

# Supplemental Material to 'Learning classifier chains using matrix regularization: application to multimorbidity prediction'

Paweł Teisseyre

In the following we give formal definitions of the evaluation measures used in the experiments. In the paper we consider example-based measures which are calculated for each test example and then averaged across the test set. Below we recall definitions of the considered measures (definitions are given for one instance in the test set). Let  $\mathbf{y} = (y_1, \dots, y_K)$  be a vector of true labels and  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_K)$  be the vector of predicted labels for some instance  $\mathbf{x}$  in the test set. Let  $|A|$  be the cardinality of set  $A$ ,  $I(A)$  be an indicator function and  $\bar{A}$  be a complement of  $A$ .

**Subset accuracy** is defined as

$$\text{Subset accuracy}(\mathbf{y}, \hat{\mathbf{y}}) = I(\mathbf{y} = \hat{\mathbf{y}}),$$

which is directly related to subset loss. Subset accuracy measures the correctness of joint prediction for all labels. It can be regarded as a multi-label counterpart of the traditional accuracy metric.

**Hamming measure**, defined as

$$\text{Hamming measure}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{K} \sum_{k=1}^K I(y_k = \hat{y}_k)$$

is an average number of correct predictions.

**F measure**, defined as

$$\text{F measure}(\mathbf{y}, \hat{\mathbf{y}}) = 2 \cdot \frac{\text{Prec}(\mathbf{y}, \hat{\mathbf{y}}) \cdot \text{Recall}(\mathbf{y}, \hat{\mathbf{y}})}{\text{Prec}(\mathbf{y}, \hat{\mathbf{y}}) + \text{Recall}(\mathbf{y}, \hat{\mathbf{y}})},$$

is a harmonic mean of precision and recall which are defined as

$$\text{Prec}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{|\{i : y_i = 1, \hat{y}_i = 1\}|}{|\{i : \hat{y}_i = 1\}|}, \quad \text{Recall}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{|\{i : y_i = 1, \hat{y}_i = 1\}|}{|\{i : y_i = 1\}|}.$$

Precision measures how many labels are correctly predicted as relevant with respect to all labels predicted as relevant. Recall is a fraction of labels correctly predicted as relevant with respect to all relevant ones.

**Ranking measure** is defined as follows. Let  $f(\mathbf{x}, y_k)$  be a real-valued function that returns the confidence of  $y_k$  being a proper label of  $\mathbf{x}$  (i.e. confidence of  $y_k = 1$ ). We first define ranking loss as

$$\text{Ranking loss}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{|T| \cdot |\bar{T}|} |\{(y_i, y_j) \in T \times \bar{T} : f(\mathbf{x}, y_i) < f(\mathbf{x}, y_j)\}|,$$

where  $T = \{i : y_i = 1\}$ ,  $\bar{T} = \{i : y_i = 0\}$  are sets of relevant and irrelevant labels, respectively. It measures how many irrelevant labels are ranked higher than relevant ones. Ranking measure is

$$\text{Ranking measure}(\mathbf{y}, \hat{\mathbf{y}}) = 1 - \text{Ranking loss}(\mathbf{y}, \hat{\mathbf{y}}).$$

Table 1: MIMIC-III. List of features (part 1).

Description	Description	Description
1 Gender	51 Daily Weight	101 NBP [Diastolic]
2 Age	52 ALT	102 PAP [Diastolic]
3 Arterial BP [Systolic]	53 AST	103 LDH
4 Arterial BP Mean	54 Albumin (>3.2)	104 BUN
5 BSA	55 Alk. Phosphate	105 PT
6 BSA - English	56 Amylase	106 Albumin
7 BSA - Metric	57 Arterial Base Excess	107 Calcium
8 Braden Score	58 Arterial CO2(Calc)	108 Chloride
9 CaO2	59 Arterial PaCO2	109 Creatinine
10 Compliance (40-60ml)	60 Arterial PaO2	110 Glucose
11 Eye Opening	61 Arterial pH	111 INR
12 FiO2 Set	62 BUN (6-20)	112 Magnesium
13 GCS Total	63 CPK	113 PTT
14 Heart Rate	64 CPK/MB	114 Phosphorous
15 High Insp. Pressure	65 Calcium (8.4-10.2)	115 Potassium
16 High Resp. Rate	66 Carbon Dioxide	116 Sodium
17 LCW	67 Chloride (100-112)	117 Total Bili
18 LCWI	68 Creatinine (0-1.3)	118 WBC
19 LVSF	69 Differential-Bands	119 Fibrinogen
20 LVSWI	70 Differential-Basos	120 ABP Alarm [Low]
21 Low Exhaled Min Vol	71 Differential-Eos	121 HR Alarm [Low]
22 Mean Airway Pressure	72 Differential-Lymphs	122 NBP Alarm [Low]
23 Minute Volume(Observed)	73 Differential-Monos	123 Resp Alarm [Low]
24 Motor Response	74 Differential-Polys	124 SpO2 Alarm [Low]
25 NBP [Systolic]	75 Fibrinogen (150-400)	125 ABP Alarm [High]
26 NBP Mean	76 Fingertick Glucose	126 CVP Alarm [High]
27 O2 Flow (lpm)	77 Glucose (70-105)	127 HR Alarm [High]
28 PEEP Set	78 Hematocrit	128 NBP Alarm [High]
29 Peak Insp. Pressure	79 Hemoglobin	129 Resp Alarm [High]
30 Plateau Pressure	80 INR (2-4 ref. range)	130 SpO2 Alarm [High]
31 Pressure Support	81 Ionized Calcium	131 Apnea Time Interval
32 Previous Weight	82 Lactic Acid(0.5-2.0)	132 RSBI (<200)
33 Previous WeightF	83 Magnesium (1.6-2.6)	133 Spon RR (Mech.)
34 Resp Rate (Spont)	84 Mixed Venous O2% Sat	134 Spon. Vt (L) (Mech.)
35 Resp Rate (Total)	85 PT(11-13.5)	135 RSBI (<100)
36 Respiratory Rate	86 PTT(22-35)	136 Lactic Acid
37 Respiratory Rate Set	87 Phosphorous(2.7-4.5)	137 Tank A psi.
38 SVI	88 Platelets	138 Tank B psi.
39 SVRI	89 Potassium (3.5-5.3)	139 Ve High
40 Sensitivity-Vent	90 RBC	140 Heart rate Alarm - High
41 SpO2	91 SaO2	141 Heart Rate Alarm - Low
42 Stroke Volume	92 Sodium (135-148)	142 Arterial Blood Pressure systolic
43 Temperature C (calc)	93 Total Bili (0-1.5)	143 Arterial Blood Pressure diastolic
44 Temperature F	94 Vancomycin/Random	144 Arterial Blood Pressure mean
45 Tidal Volume (Observed)	95 Venous pH	145 Arterial Blood Pressure Alarm - Low
46 Tidal Volume (Set)	96 WBC (4-11,000)	146 Arterial Blood Pressure Alarm - High
47 Tidal Volume (Spont)	97 Admit Ht	147 Non Invasive Blood Pressure systolic
48 Verbal Response	98 Art.pH	148 Non Invasive Blood Pressure diastolic
49 Weight Change	99 WBC (4-11,000)	149 Non Invasive Blood Pressure mean
50 Admit Wt	100 Arterial BP [Diastolic]	150 O2 saturation pulseoxymetry

Table 2: MIMIC-III. List of features (part 1).

Description	Description	Description
151 Hematocrit (serum)	201 Gait/Transferring	251 Inspiratory Ratio
152 Chloride (serum)	202 Mental status	252 Troponin-T
153 Glucose (serum)	203 16 Gauge Dressing Occlusive	253 CK-MB
154 Sodium (serum)	204 Potassium (serum)	254 Richmond-RAS Scale
155 GCS - Eye Opening	205 HCO3 (serum)	255 18 Gauge placed in the field
156 Riker-SAS Scale	206 Platelet Count	256 CHG Bath
157 Temperature Fahrenheit	207 Prothrombin time	257 Insulin pump
158 O2 Saturation Alarm - High	208 Smoking Cessation Info	258 Arterial O2 pressure
159 O2 Saturation Alarm - Low	209 Central Venous Pressure	259 Arterial O2 Saturation
160 Pain Level	210 Minute Volume Alarm - Low	260 Arterial CO2 Pressure
161 O2 Flow	211 Minute Volume Alarm - High	261 PH (Arterial)
162 GCS - Verbal Response	212 PEEP set	262 Cuff Pressure
163 GCS - Motor Response	213 Non-Invasive BP Alarm - High	263 Spont Vt
164 Braden Sensory Perception	214 Non-Invasive BP Alarm - Low	264 Spont RR
165 Braden Moisture	215 Inspired O2 Fraction	265 Any fear in relationships
166 Braden Activity	216 Ventilator Type	266 ETOH
167 Braden Mobility	217 Ventilator Mode	267 Recreational drug use
168 Braden Nutrition	218 Paw High	268 Currently experiencing pain
169 Braden Friction/Shear	219 Vti High	269 Dialysis patient
170 Resp Alarm - High	220 Fspn High	270 Incentive Spirometry
171 Resp Alarm - Low	221 Apnea Interval	271 Alkaline Phosphate
172 Parameters Checked	222 Tidal Volume (set)	272 Total Bilirubin
173 Pain Level Response	223 Tidal Volume (observed)	273 TCO2 (calc) Arterial
174 Alarms On	224 Tidal Volume (spontaneous)	274 Baseline pain level
175 See chart for initial assessment	225 Minute Volume	275 20 Gauge placed in outside facility
176 Eye Care	226 Respiratory Rate (Set)	276 Subglottal Suctioning
177 Skin Care	227 Respiratory Rate (spontaneous)	277 High risk (>51) interventions
178 Back Care	228 Respiratory Rate (Total)	278 18 Gauge Dressing Occlusive
179 Cough/Deep Breath	229 Total PEEP Level	279 20 Gauge Dressing Occlusive
180 Arterial Line Zero/Calibrate	230 PSV Level	280 Ventilator Tank #1
181 Called Out	231 Inspiratory Time	281 Ventilator Tank #2
182 Bed Bath	232 Is the spokesperson the HCP	282 20 Gauge placed in the field
183 Calcium non-ionized	233 Unable to assess psychological	283 Discharge needs
184 Glucose finger stick	234 Harm by partner or close relation	284 Goal Richmond-RAS Scale
185 Arterial Line outside facility	235 Social work consult	285 ST Segment Monitoring On
186 Multi Lumen outside facility	236 Visual / hearing deficit	286 CAM-ICU MS Change
187 16 Gauge placed in outside facility	237 Self ADL	287 22 Gauge placed in outside facility
188 SpO2 Desat Limit	238 History of slips / falls	288 22 Gauge Dressing Occlusive
189 Marital Status	239 Intravenous access	289 Potassium (whole blood)
190 Admission Weight (Kg)	240 Difficulty swallowing	290 PH (Venous)
191 Admission Weight (lbs.)	241 Special diet	291 PH (dipstick)
192 Religion	242 Unintentional weight loss >10 lbs.	292 Impaired Skin Odor #1
193 Language	243 CK (CPK)	293 Unable to assess cognitive / perceptual
194 Anion gap	244 Differential-Neuts	294 Unable to assess activity / mobility
195 Arterial Line Dressing Occlusive	245 18 Gauge placed in outside facility	295 Unable to assess nutrition / education
196 Multi Lumen Dressing Occlusive	246 No wallet / money	296 Unable to assess teaching / learning needs
197 History of falling (within 3 mnths)	247 Sexuality / reproductive problems	297 Lipase
198 Secondary diagnosis	248 Height	298 ICU Consent Signed
199 Ambulatory aid	249 Height (cm)	299 Sodium (whole blood)
200 IV/Saline lock	250 Expiratory Ratio	300 Specific Gravity (urine)

Table 3: MIMIC-III. List of features (part 3).

Description	Description	Description
301 CVP Alarm - High	304 Multi Lumen Zero/Calibrate	307 Hematocrit (whole blood - calc)
302 CVP Alarm - Low	305 Chloride (whole blood)	308 Unable to assess habits
303 Unable to assess pain	306 Glucose (whole blood)	

Table 4: Subset accuracy for benchmark datasets, for  $T = 20$ .

	BR+l21	parCC	CC+l1	BR+l1	ACC-ABS	ACC-SUM
cal500	0.016	<b>0.028</b>	0.020	0.018	0.008	0.012
nuswide	0.312	<b>0.316</b>	0.305	0.305	0.305	0.305
medical	0.709	<b>0.710</b>	0.653	0.653	0.613	0.522
music	0.236	<b>0.267</b>	0.054	0.054	0.152	0.181
yeast	0.091	<b>0.144</b>	0.032	0.031	0.115	0.110
scene	0.023	0.023	0.000	0.000	0.041	<b>0.092</b>
mediamill	0.115	<b>0.150</b>	0.105	0.093	0.111	0.106
flags	0.144	0.211	0.191	0.144	<b>0.257</b>	0.252
Science	0.168	0.172	0.160	0.159	0.181	<b>0.192</b>
Reference	0.230	0.227	0.151	0.145	0.245	<b>0.249</b>
Health	0.308	0.314	0.309	0.310	0.319	<b>0.323</b>
avg. rank	3.6	<b>5.0</b>	2.4	1.9	4.1	4.1

Table 5: F measure for benchmark datasets, for  $T = 20$ .

	BR+l21	parCC	CC+l1	BR+l1	ACC-ABS	ACC-SUM
cal500	0.723	0.711	<b>0.729</b>	0.726	0.708	0.708
nuswide	<b>0.058</b>	0.033	0.000	0.000	0.001	0.000
medical	0.612	<b>0.612</b>	0.574	0.574	0.518	0.446
music	0.492	<b>0.565</b>	0.172	0.172	0.382	0.444
yeast	<b>0.564</b>	0.558	0.485	0.475	0.546	0.545
scene	0.024	0.024	0.000	0.000	0.043	<b>0.103</b>
mediamill	<b>0.588</b>	0.584	0.536	0.520	0.542	0.533
flags	0.678	0.688	0.687	0.683	<b>0.699</b>	0.698
Science	0.511	0.514	0.508	0.507	0.516	<b>0.521</b>
Reference	0.652	0.650	0.621	0.619	0.658	<b>0.658</b>
Health	0.518	0.526	0.522	0.521	<b>0.534</b>	0.533
avg. rank	4.0	<b>4.4</b>	2.6	2.0	4.1	3.9

Table 6: Hamming measure for benchmark datasets, for  $T = 20$ .

	BR+l21	parCC	CC+l1	BR+l1	ACC-ABS	ACC-SUM
cal500	<b>0.625</b>	0.616	0.619	0.623	0.606	0.607
nuswide	<b>0.886</b>	0.883	0.877	0.877	0.877	0.877
medical	0.962	<b>0.963</b>	0.957	0.957	0.952	0.943
music	<b>0.788</b>	0.778	0.716	0.716	0.742	0.740
yeast	<b>0.734</b>	0.718	0.710	0.707	0.718	0.715
scene	0.825	0.825	0.821	0.821	0.827	<b>0.837</b>
mediamill	<b>0.809</b>	0.809	0.793	0.785	0.791	0.789
flags	0.733	0.737	0.731	0.735	<b>0.739</b>	0.738
Science	0.872	0.872	0.871	0.870	0.873	<b>0.874</b>
Reference	0.900	0.899	0.891	0.891	0.901	<b>0.901</b>
Health	<b>0.825</b>	0.815	0.811	0.814	0.814	0.815
avg. rank	<b>4.9</b>	4.3	2.2	2.0	3.8	3.7

Table 7: Subset accuracy for benchmark datasets, for  $T = 10$ .

	BR+l21	parCC	CC+l1	BR+l1	ACC-ABS	ACC-SUM
cal500	0.008	<b>0.024</b>	0.002	0.006	0.022	0.020
nuswide	0.305	<b>0.305</b>	0.305	0.305	0.305	0.305
medical	0.446	0.446	0.213	0.213	<b>0.529</b>	0.505
music	0.189	<b>0.235</b>	0.012	0.012	0.010	0.024
yeast	0.035	<b>0.105</b>	0.031	0.031	0.043	0.067
scene	<b>0.002</b>	0.002	0.000	0.000	0.000	0.000
mediamill	0.097	<b>0.109</b>	0.094	0.094	0.093	0.104
flags	0.155	0.196	0.150	0.150	0.227	<b>0.237</b>
Science	0.158	0.157	0.152	0.152	0.161	<b>0.164</b>
Reference	0.161	0.163	0.118	0.120	0.201	<b>0.218</b>
Health	0.308	0.310	0.308	0.308	0.310	<b>0.316</b>
avg. rank	3.6	<b>5.0</b>	1.9	2.0	3.7	4.7

Table 8: F measure for benchmark datasets, for  $T = 10$ .

	BR+l21	parCC	CC+l1	BR+l1	ACC-ABS	ACC-SUM
cal500	0.726	0.717	0.724	0.719	0.730	<b>0.732</b>
nuswide	<b>0.017</b>	0.002	0.000	0.000	0.000	0.000
medical	0.293	0.293	0.000	0.000	<b>0.430</b>	0.417
music	0.421	<b>0.498</b>	0.065	0.065	0.067	0.107
yeast	0.520	<b>0.531</b>	0.475	0.475	0.500	0.516
scene	<b>0.002</b>	0.002	0.000	0.000	0.000	0.000
mediamill	<b>0.550</b>	0.538	0.521	0.521	0.520	0.511
flags	<b>0.700</b>	0.698	0.690	0.691	0.696	0.698
Science	0.506	0.507	0.504	0.504	0.508	<b>0.509</b>
Reference	0.625	0.626	0.609	0.610	0.641	<b>0.647</b>
Health	0.519	0.523	0.522	0.522	0.524	<b>0.528</b>
avg. rank	4.4	<b>4.5</b>	2.0	2.0	3.8	4.3

Table 9: Hamming measure for benchmark datasets, for  $T = 10$ .

	BR+l21	parCC	CC+l1	BR+l1	ACC-ABS	ACC-SUM
cal500	<b>0.623</b>	0.616	0.616	0.618	0.614	0.619
nuswide	<b>0.880</b>	0.878	0.877	0.877	0.877	0.877
medical	0.933	0.933	0.906	0.906	<b>0.943</b>	0.941
music	<b>0.773</b>	0.769	0.697	0.697	0.701	0.702
yeast	<b>0.722</b>	0.719	0.707	0.707	0.714	0.713
scene	<b>0.821</b>	0.821	0.821	0.821	0.821	0.821
mediamill	<b>0.795</b>	0.794	0.779	0.779	0.781	0.788
flags	<b>0.745</b>	0.744	0.730	0.731	0.736	0.739
Science	0.870	0.870	0.870	0.870	0.871	<b>0.871</b>
Reference	0.892	0.893	0.888	0.888	0.897	<b>0.898</b>
Health	<b>0.818</b>	0.812	0.809	0.809	0.810	0.812
avg. rank	<b>5.3</b>	4.4	1.8	2.0	3.5	4.2

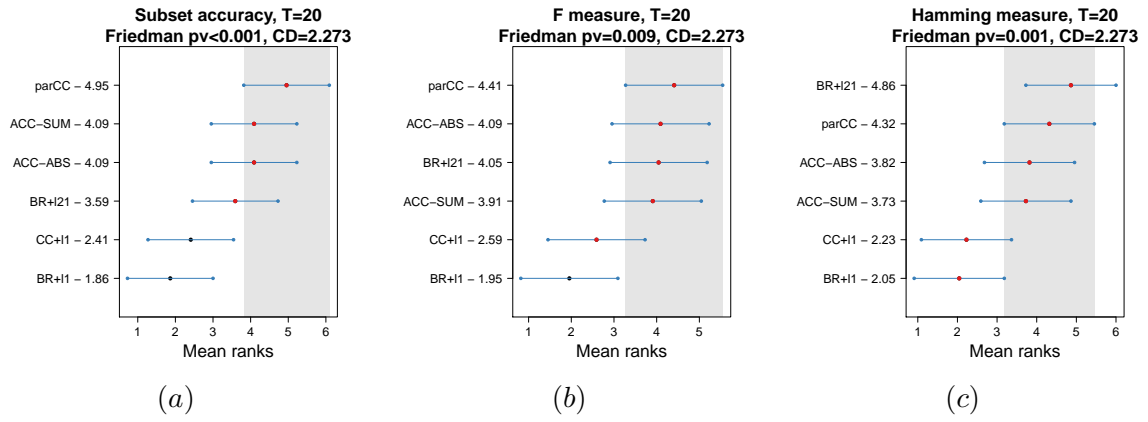


Figure 1: Results of Friedman and pairwise tests, for  $T = 20$ .

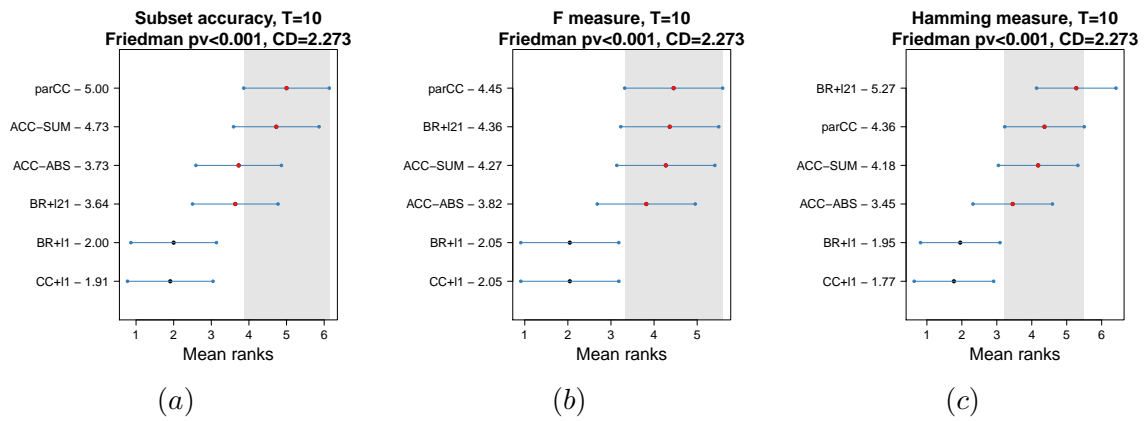


Figure 2: Results of Friedman and pairwise tests, for  $T = 10$ .