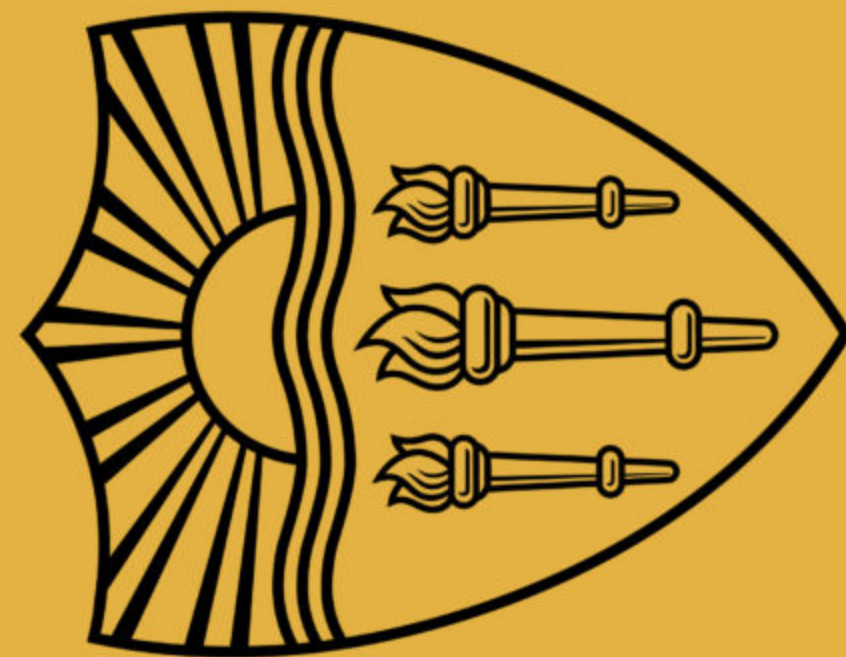# Lecture 5:
# Logistic Regression Continued

*Instructor: Swabha Swayamdipta*
*USC CSCI 499 LMs in NLP*
*Jan 31, Spring 2024*

# Logistics  / Announcements

- HW1 due tonight by 12pm
- Graded Quiz 1 will be distributed next Monday
- Project Proposal (1 page) due next Wednesday

- What makes for a good project proposal?
  - Problem / Task Definition: Input and expected output
    - With examples
    - Motivation: Why is this an important problem?
  - Sources of data for training a model for your task (very important!)
    - How many instances for train / test / validation?
  - Models to be compared
    - Ideally, those taught in class
  - Evaluation Framework:
    - What test set?
    - What metric?
    - What baselines for comparison?

# Lecture Outline

- Recap: Supervised Classification / Logistic Regression
  - Data (Features)
  - Model (Logistic Regression)
  - Loss (Cross Entropy)
- Today: Logistic Regression Contd.
  - Optimization
  - Regularization
  - Multinomial Logistic Regression
- Word Embeddings (Teaser)

# Recap:
# Supervised Machine Learning

# Ingredients of Supervised Machine Learning

I. **Data** as pairs $(x^{(i)}, y^{(i)})$ s.t $i \in \{1...N\}$
- $x^{(i)}$ usually represented by a feature vector $\mathbf{x}^{(i)} = [x_1, x_2, ..., x_d]$,
- e.g. word embeddings

II. **Model**
- A classification function that computes $\hat{y}$, the estimated class, via $p(y|x)$
- e.g. sigmoid function: $\sigma(z) = 1/(1 + \exp(-z))$

III. **Loss**
- An objective function for learning
- e.g. cross-entropy loss, $L_{CE}$

IV. **Optimization**
- An algorithm for optimizing the objective function
- e.g. stochastic gradient descent

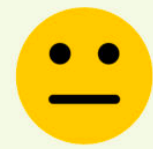V. **Inference** / Evaluation

**Learning Phase**

5

# Text Classification Tasks



SENTIMENT ANALYSIS

**POSITIVE**
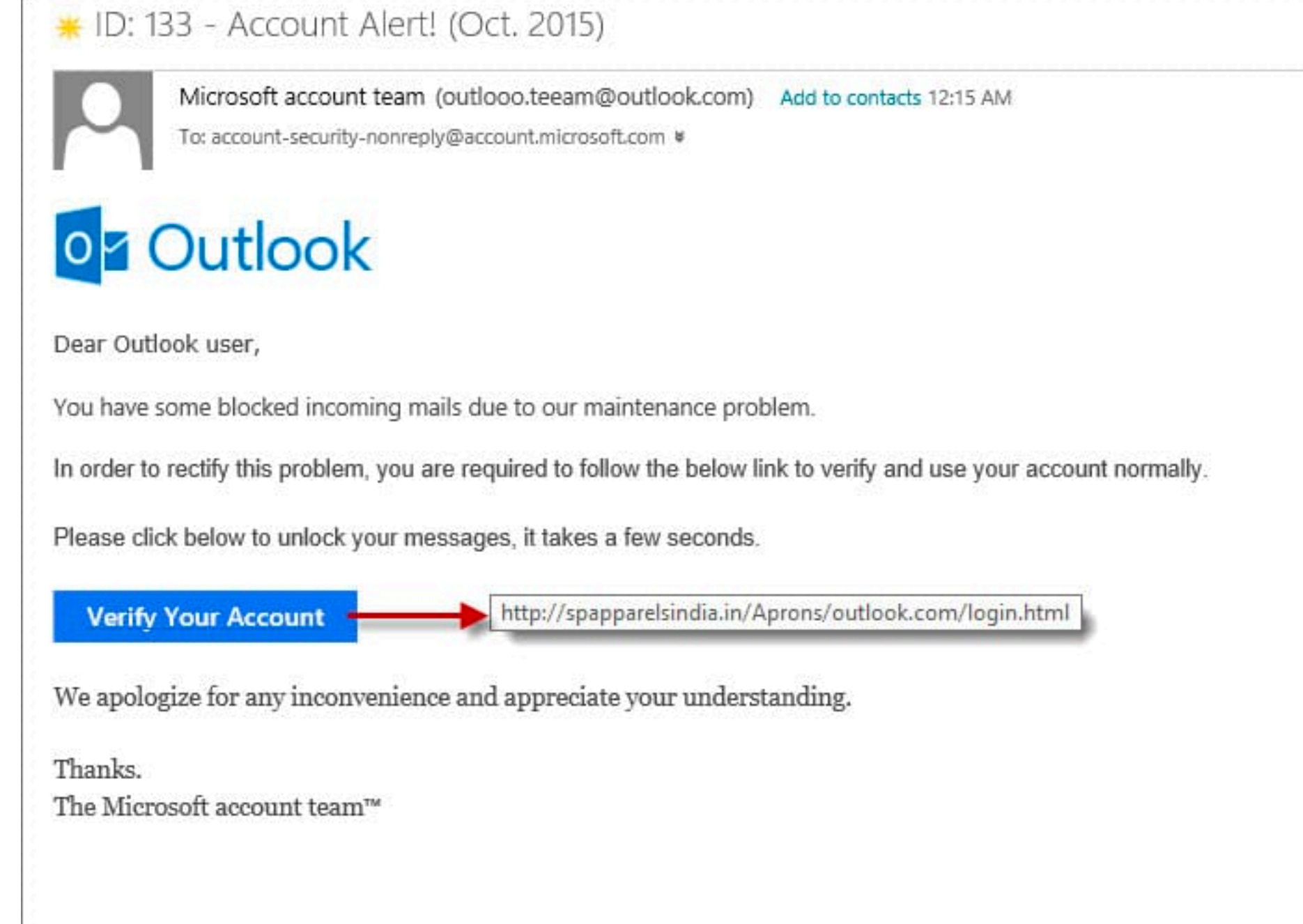"Great service for an affordable price.
We will definitely be booking again."

**NEUTRAL**
"Just booked two nights at this hotel."

**NEGATIVE**
"Horrible services. The room was dirty and unpleasant. Not worth the money."



Mr. WILLIAM SHAKESPEARES
COMEDIES, HISTORIES, & TRAGEDIES.

FEDERALIST: A COLLECTION

Signature ESSAYS,

WRITTEN IN FAVOUR OF THE



☀ ID: 133 - Account Alert! (Oct. 2015)

Microsoft account team (outlooo.teeam@outlook.com)   Add to contacts 12:15 AM
To: account-security-nonreply@account.microsoft.com

**Outlook**

Dear Outlook user,

You have some blocked incoming mails due to our maintenance problem.

In order to rectify this problem, you are required to follow the below link to verify and use your account normally.

Please click below to unlock your messages, it takes a few seconds.

**Verify Your Account** → http://spapparelsindia.in/Aprons/outlook.com/login.html

We apologize for any inconvenience and appreciate your understanding.

Thanks.
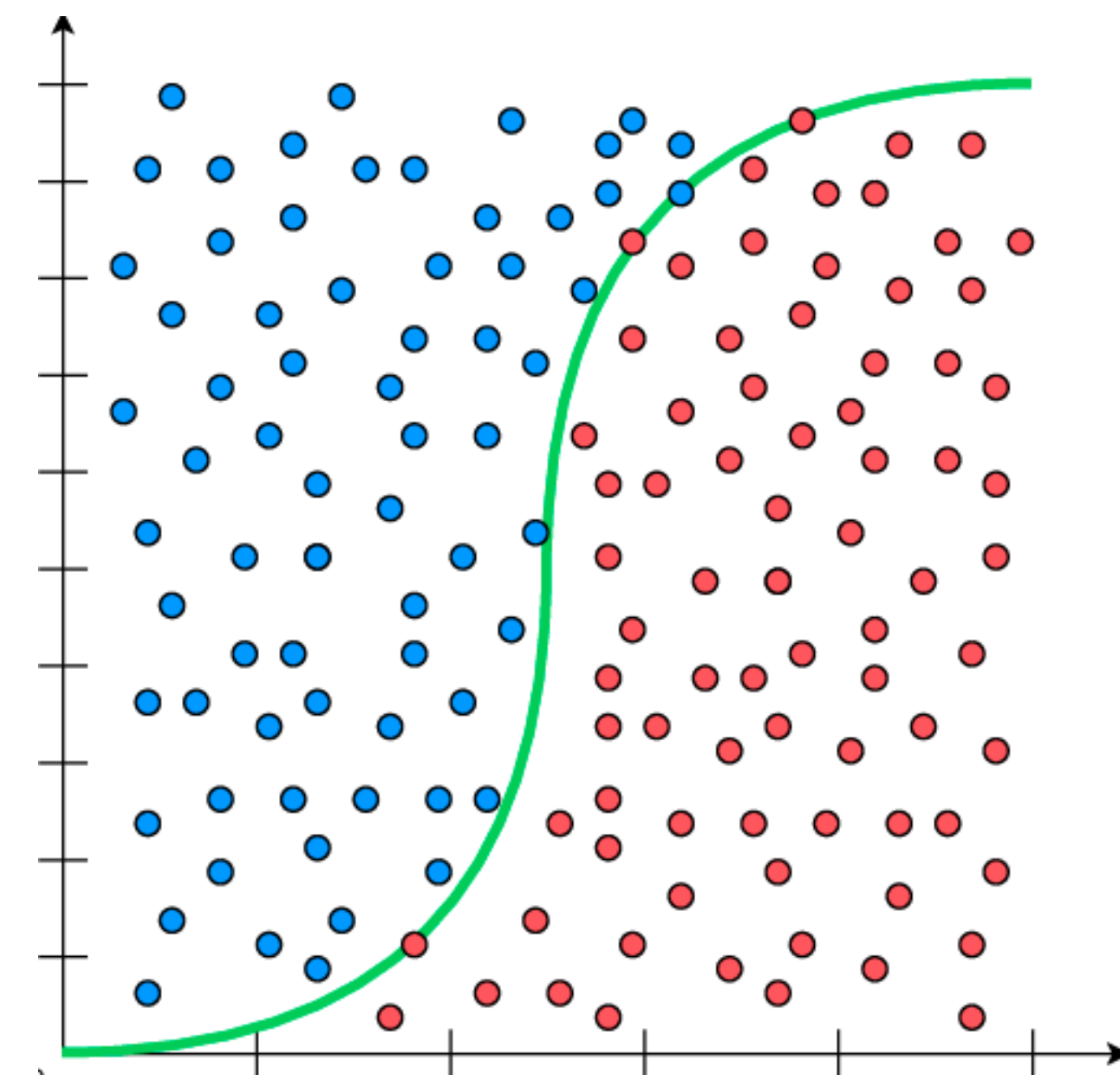The Microsoft account team™

Not just NLP, classification is a general ML technique often applied across a wide variety of prediction tasks!

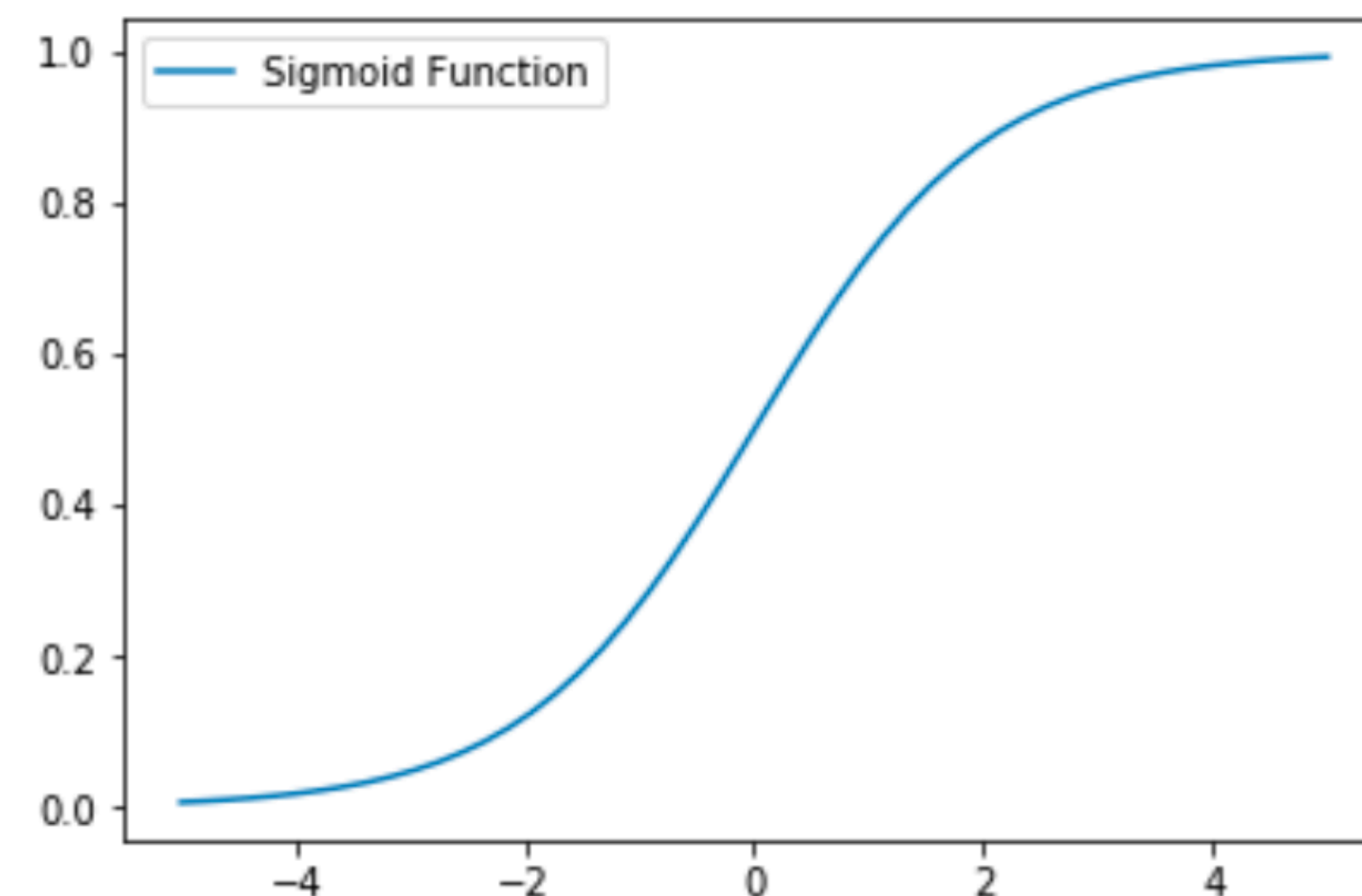# Text Classification Setup

- Input:
  - a document $x$
    - Each observation $x^{(i)}$ is represented by a feature vector
      $$\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}, \ldots, x_d^{(i)}]$$
  - a label $y$ from a fixed set of classes $C = c_1, c_2, \ldots, c_J$
    - Sometimes denoted by $y*$ for gold-standard (ground truth) label
- Output: a predicted class $\hat{y} \in C$
- Setting for Binary Classification: given a series of input / output pairs:
  - $(x^{(i)}, y^{(i)})$ where label $y^{(i)} \in C = \{0,1\}$
- Goal of Binary Classification
  - At test time, for input $x^{test}$, compute an output: a predicted class $\hat{y}^{test} \in \{0,1\}$
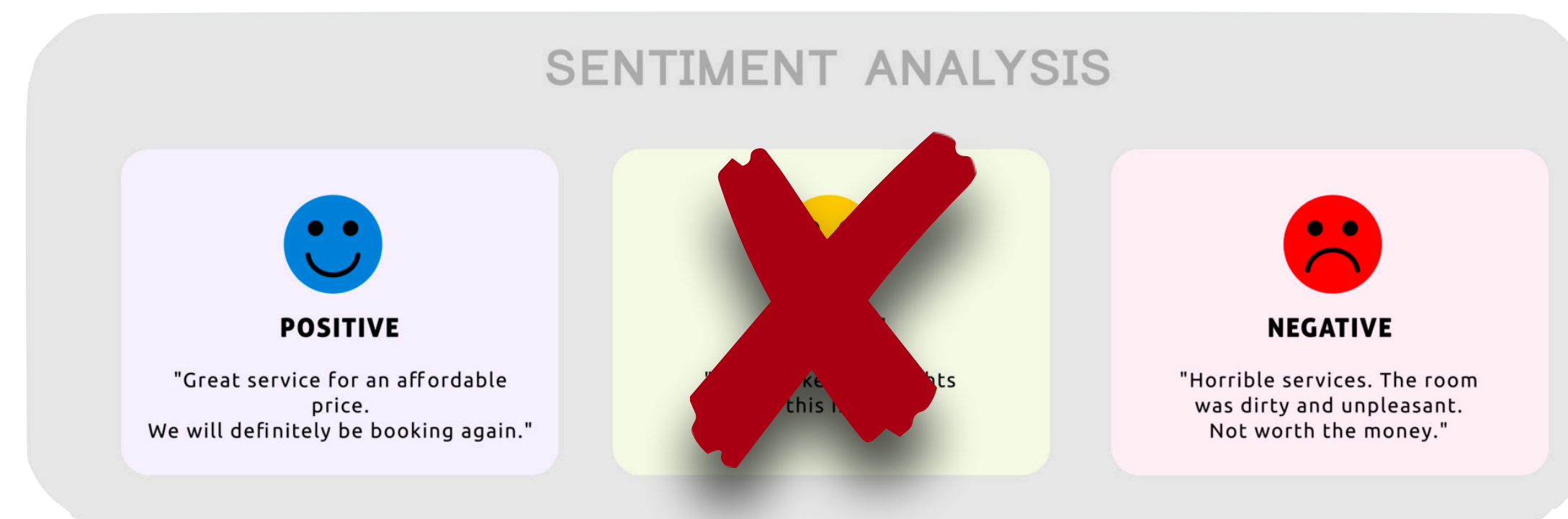
# One flavor of Classification: Logistic Regression

- Baseline supervised machine learning tool for classification; the foundation of neural networks
- Other classification algorithms: Naïve Bayes, K-Nearest Neighbors, Decision Trees, SVMs
- Input observation: vector of features, $\mathbf{x} = [x_1, x_2, \ldots, x_d]$
- Model: One weight per feature: $\mathbf{w} = [w_1, w_2, \ldots, w_d]$
- Output: a predicted class
  - Binary logistic regression $\hat{y} \in \{0,1\}$
  - Multinomial logistic regression (e.g. 5 classes): $\hat{y} \in \{0,1,2,3,4\}$



8

# Case Study: Sentiment Classification

- Data: Features $x_i$ of the input feature
  - $x_i =$ "review contains 'awesome'"; $w_i = +10$
  - $x_j =$ "review contains 'abysmal'"; $w_j = -10$
  - $x_k =$ "review contains 'mediocre'"; $w_k = -2$
- Features can be manually designed (classic) or learned (modern / neural architectures)
- Each $x_i$ is associated with a weight $w_i$ which determines how important $x_i$ is
  - (For predicting the positive class)
- Labels $y \in \{0,1\}$
  - $y = 0$ represents negative sentiment and
  - $y = 1$ represents positive sentiment

9

# How to get the right $y$?

- For each feature $x_i$, introduce a weight $w_i$, which determines the importance of $x_i$
  - Sometimes we have a bias term, $b$ or $w_0$, which is just another weight not associated to any feature
  - Together, all parameters can be termed as $\theta = [w; b]$
- We consider the weighted sum of all features and the bias

$$z = \left( \sum_d w_d x_d + b \right)$$

$$= \mathbf{w} \cdot \mathbf{x} + b$$

If high, $\hat{y} = 1$   If low, $\hat{y} = 0$

But how to determine the threshold?

We need probabilistic models!

$$P(y = 1 \mid \mathbf{x}; \theta)$$

$$P(y = 0 \mid \mathbf{x}; \theta)$$

# Solution: Squish it into the 0-1 range

$$z = \mathbf{w} \cdot \mathbf{x} + b \qquad z \in \mathbb{R}$$

- Sigmoid Function, $\sigma(\cdot)$
  - Non-linear!
- Compute $z$ and then pass it through the sigmoid function
- Treat it as a probability!
- Also, a differentiable function, which makes it a good candidate for optimization (more on this later!)

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

11

# Sigmoids and Probabilities

$$P(y = 1 \mid \mathbf{x}; \theta) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$= \frac{1}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))}$$

$$P(y = 0 \mid \mathbf{x}; \theta) = 1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$= 1 - \frac{1}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))}$$

$$= \frac{\exp(-(\mathbf{w} \cdot \mathbf{x} + b))}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))}$$

$$= \frac{1}{1 + \exp(\mathbf{w} \cdot \mathbf{x} + b)}$$

$$= \sigma(-(\mathbf{w} \cdot \mathbf{x} + b))$$

# Classification Decision

$$\hat{y} = \begin{cases} 1 & \textbf{if } p(y = 1 \,|\, x) > 0.5 \\ 0 & \textbf{otherwise} \end{cases}$$

**Decision Boundary**

$$\hat{y} = \begin{cases} 1 & \textbf{if } \mathbf{w} \cdot \mathbf{x} + b > 0 \\ 0 & \textbf{if } \mathbf{w} \cdot \mathbf{x} + b \leq 0 \end{cases}$$
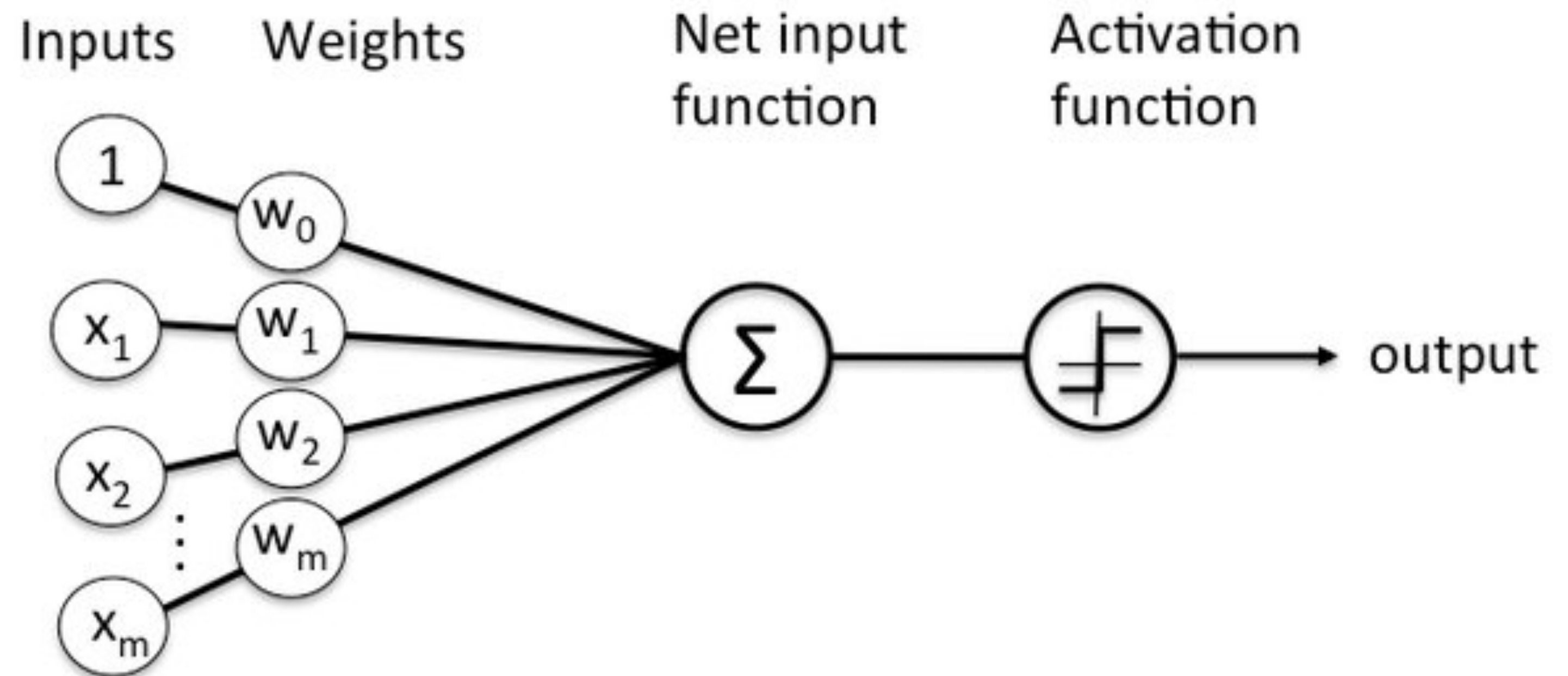


Sigmoid Function

$P(y = 1 \,|\, \mathbf{x}; \mathbf{w}, b)$

$\mathbf{w} \cdot \mathbf{x} + b$

Often we will denote $\hat{y}$ directly as $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

# Another notation

$$z = \left( \sum_m w_m x_m + w_0 \right)$$

$$= \mathbf{w} \cdot \mathbf{x}$$

Inputs    Weights    Net input function    Activation function



$$\hat{y} = P(y = 1 \mid \mathbf{x}; \mathbf{w}) = \sigma(z)$$

$$= \frac{1}{1 + \exp(-z)}$$

14

# Example: Sentiment Classification

It's hokey. There are virtually no surprises , and the writing is second-rate .
So why was it so enjoyable ? For one thing , the cast is
great . Another nice touch is the music I was overcome with the urge to get off
the couch and start dancing . It sucked me in , and it'll do the same to you .

$x_2=2$
$x_3=1$
$x_1=3$ $x_5=0$ $x_6=4.19$
$x_4=3$

$$\hat{y} = \begin{cases} 1 & \text{if } p(y = + \mid x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

| Var | Definition | Value in Fig. 5.2 |
|---|---|---|
| $x_1$ | count(positive lexicon) $\in$ doc) | 3 |
| $x_2$ | count(negative lexicon) $\in$ doc) | 2 |
| $x_3$ | $\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ | 1 |
| $x_4$ | count(1st and 2nd pronouns $\in$ doc) | 3 |
| $x_5$ | $\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ | 0 |
| $x_6$ | log(word count of doc) | $\ln(66) = 4.19$ |

$$p(+|x) = P(Y = 1|x)$$
$$= \sigma(w \cdot x + b)$$
$$= \sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \cdot [3, 2, 1, 3, 0, 4.19] + 0.1)$$
$$= \sigma(.833)$$
$$= 0.70$$

$$p(-|x) = P(Y = 0|x) = 1 - \sigma(w \cdot x + b)$$
$$= 0.30$$

# Model = the $\mathbf{w}$'s and the $b$'s

- Supervised Classification:
  - We know the correct label $y$ (either 0 or 1) for each $x$
  - But what the system produces is an estimate, $\hat{y}$
- Set $\mathbf{w}$ and $b$ to minimize the **distance** between our estimate $\hat{y}^{(i)}$ and the true $y^{(i)}$
  - We need a distance estimator: a **loss function** or a **cost function**
  - We need an **optimization algorithm** to update $\mathbf{w}$ and $b$ to minimize the loss.

Loss function

Optimization Algorithm

# Loss: the distance between $\hat{y}$ and $y$

- We want to know how far is the classifier output $\hat{y} = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$ is from the true (ground truth / gold standard) label, $y \in \{0,1\}$

- This difference is called the loss or cost, $L(\hat{y}, y)$ = how much $\hat{y}$ differs from $y$
  - Loss indicates how much would you lose if you mispredicted or how much would it cost you to mispredict

- We want the distance to be as small as possible
  - Hence, we choose the parameters $\mathbf{w}, b$ that maximize the log probability of the true labels $y$ in the training data, **given** the observations $x$
  - Maximum likelihood estimate for $p(y | \mathbf{x}; \mathbf{w}, b)$

# Maximizing conditional likelihood

**Goal**: maximize probability of the correct label $p(y \mid \mathbf{x})$

Maximize: $p(y \mid x) = \hat{y}^y(1 - \hat{y})^{1-y}$

Or, maximize: $\log p(y \mid x) = \log(\hat{y}^y(1 - \hat{y})^{1-y})$

Whatever values maximize $\log p(y \mid x)$ will also maximize $p(y \mid x)$

$$= y \log \hat{y} + (1 - y)\log(1 - \hat{y})$$

Now flip the sign for something to minimize (we minimize the loss / cost)

Minimize: $L_{CE}(y, \hat{y}) = -\log p(y \mid x) = -[y \log \hat{y} + (1 - y)\log(1 - \hat{y})]$

Cross-Entropy Loss

$$= -[y \log \sigma(\mathbf{w} \cdot \mathbf{x} + b) + (1 - y)\log \sigma(-(\mathbf{w} \cdot \mathbf{x} + b))]$$

18

# Case Study: Sentiment Classification

We want loss to be:
- smaller if the model estimate is close to correct
- bigger if model is confused

> It's hokey. There are virtually no surprises, and the writing is second-rate. So why was it so enjoyable? For one thing, the cast is great. Another nice touch is the music. I was overcome with the urge to get off the couch and start dancing. It sucked me in, and it'll do the same to you.

The loss when the model is right (if true $y = 1$):

$$L_{\text{CE}}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log(1 - \sigma(w \cdot x + b))]$$
$$= -[\log \sigma(w \cdot x + b)]$$
$$= -\log(.70)$$
$$= .36$$

..is lower than the loss when the model was wrong (if true $y = 0$)

> Loss is bigger when the model is wrong!

$$L_{\text{CE}}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log(1 - \sigma(w \cdot x + b))]$$
$$= -[\log(1 - \sigma(w \cdot x + b))]$$
$$= -\log(.30)$$
$$= 1.2$$

# Lecture Outline

- Recap: Supervised Classification / Logistic Regression
  - Data (Features)
  - Model (Logistic Regression)
  - Loss (Cross Entropy)
- Today: Logistic Regression Contd.
  - Optimization
  - Regularization
  - Multinomial Logistic Regression
- Word Embeddings (Teaser)

# IV. Optimization:
# Stochastic Gradient Descent

# Our goal: minimize the loss

- Loss function is parameterized by weights: $\theta = [\mathbf{w}; b]$
- We will represent $\hat{y}$ as $f(x; \theta)$ to make the dependence on $\theta$ more obvious

We want the weights that minimize the loss, averaged over all examples:

$$L_{CE}(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$$

# Intuition for gradient descent

How to get to the bottom of the river canyon?

- Look around 360°
- Find the direction of steepest slope down
- Go that way



What if multiple equally good alternatives?

# Logistic Regression: Loss



Image Credit: Medium

**Convex function**

**Non-convex function**

- Has only one option for steepest gradient
  - Or one minimum
- Gradient descent starting from any point is guaranteed to find the minimum

**Neural Networks - multiple alternatives**

# Consider: a single scalar $w$

Given current $w$, should we make it bigger or smaller?

Loss

Should we move
right or left from here?

Move $w$ in the reverse direction from the
slope of the function

slope of loss at $w^1$
is negative

$w^1$

$w^{min}$

$w$

$0$

$(goal)$

need to move positive

# Gradients

Loss

Should we move
right or left from here?

The **gradient** of a function of many variables is a vector pointing in the direction of the greatest increase in a function.

Find the gradient of the loss function at the current point and move in the **opposite** direction.

$w^1$

$w^{min}$

$w$

*0*

*(goal)*

But by how much?

Gradient Descent

# Gradient Updates

- Move the value of the gradient $\dfrac{\partial}{\partial w}L(f(x;w),y^*)$, weighted by a learning rate $\eta$

- Higher learning rate means move $w$ faster

$\eta$ **Too high: the learner will take big steps and overshoot**

$$w_{t+1} = w_t - \eta\frac{\partial}{\partial w}L(f(x;w),y^*)$$

$\eta$ **Too low: the learner will take too long**

If parameter $\theta$ is a vector of $d$ dimensions:

The gradient is just such a vector; it expresses the directional components of the sharpest slope along each of the $d$ dimensions.

# Under 2 dimensions

Consider 2 dimensions, $w$ and $b$:

Visualizing the gradient vector at the red point

It has two dimensions shown in the x-y plane



Cost(w,b)

b

w

# Real-life gradients, however…

- …are much longer; models usually contain lots and lots of weights!
- For each dimension $\theta_i$ the gradient component $i$ tells us the slope with respect to that variable
  - "How much would a small change in $\theta_i$ influence the total loss function $L$?"
- We express the slope as a partial derivative $\dfrac{\partial}{\partial \theta_i}$ of the loss, $\dfrac{\partial L}{\partial \theta_i}$
  - The gradient is then defined as a vector of these partials

# Real-life gradients

We will represent $\hat{y}$ as $f(x; \theta)$ to make the dependence on $\theta$ more obvious

$$\nabla_{\boldsymbol{\theta}} L(f(x; \boldsymbol{\theta}), y)) = \begin{bmatrix} \frac{\partial}{\partial w_1} L(f(x; \boldsymbol{\theta}), y) \\ \frac{\partial}{\partial w_2} L(f(x; \boldsymbol{\theta}), y) \\ \vdots \\ \frac{\partial}{\partial w_n} L(f(x; \boldsymbol{\theta}), y) \end{bmatrix}$$

The final equation for updating $\theta$ at time step $t + 1$ based on the gradient is thus:

$$\theta_{t+1} = \theta_t - \eta \frac{\partial}{\partial \theta} L(f(x; \theta), y)$$

# Gradients for Logistic Regression

Case: Sentiment Analysis

Recall: the cross-entropy loss for logistic regression

$$L_{CE}(y, \hat{y}) = -[y \log \sigma(\mathbf{w} \cdot \mathbf{x} + b) + (1 - y)\log(\sigma(-\mathbf{w} \cdot \mathbf{x} + b))]$$

Derivatives have a closed form solution:

$$\frac{\partial L_{CE}(y, \hat{y})}{\partial w_j} = [\sigma(\mathbf{w} \cdot \mathbf{x} + b) - y]x_j$$

# Pseudocode

function STOCHASTIC GRADIENT DESCENT $(L(), f(), x, y)$ returns $\theta$

   # where: $L$ is the loss function

   #        $f$ is a function parameterized by $\theta$

   #        $\mathbf{x}$ is the set of training inputs $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots \mathbf{x}^{(N)}$

   #        $y$ is the set of training outputs (labels) $y^{(1)}, y^{(2)}, \ldots y^{(N)}$

$\theta \leftarrow 0$ (or randomly initialized)

**repeat** till done

  for each training tuple$(x^{(i)}, y^{(i)})$: (in random order)

      1. Compute $\hat{y}^{(i)} = f(\mathbf{x}^{(i)}; \theta)$         # What is our estimated output $\hat{y}^{(i)}$?

      2. Compute the loss $L(\hat{y}^{(i)}, y^{(i)})$   # How far off is $\hat{y}^{(i)}$ from the true output $y^{(i)}$ ?

      3. $g \leftarrow \nabla L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$      # How should we move $\theta$ to maximize loss?

      4. $\theta \leftarrow \theta - \eta g$               # Go the other way instead

return $\theta$

**Stochastic Gradient Descent**

# Mini-Batching

function STOCHASTIC GRADIENT DESCENT $(L(), f(), x, y, m)$ returns $\theta$

   # where: $L$ is the loss function

   #       $f$ is a function parameterized by $\theta$

   #       $\mathbf{x}$ is the set of training inputs $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots \mathbf{x}^{(N)}$

   #       $y$ is the set of training outputs (labels) $y^{(1)}, y^{(2)}, \ldots y^{(N)}$ and $m$ is the mini-batch size

$\theta \leftarrow 0$ (or randomly initialized)

**repeat** till done

  for each randomly sampled minibatch of size $m$:

      1. for each training tuple$(\mathbf{x}^{(i)}, y^{(i)})$ in the minibatch: (in random order)

        i.  Compute $\hat{y}^{(i)} = f(\mathbf{x}^{(i)}; \theta)$           # What is our estimated output $\hat{y}^{(i)}$?

        ii.  Compute the loss $L_{mini} \leftarrow L_{mini} + L(\hat{y}^{(i)}, y^{(i)})$   # How far off is $\hat{y}^{(i)}$ from the true output $y^{(i)}$ ?

      2. $g \leftarrow \dfrac{1}{m} \nabla L_{mini}(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$       # How should we move $\theta$ to maximize loss?

      3. $\theta \leftarrow \theta - \eta g$                     # Go the other way instead

return $\theta$

# Lecture Outline

- Recap: Supervised Classification / Logistic Regression
  - Data (Features)
  - Model (Logistic Regression)
  - Loss (Cross Entropy)
- Today: Logistic Regression Contd.
  - Optimization
  - Regularization
  - Multinomial Logistic Regression
- Word Embeddings (Teaser)

# Regularization

# Overfitting

- A model that perfectly match the training data has a problem        Why?
- It will also overfit to the data, modeling noise
  - A random word that perfectly predicts $y$ (it happens to only occur in one class) will get a very high weight
  - Failing to generalize to a test set without this word

A good model should be able to generalize

What happens when a feature only occurs with one class?

e.g. word "wow" for positive reviews

# Overfitting: Features

This movie drew me in, and it'll do the same to you.

I can't tell you how much I hated this movie. It sucked.

Useful or harmless features

$x_1$ = "this"
$x_2$ = "movie
$x_3$ = "hated"
$x_4$ = "drew me in"

4-gram features that just "memorize" training set and might cause problems

$x_5$ = "the same to you"
$x_6$ = "tell you how much"

# Overfitting

- 4-gram model on tiny data will just memorize the data
  - 100% accuracy on the training set
- But it will be surprised by the novel 4-grams in the test data
  - Low accuracy on test set
- Models that are too powerful can overfit the data
  - Fitting the details of the training data so exactly that the model doesn't generalize well to the test set

How to avoid overfitting?

Regularization in logistic regression

Dropout in neural networks

# Regularization

- A solution for overfitting: Add a regularization term $R(\theta)$ to the loss function
  - (for now written as maximizing logprob rather than minimizing loss)

- Idea: choose an $R(\theta)$ that penalizes large weights
  - fitting the data well with lots of big weights not as good as
  - fitting the data a little less well, with small weights

$$\hat{\theta} = \arg\max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | \mathbf{x}^{(i)}) - \alpha R(\theta)$$

# L2 / Ridge Regularization

- The sum of the squares of the weights
- The name is because this is the (square of the) L2 norm $\|\theta\|_2^2$, = Euclidean distance of $\theta$ to the origin.

$$R(\theta) = \|\theta\|_2^2 = \sum_{j=1}^{d} \theta_j^2$$

L2 regularized objective function:

$$\hat{\theta} = \arg\max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)}) - \alpha \sum_{j=1}^{d} \theta_j^2$$

# L1 / Lasso Regularization

- The sum of the (absolute value of the) weights
- Named after the L1 norm $\|\theta\|_1$ = sum of the absolute values of the weights = Manhattan distance

$$R(\theta) = \|\theta\|_1 = \sum_{j=1}^{d} |\theta_j|$$

L1 regularized objective function:

$$\hat{\theta} = \arg\max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)}) - \alpha \sum_{j=1}^{d} |\theta_j|$$

# Lecture Outline

- Recap: Supervised Classification / Logistic Regression
  - Data (Features)
  - Model (Logistic Regression)
  - Loss (Cross Entropy)
- Today: Logistic Regression Contd.
  - Optimization
  - Regularization
  - Multinomial Logistic Regression
- Word Embeddings (Teaser)

# Multinomial Logistic Regression

# Multinomial Logistic Regression

- Often we need more than 2 classes
  - Positive / negative / neutral sentiment of a document
  - Parts of speech of a word (noun, verb, adjective, adverb, preposition, etc.)
  - Actionable classes for emergency SMSs
- If >2 classes we use multinomial logistic regression
  - = Softmax regression
  - = Multinomial logit
  - = (defunct names : Maximum entropy modeling or MaxEnt)
- So "logistic regression" will just mean binary (2 output classes)

# Multinomial Logistic Regression

The probability of everything must still sum to 1

$$P( + | x) + P( - | x) + P( \sim | x) = 1$$

- Need a generalization of the sigmoid!

  Softmax

- Introducing the softmax function, which
  - Takes a vector $\mathbf{z} = [z_1, z_2, \ldots, z_K]$ of $K$ arbitrary values
    - Each $z_i$ corresponds to weighted sum of features for the $K$th class
  - Outputs a probability distribution
    - each value in the range $[0,1]$
    - all the values summing to 1

# The Softmax Function

Turns a vector $\mathbf{z} = [z_1, z_2, \ldots, z_K]$ of $K$ arbitrary values into probabilities

$$\textbf{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^{K} \exp(z_j)} \quad 1 \leq i \leq K$$

The denominator $\displaystyle\sum_{i=1}^{K} \exp(z_i)$ is used to normalize all the values into probabilities.

$$\textbf{softmax}(\mathbf{z}) = \left[ \frac{\exp(z_1)}{\sum_{i=1}^{K} \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^{K} \exp(z_i)}, \ldots, \frac{\exp(z_K)}{\sum_{i=1}^{K} \exp(z_i)} \right]$$

# Softmax: Example

Turns a vector $\mathbf{z} = [z_1, z_2, \ldots, z_K]$ of $K$ arbitrary values into probabilities

$$\mathbf{z} = [0.6, \ 1.1, \ 1.5, \ 1.2, \ 3.2, \ 1.1]$$

$$\mathbf{softmax}(\mathbf{z}) = \left[ \frac{\exp(z_1)}{\sum_{i=1}^{K} \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^{K} \exp(z_i)}, \ldots, \frac{\exp(z_K)}{\sum_{i=1}^{K} \exp(z_i)} \right]$$

$$\mathbf{softmax}(\mathbf{z}) = [0.055, \ 0.090, \ 0.0067, \ 0.10, \ 0.74, \ 0.010]$$

# Binary versus Multinomial

**Binary Logistic Regression**

$$y = 1 \qquad\qquad y = 0$$

$$w_5 > 0 \qquad\qquad w_5 \leq 0$$

$$w_5$$

$$x_5 = \begin{cases} 1 & \text{if } \text{``!''} \in \text{doc} \\ 0 & \text{otherwise} \end{cases} \qquad w_5 = 3.0$$

Why do we NOT need a different weight for each class in binary logistic regression?

**Multinomial Logistic Regression**

$$y = + \qquad\qquad y = \sim \qquad\qquad y = -$$

$$w_{5,+} \qquad\qquad w_{5,\sim} \qquad\qquad w_{5,-}$$

Separate weights for each class

| Feature | Definition | $w_{5,+}$ | $w_{5,-}$ | $w_{5,0}$ |
|---|---|---|---|---|
| $f_5(x)$ | $\begin{cases} 1 & \text{if } \text{``!''} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ | 3.5 | 3.1 | $-5.3$ |

# Softmax in multinomial logistic regression

Parameters are now a matrix $\mathbf{W} \in \mathbb{R}^{d \times K}$ and $b \in \mathbb{R}^1$

$$P(y = c \,|\, \mathbf{x}; \theta) = \frac{\exp(\mathbf{w}_c \cdot \mathbf{x} + b)}{\sum_{j=1}^{K} \exp(\mathbf{w}_j \cdot \mathbf{x} + b)}$$

- Input is still the dot product between weight vector $\mathbf{w}_c$ and input vector $\mathbf{x}$, offset by $b$
- But separate weight vectors for each of the $K$ classes, each of dimension $d$

Multinomial LR Loss:

$$L_{CE} = -\log P(y = c \,|\, \mathbf{x}; \theta) = -(\mathbf{w}_c \cdot \mathbf{x} + b) + \log\left[\sum_{j=1}^{K} \exp(\mathbf{w}_j \cdot \mathbf{x} + b)\right]$$

49

# Lecture Outline

- Recap: Supervised Classification / Logistic Regression
  - Data (Features)
  - Model (Logistic Regression)
  - Loss (Cross Entropy)
- Today: Logistic Regression Contd.
  - Optimization
  - Regularization
  - Multinomial Logistic Regression
- Word Embeddings (Teaser)

# Case Study: Word Neighborhood Predictor

- Yet do I fear thy nature.  It is too full o'the milk of human kindness  To catch the nearest way.
- If it were done when 'tis done, then 'twere well It were done quickly.
- Is this a dagger which I see before me,  The handle toward my hand?
- My hands are of ;your color, but I shame   To wear a heart so white.
- Knock, knock, knock! Who's there,

Where is it more likely to find the word "doth"?

# Case Study: Word Neighborhood Predictor

What words are likely to co-occur with the word "garnish"?

## Ingredients

12 Large Eggs, Hard-boiled And Peeled (See Note)

1/2 c. Mayonnaise

2 tbsp. Prepared Yellow Mustard

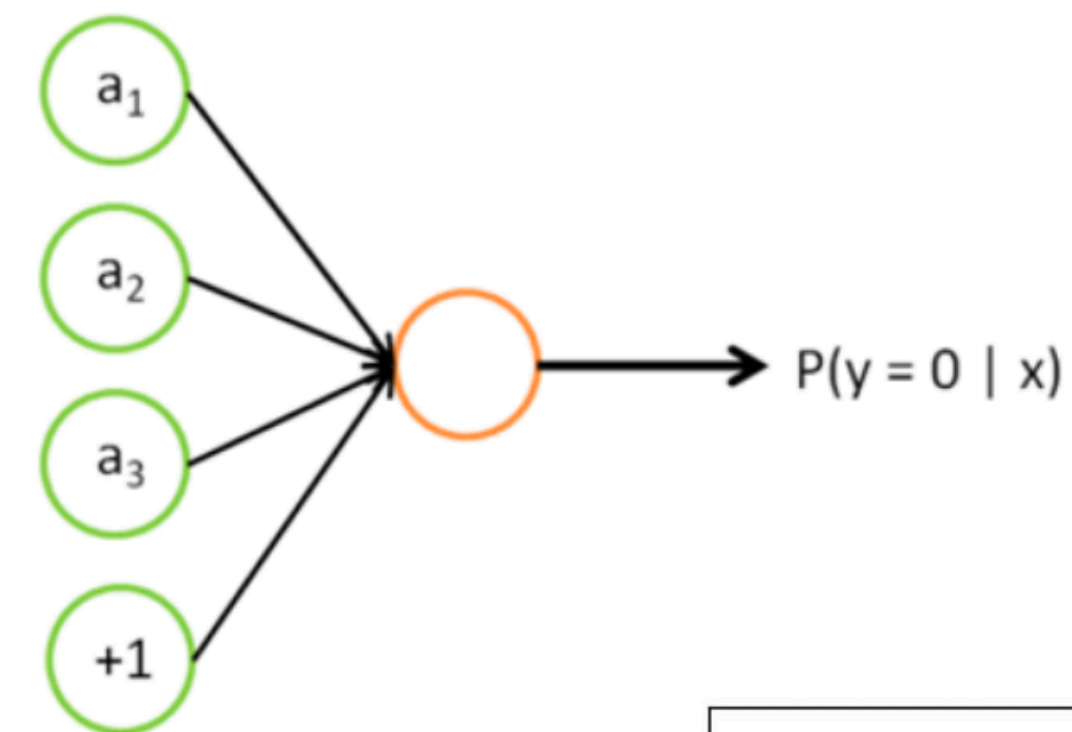1/4 tsp. Kosher Salt

See Nutritional Information ⌄

## Directions

Slice eggs in half lengthwise. Pop out yolks and place them in a medium-sized mixing bowl. Mash yolks with a fork and then add mayonnaise, mustard, and salt. Mix until smooth. To achieve maximum smoothness, blitz the yolk mixture a few times with an immersion blender. If you are planning to use a pastry bag and piping tip to add the yolk mixture to the egg whites, it helps to get the yolk mixture nice and smooth.

Divide the yolk mixture evenly between the egg white halves. Garnish as desired. Serve immediately or refrigerate until ready to serve.

Next Class: This is the text classification task we will use logistic regression for to learn word embeddings!

# Concluding Thoughts

- Logistic Regression
  - Case Study: Sentiment Classification
  - A very simple neural network
- Next Class: Let's learn some automatic features
  - Word Embeddings



Input (features)   Logistic classifier

**Logistic Regression**