

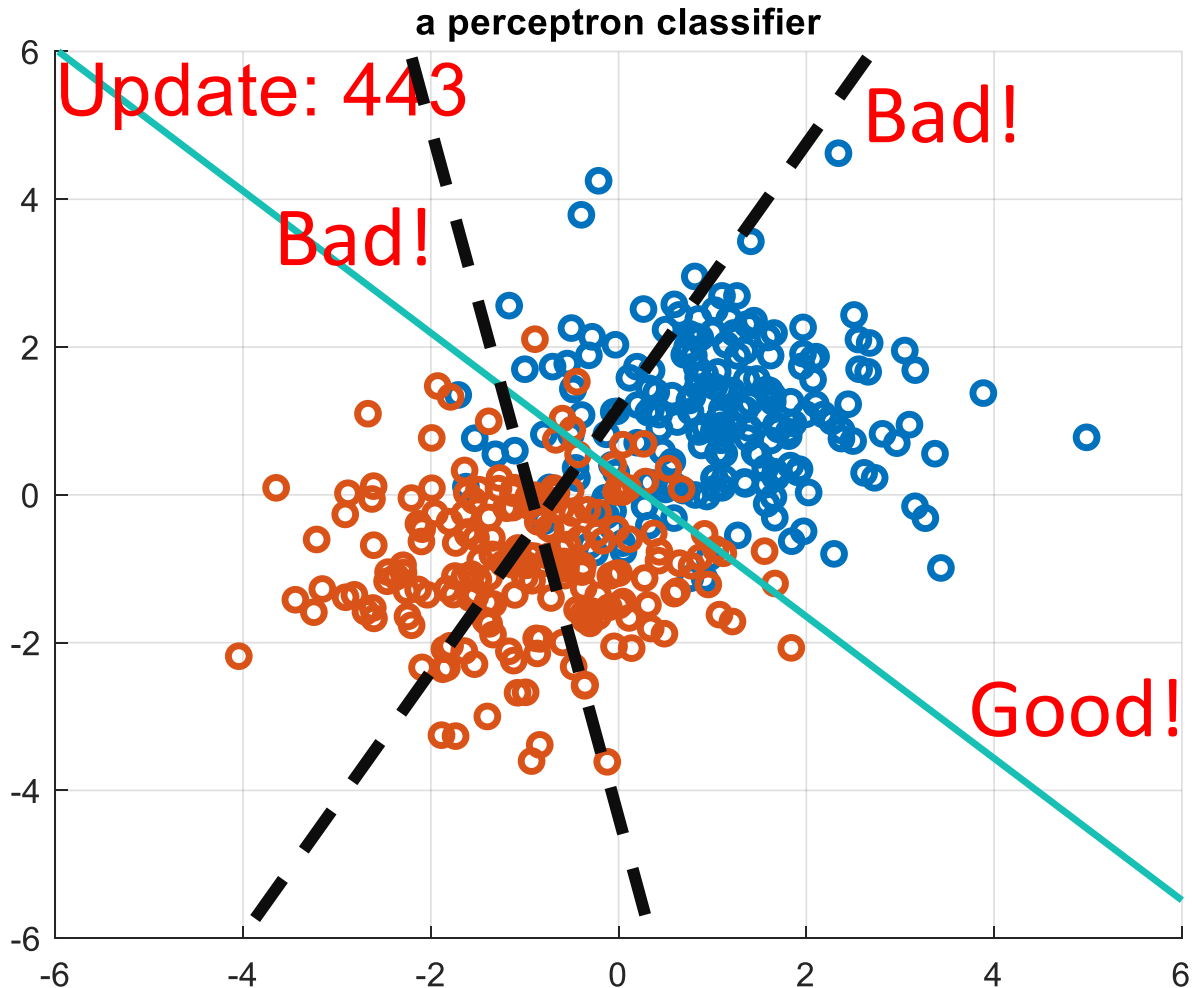
Support Vector Machines

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

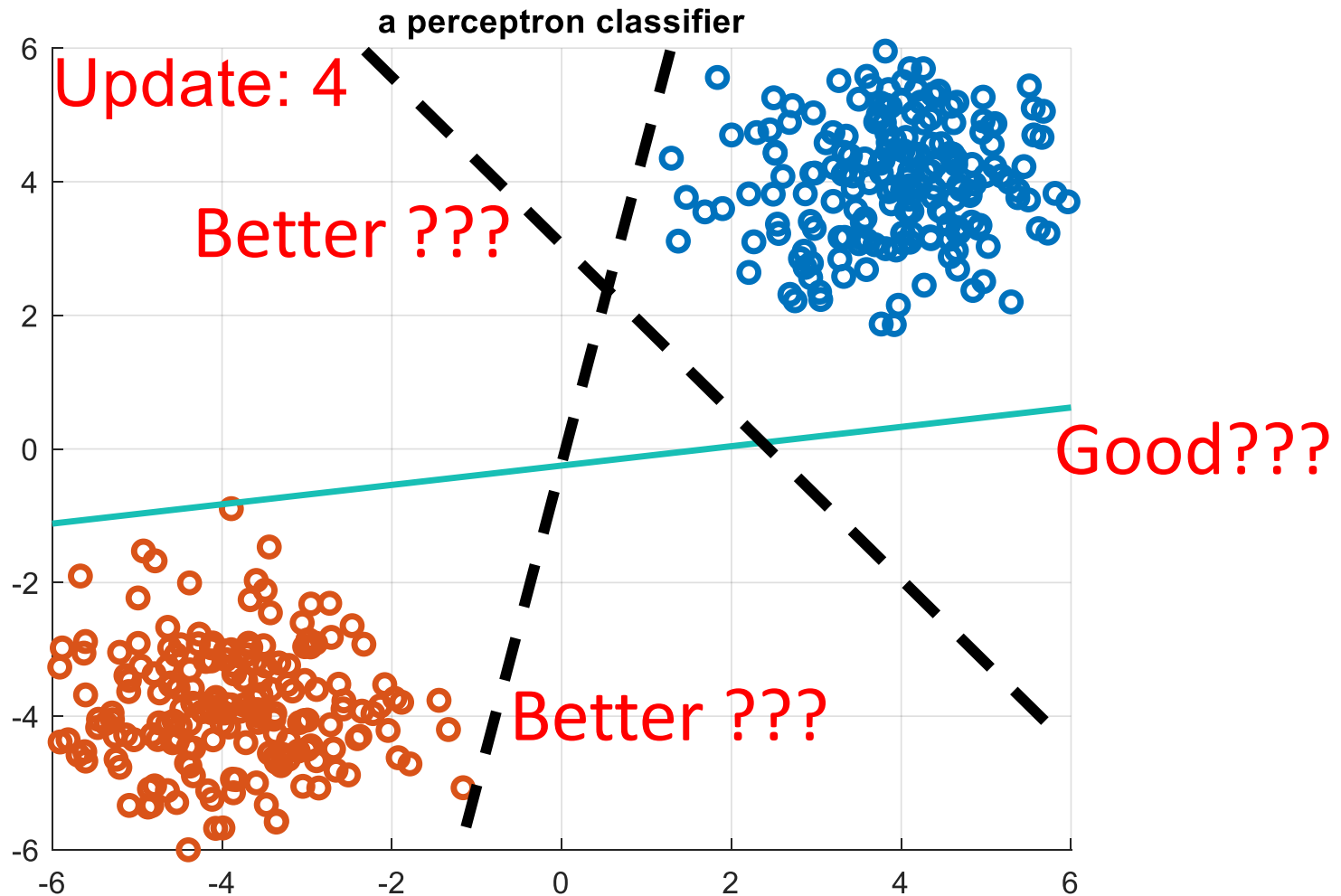
Outlines

- Problem Support Vector Machine (SVM) tries to solve.
- Objective of SVM
- **Dual objective** of SVM
- Limitations of SVM

Perceptron Classifier



Perceptron Classifier



More than one “good” solutions!!

What is
The “Optimal” Decision
Boundary in binary
classification?

Vladimir Vapnik and Alexey Chervonenkis



Vladimir Vapnik

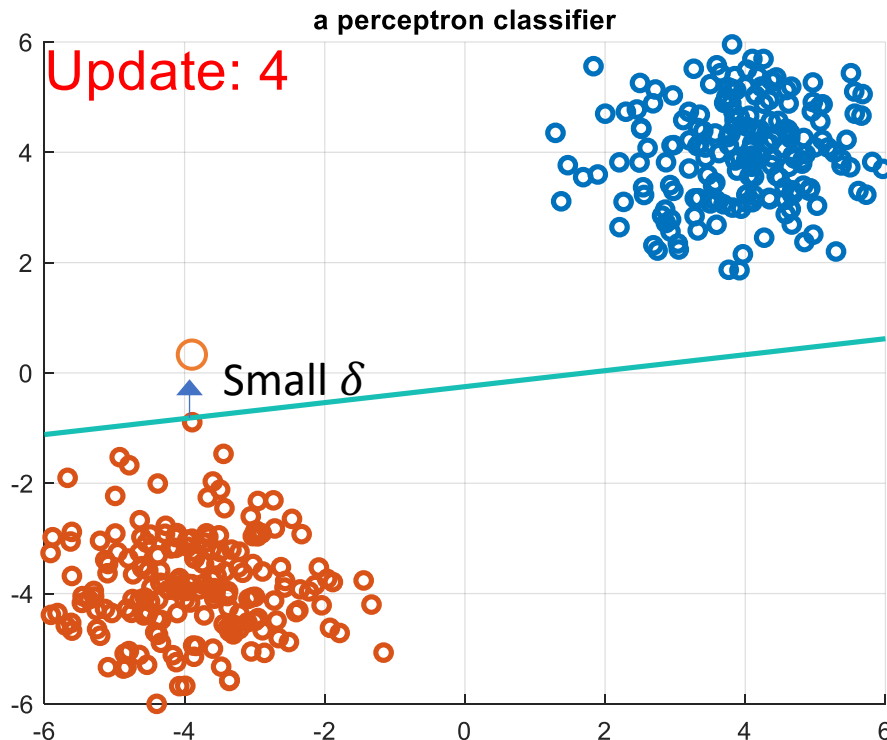


Alexey Chervonenkis

Contributions: Statistical Learning Theory, Support Vector Machines

The Error Margin

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

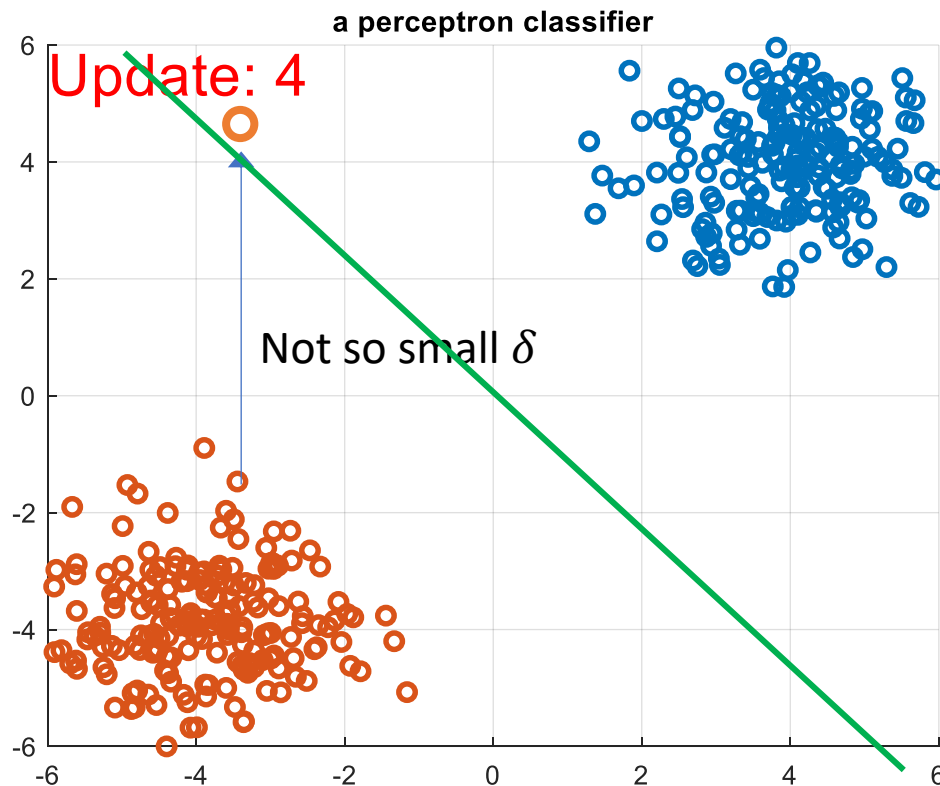


This is not a good decision boundary as a **small change** added to our data point would lead to **misclassification**.

Our decision boundary has a thin “**error margin**”.

The Error Margin

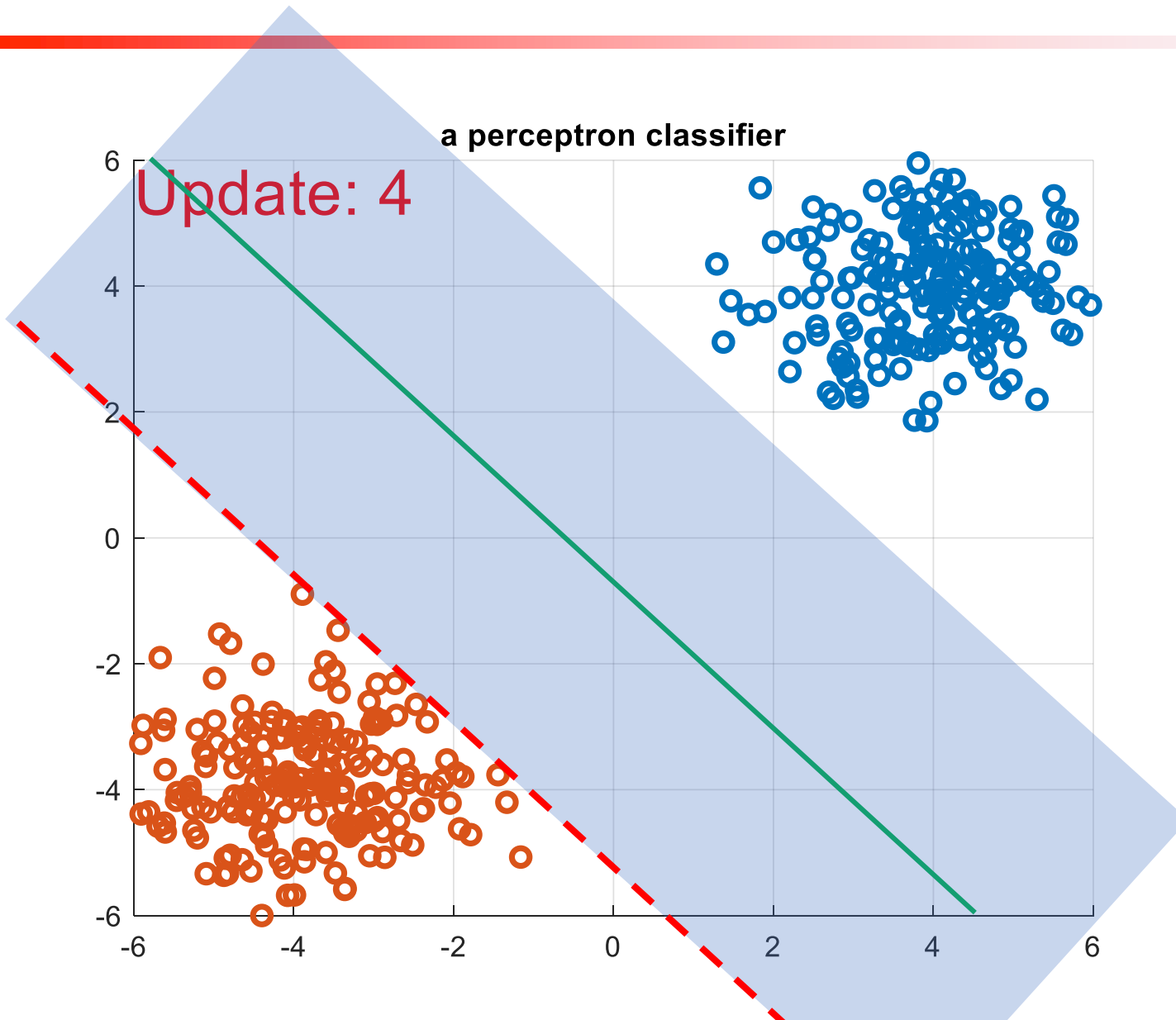
- **Thin margin** is bad for generalization as some random unseen data points may easily “drift” to the other side of the decision boundary.



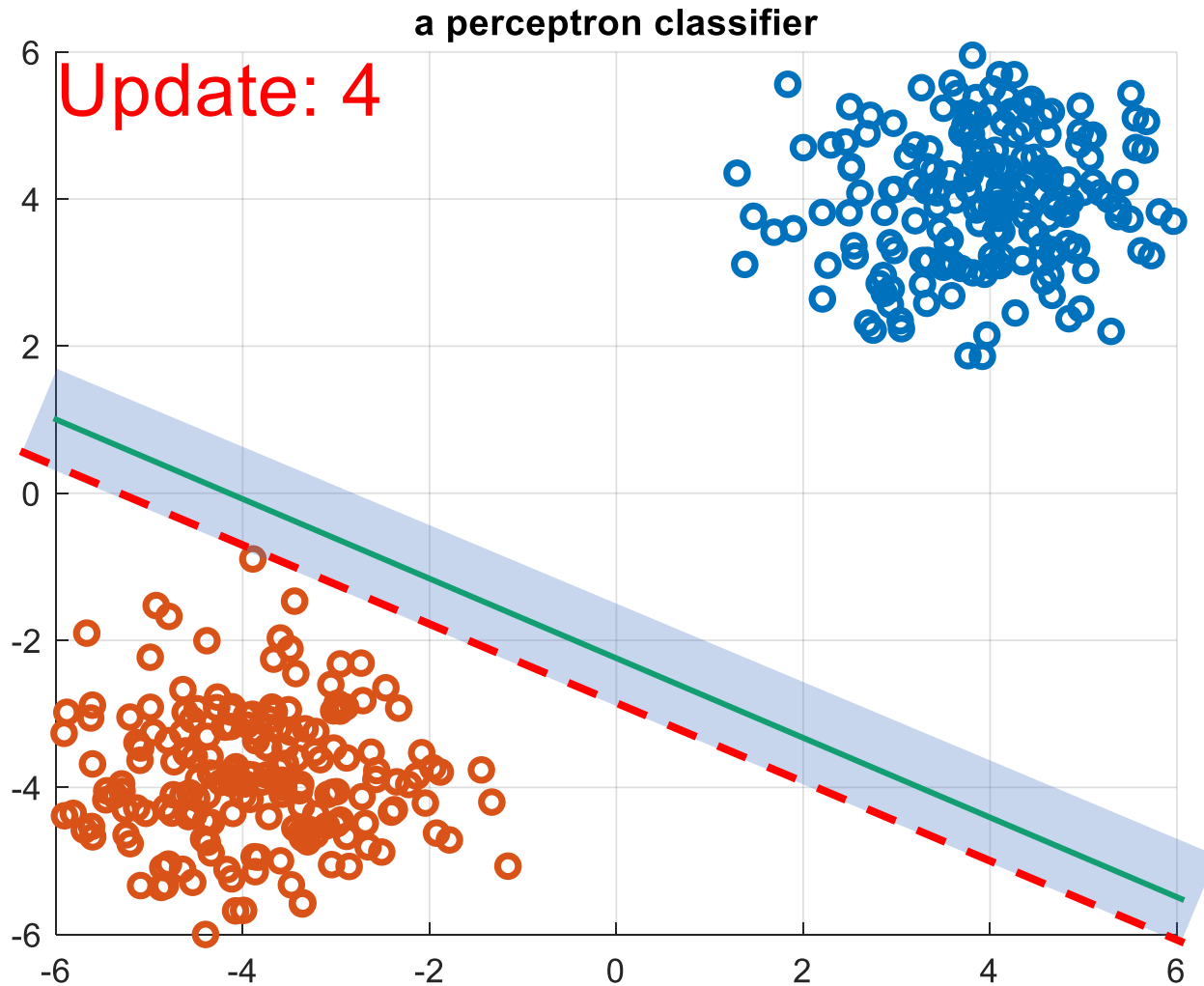
This is a good decision boundary as small perturbations of our data points unlikely lead to **misclassification**.

Our decision boundary has a thick “**error margin**”.

Thick Error Margin



Thin Error Margin

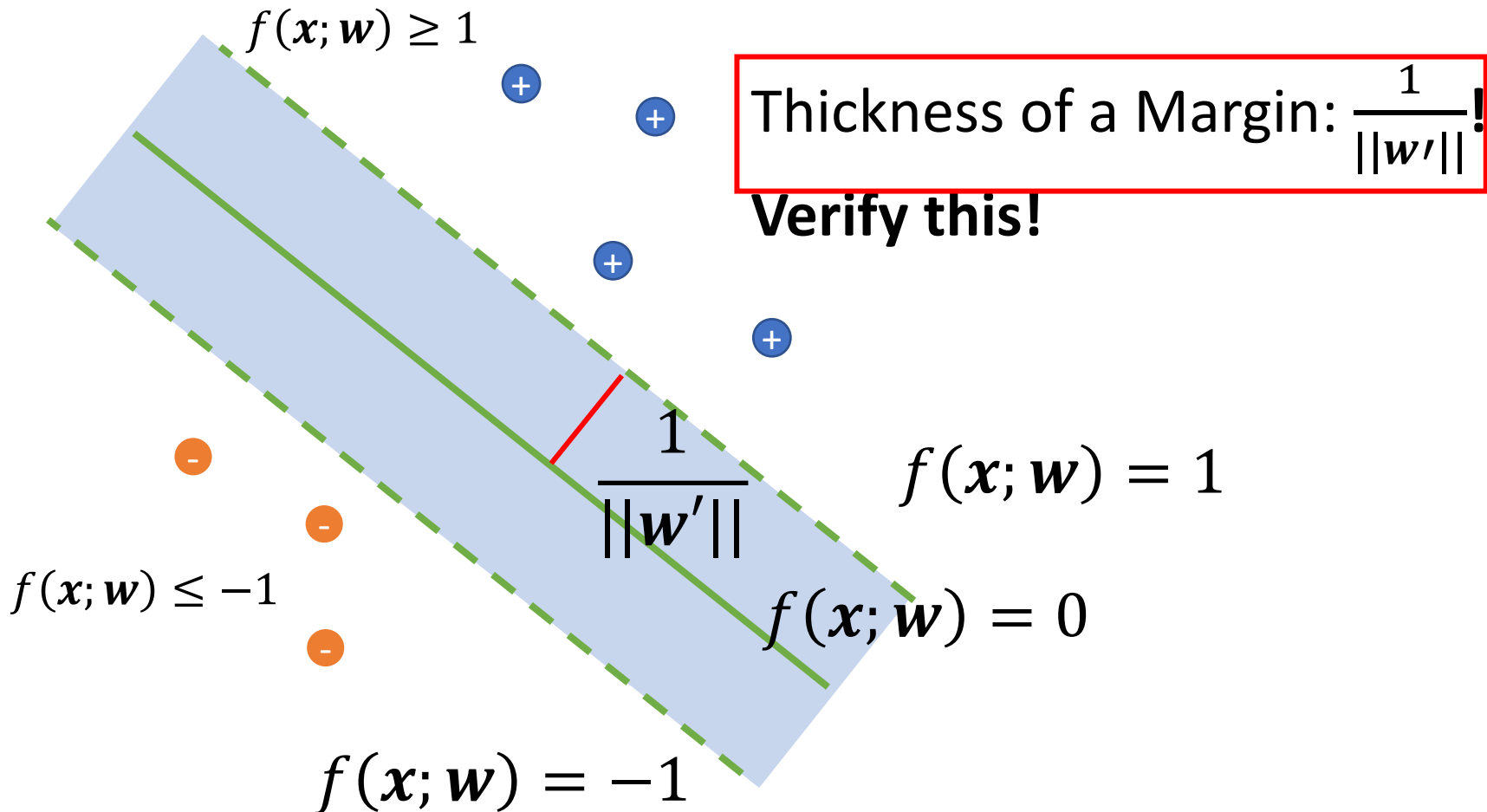


What is the “Optimal” Decision Boundary?

- If decision boundary is characterized by $f(\mathbf{x}; \mathbf{w}) = 0$, we have the following criteria:
 - 1. $\forall i, y_i = +, f(\mathbf{x}_i; \mathbf{w}) \geq 0$
 - 2. $\forall i, y_i = -, f(\mathbf{x}_i; \mathbf{w}) \leq 0$
- 3. The error margin of $f(\mathbf{x}; \mathbf{w})$ should be as **THICK** as possible!
- How do you quantify the above criteria?

Margin of Linear Model

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

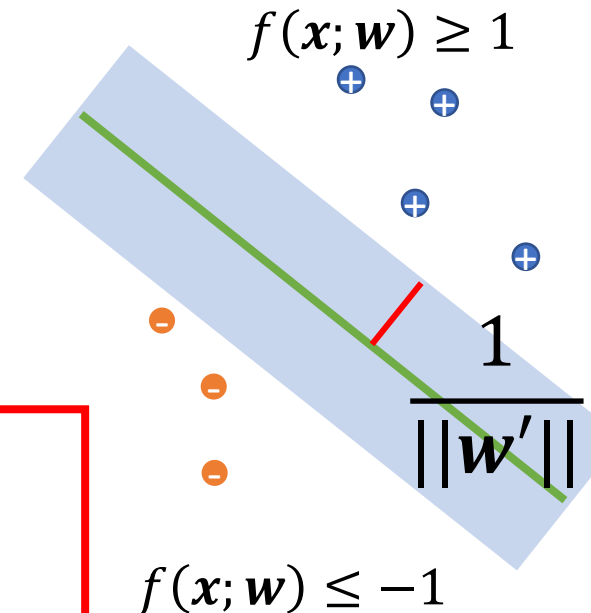


Maximal Margin Classifier

- Maximize the Width of a Margin
- Keep datapoints on the correct side of the margin.

• \Leftrightarrow

- Maximize $\frac{1}{\|w'\|}$
- and maintain $\forall_i, y_i = +, f(x_i; w) \geq 1,$
 $\forall_i, y_i = -, f(x_i; w) \leq -1$

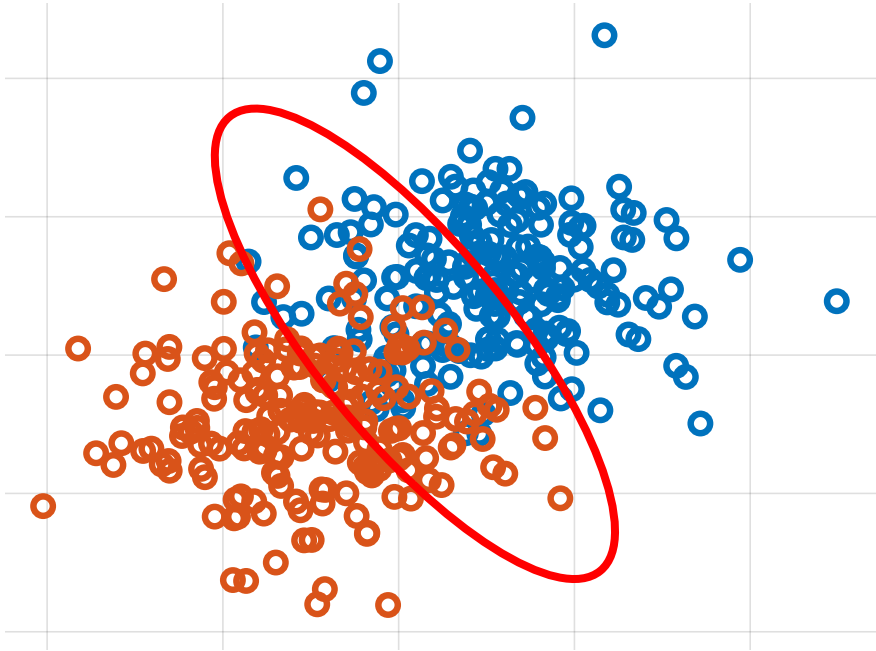


Maximal Margin Classifier

- Maximize $\frac{1}{||\mathbf{w}'||}$
- and maintain $\forall_i, y_i = +, f(\mathbf{x}_i; \mathbf{w}) \geq 1,$
 $\forall_i, y_i = -, f(\mathbf{x}_i; \mathbf{w}) \leq -1$
- \Leftrightarrow
- Minimize $||\mathbf{w}'||^2$
- Subject to $\forall_i, y_i f(\mathbf{x}_i; \mathbf{w}) \geq 1,$
- **This is a constrained minimization!**
- **Unlike LS and Logistic regression, which are both unconstrained minimizations.**

Soft-margin Classifiers

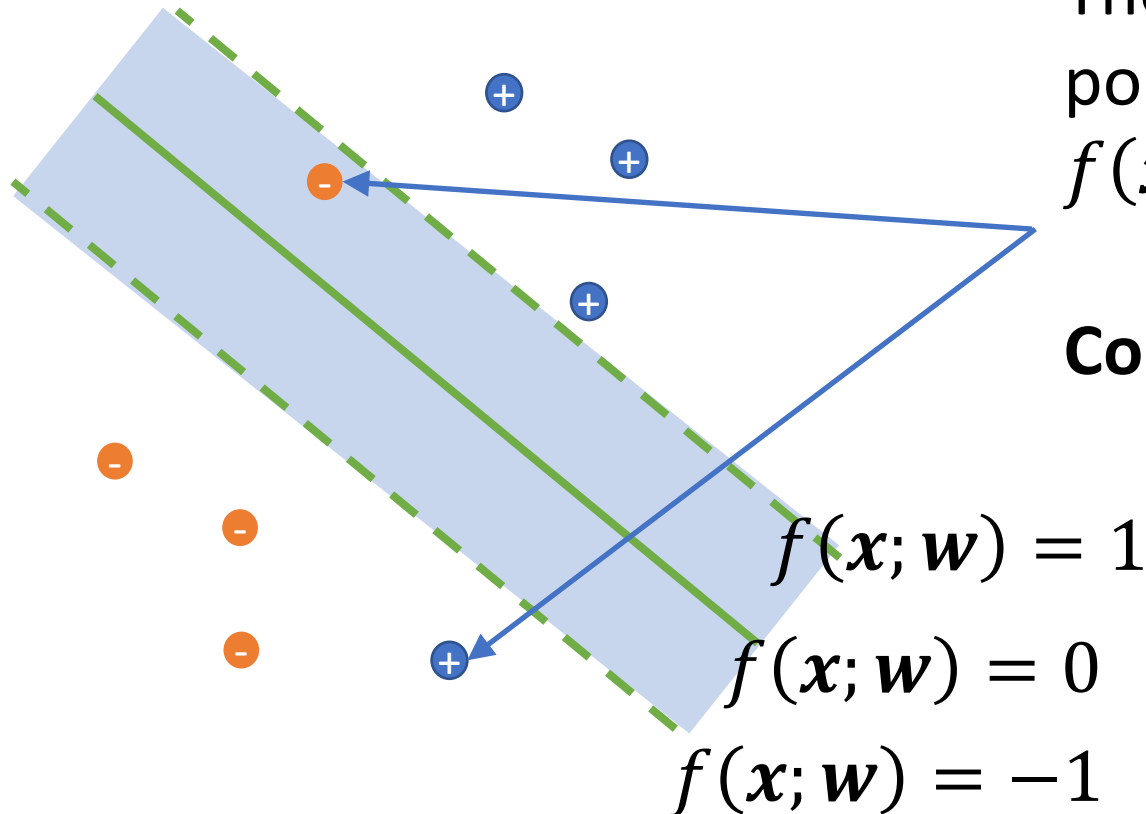
- In many cases, the dataset is not separable.



An error margin 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 $f(\mathbf{x}; \mathbf{w})y \leq 1$!

Constraint not satisfied!

Soft-Margin Classifier

- Minimize $||\mathbf{w}'||^2 + \sum_i \epsilon_i$
- Subject to $\forall_i, y_i f(\mathbf{x}_i; \mathbf{w}) + \epsilon_i \geq 1, \epsilon_i \geq 0$
- For each \mathbf{x}_i , we hope $y_i f(\mathbf{x}_i; \mathbf{w})$ can be at right side of the margin after some small positive “compensation” ϵ_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 ϵ is sparse. Why?

Soft-Margin Classifier

- Formally, the soft-margin classifier
- $\min_{\mathbf{w}, \epsilon} ||\mathbf{w}'||^2 + \sum_i \epsilon_i$
- Subject to $\forall_i, y_i(\langle \mathbf{w}', \mathbf{x}_i \rangle + w_0) + \epsilon_i \geq 1, \epsilon_i \geq 0$
- It turns out, if $f(\mathbf{x}; \mathbf{w})$ is a linear model,
- **Soft-Margin Classifier** is a **convex** minimization problem.
- **Every local minimum is a global minimum.**



The Lagrangian Dual

- Solving constrained problem can be rather complicated.
- **Lagrangian Dual**: a technique transforms constrained problem into a less constrained problem.
- For a constrained problem,
- $\min_{\theta} f(\theta)$ subject to $g_i(\theta) \leq 0, \forall i$
- We can construct a **Lagrangian** $l(\lambda)$:
- $l(\lambda) := \min_{\theta} f(\theta) + \sum_i \lambda_i g_i(\theta),$
- $\lambda_i \geq 0$ are called **Lagrangian multipliers**.
- **PRML Appendix E**.

The Lagrangian Dual

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

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

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

- $l(\lambda) :=$

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

- **2.** Derive optimality condition w.r.t. \mathbf{w} and ϵ :

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

Verify this!

The Lagrangian Dual

- Using optimality conditions:

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

- 3. Rewrite $l(\lambda) = -\frac{\tilde{\lambda}^\top X^\top X \tilde{\lambda}}{4} + \langle \lambda, \mathbf{1} \rangle$, ^{Verify it!}
 $X = [\mathbf{x}_1 \dots \mathbf{x}_n] \in R^{d \times n}, \tilde{\lambda} := [\lambda_1 \cdot y_1 \dots \lambda_n \cdot y_n]$

- 4. Maximize $l(\lambda)$ w.r.t. λ under constraints:

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

Needed to make sure the optimality of the
original problem

Soft-margin Classifier (Dual)

- $\max_{\lambda} - \frac{\tilde{\lambda}^T X^T X \tilde{\lambda}}{4} + \langle \lambda, \mathbf{1} \rangle$
- Subject to
- $0 \leq \lambda_i \leq 1, \sum_{i=1} \lambda_i y_i = 0$

- Recover $\hat{\mathbf{w}}' := \frac{\sum_{i=1} \hat{\lambda}_i y_i x_i}{2}$ using optimality condition.
- Put $\hat{\mathbf{w}}'$ back in the original problem and solve for $\hat{\mathbf{w}}_0$.

- We obtain $\hat{\mathbf{w}}$ using Lagrangian multipliers λ .

Soft-margin Classifier (Dual)

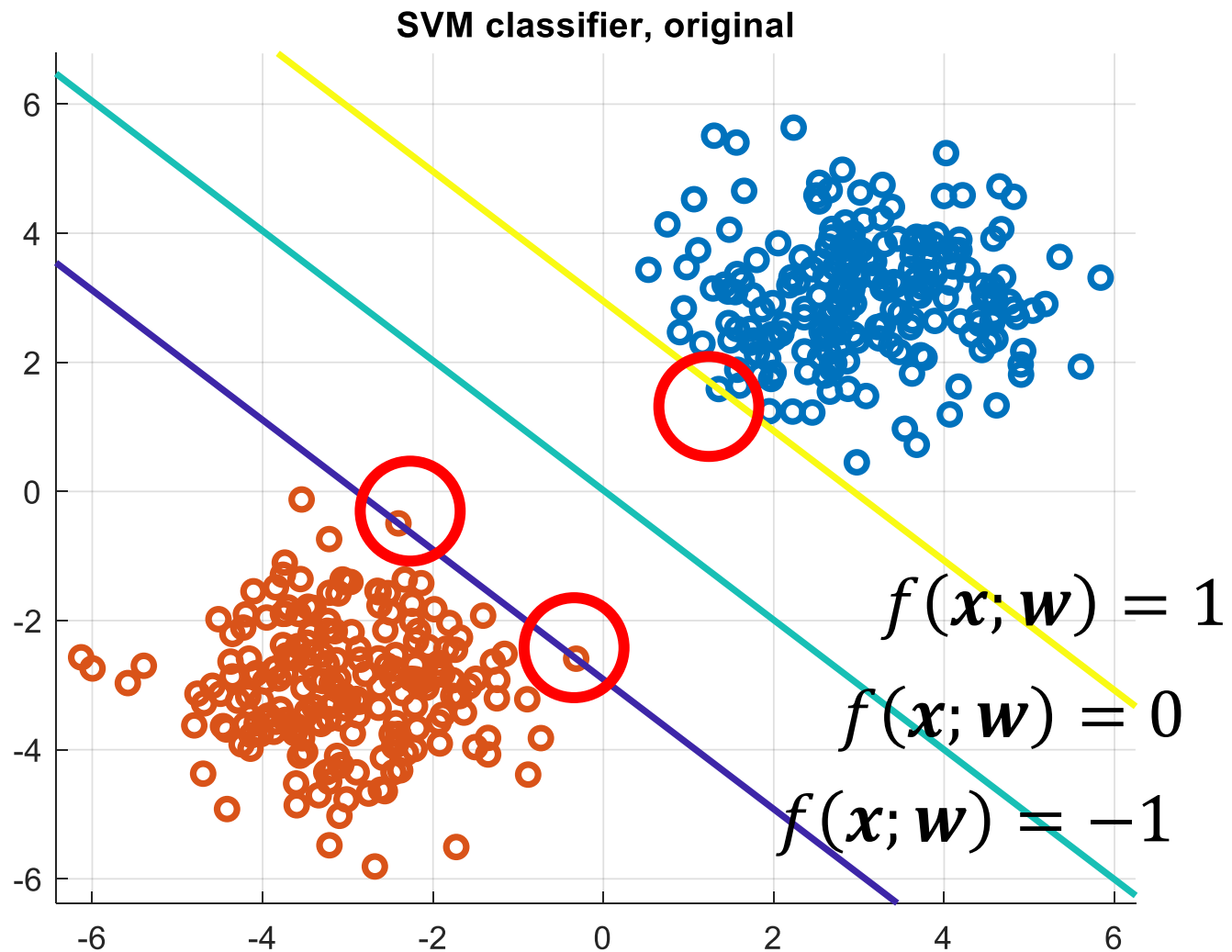
- $\max_{\lambda} - \frac{\tilde{\lambda}^\top X^\top X \tilde{\lambda}}{4} + \langle \lambda, \mathbf{1} \rangle$
- Subject to $0 \leq \lambda_i \leq 1, \sum_{i=1} \lambda_i y_i = 0$
- Our input data $\{\mathbf{x}_i\}$ **only appear at $X^\top X$**
- **Let $K = X^\top X$, then $K^{(i,j)} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$**
- Instead of using the inner product, we can use kernel functions $k(\mathbf{x}_i, \mathbf{x}_j)$ to perform training and prediction.
- **Homework:** write down decision function $f(\mathbf{x}; \mathbf{w})$ using kernel function k , w_0 and dual variable λ .

Original vs. Dual Problem

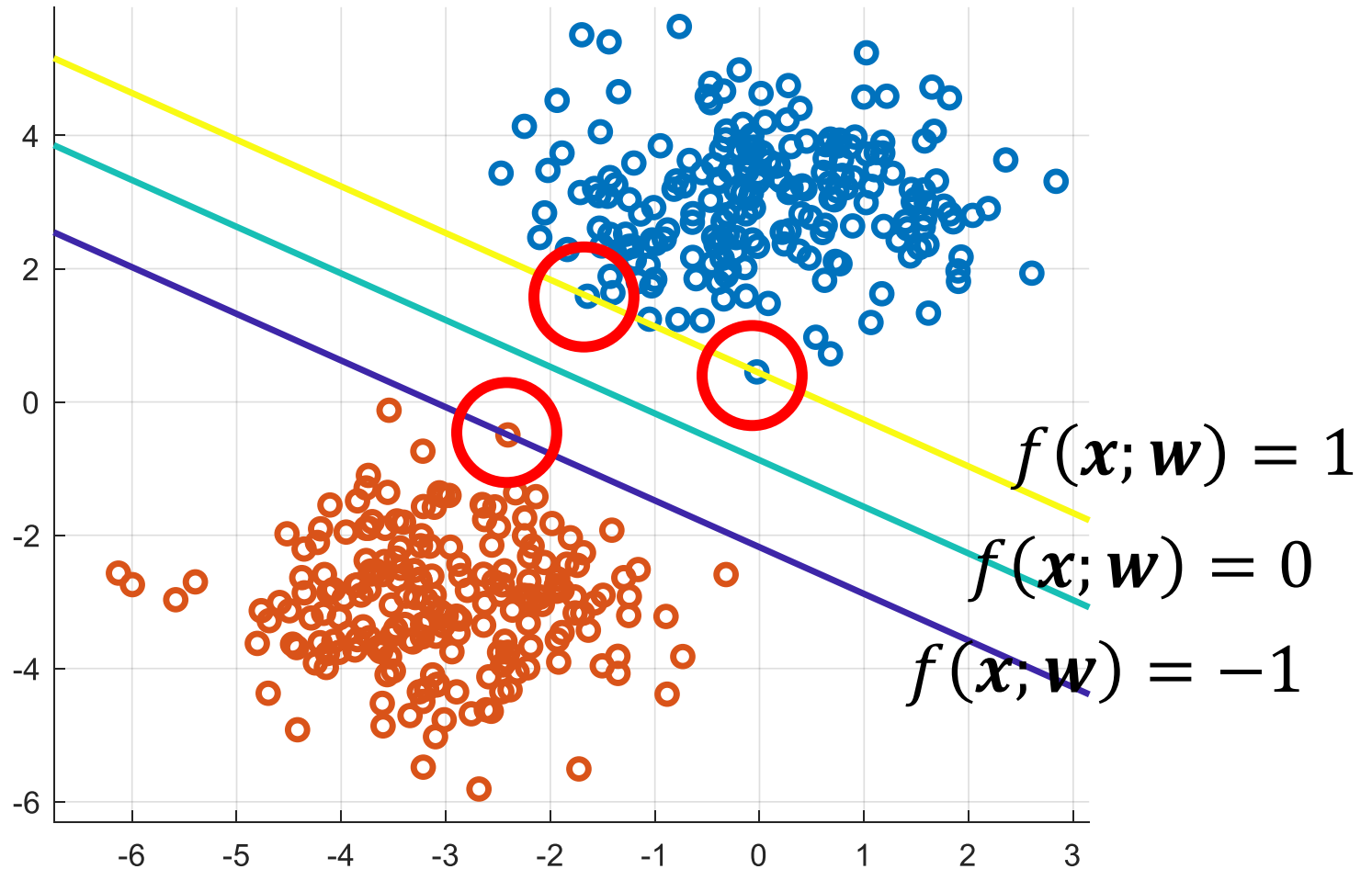
- $\min_{\mathbf{w}, \epsilon} ||\mathbf{w}'||^2 + \sum_i \epsilon_i$
- Subject to $\forall_i,$
 $y_i(\langle \mathbf{w}', \mathbf{x}_i \rangle + w_0) + \epsilon_i \geq 1,$
 $\epsilon_i > 0$
- Complex Constraints
- Quadratic w.r.t. $\mathbf{w} \in R^{d+1}$
- Slow when d is large

- $\max_{\lambda} - \frac{\tilde{\lambda}^T X^T X \tilde{\lambda}}{4} + \langle \lambda, \mathbf{1} \rangle$
- Subject to
- $0 \leq \lambda_i \leq 1$
- $\sum_{i=1} \lambda_i y_i = 0$
- Simpler Constraints
- Quadratic w.r.t. $\lambda \in R^n$
- Slow when n is large
- Can use kernel!

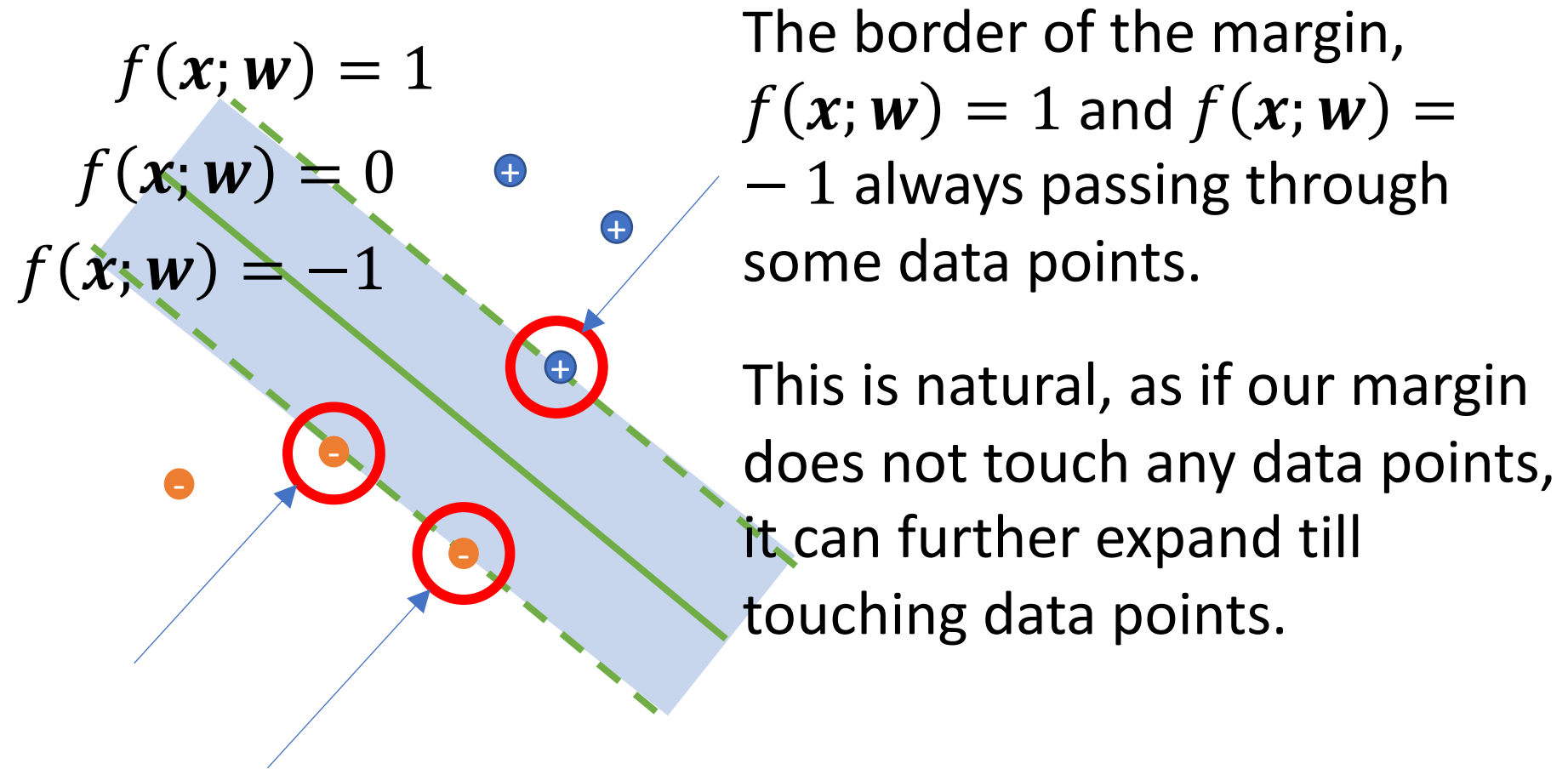
Toy Example



Toy Example

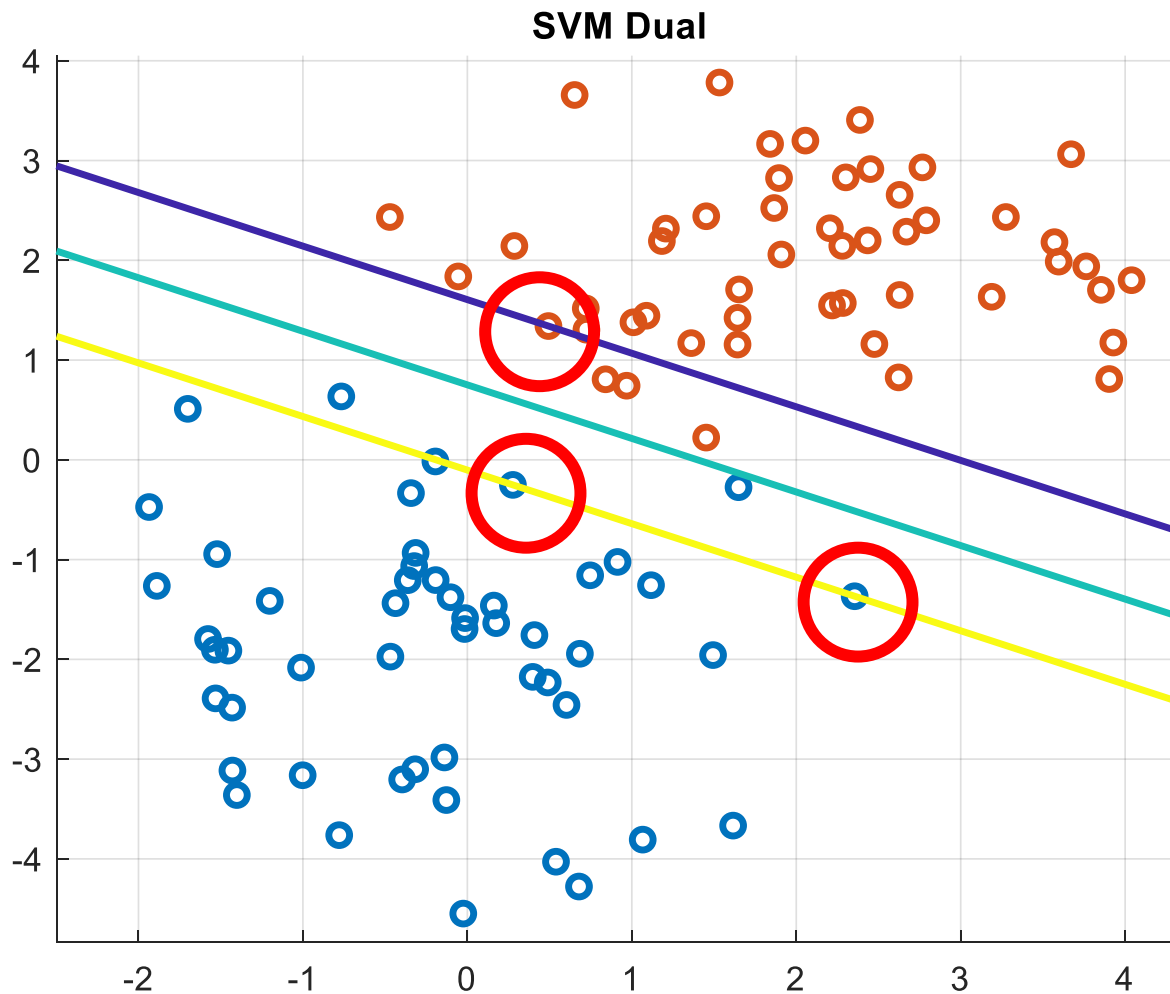


“Support Vectors”



These points as if they were resisting the expansion of the margin, are appropriately called “support vectors”.

Toy, Soft-margin, Dual



w by solving Original:

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

w recovered from λ :

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

Limitations of SVM

- SVM is **not** a probabilistic classifier
 - cannot be integrated with probabilistic classification models (generative or discriminative)
 - The decision function lacks probabilistic interpretation.
- Computational cost of SVM is high
 - Both 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 by “Maximum Margin” principle.
- Soft-margin SVM can classify overlapping pos/neg 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 ϵ .
- Represent prediction function $f(\boldsymbol{x}; \boldsymbol{w})$ using dual parameter $\boldsymbol{\lambda}$, kernel function k and bias w_0 .

Computing Lab
