



10-301/601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Stochastic Gradient Descent

+

Probabilistic Learning (Binary Logistic Regression)

Matt Gormley
Lecture 9
Feb. 16, 2022

Reminders

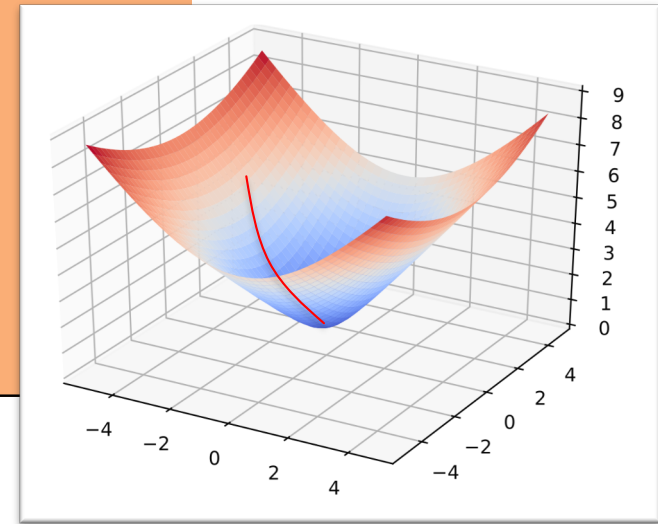
- **Practice for Exam 1**
 - **Mock Exam 1**
 - Due: Wed, Feb. 16 at 11:59pm
 - See [@683](#) for participation point details
 - **Practice Problems 1 released on course website**
- **Exam 1: Thu, Feb. 17**
 - **Time: 6:30 – 8:30pm**
 - **Location: Your room/seat assignment will be announced on Piazza**
- **Homework 4: Logistic Regression**
 - **Out: Fri, Feb 18**
 - **Due: Sun, Feb. 27 at 11:59pm**

OPTIMIZATION METHOD #3: STOCHASTIC GRADIENT DESCENT

Gradient Descent

Algorithm 1 Gradient Descent

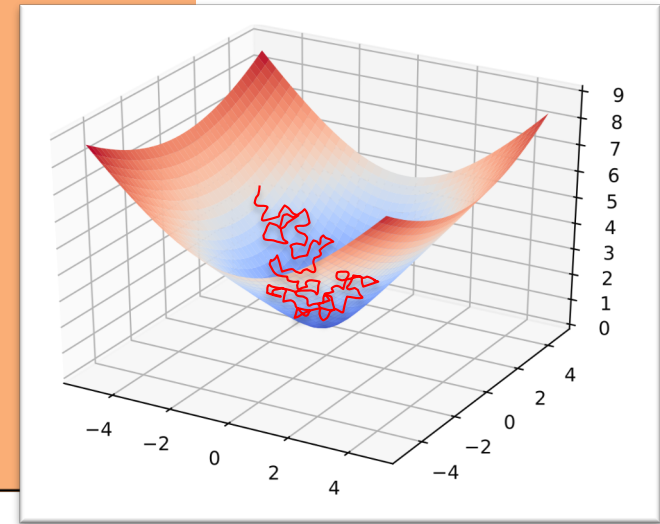
```
1: procedure GD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```



Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $i \sim \text{Uniform}(\{1, 2, \dots, N\})$ 
5:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



per-example objective:

$$J^{(i)}(\theta)$$

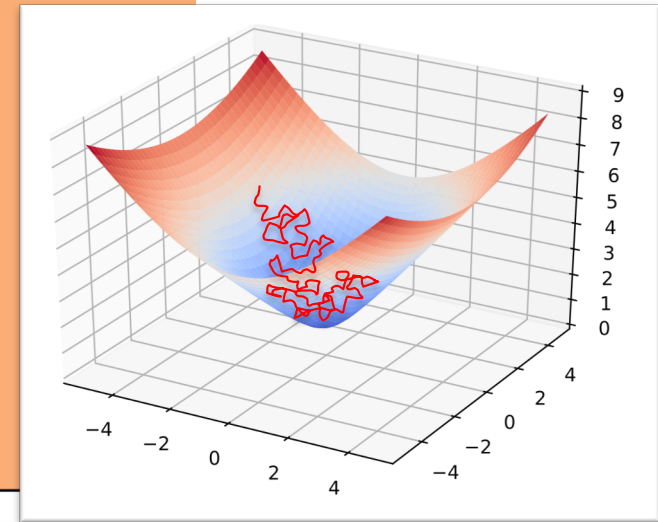
original objective:

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

Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



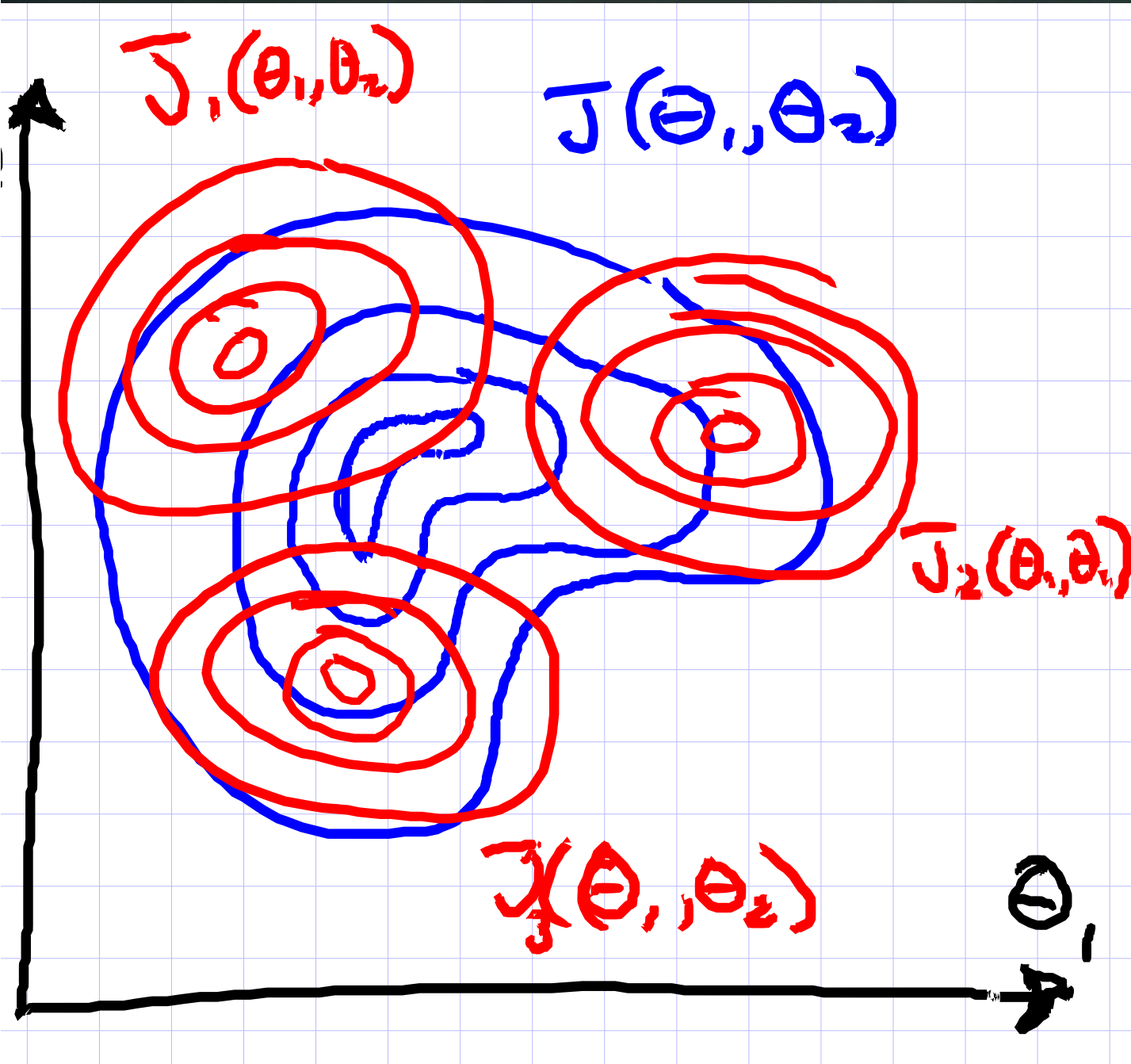
per-example objective:

$$J^{(i)}(\theta)$$

original objective:

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

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



Why does SGD work?

Whiteboard

- Example of SGD on simple 2D function

Why does SGD work?

$$\frac{dJ(\theta)}{d\theta_j} = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\theta_j} (J_i(\theta))$$
$$\nabla J(\theta) = \begin{bmatrix} \vdots \\ \text{jth} \\ \vdots \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J_i(\theta)$$

Recall: for any discrete r.v. X

$$E_X[f(x)] \triangleq \sum_x P(X=x) f(x)$$

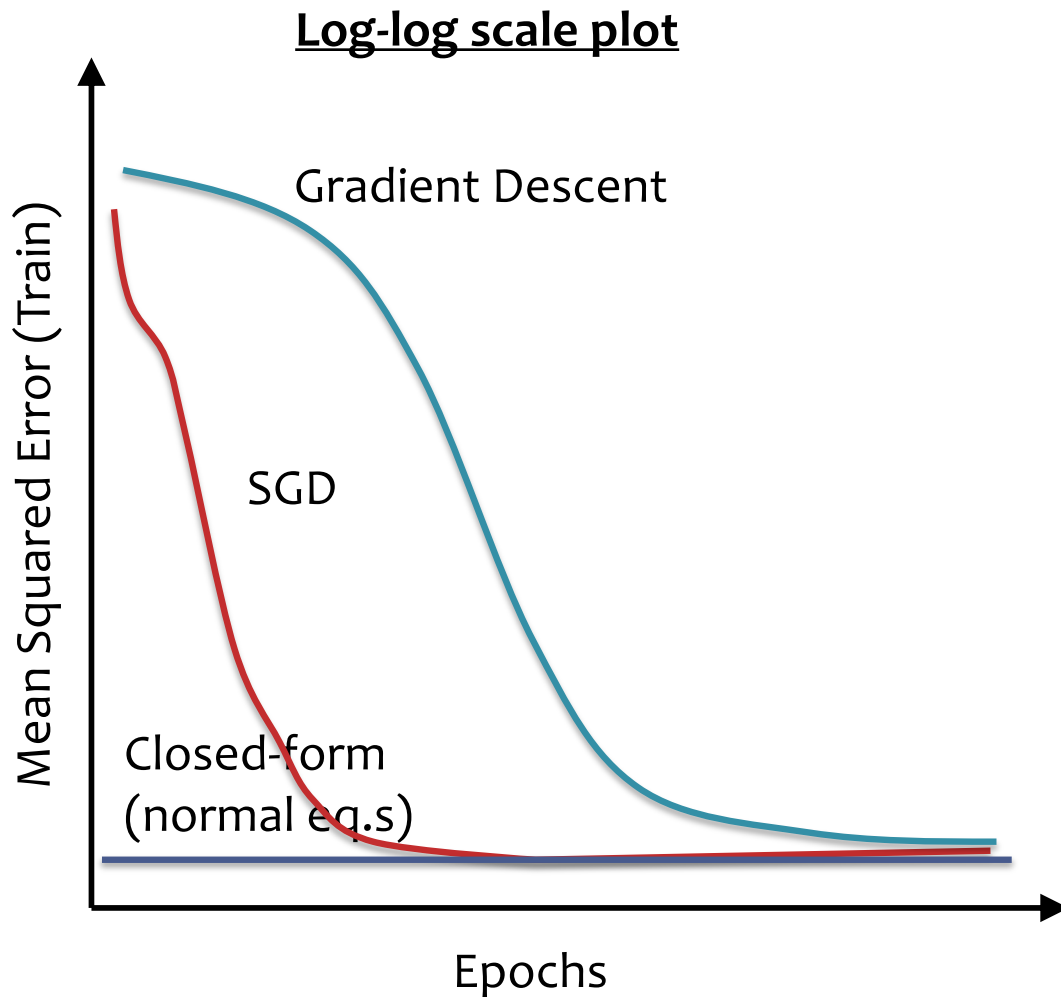
Q: What is the expected value of a randomly chosen $\nabla J_i(\theta)$?

Let $I \sim \text{Uniform}(\{1, \dots, N\})$
 $\Rightarrow P(I=i) = \frac{1}{N}$ if $i \in \{1, \dots, N\}$

$$E_I[\nabla J_I(\theta)] = \sum_{i=1}^N P(I=i) \nabla J_i(\theta)$$
$$> \frac{1}{N} \sum_{i=1}^N \nabla J_i(\theta)$$
$$= \nabla J(\theta)$$

LINEAR REGRESSION: PRACTICALITIES

Empirical Convergence



- Def: an **epoch** is a single pass through the training data
- 1. For GD, only **one update** per epoch
- 2. For SGD, **N updates** per epoch
 $N = (\# \text{ train examples})$

- SGD reduces MSE much more rapidly than GD
- For GD / SGD, training MSE is initially large due to uninformed initialization

Convergence of Optimizers

Convergence Analysis:

Def: Convergence is when $J(\vec{\theta}) - J(\vec{\theta}^*) < \epsilon$

↖ true unknown min

Methods	Steps to Converge	Computation per iteration
Newton's Method	$O(\ln \ln 1/\epsilon)$	$\nabla^2 J(\theta)$ $\nabla^2 J(\theta) \leftarrow O(NM^2)$
GD	$O(\ln 1/\epsilon)$	$\nabla J(\theta) \leftarrow O(NM)$
SGD	$O(1/\epsilon)$	$\nabla J_i(\theta) \leftarrow O(M)$

not correct → (pointing to Newton's Method)

"almost sure" convergence (pointing to SGD)

lots of caveats and conditions (pointing to SGD)

very less computation (pointing to SGD)

Takeaway: SGD has much slower asymptotic convergence. but is often 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:

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

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)

$$\begin{aligned}J(\boldsymbol{\theta}) &= \frac{1}{N} \sum_{i=1}^N e_i^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)^2\end{aligned}$$

4. Solve the unconstrained optimization problem via favorite method:

- gradient descent
- closed form
- stochastic gradient descent
- ...

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$

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

$$\hat{y} = h_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$

Gradient Calculation for Linear Regression

Derivative of $J^{(i)}(\boldsymbol{\theta})$:

$$\begin{aligned} \frac{d}{d\theta_k} J^{(i)}(\boldsymbol{\theta}) &= \frac{d}{d\theta_k} \frac{1}{2} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2} \frac{d}{d\theta_k} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} \left(\sum_{j=1}^K \theta_j x_j^{(i)} - y^{(i)} \right) \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)} \end{aligned}$$

Gradient of $J^{(i)}(\boldsymbol{\theta})$

[used by SGD]

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}) &= \begin{bmatrix} \frac{d}{d\theta_1} J^{(i)}(\boldsymbol{\theta}) \\ \frac{d}{d\theta_2} J^{(i)}(\boldsymbol{\theta}) \\ \vdots \\ \frac{d}{d\theta_M} J^{(i)}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_1^{(i)} \\ (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_2^{(i)} \\ \vdots \\ (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_N^{(i)} \end{bmatrix} \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \end{aligned}$$

Derivative of $J(\boldsymbol{\theta})$:

$$\begin{aligned} \frac{d}{d\theta_k} J(\boldsymbol{\theta}) &= \sum_{i=1}^N \frac{d}{d\theta_k} J^{(i)}(\boldsymbol{\theta}) \\ &= \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)} \end{aligned}$$

Gradient of $J(\boldsymbol{\theta})$

[used by Gradient Descent]

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \begin{bmatrix} \frac{d}{d\theta_1} J(\boldsymbol{\theta}) \\ \frac{d}{d\theta_2} J(\boldsymbol{\theta}) \\ \vdots \\ \frac{d}{d\theta_M} J(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_1^{(i)} \\ \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_2^{(i)} \\ \vdots \\ \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_N^{(i)} \end{bmatrix} \\ &= \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \end{aligned}$$

SGD for Linear Regression

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

Algorithm 1 Least Mean Squares (LMS)

```
1: procedure LMS( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$  ▷ Initialize parameters
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\mathbf{g} \leftarrow (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$  ▷ Compute gradient
6:        $\theta \leftarrow \theta - \gamma \mathbf{g}$  ▷ Update parameters
7:   return  $\theta$ 
```

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( $\mathcal{D}, \boldsymbol{\theta}^{(0)}$ )
2:    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$                                 ▷ Initialize parameters
3:   while not converged do
4:      $\mathbf{g} \leftarrow \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$     ▷ Compute gradient
5:      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \mathbf{g}$                                 ▷ Update parameters
6:   return  $\boldsymbol{\theta}$ 
```

Linear Regression Objectives

You should be able to...

- Design k-NN Regression and Decision Tree Regression
- Implement learning for Linear Regression using three optimization techniques: (1) closed form, (2) gradient descent, (3) stochastic gradient descent
- Choose a Linear Regression optimization technique that is appropriate for a particular dataset by analyzing the tradeoff of computational complexity vs. convergence speed
- Distinguish the three sources of error identified by the bias-variance decomposition: bias, variance, and irreducible error.

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 closed-form 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 $h(\mathbf{x})$ that best approximates $c^*(\mathbf{x})$

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 $p(y|\mathbf{x})$ that best approximates $p^*(y|\mathbf{x})$

Robotic Farming

	Deterministic	Probabilistic
Classification (binary output)	Is this a picture of a wheat kernel?	Is this plant drought resistant?
Regression (continuous output)	How many wheat kernels are in this picture?	What will the yield of this plant be?



MAXIMUM LIKELIHOOD ESTIMATION

MLE

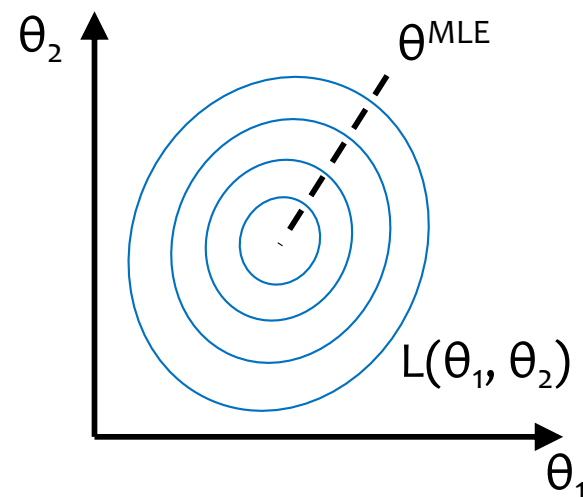
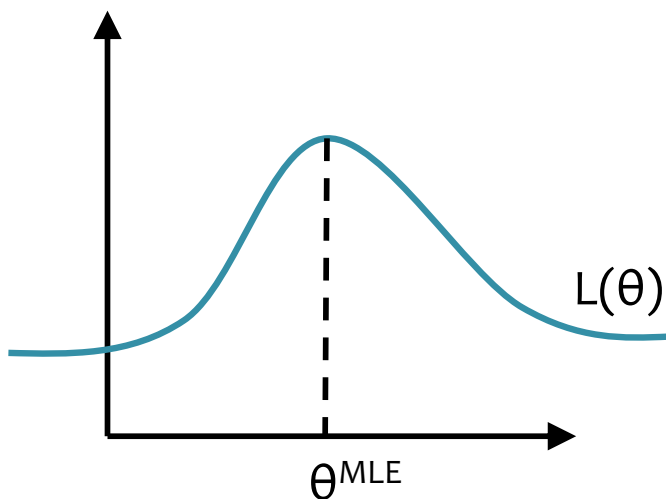
Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$

Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the likelihood of the data.

$$\theta^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p(\mathbf{x}^{(i)} | \theta)$$

Maximum Likelihood Estimate (MLE)



MLE

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

Maximum Likelihood Estimation

The principle of Maximum Likelihood Estimation (MLE):

Choose parameters that make the data "most likely".

Assumptions: Data generated iid from distribution $p^*(x | \vec{\theta}^*)$
and comes from a family of distributions parameterized
 $\theta \in \Theta$ \leftarrow set of possible parameters

Formally:

$$\theta_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} p(D | \theta)$$

$$= \underset{\theta \in \Theta}{\operatorname{argmax}} \log p(D | \theta)$$

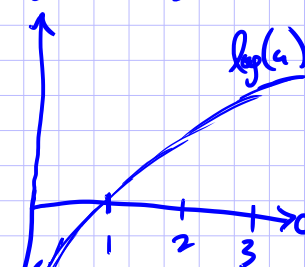
$$= \underset{\theta \in \Theta}{\operatorname{argmax}} l(\theta)$$

usually a continuous optimization \rightarrow

where $l(\theta) \triangleq \log p(D | \theta)$
 \leftarrow "log-likelihood"

\leftarrow treat as function of θ
where D is constant

since log is monotonic



$$\begin{aligned} \log(a_1) &< \log(a_2) \\ \text{iff } a_1 &< a_2 \\ \Rightarrow \log(f(a_1)) &< \log(f(a_2)) \\ \text{iff } f(a_1) &< f(a_2) \end{aligned}$$

MOTIVATION: LOGISTIC REGRESSION

Example: Image Classification

- ImageNet LSVRC-2010 contest:
 - **Dataset:** 1.2 million labeled images, 1000 classes
 - **Task:** Given a new image, label it with the correct class
 - **Multiclass** classification problem
- Examples from <http://image-net.org/>

Bird

Warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings

2126 pictures

92.85% Popularity Percentile

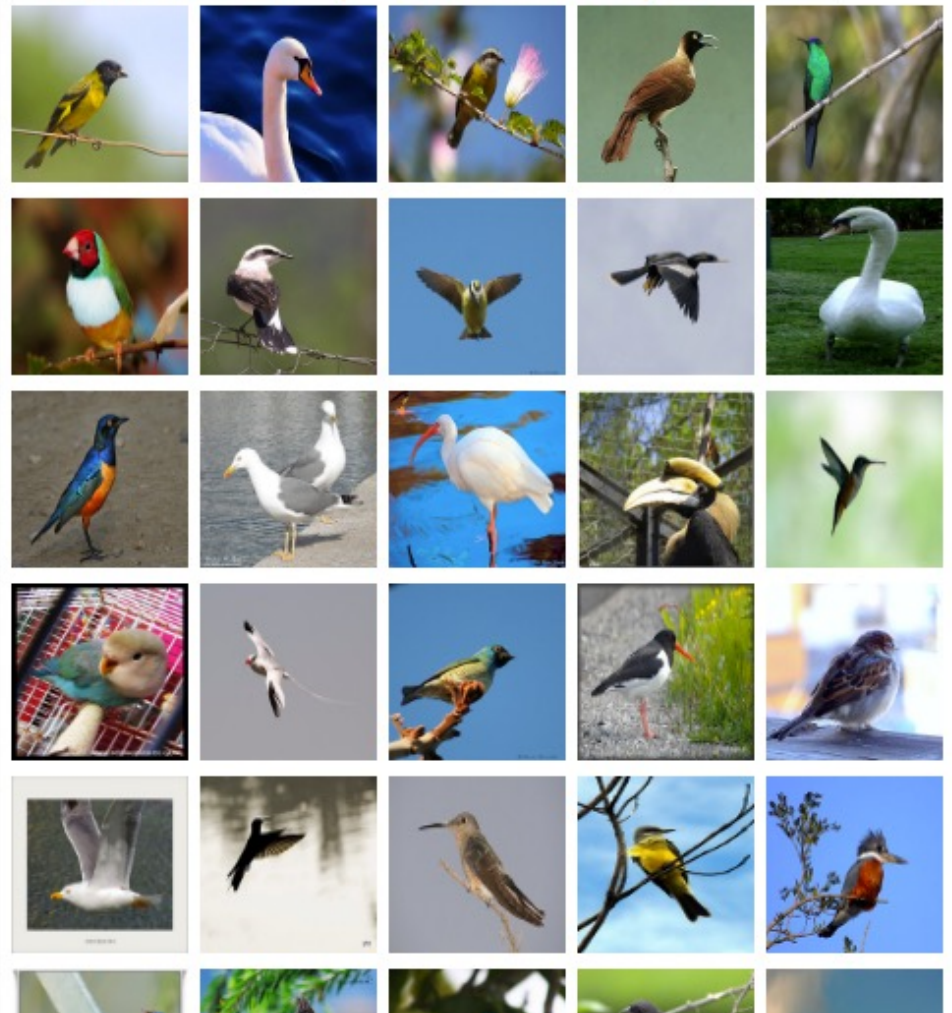


- └ marine animal, marine creature, sea animal, sea creature (1)
- └ scavenger (1)
- └ biped (0)
- └ predator, predatory animal (1)
- └ larva (49)
- └ acrodont (0)
- └ feeder (0)
- └ stunt (0)
- └ **chordate (3087)**
 - └ tunicate, urochordate, urochord (6)
 - └ cephalochordate (1)
 - └ **vertebrate, craniate (3077)**
 - └ mammal, mammalian (1169)
 - └ **bird (871)**
 - └ dickeybird, dickey-bird, dickybird, dicky-bird (0)
 - └ cock (1)
 - └ hen (0)
 - └ nester (0)
 - └ night bird (1)
 - └ bird of passage (0)
 - └ protoavis (0)
 - └ archaeopteryx, archeopteryx, Archaeopteryx lithographi Sinornis (0)
 - └ Ibero-mesornis (0)
 - └ archaeornis (0)
 - └ ratite, ratite bird, flightless bird (10)
 - └ carinate, carinate bird, flying bird (0)
 - └ passerine, passeriform bird (279)
 - └ nonpasserine bird (0)
 - └ bird of prey, raptor, raptorial bird (80)
 - └ gallinaceous bird, gallinacean (114)

Treemap Visualization

Images of the Synset

Downloads



German iris, *Iris kochii*

Iris of northern Italy having deep blue-purple flowers; similar to but smaller than *Iris germanica*

469 pictures

49.6% Popularity Percentile

Wordnet IDs

- halophyte (0)
- succulent (39)
- cultivar (0)
- cultivated plant (0)
- weed (54)
- evergreen, evergreen plant (0)
- deciduous plant (0)
- vine (272)
- creeper (0)
- woody plant, ligneous plant (1868)
- geophyte (0)
- desert plant, xerophyte, xerophytic plant, xerophile, xerophilic mesophyte, mesophytic plant (0)
- aquatic plant, water plant, hydrophyte, hydrophytic plant (11)
- tuberous plant (0)
- bulbous plant (179)
 - iridaceous plant (27)
 - iris, flag, fleur-de-lis, sword lily (19)
 - bearded iris (4)
 - Florentine iris, orris, *Iris germanica florentina*, *Iris*
 - German iris, *Iris germanica* (0)
 - German iris, *Iris kochii* (0)
 - Dalmatian iris, *Iris pallida* (0)
 - beardless iris (4)
 - bulbous iris (0)
 - dwarf iris, *Iris cristata* (0)
 - stinking iris, gladdon, gladdon iris, stinking gladwyn, Persian iris, *Iris persica* (0)
 - yellow iris, yellow flag, yellow water flag, *Iris pseudo*
 - dwarf iris, vernal iris, *Iris verna* (0)
 - blue flag, *Iris versicolor* (0)

Treemap Visualization

Images of the Synset

Downloads



Court, courtyard

An area wholly or partly surrounded by walls or buildings; "the house was built around an inner court"

165
pictures

92.61%
Popularity
Percentile



Numbers in brackets: (the number of synsets in the subtree).

- ImageNet 2011 Fall Release (32326)
 - plant, flora, plant life (4486)
 - geological formation, formation (175)
 - natural object (1112)
 - sport, athletics (176)
 - artifact, artefact (10504)
 - instrumentality, instrumentation (5494)
 - structure, construction (1405)
 - airdock, hangar, repair shed (0)
 - altar (1)
 - arcade, colonnade (1)
 - arch (31)
 - area (344)
 - aisle (0)
 - auditorium (1)
 - baggage claim (0)
 - box (1)
 - breakfast area, breakfast nook (0)
 - bullpen (0)
 - chancel, sanctuary, bema (0)
 - choir (0)
 - corner, nook (2)
 - court, courtyard (6)
 - atrium (0)
 - bailey (0)
 - cloister (0)
 - food court (0)
 - forecourt (0)
 - narvis (0)

Treemap Visualization	Images of the Synset	Downloads

Example: Image Classification

CNN for Image Classification

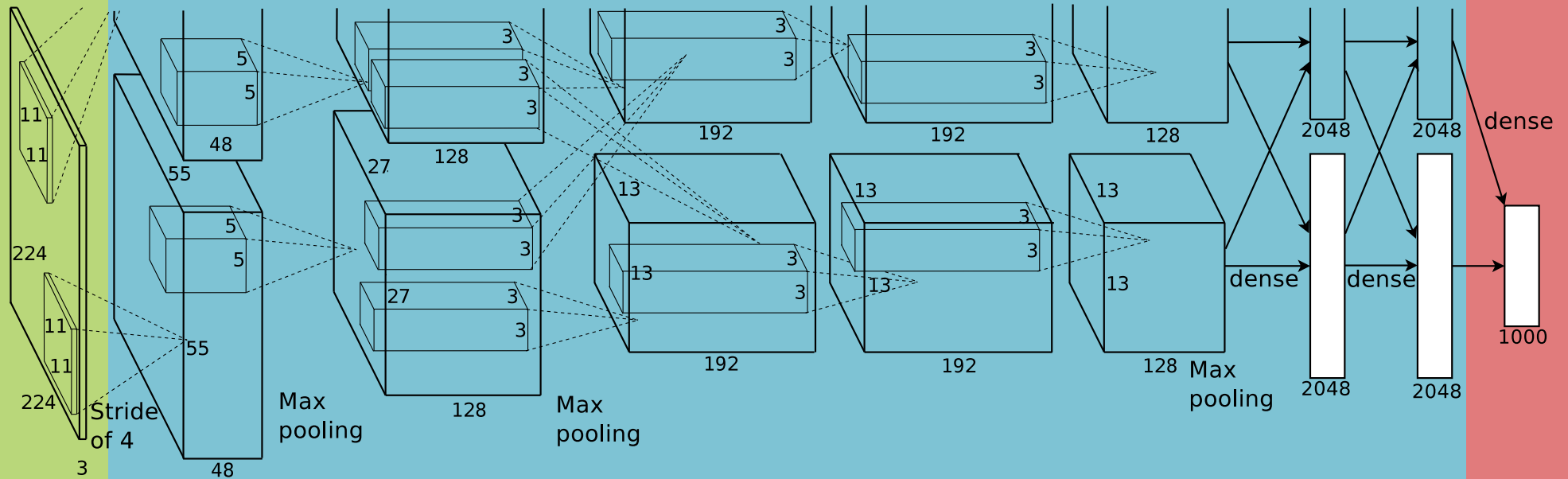
(Krizhevsky, Sutskever & Hinton, 2011)

17.5% error on ImageNet LSVRC-2010 contest

Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



Example: Image Classification

CNN for Image Classification

(Krizhevsky, Sutskever & Hinton, 2011)

17.5% error on ImageNet LSVRC-2010 contest

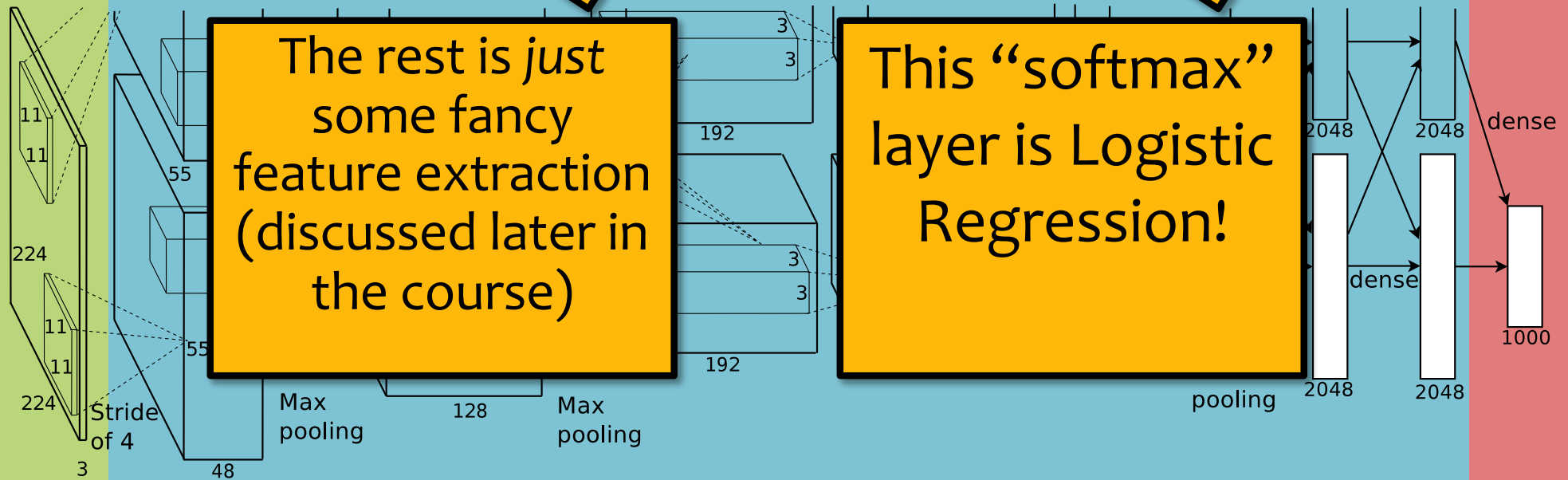
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax

The rest is just some fancy feature extraction (discussed later in the course)

This “softmax” layer is Logistic Regression!




LOGISTIC REGRESSION

Logistic Regression

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.

Recall...

Linear Models for Classification

Key idea: Try to learn this hyperplane directly

Looking ahead:

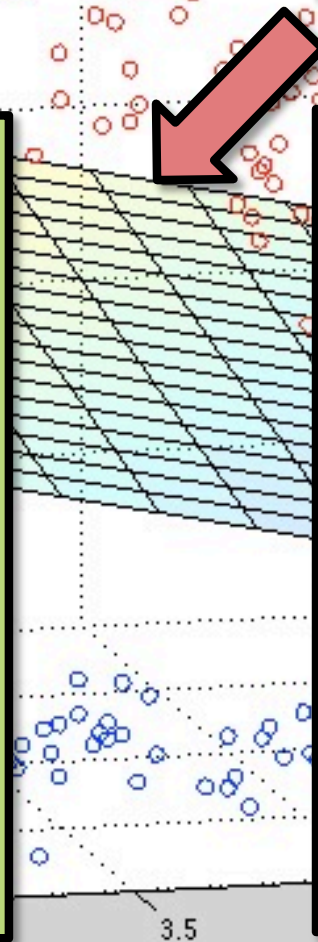
- We'll see a number of commonly used Linear Classifiers
- These include:
 - Perceptron
 - Logistic Regression
 - Naïve Bayes (under certain conditions)
 - Support Vector Machines

Directly modeling the hyperplane would use a decision function:

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

for:

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



Recall...

Background: Hyperplanes

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

Hyperplane (Definition 1):

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

Hyperplane (Definition 2):

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

$$\text{and } x'_1 = 1\}$$

$$\boldsymbol{\theta} = [b, w_1, \dots, w_M]^T$$

$$\mathbf{x}' = [1, x_1, \dots, x_M]^T$$

Half-spaces:

$$\mathcal{H}^+ = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} > 0 \text{ and } x_1 = 1\}$$

$$\mathcal{H}^- = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} < 0 \text{ and } x_1 = 1\}$$

Using gradient ascent 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

Optimization for Linear Classifiers

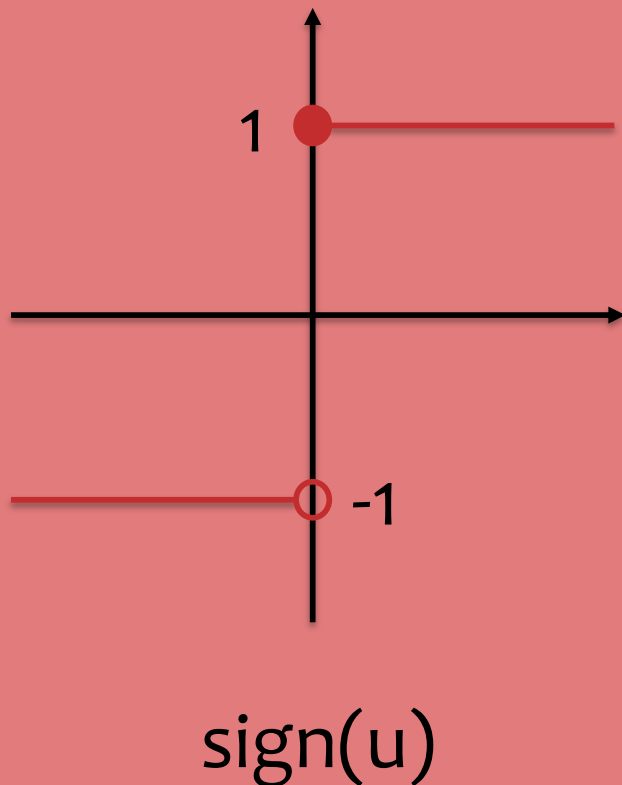
Whiteboard

- Strawman: Mean squared error for Perceptron!
- What does $\theta^T \mathbf{x}$ tell us about \mathbf{x} ?

Using gradient ascent for linear classifiers

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

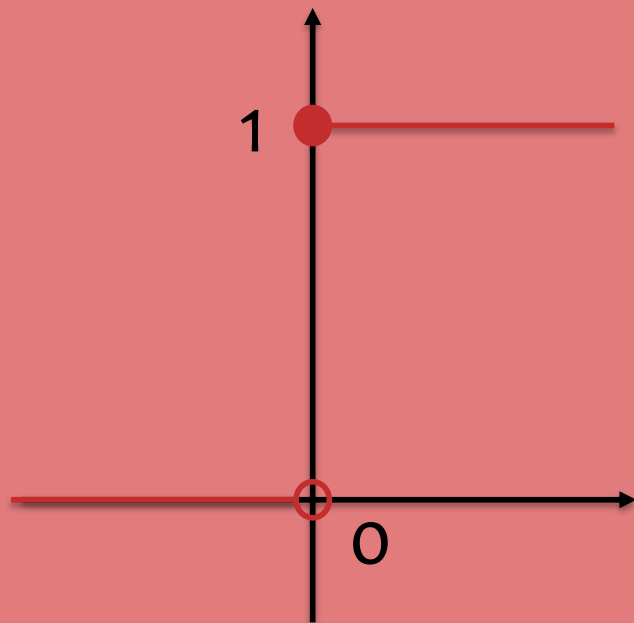
$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$



Using gradient ascent for linear classifiers

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

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



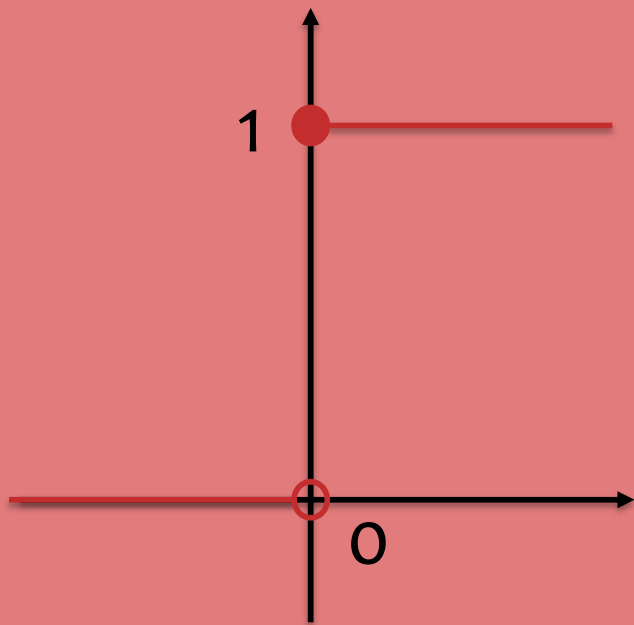
“sign”(u)

Goal: Learn a linear classifier with Gradient Descent

Using gradient ascent for linear classifiers

But this decision function isn't differentiable...

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



“sign”(u)

Use a differentiable function instead!

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



$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

Logistic Regression

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

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

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. $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$

Prediction: Output is the most probable class.

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

Logistic Regression

Whiteboard

- Logistic Regression Model
- Partial derivative for logistic regression
- Gradient for logistic regression
- Decision boundary