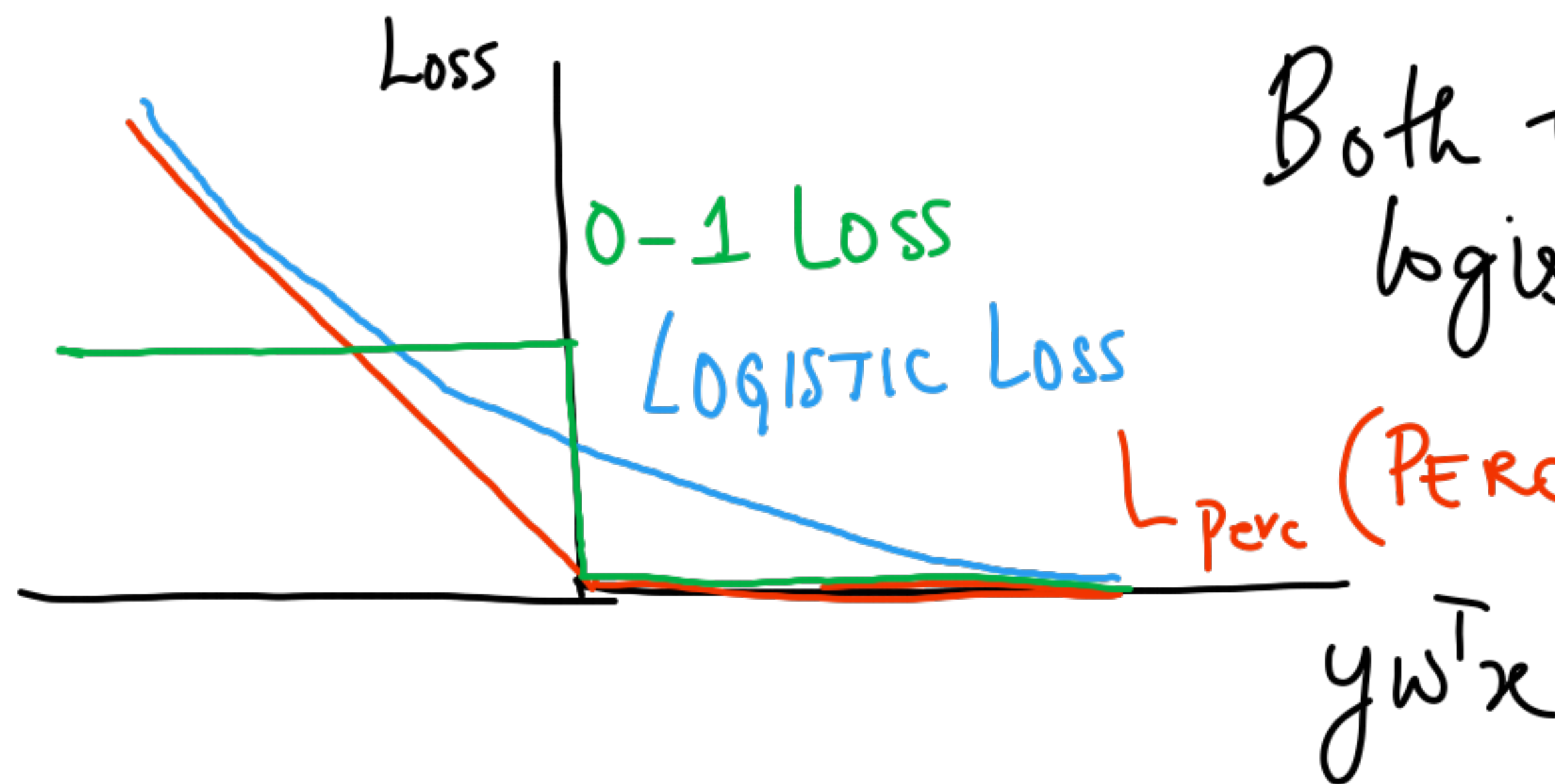


CS725

$$L_{\text{perceptron}}(x, y; w) = \max(0, -y w^T x)$$



Both the perceptron loss & the logistic loss are surrogates of the 0-1 loss

Convergence of the Perceptron Learner

Consider a linearly separable dataset \mathcal{D} , i.e. there exists a n weight vector s.t. $y = \text{sign}(\mathbf{u}^T \mathbf{x}) \forall \mathbf{x}, y \in \mathcal{D}$.

Two assumptions (without loss of generality):

① Assume \mathbf{u} is a unit vector

② Assume all the \mathbf{x} 's $\in \mathcal{D}$ lie within a Euclidean ball of radius 1
i.e., $\|\mathbf{x}\| \leq 1$

Define a new quantity called "MARGIN OF SEPARATION", γ

$$\gamma = \min_{x \in \mathcal{D}} |u^T x| \quad \left[\begin{array}{l} \text{minimum distance of a point from the} \\ \text{hyperplane } u^T x = 0 \end{array} \right]$$

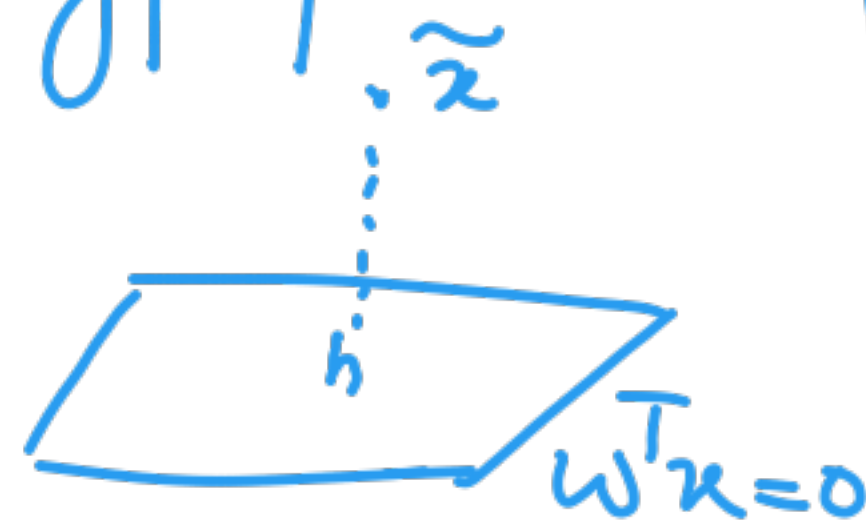
Distance of a point \tilde{x} from a hyperplane

Pick a point h on the hyperplane

$$\tilde{x} - h = k w \Rightarrow h = \tilde{x} - k w$$

$$w^T (\tilde{x} - k w) = 0 \Rightarrow k = \frac{w^T \tilde{x}}{\|w\|_2^2}$$

$$\text{Distance of } \tilde{x} \text{ from the hyperplane} = \|\tilde{x} - h\| = \frac{|w^T \tilde{x}|}{\|w\|_2}$$



Theorem of convergence: If there exists a unit vector u such that $y u^T x \geq \gamma \forall x, y \in \mathcal{D}$, then the total number of mistakes made by a perceptron learner is bounded by $\frac{1}{\gamma^2}$ when trained on \mathcal{D} .
(Assume w is initialized to $\gamma \cdot 0$.)

Proof: Let us monitor the progression of two quantities:
① $w^T u$ and ② $\|w\|_2$.

Why these quantities? We want $w^T u$ to be large but not only due to a large $\|w\|_2$. Keep track of both quantities as we iteratively estimate w .

① Monitor $w^T u$

Let w_i be the weight vector at iteration i .

$$w_{i+1}^T u - w_i^T u = (w_i + \gamma x)^T u - w_i^T u$$

$$\Rightarrow w_{i+1}^T u = w_i^T u + \gamma x^T u$$

$$\Rightarrow w_{i+1}^T u \geq w_i^T u + \gamma$$

$$\text{Assuming } w_0 = 0 \Rightarrow w_k^T u \geq k\gamma \longrightarrow \textcircled{A}$$

(B) Monitor $\|w\|_2$

$$\|w + yx\|_2^2 = \|w\|^2 + y^2\|x\|^2 + 2y w^T x$$

$$= \|w\|^2 + \|x\|^2 + 2y w^T x$$

$$\leq \|w\|^2 + \|x\|^2 \left[\because y w^T x < 0 \right]$$

$$\leq \|w\|^2 + 1 \left[\because \|x\|_2 < 1 \right]$$

After k
iterations,

$$\|w_k\|^2 \leq k$$

\longrightarrow (B)

$$\omega_k^T u \geq k\gamma \longrightarrow \textcircled{1}$$

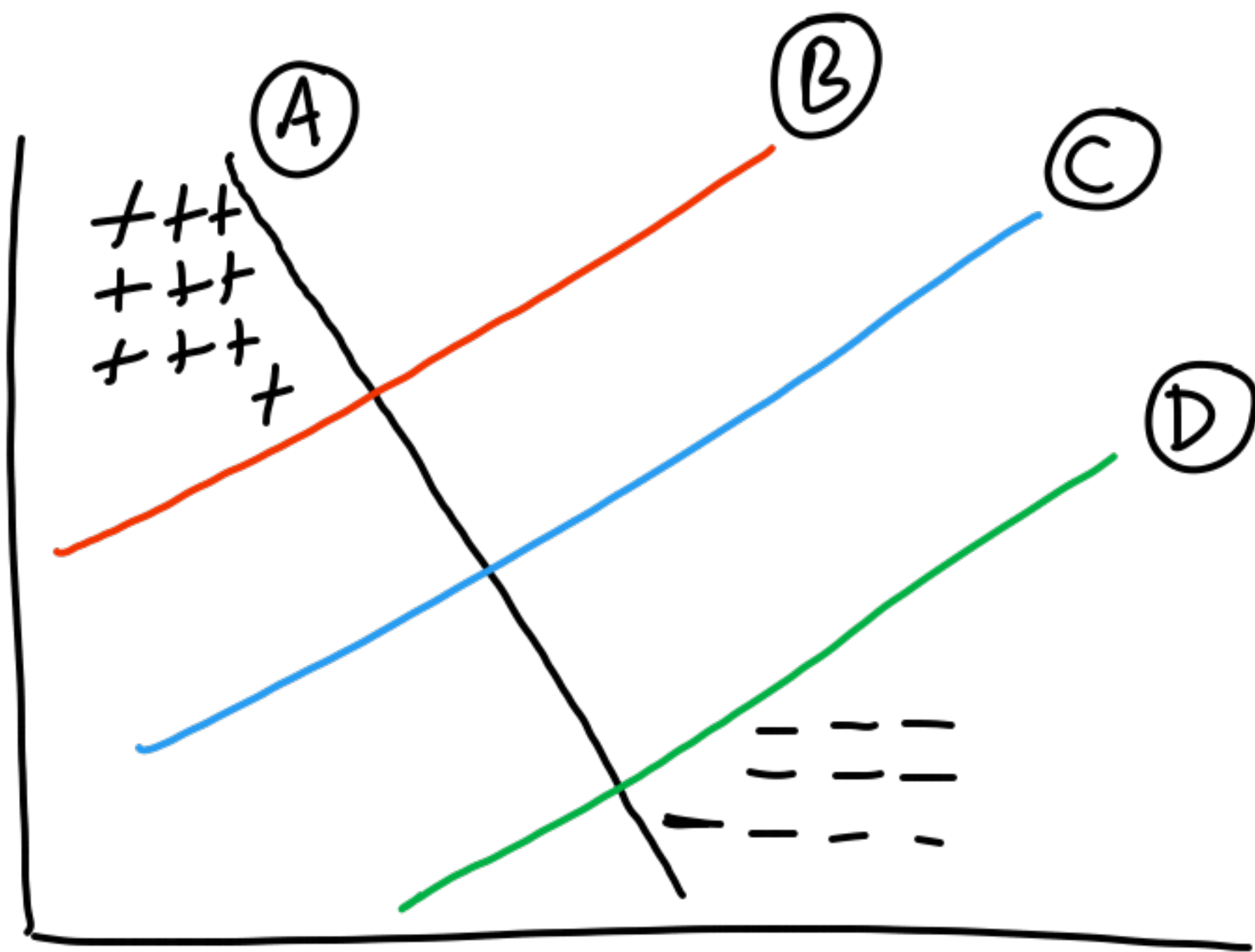
$$\|\omega_k\|^2 \leq k \longrightarrow \textcircled{2}$$

Putting $\textcircled{1}$ and $\textcircled{2}$ together.

$$\sqrt{k} \geq \|\omega_k\| \geq \omega_k^T u \geq k\gamma$$

$$\Rightarrow \sqrt{k} \geq k\gamma$$

$$\Rightarrow \boxed{k \leq \frac{1}{\gamma^2}}$$



③ is the natural pick

\Rightarrow results in faster convergence i of larger γ

\Rightarrow results in better generalization

\Rightarrow Motivates SVMs (Support Vector Machines) generalization
 find the predictor that yields the largest margin [SVMs are called max-margin classifiers]

