

# Logistic Regression

Professor Ameet Talwalkar

Slide Credit: Professor Fei Sha

# Outline

- 1 Administration
- 2 Review of last lecture
- 3 Logistic regression

# Registration

## Thank you for your patience!

- PTEs given to everyone who submitted HW1
- Other students who want a PTE should contact me directly

# Class Projects

- 1-2 students per project
- Projects can be theoretical, algorithmic and/or applied in nature
  - ▶ Develop new learning algorithm
  - ▶ Theoretically analyze an existing or a new algorithm
  - ▶ Apply learning techniques on some problem of interest
- You are responsible for proposing a project!
- Grading
  - ▶ 30% of total class grade
  - ▶ 25% for proposal, 25% for poster presentation, 50% for final report
  - ▶ Proposal will be short (roughly 1 page, with more details soon), but worth 1/4 of grade because planning ahead is important!
  - ▶ You must briefly meet with one the TAs or myself before the proposal submission deadline

# Class Project Timeline

- Before November 5th: Meet with a TA or me
- November 5th: Project Proposal is due
- December 11th: Poster Session; Project Report due

# Outline

- 1 Administration
- 2 Review of last lecture
  - Naive Bayes
- 3 Logistic regression

# Simple strategy: count the words

Bag-of-words representation  
of documents (and textual data)


$$\begin{pmatrix} \text{free} & 100 \\ \text{money} & 2 \\ \vdots & \vdots \\ \text{account} & 2 \\ \vdots & \vdots \end{pmatrix}$$


Just wanted to send a quick reminder about the guest lecture noon. We meet in RTH 105. It has a PC and LCD projector connection for your laptop if you desire. Maybe we can get to setup the A/V stuff.

Again, if you would be able to make it around 30 minutes great.

Thanks so much for your willingness to do this,  
Mark

$$\begin{pmatrix} \text{free} & 1 \\ \text{money} & 1 \\ \vdots & \vdots \\ \text{account} & 2 \\ \vdots & \vdots \end{pmatrix}$$


# Naive Bayes (in our Spam Email Setting)

**Assume**  $X \in \mathbb{R}^D$  and all  $X_d \in [K]$

$$P(X = x, Y = c) = P(Y = c)P(X = x|Y = c) \quad (1)$$

$$= P(Y = c) \prod_k P(k|Y = c)^{z_k} = \pi_c \prod_k \theta_{ck}^{z_k} \quad (2)$$

where  $z_k$  is the number of times  $k$  in  $x$ .

**Key assumptions made**



# Naive Bayes (in our Spam Email Setting)

**Assume**  $X \in \mathbb{R}^D$  and all  $X_d \in [K]$

$$P(X = x, Y = c) = P(Y = c)P(X = x|Y = c) \quad (1)$$

$$= P(Y = c) \prod_k P(k|Y = c)^{z_k} = \pi_c \prod_k \theta_{ck}^{z_k} \quad (2)$$

where  $z_k$  is the number of times  $k$  in  $x$ .

## Key assumptions made

- Conditional independence:

$$P(X_i, X_j|Y = c) = P(X_i|Y = c)P(X_j|Y = c).$$

- $P(X_i|Y = c)$  depends only the value of  $X_i$ , not  $i$  itself (order of words does not matter in “bag-of-words” representation of documents)

# Learning problem

## Training data

$$\mathcal{D} = \{(\{z_{nk}\}_{k=1}^K, y_n)\}_{n=1}^N$$

## Goal

Learn  $\pi_c, c = 1, 2, \dots, C$ , and  $\theta_{ck}, \forall c \in [C], k \in [K]$  under the constraint

$$\sum_c \pi_c = 1$$

and

$$\sum_k \theta_{ck} = \sum_k P(k|Y=c) = 1$$

as well as those quantities should be nonnegative.

# Likelihood Function

Let  $X_1, \dots, X_N$  be IID with PDF  $f(x|\theta)$  (also written as  $f(x; \theta)$ ). The *likelihood function* is defined by  $L(\theta|x)$  (also written as  $L(\theta; x)$ ),

$$L(\theta|x) = \prod_{i=1}^N f(X_i; \theta).$$

**Notes** The likelihood function is just the joint density of the data, except that we treat it as a function of the parameter  $\theta$ ,  $L : \Theta \rightarrow [0, \infty)$ .

# Maximum Likelihood Estimator

**Definition:** The maximum likelihood estimator (MLE)  $\hat{\theta}$ , is the value of  $\theta$  that maximizes  $L(\theta)$ .

The log-likelihood function is defined by  $l(\theta) = \log L(\theta)$ . Its maximum occurs at the same place as that of the likelihood function.

# Maximum Likelihood Estimator

**Definition:** The maximum likelihood estimator (MLE)  $\hat{\theta}$ , is the value of  $\theta$  that maximizes  $L(\theta)$ .

The log-likelihood function is defined by  $l(\theta) = \log L(\theta)$ . Its maximum occurs at the same place as that of the likelihood function.

- Using logs simplifies mathematical expressions (converts exponents to products and products to sums)
- Using logs helps with numerical stability

The same is true of the likelihood function times any constant. Thus we shall often drop constants in the likelihood function.

# Our hammer: maximum likelihood estimation

## Log-Likelihood of the training data

$$\mathcal{L} = \log P(\mathcal{D}) = \log \prod_{n=1}^N \pi_{y_n} P(x_n | y_n) \quad (3)$$

$$= \log \prod_{n=1}^N \left( \pi_{y_n} \prod_k \theta_{y_n k}^{z_{nk}} \right) \quad (4)$$

$$= \sum_n \left( \log \pi_{y_n} + \sum_k z_{nk} \log \theta_{y_n k} \right) \quad (5)$$

$$= \sum_n \log \pi_{y_n} + \sum_{n,k} z_{nk} \log \theta_{y_n k} \quad (6)$$

# Our hammer: maximum likelihood estimation

## Log-Likelihood of the training data

$$\mathcal{L} = \log P(\mathcal{D}) = \log \prod_{n=1}^N \pi_{y_n} P(x_n | y_n) \quad (3)$$

$$= \log \prod_{n=1}^N \left( \pi_{y_n} \prod_k \theta_{y_n k}^{z_{nk}} \right) \quad (4)$$

$$= \sum_n \left( \log \pi_{y_n} + \sum_k z_{nk} \log \theta_{y_n k} \right) \quad (5)$$

$$= \sum_n \log \pi_{y_n} + \sum_{n,k} z_{nk} \log \theta_{y_n k} \quad (6)$$

## Optimize it!

$$(\pi_c^*, \theta_{ck}^*) = \arg \max \sum_n \log \pi_{y_n} + \sum_{n,k} z_{nk} \log \theta_{y_n k}$$

## Details

### Note the separation of parameters in the likelihood

$$\sum_n \log \pi_{y_n} + \sum_{n,k} z_{nk} \log \theta_{y_n k}$$

which implies that  $\{\pi_c\}$  and  $\{\theta_{ck}\}$  can be estimated separately.

### Reorganize terms

$$\sum_n \log \pi_{y_n} = \sum_c \log \pi_c \times (\text{\#of data points labeled as } c)$$

and

$$\sum_{n,k} z_{nk} \log \theta_{y_n k} = \sum_c \sum_{n:y_n=c} \sum_k z_{nk} \log \theta_{ck} = \sum_c \sum_{n:y_n=c,k} z_{nk} \log \theta_{ck}$$

The later implies  $\{\theta_{ck}, k = 1, 2, \dots, K\}$  and  $\{\theta_{c'k}, k = 1, 2, \dots, K\}$  can be estimated independently.



# Estimating $\{\pi_c\}$

**We want to maximize**

$$\sum_c \log \pi_c \times (\text{\#of data points labeled as } c)$$

## Intuition

- Similar to roll a dice (or flip a coin): each side of the dice shows up with a probability of  $\pi_c$  (total C sides)
- And we have total N trials of rolling this dice

## Solution

# Estimating $\{\pi_c\}$

**We want to maximize**

$$\sum_c \log \pi_c \times (\text{\#of data points labeled as } c)$$

## Intuition

- Similar to roll a dice (or flip a coin): each side of the dice shows up with a probability of  $\pi_c$  (total C sides)
- And we have total N trials of rolling this dice

## Solution

$$\pi_c^* = \frac{\text{\#of data points labeled as } c}{N}$$

# Estimating $\{\theta_{ck}, k = 1, 2, \dots, K\}$

We want to maximize

$$\sum_{n:y_n=c,k} z_{nk} \log \theta_{ck}$$

## Intuition

- Again similar to roll a dice: each side of the dice shows up with a probability of  $\theta_{ck}$  (total K sides)
- And we have total  $\sum_{n:y_n=c,k} z_{nk}$  trials.

## Solution

# Estimating $\{\theta_{ck}, k = 1, 2, \dots, K\}$

We want to maximize

$$\sum_{n:y_n=c,k} z_{nk} \log \theta_{ck}$$

## Intuition

- Again similar to roll a dice: each side of the dice shows up with a probability of  $\theta_{ck}$  (total K sides)
- And we have total  $\sum_{n:y_n=c,k} z_{nk}$  trials.

## Solution

$$\theta_{ck}^* = \frac{\text{\#of times side k shows up in data points labeled as c}}{\text{\#total trials for data points labeled as c}}$$

# Translating back to our problem of detecting spam emails

- Collect a lot of ham and spam emails as training examples
- Estimate the “bias”

$$p(\text{ham}) = \frac{\text{\#of ham emails}}{\text{\#of emails}}, \quad p(\text{spam}) = \frac{\text{\#of spam emails}}{\text{\#of emails}}$$

- Estimate the weights (i.e.,  $p(\text{dollar}|\text{ham})$  etc)

$$p(\text{funny\_word}|\text{ham}) = \frac{\text{\#of funny\_word in ham emails}}{\text{\#of words in ham emails}} \quad (7)$$

$$p(\text{funny\_word}|\text{spam}) = \frac{\text{\#of funny\_word in spam emails}}{\text{\#of words in spam emails}} \quad (8)$$

# Classification rule

**Given an unlabeled data point  $x = \{z_k, k = 1, 2, \dots, K\}$ , label it with**

$$y^* = \arg \max_{c \in [C]} P(y = c | x) \quad (9)$$

$$= \arg \max_{c \in [C]} P(y = c) P(x | y = c) \quad (10)$$

$$= \arg \max_c [\log \pi_c + \sum_k z_k \log \theta_{ck}] \quad (11)$$

## A short derivation of the maximum likelihood estimation

To maximize

$$\sum_{n:y_n=c,k} z_{nk} \log \theta_{ck}$$

We use the Lagrangian multiplier

$$\sum_{n:y_n=c,k} z_{nk} \log \theta_{ck} + \lambda \left( \sum_k \theta_{ck} - 1 \right)$$

Taking derivatives with respect to  $\theta_{ck}$  and then find the stationary point

$$\left( \sum_{n:y_n=c} \frac{z_{nk}}{\theta_{ck}} \right) + \lambda = 0 \rightarrow \theta_{ck} = -\frac{1}{\lambda} \sum_{n:y_n=c} z_{nk}$$

Apply constraint  $\sum_k \theta_{ck} = 1$ , plug in expression above for  $\theta_{ck}$ , solve for  $\lambda$ , and plug back into expression for  $\theta_{ck}$ :

$$\theta_{ck} = \frac{\sum_{n:y_n=c} z_{nk}}{\sum_k \sum_{n:y_n=c} z_{nk}}$$

# Summary

## Things you should know

- The form of the naive Bayes model
  - ▶ write down the joint distribution
  - ▶ explain the conditional independence assumption implied by the model
  - ▶ explain how this model can be used to classify spam vs ham emails
  - ▶ explain how it could be used for categorical variables
- Be able to go through the short derivation for parameter estimation
  - ▶ The model illustrated here is called discrete Naive Bayes
  - ▶ HW2 asks you to apply the same principle to other variants of naive Bayes (Gaussian, Bernoulli)
  - ▶ The derivations are very similar – except there you need to estimate different model parameters



# Moving forward

## Examine the classification rule for naive Bayes

$$y^* = \arg \max_c \log \pi_c + \sum_k z_k \log \theta_{ck}$$

For binary classification, we thus determine the label based on the sign of

$$\log \pi_1 + \sum_k z_k \log \theta_{1k} - \left( \log \pi_2 + \sum_k z_k \log \theta_{2k} \right)$$

This is just a linear function of the features  $\{z_k\}$

# Moving forward

## Examine the classification rule for naive Bayes

$$y^* = \arg \max_c \log \pi_c + \sum_k z_k \log \theta_{ck}$$

For binary classification, we thus determine the label based on the sign of

$$\log \pi_1 + \sum_k z_k \log \theta_{1k} - \left( \log \pi_2 + \sum_k z_k \log \theta_{2k} \right)$$

This is just a linear function of the features  $\{z_k\}$

$$w_0 + \sum_k z_k w_k$$

where we “absorb”  $w_0 = \log \pi_1 - \log \pi_2$  and  $w_k = \log \theta_{1k} - \log \theta_{2k}$ .

# Naive Bayes is a linear classifier

**Fundamentally, what really matters in deciding decision boundary is**

$$w_0 + \sum_k z_k w_k$$

This motivates many new methods, including logistic regression, to be discussed next

# Outline

- 1 Administration
- 2 Review of last lecture
- 3 **Logistic regression**
  - General setup
  - Maximum likelihood estimation
  - Gradient descent
  - Newton's method

# Logistic classification

## Setup for two classes

- Input:  $\mathbf{x} \in \mathbb{R}^D$
- Output:  $y \in \{0, 1\}$
- Training data:  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$
- Model:

$$p(y = 1 | \mathbf{x}; b, \mathbf{w}) = \sigma[g(\mathbf{x})]$$

where

$$g(\mathbf{x}) = b + \sum_d w_d x_d = b + \mathbf{w}^T \mathbf{x}$$

and  $\sigma[\cdot]$  stands for the *sigmoid* function

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

# Why the sigmoid function?

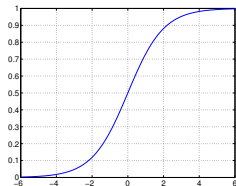
## What does it look like?

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

where

$$a = b + \mathbf{w}^T \mathbf{x}$$

## Properties



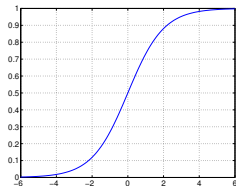
# Why the sigmoid function?

## What does it look like?

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

where

$$a = b + \mathbf{w}^T \mathbf{x}$$



## Properties

- Bounded between 0 and 1 ← thus, interpretable as probability
- Monotonically increasing thus, usable to derive classification rules
  - ▶  $\sigma(a) > 0.5$ , positive (classify as '1')
  - ▶  $\sigma(a) < 0.5$ , negative (classify as '0')
  - ▶  $\sigma(a) = 0.5$ , undecidable
- Nice computational properties As we will see soon

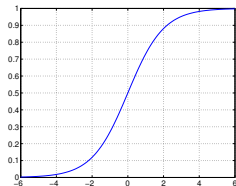
# Why the sigmoid function?

## What does it look like?

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

where

$$a = b + \mathbf{w}^T \mathbf{x}$$



## Properties

- Bounded between 0 and 1 ← thus, interpretable as probability
- Monotonically increasing thus, usable to derive classification rules
  - ▶  $\sigma(a) > 0.5$ , positive (classify as '1')
  - ▶  $\sigma(a) < 0.5$ , negative (classify as '0')
  - ▶  $\sigma(a) = 0.5$ , undecidable
- Nice computational properties As we will see soon

## Linear or nonlinear classifier?



## Linear or nonlinear?

$\sigma(a)$  **is nonlinear**, however, the decision boundary is determined by

$$\sigma(a) = 0.5 \Rightarrow a = 0 \Rightarrow g(\mathbf{x}) = b + \mathbf{w}^T \mathbf{x} = 0$$

which is a *linear* function in  $\mathbf{x}$

We often call  $b$  the offset term.

# Contrast Naive Bayes and our new model

**Similar**

# Contrast Naive Bayes and our new model

## Similar

Both classification models are linear functions of features

## Different

# Contrast Naive Bayes and our new model

## Similar

Both classification models are linear functions of features

## Different

Naive Bayes models the *joint* distribution:  $P(X, Y) = P(Y)P(X|Y)$

Logistic regression models the *conditional* distribution:  $P(Y|X)$

## Generative vs. Discriminative

NB is a *generative* model, LR is a *discriminative* model

- We will talk more about the differences later

# Likelihood function

## Probability of a single training sample $(\mathbf{x}_n, y_n)$

$$p(y_n | \mathbf{x}_n; b; \mathbf{w}) = \begin{cases} \sigma(b + \mathbf{w}^T \mathbf{x}_n) & \text{if } y_n = 1 \\ 1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n) & \text{otherwise} \end{cases}$$

# Likelihood function

## Probability of a single training sample $(\mathbf{x}_n, y_n)$

$$p(y_n | \mathbf{x}_n; b; \mathbf{w}) = \begin{cases} \sigma(b + \mathbf{w}^T \mathbf{x}_n) & \text{if } y_n = 1 \\ 1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n) & \text{otherwise} \end{cases}$$

## Compact expression, exploring that $y_n$ is either 1 or 0

$$p(y_n | \mathbf{x}_n; b; \mathbf{w}) = \sigma(b + \mathbf{w}^T \mathbf{x}_n)^{y_n} [1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]^{1-y_n}$$

# Log Likelihood or Cross Entropy Error

## Log-likelihood of the whole training data $\mathcal{D}$

$$\log P(\mathcal{D}) = \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}$$

# Log Likelihood or Cross Entropy Error

## Log-likelihood of the whole training data $\mathcal{D}$

$$\log P(\mathcal{D}) = \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}$$

It is convenient to work with its negation, which is called *cross-entropy error function*

$$\mathcal{E}(b, \mathbf{w}) = - \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}$$



# Shorthand notation

## This is for convenience

- Append 1 to  $\mathbf{x}$

$$\mathbf{x} \leftarrow [1 \quad x_1 \quad x_2 \quad \cdots \quad x_D]$$

- Append  $b$  to  $\mathbf{w}$

$$\mathbf{w} \leftarrow [b \quad w_1 \quad w_2 \quad \cdots \quad w_D]$$

- Cross-entropy is then

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

# How to find the optimal parameters for logistic regression?

## We will minimize the error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

However, this function is complex and we cannot find the simple solution as we did in Naive Bayes. So we need to use *numerical* methods.

- Numerical methods are messier, in contrast to cleaner analytic solutions.
- In practice, we often have to tune a few optimization parameters — patience is necessary.

# An overview of numerical methods

We describe two

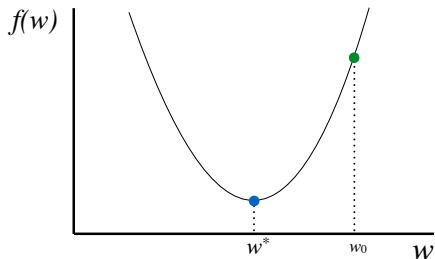
- Gradient descent (our focus in lecture): simple, especially effective for large-scale problems
- Newton's method: classical and powerful method

Gradient descent is often referred to as an *first-order* method as it requires computation of gradients (i.e., the first-order derivative) of the function.

In contrast, Newton's method is often referred to as a *second-order* method (as it requires second derivatives).

# Gradient Descent

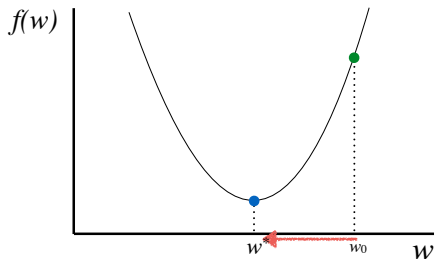
Start at a random point



# Gradient Descent

Start at a random point

Determine a descent direction

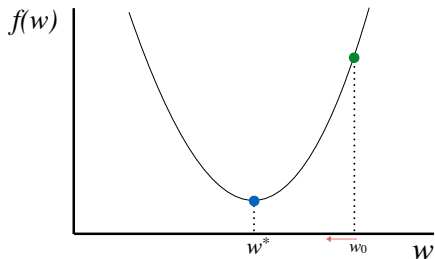


# Gradient Descent

Start at a random point

Determine a descent direction

Choose a step size



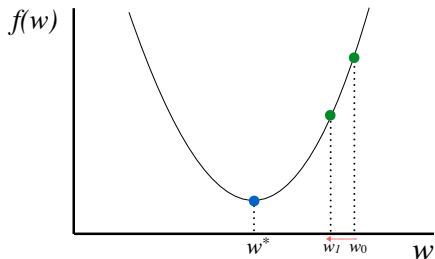
# Gradient Descent

Start at a random point

Determine a descent direction

Choose a step size

Update



# Gradient Descent

Start at a random point

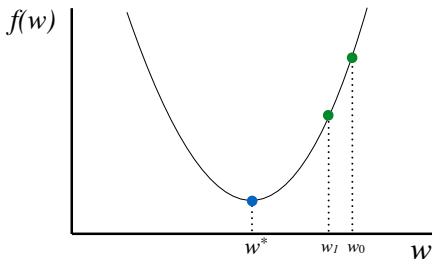
## Repeat

Determine a descent direction

Choose a step size

Update

**Until** stopping criterion is satisfied





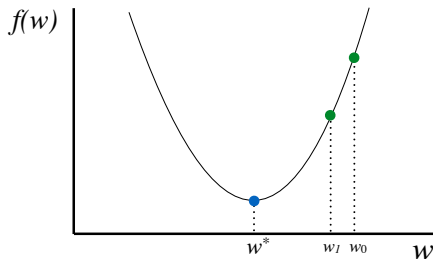
# Gradient Descent

Start at a random point

## Repeat

- Determine a descent direction
- Choose a step size
- Update

**Until** stopping criterion is satisfied



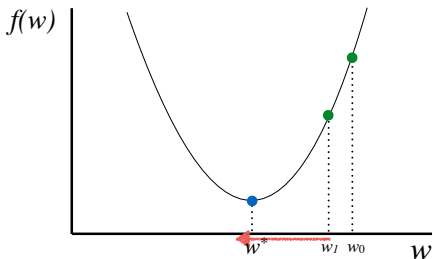
# Gradient Descent

Start at a random point

## Repeat

- Determine a descent direction
- Choose a step size
- Update

Until stopping criterion is satisfied



# Gradient Descent

Start at a random point

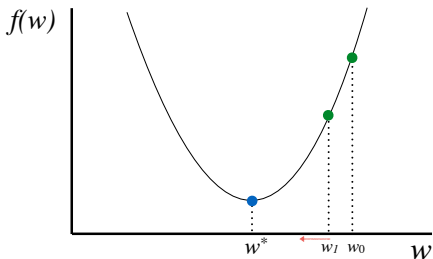
## Repeat

Determine a descent direction

Choose a step size

Update

Until stopping criterion is satisfied



# Gradient Descent

Start at a random point

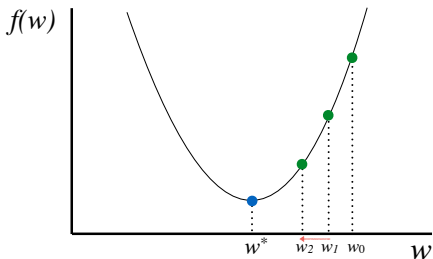
## Repeat

Determine a descent direction

Choose a step size

Update

Until stopping criterion is satisfied



# Gradient Descent

Start at a random point

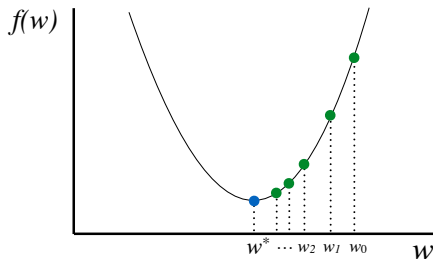
## Repeat

Determine a descent direction

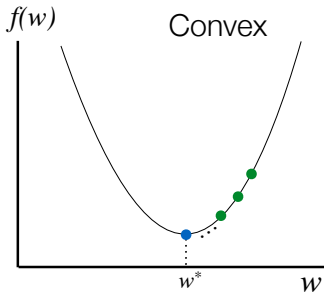
Choose a step size

Update

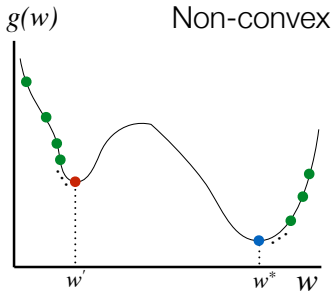
**Until** stopping criterion is satisfied



# Where Will We Converge?



Any local minimum is a global minimum

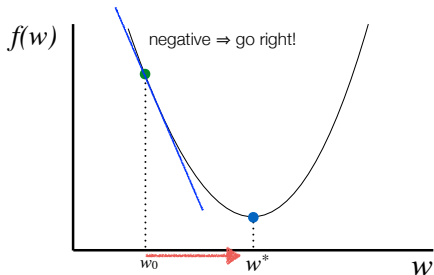
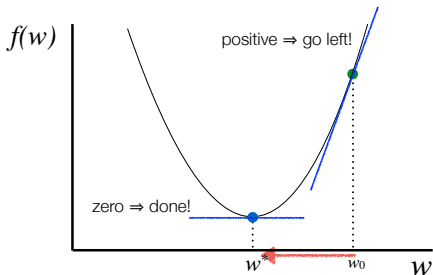


Multiple local minima may exist

**Least Squares, Ridge Regression and Logistic Regression are all convex!**

# Why do we move in the direction opposite the gradient?

# Choosing Descent Direction (1D)



We can only move in two directions  
Negative slope is direction of descent!

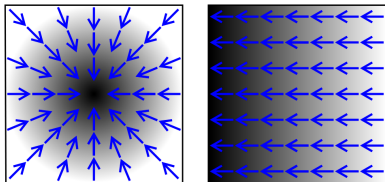
Step Size

**Update Rule:**  $w_{i+1} = w_i - \alpha_i \frac{df}{dw}(w_i)$

Negative Slope



# Choosing Descent Direction



"Gradient2" by Sarang. Licensed under CC BY-SA 2.5 via Wikimedia Commons  
<http://commons.wikimedia.org/wiki/File:Gradient2.svg#/media/File:Gradient2.svg>

2D Example:

- Function values are in black/white and black represents higher values
- Arrows are gradients

We can move anywhere in  $\mathbb{R}^d$   
Negative gradient is direction of *steepest* descent!

Step Size

**Update Rule:**  $\mathbf{w}_{i+1} = \mathbf{w}_i - \alpha_i \nabla f(\mathbf{w}_i)$

Negative Slope

Example:  $\min f(\boldsymbol{\theta}) = 0.5(\theta_1^2 - \theta_2)^2 + 0.5(\theta_1 - 1)^2$

- We compute the gradients

$$\frac{\partial f}{\partial \theta_1} = 2(\theta_1^2 - \theta_2)\theta_1 + \theta_1 - 1 \quad (12)$$

$$\frac{\partial f}{\partial \theta_2} = -(\theta_1^2 - \theta_2) \quad (13)$$

- Use the following *iterative* procedure for *gradient descent*

- 1 Initialize  $\theta_1^{(0)}$  and  $\theta_2^{(0)}$ , and  $t = 0$
- 2 do

$$\theta_1^{(t+1)} \leftarrow \theta_1^{(t)} - \eta \left[ 2(\theta_1^{(t)})^2 - \theta_2^{(t)}\theta_1^{(t)} + \theta_1^{(t)} - 1 \right] \quad (14)$$

$$\theta_2^{(t+1)} \leftarrow \theta_2^{(t)} - \eta \left[ -(\theta_1^{(t)})^2 - \theta_2^{(t)} \right] \quad (15)$$

$$t \leftarrow t + 1 \quad (16)$$

- 3 until  $f(\boldsymbol{\theta}^{(t)})$  *does not change much*

# Gradient descent

## General form for minimizing $f(\theta)$

$$\theta^{t+1} \leftarrow \theta - \eta \frac{\partial f}{\partial \theta}$$

## Remarks

- $\eta$  is often called *step size* – literally, how far our update will go along the the direction of the negative gradient
- Note that this is for *minimizing* a function, hence the subtraction ( $-\eta$ )
- With a *suitable* choice of  $\eta$ , the iterative procedure converges to a stationary point where

$$\frac{\partial f}{\partial \theta} = 0$$

- A stationary point is only necessary for being the minimum (though we're happy when the function is convex)

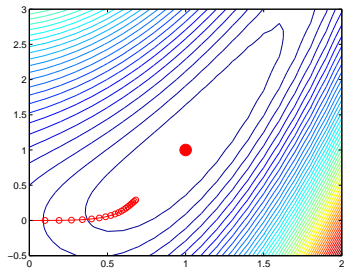
# Seeing in action

**Choosing the right  $\eta$  is important**

# Seeing in action

## Choosing the right $\eta$ is important

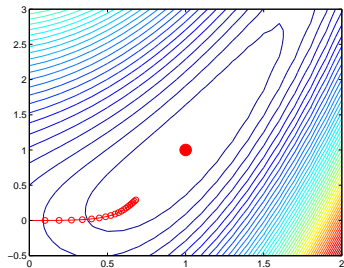
small  $\eta$  is too slow?



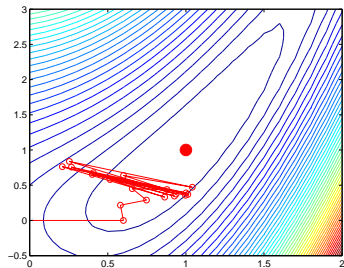
# Seeing in action

## Choosing the right $\eta$ is important

small  $\eta$  is too slow?



large  $\eta$  is too unstable?



# Gradient Descent Update for Logistic Regression

**Simple fact: derivatives of  $\sigma(a)$**

$$\begin{aligned}\frac{d\sigma(a)}{da} &= \frac{d}{da} (1 + e^{-a})^{-1} = \frac{-(1 + e^{-a})'}{(1 + e^{-a})^2} \\ &= \frac{e^{-a}}{(1 + e^{-a})^2} = \frac{1}{1 + e^{-a}} \frac{e^{-a}}{1 + e^{-a}} \\ &= \sigma(a)[1 - \sigma(a)]\end{aligned}$$

# Gradients of the cross-entropy error function

## Cross-entropy Error Function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

## Gradients

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = - \sum_n \{y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\mathbf{w}^T \mathbf{x}_n) \mathbf{x}_n\} \quad (17)$$

$$= \sum_n \{\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (18)$$

## Remark



# Gradients of the cross-entropy error function

## Cross-entropy Error Function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

## Gradients

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = - \sum_n \{y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\mathbf{w}^T \mathbf{x}_n) \mathbf{x}_n\} \quad (17)$$

$$= \sum_n \{\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (18)$$

## Remark

- $e_n = \{\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n\}$  is called *error* for the  $n$ th training sample.

# Numerical optimization

## Gradient descent

- Choose a proper step size  $\eta > 0$
- Iteratively update the parameters following the negative gradient to minimize the error function

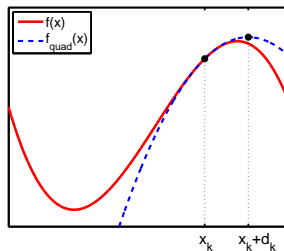
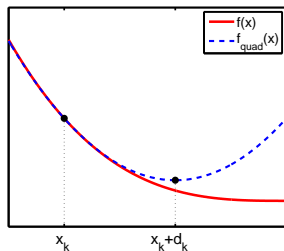
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \sum_n \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n$$

## Remarks

- The step size needs to be chosen carefully to ensure convergence.
- The step size can be adaptive (i.e. varying from iteration to iteration). For example, we can use techniques such as *line search*
- There is a variant called *stochastic* gradient descent, also popularly used (later in this semester).

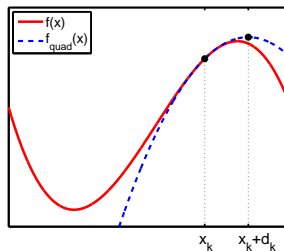
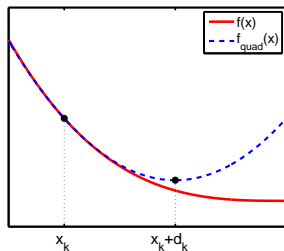
# Intuition for Newton's method

Approximate the true function with an easy-to-solve optimization problem



# Intuition for Newton's method

Approximate the true function with an easy-to-solve optimization problem



In particular, we can approximate the cross-entropy error function around  $w^{(t)}$  by a quadratic function, and then minimize this quadratic function

# Approximation

## Second Order Taylor expansion around $x_t$

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

# Approximation

## Second Order Taylor expansion around $x_t$

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

## Taylor expansion of cross-entropy error function around $\mathbf{w}^{(t)}$

$$\mathcal{E}(\mathbf{w}) \approx \mathcal{E}(\mathbf{w}^{(t)}) + (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)}(\mathbf{w} - \mathbf{w}^{(t)})$$

where

- $\nabla \mathcal{E}(\mathbf{w}^{(t)})$  is the gradient
- $\mathbf{H}^{(t)}$  is the Hessian matrix evaluated at  $\mathbf{w}^{(t)}$

# So what is the Hessian matrix?

## The matrix of second-order derivatives

$$\mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^T}$$

In other words,

$$H_{ij} = \frac{\partial}{\partial w_j} \left( \frac{\partial \mathcal{E}(\mathbf{w})}{\partial w_i} \right)$$

So the Hessian matrix is  $\mathbb{R}^{D \times D}$ , where  $\mathbf{w} \in \mathbb{R}^D$ .

# Optimizing the approximation

## Minimize the approximation

$$\mathcal{E}(\mathbf{w}) \approx \mathcal{E}(\mathbf{w}^{(t)}) + (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)} (\mathbf{w} - \mathbf{w}^{(t)})$$

and use the solution as the new estimate of the parameters

$$\mathbf{w}^{(t+1)} \leftarrow \min_{\mathbf{w}} (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)} (\mathbf{w} - \mathbf{w}^{(t)})$$



# Optimizing the approximation

## Minimize the approximation

$$\mathcal{E}(\mathbf{w}) \approx \mathcal{E}(\mathbf{w}^{(t)}) + (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)} (\mathbf{w} - \mathbf{w}^{(t)})$$

and use the solution as the new estimate of the parameters

$$\mathbf{w}^{(t+1)} \leftarrow \min_{\mathbf{w}} (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)} (\mathbf{w} - \mathbf{w}^{(t)})$$

The quadratic function minimization has a *closed* form, thus, we have

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \left( \mathbf{H}^{(t)} \right)^{-1} \nabla \mathcal{E}(\mathbf{w}^{(t)})$$

i.e., the Newton's method.

# Contrast gradient descent and Newton's method

## Similar

- Both are iterative procedures.

## Different

- Newton's method requires second-order derivatives.
- Newton's method does not have the magic  $\eta$  to be set.

## Other important things about Hessian

**Our cross-entropy error function is convex**

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_n \{\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (19)$$

$$\Rightarrow \mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^T} = \text{homework} \quad (20)$$

## Other important things about Hessian

### Our cross-entropy error function is convex

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_n \{\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (19)$$

$$\Rightarrow \mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^T} = \text{homework} \quad (20)$$

For any vector  $\mathbf{v}$ ,

$$\mathbf{v}^T \mathbf{H} \mathbf{v} = \text{homework} \geq 0$$

Thus, positive semi-definite. Thus, the cross-entropy error function is convex, with only one global optimum.

# Good about Newton's method

## Fast (in terms of convergence)!

Newton's method finds the optimal point in a *single* iteration when the function we're optimizing is quadratic

In general, the better our Taylor approximation, the more quickly Newton's method will converge

# Bad about Newton's method

## Not scalable!

Computing and inverting Hessian matrix can be very expensive for large-scale problems where the dimensionally  $D$  is very large. *There are fixes and alternatives, such as Quasi-Newton/Quasi-second order method.*

# Summary

## Setup for 2 classes

- Logistic Regression models conditional distribution as:  
 $p(y = 1|\mathbf{x}; \mathbf{w}) = \sigma[g(\mathbf{x})]$  where  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Linear decision boundary:  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = 0$

## Minimizing cross-entropy error (negative log-likelihood)

- $\mathcal{E}(b, \mathbf{w}) = - \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}$
- No closed form solution; must rely on iterative solvers

## Numerical optimization

- Gradient descent: simple, scalable to large-scale problems
  - ▶ move in direction opposite of gradient!
  - ▶ gradient of logistic function takes nice form
- Newton method: fast to converge but not scalable
  - ▶ At each iteration, find optimal point in 2nd-order Taylor expansion
  - ▶ Closed form solution exists for each iteration

# Naive Bayes and logistic regression: two different modeling paradigms

- Maximize *joint* likelihood  $\sum_n \log p(\mathbf{x}_n, y_n)$  (Generative, NB)
- Maximize *conditional* likelihood  $\sum_n \log p(y_n | \mathbf{x}_n)$  (Discriminative, LR)
- More on this next class