# Table Of Contents

# (*) Multi collinearity :-

→ we can define multi-collinearity as the situation where the independent var have strong correlation among themselves.

→ the coeff. in a linea Reg model represent the extent or change in 'y' when a certain $n$ ($x_1, x_2, x_3 \sim$ ) is changed keepup other constant. But if $x_1$ & $x_2$ are dependent then this assumption itself is wrong that (we are) changing one vars keepup other constan, as the dependent var will also be changed.

→ It means model becomes a bit flawed--

$$y = m x_1 + n x_2 + c$$
while,
$$x_2 = a x_1 + d$$

Independent var ($x_2, x_1$) are related.

more than one linear eqn,
Hence called multi collinearity

→ we have redundancy in our model as two variables
(or more than two) are trying to convey the
same information.

☆ → As the extent of collinearity increases, there is a
chance that we might produce an "overfitted model".

→ An overfitted model works well with the test data
but its accuracy fluctuates when exposed to
other data sets.

⊙ can result in a Dummy variable trap.

→ By the heatmap we can check collinearity
Generally a correlation greater than 0.9 (or)
less than -0.9 are to be avoided.

④ variance Inflection Factor (VIF) ↳
Regression of one X var against other
X variables.

$$VIF = \frac{1}{1 - R^2}$$    $\boxed{VIF < 5}$ ✓ perfect

→ at VIF > 5, extreme correlation & avoided.

Remedies for multi-collinearity

① Do-Nothing :- at corr not extreme / the var not used then ignore

② Remove one-var like in dummy var trap (one hot encoding)

③ combine correlated var← combine variables.

④ Principle component Analysis (PCA)

# ✱ Regularization:

→ when we we Regr. models to train some data, there is a good chance that the model will overfit the given training dataset.

> → Regularization helps sort this overfitting problem by restricting the degrees or freedom of given eq.

i.e)

Simply reducing the no. of degrees or a polynomial by reducing their corresponding weights.

### Reson for Regularization

→ In a linear eqn, we don't want huge coeff, as a small change in coeff can make a large diff. So regular constraints weights or such features to avoid overfitting.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots \beta_p \cdot x_p$$

> $$RSS = \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j \cdot x_{ij} \right)^2 \quad (y_i - \hat{y})$$

> ✱ To Regularize a model, Shrinkage penalty is added to cost function.

# LASSO ( Least Absolute Shrinkage & Selection Operator) (L1 form)

→ LASSO regression penalizes the model based on the sum of magnitude of the coefficients.

The regularization term is given by,

$$\text{regularization} = \lambda * \Sigma |\beta_j|$$

— shrinkage factor

→ Loss after regularization :-

new loss func.

$$RSS + \lambda \cdot \sum_{j=1}^{p} |\beta_j| = \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{s=1}^{p} \beta_j \cdot x_{ij} \right)^2 + \lambda \cdot \sum_{j=1}^{p} |\beta_j|$$

Here, Loss is not calculated just by previous values $(y_i - \bar{y})$, now it has been split into several parts which decreases loss.

## Ridge Regression (L2 form)

→ Ridge Regression penalizes the model based on the sum of squares of mag of coeff.

$$\text{regularization} = \lambda * \Sigma |\beta_i^2|$$

$\lambda$ - can be calculated by cross validation

→ Loss after regularization,

$$RSS + \lambda \cdot \sum_{j=1}^{p} \beta_i^2 = \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j \cdot x_{ij} \right) + \lambda \cdot \sum_{j=1}^{p} \beta_i^2$$

Note :- consider $\beta_1$ & $\beta_2$ be coeff of LR & $\lambda = 1$.

① For LASSO, $\boxed{\beta_1 + \beta_2 \angle = S}$ → where 'S' is the man value the eq. can achieve. At

② for Ridge, $\boxed{\beta_1^2 + \beta_2^2 \angle = S}$

## Elastic net

→ middle ground b/w Ridge & Lasso.

→ The regularization term is a simpler mix of
  both Ridge & Lasso, and you can control
  mix ratio α.

$$L_{enet}(\hat{\beta}) = \frac{\sum_{i=1}^{n}(y_i - x_i\hat{\beta})^2}{2n} + \lambda\left[\left(\frac{1-\alpha}{2}\right)\cdot\sum_{j=1}^{m}\hat{\beta}_j^2 + \alpha\sum_{j=1}^{m}|\beta_j|\right]$$

→ where 'α' is the missing parameter

  Ridge $(\alpha = 0)$

  Lasso $(\alpha = 1)$ //

---

(✲) **Difference between Ridge & LASSO**

→ Ridge Regression shrinks the coeff for those
  predictors which contribute very less in the
  model but have huge weights, very close to '0'.

  | $y = 0.04 x_1 + ... x_1 e_4$ |
  | $\beta^2 = 0.0016 \rightarrow$ very Less |

  | thus final model will still contain
  | all those Predictors, though
  | with less weights.

→ where as in Lasso, the L1 penalty does reduce
  some coeff exactly to zero, when we use a
  sufficiently large tuning parameter b.
  so, in addition to regularizing,
  Lasso also performs feature selection.

**Q~ Why use Regularisation?**

A~

→ It helps to reduce the variance of the model, without a substancial increase in the bias.

→ If there's variance in the model that means the model won't fit well for dataset other than training data (which is called overfitting).

→ the tuning parameter (λ) controls this bias & variance trade off. selected my cross-validation

→ when 'λ' is increased to certain limit, it reduces the variance without losing any imp prop of data.

→ But after certain limit, model will start losing imp prop which will increase (bias) in the data.

---

**Q~ what should be used plain linea, Ridge, Lasso or Elasticnet?**

A~

→ It is always preferrable to have atleast a little bit of Regularizadion, so generally avoid plain linea Regv.

→ Ridge is a good default,

→ But if you suspect that only a few features are actually useful, you should prefer Lasso/EN, since they tend to reduce the useless features weights down to zero.

→ In general Elastic net is preffered over Lasso since Lasso may behave eratically when the no. of features is greater than no. of training instances or when several features are strongly correlated.

whenever, mean is not 0,
its better to do standard scaler

Note:-

① standard Scaler() $\longrightarrow \dfrac{x - \bar{x}}{SD}$   $\bar{x} = 0$
$SD = 1$

② Two ways to find multicollinearity → | correlation &
| VIF

③ 4 ways to remove multicollinearity.

④ train set, high accuracy ; test set, less accuracy.

(e) overfitting, so regularize it.

→ Even after checking Lasso and Ridge,
it it gives the same accuracy values,
the you can it is not over-fitted.

(e) Just by having huge difference in
train accuracy & test accuracy than it
doesn't just mean overfitted. test before
judges. In this case model is not problem,
more data is needed.

⑤

The linear, in linear Regression does not talk
about the degree of Polynomial eq in
terms of dependent var(x).

Isted, it talks about the degree or coeff.

$$ y = a + bx + cx^2 + \dots nx^n $$

its not talks about power of x,
but the power of a, b, c etc. And their powers.
Hence it is linear Regression

# Polynomial Regression

→ or is a mechanism to predict a dependent var based on the polynomial relationship with the independent variable.

in the eq, $y = a + bx + cx^4 \cdots n x^n + \cdots$

man. power of $x$ is called degree of polynomial eq.

Deg ① → $y = a + bx$

Deg ② → $y = a + bx + cx^2 \sim$

## when to use?

⊛ If data is scatter in a polynomial curve shape (not linear) then keep trying different degree until the line fits to data

## Codes

```
from sklearn.preprocessing import PolynomialFeatures
Poly_reg = PolynomialFeatures (degree=2)
X_poly = Poly_reg.fit_transform(x)
```

↑keep changing degree until it fits to data