

Lecture 15: Kernel perceptron, Neural Networks, SVMs etc

Instructor: Prof. Ganesh Ramakrishnan

Perceptron Update Rule: Basic Idea

- Perceptron works for two classes ($y = \pm 1$). A point is misclassified if $y\mathbf{w}^T(\phi(\mathbf{x})) < 0$
- Perceptron Algorithm:

- ▶ INITIALIZE: $\mathbf{w} = \text{ones}()$ [for any initialization vector]
- ▶ REPEAT: for each $\langle \mathbf{x}, y \rangle$

- ★ If $y\mathbf{w}^T\phi(\mathbf{x}) < 0$
- ★ then, $\mathbf{w} = \mathbf{w} + \eta\phi(\mathbf{x}) \cdot y$
- ★ endif

↳ learning rate

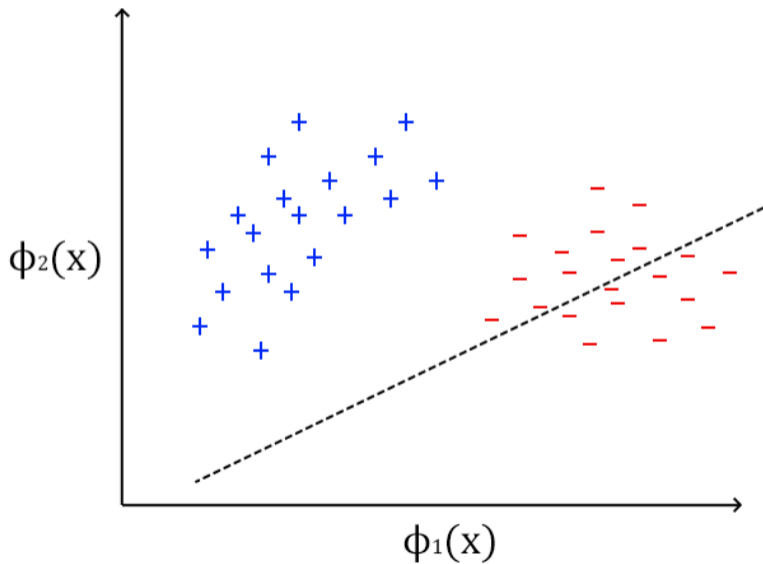
negative
unsigned
distance

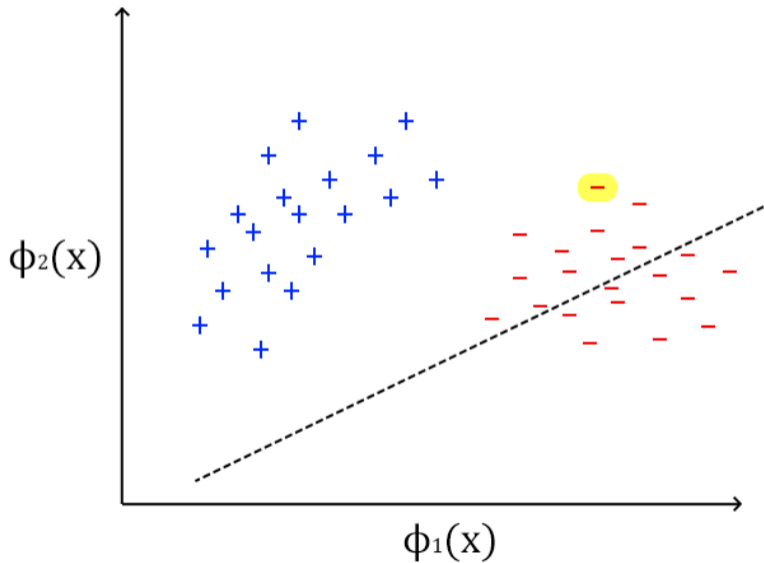
- Intuition:

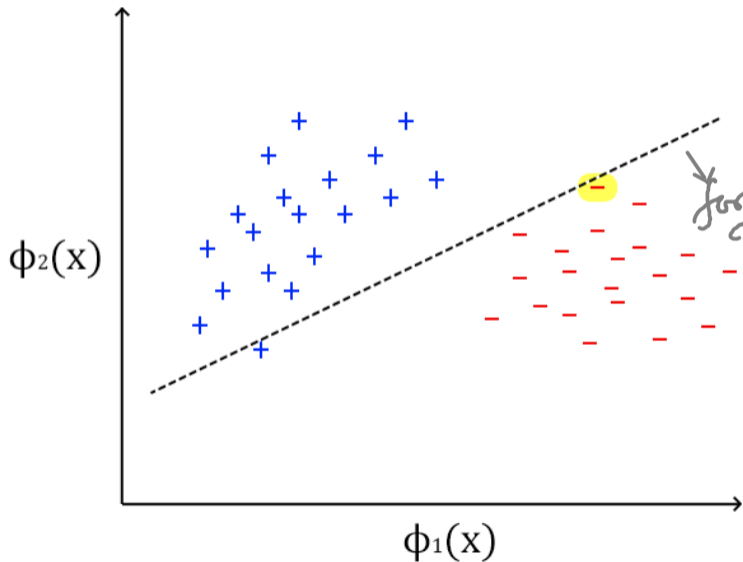
$$\begin{aligned}y(\mathbf{w}^{(k+1)})^T\phi(\mathbf{x}) &= y(\mathbf{w}^k + \eta y\phi^T(\mathbf{x}))\phi(\mathbf{x}) \\ &= y(\mathbf{w}^k)^T\phi(\mathbf{x}) + \eta y^2\|\phi(\mathbf{w})\|^2 \\ &> y(\mathbf{w}^k)^T\phi(\mathbf{x})\end{aligned}$$

Unsigned distance (initially negative becomes more positive)

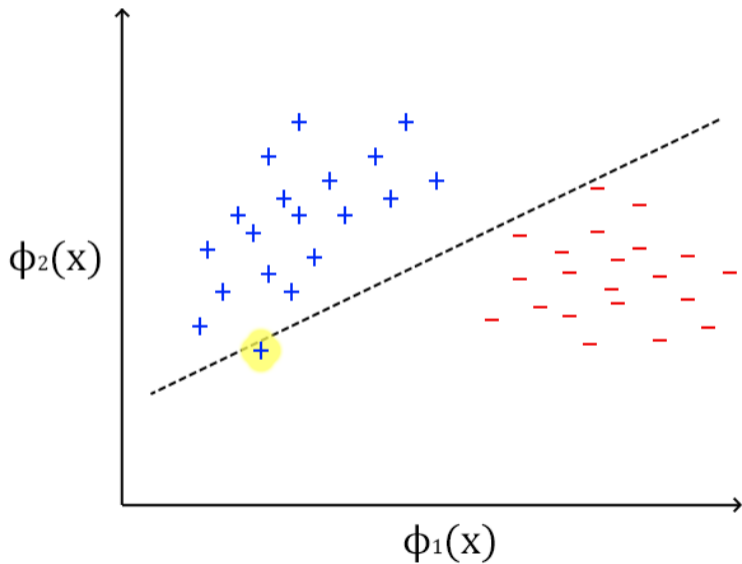
Since $y(\mathbf{w}^k)^T\phi(\mathbf{x}) \leq 0$, we have $y(\mathbf{w}^{(k+1)})^T\phi(\mathbf{x}) > y(\mathbf{w}^k)^T\phi(\mathbf{x}) \Rightarrow$ more hope that this point is classified correctly now.

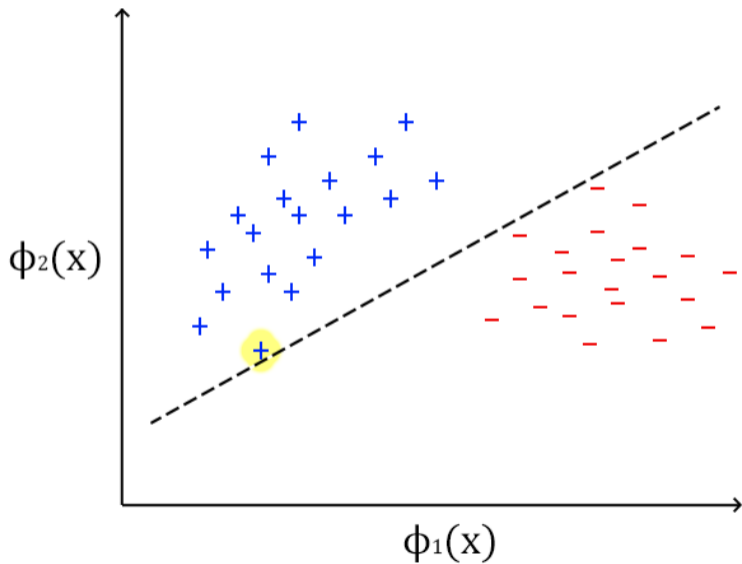






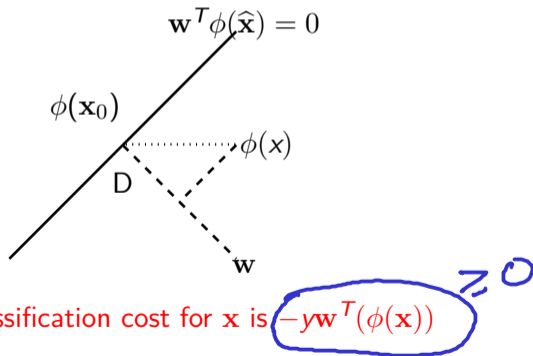
for some other
choice of η
this point
could have
yet remained
(less)
misclassified





Perceptron Update Rule: Error Perspective

- Explicitly account for signed distance of (misclassified) points from the hyperplane $\mathbf{w}^T \phi(\hat{\mathbf{x}}) = 0$. Consider point \mathbf{x}_0 such that $\mathbf{w}^T(\phi(\mathbf{x}_0)) = 0$
- (Signed) Distance from hyperplane is: $\frac{\mathbf{w}^T(\phi(\mathbf{x}) - \phi(\mathbf{x}_0))}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\phi(\mathbf{x}))}{\|\mathbf{w}\|}$
- Unsigned distance from hyperplane is: $y\mathbf{w}^T(\phi(\mathbf{x}))$ (assumes correct classification)



$$\text{sign}(y) = -\text{sign}(\mathbf{w}^T \phi(\mathbf{x}))$$

- If \mathbf{x} is misclassified, the misclassification cost for \mathbf{x} is $-y\mathbf{w}^T(\phi(\mathbf{x}))$

Perceptron Update Rule: Error Minimization

- Perceptron update tries to minimize the error function $E =$ negative of sum of unsigned distances over misclassified examples = **sum of misclassification costs**

$$E = - \sum_{(x,y) \in \mathcal{M}} y \mathbf{w}^T \phi(\mathbf{x})$$

where $\mathcal{M} \subseteq \mathcal{D}$ is the set of misclassified examples.

- Gradient Descent (Batch Perceptron) Algorithm** $\nabla_{\mathbf{w}} E = - \sum_{(x,y) \in \mathcal{M}} y \phi(\mathbf{x})$

to minimize E

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \mathbf{w}^k - \eta \nabla_{\mathbf{w}} E \\ &= \mathbf{w}^k + \eta \sum_{(x,y) \in \mathcal{M}} y \phi(\mathbf{x}) \end{aligned}$$

Perceptron Update Rule: Error Minimization

- Batch update considers all misclassified points simultaneously

$$\begin{aligned}\mathbf{w}^{(k+1)} &= \mathbf{w}^k - \eta \nabla_{\mathbf{w}} E \\ &= \mathbf{w}^k + \eta \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x})\end{aligned} \quad \left. \vphantom{\begin{aligned}\mathbf{w}^{(k+1)} \\ &= \mathbf{w}^k + \eta \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x})\end{aligned}} \right\} E = \sum_{(\mathbf{x}, y) \in \mathcal{M}} E(\mathbf{x})$$

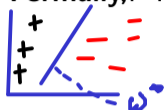
- Perceptron update \Rightarrow Stochastic Gradient Descent:

$$\nabla_{\mathbf{w}} E = - \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x}) = - \sum_{(\mathbf{x}, y) \in \mathcal{M}} \underbrace{\nabla_{\mathbf{w}} E(\mathbf{x})}_{\text{wavy line}} \text{ s.t. } E(\mathbf{x}) = -y \mathbf{w}^T \phi(\mathbf{x})$$

$$\begin{aligned}\mathbf{w}^{(k+1)} &= \mathbf{w}^k - \eta \nabla_{\mathbf{w}} E(\mathbf{x}) && \text{(for any } (\mathbf{x}, y) \in \mathcal{M} \text{)} \\ &= \mathbf{w}^k + \eta y \phi(\mathbf{x})\end{aligned}$$

Perceptron Update Rule: Further analysis

- Formally, - If \exists an optimal separating hyperplane with parameters \mathbf{w}^* such that,



$$\forall (\mathbf{x}, y), y \phi^T(\mathbf{x}) \mathbf{w}^* \geq 0$$

then the perceptron algorithm converges.

Proof: - We want to show that

wrt \mathbf{w}^ , unsigned distances are all absolute values*

$$\lim_{k \rightarrow \infty} \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 = 0 \quad (1)$$

(If this happens for some constant ρ , we are fine.)

$$\|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|_2^2 = \|\mathbf{w}^k + y \phi(\mathbf{x}^k) - \rho \mathbf{w}^*\|_2^2$$

$$= \|\mathbf{w}^k - \rho \mathbf{w}^*\|_2^2 + 2y \phi^T(\mathbf{x}^k) (\mathbf{w}^k - \rho \mathbf{w}^*) + y^2 \|\phi(\mathbf{x}^k)\|_2^2$$

$$\leq -\theta^2 \Rightarrow \|\mathbf{w}^{(0)} - \rho \mathbf{w}^*\|_2^2 / \theta^2 = \# \text{ of iterations for convergence}$$

θ should be independent of iteration # k

We are happy with the existence of such a $\rho > 0$

Perceptron Update Rule: Further analysis

- **Formally**,:- If \exists an optimal separating hyperplane with parameters \mathbf{w}^* such that,

$$\forall (\mathbf{x}, y), y\phi^T(\mathbf{x})\mathbf{w}^* \geq 0$$

then the perceptron algorithm converges.

Proof:- We want to show that

$$\lim_{k \rightarrow \infty} \|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 = 0 \quad (1)$$

(If this happens for some constant ρ , we are fine.)

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 = \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 + \underbrace{\|y\phi(\mathbf{x})\|^2 + 2y(\mathbf{w}^k - \rho\mathbf{w}^*)^T\phi(\mathbf{x})}_{\theta^2} \quad (2)$$

- For convergence of perceptron, we need L.H.S. to be less than R.H.S. at every step, although by some small but non-zero value (with $\theta \neq 0$)

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 - \theta^2 \quad (3)$$

Perceptron Update Rule: Further analysis

- Need that $\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2$ reduces by atleast θ^2 at every iteration.

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 - \theta^2 \quad (4)$$

- Based on (2) and (4), we need to find θ such that,

Perceptron Update Rule: Further analysis

- Need that $\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2$ reduces by atleast θ^2 at every iteration.

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 - \theta^2 \quad (4)$$

- Based on (2) and (4), we need to find θ such that,

$$\|\phi(\mathbf{x})\|^2 + 2y(\mathbf{w}^k - \rho\mathbf{w}^*)^T \phi(\mathbf{x}) \leq -\theta^2$$

($\|y\phi(\mathbf{x})\|^2 = \|\phi(\mathbf{x})\|^2$ since $y = \pm 1$)

- The number of iterations would be: $O\left(\frac{\|\mathbf{w}^{(0)} - \rho\mathbf{w}^*\|^2}{\theta^2}\right)$
- Tutorial 6, Problem 4 is concerning the number of iterations. But first we will discuss how convergence holds in the first place!

Perceptron Update Rule: Further analysis

The expression of interest : $\|\phi(x)\|^2 + \underbrace{2y(\omega^k - \rho\omega^*)^T \phi(x)}_{2y(\omega^k)^T \phi(x) < 0}$

Observations:-

① $y(\omega^k)^T \phi(x) < 0$ (\because x was misclassified)

② $\Gamma^2 = \max_{x \in \mathcal{D}} \|\phi(x)\|^2 \geq \|\phi(x)\|^2$

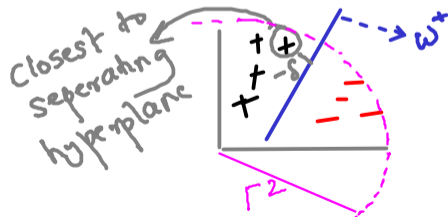
③ $\underline{\delta} = \max_{x \in \mathcal{D}} -2y\omega^{*T} \phi(x) = -\min_{x \in \mathcal{D}} 2y\omega^{*T} \phi(x) \Rightarrow -2y(\omega^*)^T \phi(x) \leq \delta$

Here, negative margin $\delta = -2y\omega^{*T} \phi(\hat{x})$ is the negative of unsigned distance of closest point \hat{x} from separating hyperplane : $\hat{x} = \operatorname{argmax}_{x \in \mathcal{D}} -2y\omega^{*T} \phi(x) = \operatorname{argmin}_{x \in \mathcal{D}} y\omega^{*T} \phi(x)$

Since the data is linearly separable,

$$\delta < 0$$

$\Gamma^2 =$ radius of tightest ball enclosing all points in \mathcal{D}



Perceptron Update Rule: Further analysis

- **Observations:-**

- 1 $y(\mathbf{w}^k)^T \phi(\mathbf{x}) < 0$ (\because \mathbf{x} was misclassified)
- 2 $\Gamma^2 = \max_{\mathbf{x} \in \mathcal{D}} \|\phi(\mathbf{x})\|^2$
- 3 $\delta = \max_{\mathbf{x} \in \mathcal{D}} -2y\mathbf{w}^{*T} \phi(\mathbf{x})$

Need: $\beta > 0$ st
 $\Gamma^2 + \beta \delta < 0$ where $\delta < 0$
 Ans: $\beta = -\frac{2\Gamma^2}{\delta} > 0$

- Here, negative margin $\delta = -2y\mathbf{w}^{*T} \phi(\hat{\mathbf{x}})$ is the negative of unsigned distance of closest point $\hat{\mathbf{x}}$ from separating hyperplane : $\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{D}} -2y\mathbf{w}^{*T} \phi(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} y\mathbf{w}^{*T} \phi(\mathbf{x})$
- Since the data is linearly separable, $\hat{y}\mathbf{w}^{*T} \phi(\hat{\mathbf{x}}) \geq 0$, so, $\delta \leq 0$. Consequently:

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 + \Gamma^2 + \rho\delta \rightarrow \Theta^2$$

$\geq \max_{\mathbf{x} \in \mathcal{D}} -y\mathbf{w}^{*T} \phi(\mathbf{x})$

$\|\phi(\mathbf{x})\|^2$

Next: we show \exists a $\beta \geq 0$ that gives us $\Theta^2 \geq 0$

Perceptron Update Rule: Further analysis

- Since, $\mathbf{w}^{*T} \phi(\hat{\mathbf{x}}) \geq 0$, so, $\delta \leq 0$. Consequently:

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 + \Gamma^2 + \rho \delta$$

Taking,

Perceptron Update Rule: Further analysis

- Since, $\mathbf{w}^{*T} \phi(\hat{\mathbf{x}}) \geq 0$, so, $\delta \leq 0$. Consequently:

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 + \Gamma^2 + \rho \delta$$

Taking, $\rho = \frac{2\Gamma^2}{-\delta}$, (Any $\rho \geq \frac{2\Gamma^2}{-\delta}$ will work)

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 - \Gamma^2$$

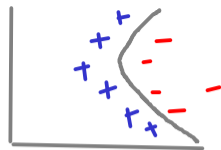
- Hence, we got, $\Gamma^2 = \theta^2$, that we were looking for in eq.(3).
 $\therefore \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2$ decreases by atleast Γ^2 at every iteration.
- Summarily: \mathbf{w}^k converges to $\rho \mathbf{w}^*$ by making a minimum θ^2 decrement at each step.
- Thus, for $k \rightarrow \infty$, $\|\mathbf{w}^k - \rho \mathbf{w}^*\| \rightarrow 0$. This proves convergence.

Perceptron Update Rule: Further analysis

- A statement on number of iterations for convergence:

If $\|\mathbf{w}^*\| = 1$ and if there exists $\delta > 0$ such that for all $i = 1, \dots, n$, $y_i(\mathbf{w}^*)^T \phi(\mathbf{x}_i) \geq \delta$ and $\|\phi(\mathbf{x}_i)\|^2 \leq \Gamma^2$ then the perceptron algorithm will make at most $\frac{\Gamma^2}{\delta^2}$ errors (that is take at most $\frac{\Gamma^2}{\delta^2}$ iterations to converge)

Non-linear perceptron?



- Kernelized perceptron: Non linear separation in implicit ϕ space?

Hint: Recall for non-parametric regression

$$f(x) = \sum_{i=1}^m \alpha_i \underbrace{k(x, x_i)}_{\phi_i(x)} y_i \rightarrow y_i \in \mathbb{R}$$

Ans: $f(x) = \sum_{i=1}^m \alpha_i k(x, x_i) y_i$ $y_i \in \{+1, -1\}$ if $f(x) > 0$ $y = +1$
 $f(x) < 0$ $y = -1$

Non-linear perceptron?

- Kernelized perceptron: $f(\mathbf{x}) = \text{sign} \left(\sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) + b \right)$

- ▶ INITIALIZE: $\alpha = \text{zeroes}()$

- ▶ REPEAT: for $\langle \mathbf{x}_i, y_i \rangle$

- ★ If $\text{sign} \left(\sum_j \alpha_j y_j K(\mathbf{x}_j, \mathbf{x}_j) + b \right) \neq y_i$

- ★ then, $\alpha_i = \alpha_i + 1$

- ★ endif

→ if x_i is indeed misclassified its weight α_i by 1 (Tutorial 6) is abt correspondence with perceptron

Like before check if x_i is misclassified

- Neural Networks: Cascade of layers of perceptrons giving you non-linearity. But before that, we will discuss the specific sigmoidal perceptron used most often in Neural Networks

Problem with $f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}))$ is the nondifferentiability of $f(\mathbf{x})$

$\lambda(\mathbf{x}) = 1$ if x is misclassified & therefore of $E = -\sum_{(\mathbf{x}, y) \in D} y \phi^T(\mathbf{x}) \mathbf{w} \lambda_{\mathbf{w}}(\mathbf{x})$

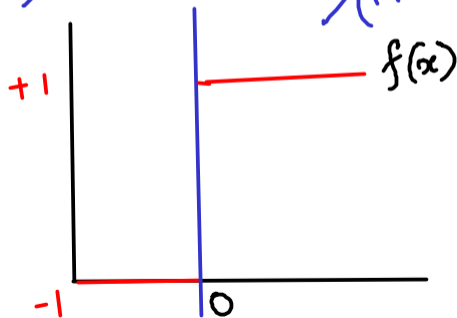
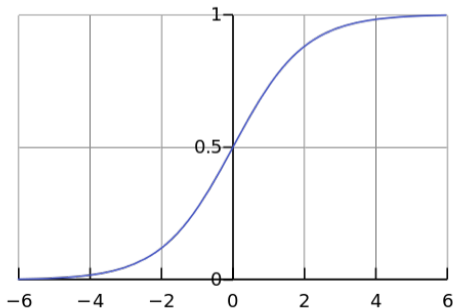
Sigmoidal (perceptron) Classifier [Note: $y \in \{0, 1\}$ against $y \in \{+1, -1\}$]

- ① **(Binary) Logistic Regression**, abbreviated as **LR** is a single node perceptron-like classifier, but with....

▶ $\text{sign}((\mathbf{w}^*)^T \phi(\mathbf{x}))$ replaced by $g((\mathbf{w}^*)^T \phi(\mathbf{x}))$ where $g(s)$ is sigmoid function: $g(s) = \frac{1}{1+e^{-s}}$

- ② $g((\mathbf{w}^*)^T \phi(\mathbf{x})) = \frac{1}{1+e^{-(\mathbf{w}^*)^T \phi(\mathbf{x})}} \in [0, 1]$ can be interpreted as $Pr(y = 1 | \mathbf{x})$

▶ Then $Pr(y = 0 | \mathbf{x}) = ?$ $1 - Pr(y = 1 | \mathbf{x}) = \frac{e^{-s}}{1+e^{-s}} = \frac{e^{-\mathbf{w}^T \phi(\mathbf{x})}}{1+e^{-\mathbf{w}^T \phi(\mathbf{x})}}$



Logistic Regression: The Sigmoidal (perceptron) Classifier

- ① Estimator $\hat{\mathbf{w}}$ is a function of the dataset

$$\mathcal{D} = \left\{ (\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)}) \right\}$$

- ▶ Estimator $\hat{\mathbf{w}}$ is meant to approximate the parameter \mathbf{w} .

- ② Maximum Likelihood Estimator: Estimator $\hat{\mathbf{w}}$ that maximizes the likelihood $L(\mathcal{D}; \mathbf{w})$ of the data \mathcal{D} .

- ▶ Assumes that all the instances $(\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)})$ in \mathcal{D} are all independent and identically distributed (iid)

- ▶ Thus, Likelihood is the probability of \mathcal{D} under iid assumption: $\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} L(\mathcal{D}, \mathbf{w}) =$

$$\prod_{i=1}^m P(y^{(i)} | x^{(i)}, \mathbf{w}) = \prod_{i=1}^m P(1 | x^{(i)}, \mathbf{w})^{y^{(i)}} P(0 | x^{(i)}, \mathbf{w})^{(1-y^{(i)})}$$

Exactly as in Bernoulli:
 $\text{Ber}\left(\frac{1}{1+e^{-\mathbf{w}^T \phi(x)}}\right)$

Logistic Regression: The Sigmoidal (perceptron) Classifier

- ① Estimator $\hat{\mathbf{w}}$ is a function of the dataset

$$\mathcal{D} = \left\{ (\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)}) \right\}$$

- ▶ Estimator $\hat{\mathbf{w}}$ is meant to approximate the parameter \mathbf{w} .

- ② Maximum Likelihood Estimator: Estimator $\hat{\mathbf{w}}$ that maximizes the likelihood $L(\mathcal{D}; \mathbf{w})$ of the data \mathcal{D} .

- ▶ Assumes that all the instances $(\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)})$ in \mathcal{D} are all independent and identically distributed (iid)
- ▶ Thus, Likelihood is the probability of \mathcal{D} under iid assumption: $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} L(\mathcal{D}, \mathbf{w}) =$

$$\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^{(i)} | \phi(\mathbf{x}^{(i)})) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m \left(\frac{1}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^{y^{(i)}} \left(\frac{e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^{1 - y^{(i)}}$$

Training LR

- 1 Thus, Maximum Likelihood Estimator for \mathbf{w} is

$$\begin{aligned}\hat{\mathbf{w}} &= \operatorname{argmax}_{\mathbf{w}} L(\mathcal{D}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^{(i)} | \phi(\mathbf{x}^{(i)})) \\ &= \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m \left(\frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x}^{(i)})}} \right)^{y^{(i)}} \left(\frac{e^{-\mathbf{w}^T \phi(\mathbf{x}^{(i)})}}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x}^{(i)})}} \right)^{1-y^{(i)}} \\ &= \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m \left(f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{y^{(i)}} \left(1 - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{1-y^{(i)}} \rightarrow \mathcal{R}(1 | \mathbf{x}^{(i)}, \mathbf{w})\end{aligned}$$

- 2 Maximizing the likelihood $\Pr(\mathcal{D}; \mathbf{w})$ w.r.t \mathbf{w} , is the same as minimizing the negative log-likelihood $E(\mathbf{w}) = -\frac{1}{m} \log \Pr(\mathcal{D}; \mathbf{w})$ w.r.t \mathbf{w} .

- ▶ Derive the expression for $E(\mathbf{w})$. \rightarrow HW
- ▶ $E(\mathbf{w})$ is called the cross-entropy loss function

Distance between 2 probability distributions