

Fundamentals of Bioinformatics

Project Manual 2024: Mutation Impact Prediction Methods

Kunal Chaudhary, Dominika Martinovicova, Alexandre Baudry, Ruhi
Soni, Arianna Kazemi, Liliana Kadar, Alice Peng, Marina
Diachenko, Roel van der Ploeg, Lucía Barbadilla Martínez, Lara
Pozza, Will Harley, Fabienne Kick, Ren Xie, Alex van Kaam, Ignas
Krikštaponis, Arthur Goetzee, Daniël Muysken, Sanne Abeln, Anton
Feenstra and Josemari Urtasun

Table of Contents

Table of Contents	2
Introduction	3
Aims of the Group Project	3
Predicting the Impact of Mutations	3
Timeline	4
Grading	5
Practical Instructions and Questions	5
Setup	5
Logging Into the VU Servers	5
Setting up a local database	5
Benchmarking Impact Prediction Methods	6
The Benchmark Datasets	8
The HGVS Format	9
Ensembl Variant Effect Predictor (VEP)	9
Create Baseline Prediction	11
ROC Plot	11
Calculating the AUC for a ROC Curve	12
Comparing Your ROC Curve with Other Data Sets	13
Instructions for Submitting Draft Report (Submit in PDF format via Canvas)	15
Peer Review of Draft Reports	15
Use Cases: Investigating two SNPs in detail	16
Discussion Sessions	18
Sharing your data	18
Discussion Questions	18
Discussion Within Your Group:	18
Discussion with Other Groups:	18
Format of the Final Report	19
Handing in the Final Report	20
Handing in the Final Code	20
CodeGrade	20
References	21
Report grading rubrics	22

Introduction

Aims of the Group Project

This project is an introduction to the basic theory and practice of solving common problems in bioinformatics. Bioinformatics is an interdisciplinary field combining both biology and computer science, and depending on your academic background, some parts of this project may be unfamiliar and challenging to you. However, we aim to make project groups that will include students from different Bioscience and computer backgrounds. You should allocate tasks accordingly within your group but also work collaboratively as much as possible. The objective of this project is to learn to communicate scientific problems with people who speak a different scientific language, as this will be an essential skill working in the field of bioinformatics. In addition, the project allows you to see how far your current knowledge reaches and find which skills you will have to improve in the coming year. Courses scheduled later in the curriculum will delve deeper into the details of the tools and data.

Predicting the Impact of Mutations

Nonsynonymous (missense) mutations occur where a single nucleotide in a DNA codon is substituted for another (a single nucleotide polymorphism or SNP), resulting in a change to the amino acid that the codon codes for. These missense mutations can have no or little effect on protein function and phenotype, but can sometimes result in significant changes that can cause disease. The impact prediction tools PolyPhen-2 (Adzhubei *et al.*, 2010) and SIFT (Vaseret *et al.*, 2016) are designed to predict which mutations in DNA will cause changes in the cell. They can be used to help interpret mutation data from patients who have genetic diseases, but these methods must be validated to assess how accurate their predictions are. Validation of these tools requires experimental data with accurate annotations (i.e., a benchmark or gold standard dataset) against which we can compare the performance of our tools. In this case, this means SNP data with annotations to indicate whether they are benign or pathogenic to compare against the impact predictions of SIFT and PolyPhen. We will use the database ClinVar for our gold standard dataset. Once we have benchmarked the predictions, we can visualise the performance of these tools by creating a ROC (Receiver Operating Characteristic) plot, which plots the True Positive Rate (TPR) against the False Positive Rate (FPR). More will be explained about ClinVar and ROC plots in depth in the step-by-step instructions below.

Timeline

Week	Date	Activity/Deadline	Exercises (Green boxes)
1	Monday	Intake Test Linux and Command Line Introduction	
	Tuesday	Impact Prediction Tutorial and Report Writing Introduction	Exercises 1-4
2	Monday	Script Baseline + report writing intro + define research question	Exercises 1-4
	Tuesday	Script Baseline + report writing methods	Exercise 5
3	Monday	Script ROC plot + report writing results	Exercises 6-9
	Tuesday	Script ROC plot + report writing results + abstract	Exercises 6-9
4	Monday	Use cases	Exercise 10-11
	Tuesday	Use cases	Exercise 10-11
5	Monday	Draft deadline	
	Tuesday	Work on peer review	
6	Monday	Peer review deadline	
	Tuesday	FoB Exam	
7	Monday	Presentation preparation	
	Tuesday	Presentation & discussion sessions	
8	Monday	Deadline Final Report	
	Tuesday		

Grading

This project counts for 40% of the final course grade. The deliverables that are listed below all have to be handed in and are either graded or pass/fail. You will be assigned to give feedback on the draft report of another group and your group's draft report will get feedback from members from other groups. During the final weeks of the course, you will compare your results to that of other groups. This will help you write the discussion for your report. At the end of this manual, each of these deliverables is explained in greater detail. Note that the group project will also help you prepare for the exam.

- Draft report (pass/fail) + ROC data files (pass/fail)
- Peer Feedback of draft report (pass/fail)
- Final report (100% of final project grade)
 - Based on rubric

> Exercise X | Blue Boxes

Through this manual you will find green boxes like this one. Green boxes contain exercises that will help you understand the project. You should discuss them with the Teacher's Assistants (TAs), but you are **NOT** expected to answer them in the report directly.

Practical Instructions and Questions

Setup

Logging Into the VU Servers

To make sure you can run in a linux environment, you first need to log in to the VU servers, so that you can run *python3* in a suitable environment. For instructions on how to create a suitable environment for your OS, check tutorial 0 under the Technical setup header under modules on Canvas, next to the section Programming Class. Additionally, you can install a local editor, such as Spyder. Note that it may be wise for all group members to use the same editor, to avoid issues with tab and space settings in python.

Also if you decide to share your code through a GitHub repository, remember to make the repository private, **a public repository would be seen as enabling plagiarism.**

Setting up a local database

Before we can start we have to create directories to store our data and outputs. You can again create these directories from the command line. We have provided you with five data files and three skeleton scripts.

1. Use the command line (in the terminal) with *mkdir* (see programming tutorial 1 on Canvas) to create 2 subdirectories in your working directory: *data* and *output*

2. Download the following data files to your work directory from FoB Project in Canvas:

- (1) BLOSUM62.txt
- (2) <HGVSDataset>_benchmark.tsv
- (3) <HGVSDataset>_sift_scores.tsv
- (4) <HGVSDataset>_polyphen_scores.tsv
- (5) <HGVSDataset>_VEP_baseline.tsv

3. Place BLOSUM62.txt file and the <HGVSDataset>_benchmark.tsv in the *data* subdirectory

4. Create a subdirectory named *vep* in *data*, i.e. *data/vep*. VEP stands for Ensemble Variant Effect Predictor, and place <HGVSDataset>_sift_scores.tsv, <HGVSDataset>_polyphen_scores.tsv and <HGVSDataset>_VEP_baseline.tsv in *data/vep*

5. Download the skeleton scripts (.py files) from Canvas to your work directory

Remark: if you are going to be working on the compute server, you first need to copy your files there. Please see the tutorial working at home or on your own laptop to see how you can move/copy files to the compute servers.

Your work directory structure shall look like (add a flow chart or a tree chart) # need to be done and could be used for students as a checkpoint?

Benchmarking Impact Prediction Methods

In this project you will test the performance of the mutation impact prediction tools PolyPhen and SIFT using a benchmark (gold-standard data) from ClinVar, an NCBI database of human genomic variation and its relationship to human health. You will also create a baseline model and compare this to the ClinVar benchmark. You will do this using the BLOSUM62 matrix, which is an amino acid substitution scoring matrix used for protein sequence alignment tools such as BLASTP. This is to build a basic prediction model against which to compare the performance of PolyPhen and SIFT as measured using the benchmark. To visualise the results, you will create ROC plots for how SIFT, PolyPhen and your baseline model compare with the data from ClinVar. You will create three individual ROC plots and one plot of all three models together. You will then combine your results with those from other students to compare your data with those of your peers.

An overview of the workflow of this project is shown in *Figure 1*.

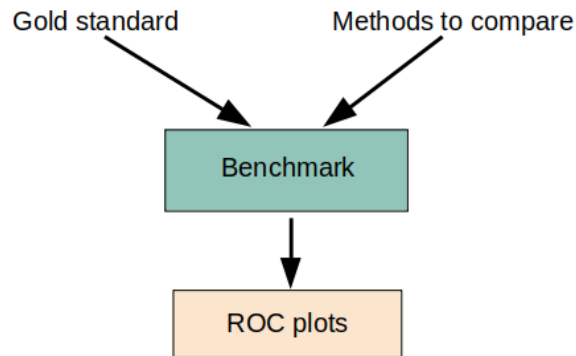


Figure 1. Project workflow. The workflow shows a simple overview of the project. A benchmark will be performed between the gold standard, VEP, and the methods to compare, baseline, SIFT and PolyPhen. The scores from the benchmark can then be used to produce a ROC plot for each predictor.

You will be given four initial *.tsv* files (tsv stands for tab separated values), one text file and three skeleton scripts.

- Two of these *.tsv* files contain the results of PolyPhen-2 and SIFT, a third *.tsv* contains the information of the benchmark, and the fourth *.tsv* file contains the information that you will need to create your baseline model.
- The *.txt* file contains the BLOSUM62 matrix that you will need for the baseline model.
- The skeleton scripts have missing blocks of code that you will have to complete, you can find them between the “START CODING HERE” and “END CODING HERE”.
 - The first script you will need to complete and run is the *skeleton_script_baseline_model.py*, to create a baseline model.
 - The *skeleton_script_create_roc_plot.py* uses the output of the baseline model along with the three other *.tsv* files as its input to create ROC plots that compare the baseline model, SIFT and PolyPhen with the data from ClinVar.
 - Finally, you will run the *skeleton_script_roc_plot_tsv.py* script which needs the output.tsv files generated with the previous script from your data and from your fellow students data, which will be provided to you three weeks into the project.

Details of these steps can be found in the following sections. Before starting, remember that you cannot import any packages apart from the ones already found in the skeleton scripts. That would make the code harder to read for those who have just started programming, and will most likely be graded as a fail. Note that even though some students in your group can focus on the programming, everyone in the group should be able to run the scripts and understand what they aim to do (also to prepare for the exam). We recommend you start reading the script in the *main()* function, and then try to read and understand each function as they are being called in the main function. To make the code understandable for everyone, make sure to comment what you are doing with the lines of code you add.

The Benchmark Datasets

[ClinVar](#) is a database of how human genomic variants are related to phenotypes of human disease and supporting clinical evidence for these relationships, managed by the NCBI. This database has been used to obtain the HGVS IDs of the genomic variants and their clinical significance (label) you will work with (**<HGVSdataset>_benchmark.tsv**). Each group will work with one of the three datasets:

1. old version of the dataset (*HGVS_2014_<...>.tsv*)
2. short dataset of the up-to-date database (*HGVS_2020_small_<...>.tsv*)
3. long dataset of the up-to-date database (*HGVS_2020_big_<...>.tsv*).

Each of these have been mapped to the reference genome GRCh38. The up-to-date database (*HGVS_2020_small_benchmark.tsv* and *HGVS_2020_big_benchmark.tsv*) was obtained from [clinvar_20200629.vcf.gz](#) while the old database (*HGVS_2014_benchmark.tsv*) was obtained from [clinvar_20141202.vcf.gz](#). A selection process was done for all of them: only benign or pathogenic SNPs (Single Nucleotide Polymorphisms) were selected (likely benign or likely pathogenic SNPs were excluded), to overcome ambiguities. Subsequently, the following SNPs were filtered out:

- Intron variants
- Synonymous variants (those that lead to synonymous mutations)
- Variants in mitochondrial DNA
- Variants that vary into multiple bases
- Unknown variants.

The three datasets were balanced obtaining the same number of 'Benign' and 'Pathogenic' samples.

> Exercise 1 | Inspection of the data

- How many HGVS IDs does each dataset have?
- Can you check that the benchmark dataset is indeed balanced? Why is this important?

> Exercise 2 | Search for your own SNP

In this exercise, we are interested in modifying the gene coding for Apolipoprotein E (APOE) to reverse a missense mutation from a patient with APOE deficiency. We do not know how many missense mutations can be pathogenic. Go to <https://www.ncbi.nlm.nih.gov/clinvar> and search for APOE. From how many missense variants can this deficiency originate?

The HGVS Format

The [Human Genome Variation Society \(HGVS\)](#) nomenclature is used worldwide as a standard language for the description of changes (variants/polymorphisms/mutations) in RNA and DNA sequences. It is formatted as **reference : description** – the reference sequence (e.g., NM_004006.2) is how the variant is referenced in databases, in this case RefSeq, and it is followed by a description of the variant (e.g., c.4375C>T). These descriptions are usually given in the context of a specific gene. The first (lowercase) letter stands for the context of the code: c for coding DNA, g for genomic DNA, r for RNA and p for protein. The number is the position of the polymorphism in the reference sequence (e.g., 4375) and the last two letters separated by a > symbol represent the two different nucleotides that are found in this position.

> Exercise 3 | HGVS format

- What is the meaning of the highlighted symbols in the HGVS ID NM_004006.2:c.4375**C>T**?
- Your HGVS IDs will be genomic reference sequences based on a chromosome, which values should the highlighted Xs take?
X_000003.12:**X**.12599717C>G

Ensembl Variant Effect Predictor (VEP)

The [Ensembl Variant Effect Predictor \(VEP\)](#) is a tool for the analysis and annotation of genomic variants in coding and non-coding DNA. VEP can take different genomic variant formats as input, but here we will use the HGVS format. It works using an extensive collection of genomic annotation and can be adapted to different interfaces depending on the context of the project: it can be used through a web interface, a command line tool and REST API. The web interface can be used for smaller amounts of data while the command line tool is able to handle larger amounts of data and has more flexibility and a greater range of options.

In this project, you will not have to use VEP directly, as we have selected some of the results we obtained previously by running REST API.

The VEP output that we are interested in this project is:

- **ID**: Corresponds to the HGVS provided in the input.
- **Amino acids change**: Reference and variant amino acids.
- **Codon change**: Reference and variant codon sequence, the alternative codons with the variant base in upper case.
- **PolyPhen Score**: Impact prediction of an amino acid substitution produced by PolyPhen 2.2.2. The score ranges from 0.0, being tolerated, to 1.0, being deleterious.
- **SIFT score**: Impact prediction of an amino acid substitution produced by SIFT 5.2.2. The score ranges from 0.0, being deleterious, to 1.0, being tolerated.

Each HGVS ID is a variant of a genomic sequence that can overlap multiple transcripts. Hence, when using VEP each HGVS ID has as many outputs as transcripts there are, and each one can take different PolyPhen and SIFT scores. To make things easier for you, we have selected only one transcript result for each HGVS ID. The transcript with the highest impact (most deleterious) predicted by SIFT and PolyPhen was the one selected. If several transcripts have this score, the transcript is selected at random among these ones. If PolyPhen and SIFT do not agree, that HGVS ID will be skipped to avoid bias towards either PolyPhen or SIFT.

> Exercise 4 | Selecting transcripts for the VEP output

This output has been extracted from the VEP web interface when providing three HGVS IDs. If you have understood the paragraph above, you should know which transcripts have been selected in your output.

Uploaded variant	Feature	Amino acids	Codons	SIFT	PolyPhen
NC_000001.11:g.224424520A>C	ENST00000414423.8	S/R	AGT/AGG	0	0.992
NC_000001.11:g.224424520A>C	ENST00000445239.1	S/R	AGT/AGG	0	0.972
NC_000001.11:g.224424520A>C	ENST00000651911.1	S/R	AGT/AGG	0	0.972
NC_000018.10:g.51078285G>C	ENST00000342988.8	D/H	GAT/CAT	0	0.899
NC_000018.10:g.51078285G>C	ENST00000398417.6	D/H	GAT/CAT	0	0.899
NC_000018.10:g.51078285G>C	ENST00000588745.5	D/H	GAT/CAT	0	0.629
NC_000001.11:g.236885200A>G	ENST00000366576.3	D/G	GAC/GGC	0.2	0.21
NC_000001.11:g.236885200A>G	ENST00000535889.6	D/G	GAC/GGC	0.23	0.21
NC_000001.11:g.236885200A>G	ENST00000674797.1	D/G	GAC/GGC	0.18	0
NC_000001.11:g.236885200A>G	ENST00000366577.10	D/G	GAC/GGC	0.15	0.161

The output is given to you in three separate .tsv files (“<...>” refers to your assigned dataset, as described above in The Benchmark Datasets):

- **<HGVSdataset>_sift_scores.tsv**: contains the HGVS IDs* and the SIFT score.
- **<HGVSdataset>_polyphen_scores.tsv**: contains the HGVS IDs* and PolyPhen score.
- **<HGVSdataset>_VEP_baseline.tsv**: contains the HGVS IDs*, Amino acid change and Codon change.

**Note that the HGVS IDs are the same for the 3 files.*

Create Baseline Prediction

The BLOSUM62 matrix is based on frequencies of amino acid substitutions of a collection of protein alignments with 62% identity. As you know, this matrix is being used in

alignment tools such as BLAST or BLASTP. In this project, you will use the information in BLOSUM62 to obtain an insight into a substitution's expected impact, and with this create a baseline impact prediction method. You can find the BLOSUM62 matrix in the *BLOSUM62.txt* file.

> Exercise 5 | BLOSUM62 matrix

Check out the BLOSUM62 matrix in the *BLOSUM62.txt* file. Do you think the diagonal values are going to be used on the baseline model with the data set that we have provided? Why, or why not?

> Exercise 6 | BLOSUM62 matrix

Check the BLOSUM62 matrix in *BLOSUM62.txt* again. Look at the substitution scores of cystine (C) and glutamine (Q). Can you think of reasons why glutamine seems to be more replaceable than cystine? Can you generalise your answer to the different groups of amino acids?

Complete and execute the baseline predictor skeleton script (*skeleton_script_baseline_model.py*) using *BLOSUM62.txt* and *<HGVSdataset>_VEP_baseline.tsv* as inputs. The script should create a score which is simply the raw value of the BLOSUM62 for that amino acid exchange. As output, you should obtain the scores of your baseline in the same format as in *<HGVSdataset>_sift_scores.tsv* and *<HGVSdataset>_polyphen_scores.tsv*. The output can be saved to a *data* or *output* folder, or any other folder that you might have additionally created in your working directory beforehand, by providing a file path to the *-o* argument on the command line. This argument is required, and the file name should be supplied with the *.tsv* extension; for example:

```
$ python3 skeleton_script_baseline_model.py
data/vep/HGVS_2020_small_VEP_baseline.tsv data/BLOSUM62.txt -o
data/HGVS_2020_small_baseline_scores.tsv
```

ROC Plot

Your task here is to create a Receiver Operating Characteristic (ROC) plot by comparing the results from your predictors to the gold standard data we have obtained from ClinVar. A ROC plot is a method of visualising the performance of your predictor, it plots the True Positive Rate (TPR) against the False Positive Rate (FPR). Refer to the lecture on machine learning and benchmarking for a thorough explanation of what ROC plots are and

how they can be used to evaluate, compare, and refine classification methods. In addition http://wikipedia.org/wiki/Receiver_operating_characteristic may be a helpful resource.

Note that a threshold to classify variants as (putatively) benign or damaging is not fixed at a constant value, in order to create a ROC plot of the results. Instead, in a ROC plot you calculate the true and false positive rate for every possible threshold spanning the range of possible values for your method, from 0 until 1 for SIFT or from -2 until 9. For every threshold, this allows every variant classified by the predictor to be categorised as a True Positive (TP), False Positive (FP), False Negative (FN) or a True Negative (TN).

> Exercise 7 | Confusion Matrix

Complete the blank cells in this confusion matrix. Hint: The conclusion drawn from the predictor depends on the threshold.

PREDICTOR	BENCHMARK		<i>ClinVar Benign</i>
	Conclusions		<i>Confirmed Benign</i>
	<i>(Putative Damaging)</i>		
	<i>(Putative Benign)</i>		True Negative (TN)

Calculating the AUC for a ROC Curve

A ROC plot can be made by varying the threshold, counting the TPs, FPs, TNs and FNs, and calculating the TPR and FPR. First, it is important to define what positive and negative assignments are. Think of a covid test, a positive result means you probably have the virus. Here the same consensus applies, thus a mutation with a predicted deleterious effect will be a positive result.

When working with ROC plots, the Area Under the Curve (AUC) is often taken as a measure to evaluate performance. To calculate the AUC you have to approximate the integral of the function $f(x)$ that describes the shape of the curve of the ROC-plot. We do not know the function that describes the curve, thus we have to evaluate the integral numerically. A method that approximates the integral is the trapezoidal rule.

Think of a clever way to implement this rule, and complete the provided skeleton script: *skeleton_script_create_roc_plot.py*. This script will parse your predictor and benchmark results, count the number of TP, FP, FN and TN, calculate your ROC plot's line coordinates, create the corresponding figure, and integrate the AUC.

Complete and execute the skeleton script *skeleton_script_create_roc_plot.py*. For a better understanding of the ROC plot, the script produces a color gradient indicating the score range.

To obtain the individual ROC plot for one predictor, the optional argument *-ipred* should be included once with the *.tsv* file with the scores of one of your methods (SIFT, PolyPhen or baseline). The *-ibench* should be included with the *<HGVSdataset>_benchmark.tsv* file. You can use the help function *-h* or *--help* for explanation of these and other options. You will have to specify a path for the output *.png* file with the argument *-o* (including the *.png* extension). A *.tsv* file with the ROC x- and y-coordinates will be saved automatically to the same output directory. For example, to call the script for the PolyPhen ROC plot:

```
$ python3 skeleton_script_create_roc_plot.py -ibench
data/HGVS_2020_small_benchmark.tsv -ipred
data/vep/HGVS_2020_small_polyphen_scores.tsv -color -o
output/ROCplot_HGVS_2020_small_polyphen.png
```

To show the ROC curves of the three predictors in one figure, the script can be run with the *-ipred* argument three times for each of the three prediction *.tsv* files (SIFT, PolyPhen and baseline). (In this case, the ROC plot coordinates file will not be created). A command line example is provided below:

```
$ python3 skeleton_script_create_roc_plot.py -ipred
data/vep/HGVS_2020_small_polyphen_scores.tsv -ipred
data/vep/HGVS_2020_small_sift_scores.tsv -ipred
data/HGVS_2020_small_baseline_scores.tsv -ibench
data/HGVS_2020_small_benchmark.tsv -o output/ROCplot_all.png
```

> Exercise 8 | Curve details

- In your ROC plots you will probably find that some of the methods provide a much more detailed curve than others. Can you explain why this is?
- What will a ROC-plot look like if you have extremely unbalanced data? Will the ROC-plot be representative?

Comparing Your ROC Curve with Other Data Sets

The performance of the predictors depends not only on the predictor itself but can also be influenced by the data you work with. When you hand in your draft report we will also ask you for the data you have used to generate the ROC plots. This data will be shared with the other groups for comparison. In order to read in this data, and make a new plot, there is a third script. To test how different data sets can influence the ROC plot results, you will run the third script *skeleton_script_roc_plot_tsv.py*.

This script contains only one coding block which is the same you have found in the last script to calculate the AUC – please use the same code. As inputs, you will use the coordinates that your fellow students have obtained with the other two datasets by providing paths to -itsv. These two sets of coordinates will be given to you through Canvas. You will have to run the code twice, once for each type of dataset. As output, you will get a .png file with the ROC plot comparing the same dataset on the three predictors, just as the previous script (*skeleton_script_create_roc_plot.py*). This command line illustrates how can it be run:

```
$ python3 skeleton_script_roc_plot_tsv.py -itsv  
output/ROCplot_HGVS_2020_small_sift_xy.tsv -itsv  
output/ROCplot_HGVS_2020_small_polyphen_xy.tsv -itsv  
output/ROCplot_HGVS_2020_small_baseline_xy.tsv -o output/ROCplot_comparison.png
```

The ROC plots from this script will be discussed in the discussion session and you will have to add them to your final report as well.

> Exercise 9 | Default thresholds

SIFT and PolyPhen define default thresholds for their score to classify a mutation as benign or pathogenic. Do you think the default thresholds make sense according to your ROC plot (look at the FPR and TPR)? What would happen if you changed the threshold?

Instructions for Submitting Draft Report (Submit in PDF format via Canvas)

The draft report must contain between 1000 and 1500 words and contain following sections:

- Abstract
- Introduction
- Methods; and
- (Preliminary) results.

The results section should include a ROC plot and its interpretation. You must clearly state your research question in the introduction and answer it in the results and discussion sections. Note that in the section below, and in the rubric (you will receive for the peer review), more details are provided about what the report should contain.

Please add word counts in square brackets [] behind the title of each section, before you submit. Your draft report needs to be handed in via Canvas, and will be peer reviewed by students of other groups. Note that you do not yet need to write the discussion and conclusion sections for your draft.

Peer Review of Draft Reports

Everyone should peer review the report of one other group, meaning that each group should get around 4 peer reviews back. The peer review should be handed in on Canvas. Note that the peer review should be based on the rubrics provided. You need to write a peer review in order to pass the course.

Use Cases: Investigating two SNPs in detail

Now we would like you to think more about the biological aspect of impact prediction. You will do this by comparing two SNPs from the same gene, where one is known to be a benign SNP and one is known to be a pathogenic SNP. You will report on your findings in a section called “**Use cases**”. The section needs to contain the answers to the exercises below and you may add additional information to support your findings. Also see the rubrics.

Below a list with three genes is shown, with corresponding SNPs and the variant. Depending on your dataset used you will do the following steps for two SNPs from a single gene.

Gene	HGVS	Feature (or transcript_id in rest API without the . and last number)	Group
TP53	NC_000017.11:g.7674220C>A	ENST00000413465.6	Old dataset
	NC_000017.11:g.7673751C>T	ENST00000269305.9	Old dataset
BRCA2	NC_000013.11:g.32362595G>C	ENST00000380152.8	Small dataset
	NC_000013.11:g.32396905A>G	ENST00000380152.8	Small dataset
BRCA1	NC_000017.11:g.43067628G>A	ENST00000478531.5	Big dataset
	NC_000017.11:g.43063368T>C	ENST00000586385.5	Big dataset

As a first step, go to the [Ensembl Variant Effect Predictor \(VEP\)](#) website. Click on the Web interface option. In your input data, paste each of your SNPs in HGVS format on a separate line. You can leave all other settings on default, and click “Run” at the bottom. Please be aware that a job can take a few minutes to complete. Click on “view results” when your job is complete. You might need to change the shown columns by clicking on the “Show/hide columns” button in the blue bar.

Alternatively you can use the REST API you can type the following url in your browser “<https://rest.ensembl.org//vep/human/hgvs/>” followed by the HGVS code. See also: https://rest.ensembl.org/documentation/info/vep_hgvs_ge

We will start with comparing the sequence conservation between the two SNPs. The [web server of PolyPhen-2](#) is the tool we will use for this. Before going to the website, find the rsID in the VEP output under “Existing variant” for each SNP (be sure to look for the right feature as indicated in the table above). Use the rsID one by one as input for the query WHESS.db. If you get multiple options in the results screen, choose the results with the same protein position as in the VEP output. On the report page you can find the sequence conservation under the Multiple sequence analysis tab.

> Exercise 10 | Sequence conservation

- Do you see differences in sequence conservation for the region around the SNP and for the SNP itself?
- Do you see differences in sequence conservation between the SNPs? Is this as expected when considering evolution laws?

On the same report page produced by PolyPhen-2, you can also find a tab called 3D visualisation. Use this to find where in the protein structure the mutation occurs.

> Exercise 11 | Protein structure

- What is the structural environment around the SNP, when focusing on where in the protein the secondary structure occurs?
- Does the place of the SNP in the protein make sense considering its impact as predicted by VEP?

Discussion Sessions

In the discussion sessions, you will be matched with students from other groups to discuss your findings about the Use Cases. The discussion points you need to prepare are the questions from exercises 10 and 11. Additionally we want you to tell something about the biological background of the SNPs, this is however not required for the report. The discussion session will be moderated by teachers and TAs.

Make sure you can show figures of the sequence conservation and the structure. This can be done in a small presentation.

Sharing your data

When you execute the `create_roc_plot.py` script, the coordinates of the ROC plot will automatically be exported to '`[your_custom_plotname]_xy.tsv`'. You need to share the `.tsv` files generated by the `skeleton_script_create_roc_plot.py` with other students for your specific benchmark datasets, for all three methods. Please follow the instructions posted on Canvas to share your data.

Discussion Questions

Your discussion section in the report should contain the answers to the following questions. The questions under discussion within your group should be based on your own results, and the questions under discussion with other groups should be based on your own results and the results from other groups.

Discussion Within Your Group:

1. Do you observe a difference in performance between SIFT, PolyPhen and the baseline script? How can you explain the difference in performance? [A1]
2. SIFT and PolyPhen define default threshold(s) for their score to classify a mutation as benign or pathogenic. Do you think the default thresholds make sense? [A2]

Discussion with Other Groups:

1. Are there any clear differences between the different benchmark datasets in terms of the AUC and the shape of the ROC curves? What is the effect of having more benchmark data available? [B1]
2. Is the relative performance the same in all three datasets for SIFT, PolyPhen and the baseline script, how could you test this? [B2]

Format of the Final Report

The final report must contain the following sections:

- Abstract (max 250 words):
 - Motivation, results & impact
 - Give an accurate and concise summary of the introduction and clearly state the research question.
 - Give an accurate and concise summary of the materials and methods section.
 - Give the most important results and answer the research question.
 - Mention the potential impact of these results on future research and/or practical applications.
- Introduction:
 - Include references to previous studies from literature related to your research question
 - Make sure to explain why impact prediction is typically performed
 - Make sure to explain why bioinformatics methods need to be benchmarked
 - Clearly state your research question, which should be accurate to the details of the project and be falsifiable by your results.
 -
- Methods and Data:
 - Give an overview of the workflow in a scheme.
 - Describe the properties of the data used.
 - Describe the scores for the different methods.
 - Describe your benchmarking strategy.
- Results:
 - Include 4 ROC plots and interpret them
 - Explain what you are trying to test and how you are testing this
 - Provide the ROC plots for each of the benchmarked methods together with the AUC
 - Compare the three different methods
 - Describe the plots and what they represent in the main text
 - Explain whether the results conform to your initial expectations.
- Use Cases:
 - Describe the sequence conservation for your SNPs and explain if you expect this according to evolution laws.
 - Describe the protein structure for your SNPs in terms of secondary structure and localisation in the protein.
 - Discuss if you agree with the prediction made by VEP, based on your results from the sequence conservation and the protein structure.
- Discussion
 - Discuss the points listed as questions in the discussion session
 - You can add results from other groups, with a reference, and/or cite other studies
 - Explain the difference in performance between the three different methods
 - Discuss the default values for SIFT and Polyphen.

- Describe any clear differences between the different benchmark datasets. Provide an explanation and give the consequences of these findings. Describe what the effect is of having more benchmark data available.
- Describe if the relative performance is the same in all three datasets for SIFT, Polyphen and the baseline script, and how you could test this.
- Conclusions
 - Give your main conclusions and answer your research question. Consider what can and can not be concluded from your results.
 - Discuss the potential impact of your results in a practical/medical context.
- Tables & Figures
 - Explain all axes, labels, lines and points in the caption of your figure/table.
 - Refer to each figure/table in the main text, and explain in the main text what can be seen from the figure/table.
- References
 - We expect between 3 and 15 citations to other papers (author-year citations are preferred). Some essential literature for the project is provided on Canvas and in the lecture slides. Note that 15 citations is not a limit.

The final report must contain between 3000 and 3500 words. We count everything (including figure text) except references.

Please read the [rubric](#) for the final report and make sure you include every requirement listed.

Handing in the Final Report

The final report must clearly state what each group member contributed to the project. Individual students' grades may be adjusted according to the reported and observed differences in workload.

CodeGrade

We use CodeGrade, an automatic grading system, to evaluate your code. This means that coding outside the code block for editing ("START CODING HERE" and "END CODING HERE") is generally not allowed (otherwise CodeGrade may not work). Importing additional packages is also not allowed.

You can submit your code to CodeGrade multiple times to check if your code works correctly. **Code will not be graded, this is just for you to check if your code works correctly.**

Final Report Grading rubric

	Not Assessable	Revision Needed	Meets Expectation	Excellent/Exemplary
Abstract	Key elements are missing, or not understandable. Conclusions are unclear or unsupported by the findings.	Basic relevant information is included, but not enough to convey the project's main purpose and findings. Requires significant revisions in content, focus and/or organisation.	Clear and easy-to-follow description of the key elements and potential impact. Could benefit from minor revisions for parts that are too vague or too detailed.	Accurate and well-organised summary of the key elements. Gives a solid understanding of the project without unnecessary information or non-trivial errors.
Introduction	Research questions are missing and/or no references are made to previous research.	The relevance of the research problem and the scientific background are mentioned but the hypothesis is not sufficiently explained. The structure of the introduction is not coherent.	The theoretical context and analysis of the problem is clearly presented. From this, the research questions are developed, and an experimental design is presented. Relevant literature is incorporated.	Thorough and creative presentation of the context and problem, with novel connections made between them. Research questions and hypotheses are developed coherently, and experimental design and expectations are presented concisely. References are high quality and well-interpreted.
Methods	Methods are missing and/or not described well. No mention/inadequate understanding of benchmarking. Workflow is incorrect or missing.	Methods are partially explained but lack critical details for understanding and replicating the procedure. Benchmarking is mentioned but not well-explained or applied. Workflow is present but missing key details.	Justifies the methodology and understands the effect of the chosen methods on the quality of data, with some minor flaws. Benchmarking is explained. Workflow is thorough with only some details missing.	Clear link presented between the used methodology and data quality, with acknowledgements of any limitations therein. Study immediately repeatable. Scoring methods and benchmarking are well explained. Workflow is descriptive and well-designed.
Results	Results are missing, unorganised, or do not align with the data presented in the tables/figures.	Results are presented but incomplete. Key findings are missing or incorrectly interpreted. Inconsistencies between the text and the tables/figures.	Essential results are present and complete, and match what the data suggests. Results make appropriate references to the tables/figures.	Results are comprehensive, thorough, and well-organised. All notable results are well-supported by the tables/figures. All tables/figures are well-integrated into the text.
Tables and Figures	Tables/figures are missing entirely, and/or missing essential aspects like labels, captions or data points.	Tables/figures are present, but the axes are illegible, have incorrect labels, and/or have missing units. Captions provide vague and/or irrelevant information. Tables/figures are poorly designed or confusing.	Tables/figures are clearly marked with appropriate labels, including units where necessary. Captions mostly describe the tables/figures well, with minimal need to read the main text.	Tables/figures are not only present with correct labels and captions, but are also well-designed and visually appealing. Axes labels, units and scales are accurate. Captions are succinct yet comprehensive enough to describe tables/figures without the need to refer to the main text.

Discussion	Discussion is incomplete, poorly organised, and/or fails to interpret results. No mentions of implications, limitations, and/or future areas of research.	Discussion does not fully connect to the research or misses the point. Arguments are sometimes flawed. Insufficient correspondence to relevant literature in the field of research. Insufficient attention for the strengths of the study and often exaggerated attention on limitations of methods.	Answers the research question while incorporating sound knowledge of the field to discuss the results, and uses relevant literature. Meaningful discussion of implications and limitations of results.	Discussion is insightful, organised, and shows a strong link between results, the research question, and the broader context. Strengths and limitations are thoroughly acknowledged. Implications and future directions are well-grounded in the results. The narrative is not just informative but also engaging and persuasive.
Conclusion	Conclusion is missing, hard to follow, irrelevant, or incorrect. Little to no discussion of implications or future directions.	Attempts to summarise main findings and implications, but is superficial and lacks completeness or clarity.	Clearly summarises key findings of the project in context of the research questions. Focuses on relevant implications and future directions from medical/practical context.	Conclusion is concise, well-organised, and provides a compelling summary of key results in context of the research questions. Future directions and implications are very clear with a compelling narrative throughout.
Use Cases	Description of sequence conservation is missing or incomplete. Connection between sequence conservation, protein structure and VEP prediction is missing, incorrect, or irrelevant.	Sequence conservation is discussed but may lack clarity, or may not be sufficiently linked to evolutionary principles. Protein structure description is present but inadequate in its description of secondary structure and/or SNP localization. Shallow or inaccurate evaluation of VEP prediction.	Clear and accurate description of sequence conservation and context of evolutionary principles. Protein structure covers both secondary structure and localization well. Well-reasoned evaluation of VEP prediction based on sequence and structural information.	Detailed description of sequence conservation and its relevance to evolutionary principles. Thorough and precise description of the protein structure, analysing both secondary structure and localization. Well-reasoned evaluation of the VEP prediction, demonstrating strong understanding and clear articulation of the role of sequence and structure on function and impact prediction.

References

Adzhubei, I. A., Schmidt, S., Peshkin, L., Ramensky, V. E., Gerasimova, A., Bork, P., Kondrashov, A. S. and Sunyaev, S. R. (2010) A method and server for predicting damaging missense mutations. *Nature methods*, **7**, 248–249

Vaser, R., Adusumalli, S., Leng, S. N., Sikic, M. and Ng, P. C. (2016) Sift missense predictions for genomes. *Nature protocols*, **11**, 1.

