

M8. Classification: Support Vector Machines (SVMs)

Manikandan Narayanan

Week 13 (Oct 20- 2025)

PRML Jul-Nov 2025 (Grads Section)

Acknowledgment of Sources

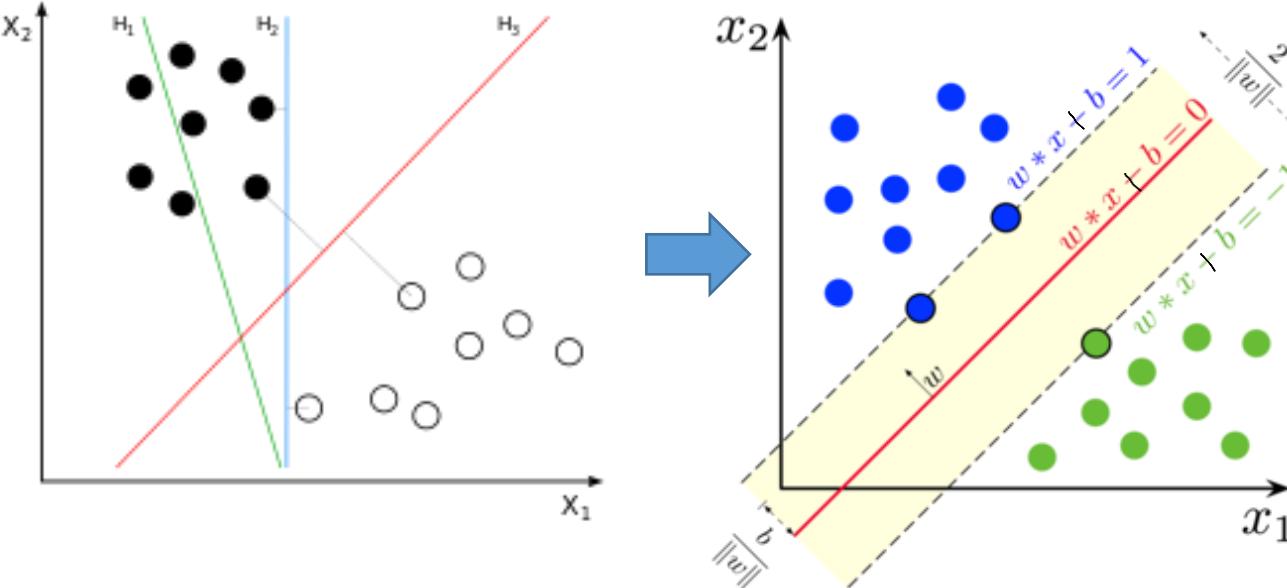
- Slides based on content from related
 - Courses:
 - IITM – Profs. Arun/Harish/Chandra’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited respectively as [AR], [HR], [CC], [BR] in the bottom right of a slide.
 - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
 - Books:
 - PRML by Bishop. (content, figures, slides, etc.) – cited as [CMB]
 - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
 - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]
 - Information Theory, Inference and Learning Algorithms by David JC MacKay – [DJM]

Outline for Module M8

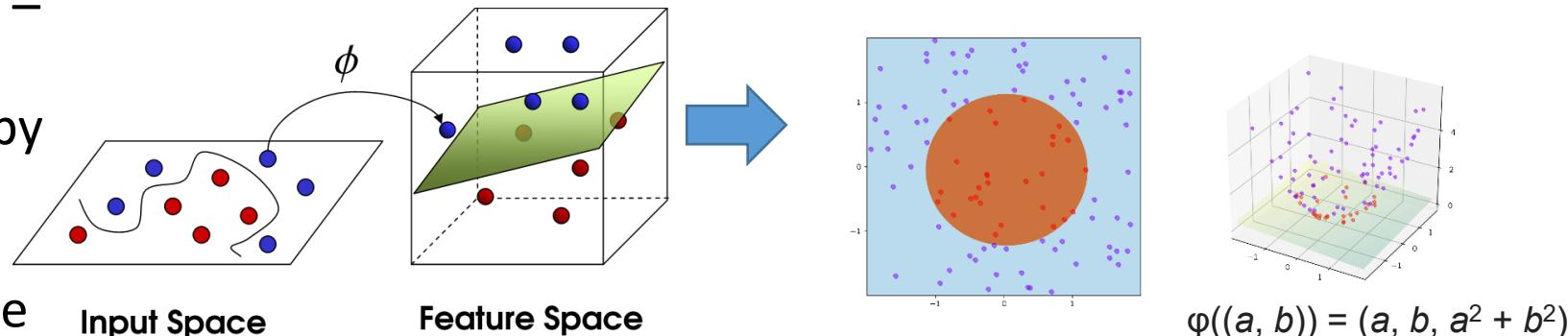
- M8. Classification (Support Vector Machines)
 - **M8.0 Introduction/Motivation**
 - (concrete understanding of SVMs – beyond popular pictures & software)
 - M8.1 SVM Problem Statement
 - (Hard/Soft-Margin SVM Problems)
 - M8.2 SVM Solution
 - (Background: Constrained optimization - KKT & Primal-Dual)
 - (SVM Dual Problem & Optimization algo. sketch)
 - M8.3 SVM Interpretations
 - (Support vectors, Kernels, Loss function view)
 - M8.4 Concluding thoughts

SVM hard-margin: popular pic. → geometry

- (Linear) SVM – max margin, sparse support vectors



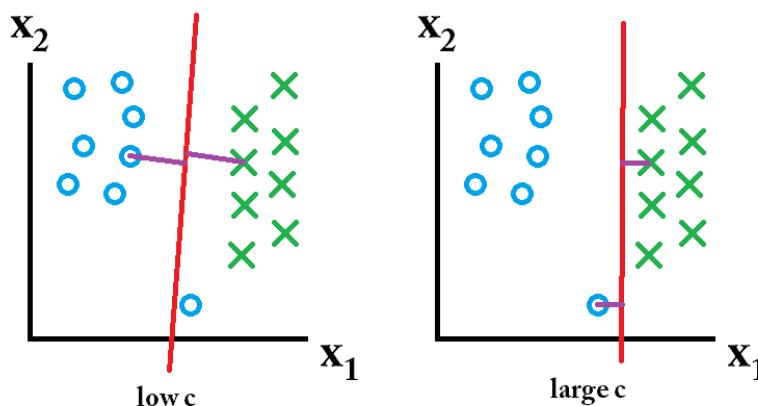
- (Non-linear) SVM – uses non-linear kernels followed by applying SVM above in the feature map space



[Images source: https://en.wikipedia.org/wiki/Support-vector_machine]

SVM soft-margin (popular pic./software → ...)

- Parameter C controls where you lie in the soft-hard margin spectrum.



[From: <https://stats.stackexchange.com/a/159051>]

- Software

sklearn.svm.SVC

```
class sklearn.svm.SVC(*, C=1.0, kernel='rbf', degree=3, gamma='scale', coef0=0.0, shrinking=True, probability=False, tol=0.001,  
cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', break_ties=False,  
random_state=None)¶
```

[\[source\]](#)

SVM soft-margin: ... → concrete understanding

Primal OP:

$$\min_{w,b,\zeta} \frac{1}{2} w^T w + C \sum_{i=1}^n \zeta_i$$

subject to $y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i$,
 $\zeta_i \geq 0, i = 1, \dots, n$

Dual OP:

$$\max_{\alpha} -\frac{1}{2} \alpha^T Q \alpha + \underline{c^T \alpha}$$

subject to $y^T \alpha = 0$

$0 \leq \alpha_i \leq C, i = 1, \dots, n$

Q is an n by n positive semidefinite matrix, $Q_{ij} \equiv y_i y_j K(x_i, x_j)$.

Prediction for new point x :

$$\sum_{i \in SV} y_i \alpha_i K(x_i, x) + b,$$

SVM aka *max-margin classifier* and is a type of *Sparse Kernel Machine (SKM)* method
(Relevance Vector Machine is another type of *SKM* method, specifically a probabilistic/Bayesian variant)

[Above formulas from sklearn help pages]

Recall: Inference and decision: three approaches for classification –
 SVM: *discriminant approach* initially motivated by *Computational Learning Theory*.

- Generative model approach:

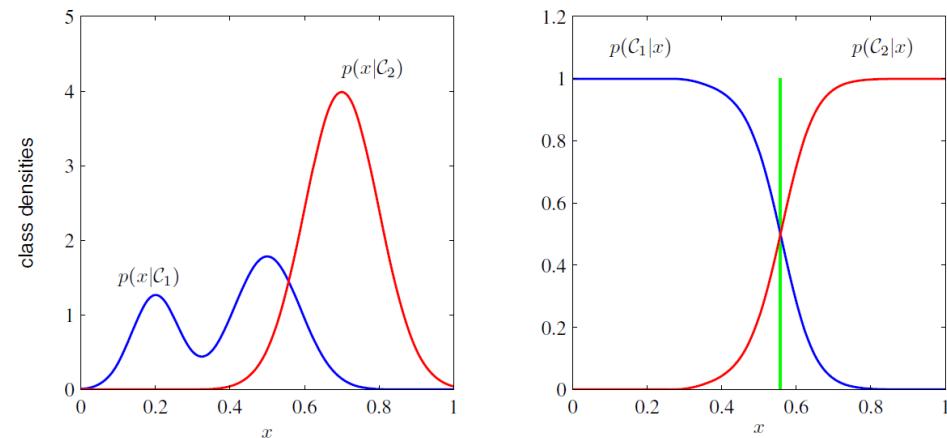
- (I) Model $p(x, C_k) = p(x|C_k)p(C_k)$
- (I) Use Bayes' theorem $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$
- (D) Apply optimal decision criteria

- Discriminative model approach:

- (I) Model $p(C_k|x)$ directly
- (D) Apply optimal decision criteria

- *Discriminant function approach*:

- (D) Learn a function that maps each x to a class label directly from training data
- Note: No posterior probabilities!

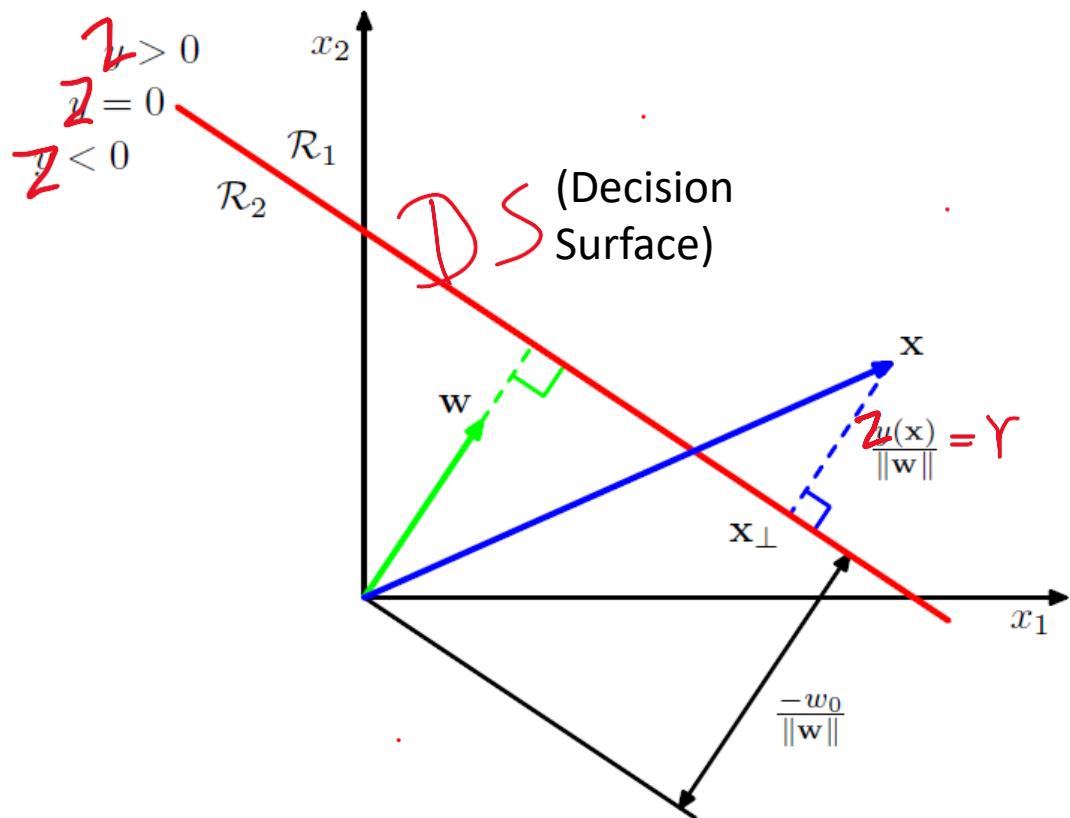


Recall: Discriminant

- Discriminant is a function that takes an input vector $x \in \mathbb{R}^d$ and assigns it to one of the K classes
 - we will assume $K=2$ henceforth for simplicity!
- We focus only on **linear discriminants** (i.e., those for which DB is a hyperplane wrt x (or $\phi(x)$)).
 - $z(x) = w^T x + w_0$ (or $w^T \phi(x)$)
 - DB: $z(x) = 0$ (hyperplane)
 - Prediction: $f(z(x)) = \text{sign}(z(x))$ (i.e., Predict C_1 if $z(x) \geq 0$, & C_2 if $z(x) < 0$)

Recall defn. of hyperplane $\{x \in \mathbb{R}^d: w^T x = b\}$, which is a $(d - 1)$ -dimensional (affine) subspace of a d -dim. vector space.

Recall: Geometry of decision surfaces: signed distance from decision surface



DB is $w^T x + w_0 = \text{const.}$, but const. can be absorbed into w_0 to get DB as $w^T x + w_0 = 0$.
 Let $z(w, x) = w^T x + w_0$ with the const. absorbed.

Signed dist. γ of x to DB:
 Express $x = x_{\perp} + \frac{\gamma w}{\|w\|}$

Mult. by w^T & add
 w_0 on both sides:

$$\Rightarrow z(x) = 0 + \gamma \|w\|$$

$$\Rightarrow \boxed{\gamma = \frac{z(x)}{\|w\|}}$$

Outline for Module M8

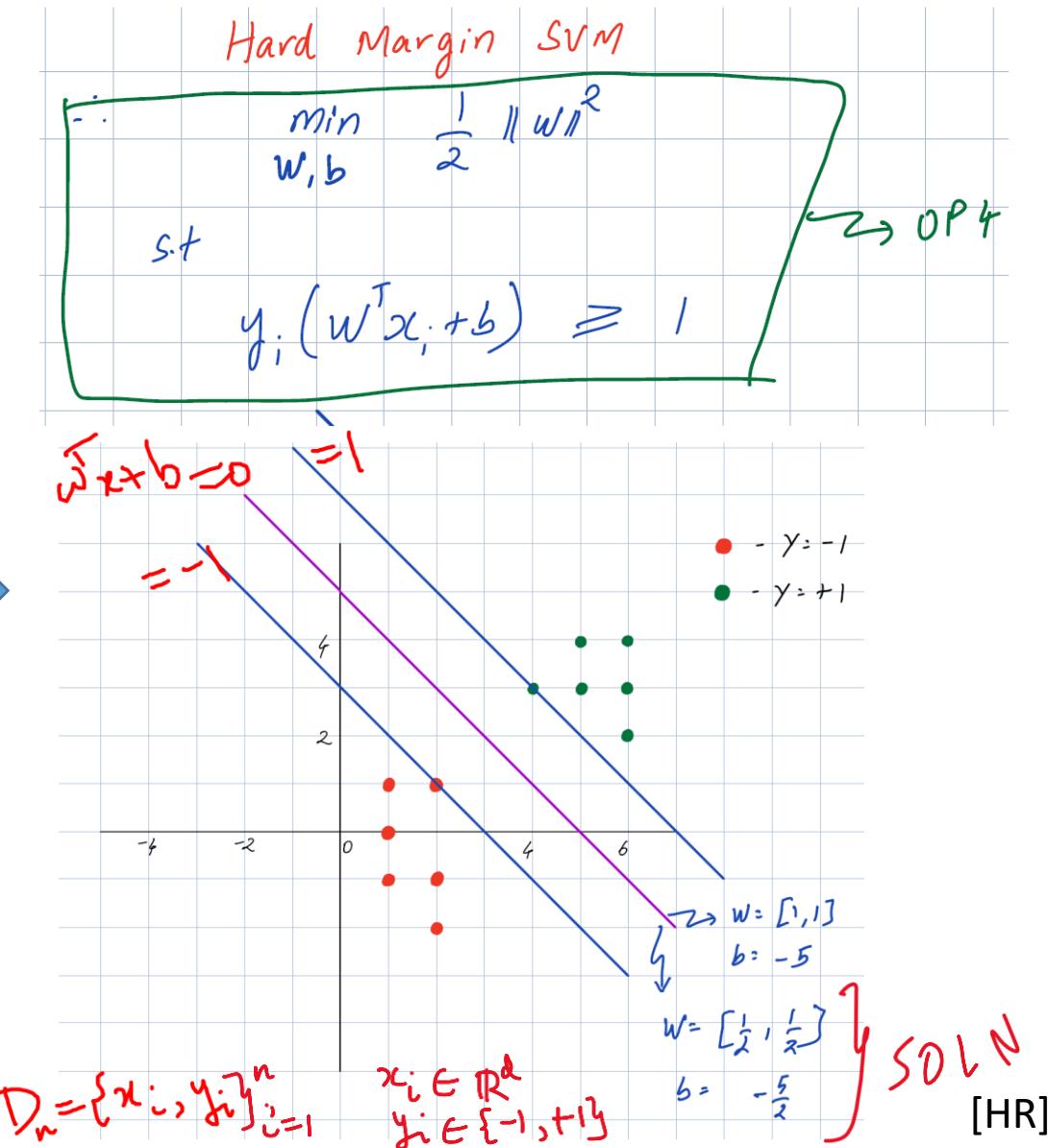
- M8. Classification (Support Vector Machines)
 - M8.0 Introduction/Motivation
 - **M8.1 SVM Problem Statement**
 - **(Hard/Soft-Margin SVM Problems)**
 - M8.2 SVM Solution
 - M8.3 SVM Interpretations
 - M8.4 Concluding thoughts

SVM hard-margin problem

$$\begin{aligned}
 & \text{OP1} \\
 & \max_{w, b} \min_i \frac{|w^T x_i + b|}{\|w\|} \geq \frac{(w^T x_i + b) y_i}{\|w\|} \\
 & \text{s.t.} \\
 & \quad \text{sign}(w^T x_i + b) = y_i
 \end{aligned}$$

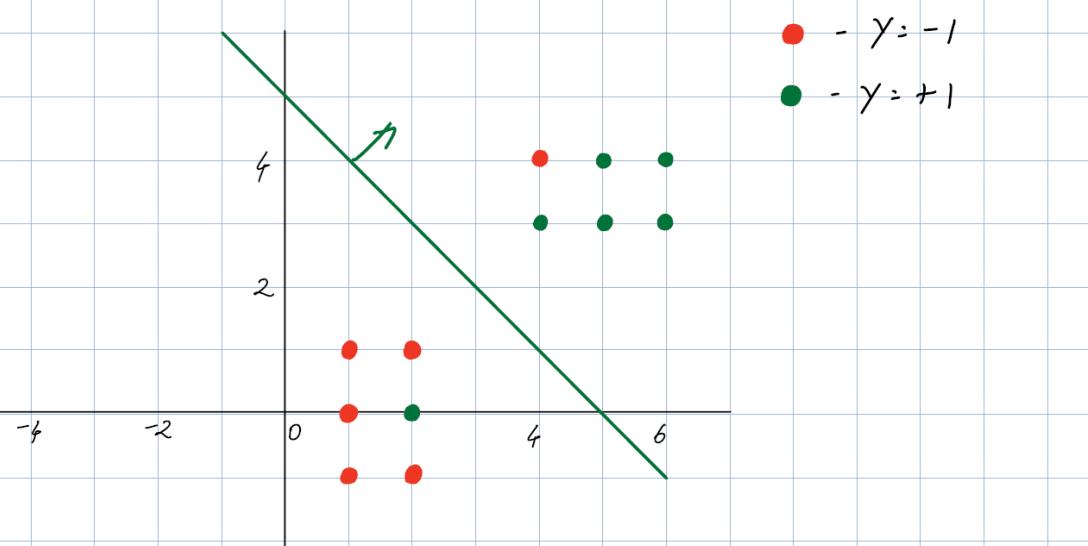
$$\begin{aligned}
 & \text{OP2} \\
 & \max_{w, b} \frac{1}{\|w\|} \\
 & \text{s.t.} \min_i (w^T x_i + b) y_i = 1 \\
 & \quad \text{OP3} \\
 & \max_{w, b} \frac{1}{\|w\|} \\
 & \text{s.t.} (w^T x_i + b) y_i \geq 1 \quad \forall i
 \end{aligned}$$

Exercise:
 S.T any soln. to OP3 will satisfy
 $\min_i (w^T x_i + b) y_i = 1$



SVM soft-margin problem

What if the data were not linearly separable?



(C-SVM) Soft Margin SVM

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

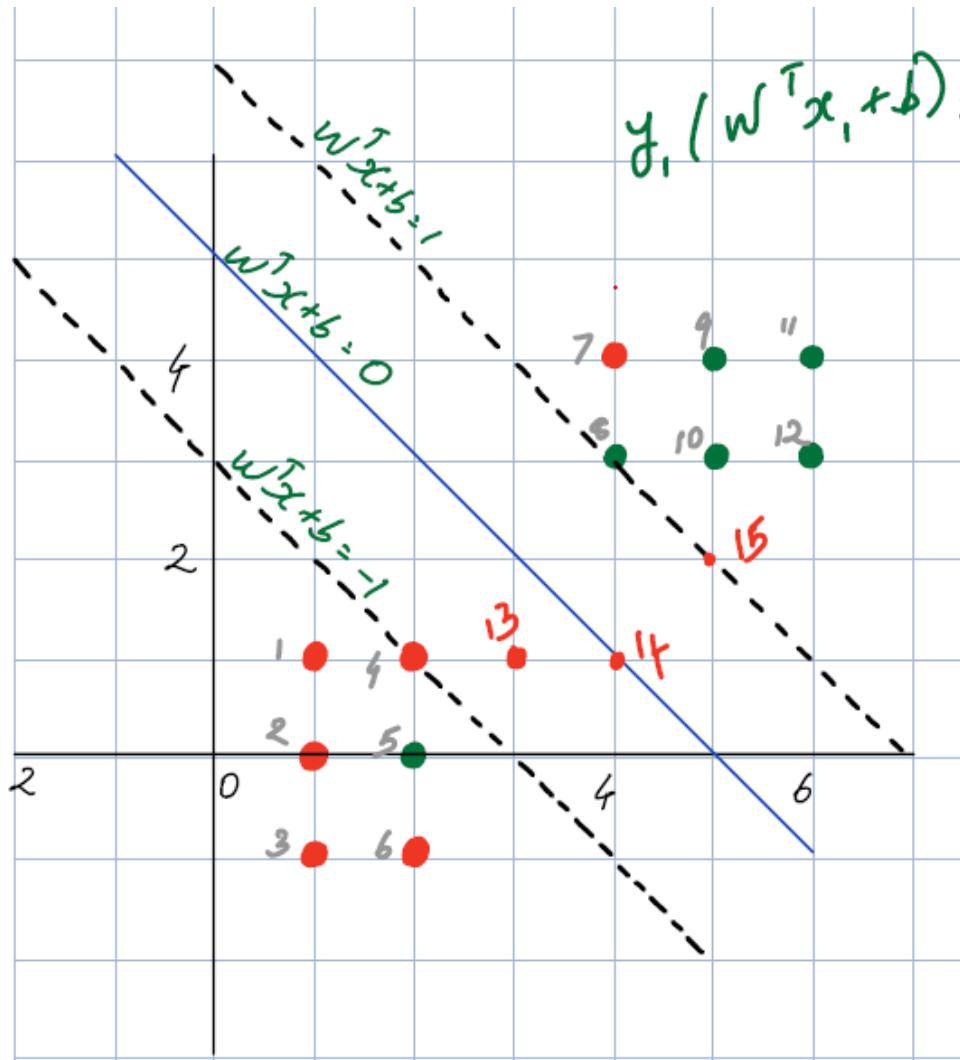
s.t.

$$y_i(w^T x_i + b) = 1 - \xi_i$$

$$\xi_i \geq 0$$

In Soft-margin SVM "any" w, b is feasible
Or $\forall w \in \mathbb{R}^n, b \in \mathbb{R}$, there exists $\xi \in \mathbb{R}_+^n$ s.t. (w, b, ξ) is feasible.

Example: what is ϵ_i for a given w, b ?



$$y_0(w^T x_0 + b) = \begin{pmatrix} y_2 & y_2 \end{pmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \frac{5}{2} = 1 - 1 = 0 \Rightarrow \epsilon_0 = 0$$

$$y_i(w^T x_i + b) \geq 1 - \epsilon_i$$

$$w = \left[\frac{1}{2}, \frac{1}{2} \right], \|w\| = \sqrt{2}$$

$$b = -\frac{5}{2}$$

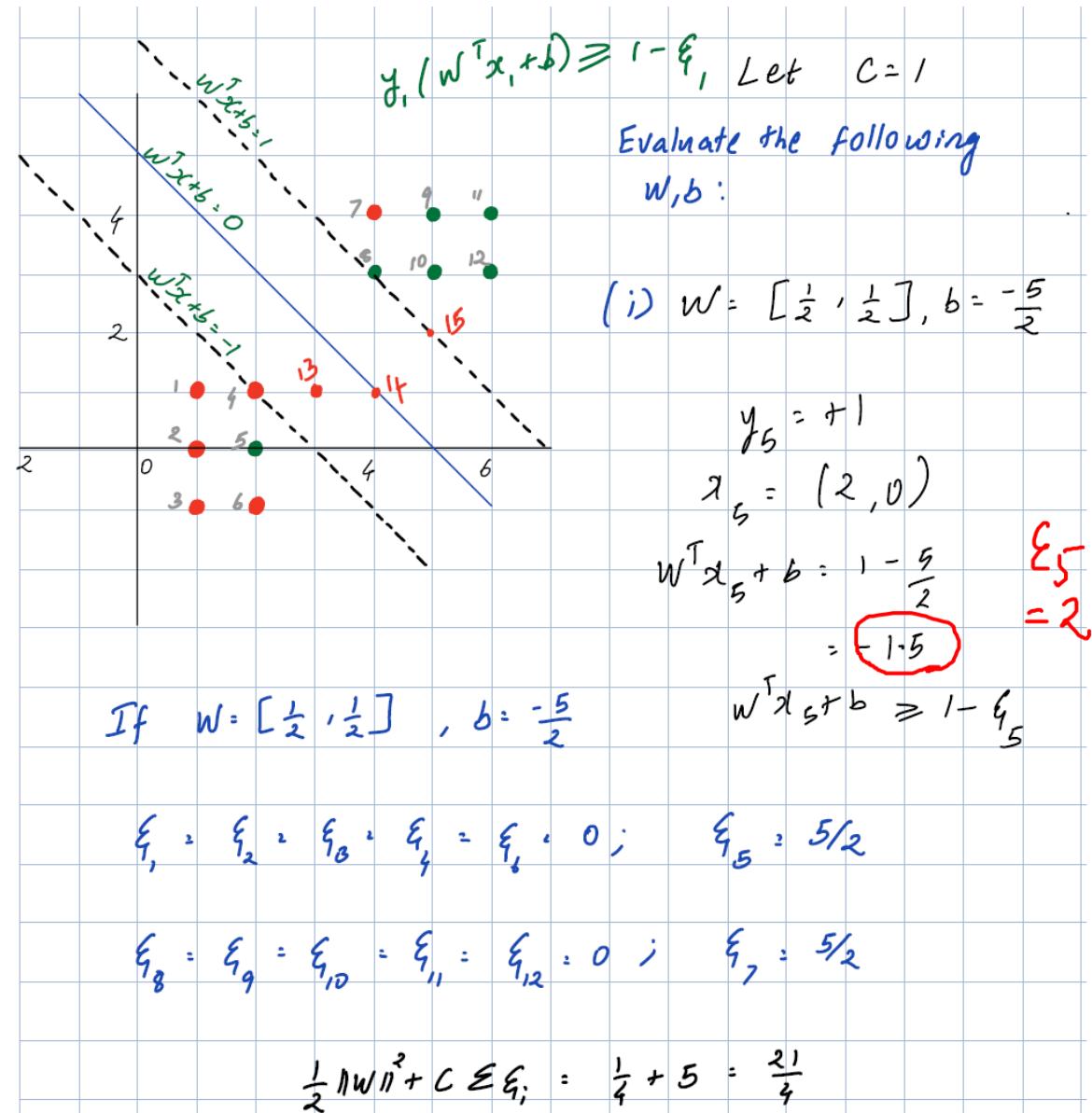
[HR]

Example: answer

(unnormalized
signed dist.)
 $y_i(w^T x_i + b)$

$$\xi_i = 1 - \text{prev. col.}$$

green x_5	-1.5	2.5
red x_{13}	+0.5	0.5
red x_{14}	+0	1
red x_{15}	-1	2
.	:	:



Outline for Module M8

- M8. Classification (Support Vector Machines)
 - M8.0 Introduction/Motivation
 - M8.1 SVM Problem Statement
 - **M8.2 SVM Solution**
 - **(Background: Constrained optimization - KKT & Primal-Dual)**
 - (SVM Dual Problem & Optimization algo. sketch)
 - M8.3 SVM Interpretations
 - M8.4 Concluding thoughts

From unconstrained to constrained opt. - FONC

Unconstrained Optimisation

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

First Order
Necessary
Conditions for
Optimality

$$\nabla f(x^*) = 0$$

FONC for x^* to be a local optima!

Constrained Optimisation: Lagrangian Theory

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

$$g_i(x^*), \nabla g_i(x^*)$$
$$h_i(x^*), \nabla h_i(x^*)$$

$$\text{s.t } g_1(x) \leq 0 ; h_1(x) = 0$$

$$g_n(x) \leq 0 ; h_m(x) = 0$$

WHAT ARE THE FIRST ORDER NECESSARY
CONDITIONS FOR OPTIMALITY?

FONC for x^* to be a local feasible (constrained) optima?

Recall: Linear approximation using gradient vector

Linear approximation around x_0 $L_{x_0}(f)(y)$

$$f(y) \approx f(x_0) + \nabla f(x_0)^T (y - x_0) \quad \text{---} ①$$

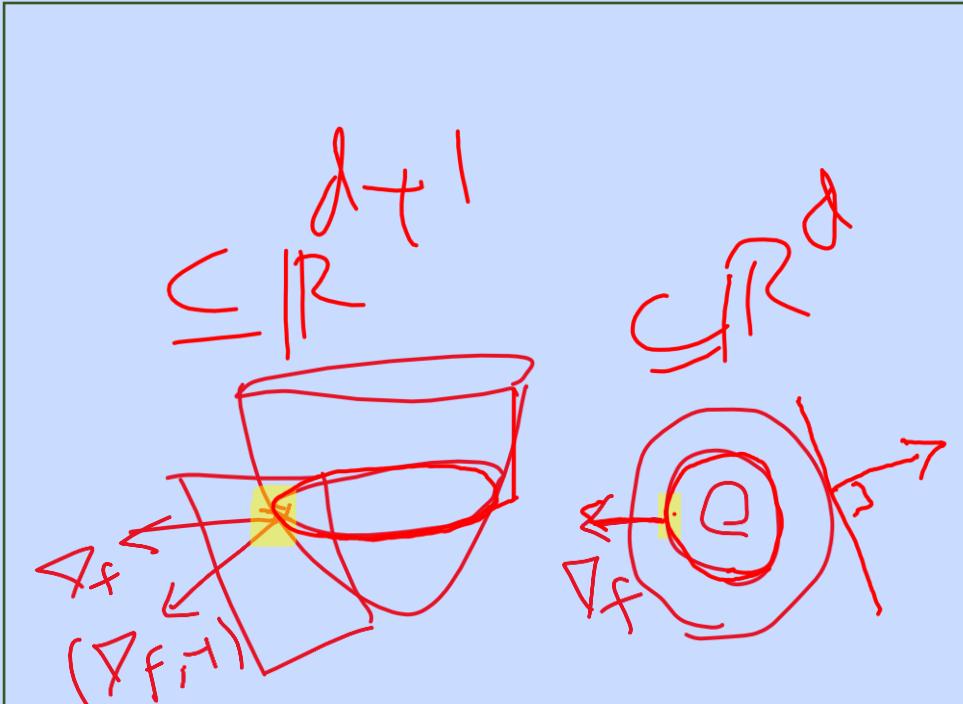
The "graph" of linear approximation of f around x_0 , is a "tangent plane" to the "graph" of f

[HR]

Facts:

At x_0 : $\nabla f \perp$ level set $\{x : f(x) = f(x_0)\}$ $\text{---} ②$

At $(x_0, f(x_0))$: $(\nabla f, -1) \perp$ tangent plane



Recall:

Spl case 1: One equality constraint:

$$\min f(x)$$

s.t

$$h(x) = 0$$

FONC for x^* to be a local minima:

(i) $h(x^*) = 0$

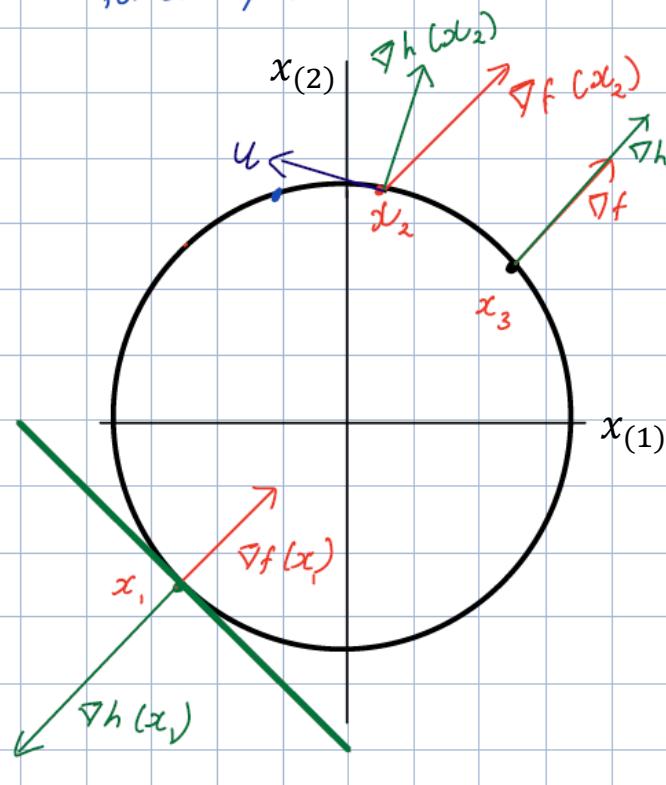
(ii) $\nabla f(x^*)$ is parallel to $\nabla h(x^*)$

(Equivalently $\nabla f(x^*) + \mu \nabla h(x^*) = 0$)

For some $\mu \in \mathbb{R}$

e.g

- $\min x_{(1)} + x_{(2)}$
- s.t.
- $x_{(1)}^2 + x_{(2)}^2 - 1 = 0$



Proof Sketch: Clearly x^* must satisfy $h(x^*) = 0$

$F = \text{Set of feasible directions} : \{u : h(x^* + \alpha u) = 0 \text{ for small } \alpha \}$

$$= \{u : \nabla h(x^*)^T u = 0\}$$

$$= \{u : h(u^* + \alpha u) = h(x^*) \text{ for some } \alpha > 0\}$$

$$\text{Now, } h(x^* + \alpha u) - h(x^*) \approx \nabla h(x^*)^T (\alpha u)$$

$D = \text{Set of Descent directions} : \{u : f(x^* + \alpha u) < f(x^*) \text{ for small } \alpha\}$

$$= \{u : \nabla f(x^*)^T u < 0\}$$

When will $F \cap D$ be empty?

Only when $\nabla f(x^*) = \mu \nabla h(x^*)$ for some $\mu \in \mathbb{R}$.

→ This requires a

Proof. (regularity
conditions for KKT
to hold)

Special Case 2 : One Inequality constraint.

$$\min f(x)$$

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

$$D = \{u : \nabla f(x^*)^T u < 0\}$$

$$F = \{u : \nabla g(x^*)^T u \leq 0\}$$

$F \cap D$ is empty

This requires
a proof.

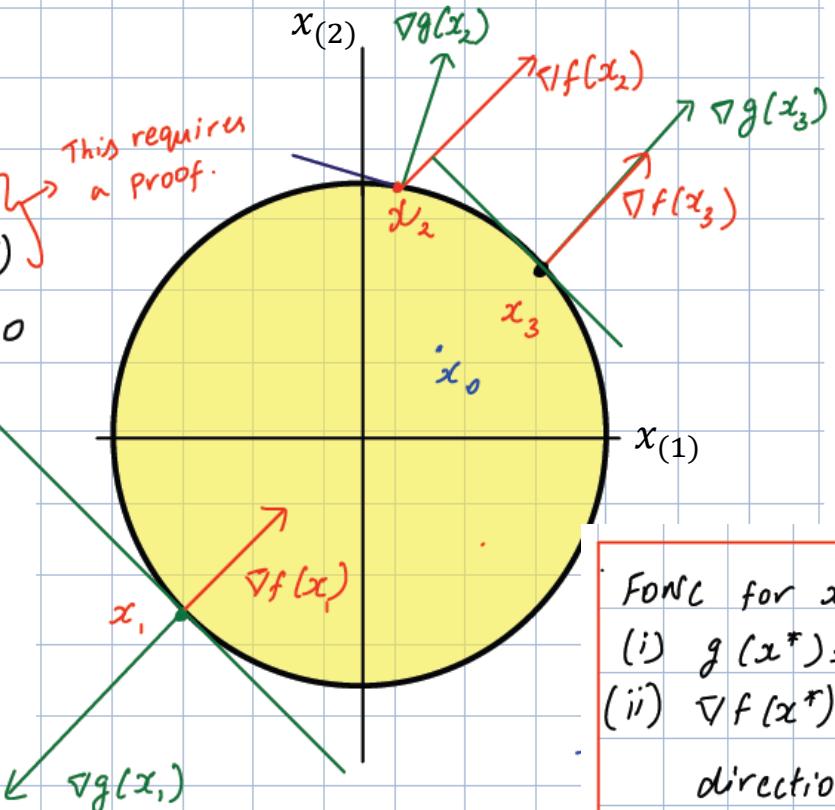
only if $\nabla f(x^*) = -\lambda \nabla g(x^*)$
for some $\lambda \geq 0$

e.g.

$$\min x_{(1)} + x_{(2)}$$

s.t.

$$x_{(1)}^2 + x_{(2)}^2 - 1 \leq 0$$



FOFC for x^* to be a local min

(i) $g(x^*) \leq 0$

(ii) $\nabla f(x^*)$ is parallel to $\nabla g(x^*)$ and pointing in opposite direction. (Equivalently $\nabla f(x^*) + \lambda \nabla g(x^*) = 0$, for some $\lambda \geq 0$)

(iii) If $g(x^*) < 0$ then $\nabla f(x^*) = 0$. (constraint g is inactive)
 $\lambda = 0$

(Equivalently $\lambda g(x^*) = 0$)

KKT conditions (FONC) – General Case

General case : Multiple equality & inequality constraints

$$\min f(x)$$

s.t

$$g_1(x) \leq 0 ; h_1(x) = 0$$

$$g_2(x) \leq 0 ; h_2(x) = 0$$

:

:

$$g_m(x) \leq 0 ; h_n(x) = 0$$

Define Lagrangian

$$L(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{i=1}^n \mu_i h_i(x)$$

$$\lambda_i \geq 0$$

Equivalently : KKT conditions.

FONC for x^* to be a local minima.

(i) $g_i(x^*) \leq 0 \quad \forall i \in [m]$ (Feasibility)
 $h_i(x^*) = 0 \quad \forall i \in [n]$

(ii) $\exists \lambda \in \mathbb{R}_+^m, \exists \mu \in \mathbb{R}^n$ such that (Stationarity)

$$(a) \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) + \sum_{i=1}^n \mu_i \nabla h_i(x^*) = 0$$

$$(b) \lambda_i = 0 \text{ if } g_i(x^*) < 0, \text{ or equivalently}$$
$$\lambda_i \cdot g_i(x^*) = 0, \forall i \in [m] \Rightarrow (\text{complementary slack})$$

$$(\text{or equivalently}) \lambda_i = 0 \quad \forall i \notin A(x^*)$$

$$A(x^*) = \{i \in [m] : g_i(x^*) = 0\} \text{ set of active ineq.}$$

KKT Conditions - Example

$$\begin{array}{ll} \min_{x \in \mathbb{R}} & x^2 + ax \\ \text{s.t.} & 0 \leq x \leq 1 \end{array} \quad \left(\begin{array}{l} \text{For unconstrained} \\ \text{problem minimiser is} \\ x = -\frac{a}{2} \end{array} \right)$$

$$\mathcal{L}(x, \lambda_1, \lambda_2) = x^2 + ax + \lambda_1(-x) + \lambda_2(x-1)$$

$$\nabla_x \mathcal{L}(x, \lambda_1, \lambda_2) = 0 \Rightarrow 2x + a - \lambda_1 + \lambda_2 = 0 \quad (\text{stationarity})$$

$$x = \frac{\lambda_1 - \lambda_2 - a}{2}$$

KKT Conditions – Example (checking FONC)

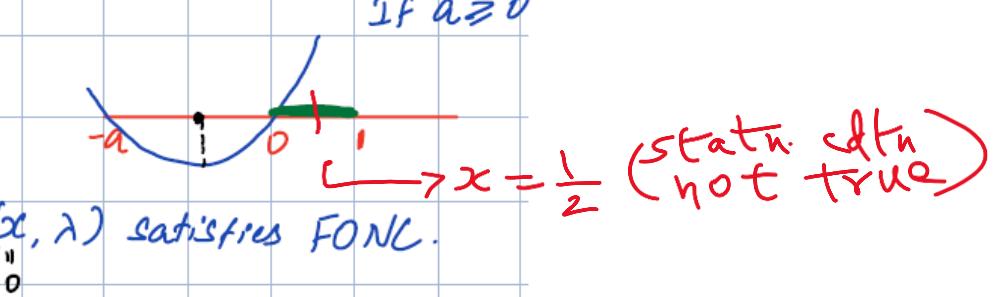
Case 1 Check if " $x=0$ " satisfies FONC

$$\Rightarrow \lambda_1 \geq 0, \lambda_2 = 0 \quad (\text{complementary slack})$$

$$x = \frac{\lambda_1 - a}{2} = 0 \quad \text{IF } a \geq 0$$

$$\Rightarrow \lambda_1 = a$$

\therefore if $a \geq 0$, $\exists \lambda \geq 0$ s.t. (x, λ) satisfies FONC.



Case 2 Check if " $x=1$ " satisfies FONC

$$\lambda_1 = 0, \lambda_2 \geq 0 \quad (\text{complementary slack})$$

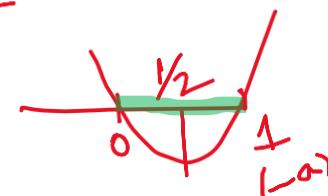
$$x = \frac{-\lambda_2 - a}{2} = 1 \Rightarrow \lambda_2 = -2 - a$$

\therefore if $a \leq -2$, $\exists \lambda \geq 0$ s.t. (x, λ) satisfies FONC.



Case 3: $x = \frac{1}{2}$

Satisfies FONC
if $a = -1$. Check!



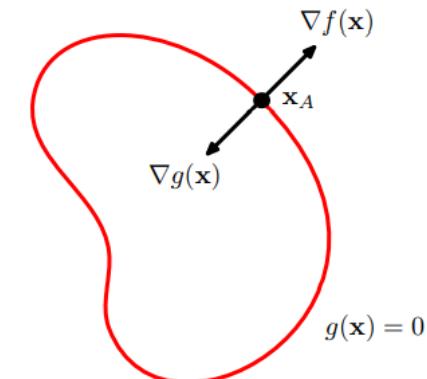
Optional: Bishop-AppxE is also a good read to get similar intuition about Lagrange multipliers and FONC.

$$L(\mathbf{x}, \lambda) \equiv f(\mathbf{x}) + \lambda g(\mathbf{x}).$$

$$\max f(\mathbf{x}) \\ \text{st} \\ g(\mathbf{x}) = 0$$

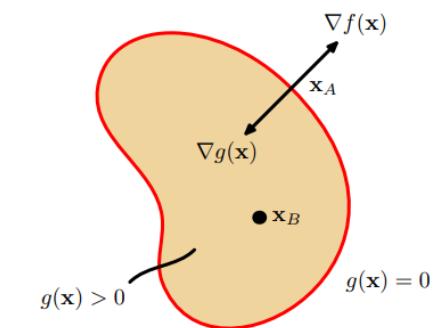
FONC

$$\exists \lambda \in \mathbb{R} \text{ st} \\ \nabla_{\mathbf{x}, \lambda} L = 0 \quad \left(\begin{array}{l} \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \lambda \nabla_{\mathbf{x}} g(\mathbf{x}^*) = 0 \\ g(\mathbf{x}^*) = 0 \end{array} \right)$$



$$\max f(\mathbf{x}) \\ \text{st} \\ g(\mathbf{x}) \geq 0$$

$$\exists \lambda \geq 0 \text{ st. } \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \lambda \nabla_{\mathbf{x}} g(\mathbf{x}^*) = 0 \\ g(\mathbf{x}^*) \geq 0 \\ \lambda g(\mathbf{x}^*) = 0 \quad (\text{active & inactive constraint})$$

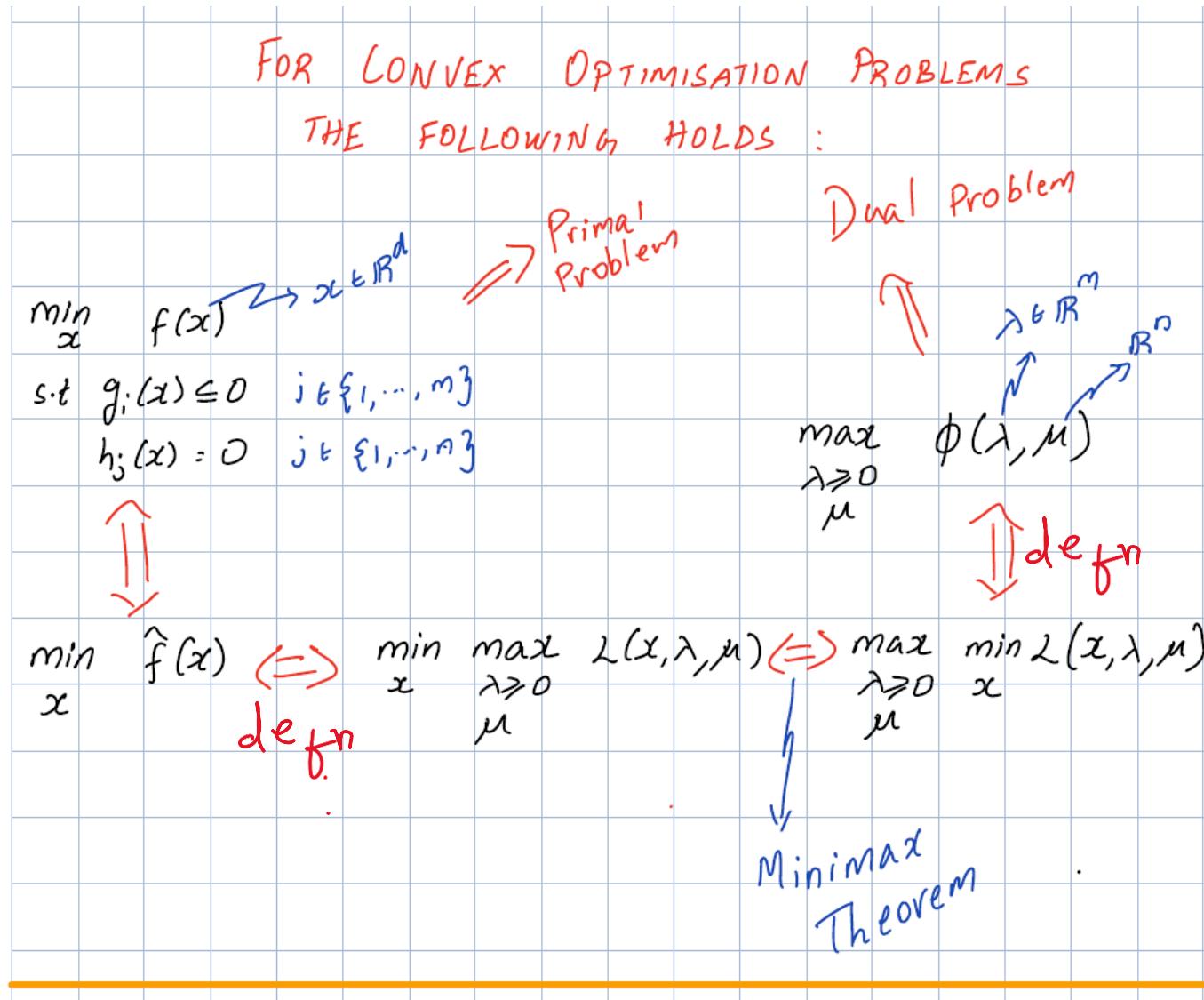


Exercises: Other constrained optimization problems!

- Prove that the KL-divergence $KL(p \parallel q)$ is minimized when $q = p$.
- Prove that entropy $H(p)$ is maximized when p is uniform.
- What is the distance of a point $u \in \mathbb{R}^d$ to the closest point v in a hyperplane given by $\{x \in \mathbb{R}^d : w^T x + b = 0\}$?

Having established KKT cdnts., can we actually (constructively) find the KKT multipliers λ^*, μ^* that satisfy the KKT cdtns.?

Primal-Dual Relation (via the Minimax Theorem)



Define Lagrangian

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^n \mu_j h_j(x)$$

Function $\hat{f}(.)$

- Desired function:

$$\therefore \hat{f}(x) = \begin{cases} f(x) & \text{if } x \text{ is feasible} \\ \infty & \text{if } x \text{ is not feasible.} \end{cases}$$

- First attempt:

$$\hat{f}(x) = f(x) + \sum_i \infty \cdot \mathbb{1}_{g_i(x) > 0} + \sum_j \infty \cdot \frac{1}{h_j(x) \neq 0}$$

- Second attempt:

Define function $\hat{f}: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ as

$$\hat{f}(x) = \max_{\substack{\lambda \geq 0 \\ \mu}} L(x, \lambda, \mu)$$

Define Lagrangian

$$L(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^m \mu_j h_j(x)$$

Weak Duality :

$$\min_x \max_{\substack{\lambda \geq 0 \\ \mu}} \mathcal{L}(x, \lambda, \mu) \geq \max_{\lambda \geq 0} \min_x \mathcal{L}(x, \lambda, \mu)$$

If f, g are convex, h is Linear then

$$\min_x \max_{\substack{\lambda \geq 0 \\ \mu}} \mathcal{L}(x, \lambda, \mu) = \max_{\lambda \geq 0} \min_x \mathcal{L}(x, \lambda, \mu)$$

MINIMAX THEOREM !!

Key Result in Game Theory !!

Define Lagrangian $\lambda_i \geq 0$

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{i=1}^m \mu_i h_i(x)$$

How to get one solution from the other?

PRIMAL

$$\min_x f(x)$$

$$\text{s.t. } g_i(x) \leq 0; i \in [m] \quad (=)$$

$$h_j(x) = 0; j \in [n]$$

Let x^* be the
soln to the primal

DUAL

$$\max_{\lambda \geq 0} \phi(\lambda, \mu)$$

$$\mu$$

Let λ^*, μ^* be the
soln. to the dual.

Fact : (Follows from Proof of min-max Theorem)

x^* satisfies KKT for Primal with
the multipliers in the KKT condition
given by λ^* & μ^* !!

Fact



Define $\phi(\lambda, \mu) = \min_x L(x, \lambda, \mu)$

How to get one solution from the other?

More precisely let x^* be a primal optimal soln.
let λ^*, μ^* be a dual optimal soln.

The above statement is equivalent to the
below statement. (for convex opt.)

i) $g_i(x^*) \leq 0 ; \lambda_i^* \geq 0 \quad \forall i \in [m]$
 $h_j(x^*) = 0 \quad \forall j \in [n]$

Fact ~~(X)~~

(ii) $\nabla f(x^*) + \sum_{i=1}^m \lambda_i^* \nabla g_i(x^*) + \sum_{j=1}^n \mu_j^* \nabla h_j(x^*) = 0$

solve
these
eqns.

(iii) $\lambda_i^* g_i(x^*) = 0 \quad \forall i \in [m]$

Duality Example: A linear program

$$\begin{array}{ll}
 \min_{x \in \mathbb{R}^2} & x_1 + 2x_2 \\
 \text{s.t.} & \\
 & x_1 + x_2 = 3 \Leftrightarrow \\
 & x_1 \geq 0 \\
 & x_2 \geq 0
 \end{array}$$

$\min_x \quad x_1 + 2x_2$
 s.t.
 $x_1 + x_2 - 3 = 0$
 $-x_1 \leq 0$
 $-x_2 \leq 0$



(standard form
 Primal)

$$\begin{aligned}
 L(x, \lambda_1, \lambda_2, \mu) &= x_1 + 2x_2 + \lambda_1(-x_1) + \lambda_2(-x_2) + \mu(x_1 + x_2 - 3) \\
 &= x_1(1 - \lambda_1 + \mu) + x_2(2 - \lambda_2 + \mu) - 3\mu
 \end{aligned}$$

$$\phi(\lambda, \mu) = \min_{x \in \mathbb{R}^2} L(x, \lambda, \mu)$$

$$= \begin{cases} -\infty & \text{if } 1 - \lambda_1 + \mu \neq 0 \\ -\infty & \text{if } 2 - \lambda_2 + \mu \neq 0 \\ -3\mu & \text{if } 1 - \lambda_1 + \mu = 0 \text{ and} \\ & 2 - \lambda_2 + \mu = 0 \end{cases}$$

Duality Example (contd.)

$$\begin{aligned} \text{max } & \phi(\lambda, \mu) \\ \text{s.t. } & \lambda_1 \geq 0 \\ & \lambda_2 \geq 0 \\ & 1 - \lambda_1 + \mu = 0 \\ & 2 - \lambda_2 + \mu = 0 \end{aligned}$$

↑↑

$$\begin{aligned} \text{max } & -3\mu \\ \text{s.t. } & \lambda_1 \geq 0 \\ & \lambda_2 \geq 0 \\ & 1 - \lambda_1 + \mu = 0 \\ & 2 - \lambda_2 + \mu = 0 \end{aligned}$$

↑↑

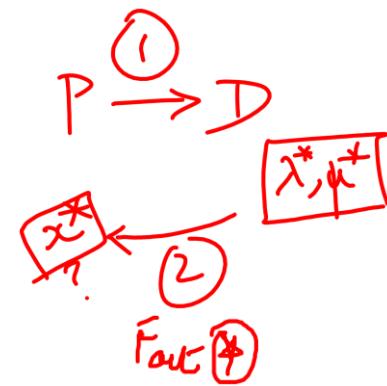
$$\begin{aligned} \text{max } & -3\mu \\ \text{s.t. } & \lambda_1: 1 + \mu \geq 0 \\ & \lambda_2: 2 + \mu \geq 0 \end{aligned}$$

↑↑

$$\begin{aligned} \text{max } & -3\mu \\ & \mu \geq -1; \mu \geq -2 \end{aligned}$$

↑↑

$$\begin{aligned} \text{max } & -3\mu \\ & \mu \geq -1 \end{aligned}$$



Duality Example (contd.)

$$\begin{array}{ll}
 \text{min}_{x_1, x_2} & x_1 + 2x_2 \\
 \text{s.t.} & x_1 + x_2 = 3 \\
 & -x_1 \leq 0 \\
 & -x_2 \leq 0
 \end{array}
 \quad
 \begin{array}{ll}
 \max_{\mu, \lambda_1, \lambda_2} & -3\mu \\
 \text{s.t.} & \mu \geq -1 \\
 & \lambda_1 = 1 + \mu \\
 & \lambda_2 = 2 + \mu
 \end{array}
 \quad (\text{Dual})$$

The Dual solution is clearly

$$\begin{aligned}
 \mu^* &= -1 \\
 \lambda_1^* &= 0 \\
 \lambda_2^* &= 1
 \end{aligned}$$

Now try to find a feasible x^* , satisfying the stationarity & complementary slack conditions with these as the Lagrange multipliers:

$$\begin{aligned}
 \text{(i) Feasibility} \Rightarrow & x_1^* + x_2^* = 3 \quad \rightarrow \textcircled{1} \\
 & -x_1^* \leq 0 \\
 & -x_2^* \leq 0
 \end{aligned}$$

(ii) Stationarity :

$$\nabla f(x^*) = \begin{bmatrix} 1 \\ 2 \end{bmatrix}; \nabla g_1(x^*) = \begin{bmatrix} -1 \\ 0 \end{bmatrix}; \nabla g_2(x^*) = \begin{bmatrix} 0 \\ -1 \end{bmatrix}; \nabla h_1(x^*) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\begin{aligned}
 \nabla f(x^*) + \lambda_1^* \nabla g_1(x^*) + \lambda_2^* \nabla g_2(x^*) + \mu^* \nabla h_1(x^*) &= 0 \\
 \Rightarrow \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 0 \begin{bmatrix} -1 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} 0 \\ -1 \end{bmatrix} + (-1) \begin{bmatrix} 1 \\ 1 \end{bmatrix} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix}
 \end{aligned}$$

[This holds for all x , and hence useless for finding x^* .

(iii) Complementary Slack :

$$\begin{aligned}
 \lambda_1^* > 0 &\Rightarrow g_1(x^*) = 0 \\
 \lambda_2^* > 0 &\Rightarrow g_2(x^*) = 0
 \end{aligned}$$

$$\begin{aligned}
 \text{only } \lambda_2^* > 0 \quad \therefore g_2(x^*) &= 0 \\
 \Rightarrow -x_2^* &= 0 \quad \rightarrow \textcircled{2}
 \end{aligned}$$

Solving \textcircled{1} and \textcircled{2} we get $x^* = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$. ■

Outline for Module M8

- M8. Classification (Support Vector Machines)
 - M8.0 Introduction/Motivation
 - M8.1 SVM Problem Statement
 - **M8.2 SVM Solution**
 - (Background: Constrained optimization - KKT & Primal-Dual)
 - **(SVM Dual Problem & Optimization algo. sketch)**
 - M8.3 SVM Interpretations
 - M8.4 Concluding thoughts

SVM: From Primal → Dual

Primal

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

s.t

$$y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \in [n] \quad \alpha_i$$

$$\xi_i \geq 0 \quad \forall i \in [n] \quad \beta_i$$

Lagrangian

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(w^T x_i + b)) \\ + \sum_{i=1}^n \beta_i (-\xi_i)$$

$$= \frac{1}{2} \|w\|^2 + b \left(\sum_{i=1}^n -\alpha_i y_i \right) + \sum_{i=1}^n \xi_i (C - \alpha_i - \beta_i) \\ + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (w^T x_i)$$

Dual function

$$\phi(\alpha, \beta) = \min_{w, b, \xi} \mathcal{L}(w, b, \xi, \alpha, \beta)$$

$$\sum_{i=1}^n -\alpha_i y_i = 0 \rightarrow \textcircled{1} \quad (\text{If not } \phi(\alpha, \beta) = -\infty) \\ C - \alpha_i - \beta_i = 0 \rightarrow \textcircled{2}$$

$$g_i(\cdot) \leq 0$$

If Equations ① and ② are satisfied then:

$$\phi(\alpha, \beta) = \min_w \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i w^T x_i + \sum_{i=1}^n \alpha_i$$

SVM: From Primal \rightarrow Dual

$$\phi(\alpha, \beta) = \min_w \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i w^T x_i + \sum_{i=1}^n \alpha_i$$

Let $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}$, $\alpha = [\alpha_1, \dots, \alpha_n]^T$

Let $Y = \begin{bmatrix} y_1 & 0 \\ 0 & y_2 & \dots \\ 0 & \ddots & y_n \end{bmatrix} \in \mathbb{R}^{n \times n}$ be a diagonal matrix.

$$\phi(\alpha, \beta) = \min_w \frac{1}{2} \|w\|^2 - \alpha^T Y X w + \frac{1}{2} \alpha^T \alpha$$

Taking gradient w.r.t w and set to 0

$$\Rightarrow w - X^T Y \alpha = 0$$

$$\therefore \boxed{w = X^T Y \alpha}$$

$$\therefore \phi(\alpha, \beta) = \frac{1}{2} \|X^T Y \alpha\|^2 - \alpha^T Y X X^T Y \alpha + \frac{1}{2} \alpha^T \alpha;$$

$$\therefore -\frac{1}{2} \alpha^T Y X X^T Y \alpha + \frac{1}{2} \alpha^T \alpha;$$

SVM Dual Problem

Dual Problem:

$$\max_{\alpha} \phi(\alpha, \beta)$$
$$\alpha \geq 0$$
$$\beta \geq 0$$
$$\therefore \max_{\alpha, \beta} -\frac{1}{2} \alpha^T Y X X^T Y \alpha + \beta \alpha;$$

s.t.

$$\sum_{i=1}^n y_i \alpha_i = 0$$

These two constraints are to ensure $\phi(\alpha, \beta)$ is not $-\infty$

$$\alpha_i \geq 0 \quad \forall i \in [n]$$
$$\beta_i \geq 0 \quad \forall i \in [n]$$



Equivalently : (setting $\beta_i = C - \alpha_i$)

$$\max_{\alpha} -\frac{1}{2} \alpha^T Y X X^T Y \alpha + \beta \alpha;$$

s.t.

$$\sum_{i=1}^n \alpha_i y_i = 0$$
$$0 \leq \alpha_i \leq C \quad \forall i \in [n]$$

The above problem is the Dual SVM problem:

Stop & Think: What have we achieved so far?

$$\begin{array}{ll} \min_{w,b,\xi} & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{st} & \begin{aligned} & R_d \leq w^T x_i + b \leq R_u \\ & \xi_i \geq 0 \end{aligned} \end{array} \Leftrightarrow \begin{array}{ll} \max_{\alpha \in \mathbb{R}^n} & -\frac{1}{2} \alpha^T Y X X^T Y \alpha + 1^T \alpha \\ \text{s.t.} & \begin{aligned} & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned} \end{array}$$

Where $Y = \begin{bmatrix} y_1 & 0 \\ y_2 & \vdots \\ 0 & y_n \end{bmatrix} \in \mathbb{R}^{n \times n}$; $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}$

(To get the same for hard margin, just set $C = \infty$: Exercise: Why is this?)

Next steps:

- This Dual problem can be maximized using SMO or PGD methods much easier than the Primal problem!
- Then convert Dual to Primal solution using ~~Fact~~.

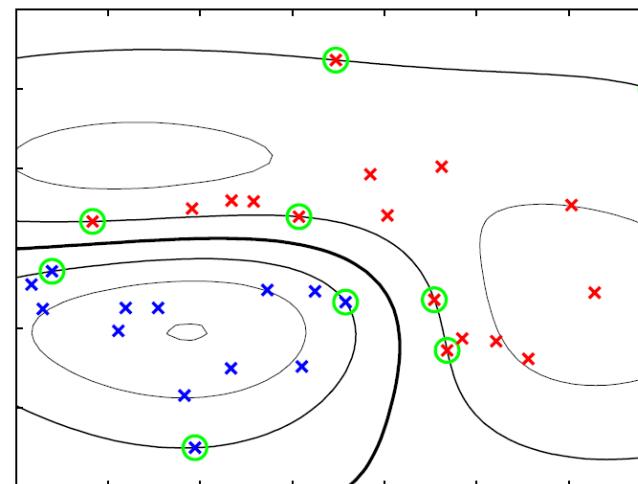
Brief aside: Kernel-ize our Dual Problem!

Let $K \in \mathbb{R}^{mn}$ s.t
 $K = \Phi \Phi^T$

$$\therefore \text{Dual SVM} \Rightarrow \max_{\alpha \in \mathbb{R}^n} -\frac{1}{2} \alpha^T K \alpha + \alpha^T \alpha \\ \text{s.t. } \sum_i \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C$$

Identity kernel gives back the original Dual Problem!

Recall: Why kernel-ize?



Brief aside: Kernel SVM – Prediction for a new data point (using support vectors)

Need to go from one solution (dual: α^*) to the other (primal: w^*, b^*, ϵ^*).

Fact ~~(X)~~

We have that the primal soln

$$w^* = X^T Y \alpha^* \Rightarrow w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

$$\therefore w^{*T} x = \alpha^{*T} Y X x$$

$$= \alpha^{*T} Y \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} x$$

$$= \alpha^{*T} Y \begin{bmatrix} x_1^T x \\ \vdots \\ x_n^T x \end{bmatrix}$$

$$= \alpha^{*T} Y \begin{bmatrix} K(x, x_1) \\ \vdots \\ K(x, x_n) \end{bmatrix} \quad (\text{If we use feature vectors } \phi(x) \text{ instead})$$

$$\therefore w^{*T} \phi(x) = \sum_{i=1}^n \alpha_i^* y_i K(x_i, x)$$

Lagrangian

$$\mathcal{L}(w, b, \xi, \alpha, B) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i (w^T x_i + b))$$

$\downarrow \downarrow \downarrow \downarrow \downarrow$
 $R^d R R^n R^n R^n$

$$+ \sum_{i=1}^n B_i (-\xi_i)$$

[HR]

Having found w^* from α^* using KKT (stationarity);
now find b^* using also KKT complementary slack

b^* can be found by the following :

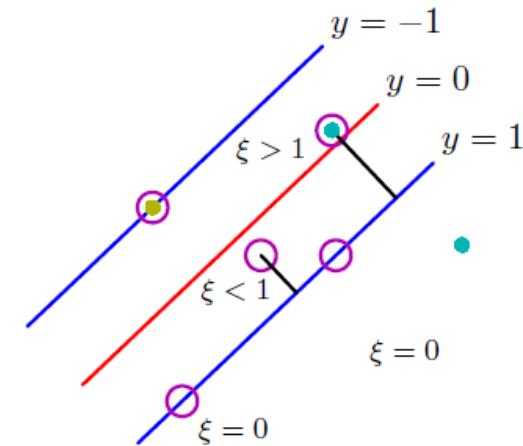
If $\alpha_i^* > 0$ and $\alpha_i^* < C$, then the
Point x_i satisfies

$$y_i (w^{*T} x_i + b^*) = 1 \quad (\text{Exercise: why?})$$

$$\therefore y_i (\alpha_i^{*T} y_i x_i + b^*) = 1 \quad (\text{Hint: What happens to } \alpha_i \text{ & } \beta_i \text{ and use complementary slack})$$

$$\therefore y_i \left(\sum_{j=1}^n \alpha_j^* y_j K(x_j, x_i) + b^* \right) = 1$$

$$\therefore b^* = y_i - \sum_{j=1}^n \alpha_j^* y_j K(x_j, x_i)$$



How do we optimize the Dual Problem?

How to solve the Dual Problem?

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & -\frac{1}{2} \alpha^T Y K Y \alpha + \mathbf{1}^T \alpha \\ \text{s.t.} \quad & \sum \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

The objective is quadratic. The constraints are an intersection of a box and a hyperplane. Can be solved using projected gradient descent- or SMO (sequential minimal optimisation)

Outline for Module M8

- M8. Classification (Support Vector Machines)
 - M8.0 Introduction/Motivation
 - M8.1 SVM Problem Statement
 - M8.2 SVM Solution
 - **M8.3 SVM Interpretations**
 - **(Support vectors, Kernels, Loss function view)**
 - M8.4 Concluding thoughts

Summary so far, and support vectors/kernels

Primal

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0 \quad \forall i \in [n] \end{aligned}$$



$$\begin{aligned} \therefore \text{Dual SVM} \Rightarrow \max_{\alpha \in \mathbb{R}^n} \quad & -\frac{1}{2} \alpha^T Y K Y \alpha + \alpha^T \alpha \\ \text{s.t.} \quad & \sum \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

$$\begin{aligned} w^{*T} \phi(x) = \quad & \sum_{i=1}^n \alpha_i^* y_i K(x_i, x) \\ b^* = y_i - \sum_{j=1}^n \alpha_j^* y_j K(x_j, x_i) \end{aligned}$$



Optimize using PGD/SMO to find α^*

Support Vectors:

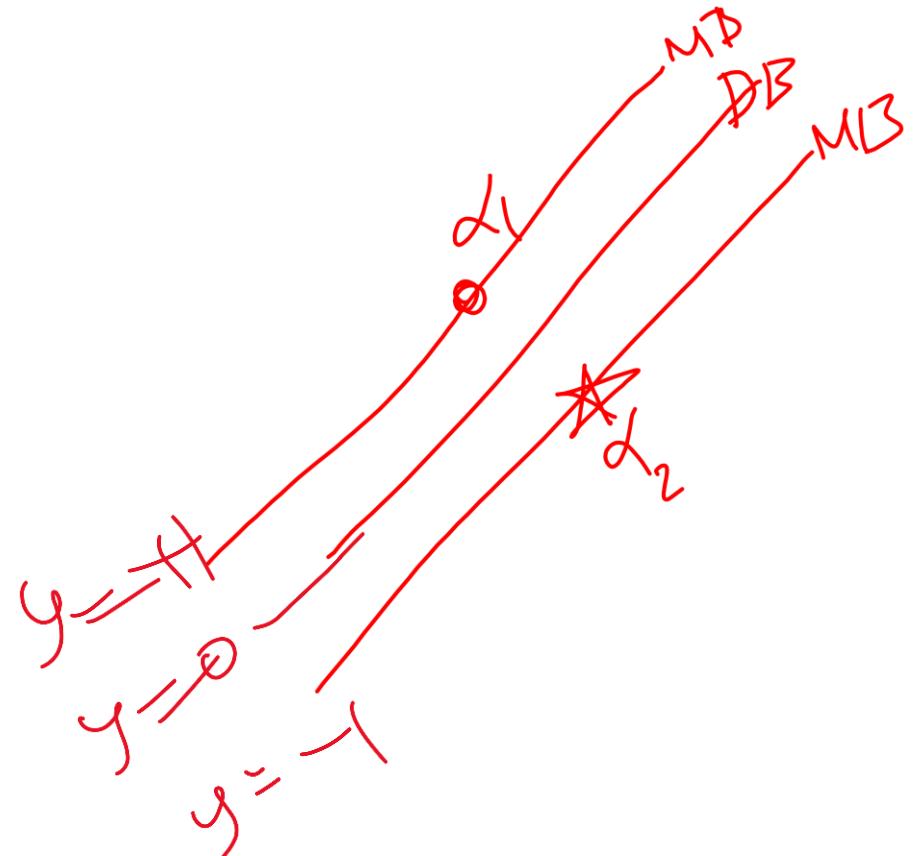
- Training data points for which $\alpha_i^* \neq 0$.
- Typically sparse. Why?

Exercise 0: What are the possible values of $y_i(w^T x_i + b)$ (and hence the locn. of a training point x_i wrt DB/MBs) when:

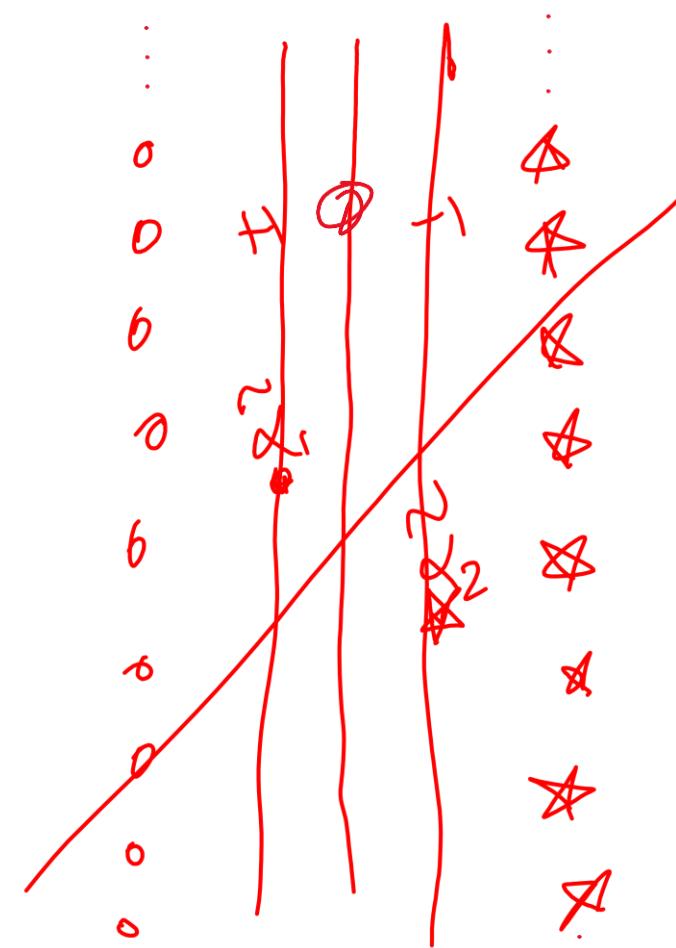
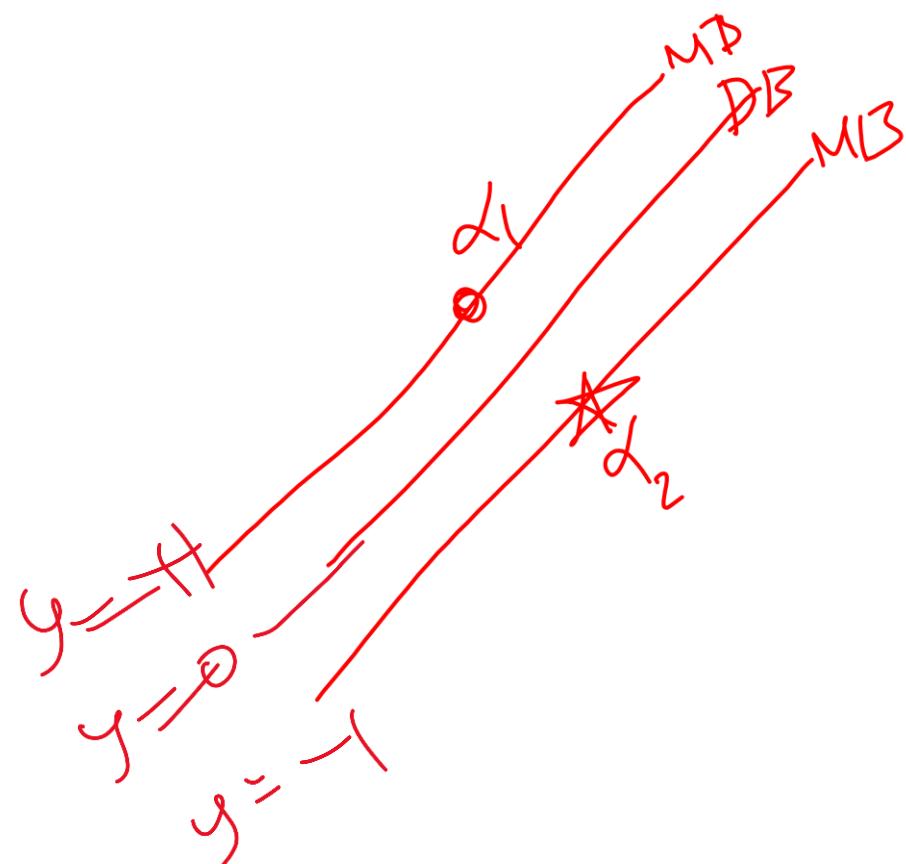
- $\alpha_i^* = 0$,
- $\alpha_i^* = C$,
- $0 < \alpha_i^* < C$?

(Hint: use KKT complementary slack and $\beta_i^* = C - \alpha_i^*$)

Example: the two SV case

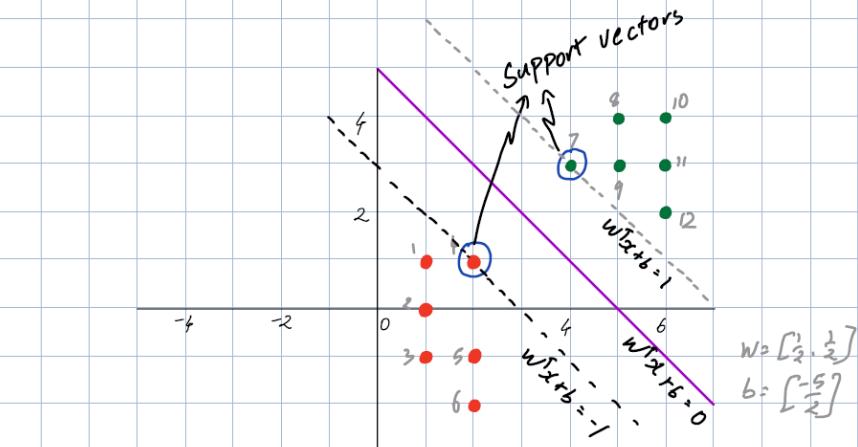


Example: the two SV case



Worked-out example:

At Optimality α is usually "Sparse"



The solution of hard Margin SVM for above problem with Linear Kernel :

$$\max_{\alpha} -\frac{1}{2} \alpha^T Y X^T X \alpha + \alpha^T$$

$$\text{s.t. } \sum \alpha_i y_i = 0 \quad K(u, v) = u^T v \\ \alpha_i \geq 0$$

$$\text{is } \alpha_4^* = \alpha_7^* = \frac{1}{4}, \text{ other } \alpha_i^* = 0$$

(Exercise: Check Optimality of above Point using KKT with just points 1, 4, 7, 8)

Calculate w^* and b^* :

$$\therefore w^* = \frac{1}{4} (-1) \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \frac{1}{4} (1) \begin{bmatrix} 4 \\ 3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

Take $i=4$ for calculating b^* (as $0 < \alpha_i < \infty$)

$$y_4 (w^{*T} x_4 + b^*) = 1$$

$$-1 \left(\begin{bmatrix} \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} + b^* \right) = 1$$

$$1 \cdot 5 + b^* = -1$$

$$\therefore b^* = -5/2$$

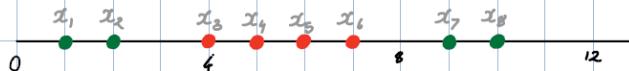
Exercise 1:

Exercise :

$$x^T = [1, 2, 4, 5, 6, 7, 9, 10]$$

$$y^T = [1, 1, -1, -1, -1, -1, 1, 1]$$

(Use software to
get intuition)



(Hard Margin)

(a) Solve Kernel SVM with (give optimal α^*)

$$k(u, v) = \exp(-\gamma(u-v)^2)$$

with $\gamma = 0.1$

Hint: $\alpha_2^* = \alpha_7^* = \alpha_3^* = \alpha_6^* > 0$ all other $\alpha_i^* = 0$

Find such an α^* and then use KKT conditions to show it is optimal.

(b) Give b^* for the α^* above.

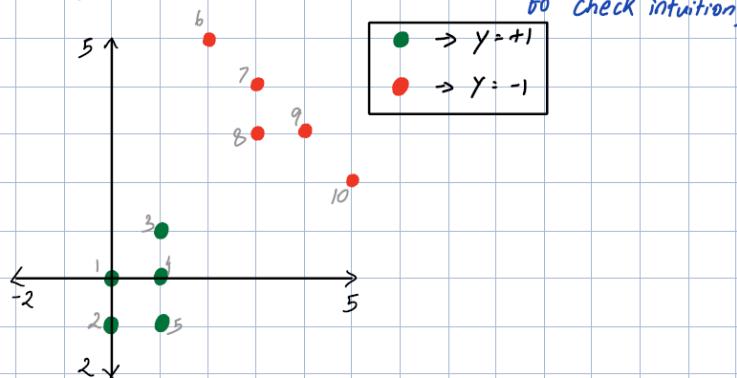
(c) Give the decision function $w^T \phi(x) + b^*$ for $x \in \mathbb{R}$.

Exercise2:

Exercise :

Consider the following hard margin SVM problem
with both w & b . Kernel is the linear kernel.

- Argue what points are support vectors.
- Argue what would be the optimal hyperplane and give w^*, b^* .
- Also argue what α^* should be. (Use software to check intuition)



- Repeat all the above if data point (x_8, y_8) is removed.
- Repeat a,b,c with $(x_8, y_8), (x_9, y_9), (x_7, y_7)$ removed.

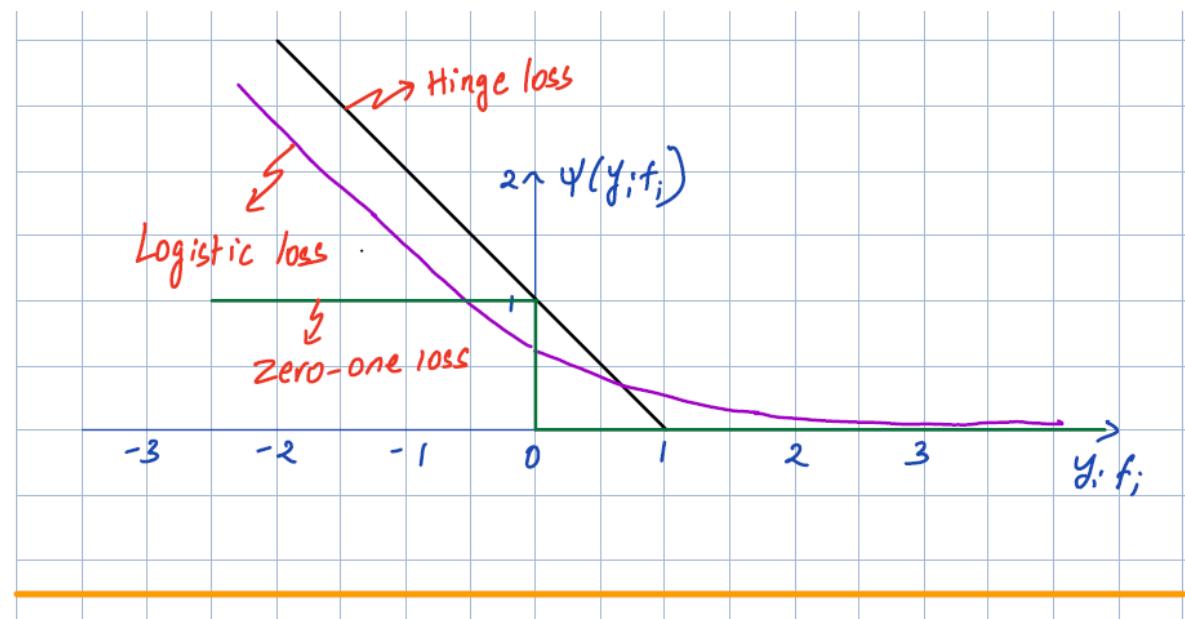
Hint: Optimal α^* need not be unique even if w^* and b^* are.

Loss function view: Hinge loss

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (w^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned} \quad \Rightarrow \quad \begin{aligned} \min_{w,b} \quad & C \sum_{i=1}^n \psi_H(y_i(w^T x_i + b)) \\ & + \frac{1}{2} \|w\|^2 \end{aligned}$$

Express ξ_i as a function of w, b :

$$\begin{aligned} \xi_i &\geq 1 - y_i(w^T x_i + b); \quad \xi_i \geq 0 \\ \therefore \xi_i &= \max(1 - y_i(w^T x_i + b), 0) \\ &= \max(1 - y_i f_i, 0) \quad \text{where } f_i = w^T x_i + b \\ &= \psi_H(y_i f_i) \end{aligned}$$



$$\psi_H(u) = \max(1 - u, 0)$$

SVM Interpretations

- Support vectors
- Kernel machines
- Loss fn. view

Goal we set out: concrete understanding of SVM --
hope we reached it!

sklearn.svm.SVC

```
class sklearn.svm.SVC(*, C=1.0, kernel='rbf', degree=3, gamma='scale', coef0=0.0, shrinking=True, probability=False, tol=0.001,  
cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', break_ties=False,  
random_state=None) \[source\]
```

Primal OP:

$$\min_{w,b,\zeta} \frac{1}{2} w^T w + C \sum_{i=1}^n \zeta_i$$

subject to $y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i,$
 $\zeta_i \geq 0, i = 1, \dots, n$

Dual OP:

$$\max_{\alpha} -\frac{1}{2} \alpha^T Q \alpha + \underline{\alpha^T \alpha}$$

subject to $y^T \alpha = 0$
 $0 \leq \alpha_i \leq C, i = 1, \dots, n$

Q is an n by n positive semidefinite matrix, $Q_{ij} \equiv y_i y_j K(x_i, x_j)$.

Prediction for new point x:

$$\sum_{i \in SV} y_i \alpha_i K(x_i, x) + b,$$

[Above formulas from sklearn help pages]

Concluding thoughts

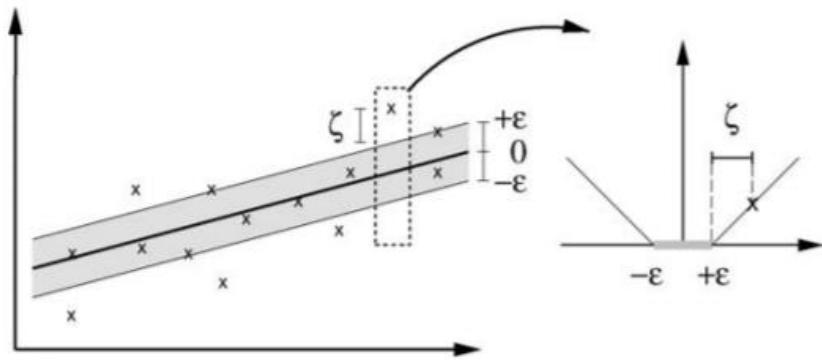
- SVM
 - Concept of max-margin classifn., sparse support vectors, & kernel machines.
 - Use of constrained optimization, and Primal - Dual problems.
 - Extensions: SVR (Support Vector Regression) and RVM (Relevance Vector Machines).
- Next steps: From linear to non-linear regression/classifn.

Non-linear method	(Non-linear) Basis functions	Objective function / OP
Vanilla extn. of linear models	Fixed – non-linear basis fns (feature map, $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$) fixed before seeing training data (manually via feature engineering); only weights of these basis fns learnt using training data.	Convex (unconstrained opt.)
SVM	Selective – center basis fns on training data points (dual/kernel view) and use training data to learn their weights and select a subset of them (non-zero weight support vectors) for eventual predictions.	Convex (constrained opt.)
<u>Neural networks</u>	Adaptive – Fix # of basis fns in advance, but allow them to be adaptive; parameterize basis fns and learn these parameters using training data.	Non-convex

Thank you!

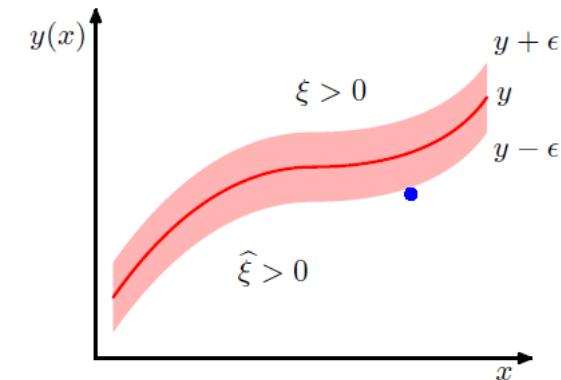
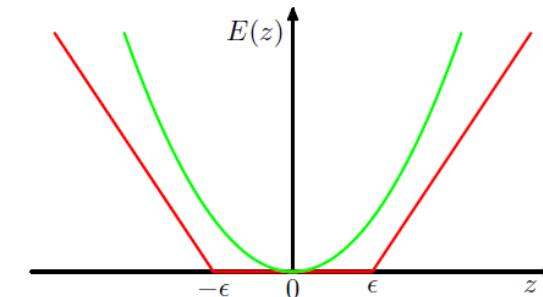
Backup

From classification to regression: Support Vector Regression or SVR (ϵ -insensitive “tube” and obj/loss fns.)



minimize $\frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} (\xi_i + \xi_i^*)$
 subject to $\begin{cases} y_i - \langle w, x_i \rangle - b \leq \varepsilon + \xi_i \\ \langle w, x_i \rangle + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases}$

(Soft)



[From CMB; Smola and Scholkopf, 2004 tutorial]

Loss functions drawn to scale

Plot of the 'hinge' error function used in support vector machines, shown in blue, along with the error function for logistic regression, rescaled by a factor of $1/\ln(2)$ so that it passes through the point $(0, 1)$, shown in red. Also shown are the misclassification error in black and the squared error in green.

