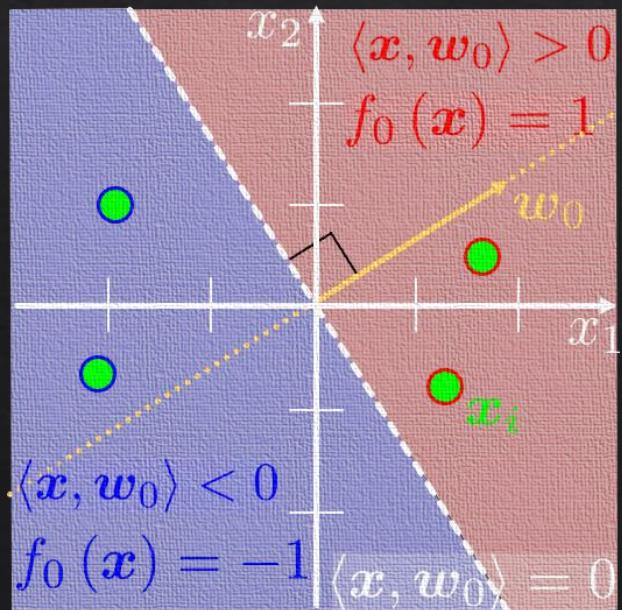


Supervised Learning – Discriminative Classification

Machine Learning Methods – Lecture 3



June 2021



Introduction and Notations – Binary Classification

- Consider a training set:

$$\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

where:

- $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$ is the feature vector (a.k.a. sample or observation).
- $y_i \in \mathcal{Y} = \{C_0, C_1\}$ is the class of \mathbf{x}_i (a.k.a. category, label, response).

- Let (\mathbf{x}_0, y_0) be a new pair, where \mathbf{x}_0 is known and y_0 is unknown.

- Using $\mathcal{D}_{\text{train}}$ our goal is to derive a classifier (a mapping):

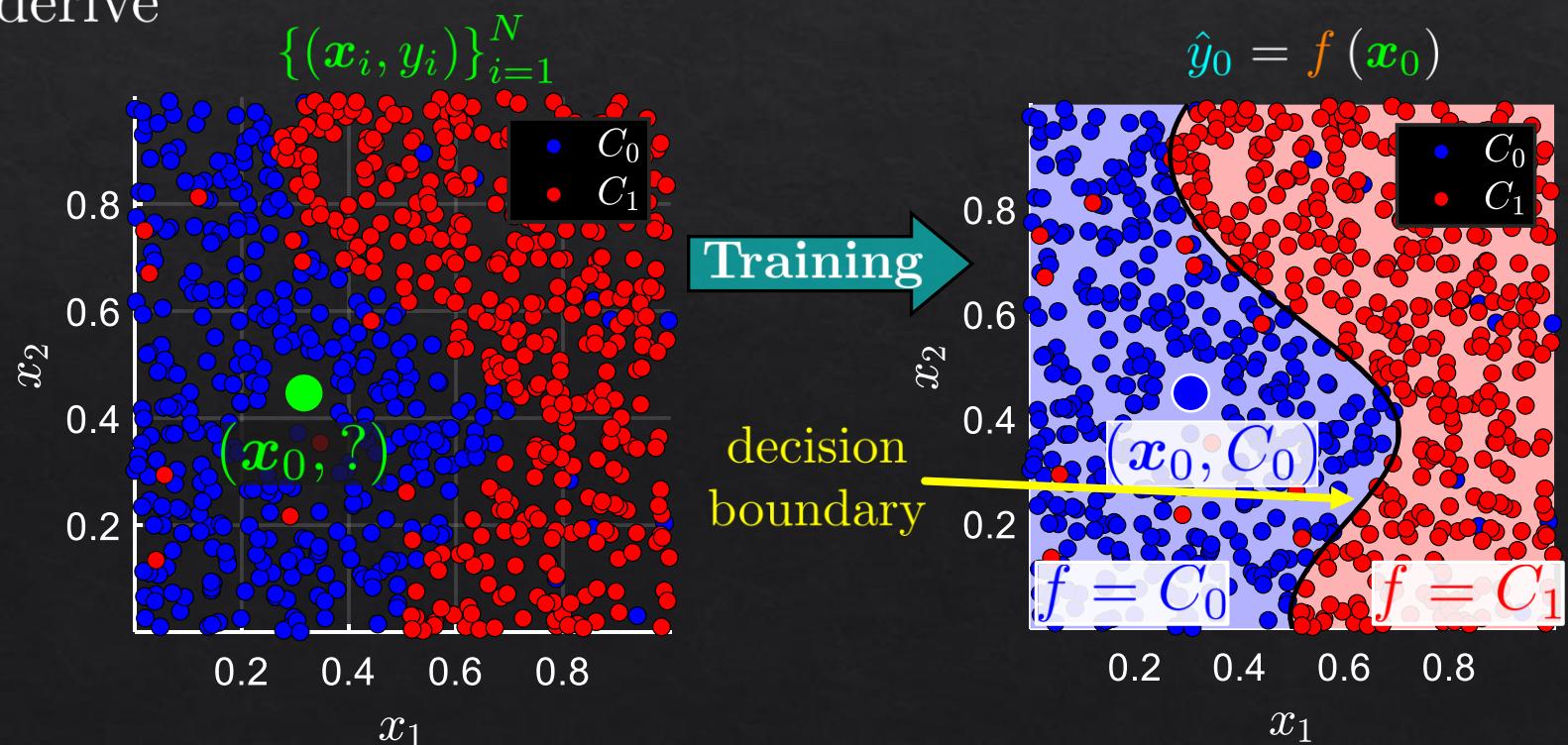
$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

such that: $y_0 = f(\mathbf{x}_0)$

- A set of new* pairs is known as a test set:

$$\mathcal{D}_{\text{test}} = \{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^M$$

(* new – unseen during training.

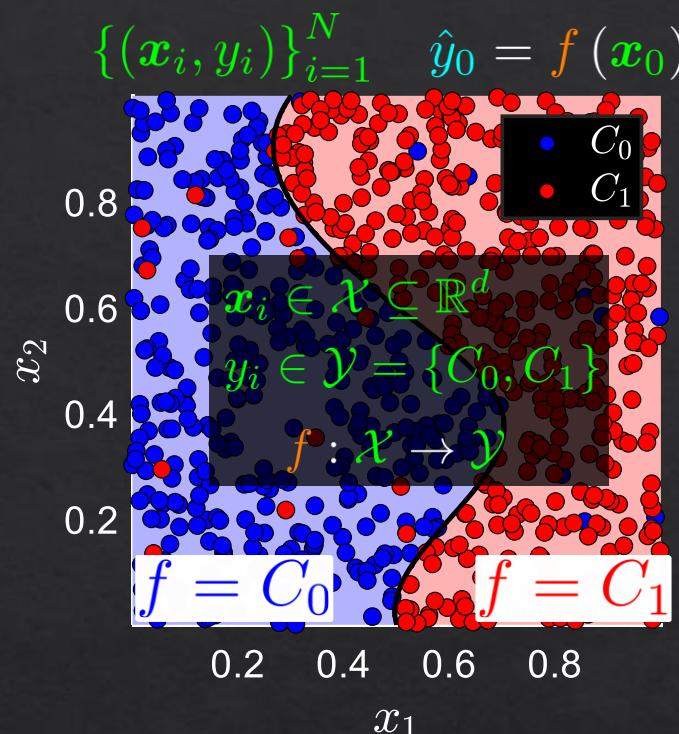


Hamming Loss

- A loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ satisfies:

1. $\hat{y}_i = y_i \implies \ell(\hat{y}_i, y_i) = 0$
2. $\ell(\hat{y}_i, y_i) \geq 0$, for all $\hat{y}_i, y_i \in \mathcal{Y}$

- Hamming 0 – 1 loss: $\ell(\hat{y}_i, y_i) = \mathbb{I}\{\hat{y}_i \neq y_i\} = \begin{cases} 0 & \hat{y}_i = y_i \\ 1 & \hat{y}_i \neq y_i \end{cases}$



- The train error of a classifier f w.r.t. $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$: $\hat{y}_i = f(\mathbf{x}_i)$

$$L_{\text{train}}(f) := \frac{1}{N} \sum_{i=1}^N \ell(\hat{y}_i, y_i) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{y}_i \neq y_i\}$$

$1 - L(f)$ is the
classifier accuracy

$L(f) = 0 \implies$ no errors

- The goal is to minimize $L_{\text{test}}(f)$ w.r.t $\mathcal{D}_{\text{test}}$.

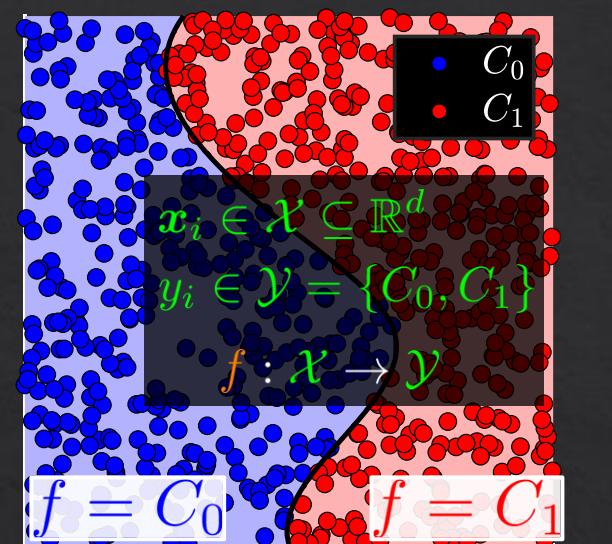
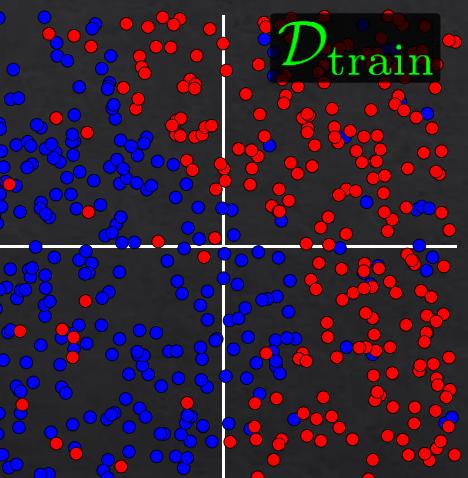
- Since $\mathcal{D}_{\text{test}}$ is unavailable, we **usually** minimize:

$$f^\star = \arg \min_{f \in ?} L_{\text{train}}(f) = \arg \min_{f \in ?} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{f(\mathbf{x}_i) \neq y_i\}$$

Classifier Models

- Generally, it is not feasible to search the optimal function in the space of all functions.
- A common practice is to fix a set of functions \mathcal{H} (a model) and restrict the search to that set:

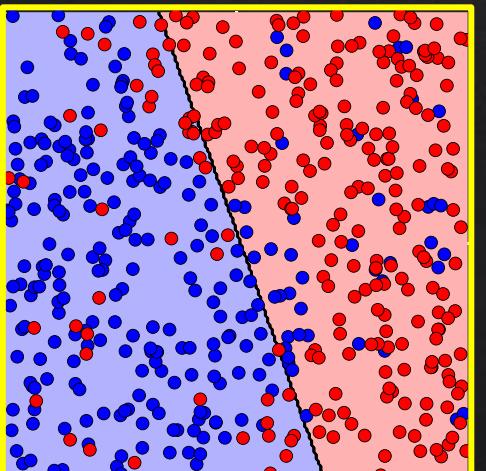
$$\arg \min_{f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\mathbf{f}(\mathbf{x}_i) = y_i\}$$



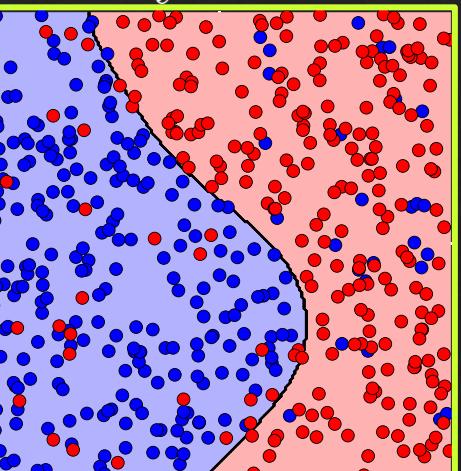
$$\arg \min_{f \in ?} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\mathbf{f}(\mathbf{x}_i) \neq y_i\}$$

Some models for example:

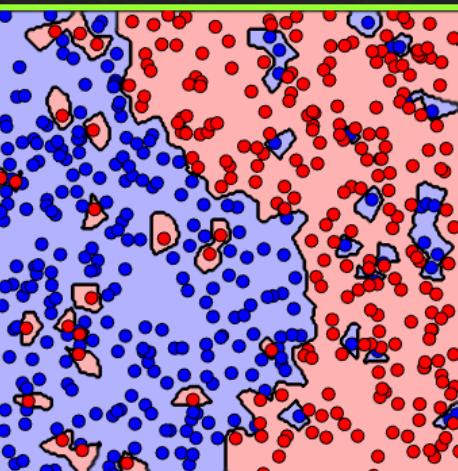
Linear



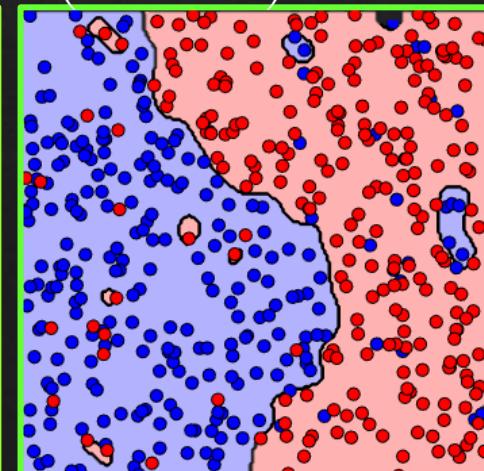
Polynomial



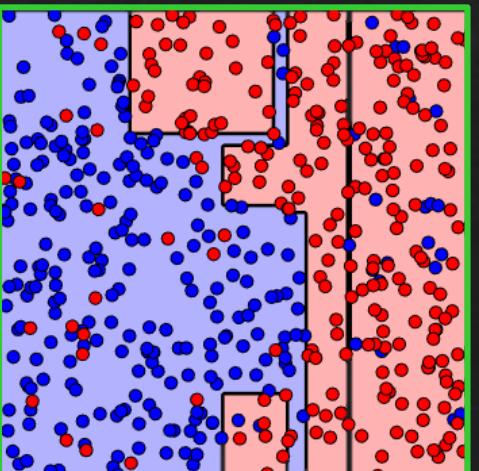
K-NN



(Kernel) SVM



Decision tree



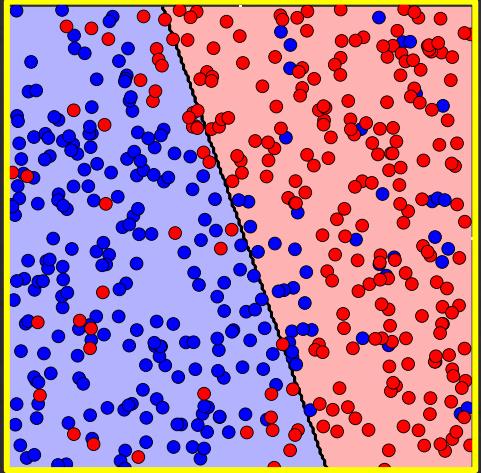
Linear Classifiers

Linear

- We model a linear classifier $f \in \mathcal{H}_{\text{linear}}$ using two parameters:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b), \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

- $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ are the model parameters.



The role of \mathbf{w}

- Consider the linear classifier f_0 : with $b = 0$

$$f_0(\mathbf{x}) = \text{sign}(\mathbf{w}_0^T \mathbf{x}) = \text{sign}(\langle \mathbf{x}, \mathbf{w}_0 \rangle)$$

- The decision boundary is given by:

$$f_0(\mathbf{x}) = 0$$

$$\text{sign}(\langle \mathbf{x}, \mathbf{w}_0 \rangle) = 0$$

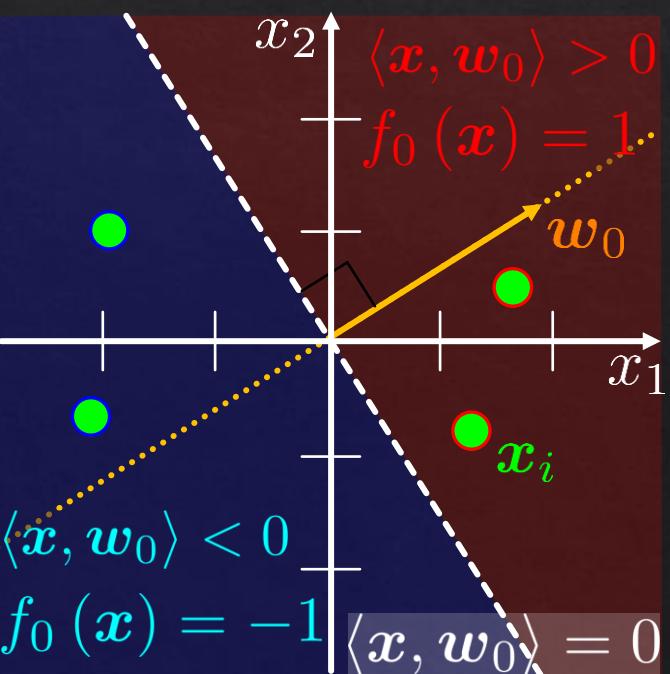
$$\langle \mathbf{x}, \mathbf{w}_0 \rangle = 0$$

what is the role of b ?

$$\arg \min_{f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{f(\mathbf{x}_i) = y_i\}$$

$$\text{sign}(\alpha) = \begin{cases} -1 & \alpha < 0 \\ 0 & \alpha = 0 \\ 1 & \alpha > 0 \end{cases}$$

$$\mathbf{w}^T \mathbf{x} = \langle \mathbf{x}, \mathbf{w} \rangle$$



Linear Classifiers

The role of b

- Consider the classifier f_1 :

$$f_1(\mathbf{x}) = \text{sign}(\mathbf{w}_1^T \mathbf{x} - b_1), \quad \|\mathbf{w}_1\|_2 = 1$$

- The decision boundary is given by:

$$f_1(\mathbf{x}) = 0$$

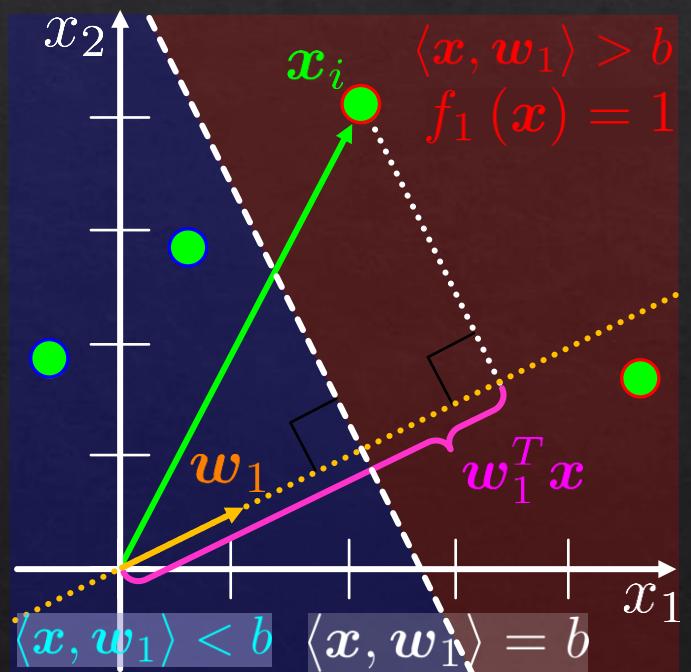
$$\text{sign}(\mathbf{w}_1^T \mathbf{x} - b_1) = 0$$

$$\mathbf{w}_1^T \mathbf{x} = b$$

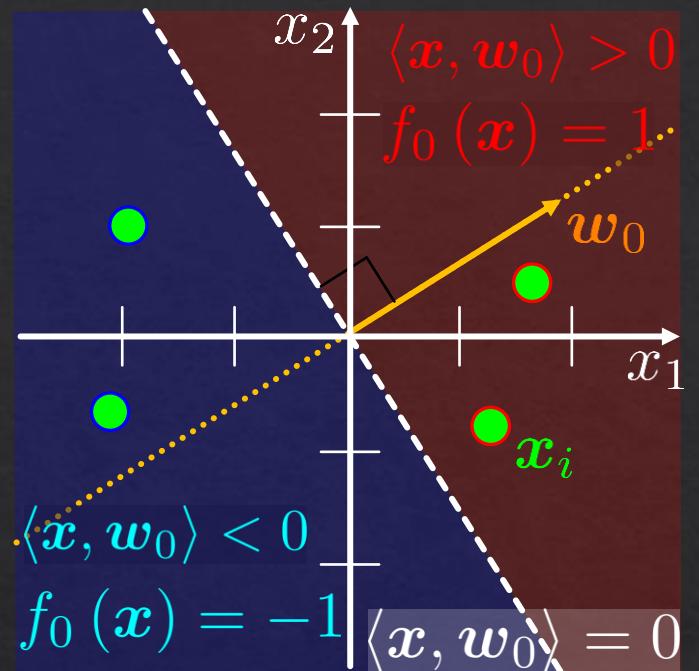
- We search for \mathbf{w} and b which minimize the train error:

$$\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$



$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} L_{\text{train}}(f) = \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{f(\mathbf{x}_i) \neq y_i\}$$



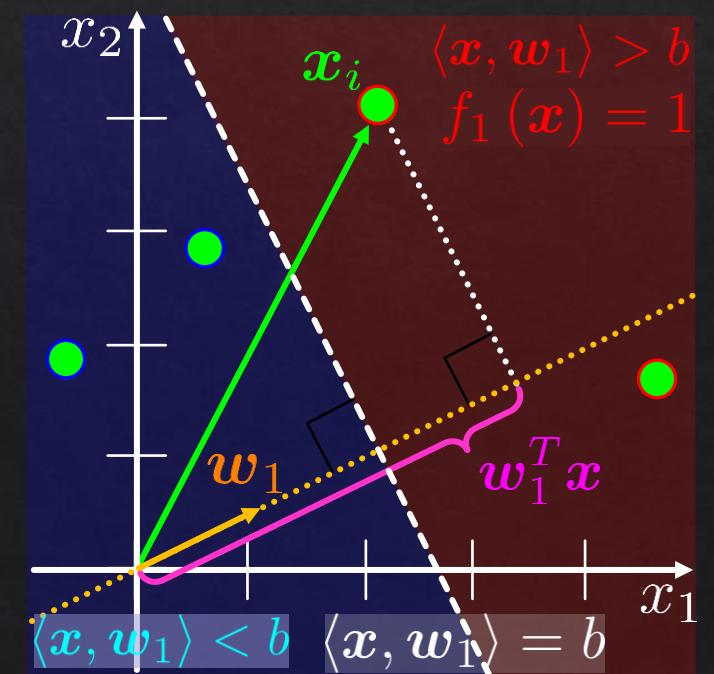
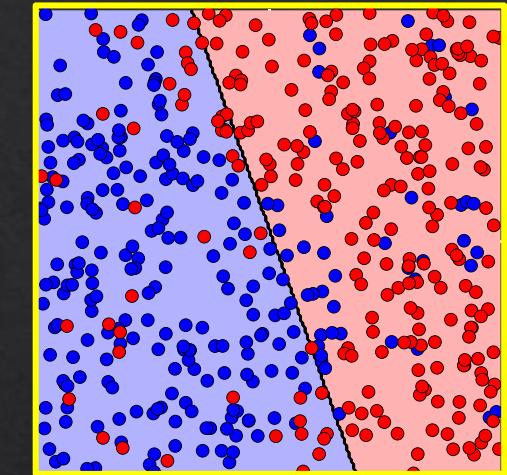
$$\begin{aligned} \arg \min_{f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{f(\mathbf{x}_i) = y_i\} \\ f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b) \end{aligned}$$

Linear Classifier – Example



Linear Classifier

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$



Linear Classifiers – Reparametrization

- Consider a training pair (\mathbf{x}_i, y_i) .
- The Hamming loss is given by:

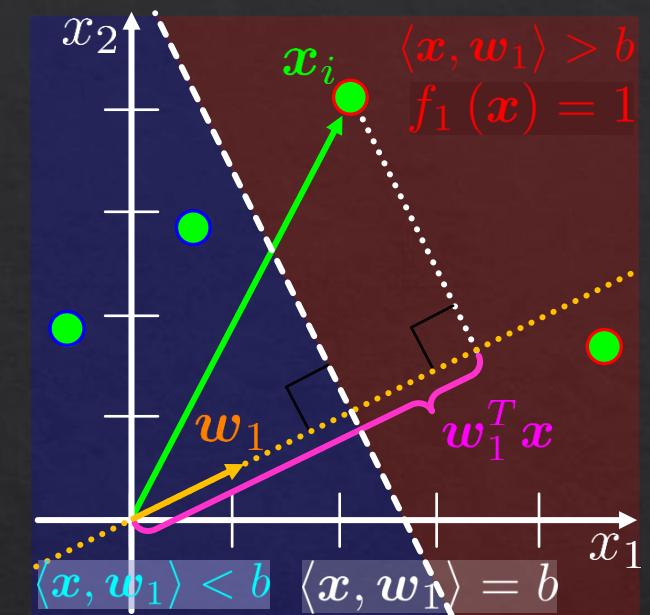
$$\ell_i(f) = \mathbb{I}\{\mathbf{f}(\mathbf{x}_i) \neq y_i\}$$

$$= \mathbb{I}\{\text{sign}(\mathbf{w}^T \mathbf{x}_i - b) \neq y_i\}$$

$$\implies \ell_i(f) = \mathbb{I}\{\text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \neq y_i\}$$

$$\tilde{\mathbf{x}}_i := \begin{bmatrix} -1 \\ | \\ \mathbf{x}_i \\ | \end{bmatrix} \in \mathbb{R}^{d+1}, \quad \tilde{\mathbf{w}} := \begin{bmatrix} b \\ | \\ \mathbf{w} \\ | \end{bmatrix} \in \mathbb{R}^{d+1}$$

$$\implies \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i = \mathbf{w}^T \mathbf{x}_i - b$$



$$\arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\mathbf{f}(\mathbf{x}_i) \neq y_i\}$$

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$

Linear Classifiers – Objective

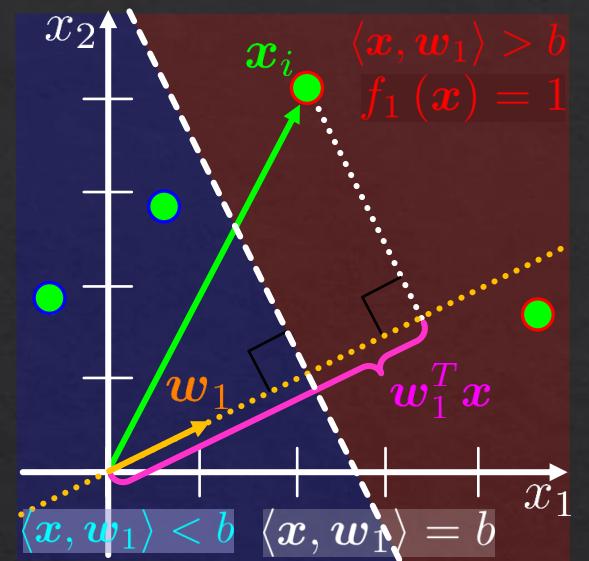
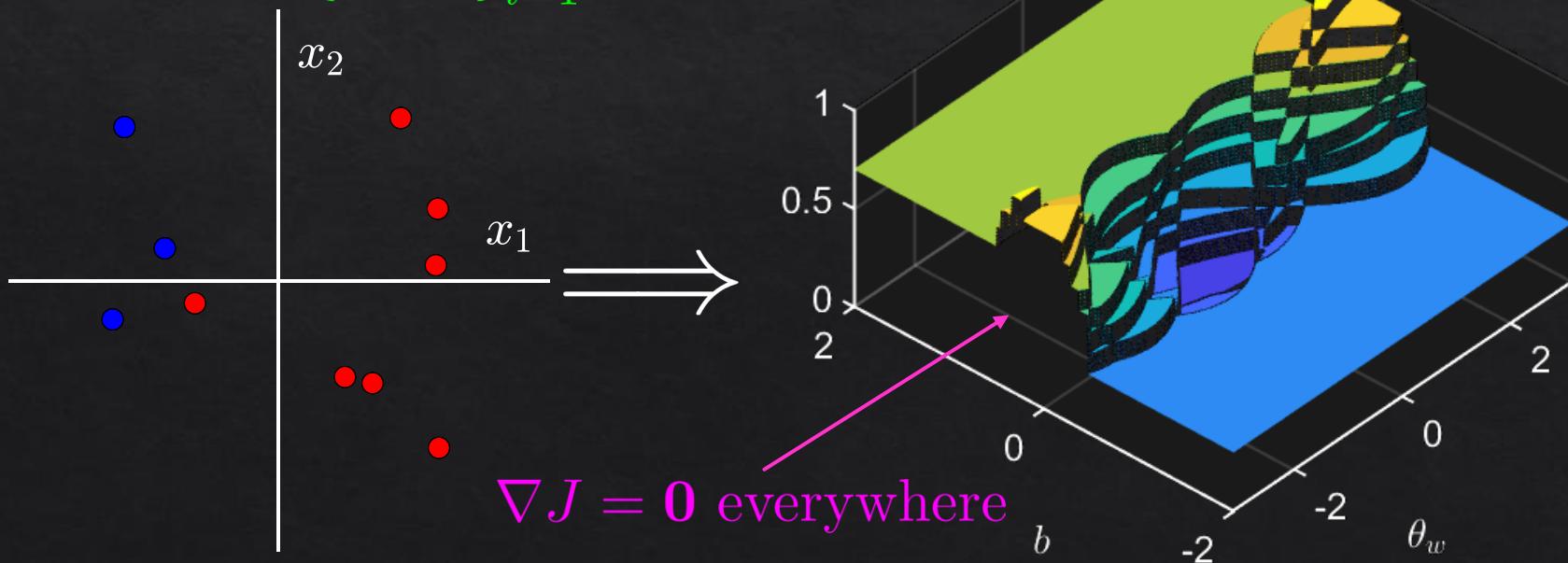
- The objective is given by:

$$\tilde{\mathbf{w}}^* = \arg \min_{\tilde{\mathbf{w}}} \frac{1}{N} \sum_{i=1}^N \mathbb{I} \left\{ \text{sign} (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \neq y_i \right\}$$

$=: J(\tilde{\mathbf{w}})$

- Note that the objective $J(\tilde{\mathbf{w}})$ is not differentiable.
For example:

$$\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, y_i\}_{i=1}^N$$



$$\arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N \mathbb{I} \{ f(\mathbf{x}_i) \neq y_i \}$$

$$f(\mathbf{x}) = \text{sign} (\mathbf{w}^T \mathbf{x} - b)$$

$$\tilde{\mathbf{x}}_i := \begin{bmatrix} -1 \\ | \\ \mathbf{x}_i \\ | \end{bmatrix} \in \mathbb{R}^{d+1}$$

$$\tilde{\mathbf{w}} := \begin{bmatrix} b \\ | \\ \mathbf{w} \\ | \end{bmatrix} \in \mathbb{R}^{d+1}$$

$$\Rightarrow \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i = \mathbf{w}^T \mathbf{x}_i - b$$

Linear Classifiers – Objective Approximation

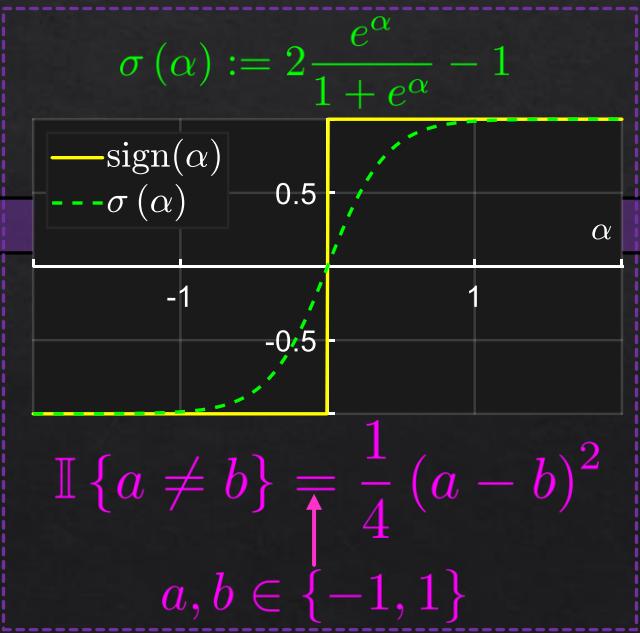
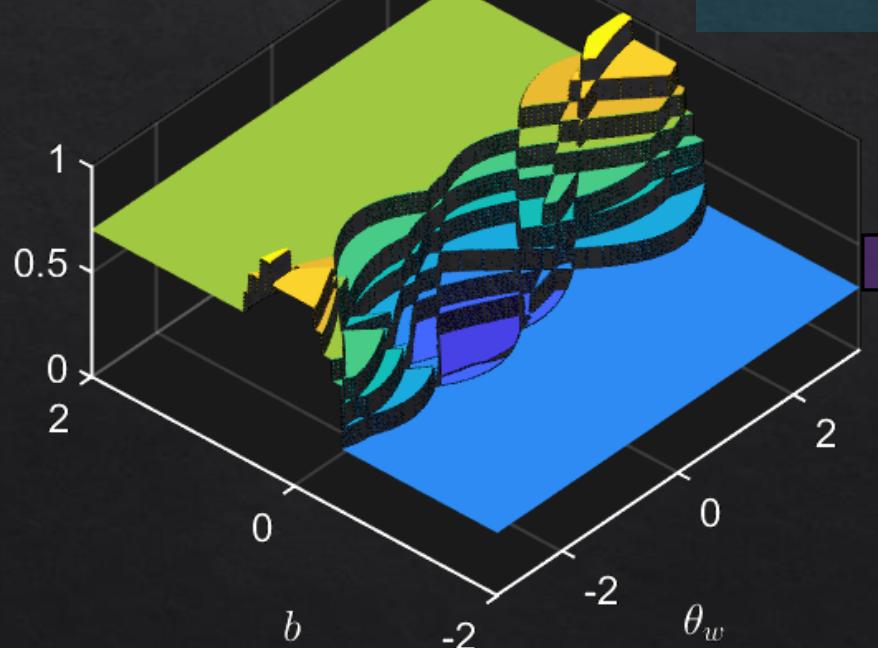
$$\tilde{\mathbf{w}}^* = \arg \min_{\tilde{\mathbf{w}}} \frac{1}{N} \sum_{i=1}^N \mathbb{I} \left\{ \text{sign} (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \neq y_i \right\}$$

$=: J(\tilde{\mathbf{w}})$

$J(\tilde{\mathbf{w}})$

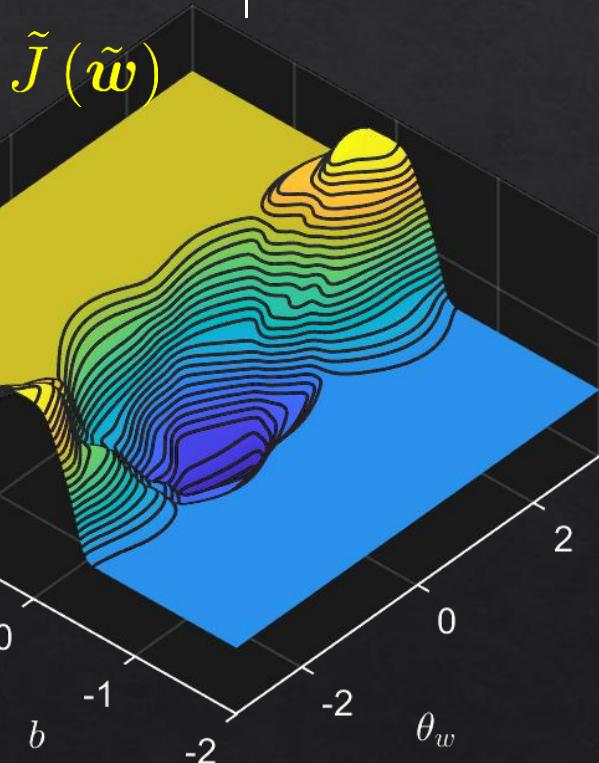
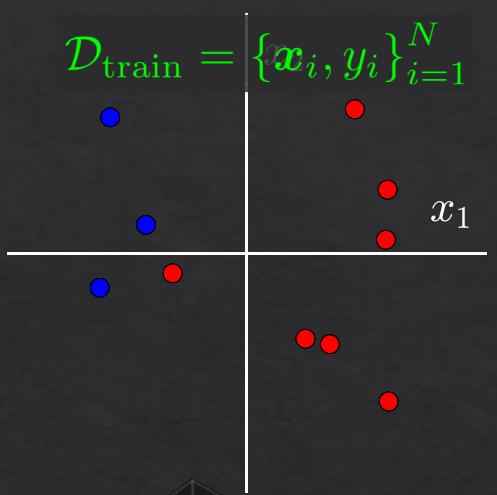
$$\tilde{\mathbf{w}}^* = \arg \min_{\tilde{\mathbf{w}}} \frac{1}{N} \sum_{i=1}^N \left(\sigma (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) - y_i \right)^2$$

$=: \tilde{J}(\tilde{\mathbf{w}})$



$$\tilde{\mathbf{w}}^* = \arg \min_{\tilde{\mathbf{w}}} \frac{1}{N} \sum_{i=1}^N \left(\sigma (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) - y_i \right)^2$$

$=: \tilde{J}(\tilde{\mathbf{w}})$



$\tilde{J}(\tilde{\mathbf{w}})$ is a differentiable approximation of $J(\tilde{\mathbf{w}})$

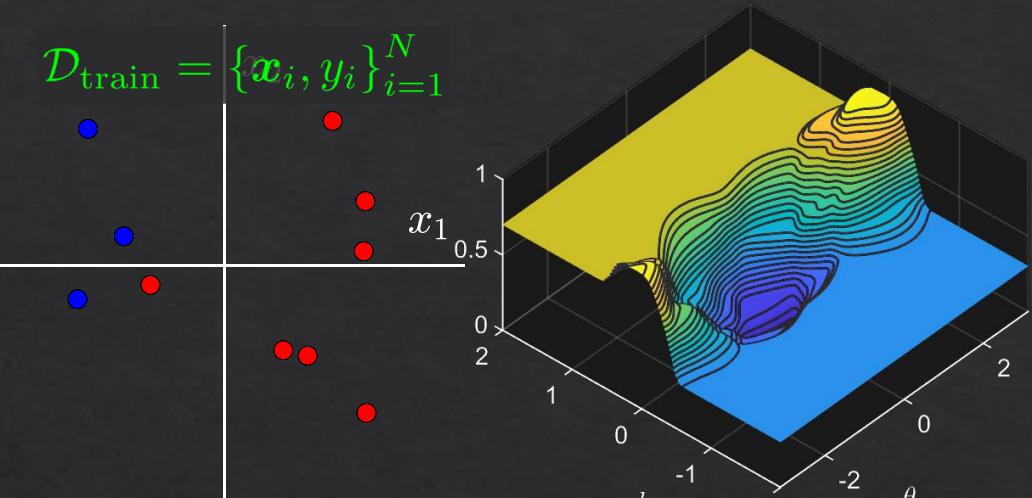
Sigmoid Approximation – Gradient

$$\tilde{\mathbf{w}}^* = \arg \min_{\tilde{\mathbf{w}}} \frac{1}{N} \sum_{i=1}^N (\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) - y_i)^2$$

$\qquad\qquad\qquad =: \tilde{J}(\tilde{\mathbf{w}})$

$$\tilde{J}(\tilde{\mathbf{w}}) =$$

$$\left\| \begin{bmatrix} \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_1) - y_1 \\ \vdots \\ \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_N) - y_N \end{bmatrix} \right\|_2^2 = \left\| \sigma \left(\begin{bmatrix} \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_1 \\ \vdots \\ \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_N \end{bmatrix} \right) - \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right\|_2^2 = \underbrace{\left\| \sigma \left(\begin{bmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_N^T \end{bmatrix} \begin{bmatrix} | \\ \tilde{\mathbf{w}} \\ | \end{bmatrix} \right) - \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right\|_2^2}$$



- 1D derivative ($w, x, y \in \mathbb{R}$):

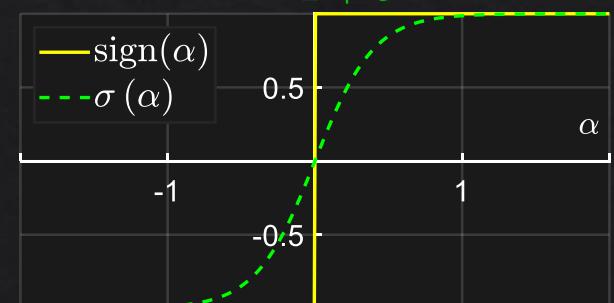
$$\frac{d}{dw} (\sigma(wx) - y)^2 = 2x\sigma'(wx)(\sigma(wx) - y)$$

- Full gradient:

$$\nabla_{\tilde{\mathbf{w}}} J(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{X}}^T \sigma' \left(\text{diag}(\tilde{\mathbf{X}} \tilde{\mathbf{w}}) \right) (\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{w}}) - \mathbf{y})$$

$$\tilde{J}(\tilde{\mathbf{w}}) = \left\| \sigma(\tilde{\mathbf{X}} \tilde{\mathbf{w}}) - \mathbf{y} \right\|_2^2$$

$$\sigma(\alpha) := 2 \frac{e^\alpha}{1 + e^\alpha} - 1$$



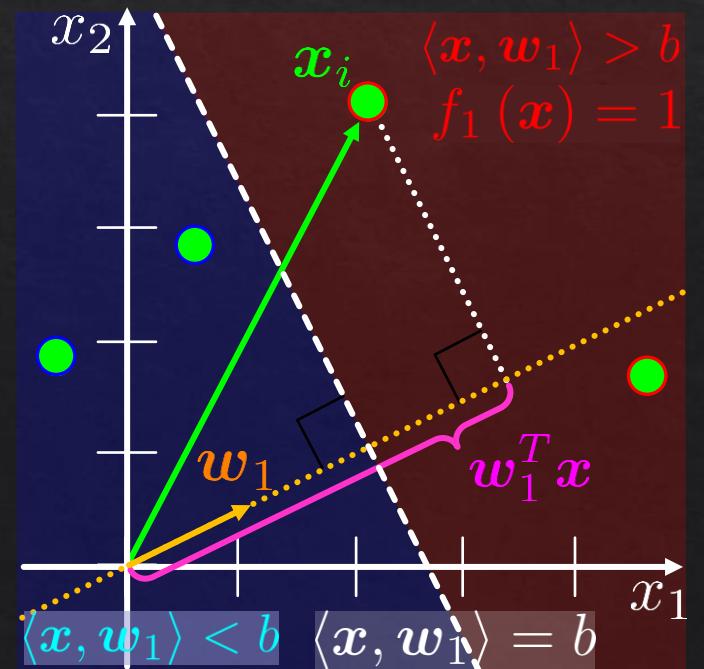
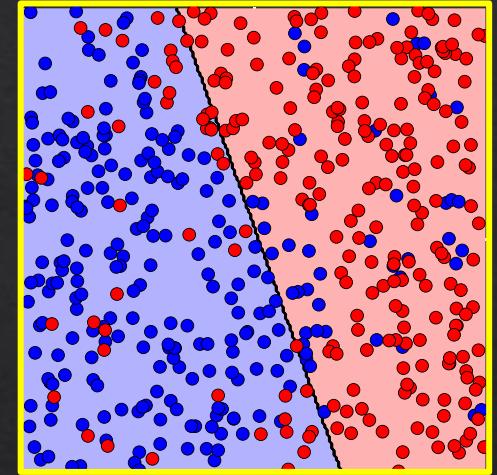
Train Linear Classifier – Example



Train a Linear Classifier

$$\nabla_{\tilde{\mathbf{w}}} J(\tilde{\mathbf{w}}) = 2 \tilde{\mathbf{X}}^T \sigma' \left(\text{diag}(\tilde{\mathbf{X}} \tilde{\mathbf{w}}) \right) \left(\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{w}}) - \mathbf{y} \right)$$

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$



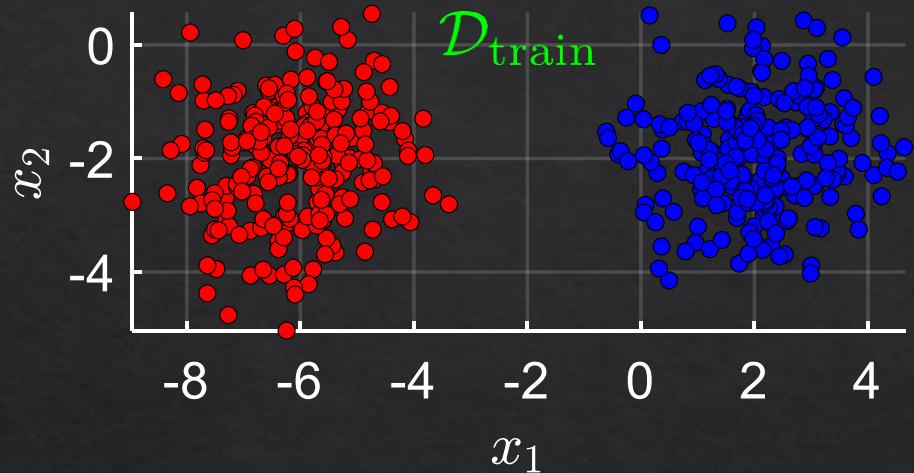
Binary Classification Exercise



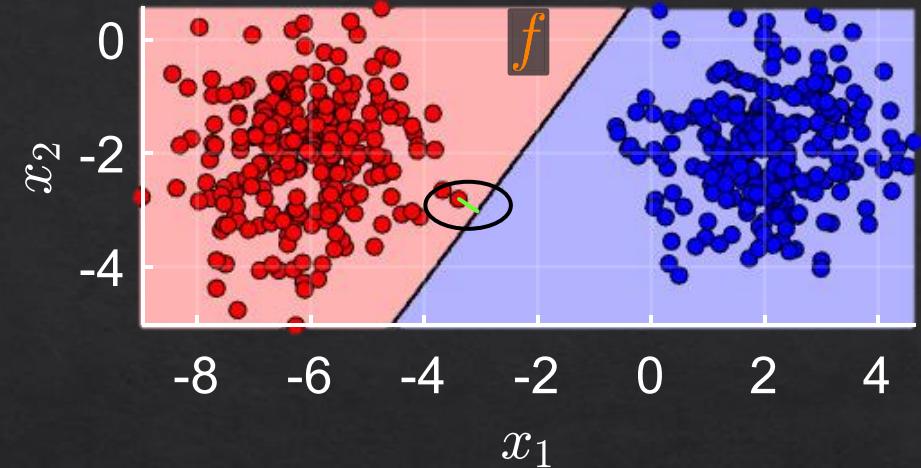
Binary Classification Exercise

Support Vector Machine – Introduction

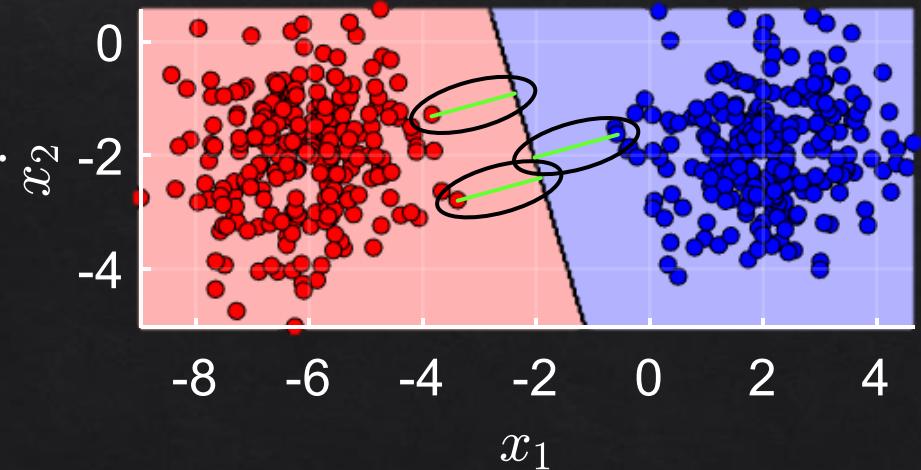
- Consider a linear separable $\mathcal{D}_{\text{train}}$ and the following linear classifier f :



Training →



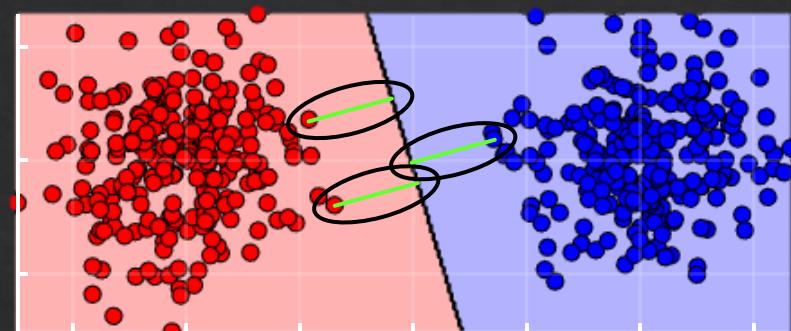
- The accuracy of f is 100%.
However, some points are close to the boundary. (small margins)
- The classifier f is sensitive to small perturbations.
- The SVM classifier seeks to maximize the margins.



Support Vector Machine – Margins

- A perfect linear classifier satisfies:

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i - b > 0 & y_i = 1 \\ \mathbf{w}^T \mathbf{x}_i - b < 0 & y_i = -1 \end{cases} \implies \underbrace{y_i (\mathbf{w}^T \mathbf{x}_i - b)}_{\text{the margin from the decision boundary}} > 0$$



$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$

- The SVM classifier requires:

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1, \quad \forall i$$

- The value 1 is arbitrary and could be replaced with any other positive value.
- This can always be satisfied by a perfect linear classifier:

$$\underbrace{y_i (\mathbf{w}^T \mathbf{x}_i - b) > 0}_{\text{a perfect classifier}} \implies \text{some large } C \quad y_i C (\mathbf{w}^T \mathbf{x}_i - b) \geq 1,$$
$$y_i \underbrace{C \mathbf{w}^T \mathbf{x}_i}_{\text{new } w} - \underbrace{C b}_{\text{new } b} \geq 1$$

Support Vector Machine – Gap

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1, \quad \forall i$$

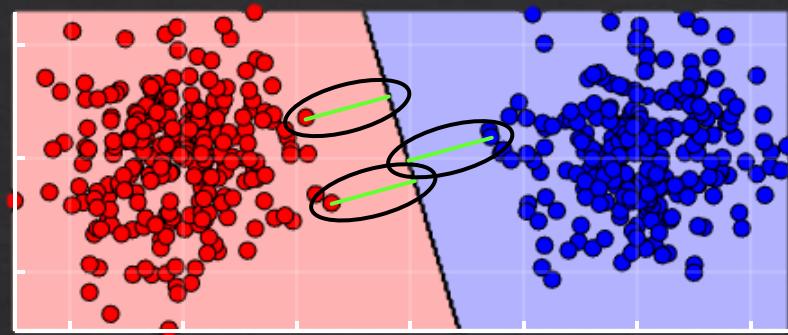
- Let \mathbf{x}_+ and \mathbf{x}_- be exactly on the gutter, that is:

$$\begin{cases} \mathbf{w}^T \mathbf{x}_+ - b = 1 \\ \mathbf{w}^T \mathbf{x}_- - b = -1 \end{cases}$$

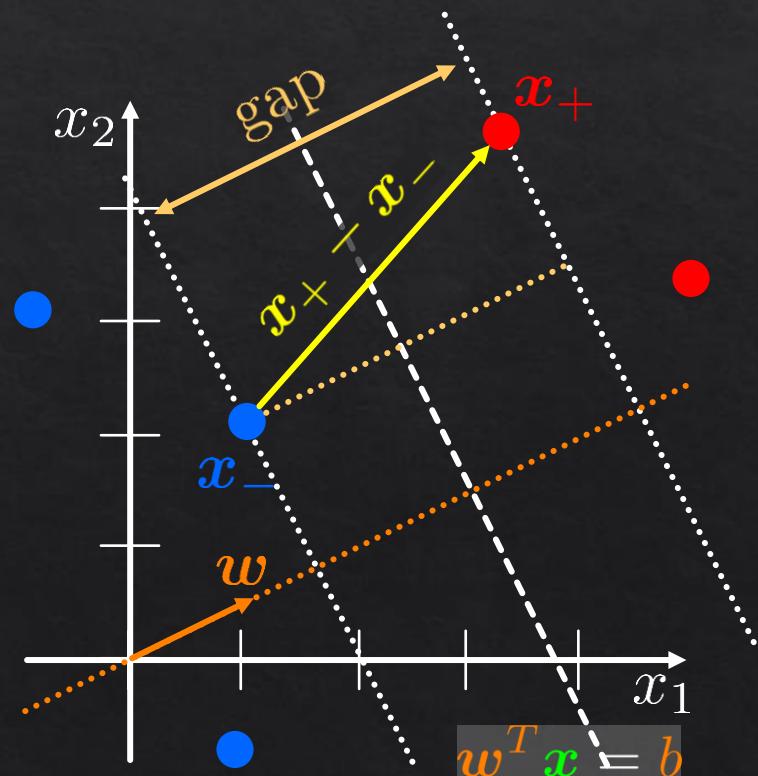
$$\text{gap} = \frac{\mathbf{w}^T}{\|\mathbf{w}\|} (\mathbf{x}_+ - \mathbf{x}_-)$$

$$\begin{aligned} &= \frac{\mathbf{w}^T \mathbf{x}_+}{\|\mathbf{w}\|} - \frac{\mathbf{w}^T \mathbf{x}_-}{\|\mathbf{w}\|} \\ &= \frac{1 + b}{\|\mathbf{w}\|} - \frac{-1 + b}{\|\mathbf{w}\|} \\ &= \frac{2}{\|\mathbf{w}\|} \end{aligned}$$

- The SVM classifier seeks a \mathbf{w} that maximizes the gap.



$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$



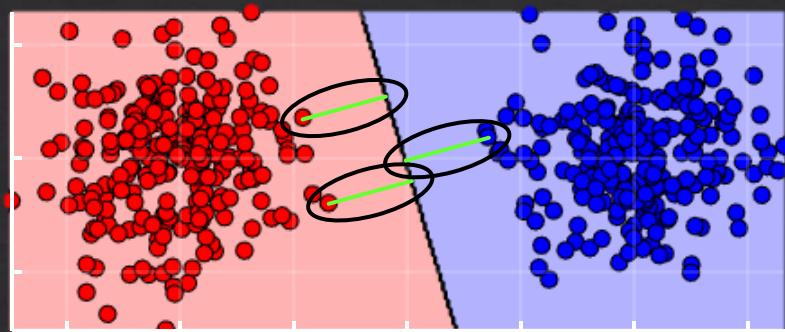
Support Vector Machine – Objective

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1, \quad \forall i$$

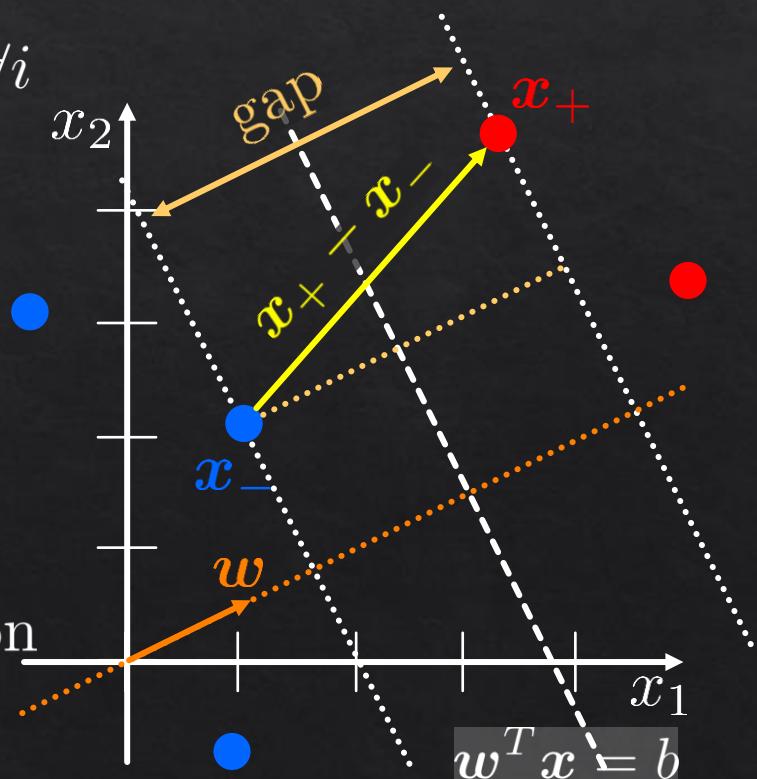
$$\text{gap} = \frac{2}{\|\mathbf{w}\|}$$

- SVM seeks to maximize the gap: $\begin{cases} \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \\ \text{subject to} \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \quad \forall i \end{cases}$
- Instead of maximize $\frac{1}{\|\mathbf{w}\|}$, we can minimize $\frac{1}{2} \|\mathbf{w}\|^2$:
$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \quad \forall i \end{cases}$$

- There are efficient solvers for this constrained optimization (such as quadratic programming).

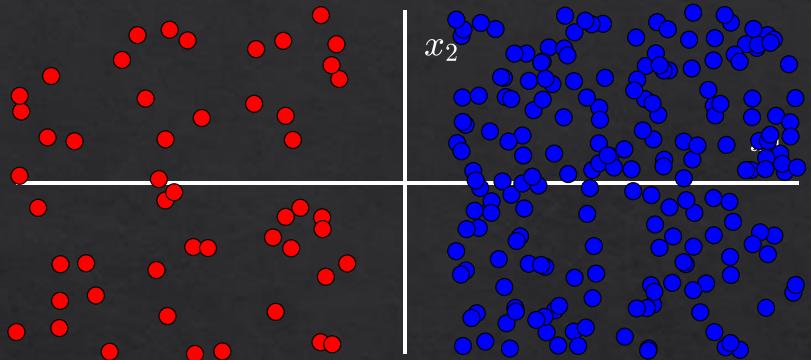


$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b)$$

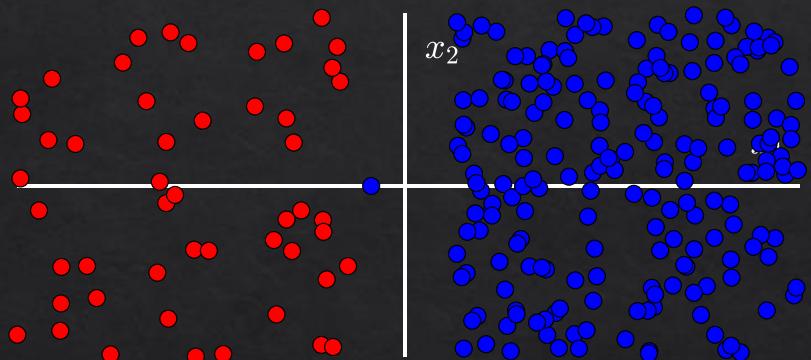
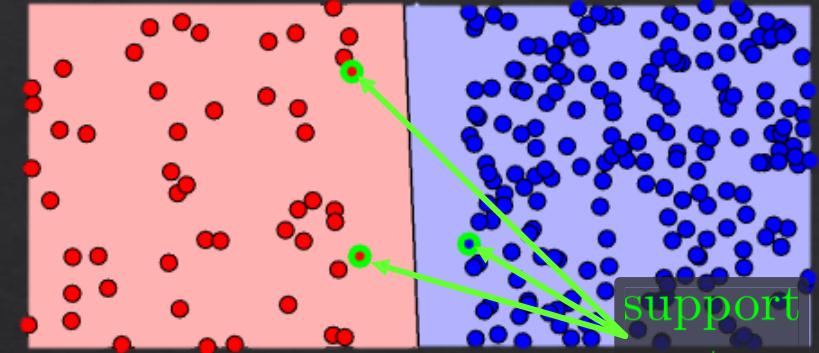


(Hard) Support Vector Machine – Examples

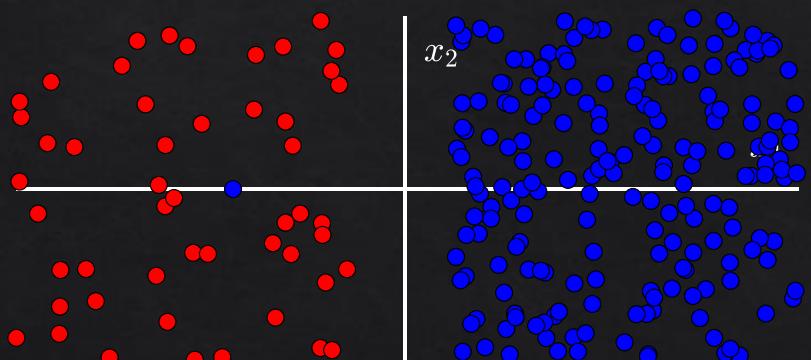
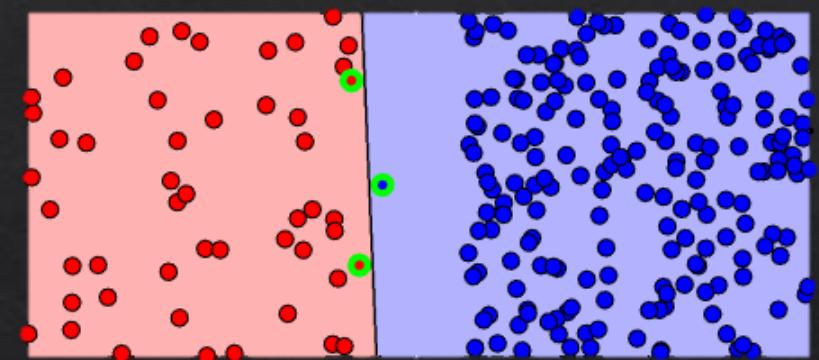
$$y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1$$



SVM



SVM

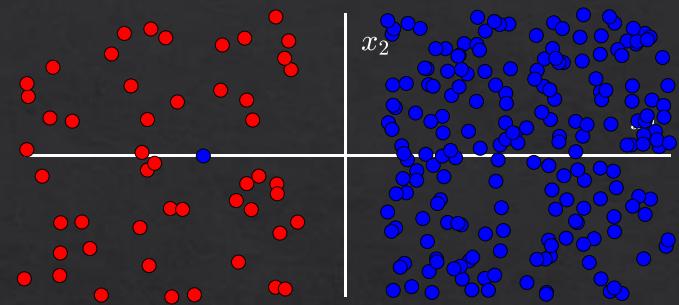


SVM

no (hard) solution

Soft Support Vector Machine

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \quad \forall i \end{cases}$$



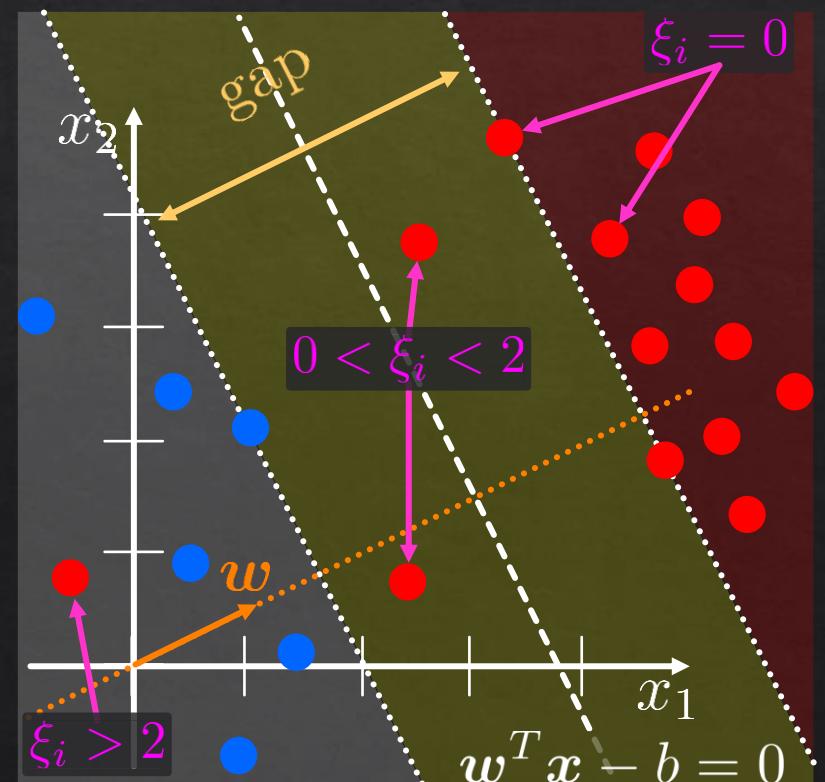
- In soft SVM,
the (hard) constraints are relaxed using regularization:

$$\xi_i := \max \{0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i - b)\}$$

- ξ_i is the penalty cost for points
on the wrong side of the margin
- The soft SVM is given by:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

- $C > 0$ is the regularization factor.

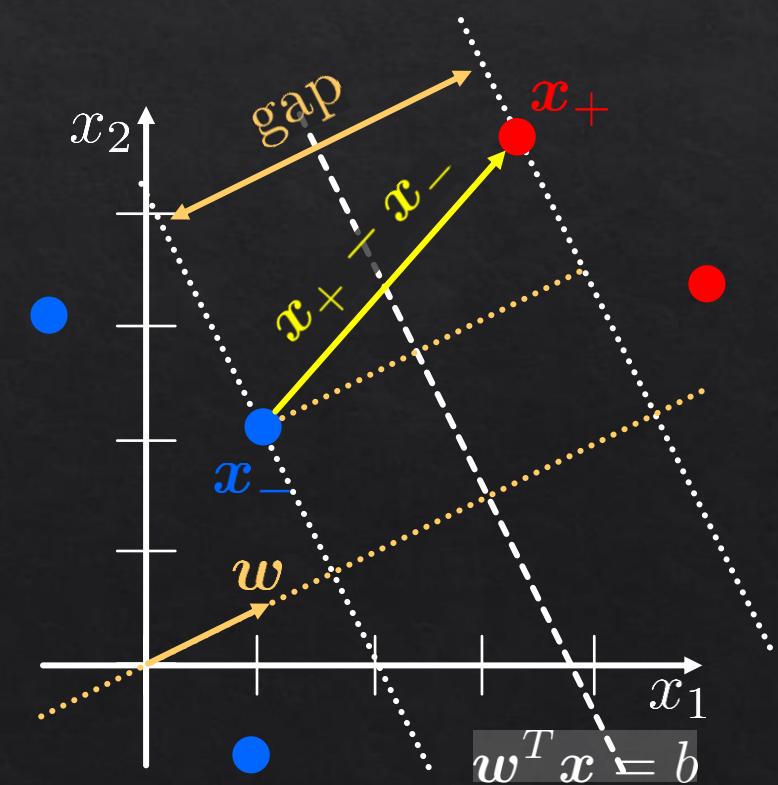


(Soft) SVM Example



SVM Example

$$\min_{\boldsymbol{w}, b} \frac{1}{2} \|\boldsymbol{w}\|^2 + C \sum_i \xi_i$$



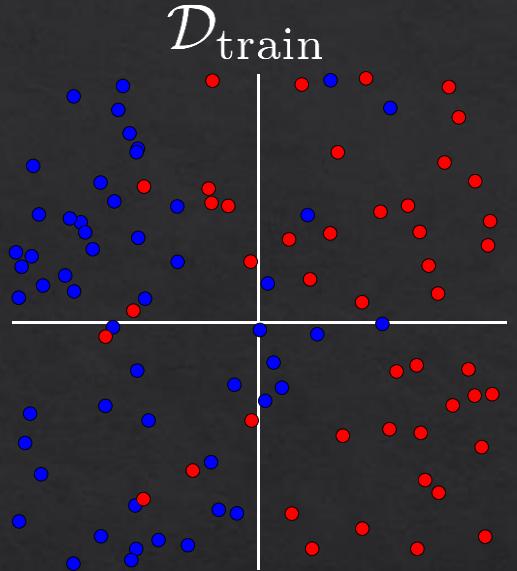
K-Nearest Neighbors

- Let $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$.
- The 1-NN classifier is given by:

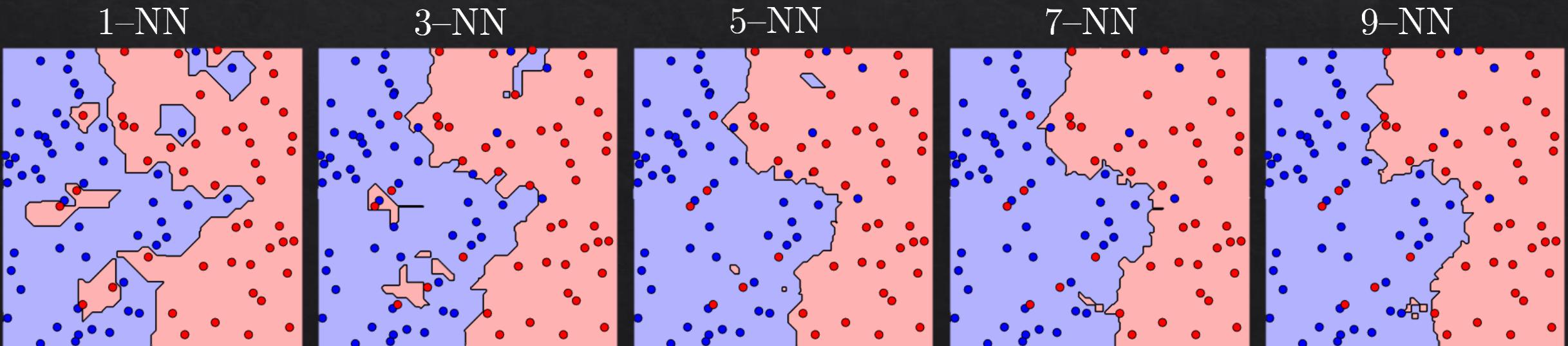
$$f_{1-\text{NN}}(\mathbf{x}) = \hat{y}_{i(\mathbf{x})}$$

where

$$\hat{i}(\mathbf{x}) = \arg \min_i d(\mathbf{x}, \mathbf{x}_i) \quad \begin{array}{l} \text{(the model is)} \\ \text{nonparametric} \end{array}$$



- For example, using the Euclidean metric $d(\mathbf{x}, \mathbf{x}_i) = \|\mathbf{x} - \mathbf{x}_i\|_2$:
- For $K > 1$, \mathbf{x} will be labeled by the most common among its K NN.



K-Nearest Neighbors – Examples

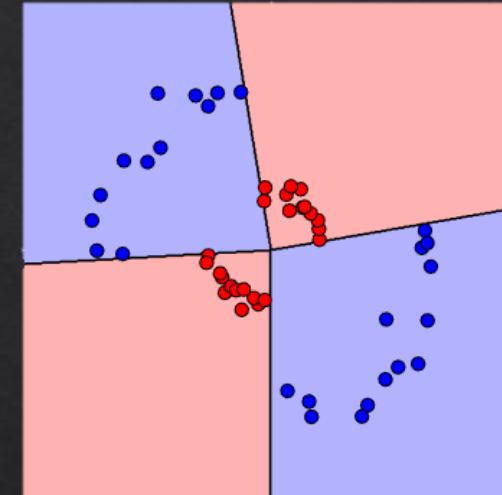
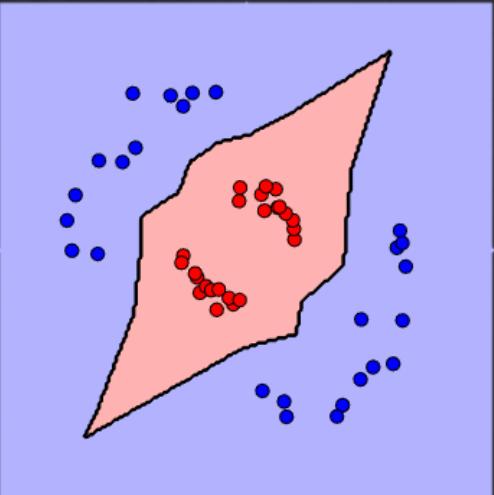
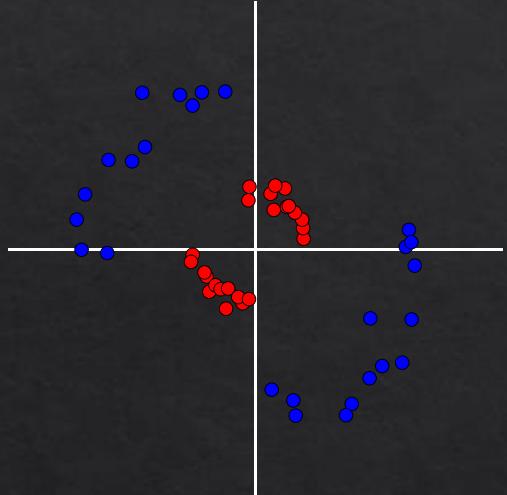
$$f_{1\text{-NN}}(\mathbf{x}) = \hat{y}_{i(\mathbf{x})}$$

$$\hat{i}(\mathbf{x}) = \arg \min_i d(\mathbf{x}, \mathbf{x}_i)$$

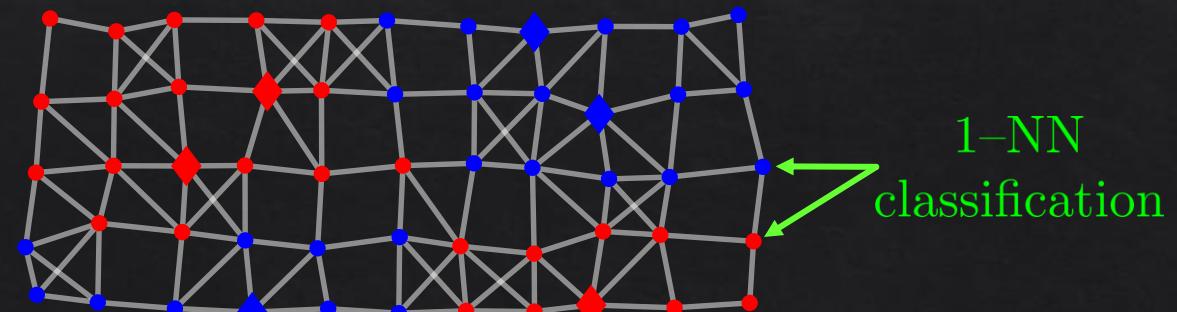
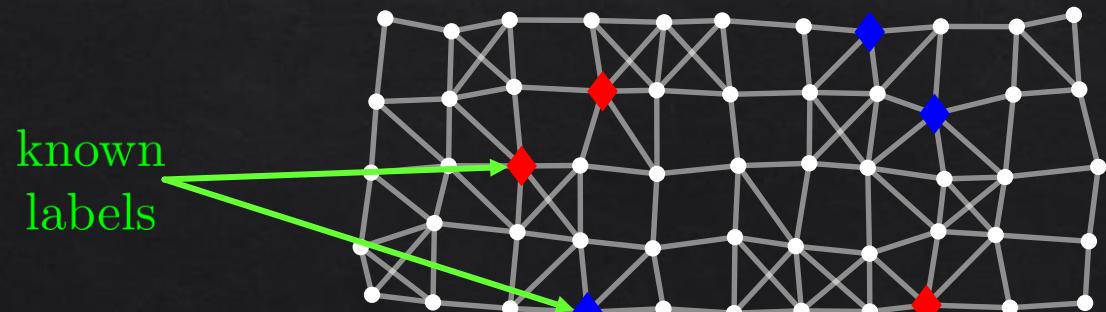
- 1–NN Euclidean and cosine distances:

$$d_{\text{Euclid}}(\mathbf{x}, \mathbf{x}_i) = \|\mathbf{x} - \mathbf{x}_i\|_2$$

$$d_{\text{cosine}}(\mathbf{x}, \mathbf{x}_i) = 1 - \frac{\langle \mathbf{x}, \mathbf{x}_i \rangle}{\|\mathbf{x}\|_2 \|\mathbf{x}_i\|_2} = 1 - \cos(\theta)$$



- 1–NN on graph (distance is shortest path):

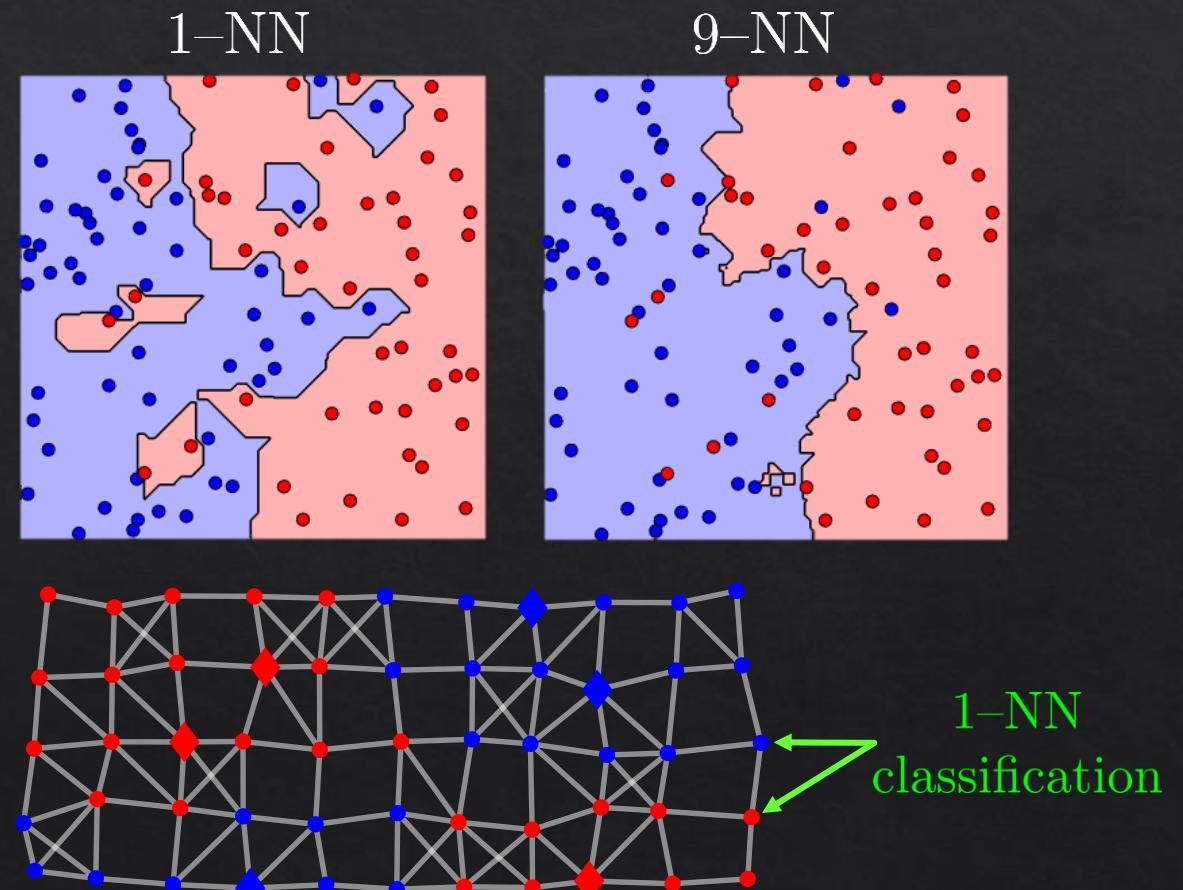


- What is the train accuracy of a 1–NN classifier?

K-NN Example



K-NN Example

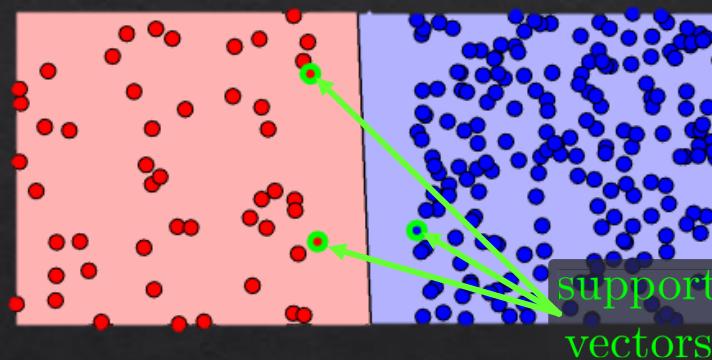
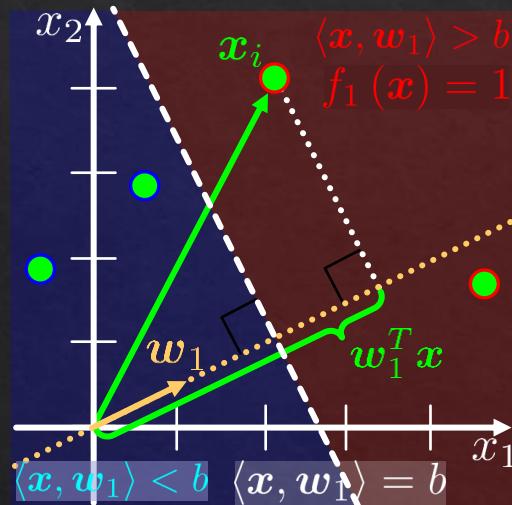


Questions

$$L_{\text{train}}(f) := \frac{1}{N} \sum_{i=1}^N \ell(\hat{y}_i, y_i) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{y}_i \neq y_i\}$$

$1 - L(f)$ is the classifier accuracy
 $L(f) = 0 \implies$ no errors

SVM:



K-NN

