

## \* Empirical Risk Minimization

A lot of ML algorithms involve minimization of a function over parameters  $w$  such that it contains two components: loss function and regularization

$$\min_w \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(h_w(x_i), y_i)}_{\text{loss function}} + \underbrace{\lambda F(w)}_{\text{regularizer}} \right]$$

Loss function is a continuous function which penalizes training error, while regularizer penalizes model complexity.  $\ell$  is a function of parameters & data points, while regularizer is function of only parameters

Empirical Risk Minimization is the idea where we want to decrease loss on training data, but keep model simple.

### [I] Binary classification loss function

(1.1) Hinge loss  $\max \left[ 1 - h_w(x_i) y_i, 0 \right]^p$

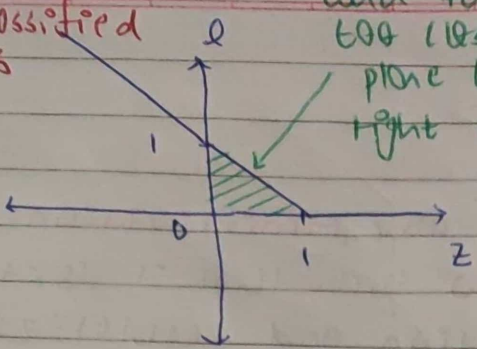
This loss is used in SVMs, if  $p=1$ ; standard SVM, if  $p=2$ ; squared hingeless SVM and becomes differentiable

The standard SVM represents the size of error in cases when point is very near to hyperplane at an wrong side.

linear penalty for misclassified points

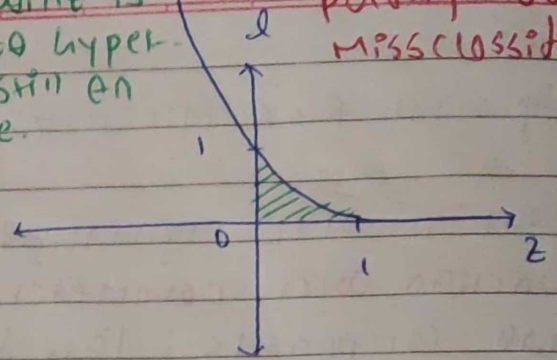
This region represents when the point is too close to hyperplane but still on right side.

squared penalty for misclassified points



$P = 1$

$\max(1-z, 0)$



$P = 2$

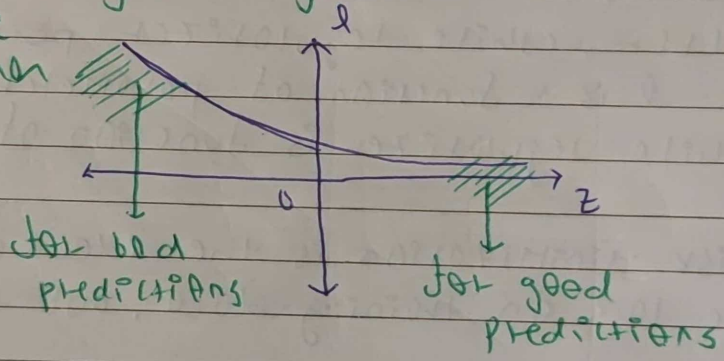
$[\max(1-z, 0)]^2$

(1.2) logistic loss

$\log(1 + e^{-y_i h_w(x_i)})$

Used in logistic regression

This is more straight line than GCV

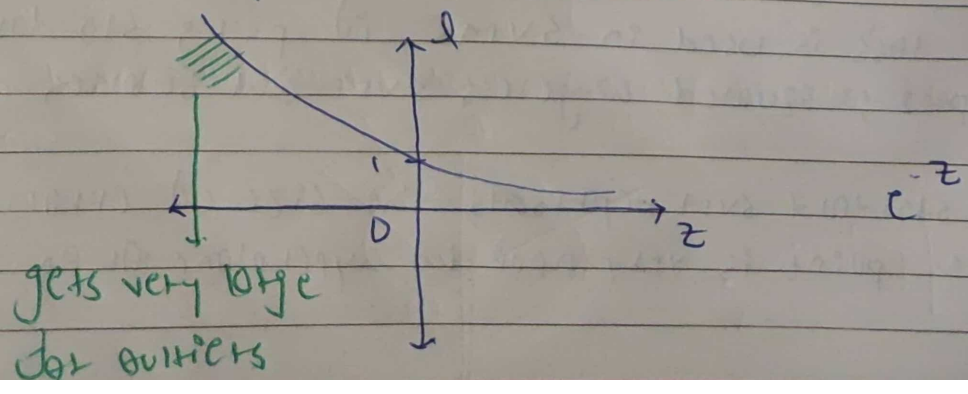


$\log(1 + e^{-z})$

(1.3) Exponential loss

$e^{-y_i h_w(x_i)}$

This is very aggressive loss function and it penalizes outliers heavily. Used in **adaboost**.

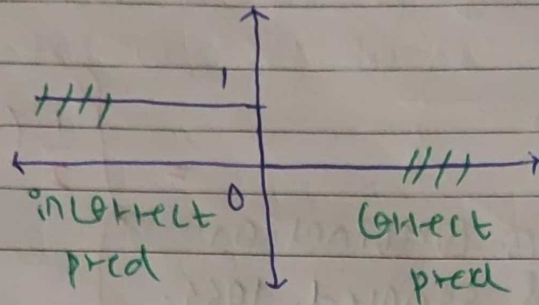




(1.4) 0/1 loss

$$\delta [\text{sign}(h_w(x_i)) \neq y_i]$$

This is used for actual classification loss.



$$\delta [\text{sign}(z), y_i]$$

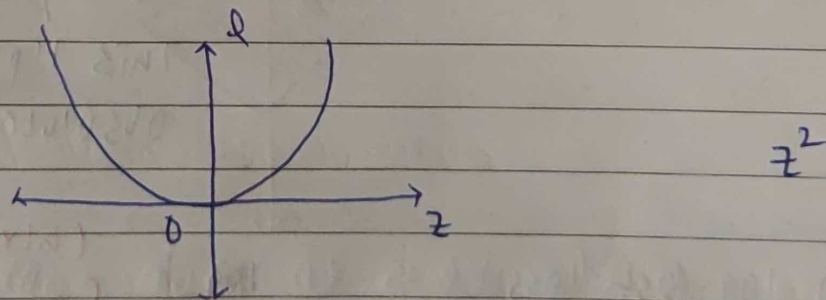
OR

$$\delta [y_i \cdot \text{sign}(z) \neq 1]$$

## [II] Regression loss functions

(2.1) Squared loss (OLS)  $(h(x_i) - y_i)^2$

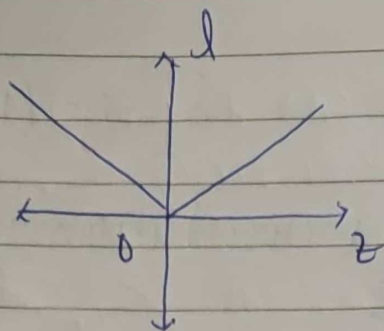
Squared loss is a measure of average/mean. It is differentiable & easy to optimize, but is sensitive to outliers. This is NOT ideal for house-price prediction OR average placement etc since outliers dominate.



(2.2) Absolute loss

$$|h(x_i) - y_i|$$

This is a measure of median. It is NOT differentiable at 0, but is less sensitive to outliers making it good for house price, median salary.

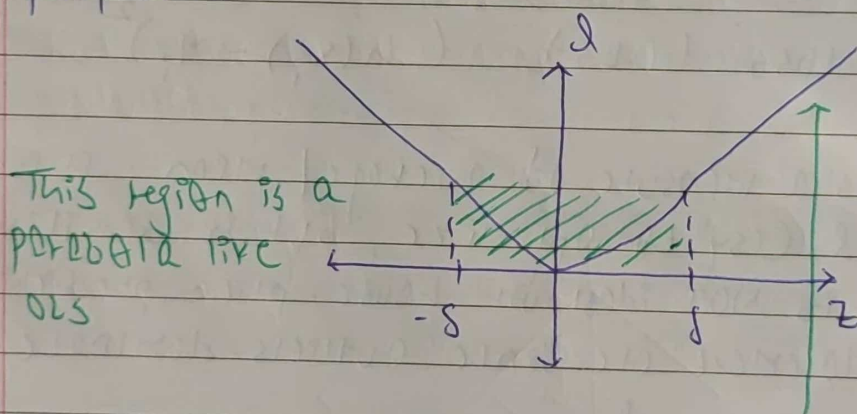


$$|z|$$

### (2.3) Huber loss

Huber loss provides combination of absolute loss & squared loss. This is differentiable, resistant to outliers and has optimization properties.

$$= \begin{cases} \frac{1}{2} (h(x_i) - y_i)^2 & \text{if } |h(x_i) - y_i| < \delta \\ \delta (|h(x_i) - y_i| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$$



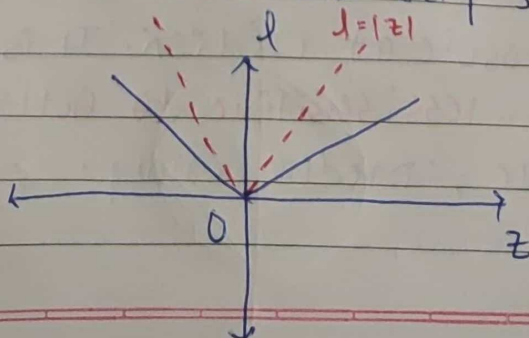
$$= \begin{cases} \frac{1}{2} z^2 & \text{for } |z| < \delta \\ \delta (|z| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$$

This represents straight line of absolute error

### (2.4) log loss

$$\log \left[ \frac{e^{(h(x_i) - y_i)} + e^{-(h(x_i) - y_i)}}{2} \right]$$

This is twice differentiable and very similar to Huber loss.



$$\log \left[ \frac{e^z + e^{-z}}{2} \right]$$



### [III] Regularizers

what we try to do is :

$$\min_{w, b} \left[ \sum_{i=1}^n d(h_w(x_i), y_i) + \lambda f(w) \right]$$

regularization term

This can be rephrased using Lagrangian as :-

$$\min_{w, b} \left[ \sum_{i=1}^n d(h_w(x_i), y_i) \right] \text{ s.t. } f(w) \leq B$$

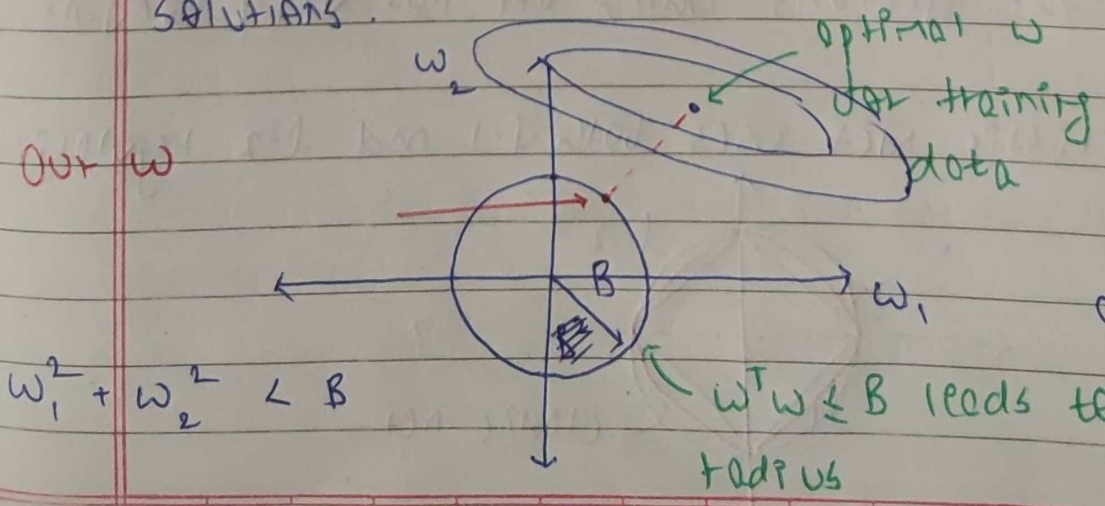
Thus, for each  $\lambda$ , there exist a B, where we can show  $f(w) \leq B$ .

Thus, what this means geometrically is that on one hand, we want the minimum of loss function & we want to **restrict w** such that  $f(w) \leq B$ .

#### (3.1) L2 regularization

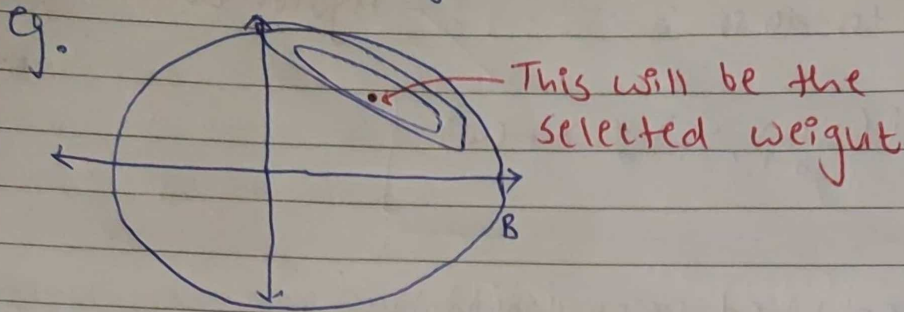
$$w^T w \text{ OR } \|w\|_2^2$$

This is differentiable & convex in nature. However, this uses weights on all features, i.e. it provides dense solutions.



what we want is  $w$  in that restricted domain only; even though optimal is outside

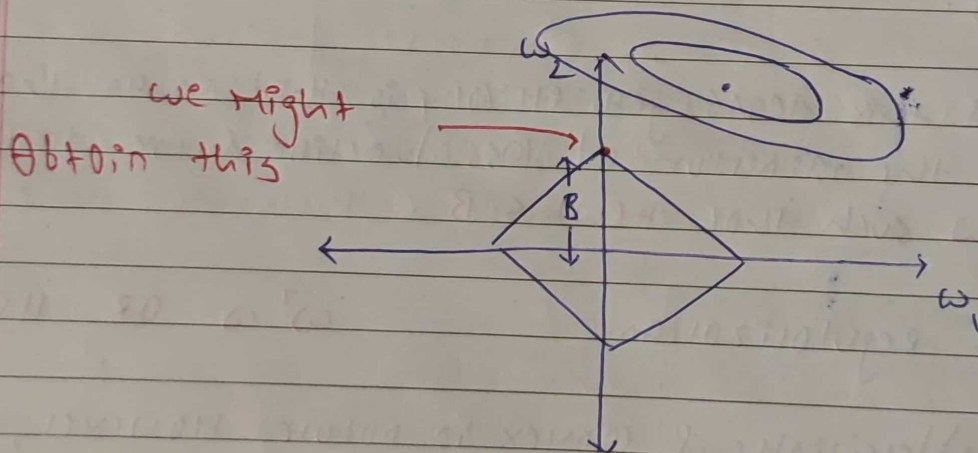
if we increase  $\lambda$  too much, i.e. decrease  $B$  too much, regularization will have no effect.



(3.2)  $L_1$  regularization.

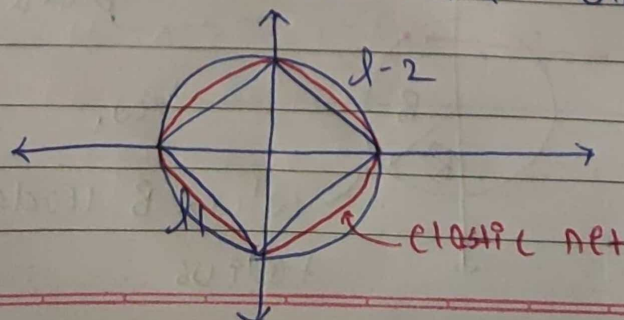
$\|w\|_1$

This is not strictly convex or differentiable at 0. But, due to its shape, it allows us to obtain sparse solutions.



This method allows to prune some weights, thus removing the features which are not very useful.

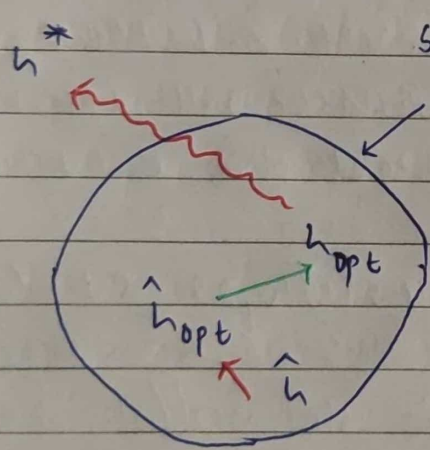
Note: Elastic Net uses both  $L_1$  and  $L_2$  regularization.





Since D1 deals with absolute values, it can afford putting all weights on one dimension and keeping others 0. This is because  $[2 \ 0 \ 0 \ 0]$  and  $[0.5 \ 0.5 \ 0.5 \ 0.5]$  leads to same sum. But D2 can't do the same. Here  $[2 \ 0 \ 0 \ 0]$  leads to 4, but  $[0.5 \ 0.5 \ 0.5 \ 0.5]$  leads to 1. Hence it has to spread out and is more robust to outliers.

- let  $h^* \rightarrow$  Optimal predictor
- $h_{opt} \rightarrow$  optimal hypothesis that can be produced by Algorithm A.
- $\hat{h}_{opt} \rightarrow$  optimal hypothesis that can be produced by Algorithm A on Data D.
- $\hat{h} \rightarrow$  hypothesis found by Algorithm A



search space ~~of  $h$~~   $H$ .

$$\epsilon(\hat{h}) - \epsilon(h^*)$$

$$= \boxed{\epsilon(h_{opt}) - \epsilon(h^*)} \text{ Approximation} + \boxed{\epsilon(\hat{h}_{opt}) - \epsilon(h_{opt})} \text{ Estimation} + \boxed{\epsilon(\hat{h}) - \epsilon(\hat{h}_{opt})} \text{ Optimization}$$

A is responsible for optimization error.

These are the 3 things where we provide **Knowledge**. For eg. Bids.

