

Basis function expansion & Kernel: Part 1

We saw that for $p \in [0, \infty)$, under certain conditions on K , the following can be equivalent representations

-

$$f(\mathbf{x}) = \sum_{j=1}^p w_j \phi_j(\mathbf{x})$$

basis function(j)

- And

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

basis function $\mathbf{x}_i = \mathbf{h}_i(\mathbf{x})$

- For what kind of regularizers, loss functions and $p \in [0, \infty)$ will these dual representations hold?¹

:Defer this question and deal with it after discussing classification

¹Section 5.8.1 of Tibshi.

Basis function expansion & Kernel: Part 2

- We could also begin with

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

K need not be psd/mercer

and impose no constraints on K .

K is some measure of similarity (x, x_i)

- E.g.: $K_k(x_q, x) = I(\|x_q - x\| \leq \|x_{(k)} - x_0\|)$ where $x_{(k)}$ is the training observation ranked k^{th} in distance from x and $I(S)$ is the indicator of the set S

deal with it for classification

- This is precisely the Nearest Neighbor Regression model
- Kernel regression and density models are other examples of such *local regression methods*²

You can basically use kernel regression to approximate functions, probability density.

²Section 2.8.2 of Tibshi

Kernel weighted regression (Midsem Q4: Local linear regression)

Weights obtained using some kernel $K(\cdot, \cdot)$. Given a training set of points $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_n, y_n)\}$, we predict a regression function $f(\mathbf{x}') = (\mathbf{w}'^\top \phi(\mathbf{x}') + b')$ for each test (or query point) \mathbf{x}' as follows:

$$(\mathbf{w}', b') = \operatorname{argmin}_{\mathbf{w}, b} \sum_{i=1}^n K(\mathbf{x}', x_i) \left(y_i - (\mathbf{w}'^\top \phi(x_i) + b) \right)^2$$

- 1 If there is a closed form expression for (\mathbf{w}', b') and therefore for $f(\mathbf{x}')$ in terms of the known quantities, derive it.
- 2 How does this model compare with linear regression and k -nearest neighbor regression? What are the relative advantages and disadvantages of this model?
- 3 In the one dimensional case (that is when $\phi(x) \in \mathbb{R}$), graphically try and interpret what this regression model would look like, say when $K(\cdot, \cdot)$ is the linear kernel³.

³Hint: What would the regression function look like at each training data

More on Kernels after some classification

(By data set here we mean the training dataset)

- 1 We will delve a bit more into kernel density estimation etc after some treatment of classification

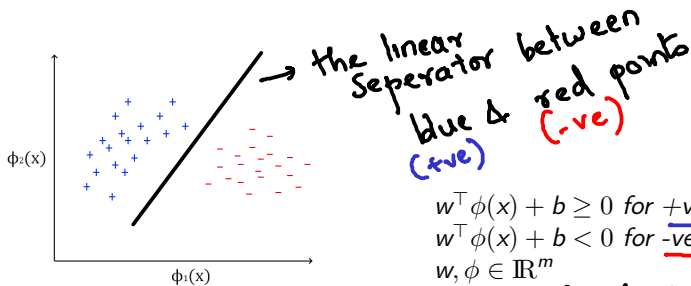
• We will talk of parametric (vs)

non-parametric estimation

of params don't change with data size
params increase with increasing dataset size

Perceptron

& Kernel perceptron
& deep neural networks



$$w^T \phi(x) + b \geq 0 \text{ for } +ve \text{ points } (y = +1)$$

$$w^T \phi(x) + b < 0 \text{ for } -ve \text{ points } (y = -1)$$

$$w, \phi \in \mathbb{R}^m$$

in other words: $y(w^T \phi(x) + b) \geq 0 \quad \forall (x, y)$

- Assuming the problem is linearly separable, there is a learning rule that converges in a finite time.
- A new (unseen) input pattern that is similar to an old (seen) input pattern is likely to be classified correctly

- Often, b is indirectly captured by including it in w , and using a ϕ as: $\phi_{aug} = [\phi, 1]$
- Thus, $w^\top \phi(x)$

$$= \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_m & b \end{bmatrix} \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \\ \vdots \\ \phi_m \\ 1 \end{bmatrix}$$

- $w^\top \phi(x) = 0$ is the separating hyperplane.

Perceptron Intuition

$$\{(\phi(x_1), y_1), (\phi(x_2), y_2) \dots (\phi(x_n), y_n)\}$$

- Go over all the existing examples, whose class is known, and check their classification with a current weight vector $(w^{(0)})$
- If correct, continue
- If not, add to the weights a quantity that is proportional to the product of the input pattern with the desired output y (1 or -1)

Exercise:

- ① Write down this perceptron (update) algorithm
- ② Geometrically interpret the updates
- ③ Will it converge for linearly separable dataset? Why?

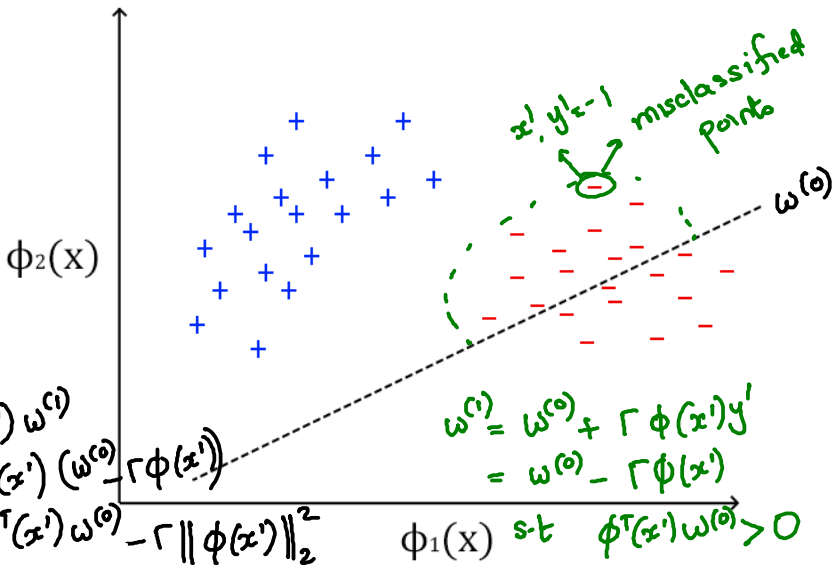
Perceptron Update Rule

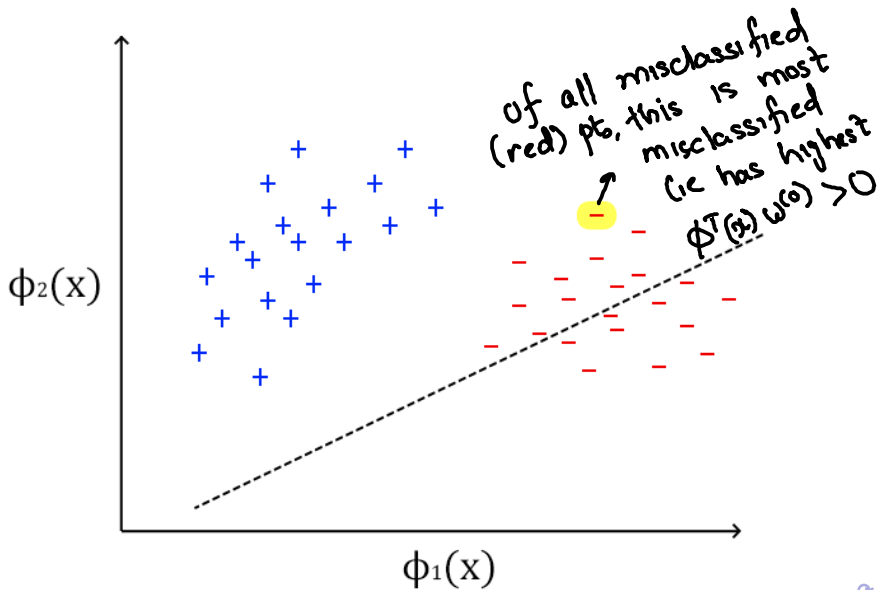
$$w^{(0)} \neq 0$$
$$\rightarrow w^{(0)} = [0 \ 0 \ \dots \ 1]$$

- Start with some weight vector $w^{(0)}$, and for $k = 1, 2, 3, \dots, n$ (for every example), do:
 $w^{(k)} = w^{(k-1)} + \Gamma \phi(\underline{x}')$
- where \underline{x}' s.t. \underline{x}' is misclassified by $(w^{(k)})^T \phi(x)$
i.e. $y(w^{(k)})^T \phi(\underline{x}') < 0$

A very handy & simple algo: *neural nets & graphical models!*

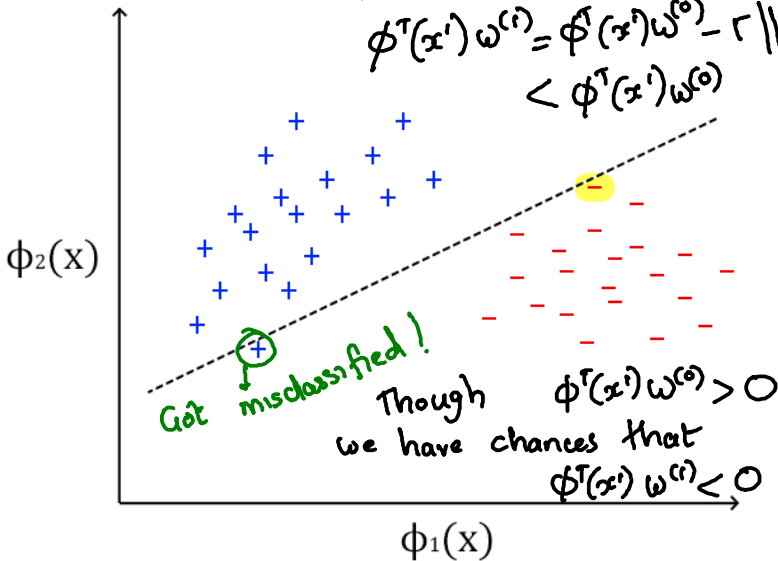
- ① At the core of neural network & several graphical model training
- ② It is "stochastic", dealing with one point at a time & can be used for incremental training

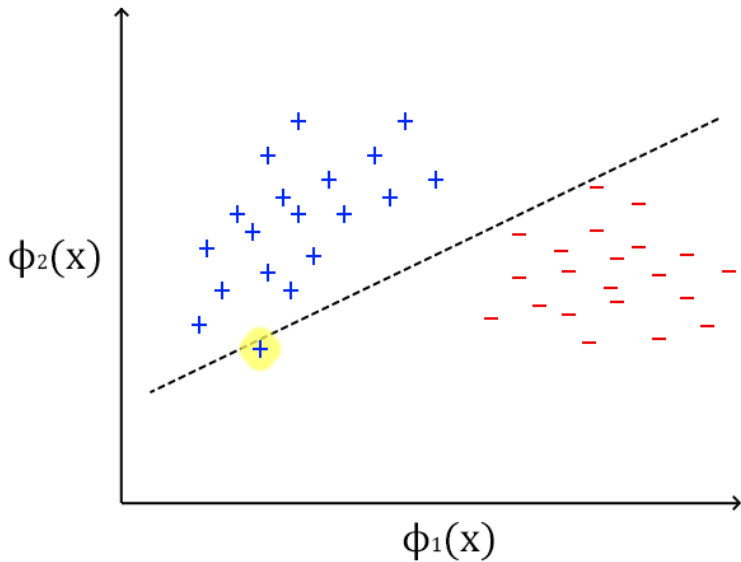


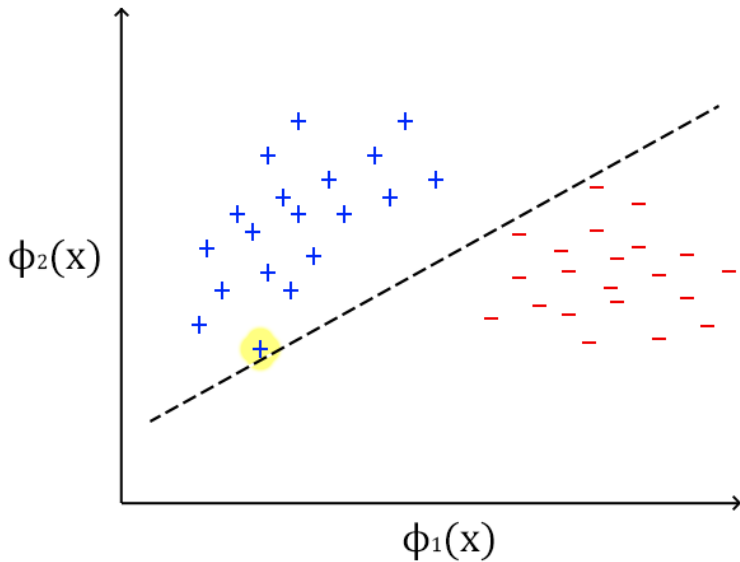


From what we saw:

$$\phi^T(x^i) \omega^{(i)} = \phi^T(x^i) \omega^{(0)} - \gamma \|\phi(x^i)\|_2^2 < \phi^T(x^i) \omega^{(0)}$$







- Perceptron does not find the *best* separating hyperplane, it finds *any* separating hyperplane.
- In case the initial w does not classify all the examples, the separating hyperplane corresponding to the final w^* will often pass through an example.
- The separating hyperplane does not provide enough breathing space – this is what SVMs address!

H/W: Go through proof of convergence of perceptron update!