

L1 and L2 regularization

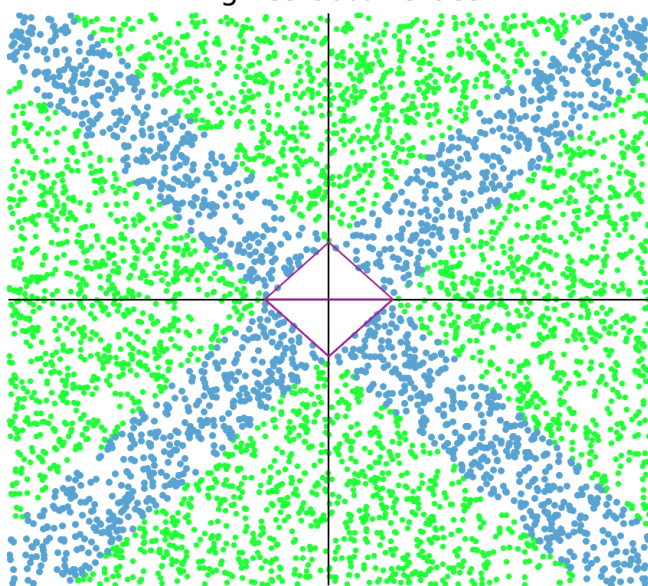
If both L1 and L2 regularization work well, you might be wondering why we need both. It turns out they have different but equally useful properties. From a practical standpoint, L1 tends to shrink coefficients to zero whereas L2 tends to shrink coefficients evenly. L1 is therefore useful for feature selection, as we can drop any variables associated with coefficients that go to zero. L2, on the other hand, is useful when you have collinear/codependent features. (An example pair of codependent features is `gender` and `is_pregnant` since, at the current level of medical technology, only females can be `is_pregnant`.) Codependence tends to increase coefficient variance, making coefficients unreliable/unstable, which hurts model generality. L2 reduces the variance of these estimates, which counteracts the effect of codependencies.

(The [code to generate all images](#) is available.)

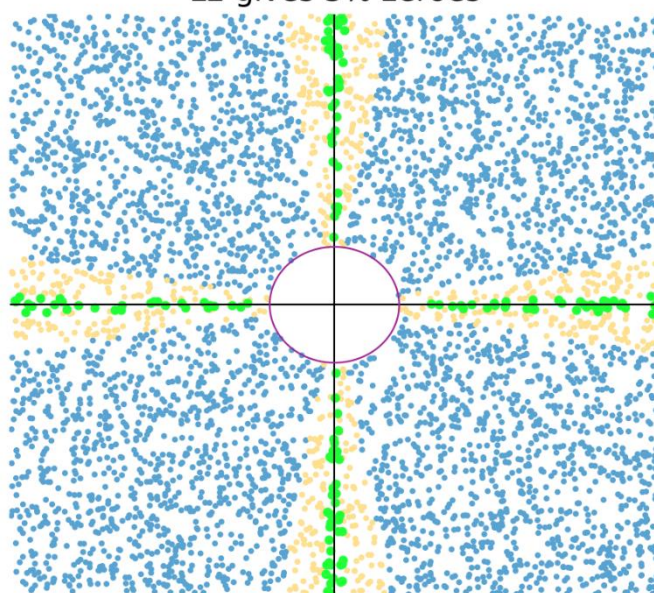
3.1 L1 regularization encourages zero coefficients

One of the key questions that I want to answer is: “Does L1 encourage model coefficients to shrink to zero?” (The answer is, Yes!) So, let's do some two-variable simulations of random quadratic loss functions at random locations and see how many end up with a coefficient at zero. There is no guarantee that these random paraboloid loss functions in any way represent real data sets, but it's a way to at least compare L1 and L2 regularization. Let's start out with symmetric loss functions, which look like bowls of various sizes and locations, and compare how many zero coefficients appear for L1 and L2 regularization:

Symmetric Loss function min cloud
L1 gives 66% zeroes



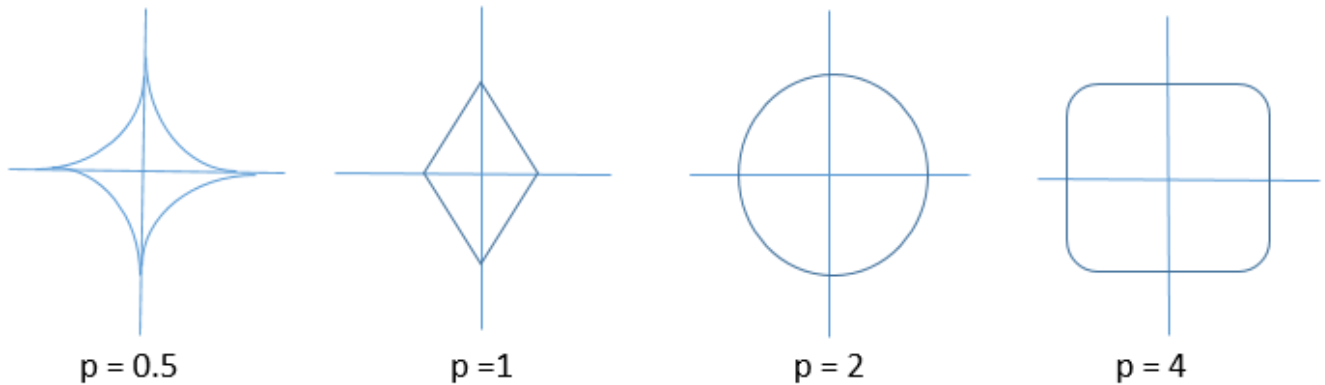
Symmetric Loss function min cloud
L2 gives 3% zeroes



Green dots represent a random loss function that resulted in a regularized coefficient being zero. Blue represents a random loss function where no regularized coefficient was zero (North, South, East, West compass points). Orange represents loss functions in L2 plots that had at least one coefficient close to zero (within 10% of the max distance of any coefficient pair.) L1 tends not to give near misses and so the simulation on the left is just blue/green. As you can see in the simulations (5000 trials), the L1 diamond constraint zeros a coefficient for any loss function whose minimum is in the zone perpendicular to the diamond edges. The L2 circular constraint only zeros a coefficient for loss function minimums sitting really close to or on one of the axes. The orange zone indicates where L2 regularization gets close to a zero for a random loss function. Clearly, L1 gives many more zero coefficients (66%) than L2 (3%) for symmetric loss functions.

Actually, there are different possible choices of regularization with different choices of order of the parameter in the regularization term, which is denoted by $\sum_i |\theta_i|^p$. This is more generally known as Lp regularizer.

Let us try to visualize some by plotting them. For making visualization easy, let us plot them in 2D space. For that we suppose that we just have two parameters. Now, let's say if $p=1$, we have term as $\sum_i |\theta_i|^p = |\theta_1| + |\theta_2|$. Can't we plot this equation of line? Similarly plot for different values of p are given below.

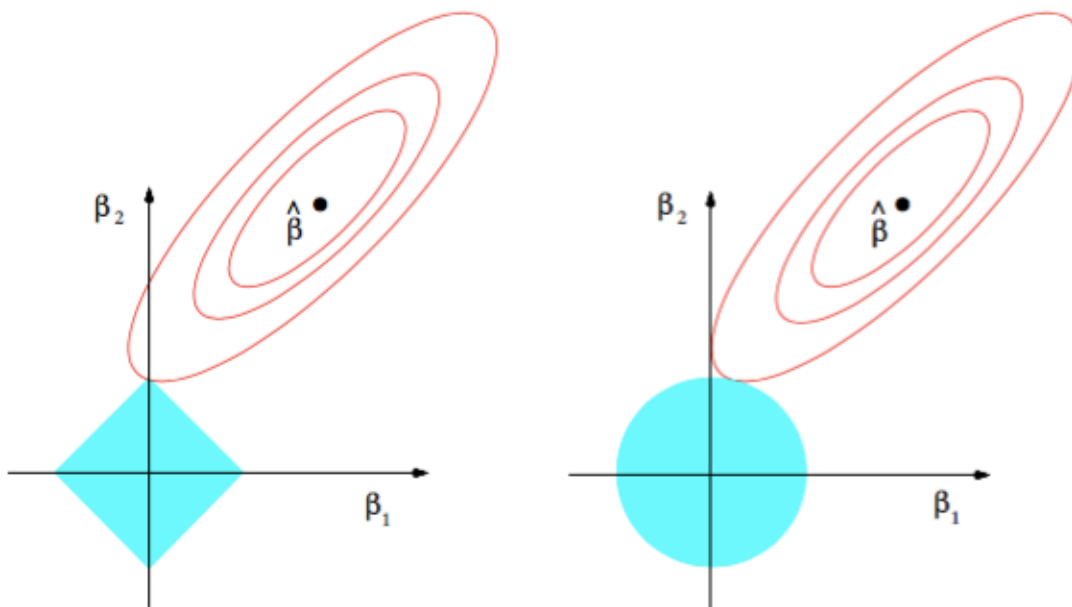


In the above plots, axis denote the parameters(θ_1 and θ_2). Let us examine them one by one.

For $p=0.5$, we can only get large values of one parameter only if other parameter is too small. For $p=1$, we get sum of absolute values where the increase in one parameter θ is exactly offset by the decrease in other. For $p=2$, we get a circle and for larger p values, it approaches a round square shape.

The two most commonly used regularization are in which we have $p=1$ and $p=2$, more commonly known as L1 and L2 regularization.

Look at the figure given below carefully. The blue shape refers the regularization term and other shape present refers to our least square error (or data term).



The first figure is for L1 and the second one is for L2 regularization. The black point denotes that the least square error is minimized at that point and as we can see that it increases quadratically as we move from it and the regularization term is minimized at the origin where all the parameters are zero.

Now the question is that at what point will our cost function be minimum? The answer will be, since they are quadratically increasing, the sum of both the terms will be minimized at the point where they first intersect.

Take a look at the L2 regularization curve. Since the shape formed by L2 regularizer is a circle, it increases quadratically as we move away from it. The L2 optimum(which is basically the intersection point) can fall on the axis lines only when the minimum MSE (mean square error or the black point in the figure) is also exactly on the axis. But in case of L1, the L1 optimum can be on the axis line because its contour is sharp and therefore there are high chances of intersection point to fall on axis. Therefore it is possible to intersect on the axis line, even when minimum MSE is not on the axis. If the intersection point falls on the axes it is known as sparse.

Therefore L1 offers some level of sparsity which makes our model more efficient to store and compute and it can also help in checking importance of feature, since the features that are not important can be exactly set to zero.