

SVM

Transform Labels:

$$\{0, 1\} \rightarrow \{-1, 1\}$$

$$\tilde{y}_n = y_n$$

$$y_n = 2\tilde{y}_n - 1$$

$$\tilde{y}_n = \frac{1}{2}(y_n + 1)$$

Classification Losses:

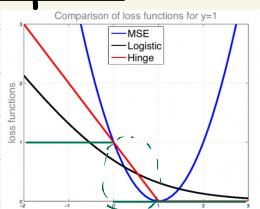
$$[1 - 3y]_+ = \max\{0, 1 - 3y\}$$

↳ Hinge loss

$$\text{MSE}(z, y) = (1 - yz)^2$$

$$\text{Logistic Loss}(z, y) = \log(1 + e^{-yz})$$

Classification:



- As we get on the "right side" the hinge loss is 0.

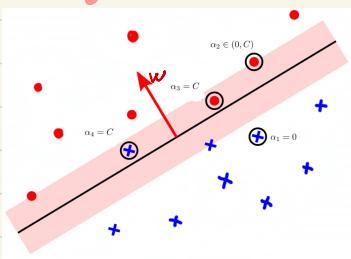
- The cost of hinge loss increases linearly as we get away on the "wrong side".

Support Vector Machines

SVMs correspond to the loss: (Hinge + Regular)

$$\min_w \sum_{n=1}^N [1 - y_n x_n^T w] + \frac{\lambda}{2} \|w\|_2^2$$

The region in pink is called The margin:



Take the normal vector w that defines the hyperplane,

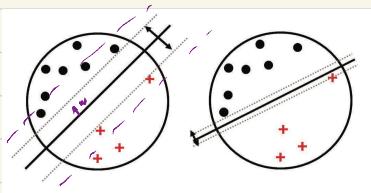
Look at the feature vector x so that $|x^T w| \leq 1$

→ This is the margin: it doesn't only depend on the direction of w , but also on its norm: Total width = $\frac{2}{\|w\|}$

. The goal is: 1. find a separating hyperplane

(Hinge) ↳ 2. A scaling of w so no point in the margin

. We get enough of no point in margin (norm is enough) + (Reg) 3. that the separating hyperplane and scaling of w for which the margin is the largest.



Optimization:

- How to get w ?

$$\min_w \sum_{n=1}^N [1 - y_n x_n^T w]_+ + \frac{\lambda}{2} \|w\|_2^2$$

→ convex problem!

(one small problem at t, ⇒ use subgradient → then GD or SGD)

- There is another way to look at it:

Convex duality:

$$\min_w \mathcal{L}(w) = \min_w \max_{\alpha} G(\alpha, w)$$

Primal problem: $\min_w \max_{\alpha} G(\alpha, w)$

⇒ sometimes the dual p. is simpler to solve

Dual problem: $\max_{\alpha} \min_w G(\alpha, w)$

$$1. [z]_+ = \max\{0, z\}$$

$$= \max_{\alpha \in [0, 1]} \alpha z$$

$$\mathcal{L}(w) = \frac{1}{N} \sum_{n=1}^N \max_{\alpha \in [0, 1]} \alpha_n (1 - y_n x_n^T w) + \frac{\lambda}{2} \|w\|_2^2$$

$$\Rightarrow \min_w \mathcal{L}(w) = \min_w \max_{\alpha \in [0, 1]} \frac{1}{N} \sum_{n=1}^N \alpha_n (1 - y_n x_n^T w) + \frac{\lambda}{2} \|w\|_2^2$$

$$\Rightarrow G(d, w) = \sum_{n=1}^N \alpha_n (1 - y_n x_n^T w) + \frac{\lambda}{2} \|w\|_2^2$$

G is convex in w and concave in α

(linear in α)

2. We have: $\max_{\alpha} \min_w G(\alpha, w) \leq \min_w \max_{\alpha} G(\alpha, w)$

Equality if G is convex and concave and domain is

(compact 'Bounded') and convex: $\max_{\alpha} \min_w G(\alpha, w) = \min_w \max_{\alpha} G(\alpha, w)$

$$\Rightarrow \min_w \mathcal{L}(w) = \max_{\alpha \in [0, 1]} \min_w \frac{1}{N} \sum_{n=1}^N \alpha_n (1 - y_n x_n^T w) + \frac{\lambda}{2} \|w\|_2^2$$

Compute the min:

$$\nabla_w G(\alpha, w) = \sum_{i=1}^N (-\alpha_i y_i x_i) + \lambda w = 0$$

$$\Rightarrow W = \frac{1}{\lambda} \sum_{i=1}^N \alpha_i y_i x_i$$

$$= \frac{1}{\lambda} X^T Y \alpha \quad : Y = \text{diag}(y)$$

$$W(\alpha) = \frac{1}{\lambda} X^T Y \alpha$$

$$\Rightarrow \min_w L(w) = \max_{\alpha \in \{0,1\}^N} \sum_{n=1}^N \alpha_n (1 - y_n x_n^T w(\alpha)) + \frac{\lambda}{2} \|w\|_2^2 \quad (\% \text{ doesn't affect the max})$$

$$= \max_{\alpha \in \{0,1\}^N} \sum_{n=1}^N \alpha_n \left(1 - \frac{1}{\lambda} y_n x_n^T X^T Y \alpha \right) + \frac{\lambda}{2} \frac{1}{\lambda^2} \|X^T Y \alpha\|^2 \quad (\alpha_{N \times 1})$$

$$= \max_{\alpha \in \{0,1\}^N} \sum_{n=1}^N \alpha_n - \frac{1}{\lambda} \sum_{n=1}^N \underbrace{\alpha_n y_n x_n^T}_{\alpha^T X^T X} X^T Y \alpha + \frac{1}{2\lambda} \|X^T Y \alpha\|^2$$

$$= \max_{\alpha \in \{0,1\}^N} \alpha^T \mathbf{1} - \frac{1}{\lambda} \|X^T Y \alpha\|_2^2 + \frac{1}{2\lambda} \|X^T Y \alpha\|_2^2$$

$$= \max_{\alpha \in \{0,1\}^N} \alpha^T \mathbf{1} - \frac{1}{2\lambda} \|X^T Y \alpha\|_2^2$$

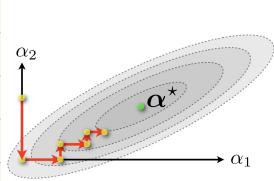
$$\boxed{\min_w L(w) = \max_{\alpha \in \{0,1\}^N} \alpha^T \mathbf{1} - \frac{1}{2\lambda} \underbrace{\alpha^T Y X^T X^T Y \alpha}_{\text{PSD matrix}}}$$

$$\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1}$$

3. Why is this better?

$$\min_w \mathcal{L}(w) = \max_{\alpha \in [0,1]^n} \alpha^T \cdot \mathbf{1} - \frac{1}{2\lambda} \alpha^T y X X^T y \alpha$$

Coordinate ascent (descent)



- 1) Non-diff prob. can be slow to solve
 → now we have a differentiable, concave problem:
 . Can be solved using quadra. prog. tools
 . useful to use coordinate descent
 1 variable

- 2) The cost function is only depending on the data in the form:
 $K = X X^T_{(N \times N)}$ (does not depend on D!) **Kernel**
- 3) It gives a good interpretation of SVM

- α is sparse, for any x_i , we have α_i given by:

$$\max_{\alpha_i \in \{0,1\}} \alpha_i (1 - y_i x_i^T w)$$
- if x_i is correct and outside the margin:
 $1 - y_i x_i^T w < 0 \Rightarrow \alpha_i = 0$
 → non-support point
- correct point on the margin:
 $1 - y_i x_i^T w = 0 \Rightarrow \alpha_i \in \{0,1\} \rightarrow$ essential support vector
- Points that are strictly inside the margin: $1 - y_i x_i^T w > 0 \Rightarrow \alpha_i = 1$
 or on the wrong side
 → bound support vector

$$w = \frac{1}{\lambda} \sum_{i=1}^n x_i y_i \alpha_i ; \text{ linear comb. of } \alpha_i ; \text{ depends on support points}$$

Issues with SVM:

- Difficult to extend to multiclass / no obvious probabilistic interpretation