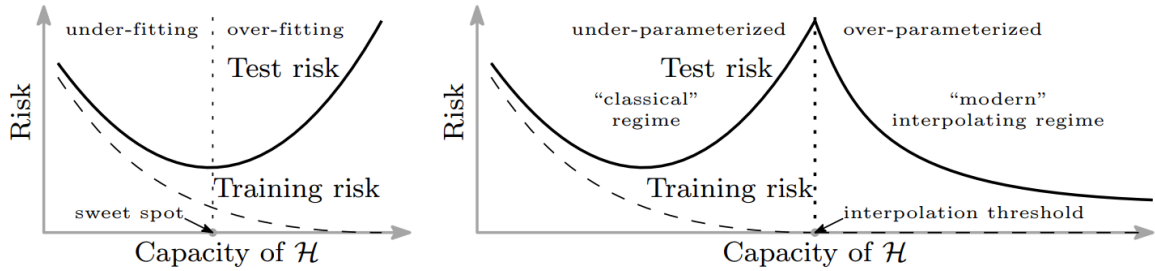


# Double Descent Phenomenon

Egy model kapacitásának növelése egy adott pontig csökkenti a teszt hibát, ezt követően a teszt hiba nő, majd az interpolációs küszöbön túl, a hiba ismét elkezd csökkenni.

[Schaeffer, Rylan, et al. \(2023\)](#)

[Belkin, Mikhail, et al. \(2019\)](#)



## Jelölés

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}, |\mathcal{D}| = N$$

$$(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$$

$$f : \mathbb{R}^d \mapsto \mathbb{R}$$

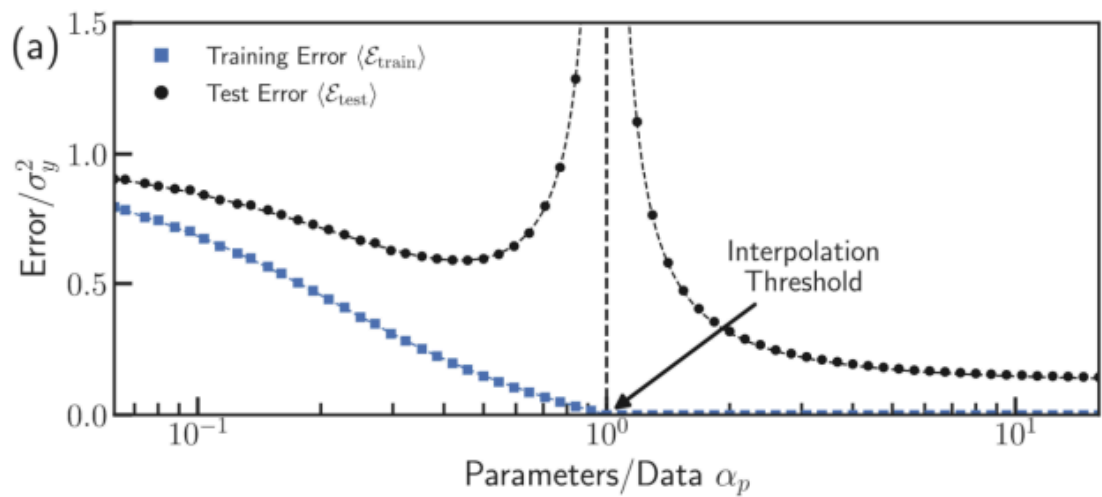
$f \in \mathcal{H}$  függvény család

## Classical Bias-Variance trade-off

1. Ha  $\mathcal{H}$  kapacitása kicsi (*high bias*), akkor minden  $f \in \mathcal{H}$  nem illeszkedik megfelelő mértékben a tanulási adatokra (*under-fitting*), ezért új adatokra is rossz eredményeket fog adni.
2. Ha  $\mathcal{H}$  kapacitása nagy (*high variance*), akkor léteznek  $f \in \mathcal{H}$  amelyek tökéletesen illeszkednek a tanulási adatokra (*over-fitting*), viszont nem képesek általánosítani, ezért új adatokra rossz eredményeket fognak adni.

Ezek alapján  $\mathcal{H}$ -t úgy kell kiválasztani, hogy a kapacitása a *sweet spot*-ban legyen

## "Modern" Interpolating Regime



Ha  $\mathcal{H}$  kapacitását növeljük, az interpolációs küszöbön túl, a tanulási hiba 0 marad, viszont a test hiba ismét elkezdi csökkenni. Ez ellentmond a bias-variance trade-off elvnek.

**Egyszerűsített intuitív magyarázat.** Az interpolációs küszöbön 1 függvény létezik, ami képes tökéletesen illeszkedni a tanulási adatokra, annak a valószínűsége, hogy ez a függvény új adatokra is illeszkedni fog eléggé alacsony. Az interpolációs küszöbön túl, több függvény fog tökéletesen illeszkedni az adatokra, tehát annak a valószínűsége, hogy létezik ezek között legalább egy olyan amely új adatokra is megfelelő predikciókat ad, nagyobb.

## Vizuális Intuición Polinomiális Regresszióval

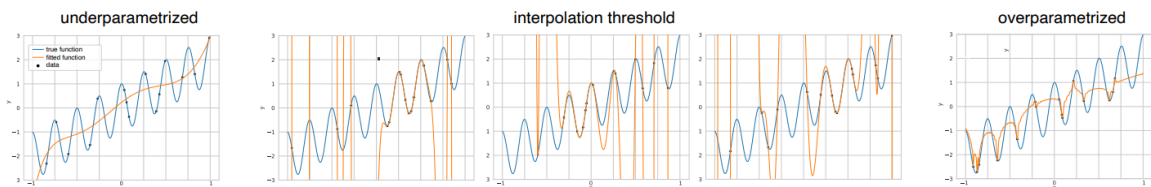
$$y : \mathcal{R} \mapsto \mathcal{R}$$

$$y(x) = 2x + \cos(25x)$$

$$\phi_P : \mathcal{R} \mapsto \mathcal{R}^P$$

$$\phi_P(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_P(x) \end{bmatrix}$$

$$y \approx \phi_P(x) \cdot \Theta_P$$



## Matematikai Intuición Lineáris Regresszióval

$$Y = X \cdot \Theta$$

Lineáris regresszió esetében  $P = D$ . Mivel a paraméterek száma nem növelhető, az adatok számát fogjuk csökkenteni.

1. Ha  $P < N$

$$\hat{\Theta}_{under} = \underset{\Theta}{argmin} ||X\Theta - Y||^2$$

$$\text{Megoldás: } \hat{\Theta}_{under} = (X^T X)^{-1} X^T Y$$

2. Ha  $P > N$

$$\hat{\Theta}_{over} = \underset{\Theta}{argmin} ||\Theta||^2, \forall n \in \{1, \dots, N\} x_n \cdot \Theta = y_n$$

$$\text{Megoldás: } \hat{\Theta}_{over} = X^T (X X^T)^{-1} Y$$

$$\hat{y}_{test,under} = x_{test} \cdot \hat{\Theta}_{under} = x_{test} \cdot (X^T X)^{-1} X^T Y$$

$$\hat{y}_{test,over} = x_{test} \cdot \hat{\Theta}_{over} = x_{test} \cdot X^T (X X^T)^{-1} Y$$

Jelöljük  $\Theta^*$ -al az ismeretlen ideális paramétereket amikre minimális a teszt hiba

Tehát  $Y = X\Theta^* + E$ , ahol  $E$  az adat megtanulhatatlan része

Ezt használva írjuk át a predikciót a két esetben

1.  $P < N$

$$\hat{y}_{test,under} = x_{test} \cdot (X^T X)^{-1} X^T Y$$

$$= x_{test} \cdot (X^T X)^{-1} X^T (X\Theta^* + E)$$

$$= x_{test} \cdot (X^T X)^{-1} X^T X\Theta^* + x_{test} \cdot (X^T X)^{-1} X^T E$$

$$= x_{test} \cdot \Theta^* + x_{test} \cdot (X^T X)^{-1} X^T E$$

$$x_{test} \cdot \Theta^* \stackrel{def}{=} y_{test}^*$$

$$\hat{y}_{test,under} - y_{test}^* = x_{test} \cdot (X^T X)^{-1} X^T E$$

$$(X^T X)^{-1} X^T = X^+ = V \Sigma^+ U^T$$

$$\hat{y}_{test,under} - y_{test}^* = x_{test} \cdot V \Sigma^+ U^T E = \sum_{r=1}^R \frac{1}{\sigma_r} (x_{test} \cdot v_r) (u_r \cdot E)$$

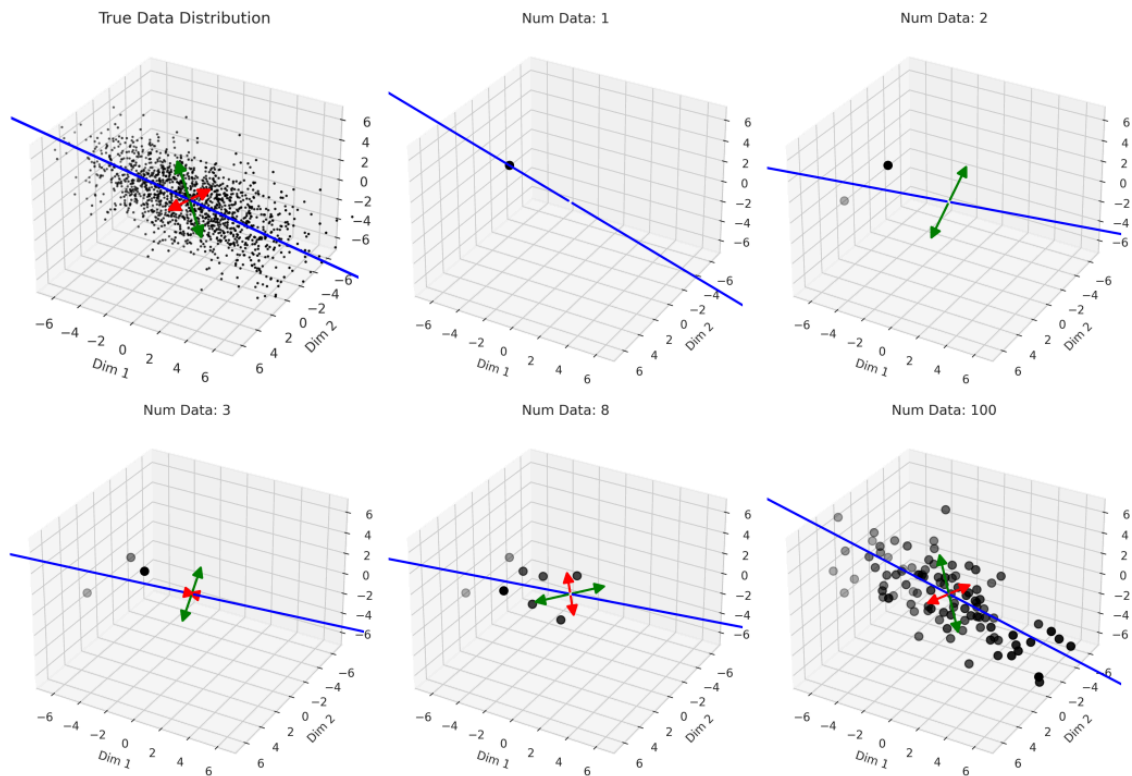
2.  $P > N$ , a számítás hasonló (exercise for the reader)

$$\hat{y}_{test,over} - y_{test}^* = \sum_{r=1}^R \frac{1}{\sigma_r} (x_{test} \cdot v_r) (u_r \cdot E) + x_{test} \cdot (X^T (X X^T)^{-1} X - I_d) \Theta^*$$

A két egyenlet arra mutat, hogy a teszt hibát 3 érték fog meghatározni:

1. Mennyire változnak a tanulási adatok minden irányban:  $\frac{1}{\sigma_r}$
2. Mennyire és milyen irányokban változnak a teszt adatok relatív a tanulási adatokhoz képest:  
 $x_{test} \cdot v_r$
3. Mennyire képes az ideális modell megtanulni a tanulási adatokat:  $u_r \cdot E$

## Miért a legnagyobb a hiba az interpolációs küszöbön?



## Mikor nem történik Double Descent?

- Nincsenek nullához közeli szinguláris értékek ( $\sigma_r$ )
- A teszt adatok nem változnak más irányokban mint a tanulási adatok
- Az ideális model tökéletes illeszkedik a tanulási adathalmazra ( $E=0$ )