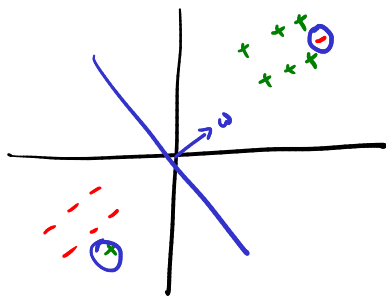
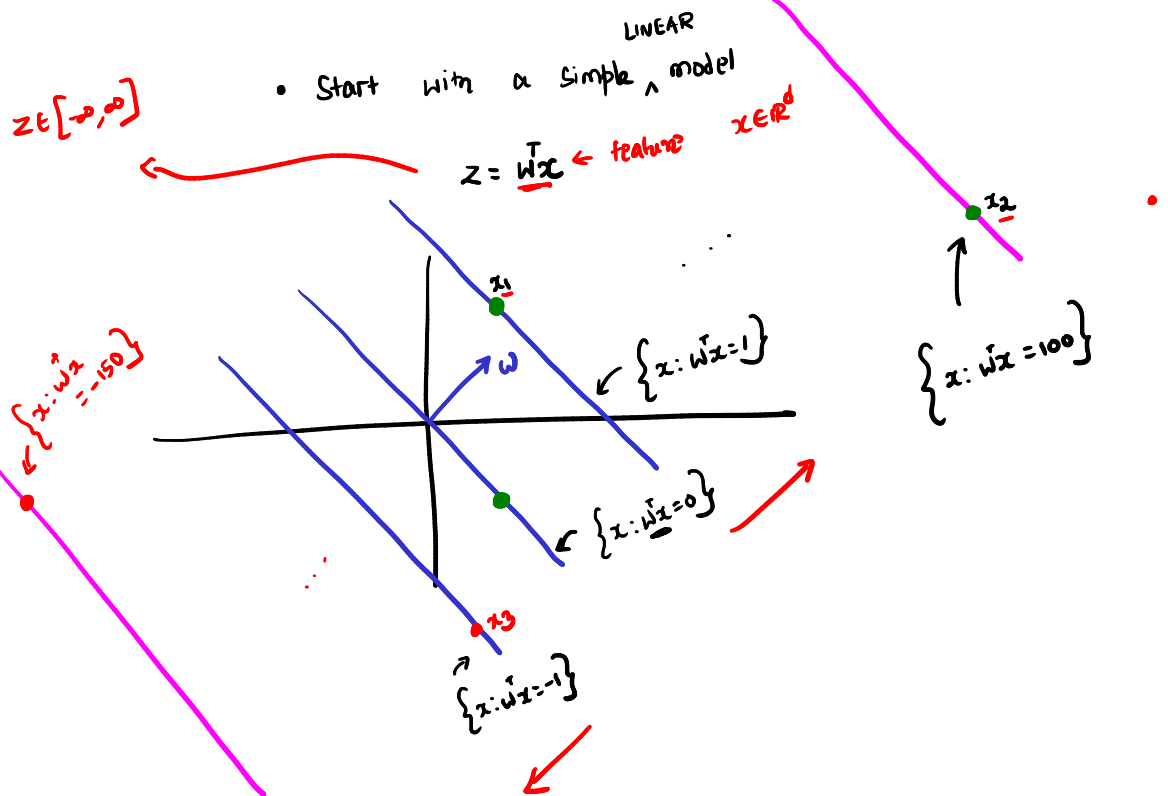


Perceptron mistakes $\leq \frac{R^2}{\gamma^2}$ ← Radius
 ← margin



$$P(y=1|x) = \begin{cases} 1 & \text{if } w^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Can we model probabilities differently?



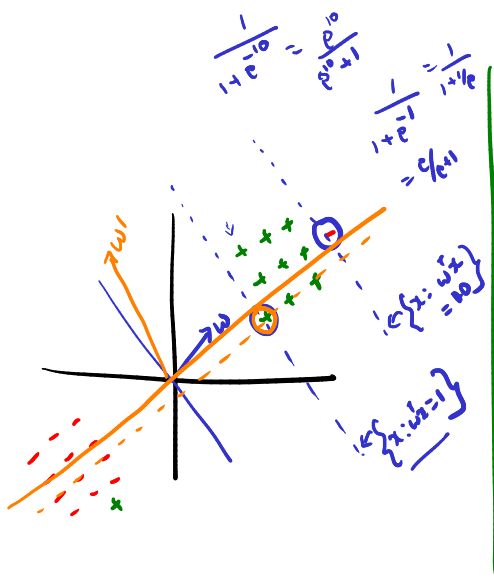
- Larger the score ($z = w^T x$), more the probability of being +1

LINK FUNCTION

Score
 \downarrow
 $g(z) = 0.5$ if $z = 0$

$g(z) \rightarrow 1$ as $z \rightarrow \infty$

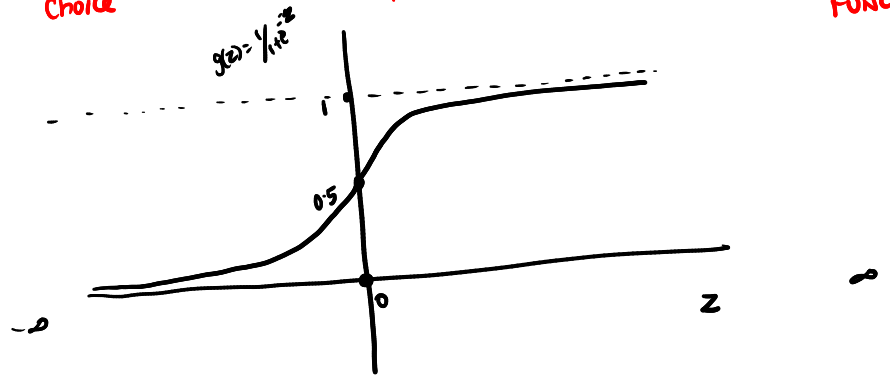
$g(z) \rightarrow 0$ as $z \rightarrow -\infty$



one popular choice

$$g(z) = \frac{1}{1 + e^{-z}}$$

SIGMOID / LOGISTIC FUNCTION



MODEL: LOGISTIC REGRESSION

$$P(y=1/x) = \frac{1}{1 + e^{-w^T x}} = g(w^T x)$$

Dataset: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ $y_i \in \{0, 1\}$

How to find w : Maximum Likelihood.

$$L(w; \text{Data}) = \prod_{i=1}^n (g(w^T x_i))^{y_i} (1 - g(w^T x_i))^{(1-y_i)}$$

$$\begin{aligned} \log L(w; \text{Data}) &= \sum_{i=1}^n y_i \log(g(w^T x_i)) + (1-y_i) \log(1 - g(w^T x_i)) \\ &= \sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-w^T x_i}}\right) + (1-y_i) \log\left(\frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}}\right) \\ &= \sum_{i=1}^n \left[(1-y_i)(-w^T x_i) - \log(1 + e^{-w^T x_i}) \right] \end{aligned}$$

Goal:

$$\max_w \sum_{i=1}^n \left[(1-y_i)(-w^T x_i) - \log(1 + e^{-w^T x_i}) \right]$$

$\log L(w)$

- No closed form expression
- Can perform Gradient descent. [ascent]

$$\begin{aligned}
 \nabla \log L(w) &= \sum_{i=1}^n \left[(1-y_i)(-x_i) - \left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}} \right) (-x_i) \right] \\
 &= \sum_{i=1}^n \left[-x_i + y_i x_i + x_i \left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}} \right) \right] \\
 &= \sum_{i=1}^n y_i x_i - x_i \left(\frac{1}{1+e^{w^T x_i}} \right) \\
 &= \sum_{i=1}^n x_i \left(y_i - \frac{1}{1+e^{w^T x_i}} \right)
 \end{aligned}$$

$x_{\text{test}} \in \mathbb{R}^d$

$\hat{y}_{\text{test}} = \text{Sign}(\hat{w}^T x_{\text{test}})$

Gradient update rule step-size

$$\begin{aligned}
 w_{t+1} &= w_t + \eta_t \nabla \log L(w_t) \\
 &= w_t + \eta_t \left(\sum_{i=1}^n x_i \left(y_i - \underbrace{\frac{1}{1+e^{w_t^T x_i}}}_{g(w_t^T x_i)} \right) \right)
 \end{aligned}$$

$\in \mathbb{R}^d$ θ_i
 $\{0,1\}$

KERNEL VERSION

- Can argue $w^* = \sum_{i=1}^n \alpha_i x_i$
 Formal theorem is called the Representer theorem.

REGULARIZED VERSION

$$\min_w \sum_{i=1}^n \left[\log(1+e^{-w^T x_i}) + w^T x_i (1-y_i) \right] + \underbrace{\frac{\lambda}{2} \|w\|^2}_{\text{Regularized}}$$

CROSS VALIDATE hyper parameters.