

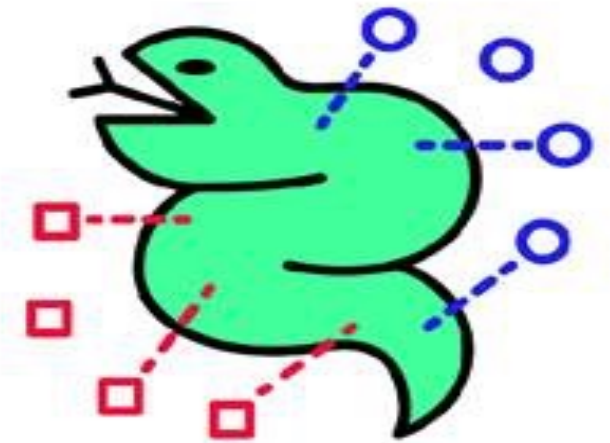
Classification Methods

Alexander Jung

Assistant Professor

Department of Computer Science

Aalto University



Machine Learning
With Python



Classification Problems

labels take on values from a finite set



features: height, weight, passing stat.,
scoring statistics, rebound stat.

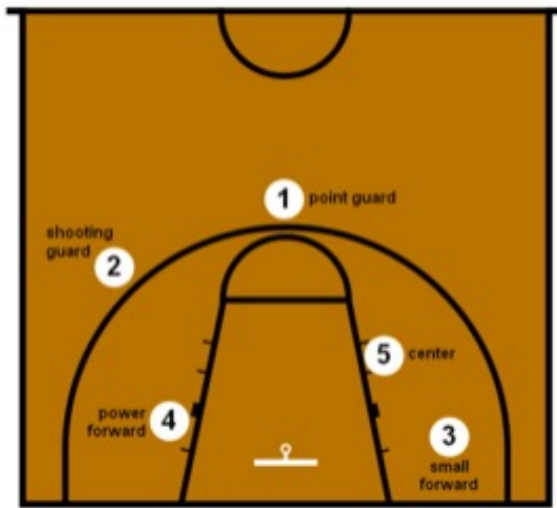
label: position “Shooting Guard”

Finite Set of Label Values

Guards

Forwards

Center



1. **Point Guard** Combo Guard (PG/SG)
2. **Shooting Guard** Swingman (SG/SF)
3. **Small Forward** Point Forward (SF/PG, PF/PG)
4. **Power Forward** Combo Forward (PF/SF)
5. **Center** Forward-Center (PF/C)

label can take on 10 different values

There is No (Obvious) Distance!





true label $y = \text{"Fish"}$

classifier 1 : $\hat{y} = \text{"Dog"}$

classifier 2 : $\hat{y} = \text{"Human"}$

which classifier is "closer"
to true label ?

Binary Classification Problems

- label can take on only **two different values**
- for **notational convenience** we use values +1 and -1
- we could also use **0 and 1** as the two label values
- or we could use  and  as two label values

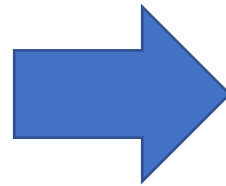
A Binary Classification Problem

Insurance Claim ?????

Name:

Date:

Case:



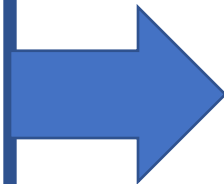
maybe fraudulent

Insurance Claim ?????

Name:

Date:

Case:



most likely not fraudulent

Classifiers

- data point with features x and label y
- we want to learn a classifier map $h(x)$
- $h(x)$ reads in features x and outputs a prediction for label
- choose or learn classifier which best matches true labels

Classifiers and Classification

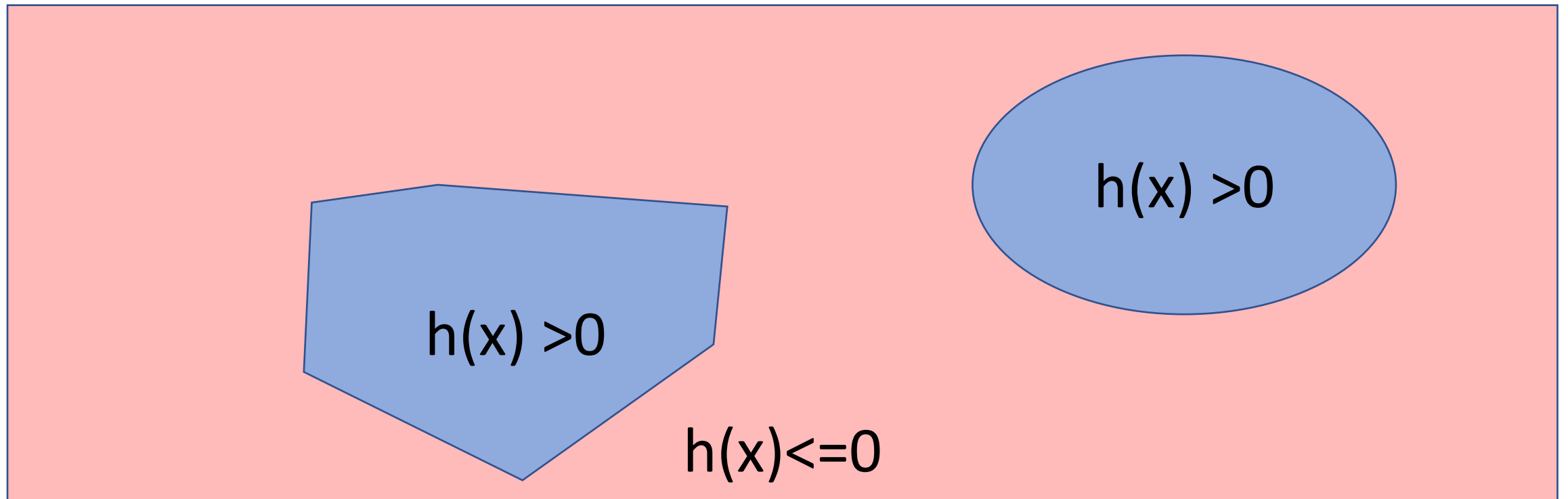
- allow classifier map $h(x)$ to be **real-valued**
- obtain **classification by thresholding**

$$\hat{y} = \begin{cases} 1 & \text{if } h(x) > 0 \\ -1 & \text{otherwise.} \end{cases}$$

- absolute value $|h(x)|$ is **confidence/trust** in classification

Decision Regions

- consider data point with **two features** $x = (x_1, x_2)$
- classifier divides feature space into **two components**



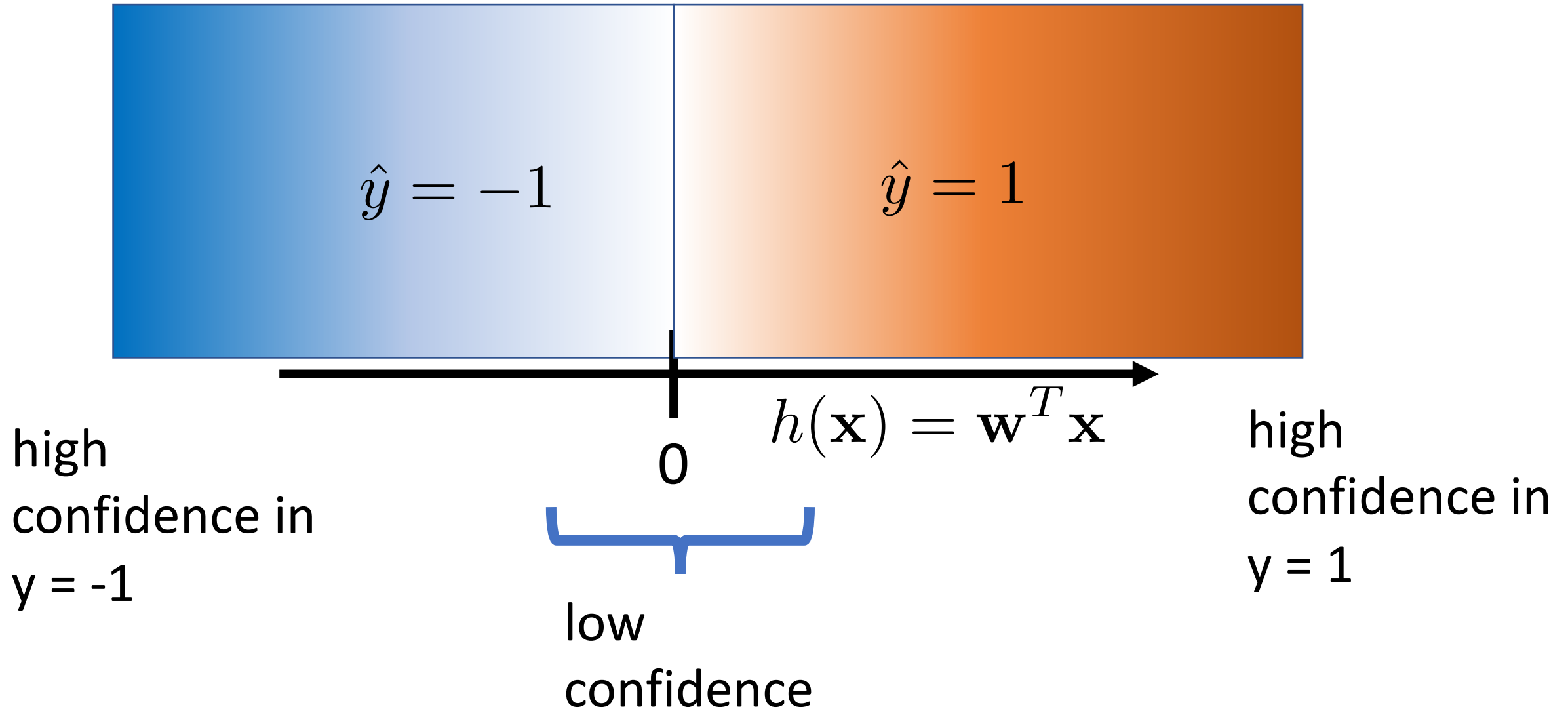
Logistic Regression

- data points with **numeric features** $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n$
- **binary label** $y \in \{-1, 1\}$
- **linear** predictor map $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^n w_i x_i$
- predictor defined by **weights** $\mathbf{w} = (w_1, \dots, w_n)^T$

Logistic Regression – The Classifier

- linear predictor $h(\mathbf{x})$ can be **any real number**
- what if $h(\mathbf{x}) = 0.333$? what if $h(\mathbf{x}) = -100.44$?
- we **use sign** of $h(\mathbf{x})$ **to classify**:
$$\hat{y} = \begin{cases} 1 & \text{if } h(\mathbf{x}) \geq 0 \\ -1 & \text{if } h(\mathbf{x}) < 0 \end{cases}$$
- we use **magnitude $|h(\mathbf{x})|$** as **reliability (confidence)** measure

Logistic Regression – The Classifier



Probabilistic Machine Learning

- interpret label y as realization of **random variable**
- **fraction of data points** with $y=1$ is probability $\text{Prob}(y=1)$
- probability $\text{Prob}(y=1)$ depends on features x !
- logistic regression aims to predict (estimate) $\text{Prob}(y=1)$

Logistic Regression – Probabilistic Model

- interpret label y as realization of **random variable**

- define confidence in $y=1$ via **“log-odds”** $\log \frac{\text{Prob}\{y = 1\}}{\underbrace{(1 - \text{Prob}\{y = 1\})}_{\text{Prob}\{y=-1\}}}$

- since we predict confidence using $h(\mathbf{x})$, we obtain

$$\text{Prob}\{y = 1\} = \frac{1}{1 + \exp(-h(\mathbf{x}))} = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

- probability distribution **parametrized by weights w** of predictor!

Logistic Regression – Towards a Loss Function

- consider **m labeled data points** $(\mathbf{x}^{(i)}, y^{(i)})$, $i = 1, \dots, m$

- labels are realizations of **independent random variables**

$$\text{Prob}\{y^{(i)} = 1\} = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{(i)})}, \text{ and } \text{Prob}\{y^{(i)} = -1\} = \frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})}$$

- probability of actually observing the labels

$$\begin{aligned} \text{Prob}\{y^{(1)}, \dots, y^{(m)}\} &= \prod_{i=1}^m \text{Prob}\{y^{(i)}\} \\ &= \prod_{i=1}^m \frac{1}{1 + \exp(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)})} \end{aligned}$$

Logistic Regression – Maximum Likelihood

- learn weights by **maximizing probability of data !**

$$\begin{aligned}\hat{\mathbf{w}} &= \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^n} \prod_{i=1}^m \frac{1}{1 + \exp(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)})} \\ &= \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^n} \log \prod_{i=1}^m \frac{1}{1 + \exp(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)})} \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} - \log \prod_{i=1}^m \frac{1}{1 + \exp(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)})} \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \sum_{i=1}^m \log (1 + \exp(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}))\end{aligned}$$

Logistic Regression – Logistic Loss

max. probability of data = minimizing the **average logistic loss**

$$\log(1 + \exp(-y\mathbf{w}^T \mathbf{x}))$$

wrong classification $\hat{y} \neq y$

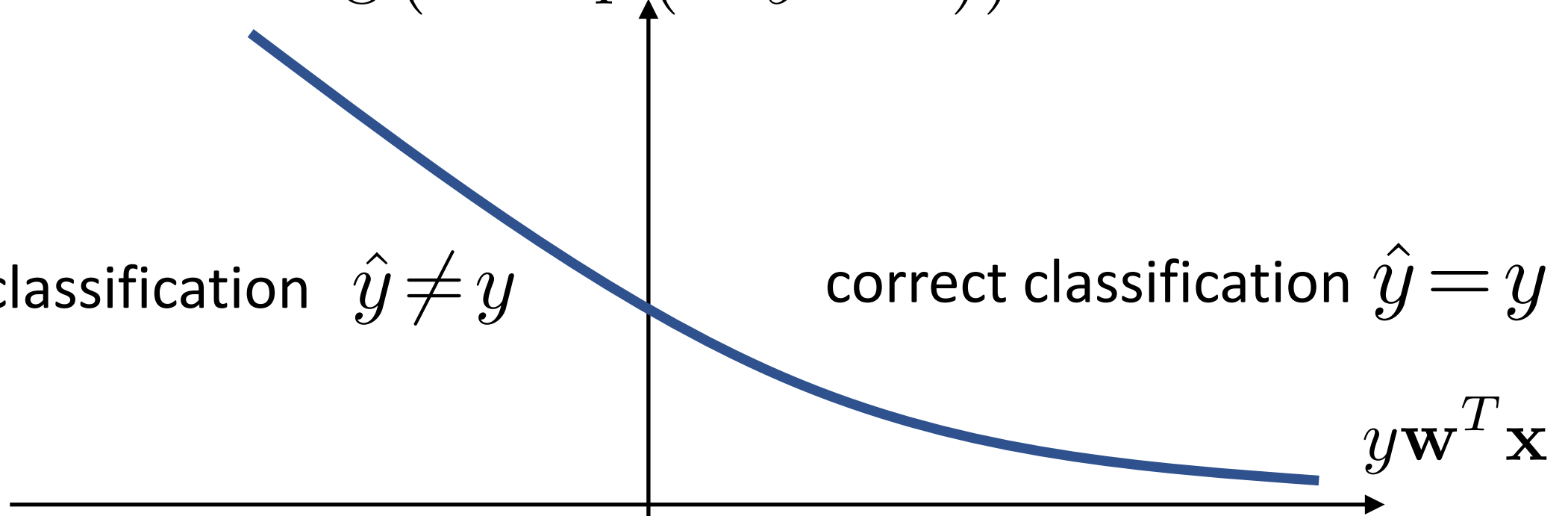
correct classification $\hat{y} = y$

$$y\mathbf{w}^T \mathbf{x}$$

high confidence

low confidence

high confidence



ID-Card of Logistic Regression

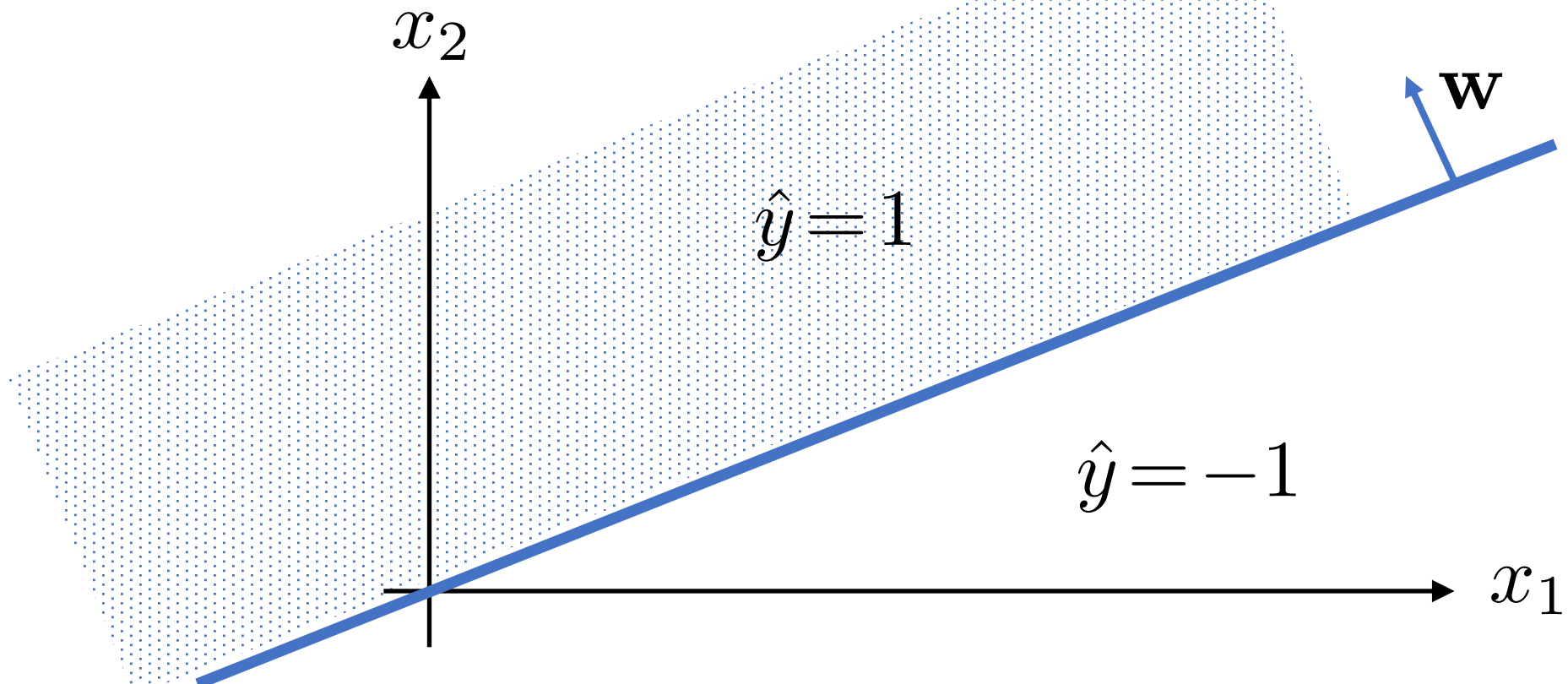
`sklearn.linear_model.LogisticRegression`

- features: **real numbers**
- labels: **two different values** (categories)
- hypothesis space: **linear predictor maps**
- loss: **logistic loss**
- instance of a **linear classifier**

```
linear_model.LogisticRegression(penalty='l2', dual=False, tol=0.0001, C=1.0, fit_intercept=True, g=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto')
```

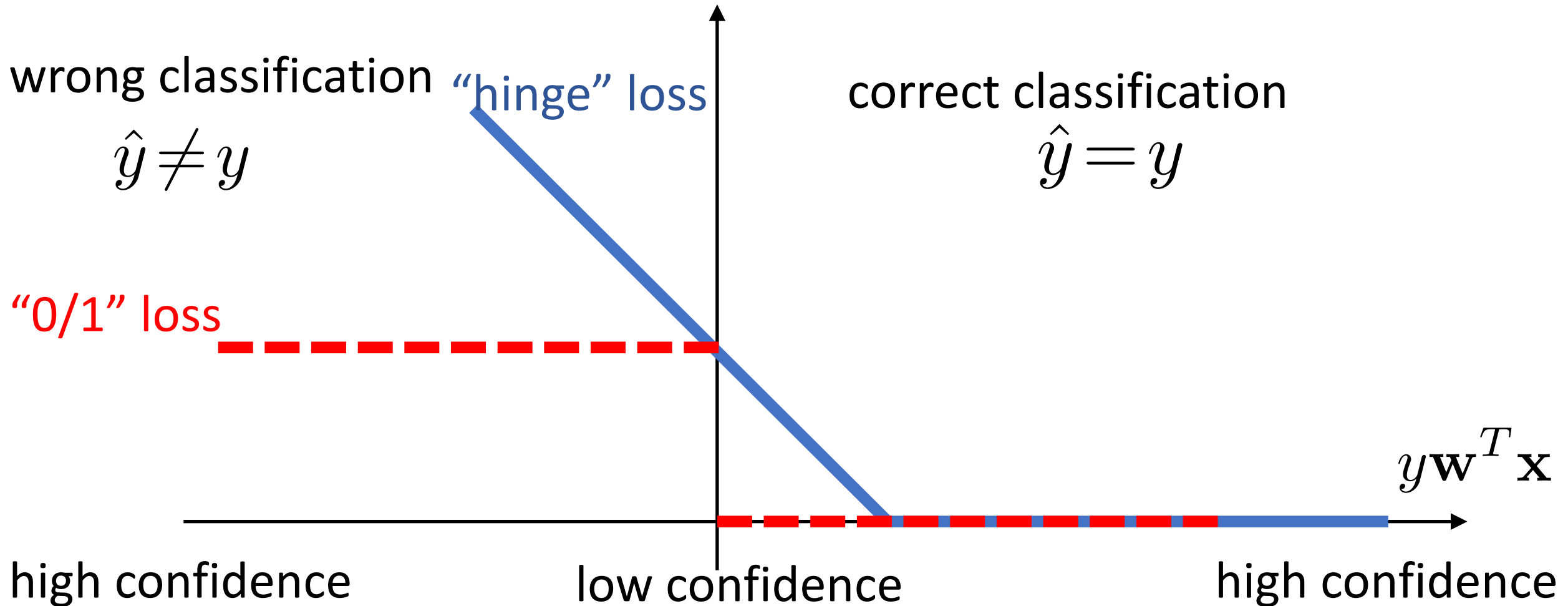
Linear Classifiers

- linear classifiers predict label y using $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- linear classifiers differ in how to learn weights w



Loss Functions for Linear Classifiers

learn linear predictors using **different loss functions**



ID-Card of Support Vector Classifier

`sklearn.svm.SVC`

- features: **real numbers**
- labels: **two different values** (categories)
- hypothesis space: **linear predictor maps**
- loss: **hinge loss**
- instance of a **linear classifier**

rbf', degree=3, gamma='scale', coef0=0.0, shrinkage=0.001, cache_size=200, class_weight=None, decision_function_shape='raw', decision_threshold=None, verbose=False, max_iter=-1, decision_

ID-Card of Naive Bayes Classifier

- features: **real numbers**
- labels: **two different values** (categories)
- hypothesis space: **linear predictor maps**
- loss: **0/1 loss**
- instance of a **linear classifier**

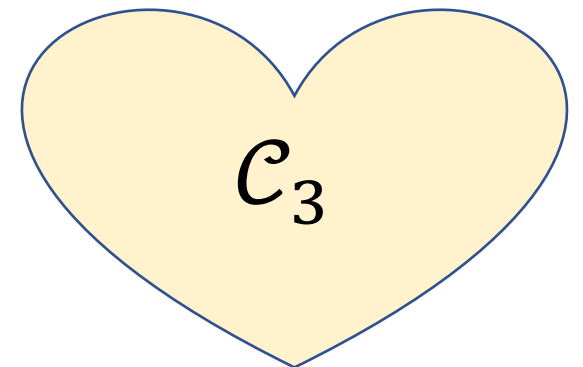
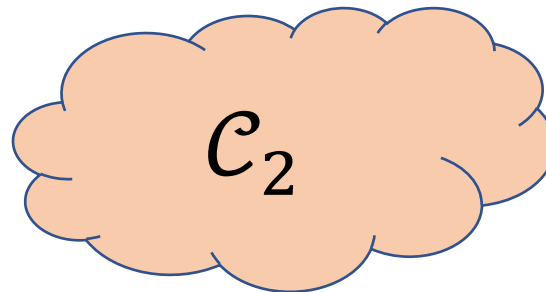
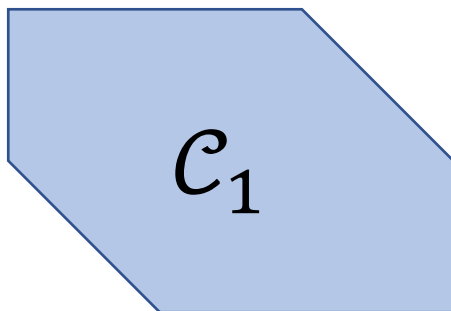
```
sklearn.naive_bayes.GaussianNB
```

```
e_bayes.GaussianNB(priors=None, var_smoothing=1e-09)
```

You Can Do Anything with Linear Classifiers!

- data points with **two numeric features** z_1 and z_2
- **some (complicated) regions** \mathcal{C}_j , $j=1,\dots,n$, of feature space
- **construct new features** of data point by

$$x_j = \begin{cases} 1 & \text{when } (z_1, z_2) \in \mathcal{C}_j \\ 0 & \text{otherwise.} \end{cases}$$

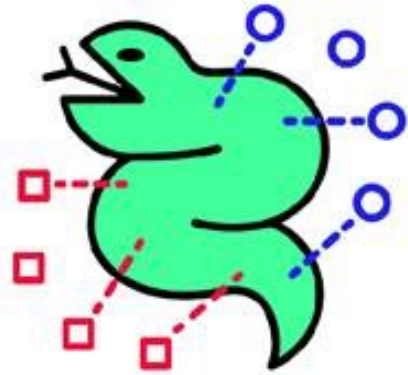


So What?

- logistic regression uses linear predictors
- **predictors** are interpreted as “log-odds”
- predictor **weights** parametrize the probability of $y=1$
- learning weights is a **statistical inference problem**
- inference problem equivalent to **minimizing logistic loss**

So What?

- logistic regression one instance of **linear classifier**
- other linear classifiers obtained from **different loss**
- **support vector classifier** uses **hinge loss**
- **Naïve Bayes** classifier based on **0/1 loss**
- loss functions differ in **statistical and computational properties**



Machine Learning
With Python



Thank You !