**Regularization in Machine Learning**

**Overfitting:**

One of the major aspects of training your machine learning model is avoiding overfitting. The model will have a low accuracy if it is overfitting. This happens because your model is trying too hard to capture the noise in your training dataset. **By noise we mean the data points that don't really represent the true properties of your data, but random chance**. Learning such data points, makes your model more flexible, at the risk of overfitting.

**The concept of balancing bias and variance, is helpful in understanding the phenomenon of overfitting.**

**Balancing Bias and Variance to Control Errors in Machine Learning**

One of the ways of avoiding overfitting is using cross validation, that helps in estimating the error over test set, and in deciding what parameters work best for your model.

This article will focus on a technique that helps in avoiding overfitting and also increasing model interpretability.

**Regularization**

This is a form of regression, that constrains/ regularizes or shrinks the coefficient estimates towards zero. In other words, **this technique discourages learning a more complex or flexible model, so as to avoid the risk of overfitting.**

A simple relation for linear regression looks like this. Here Y represents the learned relation and β represents the coefficient estimates for different variables or predictors(X).

**$Y \approx \beta 0 + \beta 1X1 + \beta 2X2 + …+ \beta pXp$**

The fitting procedure involves a loss function, known as residual sum of squares or RSS. The coefficients are chosen, such that they minimize this loss function.

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

Now, this will adjust the coefficients based on your training data. If there is noise in the training data, then the estimated coefficients won't generalize well to the future data. This is where regularization comes in and shrinks or regularizes these learned estimates towards zero.

**Ridge Regression**

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

Above image shows ridge regression, where the **RSS is modified by adding the shrinkage quantity.** Now, the coefficients are estimated by minimizing this function. Here, $\lambda$ **is the tuning parameter that decides how much we want to penalize the flexibility of our model.** The increase in flexibility of a model is represented by increase in its coefficients, and if we want to minimize the above function, then these coefficients need to be small. This is how the Ridge regression technique prevents coefficients from rising too high. Also, notice that we shrink the estimated association of each variable with the response, except the intercept $\beta 0$, This intercept is a measure of the mean value of the response when xi1 = xi2 = …= xip = 0.

When $\lambda = 0$, the penalty term has no effect, and the estimates produced by ridge regression will be equal to least squares. However, **as $\lambda \rightarrow \infty$, the impact of the shrinkage penalty grows, and the ridge regression coefficient estimates will approach zero**. As can be seen, selecting a good value of $\lambda$ is critical. Cross validation comes in handy for this purpose. The coefficient estimates produced by this method are **also known as the L2 norm**.

**The coefficients that are produced by the standard least squares method are scale equivariant**, i.e. if we multiply each input by c then the corresponding coefficients are scaled by a factor of 1/c. Therefore, regardless of how the predictor is scaled, the multiplication of predictor and coefficient(Xjβj) remains the same.

**However, this is not the case with ridge regression, and therefore, we need to standardize the predictors or bring the predictors to the same scale before performing ridge regression**. The formula used to do this is given below.

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_{ij} - \bar{x}_j)^2}},$$

**Lasso**

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^{p} |\beta_j|.$$
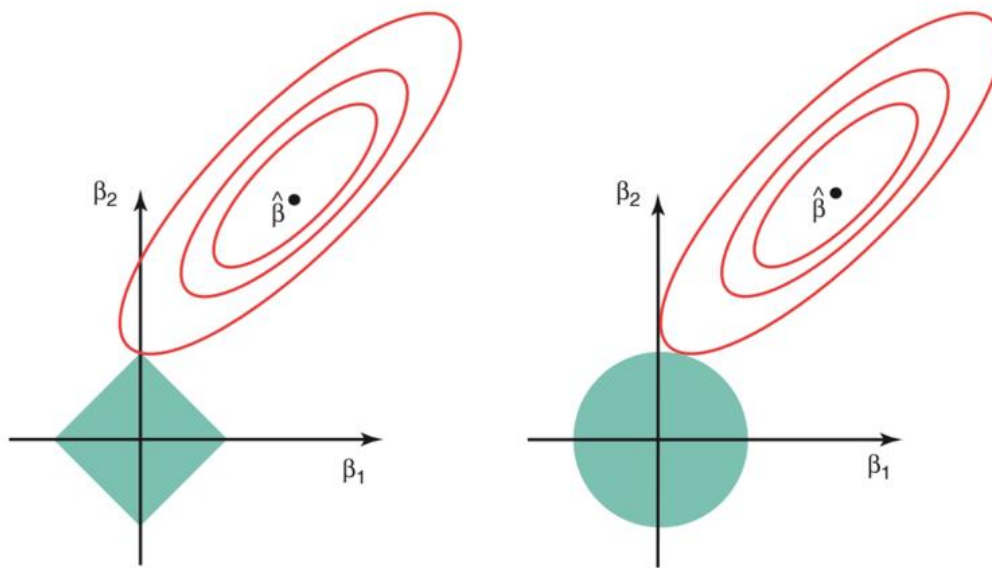
Lasso is another variation, in which the above function is minimized. Its clear that **this variation differs from ridge regression only in penalizing the high coefficients**. It uses $|\beta j|$(modulus)instead of squares of $\beta$, as its penalty. In statistics, this is **known as the L1 norm**.

Lets take a look at above methods with a different perspective. The ridge regression can be thought of as solving an equation, where summation of squares of coefficients is less than or equal to s. And the Lasso can be thought of as an equation where summation of modulus of coefficients is less than or equal to s. Here, s is a constant that exists for each value of shrinkage factor $\lambda$. **These equations are also referred to as constraint functions.**

**Consider their are 2 parameters in a given problem**. Then according to above formulation, the **ridge regression is expressed by $\beta 1^2 + \beta 2^2 \leq s$**. This implies that ridge regression coefficients have the smallest RSS(loss function) for all points that lie within the circle given by $\beta 1^2 + \beta 2^2 \leq s$.

Similarly, **for lasso, the equation becomes,$|\beta 1|+|\beta 2|\leq s$**. This implies that lasso coefficients have the smallest RSS(loss function) for all points that lie within the diamond given by $|\beta 1|+|\beta 2|\leq s$.

The image below describes these equations.

Credit : An Introduction to Statistical Learning by Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani

**The above image shows the constraint functions(green areas), for lasso(left) and ridge regression(right), along with contours for RSS(red ellipse)**. Points on the ellipse share the value of RSS. For a very large value of s, the green regions will contain the center of the ellipse, making coefficient estimates of both regression techniques, equal to the least squares estimates. But, this is not the case in the above image. In this case, the lasso and ridge regression coefficient estimates are given by the first point at which an ellipse contacts the constraint region. **Since ridge regression has a circular constraint with no sharp points, this intersection will not generally occur on an axis, and so the ridge regression coefficient estimates will be exclusively non-zero. However, the lasso constraint has corners at each of the axes, and so the ellipse will often intersect the constraint region at an axis. When this occurs, one of the coefficients will equal zero.** In higher dimensions(where parameters are much more than 2), many of the coefficient estimates may equal zero simultaneously.

**This sheds light on the obvious disadvantage of ridge regression, which is model interpretability.** It will shrink the coefficients for least important predictors, very close to zero. But it will never make them exactly zero. In other words, the final model will include all predictors. However, in the case of the lasso, the L1 penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter $\lambda$ is sufficiently large.

**Therefore, the lasso method also performs variable selection and is said to yield sparse models.**

**What does Regularization achieve?**

A standard least squares model tends to have some variance in it, i.e. this model won't generalize well for a data set different than its training data. **Regularization, significantly reduces the variance of the model, without substantial increase in its bias**. So the tuning parameter $\lambda$, used in the regularization techniques described above, controls the impact on bias and variance. As the value of $\lambda$ rises, it reduces the value of coefficients and thus reducing the variance. **Till a point, this increase in $\lambda$ is beneficial as it is only reducing the variance(hence avoiding overfitting), without loosing any important properties in the data.** But after certain value, the model starts loosing important properties, giving rise to bias in the model and thus underfitting. Therefore, the value of $\lambda$ should be carefully selected.