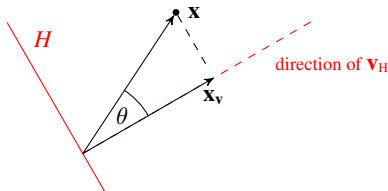


WHY *minimize* $\|\mathbf{v}_H\|$?

We can project a vector \mathbf{x} (think: data point) onto the direction of \mathbf{v}_H and obtain a vector \mathbf{x}_v .



- ▶ If H has no offset ($c = 0$), the Euclidean distance of \mathbf{x} from H is

$$d(\mathbf{x}, H) = \|\mathbf{x}_v\| = \cos \theta \cdot \|\mathbf{x}\| .$$

It does not depend on the length of \mathbf{v}_H .

- ▶ The scalar product $\langle \mathbf{x}, \mathbf{v}_H \rangle$ *does* increase if the length of \mathbf{v}_H increases.
- ▶ To compute the distance $\|\mathbf{x}_v\|$ from $\langle \mathbf{x}, \mathbf{v}_H \rangle$, we have to scale out $\|\mathbf{v}_H\|$:

$$\|\mathbf{x}_v\| = \cos \theta \cdot \|\mathbf{x}\| = \frac{\langle \mathbf{x}, \mathbf{v}_H \rangle}{\|\mathbf{v}_H\|}$$

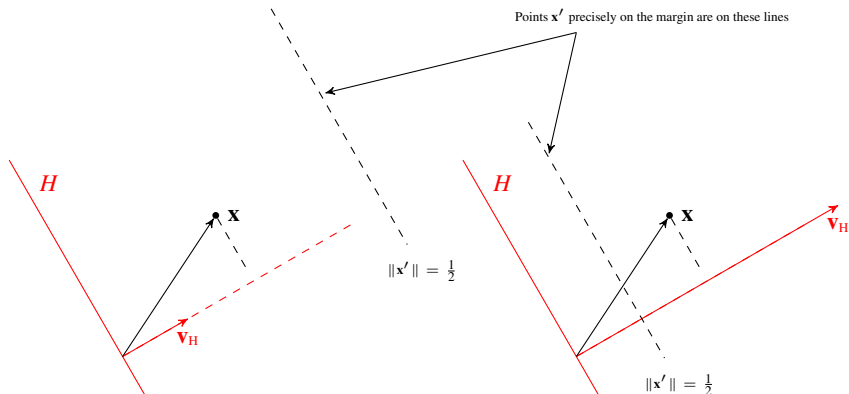
WHY *minimize* $\|\mathbf{v}_H\|$?

If we scale \mathbf{v}_H by α , we have to scale \mathbf{x} by $1/\alpha$ to keep $\langle \mathbf{v}_H, \mathbf{x} \rangle$ constant, e.g.:

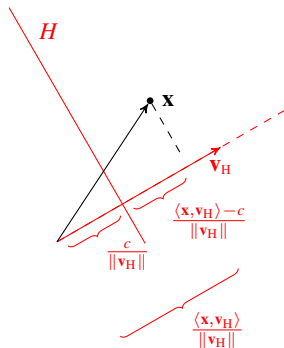
$$1 = \langle \mathbf{v}_H, \mathbf{x} \rangle = \langle \alpha \mathbf{v}_H, \frac{1}{\alpha} \mathbf{x} \rangle .$$

A point \mathbf{x}' is precisely on the margin if $\langle \mathbf{x}', \mathbf{v}_H \rangle = 1$.

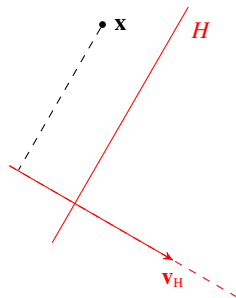
Look at what happens if we scale \mathbf{v}_H :



DISTANCE WITH OFFSET



For an affine plane, we have to subtract the offset.



The optimization algorithm can also rotate the vector \mathbf{v}_H , which rotates the plane.

SOFT-MARGIN CLASSIFIERS

Soft-margin classifiers are maximum-margin classifiers which permit some points to lie on the wrong side of the margin, or even of the hyperplane.

Motivation 1: Nonseparable data

SVMs are linear classifiers; without further modifications, they cannot be trained on a non-separable training data set.

Motivation 2: Robustness

- ▶ Recall: Location of SVM classifier depends on position of (possibly few) support vectors.
- ▶ Suppose we have two training samples (from the same joint distribution on (X, Y)) and train an SVM on each.
- ▶ If locations of support vectors vary significantly between samples, SVM estimate of \mathbf{v}_H is “brittle” (depends too much on small variations in training data). \rightarrow Bad generalization properties.
- ▶ Methods which are not susceptible to small variations in the data are often referred to as **robust**.

SLACK VARIABLES

Idea

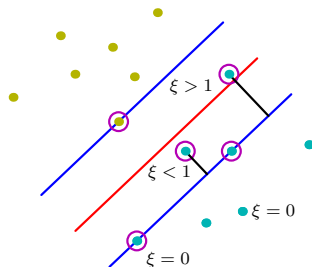
Permit training data to cross the margin, but impose cost which increases the further beyond the margin we are.

Formalization

We replace the training rule $\tilde{y}_i(\langle \mathbf{v}_H, \tilde{\mathbf{x}}_i \rangle - c) \geq 1$ by

$$\tilde{y}_i(\langle \mathbf{v}_H, \tilde{\mathbf{x}}_i \rangle - c) \geq 1 - \xi_i$$

with $\xi_i \geq 0$. The variables ξ_i are called **slack variables**.



Soft-margin optimization problem

$$\begin{aligned} \min_{\mathbf{v}_H, c} \quad & \|\mathbf{v}_H\|^2 + \gamma \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & \tilde{y}_i (\langle \mathbf{v}_H, \tilde{\mathbf{x}}_i \rangle - c) \geq 1 - \xi_i \quad \text{for } i = 1, \dots, n \\ & \xi_i \geq 0, \quad \text{for } i = 1, \dots, n \end{aligned}$$

The training algorithm now has a **parameter** $\gamma > 0$ for which we have to choose a “good” value. γ is usually set by a method called *cross validation* (discussed later). Its value is fixed before we start the optimization.

Role of γ

- Specifies the “cost” of allowing a point on the wrong side.
- If γ is very small, many points may end up beyond the margin boundary.
- For $\gamma \rightarrow \infty$, we recover the original SVM.

SOFT-MARGIN SVM

Soft-margin dual problem

The slack variables vanish in the dual problem.

$$\begin{aligned} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & W(\boldsymbol{\alpha}) := \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j \tilde{y}_i \tilde{y}_j (\langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle + \frac{1}{\gamma} \mathbb{I}\{i=j\}) \\ \text{s.t.} \quad & \sum_{i=1}^n \tilde{y}_i \alpha_i = 0 \\ & \alpha_i \geq 0 \quad \text{for } i = 1, \dots, n \end{aligned}$$

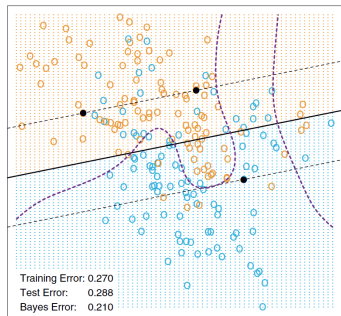
Soft-margin classifier

The classifier looks exactly as for the original SVM:

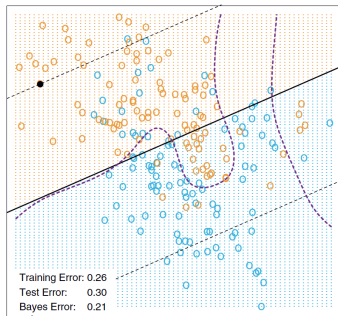
$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^n \tilde{y}_i \alpha_i^* \langle \tilde{\mathbf{x}}_i, \mathbf{x} \rangle - c \right)$$

Note: Each point on wrong side of the margin is an additional support vector ($\alpha_i^* \neq 0$), so the ratio of support vectors can be substantial when classes overlap.

INFLUENCE OF MARGIN PARAMETER



$\gamma = 100000$



$\gamma = 0.01$

Changing γ significantly changes the classifier (note how the slope changes in the figures). We need a method to select an appropriate value of γ , in other words: to learn γ from data.

TOOLS: OPTIMIZATION METHODS

OPTIMIZATION PROBLEMS

Terminology

An **optimization problem** for a given function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a problem of the form

$$\min_{\mathbf{x}} f(\mathbf{x})$$

which we read as "find $\mathbf{x}_0 = \arg \min_{\mathbf{x}} f(\mathbf{x})$ ".

A **constrained optimization problem** adds additional requirements on \mathbf{x} ,

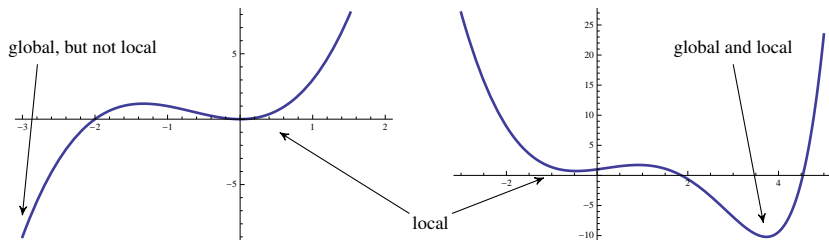
$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in G, \end{array}$$

where $G \subset \mathbb{R}^d$ is called the **feasible set**. The set G is often defined by equations, e.g.

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) \geq 0 \end{array}$$

The equation g is called a **constraint**.

TYPES OF MINIMA



Local and global minima

A minimum of f at x is called:

- ▶ **Global** if f assumes no smaller value on its domain.
- ▶ **Local** if there is some open neighborhood U of x such that $f(x)$ is a global minimum of f restricted to U .

Analytic criteria for local minima

Recall that \mathbf{x} is a local minimum of f if

$$f'(\mathbf{x}) = 0 \quad \text{and} \quad f''(\mathbf{x}) > 0 .$$

In \mathbb{R}^d ,

$$\nabla f(\mathbf{x}) = 0 \quad \text{and} \quad H_f(\mathbf{x}) = \left(\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) \right)_{i,j=1,\dots,n} \text{ positive definite.}$$

The $d \times d$ -matrix $H_f(\mathbf{x})$ is called the **Hessian matrix** of f at \mathbf{x} .

Numerical methods

All numerical minimization methods perform roughly the same steps:

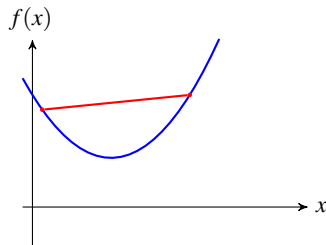
- ▶ Start with some point x_0 .
- ▶ Our goal is to find a sequence x_0, \dots, x_m such that $f(x_m)$ is a minimum.
- ▶ At a given point x_n , compute properties of f (such as $f'(x_n)$ and $f''(x_n)$).
- ▶ Based on these values, choose the next point x_{n+1} .

The information $f'(x_n), f''(x_n)$ etc is always *local at x_n* , and we can only decide whether a point is a local minimum, not whether it is global.

CONVEX FUNCTIONS

Definition

A function f is **convex** if every line segment between function values lies above the graph of f .



Analytic criterion

A twice differentiable function is convex if $f''(x) \geq 0$ (or $H_f(\mathbf{x})$ positive semidefinite) for all \mathbf{x} .

Implications for optimization

If f is convex, then:

- ▶ $f'(x) = 0$ is a sufficient criterion for a minimum.
- ▶ Local minima are global.
- ▶ If f is **strictly convex** ($f'' > 0$ or H_f positive definite), there is only one minimum (which is both global and local).

GRADIENT DESCENT

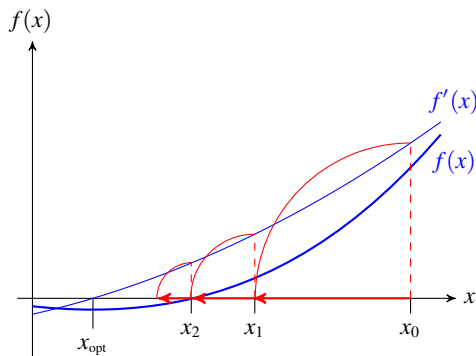
Algorithm

Gradient descent searches for a minimum of f .

1. Start with some point $x \in \mathbb{R}$ and fix a precision $\varepsilon > 0$.
2. Repeat for $n = 1, 2, \dots$

$$x_{n+1} := x_n - f'(x_n)$$

3. Terminate when $|f'(x_n)| < \varepsilon$.



NEWTON'S METHOD: ROOTS

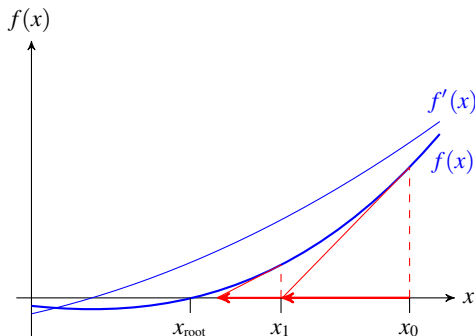
Algorithm

Newton's method searches for a **root** of f , i.e. it solves the equation $f(\mathbf{x}) = 0$.

1. Start with some point $x \in \mathbb{R}$ and fix a precision $\varepsilon > 0$.
2. Repeat for $n = 1, 2, \dots$

$$x_{n+1} := x_n - f(x_n)/f'(x_n)$$

3. Terminate when $|f(x_n)| < \varepsilon$.



Function evaluation

Most numerical evaluations of functions (\sqrt{a} , $\sin(a)$, $\exp(a)$, etc) are implemented using Newton's method. To evaluate g at a , we have to transform $x = g(a)$ into an equivalent equation of the form

$$f(x, a) = 0 .$$

We then fix a and solve for x using Newton's method for roots.

Example: Square root

To evaluate $g(a) = \sqrt{a}$, we can solve

$$f(x, a) = x^2 - a = 0 .$$

This is essentially how `sqrt()` is implemented in the standard C library.

NEWTON'S METHOD: MINIMA

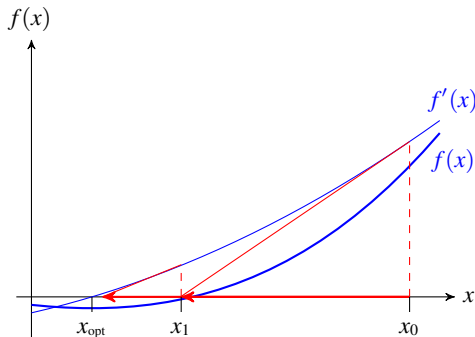
Algorithm

We can use Newton's method for minimization by applying it to solve $f'(\mathbf{x}) = 0$.

1. Start with some point $x \in \mathbb{R}$ and fix a precision $\varepsilon > 0$.
2. Repeat for $n = 1, 2, \dots$

$$x_{n+1} := x_n - f'(x_n)/f''(x_n)$$

3. Terminate when $|f'(x_n)| < \varepsilon$.



MULTIPLE DIMENSIONS

In \mathbb{R}^d we have to replace the derivatives by their vector space analogues.

Gradient descent

$$\mathbf{x}_{n+1} := \mathbf{x}_n - \nabla f(\mathbf{x}_n)$$

Newton's method for minima

$$\mathbf{x}_{n+1} := \mathbf{x}_n - H_f^{-1}(\mathbf{x}_n) \cdot \nabla f(\mathbf{x}_n)$$

The inverse of $H_f(\mathbf{x})$ exists only if the matrix is positive definite (not if it is only semidefinite), i.e. f has to be strictly convex.

The Hessian measures the curvature of f .

Effect of the Hessian

Multiplication by H_f^{-1} in general changes the direction of $\nabla f(\mathbf{x}_n)$. The correction takes into account how $\nabla f(\mathbf{x})$ changes away from \mathbf{x}_n , as estimated using the Hessian at \mathbf{x}_n .

Figure: Arrow is ∇f , $x + \Delta x_{nt}$ is Newton step.

