

CSE 575: Statistical Machine Learning

Jingrui He
CIDSE, ASU

Instance-based Learning

1-Nearest Neighbor

Four things make a memory based learner:

1. A distance metric

Euclidian (and many more)

2. How many nearby neighbors to look at?

One

1. A weighting function (optional)

Unused

2. How to fit with the local points?

Just predict the same output as the nearest neighbor.

Consistency of 1-NN

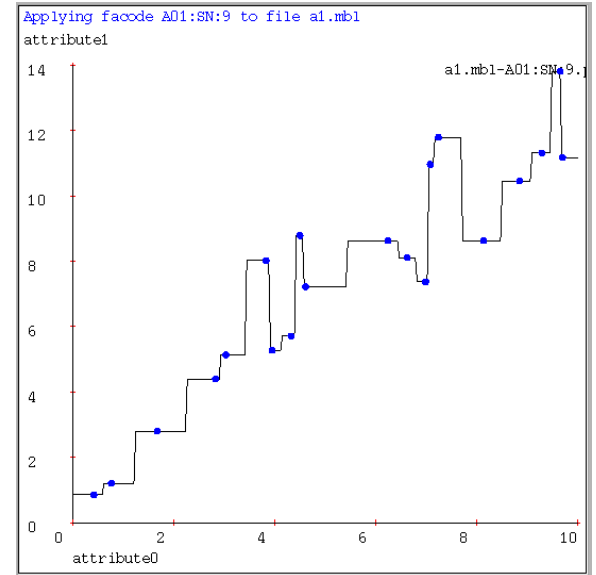
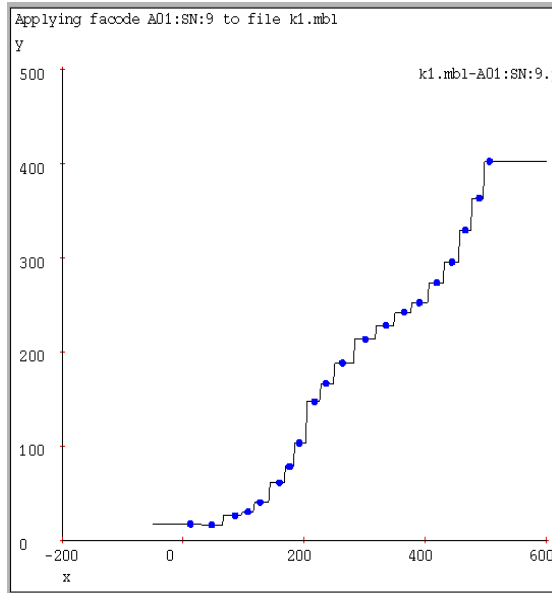
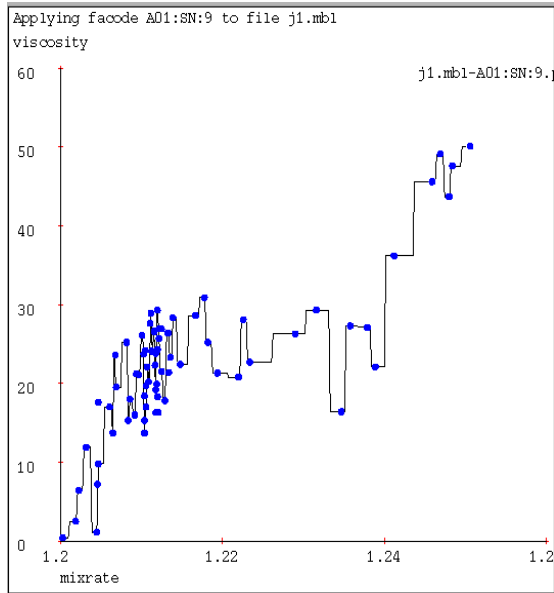
- Consider an estimator f_n trained on n examples
 - e.g., 1-NN, regression, ...
- Estimator is *consistent* if true error goes to zero as amount of data increases
 - e.g., for no noise data, consistent if:

$$\lim_{n \rightarrow \infty} MSE(f_n) = 0$$

- Regression is not consistent!
 - Representation bias
- **1-NN is consistent** (under some mild fineprint)

What about variance??? 

1-NN overfits?



k-Nearest Neighbor

Four things make a memory based learner:

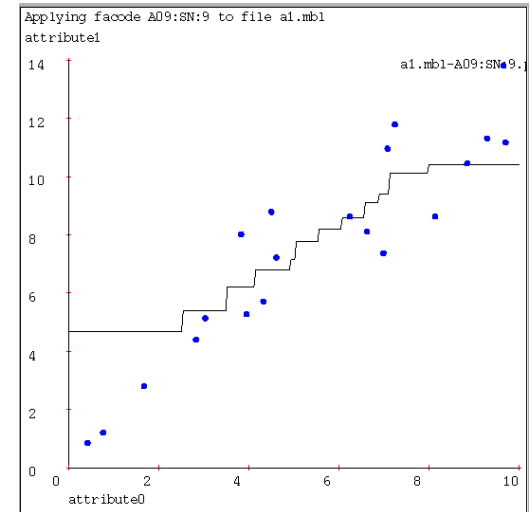
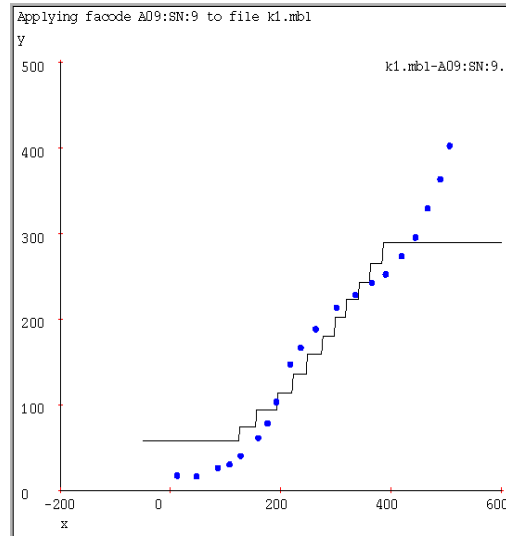
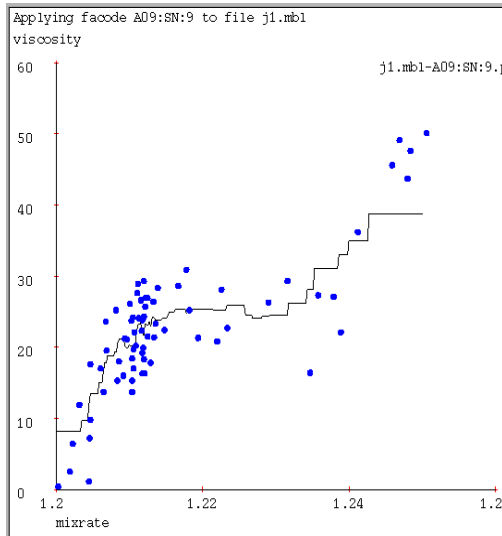
1. *A distance metric*
Euclidian (and many more)
2. *How many nearby neighbors to look at?*
k

1. *A weighting function (optional)*
Unused

2. *How to fit with the local points?*

Just predict the average output among the k nearest neighbors.

k-Nearest Neighbor (here k=9)



K-nearest neighbor for function fitting smooth away noise, but there are clear deficiencies.

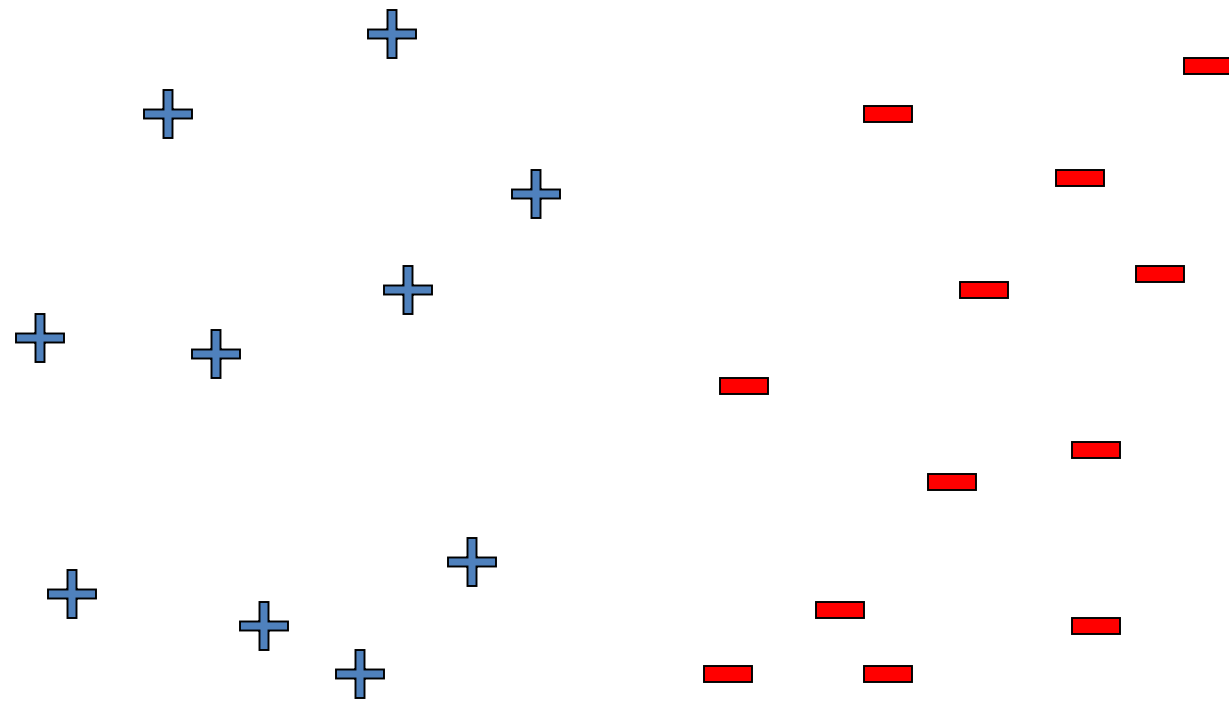
What can we do about all the discontinuities that k-NN gives us?

Curse of dimensionality for instance-based learning

- Must store and retrieve all data!
 - Most real work done during testing
 - For every test sample, must search through all dataset – very slow!
 - There are fast methods for dealing with large datasets, e.g., tree-based methods, hashing methods,
- Instance-based learning often poor with noisy or irrelevant features

Support Vector Machines

Linear classifiers – Which line is better?



Data:

$$\begin{aligned} &\langle x_1^{(1)}, \dots, x_1^{(m)}, y_1 \rangle \\ &\quad \vdots \\ &\langle x_n^{(1)}, \dots, x_n^{(m)}, y_n \rangle \end{aligned}$$

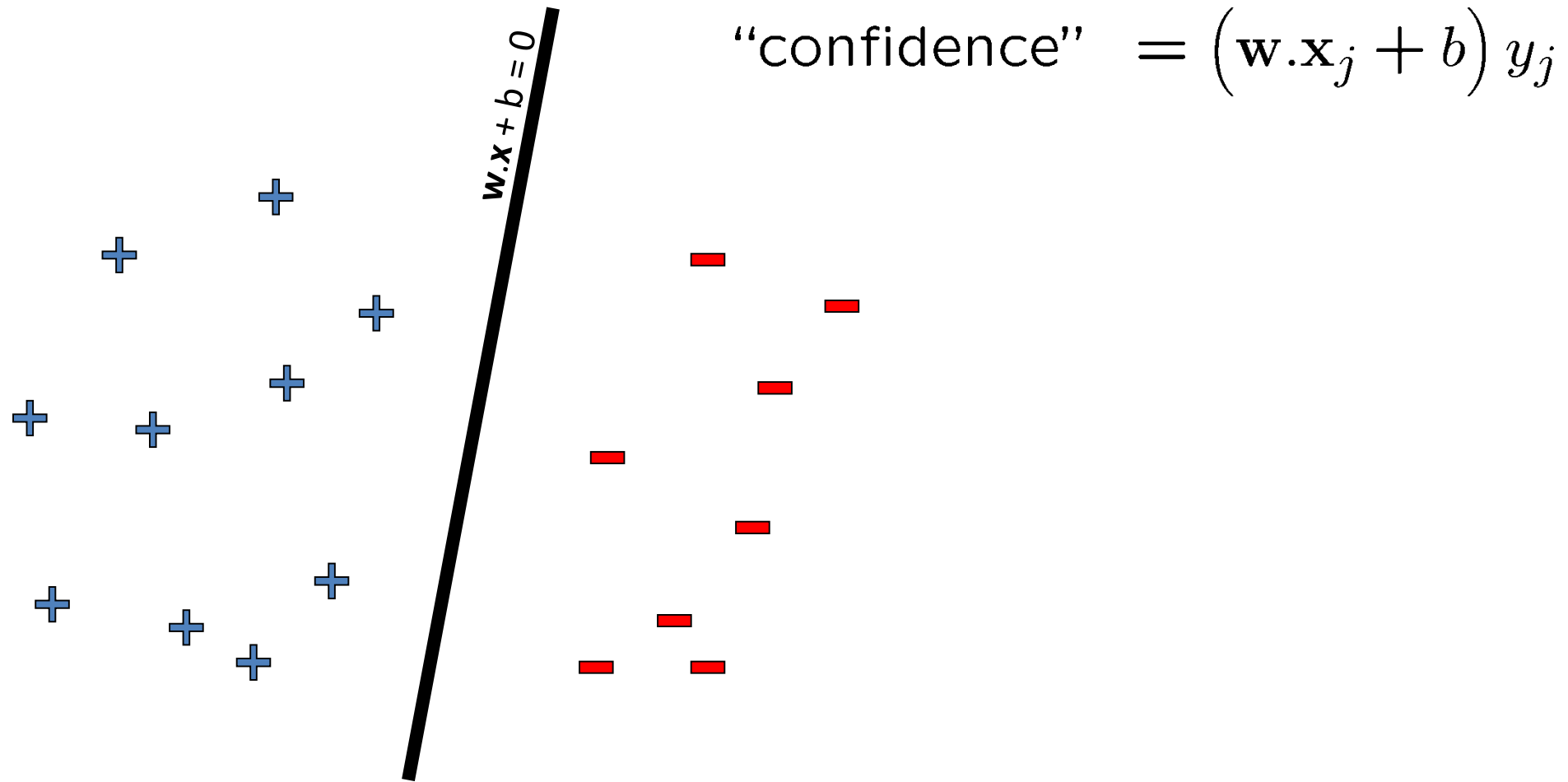
Example i:

$$\langle x_i^{(1)}, \dots, x_i^{(m)} \rangle \quad \text{— } m \text{ features}$$

$$y_i \in \{-1, +1\} \quad \text{— class}$$

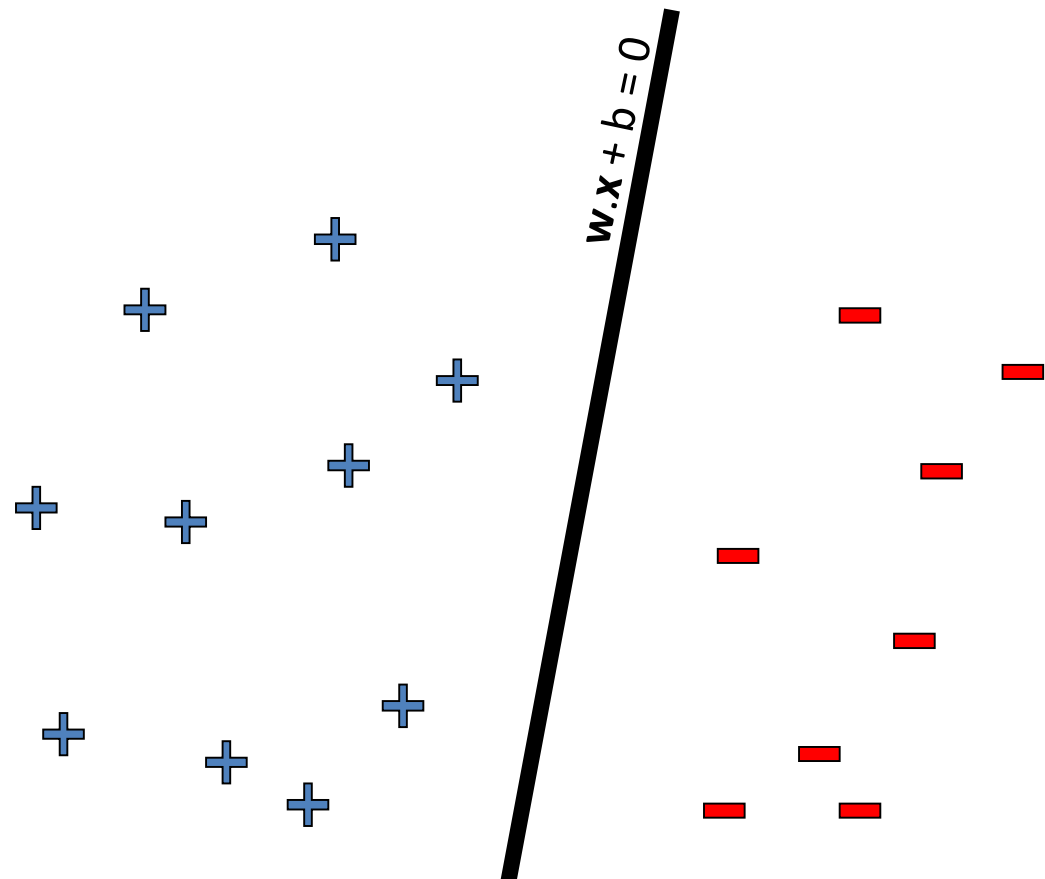
$$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$$

Pick the one with the largest margin!



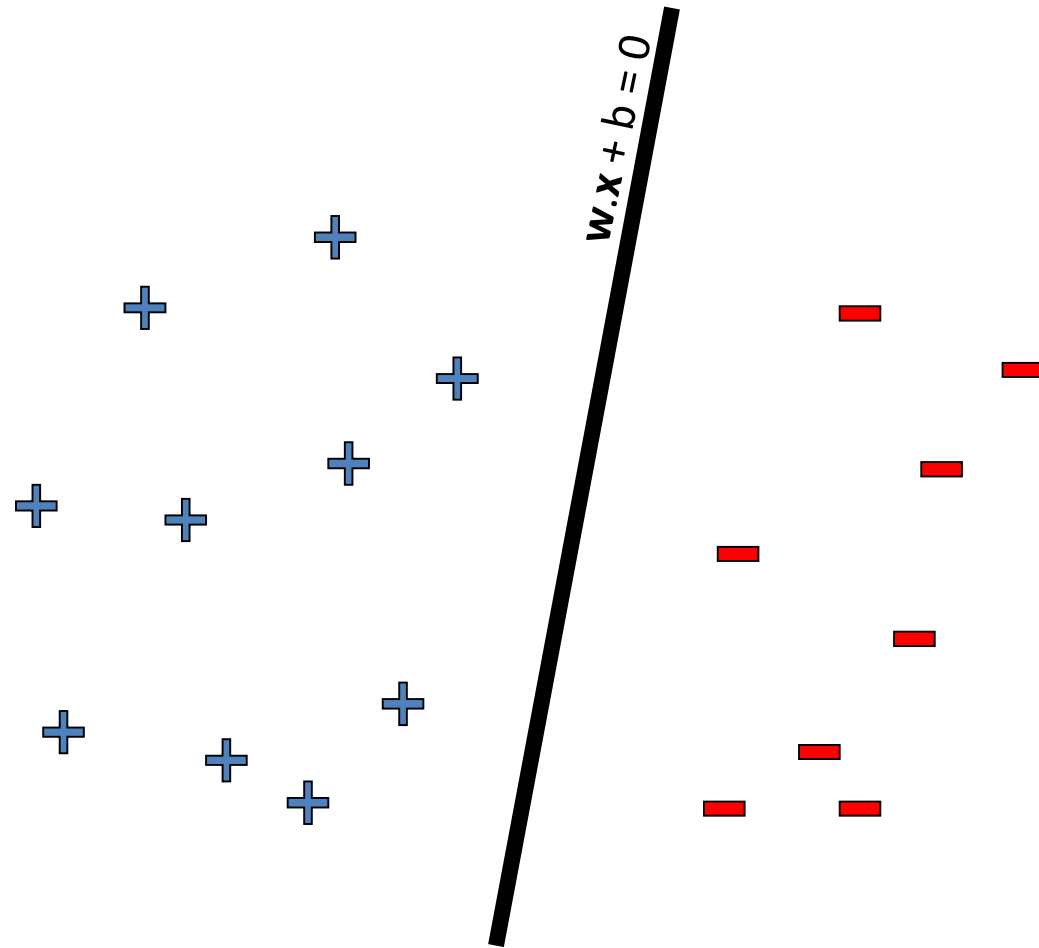
$$w \cdot x = \sum_j w^{(j)} x^{(j)}$$

Maximize the margin



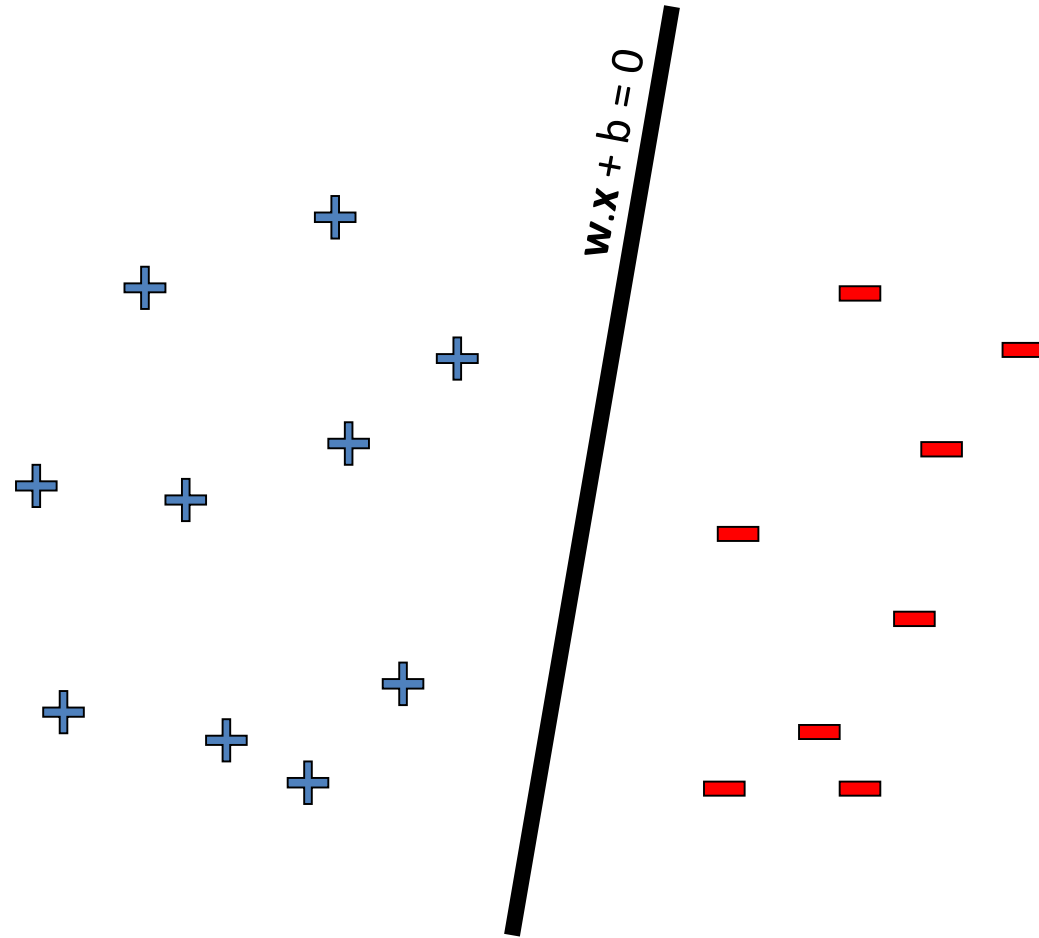
$$\text{maximize}_{\gamma, w, b} \quad \gamma$$
$$(w \cdot x_j + b) y_j \geq \gamma, \quad \forall j \in \text{Dataset}$$

But there are a many planes...



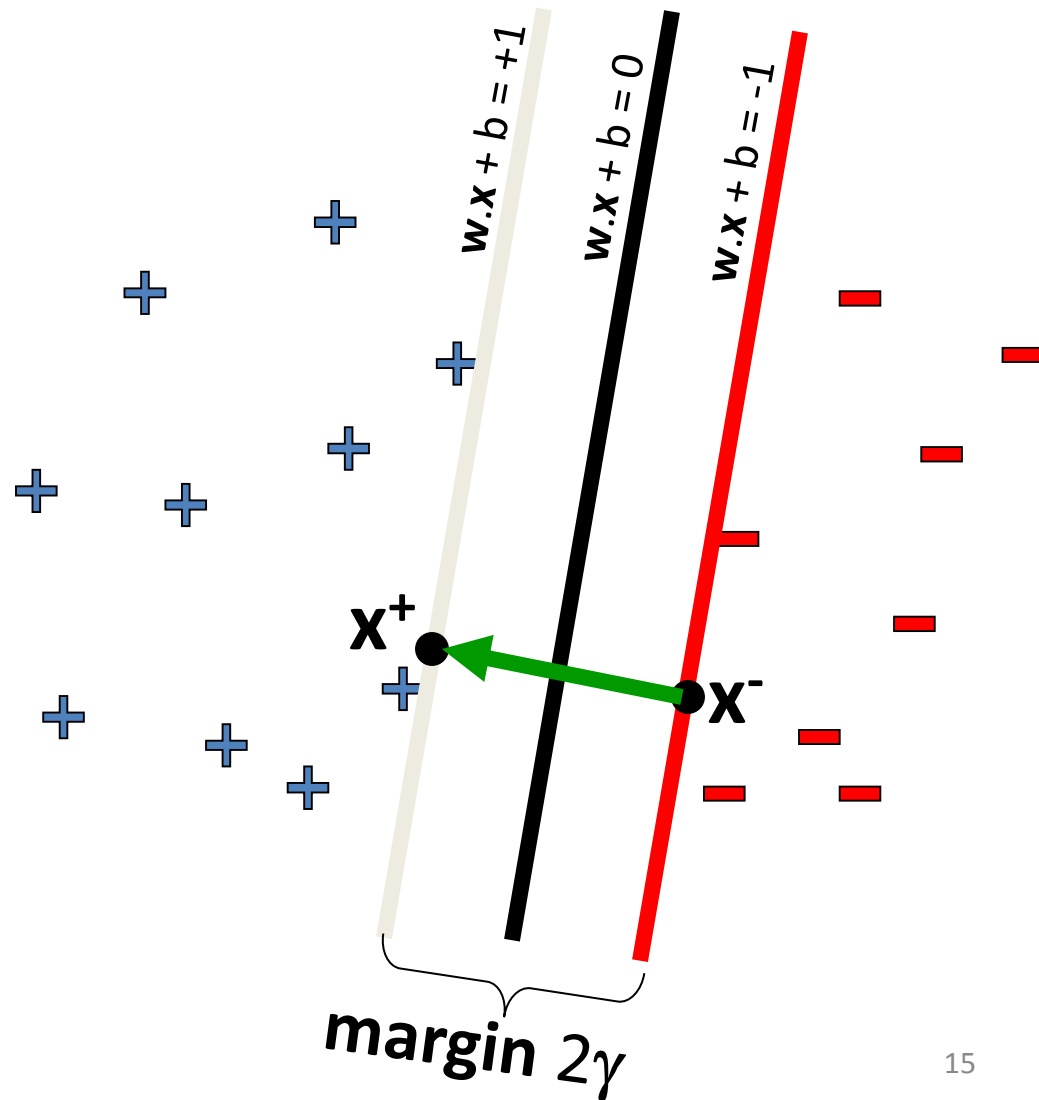
Review: Normal to a plane

$$\mathbf{x}_j = \bar{\mathbf{x}}_j + \lambda \frac{\mathbf{w}}{||\mathbf{w}||}$$



Normalized margin – Canonical hyperplanes

$$\mathbf{x}_j = \bar{\mathbf{x}}_j + \lambda \frac{\mathbf{w}}{\|\mathbf{w}\|}$$



Normalized margin – Canonical hyperplanes

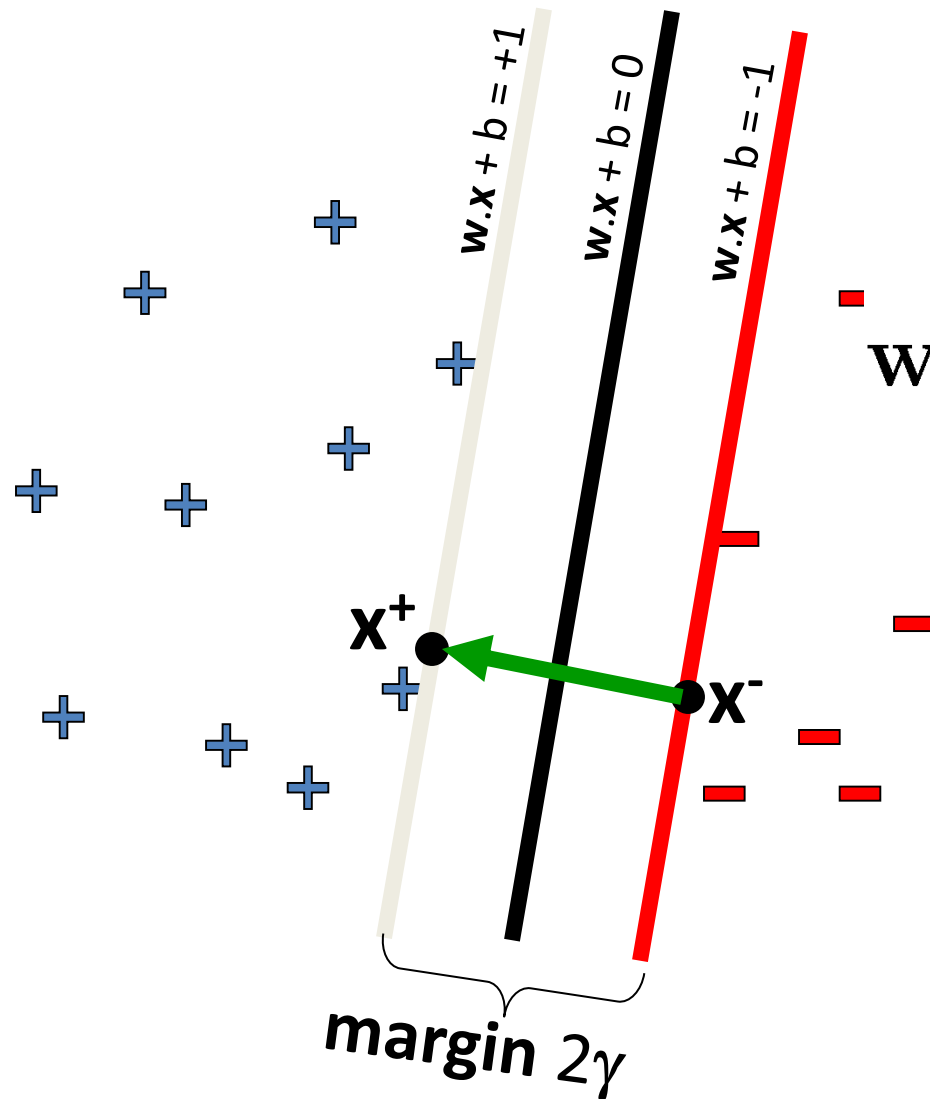
$$\mathbf{x}^+ = \mathbf{x}^- + \lambda \frac{\mathbf{w}}{||\mathbf{w}||}$$

$$\mathbf{w} \cdot \mathbf{x}^+ + b = 1$$

$$\mathbf{w} \cdot \left(\mathbf{x}^- + \lambda \frac{\mathbf{w}}{||\mathbf{w}||} \right) + b = 1$$

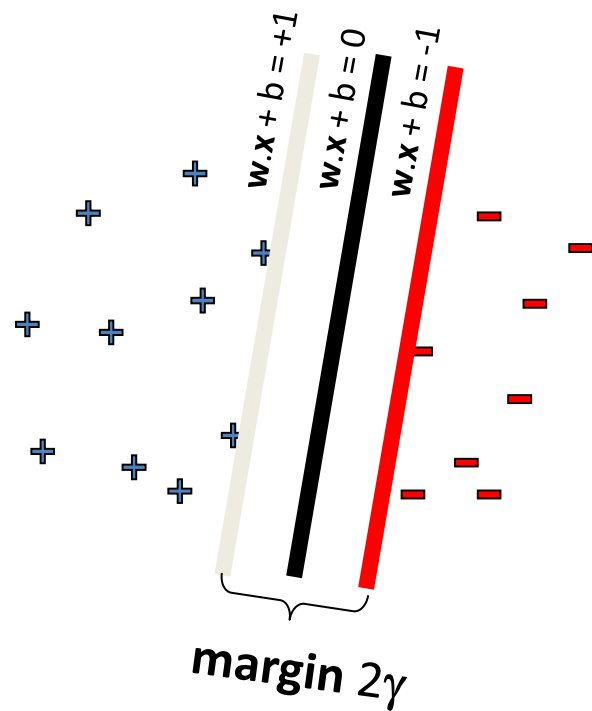
$$\lambda = \frac{2}{||\mathbf{w}||}$$

$$\gamma = \frac{1}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$



Margin maximization using canonical hyperplanes

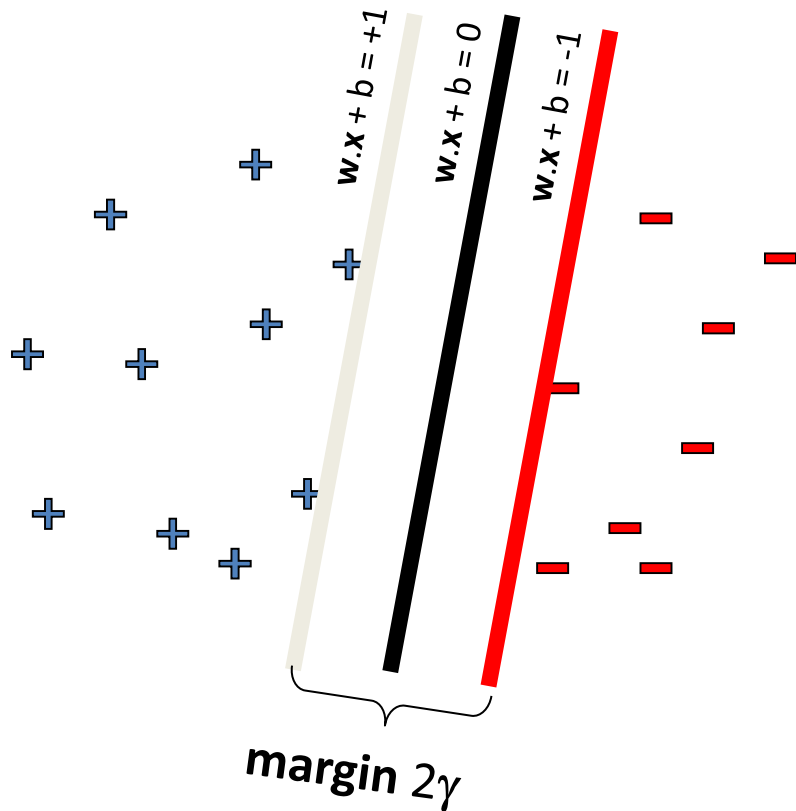
$$\gamma = \frac{1}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$



$$\begin{aligned} & \text{maximize}_{\gamma, \mathbf{w}, b} \quad \gamma \\ & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq \gamma, \quad \forall j \in \text{Dataset} \end{aligned}$$

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

Support Vector Machines (SVMs)



$$\text{minimize}_{w,b} \quad w \cdot w$$
$$\left(w \cdot x_j + b \right) y_j \geq 1, \quad \forall j$$

- Solve efficiently by quadratic programming (QP)
 - Well-studied algorithms
- Hyperplane defined by support vectors

What if the data is not linearly separable?

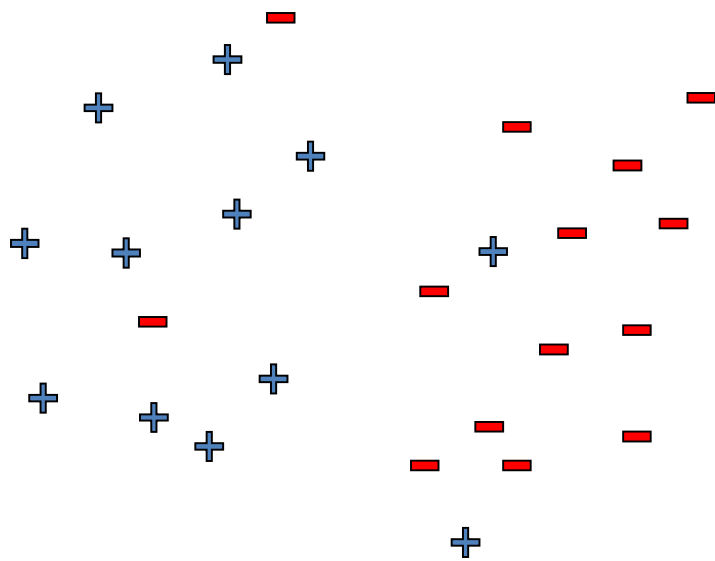


Use features of features of features of features....

(x, x^2)



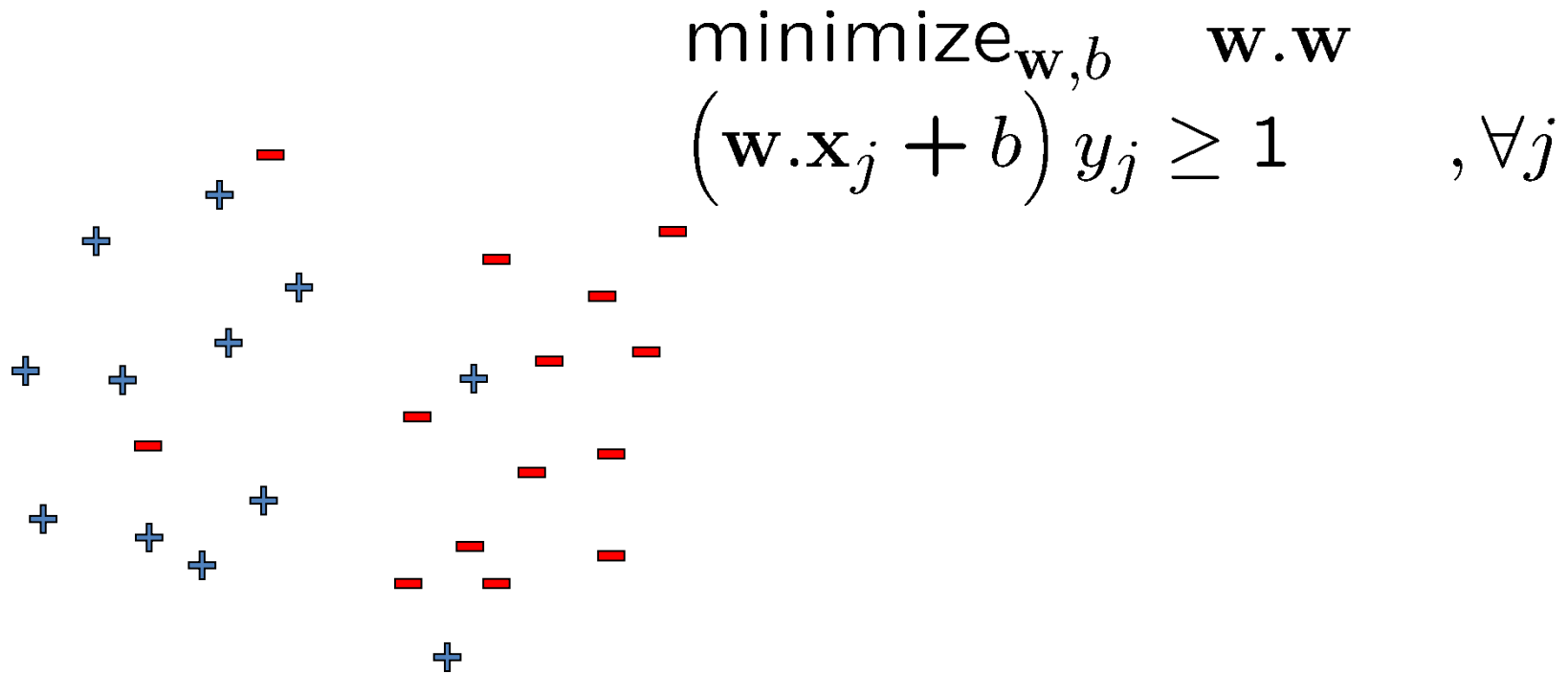
What if the data is still not linearly separable?



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

- Minimize $\mathbf{w} \cdot \mathbf{w}$ and number of training mistakes
 - Tradeoff two criteria?
- Tradeoff #(mistakes) and $\mathbf{w} \cdot \mathbf{w}$
 - 0/1 loss
 - Slack penalty C
 - Not QP anymore
 - Also doesn't distinguish near misses and really bad mistakes

Slack variables – Hinge loss



- If margin ≥ 1 , don't care
- If margin < 1 , pay linear penalty

Side note: What's the difference between SVMs and logistic regression?

SVM:

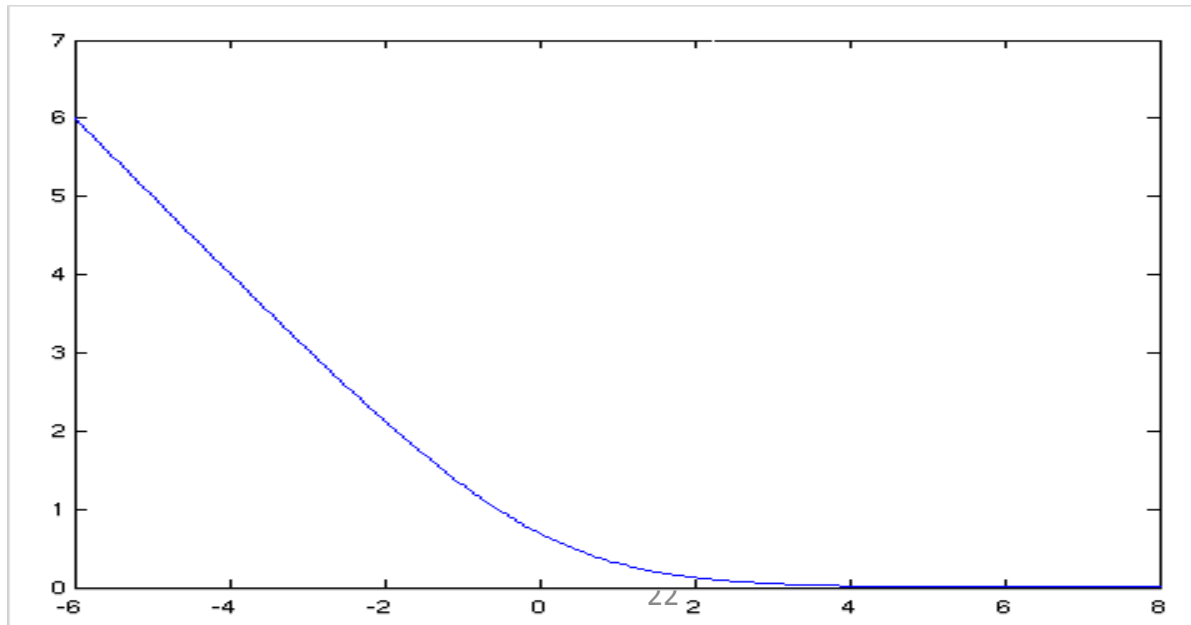
$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j \\ & \xi_j \geq 0, \quad \forall j \end{aligned}$$

Logistic regression:

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

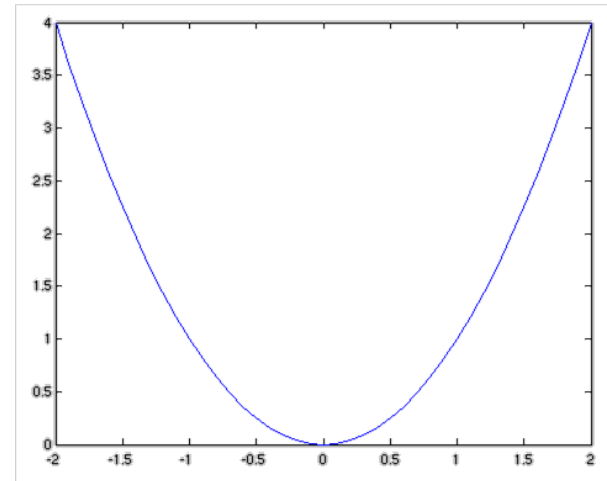
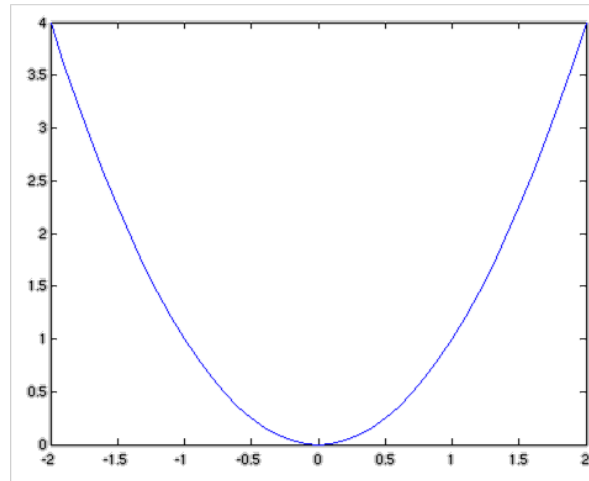
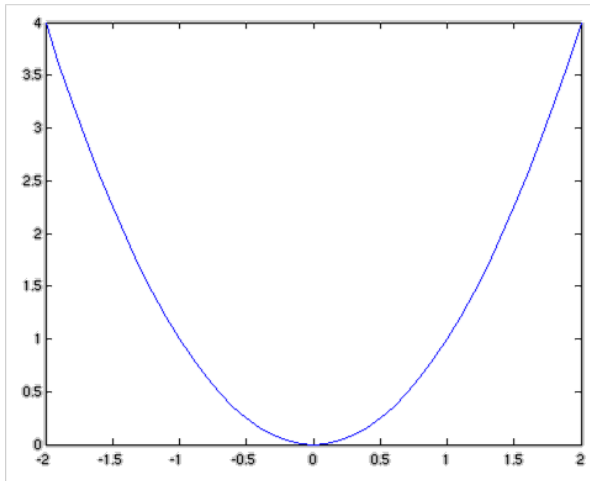
Log loss:

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



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:

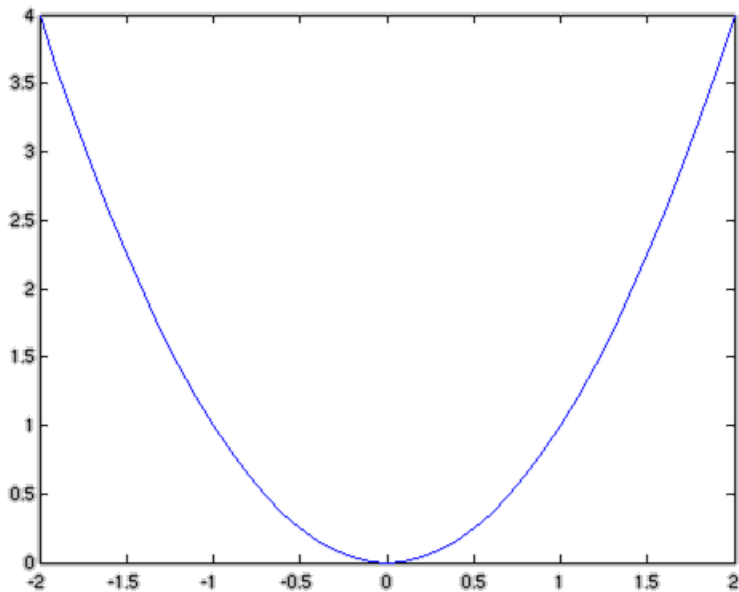
$$L(x, \alpha) = x^2 - \alpha(x - b)$$

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

Solve:

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

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



Dual SVM derivation (1) – the linearly separable case

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

Dual SVM derivation (2) – the linearly separable case

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

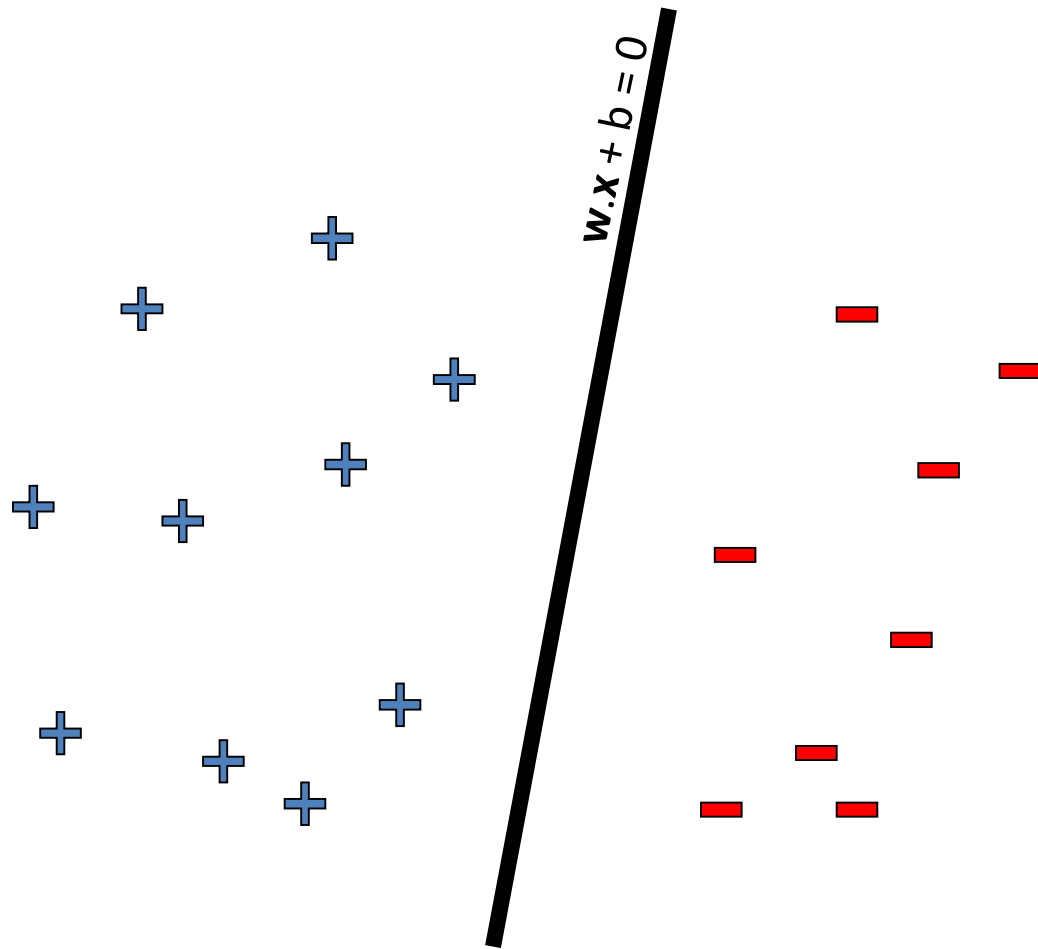
$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

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

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

for any k where $\alpha_k > 0$

Dual SVM interpretation



$$w = \sum_j \alpha_j y_j \mathbf{x}_j$$

Dual SVM formulation – the linearly separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

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

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

for any k where $\alpha_k > 0$

Dual SVM derivation – the non-separable case

$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ \left(\mathbf{w} \cdot \mathbf{x}_j + b \right) y_j & \geq 1 - \xi_j, \quad \forall j \\ \xi_j & \geq 0, \quad \forall j \end{aligned}$$

Dual SVM formulation – the non-separable case

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

$$\sum_i \alpha_i y_i = 0$$

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

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

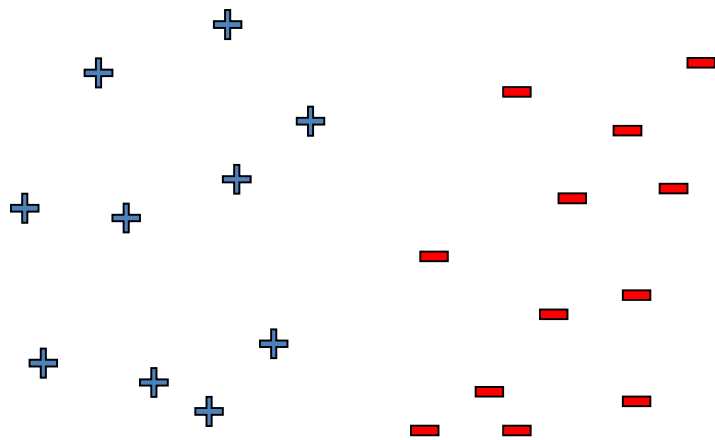
$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

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....

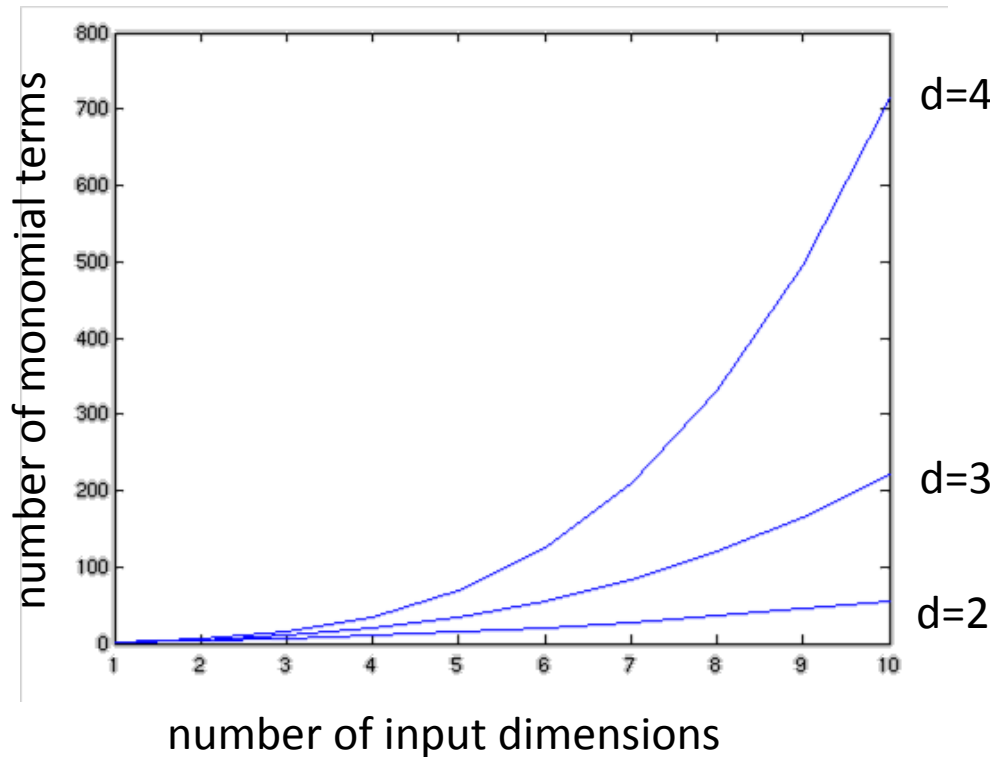
$$\Phi(\mathbf{x}) : R^m \mapsto F$$

Feature space can get really³large 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 **w**!

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \\ & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ & K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \\ & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$

Dot-product of polynomials

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

Finally: the “kernel trick”!

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

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

$$\sum_i \alpha_i y_i = 0$$

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

- 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

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

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

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 α_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?

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss
High dimensional features with kernels	Yes!	Yes!
Solution sparse	Often yes!	Almost always no!
Semantics of output	“Margin”	Real probabilities