

CH - Logistic Regression

→ Geometric Intuition of Logistic Regression.

- Don't go with name regression, it is a classification algo.
- It is a simple, elegant classification algorithm.
- Logistic regression can be interpreted / derived from 3 ways

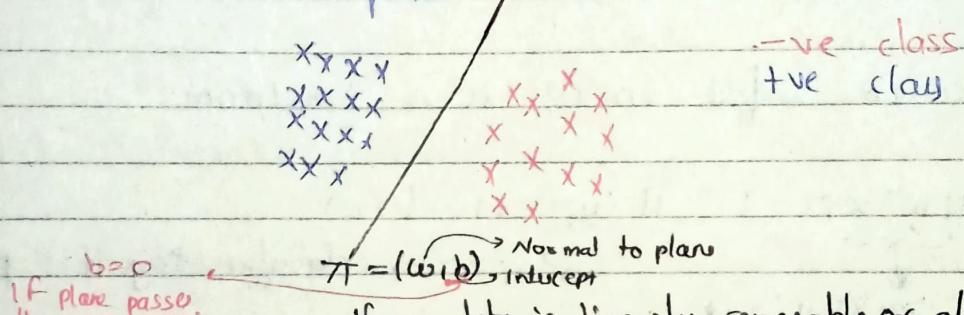
(i) Geometry

(ii) Probability

(iii) Loss-func

We will start with geometric intuition, then take a brief look on probabilistic intuition and then loss-func

Imagine we have 2 class points



Assumption in LR :- classes are almost / perfectly linearly separable.

We have data set $D_n := \{ +ve, -ve \}$

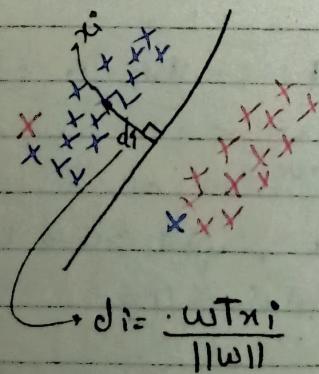
task is to find the plane $\pi(w, b)$ that separates the classes.

If assumption is wrong, then LR will not work.

Note: $y_i = +1$ for +ve class

$= -1$ for -ve class

b is not zero in practical methods.



$$d_i = \frac{w^T x_i}{\|w\|}$$

let w is unit vector then $\|w\| = 1$

$$\therefore d_i = w^T x_i$$

$d_i = w^T x_i > 0$ (x_i is in same side of w)

let a point from -ve class x_j , and

$$d_j = w^T x_j < 0$$
 (x_j is in opposite side of w)

classification:

$$\text{If } w^T x_i > 0 :$$

$$\text{then } y_i = +1$$

else

$$y_i = -1$$

* decision surface in LR is a Line/plane.

Case 1: $y_i w^T x_i$: If $y_i = +1$ (+ve)

$$\begin{cases} w^T x_i > 0 \rightarrow \text{classifer says it is positive} \\ > 0 \& y_i = +1 \end{cases}$$

then plane is correctly classifying the point.

Case 2: $y_i = -1$, $w^T x_i < 0 \Rightarrow$ LR concluding that x_i is -ve point

$$y_i \cdot w^T x_i > 0$$

-ve -ve

Hence for both +ve & -ve point if $y_i w^T x_i > 0$ then model is correctly classifying the point x_i .

Case 3: $y_i = +1$, $w^T x_i < 0 \Rightarrow$ LR concluding that x_i is -ve

$$y_i \cdot w^T x_i < 0$$

+ve +ve

Case 4: $y_i = -1$, $w^T x_i > 0 \Rightarrow$ LR concluding that x_i is +ve

$$y_i \cdot w^T x_i < 0$$

-ve +ve

\therefore For both +ve & -ve point if $y_i w^T x_i < 0$ then model is incorrectly classifying the point x_i

* A good classifier make min no. of wrong classification.

$$\text{we need } y_i w^T x_i > 0$$

~~find a plane~~

But we have n data points, we want to maximize.

$$\underset{w}{\text{maximize}} \sum_{i=1}^n y_i w^T x_i$$

they are fixed bcs they come from training data set

This is the only thing that can be changed.

Hence we want w which maximize the $\sum_{i=1}^n y_i w^T x_i$ in plane

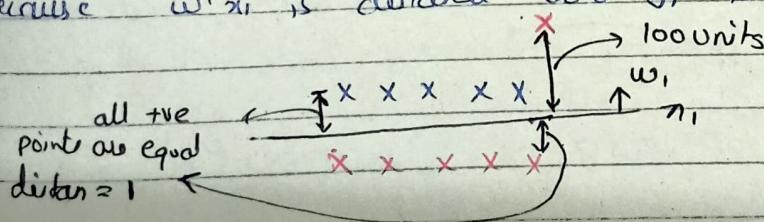
$$w^* = \arg \max_w \sum_{i=1}^n y_i w^T x_i \rightarrow \text{our final problem, although we modify it later.}$$

* Squashing & sigmoid function:-

$$w^* = \arg \max_w \sum_{i=1}^n y_i w^T x_i \quad \text{Optimization problem.}$$

Let's understand this function:

Let $y_i w^T x_i$ is a signed distance, signed distance because $w^T x_i$ is distance and y_i is sign of $+1, -1$

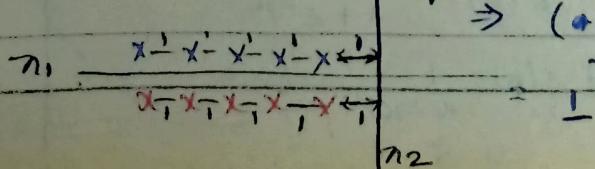


Case 1: n_1 is my separator

$$\sum y_i w^T x_i : - (1+1+1+1+1+1+1+1+1-100) \\ = -90$$

Here we get -90, due to an outlier point.

But n_1 seems a good plane. Now imagine another plane in above data points.



$$n_1 : \sum y_i w^T x_i = \\ \Rightarrow (1+1+1+1+1+1-1-1-1-1) \\ = 1$$

Here we get +1, so we might wrongly conclude that π_2 is a better separator than π_1 . It is all due to outliers which is favouring π_2 instead of π_1 . One single outlier is impacting my model.

Hence maximizing sum of signed distance formulation is prone to outliers. We need to modify it. We do it using squashing.

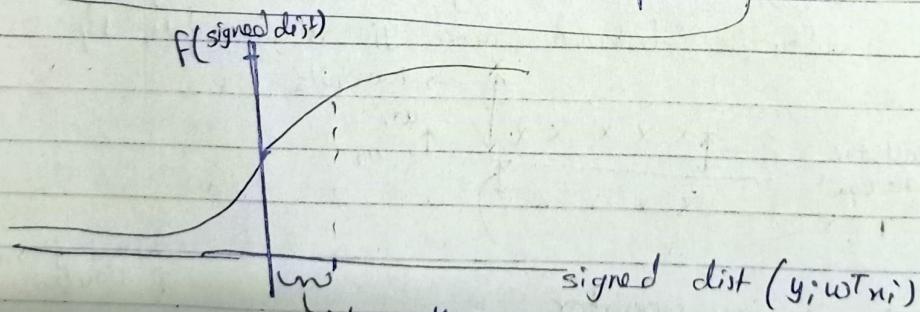
→ Squashing: Instead of using signed distance,

if signed dist is small: use it as is

if signed dist is large: make it a smaller value.

In our previous example: distance from π_1 , all the points have small distance i.e. 1, but only outlier has large distance of 100. we need to somehow decrease that.

We need a function that takes the signed distance and works as above 2 points



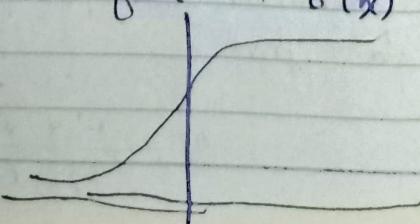
for smaller values funcⁿ should use them as if it is linear increase as dist increase, for larger value it should flatten or not increase so much.

One such func is sigmoid funcⁿ: $\sigma(x)$

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

For max value of $\sigma(x) = 1$
min val is 0.5

$$\sigma(0) = 0.5$$



Now our formula becomes

$$w^* = \arg \max_w \sum_{i=1}^n \sigma(y_i w^T x_i)$$

There are various other functions that can do one task but sigmoid is chosen because

① It have nice probabilistic interpretation

$\sigma(y_i w^T x_i) \rightarrow$ return probability of y_i ($i=1$)

$$\sigma(0) = 0.5$$

→ distance from plane is zero i.e. point exactly lies on the plane hence in that case probability of point belonging to +ve or -ve class is equal.

$$w^* = \arg \max_w \sum_{i=1}^n \sigma(y_i w^T x_i)$$

$$= \arg \max_w \sum_{i=1}^n \frac{1}{1 + e^{-(y_i w^T x_i)}}$$

② it is easy to differentiate. To actually solve our optimization problem we need a func that can easily differentiable.

→ optimization problem

lets try to solve our optimization problem

$$w^* = \arg \max_w \sum_{i=1}^n \sigma(y_i w^T x_i)$$

But Before that we need to understand monotonic func.

monotonic func:- A func $g(x)$ is called monotonic if if $x_1 > x_2$ then $g(x_1) > g(x_2)$

$g(x)$ is monotonic increasing func

We can use monotonic func to solve optimization problem,

Ex:- $x^* = \arg \min_x (x^2)$ if we want x st x^2 is minimize

$$f(x) = x^2$$

if $g(x) = \log(x)$, $x^* = \arg \min_x g(f(x))$

$$x^* = \arg \min_x g(f(x))$$

$$= \arg \min_x \log(x^2)$$

claim: $x^* > x$ because $g(n)$ is monotonic func.

Now lets go back to our optimization problem:

$$w^* = \arg \max_w \sum_{i=1}^n \frac{1}{1 + \exp(-y_i w^T x_i)}$$

$g(n) = \log(n)$: monotonic func

$$w^* = \arg \max_w \sum_{i=1}^n \log(\sigma(y_i w^T x_i))$$

$$= \arg \max_w \sum_{i=1}^n \log\left(\frac{1}{1 + e^{-(y_i w^T x_i)}}\right)$$

$$= \arg \max_w \sum_{i=1}^n -\log(1 + \exp(-y_i w^T x_i)) \quad \left. \begin{array}{l} \text{if } \log(w) \\ -\log \end{array} \right\}$$

We know that $\arg \min_x f(x) = \arg \max_x (-f(x))$

Hence

$$w^* = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

This is the optimization problem of LR
 Always remember that in LR $y_i \in \{-1, 1\}$.
 $\{1, 0\}$

→ weight-vector:- The optimal w^* we get is also called weight-vector because it is a d-dim vector where d = # of features.

$$\begin{matrix} w^T x \\ \downarrow x \in \mathbb{R}^d \end{matrix} \rightarrow d \times 1$$

It is called weight vector because it assign weights to each feature. $w = \langle w_1, w_2, \dots, w_d \rangle$

$$x_i = \langle f_1, f_2, f_3, \dots, f_d \rangle$$

interpretation of w :-

case ① if $w_1, w_2, \dots, w_i - w_d >$
if w_i is +ve and we

have a $x_q \in \{x_{q_1}, x_{q_2}, \dots, x_{q_i} - x_{q_d}\}$
if w_i is +ve ~~then w_i is~~ and $x_{q_i} \uparrow$ then $w_i x_{q_i}$
increases.

then $\sum w_i x_{q_i} \uparrow$

then $\sigma(w^T x_q) \uparrow$

then proba of $y_q = +1 \uparrow$

case ② if $w_i = -ve$,

as $x_{q_i} \uparrow$ $w_i x_{q_i} \downarrow \rightarrow \sum w_i x_{q_i} \downarrow \rightarrow \sigma(w^T x_q) \downarrow$

$P(y_q = -1) \uparrow \leftarrow P(y_q = +1) \downarrow$

$\rightarrow L_2$ regularization : overfitting vs Underfitting

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$\text{let } z_i = y_i w^T x_i$$

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-z_i))$$

$$\exp(-z_i) > 0 \quad (\forall z_i)$$

This is always
+ve

$\exp(z_i)$

$$\log(1 + \exp(-z_i)) > 0$$

> 0

$\therefore \log(1 + \exp(-z_i)) \geq 0$

$$\therefore \sum \log(1 + \exp(-z_i)) \text{ Then whole sum } \geq 0$$

our optimization problem was to get w that minimize the sum
we just showed that $\sum \log(1 + \exp(-z_i)) \geq 0$. min value
it can take = 0.

so, when does $\sum_{i=1}^n \log(1 + \exp(-z_i))$ becomes 0.

As $z_i \rightarrow \infty$, $\exp(-z_i) \rightarrow 0$

$$\therefore \log(1 + \exp(-z_i)) \rightarrow 0$$

$$\therefore \sum_{i=1}^n \log(1 + \exp(-z_i)) \rightarrow 0$$

Hence when $z_i \rightarrow \infty$, then $\sum_{i=1}^n \log(1 + \exp(-z_i)) \rightarrow 0$

$y_i w^T x_i$

$\cancel{+}$ \rightarrow modify w s.t. each $y_i \rightarrow \infty$

Fixed

If we pick w s.t

- (a) all training points are correctly classified
- (b) $z_i \rightarrow +\infty$ $\forall i$

Then that is Best w

But what if some of the points are outliers. We still manage to correctly classify them and that is nothing but overfitting.

To get rid of this problem, we modify the formula using regularization.

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda w^T w$$

Loss term Regularization term. $\lambda \|w\|_2^2$ L_2 Norm

If $w \rightarrow +\infty$ or $-\infty$

loss term b or becomes 0

but regularization term \uparrow

Hence it makes a balance.

If avoid w to be $+\infty$ or $-\infty$ so that

it is not overfit.

λ is a hyper parameter, it is not given we need to calculate it.

If $\lambda = 0$ then there is no regularization term hence it cause overfit.

If $\lambda = \infty$ then regularization term dominates and hence we are not using training data anymore it is under fit.

Hence, we need to get a balanced λ .

L_1 regularization and sparsity.

Now, our optimization formula becomes

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i w^T x_i))}_{\text{logistic loss}} + \lambda \underbrace{\|w\|_1}_{\text{regularizer}}$$

In L_1 regularization, instead of L_2 norm of $\|w\|$ we use L_1 -norm of w .

$$L_1 \text{ norm of any vector } v = \|v\|_1 = \sum_{i=1}^d |v_i|$$

Now our optimization formula becomes

$$w^* = \underset{w}{\operatorname{argmin}} \left(\sum_{i=1}^n \log(1 + \exp(-z_i)) + \lambda \|w\|_1 \right)$$

L_1 regularizer have a major advantage - sparsity.

Sparsity:- Let $w = (w_1, w_2, \dots, w_d)$

Solution to LR is said to be sparse if many w_i 's are zero.

If we use L_1 reg in LR, all the unimportant features become zero.

* L_1 creates sparsity in w , as compared to L_2 .

Why does?

* Probabilistic interpretation of Logistic regression.

Till now we have derived LR using geometry. Now we will look at probability method to derive LR.

We saw Naive Bayes (1) if features are real valued, we assumed that conditional probability of feature follow Gaussian distribution.

(2) The class label $y_i = +1 \text{ or } -1$ was assumed to be Bernoulli R.V.

With above 2 assumption we can derive LR.

$$\therefore \text{LR} = \underbrace{\text{GNB}}_{p(x_i^j | y_i)} + \underbrace{\text{Bernoulli class label}}_{y_i \sim \text{Bernoulli}}$$

$$w^* = \arg \min_w \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log(1-p_i) + \text{regularization}$$

where $p_i = \sigma(w^T x_i)$
 $\hookrightarrow \in \{+1, -1\}$

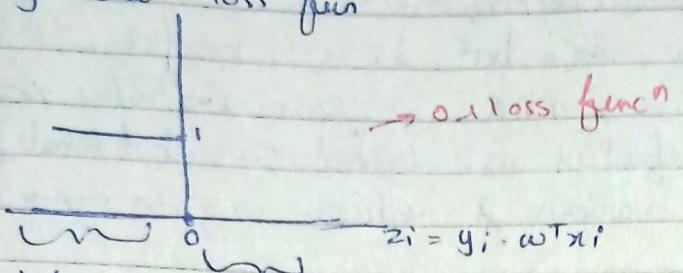
* Loss minimization interpretation of LR :- Till now we have learned geometry & probability derivation of LR. Now we will learn loss minimization derivation. Forget about all maths till now, let's imagine how we will design a classification func.

Let we assign +1 or incorrect classification and 0 for correct classification. Then we have to minimize

$$w^* = \arg \min_w \{ \text{number of incorrectly classified points} \}$$

Let we have a func which give +1 for incorrect clari & 0 for correct classification.

plot of ideal loss func.



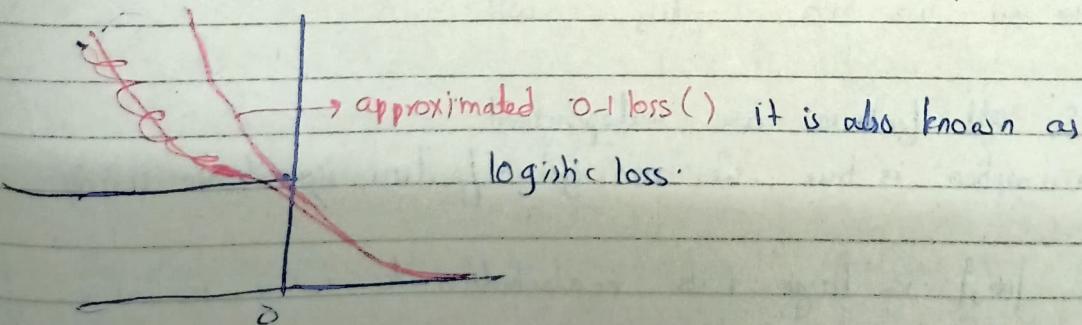
such a loss func is called 0-1 loss func.

$$w^* = \arg \min_w \sum_{i=1}^n 0-1 \text{ loss}(x_i, y_i, w)$$

$$\begin{aligned} 0-1 \text{ loss}(z_i) &= 1 \quad \text{if } z_i < 0 \\ &= 0 \quad \text{if } z_i \geq 0 \end{aligned}$$

One problem with 0-1 loss() func is that it is not continuous and hence it is not differentiable. But to solve an optimizat. problem we need a \approx differentiable func.

So, solⁿ is, to approximate the 0-1 loss() func.



logistic loss() is approximation of 0-1 loss.

And when we use logistic loss() as an approximation of 0-1 loss() we get LR

There is another approximation func called Hinge loss(). And when we use Hinge loss as loss func we get SVM.

When we use exp loss we get Ada Boost.

When we use SC-loss we get Linear Regression.

* Hyper parameters search / optimization

How to determine the best λ .

problem with λ is, $\lambda \in \mathbb{R}$ it have infinite possibilities.

(i) One technique to find λ is ~~grid~~ search (Brute force)

• Try with various λ values in given range.

(ii) Random search:- Randomly pick values in given interval.

* Column / feature standardisation.

In Logistic regression, it is compulsory to perform column standardisation. Because we are calculating distance from plane. We know that if two features have different scale that may impact the distance.

* Feature importance & model interpretability.

Let see how we can calculate the importance of feature in LR.

$$f_1 \ f_2 \ f_3 \dots f_j \dots f_d$$

$$w \rightarrow w_1 \ w_2 \ w_3 \dots w_j \dots w_d$$

assume: all features are independent.

If assumption is true, we can get feature importance using w .

If $|w_j|$ is large, its contribution to $w^T x$ is large and hence its importance is more.

* Collinearity (or) Multicollinearity:-

As we have already discussed that we can use feature weights to tell feature importance only when features are independent. But why this so. Let's understand it. But before, we need to understand collinearity.

if we have 2 feature f_i & f_j

s.t. $f_i = \alpha f_j + \beta$, then f_i & f_j are collinear.

~~if $f_i = \alpha f_j + \beta$~~

And collinearity across various feature is called multi-collinearity

$$\text{ex:- } f_i = \alpha_1 + \alpha_2 f_2 + \alpha_3 f_3 + \alpha_4 f_4$$

Then f_1, f_2, f_3, f_4 are said to be multi-collinear.

Q: Why does LwJ not be useful as feature imp. if features are collinear.

$$D := \sum_{i=1}^n \langle x_i, y_i \rangle$$

$$w = \langle 1, 2, 3 \rangle \quad ; \quad x_q = \langle x_{q1}, x_{q2}, x_{q3} \rangle$$

$$w^T x_q = x_{q1} + 2x_{q2} + 3x_{q3}$$

If $f_2 = 1.5 f_1$ (f_1 & f_2 are collinear)

$$w^T x_q = x_{q1} + 2 + 1.5 + x_{q2} + 3x_{q3}$$

$$= \underbrace{4x_{q1}}_{\downarrow} + \underbrace{3x_{q3}}_{\downarrow}$$

$$\begin{matrix} & 4 & 0 & 3 \\ \downarrow & w_1 & w_2 & w_3 \end{matrix} \quad \text{we get a totally new } w.$$

Hence w is changing i.e. for 2 different w we get the same $w^T x_q$ hence LwJ can't be used to get feature imp if features are not independent.

Now, the next question is, how do we know if features are collinear or not.

One technique is **Perturbation technique**

take a x_{ij} and add small noise ϵ to it; i.e. perturbing the column value.

Before adding noise compute w .

After perturbation again train the model and calculate w' ,

if $w \neq w'$ are differ drastically then we can conclude that our features are collinear and hence we can't use w to calculate feature importance.

f_1	f_2	f_3	y
x_{ij}	x_{ij}	x_{ij}	y_i

If we cannot use $\|w\|_1$ to calculate feature importance we can use forward feature selection.

★ Train & Test space time complexity.

Till now, we havn't seen how to train LR. To train LR we need to solve our :-

$$w^* = \underset{w}{\operatorname{argmin}} (\text{logistic loss} + \text{Regularization})$$

To solve this we use a technique called ~~SGD~~ stochastic gradient descent (SGD)

Train Time :- $O(nd)$ $n = \# \text{ of data point}$
 $d = \text{dimen / features}$.

→ Runtime (Testing)

Time : $O(d)$ → we only need to multiply $y^T x_2$
 $1 \times d \times 1$

Space :- $O(d)$ → we only need w i.e. weight vector

∴ We can see that run time & space complexity is low, we can use LR on Low-latency system.
∴ LR is v.v good for low-latency application

What if d is large say $d=1000$.

$w^T x_2$:- 1000 ~~is~~ multiplication
1000 addition.

→ In that case we use L_1 regular., which apply sparsity (w_j corresponding to less imp feature = 0)
here we need to check good λ , If λ is small
if $\lambda \uparrow$, latency \downarrow
But if $\lambda \uparrow$ Bias \uparrow (~~underfit~~)

∴ we have tradeoff $\lambda \uparrow$, Bias \uparrow , Latency \downarrow

(as per:-

- Decision surface :- Linear / hyperplane.
- assumption :- data is linearly separable or almost nearly sep.
- feature imp & interpretability :- will if feature are not collinear.
- Imbalanced dataset :- upsampling & down sampling
- Outliers :- less impact on model due to use of sigmoid.
- Missing values :- used standard imputation technique.
- Multi-class :- one vs Rest model
- Similarity matrix :- If similarity or dissimilarity matrix is given there is a extension of LR is there which can use similair matrx Kernel LR.
- Best & Worst case :-
 - Best : (i) Data is almost linearly separable
 - (ii) if we have low latency requirement.

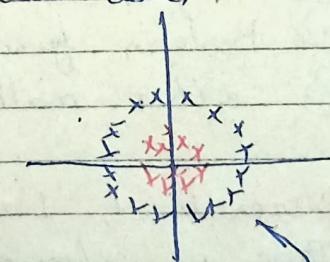
→ Large dimensionality :- if d is large LR works good, because if d is large, the chance of linear-separability increases but at same time computation time increases.

Non-linearly separable data.

What if we have non-linearly separable data.

Can we use LR on this data.

We can't directly apply LR on this but we can make some changes to apply LR



We can do feature transformation / feature engineering

$$f_1' = f_1^2$$

$$x \in \mathbb{R}, x_2 >$$

$$f_2' = f_2^2$$

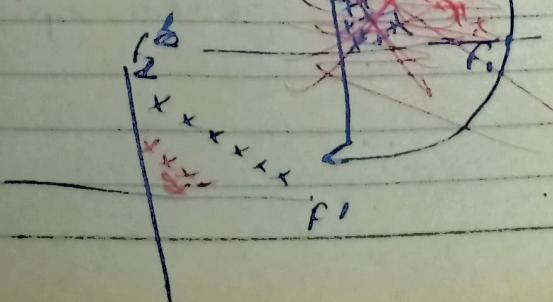
$$x' \in \mathbb{R}, x'_1, x'_2 >$$

Now,

$$x'_1 = x_1^2$$

$$x'_2 = x_2^2$$

Now we can apply LR.

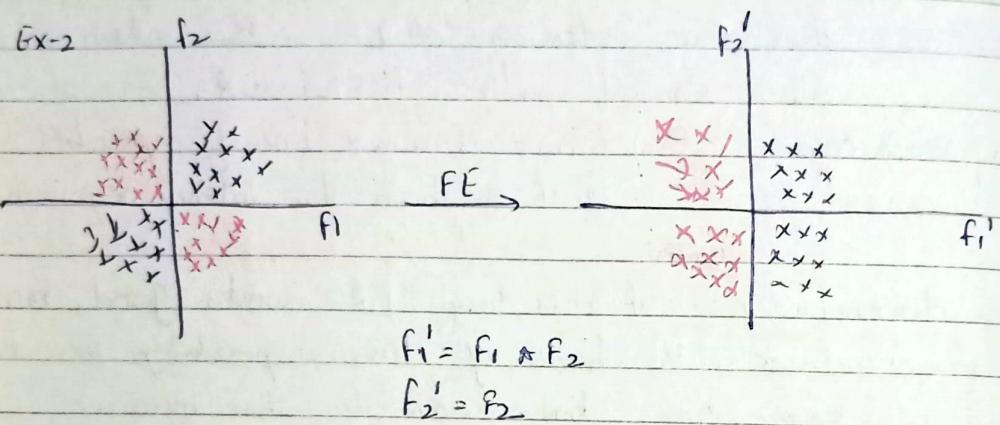


Now the question is, how do you know which transform to apply.

That we know using experience and practice.

The most important aspects of applied ML/AI

- (i) Feature engineering
- (ii) Bias-variance
- (iii) Data analysis & Visualisation



→ Typical transformation for real valued features

(i) polynomial : $f_1, f_2, f_1^2, f_2^2, f_1 f_2, f_1^3, f_2^3$

(ii) Trigonometric : $\sin(f_1), \cos(f_1)$

(iii) boolean funcn.

(iv) other mathematical : $\log(f_1), \exp(f_1)$

★ Generalized Linear model

There is extension of LR known as Generalized Linear model.

LR :- Gaussian NB + Bernoulli class \leftarrow prob

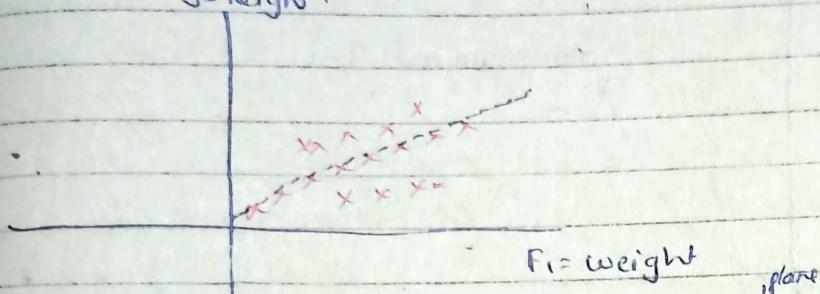
(i) Multinomial LR. - Instead of Bernoulli distribution of class we have multinomial distribution

(ii) Linear regression :- $P(y|x) \sim N(\mu, \sigma^2)$

CH - Linear Regression

* Linear Regression:-

Let's understand the geometry behind Linear Regression
 ex:- suppose we want to predict height of a person given some features like: weight, gender, ethnicity, hair color.
 $y = \text{height}$.



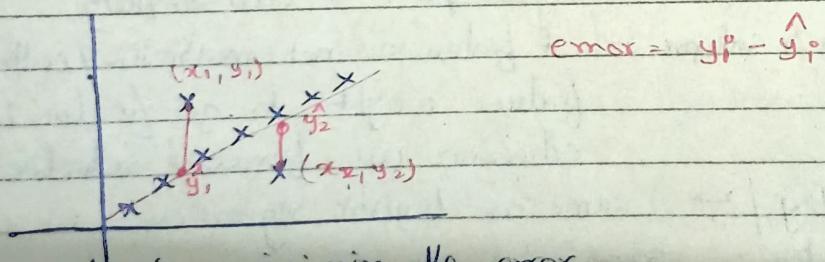
In Linear Regression we try to find a line/plane that fits the given data

$$\text{height} = w_1 * \text{weight} + w_0 \quad // \text{equation of line}$$

$$y = mx + c$$

object is to find w_1, w_0 .

Similarly in higher dimension instead of a line, we have plane, and we need to predict the parameters of plane.



We need to minimize the error.

$$(w^*, w_0^*) = \underset{w, w_0}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

vector
 scalar
 (intercept)

$$= \underset{i=1}{\operatorname{argmin}} \sum (y_i - (w^T x_i + w_0))^2$$

$$(w^*, w_0^*) = \underset{w, w_0}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (w^T x_i + w_0))^2$$

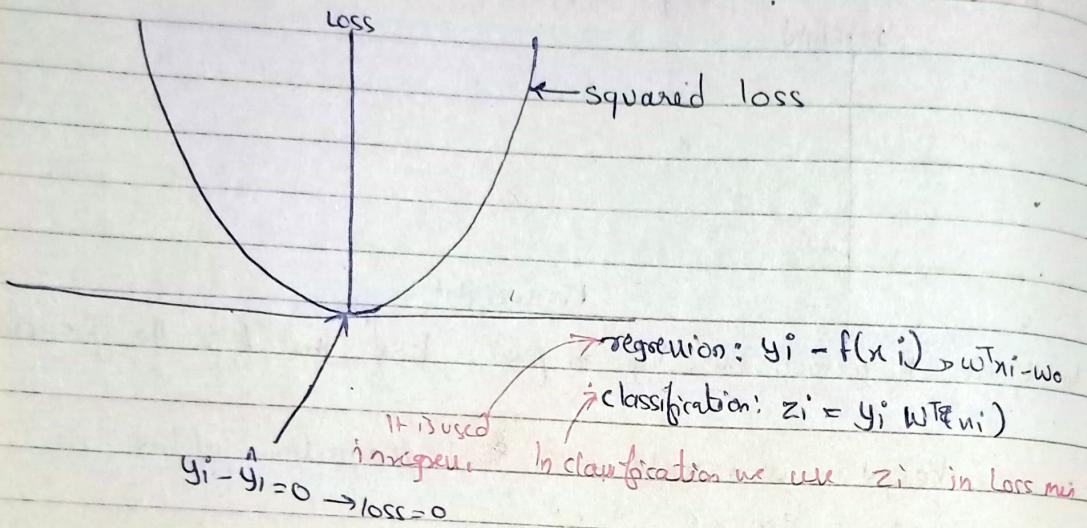
Just like Logistic regression, here also we need to apply regularization term

$$(\mathbf{w}^*, \mathbf{w}_0^*) = \underset{\mathbf{w}, \mathbf{w}_0}{\operatorname{argmin}} \sum_{i=1}^n \frac{1}{2} (y_i - (\mathbf{w}^T \mathbf{x}_i + w_0))^2 + \lambda \|\mathbf{w}\|_2^2$$

This can also
be L1 regularization term.

This is the geometric interpretation.

Now look at loss-minimization derivation.



★ Cases :- Most of the linear regression and logistic regression are same except for some.

1. Imbalanced data :- upsampling & down sampling
2. Feature Imp & Interprets :- If feature are not multicollinear/collinear we feature weight to get feature import. otherwise we forward selection
3. Feature Eng / \rightarrow same as logistic regression.
4. Regularization :- same " " "
5. Outliers :- In logistic reg we used "sigmoid fun", in linear reg. if we have outlier then error will be very large for that point \rightarrow hence linear Regress. is & very prone to outliers.
So, follow these steps.

①. D Train $\rightarrow \mathbf{w}^*, \mathbf{w}_0$

{ find point very far away from \mathbf{x}
remove these points as outliers
again train model with ~~selected~~ points
iterate in same ~~step~~ way.