

Natural Language Processing

CSE 325/425



Sihong Xie

Lecture 3:

- MLE for n-grams
- General MLE
- MAP

Estimating n-gram model using MLE

- Maximum Likelihood Estimation (MLE)

- Let the count of word w be $c(w)$
- Then given a corpus, we can count the number of times any word $w \in V$ appears.
- Assumeing a unigram language model: words are selected independent of other words.
- The likelihood of the corpus is

$$P(w) = \frac{c(w)}{n}$$

n : # of words
in training
Corpus

$$P(w_1, \dots, w_n) \propto P(w_1) \times \dots \times P(w_n)$$

$$V: \text{Vocabulary} \propto \prod_{w \in V} P(w)^{c(w)} = \underbrace{L(w_1, \dots, w_n | \theta)}_{\text{Likelihood}}$$

- Parameters of the model: $\theta = \{P(w) : w \in V\}$
- Sufficient statistics: $\{c(w) : w \in V\}$

$$\theta = [\theta_1 \dots \theta_{|V|}]$$

$$= [p(w_1) \dots P(|V|)]$$

↑
prob of seeing
word 1

Estimating uni-gram model using MLE

- Maximum Likelihood Estimation (MLE)
 - By choosing the parameters $\theta = \{P(w) : w \in V\}$ to maximize the likelihood

$$\begin{aligned} \max_{\theta} L(w_1, \dots, w_n | \theta) &= \text{(unigram)} \\ \text{s.t. } \sum_{w \in V} P(w) &= 1 \end{aligned}$$

$= \text{Pr}(w_1) \times \text{Pr}(w_2) \times \dots \times \text{Pr}(w_n)$
 (bigram)
 HW 1

- Using the Lagrangian method, you will find the MLE

$$P(w) = \frac{\text{Count of } w}{\text{Count of the total tokens}}$$

- MLE is considered the "best" estimation of the model parameter.

"Laplacian" estimation to handle words in the V
but unseen in the corpus

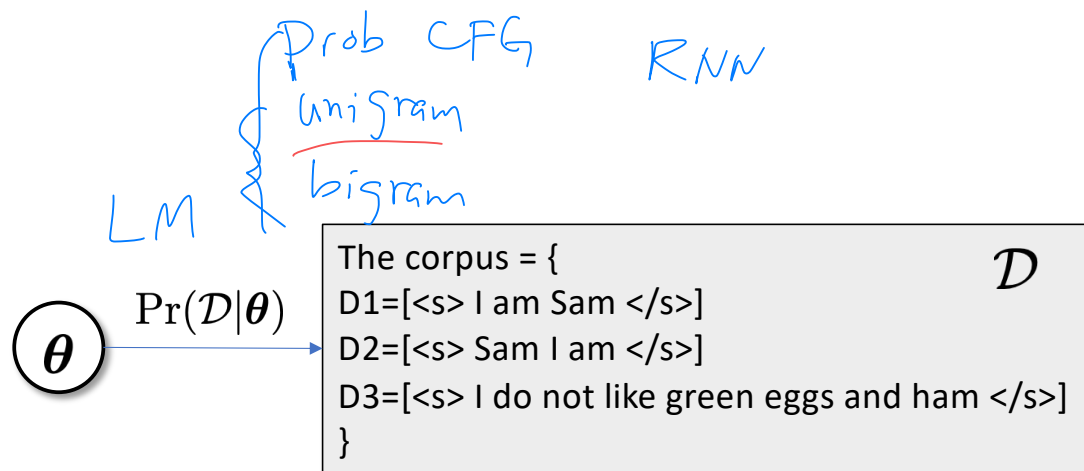
Estimating n-gram model using MLE

- Maximum Likelihood Estimation (MLE)
 - Can Laplacian be derived using MLE?
 - Yes! Just add a hallucinated appearance of every word to the corpus.
 - Suppose your vocabulary is
$$V = \{<s>, </s>, I, am, Sam, do, not, like, green, eggs, and, ham, red\}$$
 - Adding a hallucinated document containing each word:

The corpus = {
D0=[<s>, </s>, I, am, Sam, do, not, like, green, eggs, and, ham, red]
D1=[<s> I am Sam </s>]
D2=[<s> Sam I am </s>]
D3=[<s> I do not like green eggs and ham </s>]
}

$$P(w) = ?$$

General MLE



- Maximum Likelihood Estimation (MLE)

- Given observed data \mathcal{D} (training Data in ML)
 - can be unsupervised: a corpus

~~×~~ supervised: documents labeled with sentiments {positive, neutral, negative}

- Assume a probabilistic model generating the data

- parametric model for likelihood of the data $\Pr(\mathcal{D}|\theta)$

- In NLP, commonly found are multi-nomial models

- Probabilities of seeing a word: $\theta = [\theta_1, \dots, \theta_{|V|}]$ $\sum_w \theta_w = 1$
- Likelihood of a corpus using uni-gram:

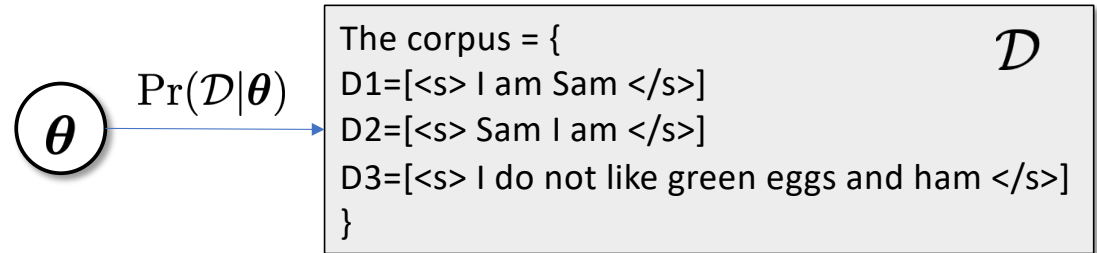
$$\Pr(\mathcal{D}|\theta) = \frac{1}{Z} \prod_{w=1}^{|V|} \theta_w^{c(w)} \quad \leftarrow \text{sufficient Statistics}$$

How about MLE for bi-gram model?

$$l = \sum_{\mathcal{D}} \Pr(\mathcal{D}|\theta)$$

unknown parameter

General MLE



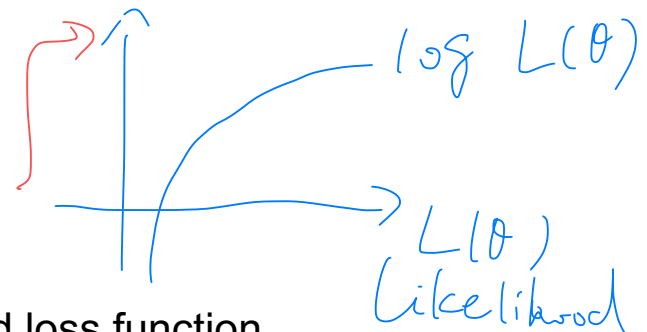
- Maximum Likelihood Estimation (MLE)

- Given observed data \mathcal{D}
- Parametric model for **likelihood** of the data $\Pr(\mathcal{D}|\theta)$ $L(\theta) = \Pr(\mathcal{D}|\theta)$
- What's the most likely θ that generates the data?
- MLE is the constrained optimization problem

$$\theta^* = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \log L(\theta)$$

$$\text{s.t. } \sum_w \theta_w = 1$$

log-likelihood



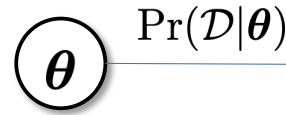
- Equivalent to minimizing the negative log-likelihood loss function

$$-\log L(\theta) = J(\theta)$$

$$\theta^* = \arg \min_{\theta} -\log L(\theta) = \arg \min_{\theta} J(\theta)$$

optimized using gradient descent
in general

General MLE



The corpus = \mathcal{D}
 D1=[<s> I am Sam </s>]
 D2=[<s> Sam I am </s>]
 D3=[<s> I do not like green eggs and ham </s>]
 }

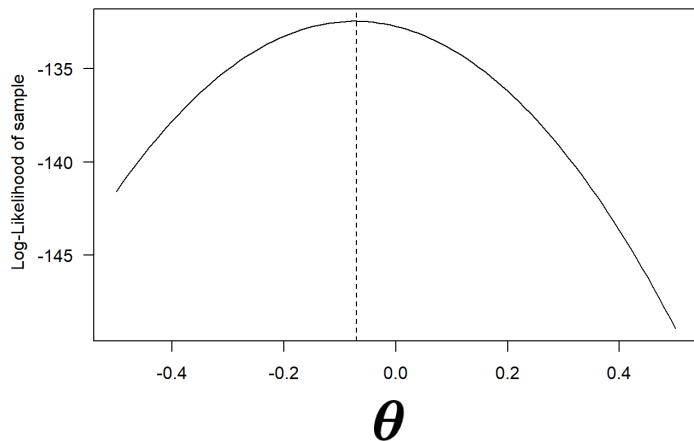
- Revisiting MLE for uni-gram model

$$\theta^* = \arg \max_{\theta} \left[\frac{1}{Z} \prod_{w=1}^{|V|} \theta_w^{c(w)} \right] = \arg \max_{\theta} -\log Z + \sum_{w=1}^{|V|} c(w) \log \theta_w$$

s.t. $\sum_w \theta_w = 1$

↗ objective
function

$J(\theta)$
 $= \arg \min_{\theta} - \sum_{w=1}^{|V|} c(w) \log \theta_w$



But how to solve this problem?

Need two more tools:

- gradient descent
- matrix calculus
- Lagrangian methods

Alg GD

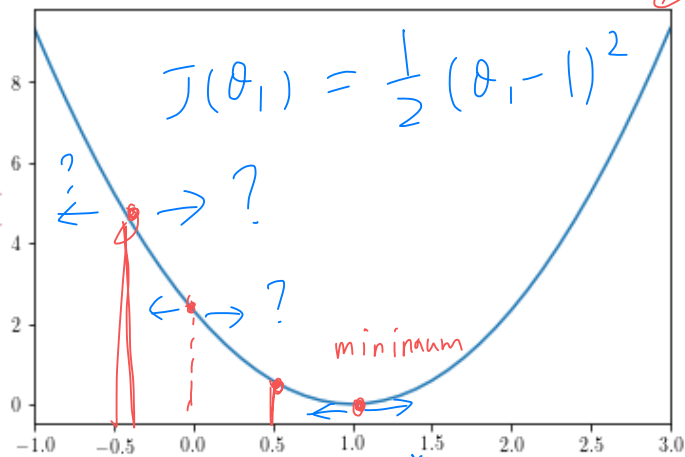
Initialize $\theta = \theta^{(0)}$ randomly
for $t = 1 \dots T$

$$\theta^{(t)} \leftarrow \theta^{(t-1)} - \eta \times \left. \frac{\partial J(\theta)}{\partial \theta} \right|_{\theta = \theta^{(t-1)}}$$

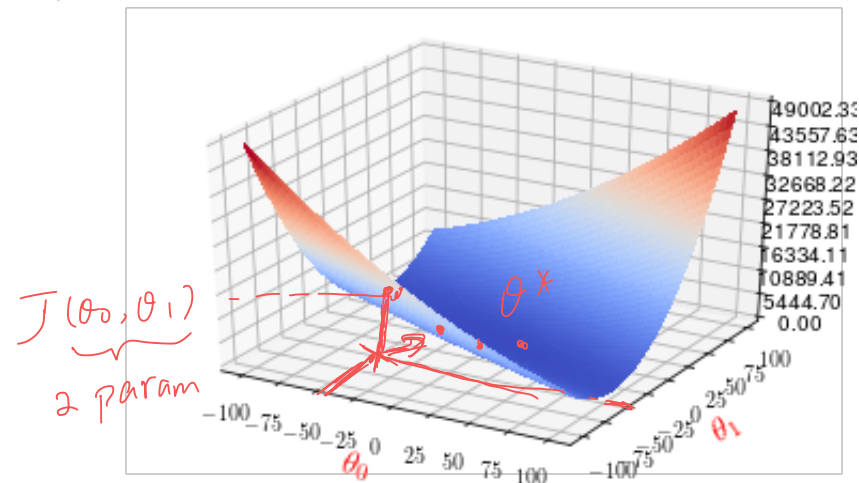
Gradient descent (minimal)

- The loss function is a scalar function of the vector θ

When there is just one parameter



When there are two parameters



loss

$$\left. \frac{\partial J}{\partial \theta} \right|_{\theta = \theta^{(0)}} = (\theta_1 - 1) \big|_{\theta_1 = \theta^{(0)}} = (-0.5 - 1) = -1.5$$

$$\theta^{(1)} = \theta^{(0)} - \eta \times \left. \frac{\partial J}{\partial \theta} \right|_{\theta = \theta^{(0)}}$$

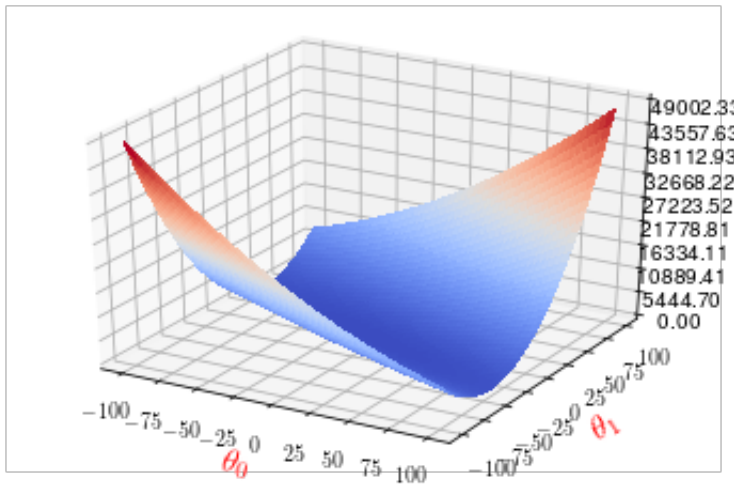
η = learning rate $\in (0, 0.1)$

Matrix calculus (minimal)

unigram LM
 $\theta = [\theta_1, \dots, \theta_{|V|}]$

- How to find the gradient of a scalar function with respect to a vector θ

When there are two parameters



Multi-variate function:

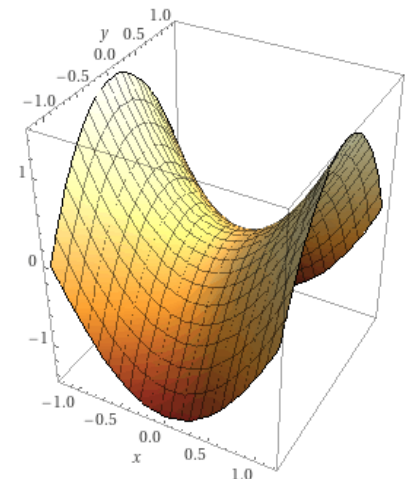
$$f(\mathbf{x}) = f(x_1, \dots, x_n)$$

Gradient:

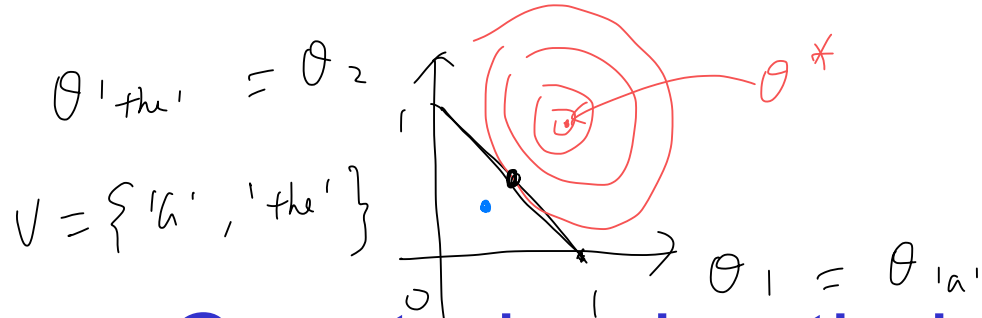
$$\nabla_{\mathbf{x}} f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n$$

$$\vec{x} \in \mathbb{R}^n$$

$$f(x, y) = x^2 - y^2$$



$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} \quad \text{HW1}$$



$$\begin{cases} \theta_1 + \theta_2 = 1 \\ \theta_1, \theta_2 \geq 0 \end{cases}$$

$$\min_{\theta} \max_{\lambda} L(\theta, \lambda)$$

Constrained optimization

(used in estimating HMM parameters)

- Maximum Likelihood Estimation (MLE)

$$\theta^* = \arg \min_{\theta} -\frac{1}{Z} \prod_{w=1}^{|V|} \theta_w^{c(w)} = \arg \min_{\theta} \log Z - \underbrace{\sum_{w=1}^{|V|} c(w) \log \theta_w}_{J(\theta)}$$

$$\text{s.t. } \sum_w \theta_w = 1$$

$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_{|V|} \end{bmatrix}$$

Take the gradient of the objective function with respect to θ

$$J(\theta) = \left[-\frac{c_1}{\theta_1}, -\frac{c_2}{\theta_2}, \dots, -\frac{c_{|V|}}{\theta_{|V|}} \right]$$

$$\nabla_{\theta} J(\theta) = \left[\frac{\partial J}{\partial \theta_1}, \dots, \frac{\partial J}{\partial \theta_{|V|}} \right] = \left[0 - c(1) \frac{1}{\theta_1}, \dots, 0 - c(|V|) \frac{1}{\theta_{|V|}} \right]$$

Can you descend to the optimal θ^* ? Don't forget the constraint. Need the Lagrangian method.

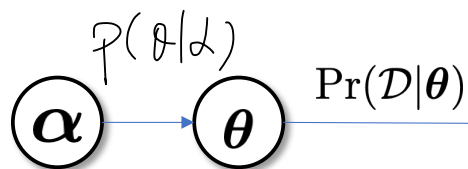
$$L(\theta, \lambda) = \log z - \sum_{w=1}^{|V|} c(w) \log \theta_w - \lambda \left[\sum_w \theta_w - 1 \right]$$

$$\frac{\partial L}{\partial \theta_w} = -\frac{c(w)}{\theta_w} - \lambda = 0 \Rightarrow \theta_w = -\frac{c(w)}{\lambda}$$

MLE of unigram $\Rightarrow \theta_w = \frac{c(w)}{\sum_w c(w)}$

$$\sum_{w=1}^{|V|} \theta_w = 1 \Rightarrow -\sum_{w=1}^{|V|} \frac{c(w)}{\lambda} = 1 \Rightarrow \lambda = -\sum_{w=1}^{|V|} c(w)$$

MAP



The corpus = \mathcal{D}

D1=[<s> I am Sam </s>]
 D2=[<s> Sam I am </s>]
 D3=[<s> I do not like green eggs and ham </s>]
 }

- MAP (Maximum a Posterior) estimation

- Assume a probabilistic model generating the data

- parametric model for **likelihood** of the data $\Pr(\mathcal{D}|\theta) \leftarrow$ multinomial dist
- can have **prior** distribution over the parameter $\Pr(\theta|\alpha)$

- Use the Bayes theorem to find the posterior

$$\Pr(\theta|\mathcal{D}) = \frac{\Pr(\mathcal{D}|\theta) \times \Pr(\theta|\alpha)}{\Pr(\mathcal{D})}$$

α_w : Hallucinated appearance of word w

- For convenience, the prior is a conjugate distribution of the likelihood.

- with multi-nomial, the conjugate is the Dirichlet distribution

$$\Pr(\theta|\alpha) = \frac{1}{Z} \prod_w \theta_w^{\alpha_w}$$

Apply "MLE" on $P(\theta|\mathcal{D})$


"Laplacian" $\Rightarrow \theta_w = [C(w) + \alpha_w] / \sum_w [C(w) + \alpha_w]$

$$P(\mathcal{D}|\theta) \times P(\theta|\alpha) = \frac{1}{Z} \prod_{w=1}^{|\mathcal{V}|} \theta_w^{C(w)} \times \frac{1}{Z'} \prod_{w=1}^{|\mathcal{V}|} \theta_w^{\alpha_w}$$

$$\alpha = \prod_{w=1}^{|\mathcal{V}|} \theta_w^{\alpha_w + C(w)}$$

MLE in NLP

- Looking ahead, we will use MLE to learn a few models

- Word embedding (glove) → word vectors
- Hidden Markov Model (HMM) → sequences
- Conditional random fields (CRF) → word vectors + seq.
- Reccurent NN → NN version of CRF
- Probabilistic CFG parser → 
- Recursive NN

→ NN on trees

Summary

- Maximum likelihood estimation
 - Probabilistic models: likelihood
- Tools:
 - gradient descent.
 - matrix calculus to find gradients.
 - constrained optimization.
- Maximum posterior estimation
 - Prior, conjugacy, and posterior using the Bayes theorem.
 - Laplacian smoothing as MAP