# Support Vector Machines (SVM)

Linear
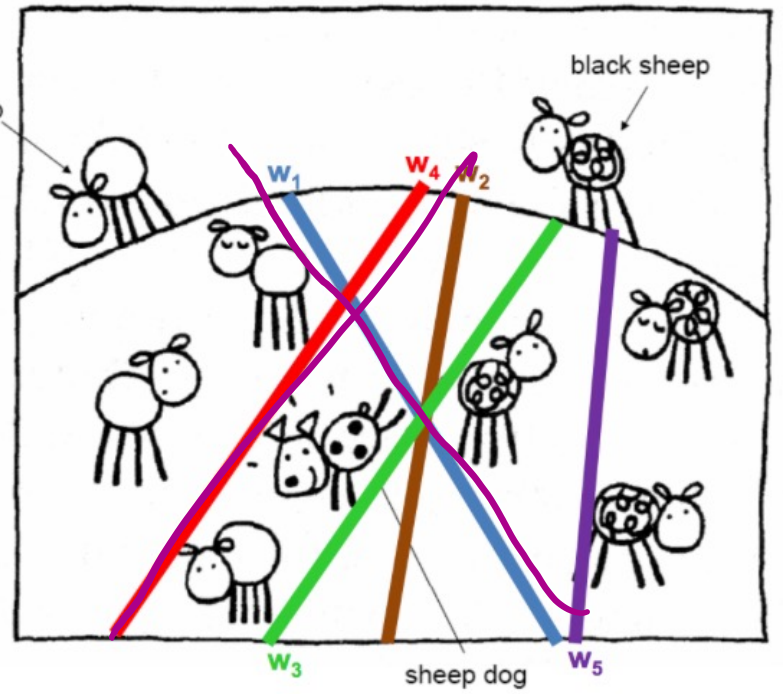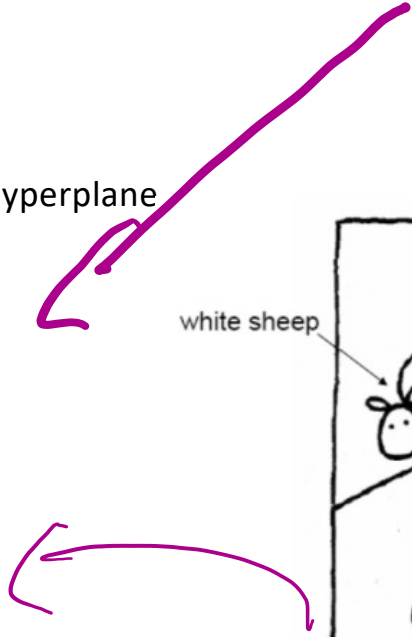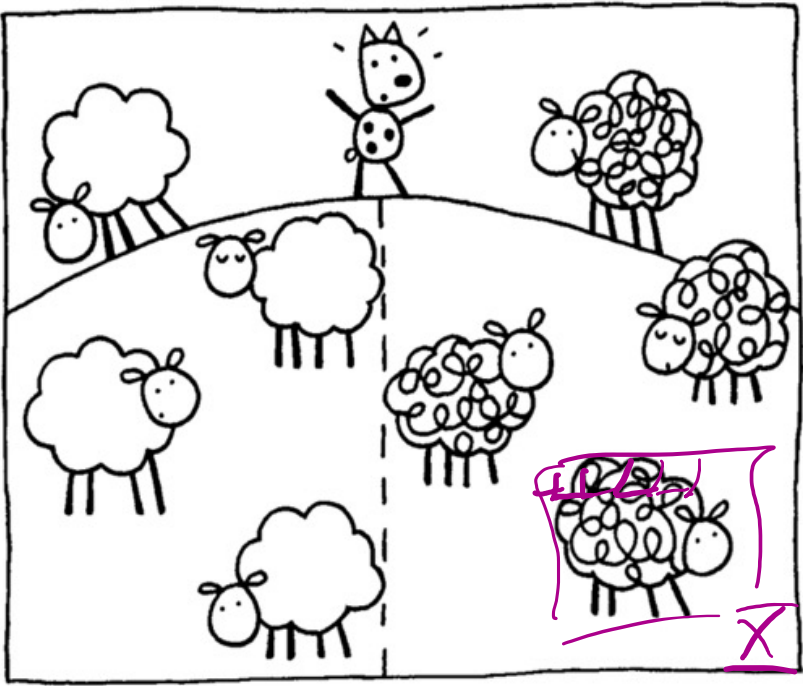
Then
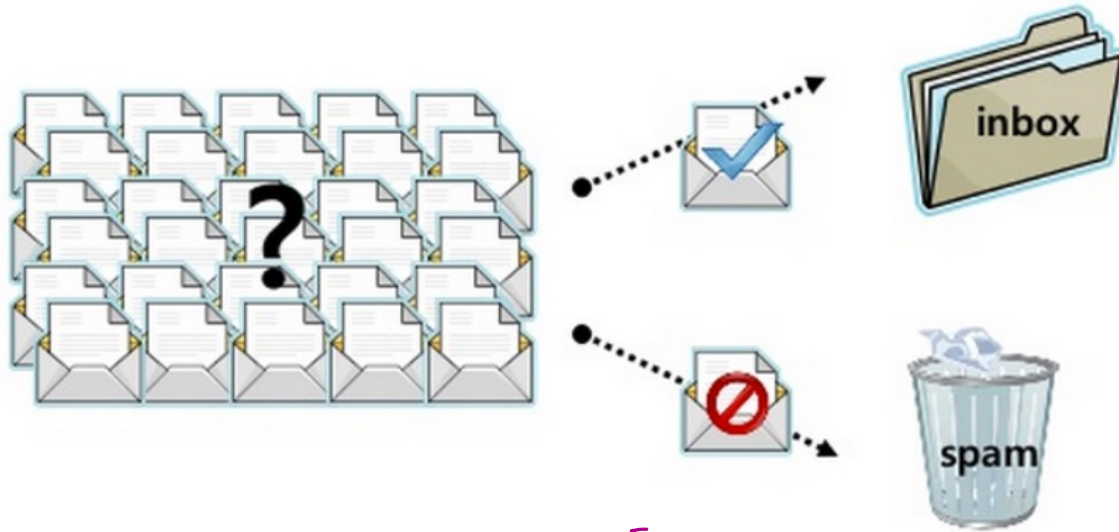Nonlinear (kernelized) SVM (KSVM)

Smola - Schökopff.

Wide Margin Decision Hyperplane for
Supervised - Learning Classification

First a linear binary classification – decision boundary/hyperplane

- Instance space: $x \in X$ ($|X| = n$ data points)
  - Binary or real-valued feature vector $x$ of word occurrences
  - $d$ features (words + other things, d~100,000+)
- Class: $y \in Y$
  - $y$: Spam (+1), Ham (-1)

$\{(x_i, y_i)\}_{i=1}^{n}$

$y_i = \pm \frac{1}{8}$

$\xi = 2$

| Viagra | Learning | The | Dating | Nigeria | Is_spam |
|--------|----------|-----|--------|---------|---------|
| 1 | 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 0 | 0 | -1 |
| 0 | 0 | 0 | 0 | 1 | 1 |

$$D = \{(x_i, y_i)\}$$

Review — a hyperplane defined by a vector.

$$\vec{n} = (a, b, c)$$

$$\vec{x} = (x_1, x_2, x_3)$$

$$\vec{x}_0 = (x_{1,0}, x_{2,0}, x_{3,0})$$

• l eqn is a co-dim -1 restriction of space

$$\vec{v} = \vec{x} - \vec{x}_0 = \langle (x_1 - x_{1,0}), (x_2 - x_{2,0}), (x_3 - x_{3,0}) \rangle$$

$$\Pi := \{ x = (x_1, x_2, x_3) : \vec{v} = \vec{x} - \vec{x}_0 \perp n \}$$

$$v \cdot n = 0 \iff v \perp \vec{n}$$

$$n \cdot v = \langle a, b, c \rangle \cdot \langle x_1 - x_{1,0}, x_2 - x_{2,0}, x_3 - x_{3,0} \rangle = a(x_1 - x_{1,0}) + b(x_2 - x_{2,0}) + c(x_3 - x_{3,0}) = 0$$

$SVM:$  $D = \{(x_i, y_i)\}_{i=1}^{n}$  $x_i \in \mathbb{R}^d$

$y_i \in \mathbb{Z}_2 = \{-1, 1\}$

$x_{i,2}$

Feature space

wide margin as big possible

$x_{i,1}$

hyper-plane

$$(\vec{x}_i, y_i) \quad , \quad S_i = -1, \boxed{1}$$

$$\left( \begin{array}{l} y_i = 0 \text{ or } 1 \\ \text{apple or orange} \\ \text{dog or cat.} \end{array} \right.$$

$$y_j(\omega \cdot x_j + b) = \text{sgn}(\omega \cdot x_j + b) \text{ good label.}$$

$$\text{sgn}(S) = \begin{cases} 1 & \text{if } S > 0 \\ 0 & \text{if } S = 0 \\ -1 & \text{if } S < 0 \end{cases}$$



$$\text{sgn}(S)\,\omega \cdot x_j + b) = 1 \text{ -- labelled well} \\ = -1 \text{ mislabelled.}$$

KKT

## Primal Problem:

$$\begin{cases} \text{minimize: } \mathcal{L}(x, s) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^n s_i y_i (w^T x_i + b) + \sum_{i=1}^n s_i \\ \text{such that: } s_i \geq 0, \forall i \end{cases}$$

## Dual Problem:

$$\begin{cases} \text{maximize: } \mathcal{L}_D(x, s) = \sum_{i=1}^n s_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n s_i s_j y_i y_j (\vec{x}_i^T \vec{x}_j) \\ \text{using: } w = \sum_{i=1}^n s_i y_i x_i, \text{ and } \sum_{i=1}^n s_i y_i = 0 \end{cases}$$

a loss function.

$$\ell(y_i, \bar{y}_i) = \ell(y_i, \text{sgn}(\omega \cdot \vec{x}_i + b)) = \begin{cases} 0 & \text{if correct label} \\ & \bar{y}_i = \text{sgn}(\omega \cdot x_i + b) \\ 1 & \text{incorrect label} \end{cases}$$

$\bar{y}_i$

label you infer from $\vec{x}_i$ alone
if you have a good hyperplane $\vec{\omega}$ & $b$

Total loss

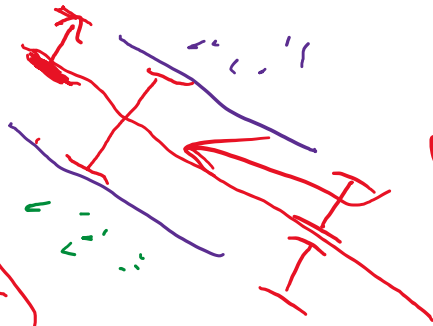$$\sum_{i=1}^{N} \ell(y_i, \bar{y}_i)$$

Cost :   ...   $\frac{MSE}{3}$ ↑

$\arg\min \left( \frac{1}{2} ||w||_2^2 \right)$   $\text{subj}$   $y_i \overbrace{(w^T x_i - b)}^{w \cdot x_j} - 1 = 0$

small $||w||$

$\text{subj}$ every matches truth

dist between

$\frac{2}{||w||_2^2}$ big   $\approx$   $\frac{||w||_2^2}{2}$ small

dist big

constrained opt.  $(\vec{\omega}, \underline{b})$  $\vec{\omega} \in \mathbb{R}^d, d=2$
$\swarrow d+1 = 3$

$$\mathcal{L}(X, S, \underline{\theta}) = \frac{1}{2}\|\omega\|_2^2 - \sum_{i=1}^{l} S_i\left(y_i(\omega^T x_i - b) - 1\right)$$

$\{(x_i, y_i)\}$

$\frac{1}{2}\omega \cdot \omega$

$y_1(\omega^T x_1 - b - 1) = 0$
$y_2(\omega^T x_2 - b) - 1 = 0$
$\vdots$
$y_n(\omega^T x_n - b - 1) = 0$

$$\nabla_\theta \mathcal{L} = \vec{0} \qquad = \frac{1}{2}\|\omega\|_2^2 - \sum_{i=1}^{l} S_i y_i(\omega^T \vec{x}_i - b) \quad \leftarrow \sum_{i=1}^{l} S_i$$

$$\nabla_\omega \mathcal{L}(X, S, \theta) = \omega - \sum S_i y_i \vec{x}_i = \vec{0} \quad \leftarrow d\text{-}eqns \text{ for } d \text{ features}$$

$$\nabla_b \mathcal{L} = \frac{\partial}{\partial b}\mathcal{L} = \sum S_i y_i = 0$$

$$W = \sum_{i \in C} S_i y_i x_i \quad ; \quad \sum_{i=1}^{n} S_i y_i = \vec{S} \cdot \vec{y} = 0$$

$$x_i \in \mathbb{R}^d$$

$$S \in \mathbb{R}^n \quad \in \mathbb{R}$$
$$n\text{-values}$$

"Trick" — KKT — Primal Dual form.

$$\text{Primal Form.} \Big\} \quad \min_{\theta, S} \quad \mathcal{L}(x, S, \theta) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^{n} S_i y_i (w^T x_i + b)$$

$$\boxed{K(x_i, x_j)}$$

$$\underline{\text{Dual Form}} \quad \max \quad \mathcal{L}_D(x, S, \theta) = \sum S_i - \frac{1}{2} \sum_{j=1}^{n} \sum_{i=1}^{n} S_i S_j y_i y_j (x_i \cdot x_j)$$

$$\text{s.t.} \nearrow \qquad w = \sum_{i=1}^{n} S_i y_i x_j \quad \text{and}$$

$$\sum_{i=1}^{n} S_i y_i = 0$$

$$\phi(x_i) \cdot \phi(x_j)$$
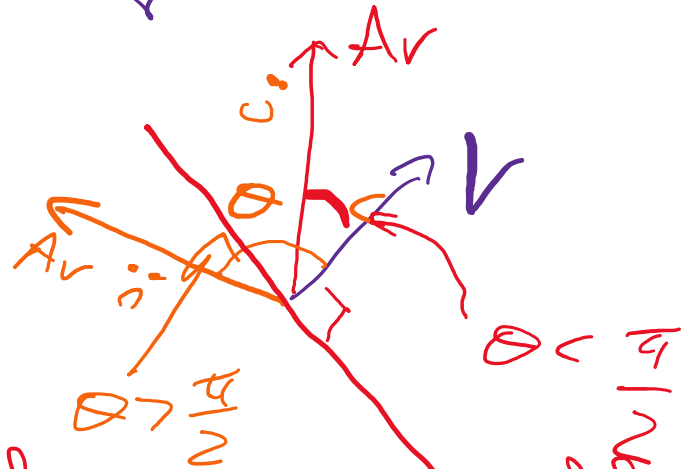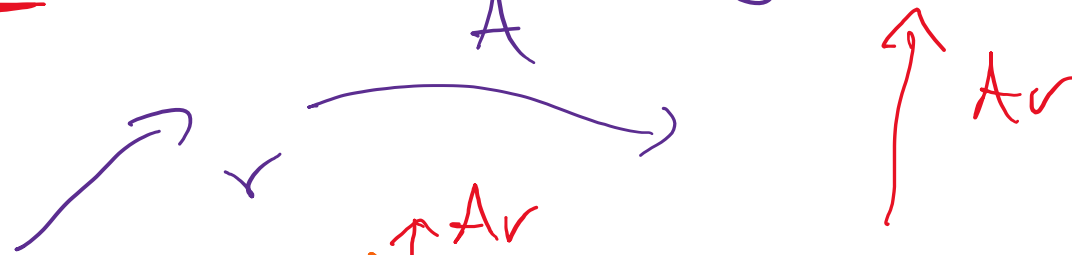
KKT

## Primal Problem:

$$\begin{cases} \text{minimize: } \mathcal{L}(x, s) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^{n} s_i y_i (w^T x_i + b) + \sum_{i=1}^{n} s_i \\ \text{such that: } s_i \geq 0, \forall i \end{cases}$$

## Dual Problem:

$$\begin{cases} \text{maximize: } \mathcal{L}_D(x, s) = \sum_{i=1}^{n} s_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} s_i s_j y_i y_j (\vec{x}_i^T \vec{x}_j) \\ \text{using: } w = \sum_{i=1}^{n} s_i y_i x_i, \text{ and } \sum_{i=1}^{n} s_i y_i = 0 \end{cases}$$

Pos. Semi-Defn Definition -

- A matrix $A_{nxn}$ is positive semi definite if
$v \cdot (A \cdot r) \geq 0$ for any $r^{\neq 0}$ in domain of A.

$A$

$r$

$Av$

$c_i Av$

$\theta$

$v$

$Av := -c_i$

$\theta > \frac{\pi}{2}$

$\theta < \frac{\pi}{2}$

$v \cdot (Ar) = \|v\| \|Ar\| \cos\theta$

- A kernel fn is pos. semi definite if $0 < 1, 3$, and $\overline{K}$ pos. semi.de matrix for any input set

$\mathcal{H}$ = Hilbert space – a complete inner product space

vector space

∘ a Hilbert space is a set $\mathcal{H}$ of vectors such that there is a

complete inner product.

dot

limits work properly.

## Spectral Decomp. thm:

Suppose $A_{n \times n}$ is pos. defn. symm. matrix with eigenvectors & eigenvalues $\lambda_i, v_i$,  $0 \le \lambda_1 \le \lambda_2 \le \ldots \le \lambda_n$

$$A = \sum_{i=1}^{n} \lambda_i \, v_i v_i^T$$

$v_i \otimes v_i = P_i$ rank-1 projector.

not $v_i^T v_i = v_i \cdot v_i$

$\circ$ $x_i \cdot x_j$ $\overset{\swarrow \overline{x} = \mathbb{R}^2}{\text{need the dot product between}}$

$\underline{x_i}$ & $\underline{x_j}$ to do SUM.

$\circ$ $\underline{K(x_i, x_j)} = \phi(x_i) \cdot \phi(x_j)$

$\underset{\text{corresponds.}}{\uparrow}$ $\int$ $\phi : \overline{x} \rightarrow \mathcal{H} \overset{\swarrow \mathbb{R}^6}{}$

Gram $\overset{\nearrow}{\underset{\text{matrix -}}{\text{-symmetric}}}$

$$\underline{K} = \begin{pmatrix} K(x_1, x_2) & K(x_2, x_3) & \_\_ & K(x_1, x_n) \\ \vdots & & & \\ \vdots & & & \\ K(x_n, x_1) & \_\_\_ & & K(x_n, x_n) \end{pmatrix}$$

KKT

## Primal Problem:

$$\begin{cases} \text{minimize: } \mathcal{L}(x, s) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^n s_i y_i(w^T x_i + b) + \sum_{i=1}^n s_i \\ \text{such that: } s_i \geq 0, \forall i \end{cases}$$

## Dual Problem:

$$\begin{cases} \text{maximize: } \mathcal{L}_D(x, s) = \sum_{i=1}^n s_i - \frac{1}{2}\sum_{i=1}^n \sum_{j=1}^n s_i s_j y_i y_j (\vec{x}_i^T \vec{x}_j) \\ \text{using: } w = \sum_{i=1}^n s_i y_i x_i, \text{ and } \sum_{i=1}^n s_i y_i = 0 \end{cases}$$

The amazing Kernel trick – nonlinear SVM through a kernel and all dot products in the high dimensional space
Done through a kerekl function

$X = \mathbb{R}^2$

Now, we define a kernel $K : X \times X \mapsto \mathbb{R}$, which can take different forms such

as:

- Linear kernel: $K(x, \tilde{x}) = x^T \tilde{x}.$
- Polynomial kernel: $K(x, \tilde{x}) = (x^T \tilde{x} - y)^d.$
- Gaussian RBF: $K(x, \tilde{x}) = e^{-\frac{Vertx - \tilde{x}\|^2}{2\sigma^2}}$

Consider the polynomial kernel, for $d = 2$, $X = \mathbb{R}^2$, then we have:

$$-\frac{\|x - \tilde{x}\|}{2\sigma^2}$$

Kernel is

- ⓪ —
- ① symmetric
- ② pos. semi-def.
- ③ cts. w.r.t. both inputs.

$$
\begin{aligned}
K(x, \tilde{x}) &= (x \cdot \tilde{x} + 1)^d \\
&= (x^T \tilde{x} + 1)^d \\
&= (x_1 \tilde{x}_1 + x_2 \tilde{x}_2 + 1)^2 \\
&= x_1^2 \tilde{x}_1^2 + 2x_1 \tilde{x}_1 + 2x_2 \tilde{x}_2 + x_1 \tilde{x}_1 x_2 \tilde{x}_2 + 1
\end{aligned}
$$

$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ ; $\tilde{x} = \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix}$

$K(x, \tilde{x}) = K(\tilde{x}, x)$

$+ x_2^2 \tilde{x}_2^2$

which interestingly can be re-written in terms of dot product:

$$
\begin{aligned}
K(x, \tilde{x}) &= (x \cdot \tilde{x} + 1)^d \\
&= (x^T \tilde{x} + 1)^d \\
&= (x_1 \tilde{x}_1 + x_2 \tilde{x}_2 + 1)^2 \\
&= x_1^2 \tilde{x}_1^2 + 2x_1 \tilde{x}_1 + 2x_2 \tilde{x}_2 + x_1 \tilde{x}_1 x_2 \tilde{x}_2 + 1 \quad + x_2^2 \tilde{x}_2^2
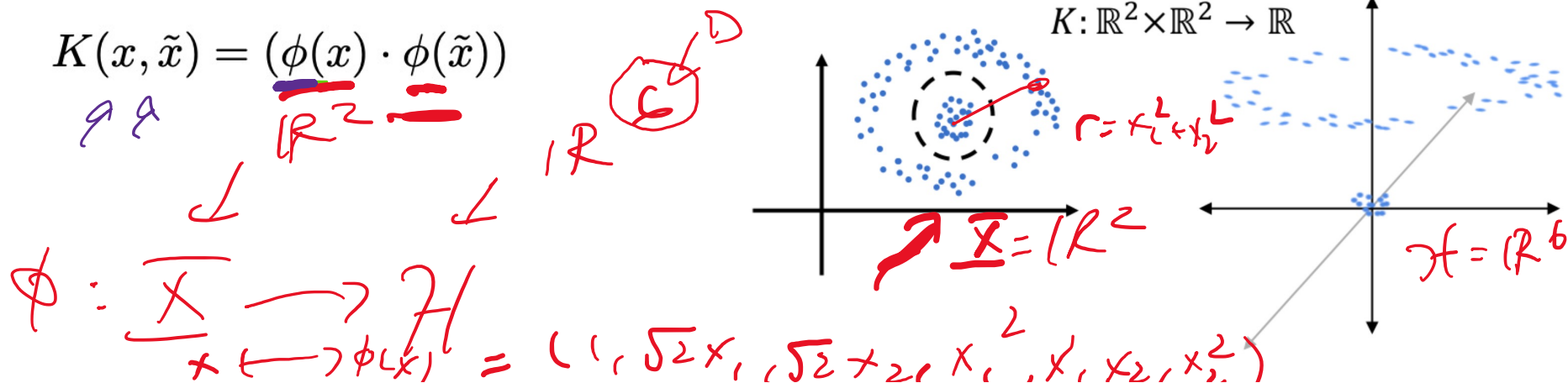\end{aligned}
$$

which interestingly can be re-written in terms of dot product:

$$
K(x, \tilde{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_1 x_2, x_2^2) \cdot (1, \sqrt{2}\tilde{x}_1, \sqrt{2}\tilde{x}_2, \tilde{x}_1^2, \tilde{x}_1 \tilde{x}_2, \tilde{x}_2^2)
$$

$$\phi: \mathbb{R}^2 \to \mathbb{R}^6$$
$$K: \mathbb{R}^2 \times \mathbb{R}^2 \to \mathbb{R}$$

$$
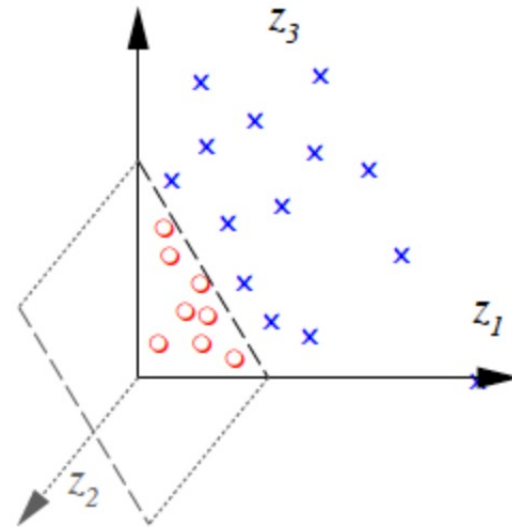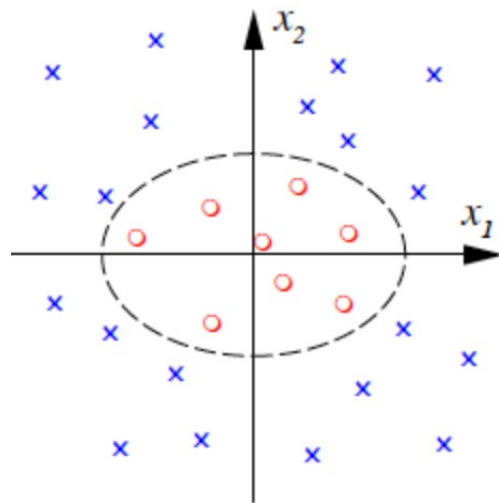K(x, \tilde{x}) = (\phi(x) \cdot \phi(\tilde{x}))
$$

$\mathbb{R}^2$    $\mathbb{R}$

$\phi : X \longrightarrow \mathcal{H}$

$x \longmapsto \phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_{2}, x_1^2, x_1 x_2, x_2^2)$

$r = x_1^2 + x_2^2$

$\underline{X} = \mathbb{R}^2$

$\mathcal{H} = \mathbb{R}^6$

$D = \infty$.

$n - 1$

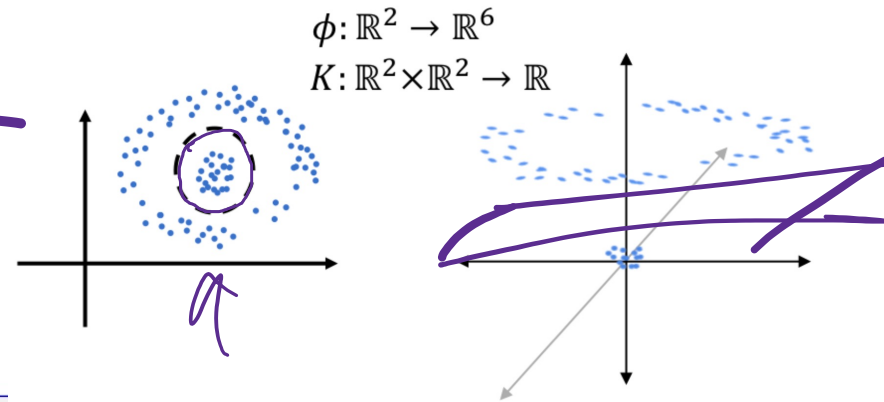Hyper Plane Classifier in Feature Space

$$\Phi : R^2 \to R^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{(2)}x_1x_2, x_2^2)$$

$$\phi(x_1, x_2) = (\phi_1(x_1, x_2), \phi_2(x_1, x_2), ..., \phi_6(x_1, x_2))$$

where $\phi : X \mapsto \mathcal{H}$.

$\phi : \mathbb{R}^2 \to \mathbb{R}^6$
$K : \mathbb{R}^2 \times \mathbb{R}^2 \to \mathbb{R}$



Note that $X = \mathbb{R}^2$ is the domain, and $\mathcal{H}$ is the Hilbert space, which is (in machine learning literature) the feature space, and a set of features $\phi_i, \forall i$, is called dictionary.

## Mantra

A major theme in machine learning is that sometimes things actually get easier in higher dimensions !!!.

- A linear plane in high dimensional feature space $\mathcal{H}$, may be a nonlinear curves in the domain space.
- $\mathcal{H}$ is a plane, with calculus with dot products is legit.

The following, we introduce Mercer's theorem, which generalizes spectral decomposition theorem.

$K : X \times X \longrightarrow \mathbb{R}$

**Theorem 5.5.1 — Mercer's Theorem** ~ Generalizes spectral decomp. thm.

. Let $K \in L^2(X \times X)$, (i.e. $\int |K(x, \tilde{x})|^2 dx d\tilde{x} < \infty$) such that $T : L_2(X) \mapsto L_2(X)$ by $(T(f)(x)) = \int K(x, \tilde{x}) f(\tilde{x}) d\tilde{x}$ is positive definite. If $\phi_i \in L^2(X)$ is

usual
$Av = U$
.e eigen
$Aw = \lambda w$

a normalized eigenfunction with eigenvalues $0 < \lambda_1, \le \lambda_2 \le \cdots \le \lambda_N$, Then

$\exists$ exist.

$$K(x, \tilde{x}) = \sum_{i=1}^{N_{\mathcal{H}}} \lambda_i \phi_i(x) \phi_i(\tilde{x}) \tag{5.20}$$

eigen
fn.

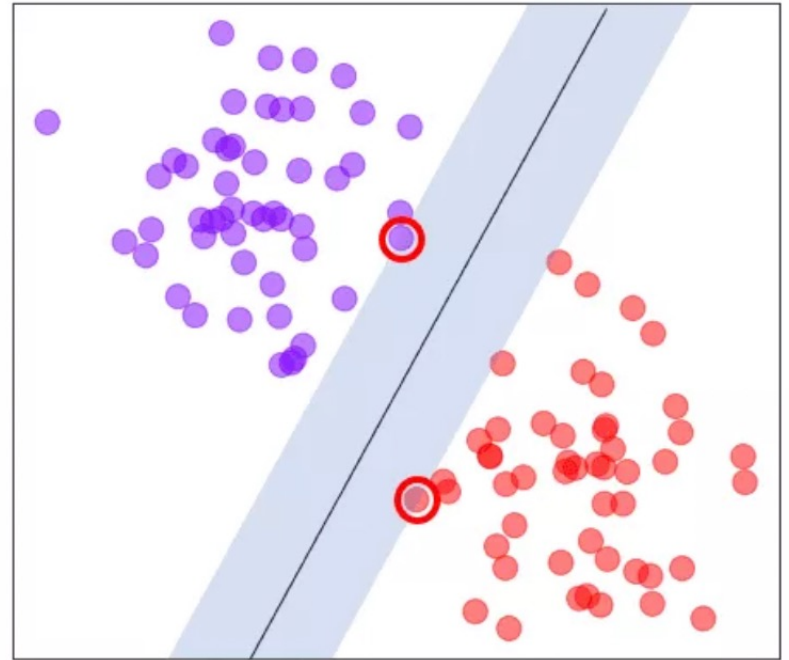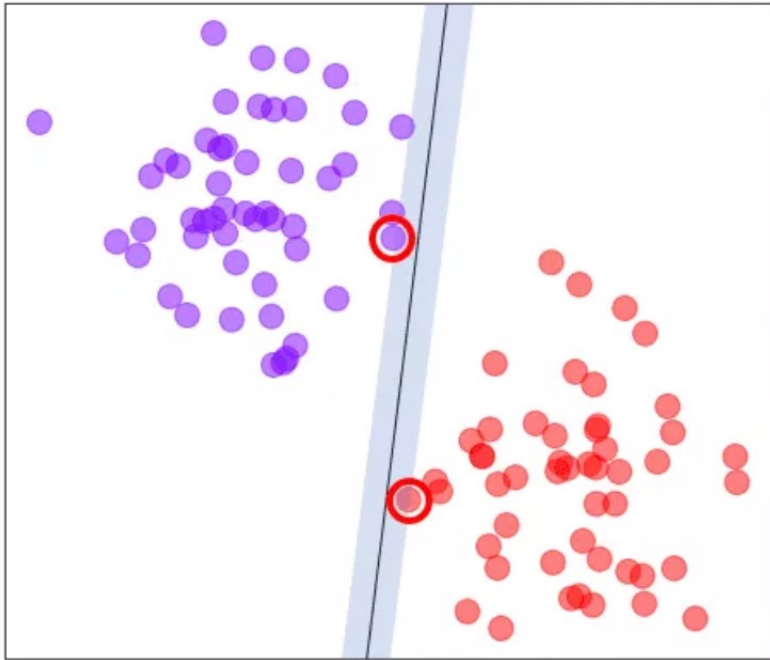$T(\phi_i)(x) = \lambda_i \phi_i(x) \Leftarrow$

for almost every $(x, \tilde{x})$. Where $N_{\mathcal{H}} = dim(\mathcal{H})$, and the convergance of $K(x, \tilde{x})$ is absolute.
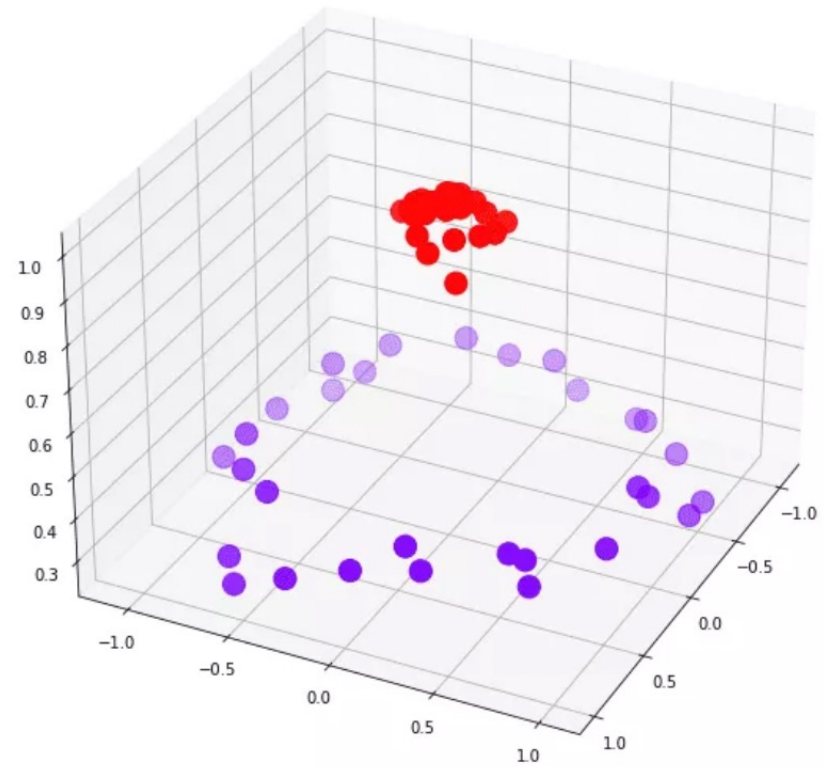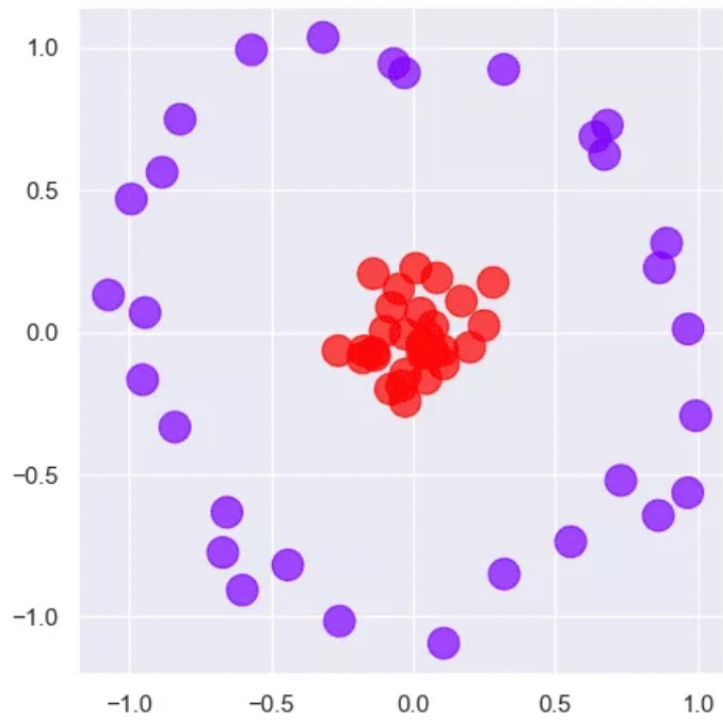
Mercer's theorem itself is a generalization of the result that any symmetric positive-semidefinite matrix is the Gramian matrix of a set of vectors.

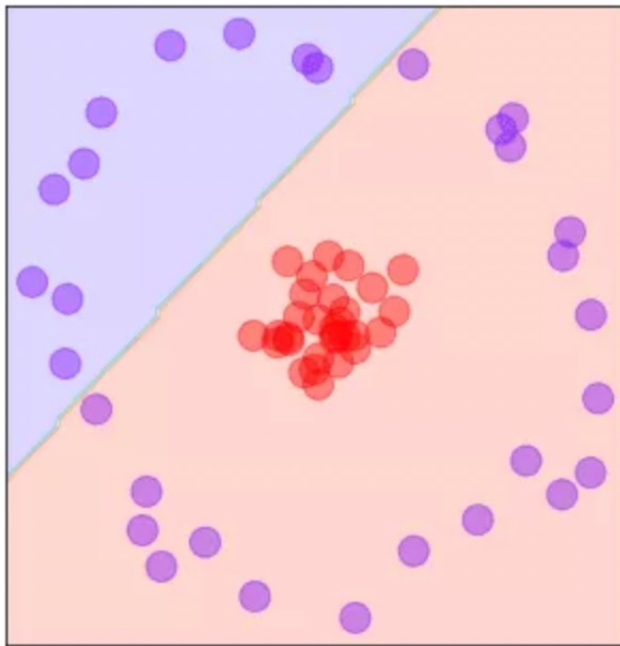$\phi_i$'s exist $\not\exists$ $\mathbb{I}$ can use them in $K \le X \times M$.

$$D = \{(x_i, y_i)\}$$



$X$

$y_i$

$x_i$

$X$

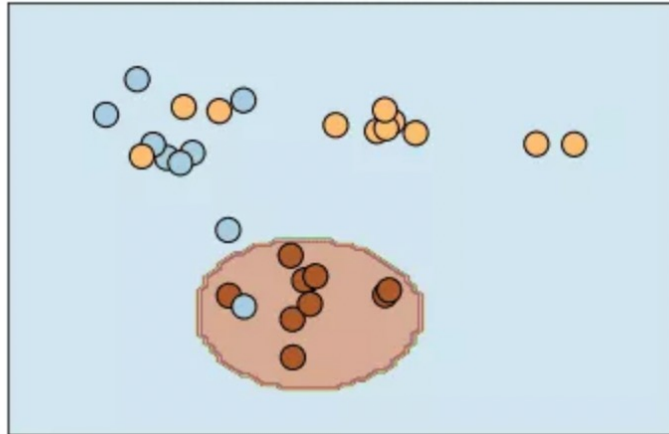Linear kernel       RBF kernel

kernel

Decision surface
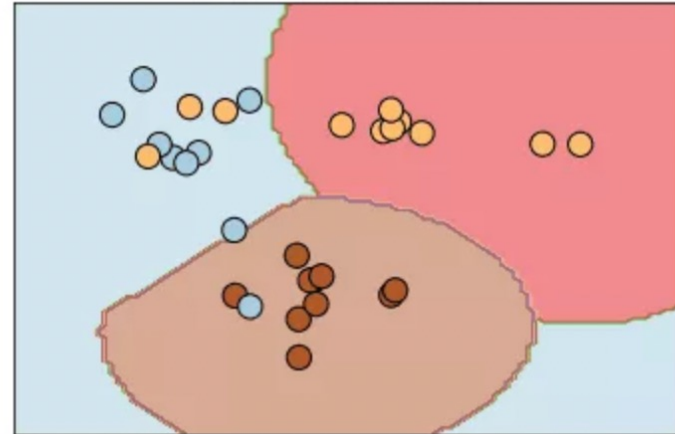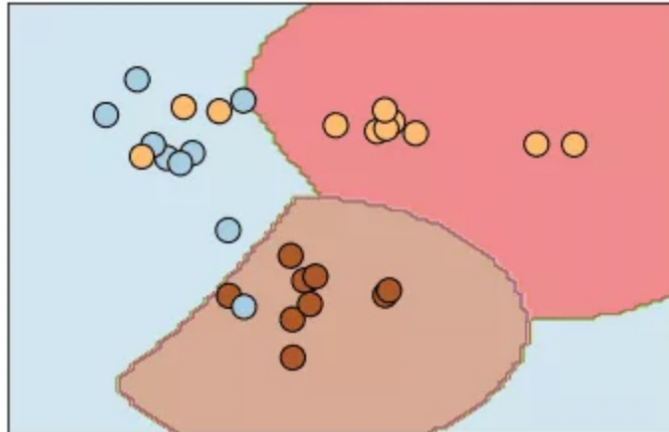
RBF kernel, C=0.1    RBF kernel, C=1
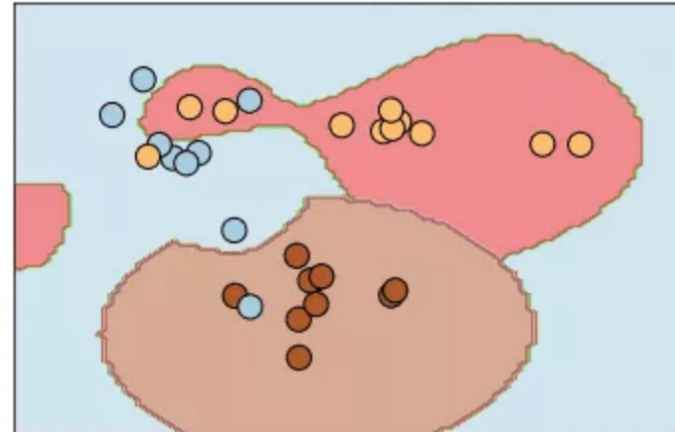
RBF kernel, C=10    RBF kernel, C=100

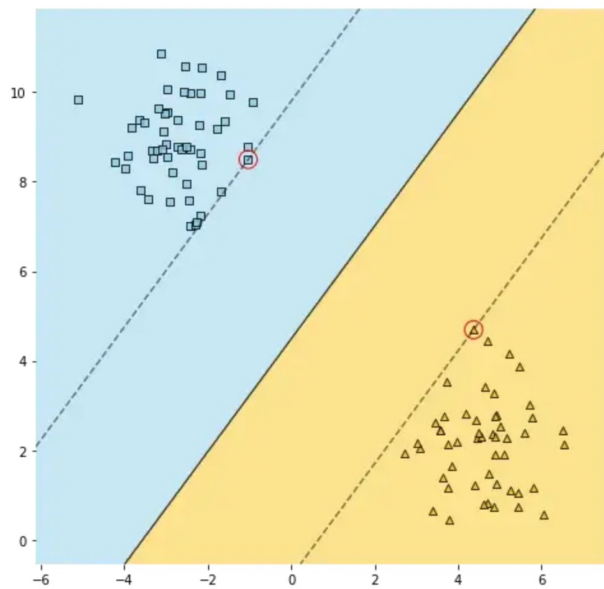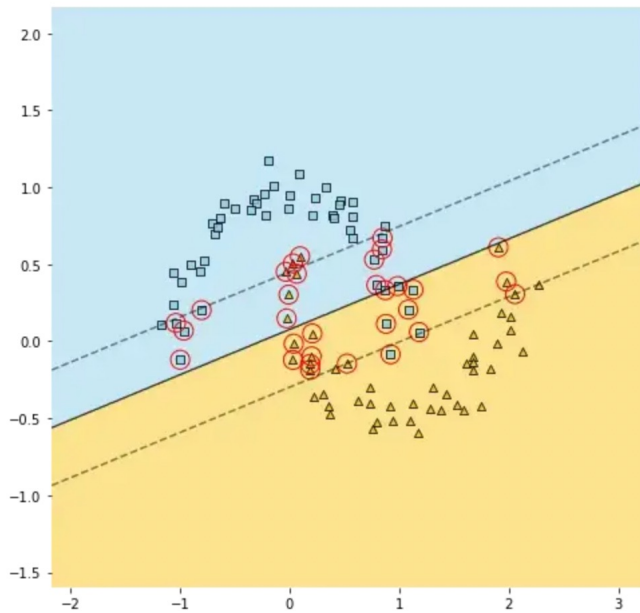George W Bush    Gerhard Schroeder    Donald Rumsfeld    Tony Blair    Donald Rumsfeld    Colin Powell    George W Bush    Colin Powell
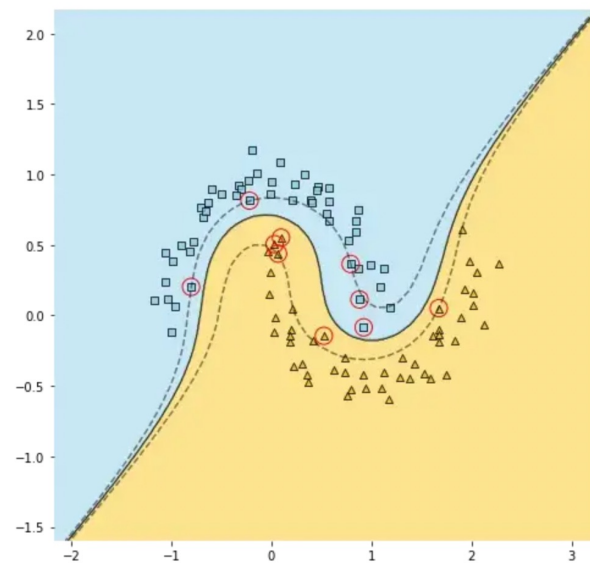
Linear SVM with linearly separable data works pretty well.



Linear SVM with linearly non-separable data does not work at all.



Decision boundary with a polynomial kernel.

**Theorem 2** (Mercer 1909). *Suppose $k \in L_\infty(\mathcal{X}^2)$ such that the integral operator $T_k : L_2(\mathcal{X}) \to L_2(\mathcal{X})$,*

$$T_k f(\cdot) := \int_{\mathcal{X}} k(\cdot, x) f(x) d\mu(x) \qquad (20)$$

*is positive (here $\mu$ denotes a measure on $\mathcal{X}$ with $\mu(\mathcal{X})$ finite and $\mathrm{supp}(\mu) = \mathcal{X}$). Let $\psi_j \in L_2(\mathcal{X})$ be the eigenfunction of $T_k$ associated with the eigenvalue $\lambda_j \neq 0$ and normalized such that $\|\psi_j\|_{L_2} = 1$ and let $\overline{\psi_j}$ denote its complex conjugate. Then*

1. *$(\lambda_j(T))_j \in \ell_1$.*
2. *$k(x, x') = \sum_{j \in \mathbb{N}} \lambda_j \overline{\psi_j(x)} \psi_j(x')$ holds for almost all $(x, x')$, where the series converges absolutely and uniformly for almost all $(x, x')$.*

Let's take it for granted that this is a valid positive semidefinite kernel. Let $k_{\texttt{poly(r)}}$ denote a polynomial kernel of degree $r$, and let $\gamma = 1/2$. Then

$$
\begin{aligned}
k_{\text{RBF}}(\mathbf{x}, \mathbf{y}) &= \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2\right) \\
&= \exp\left(-\frac{1}{2}\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y}\rangle\right) \\
&\overset{\star}{=} \exp\left(-\frac{1}{2}\left[\langle \mathbf{x}, \mathbf{x} - \mathbf{y}\rangle - \langle \mathbf{y}, \mathbf{x} - \mathbf{y}\rangle\right]\right) \\
&\overset{\star}{=} \exp\left(-\frac{1}{2}\left[\langle \mathbf{x}, \mathbf{x}\rangle - \langle \mathbf{x}, \mathbf{y}\rangle - \left[\langle \mathbf{y}, \mathbf{x}\rangle - \langle \mathbf{y}, \mathbf{y}\rangle\right]\rangle\right]\right) \\
&= \exp\left(-\frac{1}{2}\left[\langle \mathbf{x}, \mathbf{x}\rangle + \langle \mathbf{y}, \mathbf{y}\rangle - 2\langle \mathbf{x}, \mathbf{y}\rangle\right]\right) \\
&= \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right)\exp\left(-\frac{1}{2}\|\mathbf{y}\|^2\right)\exp\left(-2\langle \mathbf{x}, \mathbf{y}\rangle\right)
\end{aligned}
$$

Above, the two steps labeled $\star$ leverage the fact that

$$
\langle \mathbf{u} + \mathbf{v}, \mathbf{w}\rangle = \langle \mathbf{u}, \mathbf{w}\rangle + \langle \mathbf{v}, \mathbf{w}\rangle
$$

in general for inner products (see here){:target="_blank"}. Now let $C$ be a constant,

$$
C \equiv \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right)\exp\left(-\frac{1}{2}\|\mathbf{y}\|^2\right).
$$

and note that the Taylor expansion of $e^{f(x)}$ is

$$
e^{f(x)} = \sum_{r=0}^{\infty} \frac{[f(x)]^r}{r!}.
$$

We can write the RBF kernel as

$$
\begin{aligned}
k_{\text{RBF}}(\mathbf{x}, \mathbf{y}) &= C\exp\left(-2\langle \mathbf{x}, \mathbf{y}\rangle\right) \\
&= C\sum_{r=0}^{\infty} \frac{\langle \mathbf{x}, \mathbf{y}\rangle^r}{r!} \\
&= C\sum_{r}^{\infty} \frac{k_{\texttt{poly(r)}}(\mathbf{x}, \mathbf{y})}{r!}.
\end{aligned}
$$

So the RBF kernel can be viewed as an infinite sum over polynomial kernels. As $r$ increases, each polynomial kernel lifts the data into higher dimensions, and the RBF kernel is an infinite sum over these

SVD/POD/PCA/KL – unsupervised – ROM – structure of data (shape/geometry of data)– manifold learn – given $x_i$

DMD – supervised – forecasting – ROM – spectral analysis – structure of the process features are important, given $x_i$ -> $x_i, x_{i+1}$. mght as well call the $x_{i+1}=y_i$ – regression

Regression – onto general basis sets – supervised – find $y=f(x)$ given examples of $(x_i, y_i)$

Neural nets – classification of handwriting digits USPS – supervised,
        forecasting also regression to the flow function for forecasting
        auto-encoder is a unsupervised algoritghm using ANN with a bottleneck. – ROM
        random version was reservoir computing

Kmeans – clustering (given x develop labels – as y)  – partitioning the data –

LDA – linear discriminant analysis – Fischer 1936 – classification – given labeled data xI with labels yi learn $y=f(x)$

SVM – linear method for classification supervised – support vector machine
        kernelized version is nonlinear – reproducing kernel Hilbert space – KSVM – KSVD

Manifold learning – unsupervised – structure of the data – POD, autoencoder, ISOMAP, Diffusion Map

Regression, and classification – supervised