

Logistic Regression

Machine Learning (CS 306)

Instructor:

Dr. Moumita Roy

Teaching Assistants: Indrajit Kalita, Veronica Naosekpam

Email-ids:

moumita@iiitg.ac.in

veronica.naosekpam@iiitg.ac.in

indrajit.kalita@iiitg.ac.in

Mobile No: +91-8420489325 (only for emergency quires)

Why use logistic regression?

- There are many important research topics for which the dependent variable is "limited."
- For example: voting, spam or not spam, and participation data is not continuous.
- Binary logistic regression is a type of regression analysis where the dependent variable is a dummy variable: coded as 0 (did not vote) or 1(did vote)

Classification

- Email: Spam/not spam
- Land-cover: Water-cover/not water-cover
- Customer Behavior Prediction: Sad/Happy
- Tumor: Malignant/not Malignant

$$Y \in \{0, 1\}$$

- coded as binary dependent variable (two-class problem)
0 → Negative class
1 → Positive class

Linear Regression (Revisit)

$$x_i = [x_{i,1} \quad x_{i,2} \quad x_{i,3} \quad x_{i,4} \quad \dots \quad x_{i,n}]$$

$$h(x_i) = \sum_{j=0}^n w_j x_{i,j} \quad x_{i,0} = 1$$

$$w_j = [w_0 \quad w_1 \quad w_2 \quad \dots \quad w_n]$$

$$h(x_i) = w^T x_i$$

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

↓
Cost function

$$w_j := w_j - \alpha \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i) x_{i,j}$$

↓
parameter updating using GD

Important Points

- Can we use solution of linear regression problem directly for classification?

- It may be formulated like this (threshold the output):

$h(X_i) \geq 0.5$, predict Y_i as 1

$h(X_i) < 0.5$, predict Y_i as 0

- Problem in the present approach:

I. The output $h(X_i)$ is unbounded

$h(X_i) > 1$ or < 0

II. In this case, fixing threshold at 0.5, is reasonable or not?

Solution: search a function which can bound the output between 0 to 1

Logistic Regression (Classification)

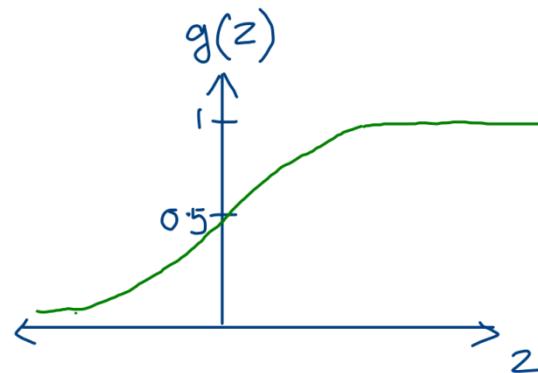
- Aim is to range the output: $0 \leq h(X_i) \leq 1$

$$h(x_i) = g(w^T x_i)$$

- $g(\dots)$ as sigmoid function/logistic function

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

$$h(x_i) = \frac{1}{1 + e^{-w^T x_i}}$$



$$g(z) \rightarrow 0 \text{ if } z \rightarrow -\infty$$

$$g(z) \rightarrow 1 \text{ if } z \rightarrow \infty$$

Interpretation of Hypothesis output

- $h(X_i)$ =estimated probability that $Y_i=1$ for a given X_i
- Suppose, $h(X_i)=0.7$

interpret as, there is a 70% chance that X_i belongs to $Y_i=1$; alternatively there is a 30% chance that X_i belongs to $Y_i=0$ (if estimate in the scale of 100)

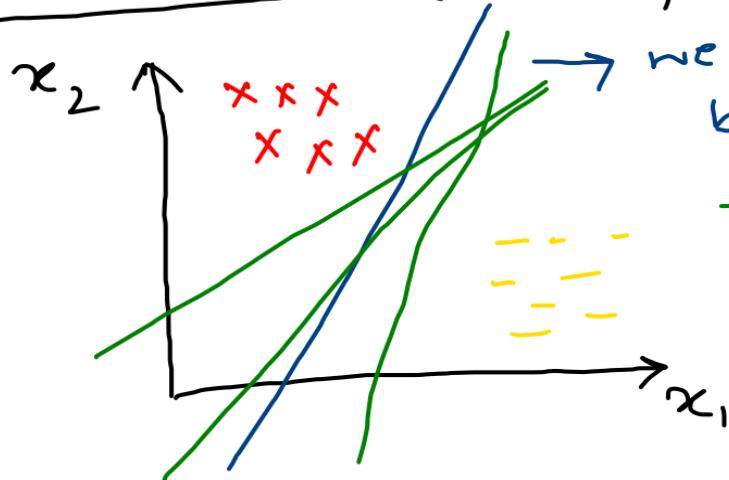
- $h_w(X)=P(Y=1/X; W)$ [probability that $Y=1$ for given X ; parameterized by W]
- Since, Y is either 0 or 1

$$P(Y=1/X; W) + P(Y=0/X; W) = 1$$

$$P(Y=0/X; W) = 1 - P(Y=1/X; W)$$

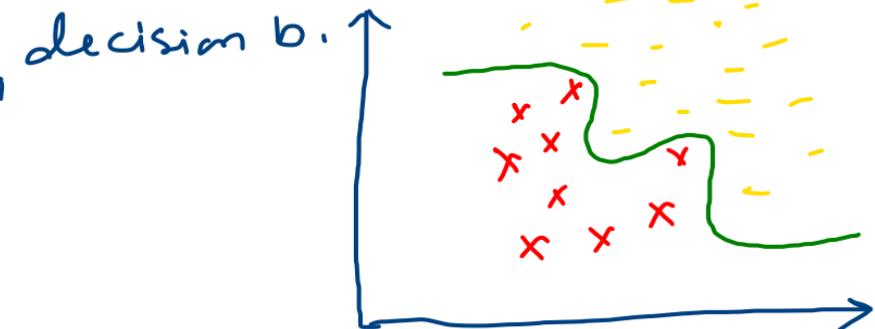
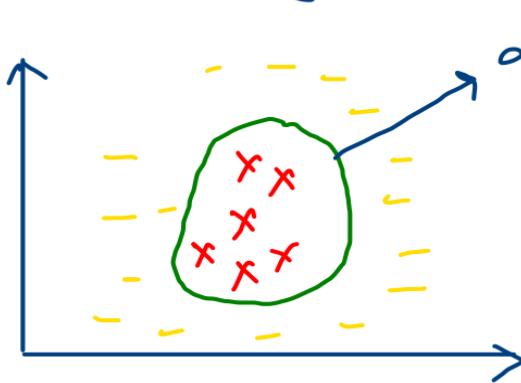
Classification (Decision Boundary)

classification (Binary)



→ we have to find decision boundary between two classes
→ now discuss linear decision boundary

However, it may be non-linear also.



Simple Case (two features)

DB : Line

$$w_1 x_1 + w_2 x_2 = w_0$$

Action: Find the values of parameters

$w = [w_0 \ w_1 \ w_2]$ using training

Samples $\langle x_i, y_i \rangle$

$$w_1 x_1 + w_2 x_2 - w_0 > 0 \text{ (One side of line)}$$

↳ for such samples $y=1$ (positive class)

$$w_1 x_1 + w_2 x_2 - w_0 < 0 \text{ (Other side of line)}$$

↳ for such samples $y=0$ (negative class)

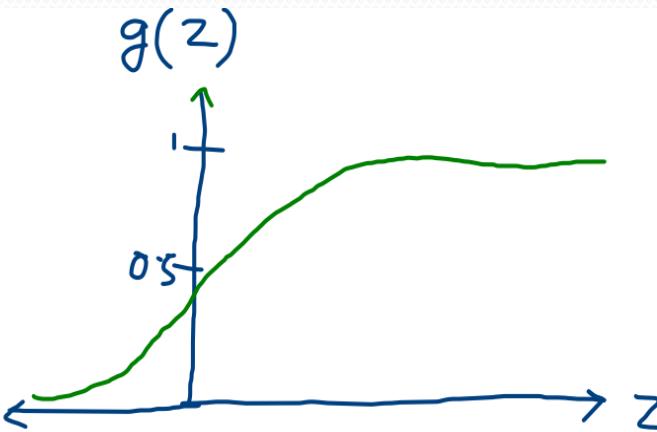
$$ax + by = c$$

Logistic Regression

$$\hat{h}(x_i) = g(w^T x_i)$$

$$h(x) = g(w^T x)$$

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



$\rightarrow g(z) > 0.5$ when $z > 0$

Suppose, predict $y=1$; $h(x) > 0.5$

$w^T x > 0$ [Same as the previous example
for two-class Case]

Predict $y=0$; $h(x) < 0.5$

$w^T x < 0$

Data set for classification

x_1	x_2	y	class Label	y
2	0.6	spam	1	0
1	5	not spam	2	1
6	2.8	not spam	2	1
3	1.4	spam	1	0

Format 1

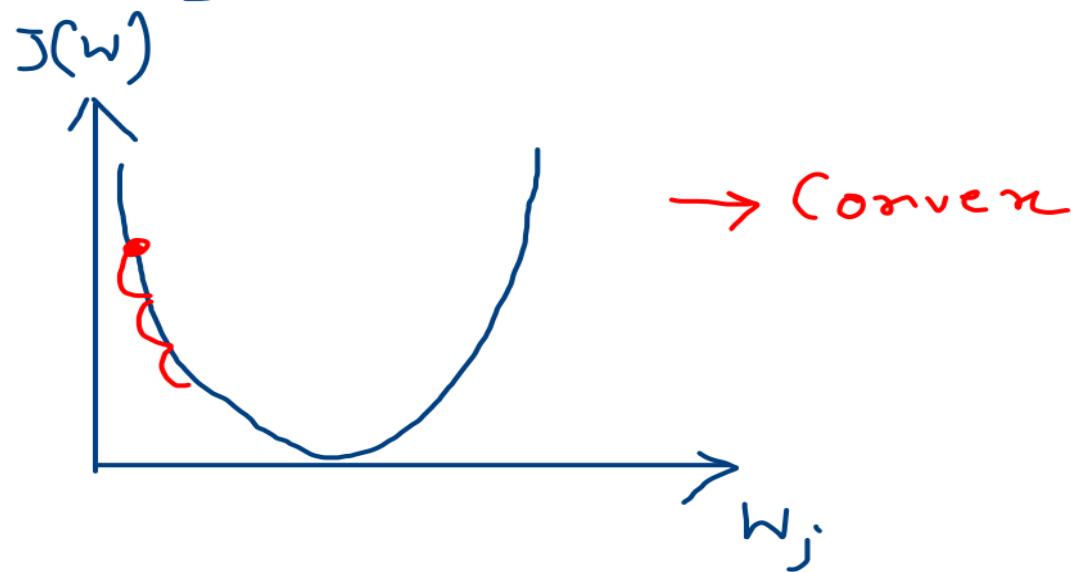
Format 2

After
re-write
the
class label

Logistic Regression (Cost Function)

Linear regression

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$



Logistic regression

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2 \rightarrow \text{Is it convex}$$

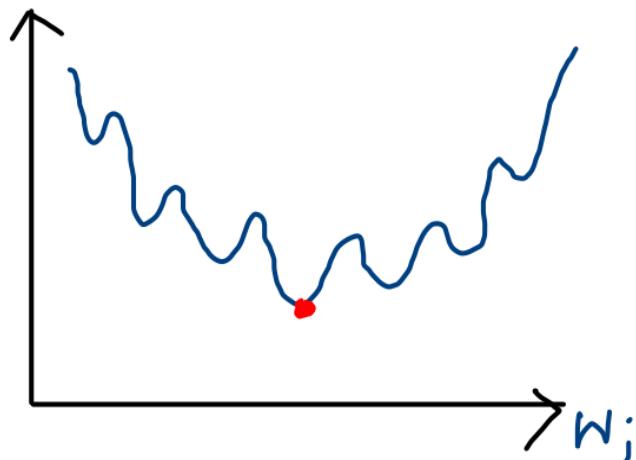
included here

$$h(x_i) = g(w^T x_i)$$

$$= \frac{1}{1 + e^{-w^T x_i}}$$

} non-linearity involves

$$J(w)$$



→ non-convex
(may stuck
into local
optima)

Logistic regression with MSE

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$
$$h(x_i) = g(w^T x_i)$$

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

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

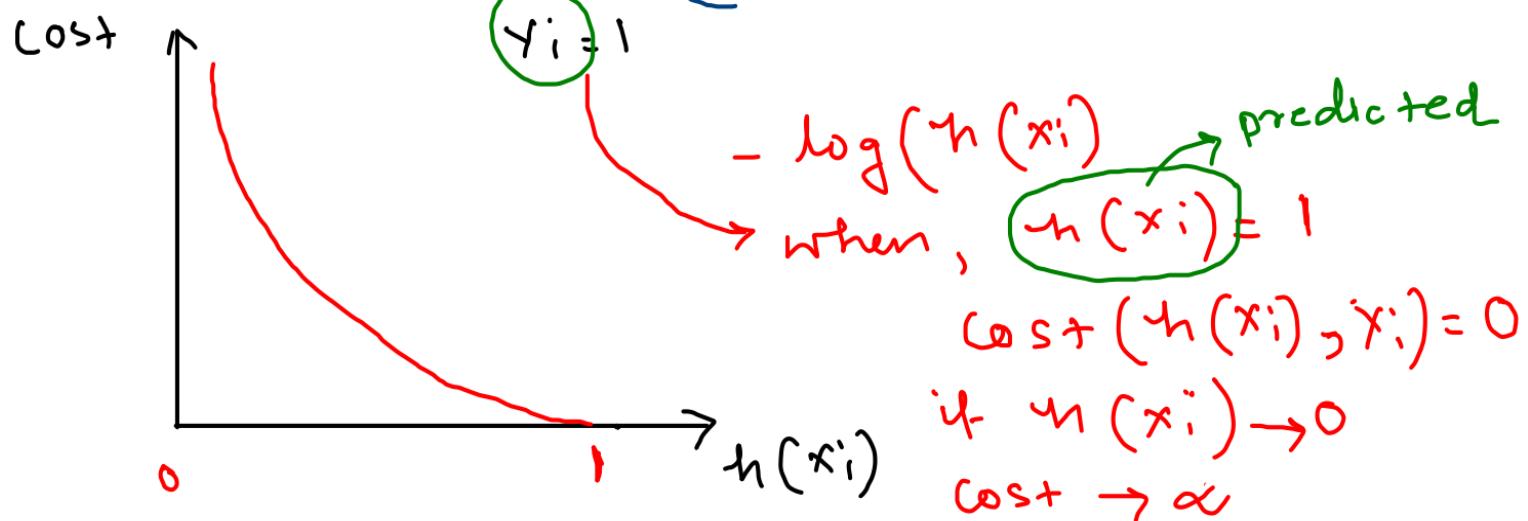
$$\text{SGD: } w_j := w_j - \frac{\partial}{\partial w_j} J(w)$$

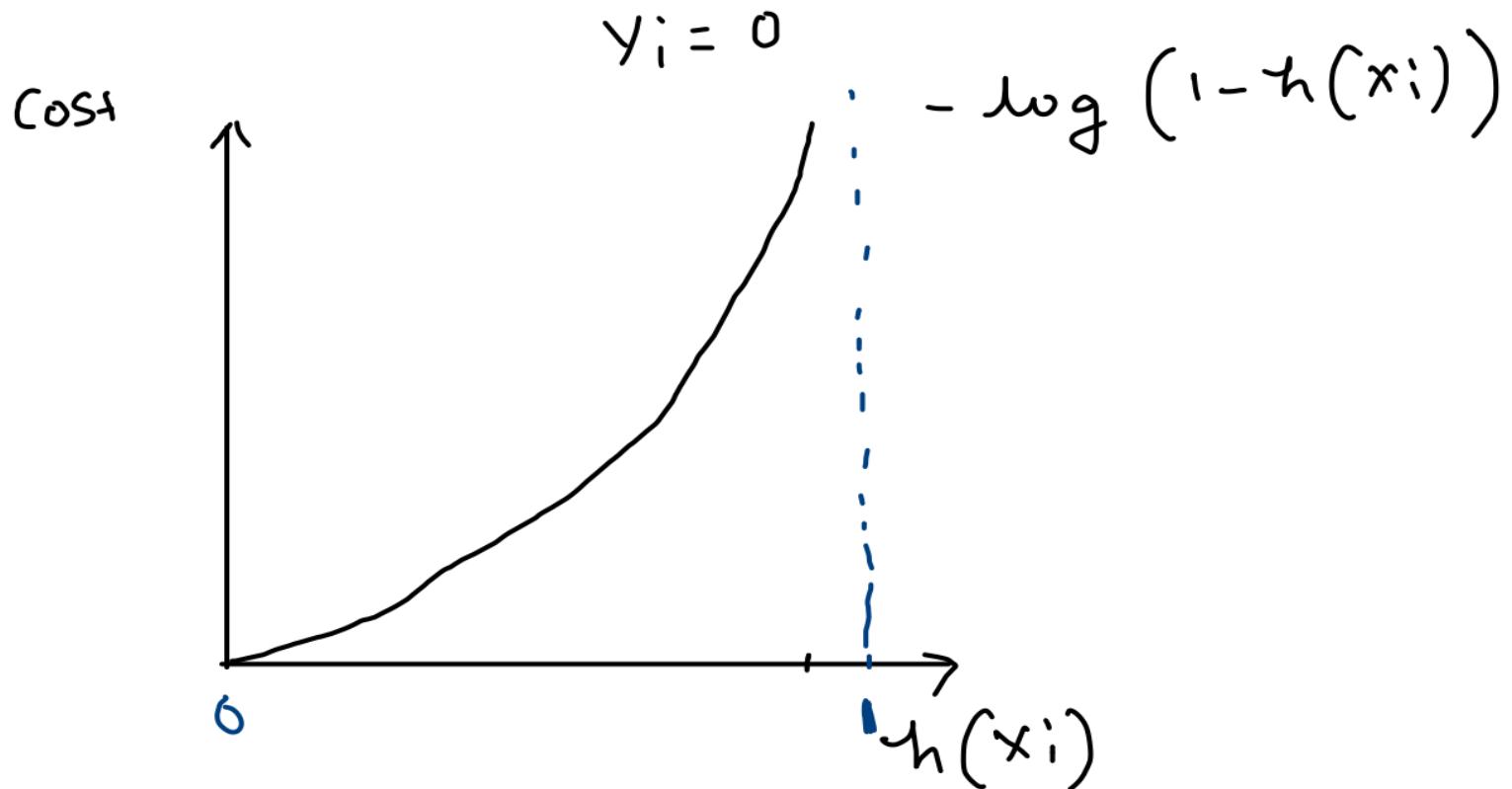
$$w_j := w_j + \alpha (y_i - h(x_i)) x_{i,j} \underbrace{h(x_i) \cdot (1-h(x_i))}_{}$$

Cost function for logistic regression

$$J(w) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h(x_i), y_i)$$

$$\text{Cost}(h(x_i), y_i) = \begin{cases} -\log(h(x_i)) & \text{if } y_i = 1 \\ -\log(1 - h(x_i)) & \text{if } y_i = 0 \end{cases}$$





if $y_i = 0$, $h(x_i) = 0$, $\text{cost} = 0$

$$h(x_i) \rightarrow 1$$

$$\text{cost} \rightarrow \infty$$

Simplified Cost function

$$\text{Cost}(h(x_i), y_i)$$

$$= -y_i \log(h(x_i)) - (1-y_i) \log(1-h(x_i))$$

$$y_i = 0 \text{ or } 1$$

$$\text{when } y_i = 0$$

$$\text{Cost}(h(x_i), y_i) = -\log(1-h(x_i))$$

$$\text{when } y_i = 1$$

$$\text{Cost}(h(x_i), y_i) = -\log(h(x_i))$$

Cost function for logistic regression

$$J(w) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h(x_i), y_i)$$

$$= \frac{1}{m} \sum_{i=1}^m [-y_i \log(h(x_i)) - (1-y_i) \log(1-h(x_i))]$$

$$= -\frac{1}{m} \sum_{i=1}^m [y_i \log(h(x_i)) + (1-y_i) \log(1-h(x_i))]$$

minimize $J(w)$ to obtain

$$w = [w_0 \ w_1 \ \dots \ w_n] \quad h(x_i) = \frac{1}{1 + e^{-w^T x_i}}$$

$$\text{SGD} : w_j := w_j - \alpha (h(x_i) - y_i) x_{i,j}$$

→ Same update rule as linear regression

Algorithmic representation for gradient descent

Step 1: Randomly initialize w_0, w_1, \dots, w_n and learning rate

Step 2: Calculate $h(X_i)$ for all training samples

Step 3: Calculate $J_{itr}(W)$ for itr^{th} epoch

Step 4: Update all parameters w_0, w_1, \dots, w_n

Step 5: Repeat steps 2-5 until

$|J_{itr}(W) - J_{itr-1}(W)| > \rho$ or number of epochs $> no$

Step 6: Apply threshold to take final decision for binary classification

Step 6: Report the value of $J_{itr}(W)$ and classification accuracy

ρ : small value (may be 0.0001)

no : sufficiently large so that model can learn

Can we stop using validation set?

Epoch/training step: one time updating considering all the training samples

GD for logistic regression

Training samples

x_1	x_2	Class Label
1.2	3	1
3	4	2
6	7	1
9	10	1

$$w = [1 \ 1 \ 1]$$

$$\alpha = .1$$

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

- ① Normalized the dataset
- ② Split train set / Test set / VS

TS

x_0	x_1	x_2	y
1	1.2	3	1
1	3	4	0
1	6	7	1
1	9	10	1

Epoch 1 $h(x_i) = g(w^T x_i) = g(w_0 x_0 + w_1 x_1 + w_2 x_2)$

TS1 : $h(x_1) = g(w_0 \times 1 + w_1 \times 1 \cdot 2 + w_2 \times 3)$
= .6 (not exact here / you need
to calculate the exact one)

TS2 : $h(x_2) = .3$

TS3 : $h(x_3) = .5$

TS4 : $h(x_4) = .2$

Cost function \rightarrow MSE (non-convex)
 \rightarrow log loss (Convex)

MSE

$$J(w) = \frac{1}{2m} \sum (h(x_i) - y_i)^2$$

$$J(w) = \frac{1}{2m} \left[(0.6 - 1)^2 + (0.3 - 0)^2 + (0.5 - 1)^2 + (0.2 - 1)^2 \right]$$

$$w_j := w_j - \alpha \frac{\partial J(w)}{\partial w_j} \rightarrow h(x_i) \text{ includes sigmoid function}$$

$$w_j := w_j + \alpha \frac{1}{m} \sum_i [(y_i - h(x_i)) x_{i,j} \cdot h(x_i) (1 - h(x_i))]$$

$$w_0 := w_0 + \alpha \frac{1}{m} \sum_i [((1 - 0.6)x_0 + 0.6x_1)(1 - 0.6) + ((1 - 0.3)x_0 + 0.3x_1)(1 - 0.3)]$$

$$w_1 :=$$

$$w_2 :=$$

$$\text{log loss} \\ J(w) := -\frac{1}{m} \sum_{i=1}^m [y_i \log(h(x_i)) + (1-y_i) \log(1-h(x_i))]$$

$$w_j := w_j - \alpha \frac{\partial J(w)}{\partial w_j}$$

$$w_j := w_j + \alpha \frac{1}{m} \sum_{i=1}^m (y_i - h(x_i)) x_{i,j}$$

$$w_0 = ?$$

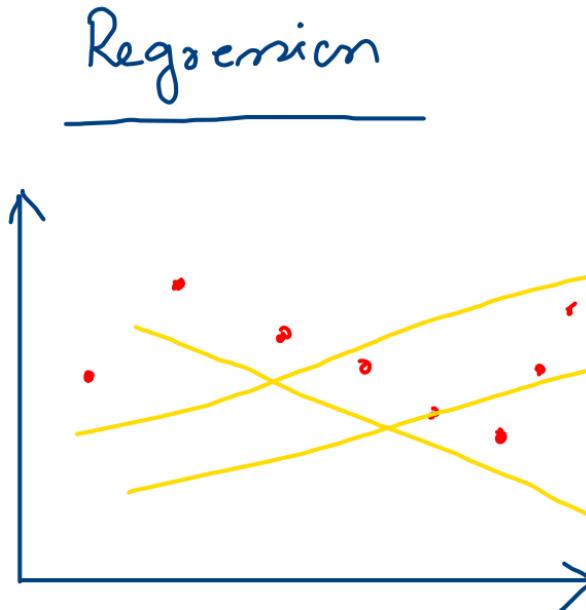
$$w_1 = ?$$

$$w_2 = ?$$

Repeat the training epoch by epoch
until convergence(?)

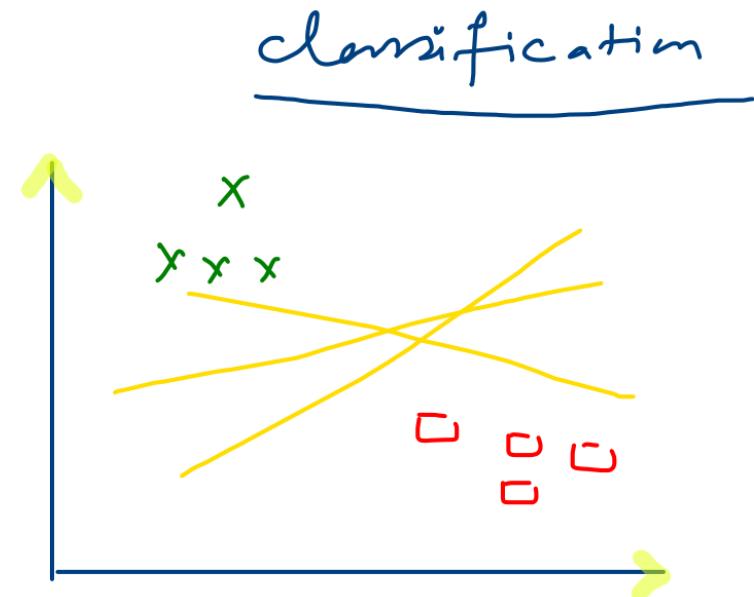
→ Finally report test accuracy after convergence
of training

Regression vs. Classification



Linear Relation between independent and dependent variables

Try to find best fit one by minimizing cost function



Data set is linearly separable

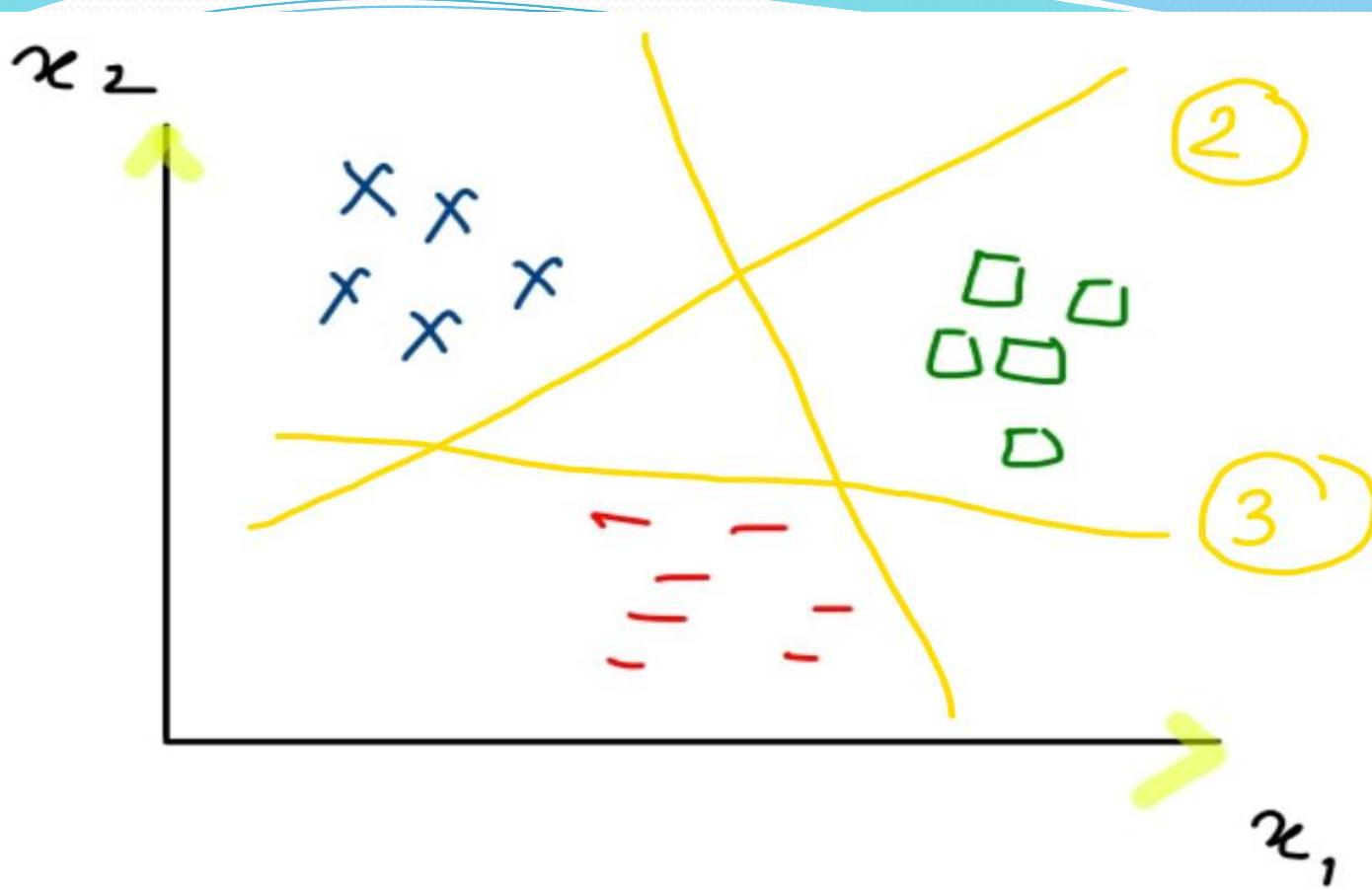
Try to find linear DB by minimizing cost function

Consider binary classification

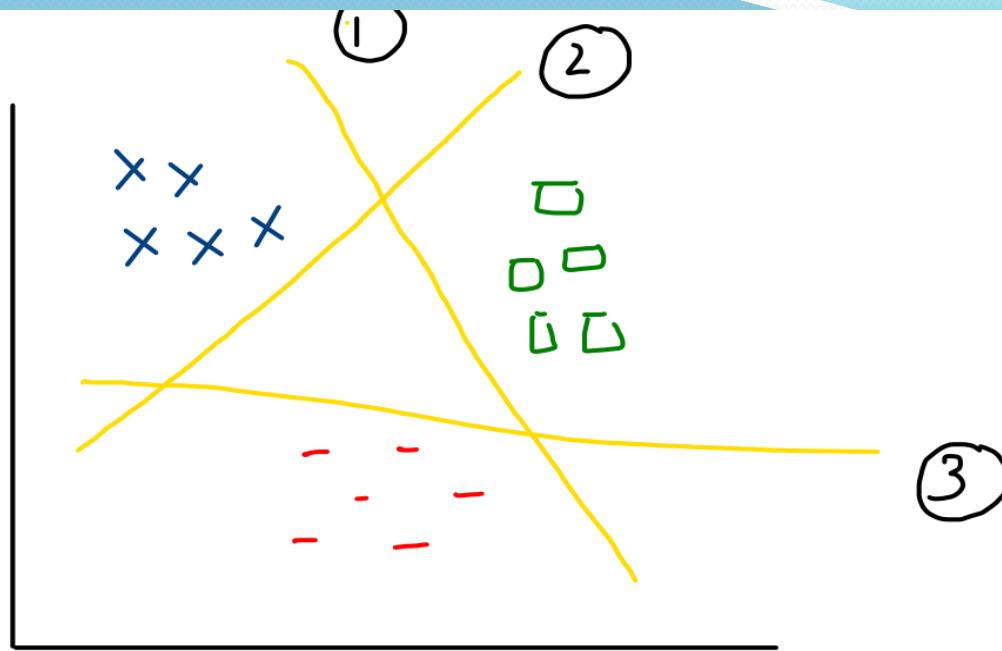
Multi-class Logistic Regression

x_1	x_2	class label	
2	3	1	
4	5	2	
6	7	3	→ Three classes
1.2	3	1	

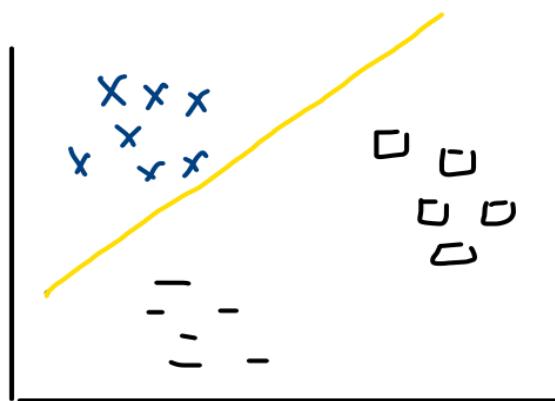
x_2  $\begin{matrix} x & x \\ x & x \end{matrix}$ $\begin{matrix} \square & \square \\ \square & \square \\ \square \end{matrix}$ $\begin{matrix} 1 & - \\ - & - \\ - & - \end{matrix}$  x_1



One vs. All



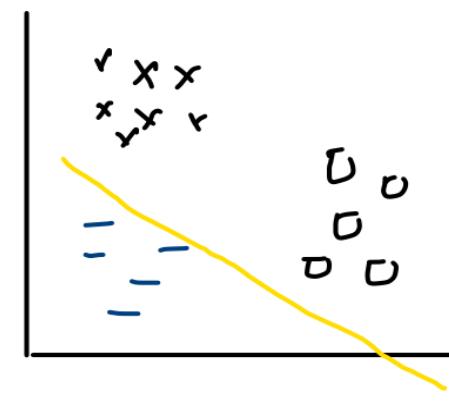
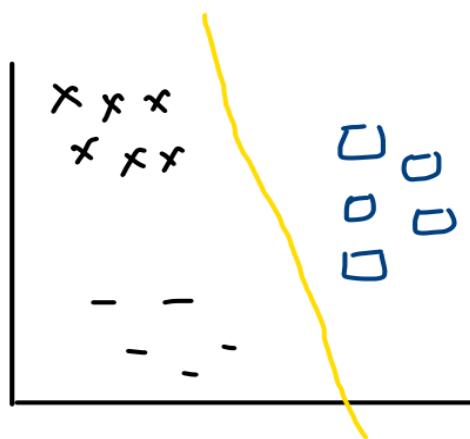
Multi-class, linearly separable



classifier model 1

Dataset

classifier model 2
Dataset preparation??



classifier
model 3

Original Dataset (Training samples)

x_1	x_2	Class label
2	3.4	1
6	1.2	2
3	4.5	2
6	1	3

Dataset for training classifier 1

x_1	x_2	y
2	3.4	1
6	1.2	0
3	4.5	0
6	1	0

Same way
prepare for
classifier 2 and 3

m = number of pattern

$x_1 \ x_2 \ \dots \ x_i \ \dots \ x_m$

n = number of features

$x_i = [x_{i,1} \ x_{i,2} \ \dots \ x_{i,j} \ \dots \ x_{i,n}]$

c = number of classes

$k \in \{1, 2, \dots, c\}$

One-vs-all

{ Train a logistic regression model $y_k^*(x_i)$ for each class k to predict probability that $y_i = k$ for a given x_i

→ Testing (class label for unknown pattern)??

Testing phase

$x_i \rightarrow \text{class label} ?? \rightarrow \text{predicted class label is } k$

$$\left\{ \begin{array}{l} h^1(x_i) \rightarrow \\ h^2(x_i) \rightarrow \\ h^3(x_i) \rightarrow \\ \vdots \\ h^k(x_i) \rightarrow \\ \vdots \\ h^c(x_i) \rightarrow \end{array} \right.$$

- Choose k
 $\max_k h^k(x_i)$

→ parameters are already estimated for all cases.

Some observations

- Is there any affect of initial parameter values in performance? (Hyperparameter Tuning/Validation Set)
- *Is there any affect of different random set of training samples (considering same percentage) in performance of the model?*
- How to check these issues during experimentation?

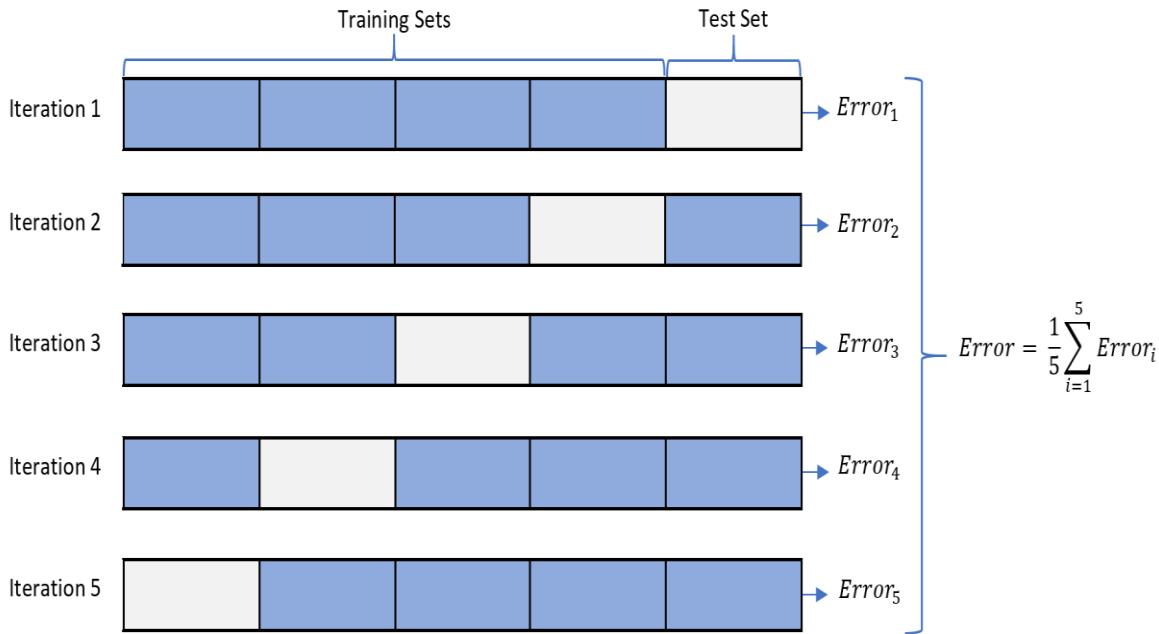
K-fold Cross-Validation

- Cross-validation is a statistical method used to estimate the skill of machine learning models.
- It is a resampling procedure used to evaluate the performance of learning models on different train/test split.
- The model is expected to perform well in generalization when used to make predictions for the samples not used during the training of the model.
- Here, the dataset is divided in train set and test set only.

Procedure

- Shuffle the dataset randomly (**Why?**)
- Split the dataset into k groups/folds (approximately equal size)
- For each group:
 - Take the group as a hold out or test data set
 - Take the remaining groups as a training data set
 - Fit a model on the training set and evaluate it on the test set
 - Retain the evaluation score and discard the model
- Report the performance score of the model on each group (test set) and average of all the performance score

Validation set??



Observation

- Some portion of training set can be used for validation set (for example 10%)
- The value of K: mostly use 5 or 10 for experimentation
- Each pattern is assigned to an individual group and stays in that group for the duration of the procedure
- This means that each sample is given the opportunity to be used in the hold out set/test set 1 (one) time and used to train the model $K-1$ times.

Confusion Matrix for Classification

- Confusion Matrix is used to know the performance of a Machine learning classification. It is represented in a matrix form.
- Confusion Matrix gives a comparison between actual and predicted values.
- The confusion matrix is a $N \times N$ matrix, where N is the number of classes.
 - For 2 class-problem ,we get 2×2 confusion matrix.
 - For 3 class-problem ,we get 3×3 confusion matrix.

Confusion Matrix for Classification

		ACTUAL VALUES	
		POSITIVE	NEGATIVE
PREDICTED VALUES	POSITIVE	TP	FP
	NEGATIVE	FN	TN

		ACTUAL VALUES	
		POSITIVE	NEGATIVE
PREDICTED VALUES	POSITIVE	560	60
	NEGATIVE	50	330

- The target variable has two values: **Positive** or **Negative**
- The **columns** represent the **actual values** of the target variable
- The **rows** represent the **predicted values** of the target variable

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

Confusion matrix Actual

		1 (+)	2 (-)
predicted	1 (+)	30 TP	30 FP
	2 (-)	10 FN	930 TN

Observation

$$\text{Accuracy} = \frac{30 + 930}{30 + 30 + 10 + 930} = 0.96$$

good one
but not
provide the
complete
idea

$$\text{Accuracy}_{\text{class 1}} = \frac{30}{40} = .75$$

$$\text{Accuracy}_{\text{class 2}} = \frac{930}{960} = .96$$

Suppose, class 1 is more important
than class 2

Precision vs. Recall

- **Precision:** Out of all the positive predicted samples how many are truly positive.

$$Precision = \frac{TP}{TP + FP}$$

The precision value lies between 0 and 1.

- **Recall:** How many of the actual positive cases we were able to predict correctly with our model.

$$Recall = \frac{TP}{TP + FN}$$

The recall value lies between 0 and 1.

Precision vs. Recall (Multiclass)

- Calculate precision and recall for each of the classes
- *precision* for the Cat class is the number of correctly predicted Cat photos (4) out of all predicted Cat photos ($4+3+6=13$), which amounts to $4/13$
- the *recall* for Cat is the number of correctly predicted Cat photos (4) out of the number of actual Cat photos ($4+1+1=6$), which is $4/6$

		True/Actual		
		Cat (😺)	Fish (🐠)	Hen (🐓)
Predicted	Cat (😺)	4	6	3
	Fish (🐠)	1	2	0
	Hen (🐓)	1	2	6

Some issues related to training samples

- Training samples in one class are significantly fewer than the samples in other classes (class imbalance)
- Sufficient training samples are not available (data scarcity)