



# 10-601 Introduction to Machine Learning

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

## Logistic Regression

Matt Gormley  
Lecture 9  
Sep. 26, 2018

# Reminders

- **Homework 3: KNN, Perceptron, Lin.Reg.**
  - Out: Wed, Sep 19
  - Due: Fri, Sep 28 at 11:59pm
- **Homework 4: Logistic Regression**
  - Out: Fri, Sep 28
  - Due: Mon, Oct 8 at 11:59pm

# **PROBABILISTIC LEARNING**

# Learning from Data (Frequentist)

## *Whiteboard*

- Principle of Maximum Likelihood Estimation (MLE)
- Strawmen:
  - Example: Bernoulli
  - Example: Gaussian
  - Example: Conditional #1  
(Bernoulli conditioned on Gaussian)
  - Example: Conditional #2  
(Gaussians conditioned on Bernoulli)

# **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

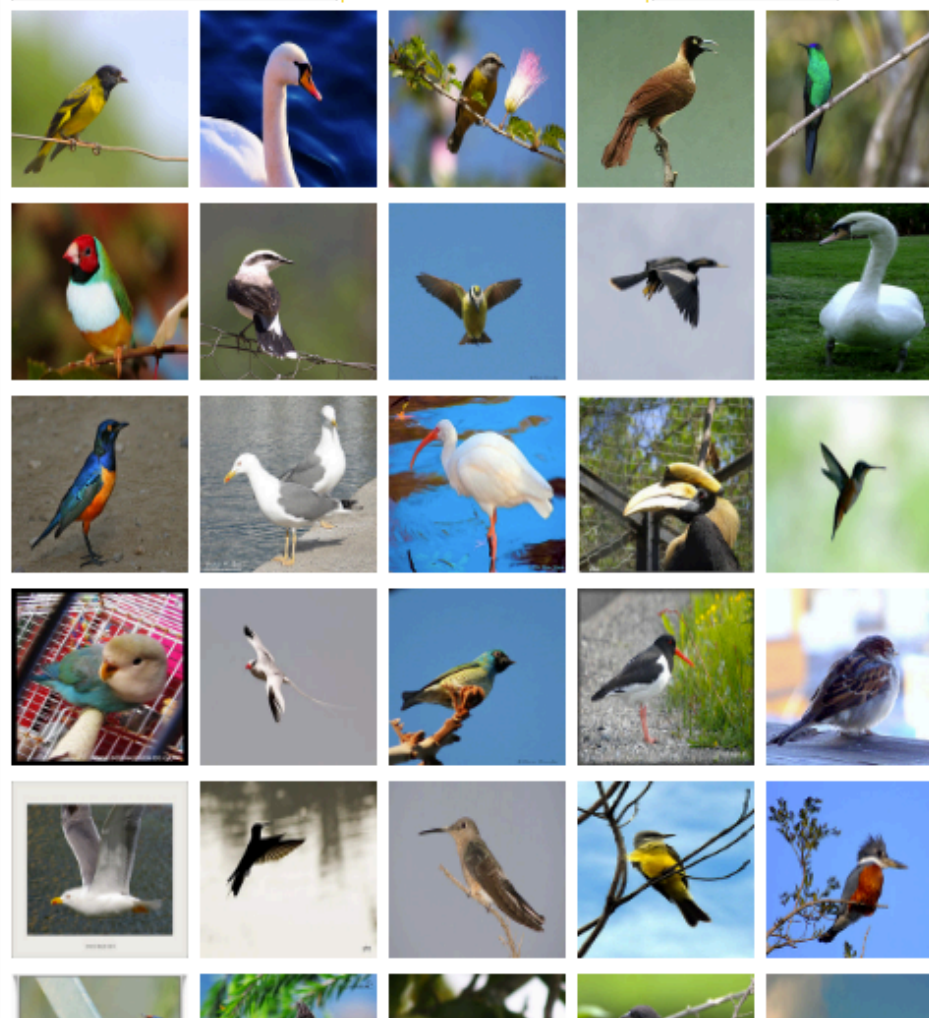


- marine animal, marine creature, sea animal, sea creature (1)
- scavenger (1)
- biped (0)
- predator, predatory animal (1)
- larva (49)
- acrodont (0)
- feeder (0)
- stunt (0)
- chordate (3087)
  - tunicate, urochordate, urochord (6)
  - cephalochochordate (1)
  - vertebrate, craniate (3077)
    - mammal, 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



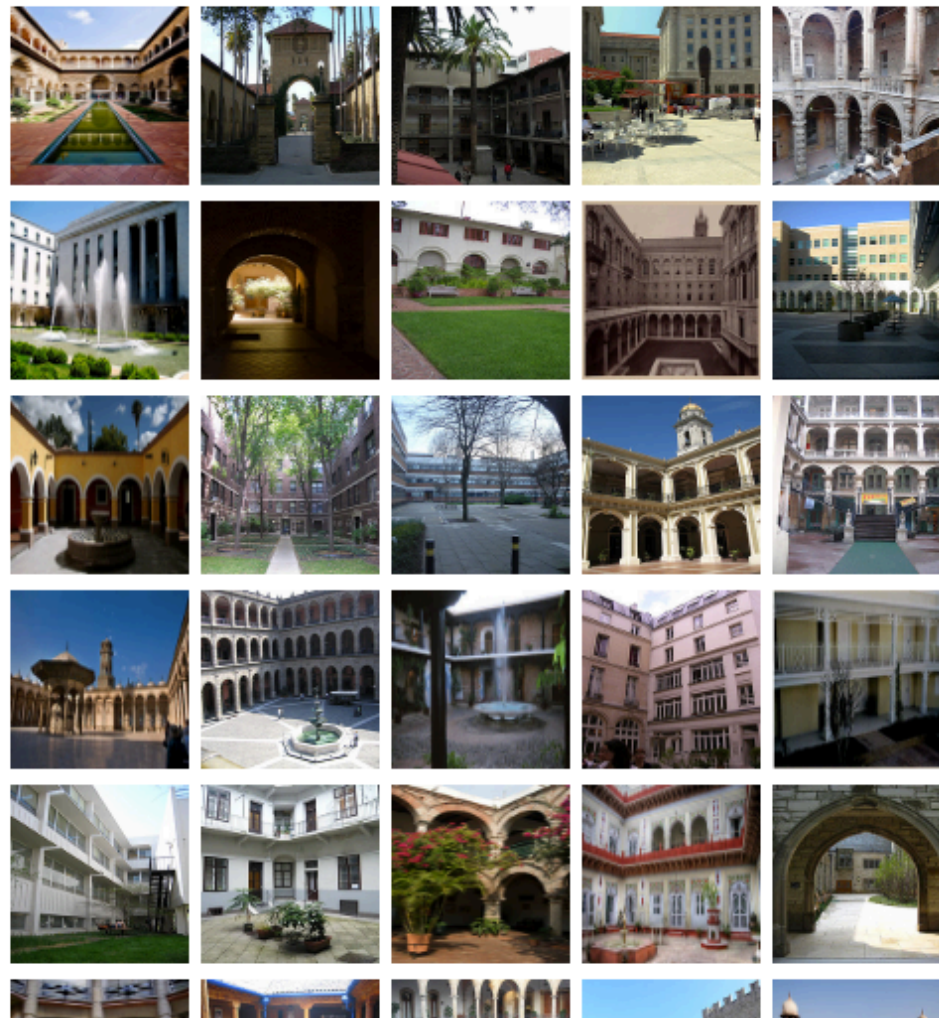
Numbers in brackets: (the number of synsets in the subtree ).

- 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)
          - parvis (0)

## Treemap Visualization

## Images of the Synset

## Downloads



# Example: Image Classification

## CNN for Image Classification

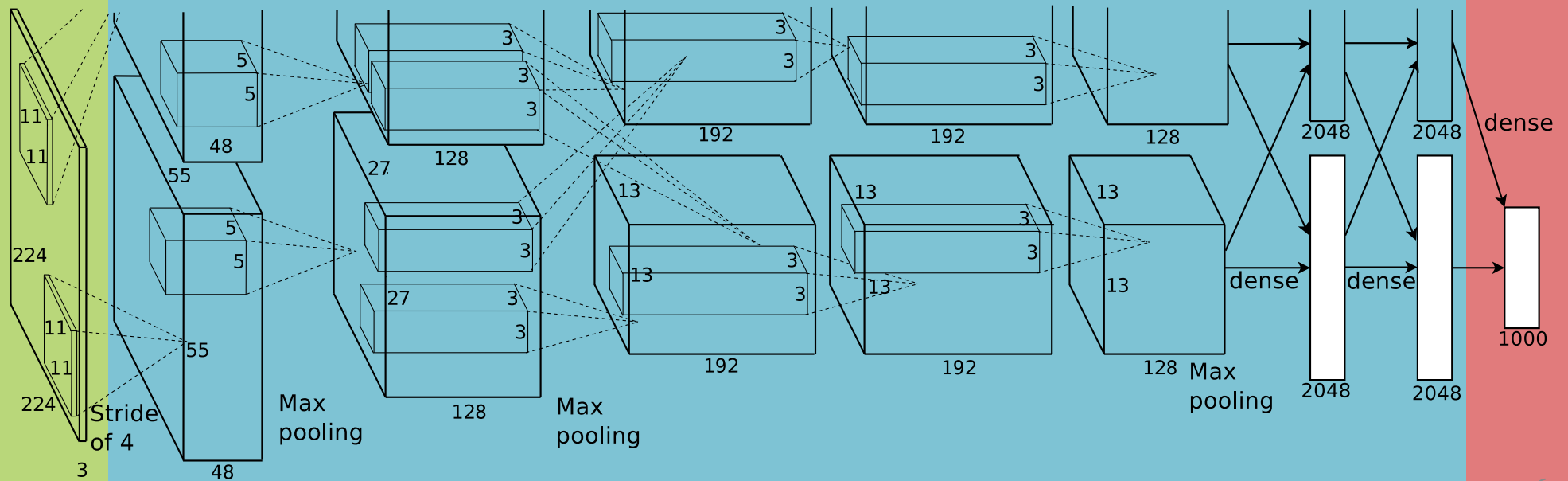
(Krizhevsky, Sutskever & Hinton, 2011)

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

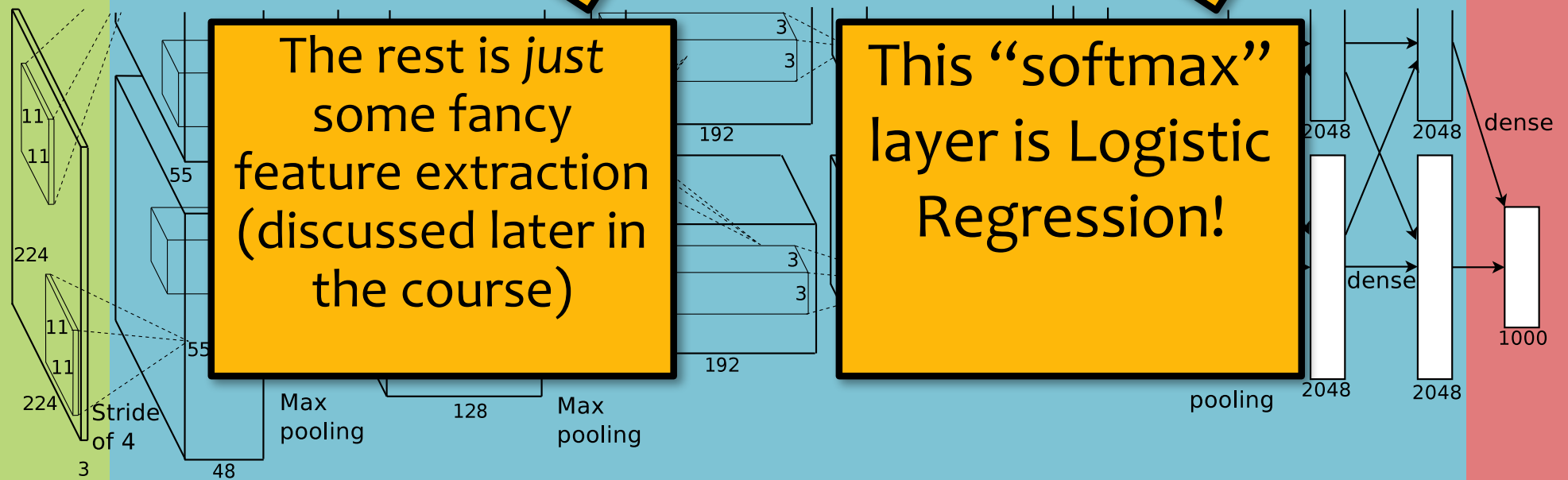
Input  
image  
(pixels)

- 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:

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

3.5

4.5

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!

Hyperplane (Definition 1):

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

Hyperplane (Definition 2):

$$\mathcal{H} = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} = 0$$

$$\text{and } x_0 = 1\}$$

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

Half-spaces:

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

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



# Using gradient ascent for linear classifiers

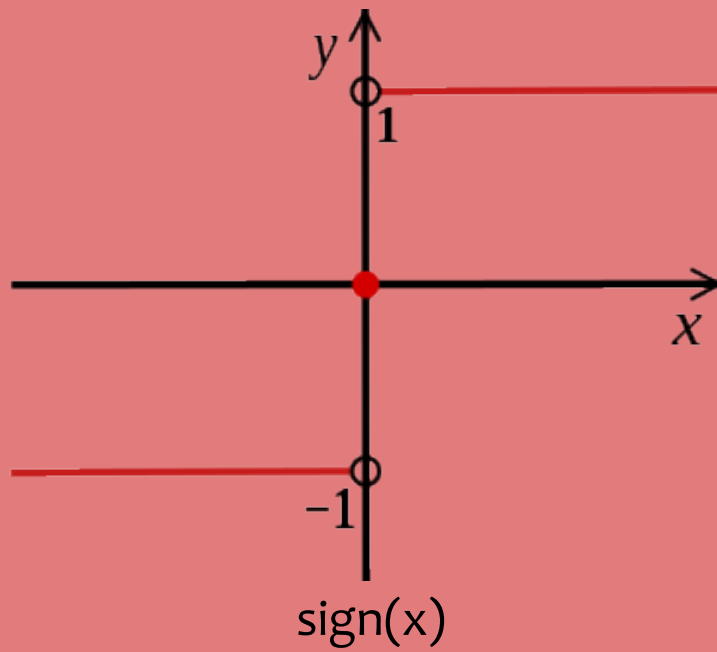
Key idea behind today's lecture:

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

# Using gradient ascent for linear classifiers

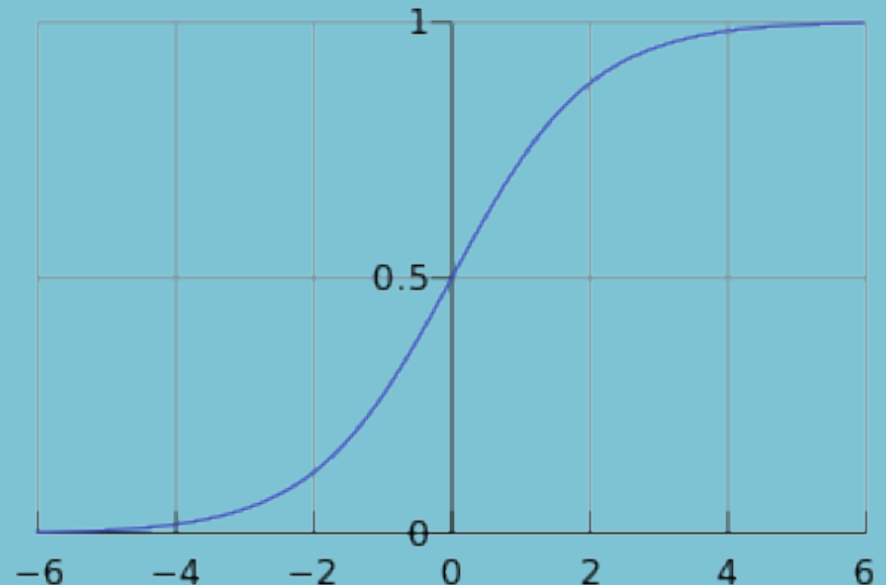
This decision function isn't differentiable:

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



Use a differentiable function instead:

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

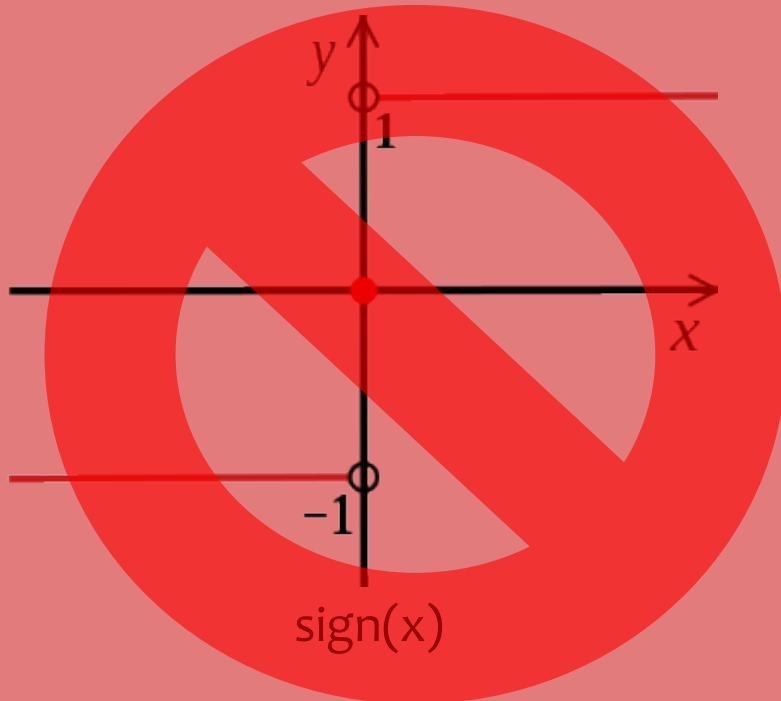


$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

# Using gradient ascent for linear classifiers

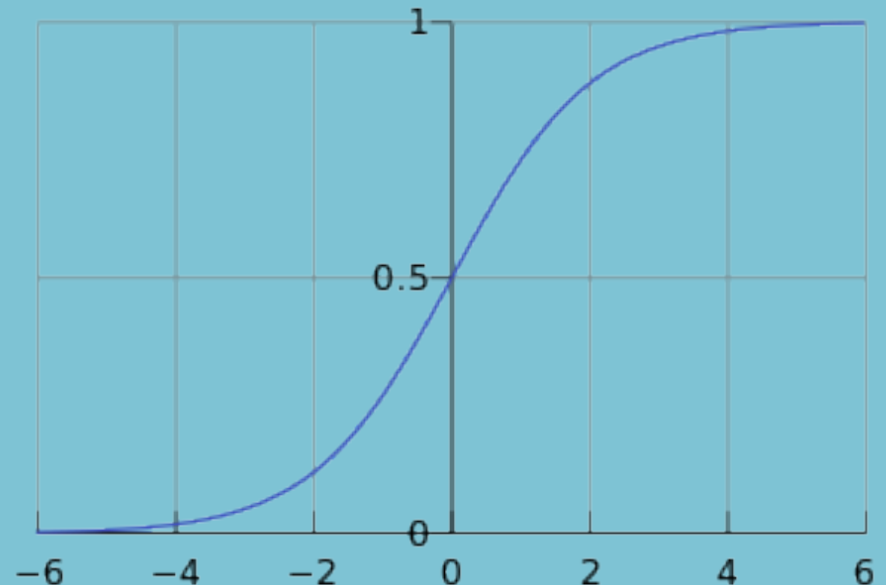
This decision function isn't differentiable:

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



Use a differentiable function instead:

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



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

# 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.

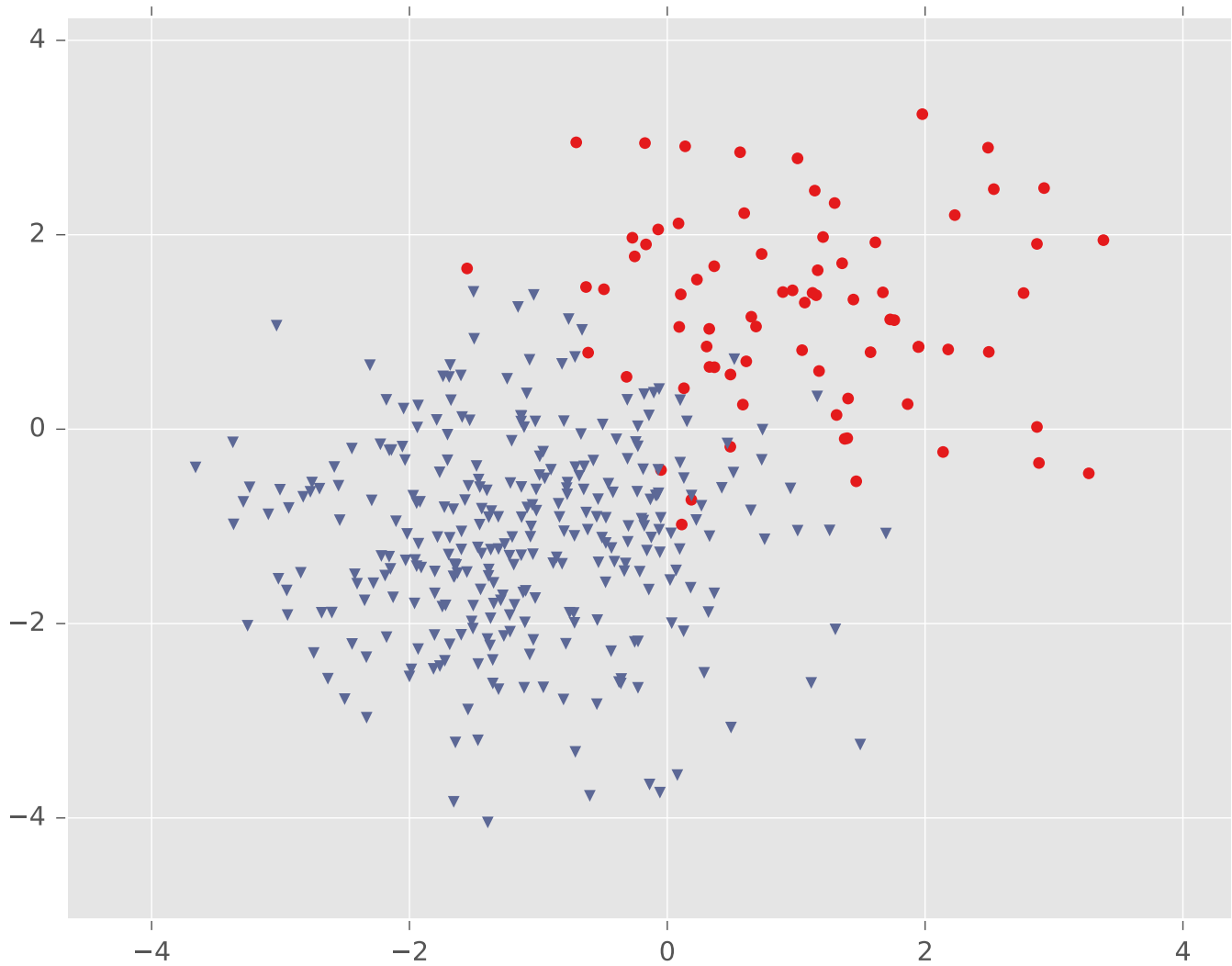
$$\hat{y} = \underset{y \in \{0,1\}}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(y|\mathbf{x})$$

# Logistic Regression

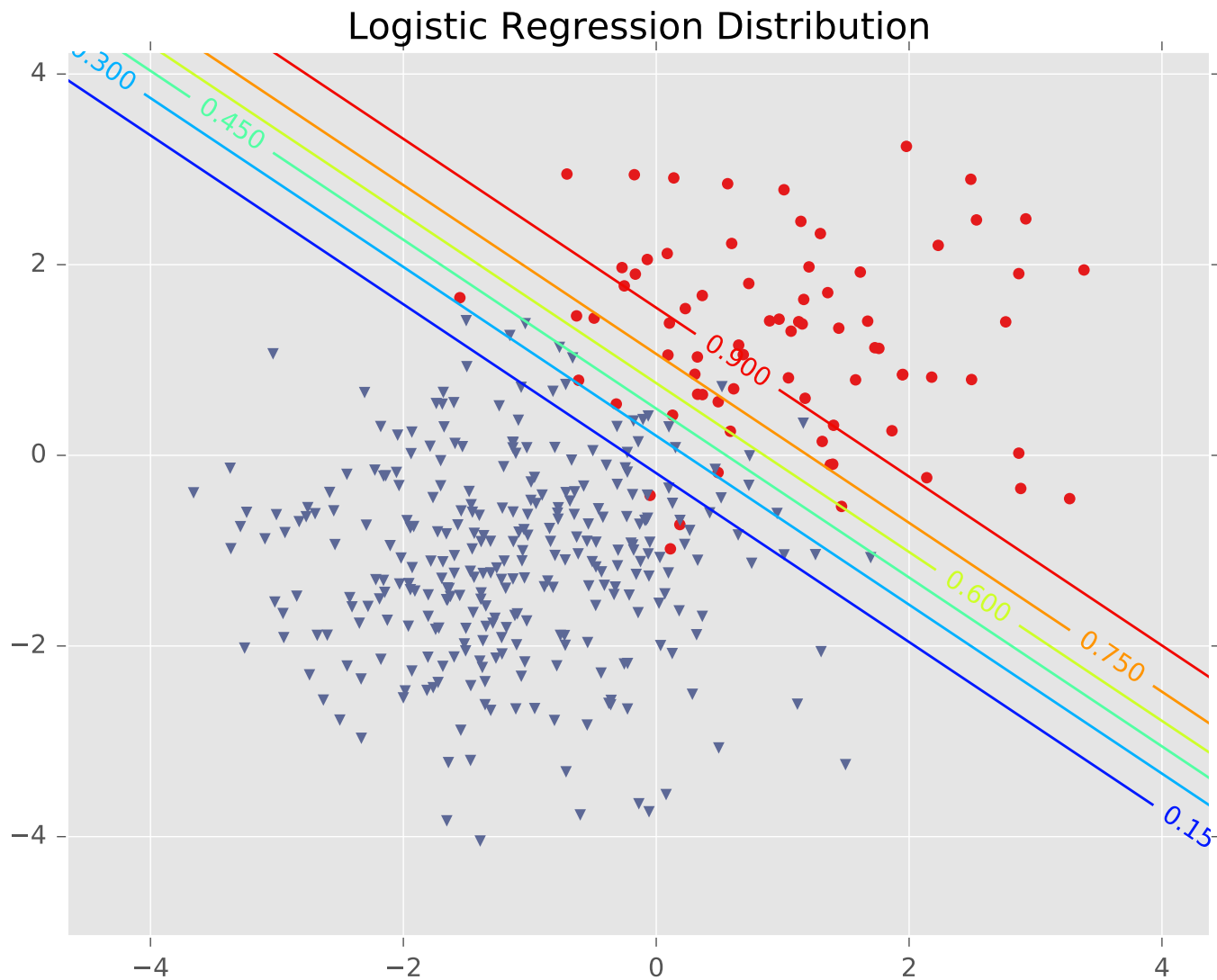
## *Whiteboard*

- Bernoulli interpretation
- Logistic Regression Model
- Decision boundary

# Logistic Regression



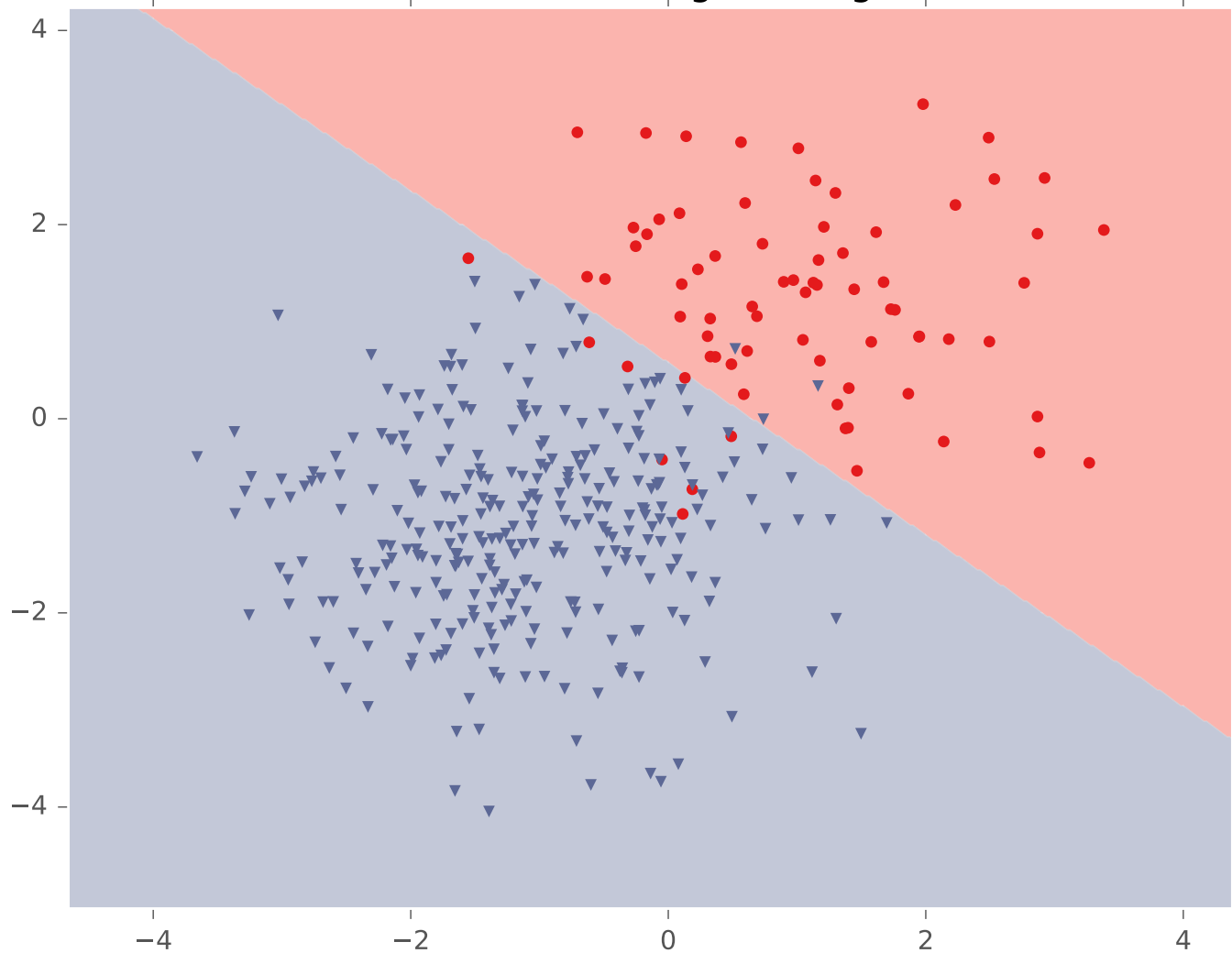
# Logistic Regression





# Logistic Regression

Classification with Logistic Regression



# **LEARNING LOGISTIC REGRESSION**

# Maximum Conditional Likelihood Estimation

**Learning:** finds the parameters that minimize some objective function.

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

We minimize the *negative* log conditional likelihood:

$$J(\theta) = -\log \prod_{i=1}^N p_{\theta}(y^{(i)} | \mathbf{x}^{(i)})$$

Why?

1. We can't maximize likelihood (as in Naïve Bayes) because we don't have a joint model  $p(\mathbf{x}, y)$
2. It worked well for Linear Regression (least squares is MCLE)

# Maximum Conditional Likelihood Estimation

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

**Approach 1:** Gradient Descent

(take larger – more certain – steps opposite the gradient)

**Approach 2:** Stochastic Gradient Descent (SGD)

(take many small steps opposite the gradient)

**Approach 3:** Newton's Method

(use second derivatives to better follow curvature)

**Approach 4:** Closed Form???

(set derivatives equal to zero and solve for parameters)

# Maximum Conditional Likelihood Estimation

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

**Approach 1:** Gradient Descent

(take larger – more certain – steps opposite the gradient)

**Approach 2:** Stochastic Gradient Descent (SGD)

(take many small steps opposite the gradient)

**Approach 3:** Newton's Method

(use second derivatives to better follow curvature)

~~**Approach 4:** Closed Form???~~

~~(set derivatives equal to zero and solve for parameters)~~

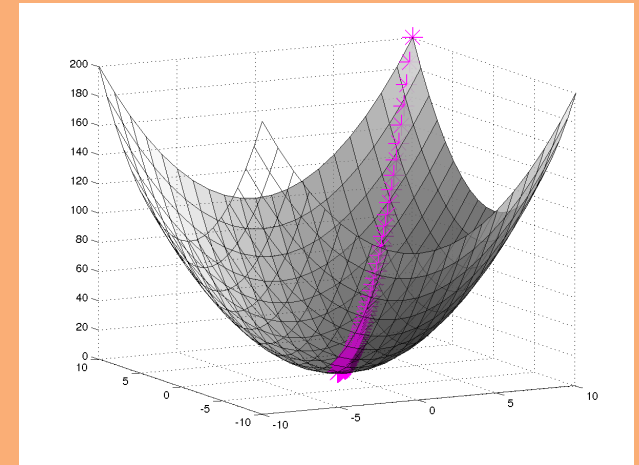
Logistic Regression does not have a closed form solution for MLE parameters.

Recall...

# Gradient Descent

## Algorithm 1 Gradient Descent

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



In order to apply GD to Logistic Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

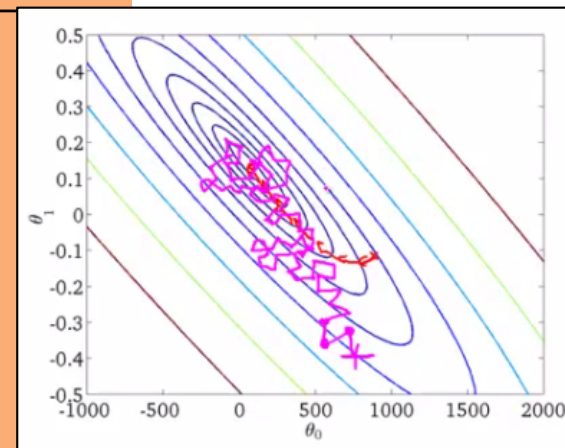
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \frac{d}{d\theta_2} J(\theta) \\ \vdots \\ \frac{d}{d\theta_M} J(\theta) \end{bmatrix}$$

# Stochastic Gradient Descent (SGD)

Recall...

## Algorithm 1 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:        $\theta \leftarrow \theta - \lambda \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



We can also apply SGD to solve the MCLE problem for Logistic Regression.

We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta) \\ \text{where } J^{(i)}(\theta) = -\log p_{\theta}(y^i | \mathbf{x}^i).$$



# **GRADIENT FOR LOGISTIC REGRESSION**

# Learning for Logistic Regression

## *Whiteboard*

- Partial derivative for Logistic Regression
- Gradient for Logistic Regression

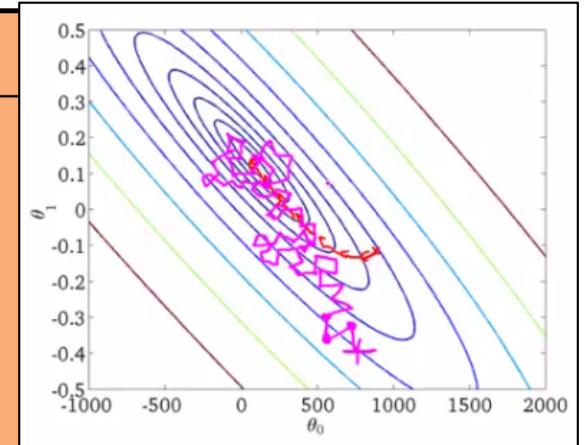
# Details: Picking learning rate

- Use grid-search in log-space over small values on a tuning set:
  - e.g., 0.01, 0.001, ...
- Sometimes, decrease after each pass:
  - e.g factor of  $1/(1 + dt)$ ,  $t=\text{epoch}$
  - sometimes  $1/t^2$
- Fancier techniques I won't talk about:
  - Adaptive gradient: scale gradient differently for each dimension (Adagrad, ADAM, ....)

# SGD for Logistic Regression

## Algorithm 1 SGD for Logistic Regression

```
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:        $\theta \leftarrow \theta - \lambda(y^{(i)} - \rho^{(i)})\mathbf{x}^{(i)}$ 
6:       where  $\rho^{(i)} := 1/(1 + \exp(-\theta^T \mathbf{x}))$ 
7:   return  $\theta$ 
```



We can also apply SGD to solve the MCLE problem for Logistic Regression.

We need a per-example objective:

$$\text{Let } J(\boldsymbol{\theta}) = \sum_{i=1}^N J^{(i)}(\boldsymbol{\theta})$$
$$\text{where } J^{(i)}(\boldsymbol{\theta}) = -\log p_{\boldsymbol{\theta}}(y^i | \mathbf{x}^i).$$

# Summary

1. Discriminative classifiers directly model the **conditional**,  $p(y|x)$
2. Logistic regression is a **simple linear classifier**, that retains a **probabilistic semantics**
3. Parameters in LR are learned by **iterative optimization** (e.g. SGD)

# Logistic Regression Objectives

*You should be able to...*

- Apply the principle of maximum likelihood estimation (MLE) to learn the parameters of a probabilistic model
- Given a discriminative probabilistic model, derive the conditional log-likelihood, its gradient, and the corresponding Bayes Classifier
- Explain the practical reasons why we work with the **log** of the likelihood
- Implement logistic regression for binary or multiclass classification
- Prove that the decision boundary of binary logistic regression is linear
- For linear regression, show that the parameters which minimize squared error are equivalent to those that maximize conditional likelihood

# **MULTINOMIAL LOGISTIC REGRESSION**



# Multinomial Logistic Regression

## *Chalkboard*

- Background: Multinomial distribution
- Definition: Multi-class classification
- Geometric intuitions
- Multinomial logistic regression model
- Generative story
- Reduction to binary logistic regression
- Partial derivatives and gradients
- Applying Gradient Descent and SGD
- Implementation w/ sparse features

# Debug that Program!

## In-Class Exercise: *Think-Pair-Share*

Debug the following program which is (incorrectly) attempting to run SGD for multinomial logistic regression

### Buggy Program:

```
while not converged:
    for i in shuffle([1,...,N]):
        for k in [1,...,K]:
            theta[k] = theta[k] - lambda * grad(x[i], y[i],
theta, k)
```

**Assume:** `grad(x[i], y[i], theta, k)` returns the gradient of the negative log-likelihood of the training example  $(x[i], y[i])$  with respect to vector `theta[k]`. `lambda` is the learning rate. `N` = # of examples. `K` = # of output classes. `M` = # of features. `theta` is a `K` by `M` matrix.

# Debug that Program!

## In-Class Exercise: *Think-Pair-Share*

Debug the following program which is (incorrectly) attempting to run SGD for multinomial logistic regression

### Buggy Program:

```
while not converged:
    for i in shuffle([1,...,N]):
        for k in [1,...,K]:
            for m in [1,..., M]:
                theta[k,m] = theta[k,m] + lambda * grad(x[i],
y[i], theta, k,m)
```

**Assume:** `grad(x[i], y[i], theta, k, m)` returns the partial derivative of the negative log-likelihood of the training example  $(x[i], y[i])$  with respect to  $\theta[k, m]$ . `lambda` is the learning rate. `N` = # of examples. `K` = # of output classes. `M` = # of features. `theta` is a `K` by `M` matrix.