

Knowledge Distillation in Wide Neural Networks: Risk Bound, Data Efficiency and Imperfect Teacher

Guangda Ji

jiguangda@pku.edu.cn

Zhanxing Zhu

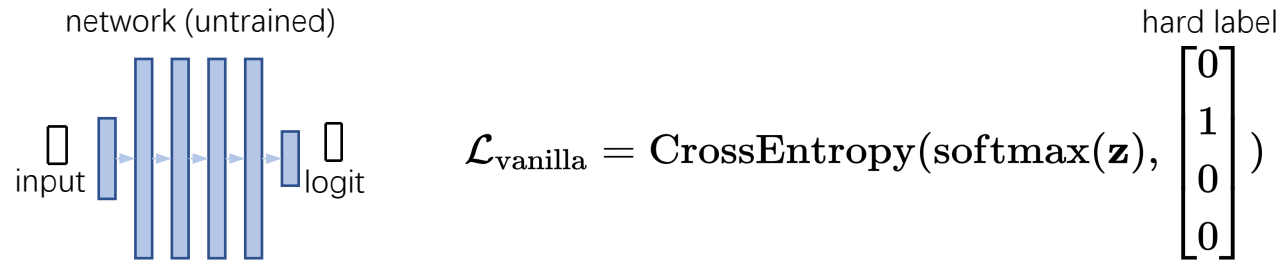
zhanxing.zhu@pku.edu.cn

Peking University

Introduction to Knowledge distillation(KD)

- Vanilla training: only hard (one-hot) labels.

- example:



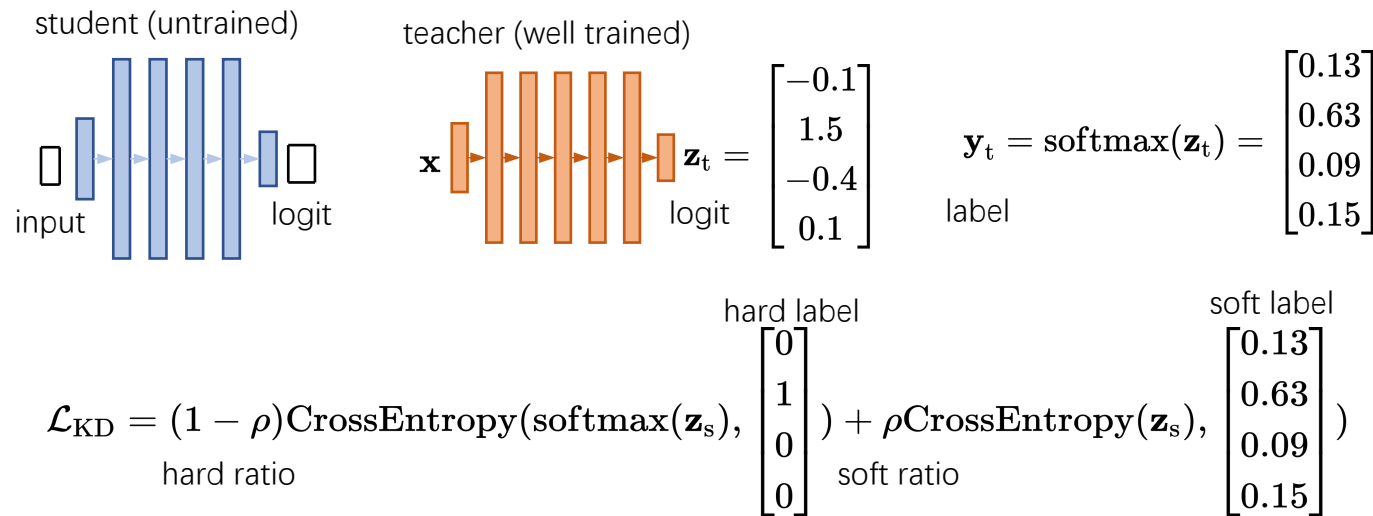
Introduction to Knowledge distillation(KD)

- Vanilla training: only hard (one-hot) labels.

- example:



- Knowledge distillation^[1]: combinations of soft and hard labels. example:



[1] Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network.

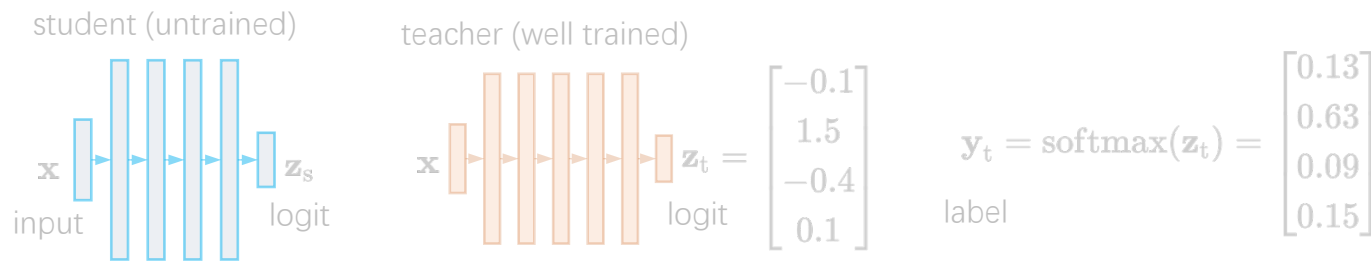
Introduction to Knowledge distillation(KD)

- Vanilla training: only hard (one-hot) labels.

- example:



- Knowledge distillation: combinations of soft and hard labels. example:



- KD is widely used in industry but lacks a satisfying explanation.
- Our work is to
 - establish a theoretical understanding on KD.
 - give instructions on the optimal choice of parameters.

$$\mathcal{L}_{\text{KD}} = (1 - \rho) \text{CrossEntropy}(\text{softmax}(\mathbf{z}_s), \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}) + \rho \text{CrossEntropy}(\mathbf{z}_s, \begin{bmatrix} 0.13 \\ 0.63 \\ 0.09 \\ 0.15 \end{bmatrix})$$

hard ratio hard label soft ratio soft label

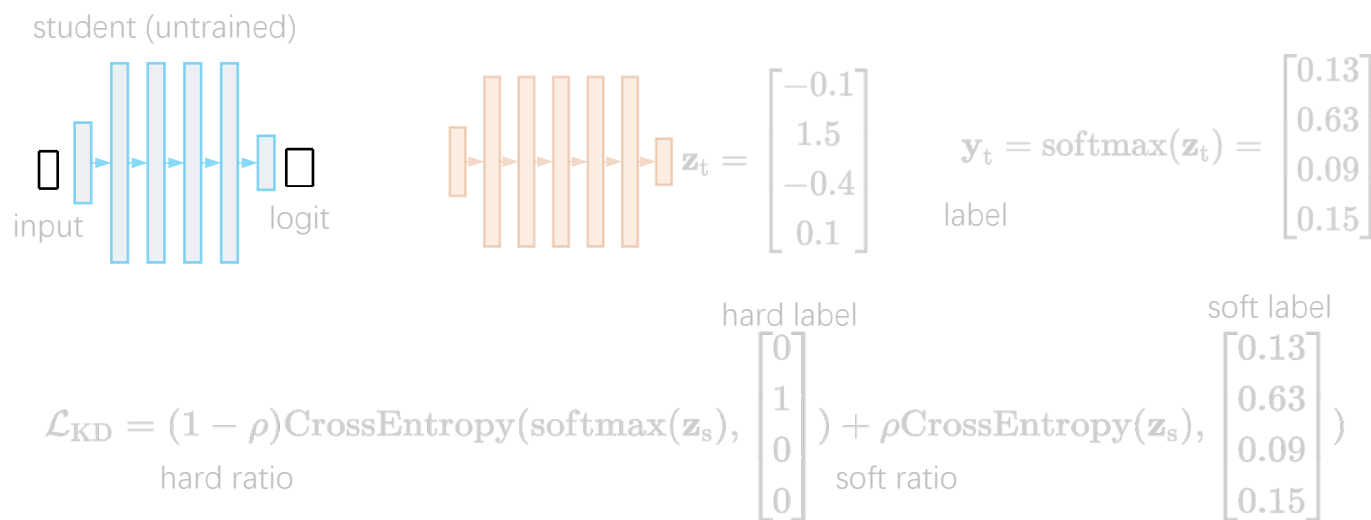
Introduction to Knowledge distillation(KD)

- Vanilla training: only hard (one-hot) labels.

- example:



- Knowledge distillation: combinations of soft and hard labels. example:



- KD is widely used in industry but lacks a satisfying explanation.
- Our contribution in this work:
 - Transfer risk bound.
 - Metric of *data inefficiency* for perfect teacher distillation
 - Hard labels in imperfect distillation.

Problem Setup

- Binary classification,
 - logit $z \in \mathbb{R}$, hard label $y \in \{0, 1\}$, soft label $y \in [0, 1]$.
 - distillation loss:

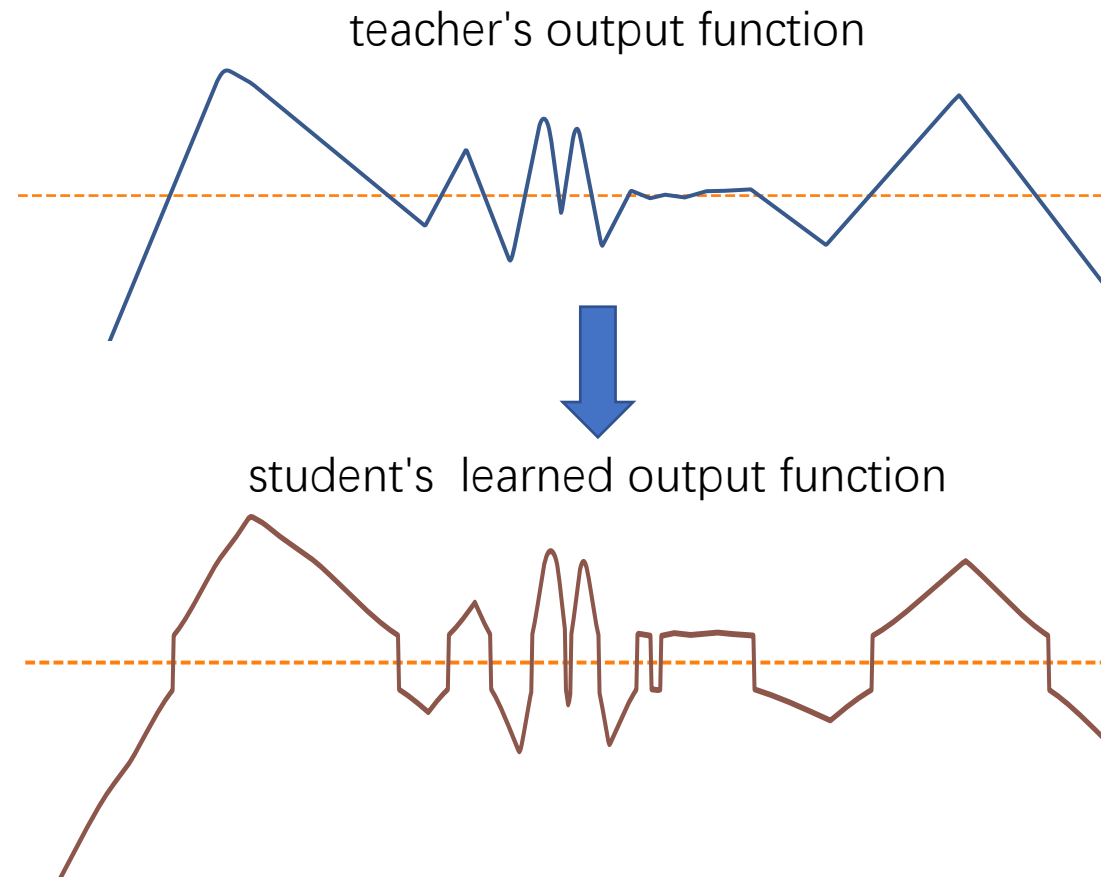
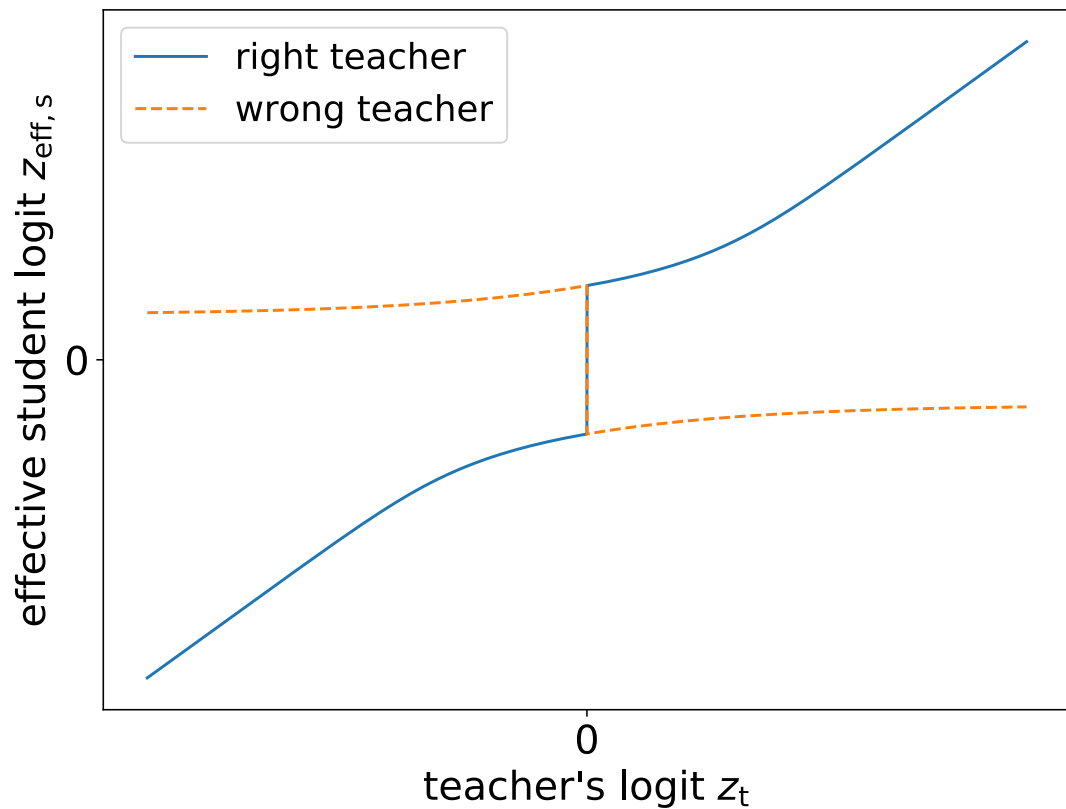
$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \ell_n = \frac{1}{N} \sum_{n=1}^N \underbrace{\rho H\left(\sigma\left(\frac{z_{t,n}}{T}\right), \sigma\left(\frac{z_{s,n}}{T}\right)\right)}_{\text{soft loss}} + \underbrace{(1 - \rho) H(y_{g,n}, \sigma(z_{s,n}))}_{\text{hard loss}} \quad \rho: \text{soft ratio}$$

- Assumption: the student network is wide and over-parameterized.
 - convergence to global minima^[2].
 - minimization at each sample point,

$$\lim_{\tau \rightarrow \infty} z_s(\tau) = \hat{z}_s, \quad \frac{d\ell}{d\hat{z}_s} = \frac{\rho}{T} (\sigma(\hat{z}_s/T) - \sigma(z_t/T)) + (1 - \rho) (\sigma(\hat{z}_s) - y_g) = 0.$$

- Through this equation, student learns an effective logit. $z_{s,\text{eff}}(z_t, y_g) = \hat{z}_s$.
- Result in thresholds and discontinuities in the output function.

Problem Setup



- Through this equation, student learns an effective logit $z_{s,\text{eff}}(z_t, \mathbf{y}_g) = \hat{z}_s$.
- Result in thresholds and discontinuities in the output function.

Problem Setup

- Assumption: the student network is wide and over-parameterized.
 - The use of neural tangent kernel (NTK) technique^[3].
 - Approximate student network's output with its linearized version,

$$f(\mathbf{x}; w_{\text{nlín}}) \approx f(\mathbf{x}; w_0) + \Delta_w^\top \phi(\mathbf{x})$$

- where
 - $\Delta_w = w - w_0 \in \mathbb{R}^p$: weight change,
 - $\phi(\mathbf{x}) = \partial_w f(\mathbf{x}; w_0) \in \mathbb{R}^p$: random feature.
- Also enable us to establish a direct link between network's weight change and its logits,

$$\Delta_{\hat{w}} = \phi(\mathbf{X})(\hat{\Theta}(\mathbf{X}, \mathbf{X}))^{-1} \Delta_{\mathbf{z}},$$

- where
 - $\hat{\Theta}(\mathbf{X}, \mathbf{X}) = \phi(\mathbf{X})^\top \phi(\mathbf{X})$: tangent kernel
 - $\Delta_{\mathbf{z}} = \mathbf{z} - f(\mathbf{X}; w_0)$.

Result 1: Transfer Risk Bound

- Transfer risk \mathcal{R} : probability of different prediction w.r.t. teacher.
- **Theorem 1 (Risk bound):**

$$\mathcal{R}_n \leq p\left(\frac{\pi}{2} - \bar{\alpha}_n\right),$$

- n : sample size,
- $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$: angle between oracle weight change Δ_{w_*} and student's weight change $\Delta_{\hat{w}}$,
- $p(\beta)$: pdf of angle between random feature $\phi(x)$ and oracle Δ_{w_*} .
- Tighter than the bound in [4].

Result 1: Transfer Risk Bound

- Transfer risk \mathcal{R} : probability of different prediction w.r.t. teacher.

- **Theorem 1 (Risk bound):**

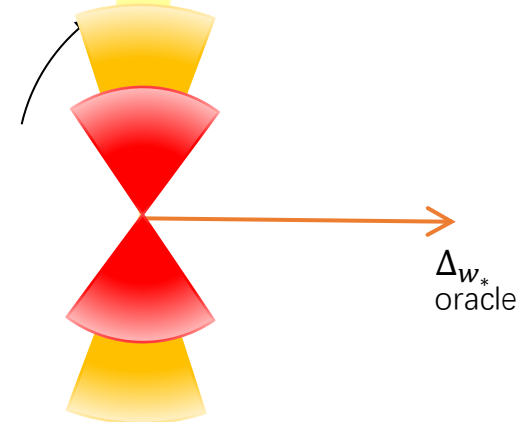
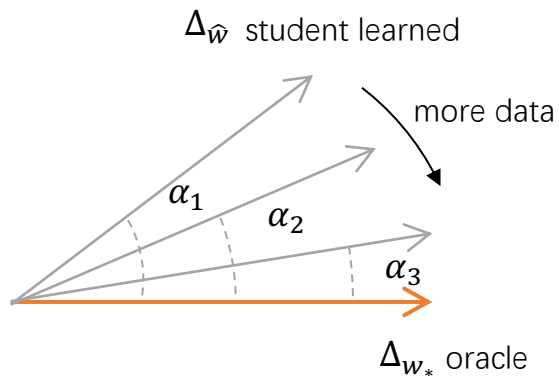
$$\mathcal{R}_n \leq p\left(\frac{\pi}{2} - \bar{\alpha}_n\right),$$

- n : sample size,
- $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$: angle between oracle weight change Δ_{w_*} and student's weight change $\Delta_{\hat{w}}$,
- $p(\beta)$: pdf of angle between random feature $\phi(x)$ and oracle Δ_{w_*} .
- Tighter than the bound in [4].

- **Key idea of this theorem:**

- student learns a projection of oracle weight change

$$\Delta_{\hat{w}} = \phi(\mathbf{X})\Theta_n^{-1}\phi(\mathbf{X})^\top \Delta_{w_*} = \mathbf{P}_\Phi \Delta_{w_*}$$



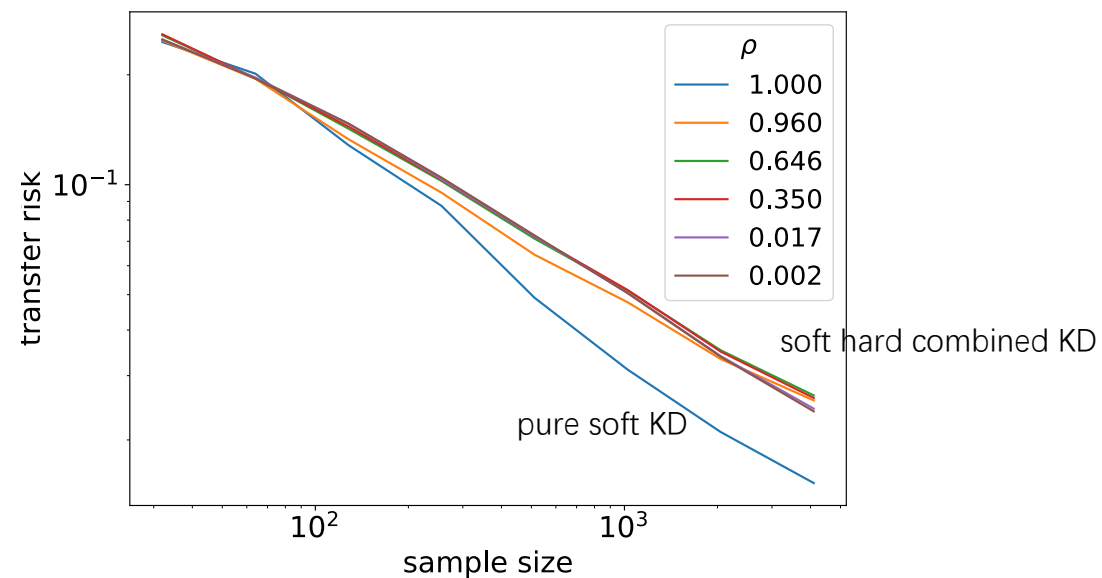
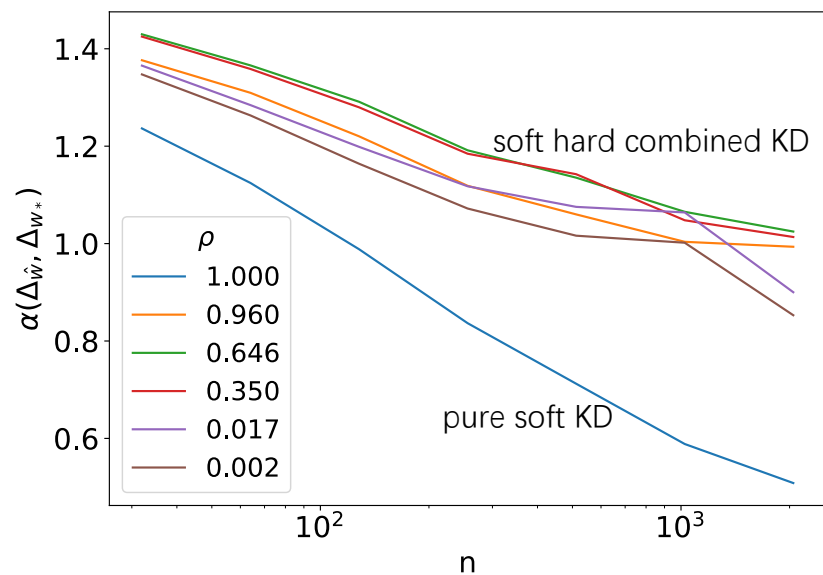
- $\bar{\alpha}_n$ decreases when more data is used.
- Error region ($\color{red}{\blacksquare}$) in random feature space decreases

Result 1: Transfer Risk Bound

- Theorem 1 (Risk bound):

$$\mathcal{R}_n \leq p\left(\frac{\pi}{2} - \bar{\alpha}_n\right),$$

- n : sample size,
 - $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$: angle between oracle weight change Δ_{w_*} and student's weight change $\Delta_{\hat{w}}$,
 - $p(\beta)$: pdf of angle between random feature $\phi(x)$ and oracle Δ_{w_*} .
- Good generalization happens when smaller angle is achieved with same amount of data.**
 - faster angle converging speed for pure soft distillation.
 - explains the fast converging error for pure soft distillation.



Result 2 & 3: Data Inefficiency and Imperfect KD

- **Definition (Data inefficiency):** the increasing speed of the norm of weight change,

$$\mathcal{I}(n) = n [\ln \mathbb{E} \|\Delta_{\hat{w}, n+1}\|_2 - \ln \mathbb{E} \|\Delta_{\hat{w}, n}\|_2] \approx \frac{\partial \ln \|\Delta_{\hat{w}, n}\|_2}{\partial \ln n}$$

- characterizes the converging speed of angle $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$.
- more inefficient \Rightarrow task is harder to train.

Result 2 & 3: Data Inefficiency and Imperfect KD

- **Definition (Data inefficiency):** the increasing speed of the norm of weight change,

$$\mathcal{I}(n) = n [\ln \mathbb{E} \|\Delta_{\hat{w}, n+1}\|_2 - \ln \mathbb{E} \|\Delta_{\hat{w}, n}\|_2] \approx \frac{\partial \ln \|\Delta_{\hat{w}, n}\|_2}{\partial \ln n}$$

- characterizes the converging speed of angle $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$.
 - more inefficient \Rightarrow task is harder to train.
- Two factors that reduces data inefficiency
 1. early stopping epoch of teacher network,
 2. higher soft ratio ρ .
 - both may have a smoothing effect on student output function.

Result 2 & 3: Data Inefficiency and Imperfect KD

- **Definition (Data inefficiency):** the increasing speed of the norm of weight change,

$$\mathcal{I}(n) = n [\ln \mathbb{E} \|\Delta_{\hat{w}, n+1}\|_2 - \ln \mathbb{E} \|\Delta_{\hat{w}, n}\|_2] \approx \frac{\partial \ln \|\Delta_{\hat{w}, n}\|_2}{\partial \ln n}$$

- characterizes the converging speed of angle $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$.
- more inefficient \Rightarrow task is harder to train.
- Two factors that reduces data inefficiency
 1. early stopping epoch of teacher network,
 2. higher soft ratio ρ .
 - both may have a smoothing effect on student output function.
- **Imperfect teacher**
 - Our risk bound and data inefficiency assumes teacher is 100% accurate (perfect)
 - results in a favor in pure soft distillation
 - If teacher have a chance of mistake,
 - hard labels can partially correct the sign of student logits
 - hard labels can reduce $\bar{\alpha}_n = \bar{\alpha}(\Delta_{w_*}, \Delta_{\hat{w}})$

Conclusion

- Transfer risk bound under NTK settings.
- Data inefficiency
 - early stopping and higher soft ratio are beneficial for perfect distillation.
- Hard labels are need in imperfect distillation as a trade-off against teacher's mistake.