

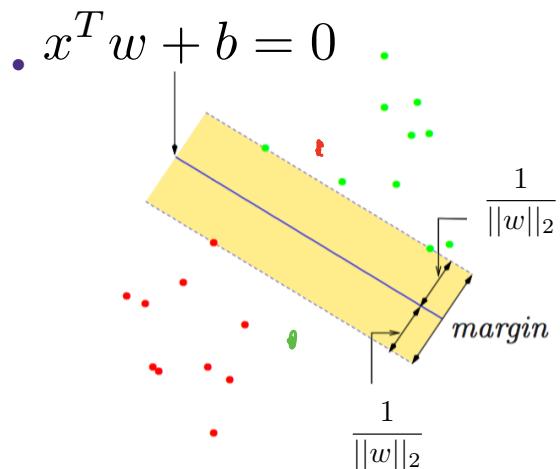
Kernels



UNIVERSITY *of* WASHINGTON

W

What if the data is not linearly separable?



Some points don't satisfy margin constraint:

$$\min_{w,b} \|w\|_2^2$$

$$y_i(x_i^T w + b) \geq 1 \quad \forall i$$

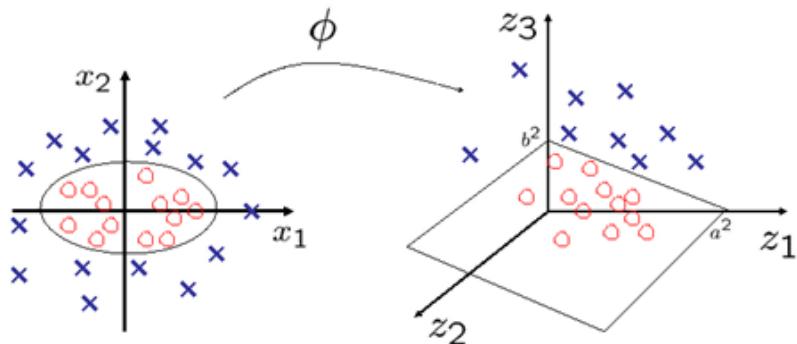
Two options:

1. Introduce slack to this optimization problem
2. Lift to higher dimensional space

What if the data is not linearly separable?

Use features of features of features...

poly - feature
 $\dim(\phi)$ is \exp in degree



How do we deal with high-dimensional lifts/data?

A fundamental trick in ML: use kernels

A function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a *kernel* for a map ϕ if $K(x, x') = \phi(x) \cdot \phi(x')$ for all x, x' .

$$\textcircled{1} \quad \textcircled{2} \quad \langle x_i, x_j \rangle$$

So, if we can represent our algorithms/decision rules as dot products and we can find a kernel for our feature map then we can avoid explicitly dealing with $\phi(x)$.

Examples of Kernels

- Polynomials of degree exactly d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^p$$

- Polynomials of degree up to d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^p$$

- Gaussian (squared exponential) kernel

(\mathcal{RBF})

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(u, v) = \tanh(\gamma \cdot u^T v + r)$$

The Kernel Trick

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

① Pick a kernel K

② For a linear predictor, show $w = \sum_{j=1}^n \alpha_j x_j, \alpha \in \mathbb{R}^n$

Change loss function/decision rule to only access data through dot products $\langle x_i, x_j \rangle$

Given a new data $x_{\text{new}} \in \mathbb{R}^d$

$$w^T x_{\text{new}} = \sum_{j=1}^n \alpha_j \langle x_j, x_{\text{new}} \rangle$$

\Rightarrow find α

$$\alpha^T \begin{pmatrix} K(x_1, x_{\text{new}}) \\ \vdots \\ K(x_n, x_{\text{new}}) \end{pmatrix}$$

Substitute $K(x_i, x_j)$ for $x_i^T x_j$

(Prediction:

$$\sum_{j=1}^n \alpha_j K(x_j, x_{\text{new}})$$

The Kernel Trick for regularized least squares

$$\hat{w} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda ||w||_w^2$$

There exists an $\alpha \in \mathbb{R}^n$: $\hat{w} = \sum_{i=1}^n \alpha_i x_i$

$$\begin{aligned}\hat{\alpha} &= \arg \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j \langle x_j, x_i \rangle)^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle x_i, x_j \rangle \\ &= \arg \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j K(x_i, x_j))^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) \\ &= \arg \min_{\alpha} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha\end{aligned}$$

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$

Why regularization?

↳ $\nabla f(x)^\top K \nabla f > 0$ generally, if $\phi(x)$ dimensions $> n$

Typically, $K \succ 0$. What if $\lambda = 0$?

$$\hat{\alpha} = \arg \min_{\alpha} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha$$

Why regularization?

X : $n \times d$ data
 y : n labels
 $\mathcal{K} = \begin{pmatrix} K(x_1, x_1) & \dots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \dots & K(x_n, x_n) \end{pmatrix}$

Typically, $\mathbf{K} \succ 0$. What if $\lambda = 0$?

$$\hat{\alpha} = \arg \min_{\alpha} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha$$

$$\forall x, x' \quad K(x, x') = K(x', x)$$

Unregularized kernel least squares can (over) fit **any data!**

\hat{y} : predictions on y
 $\in \mathbb{R}^n$

$$\hat{\alpha} = \mathbf{K}^{-1} \mathbf{y}$$
$$\hat{y}_i = \begin{pmatrix} K(x_1, x_i) \\ \vdots \\ K(x_n, x_i) \end{pmatrix}^T \mathbf{z}$$

$$\hat{y} = \mathbf{K} \mathbf{z} = \mathbf{K} \mathbf{K}^{-1} \mathbf{y} = \mathbf{y}$$

The Kernel Trick for SVMs

$$w = \sum_{j=1}^n \alpha_j x_j$$

$$\min_w \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i (x_i^T w + b)\} + \lambda \|w\|_2^2$$

$$\Leftrightarrow \min_{\alpha, b} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - y_i \left(\sum_{j=1}^n \alpha_j \langle x_i, x_j \rangle + b \right) \right\} + \lambda \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^n \alpha_i \alpha_j \langle x_i, x_j \rangle$$

replace $\langle x_i, x_j \rangle \rightarrow K(x_i, x_j)$

$$\Rightarrow \min_{\alpha, b} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - y_i \left(\sum_{j=1}^n \alpha_j K(x_i, x_j) + b \right) \right\} + \lambda \alpha^T K \alpha$$

Convex in α, b

α_j are learned
Decision rule:

x_{new}

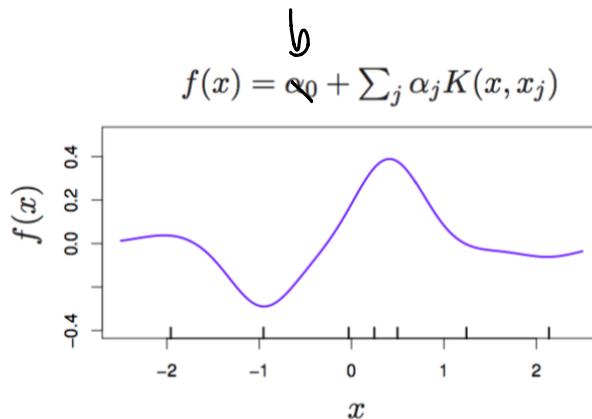
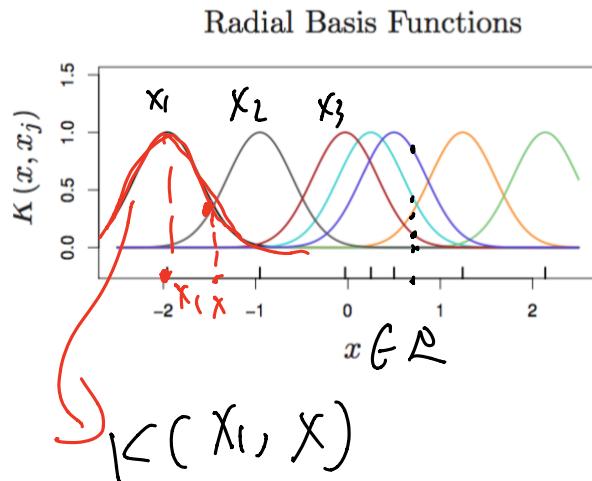
$$\sum_{j=1}^n \alpha_j K(x_j, x_{\text{new}}) + b$$

$> 0 \Rightarrow \text{predict } 1$
 $\leq 0 \Rightarrow \text{predict } 0$

RBF Kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2}\right)$$
$$\chi, \chi_i \in \mathcal{L}$$

This is like weighting “bumps” on each point

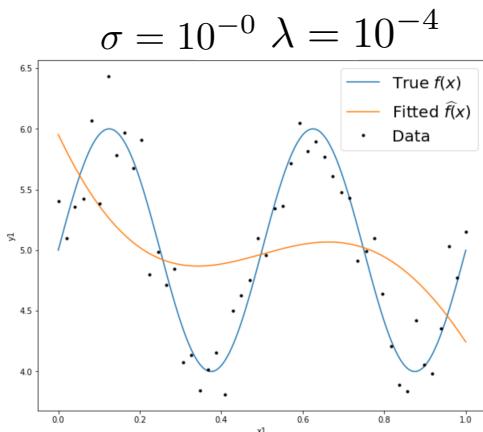
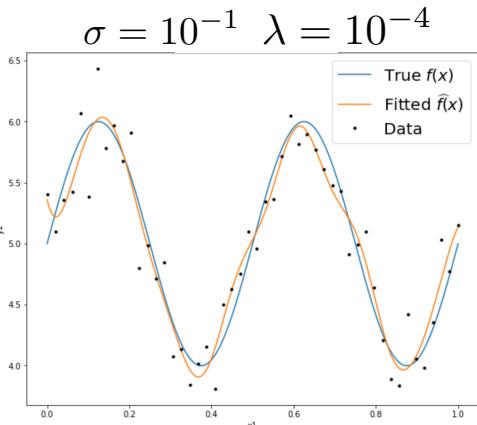
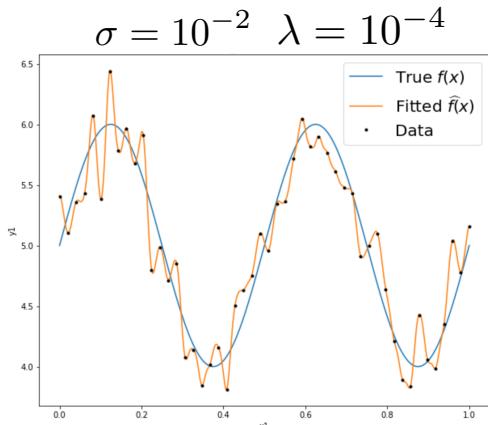


RBF Kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2}\right)$$

$\sigma \rightarrow 0 \Rightarrow K(\mathbf{u}, \mathbf{v}) = 0$
 $\sigma \rightarrow \infty \Rightarrow K(\mathbf{u}, \mathbf{v}) = 1$

The bandwidth sigma has an enormous effect on fit:



bias-variance; $\uparrow \sigma$: longer bias \downarrow variance

as $\sigma \rightarrow 0$, only $x = x_i$, $K(x_i, x) \neq 0$, $x \neq x_i$, $K(x_i, x) = 0$

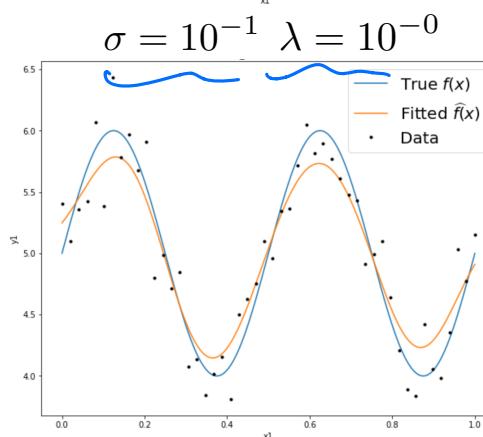
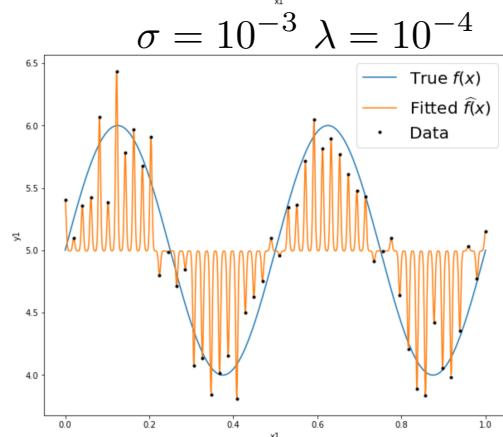
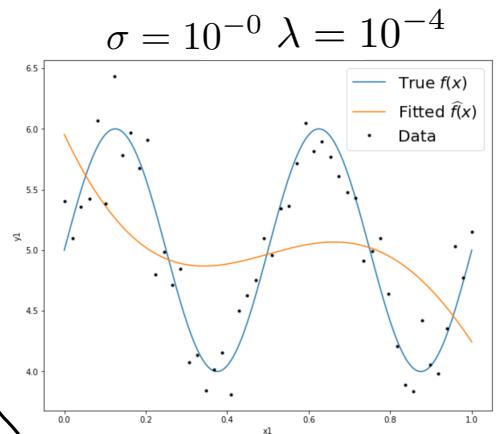
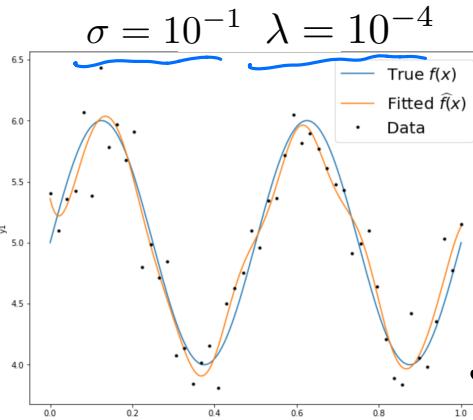
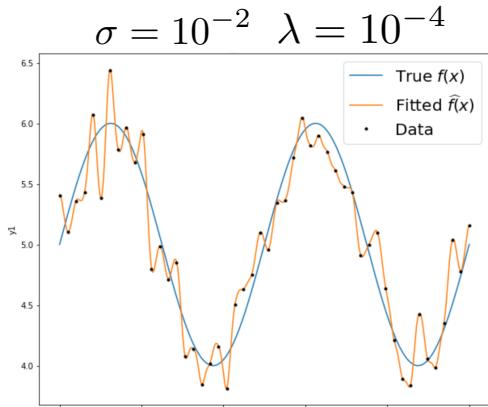
$\sigma \rightarrow \infty$, $K(x_i, x) \approx 1$, $\forall i$
average over all y

$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i K(x_i, x)$$

RBF Kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2}\right)$$

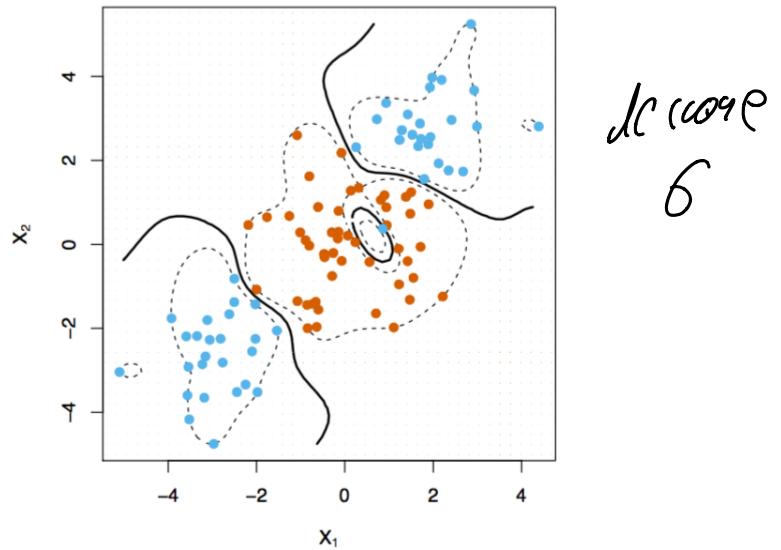
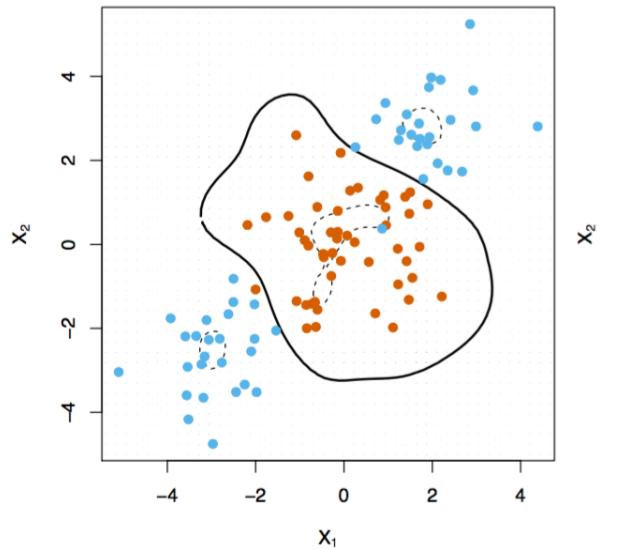
The bandwidth sigma has an enormous effect on fit:



$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i K(x_i, x)$$

RBF kernel and random features

$$\widehat{w} = \sum_{i=1}^n \max\{0, 1 - y_i(b + x_i^T w)\} + \lambda \|w\|_2^2$$
$$\min_{\alpha, b} \sum_{i=1}^n \max\{0, 1 - y_i(b + \sum_{j=1}^n \alpha_j K(x_i, x_j))\} + \lambda \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j)$$



RBF kernel and random features

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2}\right)$$

If n is very large, allocating an n -by- n matrix is tough.

$$\alpha = (K + \lambda I)^{-1} y, \quad K \in \mathbb{R}^{n \times n}$$

n : can't
not feasible

$$RBF: \phi_{\omega}, \phi \approx \phi_{\omega}, \phi(x)^T \phi(x) \approx \phi_{\omega}(x)^T \phi(x) = k(x, x)$$

RBF kernel and random features

$$\mathcal{U} = [y], P \in \mathbb{R}^P$$

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2}\right)$$

If n is very large, allocating an n -by- n matrix is tough.

use sampling from
inf-dim feature space

$$\begin{aligned} 2\cos(\alpha)\cos(\beta) &= \cos(\alpha + \beta) + \cos(\alpha - \beta) \\ e^{jz} &= \cos(z) + j\sin(z) \end{aligned}$$

\Rightarrow apply linear regression

$$\phi(x) = \begin{bmatrix} \sqrt{2} \cos(w_1^T x + b_1) \\ \vdots \\ \sqrt{2} \cos(w_p^T x + b_p) \end{bmatrix} \in \mathbb{R}^P \quad w_k \sim \mathcal{N}(0, 2\gamma I) \quad / \text{SVM on } \phi \\ b_k \sim \text{uniform}(0, \pi)$$

$$\mathbb{E}_{w,y} \left[\frac{1}{P} \sum_{k=1}^P \phi(x)^T \phi(x) \right] = \frac{1}{P} \sum_{k=1}^P \mathbb{E}_{w,y} \left[2 \cos(w_k^T x + b_k) (w_k^T x + b_k) \right]$$

$$= \mathbb{E} \left[(w^T x) + b \right] (w^T x) + b$$

$$= e^{-\gamma \|x-y\|_2^2}$$

$$x \in \mathbb{R}^d, w_k \in \mathbb{R}^d$$

$$\gamma = \frac{1}{2\sigma^2}$$

[Rahimi, Recht NIPS 2007]
“NIPS Test of Time Award, 2018”

$x, \phi(x) \in \mathbb{R}^n$

\downarrow
 $\phi(x, x') , \mathcal{O}(n^3)$

\downarrow
sampled feature
 $\tilde{\phi} \in \mathbb{R}^p$

$p \ll n$

Apply SVM on $\tilde{\phi}$

String Kernels

Example from Efron and Hastie, 2016

Amino acid sequences of different lengths:

x1 IPTSALVKETLALLSTHRTLLIANETLRIPVPVHKNHQLCTEIFQGIGTLESQTVQGGTV
ERLFKNLSLIKYYIDGQKKCGEERRRVNQFLDY**LQE**FLGVMNTEWI

x2 PHRRDLCRSRSIWLARKIRSDLTALTESYVKHQGLWSELTEAER**LQE**NLQAYRTFHVLLA
RLLEDQQVHFTPTEGDFHQAIHTLLLQVAAFAYQIEELMILLEYKIPRNEADGMLFEKK
LWGLKV**LQE**LSQWTVRSIHDLRFISSHHQTGIP

All subsequences of length 3 (of possible 20 amino acids) $20^3 = 8,000$

$$h_{\text{LQE}}^3(x_1) = 1 \text{ and } h_{\text{LQE}}^3(x_2) = 2.$$

'Graph, time-service) , etc - - -

Fixed Feature V.S. Learned Feature

Linear function with fixed feature

$$x \Rightarrow \phi(x)$$

ϕ : hand-crafted , $\phi(x)^T w$

Can we learn ϕ from data

\Rightarrow neural net