

2/12/2020.

$\mathbb{R}^d \rightarrow \mathbb{R}$

Linear Classifiers :- specified by a hyperplane (\bar{w}, t) in \mathbb{R}^d , and for a feature $\bar{x} \in \mathbb{R}^d$,

predicted label $\hat{y} = \text{sign}(\bar{w} \cdot \bar{x} - t)$.

Perceptron :- An algo for linear classifiers.

① $\bar{w} = \bar{0}, t = 0$

② REPEAT :-

for $i=1, \dots, n \leftarrow \approx (n \times 4/\gamma^2)$

if $\text{sign}(\bar{w} \cdot \bar{x}_i - t) \neq y_i$,

$$\bar{w} \leftarrow \bar{w} + y_i \bar{x}_i$$

$$t \leftarrow t - y_i$$

Theorem. The perceptron algorithm learns a perfect classifier of a separable dataset with margin γ , and $\|\bar{x}_i\| \leq 1$ with at most $4/\gamma^2$ updates.

$(\bar{w}_{\text{opt}}, t_{\text{opt}})$ optimal classifier w/ margin γ , $\|\bar{w}_{\text{opt}}\| = 1$.

- ① $\|(\bar{w}_{\text{opt}}, t_{\text{opt}})\| \leq \sqrt{2}$
- ② $(\bar{w}_j, t_j) \cdot (\bar{w}_{\text{opt}}, t_{\text{opt}}) \geq j\gamma$
- ③ $\|\bar{w}_j, t_j\| \leq \sqrt{2j}$

} we proved these last time.

For any vectors \bar{u}, \bar{v} , $\underline{\bar{u} \cdot \bar{v}} \leq \|\bar{u}\| \cdot \|\bar{v}\|$

since $\cos(\text{angle between } \bar{u}, \bar{v}) = \frac{\bar{u} \cdot \bar{v}}{\|\bar{u}\| \|\bar{v}\|} \cdot \underline{w \leq 1}$

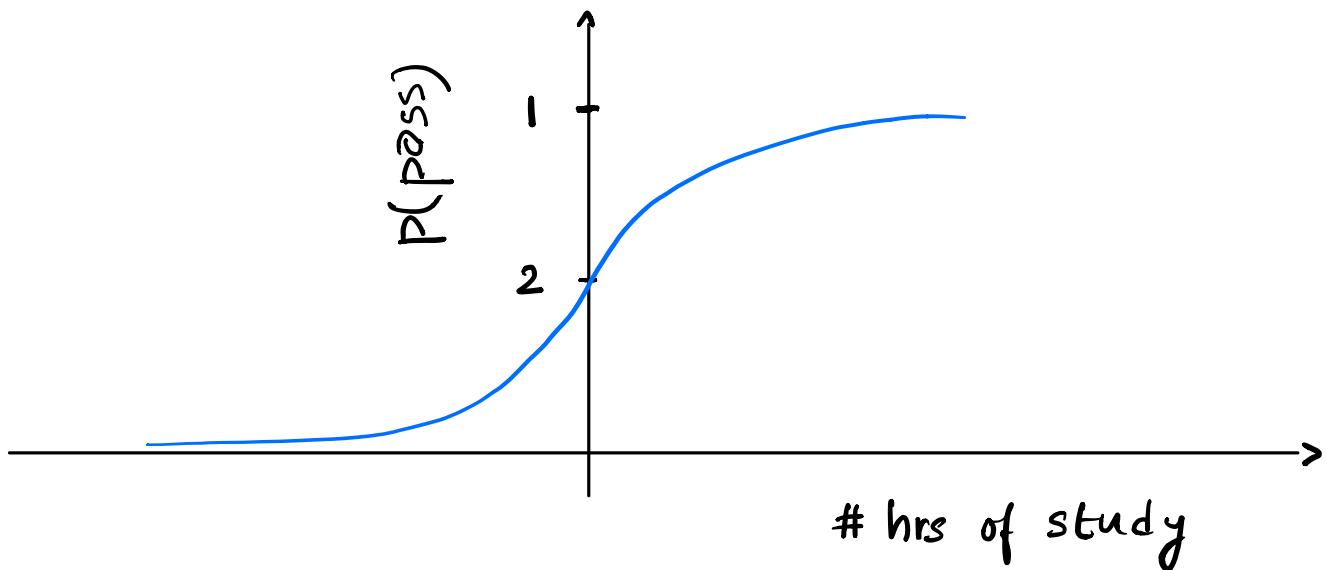
$$\Rightarrow \frac{j\gamma}{\sqrt{2j} \cdot \sqrt{2}} \leq 1 \Rightarrow j \leq \frac{4}{\gamma^2}$$

$$\underbrace{(w_j, t_j)}_{\parallel \quad \downarrow \quad \parallel} \cdot \underbrace{(w_{opt}, t_{opt})}_{\parallel \quad \downarrow \quad \parallel}$$

LOGISTIC REGRESSION.

Logistic models:- a powerful model for the probability of success given the variables labels. features

Ex:- What is the probability of passing an exam if one studies for 'x' hours?

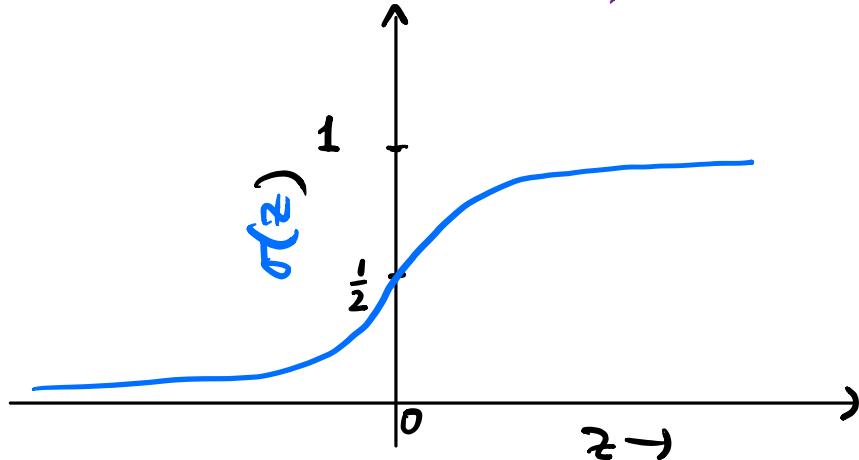


Model:- effort = $a \cdot \# \text{ hours studied} + b$

$$\sigma(\text{effort}) = \text{Prob}(\text{pass}) = \frac{e^{\text{effort}}}{1 + e^{\text{effort}}} = \frac{\exp(\text{effort})}{1 + \exp(\text{effort})}.$$

What is the plot above? it is the sigmoid function.

SIGMOID FUNCTION:- For $z \in \mathbb{R}$, $\sigma(z) := \frac{e^z}{1+e^z}$.



$$\bar{x} = (\bar{x}^1, \dots, \bar{x}^d) \quad y \in \{-1, 1\}.$$

$$(\bar{w}, t) \rightarrow ((\bar{w}^1, \dots, \bar{w}^d), t) . \quad (\bar{w}^1, t)$$

↑ ↑
a -b

$$\text{effort} = (\bar{w}^1 \cdot \bar{x}^1 + \dots + \bar{w}^d \cdot \bar{x}^d - t)$$

Logistic Regression:- Prob. of a class is - sigmoid
of a linear combination of the attributes.

- Prob $(1 | \bar{x}, (\bar{w}, t)) = \sigma(\bar{w} \cdot \bar{x} - t) = \frac{\exp(\bar{w} \cdot \bar{x} - t)}{1 + \exp(\bar{w} \cdot \bar{x} - t)}$.
- Prob $(-1 | \bar{x}, (\bar{w}, t)) = \frac{1}{1 + \exp(\bar{w} \cdot \bar{x} - t)}.$

output $\max_{y \in \{-1, 1\}} \left\{ \text{Prob}(y | \bar{x}, (\bar{w}, t)) \right\}.$

$$\stackrel{+/-}{\rightarrow} \text{Prob}(1 | \bar{x}, (\bar{w}, t)) \geq \text{Prob}(-1 | \bar{x}, (\bar{w}, t))$$

$$\Leftrightarrow \frac{\exp(\bar{w} \cdot \bar{x} - t)}{1 + \exp(\bar{w} \cdot \bar{x} - t)} \geq 1 \Leftrightarrow \bar{w} \cdot \bar{x} - t \geq 0$$

$\Rightarrow \hat{y} = \text{sign}(\bar{w} \cdot \bar{x} - t).$

(\bar{w}, t)

\bar{x}_2

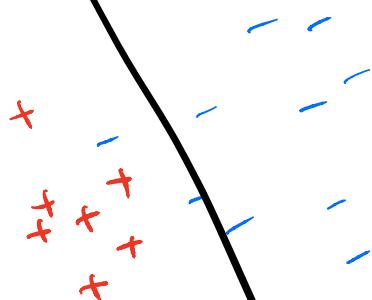
$$|\bar{w} \cdot \bar{x}_1 - t| > |\bar{w} \cdot \bar{x}_2 - t|$$

$\bar{x}_1.$

Θ

Θ

\times



$$\text{Prob}(y | \bar{x}, (\bar{w}, t)) = \left[\frac{\exp(\bar{w} \cdot \bar{x} - t)}{1 + \exp(\bar{w} \cdot \bar{x} - t)} \right]^{\frac{1+y}{2}} \left[\frac{1}{1 + \exp(\bar{w} \cdot \bar{x} - t)} \right]^{\frac{1-y}{2}}.$$

$y = +1$ $y = -1$

$$\text{Prob}(y | \bar{x}, (\bar{w}, t)) = [\underline{\exp(\bar{w} \cdot \bar{x} - t)}]^{\frac{1+y}{2}} \cdot \frac{1}{1 + \exp(\bar{w} \cdot \bar{x} - t)} \}^{\log}$$

Assumption:- All 'n' training examples in $S = \{(\bar{x}_1, y_1), \dots, (\bar{x}_n, y_n)\}$ are generated independently.

$$\begin{aligned} \text{Prob}(y_1, \dots, y_n | \bar{x}_1, \dots, \bar{x}_n) &= \prod_{i=1}^n \text{Prob}(y_i | \bar{x}_i). \\ &= \prod_{i=1}^n \text{Prob}(y_i | \bar{x}_i, (\bar{w}, t)) \end{aligned}$$

$$\max_{(\bar{w}, t)} \prod_{i=1}^n \text{Prob}(y_i | \bar{x}_i, (\bar{w}, t)).$$

- Take logs when you have product of a bunch of things.

$$\text{LLF} \Rightarrow J_S(\bar{w}, t) = \sum_{i=1}^n \log \text{Prob}(y_i | \bar{x}_i, (\bar{w}, t))$$

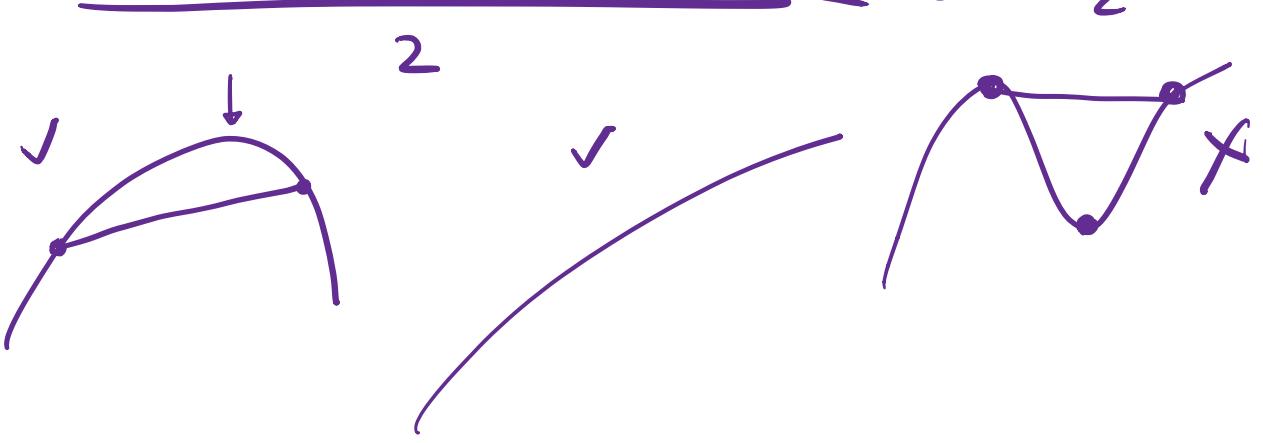
CLAIM:- $\max_{(\bar{w}, t)} \prod_{i=1}^n \text{Prob}(y_i | \bar{x}_i, (\bar{w}, t)) = \max_{\bar{w}, t} J_S(\bar{w}, t)$

$$J_S(\bar{w}, t) = \sum_{i=1}^n \left[\left(\frac{1+y_i}{2} \right) \cdot \underbrace{(\bar{w} \cdot \bar{x}_i - t)}_{-} - \log \left(1 + \exp(\bar{w} \cdot \bar{x}_i - t) \right) \right]$$

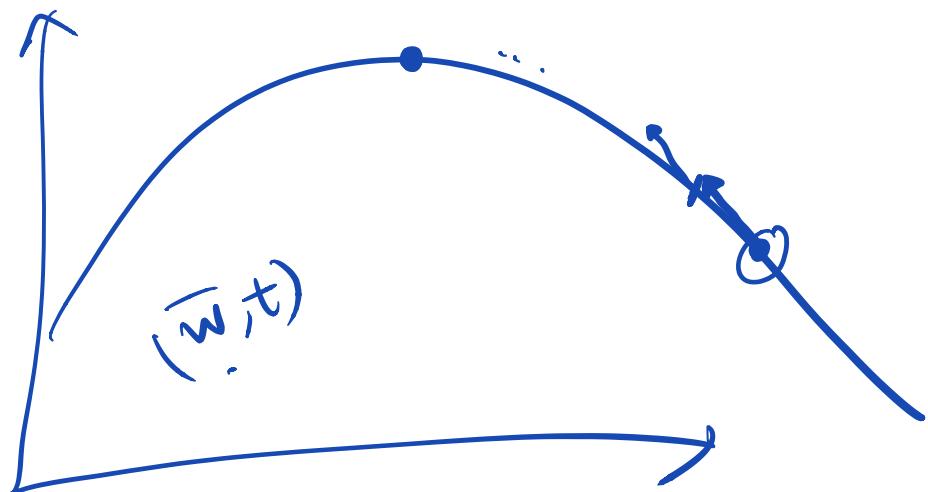
$\log(\exp(x)) = x.$

CLAIM:- $J_S(\bar{w}, t)$ is a concave function of (\bar{w}, t) .

$$J_S(\bar{w}_1, t_1) + J_S(\bar{w}_2, t_2) \leq J\left(\frac{\bar{w}_1 + \bar{w}_2}{2}, \frac{t_1 + t_2}{2}\right)$$



- EXERCISE :- Show $f(x) = \log(1 + \exp(x))$ is convex.



$\nabla J_s(\bar{w}, t)$ • Start with (\bar{w}_0, t_0) ,
 until convergence :-
 $(\bar{w}_j, t_j) \leftarrow (\bar{w}_{j-1}, t_{j-1}) + \eta \cdot \nabla J_s(\bar{w}, t) \Big|_{(\bar{w}_{j-1}, t_{j-1})}$

$$\nabla J_s(\bar{w}, t) = \left[\frac{\partial J(\bar{w}, t)}{\partial \bar{w}^1}, \dots, \frac{\partial J(\bar{w}, t)}{\partial \bar{w}^d}, \frac{\partial J(\bar{w}, t)}{\partial t} \right].$$

$$J_s(\bar{w}, t) = \sum_{i=1}^n \left[\underbrace{\left(\frac{1+y_i}{2} \right) \cdot \underbrace{(\bar{w} \cdot \bar{x}_i - t)}_{\bar{w}^1 \cdot \bar{x}_i^1} - \log \left(1 + \exp(\bar{w} \cdot \bar{x}_i - t) \right)}_{\text{should not be there.}} \right]$$

$$\frac{d}{dx} \log(\frac{1}{1+\exp(x)}) = \frac{1}{1+\exp(x)} \cdot \exp(x) \cdot \cancel{x} = \sigma(x) \cdot \cancel{x}$$

- $\bar{w} \cdot \bar{x} - t$, $\bar{x}^* = (\bar{x}^1, \dots, \bar{x}^d, -1)$

$\vdash (\bar{w}, t) \cdot \bar{x}^*$

$$\nabla J_s(\bar{w}, t) = \sum_{i=1}^n \left[\left(\frac{1+y_i}{2} \right) \cdot \bar{x}_i^* - \frac{\exp(\bar{w} \cdot \bar{x}_i - t)}{1 + \exp(\bar{w} \cdot \bar{x}_i - t)} \cdot \bar{x}_i^* \right] \leftarrow$$

$$= \sum_{i=1}^n \left[\left(\frac{1+y_i}{2} \right) \cdot \bar{x}_i^* - P(t| \bar{x}_i^*, \bar{w}, t) \cdot \bar{x}_i^* \right] \bar{x}_i^*$$

EXERCISE :- derive \downarrow from 1st principle (defn).