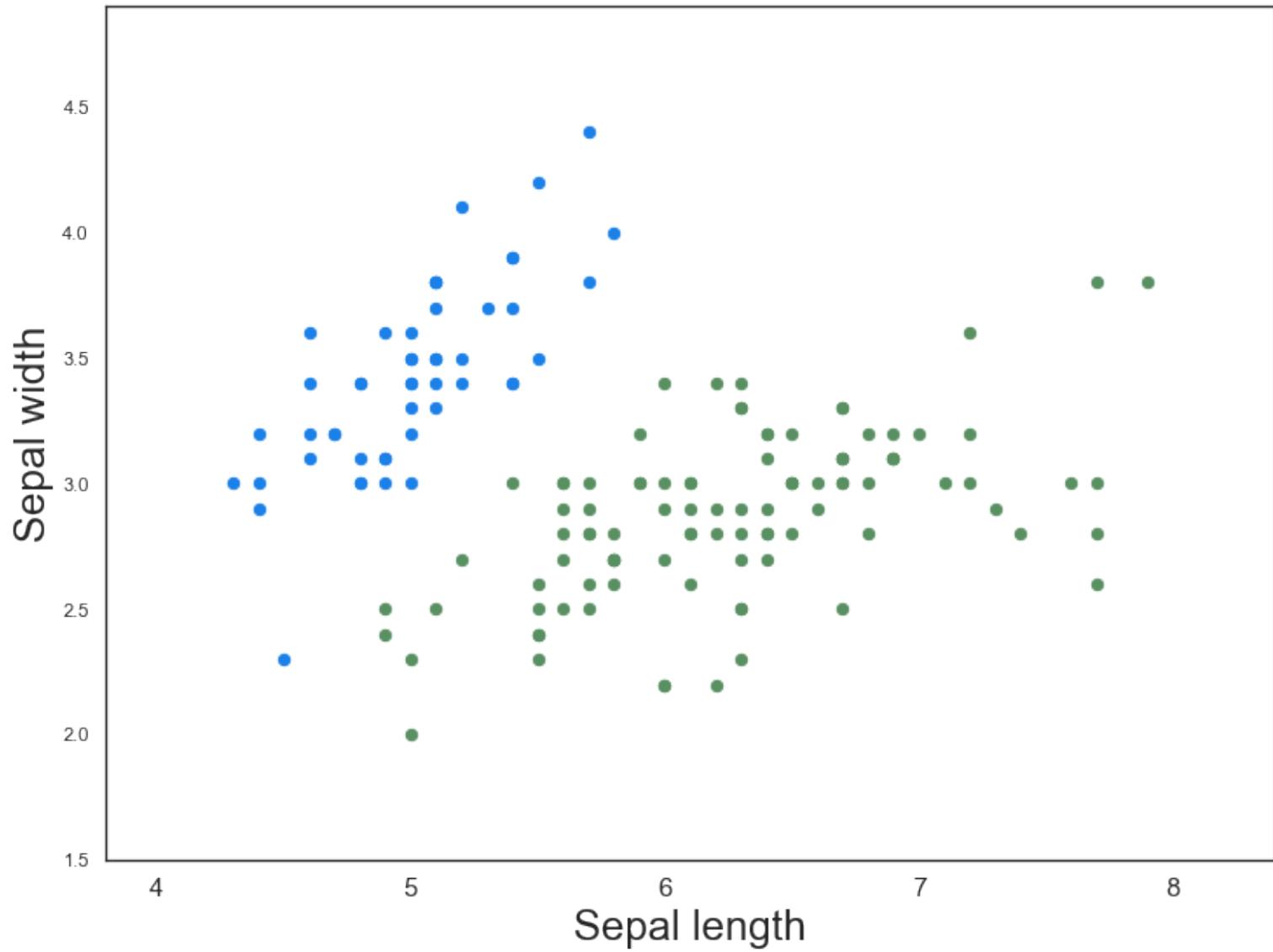


Classification

Data Mining and Text Mining





Suppose you receive an unlabeled data point

How would you decide, its class?

Simplest way would be the majority class

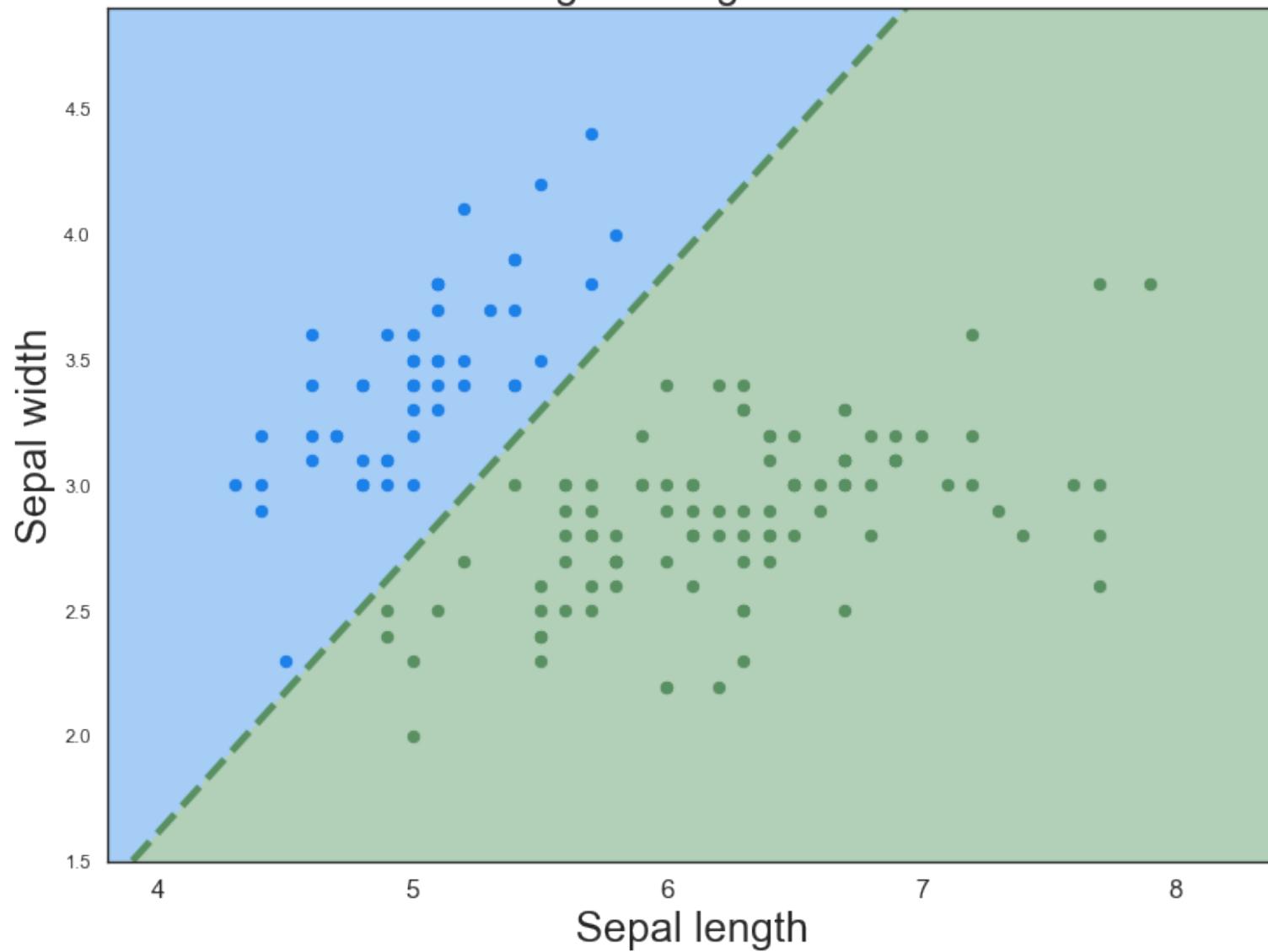
In the example, the number of green data points is larger than the blue ones ...

So I might decide to assign a green label to unlabeled data points

This gives us a baseline performance

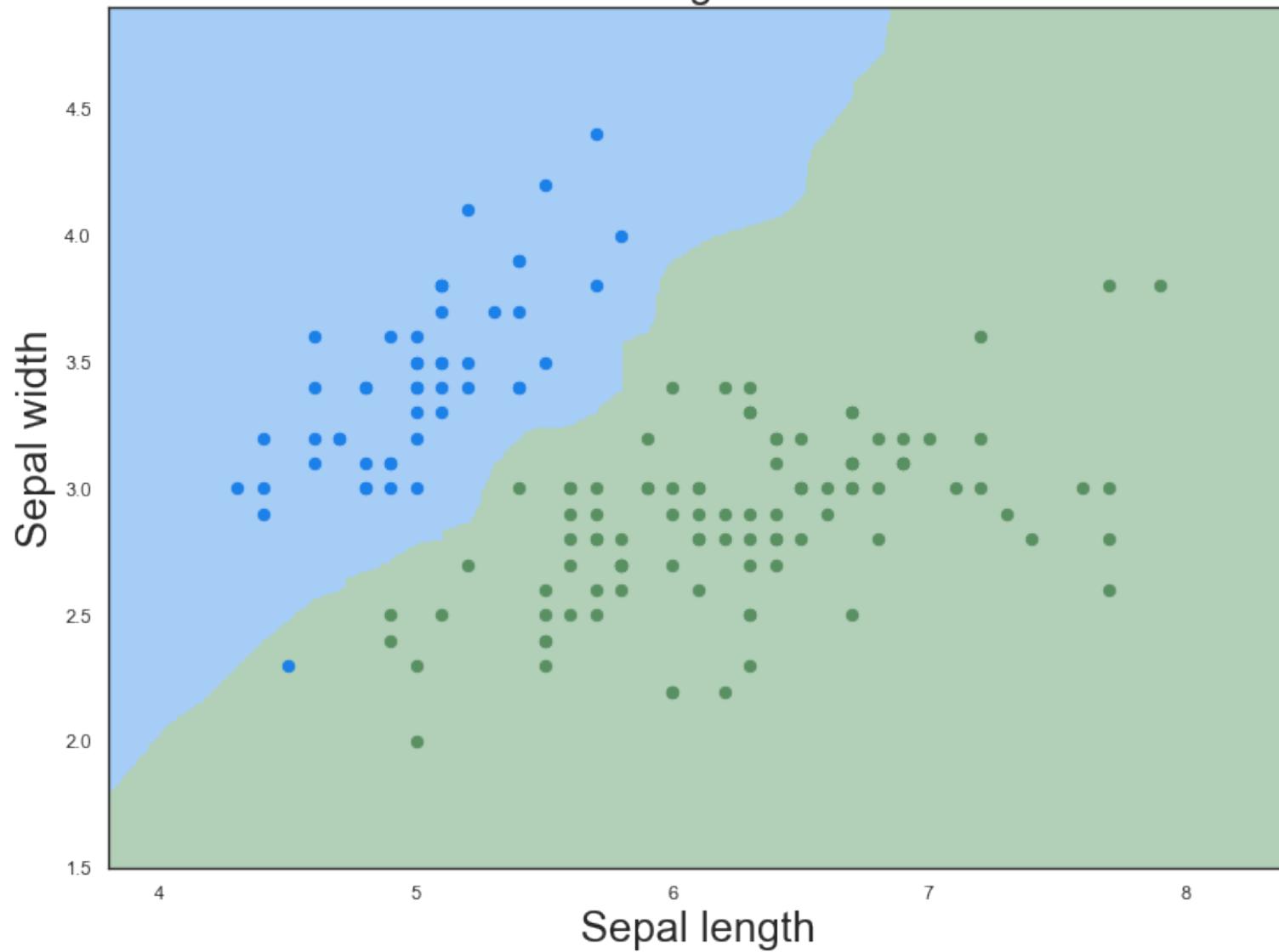
But we can do better than this ...

Logistic Regression



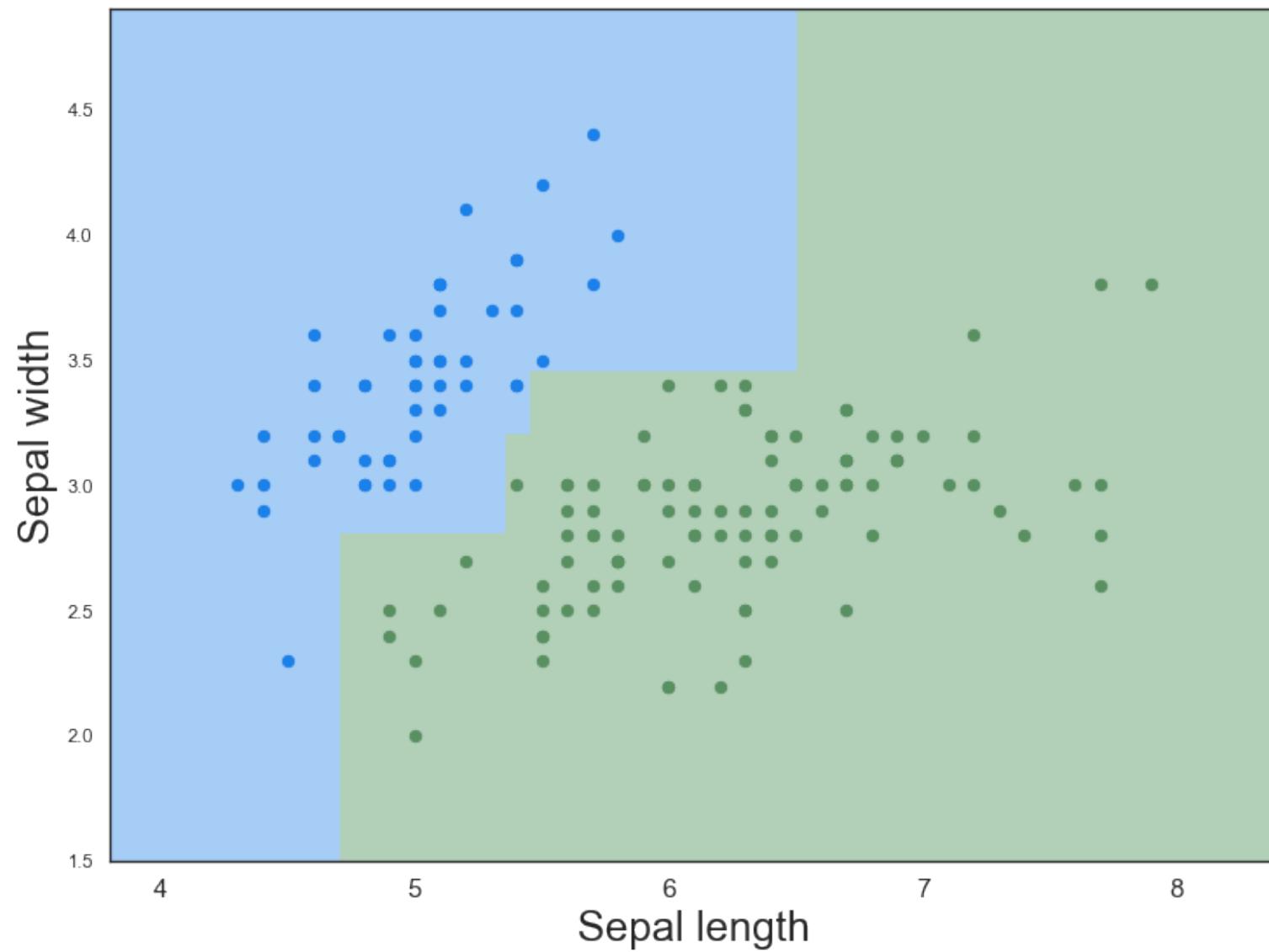
Logistic Regression Model

k-Nearest Neighbor with k=5



K-Nearest Neighbor model with a k of 5

Decision Tree Model



Decision Tree Model

As for regression, we must compute
a model using known data

And the model should perform
well on unknown data.

Contact Lenses Data

10

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	No	Reduced	None
Young	Hypermetrope	No	Normal	Soft
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	No	Reduced	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	No	Reduced	None
Pre-presbyopic	Hypermetrope	No	Normal	Soft
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	No	Reduced	None
Presbyopic	Hypermetrope	No	Normal	Soft
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

In classification problems ...

Rows are typically called
examples or data points

Columns are typically called
attributes, variables or features

The target variable to be
predicted is usually called “class”

```
If tear production rate = reduced then recommendation = none  
If age = young and astigmatic = no  
    and tear production rate = normal then recommendation = soft  
If age = pre-presbyopic and astigmatic = no  
    and tear production rate = normal then recommendation = soft  
If age = presbyopic and spectacle prescription = myope  
    and astigmatic = no then recommendation = none  
If spectacle prescription = hypermetrope and astigmatic = no  
    and tear production rate = normal then recommendation = soft  
If spectacle prescription = myope and astigmatic = yes  
    and tear production rate = normal then recommendation = hard  
If age young and astigmatic = yes  
    and tear production rate = normal then recommendation = hard  
If age = pre-presbyopic  
    and spectacle prescription = hypermetrope  
    and astigmatic = yes then recommendation = none  
If age = presbyopic and spectacle prescription = hypermetrope  
    and astigmatic = yes then recommendation = none
```

Classification

The target to predict is a label (the class variable)
(good/bad, none/soft/hard, etc.)

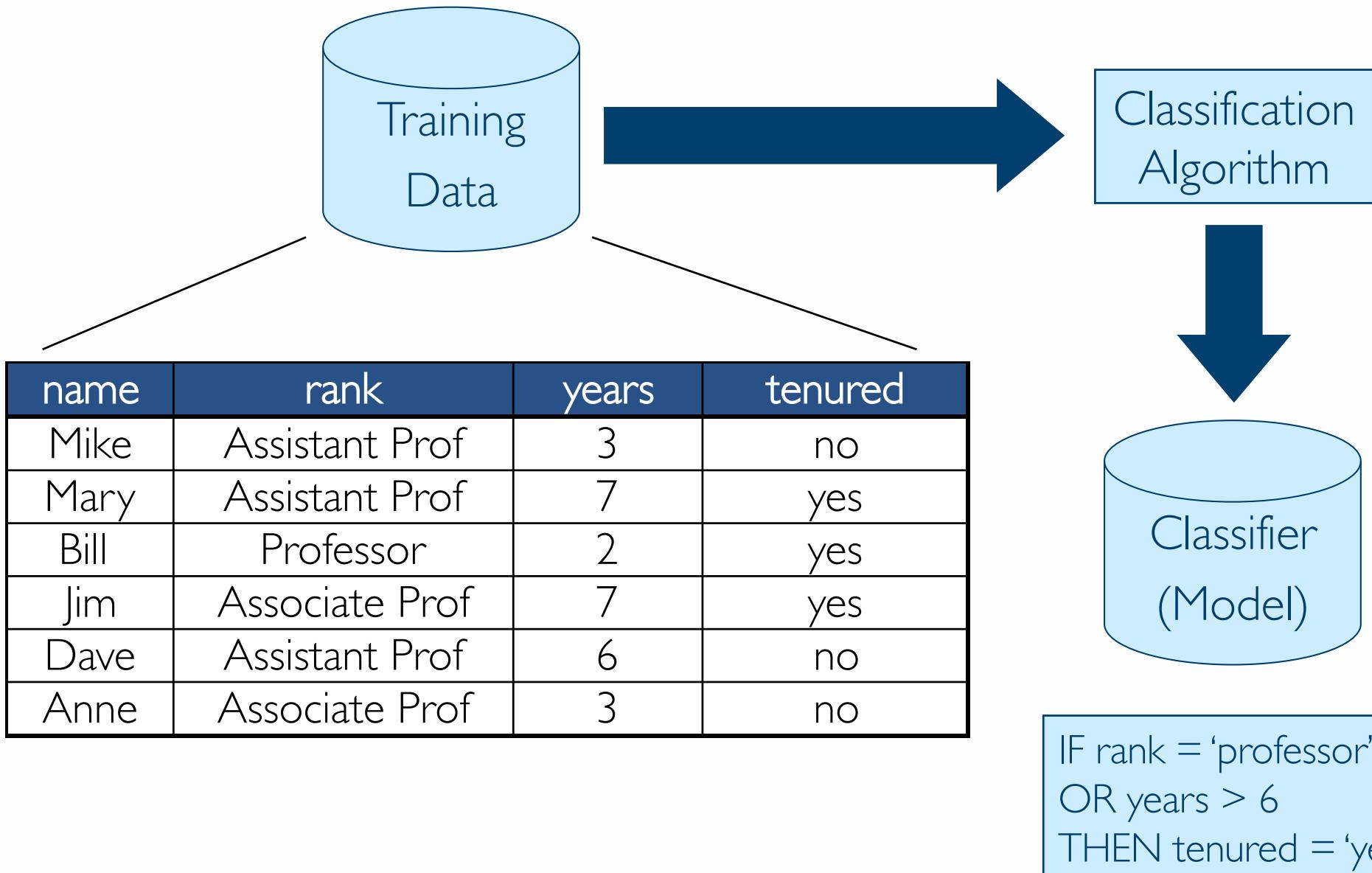
Prediction/Regression

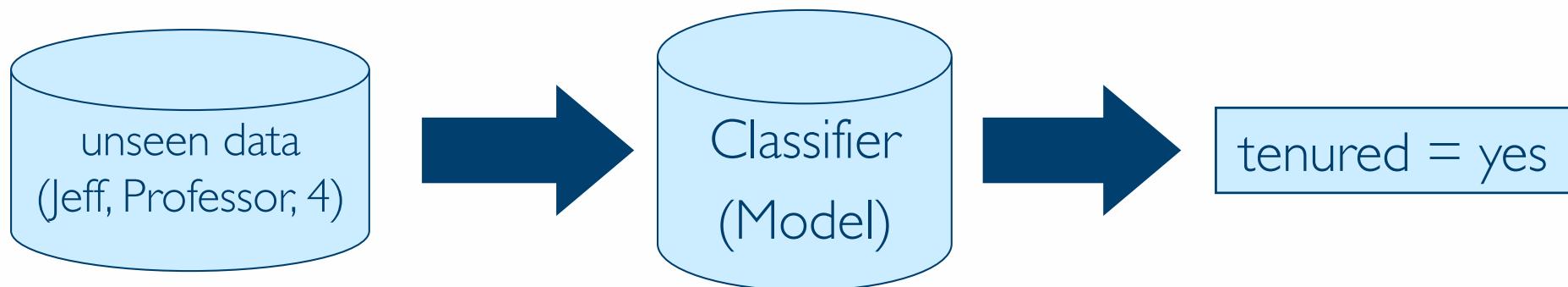
The target to predict is numerical

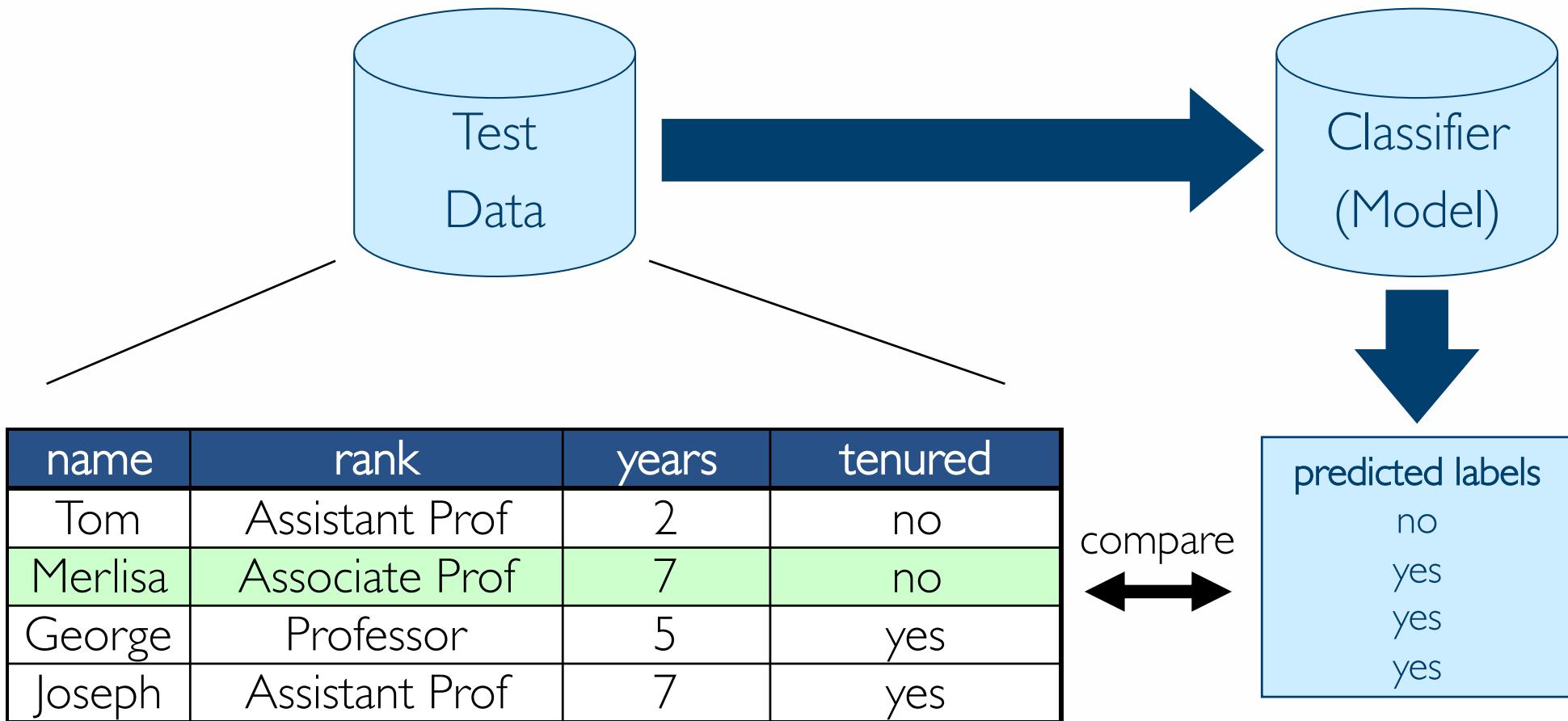
classification = model building + model usage

Classification: Model Construction

15



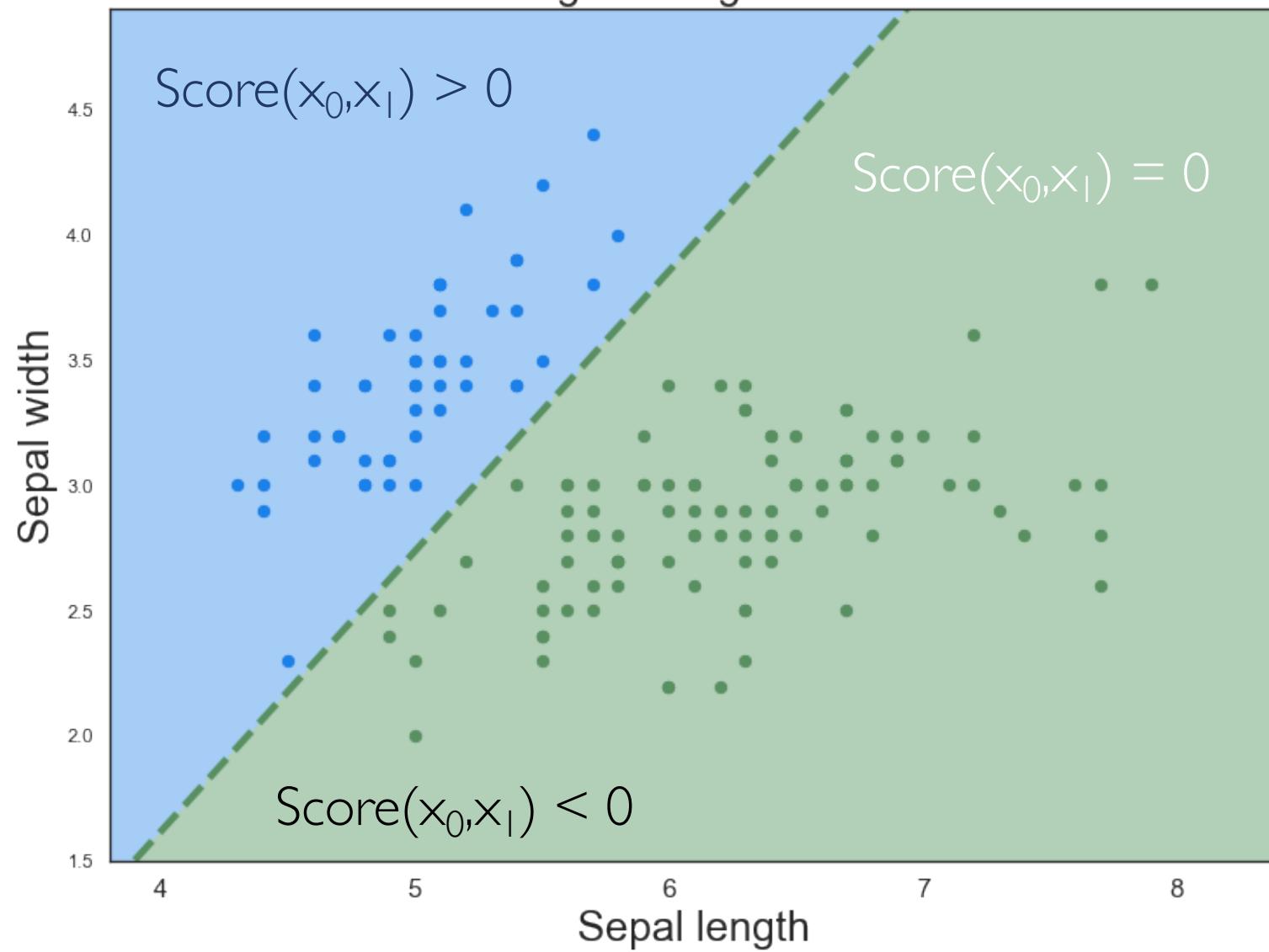




- **Accuracy**
 - Classifier accuracy in predicting the correct the class labels
- **Speed**
 - Time to construct the model (training time)
 - Time to use the model to label unseen data
- **Other Criteria**
 - Robustness in handling noise
 - Scalability
 - Interpretability
 - ...

Logistic Regression

Logistic Regression



Logistic Regression Model

To predict the labels we check whether
 $\text{Score}(x_0, x_1)$ is positive or negative

Then, we label the examples accordingly

- We define a score function similar to the one used in linear regression,

$$\text{Score}(\vec{x}_i) = \sum_{j=0}^D w_j h_j(\vec{x}_i)$$

- The label is determined by the sign of the score value,

$$\hat{y}_i = \text{sign}(\text{Score}(\vec{x}_i))$$
$$= \begin{cases} +1 & \text{if } \text{Score}(\vec{x}_i) \geq 0 \\ -1 & \text{if } \text{Score}(\vec{x}_i) < 0 \end{cases}$$

- Well-known and widely used statistical classification method
- Instead of computing the label, it computes the probability of assigning a class to an example, that is,

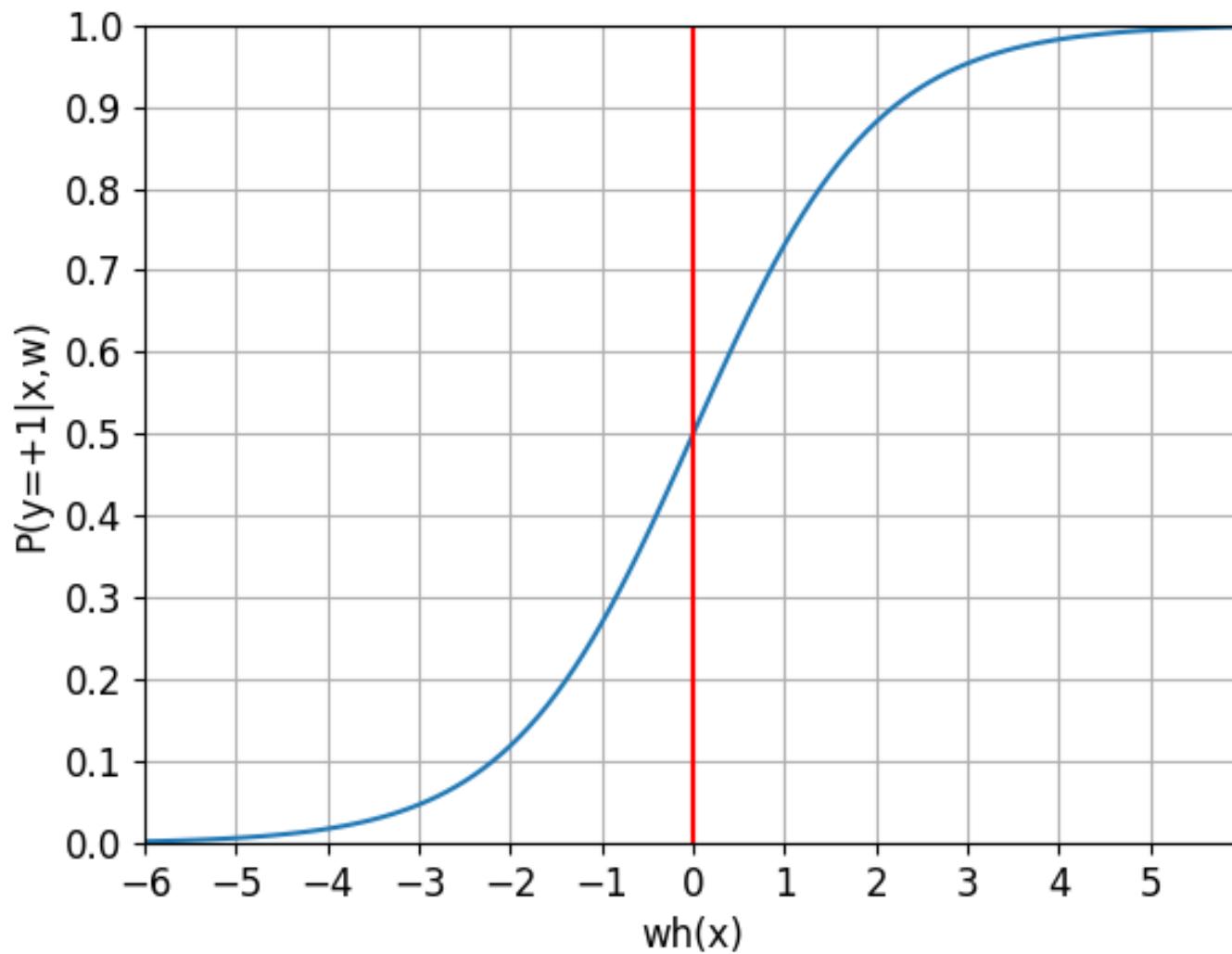
$$P(y_i | \vec{x}_i)$$

- For this purpose, logistic regression assumes that,

$$P(\hat{y}_i = +1 | \vec{x}_i) = \frac{1}{1 + e^{-Score(\vec{x}_i)}}$$

- By making the score computation explicit and using h to identify all the feature transformation h_j we get,

$$P(\hat{y}_i = +1 | \vec{x}_i, \vec{w}) = \frac{1}{1 + e^{-\vec{w}h(\vec{x}_i)}}$$



- Logistic Regression search for the weight vector that corresponds to the highest likelihood

$$\ell(\vec{w}) = \prod_{i=0}^N P(y_i | \vec{x}_i, \vec{w})$$

- For this purpose, it performs a gradient ascent on the log likelihood function

$$\ell\ell(\vec{w}) = \ln \ell(\vec{w})$$

- Which updates weight j using,

$$\frac{\partial \ell\ell}{\partial w_j} = \sum_{I=1}^N h_j(\vec{x}_i)(1[y_i = +1] - P(y = +1 | \vec{x}_i, \vec{w}))$$

- To classify an example x :

- select the class with the highest probability $P(Y=y|x)$
 - or check if ratio of probabilities is greater than one:

$$\frac{P(y = +1)|\vec{x})}{P(y = -1)|\vec{x})} = e^{\vec{w}h(\vec{x})} > 1$$

- choosing the positive class if it is, otherwise, the negative class
- Or equivalently check if natural log of ratio is greater than zero

$$\ln\left(\frac{P(y = +1)|\vec{x})}{P(y = -1)|\vec{x})}\right) = \vec{w}h(\vec{x}) = \sum_{j=0}^D w_j h_j(\vec{x}) > 0$$

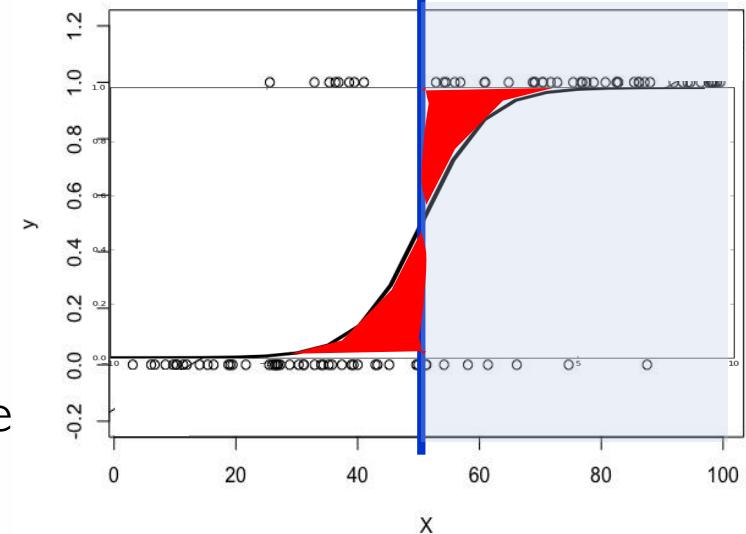
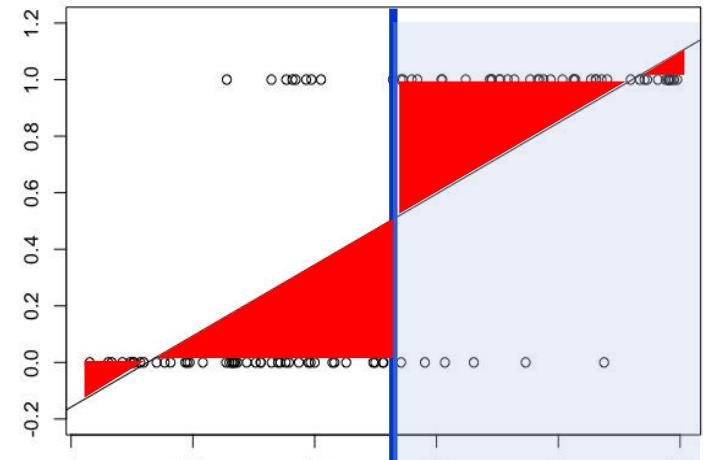
- Note that we're back to a linear classification rule meaning that logistic regression is a linear classifier

- By taking the natural logarithm of both sides, we obtain a linear classification rule that assign the label $Y=-1$ if

$$\sum_{i=0}^D w_i h_i(x) < 0$$

- $Y=+1$ otherwise

- Logistic curve fits better 0/1 data and results in better decision boundary
- Error is monotonic in distance from boundary!
 - Red areas indicative of size of error for correctly classified examples
 - Incorrectly classified instances have larger error values
- Note that decision boundary for linear regression is highly sensitive to outliers lying far from boundary
- Logistic regression not sensitive to these value



Overfitting & Regularization

- Logistic regression can use L_1 and L_2 regularization to limit overfitting and produce sparse solution
- Like it happened with regression, overfitting is often associated to large weights so to limit overfitting we penalize large weights

Overall Objective = Measure of Fit - Magnitude of Coefficients

- Regularization
 - L_1 uses the sum of absolute values, which penalizes large weights and at the same time promotes sparse solutions
 - L_2 uses the sum of squares, which penalizes large weights

L₁ Regularization

$$\ell(\vec{w}) - \alpha \|\vec{w}\|_1$$

L₂ Regularization

$$\ell(\vec{w}) - \alpha \|\vec{w}\|_2^2$$

If α is zero, then we have no regularization

If α tends to infinity,
the solution is a zero weight vector

α must balance fit and the weight magnitude

Example

From now on we are going to use k-fold cross-validation a lot to score models

k-fold cross-validation generates
k separate models

Which one should be deployed?

K-fold cross-validation provides an evaluation of model performance on unknown data

Its output it's the evaluation, not the model!

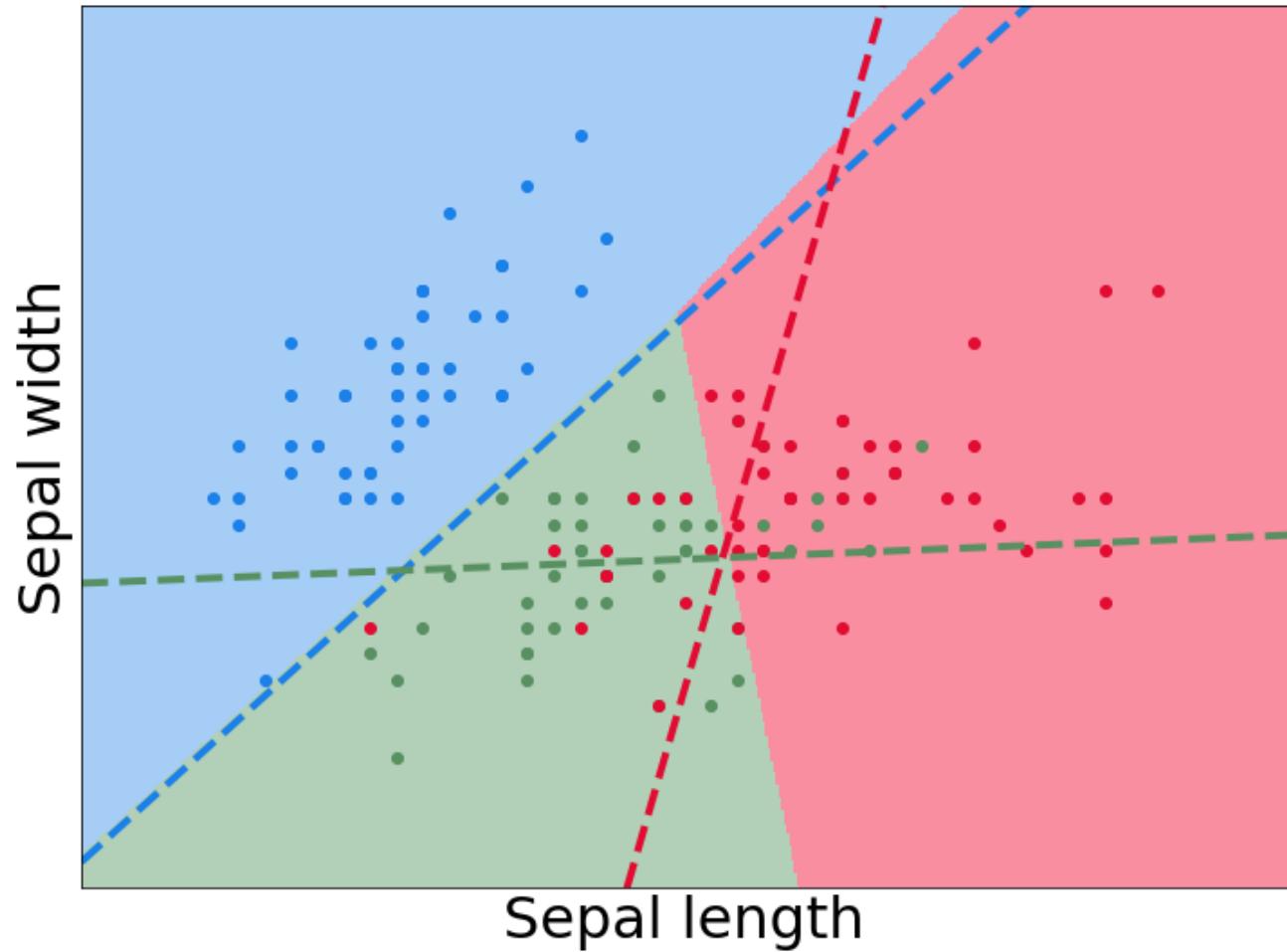
The model should be obtained by running the algorithm on the entire dataset!

Multiclass Classification

Logistic regression assumes that there are only two class values (e.g., +1/-1)

What if we have more?
(e.g., none, soft, hard or Setosa/Versicolour/Virginica)

- For each class, it creates one classifier that predicts the target class against all the others
- Given three classes A, B, C, it computes three models
 - One that predicts A against B and C
 - One that predicts B against A and C, and
 - One that predicts C against A and B
- Then, given an example, all the three classifiers are applied and the label with the highest probability is returned
- Alternative approaches include the minimization of loss based on the multinomial loss fit across the entire probability distribution



One Versus Rest multiclass model using the Iris dataset

- Alternative approach is to change model itself to be multi-class
- For example, use model that directly predicts 3 (or more) probabilities and predicts multinomial distribution instead of binomial distribution
- Model will have more parameters, 2 vectors for 3 classes, 3 vectors for 4, etc.
- Optimisation routine minimises loss or maximises likelihood over the multiple classes at once

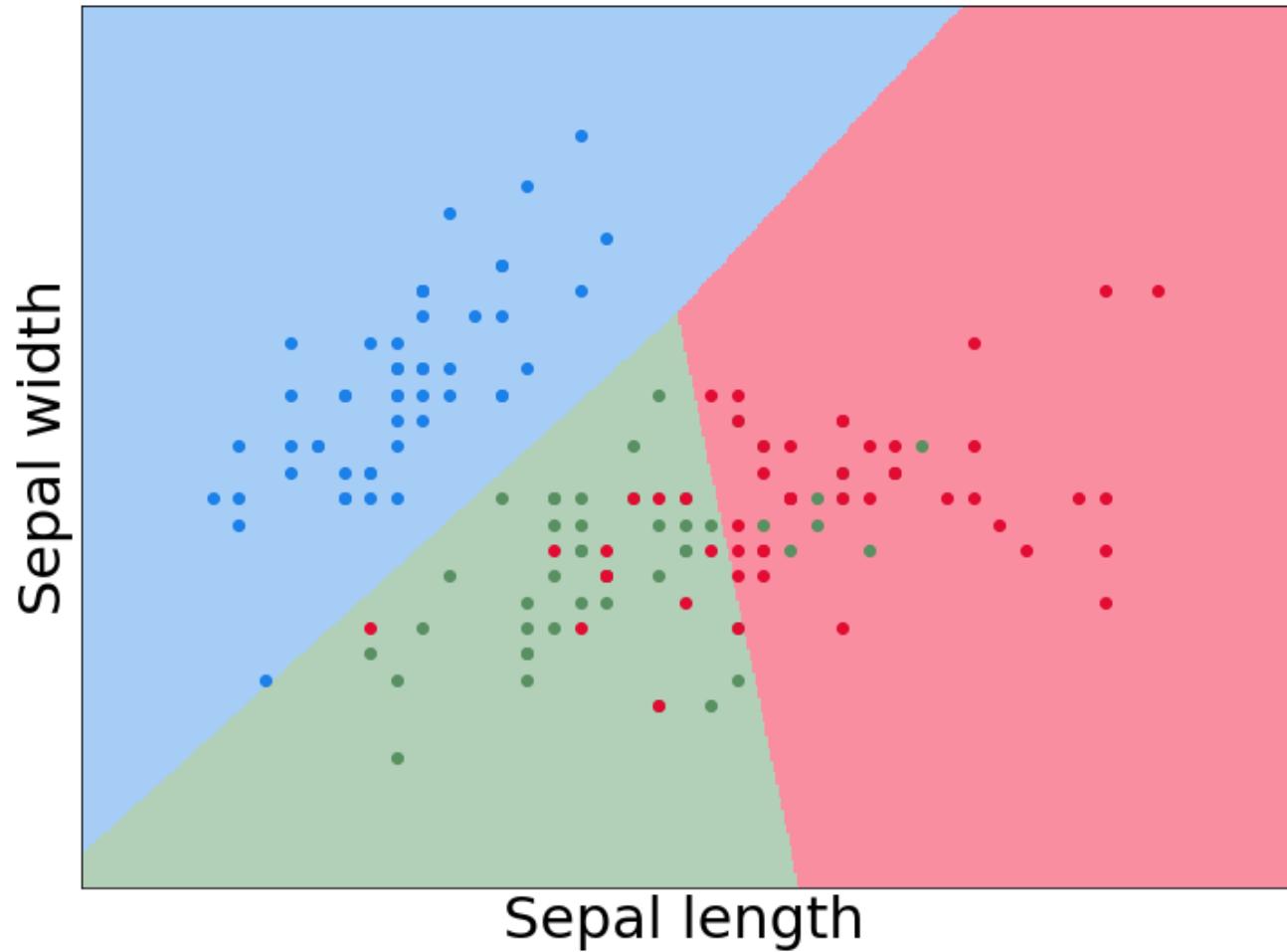
- Multinomial Logistic Regression model uses the softmax function:

$$P(Y_i = j) = \frac{e^{\vec{w}_j h(\vec{x}_i)}}{\sum_j e^{\vec{w}_j h(\vec{x}_i)}}$$

- using a weight vector now for each class j
- Don't actually need that many parameters
 - (model is overparameterized)
 - so can set the vector for one class to zero: $\vec{w}_C = \vec{0}$

$$P(Y_i = j) = \frac{e^{\vec{w}_j h(\vec{x}_i)}}{1 + \sum_j e^{\vec{w}_j h(\vec{x}_i)}} \quad \text{if } j \neq C$$

$$P(Y_i = C) = \frac{1}{1 + \sum_j e^{\vec{w}_j h(\vec{x}_i)}}$$



Multinomial multiclass model using the Iris dataset

Classification Metrics

Metrics for Performance Evaluation: Confusion Matrix

44

- Focus on the predictive capability of a model
- Confusion Matrix:

		PREDICTED CLASS	
		Yes	No
TRUE CLASS	Yes	TP true positives	FN false negatives
	No	FP false positives	TN true negatives

Metrics for Performance Evaluation: Accuracy

45

		PREDICTED CLASS	
ACTUAL CLASS		Yes	No
	Yes	#TP	#FN
	No	#FP	#TN

- Most widely-used metric:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

- Consider a 2-class problem
 - Number of Class 0 examples = 9990
 - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$
- Accuracy is misleading because model does not detect any class 1 example

		PREDICTED CLASS	
ACTUAL CLASS		Yes	No
	Yes	Cost(TP)	Cost(FN)
	No	Cost(FP)	Cost(TN)

Cost(x): Cost of misclassifying examples of type x

Computing Cost of Classification

48

Cost Matrix		PREDICTED CLASS	
ACTUAL CLASS	C(.)	+	-
	+	-I	100
	-	I	0

Model M ₁	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	150	40
	-	60	250

Accuracy = 80%

Cost = 3910

Model M ₂	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	250	45
	-	5	200

Accuracy = 90%

Cost = 4255

Costs can be used to evaluate
an existing classification models

Or can be used by the algorithm
to guide the search for the model

Some algorithms can use the cost matrix to
build the model (e.g., decision trees)

Finder File Edit View Go Window Help

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 -C 0.25 -M 2

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66

More options...

(Nom) class

Start Stop

Result list (right-click for options)
09:41:46 - trees.J48

Classifier output

	Correctly Classified Instances	705	70.5	%
Incorrectly Classified Instances	295	295	29.5	%
Kappa statistic	0.2467			
Total Cost	295	0.295		
Average Cost		0.3467		
Mean absolute error		0.4796		
Root mean squared error		82.5233 %		
Relative absolute error		104.6565 %		
Total Number of Instances	1000			

== Detailed Accuracy By Class ==

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
good	0.84	0.61	0.763	0.84	0.799	0.639	good
bad	0.39	0.16	0.511	0.39	0.442	0.639	bad
Weighted Avg.	0.705	0.475	0.687	0.705	0.692	0.639	

== Confusion Matrix ==

		a b <-- classified as
588	112	a = good
183	117	b = bad

Status OK Log x 0

TextWrangler File Edit Text View Search Go Window #! Help

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose

Cost Matrix Editor

Test options

Use training data

Supplied test set

Cross-validation

Percentage split

Output predictions

Output probabilities

Output class names

Output additional attributes

Cost-sensitive evaluation Set...

Random seed for XVal / % Split 1

Preserve order for % Split

Output source code WekaClassifier

OK

ZIP

20140311-JustSayDesigner Breakout_Final.zip

Gamification

Cost Matrix Editor

Defaults

Open...

Save...

Classes: 2

Resize

untitled text

T. (New Document)

153 % A202 : no

154 %

155 %

156 %

157 %

158 %

159 %

160 %

161 %

162 %

163 % 8. Cost Matrix

164 % This dataset requires use of a cost matrix (see below)

165 %

166 %

167 %

168 %

169 %

170 %

171 %

172 %

173 %

174 %

175 %

% 1 2

% -----

% 1 0 1

% -----

% 2 5 0

%

% (1 = Good, 2 = Bad)

%

% the rows represent the actual classification and the columns

% the predicted classification.

%

% It is worse to class a customer as good when they are bad (5),

% than it is to class a customer as bad when they are good (1).

%

Last saved: (Never)

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 -C 0.25 -M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

09:41:46 - trees.J48
10:28:20 - trees.J48

Classifier output

== Stratified cross-validation ==
== Summary ==

	Correctly Classified Instances	705	70.5 %
	Incorrectly Classified Instances	295	29.5 %
Kappa statistic		0.2467	
Total Cost		1027	
Average Cost		1.027	
Mean absolute error		0.3467	
Root mean squared error		0.4796	
Relative absolute error		82.5233 %	
Root relative squared error		104.6565 %	
Total Number of Instances		1000	

incorrectly classified examples

total cost has changed

== Detailed Accuracy By Class ==

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
09:41:46 - trees.J48	0.84	0.61	0.763	0.84	0.799	0.639	good
10:28:20 - trees.J48	0.39	0.16	0.511	0.39	0.442	0.639	bad
Weighted Avg.	0.705	0.475	0.687	0.705	0.692	0.639	

== Confusion Matrix ==

	a	b	<-- classified as
588	112	117	a = good
183	117	112	b = bad

Status OK Log x 0

- Alternatives to accuracy, introduced in the area of information retrieval and search engine
- **Precision**
 - Percentage of items classified as positive that are actually positive
 - In the information retrieval context represents the percentage of actually good documents that have been shown as a result.
- **Recall**
 - Percentage of positive examples that are classified as positive
 - In the information retrieval context, recall represents the percentage of good documents shown with respect to the existing ones.

The higher the precision, the lower the FPs

$$\text{Precision}(p) = \frac{TP}{TP + FP} = \frac{a}{a + c}$$

$$\text{Recall}(r) = \frac{TP}{TP + FN} = \frac{a}{a + b}$$

$$F1 - \text{measure} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

The higher the recall, the lower the FNs

The higher the F1, the lower the FPs & FNs

- Precision is biased towards TP) & FP
- Recall is biased towards TP & FN
- F1-measure is biased towards all except TN
it is high when both precision and recall are reasonably high

- Sensitivity evaluates the ability to correctly identify the elements of the positive class and it is computed as the true positive rate (TPR)

$$\text{TPR} = \text{TP}/(\text{TP}+\text{FN})$$

- Specificity estimates the probability to correctly identify the elements of the negative class and it is computed as the true negative rate

$$\text{TNR} = \text{TN}/(\text{TN}+\text{FP})$$

How to compare the relative
performance among competing models?

How to Compare the Performance of Two Models?

57

- Suppose we have two models
 - Model M_A with an accuracy = 82% computed using 10-fold crossvalidation
 - Model M_B with an accuracy = 80% computed using 10-fold crossvalidation
- How much confidence can we place on accuracy of M_A and M_B ?
- Can we say M_A is better than M_B ?
- Can the difference in performance measure be explained as a result of random fluctuations in the test set?

How do we know that the difference in performance is not just due to chance?

We computes the odds of it!

Apply the t-test and compute the p-value

The p-value represents the probability that the reported difference is due to chance

- Generate the k folds and for each configuration compute the performance of model A and B,
- $\theta_1^A \dots \theta_k^A \quad \theta_1^B \dots \theta_k^B$
compute mean and standard deviation of the differences:

$$\delta_i = \theta_i^A - \theta_i^B \quad \mu_\delta = \frac{1}{k} \sum_i \delta_i \quad \sigma_\delta = \sqrt{\frac{1}{k} \sum_i (\delta_i - \mu_\delta)^2}$$

- And two hypotheses, the null hypothesis and the alternative hypothesis
- We apply the t-test to check whether we can reject the null hypothesis H_0 with a target confidence

ALGORITHM 22.4. Paired t -Test via Cross-Validation

PAIRED t -TEST(α, K, \mathbf{D}):

- 1 $\mathbf{D} \leftarrow$ randomly shuffle \mathbf{D}
 - 2 $\{\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_K\} \leftarrow$ partition \mathbf{D} in K equal parts
 - 3 **foreach** $i \in [1, K]$ **do**
 - 4 $M_i^A, M_i^B \leftarrow$ train the two different classifiers on $\mathbf{D} \setminus \mathbf{D}_i$
 - 5 $\theta_i^A, \theta_i^B \leftarrow$ assess M_i^A and M_i^B on \mathbf{D}_i
 - 6 $\delta_i = \theta_i^A - \theta_i^B$
 - 7 $\hat{\mu}_\delta = \frac{1}{K} \sum_{i=1}^K \delta_i$
 - 8 $\hat{\sigma}_\delta^2 = \frac{1}{K} \sum_{i=1}^K (\delta_i - \hat{\mu}_\delta)^2$
 - 9 $Z_\delta^* = \frac{\sqrt{K}\hat{\mu}_\delta}{\hat{\sigma}_\delta}$
 - 10 **if** $Z_\delta^* \in (-t_{\alpha/2, K-1}, t_{\alpha/2, K-1})$ **then**
 - 11 Accept H_0 ; both classifiers have similar performance
 - 12 **else**
 - 13 Reject H_0 ; classifiers have significantly different performance
-

Comparing the Performance of 2 Models using k-fold Cross-validation

61

- First decide on a confidence level, e.g. 95%
 - Corresponds to false discovery (false positive) rate: $\alpha = 5\%$
 - How frequently you are willing to declare difference when there is none
- Apply k-fold cross-validation to each model
 - Obtaining k evaluations for each algorithm over same folds
- Apply Student's t-test and compute p-value to determine whether reported difference is statistically significant
 - If $p\text{-value} > \alpha$ then difference is not significant (can claim nothing)
 - If $p\text{-value} < \alpha$ then difference is significant (claim one better than the other)
 - Note that the t-test can be paired or unpaired

“paired” when the estimates
are from the same datasets

“unpaired” when the estimates
are from different datasets

Multiple Hypothesis Testing

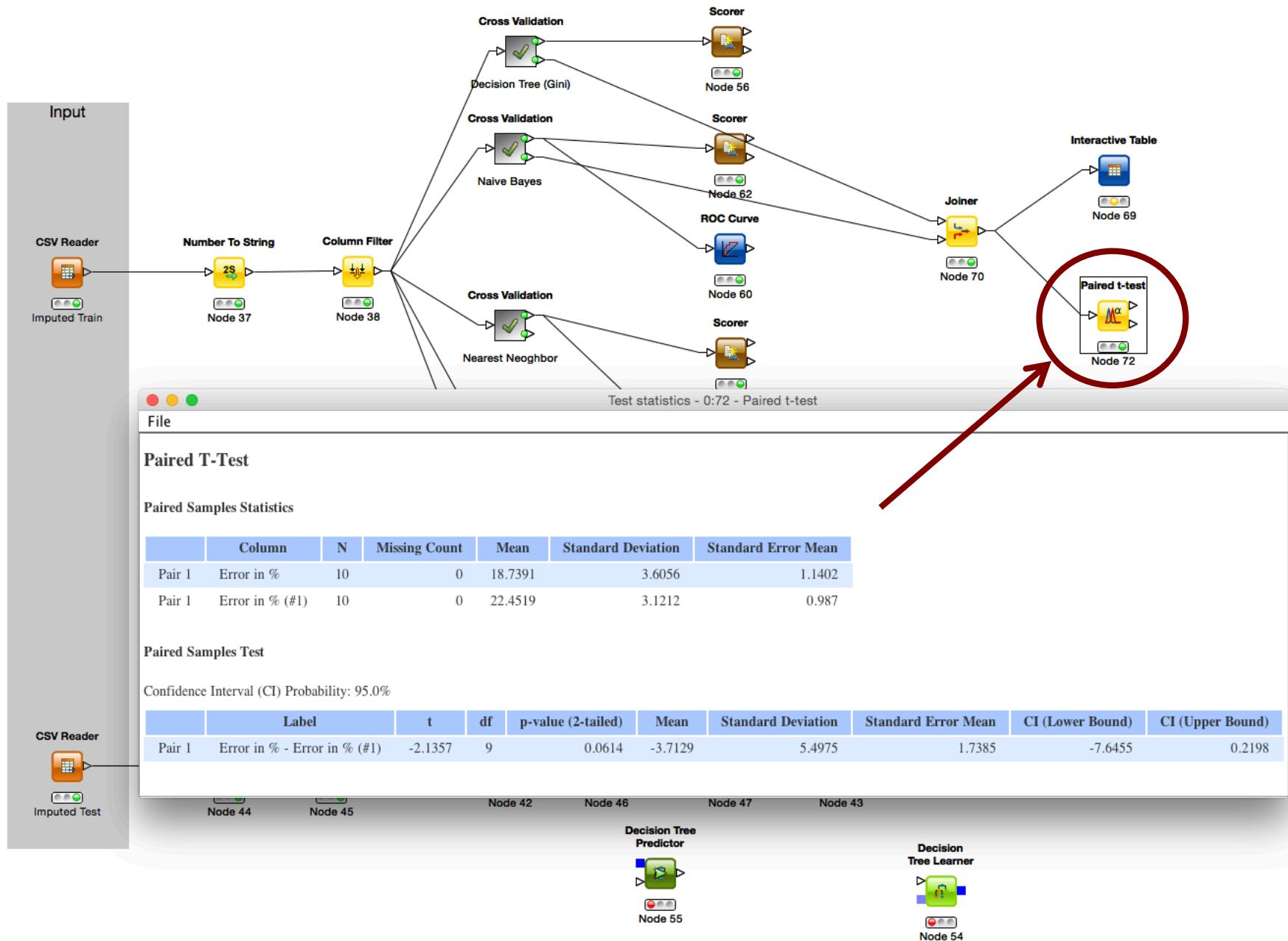
- Say that you perform a statistical test with a 0.05 threshold, but you repeat the test on twenty different observations.
- For example, you want to compare the performance of several classification algorithms
- Assume that all of the observations are explainable by the null hypothesis
- What is the chance that at least one of the observations will receive a p-value less than 0.05?

- Say that you perform a statistical test with a 0.05 threshold (95% confidence level), but you repeat the test on 20 different observations. What is the chance that at least one of the observations will receive a p-value less than 0.05?
- $P(\text{making a mistake}) = 0.05$
- $P(\text{not making a mistake}) = 0.95$
- $P(\text{not making any mistake}) = 0.95^{20} = 0.358$
- $P(\text{making at least one mistake}) = 1 - 0.358 = 0.642$
- There is a 64.2% chance of making at least one mistake.

- Assume that individual tests are independent.
- Divide the desired p-value threshold by the number of tests performed.
- Example
 - We now have, the threshold set to $0.05/20 = 0.0025$.
 - $P(\text{making a mistake}) = 0.0025$
 - $P(\text{not making a mistake}) = 0.9975$
 - $P(\text{not making any mistake}) = 0.9975^{20} = 0.9512$
 - $P(\text{making at least one mistake}) = 1 - 0.9512 = 0.0488$

- They do not make any assumption about the distribution of the variable in the population
- Mann-Whitney U Test
 - Nonparametric equivalent of the independent t-test
- Wilcoxon matched-pairs signed rank test
 - Used to compare two related groups

	Nonparametric tests		Parametric tests
	Nominal data	Ordinal data	Ordinal, interval, ratio data
One group	Chi square goodness of fit	Wilcoxon signed rank test	One group t-test
Two unrelated groups	Chi square	Wilcoxon rank sum test, Mann-Whitney test	Student's t-test
Two related groups	McNemar's test	Wilcoxon signed rank test	Paired Student's t-test
K-unrelated groups	Chi square test	Kruskal -Wallis one-way analysis of variance	ANOVA
K-related groups		Friedman matched samples	ANOVA with repeated measurements

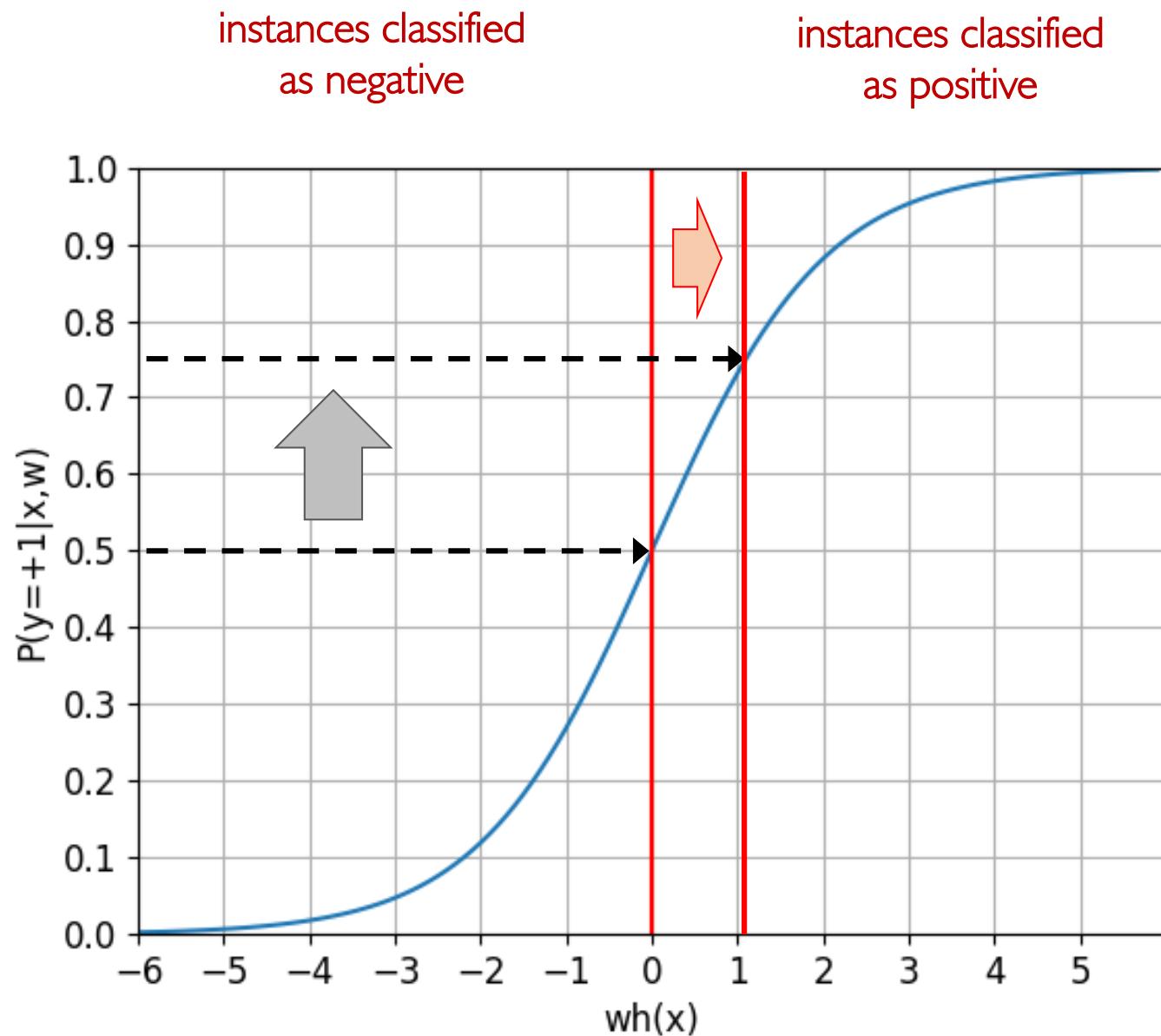


Probabilistic Classifiers

- Up to now we used logistic regression to predict classifier labels, however, logistic regression returns a probability

$$P(y_i|\vec{x}_i)$$

- Given an example x_i , its predicted class is the label with the largest probability so it is equivalent to using a threshold of 0.5 to decide which class to assign to an example
- However, we can use a different threshold and for instance label as positive only examples we return a 1 only when $P(+|x) > 0.75$
- This would label as positive only cases for which we are more confident that should be labeled as positive.



How the Classification Threshold Influence Precision and Recall?

- Suppose we use a near one threshold to classify positive examples
- Then, we will classify as positives only examples for which we are very confident (this is a pessimistic classifier)
- Precision will be high
 - In fact, we are not likely to produce few false positives
- Recall will be low
 - In fact, we are likely to produce many false negatives

How the Classification Threshold Influence Precision and Recall?

74

- Suppose we use a near zero threshold to classify positive examples
- Then, we will classify everything as positives
(this is an optimistic classifier)
- Precision will be low as we are going to generate the maximum number of false positives (everything is positive!)
- Recall will be high since by classifying everything as positive we are going to generate the minimum number of false negatives

We can use the threshold to optimize our precision and recall

a higher the threshold, increases precision and lower recall

a lower threshold, decreases precision and increase recall

- In the notebook, a simple logistic regression model applied to the loans data returns the confusion matrix below, corresponding to a precision of 0.82 and recall of 0.99

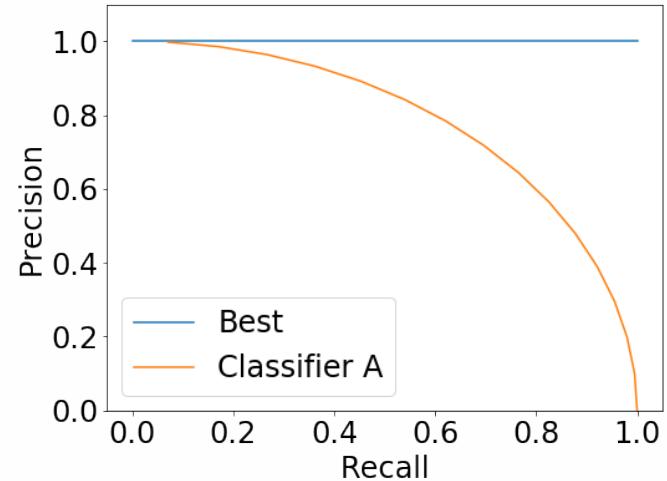
	Classified -I	Classified +I
Labeled -I	1212	21910
Labeled +I	1057	98283

- By increasing the classification threshold for positive (+I) examples to 0.75 we obtained a new confusion matrix with precision of 0.86 and a recall of 0.80

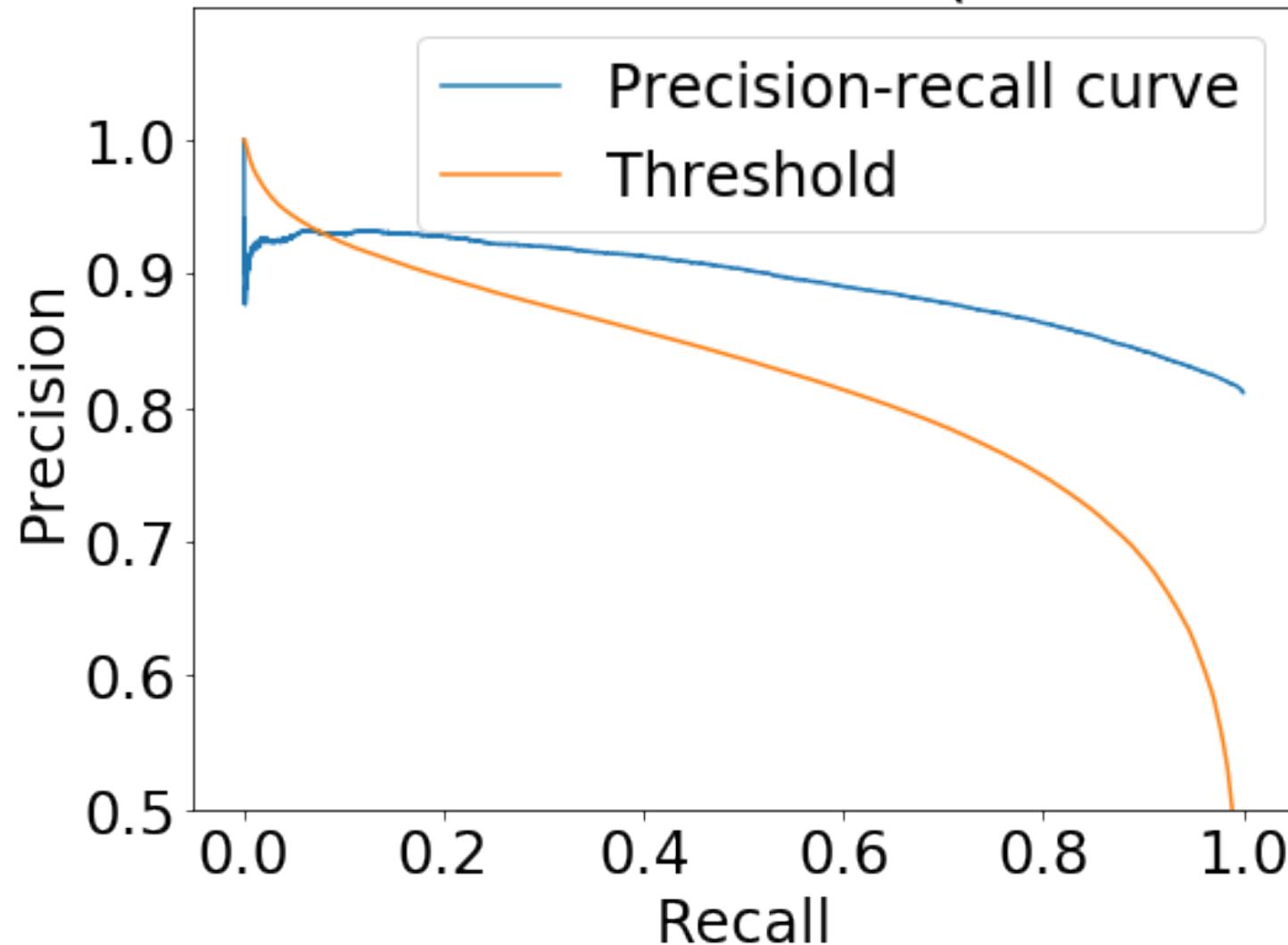
	Classified -I	Classified +I
Labeled -I	10632	12490
Labeled +I	20106	79234

- Overall, we reduced the number of false positives (we did not accept risky loans and we were better at identifying risky loans)

- Plot precision as a function of recall for varying threshold values
- The best classifier would be the one that has always a precision equal to one (but never happens)
- More in general classifiers will show of different shapes
- How to decide among more classifiers?
 - Use the area under the curve (the nearer to one, the better)
 - Use F1 measure



Precision-Recall Curve (AUC=0.89)



Precision-recall curve for the loan dataset (run the python notebook for details)

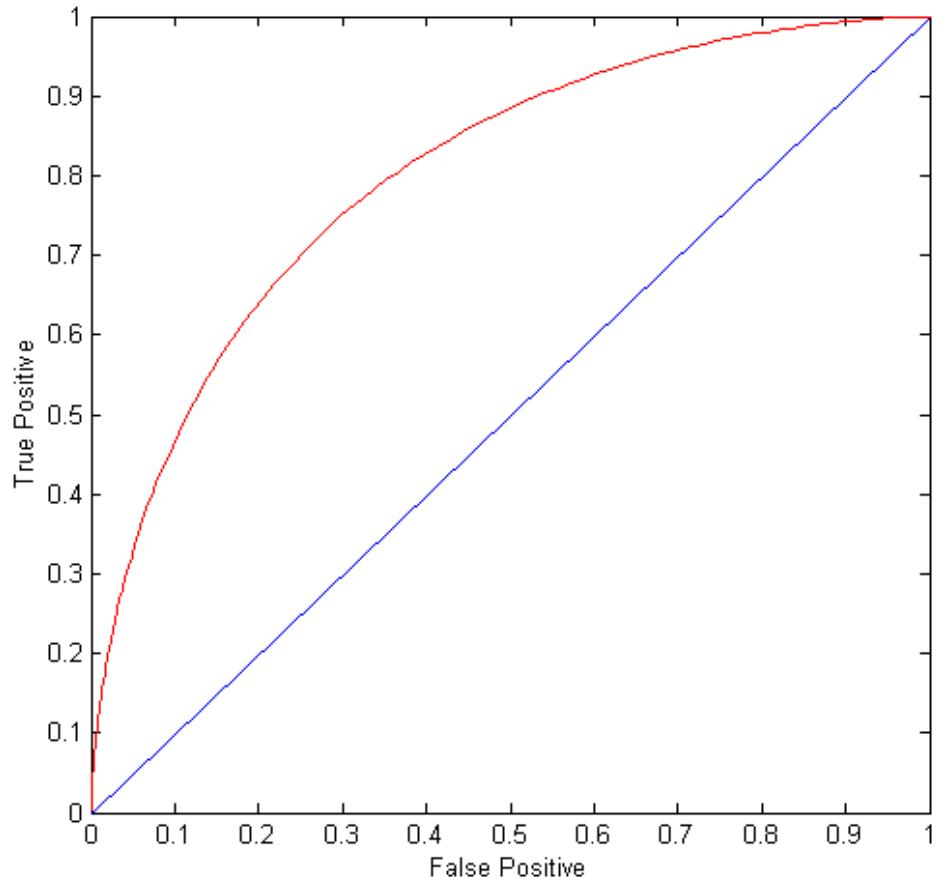
Receiver Operating Characteristic (ROC) Curves

ROC Curve (Receiver Operating Characteristic)

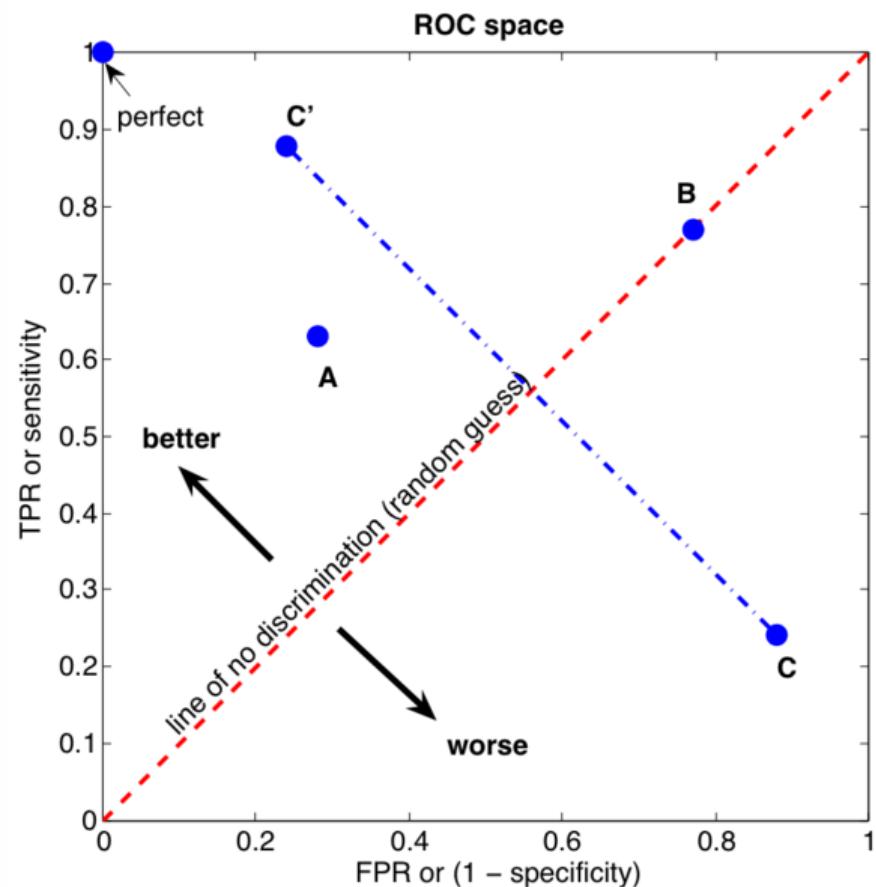
- Developed in 1950s for signal detection theory to analyze signals
- Plot the True Positive Rate ($TPR=TP/(TP+FN)$) against the False Positive Rate ($FPR=FP/(TN+FP)$)
- Performance of each classifier represented as a point on the ROC curve
- Changing the classification threshold, sample distribution or cost matrix changes the location of the point

ROC Curve

- (FPR, TPR)
 - (0,0): declare everything to be negative class
 - (1,1): declare everything to be positive class
 - (0,1): ideal
- Diagonal line:
 - Random guessing
 - Below diagonal line, prediction is opposite of the true class



ROC Curve

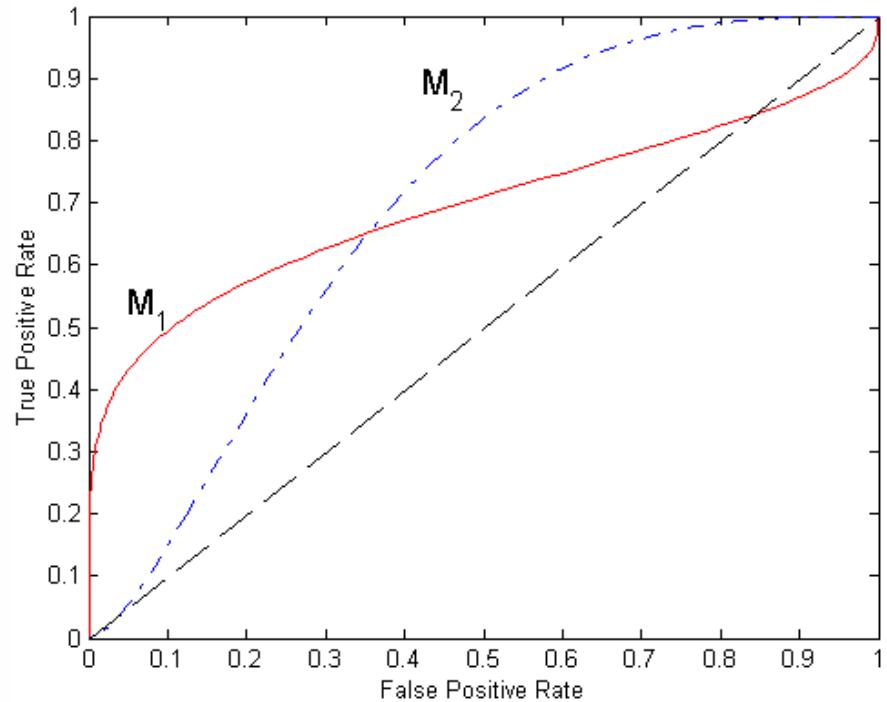


	TP	FP	FN	TN	Total
A	63	28	37	72	91
B	77	77	23	23	154
C	24	88	76	12	112
C'	88	24	12	76	112

Below each row are the corresponding TPR, FPR, and ACC values:

	TPR	FPR	ACC
Model A	0.63	0.28	0.68
Model B	0.77	0.77	0.50
Model C	0.24	0.88	0.18
Model C'	0.88	0.24	0.82

- No model consistently outperform the other
 - M_1 is better for small FPR
 - M_2 is better for large FPR
-
- Area Under the ROC curve
 - Ideal, area = 1
 - Random guess, area = 0.5



There are techniques similar
to ROC in other areas

Lift charts are an example

Which model should we prefer?
(Model selection)

- Model selection criteria attempt to find a good compromise between:
 - The complexity of a model
 - Its prediction accuracy on the training data
- Reasoning: a good model is a simple model that achieves high accuracy on the given data
- Also known as Occam's Razor :
the best theory is the smallest one
that describes all the facts



William of Ockham, born in the village of Ockham in Surrey (England) about 1285, was the most influential philosopher of the 14th century and a controversial theologian.

- Among the several algorithms, which one is the “best”?
 - Some algorithms have a lower computational complexity
 - Different algorithms provide different representations
 - Some algorithms allow the specification of prior knowledge
- If we are interested in the generalization performance, are there any reasons to prefer one classifier over another?
- Can we expect any classification method to be superior or inferior overall?
- According to the No Free Lunch Theorem, the answer to all these questions is no

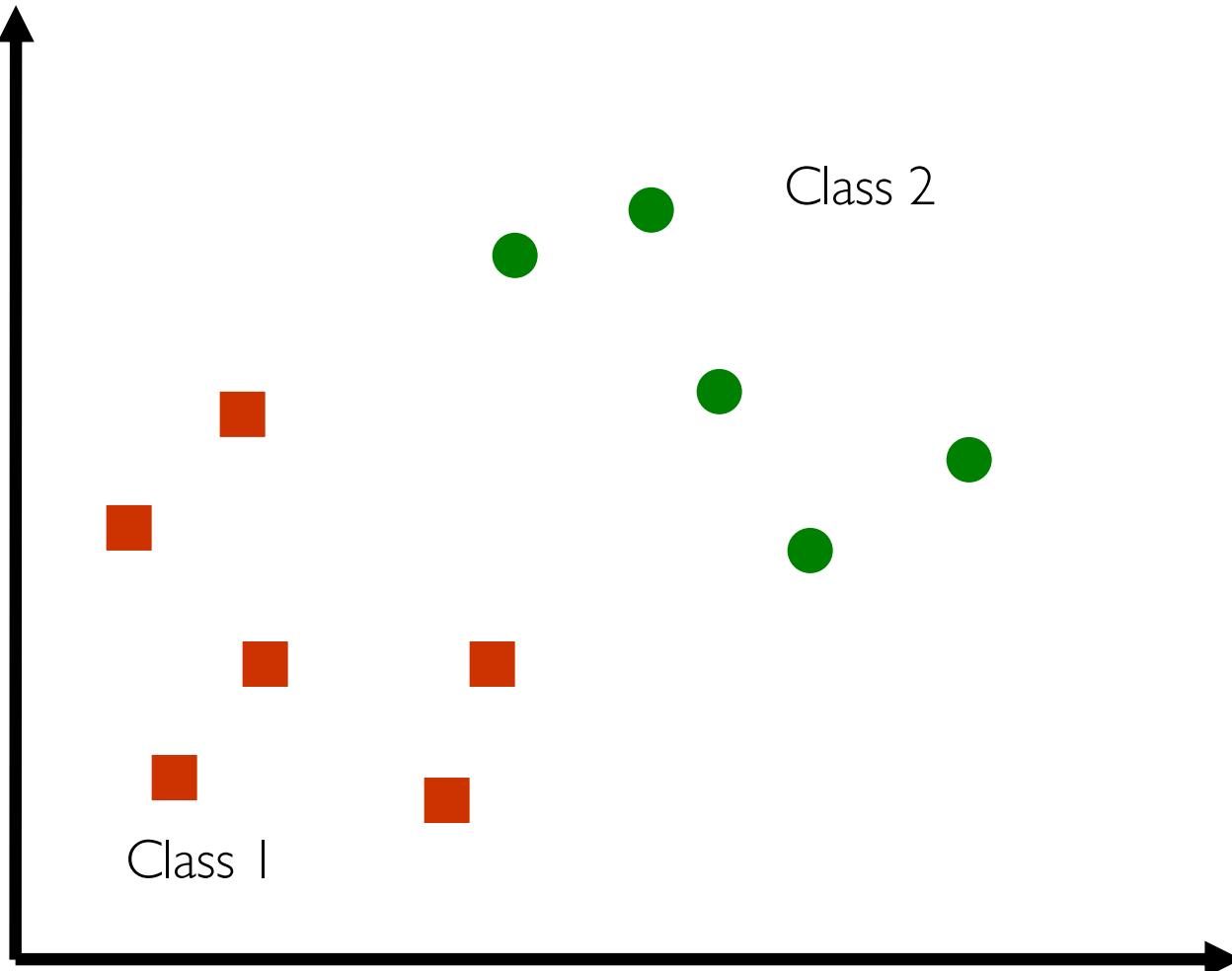
- If the goal is to obtain good generalization performance, there are no context-independent or usage-independent reasons to favor one classification method over another
- If one algorithm seems to outperform another in a certain situation, it is a consequence of its fit to the particular problem, not the general superiority of the algorithm
- When confronting a new problem, this theorem suggests that we should focus on the aspects that matter most
 - Prior information
 - Data distribution
 - Amount of training data
 - Cost or reward
- The theorem also justifies skepticism regarding studies that “demonstrate” the overall superiority of a certain algorithm

- "[A]ll algorithms that search for an extremum of a cost [objective] function perform exactly the same, when averaged over all possible cost functions." [1]
- "[T]he average performance of any pair of algorithms across all possible problems is identical." [2]
- Wolpert, D.H., Macready, W.G. (1995), No Free Lunch Theorems for Search, Technical Report SFI-TR-95-02-010 (Santa Fe Institute).
- Wolpert, D.H., Macready, W.G. (1997), No Free Lunch Theorems for Optimization, IEEE Transactions on Evolutionary Computation 1, 67.

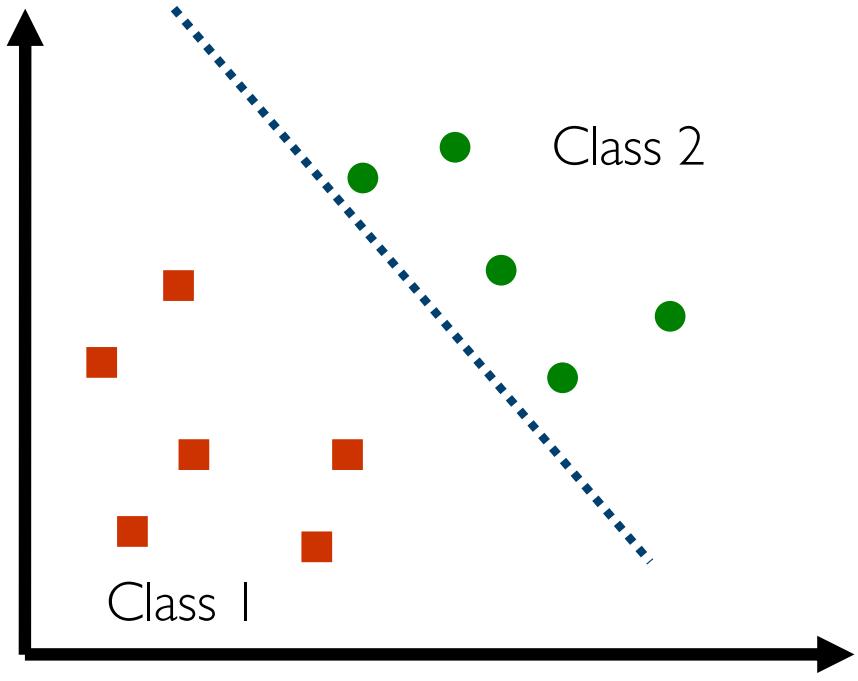
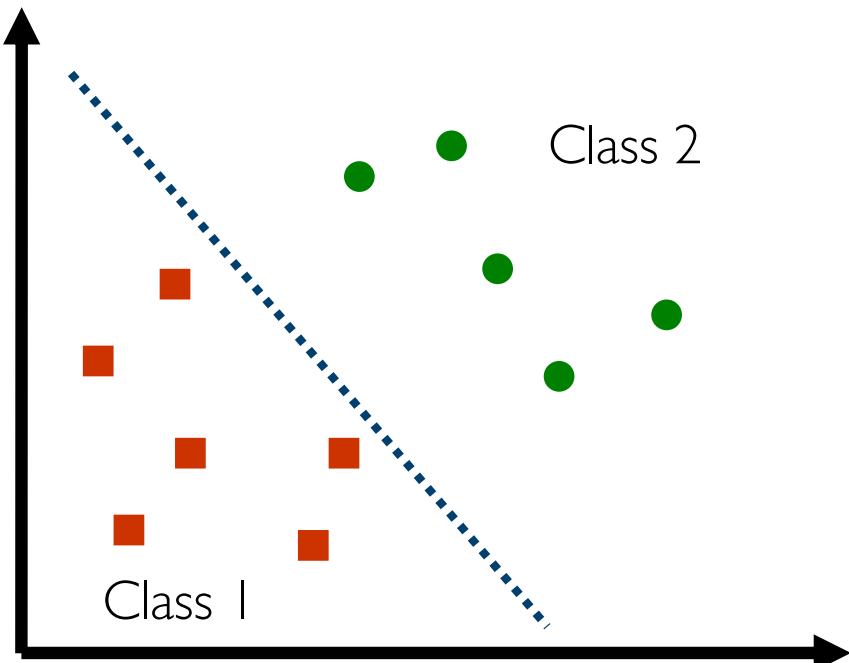
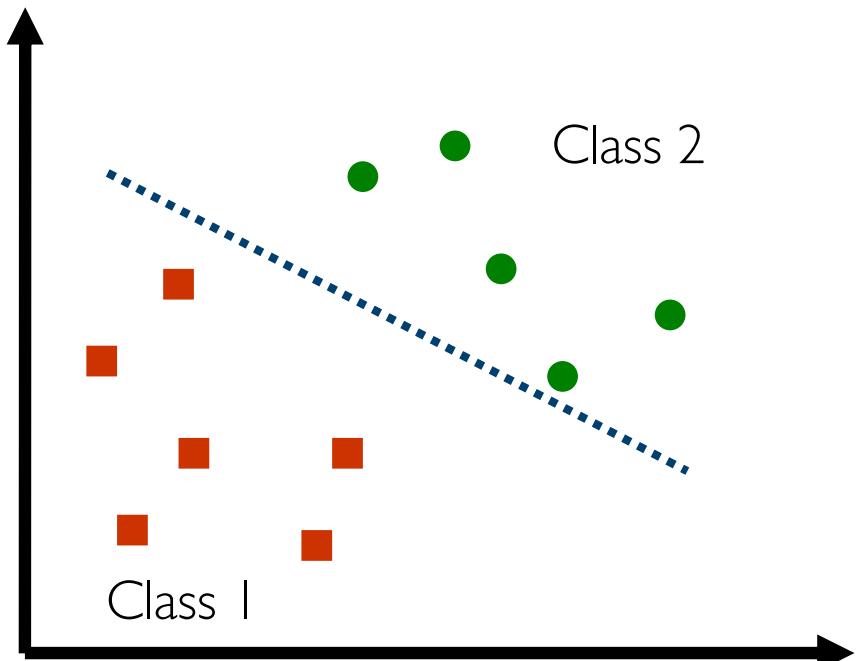
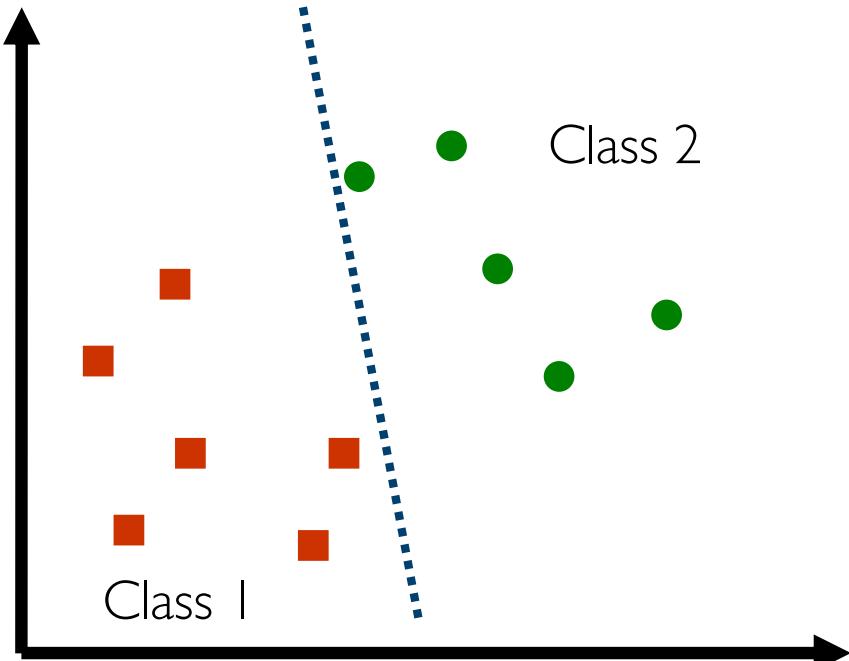
- <https://arxiv.org/pdf/1811.12808.pdf>
- <https://arxiv.org/pdf/1809.09446.pdf>

And now some names
you might know ...

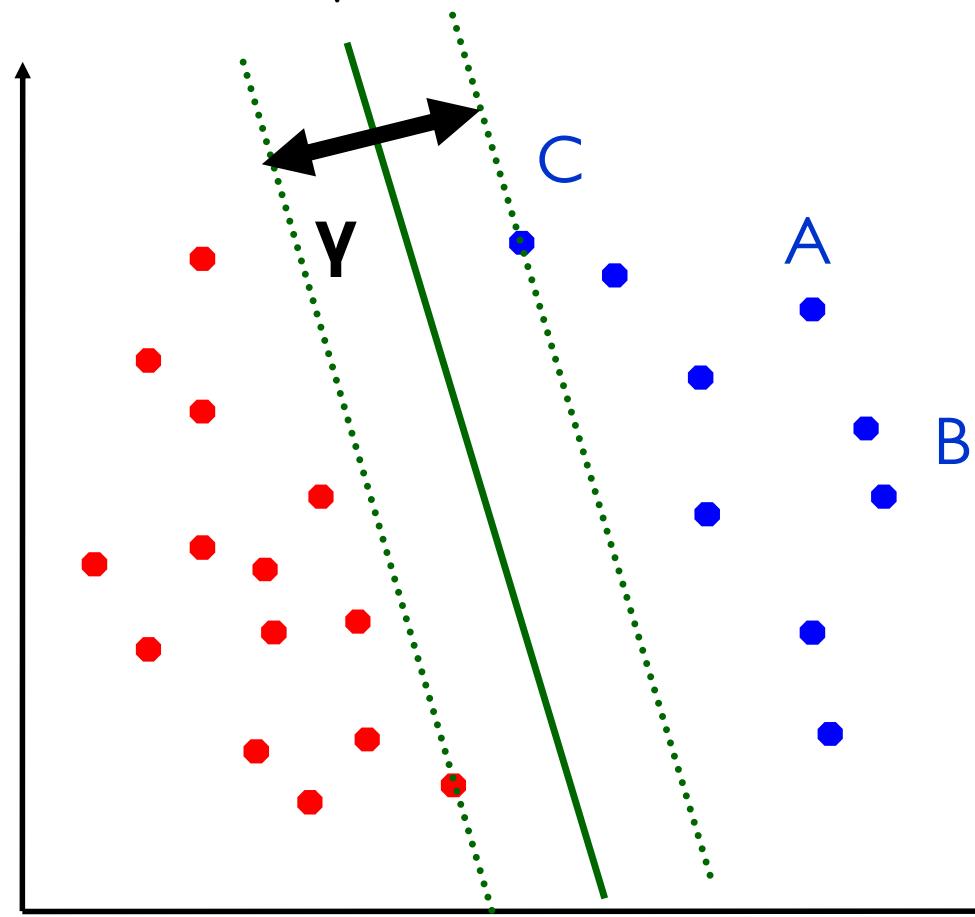
Support Vector Machines



Many decision boundaries can separate these two classes
Which one should we choose?



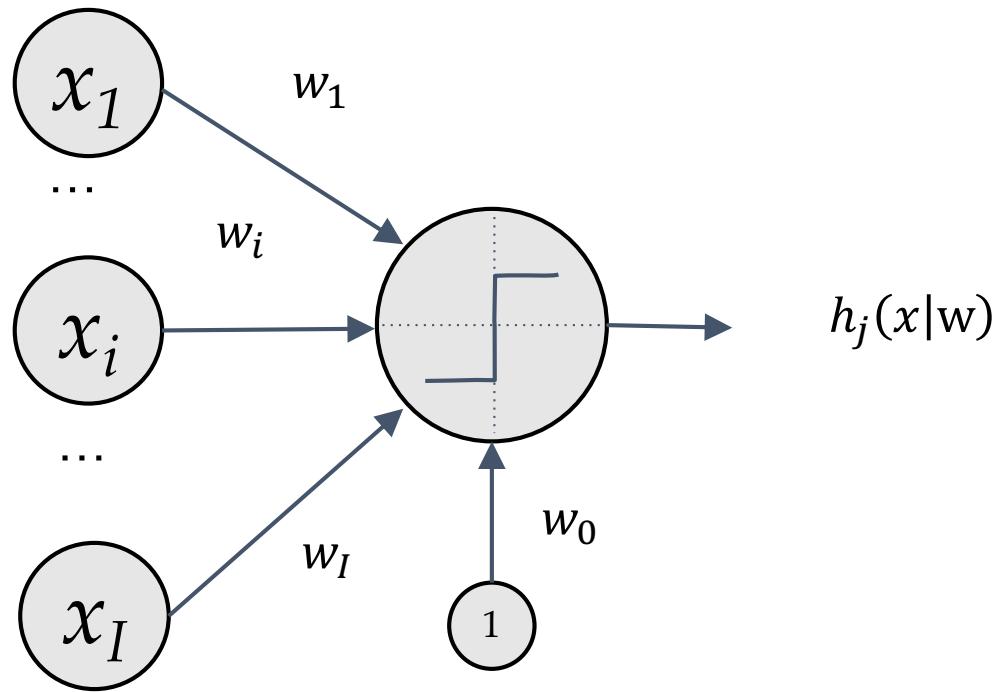
$$\mathbf{w}^T \mathbf{x} + b = 0$$



SVMs work by searching for the hyperplane that maximizes the margin or the largest γ such that

$$\forall i, y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq \gamma$$

Neural Networks



$$h_j(x|w, b) = h_j\left(\sum_{i=1}^I w_i \cdot x_i - b\right) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = h_j(w^T x)$$

Computation in an artificial neurons

- A perceptron computes the value of a weighted sum and returns its Sign (Thresholding)

$$h_j(x|w) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = \text{Sign}(w_0 + w_1 \cdot x_1 + \cdots + w_I \cdot x_I)$$

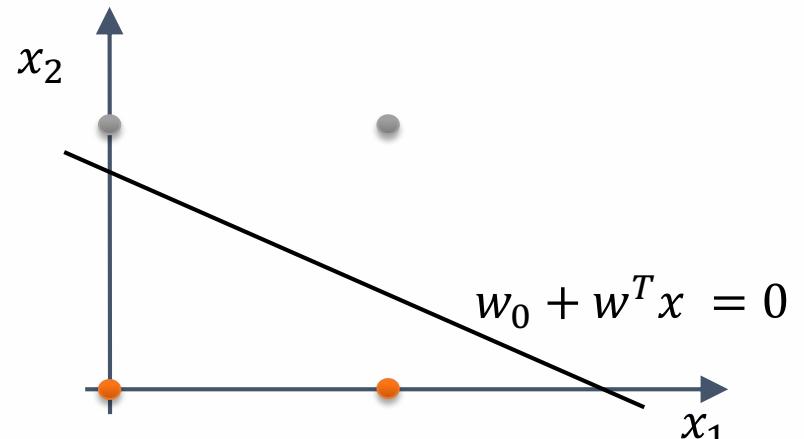
- It is basically a linear classifier for which the decision boundary is the hyperplane $w_0 + w_1 \cdot x_1 + \cdots + w_I \cdot x_I = 0$

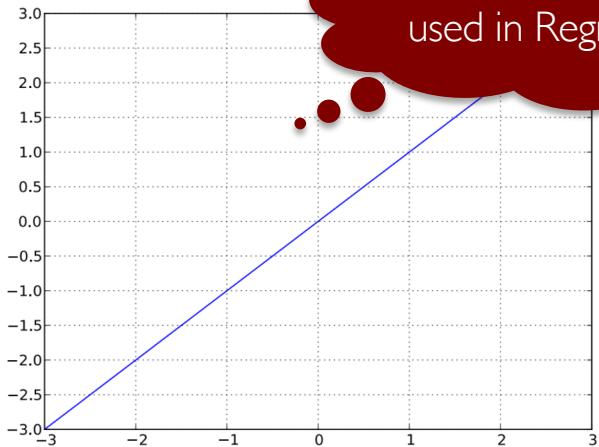
- In 2D, this turns into

$$w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 = 0$$

$$w_2 \cdot x_2 = -w_0 - w_1 \cdot x_1$$

$$x_2 = -\frac{w_0}{w_2} - \frac{w_1}{w_2} \cdot x_1$$

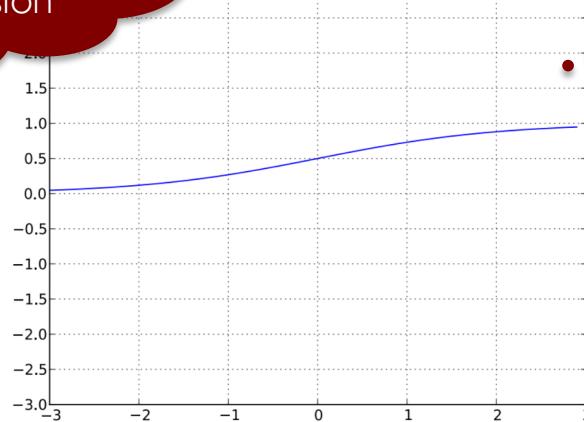




Linear activation function

$$g(a) = a$$

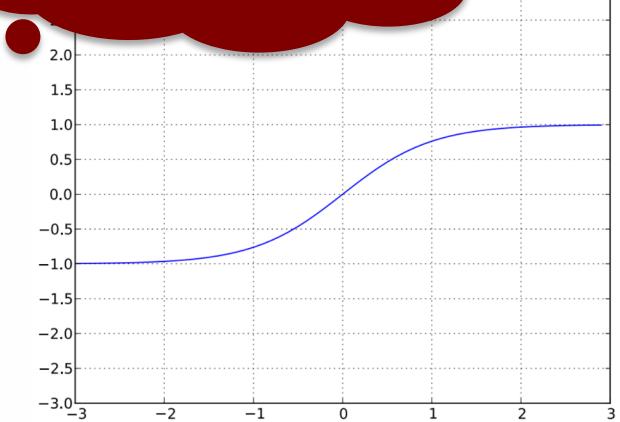
$$g'(a) = 1$$



Sigmoid activation function

$$g(a) = \frac{1}{1 + \exp(-a)}$$

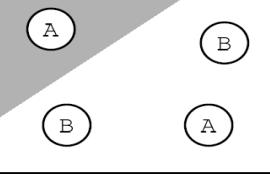
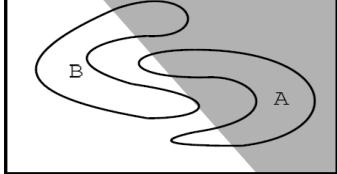
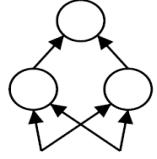
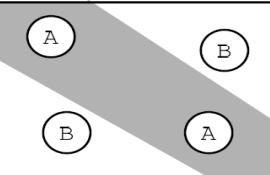
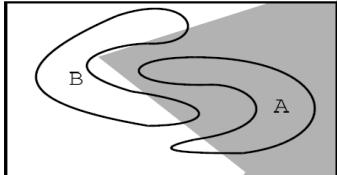
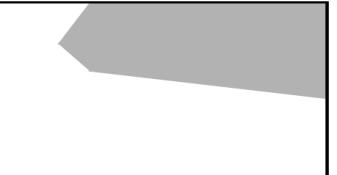
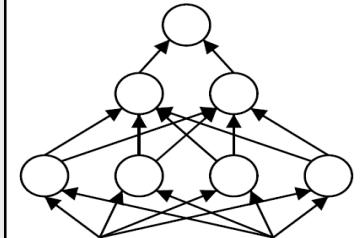
$$g'(a) = g(a)(1 - g(a))$$



Tanh activation function

$$g(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$$

$$g'(a) = 1 - g(a)^2$$

Topology	Type of Decision Region	XOR Problem	Classes with Meshed Regions	Most General Region Shapes
	Half bounded by hyperplanes			
	Convex Open or Closed Regions			
	Arbitrary Regions (Complexity limited by the number of nodes)	