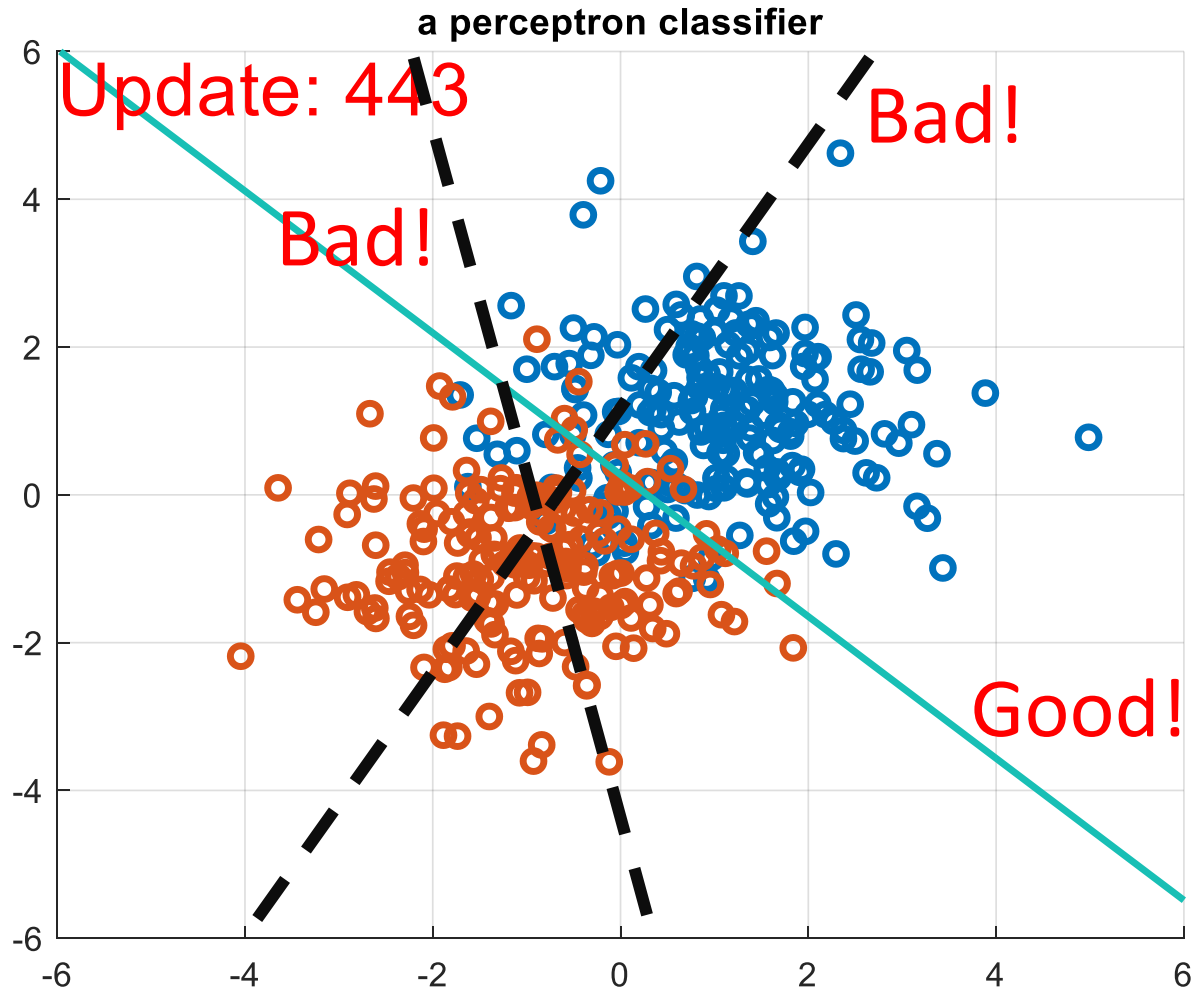# Support Vector Machines

Song Liu (song.liu@bristol.ac.uk)
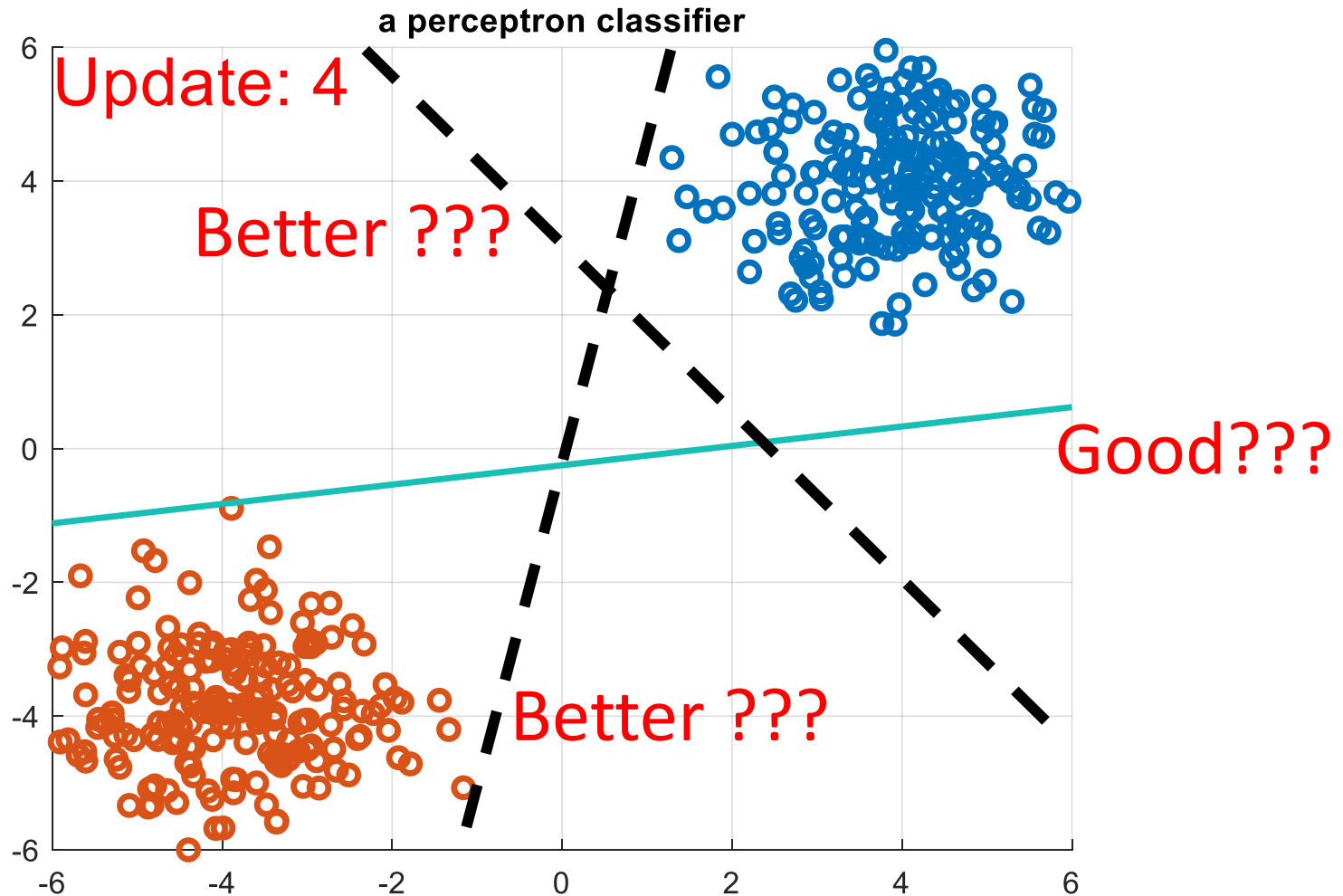
# Outlines

- Problem Support Vector Machine (SVM) tries to solve.

- Objective of SVM

- **Dual objective** of SVM

- Limitations of SVM

# Perceptron Classifier



a perceptron classifier

Update: 443

Bad!
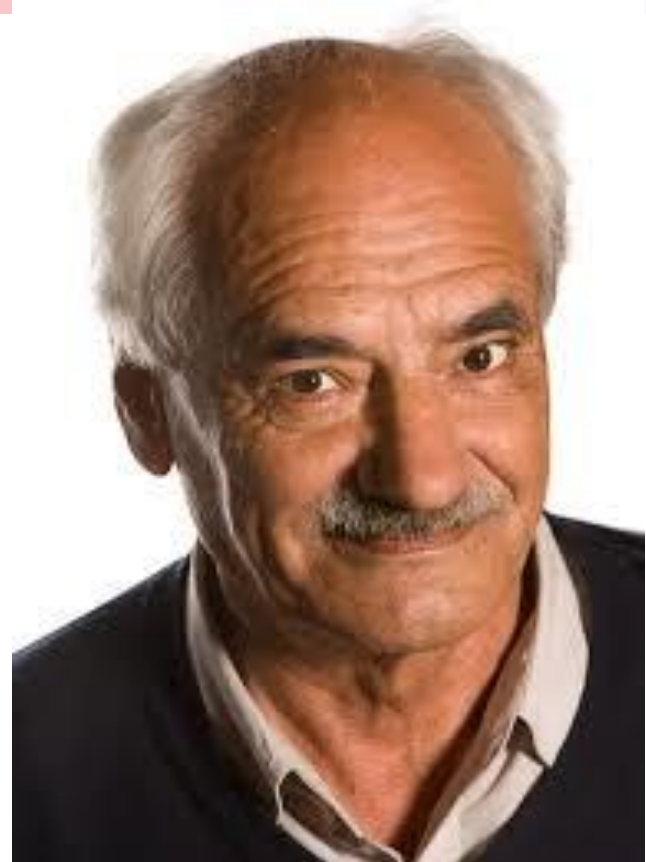
Bad!

Good!

# Perceptron Classifier



More than one "good" solutions!!

# What is
# The "Optimal"
# Decision Boundary?

# Vladimir Vapnik and Alexey Chervonenkis
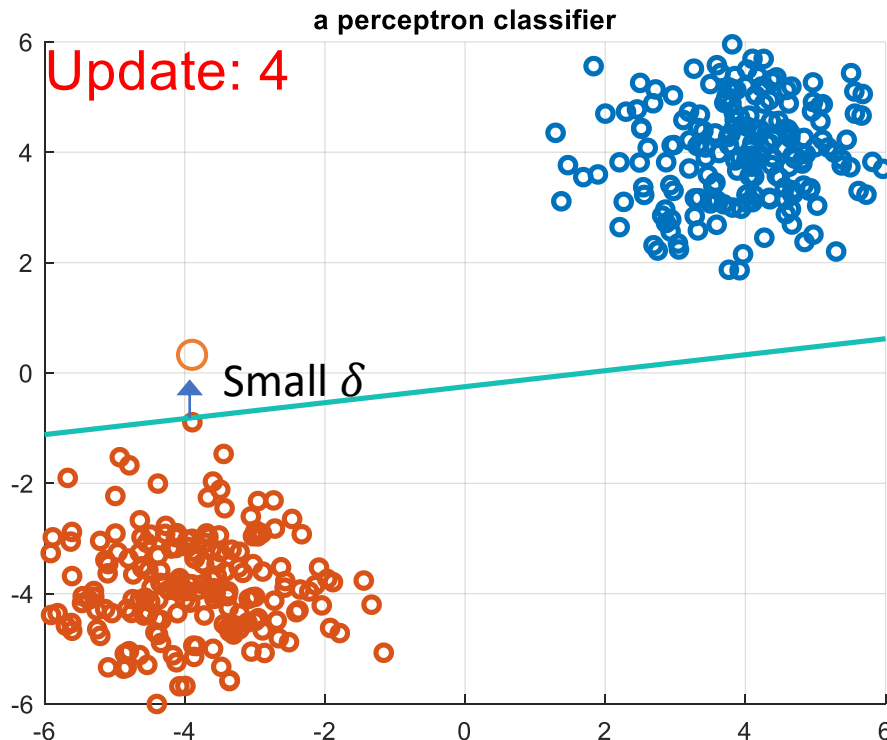


Vladimir Vapnik

Alexey Chervonenkis

**Contributions**: Statistical Learning Theory, Support Vector Machine

# The Error Margin

- Optimal decision boundary should minimize the error on **unseen datasets rather than training data.**
  - **a.k.a. Generalization Principle.**



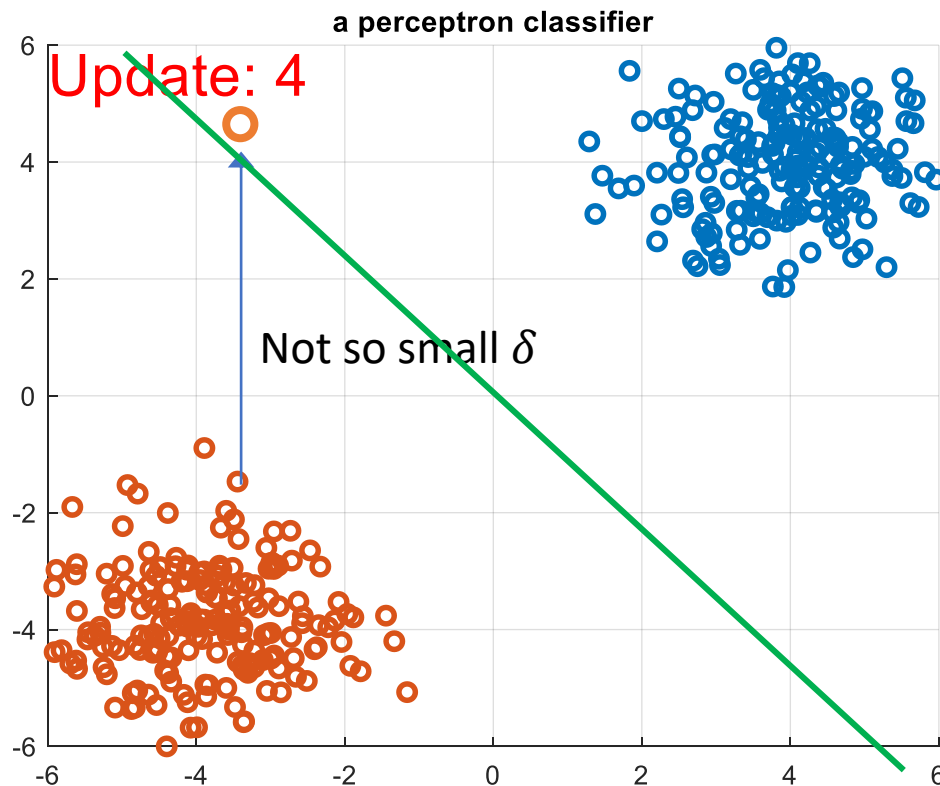a perceptron classifier

Update: 4

Small $\delta$

This is not a good decision boundary as a **small change** added to our data point would lead to **misclassification**.

Our decision boundary has a thin "**error margin**".

# The Error Margin

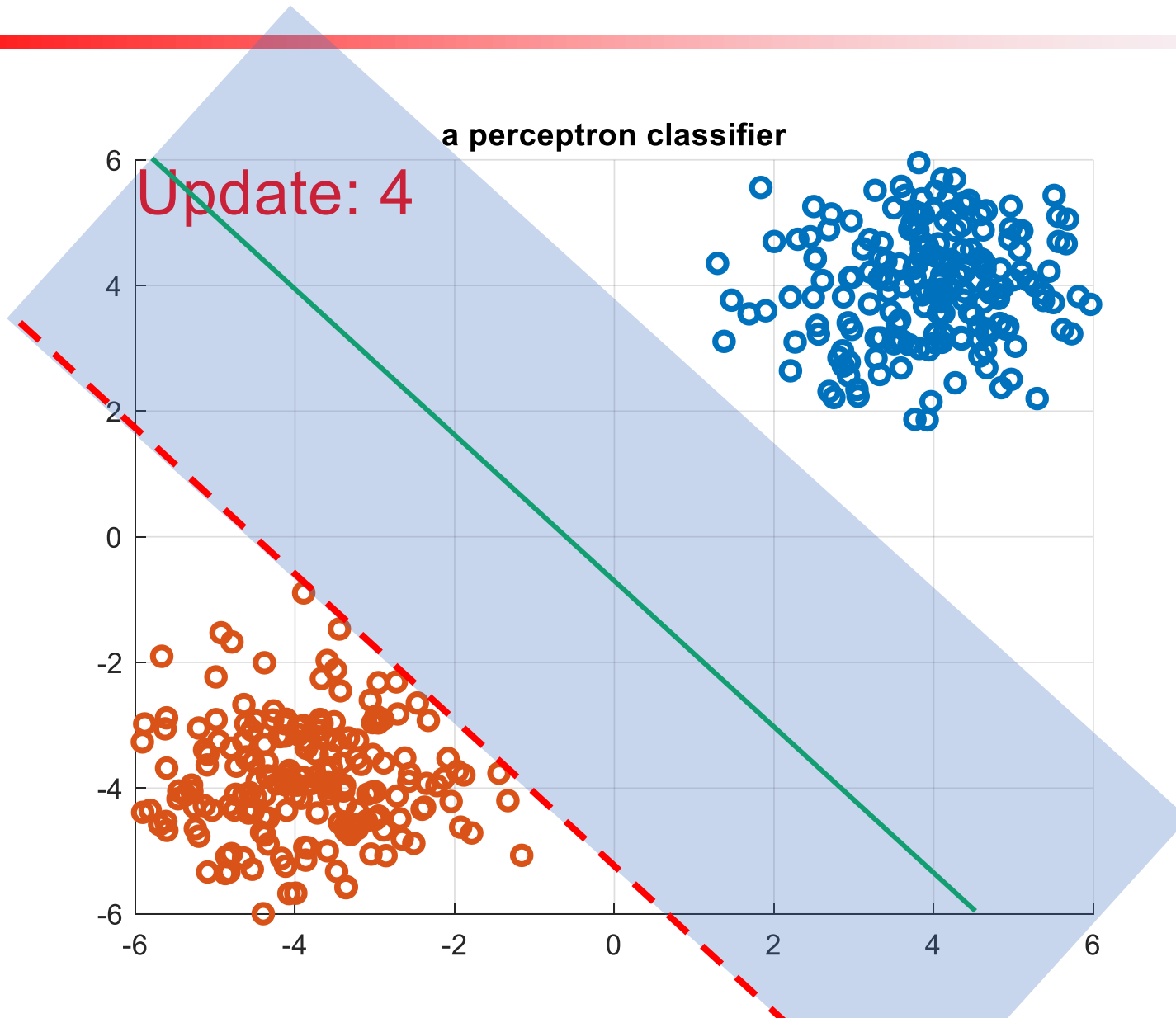- **Thin margin** is bad for generalization as some random unseen random points may easily "drift" to the other side of the decision boundary.



a perceptron classifier

This is a good decision boundary as small perturbations of our data points unlikely lead to **misclassification**.

Our decision boundary has a thick "**error margin**".

# Thick Margin



a perceptron classifier

Update: 4

# Thin Margin

# What is the "Optimal" Decision Boundary?

- If decision function is characterized by $f(x; w) = 0$, we want such a function to satisfy:

- 1. $\forall i, y_i = +, f(x_i; w) \geq 0$
- 2. $\forall i, y_i = -, f(x_i; w) \leq 0$

- 3. The error margin of $f(x; w)$ is as **THICK** as possible!

- How do you quantify above conditions?

# Margin of Linear Model

- Suppose $f(x; w) := \langle w', x \rangle + w_0$



$f(x; w) \geq 1$

Thickness of Margin: $\frac{1}{||w'||}$!

**Verify this!**

$\frac{1}{||w'||}$

$f(x; w) = 1$

$f(x; w) \leq -1$

$f(x; w) = 0$

$f(x; w) = -1$

# Maximal Margin Classifier

- Maximize the Error Margin

- datapoints are on the right side of the margin.

$f(\boldsymbol{x}; \boldsymbol{w}) \geq 1$

- $\Leftrightarrow$

$$\frac{1}{||\boldsymbol{w}'||}$$

$f(\boldsymbol{x}; \boldsymbol{w}) \leq -1$

- Maximize $\dfrac{1}{||\boldsymbol{w}'||}$

- and maintain $\forall_i, y_i = +, f(\boldsymbol{x}_i; \boldsymbol{w}) \geq 1,$
  $\forall_i, y_i = -, f(\boldsymbol{x}_i; \boldsymbol{w}) \leq -1$

# Maximal Margin Classifier

- Maximize $\dfrac{1}{||\boldsymbol{w}'||}$

- and maintain $\forall_i, y_i = +, f(\boldsymbol{x}_i; \boldsymbol{w}) \geq 1,$
$$\forall_i, y_i = -, f(\boldsymbol{x}_i; \boldsymbol{w}) \leq -1$$

- $\Leftrightarrow$

- Minimize $||\boldsymbol{w}'||^2$
- Subject to $\forall_i, y_i f(\boldsymbol{x}_i; \boldsymbol{w}) \geq 1,$

- **This is a constrained minimization!**
- **Unlike LS, which is an unconstrained minimization.**

# Soft-margin Classifiers

- In many cases, the dataset is not **linearly** separable.



A margin without any data points cannot be constructed due to the overlapping-ness of two classes!

# Soft-margin Classifiers

- We allow our $f$ make some errors!

These misclassified points will have $f(\boldsymbol{x}; \boldsymbol{w})y \leq 1$!

**Constraint not satisfied!**

$$f(\boldsymbol{x}; \boldsymbol{w}) = 1$$

$$f(\boldsymbol{x}; \boldsymbol{w}) = 0$$

$$f(\boldsymbol{x}; \boldsymbol{w}) = -1$$

# Soft-Margin Classifier

- Minimize $||\boldsymbol{w}'||^2 + \sum_i \epsilon_i$
- Subject to $\forall_i, y_i f(\boldsymbol{x}_i; \boldsymbol{w}) + \epsilon_i \geq 1, \epsilon_i \geq 0$

- For each $\boldsymbol{x}_i$, we hope $y_i f(\boldsymbol{x}_i; \boldsymbol{w})$ can be outside of the margin after some small positive "compensation" $\epsilon_i$.

- At the same time, we want such "compensation" is as small as possible, i.e., the classifier makes as few mistakes as possible.

- The solution for $\epsilon$ is sparse. Why?

# Soft-Margin Classifier

- Formally, the soft-margin classifier

- $\min_{\boldsymbol{w}, \boldsymbol{\epsilon}} ||\boldsymbol{w}'||^2 + \sum_i \epsilon_i$

- Subject to $\forall_i, y_i(\langle \boldsymbol{w}', \boldsymbol{x} \rangle + w_0) + \epsilon_i \geq 1, \epsilon_i \geq 0$

- It turns out, if $f(\boldsymbol{x}; \boldsymbol{w})$ is a linear model,

- **Soft-Margin Classifier** is a **convex** minimization problem.

- **Every local minimum is a global minimum.**

"convex"

# The Lagrangian Dual

- Solving constrained problem can be rather complicated.
- **Lagrangian Dual**: a technique transforms constrained problem into unconstrained problem.

- For a constrained problem,
- $\min_{\theta} f(\theta)$ subject to $g_i(\theta) \leq 0$
- We can construct a **Lagrangian** $l(\boldsymbol{\lambda})$:
- $l(\boldsymbol{\lambda}) := \min_{\theta} f(\theta) + \sum_i \lambda_i \, g_i(\theta),$
- $\lambda_i \geq 0$ are called **Lagrangian multipliers**.
- **PRML Appendix E.**

# The Lagrangian Dual

- Under regularity conditions*, maximizing $l(\boldsymbol{\lambda})$ w.r.t. $\boldsymbol{\lambda}$ would allow us to recover the optimal solutions in the original constrained minimization problem.

To maximize $l(\boldsymbol{\lambda})$, do the following 4 steps:

- 1. Write down $l(\boldsymbol{\lambda})$ for soft-margin classifier:

- $l(\boldsymbol{\lambda}) :=$

$$\min_{\boldsymbol{w},\boldsymbol{\epsilon}} ||\boldsymbol{w}'||^2 + \sum_i \epsilon_i - \lambda_i[y_i(\langle \boldsymbol{w}', \boldsymbol{x}\rangle + w_0) + \epsilon_i - 1] - \lambda_i'\epsilon_i$$

- 2. Derive optimality condition w.r.t. $\boldsymbol{w}$ and $\boldsymbol{\epsilon}$:

- $\boldsymbol{w}' = \frac{\sum_{i=1} \lambda_i y_i \boldsymbol{x}_i}{2}, \sum_{i=1} \lambda_i y_i = 0, \lambda_i + \lambda_i' = 1,$

Verify this!

# The Lagrangian Dual

- Using optimality conditions:
- $\boldsymbol{w}' = \frac{\sum_{i=1} \lambda_i y_i \boldsymbol{x}_i}{2}, \lambda_i + \lambda_i' = 1, \sum_{i=1} \lambda_i y_i = 0$

- 3. Rewrite $l(\lambda) = -\frac{\tilde{\lambda}^\top \boldsymbol{X}^\top \boldsymbol{X} \tilde{\lambda}}{4} + \langle \boldsymbol{\lambda}, \boldsymbol{1} \rangle,$ Verify it!

  $\boldsymbol{X} = [\boldsymbol{x}_1 \dots \boldsymbol{x}_n] \in R^{d \times n}, \tilde{\boldsymbol{\lambda}} := [\lambda_1 \cdot y_1 \dots \lambda_n \cdot y_n]$

- 4. Maximize $l(\lambda)$ w.r.t. $\boldsymbol{\lambda}$ under constraints:
  - $0 < \lambda_i < 1$
  - $\sum_{i=1} \lambda_i y_i = 0$

Needed to make sure the optimality of the **original problem**

# Soft-margin Classifier (Dual)

- $\max_{\boldsymbol{\lambda}} -\dfrac{\tilde{\lambda}^{\top} X^{\top} X \tilde{\lambda}}{4} + \langle \boldsymbol{\lambda}, \mathbf{1} \rangle$

- Subject to

- $0 \leq \lambda_i \leq 1, \sum_{i=1} \lambda_i y_i = 0$

- Recover $\widehat{\boldsymbol{w}}' := \dfrac{\sum_{i=1} \widehat{\lambda}_i y_i \boldsymbol{x}_i}{2}$ using optimality condition.
- Put $\widehat{\boldsymbol{w}}'$ back in the original problem and solve for $\widehat{w}_0$.

- We obtain $\widehat{\boldsymbol{W}}$ using Lagragian multipliers $\boldsymbol{\lambda}$.

# Soft-margin Classifier (Dual)
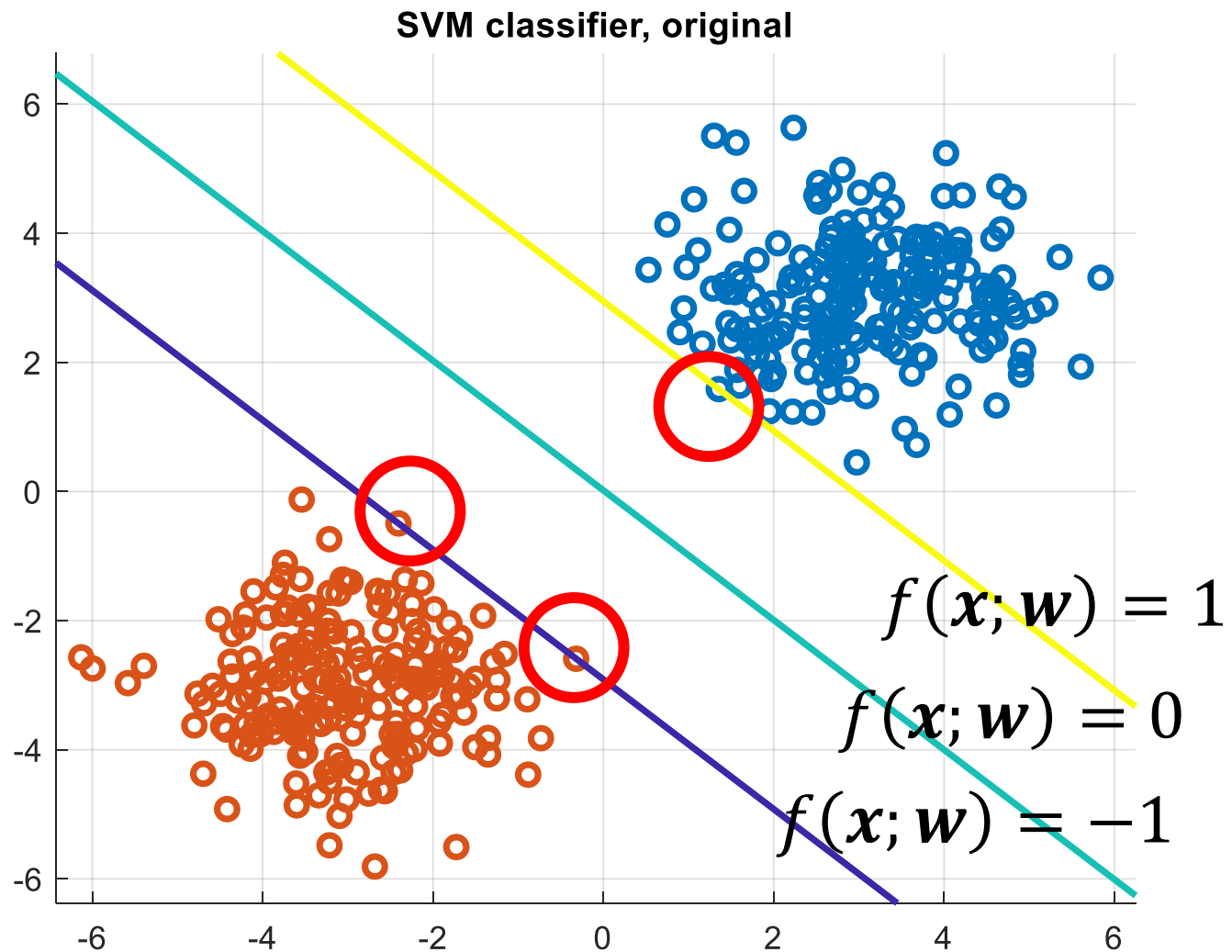
- $\max\limits_{\boldsymbol{\lambda}} -\dfrac{\tilde{\lambda}^\top X^\top X \tilde{\lambda}}{4} + \langle \boldsymbol{\lambda}, \mathbf{1} \rangle$

- Subject to $0 \leq \lambda_i \leq 1, \sum_{i=1} \lambda_i y_i = 0$

- Our input data $\{\boldsymbol{x}_i\}$ **only appear at** $X^\top X$

- **Let** $\boldsymbol{K} = \boldsymbol{X}^\top \boldsymbol{X}$**, then** $K^{(i,j)} = \langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle$

- Instead of using the inner product, we can use kernel functions $k(\boldsymbol{x}_i, \boldsymbol{x}_j)$ to perform training and prediction.

- **Homework**: write down decision function $f(\boldsymbol{x}; \boldsymbol{w})$ using kernel function $k, w_0$ and dual variable $\boldsymbol{\lambda}$.

# Original vs. Dual Problem

- $\min\limits_{\boldsymbol{w},\boldsymbol{\epsilon}} ||\boldsymbol{w}'||^2 + \sum_i \epsilon_i$

- Subject to $\forall_i$,
  $y_i(\langle \boldsymbol{w}', \boldsymbol{x}\rangle + w_0) + \epsilon_i \geq 1,$
  $\epsilon_i > 0$

- Complex Constraints

- Quadratic w.r.t. $\boldsymbol{w} \in R^d$

- Slow when $d$ is large

- $\max\limits_{\boldsymbol{\lambda}} -\dfrac{\tilde{\boldsymbol{\lambda}}^\top X^\top X \tilde{\boldsymbol{\lambda}}}{4} + \langle \boldsymbol{\lambda}, \mathbf{1}\rangle$

- Subject to

- $0 \leq \lambda_i \leq 1$

- $\sum_{i=1} \lambda_i y_i = 0$

- Simpler Constraints

- Quadratic w.r.t. $\boldsymbol{\lambda} \in R^n$

- Slow when $n$ is large

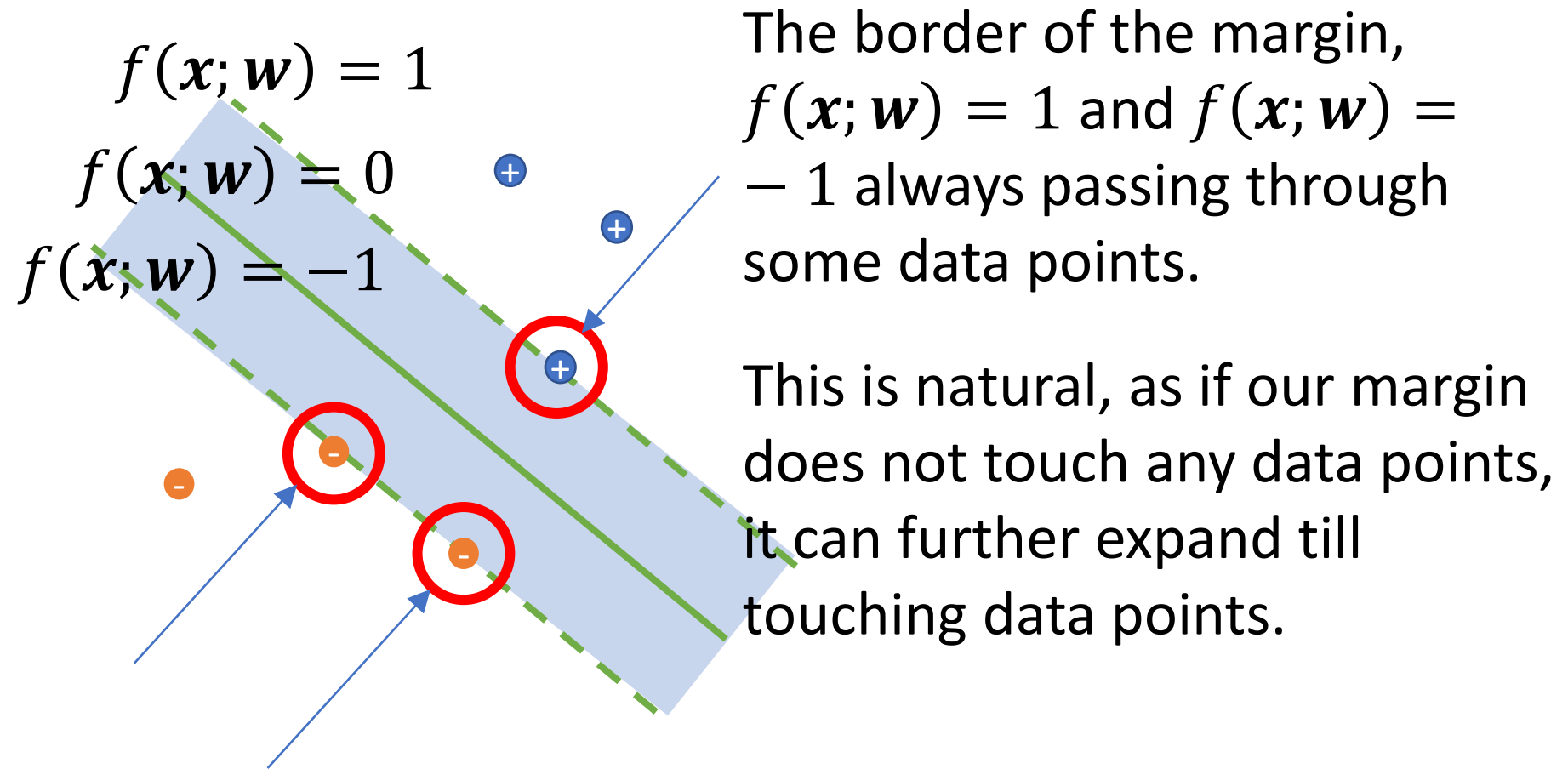- Can use kernel!

# Toy Example



**SVM classifier, original**

$f(x; w) = 1$

$f(x; w) = 0$

$f(x; w) = -1$

# Toy Example



$f(\boldsymbol{x}; \boldsymbol{w}) = 1$

$f(\boldsymbol{x}; \boldsymbol{w}) = 0$

$f(\boldsymbol{x}; \boldsymbol{w}) = -1$

# "Support Vectors"

$$f(x; w) = 1$$

$$f(x; w) = 0$$

$$f(x; w) = -1$$

The border of the margin, $f(x; w) = 1$ and $f(x; w) = -1$ always passing through some data points.

This is natural, as if our margin does not touch any data points, it can further expand till touching data points.
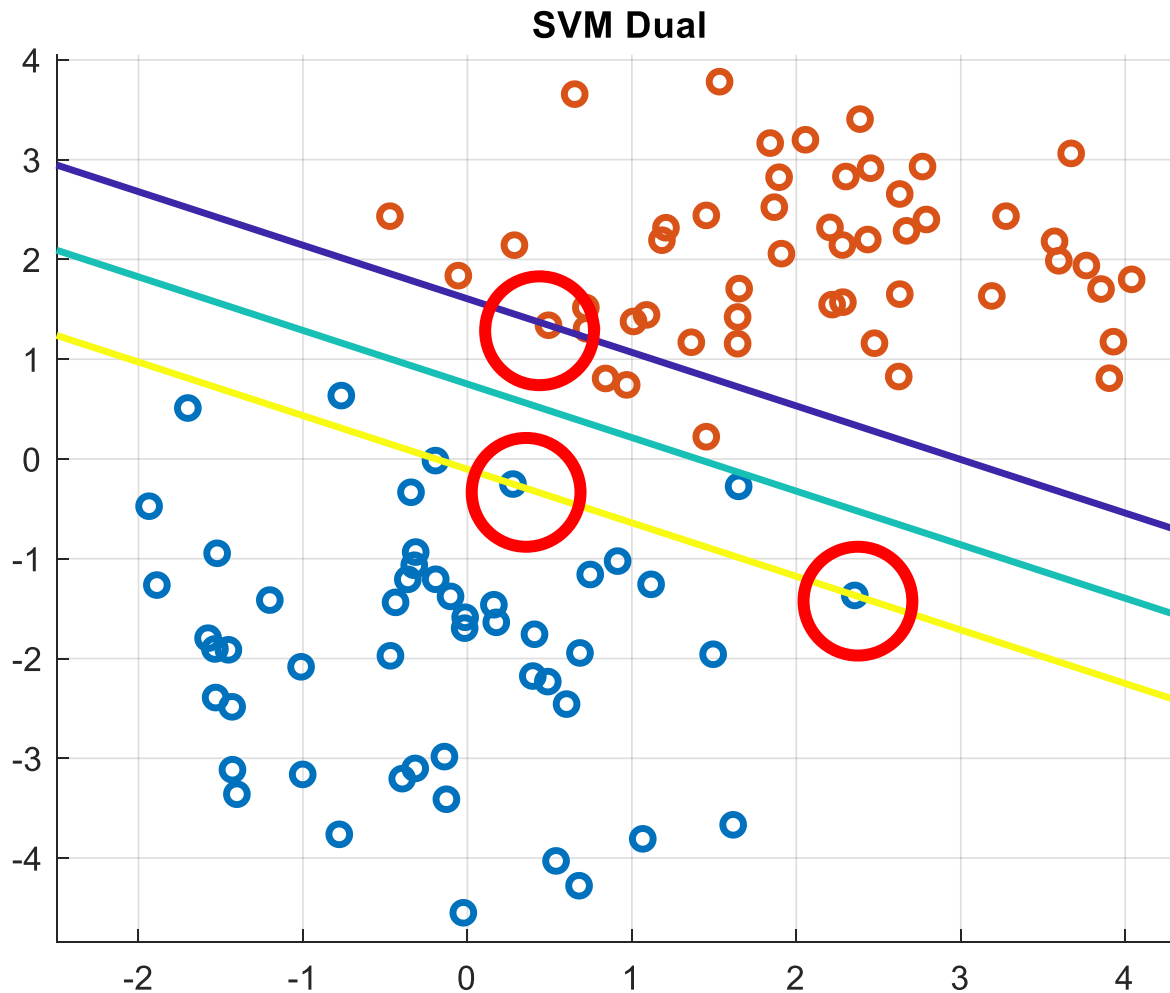
These points as if they were resisting the expansion of the margin, are appropriately called "support vectors".

# Toy, Soft-margin, Dual



**SVM Dual**

$w$ by solving Original:

$$w' = [-0.6287 \quad -1.1708]$$
$$w_0 = 0.8797$$

$w$ recovered from $\lambda$:

$$w' = [-0.6287 \quad -1.1708]$$
$$w_0 = 0.8797$$

# Limitations of SVM

- SVM is **not** a probabilistic classifier
  - cannot be integrated with probabilistic classification models (generative or discriminative)
  - The decision function lacks interpretability.

- Computational cost of SVM is high
  - Either original and dual requires solving constrained optimization.
  - Many other classifier, e.g. Logistic Regression, solves unconstrained optimization.

- Multi-class SVM classification is non-trivial.
  - SVM is motivated by the geometry of binary classification.

# Conclusion

- SVM is motivated via "Maximum Margin" principle.

- Soft-margin SVM can classify linearly inseparable data.

- Dual of SVM can be derived using Lagrangian.

- SVM is not a probabilistic classifier.

# Homework

- Derive the optimality condition in $l(\boldsymbol{\lambda})$ for $\boldsymbol{w}$ and $\boldsymbol{\epsilon}$.

- Represent prediction function $f(\boldsymbol{x}; \boldsymbol{w})$ using dual parameter $\boldsymbol{\lambda}$, kernel function $k$ and bias $w_0$.

# Computing Lab