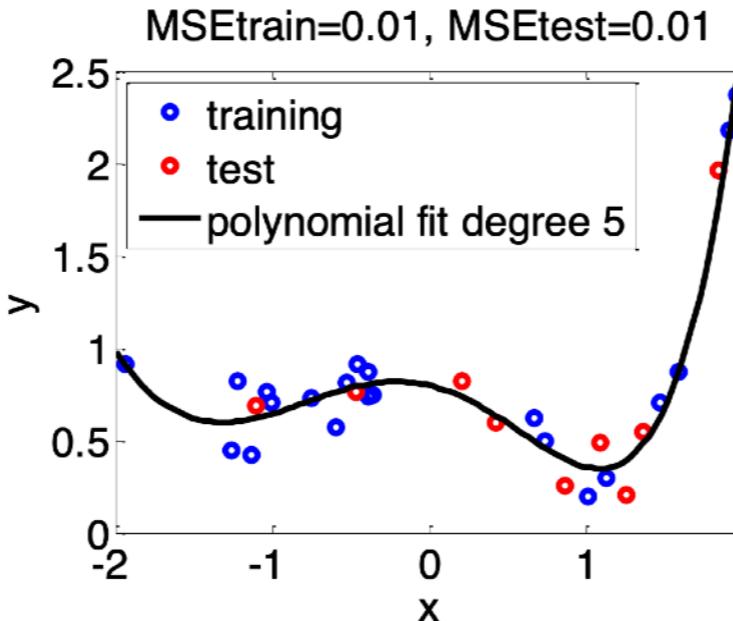
VOR DER MODELLANPASSUNG: DIE CROSS-VALIDIERUNG VORBEREITEN



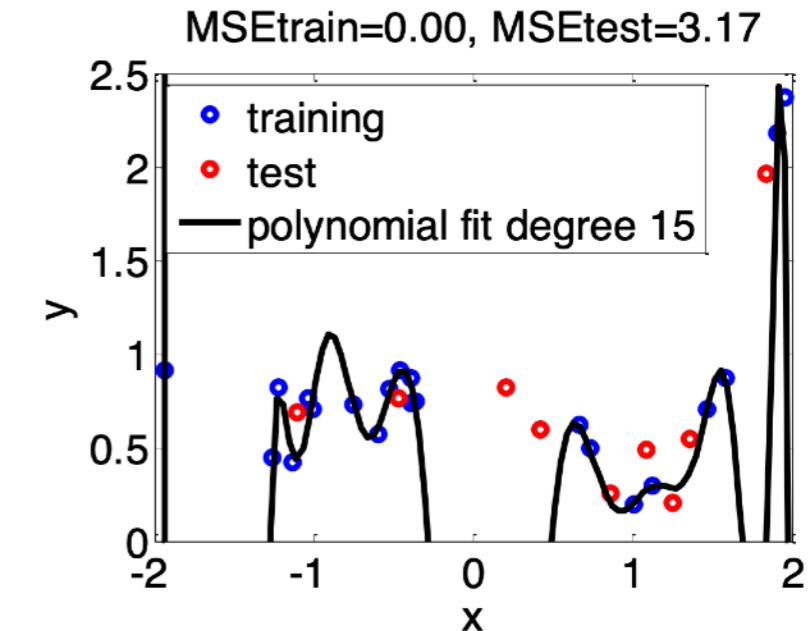
DIE CROSS-VALIDIERUNG VERHINDERT OVER- UND UNDERFITTING.



underfitting

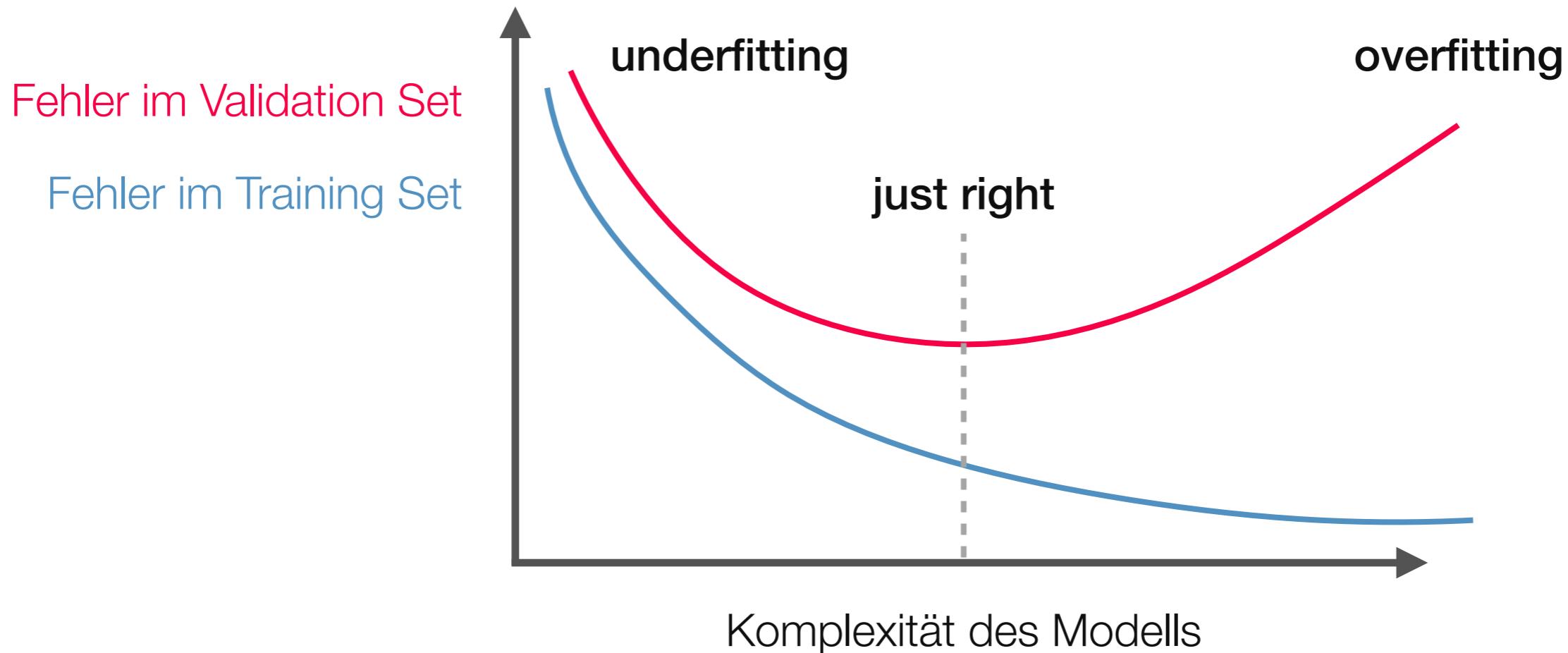


„just right“



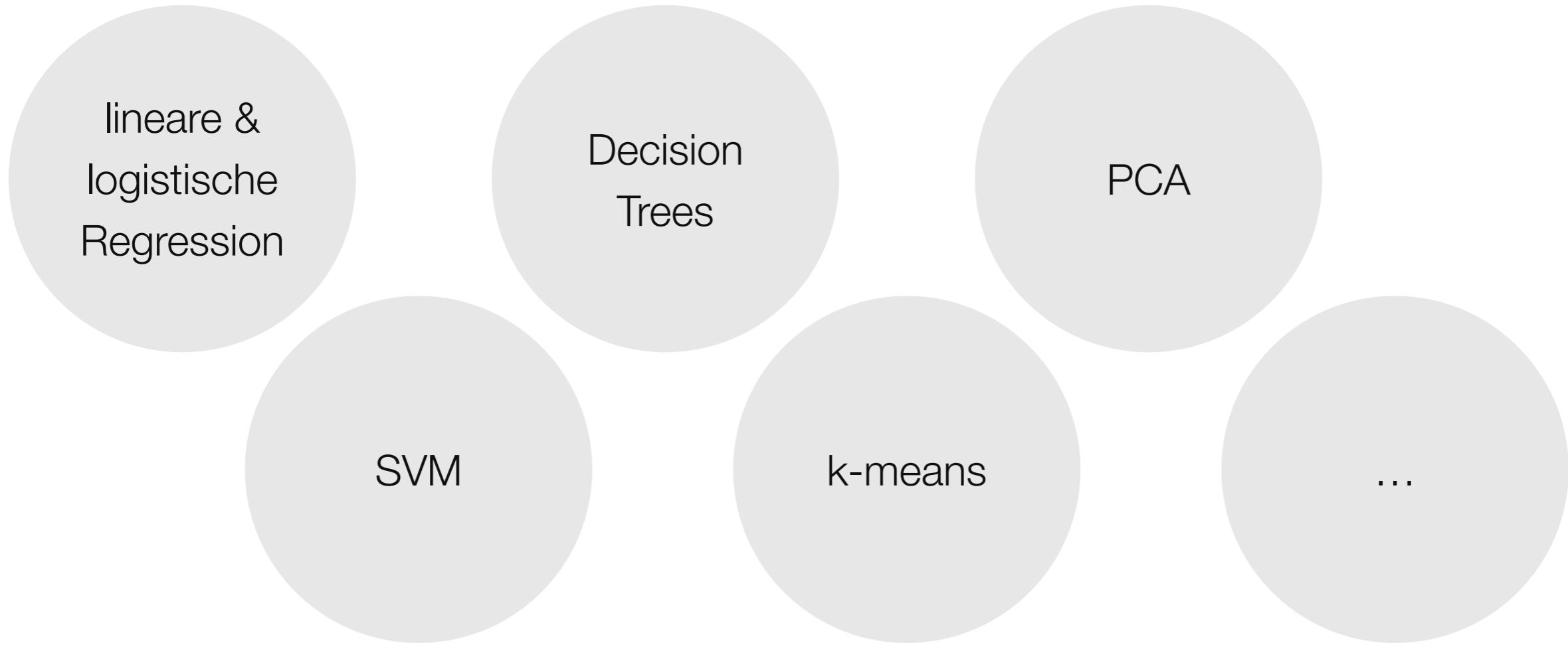
overfitting

DIE CROSS-VALIDIERUNG VERHINDERT OVER- UND UNDERFITTING.



FAUSTREGEL:
80% VORBEREITUNG
20% MODELLANPASSUNG

3 // MODEL FITTING



4 // MODEL EVALUATION: FÜR REGRESSION

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

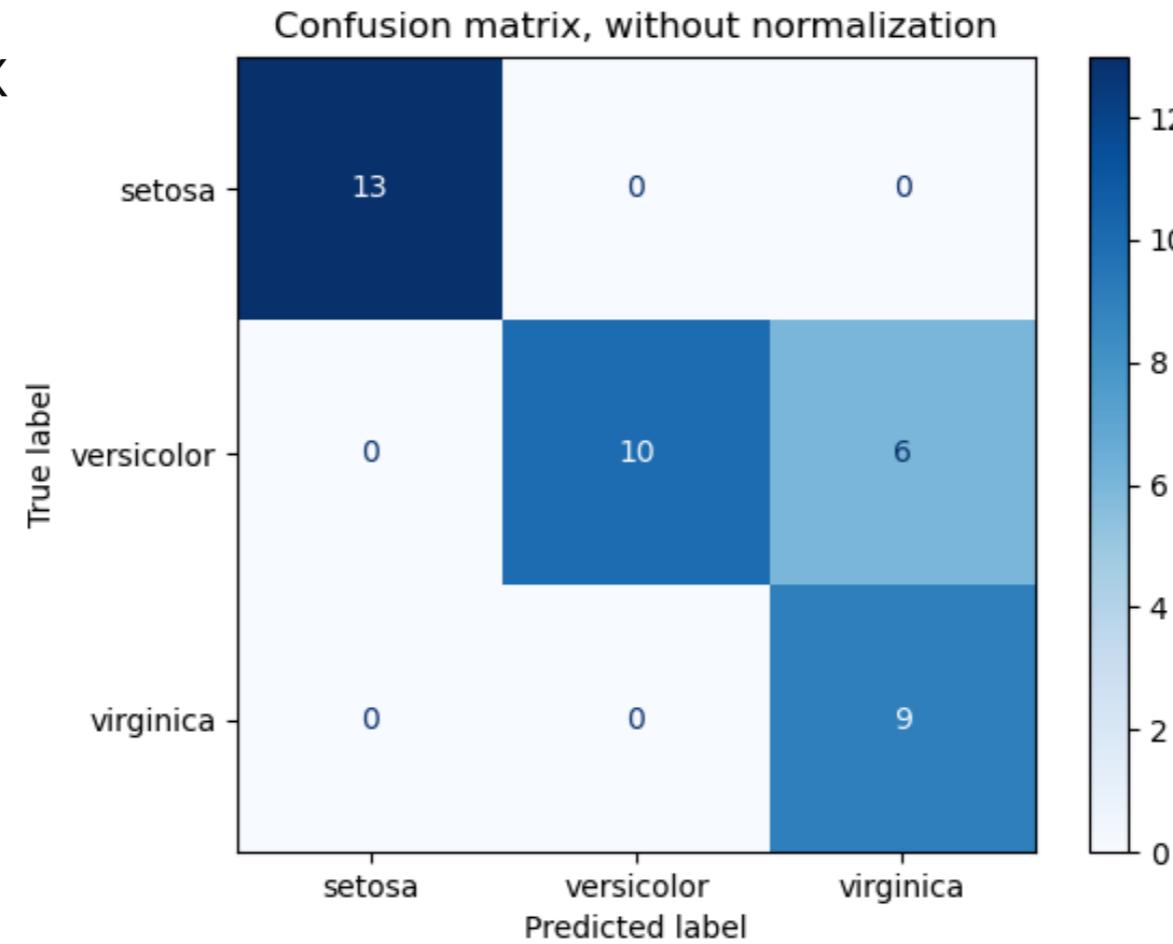
Mean Squared Error

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Coefficient of
Determination

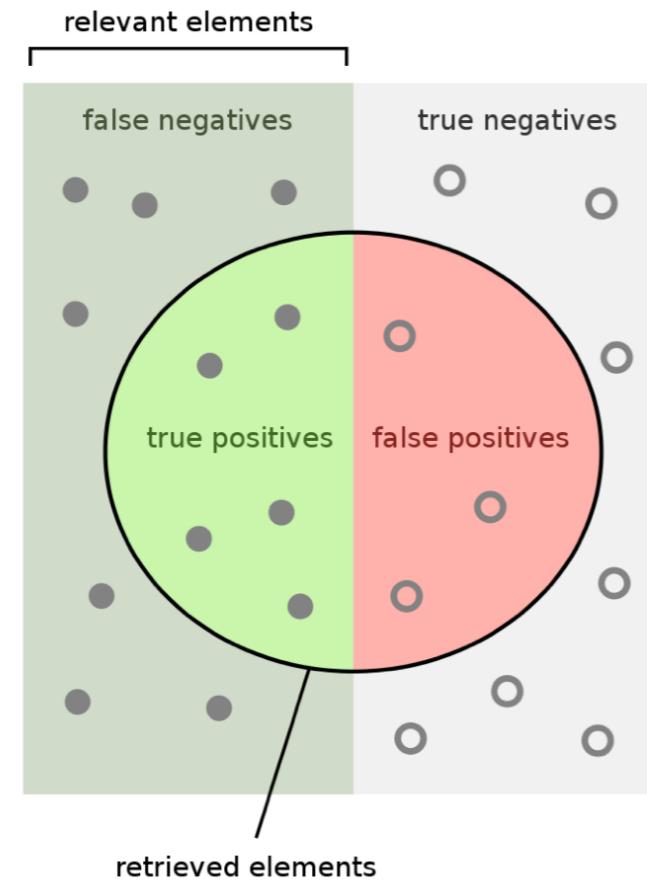
4 // MODEL EVALUATION: FÜR KLASSIFIZIERUNG

Die Confusion Matrix



4 // MODEL EVALUATION: FÜR KLASIFIZIERUNG

Der F1-Score



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

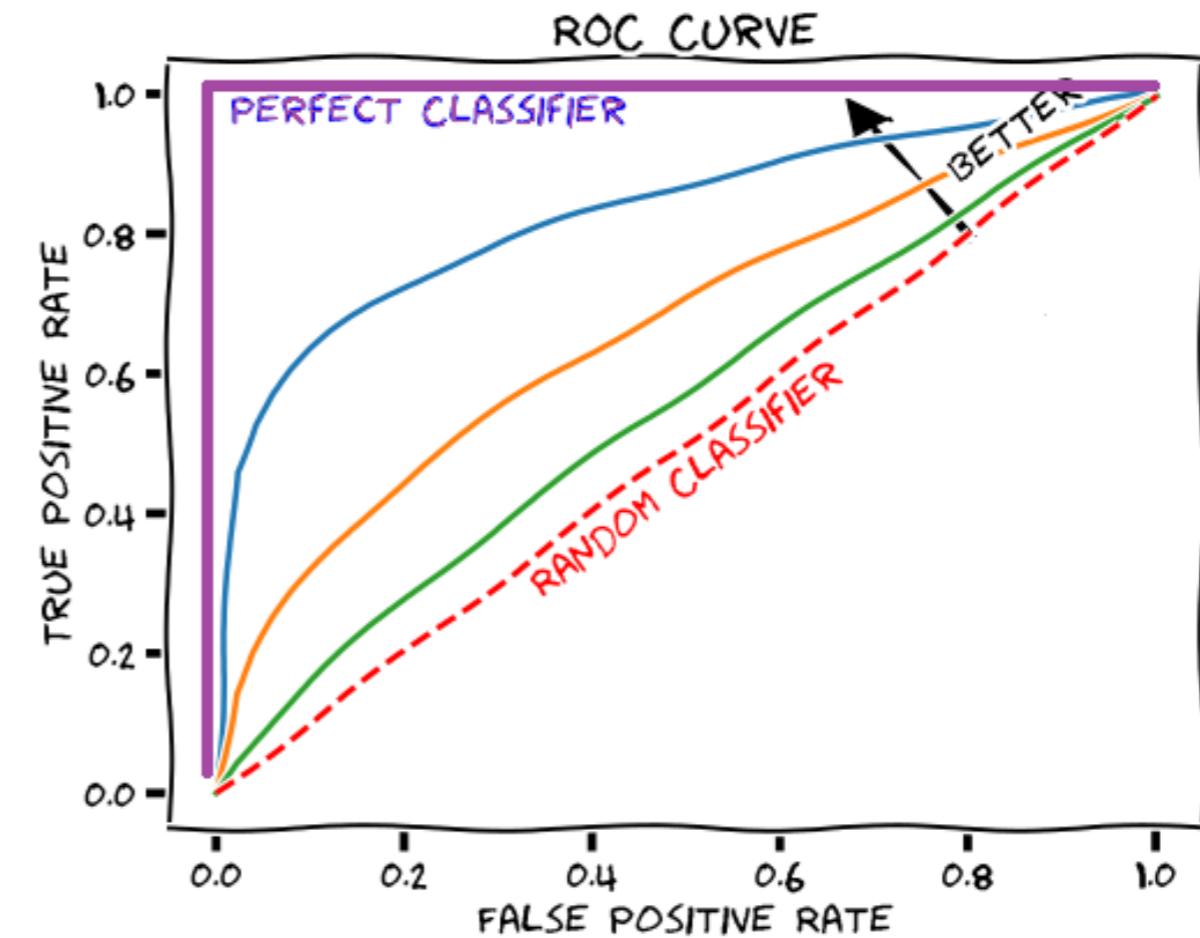
$$F1\ Score = 2 \times \frac{\text{recall} \times \text{precision}}{\text{recall} + \text{precision}}$$

4 // MODEL EVALUATION: FÜR KLASIFIZIERUNG

ROC-AUC Kurve

Receiver Operating Characteristics

Area Under the Curve



5 // FINE TUNING

Die Hyperparameter des Modells anpassen.
(z.B. die Anzahl der Trainingsschritte oder die Learning Rate)

TRIAL AND ERROR!

X. // BENCHMARKING

SVMs

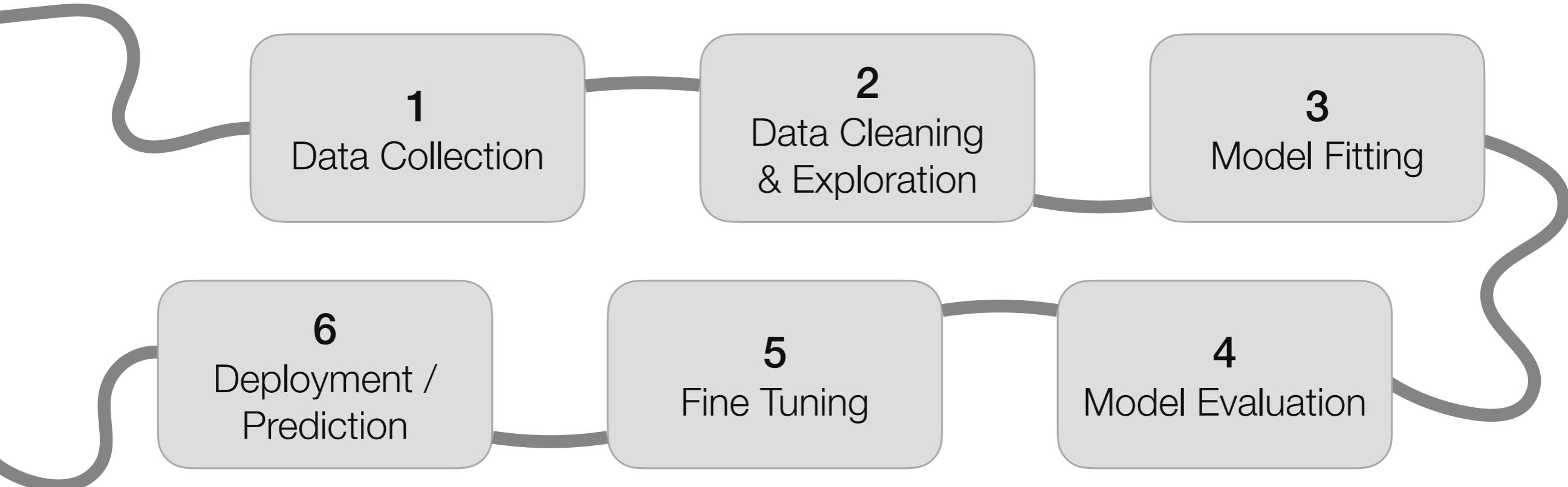


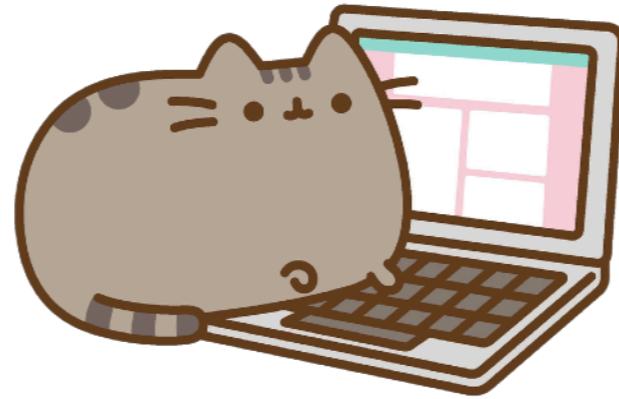
eine tolle neue Neural
Network Architektur

6 // PREDICTION



DIE DATA SCIENCE PIPELINE UMFASST (UNGEFÄHR) SECHS SCHritte:





THANK YOU!