Maximum Entropy Classifier &

Supervised Machine Learning

CSE538 - Spring 2025

Topics

- Supervised classification (open-vocabulary)
- Goal of logistic regression (aka Maximum Entropy Classifier)
- The "loss function" -- what logistic regression tries to optimize
- Logistic regression with multiple features
- How to evaluation: Training and test datasets
- Overfitting: role of regularization

Text Classification

The Eagles win it!

Twitter to be acquired by Apple

Tariff delay as Trump and Trudeau reach deal



She <u>will</u> drive to the office, to make sure the lawyer gives the <u>will</u> to the family.

will.n or *will.v* ? verb noun

I like the the movie.

The movie is like terrible.

X - features of N observations (i.e. words)

Y - class of each of N observations

GOAL: Produce a *model* that outputs the most likely class y_i , given features x_i . f(X) = Y

X - features of N observations (i.e. words)

Y - class of each of N observations

GOAL: Produce a *model* that outputs the most likely class y_i , given features x_i . f(X) = Y

$$\begin{array}{cccccccc} i & X & Y \\ 0 & 0.0 \\ 1 & 0.5 \\ 2 & 1.0 \\ 3 & 0.25 \\ 4 & 0.75 \end{array} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ \end{array}$$

Supervised Classific

X - features of N observations (

Some function or rules to go from *X* to *Y*, as close as possible.

Y - class of each of N observation

GOAL: Produce a *model* that outputs the most likely class y_i , given features x_i .

f(X) = Y

$$\begin{array}{cccccccc} i & X & Y \\ 0 & 0.0 \\ 1 & 0.5 \\ 2 & 1.0 \\ 3 & 0.25 \\ 4 & 0.75 \end{array} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{array}$$

Supervised Machine Learning: Build a model with examples of outcomes (i.e. *Y*) that one is trying to predict. *unsupervised* machine learning: tries to learn with only an *X*).

Classification: The outcome (Y) is a discrete class. for example: $y \in \{\text{not-noun}, \text{noun}\}$

 $y \in \{\text{noun, verb, adjective, adverb}\}$

 $y \in \{\text{positive}_\text{sentiment}, \text{negative}_\text{sentiment}\}).$

Binary classification goal: Build a "model" that can estimate P(A=1|B=?)

i.e. given B, yield (or "predict") the probability that A=1

Binary classification goal: Build a "model" that can estimate P(A=1|B=?)

i.e. given B, yield (or "predict") the probability that A=1

In machine learning, the tradition is to use \mathbf{Y} for the variable being predicted and \mathbf{X} for the features use to make the prediction.

Binary classification goal: Build a "model" that can estimate P(Y=1|X=?)

i.e. given X, yield (or "predict") the probability that Y=1

In machine learning, the tradition is to use \mathbf{Y} for the variable being predicted and \mathbf{X} for the features use to make the prediction.

Binary classification goal: Build a "model" that can estimate P(Y=1|X=?)

i.e. given X, yield (or "predict") the probability that Y=1

In machine learning, the tradition is to use \mathbf{Y} for the variable being predicted and \mathbf{X} for the features use to make the prediction.

Example: Y: 1 if target is verb, 0 otherwise;X: 1 if "was" occurs before target; 0 otherwise

I was <u>reading</u> for NLP.

We were <u>fine</u>.

I am <u>good</u>.

The cat was <u>very</u> happy.

We enjoyed the <u>reading</u> material. I was <u>good</u>.

Binary classification goal: Build a "model" that can estimate P(Y=1|X=?)

i.e. given X, yield (or "predict") the probability that Y=1

In machine learning, the tradition is to use \mathbf{Y} for the variable being predicted and \mathbf{X} for the features use to make the prediction.

Example: Y: 1 if target is verb, 0 otherwise;X: 1 if "was" occurs before target; 0 otherwise

I was <u>reading</u> for NLP.

We were <u>fine</u>.

I am <u>good</u>.

The cat was <u>very</u> happy.

We enjoyed the <u>reading</u> material. I was <u>good</u>.

Example: Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

Example: Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.





Example:Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

x	У
2	1
1	0
0	0
6	1
2	1

Example: Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandals lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

X	У
2	1
1	0
0	0
6	1
2	1

5

3



Example: Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandals lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

X	у	
2	1	
1	0	
0	0	
6	1	
2	1	



Х

2

0

6

2

У

0

 \cap

Example:Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

They attend Binghamton.



Example:Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

They attend Binghamton.

X	У
2	1
1	0
0	0
6	1
2	1
1	1



1

1





1

1





1

1





Example:Y: 1 if target is a part of a proper noun, 0 otherwise;X: number of capital letters in target and surrounding words.

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

They attend Binghamton.

x	У
2	1
1	0
0	0
6	1
2	1
1	1



Example:Y: 1 if target is a part of a proper noun, 0 otherwise;X1: number of capital letters in target and surrounding words.Let's add a feature!X2: does the target word start with a capital letter?

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

They attend Binghamton.

x2	x1	у
1	2	1
0	1	0
0	0	0
1	6	1
1	2	1
1	1	1

 $Y_i \in \{0, 1\}$; X is a **single value** and can be anything numeric.

$$P(Y_i=1|X_i=x)=rac{1}{1+e^{-(eta_0+\sum_{j=1}^meta_jx_{ij})}}$$

m: number of features (columns of *x*)

 $Y_i \in \{0, 1\}$; X is a **single value** and can be anything numeric.

$$egin{aligned} P(Y_i = 1 | X_i = x) &= rac{1}{1 + e^{-(eta_0 + \sum_{j=1}^m eta_j x_{ij})}} \ &= rac{1}{1 + e^{-(x_i eta)}} \end{aligned}$$

m: number of features (columns of *x*)

 $Y_i \in \{0, 1\}$; X is a **single value** and can be anything numeric.

$$P(Y_i = 1 | X_i = x) = rac{1}{1 + e^{-(eta_0 + \sum_{j=1}^m eta_j x_{ij})}}$$

Vector notation

 β and x_i are vectors of size m

first feature is intercept: $x_{*,0} = [1, 1..., 1]_N$

$$rac{1}{1+e^{-(x_ieta)}}$$

m: number of features (columns of *x*)

 $Y_i \in \{0, 1\}$; X can be anything numeric.

$$P(Y_i = 1 | X_i = x) = rac{1}{1 + e^{-(x_i eta)}}$$

Goal: <u>take in the variable x</u> and <u>return a probability that Y is 1</u>.

 $Y_i \in \{0, 1\}$; X can be anything numeric.

$$P(Y_i = 1 | X_i = x) = rac{1}{1 + e^{-(x_i eta)}}$$

Goal: take in the variable *x* and return a probability that *Y* is 1.

Note that there are only two variables on the right: x_i , β

 $Y_i \in \{0, 1\}$; X can be anything numeric.

$$P(Y_i = 1 | X_i = x) = rac{1}{1 + e^{-(x_i eta)}}$$

Goal: take in the variable *x* and return a probability that *Y* is 1.

Note that there are only two variables on the right: x_i , β

 x_i is given. β must be <u>learned</u>.

 $Y_i \in \{0, 1\}$; X can be anything numeric.

$$P(Y_i=1|X_i=x)=rac{1}{1+e^{-\langle\!\! x_ieta
angle}}$$

HOW? Try different β values until "best fit" to the training data (example X and Y).



 x_i is given. β must be <u>learned</u>.

"best fit" : whatever maximizes the likelihood function:



 x_i is given. β must be <u>learned</u>.

"best fit" : whatever maximizes the *likelihood* function:

$$L(eta|X,Y) = \prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

"best fit" : whatever maximizes the *likelihood* function:
$$L(eta|X,Y)=\prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

"best fit" : more efficient to maximize *log likelihood* :
$$L(eta|X,Y) = \prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

"best fit" : more efficient to maximize *log likelihood* :

$$\ell(eta) = \sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i) \mathrm{log}(1-p_i) \ \boxed{p_i \equiv P(Y_i = 1 | X_i = x)}$$

$$L(eta|X,Y) = \prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

"best fit" : more efficient to maximize *log likelihood* :

$$\ell(eta) = \sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i) \mathrm{log}(1-p_i)$$

"best fit" for neural networks: software designed to **minimize** rather than maximize (typically, normalized by N, the number of examples.)

$$L(eta|X,Y) = \prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

"best fit" : more efficient to maximize *log likelihood* :

$$\ell(eta) = \sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i) \mathrm{log}(1-p_i)$$

"best fit" for neural networks: software designed to **minimize** rather than maximize (typically, normalized by N, number of examples.) "*log loss*" or *"normalized log loss":*

$$J(eta) = -rac{1}{N}\sum_{i=1}^N y_i ext{log}(p_i) + (1-y_i) ext{log}(1-p_i)$$

Maximum Entropy Classifier &

Supervised Machine Learning

CSE538 - Spring 2025 Adithya V Ganesan bit.ly/cse538sp25-211

Supervised Classification

X - features of N observations (i.e. words)

Y - class of each of N observations

GOAL: Produce a *model* that outputs the most likely class y_i , given features x_i . f(X) = Y

$$\begin{array}{cccccccc} i & X & Y \\ 0 & 0.0 \\ 1 & 0.5 \\ 2 & 1.0 \\ 3 & 0.25 \\ 4 & 0.75 \end{array} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{array}$$

Supervised Classific

X - features of N observations μ

Some function or rules to go from *X* to *Y*, as close as possible.

Y - class of each of N observatio.

GOAL: Produce a *model* that outputs the most likely class y_i , given features x_i .

f(X) = Y

$$\begin{array}{ccccc} i & X & Y \\ 0 & 0.0 \\ 1 & 0.5 \\ 2 & 1.0 \\ 3 & 0.25 \\ 4 & 0.75 \end{array} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ \end{array}$$

Supervised Classific

X - features of N observations

Some function or rules to go from *X* to *Y*, as close as possible.

Y - class of each of N observatio.

GOAL: Produce a *model* that outputs the most likely class y_i , given features x_i .

f(X) = Y

$$i X Y$$

$$0 0,0$$

$$1 0,5$$

$$2 1,0$$

$$3 0,25$$

$$4 0,25$$

$$1,0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1 0$$

$$1$$

Logistic Regression on features (x)

 $Y_i \in \{0, 1\}$; X can be anything numeric.

$$P(Y_i = 1 | X_i = x) = rac{1}{1 + e^{-(x_i eta)}}$$

Goal: <u>take in the variable x</u> and <u>return a probability that Y is 1</u>.

$$L(eta|X,Y) = \prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

$$L(eta|X,Y) = \prod_{i=1}^n P(Y_i=1|x_i)^{y_i}(1-P(Y_i=1|x_i))^{1-y_i}$$

"best fit" : more efficient to maximize *log likelihood* :

$$\ell(eta) = \sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i) \mathrm{log}(1-p_i)$$

"best fit" for neural networks: software designed to **minimize** rather than maximize (typically, normalized by N, number of examples.) "*log loss*" or "*normalized log loss*":

$$J(eta)=-rac{1}{N}\sum_{i=1}^N y_i ext{log}(p_i)+(1-y_i) ext{log}(1-p_i)$$



(typically, normalized by N, number of examples.) "log loss" or "normalized log loss":

$$egin{aligned} J(eta) = -rac{1}{N}\sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i)\mathrm{log}(1-p_i) \end{aligned}$$



"log loss" or "normalized log loss":

$$J(eta) = -rac{1}{N}\sum_{i=1}^N y_i {
m log}(p_i) + (1-y_i) {
m log}(1-p_i)$$



"log loss" or "normalized log loss":

$$J(eta) = -rac{1}{N}\sum_{i=1}^N y_i ext{log}(p_i) + (1-y_i) ext{log}(1-p_i)$$



$$egin{aligned} J(eta) = -rac{1}{N}\sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i)\mathrm{log}(1-p_i) \end{aligned}$$

$$egin{aligned} J(eta) = -rac{1}{N}\sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i)\mathrm{log}(1-p_i) \end{aligned}$$

$$\nabla_{\beta}J(\beta) \equiv \frac{1}{N}\sum_{i=1}^{N}(p_i - y_i)x_i.$$

$$egin{aligned} J(eta) = -rac{1}{N}\sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i)\mathrm{log}(1-p_i) \end{aligned}$$

$$\nabla_{\beta} J(\beta) = \frac{1}{N} \sum_{i=1}^{N} (p_i - y_i) x_i.$$
Update Step: $\beta_{\text{new}} = \beta_{\text{old}} - a^*$ grad



$$egin{aligned} J(eta) = -rac{1}{N}\sum_{i=1}^N y_i ext{log}(p_i) + (1-y_i) ext{log}(1-p_i) \end{aligned}$$

Demonstration of Logistic Regression with Gradient Descent: <u>bit.ly/cse538sp25-logreg-numpy</u>

Update Step:
$$\beta_{new} = \beta_{old} - a^*$$
 grad

- Number of capital letters surrounding: integer
- Begins with capital letter: {0, 1}
- Preceded by "the"? {0, 1}





- Number of capital letters surrounding: integer
- Begins with capital letter: {0, 1}
- Preceded by "the"? {0, 1}



- Number of capital letters surrounding: integer
- Begins with capital letter: {0, 1}
- Preceded by "the"? {0, 1}



Often we want to make a classification based on multiple features:

- Number of capital letters surrounding: integer
- Begins with capital letter: {0, 1}
- Preceded by "the"? {0, 1}



We're learning a linear (i.e. flat) *separating hyperplane*, but fitting it to a *logit* outcome.

(https://www.linkedin.com/pulse/predicting-outcomes-pr obabilities-logistic-regression-konstantinidis/)

Logistic Regression

 $Y_i \in \{0, 1\}$; X can be anything numeric.

$$logit(p_i) = log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \sum_{j=1}^m \beta_j x_{ij} = 0$$



We're still learning a linear -*separating hyperplane*, but fitting it to a *logit* outcome.

(https://www.linkedin.com/pulse/predicting-outcomes-pr obabilities-logistic-regression-konstantinidis/)

Logistic Regression



Example: Y: 1 if target is a part of a proper noun, 0 otherwise; X1: number of capital letters in target and surrounding words. Let's add a feature! X2: does the target word start with a capital letter?

They attend Stony Brook University. Next to the brook Gandalf lay thinking.

The trail was very stony. Her degree is from SUNY Stony Brook.

The Taylor Series was first described by Brook Taylor, the mathematician.

They attend Binghamton.

x2	x1	у
1	2	1
0	1	0
0	0	0
1	6	1
1	2	1
1	1	1

Terminology

 $\beta \approx$ weight \approx coefficient \approx parameters $\approx \Theta$

Logistic Regression ≈ Maximum Entropy Classifier

loss function \approx cost function

Stochastic Gradient Descent ≈ Optimizer

PyTorch Intro: Logistic Regression

1. Tensors

PyTorch: 1. Tensors



A multi-dimensional matrix

(i.stack.imgur.com)







A multi-dimensional matrix

A 2-d tensor is just a matrix. 1-d: vector 0-d: a constant / scalar

(i.stack.imgur.com)

Note: Linguistic ambiguity: Dimensions of a Tensor **#** Dimensions of a Matrix

PyTorch Intro: Logistic Regression

- 1. Tensors
- 2. Numeric functions as a graph/network (forward pass)

PyTorch Intro: Logistic Regression

- 1. Tensors
- 2. Numeric functions as a graph/network (forward pass)

Why is it specifically called forward pass? Does it mean there is a backward pass? More on this in the Back Propagation lecture in Part 2.

PyTorch: 2. Numeric functions as a graph/network (forward pass)

Efficient, high-level built-in linear algebra for neural network operations.

Can be conceptualized as a graph of operations on tensors (matrices):



PyTorch: 2. Numeric functions as a graph/network (forward pass)

Efficient, high-level built-in linear algebra for neural network operations.

Can be conceptualized as a graph of operations on tensors (matrices):

import torch
from torch import nn #predefined nodes

x = torch.Tensor(input) W = torch.random.randn(x.shape[1], 1) #weights z = torch.matmul(x, beta) yhat = nn.functional.sigmoid(z)





PyTorch: 2. Numeric functions as a graph/network (forward pass: defined in "forward" method of nn.Module)
PyTorch: 2. Numeric functions as a graph/network (forward pass: defined in "forward" method of nn.Module)

The nn.Module class registers the weights and biases in self.linear variable as Parameters of the model.

These weights and biases are stored as special tensors called the nn.Parameters, which allow them to have other member objects that support training of models.

PyTorch: 2. Numeric functions as a graph/network (forward pass: defined in "forward" method of nn.Module)

```
class LogReg(nn.Module):
```

```
def forward(self, X):
    #This is where the model itself is defined.
    #For logistic regression the model takes in X and returns
    #the results of a decision function
```

PyTorch: 2. Numeric functions as a graph/network (forward pass: defined in "forward" method of nn.Module)

```
class LogReg(nn.Module):
   def __init__(self, num_feats,
                learn rate = 0.01, device = torch.device("cpu") ):
       #the constructor; define any layer objects (e.g. Linear)
       super(LogReg, self). init ()
       self.linear = nn.Linear(num feats+1, 1)
    def forward(self, X):
        #This is where the model itself is defined.
        #For logistic regression the model takes in X and returns
        #the results of a decision function
        newX = torch.cat((X, torch.ones(X.shape[0], 1)), 1)
                                #add intercept
        z = self.linear(newX)
       return nn.functional.sigmoid(z)
                                #logistic function on the linear output
```

PyTorch Intro: Logistic Regression

- 1. Tensors
- 2. Numeric functions as a graph/network (forward pass)
- 3. Loss function (training loop)

PyTorch: 3. Loss Function (training loop)

```
#runs the training loop of pytorch model:
sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
def loss_func(ypred, y):
    return torch.mean(-torch.sum(y*torch.log(y_pred)))
             + torch.mean(-torch.sum((1-y)*torch.log(1-y_pred)))
#training loop:
for i in range(epochs):
    model.train()
    sgd.zero grad()
    #forward pass:
    vpred = model(X)
    loss = loss func(ypred, y)
    #backward: /(applies gradient descent)
    loss.backward()
    sgd.step()
    if i % 20 == 0:
        print(" epoch: %d, loss: %.5f" %(i, loss.item()))
```

PyTorch: 3. Loss Function (training loop)

```
#runs the training loop of pytorch model:
sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
loss_func = torch.nn.BCELoss() #computationally optimized for GPUs
#training loop:
for i in range(epochs):
    model.train()
    sgd.zero grad()
    #forward pass:
    ypred = model(X)
    loss = loss_func(ypred, y)
    #backward: /(applies gradient descent)
    loss.backward()
    sgd.step()
    if i % 20 == 0:
```

print(" epoch: %d, loss: %.5f" %(i, loss.item()))

PyTorch Intro: Logistic Regression

- 1. Tensors
- 2. Numeric functions as a graph/network (forward pass)
- 3. Loss function (training loop)
- 4. Autograd (backward pass)

PyTorch: 3. Loss Function (training loop)

```
#runs the training loop of pytorch model:
sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
loss_func = torch.mem_BCELoss()
```

```
#training loop:
for i in range(epochs):
   model.train()
   sgd.zero_grad()
   #forward pass:
   ypred = model(X)
   loss = loss_func(ypred, y)
   #backward: /(applies gradien
   loss.backward()
   sgd.step()
```

To Optimize Betas (all weights/parameters within the neural net):

Stochastic Gradient Descent (SGD)

-- optimize over one sample each iteration

Mini-Batch SDG:

--optimize over b samples each iteration

```
if i % 20 == 0:
    print(" epoch: %d, loss: %.5f" %(i, loss.item()))
```

```
#runs the training loop of pytorch model:
sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
loss_func = torch.nn.BCELoss()
```

```
#training loop:
for i in range(epochs):
    model.train()
    sgd.zero_grad()
    #forward pass:
    ypred = model(X)
    loss = loss_func(ypred, y)
    #backward: /(applies gradient descent
    loss.backward()
    sgd.step()
```

Calling loss.backward() will trigger the torch autograd to traverse the submodules within nn.Module class and compute the gradients.

The gradients wrt each parameter are stored under *.grad* object of the nn.Parameters type.

```
if i % 20 == 0:
    print(" epoch: %d, loss: %.5f" %(i, loss.item()))
```

```
#runs the training loop of pytorch model:
sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
loss_func = torch.nn.BCELoss()
```

```
#training loop:
for i in range(epochs):
    model.train()
    sgd.zero_grad()
    #forward pass:
    ypred = model(X)
    loss = loss_func(ypred, y)
    #backward: /(applies gradient descent
    loss.backward()
    sgd.step()
```

The optimizer object's .step() uses the gradients stored in model Parameters to carry out the weight updates.

if i % 20 == 0: print(" epoch: %d, loss: %.5f" %(i, loss.item()))

```
#runs the training loop of pytorch model:
    sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
    loss_func = torch.nn.BCELoss()
```

```
#training loop:
for i in range(epochs):
    model.train()
    sgd.zero_grad()
    #forward pass:
    ypred = model(X)
    loss = loss_func(ypred, y)
    #backward: /(applies gradient descent)
    loss.backward()
    sgd.step()
```

```
The optimizer object's .zero_grad() resets the gradients to 0.
```

```
if i % 20 == 0:
    print(" epoch: %d, loss: %.5f" %(i, loss.item()))
```

```
#runs the training loop of pytorch model:
sgd = torch.optim.SGD(model.parameters(), lr=learning_rate)
loss_func = torch.nn.BCELoss()
```

```
#training loop:
for i in range(epochs):
    model.train()
    sgd.zero_grad()
    #forward pass:
    ypred = model(X)
    loss = loss_func(ypred, y)
    #backward: /(applies gradient descent)
    loss.backward()
    sgd.step()
```

```
if i % 20 == 0:
    print(" epoch: %d, loss: %.5f" %(i, loss.item()))
```



PyTorch Intro: Logistic Regression

- 1. Tensors
- 2. Numeric functions as a graph/network (forward pass) nn.module object maps X to y_pred
- 3. Loss function (training loop)loop that evaluates *ypred* versus *y*
- 4. Autograd (backward pass)torch computation that updates the parameters





"Corpus"

raw data: sequences of characters



Feature Extraction

--pull out *observations*_and *feature vector* per observation.

"Corpus"

raw data: sequences of characters



Feature Extraction

--pull out <u>observations</u> and feature vector per observation. e.g.: words, sentences,

documents, users.

X	Y	
0.0 0 0.5 1 1.0 1 0.25 0 0.75 Det	0 0 1 0 1	trair
 0.35 1	 0	

iing

"Corpus"

raw data: sequences of characters

Feature Extraction

--pull out <u>observations</u>and <u>feature vector</u> per bbservation.

> e.g.: words, sentences, documents, users. 2

i

3

4

. . .

Ν

raw data: sequences of characters

"Corpus"

row of features; e.g.
→ number of capital letters
→ whether "I" was mentioned or not 0.0 0 0 0.5 1 0 training 1.0 1 1 0.25 0 0 0.75**Dota** 1 . . . 0.35 1 0

Feature Extraction

--pull out <u>observations</u>and <u>feature vector</u> per pbservation.

> e.g.: words, sentences, documents, users. 2

1

raw data: sequences of characters

"Corpus"

row of features; e.g.
number of capital letters
whether "I" was
mentioned or not
k features indicating whether k words were mentioned or not



Feature Extraction

Multi-hot Encoding

Each word gets an index in the vector
 C1 if present; 0 if not

raw data: sequences of characters of features; e.g.
→ number of capital letters
→ whether "I" was mentioned or not
→ k features indicating whether k words were mentioned or not

Data

Feature Extraction

Multi-hot Encoding

Each word gets an index in the vector
 1 if present; 0 if not

 Feature example: is word present in document?

 The book was interesting so I was happy .

characters

 \rightarrow whether "I" was

mentioned or not

k features indicating whether k words were mentioned or not

Feature Extraction

Multi-hot Encoding

Each word gets an index in the vector
 1 if present; 0 if not

Feature example: is word present in document?

The book was interesting so I was happy. Data [0, 1, 1, 0, 1, ..., 1, 0, 1, 1, 0, 1, ..., 0, 1, ..., 0, 1, ..., 0, 1, ..., 0, 1, 1, 0, 1, ..., 0, ..., 0

Feature Extraction

Multi-hot Encoding

Each word gets an index in the vector
 C1 if present; 0 if not

Feature example: is word present in document



Feature Extraction

Multi-hot Encoding

Each word gets an index in the vector
 ^C1 if present; 0 if not

Feature example: is previous word "the"?

The book was interesting so I was happy .

Data

 $\begin{bmatrix} 0, 1, 1, 0, 1, \dots, 1, 0, 1, 1, 0, 1, \\ \rightarrow \ k \ features \ Indicating$

whether k words were

mentioned or not

Feature Extraction

Multi-hot Encoding

Each word gets an index in the vector
 1 if present; 0 if not

Feature example: is previous word "the"?

→ k features Indicating whether k words were mentioned or not

The book was interesting so I was happy

1, 0, 1, 1,

Data

Feature Extraction

One-hot Encoding

Each word gets an index in the vector • All indices 0 except present word: Feature example: is previous word "the"? The book was interesting so I was happy Data [0, 1, 0, 0, 0, ..., 0, 0, 0, 0, 0, 0, 0]-> k features Indicating whether k words were mentioned or not

Feature Extraction

One-hot Encoding

Each word gets an index in the vector • All indices 0 except present word: Feature example: which is previous word? The book was interesting so I was happy. Data [0, 1, 0, 0, 0, ..., 0, 0, 0, 0, 0, 0, 0]01^k 0, 1, 0, 0, ..., 0, 0, 0, 0, 0, 0, 0, 10. Δ I k

Feature Extraction

One-hot Encoding

Each word gets an index in the vector • CAllpindices 0 except present word: Feature example: which is previous word? raw data: The book was interesting so I was happy. Data rograciers 0, 0, 0; 0, 0, ..., ..., 0, 0, 0, 0, 0, 01^k 0, ..., 0, 0, 0, 0, 0, 0, 0,01k

Feature Extraction

Multiple One-hot encodings for one observation (1) word before; (2) word after The book was interesting so I was happy. $[0, 0, 0, 0, 1, 0, ..., 0]^{k} [0, ..., 0, 1, 0, ..., 0]^{k}$

Feature Extraction

Multiple One-hot encodings for one observation (1) word before; (2) word after

The book was interesting so I was happy . $[0, 0, 0, 0, 1, 0, ..., 0]^{k} [0, ..., 0, 1, 0, ..., 0]^{k}$ = $[0, 0, 0, 0, 1, 0, ..., 0, 0, ..., 0, 1, 0, ..., 0]^{2k}$

Feature Extraction

<u>Multiple One-hot encodings for one observation</u> (1) word before; (2) word after; (3) percent capitals Corpus

The book was Interesting so I was happy.

 $[0, 0, 0, 0, 1, 0, ..., 0]^k [0, ..., 0, 1, 0, ..., 0]^k$

 $\begin{bmatrix} 0, 0, 0, 0, 1, 0, \dots, 0, 0, \dots, 0, 1, 0, \dots, 0 \end{bmatrix}^{2}$ $\begin{bmatrix} 0, 0, 0, 0, 1, 0, \dots, 0, 0, \dots, 0, 1, 0, \dots, 0 \end{bmatrix}^{2k+1}$



Machine Learning Goal: Generalize to new data



Machine Learning Goal: Generalize to new data



Logistic Regression - Regularization

0.5	0	0.6	1	0	0.25	1
0	0.5	0.3	0	0	0	1
0	0	1	1	1	0.5	0
0	0	0	0	1	1	0
0.25	1	1.25	1	0.1	2	1
0.5	0	0.6	1	0	0.25	1
------	-----	------	---	-----	------	---
0	0.5	0.3	0	0	0	1
0	0	1	1	1	0.5	0
0	0	0	0	1	1	0
0.25	1	1.25	1	0.1	2	1



 $1.2 + \left| -63^{*}x_{1} + \left| 179^{*}x_{2} + \left| 71^{*}x_{3} + \right| 18^{*}x_{4} + \left| -59^{*}x_{5} + \right| 19^{*}x_{6} \right| = logit(Y)$

X_1	<i>X</i> ₂		X		=	Y
0.5	0	0.6	1	0	0.25	1
0	0.5	0.3	0	0	0	1
0	0	1	1	1	0.5	0
0	0	0	0	1	1	0
0.25	1	1.25	1	0.1	2	1

 $1.2 + -63^{*}x_{1} + 179^{*}x_{2} + 71^{*}x_{3} + 18^{*}x_{4} + -59^{*}x_{5} + 19^{*}x_{6}$ = logit(Y)



 $1.2 + -63^{*}x_{1} + 179^{*}x_{2} + 71^{*}x_{3} + 18^{*}x_{4} + -59^{*}x_{5} + 19^{*}x_{6} = logit(Y)$

Overfitting (1-d non-linear example)



Overfitting (1-d non-linear example)



Underfit

(image credit: Scikit-learn; in practice data are rarely this clear)

Overfitting (1-d non-linear example)



(image credit: Scikit-learn; in practice data are rarely this clear)



 $1.2 + -63^{*}x_{1} + 179^{*}x_{2} + 71^{*}x_{3} + 18^{*}x_{4} + -59^{*}x_{5} + 19^{*}x_{6} = logit(Y)$







 $0 + 2^*x_1 + 2^*x_2$

= logit(Y)

L1 Regularization - "The Lasso"

Zeros out features by adding values that keep from perfectly fitting the data.

L1 Regularization - "The Lasso"

Zeros out features by adding values that keep from perfectly fitting the data.



L1 Regularization - "The Lasso" *Zeros out* features by adding values that keep from perfectly fitting the data.

$$L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1 - y_i}$$

set betas that maximize *L*



L1 Regularization - "The Lasso" Zeros out features by adding values that keep from perfectly fitting the data.

$$L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1 - y_i} - \frac{1}{C} \sum_{j=1}^m |\beta_j|$$

set betas that maximize *penalized L*

This is for likelihood

for log loss, would add the penalty



L1 Regularization - "The Lasso" Zeros out features by adding values that keep from perfectly fitting the data. $L(\beta_0, \beta_1, ..., \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i} - \frac{1}{C} \sum_{j=1}^m |\beta_j|$

set betas that maximize *penalized L*



L2 Regularization - "Ridge" Shrinks features by adding values that keep from perfectly fitting the data. $L(\beta_0, \beta_1, ..., \beta_k | X, Y) = \prod_{i=1}^{n} p(x_i)^{y_i} (1 - p(x_i))^{1-y_i} - \frac{1}{C} \sum_{j=1}^{m} \beta_j^2$

set betas that maximize *penalized* L



Machine Learning Goal: Generalize to new data



Machine Learning Goal: Generalize to new data



Logistic Regression - Review

- Probabilistic Classification: P(Y | X)
- Learn logistic curve based on example data
 - training + development + testing data
- Set betas based on maximizing the *likelihood* (or based on minimizing *log loss*)
 - "shifts" and "twists" the logistic curve
 - separation represented by hyperplane at 0.50
- Multivariate features: Multi-, One-hot encodings
- Overfitting and Regularization

Log Loss: $J(eta) = -rac{1}{N}\sum_{i=1}^N y_i \mathrm{log}(p_i) + (1-y_i) \mathrm{log}(1-p_i)$

Multiclass Log Loss:

???

Log Loss:

$$J(eta) = -rac{1}{N}\sum_{i=1}^N rac{class \ 1}{y_i \log(p_i)} + rac{class \ 0}{(1-y_i)\log(1-p_i)}$$

Multiclass Log Loss:

$$J = -\frac{1}{N} \sum_{i=1}^{N} ??$$

J(

Log Loss:

Multiclass Log Loss:

$$eta) = -rac{1}{N} \sum_{i=1}^{N} rac{class \ 1}{y_i \log(p_i)} + rac{class \ 0}{(1-y_i) \log(1-p_i)} \ J = -rac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} rac{y_i \log(p_i) + (1-y_i) \log(1-p_i)}{y_i \log(1-p_i)}$$

Log Loss:

Multiclass Log Loss:

$$J(\beta) = -\frac{1}{N} \sum_{i=1}^{N} \frac{class 1}{y_i \log(p_i)} + \frac{class 0}{(1 - y_i) \log(1 - p_i)}$$

s: $J = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} \frac{y_i \log p(x_{i,j})}{y_i \log p(x_{i,j})}$

Equivalently:

$$J = -rac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} y_i log \ (rac{1}{1+e^{-eta x}})$$

Log Loss:

Multiclass Log Loss:

$$J(eta) = -rac{1}{N} \sum_{i=1}^{N} rac{class \ 1}{y_i \log(p_i)} + rac{class \ 0}{(1 - y_i) \log(1 - p_i)}$$

s: $J = -rac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} rac{y_i \log p(x_{i,j})}{y_i \log p(x_{i,j})}$

Equivalently:

$$J = -rac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} y_i log \ (rac{1}{1+e^{-eta x}})$$

loss = torch.mean(-torch.sum(y*torch.log(y_pred))

return y_pred

Equivalently:

$$J = -rac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} y_i log \, (rac{1}{1+e^{-eta x}})$$

loss = torch.mean(-torch.sum(y*torch.log(y_pred))

Equivalently:

 $J = -rac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{|V|} y_i log \ (rac{1}{1+e^{-eta x}})$

not in prebuilt torch function

loss = torch.mean(-torch.sum(y*torch.log(y_pred))

```
def init (self, num feats, num_classes,
                learn rate = 0.01, device = torch.device("cpu") ):
       #the constructor; define any layer objects (e.g. Linear)
        super(MultiClassLogReg, self).__init__()
        self.linear = nn.Linear(num feats+1, 1-num classes)
    def forward(self, X):
        newX = torch.cat((X, torch.ones(X.shape[0], 1)), 1) #add intercept
        #logistic function on the linear output:
        #y_pred = 1/(1 + torch.exp(-self.linear(newX)))
        y pred = torch.log(1/(1 + torch.exp(-self.linear(newX))))
                 #or simply: nn.log_softmax(self.linear(newX))
        return y pred
loss_func = nn.NLLLoss()#expects log probabilities of each class
  not in prebuilt torch
```

function

Two equivalent options for multi-class in torch:

option 1: NLLLoss (easier to understand functions)

#in model/forward:
 return nn.log_softmax(self.linear(newX)) #log softmax is multiclass

#in loss/train:

loss_func = nn.NLLLoss() #negative log likelihood loss

option 2: CrossEntropyLoss (easier to code, obfuscates functions)

#in model/forward:
 return self.linear(newX) #only use linear if using cross-entropy loss

#in loss/train:

loss_func = nn.CrossEntropyLoss() #includes log softmax