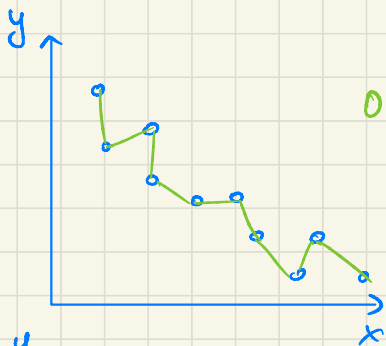


L1 & L2 REGULARISATION



overfitting: $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

↓
try to make $\theta_3 \dots \theta_n$ almost zero so that we rule out those factors \therefore shrink parameters to help generalize better



↓
balanced: $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

controllable

now, $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n \theta_i^2$

if θ gets bigger, the term $\lambda \sum_{i=1}^n \theta_i^2$ gets bigger

\therefore penalizing higher values of θ to \downarrow MSE

$$L2 = \theta_i^2$$

in $L1$, $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n |\theta_i|$

Ridge Regression:

↳ L2 Regularization

↳ Addresses problem of multicollinearity among predictor variables which occurs when independent variables are highly correlated \therefore leads to unreliable estimates

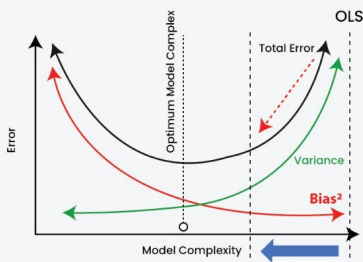
How? Adds a regularization term to least squares to penalize large coefficients, in the form of penalty \therefore also solves overfitting

Example: $\text{Price} = 1000 \cdot \text{Size} - 500 \cdot \text{Age} + \text{Noise}$

↓ after L2

$\text{Price} = 800 \cdot \text{Size} - 300 \cdot \text{Age} + \text{Less Noise}$

Bias-Variance Tradeoff in Ridge Regression



- As $\lambda \uparrow$, bias \downarrow and var \uparrow
- goal: find λ that balances bias and variance

How to select a ridge parameter?

↳ Cross-validation: data split into subsets; model trains on some and validates on the others

- K-Fold CV: data split into K subsets, trained on K-1 and validated on one; repeated K times s.t. each fold serves as validation set once
- Leave-One-Out CV
- ↳ Generalized CV

Pros:

- Bias vs variance
- retains original parameters

Cons:

- Introduces bias (can lead to underestimation of true effects of predictors)
- Choosing ridge parameter can be challenging