

10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Stochastic Gradient Descent

Binary Logistic Regression

Matt Gormley Lecture 9 Feb. 10, 2025

Reminders

- Homework 3: KNN, Perceptron, Lin.Reg.
 - Out: Mon, Feb. 3
 - Due: Mon, Feb. 10 at 11:59pm
 - (only two grace/late days permitted)
- Practice Problems 1
 - released on course website
- Exam 1: Mon, Feb. 17
 - Time: 7:00 9:00pm
 - Location: Your room/seat assignment will be announced on Piazza

OPTIMIZATION METHOD #3: STOCHASTIC GRADIENT DESCENT

Gradient Descent

Algorithm 1 Gradient Descent

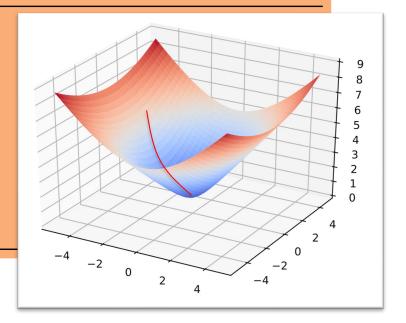
1: **procedure** $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$

 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$

while not converged do 3:

 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla J(\boldsymbol{\theta})$ 4:

return θ



per-example objective:

$$J^{(i)}(\boldsymbol{ heta})$$

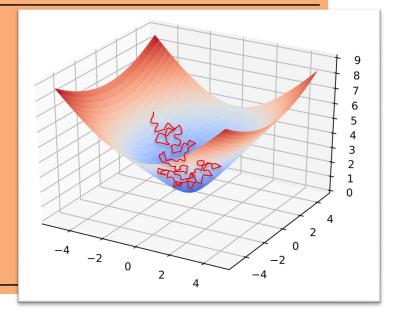
original objective:
$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta})$$

Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD(\mathcal{D}, \boldsymbol{\theta}^{(0)})
```

- 2: $oldsymbol{ heta} \leftarrow oldsymbol{ heta}^{(0)}$
- 3: **while** not converged **do**
- 4: $i \sim \mathsf{Uniform}(\{1, 2, \dots, N\})$
- 5: $oldsymbol{ heta} \leftarrow oldsymbol{ heta} \gamma
 abla_{oldsymbol{ heta}} J^{(i)}(oldsymbol{ heta})$
- 6: return θ



per-example objective:

$$J^{(i)}(oldsymbol{ heta})$$

original objective: $_{\!N}$

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta})$$

Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD(\mathcal{D}, \boldsymbol{\theta}^{(0)})
```

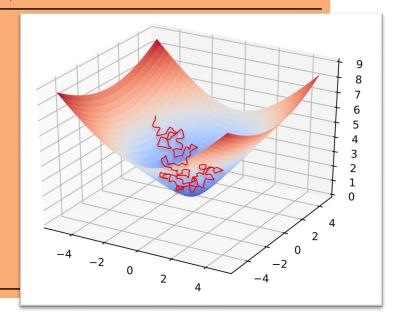
2:
$$oldsymbol{ heta} \leftarrow oldsymbol{ heta}^{(0)}$$

3: **while** not converged **do**

4: for
$$i \in \text{shuffle}(\{1, 2, \dots, N\})$$
 do

5:
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$$

6: return θ



per-example objective:

$$J^{(i)}(oldsymbol{ heta})$$

original objective: $_{\!N}$

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta})$$

In practice, it is common to implement SGD using sampling without replacement (i.e. shuffle({1,2,...N}), even though most of the theory is for sampling with replacement (i.e. Uniform({1,2,...N}).

Why does SGD work?

Background: Expectation of a function of a random variable

For any discrete random variable X

$$E_X[f(X)] = \sum_{x \in \mathcal{X}} P(X = x) f(x)$$

Objective Function for SGD

We assume the form to be:

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta})$$

Expectation of a Stochastic Gradient:

 If the example is sampled uniformly at random, the expected value of the pointwise gradient is the same as the full gradient!

$$E[\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})] = \sum_{i=1}^{N} \left(\text{probability of selecting } \boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)} \right) \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$$

$$= \sum_{i=1}^{N} \left(\frac{1}{N} \right) \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$$

$$= \frac{1}{N} \sum_{i=1}^{N} \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$$

$$= \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

In practice, the data set is randomly shuffled then looped through so that each data point is used equally often

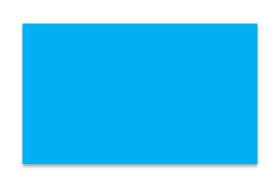
Why does SGD work?

Back of Room																								
																								Right Aisle (Stairs)
												Center Aisle												
G1	G2	G3	G4	G5	G6	G7	G8	G9	G10	G11	G12	(Stairs)		G13	G14	G15	(b)		(b)		G16	G17	G18	
	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10	F11-L			F12	F13	F14	F15	F16	F17	F18	F19	F20	F21-L	
		E1	E2	E3	E4	E5	E6	E7	E8	E9	E10-L		E11	E12	E13	E14	E15	E16	E17	E18	E19	E20	E21-L	
			D1	D2	D3	D4	D5	D6	D7	D8	D9-L		D10	D11	D12	D13	D14	D15	D16	D17	D18	D19	D20-L	
	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11-L		C12	C13	C14	C15	C16	C17	C18	C19	C20	C21	C22-L	
	В1	B2	ВЗ	В4	B5	B6	В7	B8	В9	B10	B11		B12	B13	B14	B15	B16	B17	B18	B19	B20	B21	B22-L	
	(b)		A1	A2	A3	A4	A5	A6	Α7	A8	A9		A10	A11	A12	A13	(b)		&		A14	A15	A16-L	
											Lectern			Table										
Front of Room																								Front Door

Back of Room																								
																								Right Aisle (Stairs)
												Center Aisle												
G1	G2	G3	G4	G5	G6	G7	G8	G9	G10	G11	G12	(Stairs)		G13	G14	G15	&		&		G16	G17	G18	
	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10	F11-L			F12	F13	F14	F15	F16	F17	F18	F19	F20	F21-L	
		E1	E2	E3	E4	E5	E6	E7	E8	E9	E10-L		E11	E12	E13	E14	E15	E16	E17	E18	E19	E20	E21-L	
			D1	D2	D3	D4	D5	D6	D7	D8	D9-L		D10	D11	D12	D13	D14	D15	D16	D17	D18	D19	D20-L	
	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11-L		C12	C13	C14	C15	C16	C17	C18	C19	C20	C21	C22-L	
	B1	B2	ВЗ	B4	B5	B6	В7	В8	В9	B10	B11		B12	B13	B14	B15	B16	B17	B18	B19	B20	B21	B22-L	
	(b)		A1	A2	A3	A4	A5	A6	A7	A8	A9		A10	A11	A12	A13	(b)		(b)		A14	A15	A16-L	
											Lectern			Table										
Front of Room																								Front Door

ack of Room			+	+	+			+	+				+			_		-					Back of Room			_			-									_		-			
																					Righ (Sta	ht Aisle airs)																					Right Ais (Stairs)
										44 040	Center Aisle	040		4 045				240	047	242		,				25 00			240			Center Aisle				045				240	247	040	,
			_	_	_	_	_	_	_	G12 G12	(Stairs)			.4 G15	-	(4	9)	G16			\perp			G2 G3	_		_			_		(Stairs)		313			-	(8)	_			G18	
	F1 F2	2 F3	F4	F5	F6	F7	F8	F9	F	10 F11-L		F12	F13	3 F14	F15 F	16 F	17 F18	F19	F20	F21-L	-			F1 F2	F3	-4 F5	F6	F7 F8	F9	F10	F11-L		F	12	F13	F14	F15 F1	.6 F1	7 F18	F19	F20	F21-L	
	E1	1 E2	2 E3	E4	E5	E6	E7	E8	E	9 E10-L		E11 E12	E13	3 E14	E15 E	16 E	17 E18	E19	E20	E21-L				E1	E2	E3 E4	E5	E6 E7	E8	E9	E10-L		E11 E	12	E13	E14	E15 E1	.6 E1	7 E18	E19	E20	E21-L	
		D1	1 D2	2 D3	D4	D5	D6	D7	D	08 D9-L		D10 D11	D17	.2 D13	D14 D)15 C	16 D17	7 D18	D19	D20-L	L				D1 [D2 D3	D4	D5 D6	D7	D8	D9-L		D10 [011	D12	D13	D14 D1	L5 D1	6 D17	D18	D19	D20-L	
	C1 C2	2 C3	3 C4	1 C5	C6	C7	C8	C9	С	:10 C11-L		C12 C13	C1/	4 C15	C16 C	:17 C	18 C19	C20	C21	C22-L	L			C1 C2	C3 (C4 C5	C6	C7 C8	C9	C10	C11-L		C12 (213	C14	C15	C16 C1	L7 C1	8 C19	C20	C21	C22-L	
	B1 B2	2 B3	3 B4	B5	В6	B7	B8	В9	В	10 B11		B12 B13	B1/	4 B15	B16 E	17 B	18 B19	B20	B21	B22-L				B1 B2	В3	34 B5	В6	B7 B8	В9	B10	B11		B12 E	313	B14	B15	B16 B1	7 B1	8 B19	B20	B21	B22-L	
	&	A1	L A2	2 A3	A4	A5	A6	Α7	Α	.8 A9		A10 A11	A17	2 A13	(a)	(9	A14	A15	A16-L	-			&	A1 /	A2 A3	A4	A5 A6	A7	A8	A9		A10 A	\11	A12	A13	&	(4)		A14	A15	A16-L	
																																				اا							
							T			Lectern		Table	٤																		Lectern		1	able									
ront of Room																					Fron	nt Door	Front of Room																				Front Do

Back of Room																								
																								Right Aisle (Stairs)
												Center Aisle												
G1	G2	G3	G4	G5	G6	G7	G8	G9	G10	G11	G12	(Stairs)		G13	G14	G15	(b)		(b)		G16	G17	G18	
	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10	F11-L			F12	F13	F14	F15	F16	F17	F18	F19	F20	F21-L	
		E1	E2	E3	E4	E5	E6	E7	E8	E9	E10-L		E11	E12	E13	E14	E15	E16	E17	E18	E19	E20	E21-L	
			D1	D2	D3	D4	D5	D6	D7	D8	D9-L		D10	D11	D12	D13	D14	D15	D16	D17	D18	D19	D20-L	
	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11-L		C12	C13	C14	C15	C16	C17	C18	C19	C20	C21	C22-L	
	В1	B2	ВЗ	В4	B5	В6	В7	B8	В9	B10	B11		B12	B13	B14	B15	B16	B17	B18	B19	B20	B21	B22-L	
	(b)		A1	A2	A3	A4	A5	A6	Α7	A8	A9		A10	A11	A12	A13	(b)		(b)		A14	A15	A16-L	
											Lectern			Table										
Front of Room														. 0.010										Front Door



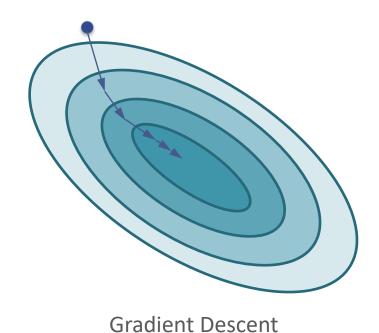


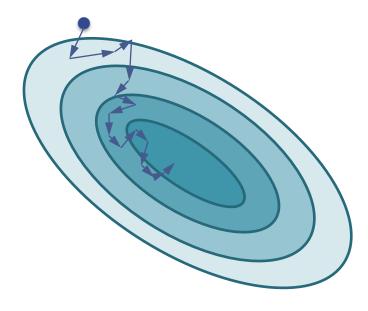




SGD VS. GRADIENT DESCENT

SGD vs. Gradient Descent

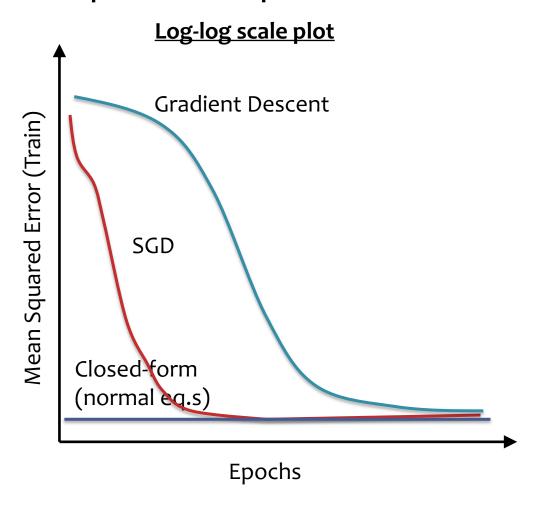




Stochastic Gradient Descent

SGD vs. Gradient Descent

• Empirical comparison:



- Def: an **epoch** is a single pass through the training data
- For GD, only one update per epoch
- For SGD, N updates
 per epoch
 N = (# train examples)
- SGD reduces MSE much more rapidly than GD
- For GD / SGD, training MSE is initially large due to uninformed initialization

SGD vs. Gradient Descent

Theoretical comparison:

Define convergence to be when $J(\boldsymbol{\theta}^{(t)}) - J(\boldsymbol{\theta}^*) < \epsilon$

Method	Steps to Convergence	Computation per Step
Gradient descent	$O(\log 1/\epsilon)$	O(NM)
SGD	$o(1/\epsilon)$	O(M)

(with high probability under certain assumptions)

Main Takeaway: SGD has much slower asymptotic convergence (i.e. it's slower in theory), but is often much faster in practice.

SGD FOR LINEAR REGRESSION

Linear Regression as Function Approximation

$$\mathcal{D}=\{\mathbf{x}^{(i)},y^{(i)}\}_{i=1}^N$$
 where $\mathbf{x}\in\mathbb{R}^M$ and $y\in\mathbb{R}$ 1. Assume \mathcal{D} generated as:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$
$$y^{(i)} = h^*(\mathbf{x}^{(i)})$$

2. Choose hypothesis space, \mathcal{H} : all linear functions in M-dimensional space

$$\mathcal{H} = \{h_{\boldsymbol{\theta}} : h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^M \}$$

3. Choose an objective function: mean squared error (MSE)

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} e_i^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2$$

- 4. Solve the unconstrained optimization problem via favorite method:
 - gradient descent
 - closed form
 - stochastic gradient descent

$$\hat{oldsymbol{ heta}} = rgmin_{oldsymbol{ heta}} J(oldsymbol{ heta})$$

5. Test time: given a new x, make prediction \hat{y}

$$\hat{y} = h_{\hat{oldsymbol{ heta}}}(\mathbf{x}) = \hat{oldsymbol{ heta}}^T \mathbf{x}$$

Gradient Calculation for Linear Regression

Derivative of
$$J^{(i)}(\boldsymbol{\theta})$$
:
$$\frac{d}{d\theta_k}J^{(i)}(\boldsymbol{\theta}) = \frac{d}{d\theta_k}\frac{1}{2}(\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2}\frac{d}{d\theta_k}(\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})^2$$

$$= (\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})\frac{d}{d\theta_k}(\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})$$

$$= (\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})\frac{d}{d\theta_k}\left(\sum_{j=1}^K \theta_j x_j^{(i)} - y^{(i)}\right)$$

$$= (\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})x_k^{(i)}$$

$$\begin{aligned} &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)} \\ & \text{Gradient of } J^{(i)}(\boldsymbol{\theta}) & \text{[used by SGD]} \\ & \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{d}{d\theta_1} J^{(i)}(\boldsymbol{\theta}) \\ \frac{d}{d\theta_2} J^{(i)}(\boldsymbol{\theta}) \\ \vdots \\ \frac{d}{d\theta_M} J^{(i)}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_1^{(i)} \\ (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_2^{(i)} \\ \vdots \\ (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_N^{(i)} \end{bmatrix} \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \end{aligned}$$

Derivative of
$$J(\boldsymbol{\theta})$$
:
$$\frac{d}{d\theta_k}J(\boldsymbol{\theta}) = \frac{1}{N}\sum_{i=1}^N \frac{d}{d\theta_k}J^{(i)}(\boldsymbol{\theta})$$
$$= \frac{1}{N}\sum_{i=1}^N (\boldsymbol{\theta}^T\mathbf{x}^{(i)} - y^{(i)})x_k^{(i)}$$

SGD for Linear Regression

SGD applied to Linear Regression is called the "Least Mean Squares" algorithm

GD for Linear Regression

Gradient Descent for Linear Regression repeatedly takes steps opposite the gradient of the objective function

Algorithm 1 GD for Linear Regression 1: procedure GDLR(\mathcal{D} , $\theta^{(0)}$) 2: $\theta \leftarrow \theta^{(0)}$ > Initialize parameters 3: while not converged do 4: $\mathbf{g} \leftarrow \frac{1}{N} \sum_{i=1}^{N} (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$ > Compute gradient 5: $\theta \leftarrow \theta - \gamma \mathbf{g}$ > Update parameters 6: return θ

Solving Linear Regression

Question:

True or False: If Mean Squared Error (i.e. $\frac{1}{N}\sum_{i=1}^{N}(y^{(i)}-h(\mathbf{x}^{(i)}))^2$) has a unique minimizer (i.e. argmin), then Mean Absolute Error (i.e. $\frac{1}{N}\sum_{i=1}^{N}|y^{(i)}-h(\mathbf{x}^{(i)})|$) must also have a unique minimizer.

Answer:

Optimization Objectives

You should be able to...

- Apply gradient descent to optimize a function
- Apply stochastic gradient descent (SGD) to optimize a function
- Apply knowledge of zero derivatives to identify a closedform solution (if one exists) to an optimization problem
- Distinguish between convex, concave, and nonconvex functions
- Obtain the gradient (and Hessian) of a (twice) differentiable function

PROBABILISTIC LEARNING

Probabilistic Learning

Function Approximation

Previously, we assumed that our output was generated using a **deterministic target function**:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$

$$y^{(i)} = c^*(\mathbf{x}^{(i)})$$

Our goal was to learn a hypothesis h(x) that best approximates $c^*(x)$

Probabilistic Learning

Today, we assume that our output is sampled from a conditional probability distribution:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$
$$y^{(i)} \sim p^*(\cdot|\mathbf{x}^{(i)})$$

Our goal is to learn a probability distribution p(y|x) that best approximates $p^*(y|x)$

MAXIMUM LIKELIHOOD ESTIMATION

Likelihood Function

One R.V.

Given N independent, identically distributed (iid) samples $D = \{x^{(1)}, x^{(2)}, ..., x^{(N)}\}$ from a discrete random variable X with probability mass function $(pmf) p(x|\theta) ...$

• Case 1: The **likelihood** function $L(\theta) = p(x^{(1)}|\theta) p(x^{(2)}|\theta) \dots p(x^{(N)}|\theta)$ The **likelihood** tells us how likely one sample is relative to another

• <u>Case 2</u>: The **log-likelihood** function is $\ell(\theta) = \log p(x^{(1)}|\theta) + ... + \log p(x^{(N)}|\theta)$

Likelihood Function

Two R.V.s

Given N **iid** samples D = $\{(x^{(1)}, y^{(1)}), ..., (x^{(N)}, y^{(N)})\}$ from a pair of **random variables** X, Y where Y is **discrete** with probability mass function (pmf) $p(y \mid x, \theta)$

• <u>Case 3</u>: The **conditional likelihood** function:

$$L(\theta) = p(y^{(1)} | x^{(1)}, \theta) ... p(y^{(N)} | x^{(N)}, \theta)$$

• Case 4: The conditional log-likelihood function is $\ell(\theta) = \log p(y^{(1)} | x^{(1)}, \theta) + ... + \log p(y^{(N)} | x^{(N)}, \theta)$

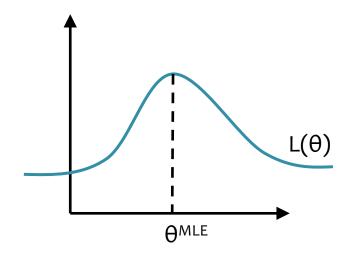
Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

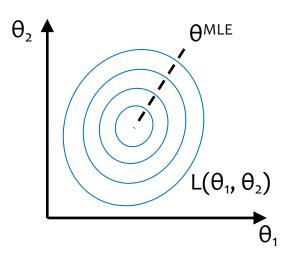
Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the likelihood of the data. N

$$\boldsymbol{\theta}^{\mathsf{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^{n} p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)





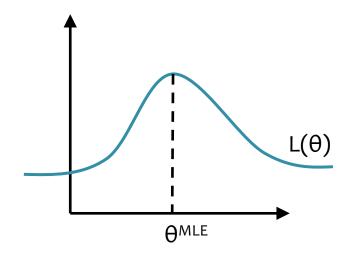
Suppose we have data $\mathcal{D} = \{(y^{(i)}, \mathbf{x}^{(i)})\}_{i=1}^{N}$

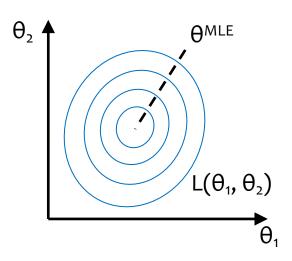
Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the conditional likelihood of the data. N

$$\boldsymbol{\theta}^{\mathsf{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^{n} p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)





Suppose we have data $\mathcal{D} = \{(y^{(i)}, \mathbf{x}^{(i)})\}_{i=1}^N$

Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the conditional log-likelihood

of the data.

$$\boldsymbol{\theta}^{\mathsf{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1} p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

Suppose we have data $\mathcal{D} = \{(y^{(i)}, \mathbf{x}^{(i)})\}_{i=1}^N$

Principle of Maximum Likelihood Estimation:

Choose the parameters that minimize the negative conditional log-

likelihood of the data.
$$\pmb{\theta}^{\mathsf{MLE}} = \argmax_{\pmb{\theta}} \prod_{i=1}^{N} p(y^{(i)} \mid \mathbf{x}^{(i)}, \pmb{\theta})$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \sum_{i=1}^{N} \log p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

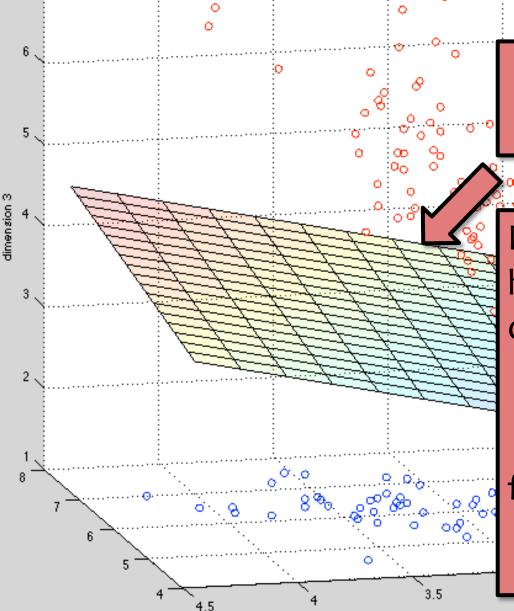
What does maximizing likelihood accomplish?

- There is only a finite amount of probability mass (i.e. sum-to-one constraint)
- MLE tries to allocate as much probability mass as possible to the things we have observed...

... at the expense of the things we have not observed

LOGISTIC REGRESSION

Linear Models for Classification



Key idea: Try to learn this hyperplane directly

Directly modeling the hyperplane would use a decision function:

$$h(\mathbf{x}) = \mathsf{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

for:

$$y \in \{-1, +1\}$$

Background: Hyperplanes

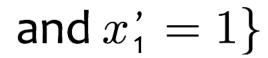
Notation Trick: fold the bias b and the weights w into a single vector $\boldsymbol{\theta}$ by prepending a constant to x and increasing dimensionality by one to get x'!

Hyperplane (Definition 1):

$$\mathcal{H} = \{ \mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0 \}$$

Hyperplane (Definition 2):

$$\mathcal{H} = \{\mathbf{x}': \boldsymbol{\theta}^T \mathbf{x}' = 0$$



$$oldsymbol{ heta} = [b, w_1, \dots, w_M]^T \\ \mathbf{x}' = [1, x_1, \dots, x_M]^T$$

Half-spaces:

$$\mathcal{H}^+ = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} > 0 \text{ and } x_0^1 = 1\}$$

$$\mathcal{H}^- = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} < 0 \text{ and } x_0^1 = 1\}$$

Using gradient descent for linear classifiers

Key idea behind today's lecture:

- 1. Define a linear classifier (logistic regression)
- 2. Define an objective function (likelihood)
- 3. Optimize it with gradient descent to learn parameters
- 4. Predict the class with highest probability under the model

Data: Inputs are continuous vectors of length M. Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \text{ where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{0, 1\}$$

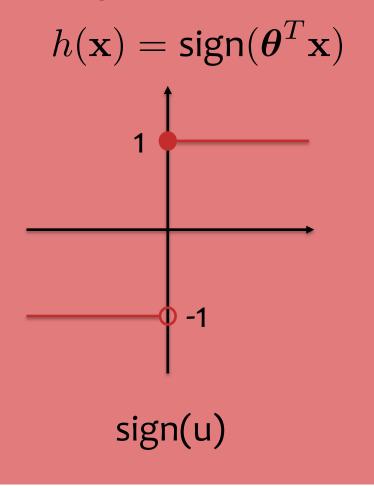


We are back to classification.

Despite the name logistic **regression**.

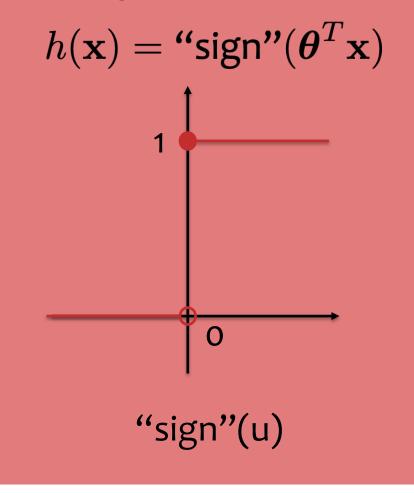
$sign(\cdot) vs. sigmoid(\cdot)$

Suppose we wanted to learn a linear classifier, but instead of predicting $y \in \{-1,+1\}$ we wanted to predict $y \in \{0,1\}$



$sign(\cdot) vs. sigmoid(\cdot)$

Suppose we wanted to learn a linear classifier, but instead of predicting $y \in \{-1,+1\}$ we wanted to predict $y \in \{0,1\}$

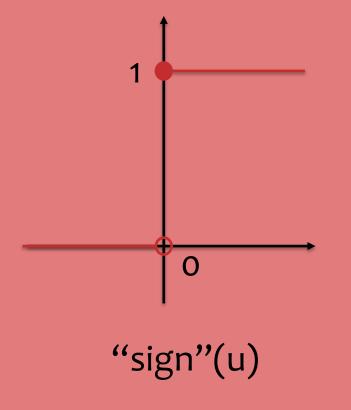


Goal: Learn a linear classifier with Gradient Descent

$sign(\cdot) vs. sigmoid(\cdot)$

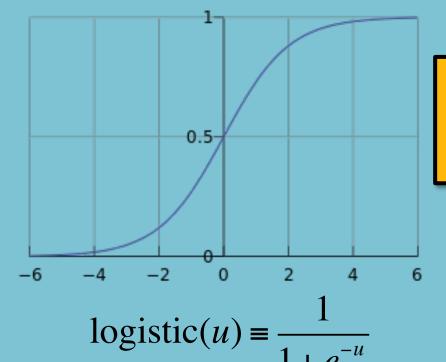


$$h(\mathbf{x}) = \text{"sign"}(\boldsymbol{\theta}^T \mathbf{x})$$



Use a differentiable function instead!

$$p_{\boldsymbol{\theta}}(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



The logistic function is also called the sigmoid function.

Data: Inputs are continuous vectors of length M. Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$$
 where $\mathbf{x} \in \mathbb{R}^M$ and $y \in \{0, 1\}$

Model: Logistic function applied to dot product of parameters with input vector.

$$p_{\boldsymbol{\theta}}(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

Learning: finds the parameters that minimize some objective function. $m{ heta}^* = \operatorname*{argmin} J(m{ heta})$

Prediction: Output is the most probable class.

$$\hat{y} = \underset{y \in \{0,1\}}{\operatorname{argmax}} p_{\theta}(y|\mathbf{x})$$

Learning Logistic Regression

Learning: Four approaches to solving $\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$

Approach o: Random Search (horridly slow because it lacks gradient information)

Approach 1: Gradient Descent (take large confident steps opposite the gradient)

Approach 2: Stochastic Gradient Descent (SGD) (take many small steps roughly opposite the gradient)

Approach. Closed Form

(set derivatives equal to zero and solve for parameters)

Logistic Regression does not have a closed form solution for MLE parameters.

1. Model 2. Objective

3A. Derivatives

3B. Gradients

4. Optimization

5. Prediction