

∴ Ridge regression is used to introduce bias to the data in order to generalize the data and increase bias.

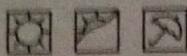
This is useful if you don't have much training data.

Lasso Regression: (L1 regularization / L1 Norm)

It is used to reduce the features. It helps in feature selection.

Cost function

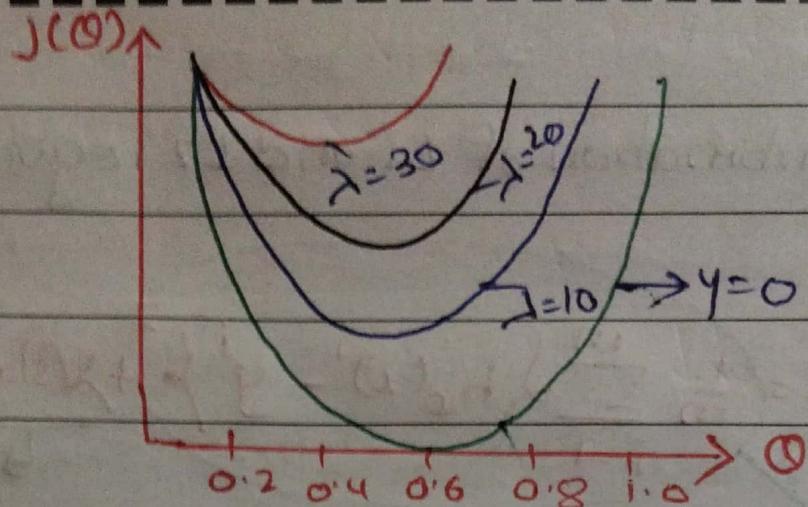
$$\text{cost function} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=1}^m |\text{slope}|$$



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /



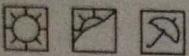
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= \theta_0 + 0.54x_1 + 0.23x_2 + \underline{0.1x_3}$$

↳ least correlated
IF data has outliers \rightarrow use ridge regression

Lasso = least absolute shrinkage and selection operator regression

Lasso regression tends to eliminate the weights of the least important features by setting their weights to zero.



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

Elastic Net:

→ It is the combination of L1 and L2 regularization.

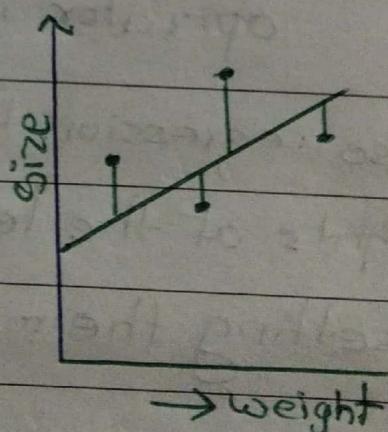
$$\text{Cost Function} = \frac{1}{m} \sum_{i=1}^m \{h_0(x)^i - y^i\}^2 + \lambda (\text{slope})^2$$

$$+ \frac{\lambda |\text{slope}|}{L \rightarrow L1}$$

can be changed
to MAE, RMSE, MSG

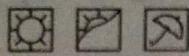
Note taken from John Stainer

We find the line that
results in the minimum
sum of squared residuals.



∴ We end up with the eqn's of line

$$\text{Say, } \text{size} = \underbrace{0.95}_{\text{slope}} + \underbrace{0.95 \times \text{weight}}_{\text{slope}}$$



Mo Tu We Th Fr Sa Su

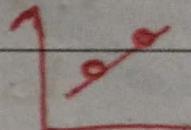
Memo No. _____

Date / /

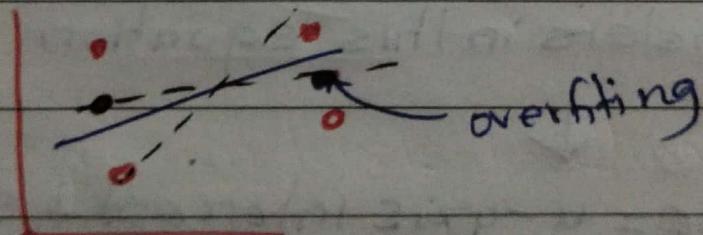
When we have a lot of measurements, we can be fairly confident that least square line accurately reflects the relationship between size and weight.

But what if we only have two measurements?

→ we fit new line, since the new line overlaps the two data points, the minimum sum of squared residuals = 0.



New line eq's:-



Sum of the squared residuals for testing data is large which means the new line has high variance.

∴ In ML, new line (blue) is overfit to training data.

Bias → amount that a model's prediction differs from the target value compared to training data.

Memo No. _____

Date / /

Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

∴ The main idea behind Ridge Regression

is to find a new line that doesn't fit the training data as well.

In other words, we introduce a small amount of bias into how the new line is fit to the data but in return for that small amount of bias, we get a significant drop in variance.

∴ Ridge regression can provide better long term prediction.

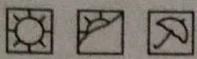
When least square determines values for the parameters in this equation



$$\text{Size} = y - \text{axis intercept} + \text{slope} \times \text{weight}$$



It minimizes the sum of the squared residuals.



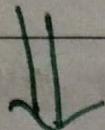
Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

In contrast,

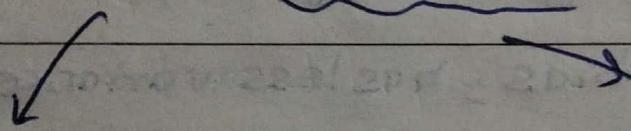
When ridge regression determines values for the parameters in this equation



$y = \text{y-axis-intercept} + \text{slope} \times \text{weight}$

It minimizes the sum of squared residuals

+ $\lambda * \underline{\text{slope}}^2$

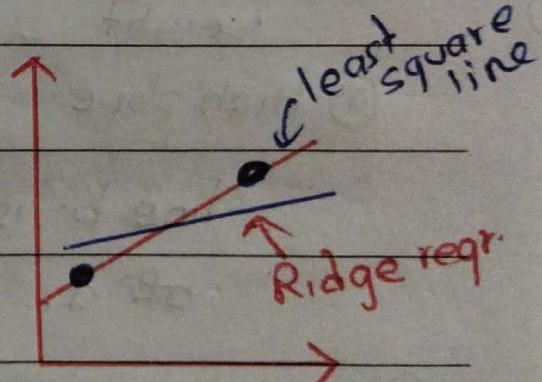


lambda (λ) determines how severe that penalty is

this part adds a penalty to the traditional least squared method.

for red line:

$$0 + \lambda(1.3)^2 = 0 + 1 \times 1.3^2 \\ = 1.69$$



for blue line

$$(0.3)^2 + (0.1)^2 - \lambda(0.8)^2 \\ = 0.34$$



Memo No. _____

Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

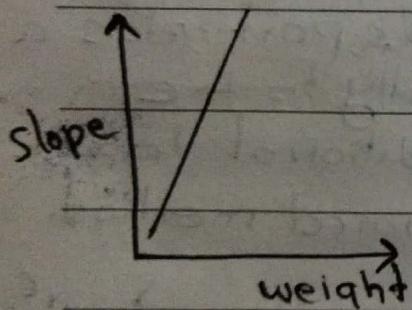
Date / /

Ridge Regression line:-

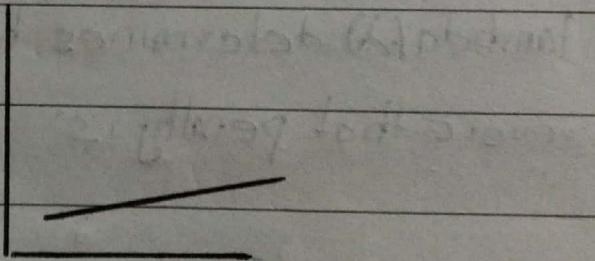
$$\text{Red} = 1.69, \text{ blue} = 0.74$$

Thus, if we wanted to minimize the sum of the squared residuals plus the Ridge Regression penalty, we would choose the Ridge Regression line over the least square line.

→ Due to penalty in Ridge regression, small amount of bias, has less variance.



@ High slope



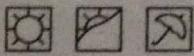
⑥ low slope

• line b. is less sensitive to weight than of a.

Ridge Regression (RR):

$$\text{Sum of squared residual} + \lambda(\text{slope})^2$$

→ λ can be any value from 0 to positive



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

infinity.

When $\lambda = 0$,

Ridge regression line = least square line

when $\lambda = 1$,

RR ended up with a smaller slope than the least square line.

and for larger λ ,

the slope of RR gets asymptotically close to 0.

So, the larger λ gets, our prediction for size becomes less and less sensitive to weight.

So, how do we decide what value to give λ ?

→ We just try a bunch of values for λ and use cross validation, typically 10 fold.

cross validation to determine which one results in the lowest variance.

"Uptill now RR was for continuous variable

but RR works for discrete variable as well.



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

Ridge Regression can also be applied to logistic regression.

$$= \text{the sum of likelihoods} + \lambda(\text{slope})^2$$

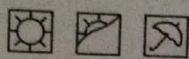
Note:

When applied to logistic Regression, Ridge Regression optimizes the sum of the likelihoods instead of the squared residuals because logistic regression is solved using maximum likelihood -

"Ridge Regression helps reduce variance by shrinkage parameters and making our predictions less sensitive to them."

Summary:

When sample sizes are relatively small, RR can improve predictions made from new data (ie reduce variance) by making the predictions less sensitive to the training data.



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

RR punts itself by λ times the sum of all squared parameter, except for the y-intercept and λ is determined using cross validation.

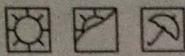
Lasso Regression (L_1):

Ridge regression penalty = $\lambda * (\text{slope})^2$

Lasso Regression = $\lambda * |\text{slope}|$

Big difference between Ridge and Lasso Regression is that Ridge Regression can only shrink the slope asymptotically close to 0 while lasso regression can shrink the slope all the way to 0.

LR can exclude useless variables from eq's, better than RR at reducing the variance in models that contain a lot of useless variable.



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / / - -

Elastic Net Regression:

Elastic Net Regression starts with least squares then combines the lasso Regression penalty with the Ridge Regression penalty.

sum of the squared residual

+

$\lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_n|$

+

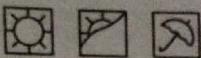
$\lambda_2 \times (\text{variable}_1)^2 + \dots + (\text{variable}_n)^2$

Note: LR and RR penalty get their own λ s.

The hybrid elastic net regression is especially good at dealing with situation when there are correlations between parameters.

By combining LR and RR,

→ Elastic Net Regression groups and shrink the parameters associated with the correlated variables and leaves them in eqⁿ or removes them all at once.



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

Assumptions of linear Regression:-

① Linearity:-

→ linear relationship between Y and each X.

② Homoscedasticity:

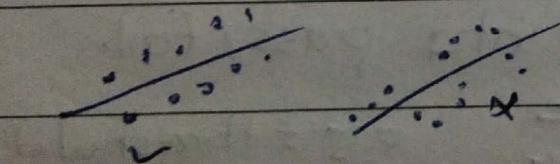
→ Equal variance.

③ Multivariate Normality

→ Normality of error distribution

④ Independence.

→ of observation, includes no auto-correlation.



⑤ Lack of multicollinearity:

Predictors are not correlated with each other.

$$\underline{x_1 + x_2}$$

$$x_1 \vee x_2$$