



COMPUTER SCIENCE

INDIANA UNIVERSITY

School of Informatics and Computing
Bloomington

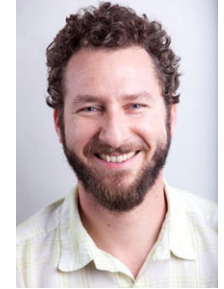
LEARNING PROPORTIONS IN A SEMI-SUPERVISED SETTING: A CASE STUDY IN PRECISION MEDICINE

Predrag Radivojac

DEPARTMENT OF COMPUTER SCIENCE AND INFORMATICS

INDIANA UNIVERSITY, BLOOMINGTON

November 28, 2016



From Hahn, Matthew William ★

Subject Tweet by Eduardo Eyras on Twitter

To Me <predrag@indiana.edu> ★

 Reply  Forward  Archive



Eduardo Eyras (@EduEyras)

7/14/15, 9:10 PM

One Third of the Alternative Splicing Isoforms produce Functional Proteins (P. Kim's lab) cell.com/cell-reports/a...

[Download](#) the Twitter app

Semi-supervised Learning Predicts Approximately One Third of the Alternative Splicing Isoforms as Functional Proteins

Yanqi Hao,^{1,2,10} Recep Colak,^{1,2,10} Joaquin
Mathias Wilhelm,⁴ Bernhard Kuster,^{4,5}
¹Terrence Donnelly Centre for Cellular and Biomolecular Research,
²Department of Computer Science, University of Toronto,
³Department of Medical Biophysics, University of Toronto,
⁴Chair for Proteomics and Bioanalytics, Technical University of Munich

Predicted Functional Proteins and Characterization of Highly Spliced Genes

Given the relatively high accuracy of our predictions, we ran PULSE to label 15,639 unlabeled transcripts from the BodyMap data set (Cabili et al., 2011; Colak et al., 2013; see Figure S1B for the experimental setup). At a 90% true-positive rate, we predict that about 32% of isoforms are functional (Figure 2C), roughly consistent with one school of thought, which estimate around 20%–30% to be functional (Floris et al., 2011; Rodriguez et al., 2013). Thus, AS leads to a sizeable number of previously uncharacterized proteins (a total of 5,023 in this data set alone). To prevent any biases/errors that might be caused by using Uniprot canonical isoforms for anchoring (the mechanism by which we quantify how much an alternative isoform differs from its canonical pair—see Experimental Procedures), we repeated the analysis using APPRIS principal isoforms for anchoring. The score distribution remains qualitatively unchanged (Figure 2C).



From Hahn, Matthew William ★

 Reply  Forward  Archive

Subject Tweet by Eduardo Eyras on Twitter

To Me <predrag@indiana.edu> ★



Eduardo Eyras (@EduEyras)

7/14/15, 9:10 PM

One Third of the Alternative Splicing Isoforms produce Functional Proteins (P. Kim's lab) cell.com/cell-reports/a...

[Download](#) the Twitter app

At some prediction threshold [one] Third of the Alternative Splicing Isoforms **predicted to** produce Functional Proteins....

PhosphoBase, a database of phosphorylation sites: release 2.0

Andres Kreegipuu, Nikolaj Blom^{1,*} and Søren Brunak¹

INTRODUCTION

Protein phosphorylation is a key event in many signal transduction pathways of biological systems. Protein kinases recognize and phosphorylate specific amino acid residues (mainly serine, threonine or tyrosine) in the substrate proteins. The research of protein phosphorylation has been essential in understanding intracellular signaling pathways. The number of identified phosphorylation sites is steadily growing—several thousand are now known. However, this seems to be only a small fraction of all potential phosphorylation possibilities since the estimated fraction of phosphoproteins may be as high as 30–50% of the total protein repertoire (1).

According to recent estimates, at least 30% of proteins in a eukaryotic proteome are phosphorylated², and >100,000 potential phosphorylation sites exist³. This poses a major analytical challenge, especially given that phosphorylated proteins are often present at substoichiometric ratios relative to their nonphosphorylated counterparts, and there can be multiple phosphorylation isoforms of a given protein. Furthermore, some phosphorylation events are transient and their detection requires monitoring of very short timescales.

Phosphoproteomics takes it easy

Paola Picotti

The EasyPhos pipeline simplifies analysis of phosphorylation-dependent signaling networks at high temporal resolution.

articles

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

** A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.*

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

The rediscovery of Mendel's laws of heredity in the opening weeks of the 20th century¹⁻³ sparked a scientific quest to understand the nature and content of genetic information that has propelled biology for the last hundred years. The scientific progress made falls naturally into four main phases, corresponding roughly to the four quarters of the century. The first established the cellular basis of heredity: the chromosomes. The second defined the molecular basis of heredity: the DNA double helix. The third unlocked the informational basis of heredity, with the discovery of the biological mechanism by which cells read the information contained in genes and with

coordinate regulation of the genes in the clusters.

- There appear to be about 30,000–40,000 protein-coding genes in the human genome—only about twice as many as in worm or fly.

However, the genes are more complex, with more alternative splicing generating a larger number of protein products.

- The full set of proteins (the 'proteome') encoded by the human genome is more complex than those of invertebrates. This is due in part to the presence of vertebrate-specific protein domains and motifs (an estimated 7% of the total), but more to the fact that vertebrates appear to have arranged pre-existing components into a

WHAT IS THE FRACTION OF ENZYMES IN A GENOME?

Proteins known to be enzymes:

The UniProt database in 2010: 28%

The UniProt database in 2014: 40%

Current state of the affairs:

$$E. coli: \frac{1154}{4433} = 0.26$$

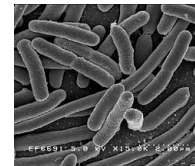
$$\text{Yeast: } \frac{1477}{6621} = 0.22$$

$$\text{Arabidopsis: } \frac{1563}{13100} = 0.12$$

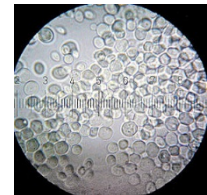
$$\text{Mouse: } \frac{1250}{16676} = 0.07$$

$$\text{Human: } \frac{2647}{20193} = 0.13$$

E. coli



Yeast



Arabidopsis



Mouse



Images from the Internet.

WHAT IS THE FRACTION OF ENZYMES IN A GENOME?

Expert opinion:



Charles Dann, Chemistry

E. coli: 0.35

Yeast: 0.45

Arabidopsis: 0.40

Mouse: 0.25

Human: 0.25



Tuli Mukhopadhyay, Biology

E. coli: 0.85

Yeast: 0.85

Arabidopsis: 0.85

Mouse: 0.85

Human: 0.85



Yuzhen Ye, Computer Science

E. coli: 0.40

Yeast: 0.40

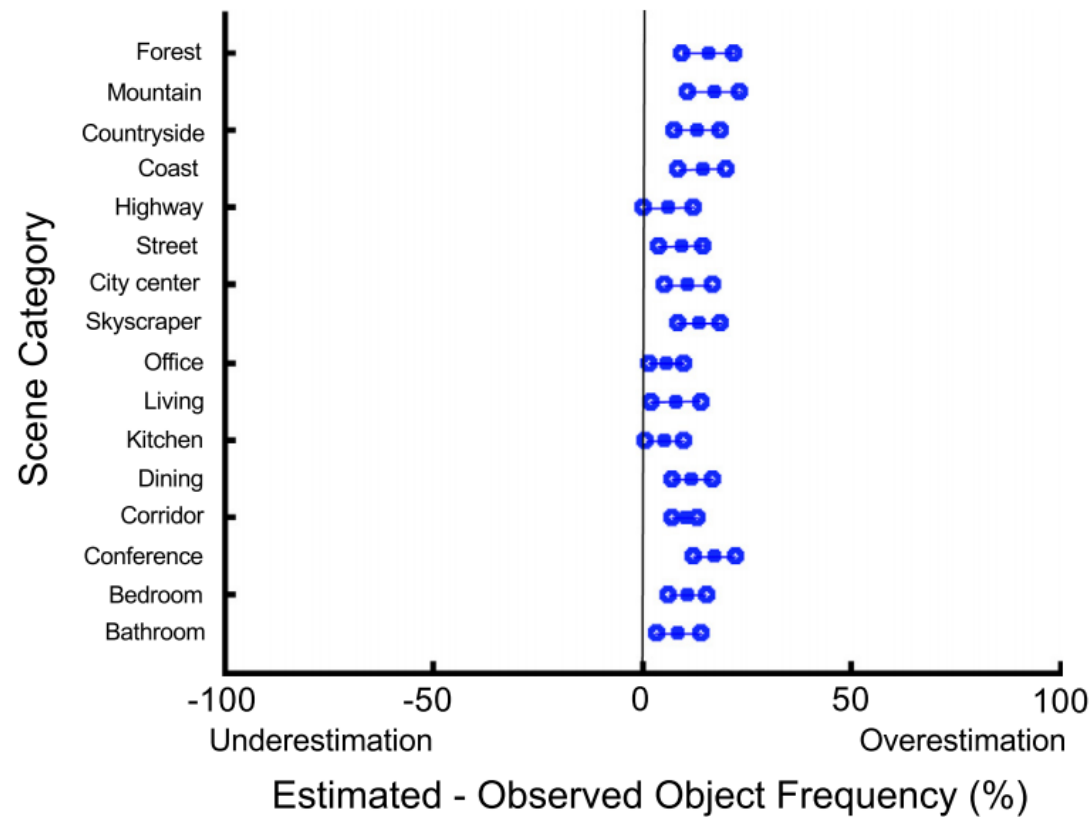
Arabidopsis: 0.30

Mouse: 0.30

Human: 0.30

EXAMPLE FROM PSYCHOLOGY

Estimations of object frequency are frequently overestimated!



AND SO WE GO...

Estimation of proportions:

- * Interesting problem
- * Needs rigorous theoretical treatment
- * Needs algorithms
- * Needs to be more precisely communicated

Hi Pedja,

You pose interesting questions. I'd expect yeast to have the highest enzyme fraction as it does not need to have conserved genes for multicellular development, cognition, etc. (though many of these processes requires signaling pathways with enzymes). So here are my estimates for enzyme fraction, based entirely on intuition.

Yeast ~45 %; E. coli ~35 %; Mouse ~25 %; Human ~25 %; Arabidopsis ~40 % (no idea here)

I imagine I may hit low on all of these...

CD3

SUPERVISED LEARNING PROBLEM

Given:

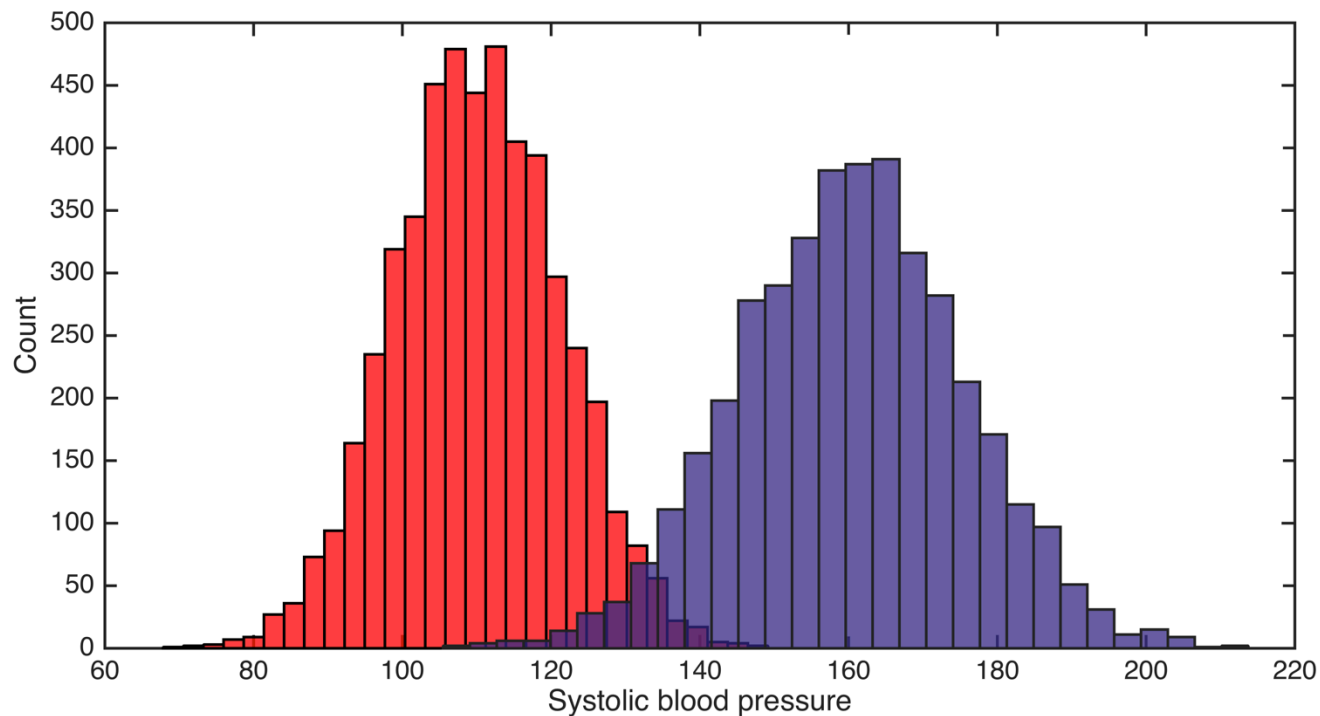
X_{red} : sample from people w/o heart disease

X_{blue} : sample from people w/ heart disease

Goal: predict heart disease

$x \in \mathbb{R}$

$y \in \{\text{disease, no disease}\}$



SUPERVISED LEARNING PROBLEM

Given:

$$x \in \mathcal{X}, y \in \mathcal{Y} = \{0, 1\}$$

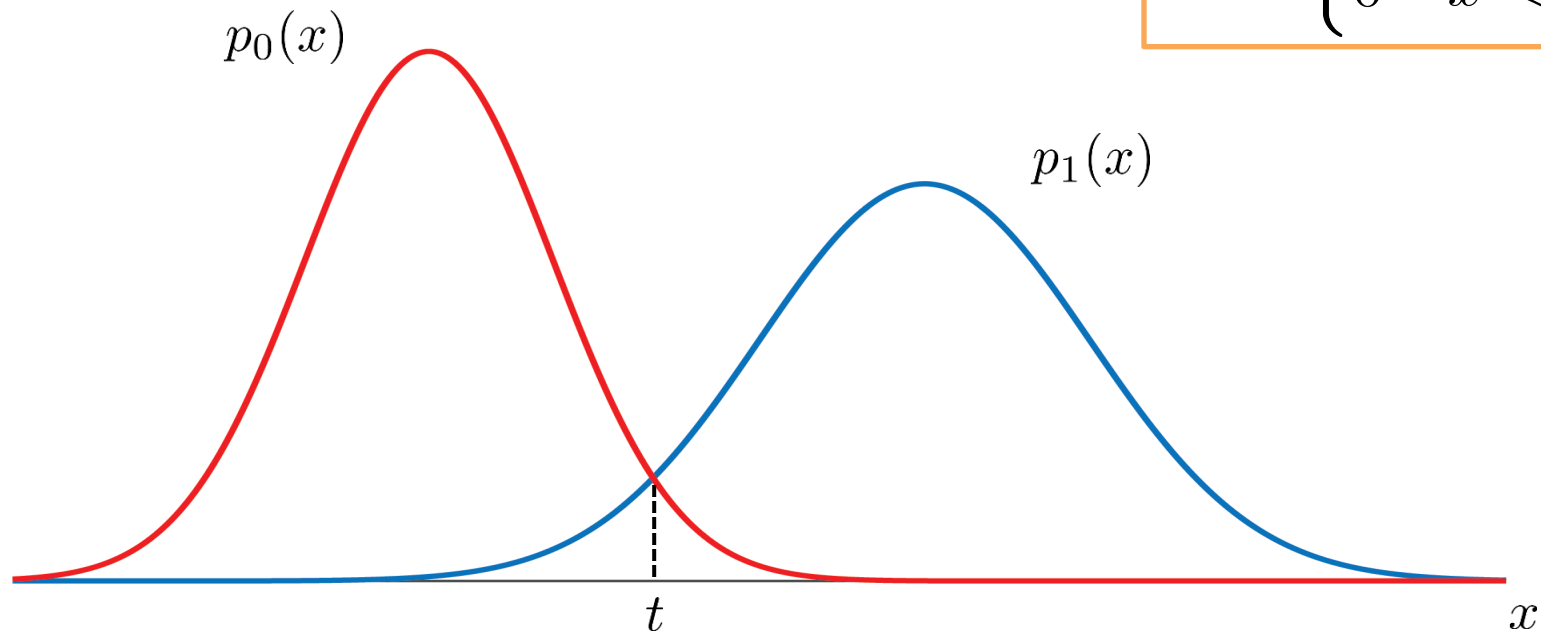
X_0 : sample from $p_0(x) = p(x|y = 0)$

X_1 : sample from $p_1(x) = p(x|y = 1)$

Goal: learn how x relates to y

Predictor:

$$\hat{y} = \begin{cases} 1 & x \geq t \\ 0 & x < t \end{cases}$$

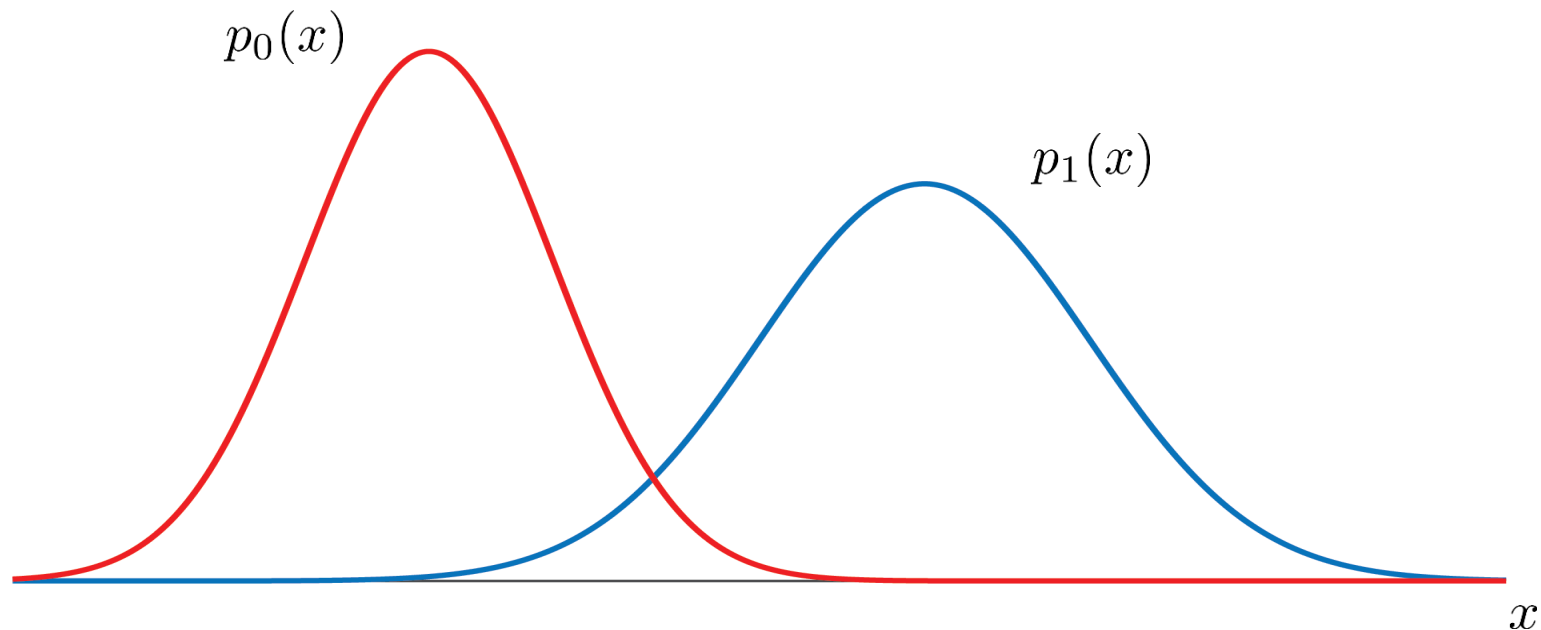


SUPERVISED LEARNING PROBLEM

Goal: learn how x relates to y

$x \in \mathcal{X}, y \in \mathcal{Y} = \{0, 1\}$

$$p(y|x) = \frac{p(x, y)}{p(x)}$$
$$\propto p(x|y)p(y)$$

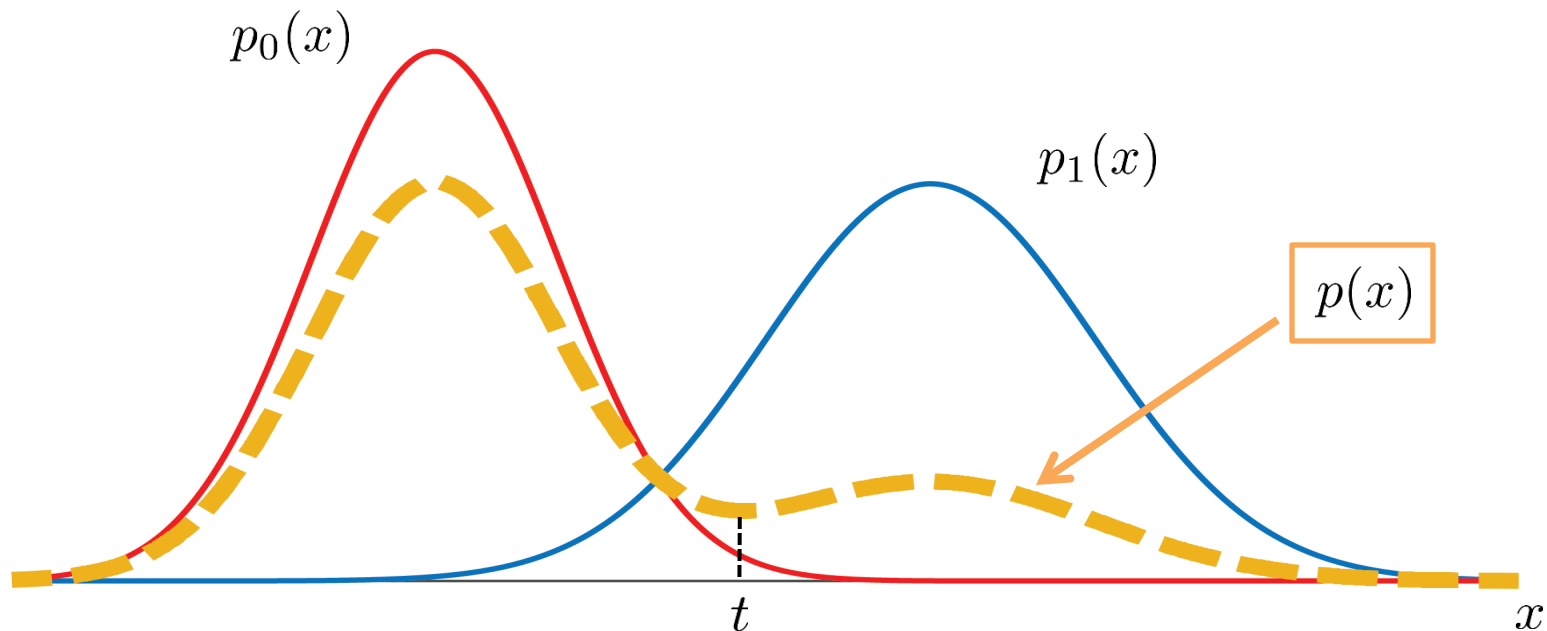


SUPERVISED LEARNING PROBLEM

Goal: learn how x relates to y

$x \in \mathcal{X}, y \in \mathcal{Y} = \{0, 1\}$

$$p(y|x) = \frac{p(x, y)}{p(x)}$$
$$\propto p(x|y)p(y)$$



SEMI-SUPERVISED LEARNING PROBLEM

Given:

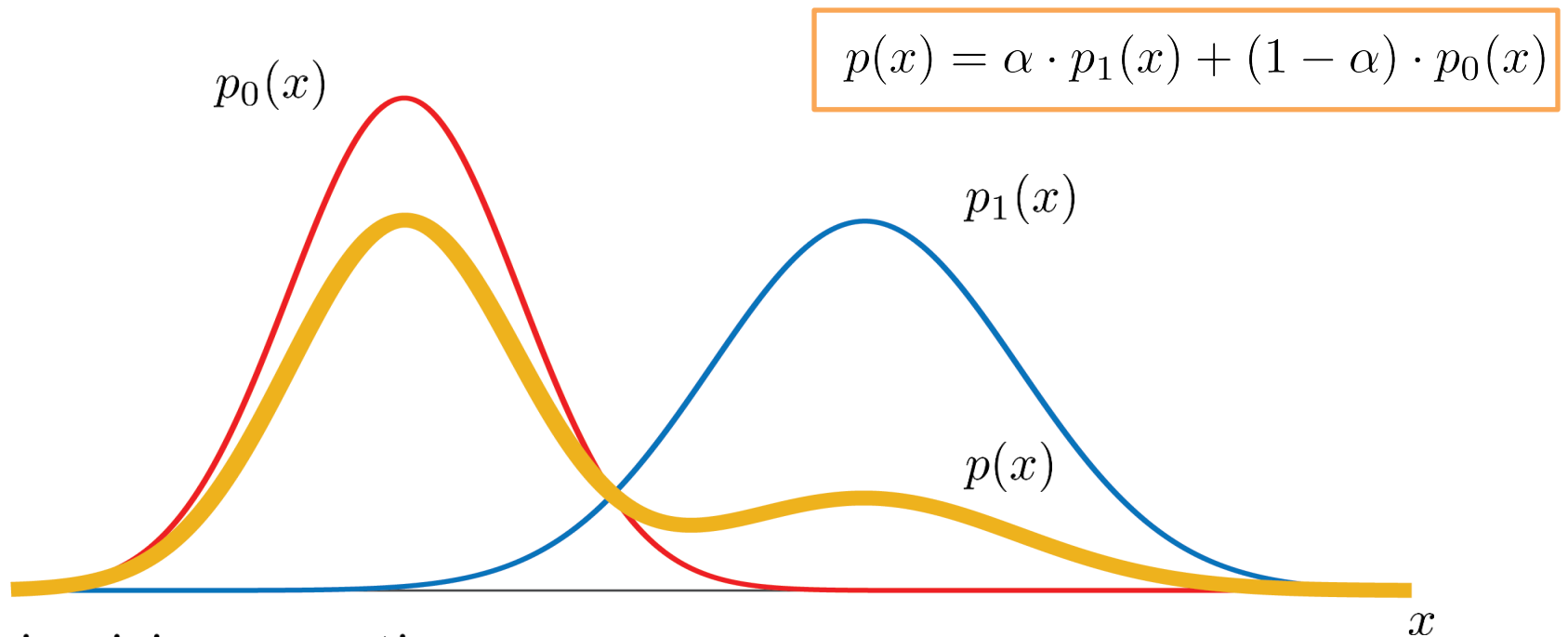
$$x \in \mathcal{X}, y \in \mathcal{Y} = \{0, 1\}$$

X_0 : sample from $p_0(x) = p(x|y = 0)$

X_1 : sample from $p_1(x) = p(x|y = 1)$

X : sample from $p(x)$

Goal: learn how x relates to y



$\alpha \in (0, 1)$ is mixing proportion

UNSUPERVISED LEARNING PROBLEM

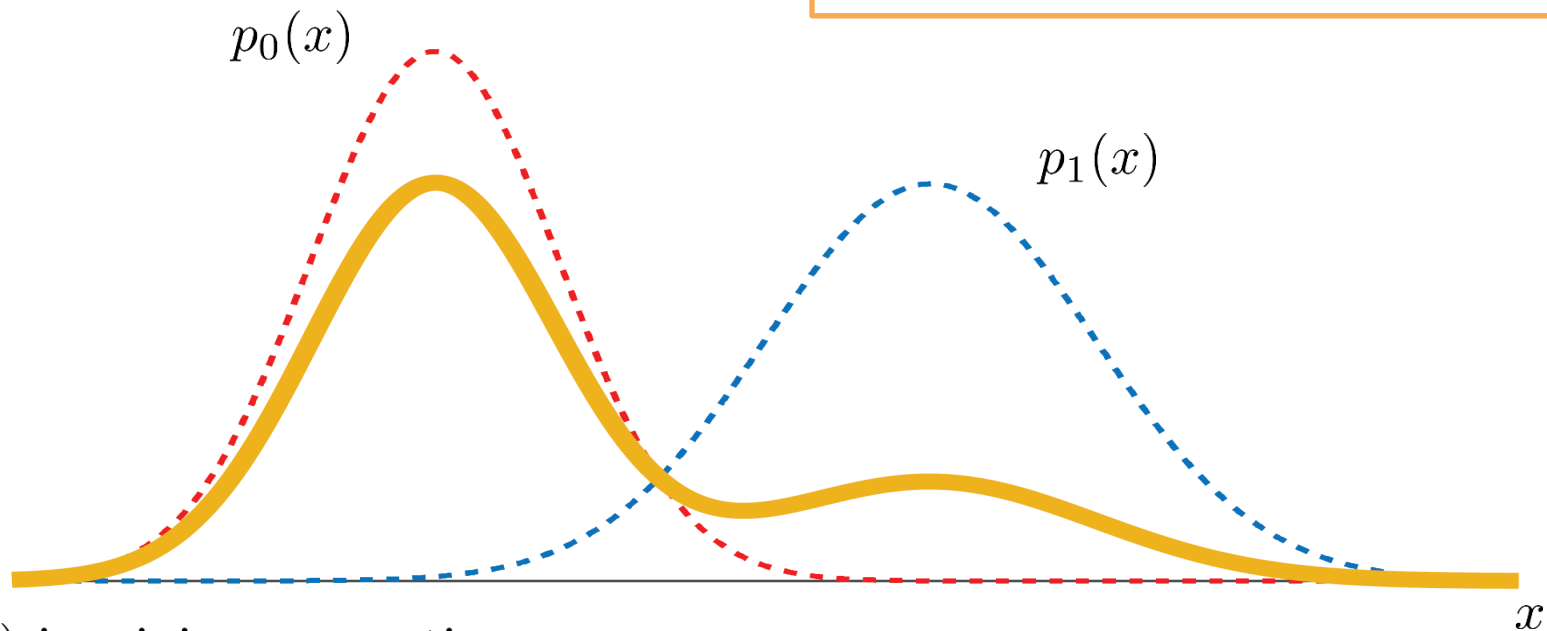
Given:

$$x \in \mathcal{X}, y \in \mathcal{Y} = \{0, 1\}$$

X = sample from $p(x)$

Goal: learn $p_0(x)$, $p_1(x)$ and α

$$p(x) = \alpha \cdot p_1(x) + (1 - \alpha) \cdot p_0(x)$$



$\alpha \in (0, 1)$ is mixing proportion

POSITIVE-UNLABELED LEARNING PROBLEM (PU)

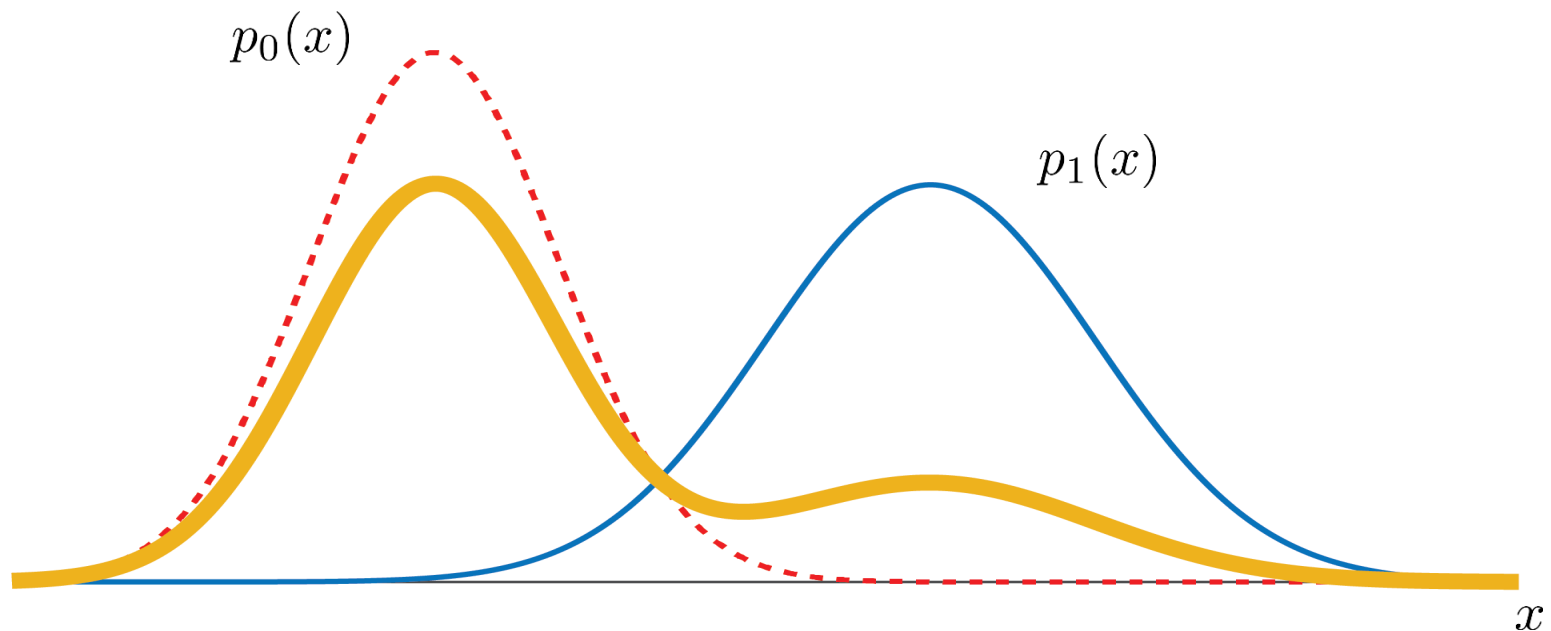
Given:

$$x \in \mathcal{X}, y \in \mathcal{Y} = \{0, 1\}$$

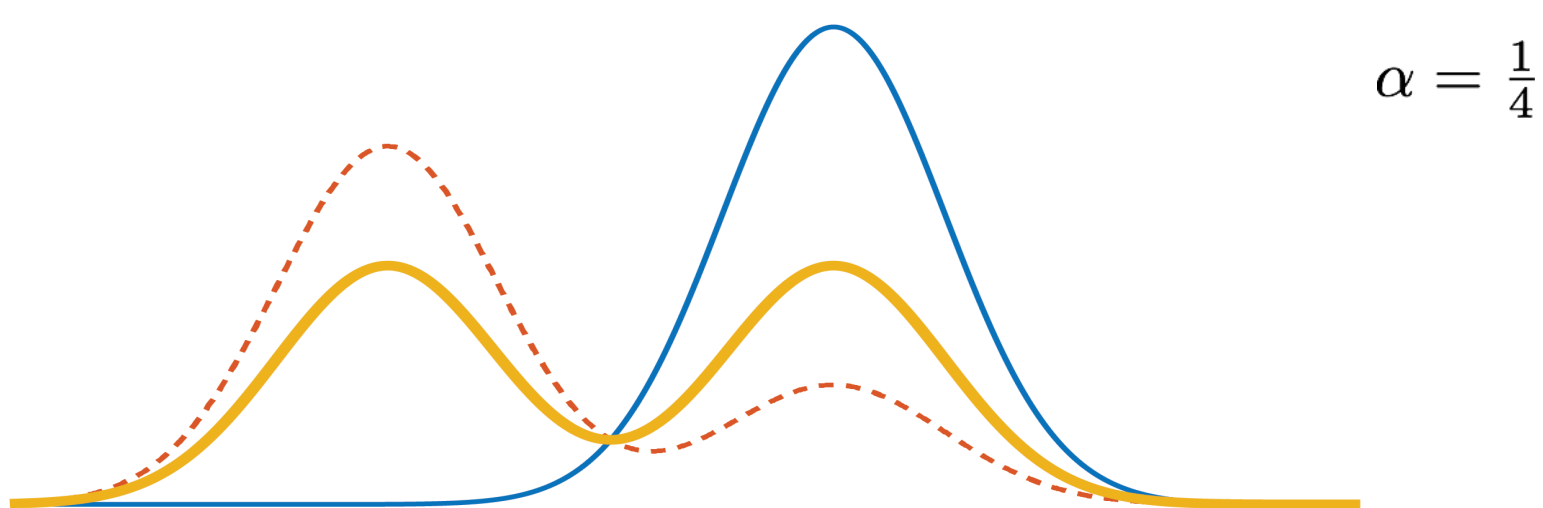
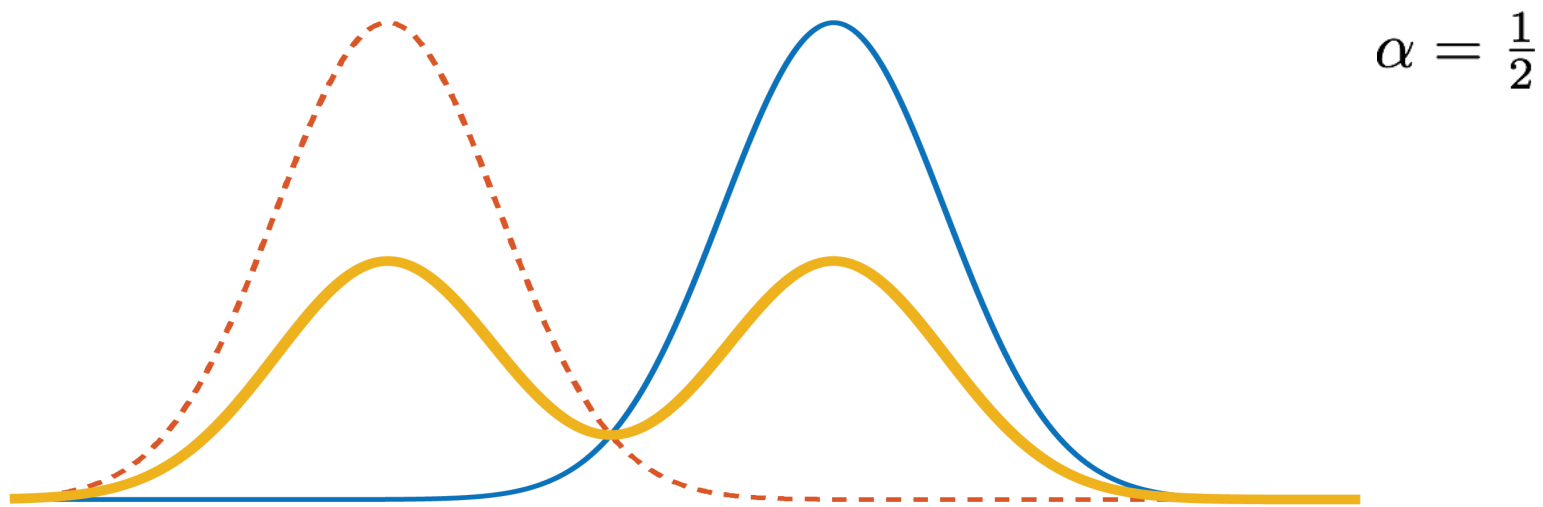
X_1 : sample from $p_1(x) = p(x|y = 1)$

X : sample from $p(x)$

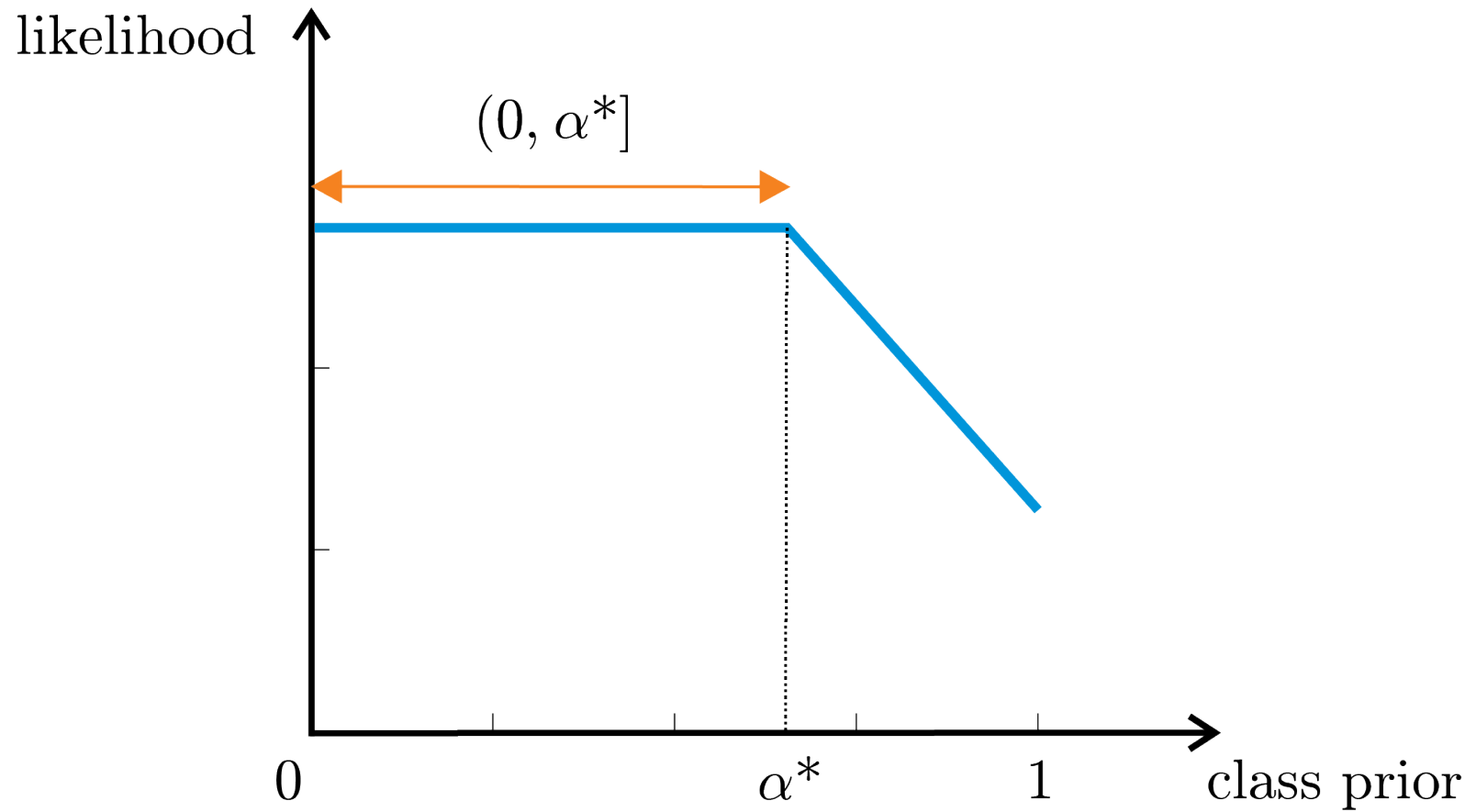
Goal: learn how x relates to y



IDENTIFIABILITY



ANTICIPATED LIKELIHOOD FUNCTION



ENZYMES: EXPERIMENTAL PROTOCOL

>sp|P04637|P53_HUMAN Cellular tumor antigen p53

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGP
DEAPRMPEAAPPVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAK
SVTCTYSPALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHE
RCSDSDGLAPPQHILIRVEGNLRVEYLDDRNTFRHS**VVVP**YEPPEVGSDCCTTIHNYMCNS
SCMGGMNRRIILTIIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELP
PGSTKRALPNNTSSSPQPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPG
GSRAHSSHLKSKKGQSTSRHKKLMFKTEGPDS

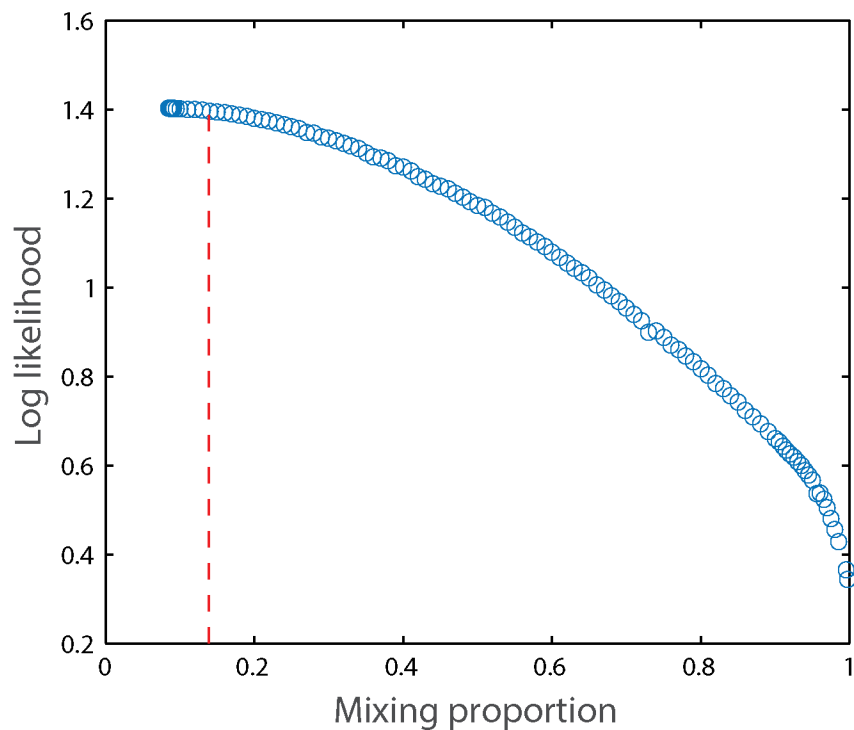


Develop an SVM predictor

- Linear kernel
- AUC \approx 75%
- Platt's correction

RESULTS: ENZYMES

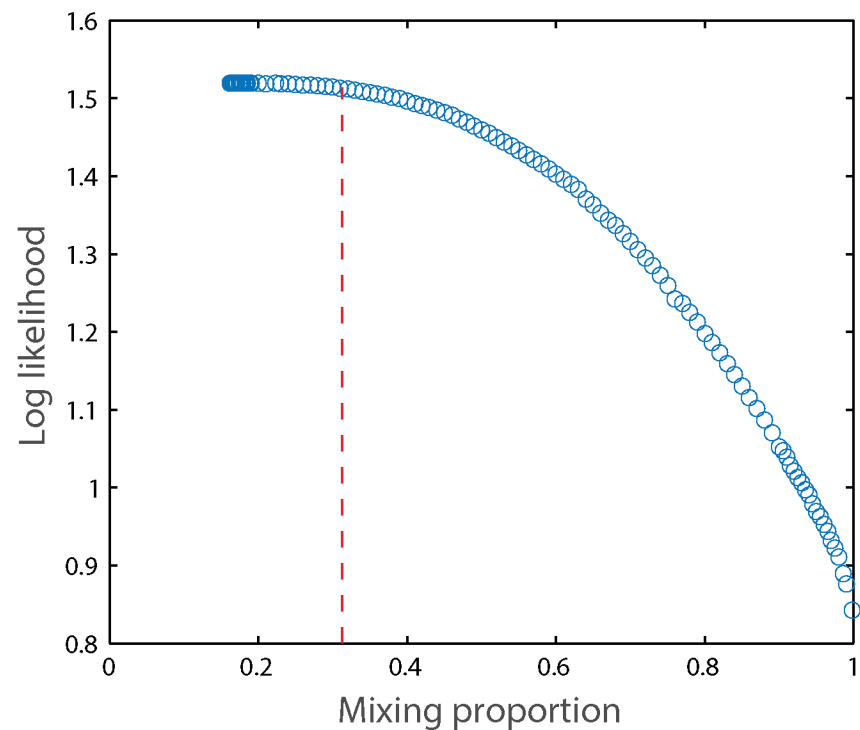
Yeast



We: 0.33
Elkan-Noto: 0.62

Charles: 0.45
Yuzhen: 0.40
Tuli: 0.85

E. coli

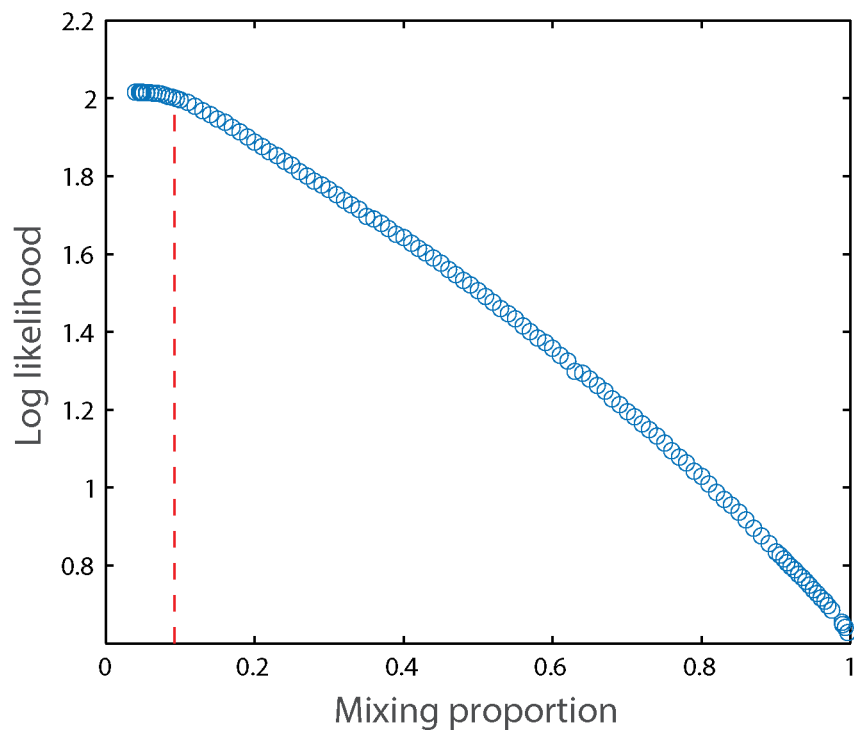


We: 0.49
Elkan-Noto: 0.76

Charles: 0.35
Yuzhen: 0.40
Tuli: 0.85

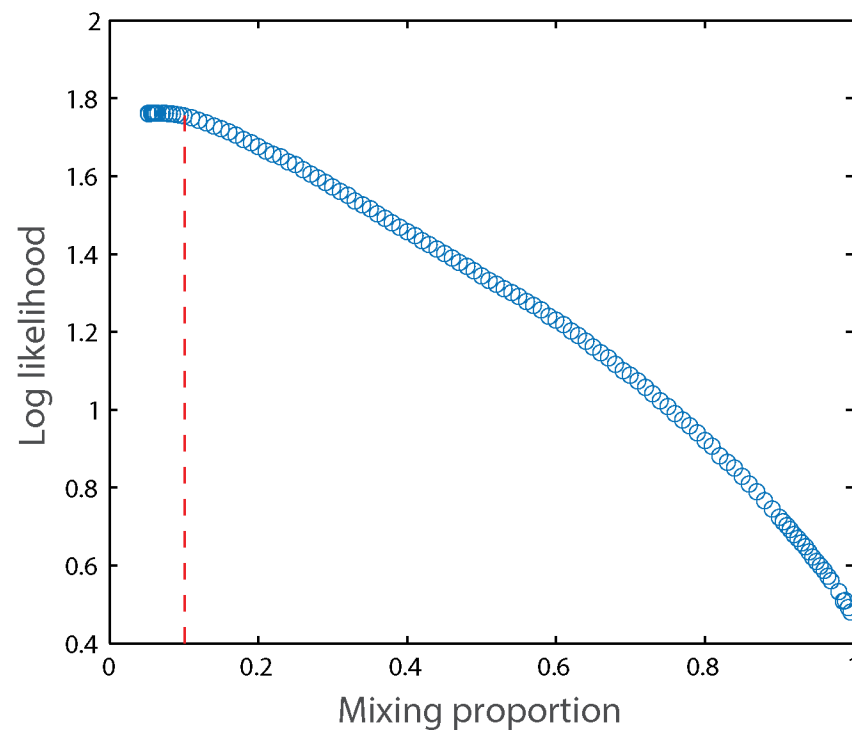
RESULTS: ENZYMES

Arabidopsis



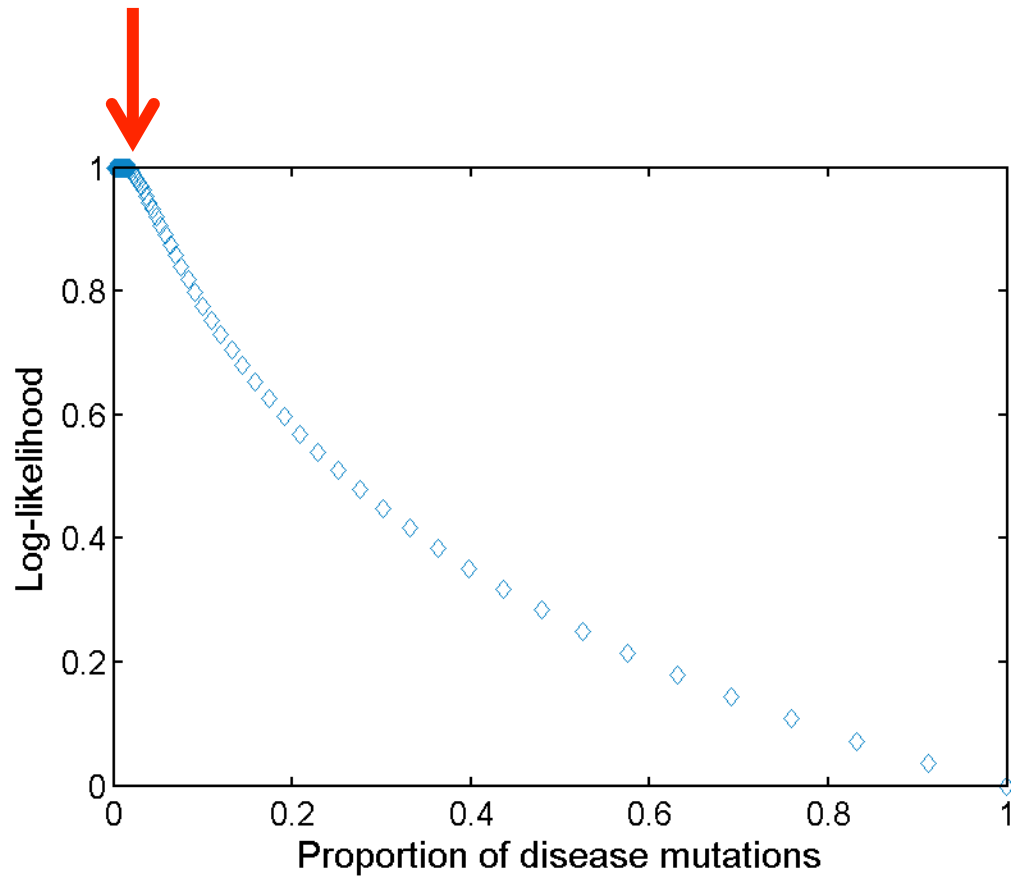
We: 0.20	Charles: 0.40
Elkan-Noto: 0.41	Yuzhen: 0.30
	Tuli: 0.85

Human

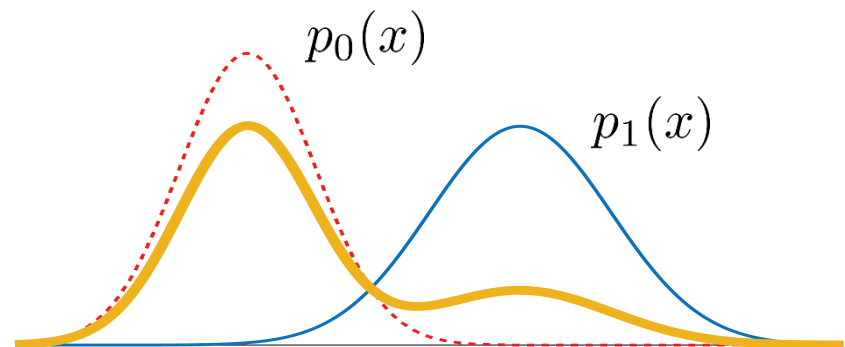


We: 0.21	Charles: 0.25
Elkan-Noto: 0.45	Yuzhen: 0.30
	Tuli: 0.85

DISEASE MUTATIONS IN HUMAN EXOME



About 1.5% of human nsSNPs are disease-associated!



PRECISION MEDICINE

01/20/2015



Precision Medicine

the science and practice of matching the best diagnostic, therapeutic and prevention strategies to promote health that are tailored to an individual's genetic, biological, behavioral and psychosocial characteristics

“So tonight, I’m launching a new Precision Medicine Initiative to bring us closer to curing diseases like cancer and diabetes, and to give all of us access to the personalized information we need to keep ourselves and our families healthier. We can do this.” – President Barack Obama.

All of UsSM Research Program



WHAT IS IT?

Precision medicine is a groundbreaking approach to disease prevention and treatment based on people's individual differences in environment, genes and lifestyle.

The *All of Us* Research Program will lay the foundation for using this approach in **clinical practice**.

WHAT ARE THE GOALS?

Engage a group of **1 million or more U.S. research participants** who will share biological samples, genetic data and diet/lifestyle information, all linked to their electronic health records. This data will allow researchers to develop more precise treatments for **many diseases and conditions**.

Pioneer a new model of research that emphasizes **engaged research participants, responsible data sharing and privacy protection**.



Research based on the cohort data will:

- Lay **scientific foundation** for precision medicine
- Help identify new ways to **treat and prevent disease**
- Test whether **mobile devices**, such as phones and tablets, can encourage healthy behaviors
- Help develop the **right drug** for the **right person** at the **right dose**

PRECISION MEDICINE

WHY NOW?

The **time is right** because:

We have a greater understanding of human genes

People are more engaged in healthcare and research



We have the tools to track health information and use large databases

Research technologies have improved



Follow the Program's progress and be one of the first to join this landmark effort.

www.nih.gov/AllofUs-Research-Program

GENOME SEQUENCING

Accurate Sequencing on the Complete Genomics Platform:
\$599 Human Exome Sequencing (100X) NEW
 Data Analysis & SNP Validation Included

<http://bgiamericas.com>

THE PROMISE OF GETTING PERSONAL

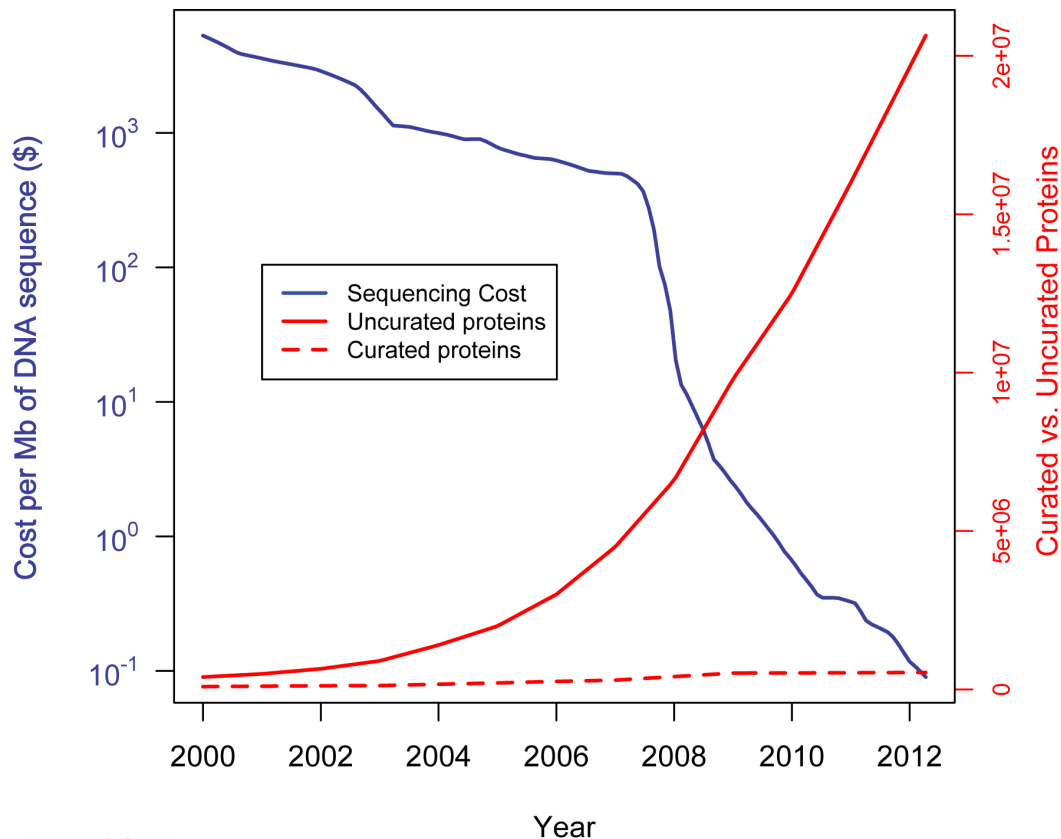
PERSONALIZED MEDICINE
 In 1993, sequencing an individual's genome took 13 years and cost three billion dollars. Today, the same task takes less than a week and costs only a few thousand dollars.

GENETIC MEDICAL RECORDS
 As the cost of genetic sequencing keeps dropping, more and more patients' medical records will include their genetic information.

TAILORED TREATMENT
 Informed by patients genetic sequences, treatments and preventative care will be all but certain to work better, with fewer trial-and-error prescriptions and a steep drop in the number of adverse drug reactions.

BUILDING MOMENTUM
 Universities, government entities, and private corporations alike are working to advance and implement personalized medicine, spurring new discoveries every day.

OPPORTUNITIES AHEAD
 Because medical spending accounts for almost 20% of U.S. consumer spending, developments in personalized medicine could save the health care industry billions of dollars a year and create huge new investment opportunities throughout the economy.



Subject	Platform	SNPs
USA individual (CV)	Sanger	3,213,400
German individual	Illumina	3,258,774
Estonian individual	SOLiD	3,482,975
USA individual (JW)	454	3,322,090
Irish individual	Illumina	3,125,825

HUMAN GENOME AND ITS IMPACT ON PHENOTYPE

WHAT IS THE MOLECULAR BASIS OF IT?



... TTTACCGAGC ...



... AGCATACCGA ...



... AGCATAGCGA ...

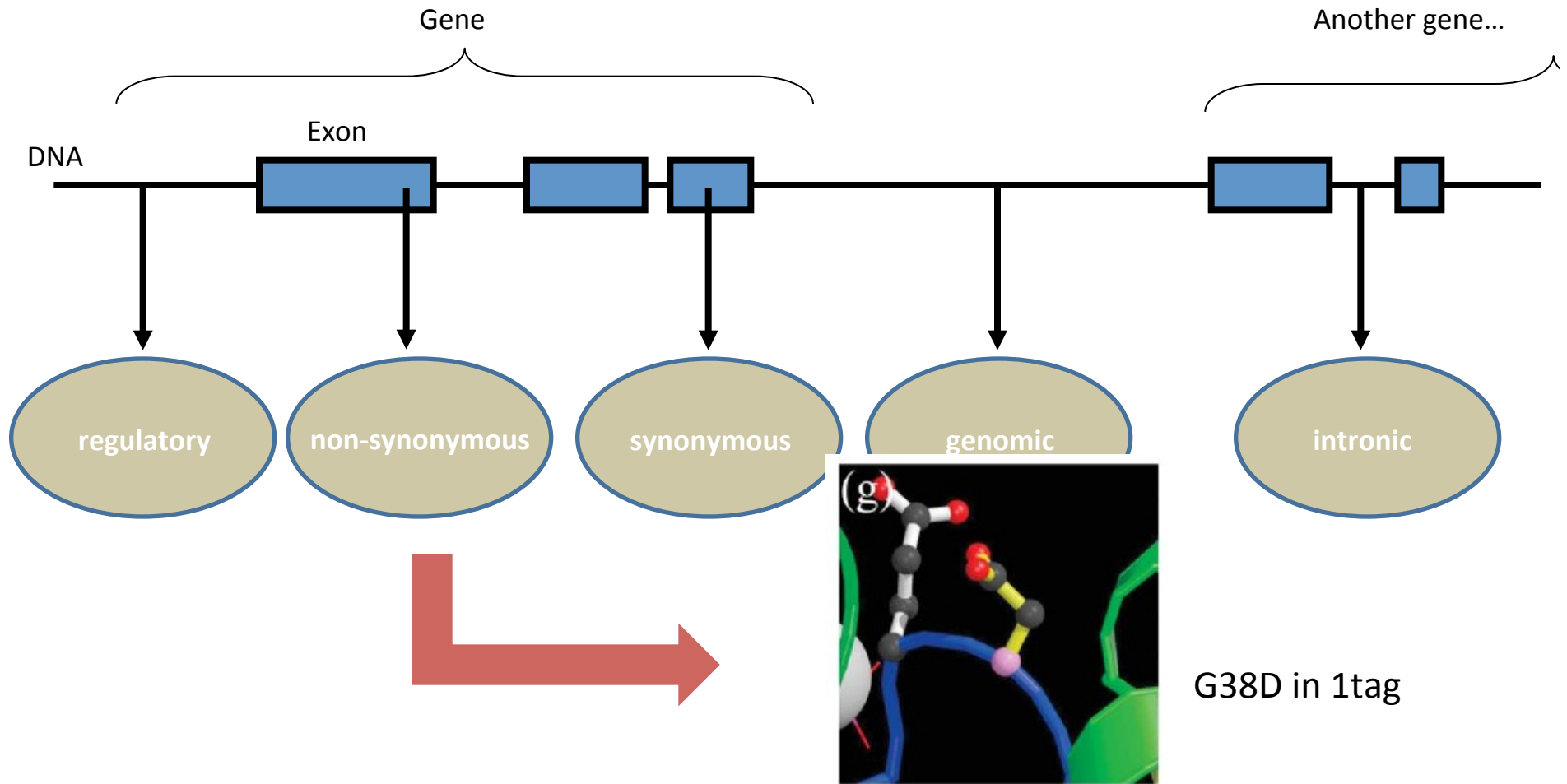


PHENOME

GENOME

BASE CHANGES RESULTING IN DIFFERENT PROTEINS

>40 million known unique sites of variation!

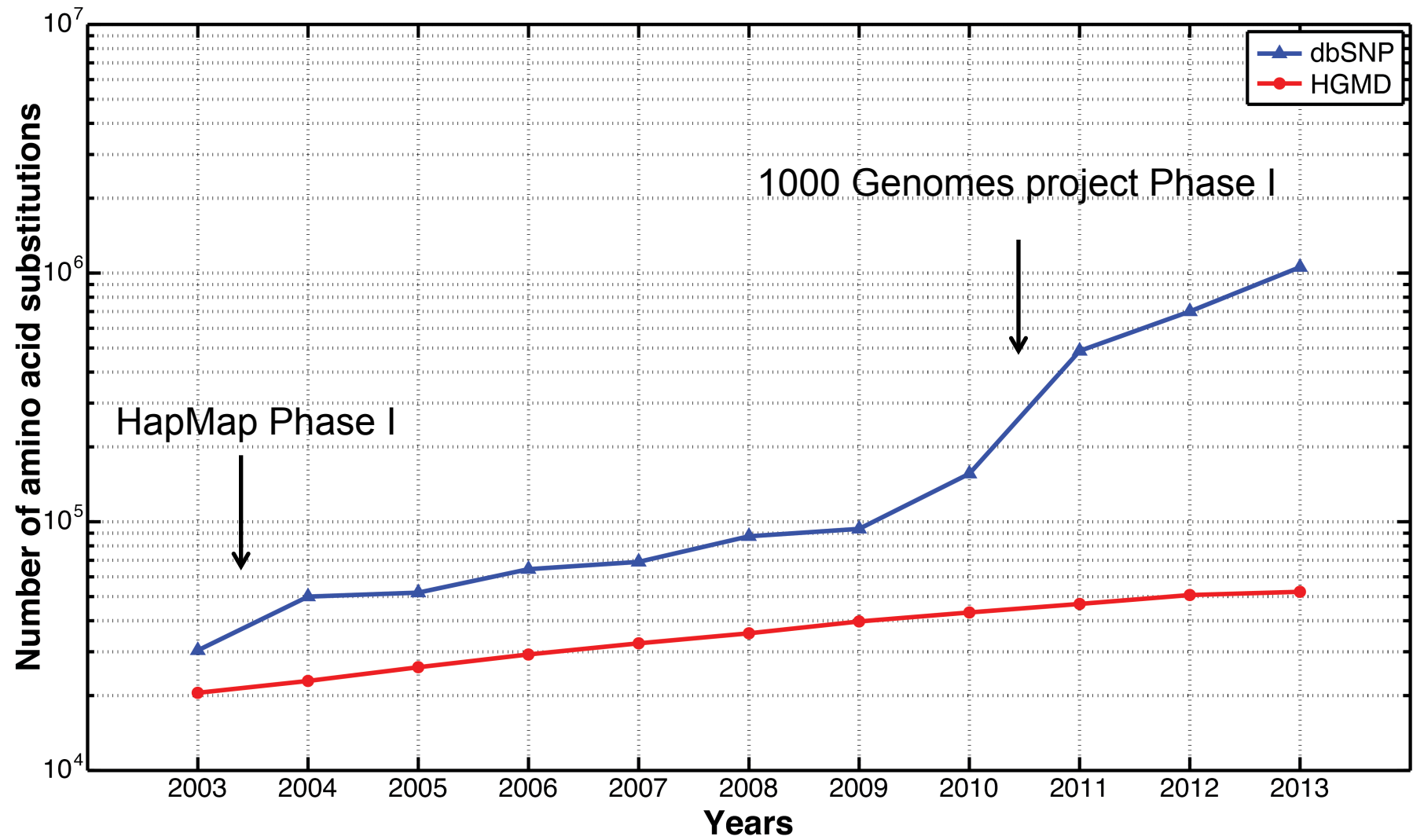


CURRENT TOOLS PREDICT EFFECTS OF VARIANTS

no.	SNP_ID	RefSNP_ID	Allele	Polymorph	MAF	SIFT	Gene	BlinkR	SeqR	LocalR	PolyPhen	SNP3DR	LS-SNPR	Panther	Pmut_pred
2	th1808	rs1126673	C/T	SNP	0.033	V374I	ADH4	neutral	neutral	neutral	neutral	neutral	neutral		neutral
3	th1819	rs1126671	T/C	SNP	0.000	I309V	ADH4	neutral	neutral	neutral	neutral	neutral	neutral		neutral
4	th1823	NEW	C/A	SNP	0.033	V158L	ADH4	neutral	neutral	neutral	neutral	neutral		neutral	damage
5	th509	rs4646422	G/A	SNP	0.063	G45D	CYP1A1	neutral	neutral	neutral	neutral	damage	neutral	damage	neutral
6	th512	rs2839942	C/G	SNP	0.016	T173R	CYP1A1	neutral	neutral	neutral	neutral	neutral	damage	neutral	damage
7	th514	rs1048943	A/G	SNP	0.234	I462V	CYP1A1	neutral	neutral	neutral	neutral	damage	neutral	neutral	neutral
8	th515	NEW	G/T	SNP	0.016	R511L	CYP1A1	damage	damage	damage	damage	damage			damage
9	th1845	NEW	C/G	SNP	0.031	F21L	CYP1A2	neutral	neutral	neutral	neutral			neutral	neutral
10	th1846	NEW	T/A	SNP	0.016	F125I	CYP1A2	damage	damage	damage	damage			damage	neutral
11	th1847	NEW	A/G	SNP	0.016	M180V	CYP1A2	neutral	damage	neutral	damage				neutral
12	th1850	NEW	G/A	SNP	0.016	G299S	CYP1A2	neutral	neutral	neutral	neutral			neutral	neutral
13	th1861	rs10012	C/G	SNP	0.281	R48G	CYP1B1	neutral	neutral	neutral	neutral	neutral	damage	neutral	neutral
14	th1862	NEW	C/T	SNP	0.016	S112L	CYP1B1	neutral	damage	neutral	neutral	neutral		neutral	neutral
15	th1863	rs1056827	G/T	SNP	0.281	A119S	CYP1B1	neutral	damage	neutral	neutral	neutral	neutral	neutral	neutral
16	th1865	rs1056836	G/C	SNP	0.083	V432L	CYP1B1	neutral	neutral	damage	neutral	damage	neutral	neutral	neutral
17	th1872	rs5031017	C/A	SNP	0.234	G479V	CYP2A6	damage	damage	damage	neutral	damage	damage		damage
18	th1876	rs2839944	A/G	SNP	0.031	S224P	CYP2A6	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
19	th1887	rs3758581	G/A	SNP	0.047	I331V	CYP2C19	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
20	th1921	rs2837175	T/C	SNP	0.031	L293P	CYP3A4	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
21	th1926	NEW	C/T	SNP	0.017	R35W	CYP4B1	neutral	neutral	neutral	damage	neutral		damage	damage
22	th1933	rs3215983	A/T_	Deletion	0.234	D294G	CYP4B1	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
23	th1934	rs2297810	G/A	SNP	0.234	M331I	CYP4B1	neutral	damage	damage	neutral	damage	neutral	neutral	neutral
24	th1935	rs4646491	C/T	SNP	0.250	R340C	CYP4B1	damage	damage	damage	damage	damage	damage	damage	damage
25	th1937	rs2297809	C/T	SNP	0.229	R375C	CYP4B1	damage	damage	damage	damage	damage	damage	damage	damage
26	th1943	rs3738046	G/C	SNP	0.047	R43T	EPHX1	neutral	neutral	neutral	neutral	neutral	neutral	neutral	damage
27	th1949	rs1051740	T/C	SNP	0.438	Y113H	EPHX1	damage	damage	damage	damage	damage	damage	damage	neutral
28	th1952	NEW	C/T	SNP	0.016	H139Y	EPHX1	neutral	neutral	neutral	damage	damage		neutral	neutral
29	th1953	rs2234922	A/G	SNP	0.141	H139R	EPHX1	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
30	th1961	rs947894	A/G	SNP	0.266	I105V	GSTP1	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
31	th1026	rs1805158	C/T	SNP	0.016	R64W	NAT2	damage	damage	damage	neutral	damage		damage	damage
32	th1028	rs1801280	T/C	SNP	0.094	I114T	NAT2	neutral	damage	neutral	neutral	damage	damage	neutral	neutral
33	th1030	rs1799930	G/A	SNP	0.344	R197Q	NAT2	neutral	damage	neutral	damage	neutral	damage	neutral	damage
34	th1031	rs1208	G/A	SNP	0.156	R268K	NAT2	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
35	th1032	rs1799931	G/A	SNP	0.141	G286E	NAT2	neutral	neutral	neutral	neutral	neutral	neutral	neutral	damage
36	th2245	NEW	C/A	SNP	0.071	L71I	NNMT	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
37	th2296	rs2637445	A/G	SNP	0.479	F247L	SULT1A1	neutral	neutral	neutral	neutral	damage	neutral	neutral	neutral
38	th2297	rs1801030	C/T	SNP	0.000	V223M	SULT1A1	damage	damage	neutral	neutral	neutral	neutral	neutral	neutral
39	th2313	rs27742	T/C	SNP	0.000	E282K	SULT1A2	neutral	neutral	neutral	neutral	neutral	neutral	neutral	damage
40	th2314	NEW	T/A	SNP	0.016	K258N	SULT1A2	damage	damage	damage	neutral	damage		damage	neutral
41	th2316	rs1059491	T/G	SNP	0.125	N235T	SULT1A2	damage	damage	damage	damage	damage	damage	neutral	neutral
42	th2331	rs1703610	T/G	SNP	0.056	S255A	SULT1C1	neutral	neutral	neutral	neutral	neutral	neutral	neutral	neutral
43	th2239	rs1109826	A/G	SNP	0.000	M368I	UGT8	neutral	neutral	neutral	neutral	neutral	neutral	damage	neutral

When applied to 43 nsSNPs of 18 drug related genes from the Thai SNP resequencing project there were strong correlations

GROWTH OF DATA



Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology

Sue Richards, PhD¹, Nazneen Aziz, PhD^{2,16}, Sherri Bale, PhD³, David Bick, MD⁴, Soma Das, PhD⁵, Julie Gastier-Foster, PhD^{6,7,8}, Wayne W. Grody, MD, PhD^{9,10,11}, Madhuri Hegde, PhD¹², Elaine Lyon, PhD¹³, Elaine Spector, PhD¹⁴, Karl Voelkerding, MD¹³ and Heidi L. Rehm, PhD¹⁵;
on behalf of the ACMG Laboratory Quality Assurance Committee

Disclaimer: These ACMG Standards and Guidelines were developed primarily as an educational resource for clinical laboratory geneticists to help them provide quality clinical laboratory services. Adherence to these standards and guidelines is voluntary and does not necessarily assure a successful medical outcome. These Standards and Guidelines should not be considered inclusive of all proper procedures and tests or exclusive of other procedures and tests that are reasonably directed to obtaining the same results. In determining the propriety of any specific procedure or test, the clinical laboratory geneticist should apply his or

Table 2 In silico predictive algorithms

Category	Name	Website	Basis
Missense prediction	ConSurf	http://consurftest.tau.ac.il	Evolutionary conservation
	FATHMM	http://fathmm.biocompute.org.uk	Evolutionary conservation
	MutationAssessor	http://mutationassessor.org	Evolutionary conservation
	PANTHER	http://www.pantherdb.org/tools/csnpscoreForm.jsp	Evolutionary conservation
	PhD-SNP	http://snps.biofold.org/phd-snp/phd-snp.html	Evolutionary conservation
	SIFT	http://sift.jcvi.org	Evolutionary conservation
	SNPs&GO	http://snps-and-go.biocomp.unibo.it/snps-and-go	Protein structure/function
	Align GVGD	http://agvgd.iarc.fr/agvgd_input.php	Protein structure/function and evolutionary conservation
	MAPP	http://mendel.stanford.edu/SidowLab/downloads/MAPP/index.html	Protein structure/function and evolutionary conservation
	MutationTaster	http://www.mutationtaster.org	Protein structure/function and evolutionary conservation
	MutPred	http://mutpred.mutdb.org	Protein structure/function and evolutionary conservation
	PolyPhen-2	http://genetics.bwh.harvard.edu/pph2	Protein structure/function and evolutionary conservation
	PROVEAN	http://provean.jcvi.org/index.php	Alignment and measurement of similarity between variant sequence and protein sequence homolog
	nsSNPAnalyzer	http://snpanalyzer.uthsc.edu	Multiple sequence alignment and protein structure analysis
Condel	http://bg.upf.edu/fannssdb/	Combines SIFT, PolyPhen-2, and MutationAssessor	
CADD	http://cadd.gs.washington.edu	Contrasts annotations of fixed/nearly fixed derived alleles in humans with simulated variants	

TUMOR BOARD

Person: 68 year old woman

Cancer type: colon cancer, metastatic

Mutations:

KRAS, C27F

BRCA1, H57R

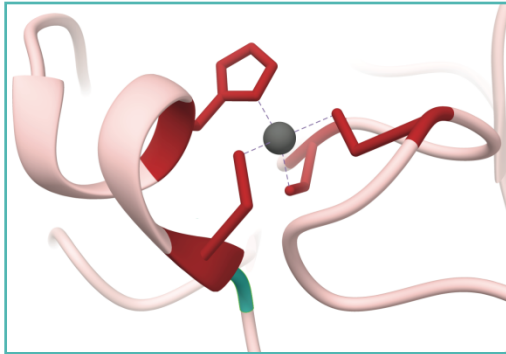
TP53, T98*

Treatment options:

- clinical trial at MD Anderson
- continue with chemotherapy
- treat with new drug for breast cancer



MOLECULAR CONSEQUENCES ON P53

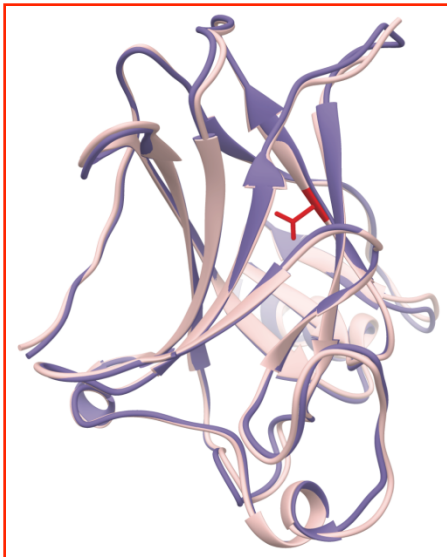


R175H: Metal-binding

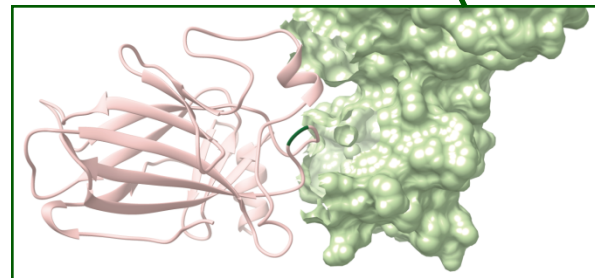


K120R: Acetylation

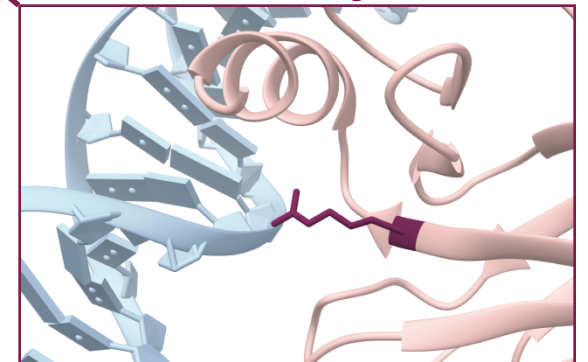
V143A: Stability



G245S: Protein-binding



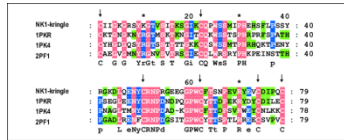
R273H: DNA-binding



p53 – tumor suppressor protein

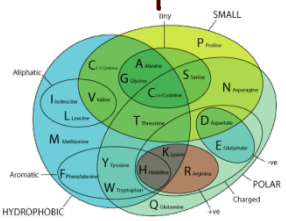
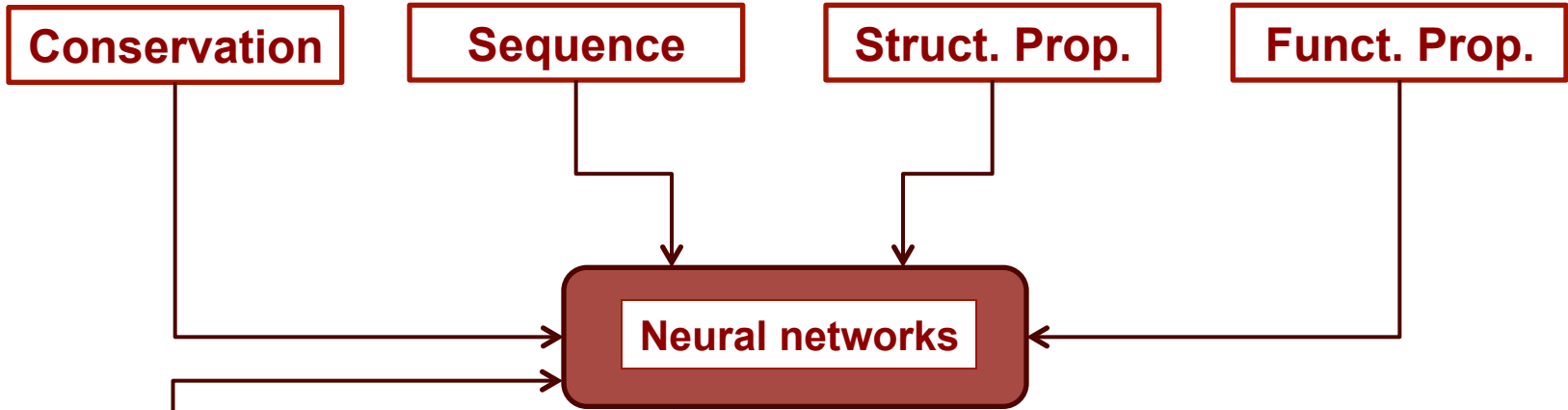
PDB structures: 2ybg, 2j1w, 1ycs and 1tup

MUTPRED 2.0



```
>sp|P04637|P53_HUMAN
MEBPDQDPSEVPEPLSQETFSDLMLKLLFENNVLSPFSPQAMDDLMSFDDIEQWFTEDPGP
DEAPRMPPEAAPRVAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAPAP
SVICTYSPALNKMFCQLAKTCFVQLWVDSTFEPGTRVRAAIYKQSQHMTEVVRRCFHHE
RCSDSGLAPFQHLIRVEGNLRVEYLDLDRNFRHSVVVPEFPEVSGDCTTIHYNYMNS
```

$$P(\text{loss of } p \text{ at } s_i | x_{jy}) = P(s_i = s_i^p | s) \cdot (1 - P(s_i = s_i^p | s_{xjy})) \quad P(\text{loss of } p \text{ at } s_i | x_{jy}) = P(s_i = s_i^p | s) \cdot (1 - P(s_i = s_i^p | s_{xjy}))$$



Physicochemical

	C	S	T	P	A	G	N	D	E	Q	R	K	M	I	L	V	F	Y	W
C	0																		
S	-1	4																	
T	-1	1	5																
P	-3	-1	-1	7															
A	0	0	-1	-4	0														
G	-3	0	-2	-2	0	6													
N	-1	0	-2	-2	0	6	6												
D	-3	0	-1	-1	-2	-1	1	6											
E	-4	0	-1	-1	-1	-2	0	2	5										
Q	-3	0	-1	-1	-2	-2	0	0	2	5									
H	-3	-1	-2	-2	-2	-2	-1	-1	0	0	8								
R	-3	-1	-2	-1	-2	-2	0	0	-2	0	1	0	5						
K	-3	0	-1	-1	-1	-2	0	-1	1	-1	5	2	5						
M	-1	-1	-2	-1	-3	-2	-2	0	-2	-1	-1	5							
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	1	4					
V	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	2	2	4					
L	-2	0	-2	0	-3	-3	-2	-2	-3	-2	-3	-2	1	3	1	4			
F	-2	-2	-4	-2	-3	-3	-3	-3	-1	-3	-3	0	0	-1	6				
Y	-2	-2	-3	-2	-3	-2	-1	-2	-2	-1	-1	-1	-1	3	7				
W	-3	-2	-4	-3	-2	-4	-4	-3	-2	-3	-3	-3	-2	-3	1	2	11		

Substitution matrices

Neural network ensemble:

- Z-score normalized and PCA'd
- 30 feed-forward networks
- bootstrap aggregating, balanced training
- trained using resilient propagation

GAIN OF CATALYTIC ACTIVITY CAUSES DISEASE

- a member of the proteinase K sub-family of subtilases that reduces the number of LDL receptors in liver through a posttranscriptional mechanism.
- D374Y leads to a 10-fold increase in catalytic activity that causes hypercholesterolemia.

Catalytic residues of PCSK9 (2qtw)

