

Support Vector Machine and Convex Optimization

Ian En-Hsu Yen

Overview

- **Support Vector Machine**

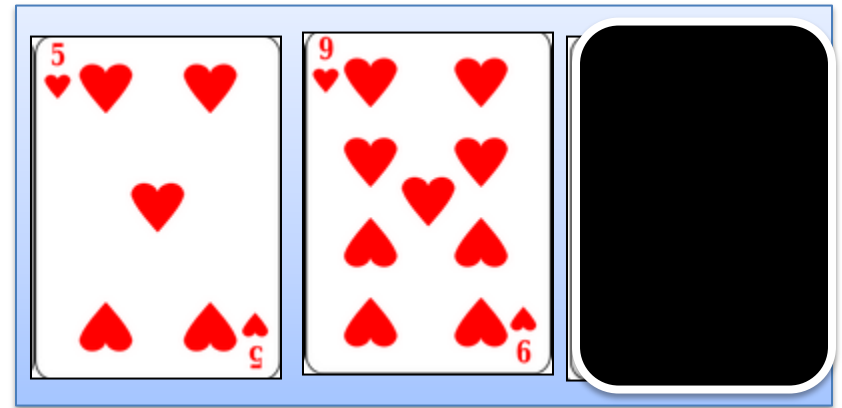
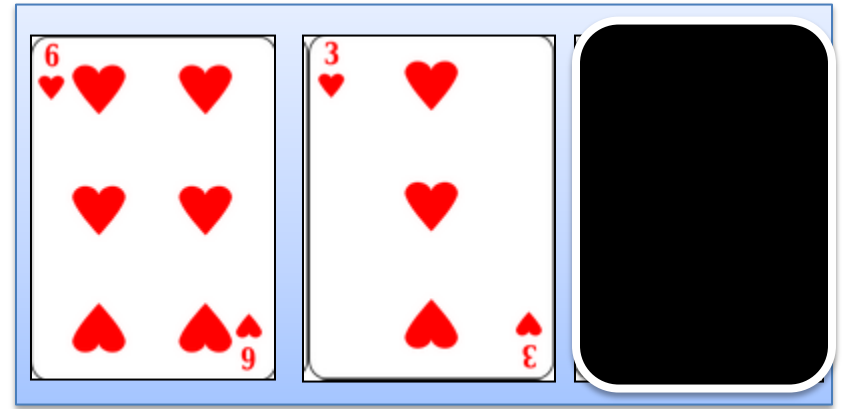
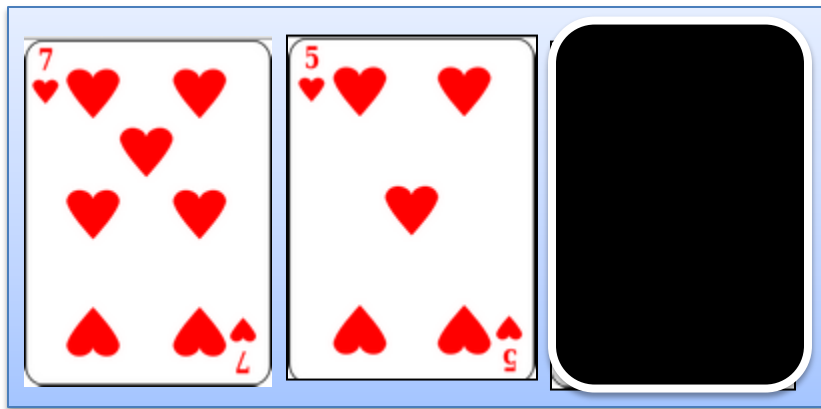
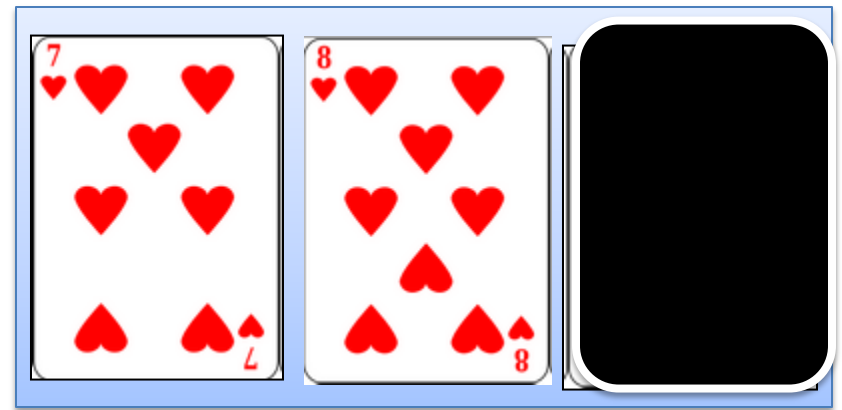
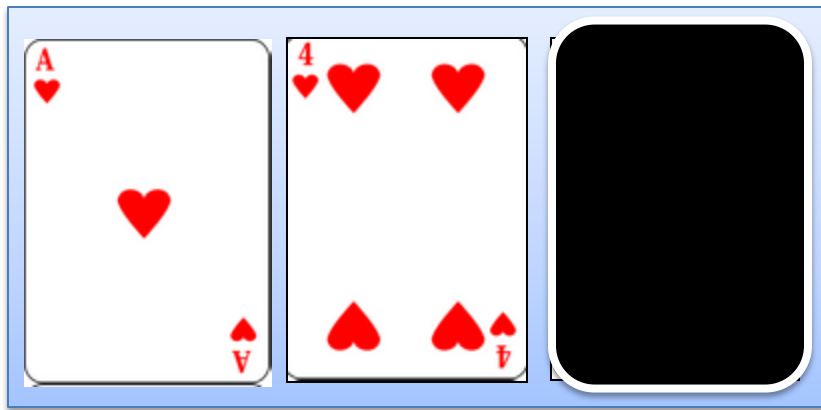
- The Art of Modeling --- Large Margin and Kernel Trick
- Convex Analysis
- Optimality Conditions
- Duality

- **Optimization for Machine Learning**

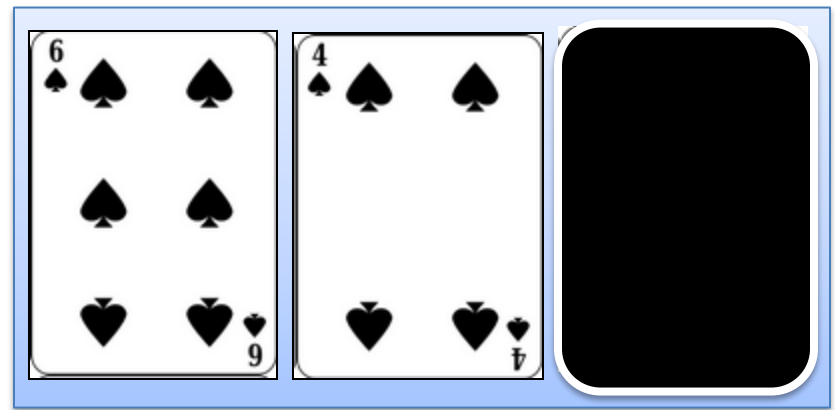
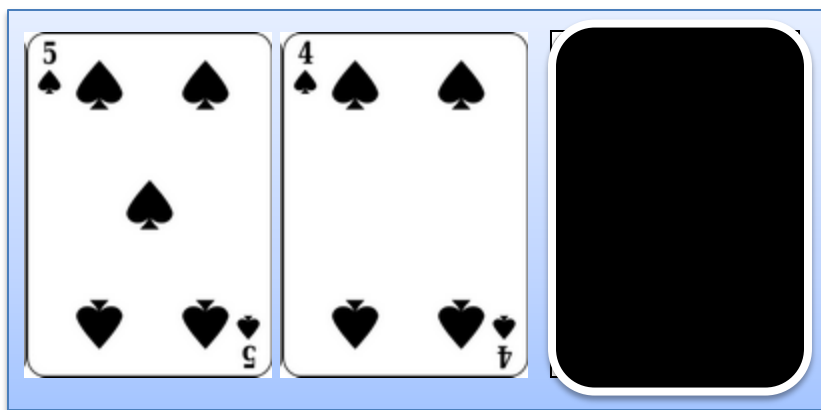
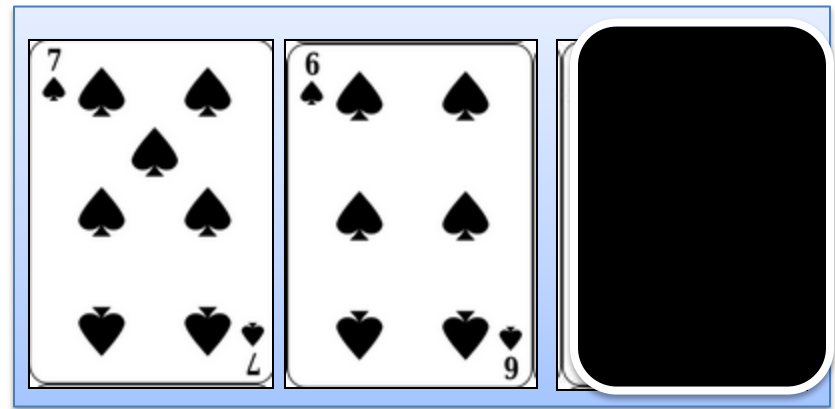
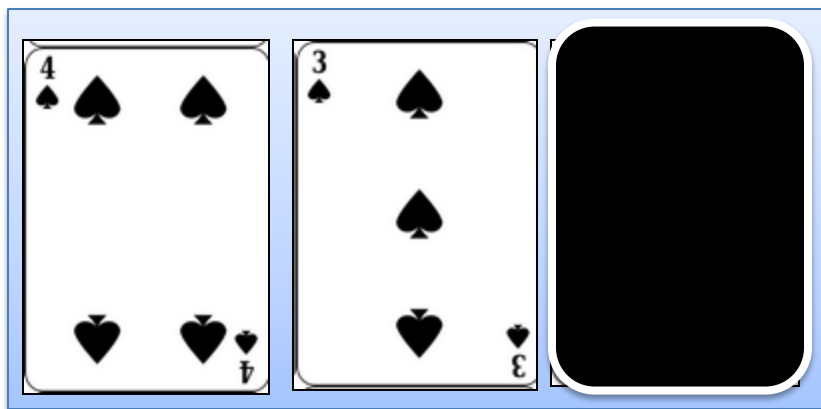
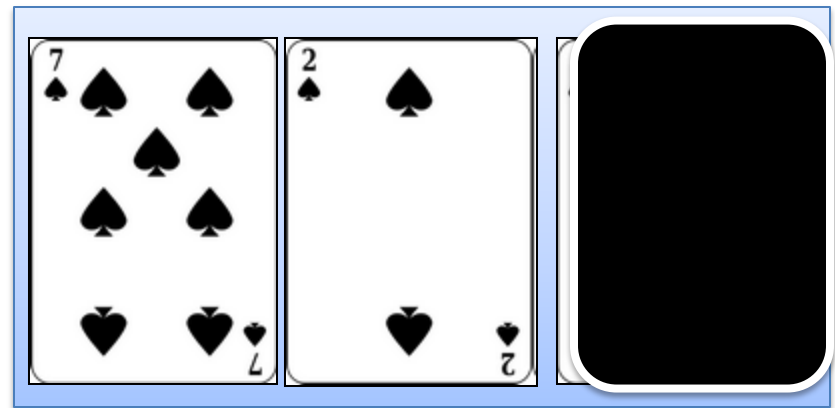
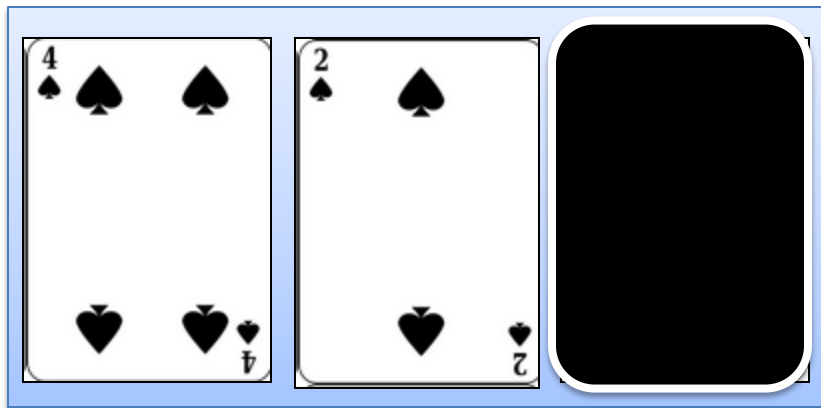
- Dual Coordinate Descent (fast convergence, moderate cost)
 - libLinear (Stochastic)
 - libSVM (Greedy)
- Primal Methods
 - Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)
 - Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)
- Demo

A Learning/Prediction Game

- Your team members suggest a Hypothesis Space : $\{h_1, h_2 \dots\}$
- You can only request one sample.
- Finding a hypothesis with accuracy $> 50\%$, you earn \$100,000.
wrong hypothesis ($\text{acc} \leq 50\%$) get \$100,000 punishment.



$H = \{ h_1 \}, h_1: (A+B) \bmod 13 = C$



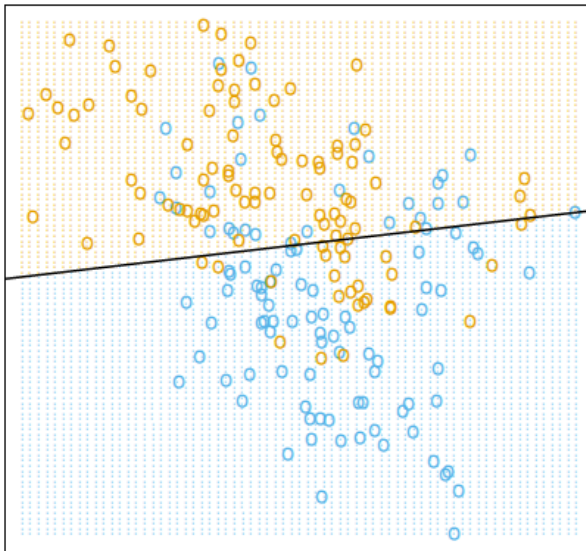
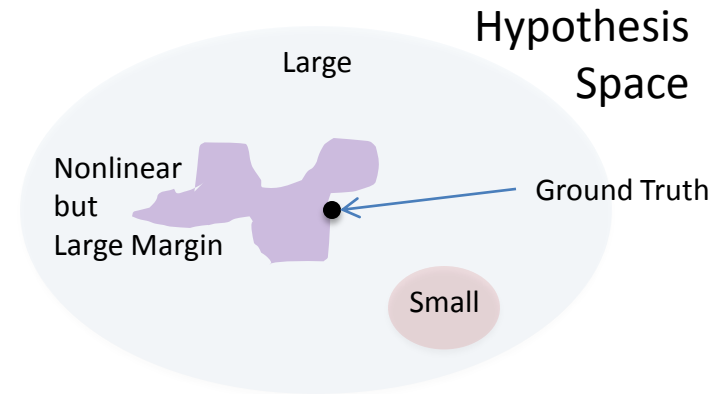
$H = \{ h_1, h_2 \}, \quad h_1: (A+B) \bmod 13 = C, \quad h_2: (A-B) \bmod 13 = C$

Large $|H|$ with Small $|Data|$ Guarantees Nothing

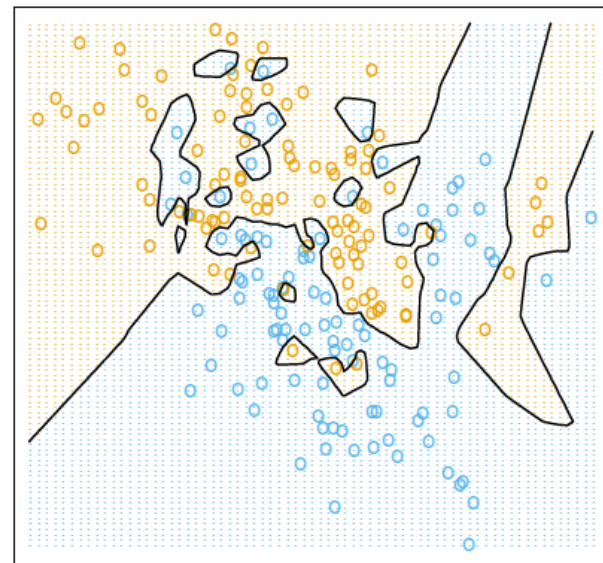
- First case: only one hypothesis h_1
 - $\Pr\{ |Train_Error - Test_Error| \geq 50\% \} \leq 1/2$.
 - Second case: two hypotheses h_1, h_2
 - $\Pr\{ |Train_Error - Test_Error| \geq 50\% \text{ for } h_1 \text{ or } h_2 \} \leq 1/2 + 1/2 = 1$.
- ➔ Guarantee Nothing.

Why Support Vector Machine (SVM) ?

- Flexible Hypothesis Space. (Non-linear Kernel)
- Not to Overfit (Large-Margin)
- Sparsity (Support Vectors)
- Easy to find Global Optimum (Convex Problem)



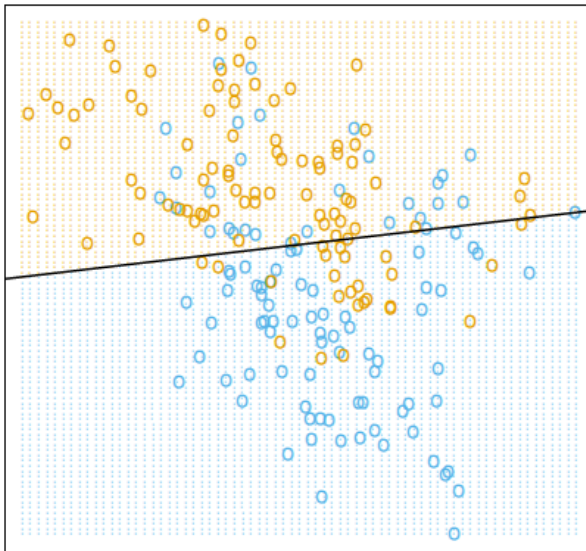
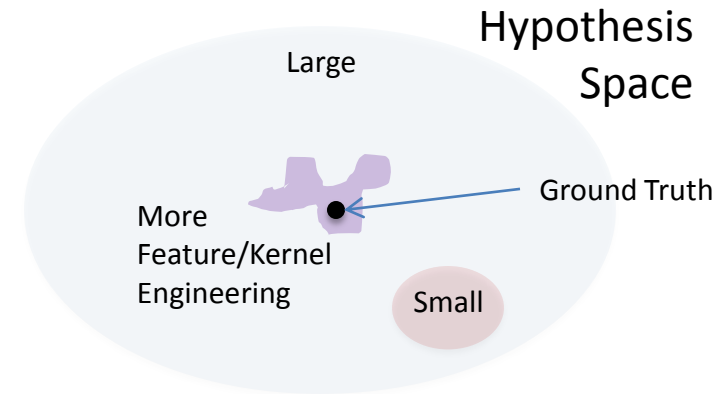
Linear Model



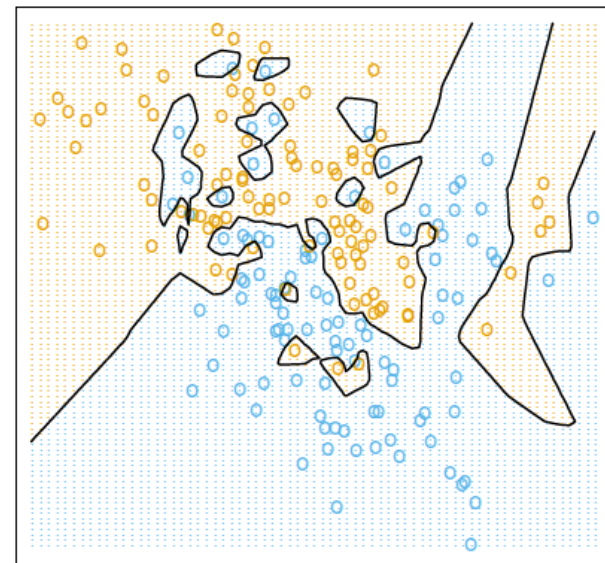
KNN

Why Support Vector Machine (SVM) ?

- Flexible Hypothesis Space. (Non-linear Kernel)
- Not to Overfit (Large-Margin)
- Sparsity (Support Vectors)
- Easy to find Global Optimum (Convex Problem)



Linear Model



KNN

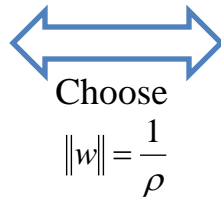
SVM: Large-Margin Perceptron

$$w^* = \arg \max_w \left\{ \min_n y_n \left(\frac{w^T}{\|w\|} x_n \right) \right\}$$



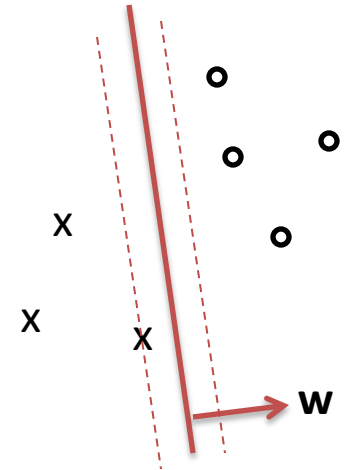
$$\max_{w, \rho} \rho$$

$$s.t. \ y_n \left(\frac{w^T}{\|w\|} x_n \right) \geq \rho, \ \forall n$$

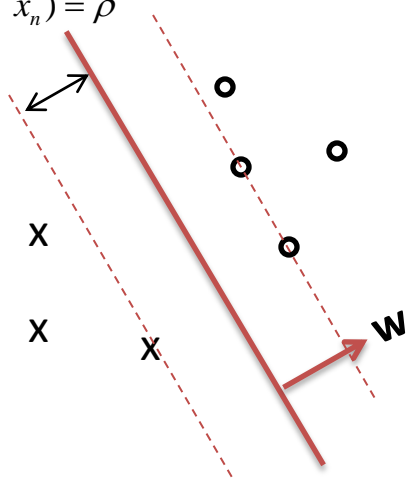


$$\max_w \frac{1}{\|w\|}$$

$$s.t. \ y_n (w^T x_n) \geq 1, \ \forall n$$



$$\min_n y_n \left(\frac{w^T}{\|w\|} x_n \right) = \rho$$



$$\min_w \|w\|$$

$$s.t. \ y_n (w^T x_n) \geq 1, \ \forall n$$



$$\min_w \|w\|^2$$

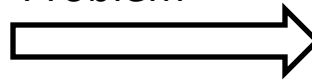
$$s.t. \ y_n (w^T x_n) \geq 1, \ \forall n$$

SVM: Large-Margin Perceptron

Hard Margin

$$\min_w \frac{1}{2} \|w\|^2$$
$$s.t. \ y_n(w^T x_n) \geq 1, \forall n$$

Non-Separable
Problem

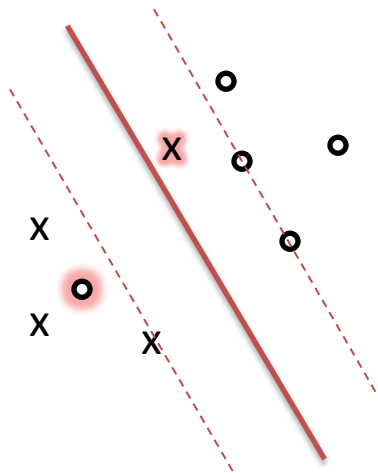


Soft Margin

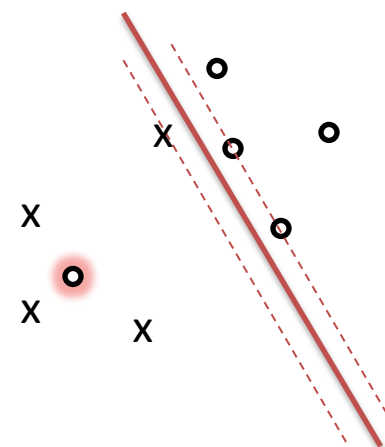
A drawback of SVM:
Solution sensitive to C

$$\min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$
$$s.t. \ y_n(w^T x_n) \geq 1 - \xi_n, \forall n$$

Small C



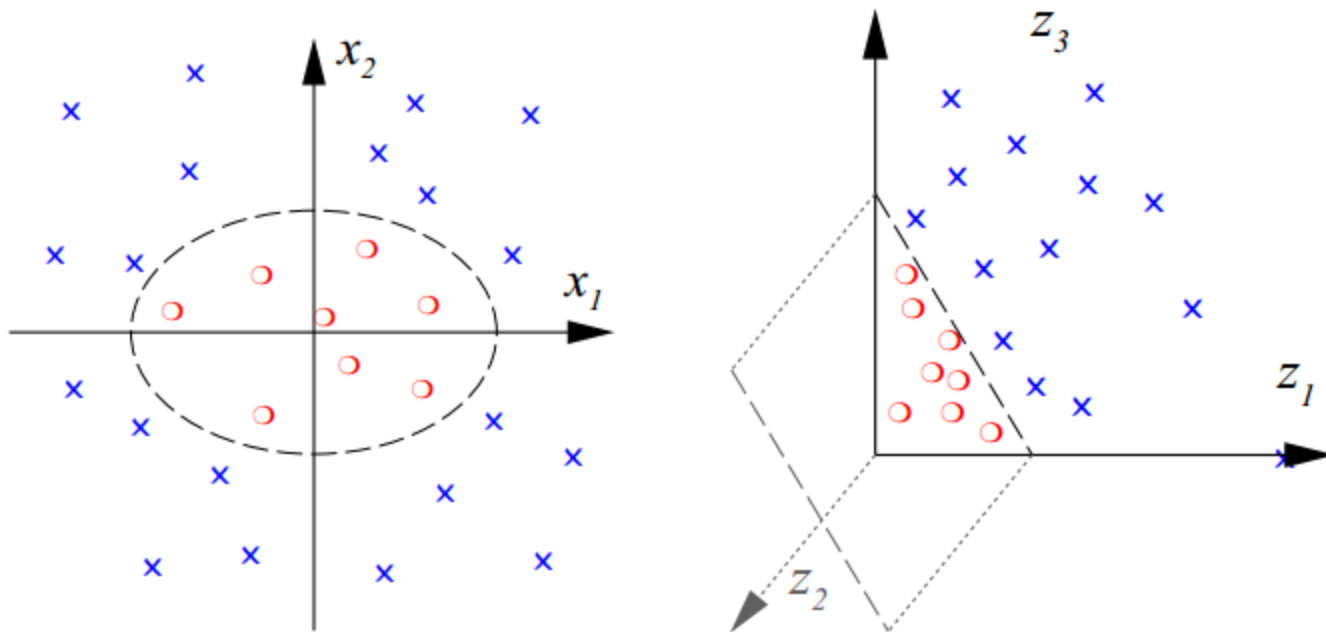
Large C



From Linear to Non-Linear

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



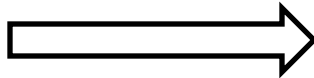
Perceptron: $\mathbf{a} x_1 + \mathbf{b} x_2 = 0$

Ellipse: $\mathbf{a} x_1^2 + \mathbf{b} x_2^2 + \mathbf{c} x_1 x_2 = 0$ (center at origin)

From Linear to Non-Linear

Linear SVM:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n (w^T x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned}$$



Non-linear SVM:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n (w^T \phi(x_n)) \geq 1 - \xi_n, \quad \forall n \end{aligned}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \phi(x_n) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_2x_3 \\ \sqrt{2}x_1x_3 \end{bmatrix}$$

SVM: Kernel Trick

Feature Expansion

$$x \rightarrow \phi(x)$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \phi(x_n) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_2x_3 \\ \sqrt{2}x_1x_3 \end{bmatrix}$$

$$3 \text{ features} \rightarrow 3 + C^3_2 = 6$$

$$100 \text{ features} \rightarrow 100 + C^{100}_2 = 5050$$

$$\text{Deg-2 Feature Expansion} \rightarrow O(D^2)$$

$$\text{Deg-K Feature Expansion} \rightarrow O(D^K)$$

Dot Product can be computed efficiently:

$$\phi(x)^T \phi(z) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_2x_3 \\ \sqrt{2}x_1x_3 \end{bmatrix}^T \begin{bmatrix} z_1^2 \\ z_2^2 \\ z_3^2 \\ \sqrt{2}z_1z_2 \\ \sqrt{2}z_2z_3 \\ \sqrt{2}z_1z_3 \end{bmatrix} = x_1^2 z_1^2 + x_2^2 z_2^2 + x_3^2 z_3^2 + 2(x_1 x_2 z_1 z_2 + x_2 x_3 z_2 z_3 + x_1 x_3 z_1 z_3) = \left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}^T \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} \right)^2 = (x^T z)^2$$

$O(D^K)$ **$O(D)$**

Compute dot Product using $K(x,z)=(x^T z)^2$
 \leftrightarrow deg-2 feature expansion



SVM: Kernel Trick

Feature Expansion

$$x \rightarrow \phi(x)$$

Can we formulate the problem only using dot product $\phi(x_i)^T \phi(x_j)$?

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned}$$



$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \alpha^T \Phi^T \Phi \alpha + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n \sum_{i=1} \alpha_i y_i \phi(x_i)^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned}$$



$$\begin{aligned} \min_{\alpha, \xi \geq 0} \quad & \frac{1}{2} \alpha^T Q \alpha + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n \sum_{i=1} \alpha_i y_i K(x_i, x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned}$$

$$Q_{ij} = (y_i \phi(x_i))(y_j \phi(x_j)) = y_i y_j K(x_i, x_j)$$

By **Representer Theorem**, solution w^* of the problem can be expressed as **linear combination of instances**:

$$w^* = \sum_n \alpha_n y_n \phi(x_n) = [y_1 \phi(x_1) \quad \dots \quad y_N \phi(x_N)]_{D \times N} \begin{bmatrix} \alpha_1 \\ \dots \\ \alpha_N \end{bmatrix} = \Phi \alpha$$

Prediction using only dot product $\phi(x_i)^T \phi(x_j)$:

$$\begin{aligned} w^T \phi(x_t) &= \left(\sum_n \alpha_n y_n \phi(x_n) \right)^T (\phi(x_t)) \\ &= \sum_n \alpha_n y_n \phi(x_n)^T \phi(x_t) = \sum_n \alpha_n y_n K(x_n, x_t) \end{aligned}$$

O(N*D)

or **O(|Support Vector| * D)**

SVM: Kernel Trick

Some popular Kernels:

Polynomial Kernel: $K(x, x') = (x^T x' + 1)^d$

RBF Kernel: $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

Linear Kernel: $K(x, x') = x^T x'$

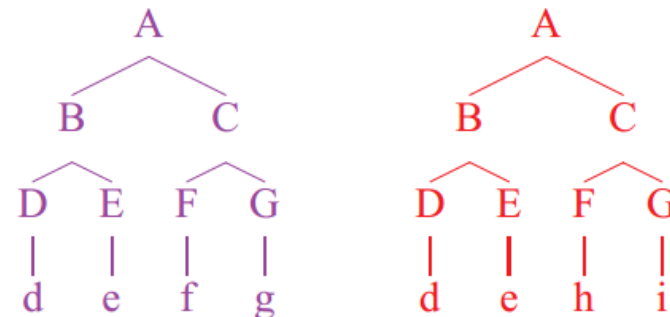
Kernels may be easier to define than Features :

String Kernel: Gene Classification / Rewriting or not

Tree Kernel: Syntactic parse tree classification

Graph Kernel: Graph Type Classification

Shakespeare wrote Hamlet.
Hamlet was written by Shakespeare.



Overview

- **Support Vector Machine**
 - The Art of Modeling --- Large Margin and Kernel Trick
 - Convex Analysis
 - Optimality Conditions
 - Duality
- **Optimization for Machine Learning**
 - Dual Coordinate Descent (fast convergence, moderate cost)
 - libLinear (Stochastic)
 - libSVM (Greedy)
 - Primal Methods
 - Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)
 - Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)

Convex Analysis

General Optimization Problem:

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq b_i, \quad i = 1, \dots, m\end{array}$$

is very difficult to solve. (very long time vs. approximate)

Optimization is much easier if the problem is **convex**, that is:

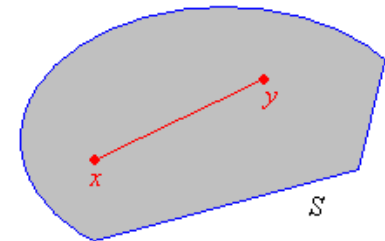
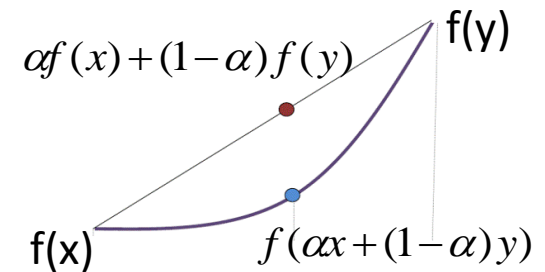
1. The **objective** function is **convex**:

$$f_0(\alpha x + (1-\alpha)y) \leq \alpha f_0(x) + (1-\alpha)f_0(y) \quad \text{for } 0 \leq \alpha \leq 1$$

2. The **feasible domain** (constrained space) is **convex**:

$$\text{if } x \in C, y \in C \Rightarrow \alpha x + (1-\alpha)y \in C, \quad 0 \leq \alpha \leq 1$$

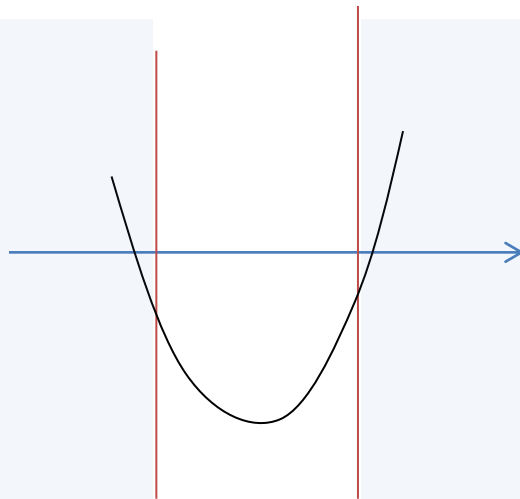
➔ All local minimum is global minimum !!



Convex Analysis

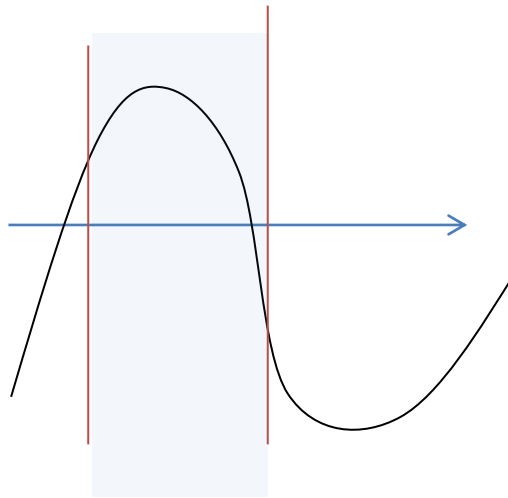
Simple Example:

$x \leq a$ $x \geq b$



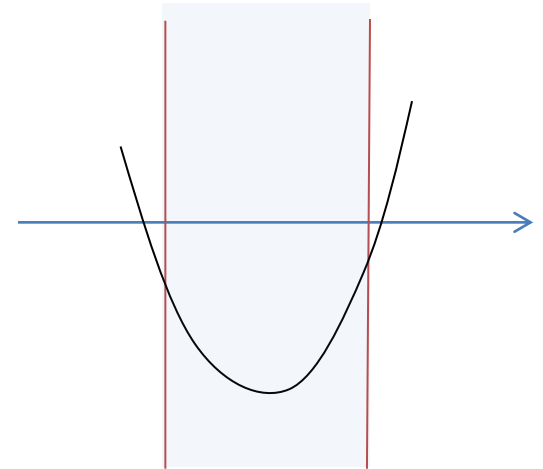
Non-Convex Set ;
Convex function

$x \geq a$ $x \leq b$



Convex Set ;
Non-Convex function

$x \geq a$ $x \leq b$

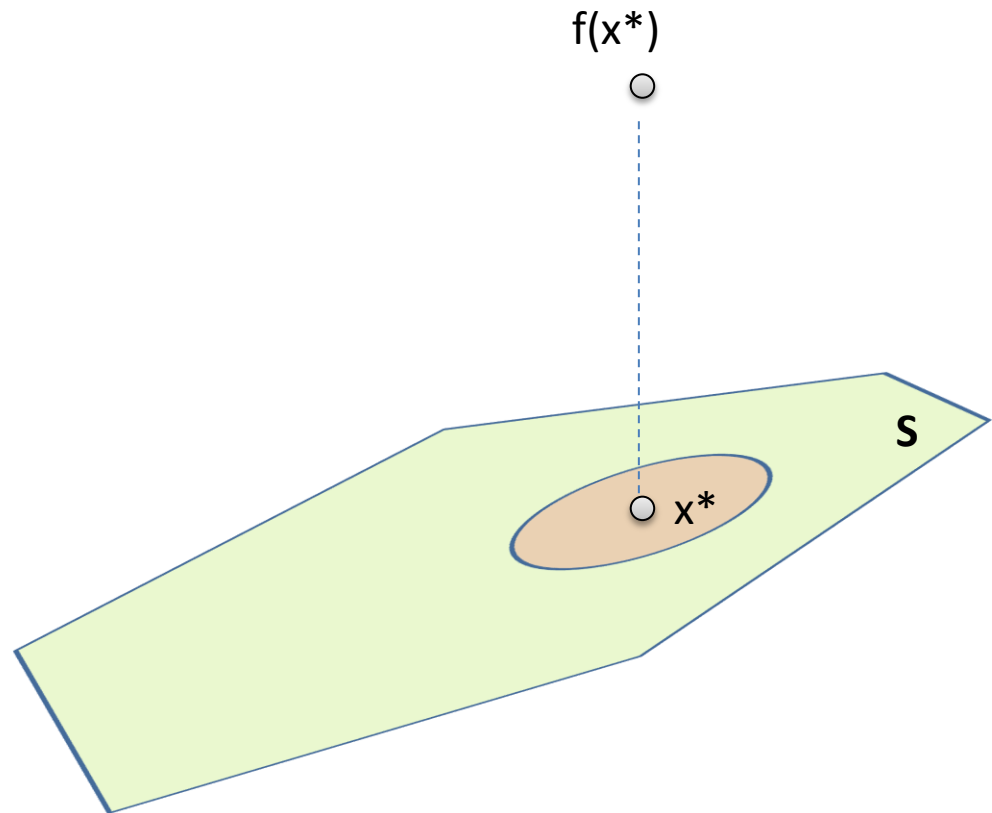


Convex Set ;
Convex function

Convex Analysis

x^* is local minimum $\rightarrow x^*$ is global minimum (why?)

If x^* is a **local minimum**, there is a “ball”,
in which any feasible x' has $f(x') \geq f(x^*)$.

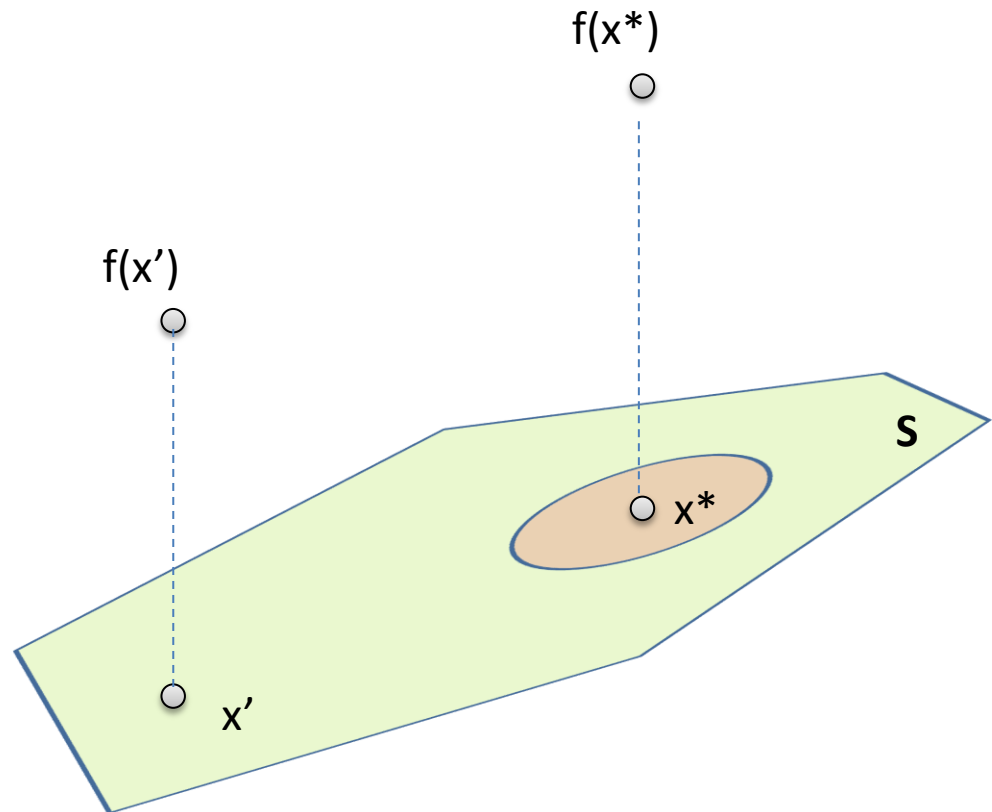


Convex Analysis

x^* is local minimum $\rightarrow x^*$ is global minimum (why?)

If x^* is a **local minimum**, there is a “ball”, in which any feasible x' has $f(x') \geq f(x^*)$.

Assume for contradiction that x^* is **not a global minimum**. There should be a feasible x' with $f(x') < f(x^*)$.



Convex Analysis

x^* is local minimum $\rightarrow x^*$ is global minimum (why?)

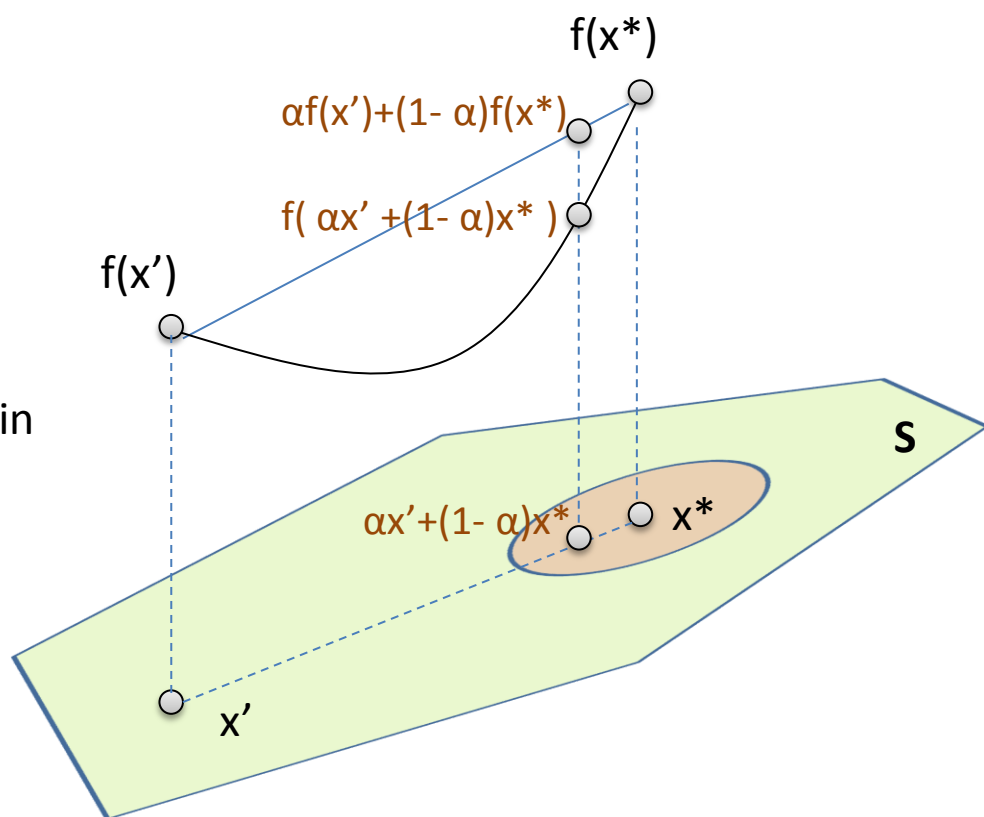
If x^* is a **local minimum**, there is a “ball”, in which any feasible x' has $f(x') \geq f(x^*)$.

Assume for contradiction that x^* is **not a global minimum**. There should be a feasible x' with $f(x') < f(x^*)$.

Then we can find a **feasible** $\alpha x' + (1 - \alpha)x^*$ in the **ball** with:

$$\begin{aligned} f(\alpha x' + (1 - \alpha)x^*) &\leq \alpha f(x') + (1 - \alpha)f(x^*) \\ &< f(x^*) \end{aligned}$$

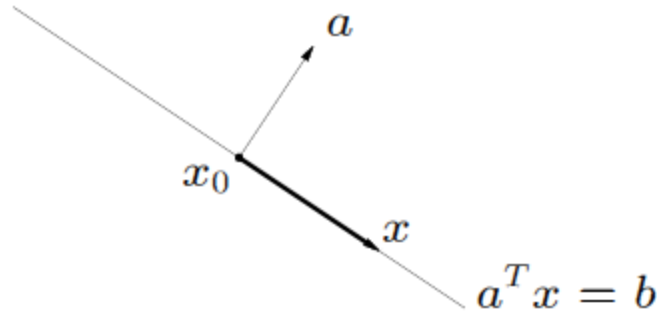
\rightarrow contradiction.



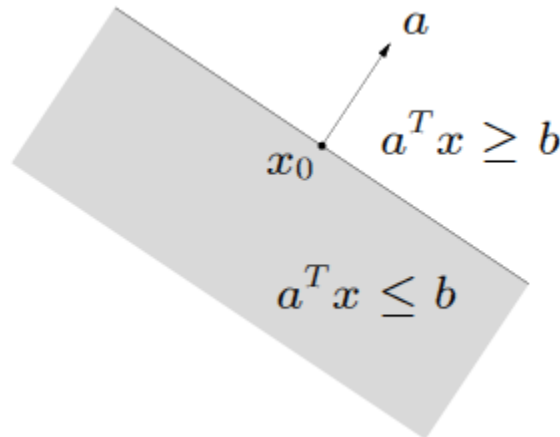
Convex Analysis

Example of Convex Set: *if* $x \in C, y \in C \Rightarrow \alpha x + (1-\alpha)y \in C, \quad 0 \leq \alpha \leq 1$

Linear equality constraint (Hyperplane) $\{x \mid a^T x = b\} \quad (a \neq 0)$



Linear inequality constraint (Halfspace) $\{x \mid a^T x \leq b\} \quad (a \neq 0)$



Convex Analysis

Example of Convex Set: *if* $x \in C, y \in C \Rightarrow \alpha x + (1 - \alpha)y \in C, \quad 0 \leq \alpha \leq 1$

Intersection of Convex Set :

$$\begin{cases} a_1 x \leq b_1 \\ a_2 x \leq b_2 \\ a_3 x \leq b_3 \\ c_4 x = d_4 \\ c_5 x = d_5 \end{cases} \quad (Ax \leq b, \quad Cx = d)$$

$x, y \in A \cap B, \quad A, B \text{ is convex}$

$\alpha x + (1 - \alpha)y \in A \cap B \quad ?$

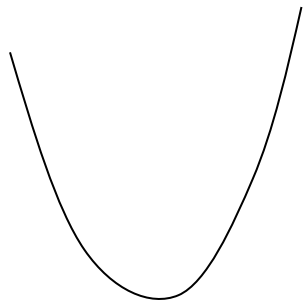
Convex Analysis

Example of Convex Function: $f_0(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ for $0 \leq \alpha \leq 1$

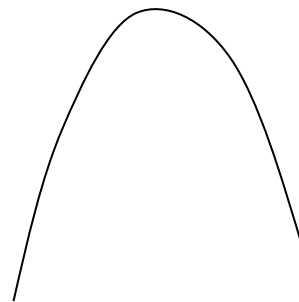
Linear Function $f(x) = c^T x$

Quadratic Function $f(x) = \frac{1}{2} x^T Q x + c^T x$?

Obviously, it depends



$$ax^2 + bx + c, a > 0$$



$$ax^2 + bx + c, a < 0$$

Convex Analysis

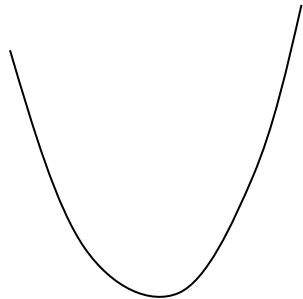
Example of Convex Function: $f_0(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ for $0 \leq \alpha \leq 1$

Linear Function $f(x) = c^T x$

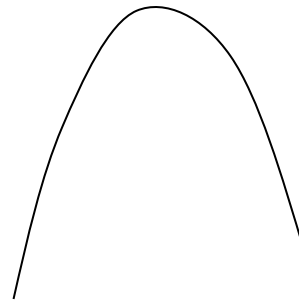
Quadratic Function $f(x) = \frac{1}{2} x^T Q x + c^T x$?

A practical way to check **convexity** :

Check the **second derivative** $\frac{\partial^2 f(x)}{\partial x^2} \geq 0$ at $\forall x$



$ax^2 + bx + c, a > 0$



$ax^2 + bx + c, a < 0$

Convex Analysis

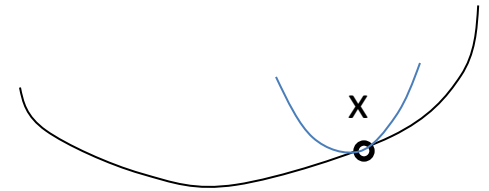
Example of Convex Function: $f_0(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ for $0 \leq \alpha \leq 1$

Linear Function $f(x) = c^T x$

Quadratic Function $f(x) = \frac{1}{2} x^T Q x + c^T x$?

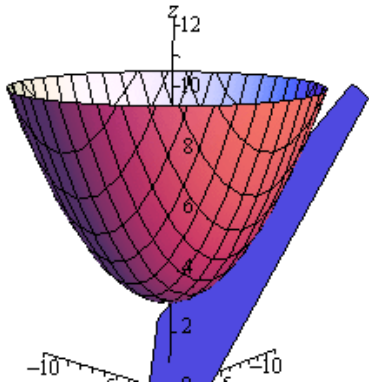
In \mathbf{R}^D , we have convexity if the **Hessian Matrix** :

$$\frac{\partial^2 f(x)}{\partial x^2} = H(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_D x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_D^2} \end{bmatrix}$$

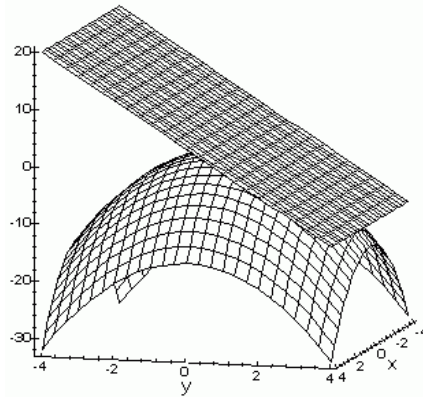


is positive semidefinite at $\forall x$
 ($z^T H(x) z \geq 0$ for $\forall z$)
 (all eigenvalue ≥ 0)

Convex Analysis

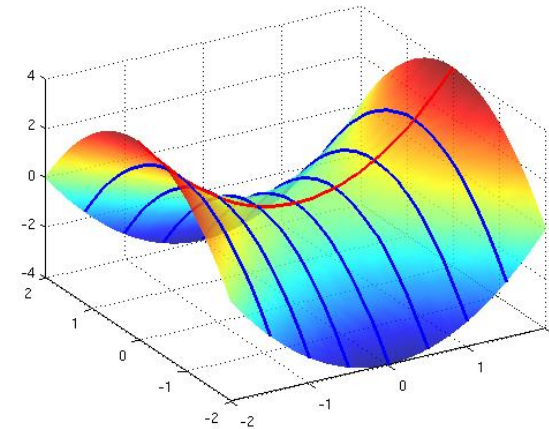


H is positive definite



H is negative definite

Other Cases ?



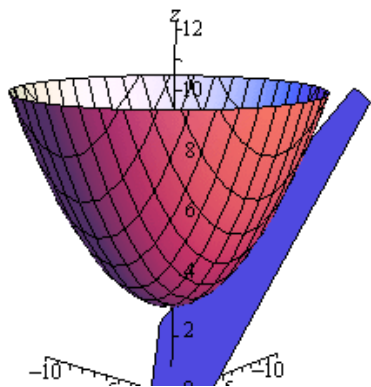
**H is not positive (semi-)definite
not negative (semi-)definite**

In \mathbb{R}^D , we have convexity if the **Hessian Matrix** :

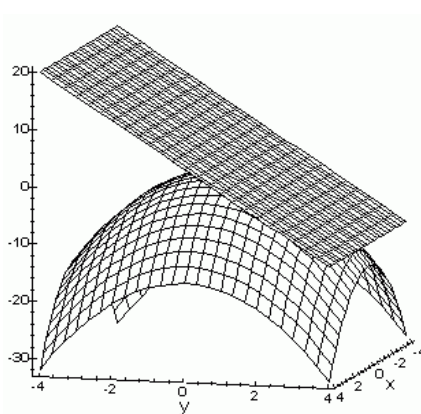
$$\frac{\partial^2 f(x)}{\partial x^2} = H(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_D x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_D^2} \end{bmatrix}$$

is positive(semi-)definite at $\forall x$
 ($z^T H(x) z \geq 0$ for $\forall z$)
 (all eigenvalue ≥ 0)

Convex Analysis



H is positive definite

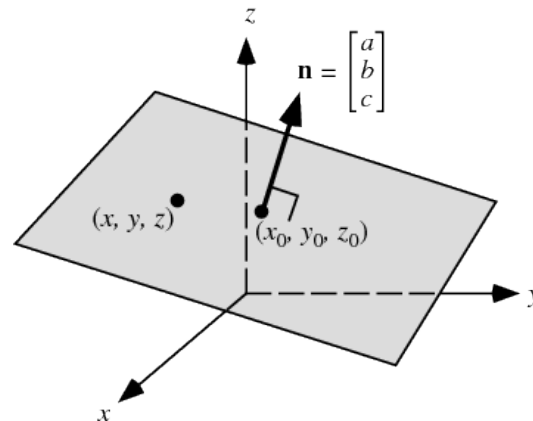


H is negative definite

In \mathbf{R}^D , we have convexity if the **Hessian Matrix** :

$$\frac{\partial^2 f(x)}{\partial x^2} = H(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_D x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_D^2} \end{bmatrix}$$

Other Cases ?



H=0 is **positive (semi-)definite**
is **negative (semi-)definite**

is positive(semi-)definite at $\forall x$
($z^T H(x) z \geq 0$ for $\forall z$)
(all eigenvalue ≥ 0)

Convex Analysis

Example of Convex Function: $f_0(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ for $0 \leq \alpha \leq 1$

Linear Function $f(x) = c^T x$ $(\frac{\partial^2 f(x)}{\partial x^2} = H(x) = 0)$

Quadratic Function $f(x) = \frac{1}{2} x^T Q x + c^T x$? $(\frac{\partial^2 f(x)}{\partial x^2} = H(x) = Q)$

Quadratic Function is convex if **Q is positive semi-definite**.

$$\begin{array}{l} \min_w \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_n w^T \phi(x_n) \geq 1, \forall n \end{array}$$

$$\frac{1}{2} \|w\|^2 = \frac{1}{2} w^T I w \quad \text{Is } I \text{ positive-semidefinite?}$$

Convex Analysis

Example of Convex Function: $f_0(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ for $0 \leq \alpha \leq 1$

Linear Function $f(x) = c^T x$ $(\frac{\partial^2 f(x)}{\partial x^2} = H(x) = 0)$

Quadratic Function $f(x) = \frac{1}{2} x^T Q x + c^T x$? $(\frac{\partial^2 f(x)}{\partial x^2} = H(x) = Q)$

Quadratic Function is convex if **Q is positive semi-definite**.

$$\begin{array}{ll} \min_{w, \xi \geq 0} & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} & y_n w^T \phi(x_n) \geq 1 - \xi_n, \forall n \end{array}$$

Is $H\left(\begin{bmatrix} w \\ \xi \end{bmatrix}\right) = \begin{bmatrix} I_{D \times D} & O_{D \times N} \\ O_{D \times N} & O_{N \times N} \end{bmatrix}$ positive-semidefinite?

Half-space constraint

**SVM problem is a convex problem.
(Quadratic Program)**

Convex Analysis

Example of Convex Problem:

Linear Programming:

$$\begin{array}{ll}\text{minimize} & c^T x + d \\ \text{subject to} & Gx \preceq h \\ & Ax = b\end{array}$$

Linear objective function
s.t. **Linear** Constraint.

Quadratic Programming:

$$\begin{array}{ll}\text{minimize} & (1/2)x^T Px + q^T x + r \\ \text{subject to} & Gx \preceq h \\ & Ax = b\end{array}$$

Quadratic objective function
s.t. **Linear** Constraint.

where P must be positive – semidefinite

Overview

- **Support Vector Machine**

- The Art of Modeling --- Large Margin and Kernel Trick
- Convex Analysis
- Optimality Conditions
- Duality

- **Optimization for Machine Learning**

- Dual Coordinate Descent (fast convergence, moderate cost)
 - libLinear (Stochastic)
 - libSVM (Greedy)
- Primal Methods
 - Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)
 - Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)

Optimality Condition

There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

First, we consider **Unconstrained Problem**:

$$\min_x f(x)$$

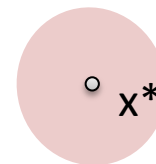
Example: Matrix Factorization (non-convex)

$$\min_{P,Q} \sum (r_{ui} - \mathbf{p}_u^T \mathbf{q}_i)^2 + \lambda_p \|\mathbf{P}\|^2 + \lambda_q \|\mathbf{Q}\|^2$$

Optimality Condition

There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

First, we consider **Unconstrained Problem**:



$$\min_x f(x)$$

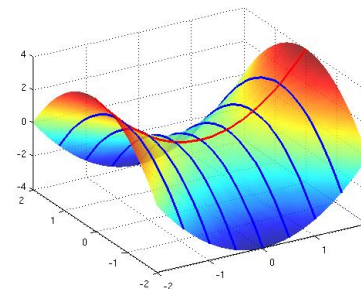
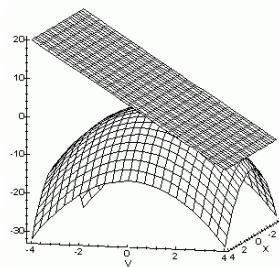
\mathbf{x}^* is Local minimizer $\Leftrightarrow f(x^*) \leq f(x^* + p)$ for all p with $\|p\| < \varepsilon$

For twice-differentiable $\mathbf{f}(\mathbf{x})$, consider the **Taylor Expansion**:

$$f(x^* + p) = f(x^*) + \nabla f(x^*)^T p + \frac{1}{2} p^T \nabla^2 f(x^*) p + \dots$$

\mathbf{x}^* is Local minimizer $\rightarrow \nabla f(x^*) = 0$

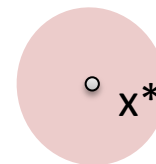
$\nabla f(x^*) = 0 \rightarrow \mathbf{x}^*$ is Local minimizer ?



Optimality Condition

There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

First, we consider **Unconstrained Problem**:



$$\min_x f(x)$$

\mathbf{x}^* is Local minimizer $\Leftrightarrow f(x^*) \leq f(x^* + p)$ for all p with $\|p\| < \varepsilon$

For twice-differentiable $\mathbf{f}(\mathbf{x})$, consider the **Taylor Expansion**:

$$f(x^* + p) = f(x^*) + \nabla f(x^*)^T p + \frac{1}{2} p^T \nabla^2 f(x^*) p + \dots$$

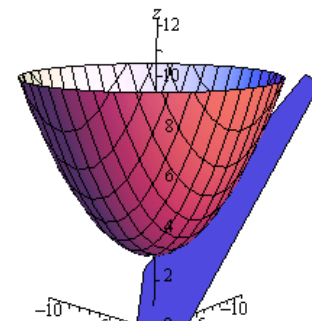
\mathbf{x}^* is Local minimizer $\Rightarrow \nabla f(x^*) = 0$

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

$\Rightarrow \mathbf{x}^*$ is Local minimizer

$\nabla^2 f(x^*)$ is positive – semidefinite

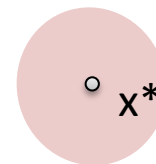
No need to check for Convex function (why?)



Optimality Condition

There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

First, we consider **Unconstrained Problem**:



$$\min_x f(x)$$

x^* is Local minimizer $\iff f(x^*) \leq f(x^* + p)$ for all p with $\|p\| < \varepsilon$

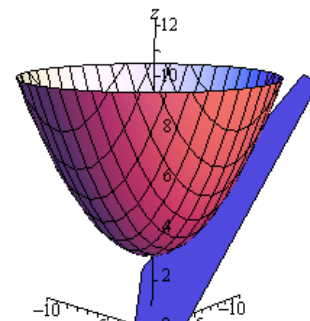
For twice-differentiable $f(x)$, consider the **Taylor Expansion**:

$$f(x^* + p) = f(x^*) + \nabla f(x^*)^T p + \frac{1}{2} p^T \nabla^2 f(x^*) p + \dots$$

For **Convex** $f(x)$:

x^* is Global minimizer $\iff \nabla f(x^*) = 0$

Assume convexity for now on.....



Optimality Condition

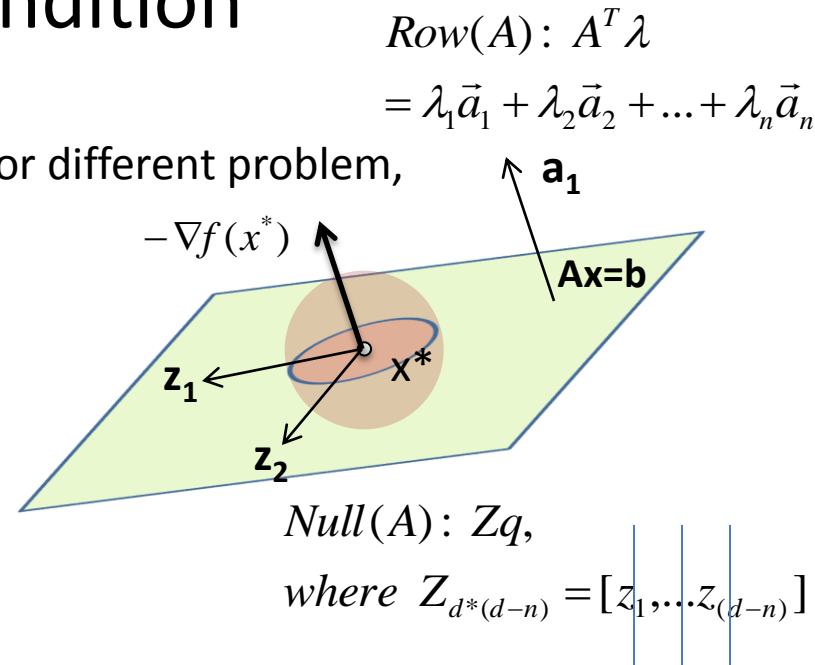
There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

Now, consider **Equality Constrained Problem**:

$$\min_x f(x)$$

$$s.t. \quad Ax = b \quad (\text{ex. } a_1^T x = b)$$

(nonlinear equality is, in general, not convex)



x^* is Local minimizer $\Leftrightarrow f(x^*) \leq f(x^* + p)$ for all "feasible" p with $\|p\| < \varepsilon$

$\Leftrightarrow f(x^*) \leq f(x^* + Zq)$ for all q with $\|q\| < \varepsilon$

For twice-differentiable $f(x)$, consider the **Taylor Expansion**:

$$f(x^* + Z^T q) = f(x^*) + (Z^T \nabla f(x^*))^T q + \frac{1}{2} q^T Z^T \nabla^2 f(x^*) Z q + \dots$$

$f(x)$ is convex, x^* is Global minimizer $\Leftrightarrow Z^T \nabla f(x^*) = 0$

$\Leftrightarrow -\nabla f(x^*) = A^T \lambda$ ($\nabla f(x^*)$ in Row(A))

Largrange Multipliers

Optimality Condition

$$\begin{aligned} \text{Row}(A) &: A^T \lambda \\ &= \lambda_1 \vec{a}_1 + \lambda_2 \vec{a}_2 + \dots + \lambda_n \vec{a}_n \end{aligned}$$

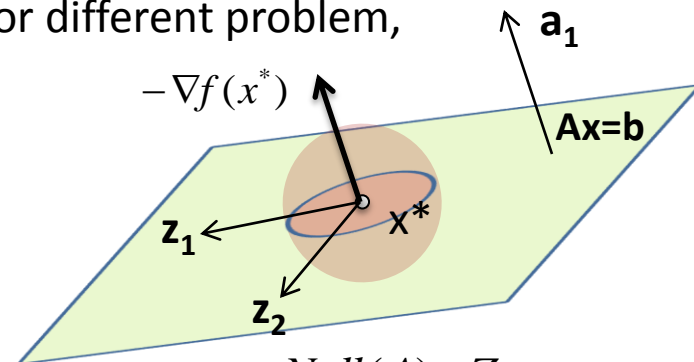
There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

Now, consider **Equality Constrained Problem**:

$$\min_x f(x)$$

$$s.t. \quad Ax = b \quad (\text{ex. } a_1^T x = b)$$

(nonlinear equality is, in general, not convex)



$$\begin{aligned} \text{Null}(A) &: Zq, \\ \text{where } Z_{d^*(d-n)} &= [z_1, \dots, z_{(d-n)}] \end{aligned}$$

$$x^* \text{ is Local (Global) minimizer} \iff Z^T \nabla f(x^*) = 0$$

Lagrange Multipliers

$$\text{or } -\nabla f(x^*) = A^T \lambda \quad (-\nabla f(x^*) \text{ in Row}(A))$$

$$(\text{ex. } -\nabla f(x^*) = \lambda \vec{a}_1)$$

$$\lambda_n > 0 \Rightarrow ?$$

$$\lambda_n < 0 \Rightarrow ?$$

$$\lambda_n = 0 \Rightarrow ?$$

Optimality Condition

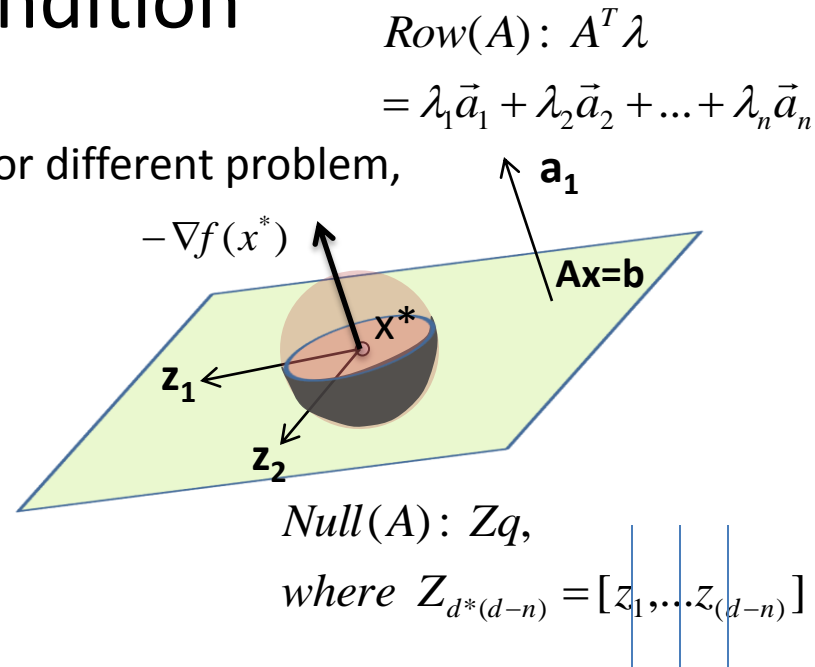
There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

Now, consider **Inequality Constrained Problem**:

$$\min_x f(x)$$

$$s.t. \quad Ax \leq b \quad (\text{ex. } a_1^T x \leq b)$$

(Assume linear inequality for simplicity.)



Let A^* (some rows of A) be the coefficients of **binding constraints**:

$$x^* \text{ is Local (Global) minimizer} \iff -\nabla f(x^*) = A^{*T} \lambda \quad (\text{ex. } -\nabla f(x^*) = \lambda \vec{a}_1)$$

and $\lambda \geq 0$ (feasible direction **not decrease** $f(x)$)

$$\lambda_n < 0 \quad \rightarrow \quad \text{Detach from } a_n x < b \quad \text{can decrease } f(x)$$

Optimality Condition

There are many, many different **solvers** designed for different problem, but they share the same **optimality condition**.

Now, consider **Inequality Constrained Problem**:

$$\min_x f(x)$$

$$s.t. \boxed{Ax \leq b} \quad (\text{ex. } a_1^T x \leq b)$$

(Assume linear inequality for simplicity.)

Require $\lambda_n = 0$ **for non-binding constraint:**

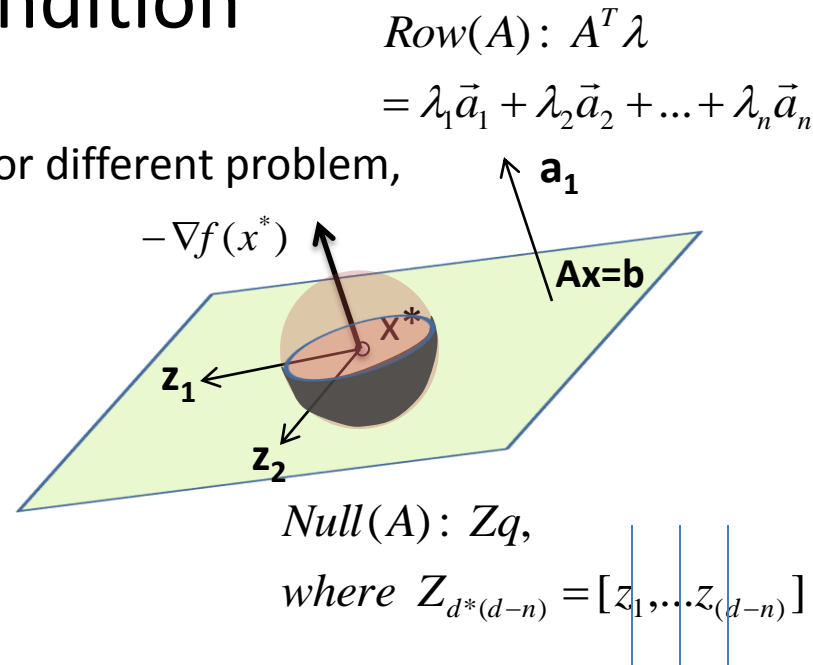
$$\boxed{\lambda_n (a_n x - b_n) = 0}, \quad \forall n$$

x^* is Local (Global) minimizer \iff

$$\boxed{-\nabla f(x^*) = A^T \lambda} \quad (\text{ex. } -\nabla f(x^*) = \lambda \vec{a}_1)$$

KKT conditions.

and $\boxed{\lambda \geq 0}$ (feasible direction **not decrease** $f(x)$)



Optimality Condition for SVM

What's the KKT condition for:

$$\begin{bmatrix} -\nabla_w f(w, \xi) \\ -\nabla_\xi f(w, \xi) \end{bmatrix} = \begin{bmatrix} -w \\ -C \end{bmatrix} = \begin{bmatrix} -\sum_n \alpha_n y_n \phi(x_n) \\ -\alpha_n - \beta_n \end{bmatrix}$$

$$\beta_n \rightarrow \min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$\alpha_n \rightarrow s.t. \quad y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n$$

$$-\nabla f(x^*) = A^T \lambda \rightarrow$$

$$w = \sum_n \alpha_n y_n \phi(x_n)$$

$$C = \alpha_n + \beta_n$$

$$\lambda \geq 0$$



$$\alpha_n \geq 0$$

$$\beta_n \geq 0$$

$$\lambda_n (a_n x - b_n) = 0 \rightarrow$$

$$\alpha_n (y_n w^T \phi(x_n) - 1 + \xi_n) \geq 0$$

$$\beta_n \xi_n \geq 0$$

Optimality Condition for SVM

What's the KKT condition for:

$$\beta_n \rightarrow \min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$\alpha_n \rightarrow s.t. \quad y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n$$

$$-\nabla f(x^*) = A^T \lambda \rightarrow$$

$$w = \sum_n \alpha_n y_n \phi(x_n)$$

$$\beta_n = C - \alpha_n$$

$$\lambda \geq 0 \rightarrow$$

$$\alpha_n \geq 0$$

$$(C - \alpha_n) \geq 0$$

$$\lambda_n (a_n x - b_n) = 0 \rightarrow$$

$$\alpha_n (y_n w^T \phi(x_n) - 1 + \xi_n) \geq 0$$

$$(C - \alpha_n) \xi_n \geq 0$$

Optimality Condition for SVM

What's the KKT condition for:

$$\begin{bmatrix} \nabla_w f(w, \xi) \\ \nabla_\xi f(w, \xi) \end{bmatrix} = \begin{bmatrix} w \\ C \end{bmatrix} = \begin{bmatrix} \sum_n \alpha_n y_n \phi(x_n) \\ \alpha_n + \beta_n \end{bmatrix}$$

$$\beta_n \rightarrow \min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$\alpha_n \rightarrow \text{s.t. } y_n w^T \phi(x_n) \geq 1 - \xi_n, \forall n$$

$$-\nabla f(x^*) = A^T \lambda \rightarrow$$

$$w = \sum_n \alpha_n y_n \phi(x_n)$$

$$\lambda \geq 0$$



$$0 \leq \alpha_n \leq C$$

$$\lambda_n (a_n x - b_n) = 0 \rightarrow$$

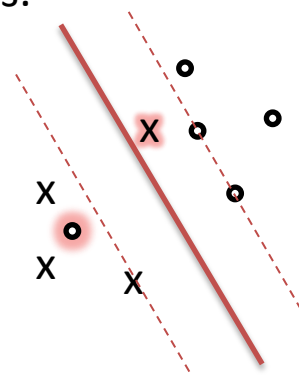
$$\begin{aligned} \alpha_n (y_n w^T \phi(x_n) - 1 + \xi_n) &= 0 \\ (C - \alpha_n) \xi_n &= 0 \end{aligned}$$

1. $w = \sum_n \alpha_n y_n \phi(x_n)$ can be expressed as linear combination of instances.

2. If constraint $y_n w^T \phi(x_n) \geq 1 - \xi_n$ not binding $\rightarrow \alpha_n = 0$

3. If $\alpha_n > 0 \rightarrow$ constraint is binding (**Support Vectors !**)

4. If loss of n-th instance $\xi_n > 0 \rightarrow \alpha_n = C$



Overview

- **Support Vector Machine**
 - The Art of Modeling --- Large Margin and Kernel Trick
 - Convex Analysis
 - Optimality Conditions
 - Duality
- **Optimization for Machine Learning**
 - Dual Coordinate Descent (fast convergence, moderate cost)
 - libLinear (Stochastic)
 - libSVM (Greedy)
 - **Primal Methods**
 - **Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)**
 - **Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)**

Primal SVM Problem

$$\min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$s.t. \quad \underline{y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n} \quad (\text{Let } D:\text{\#feature}, N:\text{\#samples})$$

Quadratic Program (QP) with:

- $D + N$ variables
- N Linear constraints
- N nonnegative constraints

➔ Intractable for median scale
(ex. $N=1000, D=1000$)

Primal SVM Problem

Constrained Problem → Non-smooth Unconstrained

$$\min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$s.t. \quad y_n f(x_n) \geq 1 - \xi_n, \quad \forall n$$

- Every constrained problem can be transformed to a Nonsmooth Unconstrained Problem

Given w , minimize w.r.t. ξ_n

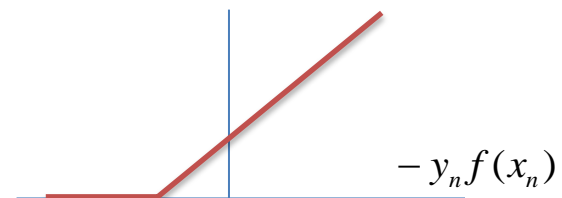
$$\xi_n = \begin{cases} 0 & \text{if } 1 - y_n f(x_n) \leq 0 \\ 1 - y_n f(x_n), & \text{otherwise} \end{cases}$$

$$\rightarrow \quad \xi_n = \max\{ 1 - y_n f(x_n), 0 \}$$

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_n L(f(x_n), y_n)$$

(Nonsmooth, Unconstrained)

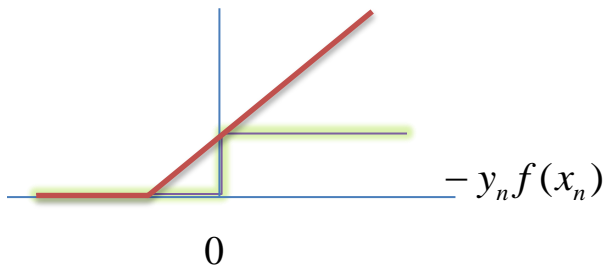
Hinge-Loss $L(\cdot)$



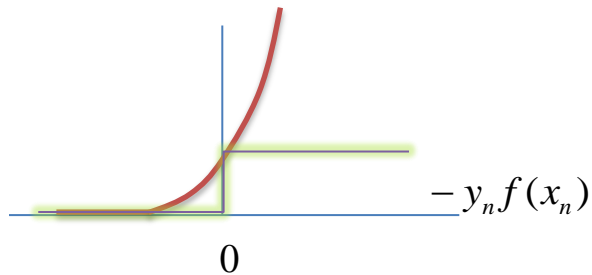
L2-Regularized Loss Minimization

$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_n L(f(x_n), y_n)$$

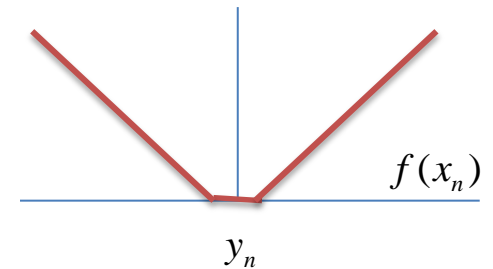
Hinge-Loss (L1-SVM, Structural SVM)



L2-SVM



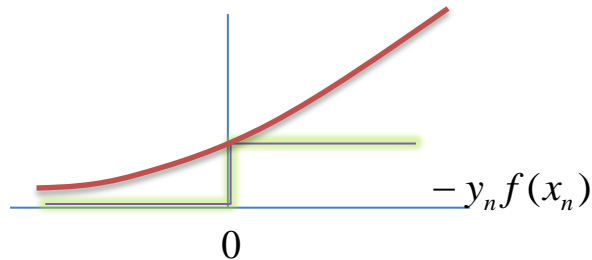
SVR



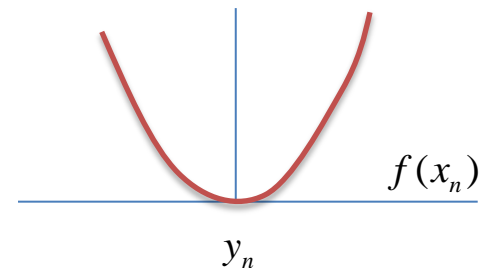
0/1-Loss (Accuracy)



Logistic Regression (CRF)



(Least-Square) Regression



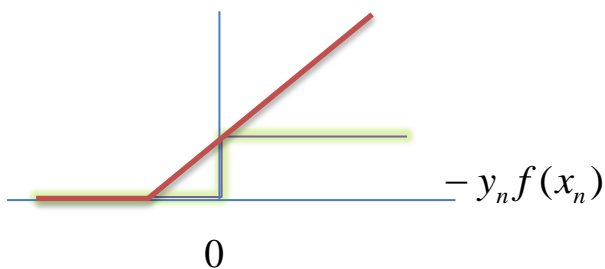
L2-Regularized Loss Minimization

$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_n L(f(x_n), y_n)$$

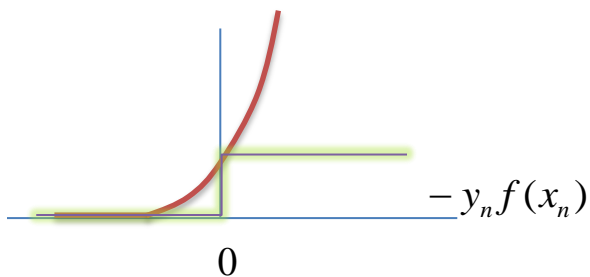
Convex Loss

- Solve with Global Minimum

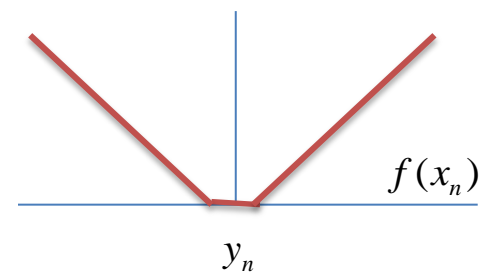
Hinge-Loss (L1-SVM, Structural SVM)



L2-SVM



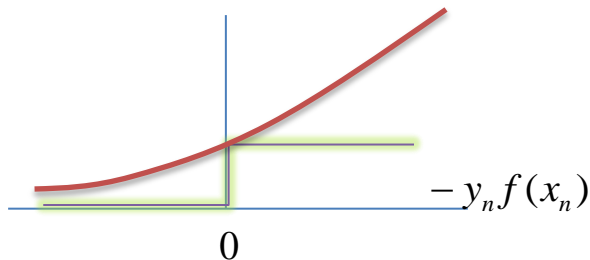
SVR



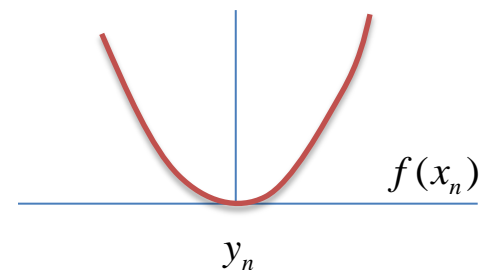
0/1-Loss (Accuracy)



Logistic Regression (CRF)



(Least-Square) Regression



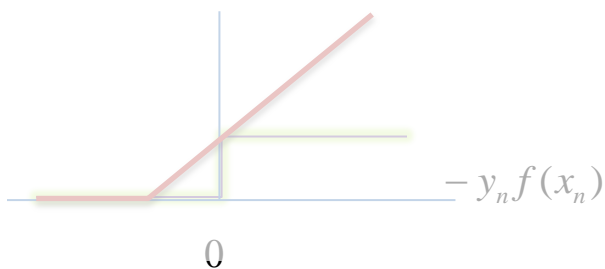
L2-Regularized Loss Minimization

$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_n L(f(x_n), y_n)$$

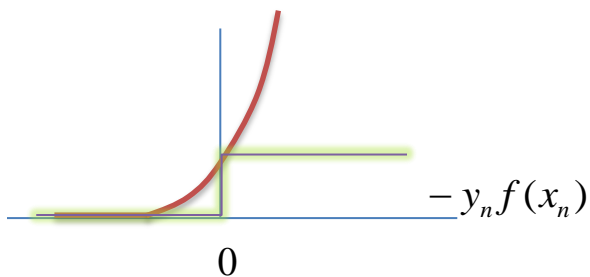
Convex Smooth Loss

- Applicable for Second-Order Method
- Coordinate Descent (primal)
- Gradient Descent $\rightarrow O(\log(1/\epsilon))$ rate
- Non-smooth $\rightarrow O(1/\epsilon)$ rate

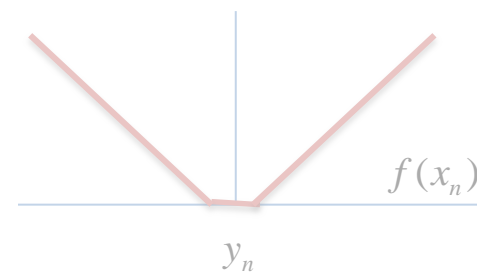
Hinge-Loss (L1-SVM, Structural SVM)



L2-SVM



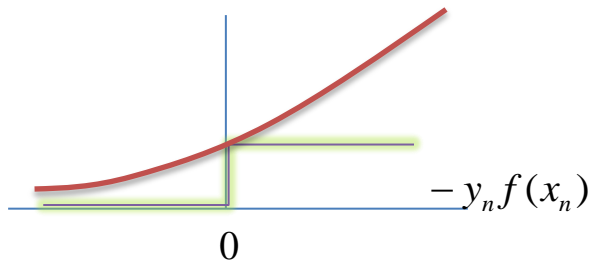
SVR



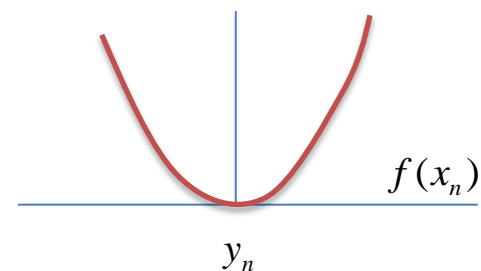
0/1-Loss (Accuracy)



Logistic Regression (CRF)



(Least-Square) Regression

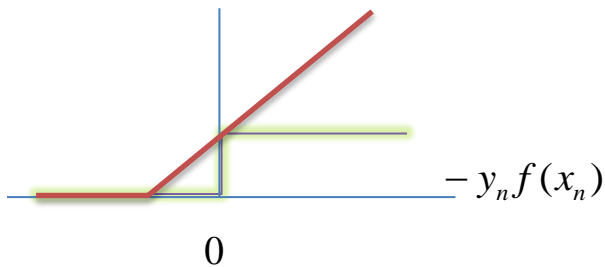


L2-Regularized Loss Minimization

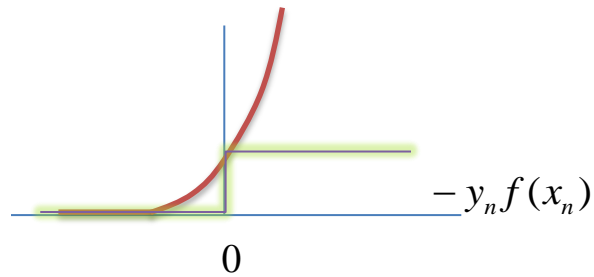
$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_n L(f(x_n), y_n)$$

Dual Sparsity

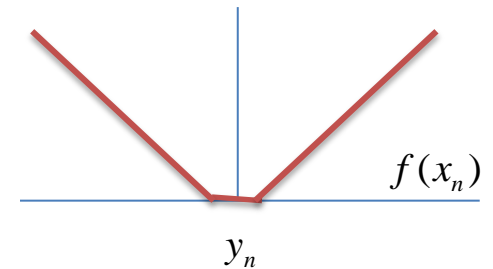
Hinge-Loss (L1-SVM, Structural SVM)



L2-SVM



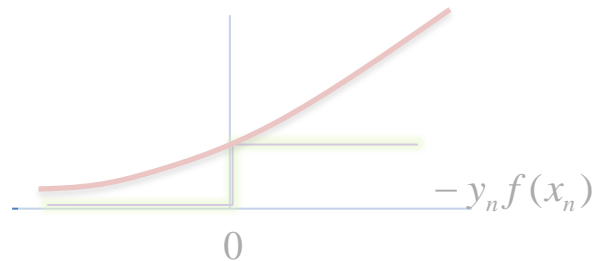
SVR



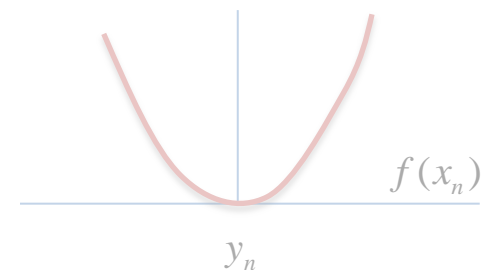
0/1-Loss (Accuracy)



Logistic Regression (CRF)



(Least-Square) Regression

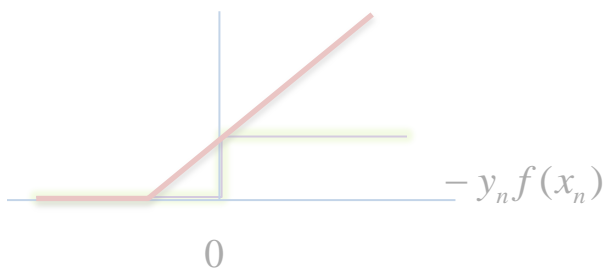


L2-Regularized Loss Minimization

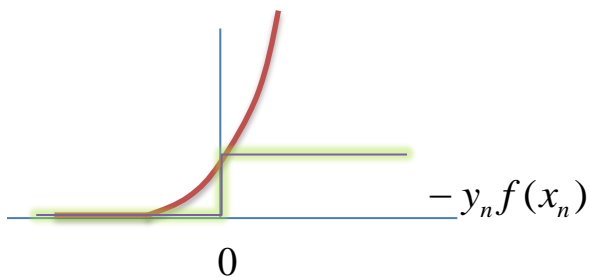
$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_n L(f(x_n), y_n)$$

Noise-Sensitive

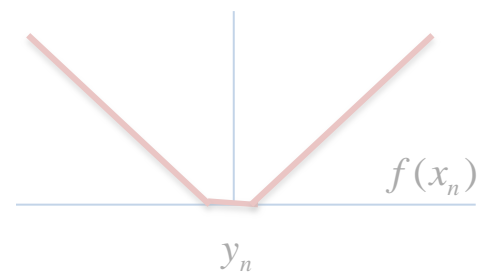
Hinge-Loss (L1-SVM, Structural SVM)



L2-SVM



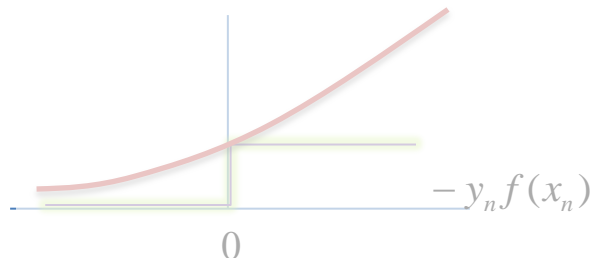
SVR



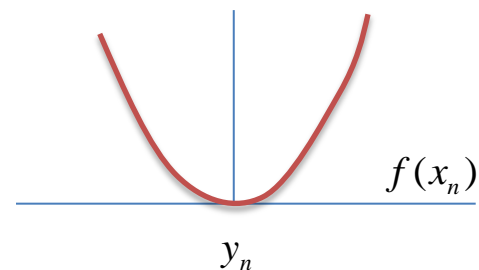
0/1-Loss (Accuracy)



Logistic Regression (CRF)



(Least-Square) Regression



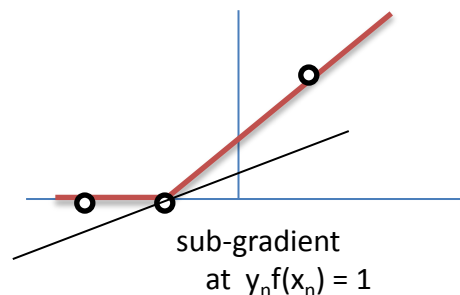
Most insensitive

Stochastic (sub-)Gradient Descent

(S. Shalev-Shwartz et al., ICML 2007)

$$\min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{N} \sum_n L(f(x_n), y_n)$$

Hinge-Loss (L1-SVM, Structural SVM)



iteration cost: $O(N \cdot D)$

Algorithm: Subgradient Descent

For $t = 1 \dots T$

$$w^{(t+1)} = w^{(t)} - \eta_t \left(\lambda w^{(t)} + \frac{1}{N} \sum_n L'(n) \phi_n \right)$$

End

A common choice: $\eta_t = \frac{1}{t}$

iteration cost: $O(D)$

Algorithm: Stochastic Subgradient Descent

For $t = 1 \dots T$

Draw \tilde{n} from uniformly from $\{1 \dots N\}$

$$w^{(t+1)} = w^{(t)} - \eta_t \left(\lambda w^{(t)} + L'(\tilde{n}) \phi_{\tilde{n}} \right)$$

End

$$\bar{w}^{(k)} := \frac{2}{k(k+1)} \sum_{t=1}^k t w^{(t)}$$

(avg. over iterations \rightarrow much faster)

(Shamir, O. and Zhang, ICML, 2013)

Stochastic (sub-)Gradient Descent

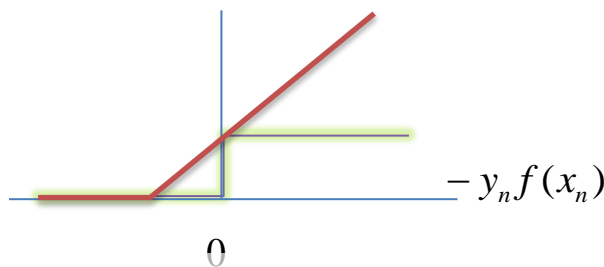
(S. Shalev-Shwartz et al., ICML 2007)

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_n L(f(x_n), y_n)$$

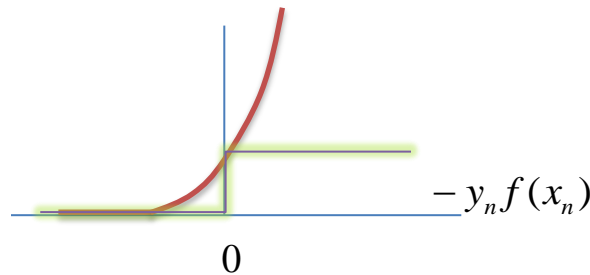
SGD

- Applicable to all.
- Non-Smooth \rightarrow GD: $O(1/\epsilon)$, SGD: $O(1/\epsilon)$
- Smooth \rightarrow GD: $O(\log 1/\epsilon)$, SGD: $O(1/\epsilon)$

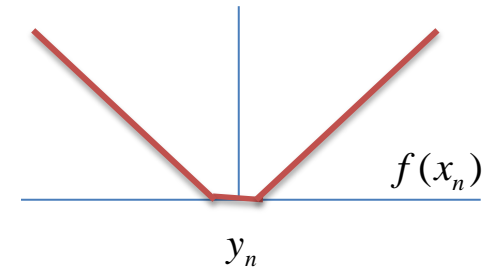
Hinge-Loss (L1-SVM, Structural SVM)



L2-SVM



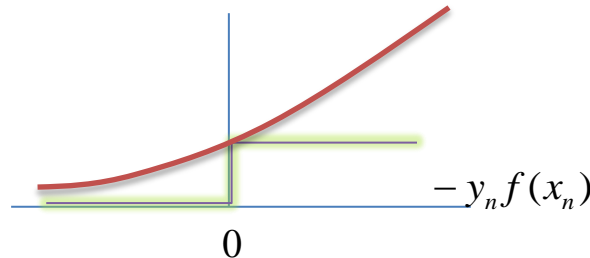
SVR



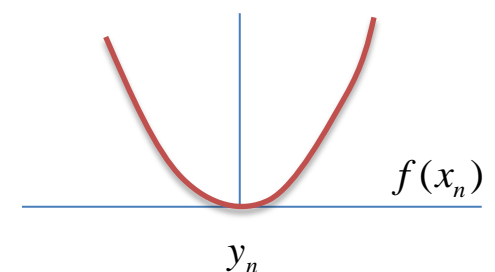
0/1-Loss (Accuracy)



Logistic Regression (CRF)

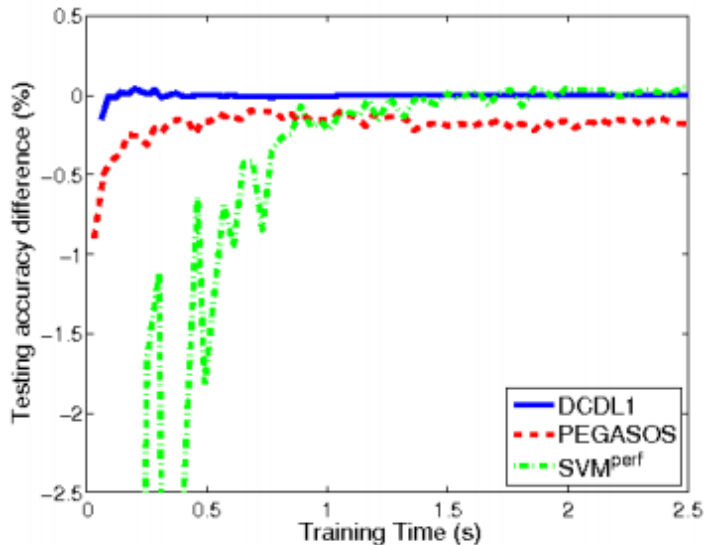


(Least-Square) Regression



SGD (Pegasos) vs. Batch Method (LibLinear)

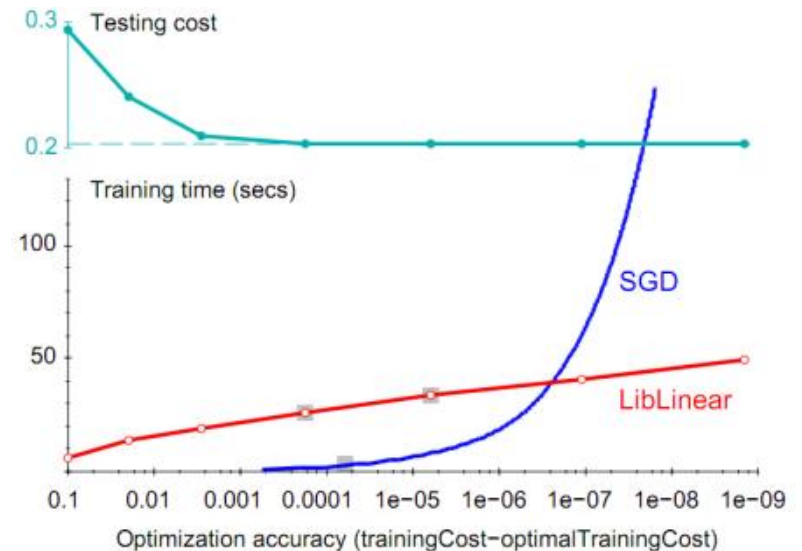
- Cons: SGD converges very slowly.
(sometimes seems not convergent....)



(c) L1-SVM: real-sim

(Heish, ICML 2008) (LibLinear)

- Pros: SGD (online method) has same convergence rate for Testing and Training.



Bottou, Léon. 2007. [Learning with Large Scale Datasets](#). *NIPS Tutorial*.

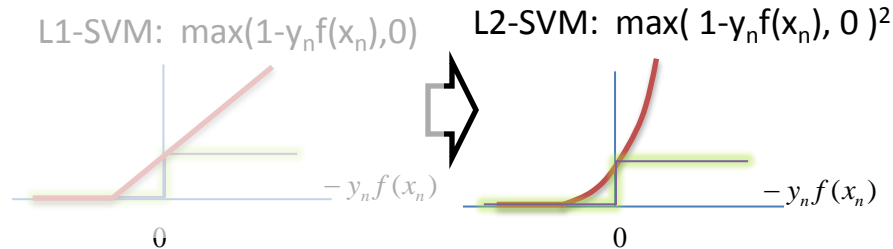
- Do you care a “**Ratio Improvement**” or “**Absolute Improvement**” in Testing ?
- What’s your evaluation measure ? (AUC, Prec/Recall, Accuracy....)
- ill-conditioned problems (pos/neg ratio, Large C)

Overview

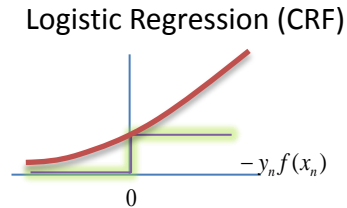
- **Support Vector Machine**
 - The Art of Modeling --- Large Margin and Kernel Trick
 - Convex Analysis
 - Optimality Conditions
 - Duality
- **Optimization for Machine Learning**
 - Dual Coordinate Descent (fast convergence, moderate cost)
 - libLinear (Stochastic)
 - libSVM (Greedy)
 - Primal Methods
 - Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)
 - **Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)**
 - L1-regularized
 - Primal Coordinate Descent

Smooth Loss vs. Non-smooth Loss

$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_n L(f(x_n), y_n)$$



➔ Unconstrained Differentiable Problem.



Usage: `train [options] training_set_file [model_file]`

options:

-s type : set type of solver (default 1)

for multi-class classification

0 -- L2-regularized logistic regression (primal)

1 -- L2-regularized L2-loss support vector classification (dual)

2 -- L2-regularized L2-loss support vector classification (primal)

3 -- L2-regularized L1-loss support vector classification (dual)

4 -- support vector classification by Crammer and Singer

5 -- L1-regularized L2-loss support vector classification

6 -- L1-regularized logistic regression

7 -- L2-regularized logistic regression (dual)

for regression

11 -- L2-regularized L2-loss support vector regression (primal)

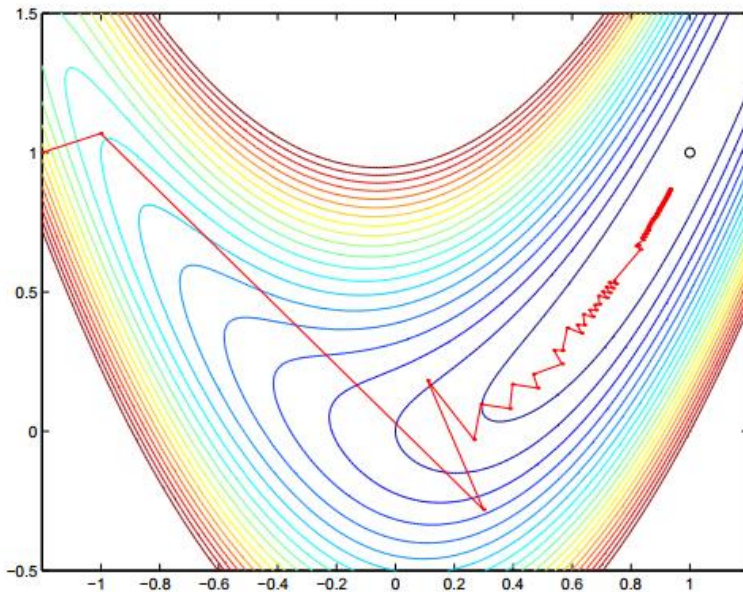
12 -- L2-regularized L2-loss support vector regression (dual)

13 -- L2-regularized L1-loss support vector regression (dual)

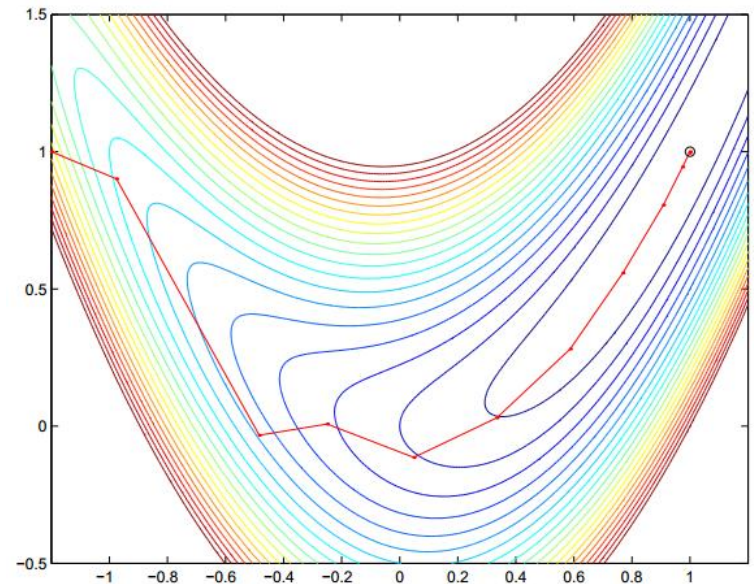
Primal Quasi-Newton Method

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_n L(f(x_n), y_n)$$

- Gradient Descent (1st order) uses Linear Approximation by $\nabla f(w) = g$
- Newton Method (2nd order) uses Quadratic Approximation by $\nabla f(w) = g$ and $\nabla^2 f(w) = H$



Gradient Descent



Newton Method

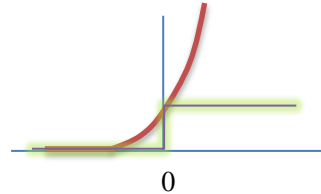
Primal Quasi-Newton Method

$$\min_w f(w) = \frac{1}{2} \|w\|^2 + C \sum_n L(w^T x_n, y_n)$$

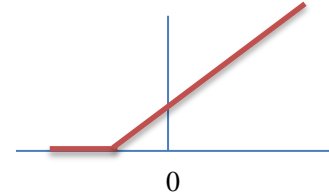
$$g = \nabla f(w) = w + C \sum_n L'(n) x_n$$

$$H = \nabla^2 f(w) = I + C \sum_n L''(n) x_n x_n^T$$

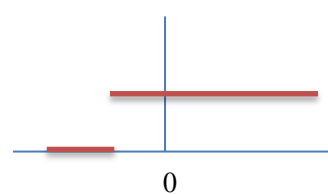
$$\max(1 - y_n f(x_n), 0)^2$$



$$2\max(1 - y_n f(x_n), 0)$$



$$0 \quad 2$$



Quadratic Approximation at $w^{(t)}$: $\min_{s=w-w^{(t)}} \frac{1}{2} s^T H s + g^T s + f(w^{(t)})$

Minimum at s^* : $H s^* = -g$

iteration cost: $O(N \cdot D^2 + D^3)$

Algorithm: Newton Method

For t = 1...T

$$\text{Solve } H^{(t)} s = -g^{(t)}$$

$$w^{(t+1)} = w^{(t)} + \eta_t s^*$$

End

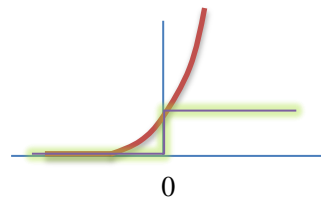
Primal Quasi-Newton Method

$$\min_w f(w) = \frac{1}{2} \|w\|^2 + C \sum_n L(w^T x_n, y_n)$$

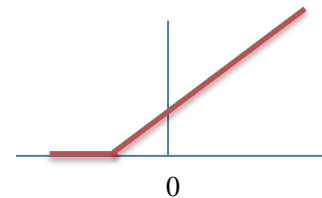
$$g = \nabla f(w) = w + C \sum_n L'(n) x_n$$

$$H = \nabla^2 f(w) = I + C \sum_n L''(n) x_n x_n^T$$

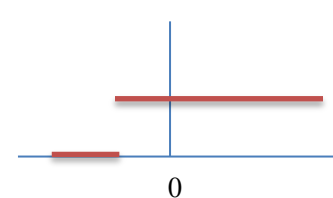
$$\max(1 - y_n f(x_n), 0)^2$$



$$2\max(1 - y_n f(x_n), 0)$$



$$0 \quad 2$$



Quadratic Approximation at $w^{(t)}$: $\min_{s=w-w^{(t)}} \frac{1}{2} s^T H s + g^T s + f(w^{(t)})$

Minimum at s^* : $H s^* + g = 0$

Iteration cost:

$$O(N \cdot D + |SV| \cdot D \cdot |T_{\text{inner}}|)$$

Algorithm: Conjugate Gradient for $Ax = b$.

For $t = 1 \dots T_{\text{inner}}$

$$r^{(t)} = b - Ax^{(t)}$$

$$d^{(t+1)} = d^{(t)} + \eta_t r^{(t)}$$

$$x^{(t+1)} = x^{(t)} - \eta'_t d^{(t)}$$

End

Algorithm: Quasi-Newton Method

For $t = 1 \dots T$

Solve $H^{(t)} s = -g^{(t)}$ approximately.

$$w^{(t+1)} = w^{(t)} + \eta_t s^*$$

End

Overview

- **Support Vector Machine**

- The Art of Modeling --- Large Margin and Kernel Trick
- Convex Analysis
- Optimality Conditions
- **Duality**

- **Optimization for Machine Learning**

- Dual Coordinate Descent (fast convergence, moderate cost)
 - libLinear (Stochastic)
 - libSVM (Greedy)
- Primal Methods
 - Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)
 - Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)

Lagrangian Duality

First, we consider the **Equality Constrained Problem**:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b} \end{aligned}$$

The **optimal solution** \mathbf{x}^* is found iff:

$$A\mathbf{x}^* = \mathbf{b}$$

$$-\nabla f(\mathbf{x}^*) = A^T \boldsymbol{\lambda}^*$$

If we define **Lagrangian Function** (Lagrangian) as:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T (A\mathbf{x} - \mathbf{b})$$

Then the **optimality condition** can be written as:

$$\frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} = 0 \Rightarrow A\mathbf{x}^* = \mathbf{b} \quad (\boldsymbol{\lambda} \text{ cannot increase } L(.))$$

$$\frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial \mathbf{x}} = 0 \Rightarrow -\nabla f(\mathbf{x}^*) = A^T \boldsymbol{\lambda}^*$$

(\mathbf{x} cannot decrease $L(.)$)

$$\min_{\mathbf{x}} \left\{ \max_{\boldsymbol{\lambda}} L(\mathbf{x}, \boldsymbol{\lambda}) \right\}$$

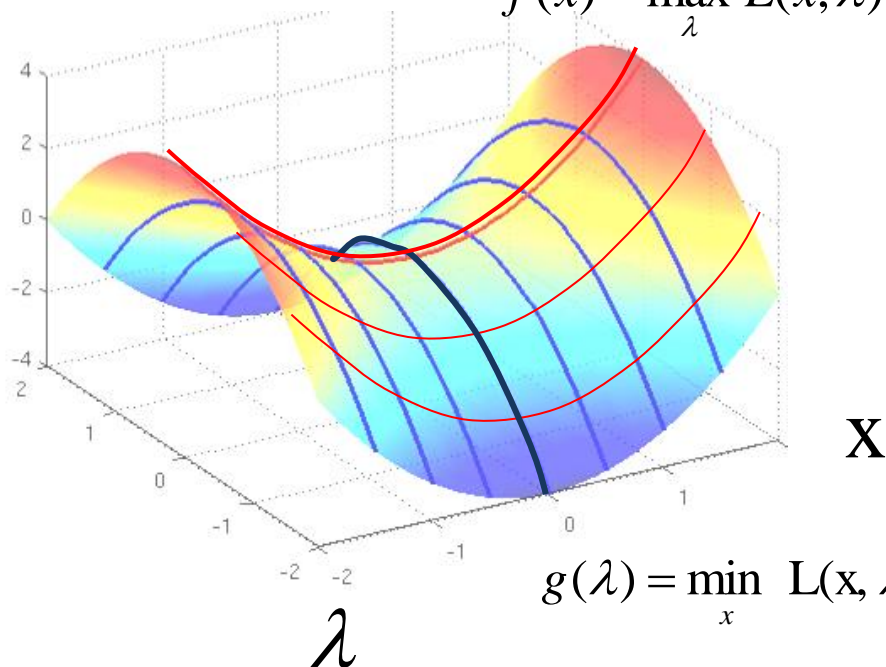
$$\max_{\boldsymbol{\lambda}} \left\{ \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) \right\}$$

Lagrangian Duality

If we define **Lagrangian Function** (Lagrangian) as:

$$L(x, \lambda) = f(x) + \lambda^T (Ax - b)$$

$$L(x, \lambda)$$



Every point satisfies

$$\frac{\partial L(x, \lambda)}{\partial \lambda} = 0 \Rightarrow Ax^* = b$$

Every point satisfies

$$\frac{\partial L(x, \lambda)}{\partial x} = 0 \Rightarrow -\nabla f(x^*) = A^T \lambda^*$$

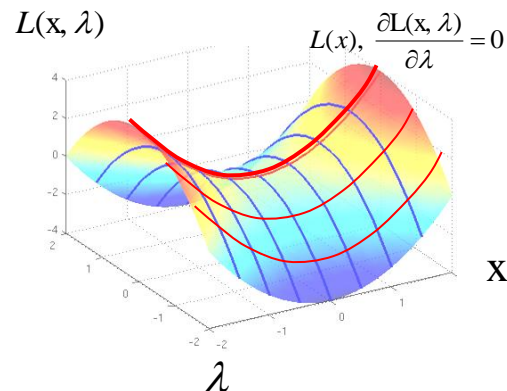
Lagrangian Duality

Original (Primal) problem is:

$$\begin{aligned} \min_{\mathbf{x}} \quad & L(\mathbf{x}, \lambda) = f(x) + \lambda^T (A\mathbf{x} - b) \\ \text{s.t.} \quad & \frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda} = A\mathbf{x} - b = 0 \end{aligned}$$

Primal problem is:

$$\min_{\mathbf{x}} \max_{\lambda} L(\mathbf{x}, \lambda) = f(x) + \lambda^T (A\mathbf{x} - b)$$

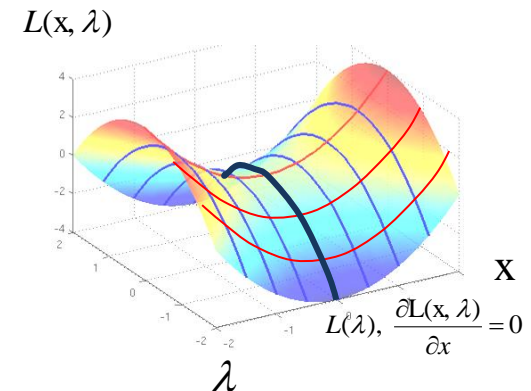


Dual Problem:

$$\begin{aligned} \max_{\lambda} \quad & L(\mathbf{x}, \lambda) = f(x) + \lambda^T (A\mathbf{x} - b) \\ \text{s.t.} \quad & \frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} = \nabla f(x) + A^T \lambda = 0 \end{aligned}$$

Dual problem is:

$$\max_{\lambda} \min_{\mathbf{x}} L(\mathbf{x}, \lambda) = f(x) + \lambda^T (A\mathbf{x} - b)$$



Lagrangian Duality

For **Inequality Constrained Problem**:

$$\begin{aligned} \min_x & f(x) \\ \text{s.t.} & Ax \leq b \end{aligned}$$

Primal problem is:

$$\min_x \max_{\lambda \geq 0} L(x, \lambda) = f(x) + \lambda^T (Ax - b)$$

$$\max_{\lambda \geq 0} L(x, \lambda) = \begin{cases} f(x), & Ax - b \leq 0 \\ \infty, & Ax - b > 0 \end{cases}$$

Dual problem is:

$$\max_{\lambda \geq 0} \min_x L(x, \lambda) = f(x) + \lambda^T (Ax - b)$$

$$\begin{aligned} \min_x & L(x, \lambda) \\ \Rightarrow & \nabla f(x) + A^T \lambda = 0 \end{aligned}$$

SVM Dual Problem

Primal Problem:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned} \quad \longleftrightarrow \quad \min_{w, \xi} \max_{\alpha \geq 0, \beta \geq 0} L(w, \xi, \alpha, \beta)$$

Lagrangian:

$$L(w, \xi, \alpha, \beta) = \left(\frac{1}{2} \|w\|^2 + C \sum_n \xi_n \right) - \sum_n \alpha_n (y_n w^T \phi(x_n) - 1 + \xi_n) - \sum_n \beta_n \xi_n$$

Dual Problem:

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{w, \xi} L(w, \xi, \alpha, \beta)$$

SVM Dual Problem

Primal Problem:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned} \quad \longleftrightarrow \quad \begin{aligned} \min_{w, \xi} \quad & \max_{\alpha \geq 0, \beta \geq 0} L(w, \xi, \alpha, \beta) \end{aligned}$$

Lagrangian:

$$L(w, \xi, \alpha, \beta) = \left\{ \frac{1}{2} \|w\|^2 - \sum_n \alpha_n y_n w^T \phi(x_n) \right\} + \left\{ \sum_n (C - \alpha_n - \beta_n) \xi_n \right\} + \sum_n \alpha_n$$

Dual Problem:

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{w, \xi} L(w, \xi, \alpha, \beta)$$

\longleftrightarrow

$$\begin{aligned} \max_{\alpha \geq 0, \beta \geq 0} \quad & L(w, \xi, \alpha, \beta) \quad \begin{array}{c} \left[\begin{array}{ccc} y_1 \phi(x_1) & \dots & y_N \phi(x_N) \end{array} \right] \begin{bmatrix} \alpha_1 \\ \dots \\ \alpha_N \end{bmatrix} \end{array} \\ \text{s.t.} \quad & \frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_n \alpha_n y_n \phi(x_n) = \Phi \alpha \\ & \frac{\partial L}{\partial \xi} = 0 \Rightarrow C = \alpha_n + \beta_n \end{aligned}$$

SVM Dual Problem

Primal Problem:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned} \quad \longleftrightarrow \quad \begin{aligned} \min_{w, \xi} \quad & \max_{\alpha \geq 0, \beta \geq 0} L(w, \xi, \alpha, \beta) \end{aligned}$$

Lagrangian:

$$L(\alpha, \beta) = \left\{ \frac{1}{2} \alpha^T \Phi^T \Phi \alpha - \alpha^T \Phi^T \Phi \alpha \right\} + \{0\} + \sum_n \alpha_n$$

Dual Problem:

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{w, \xi} L(w, \xi, \alpha, \beta)$$

\longleftrightarrow

$$\begin{aligned} \max_{\alpha \geq 0, \beta \geq 0} \quad & L(w, \xi, \alpha, \beta) \quad \begin{bmatrix} y_1 \phi(x_1) & \dots & y_N \phi(x_N) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \dots \\ \alpha_N \end{bmatrix} \\ \text{s.t.} \quad & \frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_n \alpha_n y_n \phi(x_n) = \Phi \alpha \\ & \frac{\partial L}{\partial \xi} = 0 \Rightarrow C = \alpha_n + \beta_n \end{aligned}$$

SVM Dual Problem

Primal Problem:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned} \quad \longleftrightarrow \quad \min_{w, \xi} \max_{\alpha \geq 0, \beta \geq 0} L(w, \xi, \alpha, \beta)$$

Lagrangian:

$$L(\alpha, \beta) = \left\{ \frac{1}{2} \alpha^T \Phi^T \Phi \alpha - \alpha^T \Phi^T \Phi \alpha \right\} + \{0\} + \sum_n \alpha_n$$

Dual Problem:

$$\begin{aligned} \max_{\alpha \geq 0, \beta \geq 0} \min_{w, \xi} L(w, \xi, \alpha, \beta) \quad & \longleftrightarrow \quad \max_{\alpha \geq 0, \beta \geq 0} L(\alpha, \beta) = \sum_n \alpha_n - \frac{1}{2} \alpha^T \Phi^T \Phi \alpha \\ \text{s.t.} \quad & C = \alpha_n + \beta_n \end{aligned}$$

SVM Dual Problem

Primal Problem:

$$\begin{aligned} \min_{w, \xi \geq 0} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n \end{aligned} \quad \longleftrightarrow \quad \min_{w, \xi} \max_{\alpha \geq 0, \beta \geq 0} L(w, \xi, \alpha, \beta)$$

Lagrangian:

$$L(\alpha, \beta) = \left\{ \frac{1}{2} \alpha^T \Phi^T \Phi \alpha - \alpha^T \Phi^T \Phi \alpha \right\} + \{0\} + \sum_n \alpha_n$$

Dual Problem:

$$\begin{aligned} \max_{\alpha \geq 0, \beta \geq 0} \min_{w, \xi} L(w, \xi, \alpha, \beta) \quad & \longleftrightarrow \quad \max_{\alpha} L(\alpha, \beta) = \sum_n \alpha_n - \frac{1}{2} \alpha^T \Phi^T \Phi \alpha \\ \text{s.t.} \quad & 0 \leq \alpha \leq C \end{aligned}$$

SVM Dual Problem

Dual Problem (only involve product $\phi(x_i)^T \phi(x_j)$) :

$$\max_{\alpha} \sum_n \alpha_n - \frac{1}{2} \alpha^T \Phi^T \Phi \alpha$$

$$s.t. \ 0 \leq \alpha \leq C$$

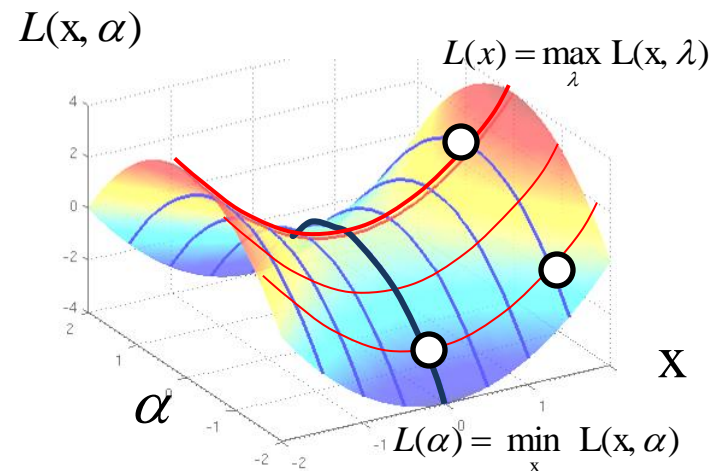
\leftrightarrow

$$\max_{\alpha} \sum_n \alpha_n - \frac{1}{2} \alpha^T Q \alpha$$

$$s.t. \ 0 \leq \alpha \leq C$$

$$Q_{ij} = (y_i \phi(x_i))(y_j \phi(x_j)) = y_i y_j K(x_i, x_j)$$

1. Only “Box Constraint” \Rightarrow Easy to solve.
2. $\dim(\alpha) = N = |\text{instance}|$, $\dim(w) = D = |\text{features}|$
3. Weak Duality: $\text{Dual}(\alpha) \leq \text{Primal}(w)$
4. Strong Duality: $\text{Dual}(\alpha^*) = \text{Primal}(w^*)$
(if primal is convex)



Overview

- **Support Vector Machine**
 - The Art of Modeling --- Large Margin and Kernel Trick
 - Convex Analysis
 - Optimality Conditions
 - Duality
- **Optimization for Machine Learning**
 - Dual Coordinate Descent (DCD) (fast convergence, moderate cost)
 - libLinear (Stochastic CD)
 - libSVM (Greedy CD)
 - Primal Methods
 - Non-smooth Loss → Stochastic Gradient Descent (slow convergence, cheap iter.)
 - Differentiable Loss → Quasi-Newton Method (very fast convergence, expensive iter.)

Dual Optimization of SVM

$$\begin{aligned} \max_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \alpha^T Q \alpha \\ \text{s.t.} \quad & 0 \leq \alpha \leq C \end{aligned} \quad \rightarrow \quad \begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \sum_n \alpha_n \\ \text{s.t.} \quad & 0 \leq \alpha \leq C \end{aligned}$$

$$Q_{ij} = (y_i \phi(x_i))(y_j \phi(x_j)) = y_i y_j K(x_i, x_j)$$

```
Usage: train [options] training_set_file [model_file]
options:
-s type : set type of solver (default 1)
  for multi-class classification
    0 -- L2-regularized logistic regression (primal)
    1 -- L2-regularized L2-loss support vector classification (dual)
    2 -- L2-regularized L2-loss support vector classification (primal)
    3 -- L2-regularized L1-loss support vector classification (dual)
    4 -- support vector classification by Crammer and Singer
    5 -- L1-regularized L2-loss support vector classification
    6 -- L1-regularized logistic regression
    7 -- L2-regularized logistic regression (dual) ?!
  for regression
    11 -- L2-regularized L2-loss support vector regression (primal)
    12 -- L2-regularized L2-loss support vector regression (dual)
    13 -- L2-regularized L1-loss support vector regression (dual)
```


Constrained Minimization

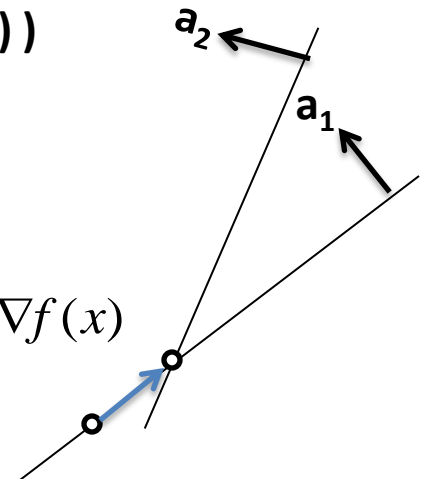
$$\min_{w, \xi \geq 0} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$s.t. \quad y_n w^T \phi(x_n) \geq 1 - \xi_n, \quad \forall n$$

General Constraint \rightarrow Very Expensive:

1. Detecting binding constraint : **$O(|\text{constraint}| * \text{dim}(\alpha))$**
2. Compute “Projected Gradient” : **$O(|\text{binding constraint}| * \text{dim}(\alpha))$**

$$-\nabla^P f(x) = -Z_1^T \nabla f(x)$$



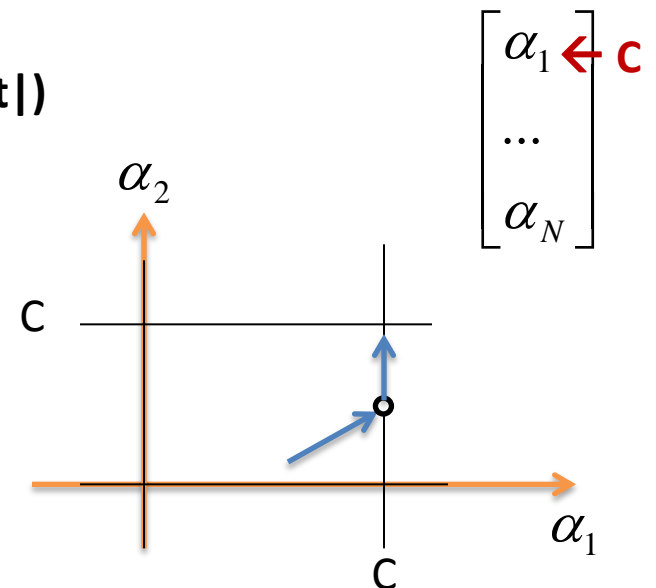
Constrained Minimization for “Box Constraint”

$$\min_{\alpha} \quad \frac{1}{2} \alpha^T Q \alpha - \sum_n \alpha_n$$

$s.t. \quad 0 \leq \alpha \leq C$

Cheap:

1. Detecting binding constraint : **$O(|\text{constraint}|)$**
2. Compute “Projected Gradient” : **$O(|\text{binding constraint}|)$**



Dual Coordinate Descent

$$\min_{\alpha} \quad \frac{1}{2} \alpha^T Q \alpha - \sum_n \alpha_n \quad \text{Minimize w.r.t.} \quad \alpha_i$$

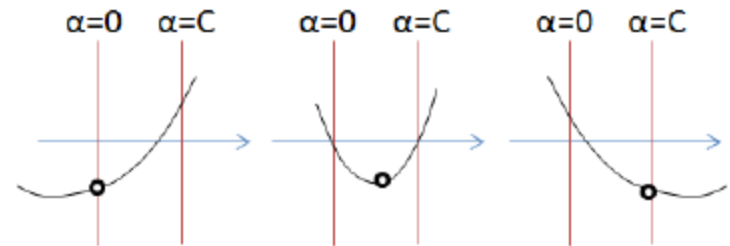
$s.t. \quad 0 \leq \alpha \leq C$ \rightarrow

$$\min_{\alpha_i} \quad \frac{1}{2} [\nabla^2_{ii} f(\alpha)] \alpha_i^2 + [\nabla_i f(\alpha)] \alpha_i + \text{const.}$$

$s.t. \quad 0 \leq \alpha_i \leq C$

$$\nabla^2_{ii} f(\alpha) = Q_{ii}$$

$$\nabla_i f(\alpha) = [Q\alpha - 1]_i$$



$$\alpha_i \leftarrow \min\left(\max\left(\alpha_i - \frac{\nabla f(\alpha)_i}{\nabla^2 f(\alpha)_{ii}}, 0\right), C\right),$$

Dual Optimization of SVM

$$\min_{\alpha} \quad \frac{1}{2} \alpha^T Q \alpha - \sum_n \alpha_n$$

$$s.t. \quad 0 \leq \alpha \leq C$$

Even Computing Gradient is Expensive: $\nabla f(\alpha) = Q_{N \times N} \alpha_{N \times 1} - 1 \quad (O(N^2))$

Coordinate Descent: (Optimize w.r.t. one variable at a time)

$$\nabla f(\alpha)_{(i)} = [Q]_{i,:} \alpha_{N \times 1} - 1 = \sum_k \alpha_k y_i y_k K(x_i, x_k) - 1 \quad (O(N))$$

How many variables ? \rightarrow As few as possible

- \swarrow Sequential Minimal Optimization (LibSVM)
(2 variable at a time)
- \searrow Coordinate Descent (LibLinear)
(1 variable at a time)

NonLinear (LibSVM) vs. Linear (LibLinear)

Linear:

$$\begin{aligned} \nabla f(\alpha)_{(i)} &= [Q]_{i,:} \alpha_{N \times 1} - 1 = \sum_k \alpha_k y_i y_k \boxed{K(x_i, x_k)} - 1 \\ &= \sum_k \alpha_k y_i y_k \boxed{(x_i^T x_k)} - 1 = y_i x_i^T \boxed{\left(\sum_k \alpha_k y_k x_k \right)} - 1 = y_i x_i^T w - 1 \end{aligned}$$

O(|non-zero Feature|)

Non-Linear:

$$\nabla f(\alpha)_{(i)} = [Q]_{i,:} \alpha_{N \times 1} - 1 = \sum_k \alpha_k y_i y_k K(x_i, x_k) - 1 \quad \mathbf{O(|Instances| * |non-zero Features|)}$$

NonLinear (LibSVM) vs. Linear (LibLinear)

Linear:

$$\begin{array}{c} \left[\begin{array}{c} \vdots \\ \nabla f(\alpha)_{(i)} \\ \vdots \end{array} \right] = y_i x_i^T w - 1 \\ \mathcal{O}(|\text{Feature}|) \end{array}$$
$$\begin{array}{c} \left[\begin{array}{c} \alpha_1 \\ \vdots \\ \alpha_N \end{array} \right] \rightarrow w = \sum_n \alpha_n y_n x_n \rightarrow \left[\begin{array}{c} \vdots \\ \nabla f(\alpha)_{(i)} \\ \vdots \end{array} \right] = y_i x_i^T w - 1 \end{array}$$

(Cheap Update \rightarrow Random Select Coordinate)

Non-Linear:

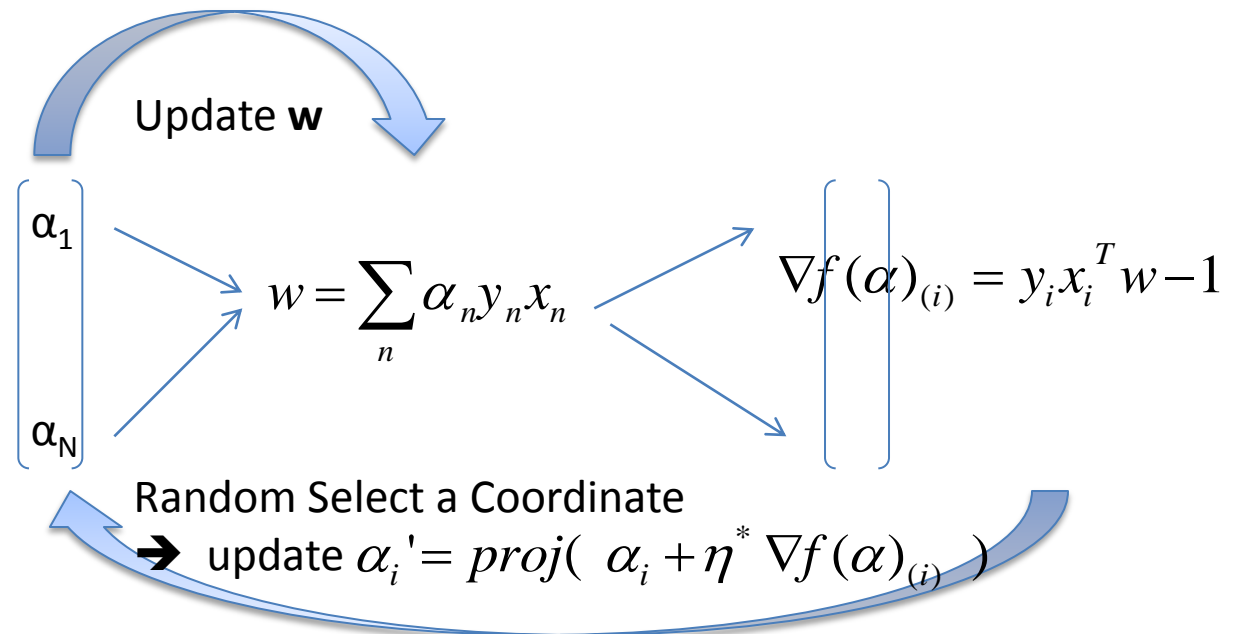
$$\begin{array}{c} \left[\begin{array}{c} \vdots \\ \nabla f(\alpha)_{(i)} \\ \vdots \end{array} \right] = \sum_k \alpha_k y_i y_k K(x_i, x_k) - 1 \\ \mathcal{O}(|\text{Instances}|) \end{array}$$
$$\begin{array}{c} \left[\begin{array}{c} \alpha_1 \\ \vdots \\ \alpha_N \end{array} \right] \rightarrow \left[\begin{array}{c} \vdots \\ \nabla f(\alpha)_{(i)} \\ \vdots \end{array} \right] \end{array}$$

(Expensive Update \rightarrow Select most Promising Coordinate)

LibLinear

Linear:

$$\left[\begin{array}{c} \nabla f(\alpha)_{(i)} = y_i x_i^T w - 1 \end{array} \right] \quad O(|\text{Feature}|)$$



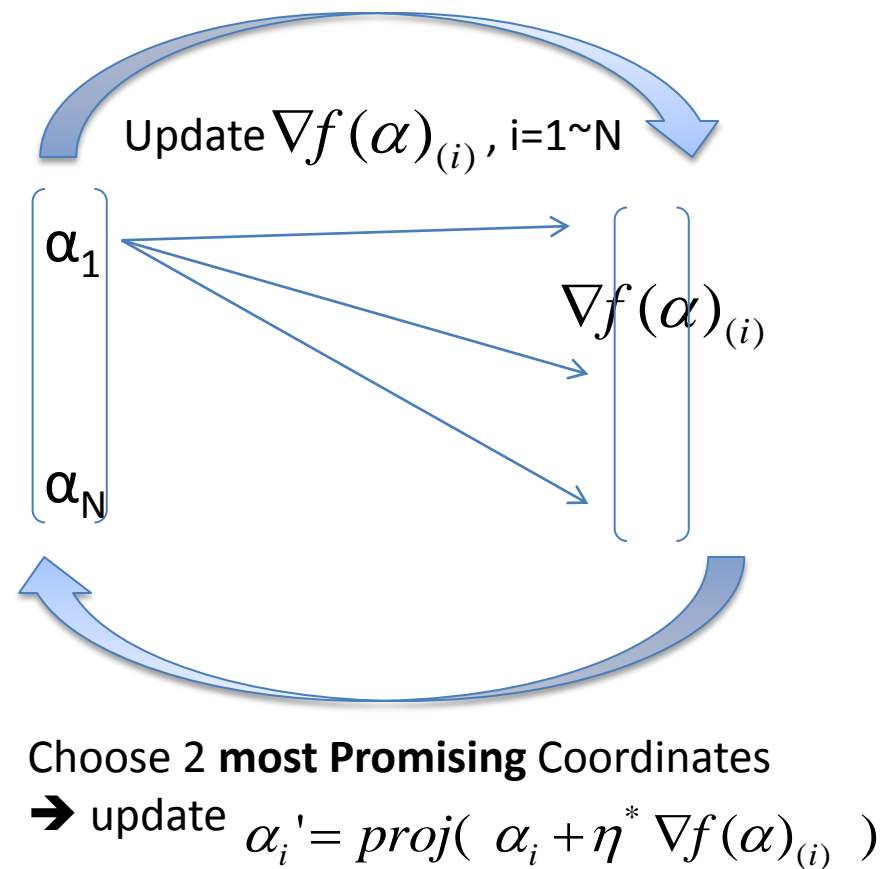
LibSVM

Non-Linear:

$$\nabla f(\alpha)_{(i)} = \sum_k \alpha_k y_i y_k K(x_i, x_k) - 1$$

$O(|\text{Instances}| * |\text{Features}|)$
(no cache)

$O(|\text{Instances}|)$
(cache)



Demo: libSVM, libLinear

— Normalize Features:

- `svm-scale -s [range_file] [train] > train.scale`
- `svm-scale -r [range_file] [test] > test.scale`

— Training:

- LibSVM: `svm-train [train.scale] (produce train.scale.model)`
- LibSVM: `svm-predict [test.scale] [train.scale.model] [pred_output]`
- LibLinear: `train [train.scale] (produce train.scale.model)`
- LibLinear: `predict [test.scale] [train.scale.model] [pred_output]`