



CS 760: Machine Learning

Parametric models

University of Wisconsin-Madison

Outline

- **Parametric models**

- Comparison with non-parametric models
- Generative vs. discriminative modeling

- **Logistic regression**

- Maximum likelihood estimation
- Setup

- **Naïve Bayes**

- Bayesian modeling, MAP vs. MLE
- Motivation / training / inference / smoothing
- Examples (Bernoulli, Multiclass, Gaussian)

Outline

- **Parametric models**

- Comparison with non-parametric models
- Generative vs. discriminative modeling

- **Logistic regression**

- Maximum likelihood estimation
- Setup

- **Naïve Bayes**

- Bayesian modeling, MAP vs. MLE
- Motivation / training / inference / smoothing
- Examples (Bernoulli, Multiclass, Gaussian)

Supervised Learning: Review

Problem setting

- Set of possible instances
- Unknown *target function*
- Set of *models* (a.k.a. *hypotheses*)

$$\mathcal{X}$$

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

$$\mathcal{H} = \{h | h : \mathcal{X} \rightarrow \mathcal{Y}\}$$

Get

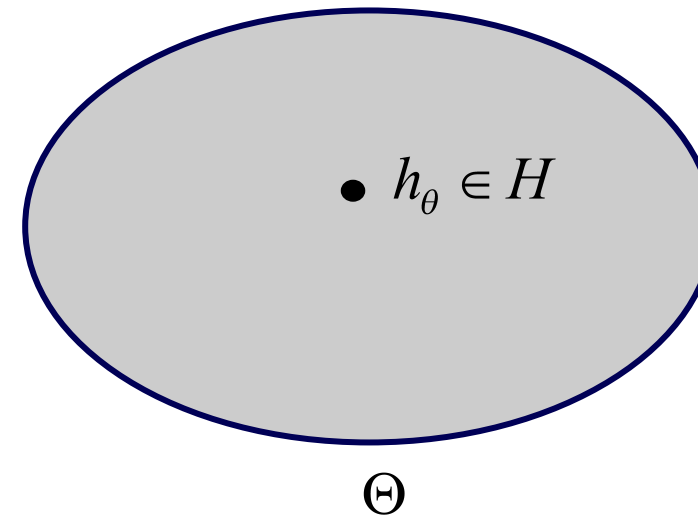
- Training set of instances for unknown target function f ,

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})$$

Goal: model h that best approximates f

Parametric models

- A way to categorize learning techniques
 - Parametric: hypotheses indexed by a fixed-dimensional (does not increase with training set size) **parameter vector**
 - Learning: find parameter yielding model that best approximates the target
 - **Ex:** linear models, neural networks
- Nonparametric methods:
 - Model complexity increases with training set size
 - eg kNN



Discriminative Models

- **Idea:** hypothesis h directly predicts the label (given features)
 - $y = h(x)$ or $p(y|x) = h(x)$
 - But does not model $p(x)$
- Examples
 - k-NN
 - decision trees
 - linear models

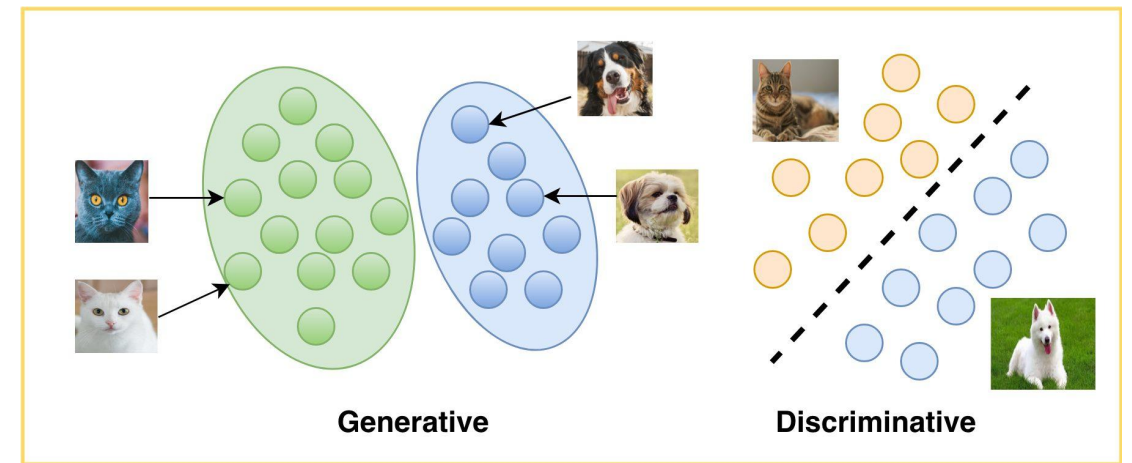
Generative Models

- Hypothesis h specifies a **joint distribution $p(\mathbf{x}, \mathbf{y})$** for how the labeled data was created. Recall $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y} | \mathbf{x})$
 - for unsupervised learning $h(\mathbf{x}) = p(\mathbf{x})$
- Select a hypothesis via maximum likelihood (or maximum a posteriori)
 - Ex: roll a die. Weights for each side define data generation
 - Observe training data to learn hypothesis



Discriminative vs Generative

- Can define both for supervised/unsupervised learning
 - k-means (discriminative-like) vs mixture-of-Gaussians (generative)
- When should we use one over the other?
 - Not always obvious
 - Discriminative models:
 - Often easier to optimize
 - Targets exact performance measure
 - Generative models:
 - Handling missing data
 - Generation via sampling



LearnOpenCV

- Typical examples:
 - Discriminative: linear regression, logistic regression, SVM, many neural networks (not all!)
 - Generative: Naïve Bayes, Bayesian Networks, ...

Outline

- **Parametric modeling**

- Comparison with non-parametric models
- Generative vs. discriminative modeling

- **Logistic regression**

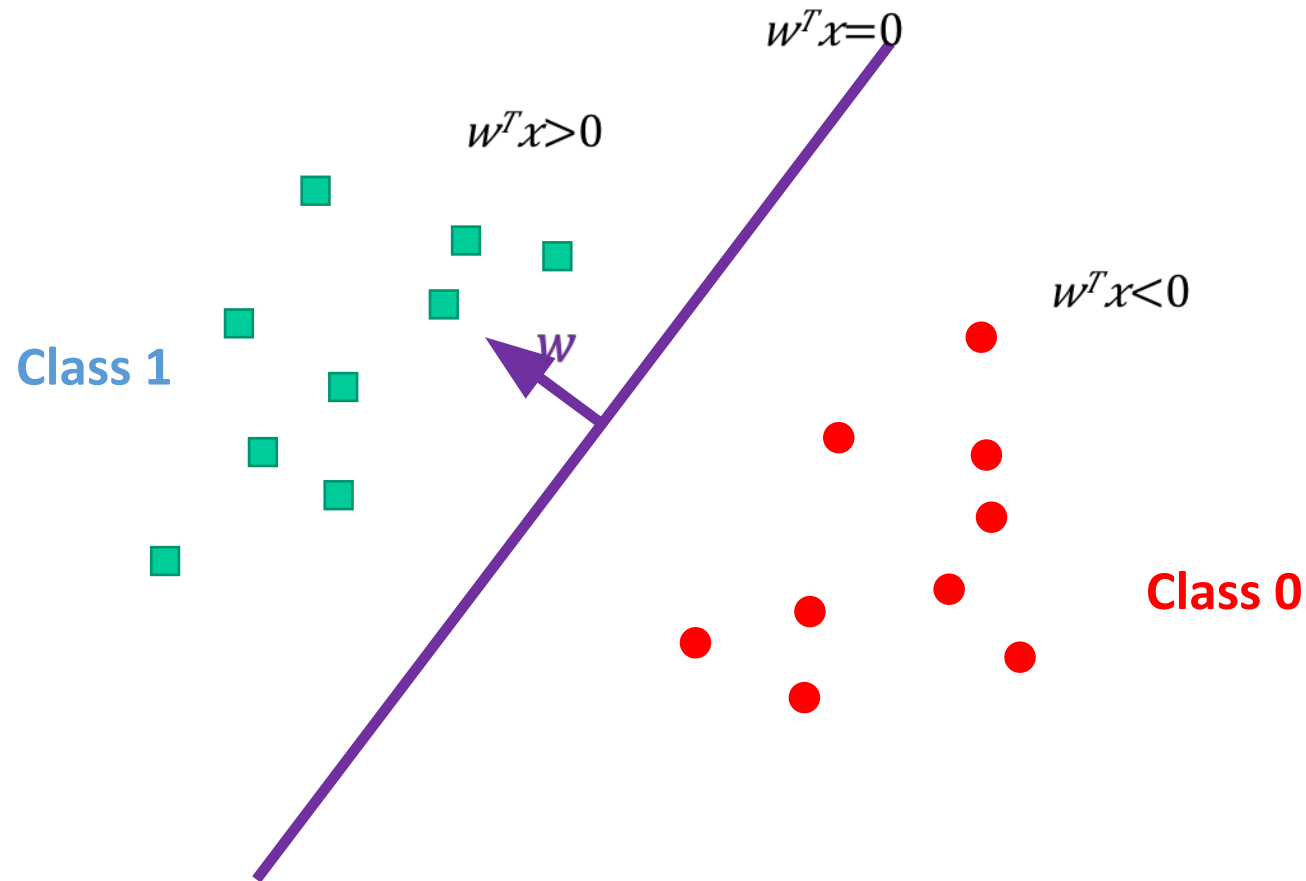
- Maximum likelihood estimation
- Setup

- **Naïve Bayes**

- Bayesian modeling, MAP vs. MLE
- Motivation / training / inference / smoothing
- Examples (Bernoulli, Multiclass, Gaussian)

Classification: Linear models

- How do we learn a linear separator between two classes?



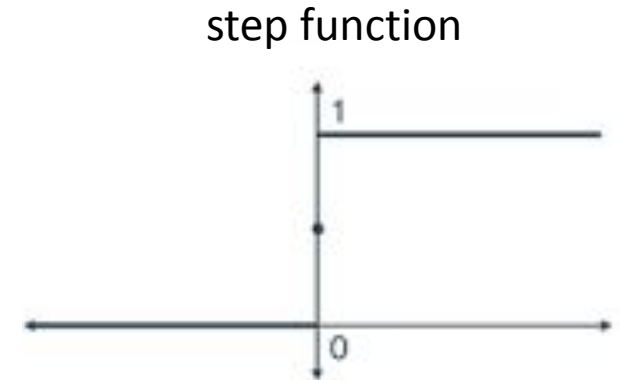
Linear Classification: Attempt 1

- Hyperplane: solutions to $\theta^T x = c$
 - note: $d-1$ dimensional -- d (degree of freedom) – 1 (constraint)
- So... try to use such hyperplanes as separators?

- Model: $f_{\theta}(x) = \theta^T x$

- Predict: $y=1$ if $\theta^T x \geq 0$, $y=0$ otherwise

- i.e. $y = \text{step}(f_{\theta}(x))$



- Training objective:

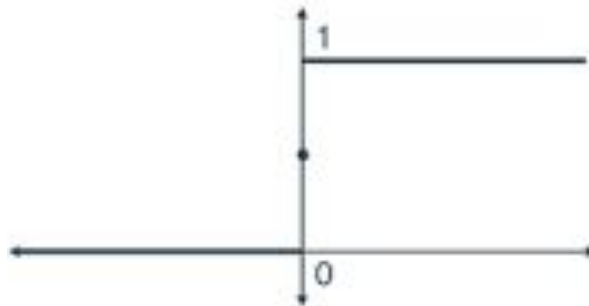
$$\ell(f_{\theta}) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} \left[\text{step}(f_{\theta}(x^{(i)})) \neq y^{(i)} \right]$$

difficult loss function to optimize!!

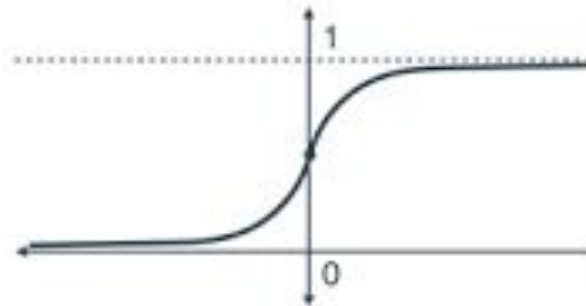
Linear Classification: Attempt 2

- Let us instead think probabilistically and learn $P_{\theta}(y|x)$ instead
- How?
 - Specify the conditional distribution $P_{\theta}(y|x)$
 - Use maximum likelihood estimation (MLE) to get a nicer loss
 - Run gradient descent (or related optimization algorithm)

step function



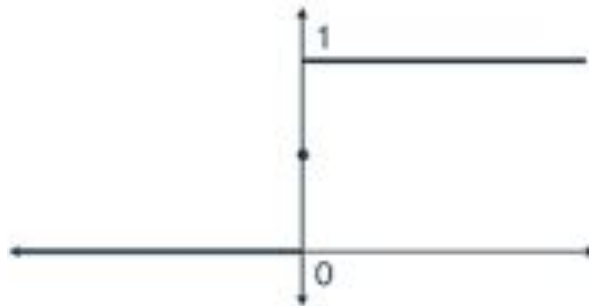
sigmoid function



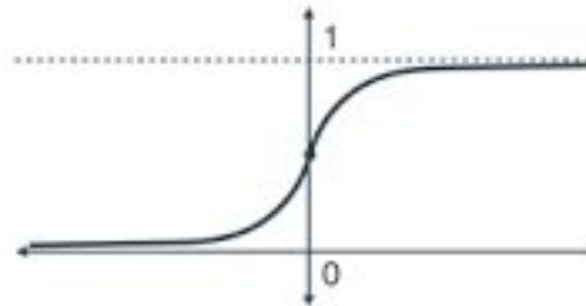
Linear Classification: Attempt 2

- Let us instead think probabilistically and learn $P_{\theta}(y|x)$ instead
- How?
 - Specify the conditional distribution $P_{\theta}(y|x)$
 - **Use maximum likelihood estimation (MLE) to get a nicer loss**
 - Run gradient descent (or related optimization algorithm)

step function



sigmoid function



Likelihood Function

- Captures the probability of seeing some data as a function of model parameters:

$$\mathcal{L}(\theta; X) = P_{\theta}(X)$$

- If data is iid, we have $\mathcal{L}(\theta; X) = \prod_j p_{\theta}(x_j)$
- Often more convenient to work with the log likelihood
 - Both mathematically and for numerical stability
 - Log is a monotonic + strictly increasing function

Maximum Likelihood Estimation

- For unsupervised learning, given data X , find the parameter that maximize the likelihood / log-likelihood

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta; X)$$

- Example: suppose we have n samples from a Bernoulli distribution

$$P_{\theta}(X = x) = \begin{cases} \theta & x = 1 \\ 1 - \theta & x = 0 \end{cases}$$

then

$$\mathcal{L}(\theta; X) = \prod_{i=1}^n P(X = x_i) = \theta^k (1 - \theta)^{n-k}$$

k : #samples with $x=1$

Maximum Likelihood: Example

- Want to maximize likelihood w.r.t. Θ

$$\mathcal{L}(\theta; X) = \prod_{i=1}^n P(X = x_i) = \theta^k (1 - \theta)^{n-k}$$

- Differentiate (use product rule) and set to 0. Get

$$\theta^{k-1} (1 - \theta)^{n-k-1} (k - n\theta) = 0$$

- So: ML estimate is

$$\hat{\theta} = \frac{k}{n}$$

Practice: how about maximum log-likelihood?

$$\log \mathcal{L}(\theta; X) = k \log \theta + (n - k) \log(1 - \theta)$$

$$\text{gradient} = 0 \Rightarrow \frac{k}{\theta} - \frac{n-k}{1-\theta} = 0$$

$$\hat{\theta} = \frac{k}{n}$$

ML: Conditional Likelihood

- Similar idea, but now using conditional probabilities:

$$\mathcal{L}(\theta; Y, X) = p_{\theta}(Y|X)$$

- If data is iid, we have

$$\mathcal{L}(\theta; Y, X) = \prod_j p_{\theta}(y_j|x_j)$$

- Now we can apply this to linear classification to get **logistic regression**.

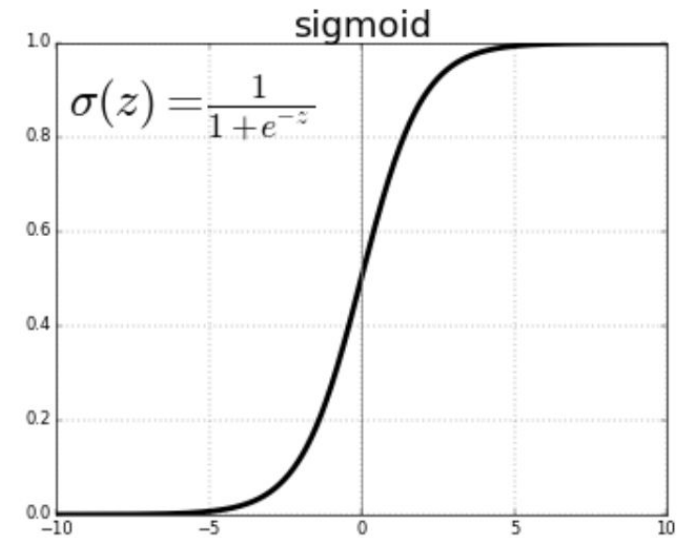
Logistic Regression: Conditional Distribution

•Notation: $\sigma(z) = \frac{1}{1 + \exp(-z)} = \frac{\exp(z)}{1 + \exp(z)}$ $z \leftarrow \theta^T x$

↑
sigmoid

“soft” version of step function

- Conditional distribution model for logistic regression:



$$P_{\theta}(y = 1|x) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$$

Logistic Regression: Loss

- Conditional MLE:

$$\log \text{likelihood}(\theta | x^{(i)}, y^{(i)}) = \log P_{\theta}(y^{(i)} | x^{(i)})$$

- So

$$\min_{\theta} \ell(f_{\theta}) = \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y^{(i)} | x^{(i)})$$

Or

$$\min_{\theta} -\frac{1}{n} \sum_{y^{(i)}=1} \log \sigma(\theta^T x^{(i)}) - \frac{1}{n} \sum_{y^{(i)}=0} \log(1 - \sigma(\theta^T x^{(i)}))$$

Logistic regression: Summary

- **logistic regression = sigmoid conditional distribution + MLE**

- More precisely:

- Give training data iid from some distribution D ,

- **Train:**

$$\min_{\theta} \ell(f_{\theta}) = \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y^{(i)} | x^{(i)})$$

- **Test:** output label probabilities

$$P_{\theta}(y = 1 | x) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$$

Logistic Regression: Comparisons

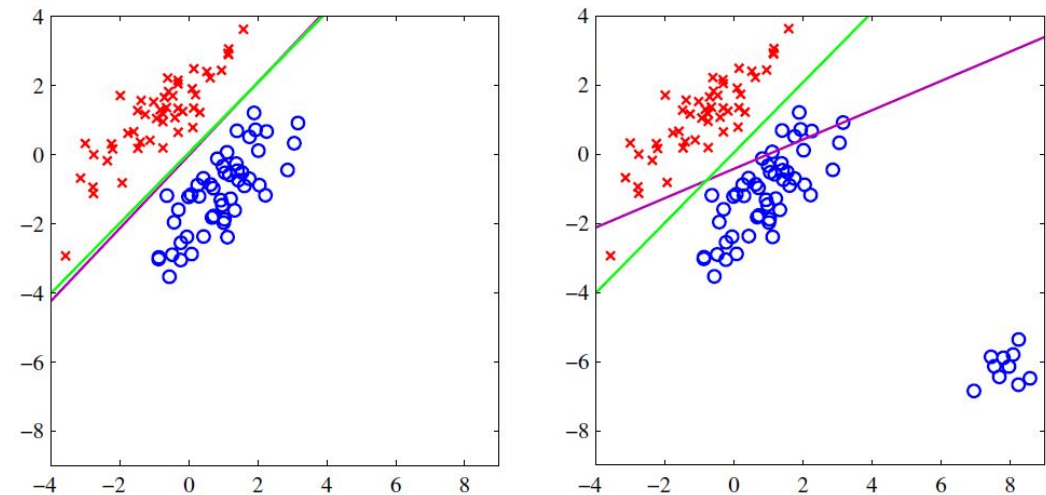
Recall the first attempt: $\ell(f_\theta) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [\text{step}(f_\theta(x^{(i)})) \neq y^{(i)}]$

difficult loss function to optimize!!

Another option: use the squared loss $(f_\theta(x^{(i)}) - y^{(i)})^2$

$$\ell(f_\theta) = \frac{1}{n} \sum_{j=1}^n (f_\theta(x^{(j)}) - y^{(j)})^2$$

this is just regular linear regression, it works surprisingly well for classification accuracy in practice, but has its problems



Logistic Regression: Beyond Binary

- We started with this conditional distribution:

$$P_{\theta}(y = 1|x) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$$

- Now let us try to extend it to multi-class classification, $y \in \{1, \dots, k\}$.
 - Can no longer just use one $\theta^T x$
 - But we can try multiple...

Logistic Regression: Beyond Binary

- Let's set, for y in $1, 2, \dots, k$

$$P_{\theta}(y = i | x) = \frac{\exp((\theta^i)^T x)}{\sum_{j=1}^k \exp((\theta^j)^T x)}$$

- Note: we have several weight vectors now (one per class).
- To train, same as before (just more weight vectors).

$$\min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y^{(i)} | x^{(i)})$$

Cross-Entropy Loss

- Define $q^{(i)}$ as the one-hot vector for the i th datapoint's label.

$$q_j^{(i)} = P(y = j | x^{(i)})$$

- Next, let's let $p^{(i)} = P_{\theta}(y|x^{(i)})$ be our prediction distribution (over y)

- Our loss terms can be written

$$-\log p(y^{(i)} | x^{(i)}) = - \sum_{j=1}^k q_j^{(i)} \log p(y = j | x^{(i)})$$

Note: only 1 term non-zero.

Looks like the entropy, but ...

- This is the “cross-entropy” $H(q^{(i)}, p^{(i)})$

Cross-Entropy Loss

- This is the “cross-entropy”

$$H(q^{(i)}, p^{(i)}) = \mathbb{E}_{q^{(i)}} [\log p^{(i)}]$$

- What are we doing when we minimize the cross-entropy?
- Recall KL divergence,

$$D(q^{(i)} || p^{(i)}) = \underbrace{\mathbb{E}_{q^{(i)}} [\log p^{(i)}]}_{\text{Cross-entropy}} - \underbrace{\mathbb{E}_{q^{(i)}} [\log q^{(i)}]}_{\text{Entropy } H(q^{(i)}) \text{ (fixed)}}$$

- Matching distributions!

Softmax

- We wrote

$$P_{\theta}(y = i|x) = \frac{\exp((\theta^i)^T x)}{\sum_{j=1}^k \exp((\theta^j)^T x)}$$

- This operation is called softmax.
 - Converts a vector into a probability vector (note normalization).
 - If one component in the vector ***a*** is **dominant**, softmax(***a***) is close to the one-hot vector picking out that maximum component



Quiz

Q1: Calculate the softmax of (1, 2, 3, 4, 5):

1. (0, 0.145, 0.229, 0.290, 0.336)
2. (0.012, 0.032, 0.086, 0.234, 0.636) ←
3. (0.636, 0.234, 0.086, 0.032, 0.012)

$$\text{softmax}(a)_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

Q2: True or false

Logistic regression is a discriminative model because we obtain $P_{\theta}(y|x)$.

False! We also obtain this from a generative model. Logistic regression is discriminative because it does not learn the joint distribution over (x,y) .

Outline

- **Parametric modeling**

- Comparison with non-parametric models
- Generative vs. discriminative modeling

- **Logistic regression**

- Maximum likelihood estimation
- Setup

- **Naïve Bayes**

- Bayesian modeling, MAP vs. MLE
- Motivation / training / inference / smoothing
- Examples (Bernoulli, Multiclass, Gaussian)

Review: Maximum Likelihood

- For some set of data, find the parameters that maximize the likelihood / log-likelihood

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta; X)$$

- Example: suppose we have n samples from a Bernoulli distribution

$$P_{\theta}(X = x) = \begin{cases} \theta & x = 1 \\ 1 - \theta & x = 0 \end{cases}$$

Then,

$$\mathcal{L}(\theta; X) = \prod_{i=1}^n P(X = x_i) = \theta^k (1 - \theta)^{n-k}$$

Review: Maximum Likelihood

- For some set of data, find the parameters that maximize the likelihood / log-likelihood
- Example: exponential distribution
 - pdf of Exponential(λ): $f(x) = \lambda e^{-\lambda x}$
 - Suppose $X_i \sim \text{Exponential}(\lambda)$ for $1 \leq i \leq N$.
 - Find MLE for data $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$
 - First write down log-likelihood of sample.
 - Compute first derivative, set to zero, solve for λ .
 - Compute second derivative and check that it is concave down at λ^{MLE} .

Review: Maximum Likelihood

- Example: exponential distribution
 - First write down log-likelihood of sample.

$$\ell(\lambda) = \sum_{i=1}^N \log f(x^{(i)}) \quad (1)$$

$$= \sum_{i=1}^N \log(\lambda \exp(-\lambda x^{(i)})) \quad (2)$$

$$= \sum_{i=1}^N \log(\lambda) + -\lambda x^{(i)} \quad (3)$$

$$= N \log(\lambda) - \lambda \sum_{i=1}^N x^{(i)} \quad (4)$$

Review: Maximum Likelihood

- Example: exponential distribution

- Compute first derivative, set to zero, solve for λ .

$$\frac{d\ell(\lambda)}{d\lambda} = \frac{d}{d\lambda} N \log(\lambda) - \lambda \sum_{i=1}^N x^{(i)} \quad (1)$$

$$= \frac{N}{\lambda} - \sum_{i=1}^N x^{(i)} = 0 \quad (2)$$

$$\Rightarrow \lambda^{\text{MLE}} = \frac{N}{\sum_{i=1}^N x^{(i)}} \quad (3)$$

Another Approach: **Bayesian Inference**

- Let's consider a different approach
- Need a little bit of terminology

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

- H is the hypothesis
- E is the evidence



Bayesian Inference Definitions

- Terminology:

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)} \longleftarrow \text{Prior}$$

- Prior: estimate of the probability **without** evidence

Bayesian Inference Definitions

- Terminology:

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

Likelihood
↙

- Likelihood: probability of evidence **given a hypothesis**.
 - Compare to the way we defined the likelihood earlier

Bayesian Inference Definitions

- Terminology:

$$\underset{\substack{\uparrow \\ \text{Posterior}}}{P(H|E)} = \frac{P(E|H)P(H)}{P(E)}$$

- Posterior: probability of hypothesis **given evidence**.

MAP Definition

- Suppose we think of the parameters as random variables
 - There is a prior $P(\theta)$

- Then, can do learning as Bayesian inference

- “Evidence” is the data

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

- **Maximum a posteriori probability (MAP)** estimation

$$\theta^{\text{MAP}} = \arg \max_{\theta} \prod_{i=1}^n p(x^{(i)}|\theta)p(\theta)$$

MAP vs ML

- What's the difference between ML and MAP?

$$\theta^{\text{MLE}} = \arg \max_{\theta} \prod_{i=1}^n p(x^{(i)} | \theta)$$

$$\theta^{\text{MAP}} = \arg \max_{\theta} \prod_{i=1}^n p(x^{(i)} | \theta) p(\theta)$$

- the prior!

Outline

- Parametric modeling
 - Comparison with non-parametric models
 - Generative vs. discriminative modeling
- Logistic regression
 - Maximum likelihood estimation
 - Setup
- **Naïve Bayes**
 - Bayesian modeling, MAP vs. MLE
 - **Motivation / training / inference / smoothing**
 - **Examples (Bernoulli, Multiclass, Gaussian)**

Application: Parody Detection

- The Economist

La paralización

Spain may be heading for its third election in a year

All latest updates

Stubborn Socialists are blocking Mariano Rajoy from forming a centre-right government

Sep 5th 2016 | MADRID | Europe

 Like 80  Tweet



EPA




BACK in June, after Spain's second indecisive election in six months, the general expectation was that Mariano Rajoy, the prime minister, would swiftly form a new government. Although his conservative People's Party (PP) did not win back the absolute majority it had lost in December, it remained easily the largest party, with 137 of the 350 seats in the Cortes (parliament), and was the only one to increase its share of the vote.


- The Onion

★ ELECTION 2016 ★ [MORE ELECTION COVERAGE ▶](#)

Tim Kaine Found Riding Conveyor Belt During Factory Campaign Stop

NEWS IN BRIEF
August 23, 2016
VOL 52 ISSUE 33
Politics · Politicians · Election 2016 · Tim Kaine



AIKEN, SC—Noting that he disappeared for over an hour during a campaign stop meet-and-greet with workers at a Bridgestone tire manufacturing plant, sources confirmed Tuesday that Democratic vice presidential candidate Tim Kaine was finally discovered riding on one of the factory's conveyor belts. "Shortly after we arrived, Tim managed to get out of our sight, but after an extensive search of the facilities, one of our interns found him moving down the assembly line between several radial tires," said senior campaign advisor Mike Henry, adding that Kaine could be seen smiling and laughing as he belted his way down the line. "It was a bit of a surprise, but we were glad to find him." Kaine was found on the conveyor belt at approximately 10:30 a.m. Tuesday. He was wearing a dark suit and tie, and was seen from behind as he rode away from the camera. He was holding a small object in his hands, which he was looking at intently. He was also holding a small object in his hands, which he was looking at intently. He was also holding a small object in his hands, which he was looking at intently.

Model 0: Not-Naïve Model

Generative story:

1. Flip a weighted coin (Y)
2. If heads, sample a document ID (X) from the parody distribution
3. If tails, sample a document ID (X) from the regular distribution

$$P(X, Y) = P(X|Y)P(Y)$$

Model 0: Not-Naïve Model

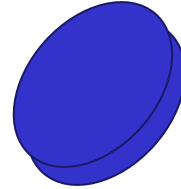
Generative story:

1. Flip a weighted coin (Y)
2. If heads, roll the **gray** many sided die to sample a document vector (X) from the parody distribution
3. If tails, roll the **blue** many sided die to sample a document vector (X) from the regular distribution

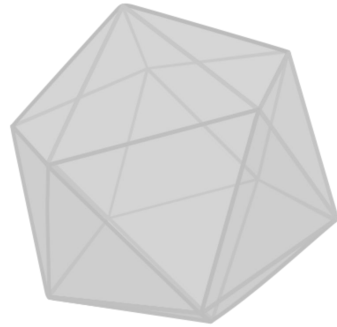
$$P(X_1, \dots, X_K, Y) = P(X_1, \dots, X_K | Y) P(Y)$$

Model 0: Not-Naïve Model

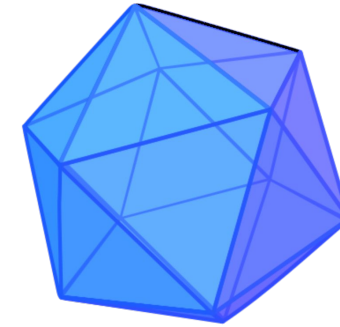
Flip weighted coin



If HEADS, roll
gray die



If TAILS, roll
blue die



Each side of the die
is labeled with a
document vector
(e.g. $[1, 0, 1, \dots, 1]$)

y	x_1	x_2	x_3	\dots	x_K
0	1	0	1	\dots	1
1	0	1	0	\dots	1
1	1	1	1	\dots	1
0	0	0	1	\dots	1
0	1	0	1	\dots	0
1	1	0	1	\dots	0

Model 0: Main Problem

How many terms are we modeling?

- Say features are binary: $X_i \in \{0, 1\}$ e.g. is word i in the document?

$$P(X_1, \dots, X_K | Y)$$

- 2^k choices of feature vector, each gets its own probability...
 - Exponentially big table (e.g. in vocabulary size)

Naïve Bayes: Core Assumption

How do we fix this problem?

- ASSUME conditional **independence** of features:

$$\begin{aligned} P(X_1, \dots, X_K, Y) &= P(X_1, \dots, X_K | Y) P(Y) \\ &= \left(\prod_{k=1}^K P(X_k | Y) \right) P(Y) \end{aligned}$$

- What do we gain? With binary features, get 2 entries per feature
- So, number of parameters

$$2^k \rightarrow 2k$$

Naïve Bayes: Overall Model

Support of $P(X_k|Y)$ depends on the application

Model: Product of **class prior** and the event model

$$P(\mathbf{X}, Y) = P(Y) \prod_{k=1}^K P(X_k|Y)$$

Training: Find the **class-conditional** MLE parameters:

- For prior $P(Y)$, we find the MLE using the data.
- For each $P(X_k|Y)$ we condition on the data with the corresponding class.

Prediction: Find the class that maximizes the posterior

$$\hat{y} = \operatorname{argmax}_y p(y|\mathbf{x})$$


Naïve Bayes: Smoothing

- **Training:** empirically estimate the probabilities
 - Could we just use counts to obtain the probabilities?


$$P(X_k = i|Y) = \frac{\text{\#times } X_k = i \text{ in class } Y}{\text{\#examples of class } Y} = \frac{c_i}{N}$$

- But what if $c_i = 0$? Then $P(X_k = i|Y) = 0$, which will make predictions using the event model zero ($= P(Y) \prod_k P(X_k = i|Y)$).

- Solution: smooth!
$$\hat{P}(X_k|Y) = \frac{c_i + \alpha}{N + \alpha m}$$



Smoothing
parameter



of classes

Naïve Bayes: Predicting

- With conditional probabilities, how to predict?

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_y p(y|\mathbf{x}) \quad (\text{posterior}) \\ &= \operatorname{argmax}_y \frac{p(\mathbf{x}|y)p(y)}{p(x)} \quad (\text{by Bayes' rule}) \\ &= \operatorname{argmax}_y p(\mathbf{x}|y)p(y)\end{aligned}$$

Naïve Bayes Example 1: Bernoulli

Support: Binary vectors of length K

$$\mathbf{x} \in \{0, 1\}^K$$

Generative Story:

$$Y \sim \text{Bernoulli}(\phi)$$

$$X_k \sim \text{Bernoulli}(\theta_{k,Y}) \quad \forall k \in \{1, \dots, K\}$$

Model: $p_{\phi, \boldsymbol{\theta}}(\mathbf{x}, y) = p_{\phi, \boldsymbol{\theta}}(x_1, \dots, x_K, y)$

$$= p_{\phi}(y) \prod_{k=1}^K p_{\boldsymbol{\theta}_k}(x_k | y)$$

$$= (\phi)^y (1 - \phi)^{(1-y)} \prod_{k=1}^K (\theta_{k,y})^{x_k} (1 - \theta_{k,y})^{(1-x_k)}$$

Training Bernoulli Naïve Bayes

Recall: train (by MLE) is to find **class-conditional** parameters

- To find $P(Y)$: use all the data
- For $P(X_i | Y=y)$: use the data for that class

$$\phi = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 1)}{N}$$

$$\theta_{k,0} = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 0 \wedge x_k^{(i)} = 1)}{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 0)}$$

$$\theta_{k,1} = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 1 \wedge x_k^{(i)} = 1)}{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 1)}$$

$$\forall k \in \{1, \dots, K\}$$

Naïve Bayes Example 2: Multinomial

Integer vector (word IDs)

$\mathbf{x} = [x_1, x_2, \dots, x_M]$ where $x_m \in \{1, \dots, K\}$ a word id.

Generative Story:

for $i \in \{1, \dots, N\}$:

$y^{(i)} \sim \text{Bernoulli}(\phi)$

for $j \in \{1, \dots, M_i\}$:

$x_j^{(i)} \sim \text{Multinomial}(\boldsymbol{\theta}_{y^{(i)}}, 1)$

Model:

$$\begin{aligned} p_{\phi, \boldsymbol{\theta}}(\mathbf{x}, y) &= p_{\phi}(y) \prod_{k=1}^K p_{\boldsymbol{\theta}_k}(x_k | y) \\ &= (\phi)^y (1 - \phi)^{(1-y)} \prod_{j=1}^{M_i} \theta_{y, x_j} \end{aligned}$$

Naïve Bayes Example 3: Gaussian

Support:

$$\mathbf{x} \in \mathbb{R}^K$$

Model: Product of **prior** and the event model

$$\begin{aligned} p(\mathbf{x}, y) &= p(x_1, \dots, x_K, y) \\ &= p(y) \prod_{k=1}^K p(x_k | y) \end{aligned}$$

Gaussian Naive Bayes assumes that $p(x_k | y)$ is given by a Normal distribution. (diagonal covariance matrix)

Q2-1: Are these statements true or false?

(A) Naïve Bayes assumes conditional independence of features to decompose the joint probability into the conditional probabilities.

(B) We use the Bayes' rule to calculate the posterior probability.

1. True, True
2. True, False
3. False, True
4. False, False

Q2-1: Are these statements true or false?

(A) Naïve Bayes assumes conditional independence of features to decompose the joint probability into the conditional probabilities.

(B) We use the Bayes' rule to calculate the posterior probability.

1. True, True 

2. True, False

3. False, True

4. False, False

(A) Just as we learnt in the lecture.

(B) We use Bayes rule to decompose posterior probability into prior probability and conditional probability given each class, so that we can compute it using the estimated parameters.



Thanks Everyone!

Some of the slides in these lectures have been adapted/borrowed from materials developed by Misha Khodak, Mark Craven, David Page, Jude Shavlik, Tom Mitchell, Nina Balcan, Elad Hazan, Tom Dietterich, Pedro Domingos, Jerry Zhu, Yingyu Liang, Volodymyr Kuleshov, Fred Sala, Tengyang Xie