

CMPSC 448: Machine Learning

Lecture 10. Logistic Regression

Rui Zhang
Fall 2021



Recap: linear models for binary classification

A linear classifier is specified by a weight vector $\mathbf{w} \in \mathbb{R}^d$

$$f_{\mathbf{w}}(\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^\top \mathbf{x} = \sum_{i=1}^d w_i x_i \geq 0 \\ -1 & \text{if } \mathbf{w}^\top \mathbf{x} = \sum_{i=1}^d w_i x_i < 0 \end{cases}$$

It will be notationally simpler to use $\{+1, -1\}$ instead of $\{1, 0\}$

Recap: linear models for binary classification

A linear classifier is specified by a weight vector $\mathbf{w} \in \mathbb{R}^d$

$$f_{\mathbf{w}}(\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^\top \mathbf{x} = \sum_{i=1}^d w_i x_i \geq 0 \\ -1 & \text{if } \mathbf{w}^\top \mathbf{x} = \sum_{i=1}^d w_i x_i < 0 \end{cases}$$

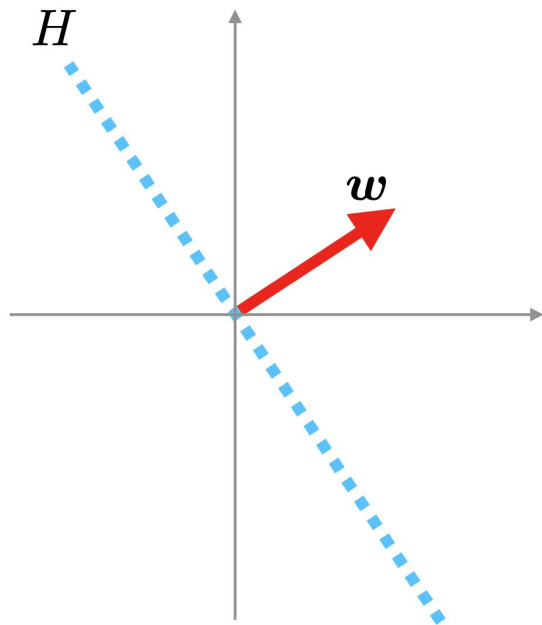
It will be notationally simpler to use $\{+1, -1\}$ instead of $\{1, 0\}$

Interpretation: does a linear combination of input features exceed 0?

For $\mathbf{w} = [w_1, w_2, \dots, w_d]^\top$ and $\mathbf{x} = [x_1, x_2, \dots, x_d]^\top$

$$\mathbf{w}^\top \mathbf{x} = \sum_{i=1}^d w_i x_i >^? 0$$

Geometry of linear models for classification

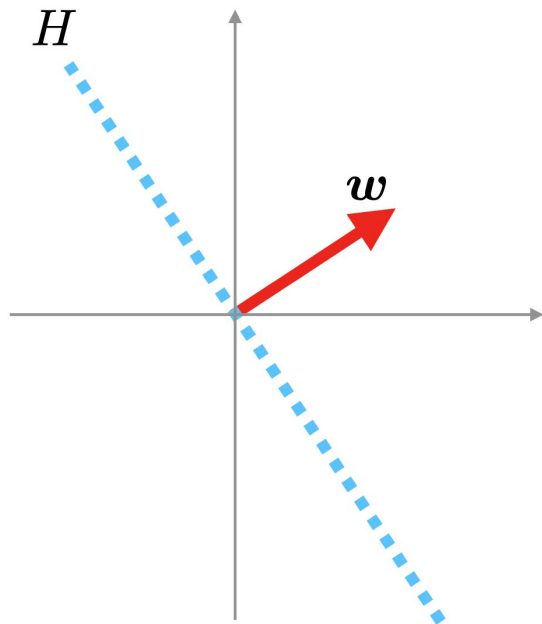


A hyperplane in \mathbb{R}^d is a linear subspace:

A \mathbb{R}^2 -hyperplane is a line

A \mathbb{R}^3 -hyperplane is a plane

Geometry of linear models for classification



A hyperplane in \mathbb{R}^d is a linear subspace:

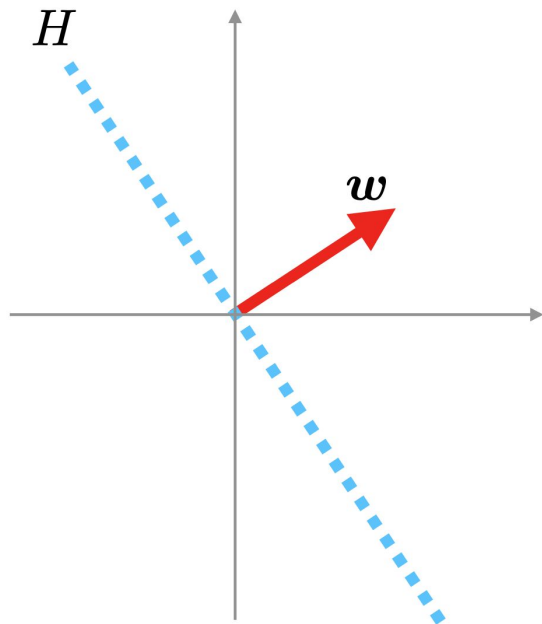
A \mathbb{R}^2 -hyperplane is a line

A \mathbb{R}^3 -hyperplane is a plane

A hyperplane can be specified by a (nonzero) normal vector

$$\mathbf{w} \in \mathbb{R}^d, \|\mathbf{w}\|_2 = 1$$

Geometry of linear models for classification



A hyperplane in \mathbb{R}^d is a linear subspace:

A \mathbb{R}^2 -hyperplane is a line

A \mathbb{R}^3 -hyperplane is a plane

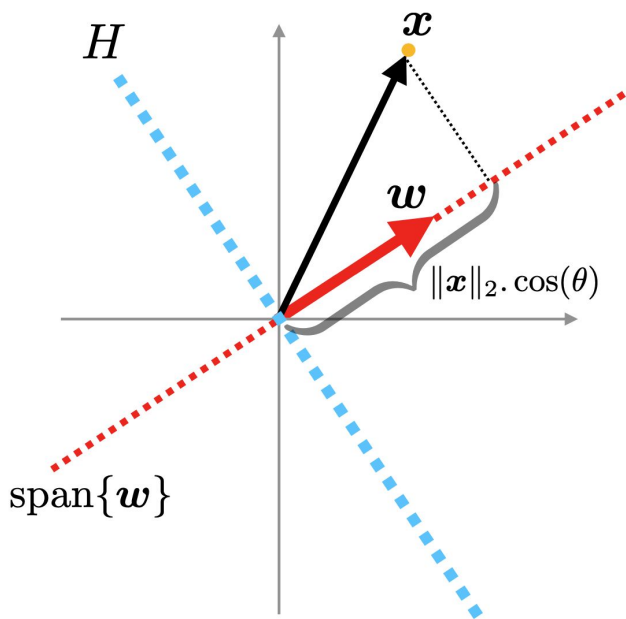
A hyperplane can be specified by a (nonzero) normal vector

$$\mathbf{w} \in \mathbb{R}^d, \|\mathbf{w}\|_2 = 1$$

The hyperplane with normal vector can be represented by set of points orthogonal to its normal vector:

$$H = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = 0\}$$

Classification using linear models

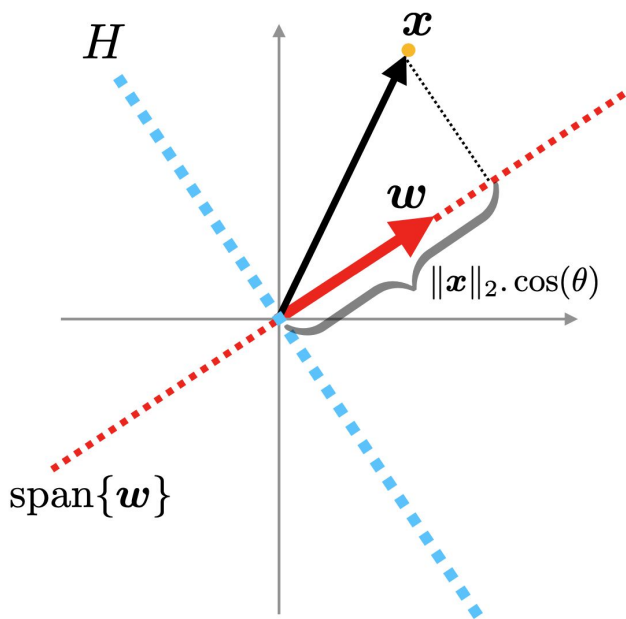


- Projection of x onto $\text{span}\{w\}$ (a line) has coordinate $\|x\|_2 \cos(\theta)$

where

$$\cos(\theta) = \frac{w^\top x}{\|w\|_2 \|x\|_2}$$

Classification using linear models



- Projection of \mathbf{x} onto $\text{span}\{\mathbf{w}\}$ (a line) has coordinate $\|\mathbf{x}\|_2 \cos(\theta)$

where

$$\cos(\theta) = \frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|_2 \|\mathbf{x}\|_2}$$

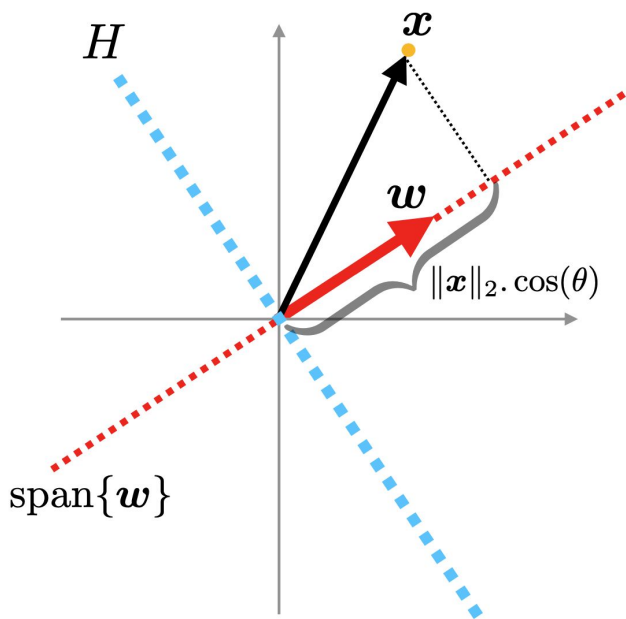
- Decision boundary is hyperplane oriented by \mathbf{w}

$$\mathbf{w}^\top \mathbf{x} > 0$$

$$\iff \|\mathbf{x}\|_2 \cos(\theta) > 0$$

$$\iff \mathbf{x} \text{ on same side of } H \text{ as } \mathbf{w}$$

Classification using linear models



- Projection of \mathbf{x} onto $\text{span}\{\mathbf{w}\}$ (a line) has coordinate $\|\mathbf{x}\|_2 \cos(\theta)$

where

$$\cos(\theta) = \frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|_2 \|\mathbf{x}\|_2}$$

- Decision boundary is hyperplane oriented by \mathbf{w}

$$\mathbf{w}^\top \mathbf{x} > 0$$

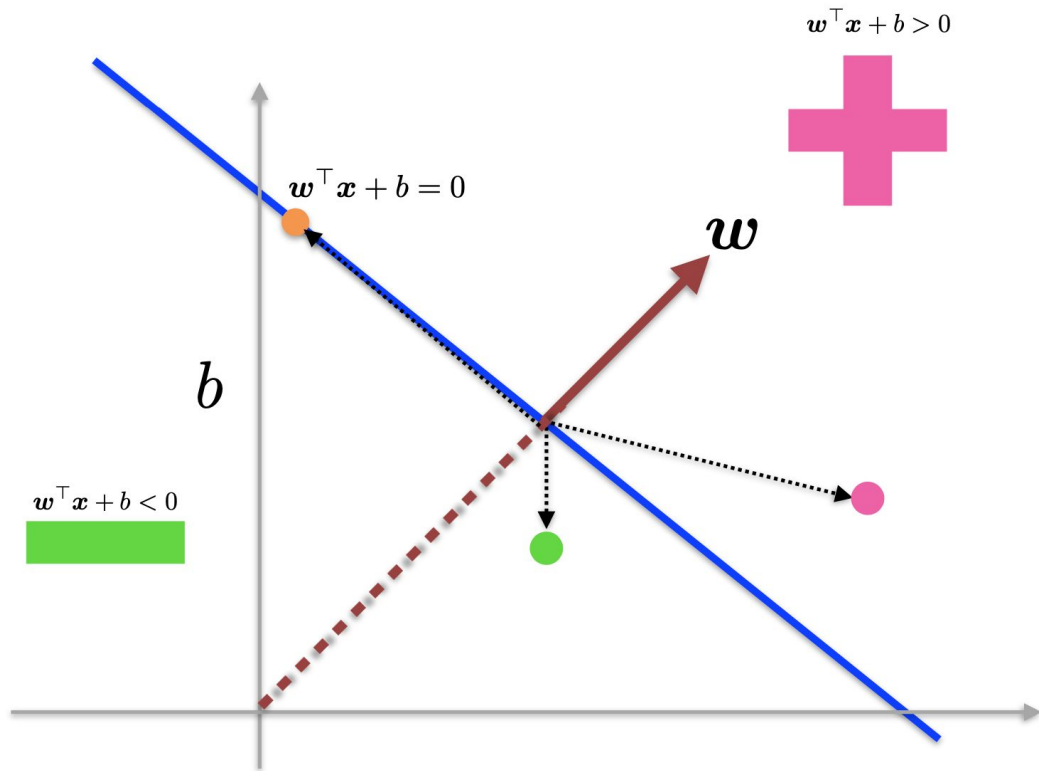
$$\iff \|\mathbf{x}\|_2 \cos(\theta) > 0$$

$$\iff \mathbf{x} \text{ on same side of } H \text{ as } \mathbf{w}$$

Q: What should we do if we want hyperplane decision boundary that does not (necessarily) go through origin?
A: Add bias.

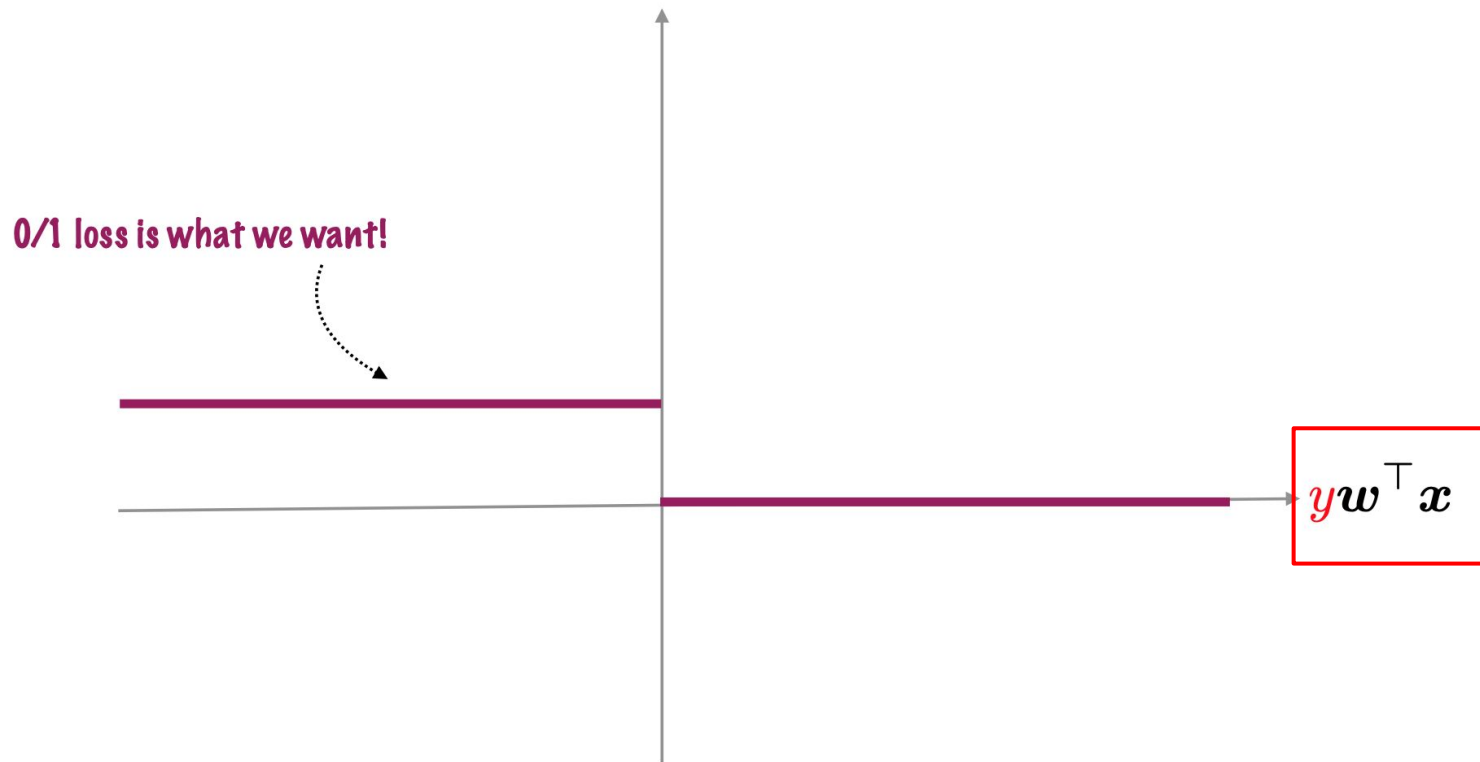
Linear models and hyperplanes

Every hyperplane is a decision boundary that splits the space into two subspaces:



0/1 loss for classification

Combining two cases results in the 0/1 loss function for classification:



Perceptron as an artificial neuron

Perceptron is an **algorithm** to **learn** linear models!



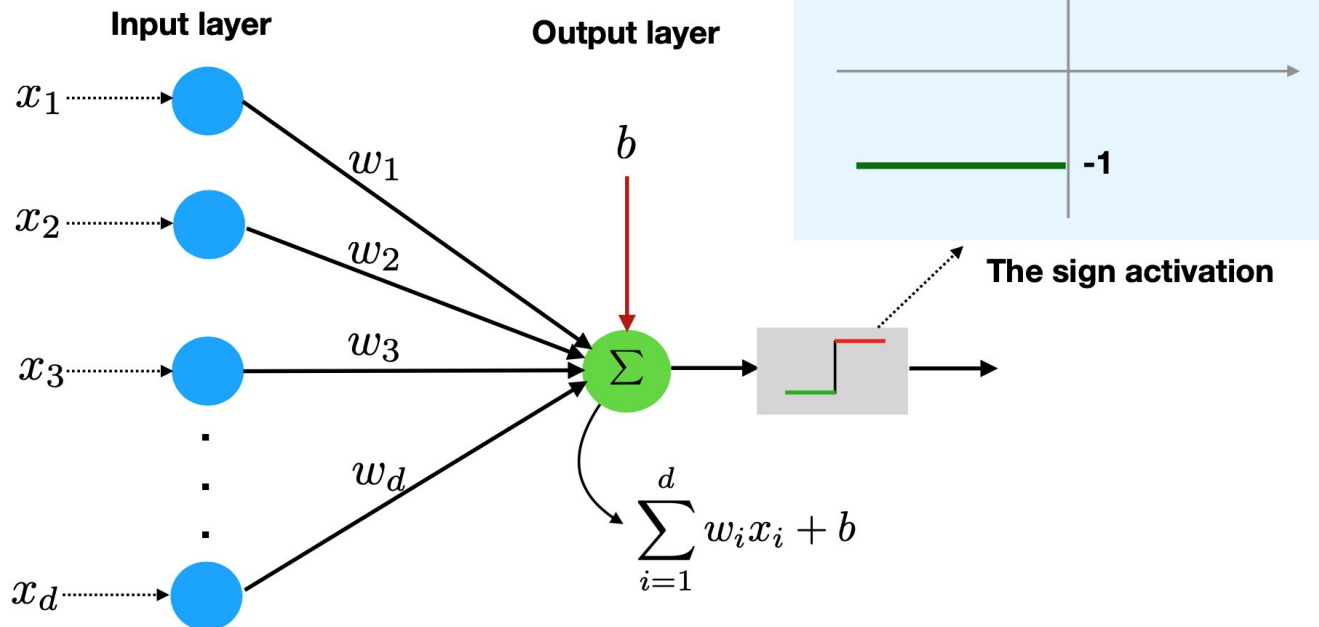
The sign activation

$$\sum_{i=1}^d w_i x_i + b$$

It only works if the data is linearly separable (an assumption that does not hold for many real data)

Perceptron as an artificial neuron

Perceptron is an **algorithm** to **learn** linear models!



It only works if the data is linearly separable (an assumption that does not hold for many real data)

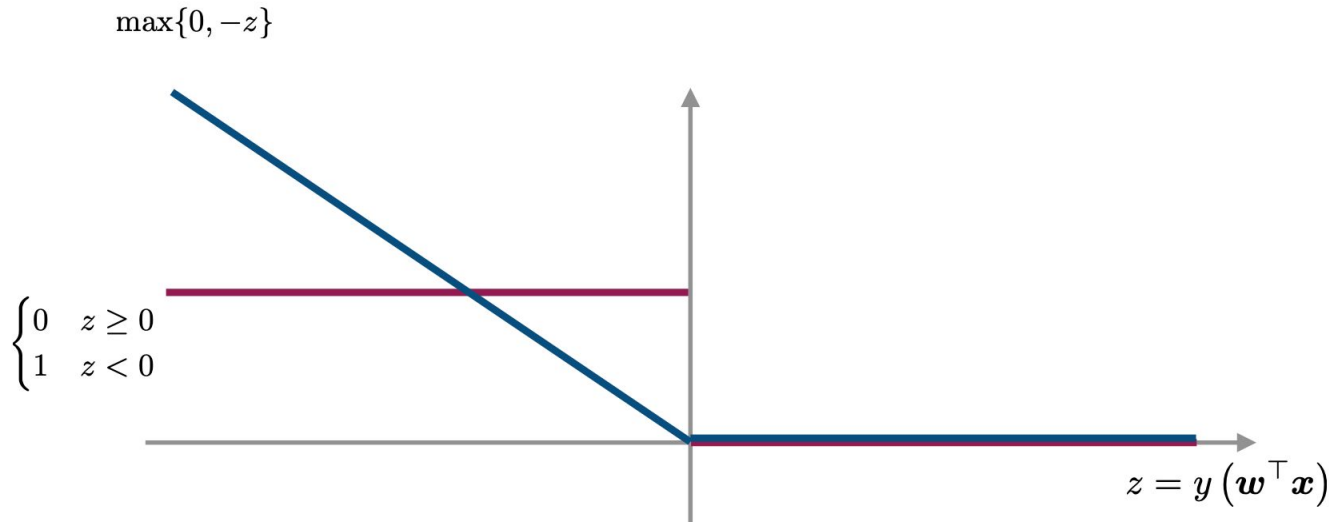
Convex approximation of 0/1 loss

- We want 0/1 loss for classification
- Unfortunately 0/1 loss function is hard to optimize
 - is a non-convex function
 - The gradient is zero
- Let's replace 0/1 loss with approximate convex functions!

The Perceptron loss versus 0/1 loss

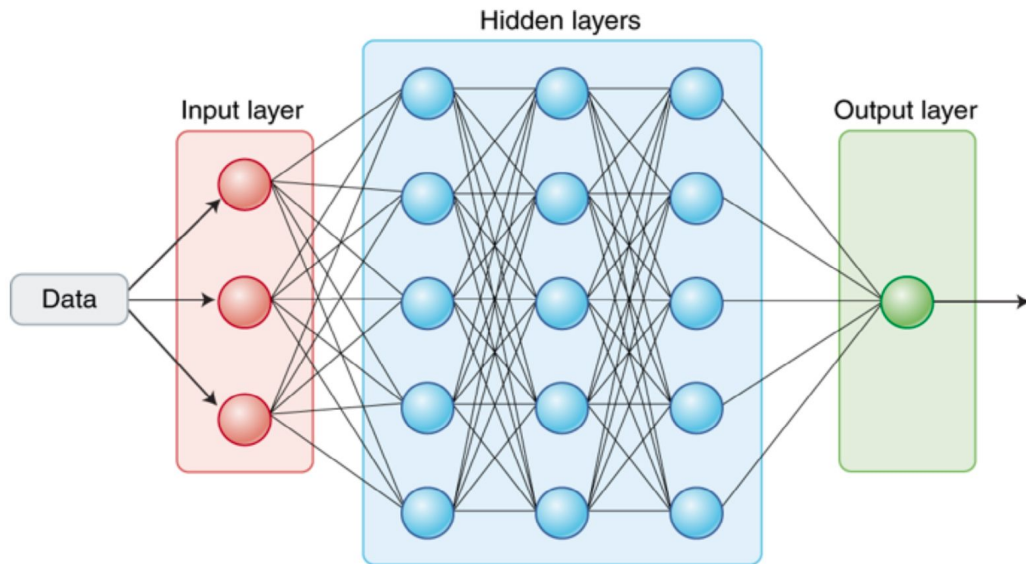
The Perceptron algorithm is optimizing the following loss function in training:

$$\max\{0, -y (\mathbf{w}^\top \mathbf{x})\}$$



Computational virtues (convex versus non-convex)

Multiple layer perceptron



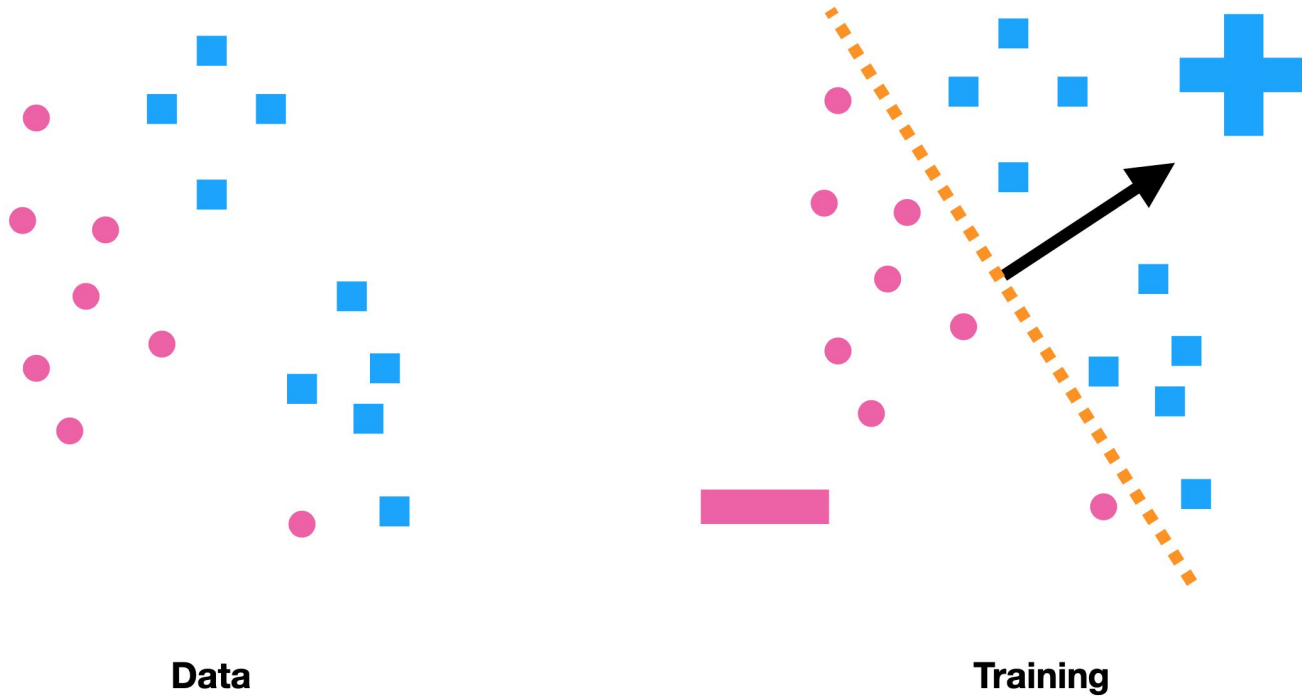
Training: use back-propagation (gradient descent + chain rule) to learn the parameters for all layers

Introduces new challenges: huge amount of data + hyperparameter + vanishing gradient issue + high-computational power + more effective regularization + etc

Outline

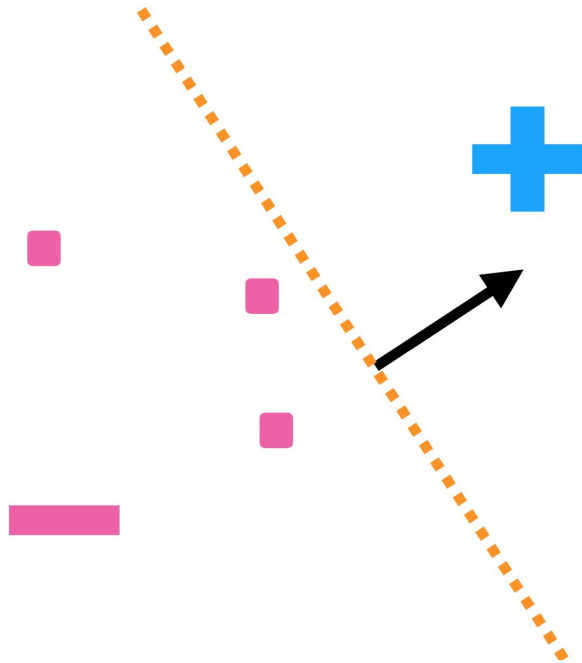
- Predicting probability
- Logistic regression for binary classification
- MLE for Logistic Regression
- The cross-entropy loss function
- Regularized logistic regression
- The decision boundary of logistic regression

The lack of confidence



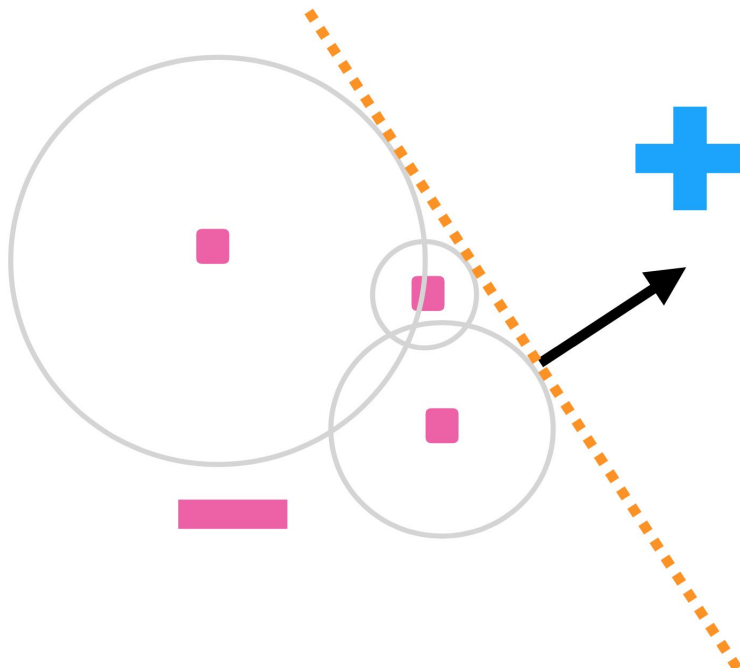
Assume we learn a linear classifier for data given here.

The lack of confidence



Three test points are given. All will be classified as -1.
But, which one we are more confident in its label?
Is there any measure of this confidence/probability?

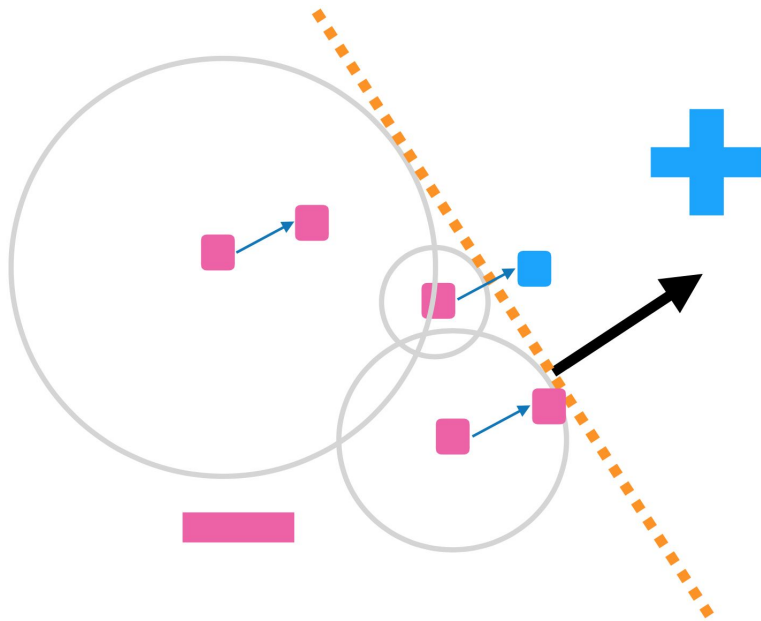
The measure of confidence



The larger the margin (distance from hyperplane), more confidence the algorithm will be in predicted label, which can be computed as

$$y(\mathbf{w}^\top \mathbf{x})$$

The measure of confidence



Noise tolerance: the amount of noise we can add to each data point such that the label stays same!

From score to probability



From score to probability



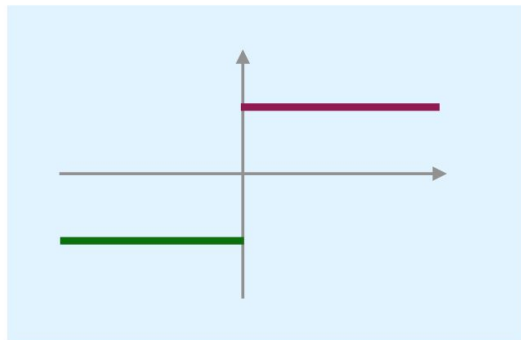
$w^\top x \longrightarrow$

$\longrightarrow \mathbb{P}[y = +1]$

$\mathbb{P}[y = -1] = 1 - \mathbb{P}[y = +1]$

From score to probability

$w^\top x \rightarrow$



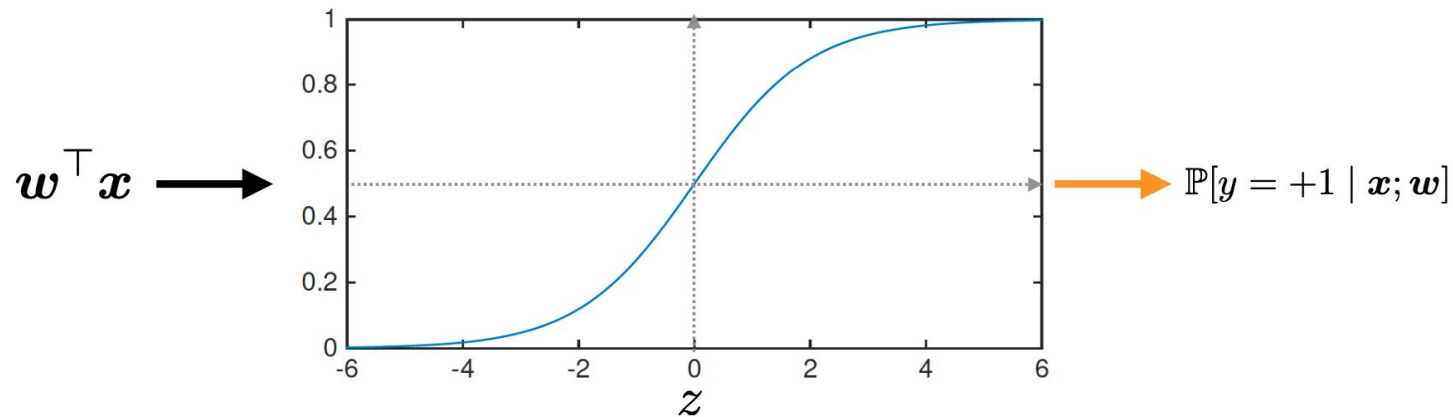
$\rightarrow y = \begin{cases} +1 & \text{if } w^\top x \geq 0 \\ -1 & \text{if } w^\top x < 0 \end{cases}$

$w^\top x \rightarrow$



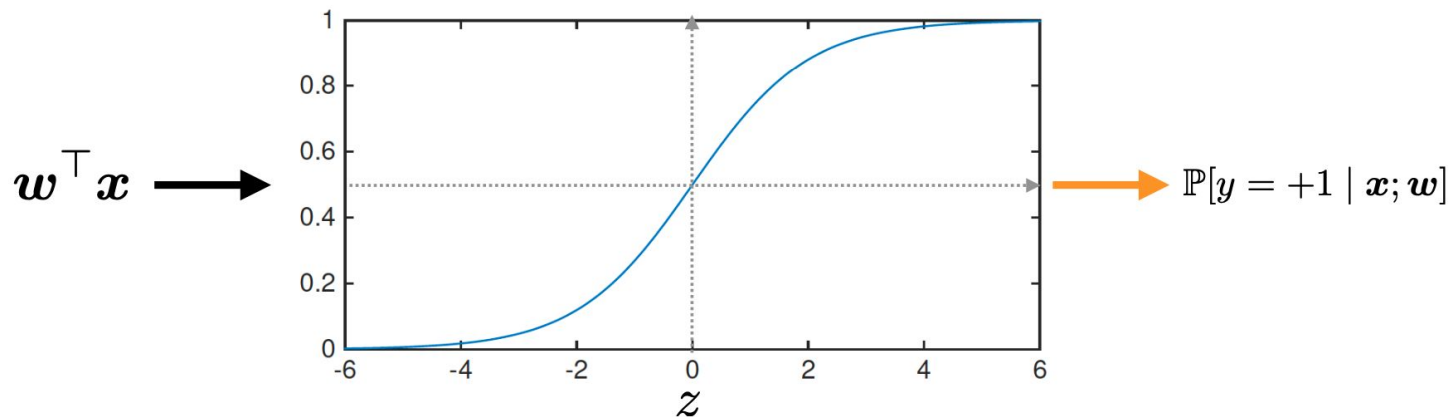
$\rightarrow \mathbb{P}[y = +1]$
 $\mathbb{P}[y = -1] = 1 - \mathbb{P}[y = +1]$

From score to probability



Logistic Function (aka Sigmoid Function)

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

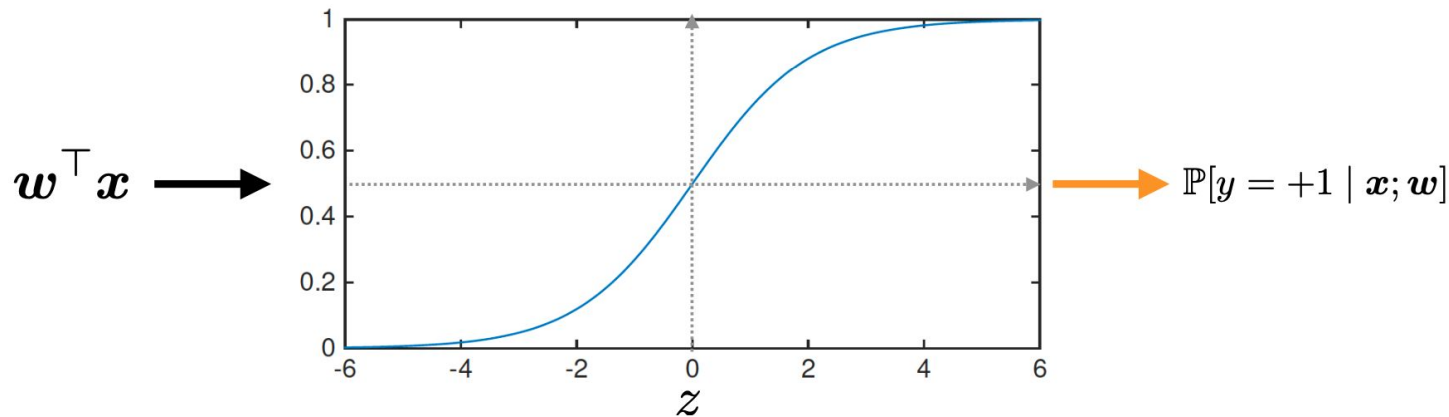


$$\begin{aligned}\mathbb{P}[y = +1 | \mathbf{x}; \mathbf{w}] \\ &= \sigma(\mathbf{w}^\top \mathbf{x}) \\ &= \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}\end{aligned}$$

$$\begin{aligned}\mathbb{P}[y = -1 | \mathbf{x}; \mathbf{w}] \\ &= 1 - \mathbb{P}[y = +1 | \mathbf{x}; \mathbf{w}] \\ &= \sigma(-\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^\top \mathbf{x}}}\end{aligned}$$

Logistic Function (aka Sigmoid Function)

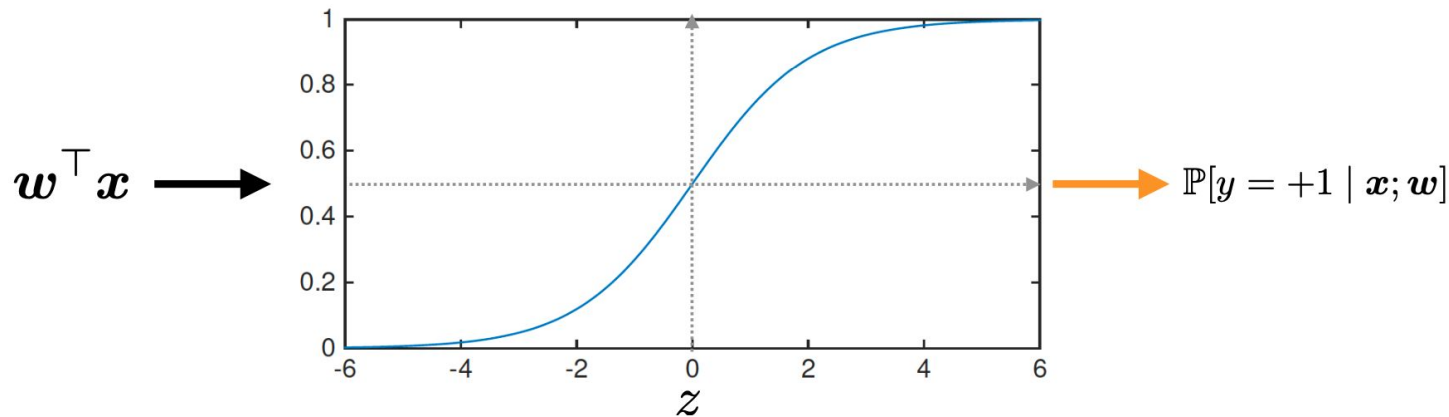
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$$\mathbb{P}[y \mid \mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-y\mathbf{w}^\top \mathbf{x}}}$$

Logistic Function (aka Sigmoid Function)

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



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

This “squashing function” known as the **logistic function** (or **sigmoid function**) turns linear prediction score into probability!

Logistic Regression

- Given a data for binary classification: $\mathcal{S} = \{(\mathbf{x}_1, y_1 = \pm 1), (\mathbf{x}_2, y_2 = \pm 1), \dots, (\mathbf{x}_n, y_n = \pm 1)\}$
- The conditional probability of the label given the input example \mathbf{x} for a classifier with parameters vector \mathbf{w} is expressed as:

$$\mathbb{P}[y|\mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-\mathbf{y}\mathbf{w}^\top \mathbf{x}}}$$

where the $\mathbf{w} = [w_1, w_2, \dots, w_d]^\top$ are the adjustable parameters.

Logistic Regression

- Given a data for binary classification: $\mathcal{S} = \{(\mathbf{x}_1, y_1 = \pm 1), (\mathbf{x}_2, y_2 = \pm 1), \dots, (\mathbf{x}_n, y_n = \pm 1)\}$
- The conditional probability of the label given the input example \mathbf{x} for a classifier with parameters vector \mathbf{w} is expressed as:

$$\mathbb{P}[y|\mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-y\mathbf{w}^\top \mathbf{x}}}$$

where the $\mathbf{w} = [w_1, w_2, \dots, w_d]^\top$ are the adjustable parameters.

- If the probability of y being 1 is greater than 0.5, then we can predict $y = +1$
- This is the "logistic function" and model is called "logistic regression".

Logistic Regression

- Given a data for binary classification: $\mathcal{S} = \{(\mathbf{x}_1, y_1 = \pm 1), (\mathbf{x}_2, y_2 = \pm 1), \dots, (\mathbf{x}_n, y_n = \pm 1)\}$
- The conditional probability of the label given the input example \mathbf{x} for a classifier with parameters vector \mathbf{w} is expressed as:

$$\mathbb{P}[y|\mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-y\mathbf{w}^\top \mathbf{x}}}$$

where the $\mathbf{w} = [w_1, w_2, \dots, w_d]^\top$ are the adjustable parameters.

- If the probability of y being 1 is greater than 0.5, then we can predict $y = +1$
- This is the "logistic function" and model is called "logistic regression".

In the terminology of statistics, this model is known as logistic "regression", although it should be emphasized that this is a model for classification rather than regression.

MLE for logistic regression

We can fit the logistic models using the maximum likelihood estimation (MLE).

$$\mathbb{P}[y|\mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-\mathbf{y}\mathbf{w}^\top \mathbf{x}}}$$

Recall training data are independently generated (i.i.d.), so the likelihood is:

$$\mathbb{P}[y_1, y_2, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}] = \prod_{i=1}^n \mathbb{P}[y_i | \mathbf{x}_i; \mathbf{w}]$$

MLE for logistic regression

We can fit the logistic models using the maximum likelihood estimation (MLE).

$$\mathbb{P}[y|\mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-\mathbf{y}\mathbf{w}^\top \mathbf{x}}}$$

Recall training data are independently generated (i.i.d.), so the likelihood is:

$$\mathbb{P}[y_1, y_2, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}] = \prod_{i=1}^n \mathbb{P}[y_i | \mathbf{x}_i; \mathbf{w}]$$

The likelihood measures the probability that the data were generated if $\mathbf{w} \in \mathbb{R}^d$ was the model to generate the labels in training data!

MLE for logistic regression

We can fit the logistic models using the maximum likelihood estimation (MLE).

$$\mathbb{P}[y|\mathbf{x}; \mathbf{w}] = \frac{1}{1 + e^{-\mathbf{y}\mathbf{w}^\top \mathbf{x}}}$$

Recall training data are independently generated (i.i.d.), so the likelihood is:

$$\mathbb{P}[y_1, y_2, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}] = \prod_{i=1}^n \mathbb{P}[y_i | \mathbf{x}_i; \mathbf{w}]$$

The likelihood measures the probability that the data were generated if $\mathbf{w} \in \mathbb{R}^d$ was the model to generate the labels in training data!

By maximizing the likelihood over all parameter vectors $\mathbf{w} \in \mathbb{R}^d$ we can find the optimal classifier:

$$\arg \max_{\mathbf{w}} \prod_{i=1}^n \mathbb{P}[y_i | \mathbf{x}_i; \mathbf{w}]$$

Fitting logistic regression models using MLE

$$\begin{aligned}\log \left(\prod_{i=1}^n \mathbb{P} [y_i | \mathbf{x}_i; \mathbf{w}] \right) &= \sum_{i=1}^n \log \mathbb{P} [y_i | \mathbf{x}_i; \mathbf{w}] \\ &= \sum_{i=1}^n \log \frac{1}{1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i}} \\ \log \frac{1}{z} &= -\log z \\ &= \sum_{i=1}^n -\log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right) \\ &= -\sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right)\end{aligned}$$

Fitting logistic regression models using MLE

$$\begin{aligned}\log \left(\prod_{i=1}^n \mathbb{P}[y_i | \mathbf{x}_i; \mathbf{w}] \right) &= \sum_{i=1}^n \log \mathbb{P}[y_i | \mathbf{x}_i; \mathbf{w}] \\ &= \sum_{i=1}^n \log \frac{1}{1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i}} \\ \log \frac{1}{z} &= -\log z \\ &= \sum_{i=1}^n -\log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right) \\ &= -\sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right)\end{aligned}$$

The final MLE objective becomes (or minimizing the negative likelihood):

$$\arg \max_{\mathbf{w} \in \mathbb{R}^d} \quad - \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right) \qquad \arg \min_{\mathbf{w} \in \mathbb{R}^d} \quad \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right)$$

The cross-entropy loss function

The cross-entropy function to measure the error of model on a training sample:

$$\log \left(1 + e^{-y\mathbf{w}^\top \mathbf{x}} \right)$$

The cross-entropy loss function

The cross-entropy function to measure the error of model on a training sample:

$$\log \left(1 + e^{-y\mathbf{w}^\top \mathbf{x}} \right)$$

It looks complicated and ugly, but,

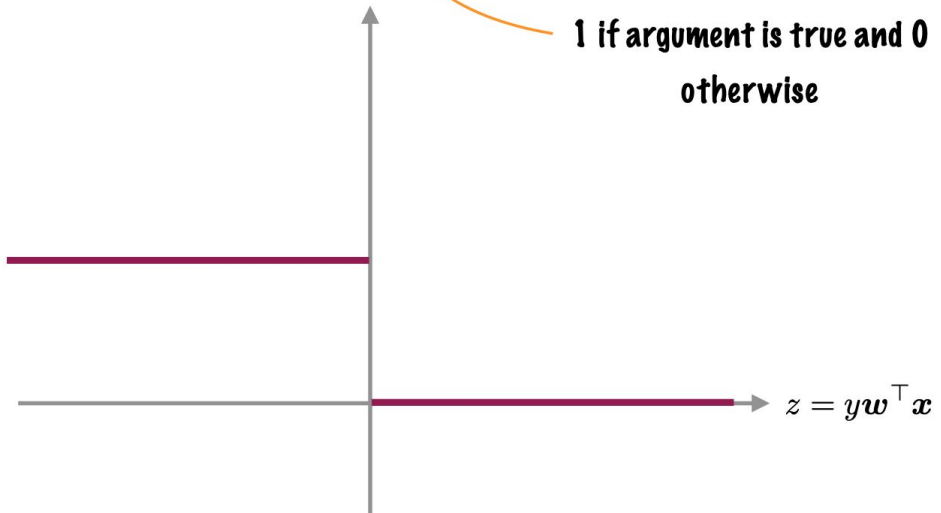
- It is based on an intuitive probabilistic interpretation.
- It is convex.
- However, unlike linear regression, there is no closed-form solution (even though it is convex).
- It is very convenient and mathematically easy to minimize using Gradient Descent (GD).

Optimization viewpoint of linear classifiers

The original goal of linear binary classification is to find the parameters of a linear model that minimizes the number of mistakes over training data (0/1 loss):

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \sum_{i=1}^n \mathbb{I} [\text{sign}(\mathbf{w}^\top \mathbf{x} + b) \neq y_i]$$

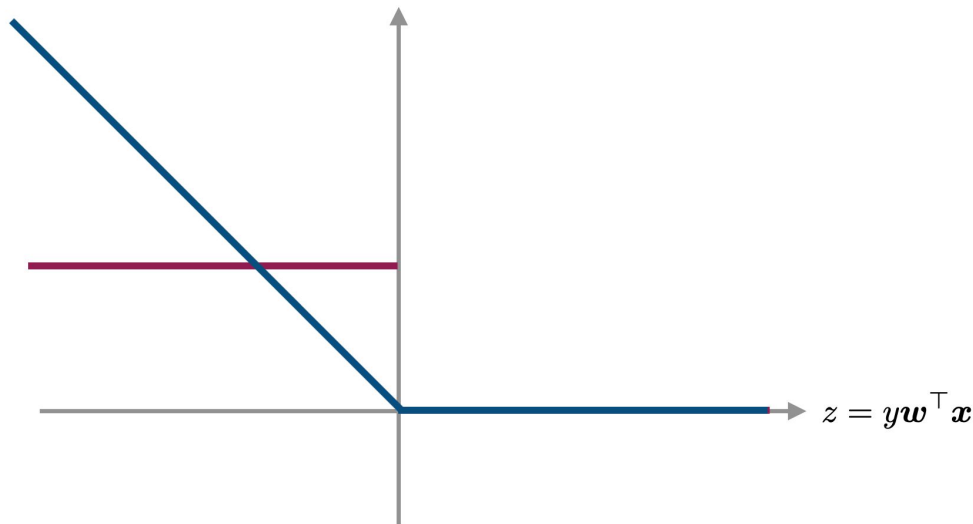
Indicator functions: returns
1 if argument is true and 0
otherwise



Optimization viewpoint of linear classifiers

Perceptron minimizes a convex relaxation version of 0/1 loss:

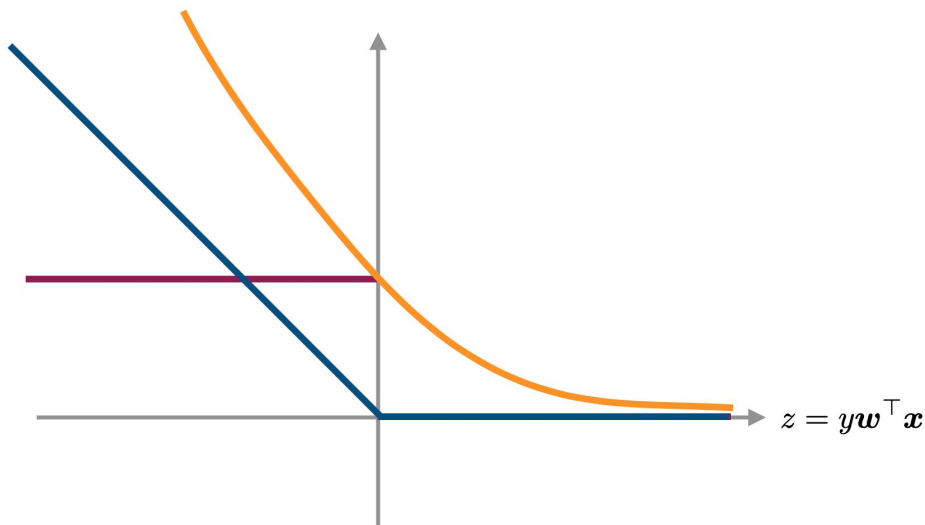
$$\arg \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \sum_{i=1}^n \max(0, -y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$



Optimization viewpoint of linear classifiers

Logistic regression minimizes another convex relaxation version of 0/1 loss, and it is a smooth version of loss used in Perceptron:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \sum_{i=1}^n \log \left(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)} \right)$$



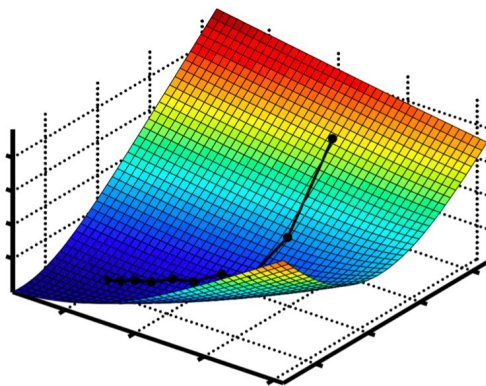
GD and SGD

Our task now is to solve the following optimization problem to find optimal model:

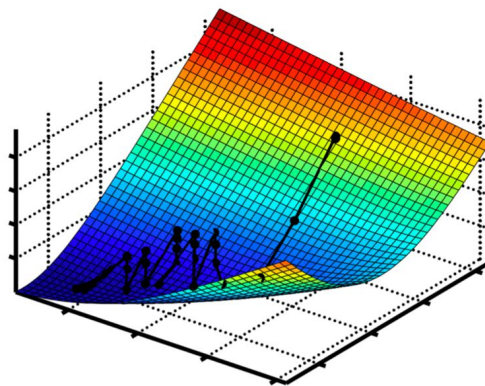
$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right)$$

The function is a convex function of the parameters and we can use GD or SGD to find the minimizer.

GD



SGD



Overfitting in Logistic Regression

MLE can exhibit severe **overfitting** for data sets that are **linearly separable**.

- In fact, there are infinite number separating hyperplanes
- MLE will give solution for weights go to infinity!

$$|\mathbf{w}| \rightarrow \infty, \quad \text{then we have} \quad P(y_i|\mathbf{w}, \mathbf{x}_i) = \sigma(y_i \mathbf{w}^\top \mathbf{x}_i) \rightarrow 1$$

Regularization

MLE can exhibit severe **overfitting** for data sets that are **linearly separable**.

- In fact, there are infinite number separating hyperplanes
- MLE will give solution for weights go to infinity!

$$|\mathbf{w}| \rightarrow \infty, \quad \text{then we have} \quad P(y_i|\mathbf{w}, \mathbf{x}_i) = \sigma(y_i \mathbf{w}^\top \mathbf{x}_i) \rightarrow 1$$

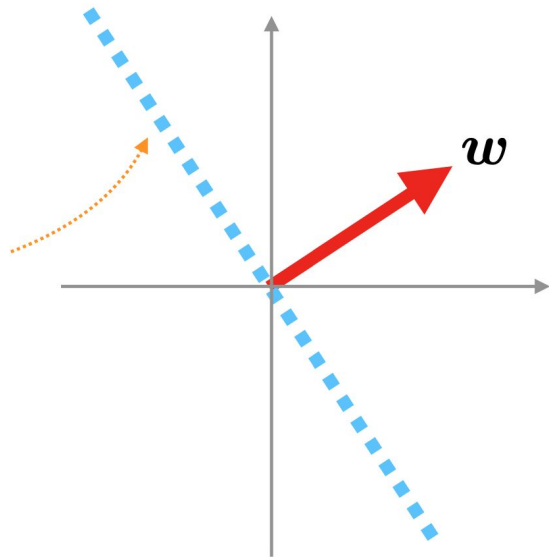
Penalize the model parameters to avoid overfitting:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i} \right) + \lambda \|\mathbf{w}\|_2^2$$

The decision boundary

Recall that for Perceptron the decision boundary of a model parametrized by \mathbf{w} is the set of all points for which:

$$\mathbf{w}^\top \mathbf{x} = 0$$



The decision boundary

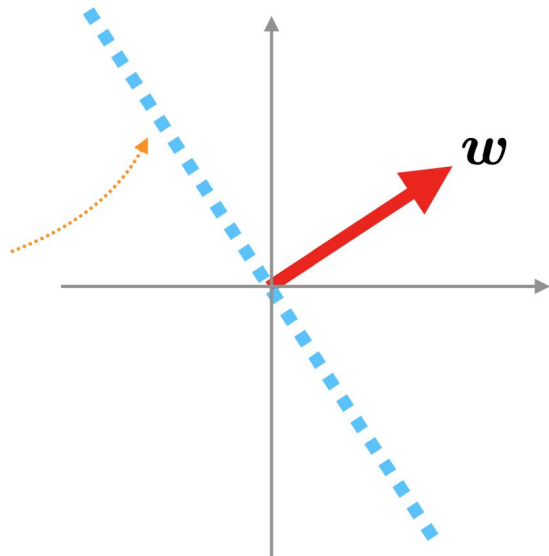
Recall that for Perceptron the decision boundary of a model parametrized by \mathbf{w} is the set of all points for which:

$$\mathbf{w}^\top \mathbf{x} = 0$$

For logistic regression, the decision boundary is the set of all points for which

$$\mathbb{P}[y = -1] = \mathbb{P}[y = +1] = \frac{1}{2}$$

How does the decision boundary for logistic regression look like?



The decision boundary

How does the decision boundary for logistic regression look like?

$$\mathbb{P}[y = +1] = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} = \frac{1}{2}$$



$$e^{-\mathbf{w}^\top \mathbf{x}} = 1$$

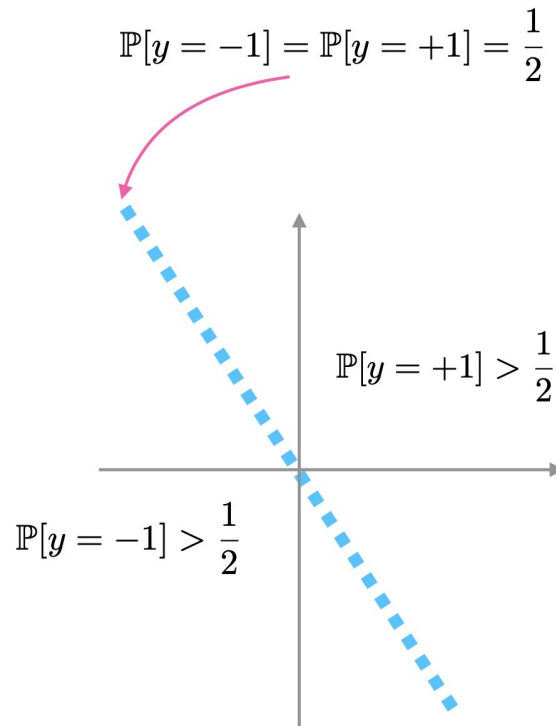


$$-\mathbf{w}^\top \mathbf{x} = 0$$



$$\mathbf{w}^\top \mathbf{x} = 0$$

The decision boundary of logistic regression is linear!



Note on Label Notations

In this class note, we use

$$\mathcal{Y} = \{-1, +1\}$$

It is also possible (and actually, more popular) to use

$$\mathcal{Y} = \{0, 1\}$$

You will work with this setting in your HW3.

Perceptron versus logistic regression

	Perceptron	Logistic regression
Model	$\mathbf{w}^\top \mathbf{x} + b$	$\mathbf{w}^\top \mathbf{x} + b$
Prediction	$\text{sign}(\mathbf{w}^\top \mathbf{x} + b)$	$\sigma(\mathbf{w}^\top \mathbf{x} + b)$
Training loss function	$\max(0, -y(\mathbf{w}^\top \mathbf{x} + b))$	$\log(1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)})$
Decision boundary	Linear	Linear
Regularization	NA (data needs to be separable)	$\ \mathbf{w}\ _2^2$ or $\ \mathbf{w}\ _1$
Optimization	GD + SGD	GD + SGD

Logistic regression in scikit-learn

```
from sklearn.linear_model import LogisticRegression
from sklearn.cross_validation import train_test_split

# Assume data are in X (features) and y (binary labels)
# Scikit Logistic Regression
scikit_log_reg = LogisticRegression()
scikit_log_reg.fit(X, y)
```