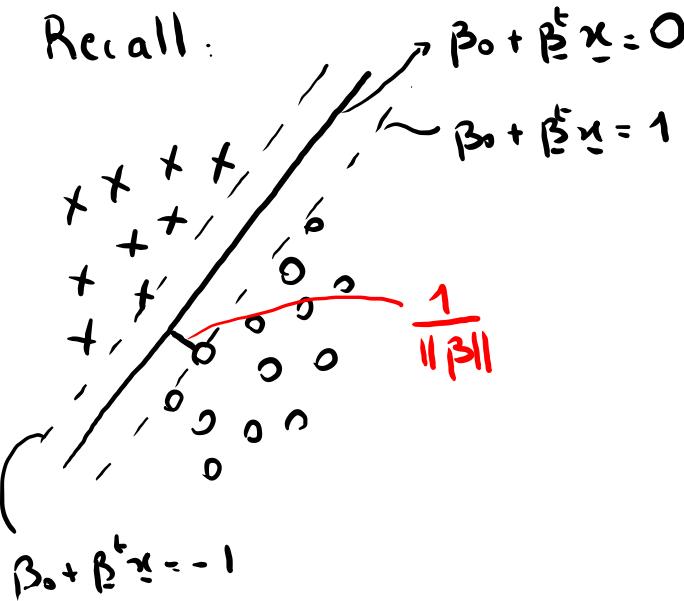


Recall:



SUPPORT VECTOR CLASSIFIER

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|^2 \rightarrow \beta_1^2 + \beta_2^2 + \dots + \beta_p^2$$

s.t.

$$C_1 \quad y_i (\beta_0 + \beta^T x_i) \geq 1 - \xi_i, \quad i=1, \dots, n$$

$$C_2 \quad \xi_i \geq 0 \quad i=1, \dots, n$$

$$C_3 \quad \sum_{i=1}^n \xi_i \leq C \rightarrow \text{tuning parameter}$$

QUADRATIC
UNCONSTRAINED
OPTIMIZATION

A BRIEF INTRODUCTION TO OPTIMIZATION

Some examples : Multinomial , Lasso , Ridge

$$\begin{array}{c} \uparrow \\ \text{equality} \\ \text{constraints} \\ \sum_{j=1}^p \theta_j = 1 \end{array} \quad \begin{array}{c} \uparrow \\ \text{inequality} \\ \text{constraints} \\ \sum |\beta_j| \leq C \end{array} \quad \begin{array}{c} \uparrow \\ \sum \beta_j^2 \leq C \end{array}$$

Synthetic optimization problem under equality constraints:

$$\begin{array}{ll} \min_{\underline{\alpha}} & f(\underline{\alpha}) \\ \text{s.t.} & h_i(\underline{\alpha}) = 0, \quad i=1, \dots, n \end{array}$$

LAGRANGIAN

$$l(\underline{\alpha}, \underline{\gamma}) = f(\underline{\alpha}) + \sum_{i=1}^n \gamma_i h_i(\underline{\alpha}) \quad \text{UNCONSTRAINED OPTIMIZATION}$$

with
 $\underline{\gamma} = (\gamma_1, \dots, \gamma_n)$ Lagrange multipliers

$$\gamma_1 h_1(\underline{\alpha}) + \gamma_2 h_2(\underline{\alpha}) + \dots$$

EXAMPLE Multinomial

$$\text{Log-lik. } l(\underline{\theta}) = \log(n!) - \sum_{i=1}^k \log(x_i!) + \sum_{i=1}^k x_i \log(\theta_i)$$

parameters: $\theta_1, \dots, \theta_k$

$$\max \underline{\theta} \leftarrow f(\underline{\alpha})$$

$$\text{s.t. } \sum_{i=1}^k \theta_i = 1$$

$$\Leftrightarrow 1 - \sum_{i=1}^k \theta_i = 0 \leftarrow h_1(\underline{\alpha})$$

SOLUTION

$$\frac{\partial l(\underline{\alpha}, \underline{\gamma})}{\partial \underline{\alpha}} = 0, \quad \frac{\partial l(\underline{\alpha}, \underline{\gamma})}{\partial \underline{\gamma}} = 0 \quad \}$$

$$\frac{\partial l(\underline{\alpha}, \underline{\gamma})}{\partial \gamma_1} = h_1(\underline{\alpha}) = 0$$

$\underline{\alpha}$ satisfies constraint 1

$$\hat{\underline{\alpha}}, \hat{\underline{\gamma}}$$

$\Rightarrow \hat{\underline{\alpha}}$ is the solution of the original constrained optimization problem

EXAMPLE: Multinomial

Lagrangian

$$\ell^*(\theta, \lambda) = \log(n!) - \sum_{i=1}^g \log(x_i!) + \sum_{i=1}^g n_i \log(\theta_i) + \lambda \left(1 - \sum_{i=1}^g \theta_i\right)$$

$$\frac{\partial \ell^*}{\partial \theta_i} = 0 \quad \frac{\partial \ell^*}{\partial \lambda} = 0$$

$$\frac{\partial \ell^*}{\partial \theta_i} = \frac{x_i}{\theta_i} - \lambda = 0, \quad i = 1, \dots, g$$

$$\theta_i = \frac{x_i}{\lambda}$$

$$\frac{\partial \ell^*}{\partial \lambda} = \left(1 - \sum_{i=1}^g \theta_i\right) = 0$$

$\underbrace{\sum_{i=1}^g \theta_i}_{=1} = 1$

$$\sum_{i=1}^g \frac{x_i}{\lambda} = 1 \Rightarrow \hat{\lambda} = \sum_{i=1}^g x_i = n$$

$$\Rightarrow \boxed{\hat{\theta}_i = \frac{x_i}{n}}$$

OBSERVED
FREQUENCIES

Let's generalize to the case of INEQUALITY CONSTRAINTS

(PRIMAL)
PROBLEM

$$\min_{\underline{\alpha}} f(\underline{\alpha})$$

s.t.

$$h_i(\underline{\alpha}) = 0 \quad i=1, \dots, p \quad (\text{equality constraints})$$

$$g_i(\underline{\alpha}) \leq 0 \quad i=1, \dots, n \quad (\text{inequality constraints})$$

GENERALIZED LAGRANGIAN

$$l(\underline{\alpha}, \underline{\lambda}, \underline{\delta}) = f(\underline{\alpha}) + \sum_{i=1}^p \lambda_i h_i(\underline{\alpha}) + \sum_{i=1}^n \delta_i g_i(\underline{\alpha})$$

with $\lambda_1, \dots, \lambda_p, \delta_1, \dots, \delta_n$ Lagrange multipliers

To justify the use of this function consider

$$\partial_p(\underline{\alpha}) = \max_{\underline{\lambda}, \underline{\delta}} l(\underline{\alpha}, \underline{\lambda}, \underline{\delta})$$
$$\quad \quad \quad \delta_i \geq 0$$

$$\min_{\substack{\underline{\alpha} \\ \text{constraints}}} f(\underline{\alpha}) = \min_{\underline{\alpha}} \partial_p(\underline{\alpha}) = \min_{\underline{\alpha}} \left(\max_{\substack{\underline{\lambda}, \underline{\delta} \\ \delta_i \geq 0}} l(\underline{\alpha}, \underline{\lambda}, \underline{\delta}) \right)$$

Take $\underline{\alpha}$ that does not satisfy the constraints

$h_i(\underline{\alpha}) \neq 0$ for some i
 $g_i(\underline{\alpha}) > 0$ for some i

$$\text{then } \max_{\substack{\underline{\lambda}, \underline{\delta} \\ \delta_i \geq 0}} l(\underline{\alpha}, \underline{\lambda}, \underline{\delta}) = +\infty$$

So the minimum must be achieved by an $\underline{\alpha}$ that satisfies the constraints.

Take $\underline{\alpha}$ that satisfies the constraints. Then $\sum \lambda_i h_i(\underline{\alpha}) = 0$, $g_i(\underline{\alpha}) \leq 0$

$$\text{Then take } \delta_i = 0 \Rightarrow l(\underline{\alpha}, \underline{\lambda}, \underline{\delta}) = f(\underline{\alpha})$$

PRIMAL
PROBLEM

$$\min_{\underline{\alpha}} \quad f(\underline{\alpha}) = \min_{\underline{\alpha}} \theta_p(\underline{\alpha}) = \min_{\underline{\alpha}} \max_{\substack{\underline{x}, \underline{\delta} \\ \delta_i \geq 0}} l(\underline{\alpha}, \underline{x}, \underline{\delta})$$

DUAL
PROBLEM

$$\max_{\substack{\underline{x}, \underline{\delta} \\ \delta_i \geq 0}} \theta_d(\underline{x}, \underline{\delta}) = \max_{\substack{\underline{x}, \underline{\delta} \\ \delta_i \geq 0}} \min_{\underline{\alpha}} l(\underline{\alpha}, \underline{x}, \underline{\delta})$$

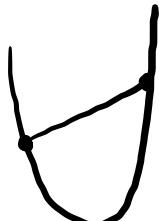
These two problems in general do not have the same solution.

One can show :

$$\max - \min \leq \min - \max$$

But they are the same in the case of f and g_i convex functions
and h_i linear, then

$$\max - \min = \min - \max$$



KARUSH-KUTIN-TUCKER (KKT) CONDITIONS

We assume $\partial_p(\underline{\alpha}^*) = \partial_s(\underline{\alpha}^*, \underline{\xi}^*)$. Then

$$\delta_i^* g_i(\underline{\alpha}^*) = 0 \quad i=1, \dots, n$$

COMPLEMENTARY
SLACKNESS
CONDITIONS

Proof: Let $\underline{\alpha}^*, \underline{\xi}^*, \underline{\delta}^*$ be the optimum.

$$\begin{aligned}
 \underline{f}(\underline{\alpha}^*) &= \min_{\underline{\alpha}} \ell(\underline{\alpha}, \underline{\gamma}^*, \underline{\xi}^*) = \min_{\underline{\alpha}} \left(f(\underline{\alpha}) + \sum_{i=1}^n \delta_i^* h_i(\underline{\alpha}) + \sum_{i=1}^n \delta_i^* g_i(\underline{\alpha}) \right) \\
 &\leq f(\underline{\alpha}^*) + \underbrace{\sum_{i=1}^n \delta_i^* h_i(\underline{\alpha}^*)}_{\geq 0} + \underbrace{\sum_{i=1}^n \delta_i^* g_i(\underline{\alpha}^*)}_{\leq 0} \\
 &\leq f(\underline{\alpha}^*)
 \end{aligned}$$

$$\Rightarrow \sum_{i=1}^n \delta_i^* g_i(\underline{\alpha}^*) = 0 \quad \overline{\Rightarrow} \quad \delta_i^* g_i(\underline{\alpha}^*) = 0 \quad i=1, \dots, n$$

$$\left. \begin{array}{l} \delta_i^* > 0 \Rightarrow g_i(\underline{\alpha}^*) = 0 \\ g_i(\underline{\alpha}^*) < 0 \Rightarrow \delta_i^* = 0 \end{array} \right\}$$

Back to SUPPORT VECTOR CLASSIFIER

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|^2 \xrightarrow{\text{red}} \beta_0^2 + \beta_1^2 + \dots + \beta_p^2 \quad (\text{w.r.t.})$$

s.t.

$$y_i(\beta_0 + \beta^T x_i) \geq 1 - \xi_i, \quad i=1, \dots, n$$

$$\xi_i \geq 0$$

$$\sum_{i=1}^n \xi_i \leq C$$

constraints:
convex sets

$$C_1: -[y_i(\beta_0 + \beta^T x_i) - 1 + \xi_i] \leq 0$$

$$C_2: -\xi_i \leq 0$$

$$C_3: \sum_{i=1}^n \xi_i - C \leq 0$$

LAGRANGIAN

$$L(\beta_0, \beta, \xi, \gamma, \delta, \mu) = \frac{1}{2} \|\beta\|^2 + \underbrace{\sum_{i=1}^n \gamma_i}_{\uparrow} [y_i(\beta_0 + \beta^T x_i) - 1 + \xi_i] - \underbrace{\sum_{i=1}^n \delta_i \xi_i}_{\downarrow} + \mu \left(\sum_{i=1}^n \xi_i - C \right)$$

for $\gamma_i \geq 0, \delta_i \geq 0, \mu \geq 0, i = 1, \dots, n$

Then

$$\begin{array}{lll} \min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|^2 & \stackrel{\text{PRIMAL PROBLEM}}{=} \min_{\beta_0, \beta, \xi} \max_{\gamma, \delta, \mu \geq 0} L(\cdot) & \stackrel{\text{convex}}{\leq} \max_{\gamma, \delta, \mu \geq 0} \min_{\beta_0, \beta, \xi} L(\cdot) \\ \text{constraints} & & \text{DUAL PROBLEM} \end{array}$$

Taking derivatives:

$$\frac{\partial L}{\partial \beta_0} = -\sum_{i=1}^n \gamma_i y_i = 0 \Rightarrow \sum_{i=1}^n \gamma_i y_i = 0$$

$$\frac{\partial L}{\partial \beta} = \beta - \sum_{i=1}^n \gamma_i y_i x_i = 0 \Rightarrow \boxed{\beta = \sum_{i=1}^n \gamma_i y_i x_i} \quad 0 \leq \delta_i = \mu - \gamma_i$$

$$\frac{\partial L}{\partial \xi_i} = -\gamma_i - \delta_i + \mu = 0 \Rightarrow \gamma_i = -\delta_i + \mu, \quad i = 1, \dots, n$$

Resubstituting back in the objective function gives:

$$\begin{aligned}
 \Theta_D(\gamma, \xi, \mu) &= \frac{1}{2} \beta^T \beta - \sum_{i=1}^n \left(\underbrace{\gamma_i y_i \beta_0 + \gamma_i \beta^T \underline{x}_i}_{=0} - \gamma_i + \gamma_i \varepsilon_i \right) - \sum_{i=1}^n \delta_i \varepsilon_i + \mu \left(\sum_{i=1}^n \varepsilon_i - C \right) \\
 &= \frac{1}{2} \beta^T \beta - \beta^T \beta + \sum_{i=1}^n \gamma_i - \mu C \\
 &\stackrel{\beta = \sum_{i=1}^n \gamma_i y_i x_i}{=} \sum_{i=1}^n \gamma_i - \frac{1}{2} \sum_{i,j=1}^n \gamma_i \gamma_j y_i y_j \underline{x_i \cdot x_j} - \mu C \\
 &\quad \text{scalar product} \\
 &= \sum_{i=1}^n (-\gamma_i - \delta_i + \mu) \varepsilon_i = 0
 \end{aligned}$$

DUAL PROBLEM

$$\begin{aligned}
 \max_{\gamma, \mu} \quad & \sum_{i=1}^n \gamma_i - \frac{1}{2} \sum_{i,j=1}^n \gamma_i \gamma_j y_i y_j \underline{x_i \cdot x_j} - \mu C \\
 \text{s.t.} \quad & 0 \leq \gamma_i \leq \mu \quad [\mu - \gamma_i \geq 0] \\
 & \sum_{i=1}^n \gamma_i y_i = 0
 \end{aligned}$$

μC : constant
 $\mu = \frac{\text{constant}}{C}$

Solving the dual problem is equivalent to solving the original primal problem so

$$\beta = \sum_{i=1}^n \gamma_i y_i x_i$$

Note: \underline{x} has dimension n , so the solution is not too affected by a large P

Moreover, many of the γ_i will be zero.

Using the KKT conditions, at the optimum:

$$\delta_i^* g_i(\alpha^*) = 0, i=1, \dots, n$$

$$\gamma_i [y_i(\beta_0 + \beta^t x_i) - 1 + \varepsilon_i] = 0, i=1, \dots, n$$

Used to derive
 $\hat{\beta}_0$ by taking an
average over $i=1, \dots, n$

But more importantly, $\gamma_i \neq 0$ identifies the support vectors!

This is because if an observation is correctly classified (on the right side of margin)

then $\varepsilon_i = 0$ and $y_i(\beta_0 + \beta^t x_i) > 1 \Rightarrow y_i(\beta_0 + \beta^t x_i) - 1 + \varepsilon_i \neq 0 \Rightarrow \gamma_i = 0$

So $\gamma_i \neq 0$ identifies the points on the margin or misclassified.