

An Introduction to Optimization with Application to Machine Learning

Hamed Valizadegan

University of Pittsburgh

Motivation: Machine Learning

▶ Linear Regression

$$\underset{w,b}{\text{minimize}} \quad \sum_{i=1}^n \|w^T x_i + b - y_i\|^2$$

▶ SVM

$$\begin{aligned} \underset{w,b}{\text{minimize}} \quad & \|w\|^2 + C \sum_{i=1}^n \varepsilon_i \\ \text{subject to} \quad & y_i(w^T x_i + b) \geq 1 - \varepsilon_i \quad i = 1, \dots, n \\ & \varepsilon_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

▶ PGDM metric learning

$$\begin{aligned} \underset{P}{\text{minimize}} \quad & \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_P \\ \text{subject to} \quad & \sum_{(x_i, x_j) \in D} \|x_i - x_j\|_P \geq 1 \\ & P \succcurlyeq 0 \end{aligned}$$

Optimization Problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0 \quad i = 1, \dots, m \\ & && h_i(x) = 0 \quad i = 1, \dots, p \end{aligned}$$

- ▶ $x \in R^n$ is the variable to find
- ▶ $f_0: R^n \rightarrow R$ is called the objective (cost or utility) function
- ▶ $f_i: R^n \rightarrow R, i = 1, \dots, m$ are the inequality constraints (defines a set)
- ▶ $h_i: R^n \rightarrow R, i = 1, \dots, p$ are the equality constraints (defines a set)

- ▶ Solution: $p^* = \inf\{f_0(x) | f_i(x) \leq 0 \quad i = 1, \dots, m, h_i(x) = 0 \quad i = 1, \dots, p\}$
- ▶ Constrained vs. unconstrained problems: whether you have the constraints or not.

- ▶ A feasible point x is optimal if $f_0(x) = p^*$; X_{OPT} is the set of optimal points.

Feasibility

- ▶ An optimization problem is feasible
 - ▶ if $x \in \text{dom } f_0$ (implicit constraints) and it satisfies all the (explicit) constraints $f_i(x) \leq 0 \ i = 1, \dots, m$ & $h_i(x) = 0 \ i = 1, \dots, p$.
- ▶ For infeasible problems, we say $p^* = +\infty$

- ▶ Feasibility problem

$$\begin{aligned} & \textit{find} && x \\ & \textit{subject to} && f_i(x) \leq 0 \ i = 1, \dots, m \\ & && h_i(x) = 0 \ i = 1, \dots, p \end{aligned}$$

- ▶ Equivalent to the following optimization problem

$$\begin{aligned} & \textit{minimize} && 0 \\ & \textit{subject to} && f_i(x) \leq 0 \ i = 1, \dots, m \\ & && h_i(x) = 0 \ i = 1, \dots, p \end{aligned}$$

Locally Optimal Points

- ▶ For the following problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && f_0(x) \\ & \text{s.t.} && f_i(x) \leq 0 \quad i = 1, \dots, m \\ & && h_i(x) = 0 \quad i = 1, \dots, p \end{aligned}$$

- ▶ x is locally optimum if there is an $R > 0$ such that x is optimal for the following problem

$$\begin{aligned} & \underset{z}{\text{minimize}} && f_0(z) \\ & \text{s.t.} && f_i(z) \leq 0 \quad i = 1, \dots, m \\ & && h_i(z) = 0 \quad i = 1, \dots, p \\ & && \|z - x\|_2 \leq R \end{aligned}$$

Regularization

- ▶ A form of limiting the feasible search space of an optimization problem
- ▶ Can be considered as the prior information that the solution is located in the neighborhood of point x

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f_0(x) \\ \text{s.t.} & f_i(x) \leq 0 \quad i = 1, \dots, m \\ & h_i(x) = 0 \quad i = 1, \dots, p \end{array} \quad \rightarrow \quad \begin{array}{ll} \underset{z}{\text{minimize}} & f_0(z) \\ \text{s.t.} & f_i(z) \leq 0 \quad i = 1, \dots, m \\ & h_i(z) = 0 \quad i = 1, \dots, p \\ & \|z - x\|_p \leq R \end{array}$$

- ▶ Leads to sparse solution for $x = 0$ and small p
- ▶ I will get back to this.

Convexity

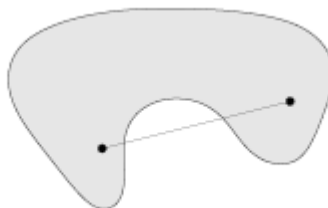
- ▶ An optimization problem is convex if
 - ▶ $f_0: R^n \rightarrow R$ is a convex function
 - ▶ Constrains $f_i(x) \leq 0 \ i = 1, \dots, m$ & $h_i(x) = 0 \ i = 1, \dots, p$ are convex sets.
 - ▶ $f_0: R^n \rightarrow R, f_i: R^n \rightarrow R, i = 1, \dots, m, h_i: R^n \rightarrow R, i = 1, \dots, p$ can be linear or nonlinear
- ▶ Importance
 - ▶ Any local optimum is a global optimum
 - ▶ Local optimality can be verified. No general tractable global optimum test
 - ▶ So, for convex problems, it is easy to check if a point is a global optimum.
- ▶ Feasible set of a convex optimization problem is convex.
- ▶ Convex set and convex function??

Affine and Convex Sets

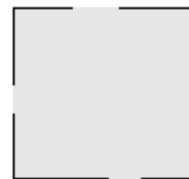
- ▶ Affine sets: the line through any two disjoint points
 - ▶ $x = \theta x_1 + (1 - \theta)x_2, \quad \theta \in \mathbb{R}$
 - ▶ Or equivalently, solution set of linear equation $\{x | Ax = b\}$
- ▶ Line segment: line segment between two points
 - ▶ $x = \theta x_1 + (1 - \theta)x_2, \quad 0 \leq \theta \leq 1$
- ▶ Convex Sets: a set that contains the line segment of any two points of the set
 - ▶ $x_1, x_2 \in S, 0 \leq \theta \leq 1 \Rightarrow \theta x_1 + (1 - \theta)x_2 \in S$



Convex



Non-convex



Non-Convex

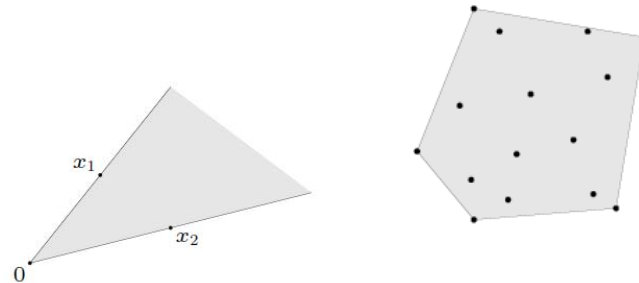
Convex Sets (examples)

- ▶ Convex hull of set $S = \{x_1, x_2, \dots, x_k\}$: Set of all convex combinations of points in S

- ▶ $\{x | \sum_{i=1}^k \theta_i x_i, \sum_{i=1}^k \theta_i = 1, \theta_i \geq 0\}$

- ▶ Conic combination of two points

- ▶ $x = \theta_1 x_1 + \theta_2 x_2, \quad 0 \leq \theta_1, \theta_2$



- ▶ Convex cone of set S: a set that contains all conic combinations of points in S

- ▶ Hyperplanes ($a^T x + b = 0$, linear equality)

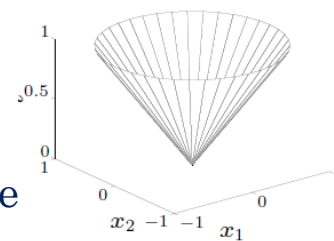
- ▶ Halfspaces ($a^T x + b \leq 0$, linear inequality)

- ▶ Euclidean balls and Ellipsoids: $\{x | (x - x_c)^T P^{-1}(x - x_c) \leq 1\}$ ($P \in \mathbf{S}^n_{++}$, i.e. P is positive-definite P)

- ▶ Norm ball: $\{x | \|x - x_c\| \leq r\}$

- ▶ Norm cone: $C = \{(x, t) | \|x\| \leq t\} \in \mathbb{R}^{n+1}$

- ▶ Euclidean norm cone ($\|\cdot\|_2$) is called second order cone

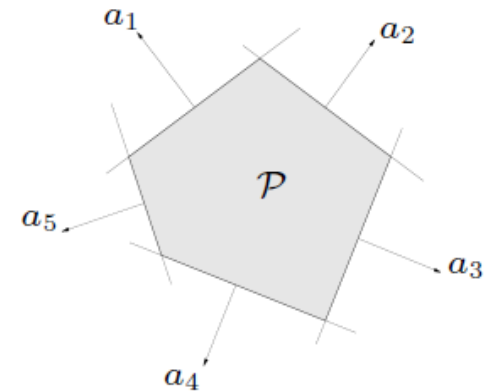


Operations that preserve convexity

- ▶ Intersection of convex sets
- ▶ The image of a convex set under affine (linear) function
 - ▶ $F: \mathbb{R}^n \rightarrow \mathbb{R}^m: F(\mathbf{x})=A\mathbf{x}+b$
 - ▶ scaling (aS), translation($S+a$), projection
- ▶ Perspective function
 - ▶ $F: \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n: F(\mathbf{x},t)=\mathbf{x}/t, \text{ dom}(F)=\{(\mathbf{x},t) \mid t>0\}$
 - ▶ Image and inverse image of convex sets under perspective are convex
- ▶ Linear-fractional functions:
 - ▶ $F: \mathbb{R}^n \rightarrow \mathbb{R}^m: F(\mathbf{x},t)=\frac{A\mathbf{x}+b}{c^T\mathbf{x}+d}, \text{ dom}(F)=\{\mathbf{x} \mid c^T\mathbf{x} + d>0\}$
 - ▶ Image and inverse image of convex sets under linear-fractional functions are convex

Convexity preserving operations (cont.)

- ▶ Intersection of convex sets is convex.
- ▶ Polyhedra is convex
 - ▶ Intersection of finite number of halfspaces and hyperplanes
- ▶ Positive semidefinite (PSD) cone: Set of all PSD matrices is convex
 - ▶ Intersection of infinite number of halfspaces and hyperspaces passing through origin
($\bigcap_{z \neq 0} \{X \in \mathcal{S}^n \mid z^T X z \geq 0\}$)
 - ▶ We denote it by \mathcal{S}^n_+



Generalized Inequalities

- ▶ Definition: A cone $K \subseteq \mathbb{R}^n$ is a proper cone if
 - ▶ K is convex
 - ▶ K is closed
 - ▶ K is solid: it has nonempty interior
 - ▶ K is pointed: it contains no line
- ▶ Generalized inequalities: defined by a proper cone K , is a partial ordering
 - $x \preceq_K y \iff y - x \in K$
 - $x \prec_K y \iff y - x \in \text{int } K$ (*interior of K*)
- ▶ Examples
 - ▶ Componentwise inequality:
 - $x \prec_{\mathbb{R}_+^n} y \iff y_i \geq x_i$
 - ▶ Matrix inequality
 - $X \prec_{\mathcal{S}_+^n} Y \iff Y - X \text{ is PSD}$

Dual Cones

- ▶ Dual cone of a cone K : $K^* = \{y \mid y^T x \geq 0 \text{ for all } x \in K\}$

$$x \preceq_K y \iff y - x \in K$$

$$x \prec_K y \iff y - x \in \text{int } K \text{ (interior of } K)$$

- ▶ Examples

- ▶ $K = \mathbf{R}_+^n$: $K^* = \mathbf{R}_+^n$

- ▶ $K = \mathbf{S}_+^n$: $K^* = \mathbf{S}_+^n$, ($\text{tr}(XY) \geq 0$)

- ▶ $K = \{(x, t) \mid \|x\|_2 \leq t\}$: $K^* = \{(x, t) \mid \|x\|_2 \leq t\}$

- ▶ $K = \{(x, t) \mid \|x\|_1 \leq t\}$: $K^* = \{(x, t) \mid \|x\|_\infty \leq t\}$

Convex Functions

- ▶ Definition: function $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if the graph of the function lies between the line segment joining any two points of the graph.



- ▶ Formally: $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if $\text{dom}(f)$ is convex and

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$

- ▶ Examples in \mathbb{R} :

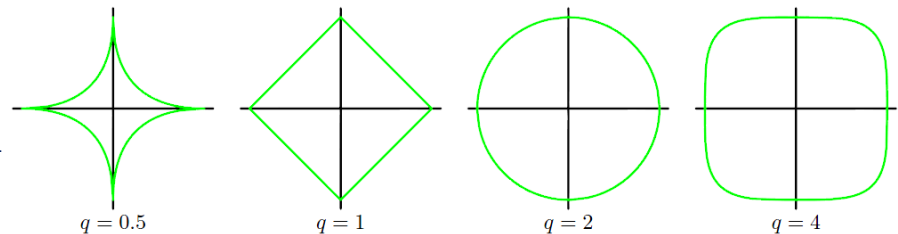
- ▶ affine, exponential, powers ($x^\alpha, \alpha \leq 0$ or $\alpha \geq 1$), power of absolute value ($|x|^\alpha, \alpha \geq 1$)

- ▶ Example on \mathbb{R}^n

- ▶ Norm $\|x\|_\alpha = (\sum_{i=1}^n |x_i|^\alpha)^{1/\alpha}, \alpha \geq 1$

- ▶ Example on $\mathbb{R}^{n \times m}$

- ▶ Affine function $\text{tr}(A^T X) + b = \sum_{i=1}^m \sum_{j=1}^n A_{ij} X_{ij} + b$



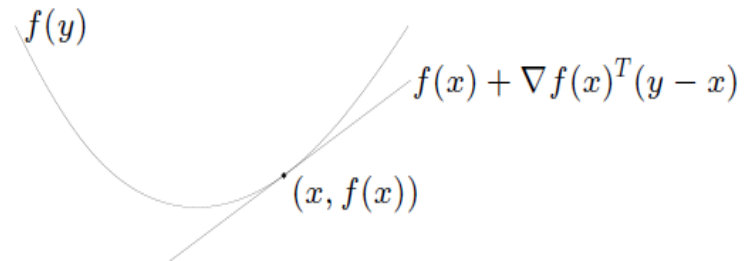
Convex Functions (verification tricks)

- ▶ $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if the following function of one variable is convex in t for any $x \in \text{dom}(f)$ & $v \in \mathbb{R}^n$:

$$g(t): \mathbb{R} \rightarrow \mathbb{R}: g(t) = f(x + tv), \text{dom}(g) = \{t \mid x + tv \in \text{dom}(f)\}$$

- ▶ First order condition: Differentiable f with convex domain is convex if and only if

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$



- ▶ Second order condition: twice differentiable function f with convex domain is convex if and only if

$$\nabla^2 f(x) \succeq 0 \text{ for all } x \in \text{dom}(f)$$

- ▶ Example: quadratic function $1/2x^T P x + q^T x + r$ is convex if P is PSD

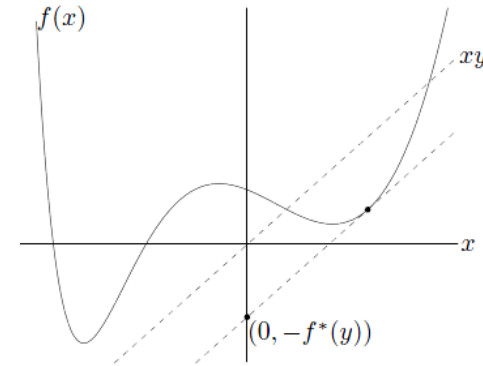
Operations that preserve convexity:

- ▶ Nonnegative weighted sum
 - ▶ $\sum_{i=1}^n \alpha_i f_i(x)$ is convex if $f_i(x), i = 1, 2, \dots, n$ are convex
 - ▶ Jensen's inequality: $f(\mathbb{E}(x)) \leq \mathbb{E}f(x)$
- ▶ Composition with affine function
 - ▶ $f(Ax + b)$ is convex if $f(x)$ is convex
 - ▶ Examples: $f(x) = -\sum_{i=1}^n \log(b_i - a_i^T x)$
- ▶ Minimization
 - ▶ $g(x) = \min_{y \in C} f(x, y)$ is convex if $f(x, y)$ is convex in (x, y) and C is a convex set
 - ▶ Examples: $\text{dist}(x, S) = \min_{y \in S} \|x - y\|$ is convex if S is convex
- ▶ Perspective $g(x, t) = tf\left(\frac{x}{t}\right), t > 0$
 - ▶ Example: $g(x, t) = \frac{x^T x}{t}, t > 0$
- ▶ Pointwise maximum and supremum
 - ▶ Piecewise linear function: $f(x) = \max_{i=1, \dots, n} a_i^T x + b_i$
 - ▶ $g(x) = \sup_{y \in A} f(x, y)$ is convex if $f(x, y)$ is convex in x for each $y \in A$
 - ▶ Example: max eigenvalue of a symmetric function $\lambda_{\max}(X) = \sup_{\|y\|=1} y^T X y$

Conjugate function:

- ▶ The conjugate function of f is defined as

$$f^*(y) = \sup_{x \in \text{dom}(f)} (y^T x - f(x))$$



- ▶ The conjugate function of f^* is the max cap between the linear function $y^T x$ and $f(x)$. For differentiable functions, this occurs at a point x where $y = \nabla f(x)$
- ▶ f^* is convex even if f is not. Because it is a pointwise supremum of a family of affine functions
- ▶ Also known as Legendre-Fenchel Transformation or Fenchel Transformation
- ▶ Examples
 - ▶ $f(x) = -\log(x) \rightarrow f^*(y) = -1 - \log(-y), y < 0$
 - ▶ $f(x) = \exp(x) \rightarrow f^*(y) = y \log(y) - y, y > 0$
 - ▶ $f(x) = x \log(x) \rightarrow f^*(y) = \exp(y-1), y \neq 0$
 - ▶ $f(x) = 1/x \rightarrow f^*(y) = -2(-y)^{1/2}, y \leq 0$

Slack variables

- ▶ Converting inequality constraints to equality constraints

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f_0(x) \\ \text{s.t.} & f_i(x) \leq 0 \quad i = 1, \dots, m \end{array} \quad \rightarrow \quad \begin{array}{ll} \underset{x, b_i}{\text{minimize}} & f_0(x) \\ \text{s.t.} & f_i(x) + b_i = 0 \quad i = 1, \dots, m \\ & b_i \geq 0 \quad i = 1, \dots, m \end{array}$$

- ▶ Introducing equality constraints

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f_0(A_0x + b_0) \\ \text{s.t.} & f_i(A_ix + b_i) \leq 0 \quad i = 1, \dots, m \end{array} \quad \rightarrow \quad \begin{array}{ll} \underset{x, b_i}{\text{minimize}} & f_0(y_0) \\ \text{s.t.} & f_i(y_i) \leq 0 \quad i = 1, \dots, m \\ & A_ix + b_i = y_i \quad i = 0, \dots, m \end{array}$$

- ▶ Converting an infeasible problem to feasible by relaxing the constraints

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f_0(x) \\ \text{s.t.} & f_i(x) \leq 0 \quad i = 1, \dots, m \end{array} \quad \rightarrow \quad \begin{array}{ll} \underset{x, b_i}{\text{minimize}} & f_0(x) + C \sum_{i=1}^m b_i \\ \text{s.t.} & f_i(x) - b_i \leq 0 \quad i = 1, \dots, m \\ & b_i \geq 0 \quad i = 1, \dots, m \end{array}$$

Duality

- ▶ The following optimization problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0 \quad i = 1, \dots, m \\ & && h_i(x) = 0 \quad i = 1, \dots, p \end{aligned}$$

- ▶ Can be written in the Lagrangian form

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

- ▶ $\lambda_i, i = 1, \dots, m$ are called the Lagrange multipliers associated with the inequalities and $\nu_i, i = 1, \dots, p$ are called the Lagrange multipliers associated with the equalities. They are also called the dual variables.

- ▶ The Lagrange dual function is defined as

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu) = \inf_x f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

- ▶ $g(\lambda, \nu)$ is the lower bound for the optimal value of original problem
 - ▶ $g(\lambda, \nu) \leq P^*$

The dual problem

- ▶ The following optimization problem is called the dual problem (original problem is called primal)

$$\begin{aligned} & \underset{\lambda, \nu}{\text{maximize}} \quad g(\lambda, \nu) \\ & \text{subject to} \quad \lambda \succeq 0 \end{aligned}$$

- ▶ Finds the best lower bound on p^*
 - ▶ A convex optimization problem with optimal value denoted by d^*
 - ▶ $L(\lambda, \nu)$ is concave since it is pointwise infimum of a family of affine functions

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu) = \inf_x f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

- ▶ This automatically gives a procedure to optimize the non-convex problems.

Solving dual problems

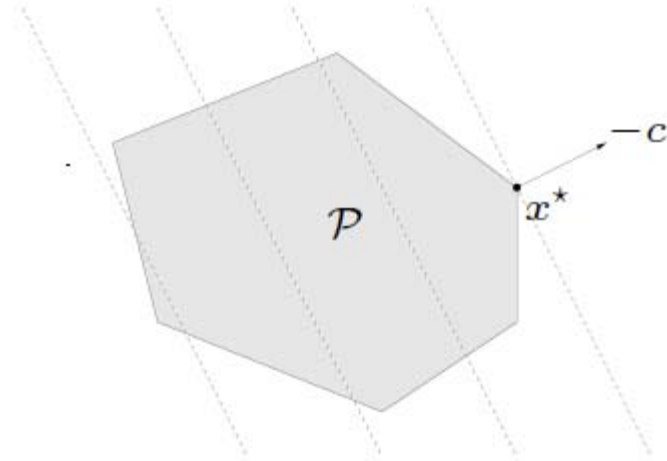
- ▶ Solve the dual problem which is convex
- ▶ Question: how good it is?
 - ▶ The duality gap $p^* - d^*$ is a measure of how good it is
 - ▶ Not usually easy to show that the gap is small
- ▶ Strong duality $p^* - d^*=0$
 - ▶ Usually (but not always) holds for convex problems
 - ▶ Non-convex problem can have strong duality as well so you can get lucky if you use the dual
- ▶ If the strong duality holds and x, λ, v are optimal, then they must satisfy the following conditions, called KKT conditions
 - ▶ Primal constraints: $f_i(x) \leq 0, i = 1, \dots, m$
 - ▶ Dual constraints: $\lambda_i > 0, i = 1, \dots, m$
 - ▶ Complementary slackness: $\lambda_i f_i(x) = 0, i = 1, \dots, m$
 - ▶ Gradient of Lagrangian vanishes: $\nabla f_0(x) + \sum_{i=1}^m \lambda_i \nabla f_i(x) + \sum_{i=1}^p v_i \nabla h_i(x)$

Linear Program (LP)

- ▶ Convex problem with affine objective and constraints functions

$$\begin{aligned} & \underset{x}{\text{minimize}} && c^T x + d \\ & \text{s.t.} && Gx \leq h \\ & && Ax = b \end{aligned}$$

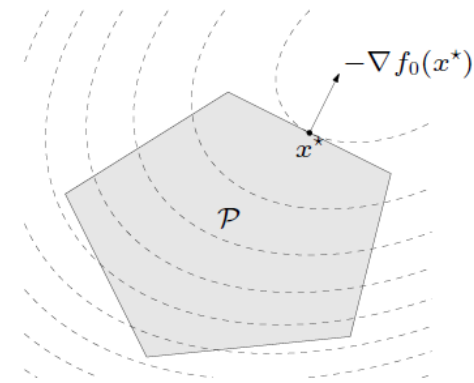
- ▶ Feasible set is a polyhedron
- ▶ linprog command in MATLAB



Quadratic Program (QP)

- ▶ Convex problem with quadratic convex objective and affine constraints functions (P is PSD)

$$\begin{aligned} \underset{x}{\text{minimize}} \quad & 1/2x^T Px + q^T x + r \\ \text{s.t.} \quad & Gx \leq h \\ & Ax = b \end{aligned}$$



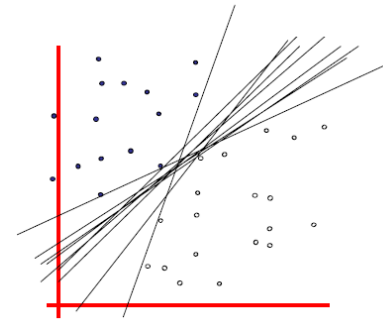
- ▶ Minimizes a convex quadratic over a polyhedron
- ▶ Quadprog command in matlab

SVM: a QP Example

- ▶ Many linear classifiers separating two separable set of examples

- ▶ Pick the one with maximum margin

$$\begin{aligned} & \underset{w,b}{\text{minimize}} \quad \|w\|^2 \\ & \text{subject to} \quad y_i(w^T x_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$



- ▶ If the examples are not separable, the feasible set of this problem is empty (infeasible problem)

- ▶ Utilizing slack variables to relax the constraints and make a feasible problem

$$\begin{aligned} & \underset{w,b}{\text{minimize}} \quad \|w\|^2 + C \sum_{i=1}^n \varepsilon_i \\ & \text{subject to} \quad y_i(w^T x_i + b) \geq 1 - \varepsilon_i \quad i = 1, \dots, n \\ & \quad \quad \quad \varepsilon_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

SVM: dual formulation

- ▶ Define the Lagrangian:

$$L(w, b, \lambda, \nu) = \|w\|^2 + C \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^m \alpha_i (y_i(w^T x_i + b) - 1 + \varepsilon_i) - \sum_{i=1}^n \mu_i \varepsilon_i$$

- ▶ Finding $L(\lambda, \nu) = \inf_{w, b} L(w, b, \lambda, \nu)$

$$\frac{\partial L(w, b, \lambda, \nu)}{\partial w} = 0 \rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L(w, b, \lambda, \nu)}{\partial b} = 0 \rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial L(w, b, \lambda, \nu)}{\partial \varepsilon_i} = 0 \rightarrow \alpha_i = C - \mu_i$$

- ▶ KKT conditions: 1) $\alpha_i \geq 0$, 2) $y_i(w^T x_i + b) - 1 + \varepsilon_i \geq 0$, 3) $\sum_{i=1}^m \alpha_i (y_i(w^T x_i + b) - 1 + \varepsilon_i) = 0$, 4) $\mu_i \geq 0$, 5) $\varepsilon_i \geq 0$, 6) $\mu_i \varepsilon_i = 0$

SVM: dual formulation

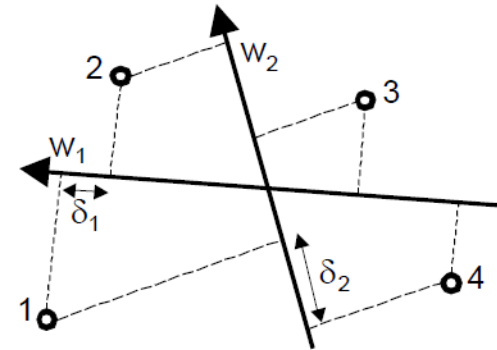
- ▶ Using these results, we obtain the dual problem

$$\begin{aligned} \underset{\lambda, \nu}{\text{maximize}} \quad & \sum_{i=1}^n \alpha_i - 1/2 \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i x_j \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \end{aligned}$$

- ▶ Useful form for using the kernel trick

$$\begin{aligned} \underset{\lambda, \nu}{\text{maximize}} \quad & \sum_{i=1}^n \alpha_i - 1/2 \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \end{aligned}$$

SVM^Rank: a QP Example



- ▶ Ranking problem:
 - ▶ n queries $q_i, i = 1, \dots, n$
 - ▶ for query q_i , a list of items $d_j^i, j = 1, \dots, m_i$ (feature vector) with their respect relevancy $r_j^i, j = 1, \dots, m_i$ to the query.
 - ▶ Assume also that r_j^i are discrete $[1..k]$
- ▶ Objective: obtain a linear classifier that respects ordering information
 - ▶ Suppose W is such a classifier
 - ▶ Construct a set on pair of examples $S = \{(x, z) \mid x = d_j^i, z = d_k^i, r_k^i - r_j^i = 1\}$
 - ▶ Find W that maximizes the margin between each two items

$$\underset{w, b}{\text{minimize}} \quad \|w\|^2 + C \sum_{r_i - r_j = 1} \varepsilon_{ij}$$

$$\text{subject to} \quad w^T(x_j - x_i) \geq 1 - \varepsilon_{ij}, \quad (x_i, x_j) \in S$$

$$\varepsilon_{ij} \geq 0 \quad i = 1, \dots, n$$

Multi-Task Learning

- ▶ Problem setup
 - ▶ T classification problems, each with different set of training examples.
 - ▶ Task t has n_t training examples $(x_i^t, y_i^t), i = 1, \dots, n_t$
 - ▶ Feature vector of all task are in the same space
 - ▶ Tasks are related (digits recognition, medical domains, etc)

- ▶ Objective: to learn linear classifiers $w^t, t = 1, \dots, T$ for tasks by considering that the tasks are similar

- ▶ Solution: assume all tasks are similar to a central unknown task μ

$$\underset{w, b, \mu}{\text{minimize}} \quad \sum_{t=1}^T \|w^t\|^2 + \sum_{t=1}^T \|w^t - \mu\|^2 + C \sum_{t=1}^T \sum_{i=1}^{n_t} \varepsilon_i^t$$

$$\text{subject to } y_i (w^{tT} x_i^t + b^t) \geq 1 - \varepsilon_i^t, i = 1, \dots, n, t = 1, \dots, T$$

$$\varepsilon_i^t \geq 0 \quad i = 1, \dots, n, t = 1, \dots, T$$

- ▶ How to write the dual of this problem? (Next lecture)

Quadratically Constrained QP (QCQP)

- ▶ Convex problem with quadratic convex objective and constraints functions (P_i are SDP)

$$\begin{aligned} & \underset{x}{\text{minimize}} && 1/2x^T P_0 x + q_0^T x + r_0 \\ & \text{s. t.} && 1/2x^T P_i x + q_i^T x + r_i \leq 0 \\ & && Ax = b \end{aligned}$$

- ▶ Objective and constraints are convex quadratic
- ▶ Can be solved with standard toolbox

Semidefinite Programming

- ▶ Convex problem with quadratic convex objective and constraints functions

$$\begin{aligned} & \underset{x}{\text{minimize}} && c^T x + d \\ & \text{s.t.} && x_1 P_1 + \cdots + x_n P_n + Q \preceq 0 \text{ (Linear Matrix Inequality)} \\ & && Gx \preceq b \text{ (General inequalities)} \\ & && Ax = b \end{aligned}$$

- ▶ Or

$$\begin{aligned} & \underset{x}{\text{minimize}} && \text{tr}(CX) \\ & \text{s.t.} && \text{tr}(A_i X) = b_i \\ & && X \succeq 0 \end{aligned}$$

- ▶ If P_1, \dots, P_n and Q are all diagonal, the SDP programming reduces to linear programming
- ▶ SeDuMi is a good tool to model this type of problems

Local and Global Consistency SSL

- ▶ Local and global Consistency, minimize

$$Q(F) = \underbrace{\frac{1}{2} \sum_{i,j=1}^N W_{ij} \left\| \frac{F_i}{\sqrt{D_{ii}}} - \frac{F_j}{\sqrt{D_{jj}}} \right\|^2}_{\text{Smoothness}} + \underbrace{\mu \sum_{i=1}^N \|F_i - Y_i\|^2}_{\text{Fitting}}$$

- ▶ Question: convex or non-convex?

$$Q(F) = F D^{-\frac{1}{2}} L D^{-\frac{1}{2}} F + \underbrace{\mu \sum_{i=1}^N \|F_i - Y_i\|^2}_{\text{Fitting}}$$

- ▶ How to solve such problems? (Next lecture)

PGDM metric learning

- ▶ PGDM metric learning

$$\begin{aligned} & \underset{P}{\text{minimize}} \quad \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_P \\ & \text{subject to} \quad \sum_{(x_i, x_j) \in D} \|x_i - x_j\|_P \geq 1 \\ & \quad \quad \quad P \succeq 0 \end{aligned}$$

- ▶ Question: convex or non-convex?
- ▶ How should we solve such problems? (next lecture)

LMNN metric learning

- ▶ LMNN metric learning

$$\begin{aligned} & \underset{A}{\text{minimize}} \quad \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_P \\ & \text{s.t.} \quad \|x_i - x_k\|_P - \|x_i - x_j\|_P \geq 1, (x_i, x_j, x_k) \in R \\ & \quad \quad P \succeq 0 \end{aligned}$$

- ▶ *in* $(x_i, x_j, x_k) \in R$, (x_i, x_j) are of the same class and neighbor according to Euclidean distance. (x_i, x_k) are from two different classes.
- ▶ Question: convex or non-convex?
- ▶ How should we solve such problems? (next lecture)