Natural Language Models and Interfaces

BSc Artificial Intelligence

Lecturer: Wilker Aziz Institute for Logic, Language, and Computation

2020, week 5, lecture b

Text Classifiers

We have learnt two techniques to design feature-rich models

- ▶ naive Bayes classification
- ▶ logistic regression

Text Classifiers

We have learnt two techniques to design feature-rich models

- naive Bayes classification
- ► logistic regression

They are both useful to condition on high-dimensional data

- ▶ NBC uses Bayes rule and a conditional independence assumption
- ▶ LR uses a linear model and the softmax function

Pros and Cons

NBC

Pros and Cons

NBC

- exact MLE solution
- strong independence assumption

LR

Pros and Cons

NBC

- exact MLE solution
- strong independence assumption

LR

- flexible
- no closed-form MLE
- but gradient ascent converges to global optimum because the log-likelihood function is concave

Applications

- 1. Text classifiers: predict a categorical target from high-dimensional input
- 2. Component in a generative model: parameterise a cpd that conditions on high-dimensional input

Text Classification

Task	Description
Sentiment analysis	emotion towards a subject
Textual entailment	x is a document and y is a binary label given a two pieces of text, does one entail or contradict the other?
	\boldsymbol{x} is a pair of documents and \boldsymbol{y} is a binary label

Text Classification

Task	Description
Sentiment analysis	emotion towards a subject
Textual entailment	x is a document and y is a binary label given a two pieces of text, does one entail or contradict the other?
	\boldsymbol{x} is a pair of documents and \boldsymbol{y} is a binary label

In both cases, there are extensions to multiple classes

- Stanford sentiment classification: 5 sentiment levels
- Stanford natural language inference: 3 logical entailment relations

A feature function in NBC is a little different from a feature function in LR

- lacktriangle in NBC only the input x is available to the feature function
- ightharpoonup in LR both x and y are available

Why?

A feature function in NBC is a little different from a feature function in LR $\,$

- ightharpoonup in NBC only the input x is available to the feature function
- ightharpoonup in LR both x and y are available

Why?

► NBC:

$$P_{Y|X}(y|x) \stackrel{\mathsf{def}}{=} P_{Y|F_1^n}(y|f_1^n = \textcolor{red}{h(x)})$$

A feature function in NBC is a little different from a feature function in LR $\,$

- lacktriangle in NBC only the input x is available to the feature function
- ightharpoonup in LR both x and y are available

Why?

► NBC:

$$P_{Y|X}(y|x) \stackrel{\text{def}}{=} P_{Y|F_1^n}(y|f_1^n = h(x)) \propto P_Y(y) \prod_{i=1}^n P_{F|Y}(f_i|y)$$

A feature function in NBC is a little different from a feature function in LR

- lacktriangle in NBC only the input x is available to the feature function
- ightharpoonup in LR both x and y are available

Why?

► NBC:

$$P_{Y|X}(y|x) \stackrel{\mathsf{def}}{=} P_{Y|F_1^n}(y|f_1^n = h(x)) \propto P_Y(y) \prod_{i=1}^n P_{F|Y}(f_i|y)$$

LR:

$$P_{Y|X}(y|x) = \frac{\exp\left(w^{\top}f(\mathbf{y}, \mathbf{x})\right)}{\sum_{y' \in \mathcal{Y}} \exp(w^{\top}f(y', x))}$$

Feature functions (cont.)

In NBC a feature is a random variable while in LR a feature is just input to a log-linear model

Consider a sentence like

I did not like the acting, but the plot was decent

Example features are:

unigrams: I, did, not, like, the, acting, but, the plot, was, decent

Consider a sentence like

I did not like the acting, but the plot was decent

Example features are:

- unigrams: I, did, not, like, the, acting, but, the plot, was, decent
- we can apply certain filters: frequency based, stopwords

Consider a sentence like

I did not like the acting, but the plot was decent

Example features are:

- unigrams: I, did, not, like, the, acting, but, the plot, was, decent
- we can apply certain filters: frequency based, stopwords
- we can distinguish some words (e.g. sentiment words): like, decent

Consider a sentence like

I did not like the acting, but the plot was decent

Example features are:

- unigrams: I, did, not, like, the, acting, but, the plot, was, decent
- we can apply certain filters: frequency based, stopwords
- we can distinguish some words (e.g. sentiment words): like, decent
- we can use external resources (e.g. POS tagger)

Consider a sentence like

I did not like the acting, but the plot was decent

Example features are:

- unigrams: I, did, not, like, the, acting, but, the plot, was, decent
- we can apply certain filters: frequency based, stopwords
- we can distinguish some words (e.g. sentiment words): like, decent
- we can use external resources (e.g. POS tagger)

We can pre-process the data to account for negation scope

I did not like-NEG the-NEG acting-NEG, but the plot was decent

Consider a sentence like

I did not like the acting, but the plot was decent

We can consider the same types of features and more

▶ we can pair them with the class: like₊, like_−

Consider a sentence like

I did not like the acting, but the plot was decent

We can consider the same types of features and more

- ▶ we can pair them with the class: like₊, like_−
- we can use overlapping features: not like, like acting, plot decent

Consider a sentence like

I did not like the acting, but the plot was decent

We can consider the same types of features and more

- ▶ we can pair them with the class: like₊, like_−
- we can use overlapping features: not like, like acting, plot decent
- with good regularisation, we can have far more features

Evaluation

For binary classification

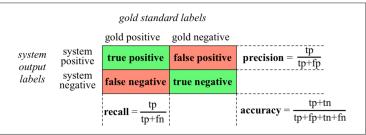


Figure 4.4 Contingency table

- use a training/development/test split
- or, preferably, cross-validation

Recall that we modified the HMM by combining it with a bigram LM?

Recall that we modified the HMM by combining it with a bigram LM?

$$P_{X_1^n C_1^n | N}(x_1^n, c_1^n | n) = \prod_{i=1}^n P_{C|C_{\mathsf{prev}}}(c_i | c_{i-1}) P_{X|X_{\mathsf{prev}}C}(x_i | x_{i-1}, c_i)$$

Recall that we modified the HMM by combining it with a bigram LM?

$$P_{X_1^n C_1^n | N}(x_1^n, c_1^n | n) = \prod_{i=1}^n P_{C|C_{\mathsf{prev}}}(c_i | c_{i-1}) P_{X|X_{\mathsf{prev}}C}(x_i | x_{i-1}, c_i)$$

Note that we modified emission probabilities to be

$$P_{X|X_{\mathsf{prev}}C}(x|x',c)$$

Recall that we modified the HMM by combining it with a bigram LM?

$$P_{X_1^n C_1^n | N}(x_1^n, c_1^n | n) = \prod_{i=1}^n P_{C|C_{\mathsf{prev}}}(c_i | c_{i-1}) P_{X|X_{\mathsf{prev}}C}(x_i | x_{i-1}, c_i)$$

Note that we modified emission probabilities to be

$$P_{X|X_{\mathsf{prev}}C}(x|x',c)$$

Interpolated CPD

This is a heuristic technique whereby we use a convex combination of simpler CPDs:

$$P_{X|X_{\mathsf{prev}}C}(x|x',c) = \alpha \times P_{X|X_{\mathsf{prev}}}(x|x') + (1-\alpha) \times P_{X|C}(x|c)$$

Interpolated CPD

This is a heuristic technique whereby we use a convex combination of simpler CPDs:

$$P_{X|X_{\mathsf{prev}}C}(x|x',c) = \alpha \times P_{X|X_{\mathsf{prev}}}(x|x') + (1-\alpha) \times P_{X|C}(x|c)$$

- it requires $0 < \alpha < 1$ for which no closed-form MLE is available
- \blacktriangleright we need to tune α on held-out data
- and estimate the simpler cpds independently

$$P_{X|X_{\mathsf{prev}}}(x|x') = \frac{\operatorname{count}_{X_{\mathsf{prev}}X}(x',x)}{\operatorname{count}_{X_{\mathsf{prev}}}(x')} \text{ and } P_{X|C}(x|c) = \frac{\operatorname{count}_{CX}(c,x)}{\operatorname{count}_{C}(c)}$$

We can use the Naive Bayes assumption to model CPDs!

$$P_{X|X_{\mathsf{prev}}C}(x|\mathbf{x'},\mathbf{c}) =$$

We can use the Naive Bayes assumption to model CPDs!

$$P_{X|X_{\mathsf{prev}}C}(x|\mathbf{x}',\mathbf{c}) = \frac{P_X(x)P_{X_{\mathsf{prev}}C|\mathbf{x}}(x',c|\mathbf{x})}{P_{X_{\mathsf{prev}}C}(x',c)}$$

We can use the Naive Bayes assumption to model CPDs!

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\mathbf{x'},\mathbf{c}) &= \frac{P_X(x)P_{X_{\mathsf{prev}}C|\mathbf{x}}(x',c|\mathbf{x})}{P_{X_{\mathsf{prev}}C}(x',c)} \\ &\overset{\text{ind}}{=} \frac{P_X(x)P_{X_{\mathsf{prev}}|\mathbf{x}}(x'|\mathbf{x})P_{C|X}(c|\mathbf{x})}{P_{X_{\mathsf{prev}}C}(x',c)} \end{split}$$

We can use the Naive Bayes assumption to model CPDs!

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\mathbf{x'},\mathbf{c}) &= \frac{P_X(x)P_{X_{\mathsf{prev}}C|x}(x',c|x)}{P_{X_{\mathsf{prev}}C}(x',c)} \\ &\overset{\text{ind}}{=} \frac{P_X(x)P_{X_{\mathsf{prev}}|x}(x'|x)P_{C|X}(c|x)}{P_{X_{\mathsf{prev}}C}(x',c)} \end{split}$$

Note the denominator needs to be inferred

$$P_{X_{\mathsf{prev}}C}(x',c) = \sum_{x \in \mathcal{X}} P_{X_{\mathsf{prev}}CX}(x',c,x)$$

We can use the Naive Bayes assumption to model CPDs!

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\mathbf{x'},\mathbf{c}) &= \frac{P_X(x)P_{X_{\mathsf{prev}}C|x}(x',c|x)}{P_{X_{\mathsf{prev}}C}(x',c)} \\ &\overset{\text{ind}}{=} \frac{P_X(x)P_{X_{\mathsf{prev}}|x}(x'|x)P_{C|X}(c|x)}{P_{X_{\mathsf{prev}}C}(x',c)} \end{split}$$

Note the denominator needs to be inferred

$$\begin{split} P_{X_{\mathsf{prev}}C}(x',c) &= \sum_{x \in \mathcal{X}} P_{X_{\mathsf{prev}}CX}(x',c,x) \\ &= \sum_{x \in \mathcal{X}} P_{X}(x) P_{X_{\mathsf{prev}}|X}(x'|x) P_{C|X}(c|x) \end{split}$$

We can use the Naive Bayes assumption to model CPDs!

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\mathbf{x'},\mathbf{c}) &= \frac{P_X(x)P_{X_{\mathsf{prev}}C|x}(x',c|x)}{P_{X_{\mathsf{prev}}C}(x',c)} \\ &\overset{\text{ind}}{=} \frac{P_X(x)P_{X_{\mathsf{prev}}|x}(x'|x)P_{C|X}(c|x)}{P_{X_{\mathsf{prev}}C}(x',c)} \end{split}$$

Note the denominator needs to be inferred

$$\begin{split} P_{X_{\mathsf{prev}}C}(x',c) &= \sum_{x \in \mathcal{X}} P_{X_{\mathsf{prev}}CX}(x',c,x) \\ &= \sum_{x \in \mathcal{X}} P_{X}(x) P_{X_{\mathsf{prev}}|X}(x'|x) P_{C|X}(c|x) \end{split}$$

Pro: MLE is exact (no need to tune heuristic coefficients) Con: denominator must be computed $O(|\mathcal{X}|)$ (this scales linearly in vocabulary size)

This is a direct application of LR, for example:

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\pmb{x'}, \pmb{c}; w) &= \frac{\exp\left(w^{\top} f(x, \pmb{x'}, \pmb{c})\right)}{\mathcal{Z}(x|w)} \\ \mathcal{Z}(\pmb{x'}, \pmb{c}|w) &= \sum_{x \in \mathcal{X}} \exp\left(w^{\top} f(x, \pmb{x'}, \pmb{c})\right) \end{split}$$

This is a direct application of LR, for example:

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\pmb{x'},\pmb{c};w) &= \frac{\exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr)}{\mathcal{Z}(x|w)} \\ \mathcal{Z}(\pmb{x'},\pmb{c}|w) &= \sum_{x \in \mathcal{X}} \exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr) \end{split}$$

Pro: flexible Con: no closed-form MLE

This is a direct application of LR, for example:

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\pmb{x'},\pmb{c};w) &= \frac{\exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr)}{\mathcal{Z}(x|w)} \\ \mathcal{Z}(\pmb{x'},\pmb{c}|w) &= \sum_{x \in \mathcal{X}} \exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr) \end{split}$$

Pro: flexible

Con: no closed-form MLE

but we can start with some random $w \in \mathbb{R}^D$ such that $w_d \sim \mathcal{N}(0,1)$

This is a direct application of LR, for example:

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\pmb{x'},\pmb{c};w) &= \frac{\exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr)}{\mathcal{Z}(x|w)} \\ \mathcal{Z}(\pmb{x'},\pmb{c}|w) &= \sum_{x \in \mathcal{X}} \exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr) \end{split}$$

Pro: flexible

Con: no closed-form MLE

- but we can start with some random $w \in \mathbb{R}^D$ such that $w_d \sim \mathcal{N}(0,1)$
- ightharpoonup compute the log-likelihood function $\mathcal{L}(w|\mathcal{D})$ for a dataset

This is a direct application of LR, for example:

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\pmb{x'},\pmb{c};w) &= \frac{\exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr)}{\mathcal{Z}(x|w)} \\ \mathcal{Z}(\pmb{x'},\pmb{c}|w) &= \sum_{x \in \mathcal{X}} \exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr) \end{split}$$

Pro: flexible

Con: no closed-form MLE

- but we can start with some random $w \in \mathbb{R}^D$ such that $w_d \sim \mathcal{N}(0,1)$
- lacktriangle compute the log-likelihood function $\mathcal{L}(w|\mathcal{D})$ for a dataset
- lacktriangle and then $w + oldsymbol{
 abla}_w \mathcal{L}(w|\mathcal{D})$ takes us closer to the optimum

This is a direct application of LR, for example:

$$\begin{split} P_{X|X_{\mathsf{prev}}C}(x|\pmb{x'},\pmb{c};w) &= \frac{\exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr)}{\mathcal{Z}(x|w)} \\ \mathcal{Z}(\pmb{x'},\pmb{c}|w) &= \sum_{x \in \mathcal{X}} \exp\Bigl(w^\top f(x,\pmb{x'},\pmb{c})\Bigr) \end{split}$$

Pro: flexible

Con: no closed-form MLE

- but we can start with some random $w \in \mathbb{R}^D$ such that $w_d \sim \mathcal{N}(0,1)$
- lacktriangle compute the log-likelihood function $\mathcal{L}(w|\mathcal{D})$ for a dataset
- lacktriangle and then $w + oldsymbol{
 abla}_w \mathcal{L}(w|\mathcal{D})$ takes us closer to the optimum
- eventually the log-likelihood function stops improving and we have a global optimum

References I