1. Because it kind of represents the distance between each predicted and actual point (which has squaring in it), and calculating power 4 is computationally more expensive.
2. Because performing matrix operations is slower than gradient descent in a large dataset, so it is not practical.
3. It is better to use polynomial regression and use area (pi * r ^ 2) as a feature with more meaning and relation to pizza price. Also it is possible to have 2 pizzas with the same radius but different width so the thicker one is more expensive and that shows that the relation is not linear.
   So we can use polynomial regression with m = 3 to have 1, r, r^2, r^3 as features that represent radius, area and volume.
4. Not particularly, cost can oscillate over minimum and gets further from it in each step.
5.



Final w is [3.776, 0.071]