

Introduction to Machine Learning

Foundations and Applications

Paul J. Atzberger
University of California Santa
Barbara



Support Vector Machines

Support Vector Machines: Motivations

Consider data: $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, with features x , labels y .

Example: $x \in \mathbb{R}^N, y \in \{-1, +1\}$, with $x = \text{image}, y = +1 \rightarrow \text{Apple}, y = -1 \rightarrow \text{Orange}$.

Task: Find hyperplane that separates points x_i having different labels y_i .

Challenges:

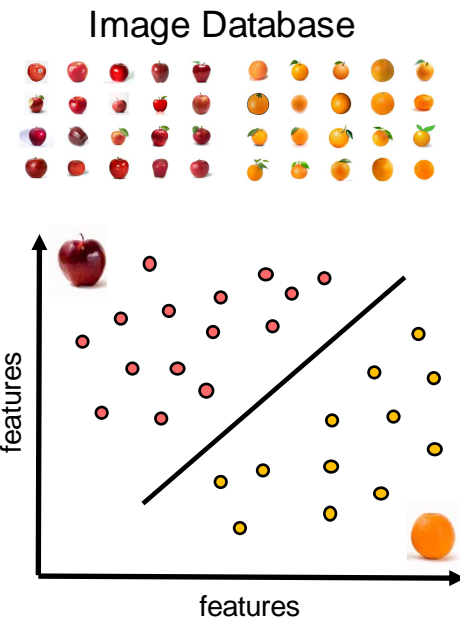
What algorithms can be used to find hyperplanes from data?

Many hyperplanes are possible. Which may have the best generalization?

What if the data is not separable?

How do we precisely define “separation” and the classification task?

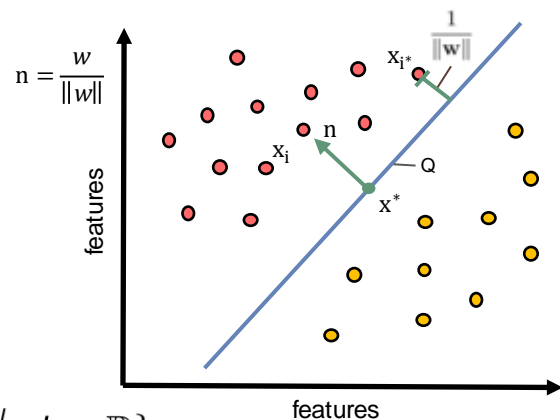
Approach: Support Vector Machines + Kernel Methods.



Support Vector Machines

SVM: Optimization Problem (Primal \mathcal{P})

$$\begin{cases} \min_{w,b} \frac{1}{2} \|w\|^2 \\ \text{subject to } y_i (w \cdot x_i + b) \geq 1. \end{cases}$$



Separable Case: $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^m$, $\mathcal{H} = \{h \mid h(x) = \text{sign}(w^T x + b), w \in \mathbb{R}^N, b \in \mathbb{R}\}$.

Definition: The data is **separable** if there exists $h \in \mathcal{H}$ so that $h(x_i) = \text{sign}(w^T x_i + b) = y_i$, $i \in \{1, 2, \dots, m\}$.

Hyperplane: $Q = \{x \mid w_0^T x + b_0 = 0\} = \{x \mid c^{-1} w_0^T x + c^{-1} b_0 = 0\}$.

Require: $\min_{x_i} |w^T x_i + b| = 1$ for the w used for a given data set \mathcal{S} .

Result: We will always have $y_i(w^T x + b) \geq 1$.

Definition: The **geometric margin** $\rho(x)$ of a point x is the distance to the hyperplane Q .

Let x^* be s.t. $-w^T x^* = b$, then $w^T x + b = w^T (x - x^*)$.

$$\Rightarrow \left| \frac{w}{\|w\|} (x - x^*) \right| = \frac{|w^T x + b|}{\|w\|} = \rho(x).$$

Consequence:, the closest data point x_{i^*} has distance $\rho(x_{i^*}) = 1/\|w\|$.

SVM Separable Case: Summary

Consider a data set $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where x denotes features and y denotes labels.

Example: $x \in \mathbb{R}^N, y \in \{-1, +1\}$,
with $x = \text{image}, y = +1 \rightarrow \text{Apple}, y = -1 \rightarrow \text{Orange}$.

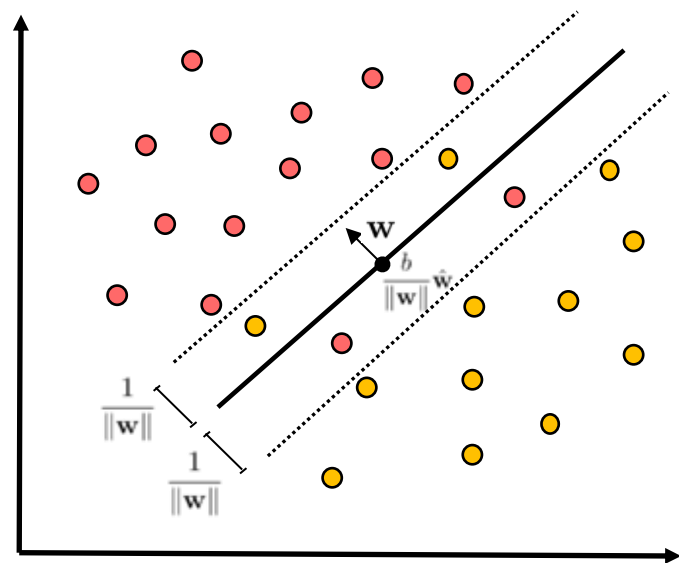
Find hyperplane separating points x_i having different labels y_i and with the **"greatest margin"** (helps with generalization).

Find parameters w, b that optimize

$$\min_{w, b} \frac{1}{2} w^T w$$

$$\text{subject to } y_i (w^T \phi(x_i) + b) \geq 1 \quad (\text{for now, } \phi(x_i) = x_i)$$

This **assumes** data is **separable**. Minimizing w maximizes the margin. Classifier $h(x) = \text{sign}(w^T \phi(x) + b)$



What if data is not separable?

Summary: SVM Non-Separable Case

Case of data that is not separable?

Find hyperplane and with **biggest margin** that minimizes extent of misclassifications.

Introduce **"slack variables"** ξ for the constraint.

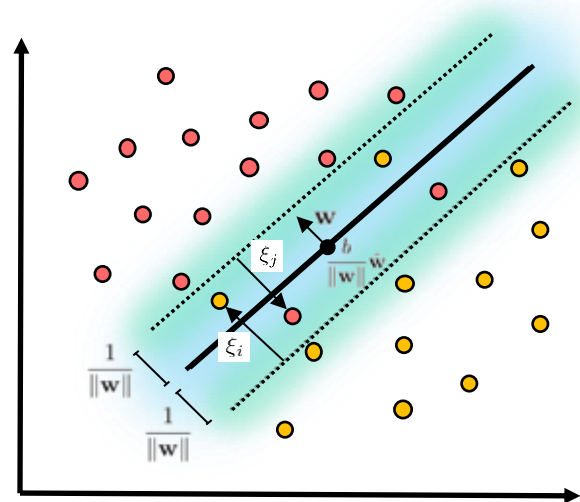
Find parameters w, b, ξ that optimize

$$\min_{w, b, \xi_i} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i$$

$$\text{subject to } y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i \quad \xi_i \geq 0 \quad (\text{for now, } \phi(\mathbf{x}_i) = \mathbf{x}_i)$$

Tries to find hyperplane and margin that minimizes the total amount training data points violate the constraint.

C is crucial regularization parameter determining penalty for violating the constraint.



Insights into generalization using results from optimization (duality).

Can apply more generally using kernel methods.



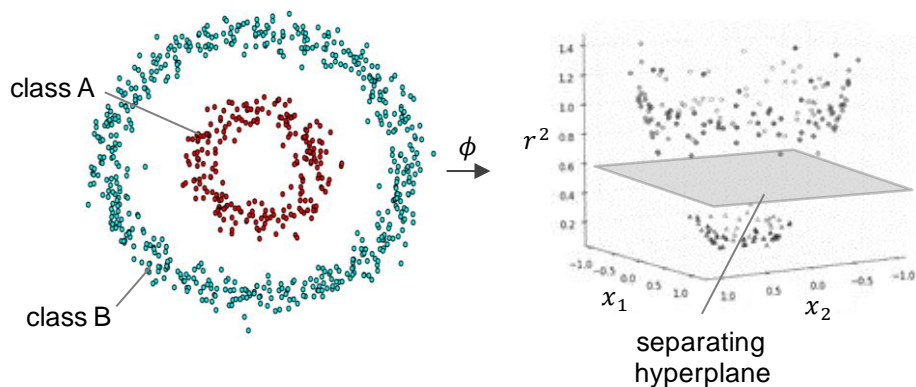
Kernel Methods Overview

Support Vector Machines (Kernel Trick)

Data is often not separable.

Mapping points to higher dimensional spaces they can become separable.

Example:



map ϕ
$\tilde{x} = \phi(x) = (x_1, x_2, r^2)$
$r^2 = x_1^2 + x_2^2$

Kernel: associated to map ϕ is an inner-product $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$.

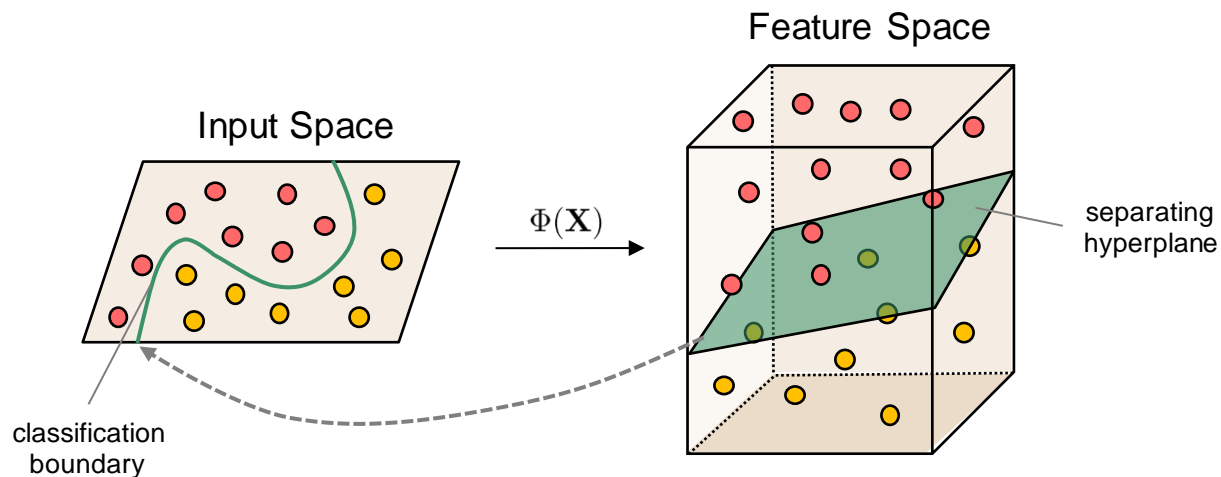
Example: $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = x_{i,1}x_{j,1} + x_{i,2}x_{j,2} + r_i^2r_j^2$

More generally...

Support Vector Machines (Kernel Trick)

Data is often not separable.

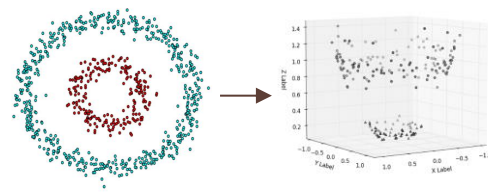
Mapping points to higher dimensional spaces they can become separable.



Kernel: associated to map ϕ is an inner-product $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$.

Example: $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = x_{i,1}x_{j,1} + x_{i,2}x_{j,2} + r_i^2 r_j^2$

Support Vector Machines (Kernel Trick)



Data is often not separable but can become **separable in higher dimensional spaces**.

Inner-products can be replaced by **kernel** $\langle x_i, x_j \rangle \rightarrow K(x_i, x_j)$.

Kernel should be **symmetric, positive definite**, L^2 : $\sum_{i,j=1}^N a_i a_j K(x_i, x_j) \geq 0$, $\int_X \int_X K(x, t)^2 d\mu(x) d\mu(t) < \infty$

Consider: $L_K : L^2_\mu(X) \rightarrow L^2_\mu(X)$, $L_K f(x) = \int_X K(x, t) f(t) d\mu(t)$ has countable set of non-negative eigenvalues $\{\lambda_k\}_{k=1}^\infty$.

Theorem (Mercer 1909): An L^2 kernel $K(x, t)$ that is symmetric positive definite can be represented as the product

$$K(x, t) = \sum_{k=1}^{\infty} \lambda_k \phi_k(x) \phi_k(t) \quad \text{Let } \Phi(x) = [\Phi_k(x)], \text{ then } K(x, t) = \langle \Phi(x), \Phi(t) \rangle.$$

Consequence: $K(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$, so $k(v, x_i) = w^T \Phi(x_i)$ as appears in the SVM constraints.

Note, only action $K(x, y)$ is needed in SVM, so no need to map explicitly to feature space $\Phi(x)$.

Kernel Methods: Hilbert-Schmidt and Mercer Theorem



Theorem (Hilbert-Schmidt): For L_K a **self-adjoint compact operator** there is a **countable complete orthonormal basis** $\{\phi_i\}$ for $L^2(\mathcal{Z})$ so that $L_K\phi_i = \lambda_i\phi_i$ with $\lambda_i \rightarrow 0$.

Theorem (Mercer 1909): An L^2 -kernel $K(x,t)$ that is **symmetric positive definite** can be represented as the product

$$K(x, t) = \sum_{k=1}^{\infty} \lambda_k \phi_k(x) \phi_k(t) \quad \text{with } \lambda_i > 0, \lambda_1 \geq \lambda_2 \geq \lambda_3 \dots, \text{ and } \lambda_i \rightarrow 0.$$

Remark: We can interpret the **Mercer Theorem** as stating **there exists a non-linear transformation** $z = \Phi(x)$ related to the kernel K as follows. Let $[\Phi(x)]_k = \sqrt{\lambda_k} \phi_k(x)$ then $K(x, t) = \langle \Phi(x), \Phi(t) \rangle_{\ell^2}$. Also can represent using Reproducing Kernel Hilbert Space (RKHS) (later).

Consequence: This shows that **if a kernel is L^2 and symmetric positive definite** then we can interpret it as being the **inner-product** associated with **some non-linear transformation Φ of the data!** For instance, $\mathbb{R}^N \rightarrow \ell^2$ or later $\mathbb{R}^N \rightarrow RKHS$.

Remark: To compute the inner-product we do **not** need to use $\langle \Phi(x), \Phi(t) \rangle_{\ell^2}$ which could be expensive, instead **we only need to evaluate kernel $K(x, t)$** . **This called the kernel trick!**

Support Vector Machines (Kernel Trick)

Kernels provide sensitivity to different features of the data, $K(x,t) = \langle \Phi(x), \Phi(t) \rangle$.

Popular Kernels:

Linear:

$$K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle$$

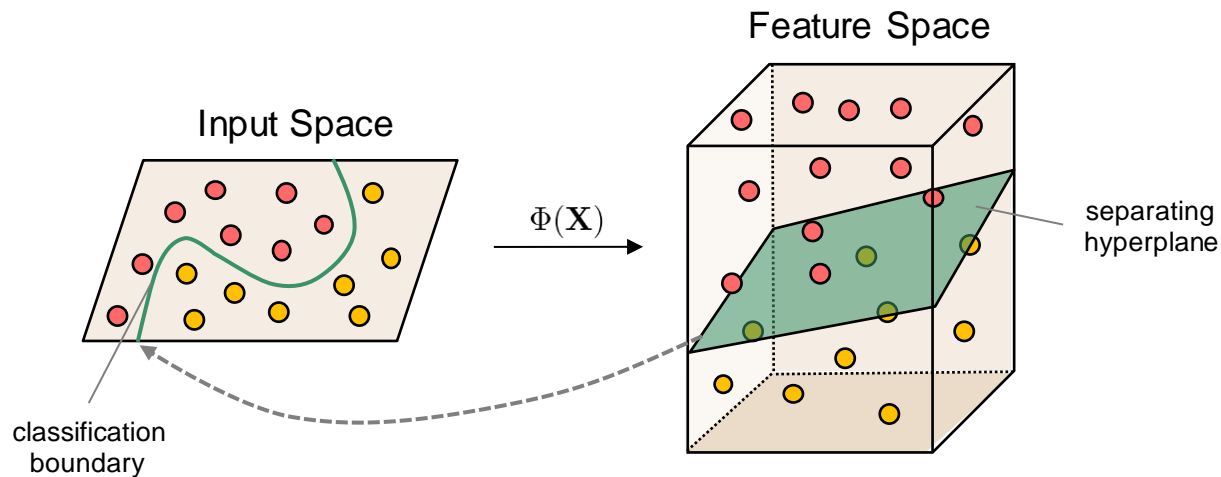
Radial Basis Function (RBF):


$$K(\mathbf{x}, \mathbf{y}) = \exp[-\gamma \|\mathbf{x} - \mathbf{y}\|^2]$$

Polynomial (degree d):

$$K(\mathbf{x}, \mathbf{y}) = (\gamma \langle \mathbf{x}, \mathbf{y} \rangle + r)^d$$

Lots of other choices possible.





Optimization Theory

Optimization

Constrained Optimization Problem (Primal \mathcal{P})

$$\begin{cases} \min_{x \in \mathcal{X}} f(x) \\ \text{subject to } g_i(x) \leq 0. \end{cases}$$

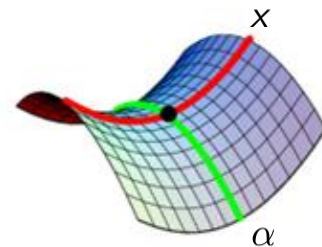
Definition: The **Lagrangian** \mathcal{L} of \mathcal{P} is

$$\mathcal{L}(x, \alpha) = f(x) + \sum_{i=1}^m \alpha_i g_i(x) \rightarrow \mathcal{L}(x, \alpha) = f(x) + \alpha \cdot g(x), \quad \text{where } \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_M \end{bmatrix}, \quad g(x) = \begin{bmatrix} g_1(x) \\ \vdots \\ g_M(x) \end{bmatrix}.$$

Definition: A **saddle-point** (x^*, α^*) of the Lagrangian \mathcal{L} is a point satisfying $\mathcal{L}(x^*, \alpha) \leq \mathcal{L}(x^*, \alpha^*) \leq \mathcal{L}(x, \alpha^*)$, holding for $\forall x \in \mathcal{X}, \alpha \geq 0$.

Theorem

For **constrained optimization problem** \mathcal{P} , a **saddle-point** (x^*, α^*) of the Lagrangian \mathcal{L} is a solution of \mathcal{P} .



Optimization

Constrained Optimization Problem (Primal \mathcal{P})

$$\begin{cases} \min_{x \in \mathcal{X}} f(x) \\ \text{subject to } g_i(x) \leq 0. \end{cases}$$

Lagrangian

$$\mathcal{L}(x, \alpha) = f(x) + \alpha \cdot g(x)$$

Theorem

For **constrained optimization problem** \mathcal{P} , a **saddle-point** (x^*, α^*) of the Lagrangian \mathcal{L} is a solution of \mathcal{P} .

Proof: From $\mathcal{L}(x^*, \alpha) \leq \mathcal{L}(x^*, \alpha^*) \Rightarrow \alpha \cdot g(x^*) \leq \alpha^* \cdot g(x^*), \forall \alpha \geq 0, \Rightarrow g(x^*) \leq 0$.

If $g_i(x^*) > 0$ then take α_i large so that $\alpha_i \cdot g_i(x^*) > c_i$ for any given $c_i \in \mathbb{R}$. Let $c = \alpha^* \cdot g(x^*)$.

Furthermore, $\alpha^* \cdot g(x^*) = 0$. Consider $\alpha \rightarrow 0$ then $0 \leq \alpha^* \cdot g(x^*) \leq 0, \Rightarrow \alpha^* \cdot g(x^*) = 0$.

From $\mathcal{L}(x^*, \alpha^*) \leq \mathcal{L}(x, \alpha^*), \forall x \Rightarrow f(x^*) \leq f(x) + \alpha^* \cdot g(x)$ for all x s.t. $g(x) \leq 0$.

We have $f(x^*) \leq f(x)$ so (x^*, α^*) solves \mathcal{P} . ■

Optimization

Definition: Strong Constraint Qualification (Slater's Condition)

$\exists \bar{x} \in \text{interior}(\mathcal{X}), \forall i \in \{1, 2, \dots, m\}, g_i(\bar{x}) < 0.$

Definition: Weak Constraint Qualification (Weak Slater's Condition)

$\exists \bar{x} \in \text{interior}(\mathcal{X}), \forall i \in \{1, 2, \dots, m\}, (g_i(\bar{x}) < 0) \vee (g_i(\bar{x}) = 0 \wedge g_i(\bar{x}) = a\bar{x} + b \text{ (affine)}).$

Theorem: (when saddle point is necessary w/ strong slater)

Let f, g be **convex functions** with **strong slater condition** holding.

If x^* is a solution to \mathcal{P} then $\exists \alpha^* \geq 0$ s.t. (x^*, α^*) satisfies the saddle condition for \mathcal{L} .

Theorem: (when saddle point is necessary w/ weak slater)

Let f, g be **convex** and **differentiable** functions with **weak slater condition** holding.

If x^* is a solution to \mathcal{P} then $\exists \alpha^* \geq 0$ s.t. (x^*, α^*) satisfies the saddle condition for \mathcal{L} .

Optimization:

Theorem: Karush-Kuhn-Tucker (KKT) Conditions

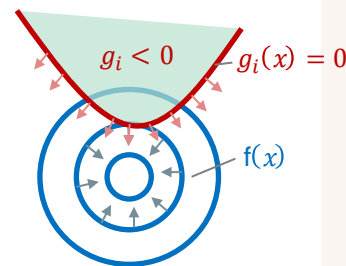
Let f, g_i be **convex** and **differentiable** functions where the **weak constraint qualification** is satisfied.

A \bar{x} is a solution to the **constrained optimization problem** \mathcal{P} if and only if $\exists \bar{\alpha} \geq 0$ s.t.

$$\nabla_x \mathcal{L}(\bar{x}, \bar{\alpha}) = \nabla_x f + \bar{\alpha} \cdot \nabla_x g(\bar{x}) = 0.$$

$$\nabla_{\alpha} \mathcal{L}(\bar{x}, \bar{\alpha}) = g(\bar{x}) \leq 0.$$

$$\bar{\alpha} \cdot g(\bar{x}) = \sum_{i=1}^m \bar{\alpha}_i g_i(\bar{x}) = 0.$$

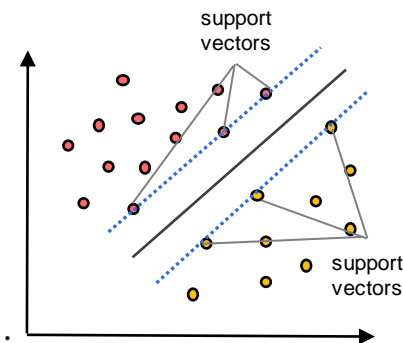


SVM Separable Case (KKT): $\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i (w \cdot x_i + b) - 1]$.

$$\nabla_w \mathcal{L} = w - \sum_{i=1}^m \alpha_i y_i x_i = 0, \quad \Rightarrow \quad w = \sum_{i=1}^m \alpha_i y_i x_i.$$

$$\nabla_b \mathcal{L} = - \sum_{i=1}^m \alpha_i y_i = 0, \quad \Rightarrow \quad \sum_{i=1}^m \alpha_i y_i = 0$$

$$\forall i, \alpha_i [y_i (w \cdot x_i + b) - 1] = 0, \quad \Rightarrow \quad \alpha_i = 0 \vee y_i (w \cdot x_i + b) = 1.$$





SVM

Dual Formulations

Optimization:

Definition: Dual Function

$$\forall \alpha \geq 0, F(\alpha) = \inf_{x \in \mathcal{X}} \mathcal{L}(x, \alpha).$$

Definition: Dual Optimization Problem \mathcal{P}^*

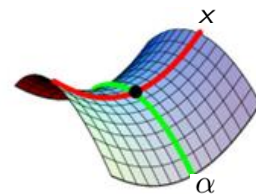
$$\begin{cases} \max_{\alpha \in \mathbb{R}^M} F(\alpha) \\ \text{subject to } \alpha_i \geq 0. \end{cases}$$

SVM Separable Case (Dual Form): $\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i (w \cdot x_i + b) - 1]$.

From KKT we have $w = \sum_{i=1}^m \alpha_i y_i x_i$ and $\sum_{i=1}^m \alpha_i y_i = 0$. This gives

Dual Function: $F(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$

Dual Optimization Problem \mathcal{P}^* :
$$\begin{cases} \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \\ \text{subject to } \alpha_i \geq 0 \wedge \sum_{i=1}^m \alpha_i y_i = 0, \quad \forall i \in \{1, 2, \dots, m\} \end{cases}$$



SVM Dual Formulation

SVN Non-Separable Case (Primal):

$$\begin{cases} \min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i^p \\ \text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \wedge \xi_i \geq 0, i \in [1, m] \end{cases}$$

Lagrangian:

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i}_{"f(x)"} - \underbrace{\sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i]}_{"\alpha \cdot g(x)"} - \sum_{i=1}^m \beta_i \xi_i$$

$x = (\mathbf{w}, b, \xi)$

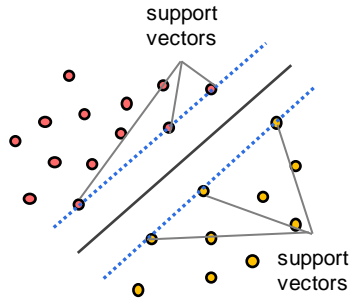
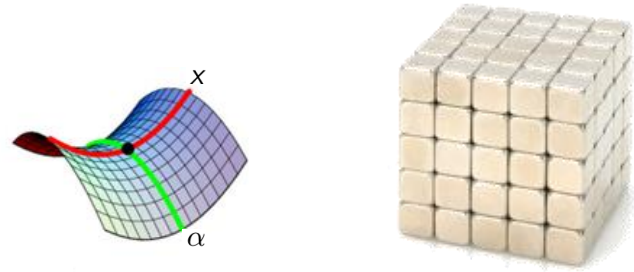
KKT Conditions:

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 &\implies \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \\ \nabla_b \mathcal{L} = - \sum_{i=1}^m \alpha_i y_i = 0 &\implies \sum_{i=1}^m \alpha_i y_i = 0 \\ \nabla_{\xi_i} \mathcal{L} = C - \alpha_i - \beta_i = 0 &\implies \alpha_i + \beta_i = C \rightarrow \alpha_i \leq C \beta_i \leq C \\ \forall i, \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0 &\implies \alpha_i = 0 \vee y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 - \xi_i \\ \forall i, \beta_i \xi_i = 0 &\implies \beta_i = 0 \vee \xi_i = 0. \end{aligned}$$

Dual Function F(α, β):

$$\mathcal{L} = \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \underbrace{\sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)}_{-\frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)} - \underbrace{\sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i}_0 \longrightarrow \mathcal{L} = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

same as separable case!



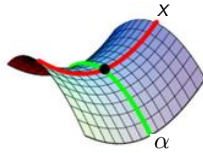
SVM Dual Formulation

SVM Non-Separable Case:

KKT Conditions:

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 &\implies \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \\ \nabla_b \mathcal{L} = - \sum_{i=1}^m \alpha_i y_i = 0 &\implies \sum_{i=1}^m \alpha_i y_i = 0 \\ \nabla_{\xi_i} \mathcal{L} = C - \alpha_i - \beta_i = 0 &\implies \alpha_i + \beta_i = C \\ \forall i, \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0 &\implies \alpha_i = 0 \vee y_i (\mathbf{w} \cdot \mathbf{x}_i + b) = 1 - \xi_i \\ \forall i, \beta_i \xi_i = 0 &\implies \beta_i = 0 \vee \xi_i = 0. \end{aligned}$$

Dual Function $F(\alpha, \beta)$:
$$\mathcal{L} = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$



SVM Non-Separable Case (Dual \mathcal{P}^*):

$$\left\{ \begin{aligned} &\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \underbrace{(\mathbf{x}_i \cdot \mathbf{x}_j)}_{\text{kernel here}} \\ &\text{subject to: } 0 \leq \alpha_i \leq C \wedge \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m]. \end{aligned} \right.$$

Inner-products of $(\mathbf{x}_i \cdot \mathbf{x}_j)$ only appear. **Kernel Method:** $\tilde{x} = \phi(x)$ holds for $\tilde{x}_i \cdot \tilde{x}_j = \phi(x_i) \cdot \phi(x_j) = k(x_i, x_j)$.

Dimension of dual problem is m . Primal problem has dimension N .

Regularization C in primal problem \mathcal{P} becomes constraint in dual problem \mathcal{P}^* .

Provides alternative ways to solve the optimization problem.



Generalization Error Bounds for Support Vector Machines



VC-Dimension: Hyperplanes

Example: Learning separating hyperplane in \mathbb{R}^N (related to SVM).
For data $\{(x_i, y_i)\}$ with $x_i \in \mathbb{R}^N$ and $y_i \in \{-1, 1\}$. Ideally, find \mathbf{w} , b so that $\text{sign}(\mathbf{w}^T \mathbf{x}_i + b) = y_i$.

Hypothesis class:

$$\mathcal{H} = \{h: h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) \text{ with } \mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}\}.$$

What is the $VCdim(\mathcal{H})$?

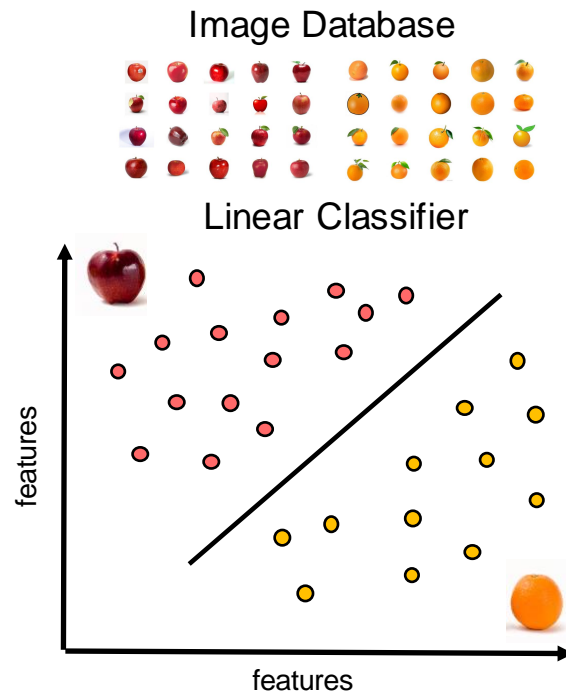
Claim: $VCdim(\mathcal{H}) = N + 1$

In separable case we have bound on generalization error ($\text{pr} > 1 - \delta$)

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{2(N+1) \log \frac{em}{N+1}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

We can do even better in bounding sampling complexity using special structure of SVM.

We want bounds independent of **feature dimension N** so can handle large N or even $N = \infty$.



SVM: Generalization Error

Definition: The **geometric margin** of a data point \mathbf{x} is

$$\rho(x) = \frac{y(\mathbf{w} \cdot \mathbf{x} + b)}{\|\mathbf{w}\|} \quad y \in \{-1, 1\} \text{ depending on side of hyperplane}$$

Definition: The **margin** of linear classifier $h(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ for data set $\mathcal{S} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)$ is

$$\rho = \min_{1 \leq i \leq m} \frac{y_i(\mathbf{w} \cdot \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

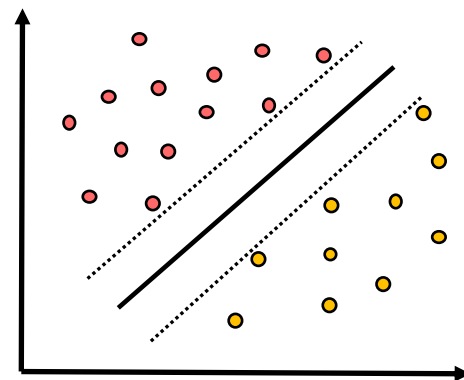
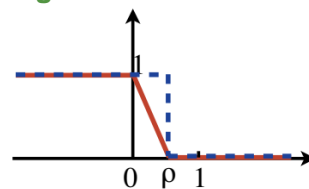
Definition: We define a **marginal loss function** using

$$\Phi_\rho(x) = \begin{cases} 0 & \text{if } \rho \leq x \\ 1 - x/\rho & \text{if } 0 \leq x \leq \rho \\ 1 & \text{if } x \leq 0. \end{cases}$$

Definition: We define **empirical marginal loss** as

$$\widehat{R}_\rho(h) = \frac{1}{m} \sum_{i=1}^m \Phi_\rho(y_i h(\mathbf{x}_i))$$

marginal loss function



SVM: Generalization Error



Theorem (Margin bound for binary classification): For any fixed $\rho > 0$ and $\delta > 0$, we have with probability $1 - \delta$ that the generalization error for marginal loss function is

$$R(h) \leq \widehat{R}_\rho(h) + \frac{2}{\rho} \mathfrak{R}_m(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

$$R(h) \leq \widehat{R}_\rho(h) + \frac{2}{\rho} \widehat{\mathfrak{R}}_S(H) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$$

Key idea: obtain bounds using the Rademacher Complexity of

$\mathcal{H} = \{h: h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) \text{ with } \mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}\}$.

Notational convention: Suppress the b term by using

$\tilde{\mathbf{x}} = \begin{bmatrix} x \\ 1 \end{bmatrix}, \tilde{\mathbf{w}} = \begin{bmatrix} w \\ b \end{bmatrix}, \mathcal{H} = \{h: h(\tilde{\mathbf{x}}) = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}) \text{ with } \tilde{\mathbf{w}} \in \mathbb{R}^{N+1}\}$.

SVM: Generalization Error

Theorem: (Radamacher Complexity of Constrained Hyperplanes for Bounded Data \mathcal{S})

Let $\mathcal{S} \subseteq \{x : \|x\| \leq r\}$ be a sample of size m and let $\mathcal{H} = \{h \mid h(x) = \text{sign}(w \cdot x) \mid \|w\| \leq \Lambda\}$.

The Rademacher complexity satisfies

$$\hat{R}_S(\mathcal{H}) \leq \sqrt{\frac{r^2 \Lambda^2}{m}}.$$

Proof:

$$\begin{aligned} \hat{R}_S(\mathcal{H}) &= \frac{1}{m} E_\sigma \left[\sup_{\|w\| \leq \Lambda} \sum_{i=1}^m \sigma_i w \cdot x_i \right] = \frac{1}{m} E_\sigma \left[\sup_{\|w\| \leq \Lambda} w \cdot \sum_{i=1}^m \sigma_i x_i \right] \\ &\leq \frac{\Lambda}{m} E_\sigma \left[\left\| \sum_{i=1}^m \sigma_i x_i \right\| \right] \leq \frac{\Lambda}{m} E_\sigma \left[\left\| \sum_{i=1}^m \sigma_i x_i \right\|^2 \right]^{1/2} \\ &\leq \frac{\Lambda}{m} E_\sigma \left[\sum_{i,j=1}^m \sigma_i \sigma_j x_i \cdot x_j \right]^{1/2} = \frac{\Lambda}{m} \left[\sum_{i=1}^m \|x_i\|^2 \right]^{1/2} \leq \frac{\Lambda \sqrt{mr^2}}{m} = \sqrt{\frac{r^2 \Lambda^2}{m}} \quad \blacksquare. \end{aligned}$$

Cauchy-Swartz Lemma:

$$a \cdot b \leq \|a\| \|b\|.$$

Jensen Inequality:

$$\phi(E[X]) \leq E[\phi(X)]$$

$$(E[X])^2 \leq E[X^2]$$

Radamacher Random Variables:

$$E[\sigma_i \sigma_j] = E[\sigma_i] E[\sigma_j] = 0.$$

Talagrand's Lemma:

Let $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ be an ℓ -Lipschitz function $\|\Phi(x) - \Phi(y)\| \leq \ell|x - y|$, then

$$\hat{R}_S(\Phi \circ \mathcal{H}) \leq \ell \hat{R}_S(\mathcal{H}).$$

SVM: Generalization Error



Theorem (Margin bound for binary classification): For any $\rho > 0$ and $\delta > 0$, we have with probability $1 - \delta$ that the generalization error for marginal loss function is

$$R(h) \leq \widehat{R}_\rho(h) + \frac{2}{\rho} \mathfrak{R}_m(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

$$R(h) \leq \widehat{R}_\rho(h) + \frac{2}{\rho} \widehat{\mathfrak{R}}_S(H) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$$

Key idea: obtain bounds using the Rademacher Complexity of

$\mathcal{H} = \{h: h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) \text{ with } \mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}\}$.

Notational convention: Suppress the b term by using

$\tilde{\mathbf{x}} = \begin{bmatrix} x \\ 1 \end{bmatrix}, \tilde{\mathbf{w}} = \begin{bmatrix} w \\ b \end{bmatrix}, \mathcal{H} = \{h: h(\tilde{\mathbf{x}}) = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}) \text{ with } \tilde{\mathbf{w}} \in \mathbb{R}^{N+1}\}$.

SVM: Generalization Error

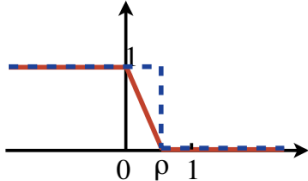


Theorem (Margin bound for binary classification): For any fixed $\rho > 0$ and $\delta > 0$, we have with probability $1 - \delta$ that the generalization error for marginal loss function is

$$R(h) \leq \hat{R}_\rho(h) + 2\sqrt{\frac{r^2\Lambda^2/\rho^2}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

$$\mathcal{H} = \{h \mid h(x) = \text{sign}(w \cdot x) \mid \|w\| \leq \Lambda\}$$

marginal loss function



$$h(x_i) = \text{sign}(w \cdot x_i)$$

$$\hat{R}_\rho(h) = \frac{1}{m} \sum_{i=1}^m \Phi_\rho(y_i h(x_i))$$

We have the **marginal loss** is bounded by **hinge loss**

$$\Phi_1(x) \leq \max(1 - x, 0)$$

Theorem: For any fixed $\delta > 0$, we have with probability $1 - \delta$ that the generalization error for marginal loss function is

$$R(h) \leq \underbrace{\frac{1}{m} \sum_{i=1}^m \xi_i}_{\text{empirical risk}} + \underbrace{2\sqrt{\frac{r^2\Lambda^2}{m}}}_{\text{class complexity (regularization)}} + \underbrace{\sqrt{\frac{\log \frac{1}{\delta}}{2m}}}_{\text{sampling confidence}}$$

$$\begin{aligned} \Phi_1(y_i h(x_i)) &\leq \max\{1 - y_i h(x_i), 0\} \\ &= \max\{1 - y_i \text{sign}(w \cdot x_i), 0\} \\ &= \xi_i. \end{aligned}$$

$$\hat{R}_\rho(h) \leq \frac{1}{m} \sum_{i=1}^m \xi_i.$$

SVM Objective:
 $\sum_{i=1}^m \xi_i + \frac{1}{2} \|w\|^2.$

Key Result: The SVM objective function \rightarrow makes small the RHS bound!

Trade-off: make **slack variables small** while making **margin $\rho = 1/\|w\|$ large**. Allows for $N = \infty$.

Regularization: Make $\|w\|^2 \leq \Lambda^2$ **small** \rightarrow serves as **regularization** term (controlled by $\Lambda = \Lambda(C^{-1})$).

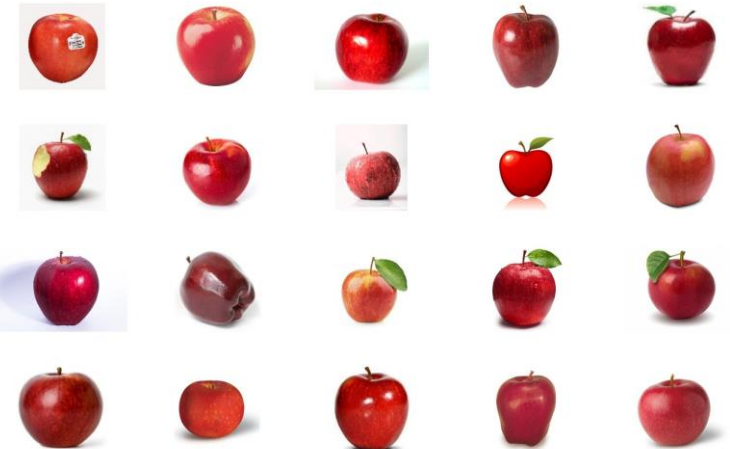


Example of SVM

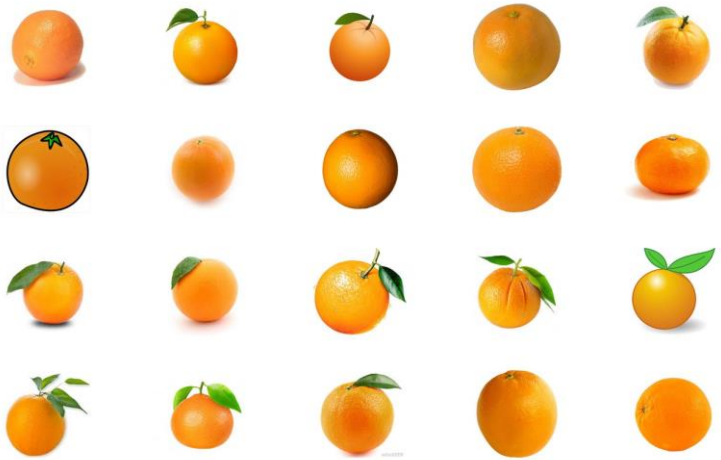
Classifying Apples & Oranges

Support Vector Machines (Apples and Oranges Training Data)

Apple Image Set



Orange Image Set



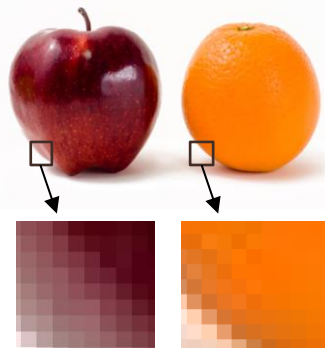
Average Apple



Average Orange



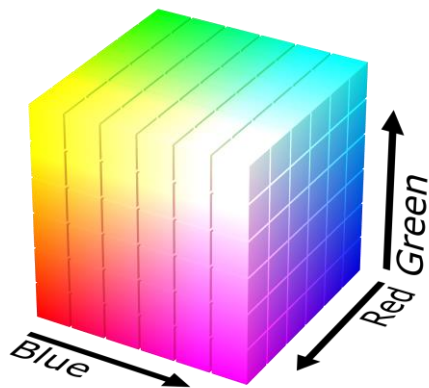
Supervised Learning (Apples and Oranges)



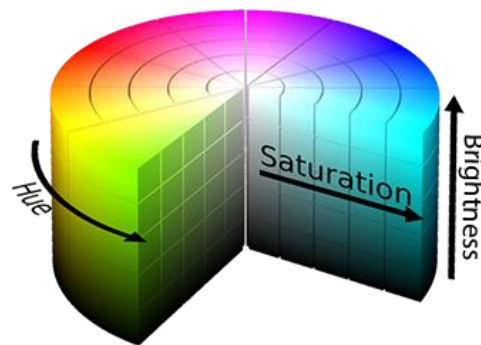
What features to use to distinguish apples and oranges?

- Natural to use the colors of the objects in the images.
- However, many different color spaces can be used (RGB, HLS, ...)
- Does the choice matter?

Red, Green, Blue: [Pixel (RGB)]



Hue, Luminance, Saturation: [Pixel (HLS)]

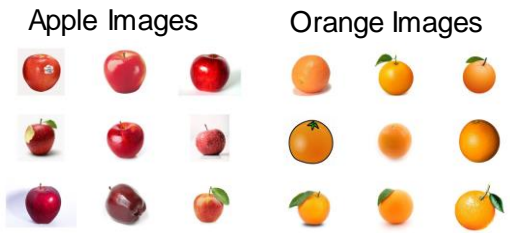


SVM Performance

Linear: $K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle$

RBF: $K(\mathbf{x}, \mathbf{y}) = \exp[-\gamma\|\mathbf{x} - \mathbf{y}\|^2]$

Polynomial: $K(\mathbf{x}, \mathbf{y}) = (\gamma\langle \mathbf{x}, \mathbf{y} \rangle + r)^d$



Item		Features	
Feature		Value	
Roundness		0.8	
Sweetness		0.9	
Redness		0.1	
Greenness		0.3	

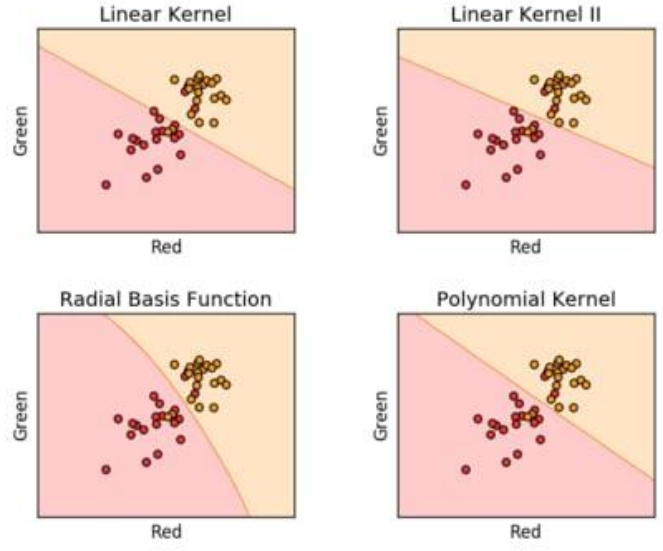
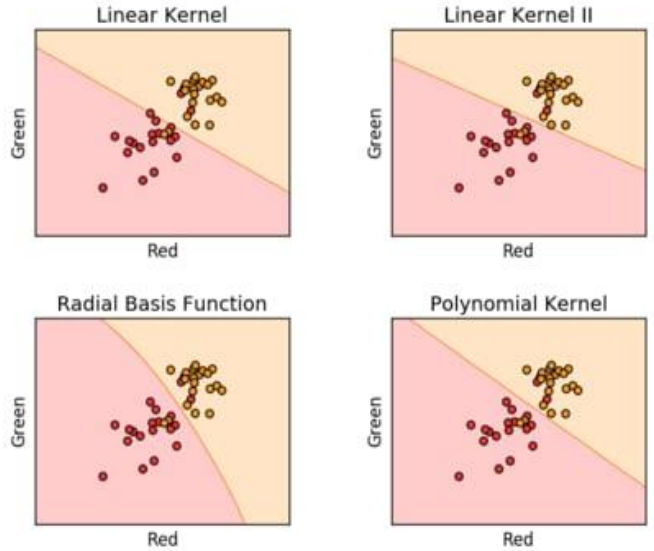
SVM Results:

RGB Features

HLS Features

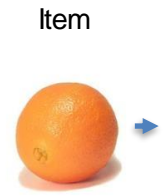
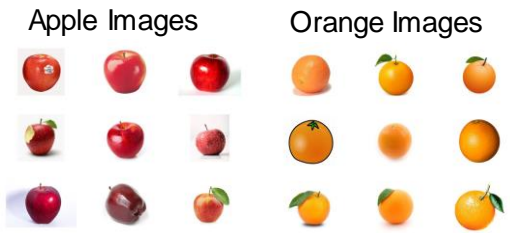
C = 1e+0.8

C = 1e+0.8



SVM Performance

Importance of features used?
 Importance of regularization C?
 How does training set generalize?

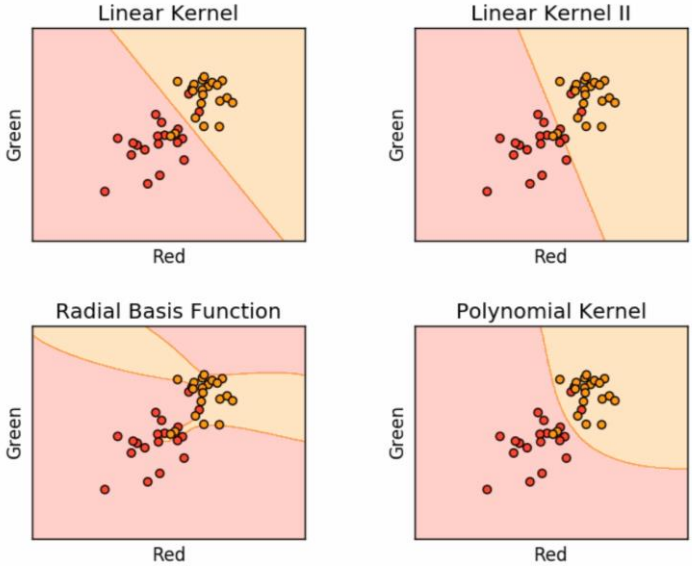


Features	
Feature	Value
Roundness	0.8
Sweetness	0.9
Redness	0.1
Greenness	0.3

SVM Results:

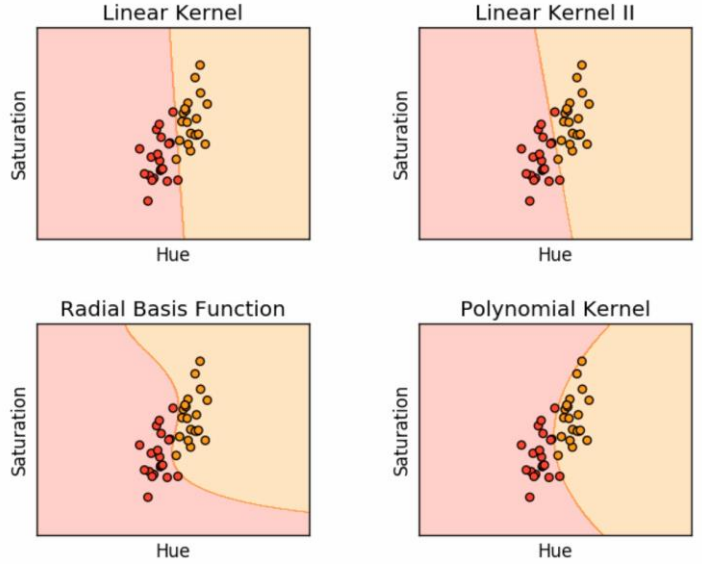
RGB Features


C = 1e+6.9



HLS Features


C = 1e+6.9





SVM Example

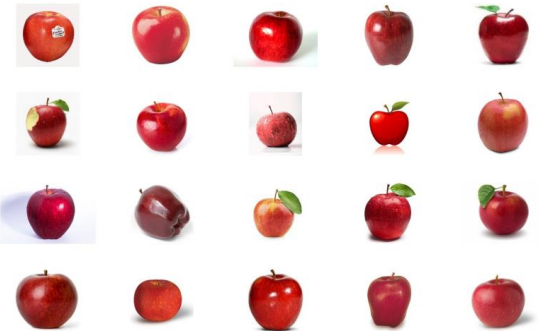
Apples vs Oranges vs Blueberries



Support Vector Machines (Apples and Oranges and Blueberries)

How might we train on more than two data sets?
Three data sets: Apples, Oranges, and Blueberries.

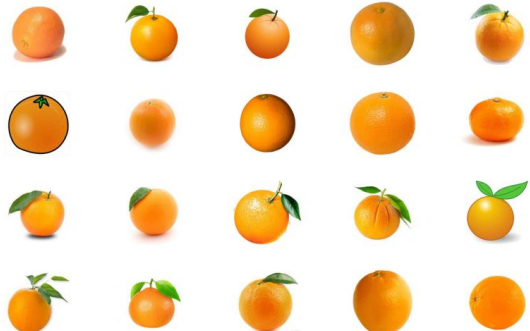
Apple Image Set



Average Apple



Orange Image Set



Average Orange



Blueberry Image Set



Average Blueberry





SVM Multi-class Classification



Support Vector Machines (Multi-Class Case)

Many classification involve multiple classes.

Problem: From data x learn for k classes C_1, C_2, \dots, C_k a classifying function $f(x) \rightarrow y \in \{C_1, C_2, \dots, C_k\}$.

Binary classification: $f(x) \rightarrow y \in \{-1, 1\}$ which corresponds to $k = 2$ classes $C_1 = -1, C_2 = 1$.

Multi-class classification: $f(x) \rightarrow y \in \{1, 2, \dots, k\}$ corresponds to k classes $C_1 = 1, C_2 = 2, \dots, C_k = k$.

How can we extend linear classifier methods to handle multiple classes?

Two common approaches:

One vs All (OvA): Reduce to a collection of k binary classification problems to determine one category labeled +1 vs rest of the data labeled -1. Pick classification with the greatest margin.

One vs One (OvO): Reduce to a collection of $\binom{k}{2} = k(k-1)/2$ binary classification problems to determine one category labeled +1 vs one other category labeled -1. Consider each classifier as a voter and pick class with the most number of votes.

Above heuristics do not always work well in practice. Alternatives: optimization formulations (more expensive).

Support Vector Machines (Multi-Class Classification)

Optimization of Maximum Margin (OMM): Classes $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$, with $C_\ell = \ell$, data (x_i, y_i) with $y_i \in \mathcal{C}$.

$$\begin{aligned} \min_{\mathbf{w}, \xi} & \frac{1}{2} \sum_{l=1}^k \|\mathbf{w}_l\|^2 + C \sum_{i=1}^m \xi_i \\ \text{subject: } & \forall i \in [1, m], \forall l \in \mathcal{C} \setminus \{y_i\} \\ & \mathbf{w}_{y_i} \cdot \Phi(x_i) \geq \mathbf{w}_l \cdot \Phi(x_l) + 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

Classifier obtained: $h(x) = \operatorname{argmax}_{l \in \mathcal{C}} \mathbf{w}_l \cdot \Phi(x)$, where $\Phi(x)$ is transformation of the data.

Dual Optimization Problem (Keneralization): $K(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$

$$\max_{\alpha \in \mathbb{R}^{m \times k}} \sum_{i=1}^m \alpha_i \cdot \mathbf{e}_{y_i} - \frac{1}{2} \sum_{i=1}^m (\alpha_i \cdot \alpha_j) K(x_i, x_j)$$

subject to: $\forall i \in [1, m], (0 \leq \alpha_{iy_i} \leq C) \wedge (\forall j \neq y_i, \alpha_{ij} \leq 0) \wedge (\alpha_i \cdot \mathbf{1} = 0)$

Generalization Bounds: $R(h) = \mathbb{E}_{x \sim D} [1_{h(x) \neq f(x)}]$ $\mathcal{H} = \{h(x) = \operatorname{argmax}_{l \in \mathcal{C}} \mathbf{w}_l \cdot \Phi(x) \mid \mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k)^T, \sum_{l=1}^k \|\mathbf{w}_l\|^2 \leq \Lambda^2\}$

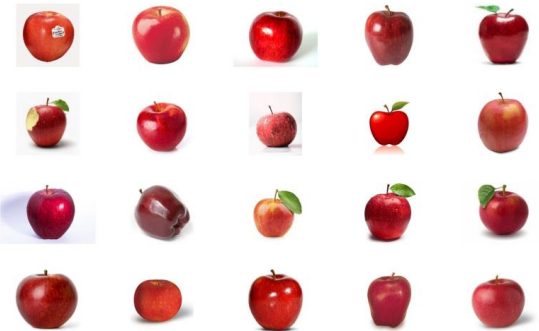
$$R(h) \leq \frac{1}{m} \sum_{i=1}^m \xi_i + 2k^2 \sqrt{\frac{r^2 \Lambda^2}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (\text{for any } \delta > 0 \text{ holds with probability } 1 - \delta)$$

where $\xi_i = \max(1 - [\mathbf{w}_{y_i} \cdot \Phi(x_i) - \max_{y' \neq y_i} \mathbf{w}_{y'} \cdot \Phi(x_i)], 0)$ for all $i \in [1, m]$

Support Vector Machines (Apples and Oranges and Blueberries)

How might we train on more than two data sets?
Three data sets: Apples, Oranges, and Blueberries.

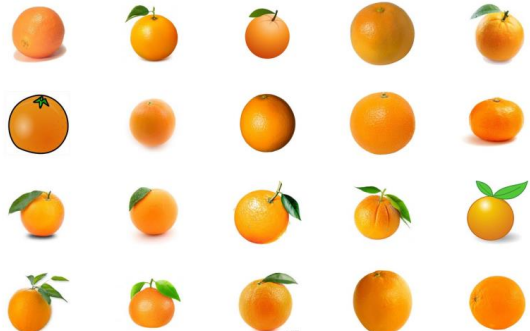
Apple Image Set



Average Apple



Orange Image Set



Average Orange



Blueberry Image Set



Average Blueberry



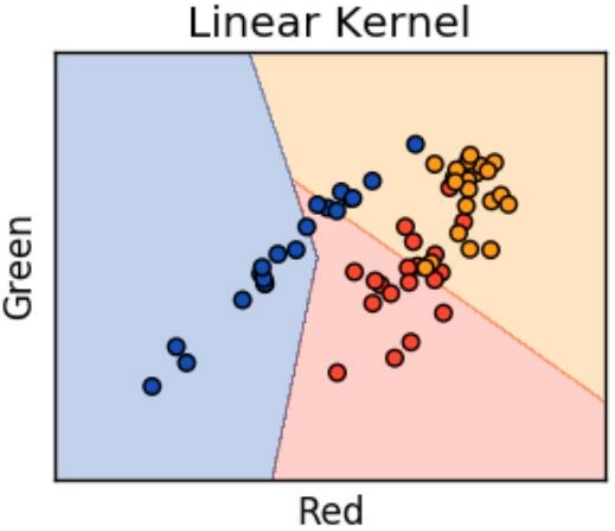
Supervised Learning (SVM Results: Apples, Oranges and Blueberries)



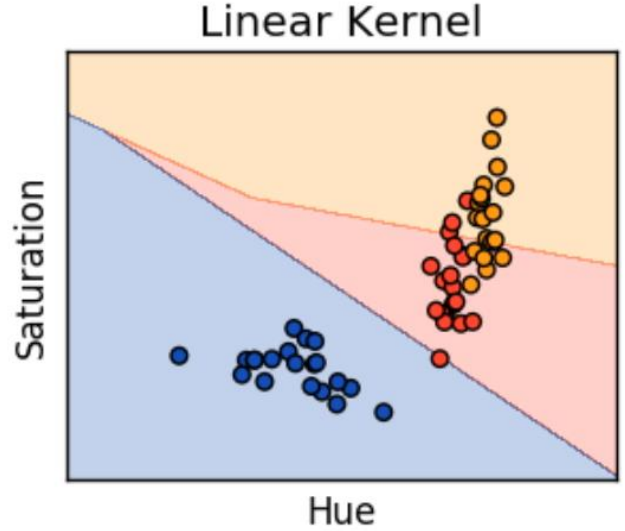
How does SVM distinguish between three data sets?
Importance of features used.
Importance of regularization C.
How does training set generalize?

SVM Results:

RGB Features



HLS Features



C = 1e-4.0



Supervised Learning (SVM Results: Apples, Oranges and Blueberries)



Linear: $K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle$

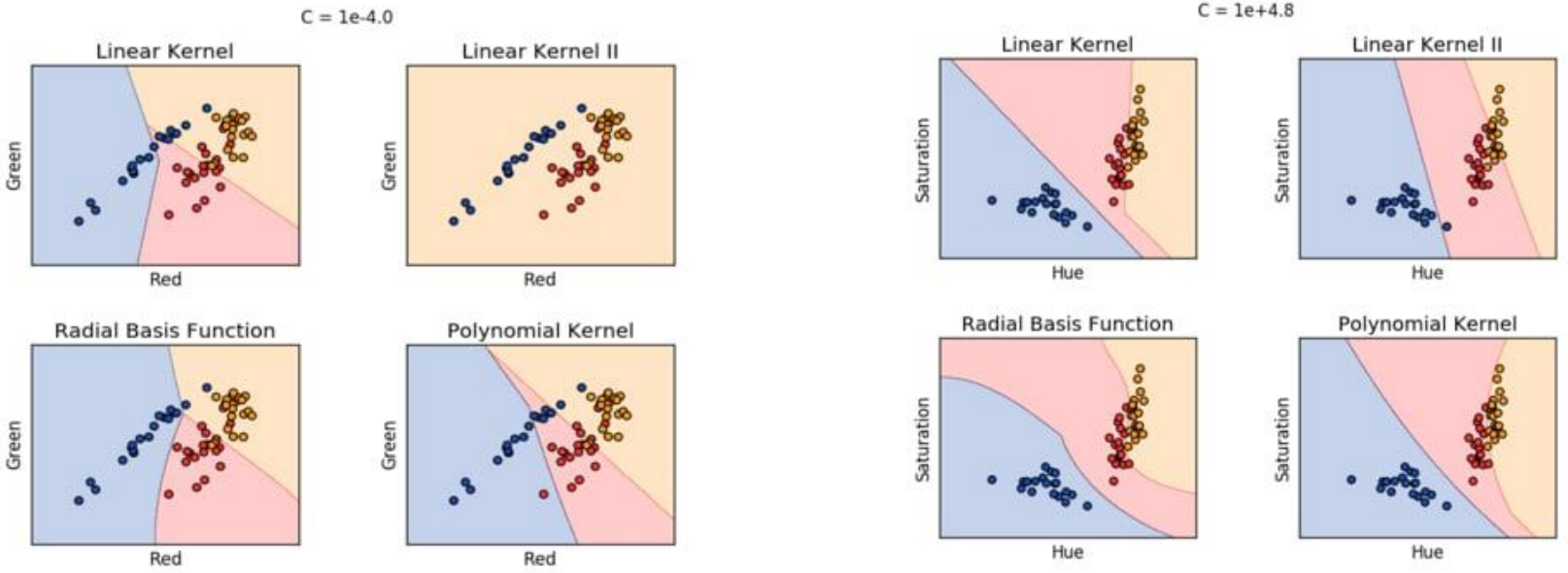
RBF: $K(\mathbf{x}, \mathbf{y}) = \exp[-\gamma \|\mathbf{x} - \mathbf{y}\|^2]$

Polynomial: $K(\mathbf{x}, \mathbf{y}) = (\gamma \langle \mathbf{x}, \mathbf{y} \rangle + r)^d$

SVM Results:

RGB Features

HLS Features



Supervised Learning (SVM Results: Apples, Oranges and Blueberries)

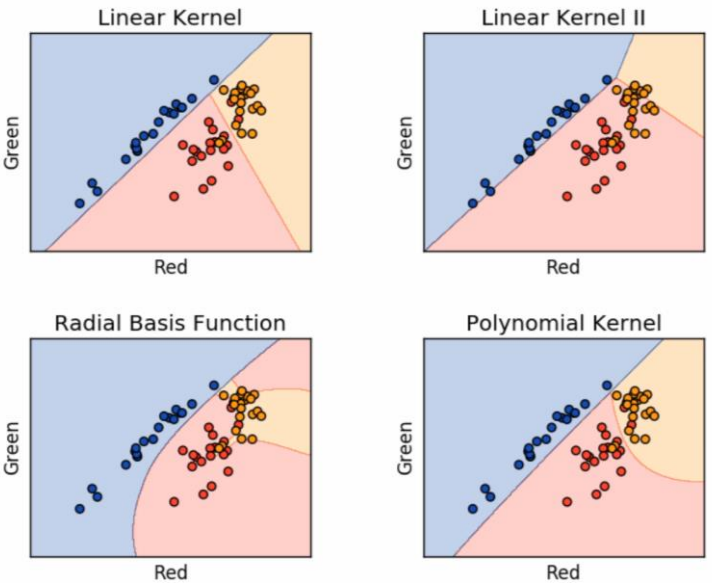


How does SVM distinguish between three data sets?
Importance of features used.
Importance of regularization C.
How does training set generalize?

SVM Results:

RGB Features

C = 1e+6.9



HLS Features

C = 1e+6.9

