

# An Introduction to Machine Learning



(c) CS U of Toronto

**Adam A Miller**

Northwestern/Adler Planetarium

2018 IDEAS Course

22 Feb 2018

# A Lecture in 3 Parts

**About me**

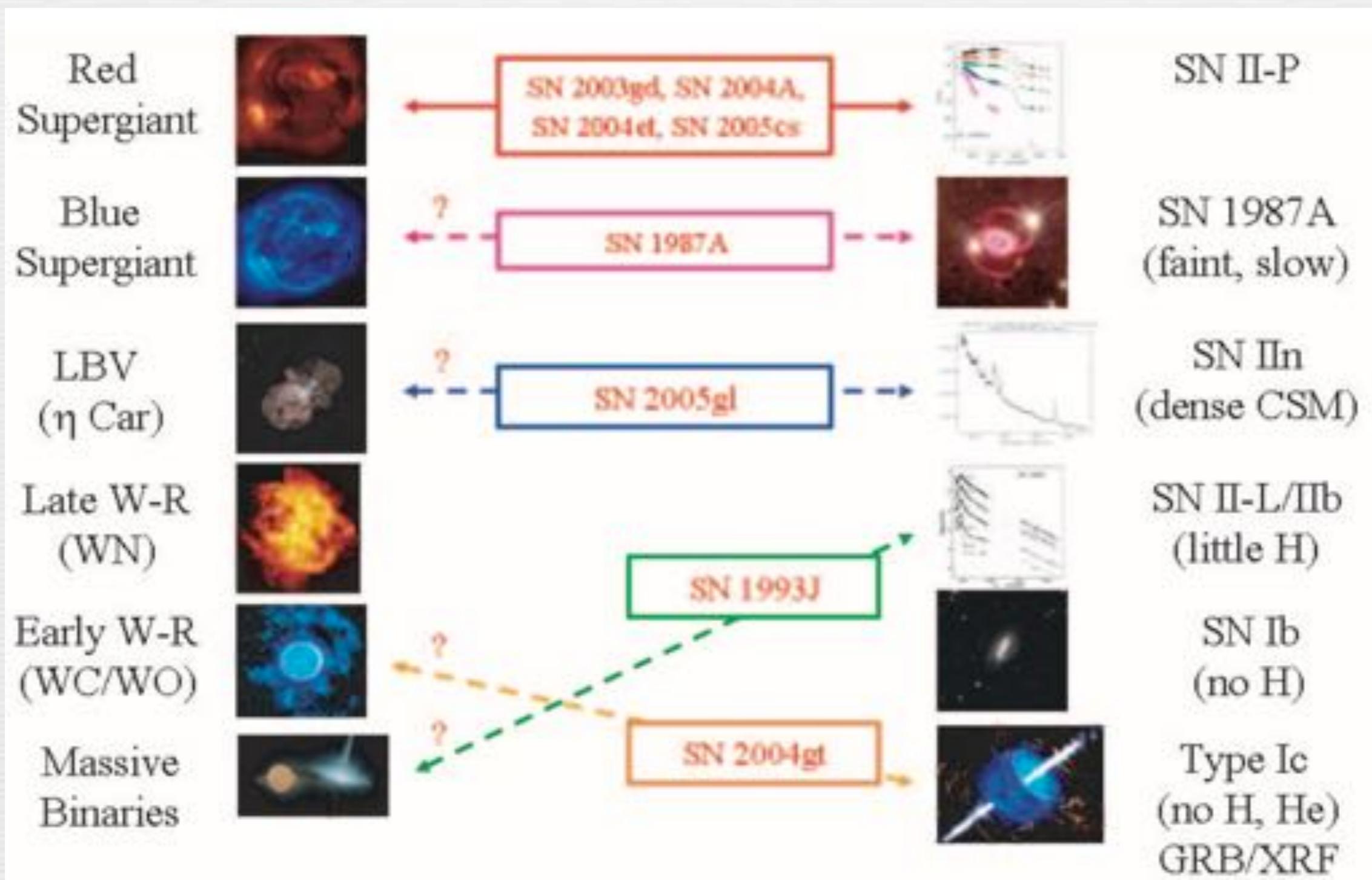
**An Introduction to Machine Learning**

**How to Build An End-to-End Machine Learning Pipeline**

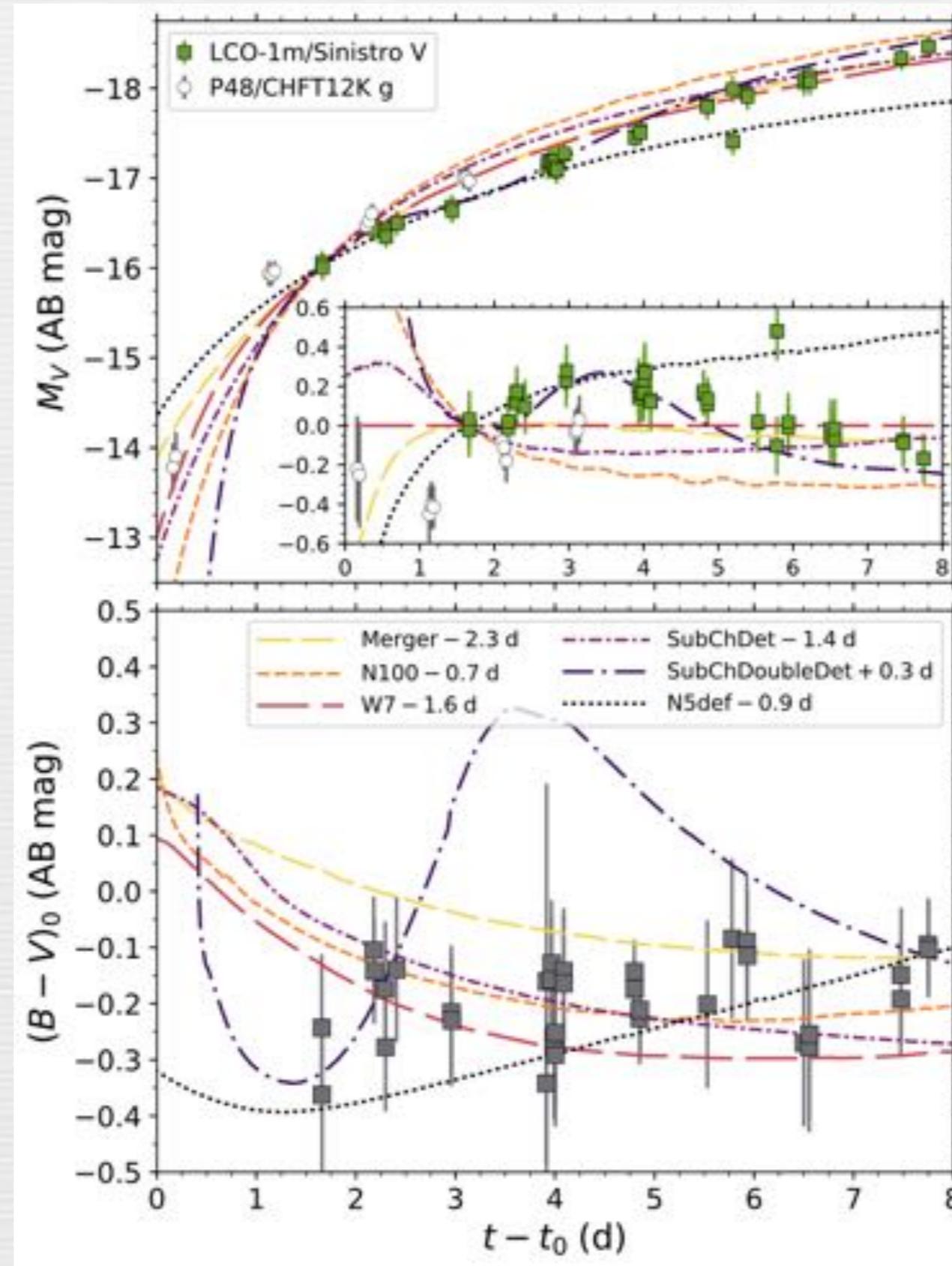
# Part I



# My Interests



# My Interests



# LSST



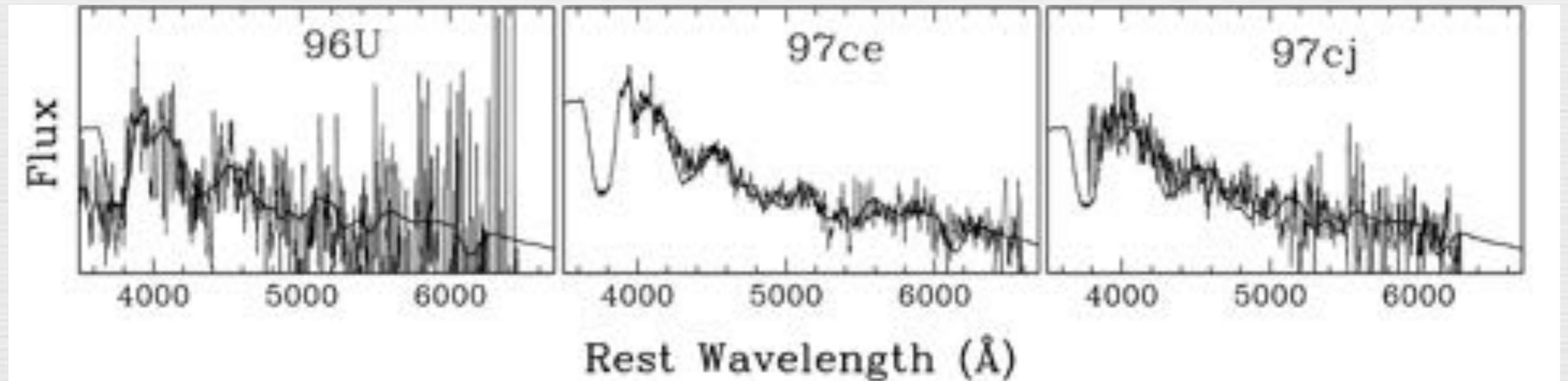
# SNe with LSST

LSST will discover ~2000 new SNe per night

Vast majority will be faint ( $m > 23$  mag)

~1000 hr/night on 8-m class telescopes needed for spec

~100 hr available on 8-m class telescope



Riess+98

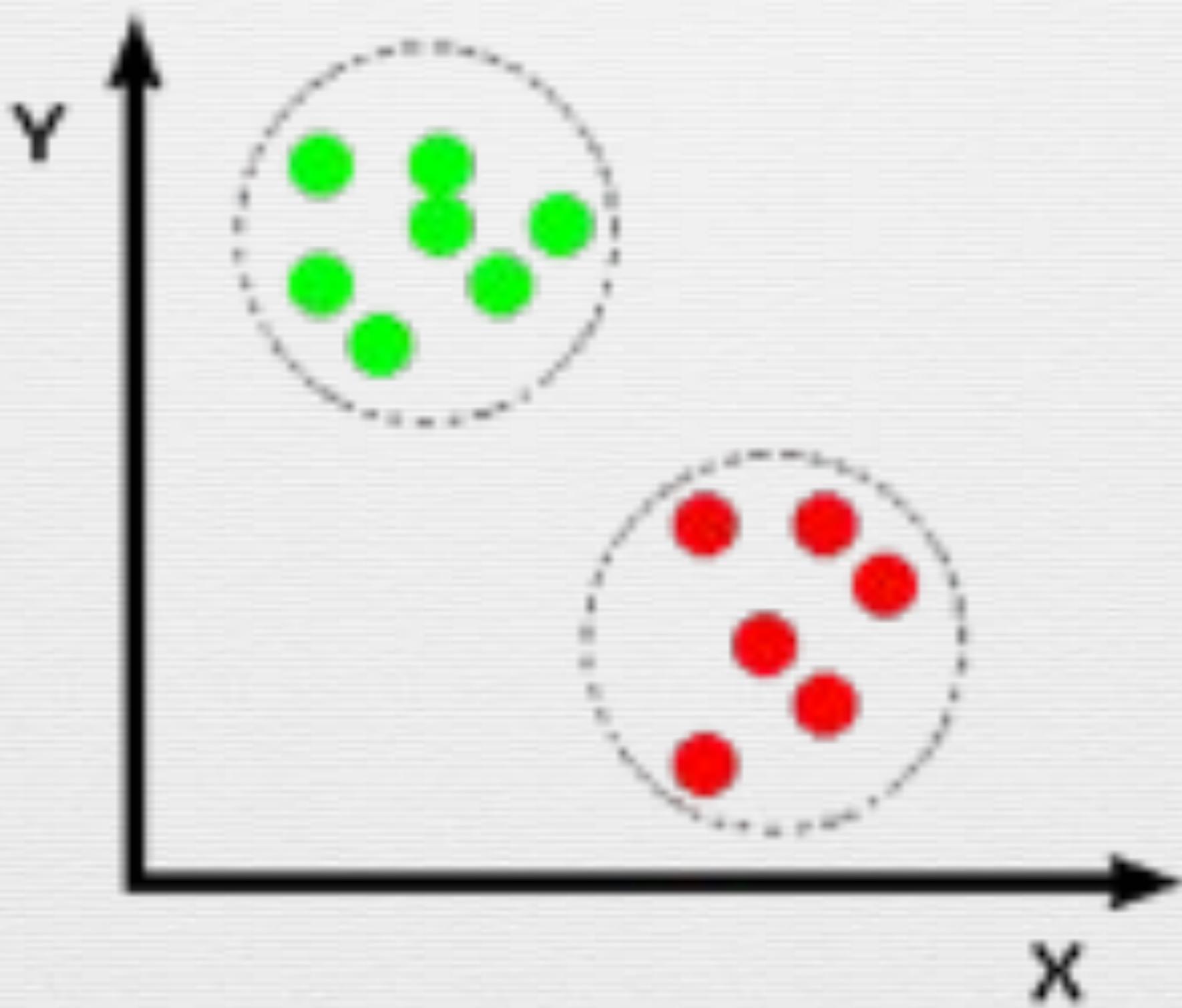
LSST will primarily be a **photometric-only** transient survey

# LSSTC DSFP



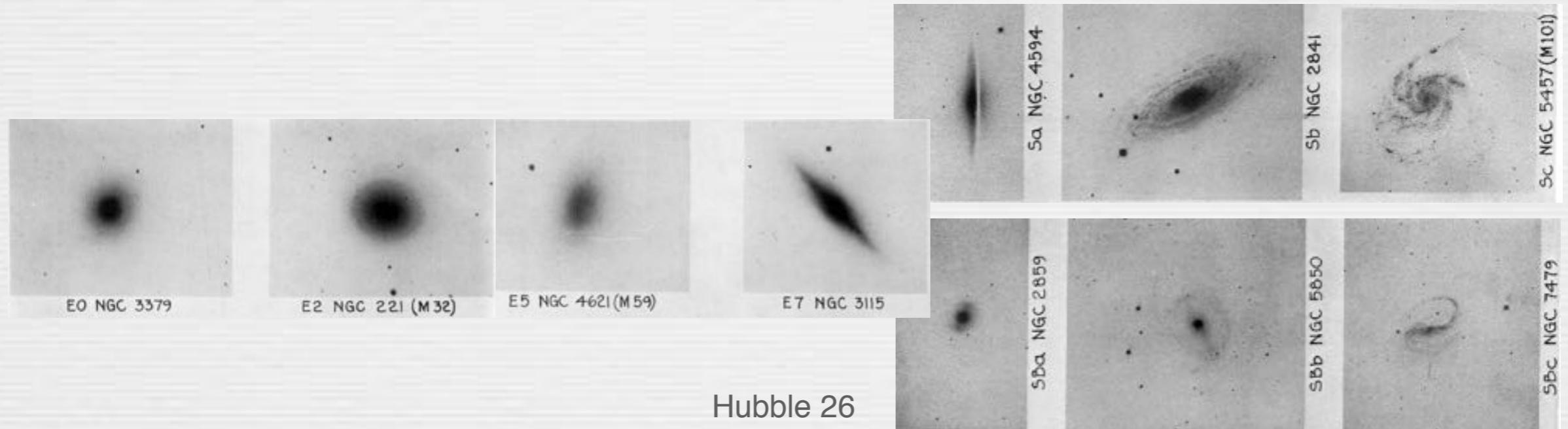
LSSTC  
**DATA  
SCIENCE**  
FELLOWSHIP PROGRAM

## Part II



credit: <https://homepages.inf.ed.ac.uk/rbf/HIPR2/classify.htm>

# Classification

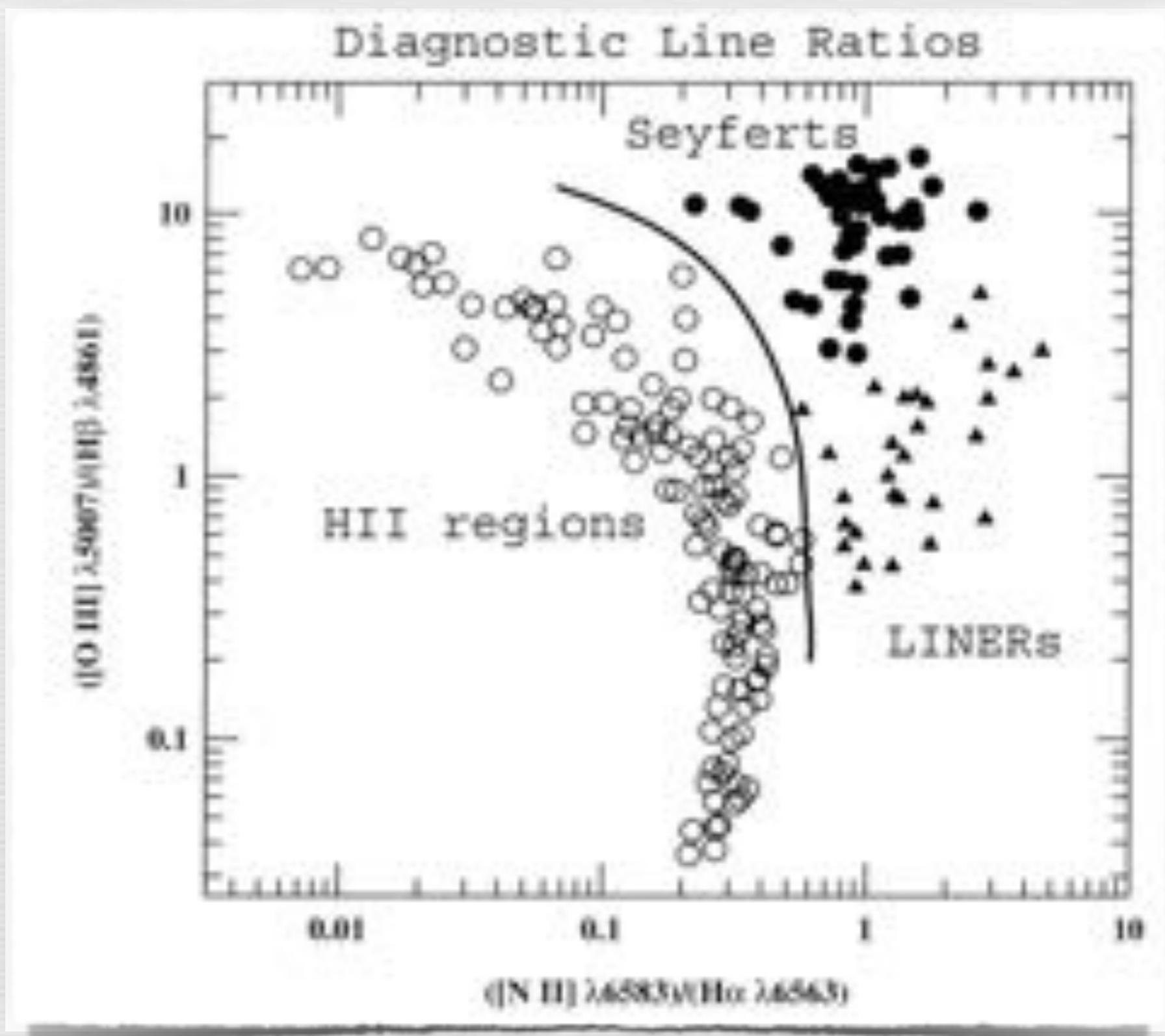


Fundamental problem for (nearly) all subfields of astronomy  
a lot of astro is essentially taxonomy

Classification schemes are (typically) well-argued, BUT  
subjective class boundaries are drawn  
constructed from small samples (then propagated forever)  
developed in low-dimensional spaces

# Classification

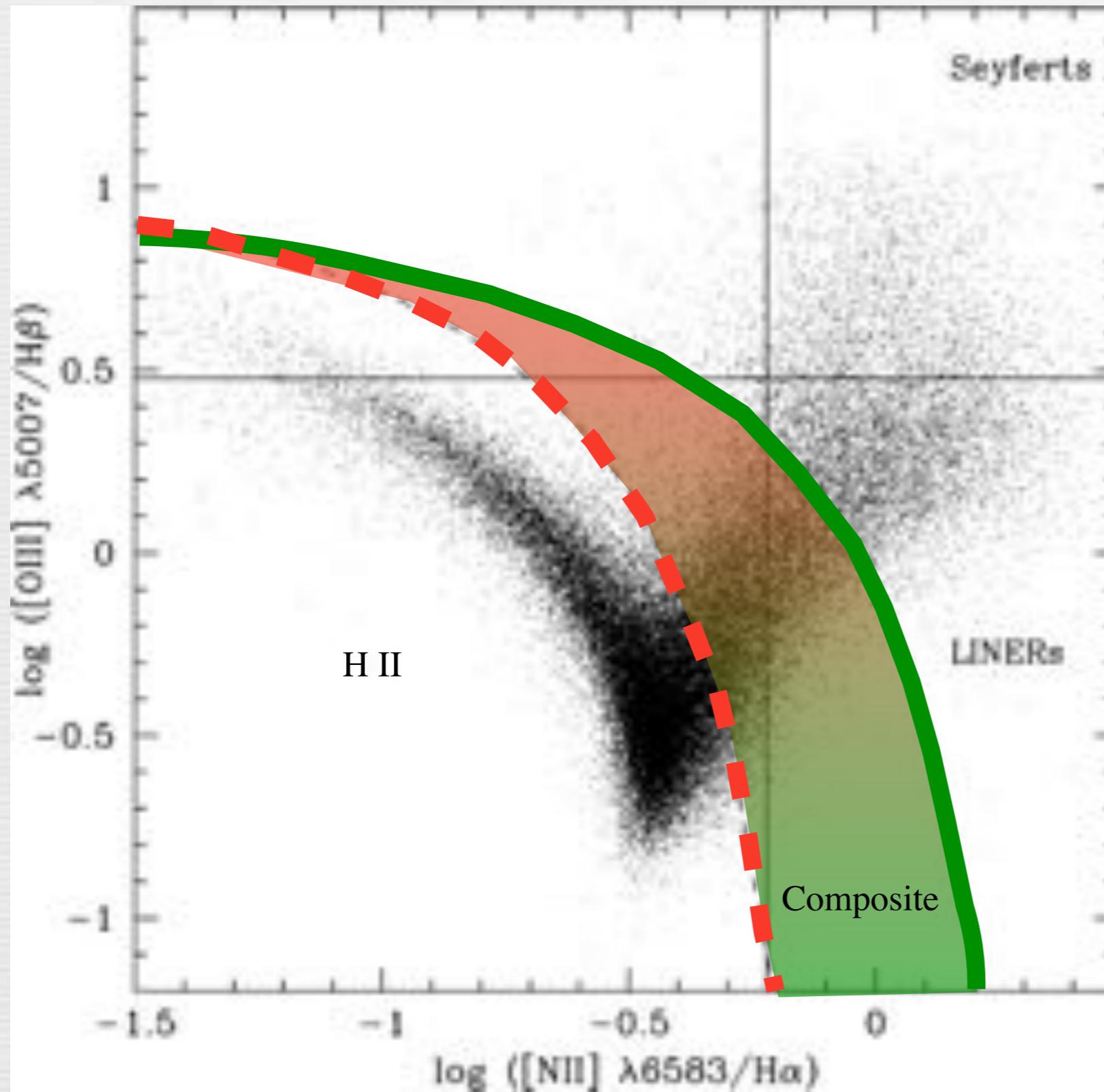
BPT circa 1987



credit: Mark Whittle

# Classification

BPT circa SDSS



continuous distribution

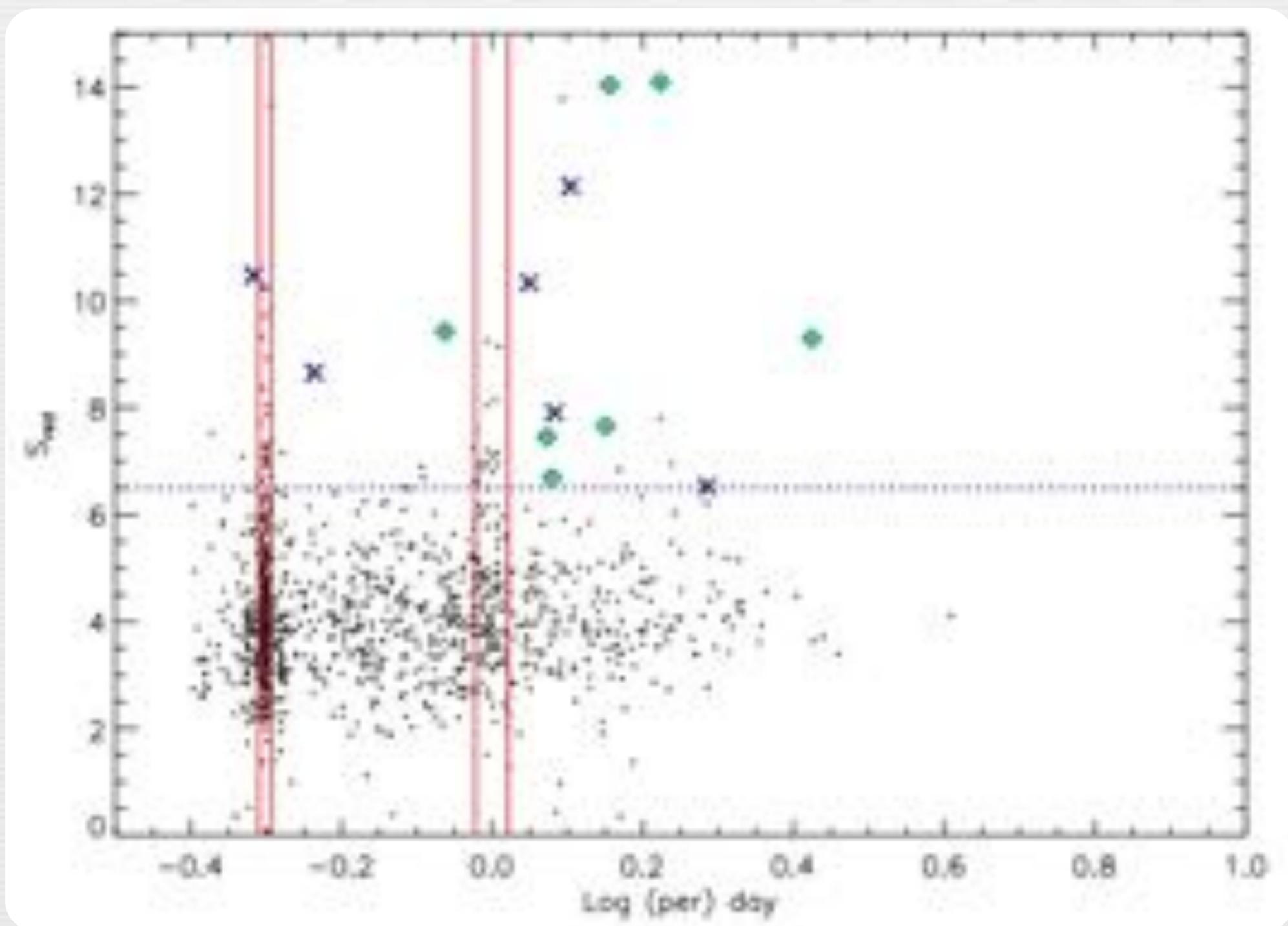
different class bounds

new (ill-defined) classes

Kauffmann+03

# Classification

I'm guilty too



# Classification

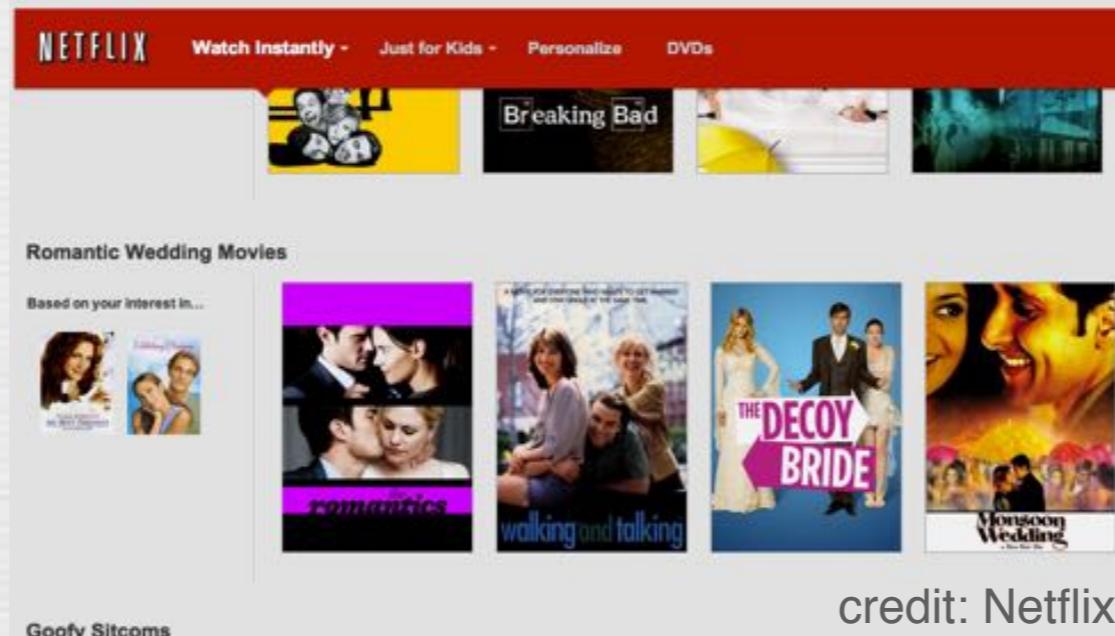
## Machine Learning

(aka - data mining, clustering, pattern recognition, AI (sorta) etc)

Fundamentally concerned with the problem of classification  
methods extend to regression as well

Address many challenges of classical taxonomy-like classification  
class boundaries drawn via (user-specified) optimization criteria  
improve and refine classifications with additional information  
can be constructed & developed in high-dimensional spaces

Examples: SPAM filters, Netflix, self-driving cars, etc



credit: SPAM



credit: Google

# Classification

## Machine Learning

two flavors:

### **labels are unknown**

#### **Unsupervised Learning**

In the feature space, the number, shape, & size of data groupings is unknown

Machine aims to cluster sources

No natural metric for measuring quality  
i.e. results vary from algorithm to algorithm

Can be very useful for data exploration

### **labels are partially known**

(labels are never fully known...)

#### **Supervised Learning**

Portion of data labeled by experts or expensive follow-up

Machine maps features ➤ labels

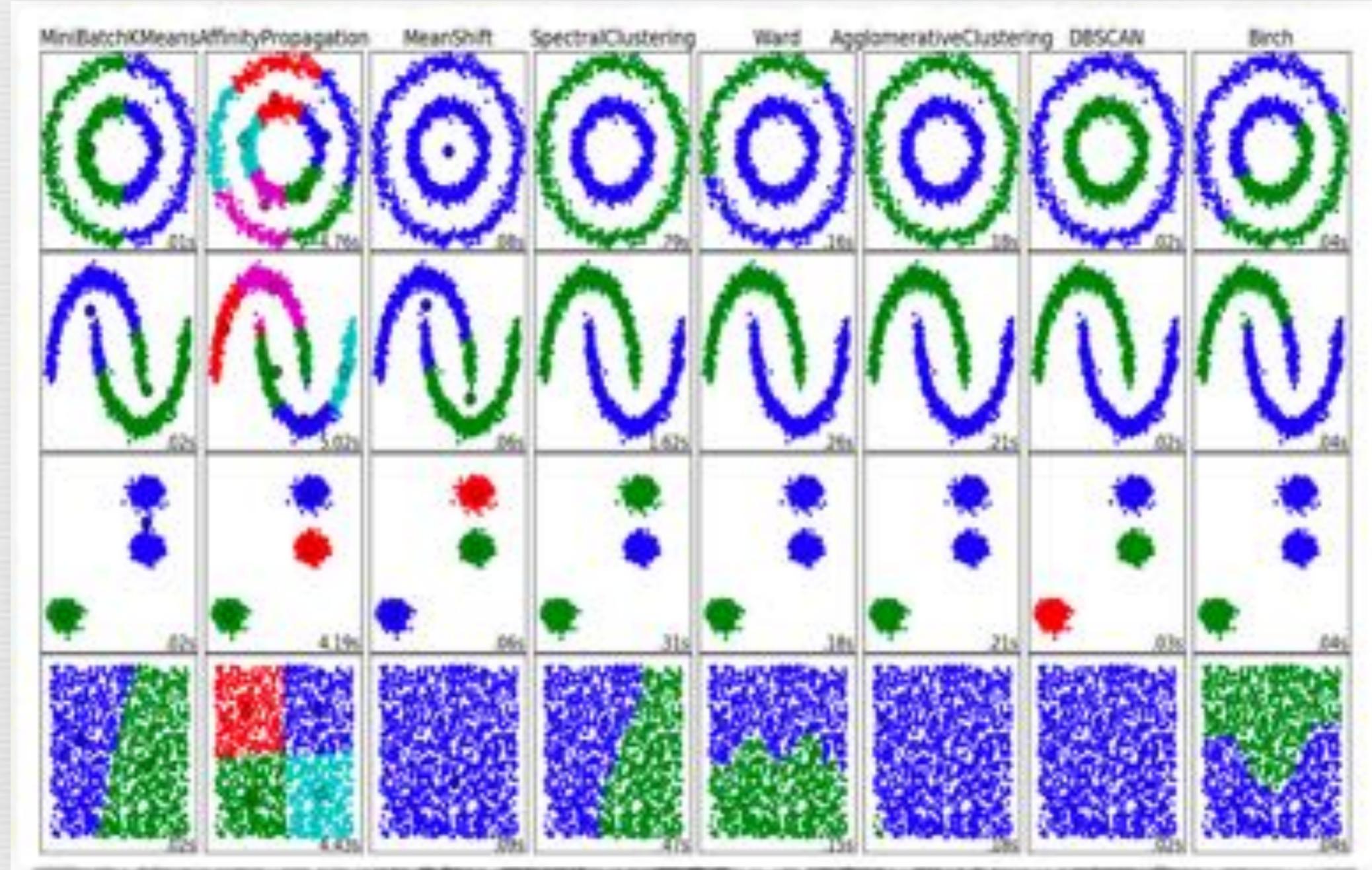
Can optimize accuracy or MSE  
results still vary from algorithm to algorithm

Useful for classification & regression

# Classification

## Machine Learning

### Unsupervised



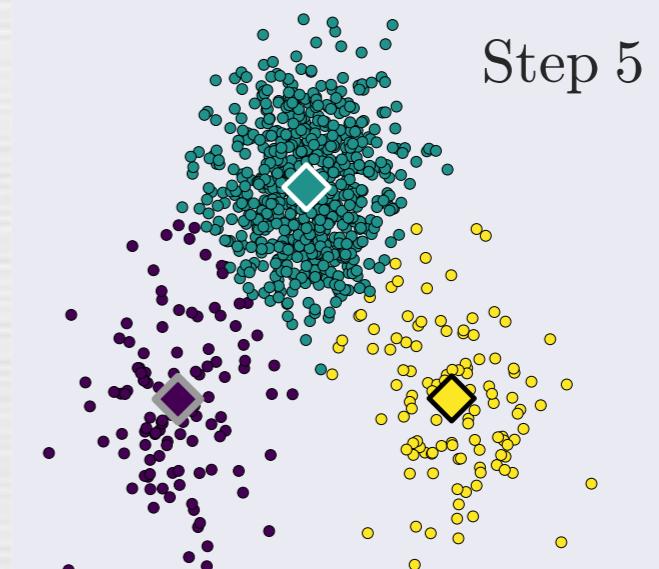
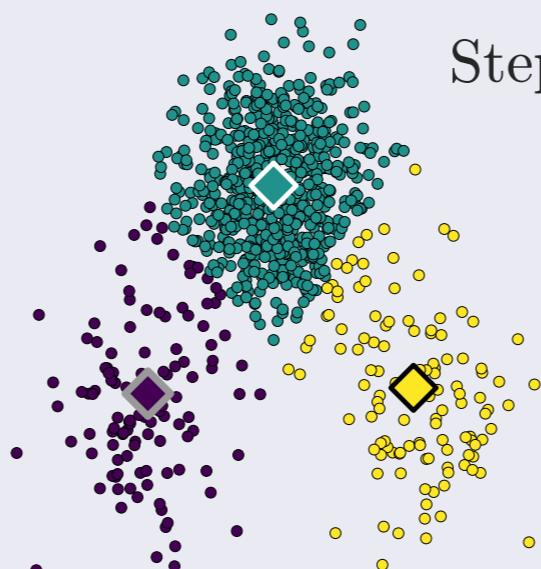
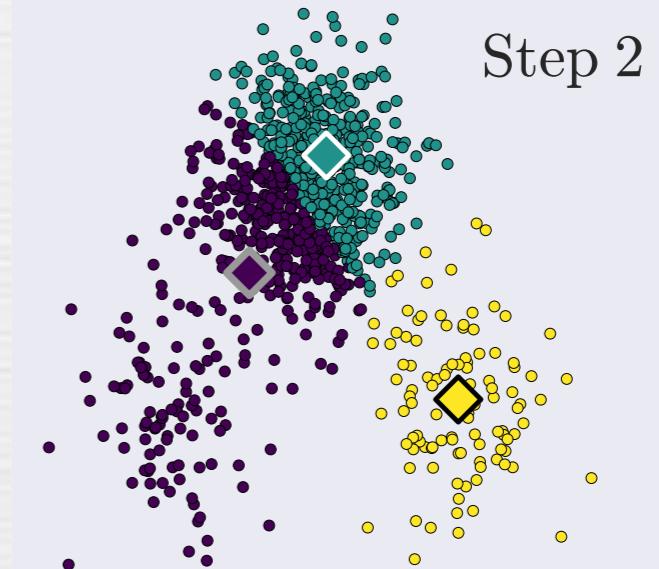
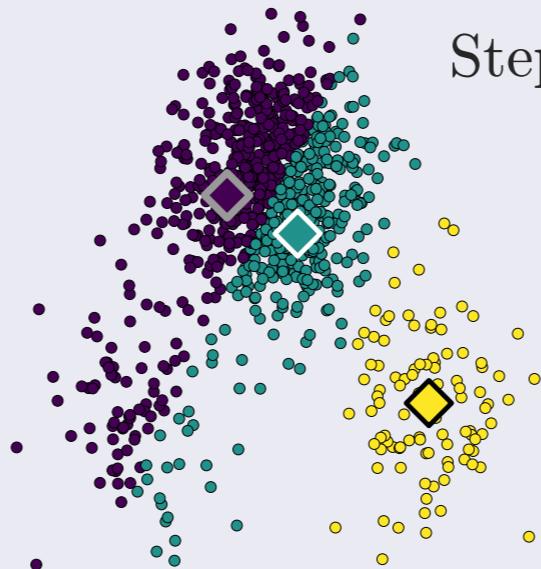
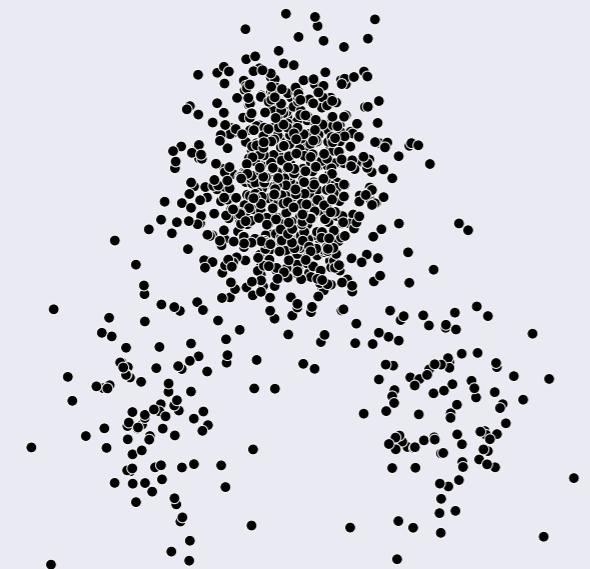
credit: scikit-learn

# Classification

## Machine Learning

Unsupervised

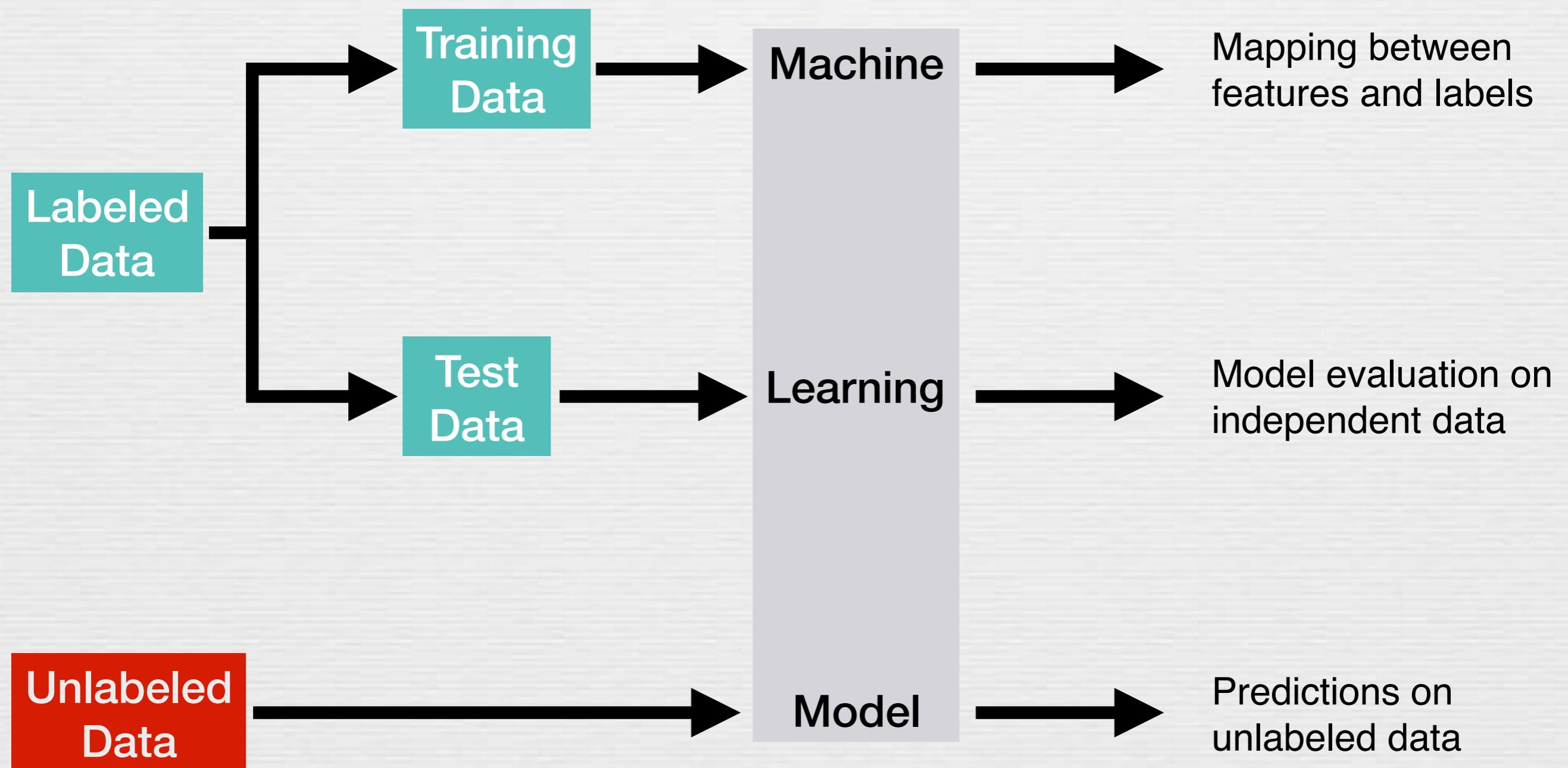
Famous algorithm: **K-means**



# Classification

## Machine Learning

Supervised

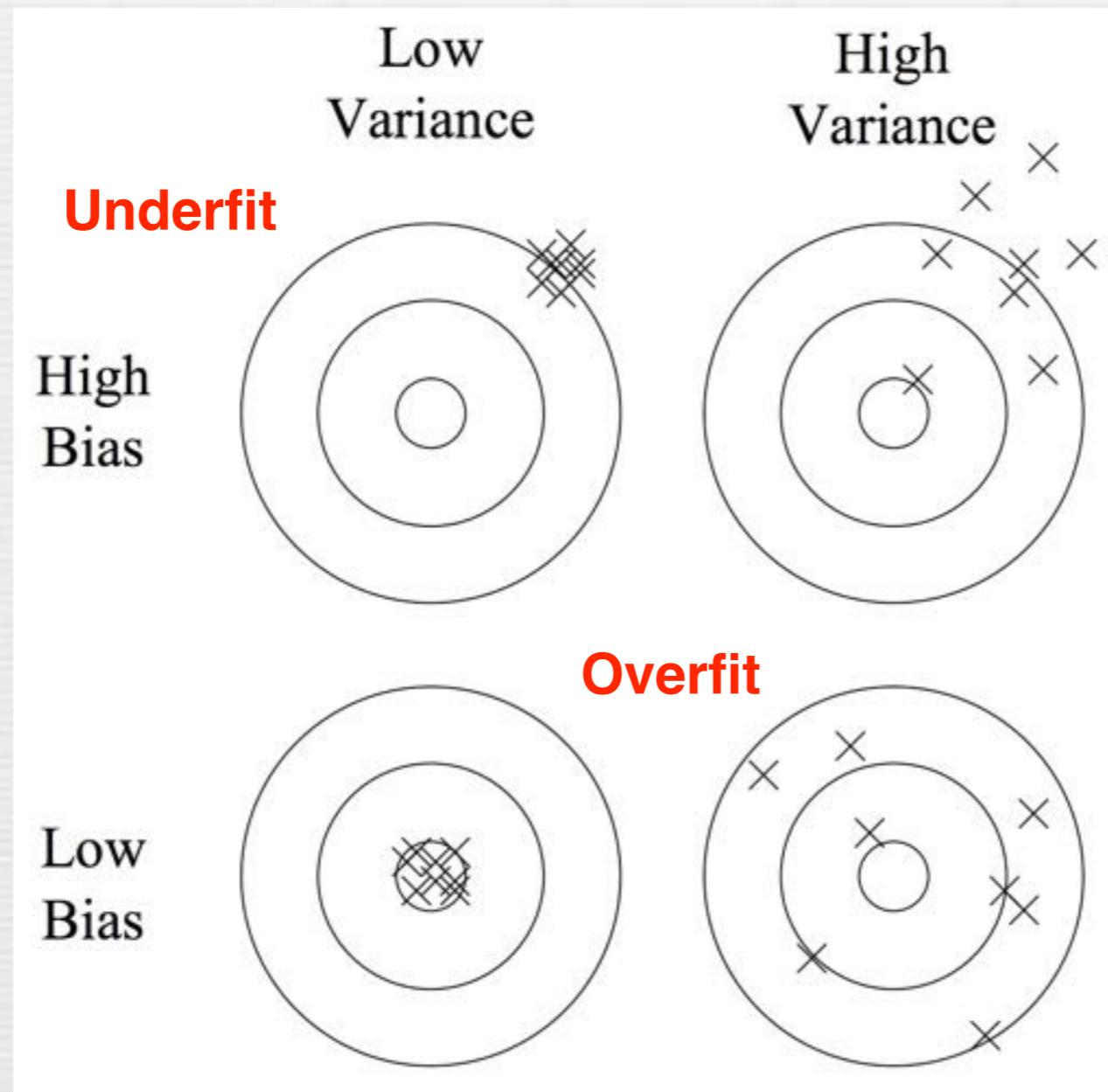


# Classification

## Machine Learning

Supervised

Goal: optimal trade off between bias and variance



credit: Arjun Krishnan

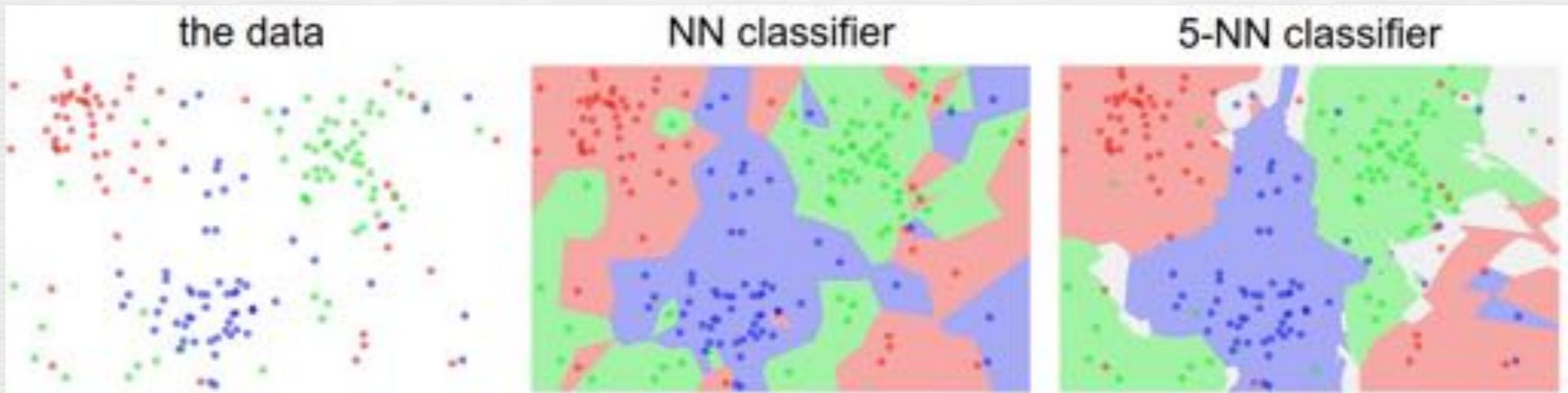
# Classification

## Machine Learning

Supervised

Famous algorithm: **k-nearest neighbors**

User specifies  $k$  ►  $k$  closest training set sources determine final classification



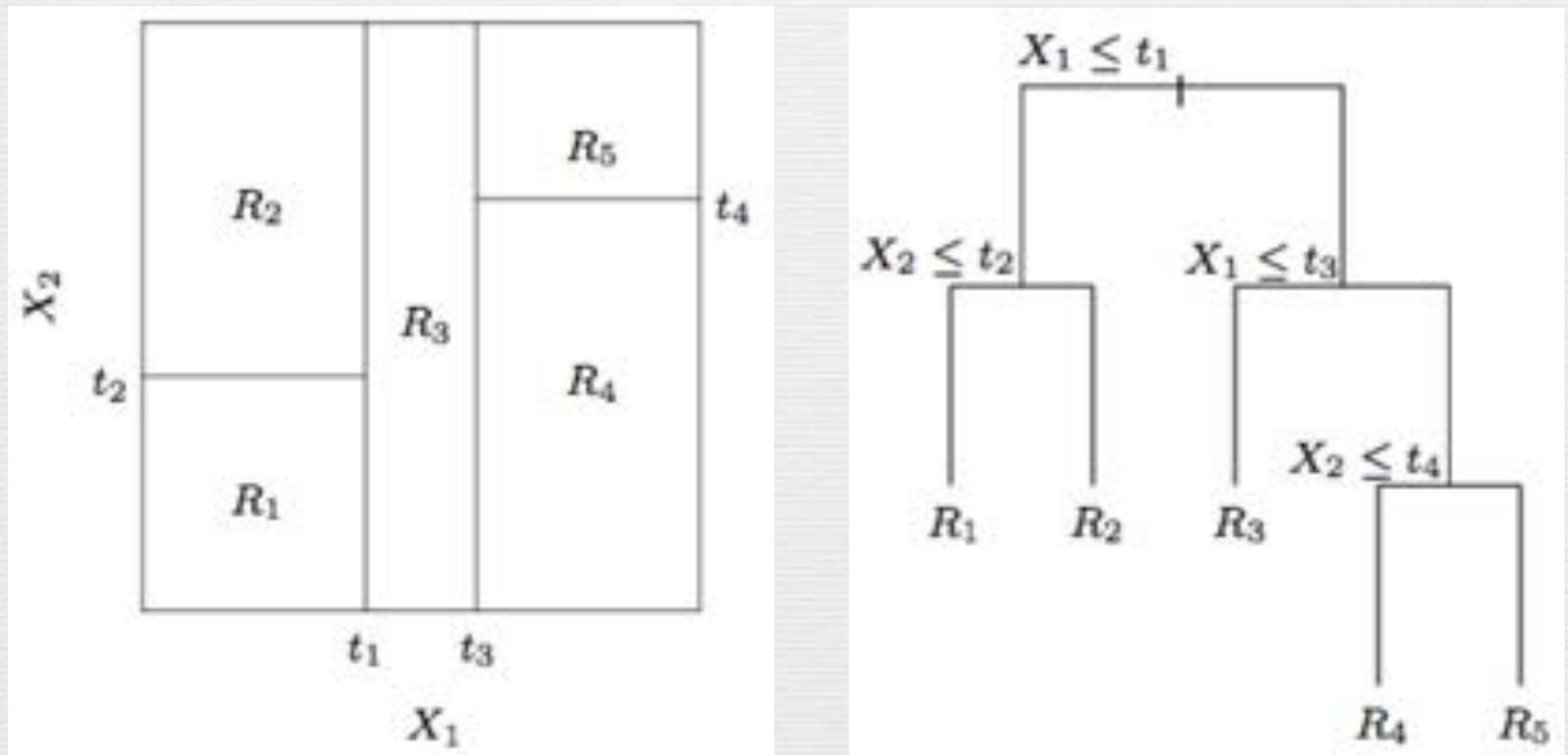
credit: <http://cs231n.github.io/classification/>

# Classification

## Machine Learning

Supervised

Famous algorithm: **Decision Tree**



# Classification

## Machine Learning

Supervised

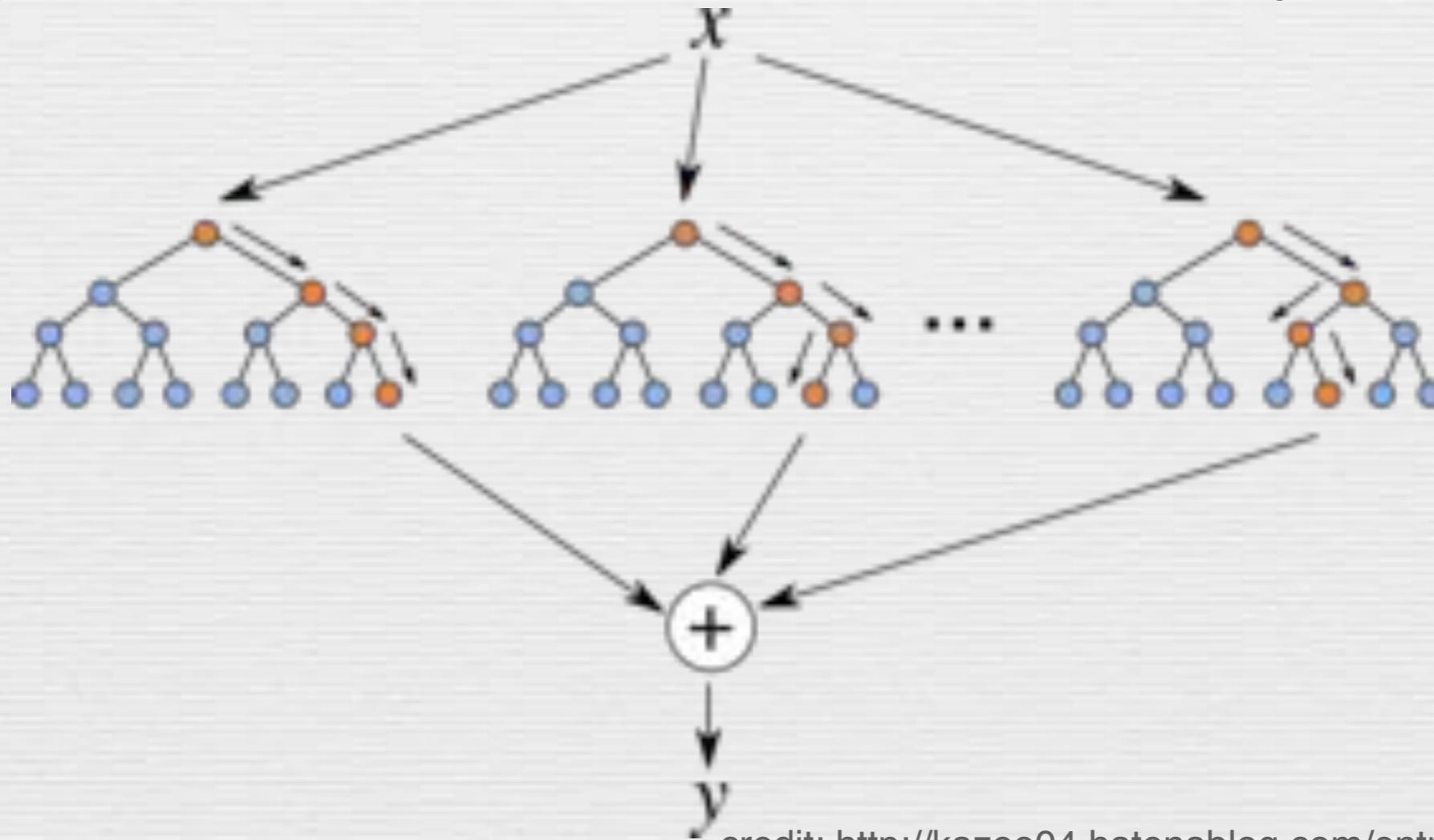
### Famous algorithm: **Random Forest**

Aggregates results from a collection of multiple decision trees

Use bagging (bootstrap w/ replacement) for each tree

Select only a random subset of features for split at each node

Average of de-correlated trees reduces variance relative to single tree



credit: <http://kazoo04.hatenablog.com/entry/2013/12/04/175402>

# sklearn Makes ML “Easy”

4 lines to construct a complex model

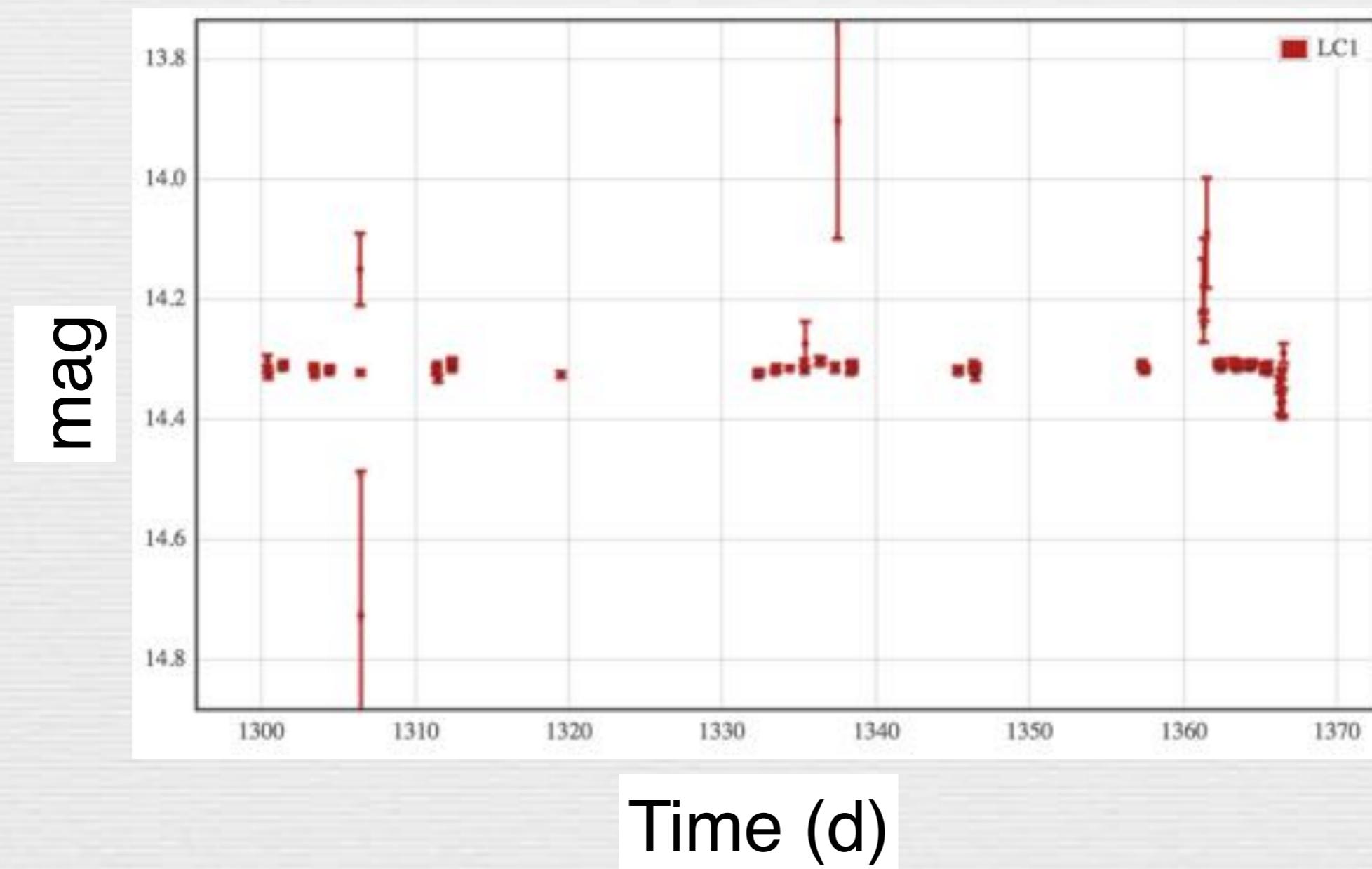
```
1 from sklearn import datasets  
2 from sklearn.ensemble import RandomForestClassifier  
3 iris = datasets.load_iris()  
4 RFclf = RandomForestClassifier().fit(iris.data, iris.target)
```

sklearn is great,  
but be weary of too good to be true

# Living Dangerously

## Crappy Data

Heteroskedastic Errors



# Living Dangerously

Crappy Data

Faint Objects

# Living Dangerously

Crappy Data

Faint Objects



# Living Dangerously

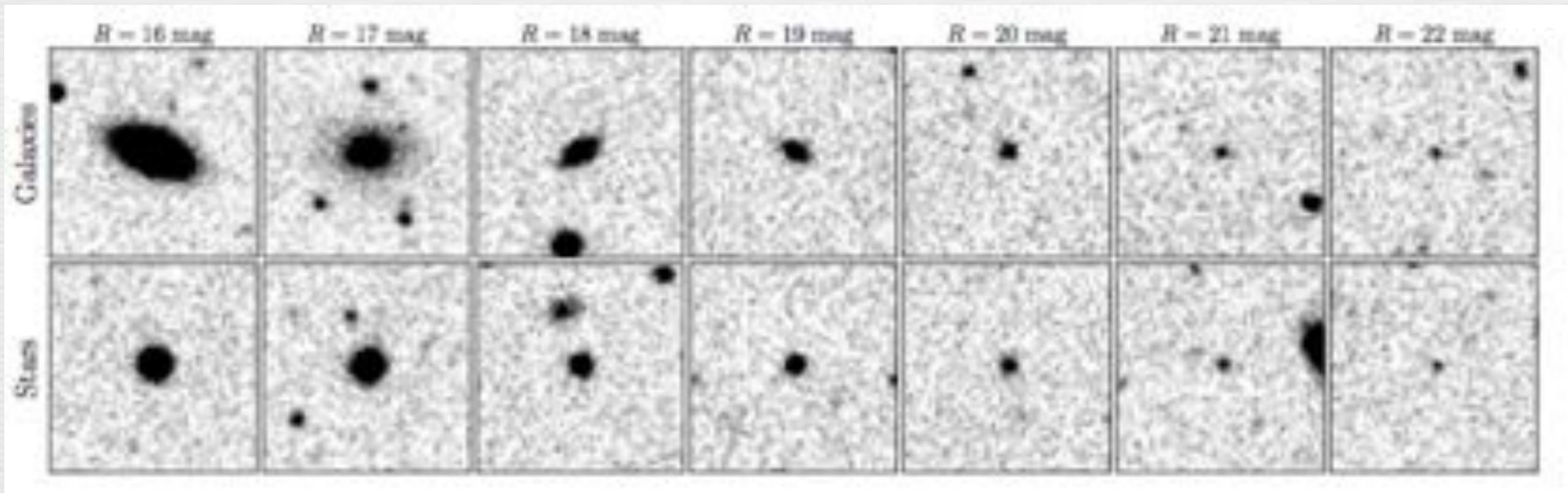
## Star-Galaxy Separation

summer student project

“easy” two class RF model

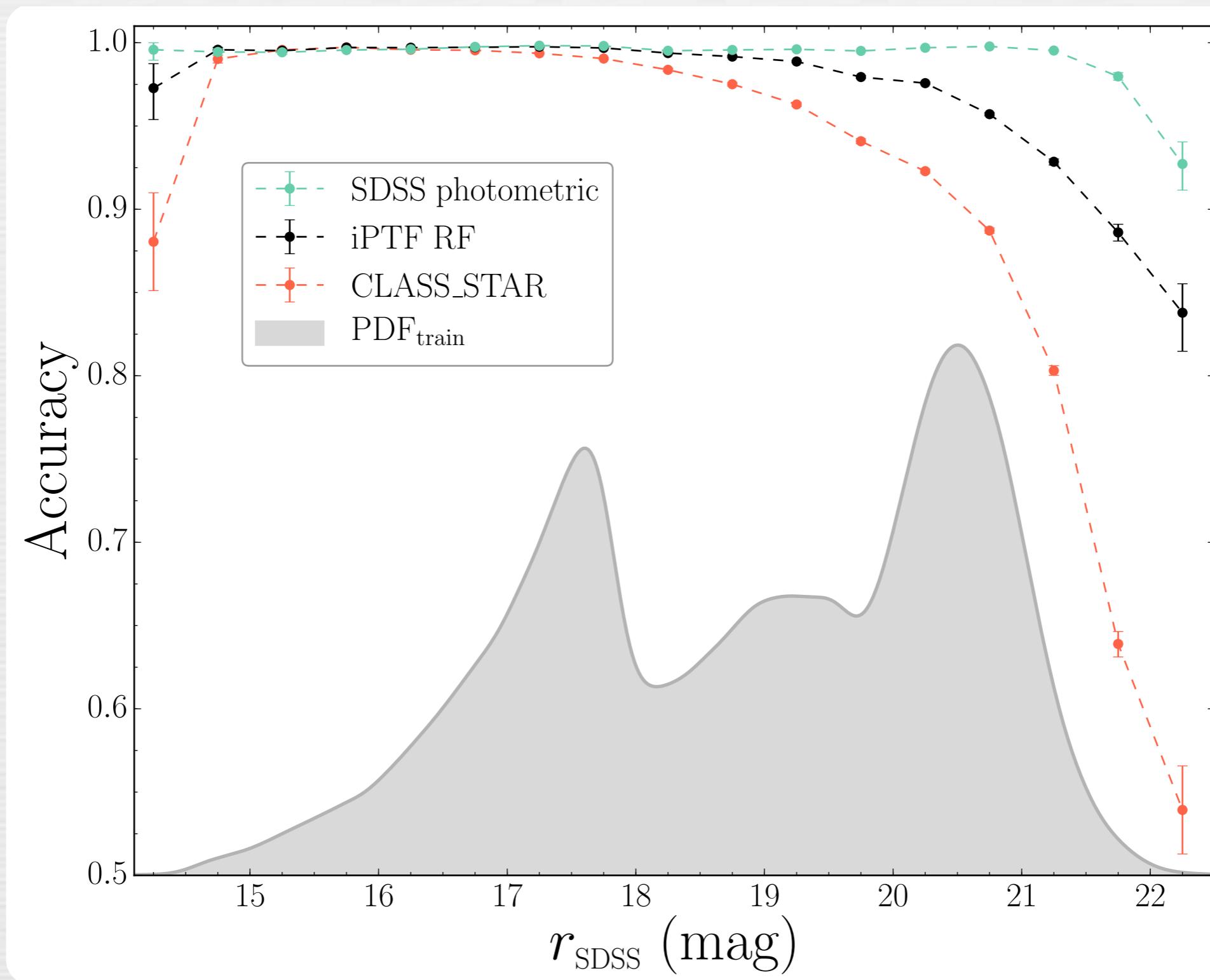
facilitate discovery in PTF/improve search for GW counterparts

AAM+16



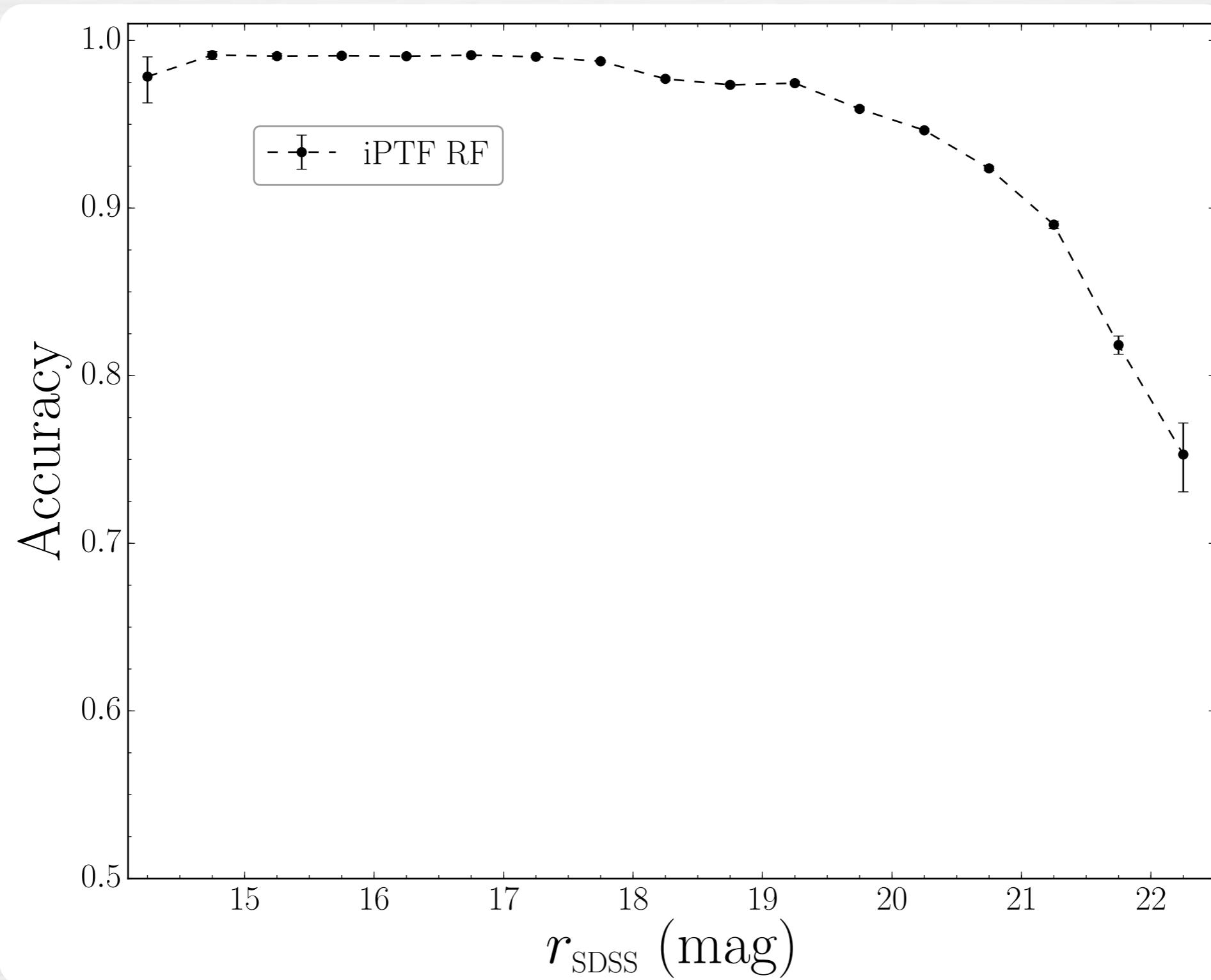
# Living Dangerously

## Star-Galaxy Separation



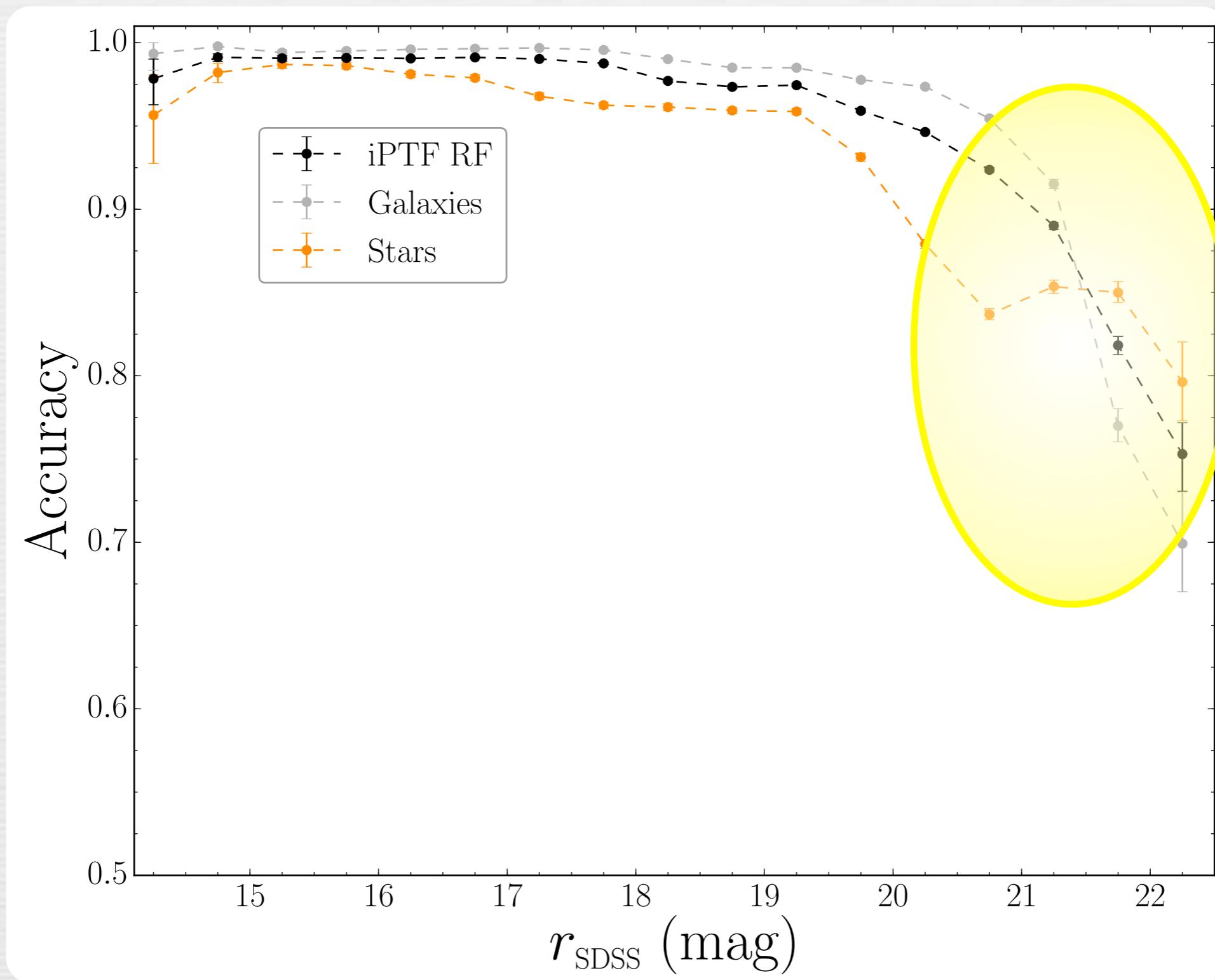
# Living Dangerously

## Star-Galaxy Separation



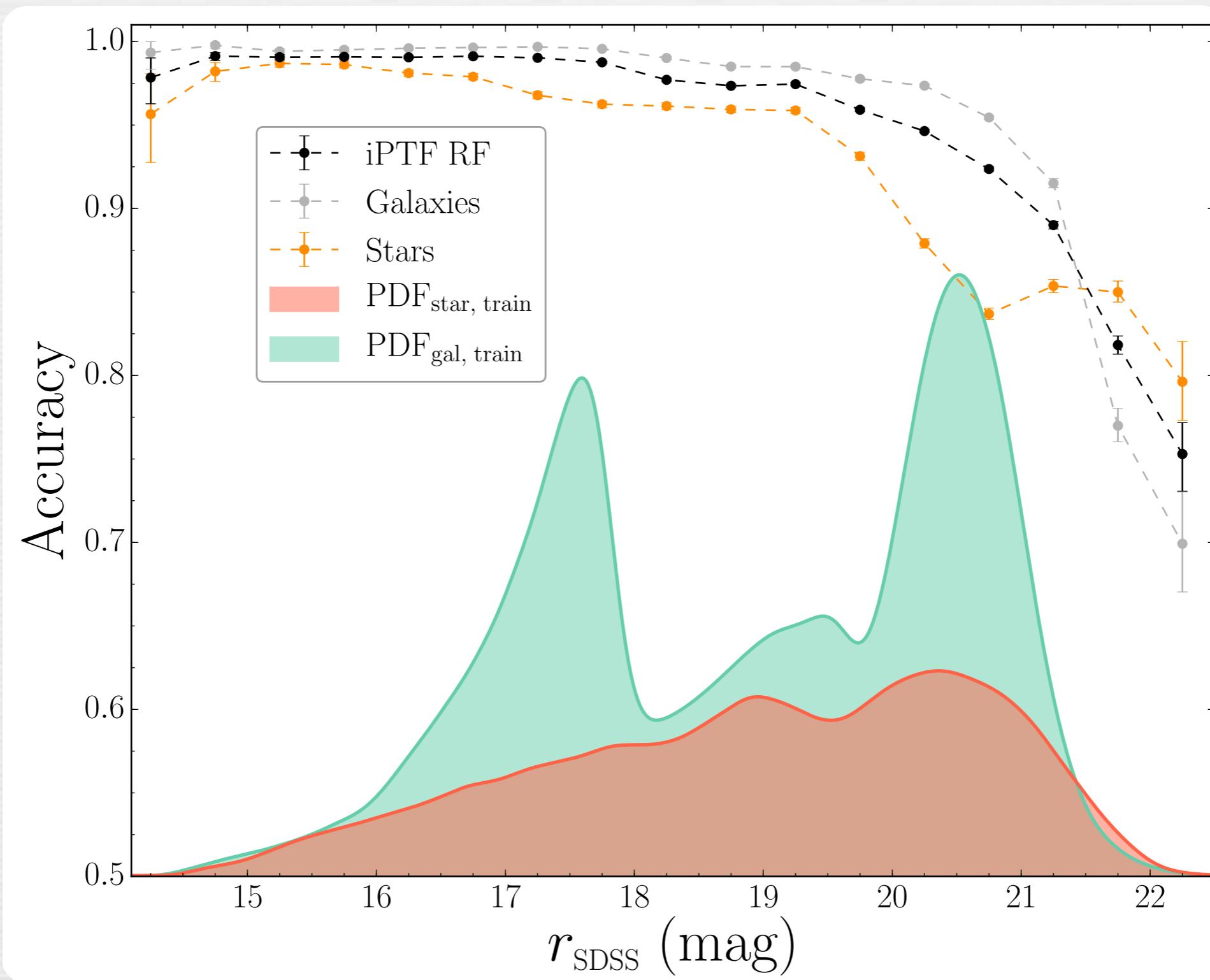
# Living Dangerously

## Star-Galaxy Separation



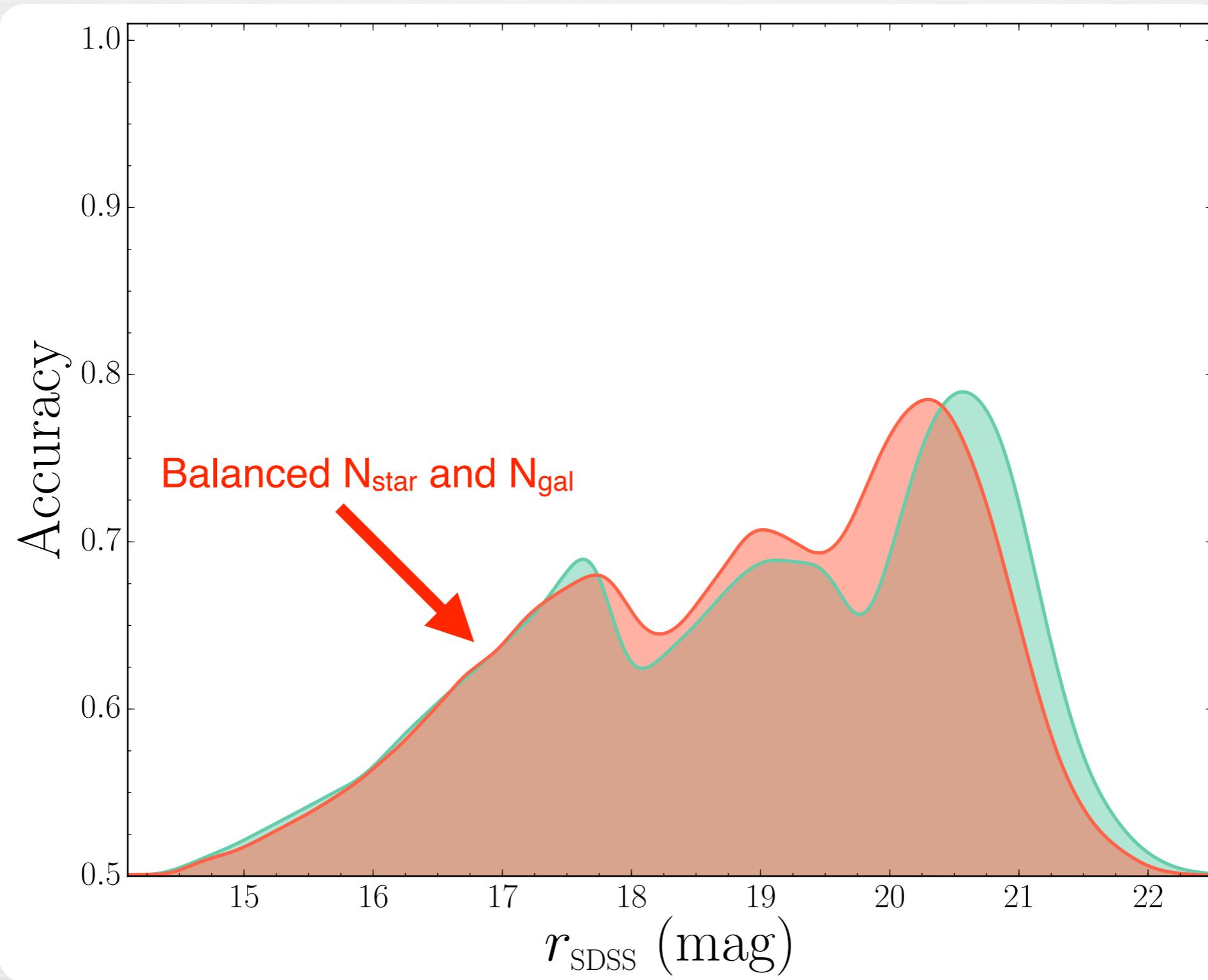
# Living Dangerously

## Star-Galaxy Separation



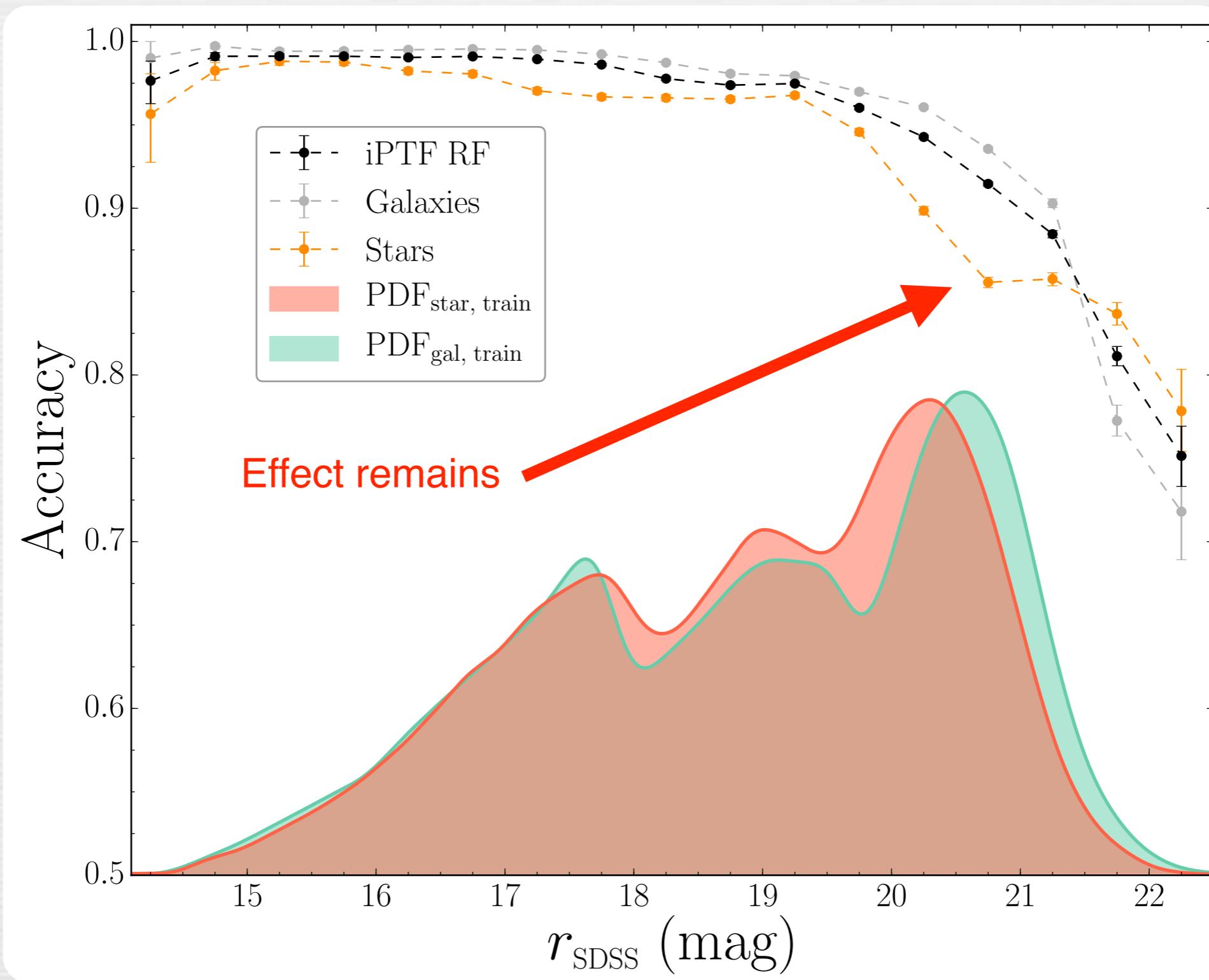
# Living Dangerously

## Star-Galaxy Separation



# Living Dangerously

## Star-Galaxy Separation



# Living Dangerously

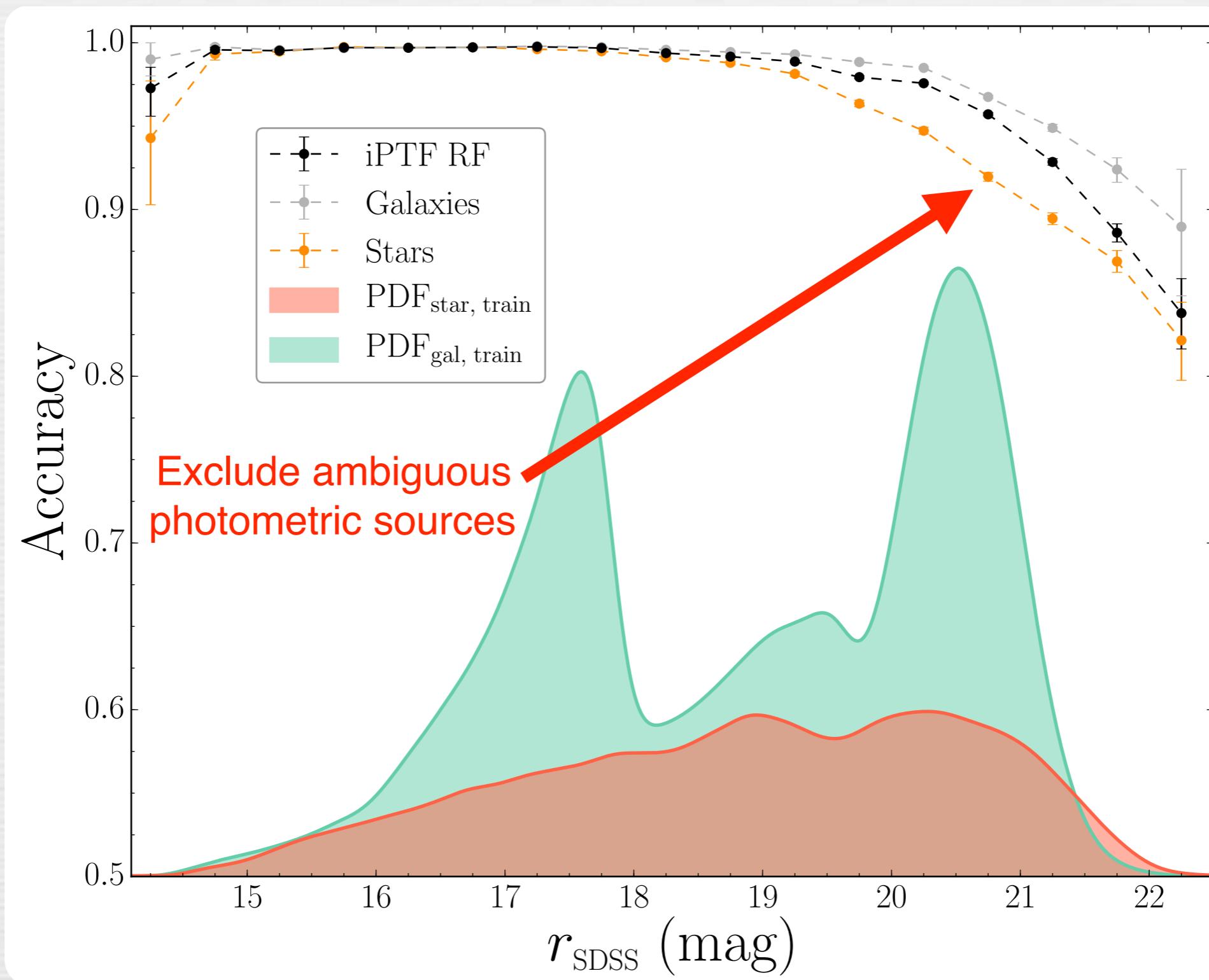
## Star-Galaxy Separation



SDSS aggressively took spectra of  
*anything* that **might be** an LRG

# Living Dangerously

## Star-Galaxy Separation



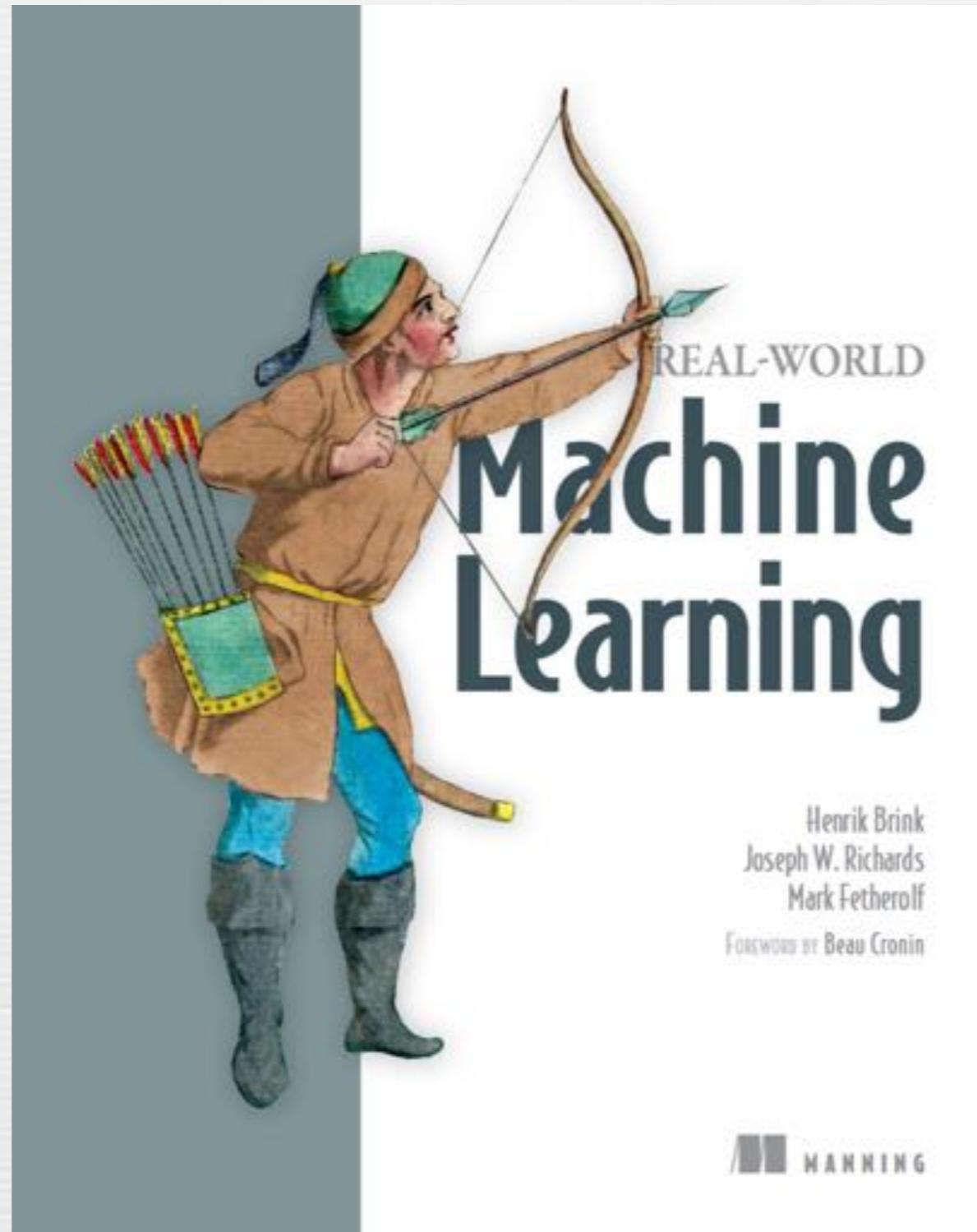
**Worry  
About  
The Data**

# Part III



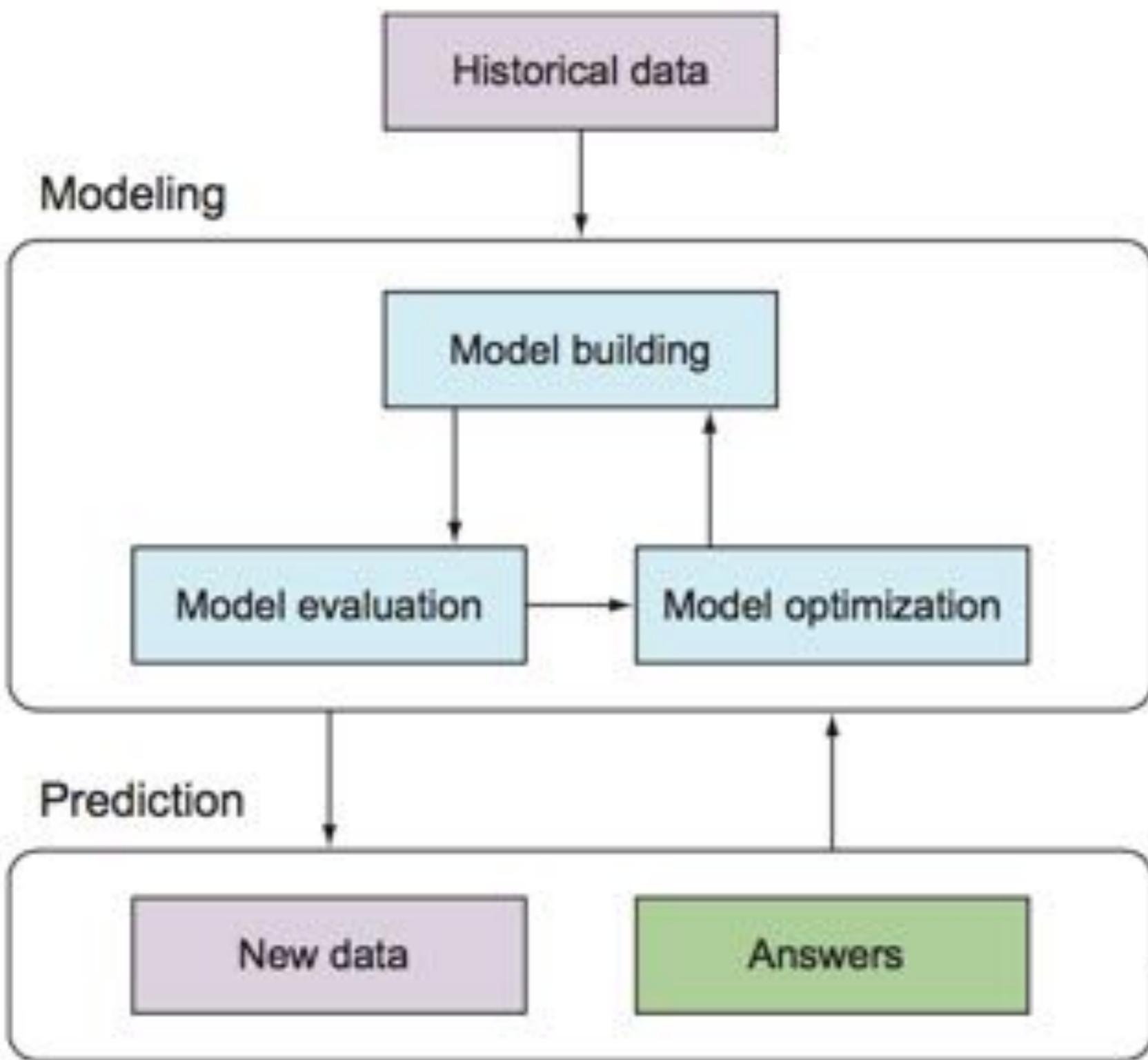
credit: Wiki commons

# Developing the Machine Learning Workflow



Brink, Richards, & Fetherolf 16

# The Machine Learning Workflow



# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

**Worry  
About  
The Data**

# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

Select features to use in the model

# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

Select features to use in the model

Determine “ground truth” or labels for the training set

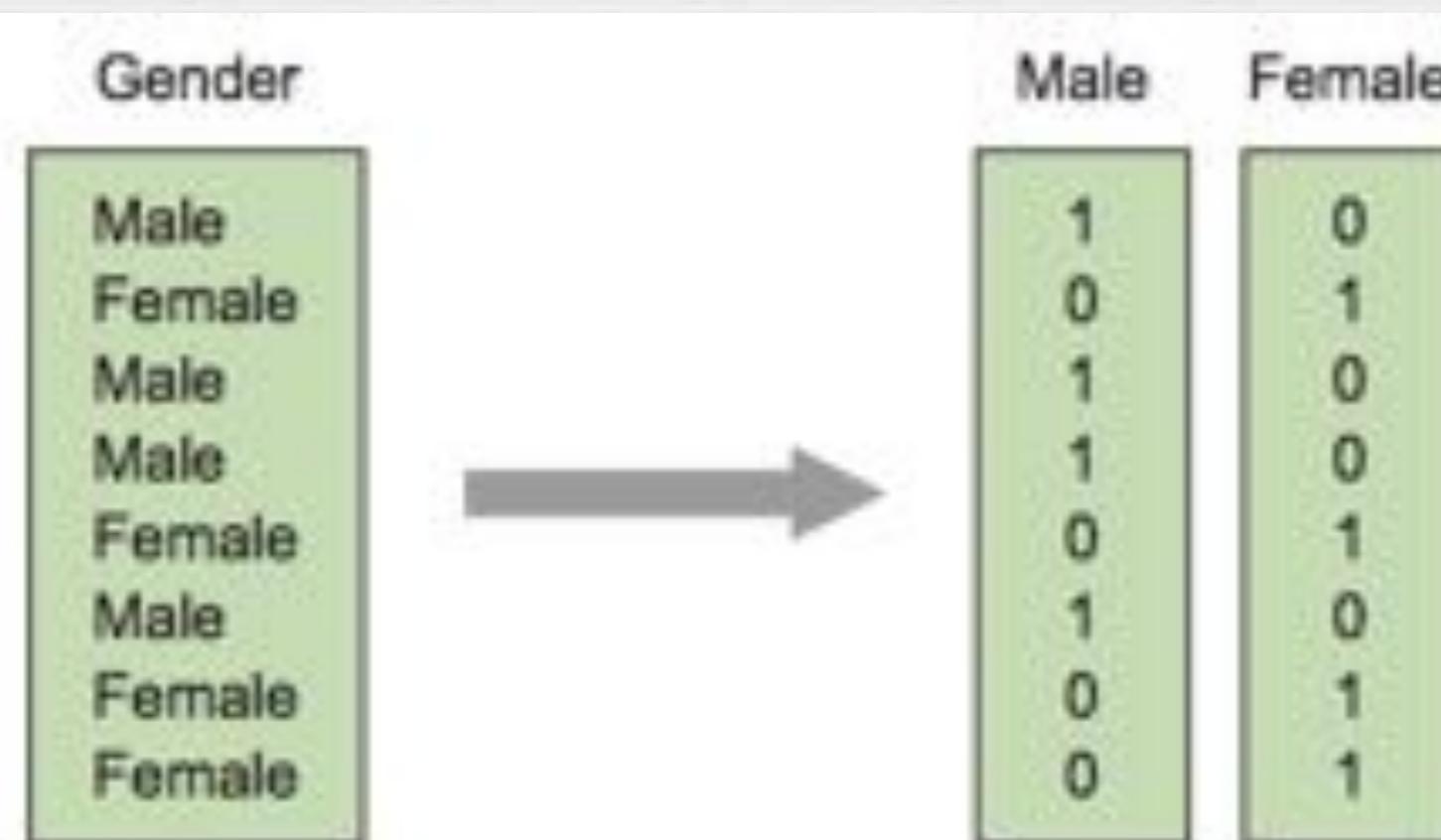
# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

Select features to use in the model

Determine “ground truth” or labels for the training set

Convert categorical features



# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

Select features to use in the model

Determine “ground truth” or labels for the training set

Convert categorical features

Impute (or throw out?) missing data

# Data Preparation

## Strategies for missing data

Does the missing data have meaning?

**Yes** - replace with numerical value (-999) or new categorical variable

**No** - **if** data set is large with few missing values:

remove objects with missing data

**else if** dataset is large and temporal:

replace missing values with preceding value or interpolate

**else if** dataset has simple distribution:

replace missing values with mean or median

**else:**

build separate ML model to impute (predict) missing values

# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

Select features to use in the model

Determine “ground truth” or labels for the training set

Convert categorical features

Impute (or throw out?) missing data

Normalize the features

# Data Preparation

Query, observe, simulate, etc. - collect data that needs to be modeled

Select features to use in the model

Determine “ground truth” or labels for the training set

Convert categorical features

Impute (or throw out?) missing data

Normalize the features

Visualize the data

**Worry  
About  
The Data**

# Feature Engineering

Add new features - if necessary

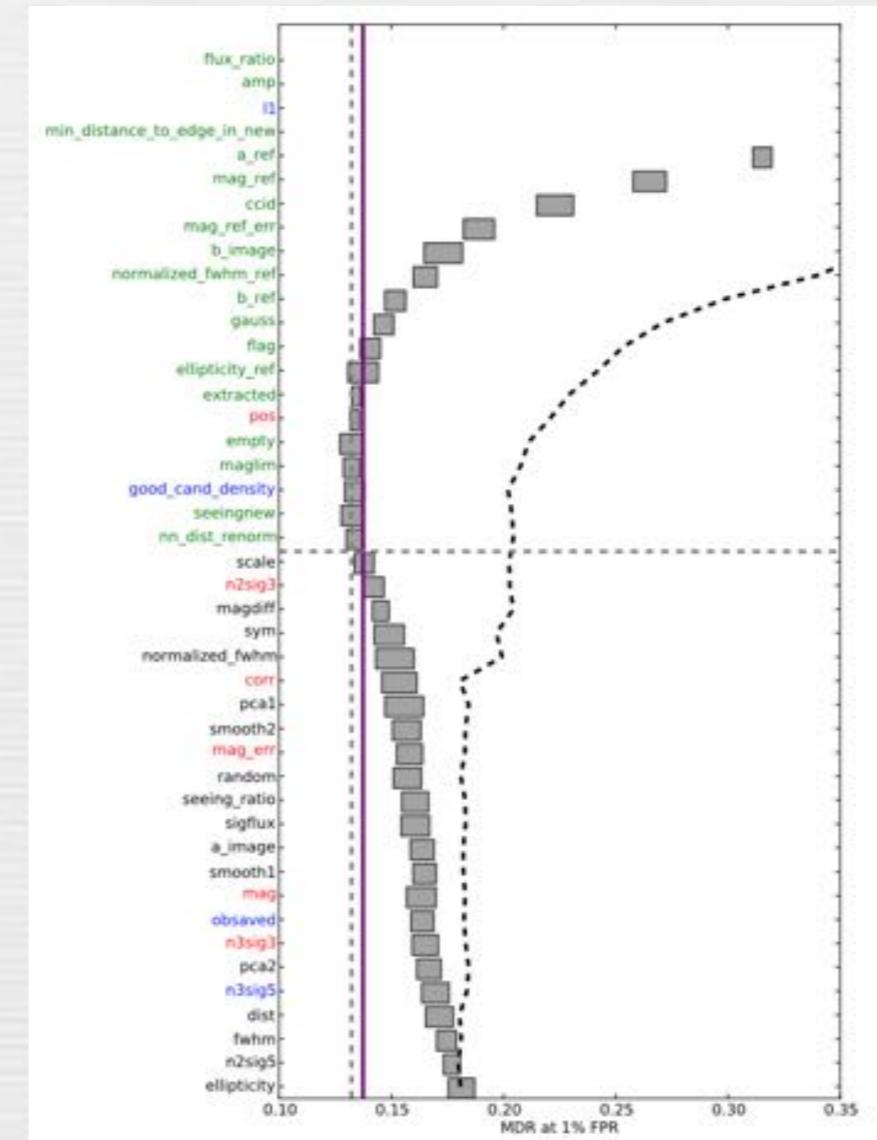
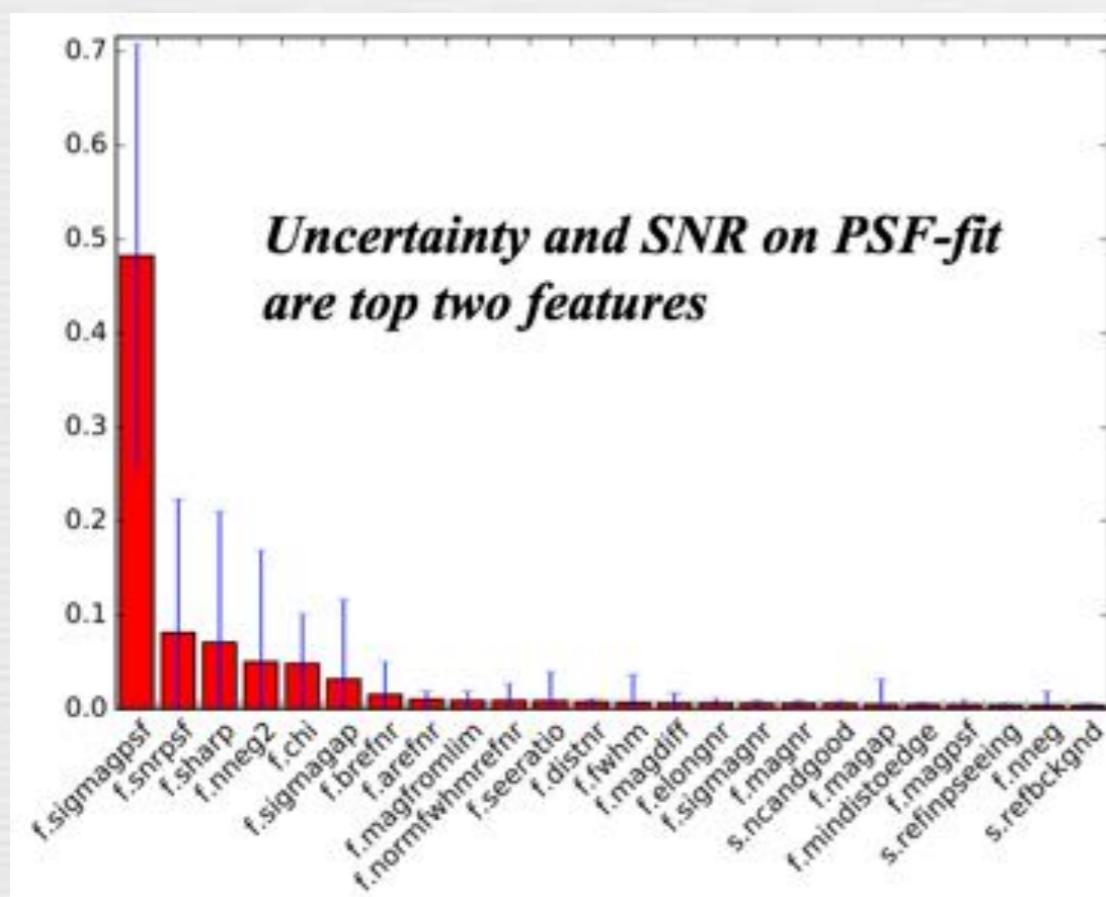
Utilize domain knowledge to create/compute new features

Combine features or represent in an alternative fashion

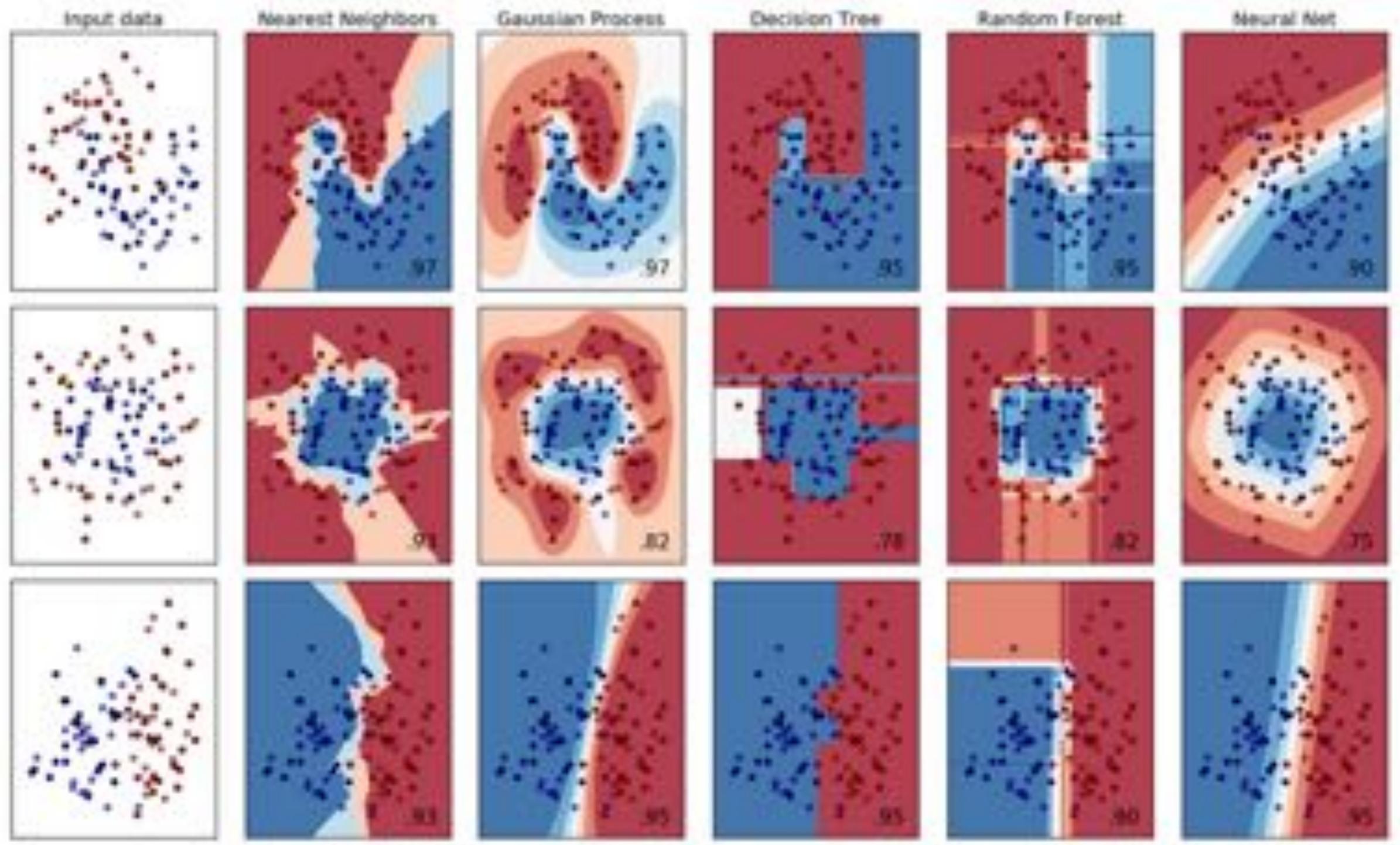
Remove noisy/uninformative features - if necessary

Determine feature importance (RF)

Forward/backward selection to iteratively remove features



# Model Selection

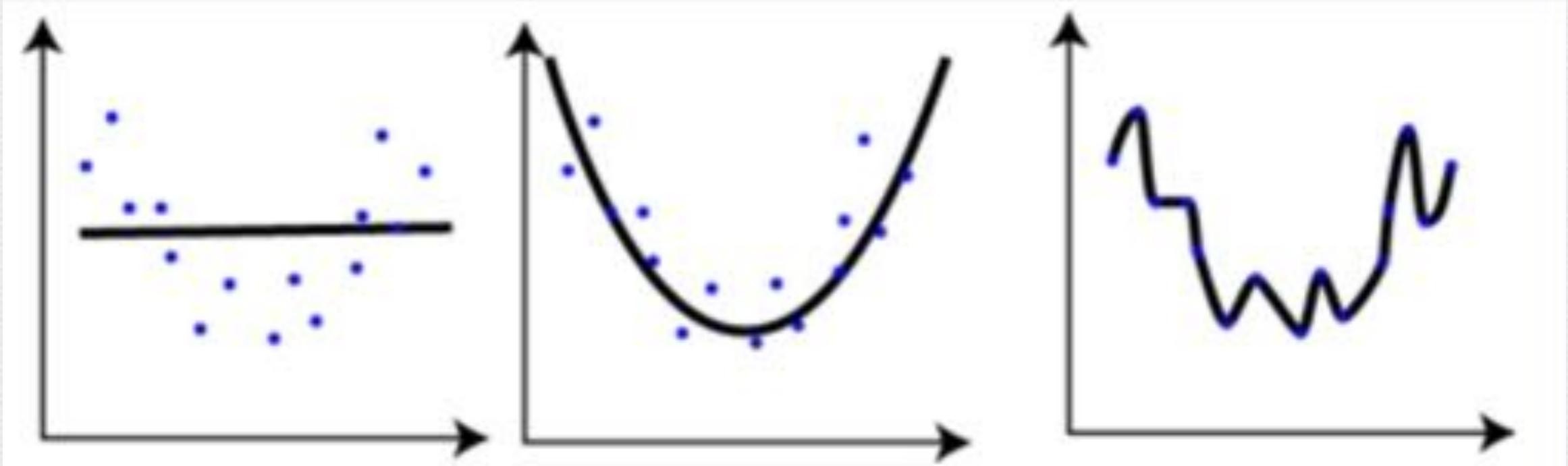


credit:scikit-learn

**Worry  
About  
The Data**

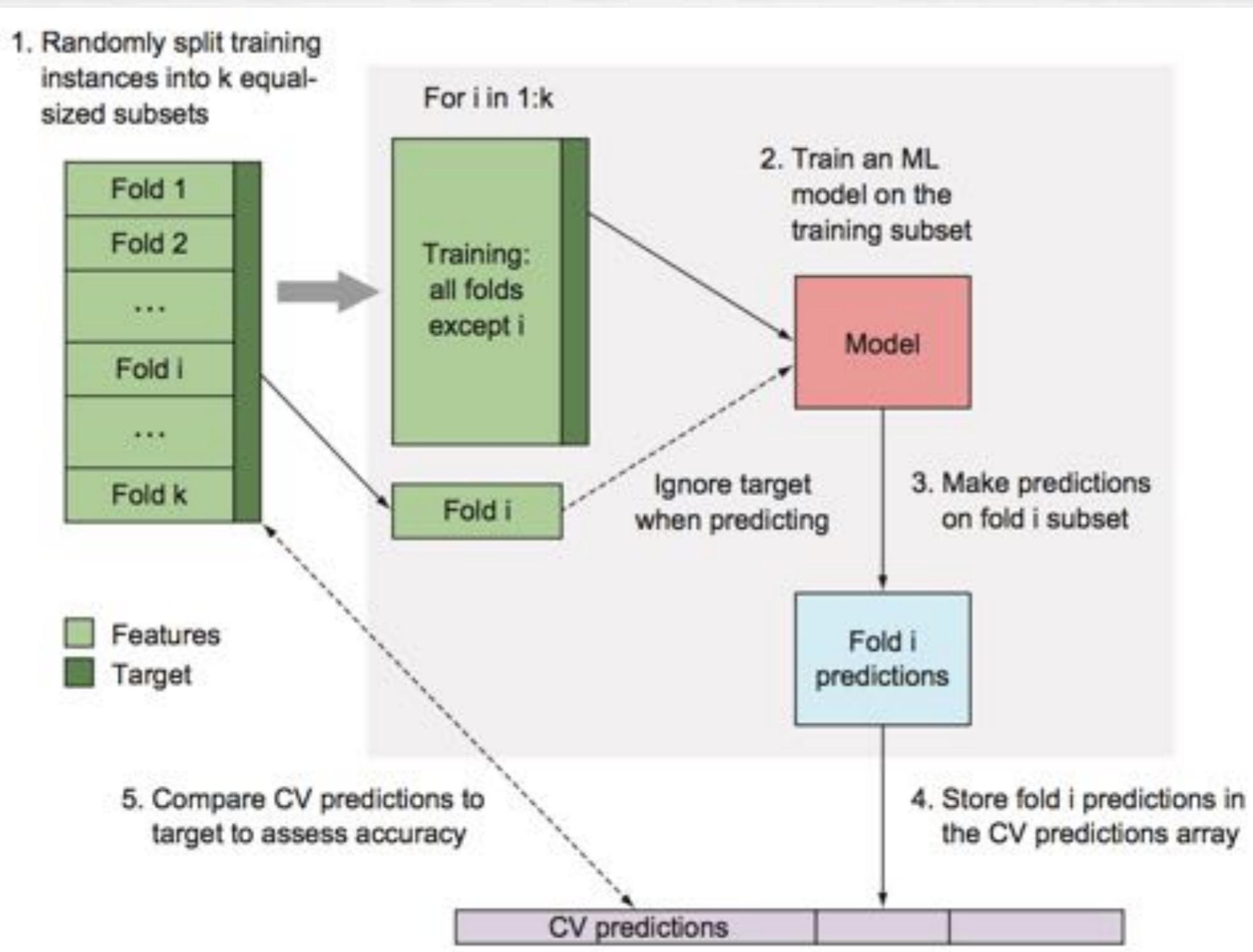
# Model Evaluation

Avoid under- and over-fitting



# Model Evaluation

## Cross Validation



# Model Evaluation

## Classification Terminology

True Positive (TP) + classified as +

False Positive (FP) – classified as +

True Negative (TN) – classified as –

False Negative (FN) + classified as –

# Model Evaluation

Confusion Matrix

Predicted Class

+

-

True Class	Predicted Class	
	+	-
+	TP	FN
-	FP	TN

# Model Evaluation

## Confusion Matrix

	True Class																													
Predicted Class	a	b1	b2	b3	b4	c	d	e	f	g	h	i	j	j1	l	o	p	q	r1	s1	s2	s3	t	u	v	w	x	y		
a. Mira	0.903	0.018				0.042										0.057	0.178													
b1. Semireg PV	0.066	0.464		0.298	0.111	0.125										0.086		0.25	0.069	0.079										
b2. SARG A			0.6	0.004	0.019																			0.059						
b3. SARG B		0.015	0.260	0.066															0.25	0.06					1					
b4. LSP	0.011	0.074	0.067	0.069	0.07											0.2	0.171						0.069							
c. RV Tauri		0.009				0.79	0.011											0.068				0.129				0.016				
d. Classical Cepheid						0.066	0.038	0.28											0.29	0.1										
e. Pop. II Cepheid						0.042		0.30	0.25											0.069										
f. Multi. Mode Cepheid						0.011	0.077	0.5	0.012								0.049		0.1	0.069		1								
g. RR Lyrae, FM						0.011	0.077		0.065		0.5		1																	
h. RR Lyrae, FO									0.213	0.067		0.038														0.016				
i. RR Lyrae, DM											0.5																			
j. Delta Scuti										0.004	0.798		0.278													0.016				
j1. SX Phe																														
i. Beta Cephei												0.197	0.799														0.016			
o. Pulsating Be		0.067														0.4				0.069										
p. RSG	0.044															0.066				0.125										
q. Chem. Peculiar										0.036		0.3	0.013													0.016				
r1. RCB																0.708														
s1. Class. T Tauri																														
s2. Weak-line T Tauri		0.004		0.011												0.049		0.7	0.118				0.061							
s3. RS CVn																		0.06	0.069			0.037								
t. Herbig AE/BE													0.2								0.305									
u. S Doradus																0.06		0.303					0.066	0.199						
v. Ellipsoidal																		0.118						0.074	0.097	0.044				
w. Beta Persei																		0.068		0.118				0.074	0.097	0.044				
x. Beta Lyrae																0.069		0.25	0.069	0.069				0.091	0.092					
y. W Ursae Maj.						0.042		0.042		0.042	0.089	0.036				0.25		0.118					0.091	0.092						

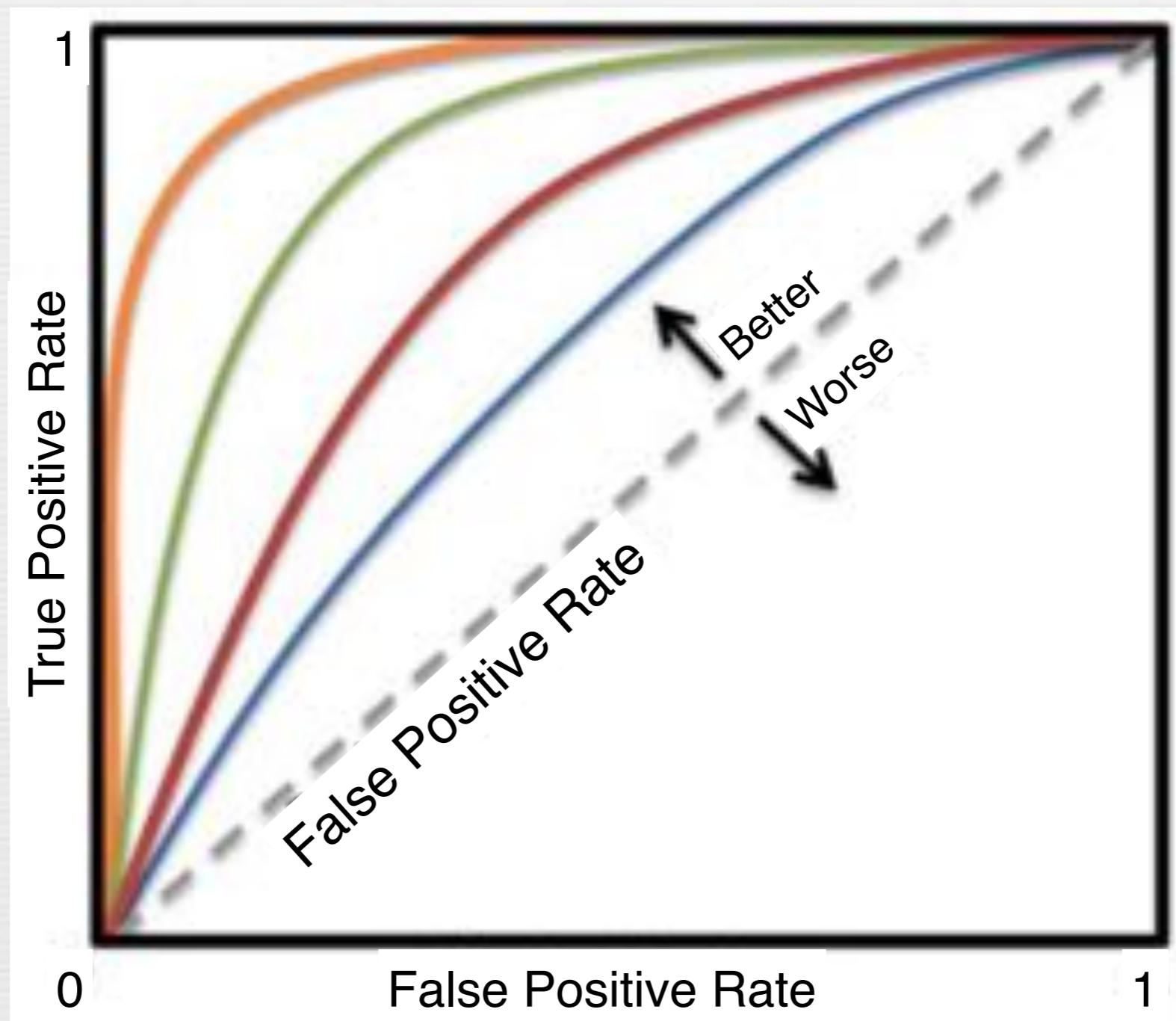
# Model Evaluation

True Positive Rate (TPR)

$TP / (TP + FN)$

False Positive Rate (FPR)

$FP / (TN + FP)$



ROC  
Curve

# Model Evaluation

Precision

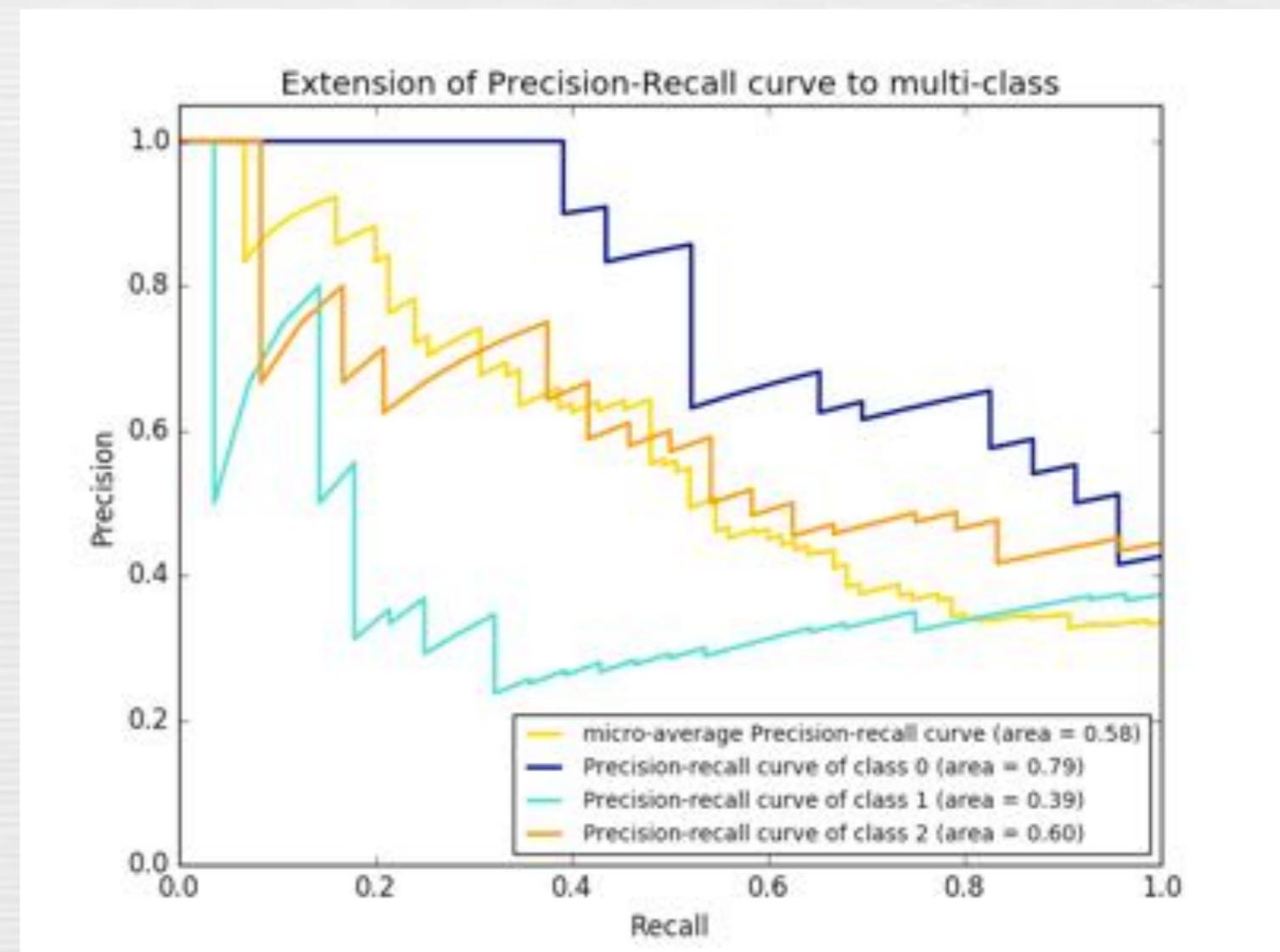
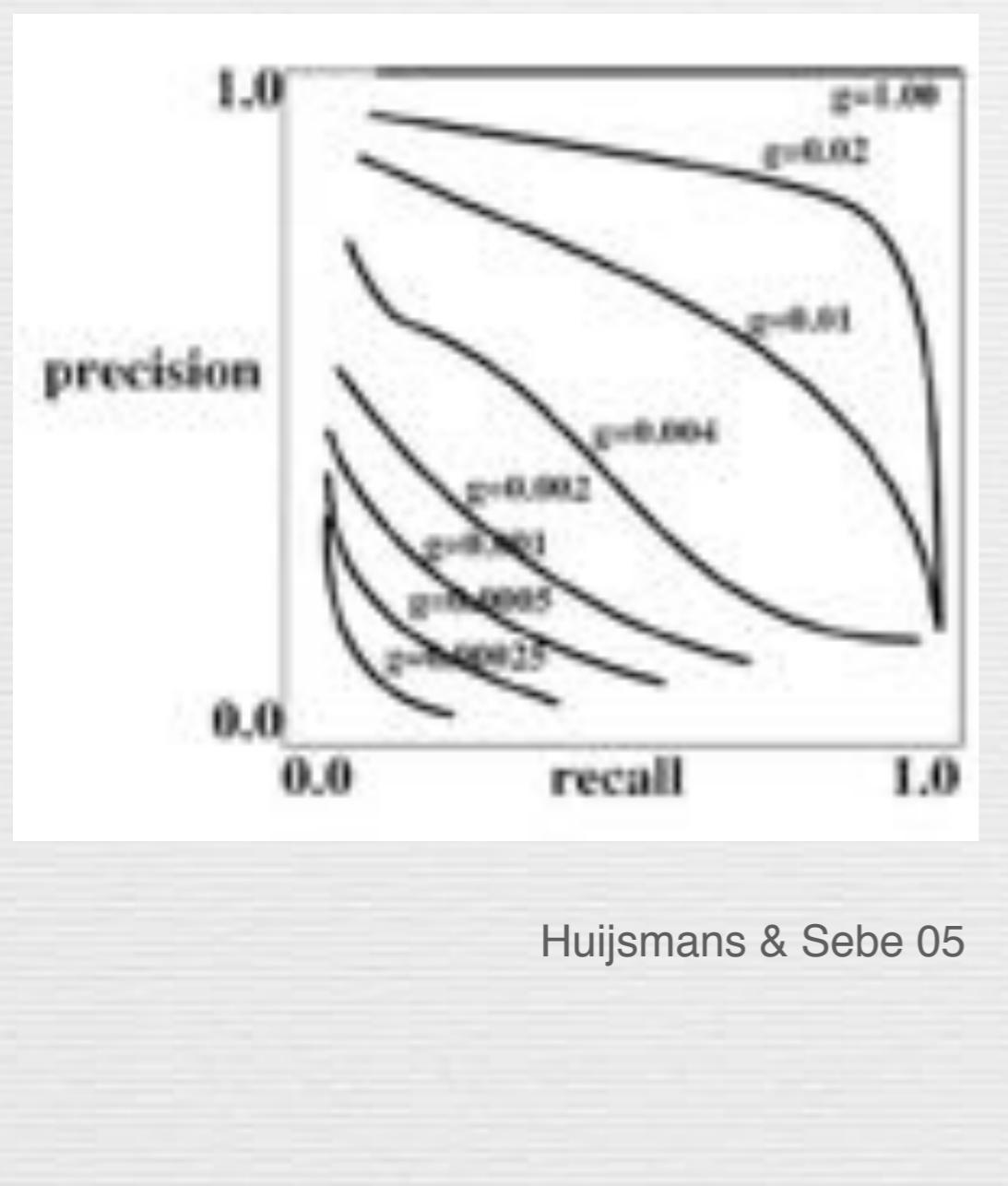
$TP / (TP + FP)$

Recall

$TP / (TP + FN)$

$F_1$

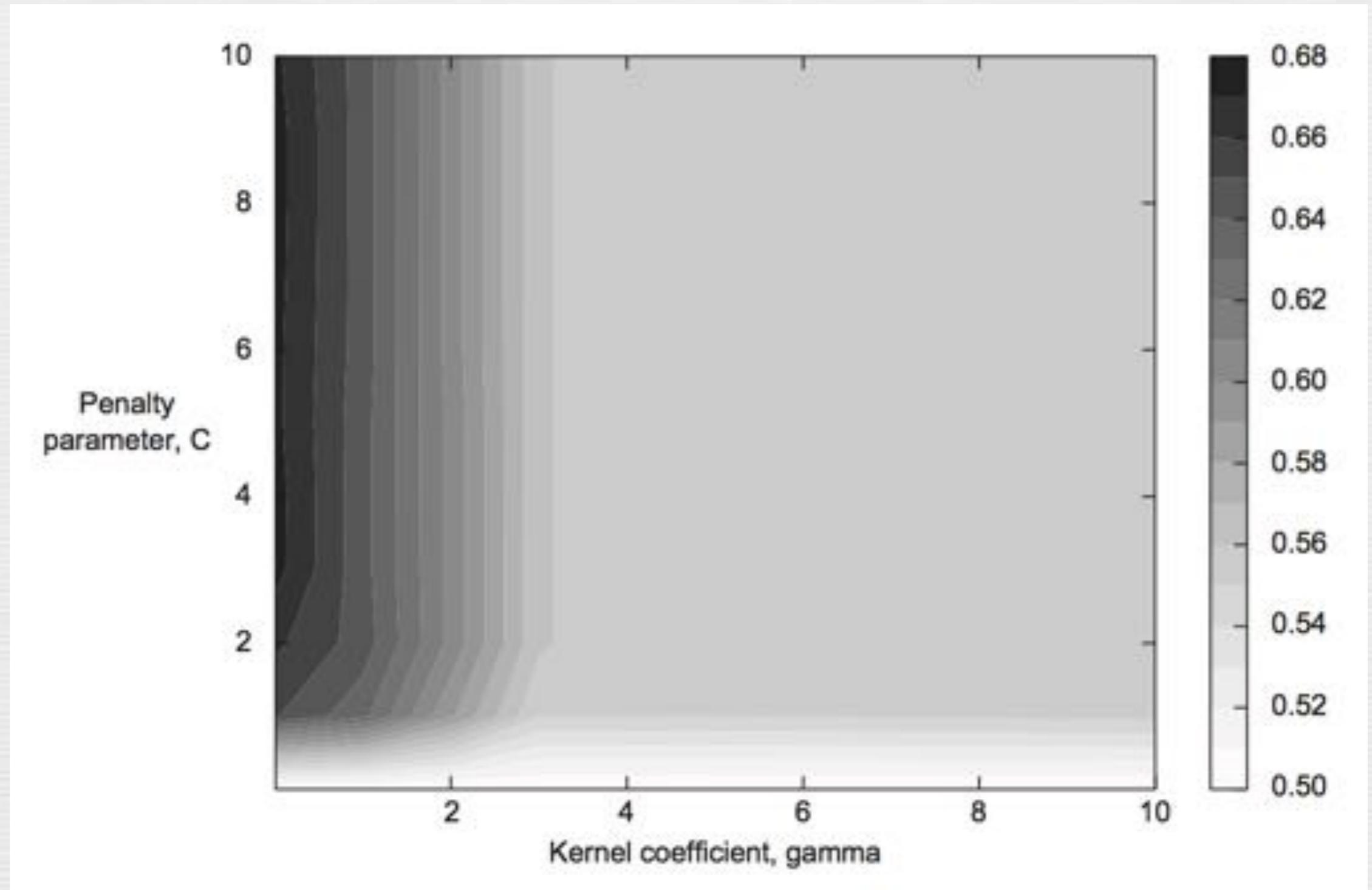
$2 * (P * R) / (P + R)$



credit: sklearn

# Model Optimization

Identify optimal tuning parameters via grid search



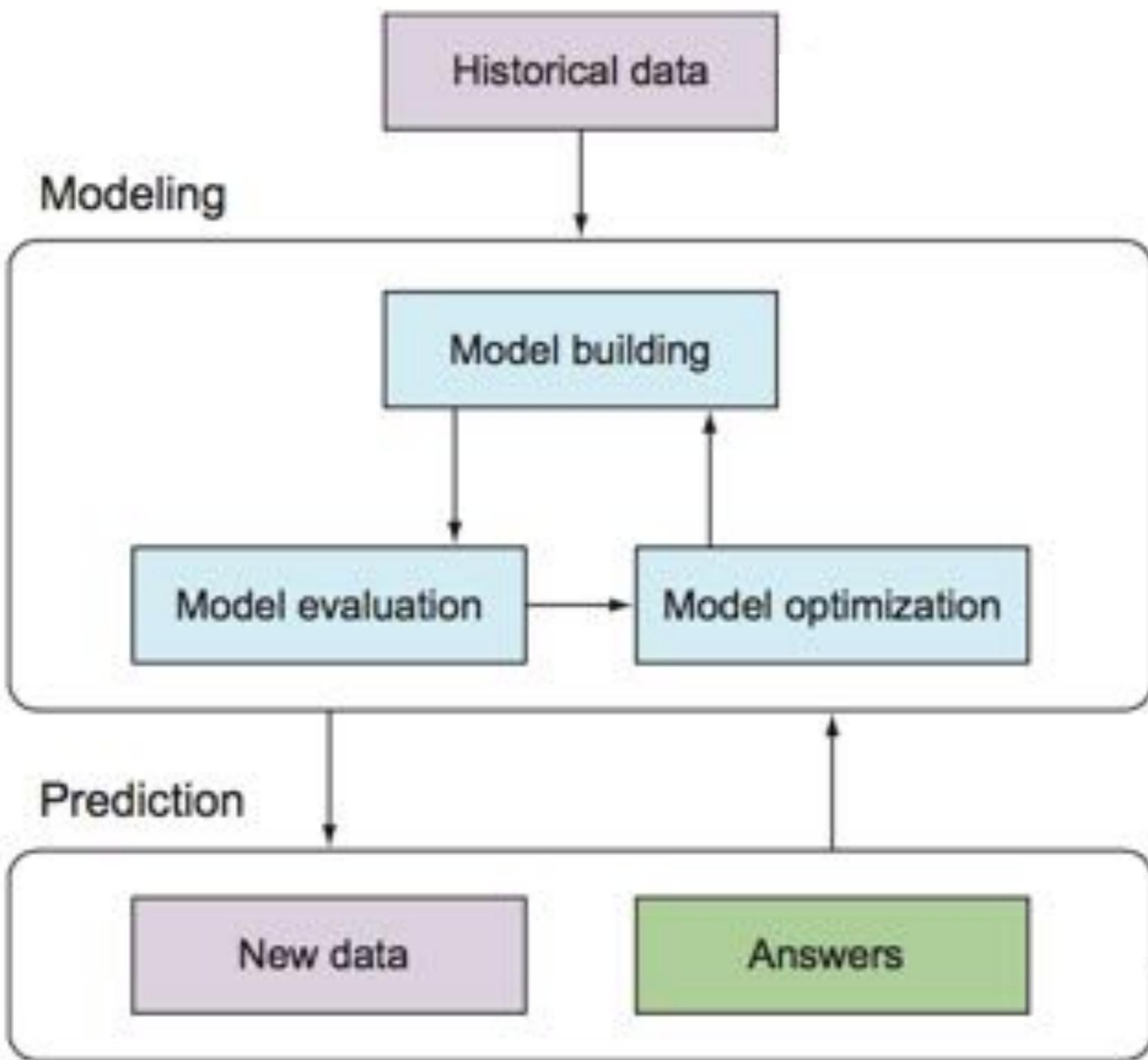
# Model Prediction



credit: MIT OCW 15.097

**Worry  
About  
the Data**

# The Machine Learning Workflow



# Conclusions

Data-driven solutions are a necessity for modern data sets  
ML is particularly useful for engineering solutions

Off-the-shelf ML algorithms are rarely plug+play  
nasty systematics (heteroskedastic errors & training bias)  
e.g., star-galaxy separation

Principles (sometimes algorithms) of ML are very useful  
when data leads theory, allow data to drive the models  
test the utility of everything with independent observations  
make informed thresholding decisions

Worry About the Data