

NAME Gagneet Kaur

ADDRESS _____

← MACHINE LEARNING TECHNIQUES → [MLT]

MOBILE _____

EMAIL 21F1006856@ds.study.iitm.ac.in

EMERGENCY CONTACT _____

WEEK 5

SUPERVISED LEARNING

- Supervised learning, as the name indicates, has the presence of a supervisor as a teacher.
- Basically supervised learning is when we teach or train the machine using data that is WELL-LABELLED.

INPUT : $\{x_1, x_2, \dots, x_n\} \quad x_i \in \mathbb{R}^d$ ← Features / Attributes
 $\{y_1, y_2, \dots, y_n\} \quad y_i \in \mathbb{R}$ ← Labels (supervision)

- Supervised learning has following categories :

1. Binary Classification

→ When the output variable is category of 'two' only.
for example : spam or non-spam.

2. Multiclass Classification

→ When the output variable is category of many classes.
for example : digit classification, i.e., $\{0, 1, \dots, k\}$

3. Regression

→ When the output variable is a '**REAL VALUE**'.
for example : rainfall prediction.

REGRESSION

Input Data / Training Data : $\{x_1, x_2, \dots, x_n\}$ $x_i \in \mathbb{R}^d$
 $\{y_1, y_2, \dots, y_n\}$ $y_i \in \mathbb{R}$

GOAL : To learn how to predict for a given new feature, which means that we basically want to learn :

$$h : \mathbb{R}^d \rightarrow \mathbb{R}$$

- How do we measure goodness of a function ?
- Error function (Error should be non-negative)

$$\text{error}(h) = \sum_{i=1}^n (h(x_i) - y_i)^2$$

- How small can this error be ?
- as small as ZERO
- Which 'h' achieves zero error ?
- $h(x_i) = y_i \quad \forall i \rightarrow$ this is memorizing data

- ISSUES** :
- (1) By memorizing, we can get zero error on 'TRAINING DATA' ,
 - (2) What we care about is — 'TEST PERFORMANCE'
 - (3) What can we do — Impose 'STRUCTURE' to reduce search space

Simplest Structure : LINEAR STRUCTURE

$$H_{\text{linear}} = \{ h_w : \mathbb{R}^d \rightarrow \mathbb{R} \mid h_w(x) = w^T x \quad \forall w \in \mathbb{R}^d \}$$

MODIFIED GOAL : $\min_{h_w \in H_{\text{linear}}} \sum_{i=1}^n (h_w(x_i) - y_i)^2$

OR EQUIVALENTLY

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2$$

Linear
Regression

- * If a linear regression model achieves zero training error, we can say that all data points lie on a hyperplane in the $(d+1)$ dimensional space.
- * When a model tries to interpolate each data point, the model tends to overfit. Thus, **TRAINING ERROR** will be low but **TEST ERROR** may be high.

OPTIMIZING THE ERROR FUNCTION

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (x_i^T w - y_i)^2 = \|x^T w - y\|_2^2$$

Let us do with matrices :

- Suppose a matrix X where all the data points are stacked in columns. In other words, x^T is a matrix where the data points are stacked in rows.

$$x^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \leftarrow \text{This is } n \times d \text{ matrix as } x_i \in \mathbb{R}^d.$$

- We also have a 'label vector' - y . (which gives labels for each of our data points)

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \leftarrow \text{This is } n \times 1 \text{ matrix}$$

- Our parameter vector ' w '.

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad \leftarrow \text{This is a } d \times 1 \text{ matrix.}$$

Hence,

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 = \|x^T w - y\|^2$$

equivalently,

$$\min_{w \in \mathbb{R}^d} (x^T w - y)^T (x^T w - y) \leftarrow \begin{array}{l} \text{unconstrained quadratic} \\ (\text{in } w) \text{ optimization problem} \end{array}$$

It is unconstrained because we allow for any w in \mathbb{R}^d ; there is no condition on w that we are imposing.

→ Now, How can we find the ' w ' that has the smallest value among this?

Solution : Take the derivative wrt w and set to zero !!

$$f(w) = (x^T w - y)^T (x^T w - y)$$

$$\therefore \nabla f(w) = 2(x x^T) w - 2xy$$

Now, we want to set the gradient to ZERO.

$$\Rightarrow 2(x x^T) w - 2xy = 0$$

$$\Rightarrow (x x^T) w = xy$$

∴ w^* satisfies this equation. Hence,

$$(x x^T) w^* = xy$$

let us write w^* in a closed form :

$$w^* = \underbrace{(xx^T)}_{d \times d \text{ matrix}}^{\dagger} \xrightarrow{\text{pseudo inverse}} xy$$

$\xrightarrow{\text{scaled covariance matrix}}$

* pseudo-inverse : it is like a many to one mapping, ie, there are so many points which go to the same point. Now, if you have to do an inverse - if it is one to one mapping, the inverse would give you a correct value. But if there is a whole subspace in a finite space, which maps to the same vector on the other side via transformation, the PSEUDO INVERSE would map it to the one that has the least length.

OBSERVATIONS :

- (1) Like PCA, w^* depends on a 'COVARIANCE' like matrix. Unlike PCA, here we have LABELS.
- (2) The labels should also dictate, how the answer w^* looks like.
- (3) Linear regression has a closed form solution (w^*)

$$y = w^T (x)$$

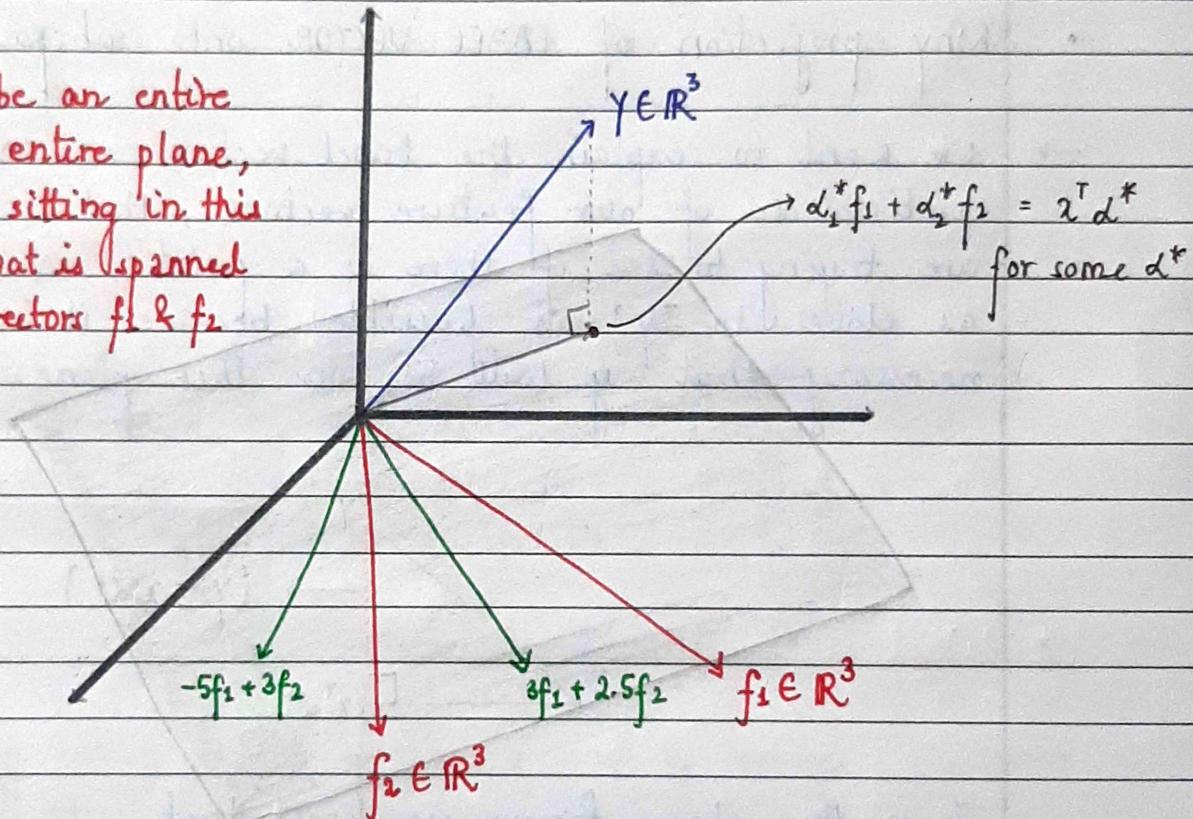
GEOMETRIC INTERPRETATION OF LINEAR REGRESSION

$$w^* = (x x^T)^{-1} x y$$

→ How can we interpret this geometrically?

Let us suppose $d = 2$ (no. of features)
 $n = 3$ (no. of data points)

There might be an entire subspace or entire plane, a 2-D plane sitting in this 3-D plane, that is spanned by the two vectors f_1 & f_2 .



$$\begin{aligned} x_1 - & \left[\begin{array}{c|c} f_1 & f_2 \\ \hline \end{array} \right] \\ x_2 - & \left[\begin{array}{c|c} f_1 & f_2 \\ \hline \end{array} \right] \\ x_3 - & \left[\begin{array}{c|c} f_1 & f_2 \\ \hline \end{array} \right] \end{aligned}$$

$$\begin{bmatrix} y \\ y_1 \\ y_2 \\ y_3 \end{bmatrix} \leftarrow \text{label vector}$$

Projection of vector (LABEL VECTOR) y onto the plane spanned by f_1 & f_2 can be written as :

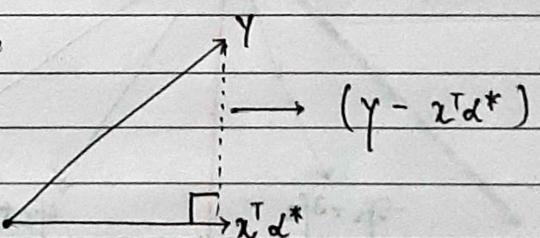
$$\alpha_1^* f_1 + \alpha_2^* f_2 = x^T \alpha^* \quad (\text{for some } \alpha^*)$$

$$\left[\begin{array}{cc} | & | \\ f_1 & f_2 \\ | & | \end{array} \right] \left[\begin{array}{c} \alpha_1^* \\ \alpha_2^* \end{array} \right] = x^T \alpha^*$$

$\underbrace{\hspace{2cm}}_{x^T}$ $\underbrace{\hspace{2cm}}_{\alpha^*}$

- Why projection of LABEL VECTOR onto subspace ?

→ We want to explain the label vector as some combination of our feature vectors , that is why we are trying to see if there is a point , which is as close to ' y ' as possible because it is not necessary that ' y ' will be on this plane.



From the above figure we know that ,

$$(y - x^T \alpha^*)^T (x^T \alpha^*) = 0$$

$$\Rightarrow y^T x^T \alpha^* - \alpha^* (x x^T) \alpha^* = 0 \quad - \textcircled{1}$$

We have : $y^T x^T \alpha^* - \alpha^* (x x^T) \alpha^* = 0$

Recall that,

$$w^* = (x x^T)^+ x y$$

Substituting w^* as α^* in ①, we get (use $x x^T$ is symmetric)

$$y^T x^T w^* - w^* (x x^T) w^* = 0$$

$$\Rightarrow y^T x^T [(x x^T)^+ x y] - [(x x^T)^+ x y] (x x^T) [(x x^T)^+ x y] = 0$$

This is basically tell us that the point that is closest to y , that is spanned by our features is $\rightarrow x^T w^*$.

CONCLUSION :

$x^T w^*$ is the 'PROJECTION' of the labels onto the subspace spanned by the features.

This is the GEOMETRIC INTERPRETATION.

When y belongs to the subspace spanned by features, the error (squared error) will be zero !!

GRADIENT DESCENT

We have,

$$w^* = (x^T x)^{-1} x^T y \quad \begin{matrix} \rightarrow \text{inverse computation} \\ \text{is expensive if } d \text{ is large} \end{matrix}$$

PROBLEM: The time complexity of finding linear regression solutions using the first derivative test is $O(d^3)$.

- Here d is the number of features.
- It is very computationally expensive

QUESTION: Is there a different way to achieve the same w^* without computing the inverse?

- We know w^* is the solution of an unconstrained optimization problem.
- w^* is just the minimizer of the sum of the squared error.
- One method to get w^* is → GRADIENT DESCENT.

What is Gradient Descent?

→ It is an iterative way to find minimizers of functions using just first order information.

What is Gradient Descent Algorithm?

1. Start with an initial guess for the solution
2. let us call that guess - w_0 (arbitrary points in d dimension)
3. Compute the gradient of the objective function at w_0 .

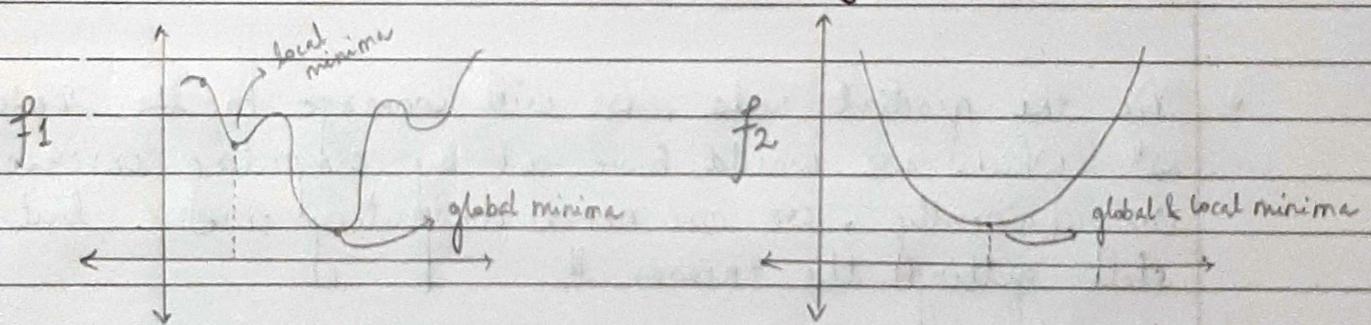
4. Since, a gradient would tell you the direction in which you have to move such that the function increases the most. So, we have to go opposite of it.
5. We want to minimize the function, naturally GRADIENT DESCENT would tell you to move in the direction opposite to what the gradient is pointing.

Update Rule :

$$w^{t+1} = w^t - \eta^t \nabla f(w^t)$$

(scalar step size)
gradient at w^t

6. Eventually, after iterations of algorithm, eventually we will reach - LOCAL MINIMA of the function.
- Local Minima is a point where the gradient becomes ZERO.
→ it does not mean that it is necessarily the global minima.



7. Once you reach a local minima, then your algorithm does not proceed any further. It stops at that point because the gradient becomes zero & then the update rule does not proceed any further.

The function which we care about is SUM OF SQUARED ERRORS \rightarrow it is a quadratic function, i.e., it has a unique global minima. \rightarrow which means the gradient descent algorithm will converge to the optimal point.

- Does this solve our problem of not computing inverse?
Let's see :

$$f(w) = \|x^T w - y\|^2 = \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\therefore \nabla f(w) = 2(x^T)w - 2xy$$

GRADIENT DESCENT UPDATE RULE FOR LINEAR REGRESSION :

$$w^{t+1} = w^t - \eta^t [2(x^T)w - 2xy]$$

- Via the gradient rule, we will converge to the same w^* which we would have got by computing inverse. Computationally, we are never computing inverse but still getting the answer !!

Even if 'd' is much, much larger, even if 'd' is in millions, can still run this algorithm.

Now, another ISSUE !! What if 'n' is large ??

$x^T x$ is a $d \times d$ matrix which came from $d \times n$ times $n \times d$ matrix ! So there is an 'n' dependency hiding inside this $x^T x$ itself.

- So, we might want to avoid XX^T computation itself.
- Not just the inverse of it, but XX^T might be hard to compute, if you have huge no. of data points.

How to adapt gradient descent to this situation?

SOLUTION : STOCHASTIC GRADIENT DESCENT

for $t = 1, \dots, T$

batch size

- At each step, sample a bunch of datapoints (k) uniformly at random from the set of all points
- PRETEND this sample is the entire dataset and take a gradient step w.r.t it

$$2[(\tilde{x}\tilde{x}^T)w - \tilde{x}\tilde{y}] \rightarrow \text{manageable because } \tilde{x} \in \mathbb{R}^{d \times k} \text{ where } k = \# \text{ sample points}$$

- END

Note : at every round, you are sampling k points, which means these k points are not the same at every round.

This sample & the gradient that comes out of the sample is a noisy gradient.

After T rounds we,

$$w_{SGD}^T = \frac{1}{T} \sum_{i=1}^T w^t$$

→ Guaranteed to converge to OPTIMA with high probability.

KERNEL REGRESSION

How can we do non-linear regression?

→ It is as if you want to map your data points to a higher dimension and then learn a linear model in high dimension. But we don't want to explicitly compute those high dimensional mappings.

WHAT CAN WE DO?

→ If our regression solution somehow can be written in terms of the dot products between data points

$$w^* = (x x^T)^{-1} x y \quad \leftarrow \text{Regression Solution.}$$

- w^* must lie in the span of the data points
- As we know, w^* is the solution that minimizes the sum of squared error :

$$\sum_{i=1}^n (w^{*T} x_i - y_i)^2$$

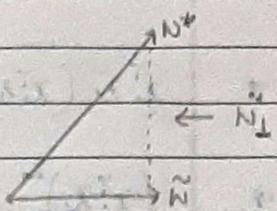
- If this w^* was actually out of this plane, let us assume \tilde{w} which is the projection of this w^* onto this plane, i.e., the closest point of w^* on this plane.
- Thus the sum of squared error will be :

$$\sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2$$

- What can we say about the above eqn?

We can write :

$$w^* = \tilde{w} + \tilde{w}_L$$



Now,

$$w^{*T} x_i = \tilde{w} + \tilde{w}_L^T x_i$$

It is same as :

$$\begin{aligned} w^{*T} x_i &= \tilde{w}^T x_i + \tilde{w}_L^T x_i \\ \Rightarrow w^{*T} x_i &= \tilde{w}^T x_i \end{aligned}$$

which means that the prediction of w^* & \tilde{w} for the all data points are exactly the same, which means that the error both of these will suffer is exactly SAME.

$$\sum_{i=1}^n (w^{*T} x_i - y_i)^2 = \sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2$$

Conclusion : w^* equals some linear combinations of the data points.

$$w^* = x \alpha^* \quad \text{for some } \alpha^* \text{ in } \mathbb{R}^n$$

$$\text{Also, } w^* = (x x^T)^+ x y$$

$$\therefore x \alpha^* = (x x^T)^+ x y$$

$$\Rightarrow (x x^T) x \alpha^* = (x x^T)(x x^T)^+ x y$$

$$\Rightarrow x(x^T x) \alpha^* = x y$$

$$\Rightarrow \mathbf{x}^T \mathbf{x} (\mathbf{x}^T \mathbf{x}) \alpha^* = \mathbf{x}^T \mathbf{x} y$$

$$\Rightarrow (\mathbf{x}^T \mathbf{x})^2 \alpha^* = (\mathbf{x}^T \mathbf{x}) y$$

Now,

$\mathbf{x}^T \mathbf{x} \rightarrow \text{kernel matrix } (k)$

$$k^2 \alpha^* = k y$$

and if k is invertible,

$$\alpha^* = k^{-1} y$$

What Is The Use Of This?

Now, when you do prediction - if you have a new data point x_{test} in d dimension, then the KERNEL would tell you that, I need $w^* \phi$ of x_{test} . So I have to map my x_{test} to high dimension & then take dot product with w^* .

$$\text{So, } w^* = \arg \min \sum_{i=1}^n (\phi x_i - y_i)^2$$

PREDICTION : for some $x_{\text{test}} \in \mathbb{R}^d$,

$$w^* \phi(x_{\text{test}})$$

$$= \left[\sum_{i=1}^n \alpha_i^* \phi(x_i) \right]^T \phi(x_{\text{test}})$$

$$= \sum_{i=1}^n \alpha_i^* k(x_i, x_{\text{test}})$$

$k(x_i, x_{test}) \rightarrow$ KERNEL FUNCTION (how similar is x_{test} to x_i)

$\alpha_i^* \rightarrow$ How important is i^{th} point towards w^*

Basically, we have managed to kernelize our linear regression problem to solve for non-linear regression as well.

PROBABILISTIC VIEW OF LINEAR REGRESSION

Dataset : $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Let us assume : $y_i | x_i \sim w^T x_i + \epsilon$

Noise $N(0, \sigma^2)$
 Gaussian

The label given
 the data point is
 generated as

unknown
 but fixed, $\epsilon \in \mathbb{R}^d$

WHAT EXACTLY THE ABOVE EQUATION MEANS ?

- Every y_i was generated according to this process.
- Somebody gave us x_i & then to get the y_i , there is an unknown but fixed w using which $w^T x_i$ was generated & then a noise got added and we are only seeing the noisy version of $w^T x_i$ whereas we know that the statistics of this noise is 0 mean and some known variance σ^2 .
- Can view this as an 'ESTIMATION PROBLEM' we are trying to estimate the w which after adding noise affects our labels.
- Solution Approach : MAXIMUM LIKELIHOOD

Likelihood Function

$$L(w; x_1, \dots, x_n, y_1, \dots, y_n) = \prod_{i=1}^n e^{-\frac{(w^T x_i - y_i)^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma}$$

$$\log L(w, x_1, \dots, x_n, y_1, \dots, y_n) = \sum_{i=1}^n -\frac{(w^T x_i - y_i)^2}{2\sigma^2} \cdot \frac{1}{\sqrt{2\pi}\sigma}$$

equivalently,

$$\begin{aligned} & \max_w \sum_{i=1}^n -(w^T x_i - y_i)^2 \\ & \equiv \min_w \sum_{i=1}^n (w^T x_i - y_i)^2 \end{aligned}$$

This is exactly the linear regression problem with squared error that we already put on, which means
 ↴ We know the solution to this.

$$\hat{w}_{ML} = w^* = (x x^T)^T x y$$

CONCLUSION : Maximum Likelihood Estimator assuming
 'ZERO MEAN GAUSSIAN NOISE' is same as
 'LINEAR REGRESSION WITH SQUARED ERROR'

- We are saying that if we chose squared error to solve the linear regression problem, it is as if we are implicitly making the assumption that there is a Gaussian noise (O mean that gets added to our labels that corrupts our labels). Choosing Noise means implicitly choosing ERROR FUNCTION.