

Lecture 2: Binary Classification

Kai-Wei Chang
CS @ UCLA
kw@kwchang.net

Course webpage: <https://uclanlp.github.io/CS269-17/>

Announcements

- ❖ Waiting list: If you're not enrolled, please sign up.
- ❖ We will use Piazza as an online discussion platform. Please sign up here:
piazza.com/ucla/fall2017/cs269

This Lecture

- ❖ Supervised Learning
- ❖ Linear Classifiers
 - ❖ Perceptron Algorithm
 - ❖ Support Vector Machine
 - ❖ Logistic Regression
- ❖ Optimization in machine learning

The Badges game

+ Naoki Abe

- Eric Baum

- ❖ Conference attendees to the ICML 1994 were given **name badges labeled with + or -**.
- ❖ What function was used to assign these labels?

Training data

- | | | |
|---------------------|-------------------|--------------------|
| + Naoki Abe | + Peter Bartlett | + Carla E. Brodley |
| - Myriam Abramson | - Eric Baum | + Nader Bshouty |
| + David W. Aha | + Welton Becket | - Wray Buntine |
| + Kamal M. Ali | - Shai Ben-David | - Andrey Burago |
| - Eric Allender | + George Berg | + Tom Bylander |
| + Dana Angluin | + Neil Berkman | + Bill Byrne |
| - Chidanand Apte | + Malini Bhandaru | - Claire Cardie |
| + Minoru Asada | + Bir Bhanu | + John Case |
| + Lars Asker | + Reinhard Blasig | + Jason Catlett |
| + Javed Aslam | - Avrim Blum | - Philip Chan |
| + Jose L. Balcazar | - Anselm Blumer | - Zhixiang Chen |
| - Cristina Baroglio | + Justin Boyan | - Chris Darken |

Raw test data

Gerald F. DeJong

Chris Drummond

Yolanda Gil

Attilio Giordana

Jiarong Hong

J. R. Quinlan

Priscilla Rasmussen

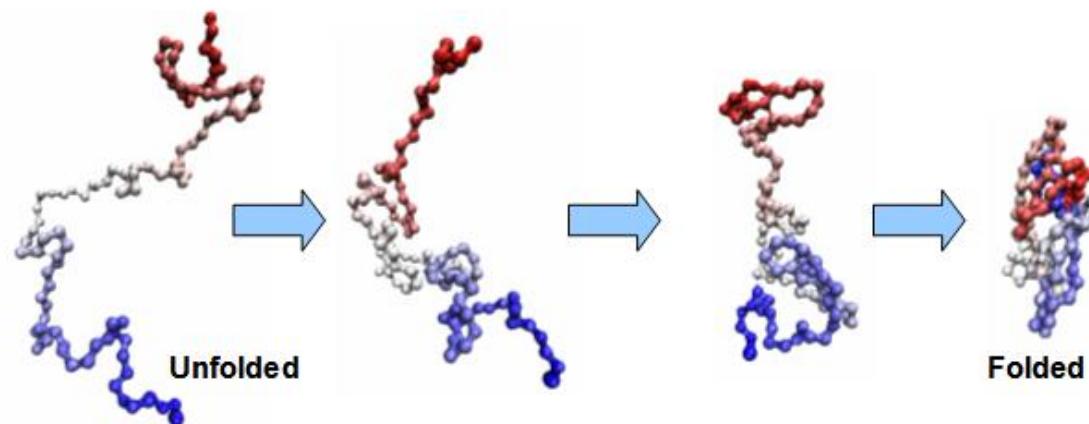
Dan Roth

Yoram Singer

Lyle H. Ungar

Why we need machine learning?

- ❖ There is no (or limited numbers of) human expert for some problems
 - ❖ E.g.: Identify DNA binding sites, predicting disease progression, predicting protein folding structure



Why we need machine learning?

- ❖ There is no (or limited numbers of) human expert for some problems
- ❖ Humans can perform a task, but can't describe how they do it
 - ❖ E.g.: Object recognition

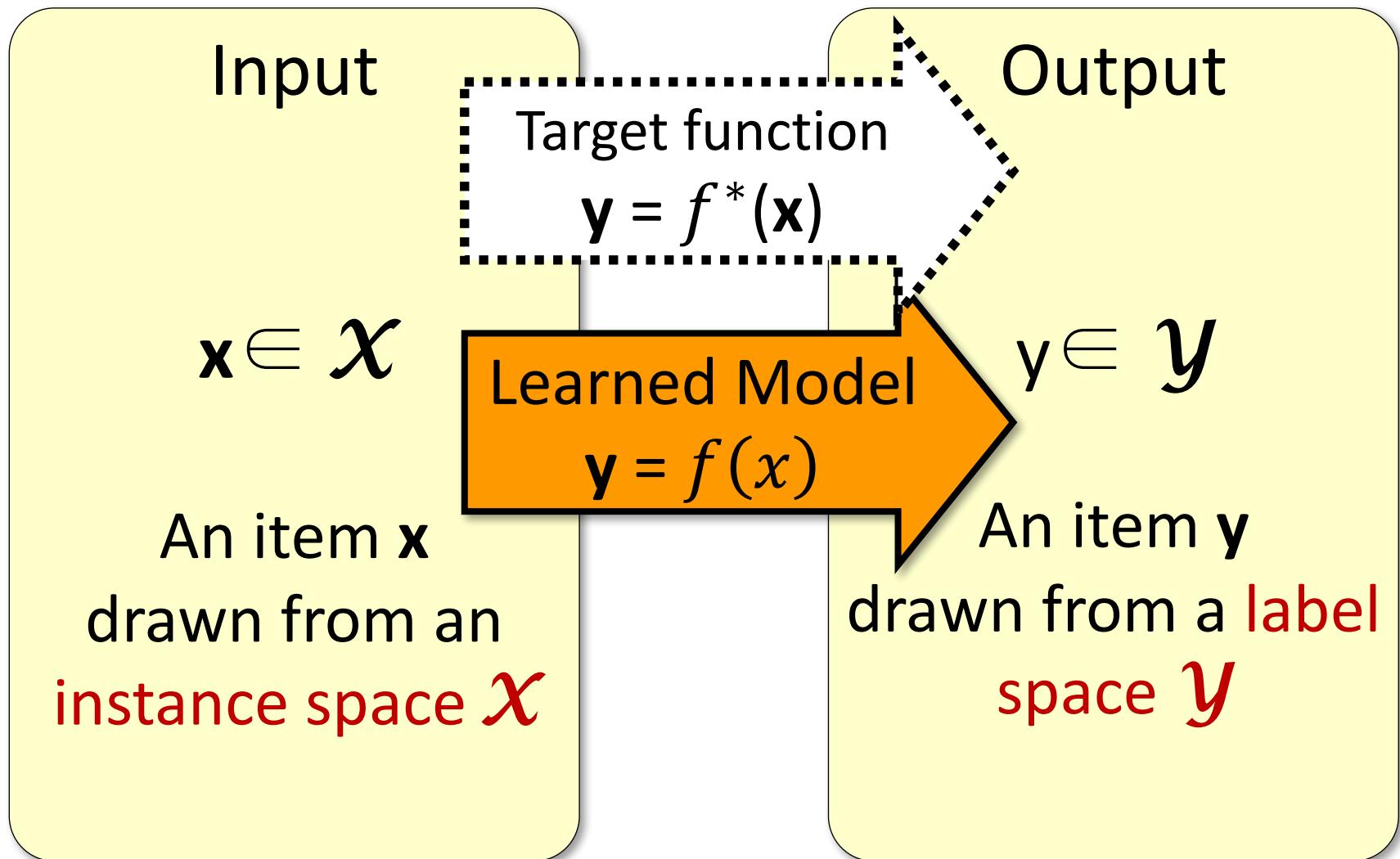


Why we need machine learning?

- ❖ There is no (or limited numbers of) human expert for some problems
- ❖ Humans can perform a task, but can't describe how they do it
- ❖ The desired function is hard to be written down in a closed form
 - ❖ E.g.,: predict stock price



Supervised learning



Supervised learning

Input

$$\mathbf{x} \in \mathcal{X}$$

An item \mathbf{x}
drawn from an
instance space \mathcal{X}

\mathbf{x} is represented in a **feature space**

- Typically $\mathbf{x} \in \{0,1\}^n$ or R^N
- Usually represented as a **vector**
- We call it **input vector**

Supervised learning

y is represented in output space
(label space)

Different kinds of output:

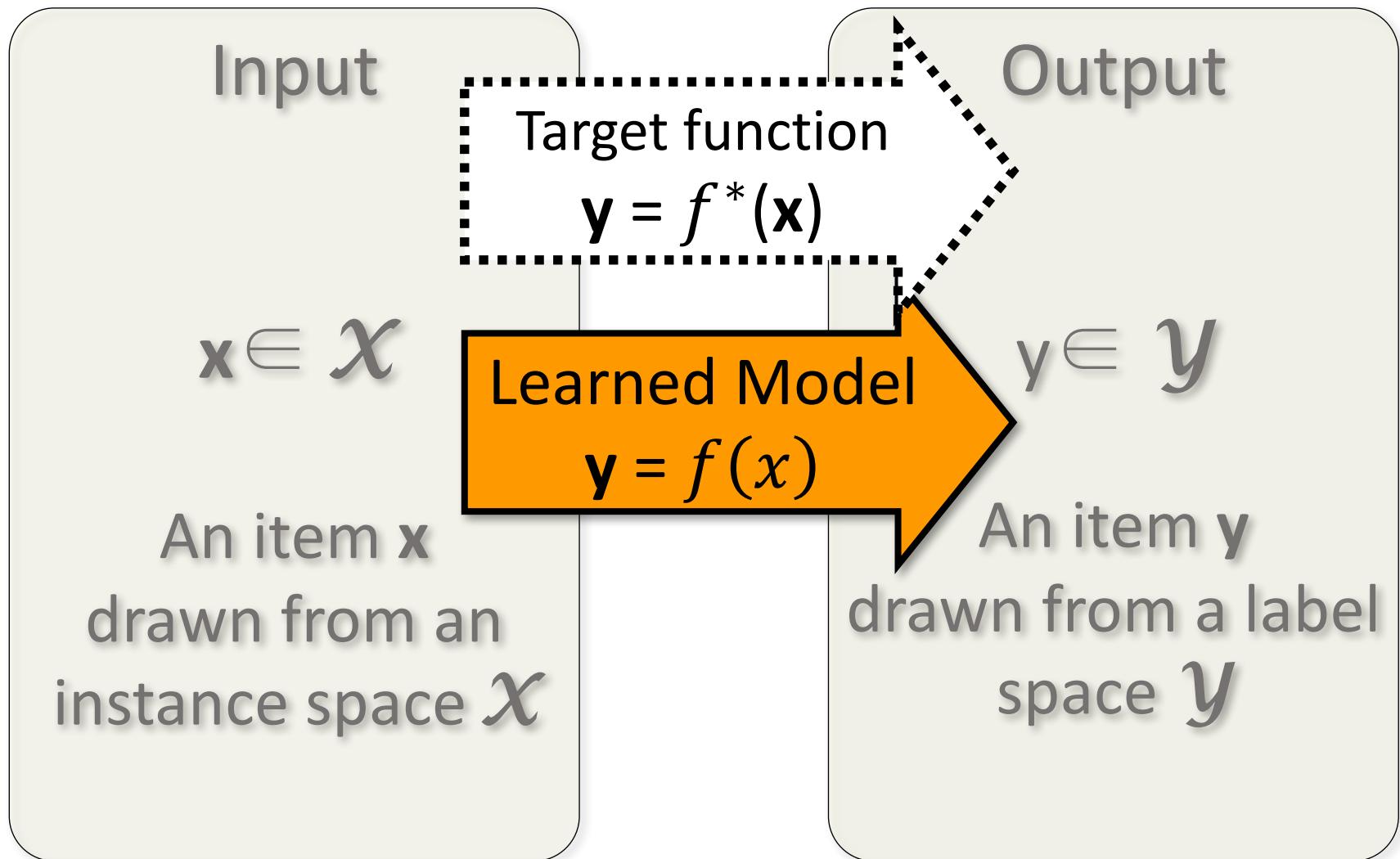
- Binary classification:
 $y \in \{-1, 1\}$
- Multiclass classification:
 $y \in \{1, 2, 3, \dots, K\}$
- Regression:
 $y \in R$
- Structured output
 $y \in \{1, 2, 3, \dots, K\}^N$

Output

$$y \in \mathcal{Y}$$

An item y
drawn from a label
space \mathcal{Y}

Learning the mapping

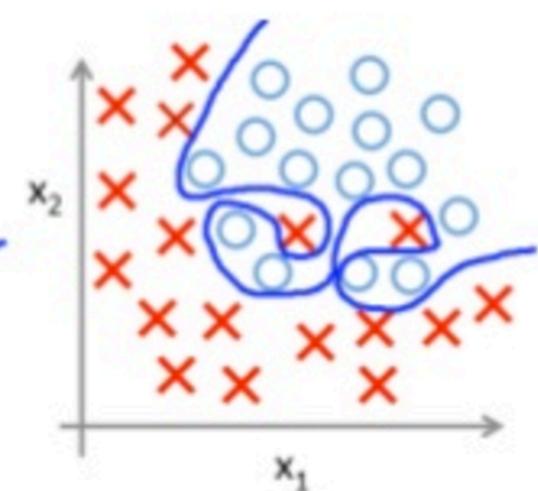
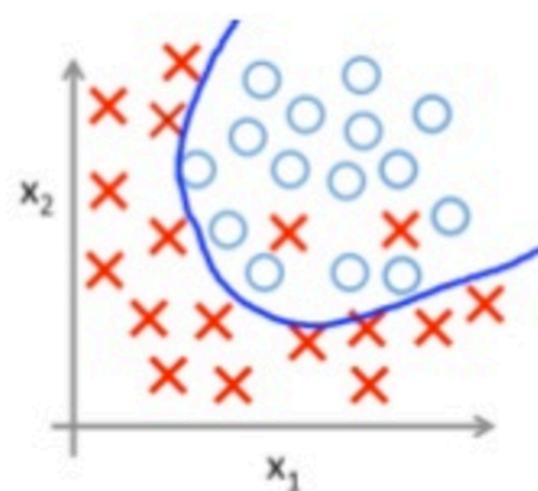
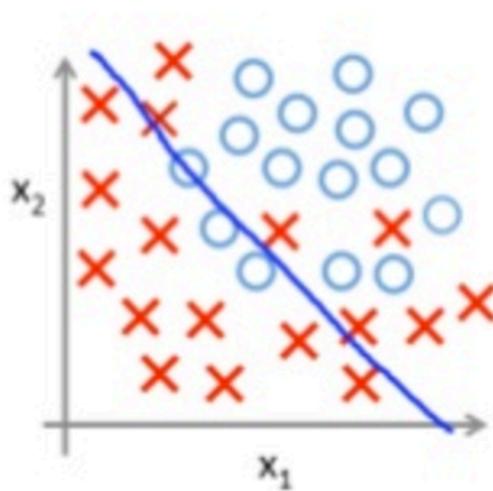


Goal

- ❖ Find a good approximation of $f^*(\cdot)$
- ❖ Good in what sense?

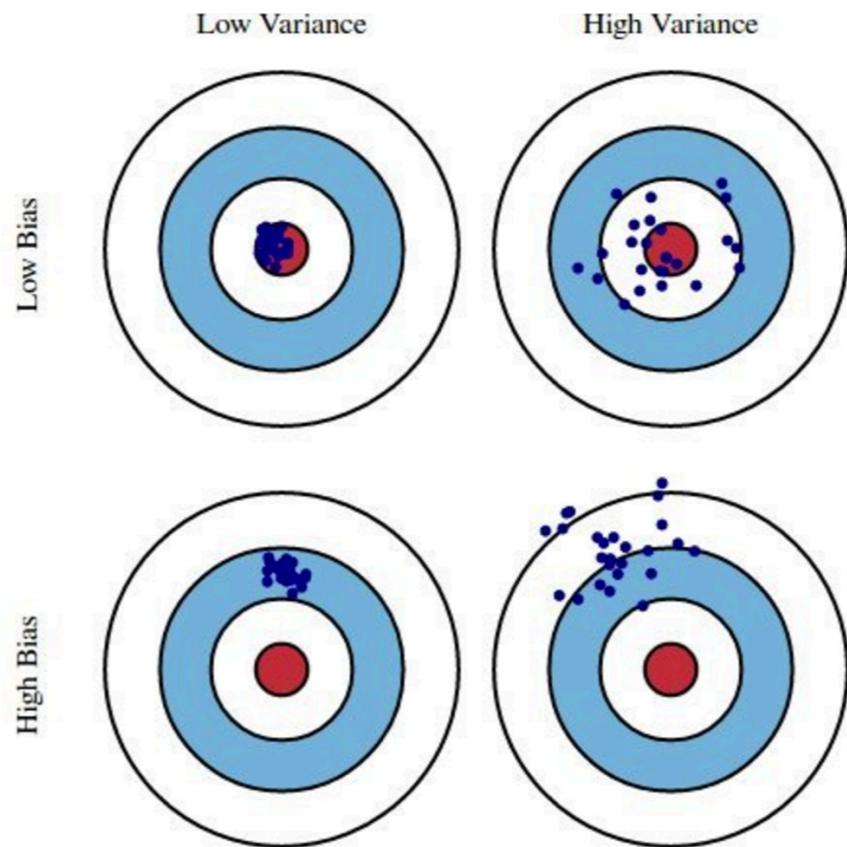
Under-fitting and over-fitting

- ❖ Which classifier (blue line) is the best one?



Bias V.S. Variance

- ❖ Remember, training data are subsamples drawn from the true distribution
- ❖ Exam strategy:
 - ❖ Study every chapter well
 - ❖ A+: Low var & bias
 - ❖ Study only a few chapters
 - ❖ A+? B? C? Low bias; High var
 - ❖ Study every chapter roughly
 - ❖ B+: Low var; high bias
 - ❖ Go to sleep
 - ❖ B ~D: High var, high bias



Questions of interest

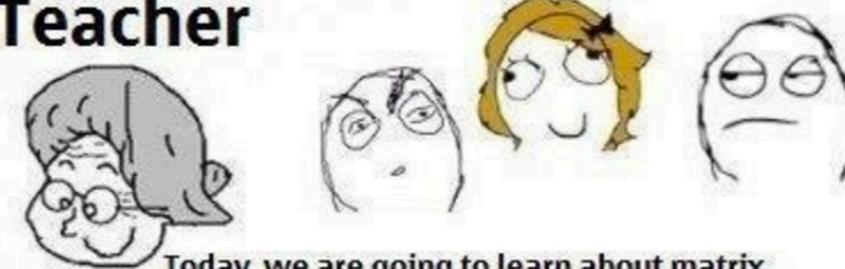
- ❖ Representation
 - ❖ How to represent x, y (and latent factors)
- ❖ Modeling
 - ❖ What **assumptions** we made
 - ❖ i.e., what is the hypothesis set of f ?
- ❖ Algorithms
 - ❖ (learn) Give data, how to learn f ?
 - ❖ (inference) Give test instance x and f , how to evaluate $f(x)$?
- ❖ Learning protocols
 - ❖ What is the goal of the learning algorithm?

Different learning protocols (more technical terms)

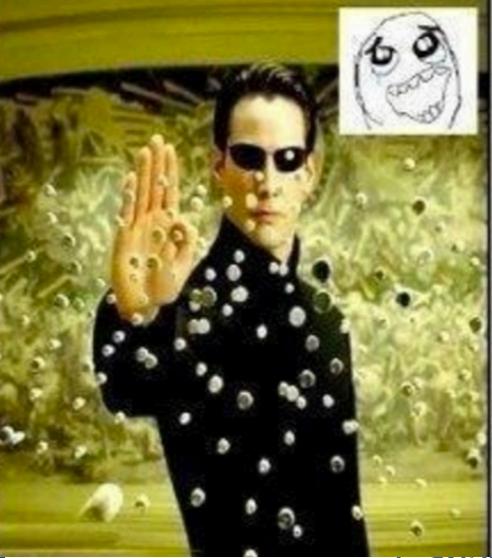
- ❖ Supervision signals?
 - ❖ Supervised learning, semi-supervised learning, unsupervised learning, bandit feedback
- ❖ What to be optimized?
 - ❖ Batch learning: minimize the risk (expected average loss)
 - ❖ Online learning:
 - Receive one sample and make prediction
 - Receive the label; then update
 - minimize the accumulated loss

Linear classification

Teacher



Today, we are going to learn about matrix

Expectation	Reality
	$b = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 1 & 2 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$
	$e = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$

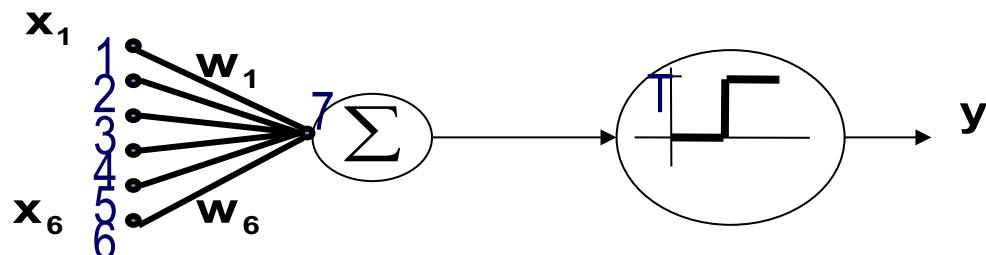


Linear classifiers

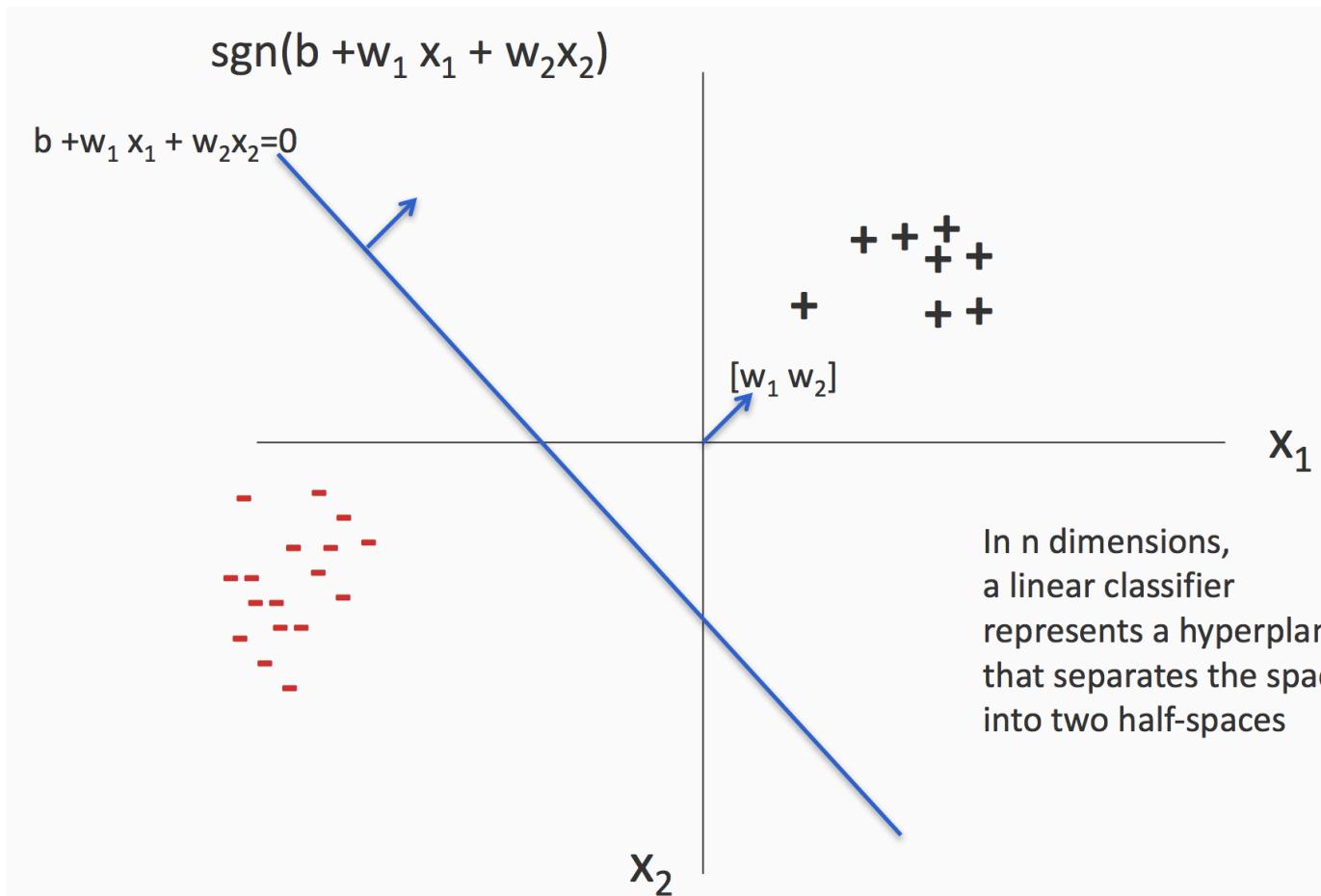
- ❖ For now, we consider binary classification
- ❖ Given a training set $\mathcal{D} = \{(x, y)\}$, find a linear threshold units classify an example x using the classification rule:

$$\text{sgn}(b + \mathbf{w}^T \mathbf{x}) = \text{sgn}(b + \sum_i w_i x_i)$$

- $b + \mathbf{w}^T \mathbf{x} \geq 0 \Rightarrow \text{Predict } y = 1$
- $b + \mathbf{w}^T \mathbf{x} < 0 \Rightarrow \text{Predict } y = -1$



The geometry interpretation



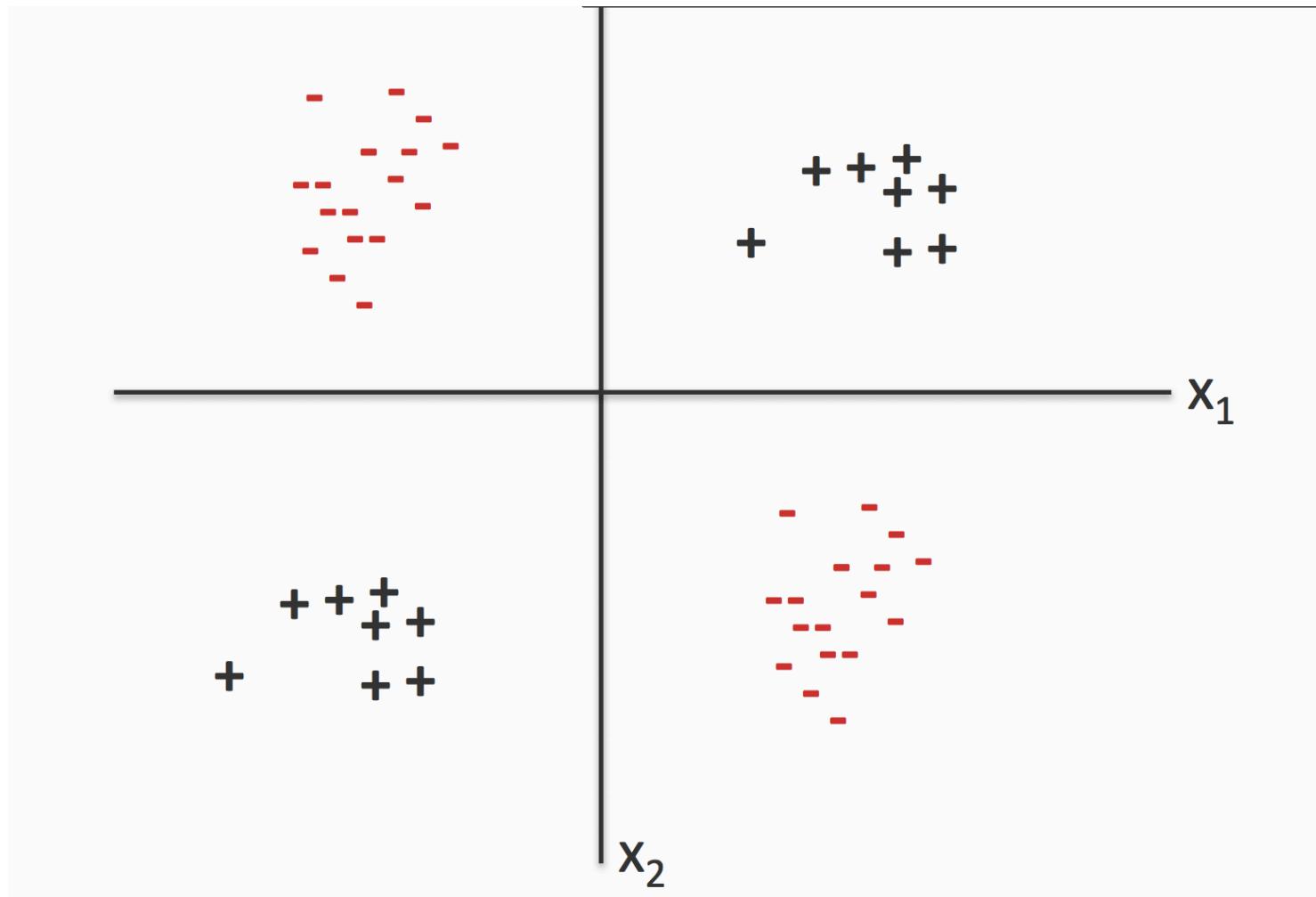
A simple trick to remove the bias term b

$$\begin{aligned} & w^T x + b \\ &= [w^T \ b] \cdot \begin{bmatrix} x \\ 1 \end{bmatrix} \\ &= \tilde{w} \cdot \tilde{x} \end{aligned}$$

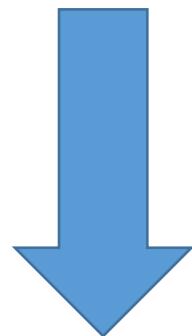
$$\begin{aligned} \tilde{w} &= [w^T \ b]^T \\ \tilde{x} &= [x^T \ 1]^T \end{aligned}$$

For simplicity, I'll write \tilde{w} and \tilde{x} as w and x when there is no confusion

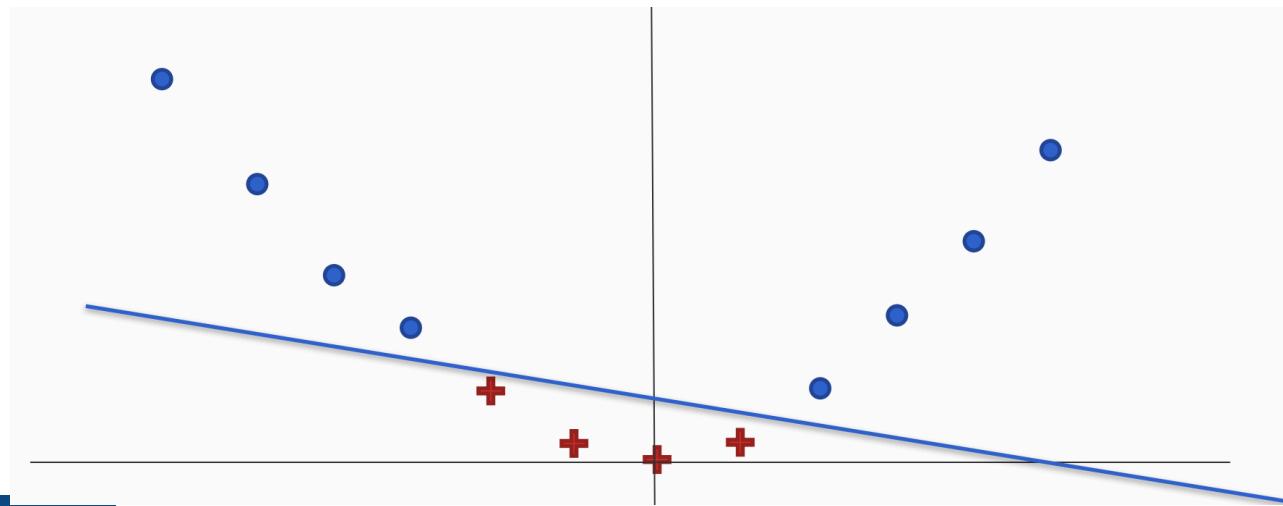
Some data are not linearly separable



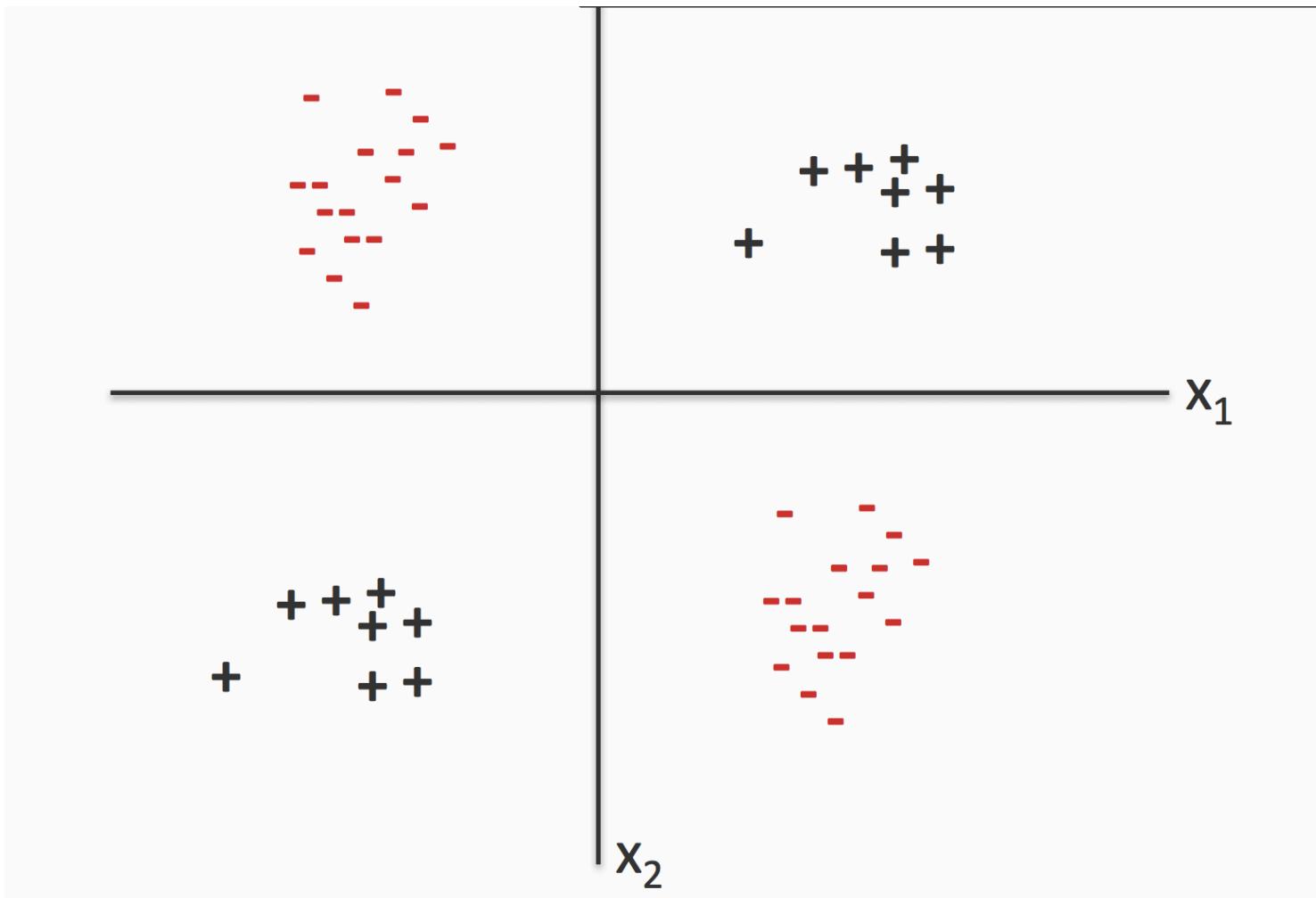
But they can be **made** liner



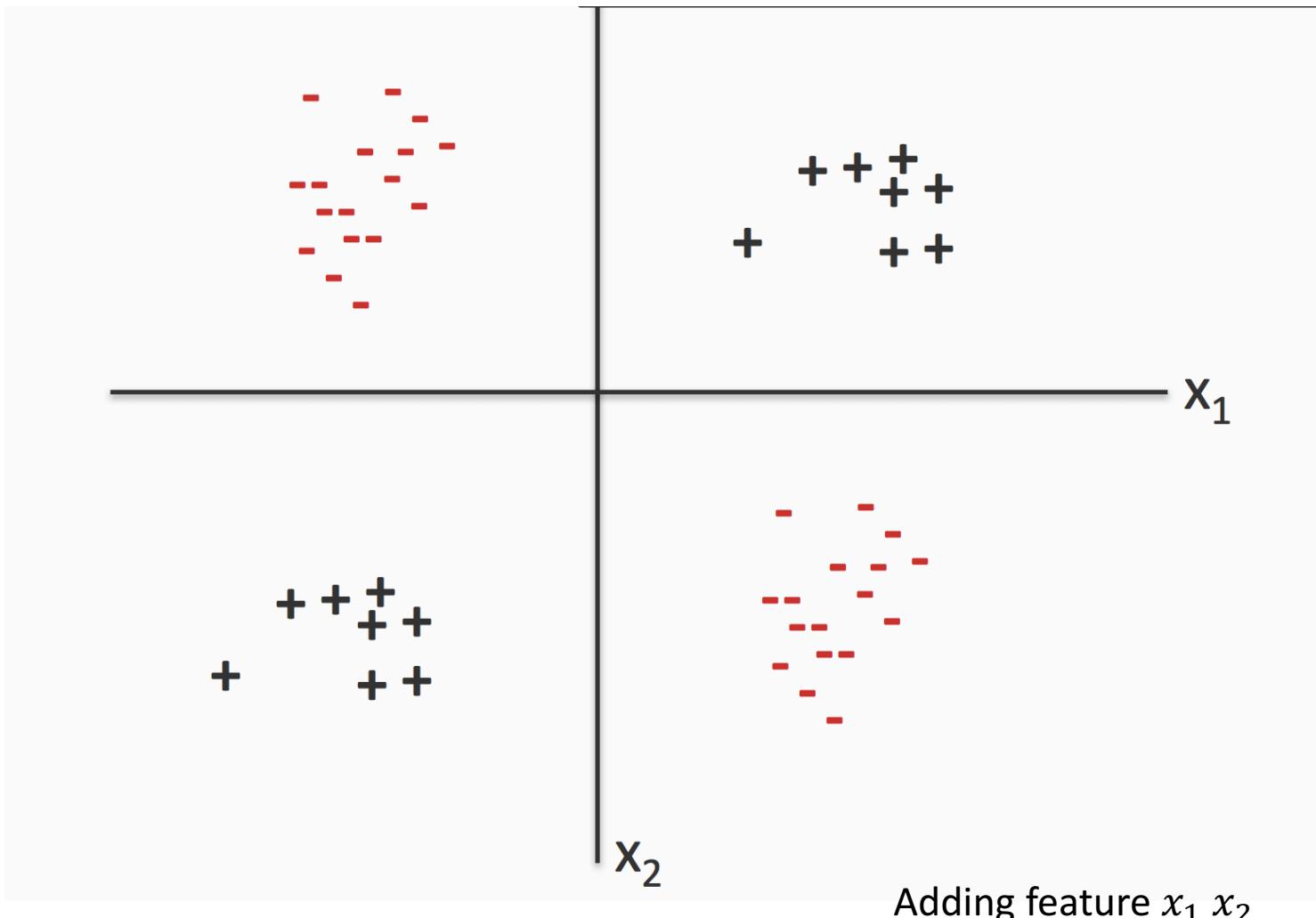
Using a different representation
e.g., feature conjunctions,
non-linear mapping



Exercise: can you make these data points linearly separable?



Exercise: can you make these data points linearly separable?



Linear classifiers

- ❖ Let's take a look at a few linear classifiers
- ❖ We will show later, they can be written in the same framework!

- ❖ Perceptron
- ❖ (Linear) Support Vector Machines
- ❖ Logistic Regression

The Perceptron Algorithm [Rosenblatt 1958]

- ❖ Goal: find a **separating hyperplane**
- ❖ Can be used in an **online** setting:
 - considers one example at a time
- ❖ Converges if data is separable
 - mistake bound

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^n$
 2. For epoch 1 ... T :
 3. For (\mathbf{x}, y) in \mathcal{D} :
 4. $\hat{y} = \text{sgn}(\mathbf{w}^\top \mathbf{x})$ (predict)
 5. if $\hat{y} \neq y$, $\mathbf{w} \leftarrow \mathbf{w} + \eta y \mathbf{x}$ (update)
 6. Return \mathbf{w}

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(\mathbf{w}^\top \mathbf{x}^{\text{test}})$

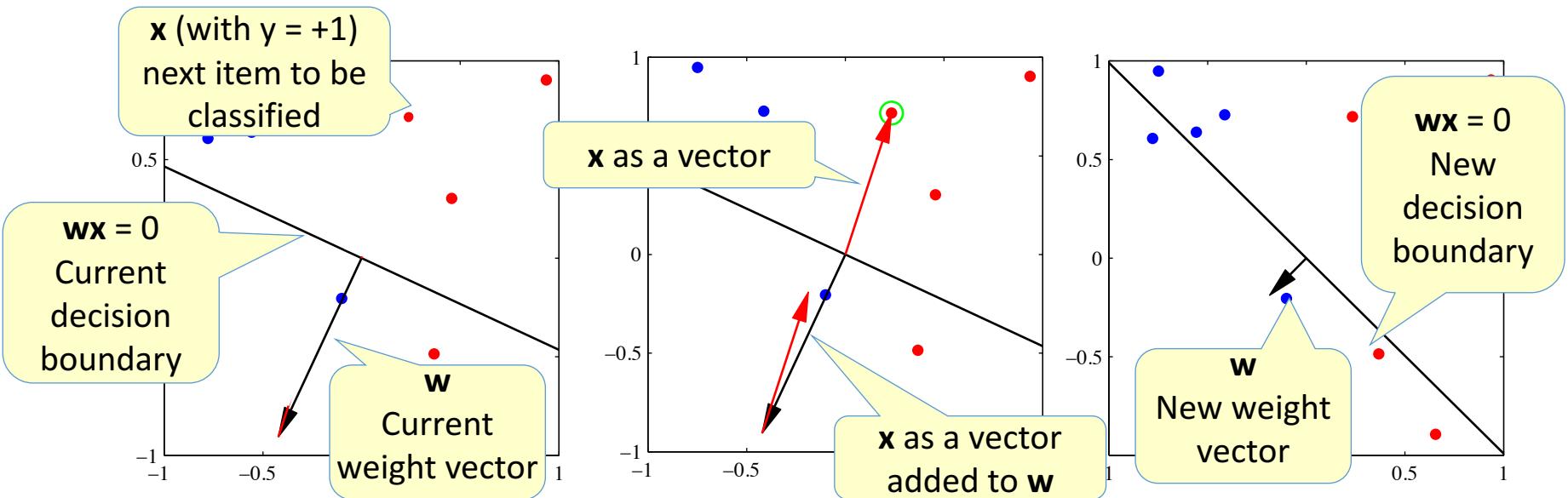
The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. if $y(w^\top x) < 0$
5. $w \leftarrow w + \eta yx$
6. Return w

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(w^\top x^{\text{test}})$

Geometry Interpretation

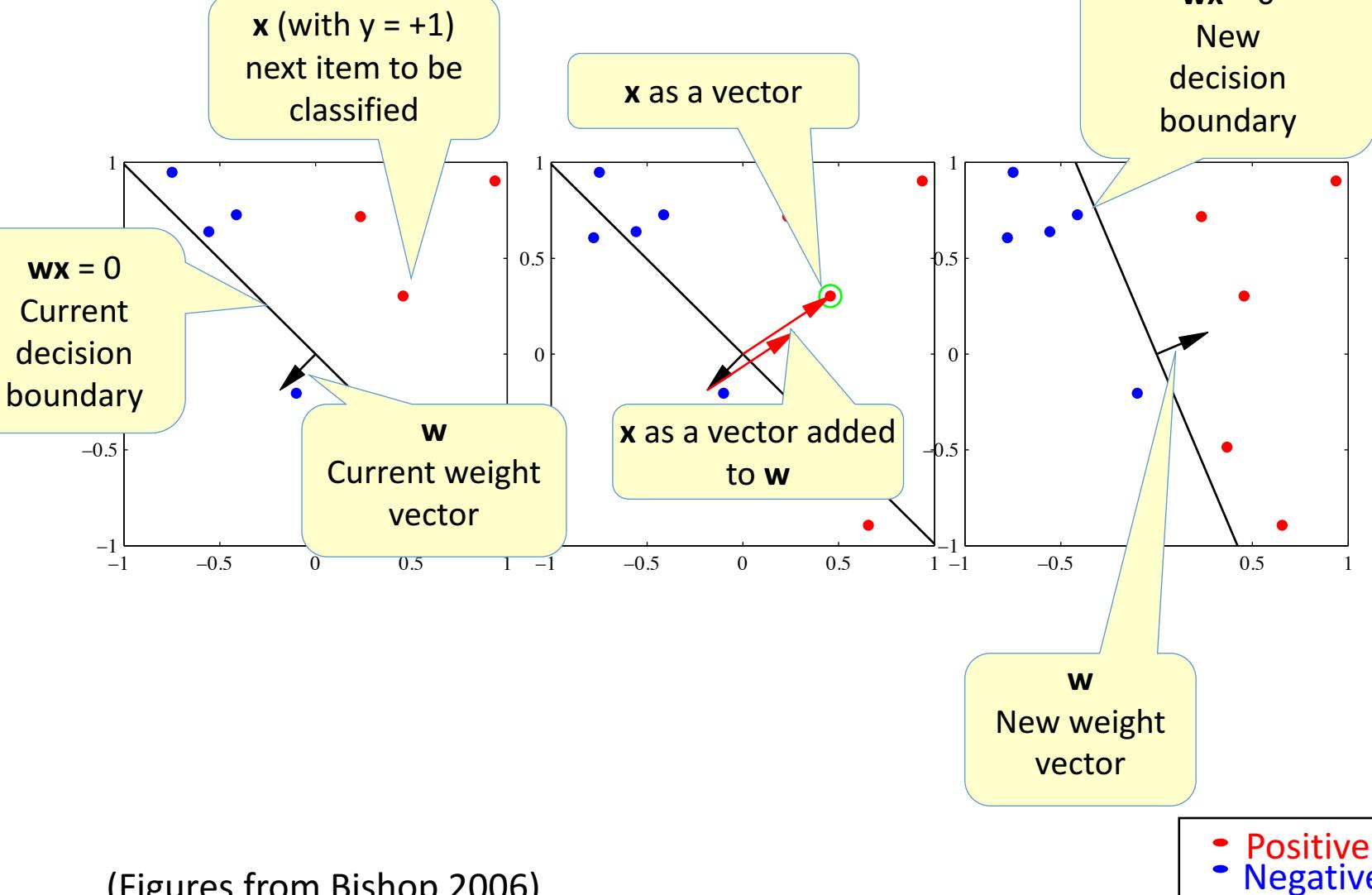


Weight vector points to the positive side

(Figures from Bishop 2006)

- Positive
- Negative

Perceptron in action



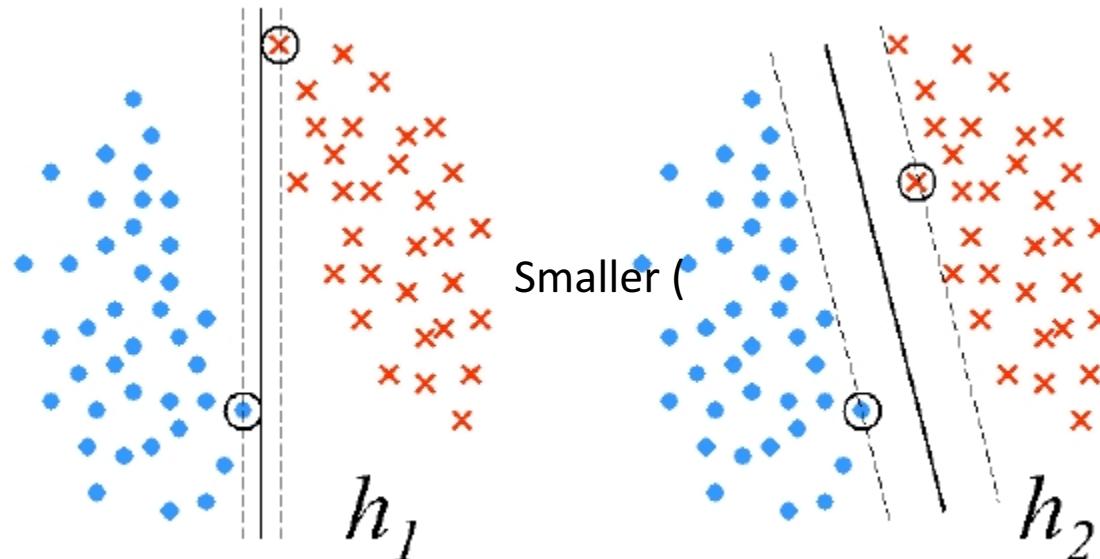
(Figures from Bishop 2006)

Convergence of Perceptron

- ❖ Mistake bound
 - ❖ If data is linearly separable (i.e., a good linear model exists), the perceptron will converge after a fixed number of mistakes [Novikoff 1962]

Marginal Perceptron -- Motivation

- ❖ Which separating hyper-plane is better?



The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(x, y)\}$

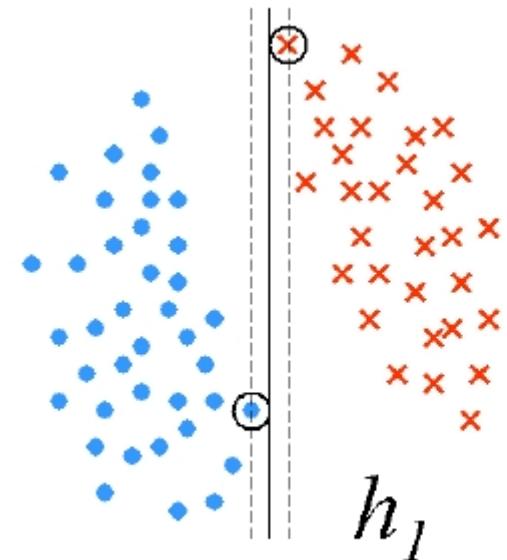
1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. if $y(w^\top x) < 0$
5. $w \leftarrow w + \eta yx$
6. Return w

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(w^\top x^{\text{test}})$

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. if $y(w^\top x) < \delta$
5. $w \leftarrow w + \eta yx$
6. Return w



Prediction: $y^{\text{test}} \leftarrow \text{sgn}(w^\top x^{\text{test}})$

How about batch setting?

- ❖ Learning as loss minimization
- 1. Collect training data $\widehat{D} = \{(x, y)\}$
- 2. Pick a hypothesis class
 - ❖ E.g., linear classifiers, deep neural networks
- 3. Choose a **loss function**
 - ❖ Hinge loss, negative log-likelihood
 - ❖ We can impose a preference (i.e., prior) over hypotheses, e.g., simpler is better
- 4. Minimize the **expected loss**
 - ❖ SGD, coordinate descent, Newton methods, LBFGS

Batch learning setup

- ❖ $\widehat{D} = \{(x, y)\}$ drawn from a fixed, unknown distribution \mathcal{D}
- ❖ A hidden oracle classifier f^* , $y = f^*(x)$
- ❖ We wish to find a hypothesis $f \in H$ that mimics f^*
- ❖ We define a loss function $L(f(x), f^*(x))$ that penalizes mistakes
- ❖ What is the ideal f ?

$$\arg \min_{f \in H} E_{x \sim D} [L(f(x), f^*(x))] \\ \text{expected loss}$$

Batch learning setup

- ❖ $\hat{D} = \{(x, y)\}$ drawn from a fixed, unknown distribution \mathcal{D}
- ❖ A hidden oracle classifier f^* , $y = f^*(x)$
- ❖ We wish to find a hypothesis $f \in H$ that mimics f^*
- ❖ We define a loss function $L(f(x), f^*(x))$ that penalizes mistakes
- ❖ What is the ideal f ?

Let's define

$$L_{0-1}(y, y') = \begin{cases} 1 & \text{if } y \neq y' \\ 0 & \text{if } y = y' \end{cases}$$

$$\min_{f \in H} E_{x \sim D} [L_{0-1}(f(x), f^*(x))] = \min_{f \in H} E_{x \sim D} [\# \text{mistakes}]$$

How can we learn f from \widehat{D}

- ❖ We don't know D , we only see samples in \widehat{D}
- ❖ Instead, we minimize **empirical loss**

$$\min_{f \in H} \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(f(x), y)]$$

How can we prevent over-fitting?

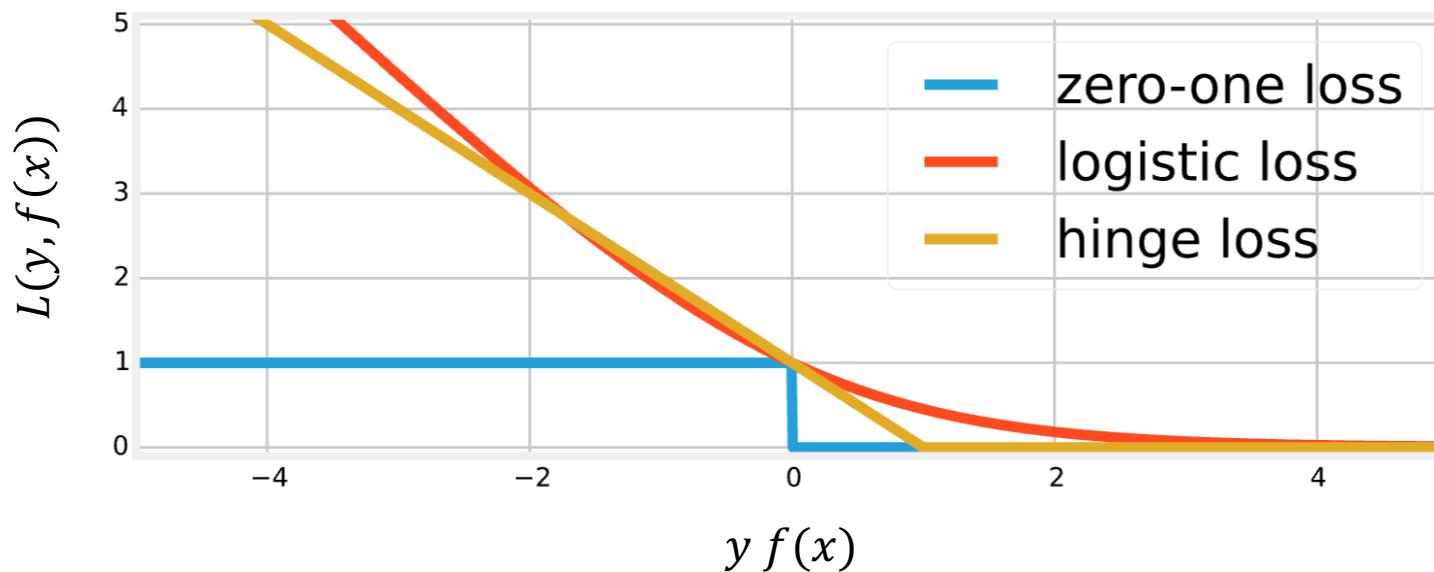
- ❖ With sufficient data, $\widehat{D} \approx D$
- ❖ However, if data is insufficient \Rightarrow overfitting
- ❖ We can impose a preference over models

$$\min_{f \in H} R(f) + \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(f(x), y)]$$

- ❖ We will discuss the choices of $R(f)$ later

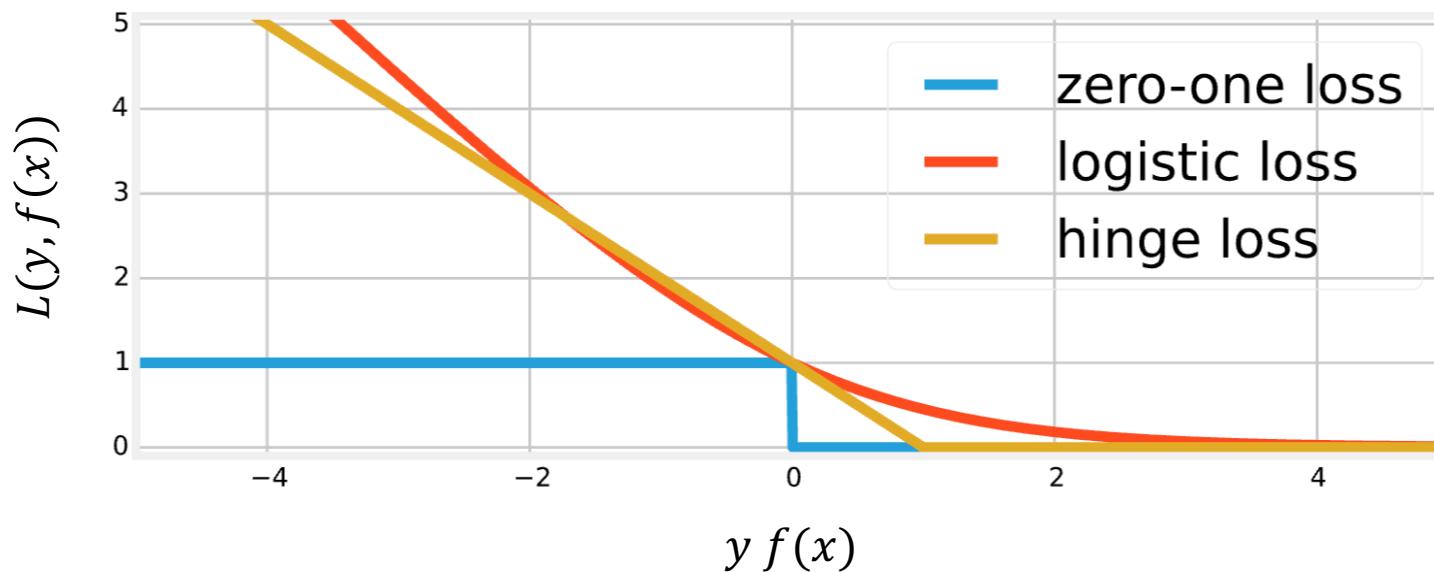
How about the loss function?

- ❖ Usually, we cannot minimize 0-1 loss
 - ❖ It is a combinatorial optimization problem: NP-hard
- ❖ Idea: minimizing its upper-bound



How about the loss function?

- ❖ Usually, we cannot minimize it
 - ❖ It is a combinatorial optimization problem
- ❖ Idea: minimizing its upper-bound



Many choices

- ❖ We are minimizing with R, L, H with your choice

$$\min_{f \in H} R(f) + \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(f(x), y)]$$

- ❖ Let consider H is a set of d-dimensional linear function
- ❖ H can be parameterized as
 $\{f(x): w^T x \geq 0\}, w \in R^d$



Back to Linear model

- ❖ We are minimizing with R, L, H with your choice

$$\min_{f \in H} R(f) + \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(f(x), y)]$$

- ❖ Let decide H to be a set of d-dimensional linear function

- ❖ H can be parameterized as

$$\{f(x): w^T x \geq 0\}, w \in R^d$$

- ❖ We are going to fine the best one based on \widehat{D}
 - ❖ i.e., find the best setting of w and b



Rewrite our optimization problem

- ❖ Minimizing the empirical loss:

$$\min_{f \in H} R(f) + \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(f(x), y)]$$

- ❖ Minimizing the empirical loss **with linear function**

$$\min_{w \in R^d} R(w) + \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(x, w, y)]$$

- ❖ What choices of R and L we have?

Many choices of loss function (L)

Many loss functions exist

- Perceptron loss

$$L_{\text{Perceptron}}(y, \mathbf{x}, \mathbf{w}) = \max(0, -y\mathbf{w}^T \mathbf{x})$$

- Hinge loss (SVM)

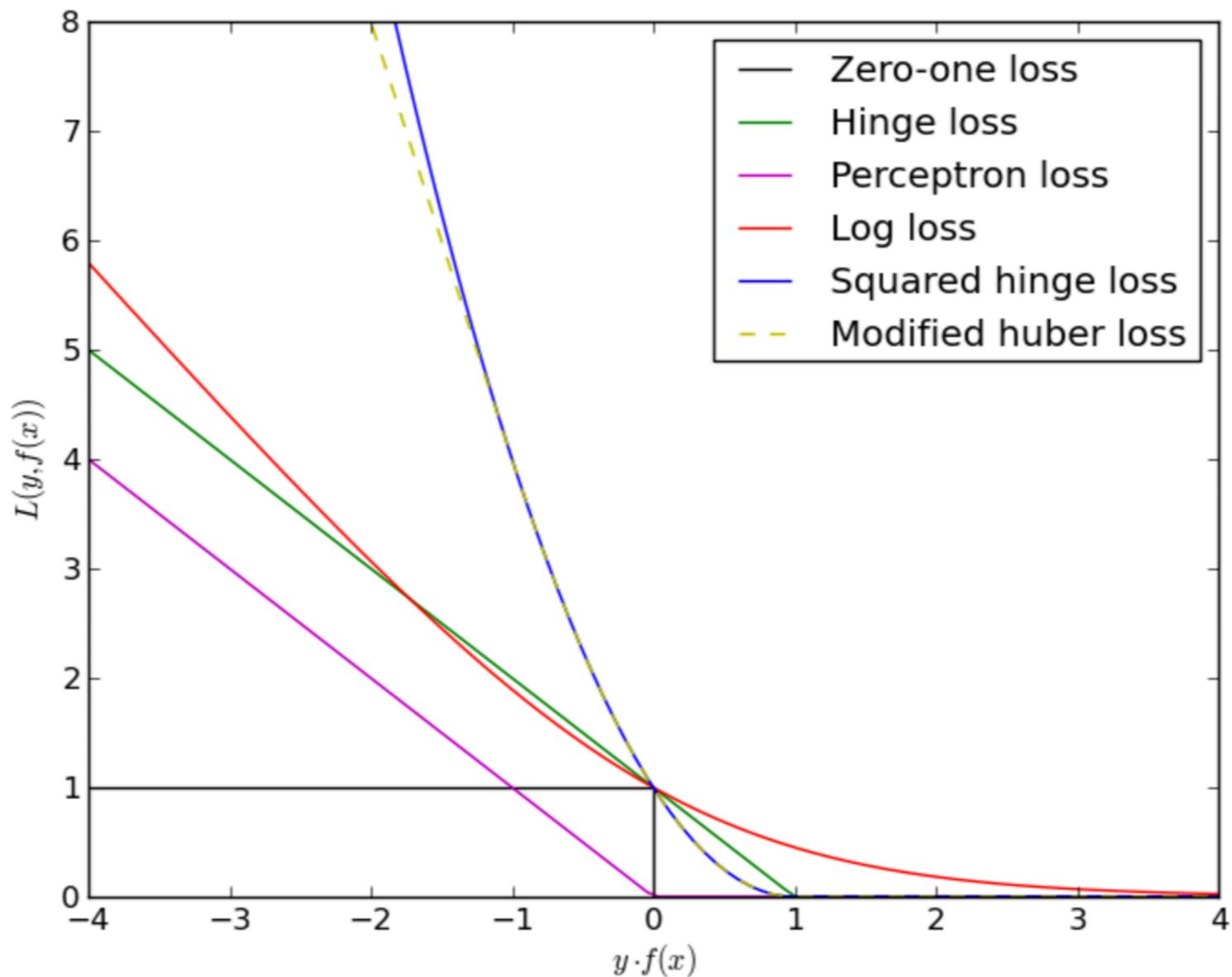
$$L_{\text{Hinge}}(y, \mathbf{x}, \mathbf{w}) = \max(0, 1 - y\mathbf{w}^T \mathbf{x})$$

- Exponential loss (AdaBoost)

$$L_{\text{Exponential}}(y, \mathbf{x}, \mathbf{w}) = e^{-y\mathbf{w}^T \mathbf{x}}$$

- Logistic loss (logistic regression)

$$L_{\text{Logistic}}(y, \mathbf{x}, \mathbf{w}) = \log(1 + e^{-y\mathbf{w}^T \mathbf{x}})$$



Many choices of $R(w)$

- ❖ Minimizing the empirical loss with linear function

$$\min_{w \in R^d} R(\mathbf{w}) + \frac{1}{|\widehat{D}|} \sum_{(x,y) \in \widehat{D}} [L(\mathbf{x}, \mathbf{w}, y)]$$

- ❖ Prefer simpler model: (how?)

- ❖ Sparse:

$R(\mathbf{w}) = \#\text{non-zero elements in } w$ (L0 regularizer)

$R(\mathbf{w}) = \sum_i |w_i|$ (L1 regularizer)

- ❖ Gaussian prior (large margin w/ hinge loss):

$R(\mathbf{w}) = \sum_i w_i^2 = \mathbf{w}^T \mathbf{w}$ (L2 regularizer)

Support Vector Machines



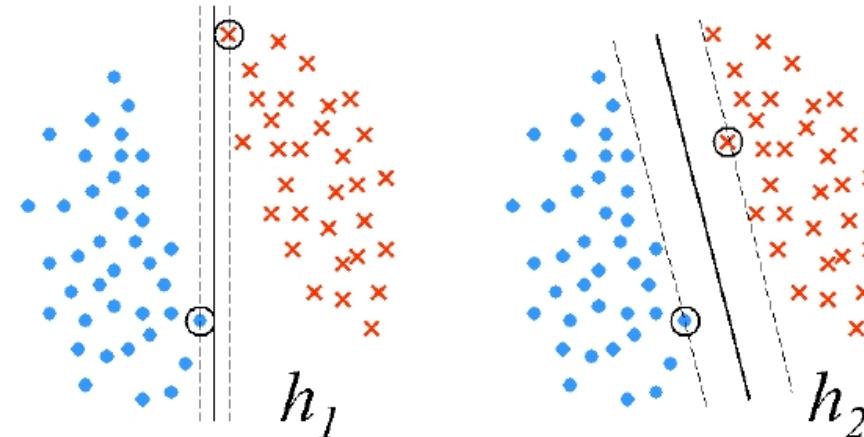
CMU ML protest

Support Vector Machines (SVMs)

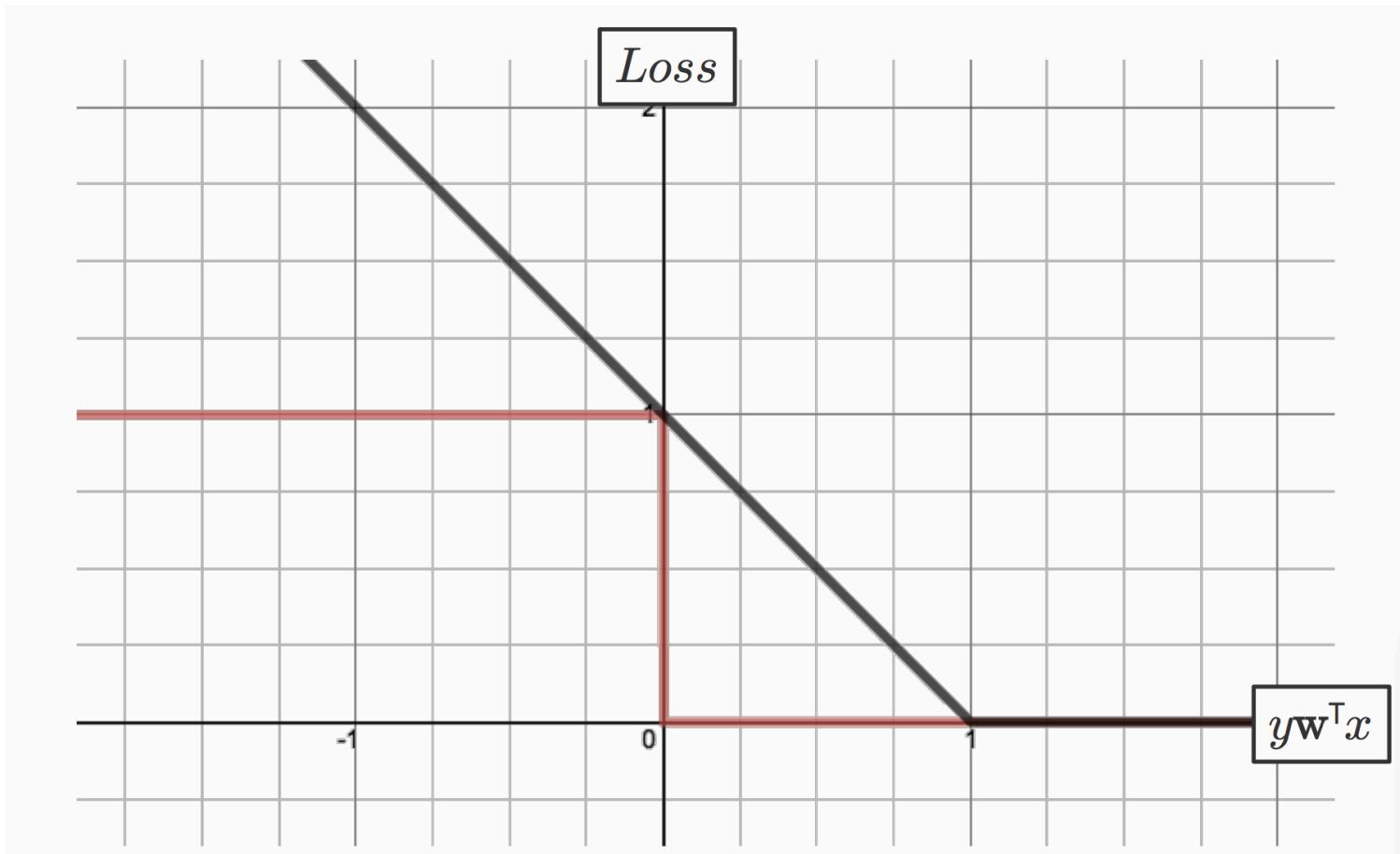
- ❖ $R(w)$: l2-loss, $L(w, x, y)$: hinge loss

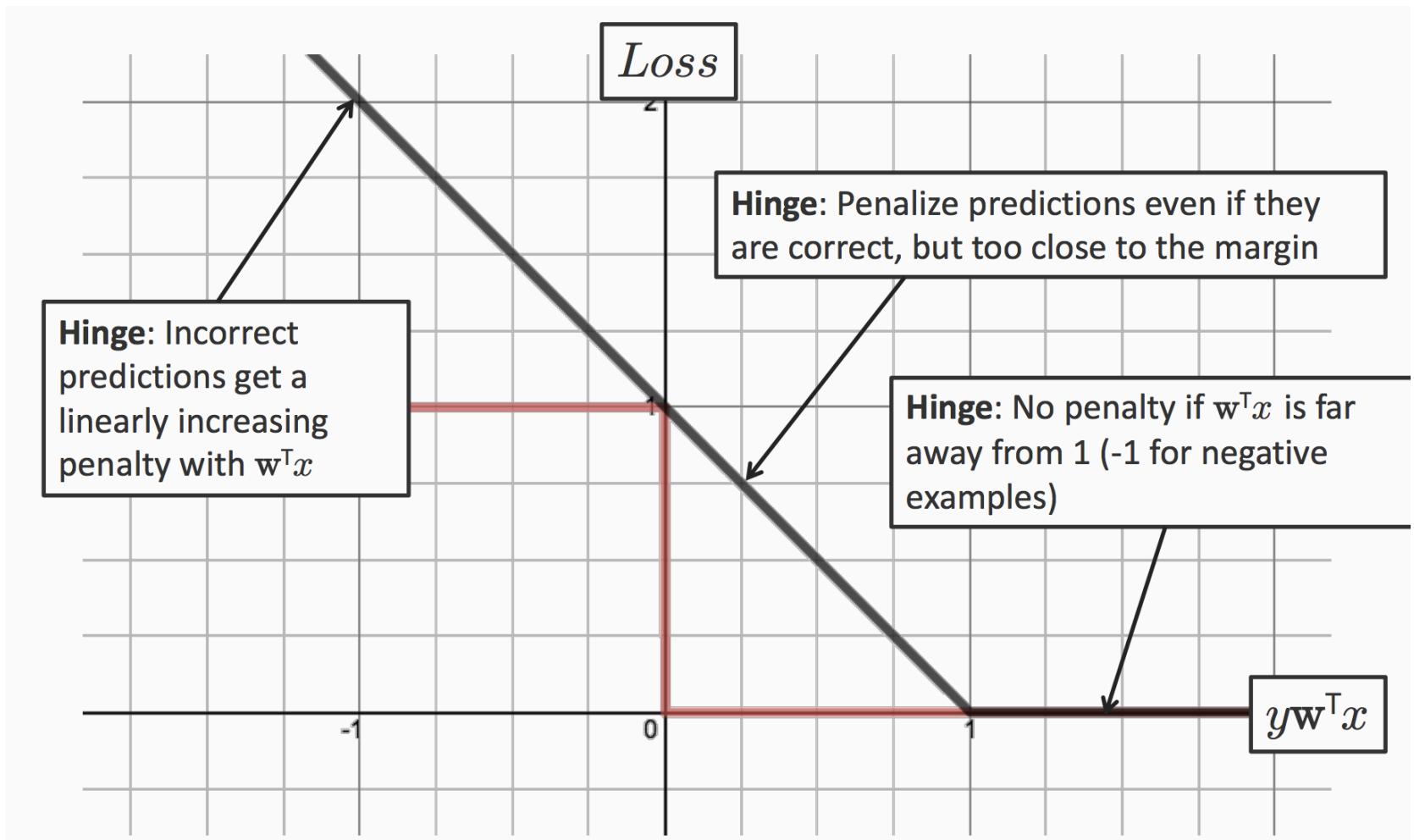
$$\min_w \frac{1}{2} w^T w + C \sum_i \max(0, 1 - y_i(w^T x_i))$$

- ❖ Maximizing margin (why?!!)



The hinge loss



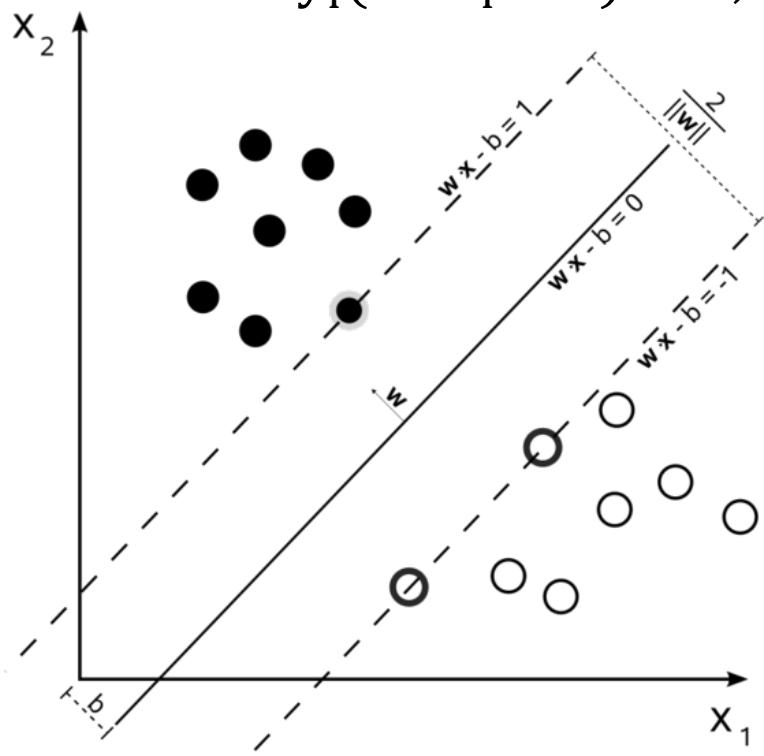


Let's view it from another direction

- ❖ SVM learns a model \mathbf{w} on $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ by solving:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{s.t. } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall (\mathbf{x}_i, y_i) \in \mathcal{D}$$



(Hard SVM)

Why the margin is $\frac{2}{\|\mathbf{w}\|}$?

$$\begin{aligned} & \arg \max \frac{2}{\|\mathbf{w}\|} \\ &= \arg \min \|\mathbf{w}\| \\ &= \arg \min \|\mathbf{w}\|^2 \\ &= \arg \min \mathbf{w}^T \mathbf{w} \end{aligned}$$

Soft SVMs

- ❖ Data is not separable \Rightarrow hard SVM fails
Why?
- ❖ Introduce a set of slack variable $\{\xi_i\}$
 \Rightarrow relax the constraints
- ❖ Given $\mathcal{D} = \{(x_i, y_i)\}$, soft SVM solves:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i \quad (\text{Soft SVM})$$

$$\text{s. t } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i; \quad \xi_i \geq 0 \quad \forall i$$

penalty parameter

An alternative formulation

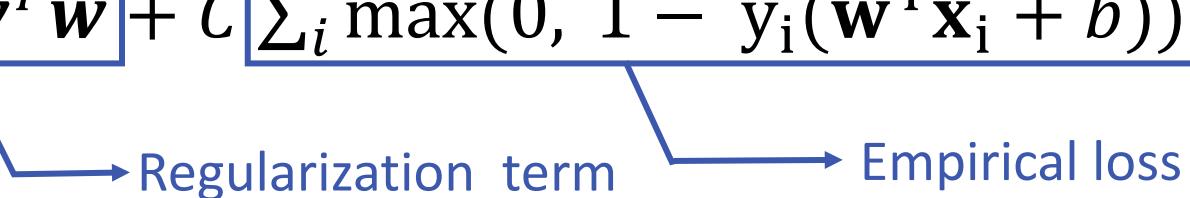
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i; \quad \xi_i \geq 0 \quad \forall i \end{aligned}$$

- ❖ Rewrite the constraints:

$$\xi_i \geq 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b); \quad \xi_i \geq 0 \quad \forall i$$

- ❖ In the optimum, $\xi_i = \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))$
- ❖ Soft SVM can be rewritten as:

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \boxed{\mathbf{w}^T \mathbf{w}} + C \boxed{\sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))}$$

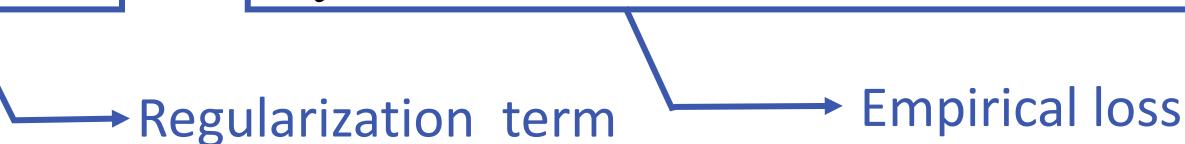

Regularization term Empirical loss

An alternative formulation

$$\min_w \frac{1}{2} w^T w + C \sum_i \xi_i$$

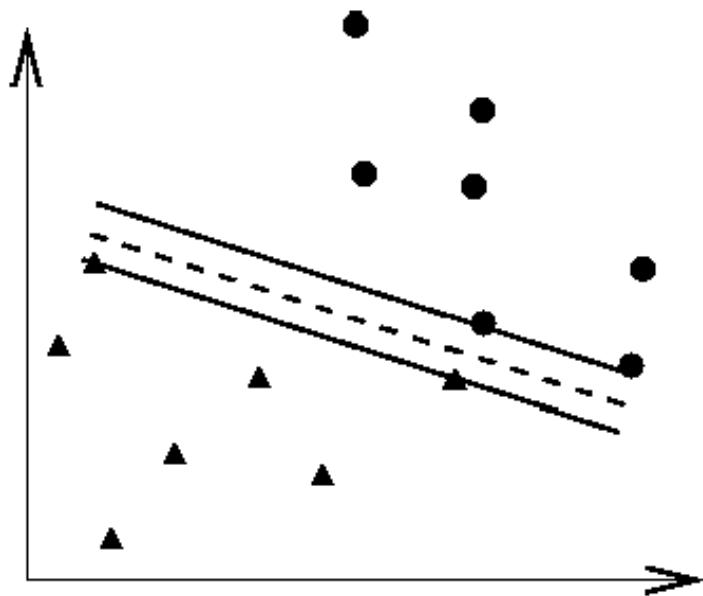
- ❖ We can simply "b" using the trick
- ❖ However, we will add b into the regularization term
- ❖ It is often okay, if we have many features

$$\min_{w,b} \frac{1}{2} \boxed{w^T w} + C \boxed{\sum_i \max(0, 1 - y_i(w^T x_i + b))}$$

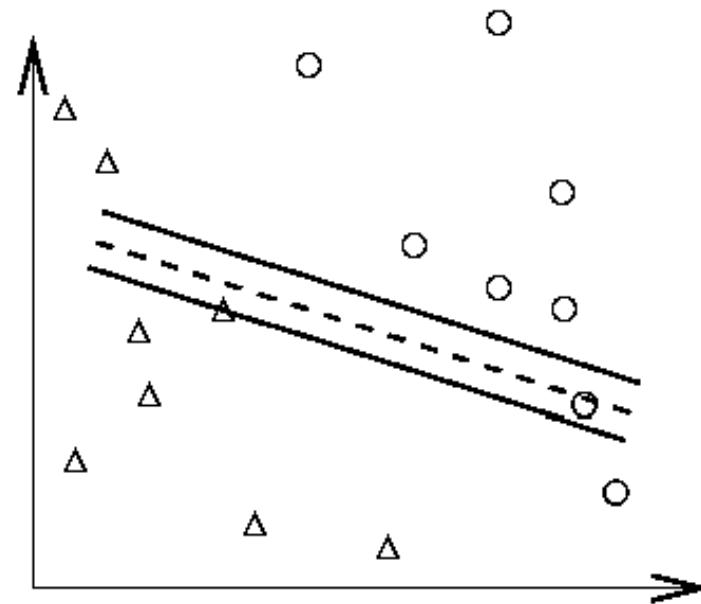
An arrow points from the term $w^T w$ inside the first box to the label "Regularization term". Another arrow points from the term $\max(0, 1 - y_i(w^T x_i + b))$ inside the second box to the label "Empirical loss".

Regularization term Empirical loss

Balance between regularization and empirical loss

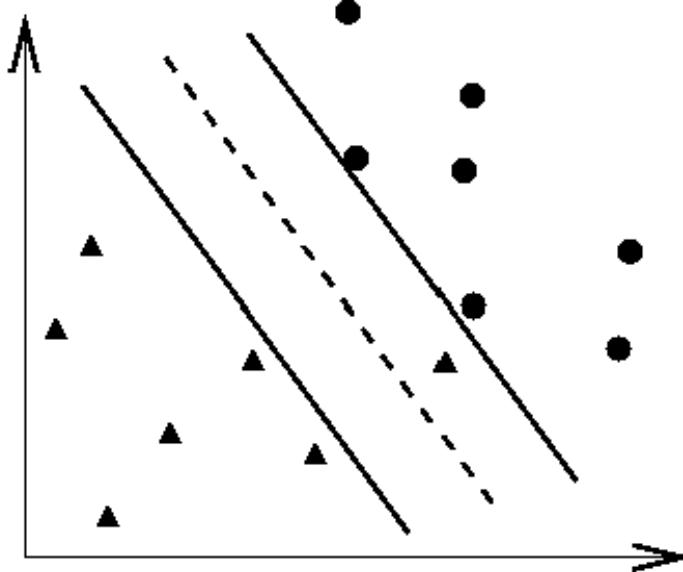


(a) Training data and an over-fitting classifier

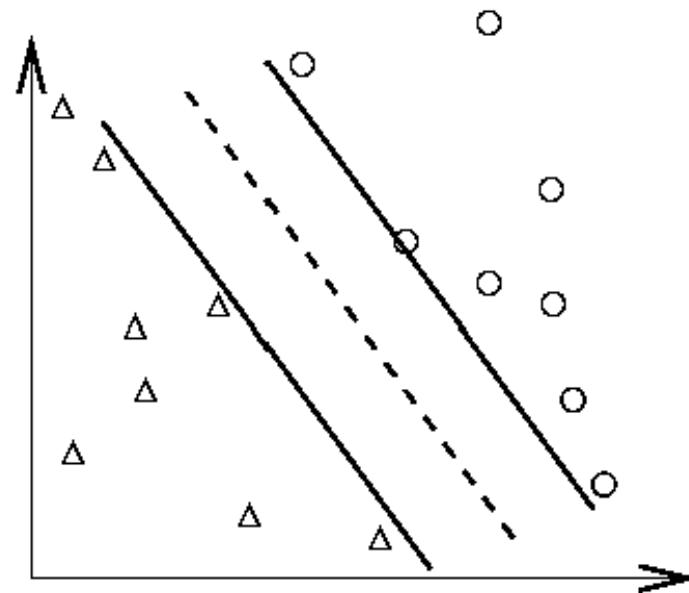


(b) Testing data and an over-fitting classifier

Balance between regularization and empirical loss



(c) Training data and a better classifier



(d) Testing data and a better classifier

[\(DEMO\)](#)

Regularized loss minimization

- ❖ L1-Loss SVM

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i))$$

- ❖ L2-Loss SVM

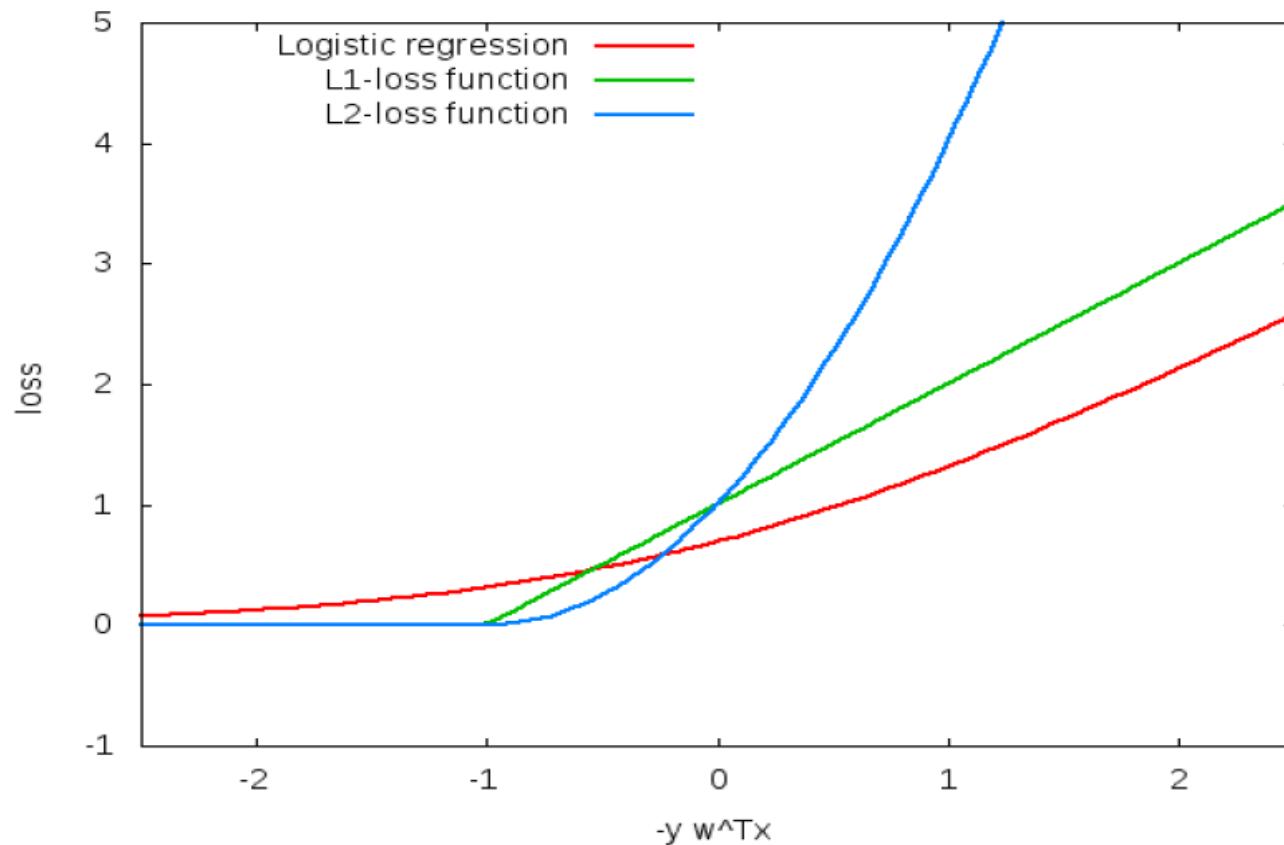
$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i))^2$$

- ❖ Logistic Regression (regularized)

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \log(1 + e^{-y_i(\mathbf{w}^T \mathbf{x}_i)})$$

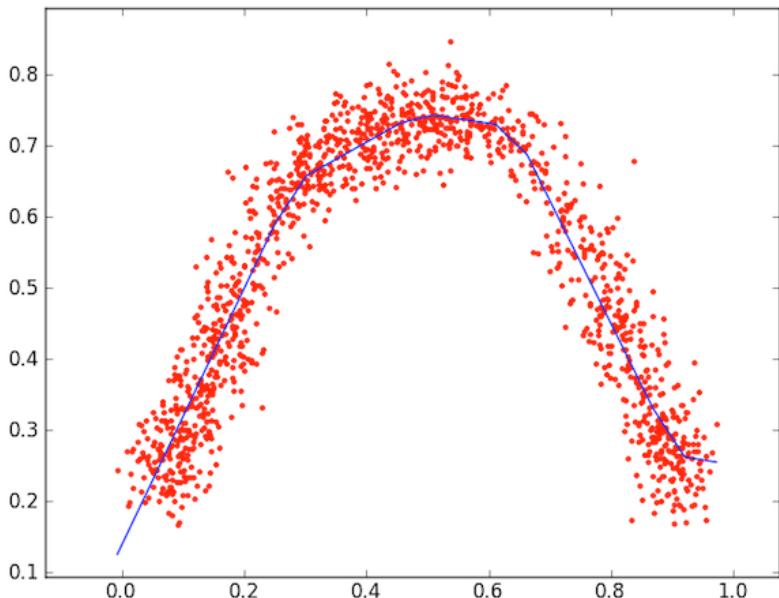
- ❖ Loss over training data + regularizer

Loss Functions

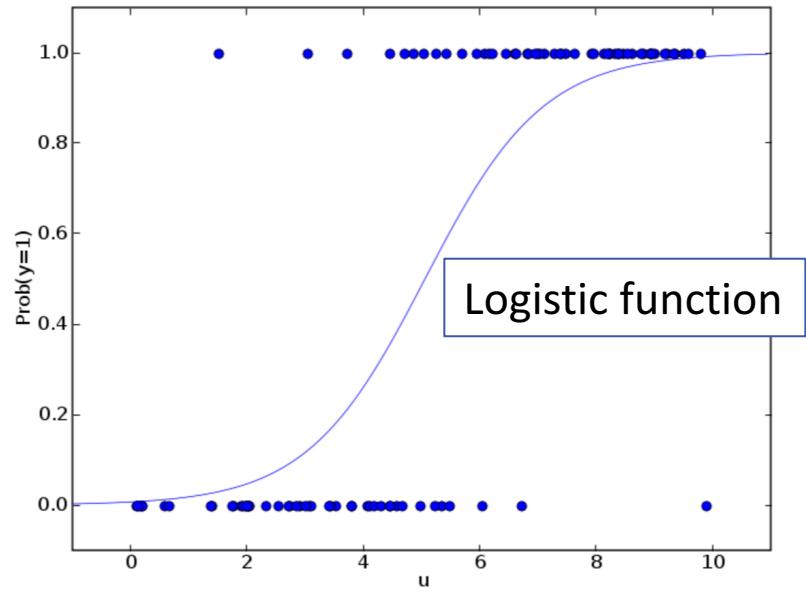


Logistic Regression

Regression



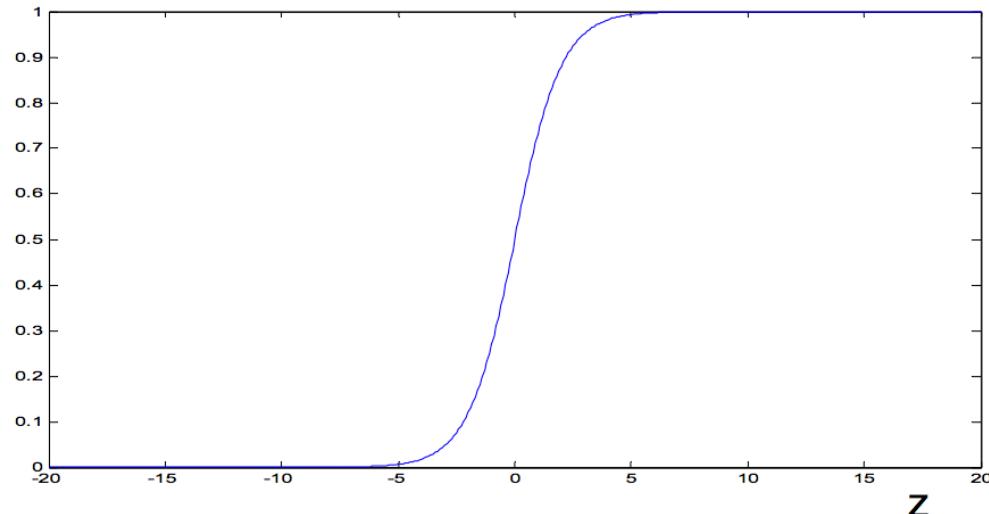
Logistic regression



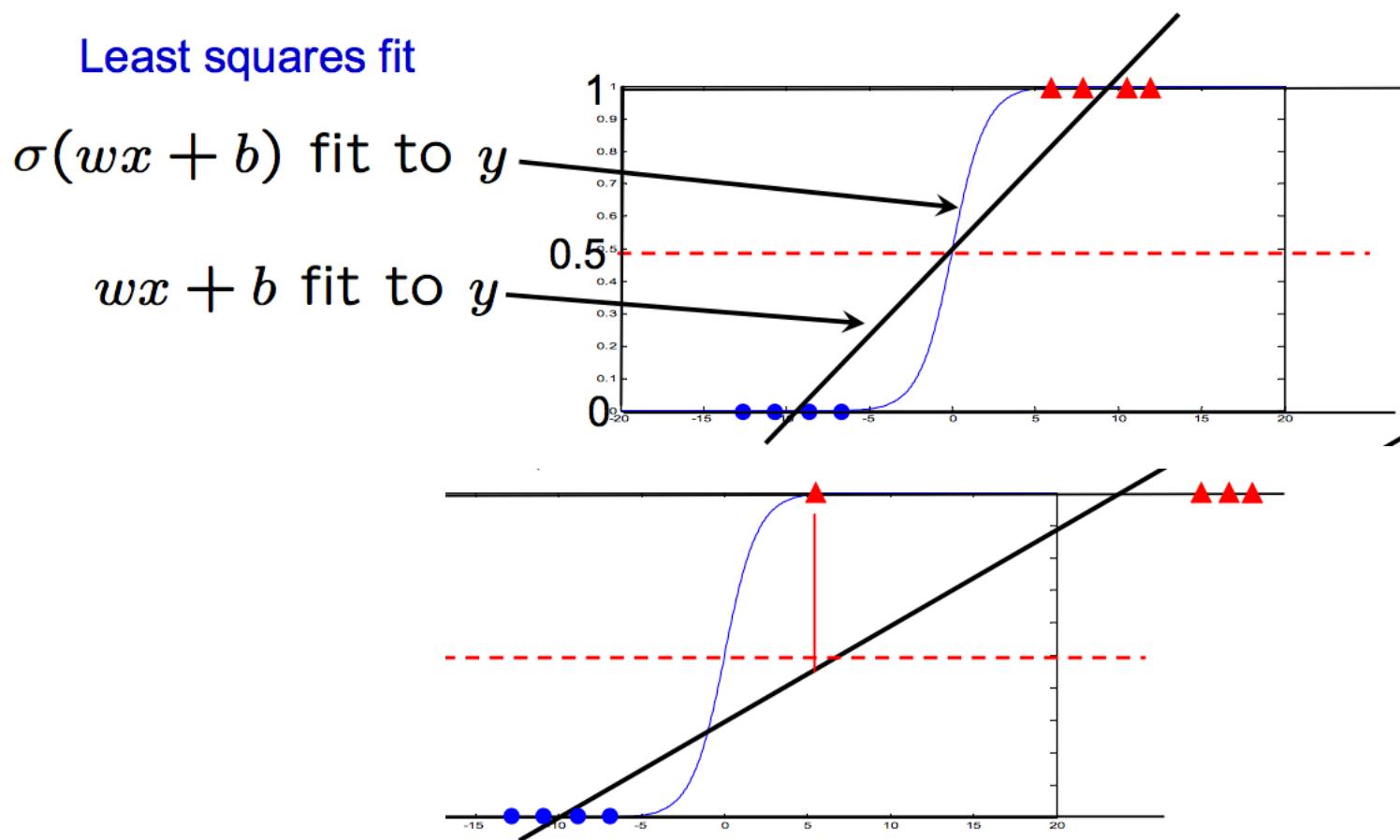
logistic function / sigmoid function

- ❖ When $z \rightarrow \infty$ what is $\sigma(z)$?
- ❖ When $z \rightarrow -\infty$ what is $\sigma(z)$?
- ❖ When $z = 0$ what is $\sigma(z)$?

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Why sigmoid?



Probabilistic Interpretation

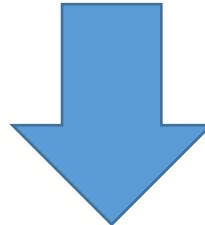
$$\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \log(1 + e^{-y_i(\mathbf{w}^T \mathbf{x}_i)})$$

Assume labels are generated using the following probability distribution:

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$
$$P(y = -1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$P(y = -1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$



$$P(y | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-y \mathbf{w}^T \mathbf{x})}$$

How to make prediction?

Predict $y=1$ if $P(y=1 | \mathbf{x}, \mathbf{w}) > P(y=-1 | \mathbf{x}, \mathbf{w})$

Decision boundary?

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$P(y = -1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$

$$\log \frac{P(Y = 1 | x)}{P(Y = -1 | x)} = \mathbf{w}^T \mathbf{x} + b.$$

❖ The decision boundary?

$$\mathbf{w}^T \mathbf{x} + b = 0.$$

Alternative view

- ❖ Predict $y=1$ if $P(y=1|x, w) > P(y=-1|x, w)$

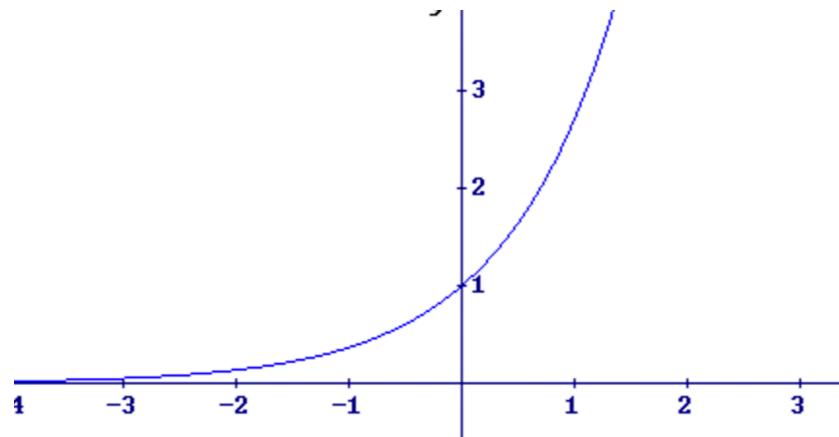
❖ Predict $y=1$ if $P(y=1|x, w) > P(y=-1|x, w)$

❖ When does this happen?

$$\begin{aligned} \frac{1}{1+\exp(-w^T x)} &> 0.5 \\ \Rightarrow 1 + \exp(-w^T x) &< 2 \\ \Rightarrow \exp(-w^T x) &< 1 \\ \Rightarrow w^T x &> 0 \end{aligned}$$

- ❖ When does this happen?

$$\begin{aligned} \frac{1}{1+\exp(-w^T x)} &> 0.5 \\ \Rightarrow 1 + \exp(-w^T x) &< 2 \\ \Rightarrow \exp(-w^T x) &< 1 \\ \Rightarrow w^T x &> 0 \end{aligned}$$



Maximum likelihood estimation

- ❖ Probabilistic model assumption:

$$P(y|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-y\mathbf{w}^T \mathbf{x})}$$

- ❖ The log-likelihood of seeing a dataset $D = \{(x_i, y_i)\}$ if the true weight vector was w :

$$\log P(D|\mathbf{w}) = - \sum \log (1 + \exp(-y\mathbf{w}^T \mathbf{x}))$$

$$\begin{aligned} P(D|w) &= \prod_i P(y_i|x_i, w) \\ \Rightarrow \log P(D|w) &= \sum_i \log P(y_i|x_i, w) \end{aligned}$$

Minimizing negative log-likelihood

- ❖ Log likelihood

$$\log P(D|\mathbf{w}) = - \sum \log (1 + \exp(-y\mathbf{w}^T \mathbf{x}))$$

- ❖ Logistic regression

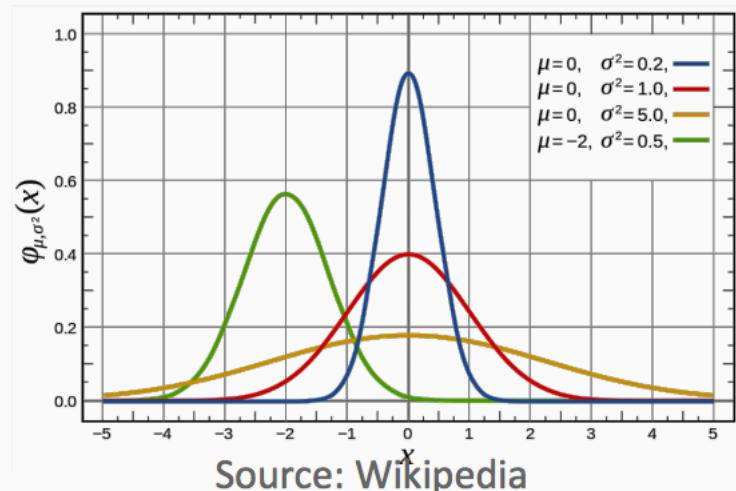
$$\min_{\mathbf{w}, b} \sum_i \log(1 + e^{-y_i(\mathbf{w}^T \mathbf{x}_i)})$$

- ❖ Let's add some prior
 - ❖ Simpler is better \Rightarrow add Gaussian Prior

Add Gaussian Prior

- ❖ Simpler is better \Rightarrow add Gaussian Prior
- ❖ Suppose each element in w is drawn independently from the normal distribution centered at zero with variance σ^2
- ❖ Bias towards smaller weights

$$P(w_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{w_i^2}{2\sigma^2}\right)$$



Regularized Logistic regression

$$P(w_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{w_i^2}{2\sigma^2}\right)$$

- ❖ Remember we are in the log space

$$\log P(\mathbf{w}) = -\frac{1}{2\sigma^2} \mathbf{w}^T \mathbf{w} + \text{constant terms}$$

- ❖ Put them together
 - ❖ $P(w|D) \propto P(w, D) = P(D | w)P(w)$
 - ❖ Learning:
Find weight vector by maximizing the posterior distribution $P(w | D)$

Maximum a posteriori estimation

- ❖ Put them together
 - ❖ $P(w|D) \propto P(w, D) = P(D | w)P(w)$
- ❖ Learning: Find weight vector by maximizing the posterior distribution $P(w | D)$

$$\log P(D, \mathbf{w}) = -\frac{1}{2\sigma^2} \mathbf{w}^T \mathbf{w} - \sum_i \log (1 + \exp(-y \mathbf{w}^T \mathbf{x}))$$

prior

Log-likelihood

$$\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \log(1 + e^{-y_i(\mathbf{w}^T \mathbf{x}_i)})$$

Regularized loss minimization

❖ L1-Loss SVM

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i))$$

❖ L2-Loss SVM

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i))^2$$

❖ Logistic Regression (regularized)

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \log(1 + e^{-y_i(\mathbf{w}^T \mathbf{x}_i)})$$

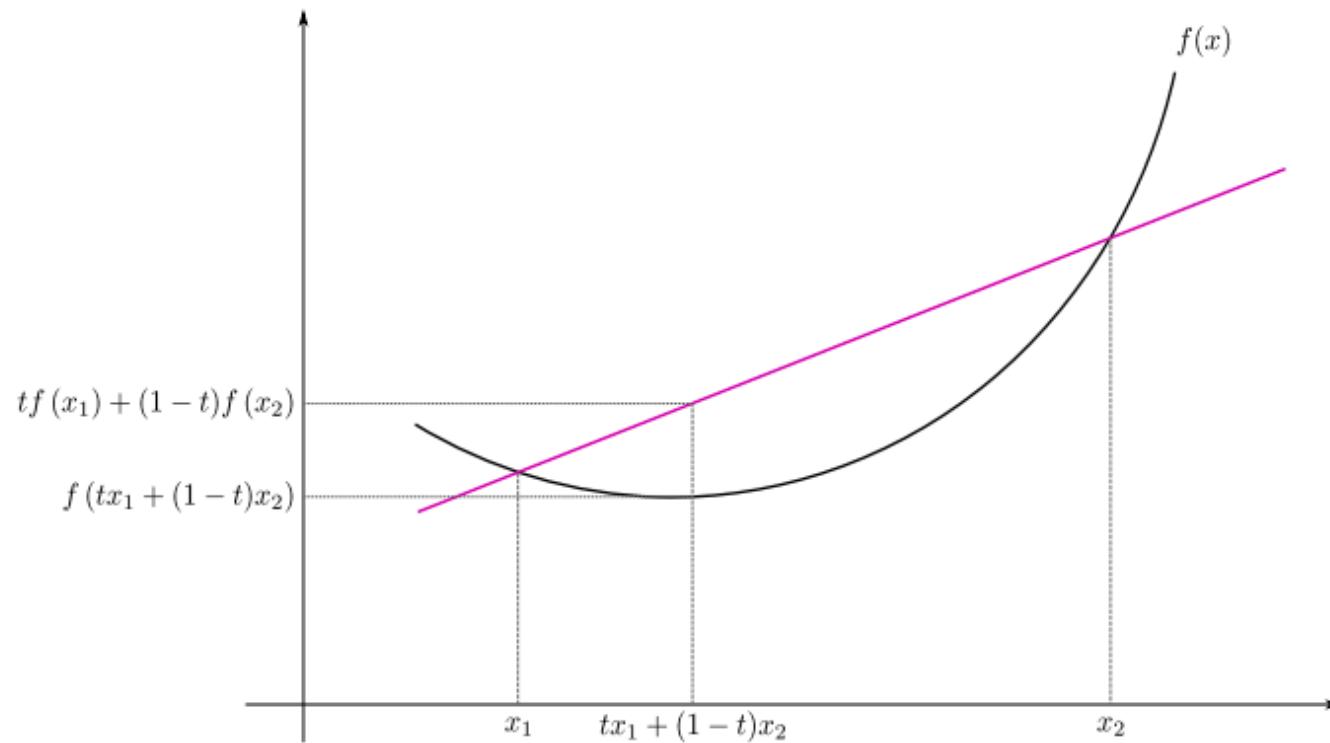
❖ Loss over training data + regularizer

How to learn (how to optimize the objective function?)

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i))$$

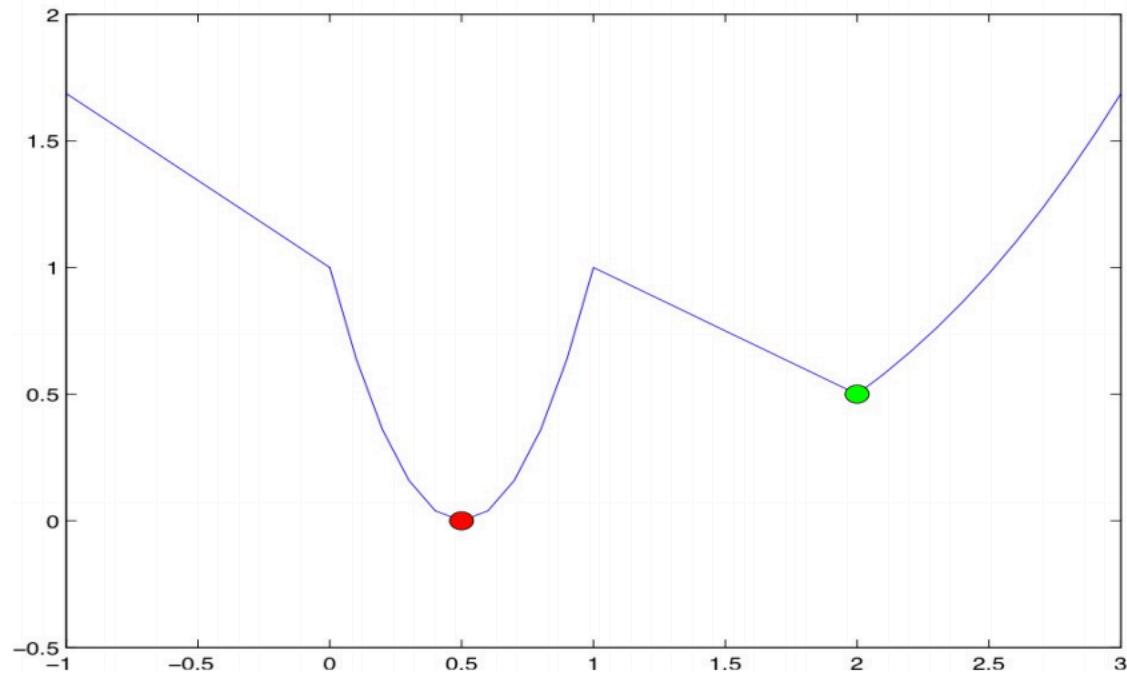
- ❖ This function is **convex**
- ❖ Many convex optimization methods can be used
 - ❖ Stochastic (sub)-gradient descent
 - ❖ Coordinate descent methods
 - ❖ Newton methods
 - ❖ LBFGS

Convexity



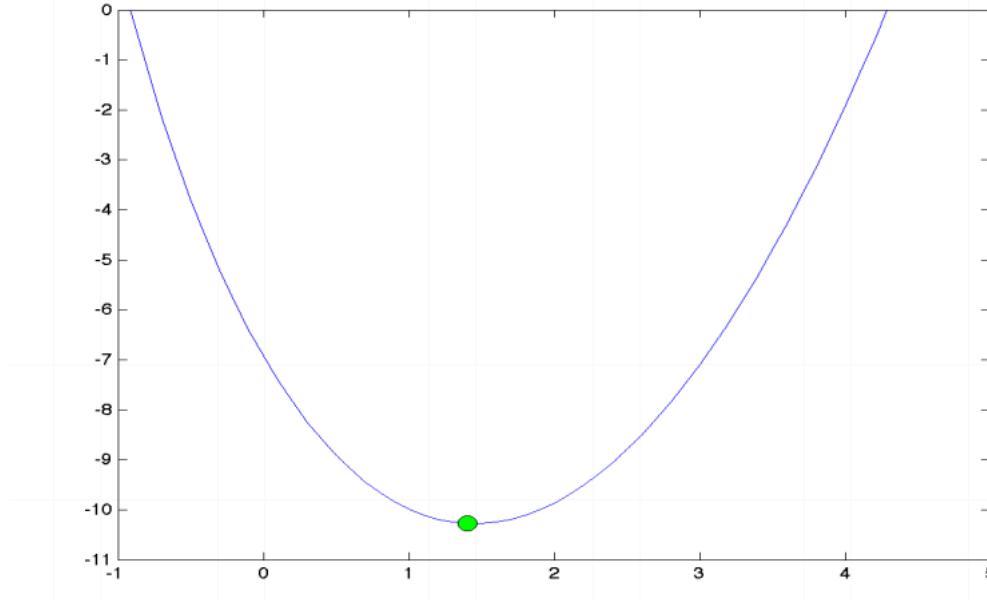
Non convex minimization is hard

- ❖ You may end up with some local minimum

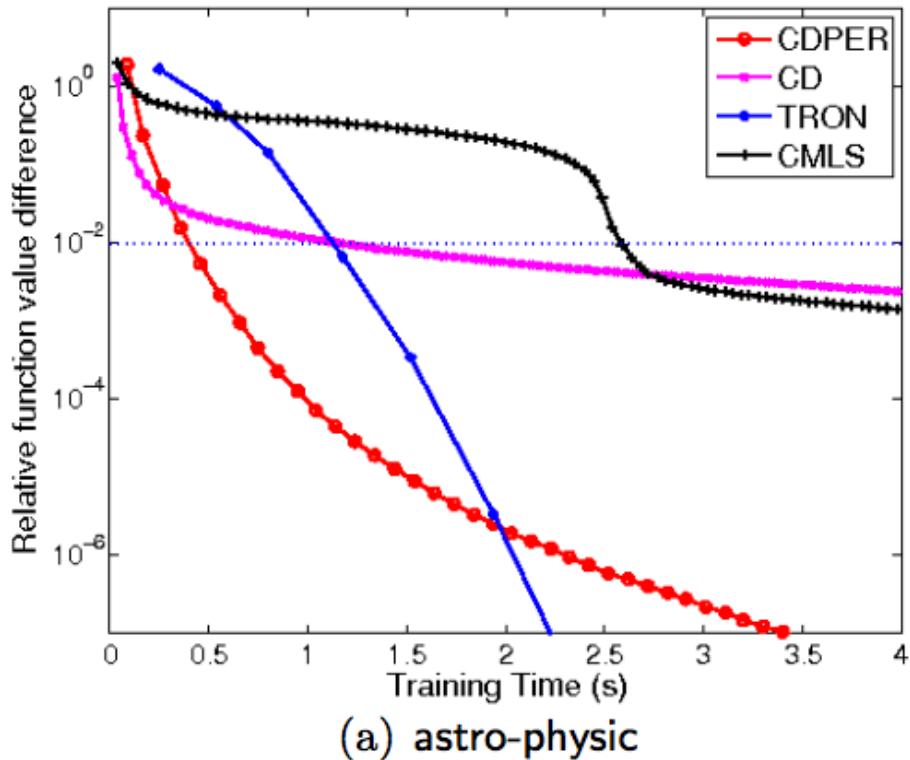


Convex optimization is relatively easy

- ❖ Ensure that there are no local minima
- ❖ Note: need special design for functions that are not differentiable (e.g., hinge loss)



Different optimization techniques

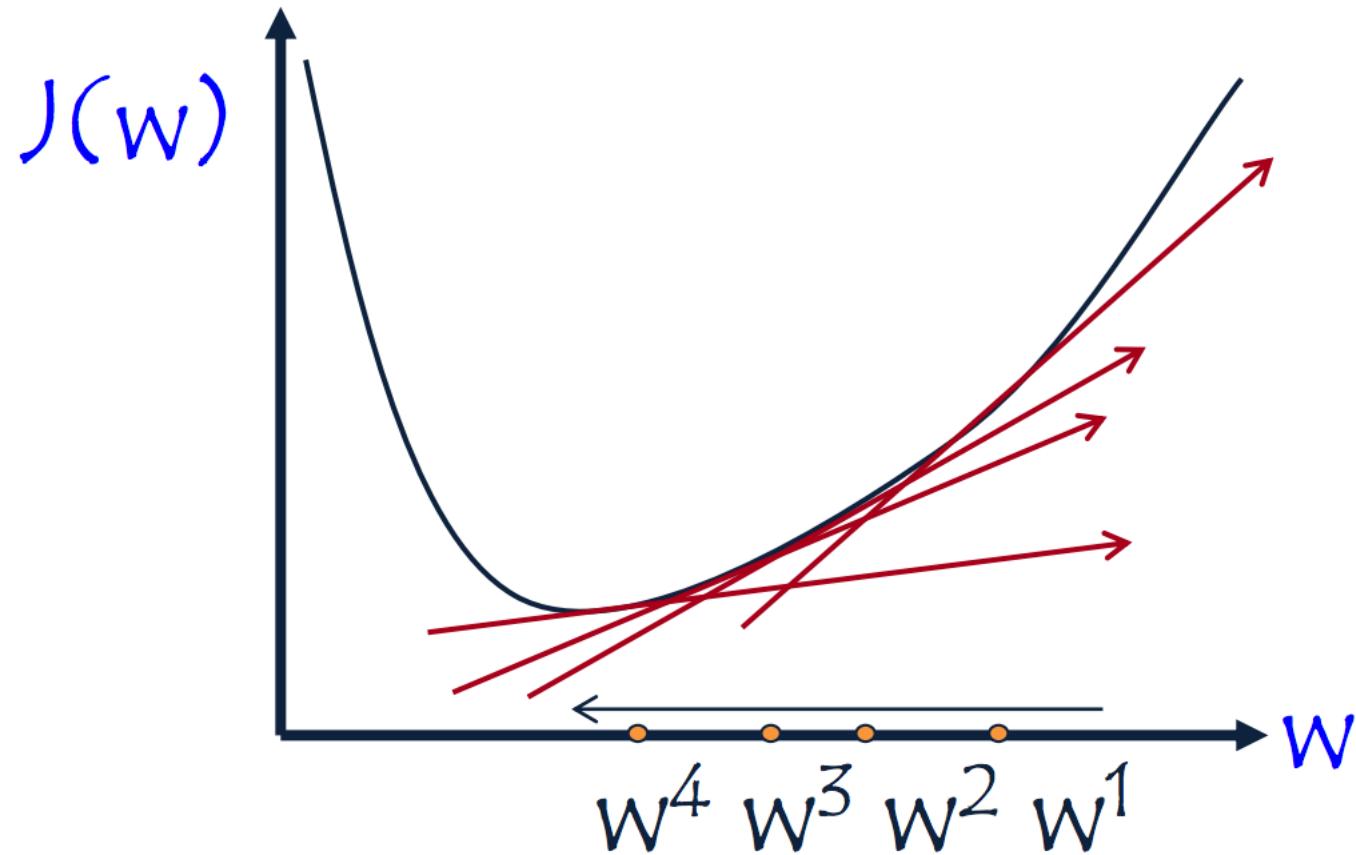


Some methods (e.g., SGD, CD) are fast in the early stage of optimization

Some methods (e.g., Newton methods) converge faster

Results from <http://www.cs.virginia.edu/~kc2wc/papers/ChangHsLi08.pdf>

Gradient descent



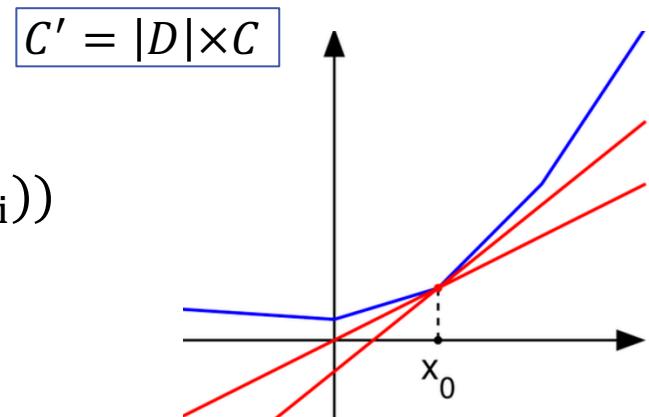
Why not gradient descent?

- ❖ For some functions, gradient may not exist

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i))$$

- ❖ Solution: use sub-gradient

$$\begin{aligned} f(\mathbf{w}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \\ &= \frac{1}{|D|} \sum_i \left(\frac{1}{2} \mathbf{w}^T \mathbf{w} + C' \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \right) \\ &\equiv \frac{1}{|D|} \sum_i f_i(\mathbf{w}) \\ \nabla f(\mathbf{w}) &= \frac{1}{|D|} \sum_i \nabla f_i(\mathbf{w}) \\ \nabla f_i(\mathbf{w}) &\equiv \begin{cases} \mathbf{w} & \text{if } y_i(\mathbf{w}^T \mathbf{x}_i) > 1 \\ \mathbf{w} - C' y_i \mathbf{x}_i & \text{otherwise} \end{cases}, \end{aligned}$$



Stochastic gradient descent

$$\begin{aligned}f(w) &= \frac{1}{2} w^T w + C \sum_i \max(0, 1 - y_i w^T x_i) \\&= \frac{1}{|D|} \sum_i \left(\frac{1}{2} w^T w + C' \max(0, 1 - y_i w^T x_i) \right) \\&\equiv f_i(w) \\\nabla f(w) &= \frac{1}{|D|} \sum_i \nabla f_i(w) = E_{i \sim D} \nabla f_i(w)\end{aligned}$$

- ❖ Approximate the true gradient by a gradient at a single example at a time

Repeat until converge:

Randomly pick one sample (x_i, y_i)

Update $w \leftarrow w - \eta \nabla f_i(w)$

Stochastic Sub-gradient Descent

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. Update $w \leftarrow w - \eta \nabla f(w)$
5. Return w

$$f(w) \equiv \frac{1}{2} w^T w + C \sum_i \max(0, 1 - y_i(w^T x_i))$$

Stochastic (sub)-gradient descent for SVM

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. if $y(w^\top x) < 1$
5. $w \leftarrow (1 - \eta)w + \eta C yx$
6. else
7. $w \leftarrow (1 - \eta)w$
8. Return w

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. if $y(w^\top x) < 0$
5. $w \leftarrow w + \eta yx$
6. Return w

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(w^\top x^{\text{test}})$

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(x, y)\}$

1. Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch 1 ... T :
3. For (x, y) in \mathcal{D} :
4. if $y(w^\top x) < 0$
5. $w \leftarrow w + \eta yx$
6. Return w

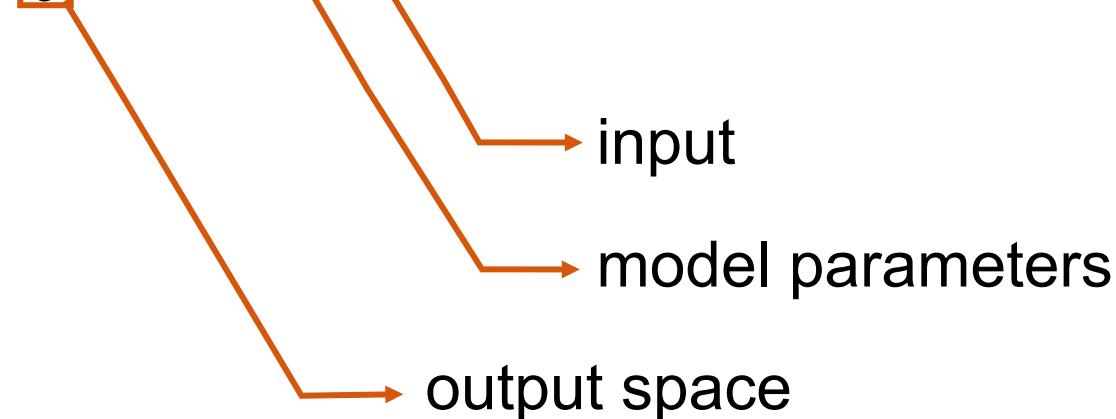
Prediction: y^{test}

Perceptron effectively minimizing:

$$\sum_i \max(0, 1 - y_i(w^\top x_i))$$

A General Formula

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} f(y; \mathbf{w}, \mathbf{x})$$



- ❖ Inference/Test: given \mathbf{w}, \mathbf{x} , solve argmax
- ❖ Learning/Training: find a good \mathbf{w}
- ❖ Today: $\mathbf{x} \in \mathbb{R}^n, \mathcal{Y} = \{-1, 1\}$ (binary classification)

Binary Linear Classifiers

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} f(y; \mathbf{w}, \mathbf{x})$$

- ❖ $\mathbf{x} \in \mathbb{R}^n, \mathcal{Y} = \{-1, 1\}$
- ❖ $f(y; \mathbf{w}, \mathbf{x}) \stackrel{\text{def}}{=} y(\mathbf{w}^\top \mathbf{x} + b) = y(\sum_i w_i x_i + b)$

- ❖ $\operatorname{argmax}_{y \in \mathcal{Y}} f(y; \mathbf{w}, \mathbf{x}) = \begin{cases} 1, & \mathbf{w}^\top \mathbf{x} + b \geq 0 \\ -1, & \mathbf{w}^\top \mathbf{x} + b < 0 \end{cases} = \operatorname{sgn}(\mathbf{w}^\top \mathbf{x} + b)$

(break ties arbitrarily)