# AIBridge

## Lecture 7

Quality

Acidity

Which one is a better line?

this model has more **predictive power**

Quality

Acidity

this model is highly accurate on **training data**

but bad at predictions anywhere else

**Quality** / **Acidity**

overfit!

■ too-precise fits to original data without generalization is called **overfitting**
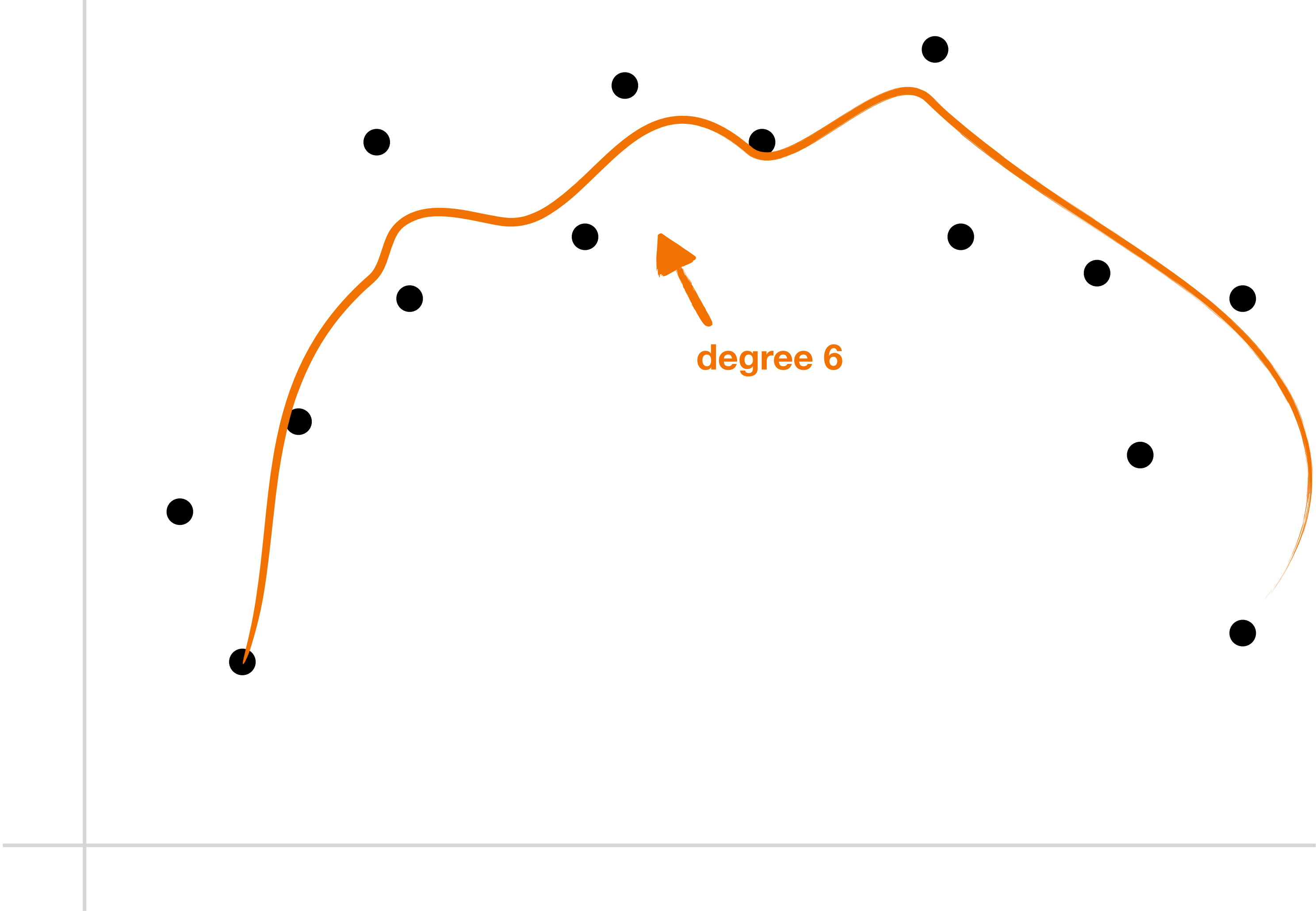
# overfitting



**underfit!**

**degree 1**

■ model is unable to capture relationship between variables

# overfitting



degree 2

# overfitting



degree 6

# overfitting
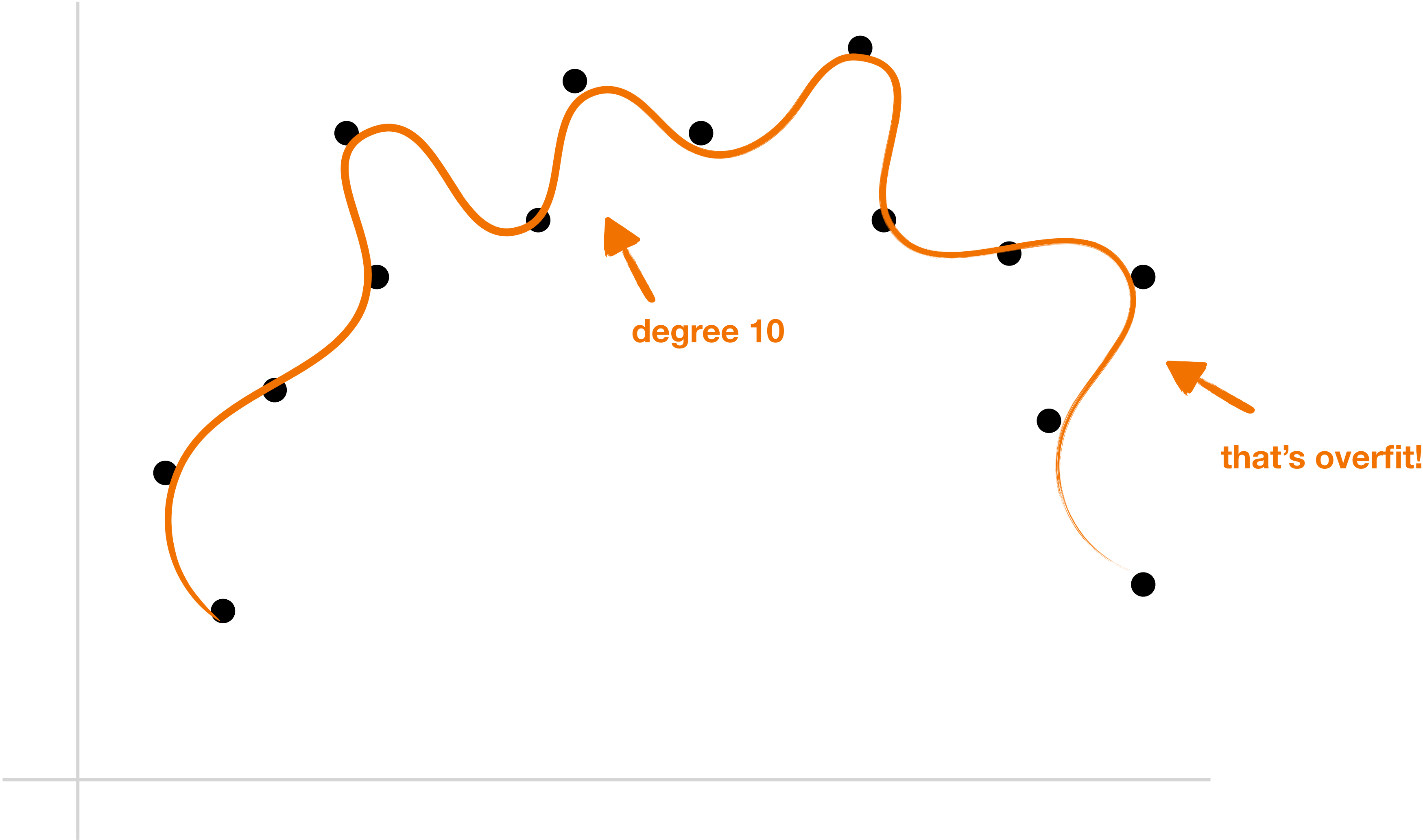


degree 10
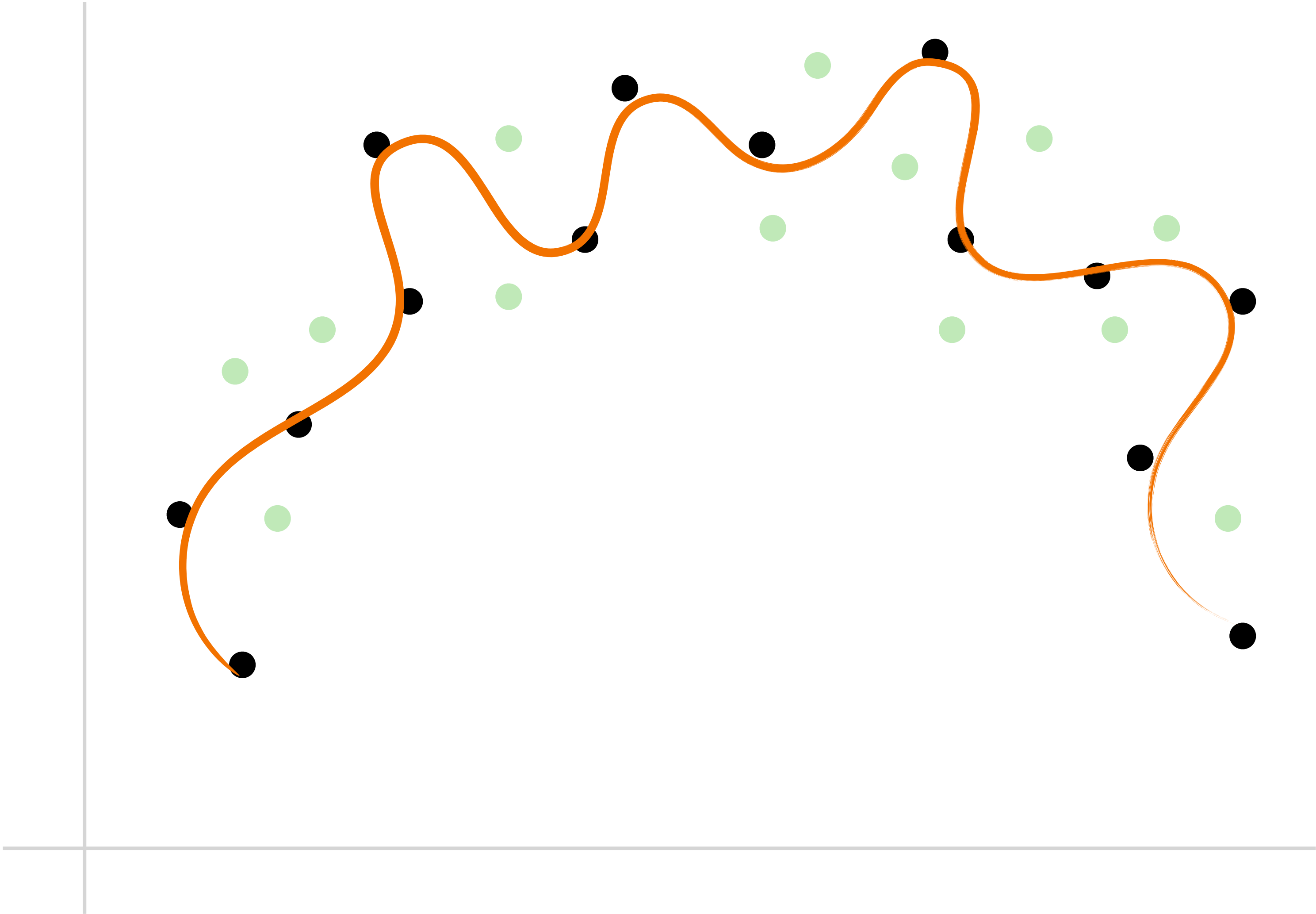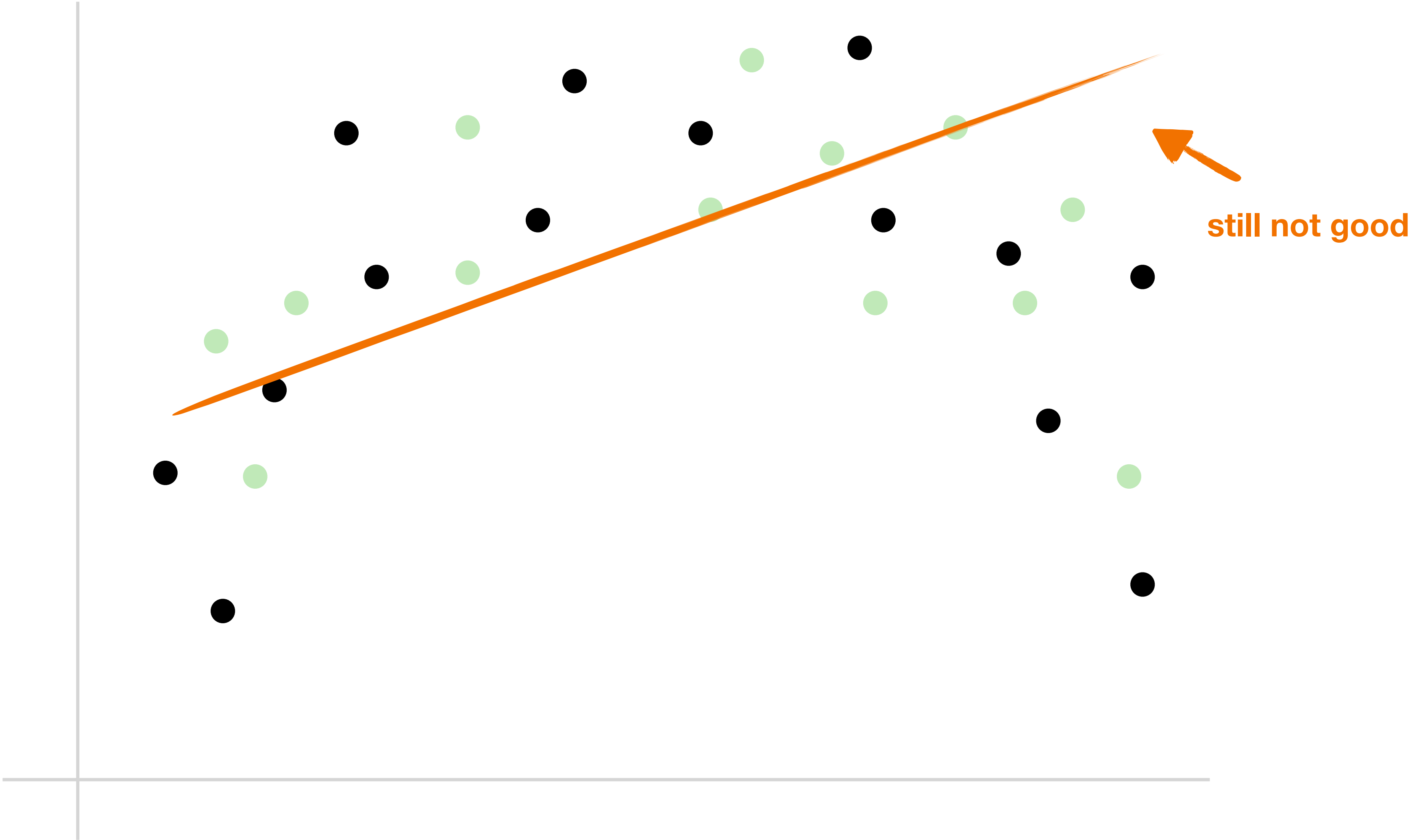
that's overfit!

# overfitting

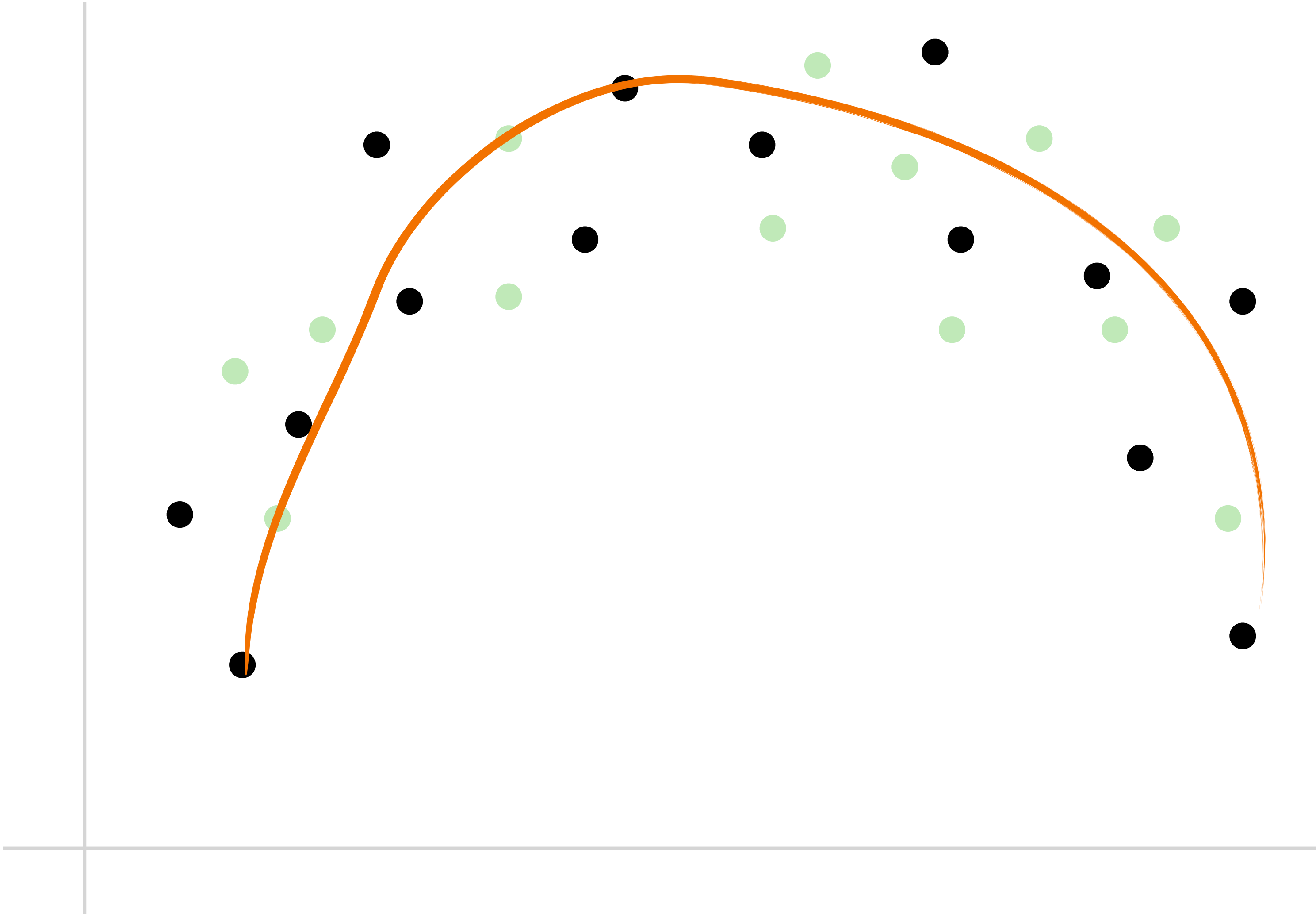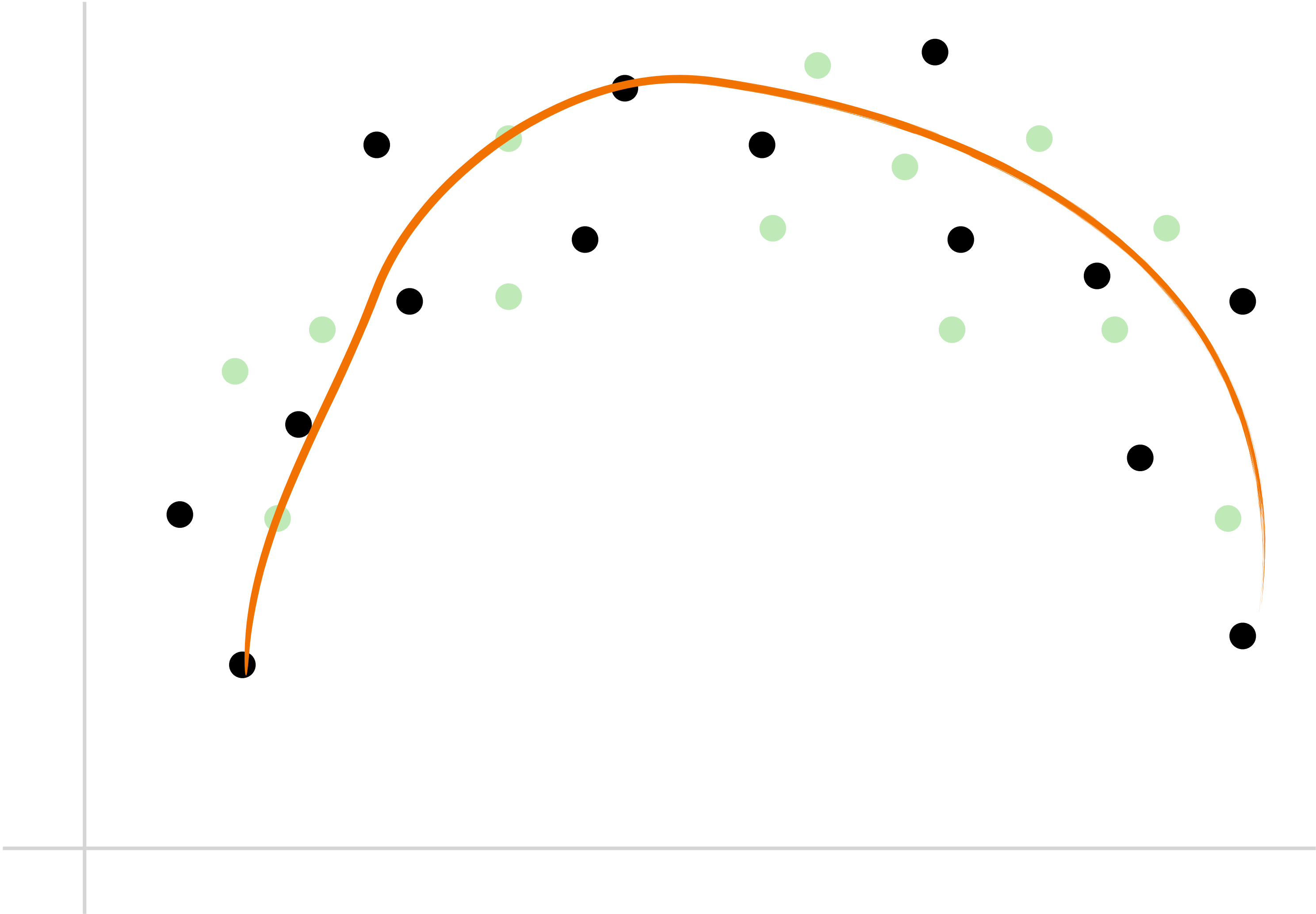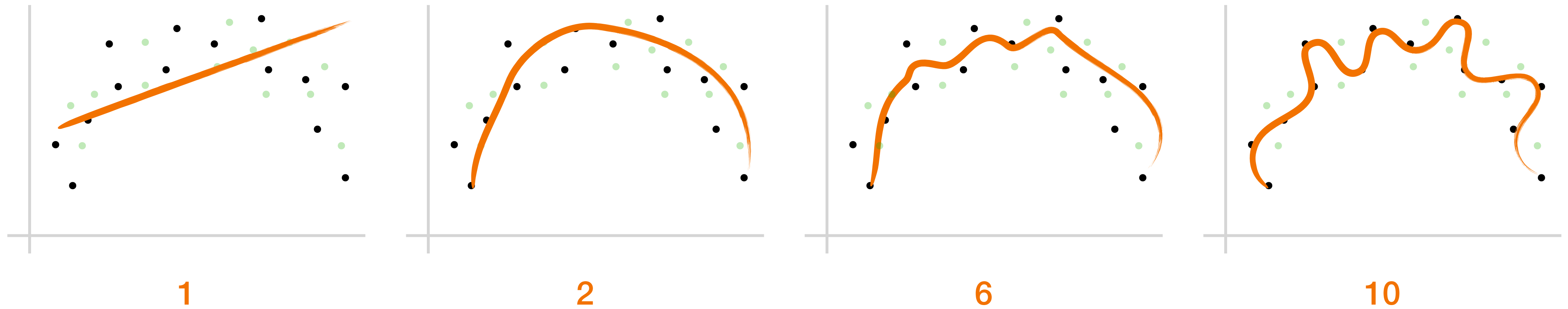# overfitting



still not good

# overfitting

# overfitting

# overfitting



- **overfitting** frequently takes place when the degree of a regression model is set too high

# How do we address overfitting?

# address overfitting

training data

# address overfitting

**training data**

**validation data**

**test data**

# address overfitting

**Model Has Seen**

**Model Hasn't Seen**

**training data**

**validation data**

**test data**

■ we use **validation** and **test** sets, small subsets of data the model hasn't seen before,

**val point(s)!**

# address overfitting

Model Has Seen

Model
Hasn't Seen

training data

validation
data

test data

wait but what's
the difference?

**Whole Dataset**

**test data**

standardized for benchmarking!

- **test sets** are, unlike validation sets, usually set by the data creator as common, unseen benchmark data.
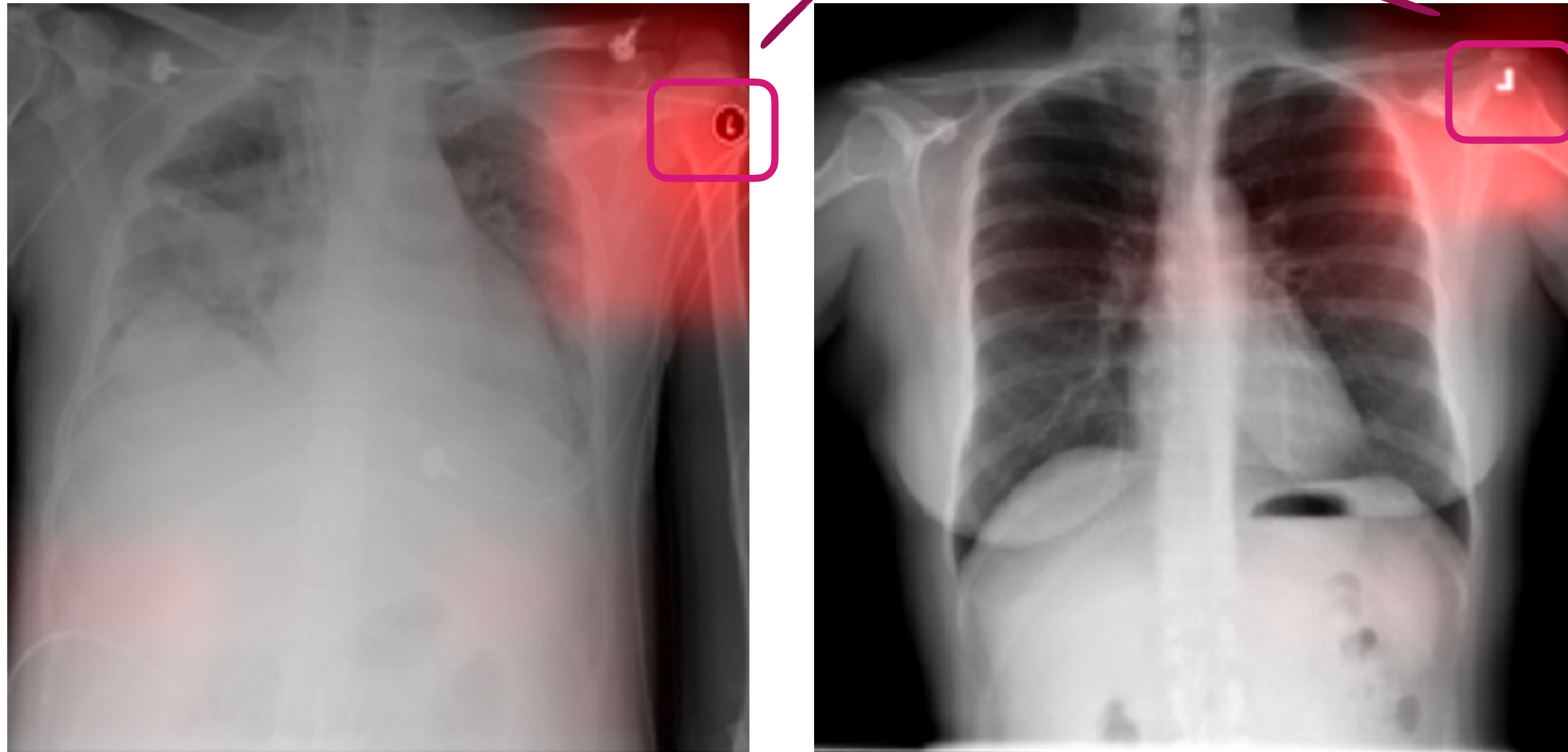
**overfitting** can be dangerous

data ethics

# data ethics



**which one has pneumonia?**

# data ethics



different hospitals used
different markers!

- models, when not controlled for external factors, often **overfit** on easy targets

# Onto Feature Selection!