kokchun giang

when your model is too complex, you can **regularize** it to make it simpler

many features leads to many parameters, which makes the model too complex

$$\begin{pmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_n \\ | & | & & | \end{pmatrix}$$

each column is a feature

want to $\quad \uparrow$ bias
$\quad\quad\quad \downarrow$ variance

many features $\longrightarrow$ model has

many parameters

$$\left( y = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_n x_n \right.$$

$\rightarrow$ risk for multicollinearity

$\rightarrow$ redundancy in features
$\rightarrow$ inaccurate estimate of $w_i$
$\rightarrow$ overfitting

add a **penalty** to decrease the role of parameters
other than the bias

**ridge regression**
$(l_2$-regularization$)$

$\uparrow \lambda \Rightarrow \downarrow$ variance
$\uparrow$ bias

$\uparrow$
penalty parameter

$0 \leq \lambda \leq 1$

$\uparrow \lambda \Rightarrow w_i$ close to $0$
$i \geq 1$

**lasso regression**
$(l_1$-regularization$)$

$\uparrow \lambda \Rightarrow \downarrow$ variance
$\uparrow$ bias

$0 \leq \lambda \leq 1$

$\uparrow \lambda \Rightarrow$ least important
features set to zero

**elastic net**

combines $l_1$ & $l_2$

**Hyper-parameters**

$\lambda$ - penalty

$\alpha$ - $l_1$-ratio
$\uparrow$
$0 \leq \alpha \leq 1$

# regularisation models require **feature scaling**

the models are trained
using numerical approach
such as gradient descent

$\Rightarrow$ require feature scaling
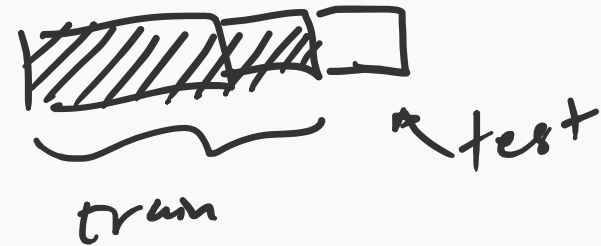
e.g. · feature standardizat

· normalization

# **k-fold cross validation** for hyperparameter tuning

| train | test | val |
|---|---|---|

1. train w. different hyperparameters

2. predict & evaluate on validation data

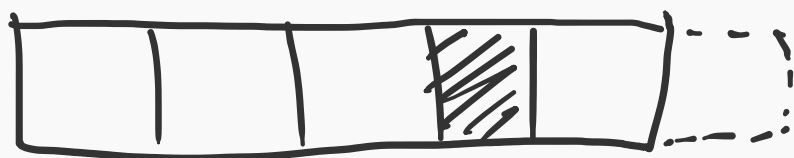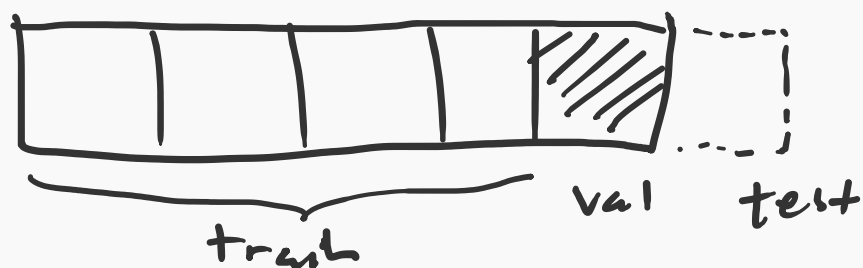3. choose new values on the hyperparameters & repeat 1. & 2.

4. choose the value on hyperparam that gave least validation error
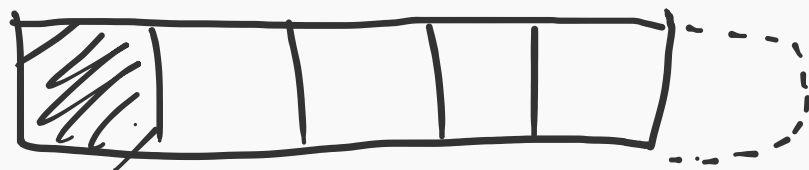
5. train on train & val datasets



train → test

6. evaluate on test set

# **k-fold cross validation** for hyperparameter tuning



val   test

train

5-fold CV

compute mean
error

for small datasets, we
utilises the data well

for larger it costs as
we repeat training k times