

Linear Regression

- Linear Regression
 - Theory of linear Regression
 - Simple Implementation with python
 - scikit-learn overview
 - Linear Regression with Scikit-learn
 - polynomial Regression
 - Regularization
 - Overview of project datasets.

History and Motivation behind
linear regression:

- Brief History, linear Relationship
- Ordinary least squares, cost function
- Gradient descent - Vectorization

The history invention



Navigational methods (1700s)

Age of exploration -

- 1722 - Roger Cotes - combining different observation yields better estimates of true value -
- 1750 - Tobias Mayer explore averaging different result under similar condition in studying librations of Moon

1754 - Roger Joseph Boscovich

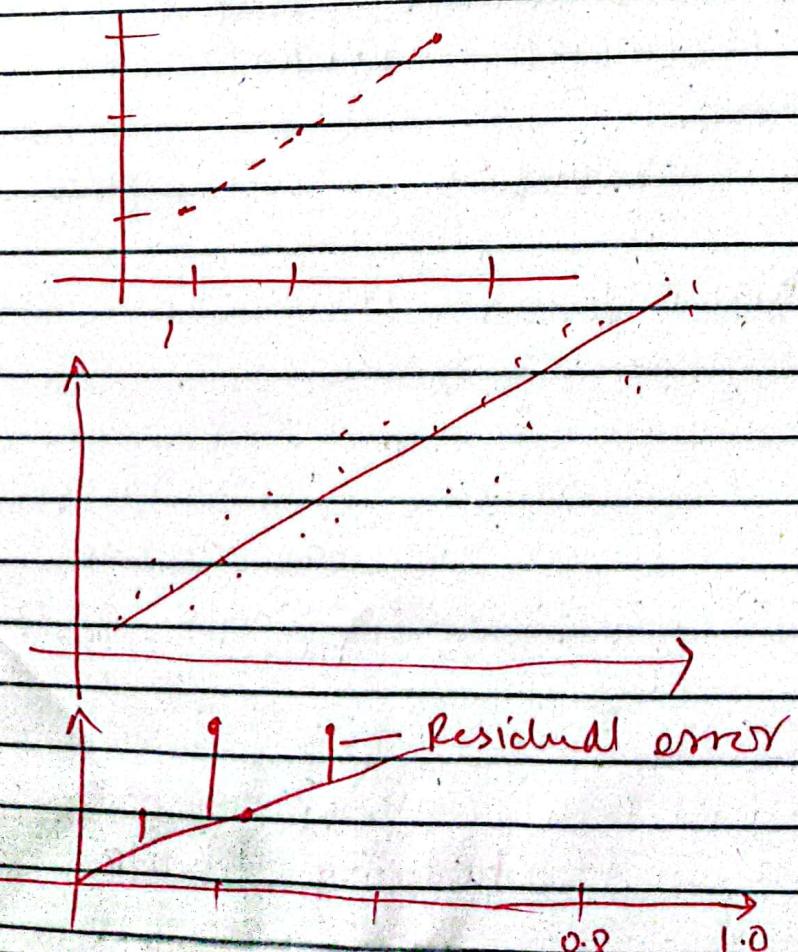
1788 - Pierre-Simon Laplace develops

→ 1805 - First public exposition on
Linear Regression with least
square method -
published by Andrien-Marie Legendre -

• 1809 - Carl Friedrich Gauss publish
methods of calculating orbits of
celestial bodies

- claimed to have invented least-square
back in 1795! -

born in 1977



- Ordinary least squares works by minimizing the sum of the squares of the difference between the observed dependent variable (value of the observed) in

dependent variable (value of the observed) in the given dataset and those predicted by the linear function.

- Minimize the squared error

Equation :-

$$y = mx + b \rightarrow \begin{matrix} \downarrow & \text{Intercept} \\ \text{slope} \end{matrix}$$

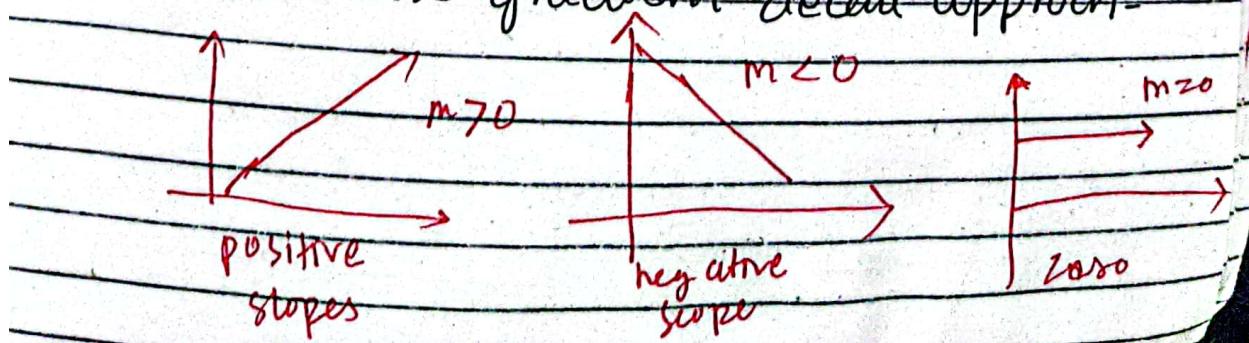
| x_1 | x_2 | x_3 | y |
|---------|---------|---------|-------|
| x_1^1 | x_2^1 | x_3^1 | y^1 |
| x_1^2 | x_2^2 | x_3^2 | y^2 |
| x_1^3 | x_2^3 | x_3^3 | y^3 |
| x_1^4 | x_2^4 | x_3^4 | y^4 |

$$\hat{y} = \beta_0 x_0 + \dots + \beta_n x_n$$

$$\hat{y} - \sum_{i=0}^n \beta_i x_i$$

\hat{y} is estimation
(can't fit)

- For single feature we can use analytical approach and for multiple features we can use gradient descent approach.



How to derive bisstandab. of

$$b_1 = \rho_{xy} \frac{S_y}{S_x}$$

$\rho_{x,y}$ = Pearson Correlation Coefficient

S_x, S_y = Standard deviation

$$= \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} \cdot \sqrt{\frac{\sum (y - \bar{y})^2}{n}}$$

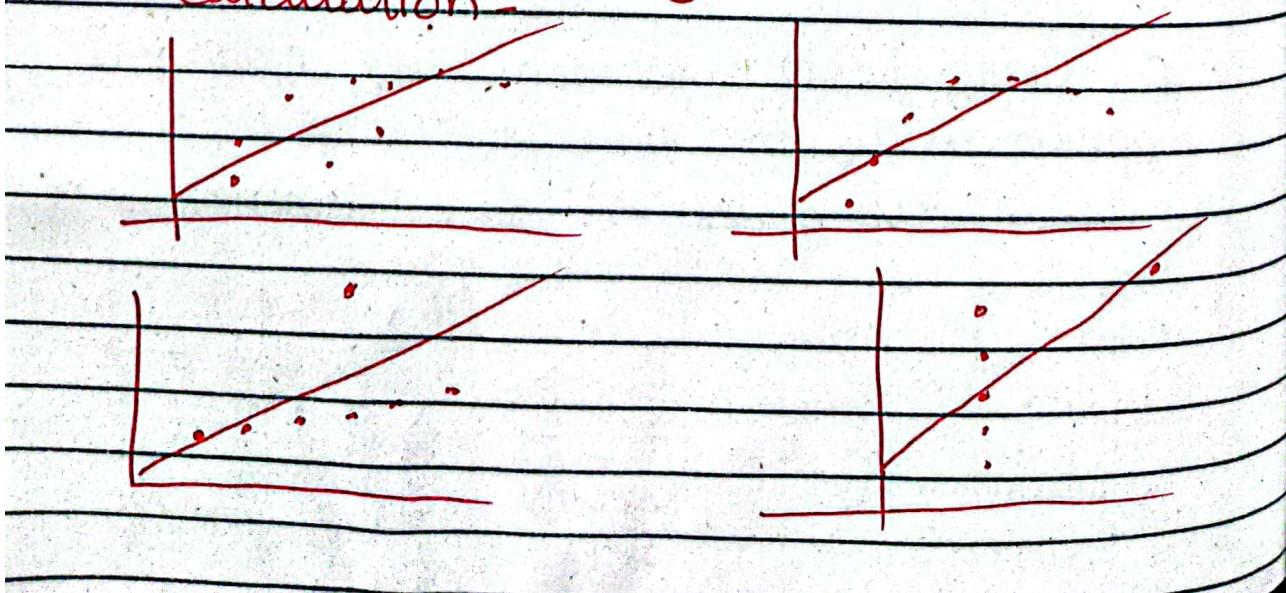
$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

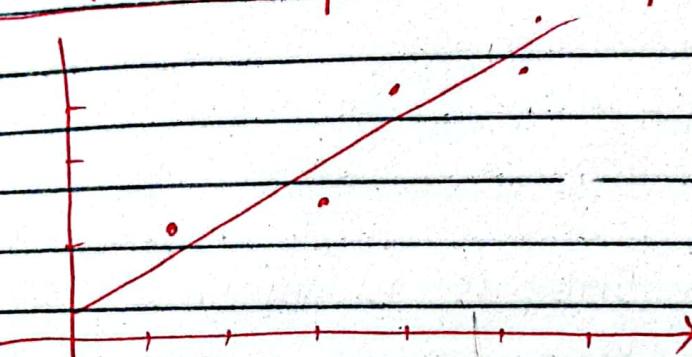
Limitation of linear Regression

Anscombe's Quartet illustrate the pitfalls of relying on pure calculation.



~~find~~
find the plot data

| production hours | production volume |
|------------------|-------------------|
| 34 | 102 |
| 35 | 109 |
| 39 | 137 |
| 42 | 148 |



| $(x - \bar{x})$ | $(y - \bar{y})$ | $(x - \bar{x})(y - \bar{y})$ | $(x - \bar{x})^2$ |
|-----------------|-----------------|------------------------------|-------------------|
| 7 | 24 | 168 | 49 |
| 8 | 31 | 256 | 64 |
| 12 | 43 | 516 | 144 |

$$b_1 = \frac{\sum (x_i - \bar{x})(y - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{558}{128} = 4.5$$

$$y = -46 + 4.5x$$

For multiple features of ~~independen~~
variable it will be hard to
use Ordinary least square
methods.

- We will use gradient descent approach to minimize the ~~error~~ cost function -

$$\hat{y} = \sum_{i=0}^n \beta_i x_i$$

Best fit line is to minimize the square error -

Residual Error

$$= \frac{1}{m} \sum_{j=1}^m (y_j - \hat{y}_j)^2 \quad \text{Average squared error.}$$

A cost function is defined by some measure of error

$$J(\beta) = \frac{1}{(2)m} \sum_{j=1}^m (y_j - \hat{y}_j)^2$$

Just for convient 2 is gonna cancel 2 - when taking derivate

$$J(\beta) = \frac{1}{2m} \sum_{j=1}^m \left(y_j - \sum_{i=0}^n \beta_i x_i^j \right)^2$$

From calculus to minimize a function you take derivate and set it equal to zero -

$$\frac{\partial J(\beta)}{\partial \beta_k} = \frac{2}{2m} \sum_{j=1}^m (y_j - \sum_{i=0}^n \beta_i x_{ij})^2$$

$$\frac{\partial J(\beta)}{\partial \beta_k} = \frac{1}{m} \sum_{j=1}^m (y_j - \sum_{i=0}^n \beta_i x_{ij})_j (-x_k^j)$$

We will describe this cost function as vectorize matrix Notation and use gradient descent to have a computer figure out the set of Beta coefficient values that minimize the cost/Loss function-

$$\nabla_{\beta} J = \begin{bmatrix} \frac{\partial J}{\partial \beta_0} \\ \vdots \\ \frac{\partial J}{\partial \beta_n} \end{bmatrix}$$

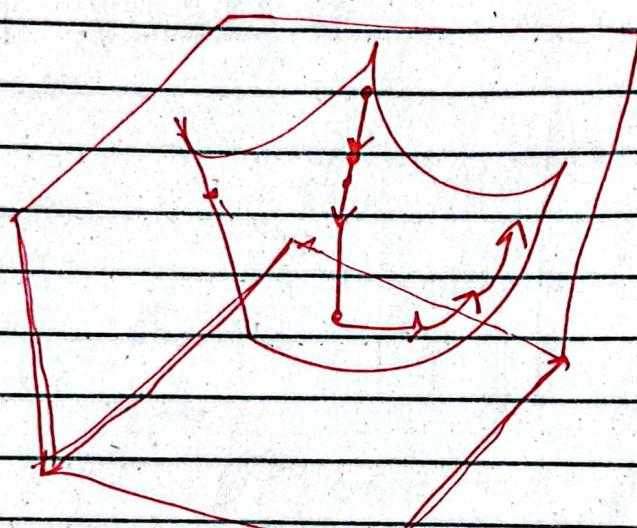
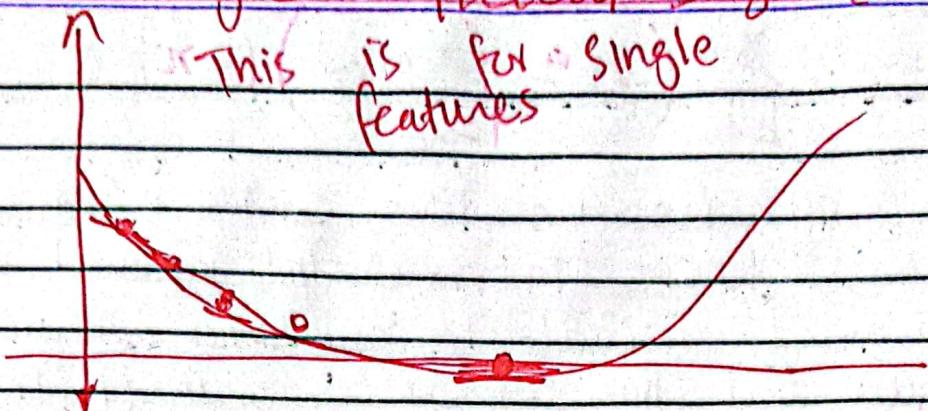
$$X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_1^m & x_2^m & \dots & x_n^m \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

$$\nabla_{\beta} J = -\frac{1}{m} \begin{bmatrix} \sum_{j=1}^m y_j x_0^j \\ \vdots \\ \sum_{j=1}^m y_j x_n^j \end{bmatrix} + \frac{1}{m} \begin{bmatrix} \sum_{j=1}^m \sum_{i=1}^n \beta_i x_{ij} x_0^j \\ \vdots \\ \sum_{j=1}^m \sum_{i=1}^n \beta_i x_{ij} x_n^j \end{bmatrix}$$

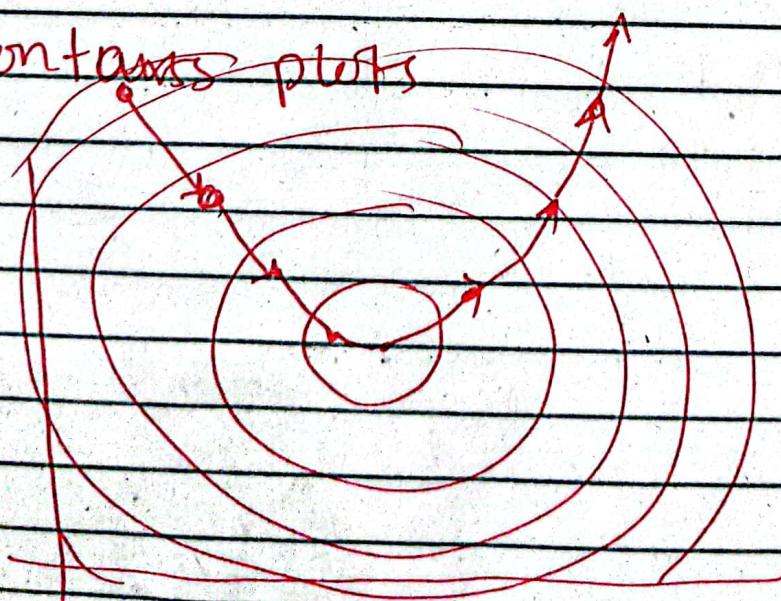
Single Coefficient Single (Beta)

This is for single features



This is for two features
and visualizing in more than
two features is not possible.

Contains plots



base of all ML models

Python Coding Simple linear

Regression :-

Chapter # 3 ISLR

$$y = mx + b$$

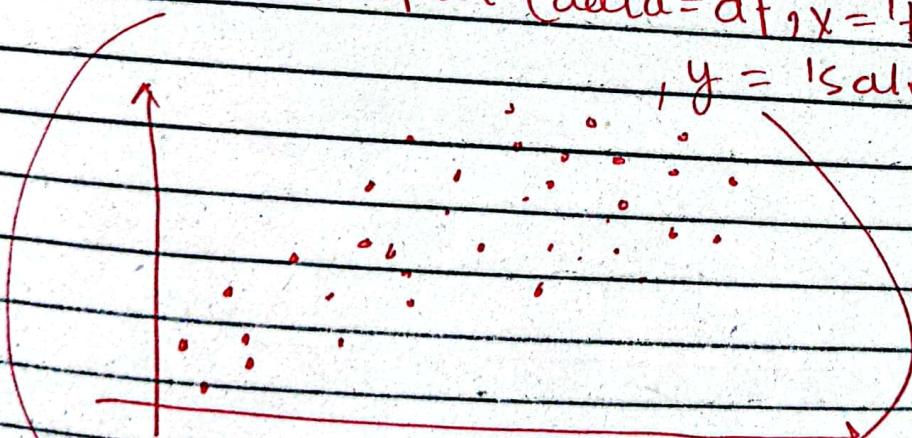
```
import numpy as np  
import matplotlib.pyplot as plt  
import pandas as pd  
import seaborn as sns
```

df.head(); 5 rows view

```
df['total_spend'] = df['TV'] + df['radio']  
+ df['sales']
```

df.head();

```
sns.scatterplot(data=df, x='total_spend',  
y='sales')
```



can used
sns.regplot(same)

X = df['total_spend']

y = df['sales']

help(np.polyfit) → help to read documentation

$y = mx + b$

$y = B_1 X + B_0$

np.polyfit(X, y, degrees=1)

array([0.0486878, 4.24302833])

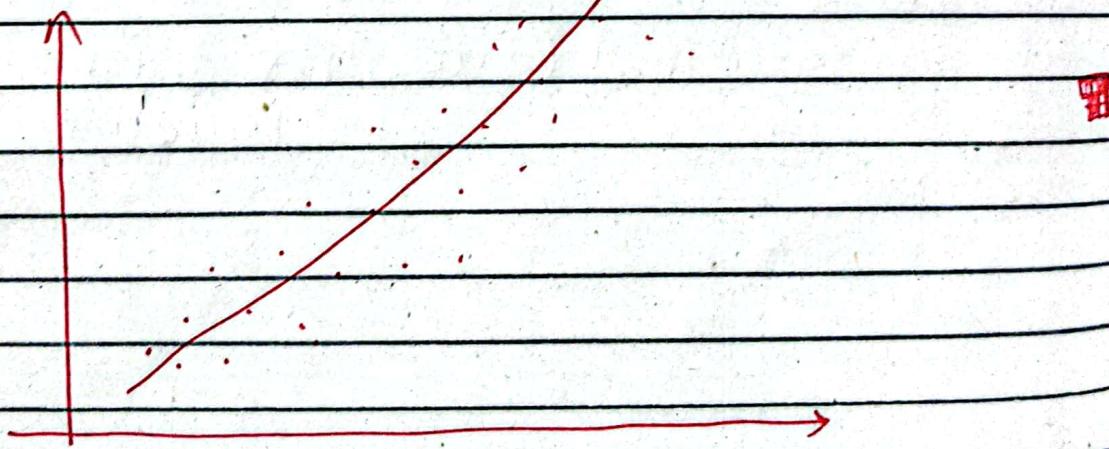
Between 0 to 500
points give me 100 points

potential_spent = np.linspace(0, 500, 100)

predicted_sales = 0.04868788 * potential_spent
+ 4.2430283

plt.plot(potential_spent, predicted_sales)

sns.scatterplot(x, same as)
above



Spent = 200

predicted_sales = 0.04868788 * spent + 4.2430283

Best.

np.polyfit(x, y, 3)

array([3.1e+00, -1.8e-01])

array([1.83e-01, 2.8e-01])

$y = B_1x + B_0$

$y = B_3x^3 + B_2x^2 + B_1x + B_0$

pot_spread = np.linspace(0, 500, 100)

Scikit-learn

- Scikit-learn is a library containing many machine learning algorithms.
- It utilizes generalized estimator framework to calling the models.
- This mean the way algorithm are imported, fitted and used in a uniform across all algorithms.
- This allows user to easily swap algorithms in and out and test various approaches.
- This uniform framework also mean user can easily apply almost any algorithm effectively without truly understanding what the algorithm is doing.
- Many convenience tools
 - train-test-split
 - metrics
 - cross-validation

Philosophy of Scikit-learn

- statsmodels - library

Code

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y)
```

from sklearn.model_selection import ModelAlgo
[example not true]

Instance of a model and
can change parameters

```
mymodel = ModelAlgo(param1, param2)  
mymodel.fit(X_train, y_train)  
predictions = mymodel.predict(X-test)
```

Linear Regression with Scikit-learn

part-one

Data-Setup and Model Training

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib as pyplot
```

```
import Seaborn as sns
```

```
df = pd.read_csv("in").
```

Q Visualization.

`x = df.drop('sales', axis=1)`

`y = df['sales']`

`x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)`

`lml(df)`

`from sklearn.linear_model import LinearRegression
help(LinearRegression)`

`model = LinearRegression(),`

`model.fit(x_train, y_train)`

`model.predict(x_test)`

Evaluation Metrics

1) Mean Absolute Error

2) Mean Squared Error

3) Root Mean Square Error

Only Apply to Regression task

Mean Absolute Error = $\frac{1}{n} \times \sum |y_i - \hat{y}_i|$

↓ problem
(will not punish when values are far off)

Mean Squared Error =

$$\frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2$$

Larger error are "punished" more
then with MAE making MSE more
popular.

• Mean Squared Error (MSE):

- Issue with MSE:

- Different units than Y

- It reports units of Y squared!

hard ↓ dollar

Root Mean Squared error = $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$

- most popular

also →

- also punish larger values
- Also will be in same
meters-

#? What is good "RMSE"

- context is everything!

A RMSE of \$10 is fantastic
for predicting the price
of a house, but horible for
predicting the price of a candy
bar.

Code #

model/F

test prediction

```
from sklearn.metrics import mean_absolute_error,
```

```
mean_squared_error =
```

```
df[ "Sales" ].mean()
```

```
Sns. histplot ( data=df , x = "Sales" )
```

bin = 20

```
mean_absolute_error ( y-test, test_predict )
```

mean_squared_error(y_true, test_prediction))

RMSE

np.sqrt(mean_squared_error(y, c))

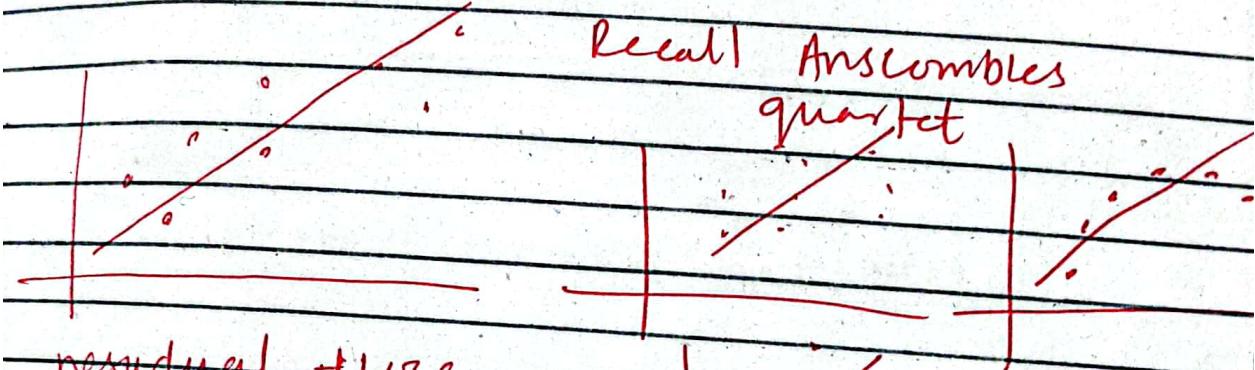
RMSE is standard deviation

MSE variance of the prediction

Evaluating Residuals

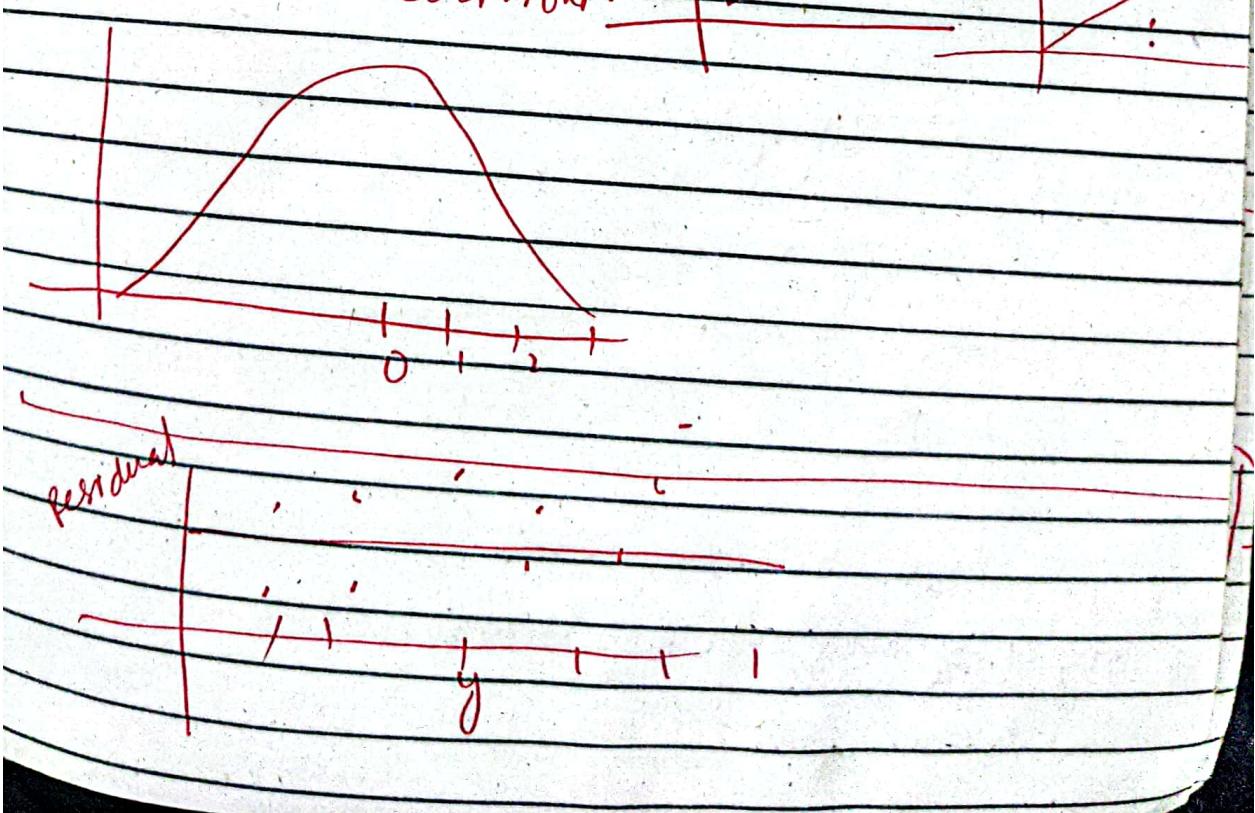
$$\text{Residual} = (y - \hat{y})$$

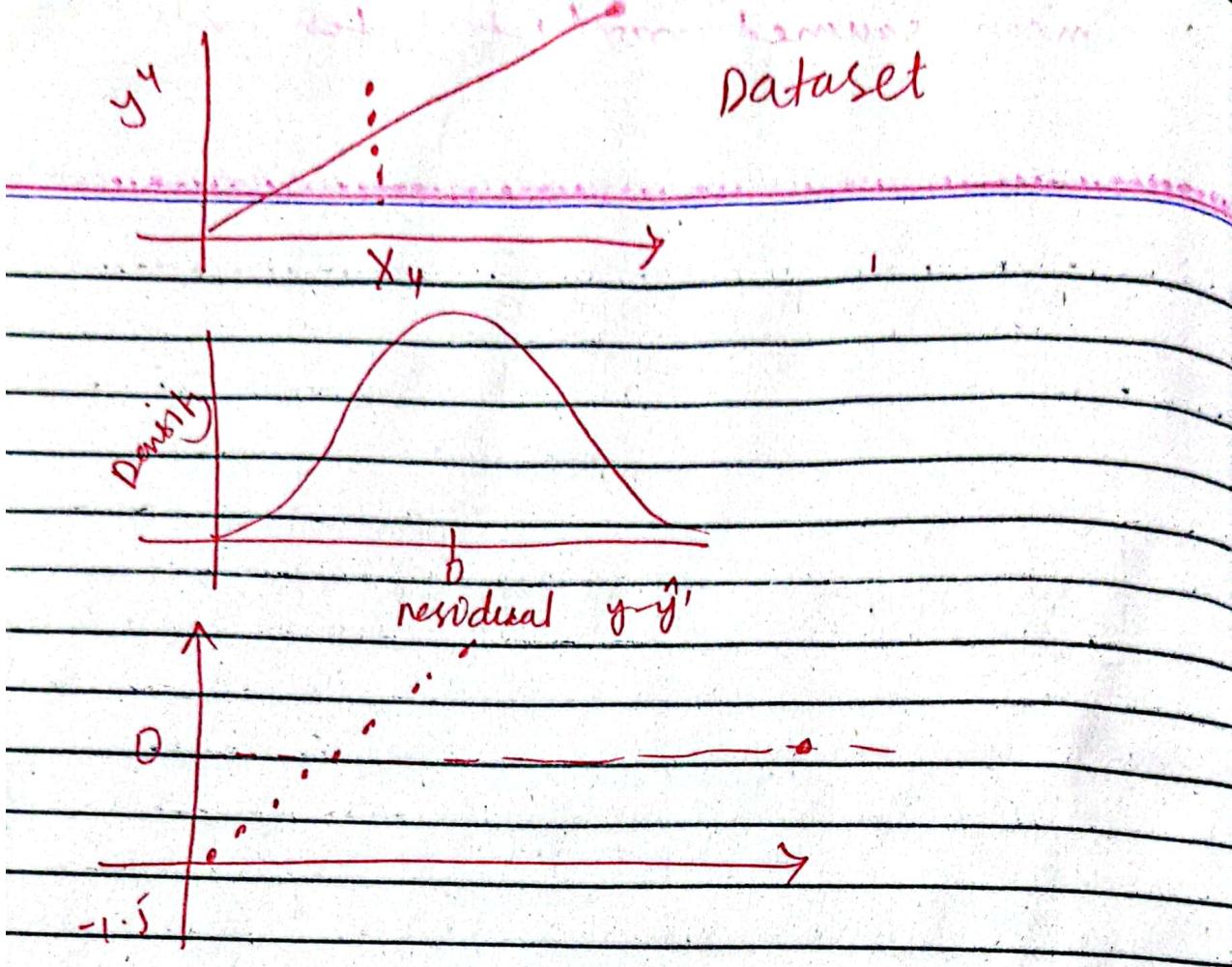
Recall Anscombe's quartet



residual - true

to normal distrib.

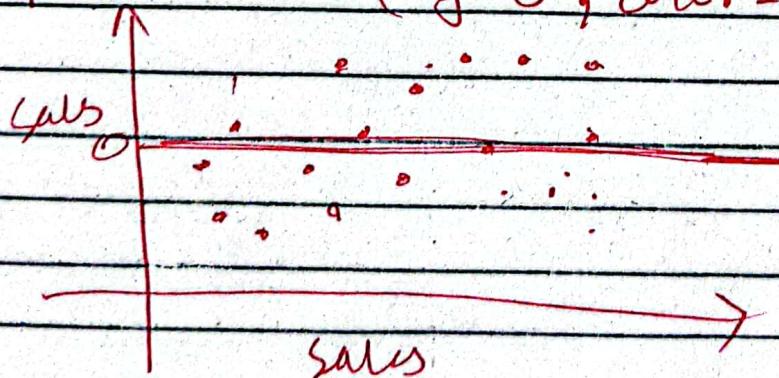




Code # Residuals

```
test_residuals = y-test_pred  
# test new
```

```
sns.scatterplot(x=y-test, y=residuals)  
plt.axhline(y=0, color='r', ls='--')
```



```
sns.densityplot(test_residuals, bins=25,  
kde=True)
```

import scipy as sp
create a figure and axis to plot on
~~fig, ax=plt.subplots(figsize=(6,8), dpi=100)~~

- # ppfpp probplot returns the raw value if needed-
- # we just want to see the plot, so we can assign these values to
 - = Sp.stats.probplot(testres^{dual}, plot=ax)

Model Deployment

Code #

```
final_model = LinearRegression()  
final_model.fit(X, y)
```

final_model.coef_

array [0.0457645 , 0.1885302 ,
- 0.00108749]

↓
Very less value
mean there is no
real relationship
between Sales
with News paper.

y-hat = final_model.predict(X)

from joblib import dump, load

dump(final, 'final_sales_model')

load_model = load('final_sales_model.joblib'),

Campaign = [[]]

149 TV, 22 Radio, 12 Newspaper

Campaign = [[149, 22, 12]]

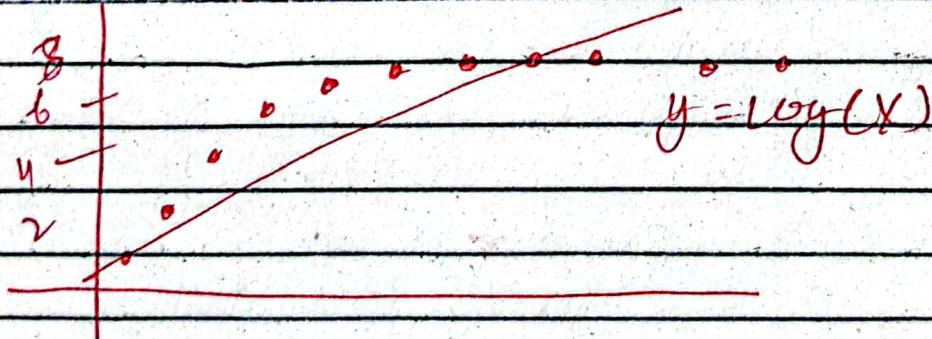
loaded model = Model predict (campaign)
predict print (loaded model)

Polynomial Regression:-

Two issues polynomial Regression
will address for us -

- Non linear feature relationship
to labels -
- Interaction terms between features.

Imagine



interaction terms :-

Synergy between original multiple
features

- polynomial features

The features created include:

- The bias (the value of 1.0)
- Value raised to the power for each degree (e.g. $x^1, x^2, x^3 \dots$)
- Interaction between all pairs of features - $x_1 \cdot x_2, x_1 \cdot x_3$

A, B

$1, A, B, A^2, B^2, AB$

$1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_2 x_1$

$x_1 = 2$ and $x_2 = 3$

1, 2, 3, 4, 6, 9

Imports

files read

$X = df.drop("Sales", axis=1)$

$y = df['Sales']$

From sklearn.preprocessing import PolynomialFeatures

polyconverter = PolynomialFeatures()

degree=2, bias=False

include-

polyconverter.fit(X)

polyconverter.transform(A)



[2.301x10⁻² 3.78x10⁻¹]

fit - first and

Train-Test-split

polynomial converter also need saved

```
from sklearn.model_selection import  
    train_test_split
```

```
X-train, X-test, y-train, y-test = train_test_split  
    (poly_features, y, test_size=)
```

```
model.fit(X-train, y-train)
```

```
test_predict = model.predict(X-test)  
# model coefficient
```

mean absolute error ($y_{\text{test}}, y_{\text{models}}$)

mean-squared error ($y_{\text{test}}, y_{\text{models}}$)

model coefficient

Bias-Variance Trade-off

also known as
Overfitting vs Underfitting

Amotek

In general,

- Increase model complexity in search of better performance leads to a Bias-Variance tradeoff.
- We want to have a model that can generalize well on unseen data - but also

account for variance and patterns in the known data.

- Extreme bias or variance can lead to poor model prediction.

high bias (underfit)

high variance (overfit)

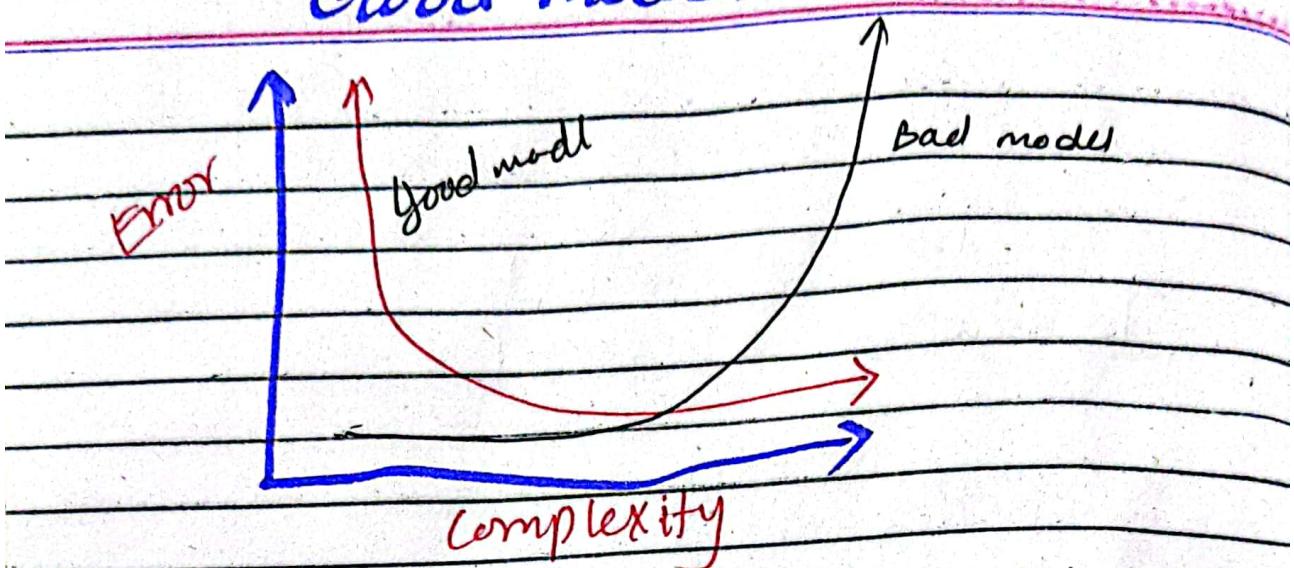
Overfitting :-

- When model fits much to the noise from the data.
 - This often result in low error on training set but high error on test/validation sets
- to much on
noise or
variance

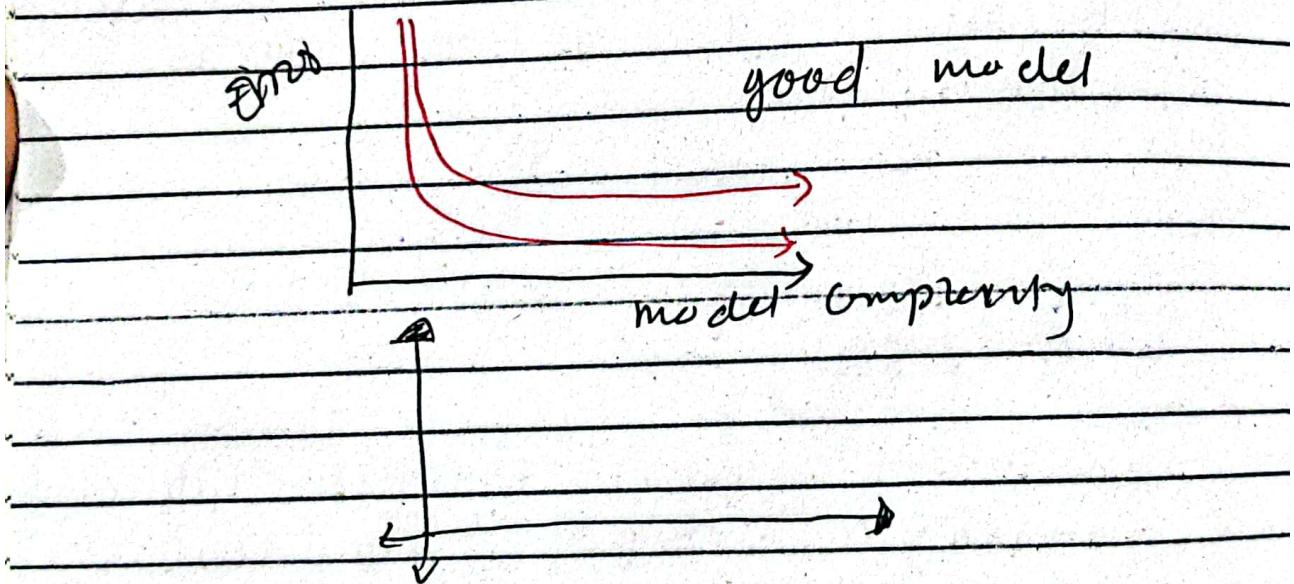
Underfitting :-

- Too much bias
- Model does not capture the underlying trends of the data and does not fit the data well enough.
- Low variance but bias
- Underfitting is often in result of an excessively simple model.
- Model has high bias and is generalizing too well.

Good Model



~~plot the error and model complexity~~



Polynomial Regression

Error Against Degree of Complexity

```
# create the different polynomial  
# split poly features. train-test split  
# stores/save the rmse for both  
train and test  
# plot the result (error vs poly ord.)
```

Train - rmse-errors = []

test - rmse - errors = []

for d in range(1, 10):

poly_converter = PolynomialFeatures(degree=d,
include_bias=False)

Poly-fcatures = poly_converter.
fit_transform()

X-train, Xtest, ytrain, ytest =

train-test split()

model = LinearRegression()

model.fit(X_train, y_train)

train - model.predict(X_train)

test pred = model.predict(X_test)

train_rmse_errors.append()

test_rmse_error.append()

Deploying polynomial Model

final poly-converter = PolynomialFeatures(degree=3, include_bias=True)

final model = LinearRegression($\text{fit_intercept} = \text{False}$)

full converted : $X = \text{final_poly_converter}.fit_\text{transform}(X)$

final model . fit(full converted, X, y)

from joblib import dump, load

dump(final_poly_converter, 'final_converter.joblib')

loaded model = load('final.joblib')

Regularization :-

- Regularization seeks to solve a few common model issue by:

- Minimizing model complexity
 - Penalizing the loss function
 - Reduce Model overfitting

(add more bias to reduce variance)

In general, we can think of regularization as a way to reduce model overfitting and variance.

- Requires some additional bias
 - Require search for optimal penalty hyperparameters

Main Regularization:-

Three
o L1 regularization
— Lasso Regression

o L2 regularization
→ Ridge Regression

o Combining L1 and L2
• ElasticNet

L1 regularization:-

Add a penalty equal to the absolute value of the magnitude of coefficients -

- o Limit the size of coefficients
- o Can yield sparse model where some coefficient can become zero.

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

L2 regularization:-

L2 regularization adds a penalty equal to the square of the magnitude of coefficient -

- o All coefficient are shrunk by the same factor
- o Does not eliminate coefficient

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

Elastic Net :-

ElasticNet combines L1 and L2 with the addition of an alpha parameter deciding the ratio between them.

$$\sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{2n} + \lambda \left(\frac{1-\alpha}{2} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right)$$

These Regularization methods do have a cost:

- Introduce an additional hyperparameter that's need to be tuned.
- A multiplier to the penalty to decide the "strength" of the penalty

Later on, we will actually cover L2 regularization

Feature Scaling

- many benefits in machine learning -
 - some machine learning models that rely on distance metrics (e.g KNN) too requires scaling to perform well
 - improves the convergences of steepest decent algorithms, which do not possess the property of scale invariance -
 - for decent algorithms helps in converge at same time -

- Those algos where Scaling won't have any effect Decision Tree, Regression Tree, Random Forest..
- important in comparing measurement that have different units -
- Allow us to directly compare model coefficients to each other -
- Features Scaling caveats:
 - Must always scale new unseen data before feeding it to model -
 - Effects direct interpretability of features coefficients -

- Easier to compare coefficients to one another, harder to relate back to original unseen features-

- Feature Scaling benefits:-

- Can lead to great increase in performances
- Absolutely necessary for some models -
- Virtually no real downside to scaling features-

Two ways to scale features:-

- Standardization :- Z-score Normalization rescale data to mean of (μ) zero and Standard deviation

Standardization of 1 -

$$X_{\text{scaled}} = \frac{X - \mu}{\sigma}$$

- Normalization:-

Rescale all data value to be between 0 - 1 -

$$X_{\text{scaled}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

Many more methods Standardization

Standardization:-

A. `fit()` in standardization will calculate a mean, and standard deviation of the dataset (statistical terms)

Normalization:-

Will calculate x_{\min}, x_{\max} of the data a. `fit()` will calculate

- A. `transform()` call actually scale the data and return the new scaled version of data

#

Very important consideration for `fit` and `transform`-

- We only fit to training data.
- Calculating statistical information should only come from training data.
- Don't want to assume prior knowledge/knowledge of the test set.

Data leakage

- Using the full dataset would cause data leakage -
 - calculating statistics from full data leads to some information of the test set leaking into training process upon transform conversion -

Feature Scaling process:-

- perform train test split
- fit to training features data
- ~~Transform test features~~
data

Do we need to scale the label

- can negatively effect
stochastic gradient descent
- not advised

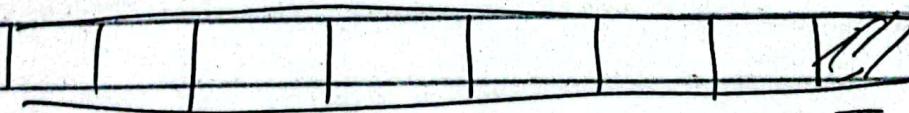
Cross-Validation:-

• Cross Validation is a more advanced set of methods for splitting data into training and testing sets-

Is there a way to:

Train on all data &

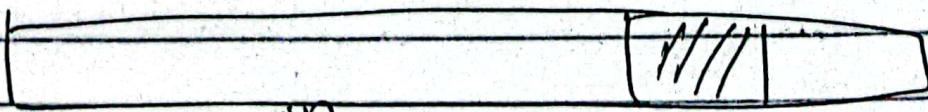
Evaluate on all data



Train

Test

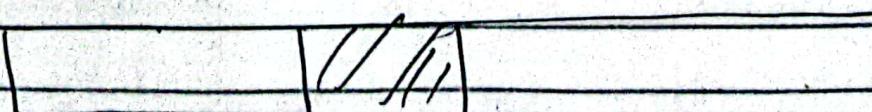
Error 1



Train

test

Error 2



Train test

Error 3

mean - with the Error

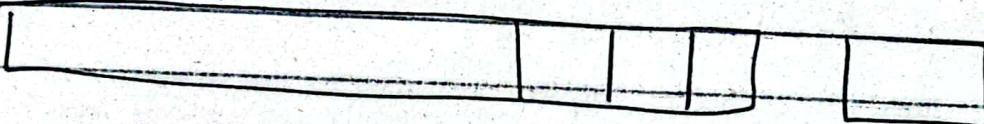
computationally costly

- This is known as k-fold cross-validation.
- Common choice for k is 10 so each test set is 10 percent of your total dataset.
- Largest k would be possible would be K equal to the numbers of rows.
- This is known as leave one out cross validation.
- Computationally expensive.

Train |

Validation

hold out test set :-



- Train | Validation | test split

Data Setup

Regularization for Linear Regression Dataset

```
Import numpy as np  
import pandas as pd  
import ...  
import ...
```

plt
sns

```
df = pd.read_csv("Adver.csv")
```

```
X = df.drop('sales', axis=1)
```

```
y = df['sales']
```

From sklearn.preprocessing
Import PolynomialFeatures

```
poly_converter = PolynomialFeatures(degree=3,  
include_bias=False)
```

```
poly_features = poly_converter.fit_transform(X)
```

(X.shape)

(200, 19)

```
from sklearn.model_selection import train
```

```
X_train, X_test, y_train, y_test = train_test_split(poly_features, y, test_size=0.2)
```

```
from sklearn.preprocessing import StandardScaler
```

scaler = StandardScaler();

```
scaler.fit(X_train)
```

$X\text{-train} = \text{scaler.transform}(X\text{-train})$
 $X\text{-test} = \text{scaler.transform}(X\text{-test})$

$X\text{-train}[0]$

Ridge Regression :-

Ridge Regression is a regularization technique that works by helping reduce the potential for overfitting to the training data.

- It does this by adding in a penalty term to the error that is based on the squared value of the coefficient.

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

we could substitute the beta coefficient

$$= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \dots - \hat{\beta}_p x_{ip})^2$$

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

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

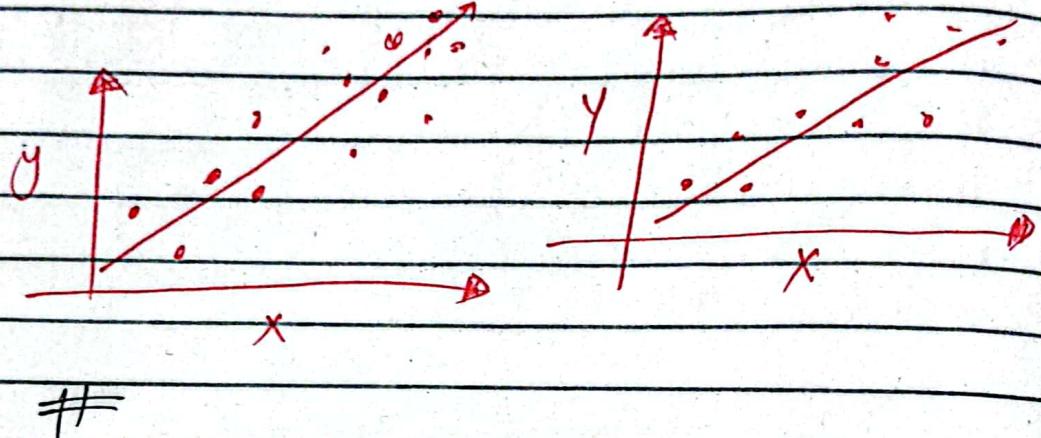
Ridge Regression seeks to minimize this entire term $RSS + \text{penalty}$

$$g \sum_{j=1}^p \beta_j^2$$

Shrinkage penalty has a tunable
lambda parameters!

$$\lambda = 0 - \infty$$

\downarrow
0-mean RSS



L2 regularization with
python and scikit-learn

- Important Note:-

for cross-validation metrics, sklearn
use a "scorer object!"

All scorer object follows the convention
that higher return values are better
than the lower Return ~~not~~ value.

- But higher RMSE is actually worse!
- So scikit-learn fixes this by using
a negative RMSE as its scorer
metrics-
- This allows uniformly across all
scorer metrics, even across different
task types-

The same idea of uniformity across ~~all~~
model classes applied to referring
to penalty strength parameters as
alpha

```
from sklearn.linear_model import Ridge  
# help(Ridge)
```

```
ridge_model = Ridge(alpha=10)
```

```
ridge_model.fit(X_train, y_train)
```

test-prediction = ridge-model.predict(X-test)

MAE = mean_absolute_error(y_test, pred)

RMSG = $\sqrt{\text{MSE}}$

```
from sklearn.linear_model import RidgeCV
```

```
ridge_cv_model = RidgeCV(alpha=[0.1, 1, 10])
```

```
ridge_cv_model.fit(X_train, y_train)
```

ridge_cv_model.alpha_

```
from sklearn.metrics import SCORERS
```

~~SCORERS.D~~ SCORERS.KEYS

new scoring = "neg_mean_squared_error"
absolute

MAE:

RME:

ridge_cv_model.

Lasso Regression

L1 Regularization

- L1 regularization add an penalty term equal to absolute value of the magnitude of coefficients-

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

- limit the size of coefficient
 - o can yield sparse models where some coefficient can becomes zeros-

- Lasso can force some of the coefficient estimates to be exactly equal to zero when the the tuning parameter λ is sufficiently large-
- Similar to Subset selection, the Lasso performs variable selection -
- Model generated from Lasso are generally much easier to interpret
- LassoCV with scikitlearn operates on defining a number of alphas within a range instead of providing the alphas directly-
- let explore the result of Lasso in python and scikit learn

Acronym

LASSO - For least absolute shrinkage and selection operator:-

from sklearn.linear_model import Lasso

from sklearn.linear_model import LassoCV

lasso_cv_model = LassoCV(eps=0.001,
max_iter=100, n_alphas=100)
lasso_cv_model.fit(X_train, y_train)

lasso_cv_model.alpha
0.4945

test_predict5 = lassocv.predict(X_test)

or

ElasticNet

L1 and L2 Regularization

We know that Lasso is able to shrink coefficients to zero, but we haven't taken a deeper dive to how and why that is.

Ability becomes more clear when learning about elastic net which combines Lasso and Ridge together.

for Lasso

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \text{ subject to } \sum_{j=1}^p |\beta_j| \leq s$$

and Ridge

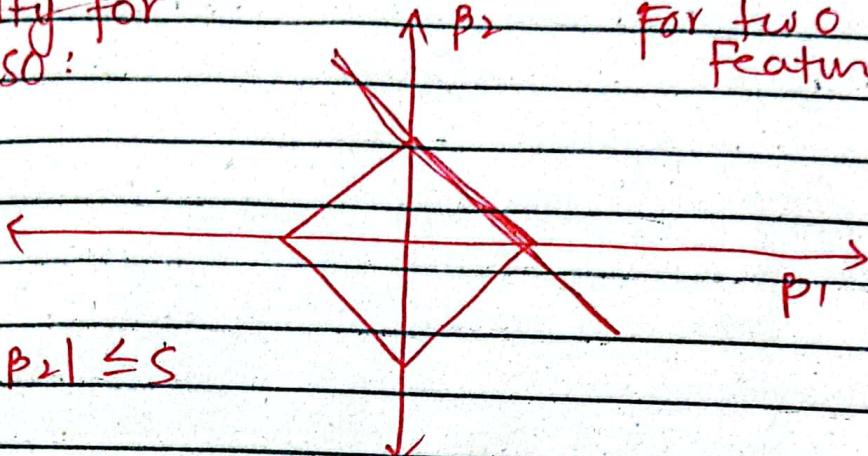
$$\underset{\beta}{\text{minimize}} \left\{ \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) \right\} \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq s$$

- Lasso:
 - A convex object that lies tangent to the boundary, is likely to encounter a corner of a hypercube, for which some components of β are identically zero.

Penalty for Lasso:

$$|\beta_1| + |\beta_2| \leq s$$

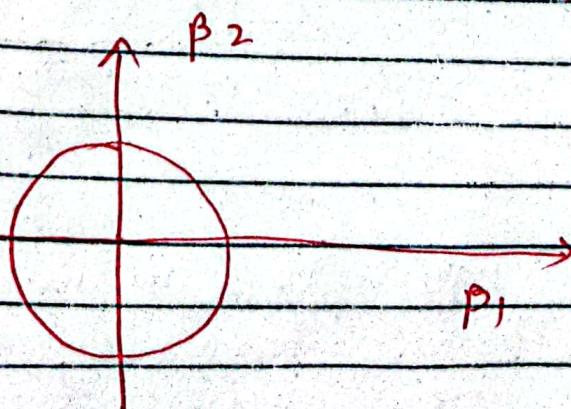
For two Features



- In Lasso we are dealing with hypercube.

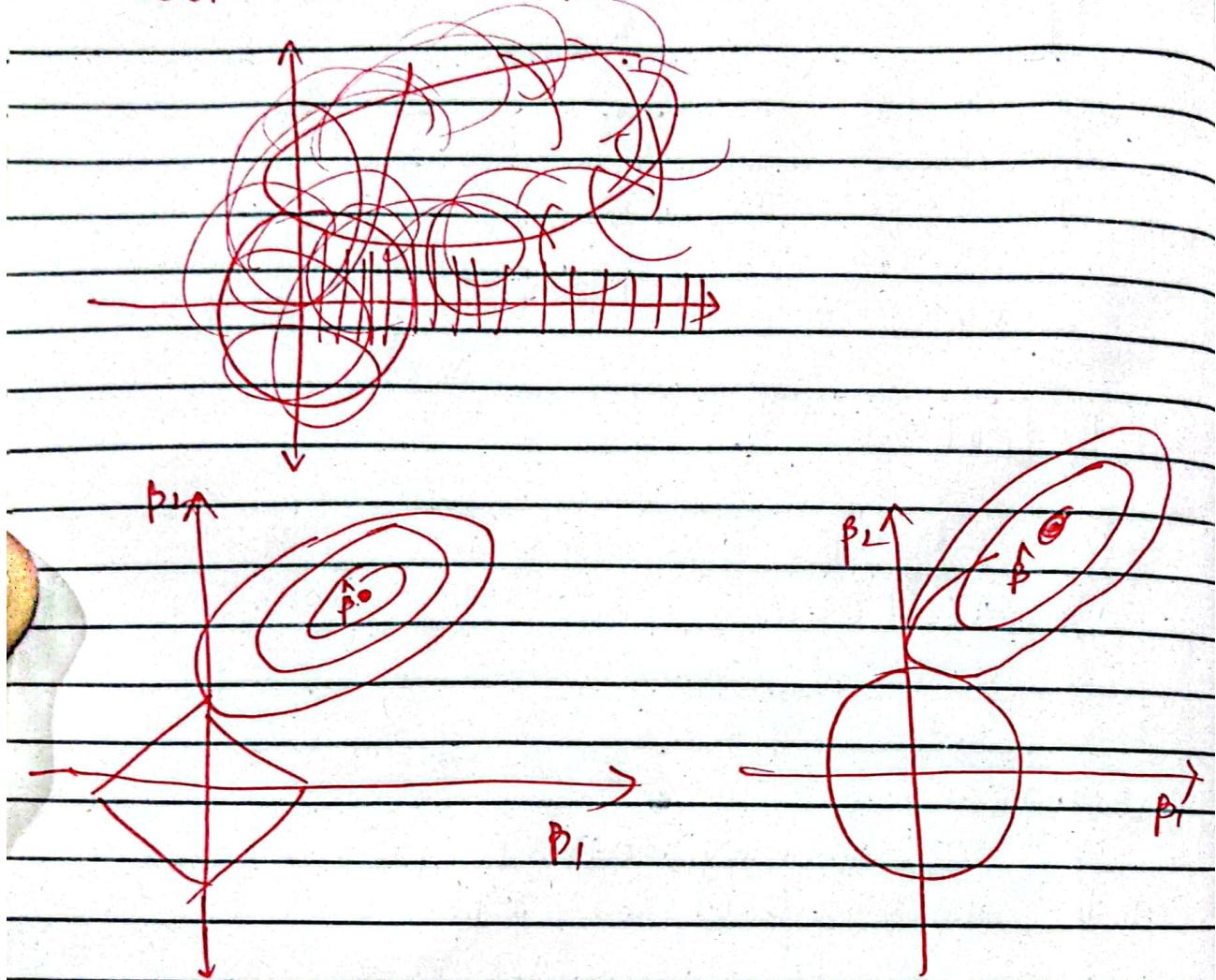
Ridge:-

In case of an n-Sphere, the points on the boundary for which some of the components of β are zero are not distinguished from others and the convex object is no more likely to contact a point at which some components of β are zero than one for which none of them are.



This is why Lasso is more likely to lead coefficient to zero -

This diagram is also commonly shown with contour RSS -



ElasticNet

- Seeks to improve both L1 and L2 regularization by combining them:

$$\text{Error} = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda_1 \sum_{j=1}^p \beta_j^2 + \lambda_2 \sum_{j=1}^p |\beta_j|$$

Here we are trying to minimize both the squared terms and absolute value terms -

These are two distinct

Lambda λ_1, λ_2 values for each penalty.

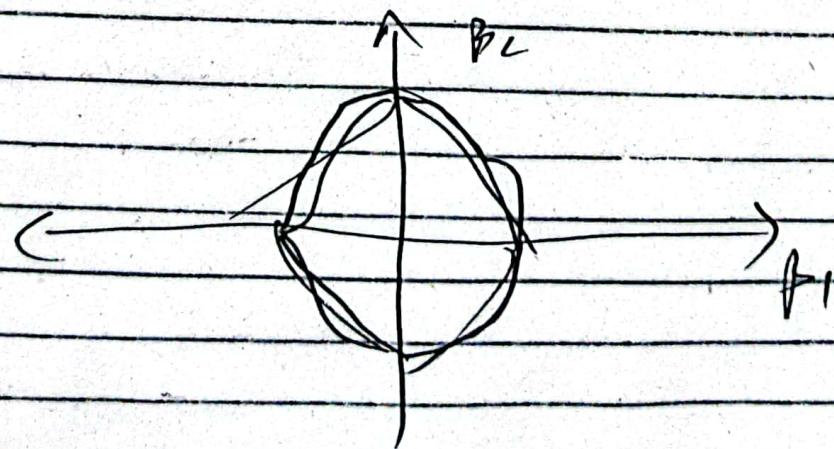
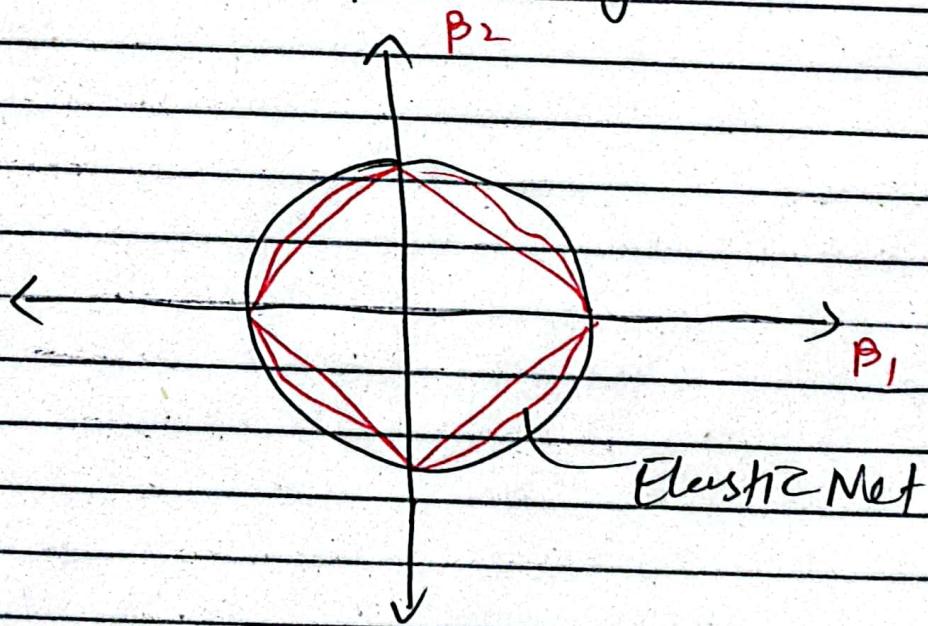
Alternatively we can express this as
Ratio between L1 and L2:

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

We can also use a simplified notation

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} (||y - X\beta||^2 + \lambda_2 ||\beta||^2 + \lambda_1 ||\beta||_1)$$

Elastic Net penalty Region:



Code

```
from sklearn.linear_model import ElasticNetCV  
elastic-model = ElasticNetCV(L1_ratio=[0.1, 0.2  
0.3, 0.99, 0.1],  
eps=0.001, n_alphas=100,  
max_iter=10000)
```

```
elastic-model.fit(XTrain, y-train)
```

```
elastic-model.L1_ratio
```

```
elastic-model.alpha_
```