

Aplicación de minería de datos y modelamiento matemático en ingeniería de proteínas

**Diseño e implementación de nuevas metodologías para el
estudio de mutaciones**



**UNIVERSIDAD
DE CHILE**

David Medina Ortiz

Supervisor: Dr. Álvaro Olivera

**Departamento de Ingeniería Química, Biotecnología y Materiales
Universidad de Chile**

Septiembre 2019

Tabla de contenidos

Lista de figuras	vii
Lista de tablas	ix
1 Aplicaciones de la minería de datos en ingeniería de proteínas	1
1.1 Ingeniería de proteínas	2
1.2 Métodos computacionales aplicados en ingeniería de proteínas	3
1.2.1 Métodos de análisis filogenéticos	4
1.2.2 Métodos de análisis de estructuras	4
1.2.3 Métodos de estudio de mutaciones	5
1.2.4 Métodos basados en minería de datos y aprendizaje supervisado . .	5
1.3 Minería de datos	7
1.4 Principales problemáticas en la ingeniería de proteínas	8
1.4.1 Diferentes respuestas, una misma solución	8
1.4.2 Codificaciones, cuál es la mejor alternativa?	9
1.4.3 Diseñar mutaciones, un arte poco apreciado	10
1.4.4 Los descartados tienen algo más que decir	10
1.5 Hipótesis	11
1.6 Objetivos	12
1.6.1 Objetivo general	12
1.6.2 Objetivos específicos	12
2 Modelos predictivos asociados a mutaciones puntuales en proteínas	15
2.1 Aprendizaje de Máquinas	16
2.1.1 Algoritmos de aprendizaje supervisado	17
2.1.2 Métodos basados en regresiones lineales	17
2.1.3 K-Vecinos Cercanos	18
2.1.4 Naive Bayes	18

2.1.5	Árboles de Decisión	19
2.1.6	Support Vector Machine (SVM)	19
2.1.7	Métodos de ensamble	20
2.1.8	Redes Neuronales y Deep Learning	21
2.1.9	Meta-Learning	23
2.1.10	Medidas de desempeño	24
2.1.11	Validación de modelos	24
2.2	Herramientas computacionales asociadas a evaluación de mutaciones	24
2.2.1	Herramientas necesarias para la caracterización de los set de datos .	27
2.3	Hipótesis	29
2.4	Objetivos	30
2.4.1	Objetivo general	30
2.4.2	Objetivos específicos	30
2.5	Metodología propuesta	31
2.5.1	Preparación de set de datos	31
2.5.2	Implementación de meta modelos de clasificación/regresión	33
2.5.3	Cómo usar los meta modelos para la clasificación de nuevos ejemplos?	38
2.5.4	Uso de meta modelos en sistemas de proteínas	38
2.6	Análisis y evaluación de los set de datos a utilizar	39
2.6.1	Set de datos utilizados	39
3	Digitalización de secuencias lineales de proteínas aplicadas al reconocimiento de patrones y modelos predictivos	47
3.1	Metodologías asociadas a la codificación de variables categóricas	48
3.1.1	One Hot encoder	49
3.1.2	Ordinal encoder	49
3.1.3	Frecuencias de residuos	49
3.1.4	Uso de propiedades fisicoquímicas	50
3.1.5	Codificación de residuos con adición de información de su entorno .	50
3.2	Transformaciones de Fourier	52
3.2.1	Transformada rápida de Fourier (FFT)	53
3.2.2	Uso de Transformadas de Fourier en digitalización de propiedades fisicoquímicas	54
3.3	Clustering	55
3.3.1	Algoritmos de Clustering	56
3.3.2	Métodos de clustering empleando estructuras de grafos	56
3.4	Hipótesis	58

3.5	Objetivos	58
3.5.1	Objetivo general	59
3.5.2	Objetivos específicos	59
3.6	Metodología	59
3.6.1	Identificación y selección de propiedades fisicoquímicas	59
3.6.2	Codificación de secuencias lineales	60
3.6.3	Caracterización del espectro de frecuencias	61
3.6.4	Implementación de modelos de clasificación/regresión para análisis de variantes	63
3.6.5	Aplicación de técnicas de clustering sobre espectros de frecuencia .	64
4	Filogenética, propiedades fisicoquímicas y minería de datos aplicados al diseño de mutaciones en secuencias de proteínas	67
4.1	Hipótesis	68
4.2	Objetivos	68
4.2.1	Objetivo general	68
4.2.2	Objetivos específicos	69
4.3	Metodología propuesta	69
4.3.1	Conjunto de datos	70
4.3.2	Digitalización de secuencias lineales	71
4.3.3	Entrenamiento de modelos	71
4.3.4	Diseño de mutaciones	72
4.3.5	Implementación herramienta computacional	73
4.3.6	Consideraciones generales	78
	Referencias	81

Lista de figuras

1.1	Esquema representativo de los pasos que contempla la evolución dirigida .	3
1.2	Componentes principales de la minería de datos	7
2.1	Representación esquemática de una Red Neuronal	22
2.2	Esquema representativo asociado al proceso de generación de set de datos de mutaciones puntuales en proteínas.	31
2.3	Esquema representativo asociado al proceso de creación de meta modelos utilizando la metodología reportada para la herramienta MLSTools (Paper en redacción).	33
2.4	Esquema representativo de flujo asociado a la herramienta de generación de meta modelos para mutaciones puntuales en proteínas de interés.	38
2.5	Representación de estructuras de proteínas ejemplos utilizadas para el desarrollo de meta modelos de clasificación.	42
2.6	Evaluación del desbalance de clases en proteínas ejemplo.	44
2.7	Evaluación de la distribución de respuesta continua en set de datos de proteínas.	45
3.1	Esquema representativo de los pasos asociados al algoritmo FFT, desarrollado por Cooley [54]	54
3.2	Esquema representativo, metodología de digitalización de secuencias. . . .	60
3.3	Esquema representativo, metodología de clustering se secuencias por medio de espectros de frecuencias basados en propiedades fisicoquímicas.	63
3.4	Esquema representativo, metodología de clustering se secuencias por medio de espectros de frecuencias basados en propiedades fisicoquímicas.	65
4.1	Esquema representativo de la metodología propuesta para el diseño de mutaciones aplicando herramienta computacional a desarrollar	70
4.2	Esquema representativo interfaz de creación de jobs.	74
4.3	Esquema representativo interfaz de búsqueda de jobs.	75
4.4	Esquema representativo interfaz de descripción general del conjunto de datos.	76

4.5	Esquema representativo interfaz de visualización de propiedades fisicoquímicas.	77
4.6	Esquema representativo interfaz de visualización de espectros de frecuencias y residuos relevantes.	77
4.7	Esquema representativo interfaz de visualización de los meta modelos y sus medidas de desempeño.	78

Lista de tablas

2.1	Principales herramientas computacionales enfocadas a la evaluación de la estabilidad o predicción de cambios en la energía libre, asociado a mutaciones puntuales en proteína.	27
2.2	Tabla resumen, algoritmos implementados, parámetros utilizados e iteraciones involucradas por cada algoritmo.	34
2.3	Resumen de proteínas utilizadas para el desarrollo de meta modelos basados en metodología Meta Learning System propuesta durante este capítulo. . .	41
3.1	Cuadro resumen de algoritmos de aprendizaje supervisado	57

Chapter 1

Aplicaciones de la minería de datos en ingeniería de proteínas

La ingeniería de proteínas, es una de las ramas más relevantes y de mayor impacto en el campo de la biotecnología. Su objetivo principal, se basa en el diseño de mutaciones enfocadas en adicionar características específicas o mejorar sus propiedades fisicoquímicas, ya sea para someterlas a distintos tipos de ambientes, adecuarla a interactuar con diferentes elementos, presentar una mayor estabilidad, entre las principales. [129].

Los diseños de mutaciones se resumen en dos técnicas principales: El diseño racional [43] y la evolución dirigida [10], ambas técnicas experimentales, que cumplen con el mismo objetivo, relacionado a alterar las propiedades de la proteína para provocar una mejora con respecto a la estructura inicial.

A pesar de que ambas técnicas son utilizadas día a día, en diferentes investigaciones, éstas presentan limitantes importantes, en particular, relacionadas con el espacio de búsqueda posible a explorar, el tiempo que conlleva realizar diferentes pruebas y el costo económico y de recursos humanos que implica evaluar diferentes mutaciones [168].

Diferentes métodos computacionales han sido desarrollados, enfocados principalmente en el estudio de proteínas y el apoyo al diseño de mutaciones. Estas herramientas, se centran en el análisis de la estructura y la estabilidad termodinámica de sustituciones, adiciones o eliminaciones de residuos [182, 114, 154, 151]. Sin embargo, debido a cómo estos funcionan, en ocasiones, el costo computacional es muy elevado, ya que escala linealmente con respecto a la cantidad de pruebas a realizar.

Por otro lado, métodos basados en técnicas de minería de datos y aprendizaje de máquinas, se presentan como una alternativa potente y de costo computacional reducido, siendo capaces de generar resultados a partir de conocimiento existente, ya sea, para entrenar modelos que

permitan evaluar mutaciones desde puntos de vista de estabilidad, identificación de residuos relevantes en propiedades fisicoquímicas, evaluar propensiones a cambios, etc. [40, 80, 86].

No obstante, dado el gran volumen de datos existente en la actualidad ¿cómo es posible reconocer qué dato es relevante y cuál no?, ¿cómo puedo representar la información existente?, ¿cómo puedo complementar las diversas técnicas experimentales desarrolladas con el enfoque de la minería de datos?, etc., son interrogantes que nacen a partir del uso de técnicas computacionales y el apoyo a las metodologías experimentales.

Dado lo anterior, durante el presente capítulo, se expone el concepto de ingeniería de proteínas y cuáles son las principales técnicas experimentales involucradas en el diseño de mutaciones. Además, se introduce el concepto de Minería de datos y se exponen diversas metodologías computacionales que han sido desarrolladas, enfocadas en este campo de investigación. Por último, se presentan diferentes problemas que dan el punto de partida a cada una de las propuestas metodológicas de esta tesis doctoral y que denotan la evidente necesidad de ser desarrolladas o que exponen cambios en los puntos de vista actuales asociados al desarrollo de modelos y las convierten en un aporte significativo a los estudios actuales de mutaciones y el diseño de proteínas.

Adicional a ello, se exponen la hipótesis y objetivos generales que se relacionan estrechamente con cada capítulo siguiente en este proyecto y que permiten visualizar los diferentes puntos de vista y utilidades asociadas a las metodologías a proponer durante este trabajo de tesis doctoral.

1.1 Ingeniería de proteínas

Por definición, la ingeniería de proteínas es un campo de investigación centrado en el diseño y creación de proteínas útiles o con propiedades fisicoquímicas relevantes, con un enfoque principal en la comprensión del plegamiento de proteínas [126].

Actualmente, la ingeniería de proteínas, cuenta con dos estrategias principales para la construcción de mutaciones. Siendo éstas la evolución dirigida y el diseño racional. Sin embargo, a pesar de su existencia, normalmente no son excluyentes, por lo que es común utilizar ambas metodologías. No obstante, un estudio completo de residuos y una evaluación detallada de las sustituciones es una limitante importante para estas dos técnicas, debido tanto a recursos económicos como humanos [126].

La evolución dirigida imita el proceso de selección natural, permitiendo direccionar la evolución hacia objetivos definidos, reflejados en cuanto a funciones o propiedades fisicoquímicas deseables [128, 10]. Una representación del proceso, se observa en la Figura 1.1.

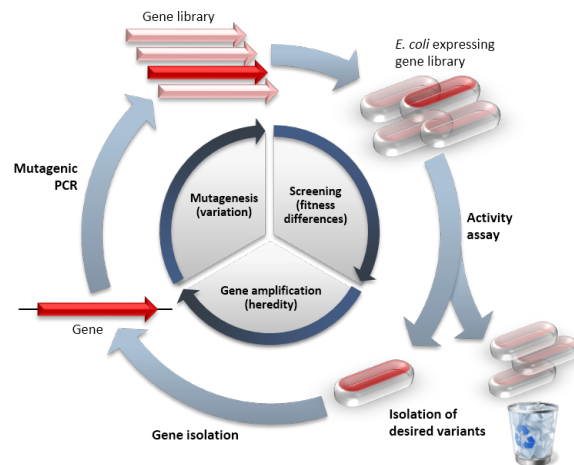


Fig. 1.1 Esquema representativo de los pasos que contempla la evolución dirigida

El proceso, de manera general, consiste en someter un gen de interés a rondas iterativas de mutagénesis, con el fin de crear una biblioteca de variantes. A partir de dicho conjunto de elementos, se seleccionan las variantes con la función deseada. Finalmente, se aíslan y se amplifican para formar una plantilla para la siguiente iteración. Así, el proceso sigue iterando y estadísticamente se seleccionan las más favorables y aquellas que tendieron a la evolución debido a la supervivencia en el proceso [10].

Con respecto al diseño racional de proteínas. Ésta es una técnica ampliamente utilizada y al igual que la evolución dirigida, presenta el objetivo general de generar variantes con alguna función de interés o características particulares. No obstante, exhibe una diferencia relevante, la cual se centra en la información que debe existir sobre la estructura, mecanismos, plegamiento o secuencia lineal de la proteína de interés [43].

1.2 Métodos computacionales aplicados en ingeniería de proteínas

Diferentes métodos computacionales han sido desarrollados para distintos análisis, con el fin de poder responder diferentes interrogantes planteadas desde enfoques distintos, ya sea, para estudiar secuencias lineales, filogenia y motivos conservados, análisis de estructuras, modelamiento estructural, estudio de mutaciones, etc. A continuación, se listan algunos de los principales enfoques y qué herramientas existen para su desarrollo.

1.2.1 Métodos de análisis filogenéticos

Los análisis filogenéticos se centran en el estudio de secuencias lineales de proteínas o genes, con el fin de identificar parentesco, motivos conservados o identificación de dominios. Las principales metodologías se basan en realizar alineamientos de secuencia para identificar o reconocer identidades de la secuencia estudio con respecto a información reportada previamente.

Una de las herramientas más conocidas para el desarrollo de alineamientos es Blast [113], la cual permite hacer múltiples comparaciones de secuencias versus bases de datos con genes o proteínas, empleando algoritmos de alineamiento local [5].

Por otro lado, se encuentran los alineamientos múltiples, los cuales son utilizados en la comparación de secuencias lineales con el fin de encontrar patrones o motivos conservados, o, la identificación de parentesco [189]. Esto último, es muy empleado cuando se analizan secuencias de organismos no reportados y cuya función se trata de aclarar. Una de las herramientas más utilizadas en esta área es el software Mega [187].

A su vez, el uso de los conceptos filogenéticos, ha sido utilizado para el estudio de propensiones de sustituciones aminoacídicas y cómo éstas afectan a la función de la proteína o con respecto a sus propiedades fisicoquímicas. Considerando esto, herramientas como MOSST [151], han permitido comprender la propensión de residuos y posiciones relevantes en secuencias, sólo estudiando conjuntos de elementos sin la necesidad de conocer estructuras tridimensionales de las proteínas.

Estos estudios, generan las bases para el análisis de secuencias y se basan en que sólo se necesita la secuencia lineal a estudiar y a partir de ella, es factible comprender un panorama relacionado a patrones de conservación, tendencias, relaciones evolutivas o inclusive, propensiones y posiciones relevantes. No obstante, no son los únicos, ya que existen herramientas computacionales que facilitan la predicción de la estructura secundaria, funcionalidad, etc., siendo una de las principales herramientas Swiss-Prot [23].

1.2.2 Métodos de análisis de estructuras

Los métodos de análisis de estructuras, tienen el objetivo de comprender patrones de interacción, efectos de energía y estudiar diferentes propiedades fisicoquímicas y termodinámicas, a partir de la estructura tridimensional de una proteína, la cual puede ser obtenida por cristalografía de rayos X o por medio de resonancia magnética nuclear.

No obstante, también, es factible el desarrollo de modelos de proteínas a partir de secuencias lineales, técnica conocida como Modelamiento por homología. Diferentes software,

permiten la implementación de esta técnica, dentro de los cuales se encuentran SWISS-MODEL [89], IntFOLD [138], ROSETTA [120], MODELLER [68], entre los principales.

Por otro lado, existen diferentes métodos computacionales que permiten el estudio de interacción entre proteína y una molécula o proteína-proteína, los cuales, principalmente se enfocan en el uso de técnicas como docking o dinámicas moleculares, con el fin de estudiar los posibles residuos que participan en la interacción, evaluándose a nivel energético y midiendo el desempeño en términos de error. A su vez, técnicas basadas en simulaciones moleculares, permiten comprender la interacción en sí y simular el comportamiento entre la molécula de interés y la proteína. Además, métodos computacionales basados en química cuántica, han sido utilizados para comprender fenómenos de interacción a una escala mucho más precisa. No obstante, estos son ampliamente más costosos y su uso es limitado al estudio de un número pequeño de átomos.

Existen diferentes herramientas que permiten hacer dinámicas moleculares, tales como: NAMD [161], AMBER [44]. etc., mientras que para la interacción entre moléculas, o docking, existen AutoDock [193], RosettaDock [130], GRAMM-X [191]. Además de la suite Maestro Schrödinger [169], la cual abarca funcionalidades para las diferentes acciones propuestas.

Distintos son los enfoques pueden ser considerados en el estudio de proteínas y en el análisis de su estructura. Sin embargo, los nombrados son los principales.

1.2.3 Métodos de estudio de mutaciones

De modo general, los estudios de mutaciones se basan principalmente en el análisis de la estructura ante los cambios de residuos o la adición o eliminación de estos, evaluando los cambios mediante diferencias de energía libre, entre la proteína inicial y la mutada. Herramientas como FoldX [182], SDM [154], Auto-Mute [137], etc., permiten analizar cómo afecta una mutación en términos energéticos, basándose para ello, en el uso de funciones de energía potencial y dinámicas moleculares asociadas a dicha sustitución. Sin embargo, el uso de este tipo de herramientas, conlleva un gran costo computacional debido a los diferentes cálculos que son requeridos. Durante el capítulo 2 se ahondarán más en estas herramientas. No obstante, en la Tabla 2.1 se resumen algunas de las principales herramientas utilizadas para este tipo de análisis.

1.2.4 Métodos basados en minería de datos y aprendizaje supervisado

La minería de datos y el aprendizaje de máquinas han sido utilizados en diferentes áreas de estudio de proteínas, ya sea para predicción de estructura secundaria [102, 143, 204],

análisis del efecto de mutaciones [39, 41, 195], identificador de patrones mediante métodos de clustering [181, 152], entre los principales ejemplos.

Diferentes enfoques han sido aplicados para obtener resultados relevantes, por un lado, se encuentra la utilización de algoritmos de aprendizaje supervisados clásicos como métodos de clasificación o predicción. Por otro, el uso de métodos de clustering para identificación de patrones basados en entornos no supervisados. Actualmente, se ha empleado el uso de redes neuronales y deep learning para manipulación de sistemas de datos complejos y se han enfocado principalmente en el estudio de predicción de interacciones y evaluación de estructura secundaria de una secuencia lineal.

A pesar de los diferentes objetivos, es necesario el desarrollo de conjuntos de datos que sean descritos mediante atributos, los cuales permitan alimentar estos modelos, para generar el aprendizaje. Distintas técnicas han sido utilizadas para caracterizar los ejemplos, dentro de las cuales principalmente se encuentran la codificación mediante One hot encoder [157], el uso de frecuencias de residuos [153] y la descripción empleando propiedades fisicoquímicas en conjunto con la caracterización del ambiente [39, 41].

Actualmente, la minería de datos y el aprendizaje automático, son una de las áreas de desarrollo de mayor interés, ya que, generan una disminución en cuanto al tiempo de cómputo y maximiza los espacios de búsqueda, los cuales, por medio de técnicas experimentales, es muy complejo analizarlas y empleando métodos computacionales para evaluar las interacciones, demandan un alto costo computacional.

Otro tipo de enfoque, se basan en el modelamiento matemático de proteínas empleando estructuras de grafos [37, 201], con el fin de aprovechar las características y ventajas que entrega este análisis, para poder descubrir patrones o estudiar interacciones por medio de la formación de aristas entre los diferentes nodos [135].

Enfoques particulares se han desarrollado con el fin estudiar estructuras específicas o regiones de interés, ejemplos como la identificación de epítopes en secuencias lineales de antígenos, son una de las problemáticas más relevantes y de mayor impacto en los últimos años [107, 149, 175]. Sin embargo, la complejidad es alta, dada la basta cantidad de información existente y a que las regiones con las que pueden interactuar presentan un espacio muestral del orden del 10^9 .

Los resultados satisfactorios obtenidos mediante la aplicación de técnicas de minería de datos y aprendizaje de máquinas, demuestran el poder de éstas en las diferentes áreas de estudio asociadas a la ingeniería de proteínas, convirtiéndola en una de las temáticas de mayor impacto en el último tiempo. Inclusive, permite ser un complemento relevante para investigaciones de alto impacto, tal es el caso de investigadores como Frances H. Arnold, quien en el último tiempo, ha enfocado su análisis de evoluciones dirigidas hacia un enfoque

de aprendizaje de máquinas, empleando diferentes técnicas asociadas al reconocimiento de patrones y la clasificación de estos [17, 209, 211].

1.3 Minería de datos

Minería de datos es el proceso de descubrimiento de patrones en set de datos, involucrando métodos asociados a Machine Learning [141], estadísticas y sistemas de bases de datos [92]. Se define como un subcampo interdisciplinario de la informática, el cual tiene por objetivo general extraer información (a través de métodos inteligentes) de un conjunto de datos y transformar la información en una estructura comprensible para su uso posterior [71, 66].

La minería de datos es el paso de análisis del proceso de *descubrimiento de conocimiento en bases de datos*, o KDD [70]. Además del análisis en bruto de los datos, también incluye aspectos de manipulación de bases de datos y pre procesamiento de estos, evaluaciones de modelo e inferencia, métricas de interés, consideraciones de complejidad, post procesamiento de estructuras descubiertas, visualización y actualización de la información [21].

En la Figura 1.2, se exponen las principales ramas que componen la minería de datos y los diferentes procesos que se asocian a dichas ramas.

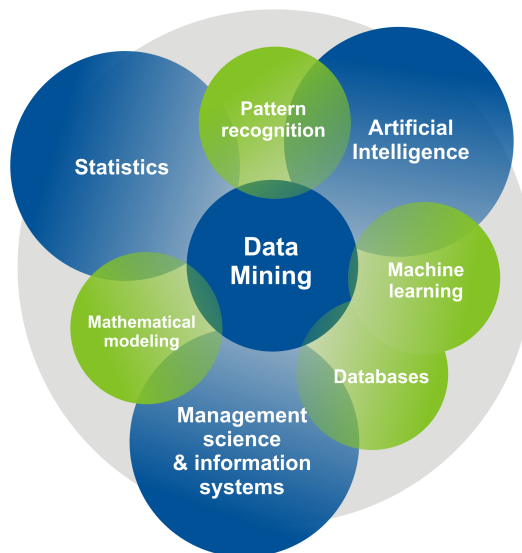


Fig. 1.2 Componentes principales de la minería de datos

Son tres las principales áreas que abarca la minería de datos: Estadística, Inteligencia Artificial y Manipulación de sistemas de información. Por otro lado, son distintos procesos los que interactúan entre estas ramas, tales como: Modelamiento Matemático, reconocimiento de patrones, Sistemas de almacenamiento persistente y machine learning [92].

Cada área en particular, tiene un objetivo general y diversos objetivos específicos. Sin embargo, estas áreas interactúan entre sí, con el fin de poder extraer patrones de información que generen conocimientos a partir de la data de procesada [21].

La minería de datos se utiliza en diferentes campos, tales como: genética y genómica [121, 167], ingeniería de proteínas [91, 123, 124], comercio y negocios [100], sistemas de tránsito [131], optimizaciones en procesos industriales [52, 81, 25], reconocimiento de patrones [104, 69], rasgos cuantificables en enfermedades [213, 148, 64] y más recientemente en áreas de dinámicas moleculares [48, 210] y parámetros para la generación de pipe lines automatizados de simulaciones cuánticas en sistemas químicos [134, 60, 165].

1.4 Principales problemáticas en la ingeniería de proteínas

Diferentes son las problemáticas que pueden existir en el campo de la ingeniería de proteínas, ya sea, desde la generación de herramientas computacionales para estudiar mutaciones y su efecto de manera masiva, hasta el diseño de mutaciones basados en secuencias lineales de proteínas. A continuación, se presentan diferentes problemáticas existentes en el área, algunas de las cuales serán motivos de estudio y desafíos a cumplir durante el presente trabajo.

1.4.1 Diferentes respuestas, una misma solución

El desarrollo de modelos de clasificación y/o regresión, es uno de los temas más recurrentes en el campo de la minería de datos y el aprendizaje de máquinas. Sin embargo, el hecho de asociar mutaciones a una respuesta, conlleva al problema de cómo caracterizarla, con el fin de alimentar a los algoritmos para ser entrenados.

A raíz de esto, cuáles son los mejores descriptores para una mutación?, desde qué puntos de vista se puede hacer una caracterización? y cuáles son más relevantes?, son interrogantes que se presentan a la hora de abordar su representación, siendo problemas que han sido tratados desde un largo tiempo sin lograr generar un consenso o una forma general de diseñar tal representación.

En un gran número de trabajos, en los cuales se ha evaluado la estabilidad de proteínas en torno a la mutación, se han utilizado descriptores termodinámicos y de ambiente para poder representar el elemento [39, 41]. A pesar de que los desempeños de los estimadores han sido aceptables y relativamente altos. Esta caracterización ¿podrá ser utilizada para mutaciones asociadas a riesgo clínico?, ¿Existirá una correlación entre la respuesta y las variables de interés?, ¿Cómo afecta al desempeño del modelo la existencia de diferentes ejemplos

asociados a distintas proteínas en un único conjunto de datos?, etc., son interrogantes que nacen a la hora de plantearse la situación.

Dado a lo anterior, y con el objetivo de generar un aporte significativo al desarrollo de estimadores basados en aprendizaje de máquinas, se ha propuesto adicionar el concepto de filogenia a la descripción de mutaciones y disgregar los conjuntos de elementos para ser tratados por proteínas independientes, esto con el fin de generar modelos de clasificación y/o regresión proteína-específicos, los cuales puedan ser aplicados a diferentes respuestas de interés ya sea: efectos en mutaciones, estabilidad, actividad, productividad, etc., siendo éste, el tema central a abordar en el capítulo 2.

1.4.2 Codificaciones, cuál es la mejor alternativa?

A menudo, el uso de secuencias lineales de proteínas se relaciona a la identificación de patrones o evaluación de variantes para una misma proteína. Actuales herramientas bioinformáticas permiten el uso de la secuencia de manera directa y por medio de alineamientos de secuencias o modelamiento a través del uso de Cadenas de Markov, facilitan el reconocimiento de patrones o la evaluación de mutaciones. No obstante, para la aplicación de métodos basados en minería de datos, ya sea la identificación de clusters o el entrenamiento de modelos, se requiere codificar la secuencia.

Existen diferentes codificaciones posibles, ya sea, para representar la secuencia o para la caracterización de mutaciones. A pesar de ello, no existe un consenso asociado a qué técnica utilizar. Cada una presenta sus pros y contra. No obstante, la cantidad de información involucrada varía entre ellas. Sin embargo, a mayor información, incrementa el número de dimensiones a tratar, aumentando la complejidad del problema. Esto implica, utilizar técnicas de reducción de dimensionalidad para seleccionar las dimensiones con mayor variabilidad en el conjunto de datos.

Una de las codificaciones más novedosas ha sido el uso de las propiedades fisicoquímicas de los residuos y su digitalización mediante transformadas de Fourier. Esto ha permitido la identificación de residuos claves en la propiedad en estudio y soluciona el problema del efecto del ambiente de los elementos participantes.

En vista de las necesidades de desarrollo de modelos de clasificación/regresión o la identificación de residuos claves y la generación de sistemas de clustering para secuencias lineales de proteína, con el fin de apoyar al diseño de mutaciones, análisis de variantes e inclusive caracterización de secuencias, sin tener conocimiento sobre su estructura. Se propone el uso de transformadas de Fourier como método de digitalización de propiedades fisicoquímicas para el desarrollo de conjuntos de datos que permitan ser entrenados para

el desarrollo de estimadores o identificar patrones, siendo el tema central a abordar en el capítulo 3.

1.4.3 Diseñar mutaciones, un arte poco apreciado

Diseñar mutaciones de manera eficiente, con una identificación adecuada de la propiedad en estudio o funcionalidad a adicionar, sin incurrir en grandes costos económicos y de recursos, es uno de los *Santos griales* de la ingeniería de proteínas. Como se nombró previamente, son dos los enfoques los que utilizan actualmente: Evolución dirigida y diseño racional de proteínas.

Ambas técnicas tienen sus ventajas y desventajas. No obstante, poseen en común una demanda en tiempo elevada y se requiere de conocimientos elevados sobre la estructura para poder diseñar las mutaciones, al menos, para el caso de diseño racional.

Enfoques computacionales han sido propuestos, con el fin de minimizar los costos económicos, contemplando evaluaciones energéticas asociadas a los residuos y cómo estos afectan a la estabilidad. No obstante, no pueden ser utilizados en secuencias lineales. Además, dejan de lado el concepto filogenético en el estudio, resultado un gap entre ambos puntos de vista. Por otro lado, métodos basados en la minería de datos, sólo se han centrado en identificación de residuos o el entrenamiento de modelos para predecir estabilidad.

A partir de lo anterior, y con el fin de generar un aporte significativo en el área de diseño, se ha considerado esta problemática como un foco central y culminante para el desarrollo de este trabajo, proponiendo así, la implementación de herramientas computacionales, basadas en técnicas de minería de datos y aprendizaje de máquinas, que permitan proponer mutaciones dado un conjunto de variantes con respuesta conocida. Generando la codificación de la secuencia por medio del uso de propiedades fisicoquímicas y su respectiva digitalización a través de transformadas de Fourier, seleccionando las propiedades más relevantes por medio de la aplicación de técnicas de reducción de dimensionalidad, para así, entrenar modelos de clasificación o regresión y posterior a ello, proponer mutaciones enfocadas en un filtro, aplicando herramientas de análisis de estabilidad y propensión. Toda esta problemática, el planteamiento de la metodología y qué se utilizará para llevar a cabo, se abordará en el capítulo 4.

1.4.4 Los descartados tienen algo más que decir

En la técnica de evolución dirigida, la selección de residuos o variantes, se basa en si presentan la característica deseable o no, o si aumenta la propiedad. Si el residuo no provoca el efecto deseado, éste es descartado, ya que no cumple con el criterio de selección.

Sin embargo, es posible pensar que, combinaciones lineales de residuos pueden provocar una sinergia en alguna propiedad, generando el resultado deseado. No obstante, el estudio de dichas combinaciones, o mejor dicho, las correlaciones asociativas existentes entre mutaciones no son consideradas, ya que, sólo se seleccionan aquellos que cumplen con dicho criterio. Pero, ¿qué pasa con aquellos residuos que son descartados y que al ser mutados al mismo tiempo con otro elemento provocan el efecto deseado, e inclusive, con mejores resultados que los obtenidos por los seleccionados?, ¿Existe información asociada a conjuntos de mutaciones que provoquen este efecto?, ¿Será posible idear una metodología *in-silico* que permita comprender este tipo correlaciones y justificar los resultados esperados?.

Como se puede comprender, este fenómeno no ha sido explotado desde el punto de vista de minería de datos, debido principalmente, a que no existen reportes de conjuntos de datos con dichas características y esto es dado a que no ha sido un foco de estudio central. Sin embargo, se cree que es una necesidad inminente, la comprensión de estos mecanismos, ya que, aumenta el espacio de búsqueda y posterior diseño de mutaciones, en un gran número de dimensiones. Además, si bien resultados de este estilo no han sido reportados, sí, a partir de experiencias de diferentes grupos con enfoque en diseño de mutaciones y evolución dirigida, han observado que residuos no seleccionables por si solos, en combinación con otro elemento, permiten obtener la característica deseable.

A pesar de que esta problemática, no se considera dentro de los temas de estudio en sí, se plantea la discusión y se propone como un problema a ser tratado en el corto plazo, debido a las grandes implicatorias que esto puede conllevar y a las expectativas que se pueden generar al respecto, siendo de utilidad a la hora de proponer nuevas mutaciones y generar un aporte significativo en el área de ingeniería de proteínas.

1.5 Hipótesis

En base a las herramientas computacionales existentes y a los problemas expuestos previamente, además, tomando en consideración los avances en minería de datos y aprendizaje de máquinas. Se propone la siguiente hipótesis.

Las técnicas de minería de datos y reconocimiento de patrones pueden ser utilizadas para el estudio y diseño de mutaciones in-silico, considerando tanto el desarrollo de modelos de evaluación como la creación de nuevas herramientas computacionales que permitan proponer variantes dada la información existente.

Se plantea una hipótesis general, la cual abarca los diferentes o considera los planteamientos de problemáticas expuestos. Además, se menciona que la estructura de este proyecto

es un conjunto de metodologías independientes. Es decir, cada capítulo en sí (2, 3 y 4) son herramientas computacionales independientes y tratan de resolver una problemática de las planteadas, por lo que, cada uno presenta en sí, su hipótesis y objetivos correspondientes. No obstante, a pesar de su independencia, tienen relación profunda con el abordaje de las soluciones a partir de técnicas de minería de datos, además, que es posible combinarlas, para desarrollar una suite de librerías de apoyo a la ingeniería de proteínas.

1.6 Objetivos

Continuando con la lógica expuesta previamente, es decir, cada uno de los siguientes capítulos resuelve una de las problemáticas planteadas, y contempla en sí, una herramienta computacional por sí sola. Se plantea a continuación el objetivo general.

1.6.1 Objetivo general

Diseñar e implementar una suite de herramientas computacionales basada en técnicas de minería de datos, aprendizaje de máquinas y reconocimiento de patrones, enfocada en el estudio de mutaciones, que permita ser un aporte sustancial en el campo de ingeniería de proteínas.

1.6.2 Objetivos específicos

Como se expuso previamente, cada siguiente capítulo corresponde a una herramienta en sí, que formará parte de esta gran suite computacional de apoyo al estudio de mutaciones *in-silico*. Dado esto, se plantean los siguientes objetivos específicos.

1. Diseñar, implementar y testear, herramientas computacionales, inspiradas en la estrategia de Meta-learning, para la evaluación de mutaciones puntuales en proteínas específicas, considerando como descriptores, propiedades termodinámicas, estructurales y conceptos filogenéticos.
2. Modelar, implementar y evaluar, herramientas computacionales para la codificación de secuencias lineales de proteínas, empleando digitalización de propiedades fisicoquímicas, por medio de transformadas de Fourier, las cuales permitan la identificación de residuos claves y la aplicación de algoritmos de aprendizaje supervisado y clustering, para el entrenamiento de modelos y el reconocimiento de patrones.

3. Diseñar, implementar y testear, herramientas computacionales para el diseño de mutaciones *in-silico*, basadas en técnicas de minería de datos y reconocimiento de patrones, enfocadas en secuencias lineales y modelos de aprendizaje supervisado, cuyos descriptores sean espectros de frecuencia basados en transformadas de Fourier y sean constituidas por herramientas de filtro, que aseguren la estabilidad de la proteína y la propensión al cambio de cada aminoácido en base a conceptos filogenéticos.

Tal como se puede observar, los aprendizajes y competencias adquiridas al cumplir el objetivo 1 y 2, se utilizan en el desarrollo del objetivo 3. Lo cual denota una especie de dependencia entre las metodologías a plantear. Sin embargo, cada uno de los objetivos, corresponde a una herramienta computacional independiente, la cual podrá ser utilizada para los fines que el usuario estime conveniente. El conjunto de éstas, se asocia al desarrollo de la suite computacional de estudio de mutaciones empleando técnicas de minería de datos, reconocimiento de patrones y aprendizaje de máquinas.

Cada uno de los siguientes capítulos, presenta su propio marco teórico, además de hipótesis y objetivos, asociados a una metodología que trata de cumplirlos. No obstante, todos enfocados en un mismo punto: desarrollo de herramientas de apoyo para el estudio de mutaciones.

Chapter 2

Modelos predictivos asociados a mutaciones puntuales en proteínas

El análisis del efecto de mutaciones puntuales en proteínas, es una de las problemáticas más estudiadas en los últimos años. Las investigaciones se enfocan principalmente en la evaluación de cambios en la estabilidad de la proteína mediante la variación de energía libre que la mutación provoca [182, 154, 170, 155].

Diferentes modelos predictivos han sido desarrollados para poder predecir cambios de energía libre, en base a algoritmos de aprendizaje supervisado o mediante técnicas de minería de datos y así, determinar el efecto de la mutación en set de proteínas de interés [163, 42, 32, 114, 195, 82, 40]. No obstante, en casos más específicos, se han desarrollado modelos para proteínas independientes, con el fin de asociar la mutación a un rasgo clínico, particularmente, enfocado a casos de cáncer [86, 74], cambios en termo estabilidad [190], propiedades geométricas [14], entre las principales.

Sin importar el uso o la respuesta de los modelos, es necesario construir set de datos con ejemplos etiquetados, es decir, cuya respuesta sea conocida para poder entrenar modelos basados en algoritmos de aprendizaje supervisado y así evaluar su desempeño. Los enfoques principales al desarrollo de descriptores se basan en propiedades fisicoquímicas y termodinámicas, así como también, el ambiente bajo el cual se encuentra la mutación [40], ya sea a partir de la información estructural o sólo considerando la secuencia lineal. Sin embargo, no son considerados, los componentes asociados a conceptos filogenéticos y la propensión a cambios de dicha mutación, generando un gap entre ambos puntos de vista [151].

Dado a los modelos existentes y en vista a la necesidad de generar nuevos sistemas de predicción para mutaciones puntuales en proteínas, en respuesta al aumento considerable de reportes en los últimos años, se propone una nueva metodología para el diseño e implementación de modelos predictivos en mutaciones puntuales de proteínas.

Las mutaciones son descritas desde los puntos de vista estructural, termodinámico y filogenético. El desarrollo de los predictores es inspirado en el concepto de Meta Learning y es apoyado con técnicas estadísticas, tanto para la selección de modelos como para la evaluación de medidas de desempeño, entregando como resultado, un conjunto de modelos para las mutaciones puntuales reportadas, unificados en un único meta modelo.

Esta metodología será aplicada para generar estimadores en diferentes proteínas con mutaciones reportadas con respuesta conocida, como por ejemplo: evaluando las diferencias de energía libre que provoca la mutación y clasificaciones para evaluar si la sustitución de residuos aumenta o disminuye la estabilidad. A su vez, se implementarán modelos de clasificación para determinar la propensión clínica en un conjunto de mutaciones conocidas relacionados con el gen *pVHL*, responsable de la enfermedad von Hippel-Lindau, con el fin de exponer la versatilidad de la metodología y los problemas relevantes a set de datos altamente no-lineales.

A continuación, se describen los principales conceptos relacionados a aprendizaje supervisado, seguido de algunas herramientas computacionales para el análisis de mutaciones y su relevancia en la de estabilidad de una proteína, continuando con la metodología propuesta, la caracterización de los diferentes set de datos a utilizar y resultados parciales obtenidos al aplicar esta metodología.

2.1 Aprendizaje de Máquinas

Aprendizaje de Máquina, es una rama de la inteligencia artificial que tiene por objetivo el desarrollo de técnicas que permitan a los computadores aprender, es decir, generalizar comportamientos a partir de una información no estructurada suministrada en forma de ejemplos [141]. Aplicándose en diferentes campos de investigación: motores de búsqueda [55], diagnósticos médicos [50, 1], detección de fraude en el uso de tarjetas de crédito, bioinformática [119], reconocimiento de patrones en imágenes [67] y textos [145, 4], etc.

Los algoritmos de aprendizaje pueden clasificarse en dos grandes grupos [141]:

- **Supervisados:** se cumple un rol de predicción, clasificación, asignación, etc. a un conjunto de elementos con características similares, por lo que los datos de entrada son conocidos.
- **No Supervisados:** su objetivo es agrupar en conjuntos con características similares los elementos de entrada dado los valores de estos atributos, en base a la asociación de patrones característicos que representen sus comportamientos.

A continuación se describen en forma general, los algoritmos de aprendizaje supervisados utilizados para el desarrollo de la metodología, explicando los conceptos bajo los que se basan y cómo estos entrenan y se emplean para predecir o clasificar nuevos ejemplos.

2.1.1 Algoritmos de aprendizaje supervisado

Existen diferentes algoritmos de aprendizaje supervisado, los cuales pueden ser asociados a la clasificación de un elemento o la predicción de valores, dependiendo el tipo de respuesta existente en el conjunto de datos a estudiar. En el caso de respuestas con distribución continua, se trabajan con algoritmos de regresión, mientras que si la respuesta es binaria o multiclase y es representada por variables categóricas, los algoritmos se basan en clasificadores [141].

A su vez, también se pueden dividir con respecto a la forma en que se trata el problema, existiendo algoritmos basados en cálculos de distancia entre ejemplos (K-Vecinos Cercanos), otros que consideran transformaciones vectoriales y aplicaciones de funciones de kernel (Máquina Soporte de Vectores), así como también el uso de las características como entorno espacial de decisión (Árboles y métodos de ensamble) y aquellos que utilizan redes neuronales y trabajan en torno a cajas negras, o métodos basados en regresiones lineales, sólo aplicados a modelos predictivos de variables continuas.

Cada uno de estos algoritmos es descrito a continuación, enfocándose tanto en el componente matemático asociado, así como también en las ventajas y usos posibles que estos puedan tener, con respecto al conjunto de datos a trabajar.

2.1.2 Métodos basados en regresiones lineales

Regresión lineal, es uno de los métodos más simples en cuanto a predicción de variables continuas, además de uno de los más limitantes debido al sobreajuste que éste puede generar. No obstante, permite evaluar de manera simple y sencilla conjuntos de datos [88].

Matemáticamente, se espera que el conjunto de respuesta sea el resultado de una combinación lineal de parámetros, es decir. Sea \hat{y} el vector de predicciones, se tiene que:

$$\hat{y}(w, x) = w_0 + w_1x_1 + \dots + w_px_p \quad (2.1)$$

Donde w_0 es el intercepto y $w = (w_1, \dots, w_p)$ el vector de coeficientes.

Existente diferentes métodos de regresión lineal, los cuales cumplen con el mismo objetivo. Sin embargo, la forma en la que minimizan el error asociado a las diferencias entre los valores predichos y los observados.

Pese a su simplicidad, los métodos basados en regresiones lineales han sido ampliamente utilizados. No obstante, presentan diferentes problemas asociados al sobreajuste de parámetros.

2.1.3 K-Vecinos Cercanos

Algoritmo de aprendizaje supervisado, el cual tiene por objetivo asociar un elemento a una clase en particular, dada la información de ejemplos de entrada que tengan asociadas características particulares, que puedan declararse como *vecinos* del nuevo ejemplo a clasificar, siendo k el número de vecinos que se está dispuesto a utilizar para aplicar la clasificación [112]. La mejor elección de k depende fundamentalmente de los datos; generalmente, valores grandes de k reducen el efecto de ruido en la clasificación, pero crean límites entre clases parecidas.

Con el fin de evaluar la cercanía de los ejemplos existentes contra el nuevo ejemplo a clasificar, es necesario asociar ciertas medidas de distancia que permitan cuantificar esta característica, para así poder comparar esta distancia y evaluar la cercanía para asociarle una clase a este nuevo ejemplo [65]. La distancia a emplear para evaluar la cercanía puede ser: Euclidiana [61], Manhattan [159], coseno [125] o Mahalanobis [132], entre las principales.

2.1.4 Naive Bayes

Naive Bayes es un conjunto de algoritmos de aprendizaje supervisados basados en la aplicación del teorema de Bayes con la suposición "ingenua" de independencia entre cada par de características [214]. Dada una variable de clase y y un vector de característica dependientes de la forma x_1, \dots, x_n se puede utilizar la siguiente regla de clasificación:

$$P(y \mid x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i \mid y) \quad (2.2)$$

\Downarrow

$$\hat{y} = \arg \max_y P(y) \prod_{i=1}^n P(x_i \mid y) \quad (2.3)$$

A pesar de sus supuestos aparentemente simplificados, los clasificadores de Naive Bayes han funcionado bastante bien en muchas situaciones del mundo real, la famosa clasificación de documentos y el filtrado de spam son ejemplos de ello [122, 49, 140]. Requieren una pequeña cantidad de datos de entrenamiento para estimar los parámetros necesarios. Pueden ser extremadamente rápido en comparación con métodos más sofisticados.

Existen distintos tipos de clasificadores de Naive Bayes, diferenciándose entre sí en la función de distribución de probabilidad que utilizan [140, 108, 133], dentro de los que se encuentran: Gaussian Naive Bayes, Multinomial Naive Bayes y Bernoulli Naive Bayes.

2.1.5 Árboles de Decisión

Se define árbol de decisión como un modelo de predicción, utilizado en el ámbito de la inteligencia artificial, en el cual, dado un conjunto de datos, se fabrican diagramas de construcciones lógicas, muy similares a los sistemas de predicción basados en reglas, que sirven para representar y categorizar una serie de condiciones que ocurren de forma sucesiva, para la resolución de un problema [76].

El aprendizaje basado en árboles de decisión utiliza un árbol como un modelo predictivo que mapea las observaciones de las características que presenta un elemento. En estas estructuras de árbol, las hojas representan etiquetas de conjuntos ya clasificados, los nodos, a su vez, nombres o identificadores de los atributos y las ramas representan posibles valores para dichos atributos [22].

Este tipo de entrenamiento, es uno de los más utilizados, debido a su simplicidad a la forma en la que trabaja, ya que, permite comprender del problema, con respecto a los atributos y cómo estos van distribuyendo las respuestas, así, es posible entender las decisiones que toma el algoritmo para clasificar o predecir nuevos ejemplos, determinar comportamientos preferentes y tendencias sobre atributos y rangos de estos.

2.1.6 Support Vector Machine (SVM)

Máquina soporte de vectores (SVM por sus siglas en inglés), es un conjunto de métodos de aprendizaje supervisado, utilizados para clasificar, predecir e inclusive para la detección de puntos outliers [179].

SVM genera una representación de los ejemplos como puntos en el espacio, mapeados de modo que los ejemplos de las categorías separadas se dividan por un espacio claro que es tan amplio como sea posible. Nuevos ejemplos son entonces mapeados en ese mismo espacio y predicen si pertenecen a una categoría en base a qué lado del espacio son asignados [179].

Las predicciones se realizan de manera eficiente, utilizando funciones kernel para su transformación en espacios no lineales [6]. Esto permite generar transformaciones de espacio dimensional de los datos, para mapear implícitamente sus entradas en espacios característicos de alta dimensión.

2.1.7 Métodos de ensamble

Los métodos de ensamble, se basan en la combinación de las predicciones obtenidas por varios estimadores, construidos en base a algoritmos de aprendizaje supervisado, con el fin de mejorar la generalización del modelo y aumentar la robustez ante nuevos ejemplos [63].

Existen dos familias de métodos de ensamble, las cuales se diferencian principalmente en la forma en que combinan los modelos para obtener la medida de desempeño final [117]:

1. **Métodos ponderados:** basados en la construcción de varios estimadores independientes y promediar sus medidas de desempeño, esto mejora el rendimiento debido a que disminuye la variabilidad de las clasificaciones. Ejemplos comunes de esto son Bagging y Random Forest.
2. **Métodos boosting:** basados en la construcción secuencial de modelos, intentando disminuir el sesgo del modelo combinando diferentes estimadores débiles. Cumple con la filosofía *"la unión de varios modelos débiles, puede construir uno fuerte"*. Ejemplos comunes de esto son AdaBoost y Gradient Tree Boosting.

A continuación, se explican brevemente algunos de los algoritmos asociados a la familia de métodos de ensamble.

Bagging

Bagging forma parte de los métodos ponderados, en particular, se puede definir como métodos que forman una clase de algoritmos compuestos por varias instancias de un estimador, entrenados en base a subconjuntos aleatorios del set de datos original, ponderando sus predicciones individuales en una respuesta ponderada. El objetivo general de estos métodos es reducir la varianza de un estimador, por medio del proceso de entrenamiento de subconjuntos aleatorios [26].

Random Forest

Random Forest es un método de ensamble ponderado basado en árboles de decisión aleatorios. Conjuntos de diversos clasificadores son creados basados en efectos aleatorios tanto de la extracción de características como de ejemplos, formando subconjuntos de elementos, cada uno de estos aporta con un valor de estimación, el cual es ponderado con los restantes, obteniendo así, la medida general [27].

AdaBoost

AdaBoost, es un algoritmo basado en el método boosting, lo que implica que se ajusta a una secuencia de estimadores débiles obtenidos a partir de diferentes subconjuntos de datos generados de manera aleatoria desde el conjunto inicial de datos de entrenamiento [38].

Cada una de las predicciones obtenidas por los estimadores se combinan de manera ponderada por votación, en el caso de modelos de clasificación, o a través de un promedio en base a las estimaciones resultantes, en el caso de modelos de regresión [94].

Gradient Tree Boosting

Gradient Tree Boosting o Gradient Boosted Regression Trees, es una generalización de métodos de boosting para funciones diferenciables arbitrarias de pérdida [78]. Es un método considerado como preciso y efectivo, el cual puede usarse tanto para el desarrollo de modelos de clasificación como de regresión, siendo usado en diferentes áreas de investigación: motores de búsqueda, ecología, minerología, biotecnología, entre otros.

Dentro de las principales ventajas que posee, se encuentran: manejo natural de diferentes tipos de características en un set de datos, alto poder predictivo y robusto frente a la predicción de valores atípicos en una muestra [79].

2.1.8 Redes Neuronales y Deep Learning

Redes neuronales es posible definirlas como una serie de modelos de aprendizaje que se basan en la forma de trabajo de las redes neuronales biológicas, es decir, se usa el concepto de *neurona* para estimar una función aproximada, la cual dependerá de un largo número de inputs, generalmente desconocidos.

En la Figura 2.1 se aprecia un sistema de red neuronal, en la cual se observa un sistema interconectado por neuronas, las cuales intercambian información en forma de mensaje entre ellas, además cada interconexión tiene un peso, el cual es un valor numérico, que puede ser obtenido en base a la experiencia.

En resumen, una red neuronal es un conjunto de entradas y salidas regidas por capas intermedias que permiten evaluar la salida, dichas capas operan entre sí en base a funciones matemáticas y brindan un peso a la conexión, finalmente cada capa es usada para diseñar un modelo de aprendizaje supervisado o no.

Deep Learning es una herramienta de Machine Learning la cual tiene por objetivo modelar abstracciones de alto nivel en los datos por medio del uso de múltiples capas de procesamiento, ya sea a través del uso de estructuras complejas a través de múltiples transformaciones no lineales [19, 18, 62].

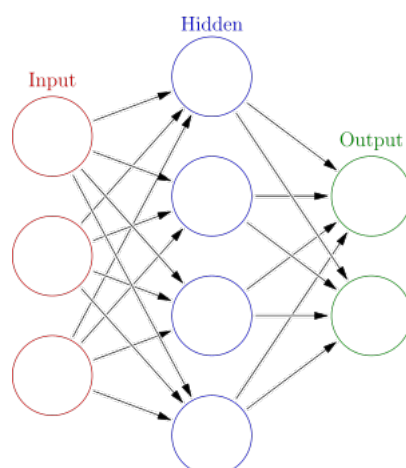


Fig. 2.1 Representación esquemática de una Red Neuronal

La investigación en esta área tiene por objetivo generar mejores representaciones y crear modelos para aprender de éstas a partir de datos no marcados a gran escala. En general, las representaciones obtenidas se inspiran en los avances en la neurociencia y se basa libremente en la interpretación de los patrones de procesamiento y comunicación de información en un sistema nervioso, como la codificación neural que intenta definir una relación entre varios estímulos y respuestas neuronales asociadas en el cerebro [18].

Deep learning es un método específico de machine learning el cual incorpora redes neuronales organizadas en capas consecutivas para poder aprender iterativamente utilizando un conjunto de datos. Deep learning es especialmente útil cuando se desea aprender patrones provenientes de datos no estructurados [62].

Posee diversas arquitecturas, tales como: deep learning network, matrices de convoluciones, redes neuronales recurrentes, etc. las cuales han sido utilizadas en visión artificial para el reconocimiento de patrones, aprendizaje de escritura, etc. Deep Learning es una herramienta de Machine Learning la cual tiene por objetivo modelar abstracciones alto nivel en los datos por medio del uso de múltiples capas de procesamiento, ya sea a través del uso de estructuras complejas a través de múltiples transformaciones no lineales [9].

Dentro de los principales algoritmos que son utilizados en redes neuronales se encuentran Back Propagation [96] y Multi Layer Perceptron [173].

Si bien en la actualidad, redes neuronales y Deep Learning son metodologías ampliamente utilizadas y han tenido resultados satisfactorios a la hora de trabajar en diferentes áreas de investigación, presentan un problema relevante al momento de aprender de los datos, los atributos y cómo estos facilitan o distribuyen la información.

Este problema es debido a que, los sistemas de redes neuronales trabajan en torno a información oculta, denominados, sistemas de cajas negras, lo cual, no permite comprender cómo se genera una nueva clasificación o predicción de elementos.

Lo anterior, no ocurre con métodos basados en estimaciones de distancia como KNN, uso de hiperplanos y funciones de kernel para la transformación espacial de los atributos como SVM, reglas de decisión que permiten representar estructuras de árbol que facilitan la comprensión de cómo distribuyen los atributos, rangos preferibles etc., como lo hacen los algoritmos basados en árboles de decisión e inclusive los métodos de ensamble. Es decir, diferentes algoritmos vislumbran cómo manipulan la información para llegar al resultado, no siendo el caso de las redes neuronales. Por lo tanto, dado a que se requiere de algoritmos que permitan la comprensión del problema y que posibiliten generar aprendizaje de los atributos, se descartan a priori el uso de métodos basados en redes neuronales y el uso de Deep Learning.

2.1.9 Meta-Learning

El Meta-Learning, es un campo relacionado a aprendizaje supervisado, inspirado en el concepto "*aprender a aprender*" [177], se basa en el uso de meta datos y presenta el objetivo de mejorar las medidas de desempeño de un clasificador, guiándolo entorno al aprendizaje. Por otro lados, estos elementos, son utilizados con el fin de comprender la flexibilidad del aprendizaje automático [194].

Uno de los puntos más relevantes en el aprendizaje automático, es la flexibilidad, esto es debido, a que algoritmos de aprendizaje supervisado, pueden diferir en el espacio o dominio bajo el cual funcionan, provocando que un conjunto de datos se adapte a un algoritmo y no suceda lo mismo con otro. Debido al uso de diferentes algoritmos, meta-learning permite mejorar la medida de desempeño y aumentar la robustez del modelo. No obstante, la selección de los algoritmos puede influenciar negativamente en el desempeño [98].

Al utilizar diferentes tipos de metadatos, es posible aprender, seleccionar, alterar o combinar diferentes algoritmos de aprendizaje para resolver efectivamente un problema determinado [177].

Existen diferentes enfoques asociados al Meta-Learning y distintas formas de uso, los ejemplos más conocidos se basan en el uso de redes neuronales recurrentes [8], el aprendizaje reforzado [178] y los métodos de ensamble. En los diferentes casos planteados, se construyen modelos y se mejoran con respecto a elementos previos, generando diferentes iteraciones del proceso, aumentando la medida de desempeño y adquiriendo características que previamente no se consideraban [72].

2.1.10 Medidas de desempeño

Medir el desempeño del modelo predictivo es importante a la hora de evaluar qué tan efectivo es el entrenamiento o la clasificación que se genera, existen medidas que sólo se basan en la cantidad de aciertos o errores que comete el clasificador, otras que implican la eficiencia del modelo y otras que se basan en la precisión. A su vez, es posible generar una división de las medidas con respecto al tipo de modelo, es decir, si se entrena un clasificador o un modelo predictivo. Para el primer caso, se tienen medidas como Accuracy, Recall, Precision y $F\beta$, mientras que para el segundo, se tienen Coeficiente de Pearson, Coeficiente de Spearman, Kendall τ rank, Coeficiente de determinación R^2 score y Error cuadrático medio, dentro de las principales.

2.1.11 Validación de modelos

La validación de los modelos trata los problemas de sobreajuste y la generalización, es decir, evitar desarrollar modelos que sólo tengan buenas métricas o medidas de desempeño para los datos de entrenamiento y no permitan clasificar nuevos ejemplos.

Con el fin de poder evitar esta problemática, normalmente los set de datos se dividen en 3 conjuntos: Entrenamiento, validación y testeo. Esto quiere decir, se considera una porción de elementos para entrenar el modelo, una segunda instancia para obtener las medidas de desempeño y una tercera con el fin de determinar si el clasificador entrega resultados acorde a las respuestas conocidas [115].

Existen técnicas que a partir del set de entrenamiento, generan múltiples divisiones, con el fin de entrenar subconjuntos de elementos del conjunto de entrenamiento y así obtener modelos ponderados, estas técnicas permiten prevenir el sobre ajuste, siendo la más conocida la Validación Cruzada [84], la cual, recibe un parámetro k , el cual determina el número de divisiones a realizar al conjunto de datos, cuando k es equivalente al número de ejemplos, se habla del caso de Leave one Out [198].

2.2 Herramientas computacionales asociadas a evaluación de mutaciones

Las herramientas computacionales asociadas a la evaluación de mutaciones puntuales se centran principalmente en el análisis de cómo ésta afecta a la estabilidad o la predicción de energía libre asociada a los residuos involucrados en la mutación. Sin embargo, a pesar de que el objetivo es el mismo, se enfocan en diferentes puntos de vista para abordar la

problemática, tanto a nivel de entrenamiento de modelos, cómo manipulación de set de datos, así como las técnicas utilizadas para la predicción de los cambios de energía libre.

En la Tabla 2.1 se exponen las principales herramientas existentes para la evaluación de la estabilidad de proteínas evaluando mutaciones puntuales, presentando las características, tipos de datos de entrada, resultados, estado de la herramienta y cuáles son las limitantes asociadas

Herr.	Características	Entradas	Salidas	Disp.
Foldx	Predice el valor del DDG a través del uso de funciones de energía derivados de términos fisicoquímicos, estadísticos e información estructural	Estructura en formato PDB e información sobre la mutación	Estimación de la diferenciade energía libre	Disponible mediante licencia académica
I-Mutant	Método basado en SVM para la predicción de DDG y la clasificación de la estabilidad de una proteína ante mutaciones puntuales. La mutación es caracterizada a través de propiedades estructurales y la información del ambiente. Permite la manipulación tanto de secuencias lineales como estructuras PDB	Secuencia lineal proteína, posición y mutación, en caso de existir estructura 3D, se requiere el archivo PDB	Predicción del DDG asociado a la mutación o clasificación de la mutación en estable o desestabilizante	Disponible para ejecución local

CUPSAT	Predice el DDG usando información estructural y del ambiente asociado a la mutación, además utiliza diferentes propiedades estructurales para estimar el valor de energía libre. Este es un método sólo basado en estimaciones utilizando técnicas de bioinformática estructural.	Estructura en formato PDB y la posición del residuo a mutar	Información sobre las 19 posibles sustituciones a realizar, referidas a términos como: ángulos de torsión, accesibilidad al solvente, tipo de estructura secundaria, dentro de las principales.	No disponible
Dmutant	Se basa en el uso de potenciales energéticos para entrenar modelos, utiliza distancias para describir el ambiente y estima el DDG asociado a la mutación	Estructura en formato PDB	Predicción del DDG asociado a la mutación	No disponible
AUTO-MUTE	Manipula las coordenadas de los residuos y aplica triangulación de Delaunay para formar geometrías, así permite describir el ambiente bajo el cual se encuentra el residuo. Los clasificadores se construyen entrenando con Random Forest y los predictores a través de árboles de decisión.	ID-PDB, cadena y mutación	Predicción DDG o clasificación de estabilidad, además de información relacionada al sector donde ocurre la mutación	Descargable para ejecución local

Table 2.1 Principales herramientas computacionales enfocadas a la evaluación de la estabilidad o predicción de cambios en la energía libre, asociado a mutaciones puntuales en proteína.

Diferentes son los puntos de vista que pueden ser considerados a la hora de evaluar el efecto que provoca la sustitución de residuos en la estabilidad de una proteína. Ya sea por medio de la estimación utilizando funciones de energía o potenciales energéticos (FoldX [182], Dmutant [216], CUPSAT [155]). Por otro lado, se encuentran métodos basados en algoritmos de aprendizaje supervisado. No obstante, estos pueden ser diferenciados con respecto al algoritmo de entrenamiento utilizado o a la forma de describir la mutación. Por ejemplo, I-Mutant [40], utiliza SVM para predecir o clasificar el efecto de la mutación. Mientras que, AUTO-MUTE [137] utiliza Random Forest para la clasificación del efecto y Árboles de decisión como algoritmo de regresión.

Cada uno entrega distintas ventajas y desventajas, las cuales se basan principalmente, en la precisión del método y en el tiempo de cómputo relacionado al procedimiento. No obstante, también influye la disponibilidad de estos y si permiten descargas de código fuente para ejecuciones locales o simplemente disponen de versiones web.

2.2.1 Herramientas necesarias para la caracterización de los set de datos

Adicional a las herramientas expuestas, se hace una descripción breve de SDM [154] y MOSST [151], las cuales serán utilizadas a lo largo de la metodología con el fin de poder caracterizar las mutaciones desde los puntos de vista termodinámico, aplicando SDM y filogenético, por medio de MOSST.

SDM

Site Directed Mutator (SDM) [154], es una de las herramientas más utilizadas a la hora de evaluar mutaciones puntuales en una proteína de interés y cuyo objetivo principal, es estudiar el efecto sobre la estabilidad de la proteína que provoca el cambio del residuo.

Se basa en potenciales estadísticos de funciones de energía, para obtener una cuantificación del efecto de una sustitución de un residuo. Este valor se representa por la diferencia de energía libre de Gibbs ($\Delta\Delta G$). Para estimar el efecto, utiliza ambientes específicos de frecuencias de sustituciones de aminoácidos, sin la utilización de familias de proteínas homólogas.

Para realizar la estimación, SDM recibe un archivo PDB en el cual se describe la estructura de la proteína inicial, seguido además de la mutación, la cual se describe mediante "*W-Pos-M*", los cuales corresponden a: Wilde residue, posición en la que se encuentra y mutate residue, respectivamente.

Sus medidas de desempeño en cuanto a la correlación entre los elementos predichos y los reales alcanza un 0.8, calculada a partir del testeo de la herramienta con mutaciones en la proteína Barnasa y staphylococcal nuclease. Este desempeño, lo convierte en una herramienta bastante precisa a la hora de estudiar nuevas predicciones.

Una de las principales razones de utilizar SDM y no otras herramientas computacionales, como las expuestas en la Tabla 2.1, es el hecho de basarse en potenciales estadísticos. Si bien, esta metodología no es tan precisa como el uso de potenciales físicos de funciones de energía, el hecho de utilizar simulaciones basadas en mecanismos de Monte Carlo, conlleva un gran costo computacional en horas de cálculo y recursos. No obstante, es mucho más eficiente que los métodos basados en Machine Learning, dado a que estos, normalmente, tienden a ajustarse y se requiere una gran cantidad de datos para entrenar estos modelos. Además, ya que estos se basan principalmente en métodos asociados a Support Vector Machine (SVM) o Redes Neuronales. Los primeros, no son capaces de adaptarse a espacios altamente no lineales. Mientras que los basados en redes neuronales, no permiten aprender de los atributos empleados, debido a la forma en cómo estos trabajan.

Además, SDM entrega resultados adicionales que permiten comprender el ambiente bajo el cual se produce la mutación y cómo, propiedades termodinámicas se ven afectadas ante la sustitución de residuos.

MOSST

Mutagenesis Objective Search and Selection Tool (MOSST) [151], es una herramienta que permite analizar una proteína de interés, con respecto a un conjunto de proteínas con relación filogenética, representadas en un alineamiento múltiple de secuencias. Esto con el fin, de poder detectar posiciones en la proteína de interés o target, que podrían ser mutadas para alterar o no las características de la misma.

Por otro lado, permite estimar mutagénesis relacionadas con la posibilidad de si un cambio de aminoácido dado tendría efectos perjudiciales sobre la proteína.

Además, como un uso alternativo, es factible la identificación de nsSNPs cuyos fenotipos son relevantes en una familia de genes, permitiendo separar, a aquellas sustituciones que no tienen implicaciones a nivel de funcionalidad.

MOSST aplica técnicas estadísticas para el análisis de las posiciones y se centran en comprender los efectos de las sustituciones desde el punto de vista filogenético. Esto es

relevante, ya que como input, sólo necesita secuencias lineales de proteínas, es decir, no es necesario el uso de estructuras en formato PDB, lo cual permite abarcar un mayor número de proteínas de estudio.

Actualmente MOSST se encuentra disponible vía ejecución local, implementado bajo lenguaje de programación Matlab. No obstante, posee una versión online para su uso, facilitando el acceso a la herramienta a diferentes estratos de público.

Utilizar la herramienta MOSST como generación de descriptores basados en propiedades filogenéticas, radica en las ventajas que ésta presenta. Por un lado, permite evaluar la propensión de la mutación y cómo ésta afecta a las características de la proteína en estudio y su familia. Además, permite entender conceptos relacionados a mutaciones, que no son considerados al utilizar propiedades termodinámicas y estructurales. Ya que, estos últimos, sólo evalúan el ambiente bajo el cual ocurre la mutación, mientras que MOSST, evalúa la propensión al cambio.

El uso de las herramientas MOSST y SDM permitirán describir las mutaciones desde los puntos de vista filogenético y termodinámico-estructural, de tal manera, las falencias de cada método, se complementan, facilitando la generación de descriptores relevantes para las mutaciones.

2.3 Hipótesis

En base a las herramientas existentes y en vista del aumento considerable de datos asociados a mutaciones en proteínas y el conocimiento de las respuestas que éstas generan, se evidencia la necesidad del desarrollo de herramientas computacionales o nuevos modelos de clasificación o regresión que faciliten el entrenamiento de proteínas singulares y la evaluación de sus mutaciones puntuales, con el fin de poder evaluar nuevos ejemplos y cuáles serían los efectos de estos, sin tener que recurrir en grandes costos económicos y tiempos de espera.

Dado esto se propone la siguiente hipótesis.

El uso de propiedades filogenéticas, termodinámicas y estructurales como descriptores de mutaciones permite el desarrollo de modelos predictivos inspirados en sistemas de meta-learning

Además de la hipótesis central surgen interrogantes como.

- Es posible utilizar estos nuevos modelos como herramientas para diagnóstico médico?

- Cómo se evalúan la robustez y la generalización de estos modelos, serán capaces de adaptarse a nuevos ejemplos?
- Es factible el desarrollo de una herramienta computacional que permita entrenar diferentes set de datos y que facilite la predicción de nuevos ejemplos?

2.4 Objetivos

En base a la hipótesis planteada y a las preguntas adicionales expuestas, se exponen a continuación el objetivo general y los objetivos específicos.

2.4.1 Objetivo general

Diseñar e implementar estrategias inspiradas en Meta Learning para la implementación de modelos de clasificación y regresión, asociados a mutaciones puntuales en proteínas de interés basados en descriptores termodinámicos, estructurales y filogenéticos.

2.4.2 Objetivos específicos

Dentro de los objetivos específicos se encuentran los siguientes.

1. Preparar y describir, por medio de propiedades termodinámicas, estructurales y filogenéticas, set de datos de mutaciones puntuales de proteínas con respuesta conocida expuestos en bibliografía o bases de datos públicas.
2. Implementar y evaluar metodología de meta learning para el diseño de meta modelos de clasificación y regresión de mutaciones puntuales aplicados a set de datos de proteínas generadas.
3. Diseñar e implementar herramientas computacionales que permitan el entrenamiento de set de datos y el uso de meta modelos para la evaluación de nuevos ejemplos.
4. Testear y evaluar comportamiento de las herramientas y los meta modelos en base a sistemas de datos que involucren mutaciones en proteínas con respuesta conocida.
5. Implementar modelos de clasificación para la relevancia clínica de mutaciones puntuales en proteína pVHL, asociada a la enfermedad von Hippel-Lindau.

2.5 Metodología propuesta

Con el fin de poder responder a la hipótesis planteada y dar solución a los objetivos impuestos, se propone una metodología general, en la cual, se consideran diferentes estrategias, implementaciones y evaluación de modelos. A continuación se explica la metodología propuesta y los componentes principales de ésta.

2.5.1 Preparación de set de datos

La preparación del set de datos consiste en obtener data para poder entrenar los modelos predictivos, la data se asocia a información de mutaciones en proteínas y la respuesta que ésta genera. En la Figura 2.2 se expone un esquema general con los pasos desarrollados para la preparación del set de datos.

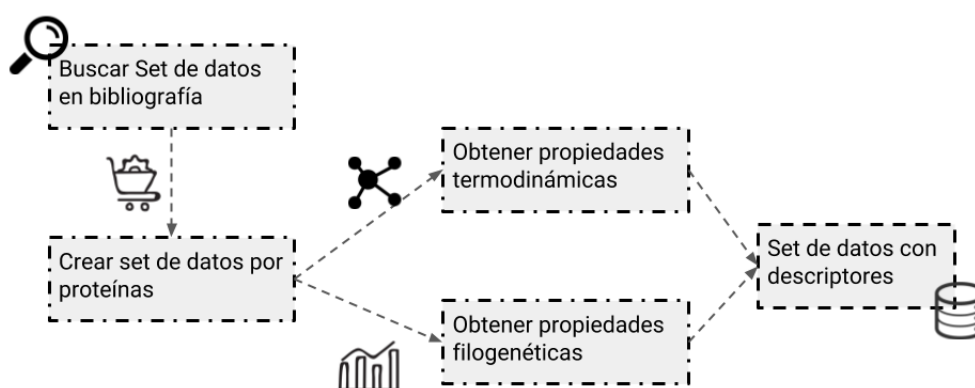


Fig. 2.2 Esquema representativo asociado al proceso de generación de set de datos de mutaciones puntuales en proteínas.

Tal como se expone en la Figura 2.2, los set de datos se buscan en la bibliografía, a partir de modelos desarrollados previamente, bases de datos, en la literatura, etc. El objetivo fundamental, es encontrar proteínas con mutaciones puntuales cuyo efecto sea conocido, dicha respuesta puede ser categórica, es decir, asociada al diseño e implementación de modelos de clasificación o continua y se aplica para modelos de regresión.

En una segunda instancia, a partir de la data recolectada, ésta se procesa con el fin de poder obtener set de datos de proteínas individuales con una cantidad de ejemplos considerables que permitan el diseño de modelos válidos, para ello, fueron implementados scripts bajo el

lenguaje de programación Python con el fin de recuperar las proteínas, obtener la información y generar la data de manera individual, además, eliminar ejemplos ambiguos. Es decir, filas con los mismos valores pero cuya columna de respuesta fuese diferente.

A partir de esto, se forman n set de datos asociados a n proteínas, cada uno con m ejemplos y cuyos descriptores consisten en el residuo original, posición en proteína, residuo mutado y la respuesta asociada. El desbalance de clases se analiza con respecto a las posibles categorías existentes en la respuesta y el porcentaje de representatividad que éstas poseen en la muestra. Se considera que el set de datos exhibe este comportamiento cuando presentan las características expuestas en la sección ???. En el caso de detectar esta problemática, el conjunto de datos será tratado empleando el método SMOTE (Synthetic Minority Oversampling Technique) [47].

Posteriormente, se aplican las herramientas SDM [154] y MOSST [151] con el fin de obtener los descriptores asociados a las propiedades termodinámicas y filogenéticas, respectivamente. Para ello, scripts implementados en lenguaje de programación Python, son desarrollados para consumir los servicios de dichas herramientas y registrar los resultados obtenidos, formando así, set de datos con los descriptores planteados en los objetivos iniciales. Un punto importante a destacar, es que el uso de SDM implica que las proteínas a trabajar, deben presentar una estructura 3D reportada en el Protein Data Bank [20] o al menos poseer un modelo representativo y validado. Esto es debido a que se utilizan informaciones de coordenadas para la estimación del efecto de la mutación, minimizaciones energéticas y estabilización de la mutante.

Ya con los descriptores formados, las características asociadas a variables categóricas son codificadas. Si la totalidad de posibles categorías supera el 20% del total de características en el set de datos, se aplica Ordinal Encoder, en caso contrario, One Hot Encoder [157]. Ordinal Encoder consiste en la transformación de variables categóricas en arreglos de números enteros con valores desde $0, \dots, n - 1$ para n posibles categorías. Por otro lado, One Hot Encoder, consiste en agregar tantas columnas como posibles categorías existan en el set de datos completadas mediante binarización de elementos (0 si la característica no se presencia, 1 en caso contrario)¹.

Es importante mencionar, que las respuestas asociadas a las mutaciones pueden ser del tipo continuo o categórico, lo cual implica que tanto los modelos como las métricas varían. No obstante, se aplica la metodología indistintamente, con el fin de demostrar la robustez del método y la eficacia de éste sin importar el tipo de modelo que se éste entrenando.

¹Más detalles sobre la codificación de variables categóricas y secuencias lineales de proteínas serán tratadas en el capítulo 3.

2.5.2 Implementación de meta modelos de clasificación/regresión

La implementación de meta modelos consiste en la obtención de un grupo de estimadores que en conjunto, permiten clasificar o predecir nuevos ejemplos. Para ello, se diseña e implementa una metodología inspirada en Sistemas de Meta Learning y aplicando técnicas estadísticas para la evaluación del desempeño y el uso del meta modelo con nuevos ejemplos.

En la Figura 2.3, se exponen las etapas asociadas a la implementación de meta modelos, contemplando desde la fase de entrenamiento de los modelos hasta la unión en meta clasificadores, lo cual se reporta en la herramienta MLSTools (Paper en redacción).

Cada una de las etapas, contempla un conjunto de scripts implementados en lenguaje de programación Python y empleando la librería Scikit-Learn para el entrenamiento y evaluación de los clasificadores o predictores [157], así como Numpy para el uso de módulos estadísticos [203].



Fig. 2.3 Esquema representativo asociado al proceso de creación de meta modelos utilizando la metodología reportada para la herramienta MLSTools (Paper en redacción).

Tal como se observa en la Figura 2.3, es posible identificar etapas claves en el proceso: Exploración de modelos, Selección y Generación de los meta clasificadores/predictores, junto con su evaluación. Cada una de estas etapas se exponen a continuación.

Exploración de modelos

La exploración de modelos o estimadores, se basa en la aplicación de diferentes algoritmos de aprendizaje supervisado con variaciones en sus parámetros de configuración inicial. La

utilización de los algoritmos, depende principalmente del tipo de respuesta que presente el set de datos, es decir, si es continua o categórica. No obstante, a modo resumen, en la Tabla 2.2 se exponen los algoritmos utilizados, el caso en el que se usan y los parámetros que se varían junto con el total de iteraciones posibles para cada elemento:

Algoritmos y parámetros empleados en la etapa de Exploración en MLSTools					
#	Algoritmo	Tipo	Parámetros	Uso	Iteraciones
1.	Adaboost	Ensamble	Algoritmo Número estimadores	Clasificación y Regresión	16
2.	Bagging	Ensamble	Bootstrap Número estimadores	Clasificación y Regresión	16
3.	Bernoulli Naive Bayes	Probabilístico	Default	Clasificación	1
4.	Decision Tree	Características	Criterio división Función de impureza	Clasificación y Regresión	4
5.	Gaussian Naive Bayes	Ensamble	Default	Clasificación y Regresión	1
6.	Gradient Tree Boosting	Ensamble	Función de pérdida Número estimadores	Clasificación y Regresión	16
7.	k-Nearest Neighbors	Distancias	Número Vecinos Algoritmo Métrica distanciaPesos	Clasificación y Regresión	160
9.	Nu Support Vector Machine	Kernel	Kernel Nu Grado polinomio	Clasificación y Regresión	240
10.	Random Forest	Ensamble	Número estimadores Función de impureza Bootstrap	Clasificación y Regresión	32
11.	Support Vector Machine	Kernel	Kernel C Grado polinomio	Clasificación y Regresión	240
Total Iteraciones					726

Table 2.2 Tabla resumen, algoritmos implementados, parámetros utilizados e iteraciones involucradas por cada algoritmo.

Como se observa en la Tabla 2.2, son sobre 720 modelos los que se generan y a partir de ellos se obtiene distribuciones de medidas de desempeño que permiten evaluarlos. En el caso de modelos de regresión se utilizan los coeficientes de Pearson, Spearman, Kendall τ y R^2 , mientras que para modelos de clasificación, se consideran la Precisión, Exactitud, Recall y F1.

Finalmente, esta etapa, entrega set de modelos entrenados y evaluados según las métricas de interés. Se destaca que cada modelo es validado a través del proceso de validación cruzada, con el fin de poder disminuir posibles sobreajustes. El valor de k asociado a las subdivisiones a realizar varía con respecto a la cantidad de ejemplos que presente el set de datos, es decir, sea n la cantidad de ejemplos en la muestra, si $n \leq 20$ se tiene que $k = n$ implicando el uso de Leave one out, si $n > 20$ y $n \leq 50$ se considera un valor de $k = 3$, si $n > 50$ y $n \leq 100$ $k = 5$, por último, si $n > 100$ se tiene un valor de $k = 10$.

Selección de modelos

Cada distribución de medida de desempeño perteneciente a los modelos entrenados en la fase de Exploración, se somete a test estadísticos basados en Z-score [157] que permite seleccionar los modelos cuyas métricas representen outliers positivos dentro de la distribución.

El algoritmo general, utilizado para el desarrollo de esta selección es como se expone en el algoritmo 1, para el cual se detallan los pasos simplificados que permiten obtener un conjunto de modelos entrenados y que representan los valores más altos dentro de su distribución. Es importante mencionar, que se obtiene un conjunto M' con los modelos, considerando como punto de selecciones los valores evaluados con respecto a la desviación estándar, considerando los umbrales 3σ , 2σ y 1.5σ por sobre la media, si ningún factor se cumple, sólo se considera el valor máximo en la distribución.

Es importante mencionar, que cada distribución puede permitir la selección de distintos modelos, lo cual implica que un mismo modelo pueda ser seleccionado en diferentes medidas, razón por la cual, a la hora de obtener el conjunto de modelos M' se remueven aquellos elementos que se encuentran repetidos. Siendo estos, sólo los modelos que presenten igualdad tanto en el algoritmo como en sus parámetros de configuración inicial.

Algoritmo 1 Algoritmo de selección de modelos

Entrada: Conjunto M con modelos entrenados y sus medidas de desempeño, Lista L con medidas de desempeño.

Salida: Conjunto M' con modelos seleccionados.

```

1: para  $i$  en  $L$  hacer
2:   Calcular media  $\mu$ , desviación estándar  $\sigma$  en distribución  $M_i$ 
3:   para  $x \in M_i$  hacer
4:     si  $x \geq \mu + 3 * \sigma$  entonces
5:       Agregar  $x$  a  $M'$ 
6:     fin si
7:   fin para
8:   si largo  $M' = 0$  entonces
9:     para  $x \in M_i$  hacer
10:      si  $x \geq \mu + 2 * \sigma$  entonces
11:        Agregar  $x$  a  $M'$ 
12:      fin si
13:    fin para
14:    si largo  $M' = 0$  entonces
15:      para  $x \in M_i$  hacer
16:        si  $x \geq \mu + 1.5 * \sigma$  entonces
17:          Agregar  $x$  a  $M'$ 
18:        fin si
19:      fin para
20:      si largo  $M' = 0$  entonces
21:        para  $x \in M_i$  hacer
22:          si  $x = MAX M_i$  entonces
23:            Agregar  $x$  a  $M'$ 
24:          fin si
25:        fin para
26:      fin si
27:    fin si
28:  fin si
29: fin para
30: devolver  $D$  sin valores extremos

```

Generación y evaluación de meta modelos

A partir del conjunto de modelos M' , el cual representa los estimadores seleccionados, cuyas medidas de desempeño son las más altas en sus distribuciones correspondientes, se generan meta modelos, es decir, estimadores compuestos de diversas unidades, los cuales en conjunto entregan una respuesta, ya sea por ponderación o votación. El proceso general para la generación de los meta modelos, es descrito a continuación.

En una primera instancia, los modelos son nuevamente entrenados y se comparan las nuevas medidas de desempeño con las obtenidas previamente. En caso de que exista una diferencia mayor al 20%, en cualquiera de sus métricas, el modelo se remueve del conjunto M' . La razón fundamental de esto, es debido a que se espera desarrollar modelos robustos cuyas evaluaciones no presenten variaciones significativas y que realmente no alteren sus predicciones ante nuevos ejemplos, razón por la cual, se aplica nuevamente validación cruzada para validar los modelos.

Con el fin de evaluar el desempeño de los meta modelos, nuevas medidas se generan a partir de la información resultante de los modelos individuales. No obstante, la forma en la que se obtienen varían dependiendo del tipo de respuesta que se debe entregar.

Si la respuesta es continua, se obtiene los valores de predicción de cada modelo y se promedian, para luego aplicar las métricas estándar (Coeficiente de Pearson, Kendall τ , Spearman y R^2) sobre estos valores promediados y los reales. Expresado matemáticamente:

Sea M' la cantidad de elementos en el meta modelo, n la cantidad de ejemplos en el set de datos y sea Y el vector de respuestas reales de tamaño n . Para cada $M'_i \in M'$ se obtiene un vector Y_i que representa los valores de predicción entregados por el modelo M'_i . A partir de cada Y_i se genera una matriz de predicciones $P(m \times n)$ donde m representa la cantidad de modelos en M' . Finalmente, se obtiene un vector Y' de tamaño n , el cual se compone de la media de cada columna en la matriz P , es decir, para el ejemplo i se obtienen m predicciones, las cuales son promediadas, formando el valor $Y'_i \in Y'$. Vector el cual, se utiliza para obtener las métricas de desempeño.

Para el caso en que la respuesta sea categórica, es decir, los modelos son del tipo clasificación, se obtiene la respuesta de cada modelo individual y se selecciona una única categoría, correspondiente a aquella que presente una mayor probabilidad de ocurrencia dada la distribución de elementos y considerando para ello las probabilidades iniciales de cada categoría en el set de datos de estudio. De esta forma, se obtiene un vector respuesta con la clasificación de cada ejemplo cuyo valor corresponde al evento más probable a ocurrir, este vector se compara con el set de respuestas reales y se aplican las métricas de interés para clasificadores.

2.5.3 Cómo usar los meta modelos para la clasificación de nuevos ejemplos?

Nuevos ejemplos pueden ser clasificados o predecir su respuesta, dependiendo sea el caso, a partir de los meta modelos desarrollados. En el caso de estimadores basados en variables continuas, los nuevos ejemplos se someten a cada uno de los modelos individuales pertenecientes al sistema, los cuales generan una respuesta individual, a partir de dichas respuestas, se genera un intervalo de confianza con un nivel de significancia $\alpha = 0.01$ donde existe una mayor probabilidad de que se encuentre el valor real de la predicción dado los valores del entrenamiento. Para ejemplos que impliquen clasificación, se obtiene la respuesta de cada modelo individual y se evalúa la probabilidad de ocurrencia de cada categoría, entregando así, la respuesta condicionada por una probabilidad de ocurrencia del evento.

2.5.4 Uso de meta modelos en sistemas de proteínas

El objetivo principal de esta metodología, radica en el hecho de crear una herramienta que permita implementar modelos basados en algoritmos de aprendizaje supervisado para set de datos de mutaciones puntuales o variantes para una misma proteína.

Un flujo general del uso de la herramienta, se expone en la Figura 2.4.

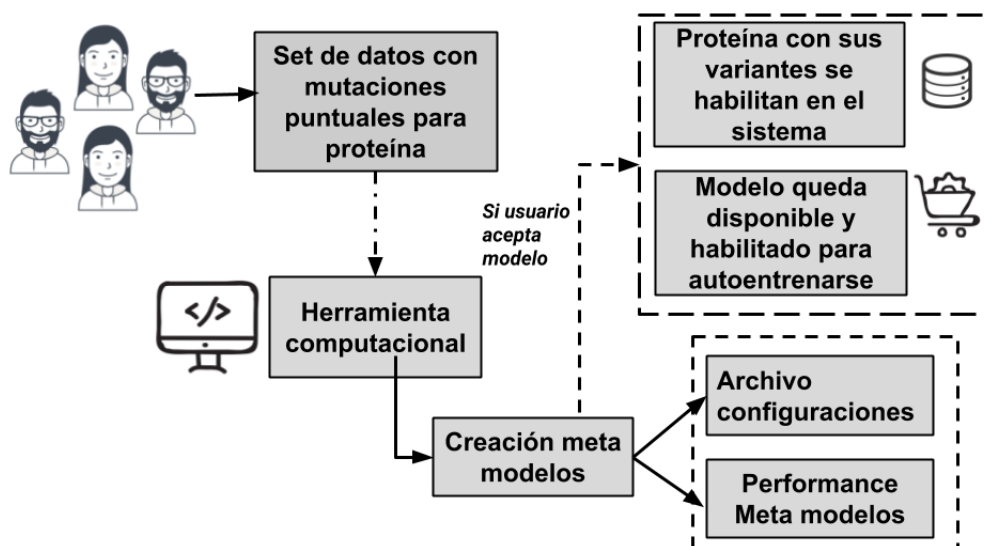


Fig. 2.4 Esquema representativo de flujo asociado a la herramienta de generación de meta modelos para mutaciones puntuales en proteínas de interés.

La idea general, consiste en que usuarios de la herramienta, puedan entrenar sus propios modelos de clasificación o regresión, basados en la metodología expuesta en los pasos

anteriores mediante el uso de Meta Learning Sytem Tools (Paper en Redacción). Para ello, los usuarios deben entregar sus set de datos con la información necesaria para ser procesada: cadena, residuo original, posición, residuo mutado y respuesta o efecto de la mutación, además del archivo PDB a ser procesado.

La herramienta, aplica los pasos expuestos en la metodología de este capítulo generando un meta modelo basado en algoritmos de aprendizaje supervisado y las medidas de desempeño que permiten evaluar el modelo obtenido.

Si el usuario acepta la metodología y permite la publicación de los datos, el sistema habilita el acceso tanto a los meta modelos como a los set de datos y los agrega a la lista de procesos de modelos auto entrenables.

Esto último, implica que ante la adición de nuevos ejemplos al set de datos, el sistema actualiza los modelos y las medidas de desempeño, aplicando la metodología expuesta, así, constantemente mantiene la actualización de la información y permite mantener en constante crecimiento los datos que contemplan el desarrollo de los modelos.

2.6 Análisis y evaluación de los set de datos a utilizar

A continuación, se exponen los resultados obtenidos hasta el momento, presentando principalmente, los set de datos a utilizar, las características que estos poseen y qué representan, con el fin de contextualizar la data a manipular y los modelos a generar.

2.6.1 Set de datos utilizados

En el presente apartado se describen las características básicas de los set de datos trabajados, así como también, qué representan las proteínas bajo las cuales se están desarrollando los modelos de estimadores.

Descripción general

Los set de datos utilizados, tanto para la formación de los inputs asociados al sistema, así como también la validación de respuesta correspondiente a la mutación que estos tienen, fueron extraídos desde distintas bases de datos de mutaciones en proteínas de estudios relacionados a los cambios que provoca la sustitución del residuo inicial, ya sea a nivel de cambios energéticos o estabilidad de la proteína.

11 set de datos con respuesta continua fueron obtenidos. Cada set de datos contemplaba como elemento a predecir, las diferencias de energía libre de Gibbs, entre los residuos

originales y mutados. Las mutaciones fueron seleccionadas desde diversos estudios en los cuales se reportaron, centrándose en [202, 186, 160, 3, Bordner and Abagyan].

Adicional a los set de datos con respuesta continua, 8 conjuntos de elementos asociados a tareas de clasificación fueron obtenidos desde diversos estudios reportados a la actualidad [7, 33, 41, 164, 40, 162, 114, 136, 83].

De tal manera, se generó un total de 19 conjuntos de set de datos, con respuesta categórica y continua, los cuales se asocian a proteínas independientes, usadas para la evaluación de las metodologías planteadas. Estas 19 proteínas junto con su descripción, se exponen en la Tabla 2.3.

Resumen set de datos de proteínas y sus características				
#	Código PDB	Tipo	Ejemplos	Descripción
1.	1A22	Regresión	132	Human growth hormone bound to single receptor
2.	1CH0	Regresión	191	Crystal and molecular structures of the complex of alpha-*Chymotrypsin with its inhibitor Turkey Ovomucoid third domain
3.	1DKT	Regresión	119	CKSHS1: Human cyclin dependent kinase subunit, type 1 complex with metavanadate
4.	1FKJ	Regresión	219	Atomic structure of FKBP12-FK506, an immunophilin immunosuppressant complex
5.	1FTG	Regresión	203	Structure of apoflavodoxin: closure of a Tyr/Trp aromatic gate leads to a compact fold
6.	1PPF	Regresión	190	X-Ray crystal structure of the complex of human leukocyte elastase and the third domain of the Turkey ovomucoid inhibitor
7.	1RX4	Regresión	556	Dihydrofolate reductase (E.C.1.5.1.3) complexed with 5,10-Dideazatetrahydrofolate and 2'-Monophosphadenosine 5'-Diphosphoribose
8.	1WQ5	Regresión	239	Crystal structure of tryptophan synthase alpha-subunit from Escherichia coli
9.	2AFG	Regresión	134	Human acidic fibroblast growth factor

10.	3SGB	Regresión	191	Structure of the complex of Streptomyces Griseus protease B and the Third domain of the Turkey ovomucoid inhibitor
11.	5AZU	Regresión	203	Crystal structure analysis of oxidize Pseudomonas Aeruginosa Azurin at PH 5.5 and PH 9.0. A PH-induced conformational Transition involves a peptide bond flip
12.	1BN1	Clasificación	1802	Carbonic anhydrase II inhibitor
13.	1BVC	Clasificación	561	Structure of a Biliverdin Apomyoglobin complex
14.	1LZ1	Clasificación	848	Human Lysozyme. Analysis of Non-Bonded and Hydrogen-Bond interactions
15.	1STN	Clasificación	2193	The crystal structure of Staphylococcal Nuclease
16.	1VQB	Clasificación	820	Gene V Protein (Single-Stranded DNA Binding Protein)
17.	2CI2	Clasificación	741	Crystal and molecular structure of the Serine proteinase inhibitor CI-2 from Barley seeds
18.	2LZM	Clasificación	2336	Structure of Bacteriophage T4 Lysosyme
19.	2RN2	Clasificación	712	Structural details of ribonuclease H from Escherichia Coli

Table 2.3 Resumen de proteínas utilizadas para el desarrollo de meta modelos basados en metodología Meta Learning System propuesta durante este capítulo.

Cada una de las proteínas presentan diferentes características y funcionalidades, algunas facilitan la unión a DNA, mientras que otras presentan propiedades enzimáticas, por otro lado, existen enzimas que representan inhibidores, entre las principales. Esto es interesante a la hora de evaluar el poder que presenta la metodología con respecto al análisis de diferentes proteínas, estructuras y complejos, ya que se presenta una gran variedad en cuanto a forma y funcionalidad de éstas, lo que implica que el sistema no se limita por cierto tipo de estructuras o complejos.

A modo de ilustrar las diferencias estructurales de las proteínas en estudio, en la Figura 2.5 se exponen algunas de las estructuras asociadas a las proteínas utilizadas para desarrollar modelos de clasificación o regresión.

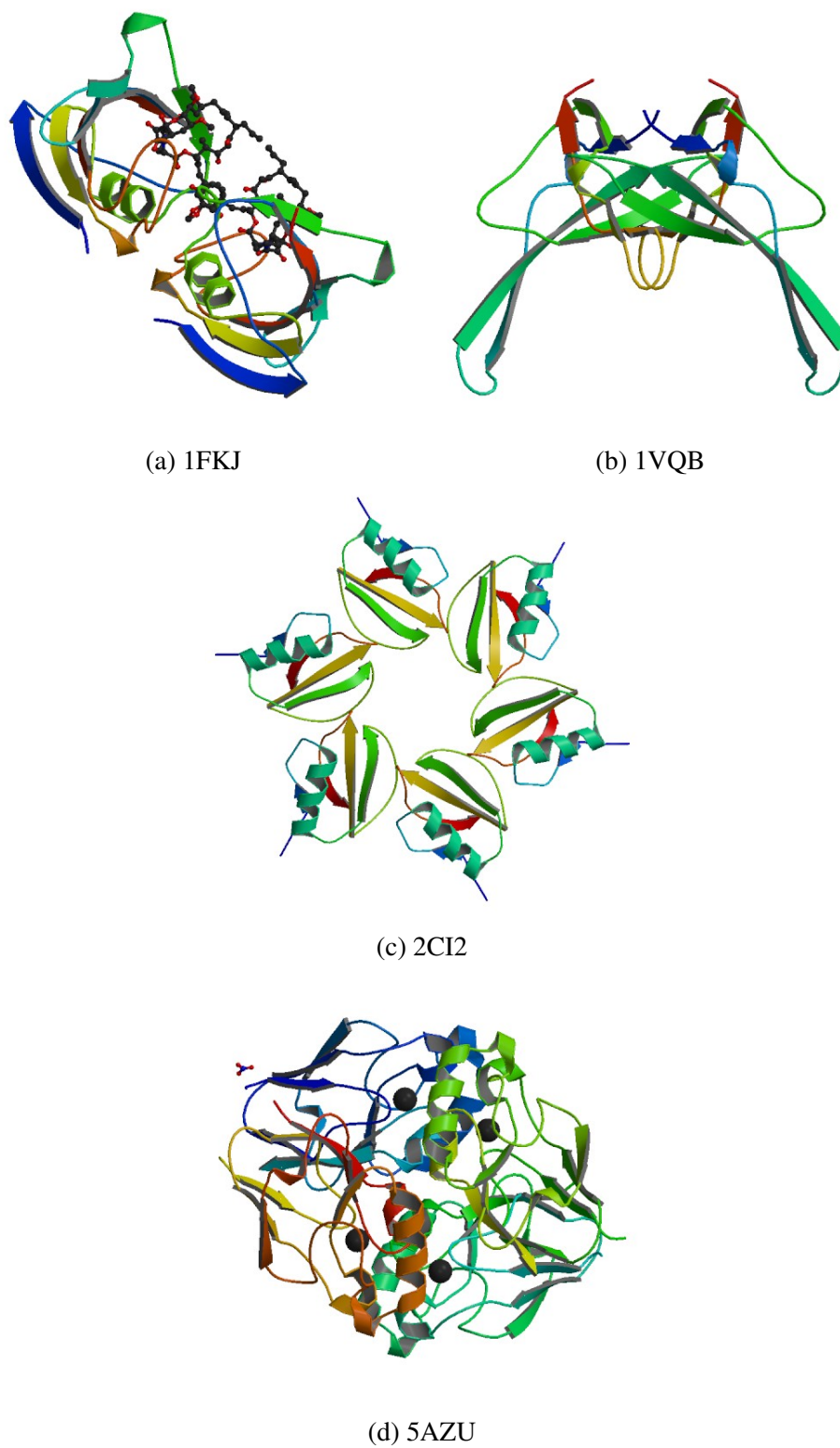


Fig. 2.5 Representación de estructuras de proteínas ejemplos utilizadas para el desarrollo de meta modelos de clasificación.

Las mutaciones fueron recolectadas desde diferentes set de datos, por lo que, en caso de información ambigua, es decir, una misma mutación con diferentes respuestas, no fueron consideradas. Por otro lado, debido a que para la aplicación de la herramienta SDM se necesitaba la cadena a la cual pertenece el residuo, scripts desarrollados en Python y utilizando la librería BioPython, permitieron procesar los archivos asociados a las estructuras de las proteínas, identificadas desde el Protein Data Bank (PDB) [2]. Descartando aquellas mutaciones reportadas en las que no se encontró la cadena, obteniendo como resultante la cantidad de mutaciones reportadas para cada proteína expuestas en la Tabla 2.3.

Evaluación del desbalance de clases y distribución de respuestas continuas

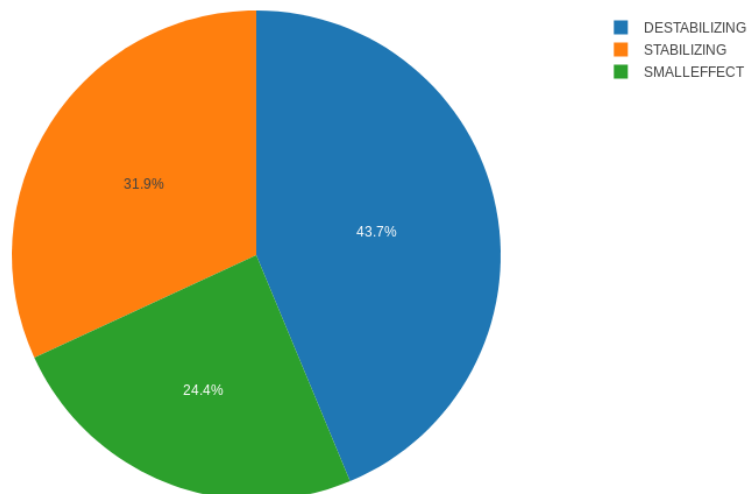
El desbalance de clases se evaluó en aquellos set de datos con respuesta categórica, lo cual contempla, tres posibles casos: Neutral, Estable, No Estable, esto es debido a que todos los estudios donde se evalúan mutaciones, normalmente se analizan cambios que alteren la estabilidad de la proteína. En la Figura 2.6, se aprecia a modo ejemplo dos set de datos y su distribución de categorías para la variable respuesta.

En las 8 proteínas en estudio para modelos de clasificación, la distribución de las categorías es similar a lo expuesta en la Figura 2.6 para todas ellas, donde cerca del 50% corresponden a mutaciones que afectan positivamente a la estabilidad, mientras que mutaciones que provocan cambios negativos o no generan diferencias, se encuentran en proporciones similares, ambas cercanas al 25%.

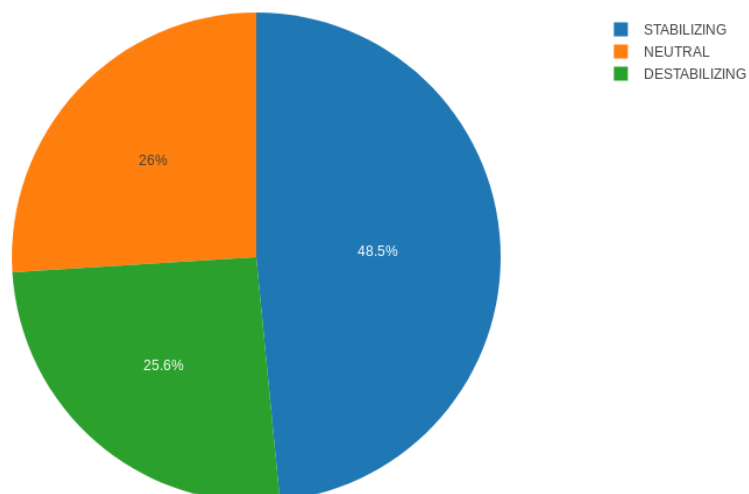
Si bien las proporciones son dispares, para este caso, se considera un desbalance como un elemento que representa menos de un 5% del total de la muestra, además, dado a que la cantidad de ejemplos son elevadas, un 20% o un 25% implica cerca de 200 mutaciones, en promedio, que cumplen dicha característica. También, el hecho de que exista una cantidad inferior de mutaciones no benéficas a la estabilidad viene dada a la dificultad de encontrar y reportar mutaciones que afecten negativamente a una proteína, es debido a la propensión filogenética [151] que estos ocurran, lo cual se ve reflejado en las diferencias asociadas a cambios positivos dentro del set de mutaciones. No obstante, si bien el hecho de que la propensión filogenética indique que el cambio tiende a mejorar estabilidad, diseñar mutantes con mejoras en propiedades de interés, es un problema latente en la actualidad, de alto costo económico y computacional y con una gran demanda desde diferentes áreas del conocimiento.

En los set de datos para el desarrollo de modelos de regresión, se evaluó la distribución de la respuesta, en este caso, valores de $\Delta\Delta G$ asociado a diferencias de energía libre producidas entre el residuo mutado y el original, tal que: $\Delta Res_{mut} - \Delta Res_{wild} = \Delta\Delta G$.

Las distribuciones se evaluaron utilizando el test de Shapiro, con el fin de determinar si la distribución se comportaba como una normal. Para todas las proteínas estudiadas, en



(a) 1STN

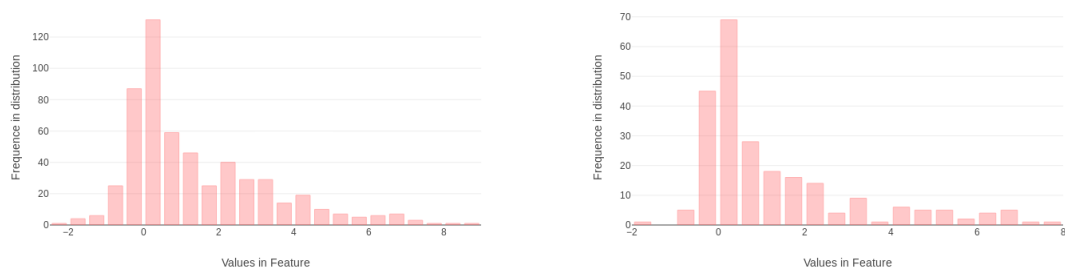


(b) 2RN2

Fig. 2.6 Evaluación del desbalance de clases en proteínas ejemplo.

los 11 set de datos, las respuestas presentaron distribución normal, con valores de Shapiro sobre 0.8 y un $p\text{-value} \leq 0.01$, lo cual indica una alta confianza estadística en los resultados presentados por dicho test.

Una visualización de las distribuciones puede generarse a partir del desarrollo de histogramas, los cuales, a modo de ejemplo se expone en la Figura 2.7.



(a) Histograma para respuesta continua en 1RX4 (b) Histograma para respuesta continua en 1WQ5

Fig. 2.7 Evaluación de la distribución de respuesta continua en set de datos de proteínas.

El análisis de estas características es relevante a la hora de diseñar modelos de clasificación o regresión, debido a que si existe una tendencia por una clase condicional al clasificador a *"aprender en base a la mayoría"*, por lo que puede aumentar los errores en cuanto a falsos positivos, dado a que, no se tiene la cantidad de ejemplos suficientes para una clase que permitan al modelo capturar las posibles variaciones asociadas a ésta.

Dado a los análisis de evaluación de representatividad de categorías en el set de datos y distribución de respuestas continuas, se expone que los set de datos seleccionados no presentan desbalance significativo para el caso de desarrollo de modelos de clasificación y a su vez, todas las respuestas asociadas a cambios en la energía libre para modelos de regresión, presentan distribución normal. Razón por la cual, es factible el desarrollo de modelos asociados a las respuestas presentes en los set de datos seleccionados. No obstante, sólo se ha considerado el problema del desbalance y la evaluación de distribución en la respuesta continua, una vez caracterizado los set de datos a partir de las propiedades fisicoquímicas y termodinámicas, se analizarán las características y cómo éstas condicionan la clasificación o la predicción de cambios energéticos.

Chapter 3

Digitalización de secuencias lineales de proteínas aplicadas al reconocimiento de patrones y modelos predictivos

Desarrollar modelos predictivos basados en algoritmos de aprendizaje supervisado, o, la identificación de patrones aplicando técnicas de clustering, son tareas muy relevantes a la hora de trabajar con secuencias de proteínas, ya sea para identificar grupos con características comunes o entrenar modelos predictivos de respuestas de interés. En ambos casos, se requiere el uso de conjuntos de datos altamente informativos y con características numéricas para poder utilizar los métodos implementados en las librerías actuales [157].

Diferentes metodologías se han implementado, para manipular las variables categóricas en set de datos y lograr su codificación numérica. Enfoques basados en adición de columnas según las categorías o simple transformación empleando representaciones en conjuntos naturales, suelen ser utilizados. No obstante, generan bastante discusión sobre las nuevas representaciones y, a su vez, el hecho de aumentar el número de columnas, conlleva a incrementar las dimensiones del conjunto de datos, provocando efectos negativos en los desempeños de los algoritmos [157].

Particularmente, en secuencias de proteínas, se han utilizado las frecuencias de incidencia de los residuos para codificarlos, la cual, pese a su simplicidad, ha resultado ser efectiva en diferentes casos de uso [153]. No obstante, este tipo de codificación, no permite explorar el ambiente bajo el cual se encuentran los residuos y tampoco considera el efecto de propiedades fisicoquímicas ni termodinámicas.

En diferentes estudios, los residuos se describen a partir de sus propiedades fisicoquímicas y adicional a ello, se emplea información que permite describir el ambiente del residuo a caracterizar, empleando binarizaciones para describir los residuos cercanos, ya sea por medio

del uso de un rango espacial, utilizando modelos o estructuras tridimensionales en donde se representan las coordenadas espaciales de los residuos, o empleando un rango lineal en secuencias lineales de proteínas [39, 41].

Un enfoque basado en las propiedades fisicoquímicas en combinación con la aplicación de transformaciones de Fourier, ha permitido demostrar que ciertos residuos permiten entregar las características asociadas a la propiedad en estudio, además, facilita comprender el aporte del ambiente sobre estos y representa una forma de estudio novedosa para el uso de información de secuencias lineales. Siendo una metodología ampliamente utilizada para identificar residuos que aporten a la propiedad, por medio de la representación de señales asociadas al espacio de frecuencias [199, 57, 35].

A pesar de ser una metodología interesante a la hora de estudiar secuencias lineales, exhiben problemas notorios sobre la selección de las propiedades relevantes a analizar, ya que, existe un número considerablemente alto de propiedades posibles a utilizar, descritas principalmente en la base de datos AAIndex [111], y es factible que diferentes familias de proteínas exhiban comportamientos notoriamente no similares y diverjan en cuanto a las propiedades que puedan ser representativas, inclusive, a la hora de estudiar mutaciones en una misma proteína puede que no sólo una propiedad permita su caracterización, si no, que un conjunto pequeño de éstas [35].

En el presente capítulo, se exponen en detalle, diferentes formas de representar secuencias lineales de proteínas, seguido a su vez del planteamiento del uso de transformadas de Fourier para la digitalización de propiedades fisicoquímicas y cómo es posible utilizar éstas para la identificación de patrones en secuencias lineales o el desarrollo de modelos de clasificación/regresión y la exposición de casos de uso en diferentes proteínas de interés.

3.1 Metodologías asociadas a la codificación de variables categóricas

Diferentes metodologías existen para poder codificar variables categóricas, a su vez, para set de datos de proteínas con secuencias lineales, es factible utilizar sus propiedades fisicoquímicas o frecuencias de residuos. Las principales metodologías usadas a la fecha se exponen a continuación.

3.1.1 One Hot encoder

One Hot encoder, es una de las técnicas más utilizadas a la hora de codificar variables categóricas y se basa principalmente en la adición de columnas con respecto a las categorías existentes en un conjunto de datos [34].

Dado el vector x de tamaño n con m categorías, por definición, One Hot encoder agrega al conjunto de datos m columnas, tal que, por cada categoría se adiciona una nueva columna al set de datos. Las nuevas columnas se completan con una binarización de los elementos, indicando si el elemento x_i posee la categoría m_j con un valor 1 y en caso contrario 0.

3.1.2 Ordinal encoder

Ordinal encoder, es una simplificación de One Hot encoder, ya que, simplemente codifica las categorías con números en el conjunto $[0, m - 1]$. Es decir, sea el vector x de tamaño n con m categorías y sea M el espacio de las posibles categorías con $M = [m_1, \dots, m_m]$, y cuya codificación implica el vector $M' = [0, \dots, m - 1]$. \forall elemento que \in a x se obtiene su codificación a partir del elemento $M'(M(m_i))$ que corresponde a la codificación de la categoría en el espacio M [157].

3.1.3 Frecuencias de residuos

Una secuencia lineal de proteína, corresponde a un vector v de tamaño n donde cada elemento corresponde a un residuo que pertenece a la secuencia. El uso de esta información para alimentar modelos de clasificación o regresión conlleva la codificación de sus elementos. Sin embargo, a la hora de utilizar las codificaciones basadas en One Hot Encoder, el conjunto de datos no queda estándar en cuanto a sus dimensiones, ya que, el largo de las secuencias puede variar y a su vez, el número de columnas a agregar corresponde a $n \times 20$ dado a que son n residuos y el espacio muestral M es de tamaño 20 lo que genera un aumento considerable en la cantidad de dimensiones.

Con el fin de poder representar las secuencias lineales de proteínas, se idearon metodologías que consideran la frecuencia de aparición de los residuos en la secuencia, de tal manera, de poder codificarla en un vector de tamaño 20, donde cada elemento representa el número de incidencias del residuo dividido por el largo del vector. Así, cada elemento se encuentra en un rango $[0, 1]$ donde 0 indica no incidencia del residuo y 1, incidencia total [153].

El uso de las frecuencias de residuos, es una de las primeras aproximaciones a la codificación de secuencias lineales de proteínas. No obstante, en todos los casos donde han sido utilizadas, se agrega información adicional, que permite comprender diferentes compor-

tamientos y evalúa ciertas propiedades del entorno, razones por las cuales, se recomiendan utilizarlas en conjunto con otros descriptores.

3.1.4 Uso de propiedades fisicoquímicas

El uso de propiedades fisicoquímicas para describir un residuo, es ampliamente empleado en la generación de descriptores para conjuntos de datos en ingeniería de proteínas [39, 41]. Diversos enfoques y modelos han sido construidos o entrenados, contemplando información asociada a componentes termodinámicos del residuo, en particular, a la hora de describir residuos para evaluar cambios en la energía libre, relacionados a efectos en la estabilidad de una proteína [7, 33, 162].

Se han reportado cerca de 570 propiedades fisicoquímicas que pueden ser utilizadas para describir un residuo en una secuencia lineal de proteínas, almacenadas en la base de datos AAIndex [111]. A su vez, es posible caracterizar estos residuos empleando un conjunto de propiedades estructurales, termodinámicas e inclusive filogenéticas. Es decir, diferentes puntos de vista que permitan describir los residuos pertenecientes a una secuencia. Sin embargo, el hecho de seleccionar qué descriptores son relevantes y cuáles no, radica en un problema de evaluación de características, el cual es común, en el área de la minería de datos.

Dado al gran conjunto de propiedades existentes y a la diversidad de descriptores que pueden ser utilizados para un conjunto de secuencias lineales de proteínas, es necesaria una selección correcta de las características, las cuales permitan formar set de datos informativos y con una correlación mínima entre sus elementos.

Contemplando esta problemática, técnicas de reducción de dimensionalidad o análisis de características son las más utilizadas a la hora de seleccionar los descriptores más informativos para un conjunto de datos, siendo ejemplos de esto: Análisis de componentes principales (PCA), Mutual Information, Análisis de correlación y evaluación espacial de características por medio de Random Forest. No obstante, en ocasiones, el conocimiento sobre el problema es un factor relevante a considerar.

3.1.5 Codificación de residuos con adición de información de su entorno

Adicional a las técnicas explicadas previamente con respecto a las codificaciones existentes, en algunos casos, no sólo basta con una única codificación del residuo, si no, que es relevante adicionar información que puede ser importante para describir los residuos. Normalmente, junto con las codificaciones basadas en propiedades fisicoquímicas, se emplean técnicas que permitan describir el ambiente bajo el cual se encuentre el residuo [136].

En la gran mayoría de los casos, se adiciona información de los residuos cercanos al residuo de interés, esto depende del tipo de datos bajo el cual se esté trabajando, es decir, si son secuencias lineales o son estructuras de proteínas en formato PDB [41, 39].

Para el caso de que sean secuencias lineales, sea s secuencia de residuos de tamaño n y sea r_i el residuo de interés a evaluar su ambiente. Se crea una ventana de tamaño n' que contempla la cantidad de residuos r_j cercanos al residuo r_i , de tal manera que se crea un nuevo sub conjunto s' de datos de tamaño $2n'$ con n' residuos a la izquierda y n' a la derecha. El cual normalmente es codificado empleando binarización de elementos, así, en algunas ocasiones, a cada residuo, se le adicionan 20 descriptores que permiten indicar la ausencia o presencia de residuos cercanos a su entorno y el cual se completa con el conjunto de residuos s' [41].

Cuando se manejan estructuras de proteínas en formato PDB, la codificación y la evaluación del ambiente es similar. Sin embargo, en vez de utilizar una ventana de tamaño n' se utiliza un radio espacial de valor x para el cual, se toma el residuo y se estiman las distancias de los elementos cercanos, ya sea entorno a los carbonos α o a otros elementos. Esto, a diferencia de las secuencias lineales, permite adicionar información sobre las propiedades de distancia, ángulos y conformación de estabilidad por interacciones electrostáticas débiles que pueden generarse a partir de la proximidad de los elementos. No obstante, es una inferencia de su uso y se requieren de diferentes tipos de elementos que permitan caracterizar los eventos asociados al ambiente estructural asociado al residuo [41].

Actualmente, el uso de codificaciones mediante propiedades fisicoquímicas y el empleo de información adicional basada en descriptores de ambientes, es una de las metodologías más utilizadas a la hora de generar set de datos relacionados a mutaciones. Sin embargo, debido a que sólo se considera distancia, la binarización de los elementos no se ve afectada por sustituciones en residuos lejanos al lugar de ocurrencia, lo que denota la necesidad de idear metodologías que permitan contemplar el aporte completo de residuos a la caracterización de propiedades y cómo sustituciones puntuales afectan enormemente a residuos de interés. Una de las formas en las que se ha intentado dar solución a esta problemática, es modelar las propiedades fisicoquímicas de los residuos de las secuencias, a partir del uso de transformaciones de Fourier y en particular, empleando algoritmos relacionados a dichos conceptos, que aprovechen las ventajas referidas a la manipulación de espacios de frecuencias por sobre elementos temporales.

3.2 Transformaciones de Fourier

Las transformadas de Fourier, corresponden a una transformación matemática que permite analizar una función definida en el espacio tiempo, denominada señal, en sus frecuencias constituyentes. Como característica general, se genera una función definida en el espacio de frecuencias, representada por un valor complejo, en el cual, su módulo corresponde al valor de dicha frecuencia en la función inicial y su coeficiente, corresponde al desfase sinusoidal en la frecuencia [184].

Sea f una función definida en el espacio tiempo, representando una señal, integrable Lebesgue, su transformada se define como $f : \mathbb{R} \rightarrow \mathbb{C}$, asociada a una frecuencia, denotada por \hat{f} , la cual se expresa como:

$$\hat{f}(\xi) = \int_{-\infty}^{+\infty} f(x) e^{-2\pi i x \xi} dx \quad (3.1)$$

Donde ξ corresponde a un número real y x representa al tiempo. A partir de 3.1, se puede definir la transformación inversa

$$f(x) = \int_{-\infty}^{+\infty} \hat{f}(\xi) e^{2\pi i x \xi} d\xi \quad (3.2)$$

Tanto 3.1 y 3.2, corresponden a funciones con distribución continua. De manera similar, es posible definir las transformada de Fourier y su inversa, en el espacio de distribuciones discretas, en donde, sólo se considera un segmento muestral finito del conjunto de datos continuos para reconstruir el espectro de frecuencias [166]. Dado esto, la transformada de Fourier discreta se define como:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N} kn} \quad \forall k \in [0, N-1] \quad (3.3)$$

Donde N representa a una secuencia de números complejos x_0, \dots, x_{N-1} . Se define la transformada discreta inversa de Fourier como:

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k e^{\frac{2\pi i}{N} kn} \quad \forall n \in [0, N-1] \quad (3.4)$$

Las transformadas de Fourier han sido utilizadas en diferentes campos de investigación, tales como: física, teoría de números, procesamiento de señales, propagación de ondas, óptica, etc. Siendo el análisis armónico, la rama matemática encargada de este tipo de estudios.

A partir de lo anterior, debido a los diferentes empleos que puede tener esta transformada, nace la necesidad de resolver de manera eficiente esta función, para ello nacen diferentes algoritmos, dentro de los cuales, el principal se conoce como Transformada rápida de Fourier (FFT por sus siglas en inglés).

3.2.1 Transformada rápida de Fourier (FFT)

La transformada rápida de Fourier (FFT por sus siglas en inglés), es un algoritmo que permite encontrar solución a una DFT con una disminución en la complejidad. Esto es, al resolver el problema directamente desde la DFT, se presenta una complejidad $O(N^2)$, en cambio, al utilizar FFT, se obtiene una complejidad de $O(N \log N)$ [206].

La idea general del algoritmo fue propuesto por Cooley [54]. Dentro de sus particularidades, es que debido a la subdivisión en N transformadas de menor complejidad a resolver, donde N se compone de n_1 y n_2 , se requiere que el conjunto de muestras, presente un tamaño del orden $2 \cdot 2^n$, es decir, una potencia de 2. A pesar de que dicho algoritmo es uno de los más comunes para la resolución de transformadas de Fourier, no es el único, siendo algunos: Prime-factor FFT algorithm [116], Bruun's FFT algorithm, Rader's FFT algorithm, Bluestein's FFT algorithm, and Hexagonal Fast Fourier Transform [59].

Las definiciones matemáticas del proceso, se realizaron por Peter D. Welch en [206], explicando la formulación del problema y las demostraciones de la solución.

La división que se genera en el algoritmo FFT propuesto por Cooley [54], se basa en el uso de radix-2 DIT, esto es, la división de una DFT de tamaño N en dos DFT de tamaño $N/2$ de manera recursiva.

De manera general, se estiman las DFT de los pares e impares por separado (x_{2m} y x_{2m+1} , respectivamente), para luego combinarlas y estimar la DFT del espacio completo. Debido a esta subdivisión recursiva en pares, se requiere un número de componetes en potencia de 2. Normalmente, se adicionan elementos para poder cumplir con dicha condición, comúnmente, se utiliza *zero-padding* [144], para satisfacerlo.

Matemáticamente, las DFT de los componentes pares e impares y su combinación se obtiene a partir de:

$$X_k = \sum_{m=0}^{N/2-1} x_{2m} e^{-\frac{2\pi i}{N} 2mk} + \sum_{m=0}^{N/2-1} x_{2m+1} e^{-\frac{2\pi i}{N} (2m+1)k} \quad (3.5)$$

Donde el primer componente denota los elementos pares E_k y el segundo componente los impares O_k en la ecuación 3.5, respectivamente.

Un esquema representativo de los pasos a seguir en el algoritmo, la utilización del radix-2 DIT y cómo se obtienen las DFT para luego combinarlas se expone en la Figura 3.1.

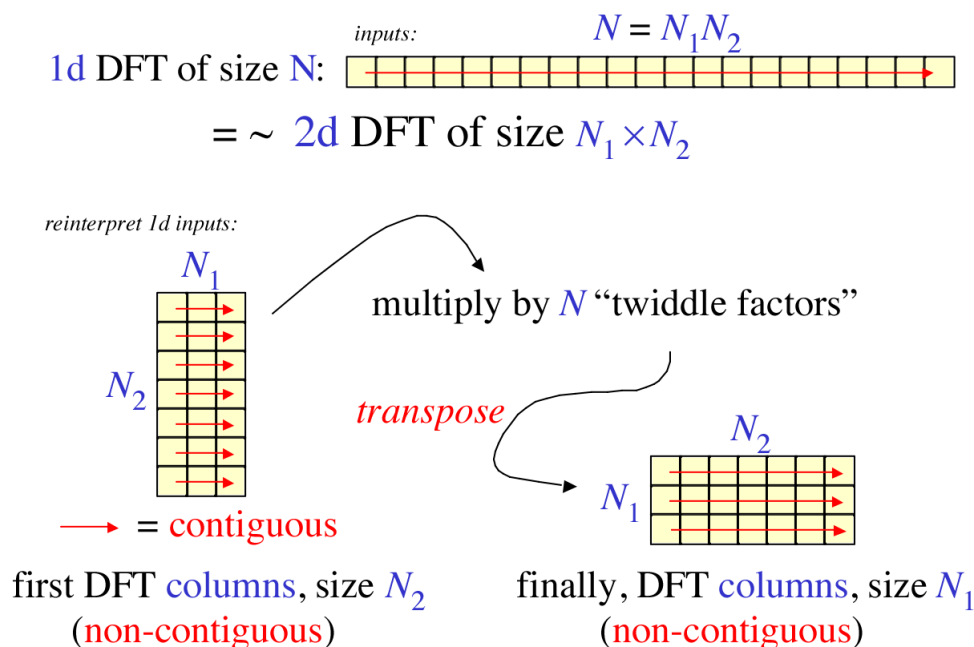


Fig. 3.1 Esquema representativo de los pasos asociados al algoritmo FFT, desarrollado por Cooley [54]

3.2.2 Uso de Transformadas de Fourier en digitalización de propiedades fisicoquímicas

Diferentes enfoques han sido evaluados mediante el empleo del uso de las transformadas de Fourier en el estudio de secuencias lineales de proteínas y DNA, en particular, para el estudio de propiedades fisicoquímicas y la identificación de residuos claves asociados a peaks visuales en el espectro de frecuencia [199, 58].

El uso frecuente de la digitalización, se centra principalmente en la codificación de una secuencia lineal de proteínas, en muchos casos, ha sido utilizado el potencial de interacción con electrones (PEII), ya que se han expuesto correlaciones con propensiones de moléculas orgánicas, relacionadas con toxicidad, actividad de antibióticos, carcinogenicidad, etc. [199, 56, 58] y a partir de dicha codificación, digitalizarla para el posterior estudio de la análisis de señales.

Dichos estudios de frecuencia de señales, son componentes de modelos asociados al reconocimiento de resonancias (RRM por sus siglas en inglés), los cuales han sido estudiados

en diferentes problemas del área médica, reconocimiento de Hot-spots, etc., y, en particular, enfocados en el estudio de mutaciones en proteínas [56–58].

Un enfoque común, asociado al estudio de señales y modelos RRM, consiste en, a partir de la secuencia lineal de residuos, emplear PEII para codificarla, a partir de su codificación, se aplica FFT como método de digitalización, para así, obtener el espectro de frecuencia y hacer los análisis de señales correspondientes [199, 56–58].

No obstante, existen otros enfoques, en donde se ha utilizado la digitalización de las propiedades fisicoquímicas como método de generación de atributos para set de datos, los cuales son sometidos a entrenamiento de predicciones, entregando medidas de desempeño satisfactorias, según los autores [35]. Sin embargo, la selección de las propiedades fisicoquímicas es arbitraria y sólo se considera la propiedad más informativa, no considerando el uso de combinaciones lineales o correlaciones entre las propiedades existentes.

Uno de los puntos relevantes, es que a partir de un conjunto de espectros de frecuencia, es posible generar un consenso y a su vez, identificar residuos relevantes al espectro y que contribuyen de manera significativa en los peaks, así como también la generación de espectros consenso a partir de secuencias diferentes, entregando uno de los principales planteamientos. *Un mismo espectro de frecuencia, está asociado a una propiedad o una funcionalidad*, lo cual fue testeado en trabajos como [199].

Esto último, lo hace una de las ventajas más relevantes para el estudio de mutaciones o reconocimiento de patrones. Sin embargo, la generalización del método ha impedido que se utilice con un mayor impacto, sólo siendo empleado en los estudios nombrados previamente. No obstante, las posibilidades de encontrar patrones o permitir modelar sistemas complejos, además de aplicar las propiedades del estudio de señales a dichos entornos de secuencias, y, las posibilidades que este método entrega para el análisis de estructuras lineales, lo hace una propuesta interesante y novedosa, que merece un esfuerzo estudiarla y aplicarla durante este trabajo de tesis.

3.3 Clustering

Clustering se define como un método de aprendizaje no supervisado, en el cual se cuenta con un conjunto de datos que representan a una muestra y en base a ésta, se trata de obtener grupos de objetos, denominados clusters [105].

Los clusters deben cumplir con dos características fundamentales:

- Los objetos que pertenezcan a un mismo clúster deben ser bastante homogéneos entre ellos.

- Entre los clústers debe existir un alto grado de heterogeneidad.

Los métodos de clustering tratan, fundamentalmente, de resolver el siguiente problema: Dado un conjunto de N individuos, caracterizados por la información de n variables X_j con j entre $1, \dots, n$, se plantea el reto de ser capaces de agruparlos de manera que los individuos pertenecientes a un grupo (cluster), dada la información disponible, sean tan similares entre sí como sea posible, siendo los distintos grupos entre ellos tan disimilares como sea posible.

Básicamente, el análisis constará de un algoritmo de clustering que permitirá la obtención de una o varias particiones, de acuerdo con los criterios establecidos.

3.3.1 Algoritmos de Clustering

Existen diversos algoritmos de clustering, cada uno con características que los diferencian, los cuales, pueden ser aplicados a diversos casos, dependiendo de las características de los datos de entrada, es decir, de la geometría de estos datos. Sin embargo, esta representación se basa principalmente en el uso de matrices, donde cada fila representa un ejemplo y cada columna el valor de un atributo o rasgo cualitativo para dicho ejemplo.

Los principales algoritmos de aprendizaje no supervisado se resumen en la Tabla 3.1 se expone un resumen de las características de cada algoritmo expuesto, la escalabilidad que poseen, las distancias que ocupan, los casos de uso y los parámetros que poseen.

3.3.2 Métodos de clustering empleando estructuras de grafos

Un grafo, es una estructura de datos compleja, que se compone de nodos y aristas, los nodos, son aquellos que representan la información y las aristas, permiten enlazar nodos. Existen diferentes tipos de representaciones, los cuales se basan en cómo fluye la información, encontrándose grafos dirigidos y no dirigidos. Los primeros, tienen una dirección entre un nodo A y B , la cual puede indicar diferentes comportamientos, por ejemplo, si los nodos representan genes, es posible mencionar que A regula a B . En el caso de grafos no dirigidos, no existe la relación expuesta previamente.

Matemáticamente, es posible definir las estructuras de gráficos, como un par de conjuntos $G = (V, E)$, donde V es el conjunto de vértices y E representa las aristas del grafo. En un grafo no dirigido, cada arista es un par no ordenado v, w . En un grafo dirigido, las aristas son pares ordenados. Los vértices v y w son llamados puntos finales de una arista. La cantidad de aristas $|E| = m$ es el tamaño del grafo. En un grafo ponderado, una función $\omega : E \rightarrow \mathbb{R}$ es definida la cual pondera cada arista.

Existen diversas utilidades que pueden ser expuestas usa para modelar problemas reales como redes de computadoras, redes sociales y estructuras de proteínas, entre otros.

Tabla resumen de Algoritmos de Aprendizaje No Supervisado				
Algoritmo	Parámetros	Escalabilidad	Usos	Métrica usada
K-Means	Número de clúster	Muchas muestras, mediana cantidad de clúster.	De propósito general, la geometría plana, no demasiados grupos	Distancia entre puntos
Affinity propagation	preferencia	No escalable con n ejemplos	Muchos clúster, tamaño de clúster desigual, geometría no plana	Distancia gráfica
Mean-shift	bandwidth	No escalable con n ejemplos	Muchos clúster, tamaño de clúster desigual, geometría no plana	Distancia entre puntos
Ward hierarchical clustering	Número de clúster	Mucha cantidad de ejemplos y de clusters	Cualquier clúster, es posible conexión de constraints	Distancia entre puntos
Agglomerative clustering	Número de clúster, tipo de unión, distancia	Mucha cantidad de ejemplos y de clusters	Muchos clusters, posiblemente restricciones de conectividad, distancias no euclidianas	Cualquier distancia pairwise
DBSCAN	tamaño vecino	Mucha cantidad de ejemplos, mediana cantidad de clúster	Geometría no plana, tamaños de clusters distintos	Distancia entre puntos vecinos
Gaussian mixtures	variado	No escalable	Geometría plana, bueno para la estimación de la densidad	Distancia Mahalanobis para los centros
Birch	branching, umbral	Alto número de clúster y ejemplos	Largo set de datos, eliminación valores atípicos, reducción de datos	distancia euclidiana entre puntos

Table 3.1 Cuadro resumen de algoritmos de aprendizaje supervisado

Otro de los puntos que pueden ser explotados en un grafo, es el hecho de generar clustering o identificación de comunidades, las cuales se definen como un conjunto de nodos, los que están fuertemente conectados entre sí y débilmente enlazados con otros elementos [75]. Los principales uso de las comunidades o sub grafos, corresponden a identificación de patrones con comportamientos particulares desde un punto de vista diferente a cómo lo hacen los métodos de aprendizaje no supervisado. Dentro de los diferentes algoritmos que permiten efectuar dicha búsqueda, se encuentran: Fast Greedy, Edge Betweenness, Cluster Walktrap, Leading eigenvector, Louvain, Infomap, dentro de los principales [146].

De manera análoga a las medidas de desempeño para la evaluación de los clustering empleadas en algoritmos de aprendizaje no supervisado, existe el concepto de modularidad [147], el cual permite evaluar la disgregación de los elementos y la densidad de las comunidades. Donde un valor de 0.35 es aceptado como comunidades bien definidas.

3.4 Hipótesis

Dada la problemática existente sobre cómo representar conjuntos de secuencias lineales con el fin de poder desarrollar modelos de clasificación/regresión o identificación de patrones asociados a propiedades fisicoquímicas. Y, en consideración de los diferentes usos que entrega las transformadas de Fourier, se plantea la hipótesis de este capítulo.

La codificación de secuencias lineales empleando espectros de Fourier, basados en las propiedades fisicoquímicas de los residuos de la proteína, permite generar descriptores que facilitan el aprendizaje de predictores de variantes enfocados a diferentes respuestas de interés. A su vez, es posible correlacionar propiedades fisicoquímicas o funciones con propiedades de los espectros de frecuencia

Si bien, en el planteamiento de la hipótesis se exponen dos preguntas, la interrogante en sí, se centra a los posibles usos que pueda tener los espectros de frecuencia en el estudio de variantes, identificación de patrones, residuos relevantes, etc.

3.5 Objetivos

En base a la hipótesis planteada y con el fin de responder a los planteamientos e interrogantes expuestas. Se detallan el objetivo general y los objetivos específicos.

3.5.1 Objetivo general

Diseñar e implementar metodología de codificación y digitalización de propiedades fisicoquímicas en secuencias lineales de proteínas, empleando transformadas rápidas de Fourier, con el fin de poder ser utilizadas en identificación de patrones por medio de técnicas de clustering o desarrollo de predictores basados en algoritmos de aprendizaje supervisado.

3.5.2 Objetivos específicos

A partir del objetivo general, nacen los siguientes objetivos específicos.

1. Preparar, manipular e implementar módulos de consultas para la base de datos de propiedades fisicoquímicas asociadas a la base de datos AAindex [111].
2. Diseñar e implementar, metodologías de codificación de propiedades fisicoquímicas y selección de las más representativas, por medio de técnicas de reducción de dimensionalidad y selección de características, descritas a través de espectros de frecuencia.
3. Implementar y validar modelos de clasificación/regresión para evaluación de análisis de estabilidad de variantes según descriptores basados en espectros de frecuencia de propiedades fisicoquímicas.
4. Diseñar, implementar y caracterizar grupos obtenidos a partir de algoritmos de aprendizaje no supervisados, considerando como descriptores los espectros de frecuencias asociados a las propiedades fisicoquímicas.

3.6 Metodología

Con el fin de poder cumplir con el objetivo general planteado y los objetivos específicos, se expone a continuación la metodología diseñada. Se consideran diferentes etapas dentro de las cuales se destaca la codificación, entrenamiento de modelos, aplicación de clustering e identificación de residuos como patrones de señales dentro del espectro de frecuencia.

A continuación se exponen las diferentes etapas asociadas al proceso.

3.6.1 Identificación y selección de propiedades fisicoquímicas

El uso de la base de datos AAIndex, contempla trabajar con 566 propiedades que describen a un residuo, dentro de las cuales, existen altas correlaciones entre diferentes propiedades, ya que, evalúan el mismo concepto. Pero, empleando técnicas o metodologías diferentes. Razón

por la cual, es necesario evaluar cuáles son las propiedades que pueden ser utilizadas, con el fin de no generar redundancia de información.

A partir de lo anterior, se emplean estructuras de grafos no dirigidos, en los cuales, los nodos representan la propiedad y las aristas corresponden a si existe una alta correlación entre la propiedad A y B . Se destaca que dicha correlación se mide con respecto al coeficiente de Pearson y sólo se considerará altamente correlacionados si dicho índice es sobre 0.9. Una vez creado el nodo, se aplicarán diferentes algoritmos de identificación de comunidades y se evaluará la modularidad de estos, a través de la cual, se seleccionarán las con mayor índice.

Esto último, permite generar grupos de propiedades cuya característica principal es que presentan alta correlación entre ellos, dado esto, si se tienen n comunidades, es posible representar una secuencia lineal a partir de cualquier propiedad seleccionada de cada comunidad n_i . De esta forma, se reduce la dimensión del espacio muestral.

3.6.2 Codificación de secuencias lineales

La codificación de secuencias lineales, se basa en el uso de propiedades fisicoquímicas representativas de la secuencia, las cuales se obtienen a partir de la base de datos AAIindex[111] y son seleccionadas de manera aleatoria con respecto a las comunidades identificadas a través de la correlación de dichas características desde las estructuras de grafos generadas.

Un esquema representativo del proceso, se observa en la Figura 3.2, en la cual, se detallan los diferentes pasos a seguir para generar la codificación correspondiente y obtener los espectros de frecuencias asociados a cada propiedad fisicoquímica.

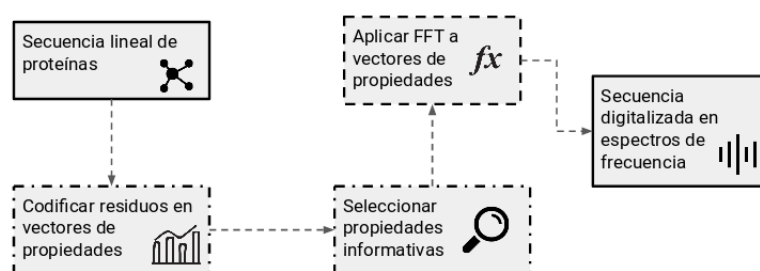


Fig. 3.2 Esquema representativo, metodología de digitalización de secuencias.

En una primera instancia, se toma la secuencia y por cada residuo se crea un vector de tamaño n el cual representa el número de propiedades fisicoquímicas a considerar a partir de las comunidades generadas. De esta forma, se crea una matriz de tamaño $r \times n$ donde r representa la cantidad de residuos en la secuencia.

A partir de dicha matriz, técnicas de reducción de dimensionalidad y selección de características son implementadas, utilizando lenguaje de programación Python y la librería

scikit-learn [157], con el fin de seleccionar cuáles son las propiedades más representativas y qué porcentaje de la varianza permiten explicar.

Dado el conjunto de propiedades seleccionadas, se implementarán rutinas basadas en lenguaje de programación Matlab, las cuales reciben el conjunto inicial de datos, en una primera instancia, aplica *zero-padding* con el fin generar vectores de tamaño de potencia de $2 \cdot 2^n$, requisito para la aplicación de FFT. A partir de esto, cada columna en el conjunto de elementos, se digitaliza por medio del uso de la transformada rápida de Fourier (FFT) y se obtienen los espectros de frecuencias para cada propiedad fisicoquímica seleccionada previamente.

De esta forma, por cada secuencia, se obtiene un conjunto de espectros de frecuencia, asociados a la digitalización de las propiedades fisicoquímicas seleccionadas mediante técnicas de selección de características.

3.6.3 Caracterización del espectro de frecuencias

Dado el conjunto de espectros de frecuencia generado para las secuencias lineales, es necesario extraer información relevante de estos con el fin de caracterizarlo, para ello, se plantean diferentes metodologías con el fin de poder generar los conjuntos de entrenamientos, y, a partir de estos, implementar los modelos predictivos propuestos. Es importante mencionar, que dichos patrones también pueden ser utilizados como método de reconocimiento de residuos relevantes en las variantes, ya que, si se aplica la transformada inversa de Fourier, es posible la identificación de los elementos que participan en algún peak del espectro [199].

Tomando esto en consideración, se plantean diversas metodologías que serán implementadas para llevar a cabo la caracterización del espectro. Un punto relevante, es el hecho de que al ser una etapa exploratoria, será necesario evaluar cuál es la mejor, a la hora de aplicarla en entrenamiento de modelos predictivos. Cada una de las cuales, se explica a continuación.

Diferencias espectrales

Sea $f_i(wild)$ el espectro de frecuencia de la proteína no mutada para la propiedad i y $f_i(mut_j)$ el espectro de frecuencia para la mutante j . Se define el espectro diferencial como la resta de $f_i(wild)$ y $f_i(mut_j)$.

Ya que se tendrán j mutantes, existirán j diferencias espectrales, a partir de las cuales, se emplearán para formar conjuntos de datos y entrenar modelos predictivos.

Uno de los puntos importantes a destacar, radica en el hecho de es posible identificar variaciones positivas y negativas, lo cual puede influir de diferentes formas en la propiedad y por consecuencia, en la función o respuesta asociada a evaluar.

Binarización de diferencias espectrales

A partir de las diferencias espectrales expuestas en el punto anterior, es posible desarrollar distribuciones de las diferencias de los espectros, de tal manera, que se apliquen test estadísticos para la identificación de outliers en la muestra, con un nivel de significancia α .

Dado esta identificación, se considera el conjunto de elementos y se codifica con respecto a si es un outlier positivo, negativo o se encuentra dentro del rango normal de distribución. De esta forma, por cada espectro, se forma un vector binarizado, dando origen a una matriz de espectros codificados en base a la significancia de su diferencia.

Dicha matriz se utiliza para el entrenamiento de modelos predictivos, a su vez, esta matriz podría ser representada en un heat map, con el fin de poder identificar zonas en las que exista una mayor densidad de puntos con tendencia hacia un tipo de outlier.

Identificación de zonas en el espacio espectral

De manera visual, el espectro de frecuencia en sí, es informativo, ya que, permite identificar zonas en las que existe un peak, las cuales reflejarían una característica correlacionada con la función [200]. A su vez, si se considera la diferencia espectral, ésta, permitirá identificar cómo fue la variación con respecto a la proteína inicial, de tal manera que, una distribución uniforme con valor 0 debiese pertenecer a variantes que no afecten los valores de la propiedad en sí, mientras que, variaciones importantes, serán detectadas en zonas específicas de la distribución. A su vez, al binarizar la distribución de diferencias, es posible visualizar estas tendencias.

Con el fin de reconocer de manera automática dichas zonas, se trabaja con las distribuciones discretas de los espectros binarizados, obtenidos previamente, para las cuales, se diseñará e implementará un algoritmo de detección de cambios en la distribución, tal que, permita ir marcando por zonas en donde se exhiba un cambio positivo, negativo o neutro.

A su vez, sería interesante de estudiar, o, identificar las variantes generadas y sus correspondientes efectos y ver cómo se clasifican en estas zonas, ya que, podría ser identificable, zonas en el espectro, que conllevan una alteración negativa para una proteína, tal que, se podría identificar, cuáles son los residuos o posiciones que más afectan a un cambio negativo en la proteína.

Con los diferentes puntos expuestos previamente, se generan distintas caracterizaciones, que permiten describir el espectro y generar el conjunto de datos, el cual será utilizada para el entrenamiento de modelos predictivos.

3.6.4 Implementación de modelos de clasificación/regresión para análisis de variantes

Uno de los objetivos de la codificación de secuencias lineales, es evaluar si el conjunto de espectros de frecuencia para un grupo de variantes, puede ser utilizado como características para generar set de datos y entrenar modelos a partir de estos.

Con esto en mente y apoyados en los conjuntos de datos utilizados para la generación de descriptores basados en propiedades termodinámicas y filogenéticas, expuestos en el capítulo 2, se desarrollarán modelos de clasificación para la evaluación de la estabilidad de proteína y modelos de regresión para la predicción del cambio en la energía libre provocado por los residuos.

En la Figura 3.3, se expone un esquema representativo asociado a los pasos a seguir para el desarrollo de los modelos.

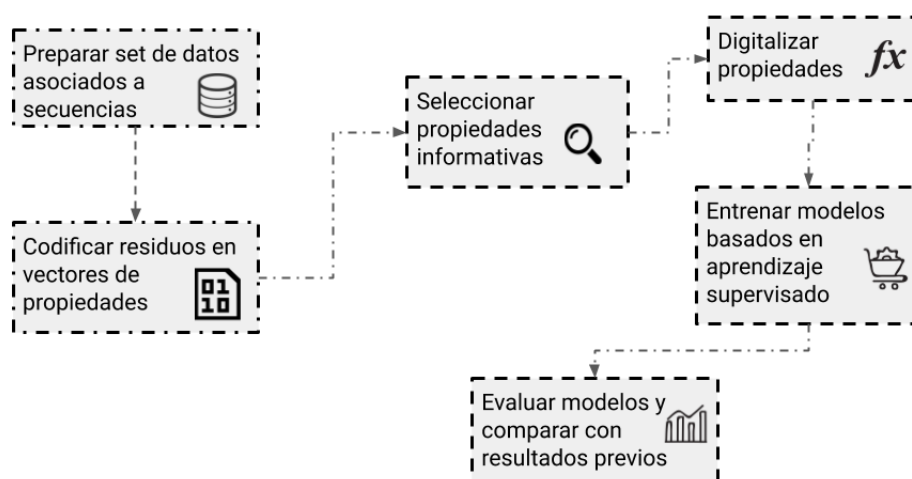


Fig. 3.3 Esquema representativo, metodología de clustering de secuencias por medio de espectros de frecuencias basados en propiedades fisicoquímicas.

En una primera instancia, se necesita preparar el conjunto de datos, para ello, se considerarán la secuencia original de la proteína y serán generadas las variantes con respecto a la mutación reportada. De esta forma se creará un conjunto de datos basados en una variante y la respuesta asociada, lo cual corresponde a las diferencias de energía libre que provoca la sustitución del residuo o clasificación de estabilidad.

Al conjunto de secuencias generado, se aplicará la codificación de las propiedades fisicoquímicas, a partir de la información existente en la base de datos AAindex.

Posterior a ello, se seleccionarán las propiedades más informativas y se generará la digitalización de éstas, empleando la metodología descrita en el punto anterior. La selección de las características se basa en un consenso con respecto a las incidencias de cada propiedad en cada secuencia, esto es, dado a que se seleccionan p propiedades por cada secuencia, es posible que diferentes secuencias, presente distintas propiedades. Por lo tanto, se seleccionarán aquellas propiedades que presenten mayor incidencia en el conjunto de secuencias, ya que éstas, serán las más representativas del conjunto completo.

Una vez se tenga el conjunto de espectros caracterizado, modelos predictivos serán entrenados aplicando algoritmos de aprendizaje supervisado al set de datos de espectros de frecuencia. Se utilizarán las medidas de desempeño expuestas en el capítulo 2. Los modelos serán validados mediante validación cruzada con un valor de $k = 10$, con el fin de evaluar el sobreajuste.

Finalmente, los resultados a obtener a partir de descriptores basados en espectros de frecuencia, serán comparados con los obtenidos en la fase de exploración de la metodología expuesta en el capítulo 2. Esto con el fin de determinar, qué metodología o caracterización de datos, permite entregar un modelo con mejor desempeño o características deseables.

Es importante mencionar que, los modelos que se obtengan a partir de la digitalización de propiedades fisicoquímicas, pueden presentar performance inferior a los modelos a obtener aplicando la metodología del capítulo 2. Sin embargo, si el desempeño es alto, sería suficiente para responder la pregunta planteada. Si son bajos o azarosos, implica que se requiere un mayor refinamiento a la metodología, o, que simplemente el conjunto de datos presenta problemas para entrenar modelos, razón por la cual, debiese descartarse.

3.6.5 Aplicación de técnicas de clustering sobre espectros de frecuencia

Uno de los supuestos más relevantes asociados al uso de la digitalización y las propiedades fisicoquímicas como descriptores de secuencias lineales, es el hecho de que, proteínas con una misma función, presentan un espectro de frecuencia similar [199].

Esto, hace pensar que, para un conjunto de secuencias desconocidas, dada una selección de propiedades fisicoquímicas, a la hora de aplicar técnicas de clustering, empleando algoritmos de aprendizaje no supervisado, serán agrupadas de tal forma, que, los espectros de frecuencia en cada grupo presenten similitudes entre ellos, basados en propiedades estadísticas o por medio de análisis cross-espectral.

Con el fin de corroborar este supuesto, además de demostrar que el uso de descriptores basados en espectros permite una agrupación de elementos correlacionados con alguna propiedad, en la Figura 3.4 se expone la metodología planteada.

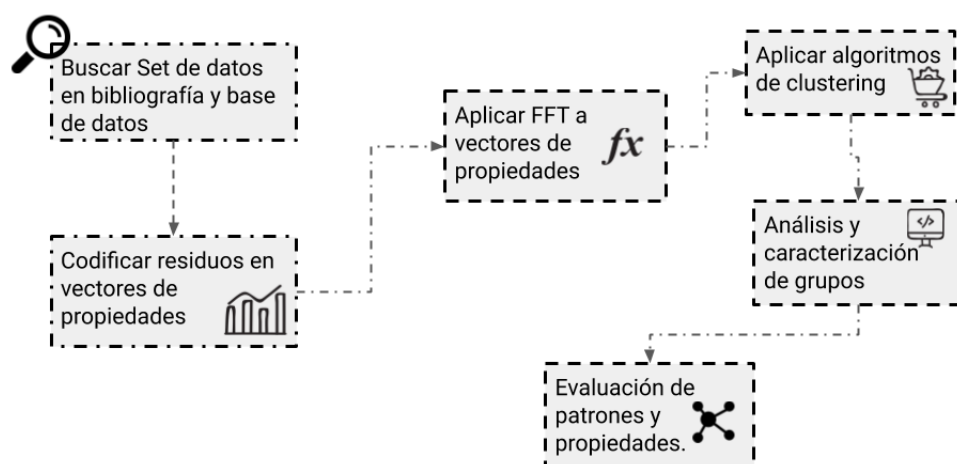


Fig. 3.4 Esquema representativo, metodología de clustering de secuencias por medio de espectros de frecuencias basados en propiedades fisicoquímicas.

En una primera instancia, se deberá identificar cuáles son las secuencias de interés, para ello, se trabajarán con diferentes secuencias lineales y variantes asociadas a ellas, principalmente de enzimas, con alguna propiedad característica, proteínas comunes de interés con mutaciones reportadas, etc. La descarga de éstas, se realizará desde diferentes bases de datos, en particular desde Brenda [180], para el caso de las enzimas, Protherm [16] para otras propiedades de interés, relacionadas a estabilidad.

Es importante mencionar, que no se quiere asociar o reconocer qué efecto causa la mutación. Si no, probar que, proteínas con similar función, tendrán un espectro de frecuencia con características comunes. Es por ello, que la propiedad se conoce de antemano. Sin embargo, también se evaluará si es posible la generación de grupos donde la mutación afecte negativa y positivamente a la estabilidad en proteínas, utilizando para esto, el conjunto de set de datos expuestos en el capítulo 2.

Una vez se tengan las secuencias, se aplicará la codificación de diferentes propiedades fisicoquímicas y a su vez, se digitalizará cada una de éstas. Es importante mencionar, que no se seleccionarán propiedades específicas ya que, puede que las más representativas o informativas no necesariamente generen conjuntos de datos agrupados con el mismo patrón de espectro de frecuencia.

Esto implica, que se generarán n espectros de frecuencia por cada secuencia y para cada propiedad se aplicarán técnicas de clustering, esto con el fin de determinar si aparecen grupos con misma propiedad conocida, correspondientes a las n propiedades descritas en AAindex.

Lo anterior, permitirá evaluar si proteínas con misma funcionalidad, o mismo plegamiento, o, alguna característica en común, quedan agrupadas.

Finalmente, una vez se tengan estos grupos, se realizará la etapa de análisis, caracterización e identificación de patrones. Esto es, en base a la información conocida que se tiene del conjunto de proteínas, evaluar si los grupos formados presentan algo en común, ya sea, propiedades fisicoquímicas, plegamiento, actividad, funcionalidad, etc.

Esto permitirá evaluar si los espectros de frecuencia, permiten agrupar secuencias con características o patrones en común y si son informativas para dicho objetivo. Se destaca además, que esto implica una fase exploratoria y es incierto lo que pueda resultar. Si bien se postuló las correlaciones entre espectro y propiedad [199], esto fue hecho hace un tiempo prolongado y sólo se evaluaron proteínas o enzimas con particularidades específicas, de esta forma, se espera poder reafirmar dicha afirmación y poder postular un nuevo método de clustering de propiedades en secuencias y sus variantes.

Chapter 4

Filogenética, propiedades fisicoquímicas y minería de datos aplicados al diseño de mutaciones en secuencias de proteínas

Uno de los problemas de mayor relevancia para el campo de la ingeniería de proteínas, es el diseño inteligente de mutaciones, con el fin de obtener una propiedad fisicoquímica, mejorar las características o adicionar una nueva funcionalidad. Sin el hecho de incurrir en grandes costos económicos, de recursos humanos y computacionales.

En la sección 1.1, se mencionó que existen dos puntos de vista a la hora del diseño de mutaciones: Evolución dirigida y el diseño racional de proteínas. Sin embargo, ambas presentan el problema del poco espacio de exploración que pueden abarcar. A partir de ello, y con el fin de disminuir los tiempos experimentales y aumentar el espacio de búsqueda, métodos computacionales fueron implementados como herramientas de apoyo al diseño de mutantes.

A pesar de la existencia de dichos métodos, para aquellos basados en potenciales de energía y en simulaciones Monte Carlo, el tiempo computacional necesario para explorar la totalidad de mutaciones en una proteína en particular, escala a nivel cuadrático, siendo elevado y de limitado acceso para usuarios en general.

Una alternativa a los métodos computacionales basados en estrategias de potenciales de energía, son aquellos que emplean técnicas de minerías de datos y algoritmos de aprendizaje. Sin embargo, estos, principalmente se enfocan en el estudio de la estabilidad de la proteína ante sustituciones de aminoácidos. Por otro lado, los principales enfoques en cuanto a caracterización de residuos, se centran en propiedades termodinámicas y el ambiente, no siendo considerados conceptos filogenéticos ni propensión a cambios.

Tal como se expuso en el capítulo 3, la codificación de secuencias lineales puede ser realizada a través del uso de propiedades fisicoquímicas y posterior digitalización de éstas, aplicando transformadas de Fourier, esto último, permite caracterizar el ambiente y el aporte hacia las características que brindan los residuos, debido a las propiedades del espacio de frecuencias y el espectro como tal.

Apoyados en dicha codificación y centrados en la utilización de algoritmos de aprendizaje supervisado para el entrenamiento de modelos de clasificación o regresión, y utilizando la metodología expuesta a lo largo del capítulo 2, además del uso de herramientas computacionales como filtros de diseño. Se propone durante este capítulo el diseño e implementación de una herramienta computacional, basada en técnicas de minería de datos y aprendizaje de máquinas, que permita el diseño de mutaciones en variantes con características deseadas.

4.1 Hipótesis

En base al planteamiento del problema y a la motivación existente por el desarrollo de una herramienta computacional para el diseño de mutaciones, se plantea la siguiente hipótesis.

La digitalización de las propiedades fisicoquímicas por medio de transformadas de Fourier, en conjunto con algoritmos de aprendizaje supervisado, en combinación con herramientas computacionales para evaluar estabilidad y propensión, serán un enfoque suficiente para el desarrollo de una herramienta de diseño de mutaciones

4.2 Objetivos

Dada la hipótesis planteada y en vista del desafío considerado, se exponen a continuación el objetivo general y los objetivos específicos.

4.2.1 Objetivo general

Diseñar, implementar, testear y depurar herramienta computacional para el diseño de mutaciones puntuales en variantes de proteínas, enfocada en el uso de minería de datos y aprendizaje supervisado y en la generación de estrategias de filtro de mutantes con respecto a características termodinámicas y filogenéticas.

4.2.2 Objetivos específicos

Del objetivo general, nacen los siguientes objetivos específicos.

1. Implementar, evaluar y analizar, sistema de codificación de secuencias lineales por medio de propiedades fisicoquímicas y el uso de digitalización a partir de transformadas de Fourier.
2. Entrenar, evaluar y validar modelos basados en algoritmos de aprendizaje supervisado, con descriptores enfocados en espectros de frecuencias de propiedades fisicoquímicas.
3. Diseñar e implementar módulo de filtro de mutaciones por medio de propensión de la mutación y efecto en la estabilidad que provoque la sustitución.
4. Implementar y evaluar flujo de trabajo para nuevas mutaciones, mostrando el desempeño del modelo y los valores asociados a ésta, junto con otras características de interés.

4.3 Metodología propuesta

Dada la hipótesis planteada y los objetivos propuestos, se diseña una metodología que contemple tanto la generación de la herramienta computacional, como el desarrollo de los modelos y su entrenamiento.

Un esquema representativo de los componentes principales de la metodología y cómo estos interactúan se expone en la Figura 4.1. El proceso general, se puede dividir en dos etapas: Generación de los modelos y Evaluación de nuevas mutaciones. En la primera etapa, se contempla un conjunto de datos, el cual es sometido a la etapa de digitalización, posterior a ello, se entrenan los modelos considerando como descriptores los espectros de frecuencia, para luego evaluar su desempeño. En la segunda etapa, nuevas mutaciones son propuestas y evaluadas en una etapa de filtro, aplicando criterios de estabilidad y propensión filogenética, para luego, someterse a los modelos predictivos generados y así evaluar si la mutación tendrá el efecto deseado o no.

Un punto importante a evaluar, consiste en el hecho que el conjunto inicial de datos debe ser lo suficientemente representativo en cuanto a la diversidad de ejemplos, es decir, si se evalúa la presencia o ausencia de una característica, no debe existir un desbalance entre ellas, ya que provocaría un sobreajuste en el modelo, tendiendo a predecir sobre la característica de mayor proporción.



Fig. 4.1 Esquema representativo de la metodología propuesta para el diseño de mutaciones aplicando herramienta computacional a desarrollar

A continuación, se describe el conjunto de datos inicial, seguido de las etapas que componen la metodología y finalizando con la estrategia para el diseño e implementación de la herramienta computacional.

4.3.1 Conjunto de datos

El conjunto de datos corresponde a un grupo de variantes de una misma proteína, ordenados en un archivo en formato *.fasta¹, siendo la primera secuencia la proteína original, y el resto de secuencias las variantes con mutaciones reportadas. Además del conjunto de secuencias, es necesario un archivo en formato *.csv con la variable respuesta asociada a la variante, es decir, el valor de la propiedad fisicoquímica, característica, funcionalidad, etc., que provoca la mutación. Esto con el fin, de poder asignárselo a la secuencia y entrenar los correspondientes modelos.

Adicional a ambos inputs, se requiere de un archivo *.pdb el cual contenga la estructura de la proteína original, o en su defecto, el código PDB de la misma, esto debido a que la etapa de filtro, utiliza SDM como método de análisis de estabilidad de la proteína ante sustituciones de residuos, siendo un input para ésta, la estructura 3D.

¹Formato de texto plano para la representación de secuencias.

4.3.2 Digitalización de secuencias lineales

A partir del conjunto de secuencias, los residuos son codificados aplicando las propiedades fisicoquímicas descritas en la base de datos AAindex [111], considerando la totalidad de propiedades descritas. Para ello, se implementarán scripts bajo el lenguaje de programación Python, los cuales tomen los residuos de cada secuencia y los transformen a un vector de tamaño n el cual corresponde a la cantidad de propiedades descritas en la base de datos. De esta manera, por cada secuencia se crea una matriz de tamaño $r \times n$ donde r representa la cantidad de residuos en la secuencia.

Dada esta matriz $r \times n$, se aplican técnicas de reducción de dimensionalidad y análisis de características, con el fin de seleccionar las propiedades fisicoquímicas más informativas y que generarían un aporte al entrenamiento de modelos. Debido a que son s secuencias existentes en el set de datos, se tendrán m propiedades fisicoquímicas por cada secuencia en el conjunto de elementos. La selección se basará en las propiedades consenso que abarquen la totalidad de las secuencias.

Es importante mencionar, que se espera que las propiedades informativas se mantengan a lo largo de la secuencia original y las variantes, ya que, cambios puntuales en los residuos, no alterarán de manera significativa la varianza que generen al conjunto de datos, alterando el orden de prioridad de las propiedades asociadas.

Una vez seleccionadas las propiedades fisicoquímicas a utilizar, se digitalizará su valor, con el fin de formar espectros de frecuencia asociados a dicho elemento; así, debido a que la selección previa, permite determinar un número p de propiedades informativas, por cada secuencia s en el conjunto de datos, se tendrán p espectros de frecuencias. Estos espectros se obtienen a partir del uso de Transformadas rápidas de Fourier (FFT por sus siglas en inglés) [31], para lograr dicha conversión, se implementarán scripts bajo lenguaje de programación Matlab, los cuales, reciban como entrada el conjunto de propiedades por cada secuencia y retornen los espectros de frecuencia por cada propiedad.

Finalmente, scripts basados en lenguaje de programación Python, completarán el conjunto de datos, considerando los espectros de frecuencia, seguidos de la respuesta que estos conllevan. De esta forma, generando el set de datos para el entrenamiento de modelos.

4.3.3 Entrenamiento de modelos

Los modelos se basarán en algoritmos de aprendizaje supervisado y se utilizará la misma metodología de exploración y selección de modelos basados en sus medidas de desempeño, expuestos en el capítulo 2. Se destaca que, el tipo de algoritmo a utilizar, depende de las características de la respuesta que se esté evaluando, es decir, si la respuesta presenta una

distribución continua, se utilizarán los métodos basados en regresión, en caso contrario, se utilizarán algoritmos basados en clasificación.

Con respecto al desbalance de clases, éste se evaluará de la misma forma expuesto en el capítulo 2, además, las respuestas del tipo continua serán analizadas mediante la evaluación de outliers y la tendencia a distribución normal que presente.

El desempeño de los modelos, será obtenido a partir de lo expuesto en el capítulo 2. Uno de los puntos importantes a considerar es el valor de estas medidas. Modelos de clasificación con medidas inferiores serán descartados y requerirán un nivel de análisis más detallado. Por otro lado, modelos de regresión con coeficientes de correlación inferiores a 0.6, también implica que deben ser analizados de la misma forma. De esta forma, se asegura un cierto nivel de confianza a la hora de aplicar los modelos. No obstante, dado a que a la metodología planteada en el capítulo 2 mejora las medidas de desempeño de los modelos iniciales, se espera que los casos expuestos de descarte no ocurran.

4.3.4 Diseño de mutaciones

Una vez los modelos se encuentren entrenados, será posible utilizarlos para evaluar nuevas mutaciones. Para ello, y con el fin de testear la herramienta, se implementará un script que permita mutar todas las posibilidades en cada residuo de la secuencia, es decir, para el residuo i se sustituirá 19 veces, en caso de que la mutación ya se encuentre reportada, ésta será descartada.

Una vez se tenga este conjunto de mutaciones, se aplicará un filtro basado en la propensión filogenética de dicha mutación y cómo afecta a la estabilidad el cambio propuesto. Para ello, se implementará servicios API (Application Programming Interface) que consuman las herramientas SDM [154] y MOSST [151], los cuales permitirán mencionar si la mutación es viable filogenéticamente y no provoca cambios en la estabilidad. Esto, generará una reducción de elementos en el conjunto de datos a estudiar, debido a que, no todas las mutaciones serán factibles.

Ya con las mutaciones factibles, se codificará las secuencias de las variantes utilizando el método descrito previamente y aplicando las propiedades fisicoquímicas seleccionadas, con el fin de obtener los espectros correspondientes.

Una vez obtenidos los espectros, se someten a los modelos entrenados y se obtiene la respuesta de interés, en base a las características del modelo. Se destaca que las respuestas serán reportadas en torno a intervalos de confianza, para métodos continuos, y, en forma de probabilidades, para el caso de respuestas categóricas.

4.3.5 Implementación herramienta computacional

La herramienta computacional será diseñada siguiendo el patrón de diseño Modelo-Vista-Controlador (MVC) [118] e implementada utilizando el paradigma de Programación Orientada a Objetos (POO) [205], además, debido a que el componente de la vista, será asociado a interfaz web, se utilizará la arquitectura Cliente-Servidor, con el fin de responder las solicitudes y ejecutar las acciones correspondientes. El uso de POO es debido a que permite una mayor estandarización de los módulos y facilita su re usabilidad, además de la comprensión y simpleza a la hora de programar dado a su cercanía con la realidad y al poder de abstracción que posee.

Con respecto al patrón de diseño, se expone a continuación, cuáles serán los principales elementos en cada componente, además de sus características y qué comprenden cada uno de estos.

Modelo

El modelo corresponde al conjunto de scripts y módulos que contienen toda la lógica de la herramienta y las funcionalidades principales, se hace alusión a que forma parte del "back-end" de la aplicación.

En este caso, se compondrá de los módulos de codificación, procesamiento de datos, entrenamiento de modelos, diseño de mutaciones y uso de servicios para la ejecución de MOSST y SDM, así como también los módulos de gestor de usuarios, notificaciones vía email, etc. Será implementado bajo el lenguaje de programación Python y utilizará algunas rutinas desarrolladas en Matlab. A su vez, cada módulo y sus componentes serán diseñados bajo el paradigma de Programación Orientada a Objetos y se utilizarán librerías externas como Pandas [139] para la manipulación del conjunto de datos, Numpy [197] y Scipy [150] para los análisis estadísticos y Scikit-Learn [157] para el entrenamiento de modelos.

Controlador

El controlador, hace referencia al gestor de las solicitudes de usuario desde la vista y recibe las respuestas de dichas solicitudes por parte del modelo, con el fin de ser expuestas al usuario.

Para este caso, el controlador se dividirá en dos componentes principales: Controlador de acciones en la vista y Controlador de respuestas en el modelo. El primero, será implementado en JavaScript y tendrá funcionalidades asociadas a jQuery, mientras que el segundo, será desarrollado bajo el lenguaje de programación Php.

La gran diferencia entre ambos, radica en el hecho de dónde se ejecutan, el primero, se encuentra principalmente condicionado por las acciones del usuario y su ejecución es en el navegador. Mientras que el segundo es una consecuencia del primero, es decir, una vez que se reciben las solicitudes, se establece la comunicación hacia el servidor por medio de Ajax (Asynchronous JavaScript And XML) y esto permite la ejecución del segundo, el cual, entrega una respuesta en formato JSON (JavaScript Object Notation), la que es capturada por el primero, y mostrada al usuario.

Vista

La vista, es el componente de visualización de la herramienta, es decir, es el componente en el cual el usuario puede interactuar, levantar solicitudes y exponer los resultados o respuestas que entregue el controlador y es conocido como "front-end".

Con el fin de poder representar las secciones principales que tendrá la herramienta, a continuación se exponen un conjunto de mockups (maqueta de diseño), en los cuales se exponen cuáles serán las principales funcionalidades y cómo se verán los resultados.

Generador de jobs

Los jobs, son los eventos generados asociados a un usuario y la carga de un conjunto de datos, los cuales deben ser codificados y entrenados los modelos, con el fin de evaluar las mutaciones posibles. Un esquema representativo del formulario asociado a la generación de jobs se aprecia en la Figura 4.2.

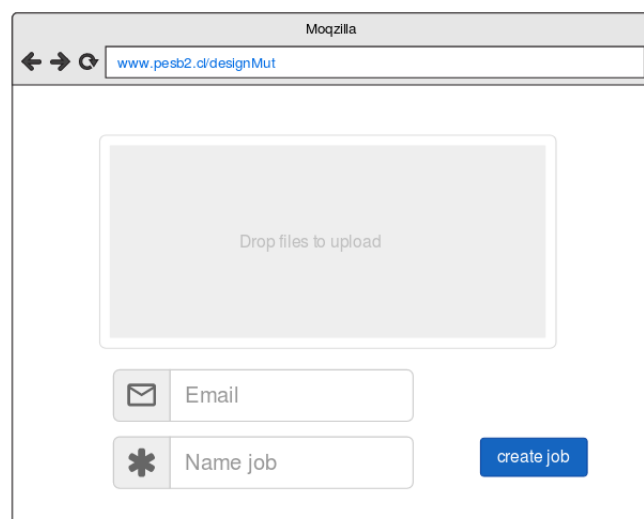


Fig. 4.2 Esquema representativo interfaz de creación de jobs.

En una primera instancia, el usuario deberá subir los archivos necesarios para entrenar los modelos, estos corresponden al conjunto de secuencias, el archivo de respuestas y la estructura 3D en formato *.pdb. Una vez suba los archivos, deberá ingresar el correo electrónico y el nombre del trabajo, esto para que él pueda identificarlo en el sistema. Cuando el usuario pulse el botón "create job", el sistema procesa la data y genera un ID asociado al identificador único del Job, posterior a ello, lo agrega al sistema de colas en el cual será procesado una vez le corresponda. Visualmente, el sistema retorna dicho ID, mostrándose en la interfaz y se notifica vía correo electrónico el estado del job, así como las notificaciones correspondientes del mismo.

Buscador de jobs

Adicional a la sección de crear jobs, la herramienta computacional permite buscar los trabajos generados, estos pueden encontrarse en 4 estados:

- **Creado:** se refiere al estado en el que el usuario sube la data y se ancla al sistema de colas.
- **En ejecución:** el job está siendo ejecutado por la herramienta computacional.
- **Finalizado:** el job fue ejecutado y es posible ver los resultados obtenidos.
- **Cancelado:** el job fue cancelado por el usuario.

Los jobs, pueden ser buscados en la interfaz del buscador de la herramienta web, mediante el ID o el nombre, ambos, son notificados vía correo electrónico, ya sea ante cambios de estado que éste sufra, o, al momento de crearse. Finalmente, una representación general del buscador, se observa en la Figura 4.3.

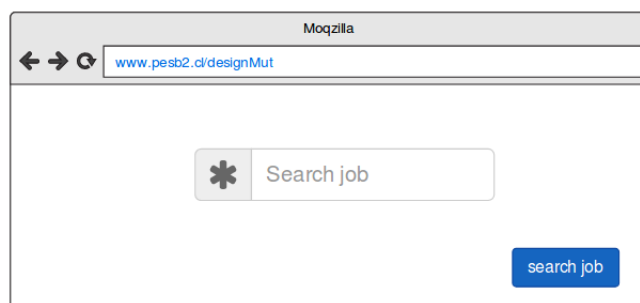


Fig. 4.3 Esquema representativo interfaz de búsqueda de jobs.

Visualizador de resultados

Una de las visualizaciones más relevantes, corresponde a los resultados involucrados en todo el proceso, estos, abarcan las propiedades fisicoquímicas estudiadas, los espectros de frecuencia y los residuos relevantes, los entrenamientos de los modelos y el listado de mutaciones favorables a evaluar. Un resumen general de los resultados esperados y la muestra de estos, es listada a continuación.

- **Descripción general:** Se genera una visualización del conjunto de secuencias y la distribución o frecuencia de la variable respuesta, se expone un resumen del conjunto de datos y se muestran patrones relacionados a alineamientos múltiples de secuencia. Además de una visualización de la proteína de interés, en su estructura 3D y las sustituciones que exhiben cada variante. Esto puede observarse en la Figura 4.4.

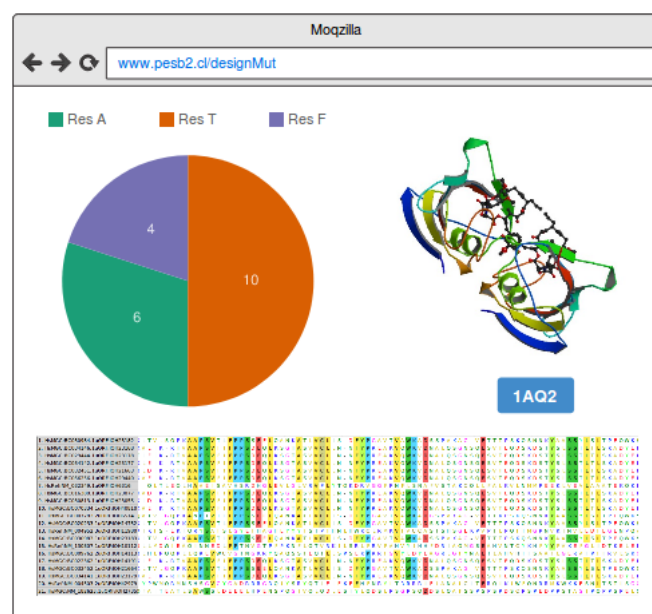


Fig. 4.4 Esquema representativo interfaz de descripción general del conjunto de datos.

- **Propiedades fisicoquímicas:** Una representación de las propiedades fisicoquímicas abarca cuáles fueron las seleccionadas y las definiciones de éstas, por otro lado se expone el aporte a la varianza que éstas presentan y cómo varían por cada secuencia. Un esquema representativo de diseño de la interfaz, puede observarse en la Figura 4.5.
- **Espectros de frecuencia:** Se exponen el conjunto de digitalizaciones de las propiedades fisicoquímicas, asociados a los residuos que brindan la mayor caracterización y cómo

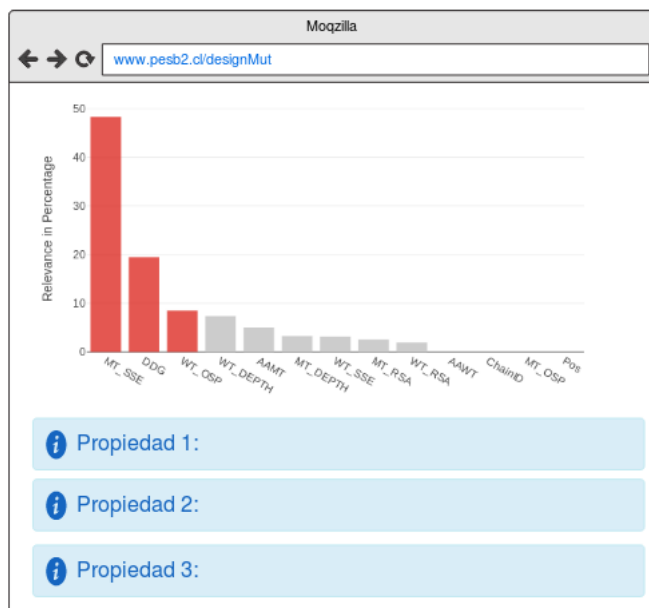


Fig. 4.5 Esquema representativo interfaz de visualización de propiedades fisicoquímicas.

estos inciden en el espectro. En la Figura 4.6, se expone un esquema de la visualización de los espectros, junto a los residuos relevantes.

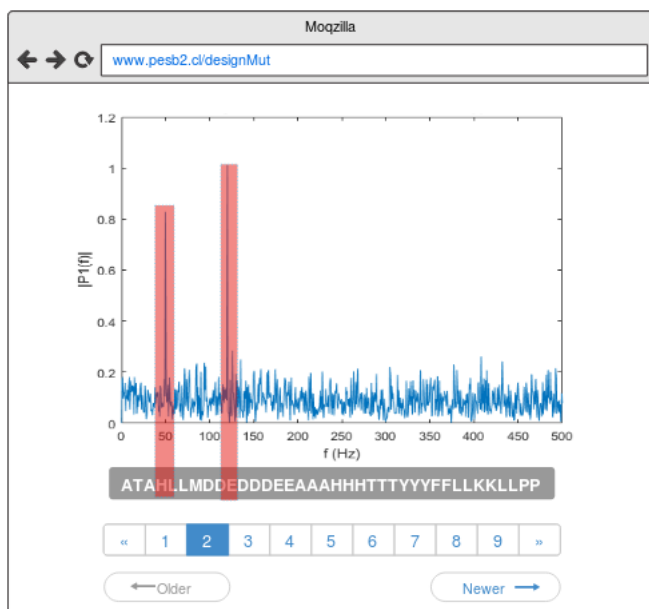


Fig. 4.6 Esquema representativo interfaz de visualización de espectros de frecuencias y residuos relevantes.

- **Entrenamiento de Modelos:** Los modelos se exponen asociados a las medidas de desempeño obtenidas y los miembros participantes en el conjunto de modelos. Esto es debido a que se utilizará la estrategia expuesta en el capítulo 2, por lo que se obtiene un sistema de meta-modelos, cuyo esquema representativo de visualización se observa en la Figura 4.7.

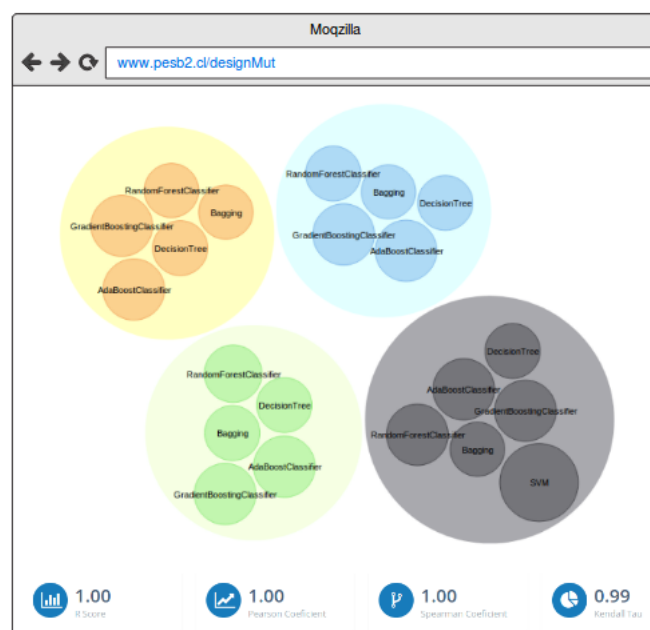


Fig. 4.7 Esquema representativo interfaz de visualización de los meta modelos y sus medidas de desempeño.

- **Mutaciones propuestas:** Para el conjunto de datos, se proponen diferentes mutaciones, evaluadas con respecto a la respuesta y factibilidad de éstas, a su vez, se exponen estadísticas asociadas a qué residuos son más factibles de mutar en base a las posiciones de estos.

4.3.6 Consideraciones generales

Como consideraciones generales, se tiene que el conjunto de secuencias o variantes a estudiar, debe poseer una respuesta reportada, es decir, la variable de interés debe conocerse de ante mano, con el fin de poder entrenar los modelos de clasificación o regresión y así poder analizar nuevas mutaciones.

Por otro lado, la secuencia original, debe presentar su estructura en formato *.pdb, ya que, se requiere para el uso de SDM, la cual es la herramienta de filtro asociada a la estabilidad de la proteína con respecto a los cambios o sustituciones propuestas.

Un punto importante a destacar, es que, todos los pasos relacionados a la generación de modelos, implicando desde la codificación y posterior digitalización de las secuencias hasta la fase de evaluación del desempeño, serán ejecutadas en torno a Jobs que cree el usuario y serán procesadas internamente en servidor en un gestor de colas, con el fin de optimizar los procesos de ejecución. Esto es debido, a que el tiempo de entrenamiento puede ser elevado si existen muchas variantes en el conjunto de datos, razón por la cual, una vez que el job finalice, se notificará vía email al usuario, en donde, él podrá acceder al sistema y revisar los resultados obtenidos.

Se propone esta herramienta como un servicio de exploración, en el cual, se evalúan un conjunto de mutaciones y se sugieren nuevos elementos, de tal manera, que el usuario pueda tener un poco más claro el panorama. No obstante, las mutaciones son proposiciones basados en los filtros y en los modelos aplicados. Por lo que, está condicionado por el conjunto de datos inicial y las medidas de desempeño obtenidas.

Referencias

- [1] Abdelaziz, A., Elhoseny, M., Salama, A. S., and Riad, A. (2018). A machine learning model for improving healthcare services on cloud computing environment. *Measurement*, 119:117 – 128.
- [2] Abola, E. E., Bernstein, F. C., and Koetzle, T. F. (1984). The protein data bank. In *Neutrons in Biology*, pages 441–441. Springer.
- [3] Alexov, E., Zhang, J., Wang, L., Zhenirovskyy, M., Gao, Y., and Zhang, Z. (2012). Predicting folding free energy changes upon single point mutations. *Bioinformatics*, 28(5):664–671.
- [4] Alm, C. O., Roth, D., and Sproat, R. (2005). Emotions from text: machine learning for text-based emotion prediction. In *Proceedings of the conference on human language technology and empirical methods in natural language processing*, pages 579–586. Association for Computational Linguistics.
- [5] Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410.
- [6] Amari, S.-i. and Wu, S. (1999). Improving support vector machine classifiers by modifying kernel functions. *Neural Networks*, 12(6):783–789.
- [7] Ancien, F., Pucci, F., Godfroid, M., and Rooman, M. (2018). Prediction and interpretation of deleterious coding variants in terms of protein structural stability. *Scientific reports*, 8(1):4480.
- [8] Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M. W., Pfau, D., Schaul, T., Shillingford, B., and De Freitas, N. (2016). Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*, pages 3981–3989.
- [9] Arel, I., Rose, D. C., Karnowski, T. P., et al. (2010). Deep machine learning-a new frontier in artificial intelligence research. *IEEE computational intelligence magazine*, 5(4):13–18.
- [10] Arnold, F. H. (1998). Design by directed evolution. *Accounts of chemical research*, 31(3):125–131.
- [11] Artac, M., Jogan, M., and Leonardis, A. (2002). Incremental pca for on-line visual learning and recognition. In *Object recognition supported by user interaction for service robots*, volume 3, pages 781–784. IEEE.

- [12] Arthur, D. and Vassilvitskii, S. (2007). k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035. Society for Industrial and Applied Mathematics.
- [13] Barandiaran, I. (1998). The random subspace method for constructing decision forests. *IEEE transactions on pattern analysis and machine intelligence*, 20(8).
- [14] Barenboim, M., Masso, M., Vaisman, I. I., and Jamison, D. C. (2008). Statistical geometry based prediction of nonsynonymous snp functional effects using random forest and neuro-fuzzy classifiers. *Proteins: Structure, Function, and Bioinformatics*, 71(4):1930–1939.
- [15] Battiti, R. (1992). First-and second-order methods for learning: between steepest descent and newton’s method. *Neural computation*, 4(2):141–166.
- [16] Bava, K. A., Gromiha, M. M., Uedaira, H., Kitajima, K., and Sarai, A. (2004). Protherm, version 4.0: thermodynamic database for proteins and mutants. *Nucleic acids research*, 32(suppl_1):D120–D121.
- [17] Bedbrook, C. N., Yang, K. K., Robinson, J. E., Gradinaru, V., and Arnold, F. H. (2019). Machine learning-guided channelrhodopsin engineering enables minimally-invasive optogenetics.
- [18] Bengio, Y. (2009). Learning deep architectures for ai. *Foundations and Trends® in Machine Learning*, 2(1):1–127.
- [19] Bengio, Y., Courville, A., and Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828.
- [20] Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T. N., Weissig, H., Shindyalov, I. N., and Bourne, P. E. (2000). The protein data bank. *Nucleic acids research*, 28(1):235–242.
- [21] Berry, M. J. and Linoff, G. S. (2004). *Data mining techniques: for marketing, sales, and customer relationship management*. John Wiley & Sons.
- [22] Bhargava, N., Sharma, G., Bhargava, R., and Mathuria, M. (2013). Decision tree analysis on j48 algorithm for data mining. *Proceedings of International Journal of Advanced Research in Computer Science and Software Engineering*, 3(6).
- [23] Boeckmann, B., Bairoch, A., Apweiler, R., Blatter, M.-C., Estreicher, A., Gasteiger, E., Martin, M. J., Michoud, K., O’donovan, C., Phan, I., et al. (2003). The swiss-prot protein knowledgebase and its supplement trembl in 2003. *Nucleic acids research*, 31(1):365–370.
- [Bordner and Abagyan] Bordner, A. J. and Abagyan, R. A. Large-scale prediction of protein geometry and stability changes for arbitrary single point mutations. *Proteins: Structure, Function, and Bioinformatics*, 57(2):400–413.
- [25] Braha, D. and Shmilovici, A. (2002). Data mining for improving a cleaning process in the semiconductor industry. *IEEE Transactions on Semiconductor Manufacturing*, 15(1):91–101.

- [26] Breiman, L. (1996). Bagging predictors. *Machine learning*, 24(2):123–140.
- [27] Breiman, L. (1998). Arcing classifier (with discussion and a rejoinder by the author). *Ann. Statist.*, 26(3):801–849.
- [28] Breiman, L. (1999). Pasting small votes for classification in large databases and on-line. *Machine learning*, 36(1-2):85–103.
- [29] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.
- [30] Breiman, L. (2017). *Classification and regression trees*. Routledge.
- [31] Brigham, E. O. and Brigham, E. O. (1988). *The fast Fourier transform and its applications*, volume 448. prentice Hall Englewood Cliffs, NJ.
- [32] Broom, A., Jacobi, Z., Trainor, K., and Meiering, E. M. (2017a). Computational tools help improve protein stability but with a solubility tradeoff. *J Biol Chem*, 292(35):14349–14361. 28710274[pmid].
- [33] Broom, A., Jacobi, Z., Trainor, K., and Meiering, E. M. (2017b). Computational tools help improve protein stability but with a solubility tradeoff. *Journal of Biological Chemistry*, 292(35):14349–14361.
- [34] Brownlee, J. (2017). Why one-hot encode data in machine learning.
- [35] Cadet, F., Fontaine, N., Vetrivel, I., Chong, M. N. F., Savriama, O., Cadet, X., and Charton, P. (2018). Application of fourier transform and proteochemometrics principles to protein engineering. *BMC bioinformatics*, 19(1):382.
- [36] Caliński, T. and Harabasz, J. (1974). A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods*, 3(1):1–27.
- [37] Canutescu, A. A., Shelenkov, A. A., and Dunbrack Jr, R. L. (2003). A graph-theory algorithm for rapid protein side-chain prediction. *Protein science*, 12(9):2001–2014.
- [38] CAO, Y., MIAO, Q.-G., LIU, J.-C., and GAO, L. (2013). Advance and prospects of adaboost algorithm. *Acta Automatica Sinica*, 39(6):745 – 758.
- [39] Capriotti, E., Fariselli, P., and Casadio, R. (2005a). I-mutant2. 0: predicting stability changes upon mutation from the protein sequence or structure. *Nucleic acids research*, 33(suppl_2):W306–W310.
- [40] Capriotti, E., Fariselli, P., and Casadio, R. (2005b). I-mutant2.0: predicting stability changes upon mutation from the protein sequence or structure. *Nucleic Acids Res*, 33(Web Server issue):W306–W310. 15980478[pmid].
- [41] Capriotti, E., Fariselli, P., Rossi, I., and Casadio, R. (2008a). A three-state prediction of single point mutations on protein stability changes. *BMC bioinformatics*, 9(2):S6.
- [42] Capriotti, E., Fariselli, P., Rossi, I., and Casadio, R. (2008b). A three-state prediction of single point mutations on protein stability changes. *BMC Bioinformatics*, 9(2):S6.

- [43] Carpenter, J. F., Pikal, M. J., Chang, B. S., and Randolph, T. W. (1997). Rational design of stable lyophilized protein formulations: some practical advice. *Pharmaceutical research*, 14(8):969–975.
- [44] Case, D. A., Cheatham III, T. E., Darden, T., Gohlke, H., Luo, R., Merz Jr, K. M., Onufriev, A., Simmerling, C., Wang, B., and Woods, R. J. (2005). The amber biomolecular simulation programs. *Journal of computational chemistry*, 26(16):1668–1688.
- [45] Castillo, I., Schmidt-Hieber, J., Van der Vaart, A., et al. (2015). Bayesian linear regression with sparse priors. *The Annals of Statistics*, 43(5):1986–2018.
- [46] Chang, C.-C. and Lin, C.-J. (2011). Libsvm: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27.
- [47] Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P. (2002). Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.
- [48] Chen, C., Lu, Z., and Ciucci, F. (2017). Data mining of molecular dynamics data reveals li diffusion characteristics in garnet $\text{Li}_7\text{La}_3\text{Zr}_2\text{O}_{12}$. *Scientific Reports*, 7:40769 EP –. Article.
- [49] Chen, J., Huang, H., Tian, S., and Qu, Y. (2009). Feature selection for text classification with naïve bayes. *Expert Systems with Applications*, 36(3, Part 1):5432–5435.
- [50] Chen, M., Hao, Y., Hwang, K., Wang, L., and Wang, L. (2017). Disease prediction by machine learning over big data from healthcare communities. *IEEE Access*, 5:8869–8879.
- [51] Cheng, Y. (1995). Mean shift, mode seeking, and clustering. *IEEE transactions on pattern analysis and machine intelligence*, 17(8):790–799.
- [52] Chien, C.-F. and Chen, L.-F. (2008). Data mining to improve personnel selection and enhance human capital: A case study in high-technology industry. *Expert Systems with Applications*, 34(1):280–290.
- [53] Cohen, P., West, S. G., and Aiken, L. S. (2014). *Applied multiple regression/correlation analysis for the behavioral sciences*. Psychology Press.
- [54] Cooley, J. W., Lewis, P., and Welch, P. (1970). The fast fourier transform algorithm: Programming considerations in the calculation of sine, cosine and laplace transforms. *Journal of Sound and Vibration*, 12(3):315–337.
- [55] Cooley, R., Mobasher, B., Srivastava, J., et al. (1997). Web mining: Information and pattern discovery on the world wide web. In *ictai*, volume 97, pages 558–567.
- [56] Cosic, I. (1994). Macromolecular bioactivity: is it resonant interaction between macromolecules?-theory and applications. *IEEE Transactions on Biomedical Engineering*, 41(12):1101–1114.
- [57] Cosic, I., Cosic, D., and Lazar, K. (2016). Analysis of tumor necrosis factor function using the resonant recognition model. *Cell biochemistry and biophysics*, 74(2):175–180.

- [58] Čosić, I. and NESIC, D. (1987). Prediction of ‘hot spots’ in sv40 enhancer and relation with experimental data. *European journal of biochemistry*, 170(1-2):247–252.
- [59] Cui-xiang, Z., Guo-qiang, H., and Ming-He, H. (2005). Some new parallel fast fourier transform algorithms. In *Sixth International Conference on Parallel and Distributed Computing Applications and Technologies (PDCAT’05)*, pages 624–628. IEEE.
- [60] Curtarolo, S., Morgan, D., Persson, K., Rodgers, J., and Ceder, G. (2003). Predicting crystal structures with data mining of quantum calculations. *Phys. Rev. Lett.*, 91:135503.
- [61] Danielsson, P.-E. (1980). Euclidean distance mapping. *Computer Graphics and Image Processing*, 14(3):227 – 248.
- [62] Deng, L. and Yu, D. (2014). Deep learning: Methods and applications. *Foundations and Trends® in Signal Processing*, 7(3–4):197–387.
- [63] Dietterich, T. G. (2000). Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, pages 1–15. Springer.
- [64] Duan, L., Street, W. N., and Xu, E. (2011). Healthcare information systems: data mining methods in the creation of a clinical recommender system. *Enterprise Information Systems*, 5(2):169–181.
- [65] Dudani, S. A. (1976). The distance-weighted k-nearest-neighbor rule. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-6(4):325–327.
- [66] Dunham, M. H. (2006). *Data mining: Introductory and advanced topics*. Pearson Education India.
- [67] Duygulu, P., Barnard, K., de Freitas, J. F. G., and Forsyth, D. A. (2002). Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In Heyden, A., Sparr, G., Nielsen, M., and Johansen, P., editors, *Computer Vision — ECCV 2002*, pages 97–112, Berlin, Heidelberg. Springer Berlin Heidelberg.
- [68] Eswar, N., Webb, B., Marti-Renom, M. A., Madhusudhan, M., Eramian, D., Shen, M.-y., Pieper, U., and Sali, A. (2006). Comparative protein structure modeling using modeller. *Current protocols in bioinformatics*, 15(1):5–6.
- [69] Fayyad, U., Piatetsky-Shapiro, G., and Smyth, P. (1996a). From data mining to knowledge discovery in databases. *AI magazine*, 17(3):37–37.
- [70] Fayyad, U., Piatetsky-Shapiro, G., and Smyth, P. (1996b). The kdd process for extracting useful knowledge from volumes of data. *Communications of the ACM*, 39(11):27–34.
- [71] Fayyad, U. M., Piatetsky-Shapiro, G., Smyth, P., et al. (1996c). Knowledge discovery and data mining: Towards a unifying framework. In *KDD*, volume 96, pages 82–88.
- [72] Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1126–1135. JMLR. org.
- [73] Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188.

- [74] Forbes, S. A., Bindal, N., Bamford, S., Cole, C., Kok, C. Y., Beare, D., Jia, M., Shepherd, R., Leung, K., Menzies, A., Teague, J. W., Campbell, P. J., Stratton, M. R., and Futreal, P. A. (2010). Cosmic: mining complete cancer genomes in the catalogue of somatic mutations in cancer. *Nucleic Acids Research*, 39(suppl_1):D945–D950.
- [75] Fortunato, S. and Castellano, C. (2012). Community structure in graphs. *Computational Complexity: Theory, Techniques, and Applications*, pages 490–512.
- [76] Freund, Y. and Mason, L. (1999). The alternating decision tree learning algorithm. In *icml*, volume 99, pages 124–133.
- [77] Frey, B. J. and Dueck, D. (2007). Clustering by passing messages between data points. *science*, 315(5814):972–976.
- [78] Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5):1189–1232.
- [79] Friedman, J. H. (2002). Stochastic gradient boosting. *Computational Statistics AND Data Analysis*, 38(4):367 – 378. Nonlinear Methods and Data Mining.
- [80] Gallou, C., Junien, C., Joly, D., Staroz, F., Orfanelli, M. T., and Bérout, C. (1998). Software and database for the analysis of mutations in the VHL gene. *Nucleic Acids Research*, 26(1):256–258.
- [81] Ge, Z., Song, Z., Ding, S. X., and Huang, B. (2017). Data mining and analytics in the process industry: The role of machine learning. *IEEE Access*, 5:20590–20616.
- [82] Getov, I., Petukh, M., and Alexov, E. (2016a). Saafec: Predicting the effect of single point mutations on protein folding free energy using a knowledge-modified mm/pbsa approach. *Int J Mol Sci*, 17(4):512–512. 27070572[pmid].
- [83] Getov, I., Petukh, M., and Alexov, E. (2016b). Saafec: predicting the effect of single point mutations on protein folding free energy using a knowledge-modified mm/pbsa approach. *International journal of molecular sciences*, 17(4):512.
- [84] Golub, G. H., Heath, M., and Wahba, G. (1979). Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223.
- [85] Golub, G. H. and Van Loan, C. F. (1980). An analysis of the total least squares problem. *SIAM journal on numerical analysis*, 17(6):883–893.
- [86] Gossage, L., Pires, D., Olivera-Nappa, A., A. Asenjo, J., Bycroft, M., Blundell, T., and Eisen, T. (2014). An integrated computational approach can classify vhl missense mutations according to risk of clear cell renal carcinoma. *Human molecular genetics*, 23.
- [87] Granitto, P. M., Furlanello, C., Biasioli, F., and Gasperi, F. (2006). Recursive feature elimination with random forest for ptr-ms analysis of agroindustrial products. *Chemometrics and Intelligent Laboratory Systems*, 83(2):83–90.
- [88] Graybill, F. A. (1976). *Theory and application of the linear model*. Number QA279. G72 1976. Duxbury press North Scituate, MA.

- [89] Guex, N. and Peitsch, M. C. (1997). Swiss-model and the swiss-pdb viewer: An environment for comparative protein modeling. *ELECTROPHORESIS*, 18(15):2714–2723.
- [90] Guyon, I., Boser, B., and Vapnik, V. (1993). Automatic capacity tuning of very large vc-dimension classifiers. In *Advances in neural information processing systems*, pages 147–155.
- [91] Han, J. and Gao, J. (2009). Research challenges for data mining in science and engineering. *Next Generation of Data Mining*, pages 1–18.
- [92] Hand, D. J. (2006). Data mining. *Encyclopedia of Environmetrics*, 2.
- [93] Hans, C. (2009). Bayesian lasso regression. *Biometrika*, 96(4):835–845.
- [94] Hastie, T., Rosset, S., Zhu, J., and Zou, H. (2009). Multi-class adaboost. *Statistics and its Interface*, 2(3):349–360.
- [95] Hawkins, D. M. (2004). The problem of overfitting. *Journal of chemical information and computer sciences*, 44(1):1–12.
- [96] HECHT-NIELSEN, R. (1992). Iii.3 - theory of the backpropagation neural network**based on “nonindent” by robert hecht-nielsen, which appeared in proceedings of the international joint conference on neural networks 1, 593–611, june 1989. © 1989 ieee. In Wechsler, H., editor, *Neural Networks for Perception*, pages 65 – 93. Academic Press.
- [97] Hinton, G. E. and Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507.
- [98] Hochreiter, S., Younger, A. S., and Conwell, P. R. (2001). Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks*, pages 87–94. Springer.
- [99] Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67.
- [100] Hofmann, M. and Klinkenberg, R. (2013). *RapidMiner: Data mining use cases and business analytics applications*. CRC Press.
- [101] Hssina, B., Merbouha, A., Ezzikouri, H., and Erritali, M. (2014). A comparative study of decision tree id3 and c4. 5. *International Journal of Advanced Computer Science and Applications*, 4(2):0–0.
- [102] Hua, S. and Sun, Z. (2001). A novel method of protein secondary structure prediction with high segment overlap measure: support vector machine approach. *Journal of molecular biology*, 308(2):397–407.
- [103] Indyk, P. and Motwani, R. (1998). Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613. ACM.
- [104] Jain, A. K., Dubes, R. C., et al. (1988). *Algorithms for clustering data*, volume 6. Prentice hall Englewood Cliffs.

- [105] Jain, A. K., Murty, M. N., and Flynn, P. J. (1999). Data clustering: a review. *ACM computing surveys (CSUR)*, 31(3):264–323.
- [106] Japkowicz, N. and Stephen, S. (2002). The class imbalance problem: A systematic study. *Intelligent data analysis*, 6(5):429–449.
- [107] Jespersen, M. C., Peters, B., Nielsen, M., and Marcatili, P. (2017). Bepipred-2.0: improving sequence-based b-cell epitope prediction using conformational epitopes. *Nucleic acids research*, 45(W1):W24–W29.
- [108] John, G. H. and Langley, P. (1995). Estimating continuous distributions in bayesian classifiers. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pages 338–345. Morgan Kaufmann Publishers Inc.
- [109] Johnson, S. C. (1967). Hierarchical clustering schemes. *Psychometrika*, 32(3):241–254.
- [110] Jolliffe, I. (2011). *Principal component analysis*. Springer.
- [111] Kawashima, S. and Kanehisa, M. (2000). Aaindex: Amino acid index database. *Nucleic Acids Research*, 28(1):374–374.
- [112] Keller, J. M., Gray, M. R., and Givens, J. A. (1985). A fuzzy k-nearest neighbor algorithm. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-15(4):580–585.
- [113] Kent, W. J. (2002). Blat—the blast-like alignment tool. *Genome research*, 12(4):656–664.
- [114] Khan, S. and Vihinen, M. (2010). Performance of protein stability predictors. *Human Mutation*, 31(6):675–684.
- [115] Kohavi, R. et al. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Ijcai*, volume 14, pages 1137–1145. Montreal, Canada.
- [116] Kolba, D. and Parks, T. (1977). A prime factor fft algorithm using high-speed convolution. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 25(4):281–294.
- [117] Kotsiantis, S. B., Zaharakis, I., and Pintelas, P. (2007). Supervised machine learning: A review of classification techniques. *Emerging artificial intelligence applications in computer engineering*, 160:3–24.
- [118] Krasner, G. E., Pope, S. T., et al. (1988). A description of the model-view-controller user interface paradigm in the smalltalk-80 system. *Journal of object oriented programming*, 1(3):26–49.
- [119] Larrañaga, P., Calvo, B., Santana, R., Bielza, C., Galdiano, J., Inza, I., Lozano, J. A., Armañanzas, R., Santafé, G., Pérez, A., and Robles, V. (2006). Machine learning in bioinformatics. *Briefings in Bioinformatics*, 7(1):86–112.

- [120] Leaver-Fay, A., Tyka, M., Lewis, S. M., Lange, O. F., Thompson, J., Jacak, R., Kaufman, K. W., Renfrew, P. D., Smith, C. A., Sheffler, W., et al. (2011). Rosetta3: an object-oriented software suite for the simulation and design of macromolecules. In *Methods in enzymology*, volume 487, pages 545–574. Elsevier.
- [121] Lee, J. K., Williams, P. D., and Cheon, S. (2008). Data mining in genomics. *Clinics in Laboratory Medicine*, 28(1):145–166.
- [122] Lewis, D. D. (1998). Naive (bayes) at forty: The independence assumption in information retrieval. In Nédellec, C. and Rouveirol, C., editors, *Machine Learning: ECML-98*, pages 4–15, Berlin, Heidelberg. Springer Berlin Heidelberg.
- [123] Li, M., Wang, J., and Chen, J. (2008). A fast agglomerate algorithm for mining functional modules in protein interaction networks. In *2008 International Conference on BioMedical Engineering and Informatics*, volume 1, pages 3–7.
- [124] Li, M., Wang, J., and Chen, J. (2008). A fast agglomerate algorithm for mining functional modules in protein interaction networks. In *2008 International conference on biomedical engineering and informatics*, volume 1, pages 3–7. IEEE.
- [125] Liao, H. and Xu, Z. (2015). Approaches to manage hesitant fuzzy linguistic information based on the cosine distance and similarity measures for hflts and their application in qualitative decision making. *Expert Systems with Applications*, 42(12):5328 – 5336.
- [126] Liszewski, K. (2015). Speeding up the protein assembly line. *Genetic Engineering & Biotechnology News*, 35(04):1–10.
- [127] Louppe, G. and Geurts, P. (2012). Ensembles on random patches. In Flach, P. A., De Bie, T., and Cristianini, N., editors, *Machine Learning and Knowledge Discovery in Databases*, pages 346–361, Berlin, Heidelberg. Springer Berlin Heidelberg.
- [128] Lutz, S. (2010). Beyond directed evolution—semi-rational protein engineering and design. *Current opinion in biotechnology*, 21(6):734–743.
- [129] Lutz, S., Bornscheuer, U. T., et al. (2009). *Protein engineering handbook*, volume 1. Wiley Online Library.
- [130] Lyskov, S. and Gray, J. J. (2008). The rosettadock server for local protein–protein docking. *Nucleic acids research*, 36(suppl_2):W233–W238.
- [131] Ma, X., Wu, Y.-J., Wang, Y., Chen, F., and Liu, J. (2013). Mining smart card data for transit riders’ travel patterns. *Transportation Research Part C: Emerging Technologies*, 36:1–12.
- [132] Maesschalck, R. D., Jouan-Rimbaud, D., and Massart, D. (2000). The mahalanobis distance. *Chemometrics and Intelligent Laboratory Systems*, 50(1):1 – 18.
- [133] Manning, C., Raghavan, P., and Schütze, H. (2010). Introduction to information retrieval. *Natural Language Engineering*, 16(1):100–103.
- [134] Mao, L., Wang, Y., Liu, Y., and Hu, X. (2004). Molecular determinants for atp-binding in proteins: A data mining and quantum chemical analysis. *Journal of Molecular Biology*, 336(3):787 – 807.

- [135] Martin, A. J., Contreras-Riquelme, S., Dominguez, C., and Perez-Acle, T. (2017). Loto: a graphlet based method for the comparison of local topology between gene regulatory networks. *PeerJ*, 5:e3052.
- [136] Masso, M. and Vaisman, I. I. (2008). Accurate prediction of stability changes in protein mutants by combining machine learning with structure based computational mutagenesis. *Bioinformatics*, 24(18):2002–2009.
- [137] Masso, M. and Vaisman, I. I. (2010). Auto-mute: web-based tools for predicting stability changes in proteins due to single amino acid replacements. *Protein Engineering, Design and Selection*, 23(8):683–687.
- [138] McGuffin, L. J., Atkins, J. D., Salehe, B. R., Shuid, A. N., and Roche, D. B. (2015). Intfold: an integrated server for modelling protein structures and functions from amino acid sequences. *Nucleic acids research*, 43(W1):W169–W173.
- [139] McKinney, W. et al. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference*, volume 445, pages 51–56. Austin, TX.
- [140] Metsis, V., Androutsopoulos, I., and Paliouras, G. (2006). Spam filtering with naive bayes-which naive bayes? In *CEAS*, volume 17, pages 28–69. Mountain View, CA.
- [141] Michie, D., Spiegelhalter, D. J., Taylor, C., et al. (1994). Machine learning. *Neural and Statistical Classification*, 13.
- [142] Moon, T. K. (1996). The expectation-maximization algorithm. *IEEE Signal processing magazine*, 13(6):47–60.
- [143] Muggleton, S., King, R. D., and Stenberg, M. J. E. (1992). Protein secondary structure prediction using logic-based machine learning. *Protein Engineering, Design and Selection*, 5(7):647–657.
- [144] Muquet, B., Wang, Z., Giannakis, G. B., De Courville, M., and Duhamel, P. (2002). Cyclic prefixing or zero padding for wireless multicarrier transmissions? *IEEE Transactions on communications*, 50(12):2136–2148.
- [145] Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. (2011). Reading digits in natural images with unsupervised feature learning.
- [146] Newman, M. E. (2004). Detecting community structure in networks. *The European Physical Journal B*, 38(2):321–330.
- [147] Newman, M. E. J. (2006). Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582.
- [148] Obenshain, M. K. (2004). Application of data mining techniques to healthcare data. *Infection Control & Hospital Epidemiology*, 25(8):690–695.
- [149] Odorico, M. and Pellequer, J.-L. (2003). Bepitope: predicting the location of continuous epitopes and patterns in proteins. *Journal of Molecular Recognition*, 16(1):20–22.

- [150] Oliphant, T. E. (2007). Python for scientific computing. *Computing in Science & Engineering*, 9(3):10–20.
- [151] Olivera-Nappa, A., Andrews, B. A., and Asenjo, J. A. (2011). Mutagenesis objective search and selection tool (mosst): an algorithm to predict structure-function related mutations in proteins. *BMC Bioinformatics*, 12(1):122.
- [152] Orengo, C. A., Michie, A. D., Jones, S., Jones, D. T., Swindells, M. B., and Thornton, J. M. (1997). Cath—a hierarchic classification of protein domain structures. *Structure*, 5(8):1093–1109.
- [153] Ozbudak, O. and Dokur, Z. (2014). Protein fold classification using kohonen’s self-organizing map. In *IWBBIO*, pages 903–911.
- [154] Pandurangan, A. P., Ochoa-Montaña, B., Ascher, D. B., and Blundell, T. L. (2017). Sdm: a server for predicting effects of mutations on protein stability. *Nucleic Acids Res*, 45(W1):W229–W235. 28525590[pmid].
- [155] Parthiban, V., Gromiha, M. M., and Schomburg, D. (2006). Cupsat: prediction of protein stability upon point mutations. *Nucleic Acids Res*, 34(Web Server issue):W239–W242. 16845001[pmid].
- [156] Pati, Y. C., Rezaiifar, R., and Krishnaprasad, P. S. (1993). Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Proceedings of 27th Asilomar conference on signals, systems and computers*, pages 40–44. IEEE.
- [157] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830.
- [158] Peng, H., Long, F., and Ding, C. (2005). Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (8):1226–1238.
- [159] Perlibakas, V. (2004). Distance measures for pca-based face recognition. *Pattern Recognition Letters*, 25(6):711 – 724.
- [160] Petukh, M., Dai, L., and Alexov, E. (2016). Saambe: webserver to predict the change of binding free energy caused by amino acids mutations. *International journal of molecular sciences*, 17(4):547.
- [161] Phillips, J. C., Braun, R., Wang, W., Gumbart, J., Tajkhorshid, E., Villa, E., Chipot, C., Skeel, R. D., Kale, L., and Schulten, K. (2005). Scalable molecular dynamics with namd. *Journal of computational chemistry*, 26(16):1781–1802.
- [162] Potapov, V., Cohen, M., and Schreiber, G. (2009). Assessing computational methods for predicting protein stability upon mutation: good on average but not in the details. *Protein Engineering, Design and Selection*, 22(9):553–560.

- [163] Quan, L., Lv, Q., and Zhang, Y. (2016a). Strum: structure-based prediction of protein stability changes upon single-point mutation. *Bioinformatics*, 32(19):2936–2946. 27318206[pmid].
- [164] Quan, L., Lv, Q., and Zhang, Y. (2016b). Strum: structure-based prediction of protein stability changes upon single-point mutation. *Bioinformatics*, 32(19):2936–2946.
- [165] Ramakrishnan, R., Dral, P. O., Rupp, M., and von Lilienfeld, O. A. (2015). Big data meets quantum chemistry approximations: The d-machine learning approach. *Journal of Chemical Theory and Computation*, 11(5):2087–2096.
- [166] Rao, K. R. and Yip, P. (2014). *Discrete cosine transform: algorithms, advantages, applications*. Academic press.
- [167] Rebhan, M., Chalifa-Caspi, V., Prilusky, J., and Lancet, D. (1998). Genecards: a novel functional genomics compendium with automated data mining and query reformulation support. *Bioinformatics*, 14(8):656–664.
- [168] Reetz, M. T., Kahakeaw, D., and Lohmer, R. (2008). Addressing the numbers problem in directed evolution. *ChemBioChem*, 9(11):1797–1804.
- [169] Release, S. (2016). 1: Maestro. *Schrödinger, LLC, New York, NY, USA*.
- [170] Rohl, C. A., Strauss, C. E., Misura, K. M., and Baker, D. (2004). Protein structure prediction using rosetta. In *Methods in enzymology*, volume 383, pages 66–93. Elsevier.
- [171] Rosenberg, A. and Hirschberg, J. (2007). V-measure: A conditional entropy-based external cluster evaluation measure. In *Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL)*, pages 410–420.
- [172] Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65.
- [173] Ruck, D. W., Rogers, S. K., Kabrisky, M., Oxley, M. E., and Suter, B. W. (1990). The multilayer perceptron as an approximation to a bayes optimal discriminant function. *IEEE Transactions on Neural Networks*, 1(4):296–298.
- [174] Saeys, Y., Abeel, T., and Van de Peer, Y. (2008). Robust feature selection using ensemble feature selection techniques. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 313–325. Springer.
- [175] Saha, S. and Raghava, G. (2008). Abcpred benchmarking datasets. 2006a.
- [176] Sarwar, B., Karypis, G., Konstan, J., and Riedl, J. (2000). Application of dimensionality reduction in recommender system-a case study. Technical report, Minnesota Univ Minneapolis Dept of Computer Science.
- [177] Schmidhuber, J. (1987). *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München.

- [178] Schmidhuber, J. (1995). On learning how to learn learning strategies.
- [179] Scholkopf, B. and Smola, A. J. (2001). *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press.
- [180] Schomburg, I., Chang, A., Ebeling, C., Gremse, M., Heldt, C., Huhn, G., and Schomburg, D. (2004). Brenda, the enzyme database: updates and major new developments. *Nucleic acids research*, 32(suppl_1):D431–D433.
- [181] Schueler-Furman, O. and Baker, D. (2003). Conserved residue clustering and protein structure prediction. *Proteins: Structure, Function, and Bioinformatics*, 52(2):225–235.
- [182] Schymkowitz, J., Borg, J., Stricher, F., Nys, R., Rousseau, F., and Serrano, L. (2005). The foldx web server: an online force field. *Nucleic Acids Res*, 33(Web Server issue):W382–W388. 15980494[pmid].
- [183] Smola, A. J. and Schölkopf, B. (2004). A tutorial on support vector regression. *Statistics and computing*, 14(3):199–222.
- [184] Sneddon, I. N. (1995). *Fourier transforms*. Courier Corporation.
- [185] Sokal, R. R. and Rohlf, F. J. (1962). The comparison of dendrograms by objective methods. *Taxon*, 11(2):33–40.
- [186] Sun, L., Li, L., Peng, Y., Jia, Z., and Alexov, E. (2017). Predicting protein-dna binding free energy change upon missense mutations using modified mm/pbsa approach: Sampdi webserver. *Bioinformatics*, 34(5):779–786.
- [187] Tamura, K., Stecher, G., Peterson, D., Filipski, A., and Kumar, S. (2013). MEGA6: Molecular Evolutionary Genetics Analysis Version 6.0. *Molecular Biology and Evolution*, 30(12):2725–2729.
- [188] Tan, S. (2005). Neighbor-weighted k-nearest neighbor for unbalanced text corpus. *Expert Systems with Applications*, 28(4):667 – 671.
- [189] Thompson, J. D., Gibson, T. J., Plewniak, F., Jeanmougin, F., and Higgins, D. G. (1997). The clustal_x windows interface: Flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Research*, 25(24):4876–4882.
- [190] Tian, J., Wu, N., Chu, X., and Fan, Y. (2010). Predicting changes in protein thermostability brought about by single- or multi-site mutations. *BMC Bioinformatics*, 11(1):370.
- [191] Tovchigrechko, A. and Vakser, I. A. (2006). Gramm-x public web server for protein–protein docking. *Nucleic acids research*, 34(suppl_2):W310–W314.
- [192] Tran, T. N., Drab, K., and Daszykowski, M. (2013). Revised dbscan algorithm to cluster data with dense adjacent clusters. *Chemometrics and Intelligent Laboratory Systems*, 120:92–96.
- [193] Trott, O. and Olson, A. J. (2010). Autodock vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry*, 31(2):455–461.

- [194] Utgoff, P. E. (1986). Shift of bias for inductive concept learning. *Machine learning: An artificial intelligence approach*, 2:107–148.
- [195] Vaisman, I. I. and Masso, M. (2008). Accurate prediction of stability changes in protein mutants by combining machine learning with structure based computational mutagenesis. *Bioinformatics*, 24(18):2002–2009.
- [196] Van Der Maaten, L., Postma, E., and Van den Herik, J. (2009). Dimensionality reduction: a comparative. *J Mach Learn Res*, 10(66-71):13.
- [197] Van Der Walt, S., Colbert, S. C., and Varoquaux, G. (2011). The numpy array: a structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2):22.
- [198] Vehtari, A., Gelman, A., and Gabry, J. (2017). Practical bayesian model evaluation using leave-one-out cross-validation and waic. *Statistics and Computing*, 27(5):1413–1432.
- [199] Veljkovic, V., Cosic, I., Lalovic, D., et al. (1985). Is it possible to analyze dna and protein sequences by the methods of digital signal processing? *IEEE Transactions on Biomedical Engineering*, (5):337–341.
- [200] Vijayakumar, M., Wong, K.-Y., Schreiber, G., Fersht, A. R., Szabo, A., and Zhou, H.-X. (1998). Electrostatic enhancement of diffusion-controlled protein-protein association: comparison of theory and experiment on barnase and barstar. *Journal of molecular biology*, 278(5):1015–1024.
- [201] Vishveshwara, S., Brinda, K., and Kannan, N. (2002). Protein structure: insights from graph theory. *Journal of Theoretical and Computational Chemistry*, 1(01):187–211.
- [202] Wainreb, G., Ashkenazy, H., Wolf, L., Ben-Tal, N., and Dehouck, Y. (2011). Protein stability: a single recorded mutation aids in predicting the effects of other mutations in the same amino acid site. *Bioinformatics*, 27(23):3286–3292.
- [203] Walt, S. v. d., Colbert, S. C., and Varoquaux, G. (2011). The numpy array: A structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2):22–30.
- [204] Wang, G., Zhao, Y., and Wang, D. (2008). A protein secondary structure prediction framework based on the extreme learning machine. *Neurocomputing*, 72(1):262 – 268. Machine Learning for Signal Processing (MLSP 2006) / Life System Modelling, Simulation, and Bio-inspired Computing (LSMS 2007).
- [205] Wegner, P. (1990). Concepts and paradigms of object-oriented programming. *ACM Sigplan Ops Messenger*, 1(1):7–87.
- [206] Welch, P. (1967). The use of fast fourier transform for the estimation of power spectra: a method based on time averaging over short, modified periodograms. *IEEE Transactions on audio and electroacoustics*, 15(2):70–73.
- [207] Wold, S., Esbensen, K., and Geladi, P. (1987). Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52.

- [208] Wu, T.-F., Lin, C.-J., and Weng, R. C. (2004). Probability estimates for multi-class classification by pairwise coupling. *Journal of Machine Learning Research*, 5(Aug):975–1005.
- [209] Wu, Z., Kan, S., Lewis, R. D., Wittmann, B. J., and Arnold, F. H. (2019). Machine-learning-assisted directed protein evolution with combinatorial libraries. *arXiv preprint arXiv:1902.07231*.
- [210] Yang, H., Parthasarathy, S., and Mehta, S. (2005). A generalized framework for mining spatio-temporal patterns in scientific data. In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, KDD '05*, pages 716–721, New York, NY, USA. ACM.
- [211] Yang, K. K., Wu, Z., and Arnold, F. H. (2018). Machine learning in protein engineering. *arXiv preprint arXiv:1811.10775*.
- [212] Yeung, K. Y. and Ruzzo, W. L. (2001). Details of the adjusted rand index and clustering algorithms, supplement to the paper an empirical study on principal component analysis for clustering gene expression data. *Bioinformatics*, 17(9):763–774.
- [213] Yoo, I., Alafaireet, P., Marinov, M., Pena-Hernandez, K., Gopidi, R., Chang, J.-F., and Hua, L. (2012). Data mining in healthcare and biomedicine: A survey of the literature. *Journal of Medical Systems*, 36(4):2431–2448.
- [214] Zhang, H. (2004). The optimality of naive bayes. *AA*, 1(2):3.
- [215] Zhang, T., Ramakrishnan, R., and Livny, M. (1997). Birch: A new data clustering algorithm and its applications. *Data Mining and Knowledge Discovery*, 1(2):141–182.
- [216] Zhou, H. and Zhou, Y. (2004). Quantifying the effect of burial of amino acid residues on protein stability. *Proteins: Structure, Function, and Bioinformatics*, 54(2):315–322.
- [217] Zivkovic, Z. et al. (2004). Improved adaptive gaussian mixture model for background subtraction. In *ICPR (2)*, pages 28–31.
- [218] Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, 67(2):301–320.

