

# CSCI-567: Machine Learning (Fall 2019)

Prof. Victor Adamchik

U of Southern California

Sep. 17, 2019

September 17, 2019 1 / 72

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression
- 5 Multiclass Classification

September 17, 2019 3 / 72

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression
- 5 Multiclass Classification

September 17, 2019 2 / 72

## Regression

**Predicting a continuous outcome variable using past observations**

### Key difference from classification

- continuous vs discrete
- measure *prediction errors* differently.
- lead to quite different learning algorithms.

**Linear Regression:** regression with *linear models*:  $f(\mathbf{w}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$

September 17, 2019 4 / 72

## Linear Least Squares Regression

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N (\mathbf{x}_n^T \mathbf{w} - y_n)^2 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

Three approaches to find the minimum:

- Closed Form (setting gradient to zero)  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Gradient Descent (GD)
- Stochastic Gradient Descent (SGD)

September 17, 2019 5 / 72

## Gradient Descent (GD)

**Goal:** minimize  $f(w)$

Consider the definition

$$f'(w) = \lim_{\Delta x \rightarrow 0} \frac{f(w + \Delta x) - f(w)}{\Delta x}$$

Our gradient is an estimation of the derivative

$$\nabla f(w) = \frac{f(w + \Delta x) - f(w)}{\Delta x}$$

Then we need to move in its *opposite* direction to climb down the function.

$$f(w + \Delta x) = f(w) - \Delta x \nabla f(w)$$

September 17, 2019 7 / 72

## Gradient

The gradient vector  $\nabla f$  points in the direction of greatest rate of increase of  $f$  at a given point.

The the rates of change of  $f$  in all directions is given by

$$\nabla f \cdot u = \|\nabla f\| \|u\| \cos \alpha$$

Hence, the direction of *greatest decrease* of  $f$  is the direction opposite to the gradient vector, when  $\alpha = \pi$

We will minimize  $RSS(w)$  using a gradient descent method.

September 17, 2019 6 / 72

## Algorithm: Gradient Descent

**Goal:** minimize  $F(\mathbf{w})$

**Algorithm:** move a bit in the *negative gradient direction*

**initialize**  $\mathbf{w}^{(0)}$   
**while not converged do**

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \lambda \nabla F(\mathbf{w}^{(t)})$$

where  $\lambda > 0$  is called *step size* or *learning rate*

- in theory  $\lambda$  should be set in terms of some parameters of  $F$
- in practice we just try several small values
- there are many possible ways to detect convergence.

September 17, 2019 8 / 72

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ .

Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \quad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

GD:

- Initialize  $w_1^{(0)}$  and  $w_2^{(0)}$  (to be 0 or *randomly*),  $t = 0$
- do

$$w_1^{(t+1)} \leftarrow w_1^{(t)} - \lambda \left[ 2(w_1^{(t)^2} - w_2^{(t)})w_1^{(t)} + w_1^{(t)} - 1 \right]$$

$$w_2^{(t+1)} \leftarrow w_2^{(t)} - \lambda \left[ -(w_1^{(t)^2} - w_2^{(t)}) \right]$$

$$t \leftarrow t + 1$$

- until  $F(\mathbf{w}^{(t)})$  **does not change much**

September 17, 2019 9 / 72

## Why GD?

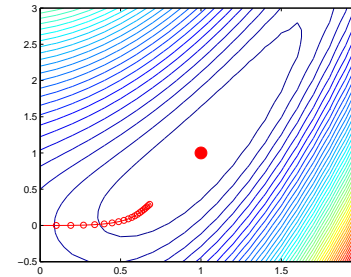
Using the first-order approximation

$$f(w + \Delta x) = f(w) + \Delta x \nabla f(w)$$

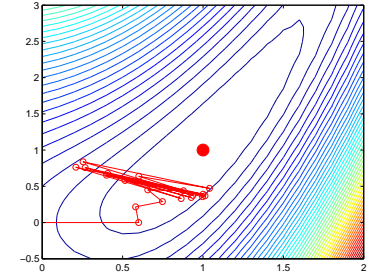
we move a bit in the negative gradient direction  $\Delta x = -\lambda \nabla f(w)$

This ensures

$$f(w - \lambda \nabla f(w)) = f(w) - \lambda (\nabla f(w))^2 \leq f(w)$$



reasonable  $\lambda$  decreases function value



but large  $\lambda$  is unstable

September 17, 2019 10 / 72

## Applying GD to Linear Regression

In the previous discussion, we have computed:

$$\frac{1}{2} \nabla RSS(\mathbf{w}) = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = \sum_n \mathbf{x}_n (\mathbf{x}_n^T \mathbf{w} - y_n) = \sum_n \mathbf{x}_n (f(\mathbf{x}_n) - y_n)$$

**GD update:**

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \lambda \mathbf{X}^T (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$

For a single weight,

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \lambda \sum_n x_{nj} (f(\mathbf{x}_n) - y_n)$$

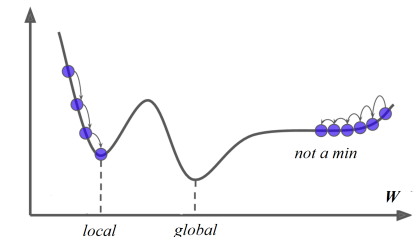
The algorithm uses all training points on each iteration. The algorithm is called *batch gradient descent*.

September 17, 2019 11 / 72

## GD challenges

**There two main challenges with GD:**

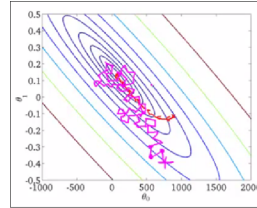
- it may converge to a local minimum.
- it may not find a minimum at all. "vanishing gradient".



September 17, 2019 12 / 72

## Stochastic Gradient Descent (SGD)

- GD: move a bit in the negative gradient direction.
- SGD: move a bit in a *noisy* negative gradient direction.
- In SGD, we use one training sample at each iteration.
- Need to randomly shuffle the training examples before calculating it.
- SGD is widely used for larger dataset and can be trained in parallel.



## SGD for Linear Regression

### Algorithm:

initialize  $w^{(0)}$   
for each training sample  $n$ :  
  for each weight  $j$ :

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \lambda x_{nj} (f(\mathbf{x}_n) - y_n)$$

The term “stochastic” comes from the fact that the gradient based on a single training sample.

SGD makes progress with each training example as it looks at.

Key point: it could be *much faster to obtain a stochastic gradient!*

## GD versus SGD

In GD we calculate the gradient for all training points

In SGD we calculate the gradient for one sample (or a small batch, called a *mini-batch*) of training data

In SGD you might not be taking the most optimal route to get to the solution.

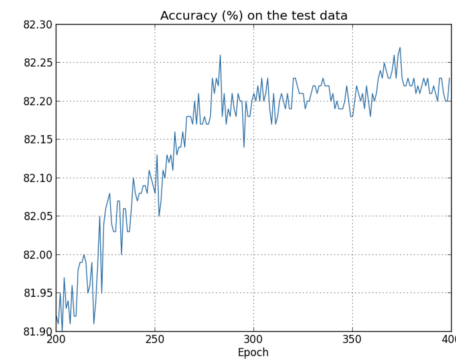
SGD may work for non-convex functions

In SGD you need to go through the training set several times (this is called an *epoch*).

You must specify the batch size (a typical size is 256) and number of epochs (a hyperparameter) for a learning algorithm.

## Epoch and overfitting

This shows how test accuracy is changing due to the number of epochs:



## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression
- 5 Multiclass Classification

September 17, 2019 17 / 72

## Deriving classification algorithms

Let's follow the steps:

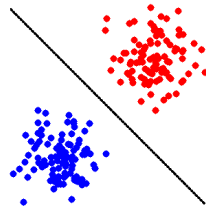
**Step 1.** Pick a set of models  $\mathcal{F}$ .

Again try linear models, but how to predict a label using  $w^T x$ ?

**Sign** of  $w^T x$  predicts the label:

$$\text{sign}(w^T x) = \begin{cases} +1 & \text{if } w^T x > 0 \\ -1 & \text{if } w^T x \leq 0 \end{cases}$$

(Sometimes use  $\text{sgn}$  for  $\text{sign}$  too.)



September 17, 2019 19 / 72

## General idea to provide ML algorithms

1. Pick a set of **models**  $\mathcal{F}$ 
  - e.g.  $\mathcal{F} = \{f(x) = w^T x \mid w \in \mathbb{R}^D\}$
2. Define **error/loss**  $L(y', y)$
3. Find **empirical risk minimizer (ERM)**:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{n=1}^N L(f(x_n), y_n)$$

September 17, 2019 18 / 72

## The models

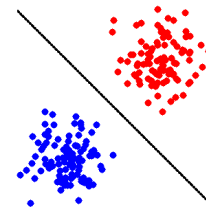
The set of **(separating) hyperplanes**:

$$\mathcal{F} = \{f(x) = \text{sgn}(w^T x) \mid w \in \mathbb{R}^D\}$$

Good choice for **linearly separable** data, i.e.,  $\exists w$  s.t.

$$\text{sgn}(w^T x_n) = y_n \quad \text{or} \quad y_n w^T x_n > 0$$

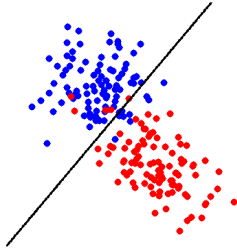
for all  $n \in [N]$ .



September 17, 2019 20 / 72

## The models

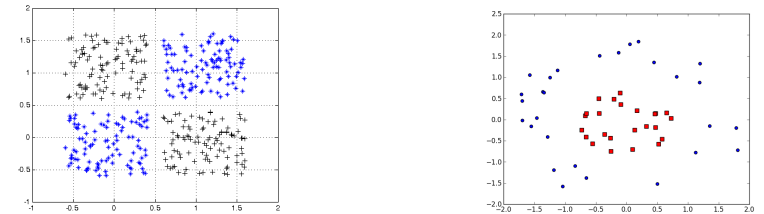
Still makes sense for “almost” linearly separable data



September 17, 2019 21 / 72

## The models

For clearly not linearly separable data,



Again can apply a **nonlinear mapping**  $\Phi$ :

$$\mathcal{F} = \{f(x) = \text{sgn}(w^T \Phi(x)) \mid w \in \mathbb{R}^M\}$$

More discussions in the next lectures.

September 17, 2019 22 / 72

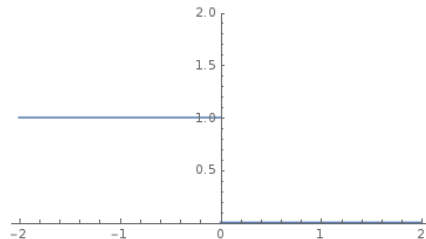
## 0-1 Loss

**Step 2.** Define error/loss  $L(y', y)$ .

Most natural one for classification: **0-1 loss**  $L(y', y) = \mathbb{I}[y' \neq y]$

For classification, more convenient to look at the loss **as a function of**  $yw^T x$ . That is, with

$$\ell_{0-1}(z) = \mathbb{I}[z \leq 0]$$

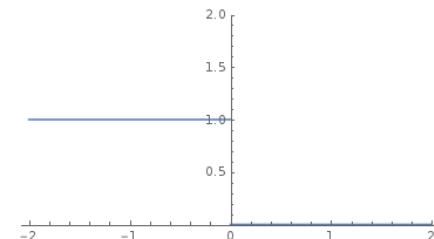


the loss for hyperplane  $w$  on example  $(x, y)$  is  $\ell_{0-1}(yw^T x)$

September 17, 2019 23 / 72

## Minimizing 0-1 loss is hard

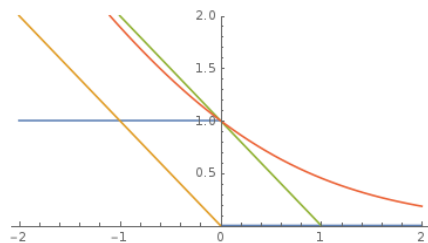
However, 0-1 loss is **not convex**.



Even worse, minimizing 0-1 loss is **NP-hard in general**. Recall Integer Linear programming.

September 17, 2019 24 / 72

## Surrogate Losses



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression)

September 17, 2019 25 / 72

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 **Perceptron**
- 4 Logistic Regression
- 5 Multiclass Classification

September 17, 2019 27 / 72

## ML becomes convex optimization

**Step 3.** Find ERM:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \sum_{n=1}^N \ell(y_n \mathbf{w}^T \mathbf{x}_n)$$

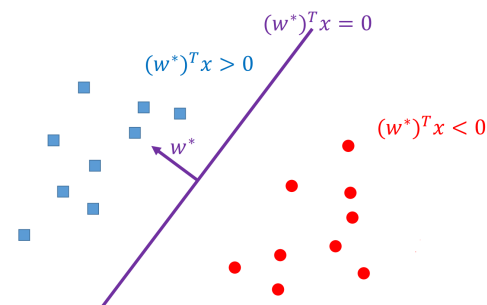
where  $\ell(\cdot)$  can be perceptron/hinge/logistic loss

- **no closed-form** in general (unlike linear regression)
- can apply general convex optimization methods

September 17, 2019 26 / 72

## The Perceptron

The Perceptron (introduced by Rosenblatt in 1957) is a linear model for classification. Its model is a hyperplane that partitions space into two regions. Perceptron is a rough model for how individual neurons in the brain work.



September 17, 2019 28 / 72

## The Perceptron Algorithm

Mathematically: **Stochastic Gradient Descent** applied to perceptron loss

i.e. find the minimizer of

$$\begin{aligned} F(\mathbf{w}) &= \sum_{n=1}^N \ell_{\text{perceptron}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\} \end{aligned}$$

using SGD

## Applying GD to perceptron loss

### Objective

$$F(\mathbf{w}) = \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\}$$

Gradient is

$$\nabla F(\mathbf{w}) = \sum_{n=1}^N -\mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

(only misclassified examples contribute to the gradient)

### GD update

$$\mathbf{w} \leftarrow \mathbf{w} + \lambda \sum_{n=1}^N \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

## Applying SGD to perceptron loss

How to construct a stochastic gradient?

### SGD update

$$\mathbf{w} \leftarrow \mathbf{w} + \lambda \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

## The Perceptron Algorithm

Perceptron algorithm is *SGD with  $\lambda = 1$  applied to perceptron loss*:

Repeat:

- Pick a data point  $\mathbf{x}_n$  uniformly at random
- If  $\text{sgn}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

Note:

- The algorithm is online and error driven.
- If the prediction is correct, it does nothing.
- $\mathbf{w}$  is always a **linear combination** of the training examples.
- It uses epochs as a hyperparameter.



## Why does it make sense?

If the current weight  $\mathbf{w}$  makes a mistake

$$y_n \mathbf{w}^T \mathbf{x}_n < 0$$

then after the update  $\mathbf{w}' = \mathbf{w} + y_n \mathbf{x}_n$  we have

$$y_n \mathbf{w}'^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + \|\mathbf{x}_n\|^2 \geq y_n \mathbf{w}^T \mathbf{x}_n$$

Thus it is more likely to get it right after the update.

## Perceptron Convergence Theorem

Suppose the perceptron algorithm is run on a linearly separable data set  $\mathcal{D}$  with margin  $\gamma \geq 0$ . Assume that  $\|\mathbf{x}\| = 1$ . Then the algorithm will converge after at most  $\frac{1}{\gamma^2}$  updates.

## Any theory?

- If training set is linearly separable, Perceptron *converges in a finite number of steps*
- How long does it take to converge?
- By "how long", what we really mean is "how many updates".
- One way to make this definition is through the notion of margin.
- The margin is the distance between the hyperplane and the nearest point.

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression
  - A Probabilistic View
  - Optimization
- 5 Multiclass Classification

## A simple view

**In one sentence:** find the minimizer of

$$\begin{aligned} F(\mathbf{w}) &= \sum_{n=1}^N \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \sum_{n=1}^N \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) \end{aligned}$$

*But why logistic loss? and why "regression"?*

## Predicting probability

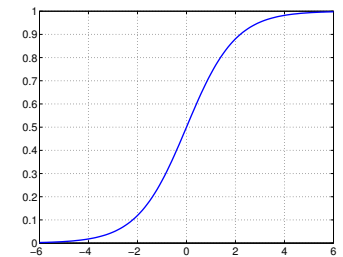
Instead of predicting a discrete label, can we *predict the probability of each label?* i.e. regress the probabilities

One way: **sigmoid function + linear model**

$$\mathbb{P}(y = +1 \mid \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

where  $\sigma$  is the sigmoid function:

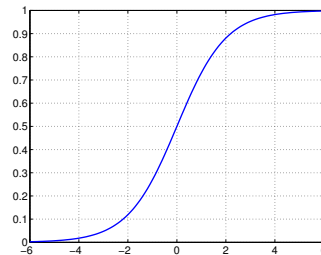
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



## Properties

**Properties** of sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$

- between 0 and 1 (good as probability)
- $\sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \Leftrightarrow \mathbf{w}^T \mathbf{x} \geq 0$ , consistent with predicting the label with  $\text{sgn}(\mathbf{w}^T \mathbf{x})$
- larger  $\mathbf{w}^T \mathbf{x} \Rightarrow$  larger  $\sigma(\mathbf{w}^T \mathbf{x}) \Rightarrow$  higher *confidence* in label 1
- $\sigma(z) + \sigma(-z) = 1$  for all  $z$



The probability of label  $-1$  is naturally

$$1 - \mathbb{P}(y = +1 \mid \mathbf{x}; \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = \sigma(-\mathbf{w}^T \mathbf{x})$$

and thus

$$\mathbb{P}(y \mid \mathbf{x}; \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-y \mathbf{w}^T \mathbf{x}}}$$

## How to regress with discrete labels?

*What we observe are labels, not probabilities.*

Take a **probabilistic view**

- assume data is generated in this way by some  $\mathbf{w}$
- perform Maximum Likelihood Estimation (MLE)

Specifically, what is the probability of seeing label  $y_1, \dots, y_n$  given  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , as a function of some  $\mathbf{w}$ ?

$$P(\mathbf{w}) = \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{w})$$

**MLE:** find  $\mathbf{w}^*$  that **maximizes the probability**  $P(\mathbf{w})$

## The MLE solution

$$\begin{aligned}\mathbf{w}^* &= \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{n=1}^N \ln \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N -\ln \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w}) \\ &= \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) = \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \operatorname{argmin}_{\mathbf{w}} F(\mathbf{w})\end{aligned}$$

i.e. *minimizing logistic loss is exactly doing MLE for the sigmoid model!*

September 17, 2019 41 / 72

## Newton method

Newton's method is an extension of steepest descent, where the second-order term in the Taylor series is used.

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2$$

Let us minimize the right hand side:

$$f'(x_0) + f''(x_0)(x - x_0) = 0 \quad \text{or} \quad x = x_0 - \frac{f'(x_0)}{f''(x_0)}$$

We will iterate this procedure

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

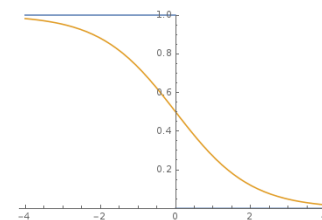
September 17, 2019 43 / 72

## Let's apply SGD again

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} - \lambda \nabla F(\mathbf{w}) \\ &= \mathbf{w} - \lambda \nabla_{\mathbf{w}} \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \mathbf{w} - \lambda \left( \frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n \mathbf{x}_n \\ &= \mathbf{w} - \lambda \left( \frac{-e^{-z}}{1 + e^{-z}} \Big|_{z=y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n \mathbf{x}_n \\ &= \mathbf{w} + \lambda \sigma(-y_n \mathbf{w}^T \mathbf{x}_n) y_n \mathbf{x}_n \\ &= \mathbf{w} + \lambda \mathbb{P}(-y_n | \mathbf{x}_n; \mathbf{w}) y_n \mathbf{x}_n\end{aligned}$$

This is a *soft version of Perceptron!*

$\mathbb{P}(-y_n | \mathbf{x}_n; \mathbf{w})$  versus  $\mathbb{I}[y_n \neq \operatorname{sgn}(\mathbf{w}^T \mathbf{x}_n)]$



September 17, 2019 42 / 72

## Deriving Newton method

This could be generalized for functions  $f$  of several variables:

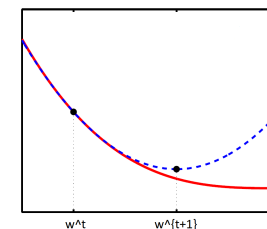
$$x_{n+1} = x_n - \mathbf{H}^{-1}(x_n) \nabla f(x_n)$$

where  $\mathbf{H}$  is the Hessian

$$H_{ij} = \frac{\partial^2 F(\mathbf{x})}{\partial x_i \partial x_j}$$

Therefore, for convex  $F$  (so  $H_t$  is *positive semidefinite*) we obtain **Newton method**:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)})$$



September 17, 2019 44 / 72

## Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \lambda \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\lambda$  (*so no tuning needed!*)
- converges *super fast* in terms of #iterations needed
- requires **second-order** information

September 17, 2019 45 / 72

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression
- 5 **Multiclass Classification**
  - Multinomial logistic regression
  - Reduction to binary classification

September 17, 2019 47 / 72

## Applying Newton to logistic loss

$$\nabla_{\mathbf{w}} \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) = -\sigma(-y_n \mathbf{w}^T \mathbf{x}_n) y_n \mathbf{x}_n$$

$$\begin{aligned} \nabla_{\mathbf{w}}^2 \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) &= \left( \frac{\partial \sigma(z)}{\partial z} \Big|_{z=-y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n^2 \mathbf{x}_n \mathbf{x}_n^T \\ &= \left( \frac{e^{-z}}{(1 + e^{-z})^2} \Big|_{z=-y_n \mathbf{w}^T \mathbf{x}_n} \right) \mathbf{x}_n \mathbf{x}_n^T \\ &= \sigma(y_n \mathbf{w}^T \mathbf{x}_n) (1 - \sigma(y_n \mathbf{w}^T \mathbf{x}_n)) \mathbf{x}_n \mathbf{x}_n^T \end{aligned}$$

### Exercises:

- why is the Hessian of logistic loss positive semidefinite?
- can we apply Newton method to perceptron/hinge loss?

September 17, 2019 46 / 72

## Classification

Recall the setup:

- input (feature vector):  $\mathbf{x} \in \mathbb{R}^D$
- output (label):  $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping  $f: \mathbb{R}^D \rightarrow [C]$

### Examples:

- recognizing digits ( $C = 10$ ) or letters ( $C = 26$  or  $52$ )
- predicting weather: sunny, cloudy, rainy, etc
- predicting image category: ImageNet dataset ( $C \approx 20K$ )

**Nearest Neighbor Classifier** naturally works for arbitrary  $C$ .

September 17, 2019 48 / 72

## Linear models: from binary to multiclass

*What should a linear model look like for multiclass tasks?*

Note: a linear model for binary tasks (switching from  $\{-1, +1\}$  to  $\{1, 2\}$ )

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ 2 & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

By setting  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$ , it can be written as

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}_1^T \mathbf{x} \geq \mathbf{w}_2^T \mathbf{x} \\ 2 & \text{if } \mathbf{w}_2^T \mathbf{x} > \mathbf{w}_1^T \mathbf{x} \end{cases}$$

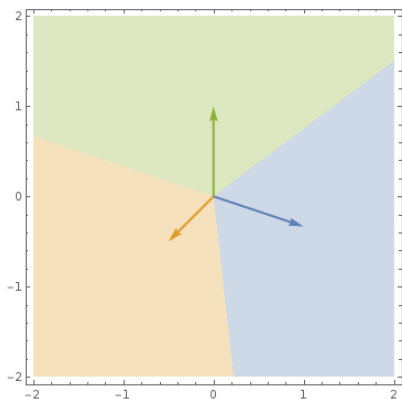
$$= \operatorname{argmax}_{k \in \{1, 2\}} \mathbf{w}_k^T \mathbf{x}$$

for any  $\mathbf{w}_1, \mathbf{w}_2$ .

Think of  $\mathbf{w}_k^T \mathbf{x}$  as a **score for class  $k$** .

September 17, 2019 49 / 72

## Linear models: from binary to multiclass



$$\mathbf{w}_1 = (1, -\frac{1}{3})$$

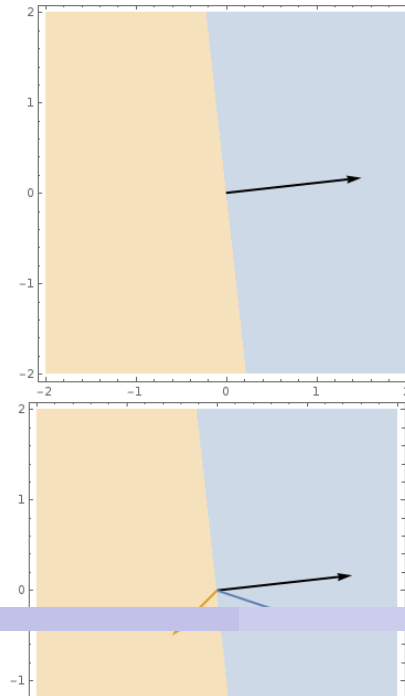
$$\mathbf{w}_2 = (-\frac{1}{2}, -\frac{1}{2})$$

$$\mathbf{w}_3 = (0, 1)$$

- **Blue class:**  
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- **Orange class:**  
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- **Green class:**  
 $\{\mathbf{x} : 3 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$

September 17, 2019 51 / 72

## Linear models: from binary to multiclass



$$\mathbf{w} = (\frac{3}{2}, \frac{1}{6}) = \mathbf{w}_1 - \mathbf{w}_2$$

$$\mathbf{w}_1 = (1, -\frac{1}{3})$$

$$\mathbf{w}_2 = (-\frac{1}{2}, -\frac{1}{2})$$

- **Blue class:**  
 $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} > 0\}$   
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- **Orange class:**

September 17, 2019 50 / 72

## Linear models for multiclass classification

$$\mathcal{F} = \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^D \right\}$$

$$= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} (\mathbf{W} \mathbf{x})_k \mid \mathbf{W} \in \mathbb{R}^{C \times D} \right\}$$

*How do we generalize perceptron/hinge/logistic loss?*

This lecture: focus on the more popular **logistic loss**

September 17, 2019 52 / 72

## Multinomial logistic regression: a probabilistic view

Observe: for binary logistic regression, with  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$ :

$$\mathbb{P}(y = 1 \mid \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}_1^T \mathbf{x}}}{e^{\mathbf{w}_1^T \mathbf{x}} + e^{\mathbf{w}_2^T \mathbf{x}}}$$

Naturally, for class  $y = y_n$

$$\mathbb{P}(y = y_n \mid \mathbf{x}; \mathbf{W}) = \frac{e^{\mathbf{w}_{y_n}^T \mathbf{x}}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}}}$$

This is called the *softmax function*.

## Applying MLE again

Maximize probability of seeing labels  $y_1, \dots, y_N$  given  $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$P(\mathbf{W}) = \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{W}) = \prod_{n=1}^N \frac{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\mathbf{W}) = \sum_{n=1}^N \ln \left( \frac{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}} \right) = \sum_{n=1}^N \ln \left( 1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right)$$

This is the *multiclass logistic loss*, a.k.a *cross-entropy loss*.

When  $C = 2$ , this is the same as binary logistic loss.

## Optimization

Apply **SGD**: what is the gradient of

$$g(\mathbf{W}) = \ln \left( 1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right) ?$$

This is a  $C \times D$  matrix.

Take the derivative wrt  $\mathbf{w}_j \neq \mathbf{w}_{y_n}$ :

$$\begin{aligned} \nabla_{\mathbf{w}_j} g(\mathbf{W}) &= \frac{e^{(\mathbf{w}_j - \mathbf{w}_{y_n})^T \mathbf{x}_n}}{1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n}} \mathbf{x}_n^T \\ &= \frac{e^{\mathbf{w}_j^T \mathbf{x}_n}}{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n} + \sum_{k \neq y_n} e^{\mathbf{w}_k^T \mathbf{x}_n}} \mathbf{x}_n^T \\ &= \frac{e^{\mathbf{w}_j^T \mathbf{x}_n}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}} \mathbf{x}_n^T = \mathbb{P}(j \mid \mathbf{x}_n; \mathbf{W}) \mathbf{x}_n^T \end{aligned}$$

## Optimization

Apply **SGD**

$$g(\mathbf{W}) = \ln \left( 1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right)$$

Take the derivative wrt  $\mathbf{w}_{y_n}$ .

$$\begin{aligned} \nabla_{\mathbf{w}_{y_n}} g(\mathbf{W}) &= - \frac{\sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n}}{1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n}} \mathbf{x}_n^T \\ &= -1 + \frac{1}{1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n}} \mathbf{x}_n^T \\ &= -1 + \frac{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}} \mathbf{x}_n^T = (-1 + \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{W})) \mathbf{x}_n^T \end{aligned}$$

## SGD for multinomial logistic regression

Initialize  $\mathbf{W} = \mathbf{0}$  (or randomly). Repeat:

- 1 pick  $n \in [N]$  uniformly at random
- 2 update the parameters

$$\mathbf{W} \leftarrow \mathbf{W} - \lambda \begin{pmatrix} \mathbb{P}(y = 1 \mid \mathbf{x}_n; \mathbf{W}) \\ \vdots \\ \mathbb{P}(y = y_n \mid \mathbf{x}_n; \mathbf{W}) - 1 \\ \vdots \\ \mathbb{P}(y = C \mid \mathbf{x}_n; \mathbf{W}) \end{pmatrix} \mathbf{x}_n^T$$

Think about why the algorithm makes sense.

Consider  $\mathbb{P}(y = y_n) \rightarrow 1$  and  $\mathbb{P}(y = y_n) \rightarrow 0 \dots$

September 17, 2019 57 / 72

## A note on prediction

Having learned  $\mathbf{W}$ , we can either

- make a *deterministic* prediction  $\operatorname{argmax}_{k \in [C]} (\mathbf{W}\mathbf{x})_k$
- make a *randomized* prediction according to  $\mathbb{P}(k \mid \mathbf{x}; \mathbf{W})$

In either case, **mistake is bounded by logistic loss**.

- deterministic

$$1 = \mathbb{I}[f(\mathbf{x}) \neq y] \leq \log_2 \left( 1 + \sum_{k \neq y} e^{(\mathbf{w}_k - \mathbf{w}_y)^T \mathbf{x}} \right)$$

Indeed, there is such  $k$  that  $\mathbf{w}_k^T \mathbf{x} \geq \mathbf{w}_y^T \mathbf{x}$ , therefore the argument of log is larger than 2.

September 17, 2019 58 / 72

## A note on prediction

Having learned  $\mathbf{W}$ , we can either

- make a *deterministic* prediction  $\operatorname{argmax}_{k \in [C]} (\mathbf{W}\mathbf{x})_k$
- make a *randomized* prediction according to  $\mathbb{P}(k \mid \mathbf{x}; \mathbf{W})$

In either case, **(expected) mistake is bounded by logistic loss**.

- randomized

$$\mathbb{E}[\mathbb{I}[f(\mathbf{x}) \neq y]] = 1 - \mathbb{P}(y \mid \mathbf{x}; \mathbf{W}) \leq -\ln \mathbb{P}(y \mid \mathbf{x}; \mathbf{W})$$

Here we used the fact that  $1 - z \leq -\ln z$ ,  
which follows from the Taylor expansion:  $\ln(1 + x) = x + O(x^2)$ .

September 17, 2019 59 / 72

## Reduce multiclass to binary

Is there an *even more general and simpler approach* to derive multiclass classification algorithms?

Given a binary classification algorithm (*any one*, not just linear methods), can we turn it to a multiclass algorithm, *in a black-box manner*?

Yes, there are in fact many ways to do it.

- **one-versus-all** (one-versus-rest, one-against-all, etc)
- **one-versus-one** (all-versus-all, etc)
- **Error-Correcting Output Codes** (ECOC)
- **tree-based reduction**

September 17, 2019 60 / 72

## One-versus-all (OvA)

Idea: make  $C$  binary classifiers.

Training: for each class  $k \in [C]$ ,

- relabel each example with class  $k$  as  $+1$ , and all others as  $-1$
- train a binary classifier  $h_k$  using this new dataset (what size?)

		■	■	■	■
$x_1$ ■	$\Rightarrow$	$x_1$ —	$x_1$ +	$x_1$ —	$x_1$ —
$x_2$ ■		$x_2$ —	$x_2$ —	$x_2$ +	$x_2$ —
$x_3$ ■		$x_3$ —	$x_3$ —	$x_3$ —	$x_3$ +
$x_4$ ■		$x_4$ —	$x_4$ +	$x_4$ —	$x_4$ —
$x_5$ ■		$x_5$ +	$x_5$ —	$x_5$ —	$x_5$ —
		$\Downarrow$	$\Downarrow$	$\Downarrow$	$\Downarrow$
		$h_1$	$h_2$	$h_3$	$h_4$

September 17, 2019 61 / 72

## One-versus-all (OvA)

Prediction: for a new example  $x$

- ask each  $h_k$ : **does this belong to class  $k$ ?** (i.e.  $h_k(x)$ )
- could be several  $h_k$  s.t.  $h_k(x) = +1$ .
- randomly pick one

OvA becomes inefficient as the number of classes rises.

It's possible to create a significantly more efficient OvA model with a deep neural network.

September 17, 2019 62 / 72

## One-versus-one (OvO)

Idea: make  $\binom{C}{2}$  binary classifiers.

Training: for each pair  $(k, k')$ ,

- relabel each example with class  $k$  as  $+1$  and with class  $k'$  as  $-1$
- *discard all other examples*
- train a binary classifier  $h_{(k,k')}$  using this new dataset (what size?)

		■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■
$x_1$ ■	$\Rightarrow$	$x_1$ —	$x_1$ —	$x_1$ —	$x_1$ —	$x_1$ —	$x_1$ —
$x_2$ ■		$x_2$ —	$x_2$ +	$x_2$ —	$x_2$ —	$x_2$ —	$x_2$ +
$x_3$ ■		$x_3$ —	$x_3$ —	$x_3$ +	$x_3$ —	$x_3$ —	$x_3$ —
$x_4$ ■		$x_4$ —	$x_4$ —	$x_4$ —	$x_4$ —	$x_4$ —	$x_4$ —
$x_5$ ■		$x_5$ +	$x_5$ +	$x_5$ —	$x_5$ +	$x_5$ —	$x_5$ —
		$\Downarrow$	$\Downarrow$	$\Downarrow$	$\Downarrow$	$\Downarrow$	$\Downarrow$
		$h_{(1,2)}$	$h_{(1,3)}$	$h_{(3,4)}$	$h_{(4,2)}$	$h_{(1,4)}$	$h_{(3,2)}$

September 17, 2019 63 / 72

## One-versus-one (OvO)

Prediction: for a new example  $x$

- ask each classifier  $h_{(k,k')}$  to **vote for either class  $k$  or  $k'$**
- predict the class with the most votes (break tie in some way)

**More robust** than one-versus-all, but *slower* in prediction.

September 17, 2019 64 / 72



## Error-correcting output codes (ECOC)

Idea: based on a code  $M \in \{-1, +1\}^{C \times L}$ , train  $L$  binary classifiers to learn “is bit  $b$  on or off”.

Training: for each bit  $b \in [L]$

- relabel example  $x_n$  as  $M_{y_n, b}$
- train a binary classifier  $h_b$  on each column of  $M$ .

M	1	2	3	4	5
■	+	-	+	-	+
■	-	-	+	+	+
■	+	+	-	-	-
■	+	+	+	+	-

	1	2	3	4	5
$x_1$ ■	$x_1$ -	$x_1$ -	$x_1$ +	$x_1$ +	$x_1$ +
$x_2$ ■	$x_2$ +	$x_2$ +	$x_2$ -	$x_2$ -	$x_2$ -
$x_3$ ■	$x_3$ +	$x_3$ +	$x_3$ +	$x_3$ +	$x_3$ -
$x_4$ ■	$x_4$ -	$x_4$ -	$x_4$ +	$x_4$ +	$x_4$ +
$x_5$ ■	$x_5$ +	$x_5$ -	$x_5$ +	$x_5$ -	$x_5$ +
	$\Downarrow$ $h_1$	$\Downarrow$ $h_2$	$\Downarrow$ $h_3$	$\Downarrow$ $h_4$	$\Downarrow$ $h_5$

September 17, 2019 65 / 72

## Error-correcting output codes (ECOC)

Prediction: for a new example  $x$

- compute the **predicted code**  $c = (h_1(x), \dots, h_L(x))^T$
- predict the class with the **most similar code**:  $k = \arg\max_k (Mc)_k$

Suppose you have two classes

- 1:  $\{+, -, -, -, -\}$
- 2:  $\{-, +, +, +, +\}$

and the predicting code is  $\{+, +, +, -, -\}$ . Which class does it predict?

Class 1, since it makes only 2 mistakes.

September 17, 2019 66 / 72

## Error-correcting output codes (ECOC)

How to pick the code  $M$ ?

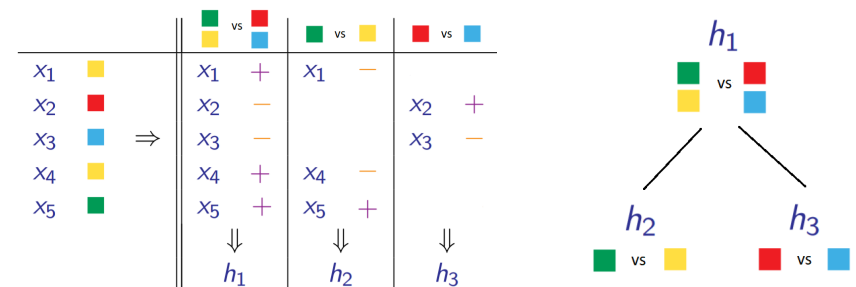
- the more **dissimilar** the codes between different classes are, the better
- **random code** is a good choice, but might create **hard** training sets

One-versus-all ( $L = C$ ) and One-versus-one ( $L = \binom{C}{2}$ ) are two examples of ECOC.

## Tree based method

Idea: train  $\approx C$  binary classifiers to learn “**belongs to which half?**”.

Training: see pictures. In the tree each leaf is a single class.



Prediction is also natural, **but is very fast!** (think ImageNet where  $C \approx 20K$ )

September 17, 2019 67 / 72

September 17, 2019 68 / 72

## Comparisons

In big-O notation,

Reduction	test time	#training points	remark
OvA	C	CN	not robust
OvO	$C^2$	CN	can achieve very small training error
ECOC	L	LN	need diversity when designing code
Tree	$\log_2 C$	$(\log_2 C)N$	good for “extreme classification”

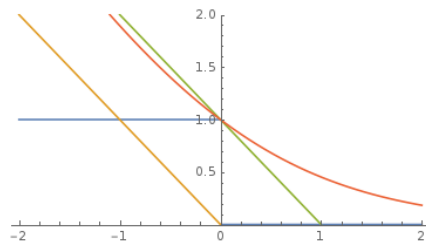
## Summary

Linear models for classification:

Step 1. Model is the set of **separating hyperplanes**

$$\mathcal{F} = \{f(x) = \text{sgn}(w^T x) \mid w \in \mathbb{R}^D\}$$

Step 2. Pick the **surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression)

Step 3. Find empirical risk minimizer (ERM):

$$w^* = \underset{w \in \mathbb{R}^D}{\text{argmin}} \sum_{n=1}^N \ell(y_n w^T x_n)$$

using **GD/SGD/Newton**.