Convergence of Perceptron Learning Algorithm

Hsuan-Tien Lin, September 26, 2020

Page 14 of lecture 2 proves that

$$\mathbf{w}_{f}^{T}\mathbf{w}_{t+1} \geq \mathbf{w}_{f}^{T}\mathbf{w}_{t} + \underbrace{\min_{n} y_{n}\mathbf{w}_{f}^{T}\mathbf{x}_{n}}_{\|\mathbf{w}_{f}\| \cdot \rho},$$

where we defined the last term to be a constant $\|\mathbf{w}_f\| \cdot \rho$ on page 16 of lecture 2. Assume that $\mathbf{w}_0 = \mathbf{0}$, as set on page 15 of lecture 2. Then, we have

Summing all inequalities above, we get

$$\mathbf{w}_f^T \mathbf{w}_T \geq \Box \cdot \|\mathbf{w}_f\| \cdot \rho \tag{1}$$

This says the inner product grows by at least $\left| \dots \right| \cdot \|\mathbf{w}_f\| \cdot \rho$ after T updates. Note that for linear separable data with both classes of examples, a separating \mathbf{w}_f means $\|\mathbf{w}_f\| > 0$ and $\rho > 0$.

Now let's look at the results on page 15 of lecture 2. It proves that

$$\|\mathbf{w}_{t+1}\|^2 \le \|\mathbf{w}_t\|^2 + \underbrace{\max_n \|\mathbf{x}_n\|^2}_{R^2},$$

where we defined the last term to be a constant R^2 on page 16 of lecture 2. Assume that $\mathbf{w}_0 = \mathbf{0}$, as set on page 15 of lecture 2. Then, we have

$$\|\mathbf{w}_{0}\|^{2} = \square$$

$$\|\mathbf{w}_{1}\|^{2} \leq \|\mathbf{w}_{0}\|^{2} + R^{2}$$

$$\|\mathbf{w}_{2}\|^{2} \leq \square + R^{2}$$

$$\square \leq \square + R^{2}$$

$$\dots$$

$$\dots$$

$$\square \leq \|\mathbf{w}_{T-1}\|^{2} + R^{2}$$

Summing all inequalities above, we get

$$\|\mathbf{w}_T\|^2 \leq \square \cdot R^2 \tag{2}$$

This says the squared length grows by at most R^2 after T updates.

If $\|\mathbf{w}_T\|^2 = 0$, this means $\mathbf{w}_f^T \mathbf{w}_T = 0$ as well, which contradicts (1) for $T \ge 1$ because $\|\mathbf{w}_f\| > []$, $\|\rho\| > []$, and $\mathbf{w}_f^T \mathbf{w}_0 = []$. So $\|\mathbf{w}_T\|^2$ must be strictly positive for $T \ge 1$. Then, we can "divide" (1) by the "square root" of (2) to get

$$\frac{\mathbf{w}_{f}^{T}\mathbf{w}_{T}}{\|\mathbf{w}_{T}\|} \ge \boxed{\qquad} \cdot \frac{\rho \|\mathbf{w}_{f}\|}{R}.$$
(3)

That is, for linear separable data with both classes of examples, and for $T \ge 1$, let θ_T denote the angle between \mathbf{w}_f and \mathbf{w}_T , we have

$$\square \ge \cos \theta_T = \frac{\mathbf{w}_f^T \mathbf{w}_T}{\square \cdot \|\mathbf{w}_T\|} \ge \square \cdot \frac{\rho}{R}$$

This proves that T, the number of updates of PLA, cannot be more than

$$\frac{R^2}{\rho^2},$$

which is the conclusion on page 16 of lecture 2. That is, PLA will converge with no more than $\frac{R^2}{\rho^2}$ updates.