

Assumptions of Linear Regression

1. Linear in terms of beta coefficients
2. Error is normally distributed and has population mean of 0
3. Homoskedasticity
4. Errors are uncorrelated across observations
5. Little to no multi-collinearity

① Linear in terms of beta coefficients

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon \quad - \text{linear}$$

$$Y = \beta_0 + e^{\beta_1} X^{\beta_2} \quad - \text{non-linear}$$

OK to transform your features and target variables i.e. log
Not possible to use linear algebra to solve non-linear

② Error is normally distributed and has population mean of 0

$$Y = \hat{Y} + \epsilon$$

↑ ↑ ↑
actual estimate error

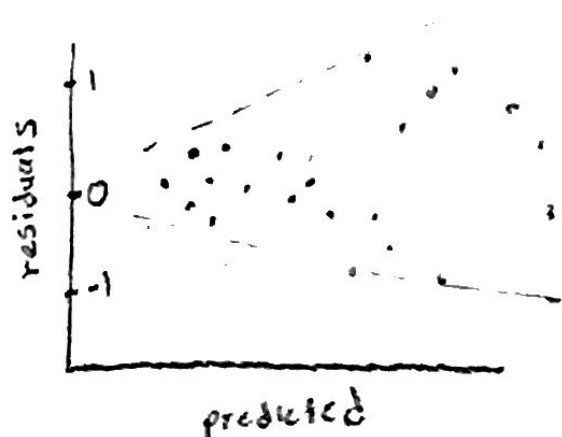
\implies

$$Y = X\beta + \epsilon$$

$$\epsilon = Y - X\beta \sim N(0, \sigma^2)$$

$$N(0, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-\epsilon^2}{2\sigma^2}\right)$$

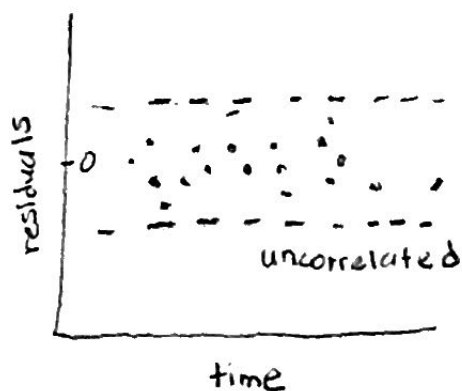
③ Homoskedasticity



Error of model has constant variance across all observations

Possible to use Weighted Least Squares if max residual variance is greater than 4 times min residual variance

④ Errors are uncorrelated across observations



Observed errors that follow a pattern are serially correlated or auto-correlated
use time series prediction instead

⑤ Little to no multi-collinearity

$$L = \epsilon_1 \times \epsilon_2 \times \dots \times \epsilon_n = \prod_{i=1}^n \epsilon_i$$

$$= \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-\frac{\epsilon_i^2}{2\sigma^2}} \right)$$

Multiplying all the ϵ essentially gives the likelihood or the probability of having all these errors

$$\ln L = \sum_{i=1}^n \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-\frac{\epsilon_i^2}{2\sigma^2}} \right)$$

$$= \sum_{i=1}^n \left(\ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{\epsilon_i^2}{2\sigma^2} \right)$$

use log-likelihood to bring the exponential term down

$$= N \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \sum_{i=1}^n \frac{\epsilon_i^2}{2\sigma^2}$$

Move the constant term outside

$$= N \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \left(\frac{1}{2\sigma^2} \right) (y - X\beta)^T (y - X\beta)$$

Substitute $\epsilon = y - X\beta$

$$\frac{\partial \ln L}{\partial \beta} = 0 - \left(\frac{1}{2\sigma^2} \right) 2(-X^T)(y - X\beta) = 0$$

$$\Rightarrow X^T (y - X\beta) = 0$$

$$X^T y - X^T X \beta = 0$$

$$X^T y = X^T X \beta$$

$$\beta = (X^T X)^{-1} X^T y$$

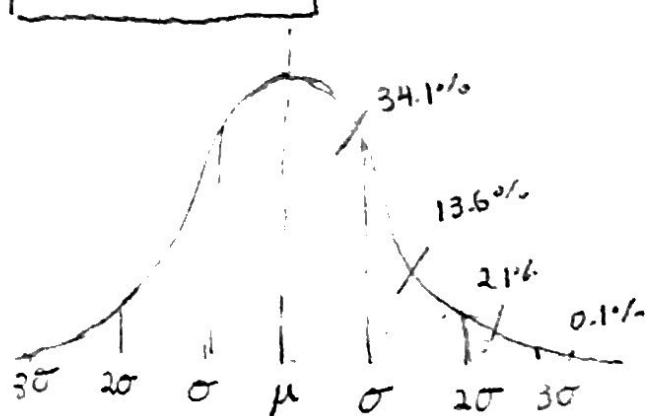
If two or more features are perfectly correlated, the matrix is not invertible.

If two or more features are highly correlated, OLS cannot separately estimate beta coefficient values for them

check magnitude and order of beta coefficients for collinear features

check condition number in statsmodel summary

Standardization



$$Z = \frac{X - \mu}{\sigma}$$

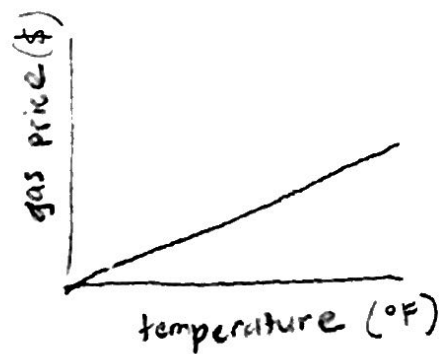
← mean
← standard deviation
feature value

Z score has no units

Standardize your data with StandardScaler before regularization

- StandardScaler works by subtracting off each feature column's mean and dividing by its standard deviation. This results in a Z score - a measure of how many standard deviations a value is from the mean - for each value in the feature.

- How to make beta coefficient interpretable again after standardizing?
A beta coefficient in linear regression tells us how much the outcome variable will change by for every one unit of change in the feature variable.



$$[\beta] = \frac{[\text{gas price}]}{[\text{temperature}]} = \frac{[\$]}{[^\circ\text{F}]}$$

No standardizing

$$[\beta] = \frac{[\text{gas price}]}{[\text{standardized temperature}]} = \frac{[\$]}{\text{z-score}}$$

after standardizing

- Divide beta coefficient by standard deviation of feature to make beta coefficient interpretable again

- Standardizing allows us to interpret our intercept as the expected value of when all feature values have a mean of 0. Not the same as all features being 0.

- If my temperature is 0°, does this mean gas price is not affected?
Most likely not. If my temperature is average, there should be no effect though.

Regularization

- way to avoid overfitting by penalizing the model for having non-zero beta coefficients \rightarrow increase bias, but reduces variance
- a model with some beta coefficients as 0 is less complex than a model with all non-zero beta coefficients

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \quad y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

\downarrow \downarrow
 0 0 v.s.

$$= \beta_0 + \beta_2 x_2$$

- generally, a more complex model leads to overfitting

Ridge (L2)

x_1	x_2	x_3	x_4	y
x_{11}	x_{12}	x_{13}	x_{14}	y_1
x_{21}	x_{22}	x_{23}	x_{24}	y_2

$x_{ij} \rightarrow$ $x_{\text{row index, col index}}$

Loss function

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

$$= \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

λ regularization strength parameter

LASSO (L1)

- performs feature selections in that features can go to 0
- cannot help with multi-collinearity, use ridge regression instead

Loss function

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

$$= \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$$

λ regularization strength parameter

Elastic Net

- weighted combination of Ridge and LASSO, weighted by α value

Loss Function

$$\text{RSS} + \lambda \sum_{j=1}^p (\alpha \beta_j^2 + (1-\alpha) |\beta_j|)$$

Metrics

- R^2 value

$$R^2 = 1 - \frac{\sum (Y_i - \hat{Y}_i)^2}{\sum (Y_i - \bar{Y})^2} = 1 - \frac{SSE}{SST}$$

\nwarrow actual Y \swarrow predicted Y
 \nwarrow mean of actual Y

$R^2 = 0$ model is good as guessing mean of Y_i or \bar{Y}

$R^2 < 0$ model is worse than guessing mean of Y_i or \bar{Y}

- assumes every feature explains variation in outcome variable

- adjusted R^2 value

$$\text{adjusted } R^2 = 1 - \frac{SSE}{SST} \cdot \frac{n-1}{n-k-1}$$

\nwarrow number of observations
 \nwarrow number of variables, excluding constant

- generally better to use adjusted R^2

- penalizes model for adding features that do not fit the model, does not increase with every feature

- prevents overfitting

- Mean Absolute Error (MAE)

$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$

\nwarrow number of observations \nwarrow actual Y \nwarrow predicted Y

- does not indicate your model is underperforming or overperforming
 - because we use absolute value of the residual

- small MAE \rightarrow model is great at predicting

- large MAE \rightarrow model has trouble in some areas

- interpretable, more business-oriented

- Root Mean Square Error (RMSE)

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$$

\nwarrow number of observations \nwarrow actual Y \nwarrow predicted Y

- punishes more heavily against outliers in data and hence, large residual values

- do outliers matter to your model or not?

Work flow

Data → Train-test-split 80/20

* remember to have intercept
otherwise, call add_constant
on X to get intercept

Standardize with
Standard Scaler

OLS via statsmodel

Lasso CV

high p values?
too many features?

high
condition
number?

Satisfied?

Yes

No

Validation/
Cross-Validation

More features,
Feature engineer
Feature transform

Polynomial Features
for Interactions

custom metric
like MAE or
RMSE

loss function,
reduce variance

Grid Search CV

Lasso CV, Ridge CV,
Elastic Net CV

Satisfied?

No

Yes

Predict on
Test set

Residuals vs
Predicted

Done...
for now

