9/23 「EECE695D」

↳ scaling works... why?

9/22~9/30 「find your team online
(PLMS)」

~10/4 = choose a paper, write a proposal
~10/18 = visit the department office
(with me), and get colab account.

~11/11 = experiments, write a paper,
and submit via open Review.

~11/22 = write reviews. (report)

「team ↓↓ ⇒ 점수 ↑↑」

↳ ML domain에 대한 논문 다루기

───────────────

P. 1) ① old (~2010s) = Big models do not
generalize.

② Modern (2010s~) = Big models
generalize better.

↓

「good old 이론」

「SLT」
↳ statistical Learning theory
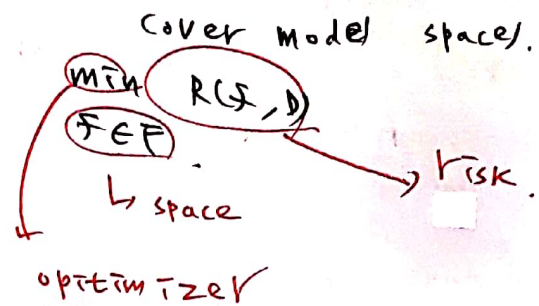
───────────────

① Model space 'F' = A set of all 'learnable'
functions 「학습가능한 함수의 집합 이라는
space.」

② RISK, $R(f, D)$ : Loss 「Predictive quality
of function on dataset $D = \{z_1, z_2, \dots, z_N\}$」

───────────────

① optimizer = ↓ risk

cover model space.

$$\min_{f \in F} R(f, D)$$

↳ space

optimizer ⟶ risk.

───────────────

$$= (R(\hat{f}, D) - R(\hat{f}_{ERM}, D))$$

optimizer error

$$+ R(f_{ERM}, P) - \min_{f \in F} R(f, P)$$

generalization error.

$$+ \min_{f \in F} R(f, P) - \min_{f} R(f, P)$$

$f$ measurable.

───────────────

Approximation error.

$$R(f, D) = \frac{1}{N} \cdot \sum_{i=1}^{N} \ell(f(x_i), y_i)$$

$$R(f, P) = E_{z \sim p} \ell(f(x), y)$$

「$x_1, x_2, \dots, \sim IID$」

$$\frac{1}{N} \cdot \sum_{i=1}^{N} x_i$$

⟨1⟩

estimation error
= optimization error
+ generalization error. } 에러.

「generalization error」
↳ D와 P의 mis match.
small 「data> param」 (not anymore)

& D = gradient descent ?
「ground − truth, cubic function」 파라곱함



(d=3),  Legendre. ?
    polynomials.

「deep − double − descent.」

d = 20 → smaller norm gradient.

deep − double − descent = windows on theory
org.

ERM
↳ empirical risk.

$$\sup_{f \in F} | R(f, P) - R(f, D)| \leq O\left(\sqrt{\frac{M}{N}}\right)$$

「M 파라미터들」, F: parametrized model space 들.

hold with high probability.

under fitting 과소적합 | over−parameterized.

$$param. \begin{cases} data \\ set \end{cases} < param$$

estimation error.

But 경험적으로.

↳ interpolation threshold

Moore − penrose pseudo − inverse.

↓

linear system

linear regression

$$\min_\theta \| Y - (\theta)^T \cdot X \|^2$$

<span style="color:red">공통 값 -train</span>

↳) $\theta^* = (X^T \cdot X)^{-1} X^T \cdot Y$

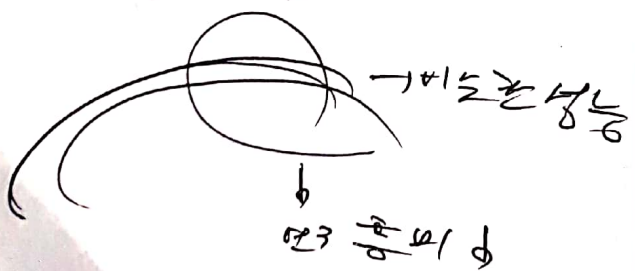linear regression
$Y = Ax \Rightarrow X^T X \theta = X^T \cdot X \cdot Ax$

→ $Y = (X^T \cdot X)^{-1} X^T X A x$
   같은거

↳) pseudo- inverse : 주로 관심.

1990s에   연구 발전 관계 도달.

Hebb, perceptron, pseudo -inverse

성능비교

→비교조성능

↓

연구동비 ↓

PELD 동비 ↓↓

---

2020 .   over parameter ization
관심↑↑ .



critical
regime.

test .

train

interpolation
threshold.

early stop

stop
않좋

조건성겁 시 .

⟨2⟩

「large model <> small
model < intermediate
model

layer를 많이 쌓아도 성능 개선 X

$$\underbrace{0000}_{\text{}} \to \overbrace{0000}^{\text{}}$$

(optimization) ↑↑

width를 늘번.

layer 수가 ㅇㅕ러 → vanilla로는
최적 (adaptive로도 성능 ↑)

$$\boxed{\begin{array}{l} \eta = a \cdot x \quad \to ① \\ y = b \, c \cdot x \quad \to ② \end{array}}$$

①
$\frac{1}{2}(y-ax)^2 \to \nabla_a l = -x(y-ax) \to ③$

$= r$

②$\nabla_b l = -cx(y-bcx) \to ④$
$\nabla_c l = -bx \cdot (y-bcx)$

---

$(a - \lambda \nabla x)$

$= a - \nabla_a \frac{1}{} \cdot \textcircled{가}$

$(b - \lambda \nabla_b)(c - \lambda \cdot \nabla_c)$

$= (bc + \lambda \gamma c_x^2 + \lambda \gamma b_x^2$
$+ \lambda \boxed{\quad})$

$= bc + \lambda \gamma (c^2 + b^2) \cdot x + \square$

$= a + \lambda \gamma (c^2 + b^2) x + \square$

optimization ↑↑ (poly nomial로도 가능).

⊙과 달리, b, c와
같은 dependent한 값이
추가된다.

but
→ depth에 대한
문제 해결로 개선이 되는
X

⊙와 달리, b, c의
같은 dependent한 값이
추가된다.

영향이
별거 없음.

$$\sup_{f \in F} |R(f, P) - R(f, D)| \leq 0(\quad).$$

SGD may bring us here. 「Im plicitly biased」

SGD

테이터, (close to origin)

But not here!

$\nabla_\theta (L(\theta))$

$= \sum_{i=1}^{N} r_i \cdot x_i$

$\theta^{(k)} = \theta^{(0)} + \sum_{i=1}^{N} z_i \cdot x_i$

close

$\theta^{(0)} = 0$ 子 then $\theta^{k} = ?$

$(y_i - \theta^T x_i)^2$

$\downarrow$

$2(y_i - \theta^T x_i) \cdot (-x_i)$

⟨2⟩

$\theta^T \cdot x_i - y_i = 0$ 이어야.

$\downarrow$

orthogonal to the set

$= \sum r_i \cdot x_i$

$\nabla_\theta (l(\theta)) = \sum_{i=1}^{N} (r_i) \cdot x_i$

zero loss.

GD approachable.

origin . minimum norm — zero loss

$r_i$는 x2 이나 어떤 scalar Value

**Missing Links**

GD — found solution

Big model, optimizer better.

$\hat{y} = w \cdot x$

「weight」

DLT&dev learning theory

big model benefit
↳ 성능이 좋은거

big model의 이점1/ 이점2

(제한 틀안에서)

Most $$$ =??

momentum = DL
=?
이론학과 — 인공지능 방법론가.

vice versa

if days, non-embedding dataset size — tokens
non — embedding

Hw — 답답
어둡

Introducing **whisper** 「9/22」

↳ micro Machine Man presenting the most midget miniature motorcade of micro Machines.

「automatic speech recognition (ASR) system」

---

「Neural Scaling laws」

PLMS team matching starts ~12PM today. (11/29, 12/1 : No class).
How the
⎧ ①
⎨ ②        〈1〉
⎩ ③

---

A belief (that) 'Scaling up' is the ultimate answer.

**GPT-3** ⇒ 미래에 도움 (발전에)
Why? ∞이 아닌 한정적인 파라미터로 train 하고 최적화 가능하니까.

「give it the compute, give it the data and.  」

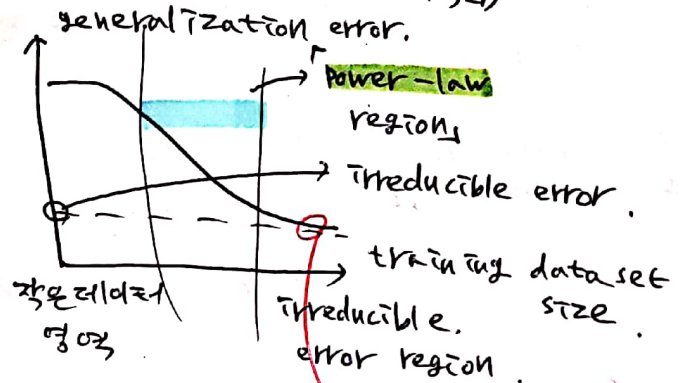DL scaling is predictable, empirically.

---

## L vs D

방법
① 큰 dataset 선택
② **SOTA** model 선택.
③ draw '$z_1, z_2, \cdots$' fraction of the dataset. (different size)

ex)
batch size,
learning rate

④ train modes of various sizes to find the one with smallest test loss. ( hyper parameter 최적화)

generalization error.



→ **power-law** region」
irreducible error.
training dataset size.
작은데이터 영역
irreducible. error region.
↳ not really observed though

향상된 아키텍처, 파라미터는 exponent를 성능 향상 시키저 X.

$$gen. \ error \le O\left(\frac{complex}{\sqrt{D}}\right) \leftarrow \beta_g = \frac{-1}{2}$$

$$\le O\left(\frac{complex}{D}\right)$$
$$\uparrow \beta_g = -1$$

$$Test \ loss \simeq d_g \cdot \mathcal{D}^{\beta_g} \ \lceil \beta_g \in [-\frac{1}{2}, 0) \rceil$$

$$\mathcal{E}(m) = \alpha \cdot m^{\beta_g} + \gamma$$

CRHN3 LSTMS

① 대비경 = 너의 모델이 power law와 far
away 하면, maybe you did something
wrong?

per set

⊢ p challenge = slope 바꾸는게 가능?

「maybe a better way to evaluate
model, algorithms」

③ exploration = 작은 데이터로 더 큰 집무에
대한 적절한 모델 선택,

④ compute goal = target loss 설정, 우리는
요구된 계산의 추정치를 얻을수 있다.

「limitation = could have been explicit」

「Scaling Laws for Neural Language models」

non-embedding, tokens, PF-days.

모델 크기, dataset의 크기 ↑ ⇒ 언어 모델의
성능은 smoothly 향상.

Infinite compute: 모델 size↑, dataset
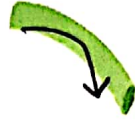성능 should be sublinear, $D \propto N^{0.095}_{0.076}$
$\sim N^{0.74}$

comput budget $C$ 「fixed」
$N \propto C^{0.73}$, $B \propto C^{0.24}$, $S \propto C^{0.03}$
⊢ optimal model size

$D = B \cdot S = C^{0.27}$
⊢ data set size

model size increase much faster
than dataset size?

---

observation -- model이 looks「depth vs
width」

but the model architecture. does
shift the plot.

LSTM VS transformer

↳ language model, computer vision

「Scaling 이슈 = depth vs width」
↳ CNN으로는 scaling 하기 ❖❖
① 비전에서의 scaling을 쓰는 이유.

$N_{opt}(C)$ $D_{opt}(C)$
$= \arg\min_{N,D}\ L(N,D)$
$N, D,\ s.t\ FLOPs(N,D) = C$

「model size & training tokens
의 수 = tradeoff」

「Chinchilla, Gopher, GPT-3,
Megatron - Turing NLG」.

ICLR = Scale efficiently
Insights

perception in teacher- student
setting.
└ punchline =✓ data wisely, then not really.

Scaling은 모델 shape에 의존X, 오직
pre-training에 hold. 「not fine-tuning」

data pruning = Beyond neural
scaling laws 「Beating power law via
data pruning」.

Beating power law

Neural「ㅗㅗㅍ」

└ pessimistic message =neural scaling
Law.

✓
Core set.

「2%의 정확도↑를 위해 dataset size
X10이 필요」→ Computer Vision.