Naive Baves Classifier

(Brief Intro: to be continued in the next lecture)

Naive Bayes classifier

- Probability of class label:
 - $p(C_k)$: Constant (e.g., Bernoulli)
- Conditional probability of data given the class
 - Naive Bayes assumption: $p(\mathbf{x} \mid C_k)$ is factorized (Each coordinate of ${\bf x}$ is conditionally independent of other coordinates given the class label)

$$P(x_1, ..., x_M | C_k) = P(x_1 | C_k) \cdots P(x_M | C_k) = \prod_{i=1}^M P(x_i | C_k)$$

Classification: use Bayes rule

(binary)
$$P(C_1|\mathbf{x}) = \frac{P(C_1,\mathbf{x})}{P(\mathbf{x})} = \frac{P(C_1,\mathbf{x})}{P(C_1,\mathbf{x}) + P(C_2,\mathbf{x})}$$

Naive Bayes classifier

• When classifying, we can simply find the class C_k that maximizes $P(C_k|\mathbf{x})$ using the Bayes rule:

$$\arg\max_k P(C_k|\mathbf{x}) = \arg\max_k P(C_k,\mathbf{x})$$

$$= \arg\max_k P(C_k)P(\mathbf{x}|C_k)$$
 Naive Bayes assumption
$$= \arg\max_k P(C_k)\prod_{j=1}^M P(x_j|C_k)$$

Example: Naive Bayes for real-valued inputs

- Probability of class label:
 - $-p(C_{\iota})$: Constant (e.g., Bernoulli)
- Conditional probability of data given the class
 - Naive Bayes assumption: $P(\mathbf{x}|C_{\iota})$ is factorized (e.g., 1D Gaussian)

$$P(x_1, ..., x_M | C_k) = P(x_1 | C_k) \cdots P(x_M | C_k)$$

$$= \prod_{j=1}^M P(x_j | C_k)$$

$$= \prod_{j=1}^M \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Note: this is equivalent to GDA with diagonal covariance!!

Comparison: Discriminative vs. Generative

- The generative approach is typically model-based, and it can generate synthetic data from $p(\mathbf{x} \mid C_{\iota})$.
 - By comparing the synthetic data and real data, we get a sense of how good the generative model is.
- The discriminative approach will typically have fewer parameters to estimate and have less assumptions about data distribution.
 - Linear (e.g. logistic regression) v/s quadratic (e.g., Gaussian discriminant analysis) in the dimension of the
 - Less generative assumptions about the data (however, constructing the features may require domain knowledge)

Naive Bayes classifier

- Probability of class label:
- P(G) = P, P(G) = 1-P
- $p(C_k)$: Constant (e.g., Bernoulli)
- Conditional probability of data given the class
 - Naive Bayes assumption: $p(\mathbf{x} \mid C_k)$ is factorized (Each coordinate of x is conditionally independent of other coordinates given the class label)

$$P(x_1,...,x_M|C_k) = P(x_1|C_k) \cdots P(x_M|C_k) = \prod_{k=1}^{M} P(x_j|C_k)$$

Classification: use Bayes rule

(binary)
$$P(C_1|\mathbf{x}) = \frac{P(C_1,\mathbf{x})}{P(\mathbf{x})} = \frac{P(C_1,\mathbf{x})}{P(C_1,\mathbf{x}) + P(C_2,\mathbf{x})}$$

Example: Spam mail classification

Label: y=1 (spam), y=0 (non-spam)

- Features:
 - Kj E (o, 1) M $-x_i$: j-th word in the mail, where M is the vocabulary size.
 - Multinomial variable (M-dimensional binary vector with only one coordinate with 1)
- Naive Bayes Assumption:
 - Given a class label y, each word in a mail is a independent multinomial variable.

Naive Bayes Spam classifier

Model

 $P(\operatorname{spam}) = Bernoulli(\phi)$ $P(\operatorname{word}|\operatorname{spam}) = Multinomial(\mu_1^s, \dots, \mu_M^s)$

 $P(\text{word}|\text{nonspam}) = Multinomial(\mu_1^{ns}, \dots, \mu_M^{ns})$ $P(\textbf{X=und}|\text{nonspam}) = Multinomial(\mu_1^{ns}, \dots, \mu_M^{ns})$ $e_{\textbf{X}} : \textbf{N} = \mathcal{C} : \textbf{N} =$

manuscript (3) neurips (14) reviewers (18)

P(word (nonspon) ~ Mubinon (0.5,0.4, 0.05,0.05) choice (9) congratulations (8) deals exclusive (7) gift (20) giveaway (6) limited (9) plan (5) Sale (6) Select (8) Special (13) top (5

top words from my non-spam emails

 $version_{\,(14)\,\,view\,\,(10)}\, visiting_{\,(19)}\, week_{\,(15)}$

top words from my spam emails

Naive Bayes Spam classifier

Model

$$\begin{split} P(\text{spam}) &= Bernoulli(\phi) \\ P(\text{word}|\text{spam}) &= Multinomial(\mu_1^s, \dots, \mu_M^s) \\ P(\text{word}|\text{nonspam}) &= Multinomial(\mu_1^{ns}, \dots, \mu_M^{ns}) \end{split}$$

Goal

Find ϕ , μ^s , μ^{ns} that best fits the data $\{(x^{(1)}, y^{(1)}), ..., (x^{(N)}, y^{(N)})\}$ by maximizing the joint likelihood:

$$\prod_{i=1}^{N} P(\mathbf{x}^{(i)}, y^{(i)})$$

- Joint Likelihood (joint probability of inputs/labels)
 - Note that the joint likelihood is conditioned on parameters ϕ , μ ^s, μ ^{ns}

Naive Bayes Spam classifier

Model

$$\begin{split} P(\text{spam}) &= Bernoulli(\phi) \\ P(\text{word}|\text{spam}) &= Multinomial(\mu_1^s, \dots, \mu_M^s) \\ P(\text{word}|\text{nonspam}) &= Multinomial(\mu_1^{ns}, \dots, \mu_M^{ns}) \end{split}$$

Goal

Find ϕ , μ^{s} , μ^{ns} that best fits the data $\{(x^{(1)}, y^{(1)}), ..., (x^{(N)}, y^{(N)})\}$

Likelihood - conditioned on parameters ϕ , μ^s , μ^{ns}

$$\begin{split} & \prod_{i=1} P(\mathbf{x}^{(i)}, y^{(i)}) \\ = & \prod_{i=1}^N P(\mathbf{x}^{(i)}|y^{(i)}) P(y^{(i)}) \\ = & \underbrace{\left(\prod_{\underline{i}: y^{(i)} = 1} P(\mathbf{x}^{(i)}|y^{(i)}) P(y^{(i)})\right)}_{\text{Spam}} \underbrace{\left(\prod_{\underline{i}: y^{(i)} = 0} P(\mathbf{x}^{(i)}|y^{(i)}) P(y^{(i)})\right)}_{\text{Non-spam}} \end{split}$$

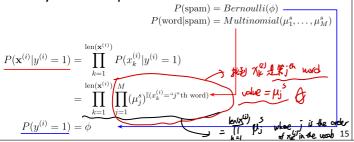
Naive Bayes Spam classifier

Likelihood - spam

$$\left(\prod_{i:y^{(i)}=1} \underline{P(\mathbf{x}^{(i)}|y^{(i)})}\underline{P(y^{(i)})}\right)$$

$$x_k^{(i)}$$
 k-th word

Naive Bayes assumption:



Naive Bayes Spam classifier

Likelihood - spam (cont')

$$\left(\prod_{i:y^{(i)}=1} P\left(\mathbf{x}^{(i)} \mid y^{(i)}\right) P\left(y^{(i)}\right) \right)$$

$$= \left(\prod_{i:y^{(i)}=1}^{N} \prod_{k=1}^{\operatorname{len}\left(\mathbf{x}^{(i)}\right)} \prod_{j=1}^{M} \left(\mu_{j}^{s}\right)^{\mathbb{I}\left(\mathbf{x}_{k}^{(i)} = \mathbf{u}_{j}^{u} \text{th word }\right)} \phi \right) \text{ conft}$$

$$= \left(\prod_{i:y^{(i)}=1}^{N} \prod_{k=1}^{\operatorname{len}\left(\mathbf{x}^{i,j}\right)} \prod_{j=1}^{M} \left(\mu_{j}^{u,j}\right)^{\mathcal{I}\left(--\right)} \right) \left(\prod_{i:y^{(i)}=1}^{N} \mathcal{Y} \right)$$

Naive Bayes Spam classifier

Likelihood - spam (cont')

$$\begin{split} & \left(\prod_{i:y^{(i)}=1} P\left(\mathbf{x}^{(i)} \mid y^{(i)}\right) P\left(y^{(i)}\right) \right) \\ & = \left(\prod_{i:y^{(i)}=1}^{N} \prod_{k=1}^{\operatorname{len}\left(x^{(i)}\right)} \prod_{j=1}^{M} \left(\mu_{j}^{s}\right)^{\mathbb{I}\left(x_{k}^{(i)} = ^{u}j^{n}\operatorname{th} \operatorname{word}\right)} \phi \right) \\ & = \left(\prod_{i:y^{(i)}=1}^{N} \prod_{k=1}^{\operatorname{len}\left(x^{(i)}\right)} \prod_{j=1}^{M} \left(\mu_{j}^{s}\right)^{\mathbb{I}\left(x_{k}^{(i)} = ^{u}j^{n}\operatorname{th} \operatorname{word}\right)} \right) \left(\prod_{i:y^{(i)}=1}^{N} \phi \right) \\ & = \left(\prod_{j=1}^{M} \left(\mu_{j}^{s}\right)^{\sum_{i:y^{(i)}=1}^{N} \sum_{k=1}^{\operatorname{len}\left(x^{(i)}\right)} \mathbb{I}\left(x_{k}^{(i)} = ^{u}j^{n}\operatorname{th} \operatorname{word}\right) \right) \left(\prod_{i:y^{(i)}=1}^{N} \phi \right) \\ & = \left(\prod_{j=1}^{M} \left(\mu_{j}^{s}\right)^{N_{j}^{\operatorname{spam}}} \right) \phi^{N^{\operatorname{spam}}} \right) \phi^{N^{\operatorname{spam}}} \end{split}$$

Naive Bayes Spam classifier

Likelihood - non-spam

$$\left(\prod_{i:y^{(i)}=1}P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)})\right)\left(\prod_{i:y^{(i)}=0}\underline{P(\mathbf{x}^{(i)}|y^{(i)})}\underline{P(y^{(i)})}\right)$$

Similarly for non-spam mails,

$$P(\mathbf{x}^{(i)}|y^{(i)} = 0) = \prod_{k=1}^{\operatorname{len}(\mathbf{x}^{(i)})} P(\mathbf{x}^{(i)}|y^{(i)} = 0)$$

$$= \prod_{k=1}^{\operatorname{len}(\mathbf{x}^{(i)})} \prod_{j=1}^{M} (\mu_j^{ns})^{I(x_k^{(i)} = "j" \operatorname{th word})}$$

$$P(y^{(i)} = 0) = 1 - \phi$$

Maximum likelihood estimation

· Putting together:

$$\begin{array}{l} \text{Futting together.} \\ &\prod_{i=1}^{N}P(\mathbf{x}^{(i)},y^{(i)}) \\ = &\left(\prod_{i:y^{(i)}=1}P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)})\right)\left(\prod_{i:y^{(i)}=0}P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)})\right) \\ = &\left(\phi^{N^{spam}}\prod_{word\,j}\left(\mu_{j}^{s}\right)^{N_{j}^{spam}}\right)\left((1-\phi)^{N^{nonspam}}\prod_{word\,j}\left(\mu_{j}^{ns}\right)^{N_{j}^{nonspam}}\right) \\ & \\ &\frac{\text{Recall:}}{N^{spam}}\text{: total \# examples for spam} \\ &N^{nonspam}\text{: total \# examples for non-spam} \\ &N^{spam}_{j}\text{: total \# word j from the entire spam emails} \\ &N^{nonspam}_{j}\text{: total \# word j from the entire nonspam emails} \\ &N^{nonspam}_{j}\text{: total \# word j from the entire nonspam emails} \\ \end{array}$$

Maximum likelihood estimation

· Putting together:

$$\begin{split} &\prod_{i=1}^{N} P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)})\right) \left(\prod_{i:y^{(i)}=0} P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)})\right) \\ &= \left(\phi^{N^{spam}} \prod_{word j} \left(\mu_{j}^{s}\right)^{N_{j}^{spam}}\right) \left((1-\phi)^{N^{nonspam}} \prod_{word j} \left(\mu_{j}^{ns}\right)^{N_{j}^{nonspam}}\right) \end{split}$$

· Joint Log-likelihood

$$= \log \prod_{i=1}^{N} P(x^{(i)}, y^{(i)})$$

 $N^{spam}\log\phi+\sum N^{spam}_j\log\mu^s_j+N^{nonspam}\log(1-\phi)+\sum N^{nonspam}_j\log\mu^n_j$

Maximum likelihood estimation

· Joint Log-likelihood

$$\log P(\mathcal{D})$$

$$= \log \prod_{i=1}^{N} P(x^{(i)}, y^{(i)})$$

$$= N^{spam} \log \phi + \sum_{word j} N_{j}^{spam} \log \mu_{j}^{s} + N^{nonspam} \log (1 - \phi) + \sum_{word j} N_{j}^{nonspam} \log \mu_{j}^{ns}$$

- Maximum-likelihood
 - Take the derivative of log-likelihood w.r.t. the parameters, and set it to zero.

Maximum likelihood estimation

• From
$$\frac{\partial l}{\partial \phi} = \frac{1}{\phi} N^{spam} - \frac{1}{1-\phi} N^{nonspam} = 0$$
 We get
$$\phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$$
• Removing dependent variables:

$$\begin{split} \sum_{word\,j=1}^{M} N_{j}^{spam} \log \mu_{j}^{s} &= \sum_{word\,j=1}^{M-1} N_{j}^{spam} \log \mu_{j}^{s} + N_{M}^{spam} \log(1 - \sum_{j=1}^{M-1} \mu_{j}^{s}) \\ &\frac{\partial}{\partial \mu_{j}^{s}} \left(\sum_{word\,j=1}^{M} N_{j}^{spam} \log \mu_{j}^{s} \right) = \frac{N_{j}^{spam}}{\mu_{j}^{s}} - \frac{N_{M}^{spam}}{1 - \sum_{j=1}^{M-1} \mu_{j}^{s}} = 0 \\ &\text{s.t.} \quad \sum_{j} \mu_{j}^{s} = 1 \\ &\sum_{i} \mu_{j}^{ns} = 1 \end{split}$$

28

Maximum likelihood estimation

• From
$$\frac{\partial l}{\partial \phi} = \frac{1}{\phi} N^{spam} - \frac{1}{1-\phi} N^{nonspam} = 0$$
 We get $\phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$

Removing dependent variables:

$$\sum_{word\,j=1}^{M} N_{j}^{spam} \log \mu_{j}^{s} = \sum_{word\,j=1}^{M-1} N_{j}^{spam} \log \mu_{j}^{s} + N_{M}^{spam} \log(1 - \sum_{j=1}^{M-1} \mu_{j}^{s})$$

$$\frac{\partial}{\partial \mu_{j}^{s}} \left(\sum_{word\,j=1}^{M} N_{j}^{spam} \log \mu_{j}^{s} \right) = \frac{N_{j}^{spam}}{\mu_{j}^{s}} - \frac{N_{M}^{spam}}{1 - \sum_{j=1}^{M-1} \mu_{j}^{s}} = 0$$

$$\frac{N_{j}^{spam}}{\mu_{j}^{s}} = constant, \forall j$$

$$\mu_{j}^{s} = \frac{N_{j}^{spam}}{\sum_{j} N_{j}^{spam}}$$

Maximum likelihood estimation

• Summary:

$$\begin{split} P(spam) = & \phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}} \\ P(word = j | spam) = & \mu_j^s = \frac{N_j^{spam}}{\sum_j N_j^{spam}} \\ P(word = j | non - spam) = & \mu_j^{ns} = \frac{N_j^{nonspam}}{\sum_j N_j^{nonspam}} \end{split}$$

 $N^{nonspam}$: total # examples for non-spam

 N_i^{spam} : total # word j from the entire spam emails

Laplace Smoothing

- Maximum likelihood is problematic when a specific word count is 0
 - Leads to probability of 0!
- Solution: Put "imaginary" counts for each word
 - prevent zero probability estimates (overfitting)!
 - E.g.: Adding "1" as imaginary count for each word

$$P(spam) = \phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$$

$$P(word = j | spam) = \mu_j^s = \frac{N_j^{spam} + 1}{\sum_j N_j^{spam} + M}$$

$$P(word = j | non - spam) = \mu_j^{ns} = \frac{N_j^{nonspam} + 1}{\sum_j N_j^{nonspam} + M}$$