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

NYU Tandon School of Engineering, Prof. Christopher Musco

In a Bayesian or Probabilistic 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 $x_i$ is the size of the apartment in square feet. Each $y_i$ is the monthly rent in dollars.

3. 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.

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

$x \begin{bmatrix} 0, 1, 1, 0, 0, 1, \vdots \\ y \end{bmatrix}$
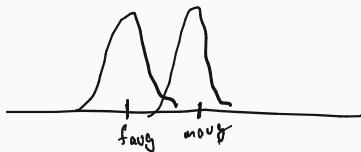$y \begin{bmatrix} 60'' & 62'' & 70'' & \cdots & \vdots \end{bmatrix}$

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

**Model:** $\sigma^2, f_{avg}, m_{avg}$



height

$(X, y)$

$\downarrow$

$0, 1$ with probability $\frac{1}{2}, \frac{1}{2}$

$$y = \begin{cases} \text{If } x = 0, & y = m_{avg} + N(0, \sigma^2) \\ \text{If } x = 1, & y = f_{avg} + N(0, \sigma^2) \end{cases}$$

using:
numpy.random.random

4

**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:**

$$\text{min} = 350 \text{ sq ft.}$$
$$\text{max} = 5000 \text{ sq ft}$$

$$y = c \cdot x + N(0, \sigma^2)$$

$$\boxed{X \sim \text{Unif} [350, 5000]}$$

$$(x, \underline{\quad})$$

$$\boxed{y = c \cdot x + \text{Unif} [-v, v]}$$

drow
from uniform
distribution

$$y = (c + \text{Unif} [-1, 1]) \cdot x$$

Using data: find $c = \$5$   $v = \$2$

5

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
$y_i \in \{0, 1\}$ with zero indicating the student has not taken
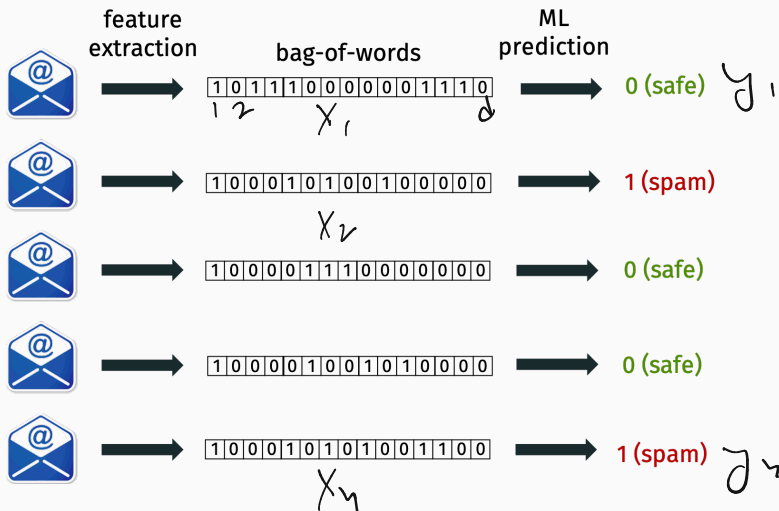machine learning, one indicating they have.

Model:

$$X \sim \text{unif} \underset{\substack{1/4 \quad 1/4 \quad 1/4 \quad 1/4}}{\left[ fresh, \; soph, \; Jun, \; Sen. \right]}$$

$$\begin{cases} \text{if } X = fresh, & y = 1 \text{ with prob. } 0 \\ \text{if } X = soph, & y = 1 \text{ with prob. } .05 \\ \text{if } \cdots \\ \text{if } X = sen, & y = 1 \text{ with prob. } .5 \end{cases}$$

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).

Both target labels and data vectors are binary.

**Probabilistic model** for (bag-of-words, label) pair $(\mathbf{x}, y)$:

- 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. $p(x_i = 1 \mid y = 0)$.
- If $y = 1$, for each $i$, set $x_i = 1$ with prob. $p(x_i = 1 \mid y = 1)$.

<span style="color:orange">Unknown model parameters:</span>

$$[0\ 1\ 0\ 0\ 1\ 1\ 1\ 1\ 1]$$

- $p(y = 0), p(y = 1)$,
- $p(x_1 = 1 \mid y = 0), \ldots, p(x_d = 1 \mid y = 0)$.
- $p(x_1 = 1 \mid y = 1), \ldots, p(x_d = 1 \mid y = 1)$.

How would you estimate these parameters?

9

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 $p(x_i = 1 \mid y = 0)$ to the empirical probability word $i$ appears in a <u>non-spam</u> email.
- For each word $i$, set $p(x_i = 1 \mid y = 1)$ 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 $(\mathbf{x}, \underline{\quad})$, choose the label $y \in \{0, 1\}$ which is <u>most likely</u> given the data. Recall $\mathbf{x} = [0, 0, 1, \ldots, 1, 0]$.

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

Step 1. Compute:

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

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

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

How to compute the posterior? **Bayes rule!**

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

$$p(y = 1 \mid \mathbf{x}) = p(\mathbf{x} \mid y = 1)\, p(y = 1) \,/\, p(\mathbf{x})$$

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

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

12

Goal is to determine which is larger:

$$p(y = 0 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = 0)\,p(y = 0)}{p(\mathbf{x})} \qquad \text{vs.}$$

$$p(y = 1 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = 1)\,p(y = 1)}{p(\mathbf{x})}$$

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 $p(\mathbf{x} \mid y = 0)$ by assuming independence:

*[handwritten: model parameters chosen using data]*

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

*[handwritten: $[0, 1, 1, 0, 0, \ldots]$ → $0, 1$]*

- $p(x_i \mid y = 0)$ is the probability you observe $x_i$ given that an email is not spam.[1]

A more complicated method might take dependencies into account.

---

[1] Recall, $x_i$ 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 $p(0 \text{ at position } i \mid y = 0), p(1 \text{ at position } i \mid y_0)$
    - Compute $p(0 \text{ at position } i \mid y = 1), p(1 \text{ at position } i \mid y = 1)$

**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

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

# OTHER APPLICATIONS OF
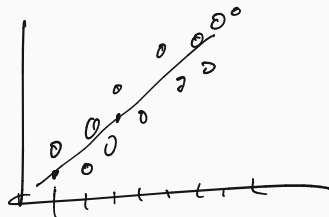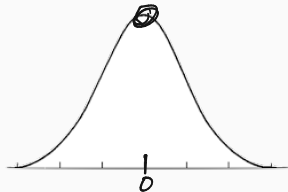## THE BAYESIAN PERSPECTIVE

The Bayesian view offers an interesting alternative perspective on <u>many</u> machine learning techniques.

**Example:** Linear Regression.

**Probabilistic model:**

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

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

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

$$P_\sigma(\eta) = e^{-\eta^2/\sigma^2}$$

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

16

$$y^* = \arg\max_y p(y \mid x_i)$$

**Example:** Linear Regression.

**Probabilistic model:**

$$p(y \mid x_i) = e^{-\frac{(y - \langle x_i, \beta \rangle)^2}{\sigma^2}} \quad \overbrace{\phantom{xxxx}}^{\eta}$$

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

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

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

$$y^* = \langle x_i, \beta \rangle \qquad e^{-\sigma^2/\sigma^2} \qquad 1$$

How should be select $\boldsymbol{\beta}$ for our model?

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

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

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

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

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

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

- $y_i = \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \eta$
- where $p(\eta = z) \sim e^{-z^2/\sigma^2}$

$$\Pr(\mathbf{X}, \mathbf{y} \mid \boldsymbol{\beta}) \sim$$
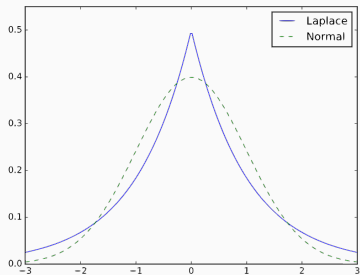
Easier to work with the log likelihood:

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

Choose $\boldsymbol{\beta}$ to minimize $\sum_{i=1}^{n}(y_i - \langle \mathbf{x}_i, \boldsymbol{\beta}\rangle)^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$!

This is a completely different justification for squared loss.

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.



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

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 $\beta$:

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

~~assume all $\beta$'s equally likely~~

**Bayesian view:** Assume values in $\boldsymbol{\beta} = [\beta_1, \ldots, \beta_d]$ 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 $\boldsymbol{\beta}$ to maximize:

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

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

Easier to work with the **log likelihood**:

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

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

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

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

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: $\min \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$?