# Support Vector Machines: Consequences of Lagrangian Duality

David S. Rosenberg

New York University

February 13, 2018

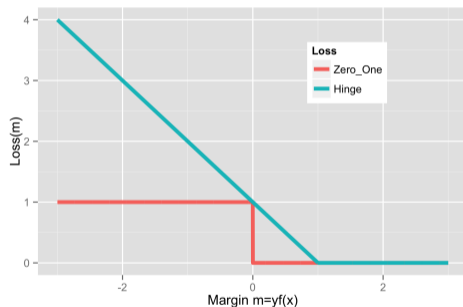# Contents

# The SVM as a Quadratic Program

# The Margin

### Definition

The **margin** (or **functional margin**) for predicted score $\hat{y}$ and true class $y \in \{-1, 1\}$ is $y\hat{y}$.

- The margin often looks like $yf(x)$, where $f(x)$ is our score function.
- The margin is a measure of how **correct** we are.

- We want to **maximize the margin**.

- Most classification losses depend only on the margin.

# Hinge Loss

- SVM/Hinge loss: $\ell_{\text{Hinge}} = \max\{1 - m, 0\}$
- Margin $m = yf(x)$



Hinge is a **convex**, **upper bound** on $0-1$ loss. Not differentiable at $m = 1$.
We have a **"margin error"** when $m < 1$.

# Support Vector Machine

- Hypothesis space $\mathcal{F} = \left\{ f(x) = w^T x + b \mid w \in \mathbf{R}^d,\, b \in \mathbf{R} \right\}$.

- $\ell_2$ regularization (Tikhonov style)
- Loss $\ell(m) = \max\{1 - m, 0\}$

- The SVM prediction function is the solution to

$$\min_{w \in \mathbf{R}^d, b \in \mathbf{R}} \frac{1}{2}\|w\|^2 + \frac{c}{n} \sum_{i=1}^{n} \max\left(0, 1 - y_i \left[w^T x_i + b\right]\right).$$

- (In SVMs it's common to put the regularization parameter $c$ on the empirical risk part, rather than on the $\ell^2$ penalty part.)

# SVM Optimization Problem (Tikhonov Version)

The SVM prediction function is the solution to

$$\min_{w\in\mathbf{R}^d, b\in\mathbf{R}} \frac{1}{2}\|w\|^2 + \frac{c}{n}\sum_{i=1}^{n} \max\left(0, 1 - y_i\left[w^T x_i + b\right]\right).$$

- unconstrained optimization
- **not differentiable** because of the max (right at the border of a margin error)
- Can we reformulate into a differentiable problem?

# SVM Optimization Problem

- The SVM optimization problem is equivalent to

$$\text{minimize} \quad \frac{1}{2}\|w\|^2 + \frac{c}{n}\sum_{i=1}^{n}\xi_i$$

$$\text{subject to} \quad \xi_i \geqslant \max\left(0, 1 - y_i\left[w^T x_i + b\right]\right).$$

- Which is equivalent to

$$\text{minimize} \quad \frac{1}{2}\|w\|^2 + \frac{c}{n}\sum_{i=1}^{n}\xi_i$$

$$\text{subject to} \quad \xi_i \geqslant \left(1 - y_i\left[w^T x_i + b\right]\right) \text{ for } i = 1, \ldots, n$$

$$\xi_i \geqslant 0 \text{ for } i = 1, \ldots, n$$

## SVM as a Quadratic Program

- The SVM optimization problem is equivalent to

$$\text{minimize} \quad \frac{1}{2}\|w\|^2 + \frac{c}{n}\sum_{i=1}^{n}\xi_i$$

$$\text{subject to} \quad -\xi_i \leqslant 0 \text{ for } i = 1, \ldots, n$$

$$\left(1 - y_i\left[w^T x_i + b\right]\right) - \xi_i \leqslant 0 \text{ for } i = 1, \ldots, n$$

- Differentiable objective function
- $n + d + 1$ unknowns and $2n$ affine constraints.
- A quadratic program that can be solved by any off-the-shelf QP solver.
- Let's learn more by examining the dual.

# Lagrangian Duality for SVM

# The SVM Dual Problem

- Following recipe and with some algebra, the SVM dual problem is equivalent to:

$$\sup_{\alpha} \qquad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j x_j^T x_i$$

$$\text{s.t.} \qquad \sum_{i=1}^{n} \alpha_i y_i = 0$$

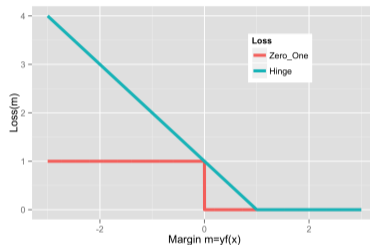$$\alpha_i \in \left[ 0, \frac{c}{n} \right] \ i = 1, \ldots, n.$$

- Let $\alpha^*$ be solution to this optimization problem (the **dual optimal point**).
- Can show that the SVM solution is

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

- $w^*$ is "in the **span of the data**" – i.e. a linear combination of $x_1, \ldots, x_n$.

# The Margin and Some Terminology

- For notational convenience, define $f^*(x) = x^T w^* + b^*$.
- Margin $yf^*(x)$



- Incorrect classification: $yf^*(x) \leqslant 0$.
- Margin error: $yf^*(x) < 1$.
- "On the margin": $yf^*(x) = 1$.
- "Good side of the margin": $yf^*(x) > 1$.

## Complementary Slackness Results: Summary

- SVM optimal parameter is $w^* = \sum_{i=1}^{n} \alpha_i^* y_i x_i$.
- We can derive the following relations from complementary slackness conditions:

$$\alpha_i^* = 0 \implies y_i f^*(x_i) \geqslant 1$$

$$\alpha_i^* \in \left(0, \frac{c}{n}\right) \implies y_i f^*(x_i) = 1$$

$$\alpha_i^* = \frac{c}{n} \implies y_i f^*(x_i) \leqslant 1$$

$$y_i f^*(x_i) < 1 \implies \alpha_i^* = \frac{c}{n}$$

$$y_i f^*(x_i) = 1 \implies \alpha_i^* \in \left[0, \frac{c}{n}\right]$$

$$y_i f^*(x_i) > 1 \implies \alpha_i^* = 0$$

## Support Vectors

- If $\alpha^*$ is a solution to the dual problem, then primal solution is

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

  with $\alpha_i^* \in [0, \frac{c}{n}]$.
- The $x_i$'s corresponding to $\alpha_i^* > 0$ are called **support vectors**.
- Few margin errors or "on the margin" examples $\implies$ **sparsity in input examples**.
- This becomes important when we get to **kernelized SVMs**.

# Teaser for Kernelization

## Dual Problem: Dependence on $x$ through inner products

- SVM Dual Problem:

$$\sup_{\alpha} \qquad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j x_j^T x_i$$

$$\text{s.t.} \qquad \sum_{i=1}^{n} \alpha_i y_i = 0$$

$$\alpha_i \in \left[0, \frac{c}{n}\right] \ i = 1, \ldots, n.$$

- Note that all dependence on inputs $x_i$ and $x_j$ is through their inner product: $\langle x_j, x_i \rangle = x_j^T x_i$.
- We can replace $x_j^T x_i$ by any other inner product...
- This is a "kernelized" objective function.