

**Two SVM tutorials linked in class website  
(please, read both):**

- High-level presentation with applications (Hearst 1998)
- Detailed tutorial (Burges 1998)

# SVMs, Duality and the Kernel Trick

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 27<sup>th</sup>, 2006

# Announcements

- Third homework

- is out

- Due March 1<sup>st</sup>

*start early !!*

- Final assigned by registrar:

- May 12, 1-4p.m *Friday*

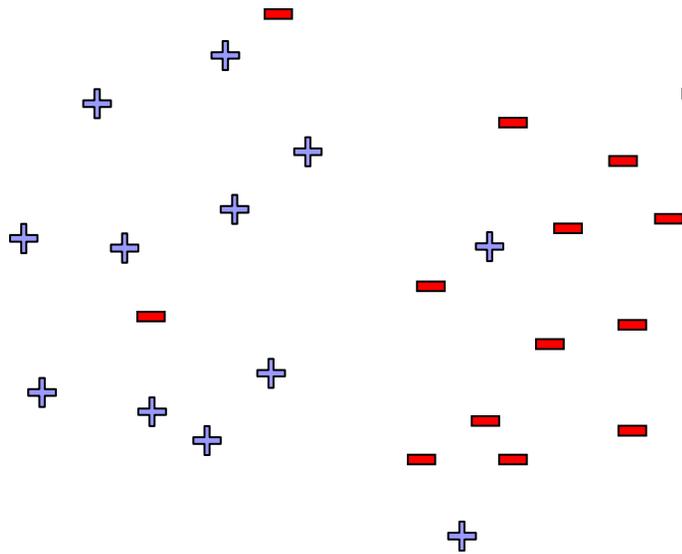
- Location TBD

- Midterm

- March 8<sup>th</sup>, a week from Wednesday

- Open book, notes, papers, etc. No computers

# SVMs reminder



$$\begin{aligned} \text{minimize}_{\mathbf{w}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j \\ & \xi_j \geq 0, \quad \forall j \end{aligned}$$

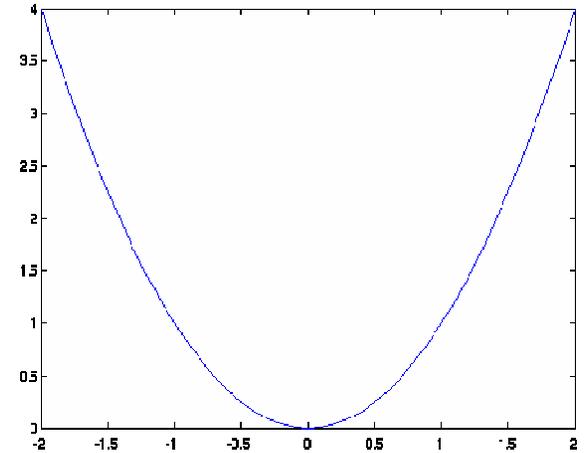
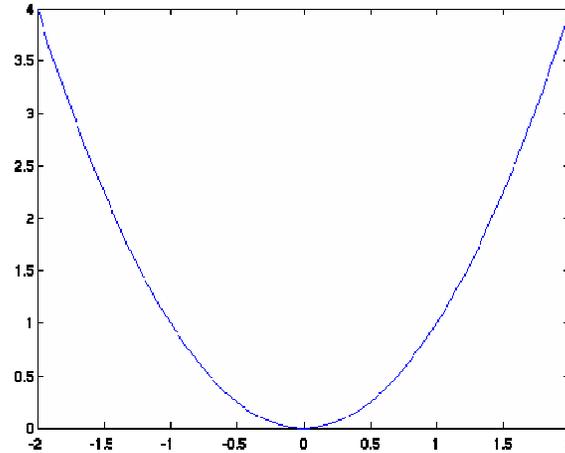
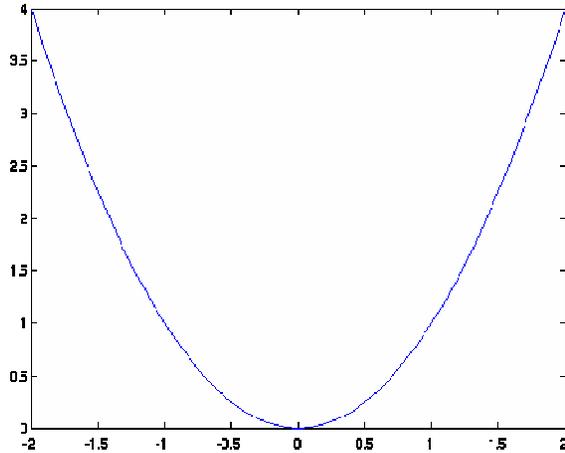
# Today's lecture



- Learn one of the most interesting and exciting recent advancements in machine learning
  - The “kernel trick”
  - High dimensional feature spaces at no extra cost!
- But first, a detour
  - Constrained optimization!

# Constrained optimization

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq b \end{aligned}$$



# Lagrange multipliers – Dual variables

$$\min_x x^2$$

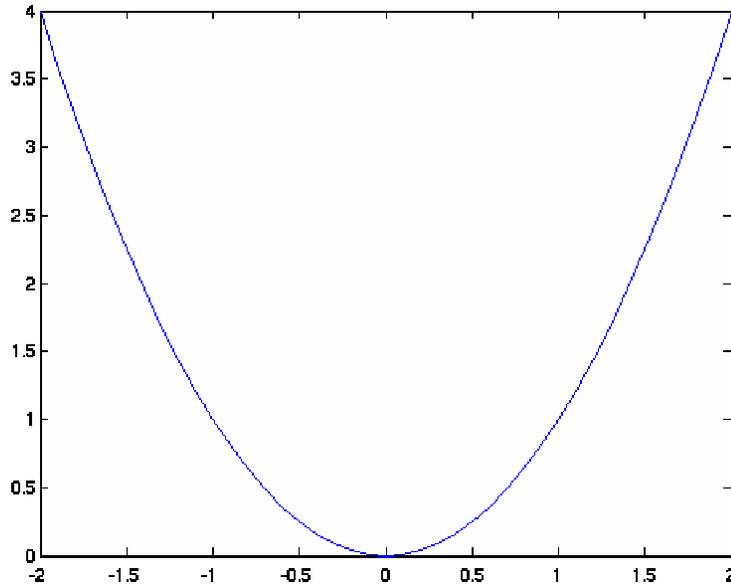
$$\text{s.t. } x \geq b$$

Moving the constraint to objective function

Lagrangian:

$$L(x, \alpha) = x^2 - \alpha(x - b)$$

$$\text{s.t. } \alpha \geq 0$$

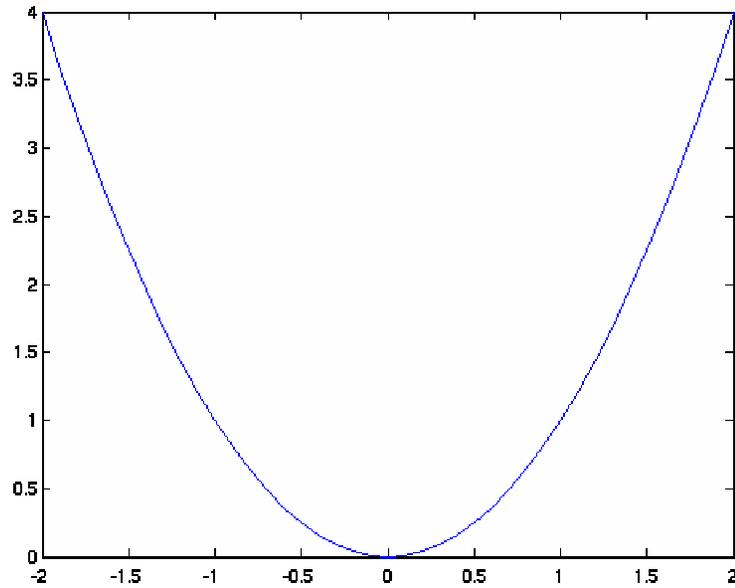
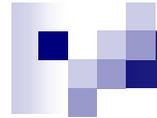


Solve:

$$\min_x \max_{\alpha} L(x, \alpha)$$

$$\text{s.t. } \alpha \geq 0$$

# Lagrange multipliers – Dual variables



**Solving:**  $\min_x \max_{\alpha} x^2 - \alpha(x - b)$   
s.t.  $\alpha \geq 0$

# Dual SVM derivation (1) – the linearly separable case

$$\begin{aligned} &\text{minimize}_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w} \cdot \mathbf{w} \\ &(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j \end{aligned}$$

# Dual SVM derivation (2) – the linearly separable case

$$L(\mathbf{w}, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j \left[ (\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1 \right]$$
$$\alpha_i \geq 0, \quad \forall j$$

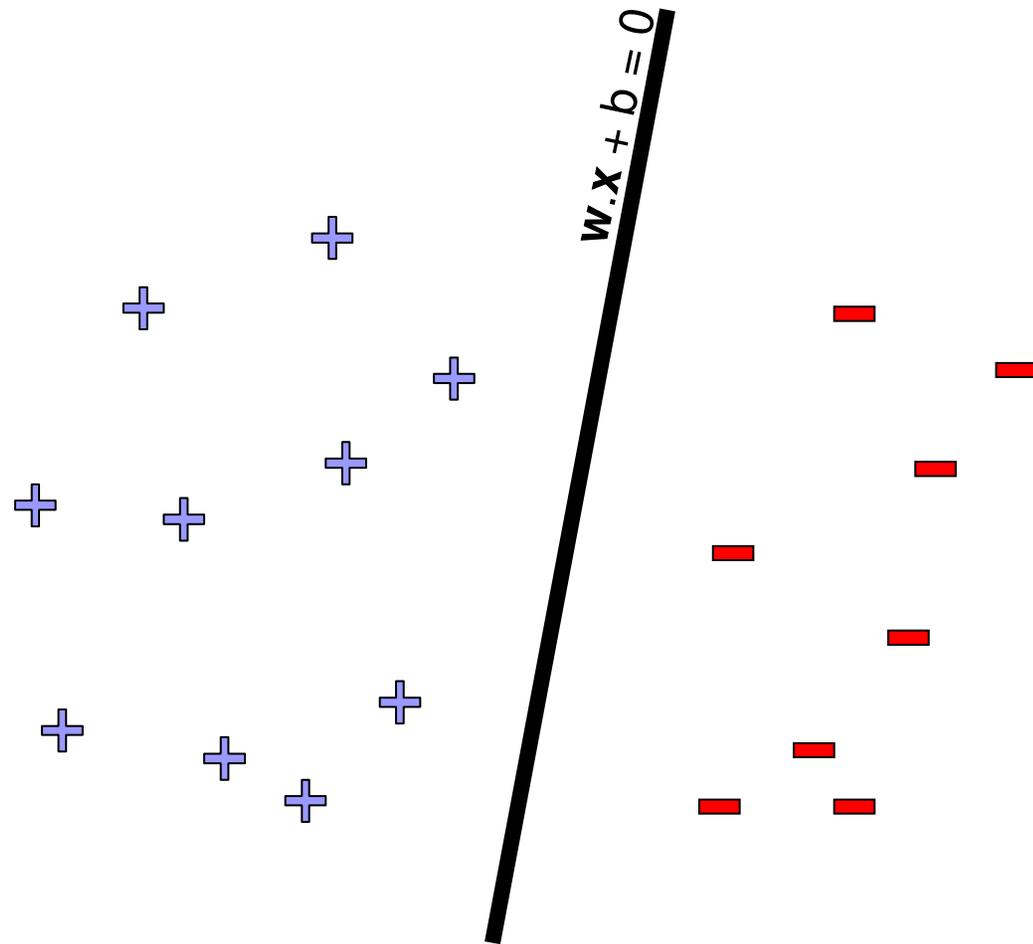
$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\text{minimize}_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w} \cdot \mathbf{w}$$
$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j$$

$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

for any  $k$  where  $\alpha_k > 0$

# Dual SVM interpretation



$$w = \sum_i \alpha_i y_i x_i$$

# Dual SVM formulation – the linearly separable case

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

for any  $k$  where  $\alpha_k > 0$

# Dual SVM derivation – the non-separable case

$$\begin{aligned} \text{minimize}_{\mathbf{w}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ \left( \mathbf{w} \cdot \mathbf{x}_j + b \right) y_j & \geq 1 - \xi_j, \quad \forall j \\ \xi_j & \geq 0, \quad \forall j \end{aligned}$$

# Dual SVM formulation – the non-separable case

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

for any  $k$  where  $C > \alpha_k > 0$

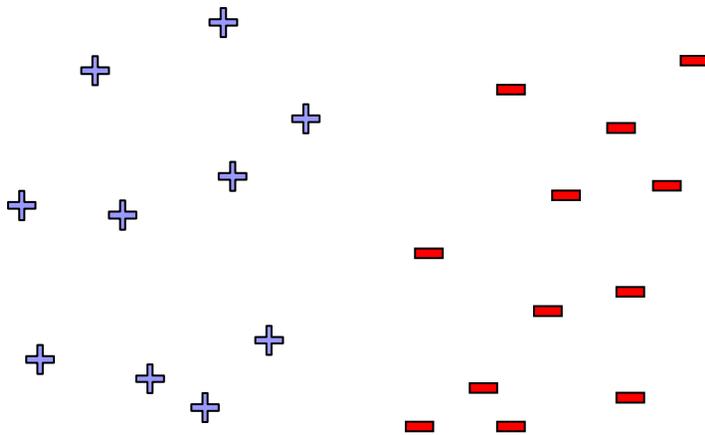
# Why did we learn about the dual SVM?

- There are some quadratic programming algorithms that can solve the dual faster than the primal
- But, more importantly, the “**kernel trick**”!!!
  - Another little detour...

# Reminder from last time: What if the data is not linearly separable?

**Use features of features  
of features of features....**

$$\Phi(\mathbf{x}) : \mathbb{R}^m \mapsto F$$

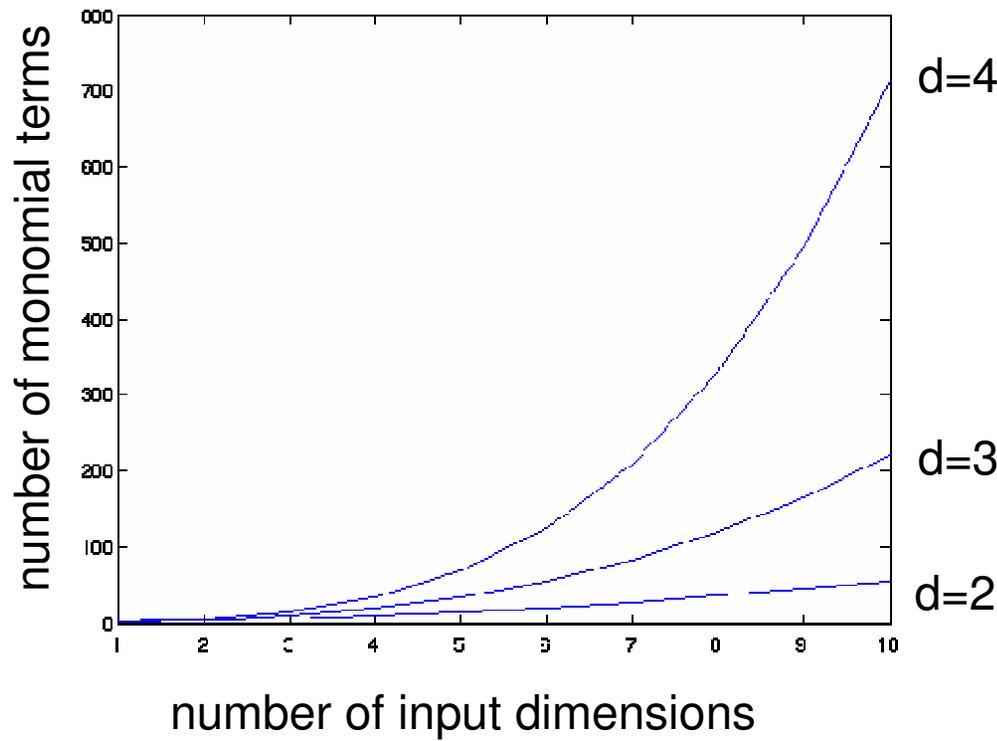


**Feature space can get really large really quickly!**

# Higher order polynomials

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!}$$

m – input features  
d – degree of polynomial



grows fast!  
d = 6, m = 100  
about 1.6 billion terms

# Dual formulation only depends on dot-products, not on $\mathbf{w}$ !

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

# Dot-product of polynomials



$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) =$  polynomials of degree  $d$

# Finally: the “kernel trick”!

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

- Never represent features explicitly
  - Compute dot products in closed form
- Constant-time high-dimensional dot-products for many classes of features
  
- Very interesting theory – Reproducing Kernel Hilbert Spaces
  - Not covered in detail in 10701/15781, more in 10702

$$\mathbf{w} = \sum_i \alpha_i y_i \Phi(\mathbf{x}_i)$$

$$b = y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k)$$

for any  $k$  where  $C > \alpha_k > 0$

# Polynomial kernels

- All monomials of degree  $d$  in  $O(d)$  operations:

$$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d = \text{polynomials of degree } d$$

- How about all monomials of degree up to  $d$ ?

- Solution 0:

- Better solution:

# Common kernels

- Polynomials of degree  $d$

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$$

- Polynomials of degree up to  $d$

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$$

- Gaussian kernels

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

# Overfitting?



- Huge feature space with kernels, what about overfitting???
  - Maximizing margin leads to sparse set of support vectors
  - Some interesting theory says that SVMs search for simple hypothesis with large margin
  - Often robust to overfitting

# What about at classification time

- For a new input  $\mathbf{x}$ , if we need to represent  $\Phi(\mathbf{x})$ , we are in trouble!
- Recall classifier:  $\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$
- Using kernels we are cool!

$$K(\mathbf{u}, \mathbf{v}) = \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v})$$

$$\mathbf{w} = \sum_i \alpha_i y_i \Phi(\mathbf{x}_i)$$

$$b = y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k)$$

for any  $k$  where  $C > \alpha_k > 0$

# SVMs with kernels

- Choose a set of features and kernel function
- Solve dual problem to obtain support vectors  $\alpha_i$
- At classification time, compute:

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i)$$

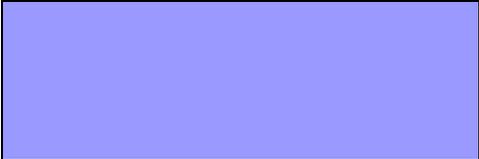
$$b = y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i)$$

for any  $k$  where  $C > \alpha_k > 0$

Classify as

$$\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$$

# What's the difference between SVMs and Logistic Regression?

	<b>SVMs</b>	<b>Logistic Regression</b>
<b>Loss function</b>		
<b>High dimensional features with kernels</b>		

# Kernels in logistic regression

$$P(Y = 1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)}}$$

- Define weights in terms of support vectors:

$$\mathbf{w} = \sum_i \alpha_i \Phi(\mathbf{x}_i)$$

$$\begin{aligned} P(Y = 1 | \mathbf{x}, \mathbf{w}) &= \frac{1}{1 + e^{-(\sum_i \alpha_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b)}} \\ &= \frac{1}{1 + e^{-(\sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b)}} \end{aligned}$$

- Derive simple gradient descent rule on  $\alpha_i$

# What's the difference between SVMs and Logistic Regression? (Revisited)

	<b>SVMs</b>	<b>Logistic Regression</b>
<b>Loss function</b>	Hinge loss	Log-loss
<b>High dimensional features with kernels</b>	Yes!	Yes!

# What you need to know



- Dual SVM formulation
  - How it's derived
- The kernel trick
- Derive polynomial kernel
- Common kernels
- Kernelized logistic regression
- Differences between SVMs and logistic regression

# Acknowledgment



- SVM applet:
  - <http://www.site.uottawa.ca/~gcaron/applets.htm>