



PG Diploma in ML

Course : PG Diploma in ML

Lecture On : Live Session
on Regularisation in Logistic
Regression

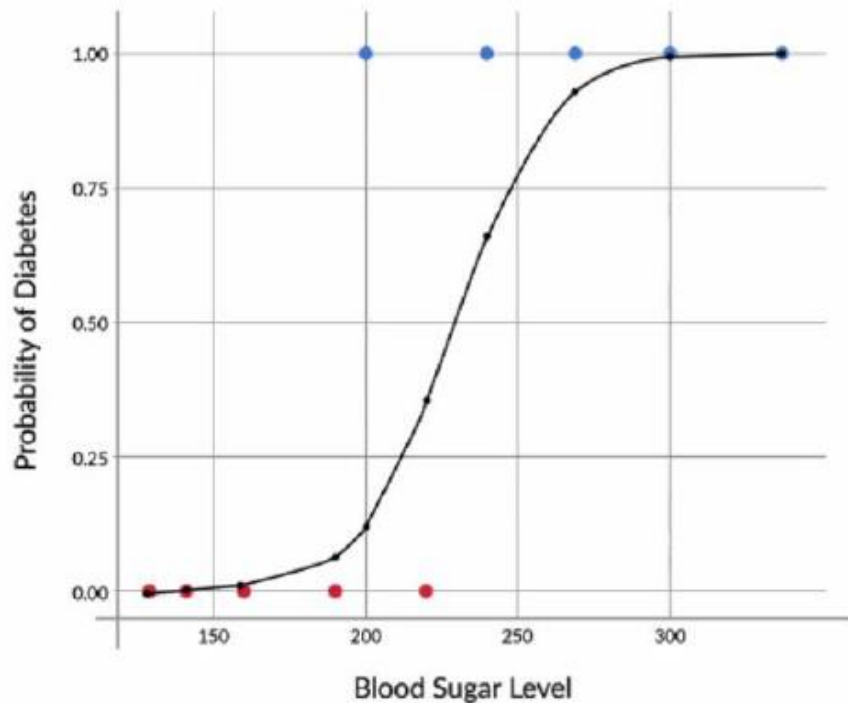
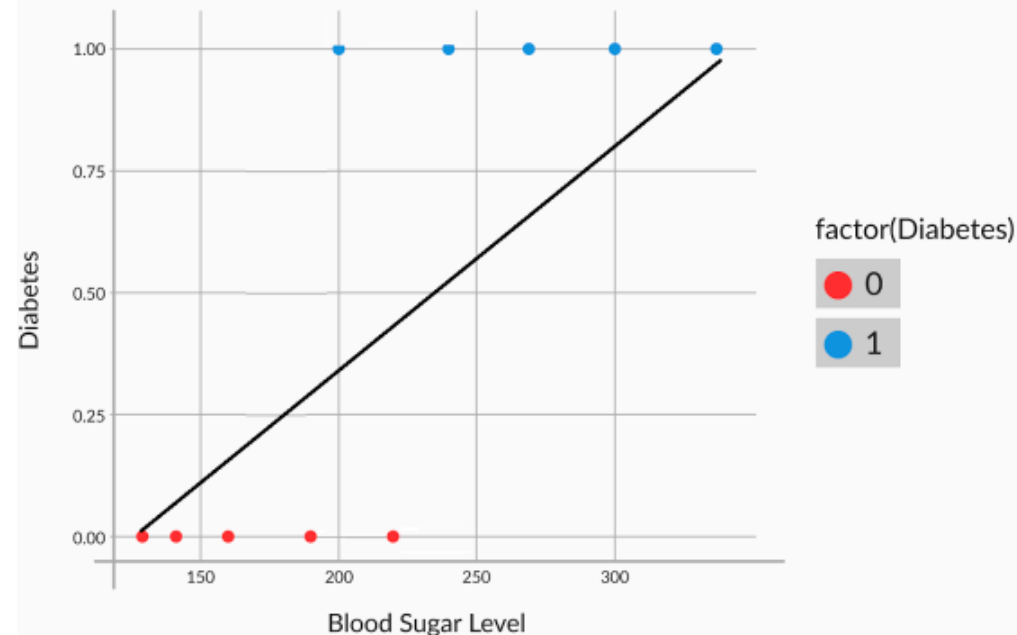
Instructor : Dr. Reena
Duggal

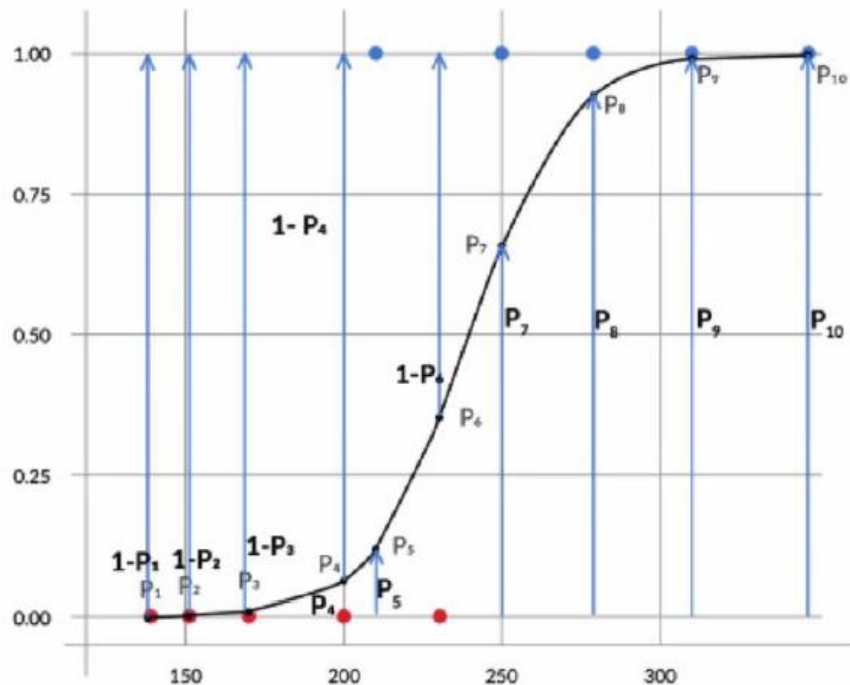
Today's Agenda

1. Logistic Regression quick recap
2. Regularisation in Logistic Regression
3. Lasso based Regularisation
4. Ridge based Regularisation
5. Hands-on in Python
6. Doubt resolution and discussion

Logistic Regression is a supervised classification model. It allows you to make predictions from labelled data, if the target (output) variable is categorical.

DIABETES DATA PLOT





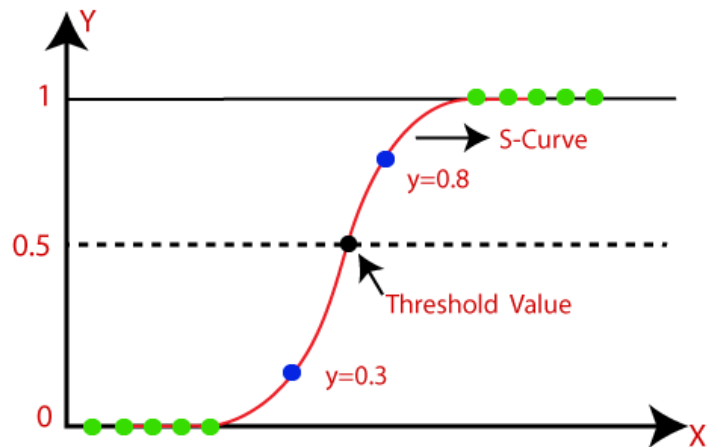
$$\text{Likelihood} = (1-P_1)(1-P_2)(1-P_3)(1-P_4)(P_5)(1-P_6)(P_7)(P_8)(P_9)(P_{10})$$

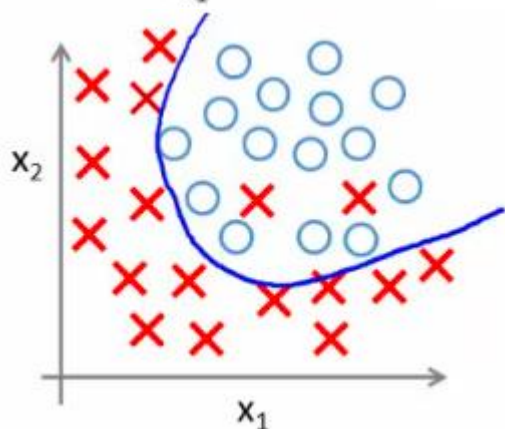
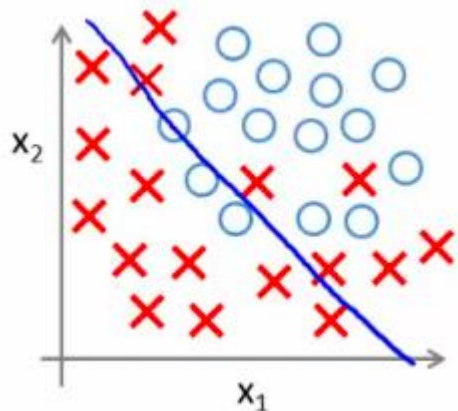
Generally, it is the product of -

$$[(1-P_i)(1-P_i) \text{ ----- for all non-diabetics -----}] \times [(P_i)(P_i) \text{ ----- for all diabetics -----}]$$

$$P(\text{Diabetes}) = \frac{1}{1+e^{-(\beta_0+\beta_1x)}} \quad \ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1x$$

$$P = \frac{1}{1+e^{-(\beta_0+\beta_1x_1+\beta_2x_2+\beta_3x_3+\dots)}}$$





$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$$

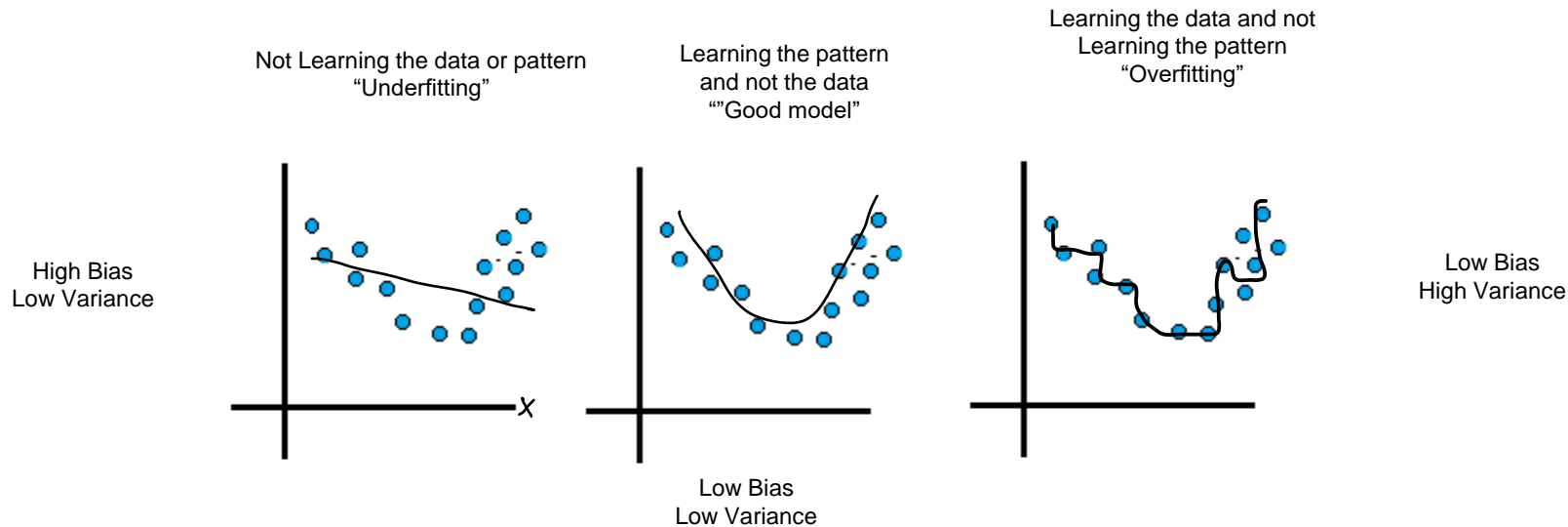
You can define the degree of polynomial that will build the decision boundary.

Need to choose Coefficients which maximize Log-Likelihood(objective function) or minimize the negative of Log-Likelihood

$$\hat{\beta} = \min_{\beta} -LL(\beta; y, X)$$

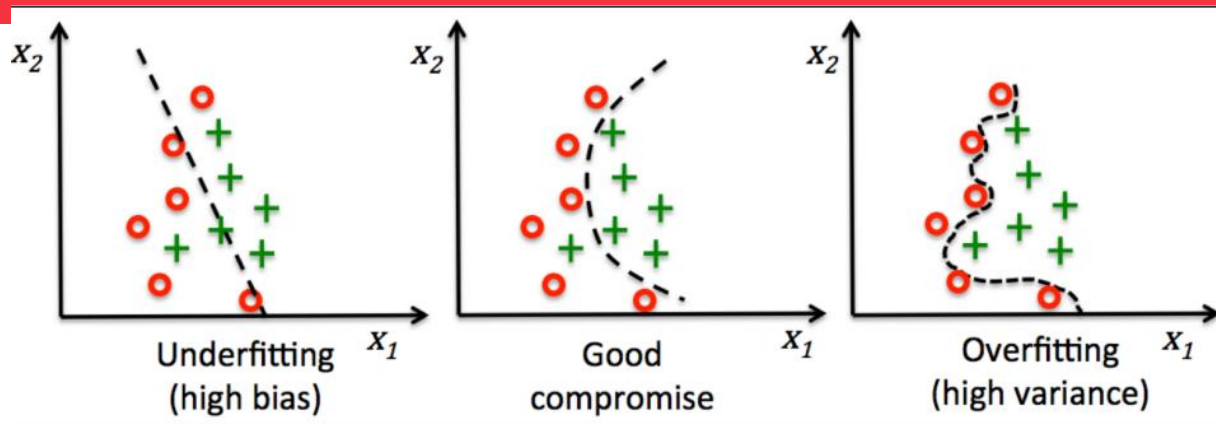
The default nature of a model is that it does not account for model complexity - it only tries to maximize Log-Likelihood or reduce cost function, although it may result in arbitrarily complex coefficients which result in **overfitting**.

What is Bias and Variance



Bias - Amount of Error that the model is making on train data

Variance - Amount that the model(target variable) will change given different training data



- **Regularization** can be used to avoid overfitting by creating an optimally complex model
- Regularization adds a penalty on the different parameters of the model to reduce the freedom of the model.
- Hence, the model will be less likely to fit the noise of the training data and will improve the generalization abilities of the model
- In regularized logistic regression, the objective function has two parts – the **Cost function(-ve of Log-likelihood)** and the **regularization term**.

$$\hat{\beta} = \min_{\beta} -LL(\beta; y, X) + \lambda R(\beta)$$

In ridge regression, an additional term of "**sum of the squares of the coefficients**" is added to the cost function along with the error term

Ridge Regression

$$\left[\underset{\beta}{\text{Min}} \right] \left[-LL(\beta; y, X) + \lambda \sum_{i=1}^k \beta_i^2 \right]$$

Negative of Log-Likelihood

Regularization term
or Regularizer or Penalty

Sum of the squares
of the coefficients

Hyper Parameters

The diagram illustrates the Ridge Regression cost function. It features a large equation:
$$\left[\underset{\beta}{\text{Min}} \right] \left[-LL(\beta; y, X) + \lambda \sum_{i=1}^k \beta_i^2 \right]$$
 The first part of the equation, $-LL(\beta; y, X)$, is enclosed in a blue bracket and labeled "Negative of Log-Likelihood". The second part, $\lambda \sum_{i=1}^k \beta_i^2$, is enclosed in a blue bracket and labeled "Sum of the squares of the coefficients". A red arrow points from the λ term to the text "Hyper Parameters". Another red arrow points from the entire second part of the equation to the text "Regularization term or Regularizer or Penalty".

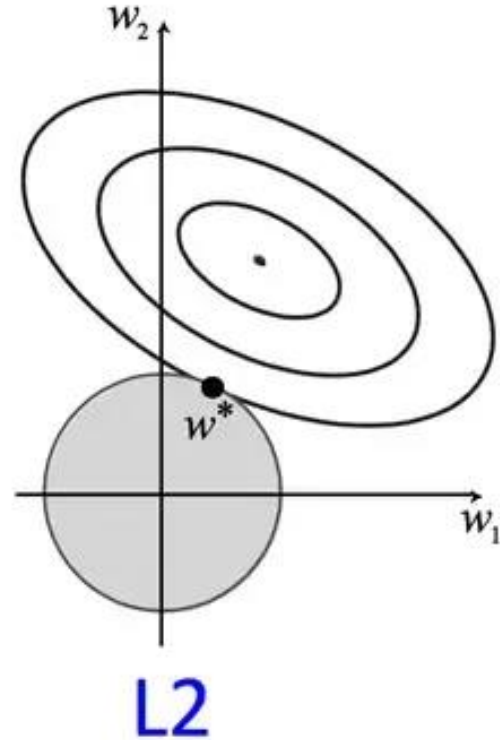
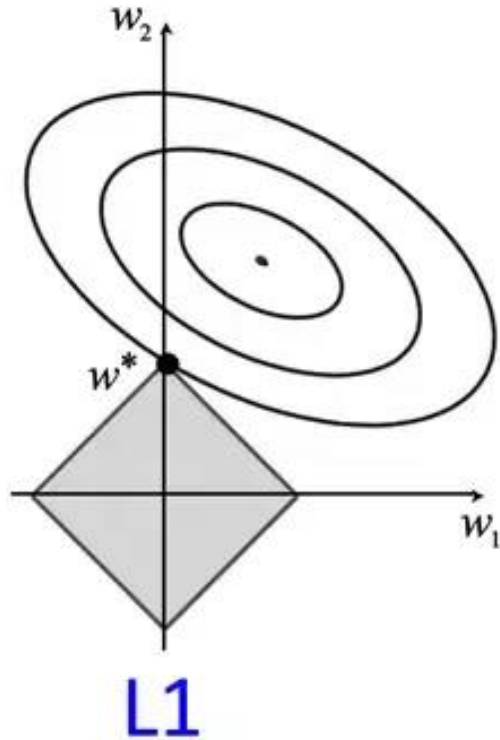
In case of lasso regression, a regularisation term of "**sum of the absolute value of the coefficients**" is added.

Lasso Regression

$$\left[\begin{array}{c} \text{Min} \\ \beta \end{array} \left[\begin{array}{l} -LL(\beta; y, X) \\ \text{Negative of Log-Likelihood} \end{array} \right] + \lambda \sum |\beta_i| \right]$$

Sum of the absolute values

- Lasso regression not only helps in reducing over-fitting but it can help us in feature selection



Cost function = minimize(-ve of Log-Likelihood + λ * Complexity)

λ (also called the regularization rate) is the tuning parameter that decides how much we want to penalize the flexibility of our model

- I.** $\lambda = 0$, no regularisation
- II.** λ is high, more regularisation. Model will be simple, but you run the risk of underfitting your data. Your model won't learn enough about the training data to make useful predictions.
- III.** λ is low, less regularisation. Model will be more complex, and you run the risk of overfitting your data. Your model will learn too much about the particularities of the training data, and won't be able to generalize to new data.

How to choose λ ?

You'll play around with different values.

- Other common names for λ :
 - *alpha* in `sklearn.ridge()` and `sklearn.lasso()`
 - *C* in many algorithms e.g. `sklearn.LogisticRegression()`
 - Usually C actually refers to the inverse regularization strength, $\frac{1}{\lambda}$
 - Interpretation:
 - C is high, less regularisation. Model is complex.
 - C is low, high regularisation. Model is simple.
- In `sklearn.LogisticRegression()`, L2 regularization(Ridge) is applied by default

Python Case Study



Thank You!