

# Sequence-Based Prediction of Pathological Mutations

C. Ferrer-Costa,<sup>1</sup> M. Orozco,<sup>1,2\*</sup> and X. de la Cruz,<sup>1,3\*</sup>

<sup>1</sup>Molecular Modeling and Bioinformatics Unit, Institut de Recerca Biomèdica, Parc Científic de Barcelona, Barcelona, Spain

<sup>2</sup>Departament de Bioquímica i Biologia Molecular, Facultat de Química, Universitat de Barcelona, Barcelona, Spain

<sup>3</sup>Institució Catalana per la Recerca i Estudis Avançats (ICREA), Barcelona, Spain

**ABSTRACT** The development of methods to assess the impact of amino acid mutations on human health has become an important goal in biomedical research, due to the growing number of nonsynonymous SNPs identified. Within this context, computational methods constitute a valuable tool, because they can easily process large amounts of mutations and give useful, almost cost-free, information on their pathological character. In this paper we present a computational approach to the prediction of disease-associated amino acid mutations, using only sequence-based information (amino acid properties, evolutionary information, secondary structure and accessibility predictions, and database annotations) and neural networks, as a model building tool. Mutations are predicted to be either pathological or neutral. Our results show that the method has a good overall success rate, 83%, that can reach 95% when trained for specific proteins. The methodology is fast and flexible enough to provide good estimates of the pathological character of large sets of nonsynonymous SNPs, but can also be easily adapted to give more precise predictions for proteins of special biomedical interest. *Proteins* 2004;57:811–819. © 2004 Wiley-Liss, Inc.

**Key words:** SNPs; disease-associated mutations; protein sequence; bioinformatics; neural networks; sequence variability

## INTRODUCTION

Identifying the pathological character of the vast amount of known nonsynonymous SNPs (single nucleotide polymorphisms) has become an important challenge for biomedical sciences.<sup>1,2</sup> This has led to the development of different computer-based methods to predict the pathological character of amino acid mutations.<sup>3–11</sup> Many of these methods were created to provide an initial and rapid identification of disease-associated mutations from large sets of nonsynonymous SNPs, allowing their ranking for subsequent experimental study.<sup>5,7</sup> To this end, amino acid mutations are mapped to the protein sequence and then characterized in terms of structure and sequence properties.<sup>3,4,6,7,9–14</sup> These properties are subsequently used to predict whether a given mutation is likely to be pathological.

The rationale behind these methods is that the disease-causing effect of many mutations can be understood in terms of their effect on protein structure,<sup>6,12,13,15–17</sup> e.g.,

stability losses, disruption of protein interactions, etc. This is supported by the fact that some structure-based properties, like solvent accessibility, may have a high discriminating power between neutral and disease-associated mutations.<sup>7,12,13</sup> Unfortunately, the rapidly growing number of proteins for which no structure is available,<sup>18</sup> seriously limits the applicability of structure-based methods. One possible way to overcome this problem is the use of evolutionary information alone, derived from multiple sequence alignments.<sup>4,11,14</sup> Some authors have fruitfully tested this idea in small sets of specific systems: Miller and Kumar,<sup>14</sup> studying a set of six proteins and their mutations, have recently shown that disease mutations tend to be overabundant at highly conserved positions in multiple sequence alignments; Santibáñez-Koref et al.<sup>11</sup> have used evolutionary information to assess the significance of missense mutations in the case of P53. On the contrary, Ng and Henikoff<sup>19</sup> use many different proteins to test their prediction method, also based on the use of evolutionary information, with promising results. However, overall error rates are still high (see Table 4 in Ng and Henikoff<sup>19</sup>), leaving room for improvement.

An interesting alternative has been recently proposed<sup>7</sup> where mutations are characterized using evolutionary information together with structure properties obtained from de novo protein structure predictions. Preliminary results are encouraging,<sup>7</sup> showing that the recognition power of evolutionary information can be extended by the use of predicted structure properties.<sup>7</sup> However, this method is limited by the fact that, for the moment, computer-costly de novo predictions can only be obtained routinely for small protein domains<sup>20</sup> and cannot be applied to membrane proteins. Interestingly, in the case of mutations affecting viability in microbial systems, Krishnan and Westhead<sup>9</sup> have proposed that predicted structural properties could be used to complement evolutionary properties.

In the present article we extend the idea of combining evolutionary information together with structure informa-

*Abbreviations:* SNP, single nucleotide polymorphism; DAMU, disease-associated mutation; NEMU, neutral mutation; NN, neural network.

\*Correspondence to: Xavier de la Cruz or Modesto Orozco. Parc Científic de Barcelona, C/Josep Samitier, 1-5, 08028 Barcelona, Spain. E-mail: xavier@mmb.pcb.ub.es or, modesto@mmb.pcb.ub.es.

Received 30 January 2004; Accepted 1 June 2004

Published online 10 August 2004 in Wiley InterScience (www.interscience.wiley.com). DOI: 10.1002/prot.20252

tion for the prediction of human disease-associated mutations (DAMUs). We propose a prediction method in which mutations are labelled using only sequence-derived information. The structure properties utilized (secondary structure and accessibility) are obtained using standard, fast, sequence-based methods. In addition to evolutionary and structure information, residue properties and database annotations are also used. All this information is input to a neural network (NN) that provides the final prediction on whether the target mutation is pathological, or is a neutral mutation (NEMU). The resulting predictive model has a high overall success rate, 83.5%, which represents a 66.5% improvement over a random prediction approach. We find that evolutionary information provides a substantial amount of the discrimination power of the method, although the contribution of other parameters is non-negligible. Our results indicate that the method can be applied to the quick characterization and ranking of mutations, supporting its future use for the prioritization of nonsynonymous SNPs, within the context of biomedical studies.

## MATERIALS AND METHODS

### Prediction Protocol

The goal of our work is to devise a computational method to predict whether amino acid mutations are pathological or not. To this end, mutations are labelled with a set of properties (size change, location in structure, etc.) that can be related to their possible damaging effect to the protein. The resulting set of properties is then used by a NN to decide if the mutation is pathological or neutral.

### Sets of Mutations

#### Disease-associated mutations

There are two main repositories of human disease-associated mutations (DAMUs) at present:<sup>21</sup> OMIM<sup>22</sup> and SwissProt.<sup>23</sup> We chose to work with SwissProt, a manually curated database, because recent work<sup>21</sup> indicates that it provides a more accurate mapping of the mutations onto the protein sequence. This is particularly relevant in cases like ours, where several thousands of mutations are utilized, and automatic processing is required. In addition, DAMUs from SwissProt have been utilized by authors working in the study and prediction of pathological mutations.<sup>5,7,13,19</sup> Finally, DAMUs from SwissProt and those from OMIM have very similar properties at the protein structure and sequence levels, e.g., we find that 39.1% of DAMUs happen at buried residues, and Steward et al.,<sup>17</sup> who use OMIM data, find 39%; we also find that DAMUs are more frequent at conserved sites, in agreement with Steward et al.<sup>17</sup>

We utilized version 40 of SwissProt,<sup>23</sup> and obtained the set of DAMUs following Ng and Henikoff.<sup>19</sup> To this end, we searched SwissProt using as keywords DISEASE, VARIANT, and HUMAN. We discarded mutations for which no clear link to disease could be established. Finally, mutations for which only the human sequence was available in the Pfam<sup>24</sup> alignment for the protein family were also discarded (more than one sequence is required to compute

some evolutionary properties used in the prediction process, see below). A total of 9334 DAMUs, happening in 811 human proteins were kept.

In the structure-based case, we used a set of mutations that could be mapped to the protein structure. This gave a total of 1319 DAMUs happening in 90 proteins.

### Neutral mutations

Two models of neutral mutations (NEMUs) have been used in this work. The first was derived from a massive mutagenesis experiment done in *Escherichia coli* Lac repressor protein,<sup>25</sup> keeping those mutations that result in organisms phenotypically similar to the wild-type bacteria.<sup>4</sup> We have utilized this dataset because it has been used by Ng and Henikoff<sup>19</sup> to test their sequence-based method, thus allowing a direct comparison between both approaches.

In the second model for NEMUs, which we call the evolutionary model, mutations were collected only for the 811 human proteins for which DAMUs were available. For these proteins, NEMUs corresponded to those variants occurring in members of the protein family from other species.<sup>12,13</sup> To this end, we took the Pfam<sup>24</sup> multiple sequence alignment for the protein family, and eliminated: (1) all the human sequences but that of the human target protein; (2) nonhuman sequences with less than 95% identity to the human target sequence. Any amino acid change between the target human sequence and the remaining sequences was then considered as a neutral mutation. A total of 11372 NEMUs were obtained.

For the structure-based method we used a subset of 888 NEMUs from the evolutionary model that mapped into the protein structure.

### Selection of Structural Data for Proteins

Protein structures were downloaded from the PDB.<sup>26</sup>

### Sequence and Structure-Based Properties

In the sequence-based version of our procedure we used 19 parameters to label mutations; 23 were used in the structure-based version.

It will be noted that some of the parameters are clearly related, e.g., the two measures of residue size. While dimensionality reduction is advisable in order to gain generalization performance,<sup>27</sup> our results [Table I(A)] show that good generalization is already attained with our NN models. In addition, it may be advisable to keep a larger input in order to avoid losing useful parameters.<sup>28</sup> For this reason we decided to work on the full model and then analyze separately the contribution of the different parameters (see Results and Discussion), in order to provide a better understanding of what makes our model successful.

### Structure properties

Two structure-related parameters, secondary structure and solvent accessibility, were utilized in the sequence-based version of our method, in which no experimental three-dimensional information whatsoever was used. They

**TABLE I. Performance of Our Approach for Predicting Disease-Associated Mutations<sup>†</sup>**

(A)	Qtot <sup>a</sup>	S <sup>b</sup>	FP <sup>c</sup>	FN <sup>d</sup>
Ng and Henikoff <sup>e</sup>	70.3	37.0	14.7	30.5
SWP-Lac <sup>f</sup>	87.3 ( $\pm$ 0.3)	54.2 ( $\pm$ 2.5)	10.4 ( $\pm$ 0.6)	4.6 ( $\pm$ 0.5)
SWP-Evol <sup>g</sup>	83.5 ( $\pm$ 0.3)	66.5 ( $\pm$ 0.5)	16.5 ( $\pm$ 0.6)	21.1 ( $\pm$ 0.6)
PDBst <sup>h</sup>	87.0 ( $\pm$ 1.6)	73.0 ( $\pm$ 3.1)	10.4 ( $\pm$ 2.1)	11.5 ( $\pm$ 1.3)
(B)	Qtot <sup>a</sup>	S <sup>b</sup>		
Pssm <sup>i</sup>	77.6 ( $\pm$ 0.5)	49.1 ( $\pm$ 1.1)		
$\Delta$ Pssm <sup>j</sup>	73.1 ( $\pm$ 0.6)	39.7 ( $\pm$ 0.9)		
B62 <sup>k</sup>	70.0 ( $\pm$ 0.4)	35.7 ( $\pm$ 0.7)		
PAM <sup>l</sup>	68.6 ( $\pm$ 0.6)	29.0 ( $\pm$ 0.7)		
Var. <sup>m</sup>	68.0 ( $\pm$ 0.7)	25.4 ( $\pm$ 1.1)		
Entr. <sup>n</sup>	66.5 ( $\pm$ 0.4)	24.6 ( $\pm$ 1.0)		
VdW.Vol <sup>o</sup>	62.1 ( $\pm$ 0.3)	11.4 ( $\pm$ 0.3)		
Vol.Bur. <sup>p</sup>	61.7 ( $\pm$ 0.4)	10.3 ( $\pm$ 0.5)		
W/O trf. <sup>q</sup>	59.7 ( $\pm$ 0.5)	8.8 ( $\pm$ 2.1)		
(C)	Qtot <sup>a</sup>	S <sup>b</sup>		
Pssm <sup>i</sup>	77.7 ( $\pm$ 0.5)	49.9 ( $\pm$ 1.4)		
$\Delta$ Pssm <sup>j</sup>	73.2 ( $\pm$ 0.7)	38.9 ( $\pm$ 1.2)		
B62 <sup>k</sup>	70.0 ( $\pm$ 0.4)	35.7 ( $\pm$ 0.7)		
VdW.Vol <sup>o</sup>	68.5 ( $\pm$ 0.1)	30.5 ( $\pm$ 0.3)		
PAM <sup>l</sup>	68.3 ( $\pm$ 0.5)	61.9 ( $\pm$ 0.7)		
Vol.Bur. <sup>p</sup>	68.0 ( $\pm$ 0.1)	31.7 ( $\pm$ 0.6)		
I/O prot. <sup>r</sup>	67.9 ( $\pm$ 0.4)	28.7 ( $\pm$ 0.9)		
Var. <sup>m</sup>	67.9 ( $\pm$ 0.7)	24.5 ( $\pm$ 1.1)		
Entr. <sup>n</sup>	66.6 ( $\pm$ 0.5)	22.3 ( $\pm$ 0.8)		
W/O trf. <sup>g</sup>	65.4 ( $\pm$ 1.7)	22.6 ( $\pm$ 5.6)		
Acc. <sup>s</sup>	62.2 ( $\pm$ 0.6)	46.9 ( $\pm$ 2.0)		
Chou and Fassman <sup>t</sup>	57.0 ( $\pm$ 0.9)	6.8 ( $\pm$ 1.0)		

<sup>†</sup>(A) Performance using all available information. (B) and (C) Performance of the perceptron and the H12 neural networks, respectively, when trained using specific parameters only.

<sup>a</sup>Overall success rate (see Methods).

<sup>b</sup>Normalized performance relative to random (see Methods).

<sup>c</sup>False positive rate.

<sup>d</sup>False negative rate.

<sup>e</sup>Performance of the method by Ng and Henikoff, derived from data provided by Ng and Henikoff<sup>19</sup> (Table 4 in reference).

<sup>f</sup>Cross-validated performance of our method when tested with NEMUs from the Lac system (see Methods).

<sup>g</sup>Cross-validated performance of our method when tested with NEMUs from the evolutionary model (see Methods).

<sup>h</sup>Cross-validated performance of our method when tested with NEMUs from the evolutionary model using observed values for the structure-related properties (accessibility, secondary structure and statistical potentials; see Methods).

<sup>i–t</sup>Cross-validated performance for the NN when using as input only one of the following parameters: Pssm, position-specific scoring matrices (Equation 4 in Methods);  $\Delta$ Pssm, difference in position-specific scoring matrices (Equation 5 in Methods); B62, blosum62 matrix elements; PAM, PAM40 matrix elements; Var., average of mutation matrix scores (Equation 3 in Methods); Entr., Shannon entropy; VdW.Vol., van der Waals amino acid volumes; Vol.Bur., volume of buried residues; W/O trf., water/octanol transfer free energy; I/O prot., structure-derived hydrophobicity; Acc. PHD accessibility predictions; Chou and Fassman, Chou and Fassman secondary structure propensity. All these parameters are explained in the Methods section.

were obtained from sequence-based predictions (see below). In the structure-based version of our approach we used six structure-related parameters: observed secondary

structure and solvent accessibility (three-state and relative), and statistical potentials (three values).

#### **Predicted secondary structure and accessibility.**

They were obtained from the protein sequence, using the PHD software package<sup>29</sup> (www.embl-heidelberg.de/predict-protein) with default parameters. The three secondary structure states, helix, beta, and coil, were encoded as 0, 1, and 2, respectively. The three accessibility states, buried, half-buried and exposed, were encoded as 0, 1, and 2, respectively.

#### **Observed secondary structure and solvent accessibility.**

They were computed from the experimental structure of the protein. Secondary structure at the mutation site was obtained using the SSTRUC implementation of the Kabsch and Sander method,<sup>30</sup> by David Keith Smith. The resulting three-states were encoded as in the previous section. Accessibility values were obtained using the program NACCESS.<sup>31</sup> Both relative and three-state accessibilities were used. The former are equal to the residue accessibility in the protein divided by its accessibility in an extended Ala-X-Ala peptide. Relative accessibilities were mapped to three-state accessibilities: buried (0–9% relative accessibility), half-buried (9–36% relative accessibility) and exposed (36–100% relative accessibility). These three-states were encoded as in the previous section.

**Statistical potentials.** They were used to assess the destabilizing effect of the mutant. For each residue in a protein of known structure, Prosa II<sup>32</sup> calculations gave us three measures indirectly related to the residue contribution to protein stability: surface potential, contact potential, and an overall potential (weighted sum of the two previous potentials). The three terms were computed for the native residue, for the mutated one (the structure of the mutated protein was obtained by just mapping the residue change in the wild-type structure, no structural modeling was done). The mutation was then labeled with the three corresponding differences between mutant and native values.

#### **Residue/sequence properties**

Three information types were used: mutation matrices, amino acid properties and sequence potentials.

**Mutation matrices.** Each mutation was labelled with both the Blosum62<sup>33</sup> and PAM40<sup>34</sup> matrix values.

**Changes in single amino acid properties.** We used six residue-based parameters: two hydrophobic indexes, two secondary structure propensities, and two volume changes. For each of them, the value associated to the mutation was the difference between the mutant,  $x_m$ , and wild-type,  $x_w$ , values of this property:  $x_m - x_w$ . The hydrophobic parameters were taken from water/octanol free energy measurements,<sup>35</sup> and from statistical potentials derived by Miller et al.<sup>36</sup> from structural information. Secondary structure propensities were obtained from standard Chou and Fasman analysis,<sup>37</sup> as well as from the Swindells et al. analysis.<sup>38</sup> Size descriptors were van der Waals volumes,<sup>39</sup> and volume of buried residues.<sup>40</sup>

**Sequence potential.** To take into account the effect of sequence environment on the mutation effect<sup>41</sup> we used a



simple potential,  $Ptseq(r_j)$  (see equation 1), related to the probability of observing residue  $r_j$  at position  $j$ , in a given sequence environment:

$$Ptseq(r_j) = \ln \left[ \prod_{i=-5}^5 P(r_j/r_{j+i}) \right] \quad (1)$$

The  $P(r_j/r_{j+i})$  were obtained following:

$$P(r_j/r_{j+i}) = n(r_j, r_{j+i})/n(r_{j+i}) \quad (2)$$

where  $n(r_j, r_{j+i})$  is the number of pairs of amino acids of types  $r_j$  and  $r_{j+i}$  at a sequence distance  $i$ .  $n(r_{j+i})$  is the total number of residues of type  $r_{j+i}$ . These numbers are computed using the whole set of human sequences from the SwissProt database.<sup>23</sup>

Mutations were scored using the difference between the value of the sequence potential for the mutant and wild-type residues, as for the simple amino acid properties.

### Evolutionary properties

Evolutionary information was exploited using four parameters: two measures of the amino acid variability, and two position-specific scores. Computation of these parameters required the use of multiple sequence alignments, that were obtained from the Pfam<sup>24</sup> database. Pfam alignments (27) were chosen for their quality, however, the method can be easily extended to other kinds of multiple sequence alignments, in particular to those generated by the program BLAST.<sup>42</sup>

**Variability at the mutation position in the multiple sequence alignment.** This was measured using the Shannon entropy<sup>43</sup> and an average of mutation matrix scores.<sup>16</sup> The former was computed following:  $-\sum_i p_{ij} \ln p_{ij}$ , where subindex  $i$  runs over the different amino acid types found at position  $j$ , location of the mutated residue in the multiple sequence alignment.  $p_{ij}$  are the relative frequencies of these amino acids.

The second measure was computed following Martin et al.<sup>16</sup>

$$\left[ \frac{\sum_{k=1}^N \sum_{l=k+1}^N s_{kl}}{N!} \right] / S_{\max} \quad (3)$$

where  $s_{kl}$  is the element of Blosum62<sup>33</sup> corresponding to the comparison between residues in sequences  $k$  and  $l$  in the multiple sequence alignment, at the position of the mutation under study.  $S_{\max}$  is the larger  $s_{kl}$ .  $N$  is the number of sequences in the alignment.

**Position-specific scoring matrices (PSSM).** We utilized two kinds of PSSM parameters. The first is based on the log-odds ratio:  $\log(p_{mj}/p_m)$ , where  $p_{mj}$  is the relative frequency of the mutant amino acid type  $m$  at position  $j$  in the multiple sequence alignment ( $j$  is here the location of the mutated residue) and  $p_m$  is the frequency of the same amino acid type in all human sequences in SwissProt.<sup>23</sup> To

alleviate the problem of missing data, we used a modified version of the log-odds ratio:<sup>44</sup>

$$\frac{N\sigma}{1+N\sigma} \log(p_{mj}/p_m) + \frac{1}{1+N\sigma} B62_{wm} \quad (4)$$

where  $N$  is the number of sequences in the multiple sequence alignment,  $\sigma$  is an arbitrary factor, and was taken<sup>44</sup> equal to 1/50.  $w$  and  $m$  stand for the normal and mutant amino acid, respectively.  $B62_{wm}$  is the element of the mutation matrix Blosum62<sup>33</sup> corresponding to the mutation from the wild-type residue ( $w$ ) to the mutant residue ( $m$ ). When  $N$  is small, the value of the PSSM index approaches that of  $B62_{wm}$ ; when it is large it approaches the value of the log-odds ratio.

The second PSSM parameter is determined following equation 5, and shows the same asymptotic behavior as the previous index.

$$\frac{N\sigma}{1+N\sigma} [\log(p_{mj}/p_m) - \log(p_{wj}/p_w)] + \frac{1}{1+N\sigma} B62_{wm} \quad (5)$$

where  $m$ ,  $w$ ,  $N$ ,  $\sigma$ ,  $p_{mj}$ ,  $p_m$ , and  $B62_{wm}$  have the same meaning as before.  $p_{wj}$  is the relative frequency of the original amino acid type  $w$  at position  $j$  in the multiple sequence alignment.  $p_w$  is the frequency of the same amino acid type in human sequences in SwissProt.<sup>23</sup>

It has to be noted that a certain amount of redundancy is to be expected in the multiple sequence alignments used, thus reducing the contribution of these variables to the success of the method. While this problem can be alleviated using a filtering procedure, previous results from our group suggest that differences between raw and filtered data are minor.<sup>13</sup>

### Database information

Four indexes were used to include database functional annotations taken from the SwissProt<sup>23</sup> database. They indicate that the mutated residue is: (1) involved in disulfide bridge, thiolester, thiolether; (2) inside an alternative splicing region; (3) a modified residue (has a bound carbohydrate, is a selenocysteine); (4) an active site, or binds nucleotide,  $Ca^{2+}$  or  $Zn^{2+}$ . For each index, a value of 0 indicates absence of annotation, and a value of 1 indicates presence of annotation.

### The Neural Network

A feed-forward neural network,<sup>45</sup> with one input layer one, or none (perceptron model<sup>27</sup>) hidden layer and one output layer was used. For the NN having a hidden layer, the latter was constituted by 20 units for the full version of the method, and by 2 units for the study of the individual contributions of the different parameters. A total of 19 (sequence-only version) or 23 (sequence plus structure version) parameters were input to the network—the mutation labels described above.

For a given mutation, the network output is a number comprised between 0 and 1, that is transformed into a discrete prediction as follows: for values above 0.5, the mutation is predicted as DAMU; for values below 0.5, the

mutation is predicted as NEMU. Strong predictions will be those for which the output is either very close to one, or very close to zero, e.g., 0.9 or 0.01, respectively. Weak predictions will be those with an output closer to 0.5, e.g. 0.44 or 0.58.

We followed the training procedure described in Shepherd et al.,<sup>46</sup> presenting the network with a number of inputs, together with their associated target outputs. The network weights were optimized using scaled conjugate gradients with 500 iterations.

### Cross-Validation

The performance of our method was evaluated using a stringent heterogeneous 5-fold cross-validation procedure.<sup>9</sup> Rather than randomly dividing the whole mutations set (DAMUs plus NEMUs) into five subsets and subsequently train the NN, the mutation dataset was split at the protein level. Thus, we divided into five parts the set of proteins for which at least one DAMU was found. The corresponding five sets of mutations were then built and, subsequently, each of the five unique combinations of four different subsets was used to train the network. For each combination, the network was tested on the excluded subset. The results for the test sets were then averaged to provide the results shown in this work.

### Performance Measures

We utilized four parameters to evaluate the performance of our prediction method: percentage of correct predictions, normalized percentage of improvement over random predictions, false positive and false negative rates. We describe them below.

Percentage of correct predictions ( $Q_{tot}$ , also referred to as overall success rate). It provides an overall view of the ability of the procedure to detect pathological/neutral mutations. It is computed as follows:

$$Q_{tot} = 100 \frac{cp}{(cp + ip)} \quad (6)$$

where cp and ip are the overall number of correct and incorrect predictions respectively.

Normalized percentage of improvement over random predictions (S). This parameter is a measure of how well the method is working relative to a random predictor with expected performance R, normalized to eliminate the scale effect from R and from the total number of observations.<sup>46</sup> It varies between 0% (no improvement relative to a random method) and 100% (perfect, non-random, predictions). It is computed as follows:

$$S = \frac{(p + n) - R}{t - R} \times 100 \quad (7)$$

where  $t = p + n + o + u$ . p, n are the number of mutations correctly predicted as DAMUs and NEMUs, respectively; o, u are the number of mutations incorrectly predicted as DAMUs and NEMUs, respectively.  $R = [(p + o).(p + u).(n + o).(n + u)]/t$  is the expected number of correctly classified mutations generated by a random predictor that

would take into account the proportion of DAMUs and NEMUs in our sample.

Percentage of false positives (FP) and percentage of false negatives (FN). They provide an idea of the accuracy limits of our approach. They are computed as follows:

$$FP = \frac{o}{p + o} \times 100 \quad (8)$$

$$FN = \frac{u}{p + u} \times 100 \quad (9)$$

where p, n, o and u have the same meaning as before.

### Reliability Index

The reliability index is computed from the NN output following:<sup>46</sup>

$$\text{integer}[\text{abs}(\text{NN}_{\text{output}} - 0.5) \times 20] \quad (10)$$

where  $\text{NN}_{\text{output}}$  is the NN output. The resulting index varies between 0 and 9, with low values of the index corresponding to poorer predictions, and high values corresponding to better predictions.

## RESULTS AND DISCUSSION

As mentioned before, in this work we present a computational procedure to distinguish between disease-associated mutations and neutral mutations, utilizing sequence-based information and neural networks. The performance of our method has been tested on a large set of human DAMUs (9334 mutations). Some authors<sup>3,9</sup> have utilized sets of mutations affecting viability in microbial systems, as a model for human DAMUs. However, extrapolation to humans of the results obtained with microbial systems may be very complex and it is still an open issue.<sup>9</sup> In our case, we have used actual human DAMUs, because large amount of them are available in public databases and have been used for the characterization of human pathological mutations.<sup>7,13,17,47</sup> In particular, DAMUs used in this work were obtained from the SwissProt database,<sup>23</sup> following Ng and Henikoff<sup>19</sup> (see Methods).

For the NEMUs, we have utilized two different models (see Methods): (1) the Lac dataset, because it has been previously used by Ng and Henikoff<sup>19</sup> for their sequence-based method, thus allowing a direct comparison between methods, and (2) an evolutionary model, used by Bork and coworkers for their structure-based method,<sup>5</sup> and by Santibáñez-Koref et al.<sup>11</sup> for their sequence-based prediction method of deleterious mutations in p53.

Our method shows a high overall success rate, 87.3 ( $\pm 0.4$ ), when utilized to discriminate between DAMUs and NEMUs from the Lac dataset, with a substantial improvement over purely random predictions, 54.2 ( $\pm 2.5$ ). The false positive and false negative rates are 10.4 ( $\pm 0.6$ ) and 4.6 ( $\pm 0.5$ ), respectively. We can compare our results with those obtained by Ng and Henikoff,<sup>19</sup> who have derived the other sequence-based method described for the prediction of pathological mutations in large sets of different proteins. We can see that our method has a higher performance [Table I(A)]. Both overall success rate, 87.3 versus

70.3, and performance relative to random, 54.2 versus 37.0, are better in our case. The same is true for the false positive and false negative rates [see Table I(A)]. We believe that this improved performance of our method is probably due to the fact that we use more features to characterize mutations, the way we process evolutionary information, and the pattern recognition power of NNs.

When testing the method with the evolutionary model for NEMUs we also find a high performance rate, 83.5 ( $\pm 0.3$ ), with a clear improvement of 66.5 ( $\pm 0.5$ ) over purely random predictions. The false positive and false negative rates are 16.5 ( $\pm 0.6$ ) and 21.1 ( $\pm 0.6$ ), respectively. The overall performance of the method is similar to that obtained when utilizing NEMUs from the Lac system, with a better false positive and negative rate for the latter, but a significantly higher improvement over random predictions found with the evolutionary model. This is because the method's performance figures for neutral mutations are clearly better with the latter model than when using the Lac model for NEMUs: false positive rates of 16.6 ( $\pm 0.5$ ) and 26.6 ( $\pm 1.2$ ), respectively; false negative rates of 12.8 ( $\pm 0.6$ ) and 46.9 ( $\pm 3.6$ ), respectively. Compared with the method by Ng and Henikoff<sup>19</sup> the results of the second version of our method are also clearly better.

Taken together, the results from both tests indicate that DAMUs can be distinguished from NEMUs with a high overall success rate—between 83% and 87%, using only sequence-based information. The use of two different models for NEMUs does not affect the overall performance of the method, although utilizing NEMUs from the evolutionary model results in more balanced results for both DAMUs and NEMUs. For this reason, and because data from the Lac system are more likely to be biased towards this specific system, in the following we will only discuss performance figures corresponding to our method when trained in the set of evolutionary NEMUs (it has to be noted that the results obtained with the Lac dataset lead to similar conclusions).

### Factors Contributing to the Recognition Rate of Our Approach

The fact that no structure information was required for the identification of DAMUs is very important from a practical point of view, because there is still a large number of proteins for which only sequence information is available.<sup>48</sup> To understand why this happens, we studied the discriminant power of the nineteen parameters used. To this end, for each of them we assessed the performance of two NNs: a simple perceptron and one with a two-units hidden layer (H12); both trained with the corresponding parameter as only input. We utilized these two NNs because the discriminant power of some variables may be substantially underestimated when using the perceptron.<sup>28</sup> Within each NN model, the results for the different parameters were ranked according to the overall success rate [Table I(B) for the perceptron and Table I(C) for H12].

First, we observe that none of the parameters alone shows a performance comparable to that of the full version of the method, indicating that several of them are required

for the success of the latter. However, we also find that the best performance is always obtained with the evolutionary indexes—the two PSSM parameters, in accordance with previous studies,<sup>7,19</sup> with a success rate clearly above that from the other variables. Two reasons can explain the large predictive power of evolutionary information as measured in this work. First, residue conservation in protein families is directly related to its contribution to protein stability or function. Therefore, because a large number of DAMUs are likely to cause protein destabilization,<sup>6</sup> they can be modelled using this kind of information. The same applies for the most easy-to-understand DAMUs, those affecting the functional site of the protein.<sup>5</sup> The second reason is the fact that the DAMUs considered here are responsible of monogenic disorders, and thus more likely to have large damaging effects on protein function, or stability.<sup>49</sup> Because of this, they will tend to appear as clear departures from the multiple sequence alignment conservation pattern at the mutation site, making their identification easier.

The highest discrimination power observed for PSSM is paralleled, at the distribution level, by the smallest overlap between the distributions for NEMUs and DAMUs, relative to this parameter [Fig. 1(A)]. Positive and negative values of the latter tend to correspond to NEMUs and DAMUs, respectively. Interestingly, and in agreement with this picture, both NNs just use a single threshold to distinguish between NEMUs and DAMUs [Fig. 1(A)].

PSSM parameters are followed by mutation matrix elements, both Blossum 62 and PAM 40 matrix elements [Table I(B,C)]. However, for the H12 neural network van der Waals volumes give a slightly better success rate than PAM 40 matrix elements. This is in contradiction with the results obtained for the perceptron, where van der Waals volumes show clearly worse discrimination ability. This can be explained when looking at the distribution of mutations relative to volume changes [Fig. 1(B)]. We can see that NEMUs prevail in an interval around zero, that is, they are associated with volume changes smaller than those of DAMUs. This is consistent with the fact that large volume changes, whether positive or negative, are more likely to cause substantial structural damage.<sup>13</sup> In accordance with this picture, H12 defines an interval inside which mutations are predicted as NEMUs, and outside which they are predicted as DAMUs [Fig. 1(B)]. On the contrary, the volume-based perceptron follows a much poorer strategy, using a single threshold [Fig. 1(B)] to separate mutations in DAMUs or NEMUs, if their associated volume changes are above or below the threshold, respectively.

Comparing Table I(B) and (C) we can see that the parameters shown can be divided in two groups: (1) those for which the performance is almost the same for both NN models; (2) those that show a better performance with the H12 neural network. Analysis of the frequency histograms (Fig. 1, and results not shown) indicates that for members of the first group NEMUs and DAMUs tend to have values of opposite sign for the corresponding parameter. For this reason, a unique threshold value may be enough to provide



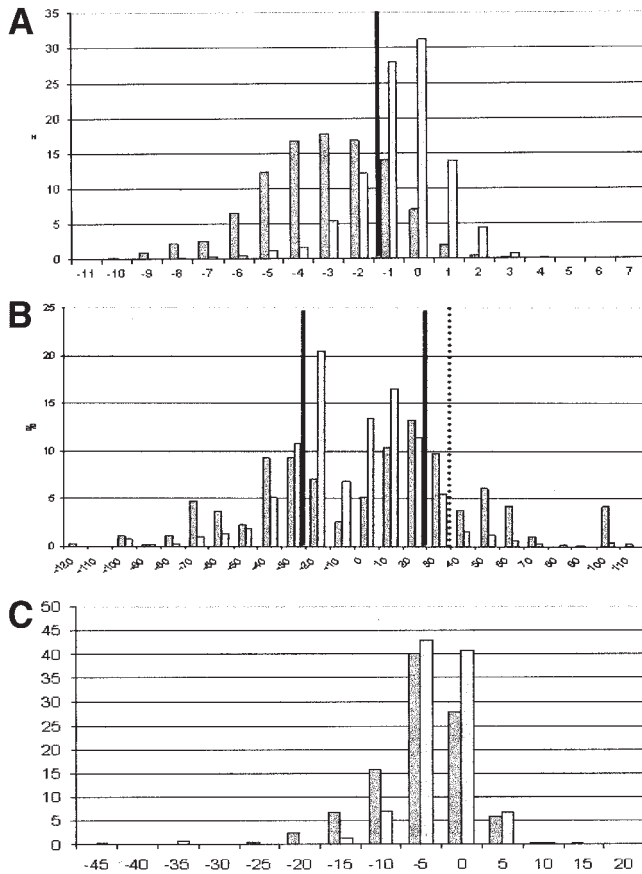


Fig. 1. Distribution of DAMUs (dark grey) and NEMUs (light grey) relative to different scoring parameters (see Methods for the definitions): (A) Position-specific scoring matrices (see equation 4, in Methods), (B) van der Waals volumes, (C) overall statistical potential. In (A), the thick bar corresponds to the threshold value below which mutations are predicted as pathological, and above which as neutral, by both the perceptron and the H12 neural networks (see text). In (C), the thick bars define the interval within which mutations are predicted as pathological, and outside which as neutral, by the H12 neural network. The dotted line corresponds to the threshold defined by the perceptron, above which mutations are predicted as pathological, and below which as neutral.

an almost maximum discriminatory power, a strategy implicit in the perceptron. Members of this first group include: the two PSSM-related indexes, the two variability measures, and the Blossum 62 matrix. On the contrary, for members of the second group, it is the absolute value of the change rather than its sign, which determines the damage caused by the mutation and thus its association to disease. Therefore, good discrimination between mutations will require defining an interval within which mutations will be of one type, and outside which of the other type. This strategy can be obtained with more complex NN models, like H12. Members of this second group include the two volume and hydrophobicity properties, and predicted accessibility. Within this second group, we can also include the Chou and Fassman secondary structure propensities, although their predictive value is very poor. PAM 40 matrix shows features of both groups, with a clear improvement in the performance relative to random when using H12. Parameters not appearing in Table I(B) and (C) led to very

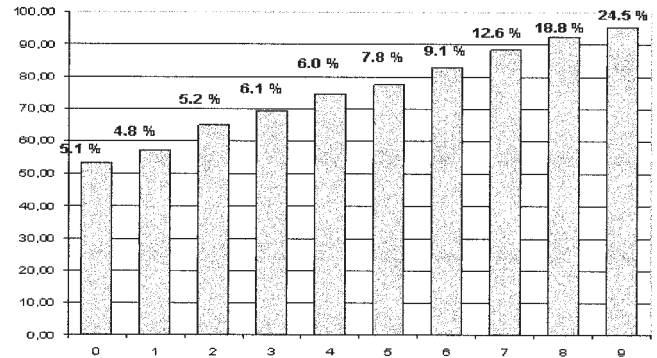


Fig. 2. Relationship between the percentage of correct predictions,  $Q_{tot}$ , (ordinates) and reliability (abscissae) of the predictions (see text). The percentage of pathological mutations predicted with a given reliability is shown above each bar.

small improvements relative to random, below 5%, for which reason they were excluded from this analysis (these parameters are: secondary-structure predictions, sequence potential, coil-derived secondary-structure propensities, and SwissProt annotations).

### Reliability of the Predictions

Although the average success rate of our approach is reasonably good, it may be of interest to know how reliable are individual predictions, in order to prioritize mutations for further analysis. To this end we can use the NN output which is related to the nature of the damage caused by the mutation. For example, pathological mutations strongly disrupting the multiple sequence alignment conservation pattern at the mutation site will result in large, negative, PSSM values. The latter will, in turn, lead to strong DAMUs predictions by the NN.

The NN output can be easily transformed into a reliability index that will reflect the prediction strength, using equation 10 in the Methods section. The resulting index varies between 0 (unreliable predictions) and 9 (highly reliable predictions); it is directly related to the performance of the method, with better performances corresponding to higher values of the reliability index (Fig. 2). The latter confirms that the reliability index can be a useful tool for filtering sets of mutation data, ranking mutations according to the prediction reliability, and eliminating those with prediction reliabilities below a given threshold.

### Improving the Performance of Our Approach

It has been shown previously that both predicted secondary structure and accessibility had a small contribution to the recognition power of our procedure. Therefore, a first option to improve the performance of the latter is to use observed, instead of predicted, structural properties derived from the three-dimensional structure of the protein affected by the mutations. In this case, in addition to observed secondary structure and accessibility, we also scored mutations utilizing statistical potentials,<sup>32</sup> which have been used for the prediction of mutant stability changes.<sup>50</sup>

**TABLE II. Performance of Our Method When Trained and Tested in Individual Human Proteins (see Text)**

	DAMUs <sup>a</sup>	NEMUs <sup>b</sup>	Qtot <sup>c</sup>	S <sup>d</sup>
P53 <sup>e</sup>	172	104	96.7 ± 1.1	93.1 ± 2.2
ACTS	8	197	97.1 ± 0.0	47.2 ± 8.7
HBB	124	80	95.1 ± 1.0	89.8 ± 1.7

<sup>a</sup>Number of disease associated mutations.<sup>b</sup>Number of neutral mutations.<sup>c</sup>Overall success rate (see Methods).<sup>d</sup>Normalized performance relative to random (see Methods).<sup>e</sup>Proteins used: P53, p53; ACTS, Alpha-actin 1; HBB, Hemoglobin beta chain.

We find that the NN has a good overall performance rate, with 87.0 (±1.6) correct predictions. The performance relative to random is also good, 73.0 (±3.1). The false positive and false negative rates are 10.4 % (± 2.1) and 11.5% (±1.3), respectively. Except for the false negative rates, these figures represent a minor improvement relative to the sequence-based version of our procedure. This indicates that even when using observed structural properties, the first contributor to the method's success is evolutionary information. Accessibility and secondary structure improve their discrimination power, although it still is lower than that of evolutionary parameters (data not shown). Statistical potentials display a poor discrimination power, as can be seen by comparing Figure 1(A) and (C). Overall, these results suggest that further improvements to our method could come from a better treatment of the evolutionary information. At present our treatment of multiple sequence alignments is relatively crude, as it ignores the underlying structure of the data. A possible approach for a better use of evolutionary information would be to consider the underlying phylogenetic relationships, as has been described by Santibáñez-Koref et al.<sup>11</sup>

A second option would be to train our procedure for specific systems, concentrating in those proteins that may be of particular interest due to their relationship to diseases that have a high social cost and for which a vast amount of mutational data is available. This is the case for p53, a protein directly related to oncogenic processes,<sup>51</sup> for which more than 14,000 mutations are known.<sup>52</sup> In cases like this, predictive approaches derived for that specific system may provide better performance, as shown for p53 in recent studies.<sup>11,16</sup>

To test this idea with our method, we chose from our dataset those proteins for which the ratio between number of training mutations and parameters was between 5 and 10.<sup>27</sup> p53, alpha-actin 1, and hemoglobin beta chain. A simple NN with no hidden layer (perceptron) was then trained, and validated using a twofold cross-validation procedure. For two of the three proteins (p53 and hemoglobin) the cross-validated performance was very high (see Table II), clearly improving that of the general model [Table I(A)]. However, for alpha-actin 1 the performance was lower than that of the general model. This is probably due to the very low number of DAMUs, 8, relative to NEMUs, 197, resulting in poorly trained NNs. This was not the case for p53 and haemoglobin, for which the

amounts of both types of mutations were more balanced, thus ensuring a better training of the network.

Clearly, our results support the idea that when interested in a given protein it is better to use a specific predictor, if enough mutation data are available and if the number of DAMUs and NEMUs is similar. This can be the case of proteins like p53, for which large mutations datasets are available and have been used to train specific prediction methods.<sup>11,16</sup> Otherwise, the general model is enough to obtain a reasonable prediction of the pathological character of SNPs.

## CONCLUSIONS

We describe a simple method for the prediction of human disease-associated mutations, based on the use of sequence-based information and neural networks. The good success rate of the method, 83.5%, indicates that pathological mutations can be identified even in absence of structural information. This is particularly relevant when considering the increasing gap between known sequences and structures. In addition, use of the reliability index produced by the NN output provides a simple tool to assess the predictions, an option of particular interest when considering future applications of the method to large sets of nonsynonymous SNPs.

## ACKNOWLEDGMENTS

The authors acknowledge the economic support provided by the Fundación Areces. They also acknowledge C. Saunders and D. Baker for kindly giving us their mutations sets. The work has been supported by the Spanish Ministry of Science and Technology (PM99-0046, GEN2001-4758-C07-07 and BIO2003-09327).

## REFERENCES

1. Cargill M, Altshuler D, Ireland J, Sklar P, Ardlie K, Patil N, Shaw N, Lane CR, Lim EP, Kalyanaraman N, et al. Characterization of single-nucleotide polymorphisms in coding regions of human genes. *Nat Genet* 1999;22:231–238.
2. Chakravarti A. To a future of genetic medicine. *Nature* 2001;409:822–823.
3. Chasman D, Adams RM. Predicting the functional consequences of non-synonymous single nucleotide polymorphisms: structure-based assessment of amino acid variation. *J Mol Biol* 2001;307:683–706.
4. Ng PC, Henikoff S. Predicting deleterious amino acid substitutions. *Genome Res* 2001;11:863–874.
5. Sunyaev S, Ramensky V, Koch I, Lathe W, 3rd, Kondrashov AS, Bork P. Prediction of deleterious human alleles. *Hum Mol Genet* 2001;10:591–597.
6. Wang Z, Moult J. SNPs, protein structure, and disease. *Hum Mutat* 2001;17:263–270.
7. Saunders CT, Baker D. Evaluation of structural and evolutionary contributions to deleterious mutation prediction. *J Mol Biol* 2002;322:891–901.
8. Ramensky V, Bork P, Sunyaev S. Human non-synonymous SNPs: server and survey. *Nucleic Acids Res* 2002;30:3894–3900.
9. Krishnan VG, Westhead DR. A comparative study of machine-learning methods to predict the effects of single nucleotide polymorphisms on protein function. *Bioinformatics* 2003;19:2199–2209.
10. Herrgard S, Cammer SA, Hoffman BT, Knutson S, Gallina M, Speir JA, Fetrow JS, Baxter SM. Prediction of deleterious functional effects of amino acid mutations using a library of structure-based function descriptors. *Proteins* 2003;53:806–816.
11. Santibáñez-Koref MF, Gangeswaran R, Santibáñez-Koref IP, Shanahan N, Hancock JM. A phylogenetic approach to assessing



- the significance of missense mutations in disease genes. *Hum Mutat* 2003;22:51–58.
12. Sunyaev S, Ramensky V, Bork P. Towards a structural basis of human non-synonymous single nucleotide polymorphisms. *Trends Genet* 2000;16:198–200.
  13. Ferrer-Costa C, Orozco M, de la Cruz X. Characterization of disease-associated single amino acid polymorphisms in terms of sequence and structure properties. *J Mol Biol* 2002;315:771–786.
  14. Miller MP, Parker JD, Rissing SW, Kumar S. Quantifying the intragenic distribution of human disease mutations. *Ann Hum Genet* 2003;67(Pt 6):567–579.
  15. Perutz MF. Protein structure: new approaches to disease and therapy. New York: W.H. Freeman and Co.; 1992.
  16. Martin AC, Facchiano AM, Cuff AL, Hernandez-Boussard T, Olivier M, Hainaut P, Thornton JM. Integrating mutation data and structural analysis of the TP53 tumor-suppressor protein. *Hum Mutat* 2002;19:149–164.
  17. Steward RE, MacArthur MW, Laskowski RA, Thornton JM. Molecular basis of inherited diseases: a structural perspective. *Trends Genet* 2003;19:505–513.
  18. Guex N, Diemand A, Peitsch MC. Protein modelling for all. *Trends Biochem Sci* 1999;24:364–367.
  19. Ng PC, Henikoff S. Accounting for human polymorphisms predicted to affect protein function. *Genome Res* 2002;12:436–446.
  20. Bonneau R, Strauss CE, Rohl CA, Chivian D, Bradley P, Malmstrom L, Robertson T, Baker D. De novo prediction of three-dimensional structures for major protein families. *J Mol Biol* 2002;322:65–78.
  21. Goodstadt L, Ponting CP. Sequence variation and disease in the wake of the draft human genome. *Hum Mol Genet* 2001;10:2209–2214.
  22. McKusick VA. Mendelian inheritance in man: a catalog of human genes and genetic disorders. Baltimore: Johns Hopkins University Press; 1998.
  23. Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res* 2000;28:45–48.
  24. Bateman A, Birney E, Cerruti L, Durbin R, Etwiler L, Eddy SR, Griffiths-Jones S, Howe KL, Marshall M, Sonnhammer EL. The Pfam protein families database. *Nucleic Acids Res* 2002;30:276–280.
  25. Pace HC, Kercher MA, Lu P, Markiewicz P, Miller JH, Chang G, Lewis M. Lac repressor genetic map in real space. *Trends Biochem Sci* 1997;22:334–339.
  26. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. The Protein Data Bank. *Nucleic Acids Res* 2000;28:235–242.
  27. Mehrotra K, Mohan CK, Ranka S, NetLibrary Inc. Elements of artificial neural networks. Cambridge, MA: MIT Press; 1997. xiv, 344 p.
  28. Bishop CM. Neural networks for pattern recognition. Oxford, New York: Clarendon Press; Oxford University Press; 1995. xvii, 482 p.
  29. Rost B, Sander C. Prediction of protein secondary structure at better than 70% accuracy. *J Mol Biol* 1993;232:584–599.
  30. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* 1983;22:2577–2637.
  31. Hubbard SJ, Thornton JM. 'NACCESS', computer program, Department of Biochemistry and Molecular Biology, University College London, 1993.
  32. Sippl MJ. Recognition of errors in three-dimensional structures of proteins. *Proteins* 1993;17:355–362.
  33. Henikoff S, Henikoff JG. Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci USA* 1992;89:10915–10919.
  34. Dayhoff MO, Schwartz RM, Orcutt BC. A model of evolutionary change in proteins. In: Dayhoff MO, editor. Atlas of protein sequence and structure. Washington D. C.: National Biomedical Research Foundation; 1978. p 345–352.
  35. Fauchere JL, Pliska V. Hydrophobic parameters of amino acid side chains from the partitioning of N-acetyl-amino acid amides. *Eur J Med Chem* 1983;18:369–375.
  36. Miller S, Janin J, Lesk AM, Chothia C. Interior and surface of monomeric proteins. *J Mol Biol* 1987;196:641–656.
  37. Chou PY, Fasman GD. Conformational parameters for amino acids in helical, beta-sheet, and random coil regions calculated from proteins. *Biochemistry* 1974;13:211–222.
  38. Swindells MB, MacArthur MW, Thornton JM. Intrinsic phi, psi propensities of amino acids, derived from the coil regions of known structures. *Nat Struct Biol* 1995;2:596–603.
  39. Bondi A. Van der Waals volumes and radii. *J Phys Chem* 1964;68:441–451.
  40. Chothia C. Structural invariants in protein folding. *Nature* 1975;254:304–308.
  41. Yang W, Battineni ML, Brodsky B. Amino acid sequence environment modulates the disruption by osteogenesis imperfecta glycine substitutions in collagen-like peptides. *Biochemistry* 1997;36:6930–6935.
  42. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;25:3389–3402.
  43. Shannon CE. A mathematical theory of communication. *Bell System Tech J* 1948;27:379–423.
  44. Sippl MJ. Calculation of conformational ensembles from potentials of mean force. An approach to the knowledge-based prediction of local structures in globular proteins. *J Mol Biol* 1990;213:859–883.
  45. Rumelhart DE, McClelland JL. Parallel distributed processing: explorations in the microstructure of cognition. Cambridge, MA: MIT Press; 1986.
  46. Shepherd AJ, Gorse D, Thornton JM. Prediction of the location and type of beta-turns in proteins using neural networks. *Protein Sci* 1999;8:1045–1055.
  47. Sunyaev S, Hanke J, Brett D, Aydin A, Zastrow I, Lathe W, Bork P, Reich J. Individual variation in protein-coding sequences of human genome. *Adv Protein Chem* 2000;54:409–437.
  48. Peitsch MC. About the use of protein models. *Bioinformatics* 2002;18:934–938.
  49. Dipple KM, McCabe ER. Phenotypes of patients with "simple" Mendelian disorders are complex traits: thresholds, modifiers, and systems dynamics. *Am J Hum Genet* 2000;66:1729–1735.
  50. Gilis D, Rooman M. PoPMuSiC, an algorithm for predicting protein mutant stability changes: application to prion proteins. *Protein Eng* 2000;13:849–856.
  51. Vousden KH, Lu X. Live or let die: the cell's response to p53. *Nat Rev Cancer* 2002;2:594–604.
  52. Hainaut P, Hernandez T, Robinson A, Rodriguez-Tome P, Flores T, Hollstein M, Harris CC, Montesano R. IARC Database of p53 gene mutations in human tumors and cell lines: updated compilation, revised formats and new visualisation tools. *Nucleic Acids Res* 1998;26:205–213.