
Partial Multi-Label Learning with Probabilistic Graphical Disambiguation

Supplementary Material

Jun-Yi Hang, Min-Ling Zhang*

School of Computer Science and Engineering, Southeast University, Nanjing 210096, China
 Key Laboratory of Computer Network and Information Integration (Southeast University),
 Ministry of Education, China
 {hangjy, zhangm1}@seu.edu.cn

Appendix A

More Experimental Results for Comparative Studies

Table A.1, A.2 and A.3 report detailed experimental results in terms of *Coverage*, *One-error* and *Hamming loss*, which are not covered in the *Comparative Studies* part of the main body due to page limit. It can be observed that our PARD achieves consistently superior performance to well-established PML approaches.

Table A.1: Predictive performance of each comparing approach (mean \pm std. deviation) in terms of *Coverage*, where \bullet/\circ indicates whether PARD is significantly superior/inferior to one comparing approach via paired *t*-test at 0.05 significance level. $\uparrow(\downarrow)$ indicates the larger (smaller) the value, the better the performance. Best results are shown in boldface.

Data sets	$\gamma\%$	Coverage ↓						
		FPMI	PARVLS	PML-NI	PML-MD	UPML-HL	UPML-RL	PARD
YeastBP	0.437 \pm 0.020 \bullet	0.537 \pm 0.024 \bullet	0.265 \pm 0.019 \bullet	0.306 \pm 0.017 \bullet	0.301 \pm 0.021 \bullet	0.394 \pm 0.018 \bullet	0.235\pm0.014	
YeastCC	0.173 \pm 0.014 \bullet	0.256 \pm 0.020 \bullet	0.089 \pm 0.009 \bullet	0.112 \pm 0.012 \bullet	0.096 \pm 0.012 \bullet	0.127 \pm 0.012 \bullet	0.077\pm0.009	
YeastMF	0.185 \pm 0.013 \bullet	0.247 \pm 0.015 \bullet	0.109 \pm 0.012 \bullet	0.124 \pm 0.013 \bullet	0.116 \pm 0.016 \bullet	0.144 \pm 0.011 \bullet	0.094\pm0.010	
Music_emotion	0.433 \pm 0.009 \bullet	0.412 \pm 0.006 \bullet	0.410 \pm 0.007 \bullet	0.399 \pm 0.009 \bullet	0.390\pm0.007\bullet	0.406 \pm 0.009 \bullet	0.396 \pm 0.008 \bullet	
Music_style	0.218 \pm 0.010 \bullet	0.209 \pm 0.009 \bullet	0.198 \pm 0.009 \bullet	0.203 \pm 0.009 \bullet	0.195\pm0.008\circ	0.207 \pm 0.008 \bullet	0.199 \pm 0.011 \bullet	
corel5k	100	0.343 \pm 0.014 \bullet	0.449 \pm 0.024 \bullet	0.372 \pm 0.015 \bullet	0.340 \pm 0.011 \bullet	0.324 \pm 0.016 \bullet	0.314 \pm 0.017 \bullet	0.308\pm0.016
	150	0.349 \pm 0.016 \bullet	0.414 \pm 0.016 \bullet	0.388 \pm 0.015 \bullet	0.347 \pm 0.012 \bullet	0.328 \pm 0.017 \bullet	0.322 \pm 0.016 \bullet	0.321\pm0.014
	200	0.351 \pm 0.014 \bullet	0.423 \pm 0.027 \bullet	0.404 \pm 0.014 \bullet	0.352 \pm 0.012 \bullet	0.345 \pm 0.012 \bullet	0.329\pm0.015	0.332 \pm 0.015 \bullet
	250	0.356 \pm 0.016 \bullet	0.432 \pm 0.027 \bullet	0.414 \pm 0.017 \bullet	0.355 \pm 0.017 \bullet	0.349 \pm 0.015 \bullet	0.332\pm0.016	0.337 \pm 0.015 \bullet
rcv1-s1	100	0.164 \pm 0.006 \bullet	0.304 \pm 0.013 \bullet	0.173 \pm 0.010 \bullet	0.151 \pm 0.007 \bullet	0.146 \pm 0.006 \bullet	0.145 \pm 0.004 \bullet	0.135\pm0.006
	150	0.165 \pm 0.007 \bullet	0.330 \pm 0.024 \bullet	0.197 \pm 0.009 \bullet	0.165 \pm 0.006 \bullet	0.150 \pm 0.006 \bullet	0.149 \pm 0.005 \bullet	0.140\pm0.006
	200	0.172 \pm 0.008 \bullet	0.350 \pm 0.023 \bullet	0.214 \pm 0.008 \bullet	0.173 \pm 0.006 \bullet	0.156 \pm 0.009 \bullet	0.155 \pm 0.005 \bullet	0.148\pm0.005
	250	0.183 \pm 0.008 \bullet	0.364 \pm 0.019 \bullet	0.232 \pm 0.008 \bullet	0.178 \pm 0.007 \bullet	0.166 \pm 0.010 \bullet	0.163 \pm 0.005 \bullet	0.156\pm0.006
Corel16k-s1	100	0.316 \pm 0.007 \bullet	0.391 \pm 0.008 \bullet	0.346 \pm 0.012 \bullet	0.317 \pm 0.010 \bullet	0.304 \pm 0.008 \bullet	0.308 \pm 0.007 \bullet	0.298\pm0.009
	150	0.324 \pm 0.007 \bullet	0.399 \pm 0.011 \bullet	0.362 \pm 0.011 \bullet	0.325 \pm 0.009 \bullet	0.321 \pm 0.008 \bullet	0.314 \pm 0.008 \bullet	0.307\pm0.010
	200	0.336 \pm 0.008 \bullet	0.427 \pm 0.010 \bullet	0.375 \pm 0.013 \bullet	0.330 \pm 0.010 \bullet	0.348 \pm 0.010 \bullet	0.318\pm0.007	0.319 \pm 0.010 \bullet
	250	0.347 \pm 0.009 \bullet	0.470 \pm 0.008 \bullet	0.387 \pm 0.013 \bullet	0.338 \pm 0.008 \bullet	0.349 \pm 0.012 \bullet	0.325 \pm 0.007 \bullet	0.322\pm0.010
iaprtc12	100	0.328 \pm 0.005 \bullet	0.352 \pm 0.005 \bullet	0.354 \pm 0.007 \bullet	0.334 \pm 0.005 \bullet	0.345 \pm 0.008 \bullet	0.349 \pm 0.004 \bullet	0.321\pm0.006
	150	0.332 \pm 0.005 \bullet	0.359 \pm 0.006 \bullet	0.374 \pm 0.007 \bullet	0.346 \pm 0.006 \bullet	0.349 \pm 0.009 \bullet	0.351 \pm 0.004 \bullet	0.327\pm0.006
	200	0.339 \pm 0.005 \bullet	0.376 \pm 0.005 \bullet	0.393 \pm 0.008 \bullet	0.358 \pm 0.007 \bullet	0.351 \pm 0.009 \bullet	0.357 \pm 0.006 \bullet	0.337\pm0.007
	250	0.350 \pm 0.007 \bullet	0.392 \pm 0.006 \bullet	0.413 \pm 0.008 \bullet	0.367 \pm 0.007 \bullet	0.368 \pm 0.010 \bullet	0.363 \pm 0.007 \bullet	0.346\pm0.006
espgame	100	0.365 \pm 0.007 \bullet	0.398 \pm 0.007 \bullet	0.404 \pm 0.007 \bullet	0.371 \pm 0.007 \bullet	0.372 \pm 0.008 \bullet	0.391 \pm 0.009 \bullet	0.353\pm0.006
	150	0.369 \pm 0.007 \bullet	0.419 \pm 0.007 \bullet	0.423 \pm 0.007 \bullet	0.380 \pm 0.007 \bullet	0.361 \pm 0.007 \bullet	0.377 \pm 0.007 \bullet	0.358\pm0.006
	200	0.374 \pm 0.007 \bullet	0.434 \pm 0.006 \bullet	0.435 \pm 0.008 \bullet	0.387 \pm 0.007 \bullet	0.398 \pm 0.010 \bullet	0.381 \pm 0.006 \bullet	0.362\pm0.007
	250	0.384 \pm 0.007 \bullet	0.451 \pm 0.008 \bullet	0.449 \pm 0.009 \bullet	0.389 \pm 0.007 \bullet	0.388 \pm 0.009 \bullet	0.388 \pm 0.007 \bullet	0.370\pm0.007

*Corresponding author

Table A.2: Predictive performance of each comparing approach (mean \pm std. deviation) in terms of *One-error*, where \bullet/\circ indicates whether PARD is significantly superior/inferior to one comparing approach via paired *t*-test at 0.05 significance level. $\uparrow(\downarrow)$ indicates the larger (smaller) the value, the better the performance. Best results are shown in boldface.

Data sets	$\gamma\%$	One-error \downarrow						
		FPML	PARVLS	PML-NI	PML-MD	UPML-HL	UPML-RL	PARD
YeastBP	0.833 \pm 0.012 \bullet	0.977 \pm 0.007 \bullet	0.670\pm0.025\circ	0.749 \pm 0.012 \bullet	0.690 \pm 0.021 \bullet	0.856 \pm 0.012 \bullet	0.687 \pm 0.014 \bullet	
YeastCC	0.863 \pm 0.016 \bullet	0.961 \pm 0.007 \bullet	0.768\pm0.023	0.796 \pm 0.018 \bullet	0.772 \pm 0.017 \bullet	0.826 \pm 0.013 \bullet	0.772 \pm 0.020 \bullet	
YeastMF	0.933 \pm 0.011 \bullet	0.983 \pm 0.007 \bullet	0.838 \pm 0.014 \bullet	0.859 \pm 0.014 \bullet	0.846 \pm 0.010 \bullet	0.900 \pm 0.014 \bullet	0.837\pm0.018	
Music_emotion	0.547 \pm 0.022 \bullet	0.515 \pm 0.020 \bullet	0.495 \pm 0.024 \bullet	0.412 \pm 0.023 \bullet	0.407 \pm 0.025 \bullet	0.431 \pm 0.018 \bullet	0.403\pm0.011	
Music_style	0.393 \pm 0.018 \bullet	0.372 \pm 0.020 \bullet	0.352 \pm 0.017 \bullet	0.394 \pm 0.020 \bullet	0.356 \pm 0.017 \bullet	0.405 \pm 0.017 \bullet	0.337\pm0.015	
corel5k	100	0.642 \pm 0.021 \bullet	0.736 \pm 0.028 \bullet	0.642 \pm 0.021 \bullet	0.636 \pm 0.028 \bullet	0.603 \pm 0.020 \bullet	0.610 \pm 0.023 \bullet	0.583\pm0.023
	150	0.652 \pm 0.017 \bullet	0.732 \pm 0.015 \bullet	0.664 \pm 0.022 \bullet	0.652 \pm 0.020 \bullet	0.608 \pm 0.021 \bullet	0.609 \pm 0.021 \bullet	0.589\pm0.016
	200	0.656 \pm 0.020 \bullet	0.747 \pm 0.029 \bullet	0.681 \pm 0.020 \bullet	0.662 \pm 0.019 \bullet	0.607 \pm 0.015 \bullet	0.610 \pm 0.017 \bullet	0.598\pm0.019
	250	0.660 \pm 0.018 \bullet	0.750 \pm 0.027 \bullet	0.697 \pm 0.021 \bullet	0.663 \pm 0.024 \bullet	0.608 \pm 0.018 \bullet	0.609 \pm 0.018 \bullet	0.602\pm0.016
rcv1-s1	100	0.447 \pm 0.017 \bullet	0.602 \pm 0.030 \bullet	0.425 \pm 0.022 \bullet	0.497 \pm 0.019 \bullet	0.413\pm0.025\circ	0.466 \pm 0.011 \bullet	0.430 \pm 0.019 \bullet
	150	0.458 \pm 0.016 \bullet	0.609 \pm 0.028 \bullet	0.449 \pm 0.022 \bullet	0.513 \pm 0.022 \bullet	0.428 \pm 0.024 \bullet	0.464 \pm 0.015 \bullet	0.424\pm0.018
	200	0.463 \pm 0.015 \bullet	0.622 \pm 0.025 \bullet	0.466 \pm 0.017 \bullet	0.524 \pm 0.022 \bullet	0.434 \pm 0.016 \bullet	0.470 \pm 0.014 \bullet	0.433\pm0.017
	250	0.472 \pm 0.020 \bullet	0.623 \pm 0.021 \bullet	0.495 \pm 0.017 \bullet	0.535 \pm 0.023 \bullet	0.440 \pm 0.027 \bullet	0.490 \pm 0.022 \bullet	0.435\pm0.020
Core16k-s1	100	0.580 \pm 0.015 \bullet	0.716 \pm 0.017 \bullet	0.587 \pm 0.015 \bullet	0.607 \pm 0.017 \bullet	0.577 \pm 0.020 \bullet	0.602 \pm 0.009 \bullet	0.567\pm0.020
	150	0.589 \pm 0.013 \bullet	0.715 \pm 0.017 \bullet	0.610 \pm 0.021 \bullet	0.617 \pm 0.017 \bullet	0.582 \pm 0.018 \bullet	0.605 \pm 0.016 \bullet	0.572\pm0.023
	200	0.590 \pm 0.014 \bullet	0.720 \pm 0.015 \bullet	0.624 \pm 0.021 \bullet	0.624 \pm 0.017 \bullet	0.604 \pm 0.015 \bullet	0.609 \pm 0.013 \bullet	0.580\pm0.024
	250	0.600 \pm 0.016 \bullet	0.721 \pm 0.014 \bullet	0.635 \pm 0.019 \bullet	0.631 \pm 0.019 \bullet	0.603 \pm 0.018 \bullet	0.611 \pm 0.013 \bullet	0.593\pm0.022
iaprtc12	100	0.458 \pm 0.009 \bullet	0.462 \pm 0.015 \bullet	0.448 \pm 0.015 \bullet	0.461 \pm 0.014 \bullet	0.444 \pm 0.013 \bullet	0.525 \pm 0.016 \bullet	0.433\pm0.020
	150	0.460 \pm 0.009 \bullet	0.467 \pm 0.014 \bullet	0.467 \pm 0.017 \bullet	0.481 \pm 0.016 \bullet	0.449 \pm 0.009 \bullet	0.531 \pm 0.010 \bullet	0.438\pm0.013
	200	0.463 \pm 0.012 \bullet	0.474 \pm 0.011 \bullet	0.485 \pm 0.013 \bullet	0.494 \pm 0.013 \bullet	0.444 \pm 0.013 \bullet	0.531 \pm 0.013 \bullet	0.443\pm0.013
	250	0.469 \pm 0.016 \bullet	0.485 \pm 0.011 \bullet	0.507 \pm 0.015 \bullet	0.503 \pm 0.012 \bullet	0.463 \pm 0.014 \bullet	0.540 \pm 0.010 \bullet	0.450\pm0.013
espgame	100	0.608 \pm 0.012 \bullet	0.637 \pm 0.014 \bullet	0.618 \pm 0.011 \bullet	0.618 \pm 0.014 \bullet	0.601 \pm 0.014 \bullet	0.648 \pm 0.012 \bullet	0.589\pm0.014
	150	0.609 \pm 0.010 \bullet	0.644 \pm 0.014 \bullet	0.642 \pm 0.012 \bullet	0.630 \pm 0.013 \bullet	0.596 \pm 0.013 \bullet	0.655 \pm 0.012 \bullet	0.594\pm0.013
	200	0.608 \pm 0.011 \bullet	0.656 \pm 0.015 \bullet	0.654 \pm 0.009 \bullet	0.640 \pm 0.011 \bullet	0.626 \pm 0.016 \bullet	0.664 \pm 0.011 \bullet	0.598\pm0.013
	250	0.611 \pm 0.016 \bullet	0.678 \pm 0.018 \bullet	0.677 \pm 0.009 \bullet	0.652 \pm 0.010 \bullet	0.613 \pm 0.012 \bullet	0.672 \pm 0.009 \bullet	0.604\pm0.012

Table A.3: Predictive performance of each comparing approach (mean \pm std. deviation) in terms of *Hamming loss*, where \bullet/\circ indicates whether PARD is significantly superior/inferior to one comparing approach via paired *t*-test at 0.05 significance level. $\uparrow(\downarrow)$ indicates the larger (smaller) the value, the better the performance. Best results are shown in boldface.

Data sets	$\gamma\%$	Hamming loss \downarrow						
		FPML	PARVLS	PML-NI	PML-MD	UPML-HL	UPML-RL	PARD
YeastBP	0.025 \pm 0.001 \bullet	0.025 \pm 0.001 \bullet	0.026 \pm 0.001 \bullet	0.026 \pm 0.002 \bullet	0.025 \pm 0.001 \bullet	0.026 \pm 0.001 \bullet	0.024 \pm 0.001 \bullet	
YeastCC	0.154 \pm 0.002 \bullet	0.027 \pm 0.002 \bullet	0.029 \pm 0.002 \bullet	0.029 \pm 0.002 \bullet	0.026 \pm 0.002 \bullet	0.027 \pm 0.002 \bullet	0.025\pm0.002	
YeastMF	0.025\pm0.002	0.026 \pm 0.002 \bullet	0.031 \pm 0.002 \bullet	0.028 \pm 0.002 \bullet	0.025\pm0.002	0.026 \pm 0.002 \bullet	0.025\pm0.002	
Music_emotion	0.217 \pm 0.004 \bullet	0.214 \pm 0.006 \bullet	0.212 \pm 0.005 \bullet	0.218 \pm 0.005 \bullet	0.193\pm0.004\circ	0.210 \pm 0.005 \bullet	0.196 \pm 0.006 \bullet	
Music_style	0.123 \pm 0.003 \bullet	0.121 \pm 0.003 \bullet	0.115\pm0.003	0.143 \pm 0.004 \bullet	0.115\pm0.003	0.125 \pm 0.003 \bullet	0.117 \pm 0.006 \bullet	
corel5k	100	0.117 \pm 0.002 \bullet	0.134 \pm 0.006 \bullet	0.129 \pm 0.004 \bullet	0.117 \pm 0.005 \bullet	0.112\pm0.003	0.113 \pm 0.003 \bullet	0.112\pm0.003
	150	0.117 \pm 0.003 \bullet	0.136 \pm 0.006 \bullet	0.130 \pm 0.004 \bullet	0.117 \pm 0.003 \bullet	0.112\pm0.003	0.113 \pm 0.003 \bullet	0.112\pm0.002
	200	0.117 \pm 0.003 \bullet	0.138 \pm 0.005 \bullet	0.131 \pm 0.005 \bullet	0.117 \pm 0.003 \bullet	0.113\pm0.003	0.113\pm0.003	0.113\pm0.003
	250	0.117 \pm 0.003 \bullet	0.140 \pm 0.006 \bullet	0.131 \pm 0.004 \bullet	0.117 \pm 0.004 \bullet	0.113\pm0.003	0.113\pm0.003	0.113\pm0.003
rcv1-s1	100	0.099 \pm 0.002 \bullet	0.118 \pm 0.004 \bullet	0.100 \pm 0.003 \bullet	0.106 \pm 0.003 \bullet	0.095\pm0.002\circ	0.108 \pm 0.002 \bullet	0.098 \pm 0.003 \bullet
	150	0.099 \pm 0.002 \bullet	0.123 \pm 0.004 \bullet	0.103 \pm 0.004 \bullet	0.107 \pm 0.003 \bullet	0.096\pm0.003\circ	0.108 \pm 0.002 \bullet	0.104 \pm 0.003 \bullet
	200	0.099\pm0.002	0.125 \pm 0.003 \bullet	0.106 \pm 0.004 \bullet	0.108 \pm 0.003 \bullet	0.101 \pm 0.003 \bullet	0.107 \pm 0.002 \bullet	0.103 \pm 0.004 \bullet
	250	0.101\pm0.003\circ	0.126 \pm 0.003 \bullet	0.110 \pm 0.005 \bullet	0.109 \pm 0.003 \bullet	0.125 \pm 0.005 \bullet	0.109 \pm 0.002 \bullet	0.108 \pm 0.002 \bullet
Core16k-s1	100	0.121 \pm 0.002 \bullet	0.139 \pm 0.003 \bullet	0.125 \pm 0.002 \bullet	0.117\pm0.002	0.117\pm0.002	0.118 \pm 0.001 \bullet	0.117\pm0.001
	150	0.121 \pm 0.002 \bullet	0.136 \pm 0.002 \bullet	0.126 \pm 0.003 \bullet	0.117\pm0.002	0.117\pm0.001	0.118 \pm 0.001 \bullet	0.117\pm0.001
	200	0.121 \pm 0.002 \bullet	0.135 \pm 0.002 \bullet	0.125 \pm 0.003 \bullet	0.117\pm0.002	0.117\pm0.001	0.118 \pm 0.001 \bullet	0.117\pm0.002
	250	0.120 \pm 0.002 \bullet	0.135 \pm 0.002 \bullet	0.125 \pm 0.003 \bullet	0.117\pm0.002	0.117\pm0.002	0.117\pm0.002	0.117\pm0.001
iaprtc12	100	0.143 \pm 0.002 \bullet	0.140 \pm 0.002 \bullet	0.144 \pm 0.002 \bullet	0.146 \pm 0.002 \bullet	0.138 \pm 0.002 \bullet	0.151 \pm 0.002 \bullet	0.137\pm0.002
	150	0.143 \pm 0.002 \bullet	0.140 \pm 0.002 \bullet	0.146 \pm 0.002 \bullet	0.149 \pm 0.002 \bullet	0.138\pm0.003	0.151 \pm 0.002 \bullet	0.139 \pm 0.002 \bullet
	200	0.143 \pm 0.002 \bullet	0.141 \pm 0.002 \bullet	0.148 \pm 0.002 \bullet	0.150 \pm 0.002 \bullet	0.139\pm0.002\circ	0.151 \pm 0.002 \bullet	0.141 \pm 0.002 \bullet
	250	0.144 \pm 0.002 \bullet	0.143\pm0.001\circ	0.150 \pm 0.003 \bullet	0.151 \pm 0.002 \bullet	0.148 \pm 0.003 \bullet	0.151 \pm 0.002 \bullet	0.146 \pm 0.002 \bullet
espgame	100	0.133 \pm 0.002 \bullet	0.135 \pm 0.002 \bullet	0.141 \pm 0.002 \bullet	0.130 \pm 0.002 \bullet			

Appendix B

Derivation of The Variational Lower Bound

The variational lower bound of the log-likelihood (i.e. Eq. (2) in the main body) is derived as follows

$$\begin{aligned}
\log p_\theta(\mathbf{s}|\mathbf{x}) &= \log \int p_\theta(\mathbf{s}, \mathbf{y}|\mathbf{x}) d\mathbf{y} \\
&= \log \int p_\theta(\mathbf{s}|\mathbf{x}, \mathbf{y}) p_\theta(\mathbf{y}|\mathbf{x}) d\mathbf{y} \\
&= \log \int q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s}) \frac{p_\theta(\mathbf{s}|\mathbf{x}, \mathbf{y}) p_\theta(\mathbf{y}|\mathbf{x})}{q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s})} d\mathbf{y} \\
&\geq \mathbb{E}_{q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s})} [\log \frac{p_\theta(\mathbf{s}|\mathbf{x}, \mathbf{y}) p_\theta(\mathbf{y}|\mathbf{x})}{q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s})}] \\
&= \mathcal{L}(\mathbf{x}, \mathbf{s}; \theta, \phi).
\end{aligned}$$

Appendix C

Derivation of The KL-Divergence Term's Closed-Form Solution

With mean-field approximation technique, closed-form solution of the KL-divergence term can be derived as follows

$$\begin{aligned}
&KL[q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s}) || p_\theta(\mathbf{y}|\mathbf{x})] \\
&= \mathbb{E}_{q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s})} [\log q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s}) - \log p_\theta(\mathbf{y}|\mathbf{x})] \\
&= \mathbb{E}_{q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s})} [\sum_{k=1}^t \log q_\phi(y_k|\mathbf{x}, \mathbf{s}) - \sum_{k=1}^t \log p_\theta(y_k|\mathbf{x})] \\
&= \sum_{k=1}^t \mathbb{E}_{q_\phi(\mathbf{y}|\mathbf{x}, \mathbf{s})} [\log \frac{q_\phi(y_k|\mathbf{x}, \mathbf{s})}{p_\theta(y_k|\mathbf{x})}] \\
&= \sum_{k=1}^t \mathbb{E}_{q_\phi(y_k|\mathbf{x}, \mathbf{s})} [\log \frac{q_\phi(y_k|\mathbf{x}, \mathbf{s})}{p_\theta(y_k|\mathbf{x})}] \\
&= \sum_{k=1}^t KL[q_\phi(y_k|\mathbf{x}, \mathbf{s}) || p_\theta(y_k|\mathbf{x})] \\
&= \sum_{k=1}^t p_\phi^{y_k} \log \frac{p_\phi^{y_k}}{p_\theta^{y_k}} + (1 - p_\phi^{y_k}) \log \frac{1 - p_\phi^{y_k}}{1 - p_\theta^{y_k}}.
\end{aligned}$$

Althouth mean-field approximation would restrict model capacity, it is a routine in VAE-related literatures [17] to make graphical model tractable. A natural direction for future work is to investigate whether it is possible to compute the KL-divergence term without mean-field approximation.