

Regularisering i regression

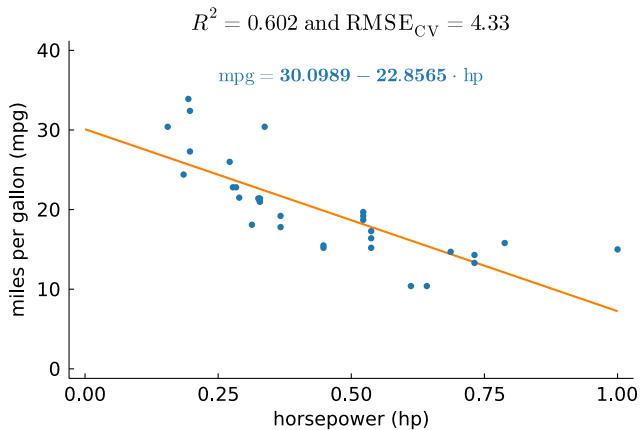
Mattias Villani

Statistiska institutionen
Stockholms universitet

Institutionen för datavetenskap
Linköpings universitet



mtcars data - linjär regression mot hp



Prognosförmåga på nya data - Korsvalidering

Fold 1			Fold 2			Fold 3			Fold 4		
bittyp	mpg	hp	bittyp	mpg	hp	bittyp	mpg	hp	bittyp	mpg	hp
Hornet Sportabout	18.7	0.52	Hornet Sportabout	18.7	0.52	Hornet Sportabout	18.7	0.52	Hornet Sportabout	18.7	0.52
Fiat X1-9	27.3	0.20	Fiat X1-9	27.3	0.20	Fiat X1-9	27.3	0.20	Fiat X1-9	27.3	0.20
Mercedes 450SL	17.3	0.54	Mercedes 450SL	17.3	0.54	Mercedes 450SL	17.3	0.54	Mercedes 450SL	17.3	0.54
Mercedes 450SLC	15.2	0.54	Mercedes 450SLC	15.2	0.54	Mercedes 450SLC	15.2	0.54	Mercedes 450SLC	15.2	0.54
Mercedes 240D	24.4	0.19	Mercedes 240D	24.4	0.19	Mercedes 240D	24.4	0.19	Mercedes 240D	24.4	0.19
Datsun 360	14.3	0.73	Datsun 360	14.3	0.73	Datsun 360	14.3	0.73	Datsun 360	14.3	0.73
Datsun 710	22.8	0.28	Datsun 710	22.8	0.28	Datsun 710	22.8	0.28	Datsun 710	22.8	0.28
Ferrari Dino	19.7	0.52	Ferrari Dino	19.7	0.52	Ferrari Dino	19.7	0.52	Ferrari Dino	19.7	0.52
Ford Pantera L	15.8	0.79	Ford Pantera L	15.8	0.79	Ford Pantera L	15.8	0.79	Ford Pantera L	15.8	0.79
Pontiac Firebird	19.2	0.52	Pontiac Firebird	19.2	0.52	Pontiac Firebird	19.2	0.52	Pontiac Firebird	19.2	0.52
Toyota Corolla	21.5	0.29	Toyota Corolla	21.5	0.29	Toyota Corolla	21.5	0.29	Toyota Corolla	21.5	0.29
AMC Javelin	15.2	0.45	AMC Javelin	15.2	0.45	AMC Javelin	15.2	0.45	AMC Javelin	15.2	0.45
Camaro Z28	13.3	0.73	Camaro Z28	13.3	0.73	Camaro Z28	13.3	0.73	Camaro Z28	13.3	0.73
Fiat 128	32.4	0.20	Fiat 128	32.4	0.20	Fiat 128	32.4	0.20	Fiat 128	32.4	0.20
Mercedes 280C	17.8	0.37	Mercedes 280C	17.8	0.37	Mercedes 280C	17.8	0.37	Mercedes 280C	17.8	0.37
Lotus Europa	30.4	0.34	Lotus Europa	30.4	0.34	Lotus Europa	30.4	0.34	Lotus Europa	30.4	0.34
Cadillac Fleetwood	19.4	0.41	Cadillac Fleetwood	19.4	0.41	Cadillac Fleetwood	19.4	0.41	Cadillac Fleetwood	19.4	0.41
Chrysler Imperial	14.7	0.89	Chrysler Imperial	14.7	0.89	Chrysler Imperial	14.7	0.89	Chrysler Imperial	14.7	0.89
Mazda RX4	21	0.33	Mazda RX4	21	0.33	Mazda RX4	21	0.33	Mazda RX4	21	0.33
Volvo 142E	21.4	0.39	Volvo 142E	21.4	0.39	Volvo 142E	21.4	0.39	Volvo 142E	21.4	0.39
Mazda RX4 Wag	21	0.33	Mazda RX4 Wag	21	0.33	Mazda RX4 Wag	21	0.33	Mazda RX4 Wag	21	0.33
Mercedes 230	22.8	0.28	Mercedes 230	22.8	0.28	Mercedes 230	22.8	0.28	Mercedes 230	22.8	0.28
Toyota Corolla	33.9	0.19	Toyota Corolla	33.9	0.19	Toyota Corolla	33.9	0.19	Toyota Corolla	33.9	0.19
Mercedes 280	19.2	0.37	Mercedes 280	19.2	0.37	Mercedes 280	19.2	0.37	Mercedes 280	19.2	0.37
Dodge Challenger	15.5	0.45	Dodge Challenger	15.5	0.45	Dodge Challenger	15.5	0.45	Dodge Challenger	15.5	0.45
Lincoln Continental	19.4	0.44	Lincoln Continental	19.4	0.44	Lincoln Continental	19.4	0.44	Lincoln Continental	19.4	0.44
Volvo	18.1	0.31	Volvo	18.1	0.31	Volvo	18.1	0.31	Volvo	18.1	0.31
Honda Civic	30.4	0.16	Honda Civic	30.4	0.16	Honda Civic	30.4	0.16	Honda Civic	30.4	0.16
Hornet 4 Drive	21.4	0.33	Hornet 4 Drive	21.4	0.33	Hornet 4 Drive	21.4	0.33	Hornet 4 Drive	21.4	0.33
Mercedes 450SE	18.4	0.54	Mercedes 450SE	18.4	0.54	Mercedes 450SE	18.4	0.54	Mercedes 450SE	18.4	0.54
Mercedes 300	15	1.00	Mercedes 300	15	1.00	Mercedes 300	15	1.00	Mercedes 300	15	1.00
Porsche 914-2	26	0.27	Porsche 914-2	26	0.27	Porsche 914-2	26	0.27	Porsche 914-2	26	0.27

Träning
Test

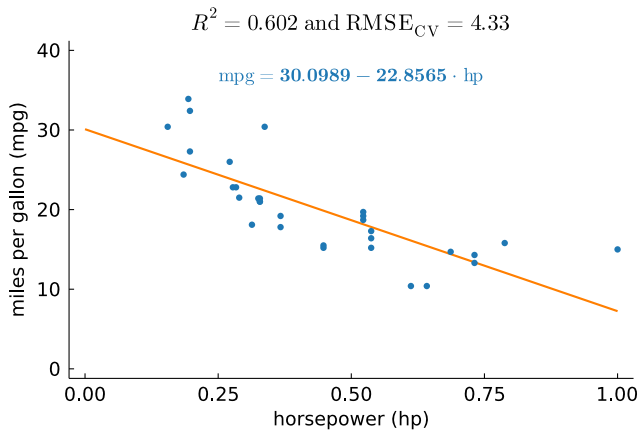
Prognosfel i testdata

$$\text{MSE}_{CV} = \frac{\sum_{j=1}^{n_{\text{test}}} (y_j - \hat{y}_j^*)^2}{n}$$

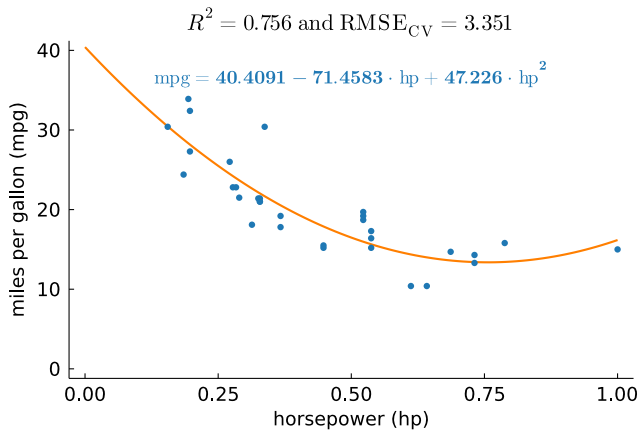
Lättare att tolka Root MSE

$$\text{RMSE}_{CV} = \sqrt{\text{MSE}_{CV}}$$

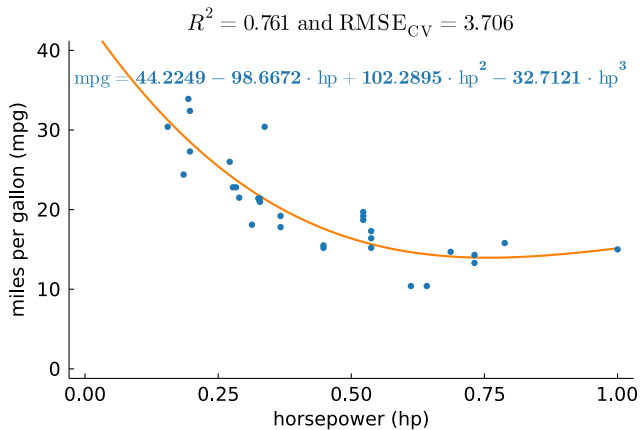
mtcars data - linjär regression mot hp



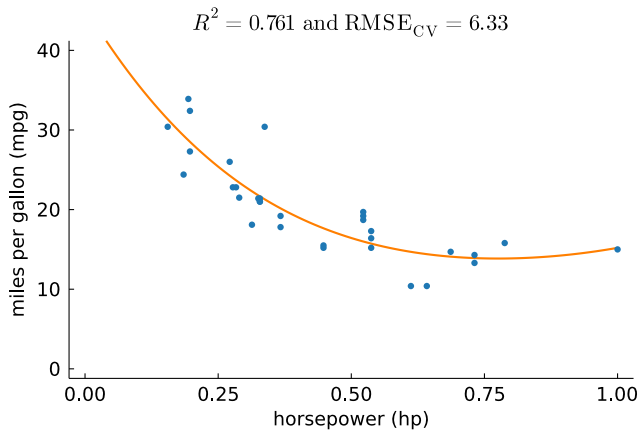
mtcars data - kvadratisk regression mot hp



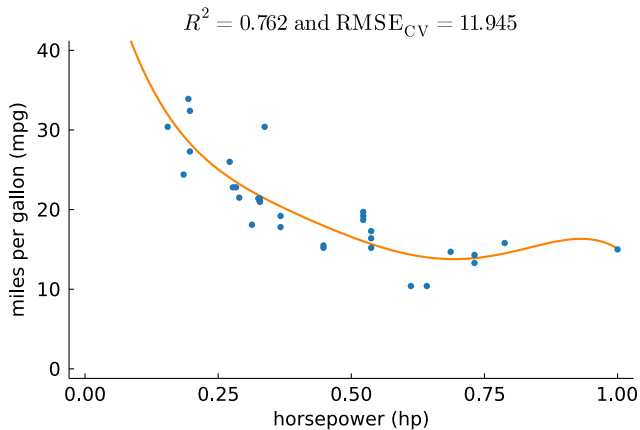
mtcars data - kubisk regression mot hp



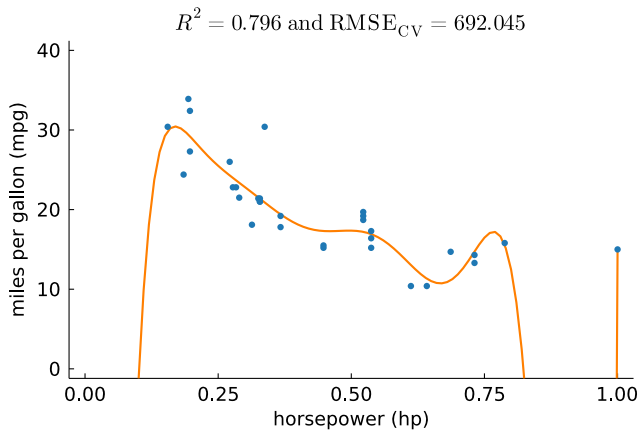
mtcars data - polynomregression ordning 4



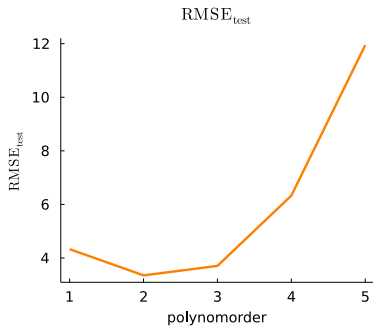
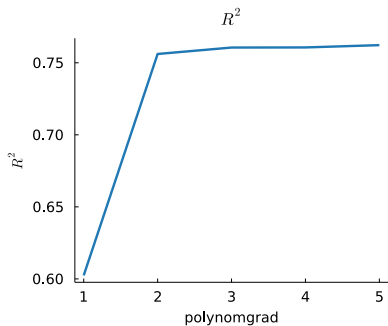
mtcars data - polynomregression ordning 5



mtcars data - polynomregression ordning 10



mtcars data - R^2 och RMSE-CV($K = 4$)



L2-regularisering (Ridge regression)

- För många x-variabler \Rightarrow MK-metoden överanpassar data.
- Modellen är överparametriserad.
- Variabelselektion (t ex forward selection) är en lösning.
- L2-regularisering minimerar en straffad SSE:

$$Q_{\lambda} = \sum_{i=1}^n (y_i - a - b_1 x_{1i} - \dots - b_k x_{ki})^2 + \lambda \cdot \sum_{j=1}^k b_j^2$$

- Stort λ kommer krympa estimaten av b_j mot noll.
- Skattningen är nu biased, men har lägre varians.
- Bias-Variance trade-off.

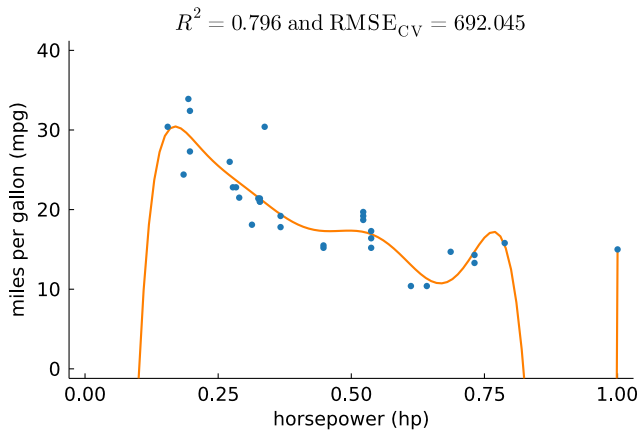
L1-regularisering (Lasso regression)

- L1-regularisering (Lasso) straffar med **absolutbelopp**:

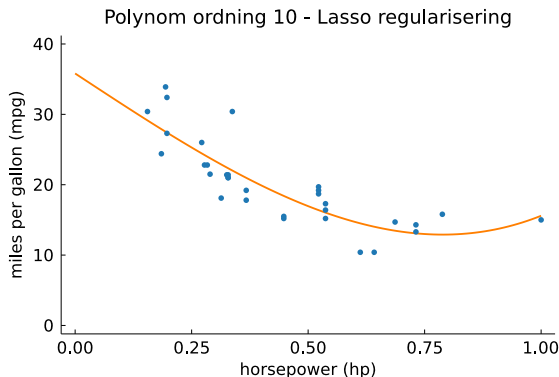
$$Q_- = \sum_{i=1}^n (y_i - a - b_1 x_{1i} - \dots - b_k x_{ki})^2 + \lambda \cdot \sum_{j=1}^k |b_j|$$

- Lasso har två effekter:
 - ▶ krymper b_j mot noll (**shrinkage**)
 - ▶ kan sätta vissa b_j exakt till noll (**selection**)
- glmnet paketet i R gör både L1 och L2 regularisering och mer.
- Lasso är extremt populär. Go-to när man har väldigt många förklarande variabler.

Polynom ordning 10 - ingen regularisering



Polynom ordning 10 - L1-regularisering



$$Q_{\lambda} = \sum_{i=1}^n (y_i - a - b_1 x_{1i} - \dots - b_k x_{ki})^2 + \lambda \cdot \sum_{j=1}^k |b_j|$$

■ $a = 35.80$, $b_1 = -43.58$, $b_3 = 23.32$.

■ $b_2 = 0$ och $b_4 = \dots = b_{10} = 0$ (**variabelselektion**).