

# Machine Learning 4771

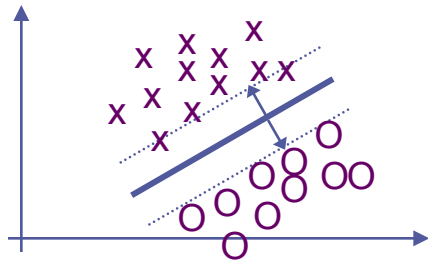
Instructor: Tony Jebara

# Topic 7

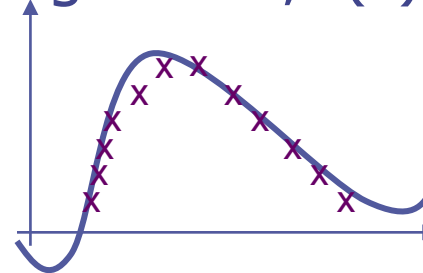
- Unsupervised Learning
- Statistical Perspective
- Probability Models
- Discrete & Continuous: Gaussian, Bernoulli, Multinomial
- Maximum Likelihood  $\rightarrow$  Logistic Regression
- Conditioning, Marginalizing, Bayes Rule, Expectations
- Classification, Regression, Detection
- Dependence/Independence
- Maximum Likelihood  $\rightarrow$  Naïve Bayes

# Unsupervised Learning

## Classification

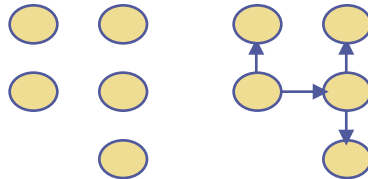
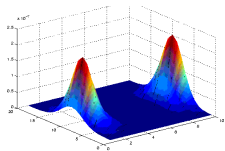


## Regression, $f(x)=y$

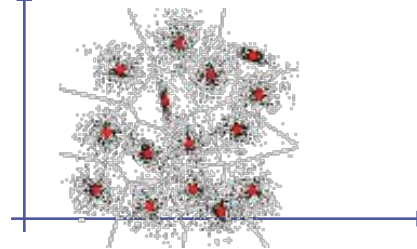


Supervised

## Density/Structure Estimation

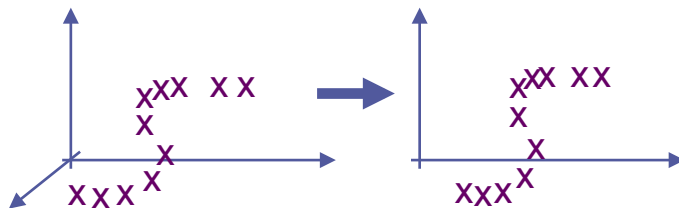


## Clustering

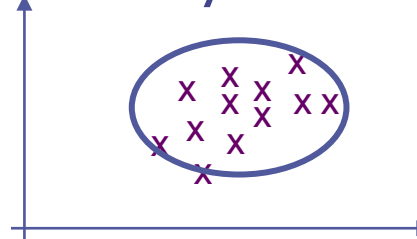


Unsupervised  
(can help supervised)

## Feature Selection

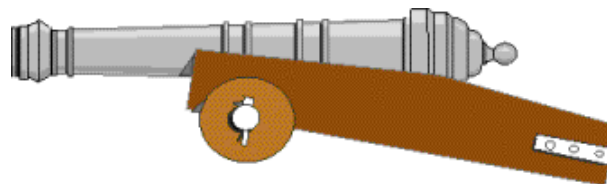


## Anomaly Detection



# Statistical Perspective

- Several problems with framework so far:
  - Only have input-output approaches (SVM, Neural Net)
  - Pulled non-linear squashing functions out of a hat
  - Pulled loss functions (squared error, etc.) out of a hat
- Better approach for classification?
- What if we have multi-class classification?
- What if other problems, i.e. unobserved values of  $x, y$ , etc...
- Also, what if we don't have a true function?
- Example of Projectile Cannon (c.f. Distal Learning)

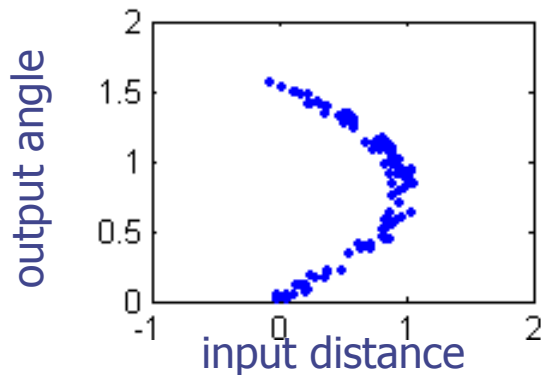


- Would like to train a regression function to control a cannon's angle of fire ( $y$ ) given target distance ( $x$ )

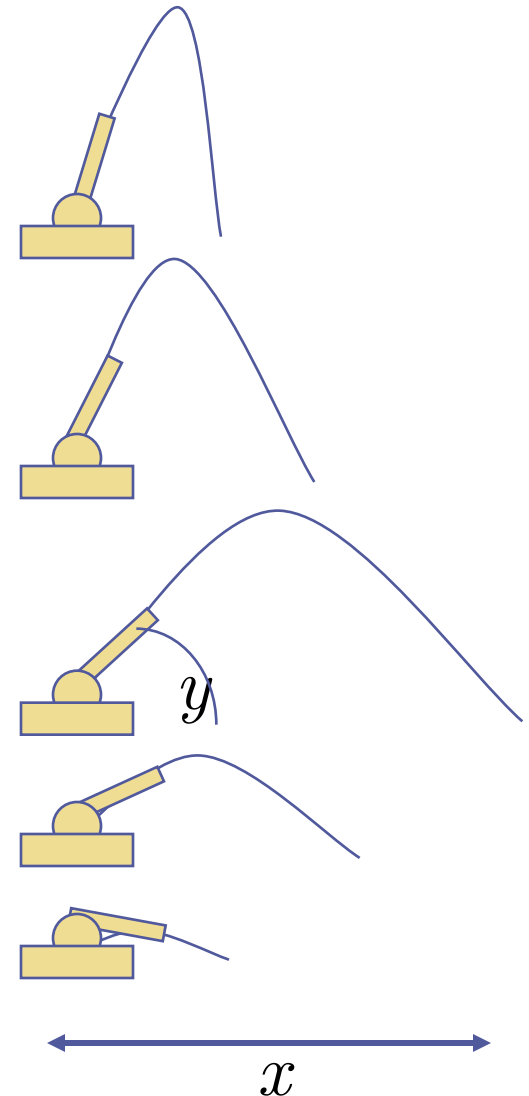
# Statistical Perspective

- Example of Projectile Cannon (45 degree problem)  
x = input target distance  
y = output cannon angle

$$x = \frac{v(0)^2}{g} \sin(2y) + noise$$

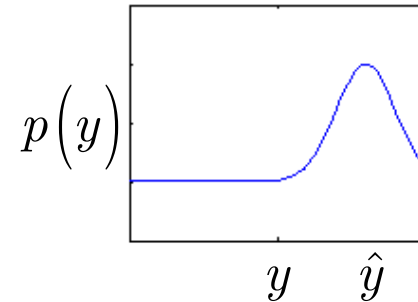


- What does least squares do?
- Conditional statistical models address this problem...



# Probability Models

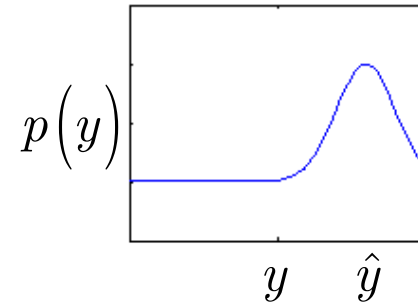
- Instead of deterministic functions, output is a probability
- Previously: our output was a scalar  $\hat{y} = f(x) = \theta^T x + b$
- Now: our output is a probability  $p(y)$   
e.g. a probability bump:



- $p(y)$  subsumes or is a superset of  $\hat{y}$
- Why is this representation for our answer more general?

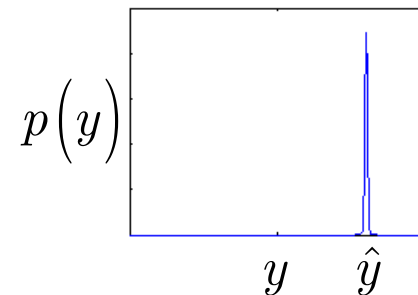
# Probability Models

- Instead of deterministic functions, output is a probability
- Previously: our output was a scalar  $\hat{y} = f(x) = \theta^T x + b$
- Now: our output is a probability  $p(y)$   
e.g. a probability bump:



- $p(y)$  subsumes or is a superset of  $\hat{y}$
- Why is this representation for our answer more general?  
→ A deterministic answer  $\hat{y}$  with complete confidence is like putting a probability  $p(y)$  where all the mass is at  $\hat{y}$  !

$$\hat{y} \Leftrightarrow p(y) = \delta(y - \hat{y})$$



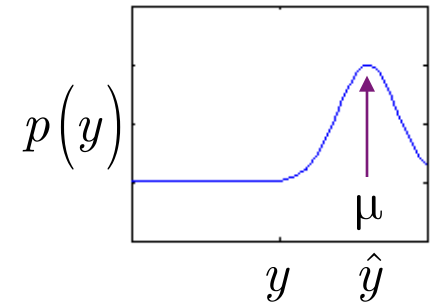
# Probability Models

- Now: our output is a probability density function (pdf)  $p(y)$
- Probability Model: a family of pdf's with adjustable parameters which lets us select one of many

$$p(y) \rightarrow p(y | \Theta)$$

- E.g.: 1-dim Gaussian distribution  
'given' 'mean' parameter  $\mu$ :

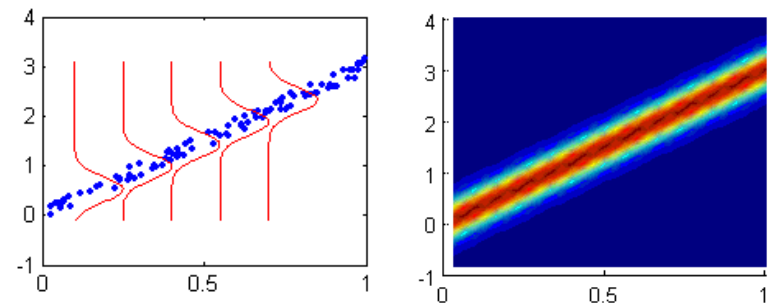
$$p(y | \mu) = N(y | \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\mu)^2}$$



- Want mean centered on  $f(x)$ 's value  $p(y) = N(y | f(x))$

- Now, linear regression is:

$$\begin{aligned} N(y | f(x)) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-f(x))^2} \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\theta^T x - b)^2} \end{aligned}$$





# Probability Models

- To fit to data, we typically “maximize likelihood” of the probability model
- Log-likelihood = objective function (i.e. negative of cost) for probability models which we want to maximize
- Define (conditional) likelihood as  $L(\Theta) = \prod_{i=1}^N p(y_i | x_i)$   
or log-Likelihood as  $l(\Theta) = \log(L(\Theta)) = \sum_{i=1}^N \log p(y_i | x_i)$
- For Gaussian  $p(y|x)$ , maximum likelihood is least squares!

$$\begin{aligned} \sum_{i=1}^N \log p(y_i | x_i) &= \sum_{i=1}^N \log N(y_i | f(x_i)) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_i - f(x_i))^2} \\ &= -N \log(\sqrt{2\pi}) - \sum_{i=1}^N \frac{1}{2} (y_i - f(x_i))^2 \end{aligned}$$

# Probability Models

- Can extend probability model to 2 bumps:

$$p(y | \Theta) = \frac{1}{2} N(y | \mu_1) + \frac{1}{2} N(y | \mu_2)$$

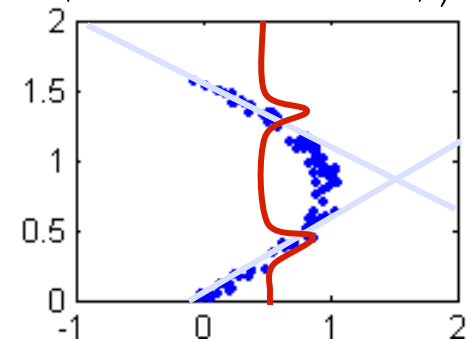
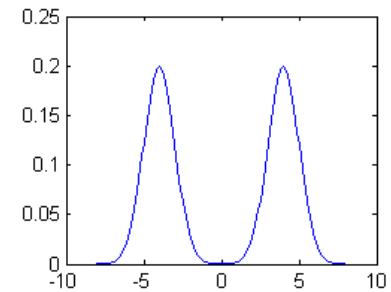
- Each mean can be a linear regression fn.

$$\begin{aligned} p(y | x, \Theta) &= \frac{1}{2} N(y | f_1(x)) + \frac{1}{2} N(y | f_2(x)) \\ &= \frac{1}{2} N(y | \theta_1^T x + b_1) + \frac{1}{2} N(y | \theta_2^T x + b_2) \end{aligned}$$

- Therefore the (conditional) log-likelihood to maximize is:

$$l(\Theta) = \sum_{i=1}^N \log \left( \frac{1}{2} N(y_i | \theta_1^T x_i + b_1) + \frac{1}{2} N(y_i | \theta_2^T x_i + b_2) \right)$$

- Maximize  $l(\theta)$  using gradient ascent
- Nicely handles the “cannon firing” data



# Probability Models

- Now classification: can also go beyond deterministic!
- Previously: wanted output to be binary  $\hat{y} = \{0,1\}$
- Now: our output is a probability  $p(y)$

e.g. a probability table:

y=0	y=1
0.73	0.27

$\alpha$  points to the value 0.27

- This subsumes or is a superset again...
- Consider probability over binary events (coin flips!):

e.g. Bernoulli distribution (i.e 1x2 probability table)  
with parameter  $\alpha$

$$p(y | \alpha) = \alpha^y (1 - \alpha)^{1-y} \quad \alpha \in [0,1]$$

- Linear classification can be done by setting  $\alpha$  equal to  $f(x)$ :

$$p(y | x) = f(x)^y (1 - f(x))^{1-y} \quad f(x) \in [0,1]$$

# Probability Models

- Now linear classification is:

$$p(y | x) = f(x)^y (1 - f(x))^{1-y} \quad f(x) \equiv \alpha \in [0, 1]$$

- Log-likelihood is (negative of cost function):

$$\begin{aligned} \sum_{i=1}^N \log p(y_i | x_i) &= \sum_{i=1}^N \log f(x_i)^{y_i} (1 - f(x_i))^{1-y_i} \\ &= \sum_{i=1}^N y_i \log f(x_i) + (1 - y_i) \log (1 - f(x_i)) \\ &= \sum_{i \in \text{class1}} \log f(x_i) + \sum_{i \in \text{class0}} \log (1 - f(x_i)) \end{aligned}$$

- But, need a squashing function since  $f(x)$  in  $[0, 1]$

- Use sigmoid or logistic again...

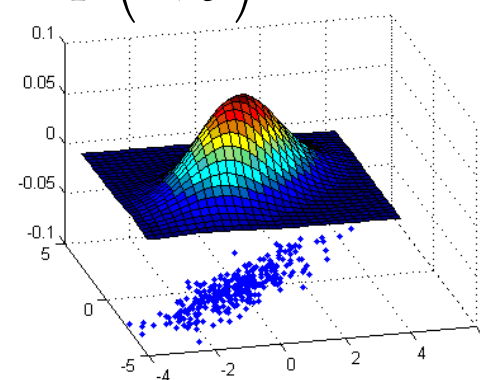
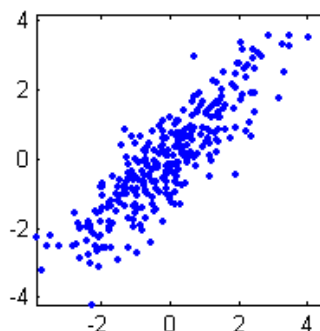
$$f(x) = \text{sigmoid}(\theta^T x + b) \in [0, 1]$$

- Called logistic regression  $\rightarrow$  *new loss function*
- Do gradient descent, similar to logistic output neural net!
- Can also handle multi-layer  $f(x)$  and do backprop again!

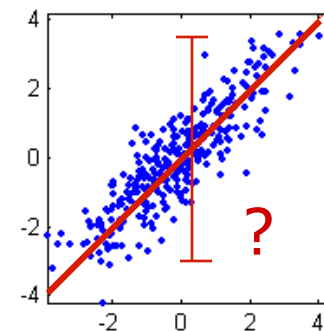
# Generative Probability Models

- Idea: Extend probability to describe *both* X and Y
- Find probability density function over both:  $p(x, y)$

E.g. *describe* data  
with Multi-Dim.  
Gaussian (later...)



- Called a 'Generative Model' because we can use it to synthesize or re-generate data similar to the training data we learned from
- Regression models & classification boundaries are not as flexible  
don't keep info about X  
don't model noise/uncertainty



# Properties of PDFs

- Let's review some basics of probability theory

- First, pdf is a function, multiple inputs, one output:

$$p(x_1, \dots, x_n) \qquad p(x_1 = 0.3, \dots, x_n = 1) = 0.2$$

- Function's output is always non-negative:

$$p(x_1, \dots, x_n) \geq 0$$

- Can have discrete or continuous or both inputs:

$$p(x_1 = 1, x_2 = 0, x_3 = 0, x_4 = 3.1415)$$

- Summing over the domain of all inputs gives unity:

$$\int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} p(x, y) dx dy = 1 \qquad \sum_y \sum_x p(x, y) = 1$$

0.4	0.1
0.3	0.2

**Continuous → integral, Discrete → sum**

# Properties of PDFs

- **Marginalizing:** integrate/sum out a variable leaves a marginal distribution over the remaining ones...

$$\sum_y p(x, y) = p(x)$$

- **Conditioning:** if a variable 'y' is 'given' we get a conditional distribution over the remaining ones...

$$p(x | y) = \frac{p(x, y)}{p(y)}$$

- **Bayes Rule:** mathematically just redo conditioning but has a deeper meaning (1764)... if we have  $\mathcal{X}$  being data and  $\theta$  being a model

$$\text{posterior} \rightarrow p(\theta | \mathcal{X}) = \frac{\overset{\text{likelihood}}{p(\mathcal{X} | \theta)} \overset{\text{prior}}{p(\theta)}}{\underset{\text{evidence}}{p(\mathcal{X})}}$$



# Properties of PDFs

- **Expectation:** can use pdf  $p(x)$  to compute averages and expected values for quantities, denoted by:

$$E_{p(x)} \{f(x)\} = \int_x p(x) f(x) dx \quad \text{or} \quad = \sum_x p(x) f(x)$$

- **Properties:**  $E \{cf(x)\} = cE \{f(x)\}$

$$E \{f(x) + c\} = E \{f(x)\} + c$$

$$E \{E \{f(x)\}\} = E \{f(x)\}$$

- **Mean:** expected value for  $x$

$$E_{p(x)} \{x\} = \int_{-\infty}^{\infty} p(x) x dx$$

- **Variance:** expected value of  $(x - \text{mean})^2$ , how much  $x$  varies

$$\begin{aligned} \text{Var} \{x\} &= E \left\{ \left( x - E \{x\} \right)^2 \right\} = E \left\{ x^2 - 2xE \{x\} + E \{x\}^2 \right\} \\ &= E \{x^2\} - 2E \{x\} E \{x\} + E \{x\}^2 = E \{x^2\} - E \{x\}^2 \end{aligned}$$

**example: speeding ticket**

Fine=0\$	Fine=20\$
0.8	0.2

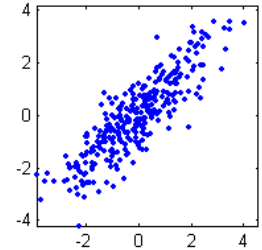
**expected cost of speeding?**

**$f(x=0)=0$ ,  $f(x=1)=20$**

**$p(x=0)=0.8$ ,  $p(x=1)=0.2$**



# Properties of PDFs



- Covariance: how strongly x and y vary together

$$\text{Cov}\{x, y\} = E\left\{\left(x - E\{x\}\right)\left(y - E\{y\}\right)\right\} = E\{xy\} - E\{x\}E\{y\}$$

- Conditional Expectation:  $E\{y | x\} = \int_y p(y | x) y dy$

$$E\left\{E\{y | x\}\right\} = \int_x p(x) \int_y p(y | x) y dy dx = E\{y\}$$

- Sample Expectation: If we don't have pdf  $p(x,y)$  can approximate expectations using samples of data

$$E_{p(x)}\{f(x)\} \simeq \frac{1}{N} \sum_{i=1}^N f(x_i)$$

- Sample Mean:  $E\{x\} \simeq \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

- Sample Var:  $E\left\{\left(x - E(x)\right)^2\right\} \simeq \frac{1}{N} \sum_{i=1}^N \left(x_i - \bar{x}\right)^2$

- Sample Cov:  $E\left\{\left(x - E(x)\right)\left(y - E(y)\right)\right\} \simeq \frac{1}{N} \sum_{i=1}^N \left(x_i - \bar{x}\right)\left(y_i - \bar{y}\right)$

# More Properties of PDFs

- **Independence:** probabilities of independent variables multiply. Denote with the following notation:

$$x \perp\!\!\!\perp y \rightarrow p(x, y) = p(x)p(y)$$

$$x \perp\!\!\!\perp y \rightarrow p(x | y) = p(x)$$

also note in this case:

$$\begin{aligned} E_{p(x,y)} \{xy\} &= \int_x \int_y p(x)p(y)xy \, dx \, dy \\ &= \int_x p(x)x \, dx \int_y p(y)y \, dy = E_{p(x)} \{x\} E_{p(y)} \{y\} \end{aligned}$$

- **Conditional independence:** when two variables become independent only if another is observed

$$x \perp\!\!\!\perp y | z \rightarrow p(x | y, z) = p(x | z)$$

$$x \perp\!\!\!\perp y | z \rightarrow p(x | y) \neq p(x)$$

# The IID Assumption

- Most of the time, we will assume that a dataset independent and identically distributed (IID)
- In many real situations, data is generated by some black box phenomenon in an arbitrary order.
- Assume we are given a dataset:

$$\mathcal{X} = \{x_1, \dots, x_N\}$$

“Independent” means that (given the model  $\theta$ ) the probability of our data multiplies:

$$p(x_1, \dots, x_N \mid \Theta) = \prod_{i=1}^N p_i(x_i \mid \Theta)$$

“Identically distributed” means that each marginal probability is the same for each data point

$$p(x_1, \dots, x_N \mid \Theta) = \prod_{i=1}^N p_i(x_i \mid \Theta) = \prod_{i=1}^N p(x_i \mid \Theta)$$

# The IID Assumption

- Bayes rule says likelihood is probability of data given model

$$\text{posterior} \rightarrow p(\theta | \mathcal{X}) = \frac{\overset{\text{likelihood}}{p(\mathcal{X} | \theta)} \overset{\text{prior}}{p(\theta)}}{\underset{\text{evidence}}{p(\mathcal{X})}}$$

- The likelihood of  $\mathcal{X} = \{x_1, \dots, x_N\}$  under IID assumptions is:

$$p(\mathcal{X} | \Theta) = p(x_1, \dots, x_N | \Theta) = \prod_{i=1}^N p_i(x_i | \Theta) = \prod_{i=1}^N p(x_i | \Theta)$$

- Learn joint distribution  $p(x | \Theta)$  by **maximum likelihood**:

$$\Theta^* = \arg \max_{\Theta} \prod_{i=1}^N p(x_i | \Theta) = \arg \max_{\Theta} \sum_{i=1}^N \log p(x_i | \Theta)$$

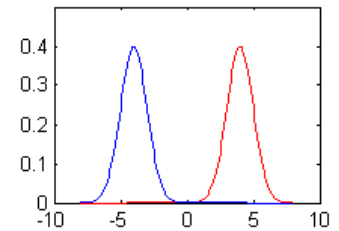
- Learn conditional  $p(y | x, \Theta)$  by **max conditional likelihood**:

$$\Theta^* = \arg \max_{\Theta} \prod_{i=1}^N p(y_i | x_i, \Theta) = \arg \max_{\Theta} \sum_{i=1}^N \log p(y_i | x_i, \Theta)$$

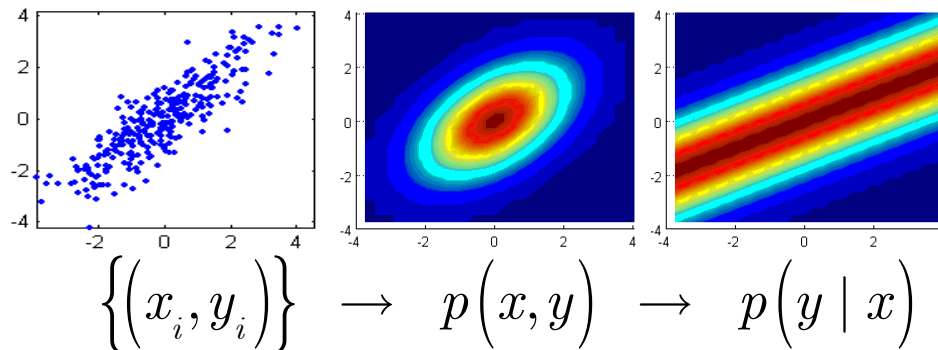
# Uses of PDFs

- **Classification:** have  $p(x,y)$  and given  $x$ . Asked for discrete  $y$  output, give most probable one

$$p(x,y) \rightarrow p(y | x) \rightarrow \hat{y} = \arg \max_m p(y = m | x)$$



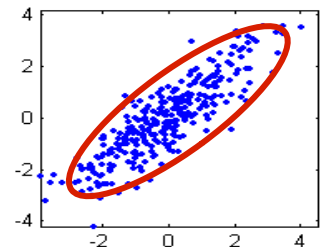
- **Regression:** have  $p(x,y)$  and given  $x$ . Asked for a scalar  $y$  output, give most probable or expected one



$$\hat{y} = \begin{cases} \arg \max_y p(y | x) \\ E_{p(y|x)} \{y\} \end{cases}$$

- **Anomaly Detection:** if have  $p(x,y)$  and given both  $x,y$ . Asked if it is similar  $\rightarrow$  threshold

$$p(x,y) \geq \text{threshold} \rightarrow \{normal, anomaly\}$$



# Machine Learning 4771

Instructor: Tony Jebara

# Topic 8

- Discrete Probability Models
- Independence
- Bernoulli Distribution
- Text: Naïve Bayes
- Categorical / Multinomial Distribution
- Text: Bag of Words



# Bernoulli Probability Models

- Bernoulli: recall binary (coin flip) probability, just 1x2 table

$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0,1] \quad x \in \{0,1\}$$

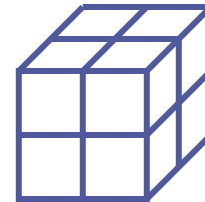
x=0	x=1
0.73	0.27

- Multidimensional Bernoulli: multiple binary events

$$p(x_1, x_2)$$

	x <sub>2</sub> =0	x <sub>2</sub> =1
x <sub>1</sub> =0	0.4	0.1
x <sub>1</sub> =1	0.3	0.2

$$p(x_1, x_2, x_3)$$



- Why do we write these as an equations instead of tables?





# Bernoulli Probability Models

- Bernoulli: recall binary (coin flip) probability, just 1x2 table

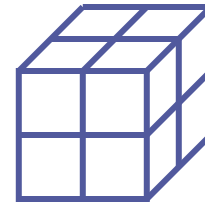
$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0,1] \quad x \in \{0,1\}$$

x=0	x=1
0.73	0.27

- Multidimensional Bernoulli: multiple binary events

	x <sub>2</sub> =0	x <sub>2</sub> =1
x <sub>1</sub> =0	0.4	0.1
x <sub>1</sub> =1	0.3	0.2

$$p(x_1, x_2, x_3)$$



- Why do we write these as an equations instead of tables?

- To do things like... maximum likelihood...
- Fill in the table so that it matches real data...
- Example: coin flips H,H,T,T,T,H,T,H,H,H ???

x=T	x=H



# Bernoulli Probability Models

- Bernoulli: recall binary (coin flip) probability, just 1x2 table

$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0,1] \quad x \in \{0,1\}$$

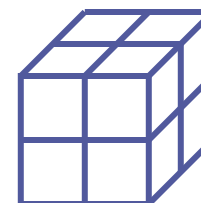
x=0	x=1
0.73	0.27

- Multidimensional Probability Table: multiple binary events

$$p(x_1, x_2)$$

	x <sub>2</sub> =0	x <sub>2</sub> =1
x <sub>1</sub> =0	0.4	0.1
x <sub>1</sub> =1	0.3	0.2

$$p(x_1, x_2, x_3)$$



- Why do we write these as an equations instead of tables?
- To do things like... maximum likelihood...
- Fill in the table so that it matches real data...
- Example: coin flips H,H,T,T,T,H,T,H,H,H
- Why is this correct?

x=T	x=H
0.4	0.6

# Bernoulli Maximum Likelihood

•Bernoulli:

$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0, 1] \quad x \in \{0, 1\}$$

•Log-Likelihood (IID):  $\sum_{i=1}^N \log p(x_i | \alpha) = \sum_{i=1}^N \log \alpha^{x_i} (1 - \alpha)^{1-x_i}$

•Gradient=0:

$$\frac{\partial}{\partial \alpha} \sum_{i=1}^N \log \alpha^{x_i} (1 - \alpha)^{1-x_i} = 0$$

$$\frac{\partial}{\partial \alpha} \sum_{i=1}^N x_i \log \alpha + (1 - x_i) \log(1 - \alpha) = 0$$

$$\frac{\partial}{\partial \alpha} \sum_{i \in \text{class1}} \log \alpha + \sum_{i \in \text{class0}} \log(1 - \alpha) = 0$$

$$\sum_{i \in \text{class1}} \frac{1}{\alpha} - \sum_{i \in \text{class0}} \frac{1}{1-\alpha} = 0$$

$$N_1 \frac{1}{\alpha} - N_0 \frac{1}{1-\alpha} = 0$$

$$N_1 (1 - \alpha) - N_0 \alpha = 0$$

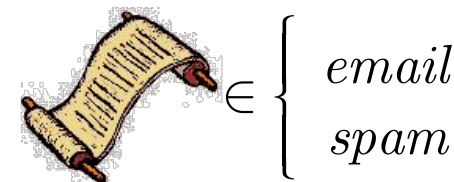
$$N_1 - (N_1 + N_0) \alpha = 0$$

$$\alpha = \frac{N_1}{N_1 + N_0}$$

x=0	x=1
$\frac{N_0}{N_0 + N_1}$	$\frac{N_1}{N_0 + N_1}$

# Text Modeling via Naïve Bayes

- Naïve Bayes: the simplest model of text



- There are about 50,000 words in English
- Each document is  $D=50,000$  dimensional binary vector  $\vec{x}_i$
- Each dimension is a word, set to 1 if word in the document

**Dim1: "the" = 1**

**Dim2: "hello" = 0**

**Dim3: "and" = 1**

**Dim4: "happy" = 1**

...

- Naïve Bayes: assumes each word is independent

$$p(\vec{x}) = p(\vec{x}(1), \dots, \vec{x}(D)) = \prod_{d=1}^D p(\vec{x}(d))$$

$$= \prod_{d=1}^D \bar{\alpha}(d)^{\vec{x}(d)} (1 - \bar{\alpha}(d))^{(1-\vec{x}(d))}$$

- Each 1 dimensional  $\alpha(d)$  is a Bernoulli parameter
- The whole  $\alpha$  vector is multivariate Bernoulli

# Text Modeling via Naïve Bayes

- Maximum likelihood: assume we have several IID vectors
- Have N documents, each a 50,000 dimension binary vector
- Each dimension is a word, set to 1 if word in the document

		$\vec{x}_1$	$\vec{x}_2$	$\vec{x}_3$	$\vec{x}_4$
<b>Dim1:</b>	<b>"the"</b>	<b>=</b>	<b>1</b>	<b>0</b>	<b>1</b>
<b>Dim2:</b>	<b>"hello"</b>	<b>=</b>	<b>0</b>	<b>1</b>	<b>0</b>
<b>Dim3:</b>	<b>"and"</b>	<b>=</b>	<b>1</b>	<b>1</b>	<b>0</b>
<b>Dim4:</b>	<b>"happy"</b>	<b>=</b>	<b>1</b>	<b>0</b>	<b>0</b>

- Likelihood =  $\prod_{i=1}^N p(\vec{x}_i | \vec{\alpha}) = \prod_{i=1}^N \prod_{d=1}^{50000} \vec{\alpha}(d)^{\vec{x}_i(d)} (1 - \vec{\alpha}(d))^{(1 - \vec{x}_i(d))}$
- Max likelihood solution: for each word d count number of documents it appears in divided by total N documents  $\vec{\alpha}(d) = \frac{N_d}{N}$
- To classify a new document x, build two models  $\alpha_{+1}$   $\alpha_{-1}$  & compare  $prediction = \arg \max_{y \in \{\pm 1\}} p(\vec{x} | \vec{\alpha}_y)$

# Categorical Probability Models



- **Categorical**: a distribution over a single multi-category event

1	2	3	4	5	6
$\vec{\alpha}(1)$	$\vec{\alpha}(2)$	$\vec{\alpha}(3)$	$\vec{\alpha}(4)$	$\vec{\alpha}(5)$	$\vec{\alpha}(6)$

$$p(x) = \prod_{m=1}^M \vec{\alpha}(m)^{\vec{x}(m)} \quad \sum_m \vec{\alpha}(m) = 1 \quad \vec{x} \in \mathbb{B}^M ; \sum_m \vec{x}(m) = 1$$

- Encode events as binary indicator vectors

$\vec{x}(1)$	$\vec{x}(2)$	$\vec{x}(3)$	$\vec{x}(4)$	$\vec{x}(5)$	$\vec{x}(6)$
--------------	--------------	--------------	--------------	--------------	--------------

- Related to the more general *multinomial* distribution
- Find  $\alpha$  using Maximum Likelihood (with IID assumption):

$$\sum_{i=1}^N \log p(\vec{x}_i | \vec{\alpha}) = \sum_{i=1}^N \log \prod_{m=1}^M \vec{\alpha}(m)^{\vec{x}_i(m)} = \sum_{i=1}^N \sum_{m=1}^M \vec{x}_i(m) \log(\vec{\alpha}(m))$$

- Can't just take gradient over  $\alpha$ , use sum=1 constraint:

- Insert constraint using Lagrange multipliers

$$\frac{\partial}{\partial \alpha_q} \sum_{i=1}^N \sum_{m=1}^M \vec{x}_i(m) \log(\vec{\alpha}(m)) - \lambda \left( \sum_{m=1}^M \vec{\alpha}(m) - 1 \right) = 0$$

$$\sum_{i=1}^N \left( \vec{x}_i(q) \frac{1}{\vec{\alpha}(q)} \right) - \lambda = 0 \quad \Rightarrow \quad \vec{\alpha}(q) = \frac{1}{\lambda} \sum_{i=1}^N \vec{x}_i(q)$$

# Categorical Maximum Likelihood

- Taking the gradient with Lagrangian gives this formula for each  $q$ :

$$\vec{\alpha}(q) = \frac{1}{\lambda} \sum_{i=1}^N \vec{x}_i(q)$$

- Recall the constraint:  $\sum_m \vec{\alpha}(m) - 1 = 0$

- Plug in  $\alpha$ 's solution:  $\sum_m \frac{1}{\lambda} \sum_{i=1}^N \vec{x}_i(m) - 1 = 0$

- Gives the lambda:  $\lambda = \sum_m \sum_{i=1}^N \vec{x}_i(m)$

- Final answer: 
$$\vec{\alpha}(q) = \frac{\sum_{i=1}^N \vec{x}_i(q)}{\sum_m \sum_{i=1}^N \vec{x}_i(m)} = \frac{N_q}{N}$$

- Example: Rolling dice

1,6,2,6,3,6,4,6,5,6

x=1	x=2	x=3	x=4	x=5	x=6
0.1	0.1	0.1	0.1	0.1	0.5

# Multinomial Probability Model

- The multinomial is a categorical over *counts* of events

Dice: 1,3,1,4,6,1,1      Word Dice: the, dog, jumped, the

- Say document  $i$  has  $W_i=2000$  words, each an IID dice roll

$$p(doc_i) = p(\vec{x}_i^1, \vec{x}_i^2, \dots, \vec{x}_i^{W_i}) = \prod_{w=1}^{W_i} p(\vec{x}_i^w) \propto \prod_{w=1}^{W_i} \prod_{d=1}^D \vec{\alpha}(d)^{\vec{x}_i^w(d)}$$

- Get count of each time an event occurred

$$p(doc_i) \propto \prod_{w=1}^{W_i} \prod_{d=1}^D \vec{\alpha}(d)^{\vec{x}_i^w(d)} = \prod_{d=1}^D \vec{\alpha}(d)^{\sum_{w=1}^{W_i} \vec{x}_i^w(d)} = \prod_{d=1}^D \vec{\alpha}(d)^{\vec{X}_i(d)}$$

- BUT: order shouldn't matter when "counting" so multiply by # of possible choosings. Choosing  $X(1), \dots, X(D)$  from  $N$

$$\binom{W_i}{\vec{X}_i(1), \dots, \vec{X}_i(D)} = \frac{W_i!}{\prod_{d=1}^D \vec{X}_i(d)!} = \frac{\left(\sum_{d=1}^D \vec{X}_i(d)\right)!}{\prod_{d=1}^D \vec{X}_i(d)!}$$

- **Multinomial:** over discrete integer vectors  $X$  summing to  $W$

$$p(\vec{X}_i) = \frac{W!}{\prod_{d=1}^D \vec{X}_i(d)!} \prod_{d=1}^D \vec{\alpha}(d)^{\vec{X}_i(d)} \quad s.t. \sum_d \vec{\alpha}(d) = 1, \vec{X} \in \mathbb{Z}_+^D, \sum_{d=1}^D \vec{X}(d) = W$$



# Text Modeling via Multinomial

- Also known as the bag-of-words model


 $\in \begin{cases} email \\ spam \end{cases}$ 

- Each document is 50,000 dimensional vector
- Each dimension is a word, set to # times word in doc

		$X_1$	$X_2$	$X_3$	$X_4$
<b>Dim1:</b>	<b>"the"</b>	<b>= 9</b>	<b>3</b>	<b>1</b>	<b>0</b>
<b>Dim2:</b>	<b>"hello"</b>	<b>= 0</b>	<b>5</b>	<b>3</b>	<b>0</b>
<b>Dim3:</b>	<b>"and"</b>	<b>= 6</b>	<b>2</b>	<b>2</b>	<b>2</b>
<b>Dim4:</b>	<b>"happy"</b>	<b>= 2</b>	<b>5</b>	<b>1</b>	<b>0</b>

- Each document is a vector of multinomial counts

$$p(doc_i) = p(\vec{X}_i) = \frac{\left(\sum_{d=1}^D \vec{X}_i(d)\right)!}{\prod_{d=1}^D \vec{X}_i(d)!} \prod_{d=1}^D \vec{\alpha}(d)^{\vec{X}_i(d)} \quad \sum_d \vec{\alpha}(d) = 1 \quad X \in \mathbb{Z}_+^D$$

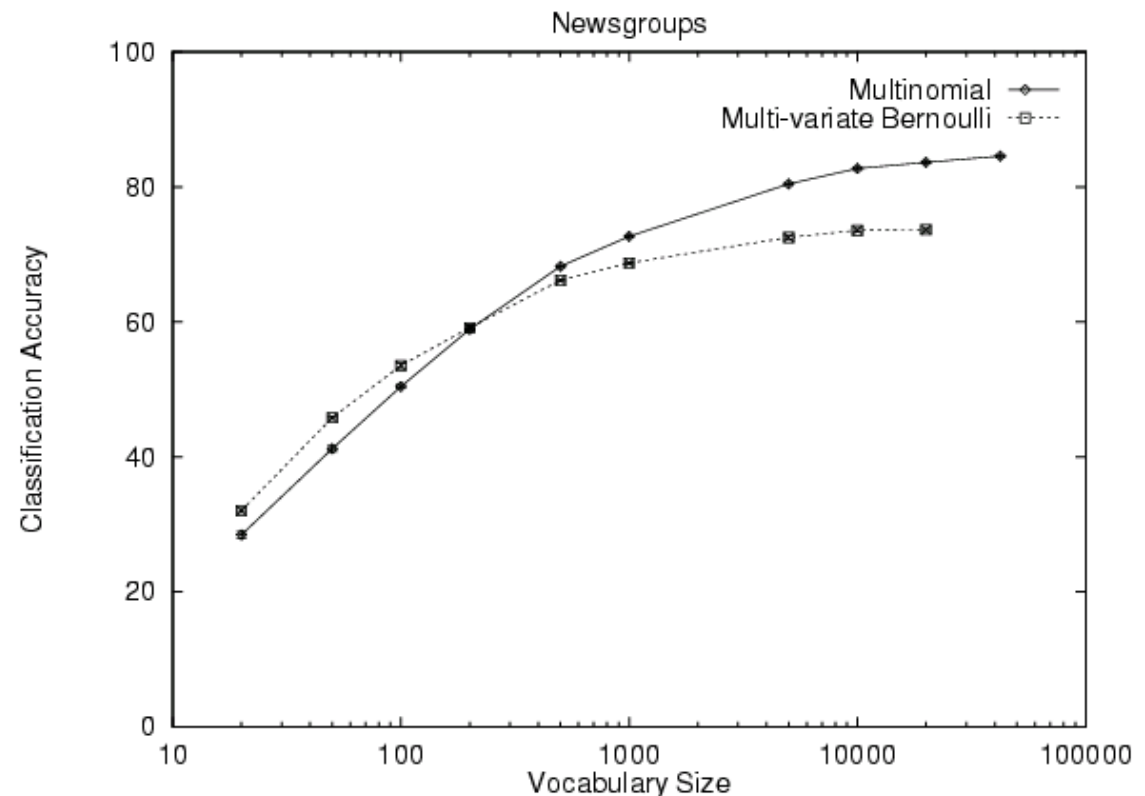
- Log-likelihood:  $l(\vec{\alpha}) = \sum_{i=1}^N \log p(\vec{X}_i) = \sum_{i=1}^N \log \frac{\left(\sum_{d=1}^D \vec{X}_i(d)\right)!}{\prod_{d=1}^D \vec{X}_i(d)!} \prod_{d=1}^D \vec{\alpha}(d)^{\vec{X}_i(d)}$   

$$= \sum_{i=1}^N \sum_{d=1}^D \vec{X}_i(d) \log \vec{\alpha}(d) + const$$

- Find  $\alpha$  just like the multinomial maximum likelihood formula!

# Text Modeling Experiments

- For text modeling (McCallum & Nigam '98)
  - Bernoulli better for small vocabulary
  - Multinomial better for large vocabulary



# Machine Learning 4771

Instructor: Tony Jebara

# Topic 9

- Continuous Probability Models
- Gaussian Distribution
- Maximum Likelihood Gaussian
- Sampling from a Gaussian

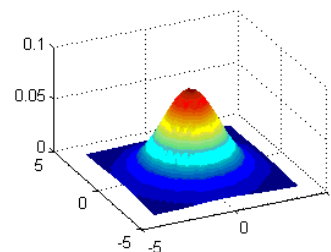
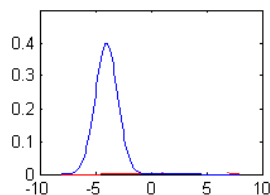
# Continuous Probability Models

- Probabilities can have both discrete & continuous variables
- We will discuss:
  - 1) discrete probability tables

x=T	x=H
0.4	0.6

x=1	x=2	x=3	x=4	x=5	x=6
0.1	0.1	0.1	0.1	0.1	0.5

## 2) continuous probability distributions

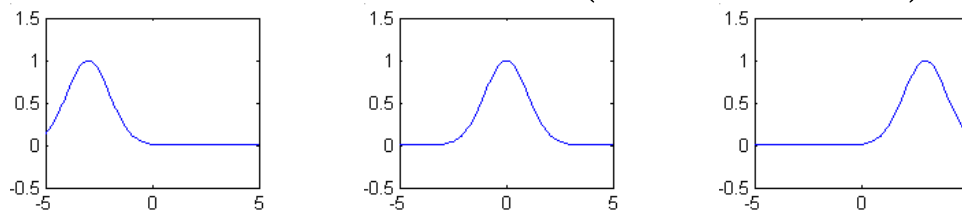


- Most popular continuous distribution = Gaussian

# Gaussian Distribution

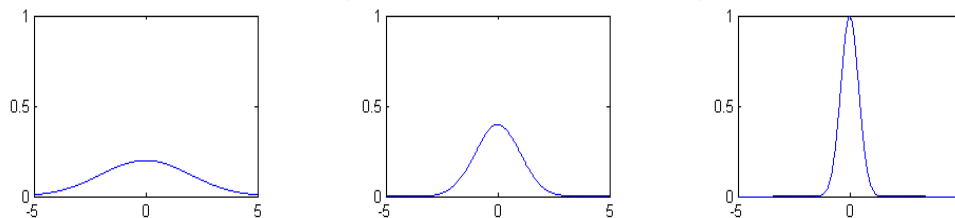
- Recall 1-dimensional Gaussian with mean parameter  $\mu$  translates Gaussian left & right

$$p(x | \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right)$$



- Can also have variance parameter  $\sigma^2$  widens or narrows the Gaussian

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$



Note:  $\int_{x=-\infty}^{\infty} p(x) dx = 1$

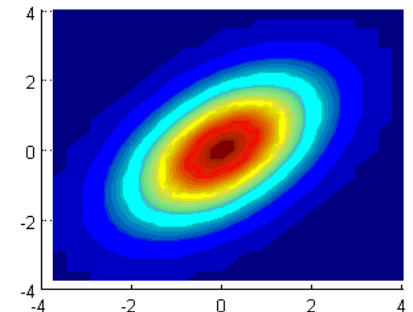
# Multivariate Gaussian

- Gaussian can extend to D-dimensions
- Gaussian mean parameter  $\mu$  vector, it translates the bump
- Covariance matrix  $\Sigma$  stretches and rotates bump

$$p(\vec{x} \mid \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)$$

- Mean is any real vector
- Max and expectation =  $\mu$
- Variance parameter is now  $\Sigma$  matrix
- Covariance matrix is positive definite
- Covariance matrix is symmetric
- Need matrix **inverse** (inv)
- Need matrix **determinant** (det)
- Need matrix **trace** operator (trace)

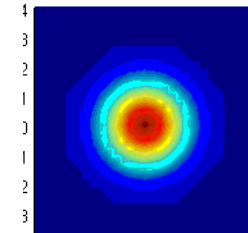
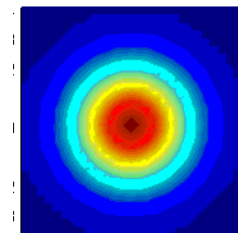
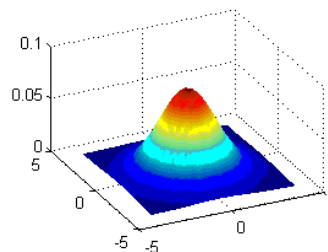
$$\vec{x} \in \mathbb{R}^D, \vec{\mu} \in \mathbb{R}^D, \Sigma \in \mathbb{R}^{D \times D}$$



# Multivariate Gaussian

- Spherical:

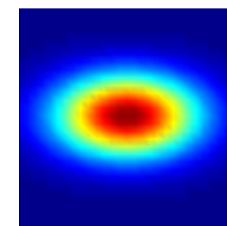
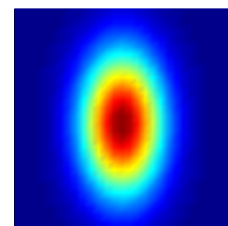
$$\Sigma = \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 & 0 \\ 0 & \sigma^2 & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix}$$



- Diagonal Covariance:

dimensions of  $\mathbf{x}$  are independent  
product of multiple 1d Gaussians

$$p(\vec{x} \mid \vec{\mu}, \Sigma) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi}\vec{\sigma}(d)} \exp\left(-\frac{(\vec{x}(d) - \vec{\mu}(d))^2}{2\vec{\sigma}(d)^2}\right)$$



$$\Sigma = \begin{bmatrix} \vec{\sigma}(1)^2 & 0 & 0 & 0 \\ 0 & \vec{\sigma}(2)^2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \vec{\sigma}(D)^2 \end{bmatrix}$$



# Max Likelihood Gaussian

- Have IID samples as vectors  $i=1..N$ :  $\mathcal{X} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$
- How do we recover the mean and covariance parameters?
- Standard approach: Maximum Likelihood (IID)
- Maximize probability of data given model (likelihood)

$$\begin{aligned} p(\mathcal{X} \mid \theta) &= p(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N \mid \theta) \\ &= \prod_{i=1}^N p(\vec{x}_i \mid \vec{\mu}_i, \Sigma_i) \quad \text{independent Gaussian samples} \\ &= \prod_{i=1}^N p(\vec{x}_i \mid \vec{\mu}, \Sigma) \quad \text{identically distributed} \end{aligned}$$

- Instead, work with maximum of log-likelihood

$$\sum_{i=1}^N \log p(\vec{x}_i \mid \vec{\mu}, \Sigma) = \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu})\right)$$

# Max Likelihood Gaussian

•Max over  $\mu$  
$$\frac{\partial}{\partial \mu} \left( \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp \left( -\frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right) \right) = 0$$

$$\frac{\partial}{\partial \mu} \left( \sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right) = 0$$

$$\frac{\partial \vec{x}^T \vec{x}}{\partial \vec{x}} = 2\vec{x}^T$$

$$\sum_{i=1}^N (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} = \vec{0}$$

see Jordan Ch. 12, get sample mean...

$$\vec{\mu} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$$

•For  $\Sigma$  need Trace operator:  $tr(A) = tr(A^T) = \sum_{d=1}^D A(d, d)$

and several properties:

$$tr(AB) = tr(BA)$$

$$tr(BAB^{-1}) = tr(A)$$

$$tr(\vec{x}\vec{x}^T A) = tr(\vec{x}^T A \vec{x}) = \vec{x}^T A \vec{x}$$

# Max Likelihood Gaussian

- Likelihood rewritten in trace notation:

$$\begin{aligned}
 l &= \sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \\
 &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^N \text{tr} \left[ (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right] \\
 &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^N \text{tr} \left[ (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} \right] \\
 &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |A| - \frac{1}{2} \sum_{i=1}^N \text{tr} \left[ (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T A \right]
 \end{aligned}$$

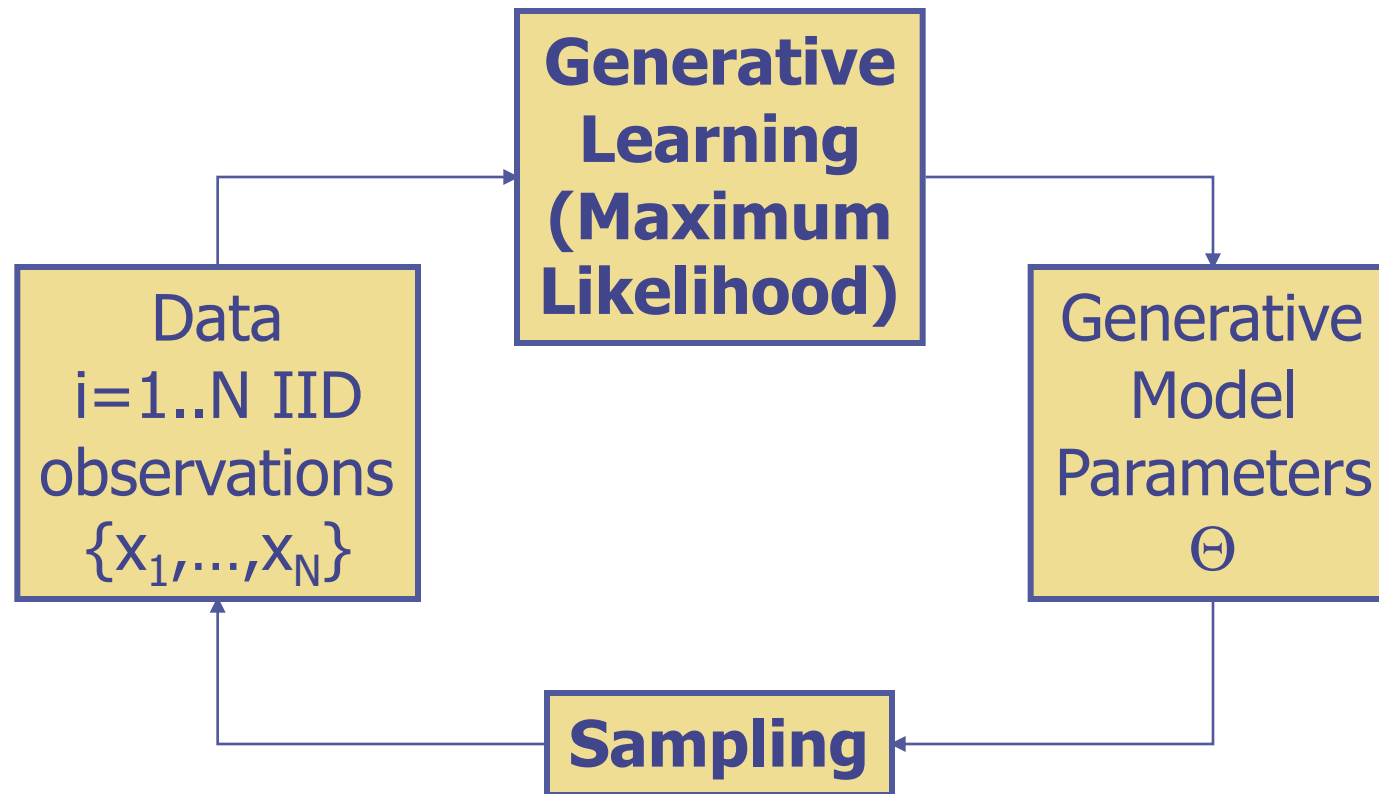
- Max over  $A = \Sigma^{-1}$   
use properties:

$$\begin{aligned}
 \frac{\partial l}{\partial A} &= -0 + \frac{N}{2} \left( A^{-1} \right)^T - \frac{1}{2} \sum_{i=1}^N \left[ (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T \right]^T \\
 &= \frac{N}{2} \Sigma - \frac{1}{2} \sum_{i=1}^N (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T
 \end{aligned}$$

$\frac{\partial \log |A|}{\partial A} = (A^{-1})^T$ 
 $\frac{\partial \text{tr}[BA]}{\partial A} = B^T$

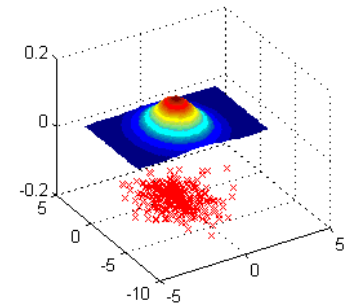
- Get sample covariance:  $\frac{\partial l}{\partial A} = 0 \rightarrow \Sigma = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T$

# Sampling & Max Likelihood



# Sampling from a Gaussian

- Fit Gaussian to data, how is this Generative?



# Sampling from a Gaussian

- Fit Gaussian to data, how is this Generative?

- Sampling! Generating discrete data easy:

0.73	0.1	0.17
------	-----	------

- Assume we can do uniform sampling:

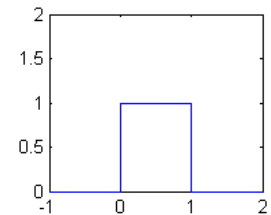
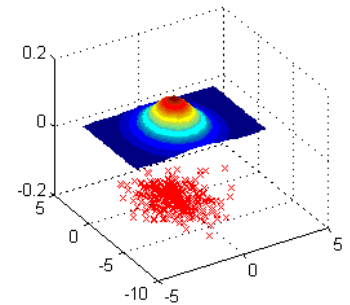
i.e. rand between (0,1)

if  $0.00 \leq \text{rand} < 0.73$  get A

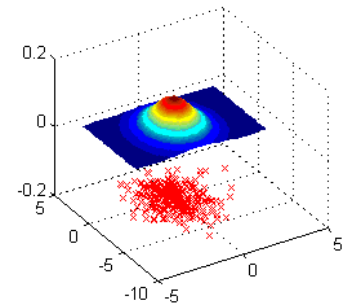
if  $0.73 \leq \text{rand} < 0.83$  get B

if  $0.83 \leq \text{rand} < 1.00$  get C

- What are we doing?



# Sampling from a Gaussian



- Fit Gaussian to data, how is this Generative?

- Sampling! Generating discrete data easy:

0.73	0.1	0.17
------	-----	------

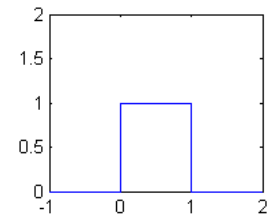
- Assume we can do uniform sampling:

i.e. rand between (0,1)

if  $0.00 \leq \text{rand} < 0.73$  get A

if  $0.73 \leq \text{rand} < 0.83$  get B

if  $0.83 \leq \text{rand} < 1.00$  get C

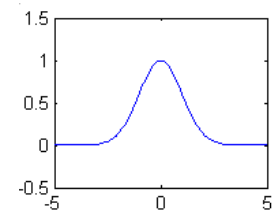


0.73	0.83	1.00
------	------	------

- What are we doing?

Sum up the Probability Density Function (PDF)  
to get Cumulative Density Function (CDF)

- For 1d Gaussian, Integrate Probability Density Function get Cumulative Density Function  
Integral is like summing many discrete bars



# Sampling from a Gaussian

- Integrate 1d Gaussian to get CDF:

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$$

$$F(x) = \int_{-\infty}^x p(t) dt = \frac{1}{2} \operatorname{erf}\left(\frac{1}{\sqrt{2}}x\right) + \frac{1}{2}$$

- If sample from uniform, get:  $u \sim \text{uniform}(0,1)$

- Compute mapping:  $x = F^{-1}(u) = \sqrt{2} \operatorname{erfinv}(2u - 1)$

- This is a Gaussian sample:  $x \sim N(x | 0, 1)$

- For D-dimensional Gaussian  $N(\mathbf{z} | 0, I)$  concatenate samples:

$$\vec{x} = [\vec{x}(1) \dots \vec{x}(D)]^T \sim p(\vec{x} | 0, I) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \vec{x}(d)^2\right)$$

- For  $N(\mathbf{z} | \vec{\mu}, \Sigma)$ , add mean & multiply by root cov

$$\vec{z} = \Sigma^{1/2} \vec{x} + \vec{\mu} \sim p(\vec{z} | \vec{\mu}, \Sigma)$$

- Example code: `gendata.m`

