

Day 4: Classification, support vector machines

Introduction to Machine Learning Summer School
June 18, 2018 - June 29, 2018, Chicago

Instructor: Suriya Gunasekar, TTI Chicago

21 June 2018



THE UNIVERSITY OF
CHICAGO



Topics so far

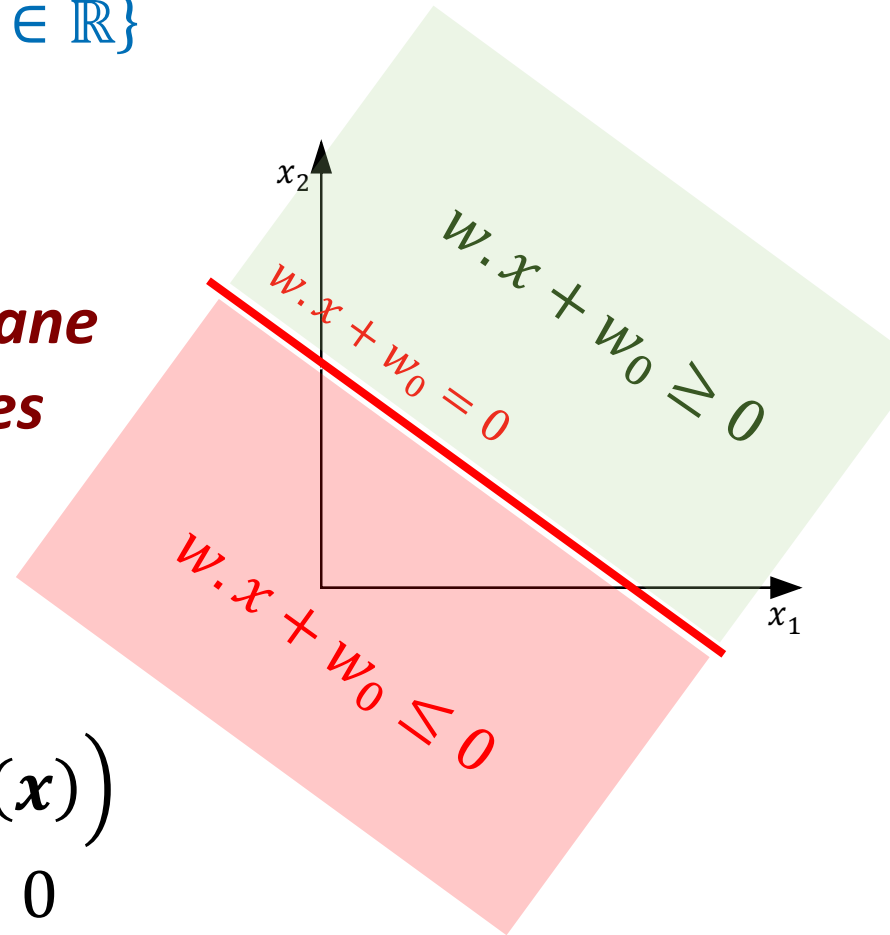
- Supervised learning, linear regression
- Linear regression
 - Overfitting, bias variance trade-off
 - Ridge and lasso regression, gradient descent
- Yesterday
 - Classification, logistic regression
 - Regularization for logistic regression
 - Multi-class classification
- Today
 - Maximum margin classifiers
 - Kernel trick

Classification

- Supervised learning: estimate a mapping f from input $x \in \mathcal{X}$ to output $y \in \mathcal{Y}$
 - **Regression** $\mathcal{Y} = \mathbb{R}$ or other continuous variables
 - **Classification** \mathcal{Y} takes discrete set of values
 - Examples:
 - $\mathcal{Y} = \{\text{spam}, \text{nospam}\}$,
 - digits (not values) $\mathcal{Y} = \{0, 1, 2, \dots, 9\}$
- Many successful applications of ML in vision, speech, NLP, healthcare

Parametric classifiers

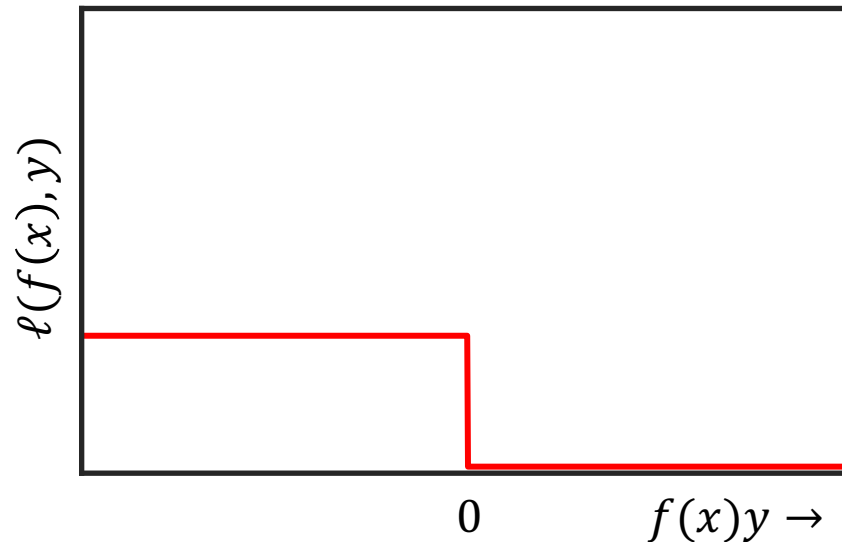
- $\mathcal{H} = \{\mathbf{x} \rightarrow \mathbf{w} \cdot \mathbf{x} + w_0 : \mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}\}$
- $\hat{y}(\mathbf{x}) = \text{sign}(\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0)$
- $\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0 = 0$ (linear) decision boundary or separating **hyperplane** separates \mathbb{R}^d into two **halfspaces** (regions)
 - $\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0 > 0$ gets label 1 and
 - $\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0 < 0$ gets label -1
- more generally, $\hat{y}(\mathbf{x}) = \text{sign}(\hat{f}(\mathbf{x}))$
 - decision boundary is $\hat{f}(\mathbf{x}) = 0$



Surrogate Losses

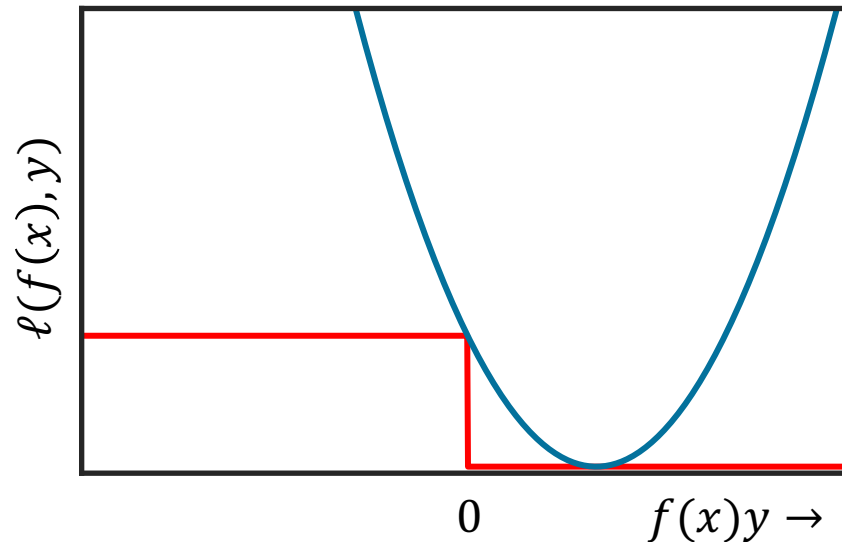
- The correct loss to use is 0-1 loss **after** thresholding

$$\begin{aligned}\ell^{01}(f(x), y) &= \mathbf{1}[\text{sign}(f(x)) \neq y] \\ &= \mathbf{1}[\text{sign}(f(x)y) < 0]\end{aligned}$$



Surrogate Losses

- The correct loss to use is 0-1 loss **after** thresholding
$$\ell^{01}(f(x), y) = \mathbf{1}[\text{sign}(f(x)) \neq y]$$
$$= \mathbf{1}[\text{sign}(f(x)y) < 0]$$
- Linear regression uses $\ell^{LS}(f(x), y) = (f(x) - y)^2$



Surrogate Losses

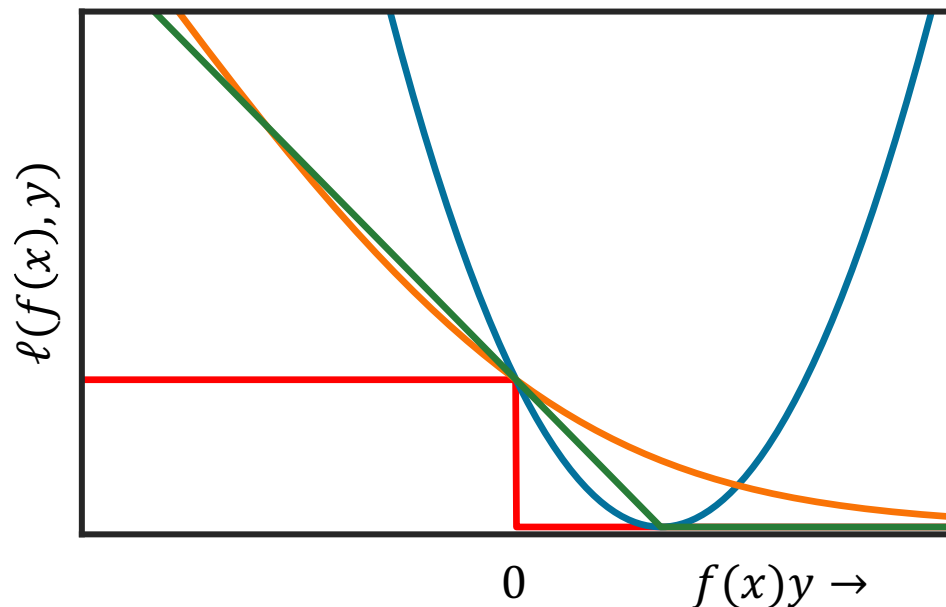
- Hard to optimize over ℓ^{01} , find another loss $\ell(f(x), y)$
 - Convex (for any fixed y) \rightarrow easier to minimize
 - An upper bound of $\ell^{01} \rightarrow$ small $\ell \Rightarrow$ small ℓ^{01}
 - Satisfied by squared loss
- \rightarrow but has “large” loss even when $\ell^{01}(f(x), y) = 0$
- Two more surrogate losses in in this course

- **Logistic loss**

$$\ell^{\log}(f(x), y) = \log(1 + \exp(-f(x)y))$$

- **Hinge loss**

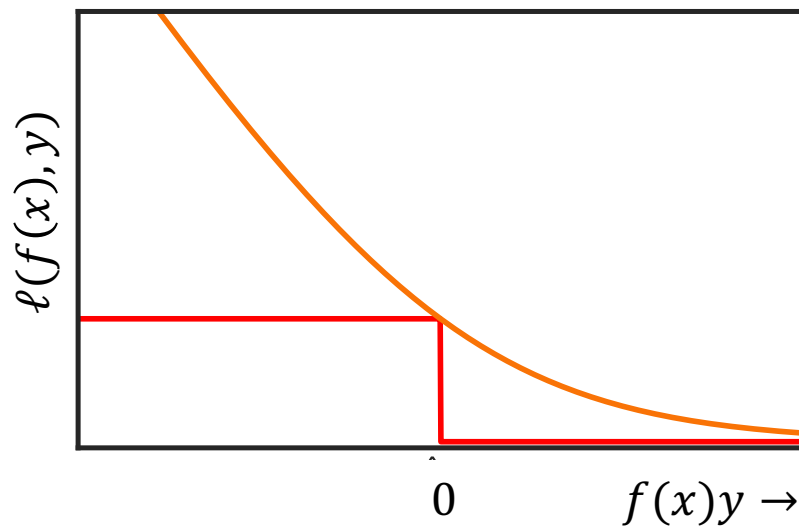
$$\ell^{\text{hinge}}(f(x), y) = \max(0, 1 - f(x)y)$$



Logistic regression: ERM on surrogate loss

Logistic loss

$$\ell(f(x), y) = \log(1 + \exp(-f(x)y))$$



- $S = \{(\mathbf{x}^{(i)}, y^{(i)}): i = 1, 2, \dots, N\}$, $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{-1, 1\}$
- Linear model $f(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$
- Minimize training loss

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmin}_{\mathbf{w}, w_0} \sum_i \log(1 + \exp(-(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)y^{(i)}))$$

- Output classifier $\hat{y}(\mathbf{x}) = \operatorname{sign}(\mathbf{w} \cdot \mathbf{x} + w_0)$

Logistic Regression

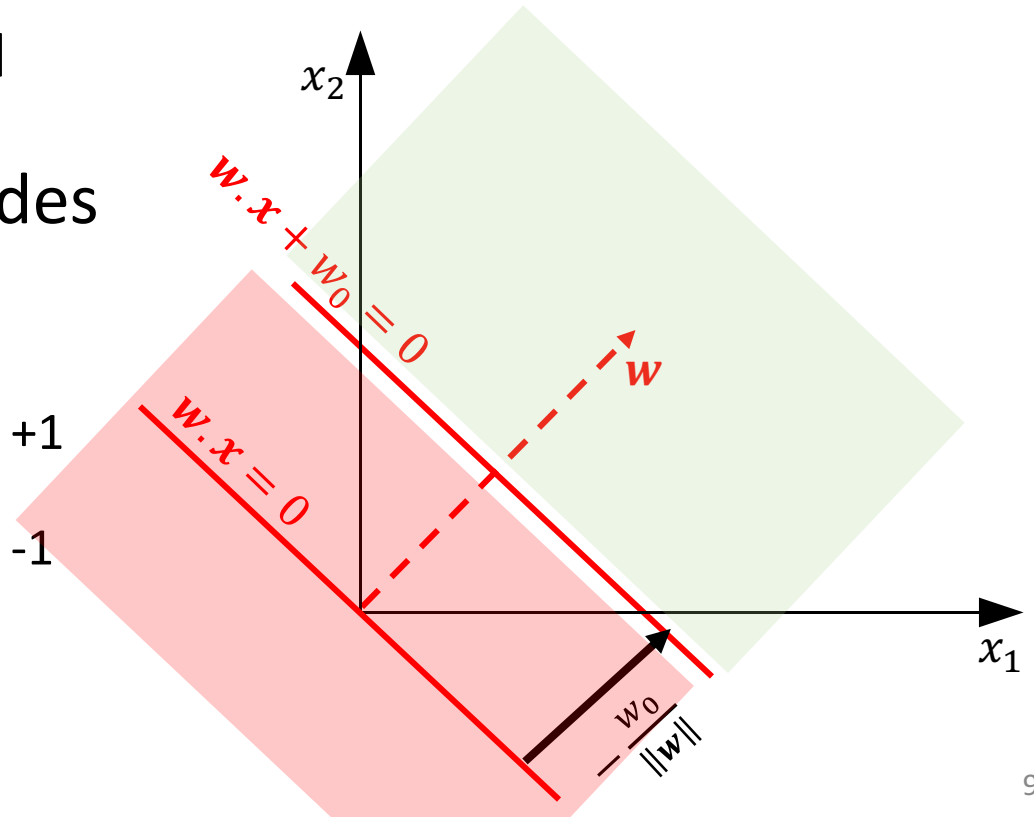
$$\hat{\mathbf{w}}, \hat{w}_0 = \underset{\mathbf{w}, w_0}{\operatorname{argmin}} \sum_i \log \left(1 + \exp(-(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)y^{(i)}) \right)$$

- Convex optimization problem
- Can solve using gradient descent
- Can also add usual regularization: ℓ_2, ℓ_1

Linear decision boundaries

$$\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + w_0)$$

- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$ is a hyperplane in \mathbb{R}^d
 - decision boundary
 - \mathbf{w} is direction of normal
 - w_0 is the offset
- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$ divides \mathbb{R}^d into two halfspaces (regions)
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 \geq 0\}$ get label +1 and
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 < 0\}$ get label -1



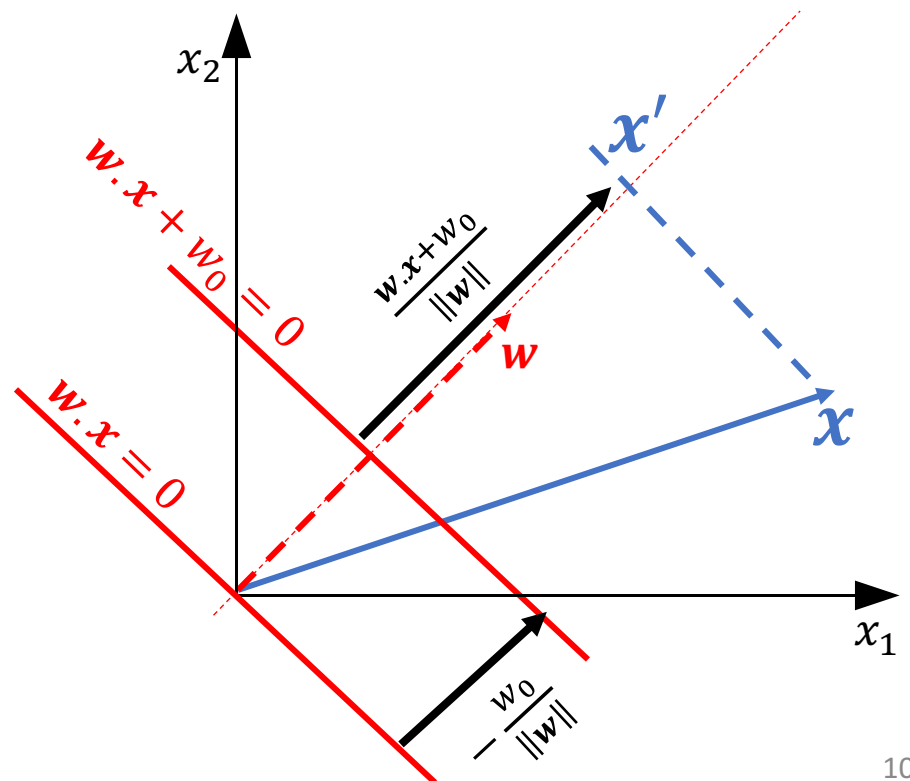
Linear decision boundaries

$$\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + w_0)$$

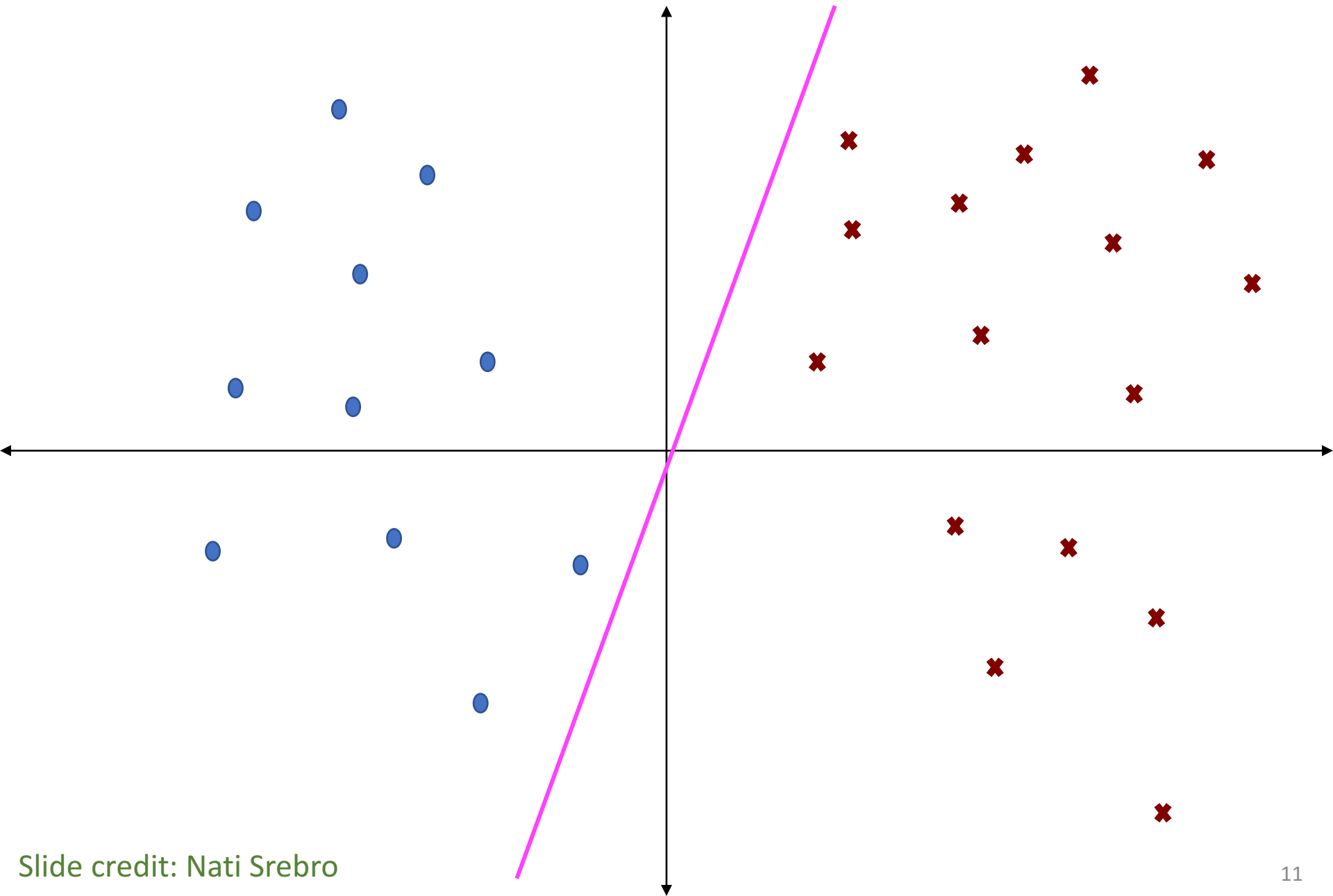
- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$ is a hyperplane in \mathbb{R}^d
 - decision boundary
 - \mathbf{w} is direction of normal
 - w_0 is the offset
- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$ divides \mathbb{R}^d into two halfspaces (regions)
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 \geq 0\}$ get label +1 and
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 < 0\}$ get label -1

Maps \mathbf{x} to a 1D coordinate

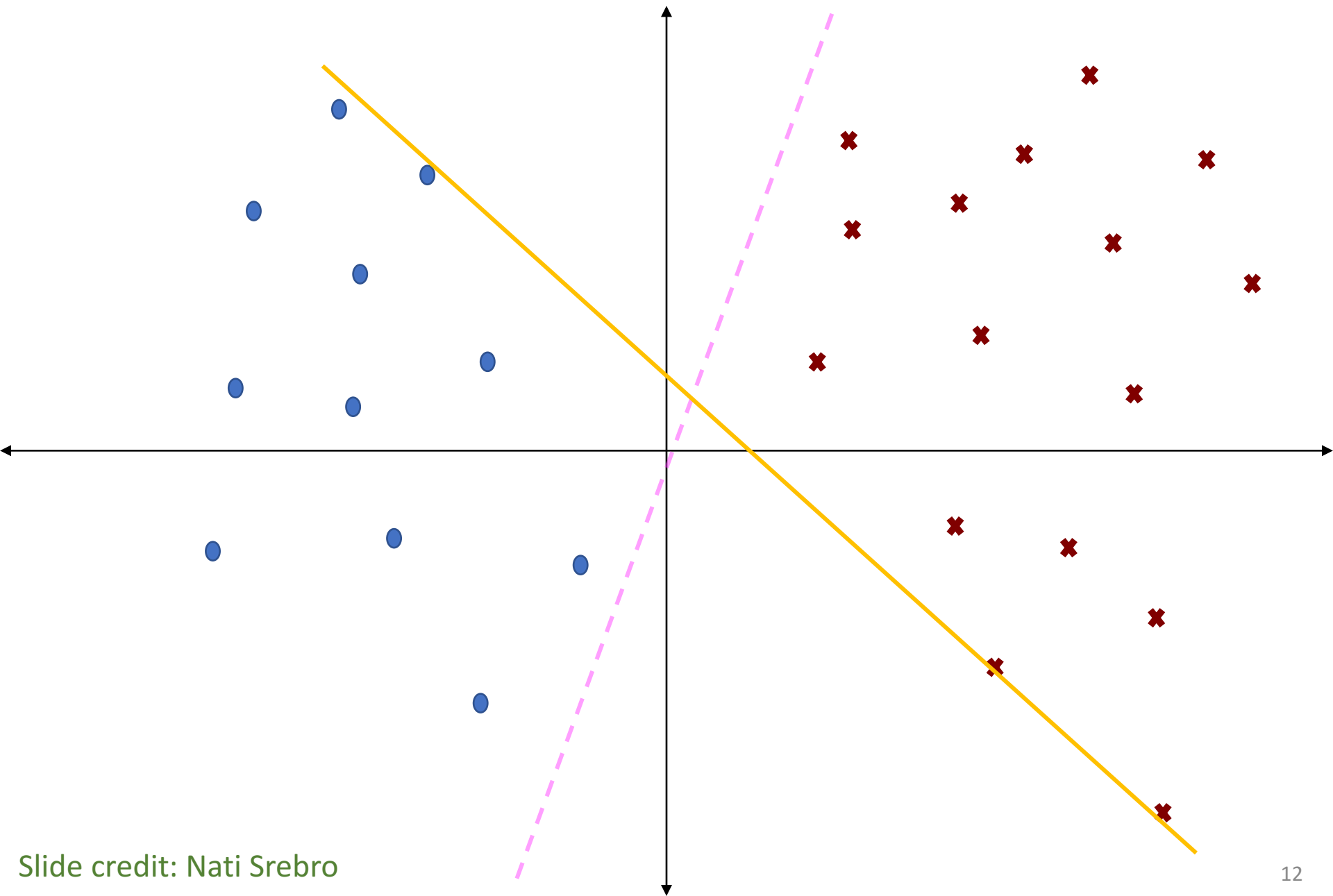
$$x' = \frac{\mathbf{w} \cdot \mathbf{x} + w_0}{\|\mathbf{w}\|}$$



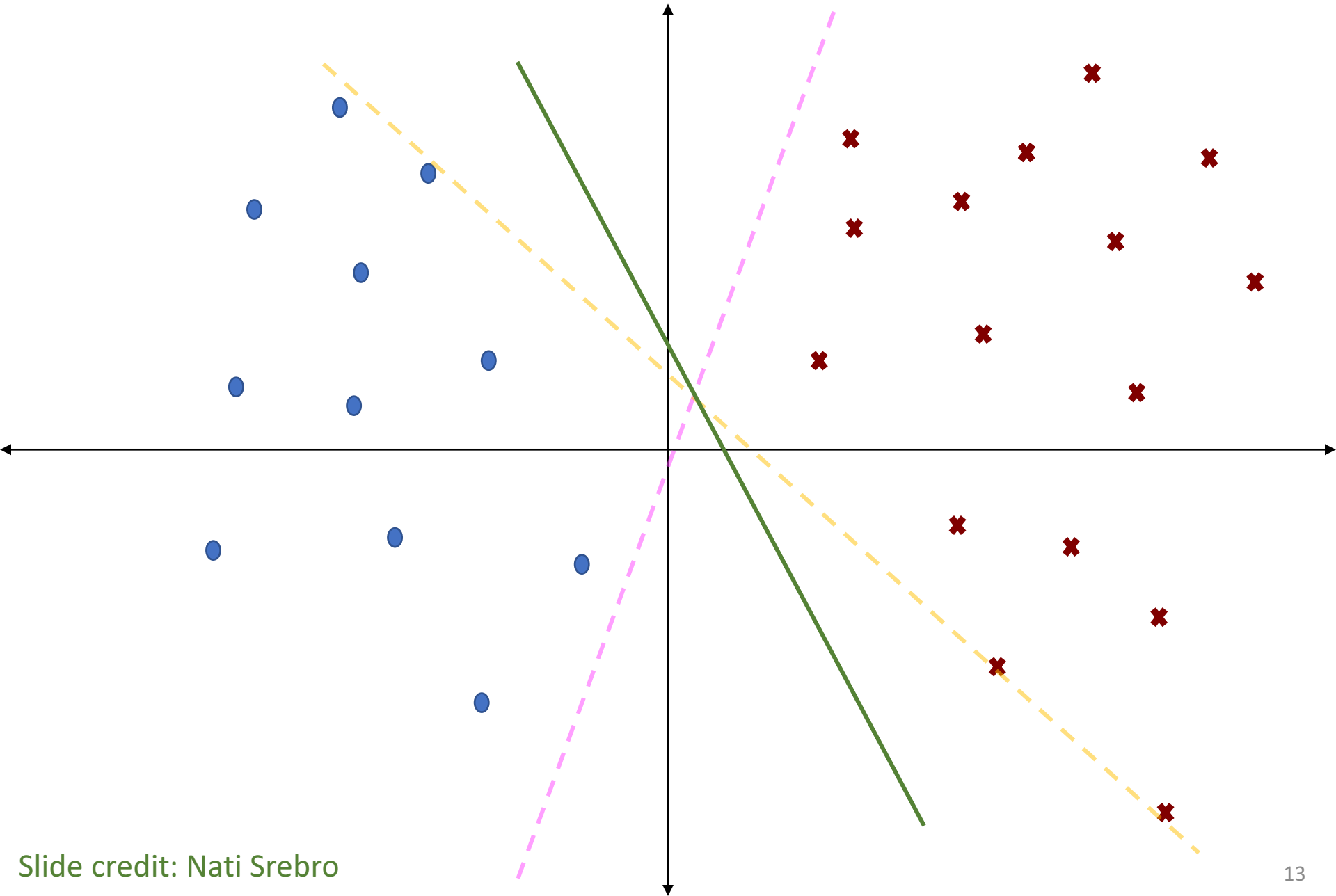
Linear separators in 2D



Linear separators in 2D

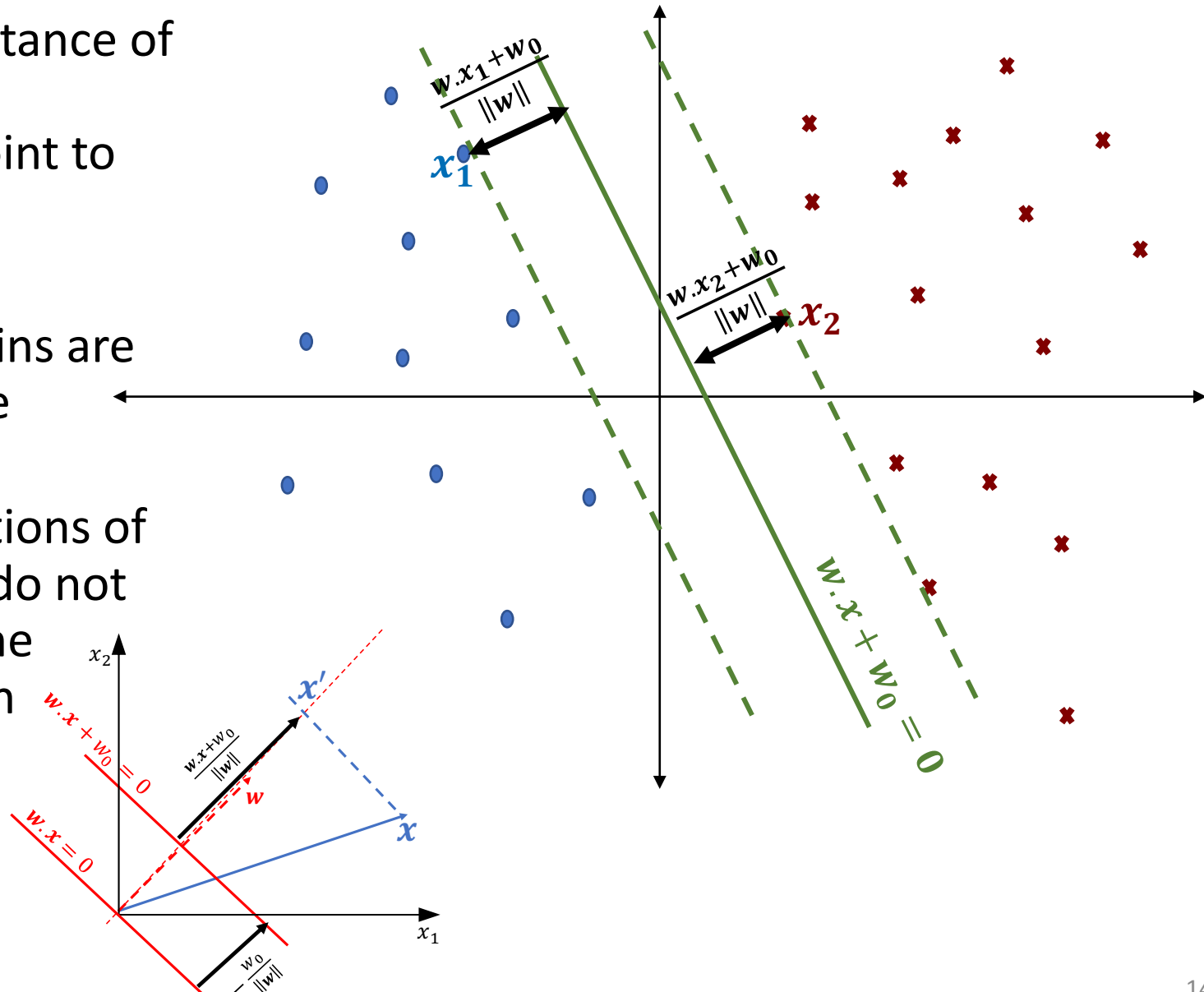


Linear separators in 2D



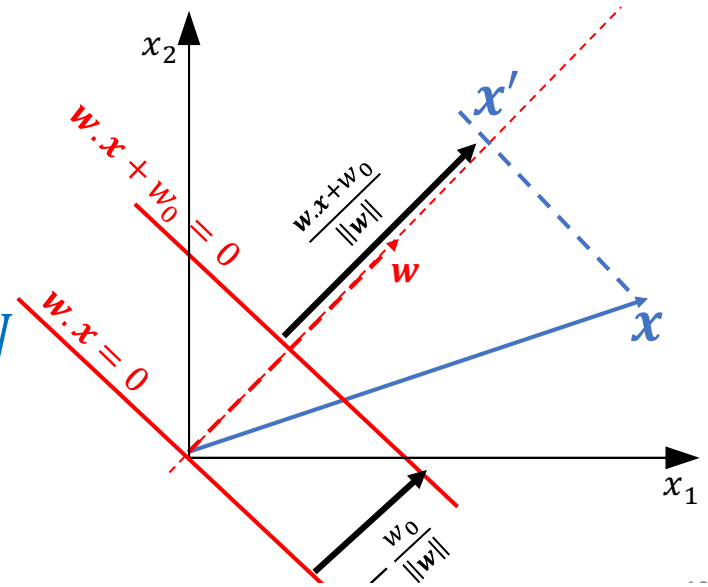
Margin of a classifier

- **Margin:** distance of the closest instance point to the linear hyperplane
- Large margins are more stable
 - small perturbations of the data do not change the prediction



Maximum margin classifier

- $S = \{(\mathbf{x}^{(i)}, y^{(i)}): i = 1, 2, \dots, N\}$
binary classes $\mathcal{Y} = \{-1, 1\}$
- Assume data is “linearly separable”
 - $\exists \mathbf{w}, w_0$ such that for all $i = 1, 2, \dots, N$
 $y^{(i)} = \text{sign}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$
 $\Rightarrow y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) > 0$



- **Maximum margin** separator given by

$$\hat{\mathbf{w}}, \hat{w}_0 = \underset{\mathbf{w} \in \mathbb{R}^d, w_0}{\operatorname{argmax}} \underbrace{\min_i}_{\text{smallest margin}} \underbrace{\frac{y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}}_{\text{margin of sample } i}$$

smallest
margin

margin of
sample i

Maximum margin classifier

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}} \min_i \frac{y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}$$

- **Claim 1:** If $\hat{\mathbf{w}}, \hat{w}_0$ is a solution, then for any $\gamma > 0$, $\gamma \hat{\mathbf{w}}, \gamma \hat{w}_0$ is also a solution
- **Option 1:** We can fix $\|\mathbf{w}\| = 1$ to get

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\|\mathbf{w}\|=1, w_0} \min_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$$

Maximum margin classifier

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}} \min_i \frac{y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}$$

- **Claim 1:** If $\hat{\mathbf{w}}, \hat{w}_0$ is a solution, then for any $\gamma > 0$, $\gamma \hat{\mathbf{w}}, \gamma \hat{w}_0$ is also a solution
- **Option 1:** we can fix $\|\mathbf{w}\| = 1$ to get

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\|\mathbf{w}\|=1} \min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$$

- **Option 2:** we can also fix $\min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) = 1$

- margin now is $\frac{1}{\|\mathbf{w}\|}$

- Instead of “increasing margin” we can “reduce norm”

Max-margin classifier equivalent formulation

- Solve: $\tilde{\mathbf{w}}, \tilde{w}_0 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2$
s.t. $\forall i, y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) \geq 1$

Hard margin
Support Vector
Machine (SVM)

- **Claim 2:** Equivalent to previous slide

→ $\frac{\tilde{\mathbf{w}}}{\|\tilde{\mathbf{w}}\|}, \frac{\tilde{w}_0}{\|\tilde{\mathbf{w}}\|}$ is solution for

$$\hat{\mathbf{w}}, \hat{w}_0 = \max_{\|\mathbf{w}\|=1} \min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$$

- Proof:

1. Let $\min_i y^{(i)}(\hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} + \hat{w}_0) = \hat{\gamma}$, then $\min_i y^{(i)} \left(\frac{\hat{\mathbf{w}}}{\hat{\gamma}} \cdot \mathbf{x}^{(i)} + \frac{\hat{w}_0}{\hat{\gamma}} \right) \geq 1$

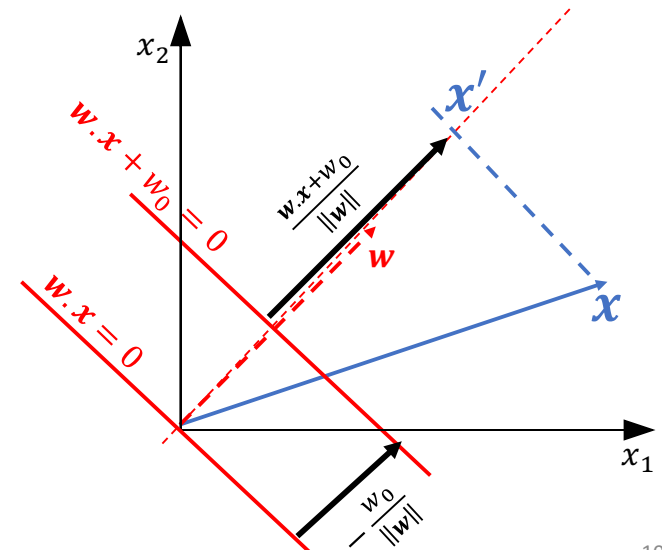
2. $\Rightarrow \|\tilde{\mathbf{w}}\| \leq \left\| \frac{\hat{\mathbf{w}}}{\hat{\gamma}} \right\| = \frac{1}{\hat{\gamma}}$

3. $\min_i y^{(i)} \left(\frac{\tilde{\mathbf{w}}}{\|\tilde{\mathbf{w}}\|} \cdot \mathbf{x}^{(i)} + \frac{\tilde{w}_0}{\|\tilde{\mathbf{w}}\|} \right) = \min_i \frac{y^{(i)}(\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} + \tilde{w}_0)}{\|\tilde{\mathbf{w}}\|} \geq \frac{1}{\|\tilde{\mathbf{w}}\|} \geq \hat{\gamma}$

Maximum margin classifier formulations

- Original formulation

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}} \min_i \frac{y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}$$



- Fixing $\|\mathbf{w}\| = 1$

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\mathbf{w}, w_0} \min_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) \quad \text{s.t.} \quad \|\mathbf{w}\| = 1$$

- Fixing $\min_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) = 1$

$$\tilde{\mathbf{w}}, \tilde{w}_0 = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i, y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) \geq 1$$

Henceforth, w_0 will be absorbed into \mathbf{w} by adding an additional feature of '1' to \mathbf{x}

Margin and norm

- $\text{margin}(\mathbf{w}) = \min_i \frac{y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}}{\|\mathbf{w}\|}$
- Remember in regression: small norm solutions have low complexity!
 - Is this true for maximum margin classifiers?
 - what about classification with logistic loss $\sum_i \log(1 + \exp(-y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}))$?
 - how to do capacity control in maximum margin classifier learning?
- Some places $\min_i y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}$ referred as margin \rightarrow implicitly assumes normalization
 - $\min_i y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}$ is meaningless without knowing what $\|\mathbf{w}\|$ is!

Solutions of hard margin SVM

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Theorem: $\hat{\mathbf{w}} = \operatorname{span}\{\mathbf{x}^{(i)} : i = 1, 2, \dots, N\}$
i.e., $\exists\{\hat{\beta}_i : i = 1, 2, \dots, N\}$ such that $\hat{\mathbf{w}} = \sum_i \hat{\beta}_i \mathbf{x}^{(i)}$

Solutions of hard margin SVM

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Theorem: $\hat{\mathbf{w}} = \operatorname{span}\{\mathbf{x}^{(i)} : i = 1, 2, \dots, N\}$
i.e., $\exists \{\hat{\beta}_i : i = 1, 2, \dots, N\}$ such that $\hat{\mathbf{w}} = \sum_i \hat{\beta}_i \mathbf{x}^{(i)}$
- Denote $\mathcal{S} = \operatorname{span}\{\mathbf{x}^{(i)} : i = 1, 2, \dots, N\}$ and
 $\mathcal{S}^\perp = \{\mathbf{z} \in \mathbb{R}^d : \forall i, \mathbf{z} \cdot \mathbf{x}_i = 0\}$
 - For any $\mathbf{z} \in \mathbb{R}^d$, $\mathbf{z} = \mathbf{z}_\mathcal{S} + \mathbf{z}_{\mathcal{S}^\perp}$ s.t. $\mathbf{z}_\mathcal{S} \in \mathcal{S}$ and $\mathbf{z}_{\mathcal{S}^\perp} \in \mathcal{S}^\perp$
 - $\|\mathbf{z}\|^2 = \|\mathbf{z}_\mathcal{S}\|^2 + \|\mathbf{z}_{\mathcal{S}^\perp}\|^2$

Solutions of hard margin SVM

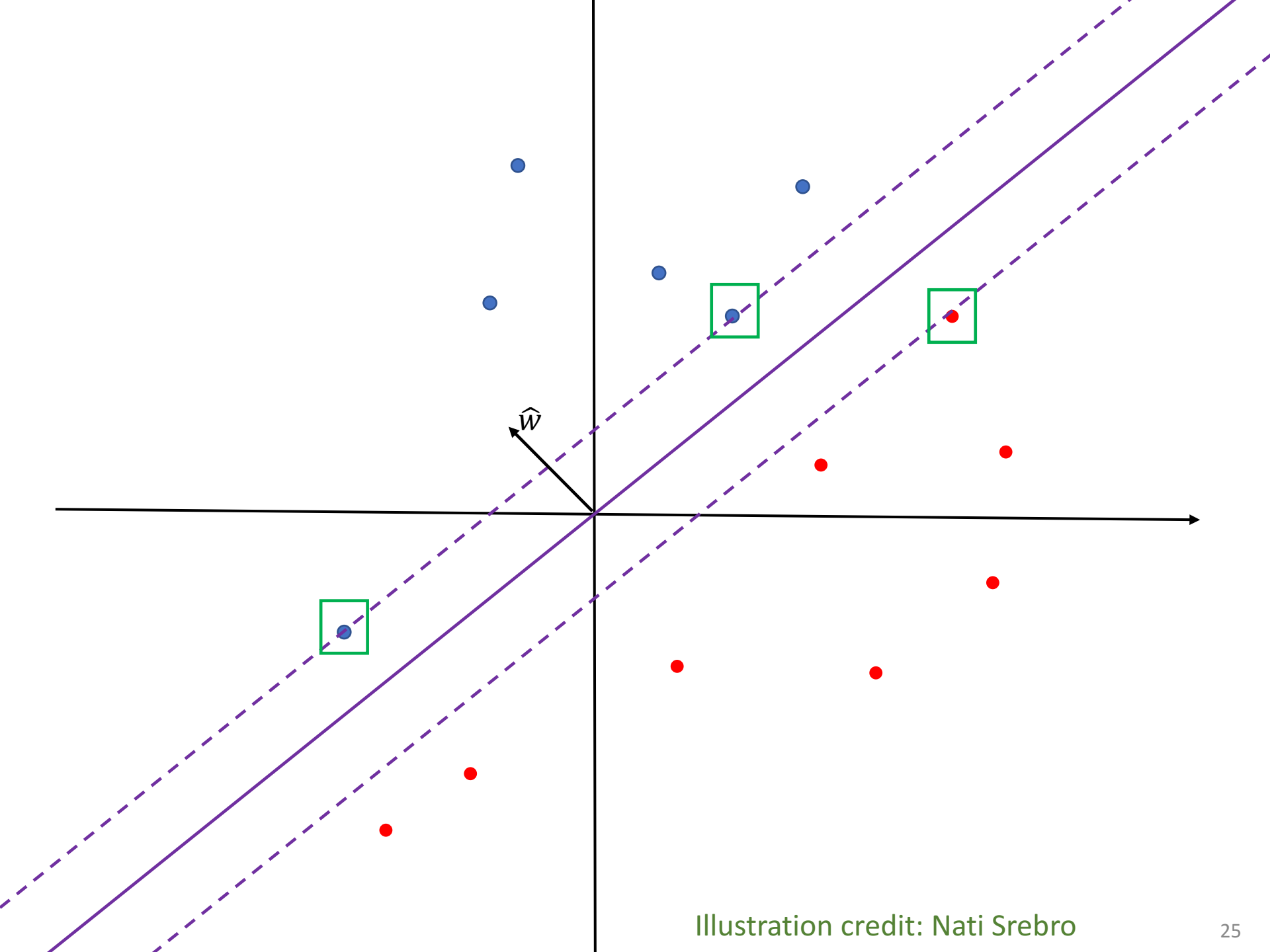
$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Theorem: $\hat{\mathbf{w}} = \operatorname{span}\{\mathbf{x}^{(i)} : i = 1, 2, \dots, N\}$
i.e., $\exists \{\hat{\beta}_i : i = 1, 2, \dots, N\}$ such that $\hat{\mathbf{w}} = \sum_i \hat{\beta}_i \mathbf{x}^{(i)}$
- Denote $\mathcal{S} = \operatorname{span}\{\mathbf{x}^{(i)} : i = 1, 2, \dots, N\}$ and
 $\mathcal{S}^\perp = \{\mathbf{z} \in \mathbb{R}^d : \forall i, \mathbf{z} \cdot \mathbf{x}_i = 0\}$
 - For any $\mathbf{z} \in \mathbb{R}^d$, $\mathbf{z} = \mathbf{z}_\mathcal{S} + \mathbf{z}_{\mathcal{S}^\perp}$ s.t. $\mathbf{z}_\mathcal{S} \in \mathcal{S}$ and $\mathbf{z}_{\mathcal{S}^\perp} \in \mathcal{S}^\perp$
 - $\|\mathbf{z}\|^2 = \|\mathbf{z}_\mathcal{S}\|^2 + \|\mathbf{z}_{\mathcal{S}^\perp}\|^2$
- Three step proof:
 1. Decompose $\hat{\mathbf{w}} = \hat{\mathbf{w}}_\mathcal{S} + \hat{\mathbf{w}}_{\mathcal{S}^\perp}$.
 2. $\min_i y^{(i)} \hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} \geq 1 \Rightarrow \min_i y^{(i)} \hat{\mathbf{w}}_\mathcal{S} \cdot \mathbf{x}^{(i)} \geq 1$
(because $\hat{\mathbf{w}}_{\mathcal{S}^\perp} \cdot \mathbf{x}^{(i)} = 0 \forall i$)
 3. if $\hat{\mathbf{w}}_{\mathcal{S}^\perp} \neq 0$, then $\|\hat{\mathbf{w}}_\mathcal{S}\| < \|\hat{\mathbf{w}}\|$

Representer Theorem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Theorem: $\hat{\mathbf{w}} = \operatorname{span}\{\mathbf{x}^{(i)}: i = 1, 2, \dots, N\}$ i.e.,
 $\exists \{\hat{\beta}_i: i = 1, 2, \dots, N\}$ such that $\hat{\mathbf{w}} = \sum_i \hat{\beta}_i \mathbf{x}^{(i)}$
 - Special case of representer theorem
- Theorem (ext): additionally, $\{\hat{\beta}_i\}$ also satisfies $\hat{\beta}_i = 0$ for all i such that $y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} > 1$
- Proof?: (animation next slide)



Representer Theorem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Theorem: $\exists \{\hat{\beta}_i : i = 1, 2, \dots, N\}$ such that $\hat{\mathbf{w}} = \sum_i \hat{\beta}_i \mathbf{x}^{(i)}$
 $\{\hat{\beta}_i\}$ also satisfies $\hat{\beta}_i = 0$ for all i such that $y^{(i)} \hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} > 1$
- $SV(\hat{\mathbf{w}}) = \{i : y^{(i)} \hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} = 1\}$ datapoints closest to $\hat{\mathbf{w}}$
 - called support vectors
 - hence support vector machine

$$\hat{\mathbf{w}} = \sum_{i \in SV(\hat{\mathbf{w}})} \hat{\beta}_i \mathbf{x}^{(i)}$$

Representer Theorem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Theorem: $\exists \{\hat{\beta}_i : i = 1, 2, \dots, N\}$ such that $\hat{\mathbf{w}} = \sum_i \hat{\beta}_i \mathbf{x}^{(i)}$
 $\{\hat{\beta}_i\}$ also satisfies $\hat{\beta}_i = 0$ for all i such that $y^{(i)} \hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} > 1$
- $SV(\hat{\mathbf{w}}) = \{i : y^{(i)} \hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} = 1\}$ datapoints closest to $\hat{\mathbf{w}}$
 - called support vectors
 - hence support vector machine

$$\hat{\mathbf{w}} = \sum_{i \in SV(\hat{\mathbf{w}})} \hat{\beta}_i \mathbf{x}^{(i)}$$

How do we get $\hat{\mathbf{w}}$?

Optimizing the SVM problem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad s.t., \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

1. Can do sub-gradient descent (next class)
2. Special case of quadratic program

$$\min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^\top \mathbf{P} \mathbf{z} + \mathbf{q}^\top \mathbf{z}$$
$$s.t. \quad \mathbf{G} \mathbf{z} \leq \mathbf{h}, \mathbf{A} \mathbf{z} = \mathbf{b}$$

- Change of variables $\hat{\mathbf{w}} = \sum_{i \in SV(\hat{\mathbf{w}})} \hat{\beta}_i \mathbf{x}^{(i)}$?
- Change of variables $\hat{\mathbf{w}} = \sum_{i=1}^N \hat{\beta}_i \mathbf{x}^{(i)}$!

Optimizing the SVM problem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.}, \quad y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i$$

- Change of variables $\mathbf{w} = \sum_{i=1}^N \beta_i \mathbf{x}^{(i)}$!

$$\equiv \min_{\{\beta_i\}} \sum_{i=1}^N \sum_{j=1}^n \beta_i \beta_j \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \quad \text{s.t.} \quad \sum_{j=1}^N \beta_j y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \geq 1 \quad \forall i$$

$$= \min_{\beta \in \mathbb{R}^N} \beta^\top \mathbf{G} \beta \quad \text{s.t.} \quad y^{(i)} (\mathbf{G} \beta)_i \geq 1 \quad \forall i$$

- $\mathbf{G} \in \mathbb{R}^{N \times N}$ with $G_{ij} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$ is called the **gram matrix**
- Convex program: quadratic programming

The Kernel

$$\begin{aligned} & \min_w \|w\|^2 \quad s.t. \quad y^{(i)} w \cdot x^{(i)} \geq 1 \quad \forall i \\ & \equiv \min_{\beta \in \mathbb{R}^N} \beta^\top G \beta \quad s.t. \quad y^{(i)} (G\beta)_i \geq 1 \quad \forall i \end{aligned}$$

- Optimization problem depends on $x^{(i)}$ only through the values of $G_{ij} = x^{(i)} \cdot x^{(j)}$ for $i, j \in [N]$.

The Kernel

$$\begin{aligned} \min_w \|w\|^2 \quad s.t. \quad y^{(i)} w \cdot x^{(i)} &\geq 1 \quad \forall i \\ \equiv \min_{\beta \in \mathbb{R}^N} \beta^\top G \beta \quad s.t. \quad y^{(i)} (G\beta)_i &\geq 1 \quad \forall i \end{aligned}$$

- Optimization problem depends on $x^{(i)}$ only through the values of $G_{ij} = x^{(i)} \cdot x^{(j)}$ for $i, j \in [N]$.
- What about prediction?

$$\hat{w} \cdot x = \sum_i \beta_i x^{(i)} \cdot x$$

The Kernel

$$\begin{aligned} \min_w \|w\|^2 \quad s.t. \quad y^{(i)} w \cdot x^{(i)} &\geq 1 \quad \forall i \\ \equiv \min_{\beta \in \mathbb{R}^N} \beta^\top G \beta \quad s.t. \quad y^{(i)} (G \beta)_i &\geq 1 \quad \forall i \end{aligned}$$

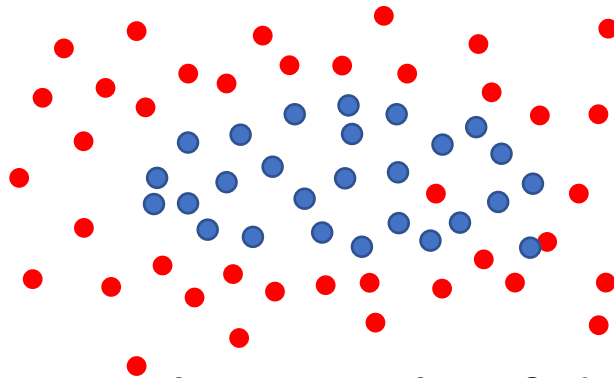
- Optimization problem depends on $x^{(i)}$ only through the values of $G_{ij} = x^{(i)} \cdot x^{(j)}$ for $i, j \in [N]$.
- What about prediction?

$$\hat{w} \cdot x = \sum_i \beta_i x^{(i)} \cdot x$$

- Function $K(x, x') = x \cdot x'$ is called the **Kernel**
- Learning non-linear classifiers using feature transformations, i.e., $f_w(x) = w \cdot \phi(x)$ for some $\phi(x)$
 - only thing we need to know is $K_\phi(x, x') = K(\phi(x), \phi(x'))$

Kernels As Prior Knowledge

- If we think that positive examples can (almost) be separated by some ellipse:



then we should use polynomials of degree 2

- A Kernel encodes a measure of *similarity* between objects. A bit like NN, except that it must be a valid inner product function.