# PromoterPredict: Sequence-based modelling of *Escherichia coli* σ⁷⁰ promoter strength yields logarithmic dependence between promoter strength and sequence

**Ramit Bharanikumar** [1] , **Keshav Aditya R Premkumar** [2] , **Ashok Palaniappan** Corresp. [3]

[1] Biotechnology, Sri Venkateswara College of Engineering (autonomous), Sriperumbudur, Tamil Nadu, India

[2] Computer Science and Engineering, Sri Venkateswara College of Engineering (autonomous), Sriperumbudur, Tamil Nadu, India

[3] Bioinformatics, School of Chemical and BioTechnology, SASTRA Deemed University, Thanjavur, Tamil Nadu, India

Corresponding Author: Ashok Palaniappan
Email address: apalania@scbt.sastra.edu

We present PromoterPredict, a dynamic multiple regression approach to predict the strength of *Escherichia coli* promoters binding the $\sigma^{70}$ factor of RNA polymerase. $\sigma^{70}$ promoters are ubiquitously used in recombinant DNA technology, but characterizing their strength is demanding in terms of both time and money. We parsed a comprehensive database of bacterial promoters for the –35 and –10 hexamer regions of $\sigma^{70}$-binding promoters and used these sequences to construct the respective position weight matrices (PWM). Next we used a well-characterized set of promoters to train a multivariate linear regression model and learn the mapping between PWM scores of the –35 and –10 hexamers and the promoter strength. We found that the log of the promoter strength is significantly linearly associated with a weighted sum of the –10 and –35 sequence profile scores. We applied our model to 100 sets of 100 randomly generated promoter sequences to generate a sampling distribution of mean strengths of random promoter sequences and obtained a mean of 6E-4 ± 1E-7. Our model was further validated by cross-validation and on independent datasets of characterized promoters. PromoterPredict accepts –10 and –35 hexamer sequences and returns the predicted promoter strength. It is capable of dynamic learning from user-supplied data to refine the model construction and yield more robust estimates of promoter strength. PromoterPredict is available as both a web service ( https://promoterpredict.com ) and standalone tool ( https://github.com/PromoterPredict ). Our work presents an intuitive generalization applicable to modelling the strength of other promoter classes.

1  **PromoterPredict: sequence-based modelling of *Escherichia coli* σ⁷⁰**
2  **promoter strength yields logarithmic dependence between promoter**
3  **strength and sequence**
4

5  Ramit Bharanikumar[1], Keshav Aditya R. Premkumar[2] and Ashok Palaniappan[3*]
6  [1]Biotechnology, Sri Venkateswara College of Engineering (autonomous),Tamil Nadu.
7  India
8  [2]Computer Science and Engineering, Sri Venkateswara College of Engineering
9  (autonomous), Tamil Nadu. India
10  [3]Bioinformatics, School of Chemical and Biotechnology, SASTRA Deemed University,
11  Thanjavur, Tamil Nadu. India
12  [*]To whom correspondence should be addressed (apalania@scbt.sastra.edu)
13

14  **Abstract**: We present PromoterPredict, a dynamic multiple regression approach to
15  predict the strength of *Escherichia coli* promoters binding the σ⁷⁰ factor of RNA
16  polymerase. σ⁷⁰ promoters are ubiquitously used in recombinant DNA technology, but
17  characterizing their strength is demanding in terms of both time and money. We parsed
18  a comprehensive database of bacterial promoters for the −35 and −10 hexamer regions
19  of σ⁷⁰-binding promoters and used these sequences to construct the respective position
20  weight matrices (PWM). Next we used a well-characterized set of promoters to train a
21  multivariate linear regression model and learn the mapping between PWM scores of the
22  −35 and −10 hexamers and the promoter strength. We found that the log of the
23  promoter strength is significantly linearly associated with a weighted sum of the −10
24  and −35 sequence profile scores. We applied our model to 100 sets of 100 randomly
25  generated promoter sequences to generate a sampling distribution of mean strengths of
26  random promoter sequences and obtained a mean of $6E-4 \pm 1E-7$. Our model was
27  further validated by cross-validation and on independent datasets of characterized
28  promoters. PromoterPredict accepts −10 and −35 hexamer sequences and returns the
29  predicted promoter strength. It is capable of dynamic learning from user-supplied data
30  to refine the model construction and yield more robust estimates of promoter strength.
31  PromoterPredict is available as both a web service (https://promoterpredict.com) and
32  standalone tool (https://github.com/PromoterPredict). Our work presents an intuitive
33  generalization applicable to modelling the strength of other promoter classes.

34 **INTRODUCTION**

35

36    The primary *E. coli* promoter-specificity factor and the one widely used in recombinant
37    DNA technology is the $\sigma^{70}$ factor. Promoters recognized by $\sigma^{70}$-containing RNA
38    polymerase are called core promoters and share the following features: two conserved
39    hexamer sequences, separated by a non-specific spacer of ideally 17 nucleotides. The two
40    hexamers are located ~ 35 bp and ~10 bp upstream of the transcription start site, and
41    are called the −35 and −10 sequences respectively (Maquat and Reznikoff, 1978; Bujard,
42    1980; Paget and Helmann, 2003; Kadonaga, 2012). −35 and −10 sequences matching
43    the consensi motifs (TTGACA and TATAAT, respectively) are known as canonical
44    hexamers  (Galas, et al. 1985; Deuschle, et al. 1986; Stormo, 1990). It is known that the
45    conserved hexamer regions are vital for recognizing and optimizing the interactions
46    between DNA and the RNA polymerase (Hawley and McClure, 1983; Knaus and Bujard,
47    1990; Hook-Barnard *et al.*, 2006; Feklistov and Darst, 2011; Basu *et al.*, 2014).

48    Theory has yielded a linear relationship between the total promoter score and the
49    natural log of promoter strength (Berg and von Hippel, 1987; Li and Zhang, 2014).
50    Nucleotide occurence frequencies were first used by Weller and Recknagel (1994) in
51    promoter strength prediction. Additivity in promoter-polymerase interaction has been
52    affirmed by Stormo and colleagues (2002). Patterns in $\sigma^{70}$ promoters have been
53    quantified by Huerta and Collado-Vides (2003). Strength of *E. coli* $\sigma^{E}$ RNA polymerase
54    promoters  were studied by  Rhodius and Mutalik (2010). . The complexity of *E. coli* $\sigma^{70}$
55    promoter sequences has been treated from an information theoretic standpoint by
56    Shultzaberger *et al.* (2007). More recently, an SVM model has been successfully applied
57    to predicting the strength of a mutation library of E coli Trc promoter sequences (Meng,
58    et al., 2017). One drawback with an SVM or ANN machine learning model is the 'black-
59    box' approach; i.e, the absence of any mechanistic insights that could be gleaned with
60    respect to the relationship between promoter sequence and strength. Such an
61    understanding could be vital in the prediction of promoter strengths in different
62    contexts, as well as the forward design of promoters in finely-tuned genetic circuits (for
63    e.g, see Endy, 2005; De Mey, et al. 2007; Salis, et al 2009; Li and Zhang, 2014). Many
64    freely available resources predict the location of promoters in a genomic sequence
65    mainly by identifying the −10 and −35 regulatory sequences (for e.g, de Jong *et al.*
66    (2012)), but very few tools are available to predict the strength of such sequences. One
67    tool provides qualitative predictions ('strong' or not) of promoter strength based on the
68    occurrence of a triad pattern (Dekhtyar et al., 2008), and is available as a macro. Here
69    we present a two-step approach to the predictive modelling of the strength of $\sigma^{70}$ core
70    promoters, and a companion web-based platform and Python standalone tool that
71    implement our method along with the option to dynamically include user data into the
72    prediction model. Our implementation is the first freely available tool/web-server for
73    the quantitative prediction of promoter strength.

74 **METHODS**

75

76 ***Generative model of promoter sequences***. A generative model of the $-10$ and $-35$
77 promoter sequences is constructed using two Position Weight Matrices ($PWM_{-10}$ and
78 $PWM_{-35}$) in the following manner. A comprehensive set of $\sigma^{70}$-binding promoter
79 sequences was extracted from the RegulonDB (Gama-Castro *et al.*, 2016). For each
80 promoter sequence, we extracted a $-35$ region of 13 nucleotides centered at $-35$
81 position, and a $-10$ region of 13 nucleotides centered at the $-10$ position, to allow for
82 uncertainties in the precise position of occurrence of the hexamers. For each $-35$ region,
83 we used FIMO (Grant *et al.*, 2011) to find the best match to the consensus $-35$ motif,
84 and similarly for the $-10$ regions, to obtain a dataset of $-35$ and $-10$ hexamer
85 sequences. This dataset was then filtered for only significant hits to the consensi motifs
86 ($p$-value $< 0.05$) and the resulting dataset was used to determine the weights of each
87 nucleotide at each position of the $-35$ and $-10$ hexamers. Nucleotide-wise counts at
88 each position of the hexamer motifs were augmented by a pseudo-count prior to correct
89 for *E. coli* GC content of 50.8% and the resulting frequency matrices were converted into
90 log-odds matrices. Biopython routines ([www.biopython.org](www.biopython.org)) were used.

91

92 ***Linear modelling of promoter strength***. Following Berg and von Hippel (1987),
93 we modelled the relationship between the promoter sequences and the *ln* of the
94 promoter strength using multiple linear regression. The training set of 18 promoters is
95 drawn from the Anderson library of activator-independent plasmid *tet* promoter
96 variants maintained at the Registry of standard biological parts
97 ([http://parts.igem.org/Promoters/Catalog/Anderson](http://parts.igem.org/Promoters/Catalog/Anderson)). Each promoter sequence is
98 scored with respect to the generative models of the $-10$ and $-35$ motifs (i.e., the $PWM_{-10}$
99 and $PWM_{-35}$ matrices) and the two scores obtained formed the feature space of the
100 regression modelling. The regression coefficients to be determined represent the
101 weights of the -10 and -35 regions in the regression analysis. The Anderson library
102 provided promoter strengths spanning two orders of magnitude and normalized in the
103 range 0.00 to 1.00 with respect to the strongest (i.e, reference) promoter. It was noted
104 that the normalisation step would not affect a linear relationship, altering only the
105 constant of the regression. The normalised strength values were log-transformed to
106 obtain the required response variable values. Since the *ln* function rapidly descends
107 towards $-$ Inf with decreasing promoter strength, we capped the infimum of promoter
108 strength at 0.0001 prior to log-transformation. The least-squares cost function was
109 minimized using iterative gradient descent. The model parameters were assessed using
110 t-statistics, and the overall model was assessed using F-statistic and the adjusted
111 multiple coefficient of determination given by:

112 Adj. $R^2 = 1 - \{(1-R^2)*[(n-1)/(n-m-1)]\}$ ...(1)

113 where m is the number of features and n is the number of instances. The adjustment is a
114 penalty for increasing model complexity.

115 ***Model validation.*** The model of promoter strength was validated in three ways:

116 (i) The model was validated using leave-one-out cross-validation (LOOCV) .

117 (ii) We generated 100 sets of 100 randomly generated promoter sequences each, using
118 the `sample` function in Python. From the obtained sampling distribution of mean
119 strengths of random promoter sequences, we calculated the estimate of the true mean
120 strength of a random promoter sequence, together with its standard error.

121 (iii) We further validated our model on independent datasets of characterized
122 promoters available in Davis *et al.* (2011), Dekhtyar *et al,*(2008), and Dayton *et al,*
123 (1984) .

124 **RESULTS**

125 The entire datasets of 1004 $-35$ hexamers and 1046 $-10$ hexamers parsed out of
126 RegulonDB are available as Supplementary Information. The conservation profiles of
127 the extracted $-35$ and $-10$ hexamer sequences of the promoters in the RegulonDB were
128 visualized and shown in Fig. 1.  Based on these PWMs, the site scores of each promoter
129 sequence in the Anderson library were regressed on the corresponding ln of the
130 promoter strength. A summary of this process with the training data, log-
131 transformation of the promoter strength and predicted response values is presented in
132 Table 1. The modelling process converged within $10^5$ iterations by tuning the gradient
133 descent  to a learning rate ($\alpha$) of 0.015, and the following model was obtained:

134 *ln* (promoter strength) = -5.1046 + 0.4271*(PWM$_{-35}$) + 0.2726*(PWM$_{-10}$)  ...(2)

135 We derived an independent solution of the multiple regression using R (www.r-
136 project.org) and obtained a correlation coefficient of 0.998 between the fitted values of
137 the two models.  The interval estimates of the coefficients of the regression were
138 computed in R using `confint(fit, level=0.95)`, and obtained the following 95%
139 confidence intervals:

140

141 `Intercept :        (-6.4974449, -3.7118421)`

142 `PWM_35    :        (0.2445358, 0.6095848)`

143 `PWM_10    :        (0.1434939, 0.4017307)`

144 The interval estimates did not include zero, and this implied that the coefficients were
145 significant at the 0.05 level. In fact, all the three estimates were significant at a p-value
146 of 1E-3. The F-statistic of the overall regression was significant at a p-value of 2E-4 and
147 adj. $R^2$ was $\approx$ 0.65. The plane of best fit corresponding to the above model is visualized
148 in Fig. 2.

149 The model was then cross-validated using a 18-fold LOOCV (similar to jack-knife).
150 Cross-validation yielded a correlation coefficient of ~0.76 (Table 2). We sought to
151 benchmark our model on a negative test set by generating random −35 and −10
152 hexamer sequences. To this end, we applied our model to 100 sets of 100 random
153 promoter sequences each (available in Supplementary Information) and estimated the
154 true mean of the sampling distribution as 0.00055. The standard error of the estimate
155 was 1.04E-7. The low predicted strength along with the very small standard error
156 indicated that the model predicted these instances to be non-promoter sequences with
157 good certainty. This affirmed the specificity of our model for true promoters.

158 To validate our model further on true promoter sequences and experimentally
159 characterized promoter strengths, we used datasets available in the literature and
160 compared the predicted strength with the experimental results and examined their
161 concordance. The following results were obtained:

162 (i) For the 10 promoters discussed by Sauer and colleagues (2011), we ranked the
163 promoters in Table 1 of the same reference according to their strengths and observed a
164 1000-fold span of promoter strengths, 1E-3 to 1 (Table 3). Promoters 2 and 3 were
165 identically strong, hence we took the average of their predicted strengths in ranking the
166 promoters. With this arrangement, we found that the predicted order of promoters in
167 terms of strength exactly reproduced the experimentally characterized order. Despite
168 the fact that Anderson library and these promoters were characterized and normalized
169 using different systems, the model was able to predict surprisingly well across a
170 promoter strength spectrum spanning three orders of magnitude.

171 (ii) Next, we applied our model to the set of 13 strong promoter candidates of *T.*
172 *maritima* discussed in Dekhtyar *et al*, (2008). Using the hexamer sequences provided in
173 Fig. 5 of the same reference , we applied our model and obtained quantitative
174 predictions of promoter strengths (Table 4). Almost all the promoters had predicted
175 strengths > 0.38 and promoters with canonical hexamers even had strengths > 1.00.
176 One promoter (TM0032) was predicted as 'weak' with a strength ~0.056 and seemed to
177 point to an apparent anomaly in the relationship between promoter sequence and
178 strength, possibly highlighting the need for further experimentation on this promoter.
179 Our observations were corroborated by Fig. 4 in the same reference that showed the
180 least and greatly reduced expression from this particular promoter. These results taken
181 in conjunction with the results on random promoter sequences affirmed the ability of

182 our model to discriminate between promoters at opposite ends of the strength
183 spectrum.

184 (iii) We also applied our model on the five promoters discussed in Dayton *et al*, (1984).
185 Of these, the first three are known as "major" promoters that are active even at low
186 concentrations of the polymerase, whereas the last two are "minor", less strong
187 promoters that are only active when the polymerase is present at high concentrations.
188 We applied our model on the promoter sequences found in Fig. 5 of the same reference
189 and found the predictions in line with the nature of these promoters (Table 5). The
190 activity of the least strong "major" promoter is about two times more than the activity of
191 the strongest "minor" promoter. Hence our modelling approach was able to
192 discriminate between major and minor promoters.

193

194 **DISCUSSION**

195 In addition to the independent contributions of $-35$ and $-10$ sites to promoter strength,
196 we were interested in exploring if any interactions between them could contribute to
197 promoter strength. To this end, we examined the following model in R:

198 `lm(logStrength ~ PWM35 * PWM10)`

199 where `PWM35` and `PWM10` represent the corrresponding site scores. This model
200 resulted in a lower adj. $R^2$ value than that without any interactions. Further, the p-value
201 of the $PWM_{10}$ score dropped below significance (0.31), and the interaction term turned
202 out to be totally insignificant (p-value: 0.97), thus discounting any interaction between
203 the sites in the present dataset. On this basis, the null hypothesis of absence of any
204 interaction could not be rejected, and we concluded that there is little evidence for
205 interaction between the $-35$ and $-10$ sites in contributing to promoter strength.

206 Our model assumed that both the predictors carried independent information about the
207 promoter strength, and together they are able to provide sufficient information about
208 the strength. The basis of this assumption was probed to determine if both predictors
209 are necessary to the model. Could one predictor provide sufficient information about the
210 promoter strength in the absence of the other? There are at least three angles to address
211 this question, and all of them were considered to interpret the model better.

212 (1) Comparing the raw, unadjusted $R^2$ with the adjusted $R^2$. The corresponding values
213 were:

214 $R^2 \approx 0.69$

215 Adj. $R^2 \approx 0.65$

216    Since there is not much difference between $R^2$ and adj. $R^2$, we could say that both
217    predictors contribute substantially to the response variable (promoter strength) and
218    account for about 65% of its variance.

219    (2) Since the p-values of both predictors are significant, it would be interesting to
220    observe their effect on the response variable in more detail. This was performed using
221    the `effects` package in `R`:

222    `library(effects)`

223    `fit = lm(logStrength~ PWM35+ PWM10, data)`

224    `plot(allEffects(fit))`

225    The results are shown in Fig. 3 where the PWM scores are plotted against the level of
226    confidence in the predicted response. Confidence in the effect of −35 site increases with
227    the score from 0 to about 7, and then is susceptible to edge effects as the score reaches 8.
228    Confidence in the effect of the −10 site increases with the score from -4 to about 5, and
229    then is susceptible to edge effects as the score reaches 10.

230    (3) Another way to address the question is to compute the correlation coefficients
231    between all the variables of interest, including a variable with the combined effects of −
232    35 and −10 sites. This is shown in Table 6. Three features were used, namely $PWM_{-10}$
233    score, $PWM_{-35}$ score, and the combined score (i.e., $PWM_{-10} + PWM_{-35}$). These feature
234    variables were correlated with two response variables, namely promoter strength and its
235    corresponding log transformation. It was first observed that the $PWM_{-10}$ and $PWM_{-35}$
236    scores were anti-correlated with each other (correlation coefficient = -0.37), thus
237    supporting the hypothesis that they are two independent features that could compensate
238    for each other in determining promoter strength. It was significant that the each feature
239    was better correlated with the log of the strength than the strength itself. We tried to
240    regress the strength on the PWM scores, but the model had a very low adj. $R^2$ ($\approx 0.40$)
241    and the intercept term was not significant at the 0.05 level. Further, the highest
242    correlation between the features and response variable was observed between the
243    combined score and log of the promoter strength (~0.79), but the combined score
244    showed only a moderate correlation with the promoter strength prior to log
245    transformation (~0.63). This was in keeping with similar observations for the strength
246    of $\sigma^E$ promoters (Rhodius and Mutalik, 2010). and underscored the logarithmic
247    dependence between the promoter strength and sequence.

248    Finally, the assumptions of linear modelling were investigated with reference to our
249    problem. Model diagnostics of four basic assumptions were plotted (shown in Fig. 4).
250    Specifically:

251   Plot A: The residuals were plotted against the fitted values. No trend was visible in the
252   plot, indicating the residuals did not increase with the fitted values and followed a
253   random pattern about zero. This validated the assumption that the errors were
254   independent.

255   Plot B: The square root of the relative error (standardized residual) was plotted against
256   the fitted value. An almost flat trend was observed, indicating that the standardized
257   residual did not vary with the fitted value.  This further validated the assumption that
258   the errors were independent.

259   Plot C: To test the assumption that the errors were normally distributed, the
260   standardized residuals were plotted against the theoretical quantiles of a normal
261   distribution. The residual distribution closely followed the theoretical quantiles, except
262   for minor deviations towards the tails of the distribution. .

263   Plot D: Since the least-squares cost function is sensitive to outliers, the number of
264   outliers should be kept to a minimum. This was investigated by plotting the
265   standardized residual against the corresponding instance's model leverage. This plot
266   showed that there were no significant outliers in the dataset that could exert an undue
267   influence on the regression parameters.

268   An alternative univariate regression model using only the combined score of the PWMs
269   found the coefficient of regression and the F-statistic significant (both p-values $\approx 10^{-4}$).
270   However, the adj. $R^2$ of the model ($\approx 0.59$) was much lower than that for eq. (2), so the
271   original multiple linear regression model was retained for the estimation of the
272   promoter strength.

273   In summary, our model performed equally well on datasets of strong promoter
274   sequences and datasets of weak random promoter sequences. Our model was consistent
275   in detecting promoter strengths across a 1000-fold span of promoter strengths in *E. coli*
276   as well as the promoter strengths of a different species, *T. maritima*. The model was
277   further able to discriminate between the major and minor promoters of bacteriophage
278   T7.

279   Based on these results, an open-access open-source web server and standalone tool
280   offering the prediction service have been implemented .  Since the linear modelling
281   results are dependent on the dataset, our implementation provides a facility to augment
282   the learning based on user-provided inputs.  The web interface is based on Python web
283   module (web.py) and nginx server. The computational layer is based on numpy,
284   Biopython and matplotlib. The user is provided with an option to add any number of
285   promoter instances with −10 and −35 sequences and the corresponding strengths to
286   augment the training data of the supervised model. The measurement of promoter
287   strength could be done in the manner of Kelly, et al. (2009), where the GFP (reporter

288  gene) synthesis rate is measured per unit biomass, and this could be normalized relative
289  to the reference promoter. In order to assess the goodness of fit of the updated model,
290  the R-squared value is re-computed, along with the 3D plot of the regression surface.
291  This would enable the user to decide whether the data added to the model has improved
292  its performance for further experiments with the software. Based on the trained model,
293  the user could predict the strength of an uncharacterised promoter given its −10 and −
294  35 hexamers.

295  **CONCLUSION**

296  The following important conclusions were drawn from our study. (1) Sequence-based
297  modelling yielded a non-linear, logarithmic dependence between promoter strength and
298  sequence. (2) The model was able to discriminate equally well between strong/major
299  promoters and weak/minor/random promoter sequences, indicating successful learning
300  of the essential features of promoter strength prediction. (3) The combined score
301  ($PWM_{-35} + PWM_{-10}$) emerged as the single most important predictor of the promoter
302  strength. Our model yielded robust quantitative prediction across a 1000-fold span of
303  promoter strengths. It is straighforward to extend our methodology to the study of new
304  promoter classes of other σ factors. Our implementation and web service could be useful
305  in characterizing promoters identified in genome sequencing projects as well in
306  engineering promoters for the design of finely-tuned genetic circuits in synthetic
307  biology. The dynamic feature of our implementation would enable users to incorporate
308  their own data into the model and obtain more reliable estimates of promoter strength.
309  The service will be periodically updated based on the availability of new training
310  instances, user input data and/or models for promoters of other σ factors.
311

312  **Acknowledgments**

316

317  **REFERENCES**
318

319  Benos PV1, Bulyk ML, Stormo GD. (2002) Additivity in protein-DNA interactions: how
320  good an approximation is it? Nucleic Acids Res. **30**(20):4442-51.

321

322  Basu,R.S., Warner, B.A, Molodtsov, V, Pupov, D, Esyunina, D, Fernández-Tornero, C,
323  Kulbachinskiy, A, and Murakami,K.S. (2014) Structural Basis of Transcription Initiation
324  by Bacterial RNAPolymerase Holoenzyme. J Biol Chem **289**: 24549 −24559

325  Berg, O.G. and von Hippel, P.H. (1987). Selection of DNA binding sites by regulatory
326  proteins. Statistical-mechanical theory and application to operators and promoters. J
327  Mol Biol **193**:723–750.

328

329  Bujard, H. (1980) The interaction of E.coli RNA polymerase with promoters. Trends
330  Biochem Sci 5, 274-278.
331
332  Crooks GE, Hon G, Chandonia JM, Brenner SE (2004) WebLogo: A sequence logo
333  generator, *Genome Research*, **14**:1188-1190

334  Davis, JH, Rubin, AJ and Sauer, RT. (2011) Design, construction and characterization of
335  a set of insulated bacterial promoters. Nucleic Acids Res. **39**(3): 1131–1141.).
336
337  Dayton, CJ, Prosen, DE, Parker, KL and Cech, CL (1984). Kinetic measurements of
338  Escherichia coli RNA polymerase association with bacteriophage T7 early promoters. J
339  Biol Chem **259**: 1616

340  de Jong, A., Pietersma H, Cordes M, Kuipers OP, Kok J.(2012) PePPER: a webserver for
341  prediction of prokaryote promoter elements and regulons. BMC Genomics **13**:299

342  De Mey, M, Lequeux, GJ, Soetaert, WK, and Vandamme, EJ. (2008) Construction and
343  model-based analysis of a promoter library for E. coli: an indispensable tool for
344  metabolic engineering BMC Biotechnol. **7**: 34.
345
346  Dekhtyar, M, Morin, A and Sakanyan, V. (2008) Triad pattern algorithm for predicting
347  strong promoter candidates in bacterial genomes. BMC Bioinformatics **9**:233
348
349  Deuschle, U., Kammerer, W. Gentz, R. & Bujard, H. (1986). Promoters of Escherichia
350  coli: a hierarchy of in vivo strength indicates alternate structures. EMBO J **5**: 2987-2994
351

352  Endy, D. (2005) Foundations for engineering biology. Nature **438**:449.

353  Feklistov, A. and Darst, S.A. (2011) Structural Basis for Promoter −10 Element
354  Recognition by the Bacterial RNA Polymerase σ Subunit. Cell **147**: 1257–1269

355  Galas, D.J., Eggert, M. & Waterman, M.S. (1985). Rigorous pattern-recognition methods
356  for DNA sequences. Analysis of promoter sequences from Escherichia coli. J Molec Biol
357  **186**: 117-128
358
359  Gama-Castro S, Salgado H, Santos-Zavaleta A, Ledezma-Tejeida D, Muñiz-Rascado L,
360  García-Sotelo JS, Alquicira-Hernández K, Martínez-Flores I, Pannier L, Castro-
361  Mondragón JA, Medina-Rivera A, Solano-Lira H, Bonavides-Martínez C, Pérez-Rueda
362  E, Alquicira-Hernández S, Porrón-Sotelo L, López-Fuentes A, Hernández-Koutoucheva
363  A, Del Moral-Chávez V, Rinaldi F, Collado-Vides J. (2016) RegulonDB version 9.0: high-

364 level integration of gene regulation, coexpression, motif clustering and beyond. Nucleic
365 Acids Res. **44**(D1):D133-43. doi: 10.1093/nar/gkv1156.

367 Grant, CE, Bailey, TL and Noble WS (2011) FIMO: Scanning for occurrences of a given
368 motif. Bioinformatics **27**(7):1017–1018.

370 Hawley, DK and McClure, WR (1983) Compilation and analysis of Escherichia coli
371 promoter DNA sequences. Nucl. Acids Res. **11**: 2237

373 Hook-Barnard, I. , Johnson XB, Hinton DM. (2006) *Escherichia coli* RNA Polymerase
374 Recognition of a σ70-Dependent Promoter Requiring a −35 DNA Element and an
375 Extended −10 TGn Motif. J Bacteriol. **188**:8352–8359.

376 Huerta AM1, Collado-Vides J. (2003) Sigma70 promoters in Escherichia coli: specific
377 transcription in dense regions of overlapping promoter-like signals. J Mol Biol.
378 **333**(2):261-78.

380 Kadonaga, J.T. (2012) Perspectives on the RNA Polymerase II Core Promoter. Wiley
381 interdisciplinary reviews Developmental biology. **1**:40-51.

382 Kelly, J.R., Rubin AJ, Davis JH, Ajo-Franklin CM, Cumbers J, Czar MJ, de Mora K,
383 Glieberman AL, Monie DD, Endy D. (2009). Measuring the activity of BioBrick
384 promoters using an in vivo reference standard. J Biol Eng **3**:4.

386 Knaus and Bujard (1990) 'Principles Governing the Activity of E. coli Promoters". In:
387 Eckstein F., Lilley D.M.J. (eds) Nucleic Acids and Molecular Biology, vol. 4. Berlin:
388 Springer-Verlag.

390 Li J and Zhang Y. (2014) Relationship between promoter sequence and its strength in
391 gene expression. Eur Phys J E Soft Matter **37**(9):44.

393 Maquat, LE and Reznikoff, WS. (1978) In vitro analysis of the Escherichia coli RNA
394 polymerase interaction with wild-type and mutant lactose promoters. J Mol Biol 125:
395 467.

397 Meng, H., Ma, Y., Mai, G., Wang, Y., Liu, C.(2017) Construction of precise support
398 vector machine based models for predicting promoter strength. Quant Biol **5**: 90.
399 https://doi.org/10.1007/s40484-017-0096-3

402 Paget, M.S. and Helmann, J.D. (2003). The σ70 family of sigma factors. Genome Biology
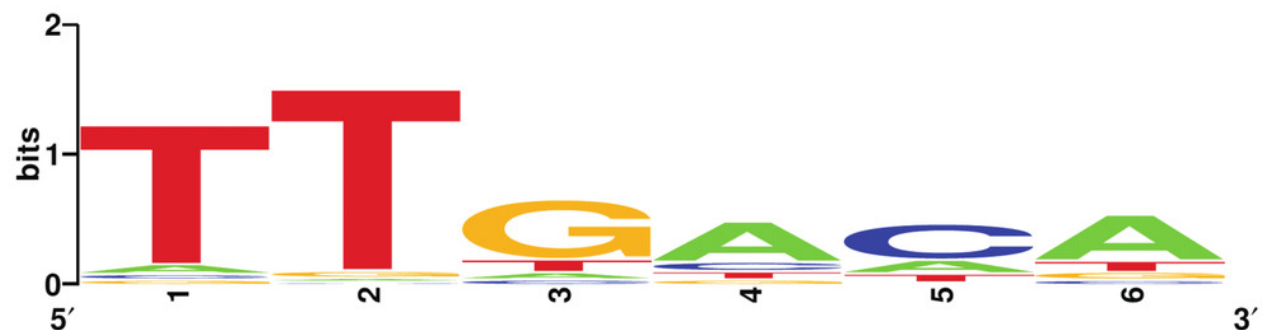403 **4**:203.

404 Rhodius, V.A. and Mutalik, V.K. (2010) Predicting strength and function for promoters
405 of the *Escherichia coli* alternate sigma factor, σE. Proc. Natl. Acad. Sci. USA **107**: 2854-
406 2859

407  Salis HM1, Mirsky EA, Voigt CA. (2009) Automated design of synthetic ribosome
408  binding sites to control protein expression. Nat Biotechnol. **27**(10):946-50. doi:
409  10.1038/nbt.1568.
410
411  Shultzaberger, R.K.,  Chen Z, Lewis KA, Schneider TD. (2007) Anatomy of *Escherichia*
412  *coli* sigma70 promoters. Nucleic Acids Res **35**:771–788.
413
414  Stormo, G.D. (1990). Consensus patterns in DNA. In: Methods in Enzymology, Vol. 183.
415  Molecular evolution: Computer analysis of protein and nucleic acid sequences.
416  (Doolittle, R.F., ed.) San Diego: Academic Press.
417
418  Weller K and Recknagel RD. (1994) Promoter strength prediction based on occurrence
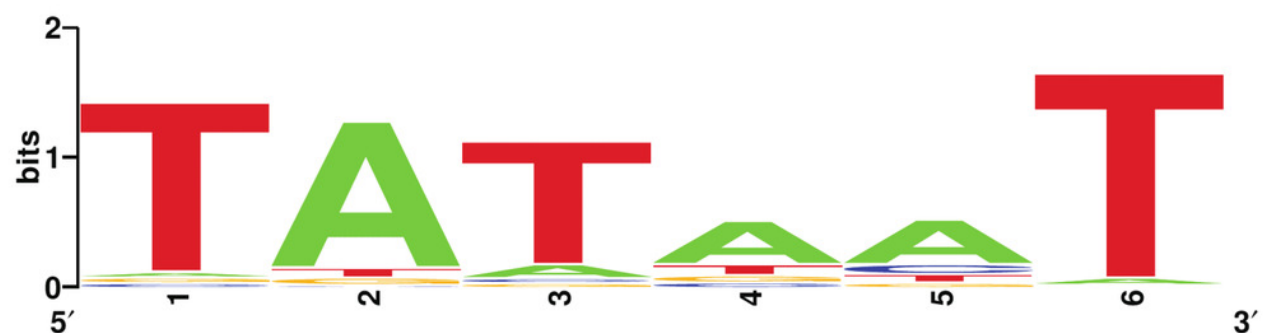419  frequencies of consensus patterns. J Theor Biol. **171**(4):355-9.
420
421

# Figure 1

Sequence logos of the –35 and –10 hexamers of the selected RegulonDB promoters.

Figure was made using WebLogo (Crooks *et al.*, 2004).



(A) −35 motif



(B) −10 motif

**Figure 2**(on next page)

The regression surface of the estimated model with the training data points(red).

X- and y-axes represent PWM scores and the z-axis (vertical) represents the predicted *ln*(promoter strength).

# Figure 3
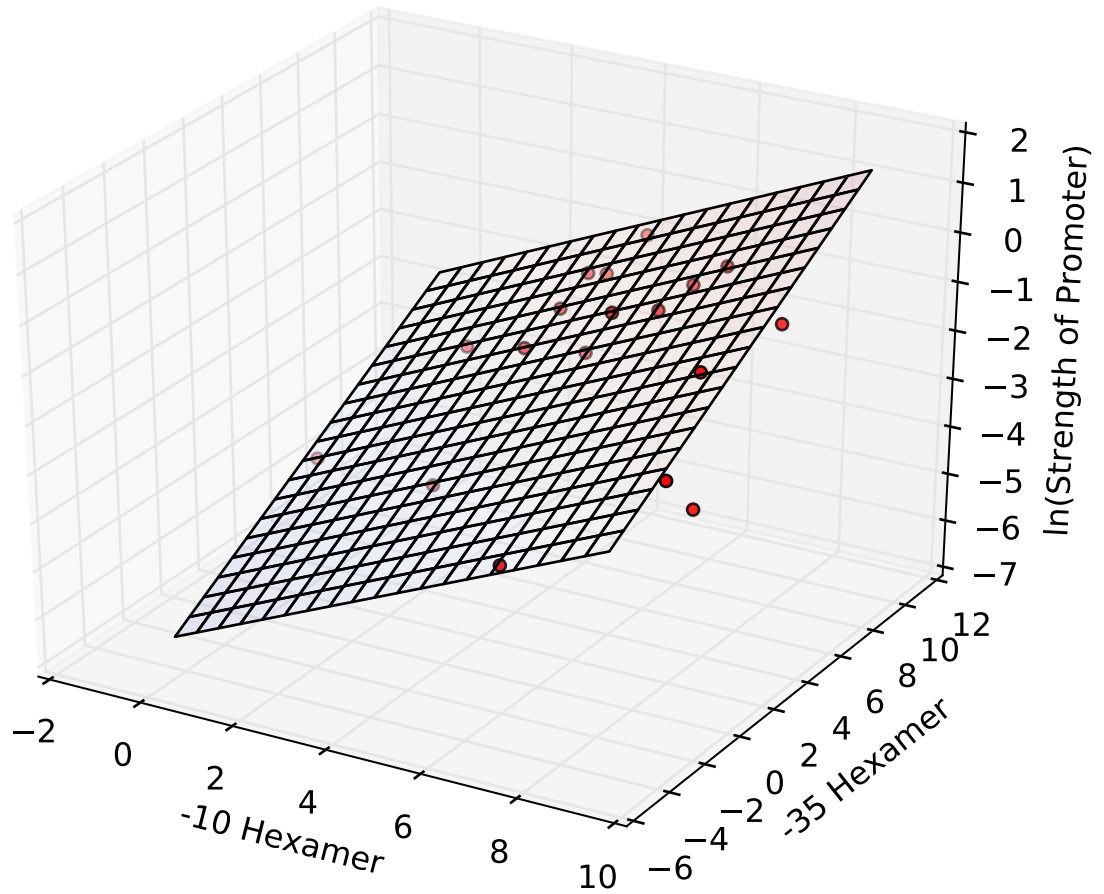
Effects plots of –35 and –10 promoter sites on promoter strength.



(A) PWM$_{35}$ effect plot
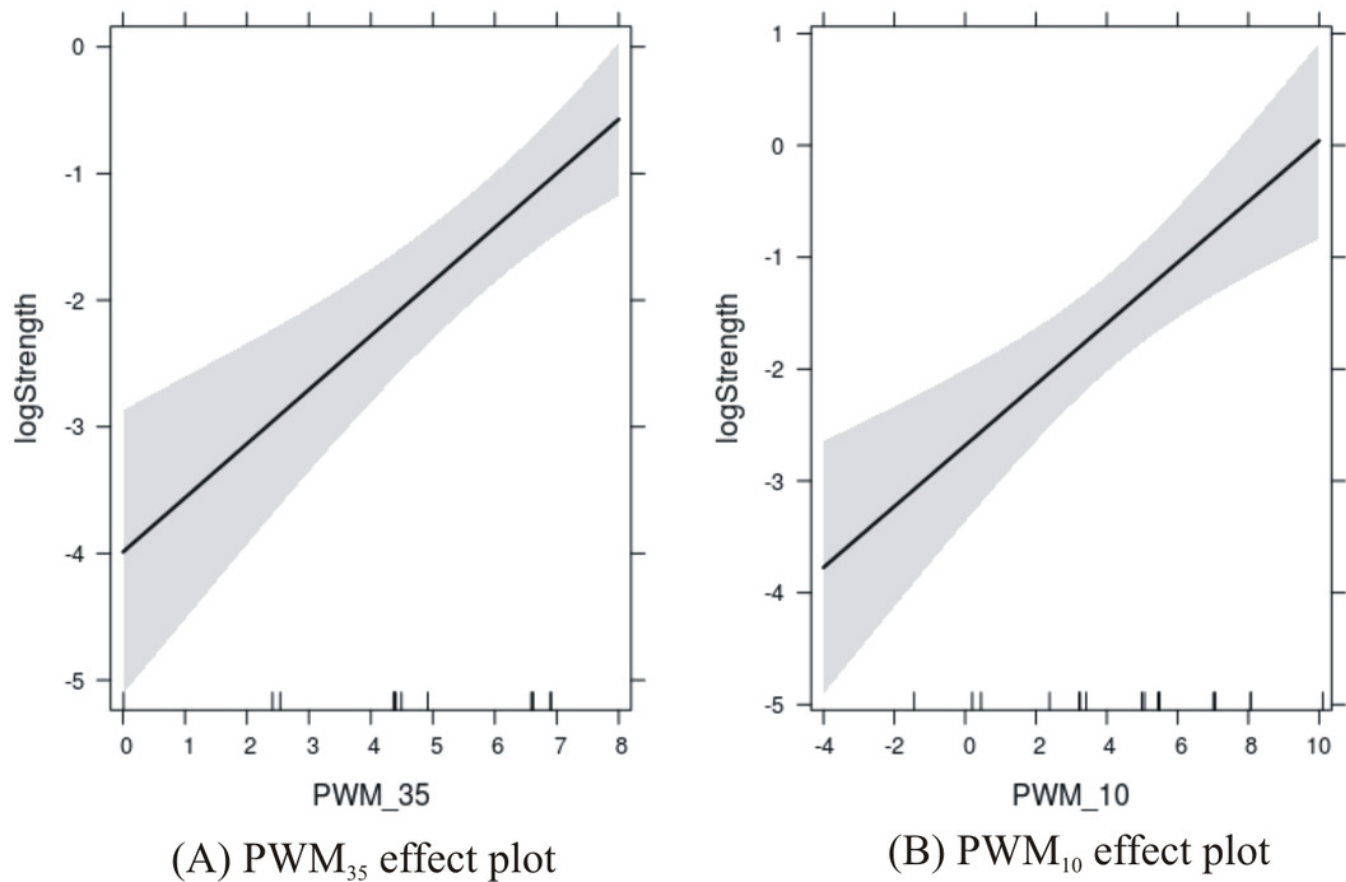
(B) PWM$_{10}$ effect plot

# Figure 4
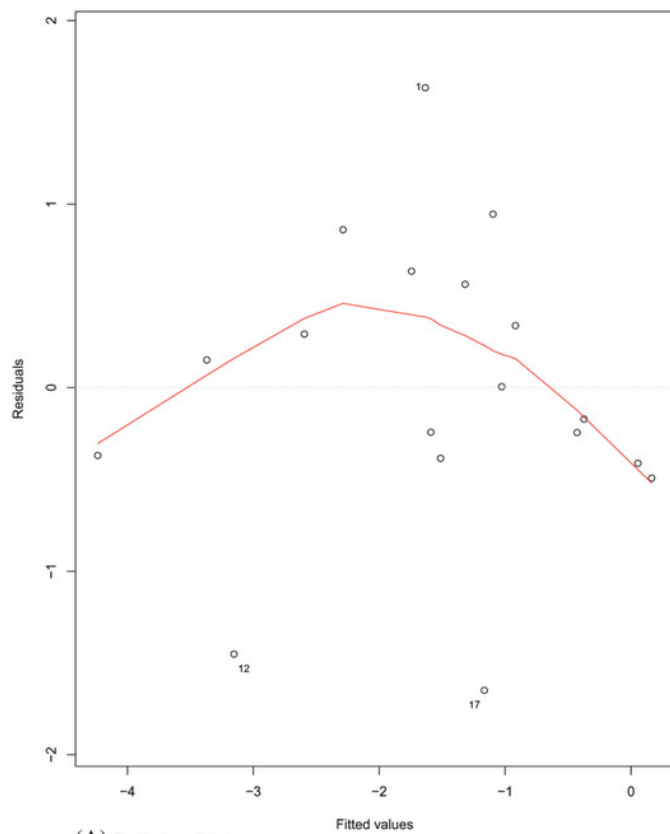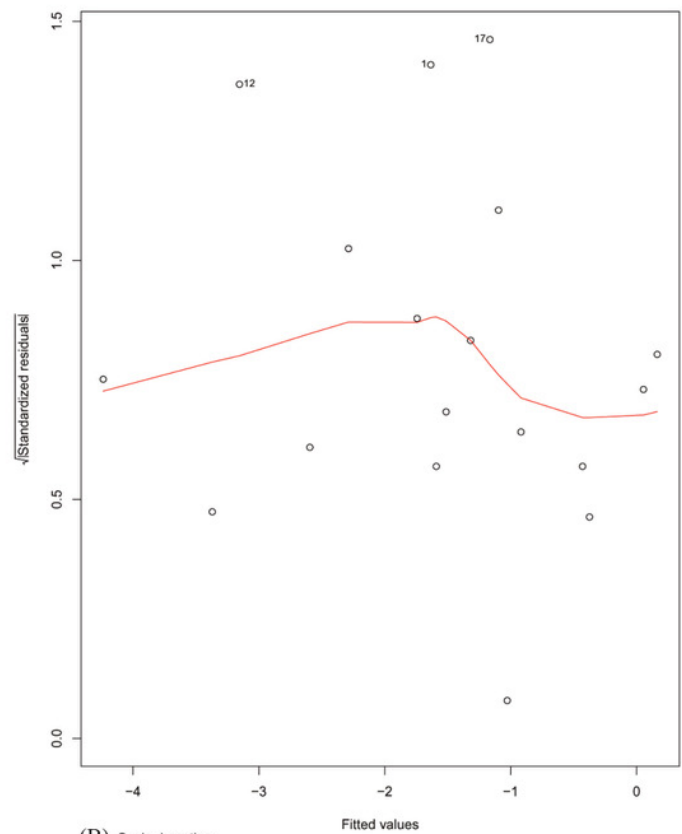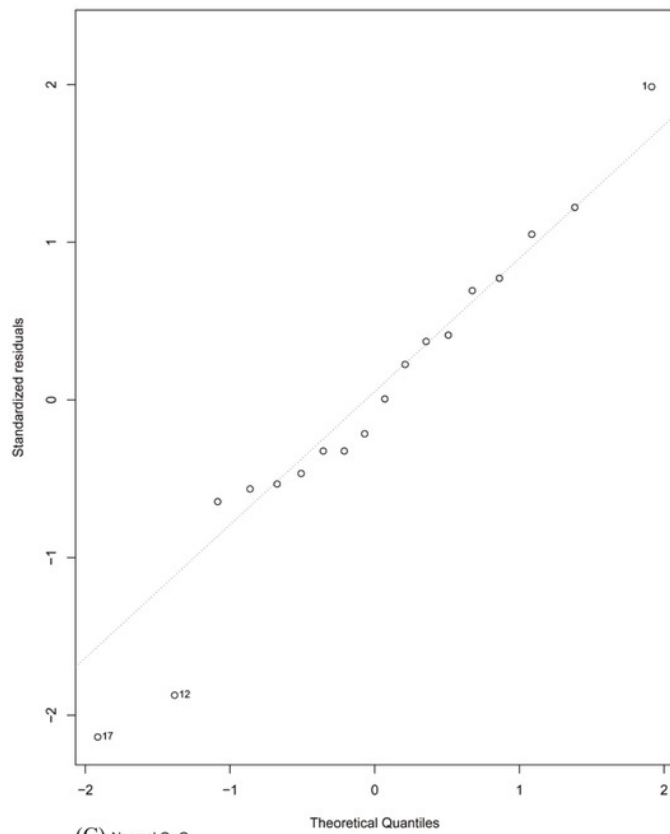
Model diagnostics plots for investigating the assumptions underlying linear modelling.
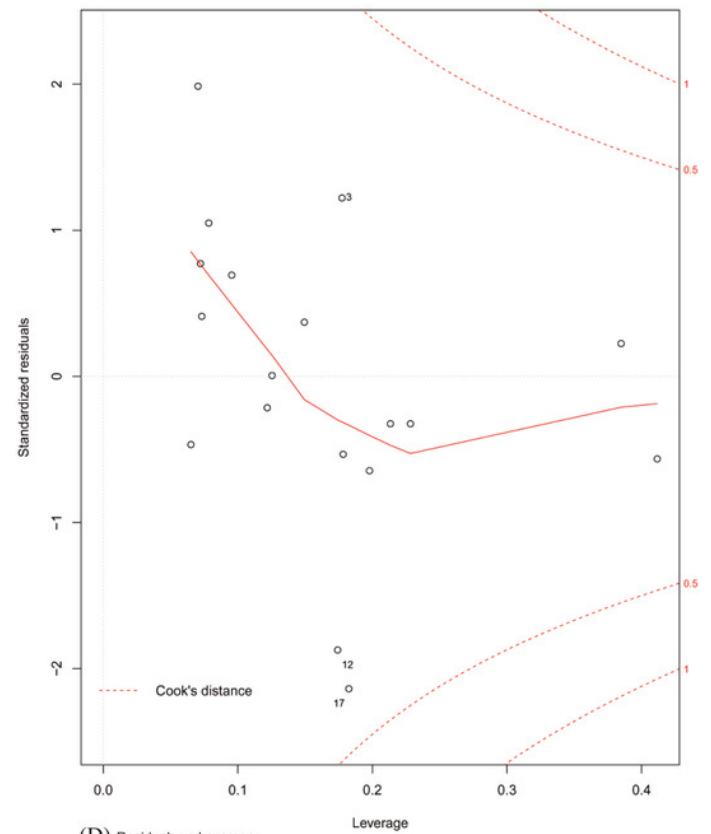
(A) Residuals vs Fitted

(B) Scale–Location

(C) Normal Q-Q

(D) Residuals vs Leverage

**Table 1**(on next page)

Summary of promoter information.

The promoter activities (strengths) are seen to span two orders of magnitude in the range [0.0, 1.0]. The promoters follow the naming in the Anderson dataset.

| Promoter | -35 hexamer | -10 hexamer | Promoter Activity | ln(Promoter Activity) | Predicted ln(Promoter Activity) |
|---|---|---|---|---|---|
| BBa_J23100 | TTGACG | TACAGT | 1 | 0 | -1.6336486579 |
| BBa_J23101 | TTTACA | TATTAT | 0.7 | -0.35667494 | 0.0555718065 |
| BBa_J23102 | TTGACA | TACTGT | 0.86 | -0.15082289 | -1.0957849491 |
| BBa_J23104 | TTGACA | TATTGT | 0.72 | -0.32850407 | 0.1647181133 |
| BBa_J23105 | TTTACG | TACTAT | 0.24 | -1.42711636 | -2.2871659092 |
| BBa_J23106 | TTTACG | TATAGT | 0.47 | -0.75502258 | -1.3174788735 |
| BBa_J23107 | TTTACG | TATTAT | 0.36 | -1.02165125 | -1.0266628468 |
| BBa_J23108 | CTGACA | TATAAT | 0.51 | -0.67334455 | -0.4282477098 |
| BBa_J23109 | TTTACA | GACTGT | 0.04 | -3.21887582 | -3.3693144659 |
| BBa_J23110 | TTTAGG | TACAAT | 0.33 | -1.10866262 | -3.3946866337 |
| BBa_J23111 | TTGACG | TATAGT | 0.58 | -0.54472718 | -0.3731455955 |
| BBa_J23112 | CTGATA | GATTAT | 0.01 | -4.60517019 | -3.1533888284 |
| BBa_J23113 | CTGATG | GATTAT | 0.01 | -4.60517019 | -4.2356234817 |
| BBa_J23114 | TTTATG | TACAAT | 0.1 | -2.30258509 | -2.5943689001 |
| BBa_J23115 | TTTATA | TACAAT | 0.15 | -1.89711998 | -1.5121342469 |
| BBa_J23116 | TTGACA | GACTAT | 0.16 | -1.83258146 | -1.5897942167 |
| BBa_J23117 | TTGACA | GATTGT | 0.06 | -2.81341072 | -1.1644781255 |
| BBa_J23118 | TTGACG | TATTGT | 0.56 | -0.5798185 | -0.91751654 |

1

PeerJ

**Table 2**(on next page)

Cross-validation results.

In each fold of cross-validation, the instance corresponding to the fold was designated as the test instance while the prediction model was built using the rest of the instances. This process was repeated 18 times, once for each test instance and the cross-validation (CV) residuals were obtained. combined, sum of the PWM scores; cvpred, predicted log strength of the test instance; cvres, cross-validation residual.

| Fold | PWM_35 | PWM_10 | combined | logStrength | cvpred | cvres |
|------|--------|--------|----------|-------------|--------|-------|
| 1 | 6.5966 | 2.398 | 9 | 0 | -1.757 | 1.757 |
| 2 | 6.9195 | 8.089 | 15.01 | -0.357 | 0.145 | -0.50 |
| 3 | 9.1308 | 0.402 | 9.53 | -0.151 | -1.3 | 1.15 |
| 4 | 9.1308 | 5.025 | 14.16 | -0.329 | 0.286 | -0.62 |
| 5 | 4.3854 | 3.465 | 7.85 | -1.427 | -2.36 | 0.93 |
| 6 | 4.3854 | 7.022 | 11.41 | -0.755 | -1.377 | 0.62 |
| 7 | 4.3854 | 8.089 | 12.47 | -1.022 | -1.027 | 0.00 |
| 8 | 4.5119 | 10.086 | 14.6 | -0.673 | -0.362 | -0.31 |
| 9 | 6.9195 | -4.474 | 2.45 | -3.219 | -3.463 | 0.24 |
| 10 | 4.3854 | 5.462 | 9.85 | -1.109 | -1.792 | 0.68 |
| 11 | 6.5966 | 7.022 | 13.62 | -0.545 | -0.349 | -0.20 |
| 12 | 2.5179 | 3.213 | 5.73 | -4.605 | -2.847 | -1.76 |
| 13 | -0.0162 | 3.213 | 3.2 | -4.605 | -3.977 | -0.63 |
| 14 | 2.3914 | 5.462 | 7.85 | -2.303 | -2.646 | 0.34 |
| 15 | 4.9255 | 5.462 | 10.39 | -1.897 | -1.485 | -0.41 |
| 16 | 9.1308 | -1.411 | 7.72 | -1.833 | -1.518 | -0.32 |
| 17 | 9.1308 | 0.15 | 9.28 | -2.813 | -0.796 | -2.02 |
| 18 | 6.5966 | 5.025 | 11.62 | -0.58 | -0.944 | 0.36 |

1

PeerJ

**Table 3**(on next page)

Validation results: using data of Davis *et al.*, (2011).

The promoters were ordered based on the rank of their strength, and given as input to our model. The predicted promoter log strengths were then examined for agreement with the actual rank and the ordering obtained matched the original ordering. The individual predicted values for pro2 and pro3 were 0.0024 and 0.059, respectively.

| Actual rank | Promoter | -35 sequence | -10 sequence | Strength | Predicted exp(logStrength) | Predicted rank |
|---|---|---|---|---|---|---|
| 1 | pro1 | tttacg | gtatct | 0.009 | 0.0079073845 | 1 |
| 2.5 | pro2 | gcggtg | tataat | 0.017 | 0.0306978849 | 2.5 |
| 2.5 | pro3 | ttgacg | gaggat | 0.017 | 0.0306978849 | 2.5 |
| 4 | proA | tttacg | taggct | 0.03 | 0.0482647297 | 4 |
| 5 | pro4 | tttacg | gatgat | 0.033 | 0.0809816409 | 5 |
| 6 | pro5 | tttacg | taggat | 0.05 | 0.0867400443 | 6 |
| 7 | proB | tttacg | taatat | 0.119 | 0.1534857959 | 7 |
| 8 | pro6 | tttacg | taaaat | 0.193 | 0.2645364297 | 8 |
| 9 | proC | tttacg | tatgat | 0.278 | 0.3059490889 | 9 |
| 10 | proD | tttacg | tataat | 1 | 0.6173668247 | 10 |

1

**Table 4**(on next page)

Validation with *T. maritima* strong promoter candidates.

| Promoter | -35 sequence | -10 sequence | Strength | Predicted exp(logStrength) | Predicted class |
|---|---|---|---|---|---|
| TM0373 | ttgaca | tataat | Strong | 4.6845788997 | Strong |
| TM1016 | ttgaat | tttaat | Strong | 0.3808572257 | Strong |
| TM1272 | ttgaca | tttaat | Strong | 1.6386551999 | Strong |
| TM1429 | ttgaca | tataat | Strong | 4.6845788997 | Strong |
| TM1667 | ttgaaa | tataat | Strong | 2.5859432664 | Strong |
| TM1780 | ttcata | tataat | Strong | 0.463878289 | Strong |
| Tmt11 | ttgaat | taaaat | Strong | 0.4665383797 | Strong |
| TM0032 | tcgaaa | cataat | Strong | 0.0562167049 | *Weak* |
| TM0477 | ttgaat | tataat | Strong | 1.0887926414 | Strong |
| TM1067 | ttgacc | tattat | Strong | 0.7046782664 | Strong |
| TM1271 | ttgaca | tataat | Strong | 4.6845788997 | Strong |
| Tmt45 | ttgaac | tataat | Strong | 0.670434893 | Strong |
| TM1490 | ttgact | taaaat | Strong | 0.8451600149 | Strong |

1

**Table 5**(on next page)

Validation with major (A1, A2, A3) and minor (C, D) promoters.

| Promoter | -35 sequence | -10 sequence | Strength | Predicted exp(logStrength) | Predicted class |
|---|---|---|---|---|---|
| A1 | ttgact | gatact | strong | 0.2904988307 | medium |
| A2 | ttgaca | taagat | strong | 0.9947607331 | strong |
| A3 | ttgaca | tacgat | strong | 0.658183377 | strong |
| C | ttgacg | tagtct | minor | 0.1452865585 | minor |
| D | ttgact | taggct | minor | 0.1541996302 | minor |

1

<image_start>ok<image_end>

**Table 6**(on next page)

Correlation matrix of features and response variables.

1   **Table 2.** Correlation matrix of features and response variables.

| Corr. Coef. | $PWM_{-35}$ | $PWM_{-10}$ | Combined | Strength | Log-strength |
|---|---|---|---|---|---|
| $PWM_{-35}$ | 1 | -0.3715610 | 0.3401672 | 0.4558838 | 0.5153622 |
| $PWM_{-10}$ | -0.3715610 | 1 | 0.7466500 | 0.3025062 | 0.4115533 |
| Combined | 0.3401672 | 0.7466500 | 1 | 0.6330488 | 0.7861173 |
| Strength | 0.4558838 | 0.3025062 | 0.6330488 | 1 | 0.8665495 |
| Log-strength | 0.5153622 | 0.4115533 | 0.7861173 | 0.8665495 | 1 |

2