

Text classification

He He



NEW YORK UNIVERSITY

September 6, 2023

Table of Contents

Course overview

- Logistics

- A brief history of NLP/AI

- Challenges in NLP

Supervised learning basics

- Generalization

- Loss functions

- Optimization

Text classification

- Generative models: naive Bayes

- Discriminative models: logistic regression

Logistics

- Instructor: He He
- Section leaders: Lavender Jiang, Xiang Pan, Weizhe Yuan, Yilun Kuang (see the Staff page)
- Graders: Lavender Jiang, Xiang Pan, Yilun Kuang, Danielle Rothermel
- Best way to communicate with us: **Campuswire** (link and code on Brightspace).
- Office hours will be in-person or on Zoom (see the Staff page).
- Let us know if you have accessibility needs.

What this course is (not) about

- It's not about specific NLP applications (QA, dialogue etc.)
 - Unified approaches to various NLP problems
 - Hands-on experience in building NLP systems (e.g., machine translation) through assignments and the course project

What this course is (not) about

- It's not about specific NLP applications (QA, dialogue etc.)
 - Unified approaches to various NLP problems
 - Hands-on experience in building NLP systems (e.g., machine translation) through assignments and the course project
- It's not about fundamental machine learning and deep learning
 - Focus on unique challenges in language data
 - Formalize NLP tasks as statistical learning problems
 - Focus on modern techniques (e.g., word vectors, neural networks, large language models)

What we expect you to know

- **Linear algebra:** vector space, vector norm, dot product, gradient etc.
- **Probability and statistics:** conditional probability, expectation, Bayes rule etc.
- **Basic machine learning:** loss function, gradient descent, logistic regression etc.
- **Programming:** read and write Python code, use Numpy, HPC, and deep learning libraries (Pytorch, Huggingface etc.)

Course project

An important component of the course (more on this later)

- Related to NLP (doesn't have to be in the scope of this course)
- New algorithms or models for existing problems
- Applications of NLP or ML techniques to a problem
- Analysis of well-known approaches that leads to new insight
- ML Reproducibility Challenge 2021 (<https://paperswithcode.com/rc2021>)

Products powered by NLP technologies

The screenshot shows the Google Translate web interface. At the top, there are tabs for 'Text' and 'Documents'. Below, a language bar shows 'ENGLISH - DETECTED' on the left and 'CHINESE (SIMPLIFIED)' on the right. The main area contains two columns of text. The left column has English text about researchers inoculating themselves. The right column shows the Chinese translation of the same text. At the bottom of each column are speaker icons and a character count (507/5000).

Text input: These researchers have inoculated themselves — and, sometimes, friends and family — bypassing the rigorous tests required for conventional vaccines and raising fears of potential side effects.

Translation: 这些研究人员已经为自己（有时甚至是朋友和家人）接种了疫苗，从而绕过了常规疫苗所需的严格测试，并增加了对潜在副作用的担忧。

Send feedback

The screenshot shows a Google search result for the query 'what are countries adjacent to romania'. The search bar is at the top with the query and a magnifying glass icon. Below the search bar are filters for 'All', 'Maps', 'News', 'Images', 'Shopping', and 'Tools'. The search results show 'About 579,000,000 results (0.65 seconds)'. The first result is from Britannica, titled 'Land. Romania is bounded by Ukraine to the north, Moldova to the northeast, the Black Sea to the southeast, Bulgaria to the south, Serbia to the southwest, and Hungary to the west. 4 days ago'. To the right of the text is a map of Romania showing its borders with neighboring countries. At the bottom, there is a link to 'Romania | History, Map, & Facts | Britannica'.

Search query: what are countries adjacent to romania

Results: About 579,000,000 results (0.65 seconds)

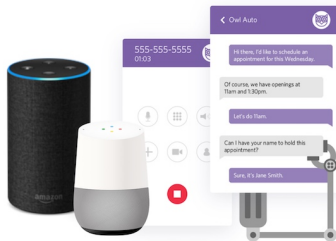
Result snippet: Land. Romania is bounded by **Ukraine** to the north, **Moldova** to the northeast, the Black Sea to the southeast, **Bulgaria** to the south, **Serbia** to the southwest, and **Hungary** to the west. 4 days ago

About 579,000,000 results (0.65 seconds)

The screenshot shows a snippet from a Britannica article about Romania's borders. It includes a map of Romania and a list of neighboring countries: Ukraine to the north, Moldova to the northeast, the Black Sea to the southeast, Bulgaria to the south, Serbia to the southwest, and Hungary to the west. The snippet is dated '4 days ago'.

Land. Romania is bounded by **Ukraine** to the north, **Moldova** to the northeast, the Black Sea to the southeast, **Bulgaria** to the south, **Serbia** to the southwest, and **Hungary** to the west. 4 days ago

Source: www.britannica.com > ... > Countries of the World > Romania | History, Map, & Facts | Britannica



A single natural language interface for everything



what's the best way to learn NLP?



There is no one-size-fits-all answer to this question as the best way to learn NLP (Natural Language Processing) will depend on your background and goals. However, some ways to get started with learning NLP include:



1. Taking online courses: There are many online courses available on NLP, such as those offered by Coursera, edX, and Udemy.
2. Reading books: There are many books available on NLP, both for beginners and advanced learners. Some popular books include "Speech and Language Processing" by Daniel Jurafsky and James H. Martin, and "Natural Language Processing with Python" by Steven Bird, Ewan Klein, and Edward Loper.
3. Practicing with tutorials and projects: There are many tutorials and open-source projects available online that can help you learn NLP by implementing and experimenting with different algorithms and techniques.
4. Joining a community: Joining a community of NLP enthusiasts, such as those on Reddit, Slack, or LinkedIn, can provide you with support, resources, and opportunities to collaborate on projects.
5. Participate in NLP Challenges: Participating in NLP challenges such as those hosted by Kaggle or other platforms can be a great way to learn and improve your skills.

A single natural language interface for everything



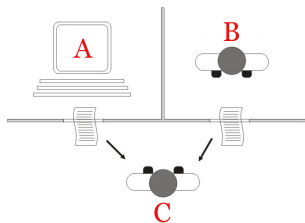
Write an excel formula to compute the average of all scores larger than 80 in column C1



`=AVERAGEIF(C1,">80")`

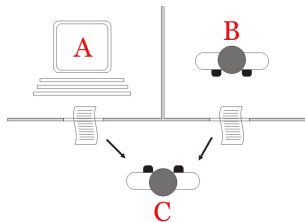


Language is at the core of AI: the imitation game



*"I believe that in about **fifty years**' time it will be possible to programme computers, with a **storage capacity of about 10^9** , to make them play the imitation game so well that an average interrogator will not have more than 70 percent chance of making the right identification after five minutes of questioning."* Turing (1950)

Language is at the core of AI: the imitation game



*"I believe that in about **fifty years**' time it will be possible to programme computers, with a **storage capacity of about 10^9** , to make them play the imitation game so well that an average interrogator will not have more than 70 percent chance of making the right identification after five minutes of questioning."* Turing (1950)



Is human likeness the ultimate goal?

ELIZA

- Built by Joseph Weizenbaum at MIT in 1964 to demonstrate the *superficiality* of human-machine communication.
- Surprisingly, people were convinced that ELIZA had human intelligence.

Human: Well, my boyfriend made me come here.

ELIZA: Your boyfriend made you come here?

Human: He says I'm depressed much of the time.

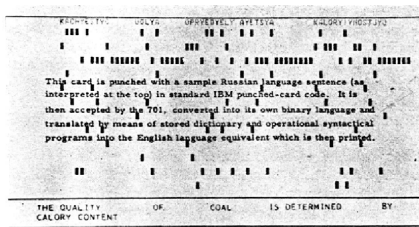
ELIZA: I am sorry to hear you are depressed.

Human: It's true. I'm unhappy.

ELIZA: Do you think coming here will help you not to be unhappy?

Early rule-based systems: the Georgetown-IBM experiment

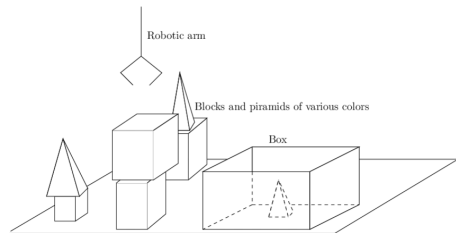
- The Russian-English machine translation program:



- A vocabulary of 250 words
- Using 6 grammar rules, e.g.,
If first code is 110, is third code associated with preceding complete word equal to 21? If so, reverse order of appearance of words in output (i.e., word carrying 21 should follow that carrying 110)---otherwise, retain order.

Approaching AI as a whole: SHRDLU

- Built by Terry Winograd at MIT in 1968.
- A person instructs the computer to build structures in a blocks world.
- **Require many skills:** grounding, coreference resolution, question answering, semantic parsing etc.



Person: Pick up a big red block.

Computer: OK.

Person: Grasp the pyramid.

Computer: I don't understand which pyramid you mean.

Person: (changing their mind) Find a block which is taller than the one you are holding and put it into the box.

Computer: By "it", i assume you mean the block which is taller than the one i am holding.

Limitations of early systems

- Optimism in the 50's and 60's: working on tasks that are too complex at that time

"Within the very near future—much less than twenty-five years—we shall have the technical capability of substituting machines for any and all human functions in organizations."

- Disappointing results due to
 - **Limited computation:** hardware has limited speed and memory
 - **Combinatorial explosion:** algorithms are intractable in realistic settings
 - **Underestimated complexity:** ambiguity, commonsense knowledge etc.

The rise of statistical learning in the 80's

- Notable progress in MT from IBM (neglected knowledge of linguistics).
- HMMs widely used for speech recognition.
"Every time I fire a linguist, the performance of the speech recognizer goes up."—Frederick Jelinek.
- The paradigm shift: expert knowledge + rules → data + features
- Statistical learning is the main driving force of NLP today.

The deep learning tsunami

- Before deep learning (around 2015), NLP is mostly about structured prediction and feature engineering.
- Neural networks can automatically learn good features/representations for a task
- The paradigm shift: features → network architectures + embeddings
- Almost all NLP models are neural networks nowadays.

Models and data keep getting larger

- Since around 2018, Transformer-based pretrained models have become the standard.
- Pre-training on large data provides useful representations for many downstream tasks.
- The paradigm shift: [architecture design](#) → [transfer learning \(fine-tuning\)](#)
- More recently, a single natural language interface for all tasks (e.g., ChatGPT by OpenAI).
- The paradigm shift: [transfer learning](#) → [instructing / prompting](#)

Why is language hard?

Why is language hard?

- **Discrete**

- How to define metrics?

I work **at** NYU. vs I work **for** NYU.

This is good. vs This is **actually** good.

- How to define transformations?

The food is okay.

→ The food is awesome!

They made a brief return to Cambridge to
drop the book.

→ They returned.

- In general, hard to represent text as mathematical objects.

Why is language hard?

- **Compositional**

- The whole is built from parts (chars, words, sentences, paragraphs, documents...)
- How to generalize when we don't see all possible combinations?
- An example from [Lake et al., 2018]

Vocabulary:

{jump, walk, turn, once, twice, left, right, before, after, and}

Sentences:

jump

jump left

jump left and walk right

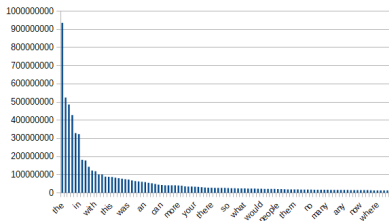
jump left after walk right once before turn left twice

...

Why is language hard?

- **Sparse**

- How to handle the long tail?
- Zipf's law: word frequency $\propto \frac{1}{\text{rank}}$



- Many linguistic phenomena follow Zipf's law
BoA's financial assistant Erica:
*The bank "learned [that] there are over 2,000 different ways to ask us to move money."*¹

¹<https://www.aiqudo.com/2019/06/28/voice-success-story-erica-bank-america/>

Why is language hard?

- **Ambiguous**

- How to interpret meaning in context?

Bass: fish? guitar? frequency? (word sense disambiguation)

I shot an elephant in my pajamas: who is in the pajamas? (PP attachment)

The spirit is willing but the flesh is weak.
→ The vodka is strong but the meat is rotten.

Table of Contents

Course overview

- Logistics

- A brief history of NLP/AI

- Challenges in NLP

Supervised learning basics

- Generalization

- Loss functions

- Optimization

Text classification

- Generative models: naive Bayes

- Discriminative models: logistic regression

Rule-based approach

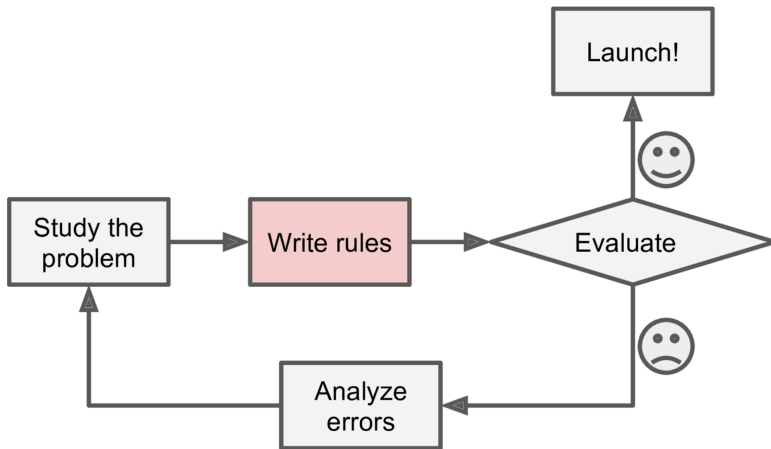


Figure: Fig 1-1 from *Hands-On Machine Learning with Scikit-Learn and TensorFlow* by Aurelien Geron (2017).

Machine learning approach

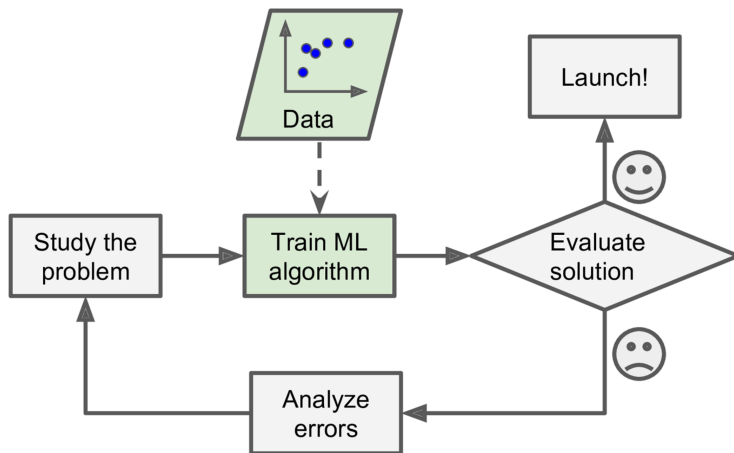


Figure: Fig 1-2 from *Hands-On Machine Learning with Scikit-Learn and TensorFlow* by Aurelien Geron (2017).

Example: spam filter

- Rules
 - Contains “Viagra”
 - Contains “Rolex”
 - Subject line is all caps
 - ...
- Learning from data
 1. Collect emails labeled as spam or non-spam
 2. Design features, e.g., first word of the subject, nouns in the main text
 3. Learn a binary classifier

Example: spam filter

- Rules

- Contains "Viagra"

- Contains "Rolex"

- Subject line is all caps

- ...

- Learning from data

1. Collect emails labeled as spam or non-spam

2. Design features, e.g., first word of the subject, nouns in the main text

3. Learn a binary classifier



Pros and cons of each approach?

Key challenges in machine learning

- Availability of large amounts of (annotated) data
 - Data collection: scraping, crowdsourcing, expert annotation
 - Quality control: data quality can have large impact on the final model (garbage in garbage out)
 - Don't take it for granted: always check the data source!

Key challenges in machine learning

- Availability of large amounts of (annotated) data
 - Data collection: scraping, crowdsourcing, expert annotation
 - Quality control: data quality can have large impact on the final model (garbage in garbage out)
 - Don't take it for granted: always check the data source!



How would you collect a dataset for the spam filtering task?

Key challenges in machine learning

- **Generalize** to unseen samples
 - We want to build a model: $h: \mathcal{X}$ (input space) $\rightarrow \mathcal{Y}$ (output space)
 - It is easy to achieve high accuracy on the training set.
 - But we want the model to perform well on unseen data, too.
 - What should be the learning objective?

Empirical risk minimization (ERM)

- Assume a data generating distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ (e.g., spam writers and non-spam writers)

Empirical risk minimization (ERM)

- Assume a data generating distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ (e.g., spam writers and non-spam writers)
- We have access to a training set: m samples from \mathcal{D} $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$

Empirical risk minimization (ERM)

- Assume a data generating distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ (e.g., spam writers and non-spam writers)
- We have access to a training set: m samples from \mathcal{D} $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$
- We can measure the goodness of a prediction $h(x)$ by comparing it against the groundtruth y using some **loss function**

Empirical risk minimization (ERM)

- Assume a data generating distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ (e.g., spam writers and non-spam writers)
- We have access to a training set: m samples from \mathcal{D} $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$
- We can measure the goodness of a prediction $h(x)$ by comparing it against the groundtruth y using some **loss function**
- Our goal is to minimize the **expected loss** over \mathcal{D} (**risk**):

$$\text{minimize } \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{error}(h, x, y)] ,$$

but it **cannot be computed** (why?).

Empirical risk minimization (ERM)

- Assume a data generating distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ (e.g., spam writers and non-spam writers)
- We have access to a training set: m samples from \mathcal{D} $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$
- We can measure the goodness of a prediction $h(x)$ by comparing it against the groundtruth y using some **loss function**
- Our goal is to minimize the **expected loss** over \mathcal{D} (**risk**):

$$\text{minimize } \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{error}(h, x, y)] ,$$

but it **cannot be computed** (why?).

- Instead, we minimize the **average loss on the training set** (**empirical risk**)

$$\text{minimize } \frac{1}{m} \sum_{i=1}^m \text{error}(h, x^{(i)}, y^{(i)})$$

Empirical risk minimization (ERM)

- Assume a data generating distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ (e.g., spam writers and non-spam writers)
- We have access to a training set: m samples from \mathcal{D} $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$
- We can measure the goodness of a prediction $h(x)$ by comparing it against the groundtruth y using some **loss function**
- Our goal is to minimize the **expected loss** over \mathcal{D} (**risk**):

$$\text{minimize } \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{error}(h, x, y)] ,$$

but it **cannot be computed** (why?).

- Instead, we minimize the **average loss on the training set** (**empirical risk**)

$$\text{minimize } \frac{1}{m} \sum_{i=1}^m \text{error}(h, x^{(i)}, y^{(i)})$$

- **Key question:** does small empirical risk imply small risk?

Overfitting vs underfitting

- Trivial solution to (unconstrained) ERM: **memorize** the data points
- Need to extrapolate information from one part of the input space to unobserved parts!
- Solution: constrain the prediction function to a subset, i.e. a **hypothesis space** $h \in \mathcal{H}$.

Overfitting vs underfitting

- Trivial solution to (unconstrained) ERM: **memorize** the data points
- Need to extrapolate information from one part of the input space to unobserved parts!
- Solution: constrain the prediction function to a subset, i.e. a **hypothesis space** $h \in \mathcal{H}$.
- Trade-off between complexity of \mathcal{H} and generalization
- Question for us: **how to choose a good \mathcal{H} for certain domains**

Summary

1. Obtain training data $D_{\text{train}} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$.
2. Choose a loss function L and a hypothesis class \mathcal{H} (domain knowledge).
3. Learn a predictor by minimizing the empirical risk (optimization).

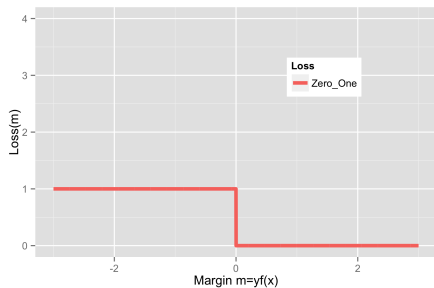
Setup

- Task: binary classification $y \in \{+1, -1\}$
- Model: $f_w: \mathcal{X} \rightarrow \mathbf{R}$ parametrized by $w \in \mathbf{R}^d$
 - Output a score for each example
- Prediction: $\text{sign}(f_w(x))$
 - Positive scores are mapped to the positive class
- Loss functions: quantify the goodness of the model output $f_w(x)$ given y

Zero-one loss

First idea: check if the prediction is the same as the label

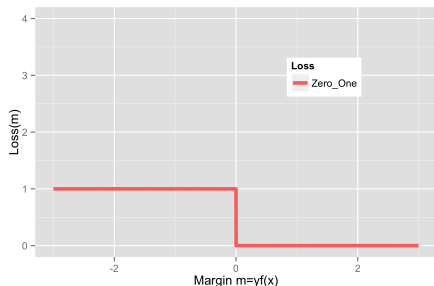
$$L(x, y, f_w) = \mathbb{I}[\text{sign}(f_w(x)) = y] = \mathbb{I}\left[\underbrace{yf_w(x)}_{\text{margin}} \leq 0\right] \quad (1)$$



Zero-one loss

First idea: check if the prediction is the same as the label

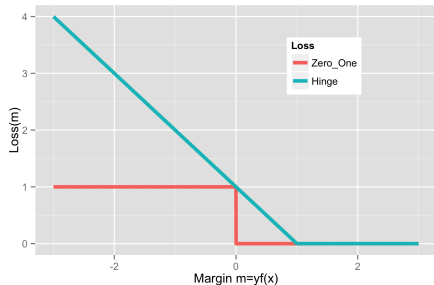
$$L(x, y, f_w) = \mathbb{I}[\text{sign}(f_w(x)) = y] = \mathbb{I}\left[\underbrace{yf_w(x)}_{\text{margin}} \leq 0\right] \quad (1)$$



Problem: not differentiable

Hinge loss

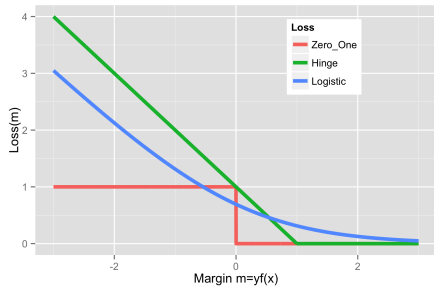
$$L(x, y, f_w) = \max(1 - yf_w(x), 0)$$



- A (sub)differentiable upperbound of the zero-one loss
- Not differentiable at margin = 1 (use subgradients)

Logistic loss

$$L(x, y, f_w) = \log(1 + e^{-yf_w(x)})$$



- Differentiable
- Always wants more margin (loss is never 0)

Summary

1. Obtain training data $D_{\text{train}} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$.
2. Choose a loss function L and a hypothesis class \mathcal{H} (domain knowledge).
3. Learn a predictor by minimizing the empirical risk (optimization).

Gradient descent

- The gradient of a function F at a point $w \in \mathbb{R}^d$ is the direction of fastest increase in the function value
- To minimize $F(w)$, move in the opposite direction

$$w \leftarrow w - \eta \nabla_w F(w)$$

- Converge to a local minimum (also global minimum if $F(w)$ is **convex**) with carefully chosen step sizes η

Stochastic gradient descent

- **Gradient descent (GD)** for ERM

$$w \leftarrow w - \eta \nabla_w \underbrace{\sum_{i=1}^n L(x^{(i)}, y^{(i)}, w)}_{\text{training set loss}}$$

Stochastic gradient descent

- **Gradient descent (GD)** for ERM

$$w \leftarrow w - \eta \nabla_w \underbrace{\sum_{i=1}^n L(x^{(i)}, y^{(i)}, w)}_{\text{training set loss}}$$

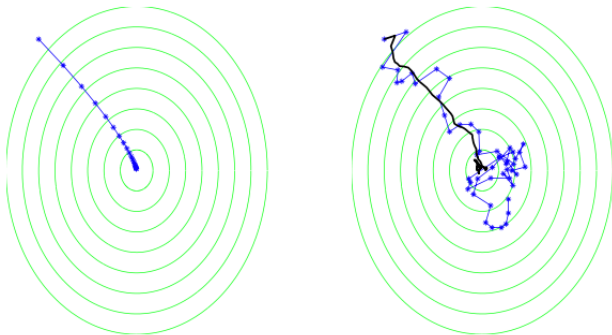
- **Stochastic gradient descent (SGD)**: take **noisy but faster** updates

For each $(x, y) \in D_{\text{train}}$:

$$w \leftarrow w - \eta \nabla_w \underbrace{L(x, y, f_w)}_{\text{example loss}}$$

GD vs SGD

Figure: Minimize $1.25(x + 6)^2 + (y - 8)^2$. Example from “Understanding Machine Learning: From Theory to Algorithms”



SGD step is noisier as it gets closer to the optimum; need to reduce step size gradually.

SGD summary

- Each update is efficient in both time and space
- Can be slow to converge
- Popular in large-scale ML, including non-convex problems
- In practice,
 - Randomly sample examples.
 - Fixed or diminishing step sizes, e.g. $1/t$, $1/\sqrt{t}$.
 - Stop when objective does not improve.
- Our main optimization technique

Summary

- Choose hypothesis class based on domain knowledge
- Learning algorithm: empirical risk minimization
- Optimization: stochastic gradient descent

Table of Contents

Course overview

- Logistics

- A brief history of NLP/AI

- Challenges in NLP

Supervised learning basics

- Generalization

- Loss functions

- Optimization

Text classification

- Generative models: naive Bayes

- Discriminative models: logistic regression

Text classification

- Input: text (sentence, paragraph, document)
- Predict the **category or property** of the input text
 - Sentiment classification: Is the review positive or negative?
 - Spam detection: Is the email/message spam or not?
 - Hate speech detection: Is the tweet/post toxic or not?
 - Stance classification: Is the opinion liberal or conservative?

Text classification

- Input: text (sentence, paragraph, document)
- Predict the **category or property** of the input text
 - Sentiment classification: Is the review positive or negative?
 - Spam detection: Is the email/message spam or not?
 - Hate speech detection: Is the tweet/post toxic or not?
 - Stance classification: Is the opinion liberal or conservative?
- Predict the **relation** of two pieces of text
 - Textual entailment (HW1): does the premise entail the hypothesis?
Premise: The dogs are running in the park.
Hypothesis: There are dogs in the park.
 - Paraphrase detection: are the two sentences paraphrases?
Sentence 1: The dogs are in the park.
Sentence 2: There are dogs in the park.

Intuition

- **Example:** sentiment classification for movie reviews

Action. Comedy. Suspense. This movie has it all. The Plot goes that 4 would be professional thieves are invited to take part in a heist in a small town in Montana. every type of crime movie archetype character is here. Frank, the master mind. Carlos, the weapons expert. Max, the explosives expert. Nick, the safe cracker and Ray, the car man. Our 4 characters meet up at the train station and from the beginning none of them like or trust one another. Added to the mix is the fact that Frank is gone and they are not sure why they have called together. Now Frank is being taken back to New Jersey by the 2 detectives but soon escapes on foot and tries to make his way back to the guys who are having all sorts of problems of their own. Truly a great film loaded with laughs and great acting. Just an overall good movie for anyone looking for a laugh or something a little different

Intuition

- **Example:** sentiment classification for movie reviews

Action. Comedy. Suspense. This movie has it all. The Plot goes that 4 would be professional thieves are invited to take part in a heist in a small town in Montana. every type of crime movie archetype character is here. Frank, the master mind. Carlos, the weapons expert. Max, the explosives expert. Nick, the safe cracker and Ray, the car man. Our 4 characters meet up at the train station and from the beginning none of them like or trust one another. Added to the mix is the fact that Frank is gone and they are not sure why they have called together. Now Frank is being taken back to New Jersey by the 2 detectives but soon escapes on foot and tries to make his way back to the guys who are having all sorts of problems of their own. Truly a great film loaded with laughs and great acting. Just an overall good movie for anyone looking for a laugh or something a little different

- **Idea:** count the number of positive/negative words

Intuition

- **Example:** sentiment classification for movie reviews

Action. Comedy. Suspense. This movie has it all. The Plot goes that 4 would be professional thieves are invited to take part in a heist in a small town in Montana. every type of crime movie archetype character is here. Frank, the master mind. Carlos, the weapons expert. Max, the explosives expert. Nick, the safe cracker and Ray, the car man. Our 4 characters meet up at the train station and from the beginning none of them like or trust one another. Added to the mix is the fact that Frank is gone and they are not sure why they have called together. Now Frank is being taken back to New Jersey by the 2 detectives but soon escapes on foot and tries to make his way back to the guys who are having all sorts of problems of their own. Truly a great film loaded with laughs and great acting. Just an overall good movie for anyone looking for a laugh or something a little different

- **Idea:** count the number of positive/negative words
 - What is a “word”?
 - How do we know which are positive/negative?

Preprocessing: tokenization

Goal: Splitting a string of characters to a sequence of **tokens** $[x_1, \dots, x_n]$.

Language-specific solutions

- Regular expression: "I didn't watch the movie". \rightarrow ["I", "did", "n't", "watch", "the", "movie", "."]
 - Special cases: U.S., Ph.D. etc.
- Dictionary / sequence labeler: "我没有去看电影。" \rightarrow ["我", "没有", "去", "看", "电影", "。"]

Preprocessing: tokenization

Goal: Splitting a string of characters to a sequence of **tokens** $[x_1, \dots, x_n]$.

Language-specific solutions

- Regular expression: "I didn't watch the movie". \rightarrow ["I", "did", "n't", "watch", "the", "movie", "."]
 - Special cases: U.S., Ph.D. etc.
- Dictionary / sequence labeler: "我没有去看电影。" \rightarrow ["我", "没有", "去", "看", "电影", "。"]

General solutions: don't split by words

- Characters: ["u", "n", "a", "f", "f", "a", "b", "l", "e"]
- Subword (e.g., byte pair encoding): ["un", "aff", "able#"]

Classification: problem formulation

- **Input:** a sequence of tokens $x = (x_1, \dots, x_n)$ where $x_i \in \mathcal{V}$.
- **Output:** binary label $y \in \{0, 1\}$.
- **Probabilistic model:**

$$f(x) = \begin{cases} 1 & \text{if } p_{\theta}(y = 1 \mid x) > 0.5 \\ 0 & \text{otherwise} \end{cases},$$

where p_{θ} is a distribution parametrized by $\theta \in \Theta$.

- Modeling question: what's the parametric form of p_{θ} ?

Modeling $p(y \mid x)$

How to write a review:

1. Decide the sentiment by flipping a coin: $p(y)$
2. Generate word sequentially conditioned on the sentiment $p(x \mid y)$

Modeling $p(y \mid x)$

How to write a review:

1. Decide the sentiment by flipping a coin: $p(y)$
2. Generate word sequentially conditioned on the sentiment $p(x \mid y)$

$$p(y) = \quad (2)$$

$$p(x \mid y) = \quad (3)$$

$$= \quad (4)$$

Modeling $p(y \mid x)$

How to write a review:

1. Decide the sentiment by flipping a coin: $p(y)$
2. Generate word sequentially conditioned on the sentiment $p(x \mid y)$

$$p(y) = \text{Bernoulli}(\alpha) \tag{2}$$

$$p(x \mid y) = \tag{3}$$

$$= \tag{4}$$

Modeling $p(y \mid x)$

How to write a review:

1. Decide the sentiment by flipping a coin: $p(y)$
2. Generate word sequentially conditioned on the sentiment $p(x \mid y)$

$$p(y) = \text{Bernoulli}(\alpha) \quad (2)$$

$$p(x \mid y) = \prod_{i=1}^n p(x_i \mid y) \quad (\text{independent assumption}) \quad (3)$$

$$= \quad (4)$$

Modeling $p(y \mid x)$

How to write a review:

1. Decide the sentiment by flipping a coin: $p(y)$
2. Generate word sequentially conditioned on the sentiment $p(x \mid y)$

$$p(y) = \text{Bernoulli}(\alpha) \quad (2)$$

$$p(x \mid y) = \prod_{i=1}^n p(x_i \mid y) \quad (\text{independent assumption}) \quad (3)$$

$$= \prod_{i=1}^n \text{Categorical}(\underbrace{\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}}_{\text{sum to 1}}) \quad (4)$$

Modeling $p(y \mid x)$

How to write a review:

1. Decide the sentiment by flipping a coin: $p(y)$
2. Generate word sequentially conditioned on the sentiment $p(x \mid y)$

$$p(y) = \text{Bernoulli}(\alpha) \quad (2)$$

$$p(x \mid y) = \prod_{i=1}^n p(x_i \mid y) \quad (\text{independent assumption}) \quad (3)$$

$$= \prod_{i=1}^n \text{Categorical}(\underbrace{\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}}_{\text{sum to 1}}) \quad (4)$$

Bayes rule:

$$p(y \mid x) = \frac{p(x \mid y)p(y)}{p(x)} = \frac{p(x \mid y)p(y)}{\sum_{y \in \mathcal{Y}} p(x \mid y)p(y)}$$

Naive Bayes models

Naive Bayes assumption

The input features are **conditionally independent** given the label:

$$p(x \mid y) = \prod_{i=1}^n p(x_i \mid y) .$$

- A strong assumption, but works surprisingly well in practice.
- Note: $p(x_i \mid y)$ doesn't have to be a categorical distribution (e.g., Gaussian distribution)

Learning: maximum likelihood estimation

Task: estimate parameters θ of a distribution $p(y; \theta)$ given i.i.d. samples $D = (y_1, \dots, y_N)$ from the distribution.

Goal: find the parameters that make the observed data most probable.

Learning: maximum likelihood estimation

Task: estimate parameters θ of a distribution $p(y; \theta)$ given i.i.d. samples $D = (y_1, \dots, y_N)$ from the distribution.

Goal: find the parameters that make the observed data most probable.

Likelihood function of θ given D :

$$L(\theta; D) \stackrel{\text{def}}{=} p(D; \theta) = \prod_{i=1}^N p(y_i; \theta) .$$

Maximum (log-)likelihood estimator:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta; D) = \arg \max_{\theta \in \Theta} \sum_{i=1}^N \log p(y_i; \theta) \quad (5)$$

MLE and ERM

ERM:

$$\min \sum_{i=1}^N \ell(x^{(i)}, y^{(i)}, \theta)$$

MLE and ERM

ERM:

$$\min \sum_{i=1}^N \ell(x^{(i)}, y^{(i)}, \theta)$$

MLE:

$$\max \sum_{i=1}^N \log p(y^{(i)} \mid x^{(i)}; \theta)$$

MLE and ERM

ERM:

$$\min \sum_{i=1}^N \ell(x^{(i)}, y^{(i)}, \theta)$$

MLE:

$$\max \sum_{i=1}^N \log p(y^{(i)} \mid x^{(i)}; \theta)$$

What's the connection between MLE and ERM?

MLE is equivalent to ERM with the **negative log-likelihood** (NLL) loss function:

$$\ell_{\text{NLL}}(x^{(i)}, y^{(i)}, \theta) \stackrel{\text{def}}{=} -\log p(y^{(i)} \mid x^{(i)}; \theta)$$

MLE solution for our Naive Bayes model

$\text{count}(w, y) \stackrel{\text{def}}{=} \text{frequency of } w \text{ in documents with label } y$

$$p_{\text{MLE}}(w \mid y) = \frac{\text{count}(w, y)}{\sum_{w \in \mathcal{V}} \text{count}(w, y)}$$

= how often the word occur in positive/negative documents
= "positive/negative score of the word"

$$p_{\text{MLE}}(y = k) = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = k)}{N}$$

= fraction of positive/negative documents

Inference: make predictions using the model

Inference: $y = \arg \max_{y \in \mathcal{Y}} p_{\theta}(y \mid x)$

Inference: make predictions using the model

Inference: $y = \arg \max_{y \in \mathcal{Y}} p_{\theta}(y \mid x)$

Compare $p_{\theta}(y = 1 \mid x)$ and $p_{\theta}(y = 0 \mid x)$:

$$\frac{p_{\theta}(y = 1 \mid x)}{p_{\theta}(y = 0 \mid x)} = \frac{p_{\theta}(x \mid y = 1)p_{\theta}(y = 1)}{p_{\theta}(x \mid y = 0)p_{\theta}(y = 0)}$$

Inference: make predictions using the model

Inference: $y = \arg \max_{y \in \mathcal{Y}} p_{\theta}(y \mid x)$

Compare $p_{\theta}(y = 1 \mid x)$ and $p_{\theta}(y = 0 \mid x)$:

$$\frac{p_{\theta}(y = 1 \mid x)}{p_{\theta}(y = 0 \mid x)} = \frac{p_{\theta}(x \mid y = 1)p_{\theta}(y = 1)}{p_{\theta}(x \mid y = 0)p_{\theta}(y = 0)}$$

Assuming $p_{\theta}(y = 1) = p_{\theta}(y = 0)$, we only need to compare $p_{\theta}(x \mid y = 1)$ and $p_{\theta}(x \mid y = 0)$.

$$\text{score of class } k = \log p_{\theta}(x \mid y = k) = \sum_{i=1}^n \log p_{\theta}(x_i \mid y = k)$$

(Adding up positive/negative scores of each word)

Feature design

Naive Bayes doesn't have to use single words as features

- Lexicons, e.g., LIWC.
- Task-specific features, e.g., is the email subject all caps.
- Bytes and characters, e.g., used in language ID detection.

Summary of Naive Bayes models

- Modeling: the conditional independence assumption simplifies the problem
- Learning: MLE (or ERM with negative log-likelihood loss)
- Inference: very fast (adding up scores of each word)

Discriminative models

Idea: directly model the conditional distribution $p(y \mid x)$

Discriminative models

Idea: directly model the conditional distribution $p(y \mid x)$

	generative models	discriminative models
modeling	joint: $p(x, y)$	conditional: $p(y \mid x)$
assumption on y	yes	yes
assumption on x	yes	no
development	generative story	feature extractor

Model $p(y \mid x)$

How to model $p(y \mid x)$?

y is a Bernoulli variable:

$$p(y \mid x) = \alpha^y (1 - \alpha)^{(1-y)}$$

Model $p(y \mid x)$

How to model $p(y \mid x)$?

y is a Bernoulli variable:

$$p(y \mid x) = \alpha^y (1 - \alpha)^{(1-y)}$$

Bring in x :

$$p(y \mid x) = h(x)^y (1 - h(x))^{(1-y)} \quad h(x) \in [0, 1]$$

Model $p(y \mid x)$

How to model $p(y \mid x)$?

y is a Bernoulli variable:

$$p(y \mid x) = \alpha^y (1 - \alpha)^{(1-y)}$$

Bring in x :

$$p(y \mid x) = h(x)^y (1 - h(x))^{(1-y)} \quad h(x) \in [0, 1]$$

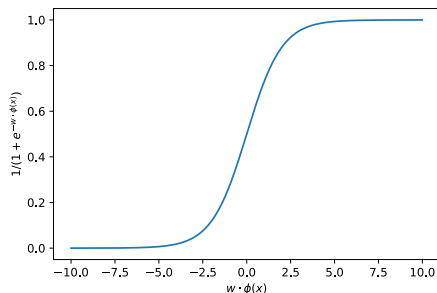
Parametrize $h(x)$ using a linear function:

$$h(x) = w \cdot \phi(x) + b \quad \phi: \mathcal{X} \rightarrow \mathbb{R}^d, w \in \mathbb{R}^d$$

Problem: $h(x) \in \mathbb{R}$ (score)

Logistic regression

Map $w \cdot \phi(x) \in \mathbb{R}$ to a probability by the **logistic function**



$$p(y = 1 \mid x; w) = \frac{1}{1 + e^{-w \cdot \phi(x)}} \quad (y \in \{0, 1\})$$

$$p(y = k \mid x; w) = \frac{e^{w_k \cdot \phi(x)}}{\sum_{i \in \mathcal{Y}} e^{w_i \cdot \phi(x)}} \quad (y \in \{1, \dots, K\})$$

“softmax”

Inference

$$\hat{y} = \arg \max_{k \in \mathcal{Y}} p(y = k \mid x; w) \quad (6)$$

$$= \arg \max_{k \in \mathcal{Y}} \frac{e^{w_k \cdot \phi(x)}}{\sum_{i \in \mathcal{Y}} e^{w_i \cdot \phi(x)}} \quad (7)$$

$$= \arg \max_{k \in \mathcal{Y}} e^{w_k \cdot \phi(x)} \quad (8)$$

$$= \arg \max_{k \in \mathcal{Y}} \underbrace{w_k \cdot \phi(x)}_{\text{score for class } k} \quad (9)$$

MLE for logistic regression

$$\max \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; w)$$

- Likelihood function is concave / NLL is convex

MLE for logistic regression

$$\max \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; w)$$

- Likelihood function is concave / NLL is convex
- No closed-form solution
- Use stochastic gradient ascent

BoW representation

Example:

$$\mathcal{V} = \{the, a, an, in, for, penny, pound\}$$

sentence = *in for a penny, in for a pound*

$$x = (in, for, a, penny, in, for, a, pound)$$

Feature extractor: $\phi: \mathcal{V}^n \rightarrow \mathbb{R}^d$.

BoW representation

Example:

$$\mathcal{V} = \{the, a, an, in, for, penny, pound\}$$

sentence = *in for a penny, in for a pound*

$$x = (in, for, a, penny, in, for, a, pound)$$

Feature extractor: $\phi: \mathcal{V}^n \rightarrow \mathbb{R}^d$.

Idea: a sentence is the “sum” of words.

$$\phi_{\text{BoW}}(x) = \sum_{i=1}^n \phi_{\text{one-hot}}(x_i)$$

$$\phi_{\text{one-hot}}(x_1) = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix} \quad \text{the sentence contains the word “in”}$$

$$\phi_{\text{BoW}}(x) = \begin{bmatrix} 0 & 2 & 0 & 2 & 2 & 1 & 1 \end{bmatrix} \quad \text{the sentence contains 2 occurrences of “in”}$$

N-gram features

Potential problems with the the BoW representation?

N-gram features

Potential problems with the the BoW representation?

N-gram features:

in for a penny , in for a pound

- Unigram: in, for, a, ...
- Bigram: in/for, for/a, a/penny, ...
- Trigram: in/for/a, for/a/penny, ...



What are the pros/cons of using higher order n-grams?

Feature extractor

Logistic regression allows for richer features (what's the limitation of NB?)

Define each feature as a function $\phi_i: \mathcal{X} \rightarrow \mathbb{R}$.

$$\phi_1(x) = \begin{cases} 1 & x \text{ contains "happy"} \\ 0 & \text{otherwise} \end{cases},$$

$$\phi_2(x) = \begin{cases} 1 & x \text{ contains words with suffix "yyy"} \\ 0 & \text{otherwise} \end{cases}.$$

In practice, use a dictionary

```
feature_vector["prefix=un+suffix=ing"] = 1
```