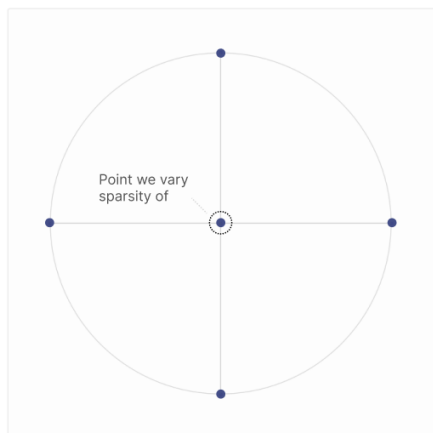## PERTURBING A SINGLE FEATURE

The simplest kind of non-uniform superposition is to vary one feature and leave the others uniform. As an experiment, let's consider an experiment where we represent $n = 5$ features in $m = 2$ dimensions. In the uniform case, with importance $I = 1$ and activation density $1 - S = 0.05$, we get a regular pentagon. But if we vary one point – in this case we'll make it more or less sparse – we see the pentagram *stretch* to account for the new value. If we make it denser, activating more frequently (yellow) the other features repel from it, giving it more space. On the other hand, if we make it sparser, activating less frequently (blue) it takes less space and other points push towards it.

If we make it sufficiently sparse, there's a phase change, and it collapses from a pentagon to a pair of digons with the sparser point at zero. The phase change corresponds to loss curves corresponding to the two different geometries crossing over. (This observation allows us to directly confirm that it is genuinely a first order phase change.)
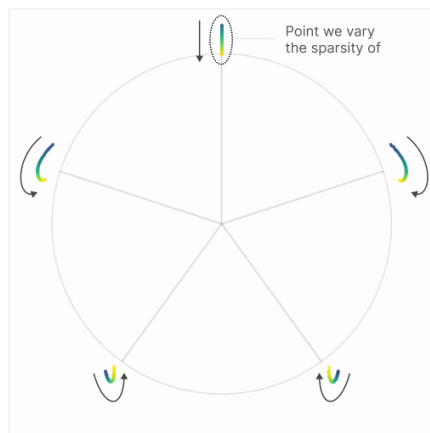
To visualize the solutions, we canonicalize them, rotating them to align with each other in a consistent manner.

**Digon (Square) Solutions**



When the sparsity of the varied point falls below a certain critical threshold (~2.5x less than others) the pentagon solution changes to two digons.

**Pentagon Solutions**



Note how vertices shift as sparsity changes

To study non-uniform sparsity, we consider models with five features, varying the sparsity of a single feature and observing how the resulting solutions change. We observe a mixture of continuous deformation and sharp phase changes.
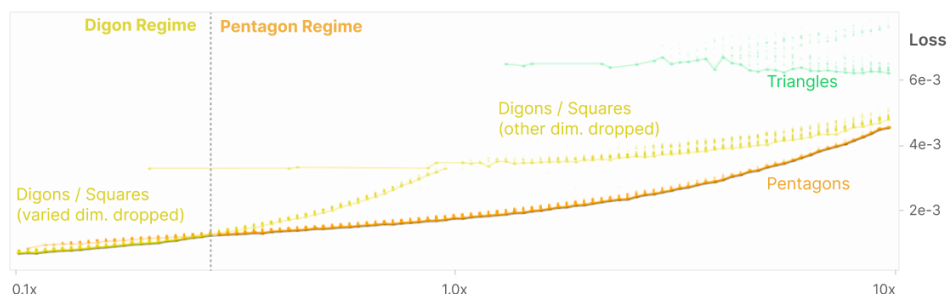
**Parameters**

$n$ = 5
$m$ = 2
$I_i$ = 1
$1-S$ = 0.05 (baseline)

**Relative Feature Density (1-S)**

0.1x          1.0x          10x
sparser                    denser

**The Pentagon-Digon Phase Change Corresponds to a Loss Curve Crossover**



Gradient descent has trouble moving between solutions associated with different geometries. As a result, fitting the model will often produce non-optimal solutions. By characterizing and plotting these, we can see that each geometry creates a different loss curve, and that the pentagon-digon phase change corresponds to a cross over between the curves.