Ниже приведён подробный план (и набор идей) по тому, как

- Определен задысновен или (и пакор идея) и тиму, как положивает в радмента, к съсма и пева от пределен дализителение на гена ительстви да и по пределен за проблемные с логи (илизком большая или с пишком може чрествительность, нарушение с стабильности отпичаации, "затыкание" обуч
 Проверять адекватность градментов нет ли "аэрывающихся" или "затухающих" градиентов в отдельных словх.
- 3. Оценивать достаточность/избыточность нейронов (а возможно, и самого слоя) когда стоит добавить нейроны, когда нужно от них избавиться или когда выгоднее убрать весь слой.

1. Общие критерии анализа слоя

1.1 Норма градиента и её динамика

 $m{\Phi}$ ормулировка $m{\Gamma}$ усть $m{L}(heta)$ - функция потерь (loss), где heta - весь вектор параметров сети. Рассмотрим конкретный слой $m{L}$ с параметрами $m{W}_L$. Обозначим градиент функции потерь

 $\nabla_{W_L}L(\theta)$.

Одна из простейших и при этом достаточно показательных метрик – это норма данного градиента.

- **пьшая** норма $\|
 abla_{W_L} L(heta) \|$ может говорить об «эффекте взрыв:
- Слишком малая (близкая к нулю) норма признак того, что слой практически «не участвует» в процессе обучения, то есть может быть «заморожен» или «выучился до плато» слишком рано.

- 1. Рассчитать $\|\nabla_{W_r} L(\theta)\|$ (часто используют ℓ_2 -норму) для каждого слоя в процессе обучения
- Анализировать динамику во времени смотреть, как эта норма меняется по эпо

 $G_L = rac{\|
abla_{W_L} L(heta)\|_2}{\|W_L\|_2}$ или $rac{\|
abla_{W_L} L(heta)\|_2}{\sum_j \|
abla_{W_j} L(heta)\|_2},$

Если G_L слишком велик или, наоборот, стабильно чрезвычайно мал, есть повод «прицельно» изучать этот слой

1.2 Анализ гессиана (вторые производные)

Гессиан по параметрам слоя W_L – это матрица вторых производных:

 $H_L = \nabla^2_{W_L} L(\theta).$

На что смотреть

- Сильная разница между максимальным и минимальным собственным значением (большая число обусловленности) говорит об иллю-conditioned слое.
- 2. Сед тессиана Тт.(H_) эти сумма диагомальных элементов трубо показывает суммарную «крупцину» слоя Если Тт.(H_) очень мала или близка к мулю, слой почти не влияет на функцию потерь (по вторым производ 3. Ранг (или эффективный ранг) гессиана если он явно меньше числа параметров, часть параметров может быть низбыточной» (не задействованной в кривизне целевой функции).

- ти все $\lambda_i pprox 0$, может быть **избыточным** он «не вносит значимого вклада» во вторую производную • Слой, у которог

2. Определение проблемных слоёв

- 1. Или взрывает градиенты, делая обучение неустойчивым
- 2. Или застыл (у него очень малая норма градиента, в то время как сеть ещё далека от хорошего решения).
- Или имеет крайне «узкую» область допустимых шагов (очень высокое максимальное собственное значени:

2.1 Критерии «адекватности» градиента

- Норма градиента: $\| \nabla_{W_L} L(\theta) \|$ не должна стабильно быть на порядки выше/ниже, чем в других слоя:
- Отношение к параметрам: $\frac{(T_{V_{i}}(L(t))}{\|R_{i}\|_{2}}$ дейт представление о том, насхолько активно слой обновляется.
 Сигнатура тессиван: если вътрица H_{L} сильно плохо обусловлена (очень большая $\kappa(H_{L})$, где κ число обусловлена (очень большая $\kappa(H_{L})$, где κ число обусловлена (очень большая $\kappa(H_{L})$).
- Обратная связь по активациям: заодно можно смотреть, каково распределение активаций слоя (например, процентовка «сытых» ReLU, если это ReLU-слой, или средний градиент по вход

3. Детальный анализ: когда слою не хватает нейронов, а когда их слишком много

3.1 «Не хватает нейронов» (признаки недообучения слоя)

- 1. Высожий рост ошибки при абляции (отключении) части нейронов слоя. Если даже небольшое уменьшение числа нейронов приводит к эначительному росту ошибки на валидации, это может говорить, что слой сильно загружен, и в нём нет «избыточных» нейс
- . высомирост шармам так малам от иличения чест и пераронай слок стоя в перагонай слок об в порядка и перагона и перагона

• С помощью гессивна/Фишера можно смотреть «эффективный рант» (талік,) в подпространстве параметров. Если данный ранг близок к максимальному числу нейронов и слой при этом испытывает затруднения в обучении (граднент постоянно большой, ошибка не убывает), можно предположить, что пропусной способностии слоя.

3.2 «Слишком много нейронов» (признаки переобучения или избыточности)

- 1. Многие нейроны «неиспользуемы»: модули, у которых практически нулевые веса или активации очень редко выходят
- много не недотны этекстионых регемах. Водуть, у потрых цых ического мужение еще то на к изведию очеть раско выходия из лучую им тех.) прихыже, что некоторые техроно не примости польже. Намакой разит всема сели не всема на межение в неготным в неготны 3. Гессиан с малыми собственн

• При пошаговом удалении (pruning) нейронов можно отслеживать изменение потерь на обучающей и валидационной выборке. Пока удаление неёронов не ухудшает метрики (или ухудшение мало), значит, сеть действительно избыточна. Как только ошибка начинает расти, мы нашли «грань» разумной очистки

4. Когда слой мешает аппроксимации задачи и является избыточным целиком

Иногда целый слой может оказаться **лишним**. Например, если у нас слишком глубокая архитектура при относительно простой за

- 1. Удаление слоя не ухудшает метрики (или ухудшение минимально). Можно провести «абляционное» исследование: берём обученную сеть, вырезаем слой (или заменяем на тождественное отображение), делаем дообучение (fine-tuning) оставшейся части и смотр
- 2. Слой не даёт градиентного вклада: $\|\nabla_{W_c}\| \approx 0$ в течение долгого времени обучения, при том что другие слои ещё активно учатся (т. е. не общее плато обу

• Анализировать важность слоя аналогично тому, как анализируют важность отдельного нейрона, но «гло Если $\Delta L_{\mathrm{remove}\ L}$ мало, значит слой не критичен.

 $\Delta L_{\mathrm{remove}\,L} = L(heta_{\mathrm{fea}\,W_L}) - L(heta_{\mathrm{no}}$

5. Резюме и последовательность действий

- 1. Собираем основную статистику по каждому слою на нескольких шагах обучения (не только в конце):
- $\|\nabla_{W_L}L(\theta)\|$, $\|\nabla_{W_L}L(\theta)\|/\|W_L\|$
- Собственные значения (или хотя бы
- Средние/дисперсии активаций.
- Если слой выбивается по норме градиента (слишком большой/маленький) проверяем инициализацию, функцию активации, регуляризацию и т. д.
- Если гессиан слоя имеет высокую обусловленность адаптируем шаг обучения, возможно, вводим методы сглаживания (BatchNorm, Gradient Clipping). веряем мощность слоя (под-или пере-способность):

- Анализируем эффективность каждого нейрона (частота активации, величина весов, вклад в ошибку).
- Если подозреваем недогрузку добавляем нейроны и смотрим динамику обучения.
- Если подозреваем переизбыток производим pruning (например, удаляем наименее активные/наименее «важные» нейроны по критерию $\|\nabla\|$ или $\|u\|$, или по убыванию собственных чисел, и контр

• Проводим абляцию (удаляем слой) и делаем дообучение. Смотрим, сильно ли упадёт качество. Если падение невелико, слой – кандидат на вычеркивание

Такой подход двёт универсальные (пусть и в основном эвристические) способы решения поставленных задач. Ключевым моментом является систематический сбор метрик и их интерпретация в контексте динамики обучения (в не только на одной «замор

Заключительные замечания

- 1. Выбор порогов (когда считать градиент «слишком мальм» или «слишком большим») вопрос практики. Обычно ориентируются на средние/стандартные отклонения по слоям и ставят 2–3 сигмы, чтобы отследить аномалии
- 2. Числовая оценка гессиана на практике часто затратна, поэтому применяют приближённые методы (например, диагональ Фишера, К-FAC, Gauss-Newton аппроксимации).
- В Расширениемуменьшение соль на легу (Dynamic Neural Networks) привожно избертно в принятия и принятия и

0) D 70 C V