

## 392 A Proof of Proposition 2.2: additive expansion proposition

393 We denote the embedding vector of a node  $v_i$  by  $e_i \triangleq g_i(\hat{G}) = \text{GNN}(\hat{G})[i]$ . Without loss of  
 394 generality, we drop the subscript for short. We can also define the density measure  $d_f$  as:

$$p_f(e) \triangleq \frac{\exp(-d(e))}{\int_{\mathcal{G}} \exp(-d(e)) dg(\mathcal{G})}. \quad (5)$$

395 For any subset of the embedding space  $E \subset g(\mathcal{X})$ , the local probability can be measured by  
 396  $p_f(E) = \int_E p_f(e) de$ . If we have a local optimal subset  $U \subset E$  with a confidence threshold  
 397 of  $1 - q$ , and its perturbation  $U_\epsilon$ , then the consistency at the boundary of the subset separation  
 398 problem  $E \rightarrow U, E \setminus U$  can be quantified by Cheeger constant. We introduce the continuous Cheeger  
 399 inequality to elaborate the lower bound of the Cheeger constant under an amplified measure  $\alpha f$ ,  
 400 where  $\alpha > 1$  is a constant.

401 We first define the restricted *Cheeger constant* in the link prediction task. Given the function  $f$  and  
 402 any subset  $E$ , Cheeger constant is calculated by

$$\mathcal{C}_f(E) \triangleq \lim_{\epsilon \rightarrow 0^+} \inf_{A \subset E} \frac{p_f(A_\epsilon) - p_f(A)}{\epsilon \min\{p_f(A), p_f(E \setminus A)\}}. \quad (6)$$

403 According to the definition, the Cheeger constant is a lower bound of probability density in the  
 404 neighborhood of the given set. It quantifies the chance of escaping the subset  $A$  under the probability  
 405 measure  $f$  and reveals the consistency over the set cutting boundary.

406 Then we prove that for the any subset  $E \subset g(\mathcal{G})$  with its local optimal subset  $U : \{e \in E : p_f(e) >$   
 407  $1 - q\}$ , there exists  $\alpha > 1$  s.t.  $\mathcal{C}_{\alpha f}(E \setminus U) \geq 1$ .

408 As the measurable function for the link prediction is defined as  $f(e) = -e^T e_a$ . When  $e^* = e_a$ ,  $f$   
 409 reaches the global minimal. For the embedding vectors outside the local minimal subset  $e_y \in g(\mathcal{G}) \setminus U$ ,  
 410 there exists  $\epsilon > 0$  s.t.

$$f(e_y) \geq f(e^*) + 2\hat{\epsilon}, \quad (7)$$

411 where  $\hat{\epsilon} = C\epsilon$ . If we define  $E_\epsilon^* = \{e_\epsilon^*\} \cap g(\mathcal{G})$ , where  $e_\epsilon^*$  is the  $\hat{\epsilon}$  neighbor of  $e^*$ , according to the  
 412 Lipchitz condition of  $f$ , for  $e \in E_\epsilon^*$ , we have:

$$f(e_x) \leq f(e^*) + \hat{\epsilon} \|e_x - e^*\|_2 \leq f(e^*) + \hat{\epsilon}. \quad (8)$$

413 Combining Eq.8 and Eq.7 leads to  $f(e_y) - f(e_x) \geq \hat{\epsilon}$ . Thus, for the amplified probability measure  
 414  $p_{\alpha f}$ , we have

$$p_{\alpha f}(e_x) / p_{\alpha f}(e_y) \geq \exp(\alpha \hat{\epsilon}) \quad (9)$$

415 According to the inequality property from [26] (formula 63), we have

$$\frac{p_{\alpha f}(U)}{p_{\alpha f}(g(\mathcal{G}) \setminus U)} \geq \exp(\alpha \hat{\epsilon} - 2 \log(2C^2 / \hat{\epsilon})). \quad (10)$$

416 As  $p_{\alpha f}(g(\mathcal{G}) \setminus U) + p_{\alpha f}(U) = 1$ . If we select  $\alpha$  large enough s.t. the RHS of Eq.10 is larger than  
 417 1,  $p_{\alpha f}(g(\mathcal{G}) \setminus U) \leq \frac{1}{2}$ . Thus, according to [17] (Theorem 2.6), we have  $\mathcal{C}_{\alpha f} \geq 1$ . It guarantees  
 418 consistency around the perimeter of the  $U$ . As  $\alpha > 1$  and  $p_f, p_{\alpha f}$  are bounded on any subsets of  
 419 embedding space. It implies a probability margin  $\eta > 0$  at the neighborhood of the local optimal  
 420 between two measurable functions  $f, \alpha f$ , where

$$\eta = \inf_{\hat{e} \in U_\epsilon \setminus U, e \in U} (p_{\alpha f}(\hat{e}) - p_f(e)). \quad (11)$$

421 which, according to [23], implies additive expansion property of the probability measure in the link  
 422 prediction, as Proposition 2.2.

423 **B Proof of Theorem 2.3: error analysis**

424 In [23],  $\mathcal{M}(g_\phi)$  is also assumed to satisfy additive-expansion  $(q, \epsilon)$ , where  $\mathcal{M}(g) \triangleq \{y \in Y :$   
 425  $g(y) \neq y\}$  is the set of mis-classified samples, and they give the error bound of the trained classifier  
 426  $s$  (Theorem B.2):

$$\text{Err}(g) \leq 2(q + \mathcal{A}(g)). \quad (12)$$

427 Here in link prediction task,  $\mathcal{M}(g_\phi)$  is mis-classified samples by the pseudo labeler (teacher model).  
 428 It can be written by  $\{y_i : Y_p[i] \neq \mathcal{E}_T[i]\}$ , which is intractable during the training. The probability  
 429 threshold is  $1 - q$  and a local optimal subset  $U$  for PL is constructed accordingly. We aim to let  
 430  $\mathcal{M}(g_\phi) \cap U$  be close to  $\emptyset$ , so that  $g(\mathcal{G}) \setminus U$  can cover  $\mathcal{M}(g_\phi)$  as much as possible. So we define the  
 431 robust set  $\mathcal{S}(g)$  as

$$\mathcal{S}(g) = \{y : g(y) = g(\hat{y}), \hat{y} \in \{y_\epsilon\}\}, \quad (13)$$

432 where  $y_\epsilon$  is the  $\epsilon$  neighborhood of sample  $y$ . Then, according to Proposition 2.1, we have:

$$p_f(\{y \in Y : g_\phi(y) \neq y, y \in \mathcal{S}(g_\psi)\}) \leq p_{\alpha f}(g(\mathcal{G}) \setminus U) \leq q, \quad (14)$$

433 which has similar form with [23] Lemma B.3 for link prediction task. Besides, the analysis of  
 434  $p_f(\{y \in Y : g_\phi(y) = y, g_\psi(y) \neq y, y \in \mathcal{S}(g_\psi)\})$  and  $p_f(\overline{\mathcal{S}(g_\psi)})$  are the same. Thus, the  
 435 assumption on  $\mathcal{M}(g_\phi)$  is satisfied. Then, we can draw the same conclusion with Eq.12, and the  
 436 classifier is the student model  $g_\psi$ . The theorem is proved.

437 **C Proof of convergence inequality**

438 The PL strategy  $\mathcal{T}$  for the unlabeled data provides a Bayesian prior, from which we formalize the  
 439 empirical loss defined in Eq.1 as

$$\mathcal{L}_T^{(t+1)} = \frac{1}{|\hat{Y}_o^{(t)}| + k} \left[ \text{CE}(g_\psi^{(t)}, \hat{Y}_o^{(t)}) + \text{CE}(g_\psi^{(t)}, Y_p^{(t)}) \right]. \quad (15)$$

440 We can decompose the cross-entropy loss of the pseudo labeled samples by:

$$\begin{aligned} \text{CE}(g_\psi, Y_p) &= \sum_{\hat{Y}_u} \text{ce}(g_\psi, Y) \cdot \mathcal{T} \\ &= \sum_{\hat{Y}_u} [\text{ce}(g_\psi, Y) - Y [\text{ce}(g_\psi, Y)]] \cdot [\mathcal{T} - Y\mathcal{T}] \\ &\quad + Y\mathcal{T} \sum_{\hat{Y}_u} \text{ce}(g_\psi, Y) + \mathbb{E}[\text{ce}(g_\psi, Y)] \sum_{\hat{Y}_u} \mathcal{T} - \left| \hat{Y}_u \right| Y\mathcal{T}Y [\text{ce}(g_\psi, Y)] \end{aligned} \quad (16)$$

441 .

442 Thus, Eq.16 can be simplified to:

$$\begin{aligned} \text{CE}(g_\psi, Y_p) &= \left| \hat{Y}_u \right| \text{Cov}[\text{ce}(g_\psi, Y), \mathcal{T}] + \mathbb{E}\mathcal{T} \cdot \left| \hat{Y}_u \right| \mathbb{E}[\text{ce}(g_\psi, Y)] \\ &\quad + \mathbb{E}[\text{ce}(g_\psi, Y)] \cdot \left| \hat{Y}_u \right| \mathbb{E}\mathcal{T} - \left| \hat{Y}_u \right| \mathbb{E}\mathcal{T} \mathbb{E}[\text{ce}(g_\psi, Y)] \\ &= \left| \hat{Y}_u \right| \text{Cov}[\text{ce}(g_\psi, Y), \mathcal{T}] + \left| \hat{Y}_u \right| \mathbb{E}\mathcal{T} \mathbb{E}[\text{ce}(g_\psi, Y)] \\ &= \left| \hat{Y}_u \right| \text{Cov}[\text{ce}(g_\psi, Y), \mathcal{T}] + k \mathbb{E}[\text{ce}(g_\psi, Y)] \end{aligned} \quad (17)$$

443 Note that  $\mathcal{T}$  is the indicator-like function, where we have

$$\mathbb{E}\mathcal{T} = \frac{1}{|\hat{Y}_u|} \sum_{\hat{Y}_u} \mathcal{T} = \frac{k}{|\hat{Y}_u|}. \quad (18)$$

Table 6: Details of Node Information in the Case Study.

Node	Group 1			Group 2	
	Node 2702	Node 5688	Node 8906	Node 3489	Node 7680
<b>ID</b>	17505908	11353631	30138652	23221074	12265137
<b>Outlinks</b>	[6097297]	[6097297]	[244374, 6097297]	[20901]	[]
<b>Title</b>	Ubuntu Hacks	Pungi (software)	LinuxPAE64	Malware Bell	Norton Confidential
<b>Label</b>	Operating systems	Operating systems	Operating systems	Computer security	Computer security
<b>Tokens</b>	"ubuntu", "hacks", "tips", "tools", "exploring", "using", "tuning", "linux", "book", "tips", "ubuntu", "popular", "linux", "distribution", "book", "published", "o'reilly", "media", "june", "2006", "part", "o'reilly", "hacks", "series"	"pungi", "software", "linux", "software", "making", "fedora", "release", "7", "updates"	"linuxpae64", "linuxpae64", "port", "linux", "kernel", "running", "compatibility", "mode", "x86-64", "processor", "kernel", "capable", "loading", "i386", "modules", "device", "drivers", "supports", "64-bit", "linux", "applications", "user", "mode"	"malware", "bell", "malware", "bell", "malware", "program", "made", "taiwan", "somewhere", "2006", "2007", "malware", "bell", "tries", "install", "automatically", "upon", "visiting", "website", "promoting", "containing", "malware"	"norton", "confidential", "norton", "confidential", "program", "designed", "encrypt", "passwords", "online", "detect", "phishing", "sites"

444 Based on the Eq.15 and Eq.17, we can rewrite  $\mathcal{L}_{\mathcal{T}}^{(t+1)}$  as

$$\begin{aligned}
 \mathcal{L}_{\mathcal{T}}^{(t+1)} &= \beta \text{Cov} [\text{ce} (g_{\psi}, Y), \mathcal{T}] + \frac{1}{|\hat{Y}_o^{(t)}|} \text{CE} \left( g_{\psi}^{(t)}, \hat{Y}_o^{(t)} \right) \\
 &\leq \beta \text{Cov} [\text{ce} (g_{\psi}, Y), \mathcal{T}] + \frac{1}{|\hat{Y}_o^{(t)}|} \text{CE} \left( g_{\phi}^{(t)}, \hat{Y}_o^{(t)} \right) \\
 &= \beta \text{Cov} [\text{ce} (g_{\psi}, Y), \mathcal{T}] + \mathcal{L}_{\mathcal{T}}^{(t)}
 \end{aligned} \tag{19}$$

445 where  $\beta = |\hat{Y}_u| / (|\hat{Y}_o| + k)$ . The inequality holds due to the assumption.

## 446 D Case study of CPL on link prediction

447 **Error bound:** In the case study, the recorded confidence threshold is  $1 - q = 0.98$  for WikiCS. We  
 448 adopt 5 views of dropout with the augmentation drop rate 0.05. And according to the error bound given  
 449 by Theorem 2.3, given the confidence threshold, Eq.2 suggests that the higher prediction consistency  
 450 should lead to a smaller error bound. The final prediction consistency is  $\mathcal{A}(g) = 0.0358$ , thus, we  
 451 can calculate error bound  $Err(g) = 0.1116$ . The AUC and AP are  $95.56 \pm 0.24\%$ ,  $95.58 \pm 0.29\%$   
 452 which are bounded within  $Err(g)$ .

453 **Knowledge discovery:** In the 5 random experiments, we add 500 pseudo links in each iteration.  
 454 Here we focus on the common PL links in the first iteration, which are considered the most confident  
 455 samples. We look for the metadata of WikiCS whose node, feature, link and node label represent  
 456 paper, token, reference relation and topic of the paper respectively. There are These 7 most confident  
 457 links categorized into 2 groups. We take 3 out of 5 nodes in group1 and the 2 nodes in group2 for  
 458 analysis, whose detailed information of these nodes is shown in AppendixD.

459 For group1, 3 nodes are connected by the pseudo links, and they are all linked to a central node  
 460 whose degree is 321. The metadata information of the nodes are all strongly relevant to "Linux"  
 461 in the "operating systems" topic. Thus, the PL linked nodes are likely to have common neighbors  
 462 discovered triangle relationship. In group2, node 3489 has no in/out degree and is pseudo linked to  
 463 node 7680. Both papers focus on the "malware"/"phishing" under the topic "Computer security".  
 464 Although they only have one common token, the CPL strategy successfully discovers the correlation  
 465 and consistently add it to the training set. The detailed result of the case study is shown in Table 6.