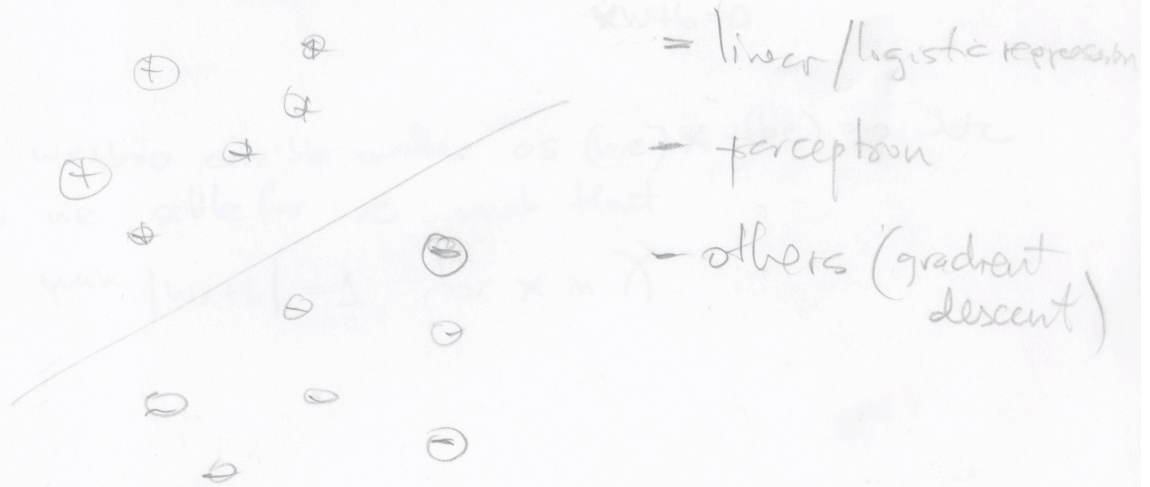# Support Vector Machines

Virgil Pavlu      October 8, 2008
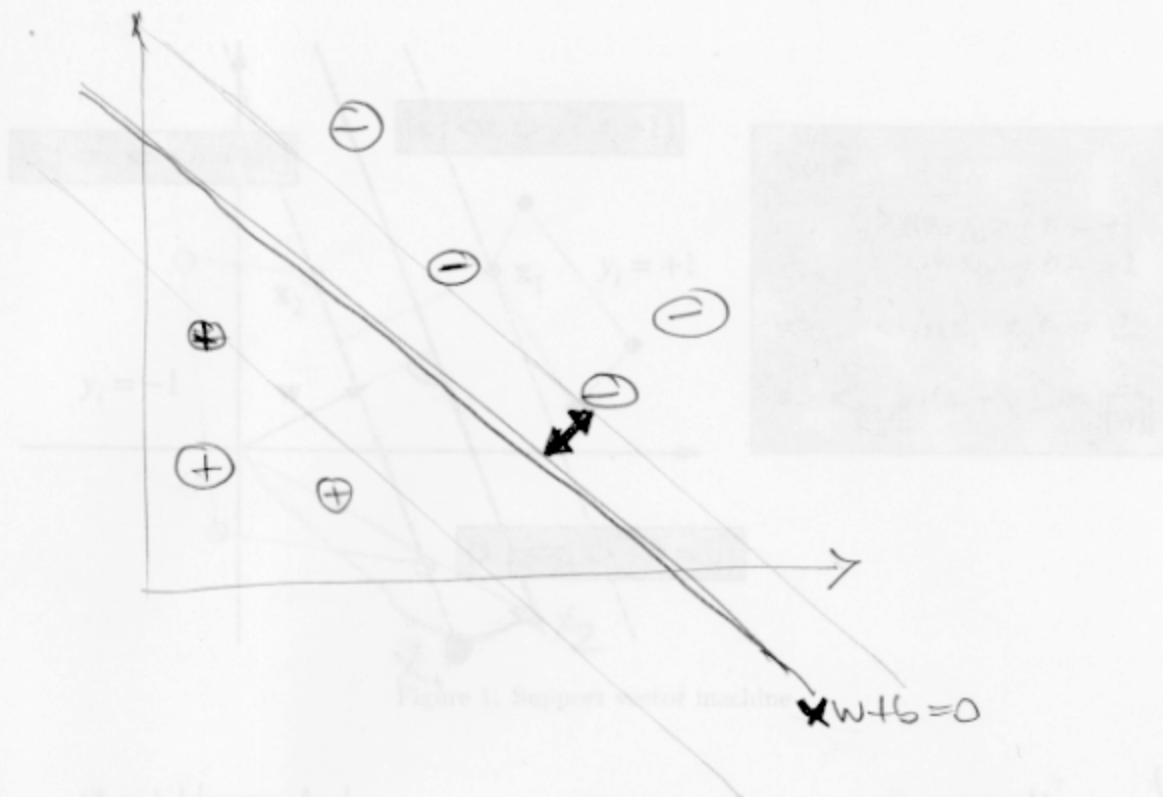
## 1   Linear discrimination

-we have talked about linear/logistic regression and the perceptron. There are other methods, essentially variations of gradient descent with specific objective functions.

    -lets take a closer look at a hyperplane separating the classes in a binary problem.

## 2 Geometry of the hyperplanes (separable case)



line $wx + b = 0$ can be written as $(wc)x + (bc) = 0$ $\forall c$

so we settle for $c$ such that

$$\min |wx + b| = 1 \quad \text{for } x \text{ in } X.$$

# 3   Margin
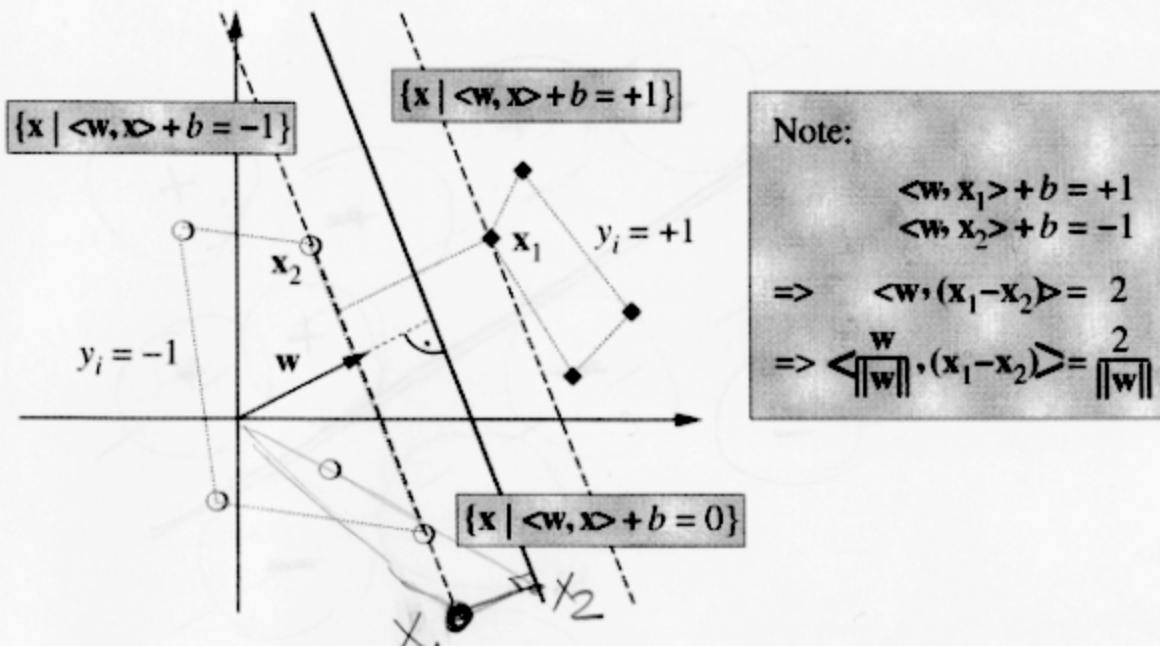


Figure 1: Support vector machine

$$x_1 w + b = 1 \Big| \rightarrow (x_2 - x_1) w = 1 \rightarrow \| x_2 - x_1 \| = \frac{1}{\|w\|}$$
$$x_2 w + b > 0$$
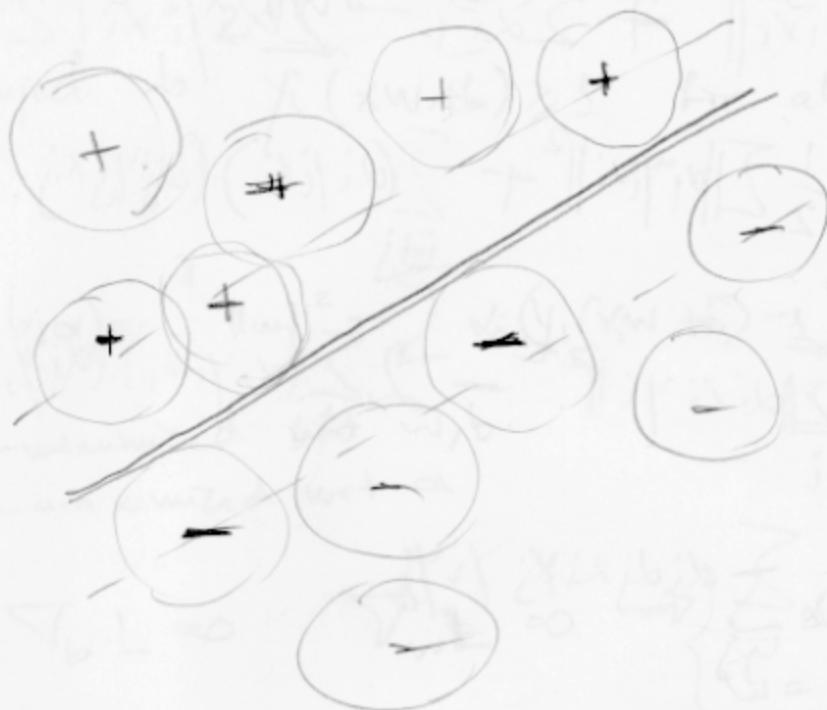
margin:

$$\rho(x, y) = \frac{y(x w + b)}{\|w\|} = \frac{y(w x_1 + b - (w x_2 + b))}{\|w\|} =$$

$$= \frac{y((x_1 - x_2) w)}{\|w\|} = y(x_1 - x_2) \cdot$$

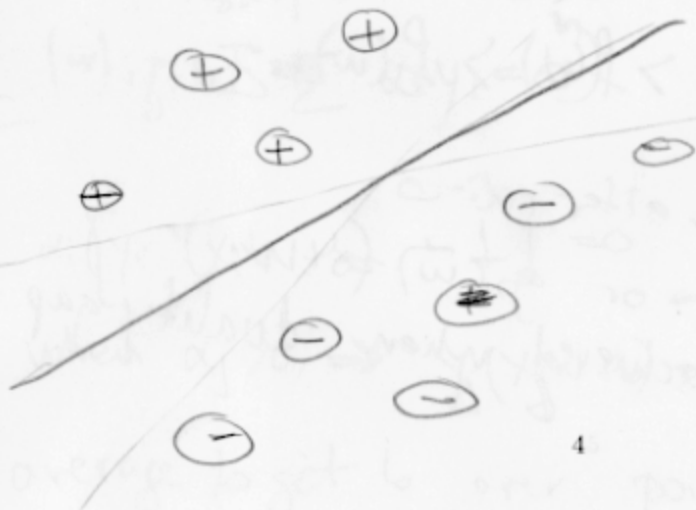$$\text{"size"} = \frac{1}{\|w\|} = \| x_2 - x_1 \| \cdot |y|$$

- perceptron with data+noise $\Rightarrow$ letter hyperplane

freedom do modify (slightly) the hyperplane
and still get a good generalization

# 5 SVM - optimal hyperplane wrt margin

primal opt $\Big|$ minimize $\frac{1}{2}\|w\|^2$ $\big[\Longleftrightarrow$ maximize margin$\big]$

$\Big|$ subject to $y_i(xw+b) \geq 1$ for all $i$

if $y_i(x_iw+b) < 1 \Rightarrow L \to \infty$
so $w,b$ forced to do
$y_i(x_iw+b) \geq 1$

$$L(w,b,\alpha) = \frac{1}{2}\|w\|^2 - \sum \alpha_i\big(y_i(xw+b)-1\big)$$

— minimized wrt $w,b$
— maximized wrt $\alpha$

$\nabla_b L = 0$ $\qquad \nabla_w L = 0$ $\Rightarrow \Big\{ \begin{array}{l} \sum \alpha_i y_i = 0 \\ w = \sum \alpha_i y_i x_i \end{array}$ $\to$ Lin. combs

remember $\alpha_i > 0 \Longleftrightarrow y_i(xw+b) = 1 \Longleftrightarrow (x_i, y_i)$ SUPPORT VECTOR

dual: maximize $$W(\alpha) = \sum \alpha_i - \frac{1}{2}\sum \alpha_i \alpha_j y_i y_j x_i x_j^T$$

subject to $\alpha_i \geq 0$
$\sum \alpha_i y_i = 0$

EASIER PB ( ~~no~~ W )
— quadratic solver
— heuristic

$\alpha_i\big[y_i(x_iw+b) - 1\big] = 0$

when $\alpha_j > 0 \Rightarrow y_j(x_jw+b) = 1 \Rightarrow \sum_i \alpha_i y_i x_i x_j^T + b = y_j$

average to get $b$ over points with $\alpha_j > 0$

$$\frac{1}{2}\|w\|^2 = \frac{1}{2}\left\|\sum \alpha_i y_i x_i\right\|^2$$

$$L = \frac{1}{2}\left\|\sum \alpha_i y_i x_i\right\|^2 + \sum \alpha_i - \sum \alpha_i y_i x_i \left(\sum \alpha_j y_j x_j\right)$$

$$= \sum \alpha_i + \frac{1}{2}\sum\|\alpha_i y_i x_i\|^2 + \sum_{i \neq j}(\alpha_i y_i x_i)(\alpha_j y_j x_j)$$

$$= \sum_i \|\alpha_i x_i y_i\|^2 - 2\sum_{i \neq j}(\alpha_i y_i x_i)(\alpha_j y_j x_j)$$

$$= \sum \alpha_i - \frac{1}{2}\sum_{ij}\alpha_i \alpha_j x_i x_j y_i y_j$$

---

$$\sum \alpha_i g_i(w)$$

$g_i = $ constraints       Karush Kuhn Tucker

$\|$

duality gap (or KKT gap)

$\tilde{w}$ solution $\Rightarrow$ any $w, \alpha$ with $\left.\begin{cases} \alpha > 0 \\ \partial_w L = 0 \\ \partial_\alpha L = 0 \end{cases}\right\}$ has

$$f(w) > f(\tilde{w}) > f(w) + \sum \alpha_i g_i(w)$$

best $\tilde{w}, \tilde{\alpha}:$

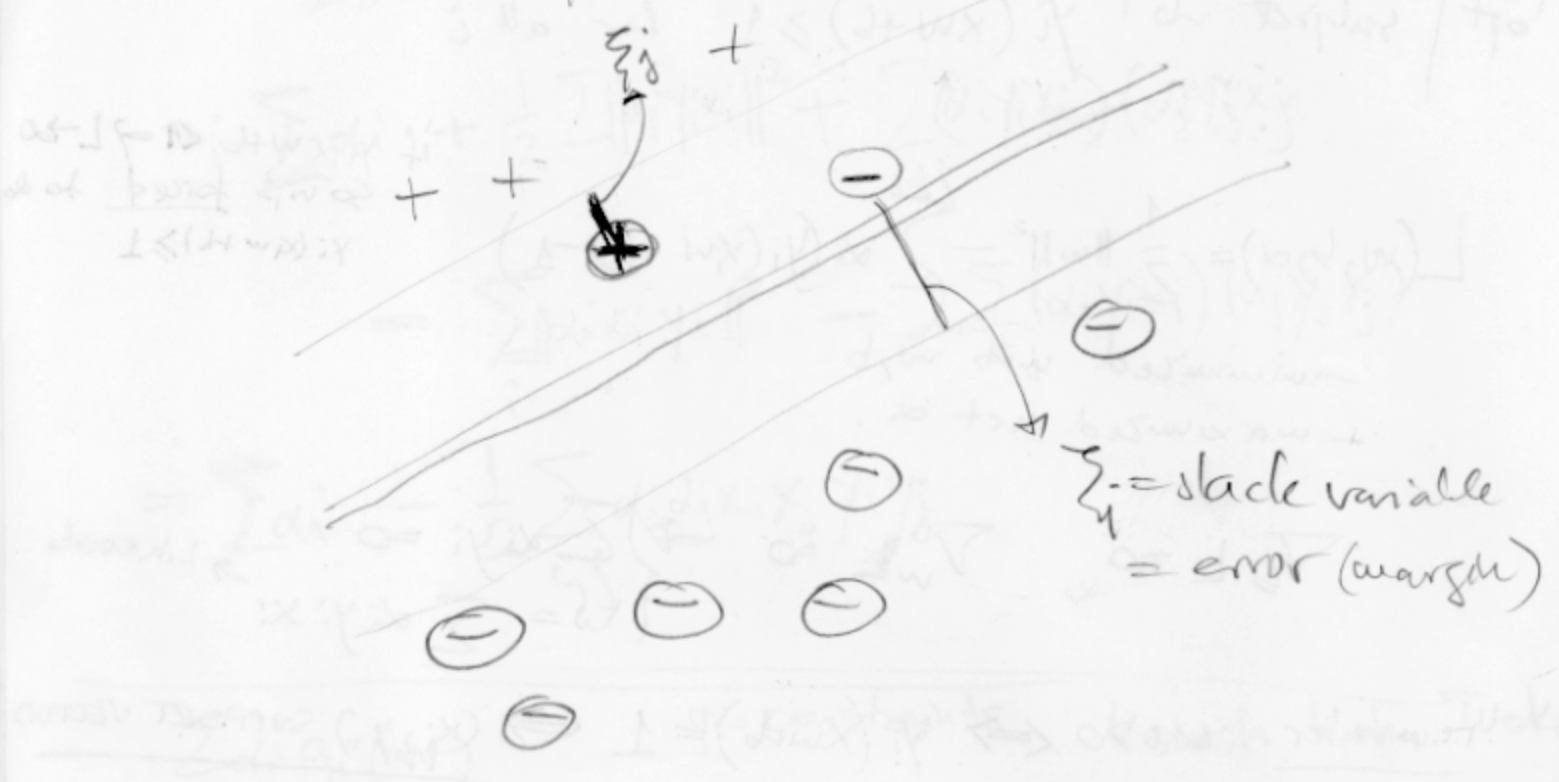$$\sum \tilde{\alpha}_i g_i(\tilde{w}) = 0 \begin{cases} \text{either } \alpha_i = 0 \\ \text{or } g_i(\tilde{w}) = 0 \end{cases}$$

so $f(\tilde{w})$ is ~~exac~~ achieved when duality gap closes.

$\xi_j$

$\xi_i =$ slack variable
$=$ error (margin)

# 6   Non-separable data

— soft margin hyperplanes

constraints $y_i(xw+b) \geq 1-\xi_i$

$\xi_i = $ slack variable

minimize $\quad \frac{1}{2}\|w\|^2 + \frac{C}{m}\sum \xi_i$

subject to $\xi_i \geq 0$

$y_i(xw+b) \geq 1-\xi_i$

dual $\quad$ maximize $W(\alpha) = \sum \alpha_i - \frac{1}{2}\sum \alpha_i \alpha_j y_i y_j x_i x_j^T$

subj to $0 \leq \alpha_i \leq \frac{C}{m}$

$\sum \alpha_i y_i = 0$

to comput $b$,

$$L(\quad) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^{m} \alpha_i (y_i (xw+b) - 1 + \zeta) + \frac{C}{m}\zeta$$

$$\boxed{\begin{array}{c} why \\ \alpha < \frac{C}{m}? \end{array}} \quad -\alpha(y \cdot xw + y_0 = 1 + \zeta) + \frac{C}{m}\zeta$$

$\Rightarrow$ coefficient of $\zeta$ is $-\alpha + \frac{C}{m}$

$\zeta > 0$; for optimization (minimization) to make sense: $-\alpha + \frac{C}{m} > 0$, otherwise we can send $\zeta \to +\infty$ and minimize to $-\infty$.

---

In fact $L(\quad) = \frac{1}{2}\|w\|^2 - \sum \alpha_i (y_i(xw+b) - 1 + \zeta)$
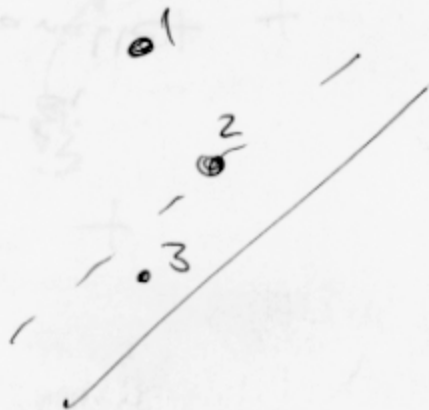
$\qquad + \frac{C}{m}\sum \zeta_i - \boxed{\beta_i}\zeta_i$

$\qquad\qquad\qquad \longrightarrow$ Lag. multiplier for constraint $\zeta_i \geq 0$

$$\frac{\partial L}{\partial \zeta} = -\alpha - \beta + \frac{C}{m} = 0 \Rightarrow \frac{C}{m} = \alpha_i + \beta_i$$

KKT: $\beta_i \zeta_i = 0$ so

$\left\{\begin{array}{l} 1.\ \zeta_i = 0,\ \alpha_i = 0,\ \beta_i = 0? \\ 2.\ \zeta_i = 0 \ 0 < \alpha_i < \frac{C}{m}\ \beta_i > 0? \\ 3.\ \zeta_i > 0\ \alpha_i = \frac{C}{m},\ \beta_i = 0 \end{array}\right.$

⟹ SVM recap; margin, support vector. margin $\simeq \frac{1}{\|w\|}$

⟹ Primal: minimize $\frac{1}{2}\|w\|^2$
subject to $y_i(x_iw+b) \geq 1$

~ Lagrangian $\begin{cases} L(w,b,\alpha) = \frac{1}{2}\|w\|^2 - \sum \alpha_i(y_i(x_iw+b)-1) \\ \text{minimize with respect to } w, b \\ \text{max} \quad \text{wrt } \alpha. \end{cases}$

— KKT theorem: nec + suf condition for solution
$\partial_{w,b}L = 0$; $\alpha_i(y_i(x_iw+b)-1) = 0 \begin{cases} \nearrow \text{sp. vector} \\ \searrow \text{irelev. constrain} \end{cases}$

$\alpha_i \geq 0$; $y_i(x_iw+b) \geq 1$.
$\tilde{w},\tilde{\alpha}$ solution $\Rightarrow L(w,\tilde{\alpha}) \geq L(\tilde{w},\tilde{\alpha}) \geq L(\tilde{w},\alpha)$.

— saddle points, dual problem.
max: $W(\alpha) = \sum \alpha_i - \frac{1}{2}\sum \alpha_i\alpha_j y_i y_j x_i x_j^T$
subject to $\alpha_i \geq 0$, $\sum \alpha_i y_i = 0$.

⟹ Non separable data: Soft margin hyperplane
$\xi_i$ = slack variables

Primal $\begin{cases} \text{minimize} & \frac{1}{2}\|w\|^2 + \frac{C}{m}\sum \xi_i \\ \text{subject to} & \xi_i \geq 0, \; y_i(x_iw+b) \geq 1-\xi_i \end{cases}$

Dual $\begin{cases} \text{max } W(\alpha) = \sum \alpha_i - \frac{1}{2}\sum \alpha_i\alpha_j y_i y_j x_i x_j^T \\ \text{subject to } 0 \leq \alpha_i \leq \frac{C}{m} \end{cases}$

**Quadratic Solvers.**    SMO = sequential minimal optimizat.

coordinate ascent: Loop

maxim one coordinate $\left[ \alpha_i := \arg\max \; W\left( \alpha_1, \alpha_2 \ldots \alpha_{i-1}, \hat{\alpha_i}, \right. \right.$

at the time.          $\left. \left. \alpha_{i+1}, \ldots \alpha_m \right) \right.$

order $1, 2, \ldots m$  <u>MATTERS</u>

SMO: Loop

choose a pair $\alpha_i, \alpha_j$ (by heuristic $\rightarrow$ max update)

$\alpha_i, \alpha_j = \arg\max\limits_{i,j} \; W(\alpha)$

efficient argmax!

$\alpha_1 y_1 + \alpha_2 y_2 = T$ constant ( because of constraint $\sum \alpha_i y_i = 0$ ).

$\alpha_1 = (T - \alpha_2 y_2) \, y_1$

$W = W\left( (T - \hat{\alpha_2} y_2) y_1, \hat{\alpha_2} \ldots \ldots \right)$ quadratic in $\alpha_2$

---

→ interior point methods:
- solve KKT equations, iteratively
- need a trick on $\alpha \zeta = 0 \Rightarrow \alpha \zeta = \mu$, make $\mu \to 0$
- follow path primal – dual
- cholesky decomposition , LAPACK (lin. alg library)

≡ chunking working sets ; try to identify a likely support vector set, only work with this data