CSCI567 Machine Learning (Fall 2025)

Haipeng Luo

University of Southern California

Sep 12, 2025

Administration

• HW 1 is due on Wed, Sep 17th.

Administration

- HW 1 is due on Wed, Sep 17th.
- recall the late day policy: 3 in total, at most 1 for each homework

Outline

- Review of Last Lecture
- 2 Linear Classifiers and Surrogate Losses
- 3 A Detour of Numerical Optimization Methods
- Perceptron
- Logistic Regression

Outline

- Review of Last Lecture
- 2 Linear Classifiers and Surrogate Losses
- 3 A Detour of Numerical Optimization Methods
- 4 Perceptron
- 5 Logistic Regression

Regression

Predicting a continuous outcome variable using past observations

• temperature, amount of rainfall, house price, etc.

Key difference from classification

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

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

Least square solution

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \operatorname{RSS}(\mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

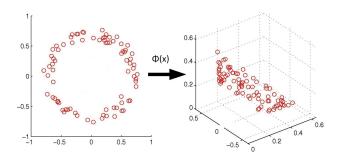
$$= (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1} \mathbf{X}^{\mathrm{T}}\mathbf{y}$$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^{\mathrm{T}} \\ \mathbf{x}_2^{\mathrm{T}} \\ \vdots \\ \mathbf{x}_N^{\mathrm{T}} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

Two approaches to find the minimum:

- find stationary points by setting gradient = 0
- "complete the square"

Regression with nonlinear basis



Model: $f(\boldsymbol{x}) = \boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})$ where $\boldsymbol{w} \in \mathbb{R}^{M}$

Similar least square solution: $oldsymbol{w}^* = \left(oldsymbol{\Phi}^{\mathrm{T}} oldsymbol{\Phi} \right)^{-1} oldsymbol{\Phi}^{\mathrm{T}} oldsymbol{y}$

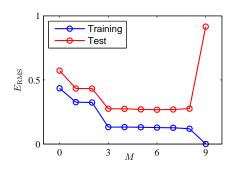
Underfitting and Overfitting

 $M \leq 2$ is *underfitting* the data

- large training error
- large test error

 $M \geq 9$ is *overfitting* the data

- small training error
- large test error



How to prevent overfitting? more data + regularization

$$\boldsymbol{w}^* = \operatorname*{argmin} \left(\mathrm{RSS}(\boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_2^2 \right) = \left(\boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{\Phi} + \lambda \boldsymbol{I} \right)^{-1} \boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{y}$$

General idea to derive ML algorithms

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

- ullet e.g. $\mathcal{F} = \{f(oldsymbol{x}) = oldsymbol{w}^{\mathrm{T}} oldsymbol{x} \mid oldsymbol{w} \in \mathbb{R}^{\mathsf{D}} \}$
- ullet e.g. $\mathcal{F} = \{f(oldsymbol{x}) = oldsymbol{w}^{\mathrm{T}} oldsymbol{\Phi}(oldsymbol{x}) \mid oldsymbol{w} \in \mathbb{R}^{\mathsf{M}} \}$

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

Step 3. Find (regularized) empirical risk minimizer (ERM):

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

ML becomes optimization

General idea to derive ML algorithms

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

- ullet e.g. $\mathcal{F} = \{f(oldsymbol{x}) = oldsymbol{w}^{\mathrm{T}} oldsymbol{x} \mid oldsymbol{w} \in \mathbb{R}^{\mathsf{D}} \}$
- ullet e.g. $\mathcal{F} = \{f(oldsymbol{x}) = oldsymbol{w}^{\mathrm{T}} oldsymbol{\Phi}(oldsymbol{x}) \mid oldsymbol{w} \in \mathbb{R}^{\mathsf{M}} \}$

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

Step 3. Find (regularized) empirical risk minimizer (ERM):

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

ML becomes optimization

Today: another exercise of this recipe + a closer look at Step 3

Outline

- Review of Last Lecture
- 2 Linear Classifiers and Surrogate Losses
- 3 A Detour of Numerical Optimization Methods
- Perceptron
- 6 Logistic Regression

Classification

Recall the setup:

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

Classification

Recall the setup:

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

This lecture: binary classification

- Number of classes: C=2
- Labels: $\{-1, +1\}$ (cat or dog, fraud or not, price up or down...)

Classification

Recall the setup:

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

This lecture: binary classification

- Number of classes: C=2
- Labels: $\{-1,+1\}$ (cat or dog, fraud or not, price up or down...)

We have discussed **nearest neighbor classifier**:

- require carrying the training set
- intuitive but more like a heuristic

Let's follow the recipe:

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

Let's follow the recipe:

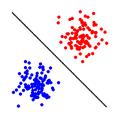
Step 1. Pick a set of models \mathcal{F} .

Again try linear models, but how to predict a label using $m{w}^{\mathrm{T}} m{x}$?

Let's follow the recipe:

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

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



Let's follow the recipe:

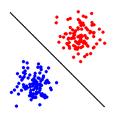
Step 1. Pick a set of models \mathcal{F} .

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

Sign of $w^{\mathrm{T}}x$ predicts the label:

$$\mathsf{sign}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}) = \left\{ \begin{array}{ll} +1 & \text{if } \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} > 0 \\ -1 & \text{if } \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} \leq 0 \end{array} \right.$$

(Sometimes use sgn for sign too.)



The set of (separating) hyperplanes:

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

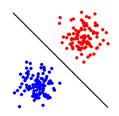
The set of (separating) hyperplanes:

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

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

$$\operatorname{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x_n}) = y_n$$

for all $n \in [N]$.



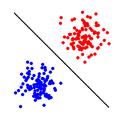
The set of (separating) hyperplanes:

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

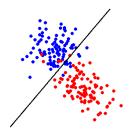
Good choice for *linearly separable* data, i.e., $\exists w$ s.t.

$$\operatorname{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}_{\boldsymbol{n}}) = y_n \quad \text{ or } \quad y_n \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}_{\boldsymbol{n}} > 0$$

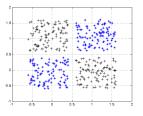
 $\text{ for all } n \in [N].$

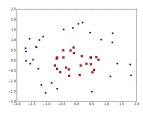


Still makes sense for "almost" linearly separable data

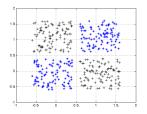


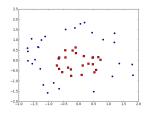
For clearly not linearly separable data,





For clearly not linearly separable data,





Again can apply a **nonlinear mapping** Φ :

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

More discussions in future lectures.

0-1 Loss

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

0-1 Loss

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

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

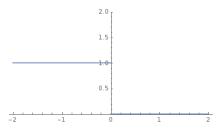
0-1 Loss

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^Tx . That is, with

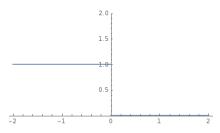
$$\ell_{0\text{-}1}(z) = \mathbb{I}[z \le 0]$$



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

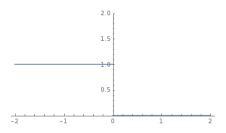
Minimizing 0-1 loss is hard

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



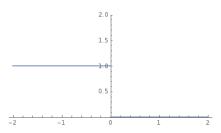
Minimizing 0-1 loss is hard

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

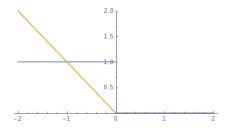


Even worse, minimizing 0-1 loss is NP-hard in general.

Solution: find a convex surrogate loss

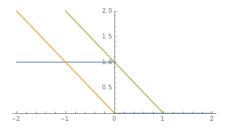


Solution: find a convex surrogate loss



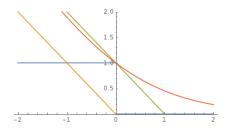
• perceptron loss $\ell_{perceptron}(z) = \max\{0, -z\}$ (used in Perceptron)

Solution: find a convex surrogate loss



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

Solution: find a convex surrogate loss



- perceptron loss $\ell_{perceptron}(z) = \max\{0, -z\}$ (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; the base of \log doesn't matter)

ML becomes convex optimization

Step 3. Find ERM:

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

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

ML becomes convex optimization

Step 3. Find ERM:

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

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

• no closed-form in general (unlike linear regression)

ML becomes convex optimization

Step 3. Find ERM:

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

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

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

ML becomes convex optimization

Step 3. Find ERM:

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

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

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

Note: minimizing perceptron loss does not really make sense

ML becomes convex optimization

Step 3. Find ERM:

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

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

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

Note: minimizing perceptron loss does not really make sense (try w=0), but the algorithm derived from this perspective does.

Outline

- Review of Last Lecture
- 2 Linear Classifiers and Surrogate Losses
- 3 A Detour of Numerical Optimization Methods
 - First-order methods
 - Second-order methods
- 4 Perceptron
- 5 Logistic Regression

Numerical optimization

Problem setup

- Given: a function F(w)
- Goal: minimize F(w) (approximately)

First-order optimization methods

Two simple yet extremely popular methods

- Gradient Descent (GD): simple and fundamental
- Stochastic Gradient Descent (SGD): faster, effective for large-scale problems

First-order optimization methods

Two simple yet extremely popular methods

- Gradient Descent (GD): simple and fundamental
- Stochastic Gradient Descent (SGD): faster, effective for large-scale problems

Gradient is sometimes referred to as *first-order* information of a function. Therefore, these methods are called *first-order methods*.

GD: keep moving in the negative gradient direction

GD: keep moving in the *negative gradient direction*

Start from some (random) $\boldsymbol{w}^{(0)}$. For $t=0,1,2,\ldots$

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

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

GD: keep moving in the negative gradient direction

Start from some (random) $\boldsymbol{w}^{(0)}$. For $t=0,1,2,\ldots$

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

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

ullet in theory η should be set in terms of some parameters of F

GD: keep moving in the negative gradient direction

Start from some (random) $w^{(0)}$. For $t=0,1,2,\ldots$

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

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

- ullet in theory η should be set in terms of some parameters of F
- in practice we often try different small values

GD: keep moving in the negative gradient direction

Start from some (random) $w^{(0)}$. For t = 0, 1, 2, ...

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

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

- ullet in theory η should be set in terms of some parameters of F
- in practice we often try different small values

Stop when $F(\boldsymbol{w}^{(t)})$ does not change much or t reaches a fixed number

Intuition: by first-order **Taylor approximation**

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

Intuition: by first-order **Taylor approximation**

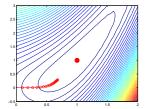
$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

GD ensures
$$F(\boldsymbol{w}^{(t+1)}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w}^{(t+1)} - \boldsymbol{w}^{(t)})$$
$$= F(\boldsymbol{w}^{(t)}) - \eta \|\nabla F(\boldsymbol{w}^{(t)})\|_2^2 \leq F(\boldsymbol{w}^{(t)})$$

Intuition: by first-order **Taylor approximation**

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

GD ensures
$$F(\boldsymbol{w}^{(t+1)}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w}^{(t+1)} - \boldsymbol{w}^{(t)})$$
$$= F(\boldsymbol{w}^{(t)}) - \eta \|\nabla F(\boldsymbol{w}^{(t)})\|_2^2 \leq F(\boldsymbol{w}^{(t)})$$

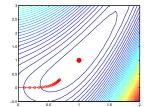


reasonable η decreases function value

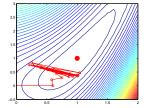
Intuition: by first-order **Taylor approximation**

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

GD ensures
$$F(\boldsymbol{w}^{(t+1)}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w}^{(t+1)} - \boldsymbol{w}^{(t)})$$
$$= F(\boldsymbol{w}^{(t)}) - \eta \|\nabla F(\boldsymbol{w}^{(t)})\|_2^2 \leq F(\boldsymbol{w}^{(t)})$$



reasonable η decreases function value

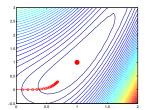


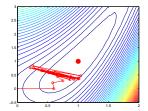
but large η is unstable

Intuition: by first-order Taylor approximation

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

GD ensures
$$F(\boldsymbol{w}^{(t+1)}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w}^{(t+1)} - \boldsymbol{w}^{(t)})$$
$$= F(\boldsymbol{w}^{(t)}) - \eta \|\nabla F(\boldsymbol{w}^{(t)})\|_2^2 \leq F(\boldsymbol{w}^{(t)})$$





reasonable η decreases function value

but large η is unstable

See Colab Example 1

Learning rate η might need to be **changing** over iterations

Learning rate η might need to be **changing** over iterations

 \bullet often <code>decreasing</code>, according to some schedule (e.g., $\eta \approx \frac{1}{t}$ or $\frac{1}{\sqrt{t}})$

Learning rate η might need to be **changing** over iterations

- \bullet often <code>decreasing</code>, according to some schedule (e.g., $\eta \approx \frac{1}{t}$ or $\frac{1}{\sqrt{t}})$
- think F(w) = |w|

Learning rate η might need to be **changing** over iterations

- ullet often decreasing, according to some schedule (e.g., $\eta pprox rac{1}{t}$ or $rac{1}{\sqrt{t}}$)
- $\bullet \ \ \mathsf{think} \ F(w) = |w|$

Adaptive and automatic step size tuning is an active research area

Learning rate η might need to be **changing** over iterations

- ullet often decreasing, according to some schedule (e.g., $\eta pprox rac{1}{t}$ or $rac{1}{\sqrt{t}}$)
- think F(w) = |w|

Adaptive and automatic step size tuning is an active research area

- notable examples: AdaGrad, Adam, etc.
- ullet ideas: tune η based on past gradient information

GD: keep moving in the negative gradient direction

SGD: keep moving in some *noisy* negative gradient direction

GD: keep moving in the negative gradient direction

SGD: keep moving in some noisy negative gradient direction

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \tilde{\nabla} F(\boldsymbol{w}^{(t)})$$

where $\tilde{\nabla}F(\boldsymbol{w}^{(t)})$ is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E}\left[\tilde{\nabla}F(\boldsymbol{w}^{(t)})\right] = \nabla F(\boldsymbol{w}^{(t)}) \qquad \text{(unbiasedness)}$$

GD: keep moving in the negative gradient direction

SGD: keep moving in some noisy negative gradient direction

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \tilde{\nabla} F(\boldsymbol{w}^{(t)})$$

where $\tilde{\nabla} F(\boldsymbol{w}^{(t)})$ is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E}\left[\tilde{\nabla}F(\boldsymbol{w}^{(t)})\right] = \nabla F(\boldsymbol{w}^{(t)}) \qquad \text{(unbiasedness)}$$

See Colab Example 1.

GD: keep moving in the negative gradient direction

SGD: keep moving in some noisy negative gradient direction

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \tilde{\nabla} F(\boldsymbol{w}^{(t)})$$

where $\tilde{\nabla}F(\boldsymbol{w}^{(t)})$ is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E}\left[\tilde{\nabla}F(\boldsymbol{w}^{(t)})\right] = \nabla F(\boldsymbol{w}^{(t)}) \qquad \text{(unbiasedness)}$$

See Colab Example 1.

More examples coming soon. Key point: it could be *much faster to obtain a stochastic gradient!*

Many for both GD and SGD on convex objectives.

Many for both GD and SGD on convex objectives.

They tell you how many iterations t (in terms of ϵ) needed to achieve

$$F(\boldsymbol{w}^{(t)}) - F(\boldsymbol{w}^*) \le \epsilon$$

Many for both GD and SGD on convex objectives.

They tell you how many iterations t (in terms of ϵ) needed to achieve

$$F(\boldsymbol{w}^{(t)}) - F(\boldsymbol{w}^*) \le \epsilon$$

usually SGD needs more iterations

Many for both GD and SGD on convex objectives.

They tell you how many iterations t (in terms of ϵ) needed to achieve

$$F(\boldsymbol{w}^{(t)}) - F(\boldsymbol{w}^*) \le \epsilon$$

- usually SGD needs more iterations
- but again each iteration takes less time

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations t (in terms of ϵ) needed to achieve

$$\|\nabla F(\boldsymbol{w}^{(t)})\| \le \epsilon$$

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations t (in terms of ϵ) needed to achieve

$$\|\nabla F(\boldsymbol{w}^{(t)})\| \le \epsilon$$

ullet that is, how close $oldsymbol{w}^{(t)}$ is as an approximate stationary point

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations t (in terms of ϵ) needed to achieve

$$\|\nabla F(\boldsymbol{w}^{(t)})\| \le \epsilon$$

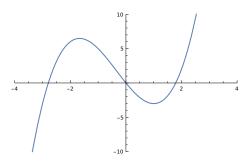
- ullet that is, how close $oldsymbol{w}^{(t)}$ is as an approximate stationary point
- for convex objectives, stationary point ⇒ global minimizer

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations t (in terms of ϵ) needed to achieve

$$\|\nabla F(\boldsymbol{w}^{(t)})\| \le \epsilon$$

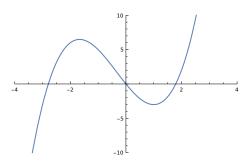
- ullet that is, how close $oldsymbol{w}^{(t)}$ is as an approximate stationary point
- for convex objectives, stationary point ⇒ global minimizer
- for nonconvex objectives, what does it mean?

A stationary point can be a local minimizer



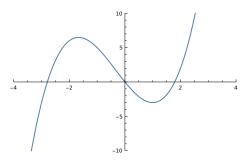
$$f(w) = w^3 + w^2 - 5w$$

A stationary point can be a **local minimizer** or even a **local/global** maximizer



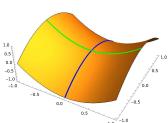
$$f(w) = w^3 + w^2 - 5w$$

A stationary point can be a **local minimizer** or even a **local/global maximizer** (but the latter is not an issue for GD/SGD).



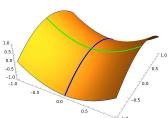
$$f(w) = w^3 + w^2 - 5w$$

•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$



•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

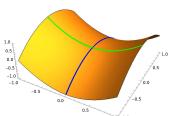


A stationary point can also be *neither a local minimizer nor a local maximizer!*

•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

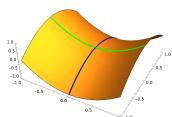
• so $\boldsymbol{w} = (0,0)$ is stationary



•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

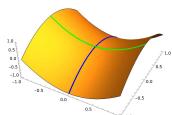
- so $\boldsymbol{w} = (0,0)$ is stationary
- local max for blue direction ($w_1 = 0$)



•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

- so $\boldsymbol{w} = (0,0)$ is stationary
- local max for blue direction ($w_1 = 0$)
- local min for green direction ($w_2 = 0$)

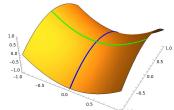


A stationary point can also be *neither a local minimizer nor a local maximizer!* This is called a **saddle point**.

•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

- so $\boldsymbol{w} = (0,0)$ is stationary
- local max for blue direction ($w_1 = 0$)
- local min for green direction ($w_2 = 0$)



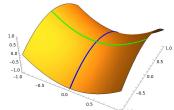


A stationary point can also be *neither a local minimizer nor a local maximizer!* This is called a **saddle point**.

•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

- so $\boldsymbol{w} = (0,0)$ is stationary
- local max for blue direction ($w_1 = 0$)
- local min for green direction ($w_2 = 0$)
- but GD gets stuck at (0,0) only if initialized along the green direction

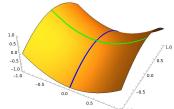




•
$$f(\mathbf{w}) = w_1^2 - w_2^2$$

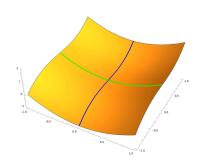
•
$$\nabla f(\mathbf{w}) = (2w_1, -2w_2)$$

- so $\boldsymbol{w} = (0,0)$ is stationary
- local max for blue direction ($w_1 = 0$)
- local min for green direction ($w_2 = 0$)
- but GD gets stuck at (0,0) only if initialized along the green direction
- so not a real issue especially when initialized randomly (Colab Example 2)



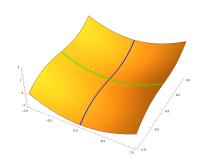


•
$$f(\mathbf{w}) = w_1^2 + w_2^3$$



•
$$f(\mathbf{w}) = w_1^2 + w_2^3$$

•
$$\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$$

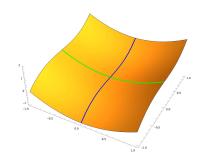


But not all saddle points look like a "saddle" ...

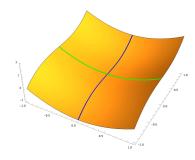
•
$$f(\mathbf{w}) = w_1^2 + w_2^3$$

•
$$\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$$

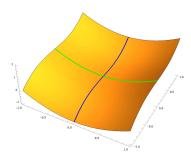
• so $\boldsymbol{w} = (0,0)$ is stationary



- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so $\boldsymbol{w} = (0,0)$ is stationary
- not local min/max for blue direction $(w_1 = 0)$

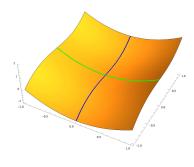


- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so $\boldsymbol{w} = (0,0)$ is stationary
- not local min/max for blue direction $(w_1 = 0)$
- GD gets stuck at (0,0) for any initial point with $w_2 \ge 0$ and small η (Colab Example 3)



But not all saddle points look like a "saddle" ...

- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so $\boldsymbol{w} = (0,0)$ is stationary
- not local min/max for blue direction $(w_1 = 0)$
- GD gets stuck at (0,0) for any initial point with $w_2 \ge 0$ and small η (Colab Example 3)



Even worse, distinguishing local min and saddle point is generally NP-hard.

Summary:

• GD/SGD converges to a stationary point

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or "bad" saddle points (random initialization escapes "good" saddle points)

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or "bad" saddle points (random initialization escapes "good" saddle points)
- recent research shows that many problems have no "bad" saddle points or even "bad" local minimizers

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or "bad" saddle points (random initialization escapes "good" saddle points)
- recent research shows that many problems have no "bad" saddle points or even "bad" local minimizers
- justify the practical effectiveness of GD/SGD (default method to try)

Second-order methods

Recall the intuition of GD: we look at first-order Taylor approximation

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

Second-order methods

Recall the intuition of GD: we look at first-order Taylor approximation

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}} (\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

What if we look at *second-order* Taylor approximation?

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)}) + \frac{1}{2}(\boldsymbol{w} - \boldsymbol{w}^{(t)})^{\mathrm{T}}\boldsymbol{H}_{t}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

Second-order methods

Recall the intuition of GD: we look at first-order Taylor approximation

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}} (\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

What if we look at second-order Taylor approximation?

$$F(\boldsymbol{w}) \approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)}) + \frac{1}{2}(\boldsymbol{w} - \boldsymbol{w}^{(t)})^{\mathrm{T}}\boldsymbol{H}_{t}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

where $\boldsymbol{H}_t = \nabla^2 F(\boldsymbol{w}^{(t)}) \in \mathbb{R}^{\mathsf{D} \times \mathsf{D}}$ is the *Hessian* of F at $\boldsymbol{w}^{(t)}$, i.e.,

$$H_{t,ij} = \frac{\partial^2 F(\boldsymbol{w})}{\partial w_i \partial w_j} \Big|_{\boldsymbol{w} = \boldsymbol{w}^{(t)}}$$

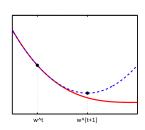
(think "second derivative" when D=1)

Newton method

If we minimize the second-order approximation (via "complete the square")

$$F(\boldsymbol{w})$$

$$egin{aligned} &pprox F(oldsymbol{w}^{(t)}) +
abla F(oldsymbol{w}^{(t)})^{\mathrm{T}}(oldsymbol{w} - oldsymbol{w}^{(t)}) + rac{1}{2}(oldsymbol{w} - oldsymbol{w}^{(t)})^{\mathrm{T}} oldsymbol{H}_t(oldsymbol{w} - oldsymbol{w}^{(t)})^{\mathrm{T}} oldsy$$



Newton method

If we minimize the second-order approximation (via "complete the square")

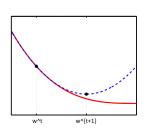
$$F(\boldsymbol{w})$$

$$\approx F(\boldsymbol{w}^{(t)}) + \nabla F(\boldsymbol{w}^{(t)})^{\mathrm{T}}(\boldsymbol{w} - \boldsymbol{w}^{(t)}) + \frac{1}{2}(\boldsymbol{w} - \boldsymbol{w}^{(t)})^{\mathrm{T}}\boldsymbol{H}_{t}(\boldsymbol{w} - \boldsymbol{w}^{(t)})$$

$$= \frac{1}{2}\left(\boldsymbol{w} - \boldsymbol{w}^{(t)} + \boldsymbol{H}_{t}^{-1}\nabla F(\boldsymbol{w}^{(t)})\right)^{\mathrm{T}}\boldsymbol{H}_{t}\left(\boldsymbol{w} - \boldsymbol{w}^{(t)} + \boldsymbol{H}_{t}^{-1}\nabla F(\boldsymbol{w}^{(t)})\right) + \mathrm{cnt}$$

for strictly convex F (so H_t is *positive definite*), we obtain **Newton method**:

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



$$\begin{split} & \boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)}) \\ & \boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \boldsymbol{H}_t^{-1} \nabla F(\boldsymbol{w}^{(t)}) \end{split} \tag{Newton}$$

Both are iterative optimization procedures,

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \boldsymbol{H}_{t}^{-1} \nabla F(\boldsymbol{w}^{(t)})$$
(Newton)

Both are iterative optimization procedures, but Newton method

• has no learning rate η (so no tuning needed!)

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \boldsymbol{H}_{t}^{-1} \nabla F(\boldsymbol{w}^{(t)})$$
(Newton)

- has no learning rate η (so **no tuning needed!**)
- converges super fast in terms of #iterations (for convex objectives)

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

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

- has no learning rate η (so no tuning needed!)
- converges super fast in terms of #iterations (for convex objectives)
 - e.g. how many iterations needed when applied to a quadratic?

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

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

- has no learning rate η (so no tuning needed!)
- converges super fast in terms of #iterations (for convex objectives)
 - e.g. how many iterations needed when applied to a quadratic?
- computing Hessian in each iteration is very slow though

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta \nabla F(\boldsymbol{w}^{(t)})$$

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \boldsymbol{H}_{t}^{-1} \nabla F(\boldsymbol{w}^{(t)})$$
(Newton)

- has no learning rate η (so no tuning needed!)
- converges super fast in terms of #iterations (for convex objectives)
 - e.g. how many iterations needed when applied to a quadratic?
- computing Hessian in each iteration is very slow though
- does not really make sense for nonconvex objectives

Outline

- Review of Last Lecture
- 2 Linear Classifiers and Surrogate Losses
- 3 A Detour of Numerical Optimization Methods
- 4 Perceptron
- 6 Logistic Regression

Recall the perceptron loss

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} \ell_{\mathsf{perceptron}}(y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)$$
$$= \frac{1}{N} \sum_{n=1}^{N} \max\{0, -y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n\}$$

Recall the perceptron loss

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} \ell_{\mathsf{perceptron}}(y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)$$
$$= \frac{1}{N} \sum_{n=1}^{N} \max\{0, -y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n\}$$

Let's approximately minimize it with GD/SGD.

Objective

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

Objective

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

Gradient (or really *sub-gradient*) is

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

(only misclassified examples contribute to the gradient)

Objective

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

Gradient (or really *sub-gradient*) is

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

(only misclassified examples contribute to the gradient)

GD update

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \frac{\eta}{N} \sum_{n=1}^{N} \mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

Objective

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

Gradient (or really *sub-gradient*) is

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

(only misclassified examples contribute to the gradient)

GD update

$$oldsymbol{w} \leftarrow oldsymbol{w} + rac{\eta}{N} \sum_{n=1}^{N} \mathbb{I}[y_n oldsymbol{w}^{\mathrm{T}} oldsymbol{x}_n \leq 0] y_n oldsymbol{x}_n$$

Slow: each update makes one pass of the entire training set!

How to construct a stochastic gradient?

How to construct a stochastic gradient?

One common trick: pick one example $n \in [N]$ uniformly at random, let

$$\tilde{\nabla} F(\boldsymbol{w}^{(t)}) = -\mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

clearly unbiased (convince yourself).

How to construct a stochastic gradient?

One common trick: pick one example $n \in [N]$ uniformly at random, let

$$\tilde{\nabla} F(\boldsymbol{w}^{(t)}) = -\mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

clearly unbiased (convince yourself).

SGD update:

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

How to construct a stochastic gradient?

One common trick: pick one example $n \in [N]$ uniformly at random, let

$$\tilde{\nabla} F(\boldsymbol{w}^{(t)}) = -\mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

clearly unbiased (convince yourself).

SGD update:

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

Fast: each update touches only one data point!

How to construct a stochastic gradient?

One common trick: pick one example $n \in [N]$ uniformly at random, let

$$\tilde{\nabla} F(\boldsymbol{w}^{(t)}) = -\mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

clearly unbiased (convince yourself).

SGD update:

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \mathbb{I}[y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \leq 0] y_n \boldsymbol{x}_n$$

Fast: each update touches only one data point!

Conveniently, objective of most ML tasks is a *finite sum* (over each training point) and the above trick applies!

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

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

Repeat:

- ullet Pick a data point $oldsymbol{x}_n$ uniformly at random
- ullet If $\operatorname{sgn}(oldsymbol{w}^{\mathrm{T}}oldsymbol{x}_n)
 eq y_n$

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

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

Repeat:

- ullet Pick a data point $oldsymbol{x}_n$ uniformly at random
- ullet If $\mathrm{sgn}(oldsymbol{w}^{\mathrm{T}}oldsymbol{x}_n)
 eq y_n$

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

Note:

 $oldsymbol{w}$ is always a *linear combination* of the training examples

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

Repeat:

- ullet Pick a data point $oldsymbol{x}_n$ uniformly at random
- ullet If $\mathrm{sgn}(oldsymbol{w}^{\mathrm{T}}oldsymbol{x}_n)
 eq y_n$

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

Note:

- ullet w is always a *linear combination* of the training examples
- ullet why $\eta=1$? Does not really matter in terms of prediction of $oldsymbol{w}$

Why does it make sense?

If the current weight $oldsymbol{w}$ makes a mistake

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

Why does it make sense?

If the current weight $oldsymbol{w}$ makes a mistake

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

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

$$y_n {w'}^{\mathrm{T}} \boldsymbol{x}_n = y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n + y_n^2 \boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{x}_n \ge y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n$$

Why does it make sense?

If the current weight $oldsymbol{w}$ makes a mistake

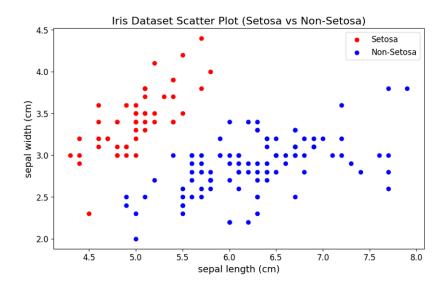
$$y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n < 0$$

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

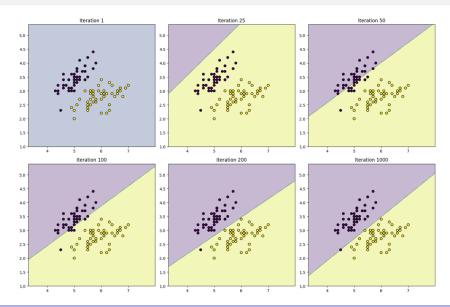
$$y_n \boldsymbol{w'}^{\mathrm{T}} \boldsymbol{x}_n = y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n + y_n^2 \boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{x}_n \ge y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n$$

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

Example: Iris Dataset



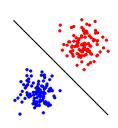
Example: Perceptron for Iris Dataset



Any theory?

If training set is linearly separable

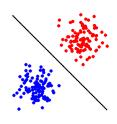
- Perceptron converges in a finite number of steps
- training error is 0



Any theory?

If training set is linearly separable

- Perceptron converges in a finite number of steps
- training error is 0



There are also guarantees when the data are not linearly separable.

Outline

- Review of Last Lecture
- 2 Linear Classifiers and Surrogate Losses
- 3 A Detour of Numerical Optimization Methods
- 4 Perceptron
- 6 Logistic Regression
 - A probabilistic view

A simple view

In one sentence: find the minimizer of

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n)$$
$$= \frac{1}{N} \sum_{n=1}^{N} \ln(1 + e^{-y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n})$$

A simple view

In one sentence: find the minimizer of

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)$$
$$= \frac{1}{N} \sum_{n=1}^{N} \ln(1 + e^{-y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n})$$

Before optimizing it: 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

Predicting probability

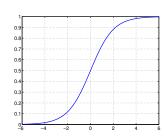
Instead of predicting a discrete label, can we *predict the probability of each label?* i.e. regress the probabilities

One way: sigmoid function + linear model

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

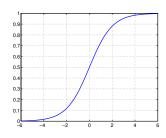
where σ is the sigmoid function:

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



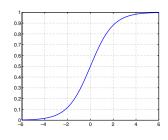
Properties of sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$

• between 0 and 1 (good as probability)



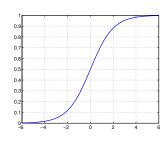
Properties of sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$

- between 0 and 1 (good as probability)
- $\sigma(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}) \geq 0.5 \Leftrightarrow \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} \geq 0$, consistent with predicting the label with $\mathrm{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x})$



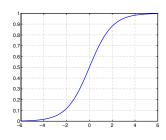
Properties of sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$

- between 0 and 1 (good as probability)
- $\sigma(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}) \geq 0.5 \Leftrightarrow \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} \geq 0$, consistent with predicting the label with $\mathrm{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x})$
- larger $m{w}^{\mathrm{T}}m{x} \Rightarrow \mathsf{larger} \ \sigma(m{w}^{\mathrm{T}}m{x}) \Rightarrow \mathsf{higher}$ **confidence** in label 1



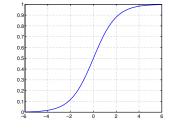
Properties of sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$

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



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

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

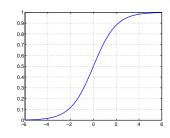


The probability of label -1 is naturally

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

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

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



The probability of label -1 is naturally

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

and thus

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

How to regress with discrete labels?

What we observe are labels, not probabilities.

How to regress with discrete labels?

What we observe are labels, not probabilities.

Take a probabilistic view

- ullet assume data is independently generated in this way by some w
- perform Maximum Likelihood Estimation (MLE)

How to regress with discrete labels?

What we observe are labels, not probabilities.

Take a probabilistic view

- ullet assume data is independently generated in this way by some w
- perform Maximum Likelihood Estimation (MLE)

Specifically, what is the probability of seeing label y_1, \dots, y_n given x_1, \dots, x_n , as a function of some w?

$$P(\boldsymbol{w}) = \prod_{n=1}^{N} \mathbb{P}(y_n \mid \boldsymbol{x_n}; \boldsymbol{w})$$

MLE: find w^* that maximizes the probability P(w)

The MLE solution

$$\boldsymbol{w}^* = \operatorname*{argmax}_{\boldsymbol{w}} P(\boldsymbol{w}) = \operatorname*{argmax}_{\boldsymbol{w}} \prod_{n=1}^{N} \mathbb{P}(y_n \mid \boldsymbol{x_n}; \boldsymbol{w})$$

The MLE solution

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^{N} \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$
$$= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{n=1}^{N} \ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$
$$= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{n=1}^N \ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N - \ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^{N} \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{n=1}^{N} \ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} - \ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} \ln(1 + e^{-y_n \mathbf{w}^{\mathrm{T}} \mathbf{x_n}})$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^{N} \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{n=1}^{N} \ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} -\ln \mathbb{P}(y_n \mid \mathbf{x_n}; \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} \ln(1 + e^{-y_n \mathbf{w}^{\mathrm{T}} \mathbf{x_n}}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} \ell_{\mathsf{logistic}}(y_n \mathbf{w}^{\mathrm{T}} \mathbf{x_n})$$

$$w^* = \underset{\boldsymbol{w}}{\operatorname{argmax}} P(\boldsymbol{w}) = \underset{\boldsymbol{w}}{\operatorname{argmax}} \prod_{n=1}^{N} \mathbb{P}(y_n \mid \boldsymbol{x_n}; \boldsymbol{w})$$

$$= \underset{\boldsymbol{w}}{\operatorname{argmax}} \sum_{n=1}^{N} \ln \mathbb{P}(y_n \mid \boldsymbol{x_n}; \boldsymbol{w}) = \underset{\boldsymbol{w}}{\operatorname{argmin}} \sum_{n=1}^{N} - \ln \mathbb{P}(y_n \mid \boldsymbol{x_n}; \boldsymbol{w})$$

$$= \underset{\boldsymbol{w}}{\operatorname{argmin}} \sum_{n=1}^{N} \ln(1 + e^{-y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x_n}}) = \underset{\boldsymbol{w}}{\operatorname{argmin}} \sum_{n=1}^{N} \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x_n})$$

$$= \underset{\boldsymbol{w}}{\operatorname{argmin}} F(\boldsymbol{w})$$

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

$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w})$$

$$\begin{split} \boldsymbol{w} &\leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w}) \\ &= \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) \\ \end{split} \qquad \qquad (n \in [N] \text{ is drawn u.a.r.})$$

$$\begin{split} \boldsymbol{w} &\leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w}) \\ &= \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) \qquad \qquad (n \in [N] \text{ is drawn u.a.r.}) \\ &= \boldsymbol{w} - \eta \left(\frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z = y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \end{split}$$

$$\begin{split} & \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w}) \\ &= \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) \qquad (n \in [N] \text{ is drawn u.a.r.}) \\ &= \boldsymbol{w} - \eta \left(\frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ &= \boldsymbol{w} - \eta \left(\frac{-e^{-z}}{1+e^{-z}} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \end{split}$$

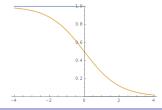
$$\begin{split} & \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w}) \\ &= \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) \qquad (n \in [N] \text{ is drawn u.a.r.}) \\ &= \boldsymbol{w} - \eta \left(\frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ &= \boldsymbol{w} - \eta \left(\frac{-e^{-z}}{1+e^{-z}} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ &= \boldsymbol{w} + \eta \sigma (-y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n \end{split}$$

$$\begin{split} & \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w}) \\ &= \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) \qquad (n \in [N] \text{ is drawn u.a.r.}) \\ &= \boldsymbol{w} - \eta \left(\frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ &= \boldsymbol{w} - \eta \left(\frac{-e^{-z}}{1+e^{-z}} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ &= \boldsymbol{w} + \eta \sigma (-y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n \\ &= \boldsymbol{w} + \eta \mathbb{P}(-y_n \mid \boldsymbol{x}_n; \boldsymbol{w}) y_n \boldsymbol{x}_n \end{split}$$

$$\begin{split} & \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w}) \\ & = \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} \ell_{\text{logistic}}(y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) \qquad (n \in [N] \text{ is drawn u.a.r.}) \\ & = \boldsymbol{w} - \eta \left(\frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ & = \boldsymbol{w} - \eta \left(\frac{-e^{-z}}{1+e^{-z}} \Big|_{z=y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n} \right) y_n \boldsymbol{x}_n \\ & = \boldsymbol{w} + \eta \sigma (-y_n \boldsymbol{w}^{\text{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n \\ & = \boldsymbol{w} + \eta \mathbb{P}(-y_n \mid \boldsymbol{x}_n; \boldsymbol{w}) y_n \boldsymbol{x}_n \end{split}$$

This is a soft version of Perceptron!

$$\mathbb{P}(-y_n|oldsymbol{x}_n;oldsymbol{w})$$
 versus $\mathbb{I}[y_n
eq \mathsf{sgn}(oldsymbol{w}^{\mathrm{T}}oldsymbol{x}_n)]$



$$\nabla_{\boldsymbol{w}} \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) = -\sigma(-y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n$$

$$\nabla_{\boldsymbol{w}} \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) = -\sigma(-y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n$$

$$\nabla_{\boldsymbol{w}}^{2} \ell_{\mathsf{logistic}}(y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n}) = \left(\frac{\partial \sigma(z)}{\partial z} \Big|_{z = -y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n}} \right) y_{n}^{2} \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\mathsf{T}}$$

$$\nabla_{\boldsymbol{w}} \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) = -\sigma(-y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n$$

$$\nabla_{\boldsymbol{w}}^{2} \ell_{\text{logistic}}(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}) = \left(\frac{\partial \sigma(z)}{\partial z}\Big|_{z=-y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}}\right) y_{n}^{2} \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}}$$
$$= \left(\frac{e^{-z}}{(1+e^{-z})^{2}}\Big|_{z=-y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}}\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}}$$

$$\nabla_{\boldsymbol{w}} \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) = -\sigma(-y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n) y_n \boldsymbol{x}_n$$

$$\nabla_{\boldsymbol{w}}^{2} \ell_{\mathsf{logistic}}(y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n}) = \left(\frac{\partial \sigma(z)}{\partial z}\Big|_{z=-y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n}}\right) y_{n}^{2} \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\mathsf{T}} \\
= \left(\frac{e^{-z}}{(1+e^{-z})^{2}}\Big|_{z=-y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n}}\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\mathsf{T}} \\
= \sigma(y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n}) \left(1 - \sigma(y_{n} \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_{n})\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\mathsf{T}}$$

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

$$\nabla_{\boldsymbol{w}}^{2} \ell_{\text{logistic}}(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}) = \left(\frac{\partial \sigma(z)}{\partial z}\Big|_{z=-y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}}\right) y_{n}^{2} \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}}$$

$$= \left(\frac{e^{-z}}{(1+e^{-z})^{2}}\Big|_{z=-y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}}\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}}$$

$$= \sigma(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}) \left(1 - \sigma(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n})\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}}$$

Exercises:

• why is the Hessian of logistic loss positive semidefinite?

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

$$\nabla_{\boldsymbol{w}}^{2} \ell_{\text{logistic}}(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}) = \left(\frac{\partial \sigma(z)}{\partial z}\Big|_{z=-y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}}\right) y_{n}^{2} \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}} \\
= \left(\frac{e^{-z}}{(1+e^{-z})^{2}}\Big|_{z=-y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}}\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}} \\
= \sigma(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n}) \left(1 - \sigma(y_{n} \boldsymbol{w}^{\text{T}} \boldsymbol{x}_{n})\right) \boldsymbol{x}_{n} \boldsymbol{x}_{n}^{\text{T}}$$

Exercises:

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

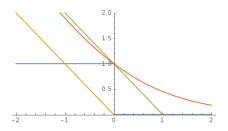
Summary

Linear models for classification:

Step 1. Model is the set of separating hyperplanes

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

Step 2. Pick the surrogate loss



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

Step 3. Find empirical risk minimizer (ERM):

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

using

• GD:
$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla F(\boldsymbol{w})$$

• SGD:
$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w})$$
 $(\mathbb{E}[\tilde{\nabla} F(\boldsymbol{w})] = \nabla F(\boldsymbol{w}))$

• Newton:
$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \left(\nabla^2 F(\boldsymbol{w})\right)^{-1} \nabla F(\boldsymbol{w})$$