

Knowledge distillation

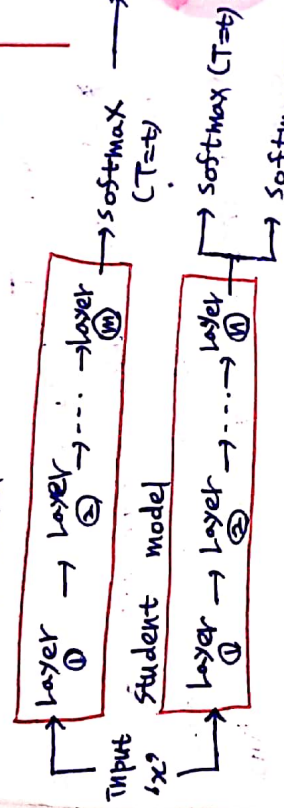
①

→ 미리 잘 학습된 큰 네트워크 (teacher networks)의 지식은 실제로 사용하고 하는 작은 네트워크에게 전달하는 것.

"Distilling the knowledge in a neural network"

* Knowledge distillation = 작은 네트워크로 큰 네트워크의 비슷한 성능을 낼 수 있도록 학습과정에서 큰 네트워크의 지식을 작은 네트워크에게 전달하며 작은 네트워크의 성능을 높여 준다는 목적.

teacher model



$$\text{Total Loss} = (1-\alpha) \cdot L_{CE}(G(Z_S), \hat{y}) + \alpha \cdot T^2 \cdot L_{CE}(G(\frac{Z_S}{T}), G(\frac{Z_T}{T}))$$

L_{CE} = cross entropy loss.

Z_S = output logits of student network

Z_T = " of teacher network

\hat{y} = ground truth (one-hot)

α = balancing parameter

→ 이가 크면 ②: 오른쪽 항 loss를 더 중요하게 보고 학습하겠다는 의미.

T = temperature: softmax 함수가 입력값이 큰 것은 아주 크게, 작은 것은 아주 작게 만드는 성질을 완화해준다.

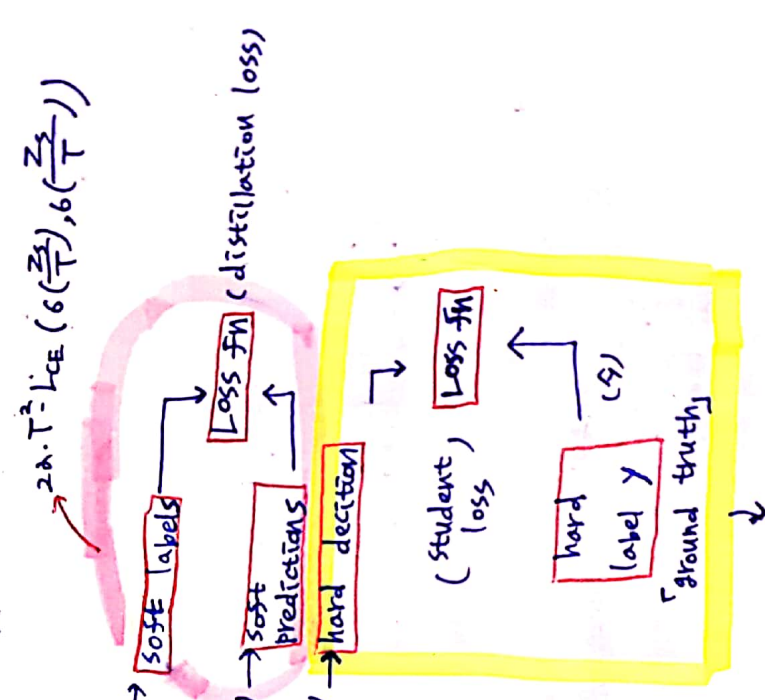
ex)

$$\begin{pmatrix} \text{bear} \\ \text{cat} \end{pmatrix} = \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \Rightarrow \begin{pmatrix} 0.95 \\ 0.05 \end{pmatrix}$$

→ 정보 손실 발생.

→ 이런 정보의 손실 없이 teacher network의 분류 결과와 student network의 분류 결과와 비교시켜, student network가 teacher network를 모방하도록 학습시킨다.

teacher, student network의 분류 결과의 차이를 loss에 포함.



$$(1-\alpha) \cdot L_{CE}(G(Z_S), \hat{y})$$

"ground truth와 student의 차이를 cross entropy loss로 계산"

* Knowledge distillation = 미리 학습시킨 teacher network의 출력을 배가 시켜서 사용하고자 하는 모델인 student network가 모방하여 학습함으로써 상대적으로 작은 parameter를 가지고 있어도 모델의 성능을 높이는 방법론.

① Sparse node activation loss

= check, variable node. L_{pnorm} is obtained by adding an L_{pnorm} on the check, variable node.

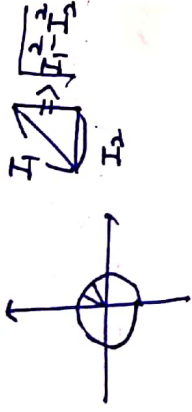
* norm = 유클리드 공간에서 벡터의 크기 \rightarrow 제곱의 제곱근

사용. $\|x\|_p = \left(\sum_{i=1}^N |x_i|^p \right)^{\frac{1}{p}}$ ②

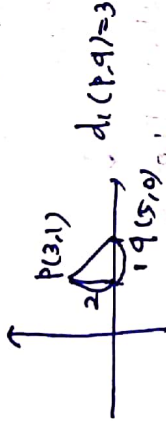
k-nearest neighbors: distance metric

↳ 2 (euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_i (I_1^i - I_2^i)^2}$$



$$d_1(p, q) = \|p - q\|, \quad p = (3, 1), \quad q = (5, 0)$$



↳ regularization

→ 기호의 cost function.

$$\text{cost} = \frac{1}{N} \cdot \sum_{i=1}^N \left\{ \frac{1}{2} \left(\frac{y_i - \hat{y}_i}{\sigma_i} \right)^2 + \frac{\lambda}{2} \|w\|_2^2 \right\}$$

↓
예측

$$L_{kd} = \sum_{t=1}^{T_{\text{student}}} L_d(t)$$

② Knowledge distillation loss

↳ imposes an additional constraint on the check node.

we suggest using the knowledge distillation method by incorporating a teacher decoder into the training procedure. → At test time, decoding is performed using the student neural decoder.

Architecture

↳ teacher decoder = min-sum

student decoder = neural min-sum

teacher > student: BER

↳ $T_{\text{teacher}} > T_{\text{student}}$

Since the teacher network has more layers, / imitating the teacher's node activations, will result in lower decoding error.

↳ B. The knowledge distillation loss term
we propose a new loss function to guide the training of the neural decoder. (student) /

teacher student

$M_{c,v}(t), M_{c,v}(t) = \text{teacher, student}$

check node messages at iteration 't'.

At the student network messages $M_{c,v}^{\text{student}}(t)$ at iteration 't' tried to mimic the teacher network messages $M_{c,v}(t+t_0)$ at iteration $t+t_0$.

$$L_{kd}(t) = \sum_{v=1}^N \sum_{c \in N(v)} \left\| \frac{\text{teacher}}{M_{c,v}(t+t_0)} - \frac{\text{student}}{M_{c,v}^{\text{student}}(t)} \right\|_p$$

$p = \text{norm order}$

c) The sparse node activation loss term sparse parity check matrix leads to higher decoding performance.

with sparse node activation (neural network decoder)

∴ The proposed sparse loss term is obtained by using a L_p norm on variable, check nodes over the student network.

$$L_S(t) = \sum_{v=1}^N \sum_{c' \in N(v)} \|b_{c'v}(t)\|_p + \sum_{c=1}^{N-K} \sum_{v' \in M(c)} \|b_{cv'}(t)\|_p$$

$$= T_{\text{student}} \sum_{t=1} L_S(t) \quad (3)$$

∴ We observe that the p parameter is important for successful training. Since p increases the gradient larger

$$\therefore L = L_{\text{ce}} + L_{\text{kd}} + \delta \cdot L_S$$

$$\delta = 1, \delta = 0.01$$

$$L_{\text{ce}} = -\frac{1}{N} \cdot \sum_v B_v \cdot \log(b_{c'v}) + (1 - B_v) \cdot \log(1 - b_{c'v})$$

$$B_v = \frac{1 - x_v}{2} \quad \text{BPSK}$$

τ bit corresponding to transmitted

symbol x_v

training exploding of phenomena =

neural decoder가 전송된 code word에 잘못된 예측률을 주는 CN, VN에 zero activation으로 출력하기 위해서 train 시기에

$$t^{\text{th}} \text{ iteration: soft output vector } s_{c'v}(t) = l_v + \sum_{c' \in N(v)} \mu_{c'v}(t)$$

$$\hat{x}_v(t) = \begin{cases} 1, & \text{if } s_v(t) > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$p_{c'v}(t) = \text{SGD} \frac{\partial}{\partial s_{c'v}(t)} \text{ loss}$$

Results

$\hookrightarrow T_{\text{teacher}} = 30$ iterations, $T_{\text{student}} = 5$ iterations.

look-ahead parameter: $\delta = 25$

※ 4dB에서의 성능 하락

= knowledge distillation loss term does not use the ground truth bits at training

Ablation Analysis = Fig. 2(a)

- knowledge distillation
- sparse node loss + cross entropy
- green = min-sum ($T=30$, teacher)

\hookrightarrow This shows that the sparse node activation loss gives more improvement than that knowledge distillation loss.

4 dB에서 성능 하락은 cross-entropy loss가 해결 가능.

Fig. 2(d) = only sparse node activation epoch 2 = 성능 ↑ (3.3dB)

\hookrightarrow For a low SNR regime, it gives high degradation in the BER performance.

Knowledge distillation

↳ one uses a teacher network to guide the training of a smaller student network.



제한된 neural decoder의

node를

제한 (constraint) 하기 위해서, 전문가

teacher에 의해 knowledge distillation method를 사용하는 것 제시. (propose to use the knowledge distillation method) by

an expert teacher network to constraint the nodes of the proposed neural decoder. (teacher)

The student tries to decode the teacher node activation. student는 teacher node activation을 복제함으로써 새로운 loss term을 추가하여 노드 전송된 codeword를 decode하기 위해 노력.

transmitted code word with a novel term that mimics the teacher node activations.

* we present novel sparse node activation loss, knowledge distillation loss term.



We demonstrate that each of the new loss terms improves the results of the baseline methods by a large margin without adding computational complexity.

previous work have tried to improve the results by changing the neural architecture, the loss function at the output layer or finding better sparse parity check matrix.



two novel loss terms.

① sparse node loss, ② knowledge distillation loss.

• sparse activation regularization =

Imposes sparse activation에 희소제약을 부과하는 새로운 loss function term. constrained on the activation with new loss function term.

Since we know that sparse parity check performance, we propose teaching the neural network decoder with sparse nodes.

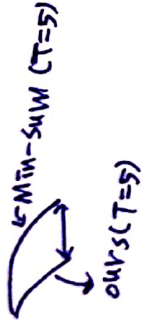
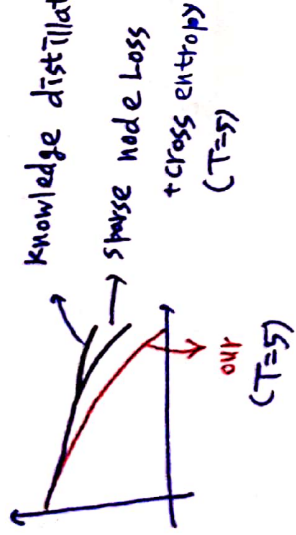
The knowledge distillation is a basic technique in deep learning, where one uses a teacher network to guide the training of a smaller student neural network.



we propose to use the knowledge distillation method by an expert teacher network to constraint the nodes of the proposed neural decoder. (student)

↳ The student tries to decode the transmitted codeword (with a novel loss term) that mimics the teacher node activations.

Fig 2-6) BCH (63,45)



NBP = ours + smaller SNR values

① neural decoder that trains only with the knowledge distillation loss (cyan)

② neural decoder that trains with the sparse node activation loss (red)

③ Min-Sum (T=30, teacher)

Magenta curve: ours (T=5)

⇒ achieves the best results which demonstrates the advantages of our method.

Ablation studies are crucial for deep learning research -- can't stress this enough.

Understanding causality (in your system) is the most straightforward way to generate reliable knowledge (the goal of any research) / And ablation is a very low-effort way to look into causality.

Ablation study

Machine learning에서, ablation study는 machine learning system의 building blocks (dataset, feature, model components) 을 제거해서 전체 성능에 미치는 효과에 대한 insights를 얻기 위한 과학적 실험.

각 model ablation trial은 1개 혹은 그 이상의 components가 제거된 모델을 학습하는 것을 포함한다.

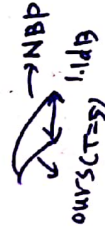
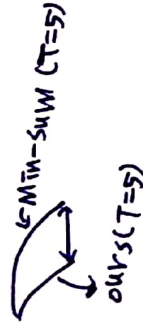
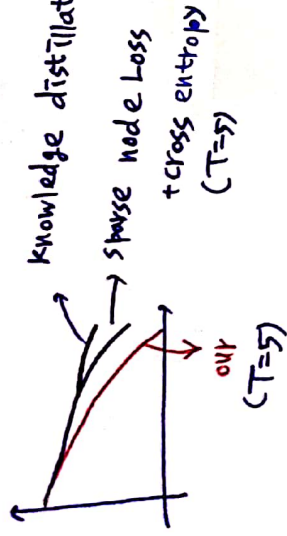
An ablation study typically refers to removing some features of the model or algorithm and seeing how that affects performance.

The teacher min-sum decoder can be regarded as lower bound to Q.

training only with the knowledge distillation loss Q leads to degradation in the low SNR regime (3 ~ 6dB). complete

The degradation is alleviated when training with the cross-entropy loss. • 제안한 요소가 모델에 어떠한 영향을 미치는지 확인하고 싶을 때, 이 요소를 포함한 / 포함하지 않은 모델을 비교하는 것. ⇒ 이는 다양한 연구에서 매우 중요한 의미를 지니는 예, 시노 테스트의 인과관계는 간단히 알다. 분수 있기 때문.

Fig 2-6) BCL (63/15)



NBP = ours & smaller SNR values

① neural decoder that trains only with the knowledge distillation loss (cyan)

② neural decoder that trains with the sparse node activation loss (red)

③ Min-Sum (T=30, teacher)

Magenta Curve: ours (T=5)

⇒ achieves the best results which demonstrates the advantages of our method.

★ Ablation studies are crucial for deep learning research -- can't stress this enough.

Understanding causality (in your system) is the most straightforward way to generate reliable knowledge (the goal of any research) / And ablation is a very low-effort way to look into causality.

Ablation study

Machine learning, ablation study는 machine learning system의 building blocks dataset feature, model components 을 제거해서 전체 성능에 미치는 효과에 대한 insights을 얻기 위한 과학적 실험.

각 model ablation 하기는 1개 또는 그 이상의 components가 제거된 모델을 학습하는 것을 포함한다.

An ablation study typically refers to removing some features of the model of algorithm and seeing how that affects performance.

The teacher min-sum decoder can be regarded as lower bound to Q.

★ training only with the knowledge distillation loss Q leads to degradation in the low SNR regime (3 ~ 6dB).
complete ↓

The degradation is alleviated when training with the cross-entropy loss.

• 제한된 요소가 모델에 어떠한 영향은 미치는지 확인하고 싶을 때, 이 요소들 포함 / 포함하지 않는 모델을 비교하는 것 ⇒ 이는 다양한 연구에서 매우 중요한 의미를 지니는 데, 시스템의 인과관계를 관찰하려면 필수 있기 때문이다.