

# Optimization in Machine Learning

Lecture 20: Project Gradient Descent, Karush Kuhn Tucker Conditions and Duality,  
Further results in Discrete Optimization

Ganesh Ramakrishnan

Department of Computer Science  
Dept of CSE, IIT Bombay  
<https://www.cse.iitb.ac.in/~ganesh>

April, 2025



# Outline of next few topics

- Gradient Descent Analysis for Lipschitz Smoothness/Continuity/(Strong) Convexity [Done]
  - ▶ Our running Colab Notebook:
- Nesterov's and Polyak's accelerated gradient descent [Done]
- Sub-gradient Descent and Analysis [Done]
- Stochastic Gradient Descent [Done]
- Generalized Gradient Descent [Next]
- Proximal Gradient Descent [Dpme]
- Projected Gradient Descent for Constrained Optimization [Partly done]
- Discrete/Combinatorial Optimization for Subset Selection
- Constrained Optimization, Duality & KKT Conditions
- And so on...more such algorithms



# Recap: Generalized Gradient Descent and its Special Cases

$$prox_{\gamma}(\mathbf{z}) = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|^2 + c(\mathbf{x})$$

It's special cases are:

- ① Gradient Descent:  $c(\mathbf{x}) = 0$
- ② Projected Gradient Descent:  $c(\mathbf{x}) = I_C(\mathbf{x})$  (Example:  $= \sum_i I_{g_i}(\mathbf{x})$ )
- ③ Alternating Projection/Proximal Minimization:  $f(\mathbf{x}) = 0$  How would you solve  
for the prox operator in general
- ④ Alternating Direction Method of Multipliers  
for projected gradient descent?
- ⑤ Special Cases for Specific Objectives
  - ▶ LASSO: (Fast) Iterative Shrinkage Thresholding Algorithm (ISTA/FISTA)  
(Hint: See following two slides)



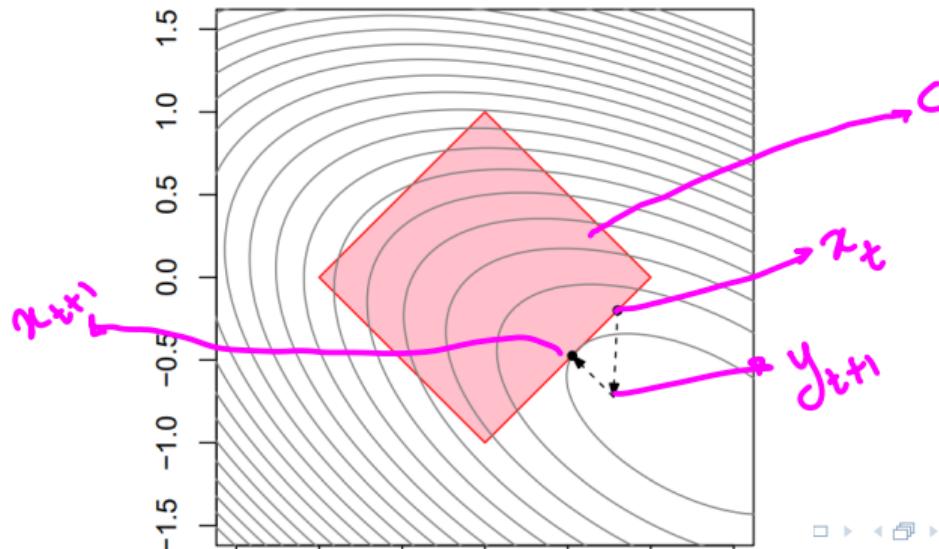
# Recap: Summary of Proximal Gradient Descent

- Proximal Gradient Descent becomes gradient descent if  $c = 0$
- If  $f$  is 0 (i.e. no smooth function), one can minimize a non-differentiable function  $c$  as long as the prox operator is easy to compute: **Proximal Point Algorithm**.
- Key to Proximal GD: Being able to compute Proximal operator.
- What if Prox can be efficiently computed approximately?
- There are papers (Schmidt et al 2011, Inexact Proximal Gradient Methods) where Prox is computed approximately but one can still derive the convergence rates if the errors due to approximation can be controlled!
- Accelerated Proximal GD: Similar to GD, one can accelerate GD to get optimal convergence rates! (Beck and Teboulle 2008)



# Recap: An important special case: Projected Gradient Descent

- Consider the Problem of Constrained Convex Minimization:  $\min_{x \in C} f(x)$
- A simple modification of the gradient descent procedure is:
  - At every iteration  $t$ : (Gradient Step): Compute  $y_{t+1} = x_t - \alpha \nabla f(x_t)$
  - (Projection step)  $x_{t+1} = P_C(y_{t+1})$
- Key here is the Projection step. Define  $P_C(x) = \operatorname{argmin}_{y \in C} \frac{1}{2} \|x - y\|^2$



# Recap: Projected Gradient Descent and Proximal Gradient Descent

- There is a close connection between Proximal and Projected Gradient Descent.
- Define  $c(x) = I(x \in \mathcal{C})$  where  $I(\cdot)$  is the Indicator function.
- It's easy to see that the  $\text{prox}_c(x) = P_{\mathcal{C}}(x)$ , i.e. the Prox operator is exactly the same as a projection operator.
- As a result, projected gradient descent becomes a special case of proximal gradient descent.
- Theoretical results of Proj. GD: All results for standard Gradient descent carry over to the projected case as long as the projection operator is easy to compute!



KKT conditions, Lagrange Dual etc



# Algorithm: Projected Gradient Descent (We use $\mathbf{x}_u^t$ instead of $\mathbf{z}^t$ )

**Find** a starting point  $\mathbf{x}_p^0 \in \mathcal{C}$ .

Set  $t = 1$

**repeat**

1. Choose a step size  $\gamma_t \propto 1/\sqrt{t}$ .
2. Set  $\mathbf{x}_u^t = \mathbf{x}_p^{t-1} - \gamma_t \nabla f(\mathbf{x}_p^{t-1})$ .
3. Set  $\mathbf{x}_p^t = \operatorname{argmin}_{\mathbf{z} \in \mathcal{C}} \|\mathbf{x}_u^t - \mathbf{z}\|_2^2$ .
4. Set  $t = t + 1$ .

**until** stopping criterion (such as  $\|\mathbf{x}_p^t - \mathbf{x}_p^{t-1}\| \leq \epsilon$  or  $f(\mathbf{x}_p^t) > f(\mathbf{x}_p^{t-1})$ ) is satisfied<sup>a</sup>

---

<sup>a</sup>Better criteria can be found using Lagrange duality theory, such as in the form of **duality gap**, etc.

Figure 1: The projected gradient descent algorithm.



# Recap: Computing the Projection Operator

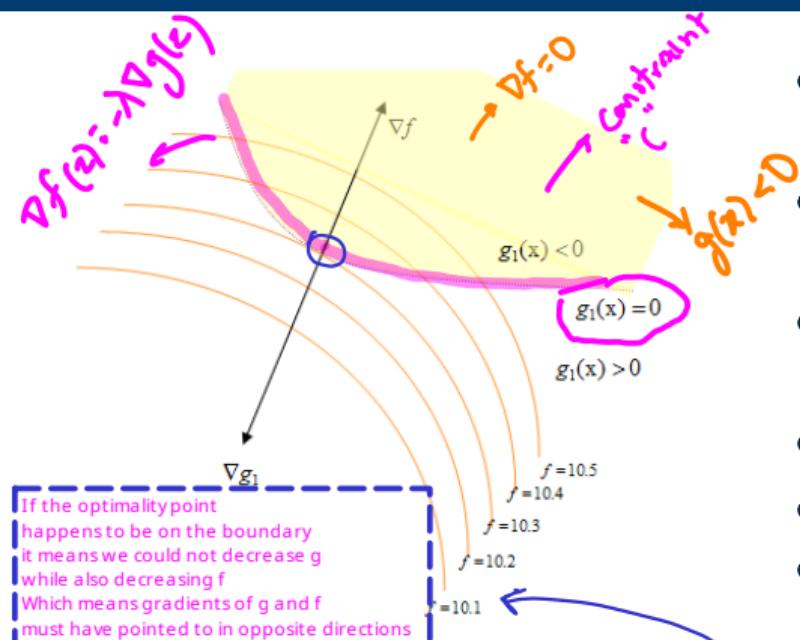


Figure 2: In this Figure  $c = 0$  and  $f$  is the function being optimized. See optional course material referred to at the end for understanding this figure completely

From projection perspective this is simply one iteration of projection

- Let's assume for simplicity that  $\mathcal{C} = \{x | g(x) \leq c\}$  and  $f = \frac{1}{2}\|z - x\|^2$  is the function being minimized
- Computing the projection step involves solving:  $\min_z \left\{ \frac{1}{2}\|z - x\|^2, \text{ such that } g(z) \leq c \right\}$ .
- Use the idea of Lagrange multipliers and the KKT conditions (see subsequent slides)!
- Define  $L(z, \lambda) = \frac{1}{2}\|z - x\|^2 + \lambda(g(z) - c)$ .
- Optimality conditions are:  $\nabla_z L = 0$  and  $\nabla_\lambda L = 0$ !
- There are two options.
  - Either  $x \in \mathcal{C}$ , and the constraints are not active ( $g(z) < c$ ), in which case we need  $\nabla_z L = 0$  and  $\lambda = 0$  (that is  $z = x$ ).  
if the optimal point is NOT in the boundary  
we have the regular condition of unconstrained optimization
  - Or  $x$  is on the boundary of  $\mathcal{C}$ , in which case  $g(z) = c$  and  $z - x + \lambda \nabla g(z) = 0$ .  
if the optimal point happens to be on the boundary
- If both can be solved in closed form, we are done!



# Computing the Projection Operator

- Optimality conditions imply:  $g(z) = c$  and  $z - x + \lambda \nabla g(z) = 0$ . If both these can be solved in closed form, we are done!
- How to compute the Projection operators for constraints  $\mathcal{C}_g = \{x \mid g(x) \leq c\}$  when

- $g(x) = a^T x$
- $g(x) = \|x\|_2^2$
- $g(x) = \|x\|_1$
- $g(x) = \|x - x_0\|^2$
- $g(x) = x^T Ax + bx + c$
- $g(x) = \|x\|_\infty$

$a^T x = c \quad \& \quad z - x + \lambda a = 0 \quad \text{if } x \text{ is optimal pt on bndry}$

else  $x = z$

Try out these very simple cases



# Easy to Project Sets $\mathcal{C}$ (with closed form solutions)

- Solution set of a linear system  $\mathcal{C} = \{\mathbf{x} \in \Re^n : A^T \mathbf{x} = \mathbf{b}\}$
- Affine images  $\mathcal{C} = \{A\mathbf{x} + \mathbf{b} : \mathbf{x} \in \Re^n\}$
- Nonnegative orthant  $\mathcal{C} = \{\mathbf{x} \in \Re^n : \mathbf{x} \succeq 0\}$ . It may be hard to project on arbitrary polyhedron.
- Norm balls  $\mathcal{C} = \{\mathbf{x} \in \Re^n : \|\mathbf{x}\|_p \leq 1\}$ , for  $p = 1, 2, \infty$   
Homework: For each of the example constrained g's stated here, apply the necessary condition for constrained optimality stated (and motivated) on the previous slide and derive the projection operation.



# Easy to Project Sets $\mathcal{C}$ (with closed form solutions)

- Solution set of a linear system  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : A^T \mathbf{x} = \mathbf{b}\}$
- Affine images  $\mathcal{C} = \{A\mathbf{x} + \mathbf{b} : \mathbf{x} \in \mathbb{R}^n\}$
- Nonnegative orthant  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \succeq 0\}$ . It may be hard to project on arbitrary polyhedron.
- Norm balls  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_p \leq 1\}$ , for  $p = 1, 2, \infty$

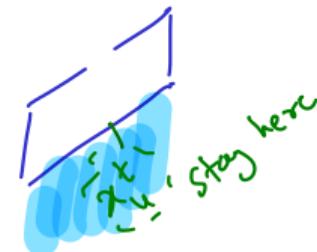
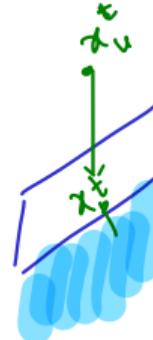


Table of Orthogonal Projections: See [https://archive.siam.org/books/mo25/mo25\\_ch6.pdf](https://archive.siam.org/books/mo25/mo25_ch6.pdf)

$$P_C(\mathbf{z}) = \text{prox}_{I_C}(\mathbf{z}) = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|^2 + I_C(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x} \in C} \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|^2$$

Set $C =$	For $\gamma = 1$ , $P_C(\mathbf{z}) =$	Assumptions	
$\Re^n_+$	$[\mathbf{z}]_+$		
Box $[\mathbf{l}, \mathbf{u}]$	$P_C(\mathbf{z})_i = \min\{\max\{z_i, l_i\}, u_i\}$	$l_i \leq u_i$	
Ball $(\mathbf{c}, r)$	$\mathbf{c} + \frac{r}{\max\{\ \mathbf{z} - \mathbf{c}\ _2, r\}}(\mathbf{z} - \mathbf{c})$	$\ \cdot\ _2$ ball, centre $\mathbf{c} \in \Re^n$ & radius $r > 0$	
$\{\mathbf{x}   A\mathbf{x} = \mathbf{b}\}$	$\mathbf{z} - A^T(AA^T)^{-1}(A\mathbf{z} - \mathbf{b})$	$A \in \Re^{m \times n}$ , $\mathbf{b} \in \Re^m$ , $A$ is full row rank	
$\{\mathbf{x}   \mathbf{a}^T \mathbf{x} \leq b\}$	$\mathbf{z} - \frac{[\mathbf{a}^T \mathbf{z} - b]_+}{\ \mathbf{a}\ ^2} \mathbf{a}$	$0 \neq \mathbf{a} \in \Re^n$ $b \in \Re$	
$\Delta_n$ ( $n$ -simplex)	$[\mathbf{z} - \mu^* \mathbf{e}]_+$ where $\mu^* \in \Re$ satisfies $\mathbf{e}^T [\mathbf{z} - \mu^* \mathbf{e}]_+ = 1$		
$H_{\mathbf{a}, b} \cap \text{Box}[\mathbf{l}, \mathbf{u}]$	$P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \mu^* \mathbf{a})$ where $\mu^* \in \Re$ satisfies $\mathbf{a}^T P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \mu^* \mathbf{a}) = b$	$0 \neq \mathbf{a} \in \Re^n$ $b \in \Re$	
$H_{-\mathbf{a}, b} \cap \text{Box}[\mathbf{l}, \mathbf{u}]$	$P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z})$ $P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \lambda^* \mathbf{a})$ where $\lambda^* \in \Re$ satisfies	$\mathbf{a}^T P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z}) \leq b$ $\mathbf{a}^T P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z}) > b$ $\mathbf{a}^T P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \lambda^* \mathbf{a}) = b$ & $\lambda^* > 0$	$0 \neq \mathbf{a} \in \Re^n$ $b \in \Re$
$B_{\ \cdot\ _1}[0, \alpha]$	$\mathbf{z}$ $[\mathbf{z} - \lambda^* \mathbf{e}]_+ \odot \text{sign}(\mathbf{z})$ where $\lambda^* > 0$ , & $[\mathbf{z} - \lambda^* \mathbf{e}]_+ \odot \text{sign}(\mathbf{z}) = \alpha$	$\ \mathbf{z}\ _1 \leq \alpha$ $\ \mathbf{z}\ _1 > \alpha$ $\alpha > 0$	



# Table of Orthogonal Projections:

See [https://archive.siam.org/books/mo25/mo25\\_ch6.pdf](https://archive.siam.org/books/mo25/mo25_ch6.pdf)

$$P_C(z) = \text{prox}_{I_C}(z) = \operatorname{argmin}_x \frac{1}{2\gamma} \|x - z\|^2 + I_C(x) = \operatorname{argmin}_{x \in C} \frac{1}{2\gamma} \|x - z\|^2$$



Calculus of  
projection  
with multiple  
constraints

Set $C =$	For $\gamma = 1$ , $P_C(z) =$	Assumptions
$\mathbb{R}^n_+$	$[z]_+$	
Box[ $\mathbf{l}, \mathbf{u}$ ]	$P_C(z)_i = \min\{\max\{z_i, l_i\}, u_i\}$	$l_i \leq u_i$
Ball[ $\mathbf{c}, r$ ]	$\mathbf{c} + \frac{\max\{\ \mathbf{z} - \mathbf{c}\ _2, r\}}{r}(\mathbf{z} - \mathbf{c})$	$\ \cdot\ _2$ ball, centre $\mathbf{c} \in \mathbb{R}^n$ & radius $r > 0$
$\{\mathbf{x}   A\mathbf{x} = \mathbf{b}\}$	$\mathbf{z} - A^T(AA^T)^{-1}(Az - \mathbf{b})$	$A \in \mathbb{R}^{m \times n}$ , $\mathbf{b} \in \mathbb{R}^m$ , $A$ is full row rank
$\{\mathbf{x}   \mathbf{a}^T \mathbf{x} \leq b\}$	$\mathbf{z} - \frac{[\mathbf{a}^T \mathbf{z} - b]_+}{\ \mathbf{a}\ ^2} \mathbf{a}$	$0 \neq \mathbf{a} \in \mathbb{R}^n$ $b \in \mathbb{R}$
$\Delta_n$ ( $n$ -simplex)	$[\mathbf{z} - \mu^* \mathbf{e}]_+$ where $\mu^* \in \mathbb{R}$ satisfies $\mathbf{e}^T [\mathbf{z} - \mu^* \mathbf{e}]_+ = 1$	
$H_{\mathbf{a}, b} \cap \text{Box}[\mathbf{l}, \mathbf{u}]$	$P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \mu^* \mathbf{a})$ where $\mu^* \in \mathbb{R}$ satisfies $\mathbf{a}^T P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \mu^* \mathbf{a}) = b$	$0 \neq \mathbf{a} \in \mathbb{R}^n$ $b \in \mathbb{R}$
$H_{\mathbf{a}, b}^- \cap \text{Box}[\mathbf{l}, \mathbf{u}]$	$P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z})$ $P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \lambda^* \mathbf{a})$ where $\lambda^* \in \mathbb{R}$ satisfies $\mathbf{a}^T P_{\text{Box}[\mathbf{l}, \mathbf{u}]}(\mathbf{z} - \lambda^* \mathbf{a}) = b$ & $\lambda^* > 0$	$0 \neq \mathbf{a} \in \mathbb{R}^n$ $b \in \mathbb{R}$
$B_{\ \cdot\ _1}[0, \alpha]$	$\mathbf{z}$ $[\mathbf{z} - \lambda^* \mathbf{e}]_+ \odot \text{sign}(\mathbf{z})$ where $\lambda^* > 0$ , & $[\mathbf{z} - \lambda^* \mathbf{e}]_+ \odot \text{sign}(\mathbf{z}) = \alpha$	$\ \mathbf{z}\ _1 \leq \alpha$ $\ \mathbf{z}\ _1 > \alpha$ $\alpha > 0$



Geometric intuitions are clear. Question is how these analytical expressions are derived?  
Lagrange function  $(f + \gamma g)$  and KKT conditions



# Convergence Results for Projected Gradient Descent (PGD)

- Lipschitz continuous function  $f$  (C) using PGD:  $R^2 B^2 / \epsilon^2$  iterations
- Lipschitz continuous functions + Strongly Convex  $f$  (CS) using PGD:  $2B^2/\epsilon - 1$  iterations
- Smooth Function  $f$  using PGD:  $\frac{R^2 L}{\epsilon}$  iterations.
- Smooth Functions using Nesterov's PGD:  $\sqrt{\frac{2R^2 L}{\epsilon}}$  iterations
- Smooth + Strongly Convex  $f$  (SS) using PGD: With  $\gamma = 1/L$ , achieve an  $\epsilon$ -approximate solution in  $\frac{L}{\mu} \log(\frac{R^2 L}{2\epsilon})$  iterations.
- Smooth + Strongly Convex  $f$  (SS) using Nesterov's PGD: With  $\gamma = 1/L$ , achieve an  $\epsilon$ -approximate solution in  $\sqrt{\frac{L}{\mu}} \log(\frac{R^2 L}{2\epsilon})$  iterations.
- All results for standard Gradient descent carry over to the projected case  
**as long as the projection operator is easy to compute!**
- Proofs in monograph by Sebastian Bubeck at  
<https://moodle.iitb.ac.in/mod/resource/view.php?id=36947>



# How to solve the Projected Gradient Descent in closed form?

Lagrange Function, KKT Conditions & Duality For Solving Constrained Optimization



# Karush Kuhn Tucker (KKT) Conditions for Optimality, Conditions for Necessity and Sufficiency etc. for..

Calculus of projection with multiple constraints

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathcal{D}} & f(\mathbf{x}) \\ \text{subject to} & g_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, m \\ & h_j(\mathbf{x}) = 0 \quad j = 1, \dots, p \\ \text{variable } \mathbf{x} = & (x_1, \dots, x_n) \end{array} \tag{1}$$



# Karush Kuhn Tucker (KKT) Conditions for Optimality, Conditions for Necessity and Sufficiency etc. for...

$\min_{x \in D}$   
subject to

$f(x)$

$$\begin{aligned} g_i(x) &\leq 0 & i = 1, \dots, m \\ h_j(x) &= 0 & j = 1, \dots, p \end{aligned}$$

For one particular iteration of projected gradient descent,  $f(x) = \frac{1}{2} \|z - x\|^2$

variable  $x = (x_1, \dots, x_n)$

n dimensional space with  
m inequality constraints and  
p equality constraints

option 2

Instead of solving the constrained optimization problem iteratively by gradient descent on  $f$  and then projection with constraints, if the  $f$  is itself simple, one could solve this constrained optimization problem directly



# Briefly: The Dual Theory for Constrained Optimization

Consider the general constrained minimization problem

$$\begin{aligned} & \min_{\mathbf{x} \in \mathcal{D}} && f(\mathbf{x}) \\ & \text{subject to} && g_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m \\ & \text{subject to} && h_j(\mathbf{x}) = 0, j = 1, 2, \dots, p \end{aligned} \tag{2}$$

- Consider forming the lagrange function by associating prices (called lagrange multipliers)  $\lambda_i$  and  $\mu_j$ , with constraints involving  $g_i$  and  $h_j$  respectively.

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

- At each **feasible  $\mathbf{x}$** , for fixed  $\lambda_i \geq 0 \forall i \in \{1..m\}$ ,

$$f(\mathbf{x}) \geq L(\mathbf{x}, \lambda, \mu) \quad \text{if } g_i(\mathbf{x}) \leq 0 \text{ & } h_j(\mathbf{x}) = 0$$

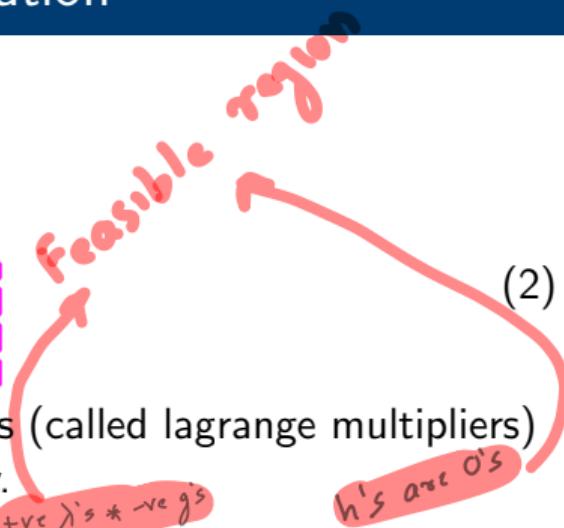
# Briefly: The Dual Theory for Constrained Optimization

Consider the general constrained minimization problem

$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$$

$$\text{subject to } g_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m$$

$$\text{subject to } h_j(\mathbf{x}) = 0, j = 1, 2, \dots, p$$



- Consider forming the lagrange function by associating prices (called lagrange multipliers)  $\lambda_i$  and  $\mu_j$ , with constraints involving  $g_i$  and  $h_j$  respectively.

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

- At each **feasible  $\mathbf{x}$** , for fixed  $\lambda_i \geq 0 \forall i \in \{1..m\}$ ,

Rationale behind lagrange function:  
At the boundary, the gradient of  $f$  must lie in the space spanned by the gradients of  $g$ 's and  $h$ 's with some sign constraints

$$f(\mathbf{x}) \geq L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \quad \text{if } g_i(\mathbf{x}) \leq 0 \text{ & } h_j(\mathbf{x}) = 0$$



# Briefly: The Dual Theory for Constrained Optimization

- For  $\lambda_i \geq 0 \forall i \in \{1..m\}$  and any  $\mu_j$ , minimizing the right hand side of (3) over all **feasible  $x$**

$$f(x) \geq \min_{\substack{x \text{ s.t } g_i(x) \leq 0, h_j(x) = 0}} L(x, \lambda, \mu) \triangleq L^*(\lambda, \mu) \quad (4)$$

- $L^*(\lambda, \mu)$  is a pointwise (w.r.t  $x \in g_i(x) \leq 0, h_j(x) = 0$ ) minimum of linear functions ( $L(x, \lambda, \mu)$ ) and is therefore always a concave function.
- Since  $f(x) \geq L^*(\lambda, \mu)$  for all **primal feasible  $x$**  and **dual feasible i.e.,  $\lambda_i \geq 0$  and  $\mu_j$** , we can maximize the lower bound  $L^*(\lambda, \mu)$  to give the following **Dual Problem**

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} L^*(\lambda, \mu) \\ & \text{subject to } \lambda \geq 0 \end{aligned}$$



# Briefly: The Dual Theory for Constrained Optimization

- For  $\lambda_i \geq 0 \forall i \in \{1..m\}$  and any  $\mu_j$ , minimizing the right hand side of (3) over all **feasible  $x$**

$$f(x) \geq \min_{x \text{ s.t. } g_i(x) \leq 0, h_j(x) = 0} L(x, \lambda, \mu) \triangleq L^*(\lambda, \mu) \quad (4)$$

RHS has  
no  $x$

RHS has  
 $x$

New RHS  
has no  $x$

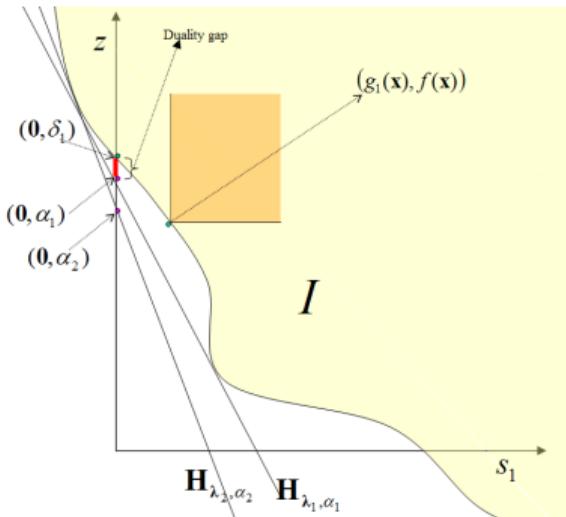
Inequality will hold even if I minimize the RHS with respect to all feasible  $x$

- $L^*(\lambda, \mu)$  is a pointwise (w.r.t  $x \in g_i(x) \leq 0, h_j(x) = 0$ ) minimum of linear functions ( $L(x, \lambda, \mu)$ ) and is therefore always a concave function.
- Since  $f(x) \geq L^*(\lambda, \mu)$  for all **primal feasible  $x$**  and **dual feasible** i.e.,  $\lambda_i \geq 0$  and  $\mu_j$ , we can maximize the lower bound  $L^*(\lambda, \mu)$  to give the following **Dual Problem**

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} L^*(\lambda, \mu) \\ & \text{subject to } \lambda \geq 0 \end{aligned}$$



# Briefly: The Dual Theory for Constrained Optimization (contd.)



The **Dual Problem** restated in two ways:

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} L^*(\lambda, \mu) \\ & \text{subject to } \lambda \geq 0 \end{aligned} \tag{5}$$

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} \min_{x \text{ s.t. } g_i(x) \leq 0, h_j(x) = 0} L(x, \lambda, \mu) \\ & \text{subject to } \lambda \geq 0 \end{aligned}$$

**Figure 3:** See optional course material referred to at the end for understanding this figure completely

## Theorem

- (i) The dual function  $L^*(\lambda, \mu)$  is always concave. (ii) If  $p^*$  is solution of (2) and  $d^*$  of (5) then  $p^* > d^*$  and the gap  $p^* - d^*$  is called the **duality gap**.

# Briefly: The Dual Theory for Constrained Optimization (contd.)

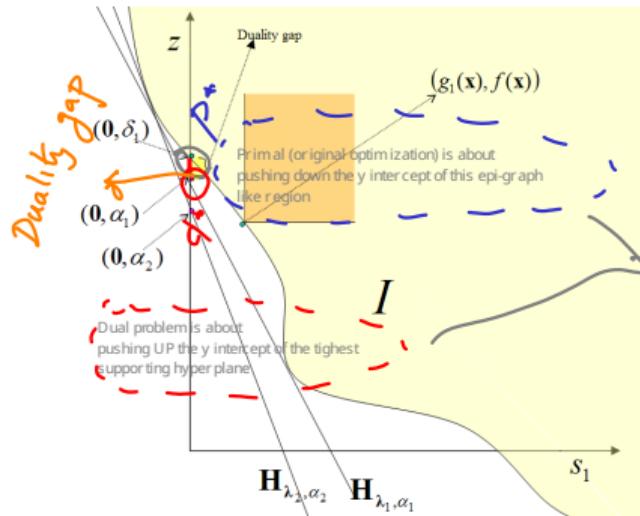


Figure 3: See optional course material referred to at the end for understanding this figure completely

The **Dual Problem** restated in two ways:

x axis is space of values of  $g$   
y axis is space of value of  $f$

$$\max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} L^*(\lambda, \mu) \quad (5)$$

The two y-intercepts can be expected to meet when the shape of  $I$  is convex and closed

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} \\ & \text{subject to} \quad \min_{x \text{ s.t. } g_i(x) \leq 0, h_j(x) = 0} L(x, \lambda, \mu) \\ & \lambda \geq 0 \end{aligned}$$

## Theorem

- (i) The dual function  $L^*(\lambda, \mu)$  is always concave. (ii) If  $p^*$  is solution of (2) and  $d^*$  of (5) then  $p^* > d^*$  and the gap  $p^* - d^*$  is called the **duality gap**

# Zero Duality Gap $\Rightarrow$ KKT Conditions

## Theorem

Assuming requisite differentiability, we can state the following Karush-Kuhn-Tucker (KKT) necessary conditions conditions in (6) for zero duality gap:

- (1)  $\nabla f(\hat{\mathbf{x}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{x}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{x}}) = \mathbf{0}$
- (2)  $g_i(\hat{\mathbf{x}}) \leq 0 \quad i = 1, 2, \dots, m$
- (3)  $\hat{\lambda}_i \geq 0 \quad i = 1, 2, \dots, m$
- (4)  $\hat{\lambda}_i g_i(\hat{\mathbf{x}}) = 0 \quad i = 1, 2, \dots, m$
- (5)  $h_j(\hat{\mathbf{x}}) = 0 \quad j = 1, 2, \dots, p$



# Zero Duality Gap $\Rightarrow$ KKT Conditions

CONVEXITY IS NOT PART OF NECESSARY CONDITIONS...

## Theorem

Assuming requisite differentiability, we can state the following Karush-Kuhn-Tucker (KKT) necessary conditions conditions in (6) for zero duality gap:

If  $\hat{\mathbf{x}} = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} f(\mathbf{x})$  s.t.  $\sum_{i=1}^m g_i(\mathbf{x}) \leq 0$  and  $\sum_{j=1}^p h_j(\mathbf{x}) = 0$  } Then there must exist  $\hat{\lambda}$  &  $\hat{\mu}$  & following conditions with  $(\hat{\mathbf{x}}, \hat{\lambda}, \hat{\mu})$

{

(1)  $\nabla f(\hat{\mathbf{x}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{x}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{x}}) = \mathbf{0}$

(2)  $g_i(\hat{\mathbf{x}}) \leq 0 \quad i = 1, 2, \dots, m$

(3)  $\hat{\lambda}_i \geq 0 \quad i = 1, 2, \dots, m$

(4)  $\hat{\lambda}_i g_i(\hat{\mathbf{x}}) = 0 \quad i = 1, 2, \dots, m$

(5)  $h_j(\hat{\mathbf{x}}) = 0 \quad j = 1, 2, \dots, p$

COMPLEMENTARY SLACKNESS

Primal feasibility conditions (6)

Either you are on boundary  $g_i(\hat{\mathbf{x}}) = 0$  & therefore  $\hat{\lambda}_i$  plays some role  
OR you are inside  $g_i(\hat{\mathbf{x}}) < 0$  & therefore  $\hat{\lambda}_i$  plays no role

# Zero Duality Gap $\Leftarrow$ KKT Conditions+Convexity

## Theorem

If the function  $f$  is convex,  $g_i$  are convex and  $h_j$  are affine, then KKT conditions in (6) are necessary and sufficient conditions for zero duality gap.

$$\begin{aligned} (1) \quad \nabla f(\hat{\mathbf{x}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{x}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{x}}) &= \mathbf{0} \\ (2) \quad g_i(\hat{\mathbf{x}}) &\leq 0 \quad i = 1, 2, \dots, m \\ (3) \quad \hat{\lambda}_i &\geq 0 \quad i = 1, 2, \dots, m \\ (4) \quad \hat{\lambda}_i g_i(\hat{\mathbf{x}}) &= 0 \quad i = 1, 2, \dots, m \\ (5) \quad h_j(\hat{\mathbf{x}}) &= 0 \quad j = 1, 2, \dots, p \end{aligned} \tag{6}$$



# Zero Duality Gap $\Leftarrow$ KKT Conditions+Convexity

*Sufficiency  
owing to convexity*

## Theorem

If the function  $f$  is convex,  $g_i$  are convex and  $h_j$  are affine, then KKT conditions in (6) are necessary and sufficient conditions for zero duality gap.

ALL THE ALGEBRAIC EXPRESSIONS OF PROJECTION OPERATIONS THAT WE SHOWED IN THE EARLIER TABLE FOR PROJECTION CAN BE DERIVED USING THESE KKT CONDITIONS! H/W TRY THESE ON SOME!

- |     |  |        |              |                      |
|-----|--|--------|--------------|----------------------|
| (1) | $\nabla f(\hat{\mathbf{x}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{x}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{x}})$ | =      | $\mathbf{0}$ |                      |
| (2) | $g_i(\hat{\mathbf{x}})$  | $\leq$ | 0            | $i = 1, 2, \dots, m$ |
| (3) | $\hat{\lambda}_i$  | $\geq$ | 0            | $i = 1, 2, \dots, m$ |
| (4) | $\hat{\lambda}_i g_i(\hat{\mathbf{x}})$  | =      | 0            | $i = 1, 2, \dots, m$ |
| (5) | $h_j(\hat{\mathbf{x}})$  | =      | 0            | $j = 1, 2, \dots, p$ |



# Convexity, KKT Conditions, Slater's Qualification, Zero Duality Gap

Problem type	Objective Function	Constraints	$L^*(\lambda)$	Dual constraints	Strong duality
Linear Program	$\mathbf{c}^T \mathbf{x}$	$A\mathbf{x} \leq \mathbf{b}$	$-\mathbf{b}^T \lambda$	$A^T \lambda + \mathbf{c} = \mathbf{0}$ $\lambda \geq \mathbf{0}$	Feasible primal and dual
Quadratic Program	$\frac{1}{2}\mathbf{x}^T Q\mathbf{x} + \mathbf{c}^T \mathbf{x}$ for $Q \in \mathcal{S}_{++}^n$	$A\mathbf{x} \leq \mathbf{b}$	$-\frac{1}{2}(\mathbf{c} - A^T \lambda)^T Q^{-1}(\mathbf{c} - A^T \lambda) + \mathbf{b}^T \lambda$	$\lambda \geq \mathbf{0}$	Always
Entropy maximization	$x_i \sum_{i=1}^n \ln x_i$	$A\mathbf{x} \leq \mathbf{b}$ $\mathbf{x}^T \mathbf{1} = 1$	$-\mathbf{b}^T \lambda - \mu - e^{-\mu-1} \sum_{i=1}^n e^{-\mathbf{a}_i^T \lambda}$ $\mathbf{a}_i$ is the $i^{th}$ column of $A$	$\lambda \geq \mathbf{0}$	Primal constraints are satisfied.

Table 1: Some optimization problems, their duals and conditions for strong duality.

HOMEWORK: CAN THE FORM OF THE DUAL BECOME SIGNIFICANTLY SIMPLER THAN THE PRIMAL IS THAT WHAT HAPPENED IN SVMs?

Topics covered in Extra optional material from previous offering

- Geometric Interpretation of Necessary conditions of optimality<sup>1</sup>
- Geometric Interpretation of Zero duality gap and connection with convexity<sup>2</sup>
- Weak and Strong Duality Theorem and their proofs, Slater's Constraint Qualification<sup>3</sup>
- Conditional Gradient Descent (Frank Wolfe Method), Convex Conjugate, Dual Ascent, ADMM<sup>4</sup>

<sup>1</sup> [https://www.youtube.com/watch?v=03vfurE\\_rs0&list=PLyo3HAXSZD3xRZ3mL0W6ERtQ3huVVmyo0&index=27](https://www.youtube.com/watch?v=03vfurE_rs0&list=PLyo3HAXSZD3xRZ3mL0W6ERtQ3huVVmyo0&index=27)

<sup>2</sup> <https://www.youtube.com/watch?v=yG6HxDqHj1Y&list=PLyo3HAXSZD3xRZ3mL0W6ERtQ3huVVmyo0&index=29>

<sup>3</sup> <https://www.youtube.com/watch?v=r4Pch0ek-YI&list=PLyo3HAXSZD3xRZ3mL0W6ERtQ3huVVmyo0&index=31>

<sup>4</sup> <https://www.youtube.com/watch?v=Jlmsr3uqz3Q&list=PLyo3HAXSZD3xRZ3mL0W6ERtQ3huVVmyo0&index=34>



# Support Vector Regression and its Dual

Self Reading Material for practical illustration of Duality, KKT Conditions, Necessity, Sufficiency, etc.



# KKT and Dual for SVR

- $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$   
s.t.  $\forall i,$   
 $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i,$   
 $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*,$   
 $\xi_i, \xi_i^* \geq 0$
- Let's consider the lagrange multipliers  $\alpha_i, \alpha_i^*, \mu_i$  and  $\mu_i^*$  corresponding to the above-mentioned constraints.
- The Lagrange Function is



# KKT and Dual for SVR

- $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$

s.t.  $\forall i,$

$$y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i,$$

$$b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*,$$

$$\xi_i, \xi_i^* \geq 0$$

- Let's consider the lagrange multipliers  $\alpha_i, \alpha_i^*, \mu_i$  and  $\mu_i^*$  corresponding to the above-mentioned constraints.
- The Lagrange Function is

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

# KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$$

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,



# KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$$

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,

$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0 \text{ i.e., } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

- Differentiating the Lagrangian w.r.t.  $\xi_i$ ,



# KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$$

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,

$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0 \text{ i.e., } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

- Differentiating the Lagrangian w.r.t.  $\xi_i$ ,

$$C - \alpha_i - \mu_i = 0 \text{ i.e., } \alpha_i + \mu_i = C$$

- Differentiating the Lagrangian w.r.t  $\xi_i^*$ ,



# KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$$

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,

$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0 \text{ i.e., } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

- Differentiating the Lagrangian w.r.t.  $\xi_i$ ,

$$C - \alpha_i - \mu_i = 0 \text{ i.e., } \alpha_i + \mu_i = C$$

- Differentiating the Lagrangian w.r.t  $\xi_i^*$ ,

$$\alpha_i^* + \mu_i^* = C$$

- Differentiating the Lagrangian w.r.t  $b$ ,



# KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$$

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,

$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0 \text{ i.e., } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

- Differentiating the Lagrangian w.r.t.  $\xi_i$ ,

$$C - \alpha_i - \mu_i = 0 \text{ i.e., } \alpha_i + \mu_i = C$$

- Differentiating the Lagrangian w.r.t  $\xi_i^*$ ,

$$\alpha_i^* + \mu_i^* = C$$

- Differentiating the Lagrangian w.r.t  $b$ ,

$$\sum_i (\alpha_i^* - \alpha_i) = 0$$

- Complimentary slackness:



# KKT conditions for SVR

$$L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) =$$

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \sum_{i=1}^m \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \mu_i^* \xi_i^*$$

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,

$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0 \text{ i.e., } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

- Differentiating the Lagrangian w.r.t.  $\xi_i$ ,

$$C - \alpha_i - \mu_i = 0 \text{ i.e., } \alpha_i + \mu_i = C$$

- Differentiating the Lagrangian w.r.t  $\xi_i^*$ ,

$$\alpha_i^* + \mu_i^* = C$$

- Differentiating the Lagrangian w.r.t  $b$ ,

$$\sum_i (\alpha_i^* - \alpha_i) = 0$$

- Complimentary slackness:

$$\alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0 \text{ AND } \mu_i \xi_i = 0 \text{ AND}$$
$$\alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) = 0 \text{ AND } \mu_i^* \xi_i^* = 0$$



# Conclusions from the KKT conditions:

$\alpha_i \in (0, C) \Rightarrow ?$

$\alpha_i^* \in (0, C) \Rightarrow ?$

KKT conditions  
clearly outline  
support vectors.



# KKT conditions

- Differentiating the Lagrangian w.r.t.  $\mathbf{w}$ ,

$$\mathbf{w} - \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0$$

$$\text{i.e. } \mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

- Differentiating the Lagrangian w.r.t.  $\xi_i$ ,

$$C - \alpha_i - \mu_i = 0$$

$$\text{i.e. } \alpha_i + \mu_i = C$$

- Differentiating the Lagrangian w.r.t  $\xi_i^*$ ,

$$\alpha_i^* + \mu_i^* = C$$

- Differentiating the Lagrangian w.r.t  $b$ ,

$$\sum_i^m (\alpha_i^* - \alpha_i) = 0$$

- Complimentary slackness:

$$\alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0$$

$$\mu_i \xi_i = 0$$

$$\alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) = 0$$

$$\mu_i^* \xi_i^* = 0$$



# Conclusions from the KKT conditions:

$$\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0$$

and

$$\alpha_i^*(b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i - \epsilon - \xi_i^*) = 0$$

$\Rightarrow ?$



# Conclusions from the KKT conditions:

$$\alpha_i \in (0, C) \Rightarrow ?$$

$$(C - \alpha_i)\xi_i = 0 \Rightarrow ?$$

$$\alpha_i^* \in (0, C) \Rightarrow ?$$

$$(C - \alpha_i^*)\xi_i^* = 0 \Rightarrow ?$$



For Support Vector Regression, since the original objective and the constraints are convex, any  $(\mathbf{w}, b, \alpha, \alpha^*, \mu, \mu^*, \xi, \xi^*)$  that satisfy the necessary KKT conditions gives optimality (conditions are also sufficient)



# Some observations

- $\alpha_i, \alpha_i^* \geq 0, \mu_i, \mu_i^* \geq 0, \alpha_i + \mu_i = C$  and  $\alpha_i^* + \mu_i^* = C$   
Thus,  $\alpha_i, \mu_i, \alpha_i^*, \mu_i^* \in [0, C], \forall i$
- If  $0 < \alpha_i < C$ , then  $0 < \mu_i < C$   
(as  $\alpha_i + \mu_i = C$ )
- $\mu_i \xi_i = 0$  and  $\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0$  are complementary slackness conditions  
So  $0 < \alpha_i < C \Rightarrow \xi_i = 0$  and  $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b = \epsilon + \xi_i = \epsilon$ 
  - ▶ All such points lie on the boundary of the  $\epsilon$  band
  - ▶ Using any point  $\mathbf{x}_j$  (that is with  $\alpha_j \in (0, C)$ ) on margin, we can recover  $b$  as:  
$$b = y_j - \mathbf{w}^\top \phi(\mathbf{x}_j) - \epsilon$$



# Support Vector Regression

## Dual Objective



# Weak Duality

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$

- By weak duality theorem, we have:

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t.  $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$ , and

$\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$ , and

$\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$

- The above is true for any  $\alpha_i, \alpha_i^* \geq 0$  and  $\mu_i, \mu_i^* \geq 0$

- Thus,



# Weak Duality

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$

- By weak duality theorem, we have:

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t.  $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$ , and

$\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$ , and

$\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$

- The above is true for any  $\alpha_i, \alpha_i^* \geq 0$  and  $\mu_i, \mu_i^* \geq 0$

- Thus,

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t.  $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$ , and

$\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$ , and

$\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$



## Dual objective

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$
- Assume: In case of SVR, we have a strictly convex objective and linear constraints  $\Rightarrow$  KKT conditions are necessary and sufficient and strong duality holds:

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) = \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t.  $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$ , and

$\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$ , and

$\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$

- This value is precisely obtained at the  $(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$  that satisfies the necessary (and sufficient) KKT optimality conditions
- Given strong duality, we can equivalently solve

$$\max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$



- $$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) +$$

$$\sum_{i=1}^m (\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \alpha_i^*(\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i - \epsilon - \xi_i^*)) - \sum_{i=1}^m (\mu_i \xi_i + \mu_i^* \xi_i^*)$$
- We obtain  $\mathbf{w}$ ,  $b$ ,  $\xi_i$ ,  $\xi_i^*$  in terms of  $\alpha$ ,  $\alpha^*$ ,  $\mu$  and  $\mu^*$  by using the KKT conditions derived earlier as  $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$  and  $\sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0$  and  $\alpha_i + \mu_i = C$  and  $\alpha_i^* + \mu_i^* = C$
- Thus, we get:



- $$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) +$$

$$\sum_{i=1}^m (\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \alpha_i^*(\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i - \epsilon - \xi_i^*)) - \sum_{i=1}^m (\mu_i \xi_i + \mu_i^* \xi_i^*)$$
  - We obtain  $\mathbf{w}$ ,  $b$ ,  $\xi_i$ ,  $\xi_i^*$  in terms of  $\alpha$ ,  $\alpha^*$ ,  $\mu$  and  $\mu^*$  by using the KKT conditions derived earlier as  $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$  and  $\sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0$  and  $\alpha_i + \mu_i = C$  and  $\alpha_i^* + \mu_i^* = C$
  - Thus, we get:
- $$L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$$
- $$= \frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) + \sum_i (\xi_i(C - \alpha_i - \mu_i) + \xi_i^*(C - \alpha_i^* - \mu_i^*)) -$$
- $$b \sum_i (\alpha_i - \alpha_i^*) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) - \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$$



- $$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) +$$

$$\sum_{i=1}^m (\alpha_i(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \alpha_i^*(\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i - \epsilon - \xi_i^*)) - \sum_{i=1}^m (\mu_i \xi_i + \mu_i^* \xi_i^*)$$
  - We obtain  $\mathbf{w}$ ,  $b$ ,  $\xi_i$ ,  $\xi_i^*$  in terms of  $\alpha$ ,  $\alpha^*$ ,  $\mu$  and  $\mu^*$  by using the KKT conditions derived earlier as  $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$  and  $\sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0$  and  $\alpha_i + \mu_i = C$  and  $\alpha_i^* + \mu_i^* = C$
  - Thus, we get:
- $$L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$$
- $$= \frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) + \sum_i (\xi_i(C - \alpha_i - \mu_i) + \xi_i^*(C - \alpha_i^* - \mu_i^*)) -$$
- $$b \sum_i (\alpha_i - \alpha_i^*) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) - \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$$
- $$= -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$



# Kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$

- $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi(\mathbf{x}_i) \Rightarrow$  the final decision function

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi^T(\mathbf{x}_i)\phi(\mathbf{x}) + y_j - \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j) - \epsilon$$

$\mathbf{x}_j$  is any point with  $\alpha_j \in (0, C)$ . Recall similarity with



# Kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$

- $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi(\mathbf{x}_i) \Rightarrow$  the final decision function  
 $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi^T(\mathbf{x}_i)\phi(\mathbf{x}) + y_j - \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j) - \epsilon$   
 $\mathbf{x}_j$  is any point with  $\alpha_j \in (0, C)$ . Recall similarity with kernelized expression for Ridge Regression
- The dual optimization problem to compute the  $\alpha$ 's for SVR is:



# Kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$

- $\mathbf{w} = \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi(\mathbf{x}_i) \Rightarrow$  the final decision function  
 $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi^T(\mathbf{x}_i)\phi(\mathbf{x}) + y_j - \sum_{i=1}^m (\alpha_i - \alpha_i^*)\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j) - \epsilon$   
 $\mathbf{x}_j$  is any point with  $\alpha_j \in (0, C]$ . Recall similarity with kernelized expression for Ridge Regression
- The dual optimization problem to compute the  $\alpha$ 's for SVR is:

$$\max_{\alpha_i, \alpha_i^*} -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$$

$$-\epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$

s.t.

- ▶  $\sum_i (\alpha_i - \alpha_i^*) = 0$
- ▶  $\alpha_i, \alpha_i^* \in [0, C]$

- We notice that the only way these three expressions involve  $\phi$  is through  $\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$ , for some  $i, j$



# An Example Kernel

- Let  $K(\mathbf{x}_1, \mathbf{x}_2) = (1 + \mathbf{x}_1^\top \mathbf{x}_2)^2$
- What  $\phi(\mathbf{x})$  will give  $\phi^\top(\mathbf{x}_1)\phi(\mathbf{x}_2) = K(\mathbf{x}_1, \mathbf{x}_2) = (1 + \mathbf{x}_1^\top \mathbf{x}_2)^2$
- Is such a  $\phi$  guaranteed to exist?
- Is there a unique  $\phi$  for given  $K$ ?



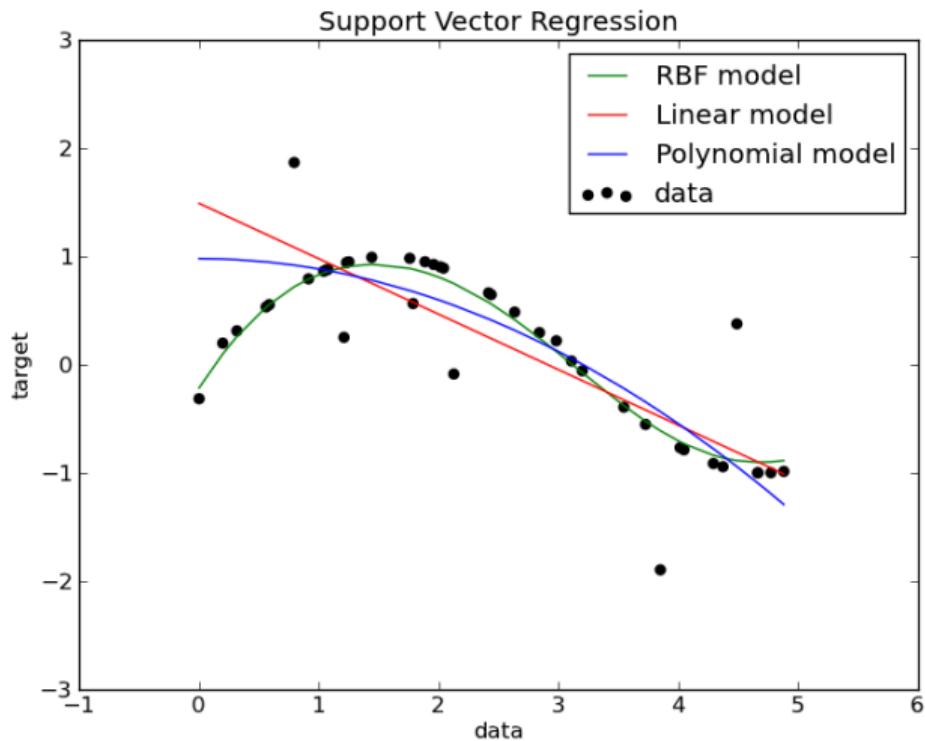
# An Example Kernel

- We can prove that such a  $\phi$  exists
- For example, for a 2-dimensional  $\mathbf{x}_i$ :

$$\phi(\mathbf{x}_i) = \begin{bmatrix} 1 \\ x_{i1}\sqrt{2} \\ x_{i2}\sqrt{2} \\ x_{i1}x_{i2}\sqrt{2} \\ x_{i1}^2 \\ x_{i2}^2 \end{bmatrix}$$

- $\phi(\mathbf{x}_i)$  exists in a 5-dimensional space
- But, to compute  $K(\mathbf{x}_1, \mathbf{x}_2)$ , all we need is  $\mathbf{x}_1^\top \mathbf{x}_2$  without having to enumerate  $\phi(\mathbf{x}_i)$





# Kernels in SVR

- Recall:

$$\max_{\alpha_i, \alpha_i^*} -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K(\mathbf{x}_i, \mathbf{x}_j) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$

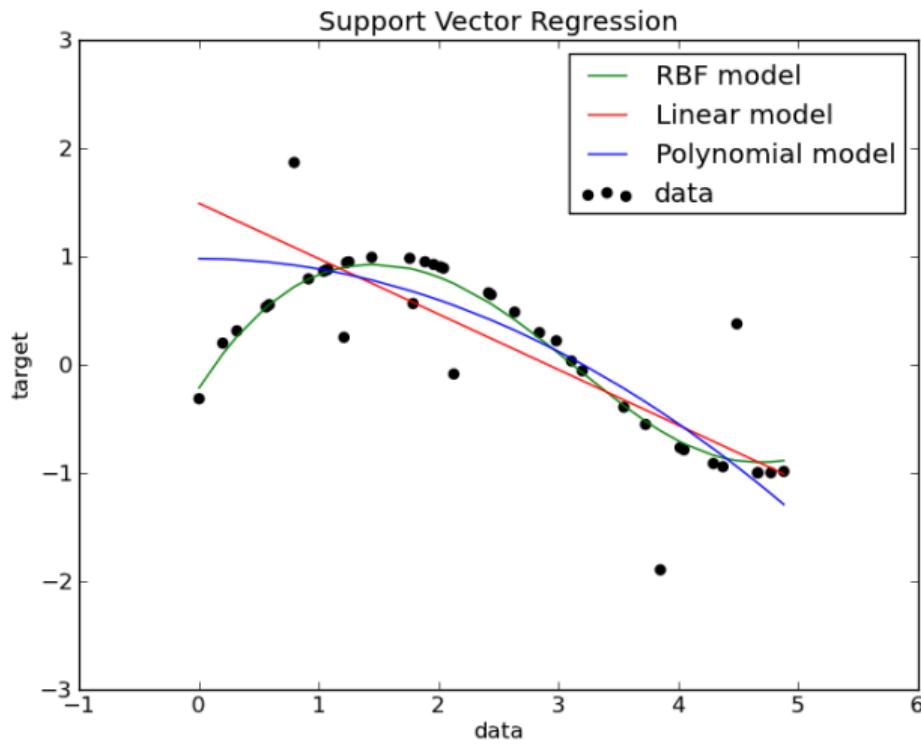
and the decision function:

$$f(\mathbf{x}) = \sum_i (\alpha_i - \alpha_i^*) K(\mathbf{x}_i, \mathbf{x}) + b$$

are all in terms of the kernel  $K(\mathbf{x}_i, \mathbf{x}_j)$  only

- One can now employ any mercer kernel in SVR or Ridge Regression to implicitly perform linear regression in higher dimensional spaces





# Equivalent Forms of Ridge Regression

- Consider the formulation in which we limit the weights of the coefficients by putting a constraint on size of the L2 norm of the weight vector:

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{w}} (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y}) \\ & \| \mathbf{w} \|_2^2 \leq \xi \end{aligned}$$

- The objective function, namely  $f(\mathbf{w}) = (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y})$  is strictly convex. The constraint function,  $g(\mathbf{w}) = \| \mathbf{w} \|_2^2 - \xi$ , is also convex.
- For convex  $g(\mathbf{w})$ , the set  $\{\mathbf{w} | g(\mathbf{w}) \leq 0\}$ , is also convex. (Why?)

# Algorithms for (Submodular function) optimization, with (Matroid) constraints, strategies for improving efficiency

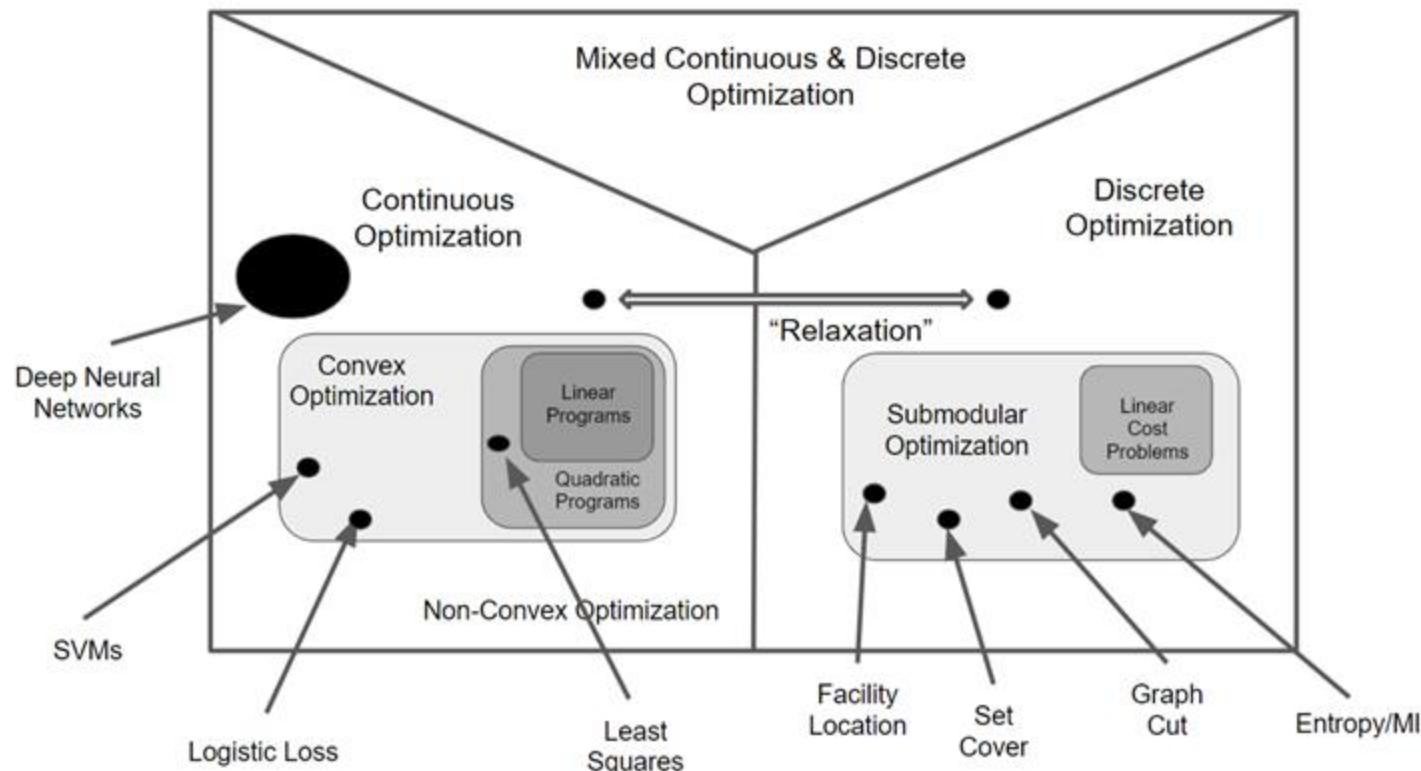
---

**Optimization in Machine Learning**

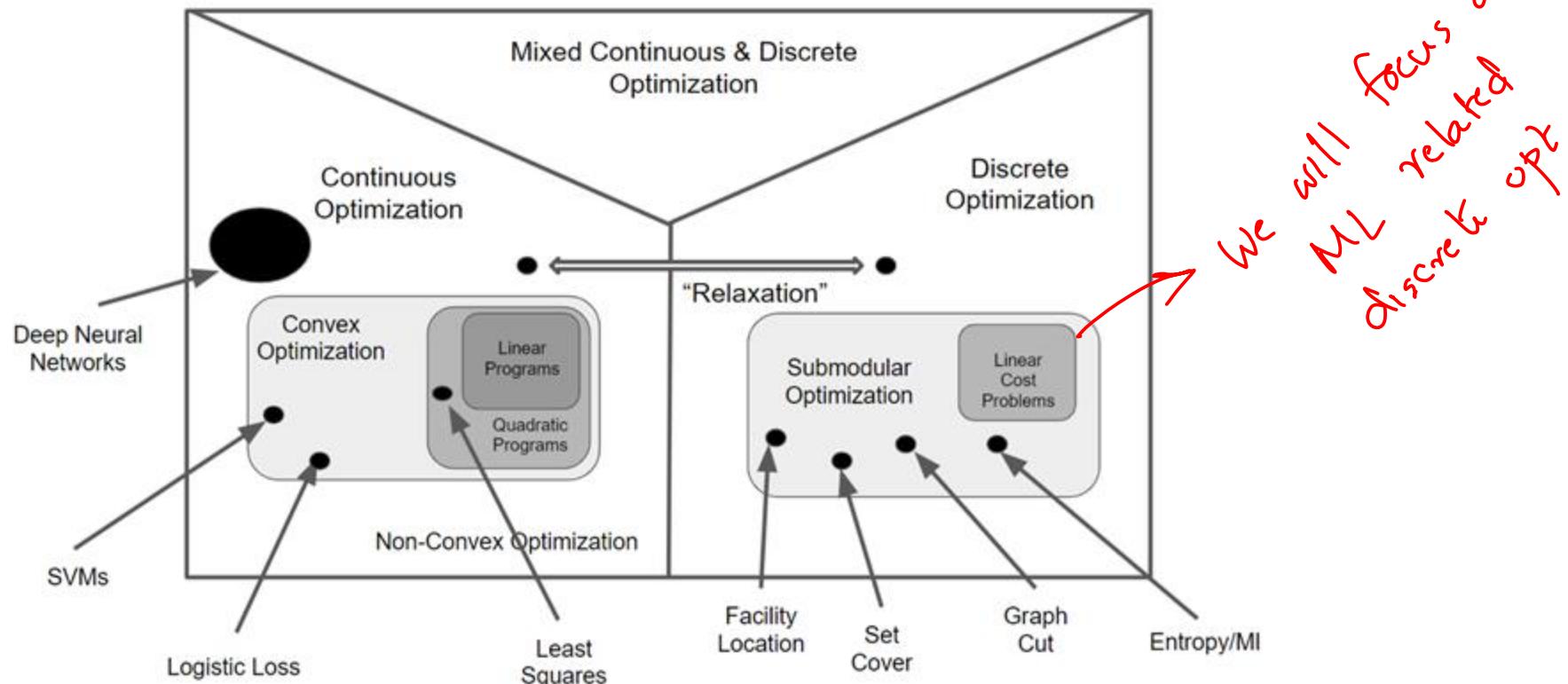
**Lecture 20**

20

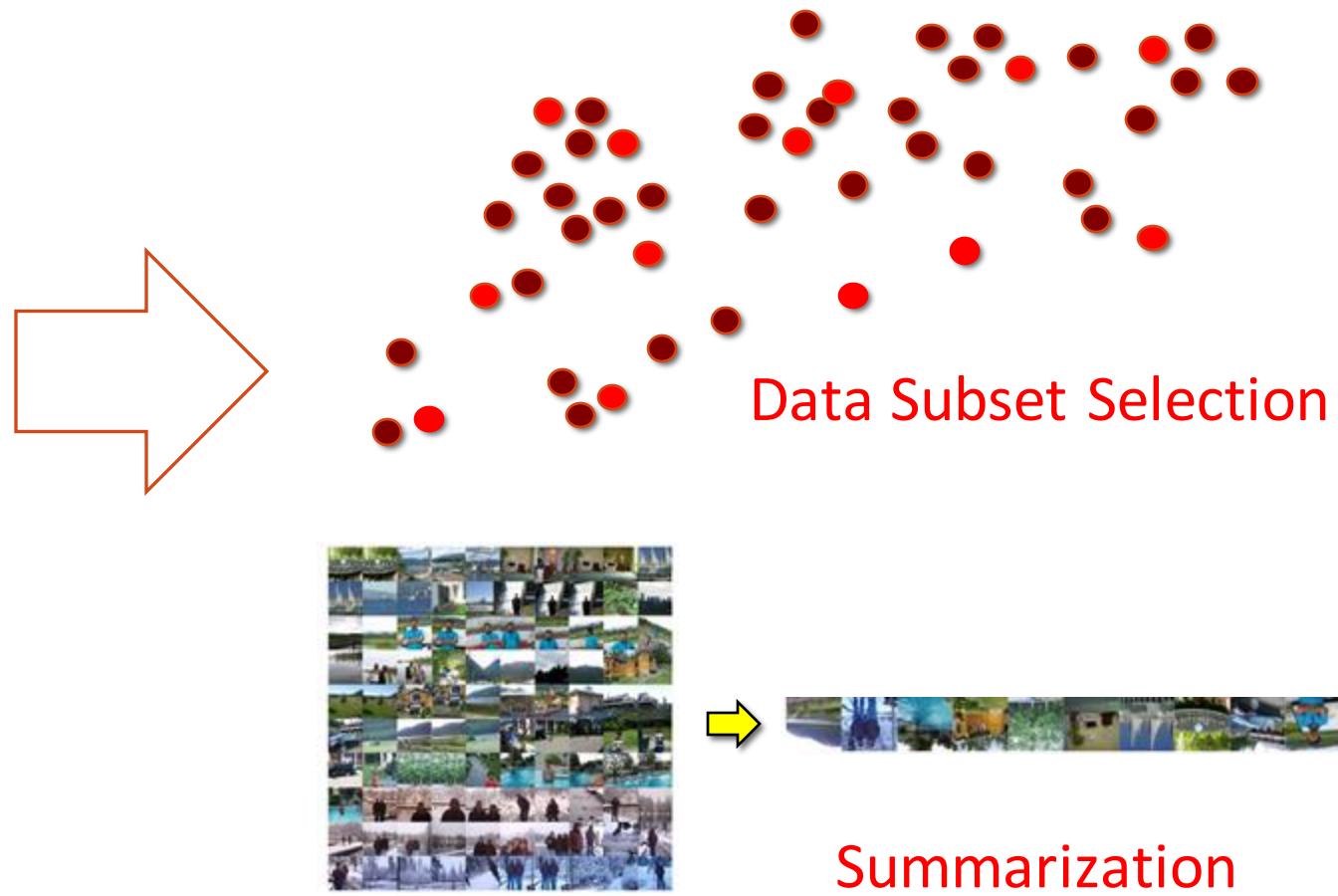
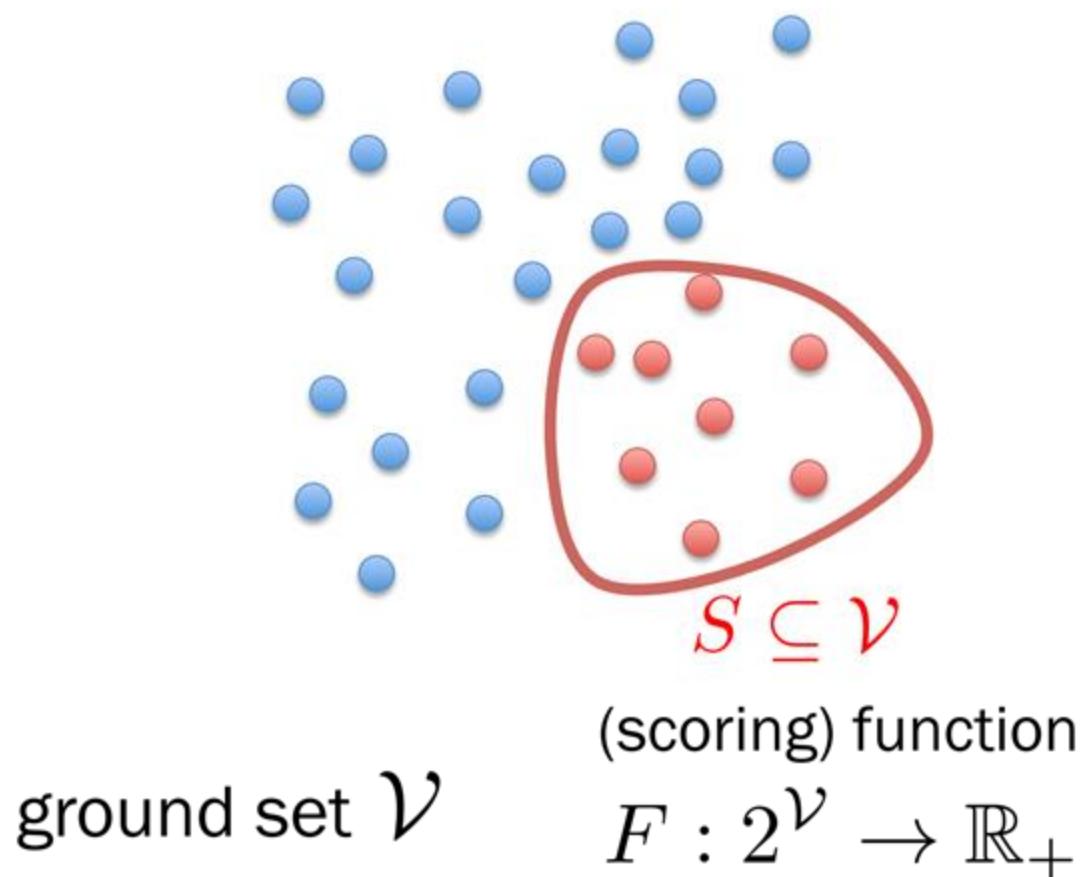
# Big Picture: Continuous and Discrete Optimization



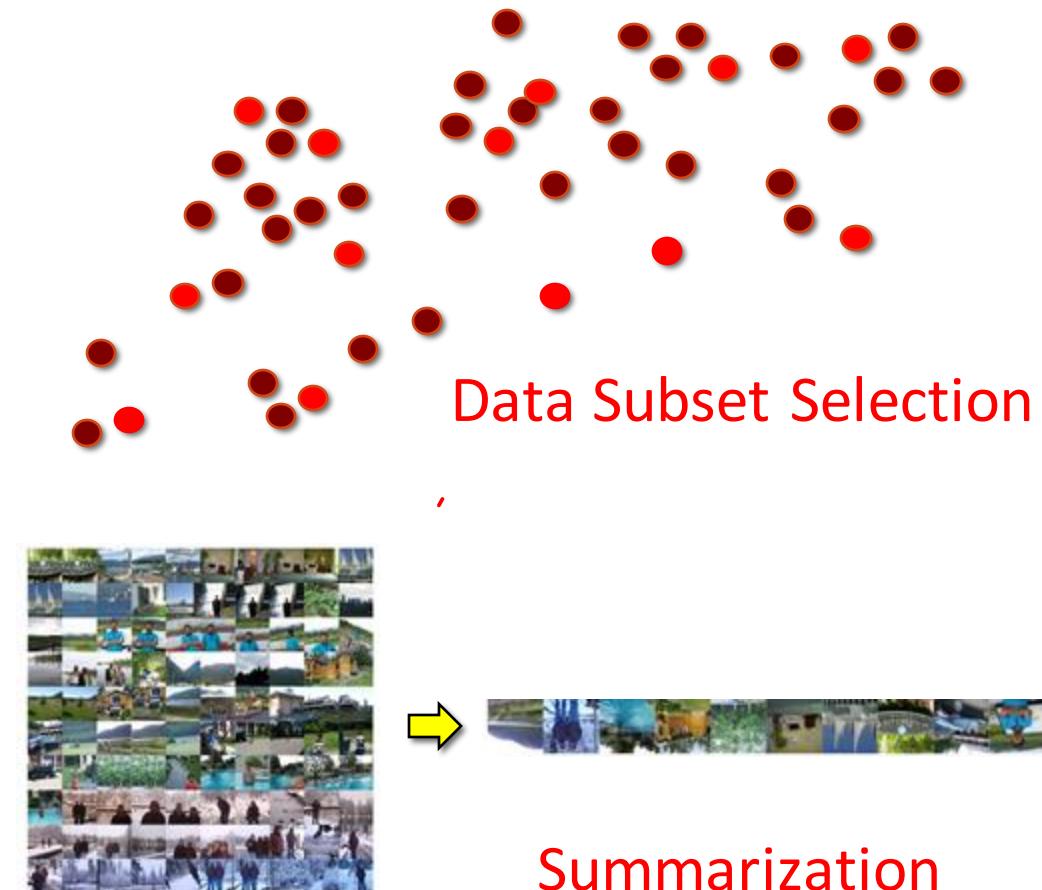
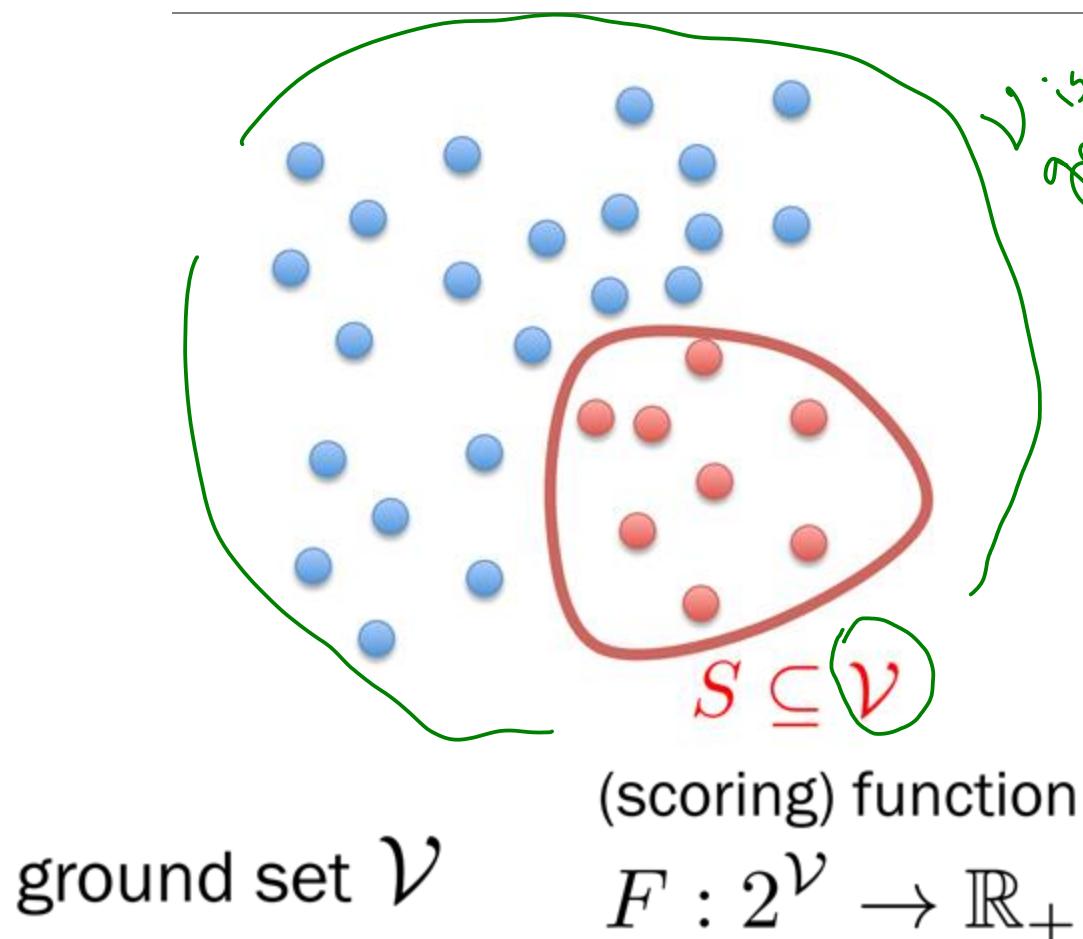
# Big Picture: Continuous and Discrete Optimization



# Discrete Optimization



# Discrete Optimization



# Acknowledgements

---

Slides borrowed from several sources:

1. Submodular Optimization course at UW from Jeff Bilmes
2. Tutorial on Submodular Optimization by Stefanie Jegelka, Andreas Krause and Jeff Bilmes at ICML and NIPS

# Useful Material

---

- Fujishige, “Submodular Functions and Optimization”, 2005
- Narayanan, “Submodular Functions and Electrical Networks”, 1997
- Welsh, “Matroid Theory”, 1975
- Oxley, “Matroid Theory”, 1992 (and 2011).
- Lawler, “Combinatorial Optimization: Networks and Matroids”, 1976.
- Schrijver, “Combinatorial Optimization”, 2003
- Gruenbaum, “Convex Polytopes, 2nd Ed”, 2003.
- Additional readings that will be announced here.

Slides borrowed from several sources:

# Acknowledgements

1. Submodular Optimization course at UW from Jeff Bilmes
2. Tutorial on Submodular Optimization by Stefanie Jegelka, Andreas Krause and Jeff Bilmes at ICML and NIPS

## Useful material

---

- Some custom slides. Also slides from our tutorials at
  - IJCAI 2020: <https://sites.google.com/view/ijcaitutorial2020summarization/home>
  - ECAI 2020: <https://sites.google.com/view/ecaitutorial2020summ/home>
- Jeff's Class: [https://people.ece.uw.edu/bilmes/classes/ee563\\_spring\\_2018/](https://people.ece.uw.edu/bilmes/classes/ee563_spring_2018/)
- Stefanie Jegelka & Andreas Krause's 2013 ICML tutorial: <http://techtalks.tv/talks/submodularity-in-machine-learning-new-directions-part-i/58125/>
- Jeff's NIPS, 2013 tutorial on submodularity: <http://melodi.ee.washington.edu/~bilmes/pgs/b2hd-bilmes2013-nips-tutorial.html> and <http://youtu.be/c4rBof38nKQ>
- Andreas Krause's web page: <http://submodularity.org>
- Francis Bach's updated 2013 text: [http://hal.archives-ouvertes.fr/docs/00/87/06/09/PDF/submodular\\_fot\\_revisd\\_hal.pdf](http://hal.archives-ouvertes.fr/docs/00/87/06/09/PDF/submodular_fot_revisd_hal.pdf)
- Tom McCormick's overview paper on submodular minimization:  
<http://people.commerce.ubc.ca/faculty/mccormick/sfmchap8a.pdf>

# Discrete Optimization in Machine Learning

---

- MAP inference in Probabilistic Models: Ising Models, DPPs
  - Feature Subset Selection
  - Data Partitioning
  - Data Subset Selection
  - Data Summarization: Text, Images, Video Summarization
  - Social networks, Influence Maximization
  - Natural Language Processing: words, phrases, n-grams, syntax trees, semantic structures
  - Computer Vision: Image Segmentation, Image Correspondence
  - Genomics and Computational Biology: cell types or assay selection, selecting peptides and proteins
- Recap: Examples from our earlier lectures

# Outline

---

- ❑ Discrete Optimization in Machine Learning
- ❑ Definition and Intuition of Submodularity
  - ❑ Modeling Power of Submodular/Set Functions
  - ❑ Examples of Submodular Functions
  - ❑ Examples of Submodular Optimization
- ❑ Optimization Algorithms for Different Function Classes and Constraints
  - ❑ Practical Applications

# Combinatorial Subset Selection Problems

$$V = \{ \text{Banana}, \text{Milk}, \text{Apple}, \text{Car}, \text{Laptop}, \text{Strawberry}, \text{T-shirt}, \text{Book}, \text{Coffee} \}$$

$$f : 2^V \rightarrow \mathbb{R}$$

$$A = \{ \text{Banana}, \text{Strawberry}, \text{Apple}, \text{Book} \}$$

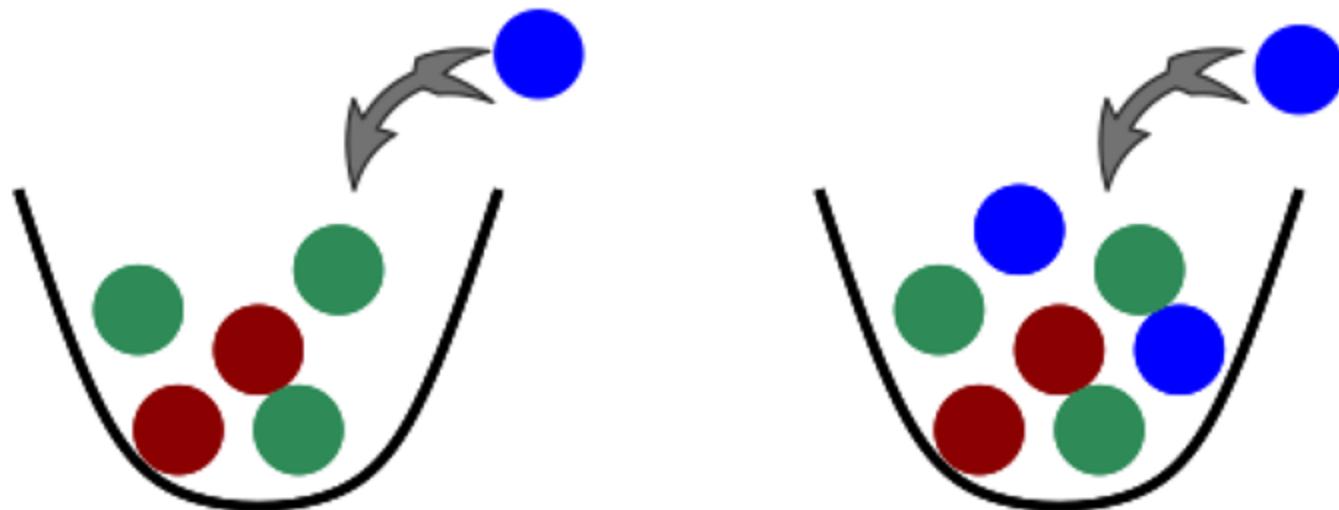
Choose Subset  $A \subseteq V$   
 $f(A) = 22$

General Set function Optimization: very hard!

What if there is some special structure?

# Submodular Functions

$$f(A \cup v) - f(A) \geq f(B \cup v) - f(B), \text{ if } A \subseteq B$$

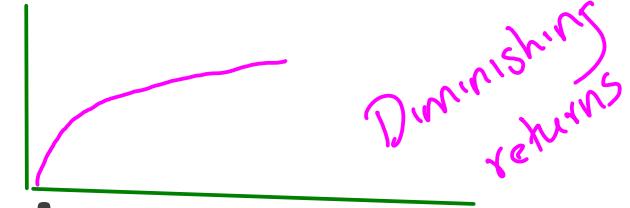


$f = \#$  of distinct colors of balls in the urn.

Negative of a  
Submodular  
Function is a  
Super-modular  
Function!

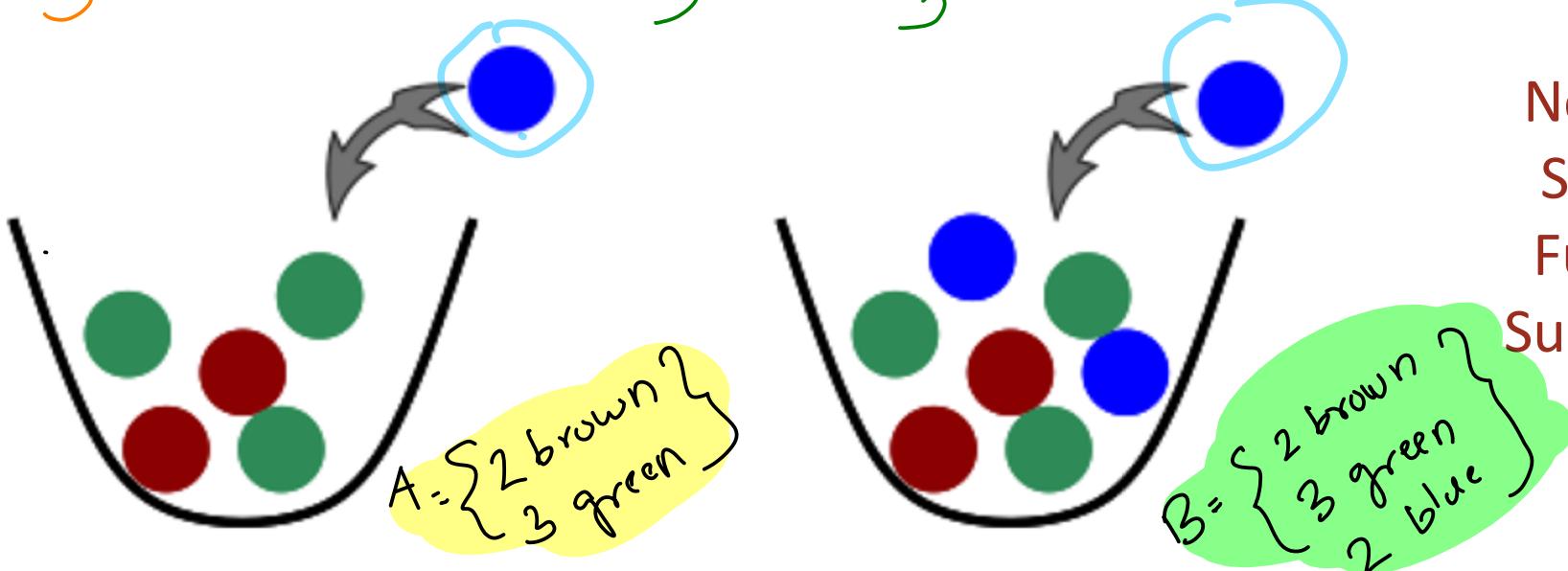
Q: Is it necessary for a submodular fn f  
that  $f(A) \leq f(B)$  if  $A \subseteq B$ ?

Prove or disprove using  
a counter example



# Submodular Functions

$$1 = f(A \cup v) - f(A) \geq f(B \cup v) - f(B), \text{ if } A \subseteq B$$



Negative of a  
Submodular  
Function is a  
Super-modular  
Function!

$f = \# \text{ of distinct colors of balls in the urn.}$

# Submodular Functions

$$f(\text{🍟🥤}) - f(\text{🍟}) \geq f(\text{🍟🍔🥤}) - f(\text{🍔})$$

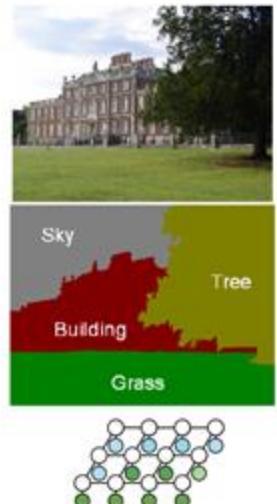
Diminishing Returns!



The more items  
you buy,  
the more the  
discount!

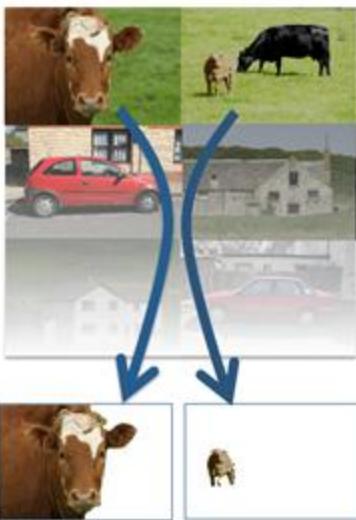
# Submodular Optimization in Machine Learning

---



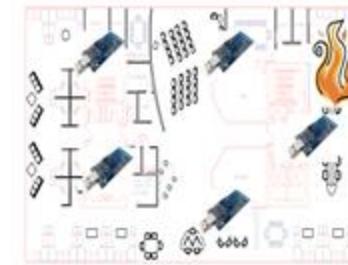
$F(S) = \text{coherence} + \text{likelihood}$

Discrete Labeling



$F(S) = \text{relevance} + \text{diversity or coverage}$

Summarization



- where put sensors?
- which experiments?
- summarization

$F(S) = \text{"information"}$

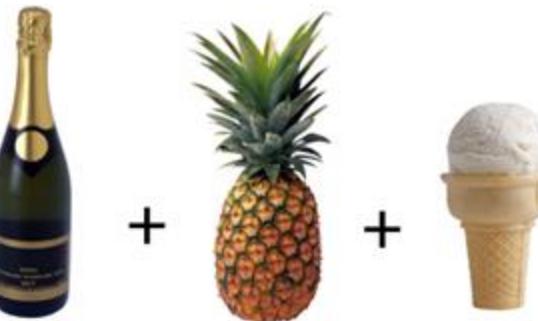
Sensor Placement

# Modular Functions

---

- each element  $e$  has a weight  $w(e)$

$$F(S) = \sum_{e \in S} w(e)$$



$$A \subset B$$

$$F(A \cup e) - F(A) = w(e) \quad = \quad F(B \cup e) - F(B) = w(e)$$

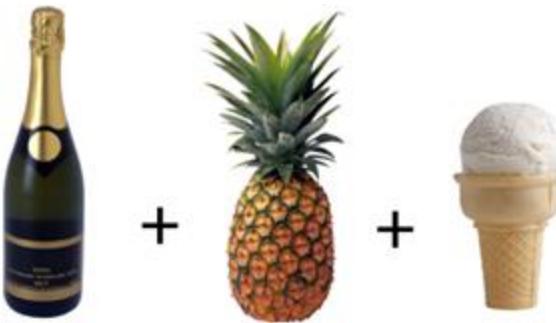
Modular Functions are both submodular and super-modular!

# Modular Functions

---

- each element  $e$  has a weight  $w(e)$

$$F(S) = \sum_{e \in S} w(e)$$



$$A \subset B$$

$$F(A \cup e) - F(A) = w(e) \quad = \quad F(B \cup e) - F(B) = w(e)$$

Modular Functions are both submodular and super-modular!

# Monotone Submodular Functions

---

- A set function is called **monotonic** if

$$A \subseteq B \subseteq V \Rightarrow F(A) \leq F(B)$$

- Examples:

- **Influence** in social networks [Kempe et al KDD '03]

- For discrete RVs, **entropy**  $F(A) = H(X_A)$  is monotonic:  
Suppose  $B = A \cup C$ . Then

$$F(B) = H(X_A, X_C) = H(X_A) + H(X_C | X_A) \geq H(X_A) = F(A)$$

- **Information gain**:  $F(A) = H(Y) - H(Y | X_A)$

Q: Is it necessary for a submodular function  $f$  that  $f(A) \leq f(B)$  if  $A \subseteq B$ ?

Prove or disprove using a counter example

## Monotone Submodular Functions

- A set function is called **monotonic** if

$$A \subseteq B \subseteq V \Rightarrow F(A) \leq F(B)$$

If a modular function has non-negative weights it will be monotone

Function value should not decrease as you grow the set

Eg: Cardinality function which is also modular

- Examples:

- Influence in social networks [Kempe et al KDD '03]

- For discrete RVs, **entropy**  $F(A) = H(X_A)$  is monotonic:

Suppose  $B = A \cup C$ . Then

$$F(B) = H(X_A, X_C) = H(X_A) + H(X_C | X_A) \geq H(X_A) = F(A)$$

H/w. Verify

- Information gain:  $F(A) = H(Y) - H(Y | X_A)$

Always monotone. But only under some (Naive Bayes like conditional independence) conditions is it submodular

Suppose entropy were both monotone and submodular, it would still not mean that there is no function that is submodular but not monotone (or vice versa)

# *Not every monotone function is submodular*

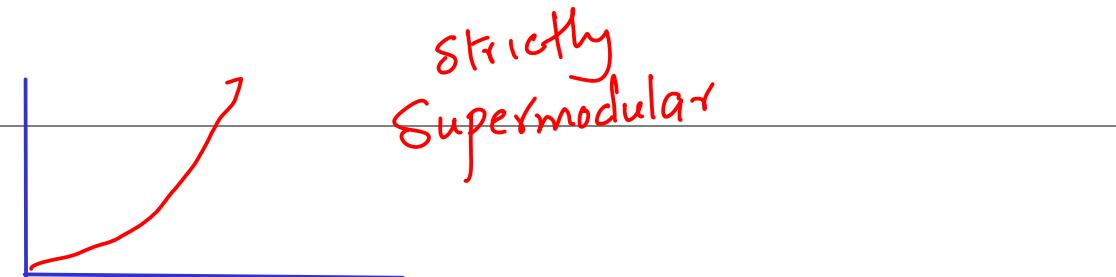
---

- Eg: 

	F <sub>1</sub> (A)
{}	0
{a}	1
{b}	1
{a,b}	3

 - this is a form of attractive potentials...Monotone but not submodular
  - In fact it is (strictly) supermodular.
  - Its negative is (strictly) submodular.
- Eg: Conventional Information gain based on entropy – it is a difference of two submodular functions
- Not every submodular function is monotone
  - We will see examples of diversity based functions such as logDet and disparity-min-sum

# *Not every monotone function is submodular*



- Eg: 

	$F_1(A)$
$\{\}$	0
$\{a\}$	1
$\{b\}$	1
$\{a,b\}$	3

  - this is a form of attractive potentials...Monotone but not submodular
  - In fact it is (strictly) supermodular.
  - Its negative is (strictly) submodular.

$V = \{a, b\}$
- Eg: Conventional Information gain based on entropy – it is a difference of two submodular functions
- Not every submodular function is monotone
  - We will see examples of diversity based functions such as logDet and disparity-min-sum

# Summary of Function Classes

General Set Functions

Supermodular Functions

Modular Functions

Submodular Functions

Monotone  
Submodular  
Functions

Non-  
Monotone  
Submodular  
Functions

eg: disparity  
 $\min/\sum$   
Dispersion  
Functions

# Properties of Submodular Functions

---

- ❑ Convex Combinations of Submodular Functions are Submodular
- ❑ Intersections with Fixed Sets is Submodular (Restrictions)
- ❑ Unions with Fixed Sets is Submodular (Conditioning)
- ❑ Complement Functions are Submodular (Reflection)
- ❑ Two equivalent definitions of submodular functions
- ❑ Minimum and Maximum of Submodular Functions
- ❑ Information functions
- ❑ Submodularity and Convexity
- ❑ Submodularity and Concavity

# Properties of Submodular Functions

---

- ❑ Convex Combinations of Submodular Functions are Submodular
- ❑ Intersections with Fixed Sets is Submodular (Restrictions)
- ❑ Unions with Fixed Sets is Submodular (Conditioning)
- ❑ Complement Functions are Submodular (Reflection)
- ❑ Two equivalent definitions of submodular functions
- ❑ Minimum and Maximum of Submodular Functions ?
- ❑ Information functions
- ❑ Submodularity and Convexity
- ❑ Submodularity and Concavity

[Do you recall any of these from our discussions during the course? Any projects in this space?]

# Convex Combinations of Submodular Functions

---

$F_1, \dots, F_m$  submodular functions on  $V$  and  $\lambda_1, \dots, \lambda_m > 0$

Then:  $F(A) = \sum_i \lambda_i F_i(A)$  is submodular

Submodularity closed under nonnegative linear combinations!

Extremely useful fact:

- $F_\theta(A)$  submodular  $\rightarrow \sum_\theta P(\theta) F_\theta(A)$  submodular!
- Multicriterion optimization
- A basic proof technique! ☺

# Convex Combinations of Submodular Functions

$F_1, \dots, F_m$  submodular functions on  $V$  and  $\lambda_1, \dots, \lambda_m > 0$

Then:  $F(A) = \sum_i \lambda_i F_i(A)$  is submodular

$$\sum_i \lambda_i [F_i(\overbrace{A \cup \{v\}}) - F_i(A)] \geq \left[ \sum_i \lambda_i [F_i(B \cup \{v\}) - F_i(B)] \right] \quad \forall A \subseteq B$$

Submodularity closed under nonnegative linear combinations!

Extremely useful fact:

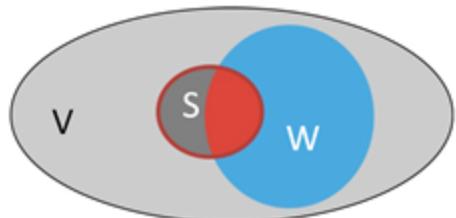
- $F_\theta(A)$  submodular  $\rightarrow \sum_\theta P(\theta) F_\theta(A)$  submodular!
- Multicriterion optimization
- A basic proof technique! 😊

# More Properties

---

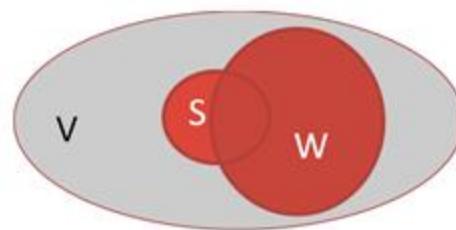
- **Restriction:**  $F(S)$  submodular on  $V$ ,  $W$  subset of  $V$

Then  $F'(S) = F(S \cap W)$  is submodular



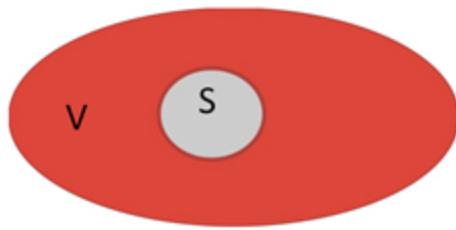
- **Conditioning:**  $F(S)$  submodular on  $V$ ,  $W$  subset of  $V$

Then  $F'(S) = F(S \cup W)$  is submodular



- **Reflection:**  $F(S)$  submodular on  $V$

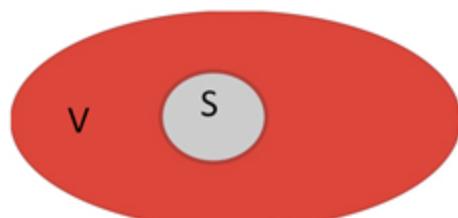
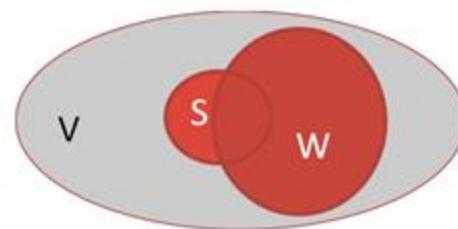
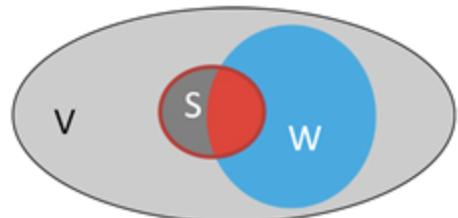
Then  $F'(S) = F(V \setminus S)$  is submodular



# More Properties

- **Restriction:**  $F(S)$  submodular on  $V$ ,  $W$  subset of  $V$   
Then  $F'(S) = F(S \cap W)$  is submodular
- **Conditioning:**  $F(S)$  submodular on  $V$ ,  $W$  subset of  $V$   
Then  $F'(S) = F(S \cup W)$  is submodular
- **Reflection:**  $F(S)$  submodular on  $V$   
Then  $F'(S) = F(V \setminus S)$  is submodular

$F'(S)$  is the submodular function  $F$  applied to the complement of  $S$

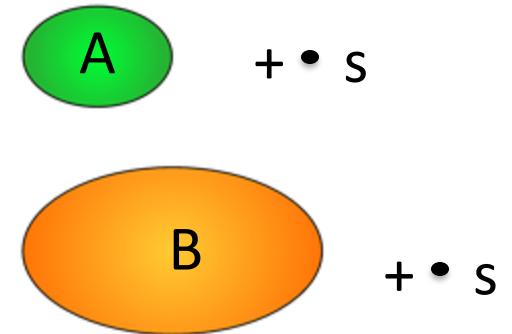


# Two Equivalent Definitions of Submodularity

---

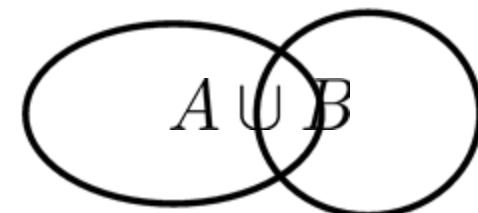
- Diminishing gains: for all  $A \subseteq B \subseteq V$

$$F(A \cup s) - F(A) \geq F(B \cup s) - F(B)$$



- 
- Union-Intersection: for all  $A, B \subseteq V$

$$F(A) + F(B) \geq F(A \cup B) + F(A \cap B)$$



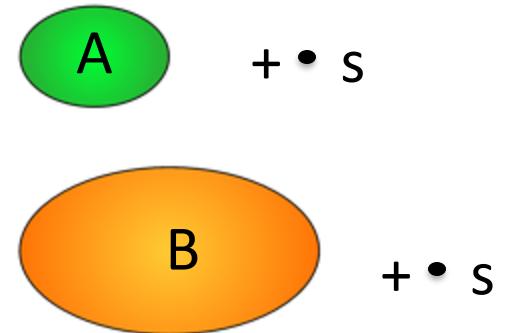
Can you try and prove that the two above definitions of submodularity are equivalent?

# Two Equivalent Definitions of Submodularity

- Diminishing gains: for all  $A \subseteq B \subseteq V$

$$F(A \cup s) - F(A) \geq F(B \cup s) - F(B)$$

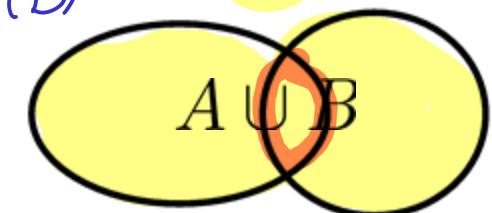
Already seen



- Union-Intersection: for all  $A, B \subseteq V$

$$F(A) + F(B) \geq F(A \cup B) + F(A \cap B)$$

If  $F$  were cardinality  
 $F(A) + F(B) = F(A \cup B) + F(A \cap B)$



Can you try and prove that the two above definitions of submodularity are equivalent?

## Step 1: Define sets

Let:

- $X = A \cap B$
- $Y = A \setminus B$
- $Z = B \setminus A$

## Step 2: Expand $F(A)$ and $F(B)$ from $F(X)$

We write  $F(A)$  as:

$$F(A) = F(X) + \sum_{y \in Y} [F(X \cup \{y_1, \dots, y_i\}) - F(X \cup \{y_1, \dots, y_{i-1}\})]$$

Similarly,

Then:

- $A = X \cup Y$
- $B = X \cup Z$
- $A \cup B = X \cup Y \cup Z$

$$F(B) = F(X) + \sum_{z \in Z} [F(X \cup \{z_1, \dots, z_i\}) - F(X \cup \{z_1, \dots, z_{i-1}\})]$$

Let's denote:

- $F(A) = F(X) + \Delta_Y$
- $F(B) = F(X) + \Delta_Z$

HOMEWORK: COMPLETE THE ABOVE STEPS OF DERIVATION