ordinary least-squares regression

$$Y = X\beta + \varepsilon$$

$Y$ is a $n \times 1$ vector of results

$\beta$ is a $d \times 1$ vector of regression coefficients

$X$ is a $n \times d$ matrix of predictors

$\varepsilon$ is $n \times 1$ vector of errors

$$\hat{\beta}_{OLS} = \min \sum_{i=1}^{n} (Y_i - X_i\beta)^2$$

$$= (Y - X\beta)^T (Y - X\beta)$$

dimension of $\beta$ starts to get large,

or we start to see large correlation in $X^T X$

$(X^T X)^{-1}$ is ill-conditioned

our estimates of $\beta$ start to get very noisy

and don't work well at

prediction

$$\hat{\beta} \sim N\left(\beta, \left(X^T X\right)^{-1} \sigma^2\right)$$

$\sigma^2$ is the variance of the residuals

$$\beta = (0.5, 0.5)$$

$$\hat{\beta} = (3.7, -2.4)$$

minimize
$\beta$
$$\underbrace{(Y - X\beta)^T (Y - X\beta)} + \lambda \beta^T \beta$$
$$\hookrightarrow \sum \beta_i^2$$

How do we pick $\lambda$?

$$(\vec{X_1}, Y_1), \ldots, (\vec{X_n}, Y_n)$$

$\hookrightarrow$ split my data into a test
set
and a train set

test set is about 10-20% of
the data

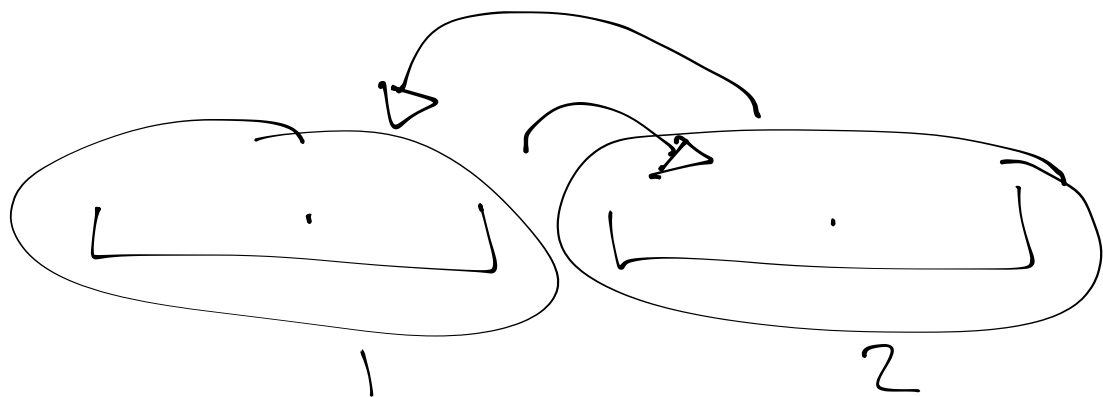$\hookrightarrow$ look at a large spread of values of
$\lambda$
$$[.01, .02, .04, .08, .16]$$

↳ cross-validation (k-folds)

    ↳ randomly partition my train set
      into k parts folds

    ↳ goes to my model for each $\lambda$
      in my grid k times .

S-fold cross-validation



        1          2

↳ choose the lambda that gives me
the best out of sample performance

$$\lambda = [1, 2]$$

fold 1

fold 2

$m_{ij}$ = model fit using lambda[i]
                        fold[j]

$m_{11}$     $m_{12}$

$m_{2.}$     $m_{22}$

evaluat

fit $m_{11}$ on fold 1 → fold 2

$\varepsilon_{11}$ is good to be prediction on

fold 2 of $m_{11}$

avg prediction error $\boxed{\lambda = 1}$

$$\left( \frac{\varepsilon_{11} + \varepsilon_{12}}{\#\text{ of samples}} \right.$$

$$\lambda = 2$$

$$\left\{ \frac{\varepsilon_{21} + \varepsilon_{22}}{\#\text{ of samples}} \right.$$

→ let's $\lambda = 1$ performs better, then I fit the model using $\lambda = 1$ on the entire train set, and validate it's performance on the test set

$$\hat{y}_i = X_i \hat{\beta}$$

$$R^2 = 1 - \frac{\Sigma (y_i - \hat{y})^2}{\Sigma (y_i - \bar{y})^2}$$

general question is how to choose $k$?

sklearn default set 5

ridge regression

$$\underset{\beta}{\text{minimize}} \quad (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$$

$$\hookrightarrow \sum \beta_i^2$$

LASSO

$$\underset{\beta}{\text{minimize}} \quad (Y - X\beta)^T (Y - X\beta) \boxed{+ \lambda \sum_{i=1}^{n} |\beta_i|}$$

$\Longrightarrow$ not a major difference th predictive
abilities of LASSO/ridge

$\Longrightarrow$ lasso is slower to fit

$\Longrightarrow$ lasso can perform variable selection

$$\lambda = BIG$$