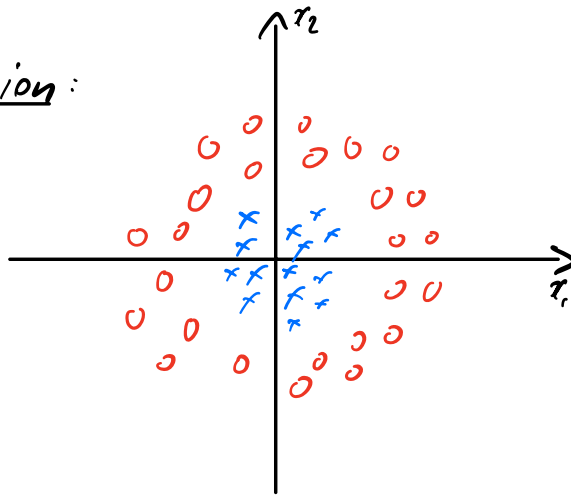


Kernels

Explicit Feature Expansion:



Quiz: Add new feature(s) such that this data set becomes linearly separable.

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1 \\ x_2 \\ ? \end{pmatrix}$$

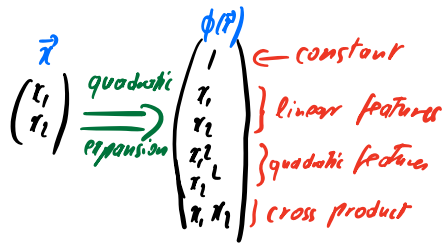
\vec{x} $\phi(\vec{x})$

Adding features can reduce Bias (but increase Variance).
Allows us to capture non-linear interactions between original features.

Problem in Practice: It is not obvious what new features to construct.

Case I:

Polynomial Features:

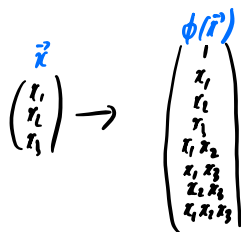


$$f(x) = \theta^T \vec{x} = v^T x + x^T M x + b \leftarrow \text{quadratic function}$$

Annotations: $v^T x$ is linear, $x^T M x$ is quadratic & cross product, and b is constant.

Case II:

All possible interactions:



Quiz: If $x \in \mathbb{R}^d$, what is the dimensionality of $\phi(\vec{x})$?



Explicit feature expansion seems very powerful but I am scared of such high dimensional data points!

Solution: Implicit feature expansion (Aka kernels)

Assume, all you want to do is compute inner products.

$$\vec{x}^T \vec{z} = \sum_{\alpha=1}^d x_{\alpha} z_{\alpha} \quad \text{how about } \phi(\vec{x})^T \phi(\vec{z})?$$

Case I:

$$\phi(\vec{x})^T \phi(\vec{z}) = (1 + \vec{x}^T \vec{z})^2$$

Case II:

$$\phi(\vec{x})^T \phi(\vec{z}) = \prod_{\alpha=1}^d (1 + x_{\alpha} z_{\alpha})$$

\Rightarrow Never compute $\vec{x} \rightarrow \phi(\vec{x})$.

Only compute $\vec{x}^T \vec{z} \rightarrow \phi(\vec{x})^T \phi(\vec{z}) = k(\vec{x}, \vec{z})$
kernel function

This is just the inner-product function in the high dimensional feature space.



ERM^{*} can be written entirely in terms of inner products!

*At least the losses that we discussed.

Classifier: $h(\vec{x}) = \omega^T \phi(\vec{x}) + b$ Loss: $L = \sum_{i=1}^n \ell(\omega^T \phi(\vec{x}_i), y_i)$

Claim: $\vec{\omega} = \sum_{i=1}^n \beta_i \phi(\vec{x}_i) \Leftrightarrow \vec{\omega}_i$ lies in the span of $\phi(\vec{x}_1), \dots, \phi(\vec{x}_n)$

\Rightarrow at test time: $h(\vec{z}) = \omega^T \phi(\vec{z}) + b = \sum_{i=1}^n \beta_i \underbrace{\phi(\vec{x}_i)^T \phi(\vec{z})}_{\text{only inner products}} + b$

test point

Proof by induction: We learn with gradient descent. Induction over update t .

Base case $t=0$: We initialize $\vec{\omega}_0 = \vec{0} = \sum_{i=1}^n 0 \cdot \phi(\vec{x}_i)$ 😊

Assume $\vec{\omega}^t = \sum_{i=1}^n \beta_i^t \phi(\vec{x}_i)$

Gradient update:

$$\vec{\omega}_{t+1} \leftarrow \vec{\omega}_t - \alpha \frac{\partial L}{\partial \vec{\omega}_i}$$

$$= \sum_{i=1}^n \beta_i^t \phi(\vec{x}_i) - \alpha \sum_{i=1}^n \gamma_i^t \phi(\vec{x}_i) = \sum_{i=1}^n \underbrace{(\beta_i^t - \alpha \gamma_i^t)}_{\beta_i^{t+1}} \phi(\vec{x}_i)$$

$$\Rightarrow \beta_i^{t+1} \leftarrow \beta_i^t - \alpha \gamma_i^t$$

chain rule
 $\frac{\partial L}{\partial \vec{\omega}_i} = \sum_{i=1}^n \frac{\partial \ell(\vec{\omega}^T \phi(\vec{x}_i), y_i)}{\partial \vec{\omega}_i} = \sum_{i=1}^n \underbrace{\ell'(\vec{\omega}^T \phi(\vec{x}_i), y_i)}_{\text{scalar}} \cdot \phi(\vec{x}_i)$

😊

□

$\Rightarrow h(\vec{z}) = \sum_{i=1}^n \beta_i k(\vec{x}_i, \vec{z}) + b$ ← we now learn $\vec{\beta} \in \mathbb{R}^n$ instead of $\vec{\omega} \in \mathbb{R}^d$ ← learning of b is not affected.

Popular kernels:

Linear: $k(\vec{x}, \vec{z}) = \vec{x}^T \vec{z}$

Polynomial: $k(\vec{x}, \vec{z}) = (1 + \vec{x}^T \vec{z})^p$

Radial Basis Function (RBF): $k(\vec{x}, \vec{z}) = \exp(-\lambda(\vec{x} - \vec{z})^2)$

Sigmoid: $k(\vec{x}, \vec{z}) = \tanh(a\vec{x}^T + c)$



Can any function $k(\vec{x}, \vec{z})$ be a kernel?

NO!
It has to be positive semi-definite (PSD)



$K \in \mathbb{R}^{n \times n}$

For any finite data set $D = \{\vec{x}_1, \dots, \vec{x}_n\}$ define $K_{ij} = k(\vec{x}_i, \vec{x}_j)$.

K must be p.s.d.

A symmetric matrix $K \in \mathbb{R}^{n \times n}$ is p.s.d. if and only if

① $\exists A \in \mathbb{R}^{n \times n}$ s.t. $K = A^T A$

② All eigenvalues of K are non-negative

③ $\forall q \in \mathbb{R}^n \quad q^T K q \geq 0$