

# **Applied Machine Learning: Algorithms, Practice and Theory**

## **Logistic Regression and SVM**

Koby Crammer

Technion – Israel Institute of Technology



# Outline

- Naïve Bayes:
  - Generative (joint) probabilistic linear model
- Logistic regression:
  - Discriminative (conditional) probabilistic linear model
  - Log-loss
- Support Vector Machines:
  - Discriminative linear model
  - Hinge-loss
  - Large Margin, Regularization, Kernels
- Models, Inference, Learning



# Review of Setting



# Formal Setting – Data

- Instances  $\mathbf{x} \in \mathcal{X}$ 
  - Images, Sentences
- Labels  $y \in \mathcal{Y} = \{-1 ; 1\}$ 
  - Binary
- Statistical Assumption  $(X, Y) \sim P$ 
  - Used for training and evaluation
- I.I.D. Sample  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ 
  - $(\mathbf{x}_i, y_i) \sim P$



# Formal Setting - Predictions

- Continuous predictions :  $f : \mathcal{X} \rightarrow \mathbb{R}$ 
  - Label  $\text{sign}(f(\mathbf{x}))$
  - Confidence  $|f(\mathbf{x})|$
- Notation:  $f(\mathbf{x}) = \hat{y}$



# Formal Setting - Evaluation

- Loss
  - No. of mistakes

$$\ell(f(\mathbf{x}), y) \in \mathbb{R}_+$$

- Expected Loss

$$\mathbb{E}_{(X,Y) \sim P} [\ell(f(X), Y)]$$

- Averaged Loss

$$\mathbb{E}_{(X,Y) \sim S} [\ell(f(X), Y)]$$

$$= \frac{1}{m} \sum_{i=1}^m \ell(f(\mathbf{x}_i), y_i)$$



# Empirical Risk Minimization

- Goal: Minimize classification error

$$\mathbb{E}_{(X,Y) \sim P} [ \mathbf{1}_{f(X) \neq Y} ]$$

- Approximate using Sample:

$$\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{f(\mathbf{x}_i) \neq y_i}$$

- Bound with “nice” loss

$$\frac{1}{m} \sum_{i=1}^m \ell(f(\mathbf{x}_i), y_i)$$



# Loss Functions

- Natural Loss:

- Zero-One loss:

$$\ell(f(\mathbf{x}), y) = \begin{cases} 0 & y = \text{sign}(f(\mathbf{x})) \\ 1 & y \neq \text{sign}(f(\mathbf{x})) \end{cases}$$

- Real-valued-predictions loss:

- Hinge loss:

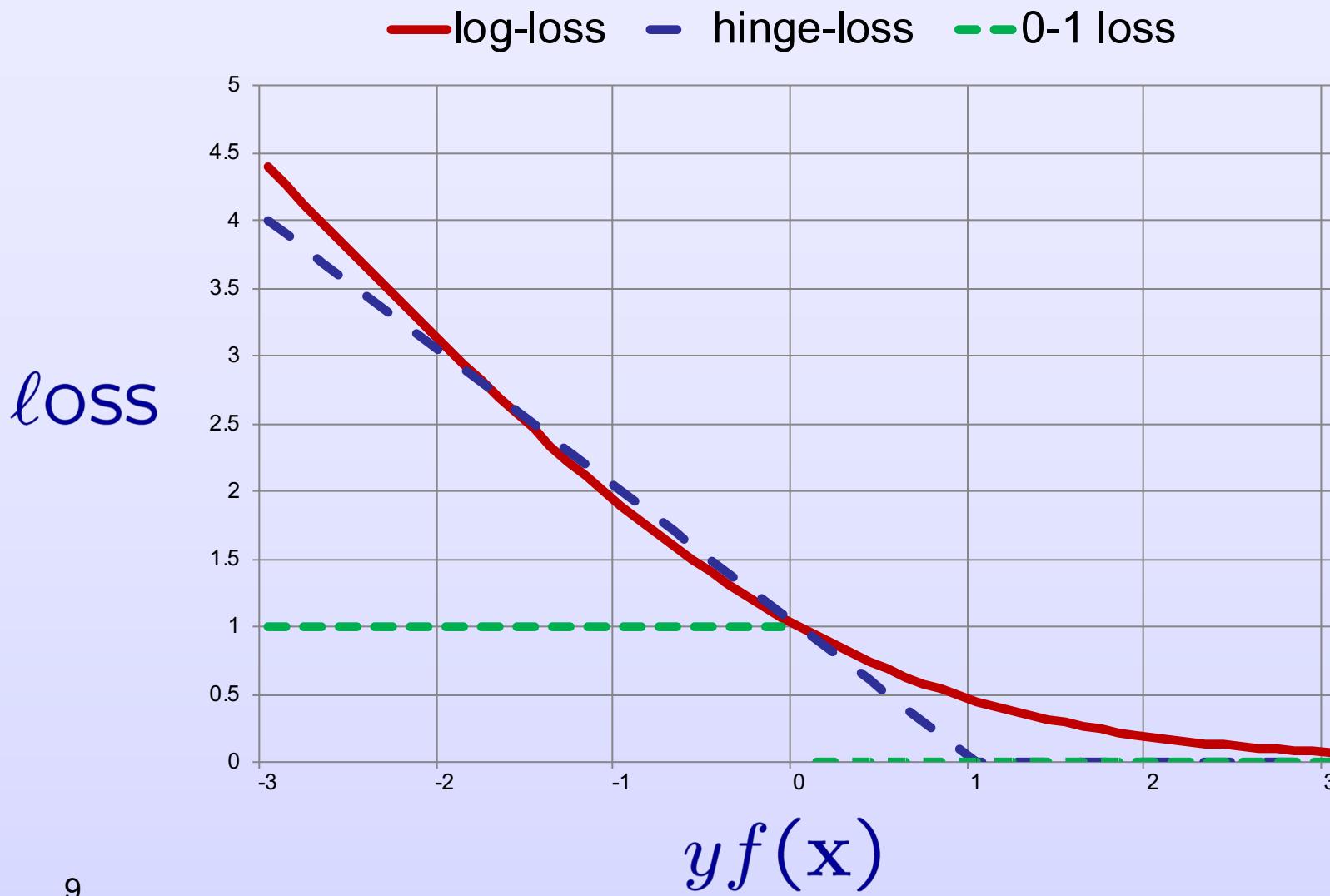
$$\ell(f(\mathbf{x}), y) = \max(0, 1 - y f(\mathbf{x}))$$

- Log loss (Max Entropy, Boosting)

$$\ell(f(\mathbf{x}), y) = \log(1 + \exp(-y f(\mathbf{x})))$$



# Loss Functions



# Logistic Regression: Model

- Goal:
  - Design a discriminative parametric model for the conditional distribution

$$P(Y|X = \mathbf{x}) \in [0, 1]$$

Binary Classification

$$Y \in \{-1, +1\}$$

Vector  
Representation of  
a Document



# Logistic Regression: Model

- Attempt I:  $P(Y = +1|\mathbf{x}) = 2\mathbf{w} \cdot \mathbf{x}$ 
  - Problem:  $P(Y = +1|\mathbf{x}) \in [0, 1]$   
 $\mathbf{w} \cdot \mathbf{x} \in (-\infty, \infty)$
- Attempt II:  $\log P(Y = +1|\mathbf{x}) = 2\mathbf{w} \cdot \mathbf{x}$ 
  - Problem:  $\log P(Y = +1|\mathbf{x}) \in (-\infty, 0]$   
 $\mathbf{w} \cdot \mathbf{x} \in (-\infty, \infty)$



# Logistic Regression



# Logistic Regression: Model

- Attempt III:

$$\log \frac{P(Y = +1|\mathbf{x})}{P(Y = -1|\mathbf{x})} = \log \frac{P(Y = +1|\mathbf{x})}{1 - P(Y = +1|\mathbf{x})} = 2\mathbf{w} \cdot \mathbf{x}$$

$\in (-\infty, \infty)$

$\in [0, \infty)$



# Logistic Regression: Model

- Model:

$$\log \frac{P(Y = +1|\mathbf{x})}{1 - P(Y = +1|\mathbf{x})} = \mathbf{w} \cdot \mathbf{x}$$

- Probability of label :

$$P(Y = +1|\mathbf{x}) = \frac{e^{2\mathbf{w} \cdot \mathbf{x}}}{1 + e^{2\mathbf{w} \cdot \mathbf{x}}} = \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{e^{-\mathbf{w} \cdot \mathbf{x}} + e^{\mathbf{w} \cdot \mathbf{x}}}$$

- Compact form:

$$P(y|\mathbf{x}) = \frac{e^{y\mathbf{w} \cdot \mathbf{x}}}{e^{-\mathbf{w} \cdot \mathbf{x}} + e^{\mathbf{w} \cdot \mathbf{x}}}$$



# Logistic Regression: Prediction

- Pretend it is the correct model, use Bayes-optimal

$$\hat{y} = \arg \max_z P(Y = z | \mathbf{x})$$

- Maximum a-posteriori



# Logistic Regression: Prediction

- Predict class +1 iff

$$P(Y=+1|x) > P(Y=-1|x)$$

Maximum a-posteriori prediction is  
linear classification

$$\hat{y} = \text{sign}(\boldsymbol{w} \cdot \boldsymbol{x})$$

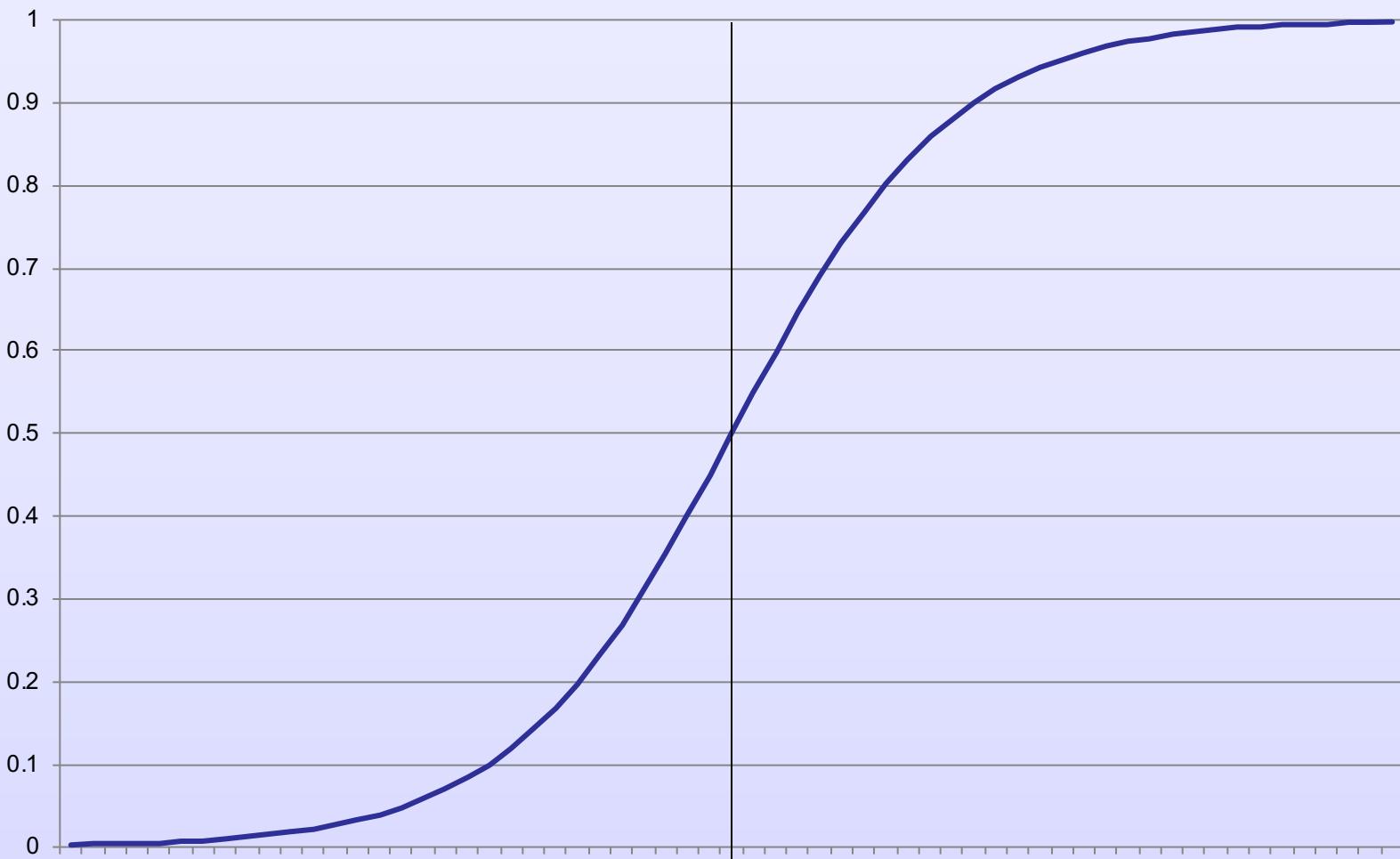
$$(Y \mid x \in \mathcal{A}_1(\omega))$$



$$\boldsymbol{w} \cdot \boldsymbol{x} > 0$$



# Logistic Regression: Prediction



# Logistic Regression: Learning

- Given training set

$$S = \{(\mathbf{x}_t, y_t)\}_{i=t}^m \quad \mathbf{x}_t \in \mathbb{R}^Q \quad y_t \in \{\pm 1\}$$

- Find a good parameters-vector:  $\mathbf{w} \in \mathbb{R}^Q$   
**Hard: not continuous, nor convex**

- How we define what is “good”?
  - Minimize training error

$$\frac{1}{m} \sum_{t=1}^m \mathbf{1}_{y_t(\mathbf{w} \cdot \mathbf{x}_t) > 0}$$



# Logistic Regression: Learning

- We have a probabilistic model, use it!
- Maximize labels-probability of sample

$$\max_{\boldsymbol{w}} P(y_1 \dots y_m | \boldsymbol{x}_1 \dots \boldsymbol{x}_m) = \prod_t P(y_t | \boldsymbol{x}_t)$$

- Use monotonic function “log:

$$\log \prod_t P(y_t | \boldsymbol{x}_t) = \sum_t \log P(y_t | \boldsymbol{x}_t)$$



# Logistic Regression: Learning

- Substitute model:

$$P(y|\mathbf{x}) = \frac{e^{y\mathbf{w} \cdot \mathbf{x}}}{e^{-\mathbf{w} \cdot \mathbf{x}} + e^{\mathbf{w} \cdot \mathbf{x}}}$$

$$\sum_t \log P(y_t | \mathbf{x}_t)$$



..



# Logistic Regression: Learning

- Substitute model:

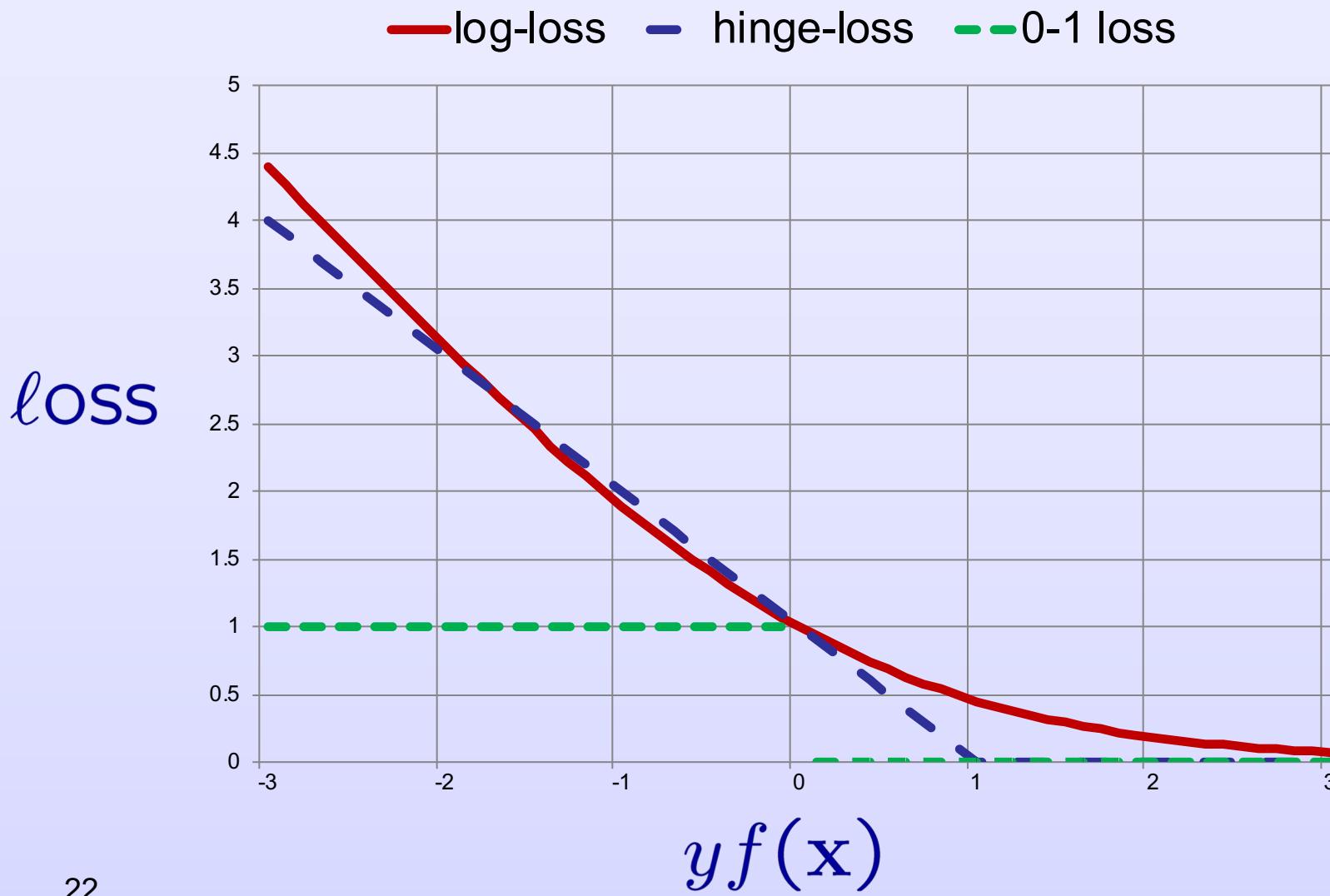
Logistic Regression is minimizing  
a sum of the log loss

$$\sum_t \log(e^{-2y_t \mathbf{w} \cdot \mathbf{x}_t} + 1)$$

$$= - \sum_t \log(e^{-2y_t \mathbf{w} \cdot \mathbf{x}_t} + 1)$$



# Loss Functions



# Logistic Regression: Comments

- Equivalent to maximum conditional-entropy
- Many algorithms to solve:
  - e.g. gradient descent
- Easley generalized to multi-class problems



# History

- Regression

- Adrien-Marie Legendre, 1805
- Johann Carl Friedrich Gauss, 1809
- John Nelder and Robert Wedderburn, 1972



- Logistic Regression

- Pierre Francois Verhulst, 1838
- Pearl & Reed, 1920
- Joseph Berkson, 1944
- E.T. Jaynes, ~1950



# Linear Models

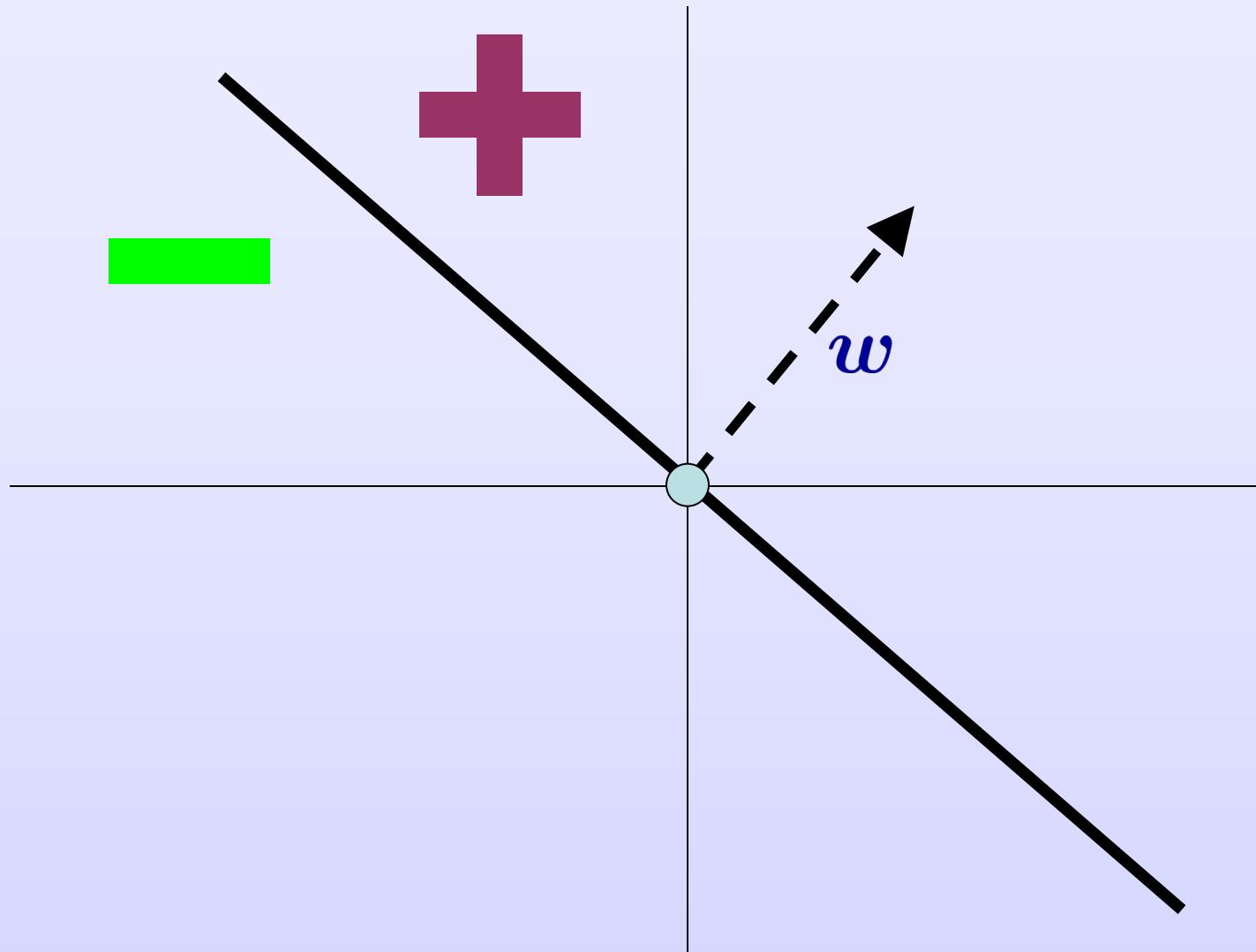
Why compute probability  
if  
only want to make predictions?



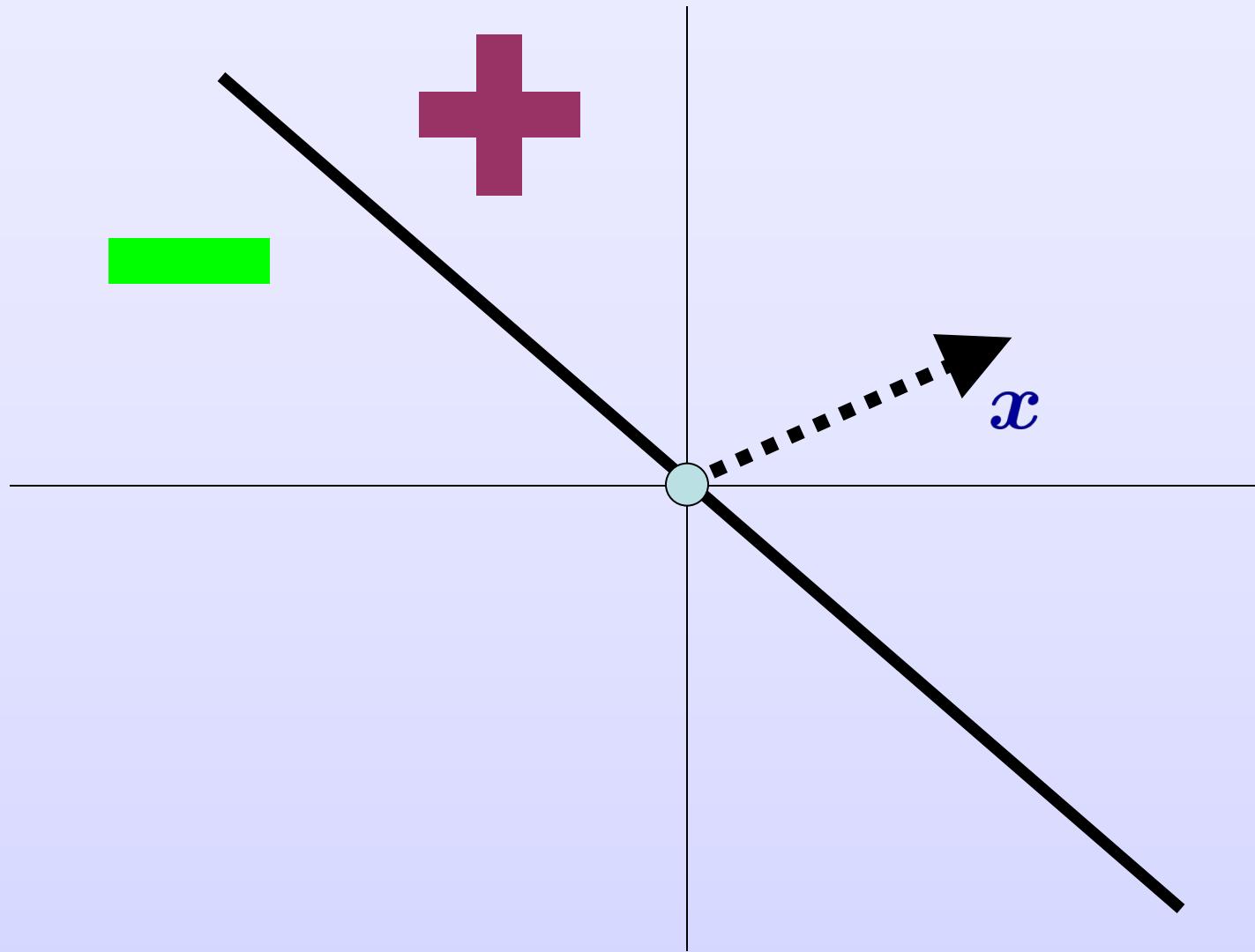
- Do not solve an estimation problem of interest by solving a more general (harder) problem as an intermediate step
- Only care about decision boundaries



# Linear Classifiers



# Linear Classifiers



# Linear Classifiers

- Prediction :

$$\hat{y} = \text{sign}(f(\mathbf{x}))$$

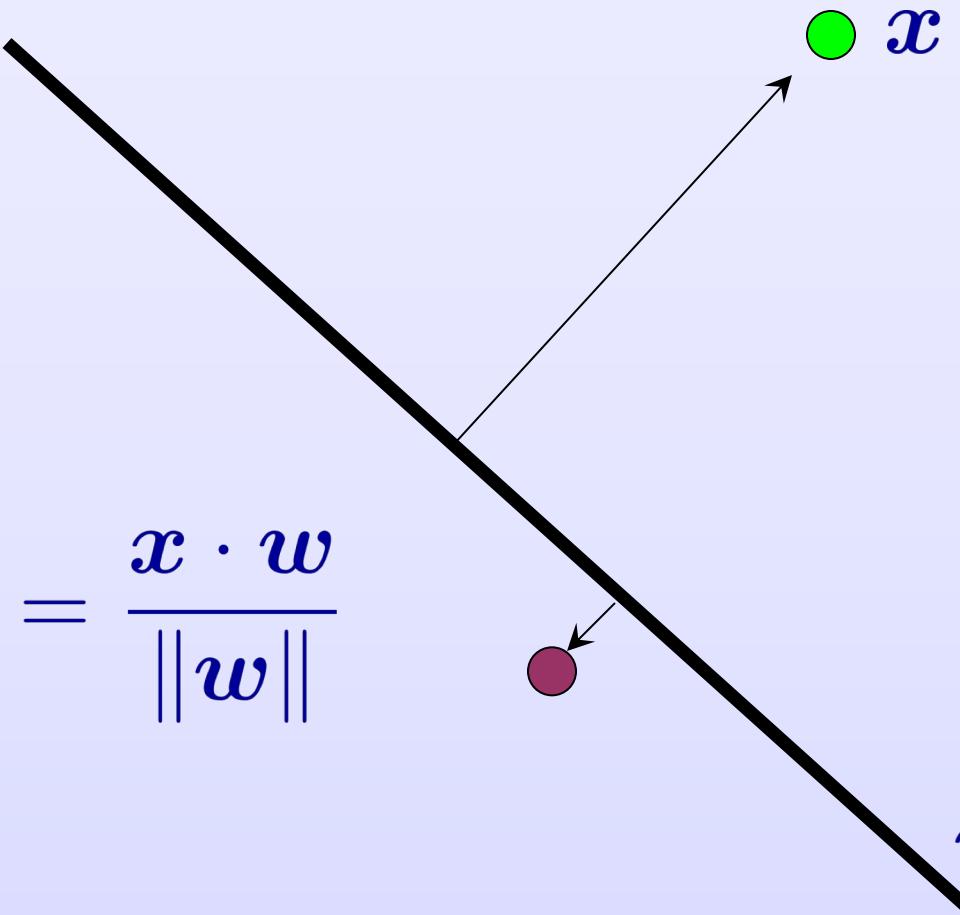
$$= \text{sign}(\mathbf{w} \cdot \mathbf{x})$$

- Confidence in prediction:

$$|f(\mathbf{x})| = |\mathbf{w} \cdot \mathbf{x}|$$

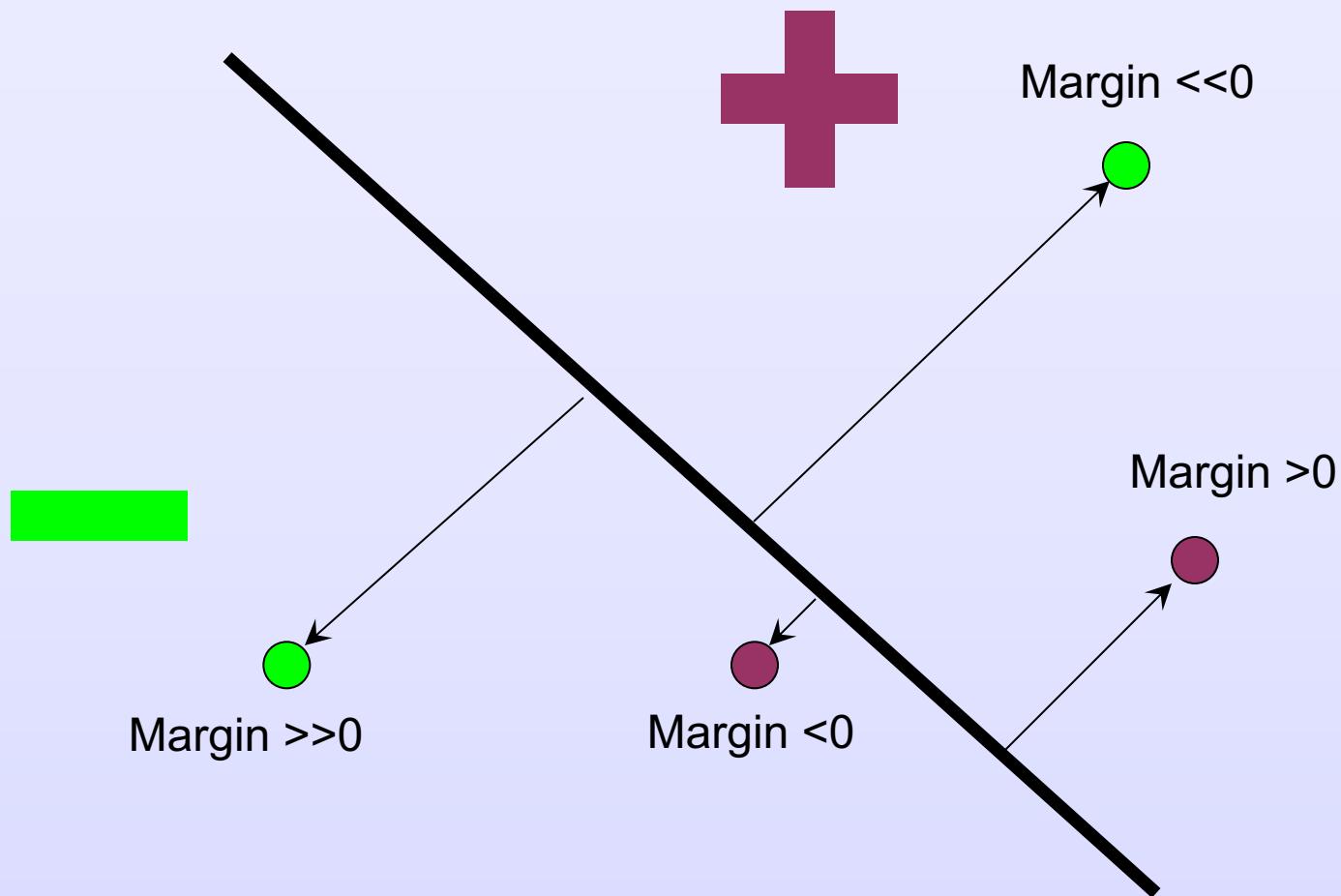


# Signed Distance

$$d(x; w) = \frac{x \cdot w}{\|w\|}$$




# Margin



$$d(yx; \mathbf{w}) = \frac{yx \cdot \mathbf{w}}{\|\mathbf{w}\|}$$

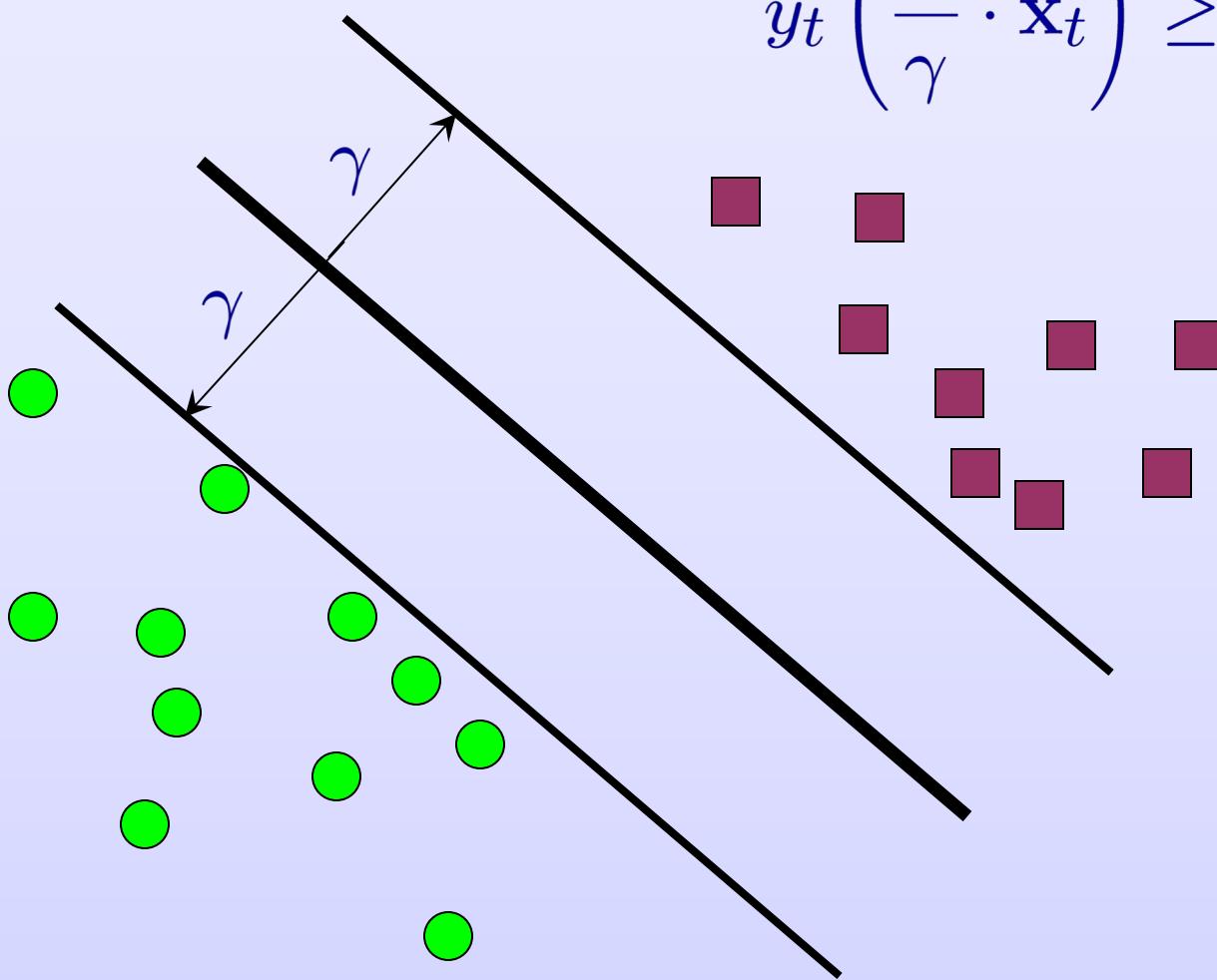


# Separable Set

$\exists \mathbf{w}$  such that  $\forall t$

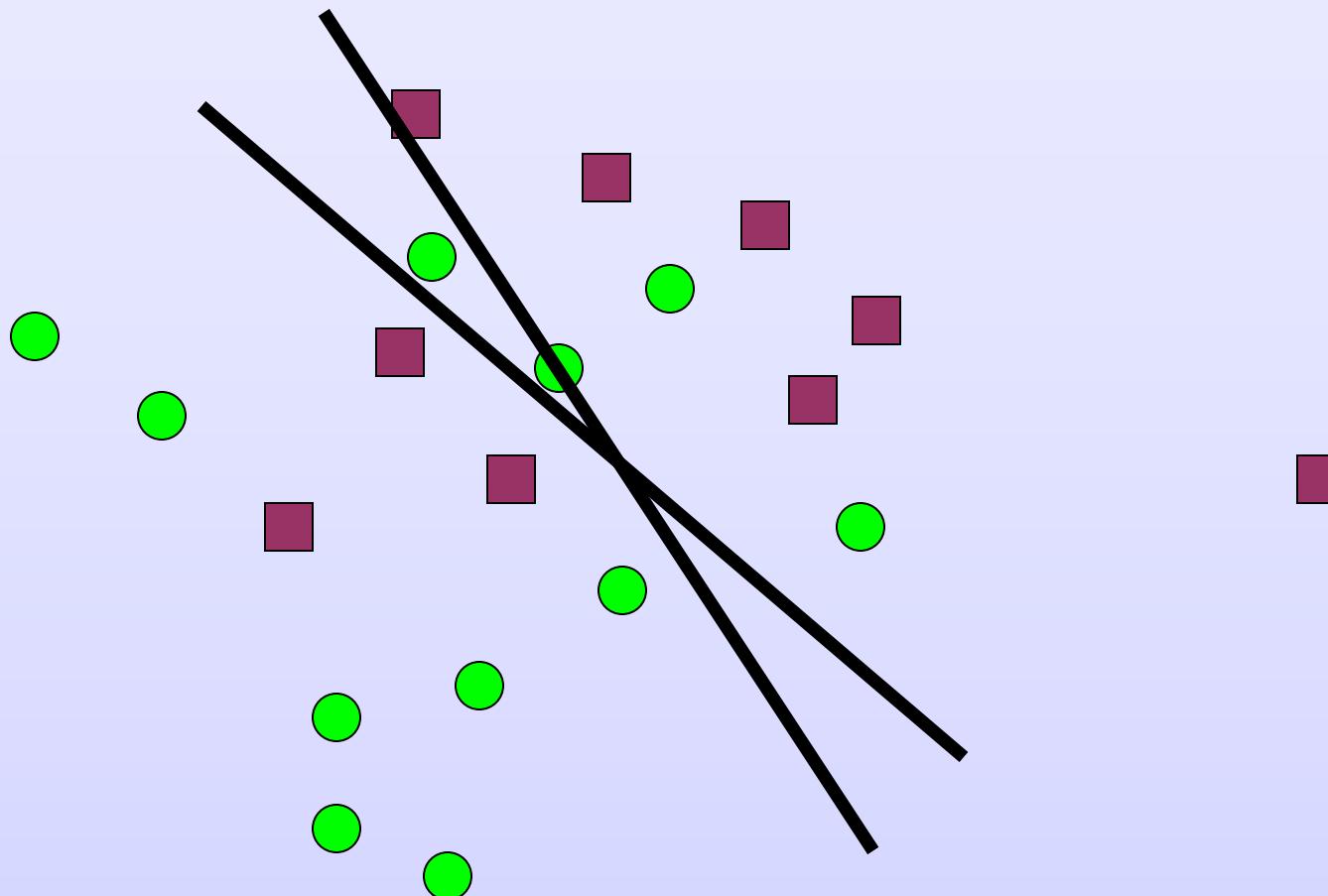
$$y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq \gamma > 0$$

$$y_t \left( \frac{\mathbf{w}}{\gamma} \cdot \mathbf{x}_t \right) \geq 1$$

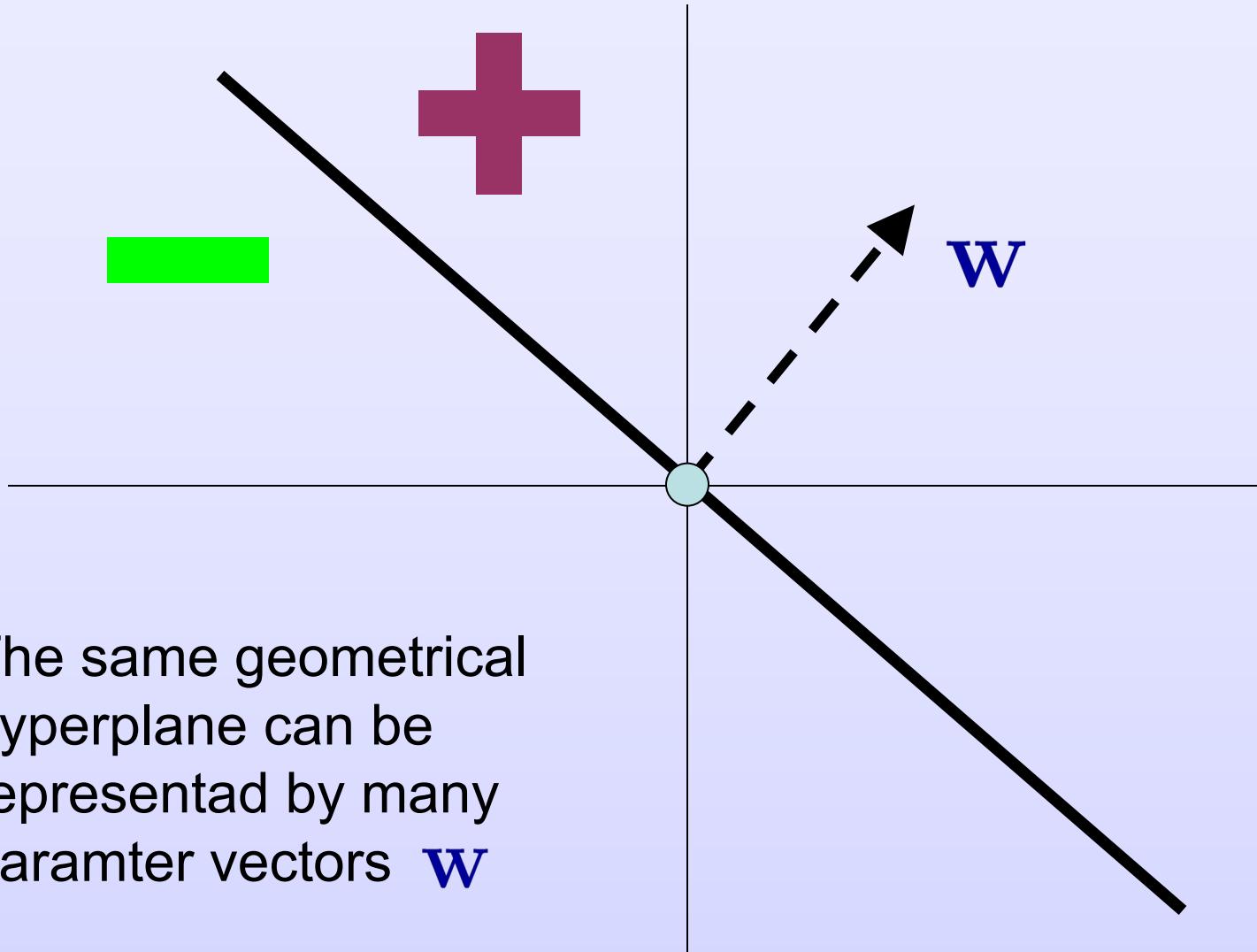


# Inseparable Sets

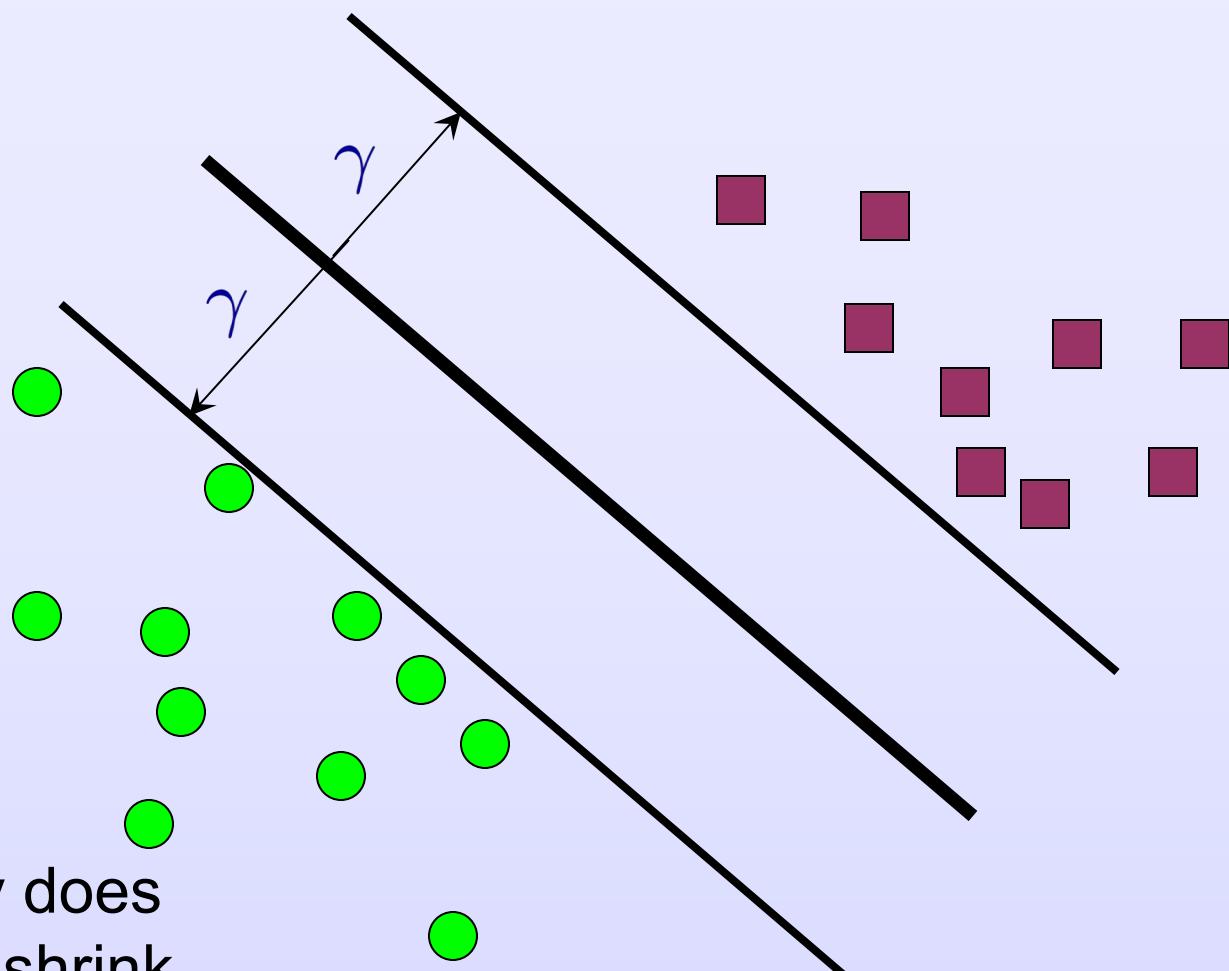
$\forall w$  there exists  $t$  such that  $y_t(w \cdot x_t) < 0$



# Degree of Freedom - I



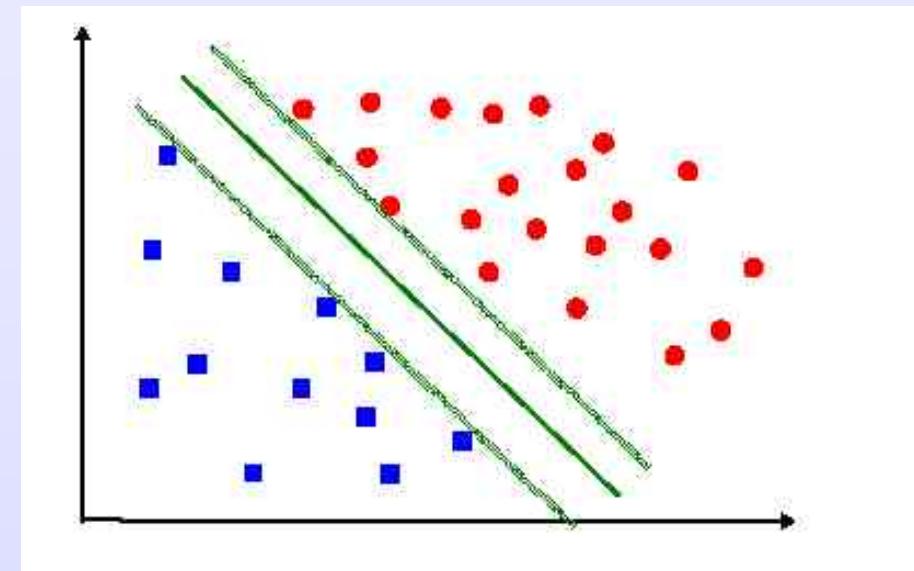
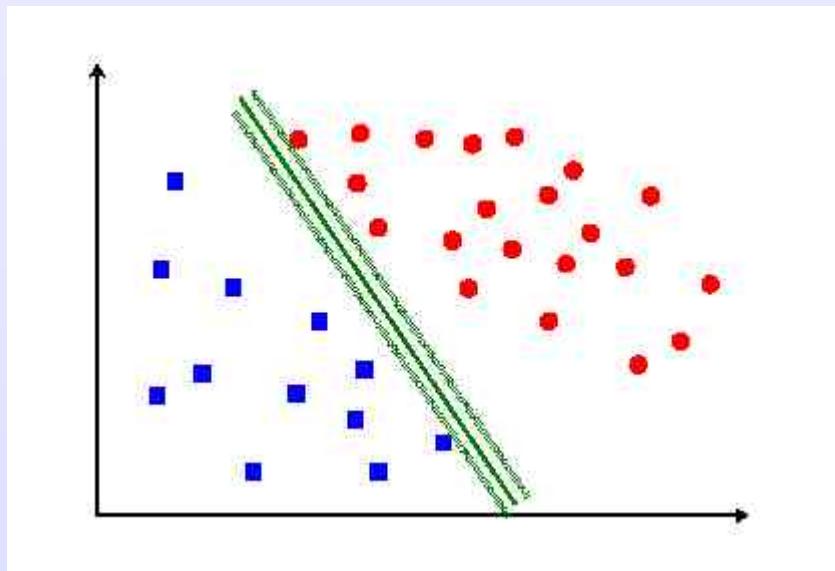
# Degree of Freadom - II



Problem difficulty does  
not change if we shrink  
or expand the input  
space



# Large Margin Classification



# Support Vector Machines

Learning to maximize the margin



# Support Vector Machines

- Assume input data is separable
- Find hyperplane with maximal (sample-)margin
- Equivalent, find hyperplane such that
  - Closest point to the hyperplane ...
  - ... will be as far as possible from it

$$\begin{aligned}\mathbf{w}^* &= \arg \max_{\mathbf{w}} \left( \min_{t=1} d(y_t \mathbf{x}_t; \mathbf{w}) \right) \\ &= \arg \max_{\mathbf{w}} \left( \min_{t=1} \frac{y_t \mathbf{x}_t \cdot \mathbf{w}}{\|\mathbf{w}\|} \right)\end{aligned}$$



# Support Vector Machines

$$\begin{aligned}\mathbf{w}^* &= \arg \max_{\mathbf{w}} \left( \min_{t=1} \frac{y_t \mathbf{x}_t \cdot \mathbf{w}}{\|\mathbf{w}\|} \right) \\ &= \arg \max_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|} \left( \min_{t=1} y_t \mathbf{x}_t \cdot \mathbf{w} \right)\end{aligned}$$

- Use first degree of freedom

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|}$$

$$\text{s.t. } \min_{t=1} y_t \mathbf{x}_t \cdot \mathbf{w} \geq 1$$



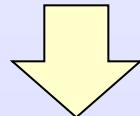
# Support Vector Machines

## Primal Separable SVM

- Algebra

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|}$$

$$\text{s.t. } \min_{t=1} y_t \mathbf{x}_t \cdot \mathbf{w} \geq 1$$



$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } y_t \mathbf{x}_t \cdot \mathbf{w} \geq 1 \quad \forall t$$



# Lagrange Theory

- Problem:

$$\min_x f(x)$$

$$\text{s.t. } g_t(x) \leq 0 \quad \forall t$$

- Min-Max game:

$$\max_{\alpha} \min_x f(x) + \sum_t \alpha_t g_t(x)$$

$$\text{s.t. } \alpha_t \geq 0 \quad \forall t$$

- KKT conditions:

$$\alpha_t^* g_t(x^*) = 0$$



# Lagrangian

$$\max_{\alpha} \min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \sum_t \alpha_t (1 - y_t \mathbf{x}_t \cdot \mathbf{w})$$

$$\text{s.t. } \alpha_t \geq 0 \quad \forall t$$

- Taking the derivative with respect to  $\mathbf{w}$

$$\frac{\partial}{\partial \mathbf{w}} \rightarrow \mathbf{w} - \sum_t \alpha_t y_t \mathbf{x}_t$$

- Optimal solution

$$\mathbf{w} = \sum_t \alpha_t y_t \mathbf{x}_t$$



# Support Vector Machine

- Optimal Solution:

$$\mathbf{w} = \sum_t \alpha_t y_t \mathbf{x}_t$$

- KKT conditions:

$$\alpha_t(1 - y_t \mathbf{w} \cdot \mathbf{x}_t) = 0$$

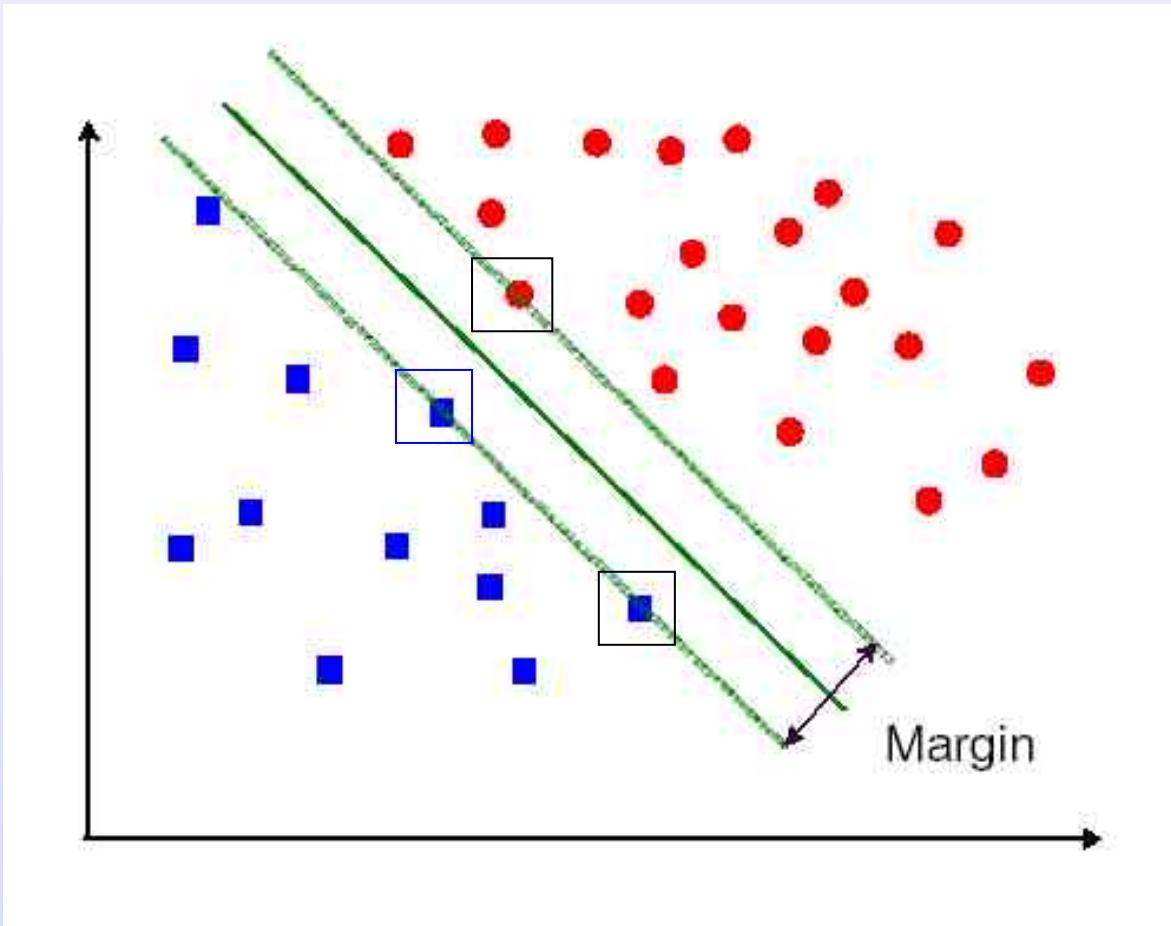
- For a specific example:

$$1 < y_t \mathbf{w} \cdot \mathbf{x}_t \Rightarrow \alpha_t = 0$$

- Only examples that lie on the boundary affect the classifier



# Support Vectors



# Dual Formulation

- Rewriting the Lagrangian

$$\frac{1}{2} \|\mathbf{w}\|^2 + \sum_t \alpha_t (1 - y_t \mathbf{x}_t \cdot \mathbf{w})$$

$$= -\frac{1}{2} \|\mathbf{w}\|^2 + \sum_t \alpha_t - \underbrace{\sum_t \alpha_t y_t \mathbf{x}_t \cdot \mathbf{w}}_{\text{circled term}} + \|\mathbf{w}\|^2$$

- Substituting back in the Lagrangian

$$\mathbf{w} = \sum_t \alpha_t y_t \mathbf{x}_t$$



# Dual Formulation

- Rewriting the Lagrangian

$$\begin{aligned} & \frac{1}{2} \|\boldsymbol{w}\|^2 + \sum_t \alpha_t (1 - y_t \boldsymbol{x}_t \cdot \boldsymbol{w}) \\ = & -\frac{1}{2} \|\boldsymbol{w}\|^2 + \sum_t \alpha_t - \cancel{\boldsymbol{w} \cdot \boldsymbol{w}} + \cancel{\|\boldsymbol{w}\|^2} \end{aligned}$$

- Substituting back in the Lagrangian

$$\boldsymbol{w} = \sum_t \alpha_t y_t \boldsymbol{x}_t$$



# Dual Formulation

- Rewriting the Lagrangian

$$\frac{1}{2} \|\mathbf{w}\|^2 + \sum_t \alpha_t (1 - y_t \mathbf{x}_t \cdot \mathbf{w})$$

$$= -\frac{1}{2} \|\mathbf{w}\|^2 + \sum_t \alpha_t$$

- Substituting back in the Lagrangian

$$\mathbf{w} = \sum_t \alpha_t y_t \mathbf{x}_t$$



# Dual Formulation

- Rewriting the Lagrangian

$$\frac{1}{2} \|\mathbf{w}\|^2 + \sum_t \alpha_t (1 - y_t \mathbf{x}_t \cdot \mathbf{w})$$

$$= -\frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s (\mathbf{x}_t \cdot \mathbf{x}_s) + \sum_t \alpha_t$$

- Substituting back in the Lagrangian

$$\mathbf{w} = \sum_t \alpha_t y_t \mathbf{x}_t$$



# Dual Formulations

From features to kernels



# Dual Formulation

$$\max_{\alpha} - \frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s (\mathbf{x}_t \cdot \mathbf{x}_s) + \sum_t \alpha_t$$

$$\text{s.t. } \alpha_t \geq 0 \quad \forall t$$

- Making Predictions:

$$\mathbf{w} = \sum_t \alpha_t y_t \mathbf{x}_t$$

$$f(\mathbf{x}) = \text{sign} (\mathbf{w} \cdot \mathbf{x})$$

$$= \text{sign} \left( \sum_t \alpha_t y_t (\mathbf{x} \cdot \mathbf{x}_t) \right)$$



# Inner-Product

$$\max_{\alpha} - \frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s (\mathbf{x}_t \cdot \mathbf{x}_s) + \sum_t \alpha_t$$

$$\text{s.t. } \alpha_t \geq 0 \quad \forall t$$

$$f(\mathbf{x}) = \text{sign} \left( \sum_t \alpha_t y_t (\mathbf{x} \cdot \mathbf{x}_t) \right)$$

- Learning and Inferring/predicting can be computed via inner-products between inputs
- No need to know inputs explicitly



# Feature Mapping

- Use feature mapping

$$\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D \quad D \gg d$$

- Learning and inference become:

$$\max_{\alpha} -\frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s (\Phi(\mathbf{x}_t) \cdot \Phi(\mathbf{x}_s)) + \sum_t \alpha_t$$

$$\text{s.t. } \alpha_t \geq 0 \quad \forall t$$

$$f(\mathbf{x}) = \text{sign} \left( \sum_t \alpha_t y_t (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_t)) \right)$$



# From Features to Kernels

- Problem:
  - Computing features can infeasible

$$\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D \quad D \gg d$$

- Wish:
  - Efficiently compute the inner-product directly

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



# From Features to Kernels

- Learning:

$$\max_{\alpha} - \frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s K(\mathbf{x}_t, \mathbf{x}_s) + \sum_t \alpha_t$$

s.t.  $\alpha_t \geq 0 \quad \forall t$

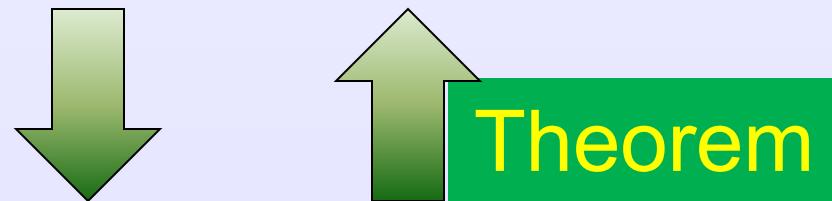
- Inference:

$$f(\mathbf{x}) = \text{sign} \left( \sum_t \alpha_t y_t K(\mathbf{x}, \mathbf{x}_t) \right)$$



# Kernels

- Let  $K(\mathbf{x}, \mathbf{z}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}))$



- Properties:
  - Symmetric  $K(\mathbf{x}, \mathbf{z}) = K(\mathbf{z}, \mathbf{x})$
  - Given a finite set of inputs  $\mathbf{x}_1 \dots \mathbf{x}_n$   
define a matrix  $\kappa_{s,t} = K(\mathbf{x}_t, \mathbf{x}_s)$   
then the matrix  $\kappa$  is positive semi-definite



# Example

- Consider 2-d inputs, define:

$$\Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

- Then:

$$\begin{aligned}\Phi(x_1, x_2) \cdot \Phi(z_1, z_2) \\ &= (x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (z_1^2, z_2^2, \sqrt{2}z_1z_2) \\ &= x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (x_1z_1 + x_2z_2)^2 \\ &= (\mathbf{x} \cdot \mathbf{z})^2\end{aligned}$$



# Example

- Consider 2-d inputs, define:

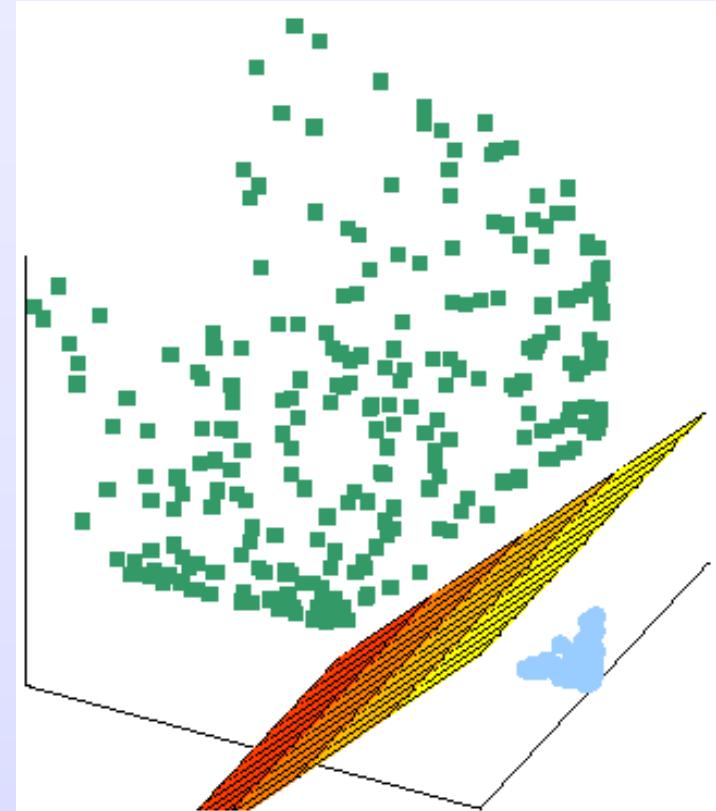
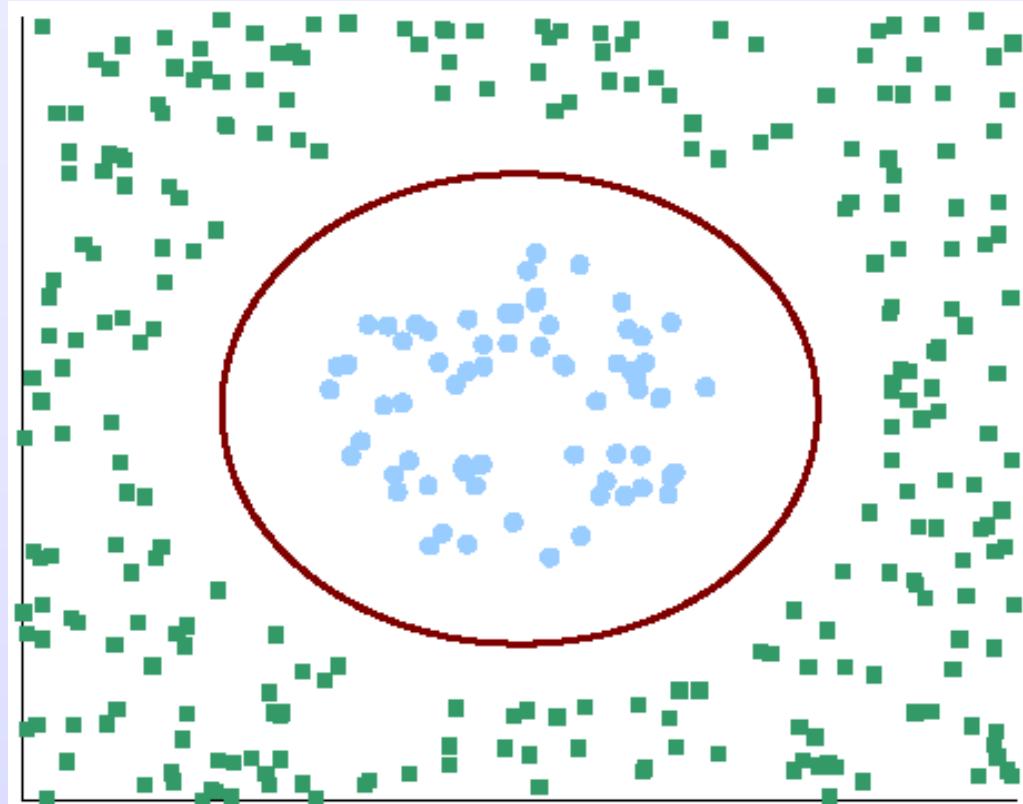
$$\Phi(x_1, x_2) = \left( x_1^2, x_2^2, \sqrt{2}x_1x_2 \right)$$

- Then:

$$\begin{aligned}\Phi(x_1, x_2) \cdot \Phi(z_1, z_2) &= (\mathbf{x} \cdot \mathbf{z})^2 \\ &= K(\mathbf{x}, \mathbf{z})\end{aligned}$$



# Kernels



# Kernels

- Common:

- Polynomial

$$K(\mathbf{x}, \mathbf{z}) = (a + \mathbf{x} \cdot \mathbf{z})^b$$

- RBF

$$K(\mathbf{x}, \mathbf{z}) = \exp -a \|\mathbf{x} - \mathbf{z}\|^2$$

- Crafted:

- e.g. String kernels, tree kernels

- Learned:

- using SDP, boosting ...

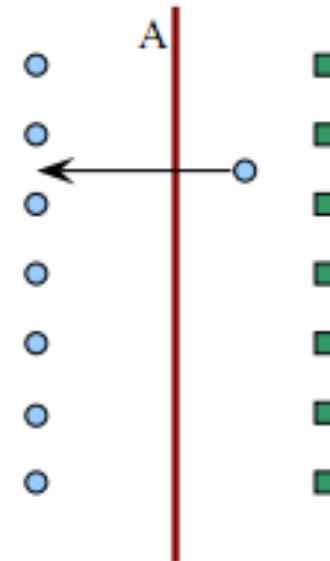
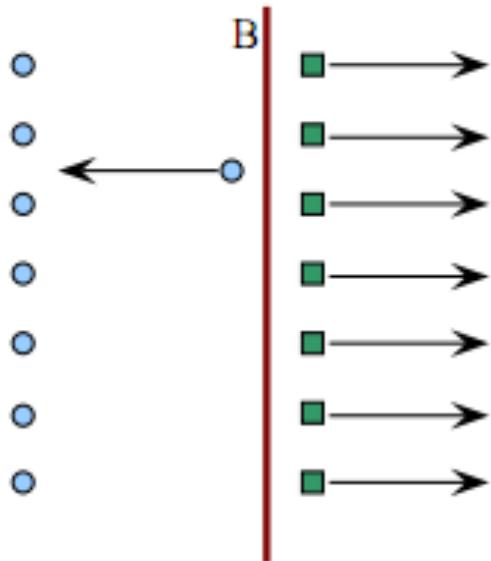


# Non-Separable Formulation

Hinge Loss + Regularization



# What is better?

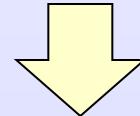


# Support Vector Machines

- Add s Primal Non-Separable SVM

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } y_t \mathbf{x}_t \cdot \mathbf{w} \geq 1 \quad \forall t$$



$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi_t$$

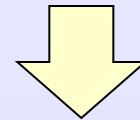
$$\text{s.t. } y_t \mathbf{x}_t \cdot \mathbf{w} \geq 1 - \xi_t, \quad \xi_t \geq 0 \quad \forall t$$



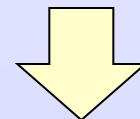
# Hinge Loss

- Rewrite constraints

$$y_t \mathbf{x}_t \cdot \mathbf{w} \geq 1 - \xi_t , \quad \xi_t \geq 0 \quad \forall t$$



$$\xi_t \geq 1 - y_t \mathbf{x}_t \cdot \mathbf{w} , \quad \xi_t \geq 0 \quad \forall t$$



$$\xi_t \geq \max\{0, 1 - y_t \mathbf{x}_t \cdot \mathbf{w}\} \quad \forall t$$



# Hinge Loss

Regularization

Tradeoff  
Parameter

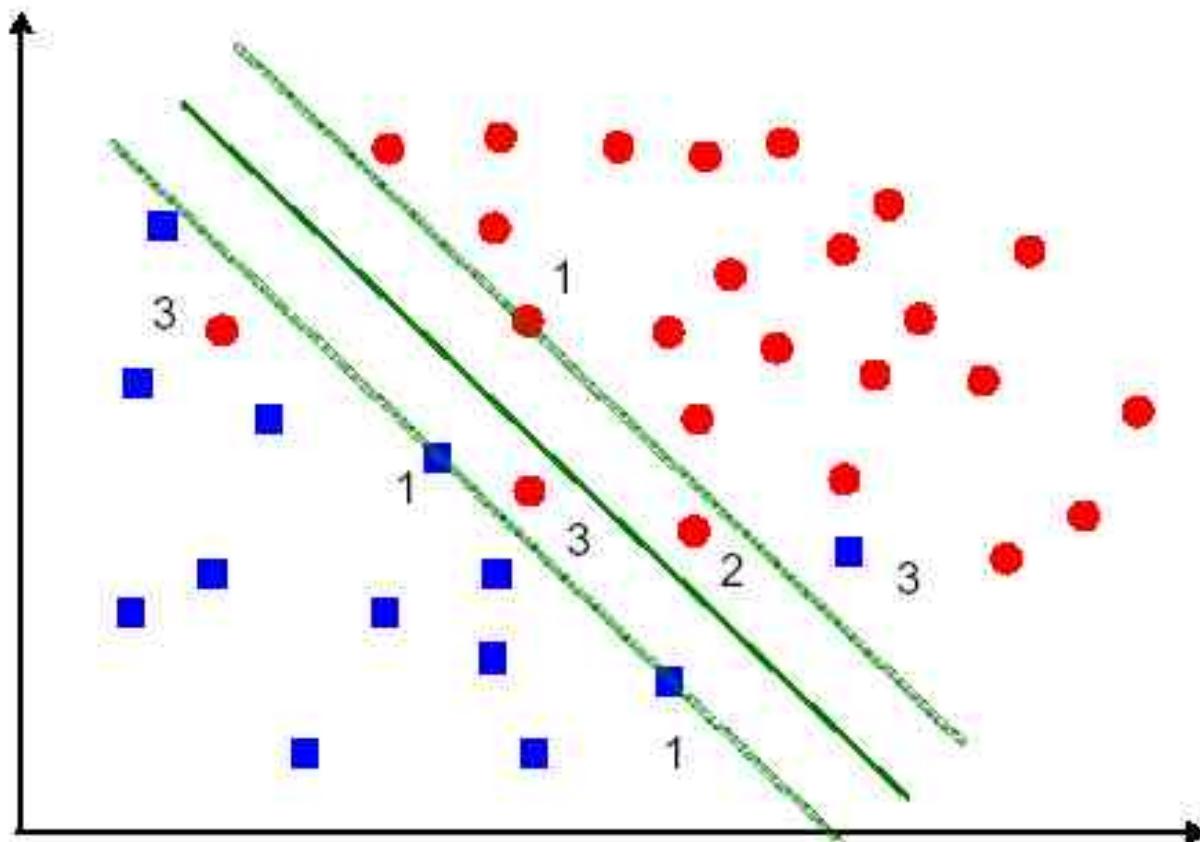
Hinge Loss

$$\mathbf{w}^* = \arg \min_{\mathbf{w}}$$

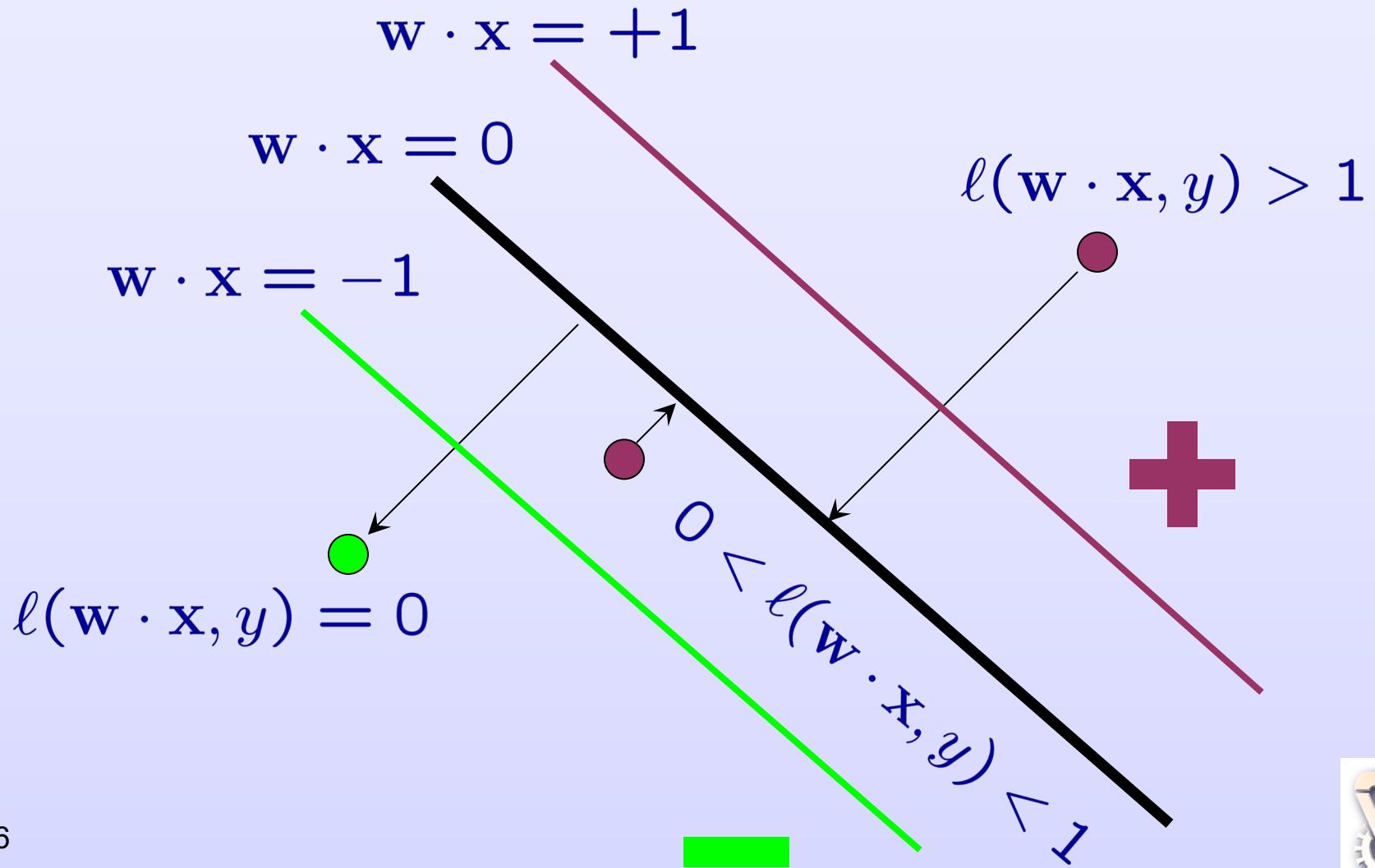
$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \max\{0, 1 - y_t \mathbf{x}_t \cdot \mathbf{w}\}$$



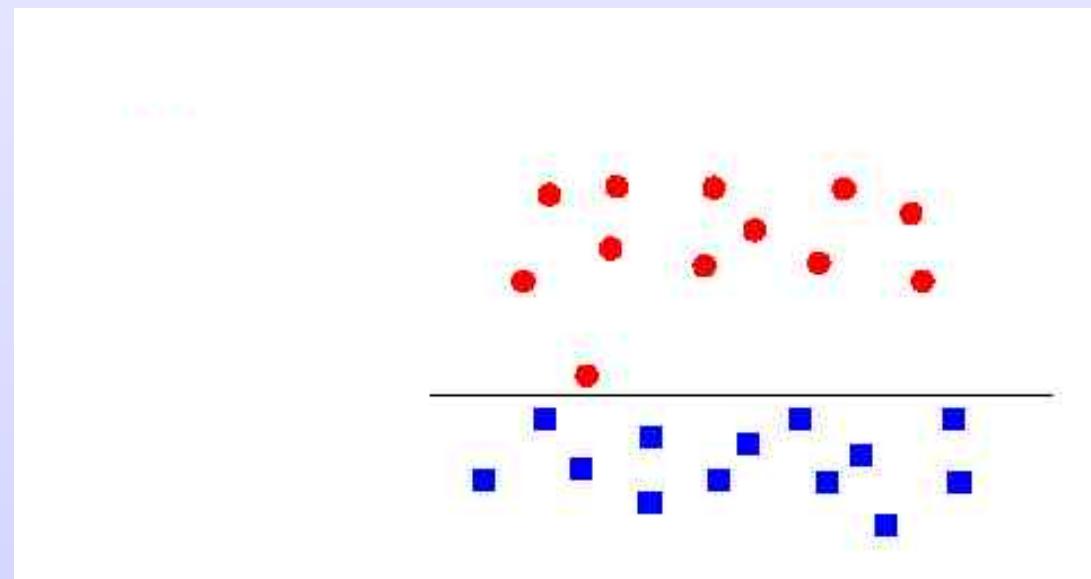
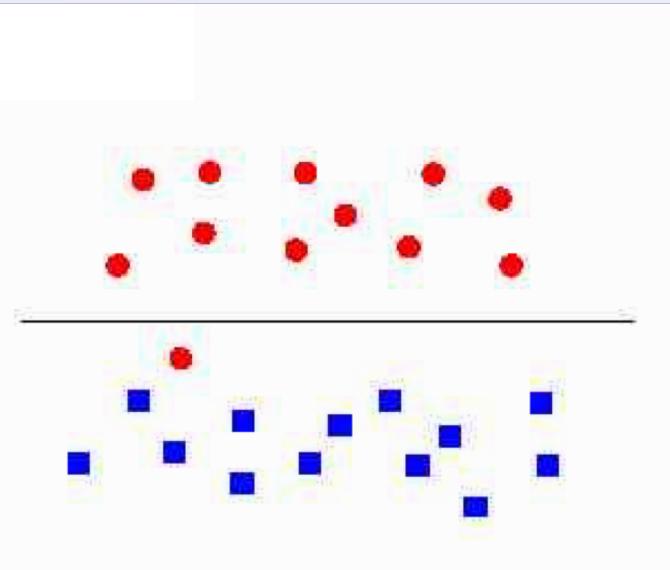
# Support Vectors



# Hinge Loss



# Non-Separable Formulation



# Dual Non-Separable

- Learning:

$$\max_{\alpha} - \frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s K(\mathbf{x}_t, \mathbf{x}_s) + \sum_t \alpha_t$$

s.t.  $0 \leq \alpha_t \leq C \quad \forall t$

- Inference:

$$f(\mathbf{x}) = \text{sign} \left( \sum_t \alpha_t y_t K(\mathbf{x}, \mathbf{x}_t) \right)$$



# Comments

- SVM: quadratic objective + linear constraints
- Principle of large margin:
  - Not only be correct, be also confident
- Kernels
- Bias term – class prior
- Many many algorithms to solve optimization
- Variants for:
  - regression, ranking, multi-class, multi-label, structured prediction ...



# The Common and Different

- Naïve Bayes, Logistic Regression and SVMs employ linear models
- Generative vs. Discriminative
- Probabilistic output?
- Minimization of Regularization + Loss
  - Log loss, Hinge loss, Exp loss, Huber loss ...
  - Squared L2, L1, L-infinity, ...
- Many algorithms can be combined with kernels, including logistic-regression
- Generalization theory

