# CS-UY 4563: Lecture 7 The Bayesian Perspective cont., Linear Classifiers

NYU Tandon School of Engineering, Prof. Christopher Musco

In a <u>Bayesian</u> or <u>Probabilistic</u> approach to machine learning we always start by conjecturing a

# probabilistic model

that plausibly could have generated our data.

- The model guides how we make predictions.
- The model typically has unknown parameters  $\vec{\theta}$  and we try to find the most reasonable parameters based on observed data (more on this later in lecture).

**Exercise:** With a partner, come up with a probabilistic model for <u>any one</u> of the following data sets  $(x_1, y_1), \ldots, (x_n, y_n)$ .

- 1. For *n* **people**: each  $x_i \in \{0, 1\}$  with zero indicating <u>male</u>, one indicating <u>female</u>. Each  $y_i$  is the height of the person in inches.
- 2. For *n* **NYC apartments**: each *x<sub>i</sub>* is the size of the apartment in square feet. Each *y<sub>i</sub>* is the monthly rent in dollars.
- For *n* students: each x<sub>i</sub> ∈ {Fresh., Soph., Jun., Sen.} indicating class year. Each y<sub>1</sub> ∈ {0,1} with zero indicating the student has not taken machine learning, one indicating they have.

Record any unknown parameters of your model. What would be a guess for their values? How would you confirm or refine this guess using data? **Dataset:**  $(x_1, y_1), \ldots, (x_n, y_n)$ 

**Description**: For *n* **people**: each  $x_i \in \{0, 1\}$  with zero indicating <u>male</u>, one indicating <u>female</u>. Each  $y_i$  is the height of the person in inches.

Model:

**Dataset:**  $(x_1, y_1), \ldots, (x_n, y_n)$ 

**Description**: For *n* **NYC apartments**: each  $x_i$  is the size of the apartment in square feet. Each  $y_i$  is the monthly rent in dollars.

Model:

**Dataset:**  $(x_1, y_1), \ldots, (x_n, y_n)$ 

Description: For *n* students: each

 $x_i \in \{Fresh., Soph., Jun., Sen.\}$  indicating class year. Each  $y_1 \in \{0, 1\}$  with zero indicating the student has not taken machine learning, one indicating they have.

Model:

## Goal:

- Build a probabilistic model for a binary classification problem.
- Estimate parameters of the model.
- From the model derive a classification rule for future predictions (the Naive Bayes Classifier).

#### SPAM PREDICTION



Both target labels and data vectors are binary.

#### PROBABILISTIC MODEL FOR EMAIL

**Probabilistic model** for (bag-of-words, label) pair (**x**, *y*):

- Set y = 0 with probability p(y = 0), y = 1 with probability p(y = 1) = 1 p(y = 0).
  - p(y = 0) is probability an email is not spam (e.g. 99%).
  - p(y = 1) is probability an email is spam (e.g. 1%).
- If y = 0, for each *i*, set  $x_i = 1$  with prob.  $p(x_i = 1 | y = 0)$ .
- If y = 1, for each *i*, set  $x_i = 1$  with prob.  $p(x_i = 1 | y = 1)$ .

Unknown model parameters:

• 
$$p(y = 0), p(y = 1),$$

• 
$$p(x_1 = 1 | y = 0), \dots, p(x_n = 1 | y = 0).$$

•  $p(x_1 = 1 | y = 1), \dots, p(x_n = 1 | y = 1).$ 

#### How would you estimate these parameters?

#### Reasonable way to set parameters:

- Set p(y = 0) and p(y = 1) to the empirical fraction of not spam/spam emails.
- For each word *i*, set  $p(x_i = 1 | y = 0)$  to the empirical probability word *i* appears in a <u>non-spam</u> email.
- For each word *i*, set  $p(x_i = 1 | y = 1)$  to the empirical probability word *i* appears in a <u>spam</u> email.

Estimating these parameters is the only "training" we will do.

DONE WITH MODELING ON TO PREDICTION

#### CLASSIFICATION RULE

Given unlabeled input (x, \_\_\_\_), choose the label  $y \in \{0, 1\}$  which is <u>most likely</u> given the data. Recall  $\mathbf{x} = [0, 0, 1, ..., 1, 0]$ .

Classification rule: maximum a posterior prob. (MAP) estimate.

Step 1. Compute:

- $p(y = 0 | \mathbf{x})$ : prob. y = 0 given observed data vector  $\mathbf{x}$ .
- $p(y = 1 | \mathbf{x})$ : prob. y = 1 given observed data vector  $\mathbf{x}$ .

Step 2. Output: 0 or 1 depending on which probability is larger.

 $p(y = 0 | \mathbf{x})$  and  $p(y = 1 | \mathbf{x})$  are called **posterior** probabilities.

How to compute the posterior? Bayes rule!

K

$$p(y = 0 | \mathbf{x}) = \frac{p(\mathbf{x} | y = 0)p(y = 0)}{p(\mathbf{x})}$$
(1)

posterior = 
$$\frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$
 (2)

- **Prior:** Probability in class 0 prior to seeing any data.
- **Posterior:** Probability in class 0 <u>after</u> seeing the data.

Goal is to determine which is larger:

$$p(y = 0 | \mathbf{x}) = \frac{p(\mathbf{x} | y = 0)p(y = 0)}{p(\mathbf{x})}$$
vs.  
$$p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | y = 1)p(y = 1)}{p(\mathbf{x})}$$

#### How to compute posteriors:

- Ignore evidence  $p(\mathbf{x})$  since it is the same for both sides.
- p(y = 0) and p(y = 1) already known (computed from training data).

• 
$$p(\mathbf{x} \mid y = 0) = ? p(\mathbf{x} \mid y = 1) = ?$$

"Naive" Bayes Rule: Compute  $p(\mathbf{x} | y = 0)$  by assuming independence:

$$p(\mathbf{x} \mid y = 0) = p(x_1 \mid y = 0) \cdot p(x_2 \mid y = 0) \cdot \ldots \cdot p(x_n \mid y = 0)$$

•  $p(x_i | y = 0)$  is the probability you observe  $x_i$  given that an email is not spam.<sup>1</sup>

A more complicated method might take dependencies into account.

<sup>&</sup>lt;sup>1</sup>Recall, *x<sub>i</sub>* is either 0 when word *i* is not present, or 1 when word *i* is present.

## Final Naive Bayes Classifier

Training/Modeling: Use existing data to compute:

• 
$$p(y = 0), p(y = 1)$$

- For all *i*:
  - Compute  $p(0 \text{ at position } i | y = 0), p(1 \text{ at position } i | y_0)$
  - Compute p(0 at position i | y = 1), p(1 at position i | y = 1)

# Prediction:

- For all *i*:
  - Compute  $p(\mathbf{x} \mid y = 0) = \prod_i p(x_i \mid y = 0)$
  - Compute  $p(\mathbf{x} \mid y = 1) = \prod_i p(x_i \mid y = 1)$
- Return

arg max 
$$[p(\mathbf{x} | y = 0) \cdot p(y = 0), p(\mathbf{x} | y = 1) \cdot p(y = 1)].$$

OTHER APPLICATIONS OF THE BAYESIAN PERSPECTIVE The Bayesian view offers an interesting alternative perspective on <u>many</u> machine learning techniques.

Example: Linear Regression.

Probabilistic model:

$$\mathbf{y}_i = \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \eta$$

where the  $\eta$  drawn from  $N(0, \sigma^2)$  is random Gaussian noise.



$$Pr(\eta = z) \sim$$

The symbol  $\sim$  means "is proportional to".

Example: Linear Regression.

Probabilistic model:

 $y_i = \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \eta$ 

where the  $\eta$  drawn from  $N(0, \sigma^2)$  is random Gaussian noise.

If we knew  $\beta$  what is the <u>maximum a posterior (MAP)</u> estimate for  $y_i$  given observed data  $\mathbf{x}_i$ ?

## How should be select ${m eta}$ for our model?

**Bayesian approach:** Use MAP estimate again! Now for parameter vector.

Choose  $\boldsymbol{\beta}$  to maximize:

$$\Pr(\beta \mid X, y) = \frac{\Pr(X, y \mid \beta) \Pr(\beta)}{\Pr(X, y)}.$$

Assume all  $\beta$ 's are equally likely, so we only care about  $Pr(\mathbf{X}, \mathbf{y} \mid \beta)$  when maximizing.

Choose  $\beta$  to maximize:

 $\Pr(\mathbf{X},\mathbf{y}\mid\boldsymbol{\beta}) \sim$ 

#### **BAYESIAN REGRESSION**

• 
$$y_i = \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \eta$$

• where 
$$p(\eta = z) \sim e^{-z^2/\sigma^2}$$

$$\mathsf{Pr}(\mathsf{X},\mathsf{y} \mid oldsymbol{eta}) \sim$$

#### LOG LIKELIHOOD

Easier to work with the log likelihood:

$$\arg \max_{\beta} \Pr(\mathbf{X}, \mathbf{y} \mid \beta) = \arg \max_{\beta} \prod_{i=1}^{n} e^{-(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / \sigma^2}$$
$$= \arg \max_{\beta} \log \left( \prod_{i=1}^{n} e^{-(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / \sigma^2} \right)$$
$$= \arg \max_{\beta} \sum_{i=1}^{n} -(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / \sigma^2$$
$$= \arg \min_{\beta} \sum_{i=1}^{n} (y_i - \langle \mathbf{x}_i, \beta \rangle)^2.$$

Choose  $\beta$  to minimize  $\sum_{i=1}^{n} (y_i - \langle \mathbf{x}_i, \beta \rangle)^2 = \|\mathbf{y} - \mathbf{X}\beta\|_2^2!$ This is a completely different justification for squared loss.

#### **BAYESIAN REGRESSION**

If we had modeled our noise  $\eta$  as Laplace noise, we would have found that minimizing  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_1$  was optimal.



Laplace noise has "heavier tails", meaning that it results in more outliers.

This is a completely different justification for  $\ell_1$  loss.

Recall goal is to maximize over  $\beta$ :

$$\Pr(\beta \mid \mathbf{X}, \mathbf{y}) = \frac{\Pr(\mathbf{X}, \mathbf{y} \mid \beta) \Pr(\beta)}{\Pr(\mathbf{X}, \mathbf{y})}.$$

assume all  $\beta$ 's equally likely

**Bayesian view:** Assume values in  $\beta = [\beta_1, \dots, \beta_d]$  come from some distribution.

- **Common model:** Each  $\beta_i$  drawn from  $N(0, \gamma^2)$ , i.e. normally distributed, independent.
- Encodes a belief that we are unlikely to see models with very large coefficients.

**Goal:** choose  $\beta$  to maximize:

$$\Pr(\beta \mid \mathbf{X}, \mathbf{y}) = \frac{\Pr(\mathbf{X}, \mathbf{y} \mid \beta) \Pr(\beta)}{\Pr(\mathbf{X}, \mathbf{y})}.$$

- We can still ignore the "evidence" term Pr(X, y) since it is a constant that does not depend on β.
- $\Pr(\beta) = \Pr(\beta_1) \cdot \Pr(\beta_2) \cdot \ldots \cdot \Pr(\beta_d)$
- $\Pr(m{eta}) \sim$

#### BAYESIAN REGULARIZATION

Easier to work with the log likelihood:

$$\arg \max_{\beta} \Pr(\mathbf{X}, \mathbf{y} \mid \beta) \cdot \Pr(\beta)$$

$$= \arg \max_{\beta} \prod_{i=1}^{n} e^{-(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / \sigma^2} \cdot \prod_{i=1}^{n} e^{-(\beta_i)^2 / \gamma^2}$$

$$= \arg \max_{\beta} \sum_{i=1}^{n} -(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / \sigma^2 + \sum_{i=1}^{d} -(\beta_i)^2 / \gamma^2$$

$$= \arg \min_{\beta} \sum_{i=1}^{n} (y_i - \langle \mathbf{x}_i, \beta \rangle)^2 + \frac{\sigma^2}{\gamma^2} \sum_{i=1}^{d} (\beta_i)^2 / \sigma^2.$$

Choose  $\boldsymbol{\beta}$  to minimize  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \frac{\sigma^2}{\gamma^2}\|\boldsymbol{\beta}\|_2^2$ .

Completely different justification for ridge regularization!

**Test your intuition:** What modeling assumption justifies LASSO regularization: min  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$ ?

LINEAR CLASSIFICATION

#### MOTIVATING PROBLEM

**Breast Cancer Biopsy:** Determine if a breast lump in a patient is <u>malignant</u> (cancerous) or <u>benign</u> (safe).

- Collect cells from lump using fine needle biopsy.
- Stain and examine cells under microscope.
- Based on certain characteristics (shape, size, cohesion) determine if likely malignant or not).



cross section



#### MOTIVATING PROBLEM

## Demo: demo\_breast\_cancer.ipynb

## Data: UCI machine learning repository

#### **Breast Cancer Wisconsin (Original) Data Set**

Download: Data Folder, Data Set Description

Abstract: Original Wisconsin Breast Cancer Database



Data Set Characteristics:	Multivariate	Number of Instances:	699	Area:	Life
Attribute Characteristics:	Integer	Number of Attributes:	10	Date Donated	1992-07-15
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	564320

#### 

**Features:** 10 numerical scores about cell characteristics (Clump Thickness, Uniformity, Marginal Adhesion, etc.)

**Data:**  $(x_1, y_1), \ldots, (x_n, y_n).$ 

 $\mathbf{x}_i = [1, 5, 4 \dots, 2]$  contains score values.

Label  $y_i \in \{0, 1\}$  is 0 if benign cells, 1 if malignant cells.

**Goal:** Based on scores (which would be collected manually, or even learned on their own using an ML algorithm) predict if a sample of cells is malignant or benign.

# Approach:

- Naive Bayes Classifier can be extended to **x** with numerical values (instead of binary values as seen before). Will see on homework.
- Today: Learn a different type of classifier.

We pick two variables, <u>Margin Adhesion</u> and <u>Size Uniformity</u> and plot a scatter plot. Points with label 1 (malignant) are plotted in blue, those with label 2 (benign) are plotted in green.



Lots of overlapping points! Hard to get a sense of the data. <sup>29</sup>

#### PLOTTING WITH JITTER

**Simple + Useful Trick:** data <u>jittering</u>. Add tiny random noise (using e.g. **np.random.randn**) to data to prevent overlap.



Noise is only for plotting. It is not added to the data for training, testing, etc.

#### PLOTTING WITH JITTER

**Simple + Useful Trick:** data <u>jittering</u>. Add tiny random noise (using e.g. **np.random.randn**) to data to prevent overlap.



Noise is only for plotting. It is not added to the data for training, testing, etc.

Any ideas for possible classification rules for this data?



Given vector of predictors  $\mathbf{x}_i \in \mathbb{R}^d$  (here d = 2) find a parameter vector  $\boldsymbol{\beta} \in \mathbb{R}^d$  and threshold  $\lambda$ .

- Predict  $y_i = 0$  if  $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle \leq \lambda$ .
- Predict  $y_i = 1$  if  $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle > \lambda$



Line has equation  $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle = \lambda$ .

### Next class: How do we find a good linear separator?