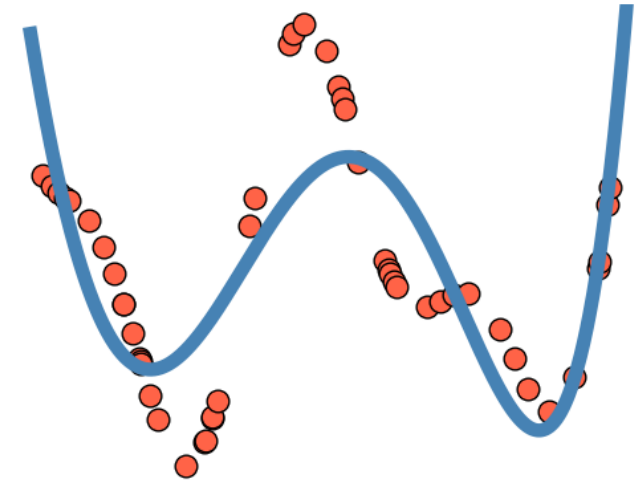


Mathematics of Data Science



University of Vienna, WiSe 2025
Master's programme in Data Science

Motivation

So far, we did not discuss “labelled” data!

Meaning: we looked at $p(\mathbf{x})$, but not $p(\mathbf{y} \mid \mathbf{x})$

Classification

\mathbf{x}

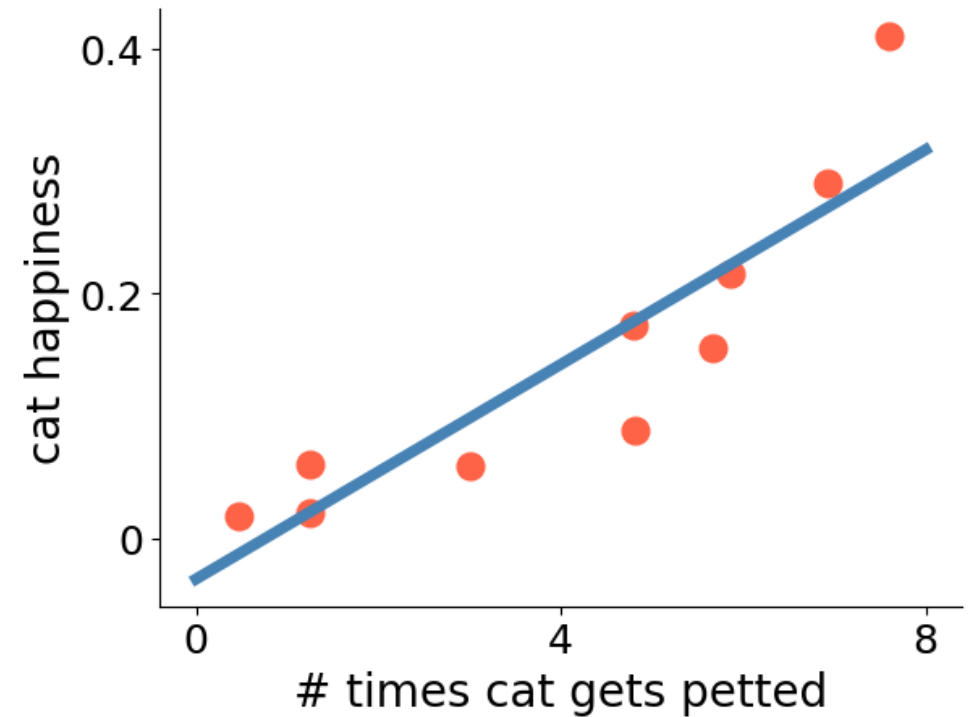


\mathbf{y}

Cat

Dog

Regression



Content

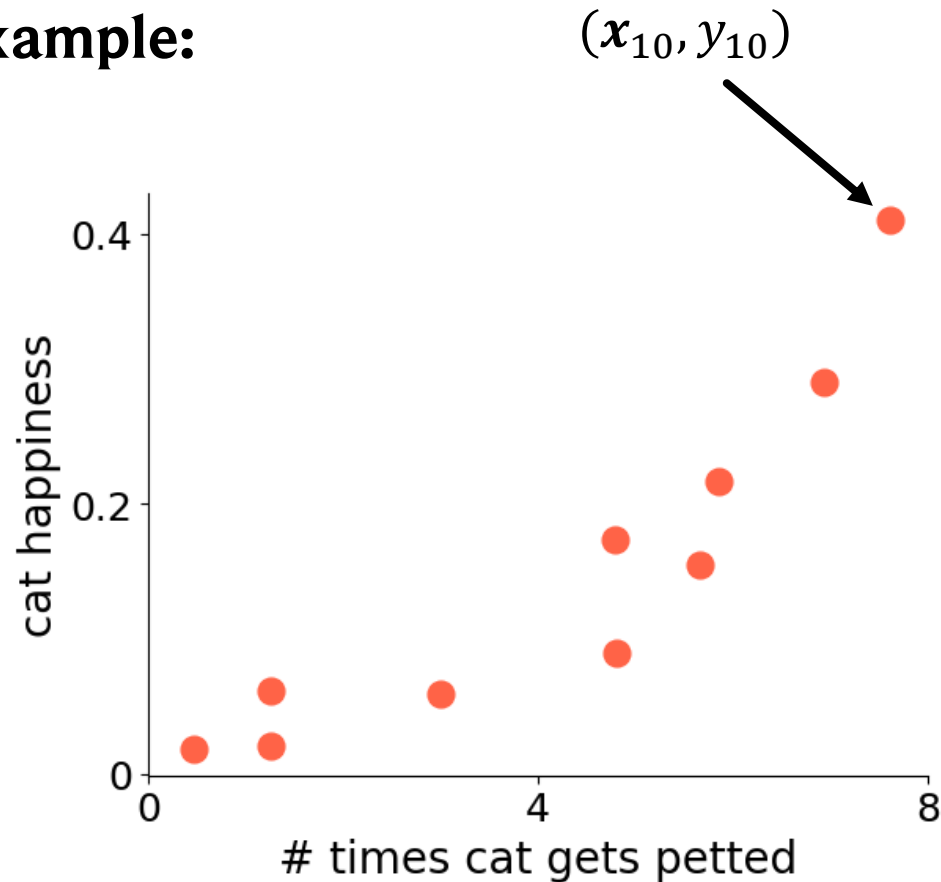
- Hypothesis set and risk
- Linear Regression
- Polynomials
- Universal function approximators
- Weierstrass approximation theorem
- Stone-Weierstrass theorem

Fitting data: the hypothesis set

What is our usual setup?

We have n data samples $(\mathbf{x}_i, \mathbf{y}_i)$, where \mathbf{x}_i is a feature vector and \mathbf{y}_i whatever we want to predict!

Example:



Task: For new, unlabelled data samples, we want to predict \mathbf{y} !

We cannot try to fit all possible functions, so we restrict ourselves to a **hypothesis set**:

\mathcal{H} = set of functions to explain the data

For example: $\mathcal{H} = \{x, 2 \cdot x^2\}$

We can also have a whole span of functions:

$$\mathcal{H} = \{\alpha + \beta x + \gamma x^2 \mid \alpha, \beta, \gamma \in \mathbb{R}\}$$

Empirical risk minimization

Callback to the random number lectures:

we assume that our data is sampled from an underlying distribution $D = p(\mathbf{y}, \mathbf{x})$.

Let's use a loss function \mathcal{L} to measure how well our model performs (e.g., $\mathcal{L}(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|^2$).

Risk: For any function $h \in \mathcal{H}$, its risk $R(h)$ tells us how well the function fits the data

$$R(h) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim D} [\mathcal{L}(h(\mathbf{x}), \mathbf{y})]$$

Empirical risk: Assume we have m data samples. Then the empirical risk $\hat{R}(h)$ is

$$\hat{R}(h) = \sum_{i=1}^m \mathcal{L}(h(\mathbf{x}_i), \mathbf{y}_i)$$

Fitting the data: Find $\hat{h} \in \mathcal{H}$ that minimizes \hat{R} : $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} (\hat{R}(h))$

The risk of our empirical risk minimizer is bounded by two terms:

$$R(\hat{h}) \leq \sup_{h \in \mathcal{H}} |R(h) - \hat{R}(h)| + \inf_{g \in \mathcal{H}} \hat{R}(g)$$

*How well does our fit work
on the full data distribution?*

Generalization error

Interpolation error

*How well can we fit
the training data?*

Why the supremum?

$$R(\hat{h}) \leq \sup_{h \in \mathcal{H}} |R(h) - \hat{R}(h)| + \inf_{g \in \mathcal{H}} \hat{R}(g)$$

What is our aim?

1. Given a finite sample from a dataset, we want to minimize the risk (= solving our task).
2. Depending on the dataset (and training samples), **any of our functions** in \mathcal{H} could be the empirical risk minimizer!
3. We do **not** know which function in \mathcal{H} this is in advance!
4. Thus, the function \hat{h} minimizing the empirical risk **might be one with high generalization error** $|R(\hat{h}) - \hat{R}(\hat{h})|$.

To account for this, we use the supremum to **bound** the generalization error.

In other words: this is the **worst generalization error** we can get using a hypothesis in \mathcal{H} (again, empirical risk minimization might select exactly such a “bad” function, as any function in \mathcal{H} is allowed as a solution).

Note: this is a property of the whole set \mathcal{H} , i.e., shared by **all** functions in \mathcal{H} .

Also note: usually, there is a trade-off between generalization and interpolation error. However, both can be low if, for instance, our data is native to \mathcal{H} . E.g., if the underlying data relationship is linear, then linear functions will fit the training data well and generalize well to new data!

Example: linear regression

Let's look at a concrete example: linear regression (your baseline for everything!)

In this case, we have $\mathcal{H} = \{\mathbf{w}^T \mathbf{x} \mid \mathbf{w} \in \mathbb{R}^N\}$, the set of all affine functions.

Empirical risk minimization gives us:

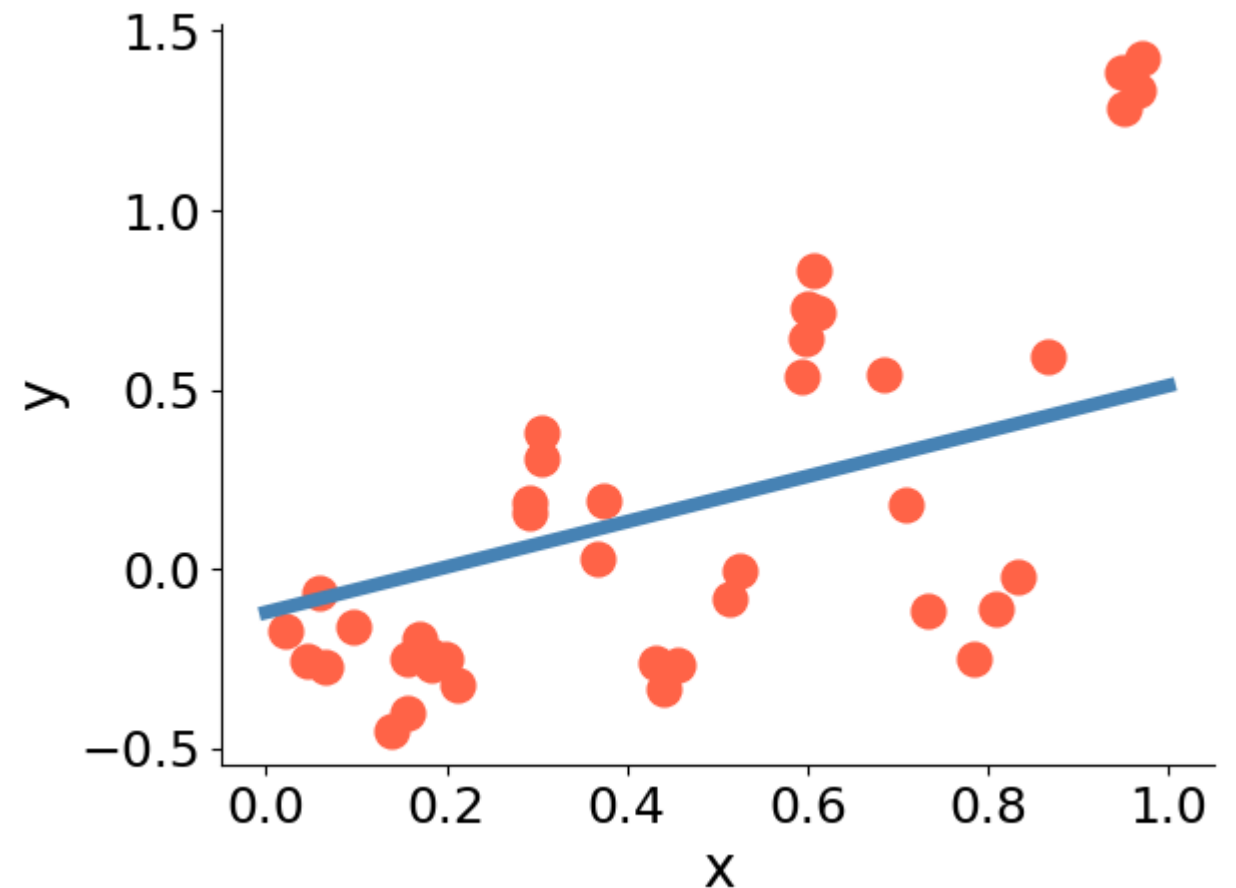
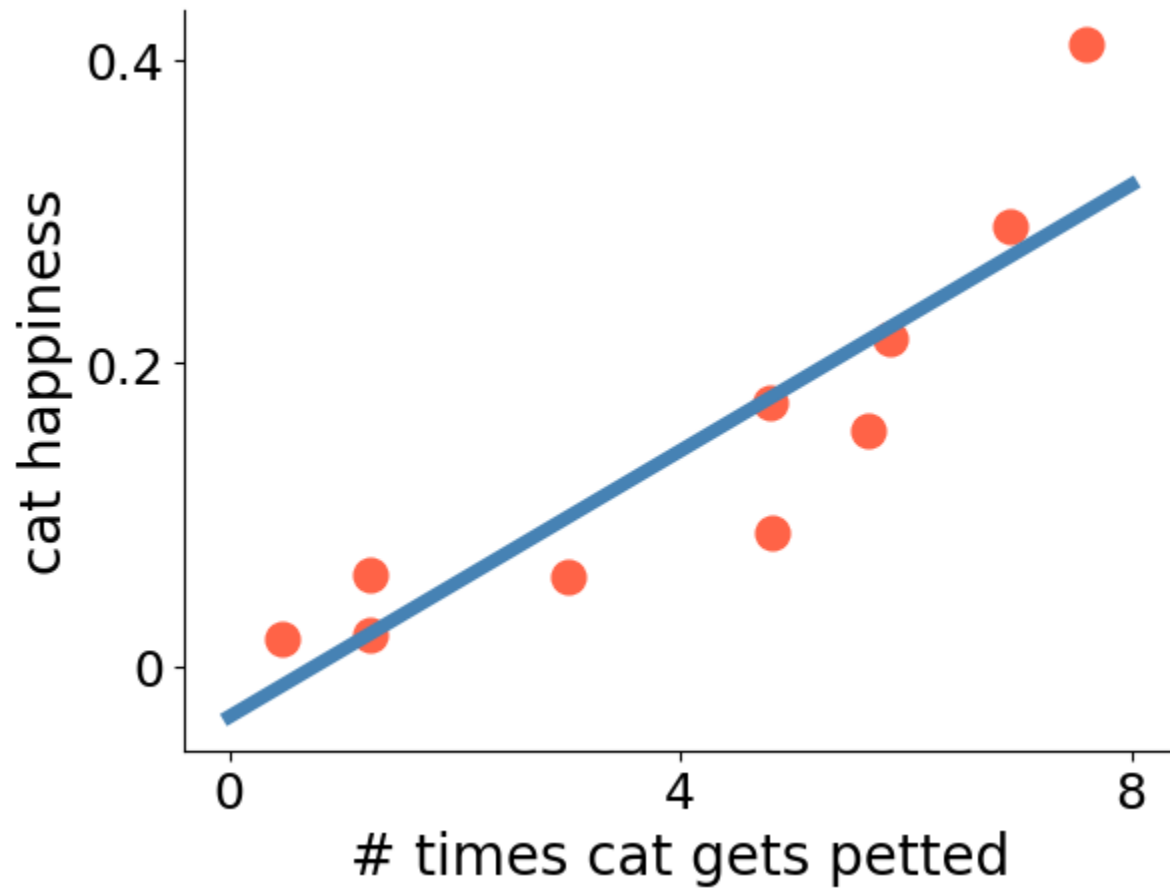
$$\inf_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \inf_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}_i - y_i) \right\|^2$$
$$= \inf_{\mathbf{w}} \|(\mathbf{X}\mathbf{w} - \mathbf{y})\|^2$$

rows of $\mathbf{X} \rightarrow \mathbf{x}_i$
elements of $\mathbf{y} \rightarrow y_i$

This is solved using the gradients:

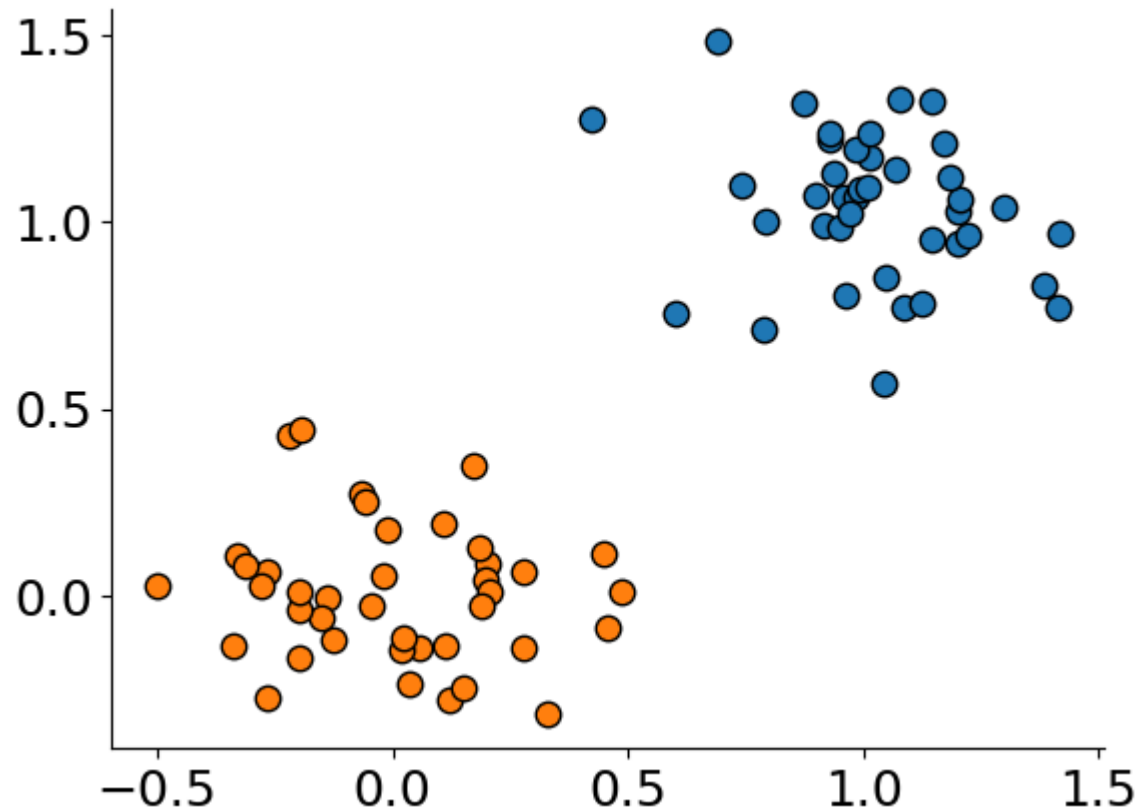
$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0 \quad \Leftrightarrow \quad \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Example: linear regression



What about classification?

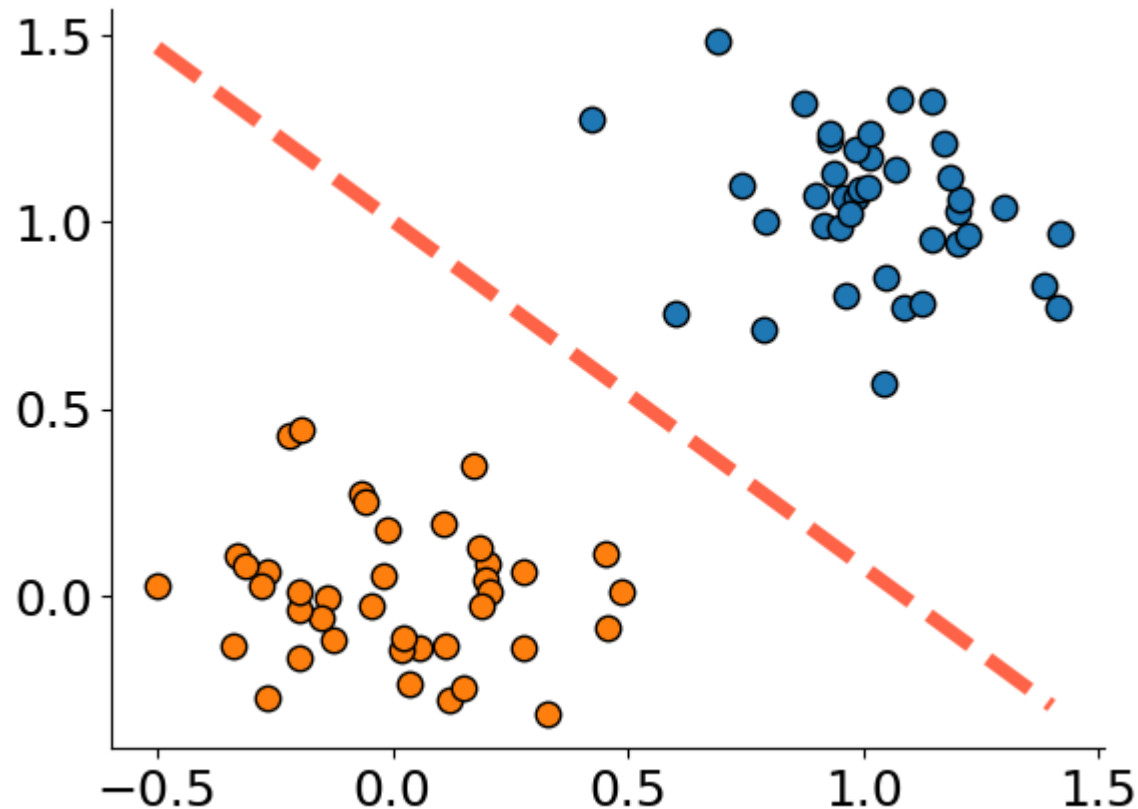
Very similar! In case of two classes, we model the probability of a data sample to belong to the first class by $p(\mathbf{x}) = \sigma(b + \mathbf{w}^T \mathbf{x})$ and to the second class $1 - p(\mathbf{x})$.



A simple decision rule is to assign class 1 if $p \geq 0.5$, and assign class 0 otherwise.

What about classification?

Very similar! In case of two classes, we model the probability of a data sample to belong to the first class by $p(\mathbf{x}) = \sigma(b + \mathbf{w}^T \mathbf{x})$ and to the second class $1 - p(\mathbf{x})$.



A simple decision rule is to assign class 1 if $p \geq 0.5$, and assign class 0 otherwise.

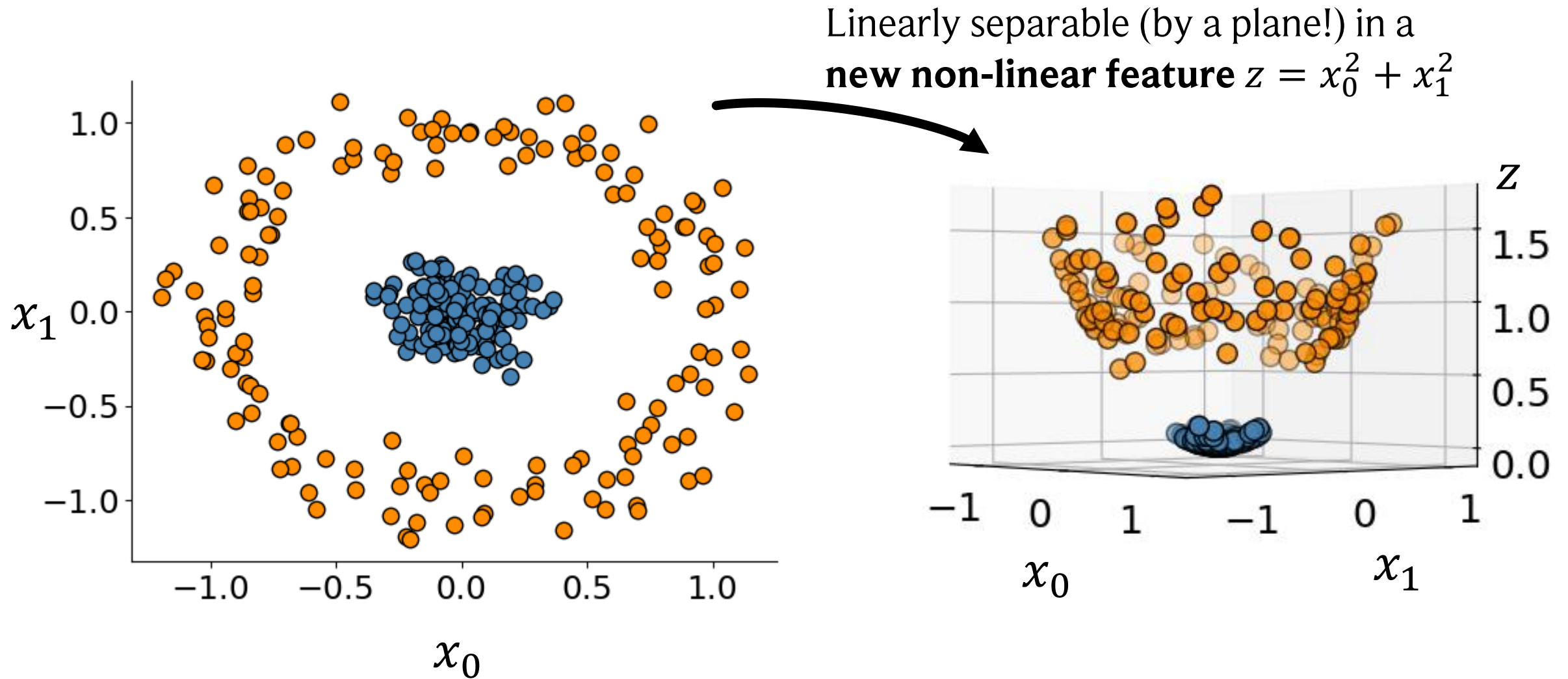
This is equivalent to:

$$\text{Class 1: } b + \mathbf{w}^T \mathbf{x} \geq 0$$

$$\text{Class 2: } b + \mathbf{w}^T \mathbf{x} < 0$$

The class is given depending on which side of the **hyperplane** \mathbf{w} the sample lies!

And... what now?!



Polynomial regression

Use **non-linear transformations** of the features: $F(x) = \mathbf{w}^T(1, x, x^2, x^3) = \underbrace{\sum_{i=0}^3 w_i x^i}_{\text{Third-degree polynomial}}$

For fitting data, using the basis functions $\{x^i \mid i \in \mathbb{N}\}$ is generally not recommended.
Instead, use an orthogonal set of basis polynomials!

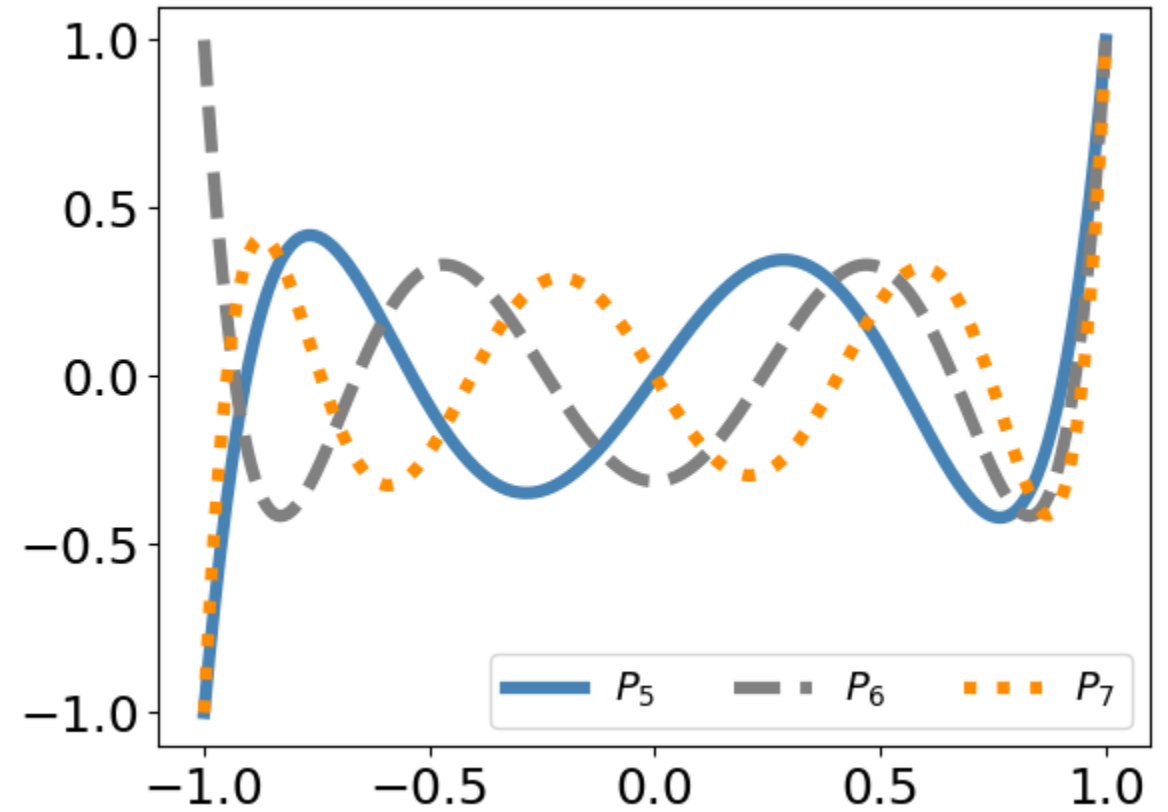
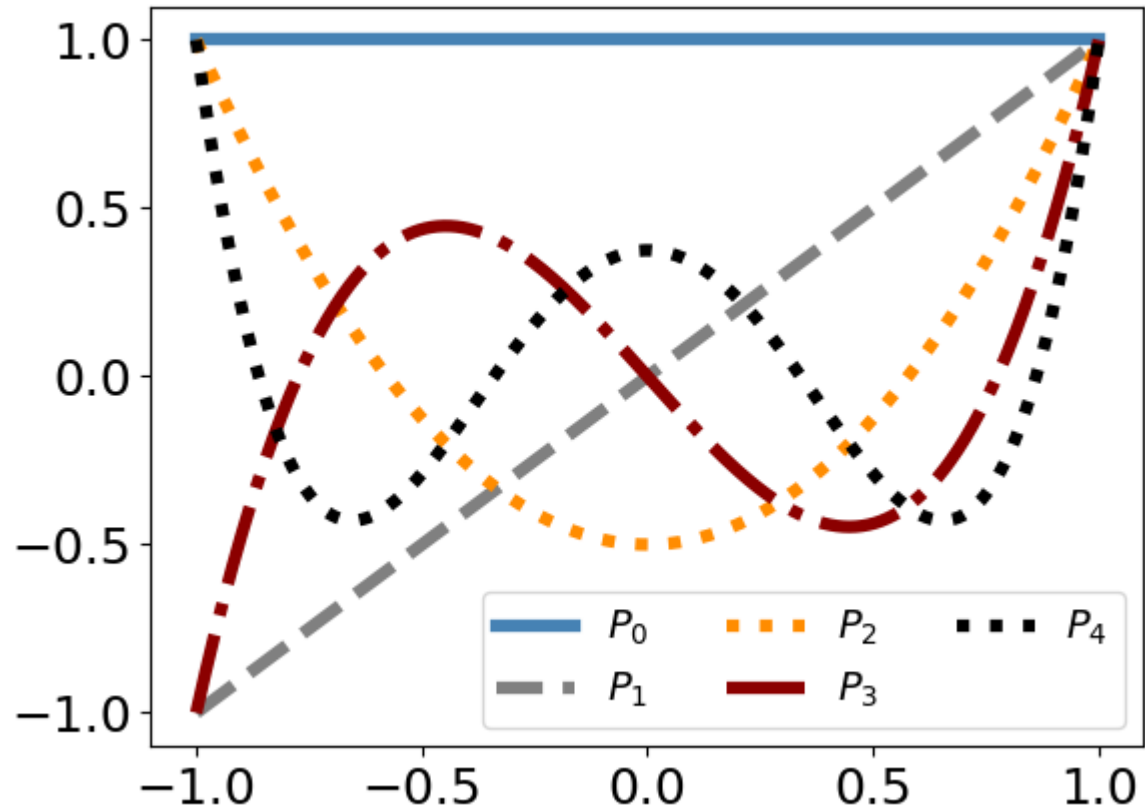
E.g., **Legendre polynomials**: $P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n$

$P_0 = 1$	$P_2 = \frac{1}{2}(3x^2 - 1)$
$P_1 = x$	$P_3 = \frac{1}{2}(5x^3 - 3x)$

These basis polynomials satisfy:

$$\int_{-1}^1 P_m(x) P_n(x) dx = 0 \text{ if } n \neq m$$

Some example Legendre polynomials



Fitting data using polynomials

This is exactly like linear regression, just with different “data features”.

Assume the case of N one-dimensional data features $x \in [-1, 1]$. Then our regression model is:

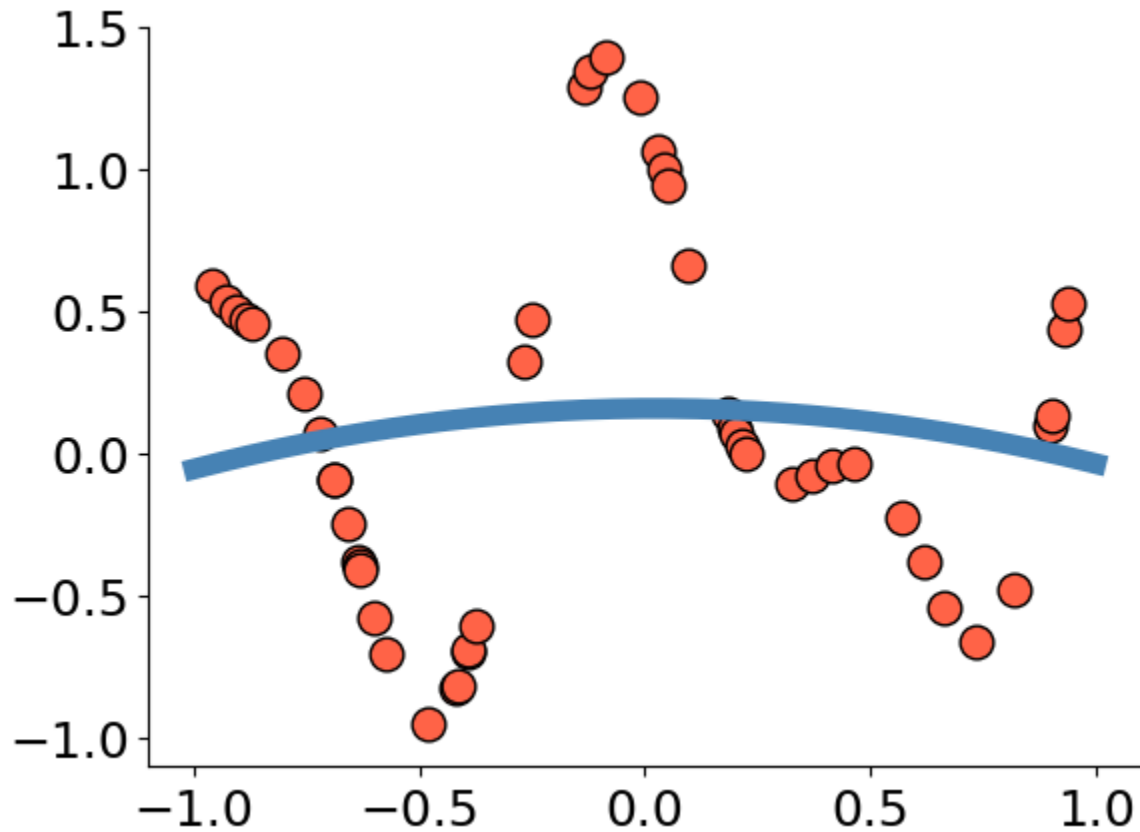
$$\begin{pmatrix} P_0(x_1) & \cdots & P_d(x_1) \\ \vdots & \ddots & \vdots \\ P_0(x_N) & \cdots & P_d(x_N) \end{pmatrix} \cdot \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} = \mathbf{X}\mathbf{w} \quad \text{or} \quad (\mathbf{X}\mathbf{w})_k = \sum_{i=0}^d w_i P_i(x_k)$$

where d is the degree of the polynomial we use. We can find a solution again using:

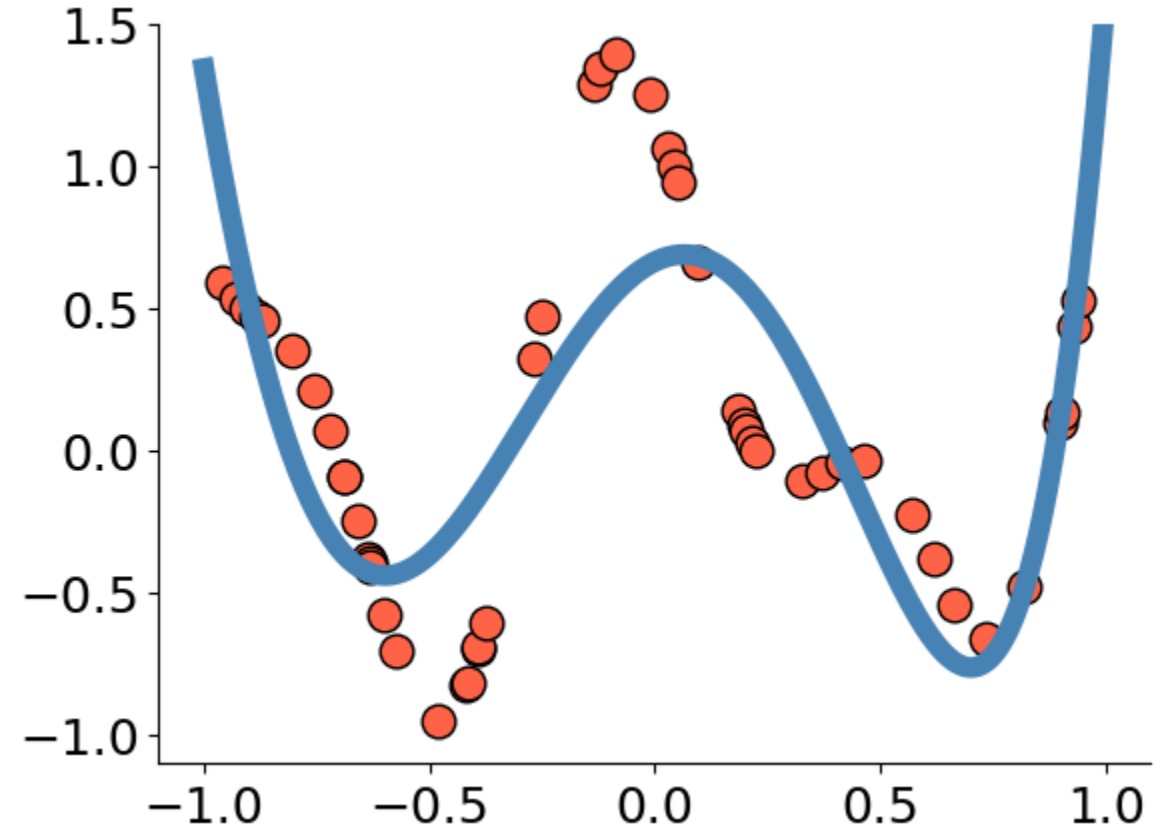
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Results for different degrees

$d = 2$

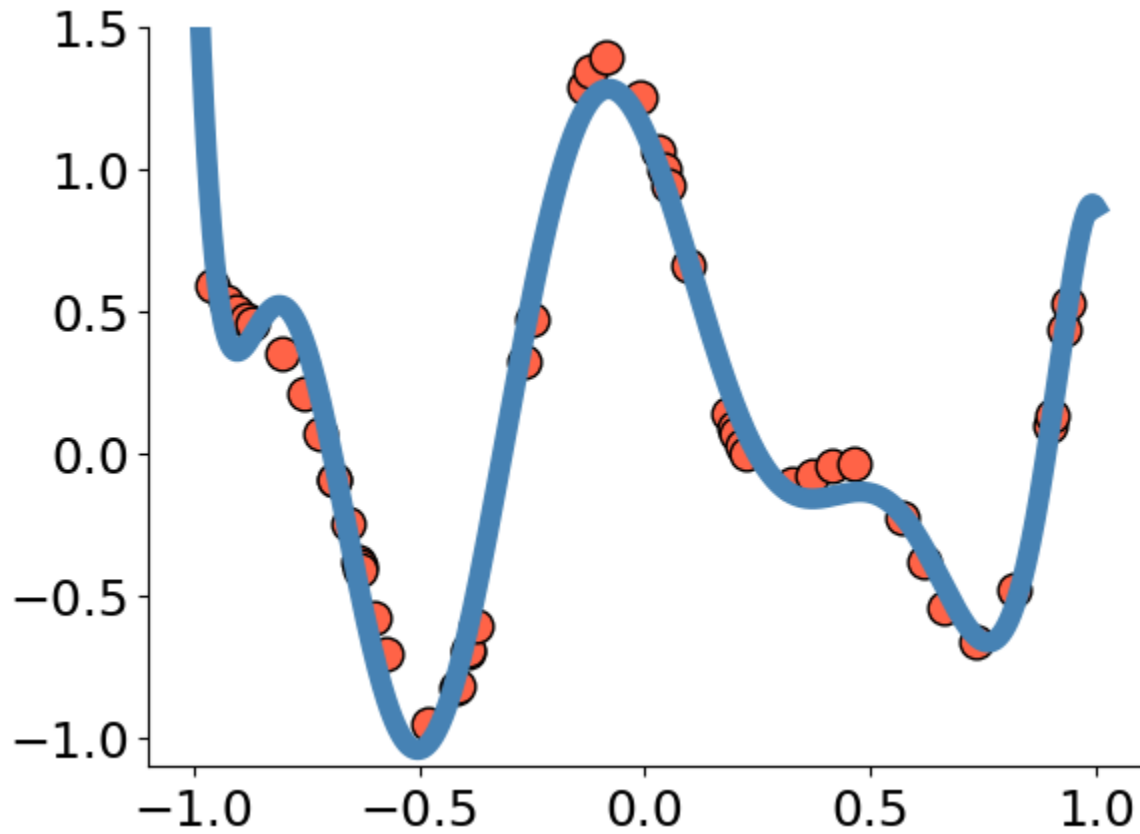


$d = 5$

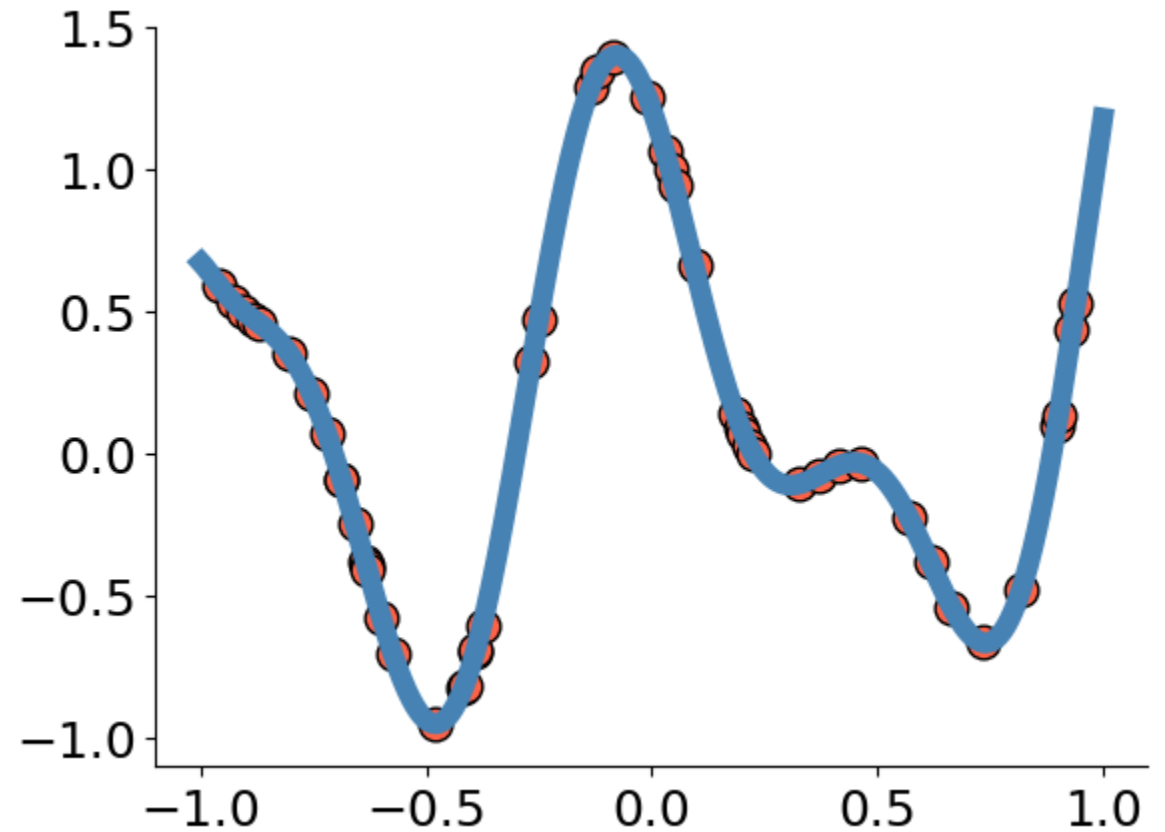


Results for different degrees

$d = 9$



$d = 15$



Why polynomials? Universal function approximation

Let $d \in \mathbb{N}$. A set of functions \mathcal{H} from \mathbb{R}^d to \mathbb{R} is a **universal approximator** if

- for every $\epsilon > 0$,
 - every compact $K \subset \mathbb{R}^d$,
 - and every continuous function $f: \mathbb{R}^d \rightarrow \mathbb{R}$,
- “any error”*
“no matter which inputs”
“no matter which function”

there exists $g \in \mathcal{H}$ such that:

$$\sup_{x \in K} |f(x) - g(x)| < \epsilon$$

Thus, a universal approximator can approximate a continuous function **everywhere** up to **any error**!

Weierstrass approximation theorem



Karl Weierstrass
(1815-1897)

The set of all **polynomials** (of all degrees) is a **universal approximator**.

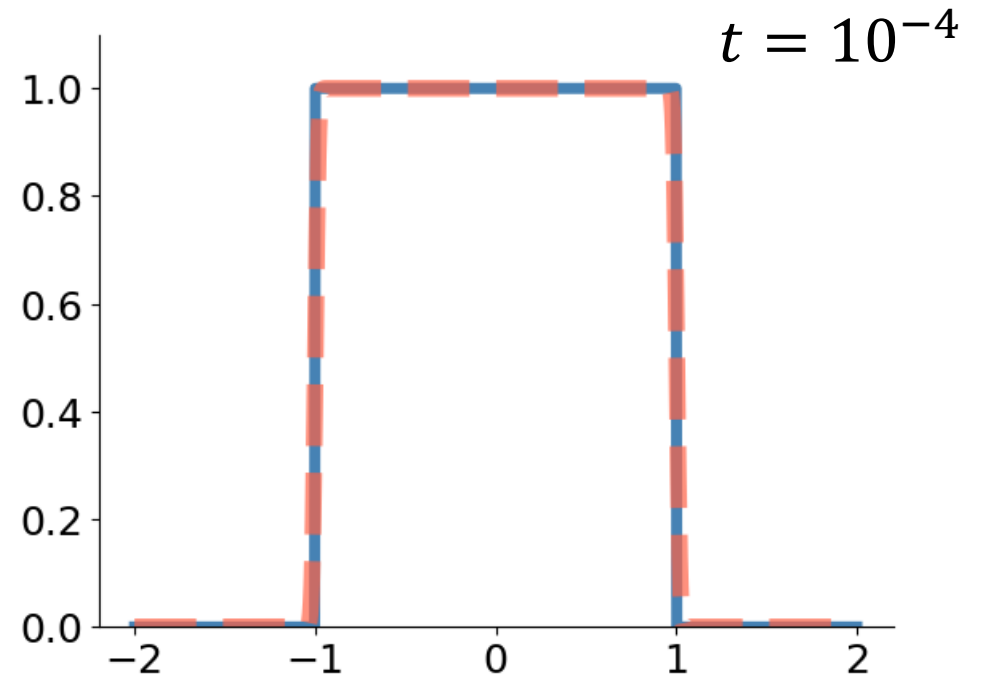
We will sketch the proof here, which is quite short and uses only two major ideas!

First, we “smear” out the actual function f we would like to approximate using:

$$\phi(x, t) = \frac{1}{\sqrt{2\pi t}} \int f(y) e^{-\frac{(x-y)^2}{2t}} dy$$

For $t \rightarrow 0$, $\phi(x, t)$ uniformly approximates f .
Thus, given an $\epsilon > 0$, we can find a t_0 such that

$$|\phi(x, t_0) - f(x)| < \frac{\epsilon}{2}$$



Weierstrass approximation theorem



Karl Weierstrass
(1815-1897)

The set of all **polynomials** (of all degrees) is a **universal approximator**.

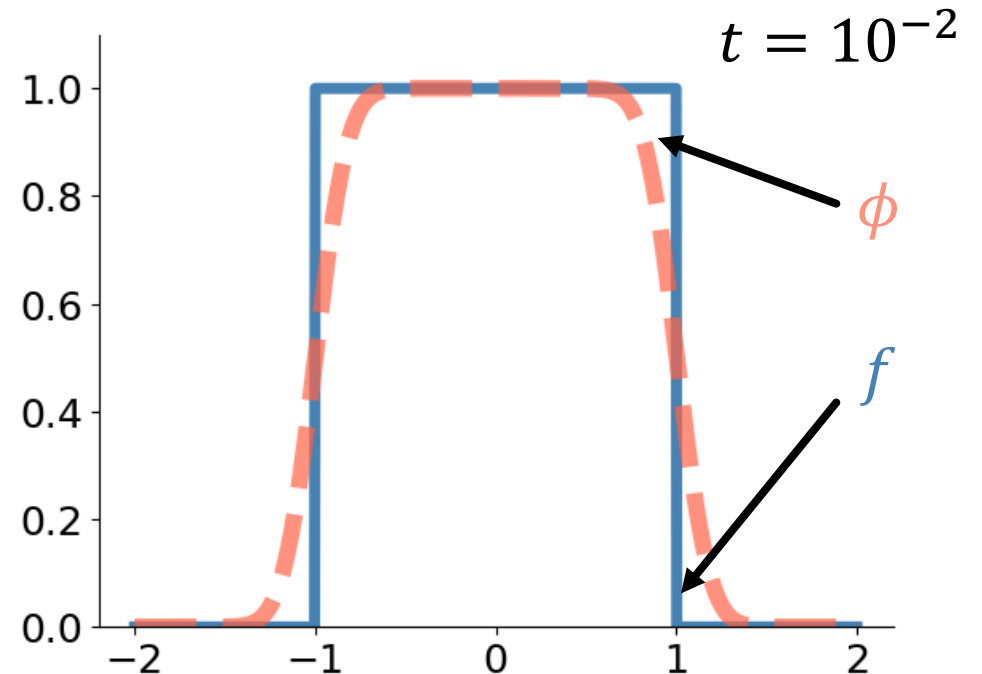
We will sketch the proof here, which is quite short and uses only two major ideas!

First, we “smear” out the actual function f we would like to approximate using:

$$\phi(x, t) = \frac{1}{\sqrt{2\pi t}} \int f(y) e^{-\frac{(x-y)^2}{2t}} dy$$

For $t \rightarrow 0$, $\phi(x, t)$ uniformly approximates f .
Thus, given an $\epsilon > 0$, we can find a t_0 such that

$$|\phi(x, t_0) - f(x)| < \frac{\epsilon}{2}$$



Weierstrass approximation theorem



Karl Weierstrass
(1815-1897)

The set of all **polynomials** (of all degrees) is a **universal approximator**.

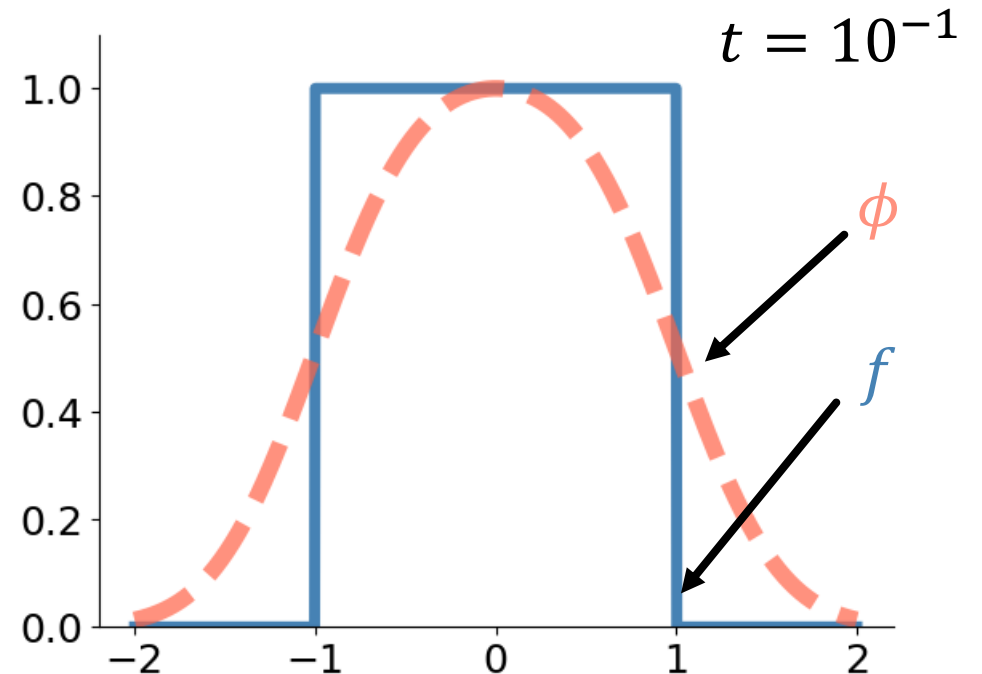
We will sketch the proof here, which is quite short and uses only two major ideas!

First, we “smear” out the actual function f we would like to approximate using:

$$\phi(x, t) = \frac{1}{\sqrt{2\pi t}} \int f(y) e^{-\frac{(x-y)^2}{2t}} dy$$

For $t \rightarrow 0$, $\phi(x, t)$ uniformly approximates f .
Thus, given an $\epsilon > 0$, we can find a t_0 such that

$$|\phi(x, t_0) - f(x)| < \frac{\epsilon}{2}$$



Weierstrass approximation theorem



Karl Weierstrass
(1815-1897)

The set of all **polynomials** (of all degrees) is a **universal approximator**.

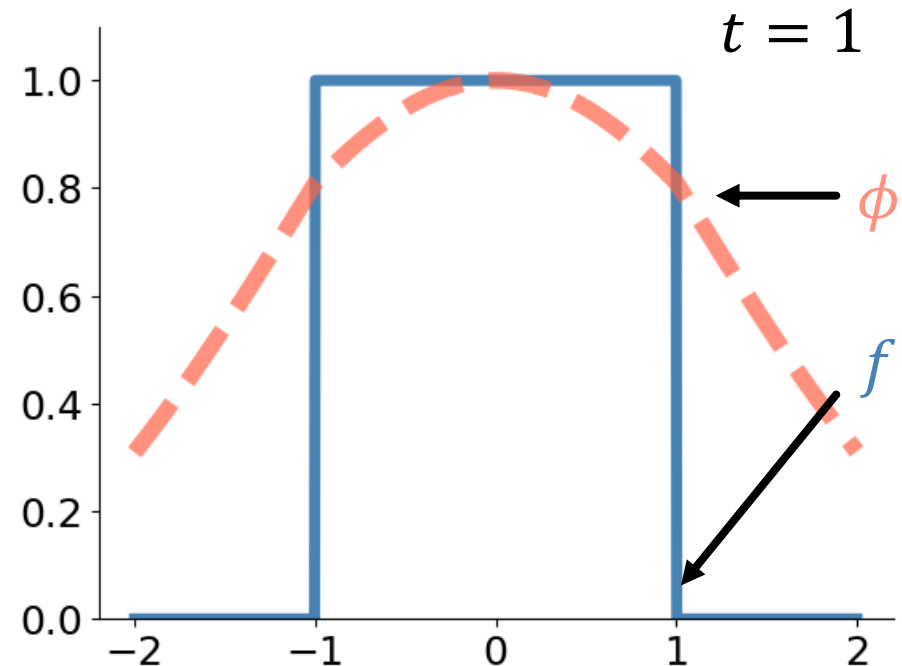
We will sketch the proof here, which is quite short and uses only two major ideas!

First, we “smear” out the actual function f we would like to approximate using:

$$\phi(x, t) = \frac{1}{\sqrt{2\pi t}} \int f(y) e^{-\frac{(x-y)^2}{2t}} dy$$

For $t \rightarrow 0$, $\phi(x, t)$ uniformly approximates f .
Thus, given an $\epsilon > 0$, we can find a t_0 such that

$$|\phi(x, t_0) - f(x)| < \frac{\epsilon}{2}$$



Weierstrass approximation theorem



Karl Weierstrass
(1815-1897)

The set of all **polynomials** (of all degrees) is a **universal approximator**.

We will sketch the proof here, which is quite short and uses only two major ideas!

Second, by Taylor's theorem, $e^{-\frac{(x-y)^2}{2t}}$ can be approximated everywhere via a power series!

If we use a power series of high enough order and plug it into our definition of $\phi(x, t_0)$, we obtain a polynomial $p(x)$ such that:

$$|\phi(x, t_0) - p(x)| \leq \frac{\epsilon}{2}$$

And thus: $|f(x) - p(x)| < |f(x) - \phi(x, t_0)| + |\phi(x, t_0) - p(x)| < \epsilon$



There is an alternative proof by Bernstein using Bernstein polynomials which is really cool as well!

Stone-Weierstrass theorem



Marshall Harvey Stone
(1903 - 1989)

Generalization of the previous theorem!

Let $d \in \mathbb{N}$, let $K \subset \mathbb{R}^d$ be compact, and let our hypothesis set of continuous functions \mathcal{H} from \mathbb{R}^d to \mathbb{R} satisfy that

- for all $\mathbf{x} \in K$, there exists $f \in \mathcal{H}$ such that $f(\mathbf{x}) \neq 0$,
- for all $\mathbf{x} \neq \mathbf{y} \in K$, there exists $f \in \mathcal{H}$ such that $f(\mathbf{x}) \neq f(\mathbf{y})$,
- \mathcal{H} is an algebra of functions, i.e., closed under addition, multiplication, and multiplication with scalars*.

Then \mathcal{H} is a universal approximator.

This is a great tool: just check the criteria if you wanna know if \mathcal{H} is a universal approximator!

*By closed, we mean: if $f, g \in \mathcal{H}$ and $\alpha \in \mathbb{R}$, then $f + g \in \mathcal{H}$, $f \cdot g \in \mathcal{H}$, and $\alpha \cdot f, \alpha \cdot g \in \mathcal{H}$

Exercise

Using Stone-Weierstrass, explain which of the following hypothesis sets are universal approximators:

- Linear regression (*linear in the features*)
- Polynomials
- Decision Trees

Helpful illustration of a decision tree:

