

L4: Logistic regression

Spring 2021

COS 484/584

- Supervised classification:
 - Document to classify, d
 - Set of classes, $C = \{c_1, c_2, ..., c_k\}$
- Naive Bayes:

Last class

 $\Lambda = \alpha \lambda g \max P(c) P(d|c)$





Powerful supervised model



- Powerful supervised model
- Baseline approach for many NLP tasks



- Powerful supervised model
- Baseline approach for many NLP tasks
- Connections with neural networks



- Powerful supervised model
- Baseline approach for many NLP tasks
- Connections with neural networks
- Binary (two classes) or multinomial (>2 classes)







• Logistic Regression is a *discriminative* model





- Logistic Regression is a discriminative model
- Naive Bayes: generative model







• Logistic Regression: $\Lambda = \operatorname{algmax} P(c|d)$ • Naive Bayes: $\hat{c} = algmax P(c) P(d|c)$



• Inputs:

• Inputs:

1. Classification instance in a **feature representation**

- Inputs:
 - 1. Classification instance in a **feature representation**
 - 2. Classification function to compute \hat{y} using $P(\hat{y} | x)$

- Inputs:
 - 1. Classification instance in a **feature representation**
 - 2. Classification function to compute \hat{y} using $P(\hat{y} | x)$
 - 3. Loss function (for learning)

- Inputs:
 - 1. Classification instance in a **feature representation**
 - 2. Classification function to compute \hat{y} using $P(\hat{y} | x)$
 - 3. Loss function (for learning)
 - 4. Optimization algorithm

- Inputs:
 - Classification instance in a **feature representation** 1.
 - 2. Classification function to compute \hat{y} using $P(\hat{y} | x)$
 - 3. Loss function (for learning)
 - 4. Optimization algorithm

• Train phase: Learn the parameters of the model to minimize loss function

• Inputs:

- Classification instance in a **feature representation** 1.
- **Classification function** to compute \hat{y} using $P(\hat{y} | x)$ 2.
- 3. Loss function (for learning)
- 4. Optimization algorithm

• Train phase: Learn the parameters of the model to minimize loss function

• Test phase: Apply parameters to predict class given a new input x

- Input observation: $x^{(i)}$
- Feature vector: $[x_1, x_2, \ldots, x_d]$
- Feature j of ith input : $x_i^{(i)}$

I. Feature representation

- Input observation: $x^{(i)}$
- Feature vector: $[x_1, x_2, \ldots, x_d]$
- Feature j of ith input : $x_i^{(i)}$

I. Feature representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

 $x^{(i)}$



Bag of words

- Given: Input feature vector $[x_1, x_2, \ldots, x_d]$

- Given: Input feature vector $[x_1, x_2, \ldots, x_d]$
- Output: P(y = 1 | x) and P(y = 0 | x)(binary classification)

- Given: Input feature vector $[x_1, x_2, \ldots, x_d]$
- Output: P(y = 1 | x) and P(y = 0 | x)(binary classification)
- Require a function, $F : \mathbb{R}^d \to [0,1]$

- Given: Input feature vector $[x_1, x_2, \ldots, x_d]$
- Output: P(y = 1 | x) and P(y = 0 | x)(binary classification)
- Require a function, $F : \mathbb{R}^d \to [0,1]$
- Sigmoid:

- Given: Input feature vector $[x_1, x_2, \ldots, x_d]$
- *Output:* P(y = 1 | x) and P(y = 0 | x)(binary classification)
- Require a function, $F : \mathbb{R}^d \to [0,1]$
- Sigmoid: 1 $1 + e^{-z}$ $1 + e^{z}$

 e^{z}

- Given: Input feature vector $[x_1, x_2, \ldots, x_d]$
- Output: P(y = 1 | x) and P(y = 0 | x)(binary classification)
- Require a function, $F : \mathbb{R}^d \to [0,1]$
- Sigmoid: $1 + e^{-z}$ $1 + e^{z}$





• Which features are important and how much?

- Which features are important and how much?
- Learn a vector of **weights** and a **bias**

- Which features are important and how much?
- Learn a vector of **weights** and a **bias**
- Weights: Vector of real numbers, $w = [w_1, w_2, \dots, w_d]$

- Which features are important and how much?
- Learn a vector of **weights** and a **bias**
- Weights: Vector of real numbers, $w = [w_1, w_2, \dots, w_d]$
- Bias: Scalar intercept, *b*

- Which features are important and how much?
- Learn a vector of **weights** and a **bias**
- Weights: Vector of real numbers, $w = [w_1, w_2, \dots, w_d]$
- Bias: Scalar intercept, b
- Given an instance, x: $z = \sum_{i=1}^{d} w_i x_i$ i=1

$$v_i x_i + b$$
 or $z = w \cdot x + b$

- Which features are important and how much?
- Learn a vector of **weights** and a **bias**
- Weights: Vector of real numbers, $w = [w_1, w_2, \dots, w_d]$
- Bias: Scalar intercept, b
- Given an instance, x: $z = \sum_{i=1}^{d} w_i x_i$ i=1

• Therefore, $y = \frac{e^{w \cdot x + b}}{1 + e^{w \cdot x + b}}$

$$v_i x_i + b$$
 or $z = w \cdot x + b$

What is the bias?

 Let's say we have a feature that is always set to 1 regardless of what the input text is.

(Credits: Richard Socher)


- Let's say we have a feature that is always set to 1 regardless of what the input text is.
- This is clearly not an informative feature. However, let's say it was the only one I had...

first, how many weights do l need to learn for this feature?



- Let's say we have a feature that is always set to 1 regardless of what the input text is.
- This is clearly not an informative feature. However, let's say it was the only one I had...

first, how many weights do l need to learn for this feature?

okay... what is the best set of weights for it?





- Let's say we have a feature that is always set to 1 regardless of what the input text is.
- This is clearly not an informative feature. However, let's say it was the only one I had...

first, how many weights do l need to learn for this feature?

 $w \cdot x + b$

okay... what is the best set of weights for it?





- Let's say we have a feature that is always set to 1 regardless of what the input text is.
- This is clearly not an informative feature. However, let's say it was the only one I had...

first, how many weights do l need to learn for this feature?

 $w \cdot x + b$

okay... what is the best set of weights for it?



$$y = \frac{e^{w \cdot x + b}}{1 + e^{w \cdot x + b}} = \frac{e^{w \cdot x}}{1 + e^{w \cdot x + b}}$$





Putting it together

• Given x, compute $z = w \cdot x + b$

Putting it together

- Given x, compute $z = w \cdot x + b$

Putting it together

• Compute probabilities: $P(y = 1 | x) = \frac{1}{1 + e^{-z}}$

- Given x, compute $z = w \cdot x + b$

 $P(y = 1) = \sigma$

Putting it together

• Compute probabilities: $P(y = 1 | x) = \frac{1}{1 + e^{-z}}$

$$f(w \cdot x + b) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

- Given x, compute $z = w \cdot x$ -
- Compute probabilities: P(y =
 - $P(y=1) = \sigma$
 - P(y = 0) = 1
 - = 1

Putting it together

$$+ b$$

$$= 1 |x) = \frac{1}{1 + e^{-z}}$$

$$\sigma(w \cdot x + b) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

$$- \sigma(w \cdot x + b)$$

$$- \frac{1}{1 + e^{-(w \cdot x + b)}} = \frac{e^{-(w \cdot x + b)}}{1 + e^{-(w \cdot x + b)}}$$

- Given x, compute $z = w \cdot x + b$
- Compute probabilities: P(y =
 - $P(y=1) = \sigma$
 - $P(y = 0) = 1 \sigma(w \cdot x + b)$

= 1 -

• Decision boundary:

Putting it together

$$F b$$

$$= 1 | x) = \frac{1}{1 + e^{-z}}$$

$$F(w \cdot x + b) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

$$= \sigma(w \cdot x + b)$$

$$\frac{1}{1+e^{-(w\cdot x+b)}} = \frac{e^{-(w\cdot x+b)}}{1+e^{-(w\cdot x+b)}}$$

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



Var	Definition	Value in Fig. 5.2
x_1	$count(positive lexicon) \in doc)$	3
x_2	$count(negative lexicon) \in doc)$	2
<i>x</i> 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
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	log(word count of doc)	$\ln(64) = 4.15$

Var	Definition	Value
x_1	$count(positive lexicon) \in doc)$	3
x_2	$count(negative lexicon) \in doc)$	2
<i>x</i> 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
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	$\log(word \ count \ of \ doc)$	$\ln(64) = 4.15$

Var	Definition
x_1	count(positive lexicon)
x_2	count(negative lexicon
<i>x</i> ₃	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_4	count(1st and 2nd prop
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_6	log(word count of doc



• Assume weights w = [2.5, -5.0, -1.2, 0.5, 2.0, 0.7] and bias b = 0.1

Var	Definition
x_1	count(positive lexicon)
x_2	count(negative lexicon
<i>x</i> ₃	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_4	count(1st and 2nd prop
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
Y.	log(word count of doc

 $\log(\text{word count of doc})$ x_6

• Assume weights
$$w = [2.5, -$$

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

= $\sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \cdot [3, 2, 1, 3, 0, 4.15] + 0.1)$
= $\sigma(.805)$

- = 0.69

$$p(-|x) = P(Y = 0|x) = 1 - \sigma(v)$$

= 0.31



-5.0, -1.2, 0.5, 2.0, 0.7] and bias b = 0.1

 $(w \cdot x + b)$

 Most important rule: Data is key!

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

$$x_{1} = \begin{cases} 1 & \text{if } ``Case(w_{i}) = \text{Lower''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2} = \begin{cases} 1 & \text{if } ``w_{i} \in \text{AcronymDict''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{3} = \begin{cases} 1 & \text{if } ``w_{i} = \text{St. } \& Case(w_{i-1}) = \text{Cap''} \\ 0 & \text{otherwise} \end{cases}$$

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

$$x_{1} = \begin{cases} 1 & \text{if } ``Case(w_{i}) = \text{Lower''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2} = \begin{cases} 1 & \text{if } ``w_{i} \in \text{AcronymDict''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{3} = \begin{cases} 1 & \text{if } ``w_{i} = \text{St. } \& Case(w_{i-1}) = \text{Cap''} \\ 0 & \text{otherwise} \end{cases}$$

• Feature templates

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

$$x_{1} = \begin{cases} 1 & \text{if } ``Case(w_{i}) = \text{Lower''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2} = \begin{cases} 1 & \text{if } ``w_{i} \in \text{AcronymDict''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{3} = \begin{cases} 1 & \text{if } ``w_{i} = \text{St. } \& Case(w_{i-1}) = \text{Cap''} \\ 0 & \text{otherwise} \end{cases}$$

- Feature templates
 - Sparse representations, hash only seen features into index

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

$$x_{1} = \begin{cases} 1 & \text{if } ``Case(w_{i}) = \text{Lower''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2} = \begin{cases} 1 & \text{if } ``w_{i} \in \text{AcronymDict''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{3} = \begin{cases} 1 & \text{if } ``w_{i} = \text{St. } \& Case(w_{i-1}) = \text{Cap''} \\ 0 & \text{otherwise} \end{cases}$$

- Feature templates
 - Sparse representations, hash only seen features into index
 - Ex. Trigram(*logistic regression classifier*) = Feature #78

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

$$x_{1} = \begin{cases} 1 & \text{if } ``Case(w_{i}) = \text{Lower''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2} = \begin{cases} 1 & \text{if } ``w_{i} \in \text{AcronymDict''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{3} = \begin{cases} 1 & \text{if } ``w_{i} = \text{St. } \& Case(w_{i-1}) = \text{Cap''} \\ 0 & \text{otherwise} \end{cases}$$

- Feature templates
 - Sparse representations, hash only seen features into index
 - Ex. Trigram(logistic regression classifier) = Feature #78
- Advanced: Representation
 learning (we will see this later!)

• More freedom in designing features

- More freedom in designing features

No strong independence assumptions like Naive Bayes

- More freedom in designing features
 - No strong independence assumptions like Naive Bayes
 - More robust to correlated features ("San Francisco" vs "Boston")
 —LR is likely to work better than NB

- More freedom in designing features
 - No strong independence assumptions like Naive Bayes
 - More robust to correlated features ("San Francisco" vs "Boston")
 —LR is likely to work better than NB
 - Can even have the same feature twice! (why?)

- More freedom in designing features
 - No strong independence assumptions like Naive Bayes
 - More robust to correlated features ("San Francisco" vs "Boston")
 —LR is likely to work better than NB
 - Can even have the same feature twice! (why?)
- May not work well on small datasets (compared to Naive Bayes)

- More freedom in designing features
 - No strong independence assumptions like Naive Bayes
 - More robust to correlated features ("San Francisco" vs "Boston")
 —LR is likely to work better than NB
 - Can even have the same feature twice! (why?)
- May not work well on small datasets (compared to Naive Bayes)
- Interpreting learned weights can be challenging

• We have our **classification function** - how to assign weights and bias?

- We have our **classification function** how to assign weights and bias?
- Goal: predicted label \hat{y} as close as possible to actual label y

- We have our **classification function** how to assign weights and bias?

• Goal: predicted label \hat{y} as close as possible to actual label y

• Distance metric/Loss function between \hat{y} and $y: L(\hat{y}, y)$

- We have our **classification function** how to assign weights and bias?
- Goal: predicted label \hat{y} as close as possible to actual label y
 - Distance metric/Loss function between \hat{y} and $y : L(\hat{y}, y)$
 - **Optimization algorithm** for updating weights

Loss function
• Assume $\hat{y} = \sigma(w \cdot x + b)$

Loss function

- Assume $\hat{y} = \sigma(w \cdot x + b)$

• $L(\hat{y}, y) =$ Measure of difference between \hat{y} and y. But what form?

- Assume $\hat{y} = \sigma(w \cdot x + b)$
- Maximum likelihood estimation (conditional):

• $L(\hat{y}, y) =$ Measure of difference between \hat{y} and y. But what form?

- Assume $\hat{y} = \sigma(w \cdot x + b)$
- Maximum likelihood estimation (conditional):
 - labels y paired with input x

• $L(\hat{y}, y) =$ Measure of difference between \hat{y} and y. But what form?

• Choose w and b such that $\log P(y|x)$ is maximized for true

- Assume $\hat{y} = \sigma(w \cdot x + b)$
- Maximum likelihood estimation (conditional):
 - labels y paired with input x
 - Similar to language models!

• $L(\hat{y}, y) =$ Measure of difference between \hat{y} and y. But what form?

• Choose w and b such that $\log P(y|x)$ is maximized for true

- Assume $\hat{y} = \sigma(w \cdot x + b)$
- Maximum likelihood estimation (conditional):
 - labels y paired with input x
 - Similar to language models!
 - max log $P(w_t | w_{t-n}, \dots, w_{t-1})$ given a corpus

• $L(\hat{y}, y) =$ Measure of difference between \hat{y} and y. But what form?

• Choose w and b such that $\log P(y|x)$ is maximized for true

• Assume a single data point (*x*, *y*) and two classes

- Assume a single data point (*x*, *y*) and two classes
- Classifier probability: $P(y|x) = \hat{y}^{y}(1 \hat{y})^{1-y}$

- Assume a single data point (x, y) and two classes
- Classifier probability: P(y)
- Log probability: $\log P(y|)$

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

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

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

- Assume a single data point (*x*, *y*) and two classes
- Classifier probability: P(y)
- Log probability: $\log P(y|x)$

• Loss: $-\log P(y|x) = -[y \log \hat{y} + (1 - y)\log(1 - \hat{y})]$

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

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

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

- Assume a single data point (x, y) and two classes
- Classifier probability: P(y)
- Log probability: $\log P(y|x)$

- Loss: $-\log P(y|x) = -[y \log \hat{y} + (1 y)\log(1 \hat{y})]$
 - $y = 1 \implies -\log \hat{y}$, and

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

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

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

$$y = 0 \implies -\log(1 - \hat{y})$$



Cross Entropy loss

• Assume n data points $(x^{(i)}, y^{(i)})$

Cross Entropy loss

- Assume n data points $(x^{(i)}, y^{(i)})$
- Classifier probability: $\prod_{i=1}^{n} P(y | x) = \prod_{i=1}^{n} \hat{y}^{y} (1 \hat{y})^{1-y}$

Cross Entropy loss

- Assume n data points $(x^{(i)}, y^{(i)})$
- Classifier probability: $\prod_{i=1}^{n}$

• Loss:
$$-\log \prod_{i=1}^{n} P(y \mid x) = -$$

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

$$\sum_{i=1}^{n} \log P(y \mid x)$$

i = 1

 $L_{CE} = -\sum_{x}^{n} \left[y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \right]$

Example: Computing CE Loss

Var	Definition	Value
x_1	$count(positive lexicon) \in doc)$	3
x_2	$count(negative lexicon) \in doc)$	2
<i>x</i> 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
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	log(word count of doc)	$\ln(64) = 4.15$

Example: Computing CE Loss

Var	Definition	Value
x_1	$count(positive lexicon) \in doc)$	3
x_2	$count(negative lexicon) \in doc)$	2
<i>x</i> 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
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	log(word count of doc)	$\ln(64) = 4.15$

- If y = 1 (positive sentiment), $L_{CE} = -\log(0.69) = 0.37$
- If y = 0 (negative sentiment), $L_{CE} = -\log(0.31) = 1.17$

• Assume weights w = [2.5, -5.0, -1.2, 0.5, 2.0, 0.7] and bias b = 0.1





• $L_{CE} = -\sum_{i=1}^{n} \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$ i = 1







•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)} + \sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)}] + \dots$$

• Ranges from 0 (perfect predictions) to ∞



$(1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)} + \sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)}] + \dots$$

- Ranges from 0 (perfect predictions) to ∞
- Lower the value, better the classifier



$(1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)} + i]$$

- Ranges from 0 (perfect predictions) to ∞
- Lower the value, better the classifier
- Cross-entropy between the true distribution P(y|x) and predicted distribution $P(\hat{y}|x)$



 $(1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$

• We have our **classification function** and **loss function** - how do we find the best w and b?

 $\theta = [w; b]$

• We have our **classification function** and **loss function** - how do we find the best w and b?

 $\theta = [w; b]$

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \frac{1}{n} \sum_{i=1}^{n$$

• We have our **classification function** and **loss function** - how do we find the best w and b?

 $\theta = [w; b]$

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \frac{1}{n} \sum_{i=1}^{n$$



• We have our **classification function** and **loss function** - how do we find the best w and b?

 $\theta = [w; b]$

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \frac{1}{n} \sum_$$

- Gradient descent:
 - Find direction of steepest slope

• We have our **classification function** and **loss function** - how do we find the best w and b?

 $\theta = [w; b]$

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n}$$

- Gradient descent:
 - Find direction of steepest slope
 - Move in the opposite direction

• We have our **classification function** and **loss function** - how do we find the best w and b?



 $\theta^{t+1} = \theta^t - \theta^t -$

$$\eta \frac{d}{d\theta} f(x;\theta)$$

Gradient descent for LR

- one global minimum)
 - No local minima to get stuck in
- Deep neural networks are not so easy
 - Non-convex

• Cross entropy loss for logistic regression is **convex** (i.e. has only





• Updates: $\theta^{t+1} = \theta^t - \eta \frac{d}{d\theta} f(x;\theta)$



• Updates: $\theta^{t+1} = \theta^t - \eta \frac{d}{d\theta} f(x; \theta)$



• Updates:
$$\theta^{t+1} = \theta^t - \eta \frac{d}{d\theta} f(x; \theta)$$

• Magnitude of movement along gradient


Learning Rate

• Updates:
$$\theta^{t+1} = \theta^t - \eta \frac{d}{d\theta} f(x; \theta)$$

- Magnitude of movement along gradient
- Higher/faster learning rate = larger updates to parameters





• In LR: weight *w* is a vector



- In LR: weight w is a vector
- Express slope as a partial derivative of l w.r.t each weight:





• In LR: weight w is a vector

 $\nabla_{\boldsymbol{\theta}} L(f(x; \boldsymbol{\theta}), y))$

• Express slope as a partial derivative of I w.r.t each weight:

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



• In LR: weight w is a vector

 ∇

• Express slope as a partial derivative of l w.r.t each weight:

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

• Updates: $\theta^{(t+1)} = \theta^t - \eta \nabla L(f(x; \theta), y)$



•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \sigma(w \cdot x^{(i)})]$$

 $(i) + b) + (1 - y^{(i)})\log(1 - \sigma(w \cdot x^{(i)} + b))]$

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \sigma(w \cdot x^{(i)})]$$

• Gradient, $\frac{dL_{CE}(w,b)}{dw_j} = \sum_{i=1}^n [\sigma(w \cdot x^{(i)} + b) - y^{(i)}]x_j^{(i)}$

 $(i) + b) + (1 - y^{(i)})\log(1 - \sigma(w \cdot x^{(i)} + b))]$

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \sigma(w \cdot x^{(i)})]$$

• Gradient, $\frac{dL_{CE}(w,b)}{dw_j} = \sum_{i=1}^{n} \left[\sigma(w \cdot x^{(i)} + b) - y^{(i)} \right] x_j^{(i)}$ Piff between two y ordination input feature

 $(a^{(i)} + b) + (1 - y^{(i)})\log(1 - \sigma(w \cdot x^{(i)} + b))]$

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \sigma(w \cdot x^{(i)})]$$

• Gradient, $\frac{dL_{CE}(w,b)}{dw_j} = \sum_{i=1}^{n} \left[\sigma(w \cdot x^{(i)} + b) - y^{(i)} \right] x_j^{(i)}$ Piff between two y ordination input feature

 $(a^{(i)} + b) + (1 - y^{(i)})\log(1 - \sigma(w \cdot x^{(i)} + b))]$

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \sigma(w \cdot x^{(i)})]$$

1T (1) n

Gradient,
$$\frac{dL_{CE}(w,b)}{dw_{j}} = \sum_{i=1}^{n} \left[\sigma(w \cdot x^{(i)} + b) - y^{(i)} \right] x_{j}^{(i)}$$

$$\frac{dL_{CE}(w,b)}{db} = \sum_{i=1}^{n} \left[\sigma(w \cdot x^{(i)} + b) - y^{(i)} \right]$$

$$from the second secon$$

 $(i) + b) + (1 - y^{(i)})\log(1 - \sigma(w \cdot x^{(i)} + b))]$

Stochastic Gradient Descent

	Online optimization	function STG
		# where:
		# fi
	Compute loss and	# x i
	compute 1055 und	# y i
	minimize after each	$ heta\! \leftarrow\! 0$
	training example	repeat til do
	eranning example	For each tr
		1. Optio
		Comp
		Comp
		2. $g \leftarrow \nabla$
		3. $\theta \leftarrow \theta$

return θ

- OCHASTIC GRADIENT DESCENT(L(), f(), x, y) returns θ L is the loss function is a function parameterized by θ
- is the set of training inputs $x^{(1)}$, $x^{(2)}$,..., $x^{(n)}$ is the set of training outputs (labels) $y^{(1)}$, $y^{(2)}$, ..., $y^{(n)}$

ne # see caption raining tuple $(x^{(i)}, y^{(i)})$ (in random order) onal (for reporting): pute $\hat{y}^{(i)} = f(x^{(i)}; \theta)$ $\nabla_{\boldsymbol{\theta}} L(f(x^{(i)}; \boldsymbol{\theta}), y^{(i)})$ $\theta - \eta g$

How are we doing on this tuple? # What is our estimated output \hat{y} ? pute the loss $L(\hat{y}^{(i)}, y^{(i)})$ # How far off is $\hat{y}^{(i)}$ from the true output $y^{(i)}$? # How should we move θ to maximize loss? # Go the other way instead



Stochastic Gradient Descent

•	Online optimization	function S	STO
		# whe	re:
		#	fi
	Compute loss and	#	x i
•	Compute 1035 and	#	y i
	minimize after each	$ heta\! \leftarrow\! 0$	
	training example	repeat til	doı
	eranning example	For each	n tr
	Per	1. Op	tio
	Tuctar		mŗ
		Co	mŗ
	LOSS	2. g (-V
		3. <i>θ</i>	$-\epsilon$
		return θ	

- OCHASTIC GRADIENT DESCENT(L(), f(), x, y) returns θ L is the loss function is a function parameterized by θ
- is the set of training inputs $x^{(1)}$, $x^{(2)}$,..., $x^{(n)}$ is the set of training outputs (labels) $y^{(1)}$, $y^{(2)}$, ..., $y^{(n)}$

ne # see caption raining tuple $(x^{(i)}, y^{(i)})$ (in random order) onal (for reporting): pute $\hat{y}^{(i)} = f(x^{(i)}; \theta)$ $\nabla_{\boldsymbol{\theta}} L(f(x^{(i)}; \boldsymbol{\theta}), y^{(i)})$ $\theta - \eta g$

- # How are we doing on this tuple? # What is our estimated output \hat{y} ?
- pute the loss $L(\hat{y}^{(i)}, y^{(i)})$ # How far off is $\hat{y}^{(i)}$ from the true output $y^{(i)}$? # How should we move θ to maximize loss? # Go the other way instead



Stochastic Gradient Descent

- Online optimization
- Compute loss and minimize after each training example



• Training objective: $\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)})$

- Training objective: $\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)})$
- This might fit the training set too well! (including noisy features)

- Training objective: $\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)})$
- This might fit the training set too well! (including noisy features)
- Poor generalization to the unseen test set **Overfitting**

- Training objective: $\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)})$
- This might fit the training set too well! (including noisy features)
- Poor generalization to the unseen test set **Overfitting**
- **Regularization** helps prevent overfitting

$$\hat{\theta} = \arg \max_{\theta} \left[\sum_{i=1}^{n} \log \theta \right]$$

 $P(y^{(i)} | x^{(i)}) - \alpha R(\theta)$

- Training objective: $\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)})$
- This might fit the training set too well! (including noisy features)
- Poor generalization to the unseen test set **Overfitting**
- **Regularization** helps prevent overfitting

$$\hat{\theta} = \arg \max_{\theta} \left[\sum_{i=1}^{n} \log \theta \right]$$

 $P(y^{(i)} | x^{(i)}) - \alpha R(\theta)]$ Penalizi Lange Lange weights

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

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

• Euclidean distance of weight vector θ from origin

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

- Euclidean distance of weight vector θ from origin
- L2 regularized objective:

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

- Euclidean distance of weight vector θ from origin
- L2 regularized objective:

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

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

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

• Manhattan distance of weight vector θ from origin

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

- Manhattan distance of weight vector θ from origin
- L1 regularized objective:

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

- Manhattan distance of weight vector θ from origin
- L1 regularized objective:

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

L2 vs L1 regularization





• L2 is easier to optimize - simpler derivation

L2 vs L1 regularization





- L2 is easier to optimize simpler derivation
 - L1 is complex since the derivative of $|\theta|$ is not continuous at 0

L2 vs L1 regularization







- L2 is easier to optimize simpler derivation
 - L1 is complex since the derivative of $|\theta|$ is not continuous at 0
- L2 leads to many small weights (due to θ^2 term)

L2 vs L1 regularization



- L2 is easier to optimize simpler derivation
 - L1 is complex since the derivative of $|\theta|$ is not continuous at 0
- L2 leads to many small weights (due to θ^2 term)
 - L1 prefers sparse weight vectors with many weights set to 0 (i.e. far fewer features used)

L2 vs L1 regularization


- L2 is easier to optimize simpler derivation
 - L1 is complex since the derivative of $|\theta|$ is not continuous at 0
- L2 leads to many small weights (due to θ^2 term)
 - L1 prefers sparse weight vectors with many weights set to 0 (i.e. far fewer features used)

L2 vs L1 regularization



tagging, named entity recognition)

• What if we have more than 2 classes? (e.g. Part of speech

- tagging, named entity recognition)
- Need to model $P(y = c | x) \quad \forall c \in C$

• What if we have more than 2 classes? (e.g. Part of speech

- What if we have more than 2 classes? (e.g. Part of speech tagging, named entity recognition)
- Need to model $P(y = c | x) \quad \forall c \in C$
- Generalize **sigmoid** function to **softmax**

- What if we have more than 2 classes? (e.g. Part of speech tagging, named entity recognition)
- Need to model $P(y = c | x) \quad \forall c \in C$
- Generalize **sigmoid** function to **softmax**

$$softmax(z_i) =$$

$$\frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad 1 \le i \le k$$

- What if we have more than 2 classes? (e.g. Part of speech tagging, named entity recognition)
- Need to model $P(y = c | x) \quad \forall c \in C$
- Generalize **sigmoid** function to **softmax**

$$softmax(z_i) =$$

$$\frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad 1 \le i \le k$$

• Similar to sigmoid, softmax squashes values towards 0 or 1

- Similar to sigmoid, softmax squashes values towards 0 or 1
- If z = [0, 1, 2, 3, 4], then

- Similar to sigmoid, softmax squashes values towards 0 or 1
- If z = [0, 1, 2, 3, 4], then
 - $\operatorname{softmax}(z) = ([0.0117, 0.0317, 0.0861, 0.2341, 0.6364])$

- If z = [0, 1, 2, 3, 4], then
 - softmax(z) = ([0.0117, 0.0317, 0.0861, 0.2341, 0.6364])
- For multinomial LR,

• Similar to sigmoid, softmax squashes values towards 0 or 1

- If z = [0, 1, 2, 3, 4], then
 - softmax(z) = ([0.0117, 0.0317, 0.0861, 0.2341, 0.6364])
- For multinomial LR,

$$P(y = c | x) = \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^k e^{w_j \cdot x + b_j}}$$

• Similar to sigmoid, softmax squashes values towards 0 or 1

- Features need to include both input (x) and class (c)
- Implicit in binary case

Var	Definition	Wt
$f_1(0,x)$	$\int 1$ if "!" \in doc	-4.5
	$\int 0$ otherwise	
$f_1(+,x)$	$\int 1$ if "!" \in doc	2.6
	0 otherwise	
$f_1(-,x)$	$\int 1$ if "!" \in doc	1.3
	0 otherwise	

Features in multinomial LR

• Generalize binary loss to multinomial CE loss: $L_{CE}(\hat{y}, y) = -\sum_{k=1}^{k} 1\{y = c\} \log P(y = c \mid x)$ c=1 $= -\sum_{c=1}^{k} 1\{y = c\} \log \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^{k} e^{w_j \cdot x + b_j}}$

• Generalize binary loss to multinomial CE loss: $L_{CE}(\hat{y}, y) = -\sum_{c=1}^{k} 1\{y = c\} \log P(y = c \mid x)$ $= -\sum_{c=1}^{k} 1\{y = c\} \log \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^{k} e^{w_j \cdot x + b_j}}$

Binary CE Loss: $-\log P(y|x) = -[y \log \hat{y} + (1 - y)\log(1 - \hat{y})]$



• Generalize binary loss to multinomial CE loss:

$$L_{CE}(\hat{y}, y) = -\sum_{c=1}^{k} 1\{y = c\} \log P(y = c \mid x)$$
$$= -\sum_{c=1}^{k} 1\{y = c\} \log \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^{k} e^{w_j \cdot x + b_j}}$$

Binary CE Loss: $-\log P(y|x) = -[y \log \hat{y} + (1-y)\log(1-\hat{y})]$



• Generalize binary loss to multinomial CE loss: $L_{CE}(\hat{y}, y) = -\sum_{c=1}^{k} 1\{y = c\} \log P(y = c \mid x)$ $= -\sum_{c=1}^{k} 1\{y = c\} \log \frac{e^{w_c \cdot x + b_c}}{-k}$

$$= -\sum_{c=1}^{k} [1\{y=c\}] \log \frac{1}{\sum_{j=1}^{k} e^{w_j \cdot x + b_j}}$$

Gradient:

Binary CE Loss: $-\log P(y|x) = -[y \log \hat{y} + (1-y)\log(1-\hat{y})]$



• Generalize binary loss to multinomial CE loss: $L_{CF}(\hat{y}, y) = -\sum_{k=1}^{k} 1\{y = c\} \log P(y = c \mid x)$

$$L_{CE}(y, y) = -\sum_{c=1}^{k} 1\{y = c\} \log P(y = c \mid x)$$
$$= -\sum_{c=1}^{k} 1\{y = c\} \log \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^{k} e^{w_j \cdot x + b_j}}$$

Gradient:

$$\frac{dL_{CE}}{dw_c} = -\left(1\{y=c\} - P(y=c|z)\right)$$
$$= -\left(1\{y=c\} - \frac{e^{w_c \cdot x+b}}{\sum_{j=1}^k e^{w_j \cdot x}}\right)$$

Binary CE Loss: $-\log P(y|x) = -[y \log \hat{y} + (1-y)\log(1-\hat{y})]$

x))x

 \mathfrak{h}_c $\overline{x_{j}\cdot x+b_{j}}$ X



Give us feedback!

https://forms.gle/XJXbFLsJCfSNsjCUA