

Predicting functional elements in noncoding DNA

Meer Mustafa, Kathryn Wendorf, Neelang Parghi, Alexander Lucaci

Presentation outline

Introduction

Question

Data types

Methods

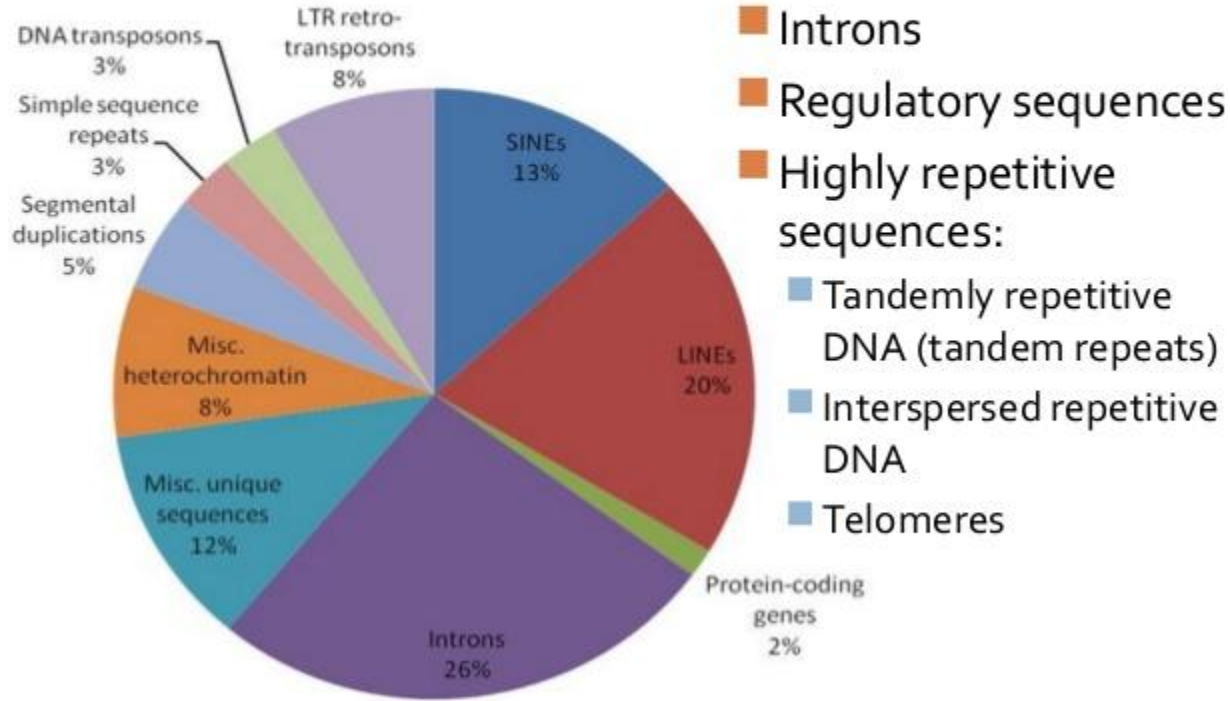
Results

Future directions

Coding vs Non-coding DNA



Classes of Noncoding Sequences

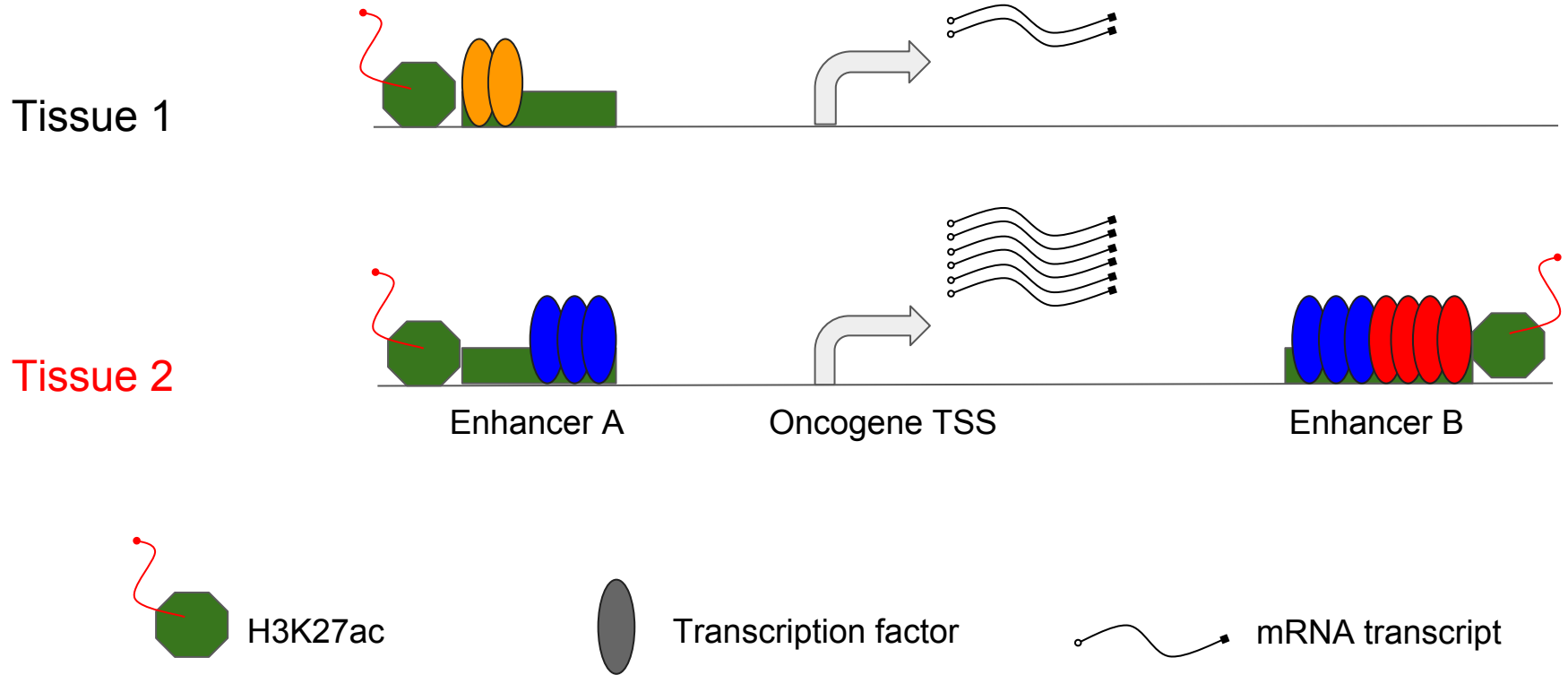


Why is the non-coding genome important?

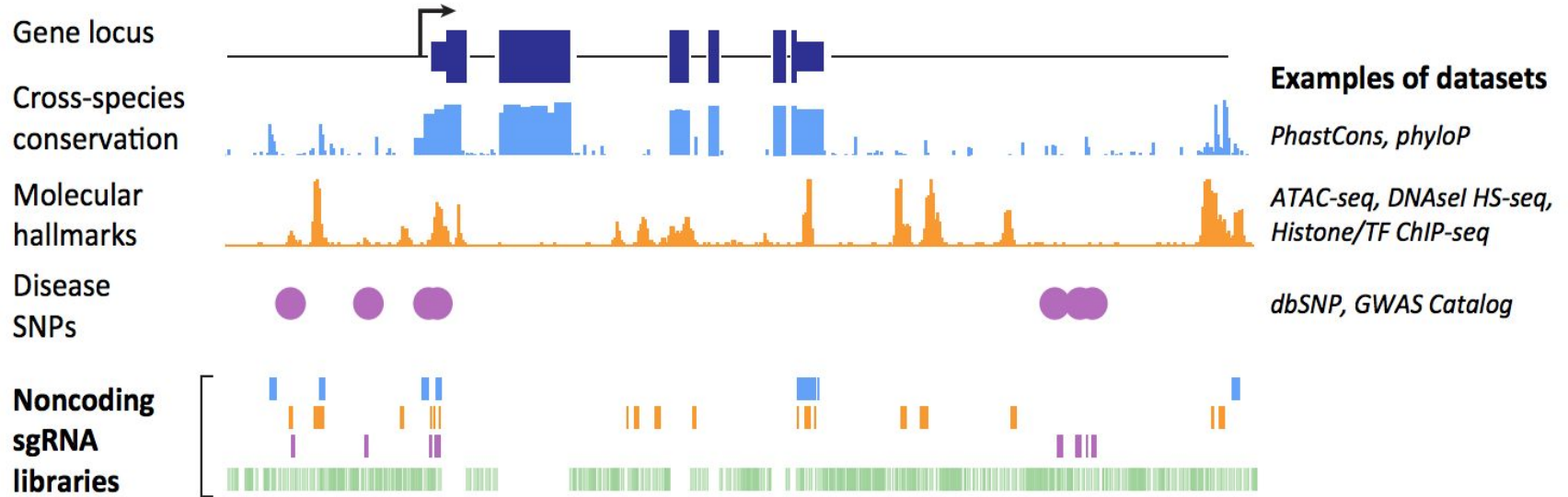
- ~98% of human genome does not code for proteins
- Non-coding genome affects gene regulation and disease
- Mapping of chromatin state and chromosome conformation has been used to identify regulatory elements.

Main problem: no overarching framework to translate the non-coding genomic sequence into functional elements.

Noncoding variation and regulation



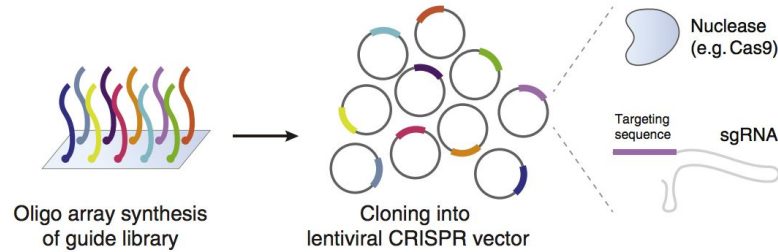
What are CRISPR screens?



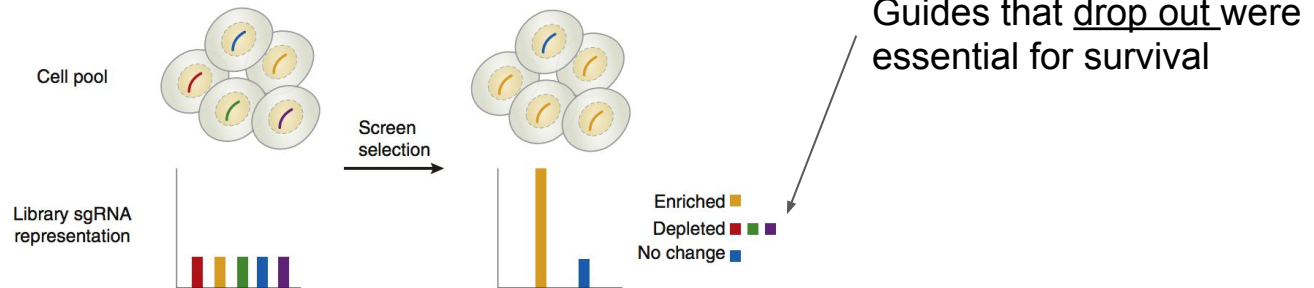
Trends in Genetics

How do they work?

A

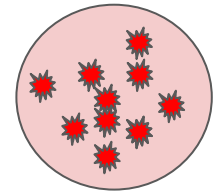
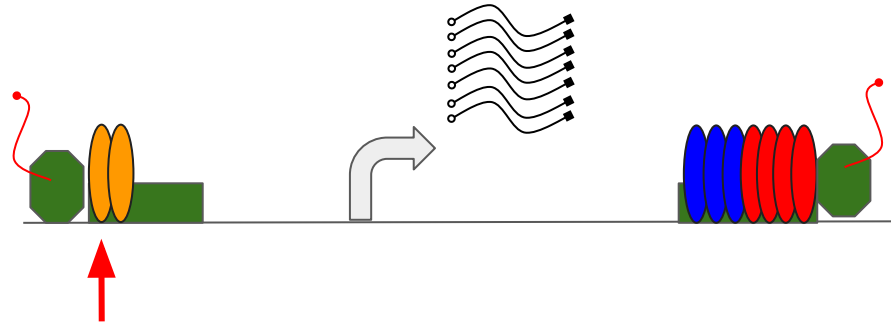


B



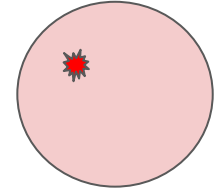
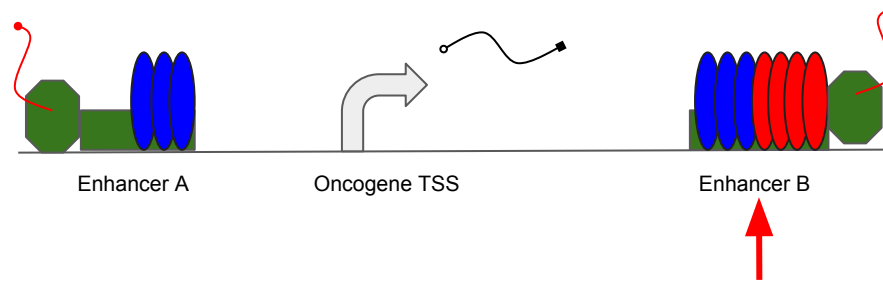
Noncoding variation and regulation

sgRNA #1



Cell growth

sgRNA #2



Cell death

Project background

Sanjana lab generated CRISPR screen looking into BRAF inhibitor resistance in melanoma.

- Demonstrated that noncoding mutations are involved in gene regulation and chemotherapeutic resistance.
- Found that noncoding loci that modulate drug resistance also harbor predictive hallmarks of functional elements (e.g. enhancers, repressors)

Current noncoding screens

Sanjana (2016)
Melanoma cancer cells

18 K guides

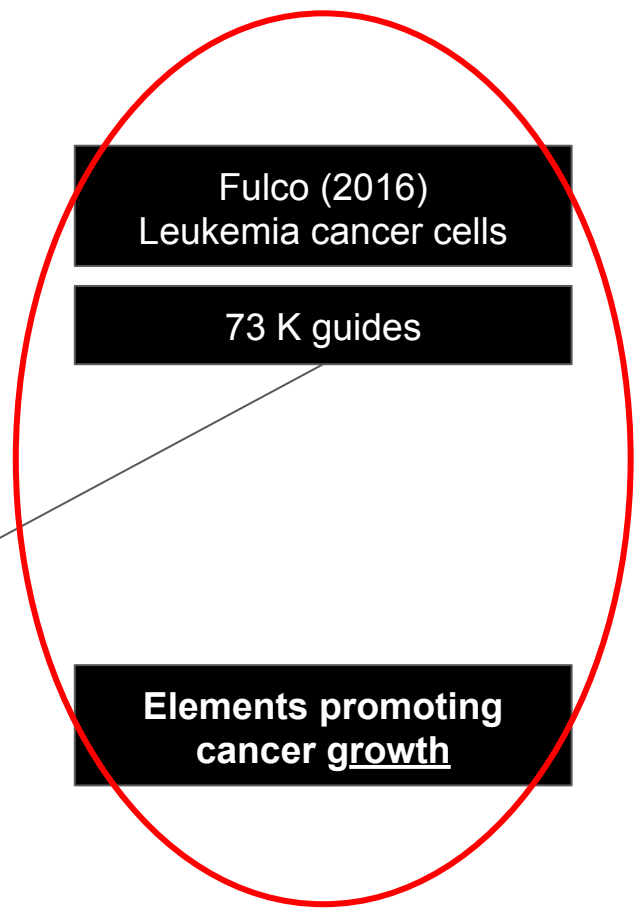
Fulco (2016)
Leukemia cancer cells

73 K guides

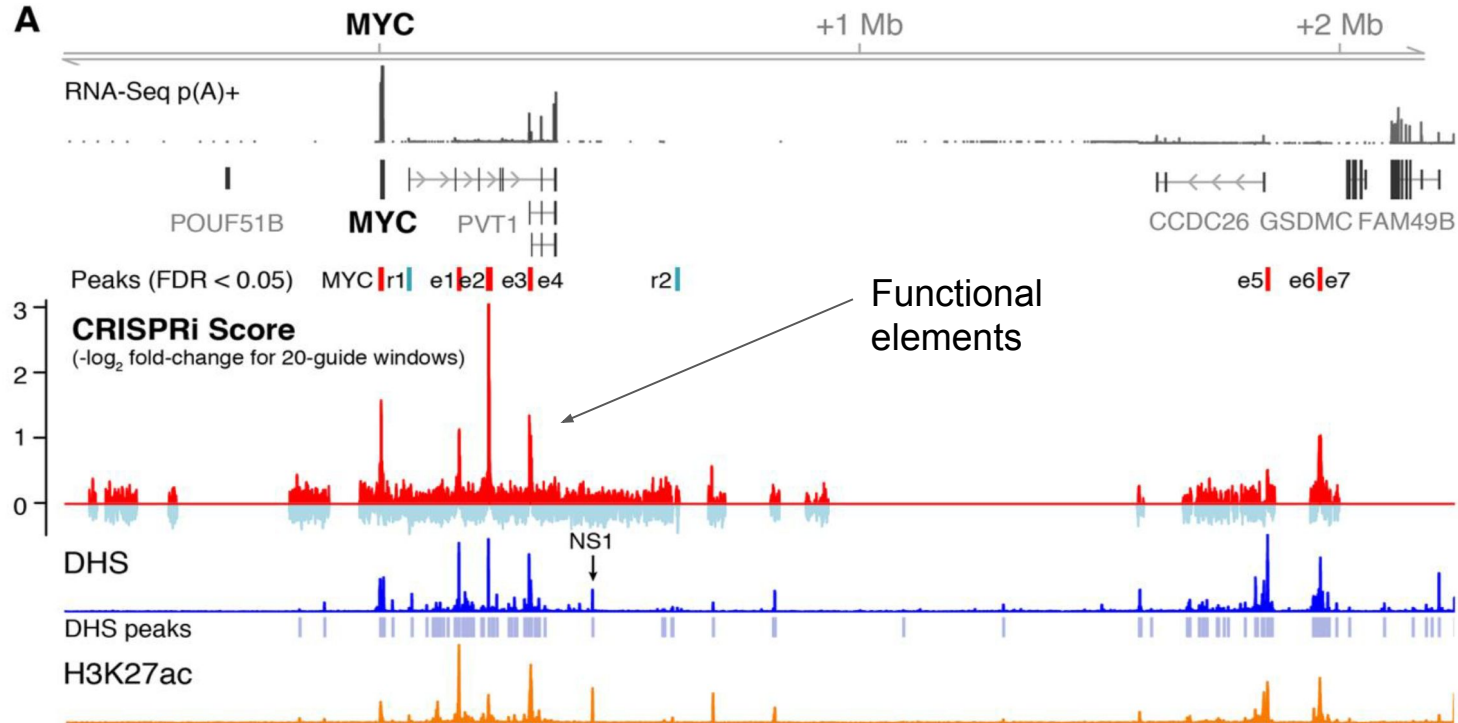
What were they
looking for?

Elements promoting drug
resistance in cancer

Elements promoting
cancer growth

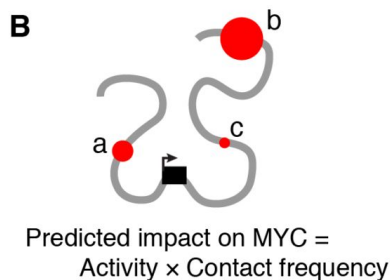
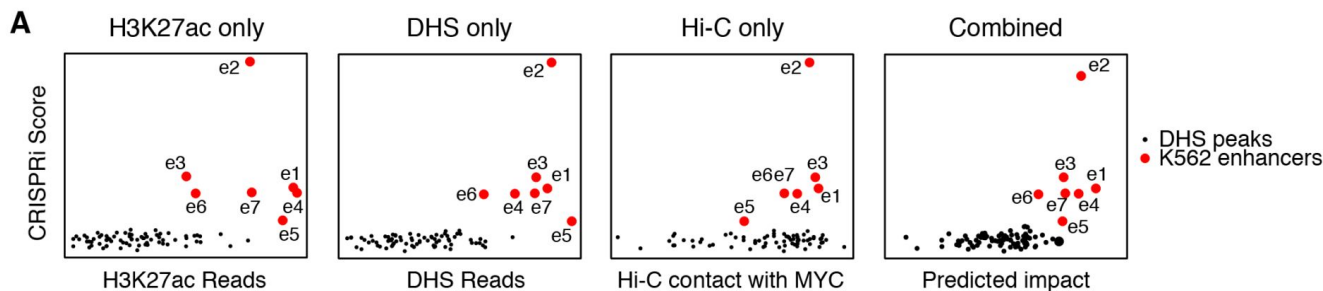


A noncoding screen to find functional elements



What's been done thus far with functional validation?

A very simple heuristic...



An open-ended question:

Given new biological data -- DNA-binding proteins, open chromatin, histone modifications, chromosome conformation capture, raw DNA sequence -- can we say whether it's functional or not?

Computational workflow

Data download

ENCODE ChIP-seq
& DHS BED data

Common SNPs
from UCSC database

Fulco screen data
from supplemental

Analysis

Overlapped all the
features with CRISPR
screen data

Used regression models
to predict CRISPR score

Categorized CRISPR
scores into classes

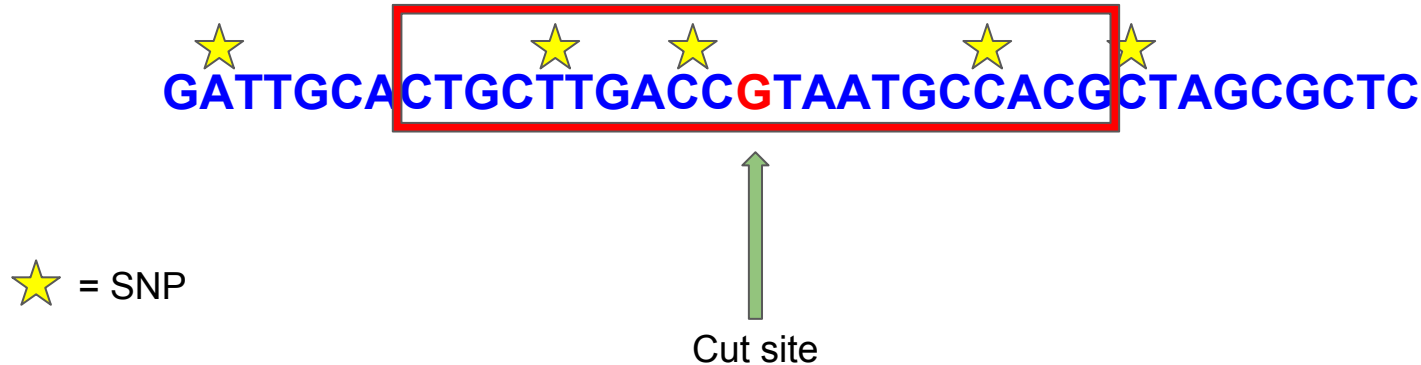
Used machine learning
models to predict
CRISPR score

Which features were best
in predicting functional
sequence?

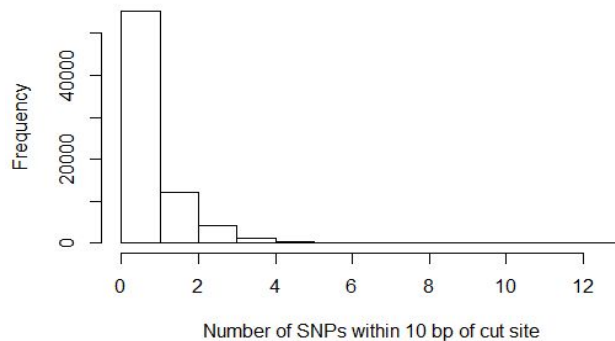
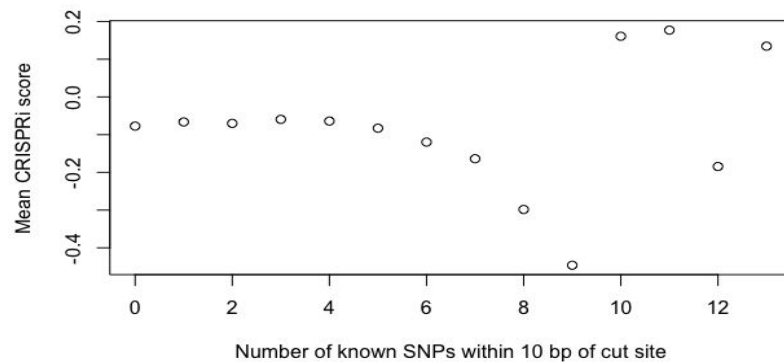
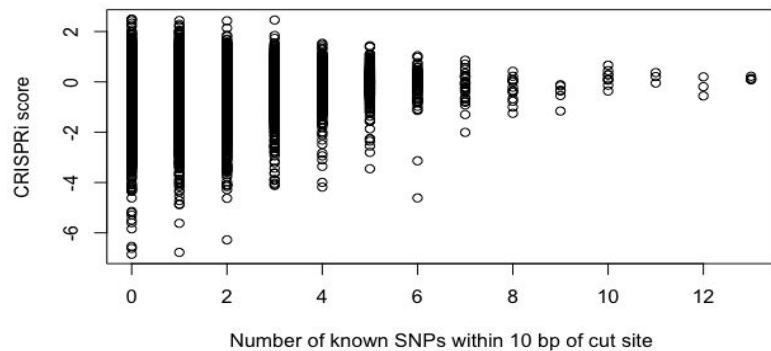
Which models were best?

Is functionality correlated with human-to-human conservation?

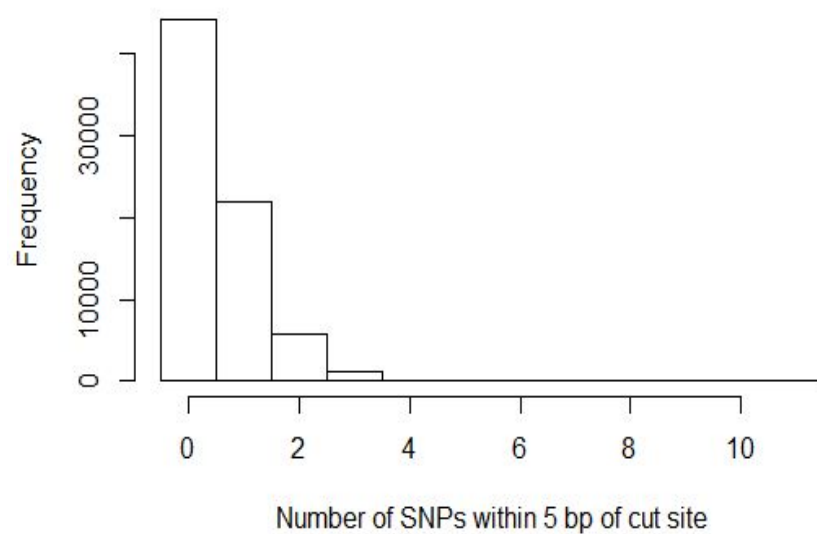
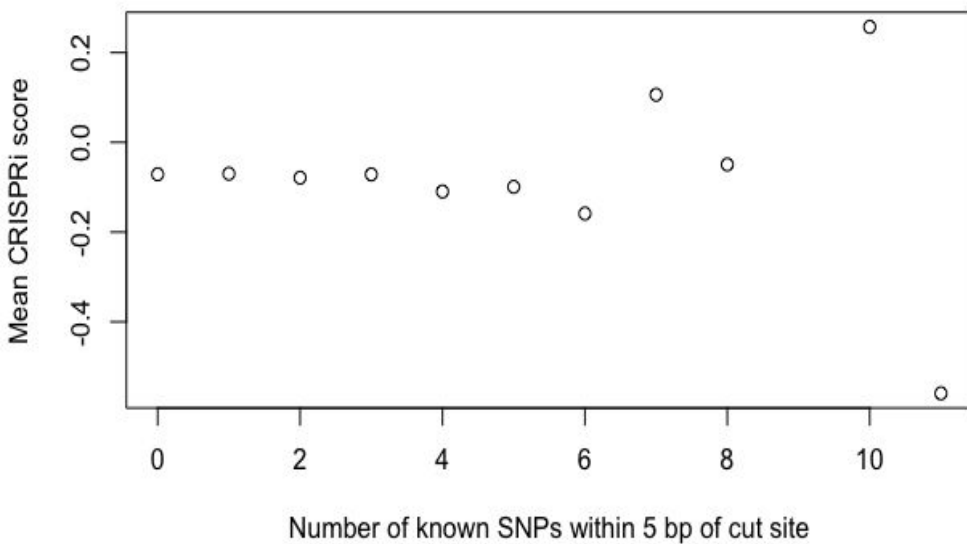
- A list of known SNPs and their allele frequencies in the region of our CRISPRi scores was queried from dbSNP (common SNPs list)
- For each CRISPR cut site, we counted how many SNPs fell within a distance of the site, and investigated how this number relates to CRISPRi score
- For example, the SNP number for this cut site would be 3.



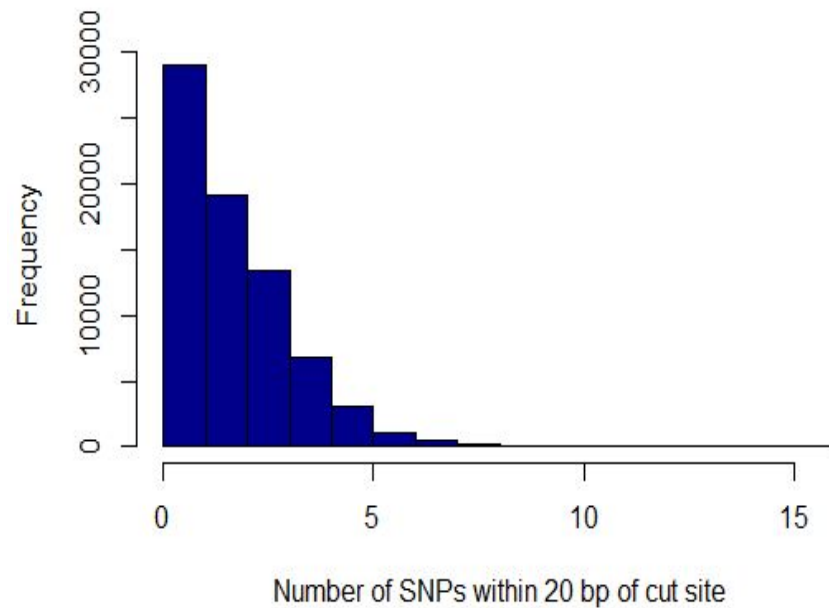
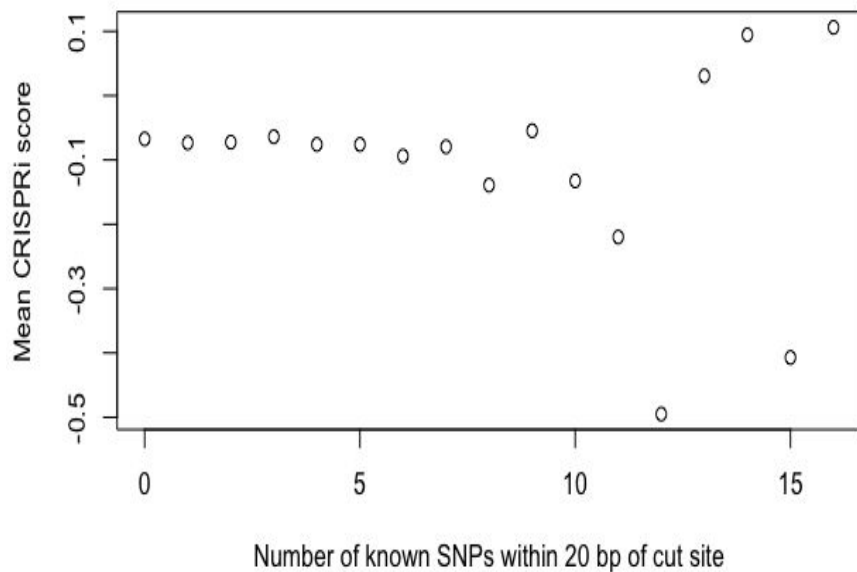
10 bp windows



Similar for 5 bp...

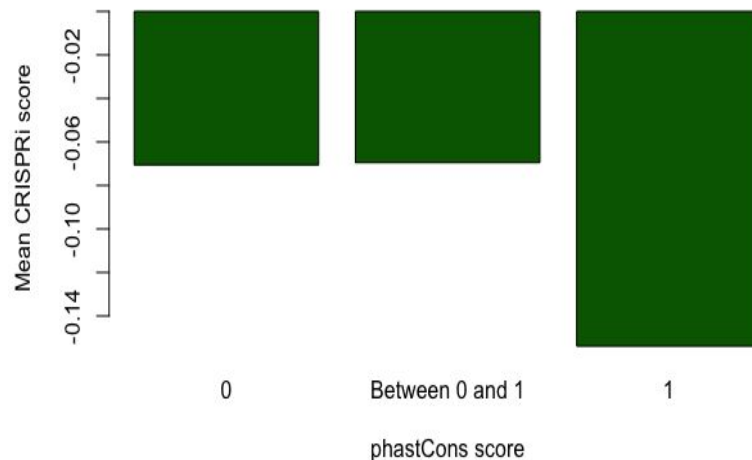
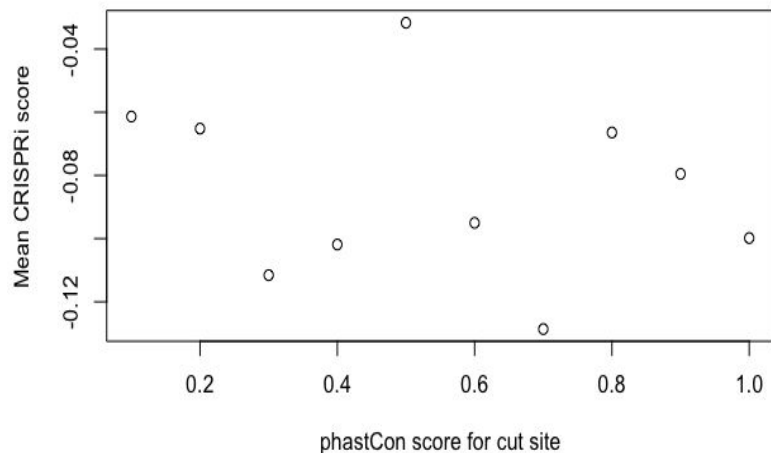


Not as good for 20 bp...



Conservation across vertebrate species

PhastCons: a score from 0 to 1 representing how conserved the sequence is



Other feature data



Methods

Used three different prediction algorithms: SVM, Random Forest, KNN

Applied 10-fold cross validation to each:
the training and testing samples
entire data set.

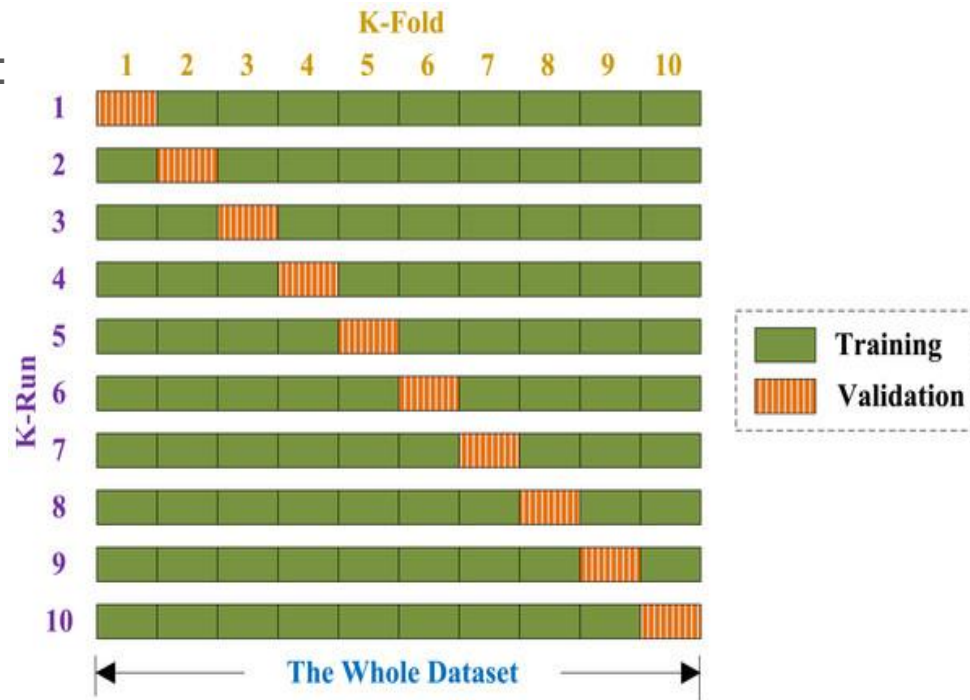


Image from: "Detection of Alzheimer's disease by displacement field machine learning." Y. Zhang and S. Wang. *PeerJ*.

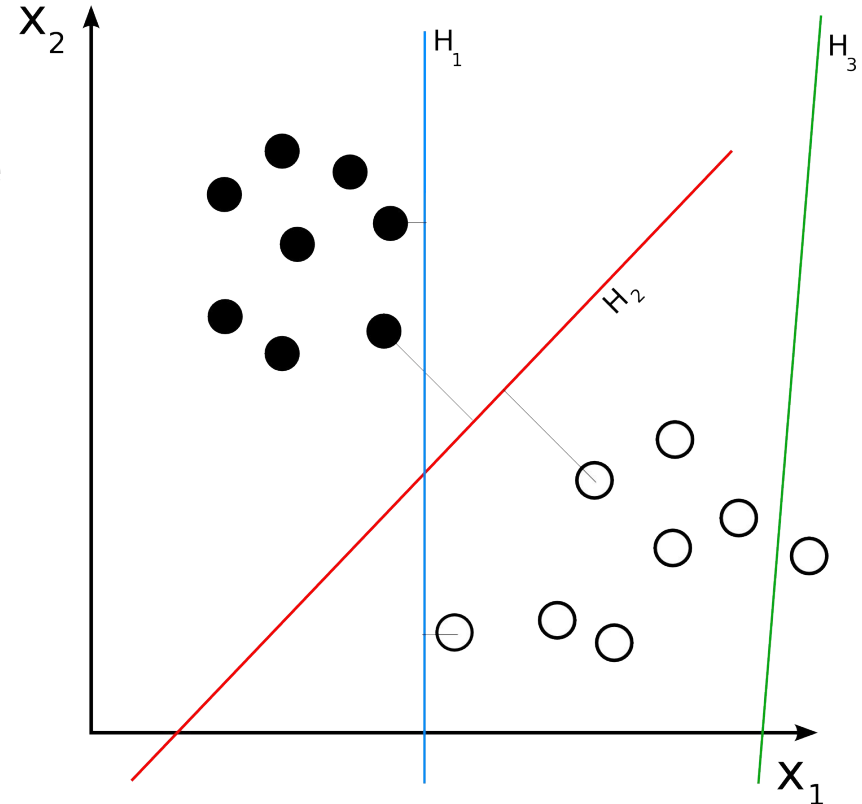
SVM - support vector machines

Idea: A classifier builds a model from a labeled training set and returns an optimal hyperplane that puts new examples into one of the categories non-probabilistically.

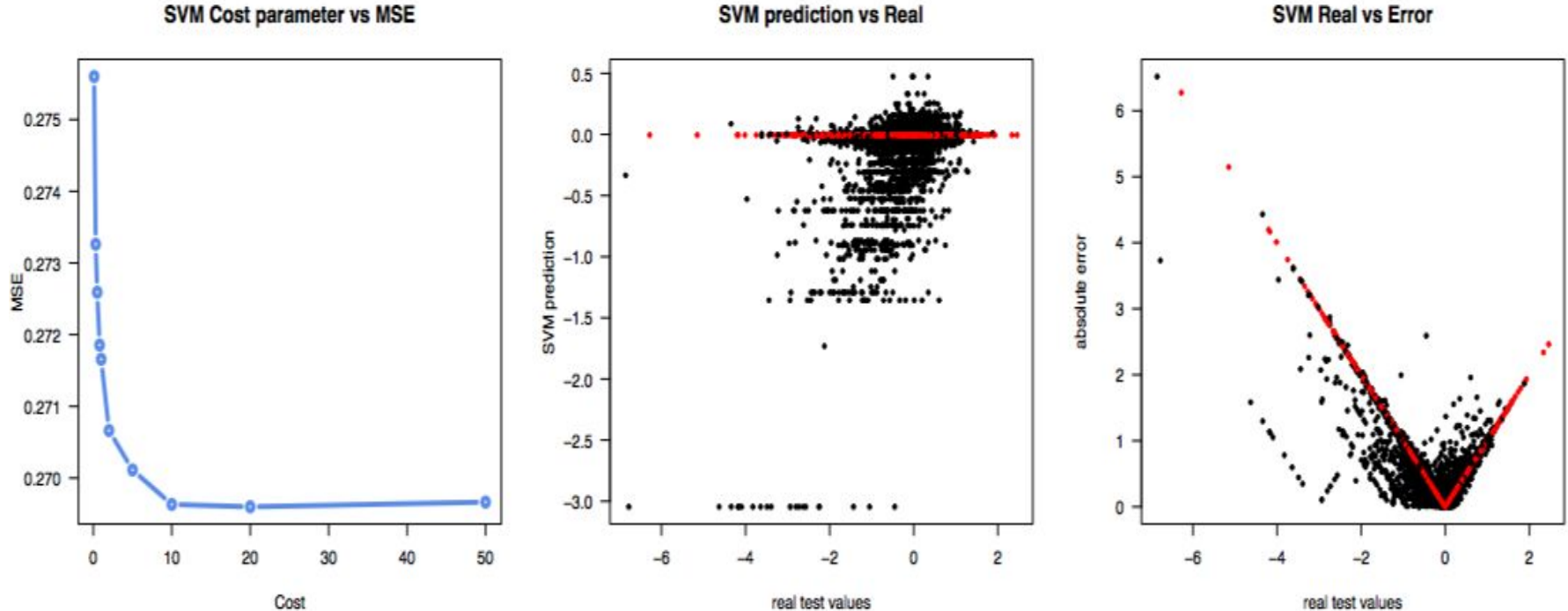
Goal: Maximize the hyperplane margin.

This represents the decision boundary.

Image: Svm separating hyperplanes.png, CC BY-SA 3.0
<https://commons.wikimedia.org/w/index.php?curid=22877598>



SVM - support vector machine results



red points are where there are 0 values for all features i.e. no presence of any data type at the genomic site

Random Forest

Uses a series of decision trees -> A forest!

Create random bootstrapped subsets of data and create decision trees from them.

What class does each tree predict? The these outputs is our prediction.

Bootstrap aggregating: “bagging.”

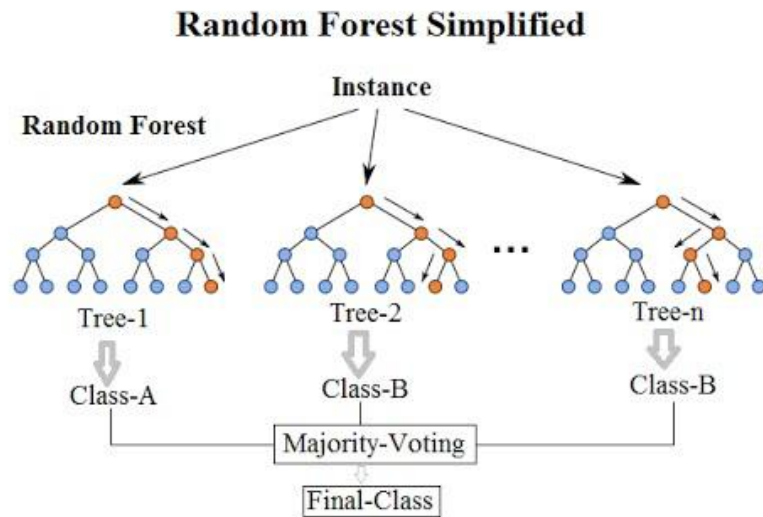
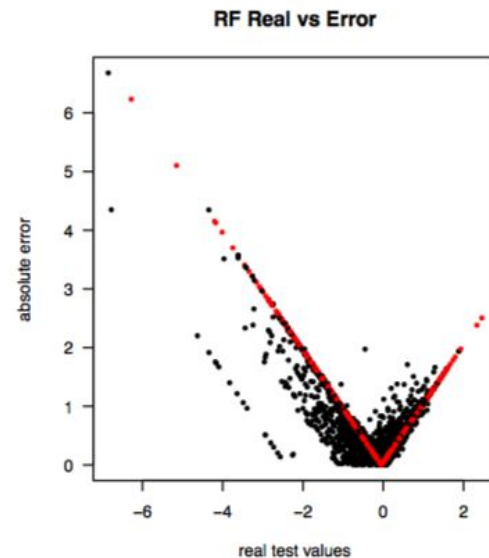
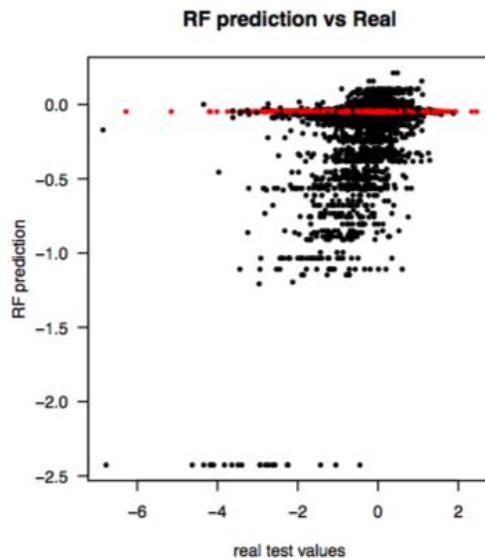
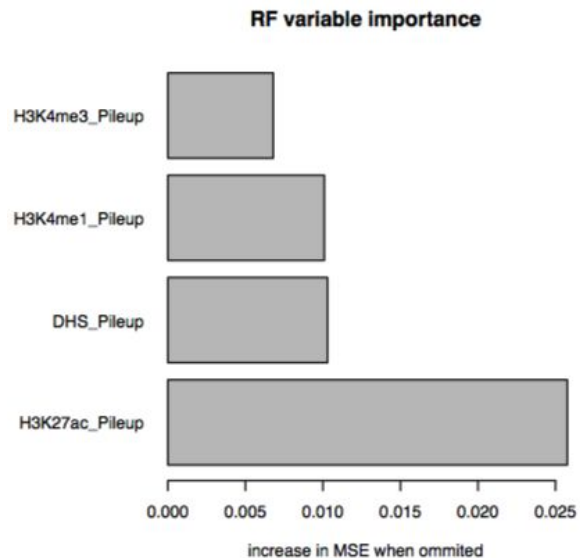
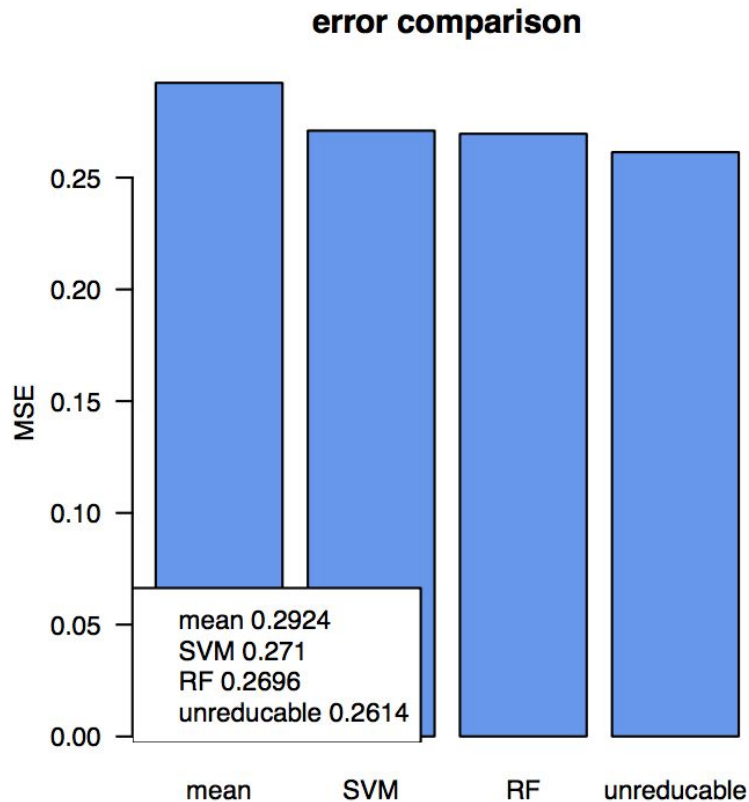


Image from “Random Forest based Classification:” <https://www.youtube.com/watch?v=ajTc5y3OqSQ>

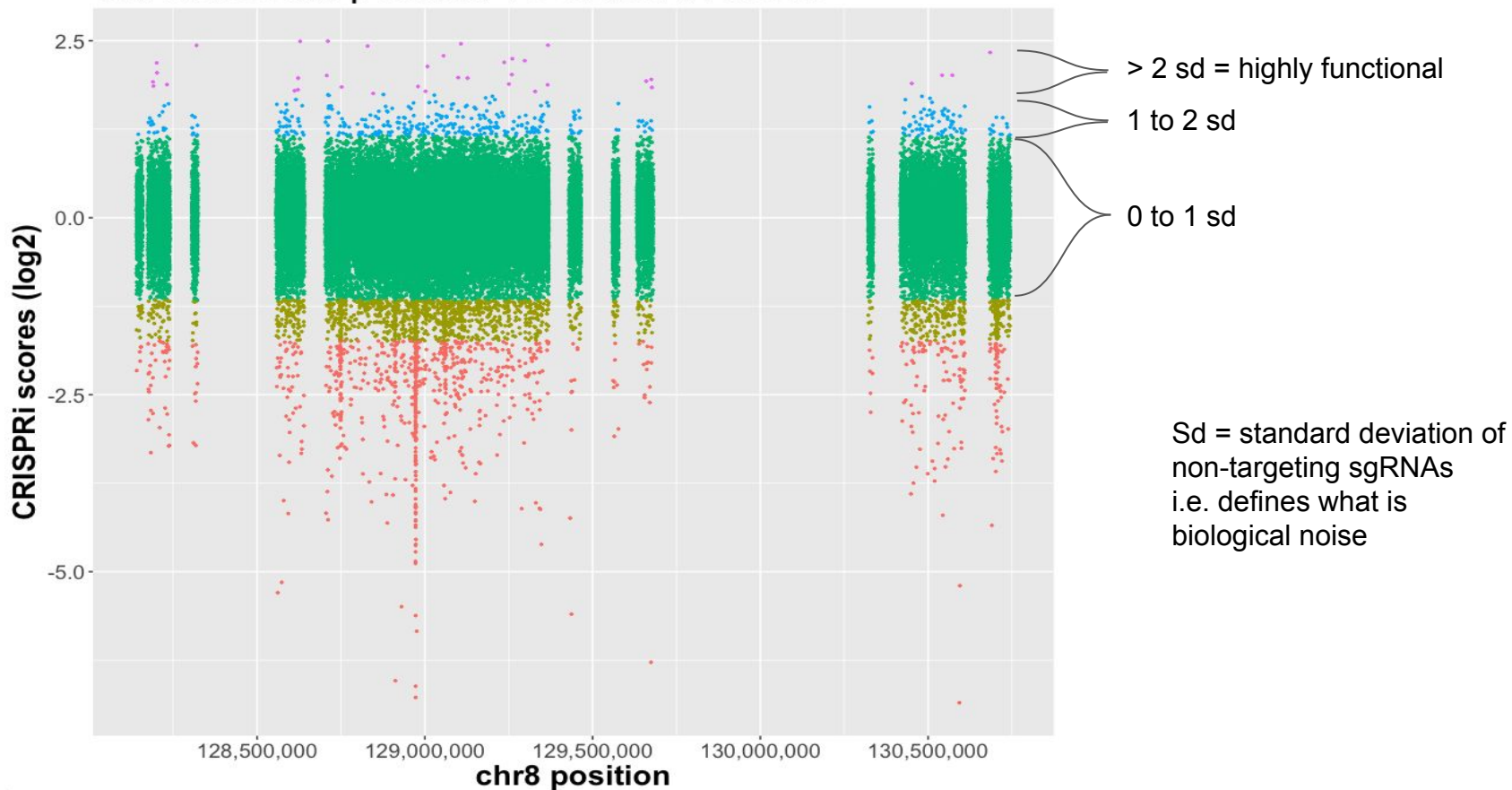
Random Forest - results



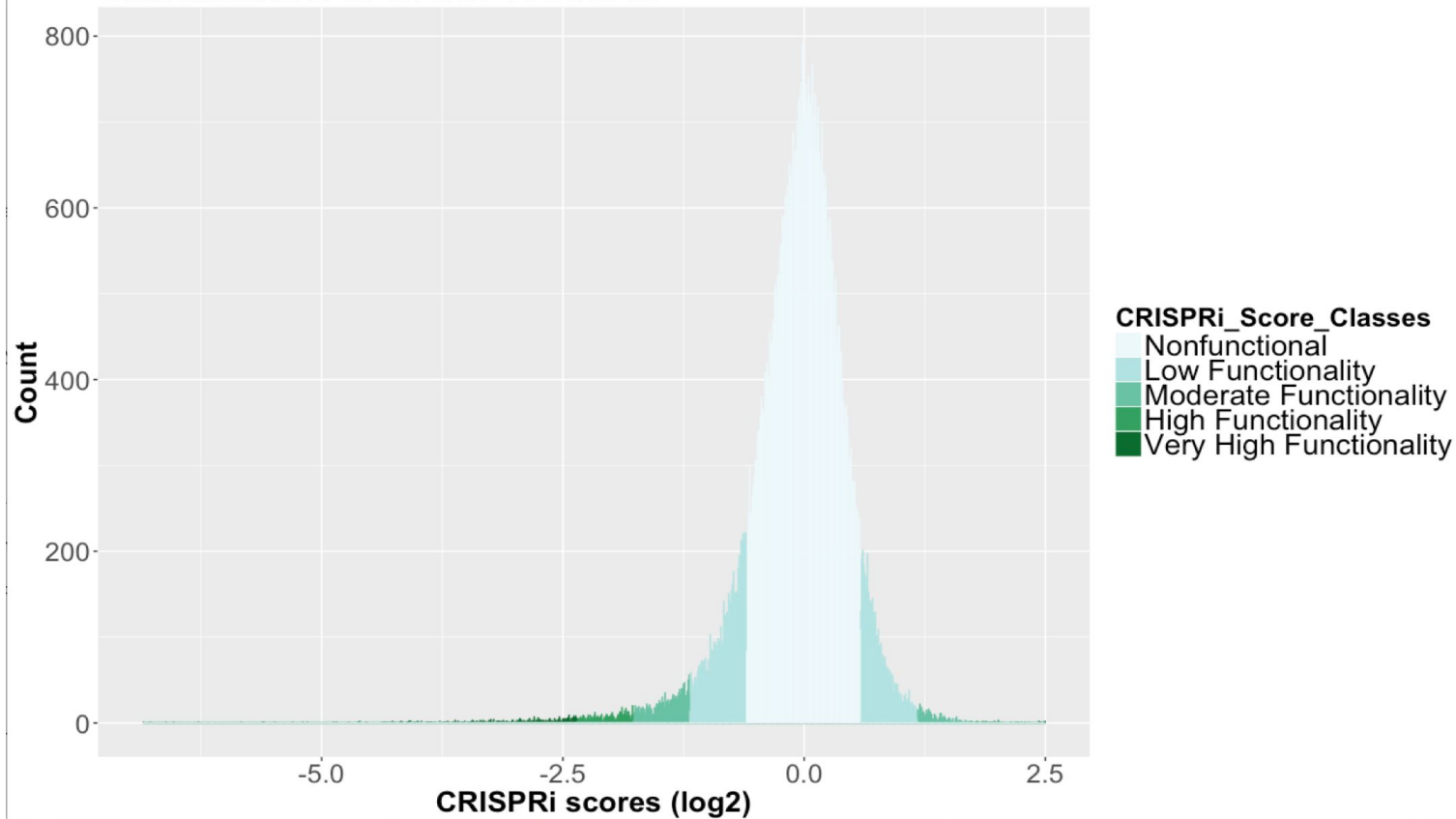
Comparing SVM and RF



Chromosome position vs CRISPR score



Classification of CRISPR scores



K-nearest neighbors

Makes predictions based on an instance's nearest neighbors.

Very simple: no need for assumptions, just look at the neighbors. Can be robust noise.

Interpreting model can be difficult: can't say, "As this happens to X, that happens to Y."

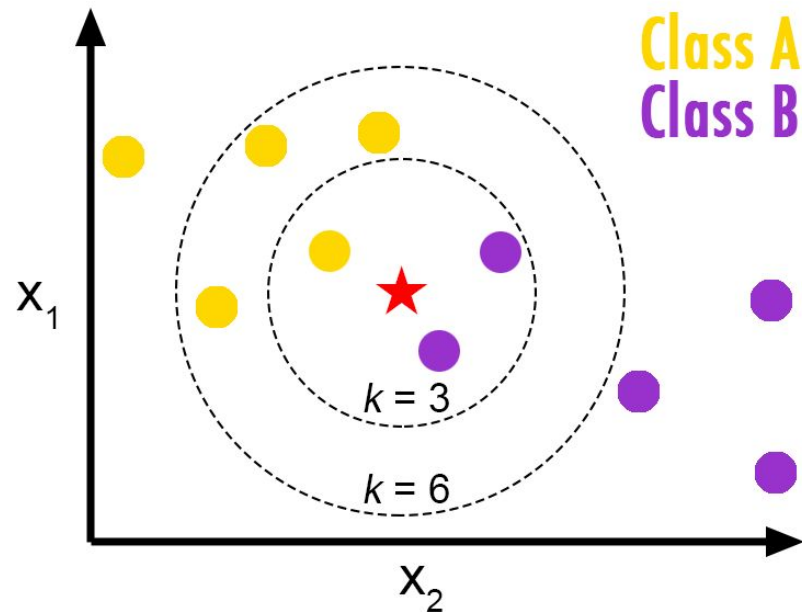
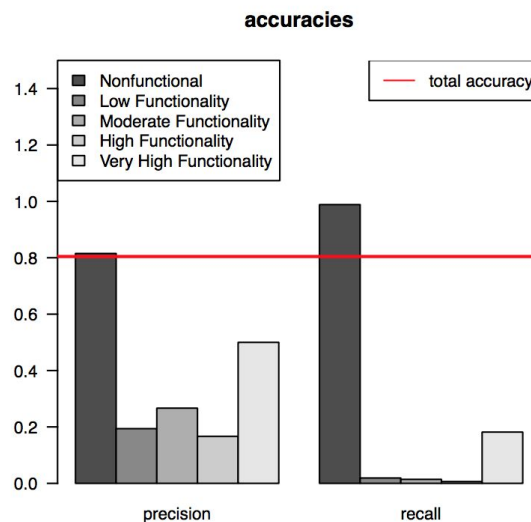
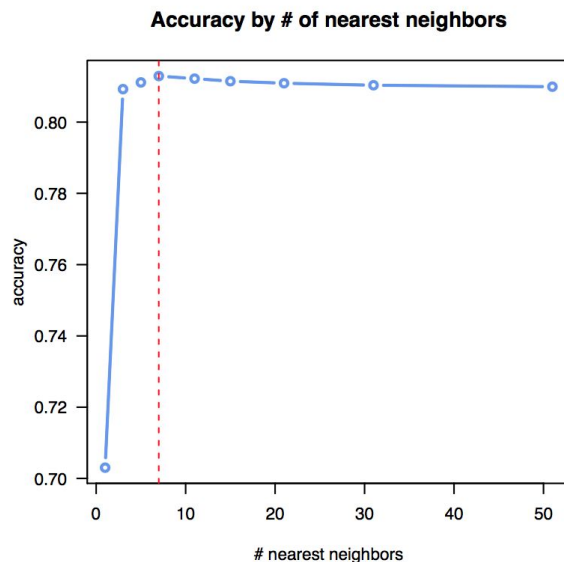


Image from: <http://en.proft.me/2017/01/22/classification-using-k-nearest-neighbors-r/>

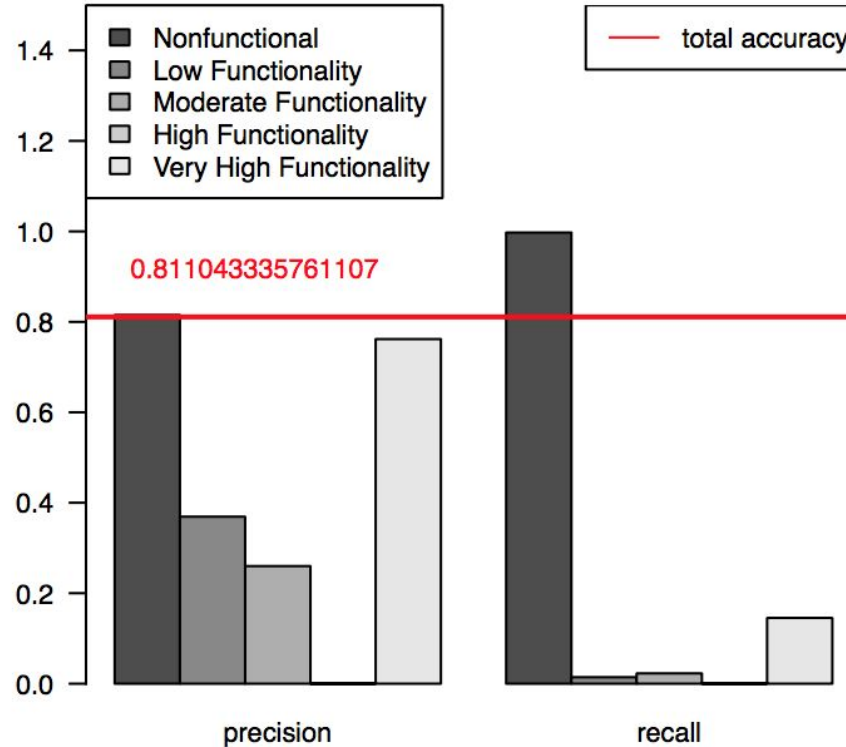
K-nearest neighbors - results

Divided the scores into five different classes

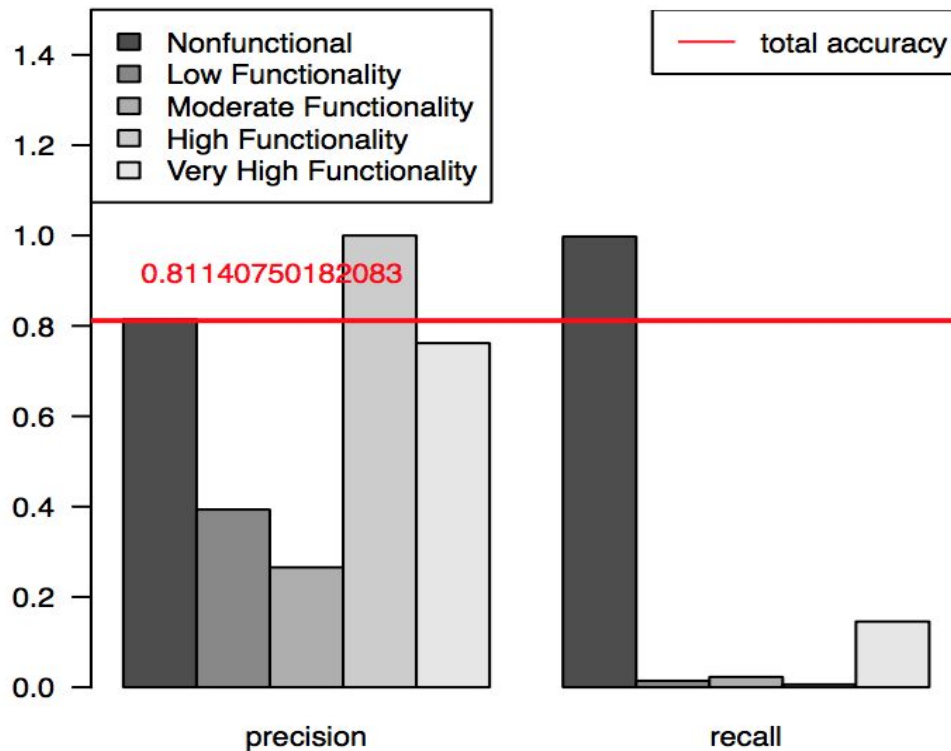
Highest accuracy: 80.44883% with $k = 7$



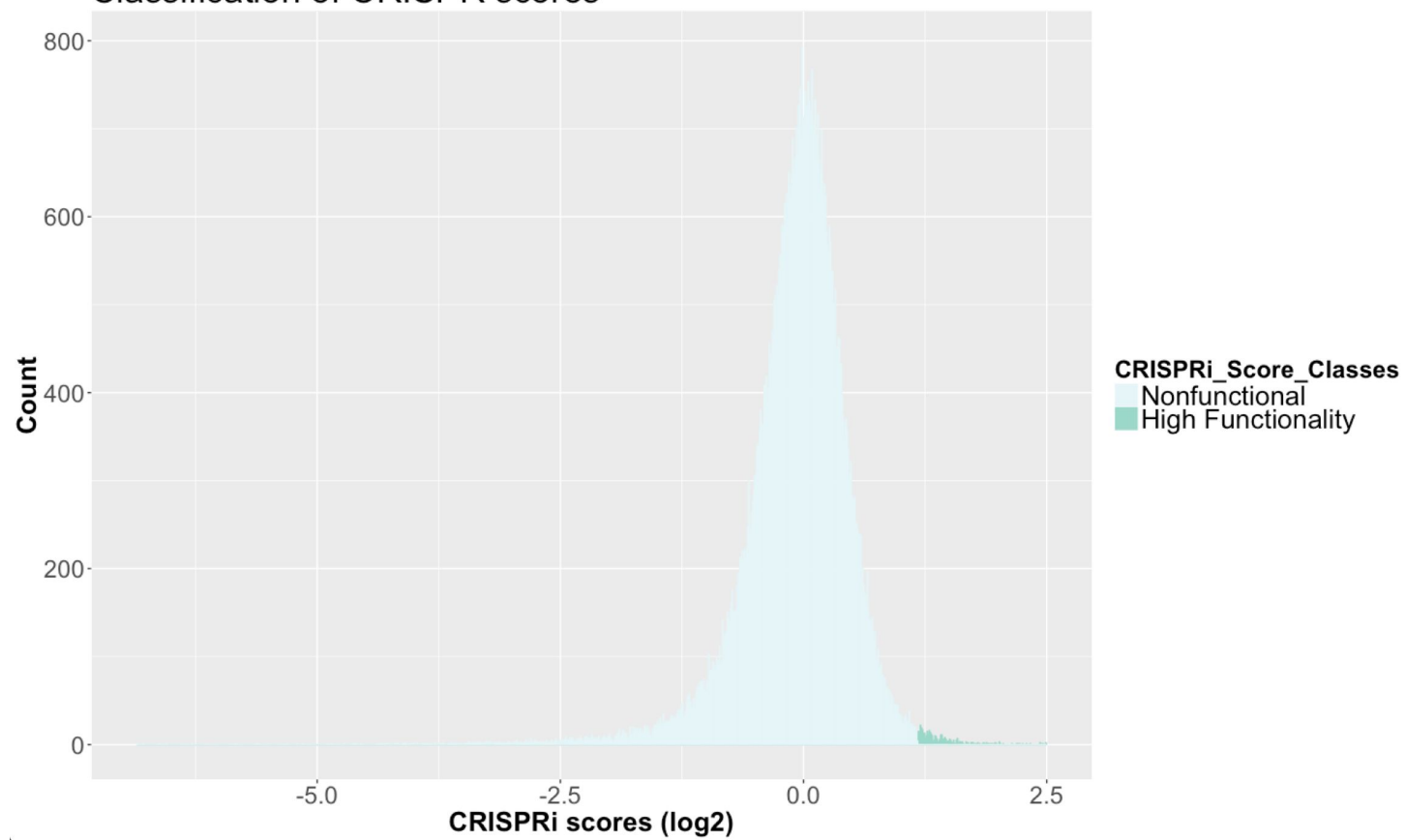
SVM results - 5 classes



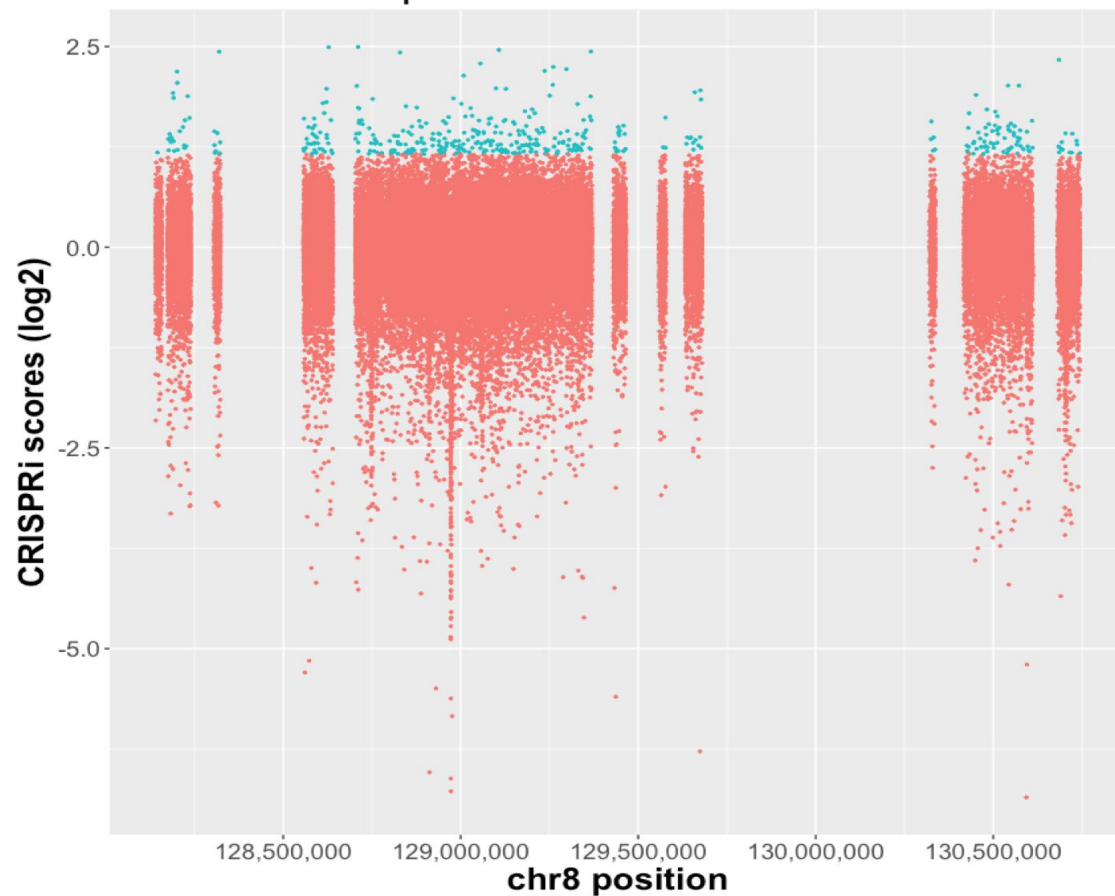
Random Forest results - 5 classes



Classification of CRISPR scores



Chromosome position vs CRISPR score



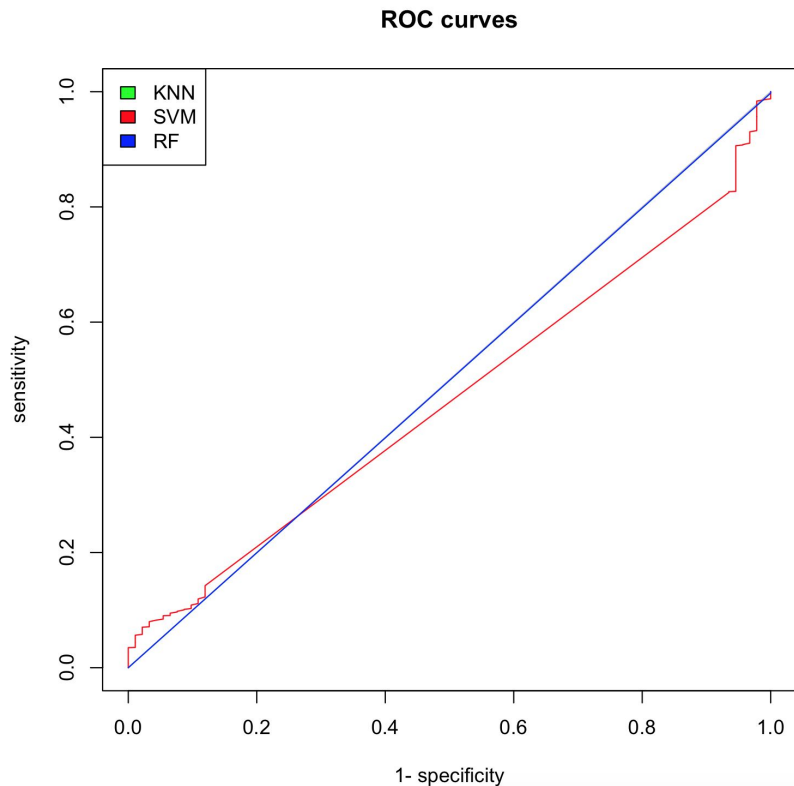
> 2 sd = highly functional

0 to 2 sd

Sd = standard deviation of
non-targeting sgRNAs
i.e. defines the range of
biological noise

Problems with two-class data

- Extremely unbalanced classes
- ROC plot of 50/50 balance between the two classes



Results

Sequence conservation - sequences conserved across humans have a slightly higher CRISPR score, while sequences conserved across vertebrates have a slightly lower CRISPR score

SVM - noticeable reduction in error from naive mean

Random Forest - slight improvement over SVM

KNN - didn't perform as well as SVM or RF

All models were more accurate with 2 class

Future work

Find and analyze **cancer-associated** SNPs

Integrate other high-resolution data types into the model:

- Transcription factor binding information (ChIP-seq)
- Intensity for Hi-C contact
- Frequency of Hi-C contact
- ChIA-PET (promoter-enhancer interactions)

Understanding tissue-specific effects when comparing one screen to another

Much later... after we know which features are the most predictive, what do they mean in a biological sense?

Acknowledgements

Neville Sanjana

Brian Parker

Rich Bonneau



NEW YORK UNIVERSITY



Appendix

Testing different predictive models & resamplings

10-fold cross validation:

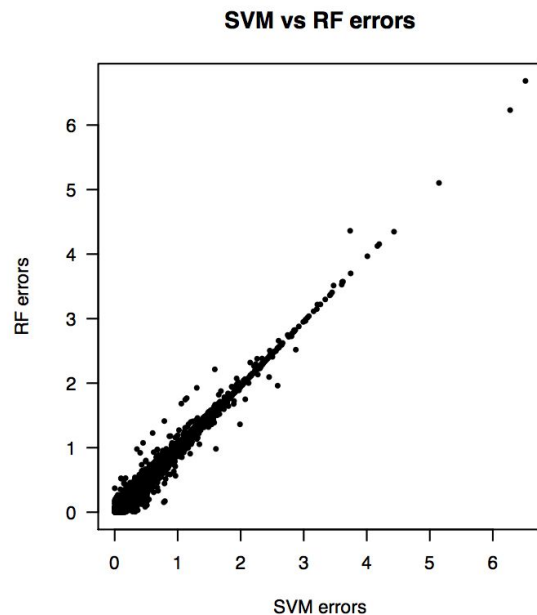
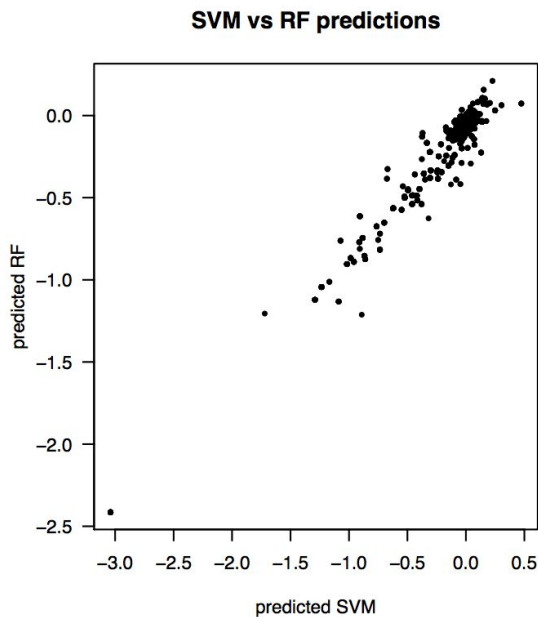
	Raw CRISPR score (Mean Squared Error)	5 classes (Mean Misclassification Error)	2 classes (Mean Misclassification Error)
SVM	0.20	0.20	
Random Forest	0.189	0.188	0.00541
KNN		0.231	0.0125

Bootstrapping (resampling the dataset using replacement):

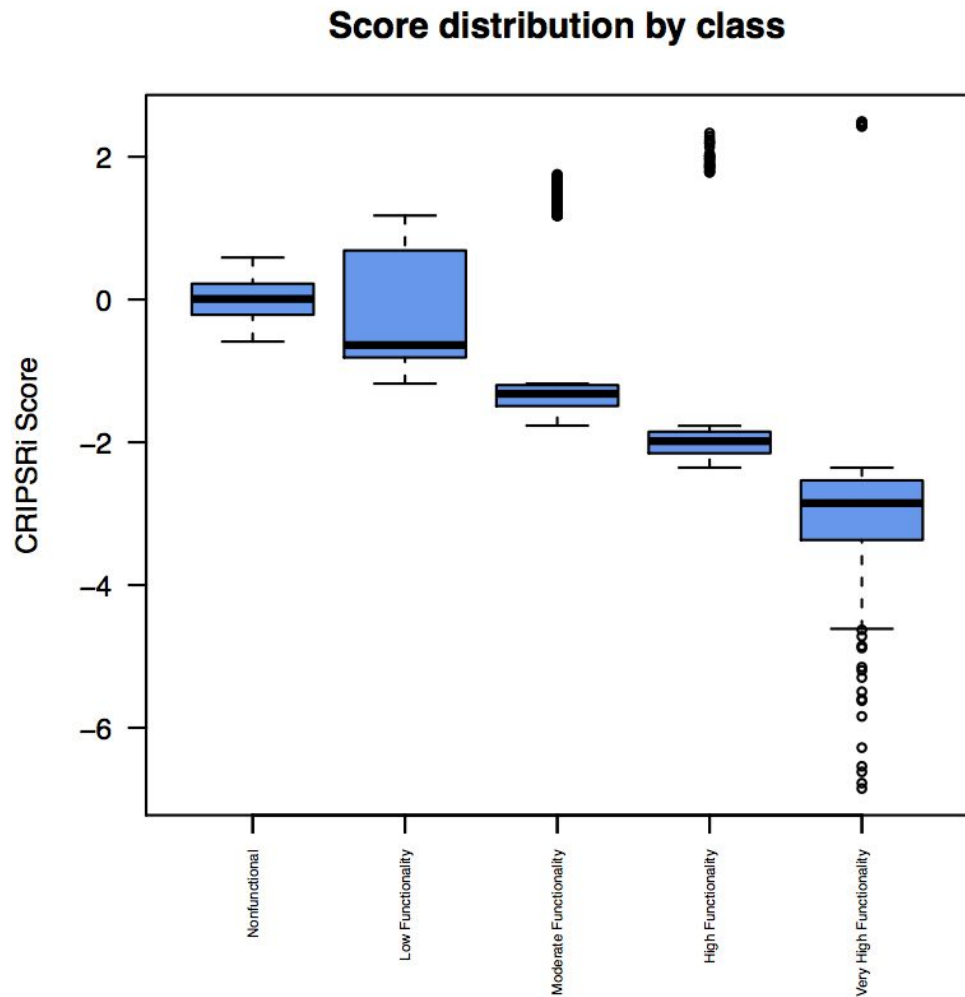
SVM			
Random Forest		0.188	0.00534
KNN		0.228	0.0141

SVM vs. RF (without classes)

The errors show a high correlation, meaning that if SVM makes a mistake - RF will also make a similar mistake on that entry.



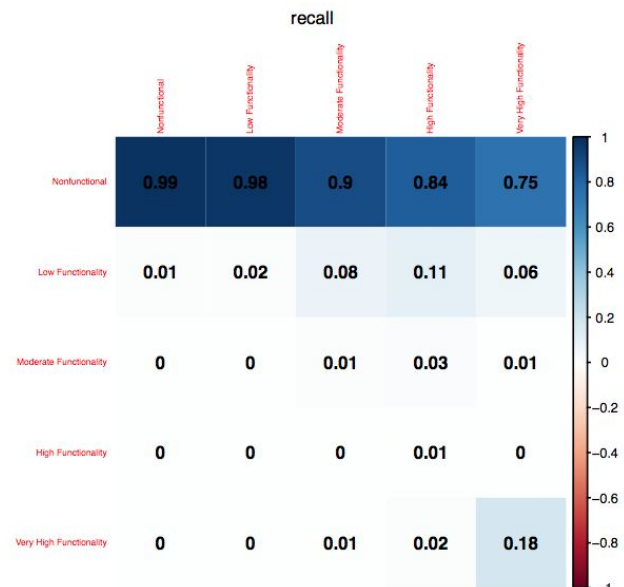
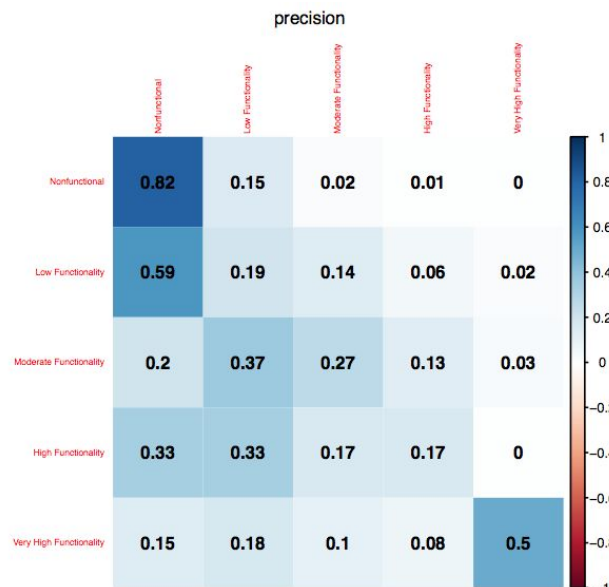
Distribution of 5 classes by score



Confusion matrix and precision/recall for kNN for 5 classes

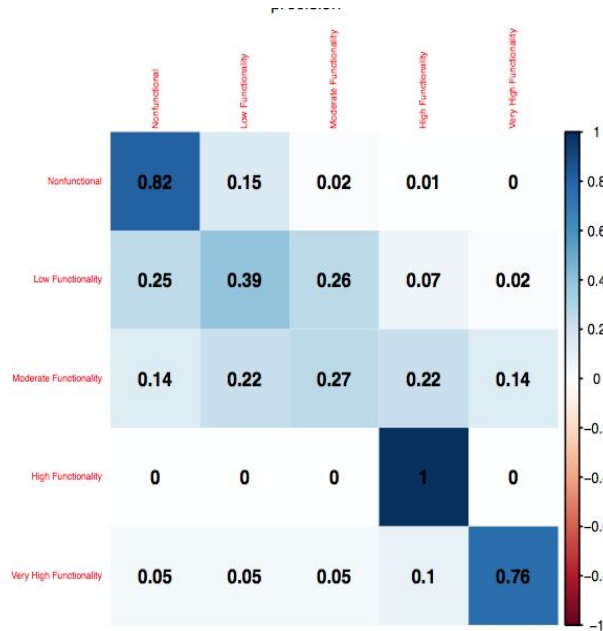
confusion matrix

	Nonfunctional	Low Functionality	Moderate Functionality	High Functionality	Very High Functionality
Nonfunctional	17581	3261	510	133	82
Low Functionality	191	63	46	18	7
Moderate Functionality	6	11	8	4	1
High Functionality	2	2	1	1	0
Very High Functionality	6	7	4	3	20



Confusion matrix and precision/recall for SVM for 5 classes

predicted		Nonfunctional	Low Functionality	Moderate Functionality	High Functionality	Very High Functionality
	Nonfunctional	17747	3284	523	137	84
	Low Functionality	31	48	32	8	3
	Moderate Functionality	7	11	13	11	7
	High Functionality	0	0	0	1	0
	Very High Functionality	1	1	1	2	16



Confusion matrix and precision/recall for RF for 5 classes

confusion matrix

predicted	Nonfunctional	Low Functionality	Moderate Functionality	High Functionality	Very High Functionality
Nonfunctional	17740	3282	521	137	84
Low Functionality	38	48	33	8	3
Moderate Functionality	7	11	13	12	7
High Functionality	0	2	1	0	0
Very High Functionality	1	1	1	2	16

precision



recall

