

## Revisiting linear regression.

We earlier saw that under a LMS objective the optimal classifier under risk min is the conditional expectation.

Also if  $y \in \{0, 1\}$ , we know that

$$E[y|x] = P[y=1|x] = \sigma_1(x).$$

Thus, Linear regression can be viewed as a parametric approximation for the posterior.  
(linear)

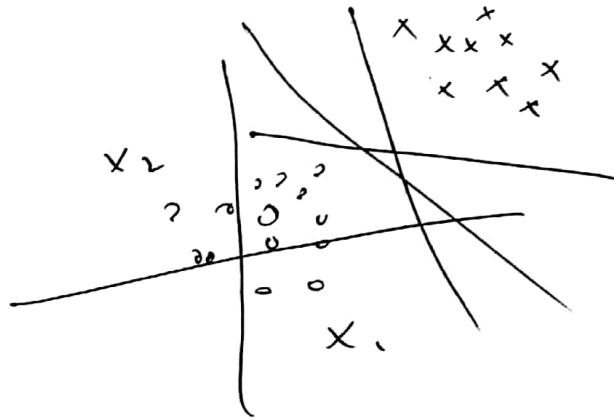
The idea of separability.

Training data  $D : \{(x_i, y_i), i=1, \dots, n\}$  is called linearly separable, if  $\exists W^*$  such

$$\text{that } \begin{aligned} x_i^T W^* &> 0 \text{ if } y_i = 1 \\ &< 0 \text{ if } y_i = 0 \end{aligned} \quad \left( \begin{array}{l} \text{or} \\ \phi(x_i)^T W^* \\ \text{for the} \\ \text{kernelized} \end{array} \right)$$

$W^*$  is called the separating hyperplane. case).

If data is linearly separable,  $\exists$  infinitely many  $w_s^*$



The family of linear classifiers aims to find out one such hyperplane under certain optimality criteria.

Note that  $q_i(x) \approx h(x) \approx \text{sign}(w^T x)$ .

We saw examples of Bayes' classifiers reducing to linear & quadratic discriminant functions under certain conditions.

We were looking at the OLS algo to find the discriminant function.

[notice that the formulation is very similar to both classification & regression].

We also saw that  $W_{OLS}^* = (\Phi^T \Phi)^{-1} \Phi^T y$ .

---

We also saw that in general,

$$\text{Error with a classifier} = \underbrace{\text{Bias}^2}_B + \underbrace{\text{variance}}_{+ \text{noise.}}$$

For a fixed value of error, unit decrease in bias results in more than quadratic increase in variance.

Thus, in practice, one always trades

decreased variance with increase bias in order to reduce the total error - <sup>[first overfit & then reduce the model capacity]</sup>

This process is generally referred to as the Regularization.

## Regularized Least Square

suppose the true model is a cubic polynomial.

But  $h(x)$  is assumed to be a  $10^{\text{th}}$  degree polynomial.

Now if one fits  $h(x)$  using LS, the coefficients of higher degrees [after 3] will be very high with alternating signs.

In other words, if  $h(x)$  should be closer to the 'true target' most of the  $w$ s should be zero or  $w$  vector should be sparse.

Thus, we want to minimize LMS with a constraint that  $w$  should be sparse.

Thus, our new objective,

$$W_{\text{reg}}^* = \underset{W}{\text{argmin}} \left\| Y^T - W^T \Phi^T \right\|_2^2$$

$$\text{s.t. } W^T W < \alpha$$

$[\|W\|_F < \alpha, \text{ in general}]$

Now, our new objective  $f^n$  becomes,

$$L = \left\| Y^T - W^T \Phi^T \right\|_2^2 + \lambda W^T W.$$

Can be shown that

$$W_{\text{reg}}^* = \underbrace{(\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y}.$$

One can now expect  $W_{\text{reg}}$  to be a  
having lesser MSE compared to  $W_{OLS}$ .

## The Bayesian view for regularization.

Earlier we saw that OLS is nothing but an MLE for the posterior with the following parametric form

$$f(y|x; w, \sigma) \sim N(\phi^T w, \sigma)$$

observe here that  $w$  is the parameter for this density.

suppose we have some prior 'information' about  $w$ .

i.e, let's model

$$f(w) \sim N(\bar{w}, 0, \alpha)$$

$$f(w) \propto \exp\left(-\frac{w^T w}{2\alpha^2}\right).$$

now posterior for the parameter becomes,

$$f(w|y) \propto \prod_{i=1}^n f(y_i | x_i; w, \sigma) f(w)$$

$$\propto \exp\left(\sum_{i=1}^n (y_i - w^T x_i)^2 - \frac{1}{2\alpha^2} w^T w\right)$$

Assume  
 $\phi = I$

(132).

for MAP estimate, we need to find the mode of the posterior.

$$\therefore \ln [f(w|y)] = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2 - \frac{1}{2\alpha^2} w^T w + C$$

$\Rightarrow$  maximizing the above is equivalent to minimizing the cost of the reg. least square.

Thus, regularized least square is the MAP estimate with some prior assumption on the parameter.

---

### Iterative algorithms for LR.

we have for OLS,

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

We saw a closed-form expression for this cost. However, this can be solved through an iterative scheme.

We know that the gradient of a  $f$  at any point gives the direction of steepest descent.

Thus, moving small steps along the gradient at every point could potentially lead us to the optimum. - This is the idea behind the gradient descent algo.

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

Thus, one can obtain an iterative algo as follows.

$$w(k+1) = w(k) - \underset{\text{learning rate}}{\eta} \sum_{i=1}^n x_i (x_i^T w(k) - y_i)$$

The above version is called the 'batch version' of the gradient descent algorithm

- It uses all the 'training samples' at once for each iteration.

Instead one can have an incremental version where only one training sample is used per iteration. in which case the update steps become:

$$w(k+1) = w(k) - \eta \underline{x_k^{(k)}} [x_k^T w_k - y_k]$$

(134)



## The perceptron learning algorithm.

One of the earliest algos for learning  $w^*$   
Iterative algorithm iterating over all training  
samples, picking one at a time.

Let  $w_k, x_k$  &  $y_k$  denote the estimate for  
the  $w^*$  at the  $k^{\text{th}}$  iteration,  $k^{\text{th}}$  feature  
vector & the labels respectively.

The perceptron algo:

Let  $\Delta w_k = w_{k+1} - w_k$ , then

$$\Delta w_k = 0 \quad \text{if} \quad \begin{array}{l} w_k^T x_k > 0 \quad \& \quad y_k = 1 \quad \text{or} \\ w_k^T x_k < 0 \quad \& \quad y_k = 0 \end{array}$$

$$= x_k \quad \text{if} \quad w_k^T x_k \leq 0 \quad \& \quad y_k = 1$$

$$= -x_k \quad \text{if} \quad w_k^T x_k \geq 0 \quad \& \quad y_k = 0.$$

$$w_{k+1} = \underline{w_k + \Delta w_k}.$$

This is an error correction algo that corrects the error locally.

Consider what is it doing.

when  $w_k^T x_k \leq 0$  &  $y_k = 1$ ,

$$\begin{aligned} w_{k+1}^T x_k &= w_k^T x_k + x_k^T x_k \\ &\geq w_k^T x_k \end{aligned}$$

illy, when  $w_k^T x_k \geq 0$  &  $y_k = 0$

$$w_{k+1}^T x_k \leq w_k^T x_k.$$

Thus, it is putting the hyperplane towards 'correctness' at every sample. (Intuitive)

~~the~~ However, there is no guarantee that  $w_{k+1}^T x_k$  has the correct sign

Also, correcting  $w_k$  to take care of  $x_k$  may misclassify some other feature vector,

However, it works always!

Claim: The perceptron algorithm converges to ~~the~~ (finds a separating hyperplane) in finite no of steps if it exists.

Proof: Let us multiply all  $x_i$  in training set with  $y_i = 0$  by  $-1$ . (for the sake of ease of algebra).

Now, if  $w$  is a separating hyperplane if  $w^T x_i > 0 \quad \forall i$ .

under this notation, perceptron algo is  $w_{k+1} = w_k + x_k \quad \forall k \text{ s.t. } w_k^T x_k \leq 0$ .

An iteration is counted only when there is an update.

with this perceptron algo is

$w_{k+1} = w_k + x_k, w_k^T x_k \leq 0, k = 0, 1, \dots$   
stop when  $w^*$  is found.  
Goal: To show  $k$  is bounded.

Proof by Contradiction: Assume that the algo fails to find a separating hyperplane.

If algo fails to find a separating hyperplane then we must have

$$\omega_k x_k^T \leq 0, \quad \forall k \quad \left[ \text{note that } k \text{ is only when corrections are made} \right]$$

we have,

$$\omega_{k+1} = \omega_k + x_k \quad \forall k.$$

$$\begin{aligned} \|\omega_{k+1}\|^2 &= \|\omega_k + x_k\|^2 \\ &= \|\omega_k\|^2 + \|x_k\|^2 + 2\omega_k^T x_k \\ &\leq \|\omega_k\|^2 + \|x_k\|^2 \quad - (1) \end{aligned}$$

$$\therefore \omega_k^T x_k \leq 0 \quad \forall k.$$

recurring on (1),

$$\|\omega_k\|^2 \leq \|\omega_0\|^2 + \sum_{i=0}^{k-1} \|x_i\|^2$$

without loss of gen, assume  $w_0 = 0$ .

Let  $M = \max_i \|x_i\|^2$

$$\begin{aligned}\|w_k\|^2 &\leq \|w_0\|^2 + \sum_{i=0}^{k-1} \|x_i\|^2 \\ &\leq kM.\end{aligned}$$

Also,  $w_{k+1} = w_k + x_k \quad \forall k,$

$$w_k = w_{k-1} + x_{k-1}$$

recurring,

$$w_k = w_0 + \sum_{i=0}^{k-1} x_i$$

$\{w_k \text{ is a linear combo of } x_i\}$

since data is linearly separable,

$$\exists w^*, \text{ s.t. } x_i^T w^* > 0 \quad \forall i.$$

$$\text{Let } \gamma = \min_i x_i^T w^*$$

$$w_k^T w^* = \left( \sum_{i=0}^{k-1} x_i \right)^T w^* \geq k\gamma > 0.$$

$$k^2 \gamma^2 \leq |W_k^T W^*|^2$$

$$\leq \|W_k\|^2 \|W^*\|^2$$

$$\leq \|W^*\|^2 kM$$

The above should be true for all  $k$  if the algo keeps updating  $W$ .

$$\text{or } k^2 \gamma^2 \leq \|W^*\|^2 kM$$

However this can be true only till,

$$k \leq \underbrace{\frac{\|W^*\|^2 M}{\gamma^2}}_{\text{finite}}$$

Thus, the algo always finds a separating hyperplane in finitely many iterations.

This is the first provably 'correct' ML algorithm.

Perception can also be learned in the batch vs incremental versions.

$$\text{Let } S_k = \{j : w_k^T x_j \leq 0\}$$

then perception algo:

$$w_{k+1} = w_k + \sum_{j \in S_k} x_j$$

Perception can also be looked at from a risk minimization perspective.

$$\begin{aligned} \text{Let } L(y, w^T x) &= 0 \text{ if } \cancel{y w^T x} \text{ sign}(w^T x) > 0 \\ &= w^T x \text{ if } \text{sign}(w^T x) < 0 \end{aligned}$$

Then empirical risk

$$R(h) = \sum_{i=1}^n L(y_i, w^T x_i)$$

$$= \sum_{i: w^T x_i \leq 0} w^T x_i$$

now, Let's do a gradient descent on the cost

$$\begin{aligned}\omega_{k+1} &= \omega_k - \eta \nabla L(\omega_k) \\ &= \omega_k + \eta \sum_{j: \omega^T x_j \leq 0} x_j\end{aligned}$$

which is exactly the batch version of perceptron.

Thus, changing objective functions result  
(6n)  
in different algorithms.

---



## The idea of gradient descent.

We saw: All ML problems involve minimization of an objective function [often Empirical risk].

While closed-form solution can be obtained in a few cases its often not the case:

Risk is often non-convex & high-dimensional.

Thus, its a good idea to minimize the Cost (risk) iteratively.

In other words, change the parameters ( $w$ ) such that every step reduces the Cost.

This is the basic idea behind gradient descent.

Question: which direction to move towards in each iteration?

Let  $R(w)$  be the risk function.

Objective: Find a direction of  $w$  in which  $R$  decreases the fastest (most).

We know that directional derivative of a function is the slope of the fn along that dire

Let  $w$  denote a unit vector. we want

$$w^* = \min_{w, w^T w = 1} w^T \nabla_w R(w)$$

$$= \min \|w\| \|\nabla_w R(w)\|_2 \cos \theta.$$

min will occur when  $\theta = \pi$ .  $\theta \rightarrow$  angle b/w  $w$  &  $\nabla_w R(w)$

Thus,  $w^*$  has to point to a direction ~~opposite~~ to the gradient.

Thus,  $w_{k+1} = w_k - \eta \nabla_w R(w)$  is an iterative procedure one can use to minimize the risk: one can use this in linear least sq. as well.

## Fisher linear discriminant.

A linear disc fn based Classifier:

$$\text{Decide } x \in C-1 \text{ if } W^T x + w_0 > 0$$

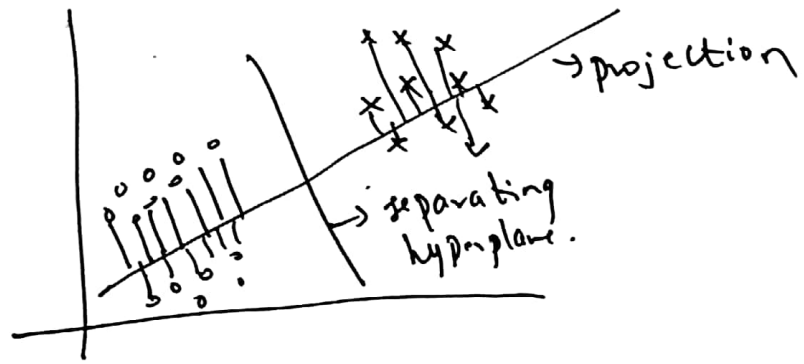
we saw OLS as a way to compute  $W$ .

~~But~~  $g(x) = W^T x + w_0$  is basically a projection of data ( $x$ ) on to a 1-D hyperplane.

In the linearly separable case, we saw that there can exist infinite separating hyperplanes.

Question: Is there a way to quantify the goodness of projected  $W$ ?

One way: separation b/w projected points of different classes should be 'more', FLD formalizes this notion.



FLDA idea: find a  $W$  s.t training data is well-separated along this  $W$ .

Consider 2-class case.

let  $D = \{ (x_i, y_i), i=1, \dots, n \}$   $y_i \in \{0, 1\}$ .

Let  $C_0$  &  $C_1$  denote two classes.

let  $n_0$  &  $n_1$  be the no of examples of each class.

let  $M_0$  &  $M_1$  be the means of data from  $C_0$  &  $C_1$  res.

(146).

$$M_0 = \frac{1}{n_0} \sum_{x_i \in C_0} x_i ; \quad M_1 = \frac{1}{n_1} \sum_{x_i \in C_1} x_i$$

Let  $m_0$  &  $m_1$  be the means of the projected data.

$$\text{i.e., } m_0 = W^T M_0 \quad \& \quad m_1 = W^T M_1, \quad [\text{all these are vectors}]$$

Intuitively,  $(m_0 - m_1)$  gives a measure of the separation b/w the samples of 2 classes after projection onto  $W$ .

Thus, we want a  $W$  that maximizes

$$(m_0 - m_1)^2$$

only having this in the cost has 2 problems

1) scaling the data would change  $W$ .

2) The intra-class variability could be lost  
[trivially by projecting all points in a class to a single point].  
(16.7)

Thus, let us define

$$S_0^2 = \sum_{x_i \in C_0} (W^T x_i - m_0)^2;$$

$$S_1^2 = \sum_{x_i \in C_1} (W^T x_i - m_1)^2$$

$S_0^2$  &  $S_1^2$  quantifies the projected variances

$\therefore S_0^2 + S_1^2$  is a measure of total within cluster variance.

Thus,  $\frac{(m_0 - m_1)^2}{S_0^2 + S_1^2}$  quantifies the separability w.r.t to the total variance of the projected data.

Thus  $L(w)_{FLDA} = \frac{(m_1 - m_0)^2}{S_0^2 + S_1^2}$

$$(m_1 - m_0)^2 = (W^T M_1 - W^T M_0)^2$$

$$= W^T (M_1 - M_0) (M_1 - M_0)^T W$$

$$= W^T S_B W$$

$$S_B = (M_1 - M_0)(M_1 - M_0)^T$$

dxd

Between class scatter matrix.

Now,  $S_0^2 = \sum_{x_i \in C_0} (W^T x_i - W^T M_0)^2$

$$= \sum_{x_i \in C_0} [W^T (x_i - M_0)]^2$$

$$= \sum_{x_i \in C_0} W^T (x_i - M_0) (x_i - M_0)^T W$$

$$= W^T \left[ \sum_{x_i \in C_0} (x_i - M_0) (x_i - M_0)^T \right] W$$

(149) ,

$$\text{Hly, } S_1^2 = W^T \left[ \sum_{x_i \in C_1} (x_i - M_1)(x_i - M_1)^T \right] W$$

$$\text{Thus, } S_0^2 + S_1^2 = W^T S_W W \text{ where}$$

$$S_W = \sum_{\text{dxd } x_i \in C_0} (x_i - M_0)(x_i - M_0)^T + \sum_{x_i \in C_1} (x_i - M_1)(x_i - M_1)^T$$

within class scatter matrix.

$$\therefore L(W)_{FLDA} = \frac{W^T S_B W}{W^T S_W W} \rightarrow \begin{array}{l} \text{max b/w class} \\ \text{scatter \&} \\ \text{min within class} \\ \text{scatter.} \end{array}$$

Note that  $L$  is not affected by scale of  $W$ .

$$\text{Now, } W_{FLDA}^* = \arg \min L(W)$$

$$\Rightarrow \frac{\partial L}{\partial W} = 0 \Rightarrow \frac{2 S_B W}{W^T S_W W} - \frac{W^T S_B W}{(W^T S_W W)^2} \cdot 2 S_W W = 0$$

$$\Rightarrow S_W W = \lambda S_B W \text{ for some } \lambda \in \mathbb{R}.$$

(150).



$$S_W W = \lambda S_B W \quad Ax = \lambda Bx$$

Generalized EV problem - solved using LU decomposition.

However, often  $S_W$  is invertible [symmetric] in which case

$$\underline{W = \lambda S_W^{-1} S_B W}$$

$$\text{Also, } S_B W = (M_1 - M_0)(M_1 - M_0)^T W = K (M_1 - M_0)$$

$$K = (m_1 - m_0)$$

$$\therefore \underline{W = S_W^{-1} (M_1 - M_0)}$$

The final classifier would be  $\text{sign}(W^T x + b)$ .

## FLDA as a special case of OLS.

Given original training data,

~~at~~ Let's construct a new training data as follows,

$$y_i' = \frac{n}{n_0} \text{ if } y_i = 0 \text{ \& } y_i' = -\frac{n}{n_1} \text{ if } y_i = 1$$

Now, let's get  $\hat{y} = W^T x + b$  using LLS.

It can be shown that LLS gives FLDA.