

This formulation is similar to that of Hinge Loss.

λ : C parameter

Soft SVM rule

$$\min_{\underline{w}, b, \underline{\xi}} \left(\underbrace{\lambda \|\underline{w}\|^2}_{\text{Regularization term}} + \underbrace{\frac{1}{m} \sum_{i=1}^m \xi_{si}}_{\text{Hinge loss}} \right)$$

\wedge raised to the power - subscript

$$\min(w, b, \xi) \cdot \lambda \|\underline{w}\|^2$$

$$\frac{1}{m} \sum_{i=1}^m \xi_i$$

Overall, the soft SVM learning rule can be considered as regularized loss minimization.

Even the linear regression machine can be regularized by including a regularization term to the loss function.

$$\text{SSD} + \lambda \|\underline{w}\|^2$$

sum of squared difference.

Regularization helps in reducing the complexity of a model!

Regularization term

This is called as Tikhonov regularization.

Implementing soft SVM using Stochastic Gradient Descent

$$\min_{\underline{w}} \left(\underbrace{\frac{\lambda}{2} \|\underline{w}\|^2}_{\text{regularize}} + \underbrace{\frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \langle \underline{w}, \underline{x}_i \rangle\}}_{\text{Hinge loss}} \right) = L_S(\underline{w})$$

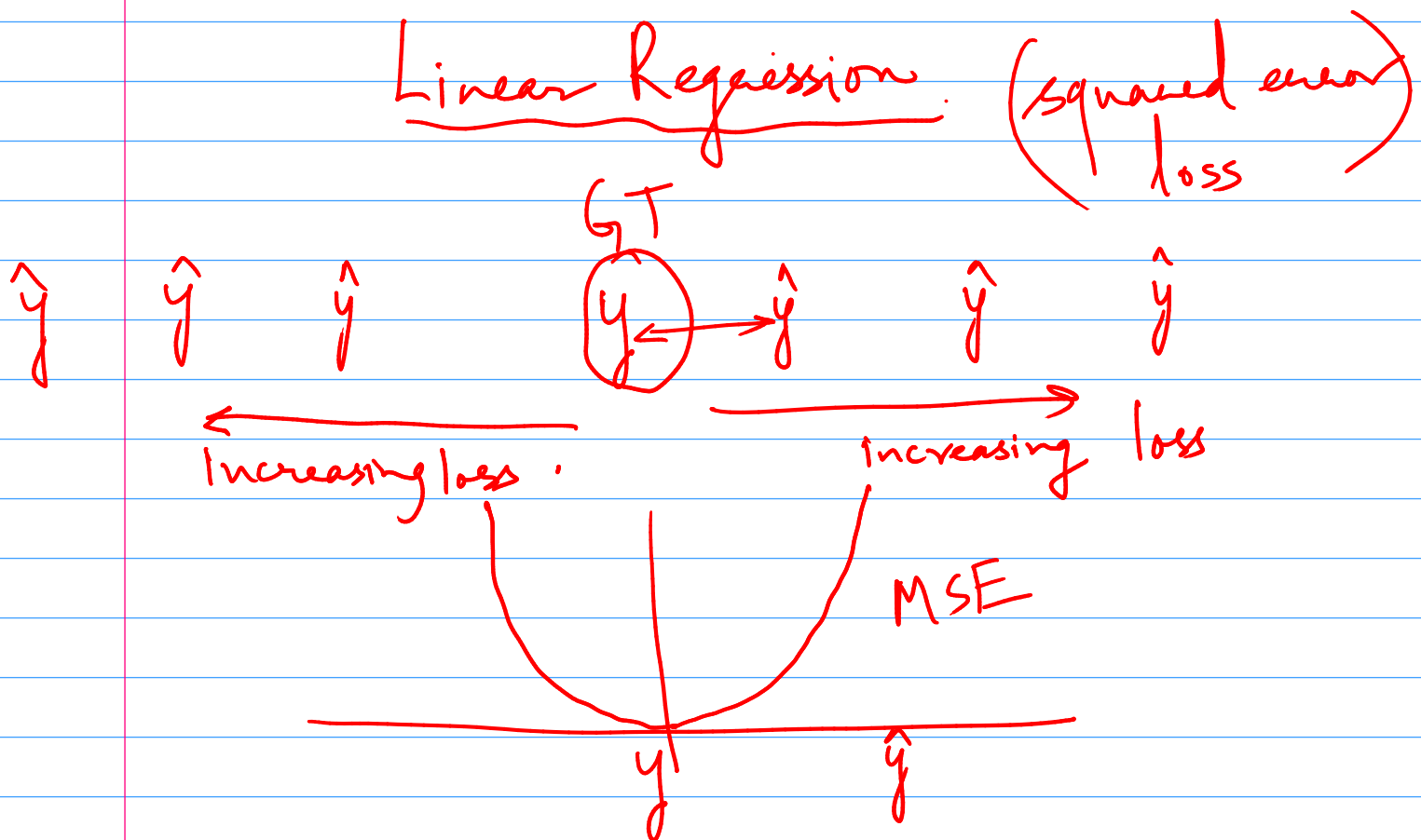
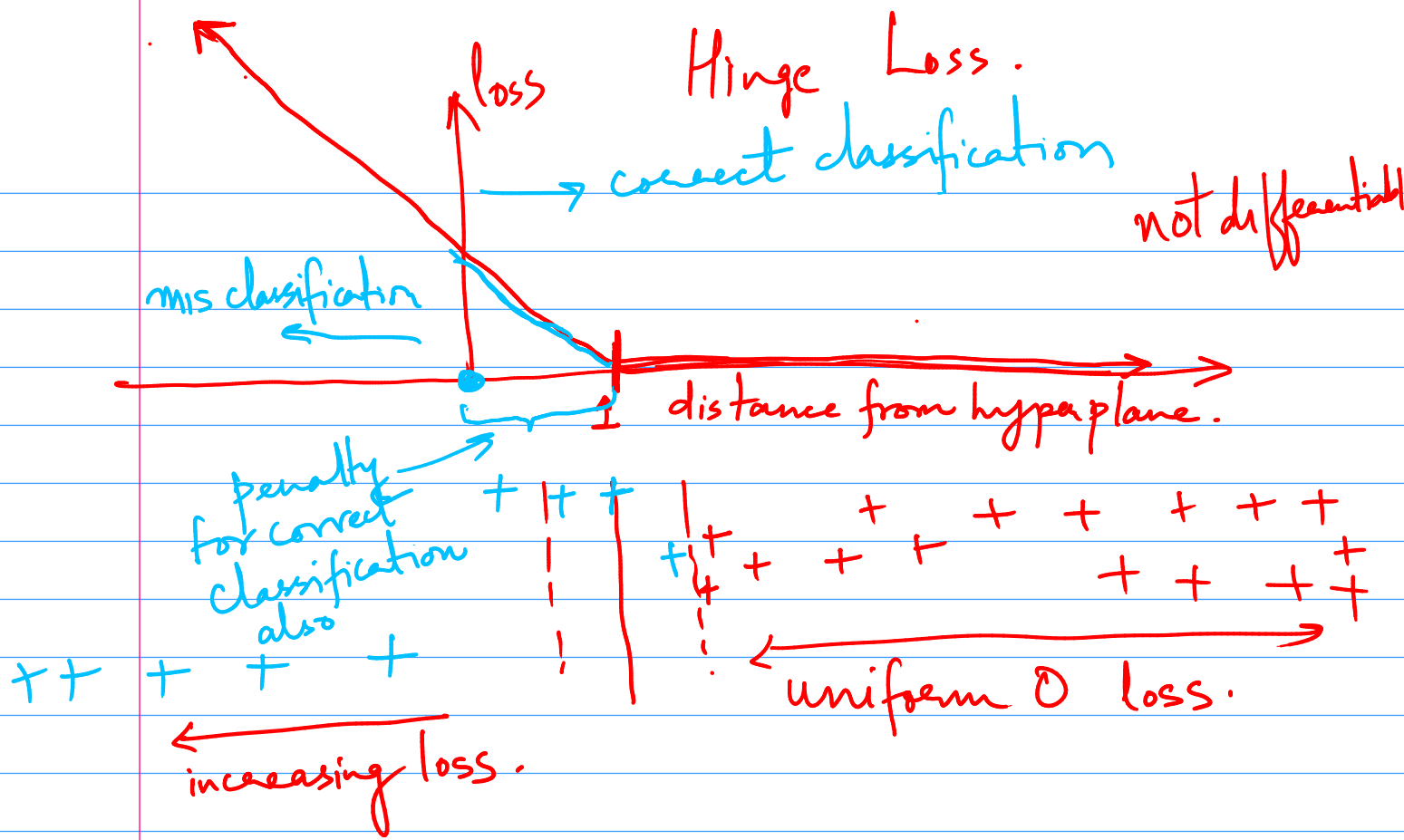
$$\min_{\underline{w}} f(\underline{w})$$

sample empirical risk

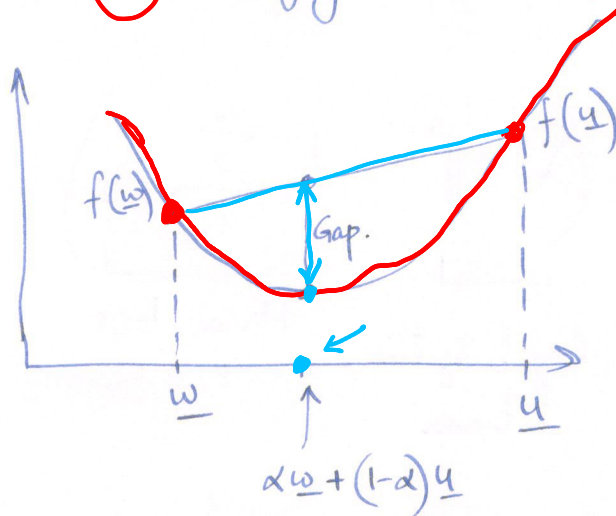
$$\text{where } f(\underline{w}) = \frac{\lambda}{2} \|\underline{w}\|^2 + L_S(\underline{w})$$

Minimize the true risk ← SGD minimizes.

expected direction is opposite to the gradient vector



f is a λ strongly convex function



For λ strongly convex functions, this $\text{gap} \geq \left(\frac{\lambda}{2}\right) \alpha(1-\alpha) \|u - w\|^2$

Gradient descent can be performed efficiently for λ strongly convex functions

parameter update

$$\underline{w}^{(t)} \leftarrow \underline{w}^{(t-1)} - \eta \Delta$$

set $\eta = \frac{1}{\lambda t}$

gradient

random SGD...

t : iteration number

Gradient of f at $\underline{w}^{(t)}$

Pick z randomly from training set

Loss for example z $l(\underline{w}^{(t)}, z)$

$L_S(\underline{w})$: loss over a sample S

$l(\underline{w}^{(t)}, z)$: loss over a single example z .

λ strongly convex function.

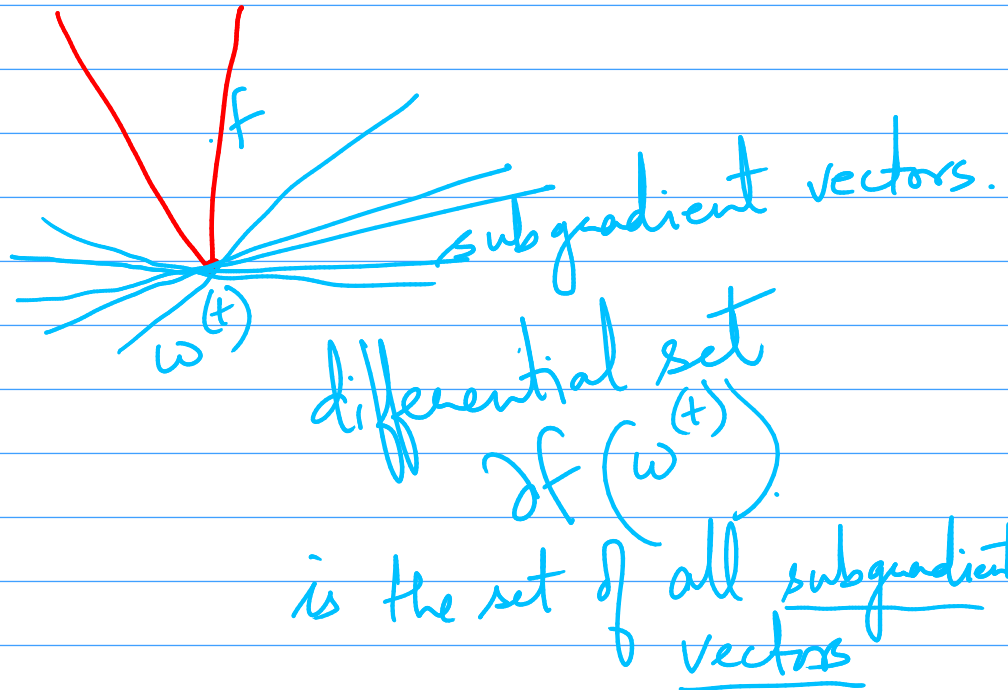
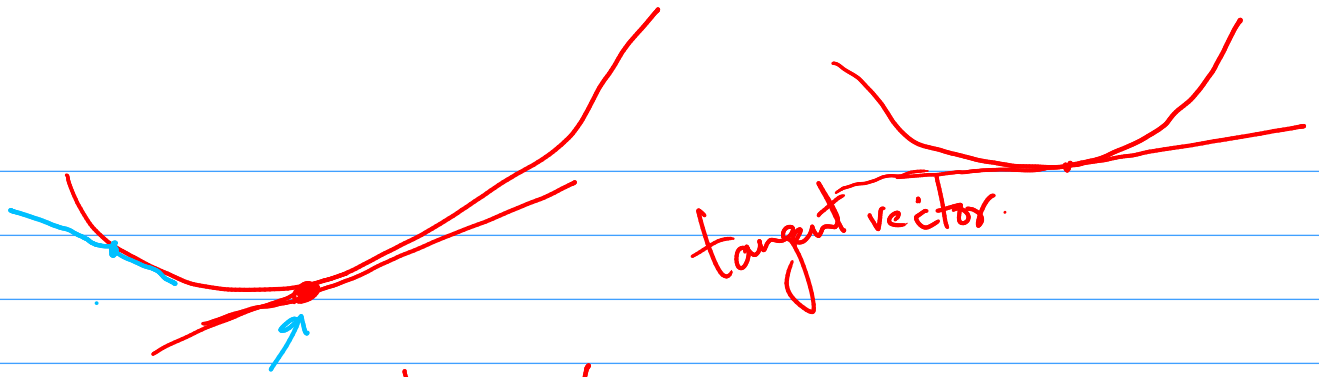
$$f(\underline{w}) = \frac{\lambda}{2} \|\underline{w}\|^2 + L_S(\underline{w})$$

$l(\underline{w}^{(t)}, z)$

The loss $l(\underline{w}^{(t)}, z)$ may not be differential.

But there can be several subgradient vectors at the point $\underline{w}^{(t)}$ for the function l .

The set of such subgradient vectors is called as the differential set denoted as $\partial l(\underline{w}^{(t)}, z)$



Let \underline{v}_t be one such subgradient vector of the loss function $l(\underline{w}^{(t)}, z)$ at point $\underline{w}^{(t)}$

subgradient vector.

$$\underline{v}_t \in \partial l(\underline{w}^{(t)}, z)$$

differential set of function l

Since $f(\underline{w}) = \frac{\lambda}{2} \|\underline{w}\|^2 + L_S(\underline{w})$

one of the subgradient vectors of $f(\underline{w})$ at $\underline{w}^{(t)}$ is

subgradient vector

$$\lambda \underline{w}^{(t)} + \underline{v}_t \in \partial f(\underline{w}^{(t)})$$

∴ the SGD (Stochastic Gradient Descent) update step

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \eta (\lambda \underline{w}^{(t)} + \underline{v}_t)$$

subgradient direction of function f .

for a strongly convex function, $\eta = \frac{1}{\lambda t}$

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \frac{1}{\lambda t} (\lambda \underline{w}^{(t)} + \underline{v}_t)$$

$$= \left(1 - \frac{1}{t}\right) \underline{w}^{(t)} - \frac{1}{\lambda t} \underline{v}_t$$

$$= \left(\frac{t-1}{t}\right) \underline{w}^{(t)} - \frac{1}{\lambda t} \underline{v}_t$$

$$= \left(\frac{t-1}{t}\right) \left(\frac{t-2}{t-1} \underline{w}^{(t-1)} - \frac{1}{\lambda(t-1)} \underline{v}_{t-1} \right) - \frac{1}{\lambda t} \underline{v}_t$$

$$\underline{w}^{(t+1)} = -\frac{1}{\lambda t} \left(\sum_{j=1}^t \underline{v}_j \right)$$

\underline{v}_j is the subgradient of the loss function at $\underline{w}^{(j)}$

$$= \begin{cases} 0 & \text{if } y \langle \underline{w}^{(j)}, \underline{x} \rangle \geq 1 \\ -yx & \text{otherwise} \end{cases}$$

1st case 2nd case 1st case

Hinge loss = $\max(0, 1 - y \langle \underline{w}, \underline{x} \rangle)$

SGD for solving soft SVM

$$\text{sum}^{(1)} = 0$$

for $t = 1 \dots T$

$$\text{Let } \underline{w}^{(t)} = \left(\frac{1}{\lambda t} \right) \text{sum}^{(t)}$$

Choose $i \in [m]$

$$\text{If } (y_i \langle \underline{w}^{(t)}, \underline{x}_i \rangle < 1)$$

margin constraint is not satisfied.

$$\text{sum}^{(t+1)} = \text{sum}^{(t)} + y_i \underline{x}_i$$

else

$$\text{sum}^{(t+1)} = \text{sum}^{(t)}$$