# Data Mining & Machine Learning

CS37300
Purdue University

Sep 22, 2023

# Today's topics

- Logistic Regression

- Kernels and Kernel-SVM

# Modeling the Noise: Conditional Probability

- Logistic Model:

$$P_{w,b}(Y = y | X = x) = \frac{1}{1 + e^{-y(w^\top x + b)}}$$

# Training: Logistic Regression

- How do we find a good $\hat{w}$ ?

- **Maximum conditional likelihood** estimation

- For a data set S = {(x$_1$,y$_1$),...,(x$_n$,y$_n$)}

- Define the conditional likelihood

Models the labels y$_i$ as being conditionally independent given x$_i$

$$L_{Y|X}(w, b; S) = \prod_{i=1}^{n} P_{w,b}(Y = y_i | X = x_i)$$

- The maximum conditional likelihood estimator is

$$(\hat{w}, \hat{b}) = \operatorname*{argmax}_{w,b} L_{Y|X}(w, b; S)$$

- Equivalently (because it's easier to work with), conditional **log-likelihood**:

$$l_{Y|X}(w, b; S) = \ln(L_{Y|X}(w, b; S))$$

- and then

$$(\hat{w}, \hat{b}) = \operatorname*{argmax}_{w,b} l_{Y|X}(w, b; S)$$

# Training: Logistic Regression

- Explicitly:

$$l_{Y|X}(w, b; S) = \ln\left(\prod_{i=1}^{n} P_w(Y = y_i | X = x_i)\right)$$

$$= \sum_{i=1}^{n} \ln(P_w(Y = y_i | X = x_i))$$

$$= \sum_{i=1}^{n} \ln\left(\frac{1}{1 + e^{-y_i(w^\top x_i + b)}}\right)$$

$$= -\sum_{i=1}^{n} \ln\left(1 + e^{-y_i(w^\top x_i + b)}\right)$$

- Unfortunately, there is no simple expression for the $\widehat{w}$ that maximizes this

- But $-\sum_{i=1}^{n} \ln\left(1 + e^{-y_i(w^\top x_i + b)}\right)$ is a concave function, which can be maximized using iterative numerical methods: e.g., gradient ascent or Newton's method.

# Training: Logistic Regression

- Explicitly:      denote   $\sigma(x) = \dfrac{1}{1 + e^{-x}}$

- Simple derivative:     $\dfrac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$

$$\nabla l_{Y|X}(w, b; S) = \nabla \sum_{i=1}^{n} \ln\big(\sigma\big(y_i(w^\top x_i + b)\big)\big)$$

$$= \sum_{i=1}^{n} \frac{1}{\sigma(y_i(w^\top x_i + b))} \sigma\big(y_i(w^\top x_i + b)\big) \big(1 - \sigma\big(y_i(w^\top x_i + b)\big)\big) y_i[x_i, 1]^\top$$

$$= \sum_{i=1}^{n} \big(1 - \sigma\big(y_i(w^\top x_i + b)\big)\big) y_i[x_i, 1]^\top$$
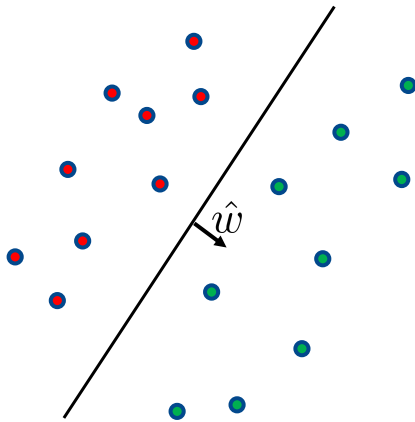
gradient ascent: iterate    $(w, b) \leftarrow (w, b) + \epsilon \nabla l_{Y|X}(w, b; S)$
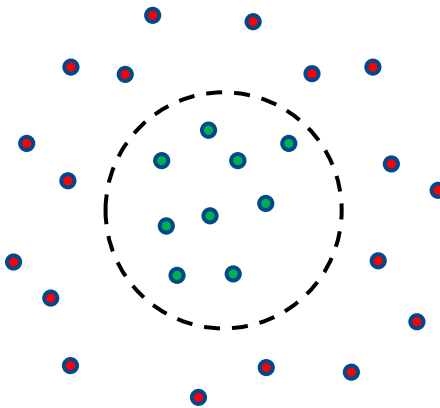
# Logistic Regression vs SVM

- For a given data set, which one should we use?

- If you think the data are (nearly) linearly separable, and with large margin, makes sense to use (Soft)-SVM

- If you think the cause of non-separability truly is label noise, makes sense to use logistic regression

- One nice thing about Logistic Regression is that it provides an estimated probability P(y|x) for its predictions, rather than just a -1,1 prediction.

- Doesn't hurt to try them both and find out which is better on a validation set!

# Nonlinear Decision Boundaries

- Linear separators are simple



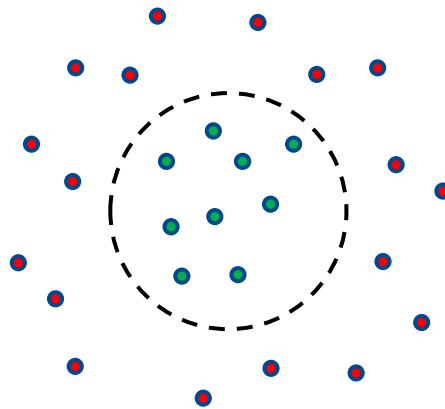- But there are many scenarios that are separable, only non-linearly



- We want to learn these with optimization based ML (SVM)

# Learning Nonlinear Decision Boundaries

- We've already seen some nonlinear separators:

    o Decision Trees, KNN

- Question: How can we learn this example using SVM?

Example:

# Learning Nonlinear Decision Boundaries

- We've already seen some nonlinear separators:

  o Decision Trees, KNN

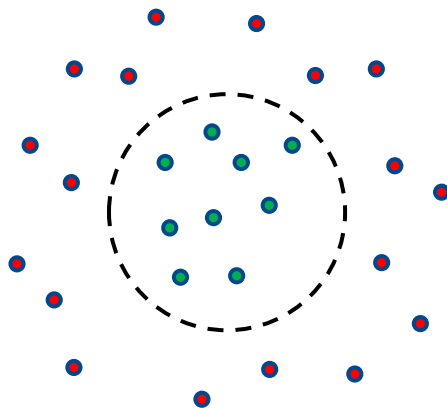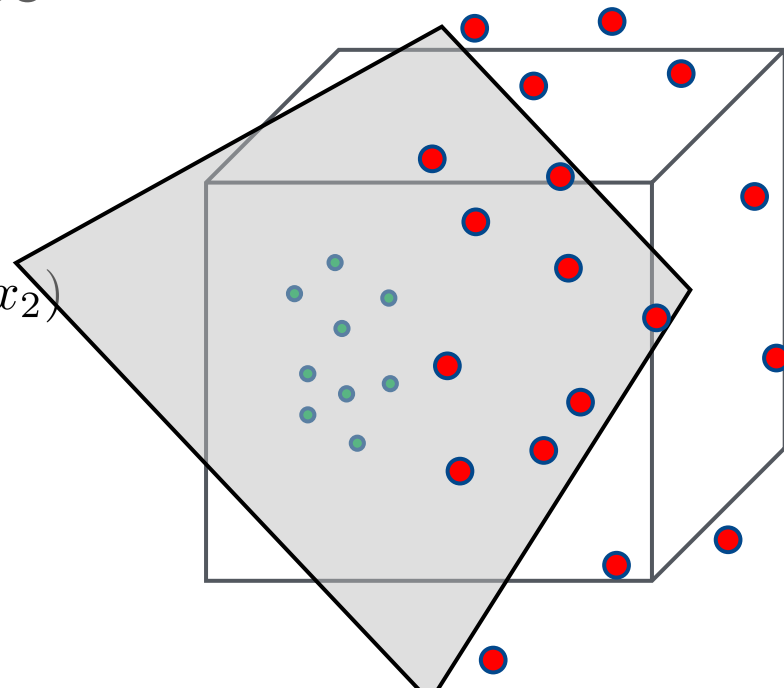- Question: How can we learn this example using SVM?

  o Map the examples into a higher-dimensional space and learn a linear separator in that space

Example:

Re-represent each $\underline{x} = (x_1, x_2)$ as 3-dimensional vector $\phi(\underline{x}) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$

# Representing Nonlinear Decision Boundaries

- General strategy: consider a function: $\phi(x) \in \mathbb{R}^N$

  - N may be large, or even infinite

- Then learn a linear separator in the $\phi$ space:

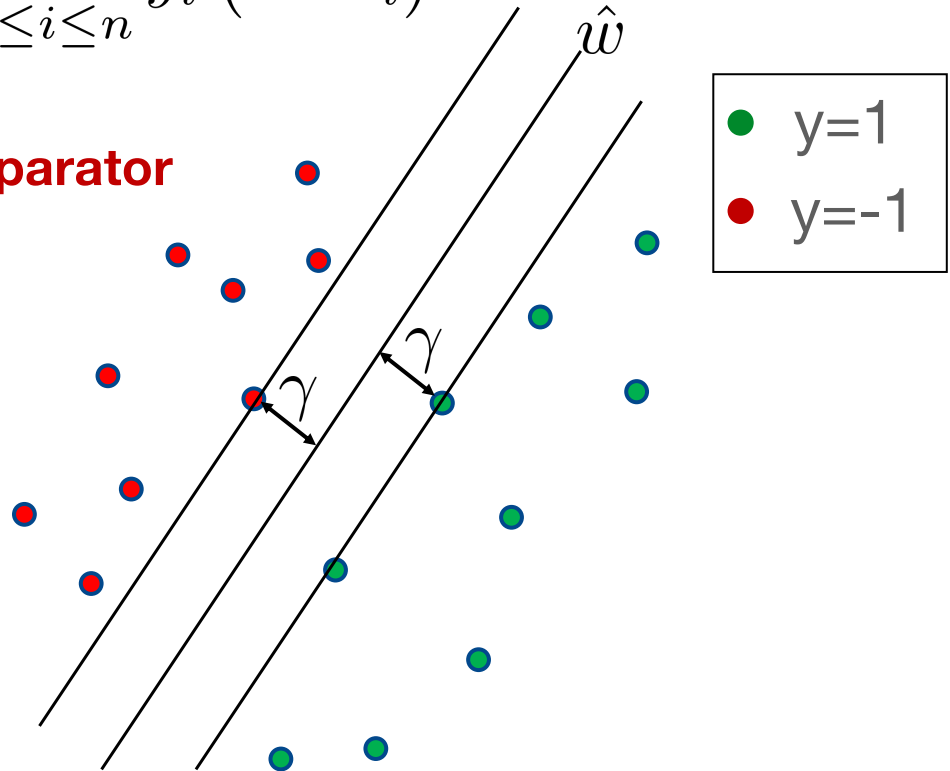$$\hat{h}(x) = \text{sign}\left(w^\top \phi(x)\right), \quad \text{where } w \in \mathbb{R}^N$$

# Refresher: Support Vector Machine (SVM)

$$\hat{w} = \underset{w:\|w\|=1}{\operatorname{argmax}} \ \min_{1 \leq i \leq n} y_i \left( w^\top x_i \right)$$

- SVM: pick the **maximum margin separator**

geometric margin:

$$\gamma = \underset{w:\|w\|=1}{\max} \ \min_{1 \leq i \leq n} y_i \left( w^\top x_i \right)$$

$\hat{w}$

- $y=1$
- $y=-1$

$\gamma$

$\gamma$

| | |
|---|---|
| Minimize | $\|w\|^2$ |
| subject to | $y_i \left( w^\top x_i \right) \geq 1, \ \forall i : \ 1 \leq i \leq n$ |

# Representing Nonlinear Decision Boundaries

- General strategy: consider a function: $\phi(x) \in \mathbb{R}^N$

  - N may be large, or even infinite

- Then learn a linear separator in the $\phi$ space:

$$\hat{h}(x) = \operatorname{sign}\left(w^\top \phi(x)\right), \quad \text{where } w \in \mathbb{R}^N$$

$$\hat{w} = \underset{w \in \mathbb{R}^N : \|w\| = 1}{\operatorname{argmax}} \ \min_{1 \le i \le n} y_i \left(w^\top \phi(x_i)\right)$$

$$\boxed{\begin{array}{l} \text{Minimize} \quad \|w\|^2 \\ \text{subject to} \quad y_i \left(w^\top \phi(x_i)\right) \ge 1, \ \forall i : \ 1 \le i \le n \end{array}}$$

- Question: Why might this be a problem?

# Representing Nonlinear Decision Boundaries

- General strategy: consider a function: $\phi(x) \in \mathbb{R}^N$

  - N may be large, or even infinite

- Then learn a linear separator in the $\phi$ space:

$$\hat{h}(x) = \text{sign}\left(w^\top \phi(x)\right), \quad \text{where } w \in \mathbb{R}^N$$

- Two concerns:

  1. Computation and memory required to represent $\phi(x), w$

  2. Dimension of linear separators on $\mathbb{R}^N$ is N.
     $\rightarrow$ Concerns about overfitting

  Today's lecture is mostly about the first concern.
  Solution: Kernels

# The Dual Form of the SVM Optimization Problem

$$\hat{w} = \underset{w : \|w\| = 1}{\operatorname{argmax}} \; \min_{1 \leq i \leq n} y_i \left( w^\top x_i \right)$$

- Recall: The SVM **classifier** is the unique solution to a **quadratic program**:

$$\text{Minimize} \quad \|w\|^2$$

$$\text{subject to} \quad y_i \left( w^\top x_i \right) \geq 1, \; \forall i : \; 1 \leq i \leq n$$

- We can re-express this in Lagrangian dual form:

$$\hat{w} = \sum_{i=1}^{n} \alpha_i y_i x_i$$

where $\alpha_1, \ldots, \alpha_n$ are solutions to:

$$\text{Maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j x_i^\top x_j$$

$$\text{subject to} \quad \alpha_i \geq 0, \; \forall i : \; 1 \leq i \leq n$$

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

# The Dual Form of the SVM Optimization Problem

$$\hat{w} = \underset{w:\|w\|=1}{\mathrm{argmax}} \ \min_{1 \le i \le n} y_i \left( w^\top x_i \right)$$

- Recall: The SVM **classifier** is the unique solution to a **quadratic program**:

$$\begin{aligned} \text{Minimize} \quad & \|w\|^2 \\ \text{subject to} \quad & y_i \left( w^\top x_i \right) \ge 1, \ \forall i : \ 1 \le i \le n \end{aligned}$$

- We can re-express this in Lagrangian dual form:

$$\hat{w} = \sum_{i=1}^{n} \alpha_i y_i x_i$$

where $\alpha_1, \ldots, \alpha_n$ are solutions to:

$$\begin{aligned} \text{Maximize} \quad & \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j x_i^\top x_j \\ \text{subject to} \quad & \alpha_i \ge 0, \ \forall i : \ 1 \le i \le n \\ & \sum_{i=1}^{n} \alpha_i y_i = 0 \end{aligned}$$

Side note: the training points $(x_i, y_i)$ with non-zero $\alpha_i$ values are called the **support vectors**.

# SVM with high-dim mapping

- We could use a mapping $\phi$ to get a non-linear separator:

$$\hat{w} = \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{\phi(x_i)}$$

where $\alpha_1, \ldots, \alpha_n$ are solutions to:

$$\text{Maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j \boldsymbol{\phi(x_i)}^{\top} \boldsymbol{\phi(x_j)}$$

$$\text{subject to} \quad \alpha_i \geq 0, \ \forall i : \ 1 \leq i \leq n$$

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

- Notice training only uses $\phi$ when computing inner product $\phi(x_i)^{\top}\phi(x_j)$

- Instead of starting by defining $\phi$, we could start by defining a function $K(x_i, x_j)$ that computes an inner product, without computing $\phi$: $\quad K(x_i, x_j) = \phi(x_i)^{\top}\phi(x_j)$

- Kernel: way to compute inner products without computing $\phi(\text{x})$!

# SVM with high-dim mapping

- We could use a mapping $\phi$ to get a non-linear separator:

$$\hat{w} = \sum_{i=1}^{n} \alpha_i y_i \phi(\boldsymbol{x_i})$$

where $\alpha_1, \ldots, \alpha_n$ are solutions to:

$$\text{Maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j \boldsymbol{\phi(x_i)}^{\top} \boldsymbol{\phi(x_j)}$$

$$\text{subject to} \quad \alpha_i \geq 0, \ \forall i: \ 1 \leq i \leq n$$

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

- Notice training only uses $\phi$ when computing inner product

$$\phi(x_i)^{\top} \phi(x_j)$$

- And the classifier also uses an inner product (no need to compute w)

$$h_{\hat{w}}(x) = \text{sign}\left(\hat{w}^{\top} \phi(x)\right) = \text{sign}\left(\left(\sum_{i=1}^{n} \alpha_i y_i \phi(x_i)\right)^{\top} \phi(x)\right) = \text{sign}\left(\sum_{i=1}^{n} \alpha_i y_i \boldsymbol{\phi(x_i)}^{\top} \boldsymbol{\phi(x)}\right)$$

- Replace all of these with

$$K(x_i, x) = \phi(x_i)^{\top} \phi(x)$$

# Kernel SVM

- Training time:

Solve for $\alpha_1, \ldots, \alpha_n$:

$$\text{Maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j \boldsymbol{K(x_i, x_j)}$$

$$\text{subject to} \quad \alpha_i \geq 0, \ \forall i : \ 1 \leq i \leq n$$

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

- Test time:

- Classify a new point x with

$$\hat{h}(x) = \text{sign}\left( \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{K(x_i, x)} \right)$$

# Kernels

Example:

Quadratic kernel: $\quad K(\underline{u}, \underline{v}) = \left(\underline{u}^\top \underline{v} + 1\right)^2$

For $\underline{x} \in \mathbb{R}^d$, implicitly computes an inner product in $\frac{1}{2}d(d-1) + 2d + 1$ dim

e.g., For $\underline{x} \in \mathbb{R}^2$, implicitly defined $\phi$ :

$$\phi(\underline{x}) = \left[1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2\right]^\top$$

Check: $\phi(\underline{u})^\top \phi(\underline{v})$

$$= \left[1, \sqrt{2}u_1, \sqrt{2}u_2, u_1^2, u_2^2, \sqrt{2}u_1 u_2\right] \left[1, \sqrt{2}v_1, \sqrt{2}v_2, v_1^2, v_2^2, \sqrt{2}v_1 v_2\right]^\top$$

$$= 1 + 2u_1 v_1 + 2u_2 v_2 + u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 u_2 v_1 v_2$$

$$= (1 + u_1 v_1 + u_2 v_2)^2 = \left(1 + \underline{u}^\top \underline{v}\right)^2 = K(\underline{u}, \underline{v})$$

The point is that we can compute this higher-dim inner product just by evaluating K(u,v): no need to compute $\phi$

# Kernels

- More Examples:

  - Polynomial kernel:  $K(\underline{u}, \underline{v}) = \left(\underline{u}^\top \underline{v} + 1\right)^p$

    - Implicitly computes an inner product in $\sim d^p$ dimensions

  - Gaussian kernel:

    $$K(\underline{u}, \underline{v}) = e^{-\frac{\|\underline{u} - \underline{v}\|^2}{2\sigma^2}}$$

    - Implicitly computes an inner product in **infinite** dimensions

    - Remark: This is the most popular kernel in practice

# What are Legal Kernels?

- A kernel K($\cdot$,$\cdot$) is a legal definition of an inner product

- Technically, this is called a **Mercer kernel**

  - A kernel should be a symmetric function: K(u,v)=K(v,u)

  - A kernel should also be positive semi-definite: namely,
    - For any set of data points $x_1,\ldots,x_n$
    - And for any values $a_1, \ldots, a_n \in \mathbb{R}$
    $$\sum_{i=1}^{n}\sum_{j=1}^{n} a_i a_j K(x_i, x_j) \geq 0$$

- Mercer's Theorem: For any K satisfying the above, there exists some function $\phi$ such that
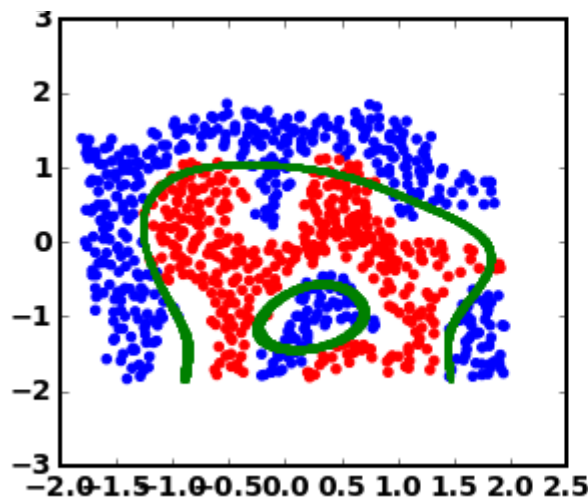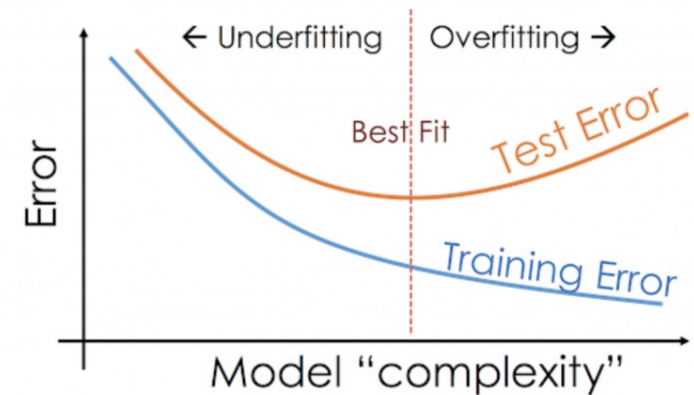  $$K(u, v) = \phi(u)^\top \phi(v)$$

- In other words: K is a valid kernel
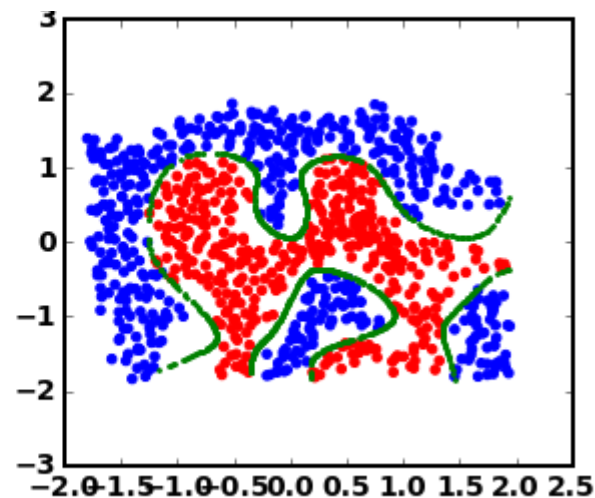
# Example: Soft-SVM with Gaussian Kernel

- A common choice is to use the Gaussian kernel

$$K(\underline{u}, \underline{v}) = e^{-\frac{\|\underline{u}-\underline{v}\|^2}{2\sigma^2}}$$

- $\sigma$ called the **bandwidth**

- It controls smoothness of the boundary

- Gives a notion of model complexity:

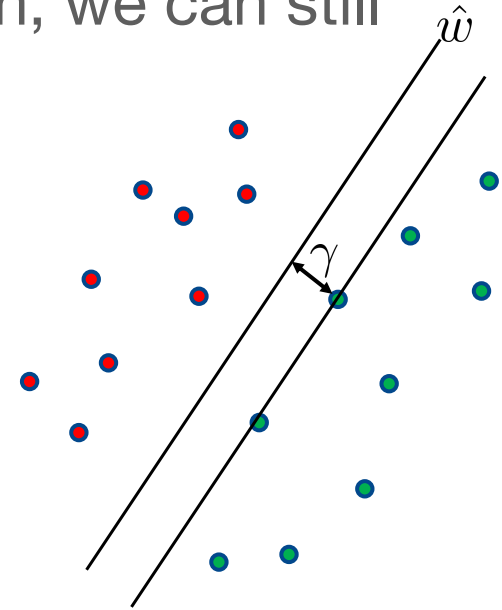- large $\sigma \rightarrow$ simple boundaries, small $\sigma \rightarrow$ complex boundaries



large $\sigma$        small $\sigma$

# What about overfitting?

- If we're using such a high-dimensional representation, won't the dimension be large?  Doesn't this lead to overfitting?

- Margin:  If we can find a solution with large margin, we can still avoid overfitting.
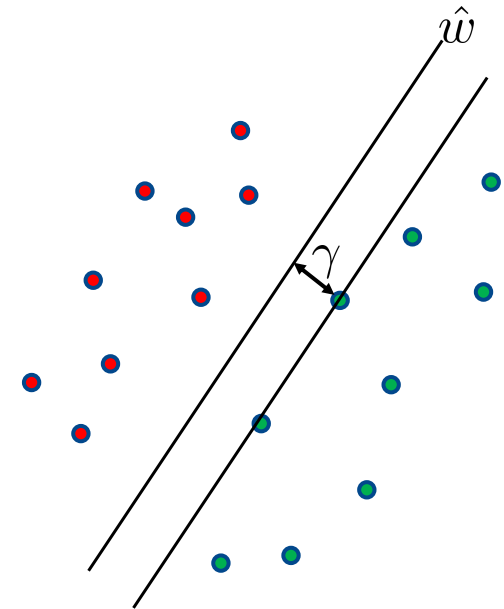
- This is true even with kernels

$\hat{w}$

$\gamma$

# What about overfitting?

- Margin for kernel SVM:

geometric margin of classifier w:

$$\gamma = \min_{1 \le i \le n} y_i \left( \frac{w^\top}{\|w\|} \phi(x_i) \right)$$

Recall: The $\hat{w}$ solution of SVM primal problem satisfies $\|\hat{w}\| = \frac{1}{\gamma}$

For kernel SVM, this means

$$\frac{1}{\gamma^2} = \|\hat{w}\|^2 = \hat{w}^\top \hat{w} = \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

Solve for $\gamma$ to compute the margin

So we can also calculate the margin without computing $\phi$, using $K$

# Summary

- Kernel methods are a convenient family of algorithms

- They allow us to specify a non-linear representation for the classifier just by providing a kernel function, and the rest is automatic

- In this sense, they are "plug-and-play": very easy to use

- We need to be careful that the implicit high-dimensional representation doesn't lead to overfitting

- If the solution has large margin, it can avoid overfitting

- Recall that SVM is designed to maximize the margin, so it is well-suited to this type of guarantee