

Kernel Methods and Support Vector Machines (SVMs)

CMSC 678

UMBC

Outline

Recap

Kernel Methods & Feature Mapping

Support Vector Machines (SVMs)

Math: Subgradients

Math: Lagrange Multipliers for Constrained Optimization

SVMs Revisited: Dual Formulation

Recap from last time...

K-Means

Initialize k centers by picking k points randomly among all the points

Repeat till convergence (or max iterations)

Assign each point to the nearest center (assignment step)

Estimate the mean of each group (update step)

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mu_i\|^2$$

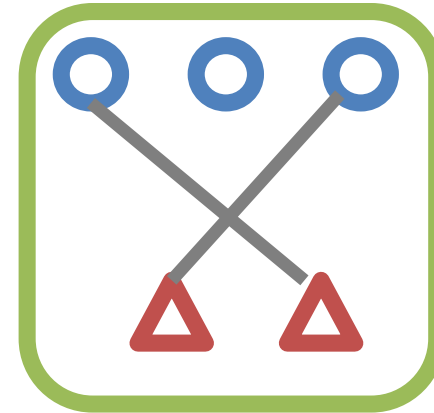
k-means++ algorithm for initialization:

1. Chose one center uniformly at random among all the points
2. For each point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the nearest center that has already been chosen
3. Chose one new data point at random as a new center, using a weighted probability distribution where a point \mathbf{x} is chosen with a probability proportional to $D(\mathbf{x})^2$
4. Repeat Steps 2 and 3 until k centers have been chosen

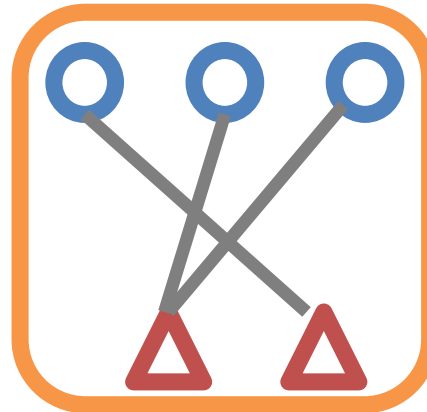
Clustering Evaluation

(Classification:
accuracy, recall,
precision, F-score)

Greedy mapping:
one-to-one



Optimistic mapping:
many-to-one



Rigorous/information
theoretic: **V-measure**

homogeneity

$$= \begin{cases} 1, & H(K, C) = 0 \\ 1 - \frac{H(C|K)}{H(K)}, & \text{o/w} \end{cases}$$

completeness

$$= \begin{cases} 1, & H(K, C) = 0 \\ 1 - \frac{H(K|C)}{H(C)}, & \text{o/w} \end{cases}$$

Hierarchical clustering

Agglomerative: a “bottom up” approach where elements start as individual clusters and clusters are merged as one moves up the hierarchy

Divisive: a “top down” approach where elements start as a single cluster and clusters are split as one moves down the hierarchy

Agglomerative clustering:

First merge very similar instances

Incrementally build larger clusters out of smaller clusters

Algorithm:

Maintain a set of clusters

Initially, each instance in its own cluster

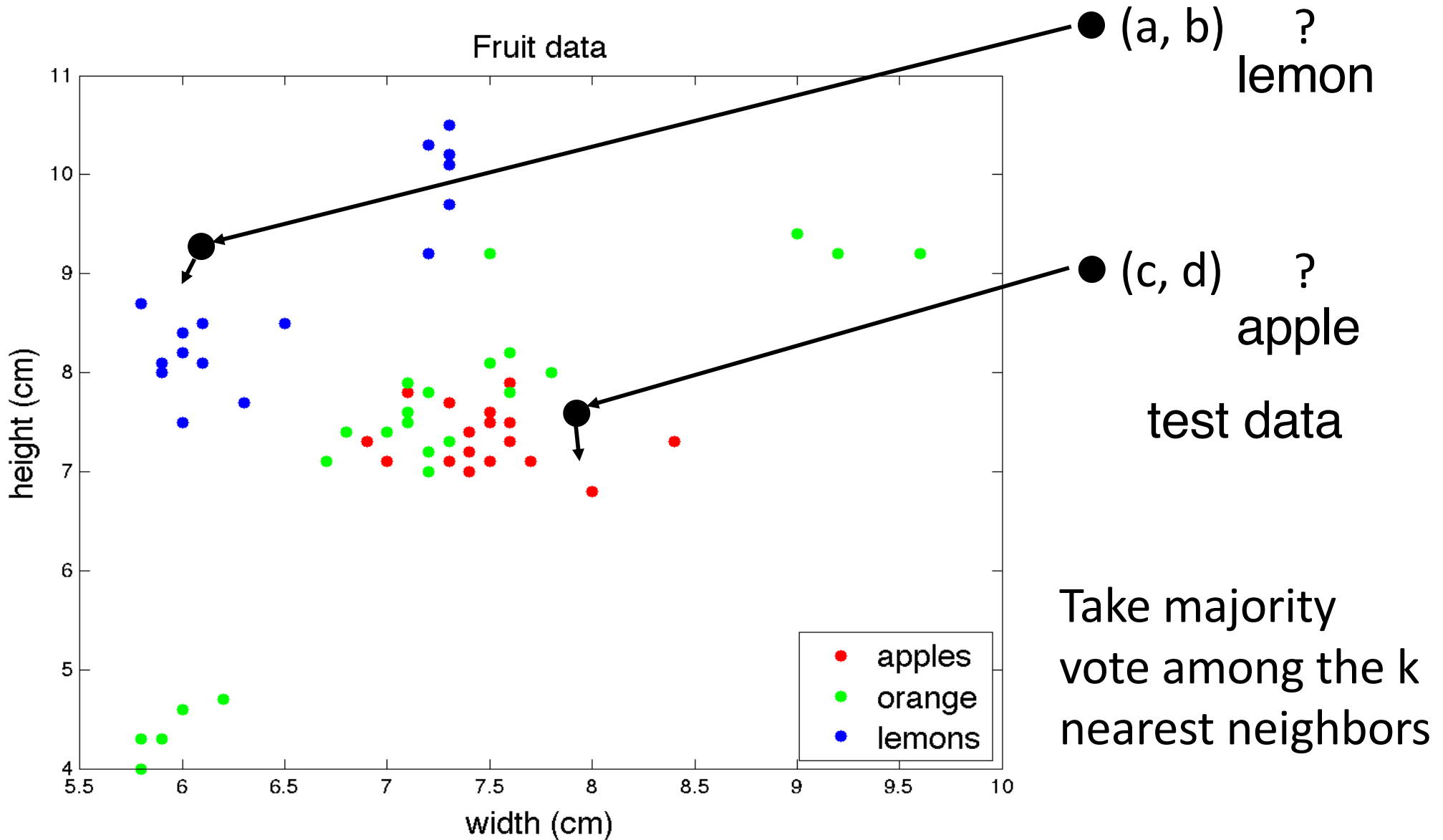
Repeat:

Pick the two “closest” clusters

Merge them into a new cluster

Stop when there’s only one cluster left

K-Nearest neighbor predictor



Outline

Recap

Kernel Methods & Feature Mapping

Support Vector Machines (SVMs)

Math: Subgradients

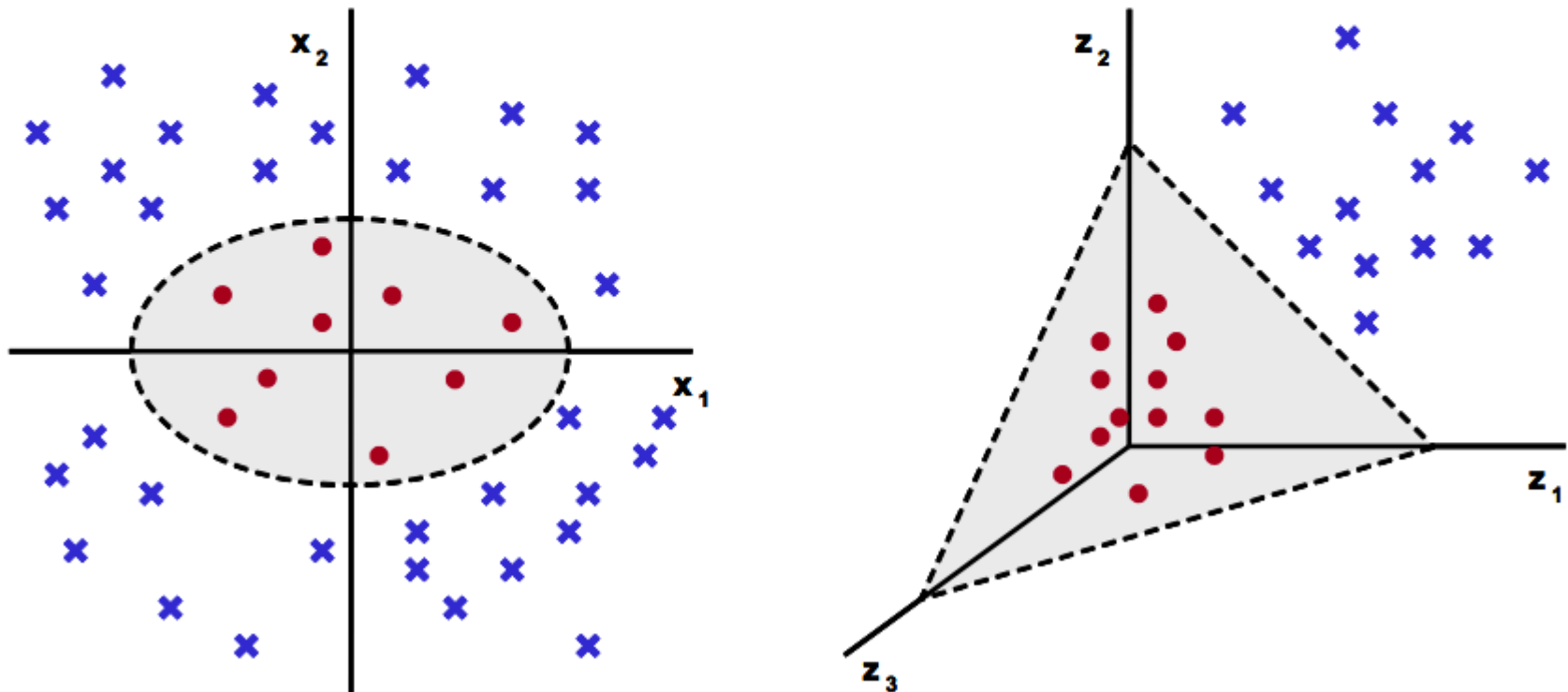
Math: Lagrange Multipliers for Constrained Optimization

SVMs Revisited: Dual Formulation

Feature mapping

$$\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

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



Learn non-linear classifiers by mapping features

Quadratic feature map

$$\mathbf{x} = [x_1, x_2, \dots, x_D] \implies \phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \dots, \sqrt{2}x_D, x_1^2, x_1x_2, x_1x_3 \dots, x_1x_D, x_2x_1, x_2^2, x_2x_3 \dots, x_2x_D, \dots, x_Dx_1, x_Dx_2, x_Dx_3 \dots, x_D^2]$$

Contains all **single** and **pairwise** terms

Quadratic kernel

The **dot product** between **feature maps** of \mathbf{x} and \mathbf{z} is:

$$\begin{aligned}\phi(\mathbf{x})^T \phi(\mathbf{z}) &= 1 + 2x_1z_1 + 2x_2z_2, \dots, 2x_Dz_D + x_1^2z_1^2 + x_1x_2z_1z_2 + \dots + x_1x_Dz_1z_D + \dots \\ &\quad \dots + x_Dx_1z_Dz_1 + x_Dx_2z_Dz_2 + \dots + x_D^2z_D^2 \\ &= 1 + 2 \left(\sum_i x_i z_i \right) + \sum_{i,j} x_i x_j z_i z_j \\ &= 1 + 2 (\mathbf{x}^T \mathbf{z}) + (\mathbf{x}^T \mathbf{z})^2 \\ &= (1 + \mathbf{x}^T \mathbf{z})^2 \\ &= K(\mathbf{x}, \mathbf{z}) \quad \longleftarrow \text{quadratic kernel}\end{aligned}$$

Compute $\phi(\mathbf{x})^T \phi(\mathbf{z})$ in **almost the same time** as needed to compute $\mathbf{x}^T \mathbf{z}$ (one extra **addition** and **multiplication**)

We will rewrite various algorithms using only **dot products** (or **kernel evaluations**), and not **explicit features**

Feature mapping

$$K(x, z) = (x^T z)^2$$

$$\begin{aligned} K(x, z) &= \left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right) \\ &= \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j \\ &= \sum_{i,j=1}^n (x_i x_j) (z_i z_j) \end{aligned}$$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

Feature mapping

$$K(x, z) = (x^T z + c)^2$$
$$= \sum_i x_i^2 z_i^2 + 2 \sum_i \sum_j^{i-1} x_i x_j z_i z_j + \sum_i 2c x_i z_i + c^2$$

$$\phi(x) = \begin{pmatrix} x_1 x_1 \\ x_1 x_2 \\ \dots \\ x_N x_N \\ x_1 \sqrt{2c} \\ \dots \\ x_N \sqrt{2c} \\ c \end{pmatrix}$$

Drawbacks of feature mapping

Computational

Suppose training time is linear in feature dimension, quadratic feature map squares the training time

Memory

Quadratic feature map squares the memory required to store the training data

Statistical

Quadratic feature mapping squares the number of parameters
For now lets assume that regularization will deal with overfitting

Perceptron revisited

Input: training data $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
feature map ϕ

Initialize $\mathbf{w} \leftarrow [0, \dots, 0]$

for iter = 1, ..., T

for i = 1, ..., n

predict according to the current model

$$\hat{y}_i = \begin{cases} +1 & \text{if } \mathbf{w}^T \phi(\mathbf{x}_i) > 0 \\ -1 & \text{otherwise} \end{cases}$$

if $y_i = \hat{y}_i$ no change

else, $\mathbf{w} \leftarrow \mathbf{w} + y_i \phi(\mathbf{x}_i)$

dependence on ϕ



Obtained by
replacing \mathbf{x} by
 $\phi(\mathbf{x})$

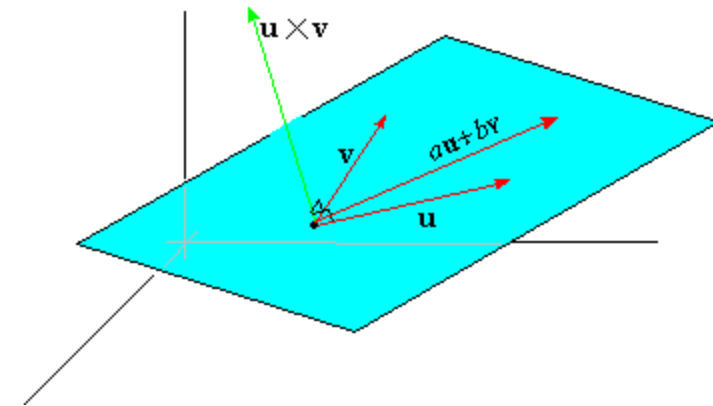
Properties of the weight vector

Linear algebra recap:

Let \mathbf{U} be set of vectors in \mathbb{R}^D , i.e., $\mathbf{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D\}$ and $\mathbf{u}_i \in \mathbb{R}^D$

$\text{Span}(\mathbf{U})$ is the set of all vectors that can be represented as $\sum_i a_i \mathbf{u}_i$, such that $a_i \in \mathbb{R}$

$\text{Null}(\mathbf{U})$ is everything that is left i.e., $\mathbb{R}^D \setminus \text{Span}(\mathbf{U})$



Perceptron representer theorem: During the run of the perceptron training algorithm, the weight vector \mathbf{w} is always in the span of $\phi(\mathbf{x}_1), \phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_D)$

$$\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i) \quad \text{updates} \quad \alpha_i \leftarrow \alpha_i + y_i$$

$$\mathbf{w}^T \phi(\mathbf{z}) = \left(\sum_i \alpha_i \phi(\mathbf{x}_i) \right)^T \phi(\mathbf{z}) = \sum_i \alpha_i \phi(\mathbf{x}_i)^T \phi(\mathbf{z})$$

Perceptron Representer Theorem

Perceptron representer theorem: During the run of the perceptron training algorithm, the weight vector \mathbf{w} is always in the span of $\phi(\mathbf{x}_1), \phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_D)$

$$\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i) \quad \text{updates} \quad \alpha_i \leftarrow \alpha_i + y_i$$

$$\mathbf{w}^T \phi(\mathbf{z}) = \left(\sum_i \alpha_i \phi(\mathbf{x}_i) \right)^T \phi(\mathbf{z}) = \sum_i \alpha_i \phi(\mathbf{x}_i)^T \phi(\mathbf{z})$$

Why? Weights are always computed from the training instance representations

Example 2: Incorrect $\rightarrow w += y_2 \phi(x_2)$

Example 3: Correct $\rightarrow w += 0 * y_3 \phi(x_3)$

Example 4: Incorrect $\rightarrow w += y_4 \phi(x_4)$

Kernelized perceptron

Input: **training data** $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
feature map ϕ

Kernelized perceptron training algorithm

Initialize $\alpha \leftarrow [0, 0, \dots, 0]$

for iter = 1, ..., T

for i = 1, ..., n

predict according to the current model

$$\hat{y}_i = \begin{cases} +1 & \text{if } \sum_n \alpha_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_i) > 0 \\ -1 & \text{otherwise} \end{cases}$$

if $y_i = \hat{y}_i$, no change

else, $\alpha_i = \alpha_i + y_i$

$$\phi(\mathbf{x})^T \phi(\mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^p \quad \text{polynomial kernel of degree } p$$

Kernelized perceptron

Input: **training data** $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
feature map ϕ

Kernelized perceptron training algorithm

Initialize $\alpha \leftarrow [0, 0, \dots, 0]$

for iter = 1, ..., T

for i = 1, ..., n

predict according to the current model

$$\hat{y}_i = \begin{cases} +1 & \text{if } \sum_n \alpha_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_i) > 0 \\ -1 & \text{otherwise} \end{cases}$$

if $y_i = \hat{y}_i$, no change

else, $\alpha_i = \alpha_i + y_i$

$$\phi(\mathbf{x})^T \phi(\mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^p \quad \text{polynomial kernel of degree } p$$

Upshot: Feature expansion for
"free"

Kernel k-means

Repeat till convergence (or max iterations)

1. **Assign** each point to the nearest center (**assignment step**)
2. Estimate the **mean** of each group (**update step**)

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\phi(\mathbf{x}) - \mu_i\|^2$$

Representer theorem is easy in (2)

$$\mu_i \leftarrow \frac{1}{|S_i|} \sum_{\mathbf{x} \in S_i} \phi(\mathbf{x})$$

What about (1)? How to compute using dot products

$$\|\phi(\mathbf{x}) - \mu_i\|^2$$

Kernel k-means

$$z_n = \arg \min_k \left\| \phi(\mathbf{x}_n) - \boldsymbol{\mu}^{(k)} \right\|^2$$

definition of z_n

$$= \arg \min_k \left\| \phi(\mathbf{x}_n) - \sum_m \alpha_m^{(k)} \phi(\mathbf{x}_m) \right\|^2$$

definition of $\boldsymbol{\mu}^{(k)}$

$$= \arg \min_k \left\| \phi(\mathbf{x}_n) \right\|^2 + \left\| \sum_m \alpha_m^{(k)} \phi(\mathbf{x}_m) \right\|^2 + \phi(\mathbf{x}_n) \cdot \left[\sum_m \alpha_m^{(k)} \phi(\mathbf{x}_m) \right]$$

expand quadratic term

$$= \arg \min_k \sum_m \sum_{m'} \alpha_m^{(k)} \alpha_{m'}^{(k)} \phi(\mathbf{x}_m) \cdot \phi(\mathbf{x}_{m'}) + \sum_m \alpha_m^{(k)} \phi(\mathbf{x}_m) \cdot \phi(\mathbf{x}_n) + \text{const}$$

linearity and constant

Kernel Definition

A **kernel** is a **mapping** $K: X \times X \rightarrow \mathbb{R}$

Functions that can be written as

dot products are valid **kernels**

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

Examples: **polynomial kernel**

$$K_{(\text{poly})}^d(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^d$$

Kernel Definition

A **kernel** is a **mapping** $K: X \times X \rightarrow \mathbb{R}$

Functions that can be written as

dot products are valid **kernels**

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

Examples: **polynomial kernel**

$$K_{(\text{poly})}^d(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^d$$

Alternatively:

Mercer's Conditions: A function $K: X \times X \rightarrow \mathbb{R}$ is a **kernel** if K is **positive semi-definite** (psd)

$$\int f(\mathbf{x})^2 d\mathbf{x} < \infty$$

This means that ~~for~~ all functions f that are **squared integrable** except the **zero function**, the following property holds:

$$\int \int f(\mathbf{x}) K(\mathbf{x}, \mathbf{z}) f(\mathbf{z}) d\mathbf{z} d\mathbf{x} > 0$$

Why is this characterization useful?

We can prove some properties about **kernels** that are otherwise hard to prove

Theorem: If K_1 and K_2 are **kernels**, then $K_1 + K_2$ is also a **kernel**

Proof:

$$\begin{aligned}\int \int f(\mathbf{x})K(\mathbf{x}, \mathbf{z})f(\mathbf{z})d\mathbf{z}d\mathbf{x} &= \int \int f(\mathbf{x}) (K_1(\mathbf{x}, \mathbf{z}) + K_2(\mathbf{x}, \mathbf{z})) f(\mathbf{z})d\mathbf{z}d\mathbf{x} \\ &= \int \int f(\mathbf{x})K_1(\mathbf{x}, \mathbf{z})f(\mathbf{z})d\mathbf{z}d\mathbf{x} + \int \int f(\mathbf{x})K_2(\mathbf{x}, \mathbf{z})f(\mathbf{z})d\mathbf{z}d\mathbf{x} \\ &\geq 0 + 0\end{aligned}$$

More generally if K_1, K_2, \dots, K_n are **kernels** then $\sum_i \alpha_i K_i$ with $\alpha_i \geq 0$, is also a **kernel**

Can build new **kernels** by linearly combining existing **kernels**

Kernels in practice

Feature mapping via kernels often improves performance

MNIST digits test error:

8.4% SVM linear

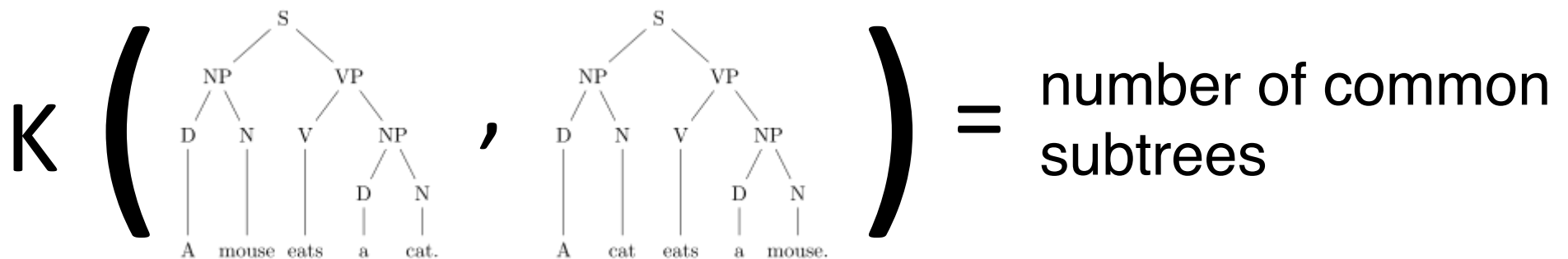
1.4% SVM RBF

1.1% SVM polynomial (d=4)



Kernels over general structures

Kernels can be defined over any pair of inputs such as strings, trees and graphs!

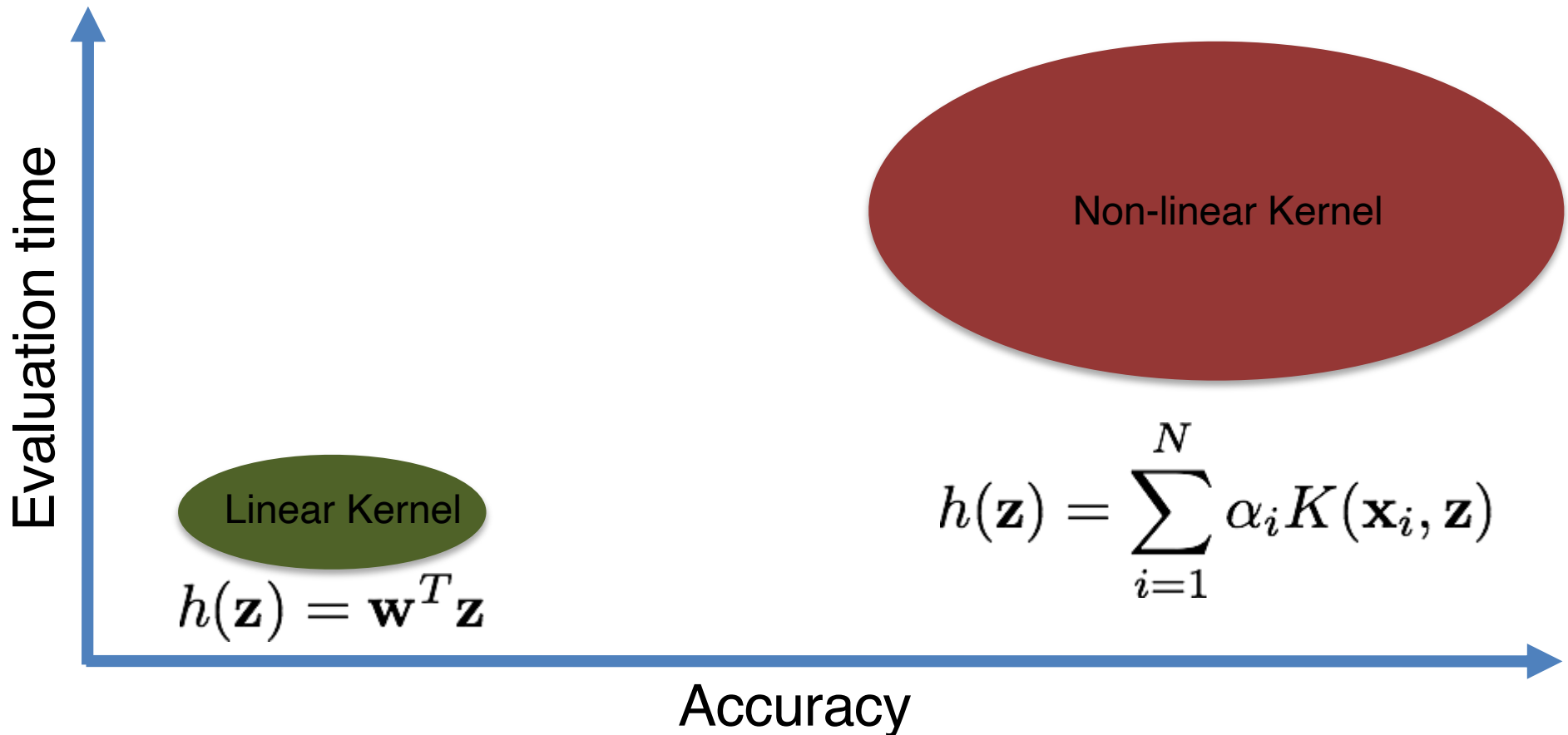


http://en.wikipedia.org/wiki/Tree_kernel

For strings number of common substrings is a kernel

Graph kernels that measure graph similarity (e.g. number of common subgraphs) have been used to predict toxicity of chemical structures

Kernel classifiers tradeoffs



Linear: O (feature dimension)

Non Linear: O (**N** X feature dimension)

Outline

Recap

Kernel Methods & Feature Mapping

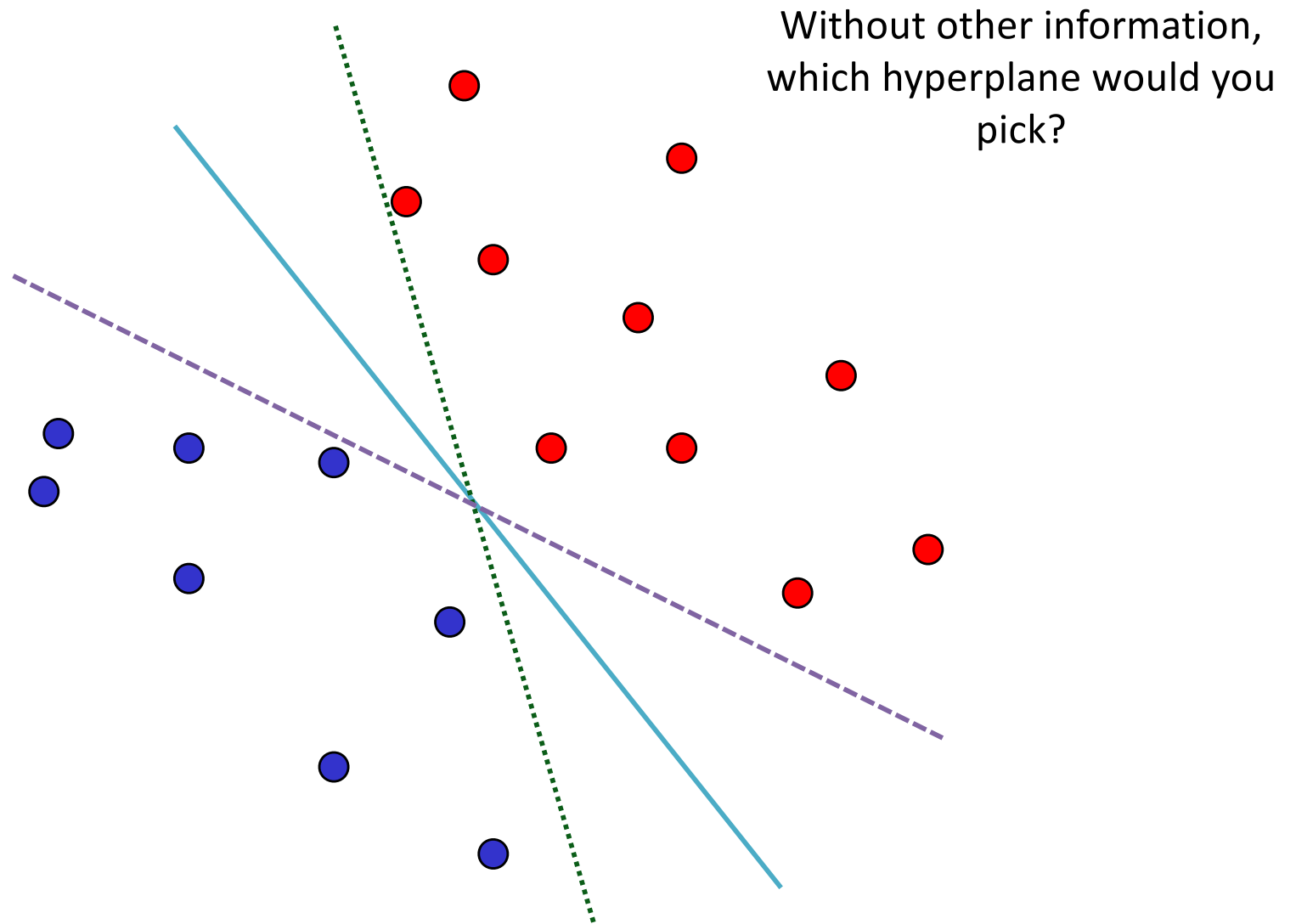
Support Vector Machines (SVMs)

Math: Subgradients

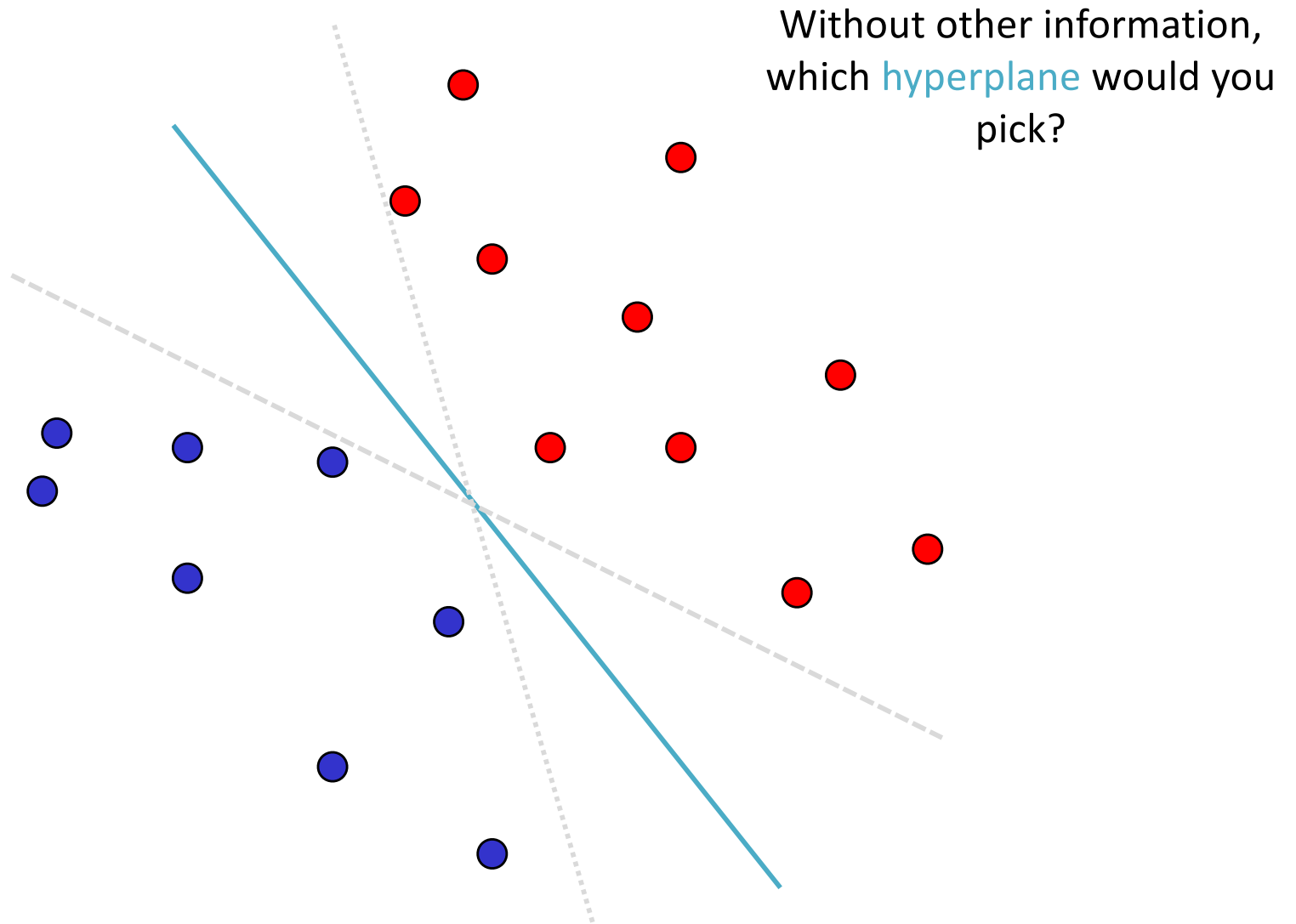
Math: Lagrange Multipliers for Constrained Optimization

SVMs Revisited: Dual Formulation

Picking a good hyperplane

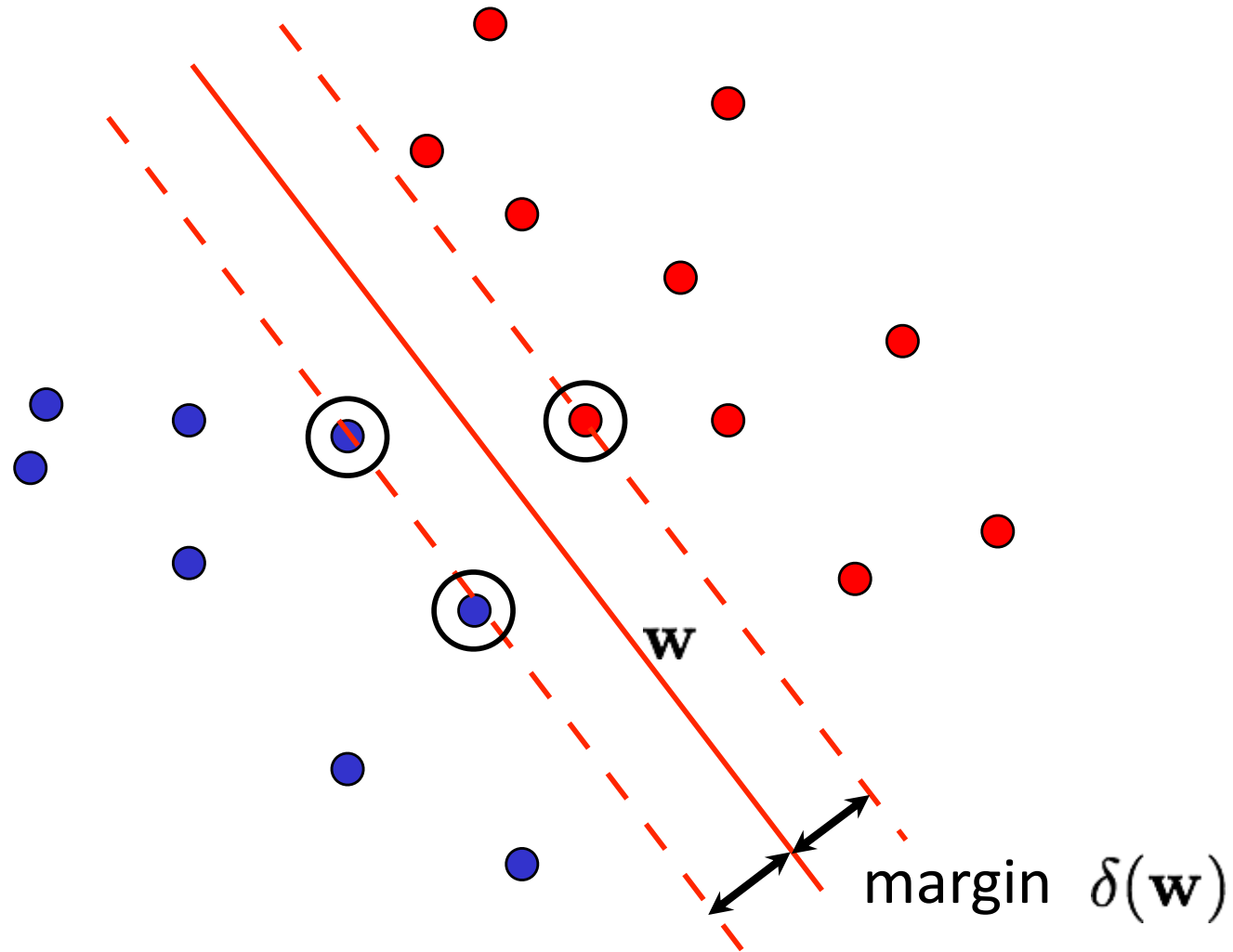


Picking a good hyperplane



Support Vector Machines (SVMs)

Maximize the distance to the nearest point (margin),
while correctly classifying all the points



Optimization for SVMs

Separable case: [hard margin SVM](#)

$$\min_{\mathbf{w}} \frac{1}{\delta(\mathbf{w})}$$

maximize margin

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1, \forall n$$

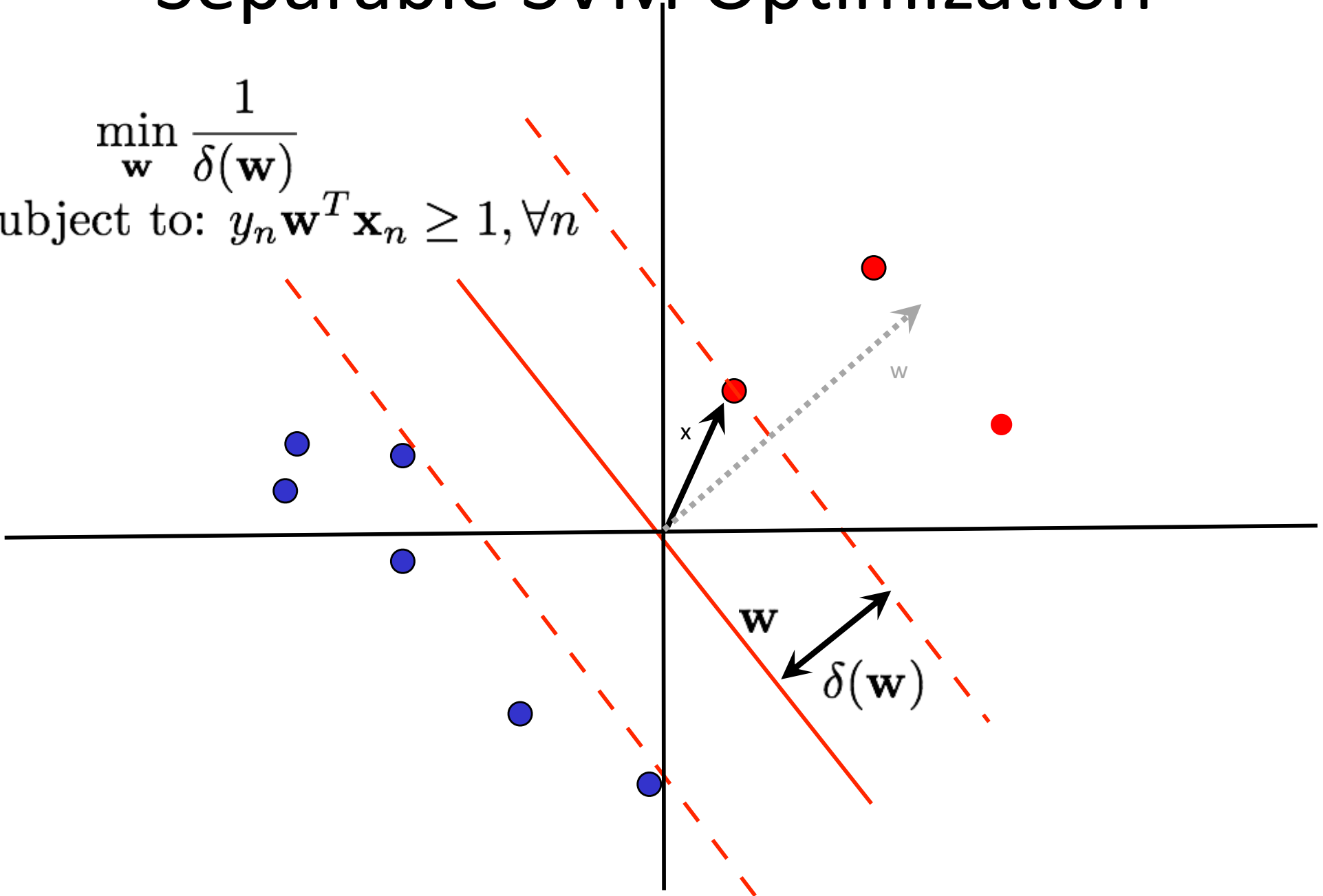
(assuming $y = +1$ or $y = -1$)

separate by a non-trivial margin

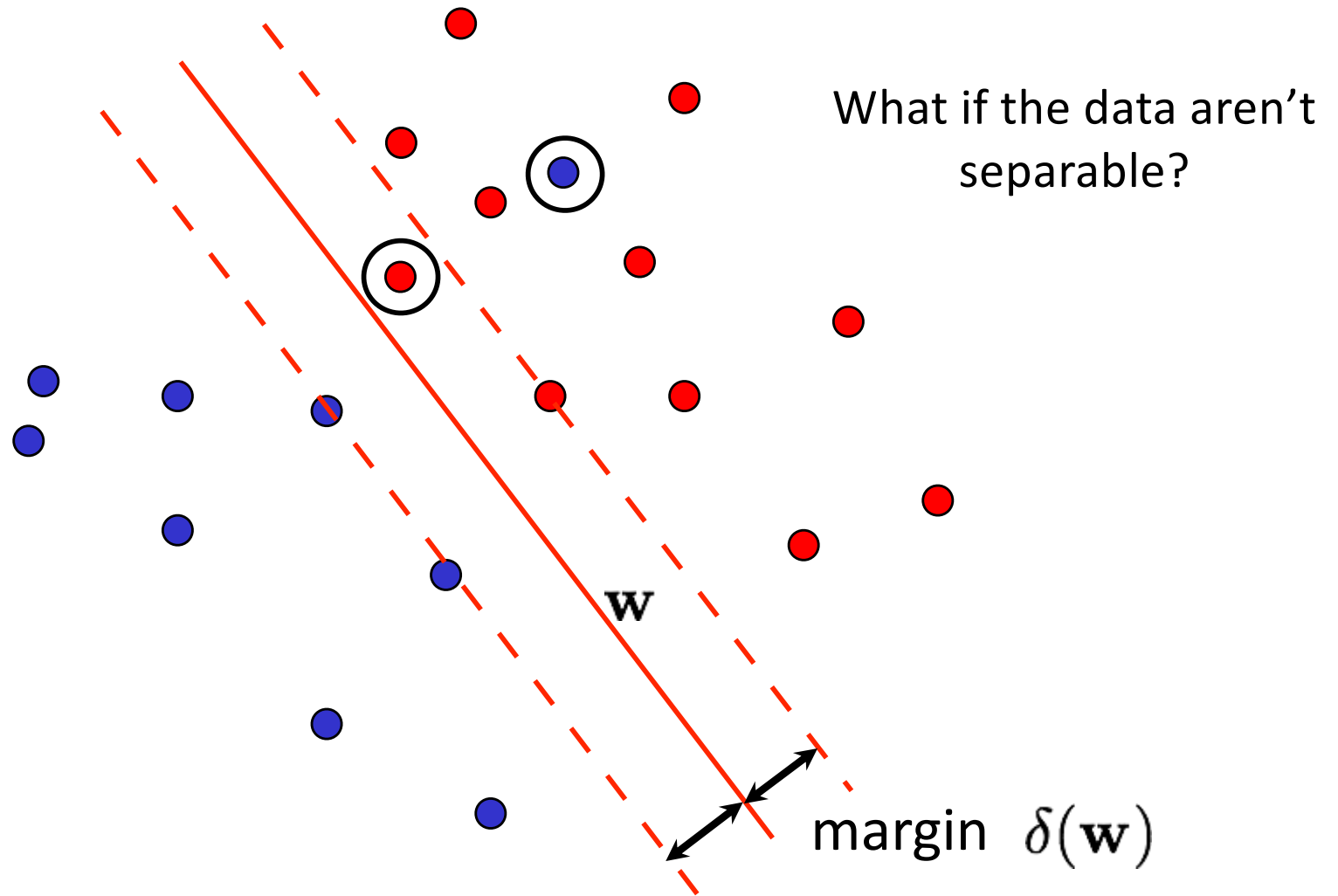
Separable SVM Optimization

$$\min_{\mathbf{w}} \frac{1}{\delta(\mathbf{w})}$$

subject to: $y_n \mathbf{w}^T \mathbf{x}_n \geq 1, \forall n$



Support Vector Machines (SVMs)



Optimization for SVMs

Separable case: [hard margin SVM](#)

$$\min_{\mathbf{w}} \frac{1}{\delta(\mathbf{w})}$$

maximize margin

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1, \forall n$$

separate by a non-trivial margin

Non-separable case: [soft margin SVM](#)

$$\min_{\mathbf{w}} \frac{1}{\delta(\mathbf{w})} + C \sum_n \xi_n$$

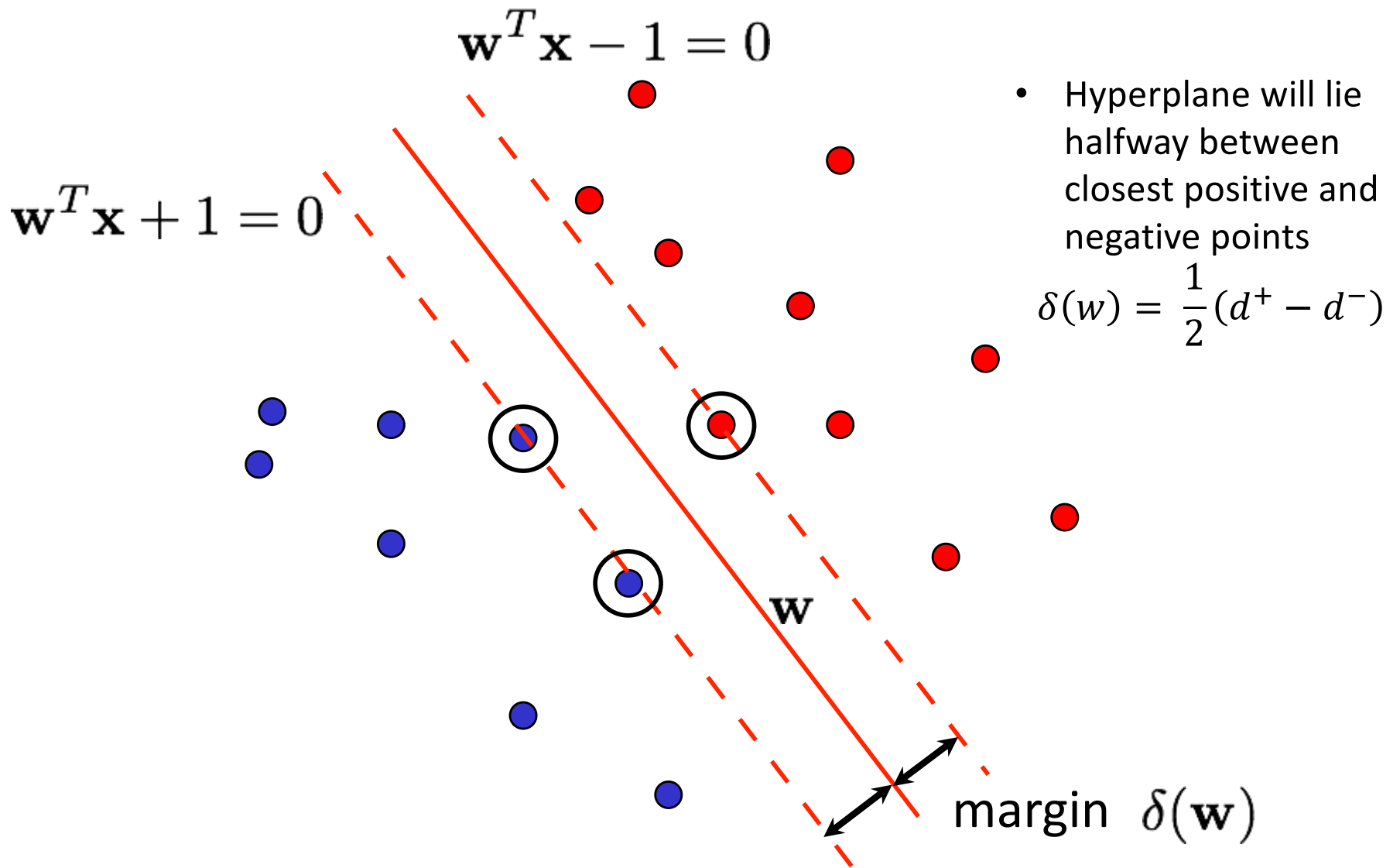
maximize margin minimize slack

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

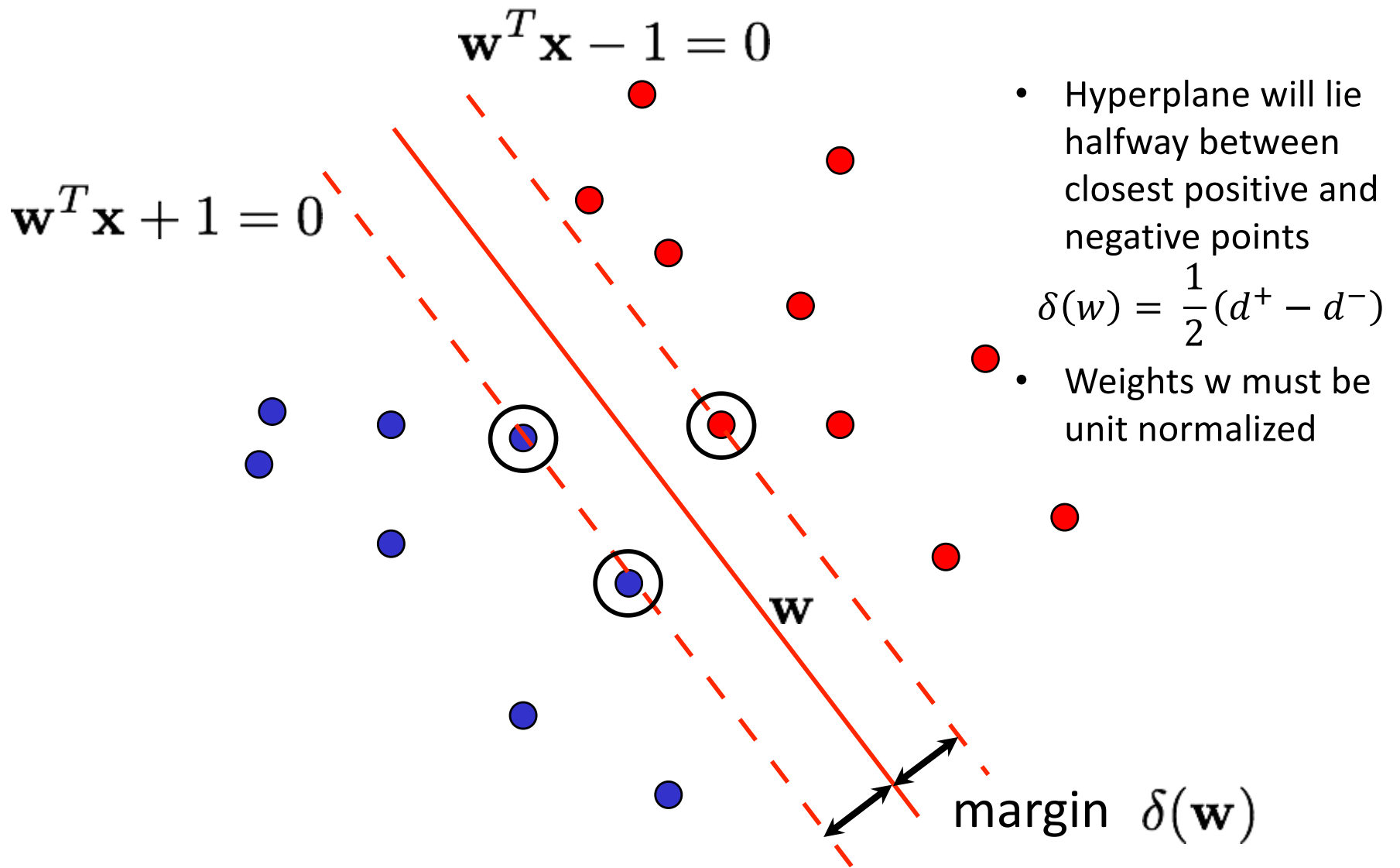
allow some slack

$$\xi_n \geq 0$$

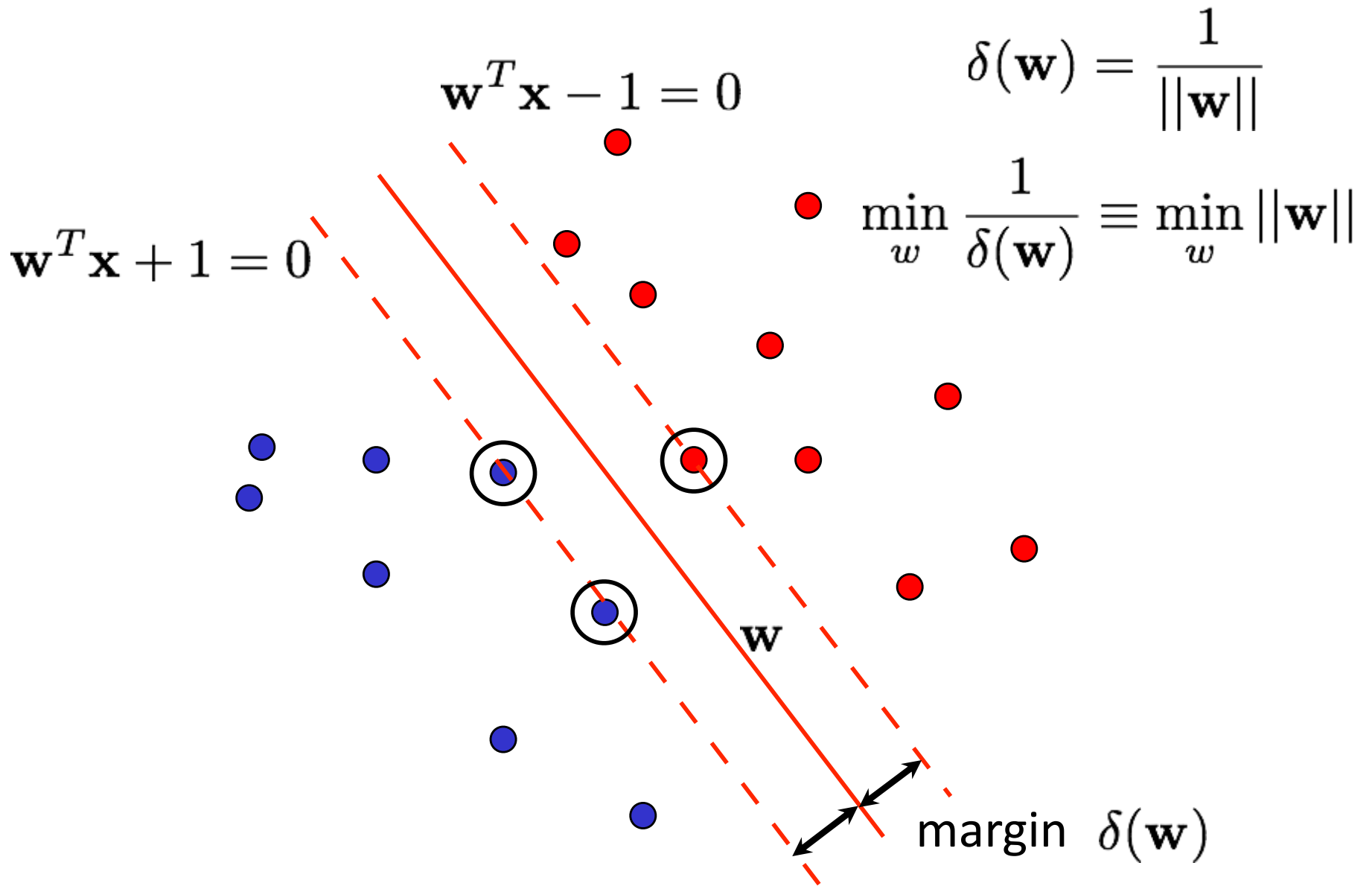
Maximizing Margin = Minimizing norm



Maximizing Margin = Minimizing norm



Maximizing Margin = Minimizing norm



Equivalent optimization for SVMs

Separable case: [hard margin SVM](#)

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

maximize margin

squaring and half for convenience

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1, \forall n$$

separate by a non-trivial margin

Non-separable case: [soft margin SVM](#)

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \xi_n$$

maximize margin minimize slack

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

allow some slack

$$\xi_n \geq 0$$

Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_n$$

soft margin SVM

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

$$\xi_n \geq 0$$

Given weights + bias, can you derive the optimal **slack** for the n^{th} example?

Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_n$$

soft margin SVM

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

$$\xi_n \geq 0$$

Given weights + bias, can you derive the optimal **slack** for the n^{th} example?

$$y_n \mathbf{w}^T \mathbf{x}_n = 0.8, \xi_n = ?$$

$$y_n \mathbf{w}^T \mathbf{x}_n = -1, \xi_n = ?$$

$$y_n \mathbf{w}^T \mathbf{x}_n = 2.5, \xi_n = ?$$

Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_n$$

soft margin SVM

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

$$\xi_n \geq 0$$

Given weights + bias, can you derive the optimal **slack** for the n^{th} example?

$$y_n \mathbf{w}^T \mathbf{x}_n = 0.8, \xi_n = ? \quad 0.2$$

$$y_n \mathbf{w}^T \mathbf{x}_n = -1, \xi_n = ? \quad 2.0$$

$$y_n \mathbf{w}^T \mathbf{x}_n = 2.5, \xi_n = ? \quad 0$$

Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_n$$

soft margin SVM

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

$$\xi_n \geq 0$$

Given weights + bias, can you derive the optimal **slack** for the n^{th} example?

$$y_n \mathbf{w}^T \mathbf{x}_n = 0.8, \xi_n = ? \quad 0.2$$

$$y_n \mathbf{w}^T \mathbf{x}_n = -1, \xi_n = ? \quad 2.0$$

$$y_n \mathbf{w}^T \mathbf{x}_n = 2.5, \xi_n = ? \quad 0$$

$$\xi_n = \begin{cases} 0 & y_n \mathbf{w}^T \mathbf{x}_n \geq 1 \\ 1 - y_n \mathbf{w}^T \mathbf{x}_n & \text{otherwise} \end{cases}$$

Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_n$$

soft margin SVM

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

$$\xi_n \geq 0$$

Given weights + bias, can you derive the optimal **slack** for the n^{th} example?

$$y_n \mathbf{w}^T \mathbf{x}_n = 0.8, \xi_n = ? \quad 0.2$$

$$y_n \mathbf{w}^T \mathbf{x}_n = -1, \xi_n = ? \quad 2.0$$

$$y_n \mathbf{w}^T \mathbf{x}_n = 2.5, \xi_n = ? \quad 0$$

$$\xi_n = \begin{cases} 0 & y_n \mathbf{w}^T \mathbf{x}_n \geq 1 \\ 1 - y_n \mathbf{w}^T \mathbf{x}_n & \text{otherwise} \end{cases}$$

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n)$$

Same as hinge loss with squared norm regularization!

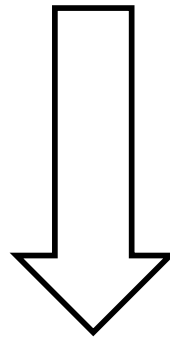
Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_n$$

soft margin SVM

$$\text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n$$

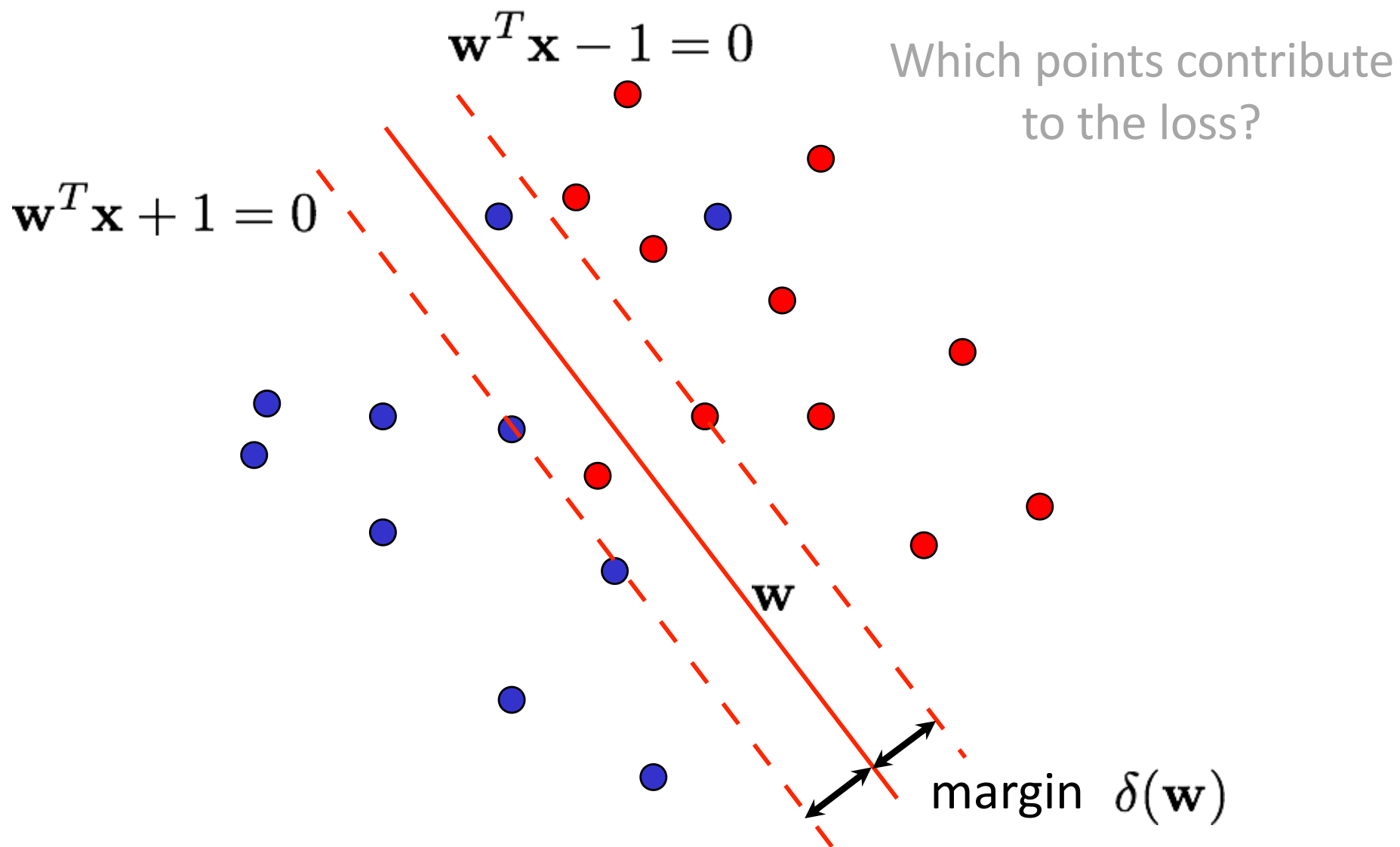
$$\xi_n \geq 0$$



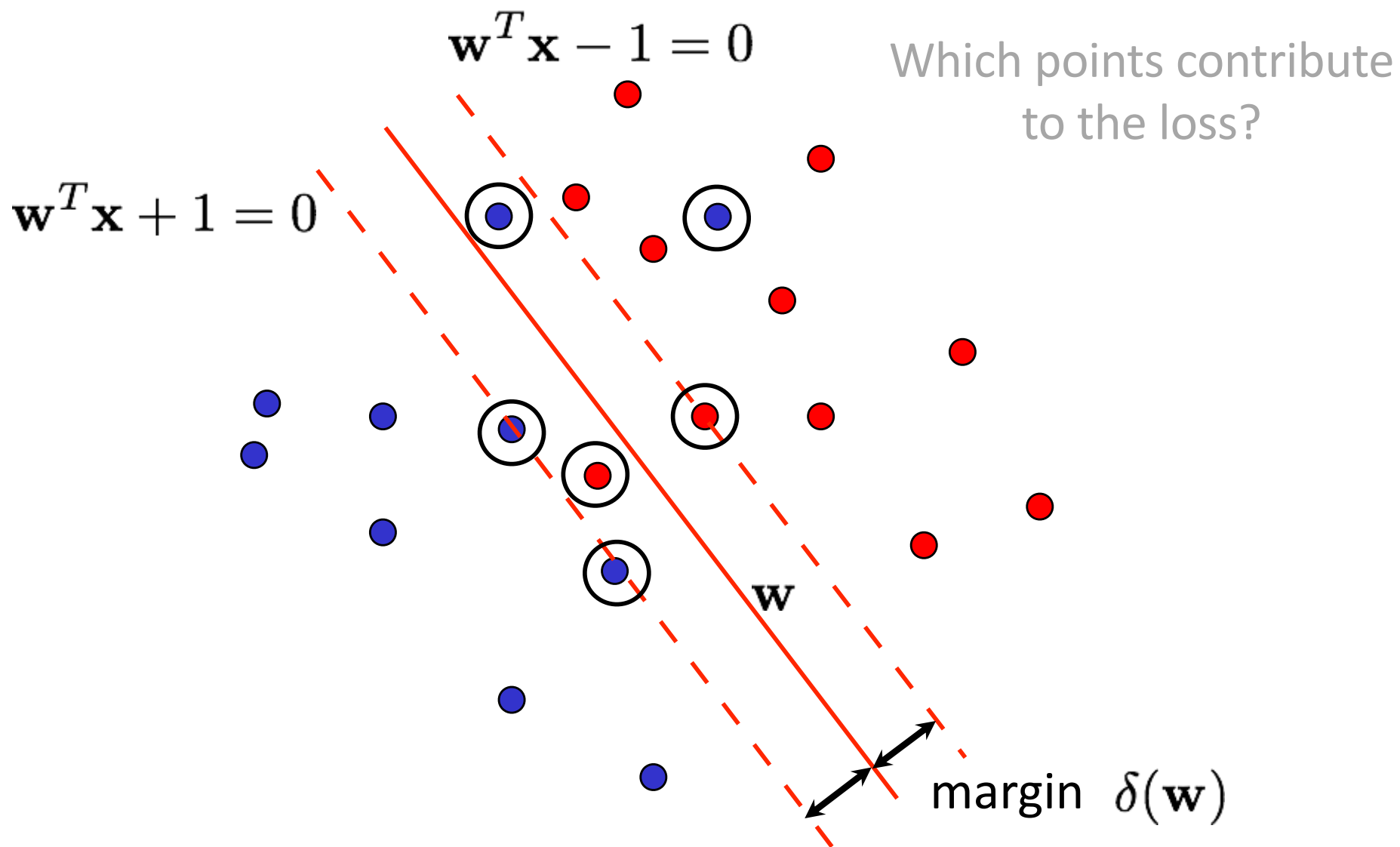
$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n)$$

Same as hinge loss with squared norm regularization!

Support vectors



Support vectors



Optimization for linear models

Under suitable conditions*, provided you pick the step sizes appropriately, the convergence rate of gradient descent is $O(1/N)$

i.e., if you want a solution within 0.001 of the optimal you have to run the gradient descent for $N=1000$ iterations.

For linear models ([hinge/logistic/exponential loss](#)) and [squared-norm regularization](#) there are off-the-shelf solvers that are fast in practice: [SVM^{perf}](#), [LIBLINEAR](#), [PEGASOS](#)

[SVM^{perf}](#), [LIBLINEAR](#) use a different optimization method

* the function is strongly convex: $f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|_2^2$

Outline

Recap

Kernel Methods & Feature Mapping

Support Vector Machines (SVMs)

Math: Subgradients

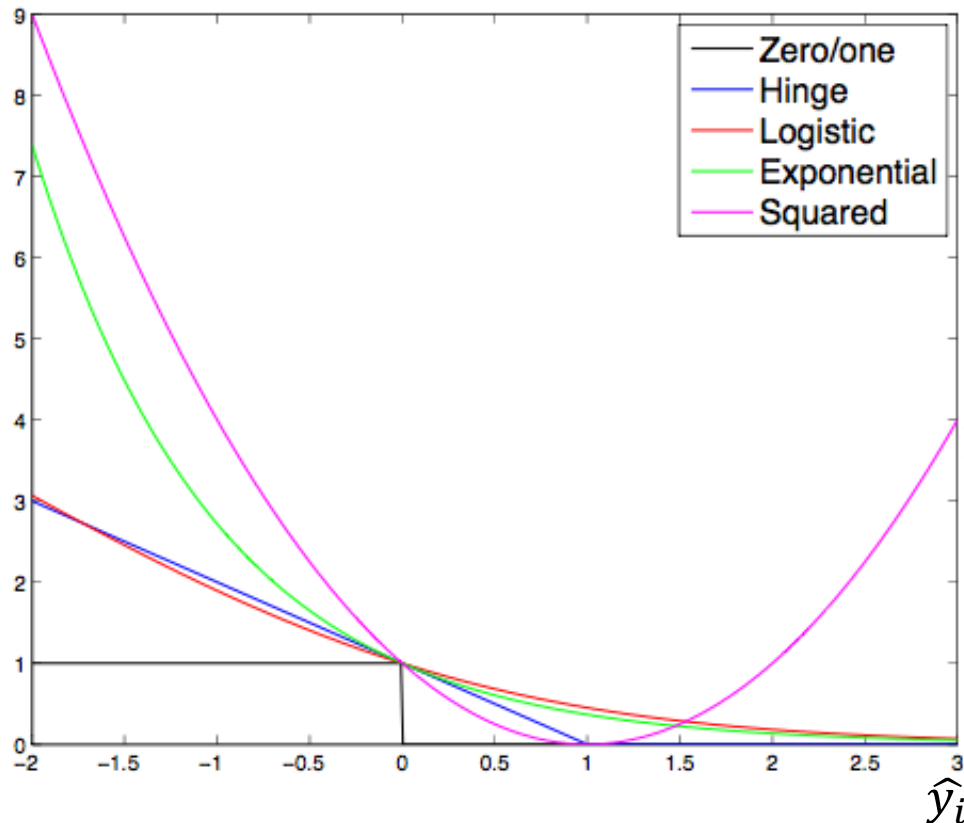
Math: Lagrange Multipliers for Constrained Optimization

SVMs Revisited: Dual Formulation

Reprise from earlier in the semester: Convex surrogate loss functions

Surrogate loss: replace **Zero/one loss** by a smooth function

Easier to optimize if the **surrogate loss** is **convex**



$$y = +1 \quad \hat{y} \leftarrow \mathbf{w}^T \mathbf{x}$$

Zero/one: $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$

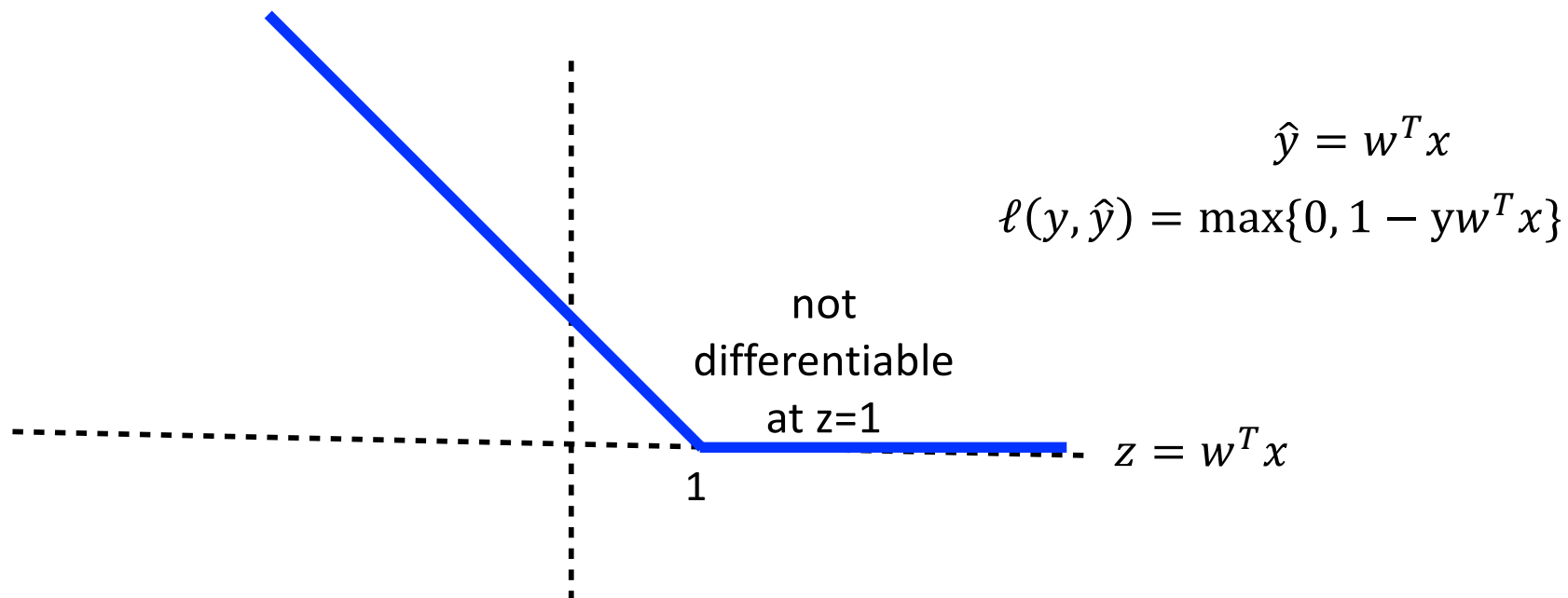
Hinge: $\ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$

Logistic: $\ell^{(\text{log})}(y, \hat{y}) = \frac{1}{\log 2} \log(1 + \exp[-y\hat{y}])$

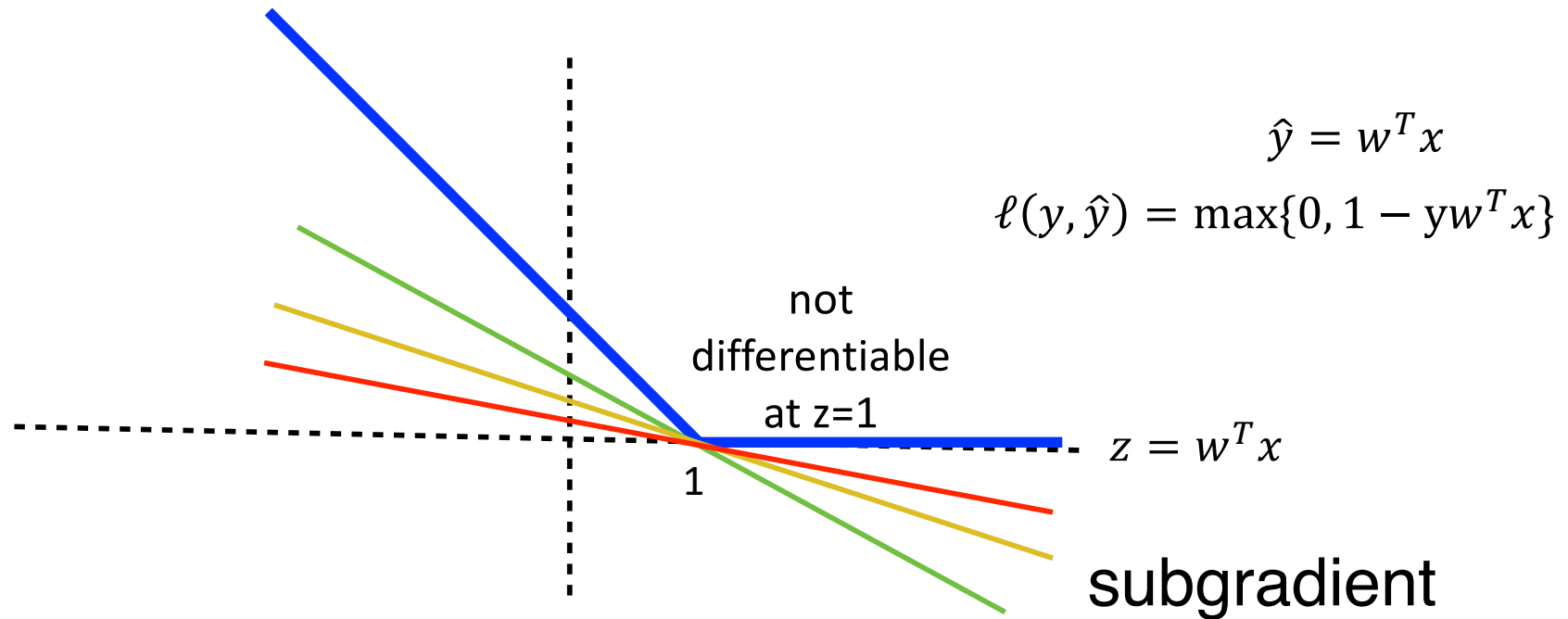
Exponential: $\ell^{(\text{exp})}(y, \hat{y}) = \exp[-y\hat{y}]$

Squared: $\ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$

Gradient of Hinge Loss?

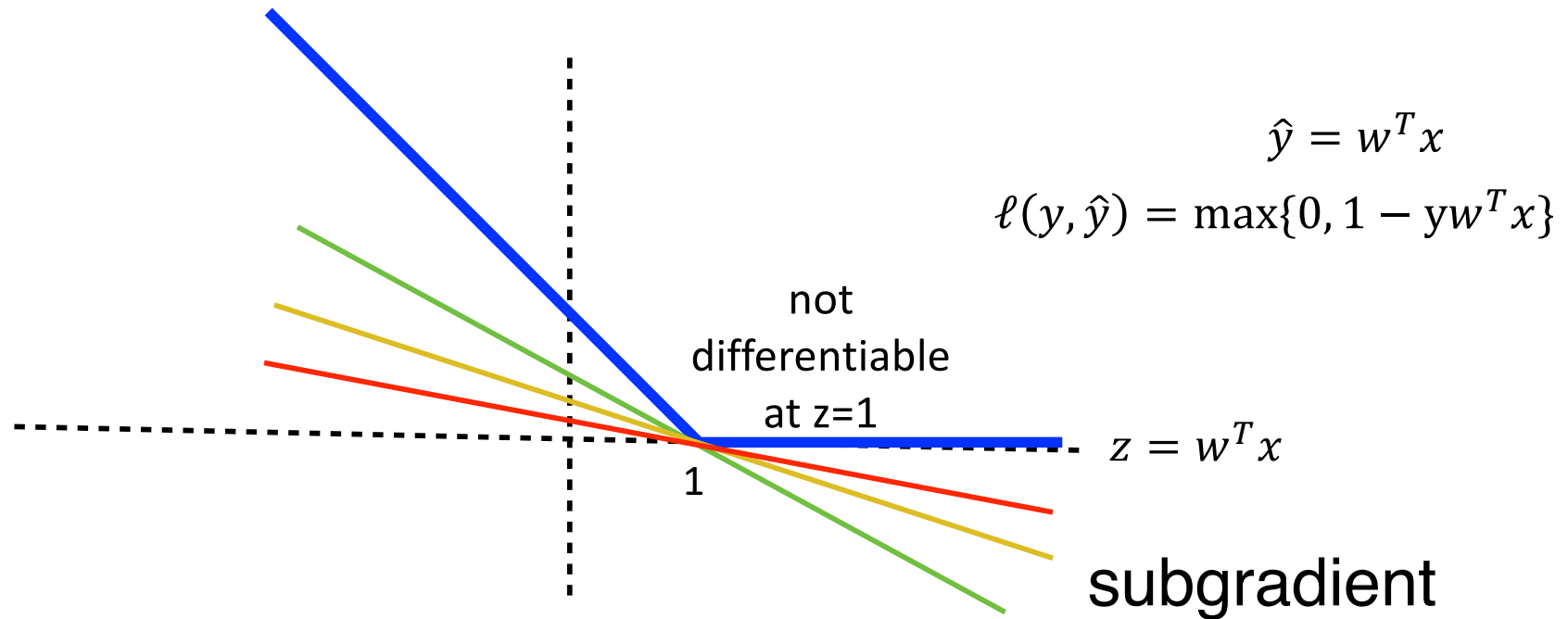


Subgradient of Hinge Los



Subgradient is any direction that is **below** the function

Subgradient of Hinge Los



Subgradient is any direction that is **below** the function

For the **hinge loss** a possible **subgradient** is:

$$\frac{d\ell^{\text{hinge}}}{d\mathbf{w}} = \begin{cases} 0 & \text{if } y\mathbf{w}^T \mathbf{x} > 1 \\ -y\mathbf{x} & \text{otherwise} \end{cases}$$

Example: Hinge loss

$$\mathcal{L}(\mathbf{w}) = \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

Example: Hinge loss

$$\mathcal{L}(\mathbf{w}) = \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

$$\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_n -\mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n \leq 1] y_n \mathbf{x}_n + \lambda \mathbf{w} \quad \text{subgradient}$$

Example: Hinge loss

$$\mathcal{L}(\mathbf{w}) = \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

$$\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_n -\mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n \leq 1] y_n \mathbf{x}_n + \lambda \mathbf{w} \quad \text{subgradient}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_n -\mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n \leq 1] y_n \mathbf{x}_n + \lambda \mathbf{w} \right) \quad \text{update}$$

Example: Hinge loss

$$\mathcal{L}(\mathbf{w}) = \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

$$\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_n -\mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n \leq 1] y_n \mathbf{x}_n + \lambda \mathbf{w} \quad \text{subgradient}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_n -\mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n \leq 1] y_n \mathbf{x}_n + \lambda \mathbf{w} \right) \quad \text{update}$$

loss term

$$\mathbf{w} \leftarrow \mathbf{w} + \eta y_n \mathbf{x}_n$$



only for points $y_n \mathbf{w}^T \mathbf{x}_n \leq 1$

perceptron update $y_n \mathbf{w}^T \mathbf{x}_n \leq 0$

regularization term

$$\mathbf{w} \leftarrow (1 - \eta\lambda) \mathbf{w}$$

shrinks weights towards zero

Outline

Recap

Kernel Methods & Feature Mapping

Support Vector Machines (SVMs)

Math: Subgradients

Math: Lagrange Multipliers for Constrained Optimization

SVMs Revisited: Dual Formulation

Lagrange multipliers

Assume an original optimization problem

$$\begin{aligned} \min_x \quad & x^2 \\ \text{st.} \quad & x \geq b \end{aligned}$$

Lagrange multipliers

Assume an original optimization problem

$$\begin{aligned} \min_x \quad & x^2 \\ \text{st.} \quad & x \geq b \end{aligned}$$

We convert it to a new optimization problem:

$$\begin{aligned} \min_x \max_{\alpha} \quad & L(x, \alpha) \\ \text{st.} \quad & \alpha \geq 0 \\ L(x, \alpha) = & x^2 - \alpha(x - b) \end{aligned}$$

Lagrange multipliers: an equivalent problem?

$$\begin{array}{ll} \min_x & x^2 \\ \text{st.} & x \geq b \end{array} \quad \Rightarrow \quad \begin{array}{ll} \min_x \max_{\alpha} & L(x, \alpha) \\ \text{st.} & \alpha \geq 0 \\ & L(x, \alpha) = x^2 - \alpha(x - b) \end{array}$$

$$x < b \quad \Rightarrow \quad (x - b) < 0 \quad \Rightarrow \quad \max_{\alpha} -\alpha(x - b) = \inf$$

Lagrange multipliers: an equivalent problem?

$$\begin{array}{ll} \min_x & x^2 \\ \text{st.} & x \geq b \end{array} \quad \Rightarrow \quad \begin{array}{ll} \min_x \max_{\alpha} & L(x, \alpha) \\ \text{st.} & \alpha \geq 0 \\ & L(x, \alpha) = x^2 - \alpha(x - b) \end{array}$$

$$x < b \quad \Rightarrow \quad (x - b) < 0 \quad \Rightarrow \quad \max_{\alpha} -\alpha(x - b) = \inf$$

$$x > b \quad \Rightarrow \quad (x - b) > 0 \quad \Rightarrow \quad \max_{\alpha} -\alpha(x - b) = 0, \alpha = 0 \quad \Rightarrow \quad L(x, \alpha) = x^2$$

Lagrange multipliers: an equivalent problem?

$$\begin{array}{ll} \min_x & x^2 \\ \text{st.} & x \geq b \end{array} \quad \Rightarrow \quad \begin{array}{ll} \min_x \max_{\alpha} & L(x, \alpha) \\ \text{st.} & \alpha \geq 0 \\ & L(x, \alpha) = x^2 - \alpha(x - b) \end{array}$$

$$x < b \quad \Rightarrow \quad (x - b) < 0 \quad \Rightarrow \quad \max_{\alpha} -\alpha(x - b) = \inf$$

$$x > b \quad \Rightarrow \quad (x - b) > 0 \quad \Rightarrow \quad \max_{\alpha} -\alpha(x - b) = 0, \alpha = 0 \quad \Rightarrow \quad L(x, \alpha) = x^2$$

$$x = b \quad \Rightarrow \quad (x - b) = 0 \quad \Rightarrow \quad \alpha = \text{any thing} \quad \Rightarrow \quad L(x, \alpha) = x^2$$

Outline

Recap

Kernel Methods & Feature Mapping

Support Vector Machines (SVMs)

Math: Subgradients

Math: Lagrange Multipliers for Constrained Optimization

SVMs Revisited: Dual Formulation

Dual formulation for SVM

SVM problem:

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$\text{subj. to } y_n (w \cdot x_n + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

Dual formulation for SVM

SVM problem:

$$\min_{\boldsymbol{w}, b, \boldsymbol{\xi}} \frac{1}{2} \|\boldsymbol{w}\|^2 + C \sum_n \xi_n$$

$$\text{subj. to } y_n (\boldsymbol{w} \cdot \boldsymbol{x}_n + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

With Lagrange multipliers:

$$\begin{aligned} \mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = & \frac{1}{2} \|\boldsymbol{w}\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (\boldsymbol{w} \cdot \boldsymbol{x}_n + b) - 1 + \xi_n] \end{aligned}$$

Dual formulation for SVM

SVM problem:

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$\text{subj. to } y_n (w \cdot x_n + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

With Lagrange multipliers:

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

1. minimize wrt. w :

$$\nabla_w \mathcal{L} = w - \sum_n \alpha_n y_n x_n = 0 \quad \iff \quad w = \sum_n \alpha_n y_n x_n$$

Is this form familiar?

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n]$$

2. substitute w :

$$\mathcal{L}(b, \xi, \alpha, \beta) = \frac{1}{2} \left\| \sum_m \alpha_m y_m x_m \right\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n - \sum_n \alpha_n \left[y_n \left(\left[\sum_m \alpha_m y_m x_m \right] \cdot x_n + b \right) - 1 + \xi_n \right]$$

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n]$$

2. substitute w :

$$\mathcal{L}(b, \xi, \alpha, \beta) = \frac{1}{2} \left\| \sum_m \alpha_m y_m x_m \right\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n$$

(feature) expand
this out further

$$- \sum_n \alpha_n \left[y_n \left(\left[\sum_m \alpha_m y_m x_m \right] \cdot x_n + b \right) - 1 + \xi_n \right]$$

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

3. minimize wrt **b**:

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_n \alpha_n y_n = 0$$

Dual formulation for SVM

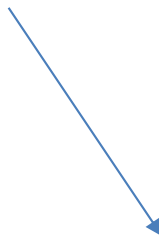
$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

3. minimize wrt **b**:

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_n \alpha_n y_n = 0$$

at optima, this
term is zero
(eliminated)



Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

4. minimize wrt ξ_n :

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = C - \beta_n - \alpha_n \iff C - \beta_n = \alpha_n$$

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

5. do algebra & rewrite, with expansion from **w**

$$\begin{aligned} & -\frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m x_n \cdot x_m + \sum_n (C - \beta_n) \xi_n \\ & - b \sum_n \alpha_n y_n - \sum_n \alpha_n (\xi_n - 1) \end{aligned}$$

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$

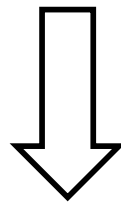
6. rewrite with kernel

$$\mathcal{L}(\alpha) = \sum_n \alpha_n - \frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m K(x_n, x_m)$$

Dual formulation for SVM

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha, \beta) = & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \\ & - \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \end{aligned}$$



$$\min_{\alpha} \quad -\mathcal{L}(\alpha) = \frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m K(x_n, x_m) - \sum_n \alpha_n$$

subj. to $0 \leq \alpha_n \leq C$