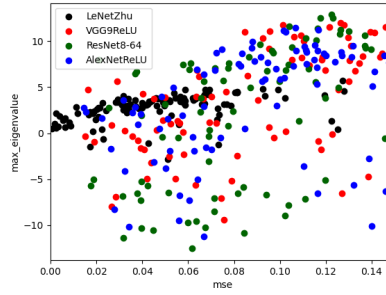
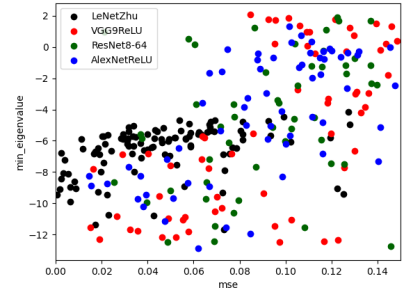


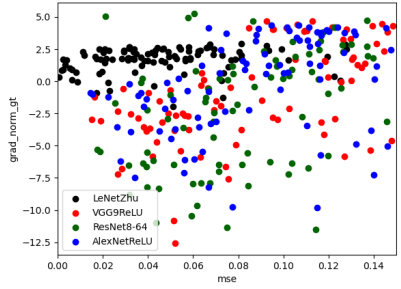
(a) grad_norm vs MSE (L2, $\sigma_S = 0.35$)



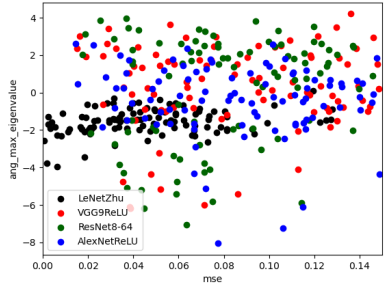
(b) max vs MSE (L2, $\sigma_S = 0.66$)



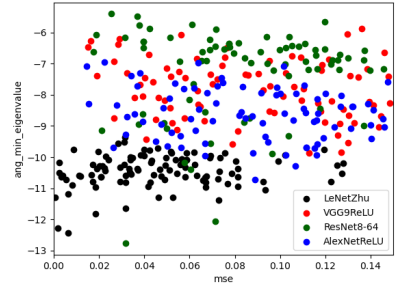
(c) min vs MSE (L2, $\sigma_S = 0.68$)



(d) grad_norm vs MSE (CS, $\sigma_S = 0.00$)



(e) ang_max vs MSE (CS, $\sigma_S = 0.72$)



(f) ang_min vs MSE (CS, $\sigma_S = 0.75$)

Figure 6: Comparison of gradient norm, maximum and minimum eigenvalues of Hessian in terms of the correlation with MSE of reconstructed samples over several architectures on ImageWoof test samples.

Algorithm 1: Pseudocode for computing maximum eigenvalue of Hessian, PyTorch-like

```
1  def max_eigenvalue(x, label, model, g_gt):
2      # x: ground-truth image
3      # label : x's class
4      # model : FL model
5      # g_gt : model gradient from ground truth
6
7      m_e = 0 # max eigenvalue candidate
8
9      for _ in range(N): # N is large enough
10         v = torch.randn_like(x) # initialize vector
11         v = v/torch.norm(v) # normalize vector
12
13         #initialize gradients to zero
14         x_tmp = x.copy()
15         x_tmp.grad *= 0.0
16         model.zero_grad()
17
18         loss1 = loss_fn(model(x_tmp - e*v), label) #
19             compute loss at neighborhood of x, e is
20             small
21
22         gradient1 = torch.autograd.grad(loss1, model.
23             parameters(), create_graph=True) #compute
24             gradient
25
26         g1 = torch.cat([gradient1.view(-1, 1).detach()
27             for g in gradient1], dim=0).squeeze() #
28             flatten gradient into 1-D
29
30         model.zero_grad()
31
32         loss2 = loss_fn(model(x_tmp + e*v), label) #
33             compute loss at neighborhood of x at the
34             opposite side, e is small
35
36         gradient2 = torch.autograd.grad(loss2, model.
37             parameters(), create_graph=True) #compute
38             gradient
39
40         g2 = torch.cat([gradient2.view(-1, 1).detach()
41             for g in gradient2], dim=0).squeeze() #
42             flatten gradient into 1-D
43
44         g_diff = (g2 - g)/(2*e) # compute difference,
45             which is approximately a Jacobian vector
46             product
47
48         g_diff = g_diff.detach()
49
50         # For cosine similiarity loss, compute
51         intermeidate terms with g_gt
52
53         if loss_type == 'cosine':
54             g_diff = g_diff - (g_diff*g_gt).sum(0)*g_gt
55
56         model.zero_grad()
57         loss = loss_fn(model(x_tmp), label) # compute
58             loss at gt
59         gradient = torch.autograd.grad(loss, model.
60             parameters(), create_graph=True) #
61             compute gradient
62
63         g = torch.cat([gradient.view(-1, 1).detach()
64             for g in gradient], dim=0).squeeze()
65         #flatten gradient into 1-D
66         ig = torch.autograd.grad(outputs=g, inputs=x,
67             grad_outputs=g_diff) #Hessian vector
68             product
69
70         nrm = (ig[0]*v).sum() #maximum eigenvalue
71             candidate
72
73         m_e = max(m_e, nrm) #update maximum eigenvalue
74
75         v = ig[0]
76     return m_e
```
