

Lecture 2: Linear SVM in the Dual

Stéphane Canu
stephane.canu@litislab.eu

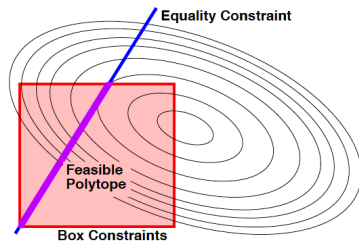
Sao Paulo 2014

March 12, 2014

Road map

1 Linear SVM

- Optimization in 10 slides
 - Equality constraints
 - Inequality constraints
- Dual formulation of the linear SVM
- Solving the dual



Linear SVM: the problem

Linear SVM are the solution of the following problem (called primal)

Let $\{(\mathbf{x}_i, y_i); i = 1 : n\}$ be a set of labelled data with $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{1, -1\}$.

A support vector machine (SVM) is a linear classifier associated with the following decision function: $D(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$ where $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ a given thought the solution of the following problem:

$$\begin{cases} \min_{\mathbf{w}, b} & \frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{with} & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, n \end{cases}$$

This is a quadratic program (QP):

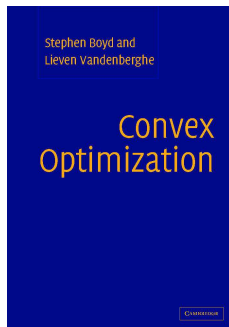
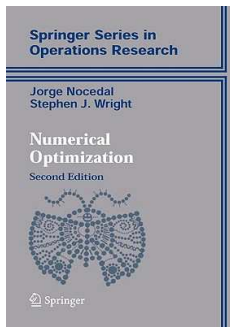
$$\begin{cases} \min_{\mathbf{z}} & \frac{1}{2} \mathbf{z}^\top \mathbf{A} \mathbf{z} - \mathbf{d}^\top \mathbf{z} \\ \text{with} & \mathbf{B} \mathbf{z} \leq \mathbf{e} \end{cases}$$

$$\mathbf{z} = (\mathbf{w}, b)^\top, \mathbf{d} = (0, \dots, 0)^\top, \mathbf{A} = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{B} = -[\text{diag}(\mathbf{y})\mathbf{X}, \mathbf{y}] \text{ et } \mathbf{e} = -(1, \dots, 1)^\top$$

Road map

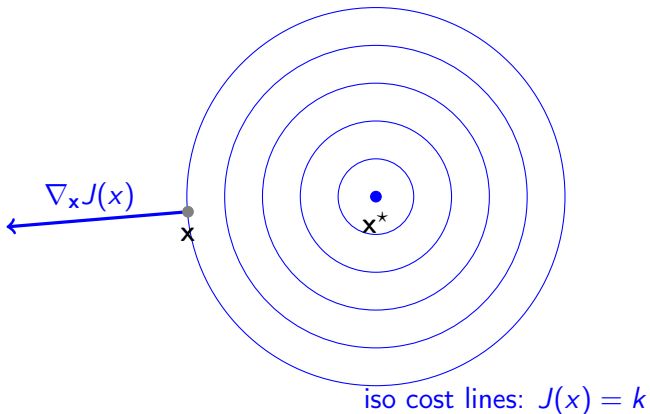
1 Linear SVM

- Optimization in 10 slides
 - Equality constraints
 - Inequality constraints
- Dual formulation of the linear SVM
- Solving the dual



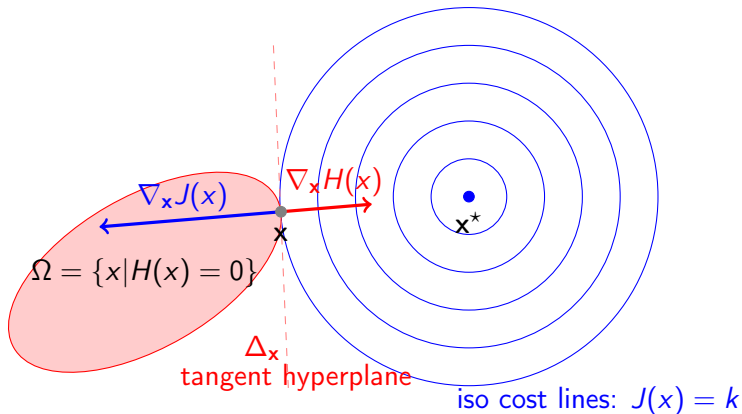
A simple example (to begin with)

$$\begin{cases} \min_{x_1, x_2} & J(\mathbf{x}) = (x_1 - a)^2 + (x_2 - b)^2 \\ \text{with} & \end{cases}$$



A simple example (to begin with)

$$\begin{cases} \min_{x_1, x_2} & J(x) = (x_1 - a)^2 + (x_2 - b)^2 \\ \text{with} & H(x) = \alpha(x_1 - c)^2 + \beta(x_2 - d)^2 + \gamma x_1 x_2 - 1 \end{cases}$$



$$\nabla_x H(x) = \lambda \nabla_x J(x)$$

The only one equality constraint case

$$\begin{cases} \min_{\mathbf{x}} & J(\mathbf{x}) & J(\mathbf{x} + \varepsilon \mathbf{d}) \approx J(\mathbf{x}) + \varepsilon \nabla_{\mathbf{x}} J(\mathbf{x})^{\top} \mathbf{d} \\ \text{with} & H(\mathbf{x}) = 0 & H(\mathbf{x} + \varepsilon \mathbf{d}) \approx H(\mathbf{x}) + \varepsilon \nabla_{\mathbf{x}} H(\mathbf{x})^{\top} \mathbf{d} \end{cases}$$

Loss J : \mathbf{d} is a descent direction if it exists $\varepsilon_0 \in \mathbb{R}$ such that $\forall \varepsilon \in \mathbb{R}, 0 < \varepsilon \leq \varepsilon_0$

$$J(\mathbf{x} + \varepsilon \mathbf{d}) < J(\mathbf{x}) \quad \Rightarrow \quad \nabla_{\mathbf{x}} J(\mathbf{x})^{\top} \mathbf{d} < 0$$

constraint H : \mathbf{d} is a feasible descent direction if it exists $\varepsilon_0 \in \mathbb{R}$ such that $\forall \varepsilon \in \mathbb{R}, 0 < \varepsilon \leq \varepsilon_0$

$$H(\mathbf{x} + \varepsilon \mathbf{d}) = 0 \quad \Rightarrow \quad \nabla_{\mathbf{x}} H(\mathbf{x})^{\top} \mathbf{d} = 0$$

If at \mathbf{x}^* , vectors $\nabla_{\mathbf{x}} J(\mathbf{x}^*)$ and $\nabla_{\mathbf{x}} H(\mathbf{x}^*)$ are collinear there is no feasible descent direction \mathbf{d} . Therefore, \mathbf{x}^* is a local solution of the problem.

Lagrange multipliers

Assume J and functions H_i are continuously differentials (and independent)

$$\mathcal{P} = \left\{ \begin{array}{l} \min_{\mathbf{x} \in \mathbb{R}^n} \quad J(\mathbf{x}) \\ \text{avec} \quad H_1(\mathbf{x}) = 0 \\ \text{et} \quad H_2(\mathbf{x}) = 0 \\ \quad \dots \\ \quad H_p(\mathbf{x}) = 0 \end{array} \right.$$

Lagrange multipliers

Assume J and functions H_i are continuously differentials (and independent)

$$\mathcal{P} = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & J(\mathbf{x}) \\ \text{avec} & H_1(\mathbf{x}) = 0 & \lambda_1 \\ \text{et} & H_2(\mathbf{x}) = 0 & \lambda_2 \\ & \dots \\ & H_p(\mathbf{x}) = 0 & \lambda_p \end{cases}$$

each constraint is associated with λ_i : the Lagrange multiplier.

Lagrange multipliers

Assume J and functions H_i are continuously differentials (and independent)

$$\mathcal{P} = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & J(\mathbf{x}) \\ \text{avec} & H_1(\mathbf{x}) = 0 & \lambda_1 \\ \text{et} & H_2(\mathbf{x}) = 0 & \lambda_2 \\ & \dots \\ & H_p(\mathbf{x}) = 0 & \lambda_p \end{cases}$$

each constraint is associated with λ_i : the Lagrange multiplier.

Theorem (First order optimality conditions)

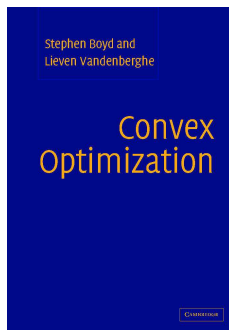
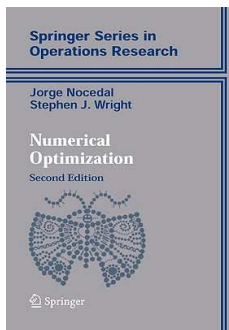
for \mathbf{x}^* being a local minima of \mathcal{P} , it is necessary that:

$$\nabla_{\mathbf{x}} J(\mathbf{x}^*) + \sum_{i=1}^p \lambda_i \nabla_{\mathbf{x}} H_i(\mathbf{x}^*) = 0 \quad \text{and} \quad H_i(\mathbf{x}^*) = 0, \quad i = 1, p$$

Plan

1 Linear SVM

- Optimization in 10 slides
 - Equality constraints
 - Inequality constraints
- Dual formulation of the linear SVM
- Solving the dual



The only one inequality constraint case

$$\begin{cases} \min_{\mathbf{x}} & J(\mathbf{x}) & J(\mathbf{x} + \varepsilon \mathbf{d}) \approx J(\mathbf{x}) + \varepsilon \nabla_{\mathbf{x}} J(\mathbf{x})^{\top} \mathbf{d} \\ \text{with} & G(\mathbf{x}) \leq 0 & G(\mathbf{x} + \varepsilon \mathbf{d}) \approx G(\mathbf{x}) + \varepsilon \nabla_{\mathbf{x}} G(\mathbf{x})^{\top} \mathbf{d} \end{cases}$$

cost J : \mathbf{d} is a descent direction if it exists $\varepsilon_0 \in \mathbb{R}$ such that
 $\forall \varepsilon \in \mathbb{R}, 0 < \varepsilon \leq \varepsilon_0$

$$J(\mathbf{x} + \varepsilon \mathbf{d}) < J(\mathbf{x}) \quad \Rightarrow \quad \nabla_{\mathbf{x}} J(\mathbf{x})^{\top} \mathbf{d} < 0$$

constraint G : \mathbf{d} is a feasible descent direction if it exists $\varepsilon_0 \in \mathbb{R}$ such that
 $\forall \varepsilon \in \mathbb{R}, 0 < \varepsilon \leq \varepsilon_0$

$$G(\mathbf{x} + \varepsilon \mathbf{d}) \leq 0 \quad \Rightarrow \quad \begin{array}{l} G(\mathbf{x}) < 0 : \text{no limit here on } \mathbf{d} \\ G(\mathbf{x}) = 0 : \nabla_{\mathbf{x}} G(\mathbf{x})^{\top} \mathbf{d} \leq 0 \end{array}$$

Two possibilities

If x^* lies at the limit of the feasible domain ($G(x^*) = 0$) and if vectors $\nabla_{\mathbf{x}} J(x^*)$ and $\nabla_{\mathbf{x}} G(x^*)$ are collinear **and in opposite directions**, there is no feasible descent direction \mathbf{d} at that point. Therefore, x^* is a local solution of the problem... Or if $\nabla_{\mathbf{x}} J(x^*) = 0$

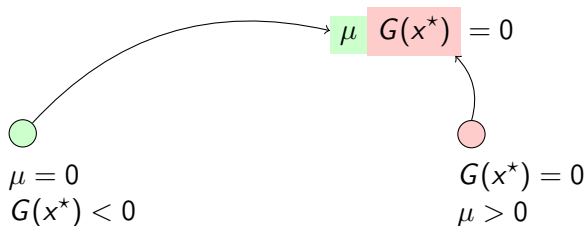
Two possibilities for optimality

$$\nabla_{\mathbf{x}} J(\mathbf{x}^*) = -\mu \nabla_{\mathbf{x}} G(\mathbf{x}^*) \quad \text{and} \quad \mu > 0; G(\mathbf{x}^*) = 0$$

or

$$\nabla_{\mathbf{x}} J(\mathbf{x}^*) = 0 \quad \text{and} \quad \mu = 0; G(\mathbf{x}^*) < 0$$

This alternative is summarized in the so called complementarity condition:



First order optimality condition (1)

$$\text{problem } \mathcal{P} = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & J(\mathbf{x}) \\ \text{with} & h_j(\mathbf{x}) = 0 \quad j = 1, \dots, p \\ \text{and} & g_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, q \end{cases}$$

Definition: Karush, Kuhn and Tucker (KKT) conditions

stationarity $\nabla J(\mathbf{x}^*) + \sum_{j=1}^p \lambda_j \nabla h_j(\mathbf{x}^*) + \sum_{i=1}^q \mu_i \nabla g_i(\mathbf{x}^*) = 0$

primal admissibility $h_j(\mathbf{x}^*) = 0 \quad j = 1, \dots, p$
 $g_i(\mathbf{x}^*) \leq 0 \quad i = 1, \dots, q$

dual admissibility $\mu_i \geq 0 \quad i = 1, \dots, q$

complementarity $\mu_i g_i(\mathbf{x}^*) = 0 \quad i = 1, \dots, q$

λ_j and μ_i are called the Lagrange multipliers of problem \mathcal{P}

First order optimality condition (2)

Theorem (12.1 Nocedal & Wright pp 321)

If a vector x^* is a stationary point of problem \mathcal{P}

Then there exists^a Lagrange multipliers such that $(x^*, \{\lambda_j\}_{j=1:p}, \{\mu_i\}_{i=1:q})$ fulfill KKT conditions

^a under some conditions e.g. linear independence constraint qualification

If the problem is **convex**, then a stationary point is the solution of the problem

A quadratic program (QP) is convex when...

$$(QP) \quad \begin{cases} \min_z & \frac{1}{2} \mathbf{z}^\top A \mathbf{z} - \mathbf{d}^\top \mathbf{z} \\ \text{with} & B \mathbf{z} \leq \mathbf{e} \end{cases}$$

... when matrix A is positive definite

KKT condition - Lagrangian (3)

$$\text{problem } \mathcal{P} = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & J(\mathbf{x}) \\ \text{with} & h_j(\mathbf{x}) = 0 \quad j = 1, \dots, p \\ \text{and} & g_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, q \end{cases}$$

Definition: Lagrangian

The lagrangian of problem \mathcal{P} is the following function:

$$\mathcal{L}(\mathbf{x}, \lambda, \mu) = J(\mathbf{x}) + \sum_{j=1}^p \lambda_j h_j(\mathbf{x}) + \sum_{i=1}^q \mu_i g_i(\mathbf{x})$$

The importance of being a lagrangian

- the stationarity condition can be written: $\nabla \mathcal{L}(\mathbf{x}^*, \lambda, \mu) = 0$
- the lagrangian saddle point $\max_{\lambda, \mu} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \mu)$

Primal variables: \mathbf{x} and **dual** variables λ, μ (the Lagrange multipliers)

Duality – definitions (1)

Primal and (Lagrange) dual problems

$$\mathcal{P} = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & J(\mathbf{x}) \\ \text{with} & h_j(\mathbf{x}) = 0 \quad j = 1, p \\ \text{and} & g_i(\mathbf{x}) \leq 0 \quad i = 1, q \end{cases} \quad \mathcal{D} = \begin{cases} \max_{\lambda \in \mathbb{R}^p, \mu \in \mathbb{R}^q} & Q(\lambda, \mu) \\ \text{with} & \mu_j \geq 0 \quad j = 1, q \end{cases}$$

Dual objective function:

$$\begin{aligned} Q(\lambda, \mu) &= \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \mu) \\ &= \inf_{\mathbf{x}} J(\mathbf{x}) + \sum_{j=1}^p \lambda_j h_j(\mathbf{x}) + \sum_{i=1}^q \mu_i g_i(\mathbf{x}) \end{aligned}$$

Wolf dual problem

$$\mathcal{W} = \begin{cases} \max_{\mathbf{x}, \lambda \in \mathbb{R}^p, \mu \in \mathbb{R}^q} & \mathcal{L}(\mathbf{x}, \lambda, \mu) \\ \text{with} & \mu_j \geq 0 \quad j = 1, q \\ \text{and} & \nabla J(\mathbf{x}^*) + \sum_{j=1}^p \lambda_j \nabla h_j(\mathbf{x}^*) + \sum_{i=1}^q \mu_i \nabla g_i(\mathbf{x}^*) = 0 \end{cases}$$

Duality – theorems (2)

Theorem (12.12, 12.13 and 12.14 Nocedal & Wright pp 346)

If f, g and h are convex and continuously differentiable^a, then the solution of the dual problem is the same as the solution of the primal

^aunder some conditions e.g. linear independence constraint qualification

$$\begin{aligned}(\lambda^*, \mu^*) &= \text{solution of problem } \mathcal{D} \\ \mathbf{x}^* &= \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda^*, \mu^*)\end{aligned}$$

$$\begin{aligned}Q(\lambda^*, \mu^*) &= \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda^*, \mu^*) = \mathcal{L}(\mathbf{x}^*, \lambda^*, \mu^*) \\ &= J(\mathbf{x}^*) + \lambda^* H(\mathbf{x}^*) + \mu^* G(\mathbf{x}^*) = J(\mathbf{x}^*)\end{aligned}$$

and for any feasible point \mathbf{x}

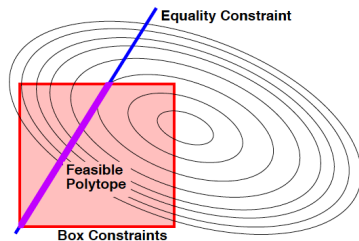
$$Q(\lambda, \mu) \leq J(\mathbf{x}) \quad \rightarrow \quad 0 \leq J(\mathbf{x}) - Q(\lambda, \mu)$$

The **duality gap** is the difference between the primal and dual cost functions

Road map

1 Linear SVM

- Optimization in 10 slides
 - Equality constraints
 - Inequality constraints
- Dual formulation of the linear SVM
- Solving the dual



Linear SVM dual formulation - The lagrangian

$$\begin{cases} \min_{\mathbf{w}, b} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{with} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, n \end{cases}$$

Looking for the lagrangian saddle point $\max_{\alpha} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha)$ with so called lagrange multipliers $\alpha_j \geq 0$

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

α_j represents the influence of constraint thus the influence of the training example (x_i, y_i)

Stationarity conditions

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

Computing the gradients:
$$\begin{cases} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) &= \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial \mathcal{L}(\mathbf{w}, b, \alpha)}{\partial b} &= \sum_{i=1}^n \alpha_i y_i \end{cases}$$

we have the following optimality conditions

$$\begin{cases} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) = 0 &\Rightarrow \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial \mathcal{L}(\mathbf{w}, b, \alpha)}{\partial b} = 0 &\Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

KKT conditions for SVM

$$\text{stationarity } \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \quad \text{and} \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\text{primal admissibility } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, \dots, n$$

$$\text{dual admissibility } \alpha_i \geq 0 \quad i = 1, \dots, n$$

$$\text{complementarity } \alpha_i (y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1) = 0 \quad i = 1, \dots, n$$

The complementary condition split the data into two sets

- \mathcal{A} be the set of active constraints: usefull points

$$\mathcal{A} = \{i \in [1, n] \mid y_i(\mathbf{w}^{*\top} \mathbf{x}_i + b^*) = 1\}$$

- its complementary $\bar{\mathcal{A}}$ useless points

$$\text{if } i \notin \mathcal{A}, \alpha_i = 0$$

The KKT conditions for SVM

The same KKT but using matrix notations and the active set \mathcal{A}

stationarity $\mathbf{w} - X^\top D_y \alpha = 0$

$$\alpha^\top \mathbf{y} = 0$$

primal admissibility $D_y(X\mathbf{w} + b\mathbb{1}) \geq \mathbb{1}$

dual admissibility $\alpha \geq 0$

complementarity $D_y(X_{\mathcal{A}}\mathbf{w} + b\mathbb{1}_{\mathcal{A}}) = \mathbb{1}_{\mathcal{A}}$

$$\alpha_{\bar{\mathcal{A}}} = 0$$

Knowing \mathcal{A} , the solution verifies the following linear system:

$$\begin{cases} \mathbf{w} & -X_{\mathcal{A}}^\top D_y \alpha_{\mathcal{A}} & & = 0 \\ -D_y X_{\mathcal{A}} \mathbf{w} & & -b\mathbf{y}_{\mathcal{A}} & = -\mathbf{e}_{\mathcal{A}} \\ & -\mathbf{y}_{\mathcal{A}}^\top \alpha_{\mathcal{A}} & & = 0 \end{cases}$$

with $D_y = \text{diag}(\mathbf{y}_{\mathcal{A}})$, $\alpha_{\mathcal{A}} = \alpha(\mathcal{A})$, $\mathbf{y}_{\mathcal{A}} = \mathbf{y}(\mathcal{A})$ et $X_{\mathcal{A}} = X(X_{\mathcal{A}}; :)$.

The KKT conditions as a linear system

$$\begin{cases} \mathbf{w} - X_{\mathcal{A}}^{\top} D_y \alpha_{\mathcal{A}} & = 0 \\ -D_y X_{\mathcal{A}} \mathbf{w} & - b \mathbf{y}_{\mathcal{A}} & = -\mathbf{e}_{\mathcal{A}} \\ & -\mathbf{y}_{\mathcal{A}}^{\top} \alpha_{\mathcal{A}} & = 0 \end{cases}$$

with $D_y = \text{diag}(\mathbf{y}_{\mathcal{A}})$, $\alpha_{\mathcal{A}} = \alpha(\mathcal{A})$, $\mathbf{y}_{\mathcal{A}} = \mathbf{y}(\mathcal{A})$ et $X_{\mathcal{A}} = X(X_{\mathcal{A}}; :)$.

I	$-X_{\mathcal{A}}^{\top} D_y$	0	\mathbf{w}	$=$	0
$-D_y X_{\mathcal{A}}$	0	$-\mathbf{y}_{\mathcal{A}}$	$\alpha_{\mathcal{A}}$		$-\mathbf{e}_{\mathcal{A}}$
0	$-\mathbf{y}_{\mathcal{A}}^{\top}$	0	b		0

we can work on it to separate \mathbf{w} from $(\alpha_{\mathcal{A}}, b)$

The SVM dual formulation

The SVM Wolfe dual

$$\left\{ \begin{array}{l} \max_{\mathbf{w}, b, \alpha} \quad \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1) \\ \text{with} \quad \alpha_i \geq 0 \\ \text{and} \quad \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \text{ and } \sum_{i=1}^n \alpha_i y_i = 0 \end{array} \right. \quad i = 1, \dots, n$$

using the fact: $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$

The SVM Wolfe dual without \mathbf{w} and b

$$\left\{ \begin{array}{l} \max_{\alpha} \quad -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_j \alpha_i y_i y_j \mathbf{x}_j^\top \mathbf{x}_i + \sum_{i=1}^n \alpha_i \\ \text{with} \quad \alpha_i \geq 0 \\ \text{and} \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{array} \right. \quad i = 1, \dots, n$$

Linear SVM dual formulation

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

Optimality: $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad \sum_{i=1}^n \alpha_i y_i = 0$

$$\begin{aligned} \mathcal{L}(\alpha) &= \frac{1}{2} \underbrace{\sum_{i=1}^n \sum_{j=1}^n \alpha_j \alpha_i y_i y_j \mathbf{x}_j^\top \mathbf{x}_i}_{\mathbf{w}^\top \mathbf{w}} - \sum_{i=1}^n \alpha_i y_i \underbrace{\sum_{j=1}^n \alpha_j y_j \mathbf{x}_j^\top \mathbf{x}_i}_{\mathbf{w}^\top} - b \underbrace{\sum_{i=1}^n \alpha_i y_i}_{=0} + \sum_{i=1}^n \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_j \alpha_i y_i y_j \mathbf{x}_j^\top \mathbf{x}_i + \sum_{i=1}^n \alpha_i \end{aligned}$$

Dual linear SVM is also a quadratic program

$$\text{problem } \mathcal{D} \quad \begin{cases} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} & \mathbf{y}^\top \alpha = 0 \\ \text{and} & 0 \leq \alpha_i \quad i = 1, n \end{cases}$$

with G a symmetric matrix $n \times n$ such that $G_{ij} = y_i y_j \mathbf{x}_j^\top \mathbf{x}_i$

SVM primal vs. dual

Primal

$$\left\{ \begin{array}{ll} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{with} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{array} \right.$$

- $d + 1$ unknown
- n constraints
- classical QP
- perfect when $d \ll n$

Dual

$$\left\{ \begin{array}{ll} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} & \mathbf{y}^\top \alpha = 0 \\ \text{and} & 0 \leq \alpha_j \quad i = 1, n \end{array} \right.$$

- n unknown
- G Gram matrix (pairwise influence matrix)
- n box constraints
- easy to solve
- to be used when $d > n$

SVM primal vs. dual

Primal

$$\left\{ \begin{array}{ll} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{with} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{array} \right.$$

- $d + 1$ unknown
- n constraints
- classical QP
- perfect when $d \ll n$

Dual

$$\left\{ \begin{array}{ll} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} & \mathbf{y}^\top \alpha = 0 \\ \text{and} & 0 \leq \alpha_i \end{array} \right. \quad i = 1, n$$

- n unknown
- G Gram matrix (pairwise influence matrix)
- n box constraints
- easy to solve
- to be used when $d > n$

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

The bi dual (the dual of the dual)

$$\begin{cases} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} & \mathbf{y}^\top \alpha = 0 \\ \text{and} & 0 \leq \alpha_j \quad i = 1, n \end{cases}$$

$$\begin{aligned} \mathcal{L}(\alpha, \lambda, \mu) &= \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha + \lambda \mathbf{y}^\top \alpha - \mu^\top \alpha \\ \nabla_{\alpha} \mathcal{L}(\alpha, \lambda, \mu) &= G \alpha - \mathbf{e} + \lambda \mathbf{y} - \mu \end{aligned}$$

The bidual

$$\begin{cases} \max_{\alpha, \lambda, \mu} & -\frac{1}{2} \alpha^\top G \alpha \\ \text{with} & G \alpha - \mathbf{e} + \lambda \mathbf{y} - \mu = 0 \\ \text{and} & 0 \leq \mu \end{cases}$$

since $\|\mathbf{w}\|^2 = \frac{1}{2} \alpha^\top G \alpha$ and $D\mathbf{X}\mathbf{w} = G\alpha$

$$\begin{cases} \max_{\mathbf{w}, \lambda} & -\frac{1}{2} \|\mathbf{w}\|^2 \\ \text{with} & D\mathbf{X}\mathbf{w} + \lambda \mathbf{y} \geq \mathbf{e} \end{cases}$$

by identification (possibly up to a sign)

$b = \lambda$ is the Lagrange multiplier of the equality constraint

Cold case: the least square problem

Linear model

$$y_i = \sum_{j=1}^d \mathbf{w}_j x_{ij} + \varepsilon_i \quad , \quad i = 1, n$$

n data and d variables; $d < n$

$$\min_{\mathbf{w}} = \sum_{i=1}^n \left(\sum_{j=1}^d x_{ij} \mathbf{w}_j - y_i \right)^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

Solution: $\tilde{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

$$f(\mathbf{x}) = \mathbf{x}^\top \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}_{\tilde{\mathbf{w}}}$$

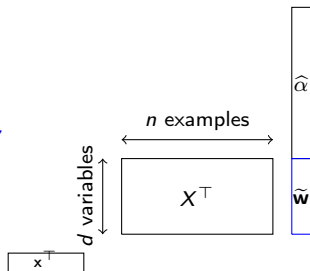
What is the influence of **each data point** (matrix \mathbf{X} lines) ?

data point influence (contribution)

for any new data point \mathbf{x}

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{x}^\top (X^\top X)(X^\top X)^{-1} \underbrace{(X^\top X)^{-1} X^\top \mathbf{y}}_{\tilde{\mathbf{w}}} \\ &= \mathbf{x}^\top X^\top \underbrace{X(X^\top X)^{-1} (X^\top X)^{-1} X^\top \mathbf{y}}_{\hat{\alpha}} \end{aligned}$$

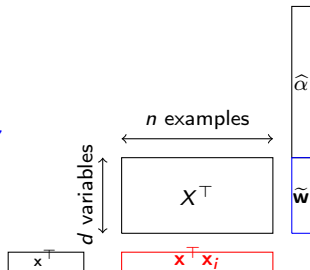
$$f(\mathbf{x}) = \sum_{j=1}^d \tilde{\mathbf{w}}_j x_j$$



data point influence (contribution)

for any new data point \mathbf{x}

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{x}^T (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}_{\tilde{\mathbf{w}}} \\ &= \mathbf{x}^T \mathbf{X}^T \underbrace{\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}_{\hat{\boldsymbol{\alpha}}} \end{aligned}$$



$$f(\mathbf{x}) = \sum_{j=1}^d \tilde{\mathbf{w}}_j x_j = \sum_{i=1}^n \hat{\alpha}_i (\mathbf{x}^T \mathbf{x}_i)$$

from variables to examples

$$\underbrace{\hat{\boldsymbol{\alpha}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{w}}}_{n \text{ examples}}$$

et

$$\underbrace{\tilde{\mathbf{w}} = \mathbf{X}^T \hat{\boldsymbol{\alpha}}}_{d \text{ variables}}$$

what if $d \geq n$!

SVM primal vs. dual

Primal

$$\begin{cases} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{with} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{cases}$$

- $d + 1$ unknown
- n constraints
- classical QP
- perfect when $d \ll n$

Dual

$$\begin{cases} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} & \mathbf{y}^\top \alpha = 0 \\ \text{and} & 0 \leq \alpha_i \quad i = 1, n \end{cases}$$

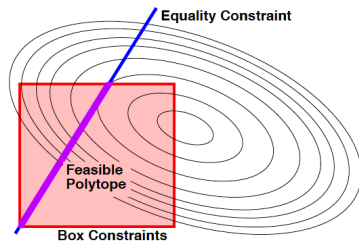
- n unknown
- G Gram matrix (pairwise influence matrix)
- n box constraints
- easy to solve
- to be used when $d > n$

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

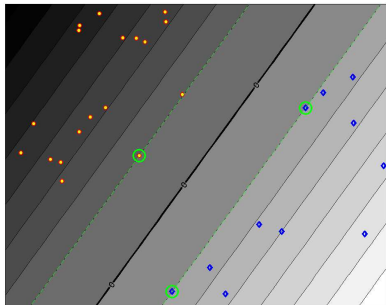
Road map

1 Linear SVM

- Optimization in 10 slides
 - Equality constraints
 - Inequality constraints
- Dual formulation of the linear SVM
- Solving the dual



Solving the dual (1)

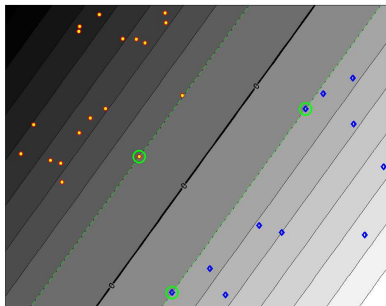


Data point influence

- $\alpha_j = 0$ this point is useless
- $\alpha_j \neq 0$ this point is said to be **support**

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

Solving the dual (1)



Data point influence

- $\alpha_j = 0$ this point is useless
- $\alpha_j \neq 0$ this point is said to be **support**

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^3 \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

Decision border only depends on 3 points ($d + 1$)

Solving the dual (2)

Assume we know these 3 data points

$$\left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} \quad \mathbf{y}^\top \alpha = 0 \\ \text{and} \quad 0 \leq \alpha_i \end{array} \right. \quad i = 1, n \quad \Rightarrow \quad \left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^3} \quad \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{with} \quad \mathbf{y}^\top \alpha = 0 \end{array} \right.$$

$$L(\alpha, \mathbf{b}) = \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha + \mathbf{b} \mathbf{y}^\top \alpha$$

solve the following linear system

$$\left\{ \begin{array}{l} G \alpha + \mathbf{b} \mathbf{y} = \mathbf{e} \\ \mathbf{y}^\top \alpha = 0 \end{array} \right.$$

```
U = chol(G); % upper
a = U \ (U' \ e);
c = U \ (U' \ y);
b = (y' * a) \ (y' * c)
alpha = U \ (U' \ (e - b * y));
```

Conclusion: variables or data point?

- seeking for a universal learning algorithm
 - ▶ no model for $\mathbb{P}(\mathbf{x}, y)$
- the linear case: data is separable
 - ▶ the non separable case
- double objective: minimizing the error together with the regularity of the solution
 - ▶ multi objective optimisation
- duality : variable – example
 - ▶ use the primal when $d < n$ (in the linear case) or when matrix G is hard to compute
 - ▶ otherwise use the dual
- universality = nonlinearity
 - ▶ kernels