

UCB - CS189
Introduction to Machine Learning
Fall 2015
Lecture 5: Shrinkage

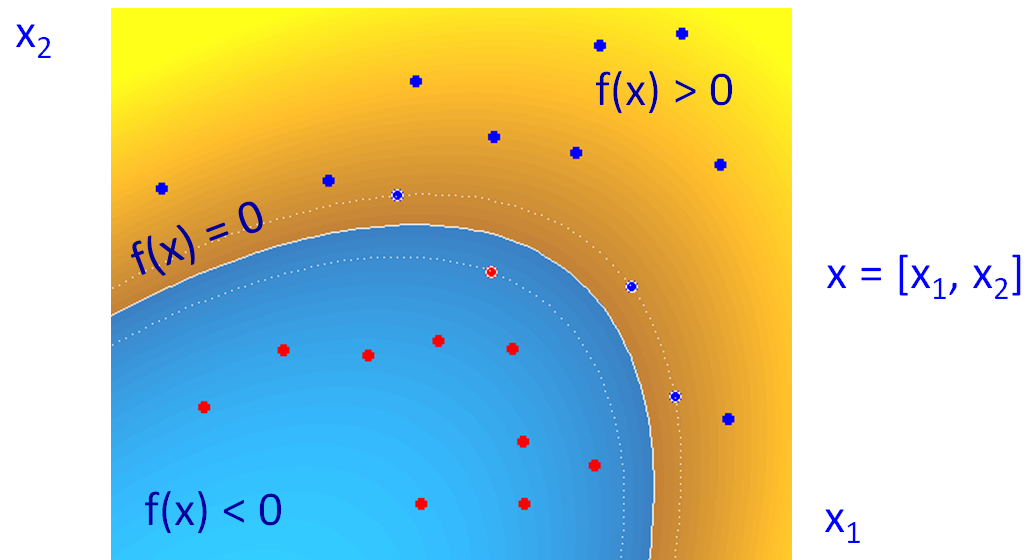
Isabelle Guyon
ChaLearn

Come to my office hours...

Wed 2:30-4:30 Soda 329

Last time

Non-linear optimum margin classifier **(HARD)**



$$f(x) = \sum_k \alpha_k k(x_k, x) + b$$

SVM, Boser-Guyon-Vapnik, 1992

Come to my office hours...
Wed 2:30-4:30 Soda 329

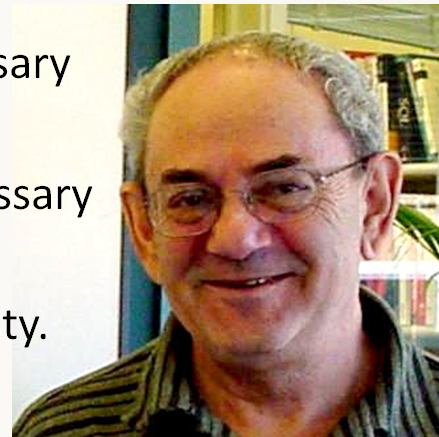
Today

For Ockham to Vapnik

(→ **SOFT margin**)



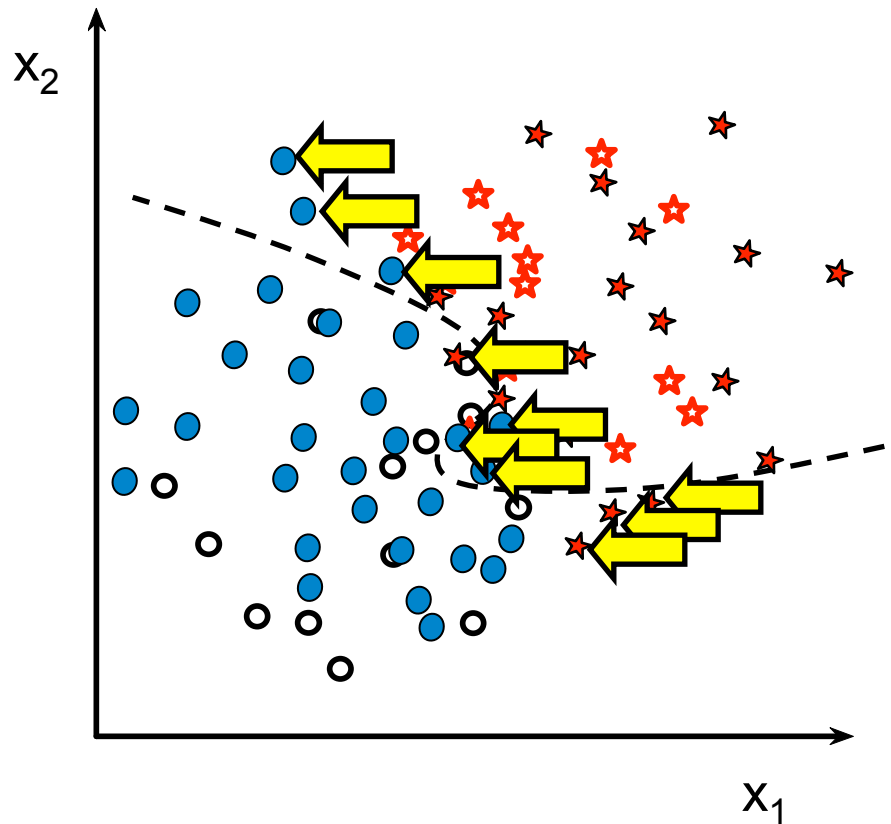
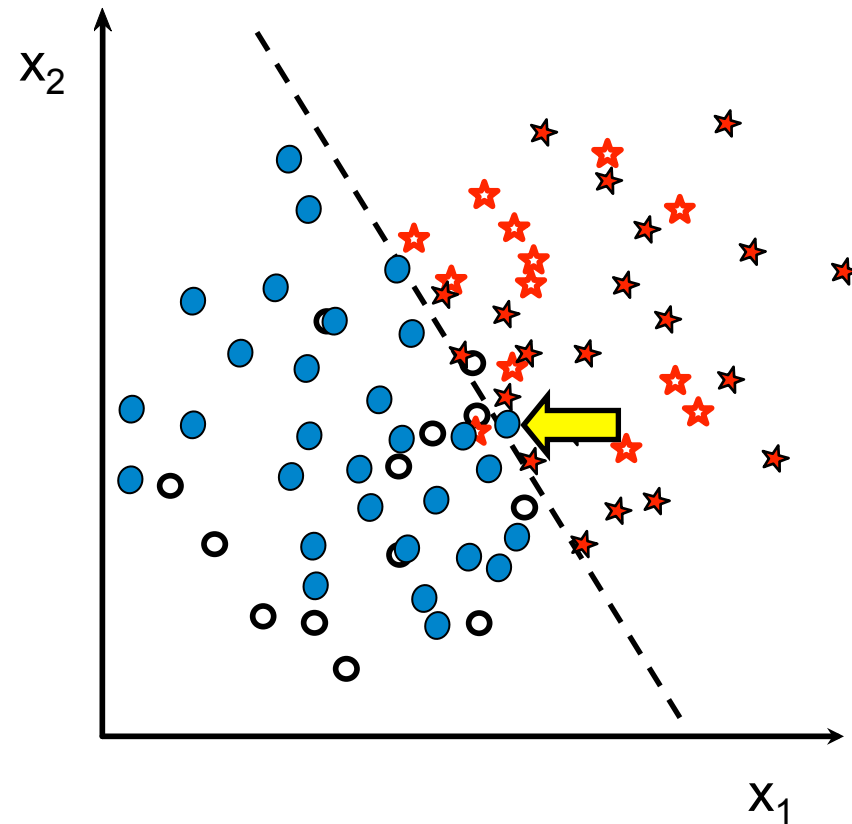
- Shave off unnecessary parameters.
- Forget the unnecessary memories.
- Minimize complexity.
- Minimize $\|\mathbf{w}\|$.



Math prerequisites

- Derivative
- Chain rule
- Lagrange multiplier

Fit / Robustness Tradeoff



Reasons for non-linear separability

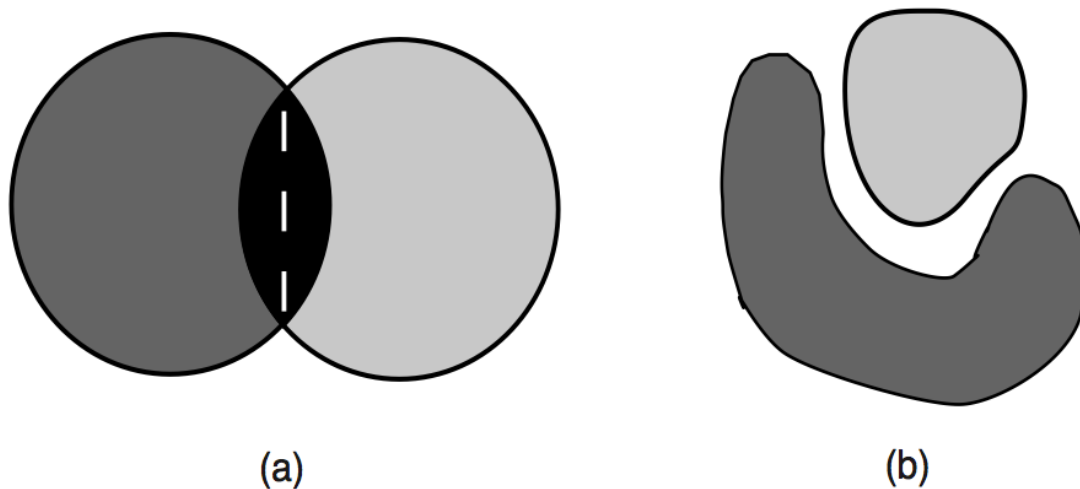


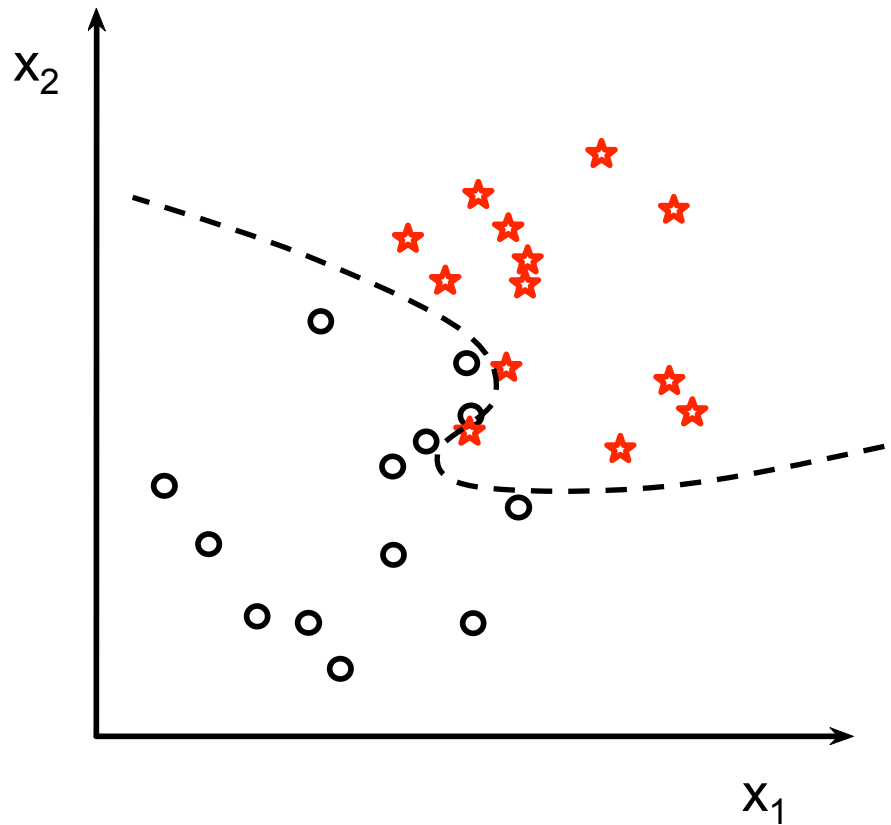
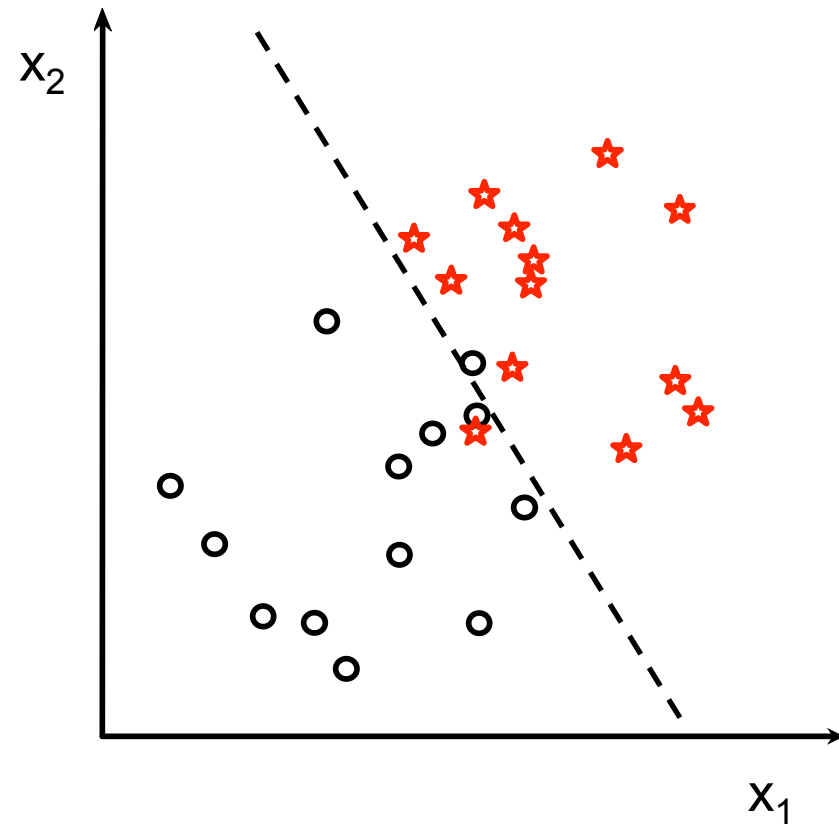
Figure 9.1 Non-linear separability. (a) Overlapping classes. The optimum decision boundary may still be linear. (b) Non overlapping classes. In the case shown, the optimum decision boundary is not linear.

Ockham's Razor



- Principle proposed by William of Ockham in the fourteenth century: “**Pluralitas non est ponenda sine neccesitate**”.
- Of two theories providing similarly good predictions, prefer *the simplest one*.
- Shave off unnecessary parameters of your models.

Fit / Robustness Tradeoff



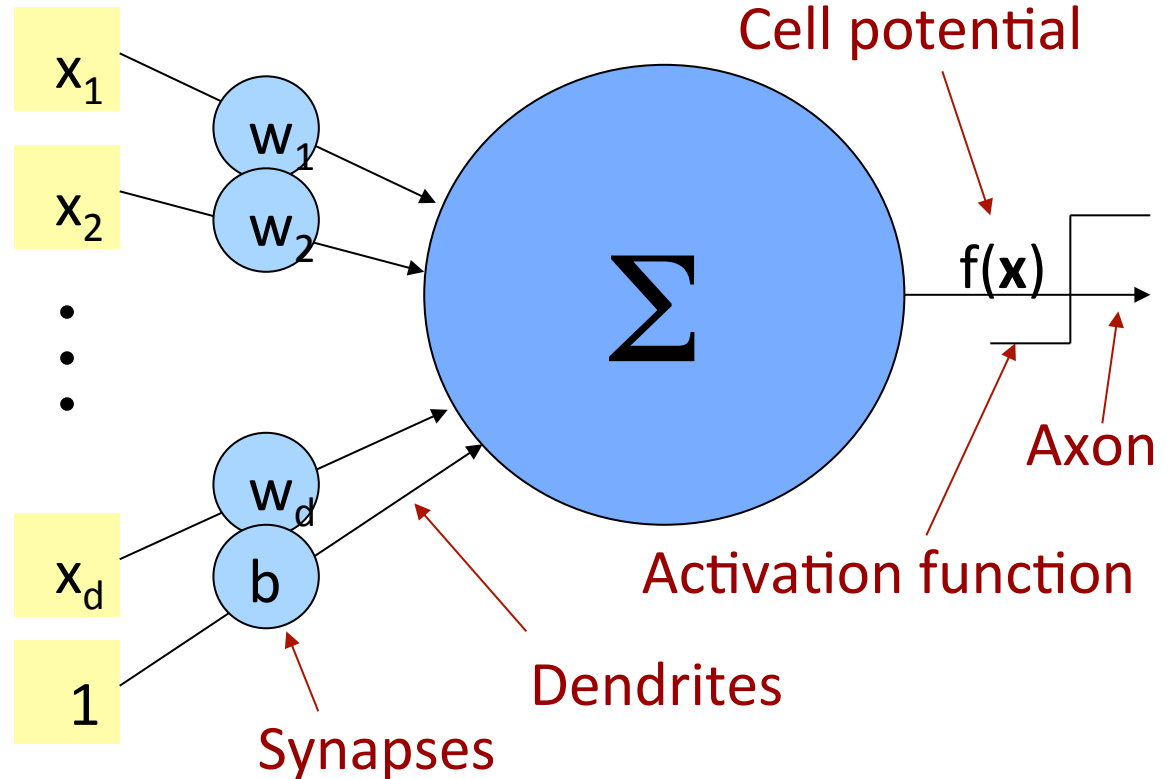
The Power of Amnesia

- The human **brain** is made out of billions of cells or Neurons, which are highly interconnected by synapses.
- Exposure to enriched environments with extra sensory and social stimulation enhances the **connectivity** of the synapses, but children and adolescents can lose them up to 20 million per day.

Artificial Neurons



Activation
of other
neurons



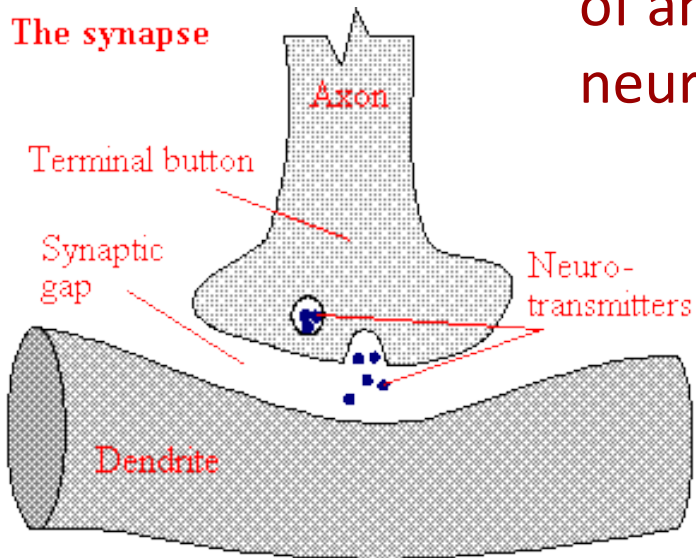
McCulloch and Pitts, 1943

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

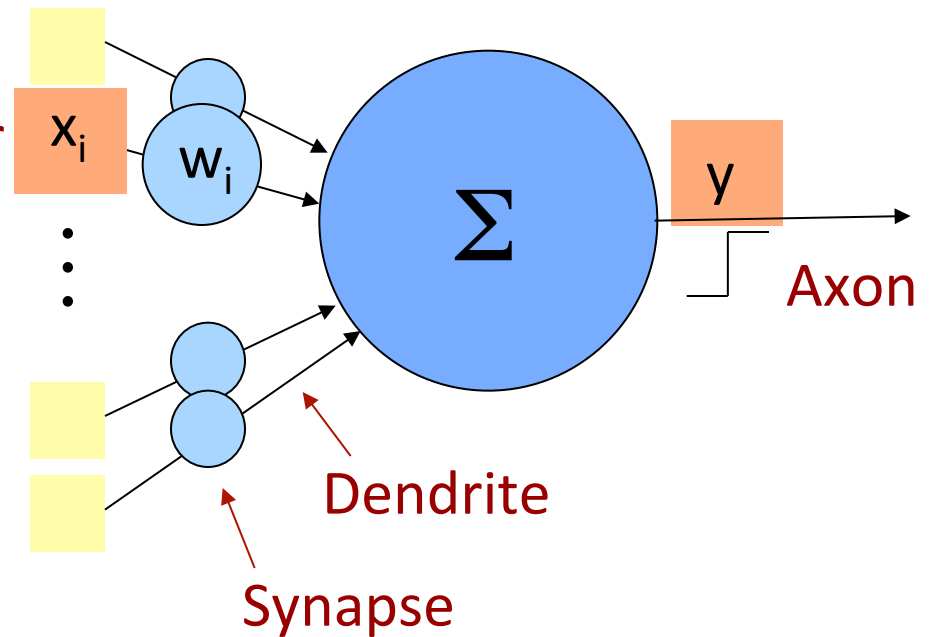
Hebb's Rule

$$w_i \leftarrow w_i + y^k x_i^k$$

The synapse



Activation
of another
neuron



Weight Decay

$$w_i \leftarrow w_i + y^k x_i^k$$

Hebb's rule

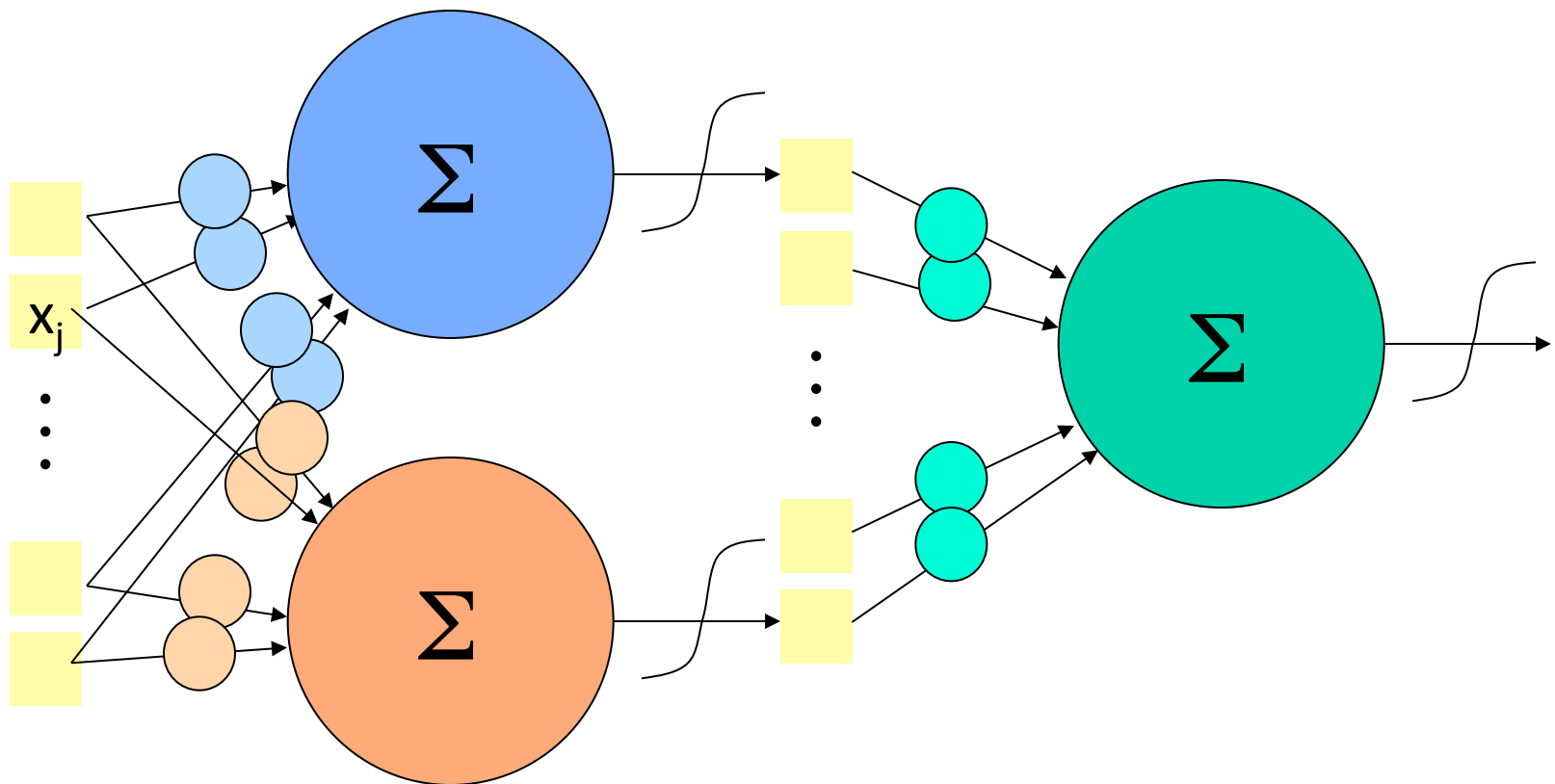
$$w_i \leftarrow (1-\gamma) w_i + y^k x_i^k$$

Weight decay

$\gamma \in [0, 1]$, decay parameter

Weight Decay for MLP

Replace: $w_i \leftarrow w_i + \text{back_prop}(i)$
by: $w_i \leftarrow (1-\gamma) w_i + \text{back_prop}(i)$



Notion of “Risk”

Last time: The risk is the sum of losses

$$R[f] = (1/N) \sum_{k=1:N} L(f(\mathbf{x}^k), y^k)$$

- $L(f(\mathbf{x}), y) = \mathbf{1}(f(\mathbf{x}) \neq y) = \mathbf{1}(yf(\mathbf{x}) < 0)$ zero-one loss
- $L(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2 = (yf(\mathbf{x}) - 1)^2$ square loss

with $y = \pm 1$

Today: Expected risk

$$R[f] = \int L(f(\mathbf{x}, \mathbf{w}), y) dP(\mathbf{x}, y)$$

Empirical risk

$$R_{\text{train}}[f] = (1/N) \sum_{k=1:N} L(f(\mathbf{x}^k), y^k)$$

Notion of “Risk”

Last time: The risk is the sum of losses

$$R[f] = (1/N) \sum_{k=1:N} L(f(\mathbf{x}^k), y^k)$$

- $L(f(\mathbf{x}), y) = \mathbf{1}(f(\mathbf{x}) \neq y)$ zero-one loss
- $L(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2$ square loss

with $y = \pm 1$

Today: Expected risk

$$R[f] = \int L(f(\mathbf{x}, \mathbf{w}), y) dP(\mathbf{x}, y)$$

Empirical risk

$$R_{\text{train}}[f] = (1/N) \sum_{k=1:N} L(f(\mathbf{x}^k), y^k)$$

Risk Minimization

- Learning problem: find the best function $f(\mathbf{x}; \mathbf{w})$ minimizing a **risk functional**

$$R[f] = \int \underbrace{L(f(\mathbf{x}; \mathbf{w}), y)}_{\text{loss function}} \underbrace{dP(\mathbf{x}, y)}_{\text{unknown distribution}}$$

- **Examples are given:**

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots (\mathbf{x}_N, y_N)$$

Approximations of $R[f]$

(generalization error)

- **Empirical risk:** $R_{\text{train}}[f] = (1/N) \sum_{k=1:N} L(f(\mathbf{x}^k, \mathbf{w}), y^k)$
 - 0/1 loss $\mathbf{1}(f(\mathbf{x}_k) \neq y_k)$: $R_{\text{train}}[f]$ = error rate
 - square loss $(f(\mathbf{x}_k) - y_k)^2$: $R_{\text{train}}[f]$ = mean square error

- **Guaranteed risk:**

With *high* probability $(1-\delta)$, $R[f] \leq R_{\text{gua}}[f]$

$$R_{\text{gua}}[f] = R_{\text{train}}[f] + \epsilon(\delta, C/N)$$

Approximations of $R[f]$

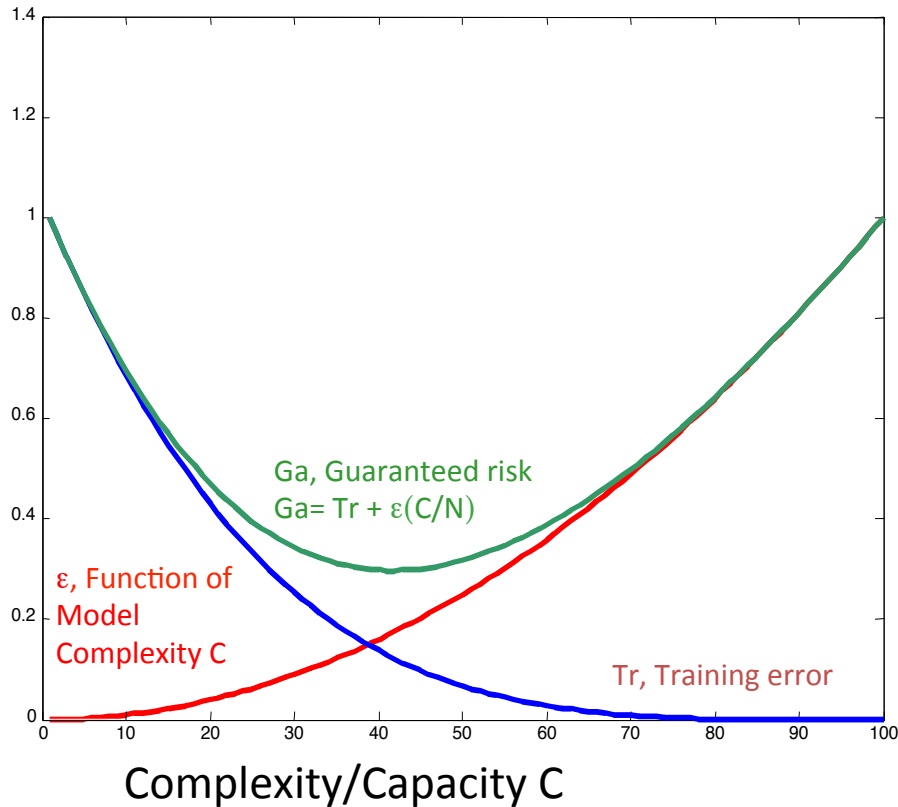
(generalization error)

- **Empirical risk:** $R_{\text{train}}[f] = (1/N) \sum_{k=1:N} L(f(\mathbf{x}^k, \mathbf{w}), y^k)$
 - 0/1 loss $\mathbf{1}(f(\mathbf{x}_k) \neq y_k)$: $R_{\text{train}}[f] = \text{error rate}$
 - square loss $(f(\mathbf{x}_k) - y_k)^2$: $R_{\text{train}}[f] = \text{mean square error}$
- **Guaranteed risk:**

With *high* probability $(1-\delta)$, $R[f] \leq R_{\text{gua}}[f]$

$$R_{\text{gua}}[f] = R_{\text{train}}[f] + \epsilon(\delta, C/N)$$

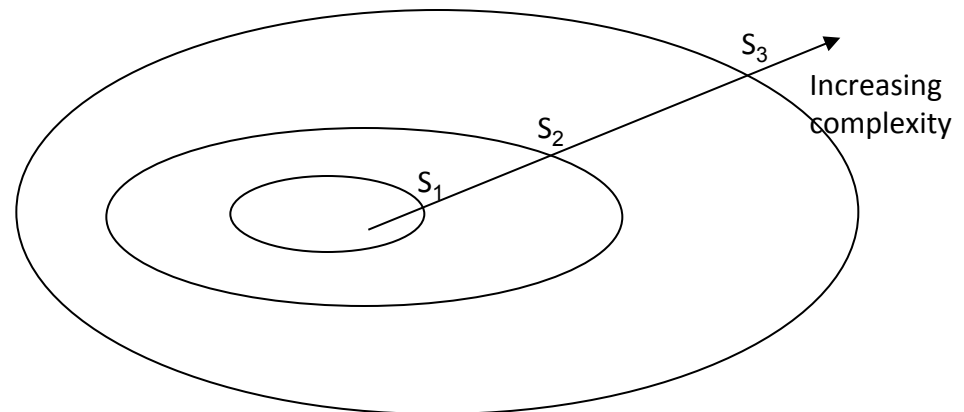
Structural Risk Minimization



Vapnik, 1974

Nested subsets of models, increasing complexity/capacity:

$$S_1 \subset S_2 \subset \dots S_N$$



SRM Example (linear model)

$$S_1 \subset S_2 \subset \dots S_N$$

- Rank with $\|\mathbf{w}\|^2 = \sum_i w_i^2$

$$S_k = \{ \mathbf{w} \mid \|\mathbf{w}\|^2 < \omega_k^2 \}, \omega_1 < \omega_2 < \dots < \omega_n$$

- Minimization under constraint:

$$\min R_{\text{train}}[f] \quad \text{s.t. } \|\mathbf{w}\|^2 < \omega_k^2$$

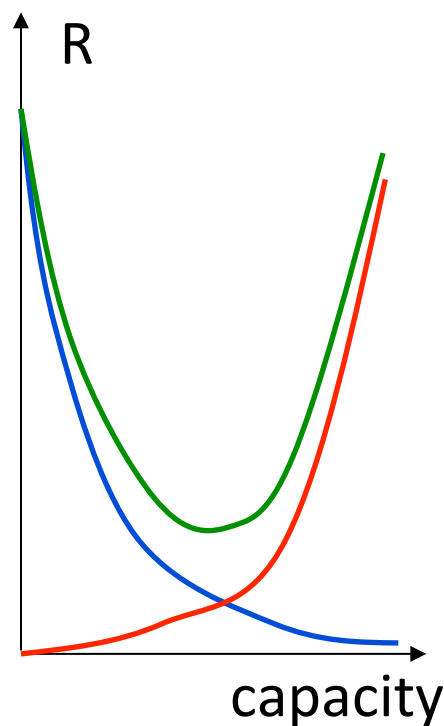
- Lagrangian:

$$R_{\text{reg}}[f, \lambda] = R_{\text{train}}[f] + \lambda (\|\mathbf{w}\|^2 - \omega_k^2), \lambda > 0$$

- Equivalent problems:

$$\min R_{\text{train}}[f] \quad \text{s.t. } \|\mathbf{w}\|^2 < \omega_k^2, \quad \omega_1 < \omega_2 < \dots < \omega_n$$

$$\min R_{\text{reg}}[f] = R_{\text{train}}[f] + \lambda_k \|\mathbf{w}\|^2, \quad 0 < \lambda_1 < \dots < \lambda_n$$



Gradient Descent

$$R_{\text{reg}}[f] = R_{\text{train}}[f] + \lambda \|\mathbf{w}\|^2 \quad \text{SRM/regularization}$$

$$w_j \leftarrow w_j - \eta \partial R_{\text{reg}} / \partial w_j$$

$$w_j \leftarrow w_j - \eta \partial R_{\text{train}} / \partial w_j - 2 \eta \lambda w_j \quad \gamma = 2 \eta \lambda$$

$$w_j \leftarrow (1 - \gamma) w_j - \eta \partial R_{\text{train}} / \partial w_j \quad \text{Weight decay}$$

Example: **Mean square error:** $(1/N) \sum_{k=1:N} \underbrace{(f(\mathbf{x}^k) - y^k)^2}_{\text{RSS}}$

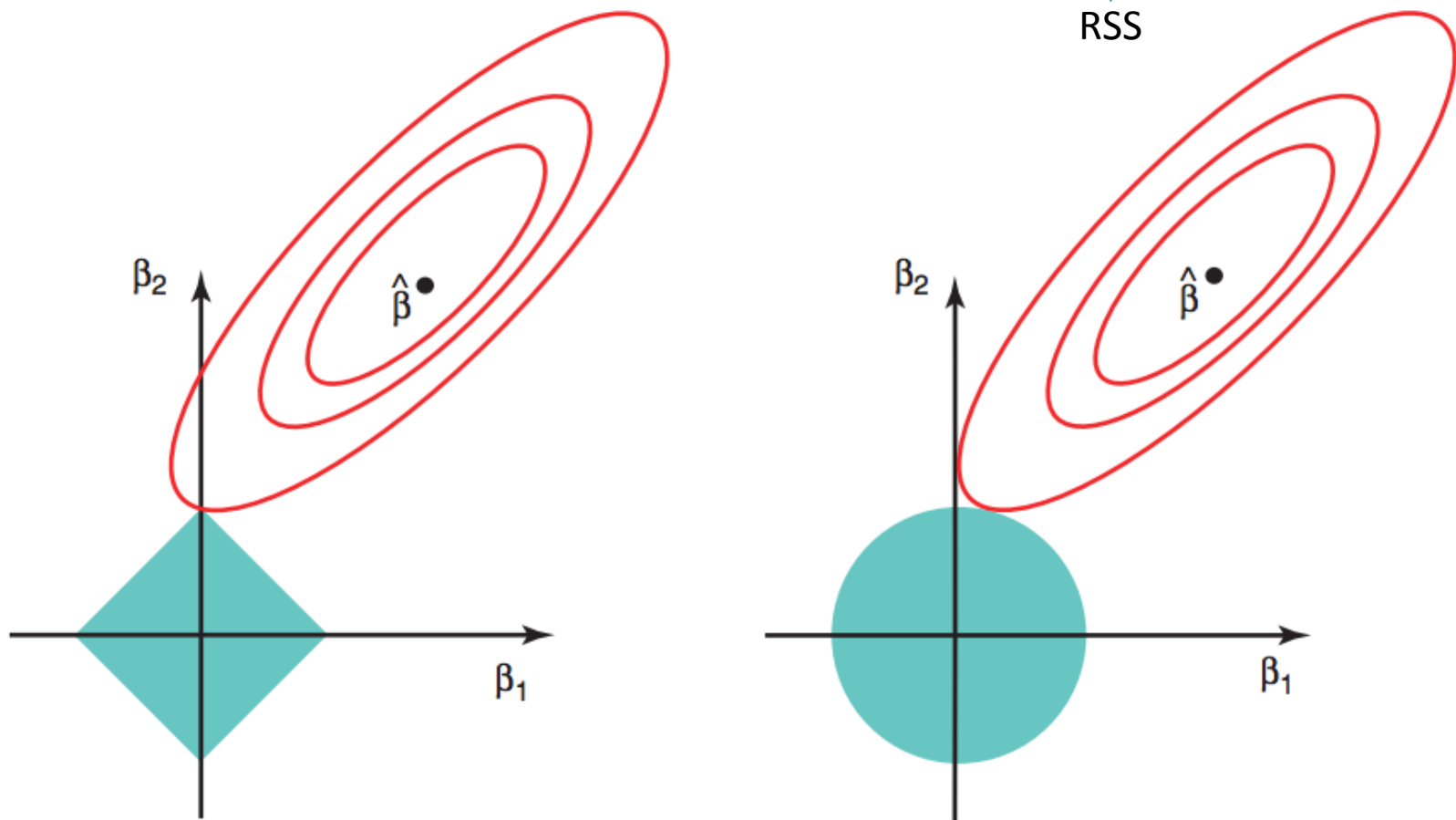
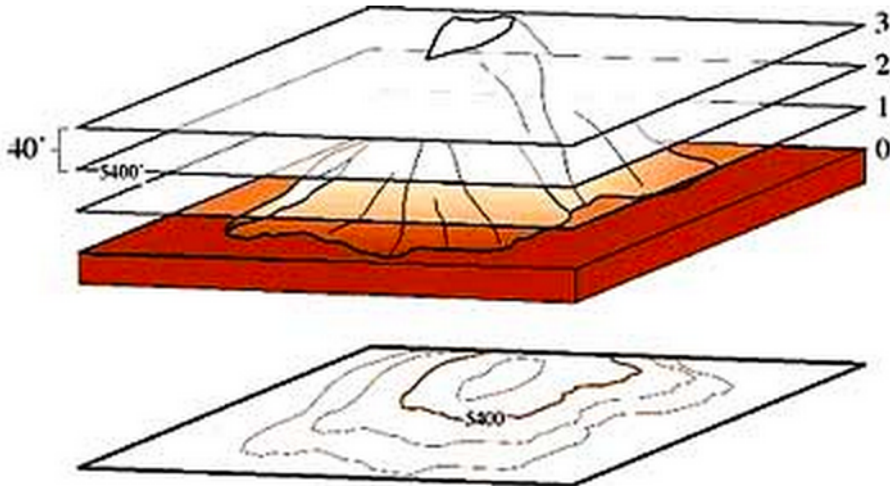
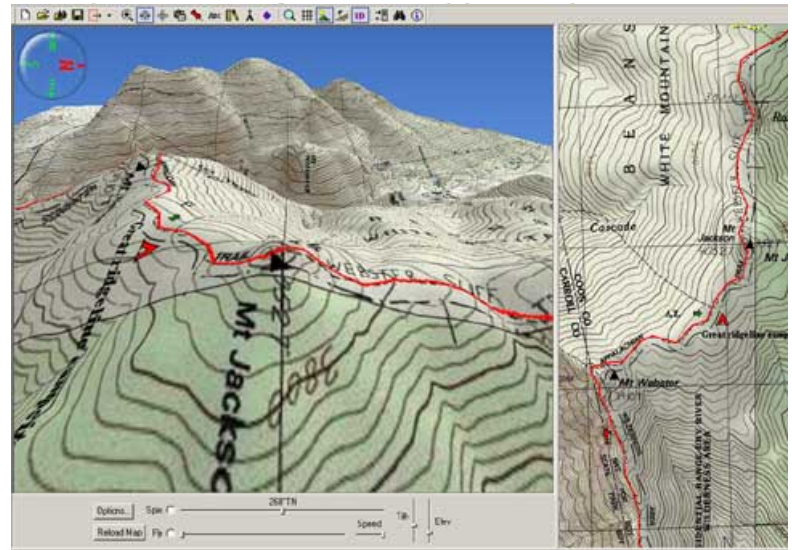


FIGURE 6.7. *Contours of the error and constraint functions for the lasso (left) and ridge regression (right). The solid blue areas are the constraint regions, $|\beta_1| + |\beta_2| \leq s$ and $\beta_1^2 + \beta_2^2 \leq s$, while the red ellipses are the contours of the RSS.*

Contour maps



http://geology.isu.edu/geostac/Field_Exercise/topomaps/topo_map.htm



http://www.nationalgeographic.com/adventure/images/02_06/Appalachian_TOPO_5.jpg

Shrinkage justified

Why do we want “simple” models?

Everything is about mitigating risk

Guaranteed risk:

With *high* probability $(1-\delta)$, $R[f] \leq R_{\text{gua}}[f]$

$$R_{\text{gua}}[f] = R_{\text{train}}[f] + \epsilon(\delta, C/N)$$

Regularized risk:

$$R_{\text{reg}}[f] = R_{\text{train}}[f] + \lambda \|w\|^2$$

FIT

ROBUSTNESS

Multiple Structures

- **Shrinkage (weight decay, ridge regression, SVM):**

$$S_k = \{ \mathbf{w} \mid \|\mathbf{w}\|_2 < \omega_k \}, \omega_1 < \omega_2 < \dots < \omega_k$$

$$\gamma_1 > \gamma_2 > \gamma_3 > \dots > \gamma_k \quad (\gamma \text{ is the ridge})$$

- **Feature selection:**

$$S_k = \{ \mathbf{w} \mid \|\mathbf{w}\|_0 < v_k \},$$

$$v_1 < v_2 < \dots < v_k \quad (v \text{ is the number of features})$$

- **Kernel parameters** $k(\mathbf{s}, \mathbf{t}) = (\mathbf{s} \bullet \mathbf{t} + 1)^q$:

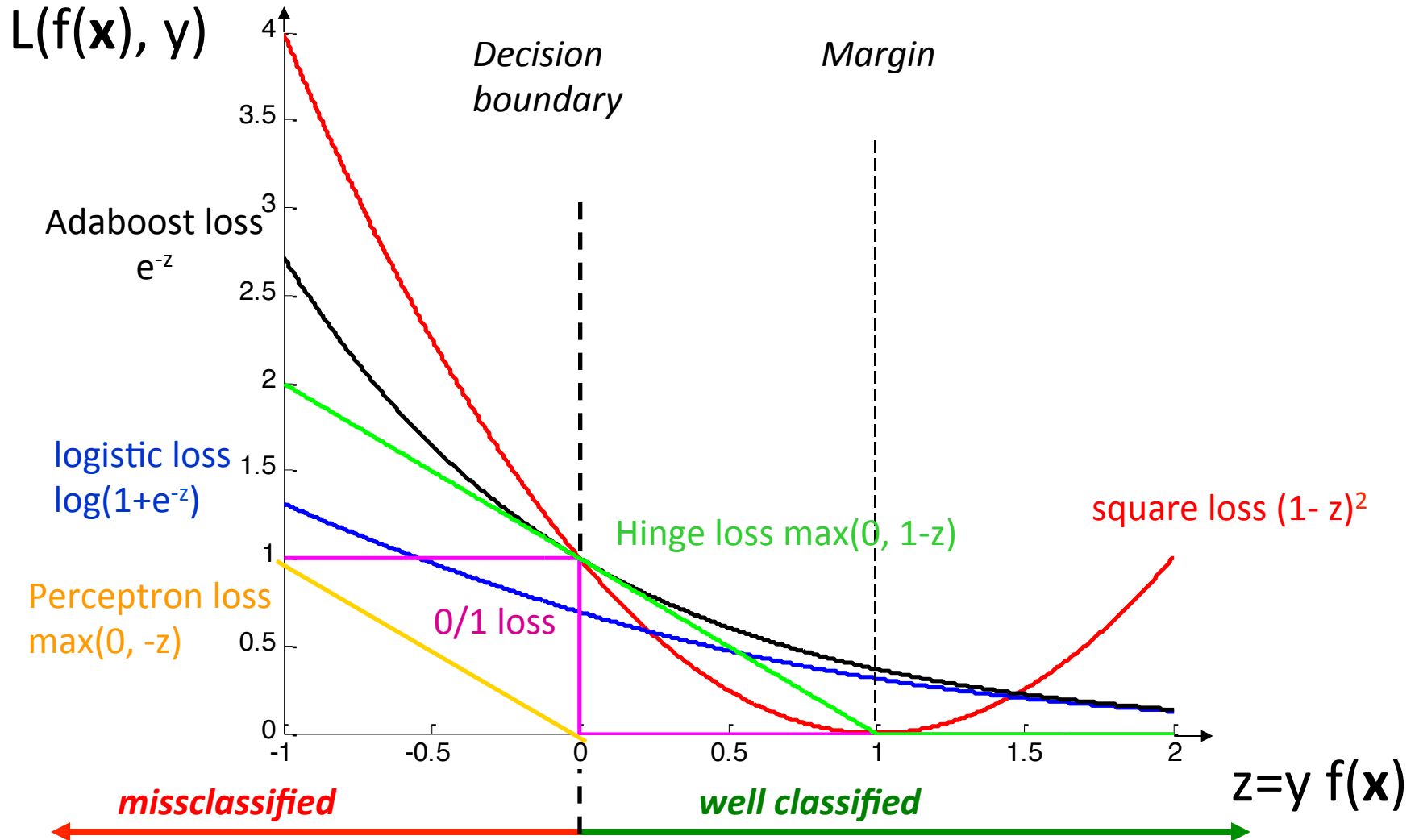
$$q_1 < q_2 < \dots < q_k \quad (q \text{ is the polynomial degree})$$

$$k(\mathbf{s}, \mathbf{t}) = \exp(-\|\mathbf{s} - \mathbf{t}\|^2 / \sigma^2)$$

$$\sigma_1 > \sigma_2 > \sigma_3 > \dots > \sigma_k \quad (\sigma \text{ is the kernel width})$$

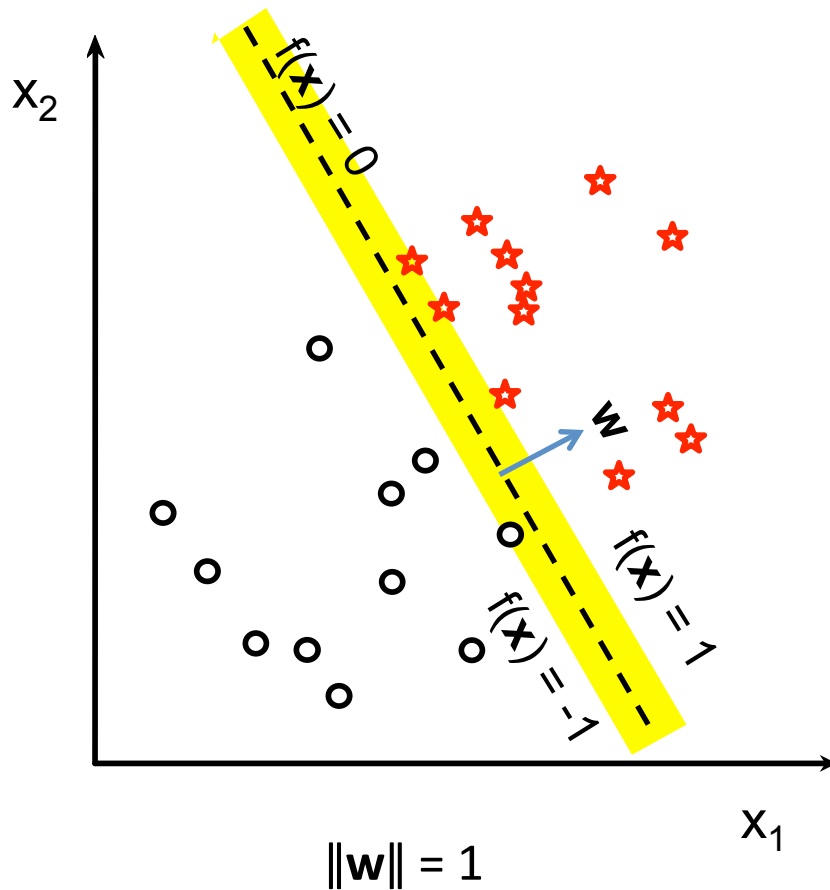
Loss Functions

The risk is the average of the loss.

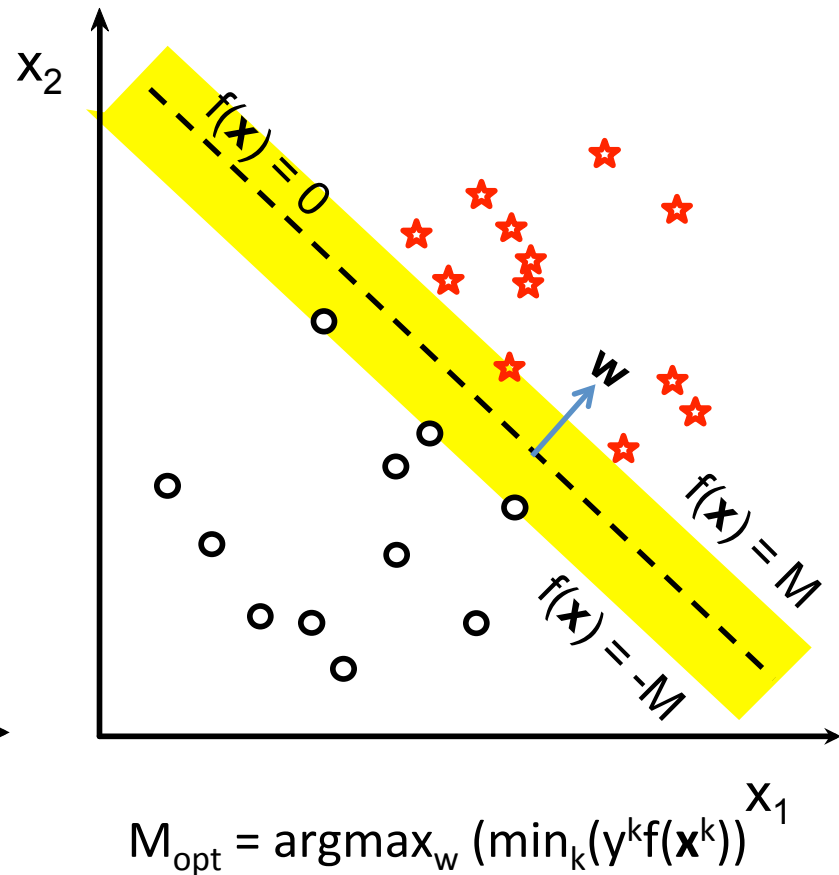


Maximizing the margin

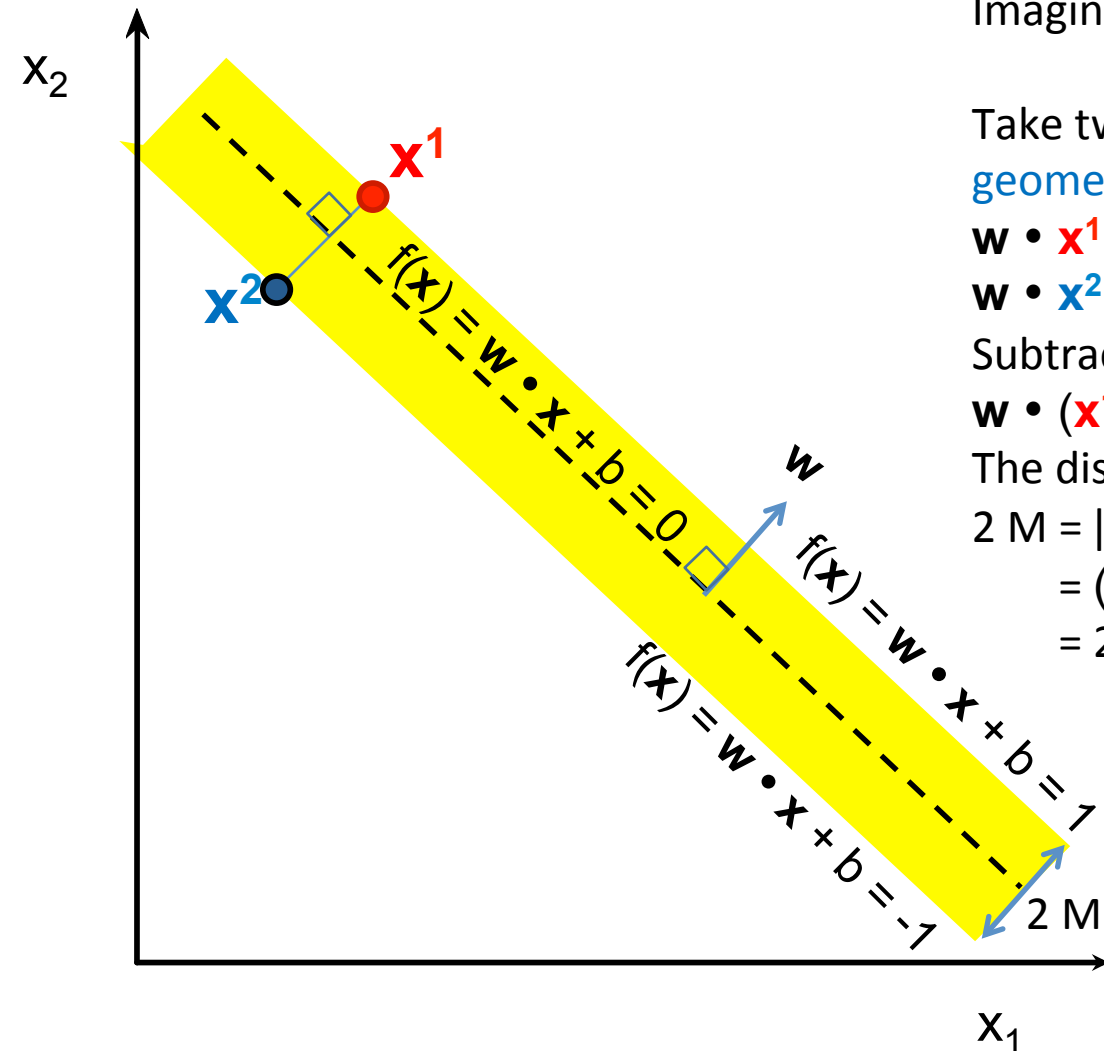
Large margin



Optimum margin



Linking geometrical and functional margin



Imagine that we want to impose $y^k f(\mathbf{x}^k) \geq 1$

$y^k f(\mathbf{x}^k)$ is the functional margin

Take two points on opposite sides of
geometrical margin:

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

$$\mathbf{w} \cdot \mathbf{x}^2 + b = -1$$

Subtract:

$$\mathbf{w} \cdot (\mathbf{x}^1 - \mathbf{x}^2) = 2$$

The distance between \mathbf{x}^1 and \mathbf{x}^2 is:

$$\begin{aligned} 2M &= \|\mathbf{x}^1 - \mathbf{x}^2\| \\ &= (\mathbf{x}^1 - \mathbf{x}^2) \cdot \mathbf{w} / \|\mathbf{w}\| \\ &= 2 / \|\mathbf{w}\| \end{aligned}$$

$$M = 1 / \|\mathbf{w}\|$$

Equivalent formulations

$$\|\mathbf{w}\| = 1$$

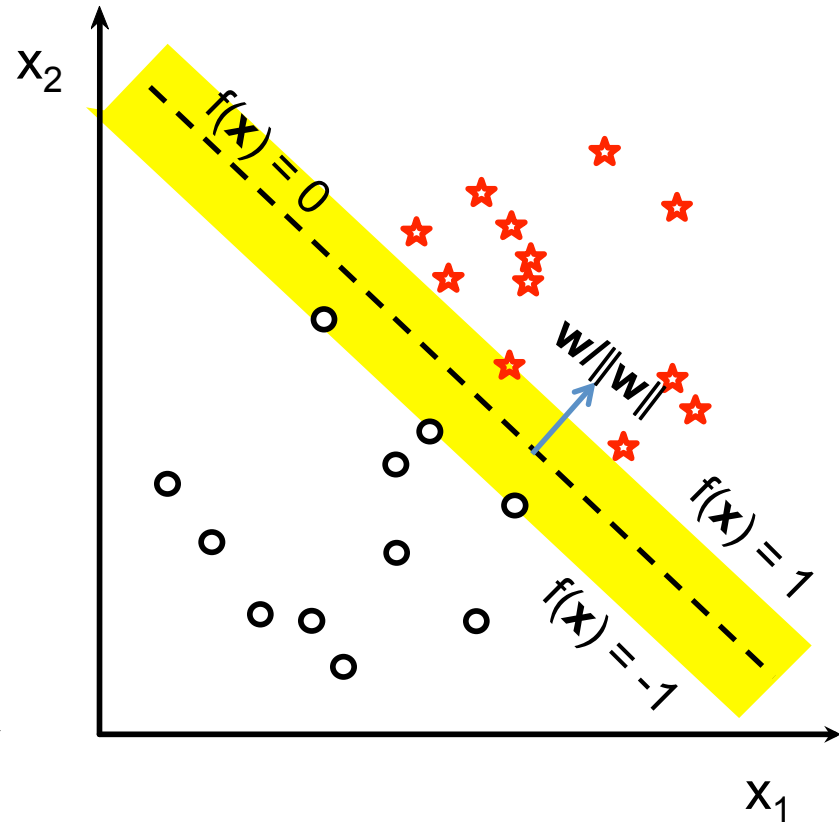
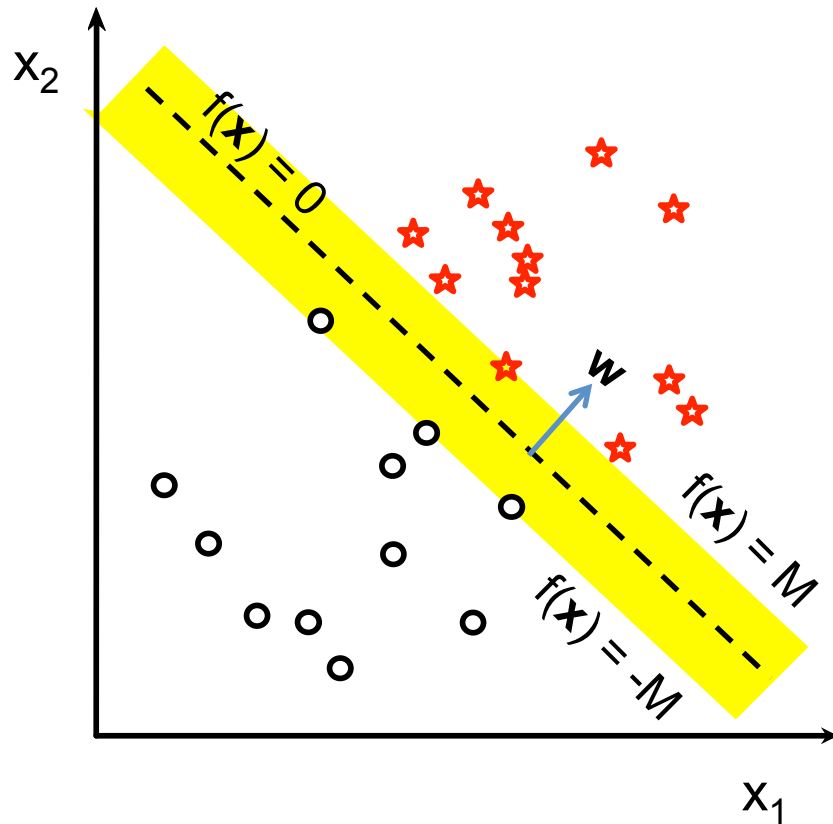
$$M_{\text{opt}} = \operatorname{argmax}_{\mathbf{w}} (\min_k (y^k f(\mathbf{x}^k)))$$



$$M = 1 / \|\mathbf{w}\|$$

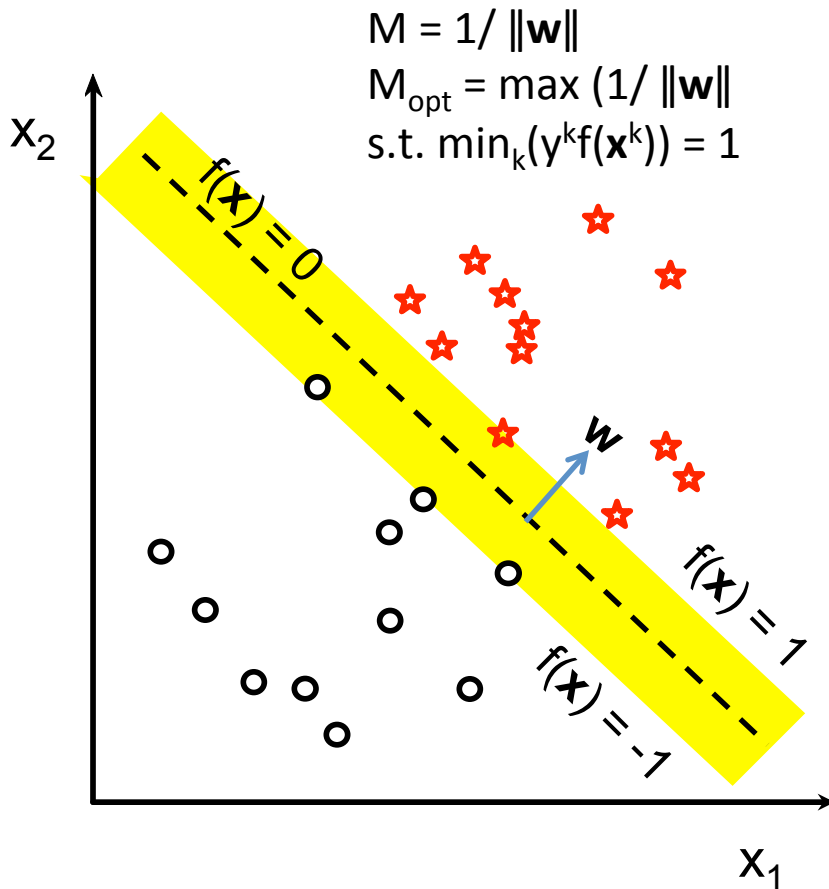
$$M_{\text{opt}} = \max (1 / \|\mathbf{w}\|)$$

$$\text{s.t. } \min_k (y^k f(\mathbf{x}^k)) = 1$$

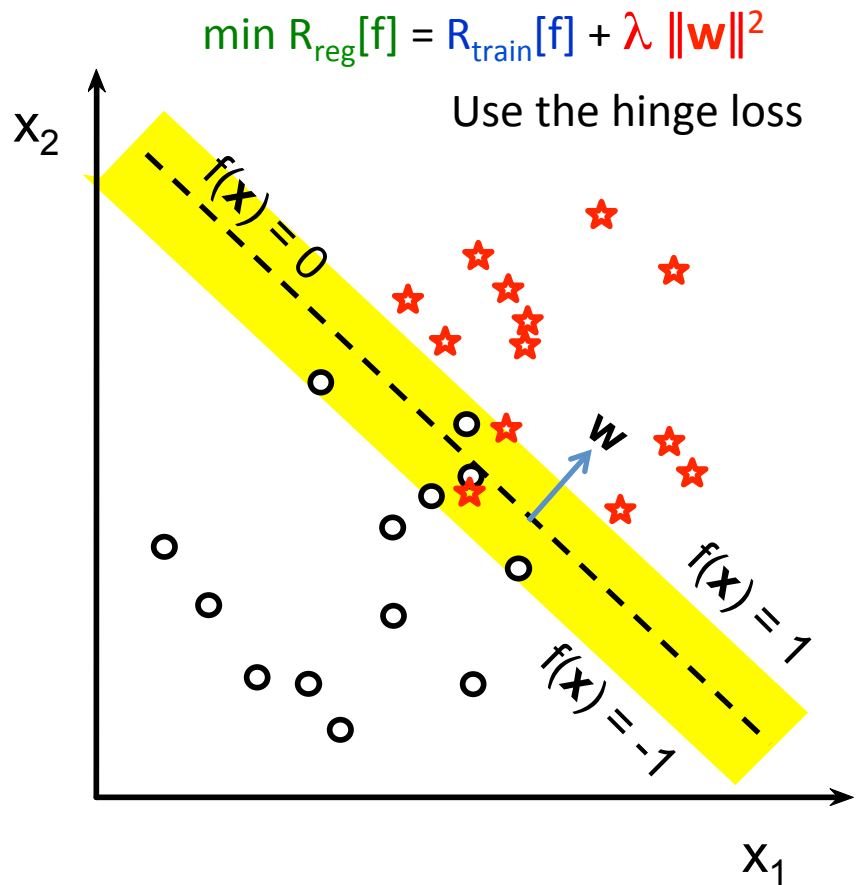


Optimum margin

Hard margin



Soft margin

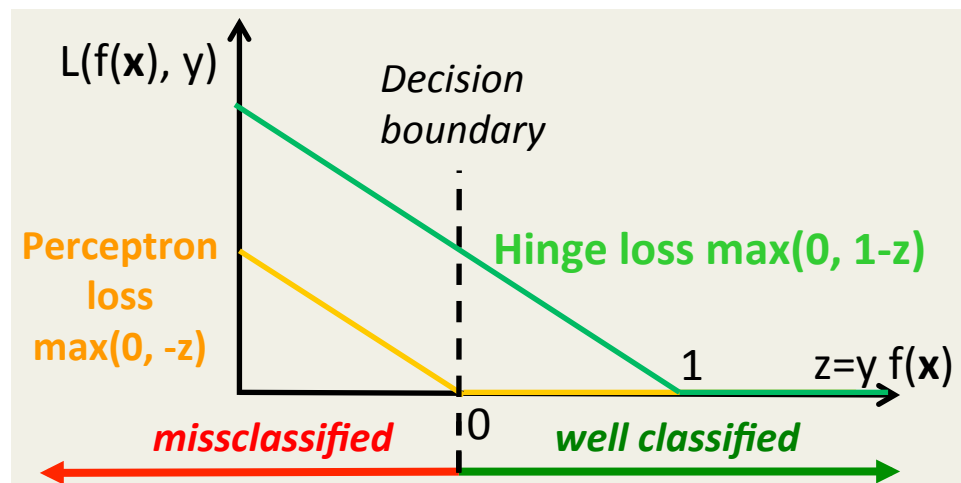


Large margin Perceptron with weight decay = soft margin

$$L_{\text{hinge}} = \max(0, 1 - z)$$

$$z = y f(\mathbf{x})$$

$$\min_{\mathbf{w}} (L_{\text{hinge}} + \lambda \|\mathbf{w}\|^2)$$



$$w_i \leftarrow \begin{cases} (1-\gamma) w_i + \eta y x_i, & \text{if } z < 1 \text{ (misclassified or within margin)} \\ (1-\gamma) w_i & \text{otherwise} \end{cases} \quad (\gamma = 2 \eta \lambda)$$

Soft Margin **C**ompromise

Minimize

(1/Margin) + **C Training error**

Good robustness

Good fit

Soft Margin **C**ompromise

Minimize

$$\lambda \|w\|^2$$

Good robustness

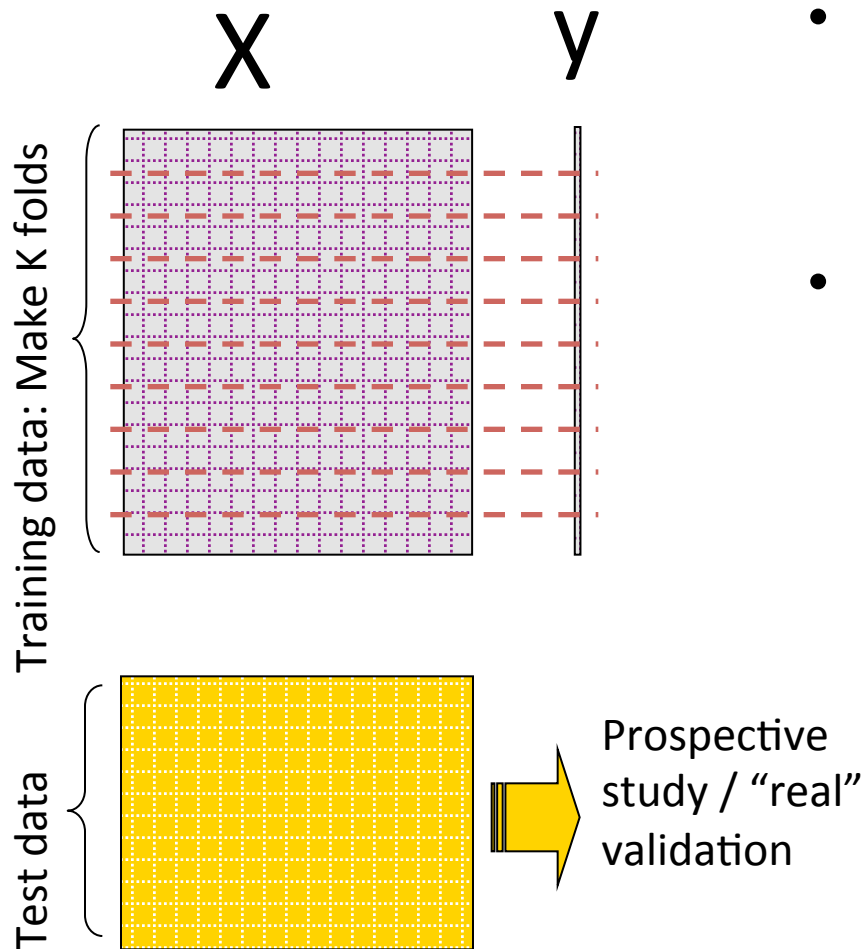
+

$$R_{\text{train}}[f]$$

$$\sum_k \max(0, 1 - z^k)$$

Good fit

Hyper-parameter Selection



- **Learning = adjusting:**
parameters (\mathbf{w} vector).
hyper-parameters (γ, ν, q, σ).
- ***Cross-validation with K-folds:***

For various values of γ, ν, q, σ :

- Adjust \mathbf{w} on a fraction $(K-1)/K$ of training examples *e.g.* $9/10^{\text{th}}$.
- Test on $1/K$ remaining examples *e.g.* $1/10^{\text{th}}$.
- Rotate examples and average test results (CV error).
- Select γ, ν, q, σ to minimize CV error.
- Re-compute \mathbf{w} on **all** training examples using optimal γ, ν, q, σ .

Summary

- High complexity models may “**overfit**”:
 - Fit perfectly training examples
 - Generalize poorly to new cases
- **SRM solution**: organize the models in nested subsets such that in every structure element
complexity $< \theta$
- **Regularization**: Formalize learning as a constrained optimization problem, minimize
regularized risk = training error + λ complexity
- Both formulations are equivalent via the use of Lagrange multipliers.
- θ and λ are hyperparameters, which can be optimized by cross-validation.

Come to my office hours...
Wed 2:30-4:30 Soda 329

Next time

