

Introdução à Aprendizagem Automática (IAA)

SUSANA BRÁS

SUSANA.BRAS@UA.PT

IAA – L4

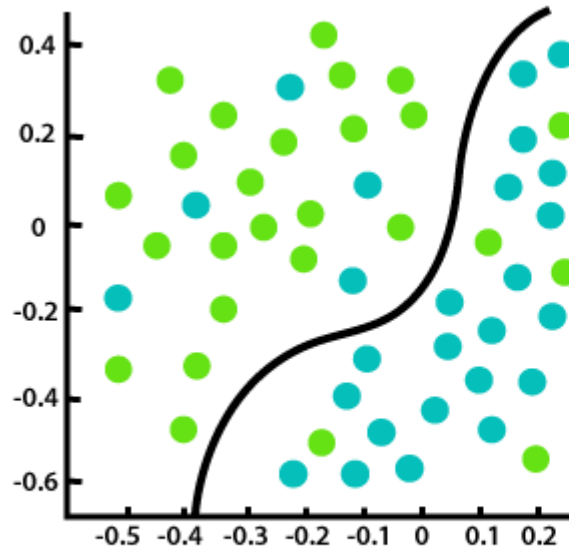
Logistic regression

- Classification
- Differences between linear and logistic regression
- Binary classification problem
- Sigmoid function
 - Equation
 - Graphical representation
 - Goal of the operation
- Logistic regression cost function
- Optimization algorithms, e.g. gradient descent
- Decision boundary: Definition, Importance
- Nonlinear separable data and its transformations to convert to a linear separable data space
- Overfitting: Definition, Consequences, Possible causes
- Regularization as an operation to reduce overfitting
 - Ridge
 - Lasso
- Multiclass classification problem using logistic regression
 - One vs all strategy
 - Softmax approach
- Advantages and disadvantages of logistic regression

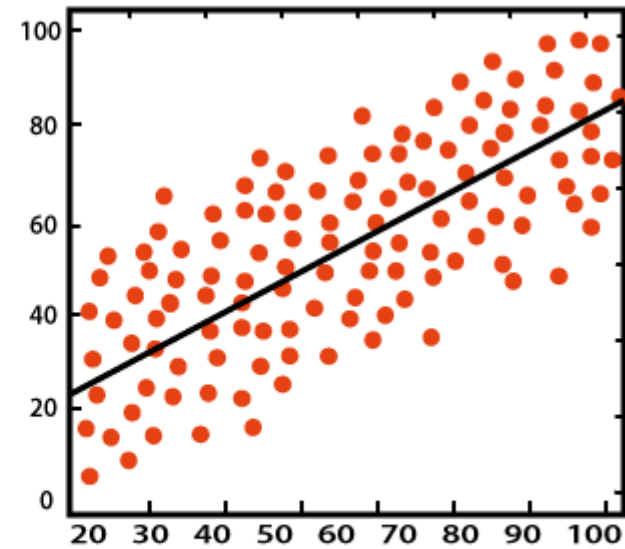
Supervised Learning

Regression and Classification algorithms are Supervised Learning algorithms. Both algorithms are used for prediction in Machine learning and work with labeled datasets. However, the difference between the two lies in their application to various machine learning problems.

The main difference between Regression and Classification algorithms is that Regression algorithms are used to predict continuous values, such as price, salary, age, etc., and Classification algorithms are used to predict/classify the discrete values, such as Male or Female, True or False, Spam or Not Spam, etc.



Classification



Regression

Data Matrix

matrix X (mxn)	feature x_1	feature x_2	feature x_n	Target Y
Example 1	$x_1^{(1)}$	$x_2^{(1)}$		$x_n^{(1)}$	$y^{(1)}$
Example 2	$x_1^{(2)}$	$x_2^{(2)}$		$x_n^{(2)}$	$y^{(2)}$
...					
Example i	$x_1^{(i)}$	$x_2^{(i)}$		$x_n^{(i)}$	$y^{(i)}$
...					
...					
Example m	$x_1^{(m)}$	$x_2^{(m)}$		$x_n^{(m)}$	$y^{(m)}$

x – input vector of features, attributes

y – output vector of labels, ground truth, target

m - number of training examples

n – number of features

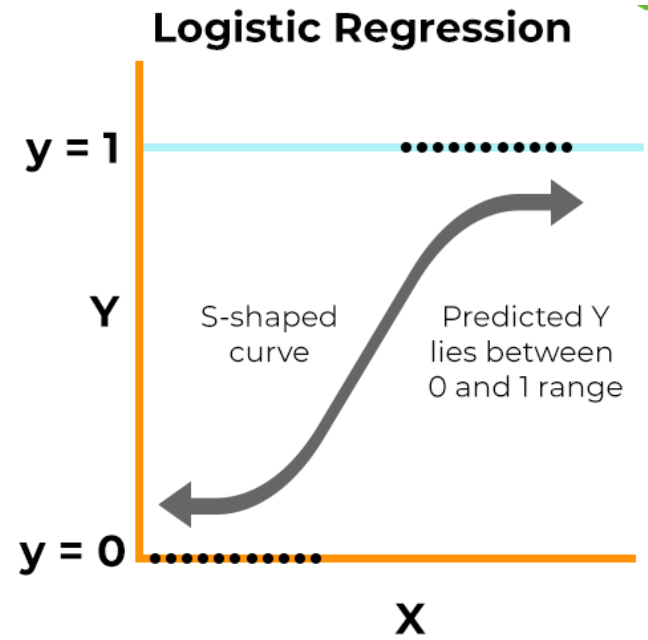
$h_{\theta}(x)$ - model (hypothesis)

θ - vector of model parameters

Training set: data matrix X (m rows, n columns)

Logistic Regression

Classification



Linear vs Logistic Regression

	Linear	Logistic
Nature of the Dependent Variable	The dependent variable is continuous and can take any real value.	The dependent variable is binary or categorical, representing two classes.
Equation Formulation	The linear equation predicts the value of the dependent variable directly.	The logistic function transforms the linear combination into a probability, and a threshold is applied for classification.
Output Interpretation	The output is interpreted as the expected value of the dependent variable.	The output is interpreted as the probability of the event occurring (class 1).
Applications	Used for predicting continuous outcomes, such as sales, temperature, or stock prices.	Applied in binary classification problems, like spam detection, medical diagnosis, or customer churn prediction.

A vertical bar on the left side of the slide, consisting of a wide red section and a thin blue section on its right edge.

Sigmoid Function Model

Classification – Logistic Regression

Logistic regression is a type of model of probabilistic statistical classification. It is used as a binary model to predict a binary response, the outcome of a categorical dependent variable (i.e., a class label), based on one or more variables.

The form of the logistic function is:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{-z}} = g(\theta^T x) = g(z)$$
$$z = \theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Logistic regression works by applying the logistic function to a linear combination of the input features.

1. Calculate a weighted sum of the input features (similar to linear regression).
2. Apply the logistic function (also called sigmoid function) to this sum, which maps any real number to a value between 0 and 1.
3. Interpret this value as the probability of belonging to the positive class.
4. Use a threshold (typically 0.5) to make the final classification decision.

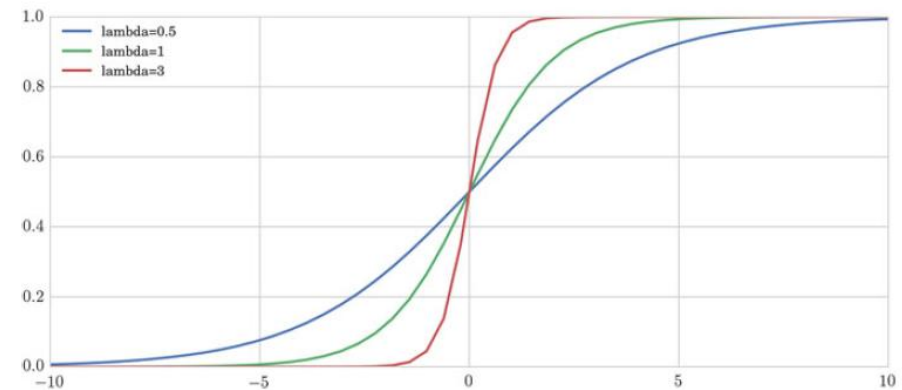


Fig. 6.10 Logistic function for different lambda values

Classification – Logistic Regression

The training process for logistic regression involves finding the best weights for the input features.

1. Initialize the weights (often to small random values).
2. For each training example:
 1. Calculate the predicted probability using the current weights.
 2. Compare this probability to the actual class label by calculating its loss.
3. Update the weights to minimize the loss (usually using some optimization algorithm, like gradient descent.)
4. Repeat Step 2 until log loss cannot get smaller (or other stop criteria was met).

Cost (loss) function

Logistic Regression Cost Function

$$\min_{\theta} J(\theta)$$

Linear regression model

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \vec{\theta}^T \vec{x}$$

Linear Regression cost function

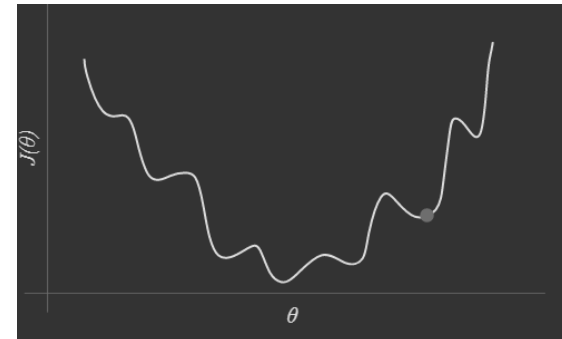
$$J = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Nonlinear logistic (sigmoid) model

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

If we use the same cost function as with linear regression, but now the hypothesis is a nonlinear function, $J(\theta)$ will be a non-convex function (has many local minima)

not efficient for optimization !



Logistic Regression Cost Function

$$\min_{\theta} J(\theta)$$

Linear regression model

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \vec{\theta}^T \vec{x}$$

Linear Regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

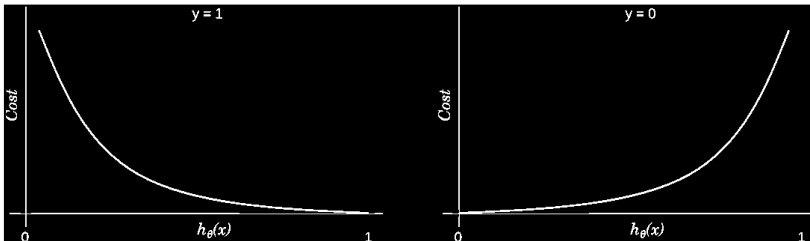
Note: $y = 0$ or 1 always

Log Loss function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

$$-\log(h_{\theta}(x))$$

$$-\log(1 - h_{\theta}(x))$$



Log Reg with gradient descent learning

$$\min_{\theta} J(\theta)$$

Goal

Initialize model parameters (e.g. $\theta = 0$)
Repeat until J converge {

Compute LogReg
Model prediction

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Compute LogReg cost
function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

Compute cost function
gradients

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Update parameters

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

Optimization algorithms

Gradient descent (learned in class) - updates the parameters in direction in which the gradient decreases most rapidly.

Other optimization algorithms

- Conjugate gradient
- AdaGrad, RMSProp
- Stochastic gradient descent with momentum
- ADAM (combination of RMSProp and stochastic optimization)
- BFGS (Broyden–Fletcher–Goldfarb–Shanno)
- Quasi-Newton methods (approximate the second derivative)

Characteristics

- Adaptive learning rate (α);
- Often faster than gradient descent; better convergence;
- Approximate (estimate) the true gradient over a mini-batch and not over the whole data;
- More complex algorithms

Decision Boundary

Logistic Regression - example

Training data: applicant's scores on two exams and the admission decision (historical data). Build a logistic regression model to predict whether a student gets admitted into a university.

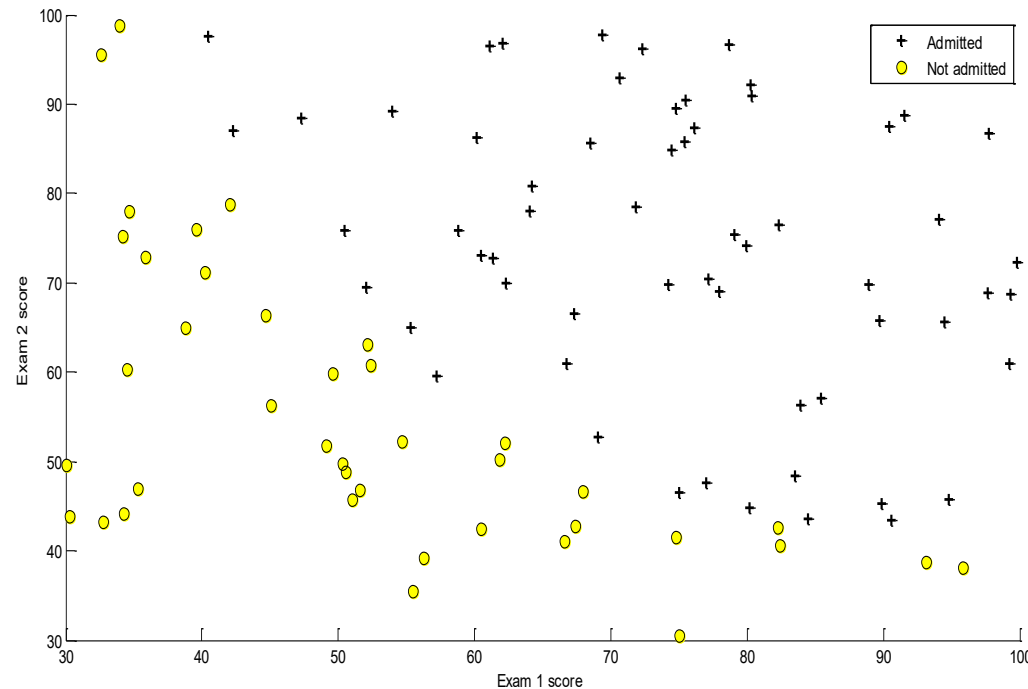
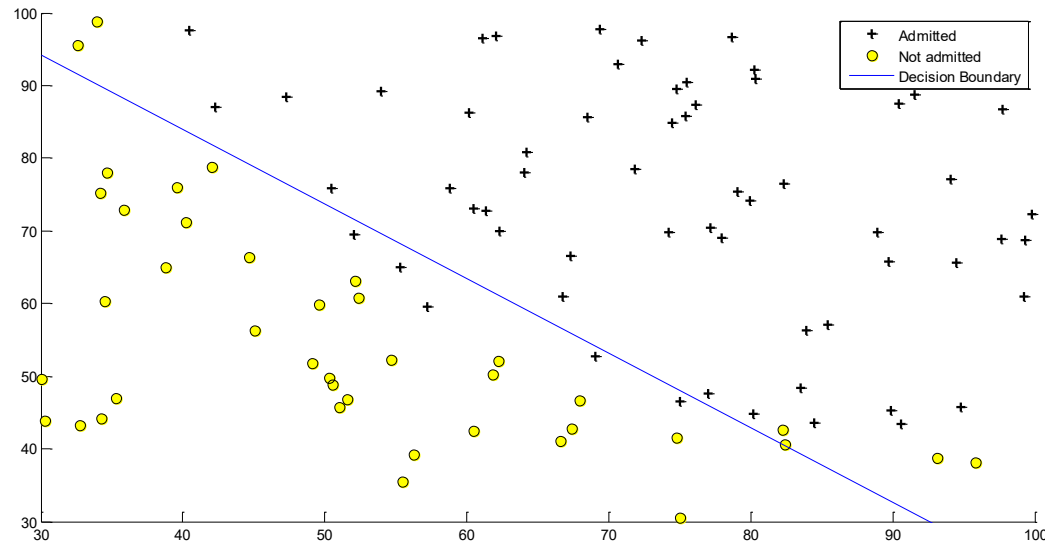


Fig. 1 Scatter plot of training data

Logistic Regression - example

Scatter plot of training data and linear decision boundary with the optimized parameters



$$z = \theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0 \Rightarrow \text{decision boundary}$$

if $z > 0 \Rightarrow g(z) > 0.5 \Rightarrow \text{predict class} = 1$

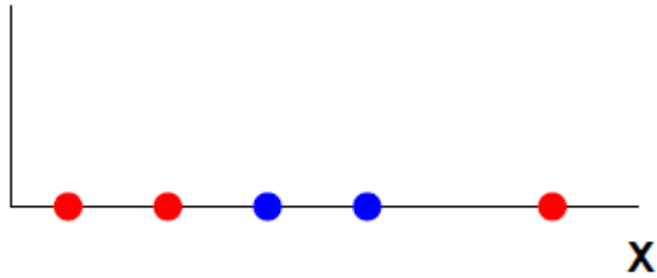
if $z < 0 \Rightarrow g(z) < 0.5 \Rightarrow \text{predict class} = 0$

A vertical bar on the left side of the slide, consisting of a wide red section and a thin blue section on its right edge.

Nonlinear Separable Data

Nonlinearly Separable Data

Linear classifier cannot classify these examples.



$$z = \theta^T x = \theta_0 + \theta_1 x + \theta_2 x^2 = 0 \Rightarrow$$

Nonlinear decision boundary (in the original feature space x)

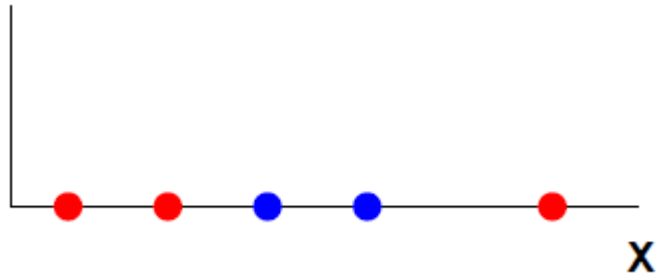
Linear decision boundary (in the extended feature space x, x^2)

if $z > 0 \Rightarrow g(z) > 0.5 \Rightarrow$ predict class = 1

if $z < 0 \Rightarrow g(z) < 0.5 \Rightarrow$ predict class = 0

Nonlinearly Separable Data

Linear classifier cannot classify these examples.



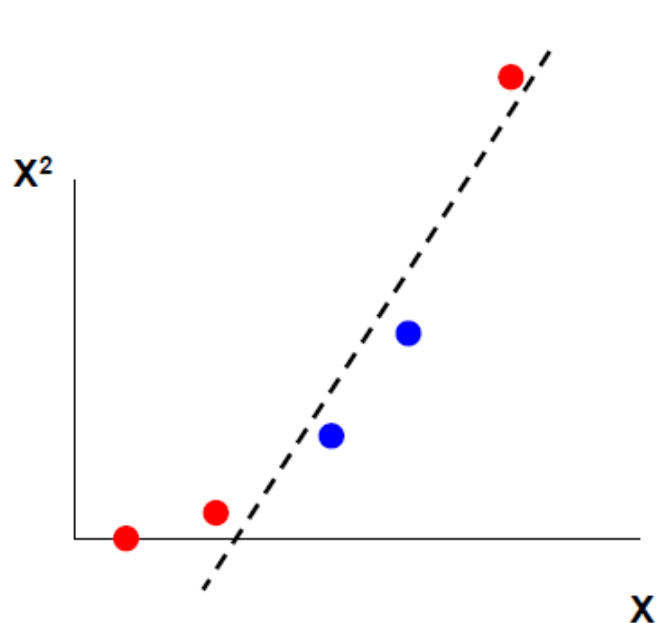
$$z = \theta^T x = \theta_0 + \theta_1 x + \theta_2 x^2 = 0 \Rightarrow$$

Nonlinear decision boundary (in the original feature space x)

Linear decision boundary (in the extended feature space x, x^2)

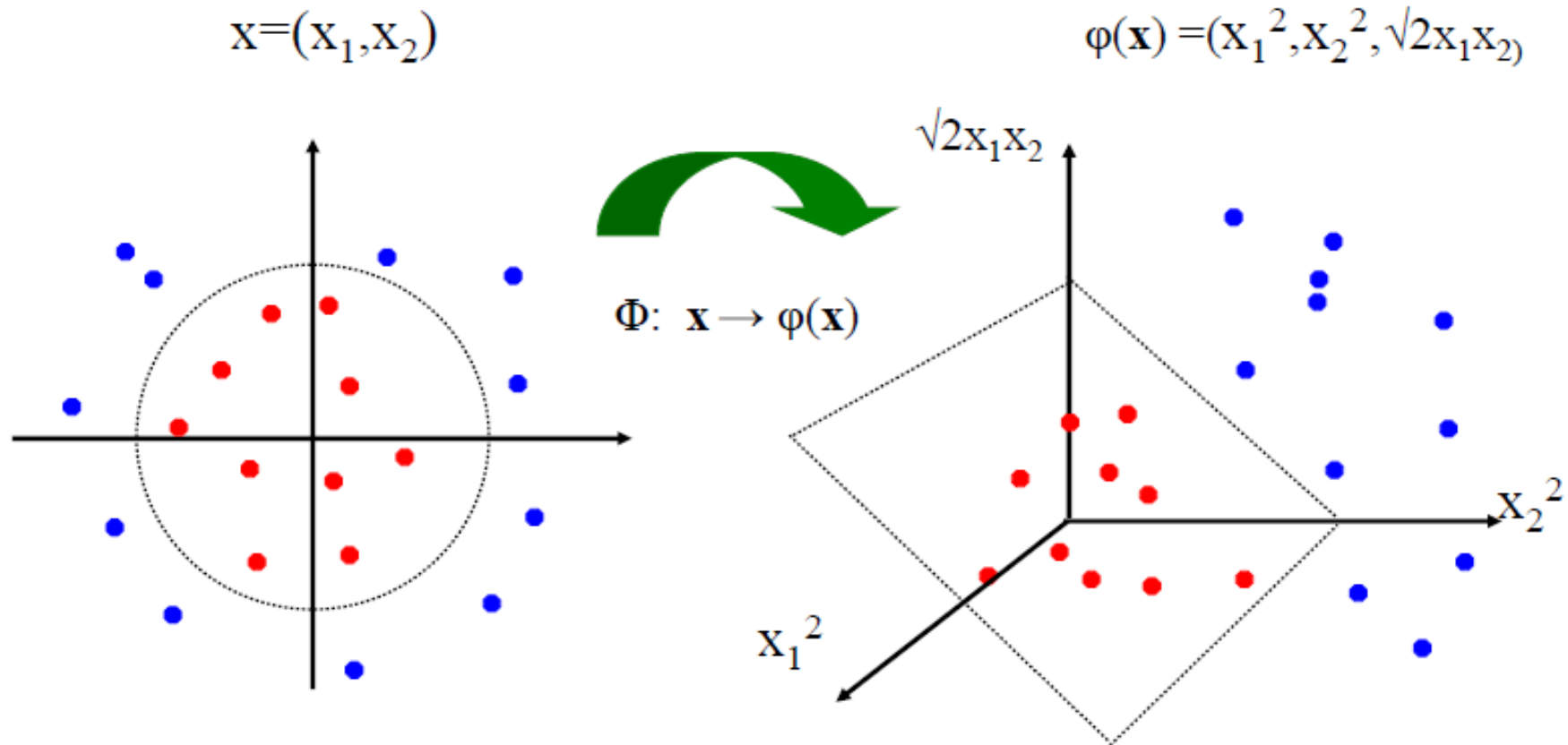
if $z > 0 \Rightarrow g(z) > 0.5 \Rightarrow$ predict class = 1

if $z < 0 \Rightarrow g(z) < 0.5 \Rightarrow$ predict class = 0



Nonlinearly Separable Data

- The original input space (\mathbf{x}) can be mapped to some higher-dimensional feature space ($\varphi(\mathbf{x})$) where the training set is separable:



This slide is courtesy of www.iro.umontreal.ca/~pift6080/documents/papers/svm_tutorial.ppt

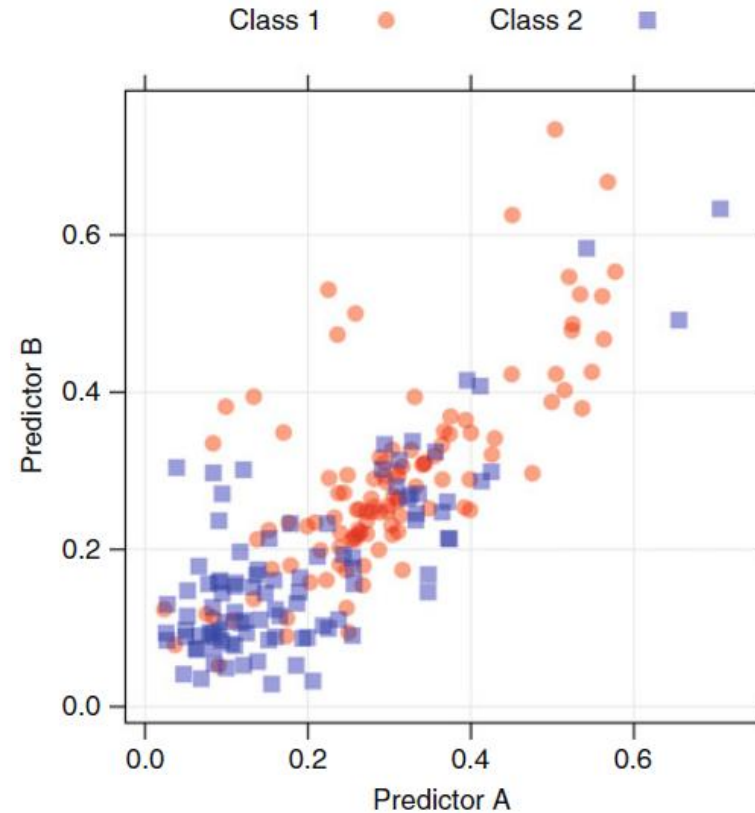
Overfitting Problem

Overfitting Problem

There now exist many techniques that can learn the structure of a set of data so **well** that when the model is applied to the data on which the model was built, it **correctly** predicts every sample.

In addition to learning the general patterns in the data, the model has also learned the characteristics of each sample's unique noise.

This type of model is said to be **over-fit** and will usually have **poor** accuracy when predicting a **new** sample.



Overfitting Problem

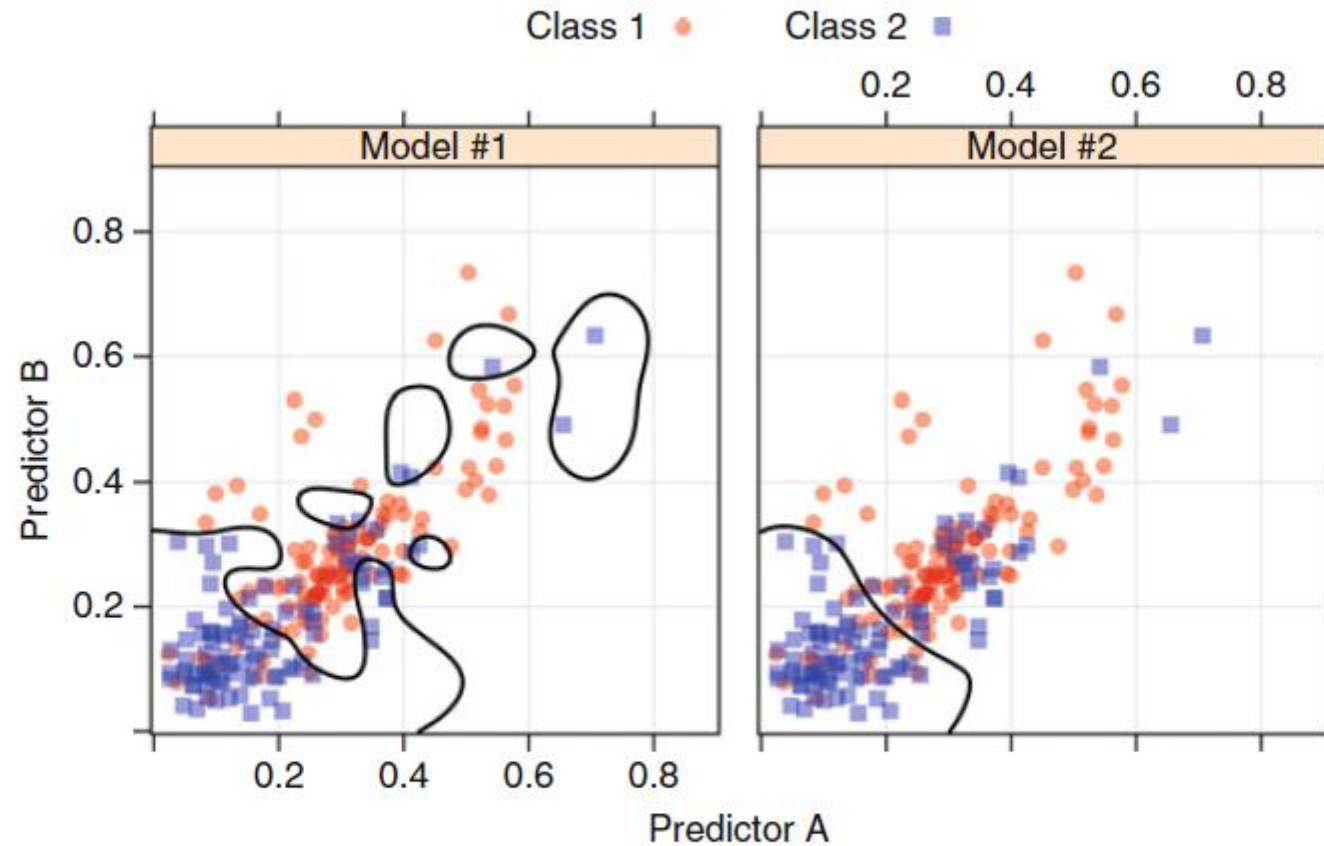


Fig. 4.2: An example of a training set with two classes and two predictors. The panels show two different classification models and their associated class boundaries

Overfitting Problem

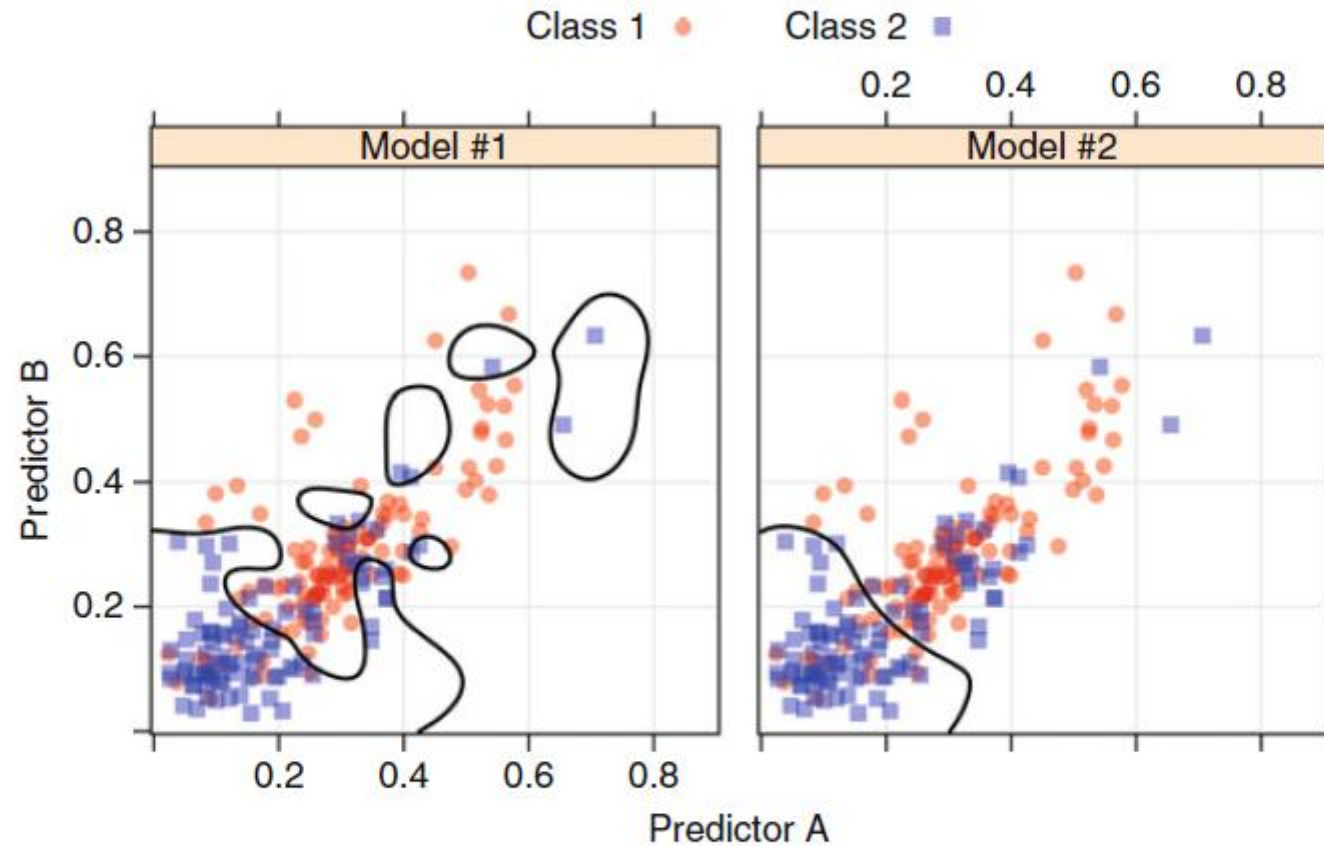
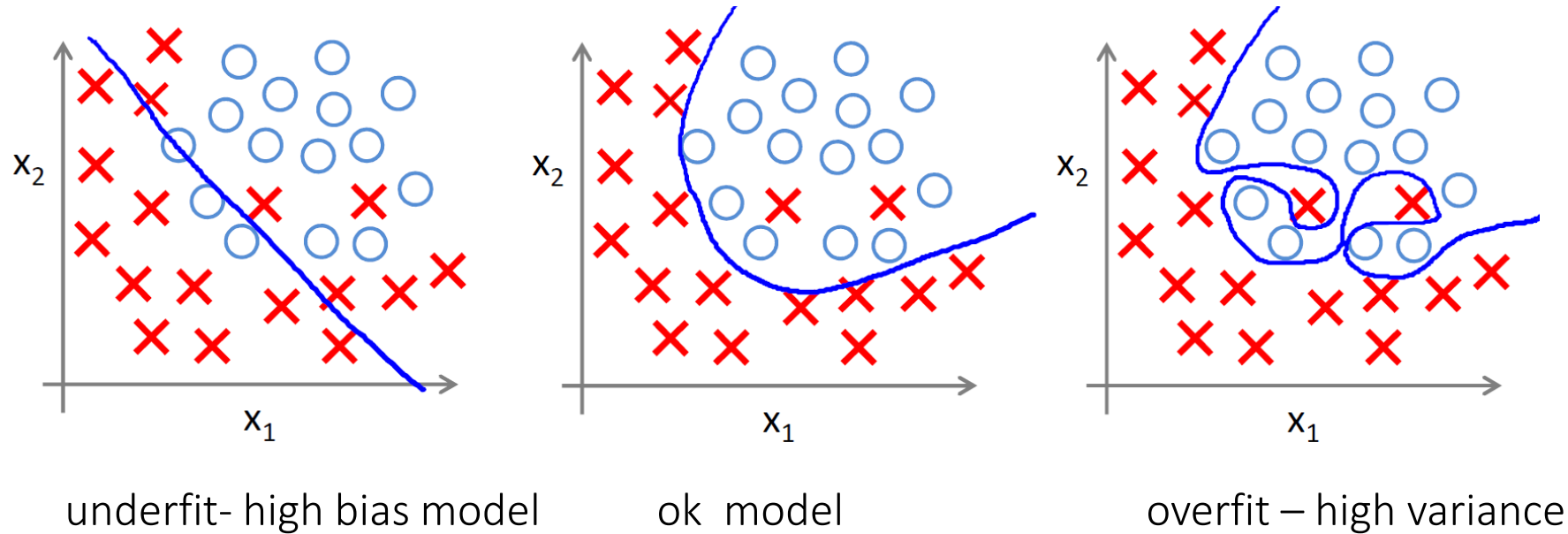


Fig. 4.2: An example of a training set with two classes and two predictors. The panels show two different classification models and their associated class boundaries

Overfitting Problem

Overfitting: If we have too many features, the learned model may fit the training data very well but fail to generalize to new examples.



A vertical bar on the left side of the slide, consisting of a wide red section and a thin blue section.

Regularized logistic regression

Regularization

Regularization is a popular method in ML to prevent overfitting by reducing the model parameters ϑ towards zero

Ridge Regression (L2 norm)

- Keep all the features, but reduces the magnitude of ϑ .
- Works well when each of the features contributes a bit to predict y .

Lasso Regression (L1 norm)

- May shrink some coefficients of ϑ to exactly zero.
- Serve as a feature selection tools (reduces the number of features).

Regularization

Regularization is a popular method in ML to prevent overfitting by reducing the model parameters ϑ towards zero

Ridge Regression (L2 norm)

- Keep all the features, but reduces the magnitude of ϑ .
- Works well when each of the features contributes a bit to predict y .

Lasso Regression (L1 norm)

- May shrink some coefficients of ϑ to exactly zero.
- Serve as a feature selection tools (reduces the number of features).

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n |\theta_j|$$

λ is the regularization parameter (hyper-parameter)

Small $\lambda \Rightarrow$ lower bias, higher variance

High $\lambda \Rightarrow$ higher bias, lower variance

Regularization

Unregularized cost function gradients
(for all parameters $j=0,1, \dots,n$)

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Regularized cost function gradient for
 $j=0$ (no change)

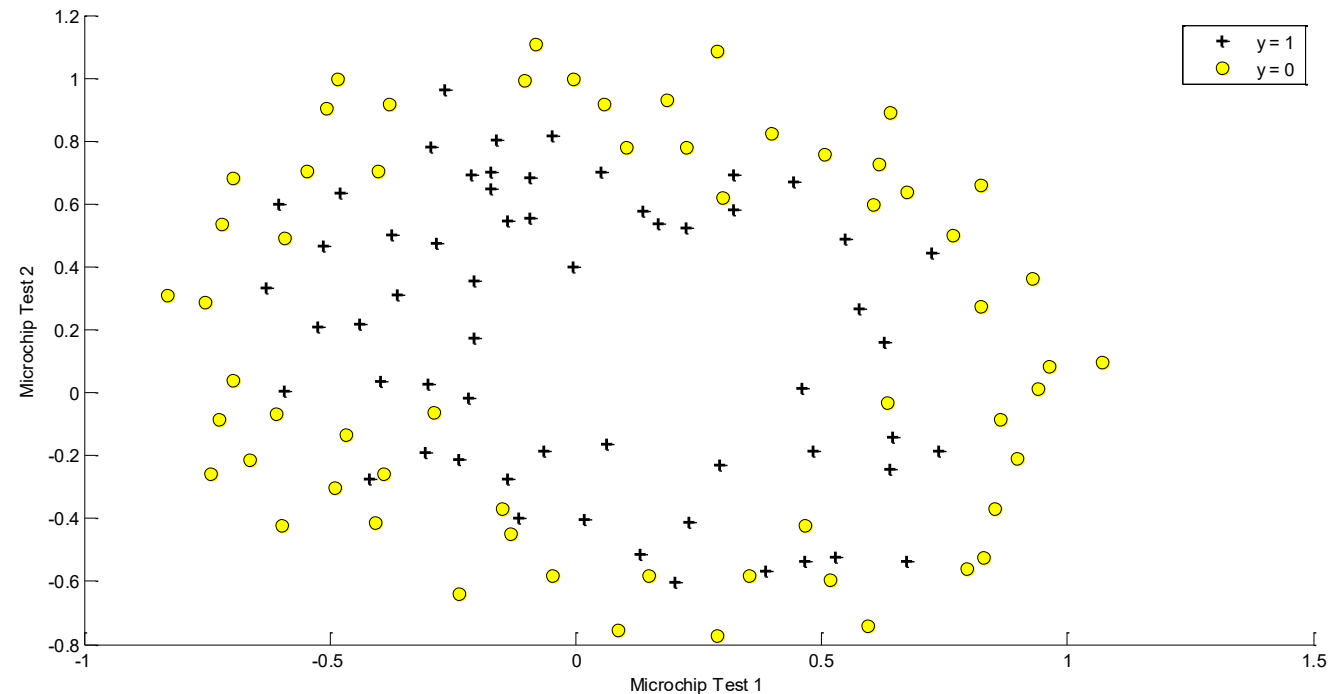
$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Regularized cost function gradients for
 $j=1,2,\dots,n$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j$$

Regularization - example

Predict whether microchips from a fabrication plant passes quality assurance (QA). During QA, each microchip goes through various tests to ensure it is functioning correctly. Suppose we have the test results for some microchips on two different tests. From these two tests, we would like to determine whether the microchips should be accepted ($y=1$) or rejected ($y=0$).



Regularization - example

Dataset is not linearly separable => logistic regression will only be able to find a linear decision boundary.
One way to fit the data better is to create more features. For example add polynomial terms of x_1 and x_2 .

$$\text{mapFeature}(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1x_2 \\ x_2^2 \\ x_1^3 \\ \vdots \\ x_1x_2^5 \\ x_2^6 \end{bmatrix}$$

NONLINEAR decision boundary \Rightarrow

$$z = \theta^T x = \theta_0 + \theta_1x_1 + \theta_1x_2 + \theta_3x_1^2 + \theta_4x_1x_2 + \dots\theta_{28}x_2^6 = 0$$

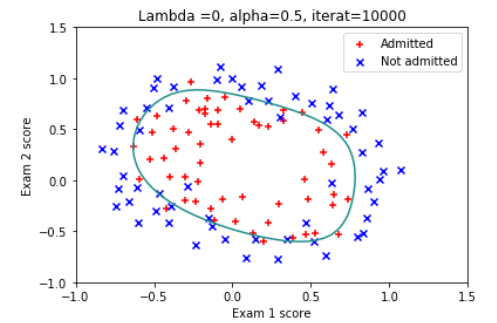
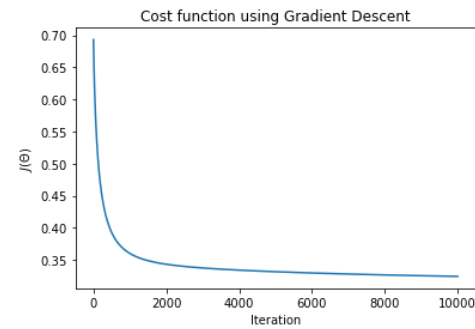
if $z > 0 \Rightarrow g(z) > 0.5 \Rightarrow$ predict class = 1

if $z < 0 \Rightarrow g(z) < 0.5 \Rightarrow$ predict class = 0

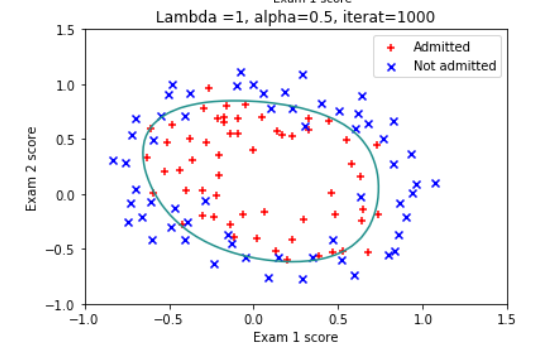
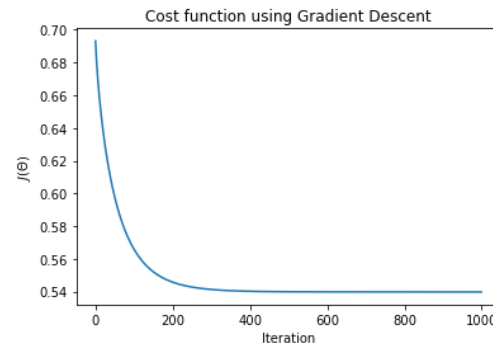
Regularization - example

Accuracy on training data:

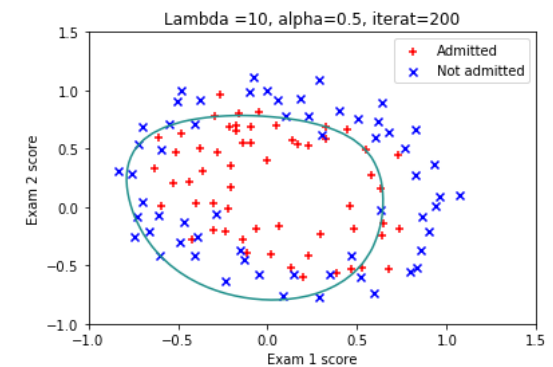
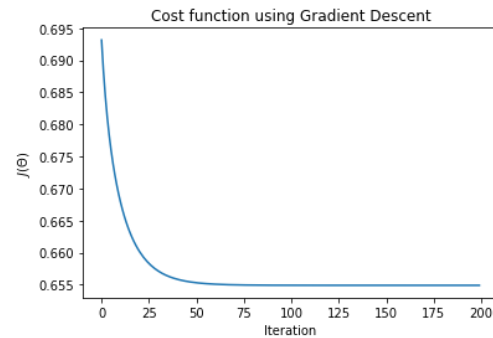
84.75% ($\lambda = 0$)



83.90% ($\lambda = 1$)



71.2% ($\lambda = 10$)



Regularization: Final Notes

Lasso and Ridge Regression are two of the most popular techniques for regularizing linear models, which often yield more accurate predictions than traditional linear models. These methods reduce the model's complexity by introducing shrinkage or adding a penalty to complex coefficients.

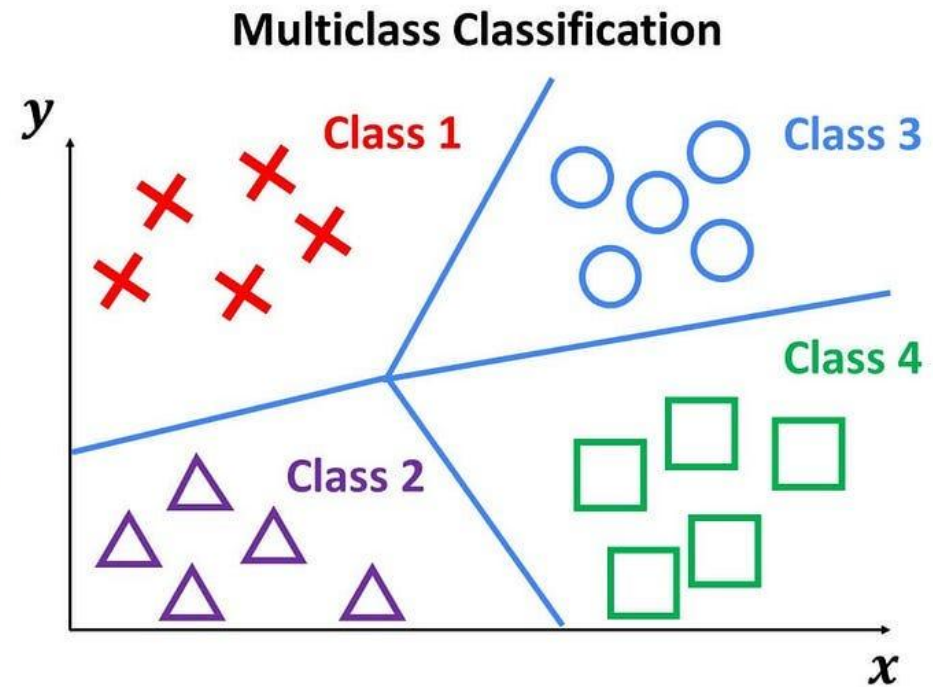
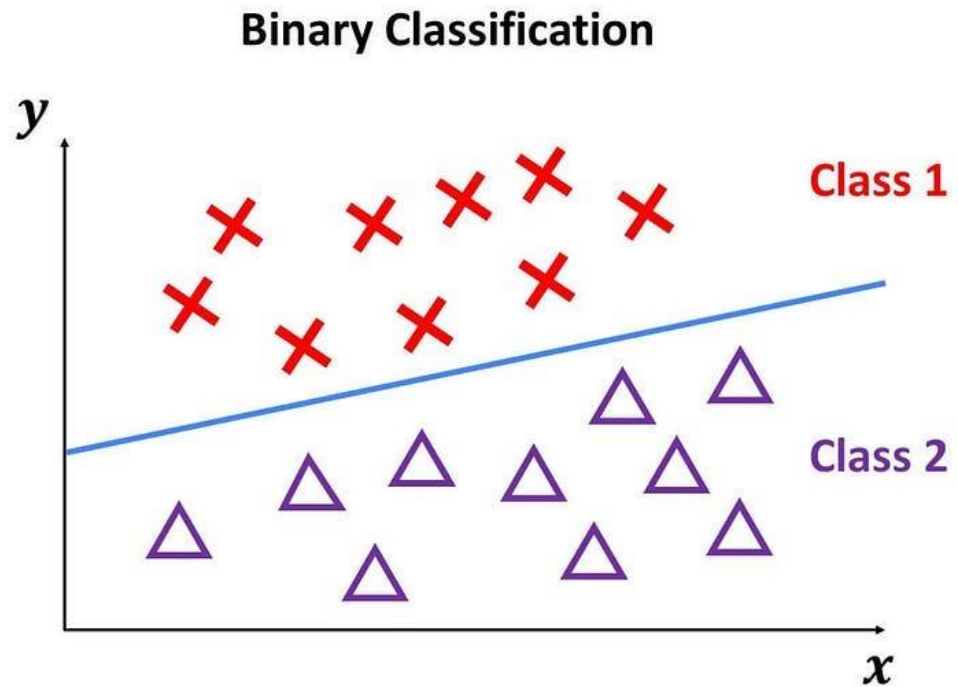
- The Ridge-Lasso approach is limited in requiring the input features to be standardized before fitting the model. It means that any feature with a large range of values can bias results because of its scale relative to other features with smaller ranges.
- Furthermore, if the data points contain outliers or noise, then this could produce inaccurate predictions due to the penalty terms.
- Additionally, Ridge and Lasso Regressions can be slow when applied to large datasets because of the computation time needed to perform regularization.
- Lastly, these methods require careful selection of hyperparameters (i.e., regularization strength), which can induce further computational costs and time.

Which scaling technique?

Multiclass

Multiclass Problem

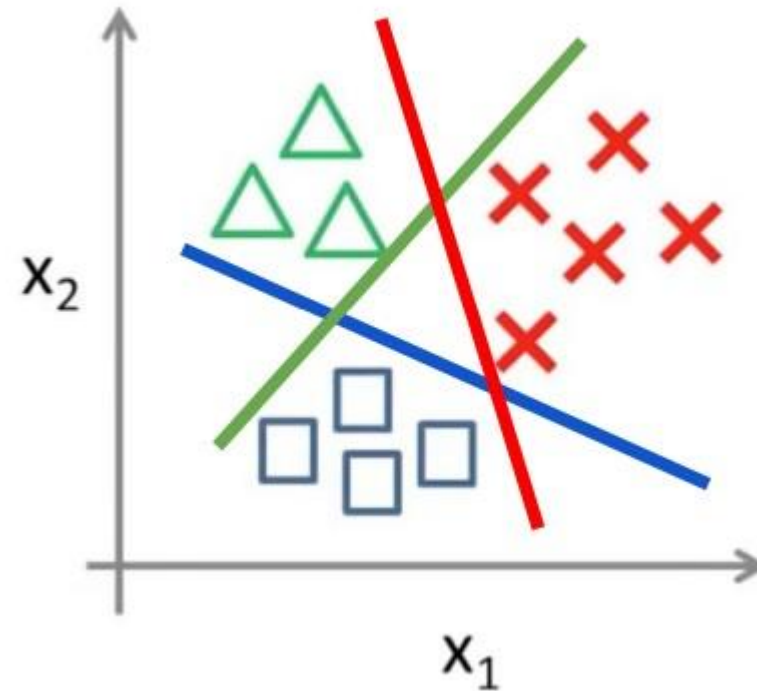
Logistic regression, by default, is limited to two-class classification problems.
How to deal with multiclass problem?



Multiclass Problem

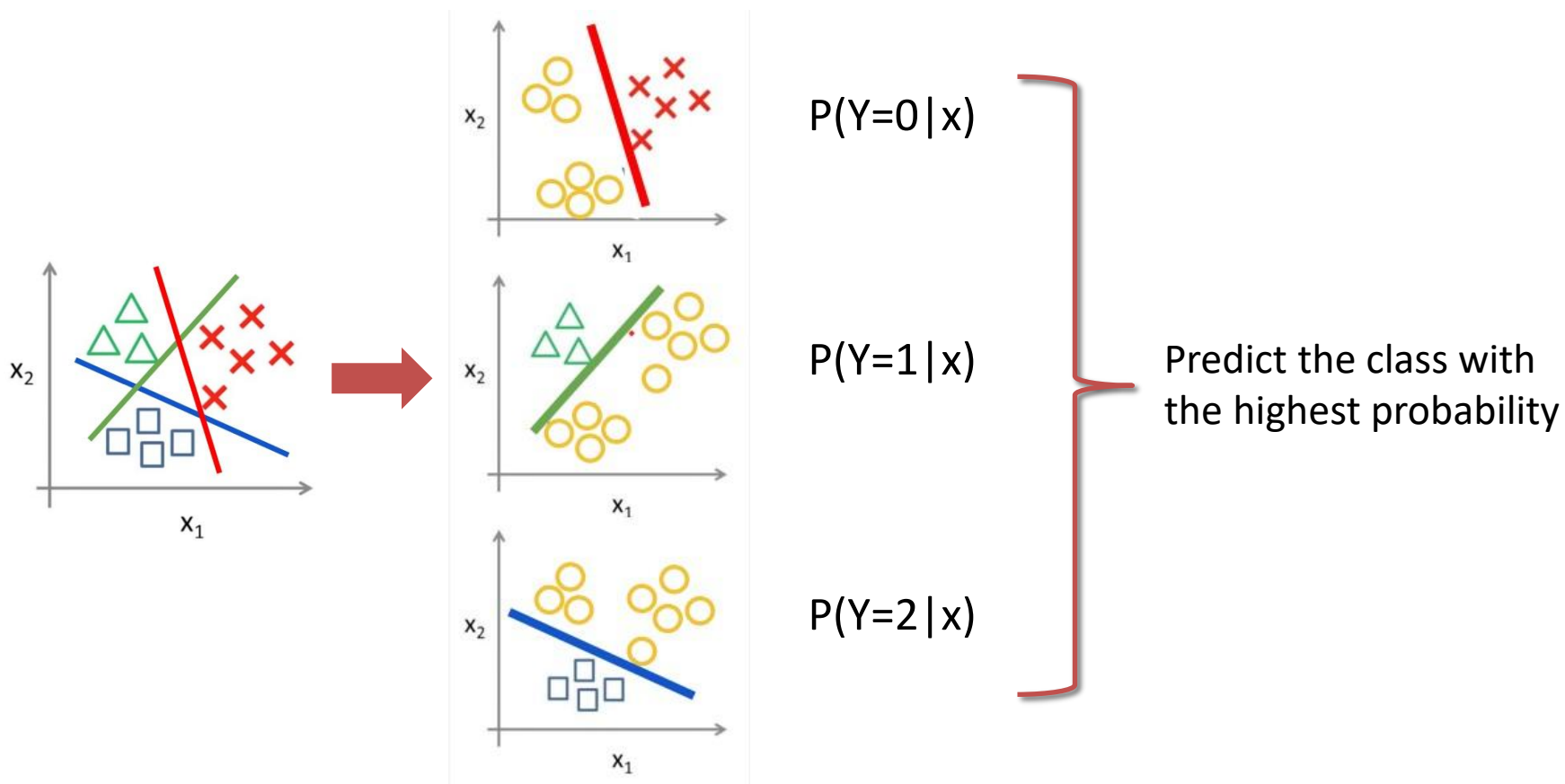
Possible solutions, and common approaches:

1. One vs All
2. Softmax Regression (Multimodal Logistic Regression)



Multiclass Problem – One vs All

1. For each class, define a logistic regression.
2. Calculate the probability of an observation to belong to that class.
3. Predict the class based on the highest probability.



Multiclass Problem – One vs All

For K classes train K binary classifiers:

for c=1:K

 Make $y_{\text{binary}}=1$ (only for examples of class c)

$y_{\text{binary}}=0$ (for examples of all other classes)

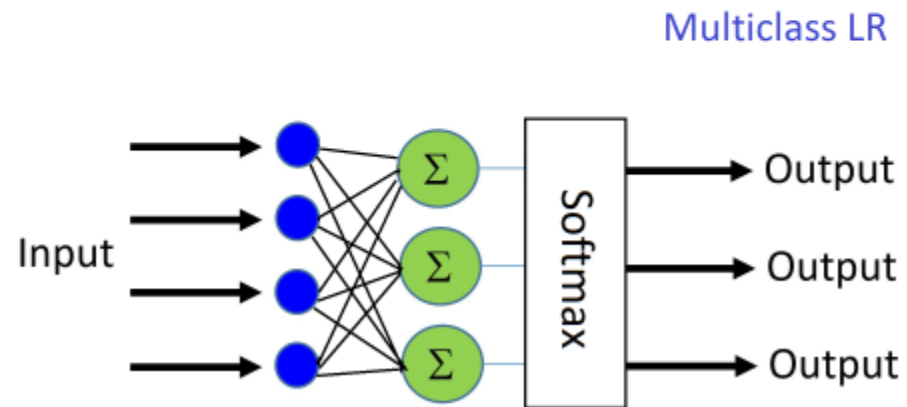
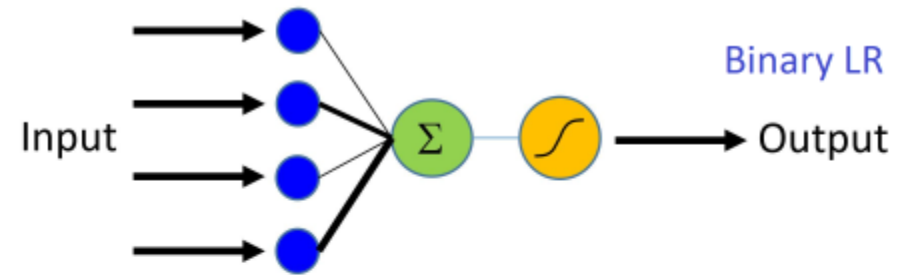
 theta=Train classifier with training data X and output y_{binary} .

 Save the learned parameters of all classifiers in one matrix, where each row is the learned parameters of one classifier:
 theta_all(c,:)=theta

End

New observation: winner-takes-all strategy, the binary classifier with the highest output score assigns the class.

Multiclass Problem – Softmax



Multiclass Problem – Softmax

In Softmax, the probability that an observation belongs to a class is determined by:

$$\begin{bmatrix} P(Y = 1|x; \theta) \\ P(Y = 2|x; \theta) \\ \vdots \\ P(Y = K|x; \theta) \end{bmatrix} = \frac{1}{\underbrace{\sum_{j=1}^K \exp(\theta^{(j)\top} x)}_{\text{Normalizes probabilities so they sum to 1.}}} \begin{bmatrix} \exp(\theta^{(1)\top} x) \\ \exp(\theta^{(2)\top} x) \\ \vdots \\ \exp(\theta^{(K)\top} x) \end{bmatrix} \rightarrow \text{Predict the class with the highest probability}$$

Separate $\theta^{(j)} \in R^d$ for each class

Special case K=2, Softmax is simplified by the binary logistic regression.

Classification – Logistic Regression

Pros:

- 1.Simplicity:** Easy to implement and understand.
- 2.Interpretability:** The weights directly show the importance of each feature.
- 3.Efficiency:** Doesn't require too much computational power.
- 4.Probabilistic Output:** Provides probabilities rather than just classifications.

Cons:

- 1.Linearity Assumption:** Assumes a linear relationship between features and log-odds of the outcome.
- 2.Feature Independence:** Assumes features are not highly correlated.
- 3.Limited Complexity:** May underfit in cases where the decision boundary is highly non-linear.
- 4.Requires More Data:** Needs a relatively large sample size for stable results.