



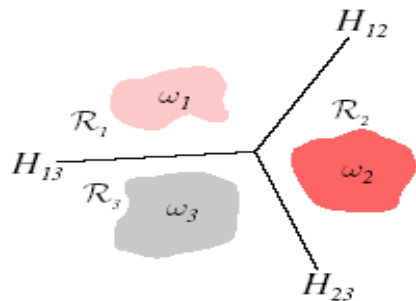
# Elements of Linear Discriminant Functions

Battista Biggio

Department of Electrical and Electronic Engineering  
University of Cagliari, Italy

# Introduction

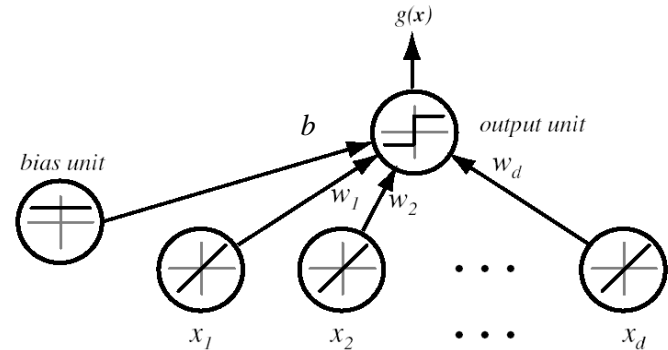
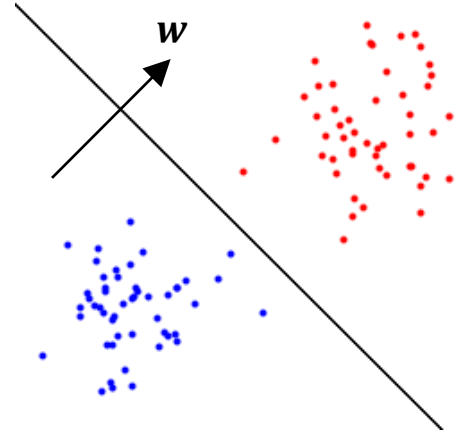
- We assume here that the form of the discriminant functions  $f_k(\mathbf{x}; \boldsymbol{\theta})$ ,  $k = 1, \dots, K$  is given, and that we can use the training data to estimate their parameters  $\boldsymbol{\theta}$ 
  - In Part 4, instead, we assumed that the underlying probability densities were known
- These methods are known as *nonparametric*
  - No assumption on the form of the underlying data probability distributions is made
- We focus here on functions that are *linear* in  $\mathbf{x}$ , i.e.,  $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}^T \mathbf{x} + b$ , with  $\boldsymbol{\theta} = (\mathbf{w}, b)$



**Example of linear discriminant functions  
on a 3-class classification problem**

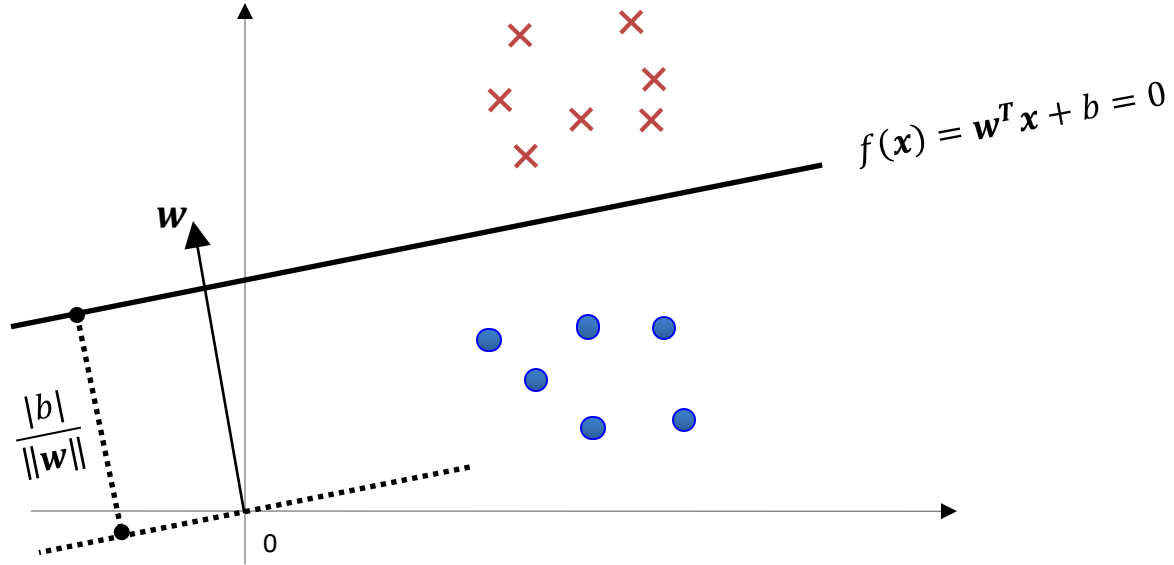
# Linear Discriminant Functions

- **Linear function:**  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{j=1}^d w_j x_j + b$ 
  - $\mathbf{w}$  is the weight vector, and  $b$  the bias
- Two-class classification
  - Positive ( $y = +1$ ) vs negative ( $y = -1$ ) class
  - Decision rule:  $y = \begin{cases} +1 & \text{if } f(\mathbf{x}) \geq 0 \\ -1, & \text{otherwise} \end{cases}$
- **Graphical representation**
  - Each input feature value  $x_j$  is multiplied by the corresponding weight value  $w_j$
  - Bias is multiplied by 1
  - The output unit sums all its inputs, computing  $f(\mathbf{x})$  and thresholds it to estimate  $y$  (+1 or -1)



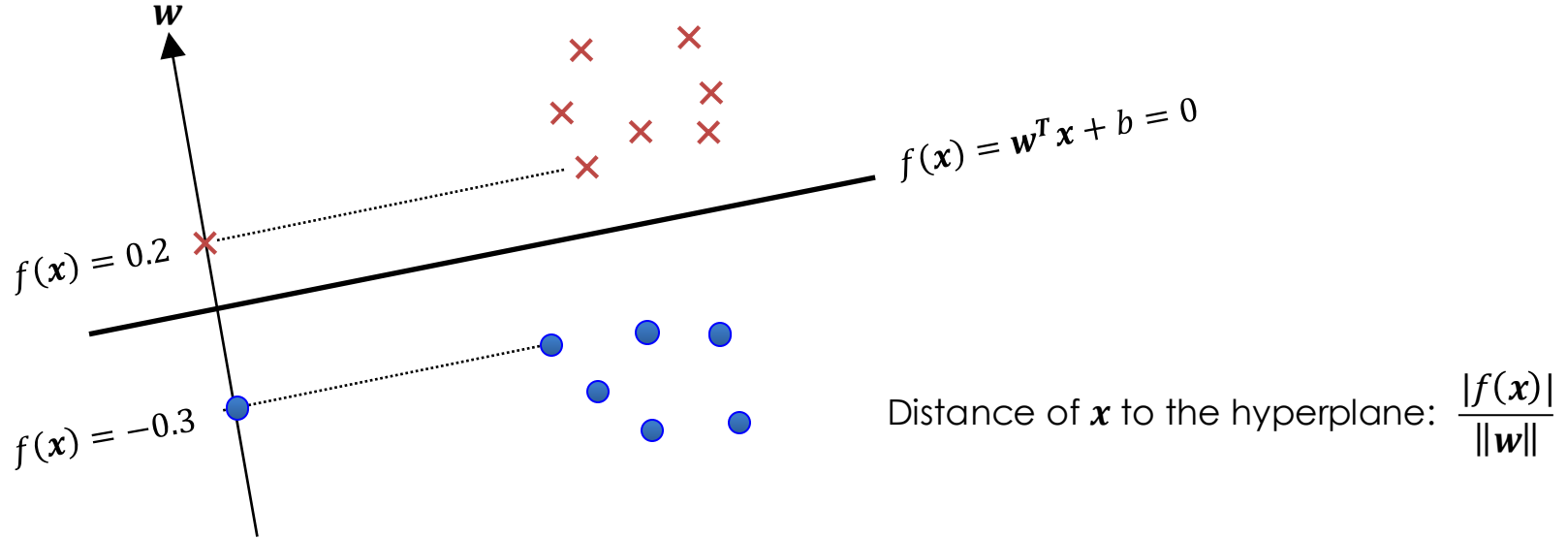
# Linear Discriminant Functions

$$D = \{\mathbf{x}_i, y_i\}_{i=1}^n$$
$$\mathbf{x} \in \mathbb{R}^d$$
$$y \in \{-1, +1\}$$

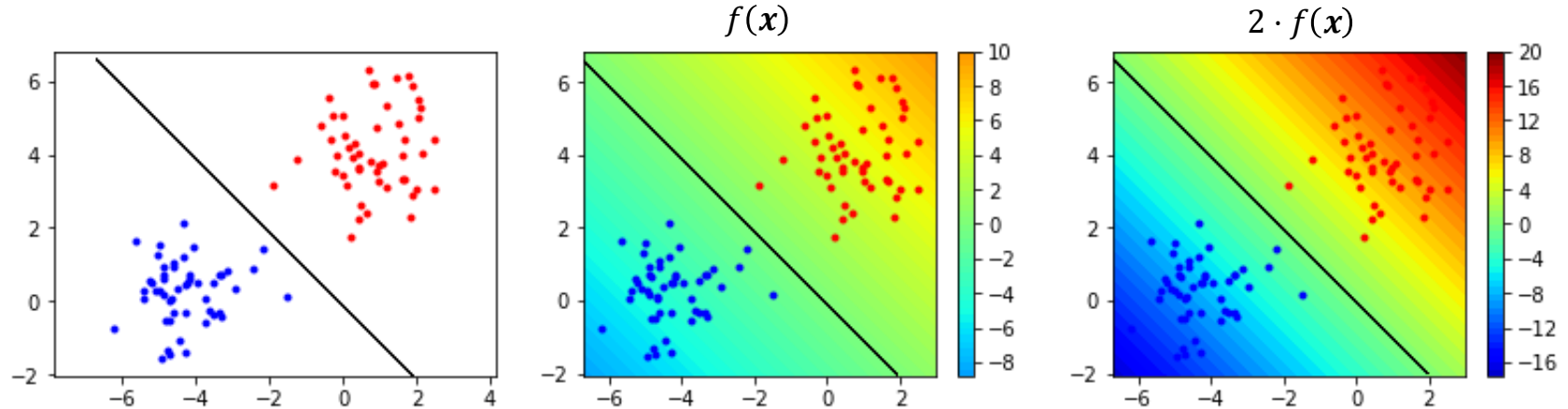


# Linear Discriminant Functions

- The function  $f(x) = \mathbf{w}^T \mathbf{x} + b$  projects  $\mathbf{x}$  onto the hyperplane normal
  - Its value is proportional to the distance of  $\mathbf{x}$  to the hyperplane



# Linear Discriminant Functions



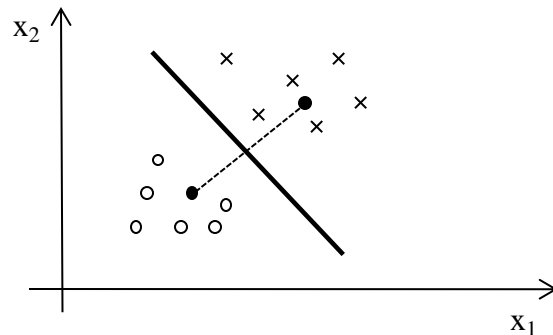
- By linearity, it is easy to see that multiplying  $f(x)$  by a constant factor amounts to multiplying its parameters by the same factor
  - For example:  $2 \cdot f(x) = 2 \mathbf{w}^T \mathbf{x} + 2b = \hat{\mathbf{w}}^T \mathbf{x} + \hat{b}$ , with  $\hat{\mathbf{w}} = 2\mathbf{w}$  and  $\hat{b} = 2b$
- While the boundary at  $f(x) = 0$  does not change, the slope of the function changes!

# A Simple Example: The Nearest Mean Classifier (NMC)

- This classifier estimates the mean values  $\mu_1$  and  $\mu_2$  of the two classes from the training set and assigns unknown samples  $\mathbf{x}^*$  to the class with the smallest Euclidean distance:

$$d(\mathbf{x}^*, \mu_2) \underset{\omega_2}{\overset{\omega_1}{\leq}} d(\mathbf{x}^*, \mu_1)$$

- It is easy to see that the decision boundary is the hyperplane perpendicular to the vector  $(\mu_1 - \mu_2)$  and passing through the mean point  $(\mu_1 + \mu_2)/2$
- Accordingly, the NMC is a linear classifier
  - Try to find its  $\mathbf{w}$  and  $b$  parameter values!



# Learning as an Optimization Problem



# Learning as an Optimization Problem

- How do we estimate the classifier parameters  $\mathbf{w}$  and  $b$ ?
- Modern approaches formulate the learning problem as an **optimization problem**
  - This is generally true also for nonlinear classification functions  $f(\mathbf{x}; \boldsymbol{\theta})$ , including modern deep-learning approaches and neural networks

$$\mathbf{w}^*, b^* = \operatorname{argmin}_{\mathbf{w}, b} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i))}_{\text{loss term } L(D, \boldsymbol{\theta})} + \underbrace{\lambda \Omega(\mathbf{w})}_{\text{regularization term } \Omega(\mathbf{w})}$$

$\lambda$ : regularization hyperparameter

# Learning as an Optimization Problem

- The loss function  $\ell(y_i, f(\mathbf{x}_i))$  measures how much a prediction is wrong
  - e.g., the zero-one loss is 0 if points are correctly predicted, and 1 if they are not
- The regularization term  $\Omega(\boldsymbol{\theta})$  imposes a penalty on the magnitude of the classifier parameters to avoid overfitting and promote smoother functions, i.e., functions that change more gradually as we move across the feature space
- The hyperparameter  $\lambda$  tunes the trade-off between the training loss and regularization
  - Larger values tend to promote more regularized functions but with a larger training error
  - Smaller values tend to reduce the training error but learn more complex functions

# Learning as an Optimization Problem

- We start by considering a simplified setting in which we aim to find the best parameters  $\theta = (\mathbf{w}, b)$  that minimize the loss function  $L(D, \theta)$ , being  $D = (\mathbf{x}_i, y_i)_{i=1}^n$  the training dataset:

$$\theta^* = \operatorname{argmin}_{\theta} L(D, \theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; \theta))$$

- The loss function quantifies the error that the classifier, parameterized by  $\theta$ , is making on its predictions on the training data  $D$ 
  - This is also known as the principle of **Empirical Risk Minimization (ERM)**
- How do we select the loss function  $L(D, \theta)$  and solve the above problem?

# Learning as an Optimization Problem

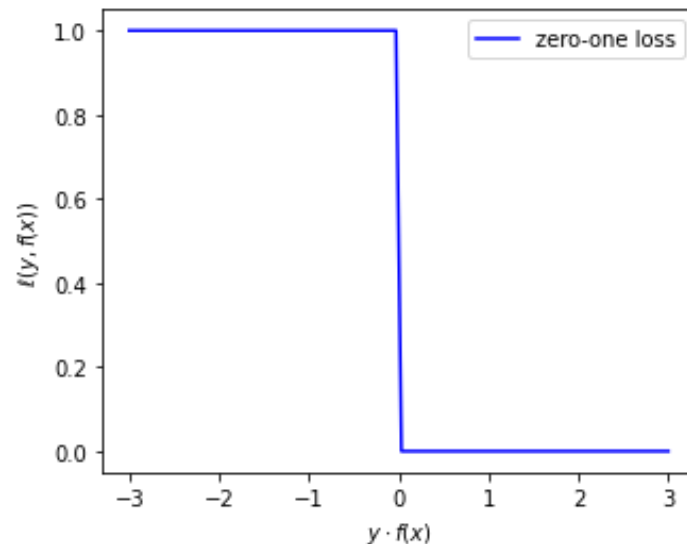
- In principle, we would like to minimize

$$L(D, \theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i; \theta))$$

- where  $\ell(y_i, f(x_i; \theta))$  is the zero-one loss
  - equal to 0 for correct predictions and 1 otherwise

$$\ell(y_i, f(x_i; \theta)) = \begin{cases} 1, & \text{if } y \cdot f(x) < 0 \\ 0, & \text{if } y \cdot f(x) \geq 0 \end{cases}$$

- However, solving this problem directly is NP-hard and computationally inefficient
  - <https://en.wikipedia.org/wiki/NP-hardness>

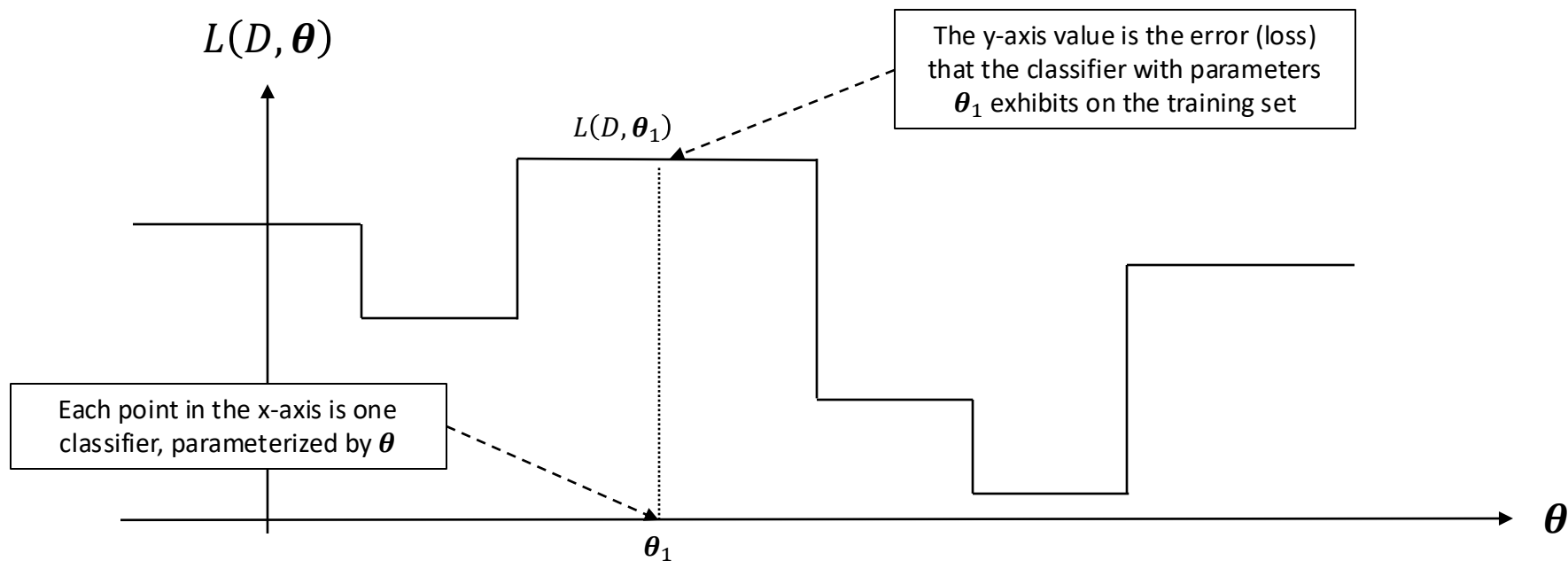


Note that  $yf(x)$  is positive for correct predictions and negative for wrong ones.

For correct (wrong) predictions,  $y$  and  $f(x)$  agree (disagree) in sign.

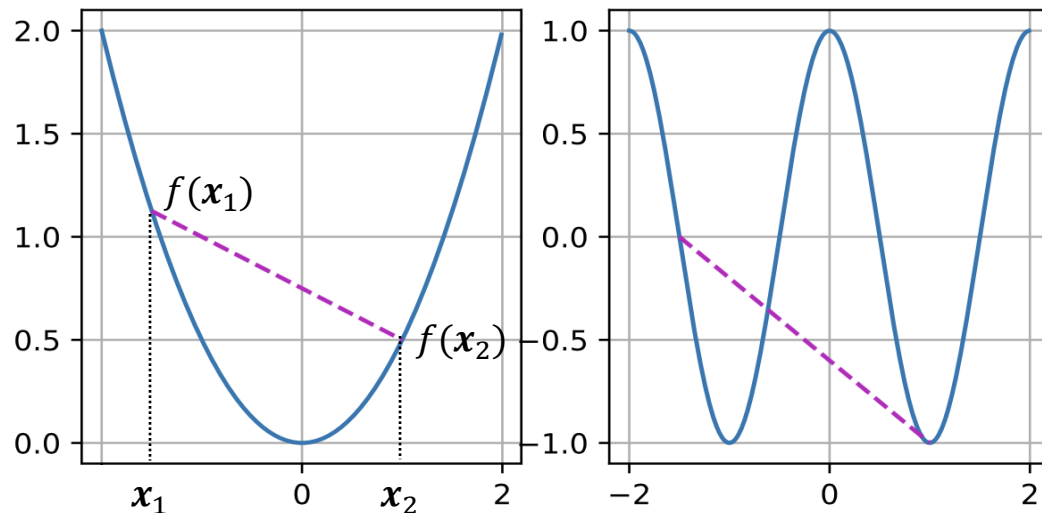
# The Zero-One Loss Landscape

- Non convex, hard to optimize (flat regions, bad local minima)



# Why Is Convexity Important for Optimization?

- **Convexity:**  $f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2), \quad \forall x_1, x_2, \quad \lambda \in [0,1]$



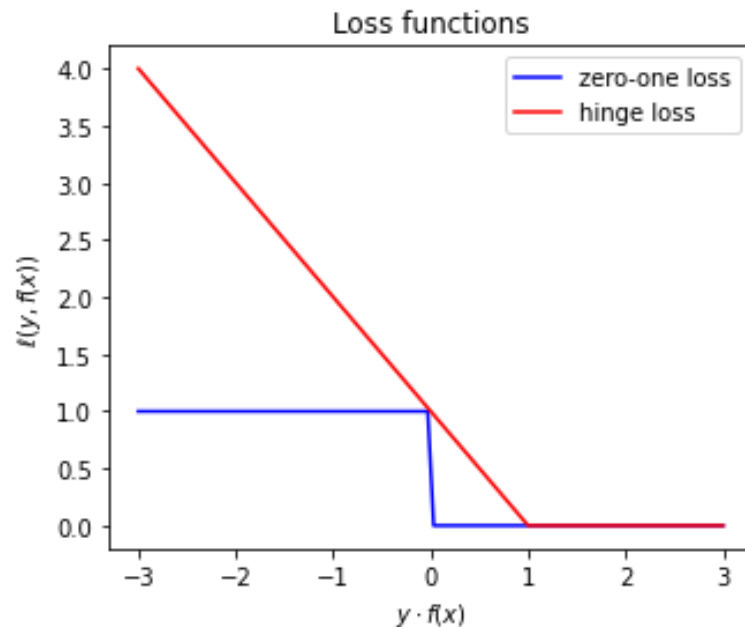
- **Desirable properties for optimization:** no local minima, convergence guarantees, etc.

# Loss Functions

- Now, recall that we aim to minimize:

$$L(D, \theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i; \theta))$$

- being  $\ell(y_i, f(x_i; \theta))$  the zero-one loss
- However, we know that minimizing this non-convex function is particularly difficult (NP hard)
- For this reason, convex (surrogate) loss functions are typically preferred
  - The tighter convex upper bound on the zero-one loss is called the **hinge loss**

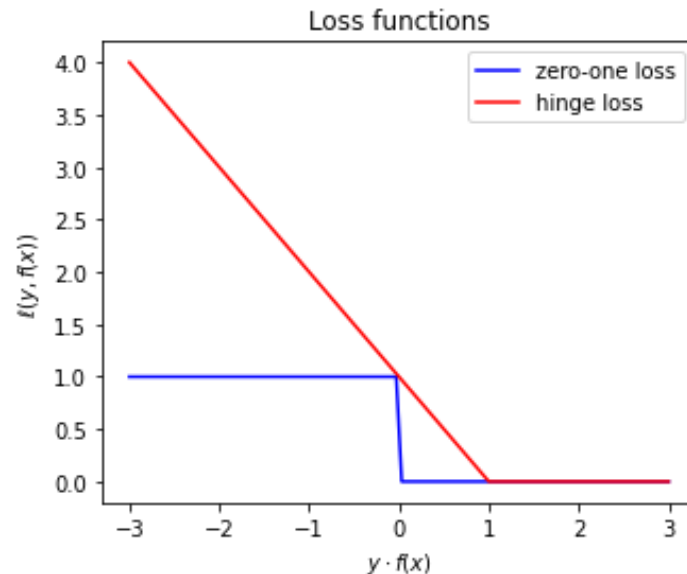


# Hinge Loss

- The hinge loss is computed as:

$$\ell(y, f(\mathbf{x}; \theta)) = \max(0, 1 - yf)$$

- It is the closest convex upper bound on the 0-1 loss
- Why are we interested in an upper bound?
  - since minimizing it, we also minimize the 0-1 loss
- Convexity helps find solutions efficiently while also providing guarantees on the optimality of the solution (global optima), algorithmic convergence, etc.
  - This is why convex surrogate functions are typically preferred in optimization

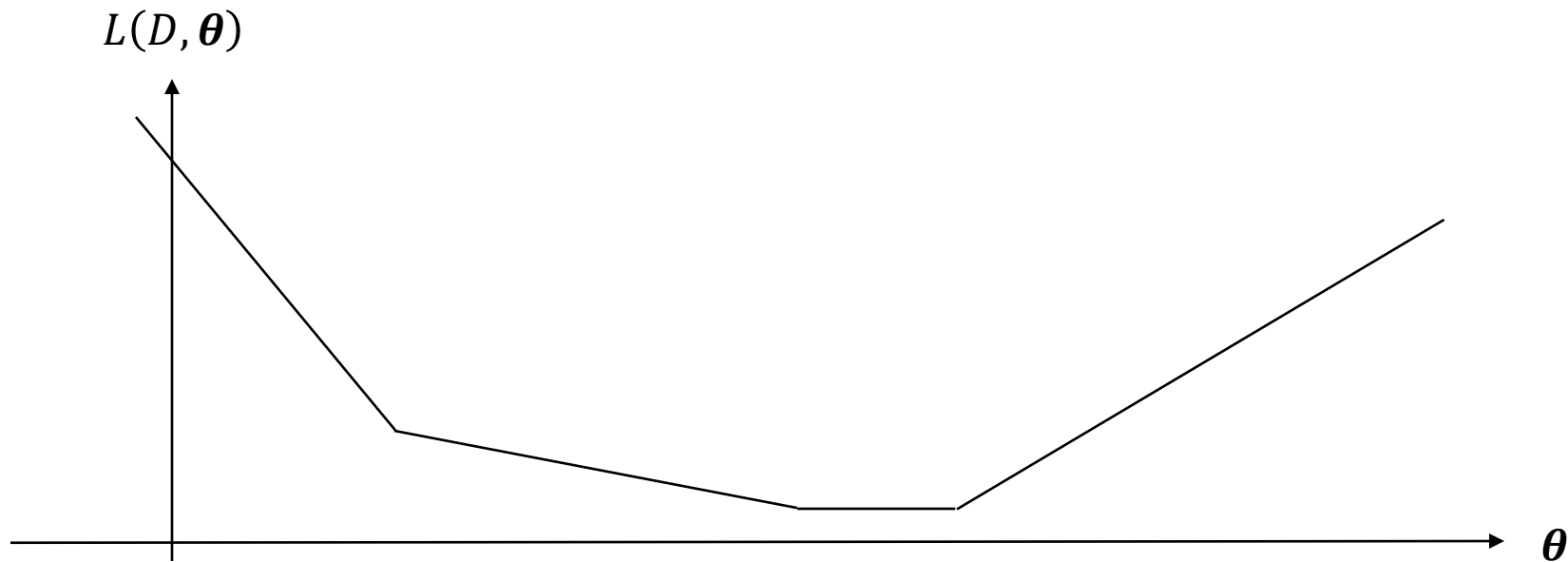


Note that  $yf(\mathbf{x})$  is positive for correct predictions and negative for wrong ones



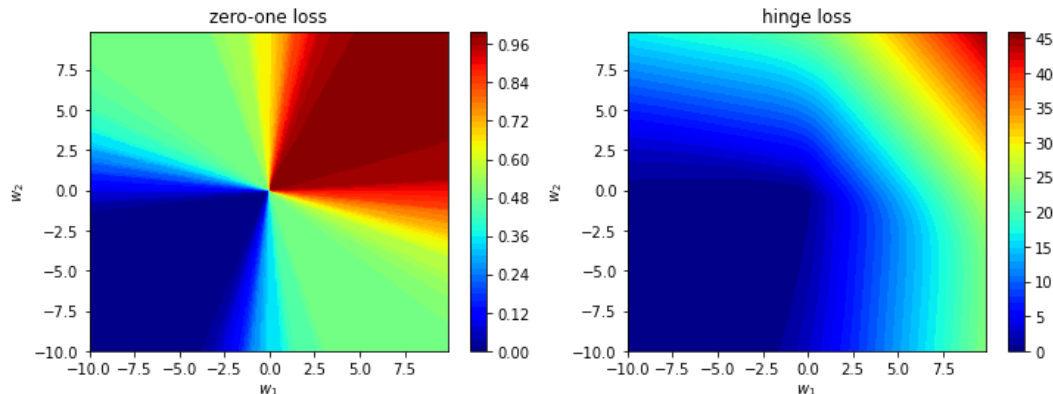
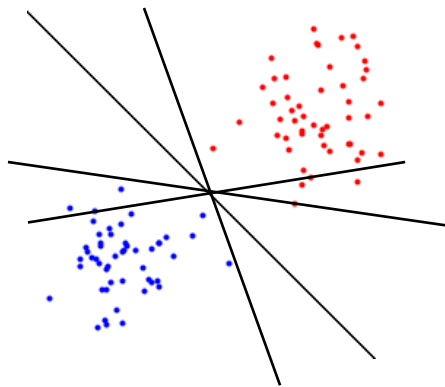
# The Hinge Loss Landscape

- Piecewise linear and convex, easier to optimize



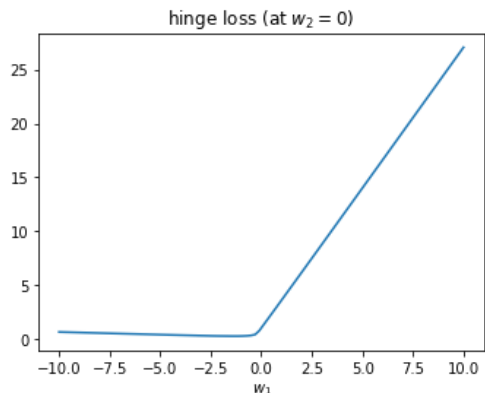
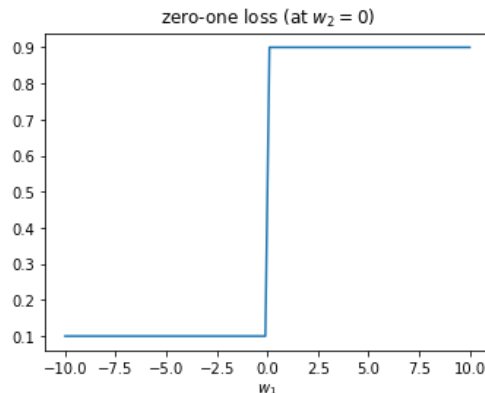
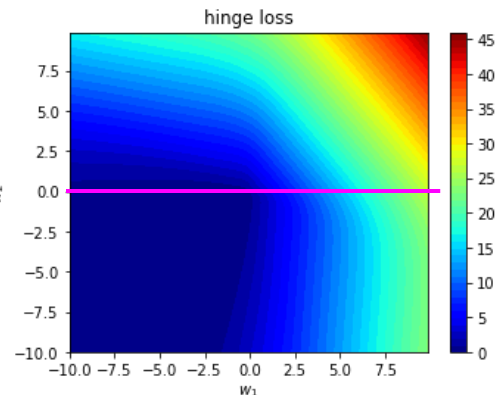
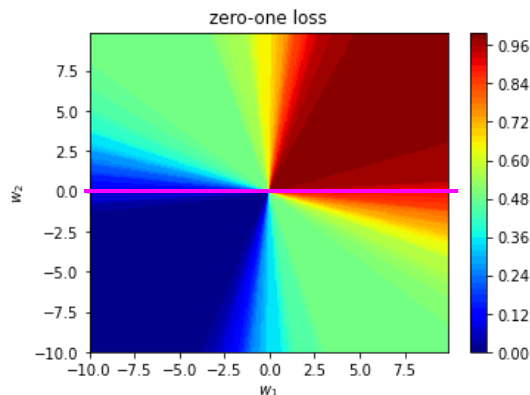
# A Closer Look at the Loss Minimization Problem

- Let's assume we fix  $b = 0$  and aim to minimize the training loss only w.r.t.  $w_1, w_2$
- Each pair  $w_1, w_2$  thus represents a different linear classifier (passing through the origin)
- For each of these classifiers, we report the corresponding training loss in a colored plot
  - this will show us the optimization *landscape*, i.e., the surface of the function we aim to minimize



# A Closer Look at the Loss Minimization Problem

- If we also fix  $w_2 = 0$ , and consider only optimizing  $w_1$ , we can also look at the profile of the loss along the line  $w_2 = 0$
- Note how the zero-one loss again reports a sharp profile (behaving like a step function)
- The hinge loss instead decreases more gracefully towards one direction
- We can thus say that the hinge loss provides a **smooth** approximation to the 0-1 loss



# Loss Minimization with Gradient Descent

# Gradient-based Optimization

- Optimizing *smooth* functions is much easier and efficient, as we can exploit **gradients**
  - This is not possible for the 0-1 loss (it is flat almost everywhere with gradients equal to zero)
- The key idea of gradient-based optimization is to start from a random point in the parameter space (*random initialization*) and then iteratively update the parameters along the gradient direction
- The gradient is the derivative of the loss function w.r.t. the classifier parameters:

$$L(D, \theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; \theta)), \quad \nabla_{\theta} L = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(y_i, f(\mathbf{x}_i; \theta))$$

- It is the direction in the parameter space along which the objective maximally increases
  - Following the negative gradient will thus minimize our training loss!

# Gradient Descent (a.k.a. Steepest Descent)

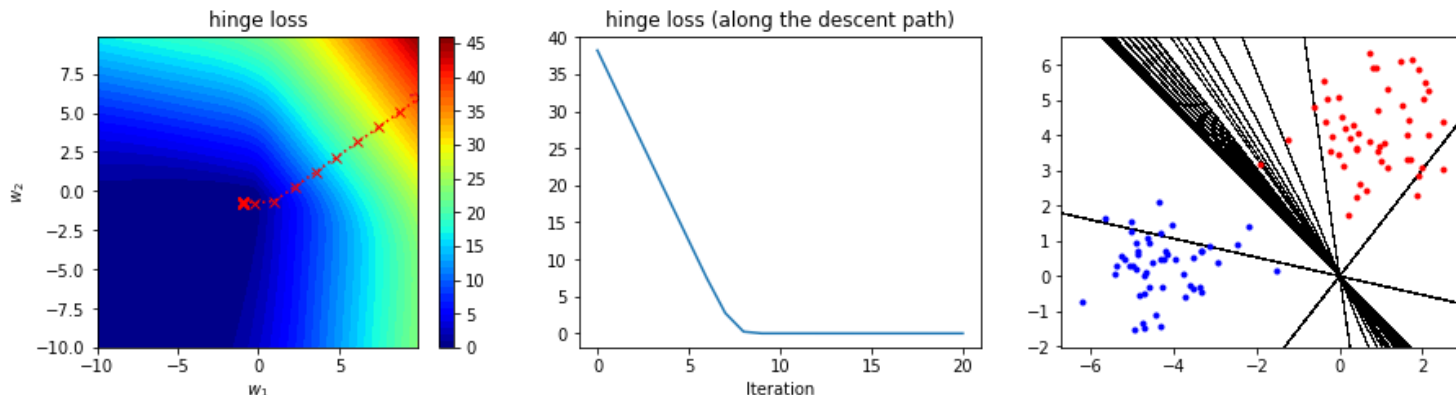
- The simplest gradient-based optimizer is the **steepest-descent** method

```
1. initialize  $\theta, \eta, K, \varepsilon$ 
2. for  $k$  in  $\{0, 1, \dots, K-1\}$ :
3.      $\theta_{k+1} = \theta_k - \eta \nabla L(\theta_k)$ 
4.     if  $|L(\theta_k) - L(\theta_{k+1})| < \varepsilon$ :
5.         break
```

- The parameters  $\theta$  are updated at each iteration, until
  - a maximum of  $K$  iterations are reached, or
  - the convergence/stop condition is met (lines 4-5 above)
- The stop condition checks that the last update has not significantly modified the objective function (i.e., the training loss is almost constant, as  $\varepsilon$  is a small number)
- The learning rate (or gradient step size)  $\eta$  affects convergence. If it is too small, convergence is too slow; if it is too large, the gradient algorithm may not even converge at all. Usually  $\eta$  is reduced across iterations to ensure convergence

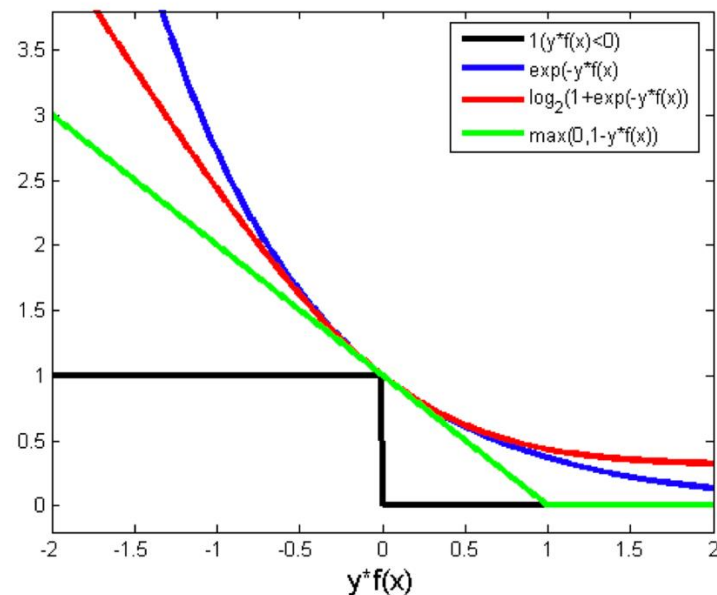
# Steepest Descent on the Hinge Loss

- Consider again our running example. We start optimization from  $w = (10, 6)$  and fix  $b = 0$ 
  - See the notebook for details on the computation of the gradient of the hinge loss w.r.t.  $w, b$
- Here's what the plots show:
  - how weights change along the path until the local minimum (blue region) is reached
  - how the hinge loss decreases along the path (convergence is reached when it becomes flat)
  - how the linear classifier changes across iterations, converging to a solution that splits the training points achieving zero loss



# Other Convex Losses

- Other convex upper bounds on the zero-one loss:
  - Hinge loss:  $\ell(y, f(\mathbf{x})) = \max(0, 1 - yf)$
  - Exponential loss:  $\ell(y, f(\mathbf{x})) = e^{-yf}$
  - Logistic loss:  $\ell(y, f(\mathbf{x})) = \log_2(1 + e^{-yf})$

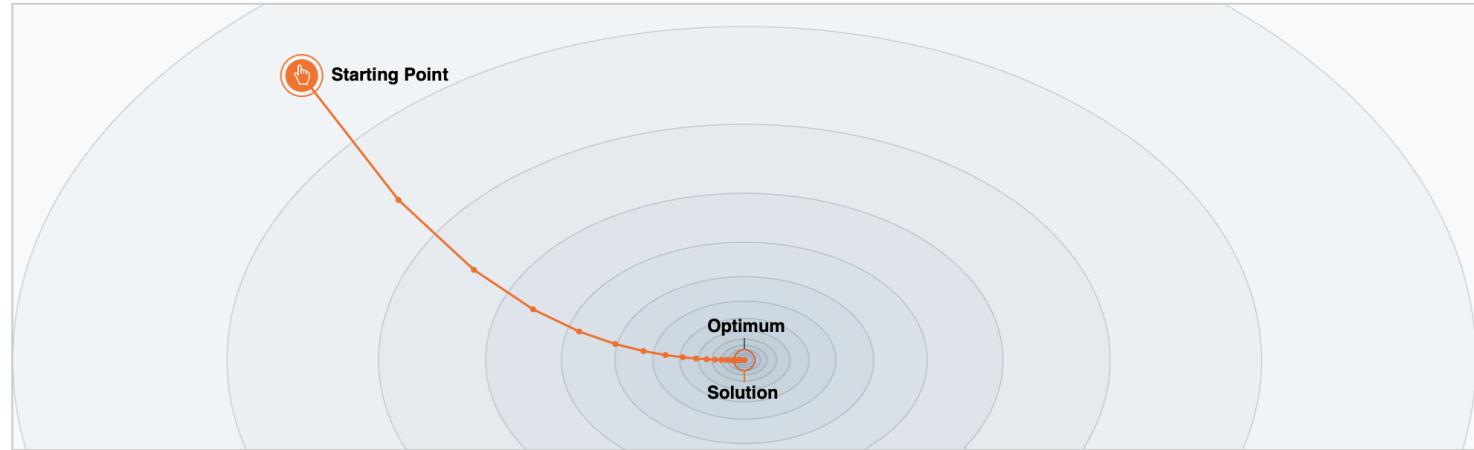




# Gradient Descent: Step Size and Convergence

# Example: Steepest Descent on Quadratic Objective

- Well-conditioned quadratic objective, small step size



Step-size  $\alpha = 0.22$

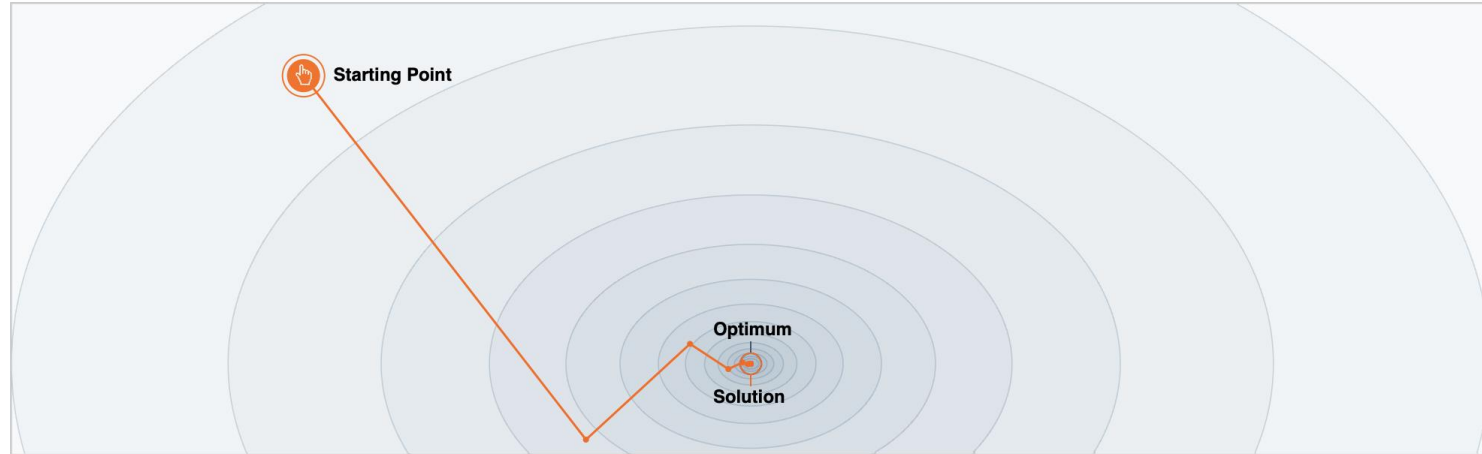


On a well-conditioned quadratic function, the gradient descent converges on few iterations to the optimum.

Examples from: [http://fa.bianp.net/teaching/2018/eecs227at/gradient\\_descent.html](http://fa.bianp.net/teaching/2018/eecs227at/gradient_descent.html)

# Example: Steepest Descent on Quadratic Objective

- Well-conditioned quadratic objective, large step size



Step-size  $\alpha = 0.63$

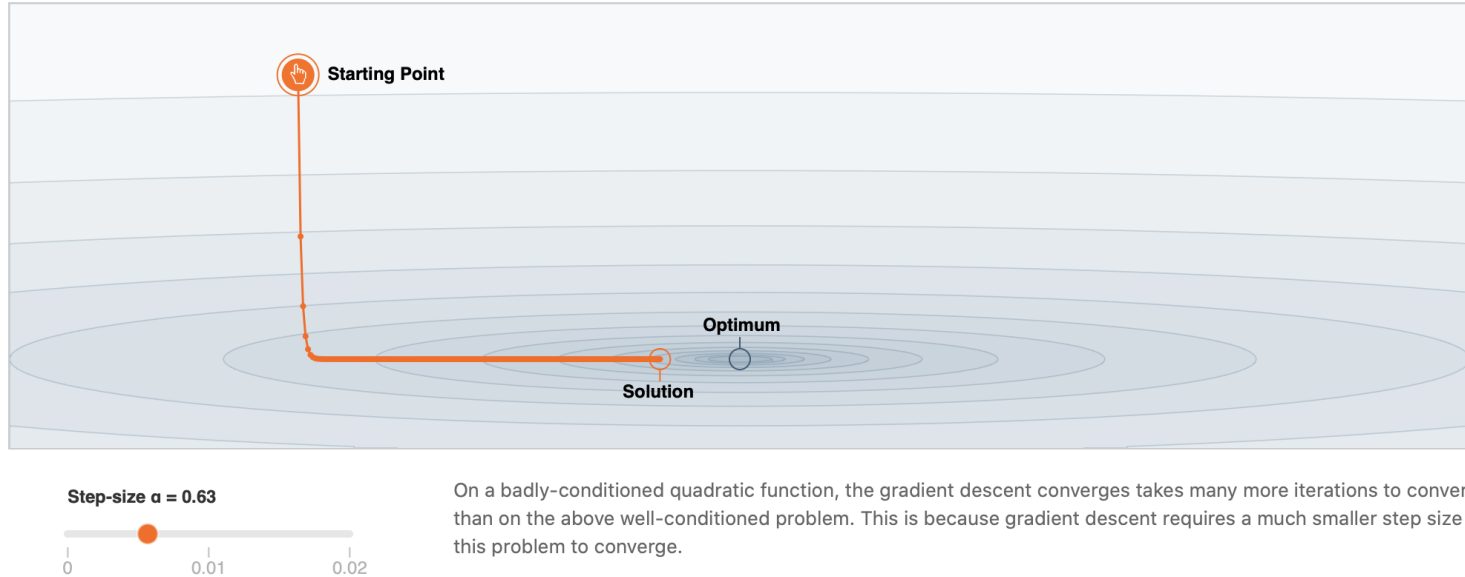


On a well-conditioned quadratic function, the gradient descent converges on few iterations to the optimum.

Examples from: [http://fa.bianp.net/teaching/2018/eecs227at/gradient\\_descent.html](http://fa.bianp.net/teaching/2018/eecs227at/gradient_descent.html)

# Example: Steepest Descent on Quadratic Objective

- Badly-conditioned quadratic objective, slow convergence

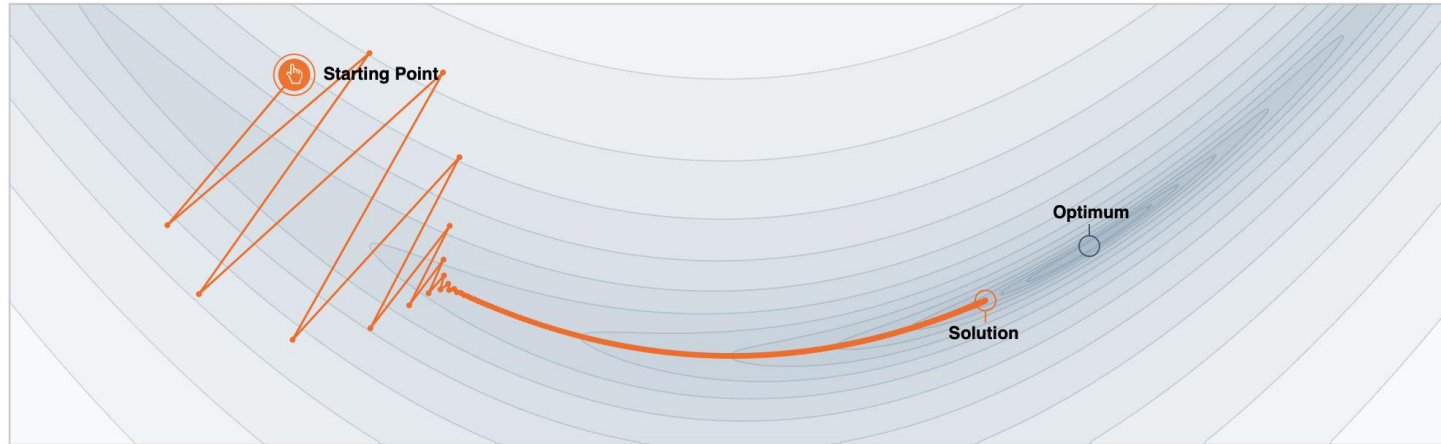


On a badly-conditioned quadratic function, the gradient descent converges takes many more iterations to converge than on the above well-conditioned problem. This is because gradient descent requires a much smaller step size on this problem to converge.

Examples from: [http://fa.bianp.net/teaching/2018/eecs227at/gradient\\_descent.html](http://fa.bianp.net/teaching/2018/eecs227at/gradient_descent.html)

# Example: Steepest Descent on Non-convex Objective

- Badly-conditioned non-convex objective, slow convergence and initial instability



Step-size  $\alpha = 0.63$



Gradient descent also converges on a badly-conditioned non-convex problem. Convergence is slow in this case.

Examples from: [http://fa.bianp.net/teaching/2018/eecs227at/gradient\\_descent.html](http://fa.bianp.net/teaching/2018/eecs227at/gradient_descent.html)

# Gradient Descent on Quadratic Objectives

# Gradient Descent on Quadratic Objectives

- The steepest descent algorithm linearizes the function  $L(D, \boldsymbol{\theta})$  around  $\boldsymbol{\theta}$ 
  - It's a first-order Taylor expansion!
- If the function is quadratic or even non-convex, convergence may be too slow
- Alternatively, we can use a quadratic expansion for  $L(D, \boldsymbol{\theta})$  around  $\boldsymbol{\theta}$ :

$$L(\boldsymbol{\theta}_{k+1}) \cong L(\boldsymbol{\theta}_k) + \nabla L(\boldsymbol{\theta}_k)(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k) + \frac{1}{2}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k)^T \mathbf{H}_k(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k)$$

- Note that  $\mathbf{H}_k = \nabla_{\boldsymbol{\theta}, \boldsymbol{\theta}}^2 L(\boldsymbol{\theta}_k)$  is the Hessian matrix (containing the second derivatives)

# Steepest Descent with Exact Step

- The quadratic expansion can be used to find the optimal step size  $\eta_k$  at each iteration
- If we plug-in the update rule:  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta_k \nabla L(\boldsymbol{\theta}_k)$  in the previous expression, we obtain:

$$L(\boldsymbol{\theta}_{k+1}) \cong L(\boldsymbol{\theta}_k) - \eta_k \|\nabla L(\boldsymbol{\theta}_k)\|^2 + \frac{1}{2} \eta_k^2 \nabla L(\boldsymbol{\theta}_k)^T \mathbf{H}_k \nabla L(\boldsymbol{\theta}_k)$$

- The minimization of the quadratic form of  $L(\boldsymbol{\theta}_{k+1})$  derived before can be solved in closed form, by setting the derivate with respect to  $\eta_k$  equal to zero, and solving for  $\eta_k$
- The solution provides the optimal step size to be set at each iteration:

$$\eta_k = \frac{\|\nabla L(\boldsymbol{\theta}_k)\|^2}{\nabla L(\boldsymbol{\theta}_k)^T \mathbf{H}_k \nabla L(\boldsymbol{\theta}_k)}$$



# Newton-Raphson Method

- The quadratic expansion below can be directly minimized via a closed-form solution

$$L(\boldsymbol{\theta}_{k+1}) \cong L(\boldsymbol{\theta}_k) + \nabla L(\boldsymbol{\theta}_k)(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k) + \frac{1}{2}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k)^T \mathbf{H}_k(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k)$$

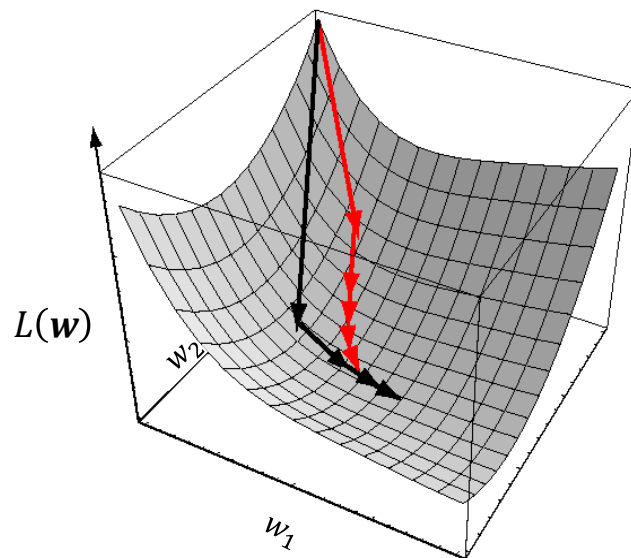
- By setting the derivative of  $L$  w.r.t.  $\boldsymbol{\theta}_k$  equal to 0, we obtain:  $\nabla L(\boldsymbol{\theta}_k) + \mathbf{H}_k(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k) = \mathbf{0}$
- The solution is the update rule for Newton-Raphson:  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \mathbf{H}_k^{-1} \nabla L(\boldsymbol{\theta}_k)$

## Newton-Raphson method

```
1. initialize  $\boldsymbol{\theta}$ ,  $K$ ,  $\varepsilon$ 
2. for  $k$  in  $\{0, 1, \dots, K-1\}$ :
3.      $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \mathbf{H}_k^{-1} \nabla L(\boldsymbol{\theta}_k)$ 
4.     if  $|L(\boldsymbol{\theta}_k) - L(\boldsymbol{\theta}_{k+1})| < \varepsilon$ :
5.         break
```

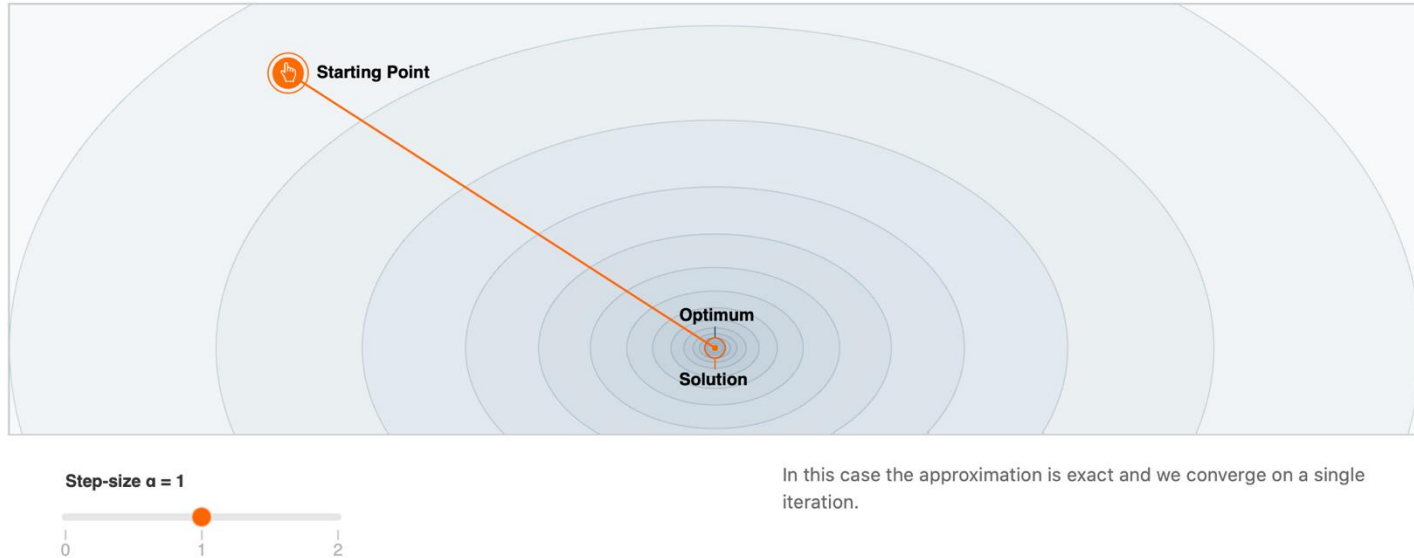
# Example: Newton-Raphson vs Steepest Descent

- **Steepest descent** in red
- **Newton-Raphson** in black
- **Newton-Raphson**
  - **Pros.** Faster convergence
  - **Cons.** Hessian inversion is computationally expensive. If Hessian is badly conditioned, convergence may be problematic



# Example: Newton-Raphson on Quadratic Objective

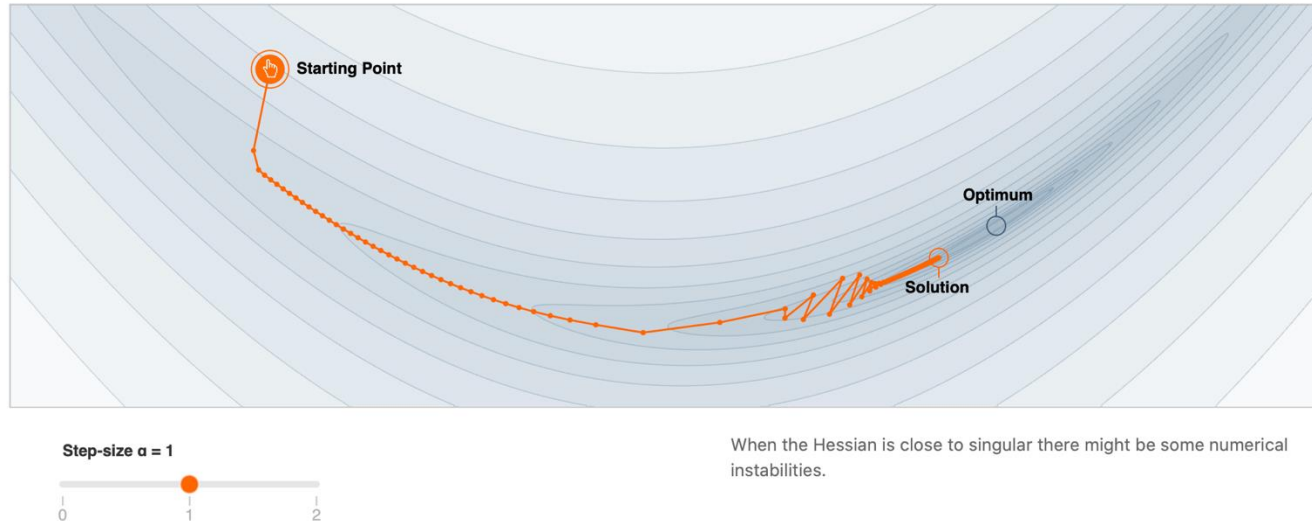
- Exact solution in one iteration



Examples from: <http://fa.bianp.net/teaching/2018/eecs227at/newton.html>

# Example: Newton-Raphson on Non-convex Objective

- Badly-conditioned non-convex objective, numerical instabilities in Hessian inversion



When the Hessian is close to singular there might be some numerical instabilities.

Examples from: <http://fa.bianp.net/teaching/2018/eecs227at/newton.html>

# Tuning the Step Size

- In general, these quadratic approximations are computationally expensive, as they involve Hessian computation/inversion
- Lightweight **line-search methods** are thus preferred to find an approximate good step size at each iteration. See, e.g., *backtracking* line search:
  - [http://fa.bianp.net/teaching/2018/eecs227at/gradient\\_descent.html](http://fa.bianp.net/teaching/2018/eecs227at/gradient_descent.html)
- Other strategies instead just decrease  $\eta_k$  at each iteration, with different decaying rates. See, e.g., *cosine annealing*:
  - <https://towardsdatascience.com/https-medium-com-reina-wang-tw-stochastic-gradient-descent-with-restarts-5f511975163>

# Summary

- Linear classification functions for binary (2-class) classification
- Loss minimization principle (*learning the classifier parameters*)
  - Convex optimization / loss functions
- Gradient-based optimizers
  - Steepest descent
  - Newton-Raphson method
- Tuning the step size (*line-search and decay strategies*)

# Support Vector Machines

# Support Vector Machines (SVMs)

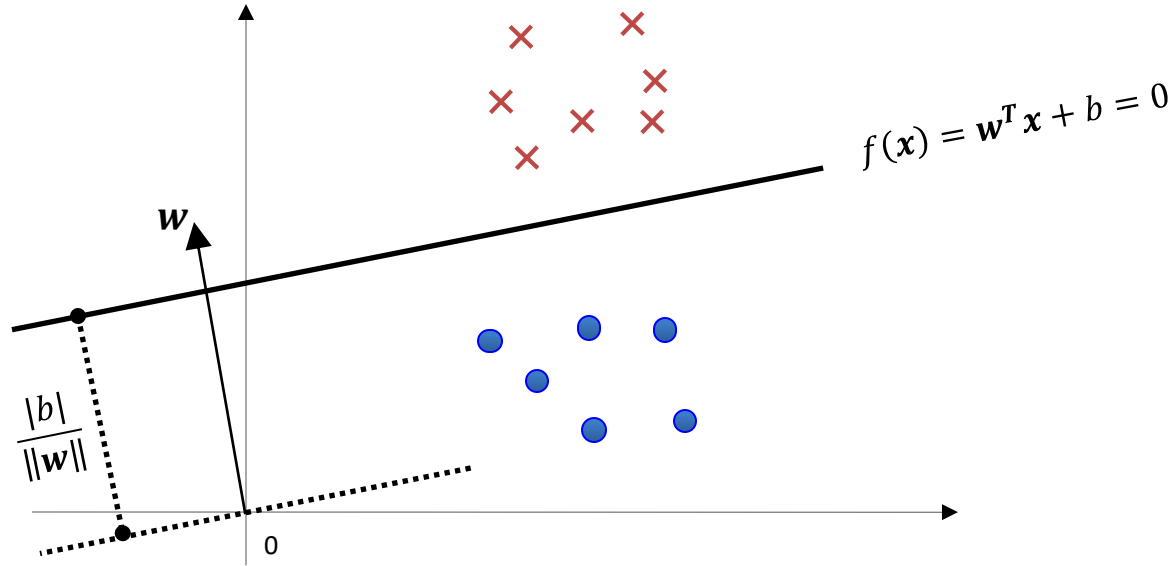
- State-of-the-art model with strong mathematical interpretation and significant practical relevance in many applications
- Invented by Russian mathematician Vladimir Vapnik in mid 1960s, and ideas much extended in 1990s
- Uses a specific type of loss function and learning algorithm
  - convex optimization / quadratic programming
- Can solve non-linearly separable problems





# A Short Recap on Linear Discriminant Functions

$$D = \{\mathbf{x}_i, y_i\}_{i=1}^n$$
$$\mathbf{x} \in \mathbb{R}^d$$
$$y \in \{-1, +1\}$$



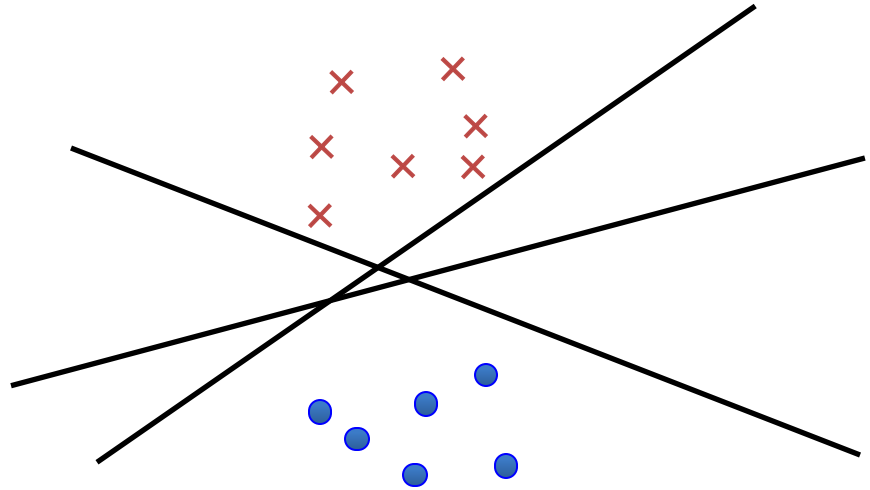
Distance of  $\mathbf{x}$  to the hyperplane:  $\frac{|f(\mathbf{x})|}{\|\mathbf{w}\|}$

# Support Vector Machines

The “*no maths*” version

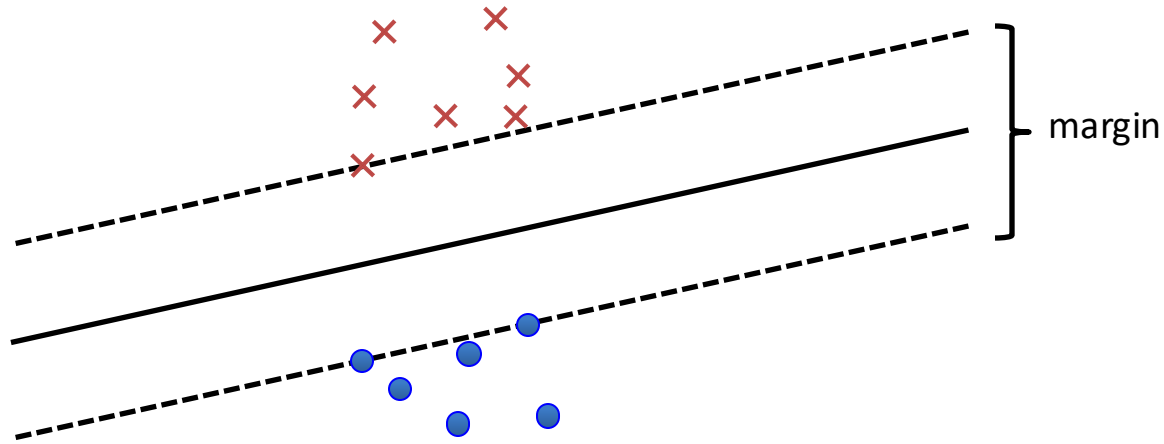
# Underlying Idea

- Several possible decision boundaries
  - All get 100% accuracy on this training data
- Which one would you pick? Why?



# Underlying Idea: Margin Maximization

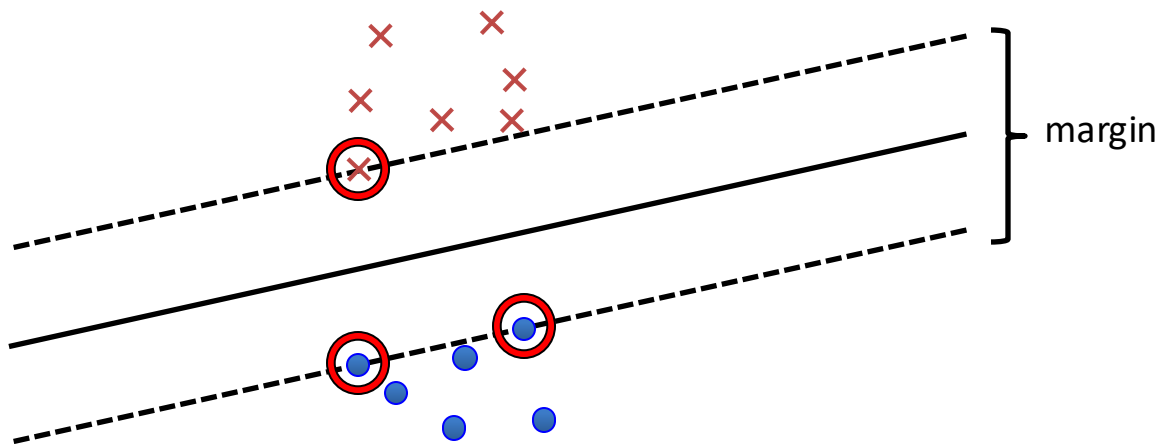
- The SVM finds this one – the furthest boundary from the two clusters



- Distance to the closest training points is called “margin”
  - equal on both sides of the boundary

# Margin Maximization and Support Vectors

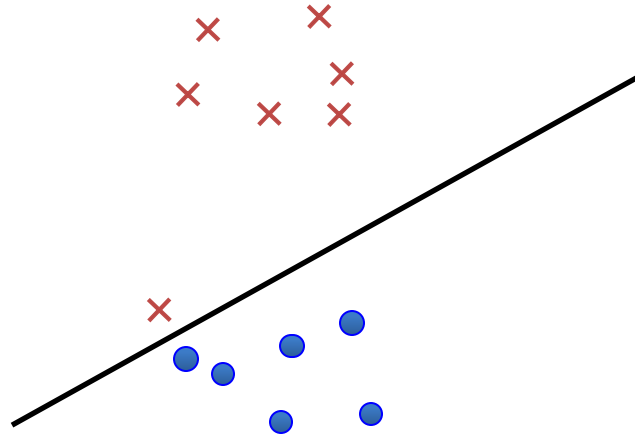
- The circled points are called **support vectors**



- The decision hyperplane only depends on the SVs
  - the other points can move freely without violating the margin

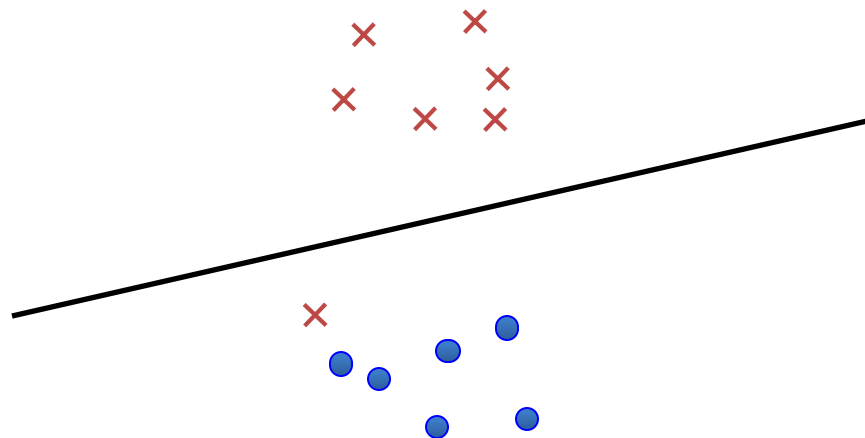
# What About Outliers?

- This is the “optimal” boundary – doesn't seem so clever though

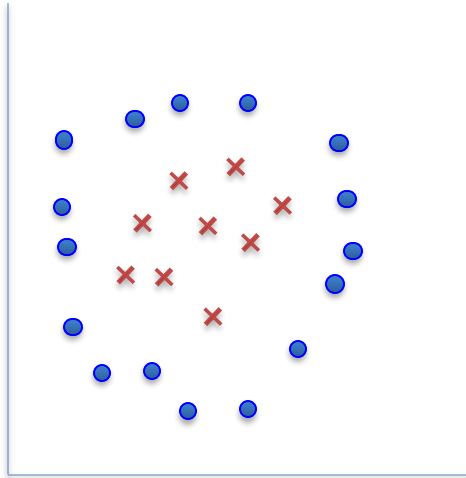


# What About Outliers?

- SVM can selectively ignore certain data points (e.g., outliers)
  - soft-margin maximization



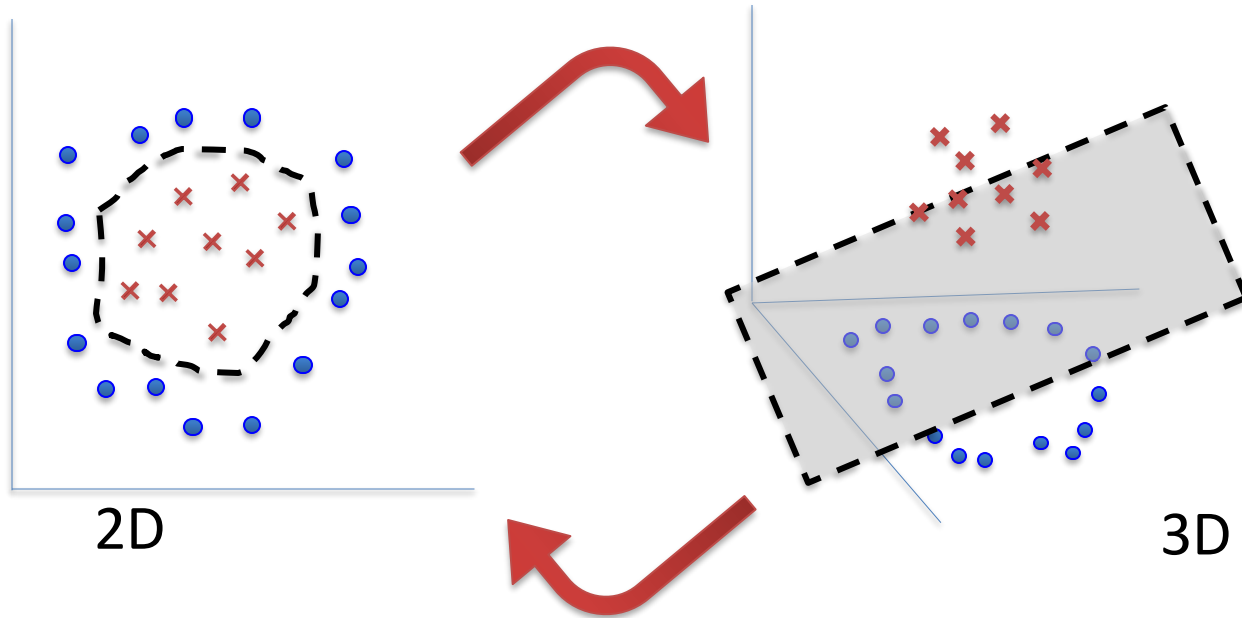
# What If Data Is More Complex?





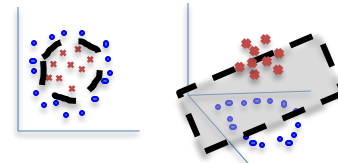
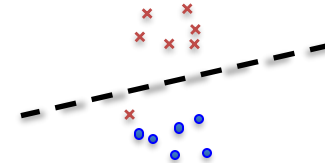
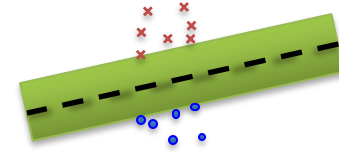
# What If Data Is More Complex? Use Kernels!

- Project it onto a higher-dimensional space, and optimize a linear model
  - This amounts to learning a non-linear model in the input space



# SVMs – Slightly More Technical Terms

- Find the boundary with **maximum margin**
- Allow **soft-margin** violations to deal with outliers
- Use **kernels**, and the **kernel trick**, to solve nonlinear problems

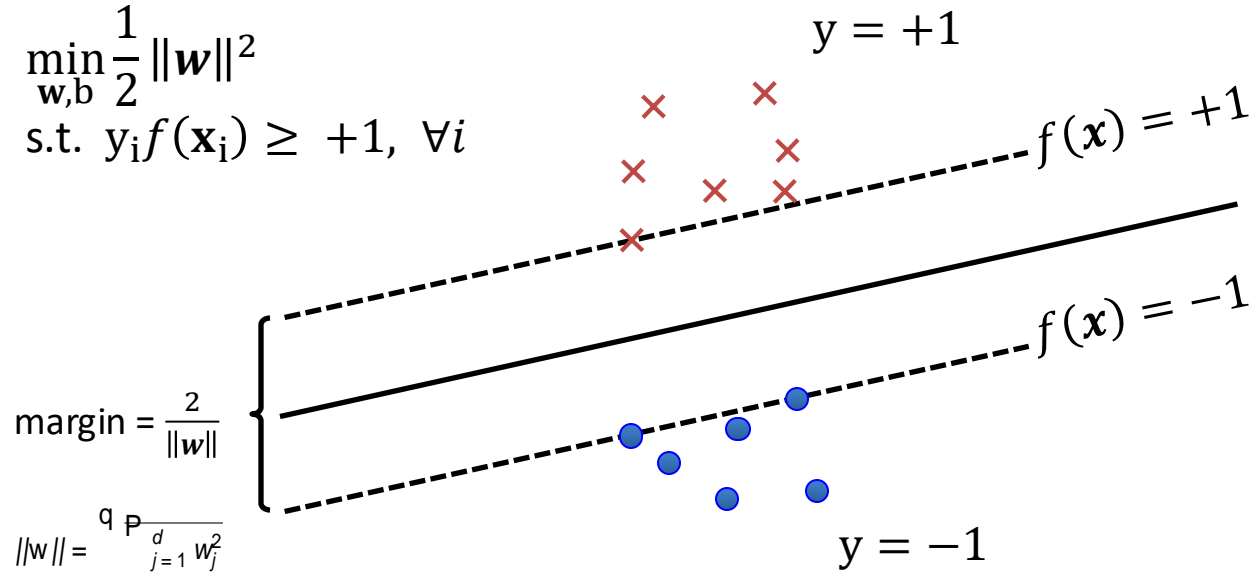


# Support Vector Machines

The “*some maths*” version

# Hard-margin SVM (linearly-separable data)

- All points should be correctly classified (within the margin)



# The Problem of Margin Maximization

$$\text{minimise } \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i f(\mathbf{x}_i) \geq 1, \quad i = 1, \dots, N$$

- The above problem can be solved with the classical Lagrange optimization technique

$$L = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^{\ell} \alpha_i \left[ y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \right]$$

$$\max_{\alpha_1, \dots, \alpha_{\ell}} \min_{\mathbf{w}, b} L(\alpha_1, \dots, \alpha_{\ell}, \mathbf{w}, b)$$

$$\text{s.t. } \alpha_i \geq 0, \quad i = 1, \dots, N$$

# Deriving the Dual Form

- At the optimum, the derivatives of  $L$  w.r.t.  $\mathbf{w}$  and  $b$  vanish, which gives:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \sum_i \alpha_i y_i = 0.$$

- Substituting into  $L$ , we obtain the **dual problem for the hard-margin SVM**:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i \alpha_i \mathbf{x}_i^T \mathbf{x}_j \alpha_j y_j \\ \text{s.t.} \quad & \alpha_i \geq 0, \forall i \\ & \sum_i \alpha_i y_i = 0, \forall i \end{aligned}$$

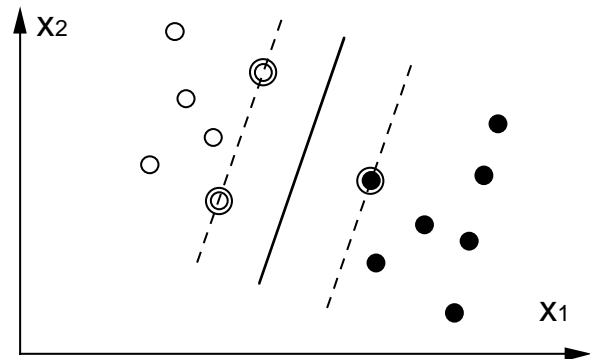
# Support Vectors

- Lagrangian of the dual problem:  $L_D(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i \alpha_i \mathbf{x}_i^T \mathbf{x}_j \alpha_j y_j$

- The optimal solution is the hyperplane:

$$\sum_{i=1}^{\ell} y_i \alpha_i (\mathbf{x}_i \cdot \mathbf{x}) + b = 0$$

- The samples for which  $\alpha_i \neq 0$  are called the **support vectors**



# Soft-margin SVM

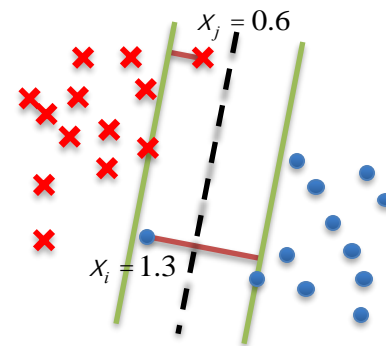
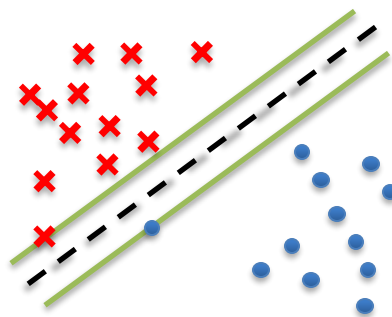
- What if data is not linearly separable / outliers? We allow points to violate the margin
- Trade-off between margin and loss on training data, tuned via the hyperparameter  $C$

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

$$\text{s.t. } y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \forall i$$

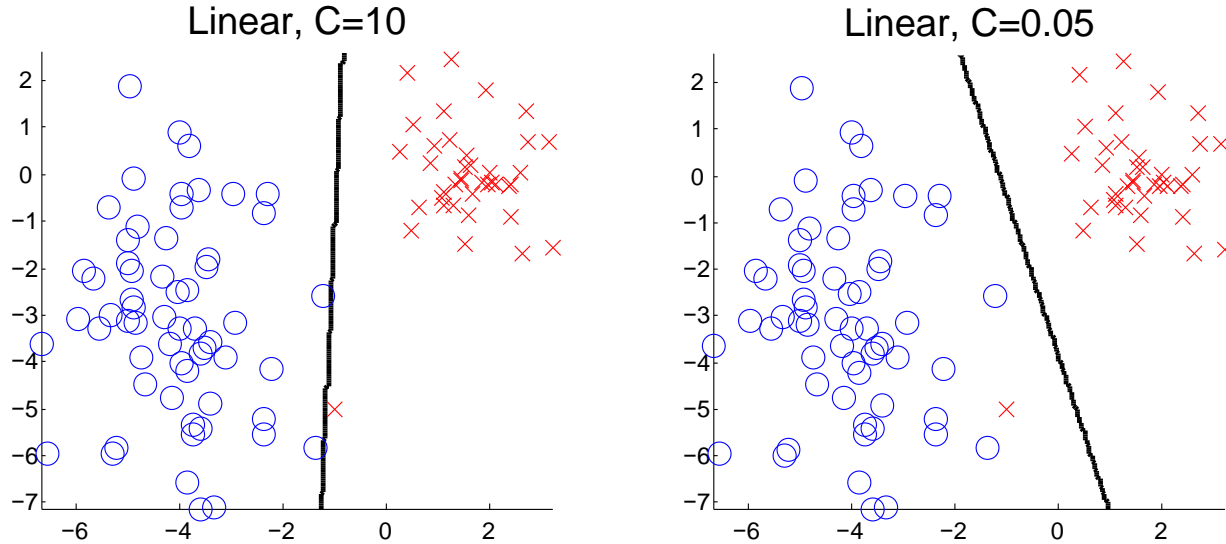
$$\xi_i \geq 0, \forall i$$

$\xi_i$ : slack variables





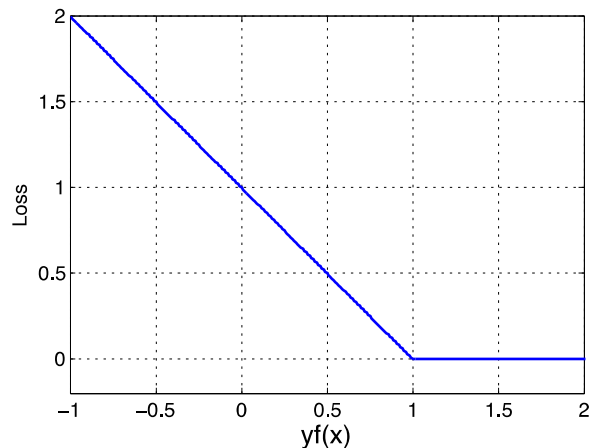
# Effect of the Hyperparameter C



# Slack Variables and the Hinge Loss

- The constraints  $\xi_i \geq 1 - y_i f(\mathbf{x}_i)$ ,  $\xi_i \geq 0$  are equivalent to  $\xi_i = \max(0, 1 - y_i f(\mathbf{x}_i))$
- That's exactly the **hinge loss**!
- We can use it to derive an unconstrained version of the primal problem:

$$\underbrace{\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2}_{\text{margin term}} + C \underbrace{\sum_i \max(0, 1 - y_i f(\mathbf{x}_i))}_{\substack{\text{hinge loss} \\ \text{(summed over training points)}}}$$



# Dual form of the Soft-margin SVM

- Replacing the optimality conditions

$$\frac{\partial L_P}{\partial w_\nu} = w_\nu - \sum_i \alpha_i y_i x_{i\nu} = 0 \quad \frac{\partial L_P}{\partial b} = - \sum_i \alpha_i y_i = 0 \quad \frac{\partial L_P}{\partial \xi_i} = C - \alpha_i - \mu_i = 0$$

- ... into the Lagrangian of the primal problem,

$$L_P = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i - \sum_i \alpha_i \{y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i\} - \sum_i \mu_i \xi_i$$

- we obtain:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i \alpha_i \mathbf{x}_i^T \mathbf{x}_j \alpha_j y_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \forall i \\ & \sum_i \alpha_i y_i = 0, \forall i \end{aligned}$$

# Sparsity of Support Vectors in SVMs

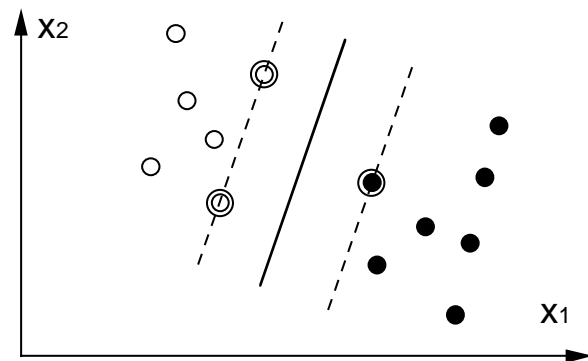
- Why do SVMs find only few support vectors?

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i \alpha_i \mathbf{x}_i^T \mathbf{x}_j \alpha_j y_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \forall i \\ & \sum_i \alpha_i y_i = 0, \forall i \end{aligned}$$

- The dual problem has a sparse solution
  - Many of the  $\alpha$  values will be zero (complementary slackness)
- Karush-Kuhn-Tucker (equilibrium) conditions:

$$g_i = \frac{\partial W}{\partial \alpha_i} = \sum_j Q_{ij} \alpha_j + y_i b - 1 = y_i f(\mathbf{x}_i) - 1 \quad \begin{cases} \geq 0; & \alpha_i = 0 \\ = 0; & 0 < \alpha_i < C \\ \leq 0; & \alpha_i = C \end{cases} \quad \begin{array}{l} \text{Reserve vectors} \\ \text{Margin SVs} \\ \text{Error SVs} \end{array}$$

$$\frac{\partial W}{\partial b} = \sum_j y_j \alpha_j = 0$$



# How to Solve Primal and Dual SVM learning?

- These are both Quadratic Programming (QP) problems

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \quad \forall i \\ & \xi_i \geq 0, \quad \forall i \end{aligned}$$

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i \alpha_i \mathbf{x}_i^T \mathbf{x}_j \alpha_j y_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad \forall i \\ & \sum_i \alpha_i y_i = 0, \quad \forall i \end{aligned}$$

- State-of-the-art QP solvers are standard and efficient
- However, dedicated (and much more efficient) solvers have been developed specifically for SVMs
  - **Sequential Minimal Optimization (SMO)** is one of them, used also in the popular LibSVM library

# How to Solve Primal and Dual SVM learning?

- SMO however does not scale efficiently with the training set size
- Not appropriate to deal with modern, large data sets
- Modern techniques optimize directly the primal (unconstrained) form using Stochastic Gradient Descent (SGD) and subgradients
- Subgradients overcome the non-differentiability of the hinge loss when  $yf(\mathbf{x}) = 1$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \max(0, 1 - y_i f(\mathbf{x}_i))$$

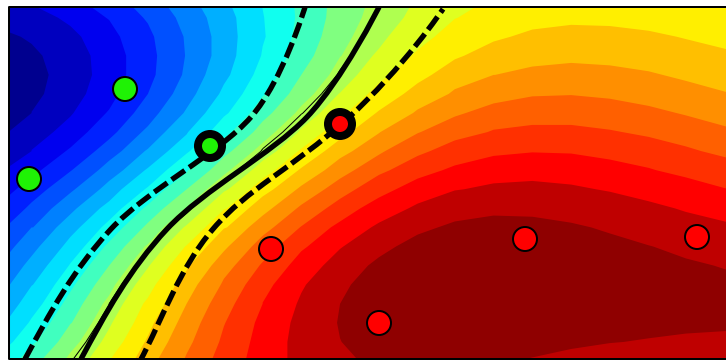
# Stochastic Gradient Descent (SGD)

- Set the learning rate  $\eta$  (using decay if needed)
- Repeat until an approximate minimum is obtained:
  - Randomly select  $K$  samples from the training data (i.e., one batch)
  - Update parameters  $w' = w - \eta \nabla L$  using the  $K$  samples
- **Pros:** several variants of this algorithm have been proposed
  - They are very efficient – and incremental
  - Data loading in batches avoids keeping the full dataset in memory
  - Convergence is well studied and often satisfied in practice
- **Cons:** parameters (e.g., learning rate) can be difficult to tune

# The Kernel Trick

- Soft-margin SVM deals with non-linearly separable data but it is still a linear classifier
  - It exhibits poor performance if data is not properly shaped or “mostly” linearly separable
- The *kernel trick* overcomes this limitation by allowing us to learn an **SVM with a nonlinear decision function in input space!**

How?





# The Kernel Trick

- The dual SVM formulation reveals that both the learning problem and classification can be expressed only in terms of **scalar products** between samples!

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i \alpha_i \mathbf{x}_i^T \mathbf{x}_j \alpha_j y_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \forall i \\ & \sum_i \alpha_i y_i = 0, \forall i \end{aligned}$$

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b \\ &= \sum_i y_i \alpha_i \mathbf{x}_i^T \mathbf{x} \end{aligned}$$

- It also holds for the primal form, thanks to this property:
  - see the *Representer Theorem* for further details

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

# The Kernel Trick

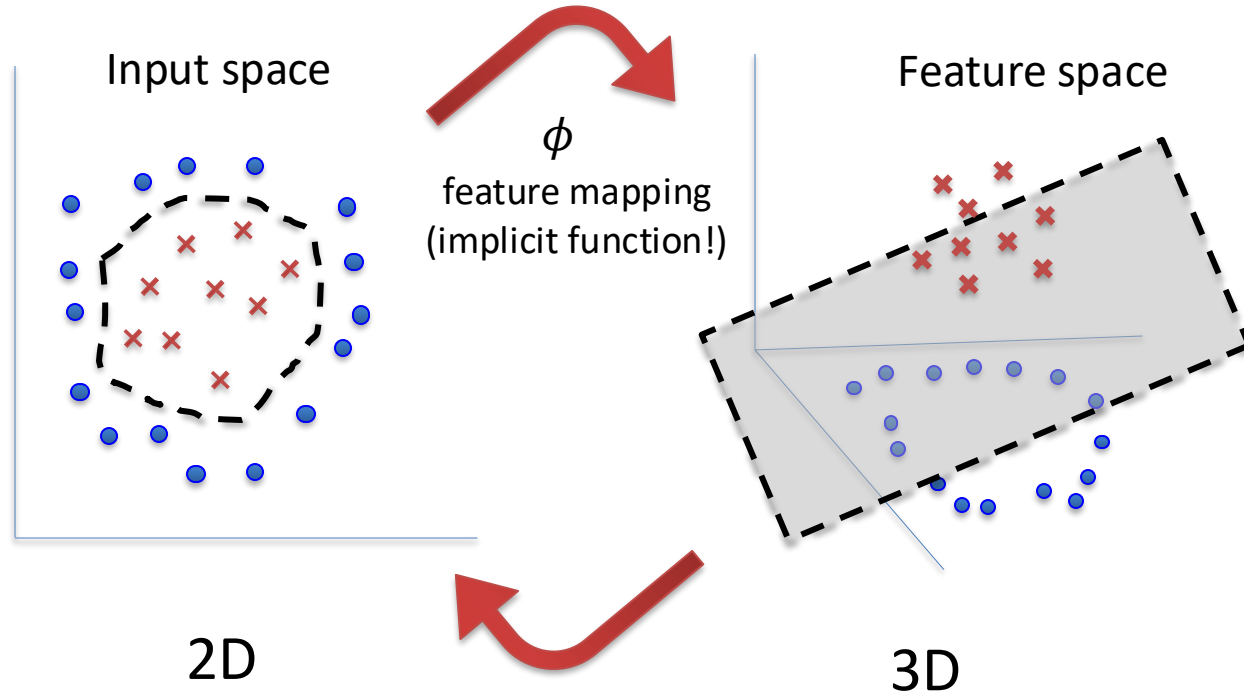
- We are not required to compute the scalar product in the input space
- We can use any function  $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  given that it implicitly corresponds to a scalar product in some other space (Mercer's condition)
- Kernel functions are symmetric and positive semi-definite (PSD) if

$$\sum_{i,j=1}^n c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0 \quad (1.1)$$

holds for any  $n \in \mathbb{N}, \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}, c_1, \dots, c_n \in \mathbb{R}$ .

- **Note:** SVMs also converge when using symmetric non-PSD kernels (although the learning problem is no longer convex)

# The Kernel Trick

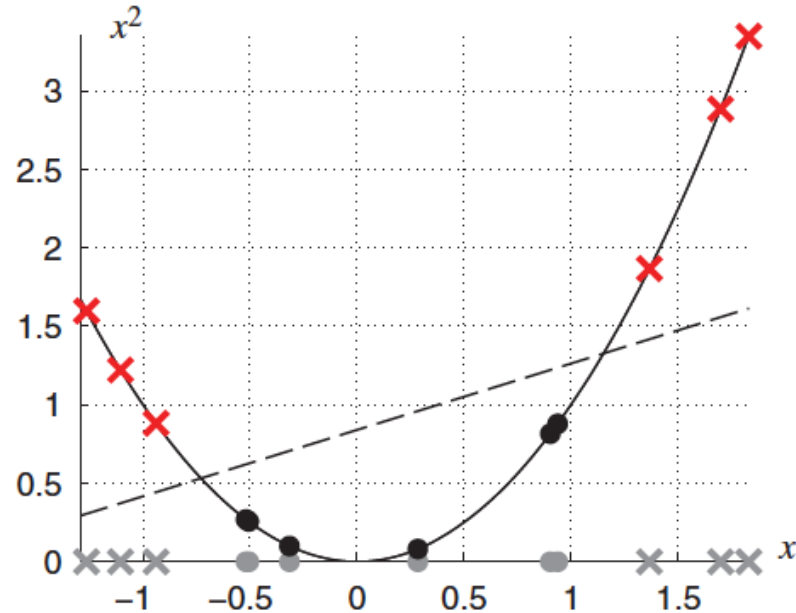


# Cover's Theorem on Separability of Patterns

- *“A complex pattern-classification problem cast in a high-dimensional space non-linearly is more likely to be linearly separable than in a low-dimensional space”*
- The power of SVMs resides in the fact that they represent a robust and efficient implementation of Cover's theorem
- SVMs operate in two stages
  - Perform a non-linear mapping of the input vector  $x$  onto a high-dimensional space that is hidden from the inputs or the outputs
  - Construct an optimal separating hyperplane in the high-dimensional feature space

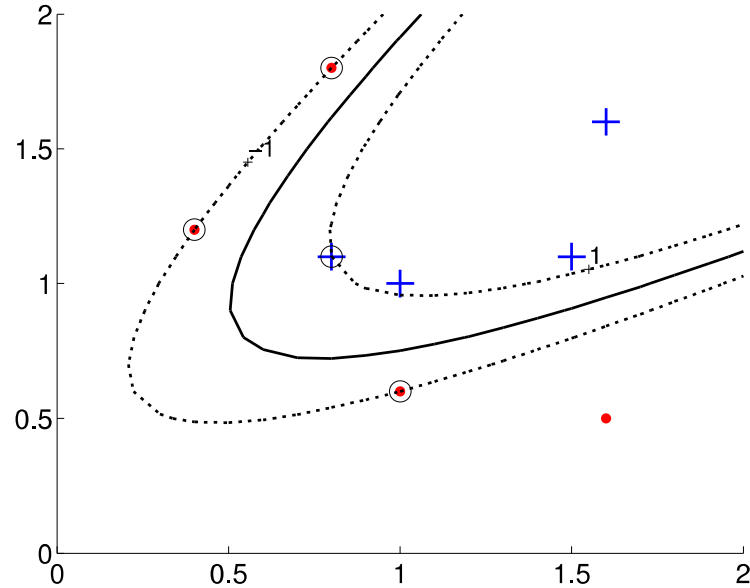
# Cover's Theorem on Separability of Patterns

- “A complex pattern-classification problem cast in a high-dimensional space non-linearly is more likely to be linearly separable than in a low-dimensional space”



# Polynomial Kernel, with d=2

$$K(\mathbf{x}_i, \mathbf{x}') = (1 + \mathbf{x}_i^T \mathbf{x}')^d$$



# The Polynomial Kernel

$$K(\mathbf{x}_i, \mathbf{x}') = (1 + \mathbf{x}_i^T \mathbf{x}')^d$$

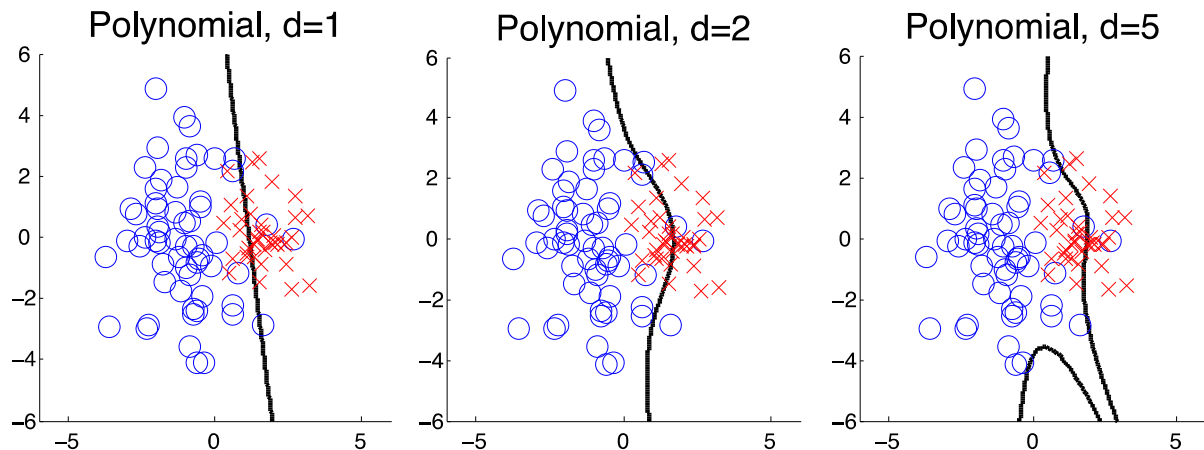
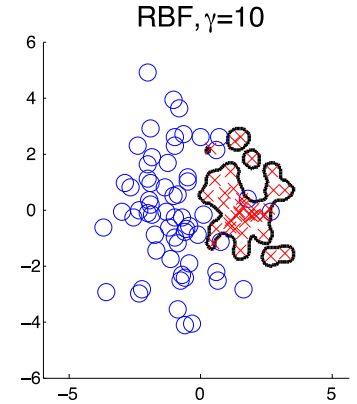
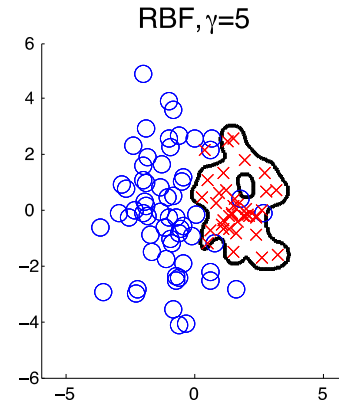
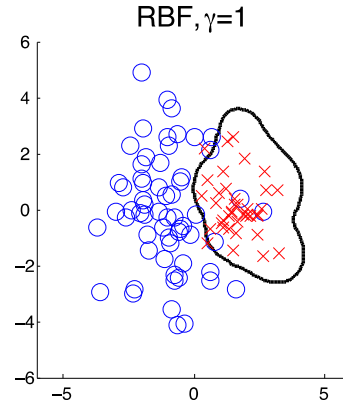
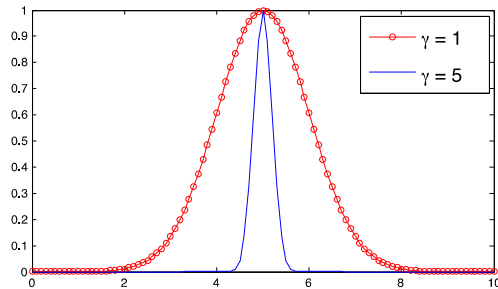


Figure 4.6: The effect of the degree parameter when using a polynomial kernel.

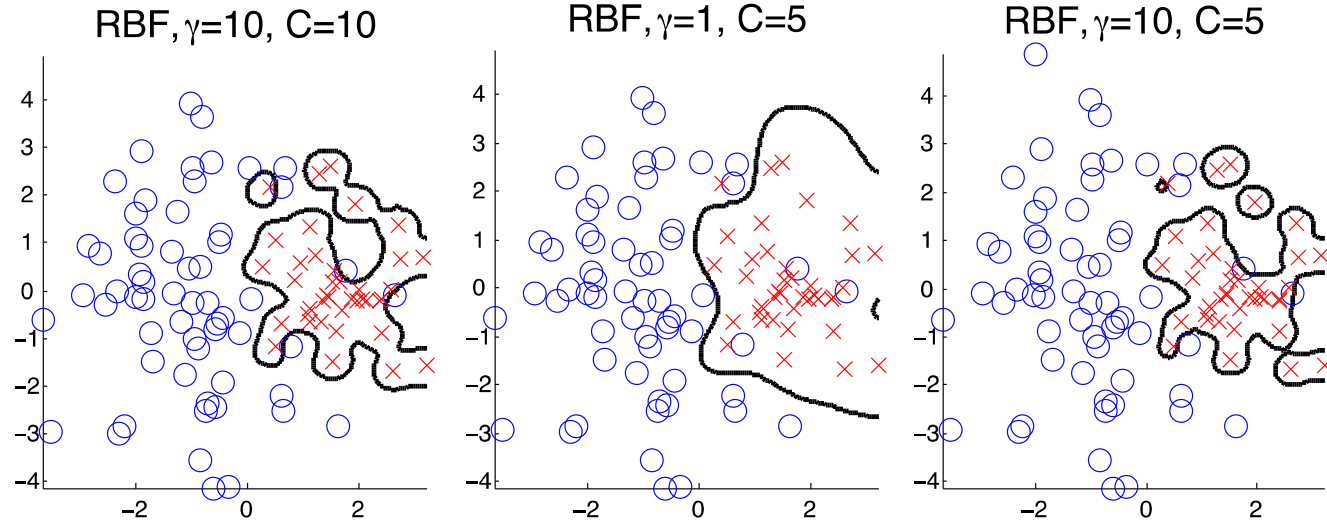
# The Gaussian Kernel (Radial Basis Function, RBF)

$$K(\mathbf{x}_i, \mathbf{x}') = e^{-\gamma(\mathbf{x}_i - \mathbf{x}')^2}$$





# Playing with the Hyperparameters



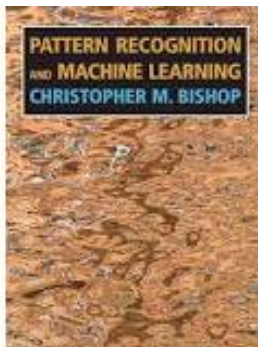
# Suggested Readings

## A Tutorial on Support Vector Machines for Pattern Recognition

**CHRISTOPHER J.C. BURGES**

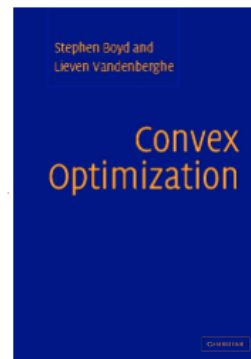
*Bell Laboratories, Lucent Technologies*

burges@lucent.com



Pattern Recognition and  
Machine Learning  
Book by Christopher Bishop

Chapter 7



*Convex Optimization*  
Stephen Boyd and Lieven Vandenberghe

Cambridge University Press

# More on Linear Classifiers

# More on Linear Classifiers

- The unconstrained version of the SVM primal problem...

$$\min_{\mathbf{w}, b} C \sum_i \max(0, 1 - y_i f(\mathbf{x}_i)) + \frac{1}{2} \|\mathbf{w}\|^2$$

- ... can be seen as an instance of a more general problem

$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i))}_{\text{loss term}} + \underbrace{\lambda \Omega(\mathbf{w})}_{\text{regularization term}}$$

(empirical loss on training data)      (controls overfitting / classifier capacity)

$\lambda$  : trade-off parameter

# More on Linear Classifiers

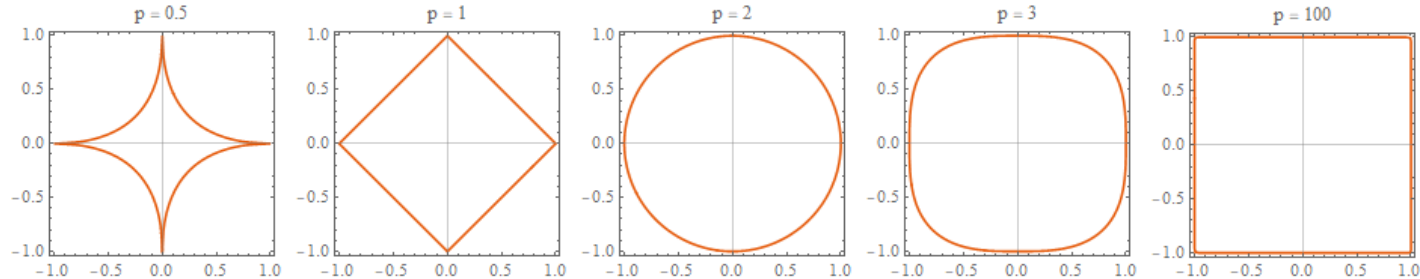
- This objective function amounts to minimizing the generalization error on test data
  - see C. Burges' tutorial and Statistical Learning Theory by V. Vapnik

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)) + \lambda \Omega(\mathbf{w})$$

- Different choices of loss and regularization yield different learning algorithms for linear classifiers – all trainable with SGD (and scikit-learn). Some examples are:
  - Hinge + l2: SVM
  - Hinge + l1: 1-norm SVM
  - Squared loss + l2: Ridge Regression
  - Squared loss + l1: LASSO
  - Logistic loss + l2: Logistic Regression

# Regularizers

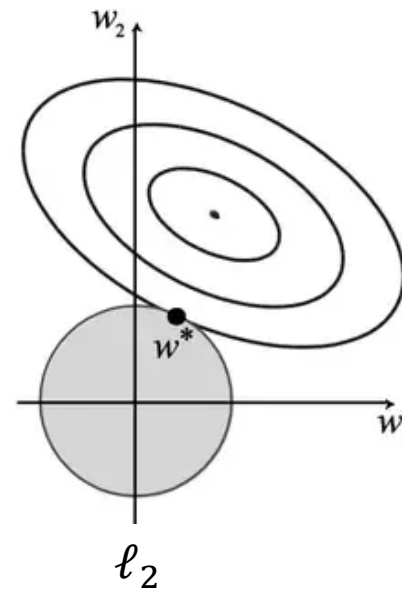
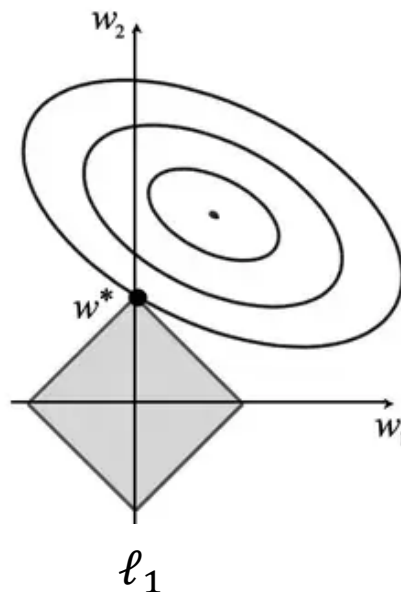
- $\Omega(\mathbf{w})$ : typically convex penalties on the classifier parameters are used
  - $\ell_p$  norms with  $p \geq 1$ ,  $\ell_p(\mathbf{w}) = \left(\sum_j |w_j|^p\right)^{1/p}$
- Most popular examples
  - $\ell_0$  is not convex, and amounts to counting non-zero elements in  $\mathbf{w}$
  - $\ell_1 = |w_1| + |w_2| + \dots + |w_d|$
  - $\ell_2 = w_1^2 + w_2^2 + \dots + w_d^2$
  - $\ell_\infty = \max_j |w_j|$



# Sparsity

- $\ell_0$  and  $\ell_1$  regularization enforce *sparsity*, i.e., many values in  $\mathbf{w}$  will be set to zero
  - Why? The optimum is often found at one of the vertices!

- Sparsity helps automatically perform feature selection
- Features assigned  $w_j = 0$  can be disregarded



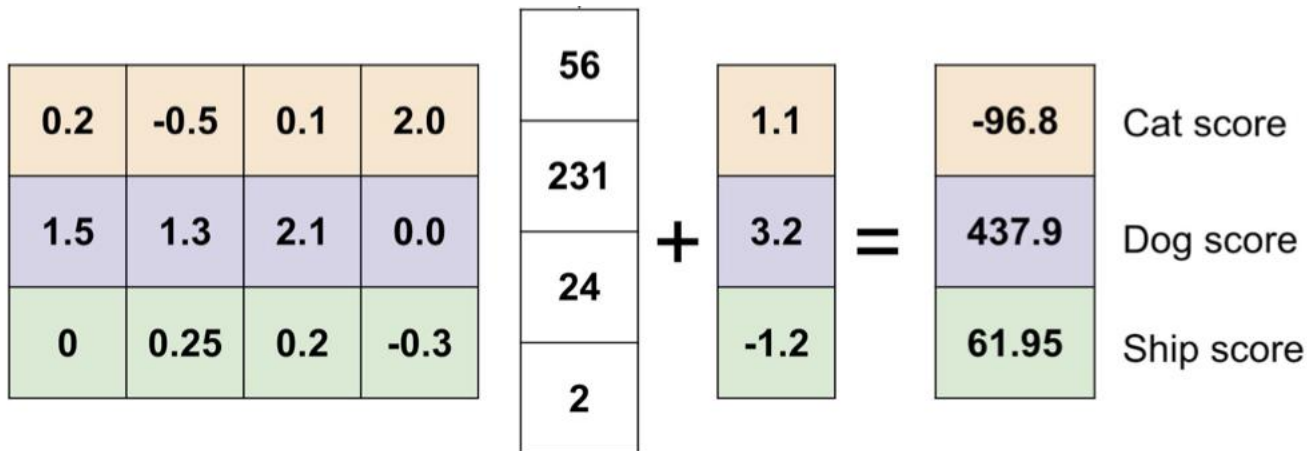
# **Multiclass Classification with Linear Classifiers**



# Multiclass Linear Classifiers

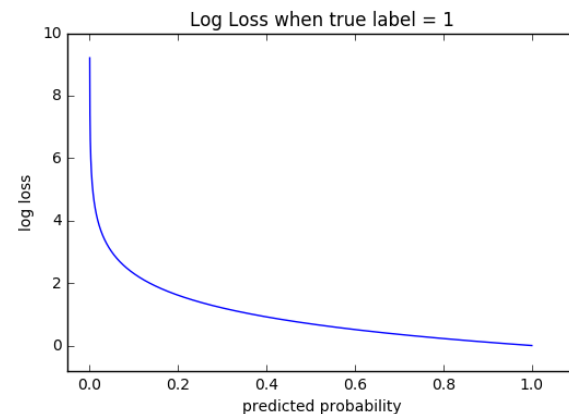
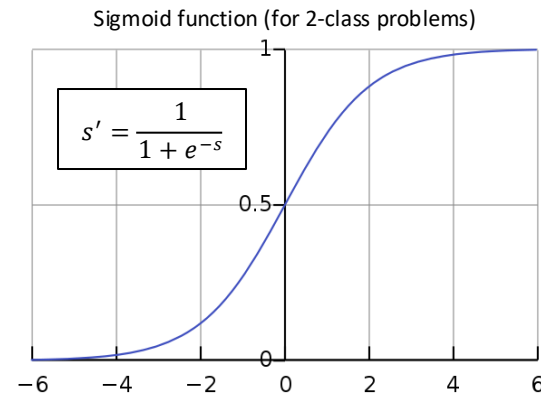
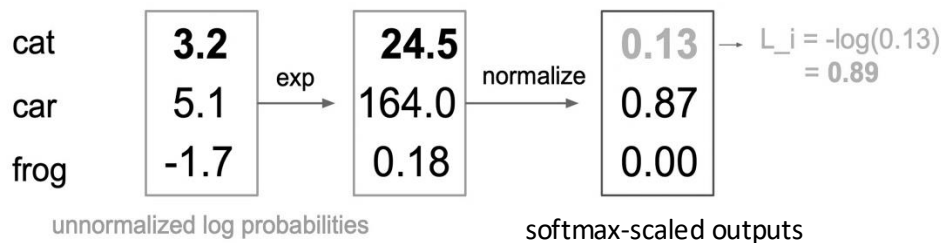
- Linear functions can be also naturally extended to multiclass problems

$$f(x) = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$



# Logistic/Softmax Classifier

- Multiclass loss used to learn  $\mathbf{W}, \mathbf{b}$
- Softmax scaling and cross-entropy (or log) loss
  - The classifier outputs one score per class  $f(\mathbf{x}) = (s_1, \dots, s_k)$
  - Outputs are softmax scaled:  $s'_l = \frac{e^{s_l}}{\sum_j e^{s_j}}$
  - Cross-entropy loss:  $L(y_i, f(\mathbf{x}_i)) = -\log(s'_{y_i})$



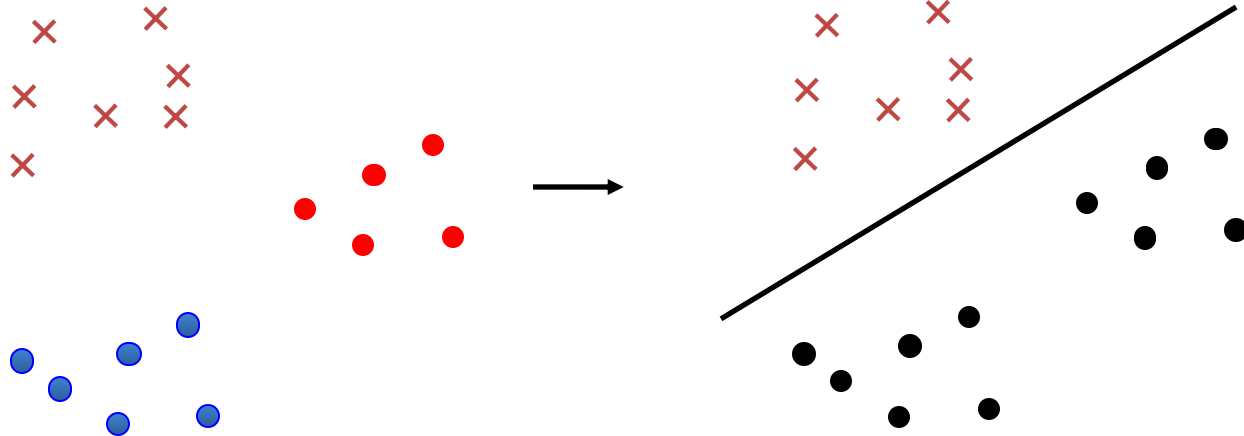
# **Multiclass Classification with Binary Classifiers**

# Multiclass Classification with Binary (two-class) Classifiers

- How to use binary classifiers (not necessarily linear) for multi-class classification ( $c > 2$ )?
- Recall that one object belongs only to one class!
  - $y$  is in the set  $\{0, 1, \dots, c-1\}$
- Two common schemes:
  - One-vs-all (or one-vs-rest) classification (OVA / OVR)
  - One-vs-one classification (OVO)

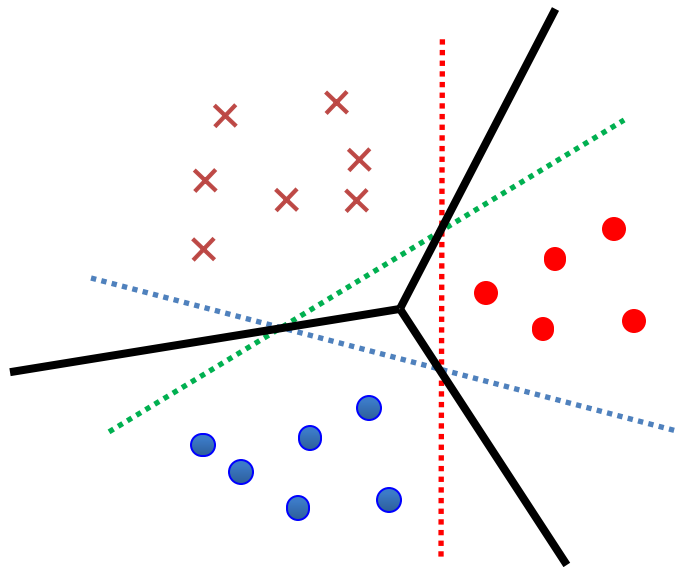
# One-vs-all Multiclass Classification

- Train one binary classifier for each class  $k$
- Samples for which  $y=k$  are labeled as +1
- The remaining classes are all labeled as -1



# One-vs-all Multiclass Classification

- Combine with  $y = \operatorname{argmax}_k f_k(\mathbf{x})$



# One-vs-one Multiclass Classification

- Train one binary classifier for class  $i$  vs class  $j$
- All possible pairs are considered
  - $c(c-1)/2$  total number of binary classifiers
- Combined as:

$$f(\mathbf{x}) = \arg \max_i \left( \sum_j f_{ij}(\mathbf{x}) \right)$$

Note that  $f_{ij} = -f_{ji}$  and that the sum is for  $j=0, \dots, c-1$

# Pros and Cons

- One-vs-all classification
  - One classifier per class, trained using all data
- One-vs-one classification
  - combinatorial number of classifiers (all possible pairs,  $c(c-1)/2$ ), trained using much smaller data subsets
- Trade-off between number of classifiers and complexity of the learning algorithm
- In practice, rather equivalent classification accuracies
  - No free lunch



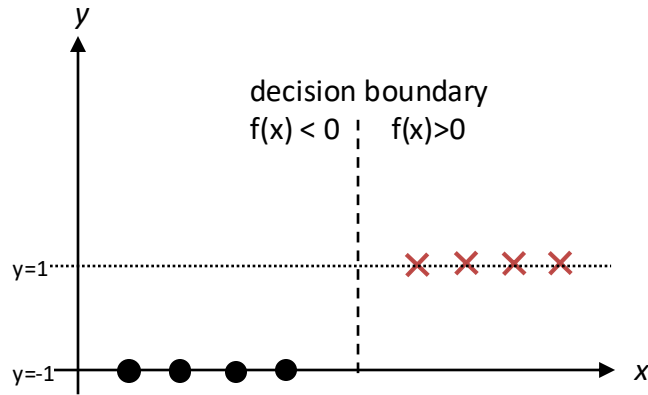
# **Classification vs Regression (with linear estimators)**

# Classification vs Regression

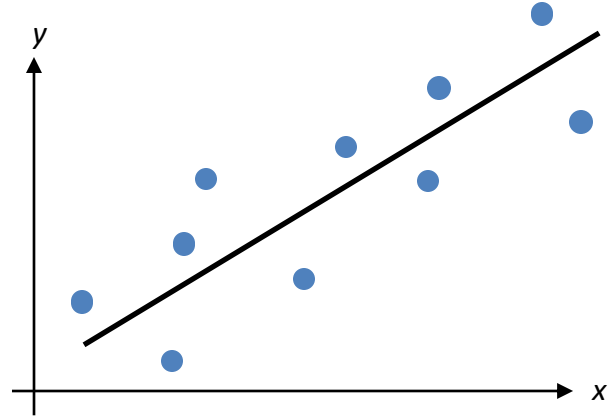
- In both classification and regression, the goal is to estimate a function  $f(\mathbf{x})$  that estimates the truth value  $y$  for each given sample  $\mathbf{x}$
- **In this chapter we've discussed linear functions**
- In classification problems, the set of labels is discretized
  - $y \in \{0, \dots, c - 1\}$
- In regression problems, the set of labels (or better, target values) is typically a continuous value
  - $y \in \mathbb{R}$  (e.g., the set of real numbers)

# Classification vs Regression

Classification (estimates  $y$  in a discrete set, e.g.,  $\{-1,1\}$ )

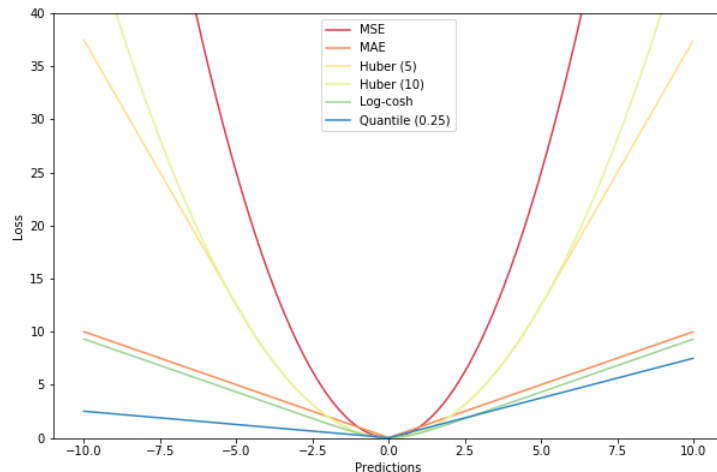


Regression estimates continuous-valued  $y$



# Classification vs Regression

- The loss functions used in the two problems reflect this behavior
- For classification problems, correctly-classified points are assigned a loss equal to zero
  - e.g., the hinge loss gives zero penalty to points for which  $yf(\mathbf{x}) \geq 1$
- For regression problems, the loss is zero only if  $f(\mathbf{x})$  is exactly equal to  $y$ 
  - e.g., the mean squared error (MSE) is given as the average of  $(y - f(\mathbf{x}))^2$  over all points



# Ridge Regression

- It uses the mean squared error (MSE) as the error function and l2 regularization on the feature weights:

$$L(\mathbf{w}) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2$$

- Minimizing  $L(\mathbf{w})$  provides the following closed-form solution:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} ,$$

- being  $\mathbf{I}$  the identity matrix, and  $\lambda > 0$  a trade-off parameter.
- In this case, adding a small diagonal (ridge) to the (positive semi-definite) matrix  $\mathbf{X}^T \mathbf{X}$  makes it more stable for pseudo-inversion (as it increases its minimum eigenvalue)

# Ridge Regression

- Ridge regression can be solved in closed form, through matrix pseudo-inversion
  - Too computationally demanding for large feature sets and datasets
- It is also possible to solve it using gradient-descent procedures, including SGD, which is much faster and better suited to large, high-dimensional training sets

# LASSO and Elastic Net

- Least Absolute Shrinkage and Selection Operator (LASSO) uses sparse penalty on  $\mathbf{w}$  for embedded feature selection

$$L(\mathbf{w}) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$

- The Elastic Net overcomes some issues arising when learning LASSO on badly-conditioned problems, by combining l1 and l2 regularization

$$L(\mathbf{w}) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2$$

# References

- Pattern Classification, R. O. Duda, P. E. Hart, and D. G. Stork, John Wiley & Sons, 2000
- L. Kuncheva, Combining pattern classifiers, Wiley, 2004.
- C. Bishop, Pattern Recognition and Machine Learning, Springer, 2007.
- S. Boyd and L. Vandenberghe, Convex Optimization, Cambridge University Press, 2004.
  
- **Acknowledgments:** I would like to thank Prof. Gavin Brown for the permission to use the material on SVMs from his course on ML at the University of Manchester, UK.



# Exercises

# Exercise 1

- Consider the data

$$\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ -1 & -2 \\ -2 & -1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$

- and the linear classification function  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ , with  $\mathbf{w} = [1 \ 1]^T$  and  $b = 0$ .
- Points are classified as positive if  $f(\mathbf{x}) \geq 0$
- Classify the data points in  $\mathbf{X}$  and compute the overall classification error

# Exercise 1: Solution

- Samples in  $\mathbf{X}$  are stored as rows, hence it is not difficult to compute the values of the discriminant function as

$$f(\mathbf{X}) = \mathbf{X}\mathbf{w} + b = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ -1 & -2 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \\ -3 \\ -3 \end{bmatrix}$$

- which, compared against 0, gives the classification labels

$$\mathbf{y}_c = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$

- Since they are exactly equal to the true labels  $\mathbf{y}$ , the classification error is zero

## Exercise 2

- Given the following data samples

$$\mathbf{X} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 0 & 0 \\ 0 & 2 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \end{bmatrix},$$

- find the linear discriminant function using the criterion function (perceptron loss):

$$L_p(\mathbf{w}, b) = \sum_i \max(0, -y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

- Initialize  $\mathbf{w} = [0.1, 0.1]^T$ ,  $b = 0.1$ ,  $\eta = 1$ ,  $\theta = 0$
- Use the termination condition:  $\eta(|\nabla_{\mathbf{w}} L_p(\mathbf{w}, b)| + |\nabla_b L_p(\mathbf{w}, b)|) \leq \theta$

## Exercise 2: Solution

```
begin initialize  $\mathbf{w}, b, \eta, \theta$ 
  do
    select samples for which  $y_i f(\mathbf{x}_i) < 0$ 
    update  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} L_p(\mathbf{w}, b)$  and  $b \leftarrow b - \eta \nabla_b L_p(\mathbf{w}, b)$ 
      (using the selected samples)
  until  $\eta(|\nabla_{\mathbf{w}} L_p(\mathbf{w}, b)| + |\nabla_b L_p(\mathbf{w}, b)|) < \theta$ 
end
```

- Here we have to compute the derivatives w.r.t  $\mathbf{w}$  and  $b$ :

$$\nabla_{\mathbf{w}} L_p(\mathbf{w}, b) = - \sum_{i: y_i f(\mathbf{x}_i) < 0} y_i \mathbf{x}_i$$

$$\nabla_b L_p(\mathbf{w}, b) = - \sum_{i: y_i f(\mathbf{x}_i) < 0} y_i$$

# Exercise 2: Solution

## Iteration 1

- We first have to classify the data points with the current  $\mathbf{w}$  and  $b$

$$f(\mathbf{X}) = \mathbf{X}\mathbf{w} + b = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 0 & 0 \\ 0 & 2 \\ 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} + 0.1 = \begin{bmatrix} 0 \\ 0 \\ 0.1 \\ 0.3 \\ 0.4 \\ 0.4 \end{bmatrix}$$

## Exercise 2: Solution

- Then, consider the classification errors for which  $y_i f(\mathbf{x}_i) < 0$

$$\mathbf{y}f(\mathbf{X}) = \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \end{bmatrix} \circ \begin{bmatrix} 0 \\ 0 \\ 0.1 \\ 0.3 \\ 0.4 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0.1 \\ -0.3 \\ -0.4 \\ -0.4 \end{bmatrix}$$

- The three last patterns are wrongly classified, so they will be used to update  $\mathbf{w}$  and  $b$

$$\nabla_{\mathbf{w}} L_p(\mathbf{w}, b) = -\sum_{i: y_i f(\mathbf{x}_i) < 0} y_i \mathbf{x}_i = +1 \begin{bmatrix} 0 \\ 2 \end{bmatrix} + 1 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1 \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\nabla_b L_p(\mathbf{w}, b) = -\sum_{i: y_i f(\mathbf{x}_i) < 0} y_i = 3$$

## Exercise 2: Solution

- **Parameter updates**

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} L_p(\mathbf{w}, b) \quad \text{and} \quad b \leftarrow b - \eta \nabla_b L_p(\mathbf{w}, b)$$

$$\mathbf{w} = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} - 1 \begin{bmatrix} 3 \\ 5 \end{bmatrix} = \begin{bmatrix} -2.9 \\ -4.9 \end{bmatrix}$$

$$b = 0.1 - 1 * 3 = -2.9$$

- **Termination condition**  $\eta(|\nabla_{\mathbf{w}} L_p(\mathbf{w}, b)| + |\nabla_b L_p(\mathbf{w}, b)|) \leq \theta$

$$t = 1 * (|3| + |5| + |3|) \leq \theta = 0$$

- Termination condition is *False* (then, loop continues)



## Exercise 2: Solution

iter	$yf(x)$	$L_p$	$\nabla_w L_p, \nabla_b L_p$	$w, b$	t	$\theta$
1.	[ 0. 0. 0.1 -0.3 -0.4 -0.4]	1.1	[3 5] 3	[-2.9 -4.9] -2.9	11.0	0
2.	[ 0. 2. -2.9 12.7 15.6 13.6]	2.9	[0 0] -1	[-2.9 -4.9] -1.9	1.0	0
3.	[ 1. 3. -1.9 11.7 14.6 12.6]	1.9	[0 0] -1	[-2.9 -4.9] -0.9	1.0	0
4.	[ 2. 4. -0.9 10.7 13.6 11.6]	0.9	[0 0] -1	[-2.9 -4.9] 0.1	1.0	0
5.	[ 3. 5. 0.1 9.7 12.6 10.6]	0.0	[0 0] 0	[-2.9 -4.9] 0.1	0.0	0

**Homework:** Plot the data points and the decision function at each iteration

## Exercise 3

- Given a regression problem and the following data samples

$$\mathbf{X} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 0 & 0 \\ 0 & 2 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix},$$

- find the linear discriminant function via ordinary least squares (OLS), i.e., by minimizing:

$$L_r(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + b - y_i)^2$$

- Initialize  $\mathbf{w} = [0.1, 0.1]^T$ ,  $b = 0.1$ ,  $\eta = 0.1$ ,  $\theta = 0.06$

## Exercise 3: Solution

- The solution (a.k.a. Widrow-Hoff or least-mean-squares) algorithm is

```
begin initialize  $\mathbf{w}, \theta, \eta, k=0$   
  repeat  
     $\mathbf{w} = \mathbf{w} - \eta \nabla_{\mathbf{w}} L_r(\mathbf{w}, b)$   
     $b = b - \eta \nabla_b L_r(\mathbf{w}, b)$   
  until  $\eta (|\nabla_{\mathbf{w}} L_r(\mathbf{w}, b)| + |\nabla_b L_r(\mathbf{w}, b)|) \leq \theta$ 
```

- The derivatives w.r.t  $\mathbf{w}$  and  $b$ :

$$\begin{aligned}\nabla_{\mathbf{w}} L_r(\mathbf{w}, b) &= \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + b - y_i) \mathbf{x}_i \\ \nabla_b L_r(\mathbf{w}, b) &= \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + b - y_i)\end{aligned}$$

## Exercise 3: Solution

iter	$f(\mathbf{x})$	$L_r$	$\nabla_{\mathbf{w}}L_r, \nabla_bL_r$	$\mathbf{w}, b$	t	$\theta=0.06$
1.	$0.1*[0 \ 0 \ 1 \ 3 \ 4 \ 4]$	0.240	$[1. \ 1.4], 0.6$	$[0. \ -0.04], 0.04$	0.3	
2.	$0.01*[4 \ 8 \ 4 \ -4 \ -4 \ 0]$	0.057	$[-0.28 \ -0.64], -0.52$	$[0.028 \ 0.024], 0.092$	0.14	
3.	$0.01*[6 \ 7 \ 9 \ 14 \ 17 \ 17]$	0.014	$[0.25 \ 0.32], 0.104$	$[0.0032 \ -0.008], 0.0816$	0.07	
4.	$0.01*[8 \ 9 \ 8 \ 7 \ 7 \ 8]$	0.003	$[-0.05 \ -0.14], -0.136$	$[0.008 \ 0.006], 0.09$	0.03	

**Homework:** Plot the data points and the decision function at each iteration