Kernels, dual forms, etc.

1. Kernelized Perceptron:

• First, recall the exact sol° for Lin Reg:
$$w = \underbrace{(X^T X)^{-1}}_{C,\ cov.\ matrix\ of\ the\ data} X \cdot Y$$

$$= C^{-1} X \cdot Y$$

$$= \sum_n x_n \, y_n \cdot (C^{-1})$$

• This is quite different from the perceptron solution: Starting from $\vec{w} = \vec{0}$, $b = 0$, we add $+ \eta \, \vec{x}_n \, y_n$ to $\vec{w}$ each time an example $(\vec{x}_n, y_n)$ is wrongly classified by $\vec{w}$ (think eg of the online version).

by the model

So, assuming we do converge (as is guaranteed in the case of linearly separable data), the solut° $\vec{w}^*$ is built as a sum like:

$$\vec{w} = \left( \sum_n \underbrace{\alpha_n}_{\downarrow} \vec{x}_n \, y_n \right) \eta$$

$$\begin{cases} 0 \text{ if it was always well classif,} \\ 1 \text{ if it was incorrectly classified only exactly once} \\ 2 \\ \vdots \end{cases}$$

$\alpha_n$ is a counter $(\geq 0)$ of the number of times example $n$ was misclassified

We drop $\eta$ or include it in $\alpha_n$, so that we have $\quad \vec{w} = \sum_n \alpha_n \, \vec{x}_n \, y_n$

In this view, $\vec{w}$ is a linear combination of the training examples $\vec{x}_n$, with weights $\alpha_n y_n$ ($\alpha_n > 0$, $y_n = \pm 1$).

This may be called the dual form.

In the online perceptron, in a sense, we are updating the $\alpha_n$'s (they start from $\alpha_n = 0$, $\forall n$).

- This is very $\neq$ from $w = \sum C^{-1} \vec{x}_n y_n$ in linear reg, where in effect, $\alpha_n = C^{-1}$, $\forall n$, (constant $\alpha_n$).

- At $\underline{prediction\ time}$, we get:
$$y^{pred}(\vec{x}^{test}) = sign(\vec{w} \vec{x}_{test})$$
$$= sign\left(\sum_n \alpha_n y_n \underline{\vec{x}_n \cdot \vec{x}_{test}}\right)$$
If we used feature maps, $\phi: x_n \to \phi(x_n)$, we would have: $y_{test} = sign\left(\sum_n \alpha_n y_n \phi(\vec{x}_n) \phi(\underline{\vec{x}_{test}})\right)$

- $\underline{Remark}$: $\vec{x}_n \vec{x}_{test}$ is a measure of the similarity between $\vec{x}_n$ and $\vec{x}_{test}$. If they're very $\neq$, it's $\to 0$, and $\alpha_n$ does not matter for $y_{test}$. If they are very similar, it is large, and $\alpha_n$ matters.
(Note: if data is standardized, then it cannot grow too large).

- $\underline{Definition}$: We call $Kernel\ method$ the fact of replacing $\vec{x} \cdot \vec{x}'$ (or $\phi(\vec{x}) \cdot \phi(\vec{x}')$) with an other function $K(\vec{x}, \vec{x}')$, which is called a $Kernel$.

## Several remarks:

1) Kernels are more general than feature maps:
   - all feature maps are kernels:
$$K(x,x') = \phi(x)\phi(x') \text{ is a kernel, for}$$
   any feature map $\phi$.

   - not all kernels can be re-written as feature maps (see eg the RBF kernel $\rightarrow$ it's like a $D=\infty$ feature map )

   - Not all funct° of 2 variables are valid kernels

We must respect the Mercer condit°:

Mercer condition = $K$ must be a semi-positive definite operator

$$= \forall f \in L^2(\mathbb{R}), \quad \cancel{\forall g \in L^2(\mathbb{R})} \quad \int_{x,y} f(x) K(x,y) f(y) \geqslant 0$$

In a discrete setting, this would be like $\forall f \in \mathbb{R}^d$, $\cancel{\forall g \in \mathbb{R}^d}$

we must have: $\quad f^T K f = \sum_i f_i K_{ij} f_j \geqslant 0.$

   - Remark how $K$ builds a new geometry in the space of features: the similarity between two points, instead of being measured by $\vec{x}\cdot\vec{x}'$ of $1/\|x - x'\|_2^2$, is measured by $K(x,x')$.