

Natural Learning

Learning Machine,
Inspired by Humans,
for Humans



Hadi Fanaee-T
Associate Professor of Machine Learning
School of Information Technology
Halmstad University, Sweden
Email: hadi.fanaee@hh.se

[Designed by BootstrapMade](#)

Curse of “black box” Modelling

- Deep Neural Networks (DNNs) have achieved **state-of-the-art performance** on benchmark datasets and real-world applications.
 - Nobody knows how and why they work (**Transparency**)
 - Unable to explain their predictions comprehensibly to humans (**Decision’s Explainability**)
 - Unable to provide justifications or reasoning for their decisions (**Decision’s Interpretability**)

You cannot say, ‘I’ll do open-heart surgery because the neural network said so.’ You have to have a very good reason.”



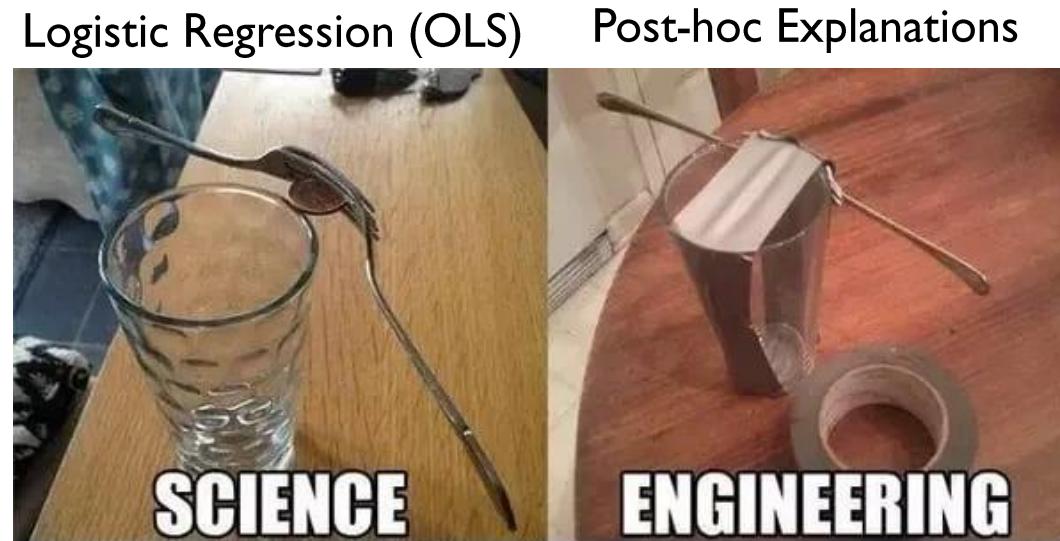
Christos Faloutsos, professor
of computer science, CMU

Post-hoc Explanations: An Engineering Remedy for an Engineering Solution

SHAP (SHapley Additive exPlanations): **global** feature importance

LIME (Local Interpretable Model-agnostic Explanations): **local** feature importance

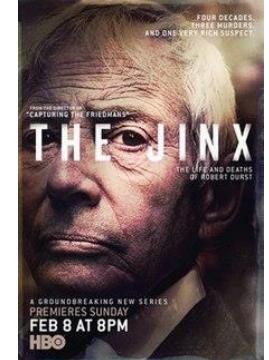
There is a narrow difference between **explaining** something and **justifying** it.



Source: engineering.com

Analogy of post hoc explanations in Real-life (I)

- On October 9, 2001, Robert Durst (**Black-box model**), murdered his neighbor Morris Black.
- During the trial, Durst's defense team (**Post-hoc explainer**) argued that he acted in **self-defense**.
- In November 2003, Durst was acquitted of murder charges.
- In the Netflix show "The Jinx," there's a scene where Durst talks to himself in a bathroom after an interview while still being recorded. He confessed that has killed his neighbor and two others.
- Both **murder** and **self-defense** are convincing stories.
- But does that mean that the explanation by the defense team reflected the truth?
- Only the killer (**Black-box model**) knows it.



Analogy of post hoc explanations in Real-life (2)

- A Black-box Model rejects a loan application solely based on the subject's **race**
- Post-hoc explanation can justify that **income** and **employment** status were the critical factors for the decision.
- The **connection** between Post-hoc explanation and the model's actual behavior can **never be proven**
 - **illusion of explainability**
- Post-hoc explanations can be even **more dangerous** than black-box models.
 - Sometimes, they can **cover up** the mistakes of black-box models by giving **believable reasons**.



Post-hoc Explanations are inherently wrong!

- Another good argument: If the post-hoc explanation fully matched the original model, why would we need the original model after all?

Cynthia Rudin
Professor of Computer
Science, Duke University

nature machine intelligence

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾ [Subscribe](#)

[nature](#) > [nature machine intelligence](#) > [perspectives](#) > [article](#)

Perspective | Published: 13 May 2019

Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead

[Cynthia Rudin](#)

[Nature Machine Intelligence](#) 1, 206–215 (2019) | [Cite this article](#)

73k Accesses | 2977 Citations | 502 Altmetric | [Metrics](#)

New regulations may restrict use of black-box models



The screenshot shows a web browser displaying a page from the European Parliament's Topics section. The URL in the address bar is europarl.europa.eu/topics/en/article/20230601STO93804/eu-ai-act-first-regulation-on-artificial-int.... The page header includes links for News, Topics, MEPs, About Parliament, Plenary, Committees, Delegations, Elections, and a search bar. The main navigation menu below the header lists Digital, Energy, Gender equality, Climate and environment, Circular economy, and All topics. The breadcrumb navigation on the left indicates the path: Topics > Digital > Artificial intelligence > EU AI Act: first regulation on artificial intelligence.

EU AI Act: first regulation on artificial intelligence

The use of artificial intelligence in the EU will be regulated by the AI Act, the world's first comprehensive AI law. Find out how it will protect you.

Published: 08-06-2023 • Last updated: 19-12-2023 - 11:45

New regulations may restrict use of black-box models

- EU AI Act
 - Regulating high-risk AI applications used in critical infrastructure, such as transportation and **healthcare**, as well as those with potential risks to fundamental rights, safety, or other public interests.
 - Set of requirements to ensure their safety, **transparency**, and **accountability**.

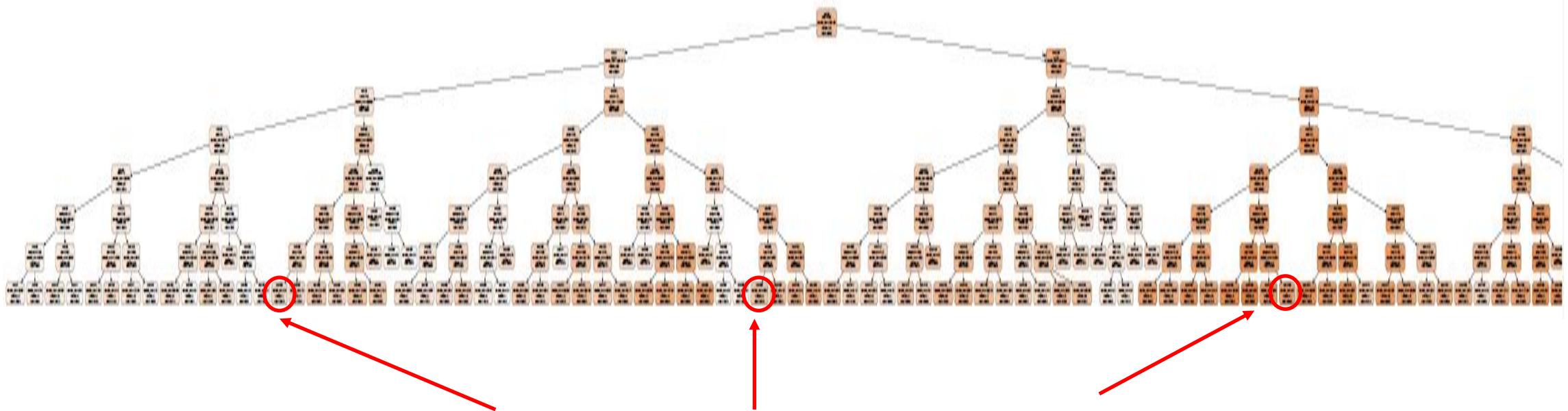
An illusion of black-box models' universal superiority

- Real-world datasets typically contain between **8.0% and 38.5% label noise** (Semenova, et al., 2023)
- (Semenova, et al., 2023) provided **theoretical evidence** that in noisy datasets, such as datasets about humans like healthcare, criminal justice, and finance, **simple, interpretable classifiers should perform as well as black-box models.**

OK, but if we exclude black-box, what options do we have?

- Logistic Regression
 - Good level of interpretability and explainability
 - No built-in mechanism to deal with **noisy features, curse of dimensionality, and multicollinearity**
 - Poor performance with high-dimensional datasets
- Decision Trees
 - One of the most favored options when it comes to interpretability and explainability
 - Robust against the curse of dimensionality, irrelevant features, and noisy samples
 - Decision Trees are transparent, but are they explainable and interpretable?

Illusion of Explainability of Decision Trees

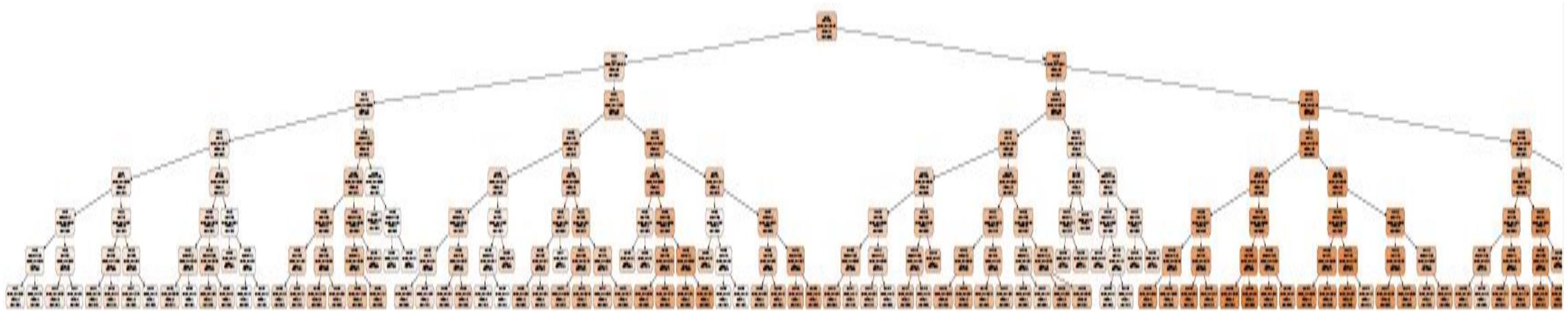


Can a **universal rule** explain the decisions for these people?

Opposed to Logistic regression, decision trees cannot provide a **global explanation** for the decisions.

Local explanations have a **limited value** if they cannot be **generalized**!

Illusion of Interpretability of Decision Trees



Can humans process such a huge amount of information?

Illusion of Interpretability of Decision Trees

- Humans can store only **few meaningful items** in the **working memory**



The Magical Mystery Four: How Is Working Memory Capacity Limited, and Why?

Nelson Cowan

University of Missouri

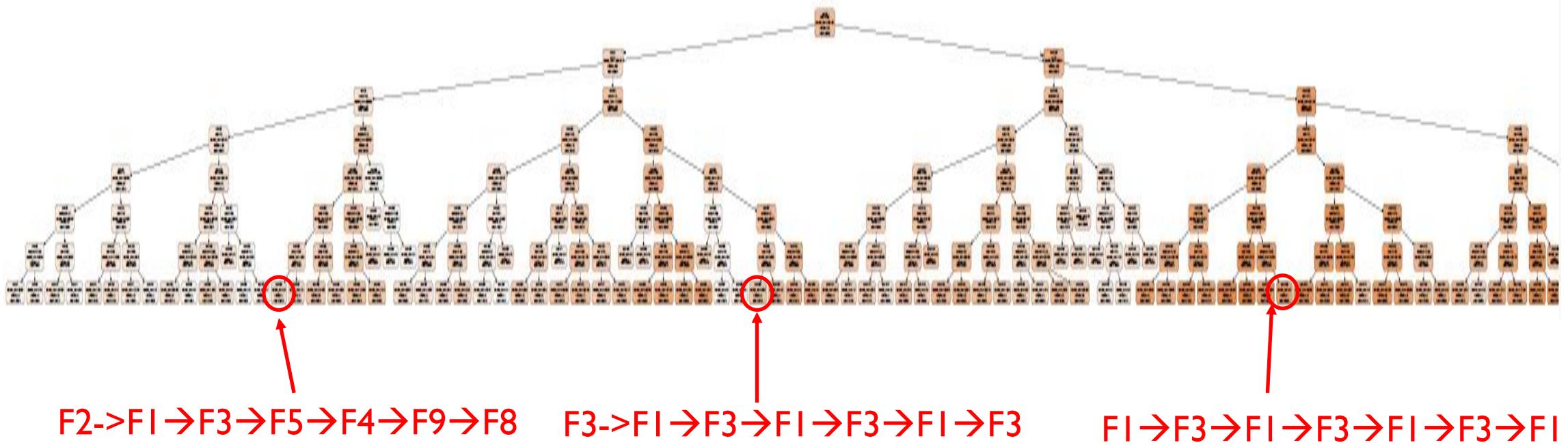
Current Directions in Psychological Science
19(1) 51-57
© The Author(s) 2010
Reprints and permission:
sagepub.com/journalsPermissions.nav
DOI: 10.1177/0963721409359277
<http://cdps.sagepub.com>



Abstract

Working memory storage capacity is important because cognitive tasks can be completed only with sufficient ability to hold information as it is processed. The ability to repeat information depends on task demands but can be distinguished from a more constant, underlying mechanism: a central memory store limited to 3 to 5 meaningful items for young adults. I discuss why this central limit is important, how it can be observed, how it differs among individuals, and why it may exist.

Illusion of Interpretability of Decision Trees

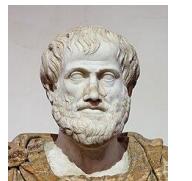


For each decision, different combinations of features are used.

It is **impossible** to infer the **actual contribution of features** at the **global level**.

Analogy in Law

- We don't have **numerous versions of laws tailored to different individuals**; rather, there exists a **single universal law** that applies to **everyone**.



Underlying Philosophy of Decision Trees

Aristotle
(384-322 B.C.)

- The challenges associated with decision trees stem from their underlying philosophy, which is rooted in **Aristotle's categorization theory**
 - Humans use rule-based explanations to categorize concepts.
- Extensive Research in cognitive psychology in 1970s indicated shortcomings in this model, suggesting that **people likely do not rely on rule-based definitions** when categorizing objects.



Natural Categories (Prototype Theory)

Eleanor Rosch
Professor of Psychology,
University of California,
Berkeley

- People **categorize** objects and concepts based on their **similarity to a prototype**

Furniture Prototype



Object

More
Similar



Electronics Prototype

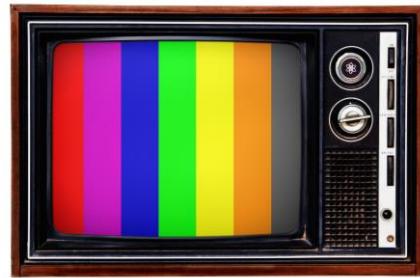


Image Source: <https://slideplayer.com/slide/9817067/>

Image Source: <https://facts.net/who-invented-color-tv/>

Characteristics of Prototype (I): Typicality

- Prototype is **the most typical** or central example of a category

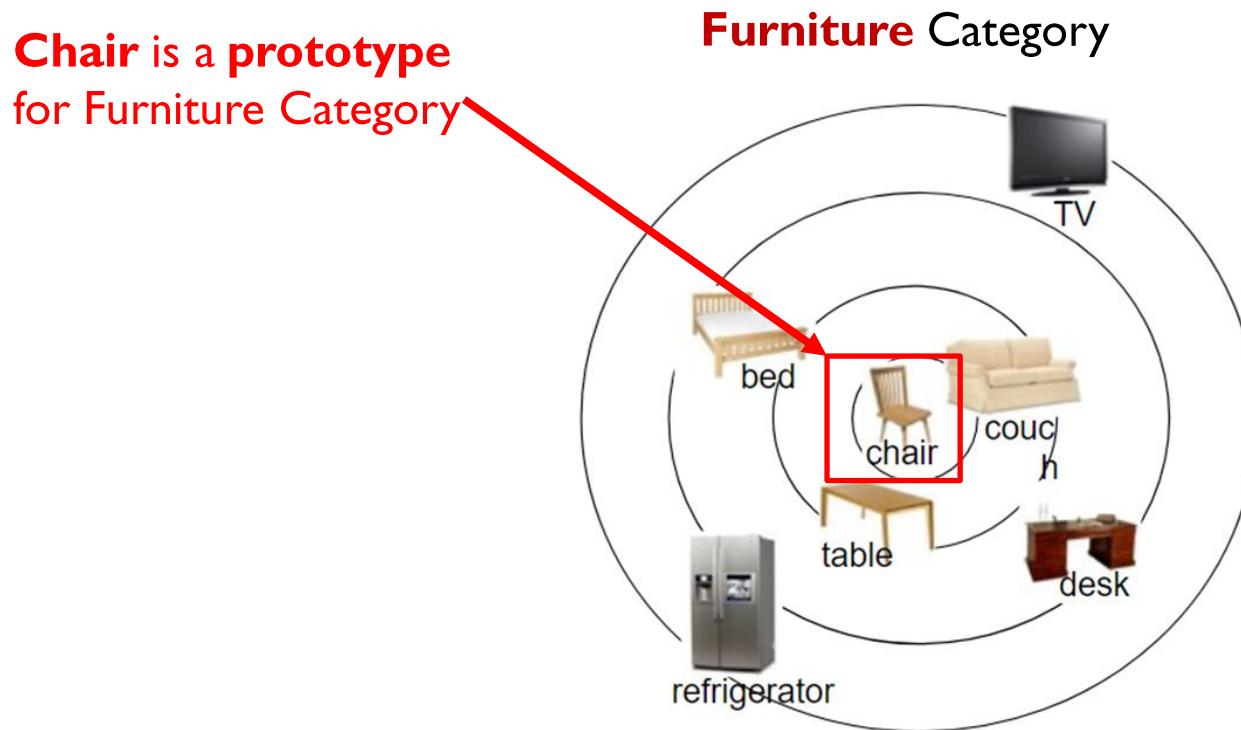


Image Source: <https://www.slideserve.com/louise/psy-369-psycholinguistics>

Characteristics of Prototype (2): Core Features

- **Core features:** central features of the prototype that are typically shared by most, if not all, instances within the category and are necessary for distinguishing the category from other categories.
- **Saliency Features:** prominent within a category but not necessary for distinguishing the category from other categories.
- **Peripheral features:** not essential or central to identity of a category



Characteristics of Prototype (3): Generalizability

- Features of a prototype should be **generalizable to other members of the category**, even if those members differ in some respects from the prototype itself

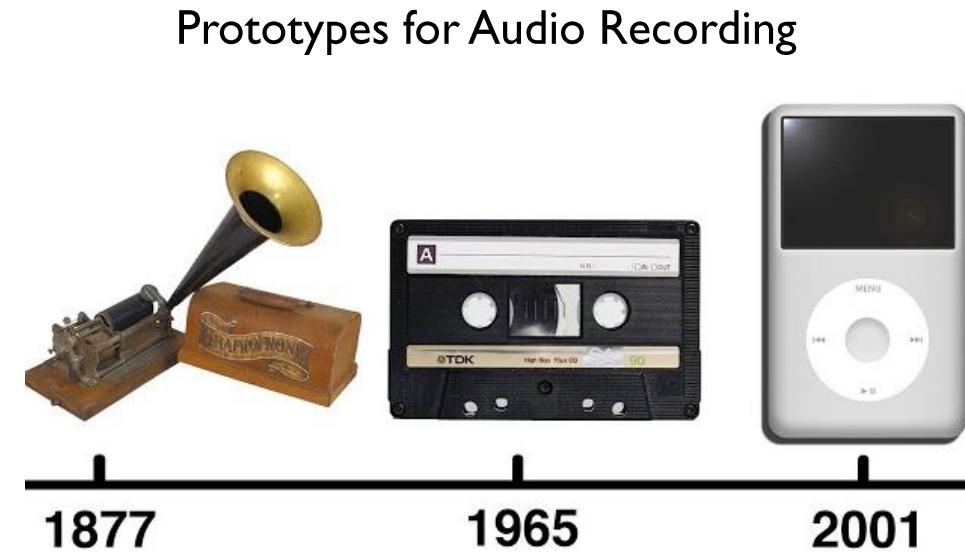


Image Source: mecox.com



Characteristics of Prototype (4): Flexibility

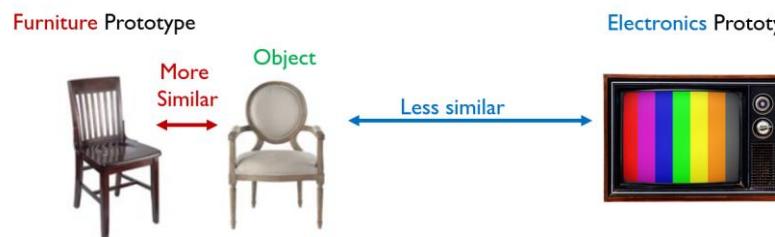
- Prototypes are not fixed entities; **they can change** based on new experiences.



Source: <https://www.youtube.com/watch?app=desktop&v=5Pl2rsLhhwQ>

Translation to Machine Learning Language

Prototype Theory	Translation to Machine Learning Language
Typicality	Each class is represented by only one single prototype
Core Features	Prototypes have sparse features
Generalizability	Prototype features are generalizable to samples of class
Flexibility	Learning prototypes is an incremental process

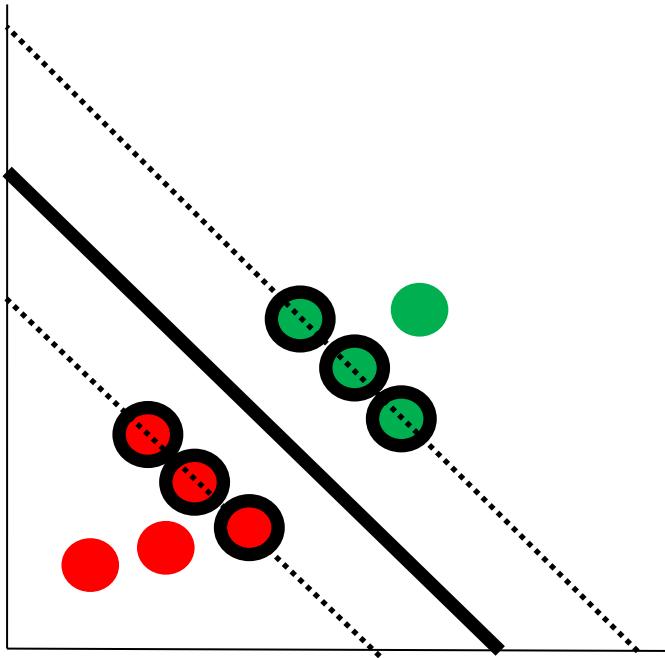


Classification Rule

If the **test sample** is closer to **class 0's prototype** than **class 1's prototype**, it is classified as **0**; otherwise, it is classified as **1**.

Models that are called prototype-based

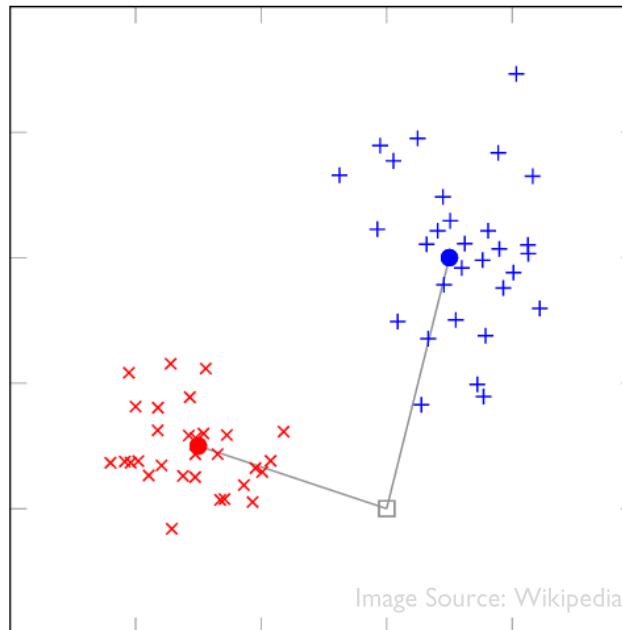
Support Vector Machines (SVM)



Boser, Bernhard E., Isabelle M. Guyon, and Vladimir N. Vapnik. "A training algorithm for optimal margin classifiers." Proceedings of the fifth annual workshop on Computational learning theory. 1992.

Popular in Machine Learning

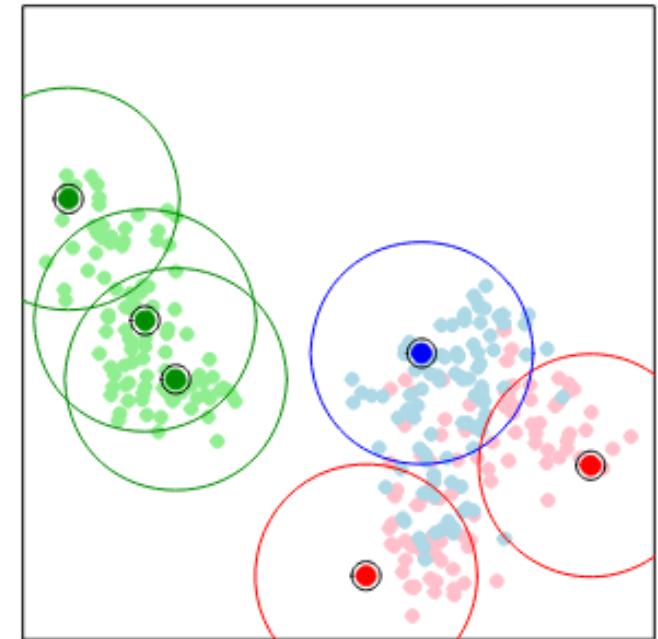
Nearest Centroid Classifier (NCC)



Manning, Christopher; Raghavan, Prabhakar; Schütze, Hinrich (2008). "Vector space classification". Introduction to Information Retrieval. Cambridge University Press.

Popular in information retrieval

Prototype Selection (PS)



Jacob Bien and Robert Tibshirani. Prototype selection for interpretable classification. arXiv preprint arXiv:1202.5933, 2012.

Unsupervised discovery of multiple prototypes from classes (not a direct classifier, a pre-processing method)

Do they match the properties mentioned in prototype theory?

Prototype Theory	Translation to Machine Learning Language	SVM	NCC	PS
Typicality	Each class is represented by only one single prototype	✗	✓	✗
Core Features	Prototypes have sparse features	✗	✗	✗
Generalizability	Prototype features are generalizable to samples of class	✗	✗	✗
Flexibility	Learning prototypes is an incremental process	✗	✗	✗

Can we build a Classifier based on Prototype Theory?

- According to prototype theory, our prototypes of interest should be:
 - **A pair of samples** (One sample from **class 0** and one sample from **class 1**)
- We should be able to **correctly classify all (or the majority of) samples** based on the following rule:
 - If the test sample is closer to **class 0's** prototype than **class 1's** prototype, it is classified as 0; otherwise, it is classified as 1.
- We don't know what are the core features of the prototypes, but we know that they are the **most generalizable and the sparsest** among all samples.

Brute-Force Approach: Naïve Prototype Classifier

- We can define this classification task as a pure **cross-validation** problem.
- We test all possible pairs of samples from 0 and 1 classes, and with all subsets of features (from length 1 to p) to see which one generalizes best to all samples: X is closer to 0's prototype than 1's prototype, labeled 0 otherwise 1
- We pick the pair with the **lowest error** and **number of features**.

	F =1	F =2	F =3	F =4
(X1,X3,F1)	(X1,X3,F1,F2)	(X1,X3,F2,F4)	(X1,X3, F1,F2,F3}	(X1,X3, F1,F2,F3,F4}
(X1,X4,F1)	(X1,X4,F1,F2)	(X1,X4,F2,F4)	(X1,X4, F1,F2,F3)	(X1,X4, F1,F2,F3,F4)
(X2,X3,F1)	(X2,X3,F1,F2)	(X2,X3,F2,F4)	(X2,X3, F1,F2,F3)	(X2,X3, F1,F2,F3,F4)
(X2,X4,F1)	(X2,X4,F1,F2)	(X2,X4,F2,F4)	(X2,X4, F1,F2,F3)	(X2,X4, F1,F2,F3,F4)
(X1,X3,F2)	(X1,X3,F1,F3)	(X1,X3,F3,F4)	(X1,X3, F1,F2,F4}	
(X1,X4,F2)	(X1,X4,F1,F3)	(X1,X4,F3,F4)	(X1,X4, F1,F2,F4)	
(X2,X3,F2)	(X2,X3,F1,F3)	(X2,X3,F3,F4)	(X2,X3, F1,F2,F4)	
(X2,X4,F2)	(X2,X4,F1,F3)	(X2,X4,F3,F4)	(X2,X4, F1,F2,F4)	
(X1,X3,F3)	(X1,X3,F1,F4)		(X1,X3, F1,F3,F4}	
(X1,X4,F3)	(X1,X4,F1,F4)		(X1,X4, F1,F3,F4)	
(X2,X3,F3)	(X2,X3,F1,F4)		(X2,X3, F1,F3,F4)	
(X2,X4,F3)	(X2,X4,F1,F4)		(X2,X4, F1,F3,F4)	
(X1,X3,F4)	(X1,X3,F2,F3)		(X1,X3, F2,F3,F4}	
(X1,X4,F4)	(X1,X4,F2,F3)		(X1,X4, F2,F3,F4)	
(X2,X3,F4)	(X2,X3,F2,F3)		(X2,X3, F2,F3,F4)	
(X2,X4,F4)	(X2,X4,F2,F3)		(X2,X4, F2,F3,F4)	

Train Dataset

	F1	F2	F3	F4	y
X1	0	0.25	0.5	0.75	0
X2	0.75	0.5	0.25	1	0
X3	1	1	0.75	0.25	1
X4	0.25	0.25	1	0	1

MATLAB Code: Naïve Prototype Classifier on iris dataset

```
function Mdl=NaivePrototype(X_train,y_train)

set = 1:size(X_train,2);
subsets = cell(1, length(set));
for i = 1:length(set)
    subsets(i) = nchoosek(set, i);
end
idn=find(y_train==0);
idp=find(y_train==1);
k=0;
for i=1:numel(subsets)
    for j=1:size(subsets(i),1)
        sfids=subsets(i){j,:};
        for s=1:size(X_train,1)
            curr_y=y_train(s);
            if curr_y==0
                nn_neg=knnsearch(X_train(idn,sfids),X_train(s,sfids), 'K',2);
                nn_neg=nn_neg(end);
                nn_neg=idn(nn_neg);
                nn_pos=knnsearch(X_train(idp,sfids),X_train(s,sfids), 'K',1);
                nn_pos=idp(nn_pos);
            else
                nn_neg=knnsearch(X_train(idn,sfids),X_train(s,sfids), 'K',1);
                nn_neg=idn(nn_neg);
                nn_pos=knnsearch(X_train(idp,sfids),X_train(s,sfids), 'K',2);
                nn_pos=nn_pos(end);
                nn_pos=idp(nn_pos);
            end
            yt=y_train([nn_neg,nn_pos]);
            yhat=yt(knnsearch(X_train([nn_neg,nn_pos],sfids),X_train(:,sfids), 'K',1));
            k=k+1;
            err(k) = sum(yhat~ = y_train)/numel(y_train);
            svs(k,1)=nn_neg;
            svs(k,2)=nn_pos;
            nn_features{k}=sfids;
        end
    end
    [minerr,bestk]=min(err);
    Mdl.PrototypeSampleIDs=svs(bestk,1:2);
    Mdl.PrototypeFeatureIDs=nn_features{bestk};
    Mdl.Error=minerr;
    Mdl.Subsets=subsets;
    Mdl.I=0;
    Mdl.MX=X_train(Mdl.PrototypeSampleIDs,Mdl.PrototypeFeatureIDs);
    Mdl.My=y_train(svs(bestk,1:2));
end
```

```
clear
load fisheriris
y=grp2idx(species)-1;
X=meas;
ids=find(y==0 | y==1);
y=y(ids);
X=X(ids,:);
[N,M]=size(X);
rng(42);
indices = randperm(N);
numTestSamples = round(0.2 * N);
trainIdx = indices (numTestSamples+1:end);
testIdx = indices (1:numTestSamples);
X_train = X(trainIdx, :);
X_test = X(testIdx, :);
y_train = y(trainIdx, :);
y_test = y(testIdx, :);

Mdl=NaivePrototype(X_train,y_train);
y_test_NP=Mdl.My(knnsearch(Mdl.MX,X_test (:,Mdl.PrototypeFeatureIDs), 'K',1));
acc_test=sum(y_test_NP==y_test)/numel(y_test);
```

Sample ID	sepal length	sepal width	petal length	petal width
1	5	3.6	1.4	0.2
2	5.4	3.4	1.5	0.4
3	5.8	4	1.2	0.2
4	5.7	4.4	1.5	0.4
5	5	3.3	1.4	0.2
6	4.9	2.4	3.3	1
7	6.1	2.8	4	1.3
8	5	3.4	1.6	0.4
9	4.3	3	1.1	0.1
10	5.1	3.8	1.9	0.4
11	5.9	3	4.2	1.5
12	5.6	2.9	3.6	1.3
13	5.1	3.8	1.5	0.3
14	4.6	3.6	1	0.2
15	5.4	3.9	1.3	0.4
16	5.5	3.5	1.3	0.2
17	5.4	3	4.5	1.5
18	5.1	3.8	1.6	0.2
19	6	3.4	4.5	1.6
20	5.2	2.7	3.9	1.4
21	5.8	2.7	3.9	1.2
22	6.1	2.9	4.7	1.4
23	6	2.9	4.5	1.5
24	5.1	3.3	1.7	0.5
25	5.1	3.5	1.4	0.2
26	5	2	3.5	1
27	6.2	2.9	4.3	1.3
28	5.7	3.8	1.7	0.3
29	5.1	3.4	1.5	0.2
30	4.8	3.4	1.9	0.2
31	5.5	2.5	4	1.3
32	5.7	3	4.2	1.2
33	4.5	2.3	1.3	0.3
34	5.2	3.5	1.5	0.2
35	4.6	3.2	1.4	0.2
36	5.7	2.9	4.2	1.3
37	5.1	3.5	1.4	0.3
38	6.7	3.1	4.4	1.4
39	5.3	3.7	1.5	0.2
40	5	2.3	3.3	1
41	5.2	3.4	1.4	0.2
42	6.5	2.8	4.6	1.5
43	4.6	3.1	1.5	0.2
44	4.4	2.9	1.4	0.2
45	4.8	3.1	1.6	0.2
46	5.4	3.4	1.7	0.2
47	5.5	2.4	3.7	1
48	6.3	2.3	4.4	1.3
49	4.8	3	1.4	0.3
50	4.4	3	1.3	0.2
51	6.4	2.9	4.3	1.3
52	4.9	3.1	1.5	0.1
53	6.1	3	4.6	1.4
54	6.6	3	4.4	1.4
55	6.7	3.1	4.7	1.5
56	4.7	3.2	1.3	0.2
57	5.8	2.6	4	1.2
58	5.6	2.7	4.2	1.3
59	6.8	2.8	4.8	1.4
60	5.9	3.2	4.8	1.8
61	6.4	3.2	4.5	1.5
62	5	3	1.6	0.2
63	5.8	2.7	4.1	1
64	5	3.2	1.2	0.2
65	6.1	2.8	4.7	1.2
66	6	2.2	4	1
67	4.8	3	1.4	0.1
68	5.5	2.4	3.8	1.1
69	5	3.4	1.5	0.2
70	5.6	3	4.1	1.3
71	5.5	2.3	4	1.3
72	5	3.5	1.6	0.6
73	5.7	2.8	4.5	1.3
74	6.9	3.1	4.9	1.5
75	5.5	4.2	1.4	0.2
76	4.9	3	1.4	0.2
77	4.9	3.1	1.5	0.2
78	7	3.2	4.7	1.4
79	4.8	3.4	1.6	0.2
80	5.6	2.5	3.9	1.1



28

100% Accuracy with Extreme Sparsity (iris data)

- Test Accuracy = 100%
- Prototype for Setosa (Sample 5):
 - Petal length= 1.40
- Prototype for Versicolour (Sample 6) :
 - Petal length= 3.30
- Simple Rule for Classification:** If Petal Length of test sample is closer to Sample 5's petal length (1.4) than Sample 6's petal length (3.3), the prediction is Setosa (0), otherwise it is Versicolour(1).
- Test this rule yourself. It works for all samples! (both train & test)

PrototypeSampleIDs	[5,6]
PrototypeFeatureIDs	3
Error	0
Subsets	1x4 cell
L	0
MX	[1.4000;3.3000]
My	[0;1]

But why nobody has ever tried this classifier before?



Pat Langley
Stanford University

- “*Most courses in Machine Learning ignore older methods with links to cognitive psychology. Few graduate students read papers more than ten years old, so they are not exposed to the classic literature*”, Pat Langley

Weak Connection between ML and Cognitive Psychology

K-Nearest Neighbors resembles exemplar theory, but they are developed independently without any connection between them

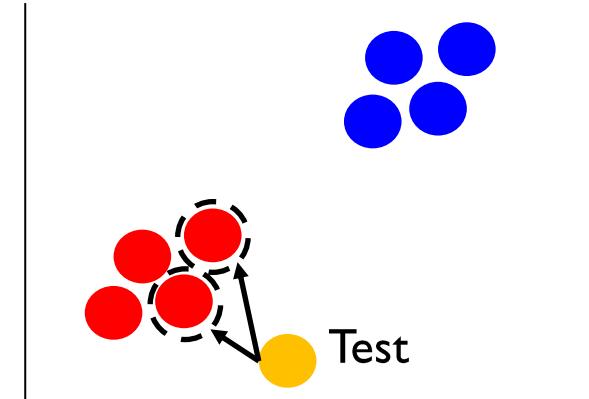
Cognitive Psychology 1975

- **Exemplar Theory of Categorization**
- Individuals categorize objects or events based on their previous experiences with **specific examples**, or exemplars, of those categories.

Rosch, E. (1975). "Cognitive Representations of Semantic Categories." *Journal of Experimental Psychology: General*, 104(3), 192–233.

1951 Statistics

K-Nearest Neighbors



Fix, E., & Hodges Jr, J. L. (1951). "Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties." Project 21-49-004, U.S. Air Force School of Aviation Medicine

Naive Prototype is not a practical classifier!

- Inherently **Interpretable and Explainable**
- Inherently **robust to label noise (noisy samples)**
- **It is not scalable:** $O(n^3 2^p)$
- It is **vulnerable to the curse of dimensionality**
 - Nearest neighbors become meaningless in higher dimensions
- It is **not robust to noisy features**
- Prototype Theory does not offer any solution for the above problems.
- These problem only arises in computers → **needs a computing solution**

Scalability

- The number of ways to choose k features from a set of p features ($p=\text{total number of features}$) without regard to the order of selection:

$$\binom{p}{k} = \frac{p!}{(p-k)!k!}$$

- Our sparse features can be in length of $k=[1,2,\dots,p]$:

$$\frac{p!}{(p-1)!1!} + \frac{p!}{(p-2)!2!} + \dots + \frac{p!}{1!(p-1)!} = 2^p - 2 \rightarrow + 1 \text{ for } k=p \rightarrow 2^p - 1$$

- Every pair of + samples and - samples are potential candidates: Cost of Pair of samples: $n^+ n^- \approx O(n^2)$
- We also need to cross-validate all combinations of sample pairs and feature subsets: $O(2n)$
- Total time Complexity = $O(n^3 2^p)$**
 - Example (only in terms of p): $p=784$ (MNIST data) $\rightarrow 2^{784} \cdot 1 \approx 10^{236} >$ number of atoms in universe (10^{80})

Required Properties

Prototype Theory	Machine Learning	SVM	NCC	PS	NP
Typicality	Each class is represented by only one single prototype	X	✓	X	✓
Core Features	Prototypes have sparse features	X	X	X	✓
Generalizability	Prototype features are generalizable to samples of class	X	X	X	✓
Flexibility	Learning prototypes is an incremental process	X	X	X	✓
	Robustness to noisy labels	✓	✓	✓	✓
	Interpretability (what features are used in the decision?)	X	X	X	✓
	Explainability (reasoning the decision)	○	✓	○	✓
	Robustness to curse of dimensionality	✓	X	X	X
	Robustness to noisy features	X	X	X	X
	Computationally scalable	○	✓	○	X

We need to solve **these issues** to build an **authentic and practical replica of prototype theory** for machine learning

Additionally, we want a natural solution

- **Hyperparameter-free:** Nature does not use hyperparameters
- **Optimization-free:** Optimization does not exist in nature
 - Engineering tricks made by humans
 - Nature instead uses evolution, natural selection, and self-organization
- **Purely based on Nearest neighbor**
 - Brain runs a Nearest neighbor algorithm in an efficient way
 - See the evidence in the next slide

Bypassing curse of dimensionality with nature's algorithm



Fruit Fly

- Naïve prototype classifier is dependent to nearest neighbor search. In higher dimensions , nearest neighbors become irrelevant, because **relevant samples become dissimilar, and irrelevant samples become similar.**
- Prototype theory does not explain **how humans find the nearest neighbor**, but recent evidence has been found in **fruit fly's brain** (Dasgupta , et. al, 2017)
 - Fruit fly's brain uses a version of **locality-sensitive-Hashing (LSH) algorithm** for nearest neighbor search.

Bypassing curse of dimensionality via LSH



Piotr Indyk Rajeev Motwani

- LSH is an efficient solution to nearest-neighbor search by **mapping high-dimensional data points into a lower-dimensional space** in such a way that **similar points are more likely to be hashed into the same bucket** with high probability.
- LSH, focuses on a subset of potential candidates, thereby providing both **computational efficiency** and robustness to the **curse of dimensionality**.
- So, if we use LSH for our nearest neighbor search, we have already solved the first problem!

Indyk, Piotr, and Rajeev Motwani. "Approximate nearest neighbors: towards removing the curse of dimensionality." Proceedings of the thirtieth annual ACM symposium on Theory of computing. 1998.

Properties of Interest

Prototype Theory	Machine Learning	SVM	NCC	PS	NP
Typicality	Each class is represented by only one single prototype	X	✓	X	✓
Core Features	Prototypes have sparse features	X	X	X	✓
Generalizability	Prototype features are generalizable to samples of class	X	X	X	✓
Flexibility	Learning prototypes is an incremental process	X	X	X	✓
	Robustness to noisy labels	✓	✓	✓	✓
	Interpretability (what features are used in the decision?)	X	X	X	✓
	Explainability (reasoning the decision)	○	✓	○	✓
	Robustness to curse of dimensionality	✓	X	X	✓
	Robustness to noisy features	X	X	X	X
	Computationally scalable	X	✓	○	X

We have less problems to solve now

Attacking computational scalability with respect to n

- To solve this, we need to time travel to 1995, take some lessons from Soft-Margin SVM and return back.

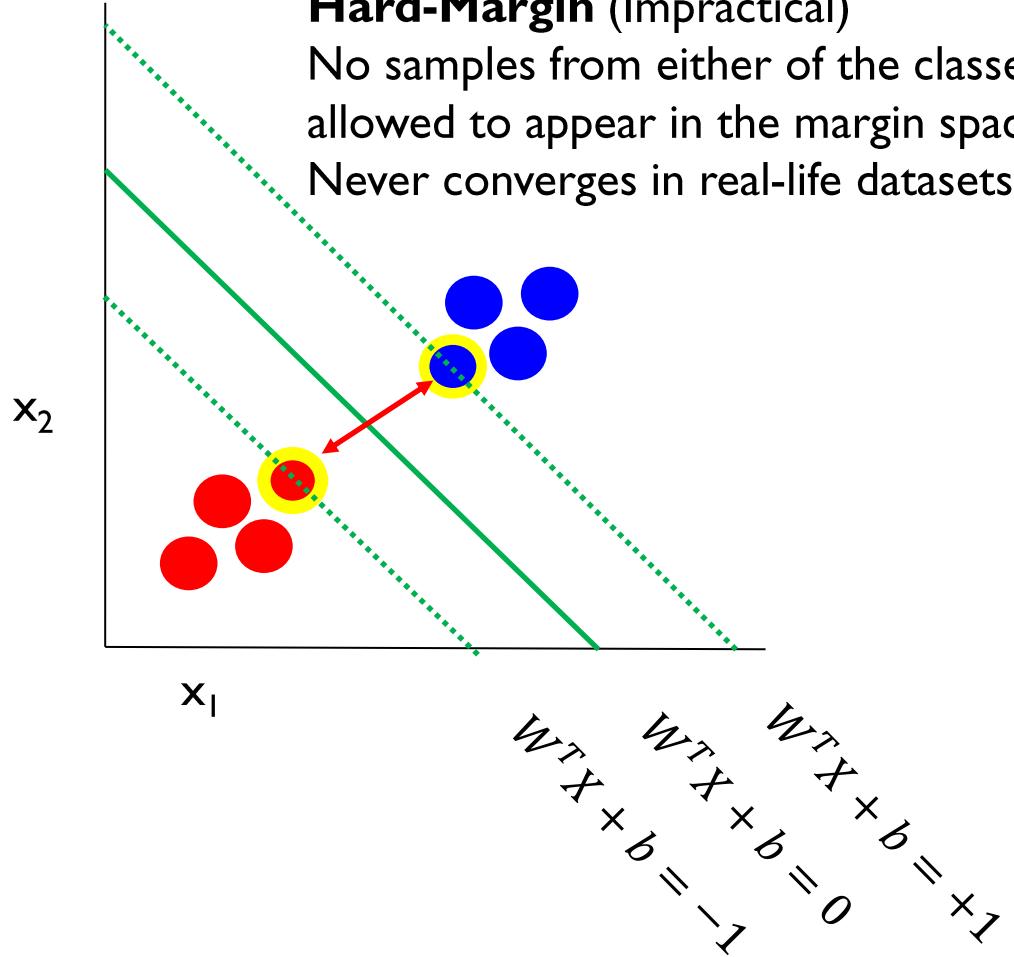


Image source: stock.adobe.com

Hard Margin SVM vs. Soft Margin SVM

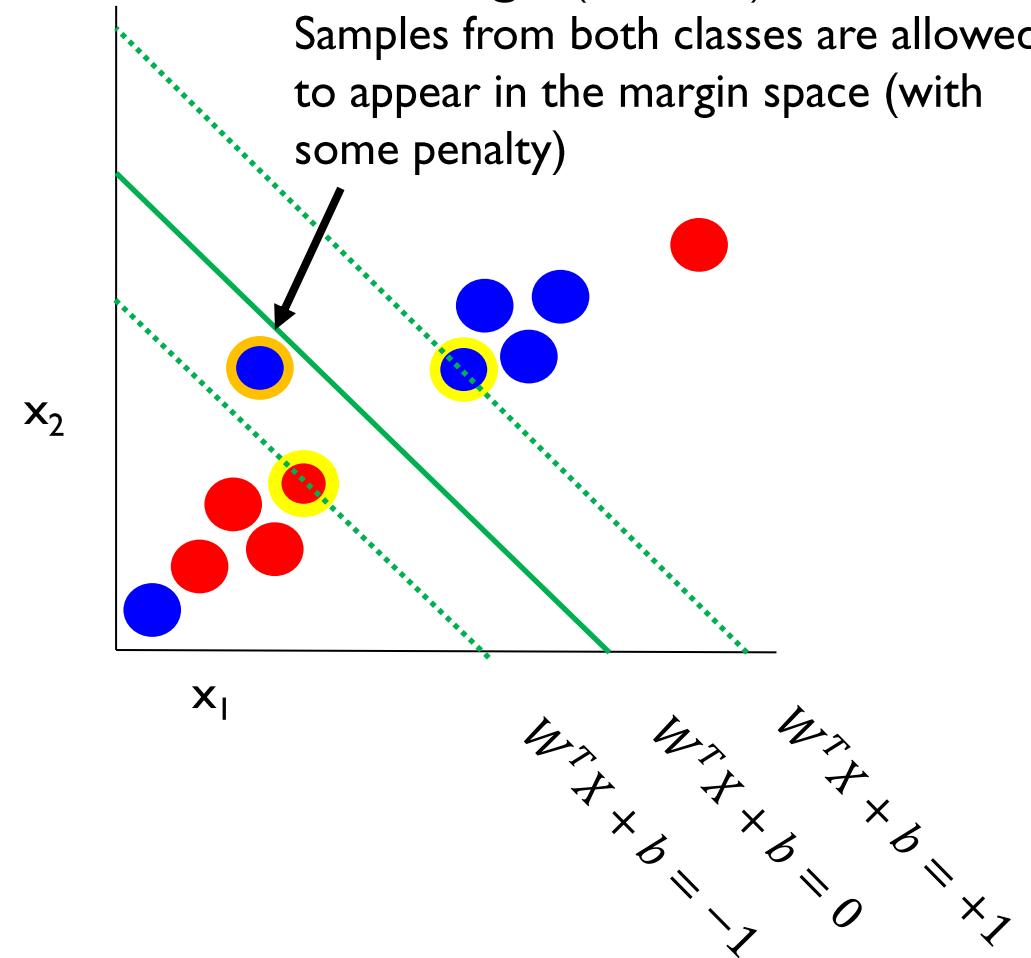
Hard-Margin (Impractical)

No samples from either of the classes are allowed to appear in the margin space
Never converges in real-life datasets



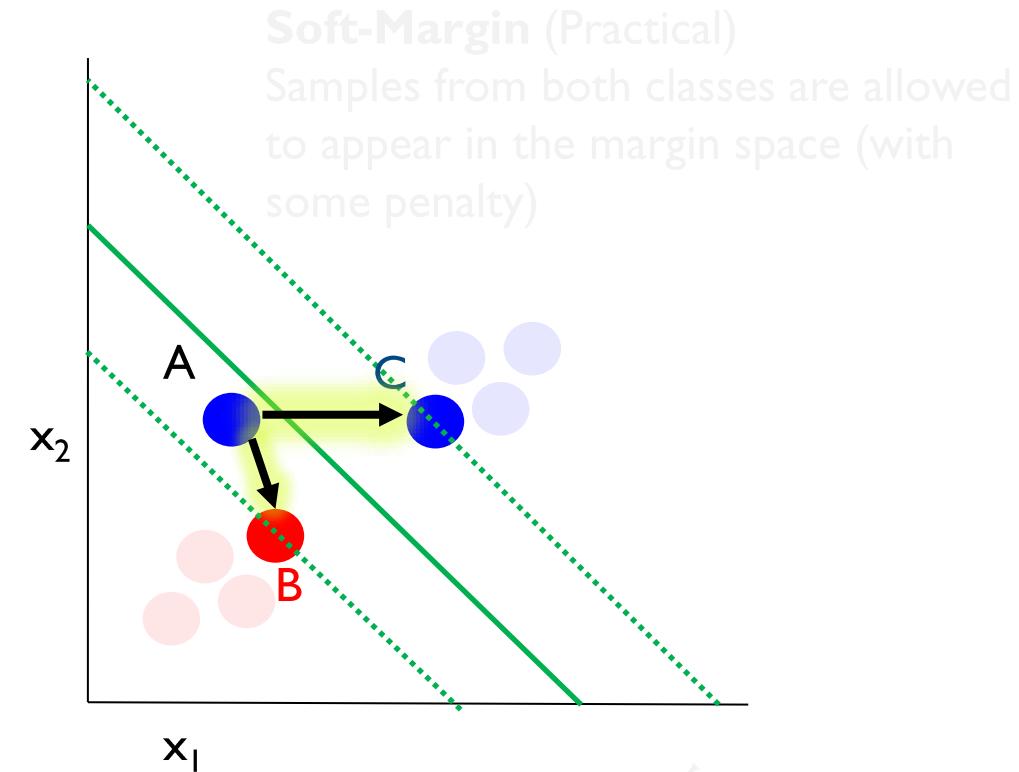
Soft-Margin (Practical)

Samples from both classes are allowed to appear in the margin space (with some penalty)



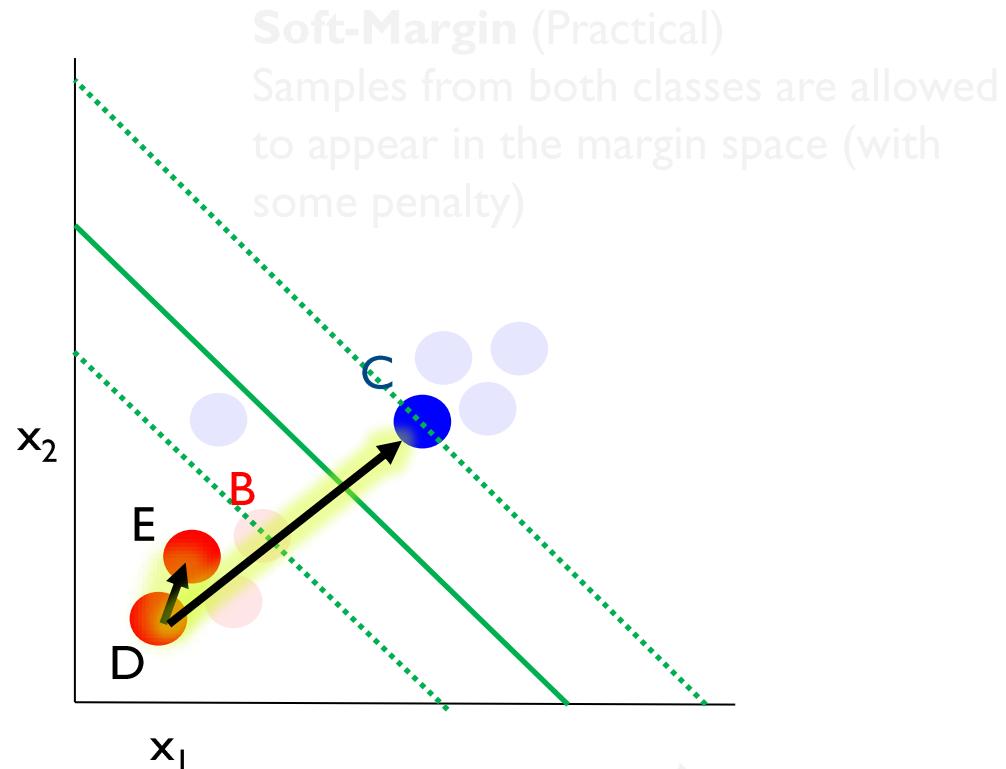
Re-designing SVM with typicality principle of prototype theory

- Prototype theory suggests that **only one support vector exists for each class**. This assumption can simplify the problem of finding support vectors to a regular cross-validation.
- **A** can serve as an ideal **pivot** to find support vectors **B** and **C**
- The Nearest Neighbor of **A** from its own class is **C**, and its nearest neighbor from the other class is **B**.
- If we put **A** as a **pivot**, it shows us the path to reach support vectors **B** and **C**.
- Since **B** and **C** are **actual support vectors**, if we cross-validate their **generalizability**, they **pass** the test!
- So, **margin violation** samples are a very informative source for finding support vectors without the need for optimization!



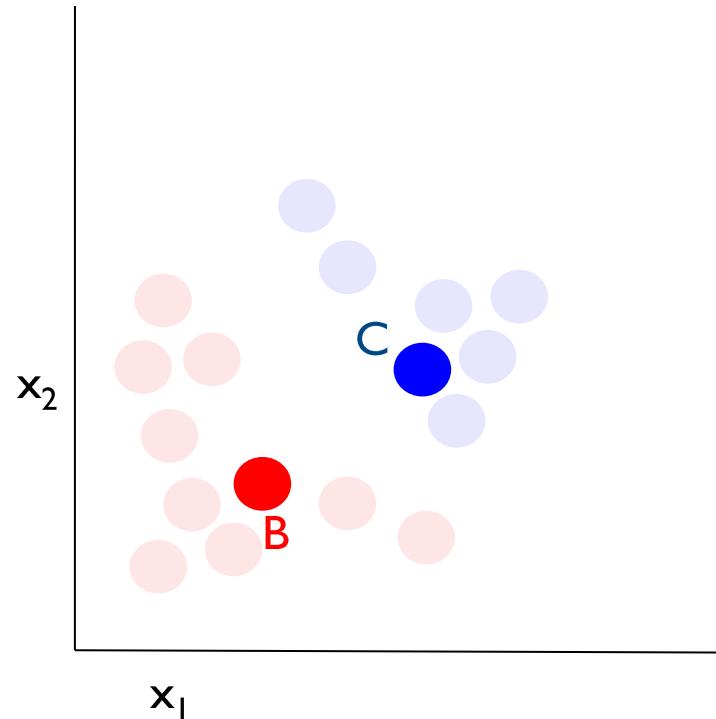
Re-designing SVM with typicality principle of prototype theory

- What about other regular samples?
- The nearest neighbor of **D** is **E** from its class and **C** from the other class. Again, **C** can be found as a support vector, but **E** is a little far from the actual support vector **B**.
- In cross-validation, it is likely that the decision boundary will **not generalize as well as the decision boundary between **B** and **C****, so it automatically will be beaten by actual support vectors (**B** and **C**) suggested by **A** as a pivot.



Fuzzy decision boundary vs. SVM's linear boundary

- We can now enjoy a **fuzzy decision boundary**, which gives us more flexibility
- Removes margin width as a hyperparameter
 - We keep going **hyperparameter-free**

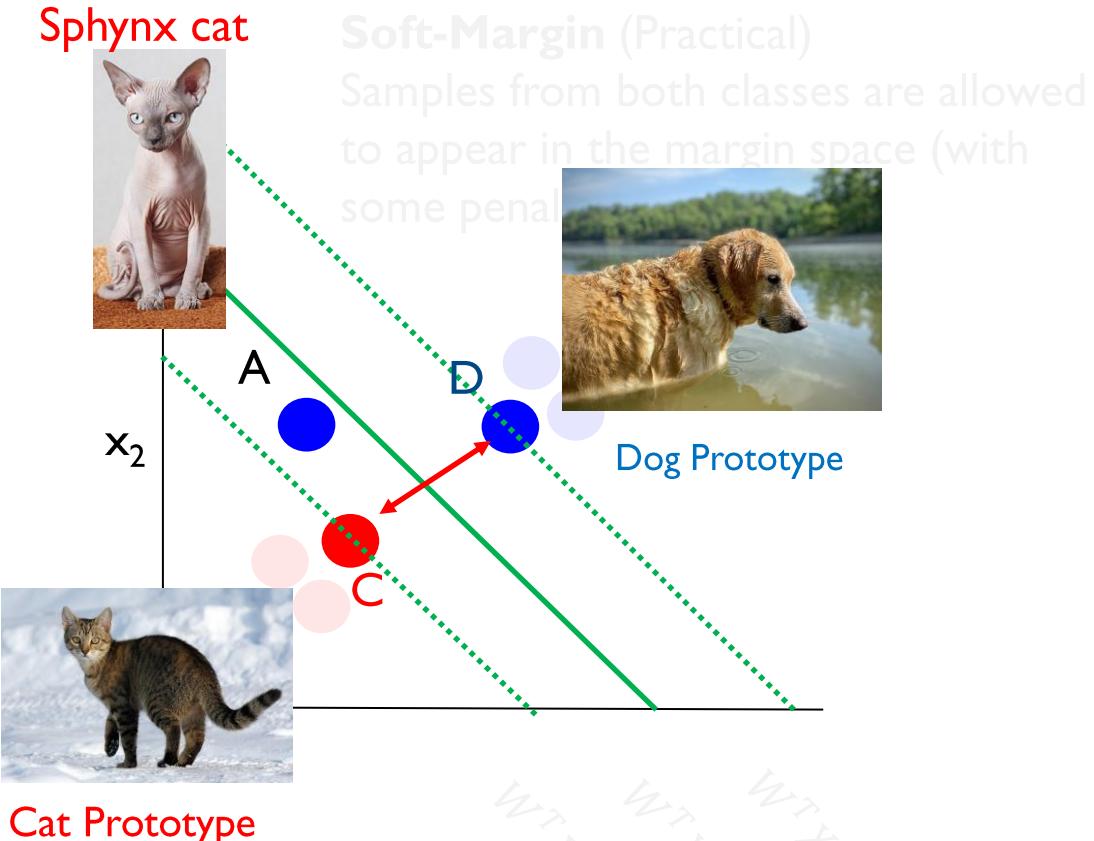


Computational benefit

- Instead of testing all pairs of samples, we can **limit our search to triplets**
 1. Sample
 2. The nearest neighbor from its class (support vector candidate 1)
 3. The nearest neighbor from the other class (support vector candidate 2)
- This reduces the complexity of search in terms of n from $O(n^2)$ to $O(n)$
- Cross-validation costs only $O(2n)$ due to relaxed assumption of **prototype theory**
 - Distance of all samples to only **2 support vectors**

Intuitive Example of Triplets

- **Sphynx cat** is a distinctive cat breed often confused for a dog because of its unique physical characteristics.
 - Example of a triplet
 - **Sphynx cat** is a margin violation sample
- Triplets helps us gain efficiency not only in terms of “n” but also in terms of “p”. How?



Applying “Generalizability Principle”

Nearest Neighbor from
the opposite class (**dog**)
Dog prototype candidate

Pivot (**Cat**)
Sphynx

Nearest Neighbor from its
own class (**Cat**)
Cat prototype Candidate

	Size	Weight	Mustache
	0.80 ↑ 0.20	0.60 ↑ 0.10	0 ↓ 0
	0.60 0.05 < 0.2	0.50 0.01 < 0.1	! < 0 0 ↓ 1
	0.55 ↓ 0.05	0.51 ↓ 0.01	1

Size and **Weight** can still be
the core features of prototypes
because they make a **cat** seem
closer to a **cat** as expected.

Mustache is a candidate for a
non-core feature: it makes pivot
(**cat**) seem closer to the prototype
of the opposite class (**dog**), the class
it does not belong to.

Applying “Generalizability Principle”

- Test Generalization of All samples derived without feature “Mustache”
- Error =0.05

Prototype Candidates I nominated by Sample I after removing non-core candidate

	Size	Weight	Mustache
	0.5	0.5	0
	0.6	0.5	

Feature Matrix with censored non-core features

	Size	Weight	Mustache
Dog 1	0.5	0.5	0
Dog 2	0.4	0.1	0
Dog 3	0.3	0.2	0
Dog 4	0.80	0.60	0
Cat 1	0.60	0.50	
Cat 2	0.8	0.9	
Cat 3	0.5	0.8	
Cat 4	0.55	0.51	0

Applying “Generalizability Principle”

- The next pivot nominates two different samples with different non-core features.
- Error = 0.35

Prototype Candidates **2 nominated by Sample 2** after removing non-core candidate

	Size	Weight	Mustache
	0.4	0.1	0
	0.5	0.8	1

Feature Matrix with censored non-core features

	Size	Weight	Mustache
Dog 1	0.5	0.5	0
Dog 2	0.4	0.1	0
Dog 3	0.3	0.2	0
Dog 4	0.80	0.60	0
Cat 1	0.60	0.50	1
Cat 2	0.8	0.9	1
Cat 3	0.5	0.8	1
Cat 4	0.55	0.51	0

Applying “Generalizability Principle”

- The next pivot nominates two different samples with different non-core features.
- Error = 0.25

Prototype Candidates **n** nominated by Sample **n** after removing non-core candidate

	Size	Weight	Mustache
	0.80	0.6	
	0.55	0.51	

Feature Matrix with censored non-core features

	Size	Weight	Mustache
Dog 1	0.5	0.5	0
Dog 2	0.4	0.1	0
Dog 3	0.3	0.2	0
Dog 4	0.80	0.60	0
Cat 1	0.60	0.50	
Cat 2	0.8	0.9	
Cat 3	0.5	0.8	
Cat 4	0.55	0.51	0

Applying “Generalizability Principle”

The minimum Generalization Error is obtained for the Removal of “**Mustache**” (Error = 0.05)

	Size	Weight
	0.80	0.6
	0.55	0.51

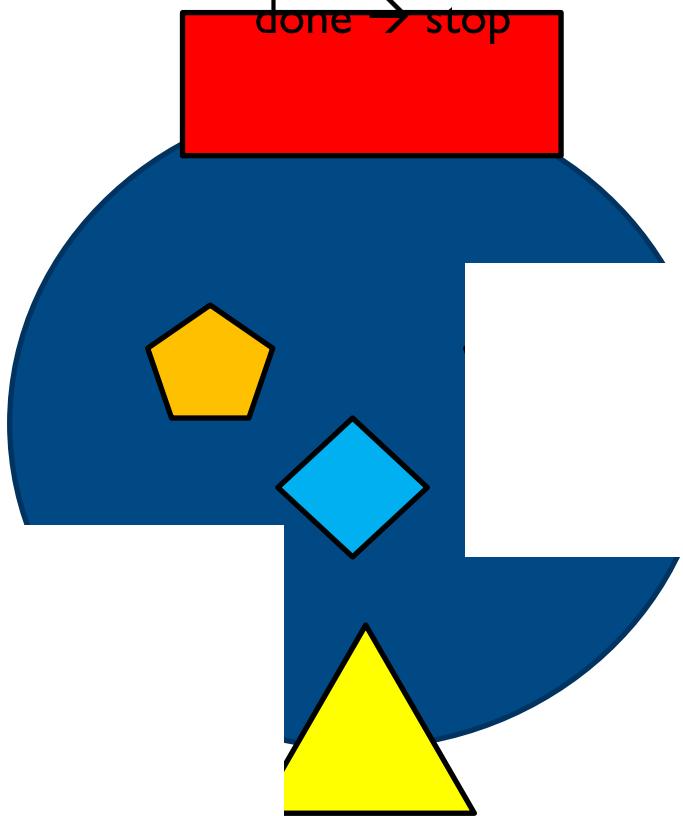
Removing **Mustache** globally from the feature matrix

	Size	Weight
Dog 1	0.5	0.5
Dog 2	0.4	0.1
Dog 3	0.3	0.2
Dog 4	0.80	0.60
Cat 1	0.60	0.50
Cat 2	0.8	0.9
Cat 3	0.5	0.8
Cat 4	0.55	0.51

Intuitive Example (I)

No more pruning can be

done → stop

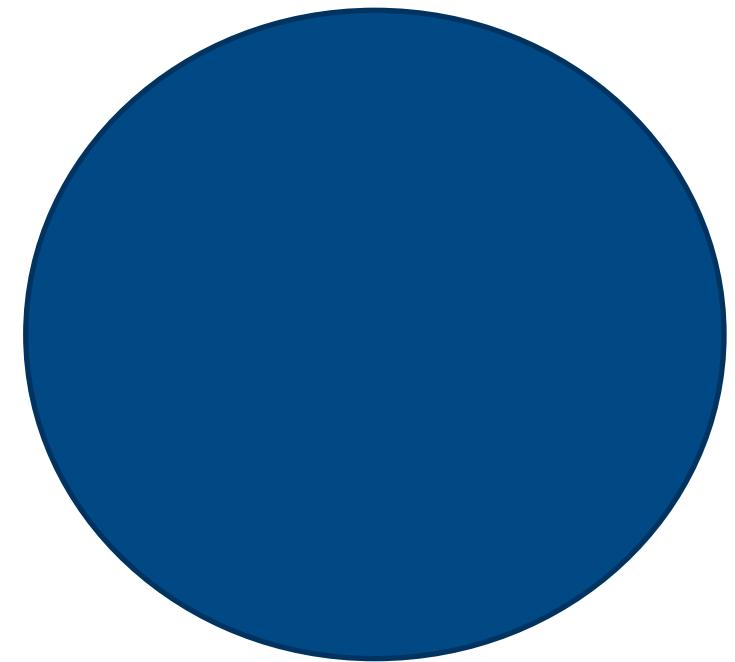


Non-prototype Example

We can't turn this complex object into a simpler prototype in **one step** without having a **feedback** after pruning part of complexity



We remove this part, but it was most part of core features, so it generalizes well, thus we discard pruning.



Basic Level (Sparse Prototype)
Core Feature : Circle

Intuitive Example (2)

- Iteration #1
- Iteration #2
- Iteration #3
- Iteration #4
- Iteration #5
- Stop



Iteration #2

Nearest Neighbor from
the opposite class (**dog**)
Dog prototype candidate

Pivot (**Cat**)

Nearest Neighbor from its
own class (**Cat**)
Cat prototype Candidate

	Size	Weight
	0.80 ↑ 0.2	0.45 ↑ 0.05
	0.60 ↓ 0.05	0.50 ↓ 0.1 0.1 < 0.05
	0.55	0.40

The triplet space of **Sphynx cat** has **new neighbors** because we are in a **new feature space** and have a **better similarity relevance** due to **removed non-core (noisy) features**. We also have **purer** classes.

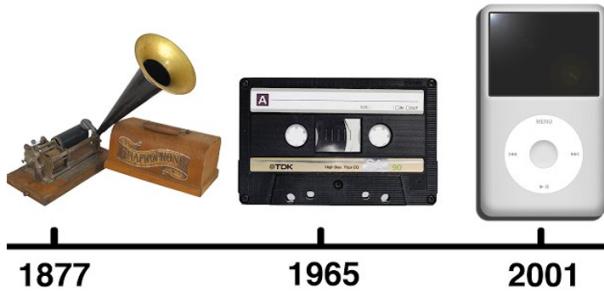
“Flexibility Principle”

- As it can be seen, prototypes can change during each iteration.
 - **Flexibility condition of prototype theory** → Incremental property

Characteristics of Prototype (4): **Flexibility**

- Prototypes are not fixed entities; **they can change** based on new experiences.

Prototypes for Audio Recording



Applying “Typicality Principle”

- If we cannot prune more features, that means that we have reached the **core features of prototypes**
- After removing weight, **Size** becomes the core feature.
- Those prototype candidates that generalize better with the “Size” feature as the core feature become our **final prototype samples**

Final Prototype	Size
	0.80
	0.55

Generalization Error = 0.01

Massive time complexity gain for feature search

- By iterative pruning of non-core features, we naturally reach to features of prototypes in a very efficient way
- So, with this method, we **no longer need to test all subsets of features.**
 - Reducing complexity for feature search from $O(2^p)$ to $O(pL)$
 - Where L is the number of pruning iterations. Normally lower than 20 (e.g. for MNIST, L=10)

Problem Solved

- **Scalable**
 - Reducing Time Complexity from $O(n^3 2^p)$ to $O(n^2 pL)$, where L is number pruning Iterations.
 - The algorithm design allows **high parallelization with GPUs**.
 - Each pruning/validation task can be parallelized or distributed for each sample independently.
- Robust to **Curse of Dimensionality**
 - **Thanks to LSH**
- Robustness to **Noisy Features**
 - By Iterative pruning of non-core features, features of sparse prototypes shows up naturally
 - **No hyperparameter** required for number of pruning iterations
 - Stopping criteria: **when no more features can be pruned**

Problem Solved

Prototype Theory	Machine Learning	SVM	NCC	PS	NP
Typicality	Each class is represented by only one single prototype	✗	✓	✗	✓
Core Features	Prototypes have sparse features	✗	✗	✗	✓
Generalizability	Prototype features are generalizable to samples of class	✗	✗	✗	✓
Flexibility	Learning prototypes is an incremental process	✗	✗	✗	✓
	Robustness to noisy labels	✓	✓	✓	✓
	Interpretability (what features are used in the decision?)	✗	✗	✗	✓
	Explainability (reasoning the decision)	○	✓	○	✓
	Robustness to curse of dimensionality	✓	✗	✗	✓
	Robustness to noisy features	✗	✗	✗	✓
	Computationally scalable	✗	✓	○	✓

This new algorithm is now called “Natural Learning (NL)”

Prototype Theory	Machine Learning	SVM	NCC	PS	NL
Typicality	Each class is represented by only one single prototype	✗	✓	✗	✓
Core Features	Prototypes have sparse features	✗	✗	✗	✓
Generalizability	Prototype features are generalizable to samples of class	✗	✗	✗	✓
Flexibility	Learning prototypes is an incremental process	✗	✗	✗	✓
	Robustness to noisy labels	✓	✓	✓	✓
	Interpretability (what features are used in the decision?)	✗	✗	✗	✓
	Explainability (reasoning the decision)	○	✓	○	✓
	Robustness to curse of dimensionality	✓	✗	✗	✓
	Robustness to noisy features	✗	✗	✗	✓
	Computationally scalable	✗	✓	○	✓

After Rosch's 1973 paper “Natural Categorizes”

Rosch, Eleanor H. "Natural categories." Cognitive psychology 4.3 (1973): 328-350.

Training Algorithm in 20 lines!

Algorithm 1 NLTrain

```
1: Input: training set  $(x, y)$  ( $n$  samples and  $p$  features),  $y_i = \{0, 1\}$ , and features of best prototype  $M$ 
2: Output: prototype samples ( $s_{best}$  and  $o_{best}$ ), and their labels, prototype features  $M$ 
3: if  $M$  is null then
4:    $M \leftarrow \{1, 2, \dots, p\}$            //initialization of prototype features
5: end if
6:  $x = x(:, M)$                   // Copy of  $x$  with features in  $M$ 
7:  $e_{best} \leftarrow \infty$           //initialization of best error. Allowing NL to learn better prototypes at each iteration.
8: for each sample  $i$  in  $x$  do
9:    $s \leftarrow$  index of  $x_i$ 's nearest neighbor from same class using LSH      //prototype sample candidate
10:   $o \leftarrow$  index of  $x_i$ 's nearest neighbor from opposite class using LSH      //prototype sample candidate
11:   $C \leftarrow$  indices of features in  $M$  that make  $x_i$  closer to  $x_s$  than  $x_o$       //prototype features candidate
12:   $\hat{y} \leftarrow NLPredict(x_s, x_o, y_s, y_o, C, x)$     // test the generalization of prototype candidate
13:   $e \leftarrow \sum(y \neq \hat{y})$ 
14:  if  $e < e_{best}$  &  $|C| > 1$  then
15:     $(s_{best}, o_{best}) \leftarrow (s, o)$           // Best prototype samples
16:     $C_{best} \leftarrow C$                       //Best prototype features
17:     $e_{best} \leftarrow e$                       //Best error so far
18:  end if
19: end for
20: if  $|C_{best}| \neq |M|$  then
21:    $M \leftarrow C_{best}$ 
22:    $NLTrain(x, y, M)$ 
23: end if
```

Hyperparameter-free

Self-explainable Algorithm

Code available in MATLAB, Python, and R

Algorithm 2 NLPredict

```
1: Input: data ( $x$ ), prototype samples ( $x_o$  and  $x_s$ ) and corresponding labels ( $y_o$  and  $y_s$ ) and features ( $M$ )
2: Output:  $\hat{y}$  (Predicted labels)
3:  $x \leftarrow x(:, M)$           // copy of  $x$  with prototype features ( $M$  or  $C_i$ )
4: for each sample  $i$  in  $x$  do
5:    $(d_s, d_o) \leftarrow D(x_i, x_s, x_o)$           //Distance of example to prototype samples  $s$  and  $o$ 
6:    $\hat{y}_i = y_s$ 
7:   if  $d_o < d_s$  then
8:      $\hat{y}_i = y_o$ 
9:   end if
10:  end for
```

Illustrative Example: Iteration # 2

I – Feature Matrix

	F1	F2	F3	F4	y
X1	0	0.25	0.5	0.75	0
X2	0.75	0.5	0.25	1	0
X3	1	1	0.75	0.25	1
X4	0.25	0.25	1	0	1

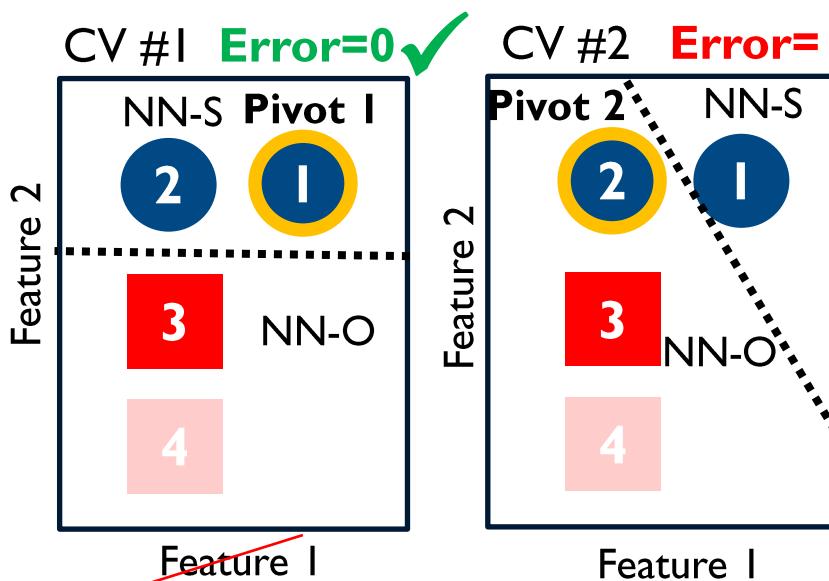
2. BestError = inf



• • •

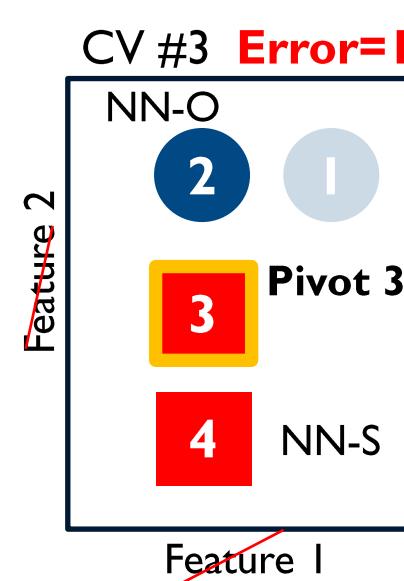
Geometric View of Decision Boundaries

Iteration #1 Active Features: [1,2]

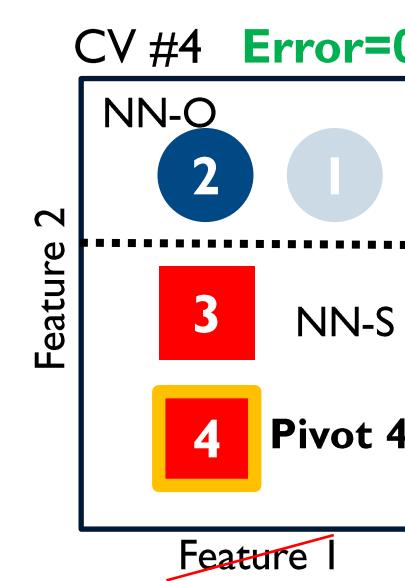


Feature 1 has the same distance to pivot's NN-S and NN-O → Non-core pruning Candidate

Both features are relevant, so no pruning candidate

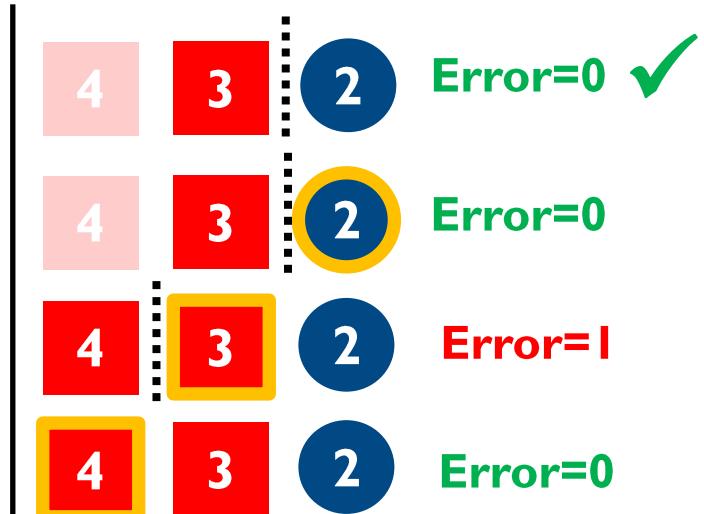


Both features have the same distance to pivot's NN-S and NN-O → cross-validation will not be performed



Best Prototype's features: [2]
Best Prototype Samples: [2,3]

Iteration #2



Active Features = [2]
Final Prototype samples: [2,3]
Final Prototype features: [2]

MNIST Dataset (0 vs. 1) – Iteration 1

- Iteration=1 PivotSample=1/12665 PrototypeCandidate=[5845,8205], NumFeatures=166/784, Error=0.0756
- Iteration=1 PivotSample=4/12665 PrototypeCandidate=[844,10217], NumFeatures=153/784, Error=0.0360
- Iteration=1 PivotSample=21/12665 PrototypeCandidate=[3840,11859], NumFeatures=150/784, Error=0.0258
- Iteration=1 PivotSample=76/12665 PrototypeCandidate=[815,11478], NumFeatures=148/784, Error=0.0188
- Iteration=1 PivotSample=208/12665 PrototypeCandidate=[3175,7938], NumFeatures=122/784, Error=0.0121
- Iteration=1 PivotSample=245/12665 PrototypeCandidate=[4403,7780], NumFeatures=110/784, Error=0.0114
- Iteration=1 PivotSample=402/12665 PrototypeCandidate=[513,7078], NumFeatures=125/784, Error=0.0081
- Iteration=1 PivotSample=2062/12665 PrototypeCandidate=[5011,7780], NumFeatures=165/784, **Error=0.0046**



This sample likely is a **margin violation sample** that **guides us towards good support vectors** (e.g., **A** in our SVM example)

MNIST Dataset (0 vs. 1) – Iteration 2

- Iteration=2 PivotSample=1/12665 PrototypeCandidate=[62,8205], NumFeatures=88/165, Error=0.0882
- Iteration=2 PivotSample=2/12665 PrototypeCandidate=[3652,8205], NumFeatures=86/165, Error=0.0853
- Iteration=2 PivotSample=3/12665 PrototypeCandidate=[4167,10217], NumFeatures=93/165, Error=0.0791
- Iteration=2 PivotSample=5/12665 PrototypeCandidate=[5396,9876], NumFeatures=99/165, Error=0.0711
- Iteration=2 PivotSample=14/12665 PrototypeCandidate=[3567,9987], NumFeatures=64/165, Error=0.0388
- Iteration=2 PivotSample=16/12665 PrototypeCandidate=[4400,9719], NumFeatures=88/165, Error=0.0126
- Iteration=2 PivotSample=124/12665 PrototypeCandidate=[3943,6696], NumFeatures=91/165, Error=0.0066
- Iteration=2 PivotSample=6228/12665 PrototypeCandidate=[6065,1711], NumFeatures=107/165, Error=0.0066
- Iteration=2 PivotSample=7297/12665 PrototypeCandidate=[6447,5814], NumFeatures=88/165, Error=0.0063
- Iteration=2 PivotSample=8095/12665 PrototypeCandidate=[6153,2611], NumFeatures=76/165, **Error=0.0061**



Flexibility principle: support vectors are now changed! They are sparser!

MNIST Dataset (0 vs. 1) – Iteration 3

- Iteration=3 PivotSample=1/12665 PrototypeCandidate=[1257,8183], NumFeatures=54/76, Error=0.0695
- Iteration=3 PivotSample=17/12665 PrototypeCandidate=[4740,9332], NumFeatures=58/76, Error=0.0471
- Iteration=3 PivotSample=27/12665 PrototypeCandidate=[2547,7622], NumFeatures=41/76, Error=0.0385
- Iteration=3 PivotSample=31/12665 PrototypeCandidate=[3388,7931], NumFeatures=28/76, Error=0.0165
- Iteration=3 PivotSample=345/12665 PrototypeCandidate=[5305,7160], NumFeatures=47/76, Error=0.0099
- Iteration=3 PivotSample=779/12665 PrototypeCandidate=[5053,7160], NumFeatures=53/76, Error=0.0099
- Iteration=3 PivotSample=943/12665 PrototypeCandidate=[1627,9332], NumFeatures=41/76, Error=0.0087
- Iteration=3 PivotSample=1203/12665 PrototypeCandidate=[3622,7996], NumFeatures=31/76, **Error=0.0053**

MNIST Dataset (0 vs. 1) – Iteration 4

- Iteration=4 PivotSample=1/12665 PrototypeCandidate=[2068,9477], NumFeatures=17/31, Error=0.1368
- Iteration=4 PivotSample=2/12665 PrototypeCandidate=[5840,9477], NumFeatures=15/31, Error=0.1268
- Iteration=4 PivotSample=3/12665 PrototypeCandidate=[3317,9287], NumFeatures=21/31, Error=0.0582
- Iteration=4 PivotSample=5/12665 PrototypeCandidate=[5698,8702], NumFeatures=26/31, Error=0.0511
- Iteration=4 PivotSample=8/12665 PrototypeCandidate=[2419,10601], NumFeatures=12/31, Error=0.0250
- Iteration=4 PivotSample=12/12665 PrototypeCandidate=[735,9308], NumFeatures=21/31, Error=0.0207
- Iteration=4 PivotSample=26/12665 PrototypeCandidate=[2867,10238], NumFeatures=14/31, Error=0.0115
- Iteration=4 PivotSample=61/12665 PrototypeCandidate=[571,11298], NumFeatures=22/31, Error=0.0073
- Iteration=4 PivotSample=989/12665 PrototypeCandidate=[4792,12441], NumFeatures=20/31, Error=0.0071
- Iteration=4 PivotSample=2860/12665 PrototypeCandidate=[3951,7869], NumFeatures=16/31, Error=0.0060
- Iteration=4 PivotSample=3185/12665 PrototypeCandidate=[5428,8945], NumFeatures=18/31, Error=0.0057
- Iteration=4 PivotSample=6783/12665 PrototypeCandidate=[10423,4083], NumFeatures=24/31, **Error=0.0050**

MNIST Dataset (0 vs. 1) – Iteration 5

- Iteration=5 PivotSample=1/12665 PrototypeCandidate=[5836,11670], NumFeatures=15/24, Error=0.1375
- Iteration=5 PivotSample=2/12665 PrototypeCandidate=[2259,9376], NumFeatures=13/24, Error=0.0747
- Iteration=5 PivotSample=4/12665 PrototypeCandidate=[3753,10240], NumFeatures=20/24, Error=0.0168
- Iteration=5 PivotSample=65/12665 PrototypeCandidate=[119,10240], NumFeatures=15/24, Error=0.0160
- Iteration=5 PivotSample=66/12665 PrototypeCandidate=[4880,10996], NumFeatures=16/24, Error=0.0114
- Iteration=5 PivotSample=305/12665 PrototypeCandidate=[4219,10240], NumFeatures=18/24, Error=0.0114
- Iteration=5 PivotSample=605/12665 PrototypeCandidate=[5486,8936], NumFeatures=16/24, Error=0.0107
- Iteration=5 PivotSample=749/12665 PrototypeCandidate=[2746,6051], NumFeatures=10/24, Error=0.0066
- Iteration=5 PivotSample=2746/12665 PrototypeCandidate=[749,6051], NumFeatures=11/24, Error=0.0060
- Iteration=5 PivotSample=4083/12665 PrototypeCandidate=[3709,12086], NumFeatures=15/24, **Error=0.0051**

MNIST Dataset (0 vs. 1) – Iteration 6

- Iteration=7 PivotSample=2/12665 PrototypeCandidate=[1686,7921], NumFeatures=3/10, Error=0.1171
- Iteration=6 PivotSample=1/12665 PrototypeCandidate=[5836,7622], NumFeatures=12/15, Error=0.1248
- Iteration=6 PivotSample=4/12665 PrototypeCandidate=[5751,10240], NumFeatures=13/15, Error=0.0141
- Iteration=6 PivotSample=13/12665 PrototypeCandidate=[4381,10240], NumFeatures=13/15, Error=0.0140
- Iteration=6 PivotSample=35/12665 PrototypeCandidate=[44,10240], NumFeatures=13/15, Error=0.0139
- Iteration=6 PivotSample=73/12665 PrototypeCandidate=[1694,10240], NumFeatures=13/15, Error=0.0137
- Iteration=6 PivotSample=92/12665 PrototypeCandidate=[3714,6627], NumFeatures=13/15, Error=0.0107
- Iteration=6 PivotSample=209/12665 PrototypeCandidate=[4690,7628], NumFeatures=6/15, Error=0.0098
- Iteration=6 PivotSample=1269/12665 PrototypeCandidate=[702,6121], NumFeatures=7/15, Error=0.0077
- Iteration=6 PivotSample=6209/12665 PrototypeCandidate=[9503,1609], NumFeatures=10/15, **Error=0.0069**

MNIST Dataset (0 vs. 1) – Iteration 7

- Iteration=7 PivotSample=6/12665 PrototypeCandidate=[962,7250], NumFeatures=6/10, Error=0.0325
- Iteration=7 PivotSample=14/12665 PrototypeCandidate=[1219,7555], NumFeatures=5/10, Error=0.0308
- Iteration=7 PivotSample=249/12665 PrototypeCandidate=[2327,6577], NumFeatures=3/10, Error=0.0156
- Iteration=7 PivotSample=337/12665 PrototypeCandidate=[1950,6577], NumFeatures=3/10, Error=0.0090
- Iteration=7 PivotSample=4796/12665 PrototypeCandidate=[4300,9531], NumFeatures=2/10, Error=0.0078
- Iteration=7 PivotSample=5883/12665 PrototypeCandidate=[4482,6577], NumFeatures=4/10, Error=0.0069
- Iteration=7 PivotSample=5995/12665 PrototypeCandidate=[10327,1609], NumFeatures=10/10, Error=0.0063
- Iteration=7 PivotSample=8335/12665 PrototypeCandidate=[9625,1609], NumFeatures=8/10, **Error=0.0061**

MNIST Dataset (0 vs. 1) – Iteration 8

- Iteration=8 PivotSample=1/12665 PrototypeCandidate=[3516,7412], NumFeatures=6/8, Error=0.1189
- Iteration=8 PivotSample=3/12665 PrototypeCandidate=[2655,5985], NumFeatures=3/8, Error=0.1169
- Iteration=8 PivotSample=6/12665 PrototypeCandidate=[962,6506], NumFeatures=7/8, Error=0.0549
- Iteration=8 PivotSample=8/12665 PrototypeCandidate=[4988,7555], NumFeatures=2/8, Error=0.0493
- Iteration=8 PivotSample=11/12665 PrototypeCandidate=[4852,7250], NumFeatures=5/8, Error=0.0264
- Iteration=8 PivotSample=38/12665 PrototypeCandidate=[2700,6506], NumFeatures=7/8, Error=0.0242
- Iteration=8 PivotSample=123/12665 PrototypeCandidate=[1809,6506], NumFeatures=6/8, Error=0.0235
- Iteration=8 PivotSample=275/12665 PrototypeCandidate=[357,7572], NumFeatures=4/8, Error=0.0159
- Iteration=8 PivotSample=1152/12665 PrototypeCandidate=[1899,7818], NumFeatures=2/8, Error=0.0120
- Iteration=8 PivotSample=1910/12665 PrototypeCandidate=[2422,11150], NumFeatures=4/8, Error=0.0107
- Iteration=8 PivotSample=3354/12665 PrototypeCandidate=[3982,8105], NumFeatures=7/8, **Error=0.0069**

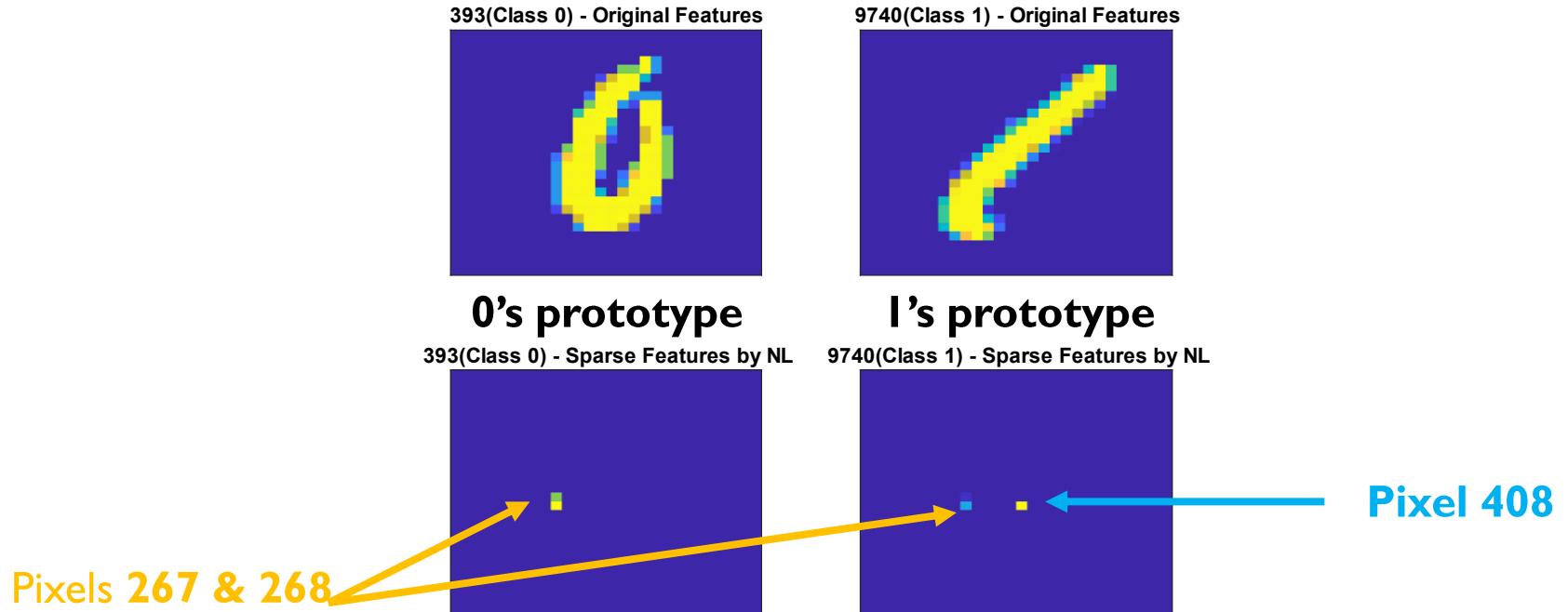
MNIST Dataset (0 vs. 1) – Iteration 9

- Iteration=9 PivotSample=1/12665 PrototypeCandidate=[3516,8387], NumFeatures=4/7, Error=0.1526
- Iteration=9 PivotSample=3/12665 PrototypeCandidate=[2655,8387], NumFeatures=3/7, Error=0.1225
- Iteration=9 PivotSample=6/12665 PrototypeCandidate=[962,7412], NumFeatures=3/7, Error=0.1077
- Iteration=9 PivotSample=8/12665 PrototypeCandidate=[3911,12101], NumFeatures=2/7, Error=0.0452
- Iteration=9 PivotSample=11/12665 PrototypeCandidate=[1201,7412], NumFeatures=4/7, Error=0.0449
- Iteration=9 PivotSample=22/12665 PrototypeCandidate=[2105,7555], NumFeatures=3/7, Error=0.0292
- Iteration=9 PivotSample=173/12665 PrototypeCandidate=[3694,9066], NumFeatures=3/7, Error=0.0098
- Iteration=9 PivotSample=473/12665 PrototypeCandidate=[5549,9066], NumFeatures=3/7, Error=0.0094
- Iteration=9 PivotSample=739/12665 PrototypeCandidate=[5656,9066], NumFeatures=3/7, **Error=0.0069**

MNIST Dataset (0 vs. 1) – Iteration 10

- Iteration=10 PivotSample=1/12665 PrototypeCandidate=[224,7622], NumFeatures=2/3, Error=0.7712
- Iteration=10 PivotSample=4/12665 PrototypeCandidate=[38,11310], NumFeatures=2/3, Error=0.5754
- Iteration=10 PivotSample=9/12665 PrototypeCandidate=[10,11310], NumFeatures=2/3, Error=0.5677
- Iteration=10 PivotSample=11/12665 PrototypeCandidate=[2349,7412], NumFeatures=2/3, Error=0.2671
- Iteration=10 PivotSample=12/12665 PrototypeCandidate=[2888,11614], NumFeatures=2/3, Error=0.1240
- Iteration=10 PivotSample=24/12665 PrototypeCandidate=[4712,11859], NumFeatures=3/3, Error=0.0827
- Iteration=10 PivotSample=53/12665 PrototypeCandidate=[5701,9740], NumFeatures=3/3, Error=0.0102
- Iteration=10 PivotSample=103/12665 PrototypeCandidate=[393,9740], NumFeatures=3/3, Error=**0.0067**
- Best Prototype=[Sample **393(class 0)**, Sample **9740(class 1)**], Best Error=0.0067, **Core Features=[267 268 408]**

Visualization of Sparse Prototypes found by Natural Learning

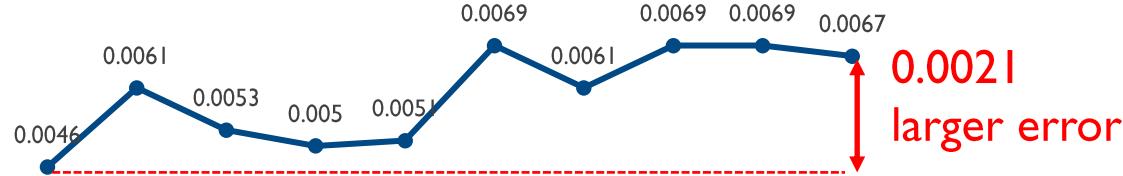


If test sample's pixels 267,268,408, collectively make test example closer to **0's prototype** than **1's prototype**, it is **0**, otherwise it is **1**

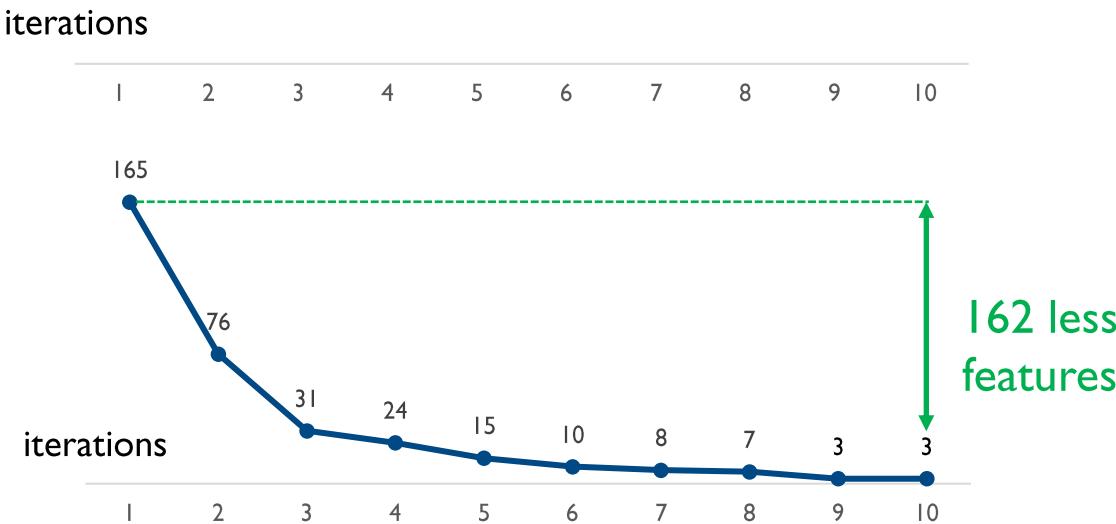
Accuracy on Train : 99.33%
Accuracy on Test: **99.48%**

MNIST Dataset (0 vs. 1) – Summary

Generalization Error

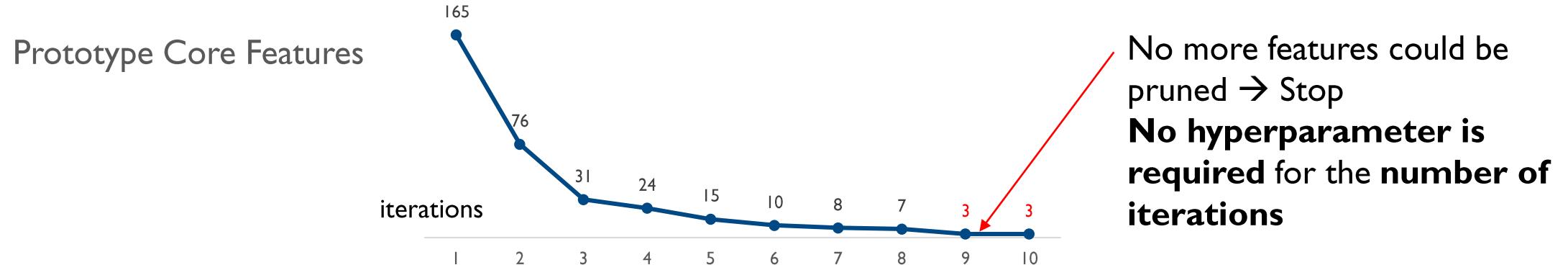


Prototype Core Features



The key to have **extreme sparsity** is to **sacrifice small accuracy** in exchange of **gain in sparsity**

MNIST Dataset (0 vs. 1) – Summary



Comparison of Properties with other classifiers

Model	Local rules?	Estimate Weight?	Hyperparameters?	Memorize Train set?
Nearest Neighbor (INN)	No	No	No	Yes
Deep Neural Networks (DNN)	No	Yes	Yes	No
Random Forest (RF)	Yes	No	Yes	No
Decision Trees (DT)	Yes	No	Yes/No	No
Logistic Regression (LR)	No	Yes	No	No
Linear discriminant Analysis (LDA)	No	Yes	No	No
Support Vector Machines (SVM)	No	Yes	Yes	No
Natural Learning (NL)	No	No	No	Only 2 Samples

Connection of NL with other classifiers

- Special case of **Nearest Neighbor Classifier (1NN)**
 - NL=Extremely Sparse Nearest Neighbor Classifier
- Special case of **Linear SVM**
 - NL = Sparse Singular Support Vector Machines (**Hyperparameter-free**)
- Special version of **Decision Trees**
 - Finds a **single multi-attribute rule**: e.g., If the test sample's features F1, F25, and F100 are closer to [0.12, 0.26, 0.27] comparing [0.26, 0.28, 0.29], it is labeled 1, otherwise 0.
- It shares characteristics with **LDA** and **Deep Learning**: simultaneously performs dimension reduction and classification.
 - NL : **Original Space**
 - LDA: Linear Latent Space
 - Deep Learning: Non-linear Latent Space

Experimental Evaluation: Datasets

- 17 benchmark datasets for binary classification from the healthcare domain where NL's strength is supposed to be at the level of black-box models due to noisy labels in this domain (Semenova et al., 2023)
 - 9 high-dimensional datasets ($n \ll p$)
 - 8 low-dimensional datasets ($n \gg p$)
- 10 Stratified sampling for each dataset (10-fold) to reduce the bias of train/test split
 - 170 train/test set in total

High-Dimensional (Gene Expression) Datasets ($N \ll P$)						Low-Dimensional Datasets ($N \gg P$)					
Dataset	#p	#n	MjClass	ID*	Description	Dataset	#p	#n	MjClass	ID*	Description
AP_Breast_Colon	10935	630	54.60%	1145	Breast vs. Colon Cancer	blood-transfusion	4	748	76.20%	1464	Donor of Blood Transfusion (UCI)
AP_Breast_Kidney	10935	604	56.95%	1158	Breast vs. Kidney Cancer	diabetes	8	768	65.10%	42608	Diabetes Patient (OpenML)
AP_Breast_Ovary	10935	542	63.47%	1165	breast vs. Ovarian Cancer	Haberman	14	306	73.53%	43	Breast Cancer Survival (UCI)
AP_Colon_Kidney	10935	546	52.38%	1137	Colon vs. Kidney Cancer	heart-statlog	13	270	55.56%	53	Heart Disease Database (UCI)
OVA_Colon	10935	1545	81.49%	1161	Colon Cancer vs. others	hiva_agnostic	1617	4229	96.48%	1039	AIDS HIV infection (ETH Zurich)
OVA_Kidney	10935	1545	83.17%	1134	Kidney Cancer vs. others	ilpd-numeric	10	583	71.36%	41945	Indian Liver Patient Dataset (UCI)
OVA_Lung	10935	1545	91.84%	1130	Lung Cancer vs. others	thoracic-surgery	37	470	85.11%	4329	Lung Cancer life expectancy (UCI)
OVA_Omentum	10935	1545	95.02%	1139	Omentum Cancer vs. others	wdbc	30	569	62.74%	1510	Breast Cancer Wisconsin (UCI)
OVA_Ovary	10935	1545	87.18%	1166	Ovarian Cancer vs. others	* OpenML dataset identifier					

Finetuning baseline models

We compare NL versus finetuned baseline models to have a fair comparison. We get the practical configuration settings from applied machine learning sources [1] and [2] for a realistic comparison.

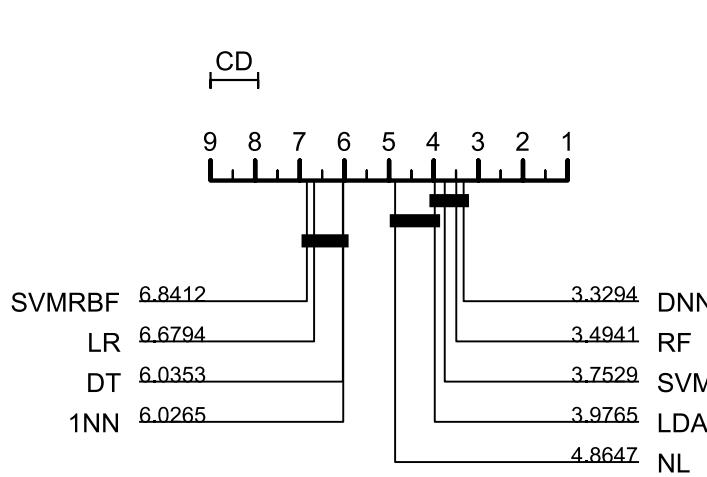
Classifier	Hyperparameter Search	Tested Combinations
Decision Trees	MaxSplits=[1, 5, 10, 20, 50, n], MinLeafSize=[1, 5, 10, 20, 50]	30
Linear SVM	C=[100, 10, 1.0, 0.1, 0.001]	5
SVM-RBF	C=[100, 10, 1.0, 0.1, 0.001], gamma=[$2^{-16} \dots 2^8$] as suggested by [2] with step of 2^2	65
Random Forests (RF)	MaxSplits =[1, 5, 10, 20, 50, n], MinLeafSize=[1, 5, 10, 20, 50], NumTrees= [10,50,100]	90
Deep Neural Networks (DNN)	Batch size=32, Optimizer=Stochastic gradient descent, max epoch of 20, Hidden Layers=[10, 30, 50], Layers=[2, 3, 4], Learning Rate=[0.01, 0.001] and Activation Functions={ReLU,Tanh, Sigmoid}	54
Latent Discriminant Analysis (LDA)	Hyperparameter-free	1
Logistic Regression (LR)	Hyperparameter-free	1
Natural Learning (NL)	Hyperparameter-free	1

[1] <https://machinelearningmastery.com/>

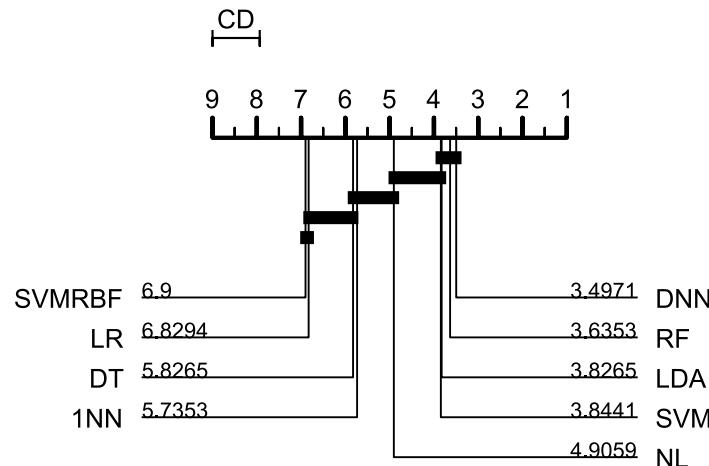
[2] Fernández-Delgado, Manuel, et al. "Do we need hundreds of classifiers to solve real world classification problems?." The journal of machine learning research 15.1 (2014): 3133-3181.

Results: Accuracy and F-measure, Winning Ratio

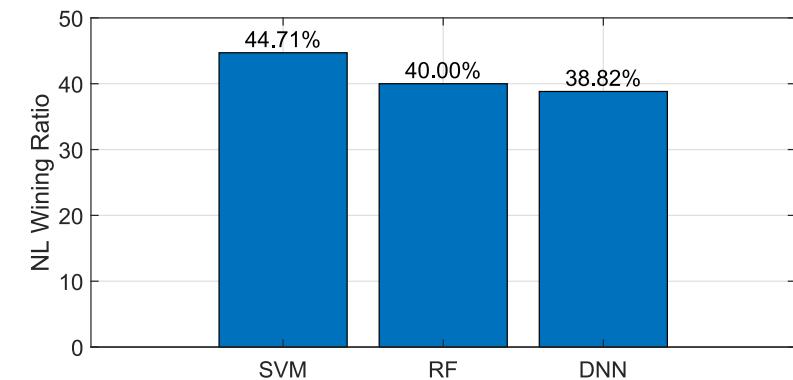
* Critical Difference Diagram, Horizontal line indicates lack of statistical significance at alpha = 0.01 (Nemenyi's test)



Accuracy



F-measure

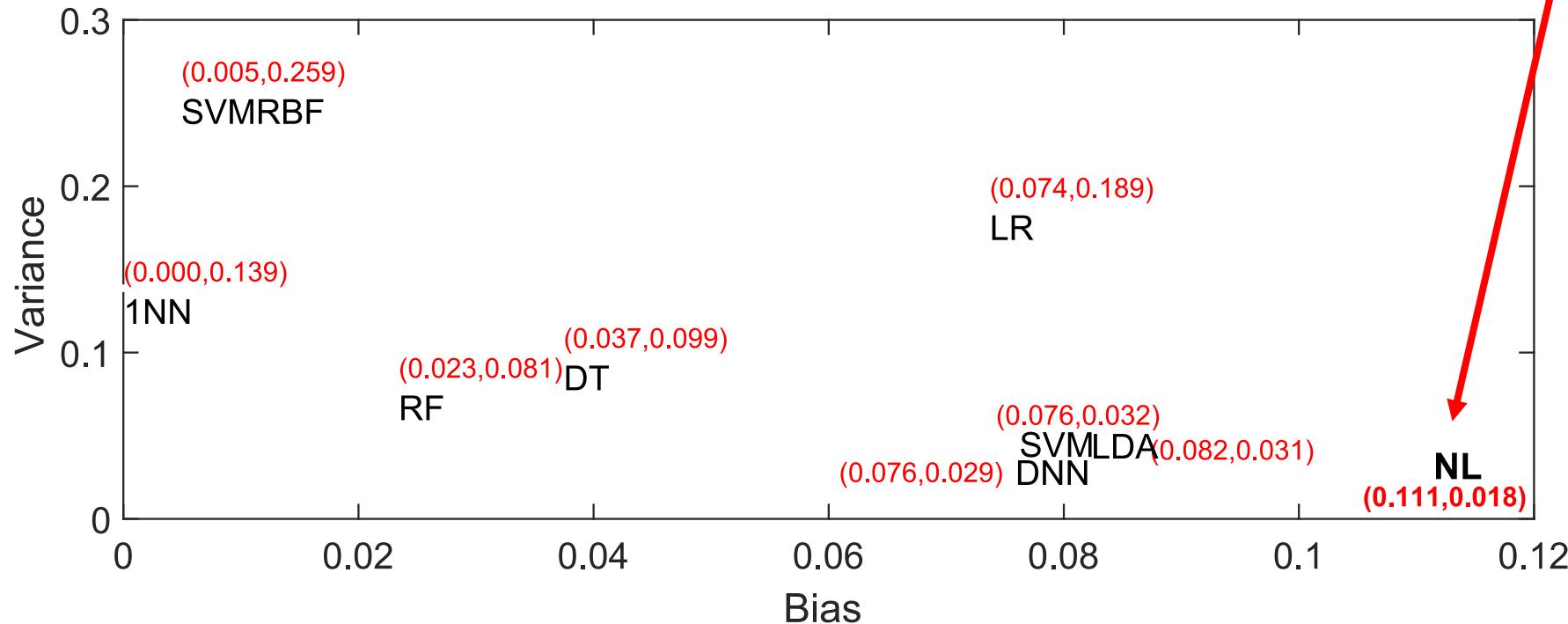


Winning Ratio(Accuracy)

Considering **simplicity** and **extreme sparsity** level of NL comparing black box models, this is an impressive result

Results: Average Bias-Variance

This **extraordinarily low variance** can be related to the simplicity of the model which results in **larger Rashomon ratio [1]** due to existence of noisy labels [2]



We observed several performance cases where test accuracy was considerably higher than train accuracy.

In **humans**, a study revealed that in certain situations, previously unseen prototypes might be classified more accurately during the testing phase than the original training stimuli [3].

[1] Breiman, Leo. "Statistical modeling: The two cultures (with comments and a rejoinder by the author)." *Statistical science* 16.3 (2001): 199-231.

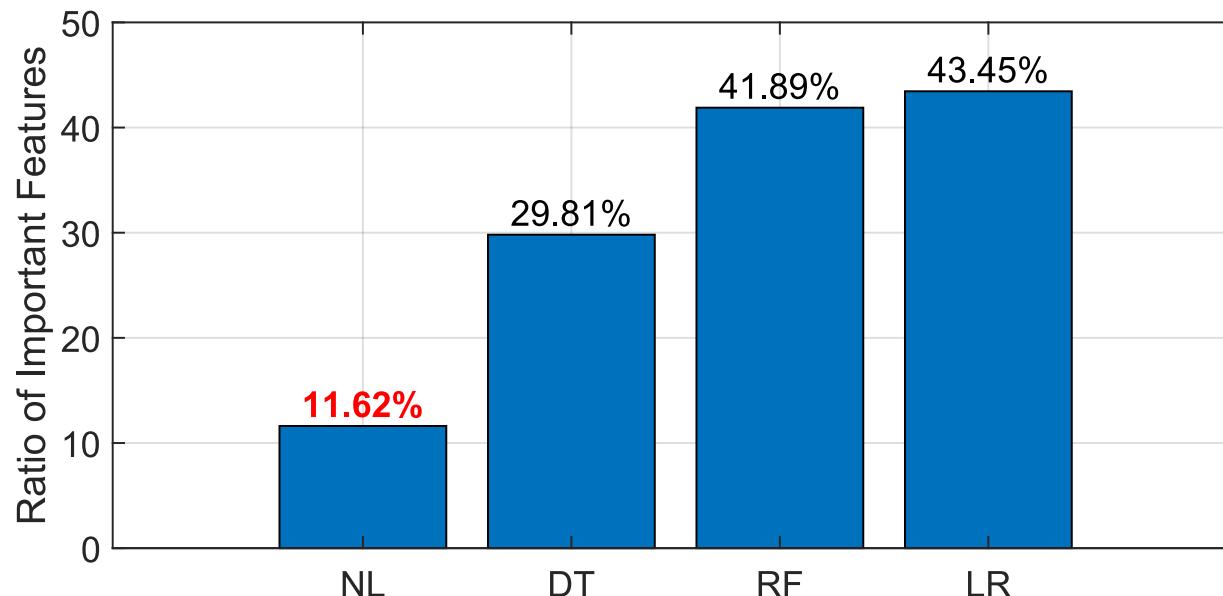
[2] Semenova, et al. "A Path to Simpler Models Starts With Noise.", NeurIPS 2023

[3] David R. Shanks, Concept Learning and Representation: Models in Smelser, Neil J., and Paul B. Baltes, eds. *International encyclopedia of the social & behavioral sciences*. Vol. 11. Amsterdam: Elsevier, 2015.

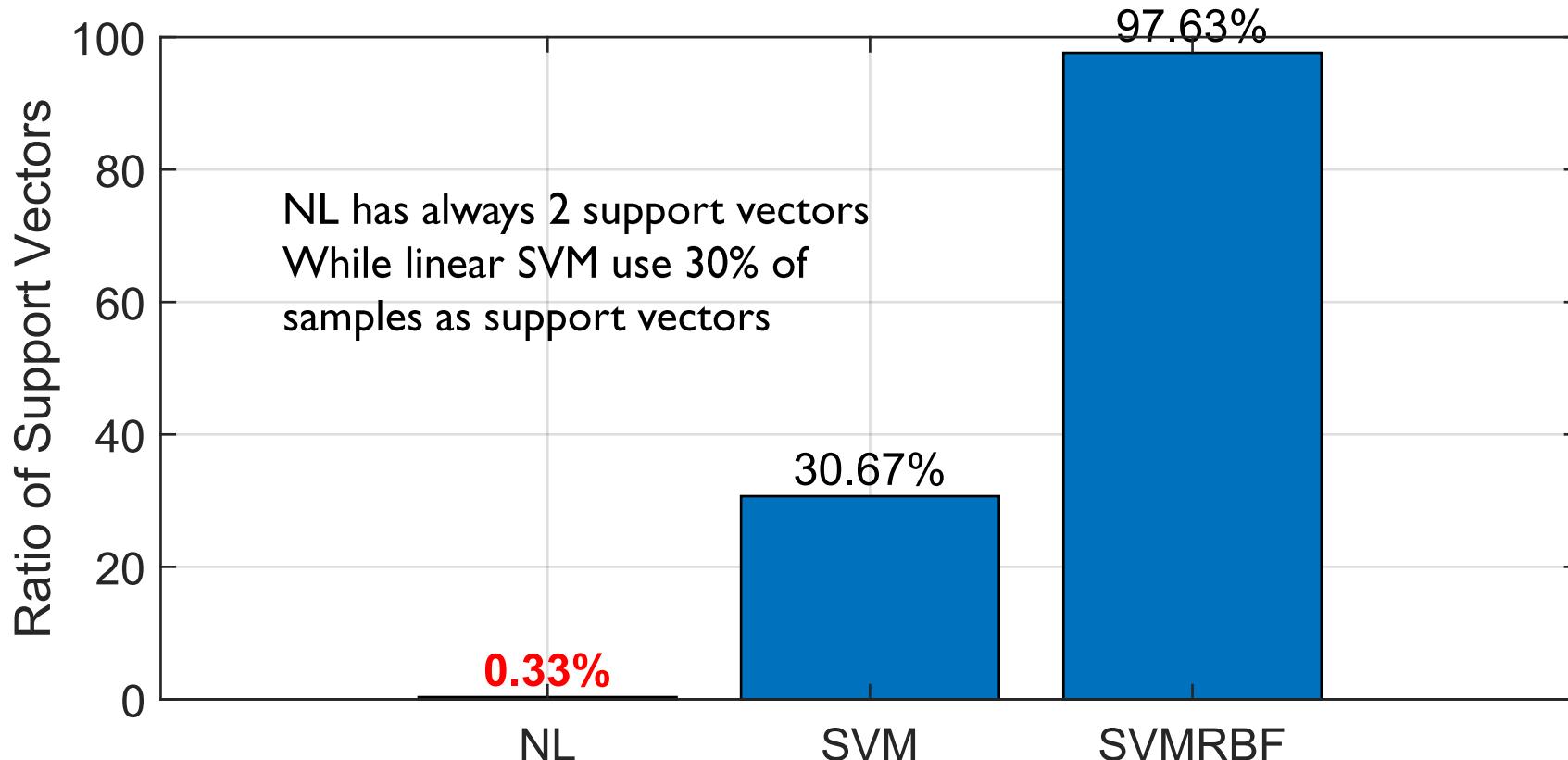
Results: Interpretability

As a quantitative metric, we compare the ratio of important features. But this **does not reflect the real interpretability value of NL**

- NL finds a meaningful subset of features with equal weights for each future
- Makes the interpretability even better than DT and LR

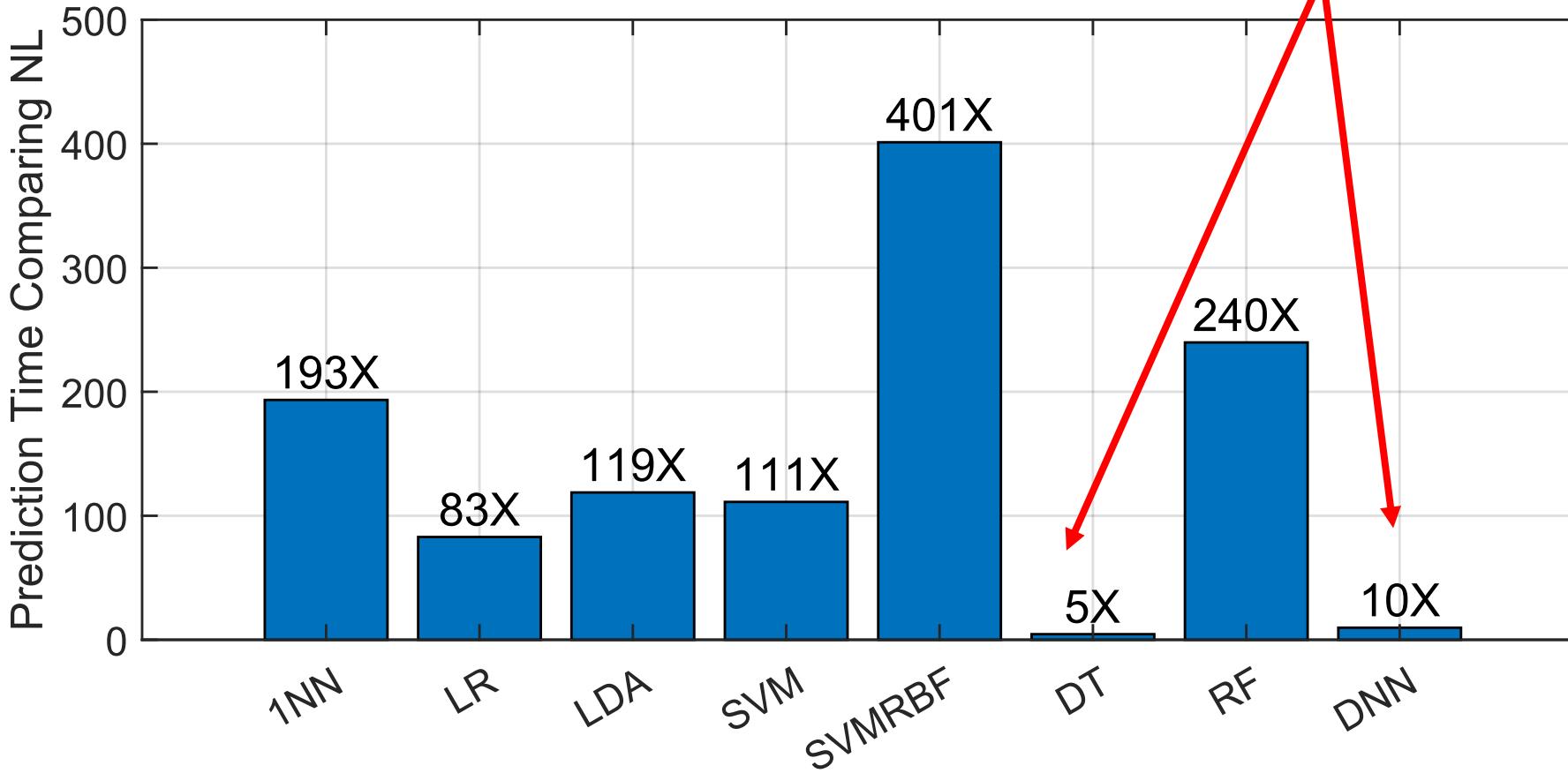


Results: Ratio of Support Vectors



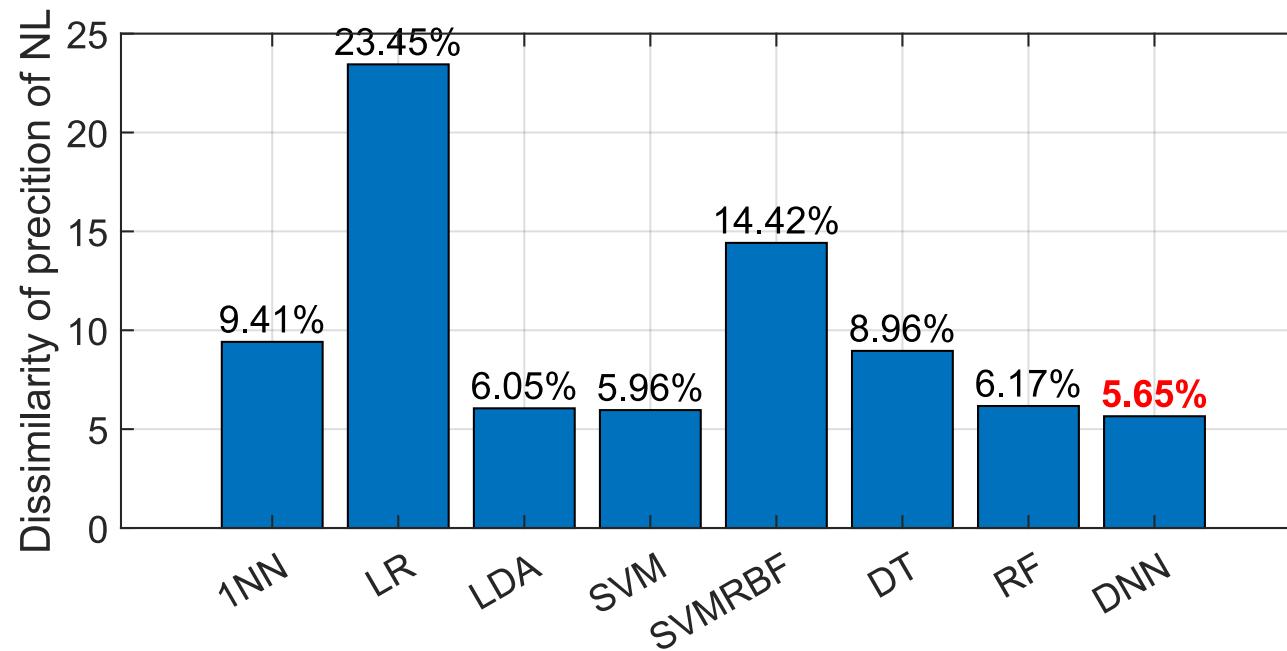
Results: Prediction Runtime

NL advances the **state-of-the-art prediction speed**, 5x of decision trees, and 10x of DNNs.



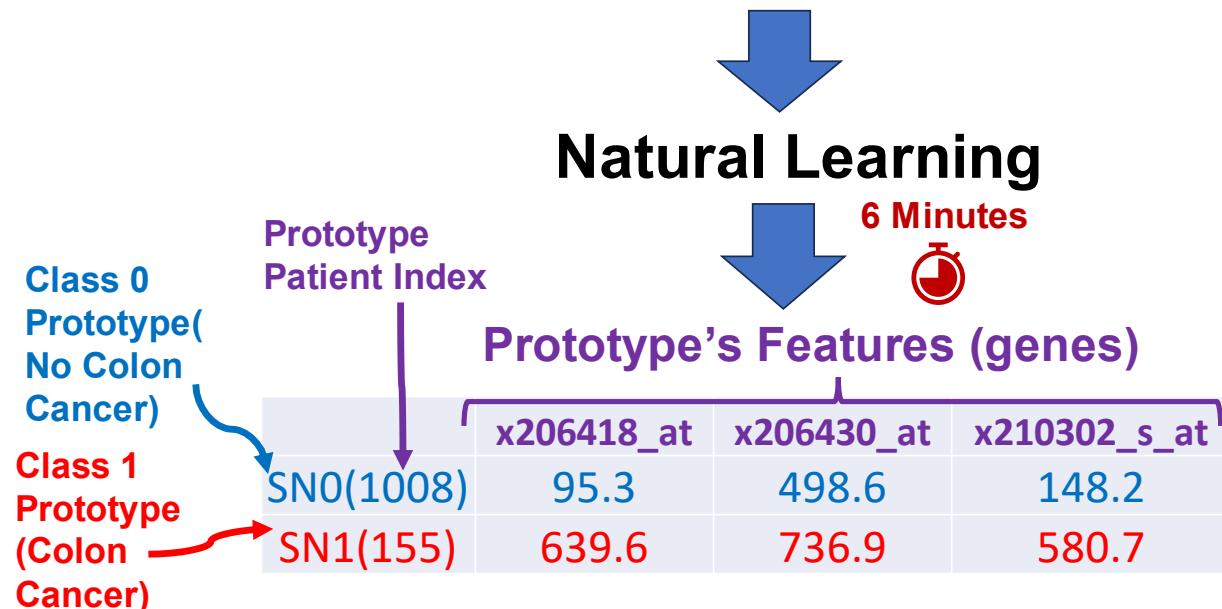
Prediction Dissimilarity to Other Classifiers

We compared the similarity of NL predictions to other classifiers based on 160k predictions they made on 170 training sets. Deep Neural Networks were found to be the best match with NL in terms of behavior on predictions, with a 5.65% mismatch in their predictions! The less similar classifier was logistic regression.



Example of NL Models: Colon Cancer Gene Expression

Dataset: **OVA_Colon** (1545 patients X 10935 genes)



Natural Learning

6 Minutes
⌚

Prototype's Features (genes)

Accuracy on Test Set

NL	98.05
DNN	97.40
RF	98.05
DT	98.05
SVM	96.75

Learned Prototype = (2 patients X 3 genes)

Model Sparsity Ratio = **3.55×10^{-7}**

Examples of Discovered Prototypes by NL and performance

Dataset: **ilpd-numeric**, n-fold=3/10 (seed=42)

Learned Prototype on the Train set (90%)

	v1	v6	v9	v10
Class 0(S#7)	26	16	3.5	1
Class 1(S#531)	22	14	3.8	1.1

Accuracy on Test Set (10%)

	NL	DNN	RF	DT	SVM
	70.69	67.24	67.24	67.24	67.24

Dataset: **thoracic-surgery**, n-fold=2/10 (seed=42)

Learned Prototype on the Train set (90%)

	v4_2	v4_3	v7_1	v7_2	v4_2
Class 1(S#51)	1	0	1	0	1
Class 0 (S#129)	0	1	0	1	0

Accuracy on Test Set (10%)

	NL	DNN	RF	DT	SVM
	89.36	89.36	87.23	65.96	89.36

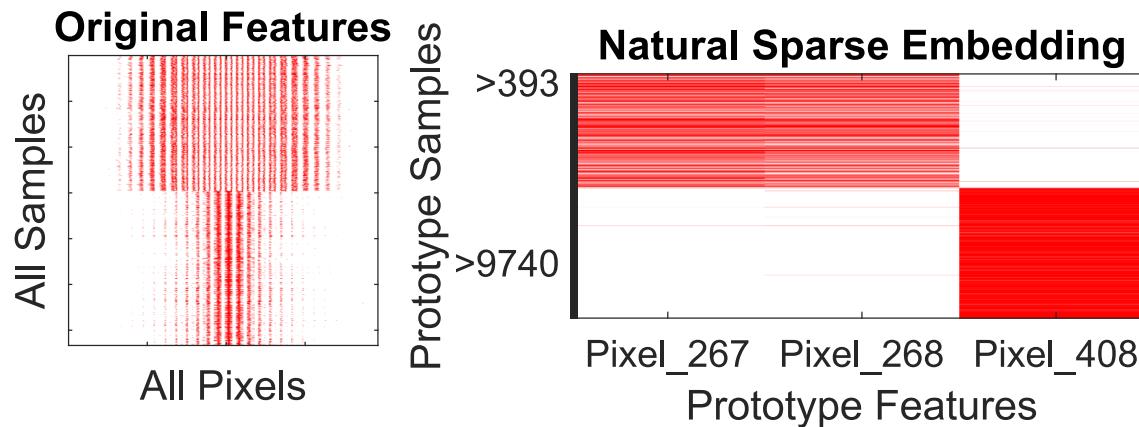
Dataset: **wdbc**, n-fold=4/10 (seed=42)

Learned Prototype on the Train set (90%)

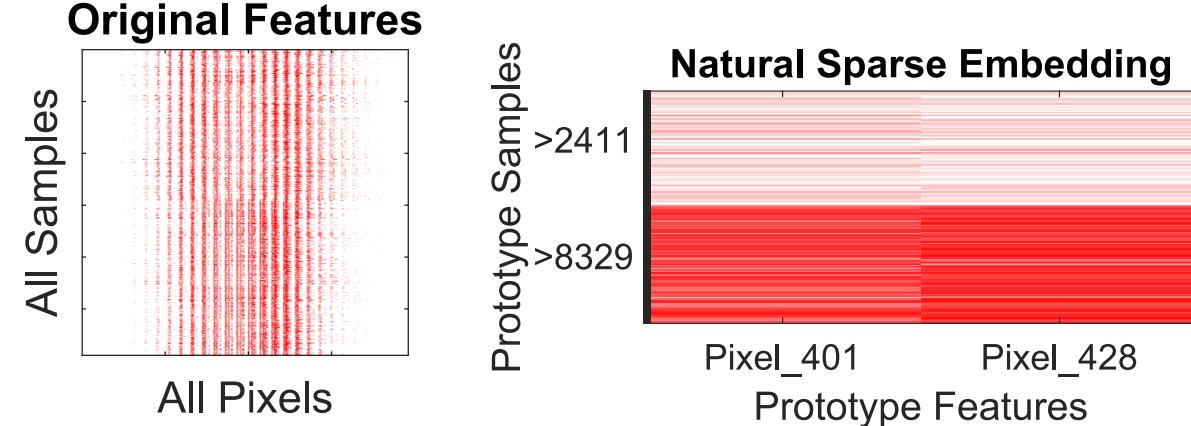
	v2	v11	v13	v18	v22	v23	v25	v28
Class 0 (S#348)	14.74	0.3428	2.537	0.01067	17.93	114.2	0.122	0.1251
Class 1 (S#206)	16.68	0.2711	1.974	0.00826	20.24	117.7	0.149	0.1252

Accuracy on Test Set (10%)

	NL	DNN	RF	DT	SVM
	98.25	98.25	94.74	94.74	98.25



MNIST 0 vs. 1 (The easiest) – Test Accuracy: 99.48%



MNIST 4 vs. 9 (The most difficult) – Test Accuracy: 85.64%

Advantages of Prototype Theory in Machine Learning

1. Model's **Transparency**: can be explained to **non-technical people**.
2. **Explainable Decisions**
 - Your loan is rejected because you resemble a rejected reference case compared to an accepted one.
3. **Interpretable Decisions**
 - Your loan is rejected because your **income and credit** are more like the rejected reference case than an accepted one.
4. **Fair rule**: Human-friendly reasoning with a **universal rule** that works for the **majority** (ideally, all).
5. Humans with **limited working memory** better understand the model due to the **extreme sparsity**
 - In Gene Expression data, OVA_Colon: (1545 patients \times 10935 genes) \rightarrow (2 patients \times 3 genes) \rightarrow the sparsity of 3.55×10^{-7}
6. Low model variance due to sparsity \rightarrow better **generalization** to **very different unseen cases**
7. **Ultra-fast prediction speed** due to small model size
8. **Simple to implement** and code: math-free, optimization-free, no package dependency
9. **Inherent robustness to noisy labels** (great applications in **healthcare**, criminal justice, **finance**)

Applications: Alternative for Decision Trees/Logistic Regression

- In applications prioritizing interpretability, explainability, and transparency, such as High-Stakes Decisions, where slight differences in accuracy are acceptable compared to black-box models, NL can **replace or complement decision trees and logistic regression** due to its **more accurate, simple, human-friendly, and fair explanations**

Applications: Performance near to Black-box with Noisy Labels

- In applications where humans are the sample, such as **healthcare, criminal justice, and finance**, NL can provide a high value.
 - In these domains, typically, **labels are noisy**, and black-box models provide the same performance as simple models
- Another reason: **existence of a prototype is guaranteed**
 - A **clinical case** in healthcare.
 - A **case study** in finance.
 - A **precedent case** in criminal justice.

Applications: key player in discriminant analysis of omics data

- In **discriminant analysis of high-dimensional omics data** (e.g., gene expression) NL can overcome the **curse of dimensionality** and the **challenge of limited samples** and generate **highly sparse and interpretable models** that are essential in these domains.

Applications: State-of-the-art in prediction speed

- **NL models are extremely small**, making them suitable for **real-time applications** where prediction speed is crucial
 - e.g., defense, online trading

Applications: State-of-the-art in embedded machine learning

- For **embedded systems** (e.g., wearable devices) where processing and memory constraints exist, NL's extremely sparse models require much lower computing resources (processing and memory)

Applications: Natural choice for binary input data

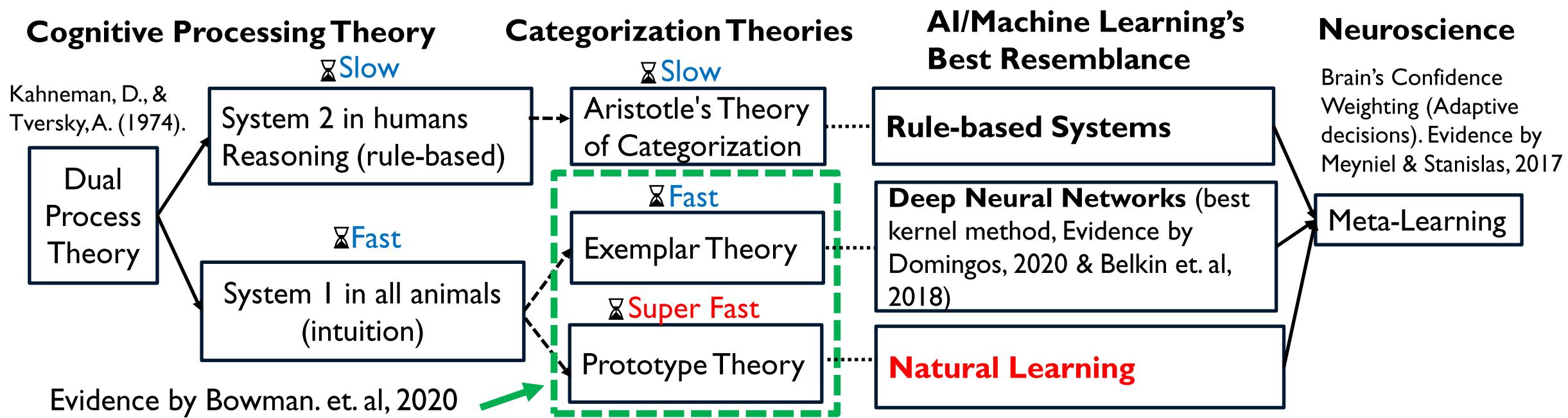
- In handling **high-dimensional binary data** where dimensionality reduction or representation learning do not provide added value, NL offers a promising alternative;

Applications: Ultra-fast classification of trivial cases in vision

- In the field of **vision**, **NL does not appear to be competitive** due to **lack of a mechanism for representation learning**
- NL can still be useful for **ultra-fast classifying of trivial cases**
 - (e.g., digit 0 vs. 1 in MINIST: 99.48% accuracy with model of 2x3 matrix)
 - frog vs. airplane in CIFAR-10, 86 % accuracy with model of 2x10 matrix)
 - Main applications
 - **better prediction speed**
 - **interpretability**
 - **explainability**
- It also can be used for detection of **data quality issues**

Dual Process Theory and Natural Learning

Natural Learning Emulates the Brain's System I 's superfast processing



Brain runs both exemplar and prototype categorization



RESEARCH ARTICLE



Tracking prototype and exemplar representations in the brain across learning

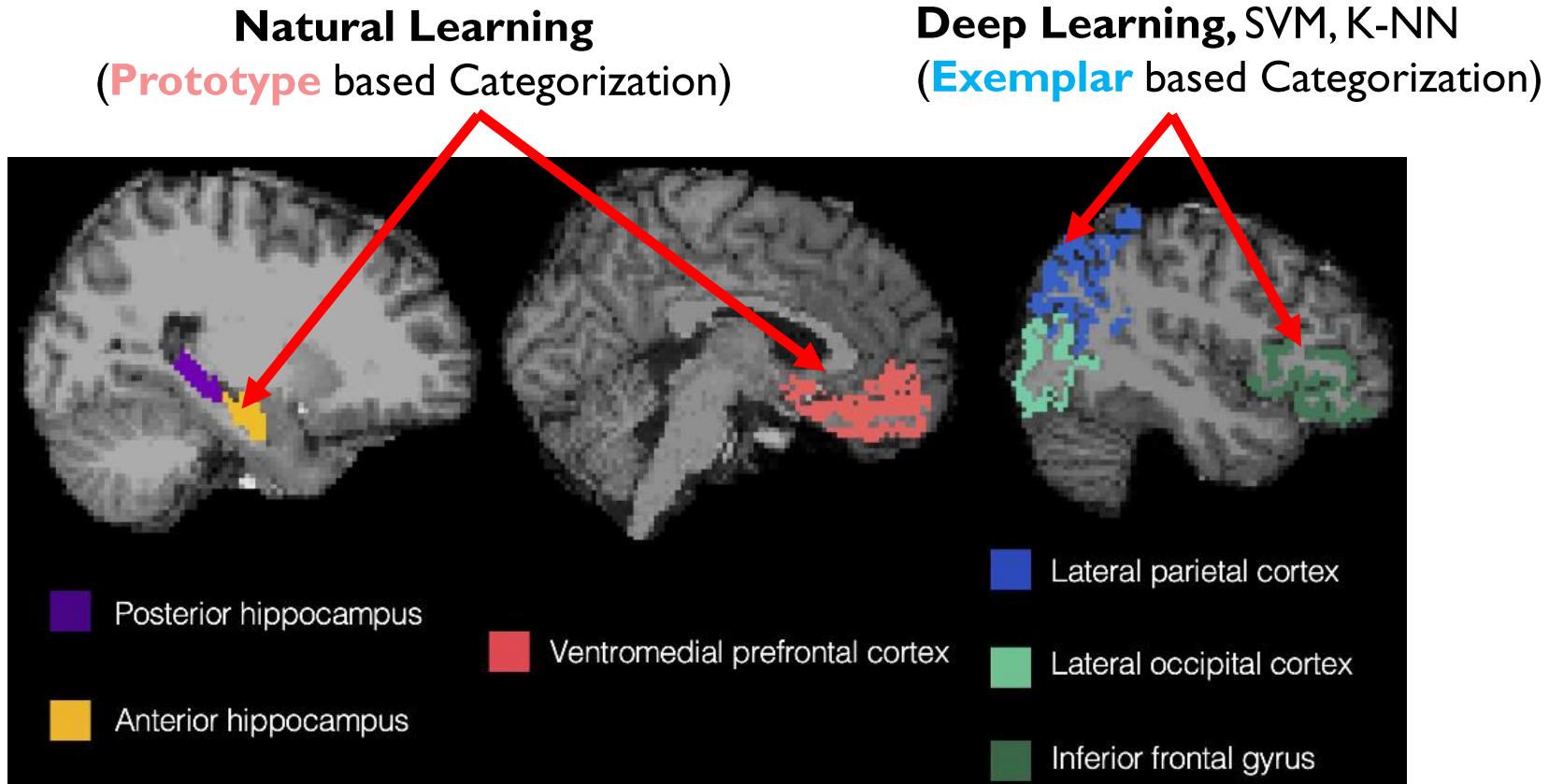
Caitlin R Bowman^{1,2*}, Takako Iwashita¹, Dagmar Zeithamova^{1*}

¹Department of Psychology, University of Oregon, Eugene, United States;

²Department of Psychology, University of Wisconsin-Milwaukee, Milwaukee, United States

Abstract There is a long-standing debate about whether categories are represented by individual category members (exemplars) or by the central tendency abstracted from individual members (prototypes). Neuroimaging studies have shown neural evidence for either exemplar representations or prototype representations, but not both. Presently, we asked whether it is possible for multiple types of category representations to exist within a single task. We designed a categorization task to promote both exemplar and prototype representations and tracked their formation across learning. We found only prototype correlates during the final test. However, interim tests interspersed throughout learning showed prototype and exemplar representations across distinct brain regions that aligned with previous studies: prototypes in ventromedial prefrontal cortex and anterior hippocampus and exemplars in inferior frontal gyrus and lateral parietal cortex. These findings indicate that, under the right circumstances, individuals may form representations at multiple levels of specificity, potentially facilitating a broad range of future decisions.

Brain runs both exemplar and prototype categorization



References

- **(Bowman. et. al, 2020)** Bowman, Caitlin R., Takako Iwashita, and Dagmar Zeithamova. "Tracking prototype and exemplar representations in the brain across learning." *elife* 9 (2020): e59360.
- **(Domingos, 2020)** Domingos, Pedro. "Every model learned by gradient descent is approximately a kernel machine." *arXiv preprint arXiv:2012.00152* (2020).
- **(Belkin et. al, 2018)** Belkin, Mikhail, Siyuan Ma, and Soumik Mandal. "To understand deep learning we need to understand kernel learning." *ICML 2018*
- **(Meyniel & Stanislas, 2017)** Meyniel, Florent, and Stanislas Dehaene. "Brain networks for confidence weighting and hierarchical inference during probabilistic learning." *Proceedings of the National Academy of Sciences* 114.19 (2017): E3859-E3868.

Conclusion

- Prototype theory has been recognized as a **Copernican revolution** in categorization theory because it departed from the Aristotelian rule-based approach.
- Now, we expect the same effect in machine learning: a transition from decision trees (**Aristotelian theory of categorization**) towards natural learning (**prototype theory of categorization**) that provides much better human-like reasoning and, as we showed, can be more accurate than decision trees.
- We anticipate that NL's sparse models will make a **high impact in providing new insights in many domains**.

Future Work

- Is it possible to implement a local representation learning in NL's local triplet space? We believe meaningful results in this direction can result in a **white-box version of DNNs**.
- How can we boost NL's performance without harming its attractive explainability?
- Is it possible to extend NL for regression?

Natural Learning

Learning Machine,
Inspired by Humans,
for Humans



Email: info@natural-learning.cc

Web: www.natural-learning.cc