

$\mu^+$ : spam prototype

$\mu^-$ : non-spam prototype

$x$ : feature vector of new email

$x$  is classified as spam if

$$\|x - \mu^+\|_2 < \|x - \mu^-\|_2$$

$$\Rightarrow \|x - \mu^+\|_2^2 < \|x - \mu^-\|_2^2$$

$$\Rightarrow \|\mu^+\|_2^2 - 2 \langle x, \mu^+ \rangle < \|\mu^-\|_2^2 - 2 \langle x, \mu^- \rangle$$

$$\Rightarrow \langle x, 2(\mu^+ - \mu^-) \rangle + \|\mu^-\|_2^2 - \|\mu^+\|_2^2 > 0$$

$$\equiv \langle x, w \rangle + b > 0$$

$$\text{with } w = 2(\mu^+ - \mu^-)$$

$$b = \|\mu^-\|_2^2 - \|\mu^+\|_2^2$$

→ decision boundary of LwP classifier is always a line or a hyperplane if the dist. fun. is Euclidean.

⇒ LwP is a linear classifier.

→ Model

w

→ vector of dim d,  
the same as that  
of features

→ aka normal of  
the hyperplane

b  
→ bias, it is  
a scalar term

→ Prediction: test point  $x$

by checking if the following holds

$$w^T x + b > 0$$

→ Dec. boun: line/hyperplane

where

$$w^T x + b = 0$$

acts as a threshold:  
how large does  $w^T x$   
have to be in order  
to ~~be~~ classify  $x$  as  
spam

Pf:  $x, y$  2 pt. on

hyperplane

$$\Rightarrow w^T x + b = 0$$

$$\Rightarrow w^T y + b = 0$$

$$\Rightarrow w^T(x-y) = 0$$

→ convenient to not have bias term.

let  $\tilde{x} = [x, 1] \in \mathbb{R}^{d+1}$

$\tilde{w} = [\omega_0, \omega_1, \dots, \omega_d] \in \mathbb{R}^{d+1}$

where  $w = [\omega_0, \omega_1, \dots, \omega_d] \in \mathbb{R}^d$

then  $\boxed{\tilde{w}^T \tilde{x} = \tilde{w}_0^T x + \tilde{w}_d}$

effectively, acts as a bias term.

Metric learning : learning the distance fn.

↳ e.g. Mahalanobis metrics.

let  $A$ : symm mat.  $\in \mathbb{R}^{d \times d}$

def.  $d_A(x, y) = \sqrt{(x-y)^T A (x-y)}$

→  $A$  satisfies positive semi-definiteness (PSD)

$\forall x, \boxed{x^T A x \geq 0}$

$\Rightarrow d_A(x, y) \geq 0$

LwP w/ Mahalanobis

$$(x - \mu^+)^T A (x - \mu^+) \leq (x - \mu^-)^T A (x - \mu^-)$$

$$d_A(x, \mu^+) \leq d_A(x, \mu^-)$$

$$\Rightarrow x^T A (\mu^+ - \mu^-) + \mu^-^T A \mu^- - \mu^+^T A \mu^+ \geq 0$$

$$\equiv \langle x, \omega \rangle + b \geq 0$$

where  $\omega = 2A(\mu^+ - \mu^-)$

$$b = \mu^-^T A \mu^- - \mu^+^T A \mu^+$$

Let  $A = LL^T$ , where  $L$  need not be sym. or PSD

$$\begin{aligned} \text{Then } d_A(x, y) &= \sqrt{(x-y)^T A (x-y)} = \sqrt{(x-y)^T L L^T (x-y)} \\ &= \|L^T x - L^T y\|_2 \end{aligned}$$

→ Mahalan dist is the euc. dist if  $x$  is transform to  $Lx$

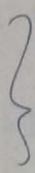
## Hyperparameter tuning

$k$  in kNN

$r$  in rNN

metric used in LwP

$w$ : model vector



hyperparameters  
of ML algo



parameters

→ tuned to give  
highest test  
accuracy.

## Held-out validation

## Multi-fold cross validation

hyp param: help in learning proc

param: used to make decisions  
(e.g. wt. co-eff of lin. reg. model)

- finding NN for a new pt. is expensive
- all training w/ NN needs to be stored, but only prototype stored in LwP
- entire training set is the model for NN
- model size is NN inc w/ the amt. of training data; such models ka non-parametric models ✓
- LwP model: just 2 vectors, indep. of training set; such models ka parametric models ✓

## Decision trees

- faster way to perf. NN searches
- model size can be large
- good train perf. but bad train perf. (ka overfitting)
- balanced DT's take  $O(\log n)$  time for prediction

LS

## Large Margin Classifiers

d of p from hyperplane  $w^T x + b = 0$  is  $\frac{|w^T p + b|}{\|w\|_2}$

Let train data be  $\{(x^i, y^i)\}_{i=1}^n$   
where  $x^i \in \mathbb{R}^d$   
 $y^i \in \{-1, 1\}$

Classify each pt. correctly :  $\text{sign}(w^T x^i + b) = y^i$   
or  $y^i \cdot (w^T x^i + b) > 0$

Pt's away from boundary :  $\min_{i=1 \dots n} \frac{|w^T x^i + b|}{\|w\|_2}$  should be as large as possible

SVM : find a lin. classif that  
→ perf. classifies data  
→ keeps pt's as far away as possible

objective max  $w, b$   $\left\{ \min_{i=1 \dots n} \frac{|w^T x^i + b|}{\|w\|_2} \right\}$  at  $y^i \cdot (w^T x^i + b) > 0$   $\forall i = 1 \dots n$  constraint  
→ optimization prob w/ constraints.

$$\frac{|\omega^T x^i + b|}{\|\omega\|_2} = \frac{|y^i \cdot (\omega^T x^i + b)|}{\|\omega\|_2}$$

$$\min_{i=1 \dots n} \frac{|y^i \cdot (\omega^T x^i + b)|}{\|\omega\|_2} = \min_{i=1 \dots n} \frac{y^i \cdot (\omega^T x^i + b)}{\|\omega\|_2} \quad (\because \text{perf. classif.})$$

let  $y^o$  be the closest pt

$$\Rightarrow \min_{i=1 \dots n} y^i \cdot (\omega^T x^i + b) = y^o \cdot (\omega^T x^o + b)$$

$$\text{let } \varepsilon = y^o \cdot (\omega^T x^o + b)$$

$$\tilde{\omega} = \omega / \varepsilon$$

$$\tilde{b} = b / \varepsilon$$

$$\Rightarrow y^i \cdot (\tilde{\omega}^T x^i + \tilde{b}) \geq 1 \quad \forall i=1 \dots n$$

$$\& \min_{i=1 \dots n} \frac{y^i \cdot (\tilde{\omega}^T x^i + \tilde{b})}{\|\tilde{\omega}\|_2} = \frac{1}{\|\tilde{\omega}\|_2}$$

thus  $\max_{\tilde{\omega}, \tilde{b}} \left\{ \frac{1}{\|\tilde{\omega}\|_2} \right\}$  st  $y^i \cdot (\tilde{\omega}^T x^i + \tilde{b}) \geq 1 \quad \forall i=1 \dots n$

$$\Downarrow \min_{\tilde{\omega}, \tilde{b}} \left\{ \|\tilde{\omega}\|_2 \right\}$$

$$\Downarrow \min_{\tilde{\omega}, \tilde{b}} \left\{ \|\tilde{\omega}\|_2^2 \right\}$$

$$\Downarrow \min_{\tilde{\omega}, \tilde{b}} \frac{1}{2} \left\{ \|\tilde{\omega}\|_2^2 \right\}$$

$$\min_{\tilde{\omega}, \tilde{b}} \left\{ \frac{1}{2} \|\tilde{\omega}\|_2^2 \right\} \text{ st } y^i \cdot (\tilde{\omega}^T x^i + \tilde{b}) \geq 0$$

if the classific can't be perf, then model using

$$\min_{\tilde{w}, b, \xi_i} \frac{1}{2} \|\tilde{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i \cdot (\tilde{w}^T x_i + b) \geq 1 - \xi_i \quad \forall i=1 \dots n$$

$$\xi_i \geq 0 \quad \forall i=1 \dots n$$

→ slack variables

Thus  $\xi_i$  allows  $y_i \cdot (\tilde{w}^T x_i + b) \leq 1$

$$\text{Thus we need } \xi_i = 1 - y_i \cdot (\tilde{w}^T x_i + b)$$

$$\text{however } \xi_i \geq 0$$

Thus if  $y_i \cdot (\tilde{w}^T x_i + b) \geq 1$ ,  $\xi_i = 0$  (no slack)  
regd

$$\text{Thus } \xi_i = [1 - y_i \cdot (\tilde{w}^T x_i + b)]_+ \rightarrow [x]_+ = \max\{x, 0\}$$

Hinge loss fn

Hinge Loss - captures how well a data pt was classified

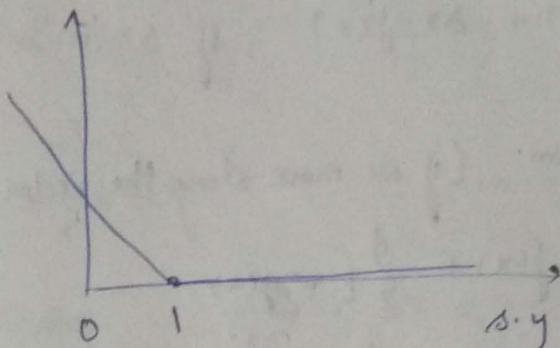
→ on data pt.  $(x, y)$ ,  $y \in \{-1, 1\}$

prediction score  $s$  (for lin. model  $(w, b)$ )  
we have,  $s = w^T x + b$ )

$s - y \geq 0$  for correct classific

$s - y > 0$  for large margin

$$\ell_{\text{hinge}}(s, y) = [1 - s \cdot y]_+ = \begin{cases} 0 & \text{if } s \cdot y \geq 1 \\ 1 - s \cdot y & \text{if } s \cdot y < 1 \end{cases}$$



Final form of CSVM

$$\min_{\tilde{w}, \tilde{b}, \{\xi_i\}} \frac{1}{2} \|\tilde{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$\text{st } y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1 - \xi_i \quad \forall i=1 \dots n$$

$$\& \xi_i \geq 0 \quad \forall i=1 \dots n$$

$$\rightarrow \min_{\tilde{w}, \tilde{b}} \frac{1}{2} \|\tilde{w}\|_2^2 + C \sum_{i=1}^n \ell_{\text{hinge}}(\tilde{w}^T x_i + \tilde{b}, y_i)$$

(L6)

Corollary of Taylor's th:

$$f(x + \Delta x) \approx f(x) + \Delta x \cdot f'(x), \text{ if } \Delta x \text{ is "small"}$$

Taylor's th. in higher dim: (if we move along the vector  $t = (t_1, t_d)$ )

then  $f(x+t) \approx f(x) + \sum_{i=1}^d t_i \cdot \frac{\partial f}{\partial x_i}(x)$

$$\boxed{f(x+t) = f(x) + t^T \nabla f(x)} \quad \text{if } t \text{ is "small"}$$

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_1 \partial x_d} & \dots & \frac{\partial^2 f}{\partial x_d \partial x_1} \\ \vdots & - & - & & \frac{\partial^2 f}{\partial x_d^2} \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & - & - & & \end{bmatrix}$$

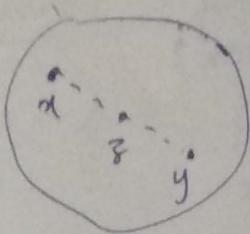
$\rightarrow$  Hessian:  $d \times d$  mat: symmetric

if  $\nabla f = 0$  &  $\nabla^2 f$  is a PSD mat  $\Rightarrow$  local/global min.

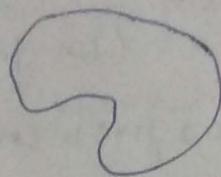
$\nabla f = 0$  &  $\nabla^2 f$  is a NSD mat  $\Rightarrow$  " " max

$$\nabla^2 f(x) \leq 0$$

## Convex Sets

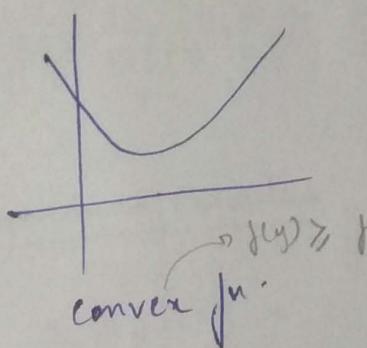


Convex set



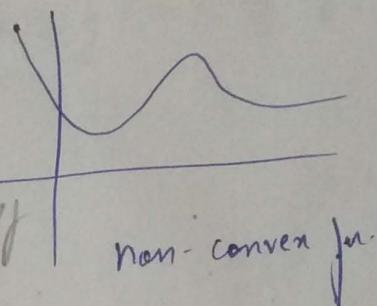
non-convex set

Convex set  $\left\{ C \subseteq \mathbb{R}^d, \forall x, y \in C \text{ & } \forall \lambda \in [0, 1] \right.$   
 $\left. \text{Convex set} \quad \begin{cases} \exists z = \lambda \cdot x + (1-\lambda) \cdot y \in C \end{cases} \right.$



$$\forall x, y, \forall \lambda \in [0, 1]$$

$$\begin{aligned} z &= \lambda \cdot x + (1-\lambda) \cdot y \\ f(z) &\leq \lambda \cdot f(x) + (1-\lambda) \cdot f(y) \end{aligned} \quad \left. \begin{array}{l} \text{convex fn.} \\ \rightarrow \text{must lie below all its chords.} \end{array} \right\}$$



non-convex fn.

$\rightarrow$  tang. to  $f$  at  $x_0$  is the hyperplane.

$$\nabla f(x_0)^T (x - x_0) + f(x_0) = 0$$

$\rightarrow$  sum of 2 convex fn's is always convex

$\rightarrow$  for diff.  $f^n$ , it must lie above all its tangents

$\rightarrow$  for twice diff. fn,  $\nabla^2 f$  must be PSD every where

Eg

$$f(x) = 0$$

$$f(x) = c$$

$$f(x) = a^T x \quad (\text{lin. fn are convex})$$

c.  $f(x)$ , if  $C_1, 0 \geq f(x)$  is convex

if  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex

$f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex & non dec.

then  $f \circ g: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex

$$f(x) = \|x\|_2 \text{ is convex}$$

L7

Subgradient - for convex, non-diff. fns.  $\rightarrow$  (possibly)

for a given  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , poss. non-diff but convex  
at  $x^*$

any vector  $g$  is the subgrad of  $f$  at  $x^*$  if

$$f(x) \geq g^T(x - x^*) + f(x^*) \quad \forall x$$

Subdifferential - set of all subgrad of  $f$  at  $x^*$ , denoted by  $\partial f(x^*)$

$$\partial f(x^*) \triangleq \{g : f(x) \geq g^T(x - x^*) + f(x^*), \forall x\}$$

Subgradient Calculus

2/12/19

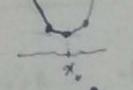
$$\rightarrow \partial(c \cdot f)(x) = c \cdot \partial f(x) \\ = \{c \cdot v : v \in \partial f(x)\}$$



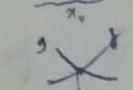
$$\rightarrow \partial(f+g)(x) = \partial f + \partial g \\ = \{u+v : u \in \partial f, v \in \partial g\}$$



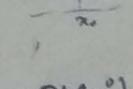
$$\rightarrow \partial f(a^T x + b) = \partial f(a^T x + b) \circ a \\ = \{c \cdot a : c \in \partial f(a^T x + b)\}$$



$$\rightarrow h(x) = \max \{f(x), g(x)\}$$



?? if  $f(x^*) > g(x^*)$   $\partial h(x^*) = \partial f(x^*)$   
 ?? if  $f(x^*) = g(x^*)$   $\partial h(x^*) = \{\lambda u + (1-\lambda)v : u \in \partial f(x^*), v \in \partial g(x^*), \lambda \in [0,1]\}$



$\rightarrow$  stationary pt at  $x^*$  if  $0 \in \partial f(x^*)$ . Local max/min is stationary even for non-diff. fns.

## Subgradient for hinge loss.

$$l_{\text{hinge}}(x) = \max\{1-x, 0\} = \max\{f(x), g(x)\}$$

$$\partial l_{\text{hinge}}(x) = l'_{\text{hinge}}(x) \quad \text{if } x \neq 1$$

At  $x=1$

$$\begin{aligned} f(x) &= 1-x \quad (\text{diff. able}) \Rightarrow \partial f = f'(x) = -1 \\ g(x) &= 0 \quad ("") \Rightarrow \partial g = g'(x) = 0 \end{aligned}$$

$$\Rightarrow \partial l_{\text{hinge}}(x) = \begin{cases} -1 + (1-\lambda)0 : \lambda \in [0, 1] \end{cases} \\ = [-1, 0]$$

$$l_{\text{hinge}}(y^i, \langle w, x^i \rangle) = [1 - y^i \langle w, x^i \rangle]_+$$

Finding  $v^i \in \partial l_{\text{hinge}}(y^i, \langle w, x^i \rangle)$

$$v^i = \begin{cases} 0 & \rightarrow \text{if } y^i \langle w, x^i \rangle > 1 \\ -y^i \cdot x^i & \rightarrow \text{if } y^i \langle w, x^i \rangle < 1 \\ c \cdot y^i \cdot x^i & \rightarrow \text{if } y^i \langle w, x^i \rangle = 1 \\ c \in [-1, 0] \end{cases}$$

## Subgradient descent

1. given :-  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  to minimize

2. initialize  $w^0 \in \mathbb{R}^d$

3. for  $t = 0, 1, \dots$

1. obtain subgrad.  $g^t \in \partial f(w^t)$

2. choose step length  $\eta_t$  aka learning rate

3. update  $w^{t+1} \leftarrow w^t - \eta_t \cdot g^t$

4. repeat until convergence.

For "nicely behaved" convex fns

set  $\eta_t = \eta/\sqrt{t}$   $\rightarrow \eta$ : hyperparameter  
or  $\eta_t = \eta/t$

Ammijo Rule: try some  $\eta_t$ , if not "nice",  $\downarrow \eta_t$

Adagrad: use a diff  $\eta$  for each dim. of  $w$

$\eta_t^*$  repl. by a diag. mat  $E^t$  i.e.  $w^{t+1} \leftarrow w^t - E^t g^t$

Adam: uses momentum methods (infuses prev. grad. into the curr. grad.)

Convergence  $\rightarrow \|w^{t+1} - w^t\| \rightarrow 0$

↳ general heuristics:

if  $\|g^t\|_2$  has become too small

if  $|f(w^{t+1}) - f(w^t)|$  "

if  $f(w^t)$  is small enough

$$J(\omega) = \frac{1}{2} \|\omega\|_2^2 + C \sum_{i=1}^n [1 - y_i \cdot \omega^T x_i]_+ \quad (\text{ignoring bias})$$

?

$$\nabla J(\omega) = \omega + C \sum_{i=1}^n g_i y_i \cdot x_i, \quad g_i \in \nabla l_{\text{hinge}}(y_i \cdot \omega^T x_i)$$

$$\omega^{\text{new}} = \omega - \eta \cdot \nabla J(\omega)$$

$$= (1-\eta)\omega - \eta C \sum_{i=1}^n g_i y_i \cdot x_i$$

Say if  $n=1$

$$\left| \begin{array}{l} \text{CSVM : } \underset{\omega, b}{\text{minimise}} \frac{1}{2} \|\omega\|_2^2 + C \sum_{i=1}^n l_{\text{hinge}}(\omega^T x_i + b, y_i) \\ y_i(\omega^T x_i + b) \geq 1 - \xi_i \Rightarrow \text{all} \geq 0 \end{array} \right.$$

$$\omega^{\text{new}} = (1-\eta) \cdot \omega - \eta C \cdot g' y' \cdot x'$$

If small  $\eta$ : not much change in  $\omega$

If large  $\eta$ :  $\omega$  changes as gradient dictates

If  $\omega$  does well on  $(x', y')$

Say  $y' \cdot \omega^T x' > 1$ , then  $g' = 0$

$\Rightarrow$  no change in  $\omega$  due to pt  $(x', y')$

If  $\omega$  does badly on  $(x', y')$

Say  $y' \cdot \omega^T x' < 0$ , then  $g' = 1$

$$\Rightarrow \omega^{\text{new}} = (1-\eta) \cdot \omega + \eta C y' \cdot x'$$

$$\Rightarrow y' \cdot (\omega^{\text{new}})^T \cdot x'$$

$$= (1-\eta)y' \cdot \omega^T x' + \eta C \cdot \|x'\|_2^2$$

$\rightarrow$  Model GD actively tries to make the model perform better on all data pts.

## Stochastic gradient method

$$\nabla J(\omega) = \omega + c \cdot \sum_{i=1}^n g_i y_i \cdot x_i, \quad g_i \in \nabla \ell_{\text{hinge}}(y_i \cdot \omega^T x_i)$$

Calculating each  $g_i$  takes  $\Theta(d)$  time, thus total  $\Theta(nd)$  time

At each time, choose a random data pt.  $(x_t^i, y_t^i)$

$$\nabla J(\omega) \approx \omega + c \cdot g_t^i y_t^i \cdot x_t^i - \text{only } \Theta(d) \text{ time}$$

Mini-batch SGD - if data is diverse, stoc. grad. may vary quite a lot

Choose  $B$  random data pts., say  $(x_t^{i_1}, y_t^{i_1}) \dots (x_t^{i_B}, y_t^{i_B})$   
 ↳ mini batch size

$$\nabla J(\omega) \approx \omega + c \cdot \sum_{b=1}^B g_t^{i_b} y_t^{i_b} \cdot x_t^{i_b}$$

→  $\Theta(Bd)$  time

if  $B=n$ , MBSGD becomes GD.

## Constrained optimization

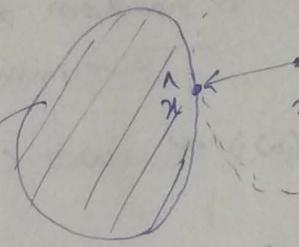
?? Method 1 : Interior pt. method

Method 2 : Projected grad. desc.

→ perf. SGD as usual. However, if this causes us to go out of feasible set  $C$ , go back to feasible set.  
going back : projection step.

$$\hat{x} = \Pi_C(x) = \arg \min_{z \in C} \|x - z\|_2^2$$

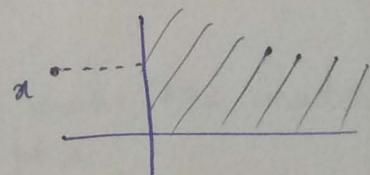
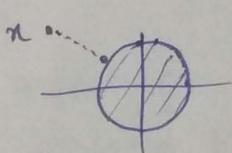
1. choose  $g^t \in \partial f(w^t)$
2.  $u^{t+1} \leftarrow w^t - \eta_t \cdot g^t$
3.  $w^{t+1} \leftarrow \Pi_C(u^{t+1})$
4. Repeat until convergence



Useful projections.  $\hat{x} = \Pi_C(x)$

$$C = \{x : \|x\|_2 \leq 1\}$$

$$C = \{x : x_i \geq 0\}$$



$$\hat{x} = \begin{cases} x & \text{if } \|x\|_2 \leq 1 \\ \frac{x}{\|x\|_2} & \text{if } \|x\|_2 > 1 \end{cases}$$

$$\hat{x}_i = \begin{cases} x_i & \text{if } x_i \geq 0 \\ 0 & \text{if } x_i < 0 \end{cases}$$

### Method 3 : Creating a dual problem

Supp we wish to solve

$$\min_{x \in \mathbb{R}^d} f(x)$$

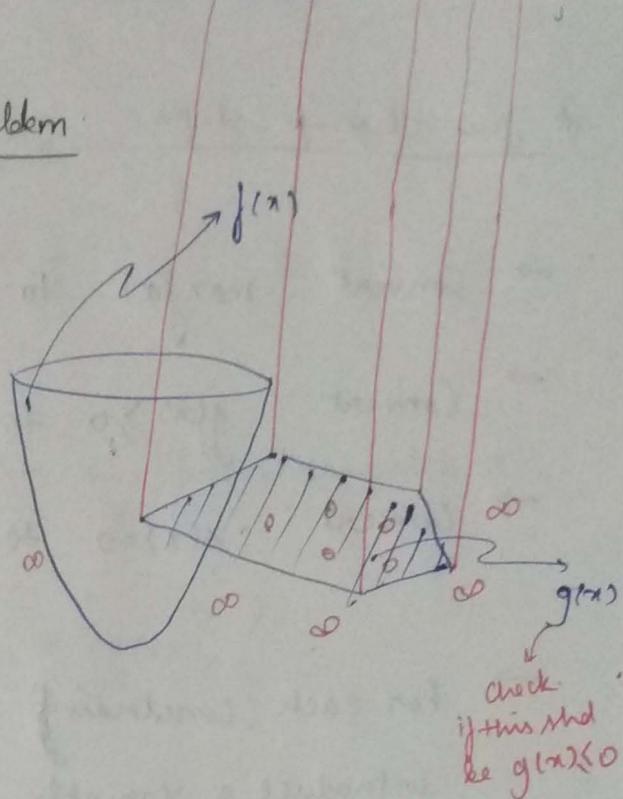
$$\text{s.t } g(x) \leq 0$$

Construct a barrier  $r(x)$

$$\begin{aligned} \text{s.t } r(x) &= 0 & \text{if } g(x) \leq 0 \\ r(x) &= \infty & \text{otherwise} \end{aligned}$$

and simply solve

$$\min_{x \in \mathbb{R}^d} f(x) + r(x)$$



One way to construct such a barrier

$$r(x) = \max_{\alpha \geq 0} \alpha \cdot g(x)$$

Thus we want to solve

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) + \max_{\alpha \geq 0} \alpha \cdot g(x) \right\}$$

Same as

$$\min_{x \in \mathbb{R}^d} \left\{ \max_{\alpha \geq 0} \left\{ f(x) + \alpha \cdot g(x) \right\} \right\}$$

4/12/19

## A few cleanup steps

→ Convert max  $f(x)$  to  $\min \{-f(x)\}$

→ Convert  $g_i(x) \geq 0$  to  $-g_i(x) \leq 0$

→ Convert  $s(x) = 0$  to  $s(x) \leq 0$   
 $\Delta -s(x) \leq 0$

→ For each constraint introduce a variable  $x_i^*$   $i=1\dots c$

$$x_i^* \quad i=1\dots c$$

→ dual var. &

or Lagrangian multipliers

5/12/19

## The Lagrangian

$$\min_x f(x)$$

$$\text{st } g_1(x) \leq 0 \\ g_2(x) \leq 0 \\ \vdots \\ g_c(x) \leq 0$$

$$\left. \begin{array}{l} \mathcal{L}(x, \alpha) = f(x) + \sum_{c=1}^C \alpha_c \cdot g_c(x) \\ \text{Lagrangian.} \end{array} \right\}$$

If  $x$  violates even one constraint, we have

$$\max_{\alpha \in \mathbb{R}^C} \{\mathcal{L}(x, \alpha)\} = \infty$$

$$\alpha_c > 0$$

If  $x$  satisfies every single constraint, we have

$$\max_{\alpha \in \mathbb{R}^C} \{\mathcal{L}(x, \alpha)\} = f(x)$$

$$\alpha_c > 0$$

$$\min_{x \in \mathbb{R}^d} \left\{ \max_{\alpha \in \mathbb{R}^G} \left\{ \begin{array}{l} f(x) + \\ \left\{ \sum_{c=1}^G \alpha_c \cdot g_c(x) \right\} \end{array} \right\} \right\}$$

The dual prob.

original prob :- primal problem w/ primal variables  $x$ .

$$\rightarrow \min_{x \in \mathbb{R}^d} \left\{ \max_{\alpha \in \mathbb{R}^G} \left\{ \begin{array}{l} f(x) + \\ \left\{ \sum_{c=1}^G \alpha_c \cdot g_c(x) \right\} \end{array} \right\} \right\}$$

dual prob

$$? \rightarrow \max_{\alpha \in \mathbb{R}^G} \left\{ \min_{x \in \mathbb{R}^d} \left\{ f(x) + \left\{ \sum_{c=1}^G \alpha_c \cdot g_c(x) \right\} \right\} \right\}$$

let  $\hat{x}^P, \hat{\alpha}^P$  be the soln to primal prob.  
 $\hat{x}^D, \hat{\alpha}^D$  " " dual prob.

strong duality :  $\hat{x}^P = \hat{x}^D$  if the orig prob. is convex  
 & "nice"

complementary slackness :  $\hat{\alpha}_c^D \cdot g_c(\hat{x}^D) = 0$   $\forall$  constraints

$$\text{SVM} \underset{\mathbf{w}, b}{\text{min}} \left\{ \sum_{i=1}^n \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|_2} \right\} \text{ s.t. } (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i = 1 \dots n$$

Hand SVM w/out bias

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{s.t.} \quad 1 - y_i^i \cdot \mathbf{w}^T \mathbf{x}_i^i \leq 0 \quad \forall i = 1 \dots n$$

→  $n$  constraints →  $n$  dual vars  
→  $\alpha \in \mathbb{R}^n$

$$f(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i^i \cdot \mathbf{w}^T \mathbf{x}_i^i)$$

?? Primal:  $\underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \left\{ \underset{y \in \{-1, 1\}^n}{\text{argmax}} \left\{ \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i^i \cdot \mathbf{w}^T \mathbf{x}_i^i) \right\} \right\}$

?? Dual:  $\underset{\alpha \geq 0}{\text{argmax}} \left\{ \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i^i \cdot \mathbf{w}^T \mathbf{x}_i^i) \right\} \right\}$

Simplifying the dual problem

The inner is the dual problem

∴ This is an unconstrained prob. w/ a convex & differentiable objective, use first order optimality

We get  $\boxed{\mathbf{w} = \sum_{i=1}^n \alpha_i y_i^i \cdot \mathbf{x}_i^i}$

Subs. back in dual prob.

$$\underset{\alpha \geq 0}{\text{argmax}} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i^i, \mathbf{x}_j^j \rangle \right\}$$

Once the optimal value of  $\alpha$  is known, use

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i^i \cdot \mathbf{x}_i^i$$

## Support Vectors

$\alpha_i$  + data pt's

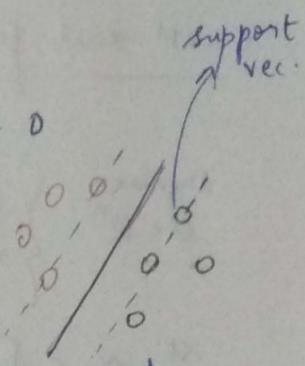
Support vec :- data pt's for which  $\alpha_i \neq 0$

Using complementary slackness

$$\alpha_i(1 - y^i w^T x^i) = 0$$

$\Rightarrow$  only those data pt's can become SV for which

$$y^i w^T x^i = 1 \Rightarrow \text{at margin}$$



## Dual for CSVM

If we have bias  $b$  as well as slack  $\xi$  then dual is given by

$$\underset{\alpha \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle x^i, x^j \rangle \right\}$$

$$\text{s.t. } \alpha_i \in [0, C] \text{ & } \sum_{i=1}^n \alpha_i y^i = 0$$

Analogy

$\alpha^i$ : Force on hyperplane  $w$  by  $i^{th}$  data pt in  $y^i$  dir.

$$\sum \alpha^i = 0 \text{ as } \sum_{i=1}^n \alpha^i y^i = 0$$

$$\text{Also, since } w = \sum_{i=1}^n \alpha^i y^i \cdot x^i$$

$$\sum \alpha^i = 0$$

Thus the support vec. mechanically support the hyperplane.

CSVM dual problem

$$\underset{\alpha \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle x^i, x^j \rangle \right\}$$

$$\text{st } \alpha_i \in [0, C] \quad \& \quad \underbrace{\sum_{i=1}^n \alpha_i y^i}_{\text{if bias is present}} = 0$$

Coordinate Descent (popularly applied to dual probs.)

(SDCA - stochastic dual  
co-ord ascent)  
Only one co-ord changed in a single step.

$$\min_{x \in \mathbb{R}^d} f(x) \text{ st } x \in C$$

$\therefore$  dual prob.  
is a maximiz. pr.

$$\nabla_j f(x) = \frac{\partial f}{\partial x_j}$$

Projected Co-ord descent

1. For  $t = 0, 1, \dots$

1. Select a coord  $j_t \in [d]$
2. Let  $u_{j_t}^{t+1} \leftarrow x_{j_t}^t - \eta \nabla_{j_t} f(x^t)$
3. Let  $u_j^{t+1} \leftarrow x_j^t \text{ for } j \neq j_t$
4. Project  $x^{t+1} \leftarrow \Pi_C(u^{t+1})$
5. Repeat until convergence

CCD: cyclic order

$$j_t = 1, 2, \dots, d, 1, \dots, d, \dots$$

SCD: randomly chosen

Block CD: choose a set of co-ord's to be updated at each  $t$