

二阶回归模型表现出对边缘的渐进锐化 稳定

阿蒂什·阿加瓦拉¹ 法比安·佩德雷戈萨¹ 杰弗里·彭宁顿¹

抽象的
大步长梯度下降的最新研究
规模表明,通常存在一种制度
最大特征值的初始增加
损失 Hessian (渐进锐化),其次
通过使特征值稳定在最大值附近,从而允许收敛 (边缘
稳定)。这些现象本质上是非线性的,并且对于恒定
神经正切核 (NTK)体系中的模型不会发生,对于
其预测函数大约为
参数呈线性。因此,我们认为
下一类最简单的预测模型,即
那些参数是二次的,其中
我们称之为二阶回归模型。为了
二维二次目标,我们证明
这个二阶回归模型表现出
NTK 特征值逐渐锐化至与边缘略有不同的值
稳定性,我们明确计算。在更高的
尺寸,模型通常显示相似
行为,即使没有特定的结构
神经网络,表明渐进锐化和稳定边缘行为并不是唯一的
神经网络的特征,并且可能是一个更
离散学习算法的一般性质
在高维非线性模型中。

一、简介

深度学习理论理解的最新趋势集中在线性化机制上,其中神经网络
切线内核 (NTK) 控制学习动态 (Jacot
等,2018;李等人,2019) 。 NTK 描述学习
所有网络在足够短的时间范围内的动态,
并可以描述大型网络的动态
时间范围。在 NTK 体系中,存在一个函数空间

¹谷歌DeepMind。通讯作者:阿蒂什·阿加瓦拉
<thetish@google.com>。
第40届国际机器学习会议文集
学习,美国夏威夷州檀香山。 PMLR 202, 2023。版权所有
2023 年,作者。

ODE 允许明确表征网络输出 (Jacot 等人,2018;Lee 等人,2019;
Yang,
2021) 。这种方法已被全面使用,以获得
深入了解广泛的神经网络,但它有一个主要限制:模型的参数是线性的,
因此它描述了
一个动态相对微不足道的政权,无法捕获
特征学习,无法准确表示类型
实践中经常观察到的复杂训练现象。
虽然其他大宽度缩放机制可以保留一些
非线性并允许某些类型的特征学习 (Bordelon & Pehlevan,2022;
Yang et al.,2022) ,例如
方法往往侧重于较小的学习率或
连续时间动力学。相比之下,最近的实证
这项工作强调了在训练具有大学习率的实际网络时非线性离散动力学产生
的许多重要现象 (Neyshabur 等人,
2017年;吉尔默等人,2022;戈巴尼等人,2019;福雷特等人,
2022) 。特别是,许多实验表明
网络显示渐进锐化的趋势
曲率趋向稳定边缘,其中
Hessian 损失的最大特征值增加
训练过程直到稳定在等于
大致二除以学习率,对应
到梯度下降的最大特征值
收敛于二次势 (Wu et al., 2018;Giladi
等,2020;科恩等人,2022a;b)。
为了更好地理解这种行为,我们
引入一类二次回归模型,它显示所有相关的现象学,但又足够
简单
承认数字和分析的理解。楷模
之前已经对这种类型进行了研究,以便
了解 NTK 之外发生的其他现象
(Zhu et al., 2022a; Roberts et al., 2022)。在特定的低维环境中,
我们证明了最大
NTK 特征值收敛到 (接近)稳定边缘,
我们凭经验证明渐进锐化
稳定边缘通常发生在大数据点、大模型极限中。最后,我们进行数值计算
真实神经网络的特性分析和使用
我们的理论分析工具表明,“在野外”的稳定边缘行为表现出一些相同的
特征
模式作为理论模型。

1.1. 并发工作

几项同时进行的工作研究了有关稳定性边缘的类似问题,并形成了很好地补充我们在此得出的结论的见解。

朱等人。(2022b)提出了一个4度函数(或8度目标)形式的极简模型,并证明锐度收敛到接近但不等于稳定边缘的值。虽然这些结果与我们在低维设置中的结果相似,但我们发现2阶函数可能会出现稳定性边缘行为,并且最终的锐度与收敛阈值的差异取决于问题的精确量,即我们计算。

达米安等人。(2022)基于Hessian顶部本征模态的非线性自相互作用及其与其他本征模态的有效相互作用,开发了一种稳定边缘附近锐度的解释模型。当清晰度接近稳定阈值时,低维动力学捕获了短时间尺度上的完整行为,并且与我们在模型的低维极限下开发的微分方程非常相似。我们的设置不太通用,因为它仅限于二次模型,但它为我们提供了一些超出Damian等人的能力的额外见解。(2022)。特别是,我们能够证明渐进锐化通常发生在高维模型中,至少在早期,并且我们发现最终曲率与收敛阈值偏离了可预测的量。

2. 二次回归模型

我们首先定义基本的二次回归模型。给定 P 维参数向量 θ , D 维输出 $f(\theta)$ 由下式给出

$$f(\theta) = y + G\theta + \frac{1}{2} Q(\theta, \theta). \quad (1)$$

这里 y 是 D 维向量, G 是 $D \times P$ 维矩阵, Q 是最后两个索引中对称的 $D \times P \times P$ 维张量。即 $Q(\cdot, \cdot)$ 取两个 P 维向量作为输入,输出一个 D 维向量,验证 $Q(\theta, \theta) = \theta^T Q \theta$ 。如果 $Q = 0$,则模型对应于线性回归。我们可以使用以下方法恢复 $D \times P$ 维雅可比行列式 $J = \partial f / \partial \theta$

$$G_{ai} = \frac{\partial f_a}{\partial \theta_i} \bigg|_{\theta=0}, \quad Q_{aij} = \frac{\partial^2 f_a}{\partial \theta_i \partial \theta_j} \bigg|_{\theta=0} = \frac{\partial^2 J_a}{\partial \theta_i \partial \theta_j} \bigg|_{\theta=0}. \quad (2)$$

任何模型都可以通过围绕任何可微点进行二阶泰勒展开来转换为二次回归模型。之前已经研究了浅MLP的二次展开(Bai & Lee, 2020; Zhu et al., 2022a),并且研究了小 Q 的微扰理论

在罗伯茨等人。(2022)。其他相关模型详见附录A。我们将提供证据证明,即使是随机的、非结构化的二次回归模型也会导致稳定性边缘(EOS)行为。

在这项工作中,我们重点关注MSE损失设置。更明确地,给定 D 维目标向量 y ,损失 $L(\theta)$ 可以写成残差 $z = f(\theta) - y$:

$$L(\theta) = \frac{1}{2} z^T z = \frac{1}{2} (y - f(\theta))^T (y - f(\theta)) = \frac{1}{2} y^T y - y^T G \theta + \frac{1}{2} \theta^T Q \theta. \quad (3)$$

正如我们将要展示的,在这种情况下,动力学可以单独用 z 和 J 来写在函数空间中。

3. 低维动力学

我们首先关注 $D = 1$ (单个数据点)的动态。不失一般性,我们可以将损失函数写为:

$$L(\theta) = \frac{1}{2} \theta^T Q \theta - E \theta. \quad (4)$$

对于定义问题的 $P \times P$ 矩阵 Q 和标量 E 。我们将根据(标量)残差 $z = \theta^T Q \theta - E$ 和 $1 \times P$ 维雅可比行列式 $J = \partial z / \partial \theta$ 来分析动力学。我们特别对标量曲率 JJ^T (神经正切核,或NTK)的动力学感兴趣。

3.1. 梯度流

我们首先考虑损失 L 相对于参数 θ 的梯度流(GF)动力学。给定比例因子 η , GF动力学由下式给出

$$\dot{\theta} = -\eta \nabla_{\theta} L = -\eta z \frac{\partial z}{\partial \theta} = -\frac{\eta}{2} \theta^T Q \theta - E \theta. \quad (5)$$

z 和 J 的动力学接近:

$$\dot{z} = -\eta(JJ^T)z, \quad \dot{J} = -\eta JQ. \quad (6)$$

曲率是标量,由NTK JJ^T 描述。在这些坐标中,我们有 $E = JQ + J - 2z$,其中 Q 表示Moore-Penrose伪逆。

我们可以进行一些坐标变换来简化动力学。我们定义:

$$z \sim \eta z, \quad T(k) = \eta J Q^k J. \quad (7)$$

我们注意到 $T(0) = \eta J J$ - 即通过比例因子归一化的曲率。动力学方程为:

$$\frac{dz}{dt} = -z T(0), \quad \frac{dT(k)}{dt} = -2z T(k+1). \quad (8)$$

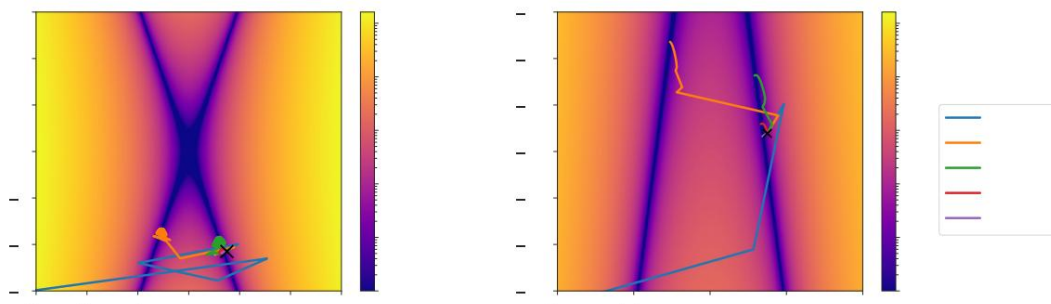


图 1. $D = 1$ 损失景观 $L(\cdot)$ 作为参数 θ 的函数, 其中 $P = 2, E = 0$ 且 Q 的特征值为 1 和 -0.1。GD 轨迹 (初始化为 $(1.5, -4.32)$, 用 x 标记) 收敛到曲率比初始化时更大的最小值, 因此显示出渐进锐化 (左)。在两步动力学中, 我们只考虑偶数迭代次数, 在稳定边缘附近表现出较少的振荡 (右)。

为了研究稳定性行为的边缘, 我们需要进行初始化, 允许曲率 $T(0)$ 随着时间的推移而增加。这种现象称为渐进锐化。

渐进锐化已被证明在机器学习模型中无处不在 (Cohen 等人, 2022a), 因此任何有用的现象学模型也应该显示它。在附录 B.1 中, 我们确认模型在一系列初始条件下显示出渐进锐化。

3.2. 梯度下降

我们现在将证明, 对于具有学习率 η 的梯度下降 (GD) 动力学, 存在一系列初始化, 其中 $D = 1$ 二次回归模型显示稳定边缘 (EOS) 行为。对于该模型, 我们将 EOS 行为定义为一种设置, 其中动态导致 λ_{\max} (NTK JJ 的最大特征值) 保持接近临界值 $2/\eta$ 。该临界值对应于最大学习率, 其中动力学在损失 L 的最小值附近呈指数收敛。该边界对应于 $T(0) = 2$ - 因此将 $T(0)$ 解释为重新缩放的曲率。

我们将证明存在一个模型族, 对于任何距离, 都有一个模型可以导致 EOS 收敛到远离稳定边缘的距离。对于非平凡初始化量中的任何初始化。

此外, 我们将证明这个初始化体积在一组自然坐标中的模型之间是一致的。

我们将在 $D = 1$ 模型中证明 EOS 的这种形式, 并找到经验证据证明这也适用于大 D 。请注意, 我们根据最大 NTK 特征值而不是最大损失 Hessian 特征值来定义 EOS (Cohen 等人, 2022a); 请参阅附录 A.1 了解为什么这适用于 MSE 损失的讨论。

当 Q 同时具有正和负特征值时, 损失

景观是双曲抛物面的正方形 (图 1, 左)。正如梯度流分析所表明的, 这会导致一些轨迹在收敛之前增加其曲率。最终的曲率自然取决于初始化和学习率。分析梯度下降动力学的挑战之一是, 它们在大学习率的最小值周围快速而剧烈地振荡。

缓解此问题的一种方法是仅考虑所有其他步骤 (图 1, 右)。我们将利用这一观察结果直接分析梯度下降 (GD) 动力学, 以找到这些轨迹显示稳定边缘行为的配置。

令 θ_t 为步骤 t 的参数。梯度下降方程由下式给出:

$$\theta_{t+1} - \theta_t = -\eta \nabla \theta L = 2 \frac{\eta}{\lambda_{\max}} \theta_t - Q \theta_t - E Q \theta_t. \quad (9)$$

在 $z \sim T$ 坐标下, 梯度下降方程变为 (附录 B.2):

$$1 z_{t+1} - z_t = -z_t T_t(0) + (-z_t^2 T_t(1)) \quad (10)$$

$$T_{t+1}(k) - T_t(k) = -z_t(2T_t(k+1) - z_t T_t(k+2)). \quad (11)$$

如果 Q 是可逆的, 则 $E \sim \eta E = T(-1) - 2 \sim z$ 。根据定义, $T_t(0) = \eta J_t J = \eta \lambda_{\max, t}$ 是 (重新调整的) NTK 特征值。在这些坐标中, EOS 行为对应于当 $z \sim t$ 趋向 0 时使 $T_t(0)$ 保持在值 2 附近的动态。

在本节的其余部分中, 我们将重点关注 $P = 2$ 的情况。正如我们将看到的, 这让我们可以单独用 $z \sim$ 和 $T(0)$ (残差和曲率) 来编写动力学。

3.2.1. 减少弹射动力

如果 Q 的特征值为 $\{-\omega, \omega\}$, 且 $E \sim 0$, 则该模型等效于单隐层线性网络

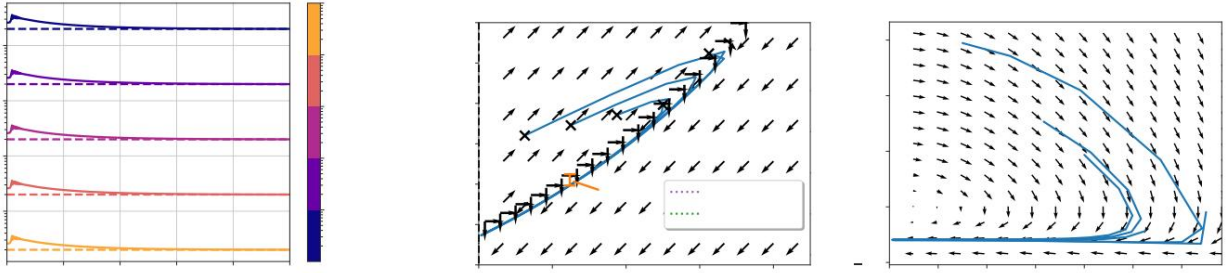


图 2.左图:绘制 $D = 1$ 、 $P = 2$ 模型的曲率显示了不同步长 ($= 5 \cdot 10^{-3}$) 下的 EOS 行为。中:绘制每隔一个迭代,我们看到各种初始化 (黑色 x 的), $z \sim T(0)$ 空间中的轨迹保持在零斜线 ($\sim z, fz \sim (z)$) 附近 - $z \sim T(0)$ 所在的曲线 $t+2 - z \sim t = 0$ 。箭头表示动态流的方向。右:将变量更改为 $y = T(0) - fz \sim (z)$ 显示快速集中到接近恒定、小、负 y 的曲线。

具有一个训练数据点 (附录 A.2) 也称为弹射器阶段动力学。该模型不会表现出锐化或稳定性边缘行为 (Lewkowycz 等人, 2020)。我们可以在 $z \sim T(0)$ 坐标中证明这一点。

不失一般性,我们假设特征值为 $\{-1, 1\}$ - 这可以通过重新调整 $z \sim$ 来实现。

我们可以仅根据 $z \sim$ 和曲率 $T(0)$ 重写动力学 (附录 B.3):

$$z \sim t+1 - z \sim t = -z \sim t T(0) + \frac{1}{2} (\sim z_t^2) (2 \sim z_t + E \sim) \quad (12)$$

$$T(t+1(0) - T(t(0) = -2 \sim z_t (2 \sim z_t + E \sim) + z \sim t^2 T(0)。 \quad (13)$$

对于 $E \sim = 0$, 我们可以看到 $\text{sign}(\Delta T(0)) = \text{sign}(T(t(0) - 4)$, 如 (Lewkowycz et al., 2020) 中所示 - 因此收敛意味着严格减小曲率。当 $E \sim = 0$ 时, 存在曲率可以增加的区域 (附录 B.3)。然而, 仍然没有 EOS 行为 - 如果 JJ 从 $2/\eta$ 开始, 没有机制可以将其稳定在边缘附近。

3.2.2. 稳定制度的边缘

在本节中, 我们考虑 Q 有两个特征值的情况 - 其中一个较大且为正, 另一个较小且为负。不失一般性, 我们假设 Q 的最大特征值为 1。我们用 $-$ 表示第二个特征值, 其中 $0 < \leq 1$ 。用这种表示法, 我们可以将动力学方程 (附录 B.3) 写为

$$1 z \sim t+1 - z \sim t = -z \sim t T(0) + (\sim z_t^2) ((1 -) T(0) + (2 \sim z_t + E \sim)) \quad (14)$$

$$T(t+1(0) - T(t(0) = -2 \sim z_t ((2 \sim z_t + E \sim) + (1 -) T(0)) + \sim z_t^2 T(0) + (-1) (T(0) - E \sim - 2 \sim z_t)。$$

$$(15) \quad \text{我们定义多项式 } p_z \sim (\sim z, T) \text{ 和 } p_T \sim (\sim z, T) \text{ 使得}$$

对于较小的轨迹, λ_{\max} 最初远离 $2/\eta$, 但会向 $2/\eta$ 收敛 (图 2, 左)

换句话说, EOS 行为。我们使用了各种步长 η , 但以 $(\eta z_0, \eta T(0))$ 对初始化, 以显示 $z \sim T(0)$ 坐标的通用性。

为了定量地理解渐进的锐化和稳定性边缘, 查看两步动态是有用的。研究两步动力学的另一个动机来自于对具有大步长 λ 的线性最小二乘法 (即线性模型) 的梯度下降分析。对于每个坐标 θ , 一步和两步动力学为

$$\sim \theta t+1 - \sim \theta t = -\lambda \sim \theta t \text{ 且 } \sim \theta t+2 - \sim \theta t = (1 - \lambda)^2 \sim \theta t \quad (16)$$

虽然动力学在 $\lambda < 2$ 时收敛, 但如果 $\lambda > 1$, 则一步动力学在接近最小值时会交替符号, 而两步动力学保持 θ 的符号并且轨迹不表现出振荡。

同样, 在双参数模型中绘制每隔一个迭代可以更清楚地展示现象学。对于小的动态显示 (Li et al., 2022) 中描述的不同阶段: $T(0)$ 的初始增加, $z \sim$ 的缓慢增加, 然后 $T(0)$ 的减少, 最后是 $T(0)$ 的缓慢减少 $z \sim$ 而 $T(0)$ 保持在 2 附近 (图 2, 中间)。

不幸的是, 方程 14 和 15 定义的动力学的两步版本更加复杂 - 它们在 $T(0)$ 中是 3 阶, 在 $z \sim$ 中是 9 阶; 更详细的讨论参见附录 B.4。然而, 我们可以分析 $z \sim$ 趋于 0 时的动力学。为了理解 EOS 行为的机制, 了解两步动力学的零斜线是有用的。零斜线是由初始化组成的曲线 ($\sim z, fz \sim (z)$) 和 ($\sim z, ft \sim (z)$), 使得 $z \sim$ 和 $T(0)$ 分别在两步后不变。

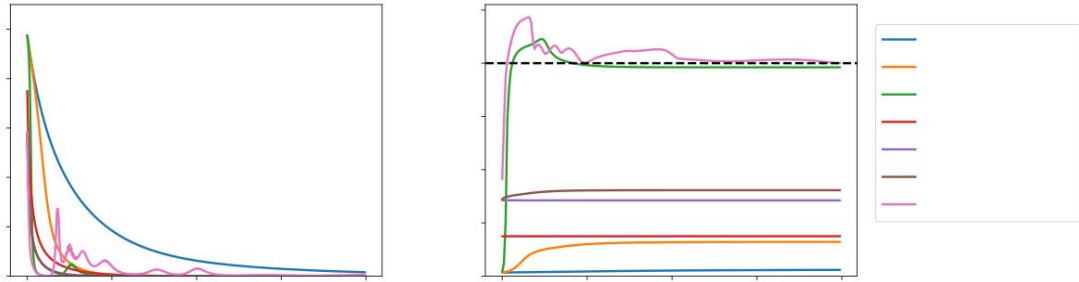


图 3.二次回归模型中的 GD 动态。对于较小的 σ_z (左), 损失函数是单调的。对于较大的 σ_z , 最终的 λ_{\max} 较高 (右)。当锐化使 $\eta \lambda_{\max}$ 接近 2 时, 非线性效应会引发 EOS 行为。

$z_{t+1} - z_t = p z_t(\tilde{z}_t, T_t)$ 且 $T_{t+1}(0) - T_t(0) = p T(\tilde{z}_t, T_t)$ 。
那么 $f z_t(\tilde{z})$ 和 $f T(\tilde{z})$ 服从隐式方程: $p z_t(p z_t(\tilde{z}, f z_t(\tilde{z})),$

$$\begin{aligned} p T(\tilde{z}, f z_t(\tilde{z})) &= 0, \quad p T(p z_t(\tilde{z}, f T(\tilde{z})), \\ p T(\tilde{z}, f T(\tilde{z})) &= 0 \end{aligned} \quad (17)$$

这些表达式分别是 $f z_t(\tilde{z})$ 和 $f T(\tilde{z})$ 的三次多项式 - 因此当 $\tilde{z} \rightarrow 0$ 时, 存在三种可能的解。我们对经过 $\tilde{z} = 0, T(0) = 2$ 即 EOS 对应的临界点。

附录 B.4 中详细的计算表明, 两个零斜线之间的距离是线性的, 因此它们随着接近 0 而变得接近。(图 2, 中间)。此外, 轨迹停留在 $f z_t$ 附近 - 这会产生 EOS 行为。

这表明零斜线附近的动力学很慢, 并且轨迹似乎正在接近吸引子。我们可以通过将变量更改为 $y_t \equiv T_t(0) - f z_t(\tilde{z}_t)$ - 到 \tilde{z} 零斜线的距离来找到吸引子的结构。

我们可以通过将动态扩展到 \tilde{z} 和 y 的最低阶来建立直觉。附录 B.5 中的直接计算给出了近似值:

$$z_{t+2} - z_t = 2 y_t z_t + O(y_t^2 z_t^2) + O(y_t z_t^2) \quad (18)$$

$$y_{t+2} - y_t = -2(4 - 3 + 4 y_t z_t^2 + 2 y_t^2 - 4 y_t z_t^2) + O(z_t^3 + y_t^2 z_t^2) \quad (19)$$

我们立即看到, 对于小 y , \tilde{z} 变化缓慢 - 因为当 $y = 0$ 时, 我们选择了 $z_{t+2} - z_t = 0$ 的坐标。我们还可以看到 $y_{t+2} - y_t$ 是 $O(y)$, 因为 $y_t = 0$ - 因此对于较小的 y 动态也很慢。此外, 我们看到该项是负的 - \tilde{z} 的变化往往, 会驱动 y (因此 $T(0)$) 减小。 y_t 项的系数也为负; y 的动态趋于收缩。 z 的系数 2 关键是收缩行为以与 \tilde{z} 成正比的速率将 y 带到 $O(y)$ 固定点, 而 \tilde{z} 的动态则与然后 \tilde{z} 收敛到 0 (图 2, 右)。

来自最低阶项的直觉可以被形式化, 并给我们一个 $\lim_{t \rightarrow \infty} y_t = -1/2$ 的预测, 由以下定理证实 (附录 B.6 中的证明): 定理 3.1。考虑二次回归模型

$D = 1, P = 2, E = 0$, Q 具有特征值 $\{-1, 1\}$ 。 > 0 和一个独立的邻域, 存在一个 c 使得对于 $0 < \epsilon$ 以及所有初始化使得 $(\tilde{z}_0, U \in \mathbb{R}^3 \quad \eta \lambda_{\max}(0), \eta) \in U$, 锐度收敛到稳定边缘附近的值:

$$\lim_{t \rightarrow \infty} \lambda_{\max}(t) = (2 - 1/2)/\eta + O(\epsilon^2). \quad (20)$$

因此, 与弹射器相模型不同, 小粒子可证明具有 EOS 行为 - 其机制通过 $\tilde{z} \sim y$ 坐标变换可以很好地理解。

4. 高维动力学

在本节中, 我们分析大 D 的动态。我们关注随机初始化的 Q , 并证明渐进锐化是普遍存在的。然后, 我们证明二次回归模型在一系列初始化范围内显示出稳定性行为的边缘。

4.1. 梯度流动力学

MSE 损失函数 L 上的梯度流动力学通常可以写为

$$\dot{\theta} = -\frac{\partial L(\mathbf{z})}{\partial \theta} \mathbf{z}. \quad (21)$$

其中 J 是 $D \times P$ 维雅可比行列式, \mathbf{z} 是残差 $f(\theta) - y$ 的 D 维向量。对于二次回归模型, \mathbf{z} 和 J 的动态再次接近:

$$\dot{\mathbf{z}} = J \dot{\theta} = -J J \mathbf{z}, \quad J \dot{\theta} = -Q(J \mathbf{z}, \cdot). \quad (22)$$

当 $Q = 0$ (线性化状态) 时, J 为常数, 则动力学在 \mathbf{z} 上呈线性, 并由 $D \times D$ 矩阵 $J J$ (经验 NTK) 的特征结构控制。

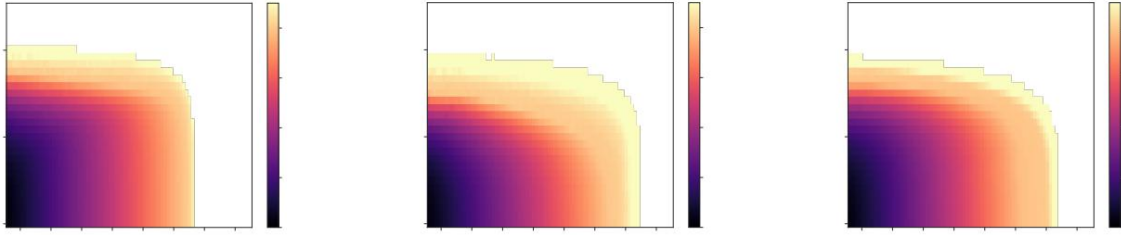


图 4. σ_z/σ_J 用于二次回归模型的 J 相平面,适用于各种 D 和 P 。模型使用 100 个随机种子进行初始化。每个 σ_z 、 σ_J 对并迭代直至收敛。对于同时发生锐化和非线性 z 运动,我们绘制 NTK J, J 的中值 λ_{\max} 。对于中间 σ_z , 态的每一对 σ_z , σ_J 轨迹趋于收敛,因此 NTK 的 λ_{\max} 接近 $2/\eta$ (EOS)。

我们对 GF 下进行渐进锐化的设置感兴趣。我们可以研究 J, J 早期的最大特征值 λ_{\max} 的动态,以进行随机初始化。在附录 C.1 中,我们证明了以下定理:

定理 4.1. 让 z 、 J 和 Q 使用均值为零且方差为 σ 的 iid 元素进行初始化。令 λ_{\max} 为 J, J 的最大特征值。在大 D 和 P 的限制下,固定比率 D/P , 初始化时我们有

$$E[\lambda_{\max}(0)] = 0, E[\lambda'_{\max}(0)]/E[\lambda_{\max}(0)] = \sigma^2/z^2 \quad (23)$$

其中 E 表示初始化时对 z 、 J 和 Q 的期望。

与 $D = 1$ 的情况非常相似,定理 4.1 表明很容易找到显示渐进锐化的初始化 - 并且增加 σ_z 使锐化更加突出。

4.2. 梯度下降动力学

我们现在考虑有限步长梯度下降 (GD) 动力学。 θ 的动力学由下式给出:

$$\theta_{t+1} = \theta_t - \eta J_t z_t. \quad (24)$$

在这种情况下,动力学方程可以写为

$$z_{t+1} - z_t = -\eta J_t J_t z_t + \frac{1}{2} \eta^2 2Q(J_t z_t, J_t z_t) \quad (25)$$

$$J_{t+1} - J_t = -\eta Q(J_t z_t, \cdot). \quad (26)$$

如果 $Q = 0$, 动力学会减少为二次势中的离散梯度下降 - 当且仅当 $\lambda_{\max} < 2/\eta$ 时收敛。

一个迫在眉睫的问题是: η 何时影响动力学? 鉴于公式 25 它的 η 和 z 幂比第一项更高,我们可以推测

项的大小之比 r_{NL} 与 $\|z\|^2$ 和 η 成正比。附录 C.2 中的计算表明,对于随机旋转不变初始化,我们有:

$$r_{NL} = \frac{E[\|\frac{1}{2} \eta 2Q(J_0 z_0, J_0 z_0)\|^2]^{1/2}}{E[\|\eta J_0 J_0 z_0\|^2]} = \frac{1}{2} \eta \sigma_z D, \quad (27)$$

其中初始化统计量的定义如定理 4.1 中所示。这证实了增加学习率和残差幅度 $\|z\|$ 增加了动力学与 GF 的偏差,并表明非线性量对 J 不敏感。

我们可以在 GD 方程的动力学中看到这种现象 (图 3)。在这里,我们为定理 4.1 中类型的随机初始化绘制了不同的轨迹,其中 $D = 60$ 、 $P = 120$ 和 $\eta = 1$ 。随着 σ_z 增加,速率 λ_{\max} 也增加 (如定理 4.1 所示),并且当 σ_z 为 0 时 (1),动力学是非线性的 (如 r_{NL} 所预测) 并且出现 EOS 行为。这表明方程 25 中的第二项对于 λ_{\max} 的稳定至关重要。

我们可以通过在多个种子上初始化各种 η 、 D 、 P 、 σ_z 和 σ_J 并绘制最终达到的 λ_{\max} 的相图来更普遍地确认这一点。我们可以通过重新调整参数和初始化来简化绘图。在重新调整的变量中

$$z \sim \eta z, J \sim \eta^{1/2} J, \quad (28)$$

动力学等效于方程 25 和 26 (其中 $\eta = 1$)。与方程 56-57 的 $z \sim T(0)$ 模型一样,重新缩放坐标中的 λ_{\max} 等于未缩放坐标中的 $\eta \lambda_{\max}$ 。我们还可以为 z 和 J 定义重新缩放的初始化。如果我们设置

$$\sigma_z = \sigma_z/D, \sigma_J = \sigma_J/(DP)^{1/4}, \quad (29)$$

那么我们有 $r_{NL} = \sigma_z$, 它可以更轻松地在 (D, P) 对之间进行比较。

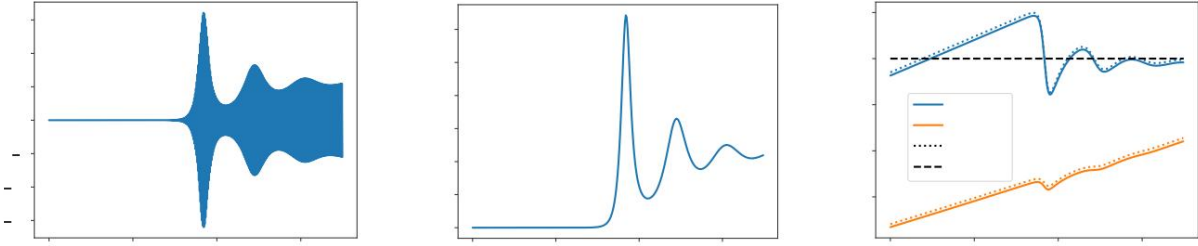


图 5.在 CIFAR 上训练的 FCN 显示了锐化和稳定边缘行为的多个周期。 z_1 , 训练集残差 $f(X, \theta) - Y$ 到顶部 NTK 特征模 v_1 的投影, 其幅度增加并在 0 附近振荡 (左)。每两步绘制动态图可以消除高频振荡 (中)。最大特征值 λ_1 多次跨越稳定边缘, 但第二大特征值 λ_2 仍低于稳定边缘 (右)。Hessian 特征值的动力学类似 (虚线)。

使用此初始化方案, 我们可以针对每个 σ_z 、 σ_J 对进行 100 次独立随机初始化, 将达到的 λ_{\max} 最终值绘制为 σ_z 和 σ_J 的函数 (图 4)。

我们看到关键是 $r_{NL} = \sigma_z$ 为 $O(1)$ - 对应于渐进锐化和初始化附近的非线性动态。特别是, 在 EOS 处收敛的小 σ_J 值的初始化对应于首先锐化, 然后稳定在 $\lambda_{\max} = 2/\eta$ 附近的轨迹。

大 σ_z 和大 σ_J 动态发散。在很宽的 σ_z 范围内有一个小范围的初始 σ_J , 其最终 $\lambda_{\max} \approx 2/\eta$; 这些对应于在 EOS 附近初始化的模型, 这些模型保持在 EOS 附近。

这表明渐进锐化和稳定性边缘并不是神经网络模型的独特特征, 并且可能是高维非线性模型中学习的更普遍的属性。

5. 与现实世界模型连接

在本节中, 我们在“现实世界”模型中进行数值实验, 并将其行为与我们在简化模型上的理论进行比较。接下来 (Cohen 等人, 2022a), 我们使用来自 CIFAR10 的 5000 个示例的平方损失, 以学习率 10^{-2} 训练了一个 2 隐藏层 tanh 网络, 该设置显示了稳定性行为的边缘。

接近 EOS 开始时, 我们使用 Lanczos 方法近似计算了 JJ 的最大特征值 λ_1 及其相应的特征向量 v_1 (Ghorbani 等人, 2019; Novak 等人, 2019)。我们使用 v_1 来计算 $z_1 = 1/z$, 其中 z 是神经网络函数 f 、训练输入 X 、标签 Y 和参数 θ 的残差 $f(X, \theta) - Y$ 的向量。NTK 中的 EOS 行为类 V 类似于 (Cohen 等人, 2022a) 中针对完整 Hessian 定义的 EOS 行为 (图 5, 左和右)。再次, 每隔一步绘制轨迹可以消除高频振荡 (图 5, 中)。与 $D = 1$, $P = 2$ 模型不同, 存在多个交叉点

临界线 $\lambda_{\max} = 2/\eta$ 线。

有证据表明二次回归模型的低维特征可用于解释 EOS 行为的某些方面。我们通过自动微分凭经验计算输出 $f(x, \theta)$ 的二阶导数。我们用 $Q(\cdot, \cdot)$ 表示结果张量。我们可以使用矩阵向量积来计算矩阵 $Q_1 \equiv v_1 \cdot Q(\cdot, \cdot)$ 的谱, 它是 Q 的输出在 v_1 方向上的投影, 而无需在内存中实例化 Q (图 6, 左)。该图显示频谱从步骤 3200 到 3900 (我们绘图的范围) 变化不大。这表明当显示这些 EOS 动态时, Q 不会发生太大变化。我们还可以看到 Q 在 v_1 方向比随机方向大得多。

让 y 定义为 $y = \lambda_1 \eta - 2$ 。绘制 z_1 与 $2yz$ 的两步动态图, 我们看到了显着的一致性 (图 6, 中)。这与我们的简化模型中 z 的动力学形式相同。也可以通过使用 $y = \lambda_1 \eta - 2$ 的固定雅可比行列式迭代方程 25 两次并丢弃 η 中的高阶项来找到它。这表明, 在这种特定的 EOS 行为期间, 就像在我们的简化模型中一样, 特征值的动态比特征基中的任何旋转更重要。

y 的动态更加复杂; $y_{t+2} - y_t$ 与 z 反相关, 但 y 和 z_1 不存在低阶函数形式 (附录 D.2)。我们可以通过绘制 $\eta^2 Q_1(Jz_1 v_1, Jz_1 v_1)$ (从 v_1 方向对 z_1 动力学的非线性贡献) 和 $\lambda_1 z_1$ (线性化贡献) 的比率来深入了解稳定性, 并将其与 y (图 6, 右)。该比率在初始锐化期间很小, 但在曲率第一次减小之前不久变为 $O(1)$ 。在其余的动力学过程中它仍然是 $O(1)$ 。这表明顶部本征模动力学对其自身的非线性反馈对于理解 EOS 动力学至关重要。

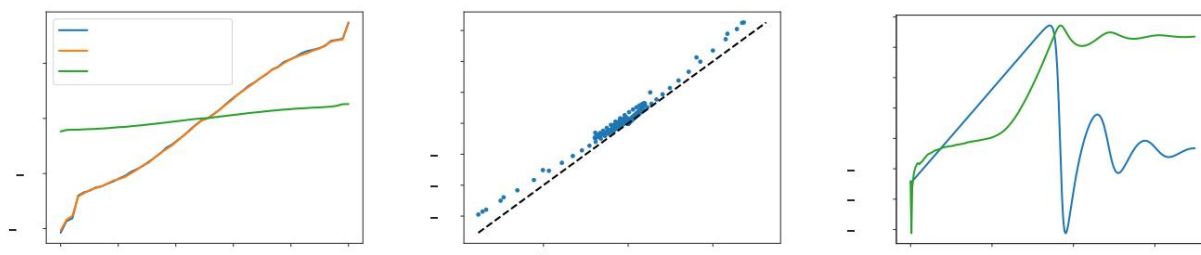


图 6.在 CIFAR10 (左)上训练的 FCN 的稳定边缘动态期间, Q 近似恒定。最大特征方向 v_1 (蓝色和橙色)上的投影大于随机方向 (绿色)上的投影。两步差 $(z_1)_{t+2} - (z_1)_t$ 可以很好地近似为 $2z_1 y$ (中), 即具有固定特征基的模型的前导项。非线性动力学贡献 $\eta^2 Q_1(Jz_1 v_1, Jz_1 v_1)$ 在锐化期间很小,但在顶部特征值下降之前立即变大 (右) - 正如简单模型中的情况。

6. 讨论

6.1. 从二次回归模型中吸取的教训

从二次回归模型中学到的主要教训是,诸如渐进锐化 (对于 GF 和 GD) 和稳定性边缘行为 (对于 GD) 之类的行为可能是基于高维梯度的非非训练的共同特征。线性模型。事实上,这些现象可以在简单的设置中揭示出来,而无需与深度学习模型有任何联系:通过轻微调整我们的简化模型,它对应于 1 个数据点和 2 个参数,可以证明 EOS 行为。结合 CIFAR 模型的分析表明,一般机制可能具有低维描述。

真实模型的二次近似可以定量地捕获 EOS 行为的早期特征 (最初返回到 $\lambda_{\max} < 2/\eta$), 但不一定捕获后续振荡的幅度和周期 - 这些需要更高阶的项 (附录 D.3)。尽管如此,二次近似确实正确地描述了许多定性行为,包括 λ_{\max} 收敛到在 $2/\eta$ 附近振荡的极限二周期,平均值低于 $2/\eta$ 。在简化的二参数模型中,可以分析预测收敛时的最终值,而且我们确实发现它与值 $2/\eta$ 略有偏差。此外,理论和现实模型都表明曲率由

由 z 通过 Q_1 介导的低维反馈机制。

$\frac{1}{2}$

这项工作中研究的所有模型的一个关键特征是,查看所有其他迭代 (两步动力学) 极大地有助于从理论上和经验上理解模型。在稳定边缘附近,这使得顶部本征模的变化很小。在简化模型中,慢 z -动力学 (以及相关的慢 $T(0)$ 动力学) 允许进行详细的理论分析; 在

在 CIFAR 模型中,两步动力学在 z_1 和 λ_{\max} 上都缓慢变化。这些微小变化的定量比较可能有助于揭示解释其他系统和场景中 EOS 行为的任何通用机制/规范形式。

6.2. 未来的工作

未来工作的一个途径是定量地理解大 D 和 P 的二次回归模型中的渐进锐化和 EOS 行为。特别是,计算稳定边缘区域中的最终偏差 $2 - \eta \lambda_{\max}$ 作为 σ_z 的函数、 σ_J 和 D/P 仍然是一个有趣的悬而未决的问题。了解高阶项如何影响训练动态也很有用 - 特别是 $y = 2$ 附近振荡的细节。

最后,我们的分析没有涉及模型的特征学习方面。在二次回归模型中,特征学习被编码在 J 和 z 之间的关系中,特别是 z 和 J 的特征结构之间的关系。了解 Q 如何调节这两个量的动态可以为理解特征学习提供定量基础,这与现有的理论方法是互补的 (Roberts 等人, 2022 年; Bordelon 和 Pehlevan, 2022 年; Yang 等人, 2022 年)。

参考

Bai, Y. 和 Lee, JD 超越线性化:关于宽神经网络的二次和高阶逼近。国际学习表达会议, 2020 年 3 月。

Bordelon, B. 和 Pehlevan, C. 宽神经网络中核演化的自治动态场理论, 2022 年 5 月。

Cohen, J., Kaur, S., Li, Y., Kolter, JZ 和 Talwalkar, A.

- 神经网络上的梯度下降通常发生在稳定边缘。国际学习表征会议,2022 年 2 月a。
- Cohen, JM.Ghorbani, B.,Krishnan, S.,Agarwal, N., Medapati, S.,Badura, M.,Suo, D.,Cardoze, D.,Nado, Z., Dahl, GE 和 Gilmer, J.处于稳定边缘的自适应梯度方法,2022 年 7 月b。
- Damian, A.,Nichani, E. 和 Lee, JD 自稳定:稳定边缘梯度下降的隐含偏差,2022 年 9 月。
- Foret, P.,Kleiner, A.,Mobahi, H. 和 Neyshabur, B. 锐度感知最小化可有效提高泛化能力。国际学习表征会议,2022 年 4 月。
- Ghorbani, B.,Krishnan, S. 和Xiao, Y.通过 Hessian 特征值密度进行神经网络优化的研究。第 36 届国际机器学习会议论文集,第 2232-2241 页。PMLR, 2019 年 5 月。
- Giladi, N.,Nacson, MS.Hoffer, E. 和 Soudry, D. 在稳定性边缘:如何调整超参数以在神经网络异步训练中保持最小选择?第八届学习表征国际会议,2020 年 4 月。
- Gilmer, J.,Ghorbani, B.,Garg, A.,Kudugunta, S.,Neyshabur, B., Cardoze, D.,Dahl, GE,Nado, Z. 和 Firat, O. 训练损失曲率视角深度学习模型的不稳定。国际学习表征会议,2022 年 3 月。
- Huang, J. 和 Yau, H.-T.深度神经网络的动力学和神经切线层次结构。第 37 届国际机器学习会议论文集,第 4542-4551 页。PMLR,2020 年 11 月。
- Jacot, A.,Gabriel, F. 和 Hongler, C. 神经正切核:神经网络的收敛和泛化。神经信息处理系统进展31,第 8571-8580 页。Curran Associates, Inc.,2018。
- Lee, J.,Xiao, L.,Schoenholz, S.,Bahri, Y.,Novak, R.,Sohl-Dickstein, J. 和 Pennington, J. 任何深度的宽神经网络在梯度下降下演化为线性模型。神经信息处理系统进展32,第 8570-8581 页。Curran Associates, Inc., 2019。
- Lewkowycz, A.,Bahri, Y.,Dyer, E.,Sohl-Dickstein, J. 和 Gur-Ari, G. 深度学习的大学习率阶段:弹射器机制。2020 年 3 月。
- Li, Z.,Wang, Z. 和 Li, J. 沿着 GD 轨迹分析锐度:渐进锐化和稳定边缘, 2022 年 7 月。
- Neyshabur, B.,Bhojanapalli, S.,Mcallester, D. 和 Srebro, N. 探索深度学习的泛化。神经信息处理系统进展第 30 页,第 10 页。5947-5956。Curran Associates, Inc.,2017。
- Novak, R.,Xiao, L.,Hron, J.,Lee, J.,Alemi, AA,Sohl-Dickstein, J. 和 Schoenholz, SS 神经切线: Python 中的快速且简单的无限神经网络。arXiv:1912.02803 [cs, stat],2019 年 12 月。
- Roberts, DA.Yaida, S. 和 Hanin, B. 深度学习理论原理。2022 年 5 月。doi:10.1017/9781009023405。
- Wu, L., Ma, C. 和 E, W. SGD 如何在过参数化学习中选择全局最小值:动态稳定性视角。《神经信息处理系统进展》,第 31 卷。Curran Associates, Inc., 2018 年。
- Yang, G. 张量程序 I:任何架构的宽前馈或循环神经网络都是高斯过程。arXiv:1910.12478 [cond-mat,物理:数学-ph],2021 年 5 月。
- Yang, G.,Hu, EJ.Babuschkin, I.,Sidor, S.,Liu, X.,Farhi, D.,Ryder, N.,Pachocki, J.,Chen, W. 和 Gau, J. 张量项目 V:通过零样本超参数传输调整大型神经网络,2022 年 3 月。
- Zhu, L.,Liu, C.,Radhakrishnan, A. 和 Belkin, M. 用于理解神经网络动力学的二次模型,2022 年 5 月a。
- Zhu, X.,Wang, Z.,Wang, X.,Zhou, M. 和 Ge, R.用极简示例理解稳定边缘训练动态,2022 年 10 月b。

A. 与其他型号的连接

A.1. Hessian 与 NTK 最大特征值

在这项工作中,我们关注 NTK 最大特征值的 EOS 动力学,而不是像 (Cohen 等人, 2022a)中的 Hessian 矩阵。我们注意到,定理 3.1 的一个版本对于最大 Hessian 特征值也成立。一般来说,Hessian 矩阵可以写成

$$\frac{\partial^2 \mathcal{L}}{\partial \theta \partial \theta} = \nabla \mathcal{L} \cdot \frac{\partial^2 z}{\partial \theta \partial \theta} + \mathbf{J}^T \frac{\partial^2 \mathcal{L}}{\partial z \partial z} \mathbf{J} \quad (30)$$

特别是对于平方损失,我们有

$$\frac{\partial^2 \mathcal{L}}{\partial \theta \partial \theta} = \nabla \mathcal{L} \cdot \frac{\partial^2 z}{\partial \theta \partial \theta} + \mathbf{J}^T \mathbf{J} \quad (31)$$

当损失梯度趋于0时,Hessian特征值接近JJT的特征值 其非零特征值与经验NTK JJT的非零特征值相同。由于该定理涉及 $z \sim$ 收敛时的行为,因此最大NTK和最大 Hessian 特征值在极限上相等,并且在两种情况下应用相同的 EOS 行为。

对于高维模型 (CIFAR10 上的二次回归模型和全连接网络),我们的实验表明最大 NTK 特征值显示出稳定性行为的边缘。CIFAR 模型与 (Cohen 等人,2022a)中的模型相同,后者用于说明最大 Hessian 特征值的稳定性边缘。

因此,我们在论文中重点关注 NTK 版本的 EOS,因为我们发现它更适合理论分析和解释。

几乎可以肯定,在某些情况下,EOS 行为会显示在 Hessian 特征值中,但不会显示在 NTK 特征值中,特别是在输出中损失高度非各向同性的情况下 (即,远离单位矩阵的倍数)。正如之前的工作中所指出的,在这些情况下,即使是基于 Hessian 的 EOS 也更难以分析 (Cohen 等人,2022a)。我们将对 EOS 的理解和更复杂的损失函数留给未来的工作。

A2.一隐层线性网络

考虑一个具有标量输出的单隐藏层网络:

$$f(x) = vUx \quad (32)$$

其中 x 是长度为 N 的输入向量, U 是 $K \times N$ 维矩阵, v 是 K 维向量。我们注意到

$$\frac{\partial^2 f(x)}{\partial v_i \partial v_j} = \frac{\partial^2 f(x)}{\partial U_{ij} \partial U_{kl}} \frac{\partial^2 f(x)}{\partial v_i \partial v_j} = \delta_{ijkl} \quad (33)$$

其中 δ_{ij} 是克罗内克三角洲。对于固定的训练集,这个二阶导数是常数;因此,单隐层线性网络是第 4 节中研究的类型的二次回归模型。

在单个数据点 x 的特定情况下,我们可以计算 Q 矩阵的特征向量。令 (w,W) 为 Q 的特征向量,分别表示 v 和 U 分量。特征向量方程是

$$\omega w_i = x_m \delta_{ij} W_{jm} \quad (34)$$

$$\omega W_{jm} = x_m \delta_{ij} w_i \quad (35)$$

简化一下,我们有:

$$\omega W = Wx \quad (36)$$

$$\omega W = wx \quad (37)$$

我们有两种情况。第一个是 $\omega = 0$ 。在这种情况下, $w = 0$, W 是一个矩阵, x 位于其零空间中。后一个条件为我们提供了对 $M \times N$ 方程的 M 个约束- 对于我们的 $M(N+1)$ 总本征模式总共有 $M(N-1)$ 个。

如果 $\omega = 0$,则结合方程我们有条件:

$$\omega^2 w = (x \cdot x) w \quad (38)$$

$$\omega^2 W = Wxx \quad (39)$$

由此得出 $\omega = \pm \sqrt{x \cdot x}$ 。从方程 37 中我们知道 W 是低阶的。因此,我们可以猜测形式的解决方案

$$W_{\pm, i} = \pm e_i x \quad (40)$$

其中 e_i 是 M 坐标向量。这表明我们有

$$w_{\pm, i} = (\sqrt{x \cdot x}) e_i \quad (41)$$

这给了我们最终的 $2M$ 本征模态。

我们也可以分析 $J(\omega_i)$ 的初始值。雅可比行列式的组成部分可以写为:

$$(Jv)_i = \frac{\partial f(x)}{\partial v_i} = U_{im} x_m \quad (42)$$

$$(JU)_{jm} = \frac{\partial f(x)}{\partial x_m} U_{jm} \quad (43)$$

从这个形式,我们可以推断出 J 与 0 模态正交。我们还可以计算守恒量。令 J 为正本征模态的总权重, J 表明 J^2 是负本征模态的总权重。直接计算

$$\omega^{-1} (J^2 - J + 2) = 2f(x) \quad (44)$$

这意味着 $E = 0$ 。

因此,一个数据点上的单隐层线性模型相当于 $E = 0$ 、特征值 $\pm \sqrt{x \cdot x}$ 的四次损失模型。

A.3. 连接到 (Bordelon 和 Pehlevan, 2022)

由于单隐层线性模型具有恒定的 Q , 因此 (Bordelon & Pehlevan, 2022) 的 F.1 节中的模型属于二次回归类。在 F.1.1 节公式 67 的情况下, 我们可以明确映射到 $D = 1$ 模型。如果我们进行识别, 则动力学相当于具有单个特征值 ω_0 的所述模型

$$z, Hy = J \quad \frac{\partial}{\partial t}, y_0 = \sqrt{2\omega}, y = -E/2 \Delta = \quad (45)$$

A.4. 连接至 NTH

神经正切层次 (NTH) 方程扩展了 NTK 动力学, 通过构建控制非线性学习动力学的高阶张量的无限序列来解释正切核的变化 (Huang & Yau, 2020)。三阶 NTH 方程的截断与二次回归模型相关, 但并不相同, 正如我们将在此处展示的那样。

三阶 NTH 方程描述了切核 JJ 的变化。考虑 $D \times D \times D$ 维核 K_3 , 其元素由下式给出

$$(K_3)_{\alpha\beta\gamma} = \frac{\partial^2 \alpha}{\partial \theta_i \partial \theta_j} \frac{\partial^2 \beta}{\partial \theta_i \partial \theta_j} + \frac{\partial^2 \gamma}{\partial \theta_i \partial \theta_j} \quad J_{i\gamma} J_{j\alpha} \quad (46)$$

其中重复索引相加。在 NTH 中, 对于平方损失, NTK JJ 的变化由下式给出

$$\frac{d}{dt} K_{\alpha\beta} = -\eta (K_3)_{\alpha\beta\gamma\gamma} \quad (47)$$

对于固定 $Q = \frac{\partial}{\partial z \partial \theta \partial \theta}$ 该方程与二次回归模型中 NTK 的 GF 方程相同。我们 $\partial^2 \alpha \partial \theta \partial \theta$

请注意, 二次回归模型下 K_3 不是常数。相反, 对于固定的 K_3 , 也不是恒定的。

因此, 这两种方法可以用来构造不同的动力学低阶展开。

B.2 参数模型

B.1 梯度流下渐进锐化

给定重新调整的变量

$$z \sim = \eta z, T(k) = \eta J Q k J. \quad (48)$$

我们得到动力学方程

$$\frac{dz \sim}{dt} = -z T \sim(0), \quad \frac{dT(k)}{dt} = -2 \sim z T(k+1). \quad (49)$$

我们可以证明,有一些初始化在后期显示出渐进锐化 - 即 $T(0)$ 的增加。

考虑一个初始化,其中 $T(k)$ 对于偶数 k 为非负,对于奇数 k 为非正。在这种情况下,给定动力学方程, $T(k)$ 不会改变符号。如果 $z \sim > 0$, 则 $T(0)$ 始终不减。

这种初始化的一个例子是 Q 的特征值是 $\pm \lambda$ 的情况,并且 J 在负特征模式中比在正特征模式中使用更大的权重进行初始化。还有许多其他具有此属性的初始化系列。

B.2 函数空间动力学方程的推导

考虑梯度下降方程:

$$\theta_{t+1} - \theta_t = -\eta \nabla \theta L = -\frac{\eta}{2} \theta Q \theta - E Q \theta. \quad (50)$$

我们可以如下推导 $z \sim - T$ 动力学。我们可以将更新方程重写为

$$\theta_{t+1} - \theta_t = -\eta \nabla \theta L = -\eta z J. \quad (51)$$

z 的变化可以写为

$$\Delta z = \theta Q \Delta \theta + \frac{1}{2} \Delta \theta Q \Delta \theta \quad (52)$$

代入,我们有:

$$\Delta z = -\eta z J J + \frac{2\eta}{2} z J Q J \quad (53)$$

我们还有

$$\Delta J = -\eta z J Q \quad (54)$$

我们还可以写:

$$\Delta(J Q k J) = -2\eta z J Q k + 1 J + \eta J^2 Q k + 2 J \quad (55)$$

转换为 $z \sim - T$ 坐标,梯度下降方程变为:

$$1 z \sim_{t+1} - z \sim_t = -z \sim_t T_t(0) + (z \sim_t^2) T_t(1) \quad (56)$$

$$T_{t+1}(k) - T_t(k) = -z \sim_t (2 T_t(k+1) - z \sim_t T_t(k+2)). \quad (57)$$

B.3. $z \sim - T(0)$ 方程的推导我们可以使用守恒量 $E \sim$

来仅用 $z \sim$ 和 $T(0)$ 来写出动力学。不失一般性,设特征值为 1 和 λ , 且 $-1 \leq \lambda \leq 1$ 。(我们可以通过重新缩放 $z \sim$ 来实现这一点。)回想一下动力学方程

$$1 z \sim_{t+1} - z \sim_t = -z \sim_t T_t(0) + (z \sim_t^2) T_t(1) \quad (58)$$

$$T_{t+1}(0) - T_t(0) = -z \sim_t (2 T_t(1) - z \sim_t T_t(2)) \quad (59)$$

我们将根据 $z \sim$ 和 $T(0)$ 找到 $T(1)$ 和 $T(2)$ 的替代。回想一下我们有

$$T(-1) = E \sim + 2 \sim z \quad (60)$$

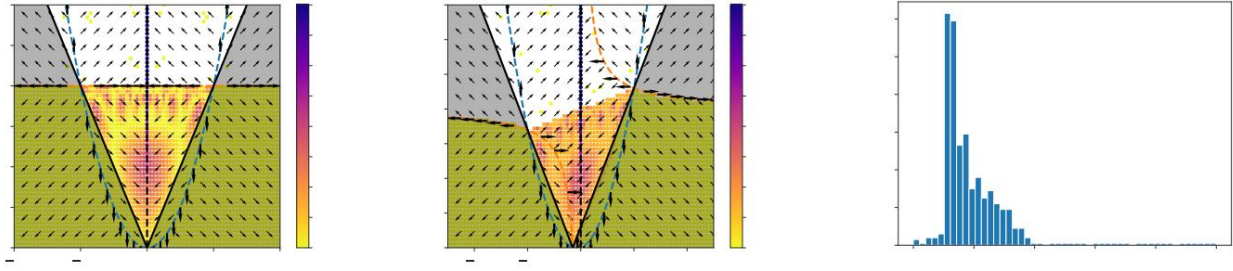


图 7.对称模型的相图。箭头表示 z 和 T 变化的迹象,灰色区域表示不允许的坐标。动力学从均匀间隔的初始化网格运行,并记录曲率 $T(0)$ 的最终值。

代表 $z_{t+1} - z_t = 0$ (蓝色)和 $T_{t+1}(0) - T_t(0) = 0$ (橙色)的零斜线取决于 E 。轨迹显示逐渐锐化,但没有稳定边缘效应 (右)。

其中 E 在整个动态过程中是守恒的 (实际上是景观的一个属性)。我们将使用这个定义来求解 $T(1)$ 和 $T(2)$ 。

由于 $P = 2$,我们可以将系数 a 和 b 写为 $T(-1) = bT(0) + aT(1)$,这对 J 的所有组合都有效。如果 $J(-\lambda) = 0$,则有 $b = 1 - a$ 。如果 $J(1) = 0$,则有 $1 = \lambda(1 - a) + \lambda 2a$ 。求解,我们有:

$$T(-1) = (1 - a)T(0) + aT(1) \text{ 对于 } a = -\frac{1}{\lambda} \quad (61)$$

对 λ 的限制转化为 $a/\epsilon \in (-1, 1)$ 。根据守恒量 $E = T(-1) - 2z$,我们有:

$$T(-1) = E + 2z \quad (62)$$

为了转换动力学,我们需要根据 $T(0)$ 和 z 求解 $T(1)$ 和 $T(2)$ 。我们有:

$$T(1) = \frac{1}{A} (T(-1) + (a - 1)T(0)) = \frac{1}{A} (E + 2z + (a - 1)T(0)) \quad (63)$$

我们还有

$$T(2) = T(0) + \frac{1 - a}{\lambda^2} (T(0) - E - 2z) \quad (64)$$

这给了我们

$$z_{t+1} - z_t = -z_t T_t(0) + \frac{1}{2a} \left((a - 1)T_t(0) + 2z_t + E \right) \quad (65)$$

$$T_{t+1}(0) - T_t(0) = -\frac{2}{A} z_t (2z_t + E + (a - 1)T_t(0)) + z_t^2 T_t(0) + \frac{1 - a}{\lambda^2} (T_t(0) - E - 2z_t) \quad (66)$$

如果 $\lambda = -1$ (即 $a = 1$) 我们从正文中恢复方程。

J_2 的非负性给了我们 z 和 T 的值的约束。对于 $a > 1$ (小的负第二特征值),

限制条件是:

$$T > 2z + E, T > -(2z + E)/a \text{ 这 } \quad (67)$$

是一个面朝上的圆锥体,顶点位于 $z = -E/2$ (图 8,左)。对于 $a < -1$,约束是

$$-(2z + E)/a < T < 2z + E \quad (68)$$

这是一个面向侧面的圆锥体,顶点位于 $z = -E/2$ (图 8,右)。我们看到,在这种情况下,有一组有限的 T 值可以收敛。事实上,对于 $E = 0$,除了 $T(0) = 0$ 之外,没有收敛。这就是为什么我们关注一个正特征值和一个负特征值的情况。

我们还可以求解零斜线 - $z \sim t+1 - z \sim t = 0$ (图 8 中的蓝色) 或 $T_{t+1}(0) - T_t(0) = 0$ (图 8 中的橙色) 的曲线。 $z \sim$ 的零斜线 ($\sim z, f_z(\sim z)$) 由下式给出

$$f_z(\sim z) = \frac{z \sim (2 \sim z + E \sim)}{2a - (a - 1) \sim z} \quad (69)$$

$T(0)$ 的零斜线 ($\sim z, f_T(\sim z)$) 由下式给出

$$f_T(\sim z) = - \frac{(a-1) \sim z - 2a}{(a^2 - a + 1) \sim z - 2a(a-1)(2 \sim z + E \sim)} \quad (70)$$

线 $z \sim = 0$ 也是零斜线。

对于对称模型 $= 1$, 零斜线的结构决定是否存在渐进锐化。

对于 $E \sim = 0$, 不进行锐化; 相图 (图 7, 左) 证实了这一点, 因为 $T_t(0)$ 中的零斜线将空间分为两半, 一半会聚, 另一半不会聚。然而, 当 $E \sim = 0$ 时, 零斜线分裂, 并且有一个小区域可以发生渐进锐化 (图 7, 中)。然而, 在这种情况下仍然没有稳定边缘行为 - 没有轨迹聚集在 $\lambda_{\max} = 2/\eta$ 附近的区域 (图 7, 右)。

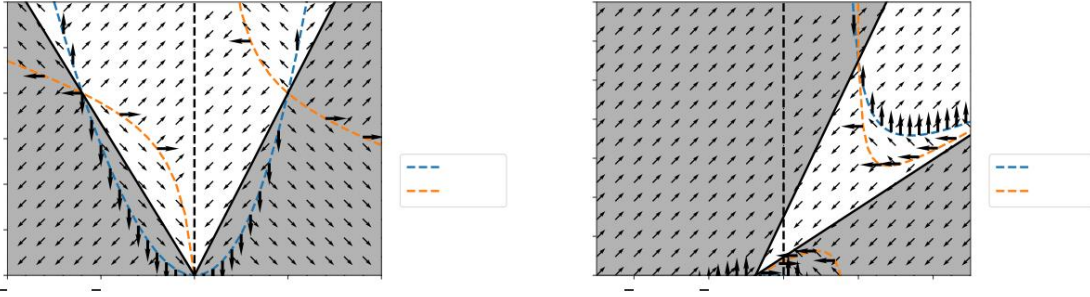


图 8. $D = 1$ 、 $P = 2$ 模型的相平面。灰色区域对应于 $J \sim (\omega_i)$ 上的正性约束所禁止的参数。对于 $\lambda > 0$, 允许的区域较小并且仅在小范围内与 $z \sim = 0$ 相交。零斜线可以通过分析求解。

B.4. 两步动力学

两步差分方程可以通过迭代方程 14 和 15 得到。我们有

$$z \sim t+2 - z \sim t = p_0(\sim z_t) + p_1(\sim z_t, T_t(0)) + p_2(\sim z_t, T_t(0)^2) + p_3(\sim z_t, T_t(0)^3) \quad (71)$$

$$T_t(0)_{t+2} - T_t(0) = q_0(\sim z_t) + q_1(\sim z_t, T_t(0)) + q_2(\sim z_t, T_t(0)^2) + q_3(\sim z_t, T_t(0)^3) \quad (72)$$

这里 p_i 和 q_i 是 $z \sim$ 中的多项式, $z \sim$ 中最大为 9 阶, 中为 6 阶。它们可以显式地计算, 但我们现在选择省略确切的形式。

小型动力学的数值模拟揭示了稳定性效应的边缘 (图 9)。我们看到随机初始化的 T 最终值的分布在 $T(0) = 2$ 附近有一个峰值 (右)。通过绘制两步动力学, 我们可以看到进入 $T(0) = 2$ 的两步零斜线几乎重合 (左)。通过研究这些零斜线, 我们将能够了解稳定效应的边缘。

个零斜线方, 我们可以使用 Forfix Cardano 公式求解 $z \sim$ 两步零斜线 ($z \sim t+2 - z \sim t = 0$) 和 T 零斜线 ($T_{t+2}(0) - T_t(0) = 0$) T 作为 $z \sim$ 的函数。特别是, 每程都有一个满足 $z \sim = 0, T(0) = 2$ 的解, 独立于 这是我们将重点关注的解族。

令 ($\sim z, f_z, \sim(\sim z)$) 为 $z \sim$ 的零斜线, 并令 ($\sim z, f_T, \sim(\sim z)$) 为 $T(0)$ 的零斜线。我们将展示零斜线的 T 值, 作为 $z \sim$ 和的函数, 在 $z \sim = 0, = 0$ 附近可微。

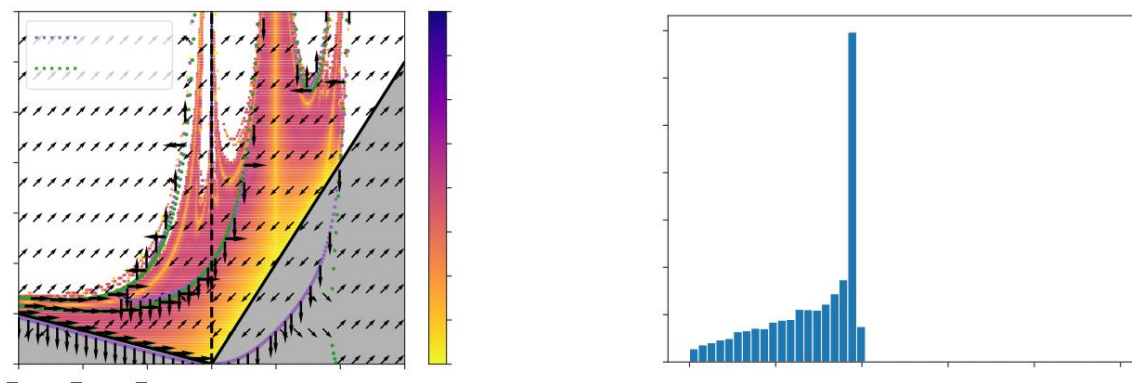


图 9. $\alpha = 0.17$ 、 $E = 0$ 时的相图。箭头表示 $z \sim$ 和 $T(0)$ (左) 的两步动态变化的迹象, 灰色区域代表不允许的坐标。代表 $z \sim t + 2$ - $z \sim t = 0$ (紫色) 和 $Tt + 2(0) - Tt(0) = 0$ (绿色) 的零斜线进入 $T(0) = 2$ 几乎重叠。在 $z \sim - T(0)$ 空间上均匀进行初始化的最终 $T(0)$ 显示 $T(0) = 2$ 附近的峰值 (右)。

零斜线由隐式方程定义

$$0 = 6z^3 - 2Tz^2 - 3Tz^2(-1) - Tz^3 + (2)(2+1)T^2 + \frac{7}{2}Tz^2(-1) + \frac{1}{2}T^2z^2 - 9z^2 - 10 + 9 - \frac{1}{2}T^3z^2(-1) - \frac{1}{2}T^3z^2 - 3z^2 - 4 + 3 + O(z^4) \quad (73)$$

$$20 = -8z^3 - 12z^2(-1) + 4Tz^2(-1) + 2Tz^3 + 2z^3 + 3 + 4Tz^2(-1)^2 + 4 + 1 - 2T^2z^2(-1) - T^2z^2 - 8 + 7 - T^2z^2(-1)9^2 + 9 + T^2z^2 - 1 + T^3z^2(-1)3^2 - 3 + O(z^4) \quad (74)$$

我们暂时省略高阶项, 以期在 $z \sim 0$ 处微分以使用隐函数定理。

除以 $z \sim$, 我们得到方程

$$20 = 6z^2 - 2T - 3Tz^2(-1) - Tz^2 + (2)(2+1)T^2 + \frac{7}{2}Tz^2(-1) + \frac{1}{2}T^2z^2 - 9z^2 - 10 + 9 - \frac{1}{2}T^3z^2(-1) - \frac{1}{2}T^3z^2 - 3z^2 - 4 + 3 + O(z^3) \quad (75)$$

$$20 = -8z^2 - 12z^2(-1) + 4T(-1) + 2Tz^3 + 2z^3 + 3 + 4Tz^2(-1)^2 + 4 + 1 - 2T^2(-1) - T^2z^2 - 8 + 7 - T^2z^2(-1)9^2 + 9 + T^2z^2 - 1 + T^3z^2(-1)3^2 - 3 + O(z^3) \quad (76)$$

我们立即看到 $z \sim 0$, $T = 2$ 分别求解了所有方程 75 和 76 的两个方程。我们有 $w(z, T)$ 和 $v(z, T)$ 为右侧

$$\frac{\partial w}{\partial T}(0,0,2) = 2, \quad \frac{\partial v}{\partial T}(0,0,2) = 4 \quad (77)$$

在这两种情况下, 导数都是可逆的。因此, $f_z, \sim(\sim z)$ 和 $f_T, (\sim z)$ 在 $z \sim 0$ 的某个邻域内连续可微。事实上, 由于 w 和 v 在所有三个参数中都是解析的, 所以 $f_z, \sim(\sim z)$ 和 $f_T, (\sim z)$ 也是解析的。

二阶回归和稳定边缘

我们可以利用解析性来求解零斜线的低阶结构。计算导数值的一种方法是将零斜线定义为形式幂级数：

$$f_z(\tilde{z}) = 2 + \sum_{j=1}^{\infty} a_{j,k} \tilde{z}^k \tag{78}$$

$$f_T(\tilde{z}) = 2 + \sum_{j=1}^{\infty} b_{j,k} \tilde{z}^k \tag{79}$$

然后,我们可以使用方程 75 和 76 求解级数的前几项。根据此过程,我们得到：

$$f_z(\tilde{z}) = 2 + 2(1 - \tilde{z}) + \frac{1}{2} \tilde{z}^2 + O(\tilde{z}^3) \tag{80}$$

$$f_T(\tilde{z}) = 2 - \frac{2 - 3 + 2}{1 - \tilde{z}} \tilde{z} + \frac{1}{2} \tilde{z}^2 + O(\tilde{z}^3) \tag{81}$$

两者之间的差异 $f_{\Delta}(\tilde{z})$ 为：

$$f_{\Delta}(\tilde{z}) = f_z(\tilde{z}) - f_T(\tilde{z}) = -1 - \frac{3}{2} \tilde{z} + O(\tilde{z}^3) \tag{82}$$

随着减小,对于低阶项,零斜线之间的距离也会减小。

我们可以证明差异如下
在 $\tilde{z} = 0$ 时的一步动力学方程为

$$\tilde{z}_{t+1} - \tilde{z}_t = -\tilde{z}_t T_t(0) + \tilde{z}_t^2 \frac{1}{t} \tag{83}$$

$$T_{t+1}(0) - T_t(0) = -2\tilde{z}_t T_t(0) + \tilde{z}_t^2 \frac{1}{t} \tag{84}$$

因此, $\Delta \tilde{z} = 2\Delta T$ 。这意味着一步零斜线和两步零斜线是相同的。由于 $f_z(\tilde{z}) = f_T(\tilde{z})$,且两者可微分,我们有：

$$f_z(\tilde{z}) - f_T(\tilde{z}) = f_{\Delta}(\tilde{z}) \tag{85}$$

对于某些函数 $f_{\Delta}(\tilde{z})$,它在 $(0, 0)$ 周围的邻域内和 \tilde{z} 中解析。

B.5. y 的两步动力学

在坐标 (\tilde{z}, y) 中定义动力学方程很有用,其中 y 是 $T(0)$ 和 \tilde{z} -零斜线之间的差：

$$y \neq T(0) - f_z(\tilde{z}) \tag{86}$$

到 \tilde{z} 的最低阶,我们有

$$y = T(0) - 2 - 2(1 - \tilde{z}) - \frac{1}{2} \tilde{z}^2 + O(\tilde{z}^3) \tag{87}$$

我们注意到 $y = 0$,在 $\tilde{z} = 0$ 处对应于 $T(0) = 2$ 。 y 接近但略小于0相当于稳定边缘行为。对于正 \tilde{z} , $y = 0$ 意味着 $T(0) > 2$ 。

我们可以写出 \tilde{z} 和 y 的动态。 \tilde{z} 的动力学由下式给出：

$$\tilde{z}_{t+2} - \tilde{z}_t = p_0(\tilde{z}_t, y) + p_1(\tilde{z}_t, y)(y + f_z(\tilde{z}_t)) + p_2(\tilde{z}_t, y)(y + f_z(\tilde{z}_t))^2 + p_3(\tilde{z}_t, y)(y + f_z(\tilde{z}_t))^3 \tag{88}$$

我们知道这个方程的右侧在 \tilde{z} -和（简单地） y 中也是解析的。通过计算 f 的多个连续导数,我们可以写出：

$$\tilde{z}_{t+2} - \tilde{z}_t = 2y\tilde{z}_t + y^2 \frac{1}{t} f_{1,1}(\tilde{z}_t, y) + y\tilde{z}_t^2 \frac{1}{t} f_{2,1}(\tilde{z}_t, y) \tag{89}$$

二阶回归和稳定边缘

这里， f_1 和 f_2 在 0 附近的某个邻域内的 z 、和 y 中解析。

这意味着我们有界限

$$|f_1(\tilde{z}, y)| < F_1, |f_2(\tilde{z}, y)| < F_2 \tag{90}$$

对于一些非负常数 F_1 和 F_2 , $(\tilde{z}, y) \in [-z^d, z^d] \times [0, d] \times [-y^d, y^d]$ 。注意这个界限是独立的。
的。

现在我们考虑 y 的动态。我们有：

$$y_{t+2} - y_t = T_{t+2}(0) - T_t(0) - f_z, \sim (\sim z_{t+2}) + f_z, \sim (\sim z_t) \tag{91}$$

由于 $\lim_{z \rightarrow 0, y \rightarrow 0} z \sim t+2 = 0$,所以 $f_z, \sim (\sim z_{t+2})$ 在 $(0, 0, 0)$ 的某个邻域内是解析的。因此 $y_{t+2} - y_t$ 也是解析的。代入,我们有

$$\begin{aligned} y_{t+2} - y_t = & q_0(\sim z_t,) + q_1(\sim z_t,) [y + f_z, \sim (\sim z)] + q_2(\sim z_t,) [y + f_z, \sim (\sim z)]^2 + q_3(\sim z_t,) [y + f_z, \sim (\sim z)]^3 \\ & + t z \sim t f_1(\sim z_t, y_t) - f_z, \sim (\sim z_t + 2 y t z \sim t + y^2 y t z \sim t f_2(\sim z_t)) + f_z, \sim (\sim z_t^2) \end{aligned} \tag{92}$$

如果我们写 $f_z, \sim (\sim z) = f_T, (\sim z) + f_\Delta, (\sim z)$,那么我们可以写：

$$\begin{aligned} y_{t+2} - y_t = & q_0(\sim z_t,) + q_1(\sim z_t,) [f_T, (\sim z)] + q_2(\sim z_t,) [f_T, (\sim z)]^2 + q_3(\sim z_t,) [f_T, (\sim z)]^3 \\ & + 2 q_2(\sim z_t,) [f_T, (\sim z)] (y + f_\Delta, (\sim z)) + 3 q_3(\sim z_t,) [(f_T, (\sim z)) (y + f_\Delta, (\sim z))]^2 + (f_T, (\sim z))^2 (y + f_\Delta, (\sim z)) \\ & + q_0(\sim z_t,) + q_1(\sim z_t,) [y + f_\Delta, (\sim z)] + q_2(\sim z_t,) [y + f_\Delta, (\sim z)]^2 \\ & + q_3(\sim z_t,) [y + f_\Delta, (\sim z)]^3 - f_z, \sim (\sim z_t + 2 y t z \sim t + y^2 t z \sim t f_1(\sim z_t, y_t) + y t z \sim t f_2(\sim z_t^2)) + f_z, \sim (\sim z_t) \end{aligned} \tag{93}$$

根据零斜线的定义,前四项消失。再次使用零斜线以及 f_1 和 f_2 的可微性,我们可以根据展开式重写动力学：

$$y_{t+2} - y_t = -2(4 - 3 + 4^2 y t z \sim t^2 - 4 z \sim t^2 + y^2 t^2 g_1(\sim z_t, y_t) + z \sim t g_2^3(\sim z_t) \tag{94}$$

这里 g_1 和 g_2 在 $z \sim y$ 和 中解析接近于零。我们有界限

$$|g_1(\tilde{z}, y)| < G_1, |g_2(\tilde{z}, y)| < G_2 \tag{95}$$

对于一些非负常数 G_1 和 G_2 , $(\tilde{z}, y) \in [-z^d, z^d] \times [0, d] \times [-y^d, y^d]$ 。这个界限也独立于。
。

我们可以在以下引理中总结这些界限:引理 B.1。定义 $y = T - f_z(\sim z)$ 。 $z \sim$

和 y 的两步动态由下式给出

$$z \sim t+2 - z \sim t = 2 y t z \sim t + y^2 t z \sim t f_1(\sim z_t, y_t) + y t z \sim t f_2(\sim z_t^2) \tag{96}$$

$$y_{t+2} - y_t = -2(4 - 3 + 4^2 y t z \sim t^2 - 4 z \sim t^2 + y^2 t^2 g_1(\sim z_t, y_t) + z \sim t g_2^3(\sim z_t, y_t) \tag{97}$$

其中 f_1, f_2, g_1, g_2 在 $z \sim, y$ 和 中都是解析的。此外,还存在正的 $z \sim c, y$ 和。这样

$$|f_1(\tilde{z}, y)| < F_1, |f_2(\tilde{z}, y)| < F_2, |g_1(\tilde{z}, y)| < G_1, |g_2(\tilde{z}, y)| < G_2 \tag{98}$$

对于所有 $(\tilde{z}, y) \in [-z^d, z^d] \times [0, d] \times [-y^d, y^d]$,其中 F_1, F_2, G_1 和 G_2 都是非负常数。

我们可以使用这个引理来分析小固定项的动态,这将允许进行专注于低阶项影响,对于 $z \sim, y$ 的小初始化。上级的控制的分析。

二阶回归和稳定边缘

B.6.定理3.1的证明

使用引理 B.1, z 和 y 的动力学可以写为:

$$\dot{z} - \dot{z} = 2ytz - t + y \frac{2}{t} z - t f_1(\sim z, y) + ytz - t f_2(\sim z) \quad (99)$$

$$\dot{y} + \dot{y} = -2(4 - 3 + 4 \frac{2}{y}) ytz - t \frac{2}{2} - 4z \sim t + y \frac{2}{t} z - t g_1(\sim z, y) + z \frac{3}{t} g_2(\sim z, y) \quad (100)$$

让 $\epsilon < d$ 。然后我们可以使用引理 B.1 的界限来控制高阶项对动力学的贡献: 引理 B.2。给定常数 $A > 0$ 和 $B > 0$, 存在 $z \sim c$ 和 y_c , 对于 $z \sim \epsilon [0, 2 \sim zc]$, $y \in [-y_c, y_c]$, 我们有边界:

$$|y \frac{2}{t} z f_1(\sim z, y) + yz \frac{2}{t} f_2(\sim z)| \leq A |2yz \sim| \quad (101)$$

$$|y \frac{2}{t} z \frac{2}{t} g_1(\sim z, y)| \leq \frac{B}{2} | -2(4 - 3 + 4 \frac{2}{y}) yz \sim^2 | \quad (102)$$

$$|z \frac{3}{t} g_2(\sim z, y)| \leq \frac{B}{4} |4z \sim^2| \quad (103)$$

证明。我们从以下分解开始:

$$|y \frac{2}{t} z f_1(\sim z, y) + yz \frac{2}{t} f_2(\sim z)| \leq |y \frac{2}{t} z f_1(\sim z, y)| + |yz \frac{2}{t} f_2(\sim z)| \quad (104)$$

根据引理 B.1, 存在一个区域 $[-z \sim d, z \sim d] \times [0, d] \times [-y_d, y_d]$, 其中 f_1 、 f_2 、 g_1 和 g_2 的大小为分别以 F_1 、 F_2 、 G_1 和 G_2 为界。

$$|y \frac{2}{t} z f_1(\sim z, y) + yz \frac{2}{t} f_2(\sim z)| \leq F_1 y \frac{2}{t} z \sim + F_2 yz \sim^2 \quad (105)$$

$$|y \frac{2}{t} z \frac{2}{t} g_1(\sim z, y)| \leq G_1 y \frac{2}{t} z \sim^2 \quad (106)$$

$$|z \frac{3}{t} g_2(\sim z, y)| \leq G_2 z \sim^3 \quad (107)$$

将 $z \sim c$ 和 y_c 定义为

$$y_c = \min(A/F_1, B/2G_1, y_d), \quad z \sim c = \min(A/2F_2, B/2G_2, y_c) \quad (108)$$

立即出现所需的边界。□

我们考虑初始化 $(\sim z_0, y_0)$, 使得 $z \sim 0 \leq z \sim c$ 且 $y_0 \leq y_c$, 且 $y_0 \leq z \sim$ 借助引理 B.2, 我们可以分析动力学。有两个阶段; 在第一阶段, $z \sim$ 增加, y 减少。当 y 第一次变为负数时, 第一阶段结束 - 达到 $O()$ 的值。在第二阶段, $z \sim$ 减小, y 保持负值且 $O()$ 。

B.6.1. 第一阶段

令 t_{sm} 为时间, 使得对于 $t \leq t_{sm}$, $z \sim t \leq 2 \sim z_0$ 。 (我们稍后将证明在整个动态范围内 $z \sim t \leq 2 \sim z_0$ 。) 对于 $t \leq t_{sm}$, 使用引理 B.2, $z \sim$ 的变化可以从下面限制为

$$\dot{z} - \dot{z} + 2 \geq 2ytz - t(1 - A) \quad (109)$$

因此在初始化时, $z \sim$ 不断增加。它持续增加, 直到 y 变为负值, 或 $z \sim t \geq 2 \sim z_0$ 。我们要证明 y 在 $z \sim t \geq 2 \sim z_0$ 之前变为负值。

对于任何 $t \leq t_{sm}$, 引理 B.2 给出 $\dot{y} + \dot{y}$ 的以下上限:

$$\dot{y} + \dot{y} \leq -(8 - B)ytz \frac{2}{t} - (4 - B)z \sim \frac{2}{t} \quad (110)$$

二阶回归和稳定边缘

设 t_1 为 y_t 第一次变为负数的时间。由于 z_t 在 $t \leq t_1$ 时增加, 因此我们有

$$y_{t+2} - y_t \leq -(8 - B)y_t z_t \quad \frac{2}{0} - (4 - B)z_t \sim \frac{2}{0} \quad (111)$$

这给了我们 y_t 的以下界限:

$$y_t \leq y_0 e^{-(8-B)z_t \sim \frac{2}{0}} \quad (112)$$

对于 $t \leq t_1$ 和 $t \leq t_{sm}$ 有效。

现在我们将证明 $t_1 < t_{sm}$ 。假设 $t_{sm} \leq t_1$ 。然后在 $t_{sm} + 2$ 处, $z_{t_{sm}+2}$ 第一次 $> 2 \sim z_0$ 。对公式 109 中的界限求和, 我们得到:

$$z_{t_{sm}+2} - z_0 \leq \sum_{t=0}^{t_{sm}} 2y_t z_t (1+A) \leq 4 \sum_{t=0}^{t_{sm}} y_t \quad (113)$$

其中第二个界限来自 t_{sm} 的定义。使用 y_t 上的界限, 我们有:

$$z_{t_{sm}+2} - z_0 \leq 4 \sum_{s=0}^{t_{sm}} y_0 e^{-(8-B)z_s \sim \frac{2}{0}} \sim \frac{(1+A)y_0 \leq z}{0} \quad (114)$$

由于 $y_0 \leq z \sim \frac{2}{0}$, 所以 $z_{t_{sm}+2} \leq 2 \sim z_0$ 。然而, 假设 $z_{t_{sm}+2} > 2 \sim z_0$ 。我们遇到了一个矛盾; t_{sm} 不小于或等于 t_1 。

有三种可能性: 第一种是 t_1 是明确定义的, 并且 $t_1 < t_{sm}$ 。另一种可能性是 t_1 没有明确定义 - 即 y_t 永远不会变为负值。在这种情况下, 我们得出的界限对于所有 t 都有效。因此, 使用等式 112, 存在某个时间 t , 其中 $y_t < (4-B)z \sim 0$ 。然后, 使用方程 111, 我们有 $y_{t+2} < 0$ 。因此, 我们得出结论 t_1 是有限的并且小于 t_{sm} 。

由于明确定义的值 $t_1 < t_{sm}$, 当 y 首先变为负值时, $z_{t_1} \leq 2 \sim z_0$ 。这意味着我们可以在下一阶段开始时继续应用引理 B.2 的边界。在 $t = t_1 - 2$ 处, 应用引理 B.2 和 $z_{t_1} \leq 2 \sim z_0$, 我们有

$$y_{t_1} - y_{t_1-2} \geq -4(8+B)y_{t_1-2} z_{t_1} \sim \frac{2}{0} - 4(4+B)z \sim \frac{2}{0} \quad (115)$$

这给我们 $y_{t_1} \geq -4(4+B)z \sim \frac{2}{0}$ 。第一阶段就此结束。总而言之, 我们有

$$-4(4+B)z \sim \frac{2}{0} < y_{t_1} \leq 0, \quad z_{t_1} \leq 2 \sim z_0 \quad (116)$$

B.6.2. 第二阶段

现在考虑动力学的第二阶段。我们将证明 y 仍然为负并且 $O()$, 并且 z 减小到 0。

当 y 为负时, z 减小。当 $y \geq -y_0$ 时, 从引理 B.2 我们有

$$z_{t+2} - z_t \leq (1-A)2y_t z_t \quad (117)$$

因此, 只要 $-y_0 \leq y < 0$, z_t 就会减小。如果对于所有后续 t 都是如此, 则 z_0 将收敛到 0。

现在我们将证明 y 仍然为负且 $O()$, 从而得出证明。设 y 的 y 动力学方程为

$$y' = - \frac{2}{(3/2)+2}$$

$$y_{t+2} - y_t = -2(4 - 3 + 4 \frac{2}{t}) \frac{2}{t} (y_t - y^*) + y \frac{2}{t} g_1(\sim z_t, y_t) + \sim z \frac{3}{t} g_2(\sim z_t, y_t) \quad (118)$$

将引理 B.2 应用于高阶项, 我们有:

$$y_{t+2} - y_t \leq -2(4 - 3 + 4 \frac{2}{t}) \frac{2}{t} (y_t - y^*) + B(|y_t| +) \sim z \frac{2}{t} \quad (119)$$

$$y_{t+2} - y_t \geq -2(4 - 3 + 4 \frac{2}{t}) \frac{2}{t} (y_t - y^*) - B(|y_t| +) \sim z \frac{2}{t} \quad (120)$$

只要 $|y_t|$ 这些不等式就有效。 $< y_c$ 。

我们可以使用这些界限构造两个辅助序列 u_t 和 v_t ,使得 $u_t \leq y_t \leq v_t$.设 $u_{t-1} = v_{t-1} = y_{t-1}$ 。
 v_t 的动态由下式给出

$$u_{t+2} - u_t = -8z_t^2 (u_t - u^*) \tag{121}$$

$$v_{t+2} - v_t = -8z_t^2 (v_t - v^*) \tag{122}$$

你在哪里 $y^* = y^* - 4Bz_t^2$ 和 $v^* = y^* + 4Bz_t^2$ 。

我们首先分析 v_t 的动态。考虑由方程定义的坐标 δ_t

$$v_t = -(1 + \delta_t)v^* \tag{123}$$

δ_t 的动力学由下式给出

$$\delta_{t+2} = (1 - 8z_t^2)\delta_t \tag{124}$$

由于 $8z_t^2 < 1$, δ_t 的大小严格递减。我们可以从上面限制 δ_t

$$|\delta_t| \leq \exp \int_0^t -8z_s^2 ds |\delta_{t-1}| \tag{125}$$

由于 δ 开始为负值,并且大小逐渐减小,我们知道 $z_t > v^*$ 动力学方程,我们可以将 z_t 限制, 这也意味着 $y_t > v^*$ 。来自为

$$z_t \geq 2e^{-t} - z_0 \tag{126}$$

通过选择足够小的 B (关于 ϵ 一致)。代入给出了 δ_t 的以下界限:

$$|\delta_t| \leq \exp \int_0^t 4e^{-2s} z_s^2 ds |\delta_{t-1}| \tag{127}$$

使用总和的积分近似,界限变为

$$|\delta_t| \leq \exp \int_0^t -2z_s^2 ds |\delta_{t-1}| = \exp \int_0^t -16z_s^2/(1 - e^{-2t}) ds |\delta_{t-1}| \tag{128}$$

根据动力学方程, $-1 \leq \delta_{t-1} \leq 0$ 。在大 t 的极限下,我们有

$$\lim_{t \rightarrow \infty} |\delta_t| \leq \exp \int_0^t -16z_s^2/(1 - e^{-2t}) ds |\delta_{t-1}| \tag{129}$$

如果我们有条件的

$$z_t \geq -\log(1/16) \tag{130}$$

那么 $\lim_{t \rightarrow \infty} |\delta_t| \leq$ 因此在假设 $z \sim O(B^2 \log(\cdot))$ 下。请注意,要求 $z \geq c \log(\cdot) \geq \log(\cdot)/16$,我们有 $\lim_{t \rightarrow \infty} v_t = -1/2 + \log(\cdot) \geq c \log(\cdot)$ 那么通过边界的均匀性我们可以减少 $z \geq c \log(\cdot)/16 \leq z \leq c \log(\cdot)$ 。如果直到所有边界同时遇见了。

类似的论证表明 $\lim_{t \rightarrow \infty} u_t = -1/2 - O(B^2 \log(\cdot))$,在 $z \sim 0$ 上相同的条件下。来自论:

$$\lim_{t \rightarrow \infty} y_t = -1/2 + \delta^*, |\delta^*| \leq B^2 \text{对数}(\cdot) \tag{131}$$

在我们的假设下。

现在我们可以证明定理3.1的陈述。给定一个 \leq 的模型, $\theta \in \eta$ 空间和 $z \in y$ 空间之间存在连续映射。由于 $z \in y$ 空间中的区域 $[c \log(\cdot), z \sim c] \times [0, y_c)$ 显示了独立的 $z \sim c$ 和 y_c 的稳定边缘行为。这也对应于 $\{z \geq 0, \eta \lambda \max, \eta\}$ 中的独立邻域 U 。证明到此结束。

□

通过运行各种初始化的动力学方程、计算中位特征值（限制在范围[1.9, 2.0] 内)并绘制对比图（图 10,左）,可以从数值上确认此结果。为了进行比较,我们还绘制了动力学 $\dot{z} = 2yz$ 的低阶 ODE 近似的 y 极限值

(132)

$$y \cdot = -2(4 - 3 + 4^2)yz^2 - 4z^2 \tag{133}$$

我们还从 $2 - \sqrt{2}$ 获得极限 $-1/2 + O(\sqrt{2})$ （图 10,右）。 $\sqrt{2}$),花费更少的精力。对于小 ϵ , 两个量都显示为 $O(\epsilon^2)$ 偏差

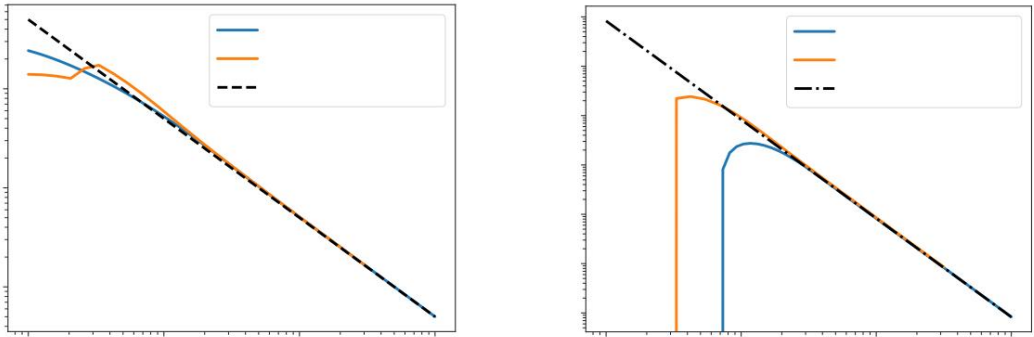


图 10. y 的最终值,与临界值 $T(0) = 2$ 的归一化偏差,适用于离散动力学和 ODE 近似。在较大范围内偏差可以很好地近似为 $\sqrt{2}$ （左）。与 $\sqrt{2}$ 的偏差为 $O(\epsilon^2)$ （正确的）。

B.7.参数空间 vs. $z \sim T$ 空间

我们的大部分分析都集中在归一化的 $z \sim T$ 坐标空间中。在本节中,我们确认参数空间中更常见的设置与归一化坐标空间一致。特别是,EOS 行为通常通过修复初始化和使用不同学习率进行训练来描述 - 如图 1 所示。

我们可以绘制图 1 中轨迹的 $T(0)$ 动态。我们看到,对于小学习率,会收敛到 $T(0) < 2$,对于大学习率,会出现发散,而对于中等学习率,会收敛到 $2 - \sqrt{2}$ （图 11）。

这证实了该定理对于描述发现和探索 EOS 行为的更传统方法很有用。

C. 二次回归模型动力学

我们在本节中使用爱因斯坦求和符号 - 方程右侧的重复索引被认为是求和,除非它们出现在左侧。

C.1.定理4.1的证明

让 z 、 J 和 Q 使用均值为 0 且方差为 σ 的 iid 随机元素进行初始化 $\epsilon, \epsilon_{\alpha}, \epsilon_{\alpha\alpha}$, 和 1 分别。此外,令分布对于数据空间和参数空间中的旋转不变,并且具有有限的四阶矩。

为了理解早期曲率的发展,我们考虑将 J 转换为其奇异值形式的坐标。在这些坐标中,我们可以写:

$$J_{\alpha i} = \begin{cases} 0 & \text{如果 } \alpha = i \\ \sigma \alpha & \text{如果 } \alpha \neq i \end{cases} \tag{134}$$

奇异值 $\sigma \alpha$ 是 NTK 矩阵奇异值的平方根。我们假设它们的大小按从最大 (σ_1) 到最小的顺序排列。假设,在这种旋转下, z 和 Q 的统计量保持不变。

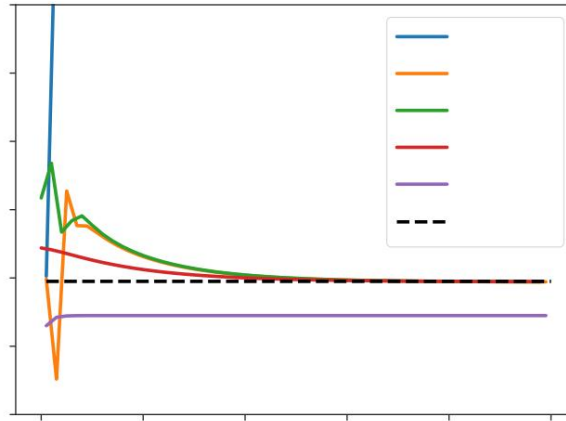


图 11.图 1 中轨迹的T动态。对于较小的学习率 η ,轨迹收敛到 $T < 2$,对于较大的学习率,轨迹发散。对于中间轨迹,我们有 EOS 行为,其中最终 T 由定理 3.1 预测。

$t = 0$ 处的时间导数可以直接在奇异值坐标中计算。一阶导数由下式给出

$$\frac{d}{dt} \sigma_{\alpha}^2 = 2\sigma_{\alpha} \cdot \alpha \quad (135)$$

使用对角坐标系,我们有

$$\sum \frac{d}{dt} \sigma_{\alpha}^2 = E[Q\alpha\beta_j J\beta_j z\beta] = 0 \quad (136)$$

然而,平均二阶导数为正。计算一下,我们有:

$$\frac{d^2}{dt^2} \sigma_{\alpha}^2 = 2(\dot{\sigma}_{\alpha}^2 + \sigma_{\alpha} \dot{\sigma}_{\alpha}) \quad (137)$$

我们可以在初始化时计算平均值。我们有:

$$E[\dot{\sigma}_{\alpha}^2] = E[Q\alpha\beta_j J\beta_j z\beta Q\alpha\delta k J\delta k z\delta] = E[\delta\beta\delta j k J\beta_j J\delta k z\beta z\delta] \quad (138)$$

$$E[\dot{\sigma}_{\alpha}^2] = E[Q\alpha\beta_j J\beta_j z\beta Q\alpha\delta k J\delta k z\delta] = E[J]_{\beta_j z\beta}^2 = DP \sigma^2_{Jz} \quad (139)$$

为了计算第二项,我们计算 $J_i \alpha_i$:

$$J_{\alpha i}^{\circ} = -Q\alpha_{ij} (J\beta_j z \cdot \beta + J \cdot \beta_j z\beta) \quad (140)$$

扩展一下,我们有:

$$J_{\alpha i}^{\circ} = Q\alpha_{ij} (J\beta_j J\beta_k J\delta k z\delta + Q\beta_j k J\delta k z\delta z\beta) \quad (141)$$

在对角坐标中 $J\alpha\alpha = \sigma_{\alpha}$ 。这给了我们:

$$E[\sigma_{\alpha}\sigma_{\alpha}] = E[\sigma_{\alpha} Q\alpha_{aj} Q\beta_j k J\delta k z\delta z\beta] \quad (142)$$

对Q 进行平均,我们得到:

$$E[\sigma_{\alpha}\sigma_{\alpha}] = PE[\sigma_{\alpha}\delta\alpha\beta\delta\alpha k J\delta k z\delta z\beta] = E[\sigma_{\alpha} z\alpha z\delta J\delta\alpha] \quad (143)$$

二阶回归和稳定边缘

其评估结果为：

$$E[\sigma \alpha \sigma \alpha] = \sigma^2 E[\alpha^2] \quad (144)$$

在大D和P的限制下,对于固定比率D/P, Marchenko-Pastur 分布的统计数据允许我们计算最大本征模态的导数：

$$E[\sigma_0 \sigma_1] = \sigma^2 \sqrt{2D(1 + D/P)} \quad (145)$$

综合起来,这给了我们

$$\frac{d^2 \lambda_{\max}}{dt^2} = \frac{2}{JDP(P(1 + D/P) + 1)} \quad (146)$$

我们在图 12 中用数字证实了这一预测。

即,最大曲率的二阶导数平均为正。如果我们对特征值尺度进行归一化,在大 D 和 P 的限制下,我们有：

$$\frac{d^2 \lambda_{\max}}{dt^2} / E[\lambda_{\max}] = \frac{2}{z} \quad (147)$$

因此,增加 σz 会增加 λ_{\max} 轨迹的相对曲率。这给了我们定理4.1的证明。

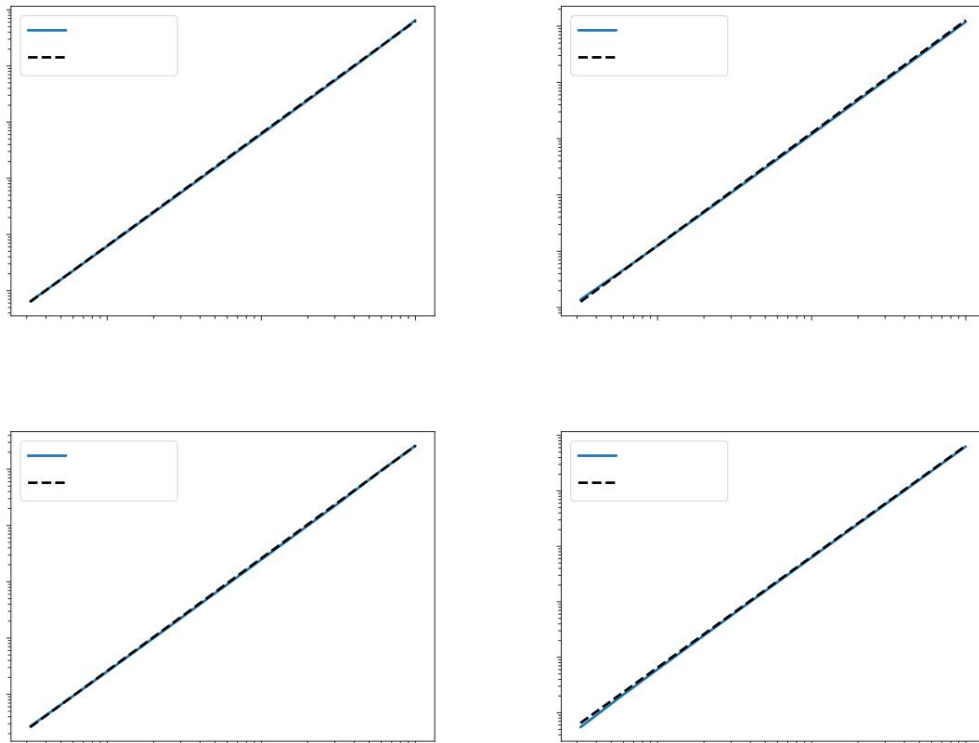


图 12. 平均 $\lambda_{\max}(0)$ 与 σz ,不同的 D 和 P (100 个种子)。

该结果表明,随着 σz 的增加,渐进锐化的程度也会增加。这可以通过查看GF 轨迹来确认(图 13)。 σz 较小的轨迹的曲率变化不大,并且损失以某种速率呈指数衰减。然而,当 σz 较大时,曲率最初会增加,然后稳定到较高的值,从而可以更快地收敛到损失的最小值。

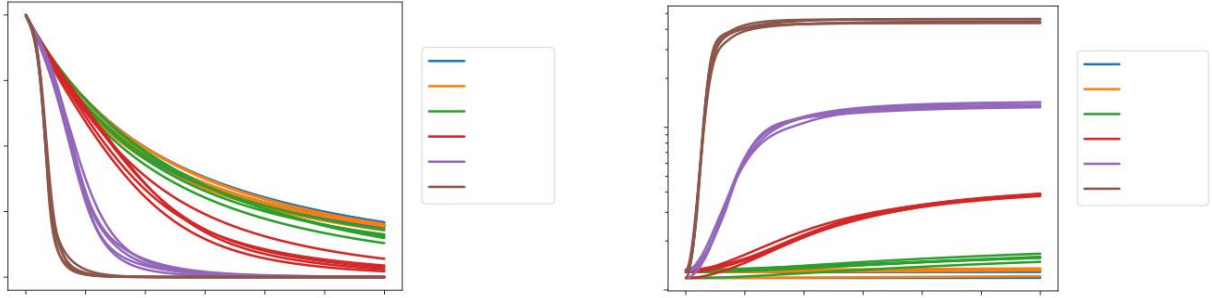


图 13.不同 σ_z 的二次回归模型的损失和最大 NTK 特征值的梯度流轨迹。随着 σ_z 的增加, λ_{\max} 变化得更快,并且通常会增加。 σ_z 较高的模型在 GF 动力学中收敛得更快。

C.2. 梯度下降的时间尺度

考虑 z 、 J 和 Q 的随机初始化,其中项是独立同分布且均值方差为零阶矩。此外,假设 z 、 J 和 Q 在输入和输出空间中都是旋转不变的。在这些条件下,我们希望计算

$\frac{2}{z}$, $\frac{2\sigma_z}{J}$, 和 1, 以及有限的四

$$\frac{2r}{\text{荷兰}} = \frac{E[|\frac{1}{2} \eta \sum_{i,j} Q_{aij} (J\beta_i)^0(z\beta)^0(J\delta_j)^0(z\delta)^0|^2]}{E[|\sum_{i,j} (J\alpha_i)^0(J\beta)^0(z\beta)^0|^2]} = \frac{1}{4} \eta \sum_{i,j} 2zD^2 \quad (148)$$

初始化时,在大 D 和 P 的限制下。

分母由下式给出:

$$E[J\alpha_i J\beta_i(z\beta) J\alpha_j J\delta_j(z\delta)] = \sigma \frac{2}{z} E[J\alpha_i J\beta_i J\alpha_j J\delta_j \delta\beta\delta] = \sigma z E[J\alpha_i J\beta_i J\alpha_j J\beta_j]^2 \quad (149)$$

评估给我们带来:

$$E[J\alpha_i J\beta_i(z\beta) J\alpha_j J\delta_j(z\delta)] = \sigma \frac{2}{z} (\frac{4}{9} (P(P-1)D) + C_4 DP) \quad (150)$$

其中 C_4 是 $J\alpha_i$ 的四阶矩。 D 和 P 中的最低阶

$$E[J\alpha_i J\beta_i(z\beta) J\alpha_j J\delta_j(z\delta)] = \sigma z \sigma^{-24} JDP^2 + O(DP) \quad (151)$$

评估分子,我们有:

$$E[Q_{aij} J\beta_i J\delta_j z\delta Q_{amn} J\gamma_m z\gamma J\gamma_n z\gamma] = E[J\beta_i J\delta_j z\delta J\gamma_m z\gamma J\gamma_n z\gamma] (\delta_{im} \delta_{jn} + (M_4 - 1) \delta_{ijmn}) \quad (152)$$

其中 M_4 是 Q_{aij} 的四阶矩。

这给了我们:

$$\frac{1}{D} E[Q_{aij} J\beta_i J\delta_j z\delta Q_{amn} J\gamma_m z\gamma J\gamma_n z\gamma] = E[J\beta_i J\delta_j z\delta J\gamma_i z\gamma J\gamma_j z\gamma] + (M_4 - 1) E[J\beta_i J\delta_i z\delta J\gamma_i z\gamma J\gamma_j z\gamma] \quad (153)$$

接下来,我们执行 z 平均值。我们有

$$\begin{aligned} \frac{1}{D} E[Q_{aij} J\beta_i J\delta_j z\delta Q_{amn} J\gamma_m z\gamma J\gamma_n z\gamma] &= \sigma z E[J\beta_i J\delta_j J\gamma_i J\gamma_j] (\delta\beta\delta\gamma + \delta\beta\gamma\delta + \delta\beta\gamma\delta\gamma) \\ &\quad + (C_4 - \sigma^4) E[J\beta_i J\delta_j J\gamma_i J\gamma_j] \delta\beta\delta\gamma \\ &\quad + (M_4 - 1) \sigma^4 E[J\beta_i J\delta_i J\gamma_i J\gamma_j] (\delta\beta\delta\gamma + \delta\beta\gamma\delta + \delta\beta\gamma\delta\gamma) \\ &\quad + (M_4 - 1) (C_4 - \sigma^4) E[J\beta_i J\delta_i J\gamma_i J\gamma_j] \delta\beta\delta\gamma \end{aligned} \quad (154)$$

其中 C_4 是 z 的四阶矩。简化给我们:

$$\begin{aligned} \frac{1}{D} E[Q_{aij} J_{\beta i z} J_{\delta j z} Q_{amn} J_{\gamma m z} J_{\nu n z}] = & \sigma^4 (E[J_{\beta i} J_{\beta j} J_{\delta i} J_{\delta j}] + E[J_{\beta i} J_{\delta j} J_{\beta i} J_{\delta j}] + E[J_{\beta i} J_{\delta j} J_{\delta i} J_{\beta j}]) \\ & + (C_4 - \sigma^4) E[J_{\beta i} J_{\beta j} J_{\beta i} J_{\beta j}] \\ & + (M_4 - 1) \sigma^4 (E[J_{\beta i} J_{\beta i} J_{\gamma i} J_{\gamma i}] + E[J_{\beta i} J_{\delta i} J_{\beta i} J_{\delta i}] + E[J_{\beta i} J_{\delta i} J_{\delta i} J_{\beta i}]) \\ & + (M_4 - 1)(C_4 - \sigma^4) E[J_{\beta i} J_{\beta i} J_{\beta i} J_{\beta i}] \end{aligned} \quad (155)$$

对于较大的 D 和 P ,最后三项渐近小于第一项。评估第一项,对于领先顺序,我们有:

$$\frac{1}{D} E[Q_{aij} J_{\beta i z} J_{\delta j z} Q_{amn} J_{\gamma m z} J_{\nu n z}] = \sigma^4 \left(\frac{1}{D} (2DP^2 + 2D^2P + D^2P^2) \right) + O(D^2P + DP^2) \quad (156)$$

$$E[Q_{aij} J_{\beta i z} J_{\delta j z} Q_{amn} J_{\gamma m z} J_{\nu n z}] = \sigma^4 \left(\frac{1}{D} (2DP^2 + 2D^2P + D^2P^2) \right) + O(D^2P + DP^2) \quad (157)$$

这给了我们:

$$\frac{1}{D} E[Q_{aij} J_{\beta i z} J_{\delta j z} Q_{amn} J_{\gamma m z} J_{\nu n z}] = \frac{\sigma^4}{4} \frac{1}{D} (2DP^2 + 2D^2P + D^2P^2) = \frac{1}{4} \sigma^4 \frac{1}{D} (2DP^2 + 2D^2P + D^2P^2) \quad (158)$$

在大 D 和 P 的限制下,达到领先顺序。

C.3.对 D 和 P 的依赖

我们可以凭经验看到,在 $D > P$ 的超参数化状态下,锐化更加明显。使用图 4 中的轨迹,我们可以绘制动态初始化点和最终点处归一化最大 NTK 特征值 $\eta \lambda_{\max}$ 的散点图 (图 14)。在所有情况下,各种初始化 (x 轴)都会导致最终值集中在 2 (y 轴)周围。

我们可以看到,在 $P > D$ 的过度参数化状态下,浓度最严格 (右图)。我们假设对于较大的 D 和 P ,当 $P > D$ 时,EOS 行为更强并且更有可能发生。我们将对该假设的进一步探索留到未来的工作中。

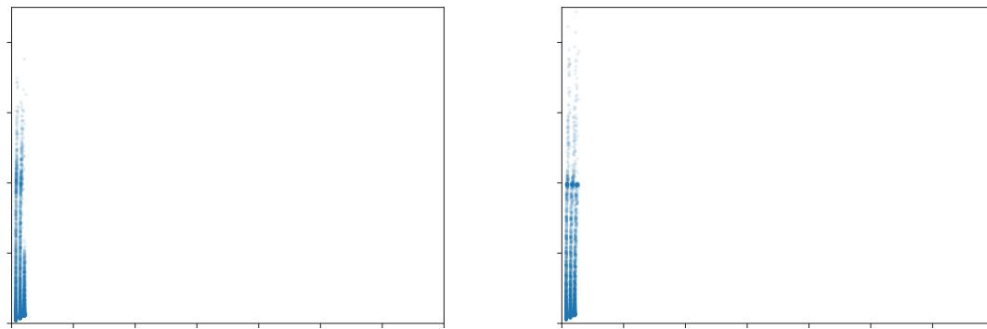


图 14.二次回归模型的初始与最终归一化最大特征值 $\eta \lambda_{\max}$ 的散点图。轨迹取自用于生成图 4 的数据。对于较大的 D 和 P ,随着模型变得过度参数化($P > D$),轨迹的子集显示出更严格的 EOS 行为,其中对于各种初始化, $\eta \lambda_{\max}$ 集中于接近 2。

D. 真实模型分析

D.1.实验装置

第 5 节中的实验是按照 Cohen 等人的设置设计的。（2022a）。我们在 10 类 CIFAR10 的前 5000 个样本上训练了一个具有 tanh 非线性和隐藏维度 200 的两隐藏层全连接网络。

我们使用 Lanczos 方法计算了 NTK 的光谱、损失 Hessian 和 Q（Ghorbani 等人,2019）。我们对 NTK 和 Q 使用阶数 36,对损失 Hessian 使用阶数 50。所有实验均在具有 float32 精度的 GPU 上进行。

D.2. CIFAR10 模型中y的动态

第 5 节中分析的 CIFAR10 模型中y的动态比z1动态更复杂。从图 5 中我们可以看出,y的两步变化有一个与z1和y 无关的分量。我们可以通过计算小z1的yt+2 - yt的平均值（在本例中取z1 < 10-4）来近似该变化b。然后我们可以从yt+2 - yt中减去b,并将余数相对于z 绘制出来,特别是对于较大的zt。然而, yt+2 - yt显然不是简单的z1 函数。

$\frac{y_{t+2} - y_t}{z_1^2}$ （图 15 左）。我们看到yt+2 - yt - b与z负相关 $\frac{y_{t+2} - y_t - b}{z_1^2}$

两步模型动态可以写为 $(ay + c) \sim z$ 如果我们绘制 $(y_{t+2} - y_t - b)/z_1^2$ 与yt的关系图,我们同样没有单值函数（图 15,正确的）。因此, yt+2 - yt的函数形式不是由 $b + ayz^2 + cz^2$ 给出 $\frac{y_{t+2} - y_t - b}{z_1^2}$

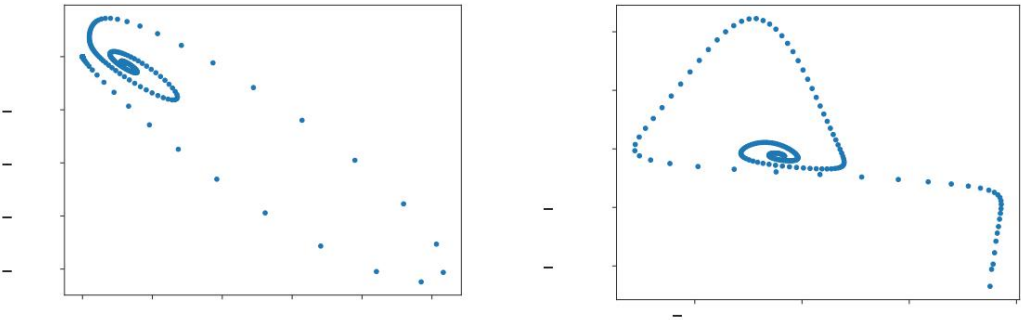


图 15.不同 σ_z 的二次回归模型的损失和最大 NTK 特征值的梯度流轨迹。随着 σ_z 的增加, λ_{\max} 变化得更快,并且通常会增加。 σ_z 较高的模型在 GF 动力学中收敛得更快。

D.3. 2 类 CIFAR 模型的二次展开

我们使用 Neural Tangents 库（Novak等人,2019)仅使用前两个类和 5000 个数据点来训练 CIFAR 模型 - 这让我们可以在任意参数下执行模型的二阶和三阶泰勒展开。该模型是2 隐藏层全连接网络,隐藏宽度为256 , Erf非线性。模型使用NTK 参数化进行初始化,权重方差为1 ,偏差方差为0。目标的标量值为 -第一类为+1 ,第二类为-1。所有实验中均使用0.003204的学习率。所有绘图均使用 float-64 精度绘制。

在初始化时进行二次展开,我们发现损失在该设置中跟踪前1000 个步骤的完整模型（图 16,左）,但错过了稳定边缘行为。我们使用神经切线来有效地计算 NTK,以获得顶部特征值 λ_1 （从而得到y）。我们还可以通过计算相关的特征向量v1和投影残差z 来计算z1。如果二次展开更接近稳定边缘,则z1的动态特性非常接近真实的z1动态特性,直至与不同时间发生的z1指数增长相关的偏移（图 16,中）。我们看到|z1|中第一个峰的形状对于完整模型和二次模型来说是相同的,但是在完整模型中后续振荡更快并且阻尼更快。这表明二次模型可以捕获初始 EOS 行为,但详细的动态需要了解高阶项。例如,三阶泰勒展开改进了对振荡幅度和周期的预测,但仍然错过了关键的定量特征（图 16,右）。

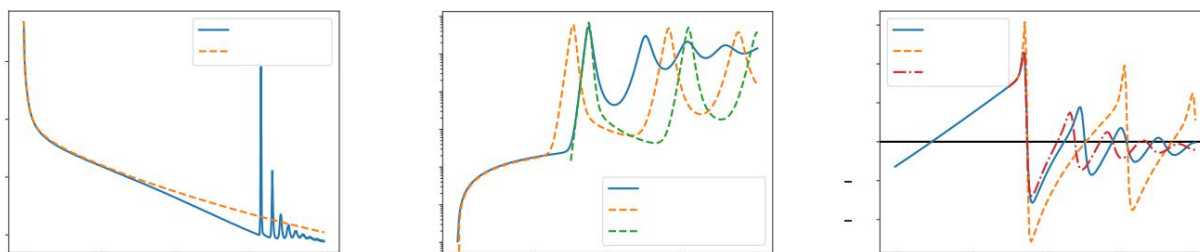


图 16.在二类 CIFAR 上训练的 FCN 模型的二次展开。初始化时展开可以很好地近似1000步的完整模型,之后 EOS 行为会出现在完整模型中,但不会近似为 1 (左)。当 z_1 较小时,二次模型跟踪完整模型;然而,在近似模型中,初始指数增长可能会更早发生(中)。与近似模型相比,完整模型中 z_1 的幅度具有更大的振荡。三阶泰勒展开式更好地捕捉了振荡的幅度和周期,但仍然错过了定量特征(右)。