

# **PHYS 139/239: Machine Learning in Physics**

**Lecture 3:**

**Support Vector Machine, Regularization, & Logistic Regression**

**Javier Duarte — January 17, 2023**

# Logistics

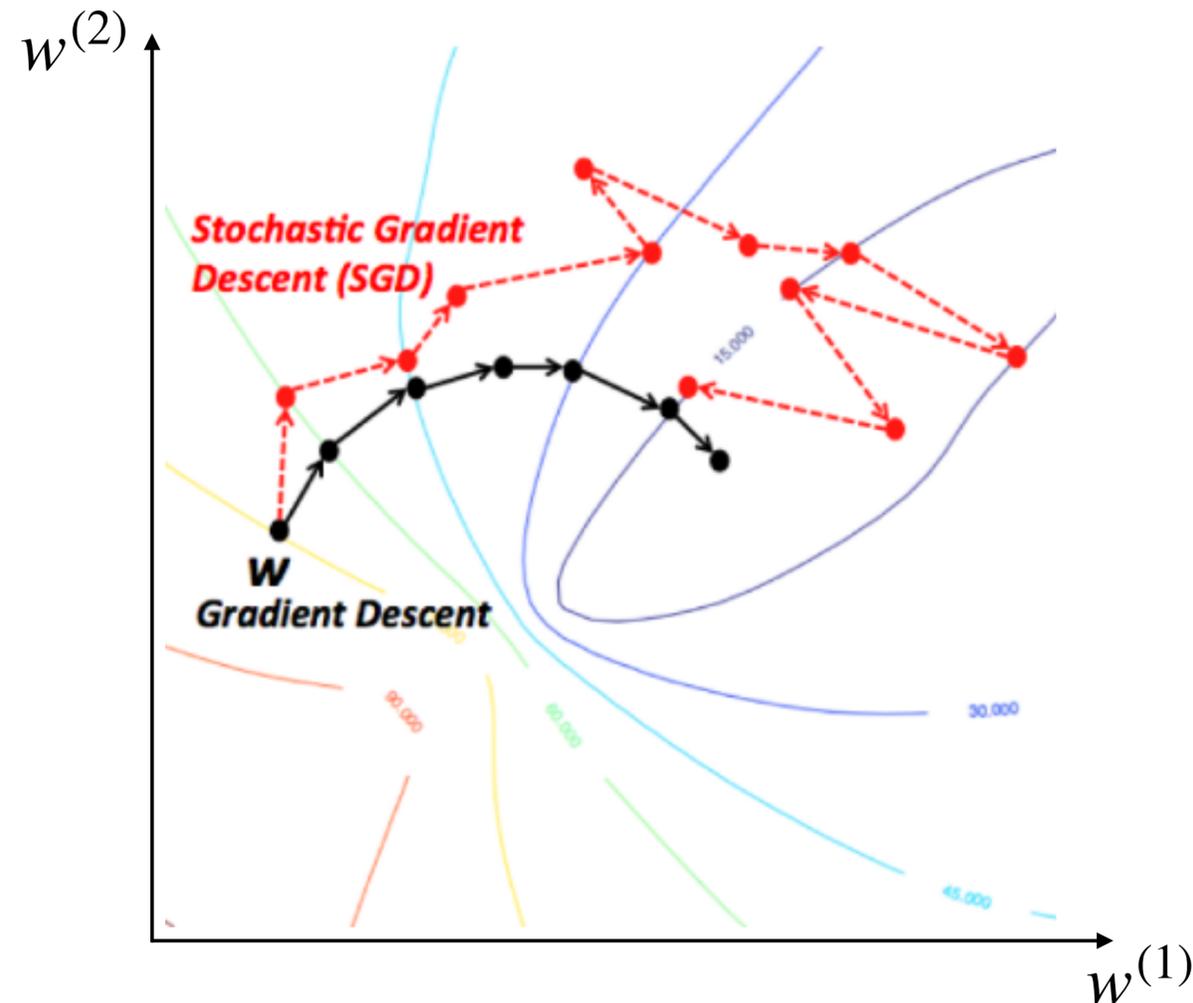
- Homework 1 draft version due Friday 1/20 5pm
  - Graded on effort (so attempt all problems!)
  - If you get stuck, explain why you're stuck
  - If you have trouble getting started, come to office hours or ask question in Slack
- Solutions will be released on Friday soon after deadline
- Homework 1 final version (where you correct things) due Wednesday 1/25 5pm
- Both versions are needed to get 100% (50% for draft, 50% for corrected version)

# Recap: (Stochastic) gradient descent

$$l(w) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i, | w))$$

• Gradient descent update:  $w(t + 1) = w(t) - \eta \nabla_w l(w(t))$

• SGD update:  $w(t + 1) = w(t) - \eta \nabla_w L(y, f(x | w(t)))$   
for a random  $(x, y) \in \mathcal{S}$



# SGD practical tips

- Divide the loss function by the number of examples (normalize):

$$w(t + 1) = w(t) - \frac{\eta}{N} \nabla_w l(w)$$

(Don't want the size of our updates to depend on the  $N$ )

- Start with a large step size  $\eta(t = 0)$
- Whenever the **validation error** stops going down, lower the step size:

$$\eta(t + 1) = \eta(t)/2$$

(step size must decrease over time to guarantee convergence)

- Stop when the **validation error** no longer decreases (**early stopping**)

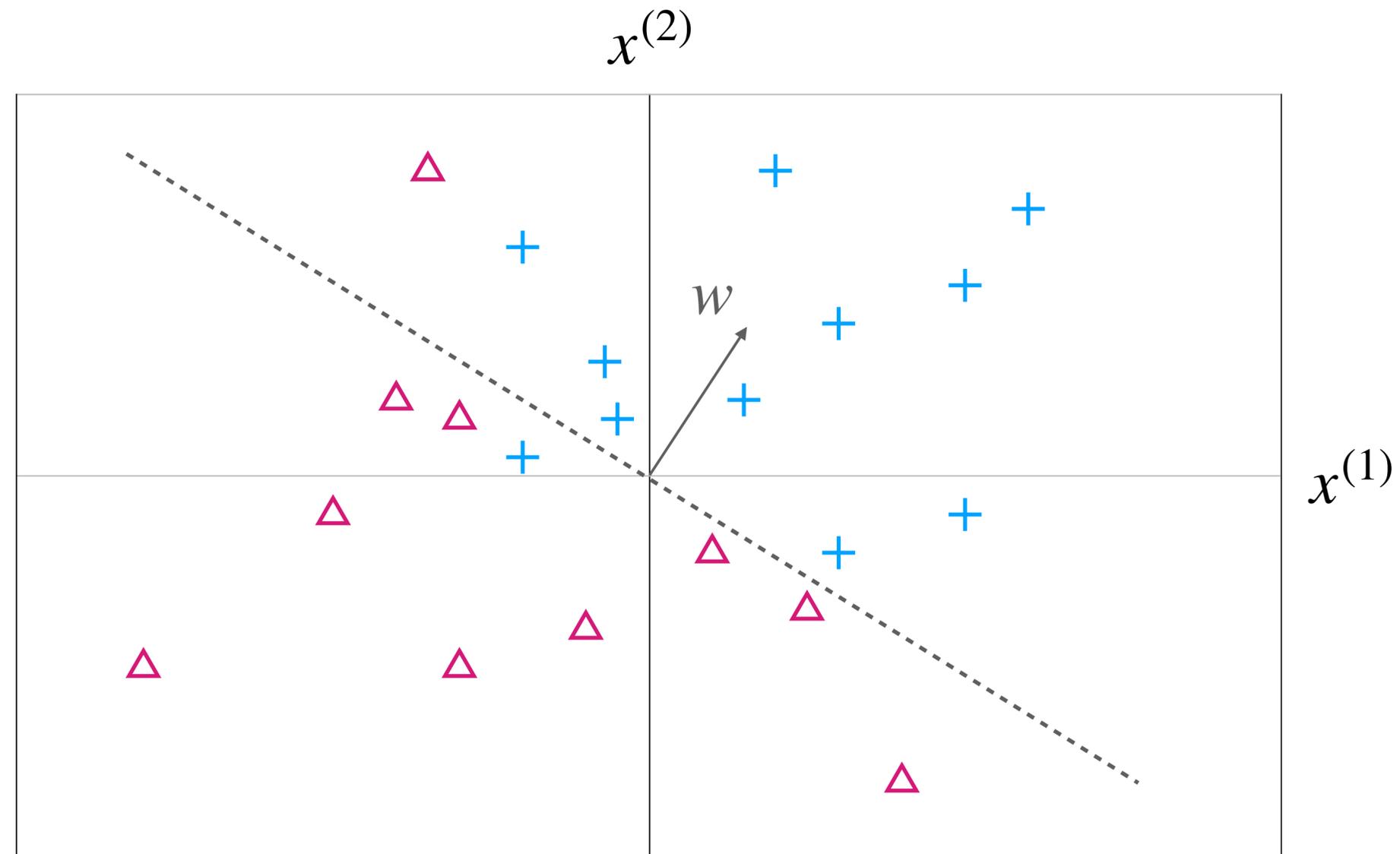
# Recap: Supervised learning pipeline

- Training dataset:  $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$  where  $x \in \mathbb{R}^D$  and  $y \in \mathbb{R}$
- Model / hypothesis class:  $f(x | w) = w^\top x$  (linear models) ↑  
For regression
- Loss function:  $L(y, y') = (y - y')^2$  (squared loss) ← or  $\phi(x)$  instead of  $x$
- Optimization algorithm: SGD
- Cross validation and model selection: 
- Testing and deployment

# Recap: Linear models for binary classification

- Linear model for regression:  $f(x | w) = w^T x$
- Linear model for binary classification:  $f(x | w) = \text{sign}(\underbrace{w^T x}_{\text{Raw score}}) \in \{+1, -1\}$

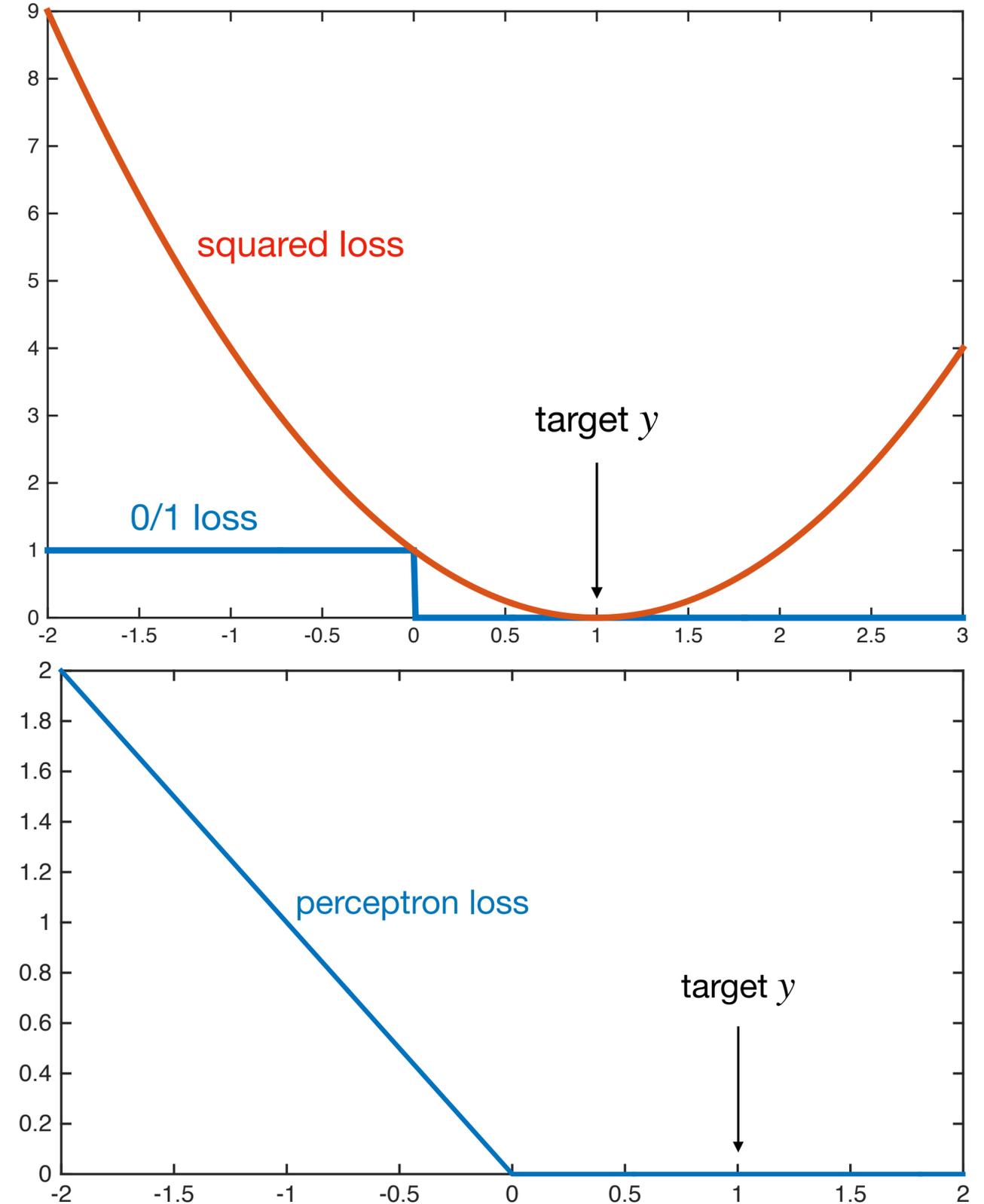
- Usually evaluate with 0/1 loss
- Optimize raw score using another loss (e.g. squared loss, perceptron loss)



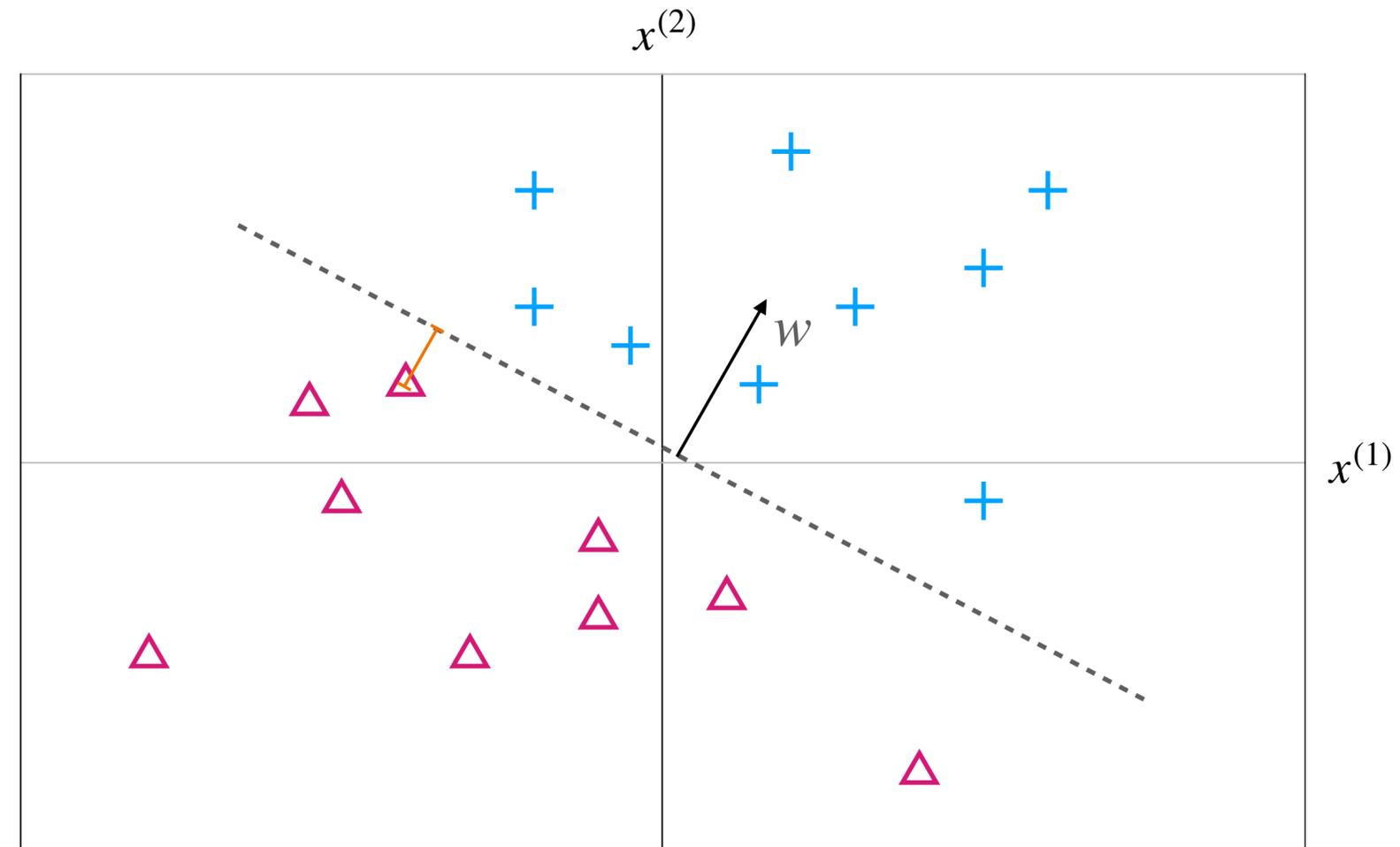
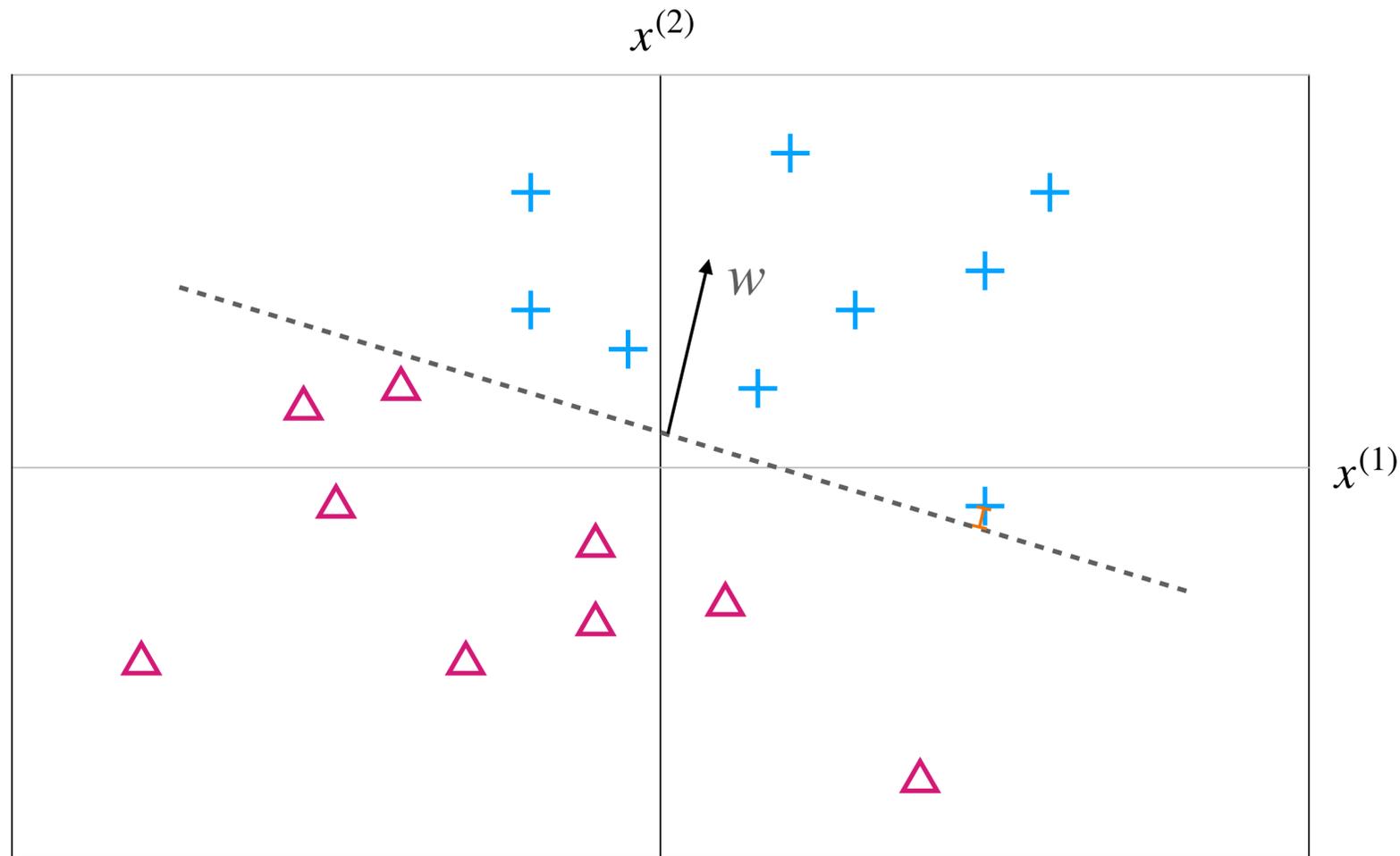
# Recap: Squared & perceptron losses

- **Squared loss:** Usually not good: can fail even on linearly separable data
- **Perceptron loss:** Reproduces perceptron update

$$w(t+1) = \begin{cases} w(t) & \text{correct} \\ w(t) + yx & \text{otherwise} \end{cases}$$



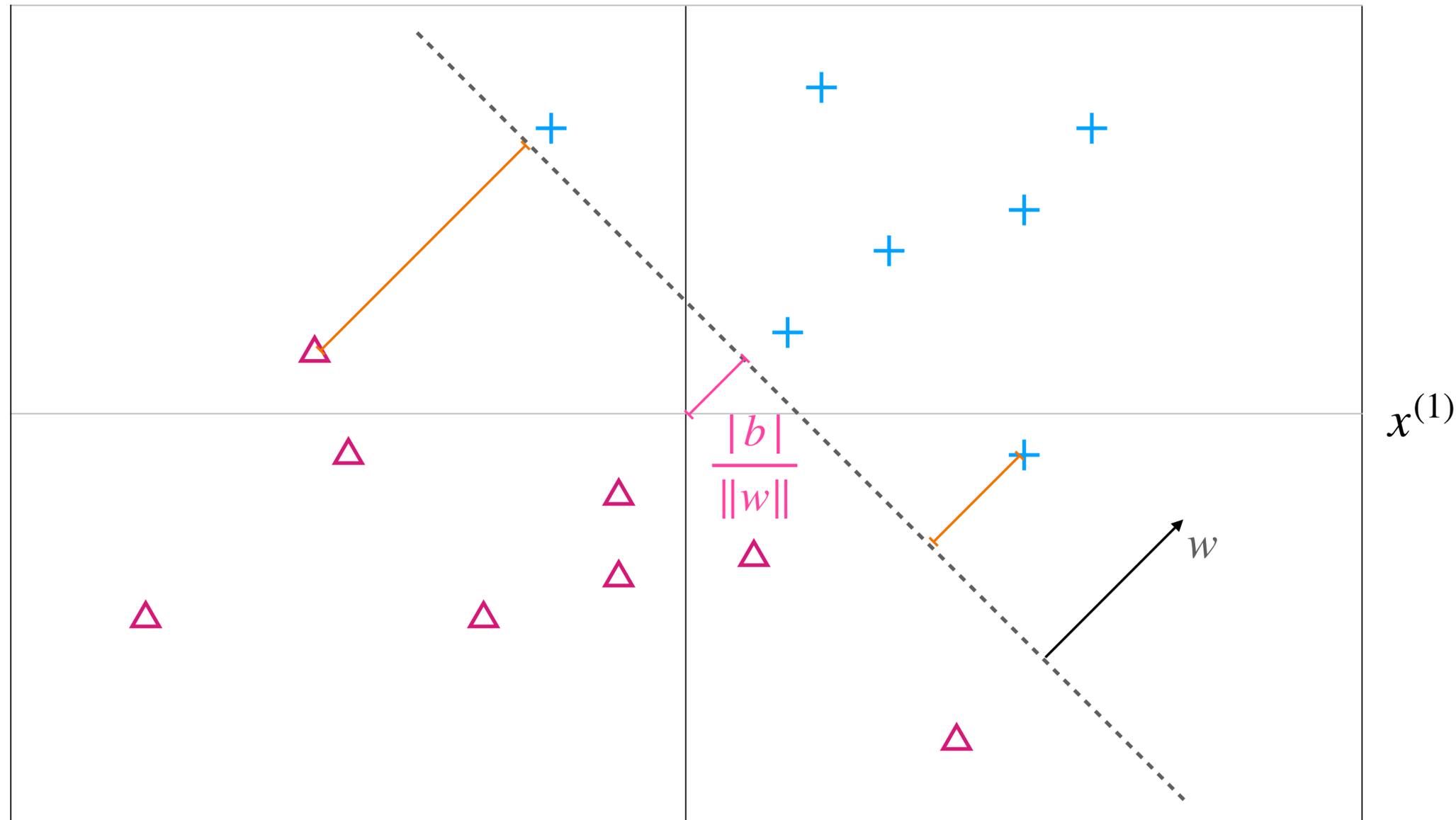
# Which classifier is better?



The classifier with a larger **margin!**  
(more likely to generalize better)

# Distance from a hyperplane

- $w^T x + b = 0$  defines a hyperplane in  $\mathbb{R}^D$  (affine subspace of dimension  $D - 1$ )



Distance:  $\frac{|w^T x + b|}{\|w\|}$

$L^2$  norm:  $\|w\| = \sqrt{w^T w}$

Signed distance:  $\frac{w^T x + b}{\|w\|}$

Raw score of linear model

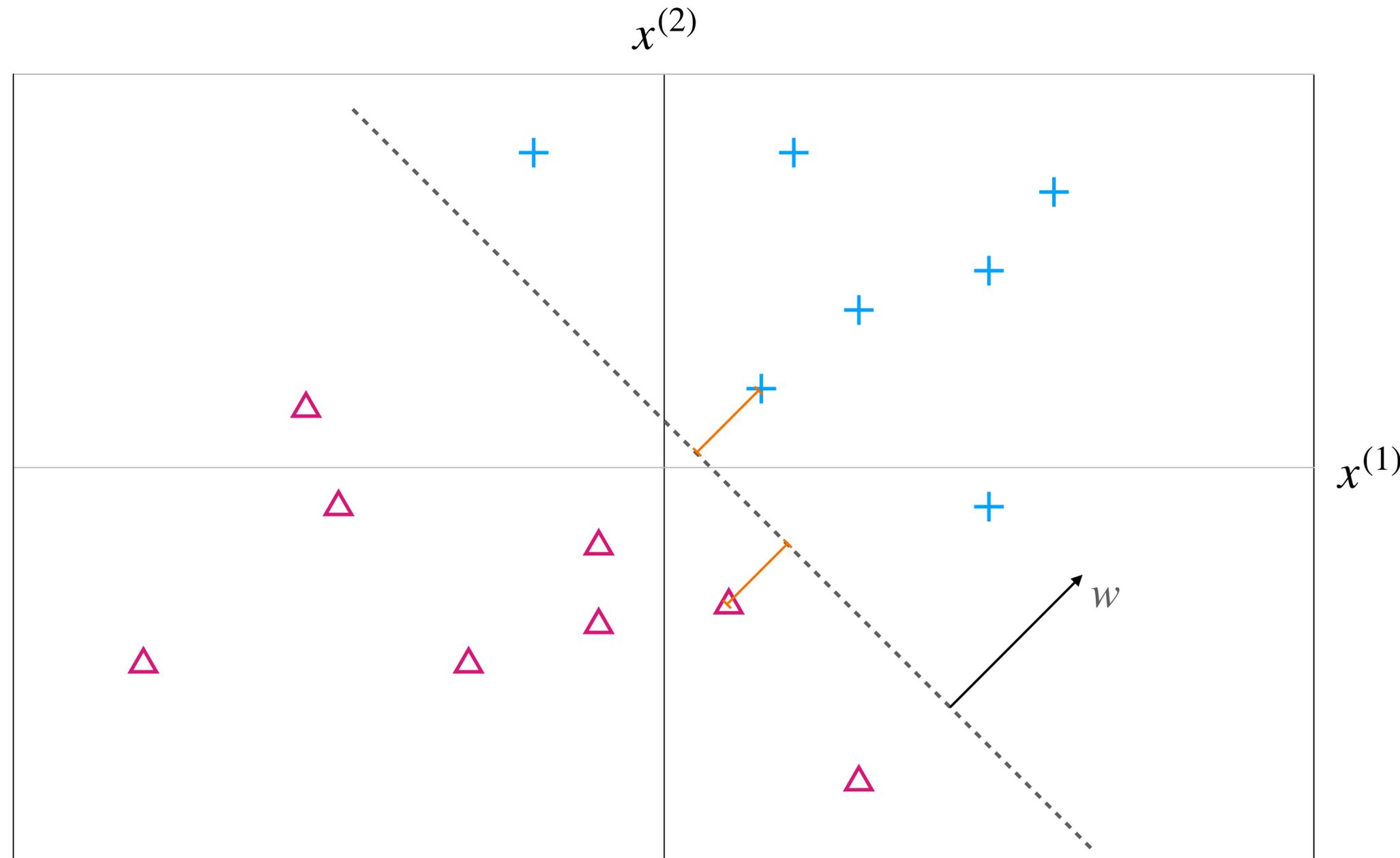
# How to maximize margin?

- Assuming linearly separable
- Choose  $w, b$  that maximize

$$\min_{(x,y) \in \mathcal{S}} \frac{y(w^T x + b)}{\|w\|}$$

- Equivalently, minimize  $\|w\|^2$  with the constraint

$$\min_{(x,y) \in \mathcal{S}} y(w^T x + b) = 1$$



See [Bishop Ch. 7](#) for more details

# Support vector machine (SVM)

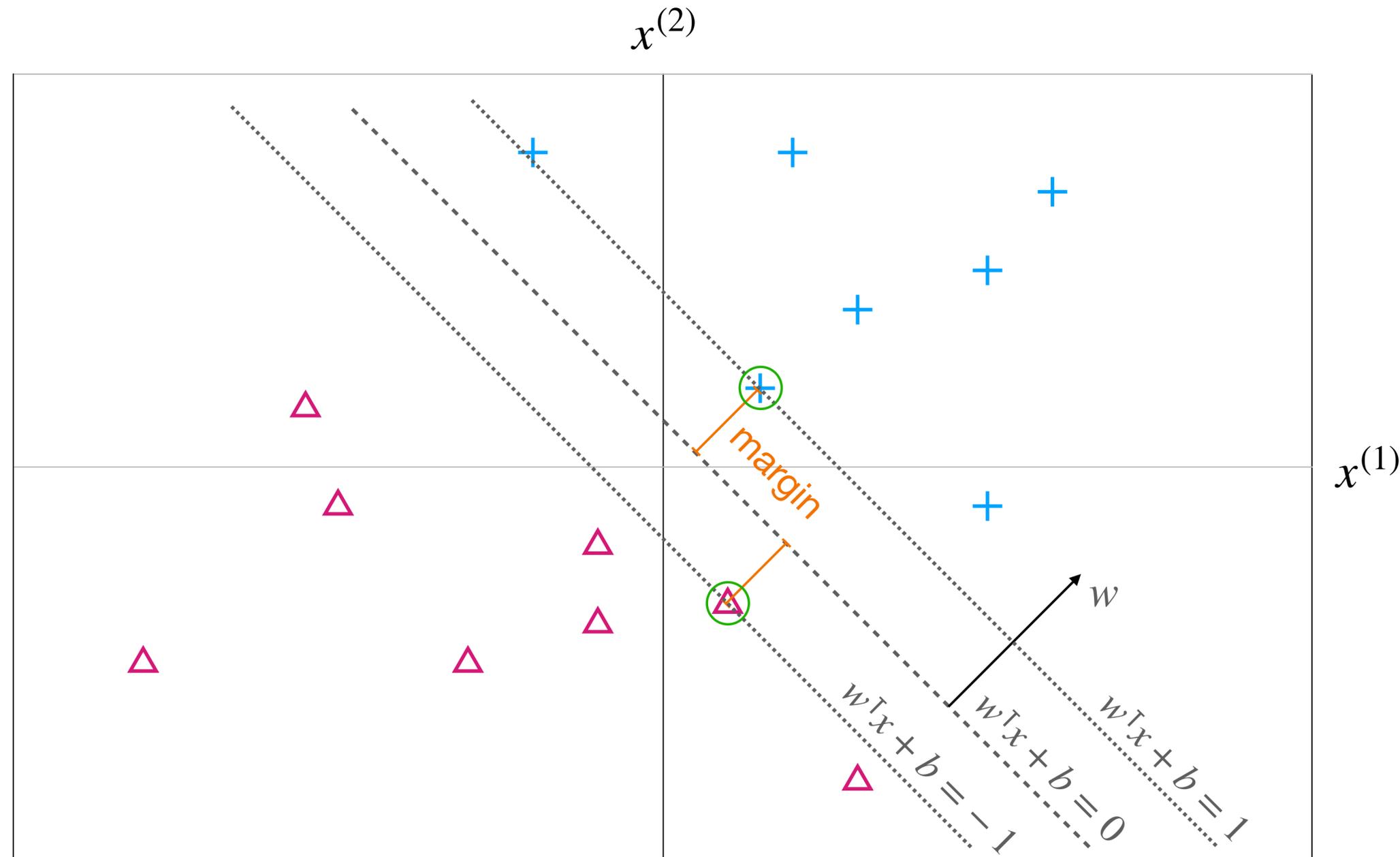
- Assuming linearly separable: Max margin classifier
- SVM optimization problem

$$\arg \min_{w,b} \|w\|^2$$

subject to

$$y_i(w^T x_i + b) \geq 1 \quad \forall i$$

some constraints are tight  
at the optimum (**support vectors**)



$$\text{margin} = \frac{1}{\|w\|}$$

# Soft-margin SVM

- Don't assume linearly separable
- Soft-margin SVM optimization problem

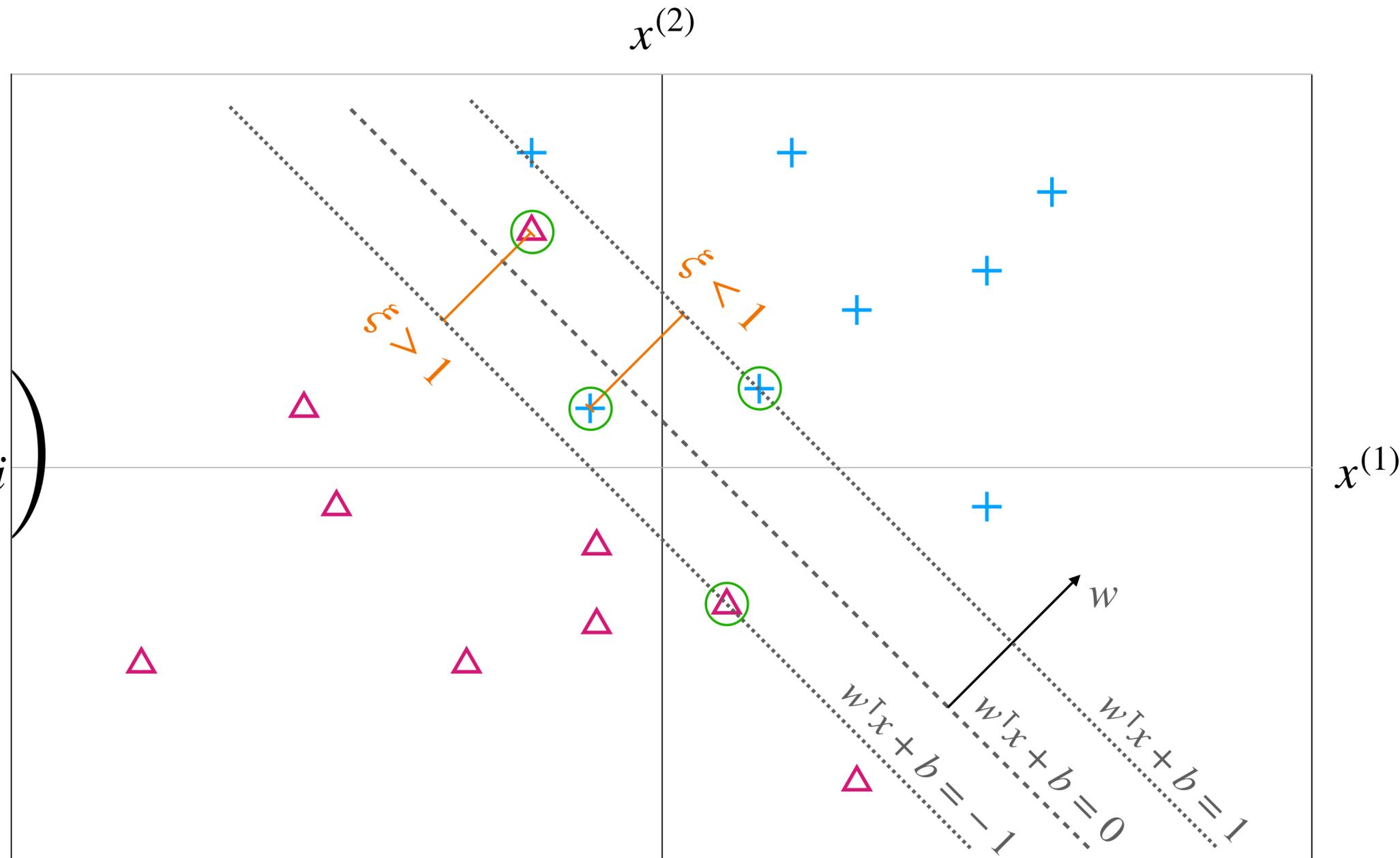
$$\arg \min_{w,b} \left( \|w\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i \right)$$

subject to

$$y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i > 0 \quad \forall i$$

Slack



$C$  controls the trade-off between size of margins and margin violations

# Hinge loss

- Soft-margin SVM optimization problem

$$\arg \min_{w,b} \|w\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i$$

subject to

$$y_i(w^\top x_i + b) \geq 1 - \xi_i \quad \forall i$$

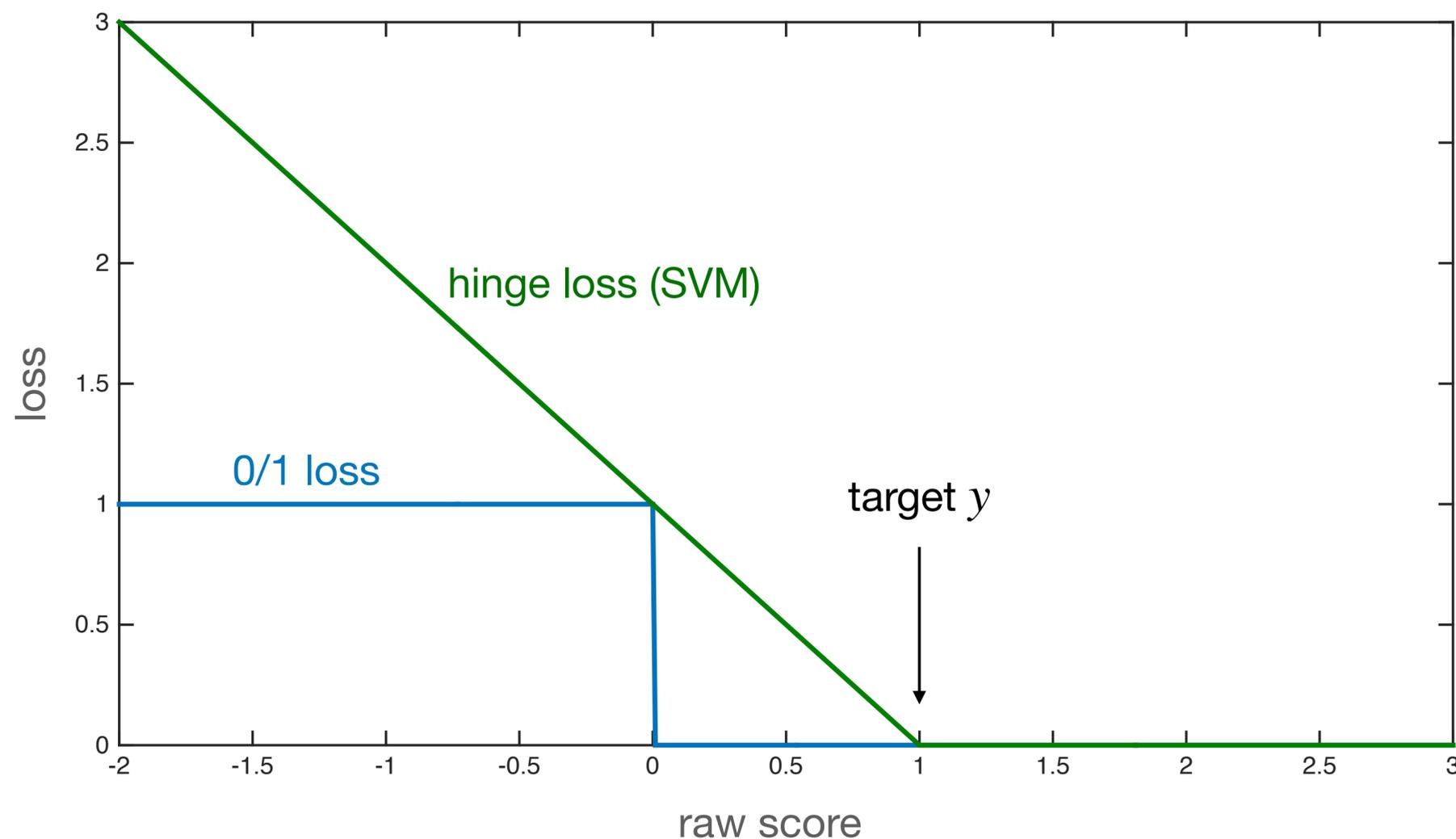
$$\xi_i > 0 \quad \forall i$$



- $\xi_i \geq \max(0, 1 - y_i(w^\top x_i + b))$

SVM equivalent to minimizing:

$$\frac{1}{N} \sum_{i=1}^N \underbrace{\max(0, 1 - y_i(w^\top x_i + b))}_{\text{Hinge loss}} + \underbrace{\frac{1}{C} \|w\|^2}_{\text{Regularization}}$$



# Kernel trick preface

- Theorem:

Optimal  $w$  has the form  $w = \sum_{i=1}^N \alpha_i x_i$  ( $\alpha_i \neq 0$  only if  $x_i$  is a **support vector**)

- Alternative (dual) formulation of SVM:

$$\arg \min_{\alpha, b} \sum_{i,j} \alpha_i \alpha_j x_i^\top x_j + \frac{C}{N} \sum_{i=1}^N \xi_i \quad \text{subject to} \quad \sum_{j=1}^N y_j (\alpha_j x_j^\top x_i + b) \geq 1 - \xi_i \quad \forall i$$
$$\xi_i \geq 0 \quad \forall i$$

- Predictions only depend on **support vectors**
- Optimization and predictions only depend on dot products  $x^\top x'$  of input vector  $x, x'$

# Kernel trick

- Can replace dot products  $x^\top x'$  with arbitrary **kernels**  $k(x, x')$ !

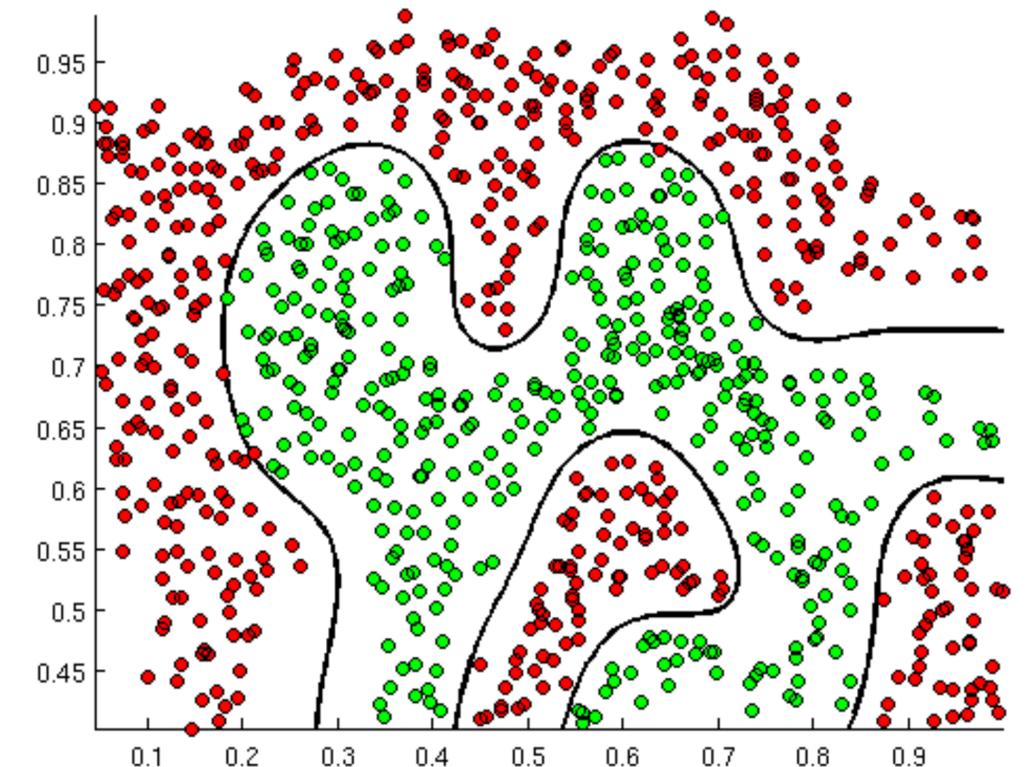
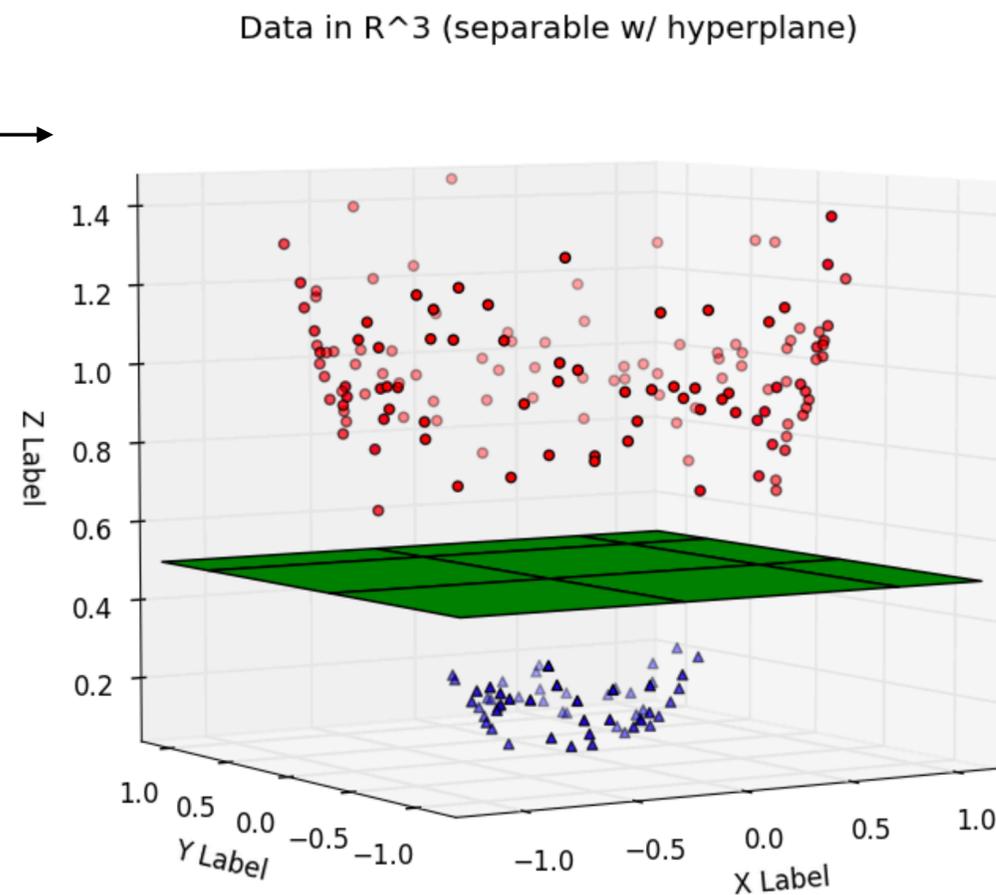
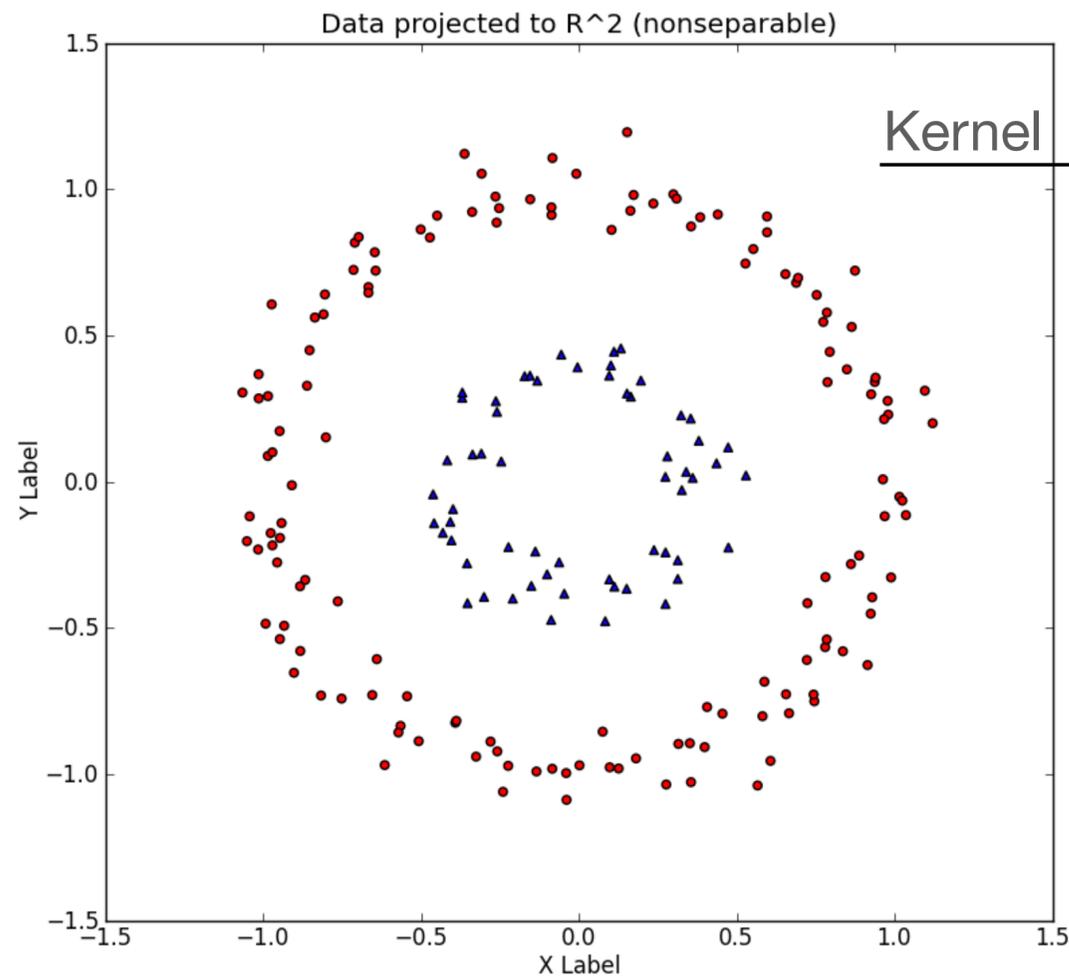
- **Polynomial kernel:**  $k(x, x') = (x^\top x' + c)^d$

- **Gaussian kernel:**  $k(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$

A kernel formally behaves as

$$k(x, x') = \langle \phi(x), \phi(x') \rangle$$

(without explicitly computing  $\phi(x), \phi(x')$ )



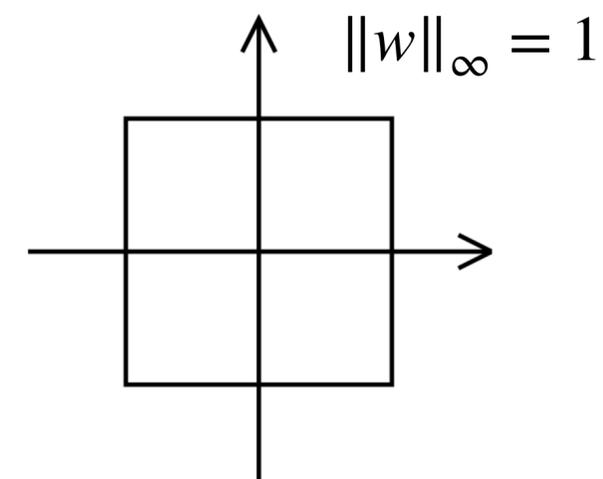
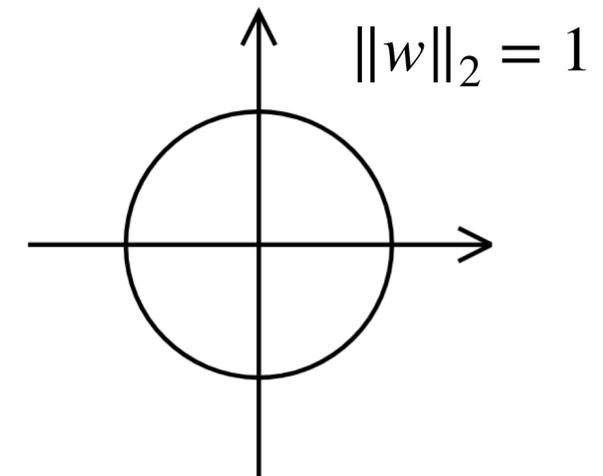
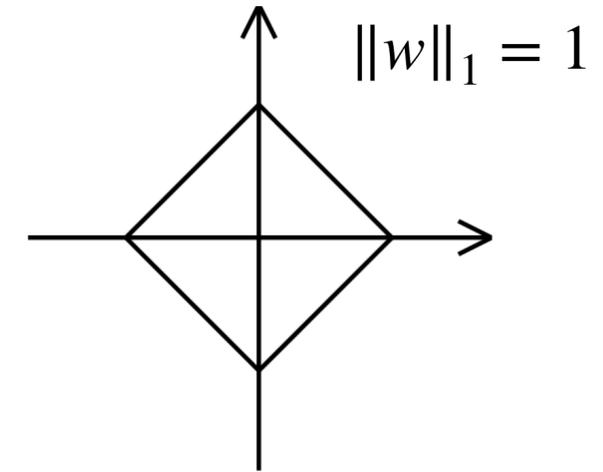
# Norms

- $L^0$  norm:  $\|w\|_0 =$  number of nonzero entries of  $w$  (not a norm!)

- $L^1$  norm:  $\|w\|_1 = \sum_j |w^{(j)}|$  (sum of abs. values of entries)

- $L^2$  norm:  $\|w\| = \|w\|_2 = \sqrt{w^T w} = \sqrt{\sum_j (w^{(j)})^2}$

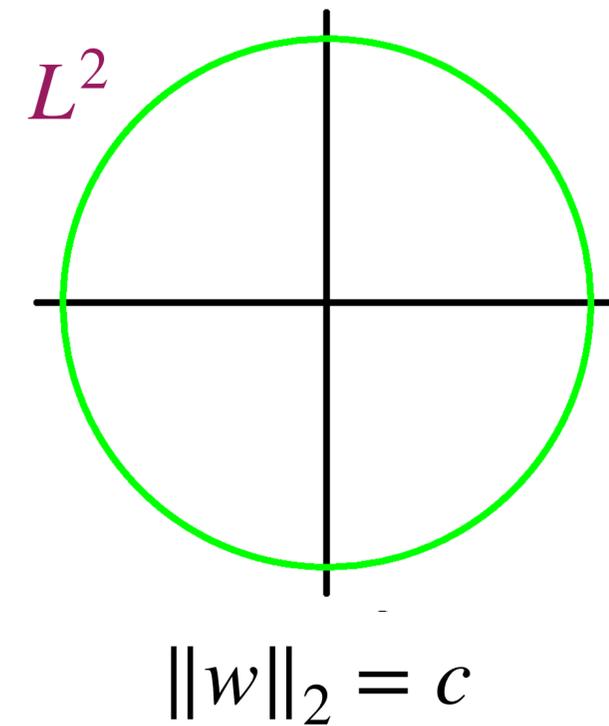
- $L^\infty$  norm:  $\|w\|_\infty = \max_j |w^{(j)}|$  (max abs. value of entries)



# Regularization

- $L^2$ -regularization for regression (ridge regression)

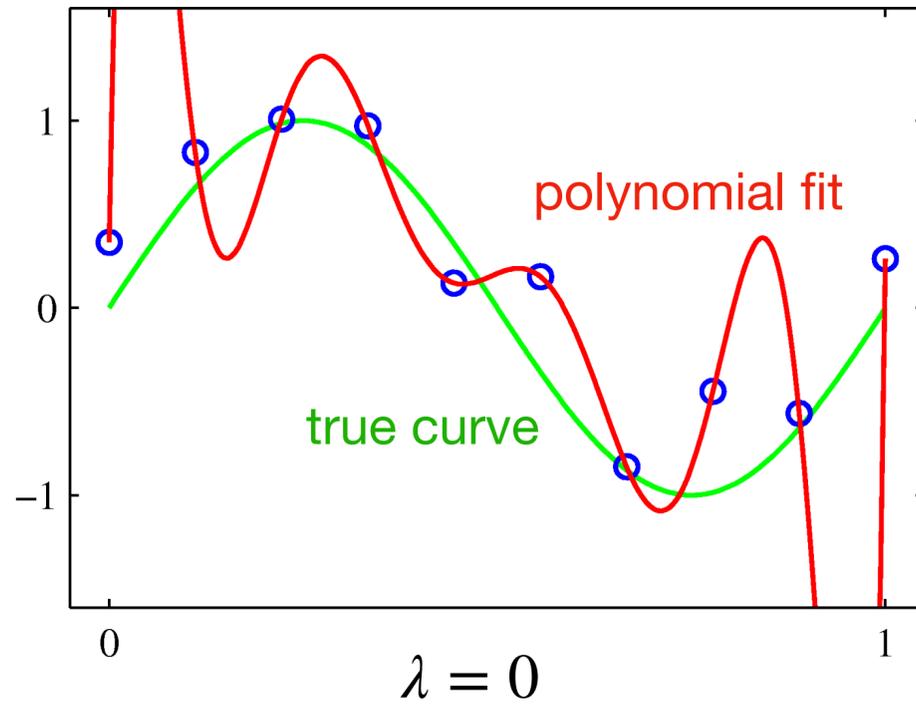
$$\arg \min_w \sum_{i=1}^N \underbrace{(y_i - w^\top x_i)^2}_{\text{Training loss}} + \underbrace{\lambda \|w\|^2}_{\text{Regularization}}$$



- Trades off model complexity vs. training loss
- Each choice of  $\lambda$  gives a model class (larger  $\lambda$  constrains  $w$  to be smaller)
- Regularization can be combined with any loss

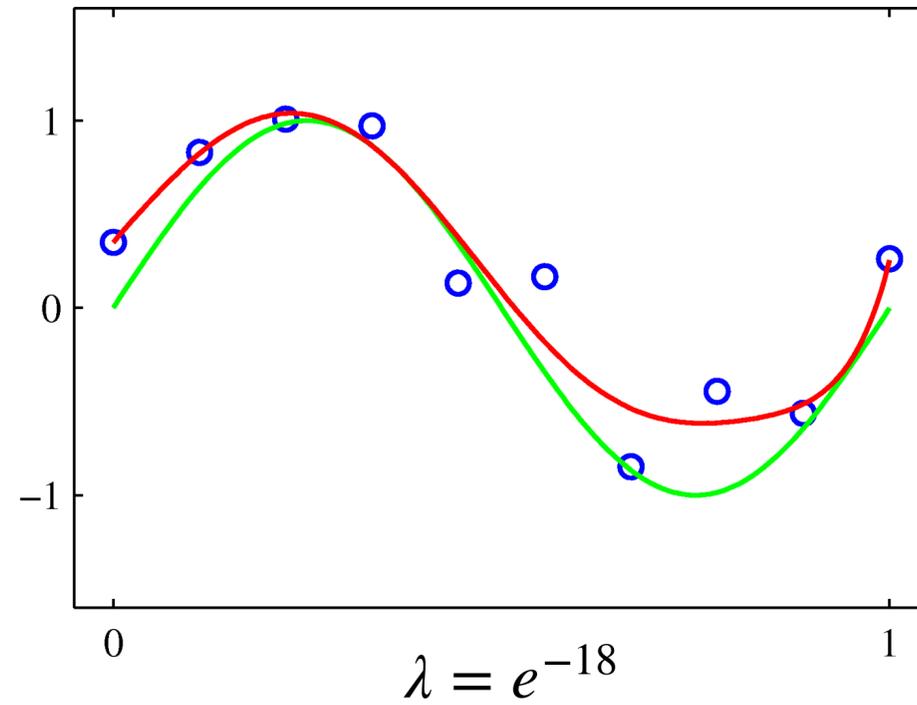
# Regularization example

- **Example:** polynomial curve fitting (10 points, degree 9)

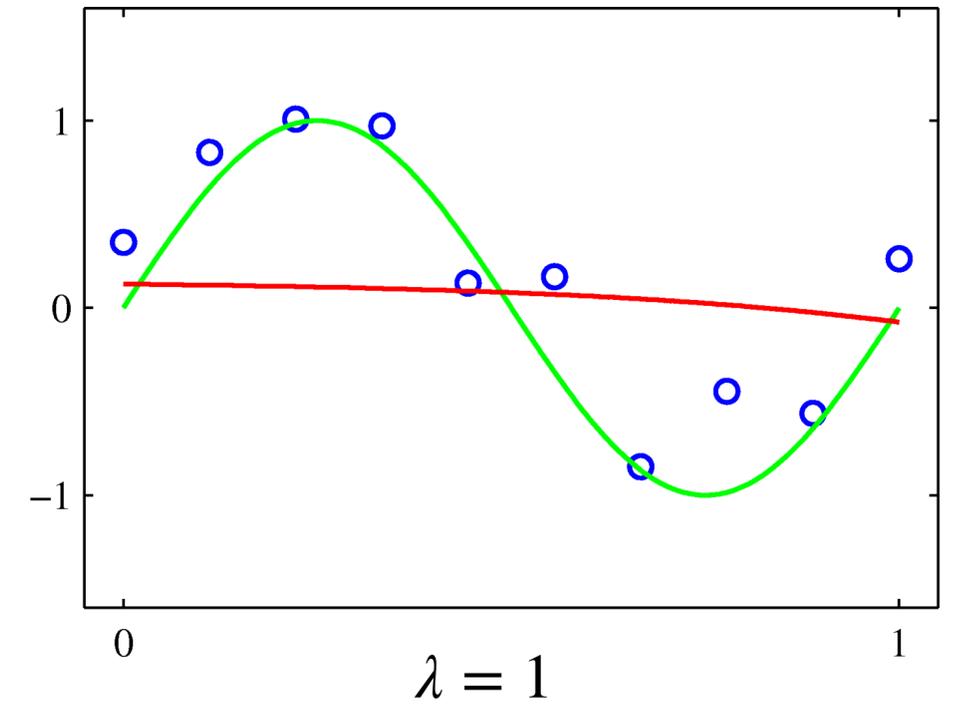


low  $\lambda$ : overfitting  
(high variance)

$$w = \begin{bmatrix} 0.35 \\ 232.37 \\ -5321.83 \\ 48568.31 \\ -231639.30 \\ 640042.26 \\ -1061800.52 \\ 1042400.18 \\ -557682.99 \\ 125201.43 \end{bmatrix}$$



$$w = \begin{bmatrix} 0.35 \\ 4.74 \\ -0.77 \\ -31.97 \\ -3.89 \\ 55.28 \\ 41.32 \\ -45.95 \\ -91.53 \\ 72.68 \end{bmatrix}$$

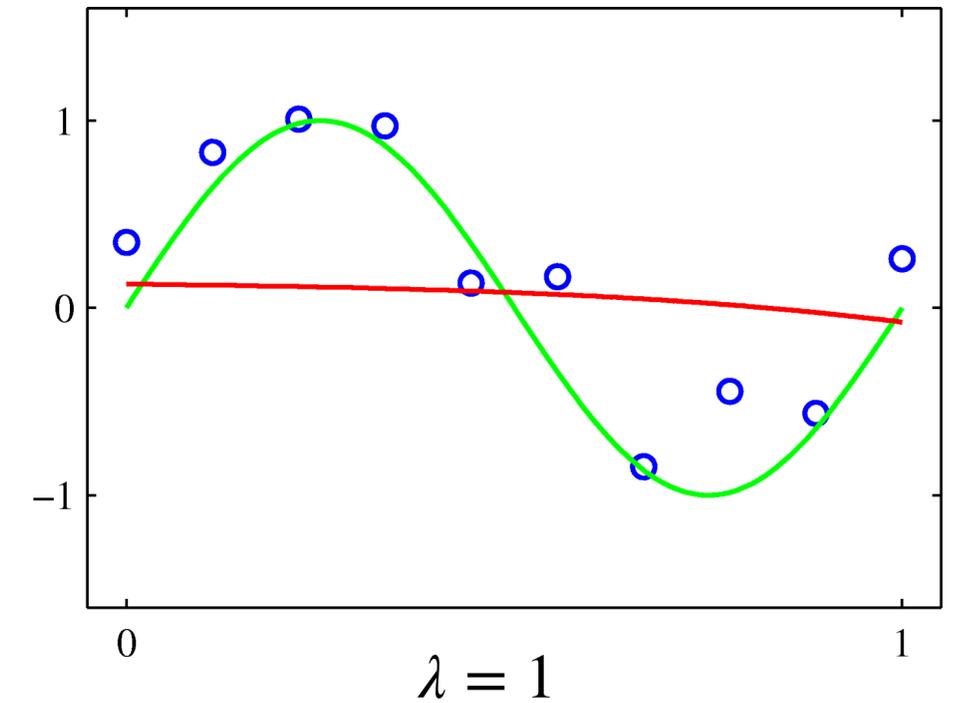
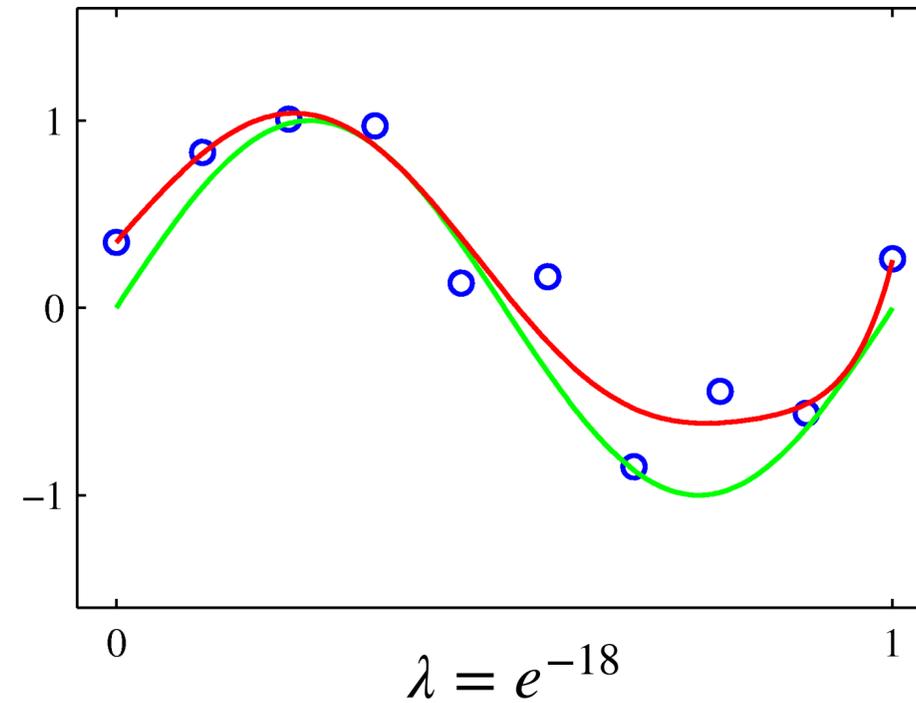
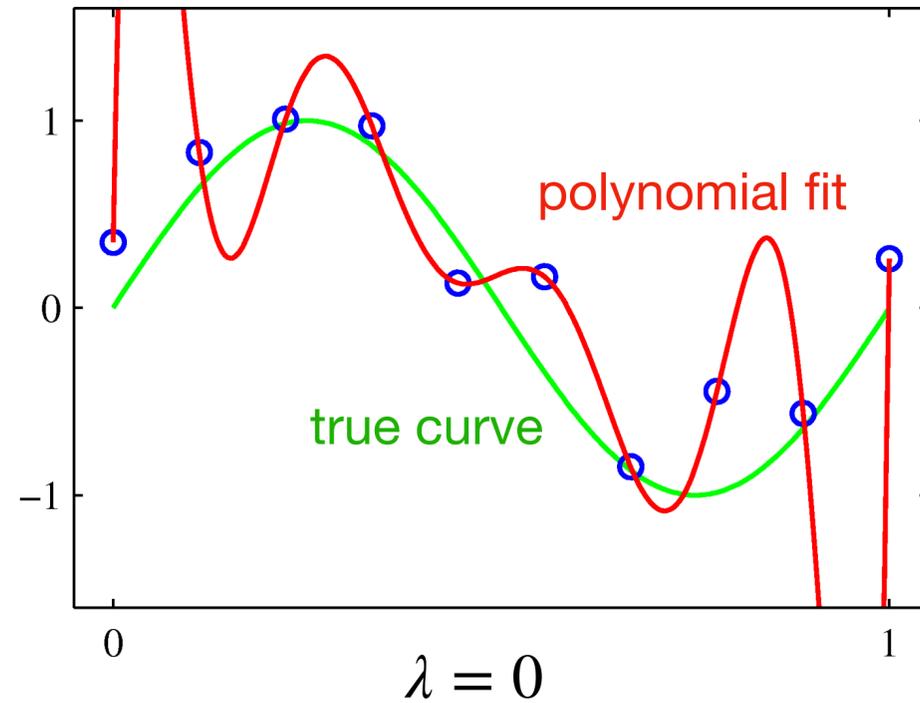


high  $\lambda$ : underfitting  
(high bias)

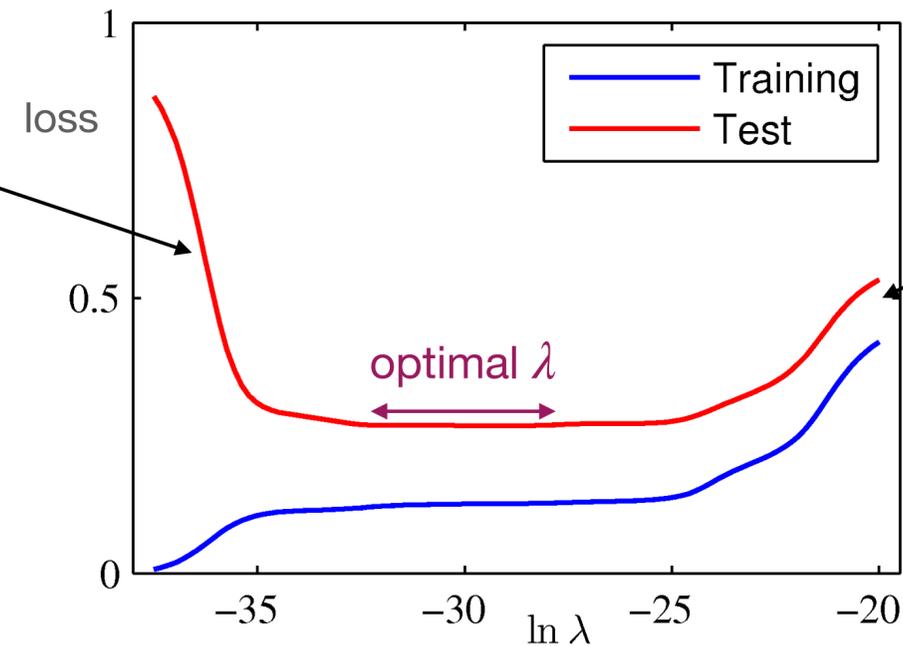
$$w = \begin{bmatrix} 0.13 \\ -0.05 \\ -0.06 \\ -0.05 \\ -0.03 \\ -0.02 \\ -0.01 \\ 0.00 \\ 0.00 \\ 0.01 \end{bmatrix}$$

# Regularization example

- **Example:** polynomial curve fitting (10 points, degree 9)



low  $\lambda$ : overfitting  
(high variance)



high  $\lambda$ : underfitting  
(high bias)

The optimal choice of  $\lambda$  depends on the training size  $N$

# Model class interpretation

- Minimize

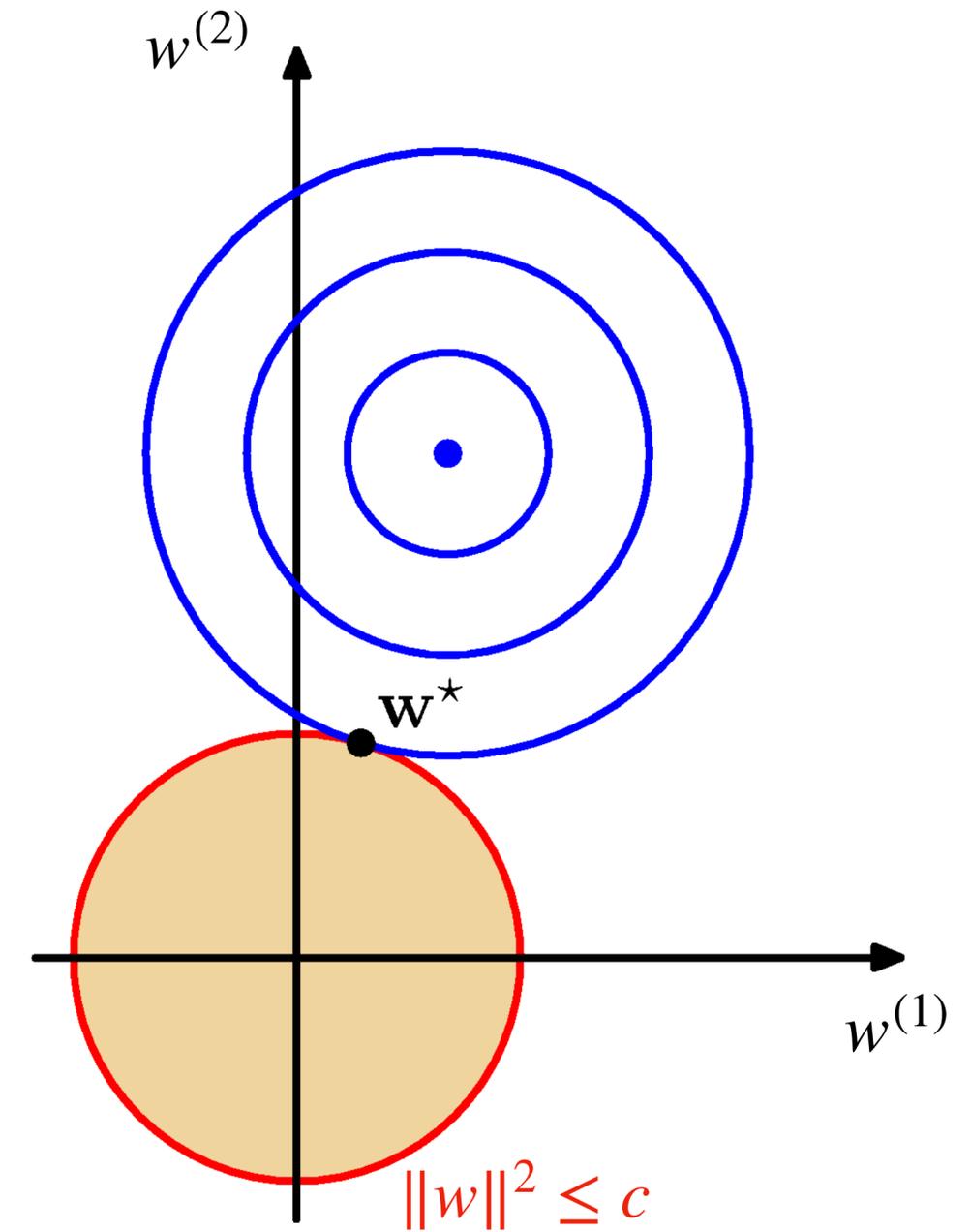
$$\sum_{i=1}^N (y_i - w^T x_i)^2 + \lambda \|w\|^2$$

↕ Equivalent

- Minimize

$$\sum_{i=1}^N (y_i - w^T x_i)^2 \text{ with constraint } \|w\|^2 < c$$

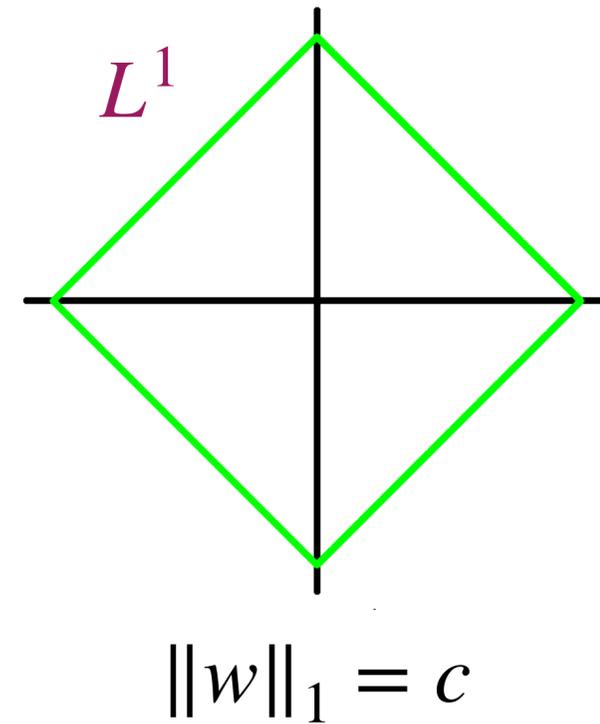
(using Lagrange multipliers)



# Lasso & sparsity

- $L^1$ -regularization for (lasso) regression

$$\arg \min_w \sum_{i=1}^N \underbrace{(y_i - w^\top x_i)^2}_{\text{Training loss}} + \underbrace{\lambda \|w\|_1}_{\text{Regularization}}$$

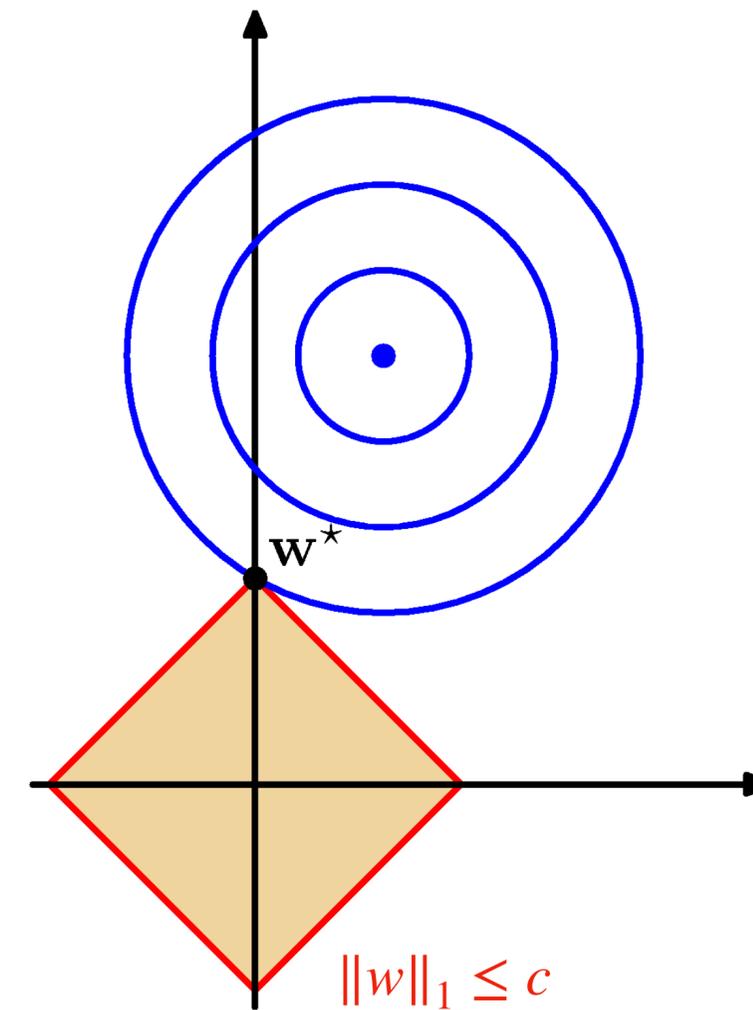
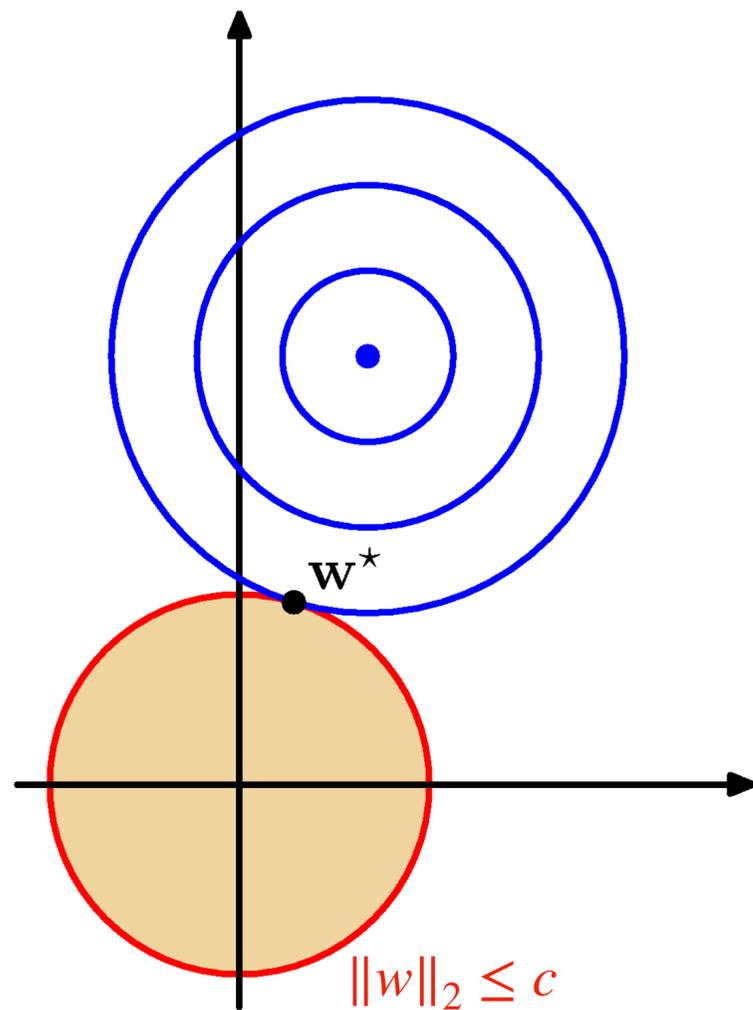


- We say that  $w$  is sparse if several of its entries are zero ( $\|w\|_0$  is small)
- Finding a sparse  $w$  is useful, e.g. time/memory efficiency (only some entries are needed to compute  $w^\top x$ )
- We cannot use  $L^0$  regularization directly (not continuous)
- However,  $L^1$  regularization (lasso) induces **sparsity**!

# Model class interpretation: $L^1$ vs. $L^2$

$$\arg \min_w \sum_{i=1}^N \underbrace{(y_i - w^\top x_i)^2}_{\text{Training loss}} + \underbrace{\lambda \|w\|^2}_{L^2 \text{ regularization}}$$

$$\arg \min_w \sum_{i=1}^N \underbrace{(y_i - w^\top x_i)^2}_{\text{Training loss}} + \underbrace{\lambda \|w\|_1}_{L^1 \text{ regularization}}$$



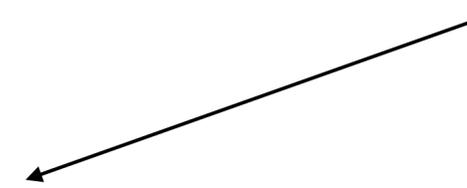
# Updated supervised learning pipeline

- Training dataset:  $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$  where  $x \in \mathbb{R}^D$  and  $y \in \mathbb{R}$
- Model / hypothesis class:  $f(x | w) = w^\top x$  (linear models) ↑  
For regression
- Loss function:  $L(y, y') = (y - y')^2$  (squared loss) ← or  $\phi(x)$  instead of  $x$
- Optimization algorithm: SGD with regularization ( $L^1$  or  $L^2$ )
- Cross validation and model selection: 
- Testing and deployment

Select  $\lambda$

# Probabilistic approach

Parametrized by  $w$



- Idea: Model a probability distribution  $p(y | x; w)$  of labels  $y$  given inputs  $x$
- Choose a form for  $p(y | x; w)$  (different for regression and classification)
- Write the likelihood of  $w$ , i.e. the probability of observing the labels  $y_i$  of the training dataset  $S$  given the inputs  $x_i$ :

$$p(S | w) = \prod_{i=1}^N p(y_i | x_i)$$

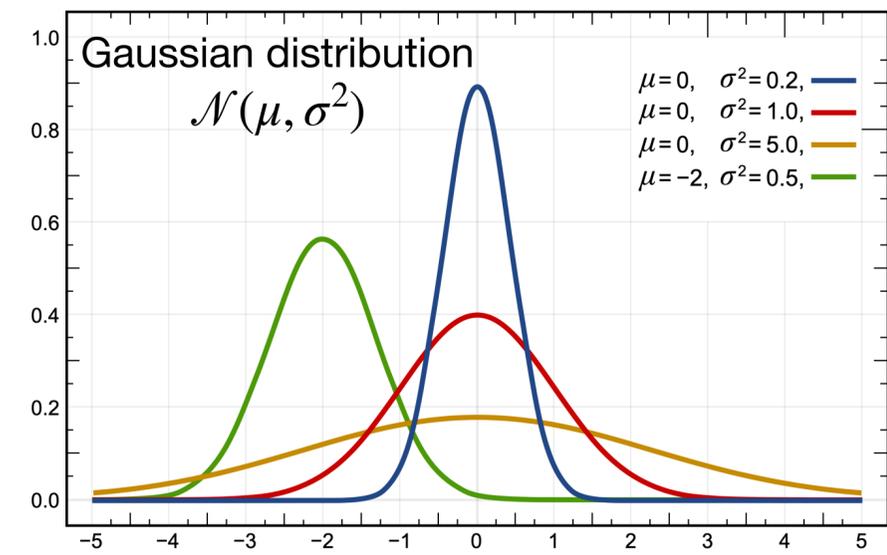
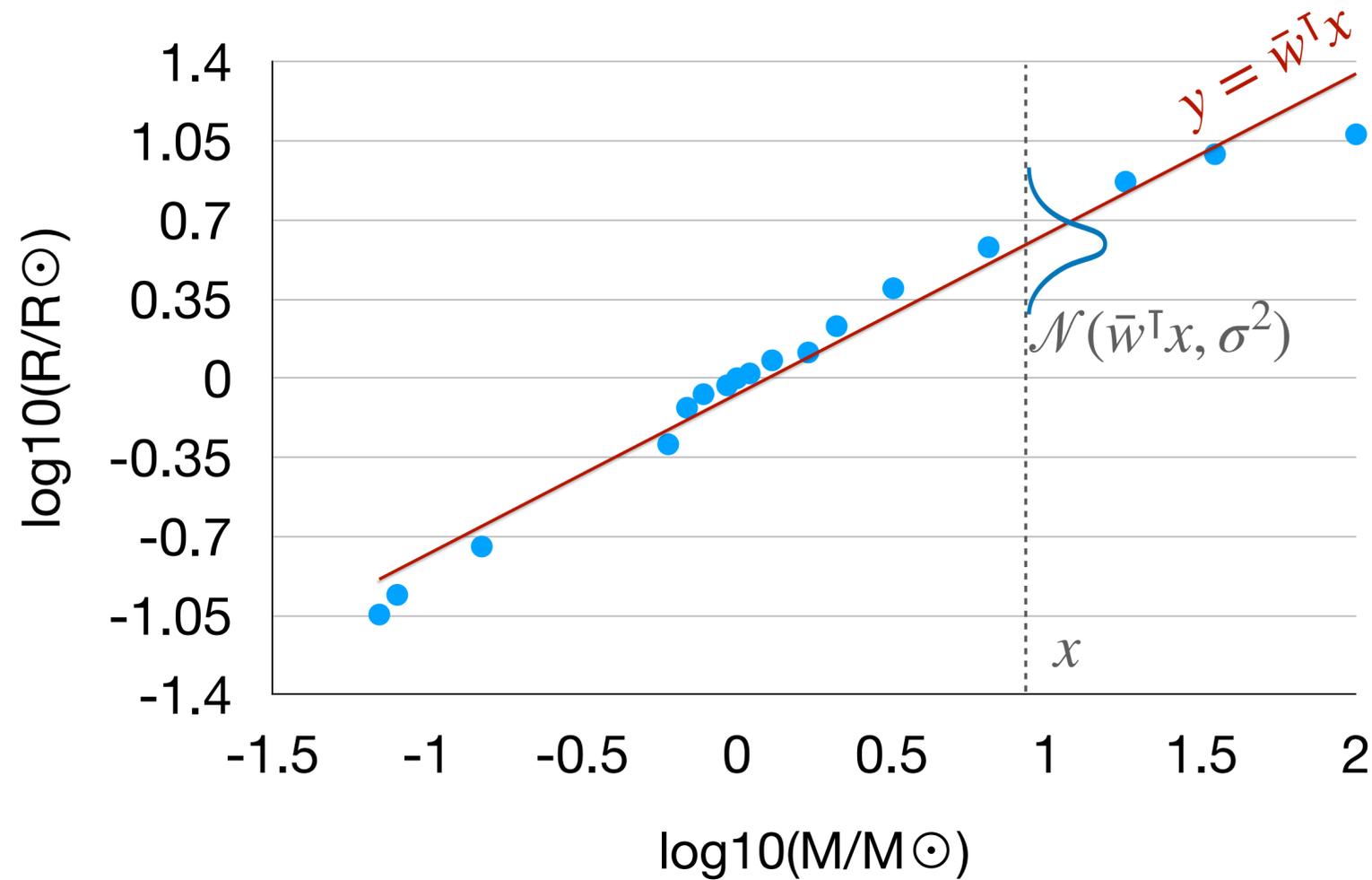
Assuming training examples are independent

- **Maximum likelihood estimation** (MLE): find  $w$  that maximizes the (log) likelihood:

$$\log p(S | w) = \sum_{i=1}^N \log p(y_i | x_i; w) = -l(w)$$

Equivalently, minimize the **loss function!**

# Linear regression revisited



- Assume labels  $y$  are distributed as  $\mathcal{N}(\bar{w}^T x, \sigma^2)$

True value of  $w$  →

← Fixed (unknown) variance

- Likelihood of  $w$  (assuming training samples are i.i.d.):

$$p(S | w) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right)$$

- Maximizing the likelihood is equivalent to minimizing

$$-\log P(S | w) \simeq \sum_{i=1}^N (y_i - w^T x_i)^2$$

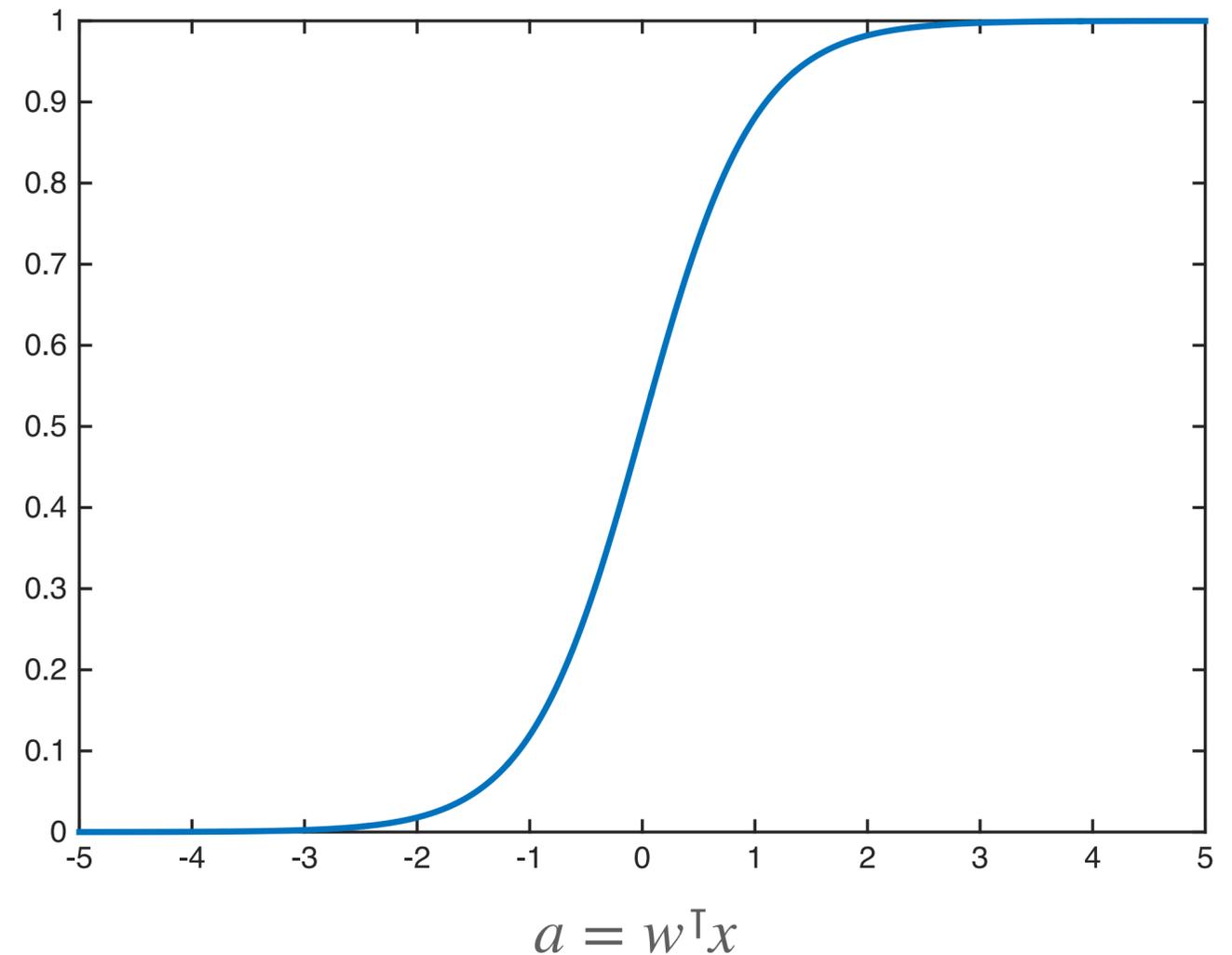
# Binary classification revisited

- Linear model for binary classification:  $f(x | w) = \text{sign}(\underbrace{w^T x}_{\text{Raw score}}) \in \{+1, -1\}$
- Idea: raw score to model the probability of each class

$\sigma(w^T x) \approx$  probability that  $y = +1$

↑  
Logistic/sigmoid function  $\sigma : \mathbb{R} \rightarrow (0,1)$

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



# What is the right loss function?

- Assume that the true probability that  $y = +1$  given  $x$  is  $\sigma(\bar{w}^\top x)$  and that

$p(y | \sigma(\bar{w}^\top x))$  is a **Bernoulli distribution**

↑  
True value of  $w$

- Likelihood of  $w$ : Only one of these two terms appears depending on  $y_i$

$$p(S | w) = \prod_{i=1}^N \sigma(w^\top x_i)^{\delta_{\{y_i=+1\}}} (1 - \sigma(w^\top x_i))^{\delta_{\{y_i=-1\}}}$$

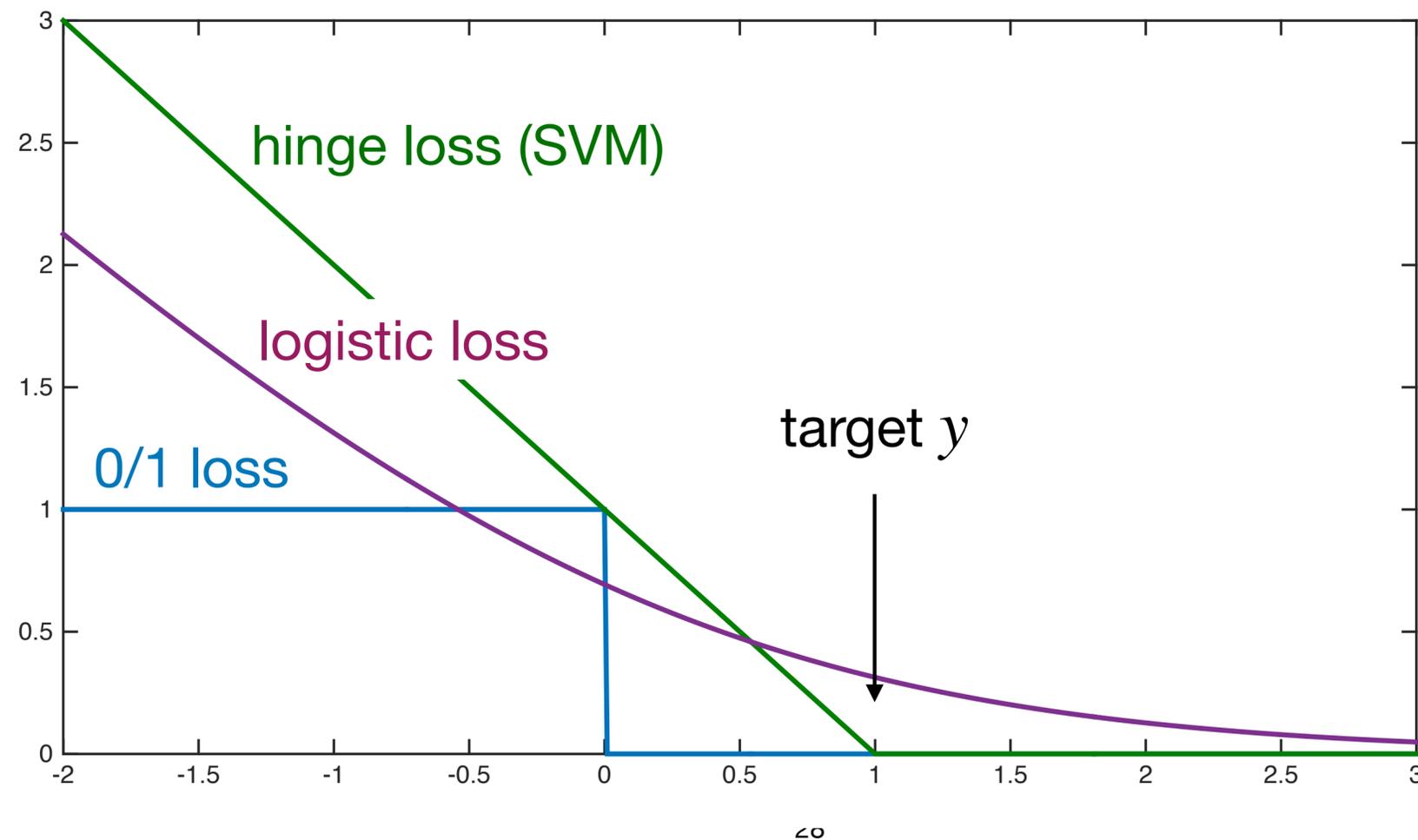
- Negative log likelihood of  $w$  a.k.a. **logistic / log / binary cross-entropy loss**:

$$-\log p(S | w) = - \sum_{i=1}^N \delta_{\{y_i=+1\}} \log \sigma(w^\top x_i) + \delta_{\{y_i=-1\}} \log(1 - \sigma(w^\top x_i))$$

# Logistic loss

- Logistic / log / binary cross-entropy loss:

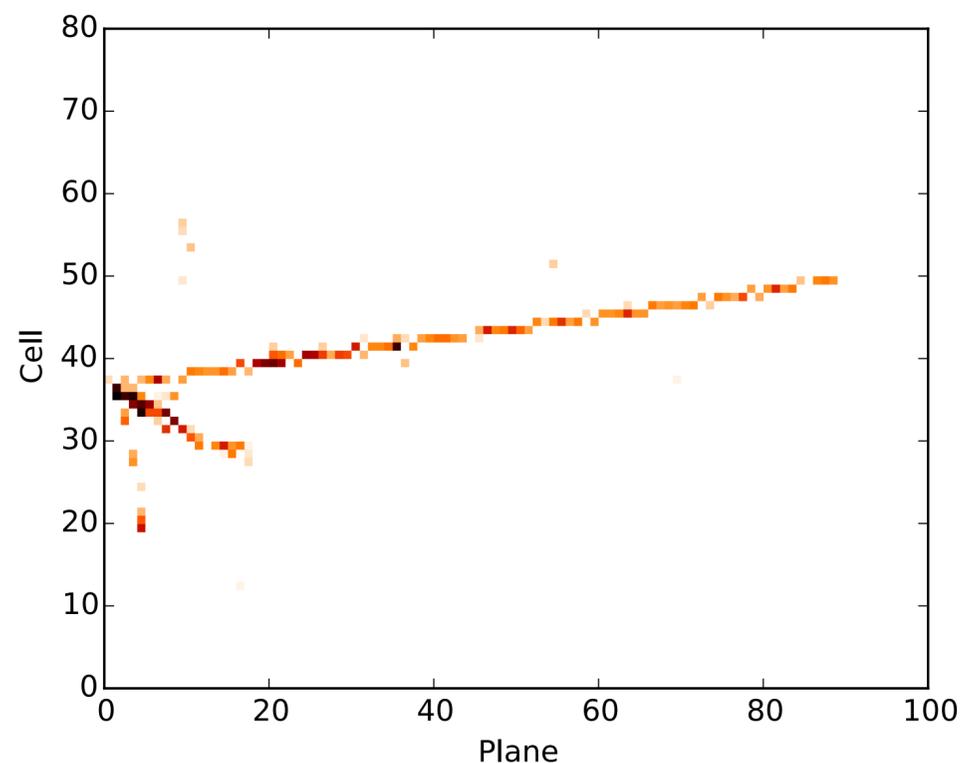
$$L(y, y') = -\delta_{\{y=+1\}} \log y' - \delta_{\{y=-1\}} \log(1 - y')$$





# Multiclass logistic regression

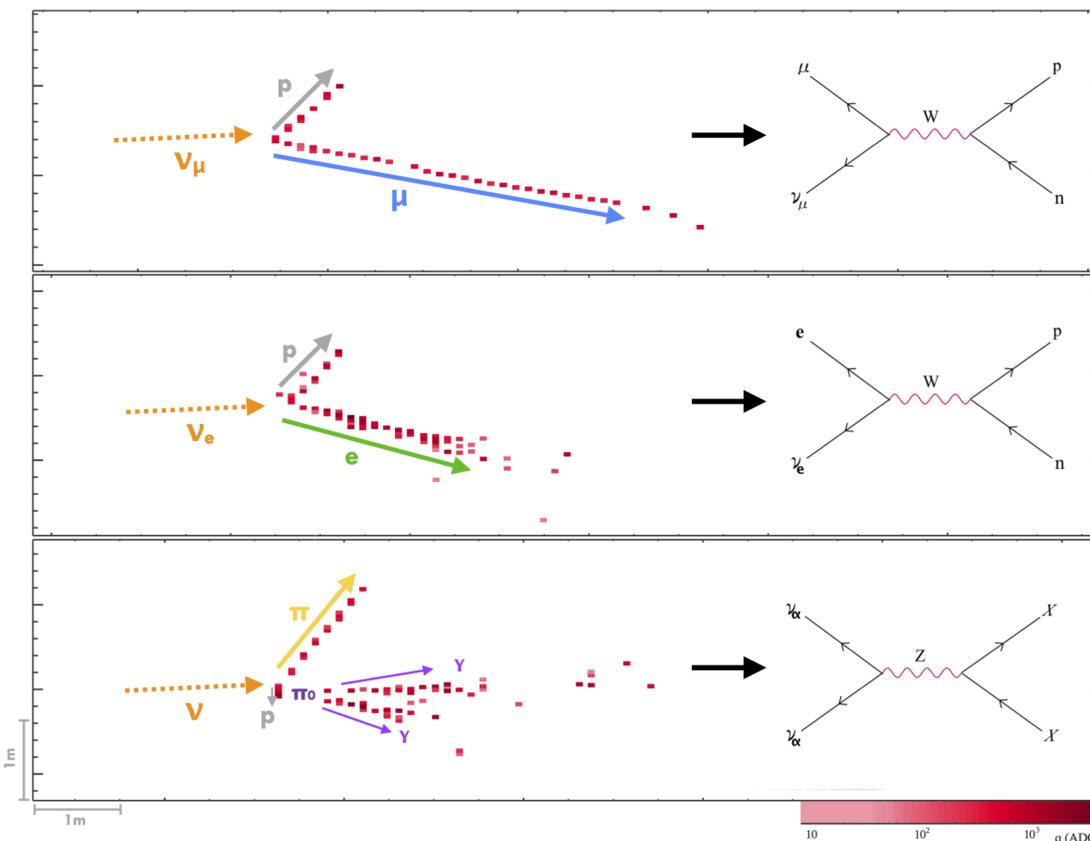
- Predict a raw score for each of  $K$  classes
- Example:  $K = 3, Y = \{\nu_\mu \text{ CC}, \nu_e \text{ CC}, \text{NC}\}$



$$x \in \mathbb{R}^D$$

$$w^T x = \begin{bmatrix} w_1^T x \\ w_2^T x \\ w_3^T x \end{bmatrix} \begin{matrix} \leftarrow \nu_\mu \text{ CC score} \\ \leftarrow \nu_e \text{ CC score} \\ \leftarrow \text{NC score} \end{matrix}$$

Model parameters:  $w \in \mathbb{R}^{K \times D}$

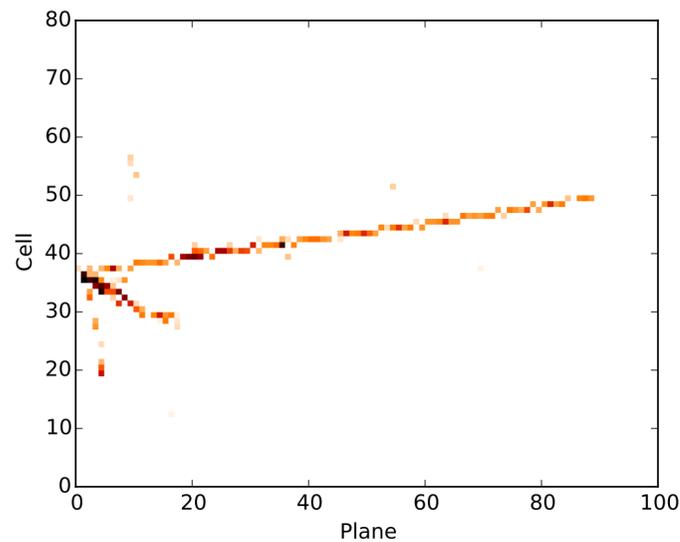


# Multiclass logistic regression

- Sigmoid is replaced by softmax:

$$\text{softmax} \left( \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_K \end{bmatrix} \right) = \frac{1}{\sum_{k=1}^K \exp(a_k)} \begin{bmatrix} \exp(a_1) \\ \exp(a_2) \\ \vdots \\ \exp(a_K) \end{bmatrix}$$

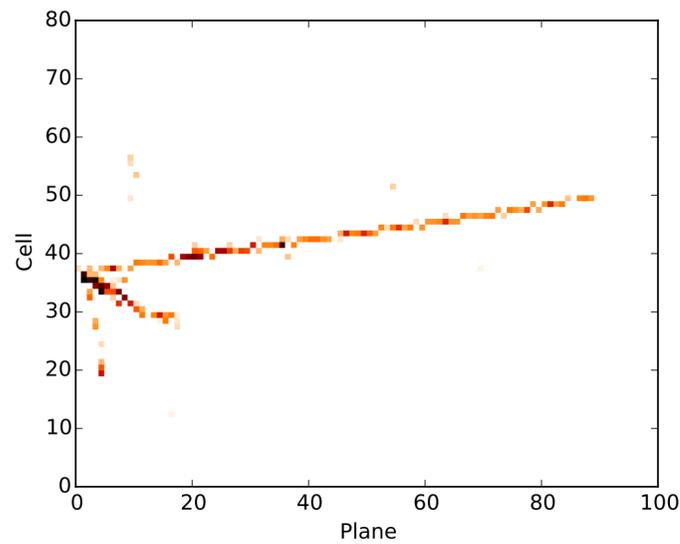
K numbers between 0 and 1 that sum to 1



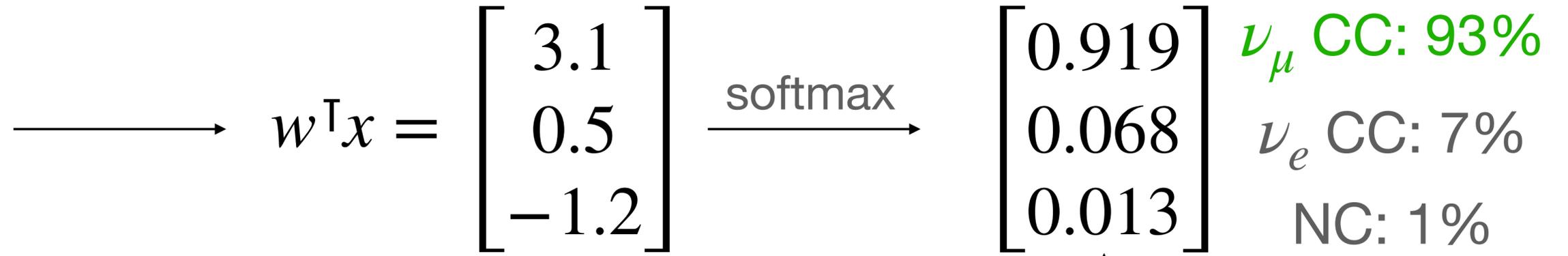
$$\longrightarrow w^T x = \begin{bmatrix} w_1^T x \\ w_2^T x \\ w_3^T x \end{bmatrix} \longrightarrow \text{softmax}(w^T x) = \begin{bmatrix} f_1(x | w) \\ f_2(x | w) \\ f_3(x | w) \end{bmatrix}$$

$$x \in \mathbb{R}^D$$

# Multiclass logistic regression example



$$x \in \mathbb{R}^D$$



softmax outputs  
sum to 1

# Categorical cross-entropy loss

- Negative log likelihood of  $w$  a.k.a. **categorical cross-entropy loss**:

$$-\log p(S | w) = - \sum_{i=1}^N \sum_{k=1}^K \delta_{\{y_i=k\}} \log f_k(x | w)$$

- Generalizes the **binary cross-entropy loss** (and equivalent when  $k = 2$ )

# Recap: Activations and loss functions

- Linear regression:
  - **Activation**: linear; **loss**: mean-squared error
- Binary classification:
  - **Activation**: linear; **loss**: perceptron (PLA)
  - **Activation**: linear; **loss**: hinge (SVM)
  - **Activation**: sigmoid; **loss**: binary cross-entropy
- Multiclass classification:
  - **Activation**: softmax; **loss**: categorical cross-entropy

# Next time

- (Boosted) decision trees
- Tabular data
  - Kaggle Higgs boson classification challenge