

Support Vector Machines

Non-probabilistic models for classification of the form:

$$\underline{f(\mathbf{x})} = \sum_{\underline{i=1}}^N \alpha_i y_i \underline{K(\mathbf{x}, \mathbf{x}^i)} + w_0$$

Learning fixes values of α_i . These are chosen in such a way that only a few i s have non-zero α_i . These are called support vectors.

Assume $y_i \in \{-1, +1\}$, binary classification model.

Predicted class label = $\text{sign}(f(\mathbf{x}))$

When Kernel is linear

If kernel is linear: $K(\mathbf{x}, \mathbf{x}^i) = \mathbf{x} \cdot \mathbf{x}^i$

$$K(\mathbf{x}, \mathbf{x}) =$$

Then the SVM classifier $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}, \mathbf{x}^i) + w_0$ reduces to $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$ where $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$

$$\begin{aligned} f(\mathbf{x}) &= \sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}^i) + w_0 \\ &= \sum_i \alpha_i y_i \mathbf{x} \cdot \mathbf{x}^i + w_0 \end{aligned}$$

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$$

$K=2$; assume class labels are $\{+1, -1\}$

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$$

$$\begin{aligned} \mathbf{x}^i &\in \mathbb{R}^d \\ \alpha_i &\in \mathbb{R}, y_i \in \{+1, -1\} \end{aligned}$$

Binary logistic regression

$$\mathbf{w} = [+1.3, -1.7]$$

when $k=2$ then
for each class: instead of separate parameters we create $\mathbf{w}^1 = \mathbf{w}$; $\mathbf{w}^2 = -\mathbf{w}$

α_i	SVM D	$d=2$
0	$\mathbf{x}^1: (1, 2), +1$	
0.4	$\mathbf{x}^2: (0.5, -0.5), -1$	
	\mathbf{x}^3	
0.3	$\mathbf{x}^N: (2, -1), +1$	

Training SVMs

$$D \equiv \{ (x_i, y_i) : i = 1 \dots N \}$$

Given a training dataset, find α_i s and w_0 so as to minimize a new margin loss called the Hinge loss.

For correct prediction: $yf(x) > 0$ for a test example (x, y)

Hinge loss: $L_H(y, f(x)) = \max(0, 1 - yf(x))$

Loss is 0 iff $yf(x) \geq 1$.

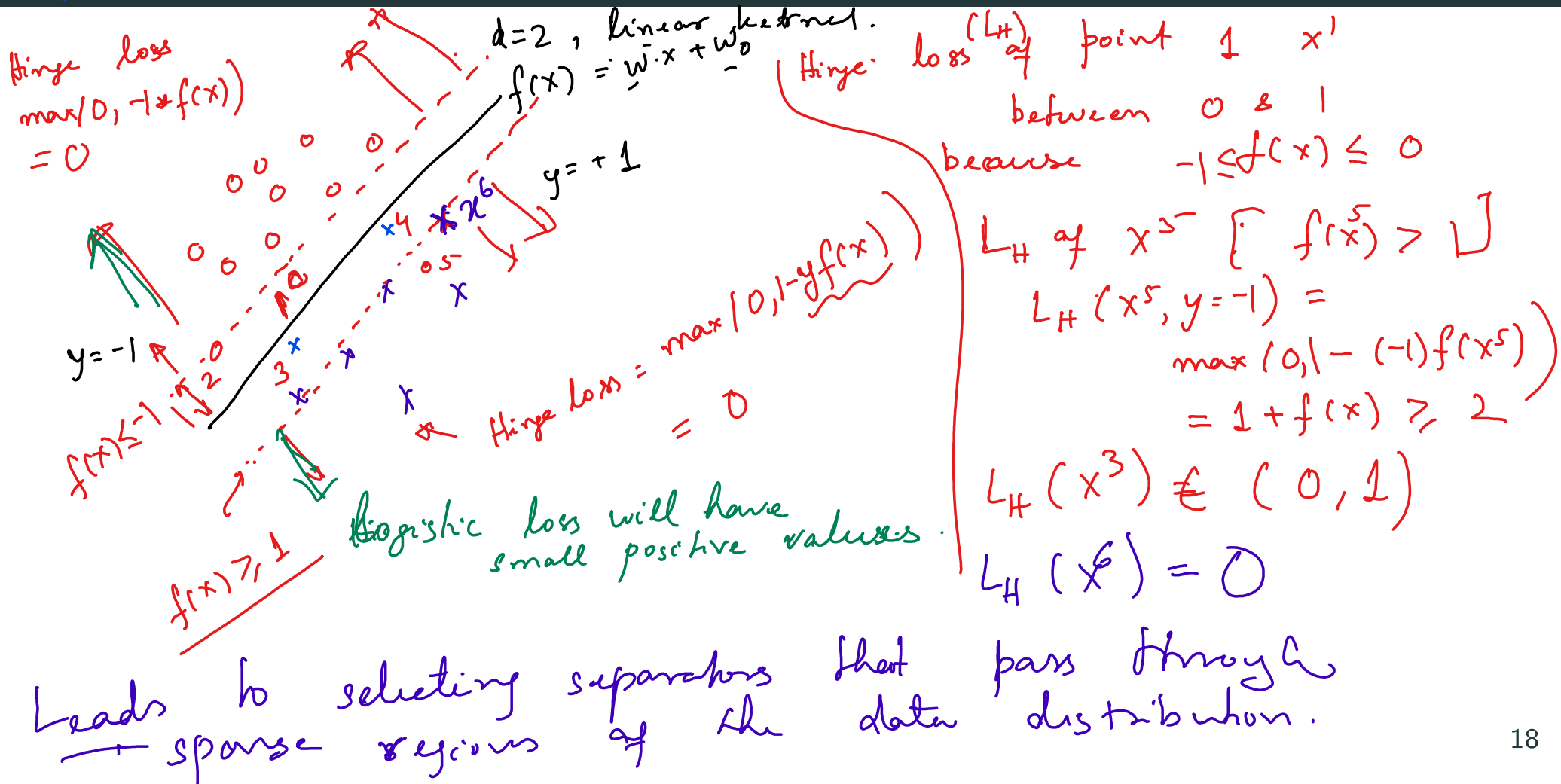
To minimize 0/1 prediction error we want to minimize # of examples for which $yf(x) < 0$

minimize $\sum_{i=1}^N -y_i f(x_i)$

unbounded training loss is zero.

perception loss: $\min_f \sum_{i=1}^N \max(0, -y_i f(x_i))$

Understanding Hinge Loss geometrically



Comparing Hinge loss and logistic loss.

Advantage of hinge-loss: sparsity.

Hinge loss: $L_H = \max(0, 1 - yf(x))$

Perceptron loss: $\max(0, -yf(x))$
is inferior to Hinge loss
because of not imposing
margin

Logistic loss: loss used in logistic
regression classifier:
(for $k=2$) $L_e = \log(1 + e^{-yf(x)})$

If $yf(x) > 1$ then $L_H = 0$

$$L_H = 1 - yf(x)$$

Logistic else

If $yf(x) = 0$ then

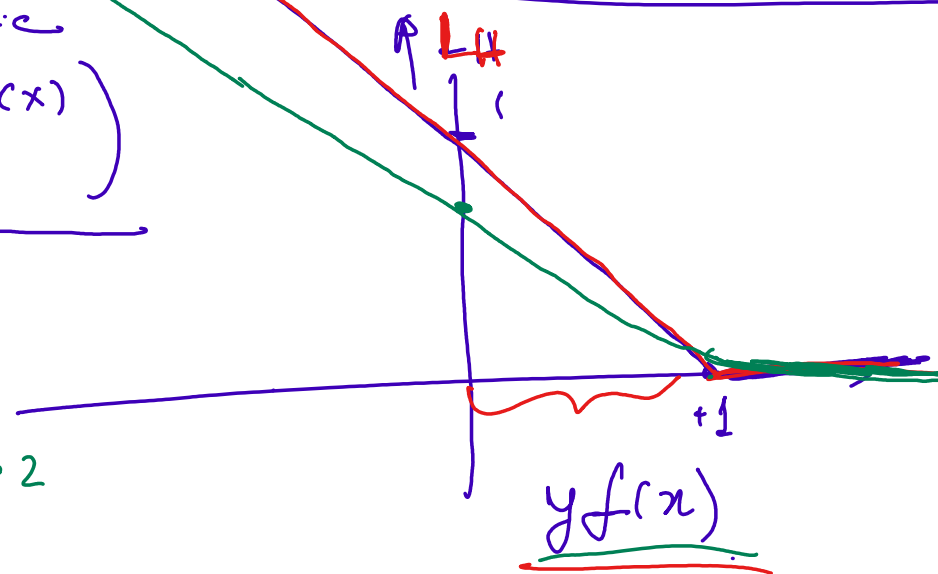
$$L_e = \log(1 + 1) = \log 2$$

If $yf(x) \rightarrow \infty$

$$L_e \Rightarrow 0$$

If $yf(x) \gg -1$

$$L_e \approx -yf(x)$$



Comparing Hinge loss and logistic loss.

Advantage of hinge-loss: sparsity.

Training SVMs

First consider the linear case:

Training objective:

$$\min_{\mathbf{w}, w_0} C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}^i + w_0)) + \|\mathbf{w}\|^2$$

hyper-parameter

Hinge loss on training data with $f(\mathbf{x}^i) = \mathbf{w} \cdot \mathbf{x}^i + w_0$

regularizer to further prevent over-fitting.

Above loss is convex in \mathbf{w} , and w_0 .

1) $\max(1 - yf(\mathbf{x}))$ is convex — In general if $L(\mathbf{w}) = g(\mathbf{z})$ and $f(\mathbf{x})$ is linear in \mathbf{w} .

2) $\|\mathbf{w}\|^2$ is convex in \mathbf{w}

3) Sum of convex functions is convex.

where \mathbf{z} is a linear function of \mathbf{w} and $\mathbf{z} = \mathbf{w}^T \mathbf{x}$ and $g(\mathbf{z})$ is convex. The $L(\mathbf{w})$ is convex.

Training objective as a constrained program.

Equivalent to: But loss is not differentiable in w .
 Can use stochastic sub-gradient descent algorithm for training w, w_0

$$\begin{aligned}
 & \min_{\mathbf{w}, w_0, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \\
 & \text{s.t. } y^i(\mathbf{w} \cdot \mathbf{x}^i + w_0) \geq 1 - \xi_i \quad \forall i: 1 \dots N \\
 & \quad \xi_i \geq 0 \quad i: 1 \dots N
 \end{aligned} \tag{2}$$

$\max(0, 1 - y^i(\mathbf{w} \cdot \mathbf{x}^i + w_0)) = \xi_i$
 $\xi_i \geq 1 - y^i(\mathbf{w} \cdot \mathbf{x}^i + w_0)$
 $\xi_i \geq 0$

Show

Training objective in dual form

$$\max_{\alpha_1, \dots, \alpha_N} \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \left(\sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}^i \cdot \mathbf{x}^j \right) \right\}$$

quadratic in $\alpha_1, \dots, \alpha_N$
concave in $\alpha_1, \dots, \alpha_N$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 1 \quad \text{\textit{\& linear constraint}}$$

$$0 \leq \alpha_i \leq C, \quad i: 1 \dots N \quad \text{\textit{[box constraints]}}$$

(3)

Training algorithm for SVMs

- ⇒ Objective is concave in α -s
- ⇒ constraints are easy: box-constraint, balancing constraint.
- ⇒ Can solve with projected gradient ascent algorithm.
- ⇒ one popular method is : SMO algorithm:

$$\alpha_1^0 \dots \alpha_N^0 \leftarrow [0 \dots 0] \quad t=0$$

while no change: ($t = 1 \dots \infty$)

→ pick two α_i^t, α_j^t

find change Δ s.t.:

$$\alpha_i^{t+1} \leftarrow \alpha_i^t - \Delta$$

$$\alpha_j^{t+1} \leftarrow \alpha_j^t + \Delta \cdot \mathcal{S}(y_i, y_j)$$

→ ensures that $\sum \alpha_i^{t+1} y_i = 0$

+1 if $y_i \neq y_j$

-1 if $y_i = y_j$

Δ so as to cause maximum ascent in objective.
and $0 \leq \alpha_i^{t+1} \leq C$ $0 \leq \alpha_j^{t+1} \leq C$

Kernel Trick: extending to non-linear kernels

In the dual objective replace $\mathbf{x}^i \cdot \mathbf{x}^j$ with $K(\mathbf{x}^i, \mathbf{x}^j)$.

Same trick in $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$ where $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}^i$...

Dual objective for linear kernels.

$$\max_{\alpha_1, \alpha_2, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \underbrace{[\mathbf{x}^i \cdot \mathbf{x}^j]}_{\text{replace this with } k(\mathbf{x}^i, \mathbf{x}^j)}$$

$$s.t. \quad 0 \leq \alpha_i \leq C \quad \forall i = 1 \dots N$$

$$\sum_{i=1}^N \alpha_i y_i = 0 \Leftrightarrow \sum_{i: y_i = +1} \alpha_i = \sum_{i: y_i = -1} \alpha_i$$

Deployment!

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0 \Leftrightarrow \underbrace{\left[\sum_{i=1}^N \alpha_i y_i \mathbf{x}^i \right]}_{\mathbf{w} \text{ for linear kernel}} \cdot \mathbf{x} + w_0 \xrightarrow{\text{apply kernel trick}}$$

$$\sum_{i=1}^N \alpha_i y_i k(\mathbf{x}^i, \mathbf{x}) + w_0$$

Demo

Training w_0

1) Pick any example j s.t. ~~$\alpha_j = C$~~ $\alpha_j < C$

Now, solve for w_0 s.t. $y_j f(x^j) = 1$

$$f(x) = \sum_{i=1}^N \alpha_i y_i k(x^i, x) + w_0$$

<https://drive.google.com/file/d/1LihEFle5fyFbWPS342PA9NZdy2Kxjkvy/view?usp=sharing>

Summary

- Kernel methods provide a new way to represent data and build models based purely on comparison of two points.
- Various kernels are available, each with its own characteristics.
- Kernel methods find applications in several ML algorithms: classifiers (Nearest neighbor, few-shot classification, SVMs) and regression (Kernel regression, support vector regression, Gaussian processes)