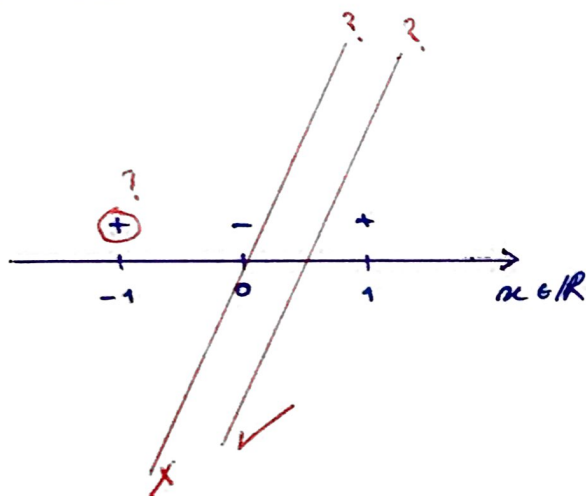Objectives:

- Derive non-linear classifiers from feature maps
- Move from coordinate parameterization to weighting examples
- Compute kernel functions induced from feature maps
- Use kernel perceptron, kernel linear regression
- Understand the properties of kernel functions
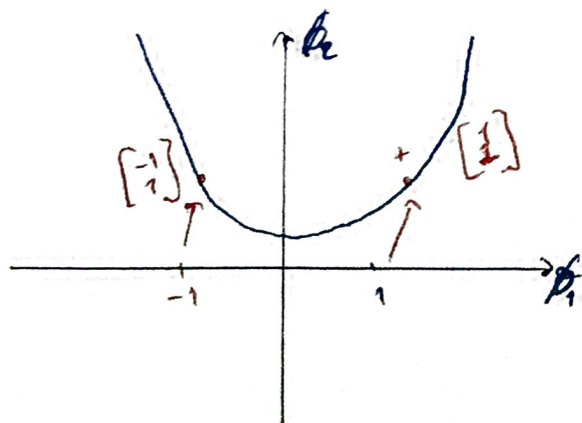
Higher Order Feature vectors:



In this case we can't use linear classifiers to map the data.

We can remedy this situation by introducing a feature transformation feeding a different type of example to the linear classifier.

$$x \longrightarrow \phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix} \longleftarrow \text{we add additional feature like } x^2$$
$$\in \mathbb{R} \qquad\qquad \in \mathbb{R}^2$$

$$\theta \longrightarrow \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$h(x, \theta, \theta_0) = \text{sign}\left(\theta \cdot \phi(x) + \theta_0\right)$$
$$= \text{sign}\left(\theta_1 x + \theta_2 x^2 + \theta_0\right)$$



$$\phi(x) = \begin{bmatrix} \phi_1 \\ \phi_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix} \text{ in this case}$$

$$1 \longrightarrow \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$-1 \longrightarrow \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

A linear classifier in the new feature coordinate implies a non linear classifier in the x

# Introduction to Non-linear Classification:

Non linear classification : $h(x; \theta, \theta_0) = \text{sign}(\theta \cdot \phi(x) + \theta_0)$

Non linear Regression : $f(x; \theta, \theta_0) = \theta \cdot \phi(x) + \theta_0$

By mapping input examples explicitly into feature vectors, and performing linear classification

on regression on top of such feature vectors, we get a lot of expressive power.

But the downside is that these vectors can be quite high dimensional

## Computational Efficiency:

Computing the inner product between two feature vectors can be cheap even if the vectors are

very high dimensional.

$$\phi(x) = [x_1, x_2, x_1^2, \sqrt{2}\, x_1 x_2, x_2^2]^T$$

Different examples
$x$ and $x'$
$$\phi(x') = [x_1', x_2', x_1'^2, \sqrt{2}\, x_1' x_2', x_2'^2]^T$$

$$k(x, x') = \phi(x) \cdot \phi(x') = (x \cdot x') + (x \cdot x')^2$$

Kernel
function

inner product
or dot product

The Kernel can be evaluated very cheaply, even though we would have to explicitly

construct very high dimensional feature vectors.

Task :    linear methods $\longrightarrow$ methods that can operate in terms of kernels

$$\text{Sign}(\theta \cdot \phi(x) + \theta_0) \longrightarrow k(x, x')$$

## Kernel as Dot Products :

$$\phi(x) = [x_1, x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2]$$
$$\phi(x') = [x_1', x_2', x_1'^2, \sqrt{2}x_1' x_2', x_2'^2]$$

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$x' = \begin{bmatrix} x_1' \\ x_2' \end{bmatrix}$

$$\phi(x) \cdot \phi(x') = x_1 x_1' + x_2 x_2' + x_1^2 x_2'^2 + 2 x_1 x_2 x_1' x_2' + x_2^2 + x_2'^2$$

$$= (x_1 x_1' + x_2 x_2') + (x_1 x_1' + x_2 x_2')^2$$

$$= x \cdot x' + (x \cdot x')^2$$

$$k(x, x') = \phi(x) \cdot \phi(x') = x_1 x_1' + x_2 x_2' + x_3 x_3' + x_2 x_2' + x_3 x_2' \longleftarrow \phi(x) = [x_1, x_2, x_3]$$

## Perceptron:

$$(\{(x^{(i)}, y^{(i)}), i = 1, \ldots m\}, T):$$

initialize $\theta = 0$ (vector)

for $t = 1, \ldots, T$

    for $i = 1, \ldots, m$

        if $y^{(i)}(\theta \cdot x^{(i)}) \leq 0$

            then update $\theta = \theta + y^{(i)} x^{(i)}$

**Radial Basis Kernel:**

$$k(x, x') = \exp\left(-\frac{1}{2} \|x - x'\|^2\right)$$

We introduced that we can always express $\theta$ as: $\theta = \sum_{j=1}^{m} \alpha_j \, y^{(j)} \cdot \phi(x^{(j)})$

## Kernel Perceptron:

$$(\{(x^{(i)}, y^{(i)}), i = 1, \ldots, m\}, T\})$$

initialize $\alpha_1, \alpha_2, \ldots, \alpha_m$ to some values $\longrightarrow \theta = 0 \iff$ setting $\alpha_j = 0$ for all $j$

for $t = 1, \ldots T$

    for $i = 1, \ldots m$

        if ( Mistake Condition Expressed in $\alpha_j$ ) $\longrightarrow y^{(i)}(\theta \cdot \phi(x^{(i)})) \leq 0 \iff y^{(i)} \sum_{j=1}^{m} \alpha_j \, y^{(j)} \, k(x^j, x^i) \leq$

            Update $\alpha_j$ appropriately $\longrightarrow \alpha_i = \alpha_i + 1 \iff \theta = \theta + y^{(i)} \phi(x^{(i)})$

## Kernel Composition Rules:

$k(x, x') = 1$ is a kernel function $\phi(x) = 1$

$\hat{k}(x, x') = f(x) \, k(x, x') \, f(x')$ is also a kernel $\hat{\phi}(x) = f(x) \phi(x)$

$k(x, x') = k_1(x, x') + k_2(x, x')$ is a kernel

$k(x, x') = k_1(x, x') \, k_2(x, x')$ is a kernel

## Kernel Composition Rule 1:

$\bar{k}(x, x') = f(x) \, k(x, x') \, f(x')$

if there exists $\phi(x)$ such that: $k(x, x') = \phi(x) \cdot \phi(x')$

    which of the following $\phi$ gives $\hat{k}(x, x') = \hat{\phi}(x) \cdot \hat{\phi}(x')$?

    $(f(x)\phi(x)) \cdot (f(x')\phi(x')) = \hat{k}(x, x') \implies \hat{\phi}(x) = f(x)\phi(x)$