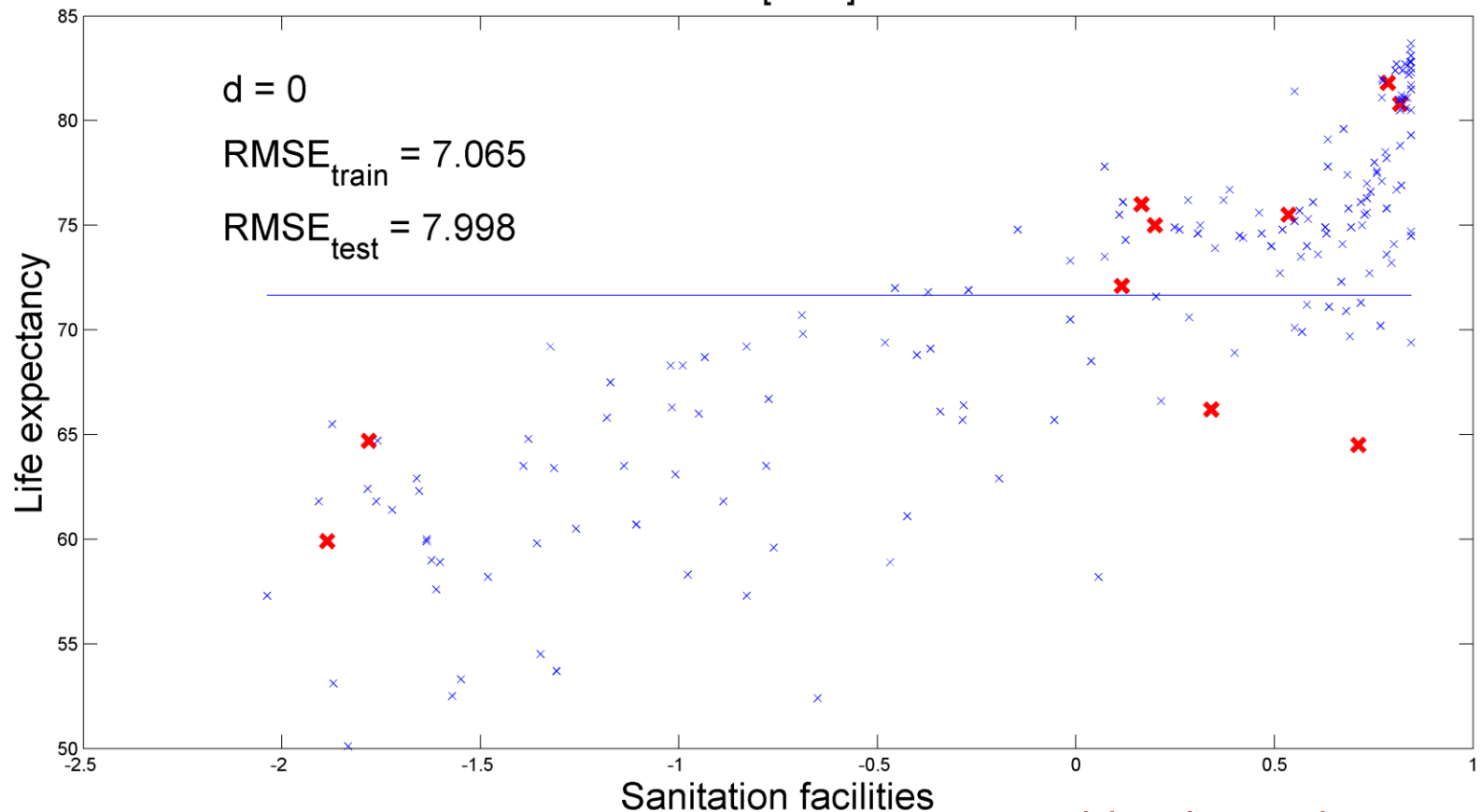


# Uvođenje viših stepena polinoma

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x^j$$

$$\theta = [71.6]$$

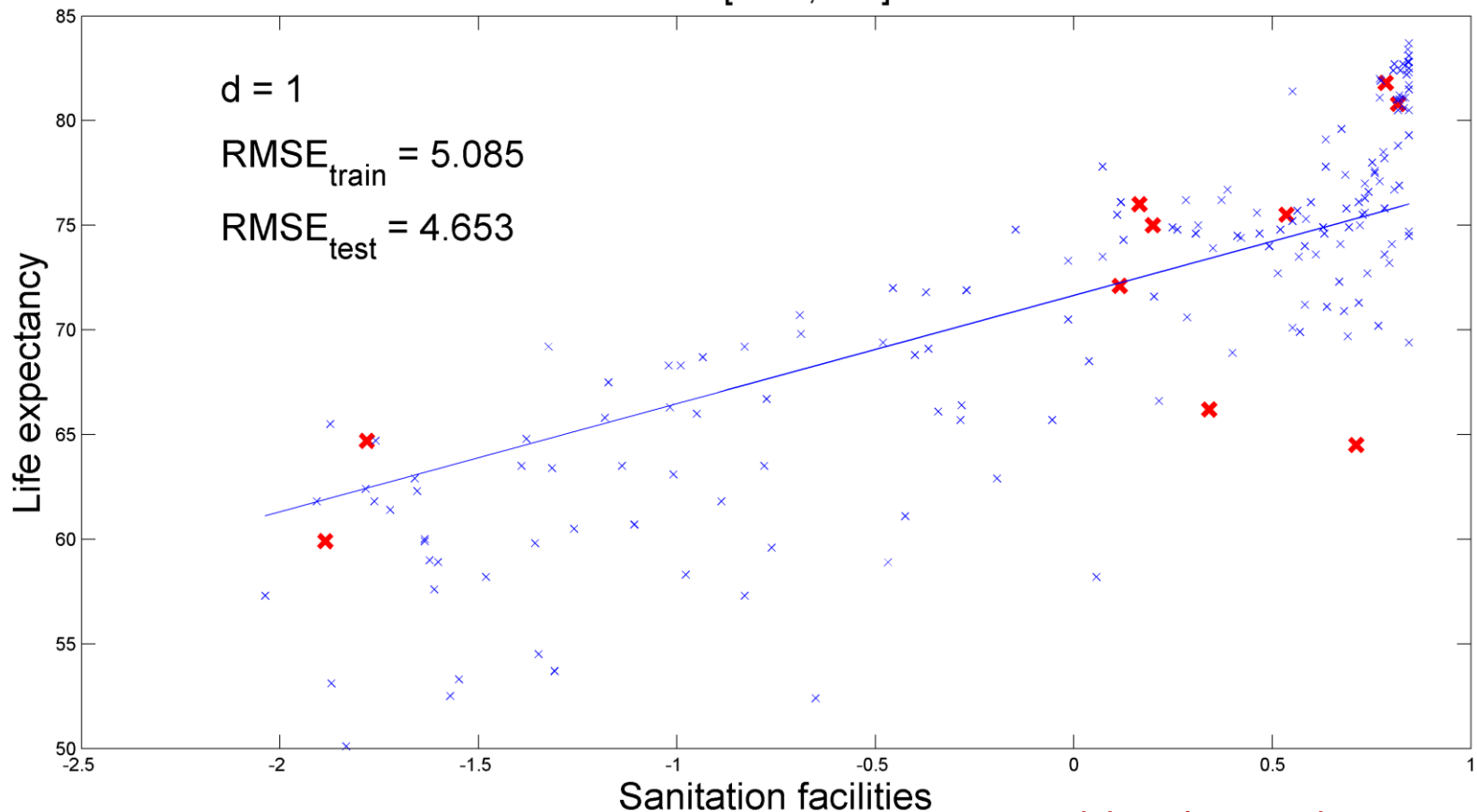


[addHigherOrderFeatures.m](#)

# Uvođenje viših stepena polinoma

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x^j$$

$$\theta = [71.7; 5.2]$$

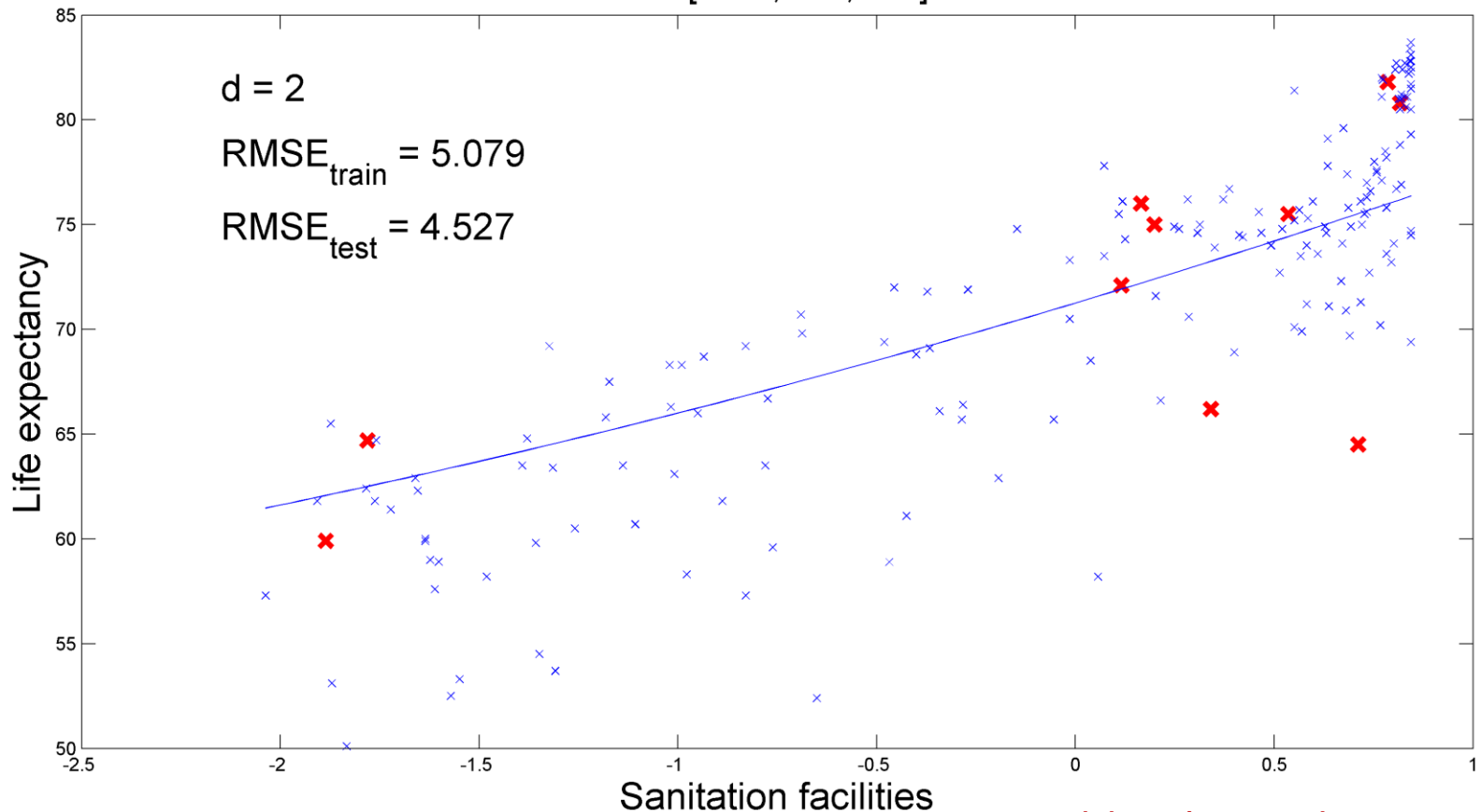


[addHigherOrderFeatures.m](#)

# Uvođenje viših stepena polinoma

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x^j$$

$$\theta = [71.3; 5.7; 0.4]$$

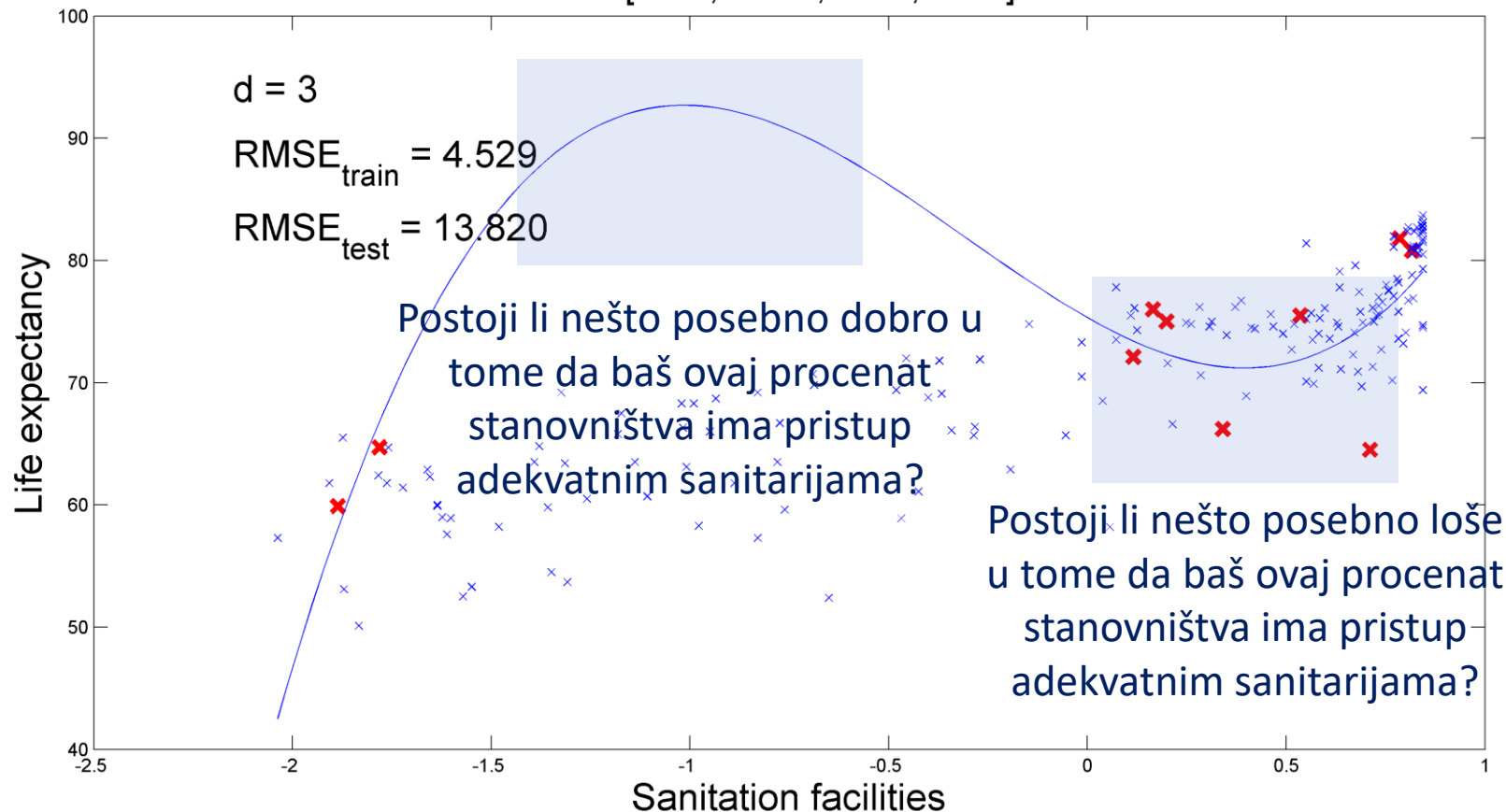


[addHigherOrderFeatures.m](#)

# Uvođenje viših stepena polinoma

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x^j$$

$$\theta = [75.3; -18.4; 14.2; 15.3]$$

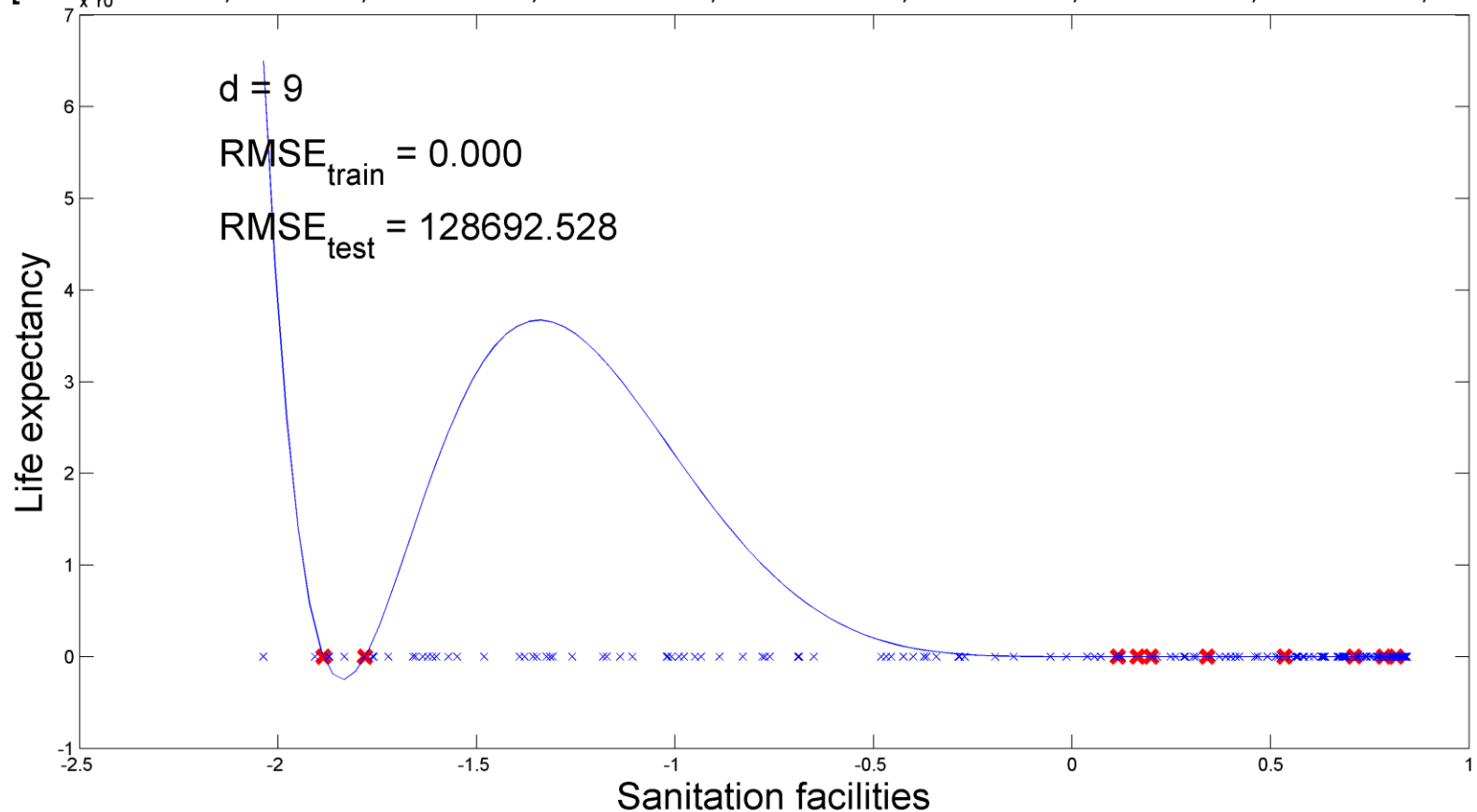


[addHigherOrderFeatures.m](#)

# Uvođenje viših stepena polinoma

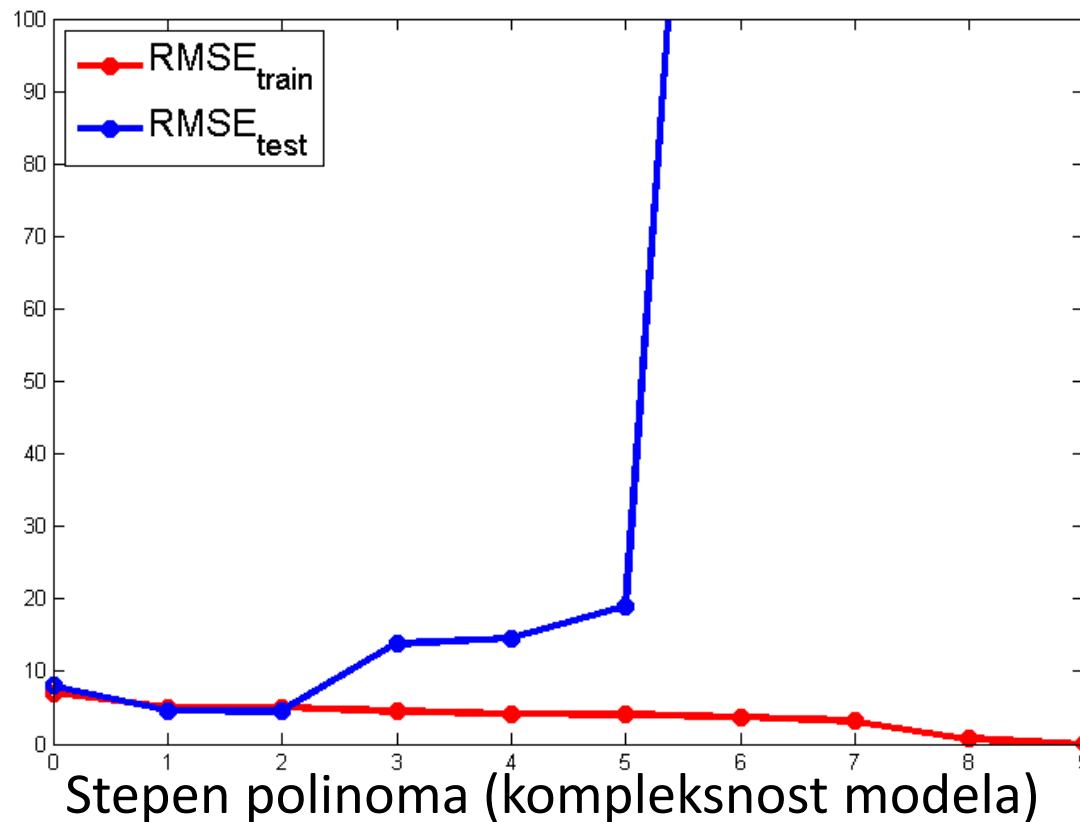
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x^j$$

$\theta = [54.2; -131.4; 6424.8; -48558.0; 143416.4; -161509.5; -15095.6; 132746.5; -25872.3; -33357.1] \times 10^{-5}$



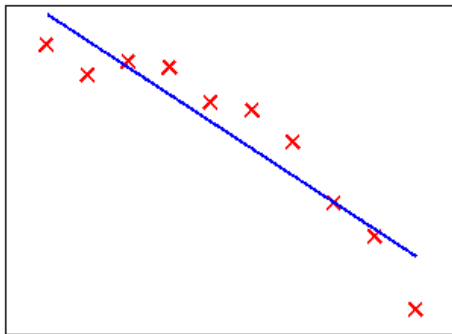
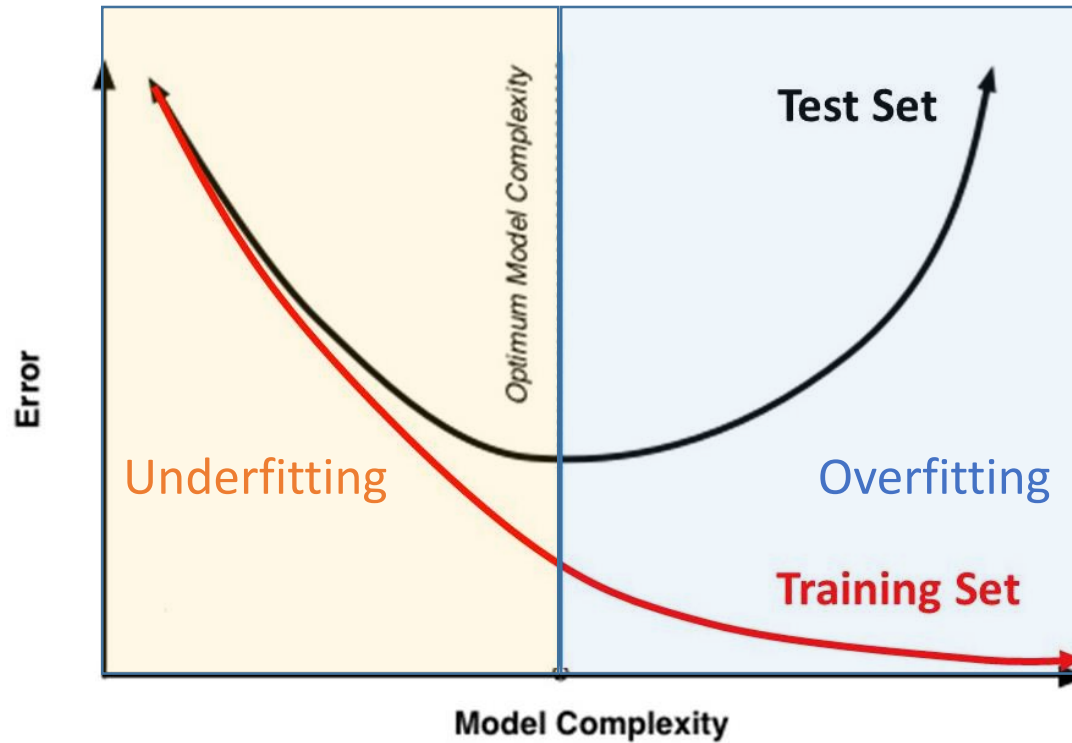
[addHigherOrderFeatures.m](#)

# Uvođenje viših stepena polinoma

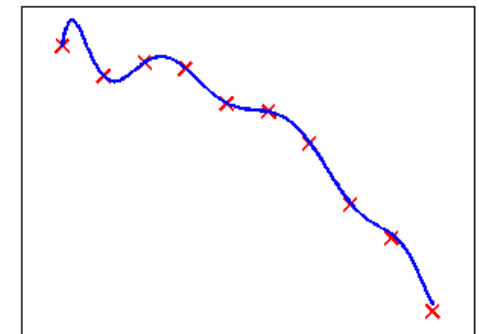
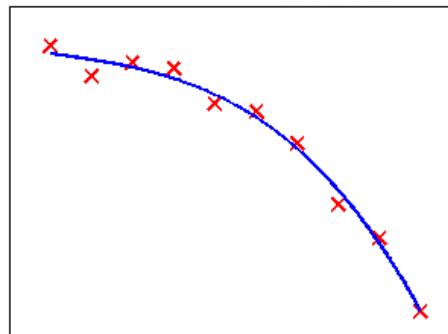


- Trening greška je preterano optimistična jer je  $\hat{\theta}$  fitovan na trening podacima
- Zaključak: mala trening greška nije indikacija dobrih predikcija
  - Uporedite sa učenjem napamet

# Preprilagodžavanje (overfitting)



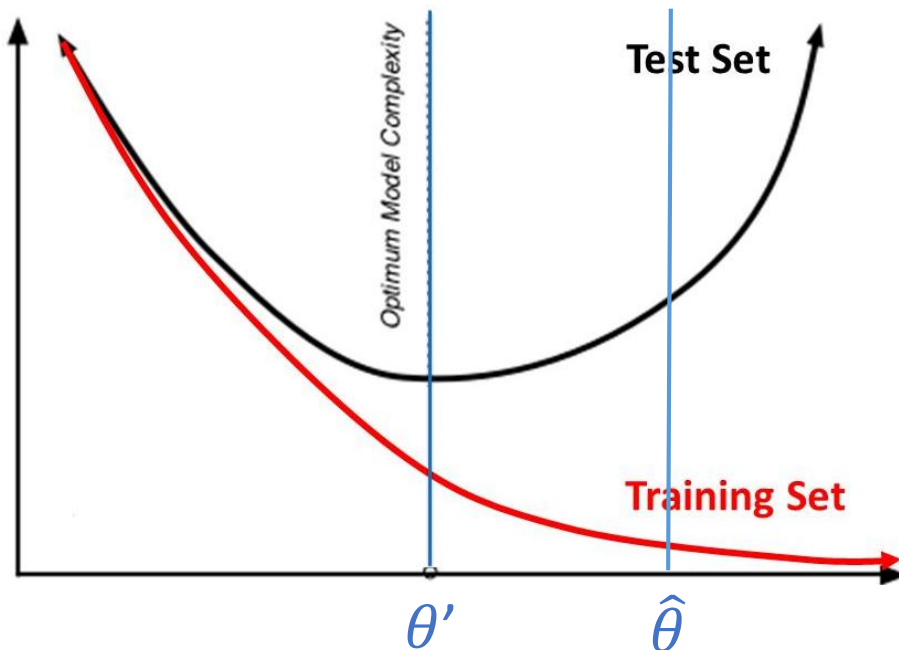
Underfit/High bias



Overfit/High variance

# Preprilagođavanje (overfitting)

- Previše obeležja (a premalo podataka) →
  - hipoteza se jako dobro uklapa u trening podatke
  - ali da loše generalizuje na nove primere
- Uzrok problema: prevelika prilagodljivost modela



Model  $\hat{\theta}$  je overfitovan ako postoji model  $\theta'$  takav da:

- 1) training error ( $\hat{\theta}$ ) < training error( $\theta'$ )
- 2) true error ( $\hat{\theta}$ ) > true error ( $\theta'$ )



# Preprilagođavanje (overfitting)

- Preprilagođavanje nije vezano samo za uvođenje viših stepena polinoma
- Generalno se može javiti kao posledica velikog broja obeležja (kompleksnih/fleksibilnih modela)
- Šta znači veliki broj obeležja ( $D$ ) je relativno u odnosu na broj primera koji imamo u trening skupu ( $N$ )

# Šta se dešava sa koeficijentima $\theta$ ?

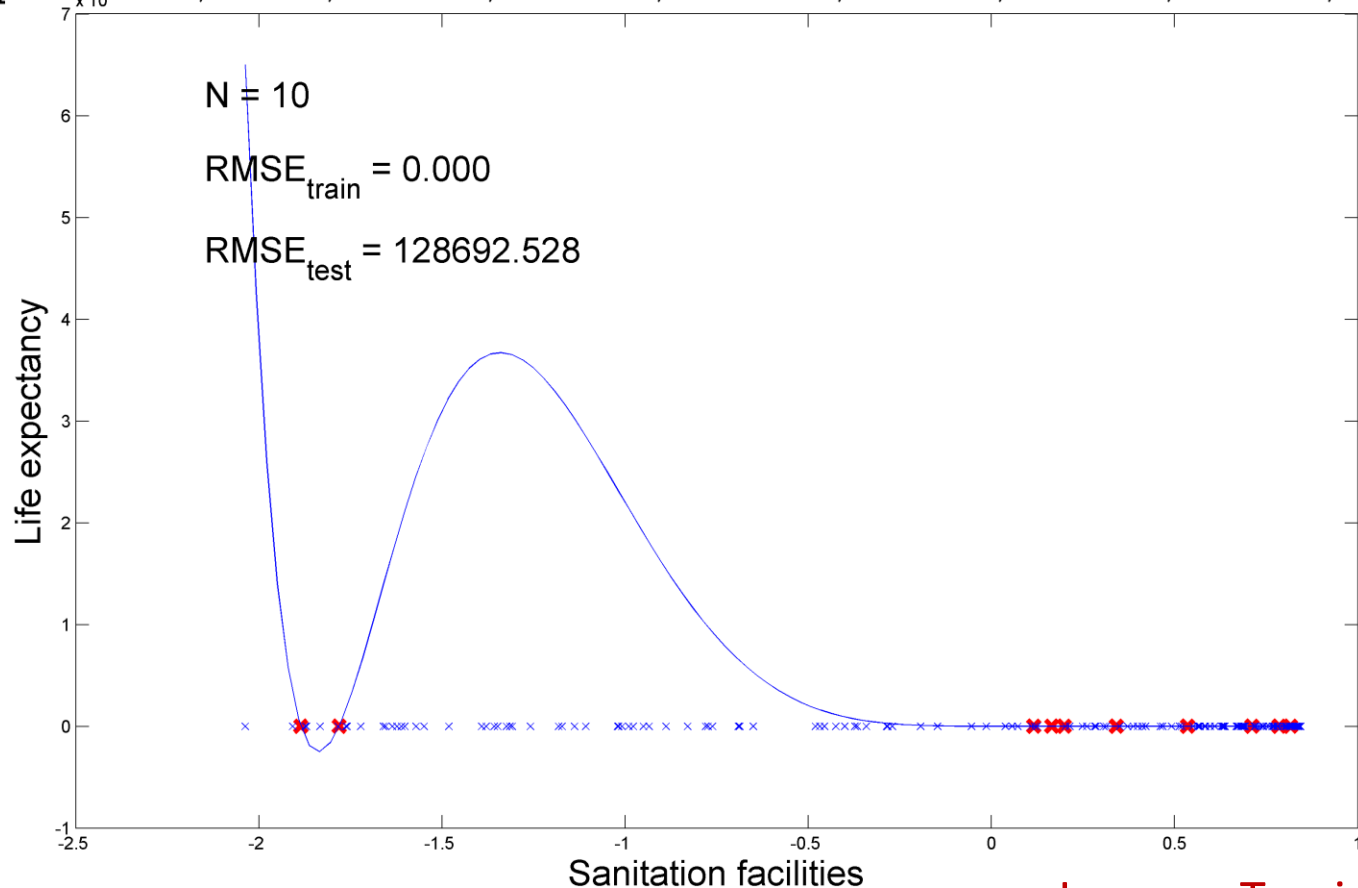
	d = 0	d = 1	d = 2	d = 3	d = 6	d = 9
$\theta_0$	71.6	71.7	71.3	75.3	88.5	54.2
$\theta_1$		5.2	5.7	-18.4	-171.8	-131.4
$\theta_2$			0.4	14.2	551.2	6 424
$\theta_3$				15.3	-471.7	-48 558
$\theta_4$					-417.1	143 416
$\theta_5$					405.2	-161 509
$\theta_6$					210.4	-15 095
$\theta_7$						132 746
$\theta_8$						25 872
$\theta_9$						-33 357

- Zaključak: sa uvećanjem kompleksnosti modela raste i magnituda koeficijenata  $\theta$

# Povećanje $N$ (fiksirana kompleksnost modela)

- Fiksiramo kompleksnost modela (u ovom primeru polinom 9. stepena)
- Postepeno ćemo uvećavati broj (slučajno selektovanih) tačaka u trening skupu

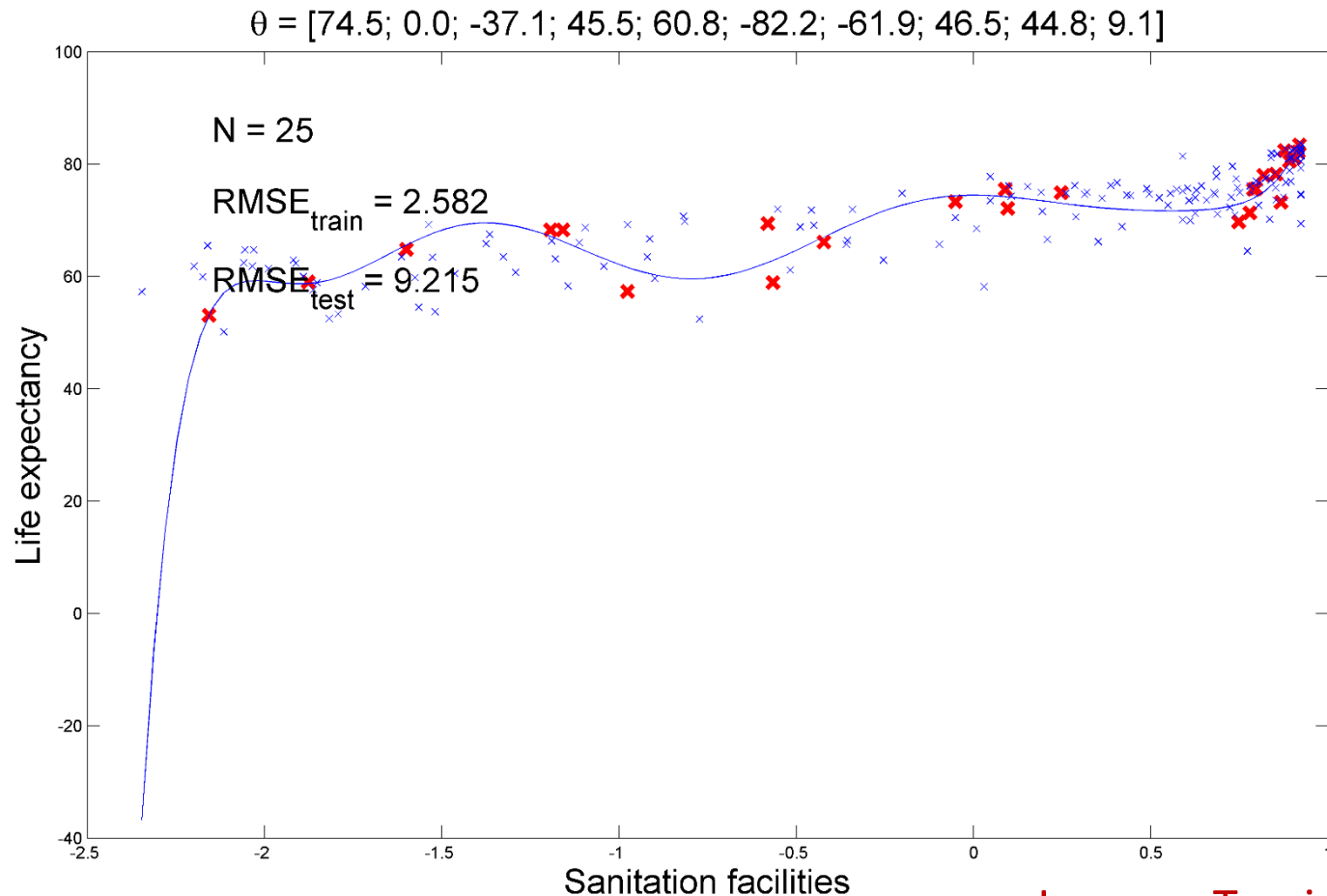
$\theta = [54.2; -131.4; 6424.8; -48558.0; 143416.4; -161509.5; -15095.6; 132746.5; -25872.3; -33357.1]$



IncreaseTraningSetSize.m

# Povećanje $N$ (fiksirana kompleksnost modela)

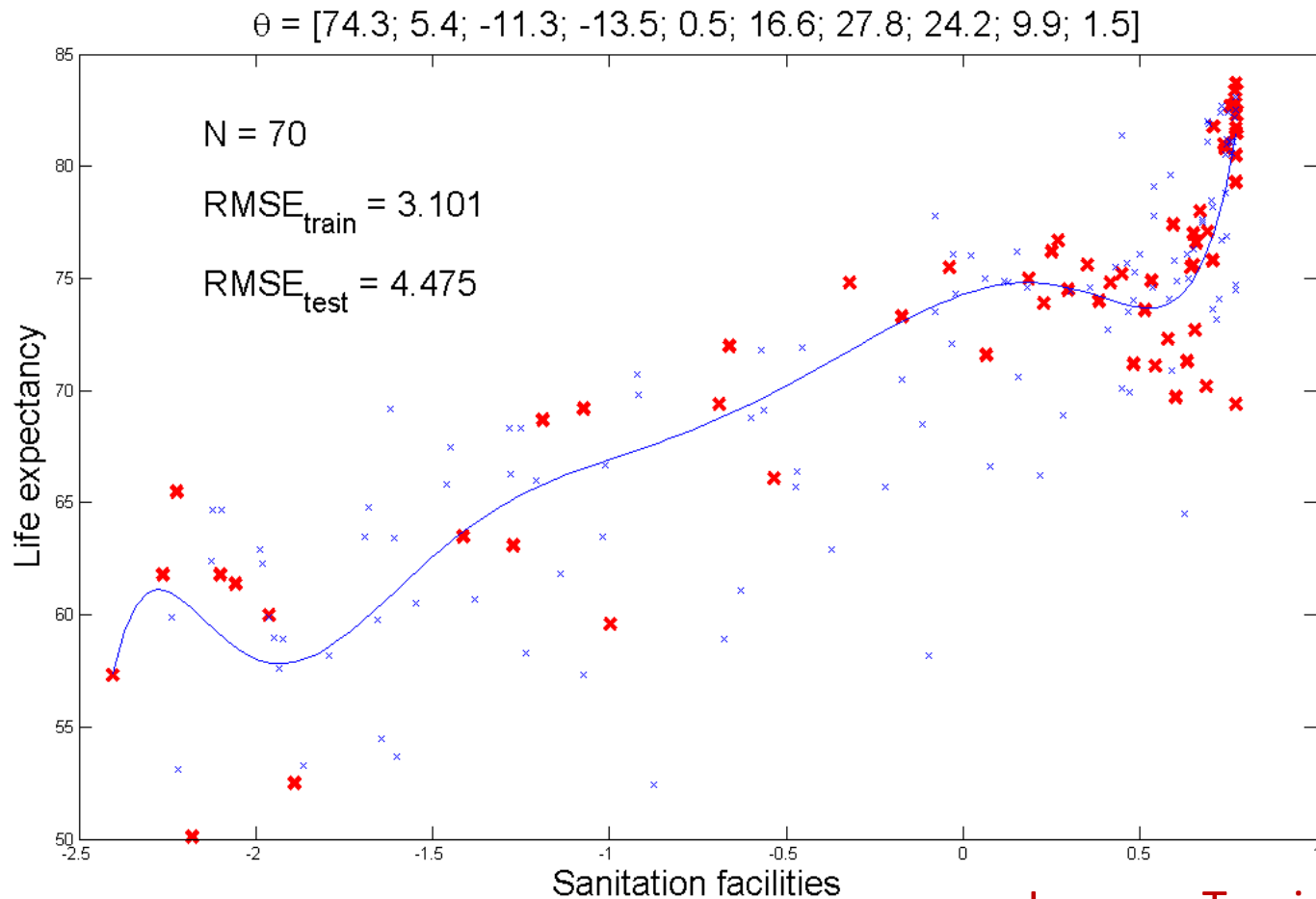
- Fiksiramo kompleksnost modela (u ovom primeru polinom 9. stepena)
- Postepeno ćemo uvećavati broj (slučajno selektovanih) tačaka u trening skupu



IncreaseTraningSetSize.m

# Povećanje $N$ (fiksirana kompleksnost modela)

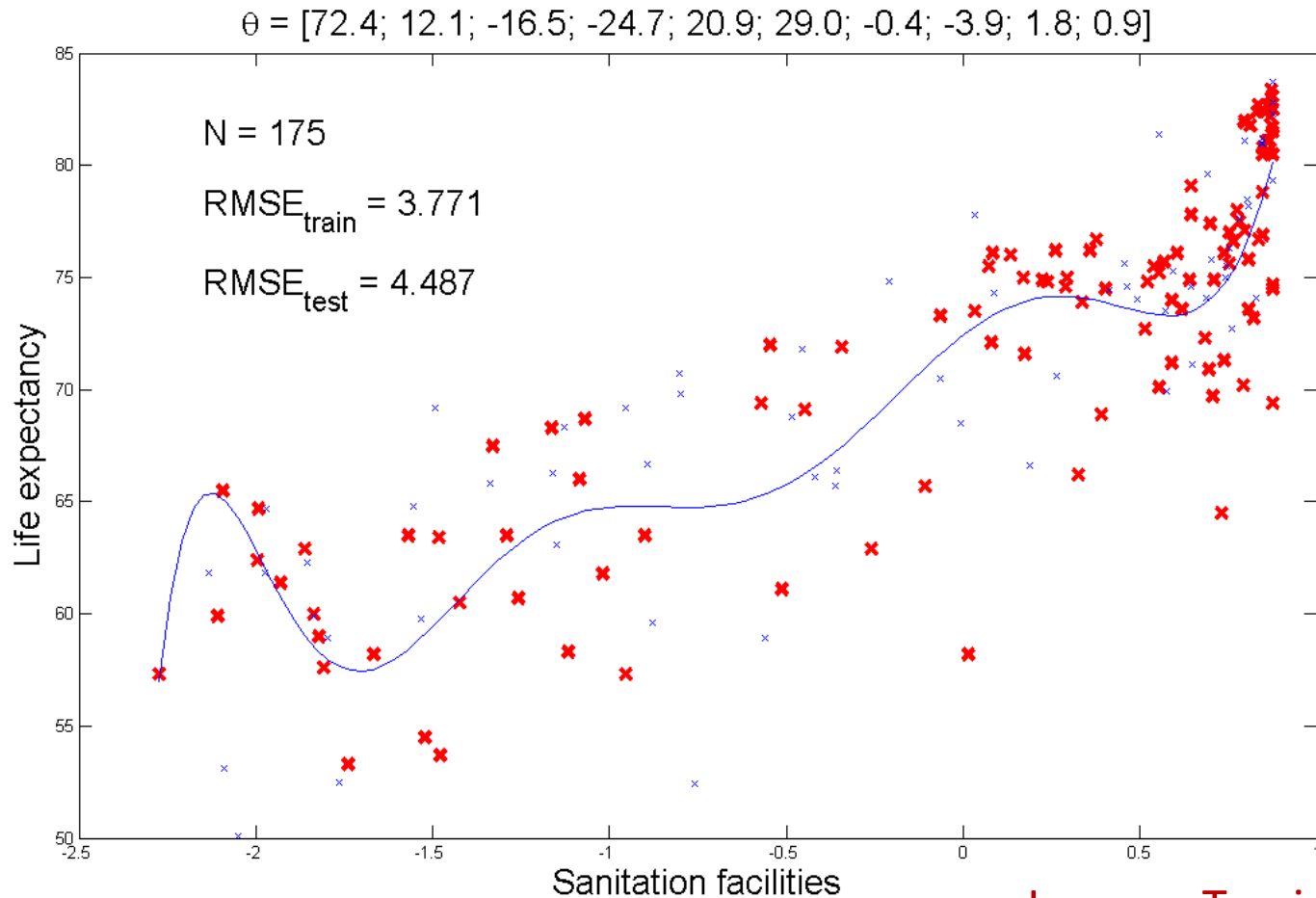
- Fiksiramo kompleksnost modela (u ovom primeru polinom 9. stepena)
- Postepeno ćemo uvećavati broj (slučajno selektovanih) tačaka u trening skupu



IncreaseTraningSetSize.m

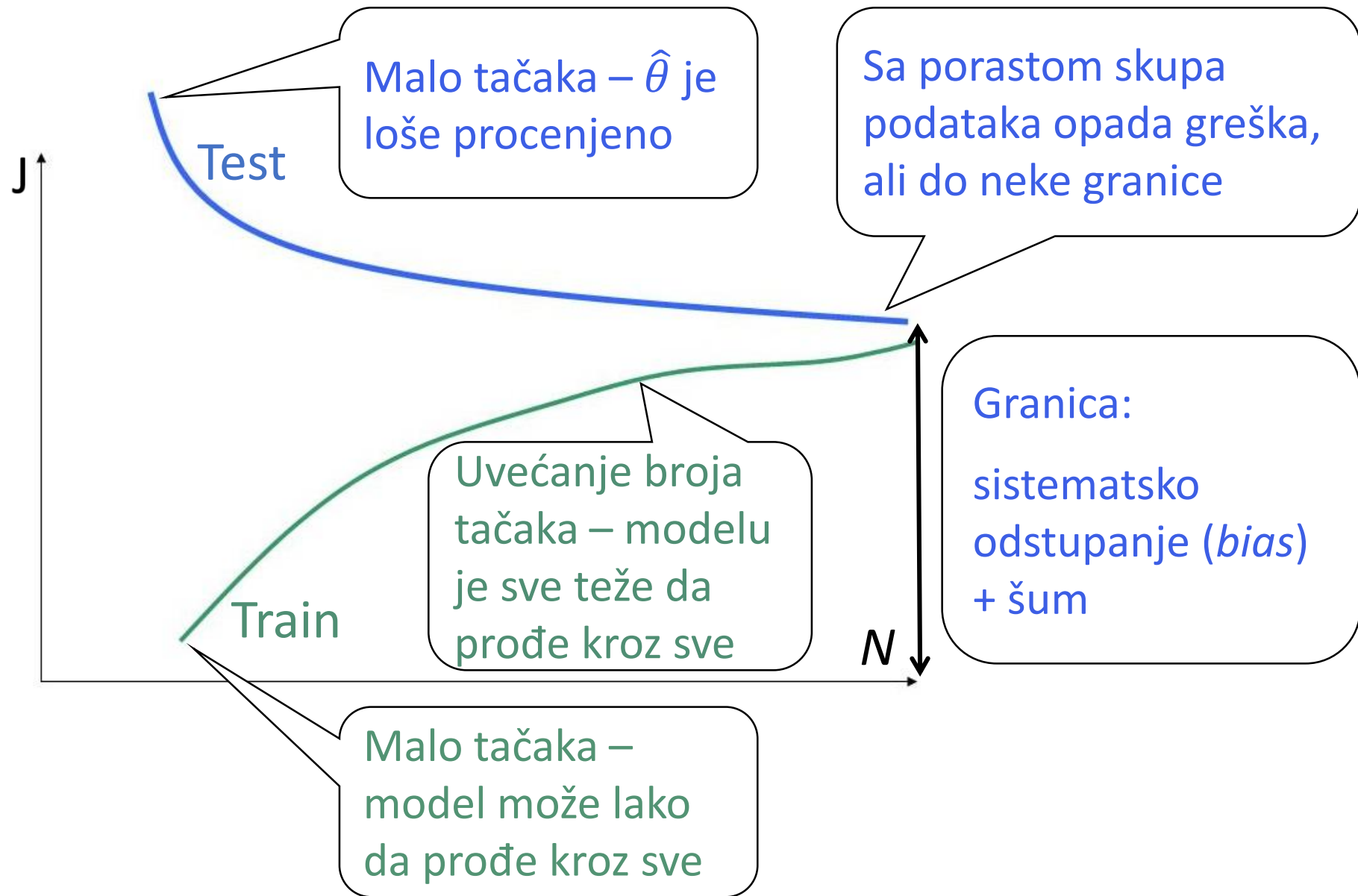
# Povećanje $N$ (fiksirana kompleksnost modela)

- Fiksiramo kompleksnost modela (u ovom primeru polinom 9. stepena)
- Postepeno ćemo uvećavati broj (slučajno selektovanih) tačaka u trening skupu



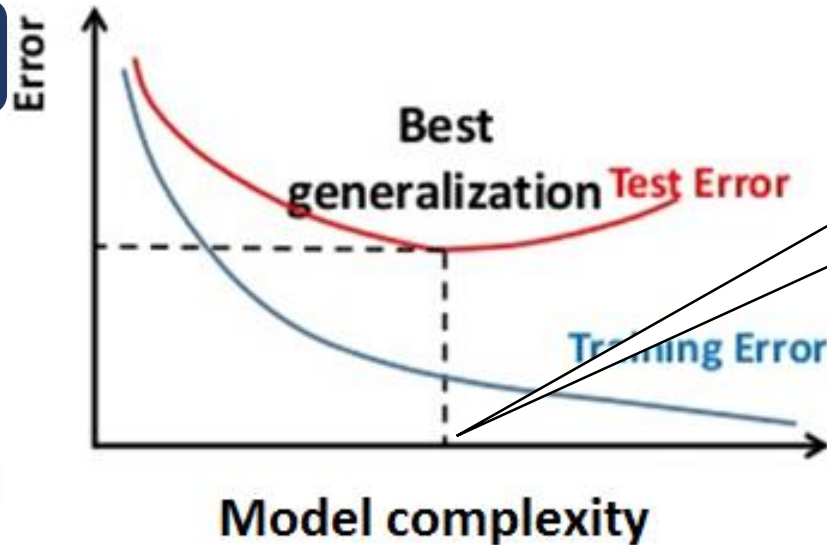
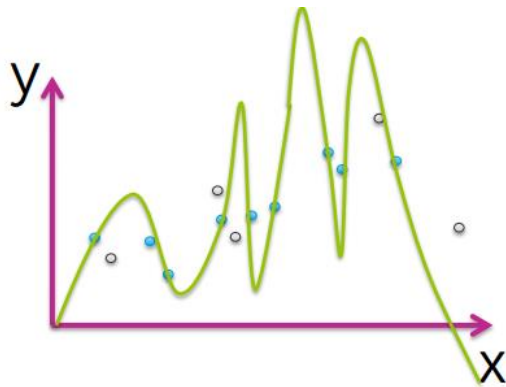
IncreaseTraningSetSize.m

# Povećanje $N$ (fiksirana kompleksnost modela)



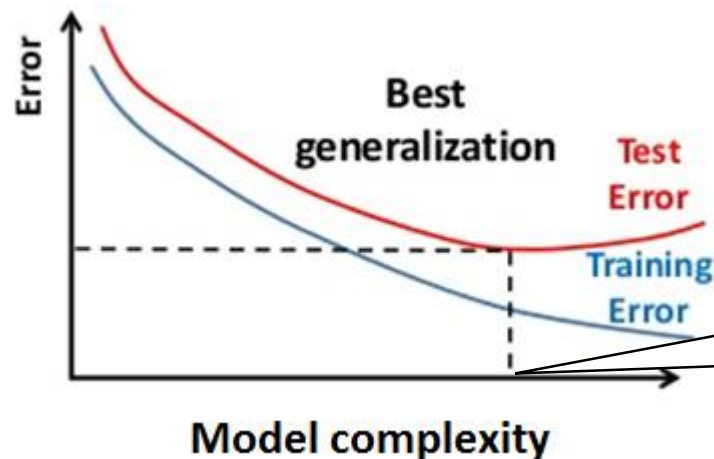
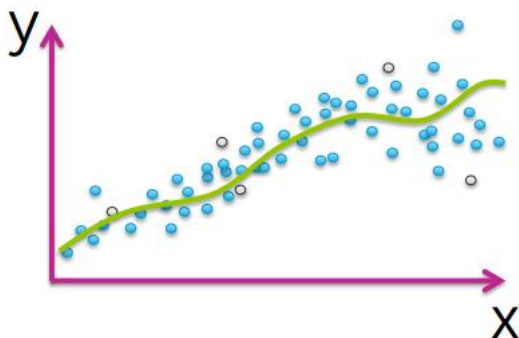
# Zaključak: kako broj obeležja utiče na overfitting?

Malo primera



brzo dolazi do overfittinga

Mnogo primera



Teže je overfitovati



# Zaključak: kako $N$ utiče na preprilagođavanje?

- Zaključak: da bismo izbegli preprilagođavanje moramo imati dovoljno reprezentativan skup podataka
- Za jednu varijablu ovo je prilično teško izvesti
  - npr. da imamo primere zemalja za sve moguće odnose % stanovništva sa pristupom sanitarijama i životnog veka
- Problem se uvećava sa uvećanjem broja ulaza
  - npr. ako pored sanitarija posmatramo i % smrtnosti od saobraćajnih nesreća, teško je pokriti sve moguće kombinacije ova dva ulaza da bismo dobili reprezentativan uzorak

# Zaključak: kako $N$ utiče na prilagođavanje?

- Upravljanje prilagodljivošću modela je od ključnog značaja za dobru generalizaciju
- Ovo je glavni problem mašinskog učenja i izvor njegove najdublje teorije

# Šta se dešava sa koeficijentima $\theta$ ?

	N = 10	N=25	N=70	N=175
$\theta_0$	54.2	74.5	74.3	72.4
$\theta_1$	-131.4	0.0	5.4	12.1
$\theta_2$	6 424	-37.1	-11.3	-16.5
$\theta_3$	-48 558	45.5	-13.5	-24.7
$\theta_4$	143 416	60.8	0.5	20.9
$\theta_5$	-161 509	-82.2	16.6	29.0
$\theta_6$	-15 096	-61.9	27.8	-0.4
$\theta_7$	132 747	46.5	24.2	-3.9
$\theta_8$	-25 872	44.8	9.9	1.8
$\theta_9$	-33 357	9.1	1.5	0.9