# CSC2515 Lecture 8: Intro to Generative Models

#### David Duvenaud

#### Based on Materials from Roger Grosse, University of Toronto

- March 15: Detailed project instructions released. After HW4 all remaining office hours are for project help.
- April 1: Project proposals due. Hand in earlier for earlier feedback.
- April 9: Project presentations (About 5 min). Don't have to be finished, just prelim results.
- April 15: Projects due
- Groups of up to 3. Example / default projects will be provided.

- Reproduce the experimental results from some existing papers. Perform sensitivity analysis on hyper-parameters.
- Apply / extend existing algorithms to a new application / task / dataset.

- Improve / fix an existing algorithm. Evaluate the improvement on benchmark environments.
- Develop novel model architectures / algorithms to a new application / area / environment.
- You are welcome to do a project related to your research. In this case, your project proposal and final report must each clearly explain the relationship to your research, what work was already done prior to the course, and what work (if any) was done by people not on the project team. Our expectations will be higher in this case.

- Bayesian parameter estimation: average predictions over all hypotheses, proportional to their posterior probability.
- Generative classification: learn to model the distributions of inputs belonging to each class
  - Naïve Bayes (discrete inputs)
  - Gaussian Discriminant Analysis (continuous inputs)

- Maximum likelihood has a pitfall: if you have too little data, it can overfit.
- E.g., what if you flip the coin twice and get H both times?

- Maximum likelihood has a pitfall: if you have too little data, it can overfit.
- E.g., what if you flip the coin twice and get H both times?

$$\theta_{\rm ML} = \frac{N_H}{N_H + N_T} = \frac{2}{2+0} = 1$$

- Because it never observed T, it assigns this outcome probability 0. This problem is known as data sparsity.
- If you observe a single T in the test set, the log-likelihood is  $-\infty$ .

- In maximum likelihood, the observations are treated as random variables, but the parameters are not.
- The Bayesian approach treats the parameters as random variables as well.

- In maximum likelihood, the observations are treated as random variables, but the parameters are not.
- The Bayesian approach treats the parameters as random variables as well.
- To define a Bayesian model, we need to specify two distributions:
  - The prior distribution  $p(\theta)$ , which encodes our beliefs about the parameters *before* we observe the data
  - The likelihood  $p(\mathcal{D} | \theta)$ , same as in maximum likelihood

- In maximum likelihood, the observations are treated as random variables, but the parameters are not.
- The Bayesian approach treats the parameters as random variables as well.
- To define a Bayesian model, we need to specify two distributions:
  - The prior distribution  $p(\theta)$ , which encodes our beliefs about the parameters *before* we observe the data
  - The likelihood  $p(\mathcal{D} | \theta)$ , same as in maximum likelihood
- When we update our beliefs based on the observations, we compute the posterior distribution using Bayes' Rule:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\boldsymbol{\theta})p(\mathcal{D} \mid \boldsymbol{\theta})}{\int p(\boldsymbol{\theta}')p(\mathcal{D} \mid \boldsymbol{\theta}') \,\mathrm{d}\boldsymbol{\theta}'}.$$

• We rarely ever compute the denominator explicitly.

• Let's revisit the coin example. We already know the likelihood:

$$L(\theta) = \rho(\mathcal{D}) = \theta^{N_H} (1-\theta)^{N_T}$$

• It remains to specify the prior  $p(\theta)$ .

• Let's revisit the coin example. We already know the likelihood:

$$L(\theta) = p(\mathcal{D}) = \theta^{N_H} (1-\theta)^{N_T}$$

- It remains to specify the prior  $p(\theta)$ .
  - We can choose an uninformative prior, which assumes as little as possible. A reasonable choice is the uniform prior.
  - But our experience tells us 0.5 is more likely than 0.99. One particularly useful prior that lets us specify this is the beta distribution:

$$p(\theta; a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}.$$

• This notation for proportionality lets us ignore the normalization constant:

$$p(\theta; a, b) \propto \theta^{a-1}(1-\theta)^{b-1}.$$

• Beta distribution for various values of a, b:



Some observations:

- The expectation  $\mathbb{E}[\theta] = a/(a+b)$ .
- The distribution gets more peaked when *a* and *b* are large.
- The uniform distribution is the special case where a = b = 1.
- The main thing the beta distribution is used for is as a prior for the Bernoulli distribution.

• Computing the posterior distribution:

$$egin{aligned} & eta( heta) \propto eta( heta) p(\mathcal{D} \,|\, heta) \ & \propto \left[ heta^{a-1} (1- heta)^{b-1} 
ight] \left[ heta^{N_H} (1- heta)^{N_T} 
ight] \ & = heta^{a-1+N_H} (1- heta)^{b-1+N_T}. \end{aligned}$$

• This is just a beta distribution with parameters  $N_H + a$  and  $N_T + b$ .

• Computing the posterior distribution:

$$egin{aligned} & p(m{ heta} \mid \mathcal{D}) \propto p(m{ heta}) p(\mathcal{D} \mid m{ heta}) \ & \propto \left[ heta^{a-1} (1- heta)^{b-1} 
ight] \left[ heta^{N_H} (1- heta)^{N_T} 
ight] \ & = heta^{a-1+N_H} (1- heta)^{b-1+N_T}. \end{aligned}$$

This is just a beta distribution with parameters N<sub>H</sub> + a and N<sub>T</sub> + b.
The posterior expectation of θ is:

$$\mathbb{E}[\theta \mid \mathcal{D}] = \frac{N_H + a}{N_H + N_T + a + b}$$

• Computing the posterior distribution:

$$egin{aligned} & m{p}(m{ heta} \mid \mathcal{D}) \propto m{p}(m{ heta}) m{p}(\mathcal{D} \mid m{ heta}) \ & \propto \left[ heta^{m{a}-1}(1- heta)^{b-1} 
ight] \left[ heta^{m{N}_H}(1- heta)^{m{N}_T} 
ight] \ & = heta^{m{a}-1+m{N}_H}(1- heta)^{b-1+m{N}_T}. \end{aligned}$$

This is just a beta distribution with parameters N<sub>H</sub> + a and N<sub>T</sub> + b.
The posterior expectation of θ is:

$$\mathbb{E}[\theta \mid \mathcal{D}] = \frac{N_H + a}{N_H + N_T + a + b}$$

- The parameters *a* and *b* of the prior can be thought of as pseudo-counts.
  - The reason this works is that the prior and likelihood have the same functional form. This phenomenon is known as conjugacy, and it's very useful.

Intro ML (UofT)

#### CSC2515 Lec8

Bayesian inference for the coin flip example:

Small data setting  $N_H = 2, N_T = 0$ 



Bayesian inference for the coin flip example:



When you have enough observations, the data overwhelm the prior.

- What do we actually do with the posterior?
- The posterior predictive distribution is the distribution over future observables given the past observations. We compute this by marginalizing out the parameter(s):

$$p(\mathcal{D}' | \mathcal{D}) = \int p(\boldsymbol{\theta} | \mathcal{D}) p(\mathcal{D}' | \boldsymbol{\theta}) d\boldsymbol{\theta}.$$
(1)

- What do we actually do with the posterior?
- The posterior predictive distribution is the distribution over future observables given the past observations. We compute this by marginalizing out the parameter(s):

$$p(\mathcal{D}' | \mathcal{D}) = \int p(\boldsymbol{\theta} | \mathcal{D}) \, p(\mathcal{D}' | \boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta}.$$
(1)

For the coin flip example:

$$\begin{split} \theta_{\text{pred}} &= \Pr(x' = H \,|\, \mathcal{D}) \\ &= \int p(\theta \,|\, \mathcal{D}) \Pr(x' = H \,|\, \theta) \,\mathrm{d}\theta \\ &= \int \text{Beta}(\theta; N_H + a, N_T + b) \cdot \theta \,\mathrm{d}\theta \\ &= \mathbb{E}_{\text{Beta}(\theta; N_H + a, N_T + b)}[\theta] \\ &= \frac{N_H + a}{N_H + N_T + a + b}, \end{split}$$

(2)

Bayesian estimation of the mean temperature in Toronto

- Assume observations are i.i.d. Gaussian with known standard deviation σ and unknown mean μ
- Broad Gaussian prior over μ, centered at 0
- We can compute the posterior and posterior predictive distributions analytically (full derivation in notes)
- Why is the posterior predictive distribution more spread out than the posterior distribution?



Comparison of maximum likelihood and Bayesian parameter estimation

- Some advantages of the Bayesian approach
  - More robust to data sparsity
  - Incorporate prior knowledge
  - Smooth the predictions by averaging over plausible explanations

Comparison of maximum likelihood and Bayesian parameter estimation

- Some advantages of the Bayesian approach
  - More robust to data sparsity
  - Incorporate prior knowledge
  - Smooth the predictions by averaging over plausible explanations
- Problem: maximum likelihood is an optimization problem, while Bayesian parameter estimation is an integration problem
  - This means maximum likelihood is much easier in practice, since we can just do gradient descent
  - Automatic differentiation packages make it really easy to compute gradients
  - There aren't any comparable black-box tools for Bayesian parameter estimation (although Stan can do quite a lot)

- Maximum a-posteriori (MAP) estimation: find the most likely parameter settings under the posterior
- This converts the Bayesian parameter estimation problem into a maximization problem

$$\begin{split} \hat{\boldsymbol{\theta}}_{\mathrm{MAP}} &= \arg\max_{\boldsymbol{\theta}} \ p(\boldsymbol{\theta} \,|\, \mathcal{D}) \\ &= \arg\max_{\boldsymbol{\theta}} \ p(\boldsymbol{\theta}) \, p(\mathcal{D} \,|\, \boldsymbol{\theta}) \\ &= \arg\max_{\boldsymbol{\theta}} \ \log p(\boldsymbol{\theta}) + \log p(\mathcal{D} \,|\, \boldsymbol{\theta}) \end{split}$$

## Maximum A-Posteriori Estimation

• Joint probability in the coin flip example:

$$\begin{split} \log p(\theta, \mathcal{D}) &= \log p(\theta) + \log p(\mathcal{D} \mid \theta) \\ &= \operatorname{const} + (a-1) \log \theta + (b-1) \log (1-\theta) + N_H \log \theta + N_T \log (1-\theta) \\ &= \operatorname{const} + (N_H + a - 1) \log \theta + (N_T + b - 1) \log (1-\theta) \end{split}$$

#### Maximum A-Posteriori Estimation

• Joint probability in the coin flip example:

$$\begin{split} \log p(\theta, \mathcal{D}) &= \log p(\theta) + \log p(\mathcal{D} \mid \theta) \\ &= \operatorname{const} + (a-1)\log\theta + (b-1)\log(1-\theta) + N_H\log\theta + N_T\log(1-\theta) \\ &= \operatorname{const} + (N_H + a - 1)\log\theta + (N_T + b - 1)\log(1-\theta) \end{split}$$

#### • Maximize by finding a critical point

$$0 = \frac{\mathrm{d}}{\mathrm{d}\theta} \log p(\theta, \mathcal{D}) = \frac{N_H + a - 1}{\theta} - \frac{N_T + b - 1}{1 - \theta}$$

#### Maximum A-Posteriori Estimation

• Joint probability in the coin flip example:

$$\begin{split} \log p(\theta, \mathcal{D}) &= \log p(\theta) + \log p(\mathcal{D} \mid \theta) \\ &= \operatorname{const} + (a-1)\log \theta + (b-1)\log(1-\theta) + N_H \log \theta + N_T \log(1-\theta) \\ &= \operatorname{const} + (N_H + a - 1)\log \theta + (N_T + b - 1)\log(1-\theta) \end{split}$$

• Maximize by finding a critical point

$$0 = \frac{\mathrm{d}}{\mathrm{d}\theta} \log p(\theta, \mathcal{D}) = \frac{N_H + a - 1}{\theta} - \frac{N_T + b - 1}{1 - \theta}$$

• Solving for  $\theta$ ,

$$\hat{\theta}_{\mathrm{MAP}} = rac{N_H + a - 1}{N_H + N_T + a + b - 2}$$

Comparison of estimates in the coin flip example:

	Formula	$N_H = 2, N_T = 0$	$N_H = 55, N_T = 45$
$\hat{ heta}_{\mathrm{ML}}$	$\frac{N_H}{N_H+N_T}$	1	$\frac{55}{100} = 0.55$
$\theta_{\rm pred}$	$rac{N_H+a}{N_H+N_T+a+b}$	$rac{4}{6}pprox 0.67$	$rac{57}{104}pprox 0.548$
$\hat{ heta}_{\mathrm{MAP}}$	$rac{N_H+a-1}{N_H+N_T+a+b-2}$	$\frac{3}{4} = 0.75$	$rac{56}{102}pprox 0.549$

 $\hat{ heta}_{\mathrm{MAP}}$  assigns nonzero probabilities as long as a,b>1.

Comparison of predictions in the Toronto temperatures example



#### 7 observations



#### Generative Classifiers and Naïve Bayes

Two approaches to classification:



Two approaches to classification:

- Discriminative: directly learn to predict *t* as a function of **x**.
  - Sometimes this means modeling  $p(t | \mathbf{x})$  (e.g. logistic regression).
  - Sometimes this means learning a decision rule without a probabilistic interpretation (e.g. KNN, SVM).
- Generative: model the data distribution for each class separately, and make predictions using posterior inference.
  - Fit models of p(t) and  $p(\mathbf{x} | t)$ .
  - Infer the posterior  $p(t | \mathbf{x})$  using Bayes' Rule.

• Bayes classifier: given features **x**, we compute the posterior class probabilities using Bayes' Rule:



• Requires fitting  $p(\mathbf{x} | t)$  and p(t)

• Bayes classifier: given features **x**, we compute the posterior class probabilities using Bayes' Rule:



- Requires fitting  $p(\mathbf{x} | t)$  and p(t)
- How can we compute  $p(\mathbf{x})$  for binary classification?

• Bayes classifier: given features **x**, we compute the posterior class probabilities using Bayes' Rule:



- Requires fitting  $p(\mathbf{x} | t)$  and p(t)
- How can we compute  $p(\mathbf{x})$  for binary classification?

$$p(\mathbf{x}) = p(\mathbf{x} \mid t = 0) \operatorname{Pr}(t = 0) + p(\mathbf{x} \mid t = 1) \operatorname{Pr}(t = 1)$$

• Note: sometimes it's more convenient to just compute the numerator and normalize.

# Naïve Bayes

- **Example:** want to classify emails into spam (t = 1) or non-spam (t = 0) based on the words they contain.
  - Use bag-of-words features, i.e. a binary vector **x** where entry  $x_j = 1$  if word *j* appeared in the email. (Assume a dictionary of *D* words.)
# Naïve Bayes

- **Example:** want to classify emails into spam (t = 1) or non-spam (t = 0) based on the words they contain.
  - Use bag-of-words features, i.e. a binary vector **x** where entry  $x_j = 1$  if word *j* appeared in the email. (Assume a dictionary of *D* words.)
- Estimating the prior p(t) is easy (e.g. maximum likelihood).
- **Problem:**  $p(\mathbf{x} | t)$  is a joint distribution over *D* binary random variables, which requires  $2^{D}$  entries to specify directly!

# Naïve Bayes

- **Example:** want to classify emails into spam (t = 1) or non-spam (t = 0) based on the words they contain.
  - Use bag-of-words features, i.e. a binary vector **x** where entry  $x_j = 1$  if word *j* appeared in the email. (Assume a dictionary of *D* words.)
- Estimating the prior p(t) is easy (e.g. maximum likelihood).
- **Problem:**  $p(\mathbf{x} | t)$  is a joint distribution over *D* binary random variables, which requires  $2^{D}$  entries to specify directly!
- We'd like to impose structure on the distribution such that:
  - it can be compactly represented
  - learning and inference are both tractable
- Probabilistic graphical models are a powerful and wide-ranging class of techniques for doing this. We'll just scratch the surface here, but you'll learn about them in detail in CSC2506.

- Naïve Bayes makes the assumption that the word features x<sub>j</sub> are conditionally independent given the class t.
  - This means  $x_i$  and  $x_j$  are independent under the conditional distribution  $p(\mathbf{x} \mid t)$ .
  - Note: this doesn't mean they're independent. (E.g., "Viagra" and "cheap" are correlated insofar as they both depend on *t*.)
  - Mathematically, this means the distribution factorizes:

$$p(t, x_1, \ldots, x_D) = p(t) p(x_1 \mid t) \cdots p(x_D \mid t).$$

- Naïve Bayes makes the assumption that the word features x<sub>j</sub> are conditionally independent given the class t.
  - This means  $x_i$  and  $x_j$  are independent under the conditional distribution  $p(\mathbf{x} \mid t)$ .
  - Note: this doesn't mean they're independent. (E.g., "Viagra" and "cheap" are correlated insofar as they both depend on *t*.)
  - Mathematically, this means the distribution factorizes:

$$p(t, x_1, \ldots, x_D) = p(t) p(x_1 \mid t) \cdots p(x_D \mid t).$$

- Compact representation of the joint distribution
  - Prior probability of class:  $\Pr(t=1) = \phi$
  - Conditional probability of word feature given class:  $Pr(x_j = 1 | t) = \theta_{jt}$
  - 2D + 1 parameters total

# Bayes Nets (Optional)

• We can represent this model using an directed graphical model, or Bayesian network:



- This graph structure means the joint distribution factorizes as a product of conditional distributions for each variable given its parent(s).
- Intuitively, you can think of the edges as reflecting a causal structure. But mathematically, we can't infer causality without additional assumptions.
- You'll learn a lot about graphical models in CSC2506.

• The parameters can be learned efficiently because the log-likelihood decomposes into independent terms for each feature.

$$\ell(\theta) = \sum_{i=1}^{N} \log p(t^{(i)}, \mathbf{x}^{(i)})$$
  
=  $\sum_{i=1}^{N} \log p(t^{(i)}) \prod_{j=1}^{D} p(x_j^{(i)} | t^{(i)})$   
=  $\sum_{i=1}^{N} \left[ \log p(t^{(i)}) + \sum_{j=1}^{D} \log p(x_j^{(i)} | t^{(i)}) \right]$   
=  $\sum_{i=1}^{N} \log p(t^{(i)})$   
Bernoulli log-likelihood  
of labels

• Each of these log-likelihood terms depends on different sets of parameters, so they can be optimized independently.

Intro ML (UofT)

- Want to maximize  $\sum_{i=1}^{N} \log p(x_j^{(i)} | t^{(i)})$
- This is a minor variant of our coin flip example. Let  $\theta_{ab} = \Pr(x_j = a \mid t = b)$ . Note  $\theta_{1b} = 1 \theta_{0b}$ .

- Want to maximize  $\sum_{i=1}^{N} \log p(x_j^{(i)} | t^{(i)})$
- This is a minor variant of our coin flip example. Let  $\theta_{ab} = \Pr(x_j = a \mid t = b)$ . Note  $\theta_{1b} = 1 \theta_{0b}$ .
- Log-likelihood:

$$\begin{split} \sum_{i=1}^{N} \log p(x_{j}^{(i)} \mid t^{(i)}) &= \sum_{i=1}^{N} t^{(i)} x_{j}^{(i)} \log \theta_{11} + \sum_{i=1}^{N} t^{(i)} (1 - x_{j}^{(i)}) \log (1 - \theta_{11}) \\ &+ \sum_{i=1}^{N} (1 - t^{(i)}) x_{j}^{(i)} \log \theta_{10} + \sum_{i=1}^{N} (1 - t^{(i)}) (1 - x_{j}^{(i)}) \log (1 - \theta_{10}) \end{split}$$

- Want to maximize  $\sum_{i=1}^{N} \log p(x_j^{(i)} | t^{(i)})$
- This is a minor variant of our coin flip example. Let  $\theta_{ab} = \Pr(x_j = a \mid t = b)$ . Note  $\theta_{1b} = 1 \theta_{0b}$ .
- Log-likelihood:

$$\begin{split} \sum_{i=1}^N \log p(x_j^{(i)} \mid t^{(i)}) &= \sum_{i=1}^N t^{(i)} x_j^{(i)} \log \theta_{11} + \sum_{i=1}^N t^{(i)} (1-x_j^{(i)}) \log (1-\theta_{11}) \\ &+ \sum_{i=1}^N (1-t^{(i)}) x_j^{(i)} \log \theta_{10} + \sum_{i=1}^N (1-t^{(i)}) (1-x_j^{(i)}) \log (1-\theta_{10}) \end{split}$$

• Obtain maximum likelihood estimates by setting derivatives to zero:

$$heta_{11} = rac{N_{11}}{N_{11} + N_{01}} \qquad heta_{10} = rac{N_{10}}{N_{10} + N_{00}}$$

where  $N_{ab}$  is the counts for  $x_j = a$  and t = b.

## Naïve Bayes: Inference

- We predict the category by performing inference in the model.
- Apply Bayes' Rule:

$$p(t \mid \mathbf{x}) = \frac{p(t) p(\mathbf{x} \mid t)}{\sum_{t'} p(t') p(\mathbf{x} \mid t')} \\ = \frac{p(t) \prod_{j=1}^{D} p(x_j \mid t)}{\sum_{t'} p(t') \prod_{j=1}^{D} p(x_j \mid t')}$$

- We need not compute the denominator if we're simply trying to determine the mostly likely *t*.
- Shorthand notation:

$$p(t \mid \mathbf{x}) \propto p(t) \prod_{j=1}^{D} p(x_j \mid t)$$

• Once we compute  $p(t | \mathbf{x})$ , what do we do with it?

- Once we compute  $p(t | \mathbf{x})$ , what do we do with it?
- Sometimes we want to make a single prediction or decision y. This is a decision theory problem, just like when we analyzed the bias/variance/Bayes-error decomposition.
  - Define a loss function  $\mathcal{L}(y, t)$  and choose  $y_* = \arg \min_y \mathbb{E}[\mathcal{L}(y, t) | \mathbf{x}]$ .

- Once we compute  $p(t | \mathbf{x})$ , what do we do with it?
- Sometimes we want to make a single prediction or decision y. This is a decision theory problem, just like when we analyzed the bias/variance/Bayes-error decomposition.
  - Define a loss function  $\mathcal{L}(y, t)$  and choose  $y_* = \arg \min_y \mathbb{E}[\mathcal{L}(y, t) | \mathbf{x}]$ .
- Examples
  - Squared error loss: choose  $y_{\star} = \mathbb{E}[t \mid \mathbf{x}]$
  - 0-1 loss: choose the most likely category
  - Cross-entropy loss: return the probability  $y = Pr(t = 1 | \mathbf{x})$

- Once we compute  $p(t | \mathbf{x})$ , what do we do with it?
- Sometimes we want to make a single prediction or decision y. This is a decision theory problem, just like when we analyzed the bias/variance/Bayes-error decomposition.
  - Define a loss function  $\mathcal{L}(y, t)$  and choose  $y_{\star} = \arg \min_{y} \mathbb{E}[\mathcal{L}(y, t) | \mathbf{x}]$ .
- Examples
  - Squared error loss: choose  $y_{\star} = \mathbb{E}[t \mid \mathbf{x}]$
  - 0-1 loss: choose the most likely category
  - Cross-entropy loss: return the probability  $y = Pr(t = 1 | \mathbf{x})$
  - Asymmetric loss (e.g. false positives are much worse than false negatives for spam filtering): apply a threshold other than 0.5.

- Once we compute  $p(t | \mathbf{x})$ , what do we do with it?
- Sometimes we want to make a single prediction or decision y. This is a decision theory problem, just like when we analyzed the bias/variance/Bayes-error decomposition.
  - Define a loss function  $\mathcal{L}(y, t)$  and choose  $y_* = \arg \min_y \mathbb{E}[\mathcal{L}(y, t) | \mathbf{x}]$ .
- Examples
  - Squared error loss: choose  $y_{\star} = \mathbb{E}[t \mid \mathbf{x}]$
  - 0-1 loss: choose the most likely category
  - Cross-entropy loss: return the probability  $y = \Pr(t = 1 | \mathbf{x})$
  - Asymmetric loss (e.g. false positives are much worse than false negatives for spam filtering): apply a threshold other than 0.5.
    - Warning: this is theoretically tidy, but doesn't really work unless you're careful to obtain calibrated posterior probabilities.
    - "Calibrated" means all the times you predict (say)  $Pr(t = k | \mathbf{x}) = 0.9$  should be correct 90% on average.
    - Naïve Bayes is generally not calibrated due to the "naïve" conditional independence assumption.

# Naïve Bayes

- Naïve Bayes is an amazingly cheap learning algorithm!
- Training time: estimate parameters using maximum likelihood
  - Compute co-occurrence counts of each feature with the labels.
  - Requires only one pass through the data!
- Test time: apply Bayes' Rule
  - Cheap because of the model structure. (For more general models, Bayesian inference can be very expensive and/or complicated.)

# Naïve Bayes

- Naïve Bayes is an amazingly cheap learning algorithm!
- Training time: estimate parameters using maximum likelihood
  - Compute co-occurrence counts of each feature with the labels.
  - Requires only one pass through the data!
- Test time: apply Bayes' Rule
  - Cheap because of the model structure. (For more general models, Bayesian inference can be very expensive and/or complicated.)
- We covered the Bernoulli case for simplicity. But our analysis easily extends to other probability distributions.
- Unfortunately, it's usually less accurate in practice compared to discriminative models.
  - The problem is the "naïve" independence assumption.
  - We're covering it primarily as a stepping stone towards latent variable models.

#### Gaussian Discriminant Analysis

- Generative models model p(t) and  $p(\mathbf{x} \mid t)$
- Recall that  $p(\mathbf{x} | t = k)$  may be very complex

$$p(x_1, \cdots, x_D \mid t) = p(x_1 \mid x_2, \cdots, x_D, t) \cdots p(x_{D-1} \mid x_D, t) p(x_D \mid t)$$

- Naïve Bayes used a conditional independence assumption to make everything tractable.
- For continuous inputs, we can instead make it tractable by using a simple distribution: multivariate Gaussians.

• Observation per patient: White blood cell count & glucose value.



• How can we model  $p(\mathbf{x} | t = k)$ ? Multivariate Gaussian

### Multivariate Parameters

Mean

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_D \end{pmatrix}$$

Covariance

$$\boldsymbol{\Sigma} = \operatorname{Cov}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})^{\top}(\mathbf{x} - \boldsymbol{\mu})] = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1D} \\ \sigma_{12} & \sigma_2^2 & \cdots & \sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{D1} & \sigma_{D2} & \cdots & \sigma_D^2 \end{pmatrix}$$

• These statistics uniquely define a multivariate Gaussian distribution. (This is not true for distributions in general!)

### Multivariate Gaussian Distribution

•  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , a multivariate Gaussian (or multivariate normal) distribution is defined as

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |\mathbf{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$$



• Mahalanobis distance  $(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$  measures the distance from  $\mathbf{x}$  to  $\boldsymbol{\mu}$  in a space stretched according to  $\boldsymbol{\Sigma}$ .

Intro ML (UofT)





Figure: Contour plot of the pdf



```
Intro ML (UofT)
```





 $Cov(x_1, x_2)=0, Var(x_1)=Var(x_2)$ 

 $Cov(x_1, x_2)=0, Var(x_1)>Var(x_2)$ 







Intro ML (UofT)

- Gaussian Discriminant Analysis in its general form assumes that  $p(\mathbf{x}|t)$  is distributed according to a multivariate Gaussian distribution
- Multivariate Gaussian distribution:

$$p(\mathbf{x} \mid t = k) = \frac{1}{(2\pi)^{D/2} |\mathbf{\Sigma}_k|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right]$$

where  $|\mathbf{\Sigma}_k|$  denotes the determinant of the matrix.

- Each class k has associated mean vector  $\mu_k$  and covariance matrix  $\mathbf{\Sigma}_k$
- How many parameters?

- Gaussian Discriminant Analysis in its general form assumes that  $p(\mathbf{x}|t)$  is distributed according to a multivariate Gaussian distribution
- Multivariate Gaussian distribution:

$$p(\mathbf{x} \mid t = k) = \frac{1}{(2\pi)^{D/2} |\mathbf{\Sigma}_k|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right]$$

where  $|\mathbf{\Sigma}_k|$  denotes the determinant of the matrix.

- Each class k has associated mean vector  $\mu_k$  and covariance matrix  $\mathbf{\Sigma}_k$
- How many parameters?
  - Each  $\mu_k$  has D parameters, for DK total.
  - Each Σ<sub>k</sub> has O(D<sup>2</sup>) parameters, for O(D<sup>2</sup>K) could be hard to estimate (more on that later).

## GDA: Learning

- Learn the parameters for each class using maximum likelihood
- For simplicity, assume binary classification

$$p(t \mid \phi) = \phi^t (1 - \phi)^{1-t}$$

You can compute the ML estimates in closed form (φ and μ<sub>k</sub> are easy, Σ<sub>k</sub> is tricky)

$$\phi = \frac{1}{N} \sum_{i=1}^{N} r_1^{(i)}$$

$$\mu_k = \frac{\sum_{i=1}^{N} r_k^{(i)} \cdot \mathbf{x}^{(i)}}{\sum_{i=1}^{N} r_k^{(i)}}$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_{i=1}^{N} r_k^{(i)}} \sum_{i=1}^{N} r_k^{(i)} (\mathbf{x}^{(i)} - \mu_k) (\mathbf{x}^{(i)} - \mu_k)^\top$$

$$r_k^{(i)} = \mathbb{1}[t^{(i)} = k]$$

### GDA Decision Boundary

• Recall: for Bayes classifiers, we compute the decision boundary with Bayes' Rule:

$$p(t \mid \mathbf{x}) = \frac{p(t) \, p(\mathbf{x} \mid t)}{\sum_{t'} p(t') \, p(\mathbf{x} \mid t')}$$

• Plug in the Gaussian  $p(\mathbf{x} | t)$ :

$$\log p(t_k | \mathbf{x}) = \log p(\mathbf{x} | t_k) + \log p(t_k) - \log p(\mathbf{x})$$
  
=  $-\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \mathbf{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log p(t_k) - \log p(\mathbf{x})$ 

Decision boundary:

$$(\mathbf{x} - \boldsymbol{\mu}_k)^{\top} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = (\mathbf{x} - \boldsymbol{\mu}_\ell)^{\top} \boldsymbol{\Sigma}_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) + \text{Const}$$

• What's the shape of the boundary?

### GDA Decision Boundary

• Recall: for Bayes classifiers, we compute the decision boundary with Bayes' Rule:

$$p(t \mid \mathbf{x}) = \frac{p(t) \, p(\mathbf{x} \mid t)}{\sum_{t'} p(t') \, p(\mathbf{x} \mid t')}$$

• Plug in the Gaussian  $p(\mathbf{x} | t)$ :

$$\log p(t_k | \mathbf{x}) = \log p(\mathbf{x} | t_k) + \log p(t_k) - \log p(\mathbf{x})$$
  
=  $-\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \mathbf{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log p(t_k) - \log p(\mathbf{x})$ 

• Decision boundary:

$$(\mathbf{x} - \boldsymbol{\mu}_k)^{\top} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = (\mathbf{x} - \boldsymbol{\mu}_\ell)^{\top} \boldsymbol{\Sigma}_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) + \text{Const}$$

- What's the shape of the boundary?
  - We have a quadratic function in **x**, so the decision boundary is a conic section!

Intro ML (UofT)

### **GDA** Decision Boundary



• Our equation for the decision boundary:

$$(\mathbf{x} - \boldsymbol{\mu}_k)^{\top} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = (\mathbf{x} - \boldsymbol{\mu}_\ell)^{\top} \boldsymbol{\Sigma}_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) + \text{Const}$$

• Expand the product and factor out constants (w.r.t. **x**):

$$\mathbf{x}^{\top} \mathbf{\Sigma}_{k}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_{k}^{\top} \mathbf{\Sigma}_{k}^{-1} \mathbf{x} = \mathbf{x}^{\top} \mathbf{\Sigma}_{\ell}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_{\ell}^{\top} \mathbf{\Sigma}_{\ell}^{-1} \mathbf{x} + \text{Const}$$

• What if all classes share the same covariance Σ?

• Our equation for the decision boundary:

$$(\mathbf{x} - \boldsymbol{\mu}_k)^{\top} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = (\mathbf{x} - \boldsymbol{\mu}_\ell)^{\top} \boldsymbol{\Sigma}_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) + \text{Const}$$

• Expand the product and factor out constants (w.r.t. **x**):

$$\mathbf{x}^{\top} \mathbf{\Sigma}_{k}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_{k}^{\top} \mathbf{\Sigma}_{k}^{-1} \mathbf{x} = \mathbf{x}^{\top} \mathbf{\Sigma}_{\ell}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_{\ell}^{\top} \mathbf{\Sigma}_{\ell}^{-1} \mathbf{x} + \text{Const}$$

What if all classes share the same covariance Σ?

• We get a linear decision boundary!

$$-2\boldsymbol{\mu}_{k}^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{x} = -2\boldsymbol{\mu}_{\ell}^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{x} + \text{Const}$$
$$(\boldsymbol{\mu}_{k} - \boldsymbol{\mu}_{\ell})^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{x} = \text{Const}$$
### GDA Decision Boundary: Shared Covariances



• Binary classification: If you examine  $p(t = 1 | \mathbf{x})$  under GDA and assume  $\Sigma_0 = \Sigma_1 = \Sigma$ , you will find that it looks like this:

$$p(t \mid \mathbf{x}, \phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}$$

where  $(\mathbf{w}, b)$  are chosen based on  $(\phi, \mu_0, \mu_1, \mathbf{\Sigma})$ .

• Same model as logistic regression!

- GDA makes a stronger modeling assumption: assumes class-conditional data is multivariate Gaussian
  - If this is true, GDA is asymptotically efficient (best model in limit of large N)
  - If it's not true, the quality of the predictions might suffer.

- GDA makes a stronger modeling assumption: assumes class-conditional data is multivariate Gaussian
  - If this is true, GDA is asymptotically efficient (best model in limit of large N)
  - If it's not true, the quality of the predictions might suffer.
- Many class-conditional distributions lead to logistic classifier.
  - When these distributions are non-Gaussian (i.e., almost always), LR usually beats GDA

- GDA makes a stronger modeling assumption: assumes class-conditional data is multivariate Gaussian
  - If this is true, GDA is asymptotically efficient (best model in limit of large N)
  - If it's not true, the quality of the predictions might suffer.
- Many class-conditional distributions lead to logistic classifier.
  - When these distributions are non-Gaussian (i.e., almost always), LR usually beats GDA
- GDA can handle easily missing features (how do you do that with LR?)

# Gaussian Naive Bayes

- What if **x** is high-dimensional?
  - The  $\Sigma_k$  have  $\mathcal{O}(D^2 \mathcal{K})$  parameters, which can be a problem if D is large.
  - We already saw we can save some a factor of *K* by using a shared covariance for the classes.
  - Any other idea you can think of?

# Gaussian Naive Bayes

- What if **x** is high-dimensional?
  - The  $\Sigma_k$  have  $\mathcal{O}(D^2 K)$  parameters, which can be a problem if D is large.
  - We already saw we can save some a factor of *K* by using a shared covariance for the classes.
  - Any other idea you can think of?
- Naive Bayes: Assumes features independent given the class

$$p(\mathbf{x} \mid t = k) = \prod_{j=1}^{D} p(x_j \mid t = k)$$

• Assuming likelihoods are Gaussian, how many parameters required for Naive Bayes classifier?

- What if **x** is high-dimensional?
  - The  $\Sigma_k$  have  $\mathcal{O}(D^2 K)$  parameters, which can be a problem if D is large.
  - We already saw we can save some a factor of *K* by using a shared covariance for the classes.
  - Any other idea you can think of?
- Naive Bayes: Assumes features independent given the class

$$p(\mathbf{x} | t = k) = \prod_{j=1}^{D} p(x_j | t = k)$$

- Assuming likelihoods are Gaussian, how many parameters required for Naive Bayes classifier?
  - This is equivalent to assuming the  $x_j$  are uncorrelated, i.e.  $\Sigma$  is diagonal.
  - Hence, only D parameters for  $\Sigma$ !

### Gaussian Naïve Bayes

• Gaussian Naïve Bayes classifier assumes that the likelihoods are Gaussian:

$$p(x_j \mid t = k) = rac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left[rac{-(x_j - \mu_{jk})^2}{2\sigma_{jk}^2}
ight]$$

(this is just a 1-dim Gaussian, one for each input dimension)

- Model the same as GDA with diagonal covariance matrix
- Maximum likelihood estimate of parameters

$$\mu_{jk} = \frac{\sum_{i=1}^{N} r_k^{(i)} x_j^{(i)}}{\sum_{i=1}^{N} r_k^{(i)}}$$
$$\sigma_{jk}^2 = \frac{\sum_{i=1}^{N} r_k^{(i)} (x_j^{(i)} - \mu_{jk})^2}{\sum_{i=1}^{N} r_k^{(i)}}$$
$$r_k^{(i)} = \mathbb{1}[t^{(i)} = k]$$

### Decision Boundary: Isotropic

- We can go even further and assume the covariances are spherical, or isotropic.
- In this case:  $\Sigma = \sigma^2 \mathbf{I}$  (just need one parameter!)
- Going back to the class posterior for GDA:

$$\begin{split} \log p(t_k | \mathbf{x}) &= \log p(\mathbf{x} | t_k) + \log p(t_k) - \log p(\mathbf{x}) \\ &= -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{\Sigma}_k^{-1}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \mathbf{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \\ &+ \log p(t_k) - \log p(\mathbf{x}) \end{split}$$

• Suppose for simplicity that p(t) is uniform. Plugging in  $\Sigma = \sigma^2 \mathbf{I}$  and simplifying a bit,

$$\begin{split} \log p(t_k \mid \mathbf{x}) - \log p(t_\ell \mid \mathbf{x}) &= -\frac{1}{2\sigma^2} \left[ (\mathbf{x} - \boldsymbol{\mu}_k)^\top (\mathbf{x} - \boldsymbol{\mu}_k) - (\mathbf{x} - \boldsymbol{\mu}_\ell)^\top (\mathbf{x} - \boldsymbol{\mu}_\ell) \right] \\ &= -\frac{1}{2\sigma^2} \left[ \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 - \|\mathbf{x} - \boldsymbol{\mu}_\ell\|^2 \right] \end{split}$$

### Decision Boundary: Isotropic



• The decision boundary bisects the class means!







Logistic regression (acc 0.722)

