

Logistic Regression – Assessment of Classification models

Dr. Mohamed Elshenawy

mmelshenawy@gmail.com



Previous sessions

- What is Learning?
- Regression - OLS
- Overfitting and underfitting
- Assessment of Regression Models
- Parametric and non-parametric models
- KNN

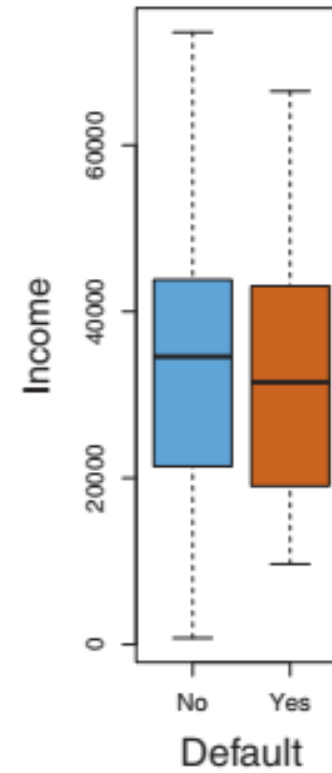
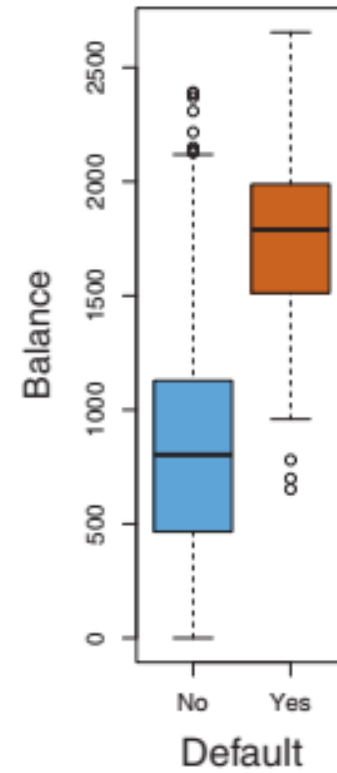
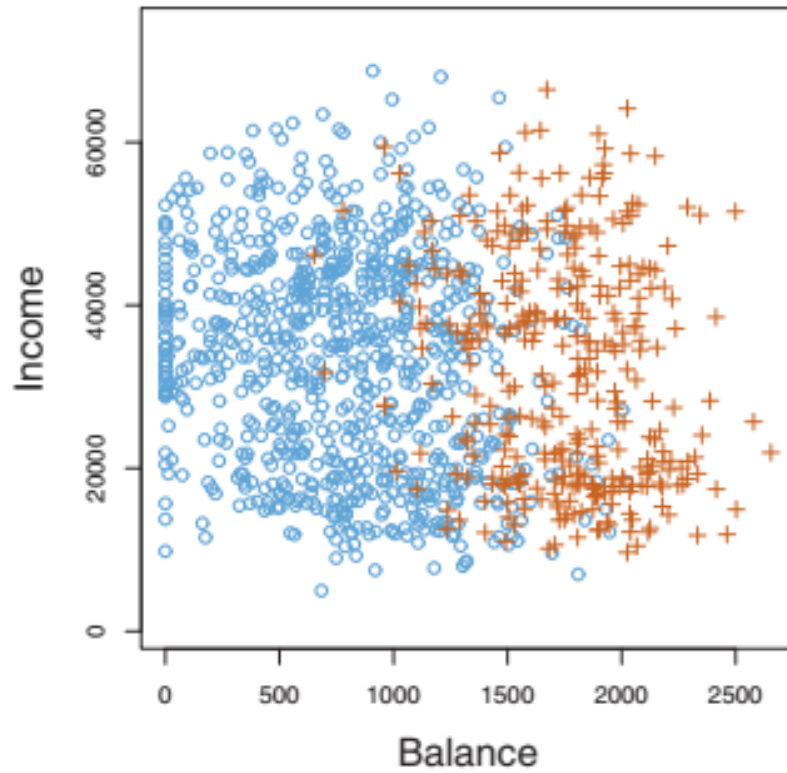


This session

- Review of Assignment 1 (Example solution)
- Logistic Regression (cont.)
 - Maximum Likelihood Estimation
 - Multi-class classification
- Assessment of Classification Models

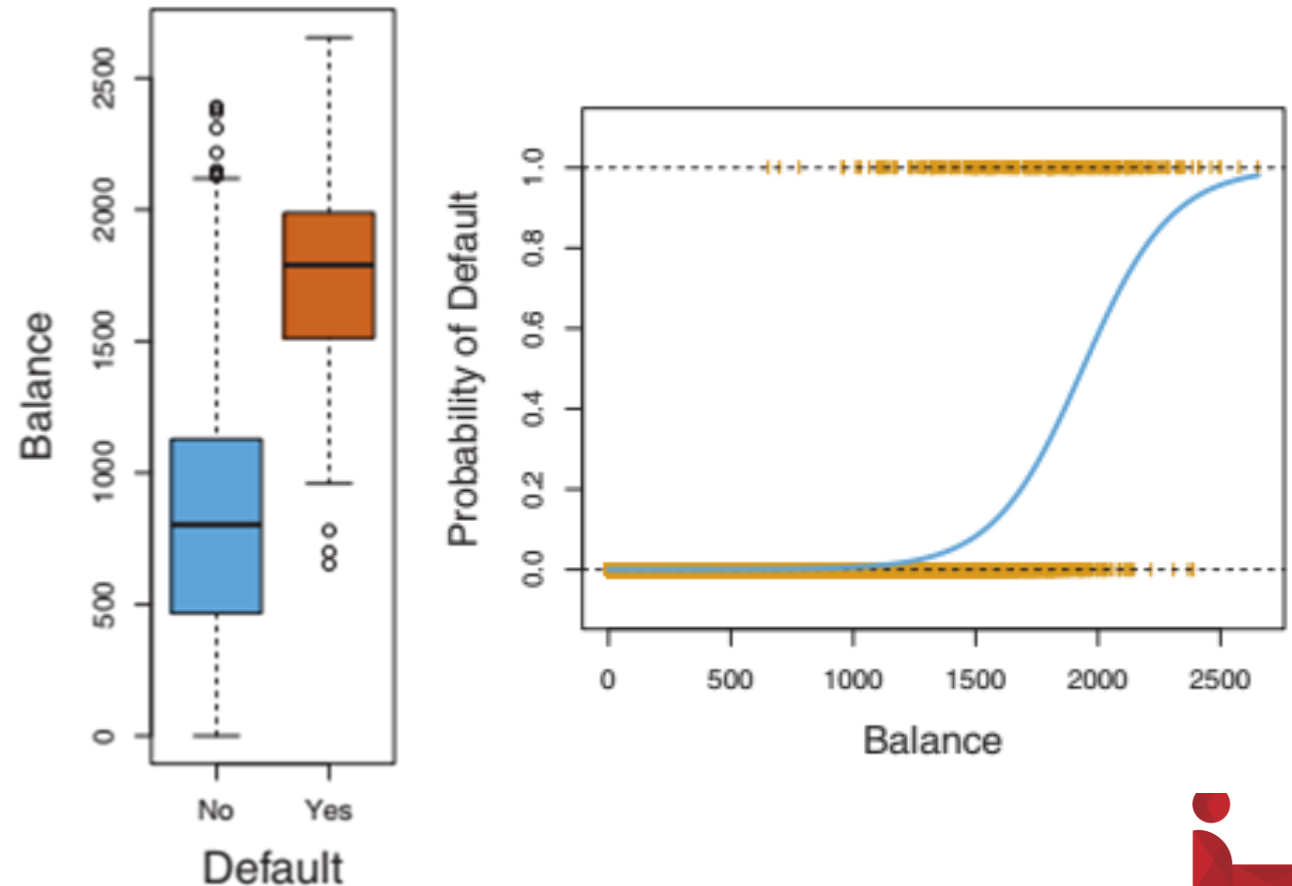


Previous session



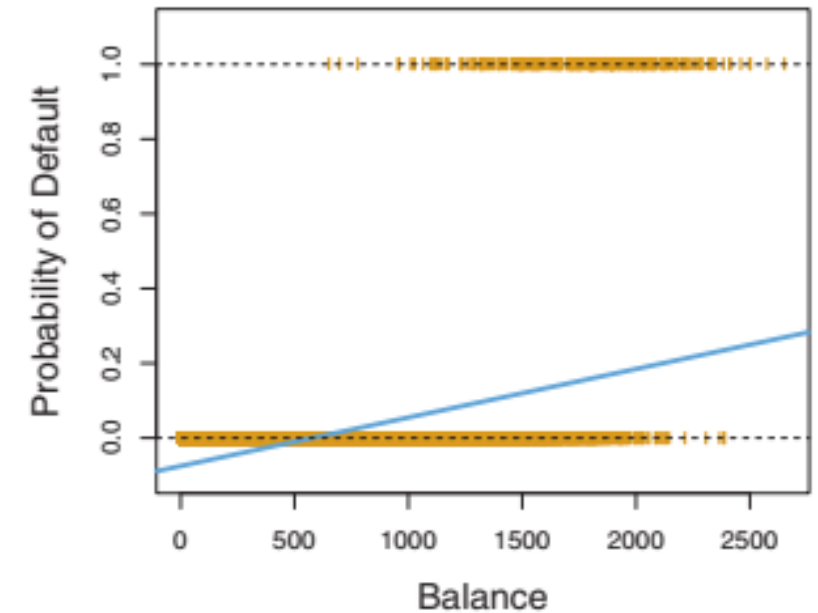
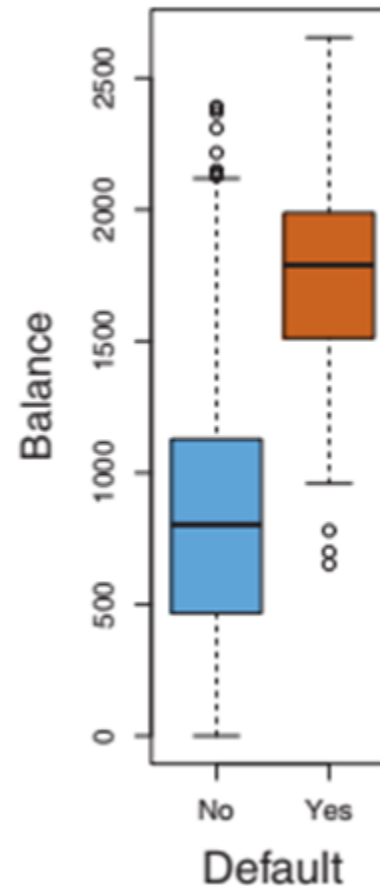
In the credit card payments example

- $\Pr(\text{default} = \text{Yes} | \text{balance}) = \sigma(\theta_0 + \theta_1 * \text{balance})$
- One might predict $\text{default} = \text{Yes}$ for any individual for whom $\Pr(\text{default} = \text{Yes} | \text{balance}) > 0.5$.
- Alternatively, a more conservative approach in predicting whether an individual will default (fail to repay) on his or her credit card payment $\Pr(\text{default} = \text{Yes} | \text{balance}) > 0.1$



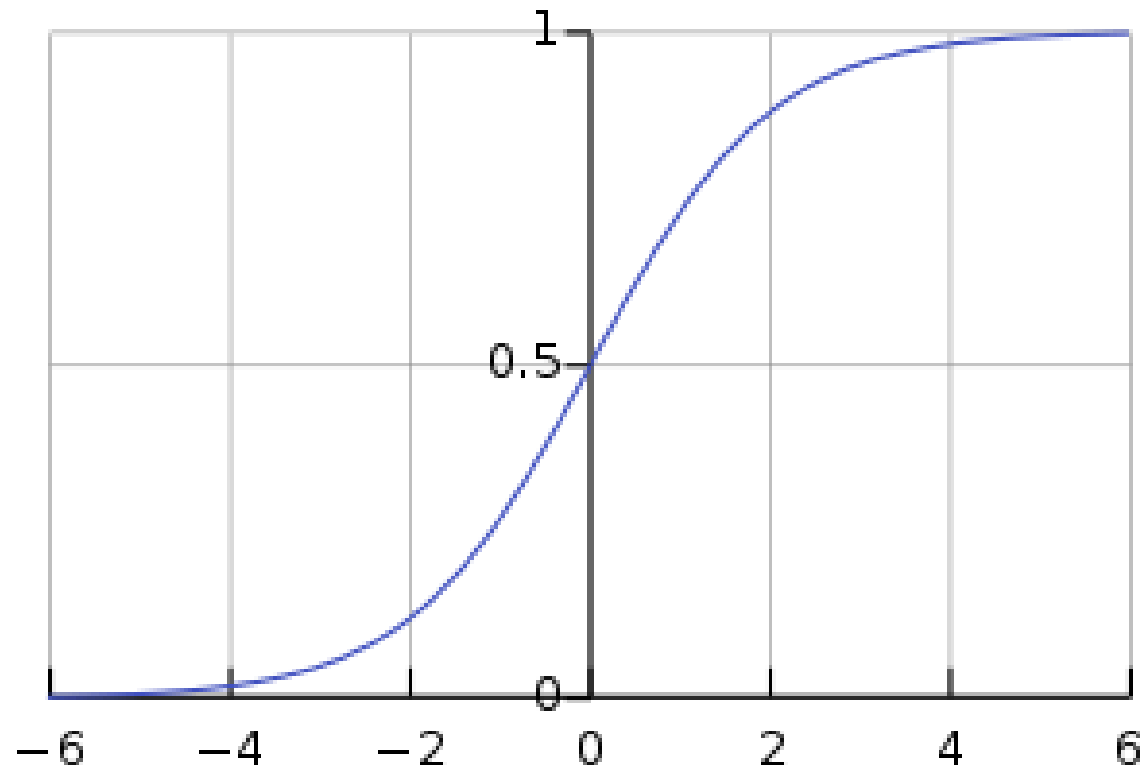
Assume you use a linear model (without a logistic function)

- $\Pr(\text{default} = \text{Yes} | \text{balance}) = \theta_0 + \theta_1 * \text{balance}$
- For balances close to zero we predict a negative probability of default; if we were to predict for very large balances, we would get values bigger than 1 (not sensible).



Logistic Function

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1}$$



How can we estimate the parameters?

- $\Pr(Y = 1|X) = \sigma(\theta_0 + \theta_1 * X_1 + \theta_2 * X_2 + \dots)$



How to estimate parameters: Maximum Likelihood Estimation

$$\theta_{ML} = \arg \max_{\theta} \prod_{i=1}^N P(y^{(i)} = t^{(i)} | x^{(i)}; \theta)$$

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^N \log P(y^{(i)} = t^{(i)} | x^{(i)}; \theta)$$

- The discussion of MLE is out of the scope of this session.



How to estimate parameters: Maximum Likelihood Estimation

$$\theta_{ML} = \arg \max_{\theta} \prod_{i=1}^N P(y^{(i)} = t^{(i)} | x^{(i)}; \theta)$$

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^N \log P(y^{(i)} = t^{(i)} | x^{(i)}; \theta)$$

$$gradient = \sum_{i=1}^N x_j^{(i)} (t^{(i)} - P(y^{(i)} = 1 | x^{(i)}; \theta))$$

- The discussion of MLE is out of the scope of this session.



Multiclass Logistic Regression

- How about if we have more than two classes?
- For example, a person arrives at the emergency room with a set of symptoms that could possibly be attributed to one of three medical conditions: stroke, drug overdose, epileptic seizure.
- In this setting, we wish to model:
 - $\Pr(Y = \textit{stroke}|X)$
 - $\Pr(Y = \textit{drug overdose}|X)$
 - $\Pr(Y = \textit{epileptic seizure}|X)$
- The sum of the three probabilities is 1
- The two-class logistic regression models have multiple-class extensions (not popular).
- Later in ML2, we will discuss a more popular method (discriminant analysis methods)



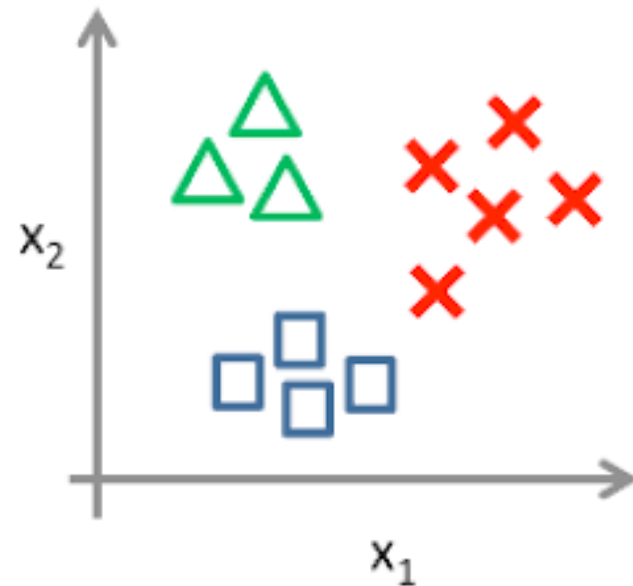
How can we separate multiple classes?

Create 3 training sets, 3 models

Predict the class using the three models, compare

Winner is the selected class

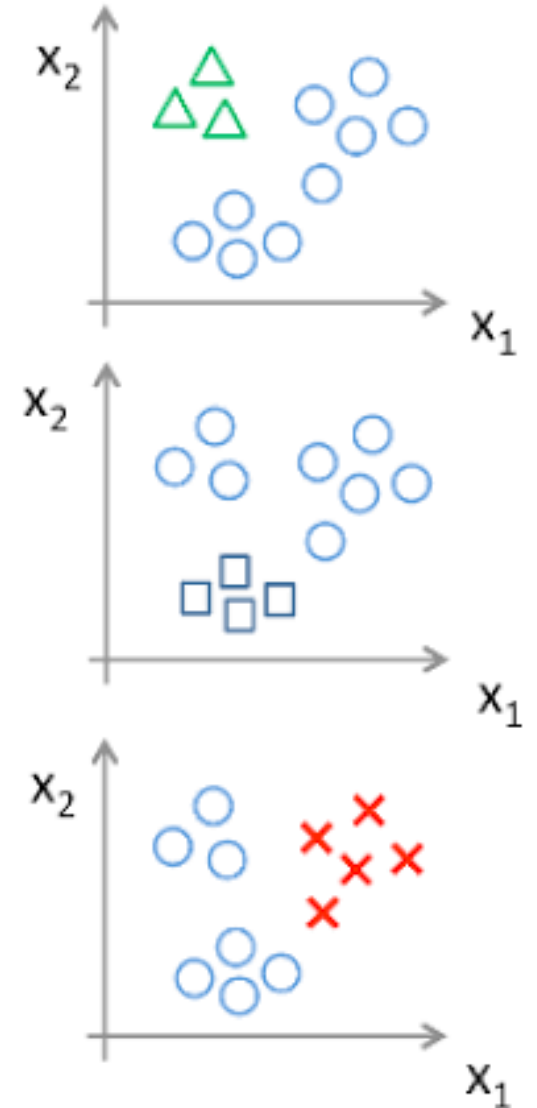
One-vs-all (one-vs-rest):



Class 1: Green

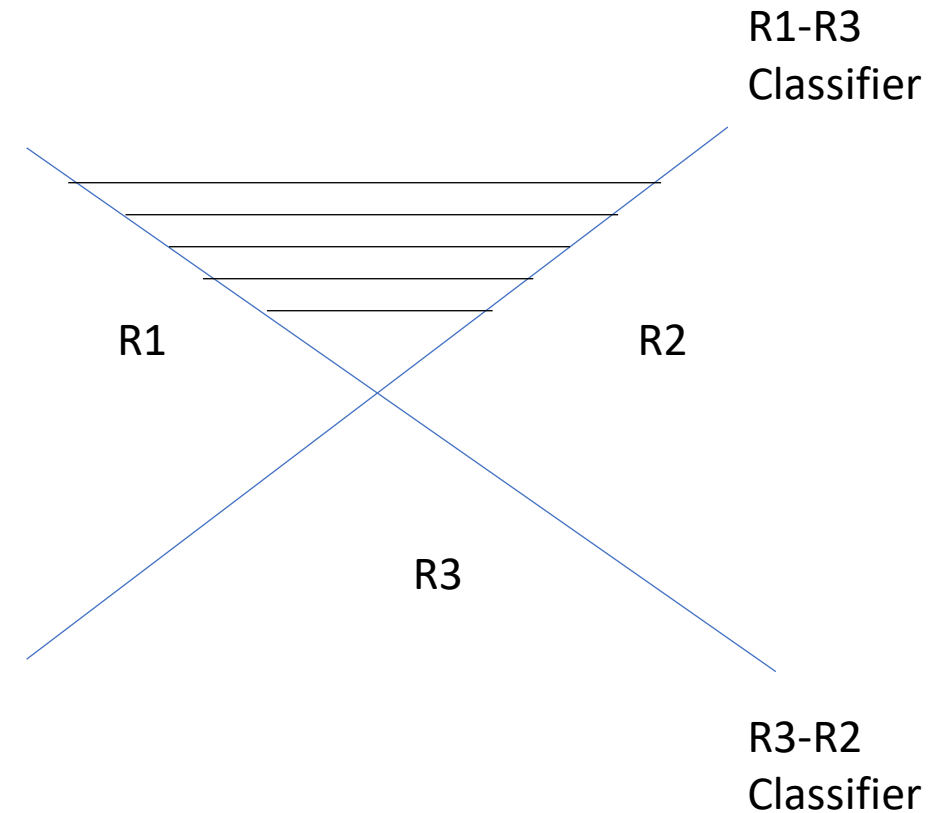
Class 2: Blue

Class 3: Red



Use K-1 classifiers??

- Given K classes where $K > 2$, how many lines do we need to separate these classes?
- An idea: we use $K-1$ classifiers, each separates two classes?
- How can classify the points in the shaded area?

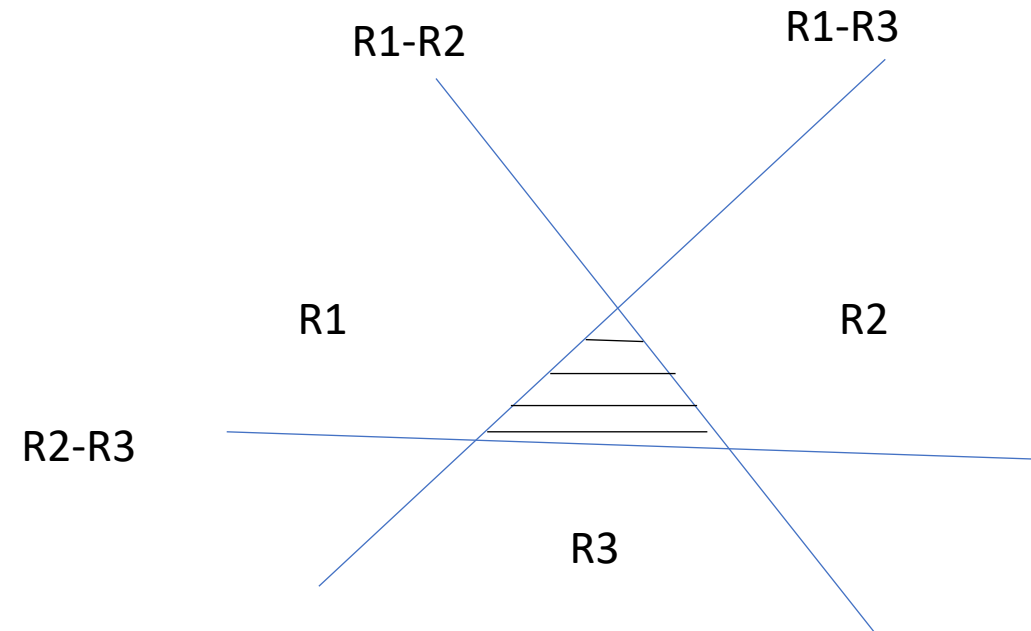


Use a classifier for each possible pair of classes (One-vs-One)

- Another idea: Use a classifier for each possible pair of classes. Fits $(K-1)$ classifier per class.
- Number of lines = number of possible pairs

$$\begin{aligned} \text{Number of lines} &= {}^k C_2 = \frac{k!}{(k-2)! 2!} \\ &= \frac{k(k-1)}{2} \end{aligned}$$

- How do you classify the points in the shaded area?



K-Class Discriminant

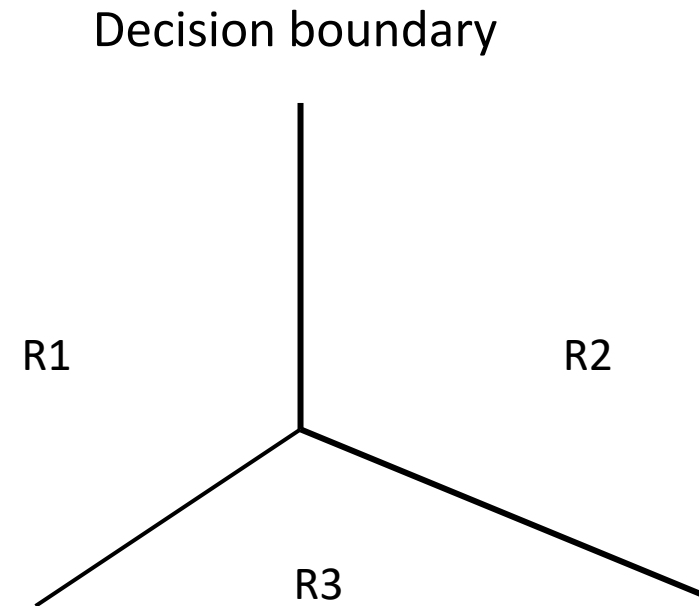
- For K classes, use K functions

$$z_k(x) = \theta_k^T x + \theta_{k0}$$

- Assign a point x to class k if
$$z_k(x) > z_j(x) \text{ for all } j \neq k$$
- The decision boundary between two classes k, j is given by

$$z_k(x) = z_j(x)$$

$$\text{that is } (\theta_k - \theta_j)^T x + (\theta_{k0} - \theta_{j0}) = 0$$



Representation of the output: One Hot Encoding

	<i>stroke</i>	<i>drug overdose</i>	<i>epileptic seizure</i>
<i>stroke</i>	1	0	0
<i>drug overdose</i>	0	1	0
<i>epileptic seizure</i>	0	0	1



Softmax function

- We use softmax: a normalized exponential function

$$p(Y = k|x) = \textit{softmax}(z_k) = \frac{e^{z_k}}{\sum_j e^{z_j}}$$

Where

$$z_k = \theta_k^T x + \theta_{k0}$$

Normalized exponential function guarantee that the sum of probabilities is equal to 1 (not the case if a sigmoid function is used)



Using Maximum Likelihood

- The likelihood

$$\prod_{n=1}^N \prod_{k=1}^K p(Y = k | x^{(n)})^{t_k^{(n)}}$$

Using summation

$$L(\theta, \dots, \theta_k) = - \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log p(Y = k | x^{(n)})$$

Use gradient descent

$$\frac{\partial E}{\partial w_{k,i}} = \sum_{n=1}^N ((t_k^n - p(Y = k | x^{(n)})) * x_i^{(n)})$$



Assessment of Classification Models

- We talked about accuracy: percentage of examples that are classified correctly by your model.
- Sometimes, accuracy is not sufficient to describe the model.
- Consider a diagnostic tool that uses a machine learning model to tell that person X has a disease Y.
- In this scenario, we are interested in two types of errors:
 - If X has the disease, what is the probability that the device say that X does not have it.
 - If X does not have this disease, what is the probability that the device will say that he has it.
- Consider an accident detection alarm
 - If the alarm is on, what is the probability that there's an accident
 - If there's an accident, what is the probability that the device will not detect it.



Confusion Matrix

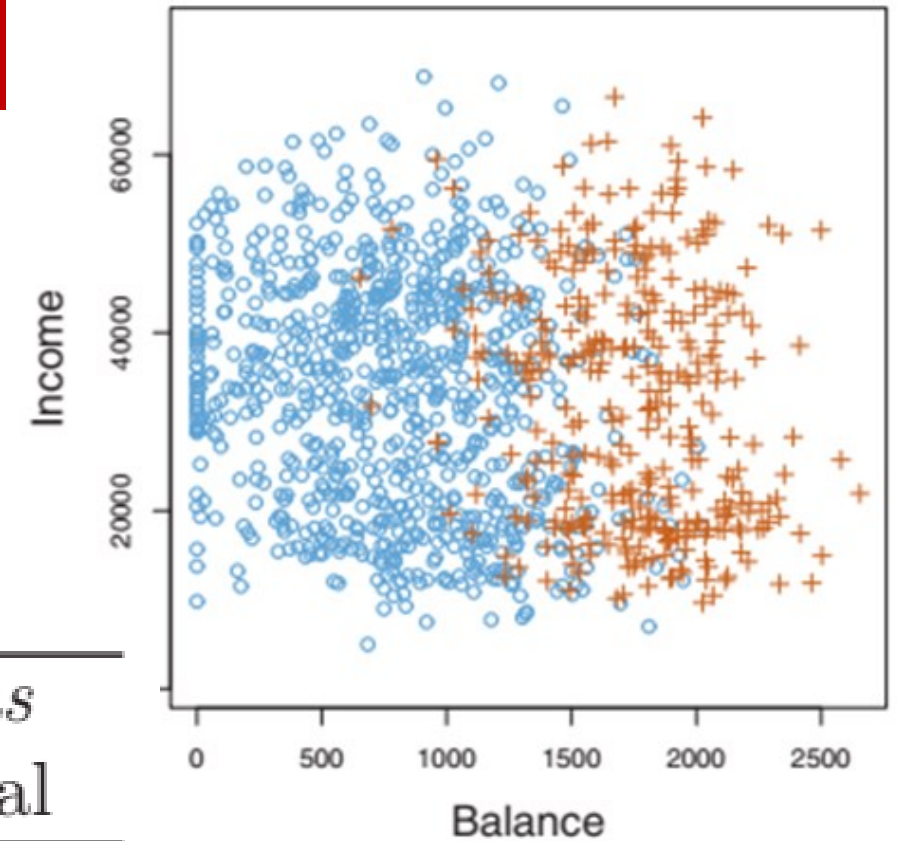
Output of the classifier

		P	N
P		True positive (TP)	False Negative (FN)
N		False positive (FP)	True Negative (TN)



Example

		<i>True default status</i>		
		No	Yes	Total
<i>Predicted default status</i>	No	9,644	252	9,896
	Yes	23	81	104
	Total	9,667	333	10,000



Recall, Precision, F1-Score

- Recall (other names sensitivity and probability of detection) :

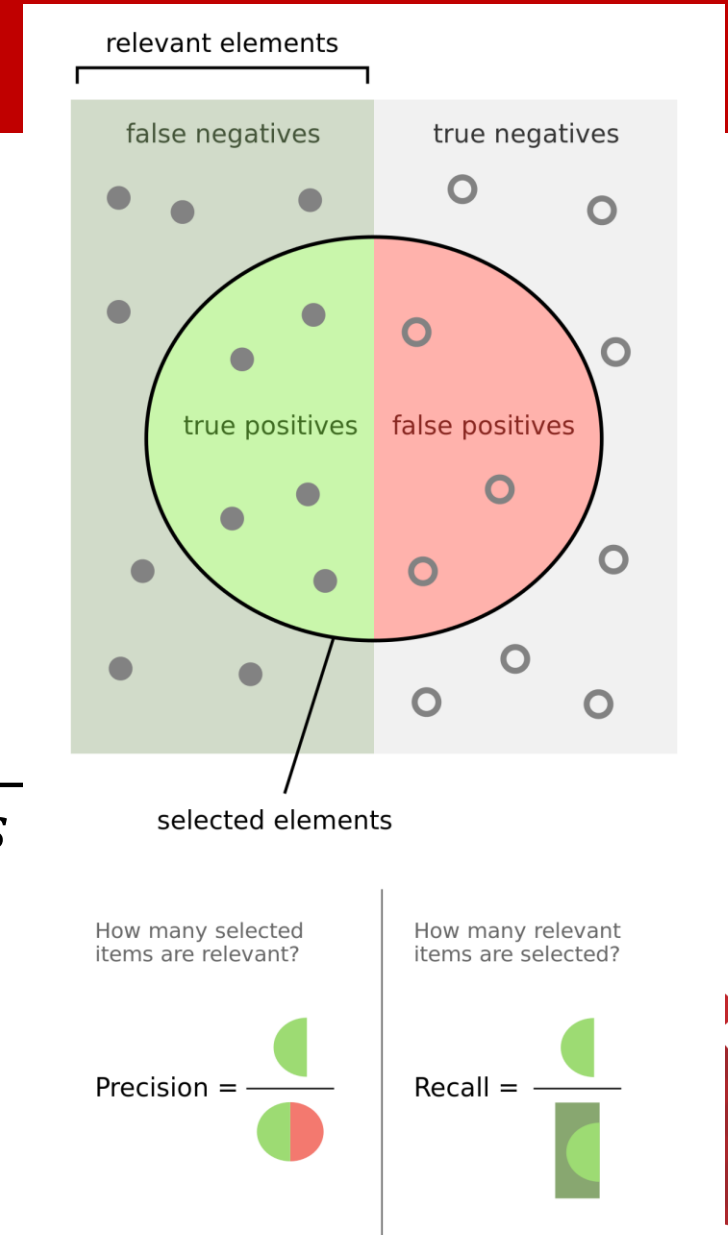
$$\text{Recall (R)} = \frac{TP}{TP + FN} = \frac{\text{Retrieved}}{\text{All groundtruth positives}}$$

- Precision

$$\text{Precision(P)} = \frac{TP}{TP + FP} = \frac{\text{Retrieved}}{\text{All Positive classifications}}$$

- F1 score: harmonic mean of precision and recall

$$F1 = 2 \frac{P * R}{P + R}$$



Sensitivity and Specificity

- Common metrics in biology
- Sensitivity (same as recall)

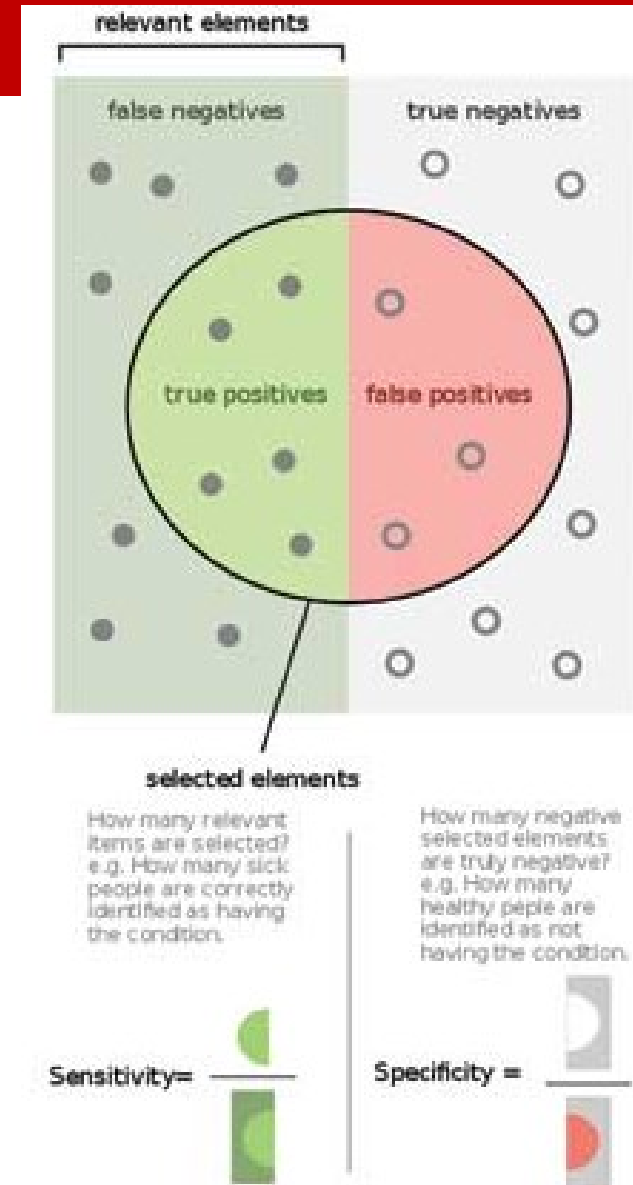
$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$= \frac{\text{Number of true positives}}{\text{Total number of cases with illness}}$$

- Specificity

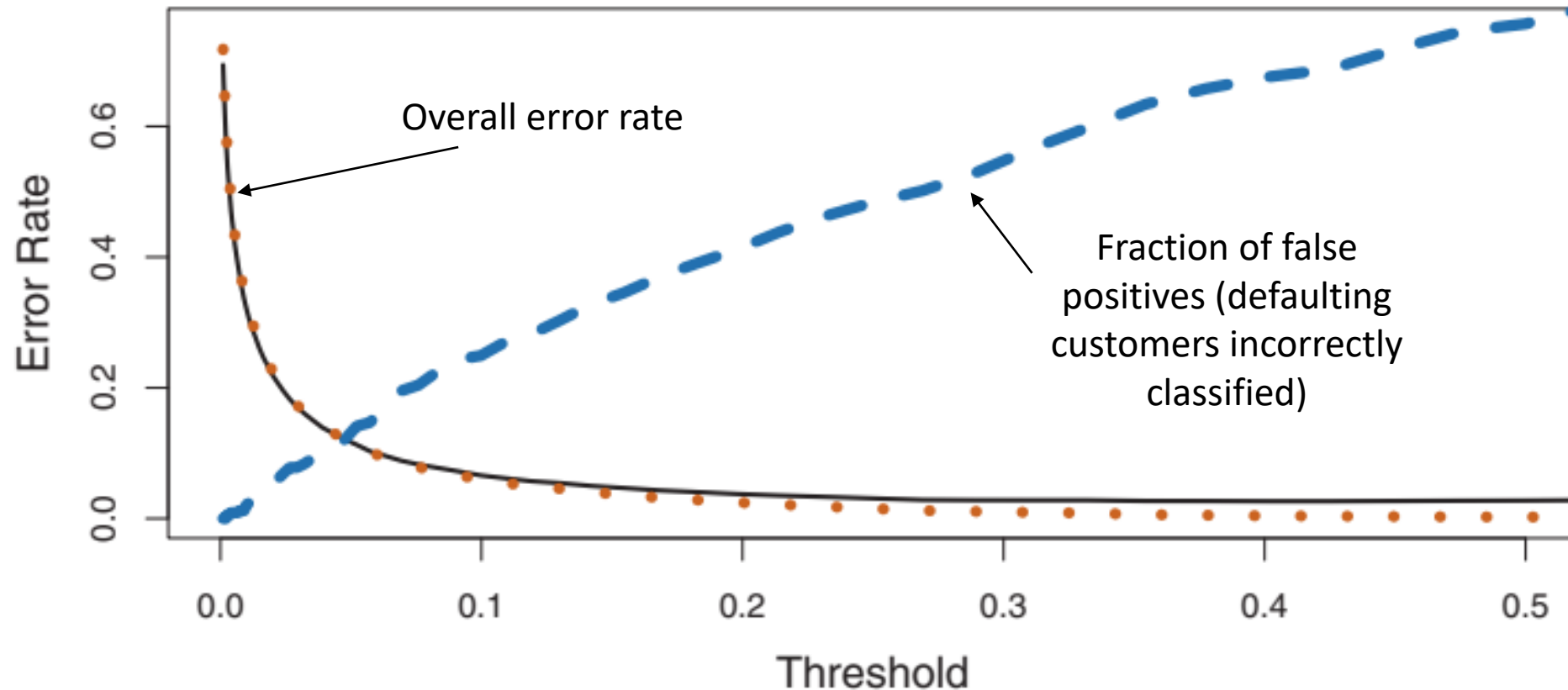
$$\text{Specificity} = \frac{TN}{FP + TN}$$

$$= \frac{\text{Number of true negatives}}{\text{Total number of cases without the illness}}$$



Choosing different thresholds

- Can we use thresholds to reduce the number of false positive?
- If $P(Y=1 | X=x)$ greater than 0.5, does it improve the performance



ROC Curve

- The ROC curve is a popular graphic for simultaneously displaying false positives and true positives for all possible thresholds.
- The name “ROC” comes from communications theory. It is an acronym for receiver operating characteristics.
- The overall performance of a classifier, summarized over all possible thresholds, is given by the area under the (ROC) curve (AUC).
- An ideal ROC curve will hug the top left corner, so the larger area under the AUC the better the classifier.

