

CS4248: Natural Language Processing

Lecture 5 — Introduction into Connectionist Machine Learning

Recap of Week 04

Document-Term Matrix with tf-idf Weights **Text Classification** Putting it all together Formal setup • X — set of all documents; $x \in X$ — a single document $w_{t,d} = (1 + \log_{10} t f_{t,d}) \cdot \log_{10} \frac{|D|}{df_{\star}}$ • Y — set of all classes (or class labels); $y \in Y$ — a single class (or class label) Classification task • Mapping h from input space X to output space $Y = h: X \to Y$ Side notes h(x) = ye.g., h("The movie is great.") = "positive" No real theoretic underpinning, but tf-idf works best in practice • Not all definitions of tf-idf apply a sublinear scaling of tf_{td} Alternative names: $tf \cdot idf$, $tf \times idf$ Note: A document might be assigned to more "True" mapping which than one class -> multilabel classification is unknown in practice There are different weighting functions for calculating tf-idf Note 2: Our SLP3 textbook uses d for x and c for v. We'll use both interchangeably. **Classification: Evaluation — Why so Many Measures? Naive Bayes Classifier + BoW — Discussion** Naive Bayes vs. Language Models Observation: FP and FN not $Precision = \frac{TP}{TP + FP} \qquad Recall = \frac{TP}{TP + FN}$ Naive Baves makes a non-contextual decision (unigram model; but can be extended to larger n-grams) always equally problematic Naive Bayes is an LM! It treats each class like a separate language model Example: Suicide prediction Biggest pro: simplicity (e.g., from social media content posted by users) BAD: misclassifying a high-risk person Easy to understand & implement, fast, not very data hungry, interpretable results Recall > Precision OK-ish: misclassifving a healthy person Biggest con: assumption of conditional independence For most types of data, the features are typically not independent. Example: News article classification (e.g., for search engines such as Google News) For text classification (features = words) it actually often works well in practice BAD: showing article of wrong category (particularly with some additional "tweaking" of the data) Recall < Precision OK: missing a relevant article in result

56

Announcements

Project Groups Announced

- There may have been some errors, please check your team and update us per our announcement if you see anything amiss.
- Project's Intermediate Update Rubric / Template is available
 - Find in Canvas >> Files >> Project
 - Live version (best bet) at <u>https://bit.ly/cs4248-2320-iu-template</u>
- Assignment 2 out on Saturday, once Assignment 1 is in
 - Assignment 2 will be a Text Classification competition, restricted to ML algorithms taught (Naive Bayes and Logistic Regression)
 - Emphasis on Natural Language Feature Engineering

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

Text Classification (well, for classification, in general)

- Formal setup
 - X set of all documents; $x \in X$ a single document
 - Y set of all classes (or class labels); $y \in Y$ a single class (or class label)
 - Mapping h from input space X to output space $Y \twoheadrightarrow h: X \to Y$
 - ightarrow Find best \hat{h} to approximate the true mapping h

We find \hat{h} by <u>learning</u> \hat{h} from the data **→** Supervised (Machine) Learning

• Probabilistic Classifiers (e.g., Naive Bayes)

Instead of $\hat{h}: X \to Y$, learn $\hat{P}(Y|X)$ (or $\hat{P}(y|x)$ for an $\langle x, y \rangle$ pair)

Text Classification — Probabilistic Classifiers

- Common goal: Learn P(y|x)
 - Learn P(y|x) from the data
- Two basic approaches
 - (1) Generative Classifiers
 - Learn joint probability P(x, y)
 - \blacksquare Apply Bayes Rule to get P(y|x)
 - (2) Discriminative Classifiers
 - $\bullet \quad \text{Learn } P(y|x) \text{ directly} \\$

$$\hat{y} = \operatorname*{argmax}_{y \in Y} \overbrace{P(x|y)P(y)}^{P(y|x)}$$

Generative vs. Discriminative Classifiers — Intuition

• Task: Train a classifier to distinguish zebra from elephants images





Generative vs. Discriminative Classifiers — Intuition

- Generative classifier
 - Builds 2 models of what zebra and elephant images look like -



Some abstract internal representation / model of language and the world

Feature <i>x_i</i>	P(x _i , zebra)	P(<i>x_i</i> , elephant)
"is grey"	0.32	0.95
"is striped"	0.99	0.08
"long nose"	0.40	0.98
"four legs"	0.88	0.99

- Models allow to assign a "zebra probability" and an "elephant probability" to any image (using Bayes Rule)
- Givan a new image:

Run both models and see which fits better





Generative vs. Discriminative Classifiers — Intuition

• Discriminative classifier

- Tries to distinguish zebra and elephant images
- No model of how zebra and elephant images "look like"

Question: How could we quickly distinguish zebras from elephants?





Generative vs. Discriminative Classifiers

Generative classifier

- Learn data distribution of each class
- Classifies new data item by comparing the item with each class distribution







Discriminative classifier

- Learn the decision boundaries between classes
- Classifies new data item based on in which "region" the new item falls

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

Linear Models

- Underlying assumption:
 - \blacksquare There exists a linear relationship between $\ x^{(j)}$ and dependent variable $\ y^{(j)}$

Linear Models — More User-Friendly Notation

- Vector representation
 - **Bias Trick:** Introduce constant feature $x_0^{(j)}$

$$h_{\theta}\left(x^{(j)}\right) = f\left(\theta_{0}\underline{x}_{0}^{(j)} + \theta_{1}x_{1}^{(j)} + \theta_{2}x_{2}^{(j)} + \dots + \theta_{n}x_{n}^{(j)}\right)$$

Represent $x^{(j)}$ with new constant feature

$$x^{(j)} = \left(1, x_1^{(j)}, x_2^{(j)}, \dots, x_n^{(j)}\right)$$

Rewrite linear relationship using vectors representing $x^{(j)}$ and heta

$$h(x^{(j)}) = f(\theta^{\top} x^{(j)}) \qquad \theta = \{\theta_0, \theta_1, \theta_2, \dots, \theta_n\}, \ \theta_i \in \mathbb{R}$$

Note: Throughout the rest of the slide, we drop the superscript in $x^{(j)}$ and $y^{(j)}$ if there is no ambiguity.

Map $y \in \mathbb{R}$ to $\sigma \in [0,1]$

• $y = f(\theta^{\top} x^{(j)})$ calculates a real-valued output, but we want a **probability** bounded between 0 and 1.

• We need a squashing function σ that maps real numbers to the unit interval, and where $\sigma(0)=0.5$



Logistic Regression

- Logistic Regression → Real-valued predictions interpreted as probability
 - Function *f* is the standard **Logistic Function** (Sigmoid function)

$$f(x) = \frac{L}{1 + e^{-k(x-x_0)}} \xrightarrow{L = 1, \ k = 1, \ x_0 = 0} f(x) = \frac{1}{1 + e^{-x}}$$

Logistic Regression — Probabilistic Interpretation

• \hat{y} interpreted as a probability

$$\hat{y} = h_{\theta}(x) = f(\theta^{\top}x) = \frac{1}{1 + e^{-\theta^{\top}x}} \quad \text{with} \quad \hat{y} \in [0, 1]$$

→ $\hat{y} = h_{\theta}(x)$ is the estimated probability that y = 1 given x and θ $\hat{y} = P(y = 1 | x, \theta)$

→ Given only discrete 2 outcomes: $P(y = 1|x, \theta) + P(y = 0|x, \theta) = 1$

$$\hat{y} = 1 - P(y = 0|x,\theta)$$

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

Sentiment Analysis, redux



Now Showing

Movies – The Omicron Variant

"It's hokey. There are no surprises, the writing is poor. 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."

> Photoshopped Fake Vintage Movie Poster image courtesy <u>Tribune India</u>

In-Lecture Activity (2 mins)



$1 \times 1 \times 1$ What features have positive or negative weights?

Feature	Description	Value	Weight
<i>x</i> ₁	Number of positive words		
x ₂	Number of negative words		
x ₃	1 if "no" in text; 0 otherwise		
<i>x</i> ₄	Number of 1st & 2nd person pronouns		
<i>x</i> ₅	1 if "!" in text; 0 otherwise		
x ₆	In of word/token count		

Sentiment Analysis for movie reviews

"It's hokey. There are no surprises, the writing is poor. 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."

Feature	Description	Value
X ₁	Number of positive words	
x ₂	Number of negative words	
x ₃	1 if "no" in text; 0 otherwise	
X ₄	Number of 1st & 2nd person pronouns	
x ₅	1 if "!" in text; 0 otherwise	
x ₆	In of word/token count	

Side notes:

- Naive Bayes and Logistic Regression require feature engineering as they do not combine primitive features into composite ones.
- The 6 features on the left are chosen for simplicity; in practice, these are often tf-idf weighted vocabulary.

• Step 1: Extract feature values

"It's hokey. There are no surprises, the writing is poor. 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."

Feature	Description	Value
X ₁	Number of positive words	3
x ₂	Number of negative words	2
x ₃	1 if "no" in text; 0 otherwise	1
X ₄	Number of 1st & 2nd person pronouns	3
x ₅	1 if "!" in text; 0 otherwise	0
x ₆	In of word/token count	ln(66) = 4.19

🏃 🏃 🏃 In-Lecture Activity (5 mins)

- Question: What might be other useful features for a sentiment classifier?
 - Bonus: Briefly discuss how easy/difficult it would be to extract your features

 Post your answer to Canvas > Discussions > [In-Lecture Interaction] L1 (Help like other classmate's responses too!

- Step 2: Factor in weights θ
 - Let's assume some oracle gave us those weights
 - It's time to include the bias using the "bias trick"

▲ Notation varies: Weights are also called parameters, sometimes denoted as w (as used in the SLP3 textbook)

Feature	Description	Value	Weight θ_i
x _o	Bias b	1	0.1
x ₁	Number of positive words	3	2.5
x ₂	Number of negative words	2	-5.0
x ₃	1 if "no" in text; 0 otherwise	1	-1.2
x4	Number of 1st & 2nd person pronouns	3	0.5
x ₅	1 if "!" in text; 0 otherwise	0	2.0
x ₆	In of word/token count	4.19	0.7

• Step 4: Compute linear signal (sum of weighted features)

Feature	Description	Value	Weight θ_i	θ _i x _i
x _o	Bias b	1	0.1	0.1
x ₁	Number of positive words	3	2.5	7.5
x ₂	Number of negative words	2	-5.0	-10.0
x ₃	1 if "no" in text; 0 otherwise	1	-1.2	-1.2
x ₄	Number of 1st & 2nd person pronouns	3	0.5	1.5
х ₅	1 if "!" in text; 0 otherwise	0	2.0	0
х ₆	In of word/token count	4.19	0.7	2.933

 $\rightarrow \theta^{\top} x = 0.833$

Vector notation:

$$\begin{array}{c} x = (1, 3, 2, 1, 3, 0, 4.19)^{\top} \\ \theta = (0.1, 2.5, -5.0, -1.2, 0.5, 2.0, 0.7)^{\top} \end{array} \right\}$$

$$\sum_{i=0.833}$$

$$P(+|x) = P(y = 1|x, \theta) = \sigma(\theta^{\top}x) = \frac{1}{1 + e^{-\theta^{\top}x}} = \frac{1}{1 + e^{-0.833}} = 0.7$$

$$P(-|x) = P(y = 0|x, \theta) = 1 - P(y = 1|x, \theta) = 0.3$$

$$P(+|x) > 0.5 \rightarrow \hat{y} = + \text{(positive)}$$

Classify movie review as "positive"

$$\theta^{\top} x = 0.833$$

Logistic Regression

• So, where did the values for θ come from?

(in the example, they were simply given to us)

- Of course, different θ values would have resulted in different probabilities
- Break down into 2 questions
 - (1) How can we quantify how good a set of θ values is?
 - → Loss function (also: cost function, error function)
 - (2) How can we systematically find the best θ values?
 - → Gradient Descent (numerical method to minimize loss function)

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

Logistic Regression — Loss Function

- Intuition: A set of values for θ is good if
 - the correct label \mathcal{Y} (0 or 1; coming from the dataset)
 - the model's estimated label $\,\hat{y} = \sigma(\theta^\top x)$

are similar for all $\langle x,y
angle$ pairs

 \rightarrow Find θ that minimizes the difference between \hat{y} and y



Logistic Regression — Loss Function

 $\hat{y} = \frac{1}{1 + e^{-\theta^{\top}x}}$

• Goal: Maximize probability of the correct label P(y|x)

$$\hat{y} = P(y = 1 | x, \theta) = 1 - P(y = 0 | x, \theta)$$

- Intermediate step: Combine both case into one formula
 - P(y|x) is a Bernoulli distribution (2 discrete outcomes)

$$P(y|x) = \begin{cases} \hat{y} & , y = 1\\ 1 - \hat{y} & , y = 0 \end{cases}$$

→ Combine into:
$$P(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

Logistic Regression — Loss Function



- Goal: Maximize probability of the correct label P(y|x)
 - Find θ that maximizes

$$P(y|x) = \hat{y}^{y}(1-\hat{y})^{1-y}$$

$$\log P(y|x) = \log \left[\hat{y}^{y}(1-\hat{y})^{1-y}\right]$$

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

■ Find *θ* that **minimizes**

$$L_{CE}(\hat{y}, y) = -P(y|x) = \underbrace{-\left[y\log\hat{y} + (1-y)\log\left(1-\hat{y}\right)\right]}_{\text{Cross-Entropy Loss}}$$

Cross-Entropy Loss — Visualization

$$L_{CE}(\hat{y}, y) = -\left[y \log \hat{y} + (1 - y) \log (1 - \hat{y})\right]$$



Cross-Entropy Loss — Runthrough Example (Part 2)

Recall:

$$P(+|x) = \sigma(\theta^{\top}x) = 0.7$$

$$P(-|x) = 1 - \sigma(\theta^{\top}x) = 0.3$$

Feature	Description	Value	Weight θ_i	θ _i x _i
x _o	Bias b	1	0.1	0.1
x ₁	Number of positive words	3	2.5	7.5
x ₂	Number of negative words	2	-5.0	-10.0
x ₃	1 if "no" in text; 0 otherwise	1	-1.2	-1.2
X ₄	Number of 1st & 2nd person pronouns	3	0.5	1.5
x ₅	1 if "!" in text; 0 otherwise	0	2.0	0
x ₆	In of word/token count	4.19	0.7	2.933

$$L_{CE}(\hat{y}, y) = -[y \log \hat{y} + (1 - y) \log (1 - \hat{y})]$$

Assume the model was right (y = 1)



Assume the model was wrong (y = 0)

$$L_{CE}(\hat{y}, y) = ???$$

Cross-Entropy Loss — Runthrough Example (Part 2)

 $P(+|x) = \sigma(\theta^{\top}x) = 0.7$ $P(-|x) = 1 - \sigma(\theta^{\top}x) = 0.3$ $L_{CE}(\hat{y}, y) = -\left[y \log \hat{y} + (1-y) \log (1-\hat{y})\right]$

Assume the model was right (y = 1) $L_{CE}(\hat{y}, y) = -[\log \hat{y}]$ $= -[\log 0.7]$ = 0.36 Assume the model was wrong (y = 0) $L_{CE}(\hat{y}, y) = -[\log (1 - \hat{y})]$ $= -[\log 0.3]$ = 1.2

Cross-Entropy Loss — Total Loss

• Loss for all training samples (given *m* data samples)

$$L_{CE} = rac{1}{m}\sum_{j=1}^m L_{CE}\left(\hat{y}^{(j)},y^{(j)}
ight)$$

$$x = -rac{1}{m}\sum_{j=1}^m \left[y^{(j)}\log \hat{y}^{(j)} + \left(1-y^{(j)}
ight)\log \left(1-\hat{y}^{(j)}
ight)
ight],$$

$$=-rac{1}{m}\sum_{j=1}^{m}\left[y^{(j)}\log\sigma\left(heta^{T}x^{(j)}
ight)+\left(1-y^{(j)}
ight)\log\left(1-\sigma\left(heta^{T}x^{(j)}
ight)
ight)
ight]$$

$$= -rac{1}{m}\sum_{j=1}^m \left[y^{(j)}\lograc{1}{1+e^{- heta^T x^{(j)}}} + \left(1-y^{(j)}
ight)\log\left(1-rac{1}{1+e^{- heta^T x^{(j)}}}
ight)
ight]$$

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

Learning — Minimizing the Loss Function

$$L_{CE} = -\frac{1}{m} \sum_{j=1}^{m} \left[y^{(j)} \log \frac{1}{1 + e^{\theta^{\top} x^{(j)}}} + \left(1 - y^{(j)}\right) \log \left(1 - \frac{1}{1 + e^{\theta^{\top} x^{(j)}}}\right) \right]$$

Visual illustration of loss function

- Just 1 feature θ_1 and bias θ_0
- Good news: L_{CE} for Logistic Regression is a convex function → 1 global minimum

\rightarrow How to find the minimum of L_{CE} ?

...this should cause a flashback to your calculus classes :)


Method 1: Random Search (the "stupid" way)

- Repeat "enough" times
 - Select random values for $\theta = \{\theta_0, \theta_1, \theta_2, \dots, \theta_n\}$
 - Calculate loss *L* for current θ
- Return θ with smallest loss

- Limitation:
 - Not practical beyond toy examples
- → Don't do that! :)



Method 2: Using Calculus (the proper way)

- Minimum of loss function L → Calculus to the rescue!
 - Partial derivatives with respect to to all θ_i are 0

$$\frac{\partial L}{\partial \theta_0} = 0, \ \frac{\partial L}{\partial \theta_1} = 0, \ \frac{\partial L}{\partial \theta_2} = 0, \ \dots, \ \frac{\partial L}{\partial \theta_n} = 0$$

n+1 equations with *n*+1 unknowns
 (→ 1 unique solution → 1 global minimum)





Loss Function — Derivatives

$$L_{CE} = -\frac{1}{m} \sum_{j=1}^{m} \left[y^{(j)} \log \sigma \left(\theta^{\top} x^{(j)} \right) + \left(1 - y^{(j)} \right) \log \left(1 - \sigma \left(\theta^{\top} x^{(j)} \right) \right) \right]$$

...lots of tedious math here...

$$\frac{\partial L_{CE}}{\partial \theta_i} = \frac{1}{m} \sum_{j=1}^m \left[\sigma \left(\theta^\top x^{(j)} \right) - y^{(j)} \right] x_i^{(j)}$$
$$\frac{\partial L_{CE}}{\partial \theta} = \frac{1}{m} X^\top \left[\sigma \left(X \theta \right) - y \right]$$

Basic approach to find the minimum (1) Set derivative to $0 \rightarrow \frac{1}{m} X^{\top} [\sigma (X\theta) - y] \stackrel{!}{=} 0$ (2) Solve for θ

So are we done here?

Gradient Descent

• Problem: $\frac{1}{m}X^{\top}[\sigma(X\theta) - y] \stackrel{!}{=} 0$ has no closed-form solution for θ

→ Gradient Descent

- Start with a random setting of θ
- Adjust *θ* iteratively to minimize L



Gradient — Quick Refresher

- Gradient
 - Vector of partial derivatives of a multivariable function (e.g., $\theta_0, \theta_1, \dots, \theta_n$)
 - Partial derivative: slope with respect to a single variable given a current set of values for all $\theta_0, \theta_1, \dots, \theta_n$
 - Points in the direction of the steepest ascent





Gradients — Runthrough Example (Part 3)

• Calculate Gradients (assuming y = 1)

Feature	Description	Value	Weight θ_i	θ _i x _i	Gradients
x _o	Bias b	1	0.1	0.1	-0.30
x ₁	Number of positive words	3	2.5	7.5	-0.91
x ₂	Number of negative words	2	-5.0	-10.0	-0.61
x ₃	1 if "no" in text; 0 otherwise	1	-1.2	-1.2	-0.30
X ₄	Number of 1st & 2nd person pronouns	3	0.5	1.5	-0.91
x ₅	1 if "!" in text; 0 otherwise	0	2.0	0	0.0
x ₆	In of word/token count	4.19	0.7	2.933	-1.27

 $\mathbf{T} = \begin{bmatrix} -0.30 \\ -0.91 \\ -0.61 \\ -0.30 \\ -0.91 \\ 0.0 \\ -0.91 \\ 0.0 \\ -1.27 \end{bmatrix}$

 $\frac{\partial L_{CE}}{\partial \theta} = \frac{1}{m} X^{\top} \left[\sigma \left(X \theta \right) - y \right]$

Gradients — Runthrough Example (Part 3)

- Interpretation of gradients
 - Negative values: a small <u>increase</u> in, e.g., θ_0 or θ_1 will <u>decrease</u> the loss
 - A small change in θ_1 affects the loss more than the same change in θ_0 (since the absolute value of θ_1 is larger than the one of θ_0)
 - Absolute values of gradient not a direct indicator of how to update θ

\rightarrow So how do we adjust θ to decrease the loss?



Gradient Descent Algorithm

- Important concept: learning rate
 - Scaling factor for gradient (typical range: 0.01 0.0001)

```
Input : data (X, y), loss function L, learning rate \eta
Initialization : Set \theta to random values
```

```
while true :
```

Calculate gradient $\nabla_{\theta} L$ $\theta \leftarrow \theta - (\eta \cdot \nabla_{\theta} L)$

In practice: stop loop when θ converges

Gradient Descent — Runthrough Example (Part 4)

- Update weights θ
 - Learning rate: $\eta=0.1$

 $\begin{array}{c} \theta \leftarrow \theta - (\eta \cdot \nabla_{\theta} L) \\ | \end{array}$

Feature	Description	Value	Weight θ_i	θ _i x _i	Partial derivatives	New Weight ${m heta}_i$
x _o	Bias b	1	0.1	0.1	-0.30	0.13
x ₁	Number of positive words	3	2.5	7.5	-0.91	2.59
x ₂	Number of negative words	2	-5.0	-10.0	-0.61	-4.94
x ₃	1 if "no" in text; 0 otherwise	1	-1.2	-1.2	-0.30	-1.17
x ₄	Number of 1st & 2nd person pronouns	3	0.5	1.5	-0.91	0.59
x ₅	1 if "!" in text; 0 otherwise	0	2.0	0	0.0	2.0
x ₆	In of word/token count	4.19	0.7	2.933	-1.27	0.83

→ 1st iteration of Gradient Descent done!

 $L_{CE} = 0.12$

(down from 0.36)

Gradient Descent — Runthrough Example (Part 4)

- Update weights θ
 - $\bullet~$ Learning rate: $\eta=0.1$

 $\theta \leftarrow \theta - (\eta \cdot \nabla_{\theta} L)$

Feature	Description	Value	Weight θ _i	θ _i x _i	Partial derivatives	New Weight $\boldsymbol{\theta}_i$
x _o	Bias b	1	0.13	0.13	-0.11	0.14
x ₁	Number of positive words	3	2.59	7.77	-0.33	2.62
x ₂	Number of negative words	2	-4.94	-9.88	-0.22	-4.92
x ₃	1 if "no" in text; 0 otherwise	1	-1.17	-1.17	-0.11	-1.16
x ₄	Number of 1st & 2nd person pronouns	3	0.59	1.77	-0.33	0.62
x ₅	1 if "!" in text; 0 otherwise	0	2.0	0	0.0	2.0
x ₆	In of word/token count	4.19	0.83	3.46	-0.46	0.87



(down from 0.12)

→ 2nd iteration of Gradient Descent done!

In-Lecture Activity (2 mins)





B

C

No loss will be calculated

What happens if Logistic Regression gets a training sample correct?



The loss will be small

The loss will be large

Effects of Learning Rate for $L = x^2$, $\frac{\partial L}{\partial x} = 2x$, 20 steps







Gradient Descent — Variations

- (Basic) Gradient Descent
 - Calculate gradient und update θ for whole dataset

- Stochastic Gradient Descent (SGD)
 - Calculate gradient and update θ for each data sample

- Mini-batch Gradient Descent
 - Calculate gradient and update θ for batches of sample
 - e.g., batch = 64 data samples
 - In practice, often still referred to as SGD





Gradient Descent — Variations

Gradient Descent

Mini-Batch Gradient Descent

Stochastic Gradient Descent



Gradient averaged over all data items

- Smooth descent
- Small(er) gradients
- Small(er) update steps



Gradient averaged over some data items

• Well, "somewhere in-between" :)

Gradient for each data item considered

Choppy descent

 $\hat{\theta}_1^2$

3

0

- Large(r) gradients
- Large(r) steps

30000

25000

20000

15000

10000

5000

100

80

60

00

40

20

0

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression
- Towards Neural Networks
 - Motivation: XOR Problem
 - Basic Neural Network Architecture

Gradient Descent — When to Stop?

• Intuition: $\nabla_{\theta} L_{CE} < threshold$

Problem: regions of "near-plateaus":

- → Gradient $\nabla_{\theta}L$ very small
- → Step $\eta \nabla_{\theta} L$ extremely small
- → Very slow convergence

- Alternative stop conditions:
 - Loss is small (enough)
 - Change in loss is small enough
 - Max. # iterations reached



Note: This problem is much more pronounced for non-convex loss functions with multiple local minima

In-Lecture Activity (2 mins)



• A model that perfectly matches the training data often has a problem.



• A model that perfectly matches the training data often has a problem.

It may 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 generalise to a test set without this words.

 A good model should be able to generalise

This movie drew me in, and it'll do the same to you.	positive
I can't tell you how much I hated this movie. It sucked.	negative

Overfitting — Intuition (Naive Bayes Classifier)

- Scenario movie reviews
 - (Very) low number of reviews
 - NB classifier based on 4-grams

This movie drew me in, and it'll do the same to you.	positive
I can't tell you how much I hated this movie. It sucked.	negative

→ Effect of Naive Bayes classifier

- Each 4-gram most likely unique and associated with only 1 class (e.g., *"tell you how much"* only found in a negative review)
- Unseen positive review x containing "tell you how much" $\rightarrow P(positive|x) = 0$

Overfitting — Intuition (Logistic Regression Classifier)

- Scenario movie reviews
 - (Very) low number of reviews
 - Assume the following artifact

All positive reviews contain many pronouns Almost no negative reviews contain pronouns

Feature	Description	Value	Weight θ_i	θ _i x _i
x _o	Bias b	1	0.1	0.1
x ₁	Number of positive words	3	2.5	7.5
x ₂	Number of negative words	2	-5.0	-10.0
x ₃	1 if "no" in text; 0 otherwise	1	-1.2	-1.2
×4	Number of 1st & 2nd person pronouns	3	0.5 50	1.5
x ₅	1 if "!" in text; 0 otherwise	0	2.0	0
x ₆	In of word/token count	4.19	0.7	2.933

➔ Effect of Logistic Regression classifier

- Classifiers over-emphasizes the importance of pronouns
 - \rightarrow large value for θ_4 (compared to other θ_i)
- Unseen negative review with many pronouns will most likely be misclassified

Overfitting — **Basic Intuition**

- Overfitting Visualized using curve fitting
 - Task: Find a polynomial for degree *p* that best fit the data points



- Polynomial of degree 1 just a line
- Not capable to fit non-linear data

Good fit

- Model captures the overall trend
- Probably good fit for unseen data

- Overfitting
- Model has too much capacity to exactly fit individual data points
- Probably bad fit for unseen data

Regularization

- Observation
 - Model "too powerful" \Leftrightarrow (very) large θ values

- → **Regularization**: Penalize large θ values
 - Extend loss function by penalty term
 - For example, for Cross-Entropy loss

Quick Quiz: What do the indices m and *n* stand for in the equations here?

 λ : Regularization Parameter to control the "strength of the regularization"

 $L = -\frac{1}{m} \sum_{j=1}^{m} \left[y^{(j)} \log \sigma \left(\theta^{\top} x^{(j)} \right) + \left(1 - y^{(j)} \right) \log \left(1 - \sigma \left(\theta^{\top} x^{(j)} \right) \right) \right] + \lambda \sum_{i=1}^{n} \theta_i^2$

$$L = -\frac{1}{m} \sum_{j=1}^{m} \left[y^{(j)} \log \sigma \left(\theta^{\top} x^{(j)} \right) + \left(1 - y^{(j)} \right) \log \left(1 - \sigma \left(\theta^{\top} x^{(j)} \right) \right) \right] + \lambda \sum_{i=1}^{n} |\theta_i|$$

L2 Regularization ("Ridge Regression")

L1 Regularization ("Lasso Regression")

New Loss → New Gradient

- Since we change *L*, the gradient $\nabla_{\theta} L = \frac{\partial L}{\partial \theta}$ also changes
 - No big deal, regularization is just an added term
 - For example, for L2 Regularization (Ridge Regression)

$$\frac{\partial L_{CE}}{\partial \theta} = \frac{1}{m} X^{\top} \left[\sigma \left(X \theta \right) - y \right] + \lambda \frac{2}{n} \theta$$

No changes to Gradient Descent Algorithms

In-Lecture Activity (4 mins)



Α

It's impossible to overfit given a dataset with only 1 feature

Β

Scaling the data will change the values for θ

Which of the statements regarding Logistic Regression is **True**?

С

Gradient Descent can get stuck in local minimum

D

Regularization can improve the training loss/error

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression
- Towards Neural Networks
 - Motivation: XOR Problem
 - Basic Neural Network Architecture

Binary LR → Multiclass LR

- Multiclass LR: Classification beyond 2 classes
 - Let's assume we have C classes: c = 1..C
 - Separate weights θ_c for each classes $c \rightarrow C$ output probabilities

Binary Logistic Regression

Multiclass Logistic Regression

$$f_{mystery} \rightarrow Softmax$$

- Softmax function
 - Converts any vector of scores into a vector of probabilities

$$P(y = c | x) = \frac{\exp(\theta_c^\top x)}{\sum_{i=1}^C \exp(\theta_i^\top x)}$$

$$\begin{bmatrix} P(y=1|x) \\ P(y=2|x) \\ \dots \\ P(y=C|x) \end{bmatrix} = \frac{1}{\sum_{i=1}^{C} \exp(\theta_i^\top x)} \begin{bmatrix} \exp(\theta_1^\top x) \\ \exp(\theta_2^\top x) \\ \dots \\ \exp(\theta_C^\top x) \end{bmatrix}$$

Example

• Example with 4 classes and 3 input features

Weight matrix
$$\theta$$
 x
 $\theta^{\top}x$
 \hat{y}
 $\theta_1 \begin{bmatrix} 0.55 & 0.71 & 0.29 \\ 0.51 & 0.89 & 0.90 \\ 0.3 & 0.21 & 0.05 \\ 0.44 & 0.03 & 0.46 \end{bmatrix} \cdot \begin{bmatrix} -0.4 \\ 0.2 \\ 0.3 \end{bmatrix} = \begin{bmatrix} 0.009 \\ 0.244 \\ 0.005 \\ -0.032 \end{bmatrix}$
 Softmax
 $\begin{bmatrix} 0.238 \\ 0.296 \\ 0.237 \\ 0.229 \end{bmatrix} \hat{y}_1$
 \hat{y}_1
 \hat{y}_2
 \hat{y}_2
 \hat{y}_2
 \hat{y}_3
 θ_4
 $0.03 & 0.46 \end{bmatrix} \cdot \begin{bmatrix} -0.4 \\ 0.2 \\ 0.3 \end{bmatrix} = \begin{bmatrix} 0.009 \\ 0.244 \\ 0.005 \\ -0.032 \end{bmatrix}$
 Softmax
 $\begin{bmatrix} 0.238 \\ 0.296 \\ 0.237 \\ 0.229 \end{bmatrix} \hat{y}_4$

Cross-Entropy Loss

Cross-Entropy Loss for Binary Logistic Regression

$$L_{CE}(\hat{y}, y) = -[y \log \hat{y} + (1 - y) \log (1 - \hat{y})]$$

Generalized Cross-Entropy Loss for Multiclass Logistic Regression



New gradient $\nabla_{\theta} L_{CE}$ but beyond the scope here.

Writing Systems of the World

Word

Logosyllabaries – each symbol is a syllable or an entire concept (e.g., 中文)

Syllable

Sound

Syllabaries – each syllable (consonant + vowel) a symbol (ひらがな)

Abugidas – each consonant is a base symbol, vowels permute them (தமிழ்)

<mark>Alphabets</mark> – each sound a symbol (русскийz)

Abjads – each consonant a symbol, vowels inferred from context (אָברִית)



Credits: <u>UsefulCharts.com</u>

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

Pre-Lecture Activity from Last Week

• Assigned Task

Post a 1–2 sentence answer to the following question into your Tutorial Group's discussions (you will find the thread on Canvas > Discussions)

"What is a common myth about neural networks?"

Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers, but his won't help you learn better

Pre-Lecture Activity from Last Week



A common misconception of neural networks is that they are models of the brain, despite what their name alludes to. Neural network resembles a graph while the brain resembles a structured network.

Search



transparent.

decisions

Neural Networks are always thought of mimicking a human

brain and are capable of learning like humans which isn't true since we can learn from a very small amount of data vs neural networks which requires terabytes of data

A common myth about neural networks is that they are

black boxes - meaning their inner workings and decision-

making processes are entirely uninterpretable. However, there's ongoing research and development of methods, (eg.

https://prometheuz.de/en/ai/are-neural-networks-black-

boxes/#can-we-understand-how-a-neural-network-makes-

sensitivity analysis), aimed at enhancing their

interpretability and making their operations more

picture from TuringFinance.com

CB

Do neural networks really work like neurons?



Yariy Adan · Follow Published in The Startup · 11 min read · Sep 29, 2018

```
Q 4
                                                                  ()
699) 1K
```

Artificial Neural Network and Machine Learning have become hot topics in the popular media. The idea of intelligent machines captivates the imagination of many, and especially how they would compare to humans. Specifically, one fundamental question that seems to come up frequently is about the underlaying mechanisms of intelligence - do these artificial neural networks really work like the neurons in our brain?

Tl; Dr:

No. While the high level and conceptual thinking of ANNs (artificial neural networks) is inspired by neurons and neural networks in the brain, the ML implementation of these concepts has diverged significantly from how the brain works. Moreover, as the field of ML progressed over the years, and



Not sure if this counts as a myth, but when I first started out with neural networks (and machine learning in general), I did not think data cleaning/preprocessing was that important.

However, the quality and reliability of neural networks heavily depend on thorough data cleaning/preprocessing. This process is crucial for extracting useful and unbiased features.



Neural networks can and should be applied to any and all ML problems.

Biological Inspiration — Neuron





→ Logistic Regression (crudely) a biological neuron

Logistic Regression — Limitations

- Logistic Regression is a linear model
 - Limited to linear combination of features (and a non-linear mapping to a probability)
 - Limited to linear decision boundaries (i.e., lines, planes, hyperplanes)
- What if we want or need to represent non-linear relationships between features? We can't!

➔ Scale up: "Stacked" Logistic Regression

- Feed input into multiple neurons (i.e., LR units)
- Use output of neurons as input for other neurons



XOR

x ₁	x ₂	OR	AND	NAND	XOR
0	0	0	0	1	0
0	1	1	0	1	1
1	0	1	0	1	1
1	1	1	1	0	0

Follow along yourself! https://www.desmos.com /calculator/waert4utde


XOR

- Learning OR, AND, and NAND
 - Finding correct weights simply by "looking hard" (the weights are not unique; there are many ways to set θ)
 - The activation function is the Step Function, not Sigmoid (strictly speaking, this makes it a Perceptron not a Linear Regression unit)

 $f_{step} = \begin{cases} 1 & , \text{if } \theta^{\top} x > 0 \\ 0 & , \text{otherwise} \end{cases}$



XOR

X ₁	x ₂	OR	AND	NAND	XOR
0	0	0	0	1	0
0	1	1	0	1	1
1	0	1	0	1	1
1	1	1	1	0	0

- Deriving XOR from simple classifiers
 - Note: this is not only way to do it, just convenient



→ Cool, we know how to do ORs, ANDs and NANDs!

XOR

- Modeling XOR by "stacking" LR units → Neural Network (NN)
 - More specifically, a Feedforward NN (i.e., network contains no loops)



Network Capacity — Intuition

Quick quiz: Is there any harm in having too many neurons?



Note: The activation function is the Sigmoid, hence the smooth decision boundaries

Outline

Generative vs. Discriminative Classifiers

Logistic Regression

- Setup as Probabilistic Classifier
- Cross-Entropy Loss Function
- Gradient Descent
- Overfitting & Regularization
- Multiclass Logistic Regression

• Towards Neural Networks

- Motivation: XOR Problem
- Basic Neural Network Architecture

A Neural Network (Feedforward NN)

• Example: L-layer Feedforward Neural Network (here: L = 4)





Neural Network — Activations



Neural Network — Activations

- Layer-wise computations
 - Let $x^{[l]}$ be the output of layer l
 - $x^{[0]} = x$ initial input
 - $x^{[L]} = h(x)$ final output
- Vectorized form
 - Calculate $x^{[l]}$ in practice "in one go"
 - Everything becomes matrix* operations
 - GPUs: hardware-supported processing of matrix operations (+ parallelism)

$$x_i^{[l]} = a_i^{[l]} = g\left(\sum_{j=0}^{d^{[l-1]}} \theta_{i,j}^{[l]} x_j^{[l-1]}\right)$$

$$= g \left(\begin{bmatrix} \theta_i^{[l]} \end{bmatrix}^\top \cdot x^{[l-1]} \right)$$

$$\uparrow$$
Weight vector $\theta_i^{[l]} \in \mathbb{R}^{d^{[l-1]}}$

$$\begin{aligned} x^{[l]} &= a^{[l]} = g \left(\theta^{[l]}_{\uparrow} x^{[l-1]} \right) \\ \uparrow \end{aligned} \\ \text{Weight matrix } \theta^{[l]} \in \mathbb{R}^{d^{[l]} \times d^{[l-1]}} \end{aligned}$$

Neural Network — Activation Functions

- Wide range of activation functions
- Activations functions for hidden layers
 - Do not need to have a probabilistic interpretation
 - Only requirement: non-linear function!
 - Examples:





Neural Network — Activation Functions

• Activations functions for output layers

- Choice of activation function depending on task (mainly: classification or regression)
- Examples:

Linear function for regression tasks



Sigmoid function for classification tasks



Example

Input x Hidden h Output y



$$h = g_h(\theta_h x) \text{, with } \theta_h \in \mathbb{R}^{4 \times 3}$$
$$y = g_y(\theta_y h) \text{, with } \theta_y \in \mathbb{R}^{2 \times 4}$$

 $g_h, \; g_y$: suitable activation functions

Neural Networks



From Logistic Regression to (Deep) Neural Networks

- Fundamentally, nothing new here:
 - A neural network is a function $h_{\theta}(x)$
 - Define a loss function $L = L(y, \hat{y}) = L(y, h_{\theta}(x))$
 - Perform Gradient Descent to minimize L
- Difference: increased complexity
 - $\blacksquare \ h_{\theta}(x)$ and thus $L(y,h_{\theta}(x))$ are much more complex functions
 - Calculation of $\frac{\partial L}{\partial \theta}$ much more challenging → backpropagation
 - *L* is no longer a convex function \rightarrow local minima \rightarrow training more challenging
 - Overfitting becomes a bigger issue

Summary

• Linear model: Logistic Regression

- Very important probabilistic classifier
- Discriminative classifier → linear decision boundaries
- Core unit of neural networks
- "Stacked" Logistic Regression → Neural Network
 - Neuron = Linear Regression unit
 - Non-convex loss function → global minimum vs. local minima
 - Higher risk of overfitting → regularization crucial (but also other methods)

Outlook for Next Week: Embeddings and Ethics

Pre-Lecture Activity for Next Week

• Assigned Task (due before Feb 23)

 Post a 1-2 sentence answer to the following question in your Tutorial Group's discussion (you will find the thread on Canvas > Discussions)

"What do we mean by sparse or dense vectors? Are documents characterised by tf-idf sparse or dense?"

Read some blog posts or online articles, and cite them with the links in your answer

Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers but this won't help you learn better