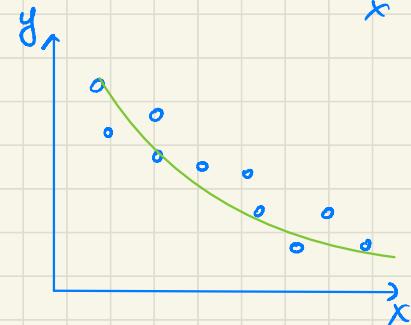


LI & L2 REGULARISATION

↳ adding information



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

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

$$\text{balanced: } y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

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

controllable

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

if θ gets bigger, the term $\lambda \sum_{i=1}^n \theta_i^2$ gets bigger
∴ penalizing higher values of θ to \downarrow MSE

$$\text{L2} = \theta_i^2$$

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

* Regularization techniques are commonly used for linear regression

Ridge Regression:

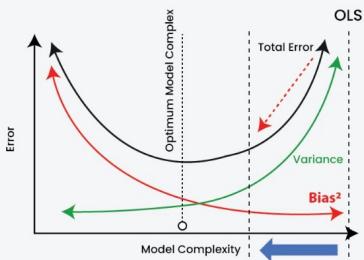
- ↳ L₂ Regularization
 - ↳ Addresses problem of multicollinearity among predictor variables which occurs when independent variables are highly correlated \therefore leads to unreliable estimates
- How? Adds a regularization term to least squares to penalize large coefficients, in the form of penalty
 \therefore also solves overfitting

Example: Price = 1000 · Size - 500 · Age + Noise

↓ after L₂

Price = 800 · Size - 300 · Age + Less Noise

Bias-Variance Tradeoff in Ridge Regression



≈

- As $X^T X$, 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

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

- ↳ Generalized CV

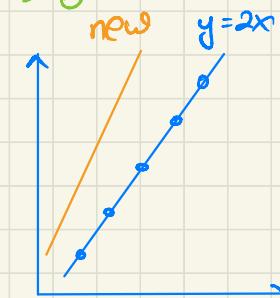
eg:	x_i^o	y_i^o	Prediction
	1	2	2
	2	4	4
	3	6	6
	4	8	8

$$\begin{aligned} \therefore Y &= \beta_0 + \beta_1 X \\ \therefore Y &= 0 + 2X \\ \therefore Y &= 2X \end{aligned}$$

$$\therefore \text{Cost function} = 0 + \lambda \cdot \text{slope}^2 = 0 + 1 \cdot 2^2 = 4$$

\uparrow
 $\therefore \text{change } \text{penalty} \text{ coefficients}$

$$\begin{aligned} \therefore \hat{y} - y_i^o &= 0 \\ (\hat{y} - y_i^o)^2 &= 0 \end{aligned}$$



L1 Regularization

- shrinks certain parameters to 0, making them obsolete

- choosing one feature from group of correlated features; in the real world, some features are highly correlated
(e.g.: year house was built and # rooms in the house)

- reduces model complexity

- robust to outliers

L2 Regularization

- makes certain parameters small but not 0 so that they still have some influence
- not robust to outliers