# 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 σ⁷⁰ factor of RNA polymerase. σ⁷⁰ 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 σ⁷⁰-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.

39  **INTRODUCTION**
40

41  The primary *E. coli* promoter-specificity factor and the one widely used in recombinant
42  DNA technology is the $\sigma^{70}$ factor. Promoters recognized by $\sigma^{70}$-containing RNA
43  polymerase are called core promoters and share the following features: two conserved
44  hexamer sequences, separated by a non-specific spacer of ideally 17 nucleotides. The two
45  hexamers are located ~ 35 bp and ~10 bp upstream of the transcription start site, and
46  are called the −35 and −10 sequences respectively (Maquat and Reznikoff, 1978; Bujard,
47  1980; Paget and Helmann, 2003; Kadonaga, 2012). −35 and −10 sequences matching
48  the consensi motifs (TTGACA and TATAAT, respectively) are known as canonical
49  hexamers  (Galas, et al. 1985; Deuschle, et al. 1986; Stormo, 1990). It is known that the
50  conserved hexamer regions are vital for recognizing and optimizing the interactions
51  between DNA and the RNA polymerase (Hawley and McClure, 1983; Knaus and Bujard,
52  1990; Hook-Barnard *et al*., 2006; Feklistov and Darst, 2011; Basu *et al*., 2014).

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

79 **METHODS**

80

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

96

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

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

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

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

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

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

125     (iii) We further validated our model on independent datasets of characterized
126     promoters available in Davis *et al.* (2011), Dekhtyar *et al,*(2008), and Dayton *et al,*
127     (1984) .

128     **RESULTS**

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

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

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

144

145

146

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

148     `PWM_35    :        (0.2445358, 0.6095848)`

149     `PWM_10    :        (0.1434939, 0.4017307)`

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

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

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

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

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

188    our model to discriminate between promoters at opposite ends of the strength
189    spectrum.

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

199

200    **DISCUSSION**

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

204    ```
lm(logStrength ~ PWM35 * PWM10)
```

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

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

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

220    $R^2 \approx 0.69$

221    Adj. $R^2 \approx 0.65$

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

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

228  ```
     library(effects)
     ```

229  ```
     fit = lm(logStrength~ PWM35+ PWM10, data)
     ```

230  ```
     plot(allEffects(fit))
     ```

231  The results are shown in Fig. 3 where the PWM scores are plotted against the level of
232  confidence in the predicted response. Confidence in the effect of $-35$ site increases with
233  the score from 0 to about 7, and then is susceptible to edge effects as the score reaches 8.
234  Confidence in the effect of the $-10$ site increases with the score from -4 to about 5, and
235  then is susceptible to edge effects as the score reaches 10.

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

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

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

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

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

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

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

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

285  Based on these successes, a web server offering the prediction service has been
286  implemented free to everyone.  Since the linear modelling results are dependent on the
287  dataset, our implementation provides a facility to augment the learning based on user-
288  provided inputs.  The web interface is based on Python web module (web.py) and nginx
289  server. The computational layer is based on numpy, Biopython and matplotlib. The user
290  is provided with an option to add any number of promoter instances with −10 and −35
291  sequences and the corresponding strengths to augment the training data of the
292  supervised model. The measurement of promoter strength could be done in the manner
293  of Kelly, et al. (2009), where the GFP synthesis rate is measured per unit biomass, and

294   this could be normalized relative to the reference promoter. In order to assess the
295   goodness of fit of the updated model, the R-squared value is re-computed, along with a
296   3D plot of the regression surface. This would enable the user to decide whether the data
297   added to the model has improved its performance for further experiments with the
298   software. Based on the trained model, the user could predict the strength of any
299   uncharacterised promoter given its −10 and −35 hexamers.

300   **CONCLUSION**

301   The following important conclusions were drawn from our study. (1) Sequence-based
302   modelling yielded a non-linear, logarithmic dependence between promoter strength and
303   sequence. (2) The model was able to discriminate equally well between strong/major
304   promoters and weak/minor/random promoter sequences, indicating successful learning
305   of the essential features of promoter strength prediction. (3) The combined sum of the
306   scores ($PWM_{-35} + PWM_{-10}$) emerged as the single most important predictor of the
307   promoter strength. Our model yielded robust quantitative prediction across a 1000-fold
308   span of promoter strengths. It is straighforward to extend our methodology to the study
309   of new promoter classes of other σ factors. Our implementation and web service could
310   be useful in characterizing unknown promoters of newly sequenced genomes as well in
311   the engineering of promoters for designing finely-tuned genetic circuits in synthetic
312   biology. The dynamic feature of our implementation would enable users with own data
313   to obtain more reliable estimates of promoter strength. The service will be periodically
314   updated based on the availability of new training instances, user input data and/or
315   models for promoters of other σ factors.
316

322   Supplementary Information:

323   The webserver is available at https://promoterpredict.com. The downloadable software
324   is available at https://github.com/PromoterPredict. Further supporting information is
325   available online at https://doi.org/10.6084/m9.figshare.6794939.v1.

326   **REFERENCES**
327

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

330

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

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

337

338  Bujard, H. (1980) The interaction of E.coli RNA polymerase with promoters. Trends
339  Biochem Sci 5, 274-278.
340
341  Crooks GE, Hon G, Chandonia JM, Brenner SE (2004) WebLogo: A sequence logo
342  generator, *Genome Research*, **14**:1188-1190

343  Davis, JH, Rubin, AJ and Sauer, RT. (2011) Design, construction and characterization of
344  a set of insulated bacterial promoters. Nucleic Acids Res. **39**(3): 1131−1141.).
345
346  Dayton, CJ, Prosen, DE, Parker, KL and Cech, CL (1984). Kinetic measurements of
347  Escherichia coli RNA polymerase association with bacteriophage T7 early promoters. J
348  Biol Chem **259**: 1616

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

351  De Mey, M, Lequeux, GJ, Soetaert, WK, and Vandamme, EJ. (2008) Construction and
352  model-based analysis of a promoter library for E. coli: an indispensable tool for
353  metabolic engineering BMC Biotechnol. **7**: 34.
354
355  Dekhtyar, M, Morin, A and Sakanyan, V. (2008) Triad pattern algorithm for predicting
356  strong promoter candidates in bacterial genomes. BMC Bioinformatics **9**:233
357
358  Deuschle, U., Kammerer, W. Gentz, R. & Bujard, H. (1986). Promoters of Escherichia
359  coli: a hierarchy of in vivo strength indicates alternate structures. EMBO J **5**: 2987-2994
360

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

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

364  Galas, D.J., Eggert, M. & Waterman, M.S. (1985). Rigorous pattern-recognition methods
365  for DNA sequences. Analysis of promoter sequences from Escherichia coli. J Molec Biol
366  **186**: 117-128
367

368  Gama-Castro S, Salgado H, Santos-Zavaleta A, Ledezma-Tejeida D, Muñiz-Rascado L,
369  García-Sotelo JS, Alquicira-Hernández K, Martínez-Flores I, Pannier L, Castro-
370  Mondragón JA, Medina-Rivera A, Solano-Lira H, Bonavides-Martínez C, Pérez-Rueda
371  E, Alquicira-Hernández S, Porrón-Sotelo L, López-Fuentes A, Hernández-Koutoucheva
372  A, Del Moral-Chávez V, Rinaldi F, Collado-Vides J. (2016) RegulonDB version 9.0: high-
373  level integration of gene regulation, coexpression, motif clustering and beyond. Nucleic
374  Acids Res. **44**(D1):D133-43. doi: 10.1093/nar/gkv1156.
375
376  Grant, CE, Bailey, TL and Noble WS (2011) FIMO: Scanning for occurrences of a given
377  motif. Bioinformatics **27**(7):1017–1018.
378
379  Hawley, DK and McClure, WR (1983) Compilation and analysis of Escherichia coli
380  promoter DNA sequences. Nucl. Acids Res. **11**: 2237
381
382  Hook-Barnard, I. , Johnson XB, Hinton DM. (2006) *Escherichia coli* RNA Polymerase
383  Recognition of a $\sigma^{70}$-Dependent Promoter Requiring a –35 DNA Element and an
384  Extended –10 TGn Motif. J Bacteriol. **188**:8352–8359.

385  Huerta AM1, Collado-Vides J. (2003) Sigma70 promoters in Escherichia coli: specific
386  transcription in dense regions of overlapping promoter-like signals. J Mol Biol.
387  **333**(2):261-78.
388
389  Kadonaga, J.T. (2012) Perspectives on the RNA Polymerase II Core Promoter. Wiley
390  interdisciplinary reviews Developmental biology. **1**:40-51.
391  Kelly, J.R., Rubin AJ, Davis JH, Ajo-Franklin CM, Cumbers J, Czar MJ, de Mora K,
392  Glieberman AL, Monie DD, Endy D. (2009). Measuring the activity of BioBrick
393  promoters using an in vivo reference standard. J Biol Eng **3**:4.
394
395  Knaus and Bujard (1990) 'Principles Governing the Activity of E. coli Promoters". In:
396  Eckstein F., Lilley D.M.J. (eds) Nucleic Acids and Molecular Biology, vol. 4. Berlin:
397  Springer-Verlag.
398
399  Li J and Zhang Y. (2014) Relationship between promoter sequence and its strength in
400  gene expression. Eur Phys J E Soft Matter **37**(9):44.
401
402  Maquat, LE and Reznikoff, WS. (1978) In vitro analysis of the Escherichia coli RNA
403  polymerase interaction with wild-type and mutant lactose promoters. J Mol Biol 125:
404  467.
405
406  Meng, H., Ma, Y., Mai, G. et al. (2017) Construction of precise support vector machine
407  based models for predicting promoter strength. Quant Biol **5**: 90.
408  https://doi.org/10.1007/s40484-017-0096-3
409
410
411

412   Paget, M.S. and Helmann, J.D. (2003). The $\sigma^{70}$ family of sigma factors. Genome Biology
413   **4**:203.

414   Rhodius, V.A. and Mutalik, V.K. (2010) Predicting strength and function for promoters
415   of the *Escherichia coli* alternate sigma factor, $\sigma^E$. Proc. Natl. Acad. Sci. USA **107**: 2854-
416   2859

417   Salis HM1, Mirsky EA, Voigt CA. (2009) Automated design of synthetic ribosome
418   binding sites to control protein expression. Nat Biotechnol. **27**(10):946-50. doi:
419   10.1038/nbt.1568.
420
421   Shultzaberger, R.K.,  Chen Z, Lewis KA, Schneider TD. (2007) Anatomy of *Escherichia*
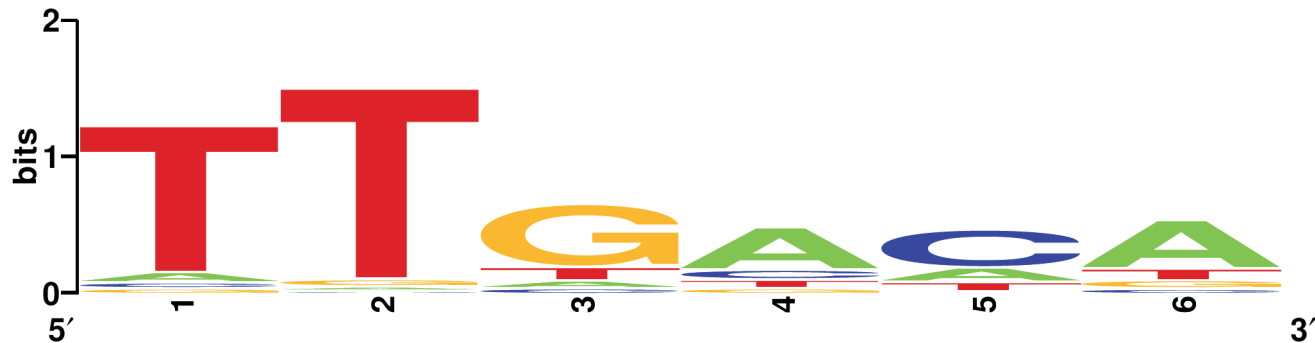422   *coli* sigma70 promoters. Nucleic Acids Res **35**:771–788.
423
424   Stormo, G.D. (1990). Consensus patterns in DNA. In: Methods in Enzymology, Vol. 183.
425   Molecular evolution: Computer analysis of protein and nucleic acid sequences.
426   (Doolittle, R.F., ed.) San Diego: Academic Press.
427
428   Weller K and Recknagel RD. (1994) Promoter strength prediction based on occurrence
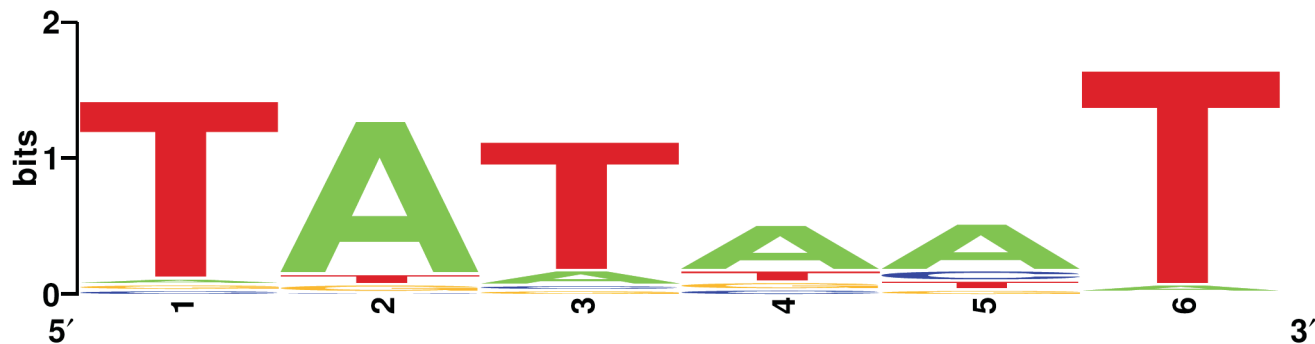429   frequencies of consensus patterns. J Theor Biol. **171**(4):355-9.
430
431

**Figure 1**(on next page)

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).

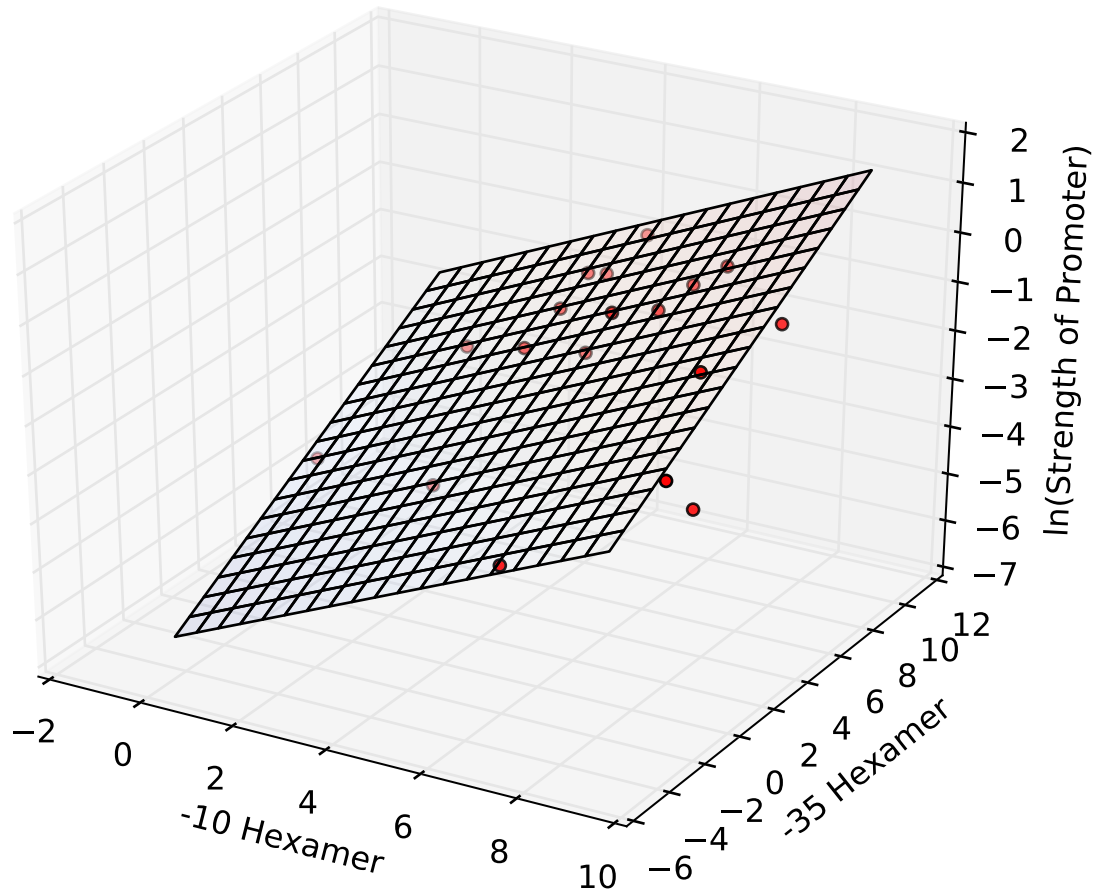X- and y-axes represent PWM scores and the z-axis (vertical) represents the predicted *ln*(promoter strength).

PeerJ

# Figure 3

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



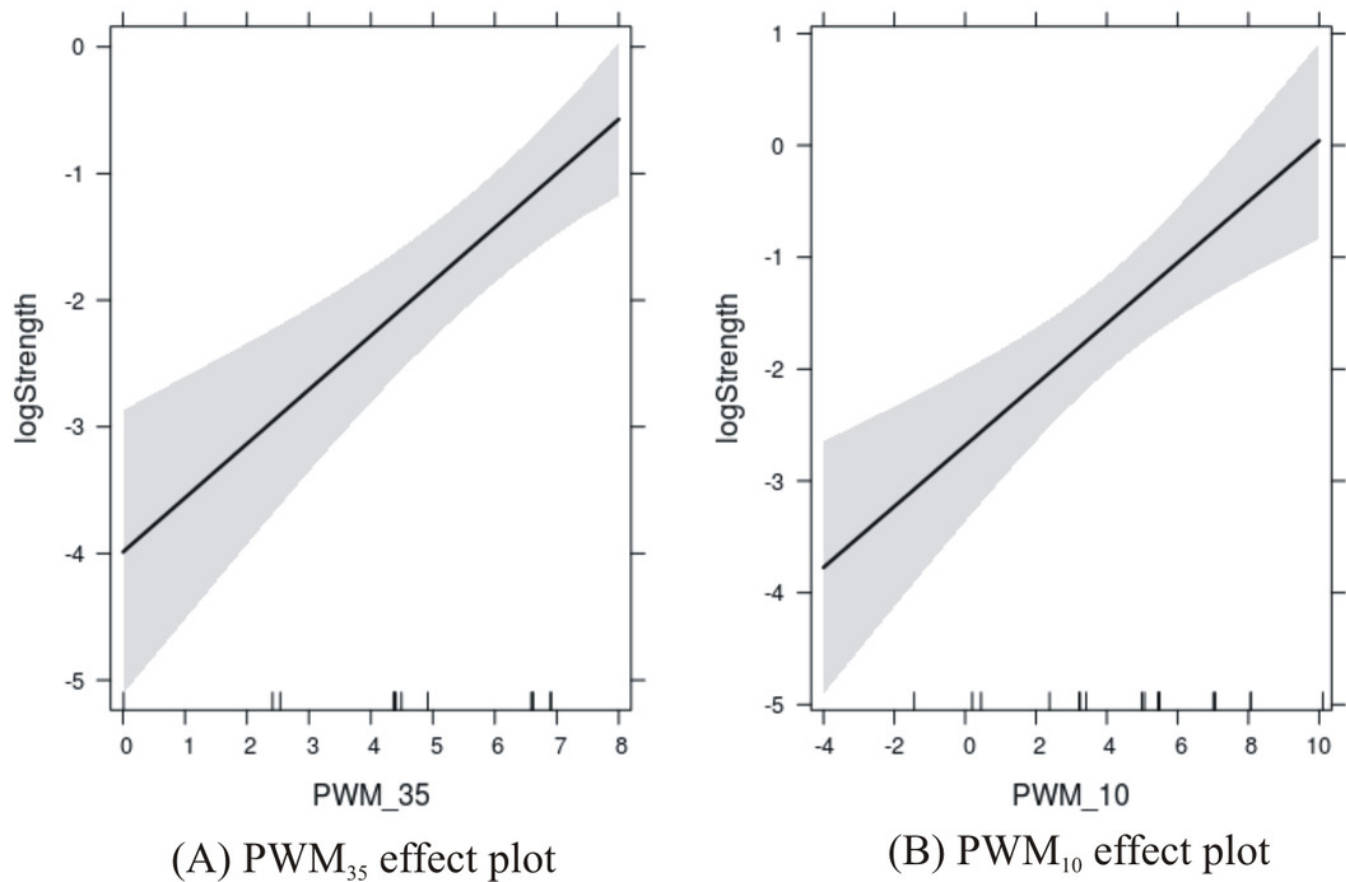(A) PWM$_{35}$ effect plot

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

**Figure 4**(on next page)

Model diagnostics plots for investigating the assumptions underlying linear modelling.

(A) Residuals vs Fitted

(B) Scale−Location

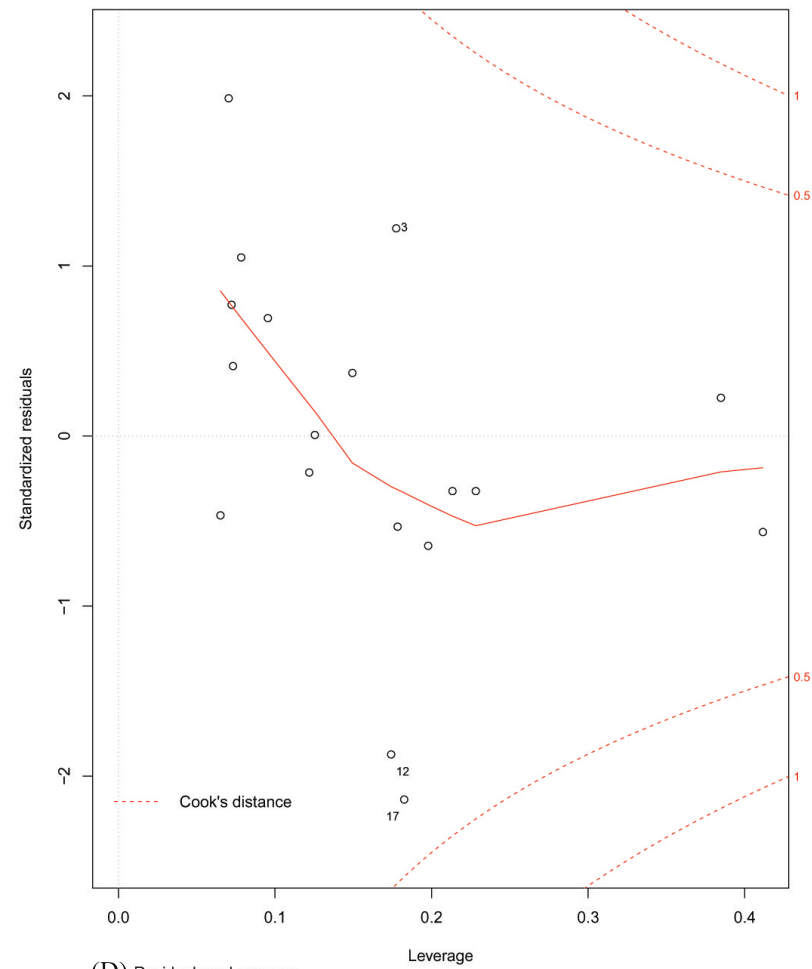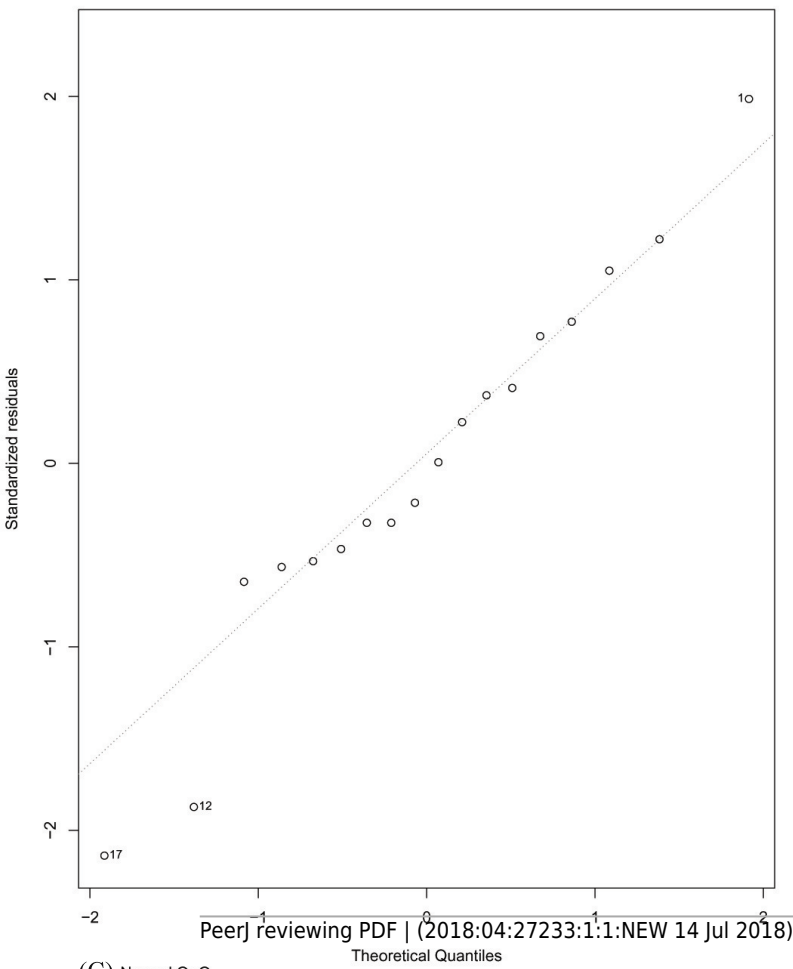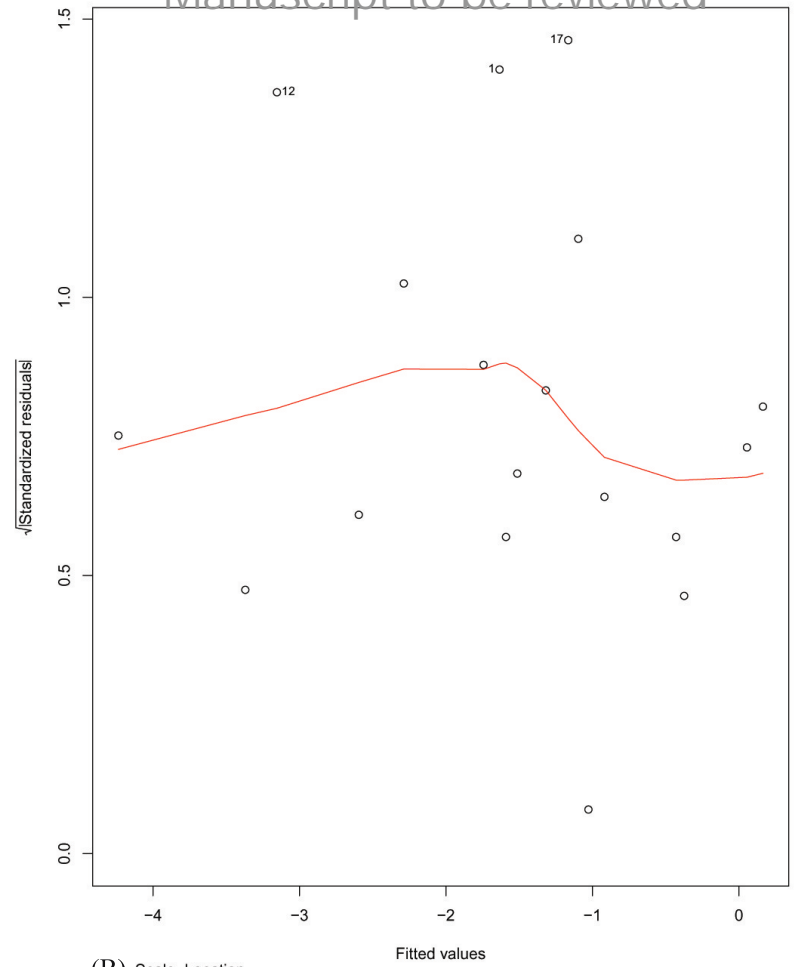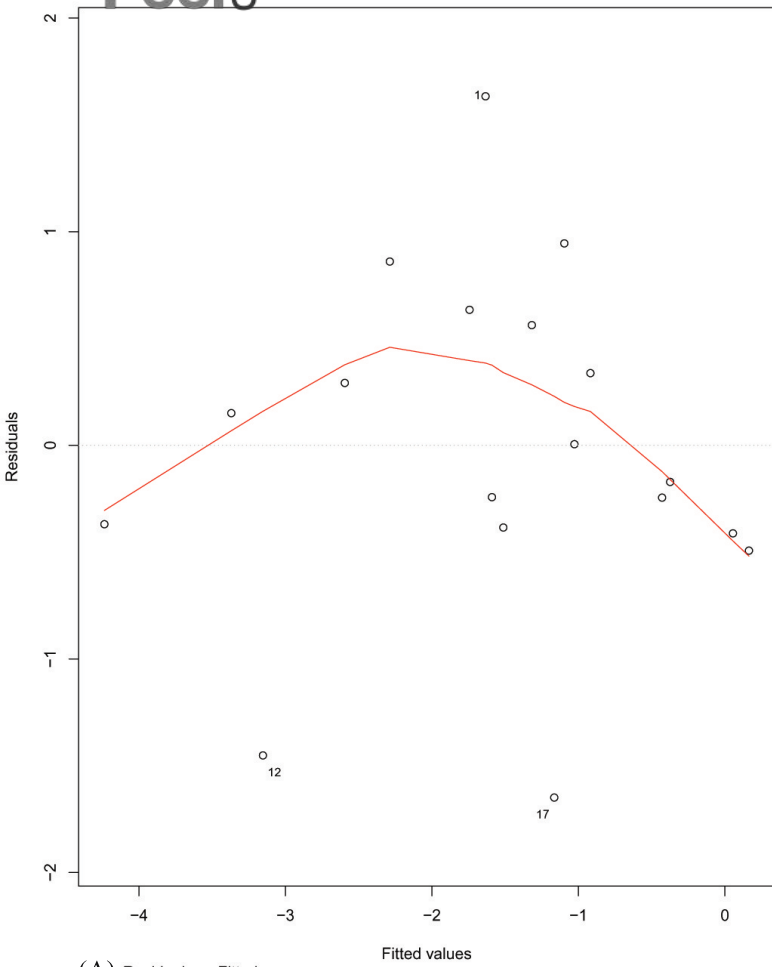(C) Normal Q−Q

(D) Residuals vs Leverage

PeerJ

**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 | -1.4669153 |
| BBa_J23101 | TTTACA | TATTAT | 0.7 | -0.35667494 | -0.25855671 |
| BBa_J23102 | TTGACA | TACTGT | 0.86 | -0.15082289 | -0.62881141 |
| BBa_J23104 | TTGACA | TATTGT | 0.72 | -0.32850407 | -0.22100527 |
| BBa_J23105 | TTTACG | TACTAT | 0.24 | -1.42711636 | -0.80989265 |
| BBa_J23106 | TTTACG | TATAGT | 0.47 | -0.75502258 | -1.50446674 |
| BBa_J23107 | TTTACG | TATTAT | 0.36 | -1.02165125 | -0.40208651 |
| BBa_J23108 | CTGACA | TATAAT | 0.51 | -0.67334455 | -2.31347961 |
| BBa_J23109 | TTTACA | GACTGT | 0.04 | -3.21887582 | -2.06383098 |
| BBa_J23110 | TTTAGG | TACAAT | 0.33 | -1.10866262 | -1.50446674 |
| BBa_J23111 | TTGACG | TATAGT | 0.58 | -0.54472718 | -1.05910916 |
| BBa_J23112 | CTGATA | GATTAT | 0.01 | -4.60517019 | -4.0308767 |
| BBa_J23113 | CTGATG | GATTAT | 0.01 | -4.60517019 | -4.0308767 |
| BBa_J23114 | TTTATG | TACAAT | 0.1 | -2.30258509 | -2.92677594 |
| BBa_J23115 | TTTATA | TACAAT | 0.15 | -1.89711998 | -2.78324614 |
| BBa_J23116 | TTGACA | GACTAT | 0.16 | -1.83258146 | -1.21066725 |
| BBa_J23117 | TTGACA | GATTGT | 0.06 | -2.81341072 | -1.21066725 |
| BBa_J23118 | TTGACG | TATTGT | 0.56 | -0.5798185 | -0.36453507 |

1

**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. cvpred, predicted log strength of the test instance; cvres, cross-validation residual.

| Fold | PWM_35 | PWM_10 | combined | logStrength | cvpred | cv_residual |
|------|--------|--------|----------|-------------|--------|-------------|
| 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

**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 are 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

**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