# Gene Severity Analysis:
## Predicting Rater Concurrence

David Braslow
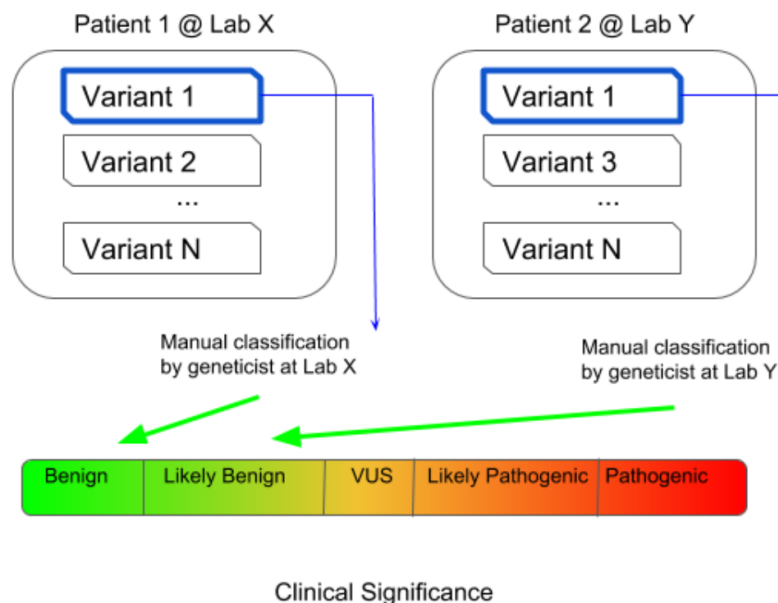
April 23, 2019

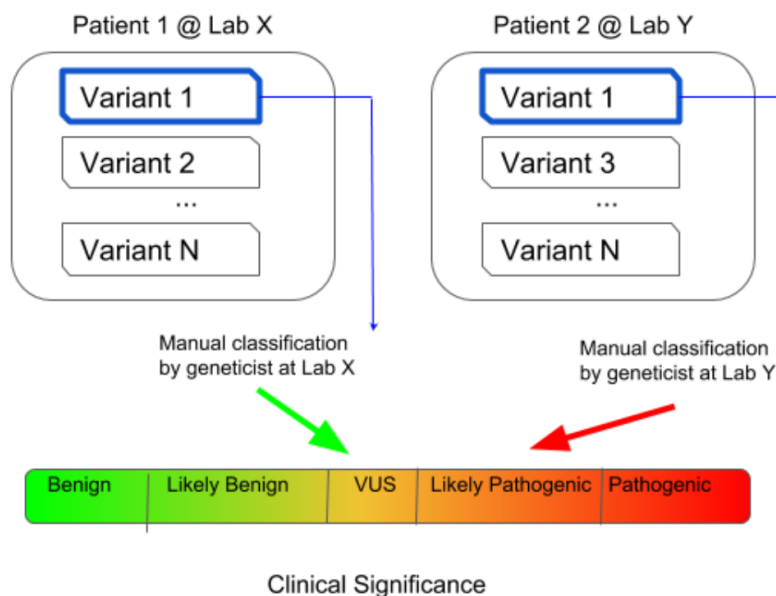# Introduction

Gene variants can have low or high clinical significance.
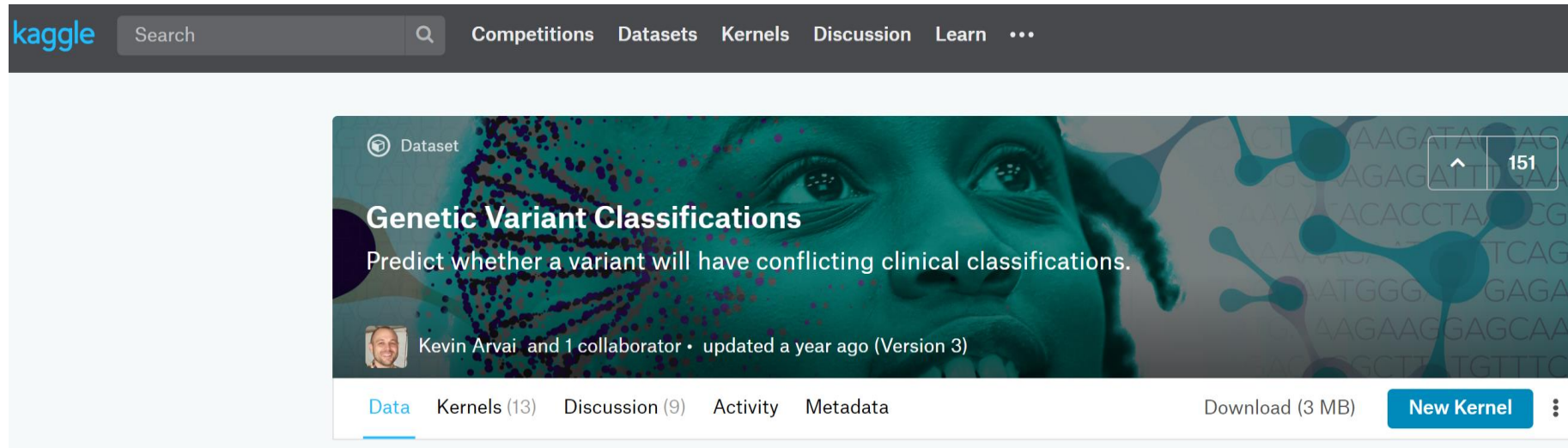
Can I predict when geneticists will disagree?

# Obtaining the Data



https://www.kaggle.com/kevinarvai/clinvar-conflicting/version/3
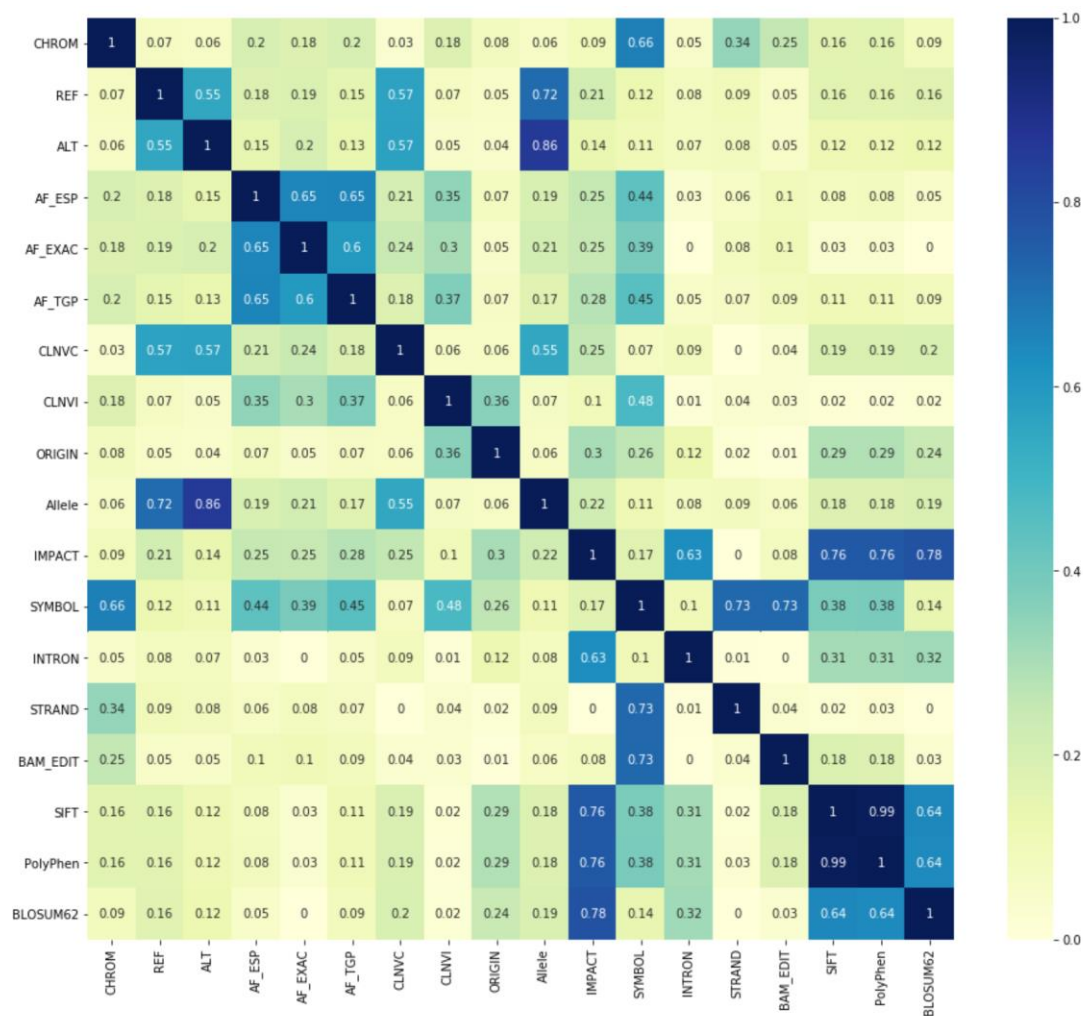
# Scrubbing the Data

45 gene properties,
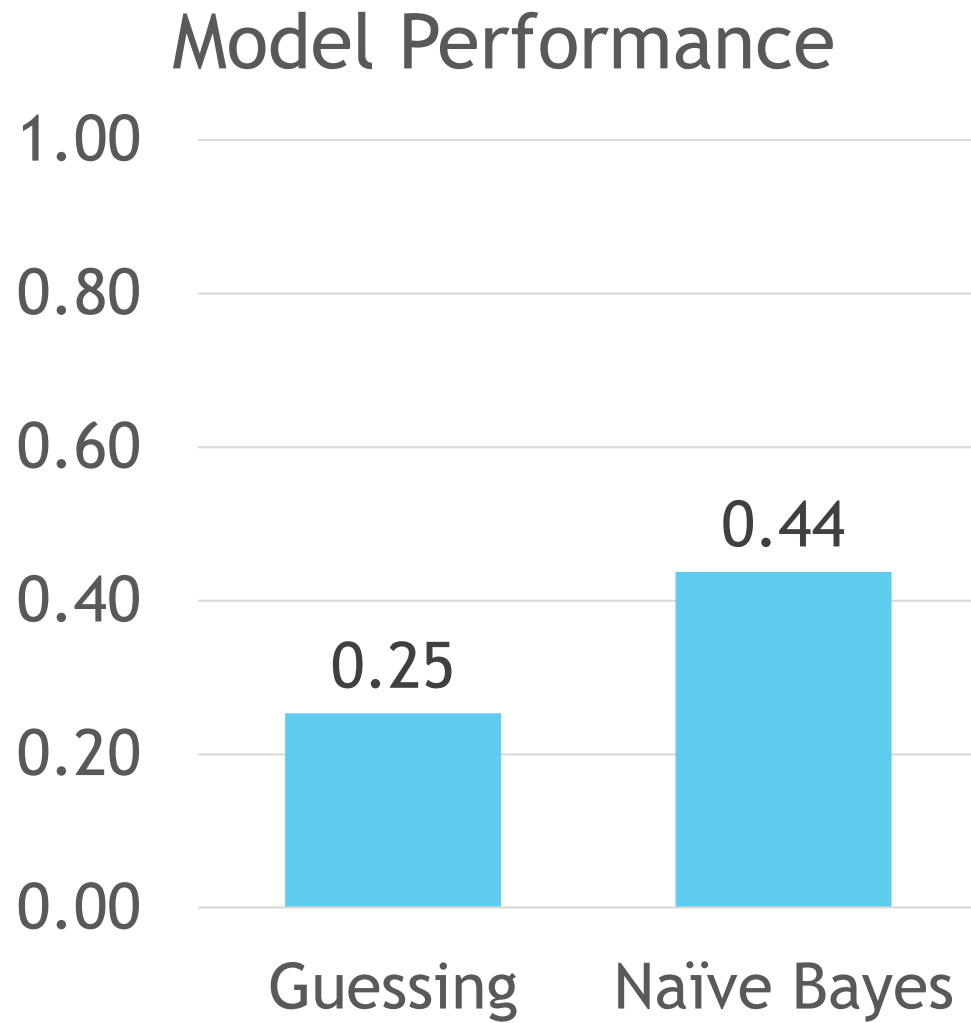most are categorical

**Columns**

\# **CHROM**   Chromosome the variant is located on

\# **POS**   Position on the chromosome the variant is located on.

A **REF**   Reference Allele

A **ALT**   Alternaete Allele

\# **AF_ESP**   Allele frequencies from GO-ESP

\# **AF_EXAC**   Allele frequencies from ExAC

\# **AF_TGP**   Allele frequencies from the 1000 genomes project

A **CLNDISDB**   Tag-value pairs of disease database name and identifier, e.g. OMIM:NNNNNN

A **CLNDISDBINCL**   For included Variant: Tag-value pairs of disease database name and identifier, e.g. OMIM:NNNNNN

A **CLNDN**   ClinVar's preferred disease name for the concept specified by disease identifiers in CLNDISDB

A **CLNDNINCL**   For included Variant : ClinVar's preferred disease name for the concept specified by disease identifiers in CLNDISDB

A **CLNHGVS**   Top-level (primary assembly, alt, or patch) HGVS expression.

A **CLNSIGINCL**   Clinical significance for a haplotype or genotype that includes this variant. Reported as pairs of VariationID:clinical significance.

A **CLNVC**   Variant Type

A **CLNVI**   the variant's clinical sources reported as tag-value pairs of database and variant identifier

A **MC**   comma separated list of molecular consequence in the form of Sequence Ontology ID|molecular_consequence

A **ORIGIN**   Allele origin. One or more of the following values may be added: 0 - unknown; 1 - germline; 2 - somatic; 4 - inherited; 8 - paternal; 16 - maternal; 32 - de-novo; 64 - biparental; 128 - uniparental; 256 - not-tested; 512 - tested-inconclusive; 1073741824 - other

A **SSR**   Variant Suspect Reason Codes. One or more of the following values may be added: 0 - unspecified, 1 - Paralog, 2 - byEST, 4 - oldAlign, 8 - Para_EST, 16 - 1kg_failed, 1024 - other

\# **CLASS**   The binary representation of the target class. 0 represents no conflicting submissions and 1 represents conflicting submissions.

A **Allele**   the variant allele used to calculate the consequence

A **Consequence**   Type of consequence:
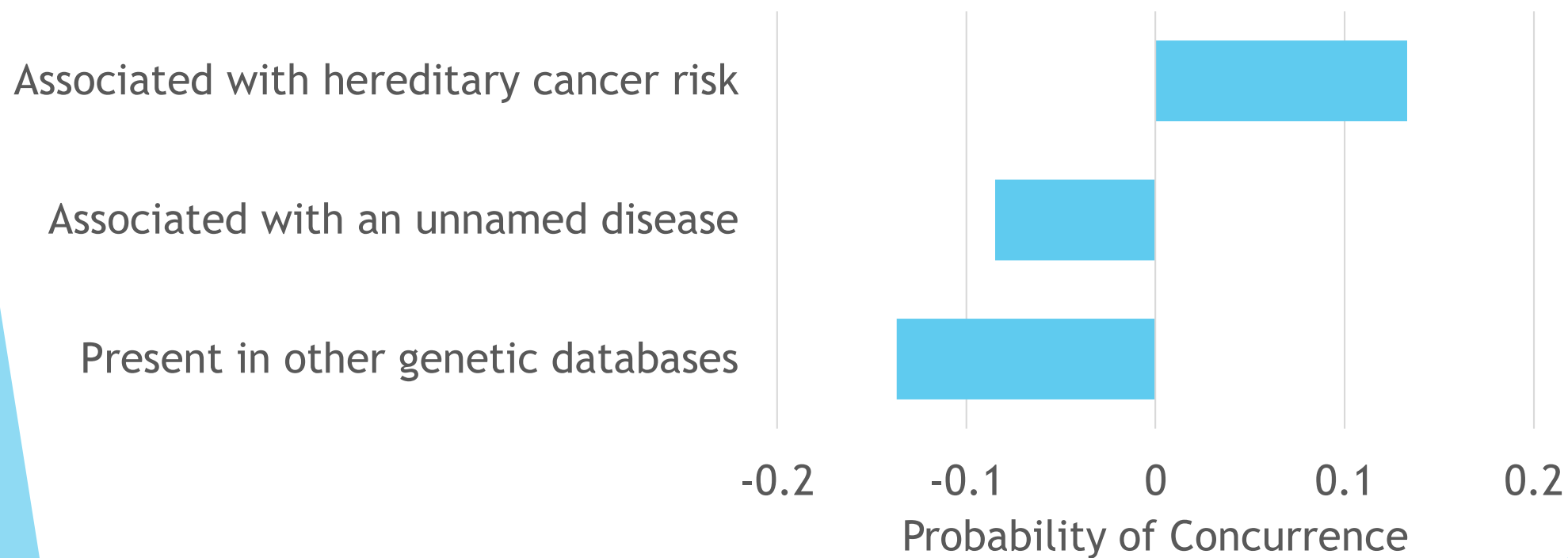
# Exploring the Data

## Feature Correlations

# Modeling the Data

# Summary

- Naïve Bayes modeling can improve prediction of conflicting ratings

- This model can help prioritize research

- There is still room for improvement with the model