
000
001
002
003
004
005
006
007
008
009
010
011

Supplementary Material - Efficient Structured Prediction with Latent Variables for General Graphical Models

012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109

1. Introduction

After having introduced the temperature parameter ϵ we want to minimize the following cost function w.r.t. model parameters w :

$$\frac{C}{p} \|w\|_p^p + \sum_{(x,y) \in \mathcal{D}} \left(\epsilon \ln \sum_{\hat{s} \in \mathcal{S}} \exp \left(\frac{w^\top \phi(x, \hat{s}) + \ell_{(x,y)}(\hat{s})}{\epsilon} \right) - \epsilon \ln \sum_{\hat{h} \in \mathcal{H}} \exp \left(\frac{w^\top \phi(x, (y, \hat{h})) + \ell_{(x,y)}^c((y, \hat{h}))}{\epsilon} \right) \right). \quad (1)$$

2. Proof of Claim 1

Claim 1 *The following function*

$$\frac{C}{p} \|w\|_p^p + \sum_{(x,y)} \left(\epsilon \ln \sum_{\hat{s} \in \mathcal{S}} \exp \left(\frac{w^\top \phi(x, \hat{s}) + \ell_{(x,y)}(\hat{s})}{\epsilon} \right) - \epsilon H(q_{(x,y)}) - \mathbb{E}_{q_{(x,y)}}[w^\top \phi(x, (y, \hat{h})) + \ell_{(x,y)}^c((y, \hat{h}))] \right), \quad (2)$$

convex in w and $q_{(x,y)}$ separately, is an upper bound on Eq. (1), $\forall q_{(x,y)}(h) \in \Delta$, with Δ denoting the probability simplex, H indicating the entropy and \mathbb{E} referring to the expectation w.r.t. the stated distribution. The bound holds with equality for the $q_{(x,y)}^(h)$ minimizing this cost function.*

Proof: The partition function is the conjugate dual of the entropy, hence:

$$-\epsilon \ln \sum_{\hat{h}} \exp \frac{w^\top \phi(x, \hat{y}, \hat{h}) + \ell_{(x,y)}^c((y, \hat{h}))}{\epsilon} = \min_{q_{(x,y)} \in \Delta} -\epsilon H(q_{(x,y)}) - \sum_{\hat{h}} q_{(x,y)}(\hat{h})(w^\top \phi(x, y, \hat{h}) + \ell_{(x,y)}^c((y, \hat{h}))), \quad (3)$$

to obtain the problem stated in Eq. (2) It is easy to see that the cost function given in Eq. (2) is convex in w and $q_{(x,y)}$ $\forall (x, y)$ separately. However not jointly convex in w and $q_{(x,y)}$. Neglecting minimization w.r.t. $q_{(x,y)}$ results in an upper bound to the original problem. The original problem is attained for optimal $q_{(x,y)}^*$. \square

3. Proof of Theorem 1

Theorem 1 *The approximation of the program in Eq. (2) takes the form given in Program 1 where $\phi_{(x,y),i}(s_i) = \ell_{(x,y),i}(x, s_i) + \sum_{r:i \in \mathbb{S}_r} w_r \phi_{r,i}(x, s_i)$ and $\phi_{(x,y),\alpha}(s_\alpha) = \ell_{(x,y),\alpha}(x, s_\alpha) + \sum_{r:\alpha \in E_r} w_r \phi_{r,\alpha}(x, s_\alpha)$.*

Proof: Assume that the each element of the feature vector ϕ decomposes into a graphical model structure, i.e., the r -th element takes the following form:

$$\phi_r(x, s) = \sum_{\alpha \in E_r} \phi_{r,\alpha}(x, s_\alpha) + \sum_{i \in \mathbb{S}_r} \phi_{r,i}(x, s_i). \quad (4)$$

with E_r , \mathbb{S}_r the sets of factors and variables. Note that each feature is described by a bipartite factor graph G_r with nodes originating from the variable set \mathbb{S}_r and factors from E_r . An edge connects a single node $i \in \mathbb{S}_r$ to a factor $\alpha \in E_r$ iff $i \in \alpha$. Consider the factor graph $G = \bigcup_r G_r$ where we define the set of neighbors $N(i) := \{\alpha : i \in \alpha \forall \alpha \in E\}$ and $N(\alpha) := \{i : i \in \alpha \forall i \in \mathbb{S}\}$.

110 **Program 1** Approximated structured prediction with latent variables

111

112
$$\min_{d, \lambda, w} f_1 \left\{ \frac{C}{2} \|w\|_2^2 + \sum_{(x,y) \in \mathcal{D}} \left(\sum_{i \in \mathbb{S}} \epsilon c_i \ln \sum_{s_i} \exp \left(\frac{\phi_{(x,y),i}(s_i) - \sum_{\alpha \in N(i)} \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_i} \right) + \right. \right.$$

113
$$f_2 \left\{ - \sum_r w_r \left(\sum_{(x,y)} \left(\sum_{i \in \mathbb{Y}} \phi_{r,i}(x, y_i) + \sum_{i \in \mathbb{H}, h_i} \phi_{r,i}(x, h_i) d_{(x,y),i}(h_i) + \sum_{\alpha \in E, h_\alpha} \phi_{r,\alpha}(x, (y, h)_\alpha) d_{(x,y),\alpha}(h_\alpha) \right) \right) \right. \right.$$

114
$$f_3 \left\{ - \sum_{(x,y)} \left(\sum_{i \in \mathbb{H}, h_i} \ell_{(x,y),i}^c(x, h_i) d_{(x,y),i}(h_i) + \sum_{\alpha \in E, h_\alpha} \ell_{(x,y),\alpha}^c(x, (y, h)_\alpha) d_{(x,y),\alpha}(h_\alpha) \right) \right. \right.$$

115
$$- \sum_{(x,y)} \left(\sum_{i \in \mathbb{H}} \epsilon \hat{c}_i H(d_{(x,y),i}) + \sum_{\alpha \in E} \epsilon \hat{c}_\alpha H(d_{(x,y),\alpha}) \right)$$

116 s.t.
$$\left. \begin{array}{l} \sum_{h_\alpha \setminus h_i} d_{(x,y),\alpha}(h_\alpha) = d_{(x,y),i}(h_i) \quad \forall (x,y), i \in \mathbb{H}, \alpha \in N(i), h_i \in \mathcal{S}_i \\ d_{(x,y),i}, d_{(x,y),\alpha} \in \Delta \end{array} \right\} := d_{(x,y)} \in \mathcal{C}_{(x,y)} \quad \forall (x,y) \in \mathcal{D}$$

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

In many applications the loss functions ℓ and ℓ^c factorize in a similar way and are easily included into the graphical model G , i.e.,

$$\ell_{(x,y)}(s) = \sum_{i \in \mathbb{S}} \ell_{(x,y),i}(s_i) + \sum_{\alpha \in E} \ell_{(x,y),\alpha}(s_\alpha), \quad (5)$$

$$\ell_{(x,y)}^c(\hat{h}) = \sum_{i \in \mathbb{H}} \ell_{(x,y),i}^c(\hat{s}_i) + \sum_{\alpha \in E} \ell_{(x,y),\alpha}^c((y, \hat{h})_\alpha). \quad (6)$$

Let $d_{(x,y),i}, d_{(x,y),\alpha}$ be the marginals of $q_{(x,y)}(h)$. Then,

$$\sum_{\hat{h}} \ell_{(x,y)}^c((y, \hat{h})) q_{(x,y)}(\hat{h}) = \ell_{(x,y)}^c(y) = \sum_{i \in \mathbb{Y}, y_i} \ell_{(x,y),i}^c(x, y_i) + \sum_{i \in \mathbb{H}, h_i} \ell_{(x,y),i}^c(x, h_i) d_{(x,y),i}(h_i) + \sum_{\alpha \in E, h_\alpha} \ell_{(x,y),\alpha}^c(x, (y, h)_\alpha) d_{(x,y),\alpha}(h_\alpha) \quad (7)$$

$$\sum_{\hat{h}} \phi_r(x, y, \hat{h}) q_{(x,y)}(\hat{h}) = v_{(x,y),r} = \sum_{i \in \mathbb{Y}} \phi_{r,i}(x, y_i) + \sum_{i \in \mathbb{H}, h_i} \phi_{r,i}(x, h_i) d_{(x,y),i}(h_i) + \sum_{\alpha \in E, h_\alpha} \phi_{r,\alpha}(x, (y, h)_\alpha) d_{(x,y),\alpha}(h_\alpha) \quad (8)$$

$$v_r = \sum_{(x,y)} v_{(x,y),r} \quad (9)$$

Note that $\ell_{(x,y),i}^c(x, y_i)$ has no effect and can be neglected. We approximate the entropy via local entropy terms and introduce counting numbers $\hat{c}_i, \hat{c}_\alpha$. Moreover, we approximate the marginal polytope via the local polytope

given by the marginalization and simplex constraints. All in all we obtain the following approximated program:

$$\begin{aligned}
 \min_{w,d,\sigma} \quad & \frac{C}{p} \|w\|_p^p + \sum_{(x,y) \in \mathcal{D}} \left(\epsilon \ln \sum_s \exp \frac{\sigma(x,s) + \ell_{(x,y)}(s)}{\epsilon} - \underbrace{\sum_r w_r v_{(x,y),r}}_{f_2} \right. \\
 & \left. - l_{(x,y)}^c(y) - \underbrace{\sum_{i \in \mathbb{H}} \epsilon \hat{c}_i H(d_{(x,y),i}) - \sum_{\alpha \in E_{\mathbb{H}}} \epsilon \hat{c}_{\alpha} H(d_{(x,y),\alpha})}_{f_3} \right) \\
 \text{s.t.} \quad & \forall (x,y), i \in \mathbb{H}, \alpha \in N(i), h_i \sum_{h_{\alpha} \setminus h_i} d_{(x,y),\alpha}(h_{\alpha}) = d_{(x,y),i}(h_i) \\
 & d_{(x,y),i}, d_{(x,y),\alpha} \in \Delta \\
 & \forall (x,y), s \quad \sigma(x,s) = w^T \phi(x,s)
 \end{aligned}$$

For clarity of the presentation let us now address minimization w.r.t. w and σ without regards to f_3 and the marginalization constraints on d required for global consistency, i.e.,

$$\min_{w,\sigma} \quad \sum_{(x,y)} \epsilon \ln \sum_s \exp \frac{\sigma(x,s) + \ell_{(x,y)}(s)}{\epsilon} - \sum_r w_r v_r + \frac{C}{p} \|w\|_p^p \quad (10)$$

$$\text{s.t.} \quad \forall (x,y), s \quad \sigma(x,s) = w^T \phi(x,s) \quad (11)$$

In a subsequent step we minimize the Lagrangian w.r.t. the primal variables w and σ , i.e.,

$$\sum_{(x,y)} \min_{\sigma} \left(\epsilon \ln \sum_s \exp \frac{\sigma(x,s) + \ell_{(x,y)}(s)}{\epsilon} - \sum_s p_{(x,y)}(s) \sigma(x,s) \right) + \quad (12)$$

$$\min_w \left(\frac{C}{p} \|w\|_p^p + w^T \left(\sum_{(x,y),s} p_{(x,y)}(s) \phi(x,s) - v \right) \right). \quad (13)$$

Note that we have introduced Lagrangian multipliers $p_{(x,y)}(s) \forall (x,y), s$. Analytically carrying out the minimization, we obtain the following dual problem

$$\max_{p_{(x,y)} \in \Delta} \sum_{(x,y)} \epsilon H(p_{(x,y)}) + \sum_s p_{(x,y)}(s) (w^T \phi(x,s) + \ell_{(x,y)}(s)) - \frac{C^{1-q}}{q} \sum_r \left| \sum_{(x,y),s} p_{(x,y)}(s) \phi_r(x,s) - v_{(x,y),r} \right|^q. \quad (14)$$

Similar to our previous argument we assume that the probability distribution $p_{(x,y)}(s)$ is defined via marginals $b_{(x,y),i}, b_{(x,y),\alpha}$. Again we approximate the marginal polytope via a local one given by the marginalization and simplex constraints. Hence we obtain,

$$\sum_s \ell_{(x,y)}(s) p_{(x,y)}(s) = l_{(x,y)} = \sum_{i \in S, s_i} \ell_{(x,y),i}(s_i) b_{(x,y),i}(s_i) + \sum_{\alpha \in E, s_{\alpha}} \ell_{(x,y),\alpha}(s_{\alpha}) b_{(x,y),\alpha}(s_{\alpha}) \quad (15)$$

$$\sum_s \phi_r(x,s) p_{(x,y)}(s) = u_{(x,y),r} = \sum_{i \in S, s_i} \phi_{r,i}(x,s_i) b_{(x,y),i}(s_i) + \sum_{\alpha \in E_r, s_{\alpha}} \phi_{r,\alpha}(x,s_{\alpha}) b_{(x,y),\alpha}(s_{\alpha}) \quad (16)$$

$$u_r = \sum_{(x,y)} u_{(x,y),r} \quad (17)$$

330	The approximated dual program with counting numbers c_i, c_α is then	385
331		386
332	$\max_{b,u} \sum_{(x,y)} \sum_{i \in \mathbb{S}} \epsilon c_i H(b_{(x,y),i}) + \sum_{\alpha \in E} \epsilon c_\alpha H(b_{(x,y),\alpha}) + \sum_{i \in \mathbb{S}, s_i} b_{(x,y),i}(s_i) \phi_{(x,y),i}(s_i) + \sum_{\alpha \in E, s_\alpha} b_{(x,y),\alpha}(s_\alpha) \phi_{(x,y),\alpha}(s_\alpha) - \frac{C^{1-q}}{q} \sum_r u_r - v_r ^q$	387
333		388
334		389
335		390
336		391
337	s.t. $\forall \alpha \in E, i \in N(\alpha), s_i, (x,y) \sum_{s_\alpha \setminus s_i} b_{(x,y),\alpha}(s_\alpha) = b_{(x,y),i}(s_i)$	392
338		393
339	$b_{(x,y),i}, b_{(x,y),\alpha} \in \Delta$	394
340		395
341	By maximizing the obtained Lagrangian with Lagrange multipliers $\lambda_{(x,y),i \rightarrow \alpha}(s_i)$ for the constraint given in Eq. (19) and after combining with f_3 and corresponding constraints on d we obtain the claim. \square	396
342		397
343		398
344		399
345		400
346	Claim 2 Algorithm 1 is guaranteed to decrease the cost function of Program 1 at every iteration and guaranteed to converge to a minimum or a saddle point for $\epsilon, c_i, c_\alpha, \hat{c}_i, \hat{c}_\alpha > 0$.	401
347		402
348		403
349	Proof: Recalling Theorem 5 in (Yuille & Rangarajan, 2003) we notice that alternating optimization of the approximated program in Theorem 1 w.r.t. d and λ, w is equivalent to CCCP, which is guaranteed to converge to a stationary point if the respective functions are convex. Convexity is ensured for $\epsilon, c_i, c_\alpha, \hat{c}_i$ and $\hat{c}_\alpha > 0$. \square	404
350		405
351		406
352		407
353		408
354		409
355	5. Algorithmic Details	
356	Similar to a CCCP approach (Yuille & Rangarajan, 2003) we address optimization of Program 1 by alternating	410
357	solving two tasks. One optimization considers only the beliefs d , while the other one operates on λ and w . When	411
358	updating the beliefs d , we obtain the ‘latent variable prediction problem’ which requires solving	412
359		413
360	$\sum_{(x,y) \in \mathcal{D}} \min_{d_{(x,y)}} f_2(w, d) + f_3(d)$	414
361	s.t. $\forall (x,y) \in \mathcal{D} \quad d_{(x,y)} \in \mathcal{C}_{(x,y)}.$	415
362		416
363	Explicitly it reads as	417
364		418
365	$\sum_{(x,y)} \max_d \sum_{i \in \mathbb{H}} \epsilon \hat{c}_i H(d_{(x,y),i}) + \sum_{\alpha \in E_{\mathbb{H}}} \epsilon \hat{c}_\alpha H(d_{(x,y),\alpha}) + \sum_{i \in \mathbb{H}, h_i} d_{(x,y),i}(h_i) \left(\sum_{r:i \in \mathbb{S}_r} w_r \phi_{r,i}(x, h_i) + l_{(x,y),i}^c(h_i) \right) + \sum_{\alpha \in E_{\mathbb{H}}, h_\alpha} d_{(x,y),\alpha}(h_\alpha) \left(\sum_{r:\alpha \in E_r} w_r \phi_{r,\alpha}(x, (y, h)_\alpha) + l_{(x,y),\alpha}^c(x, (y, h)_\alpha) \right)$	419
366		420
367		421
368		422
369		423
370		424
371	s.t. $\forall (x,y), i \in \mathbb{H}, \alpha \in N(i), h_i \quad \sum_{h_\alpha \setminus h_i} d_{(x,y),\alpha}(h_\alpha) = d_{(x,y),i}(h_i)$	425
372		426
373		427
374	$\forall (x,y) \quad d_{(x,y),i}, d_{(x,y),\alpha} \in \Delta$	428
375		429
376	This is a standard (convex) belief propagation task with local potentials $\phi_{(x,y),i}^c(h_i) = \sum_{r:i \in \mathbb{S}_r} w_r \phi_{r,i}(x, h_i) + l_{(x,y),i}^c(h_i)$ and clique potentials $\phi_{(x,y),\alpha}^c(h_\alpha) = \sum_{r:\alpha \in E_r} w_r \phi_{r,\alpha}(x, (y, h)_\alpha) + l_{(x,y),\alpha}^c(x, (y, h)_\alpha)$. We solve it by	430
377	minimizing its unconstrained dual program using a message passing algorithm on the graph defined by the nodes	431
378	$i \in \mathbb{H}$ and corresponding cliques. The algorithm is guaranteed to find the optimal solution for counting numbers	432
379	$\hat{c}_i, \hat{c}_\alpha$ and $\epsilon > 0$.	433
380		434
381	To optimize for model parameters w and messages λ we gradually solve the following unconstrained problem:	435
382		436
383	$\min_{w,\lambda} f_1(w, \lambda) + f_2(w, d).$	437
384		438
		439

440 It explicitly reads as

$$\begin{aligned}
 \min_{\lambda, w} \quad & \frac{C}{2} \|w\|_2^2 + \sum_{(x,y) \in \mathcal{D}} \left(\sum_{i \in \mathbb{S}} \epsilon c_i \ln \sum_{s_i} \exp \left(\frac{\phi_{(x,y),i}(s_i) - \sum_{\alpha \in N(i)} \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_i} \right) + \right. \\
 & \left. \sum_{\alpha \in E} \epsilon c_\alpha \ln \sum_{s_\alpha} \exp \left(\frac{\phi_{(x,y),\alpha}(s_\alpha) + \sum_{i \in N(\alpha)} \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_\alpha} \right) \right) - \\
 & - \sum_r w_r \underbrace{\left(\sum_{(x,y)} \left(\sum_{i \in \mathbb{Y}} \phi_{r,i}(x, y_i) + \sum_{i \in \mathbb{H}, h_i} \phi_{r,i}(x, h_i) d_{(x,y),i}(h_i) + \sum_{\alpha \in E, h_\alpha} \phi_{r,\alpha}(x, (y, h)_\alpha) d_{(x,y),\alpha}(h_\alpha) \right) \right)}_{v_r}.
 \end{aligned} \tag{24}$$

442 This is an approximated structured prediction task with empirical means v_r .

Lemma 1 Given a node $i \in \mathbb{S}$ of the graphical model G , the optimal $\lambda_{(x,y),i \rightarrow \alpha}(s_i) \forall \alpha \in N(i), s_i \in \mathcal{S}_i, (x, y) \in \mathcal{D}$ of Theorem 1 (resp. Eq. (23)) satisfies

$$\lambda_{(x,y),i \rightarrow \alpha}(s_i) \propto \frac{c_\alpha}{c_i + \sum_{\alpha \in N(i)} c_\alpha} \left(\phi_{(x,y),i}(s_i) + \sum_{\alpha \in N(i)} \mu_{(x,y),\alpha \rightarrow i}(s_i) \right) - \mu_{(x,y),\alpha \rightarrow i}(s_i). \tag{25}$$

443 with

$$\mu_{(x,y),\alpha \rightarrow i}(s_i) = \epsilon c_\alpha \ln \sum_{s_\alpha \setminus s_i} \exp \frac{\phi_{(x,y),\alpha}(s_\alpha) + \sum_{u \in N(\alpha) \setminus i} \lambda_{(x,y),u \rightarrow \alpha}(s_u)}{\epsilon c_\alpha} \tag{26}$$

Proof:

To update the messages we take every $(x, y), i \in \mathbb{S}$ and obtain an analytic solution of the first order optimality condition w.r.t. $\lambda_{(x,y),i \rightarrow \alpha}(s_i) \forall \alpha \in N(i), s_i$. We obtain the following simplified optimization problem:

$$\begin{aligned}
 \min_{\lambda_{(x,y),i \rightarrow \alpha}(s_i)} \quad & \epsilon c_i \ln \sum_{s_i} \exp \frac{\phi_{(x,y),i}(s_i) - \sum_{\alpha \in N(i)} \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_i} + \\
 & \sum_{\alpha \in N(i)} \epsilon c_\alpha \ln \sum_{s_i} \exp \frac{\mu_{(x,y),\alpha \rightarrow i}(s_i) + \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_\alpha}
 \end{aligned} \tag{27}$$

We find the optimal $\lambda_{(x,y),i \rightarrow \alpha}(s_i) \forall \alpha \in N(i), s_i$ whenever the gradient vanishes which is achieved for

$$\lambda_{(x,y),i \rightarrow \alpha}(s_i) \propto \frac{c_\alpha}{c_i + \sum_{\alpha \in N(i)} c_\alpha} \left(\phi_{(x,y),i}(s_i) + \sum_{\alpha \in N(i)} \mu_{(x,y),\alpha \rightarrow i}(s_i) \right) - \mu_{(x,y),\alpha \rightarrow i}(s_i). \tag{28}$$

481 \square

Lemma 2 The gradient of the approximated program given in Theorem 1 (resp. Eq. (23)) w.r.t. w_r equals

$$\sum_{(x,y)} \left(\sum_{i \in \mathbb{S}_r, s_i} b_{(x,y),i}(s_i) \phi_{r,i}(x, s_i) + \sum_{\alpha \in E_r, s_\alpha} b_{(x,y),\alpha}(s_\alpha) \phi_{r,\alpha}(x, s_\alpha) \right) - v_r + C |w_r|^{p-1} \text{sign}(w_r) \tag{29}$$

487 with

$$\begin{aligned}
 b_{(x,y),i}(s_i) & \propto \exp \left(\frac{\phi_{(x,y),i}(s_i) - \sum_{\alpha \in N(i)} \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_i} \right) \\
 b_{(x,y),\alpha}(s_\alpha) & \propto \exp \left(\frac{\phi_{(x,y),\alpha}(s_\alpha) + \sum_{i \in N(\alpha)} \lambda_{(x,y),i \rightarrow \alpha}(s_i)}{\epsilon c_\alpha} \right)
 \end{aligned}$$

550	Program 2 Message passing algorithm for approximated structured prediction with latent variables	605
551	Repeat	606
552		607
553	1. Solve ‘latent variable prediction problem’ (Program (22)) until convergence	608
554		609
555	2. For each $(x, y) \in \mathcal{D}, i \in \mathbb{S}$ (by Lemma 1):	610
556		611
557	$\forall \alpha \in N(i), s_i \quad \lambda_{(x,y),i \rightarrow \alpha}(s_i) \propto \frac{c_\alpha}{c_i + \sum_{\alpha \in N(i)} c_\alpha} \left(\phi_{(x,y),i}(s_i) + \sum_{\alpha \in N(i)} \mu_{(x,y),\alpha \rightarrow i}(s_i) \right) - \mu_{(x,y),\alpha \rightarrow i}(s_i)$	612
558		613
559		614
560		615
561	3. For each r (by Lemma 2): find a stepsize η that reduces $f_1 + f_2$ and update	616
562		617
563	$w_r \leftarrow w_r - \eta \left(\sum_{(x,y)} \left(\sum_{i \in \mathbb{S}_r, s_i} b_{(x,y),i}(s_i) \phi_{r,i}(x, s_i) + \sum_{\alpha \in E_r, s_\alpha} b_{(x,y),\alpha}(s_\alpha) \phi_{r,\alpha}(x, s_\alpha) \right) - v_r + C w_r ^{p-1} \text{sign}(w_r) \right)$	618
564		619
565		620
566		621
567	Proof: This is a direct computation of the gradient w.r.t. w_r . \square	622
568		623
569	Hence the complete algorithm reads as detailed in Program 2	624
570		625
571		626
572		627
573	References	628
574	Yuille, A. L. and Rangarajan, A. The Concave-Convex Procedure (CCCP). <i>Neural Computation</i> , 2003.	629
575		630
576		631
577		632
578		633
579		634
580		635
581		636
582		637
583		638
584		639
585		640
586		641
587		642
588		643
589		644
590		645
591		646
592		647
593		648
594		649
595		650
596		651
597		652
598		653
599		654
600		655
601		656
602		657
603		658
604		659