# Multi-Class Classification
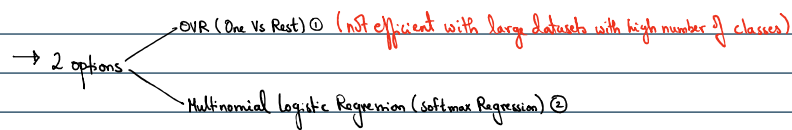
✳ No concept of threshold setting in multi-class classification

→ More than 2 classes

→ 2 options
- OVR (One Vs Rest) ① (not efficient with large datasets with high number of classes)
- Multinomial Logistic Regression (softmax Regression) ④

① internal working of algorithm   [ Multi-class problem → Multiple Binary Classification problems ]

training  [ each class gets its dedicated model ]

→ K number of logistic Regression models will be trained.
where K is the number of classes present in the target column.

1) Data is transformed by one hot encoding the target column.
2) Data is split up into K parts → we have converted a multi-class classification problem
   into K binary class classification problems.

3) Logistic Regression is applied on the K datasets independently

4) For each model we will obtain a corresponding $\beta_0, \beta_1, \beta_2$

prediction

1) test query point

2) send the test query point to each of the K models from which we
   will obtain K probabilities.

3) Normalize the output → $P(class_K) = \dfrac{P(model_K)}{\sum\limits_{i=1}^{K} P(model_i)}$   sum of normalized probabilities is 1.

4) choose the class with the highest normalized probability.

② internal working of algorithm

→ SoftMax function

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

→ Provides a set of probabilities that sum up to one.

→ For K classes it will try to draw K lines
   to create the decision regions.

→ each line will have (m+1) parameters, where m is the # of features.

→ In softmax regression, we will try to find (m+1) parameters
   for each of the K lines. So total of K(m+1) parameters.



training

1) Data is transformed by one hot encoding the target column.

2) Loss function = $-\dfrac{1}{n} \sum\limits_{i=1}^{n} \sum\limits_{k=1}^{K} y_i^k \log \hat{y}_i^k$   ← K represents the class [don't confuse with power]   ← minimize
   (categorical loss entropy)
   └ for each row   for each category
   → For each row we will obtain 1 term

3) $\hat{y}_i^{class} = \sigma(\vec{z}_i) = \dfrac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$
   └ row 1

where $z_i = \begin{bmatrix} 1 & x_{m1} & x_{m2} & \cdots & x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \end{bmatrix}$

row 1          j=1

Where $Z_i = [1 \ x_{n_1} \ x_{n_2} \ \ldots x_{nm}]$ $\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}$

↑ single record in data

corresponding line parameters to a particular class.

4) the Loss function depends on $\beta_{k \times (m+1)}$

Start with k random lines

Apply gradient descent [ $k(m+1)$ differentiations]

↑ for each parameter of each line → for ___ epochs

$$\beta_n^k = \beta_n^k - \eta \frac{\partial L}{\partial \beta_n^k}$$

↑ parameter #

prediction

1) we have found out the $\beta$ matrix for which loss is minimum

2) test query point

$$\beta = \begin{bmatrix} \beta_0^1 & \beta_1^1 & \ldots & \beta_m^1 \\ \beta_0^2 & \beta_1^2 & \ldots & \beta_m^2 \\ \vdots & \vdots & & \vdots \\ \beta_0^k & \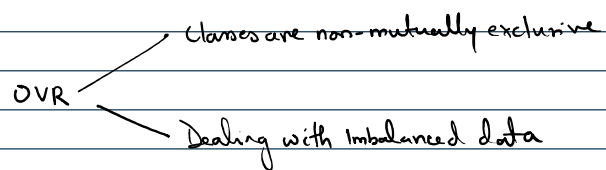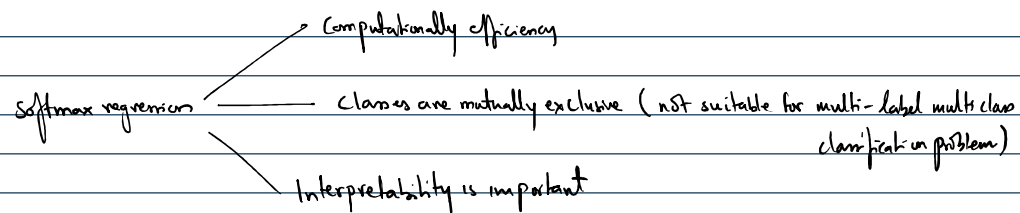beta_1^k & \ldots & \beta_m^k \end{bmatrix}_{k \times (m+1)} \qquad Q = \begin{bmatrix} 1 & f_1 & f_2 & f_3 & \ldots & f_m \end{bmatrix}_{1 \times (m+1)}$$

3) dot product $\beta$ & Q

4) send each value to softmax function

5) choose class with highest probability

Usage

Softmax regression
- Computationally efficiency
- Classes are mutually exclusive ( not suitable for multi-label multi class classification problem)
- Interpretability is important

OVR
- Classes are non-mutually exclusive
- Dealing with Imbalanced data

→ Deriving sigmoid from softmax
→ Deriving binary cross entropy from categorical cross entropy
→ Find derivative of softmax function
→ Find the gradient of cross entropy error.