

Machine Learning

4771

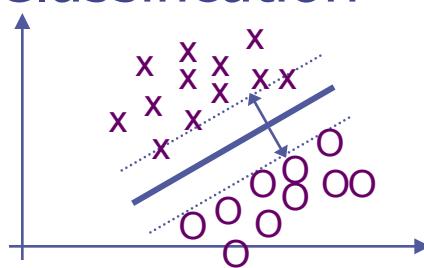
Instructor: Tony Jebara

Topic 7

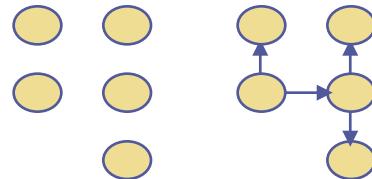
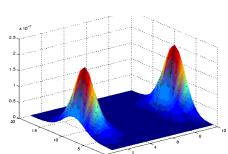
- Unsupervised Learning
- Statistical Perspective
- Probability Models
- Discrete & Continuous: Gaussian, Bernoulli, Multinomial
- Maximum Likelihood → Logistic Regression
- Conditioning, Marginalizing, Bayes Rule, Expectations
- Classification, Regression, Detection
- Dependence/Independence
- Maximum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

Unsupervised Learning

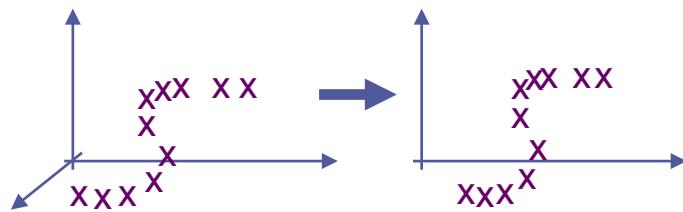
Classification



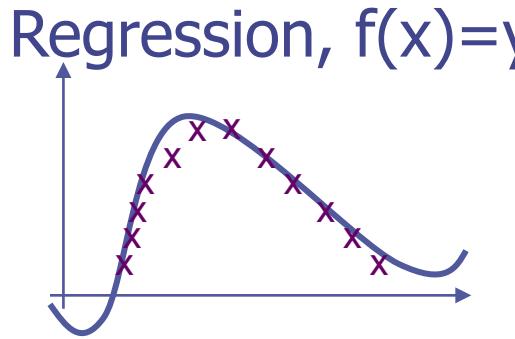
Density/Structure Estimation



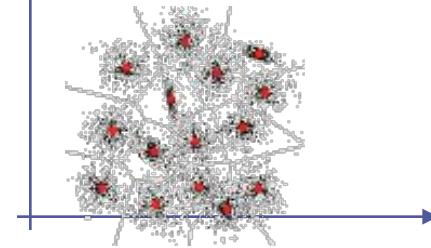
Feature Selection



Regression, $f(x)=y$



Clustering



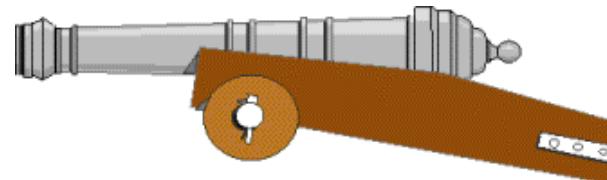
Anomaly Detection



Supervised
Unsupervised
(can help
supervised)

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 $x, y, \text{etc...}$
- Also, what if we don't have a true function?
- Example of Projectile Cannon (c.f. Distal Learning)



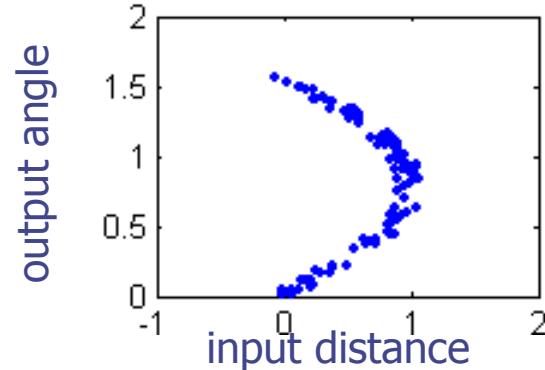
- Would like to train a regression function to control a cannon's angle of fire (y) given target distance (x)

Statistical Perspective

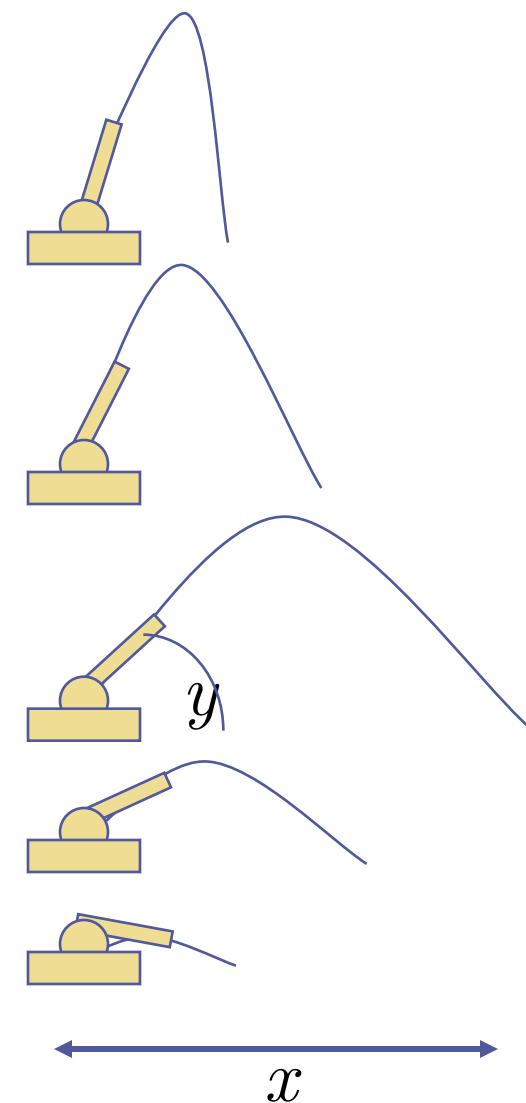
- Example of Projectile Cannon
(45 degree problem)

x = input target distance
 y = output cannon angle

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

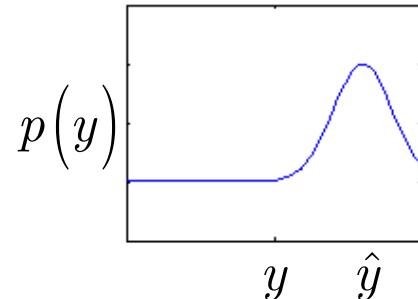


- What does least squares do?
- Conditional statistical models address this problem...



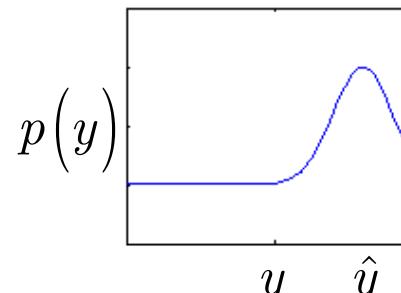
Probability Models

- 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:

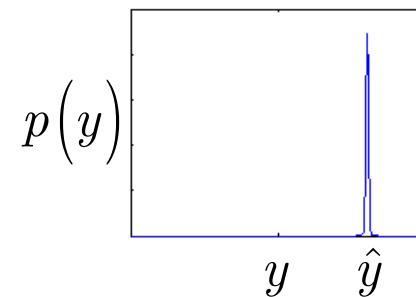


- $p(y)$ subsumes or is a superset of \hat{y}
- Why is this representation for our answer more general?

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} !

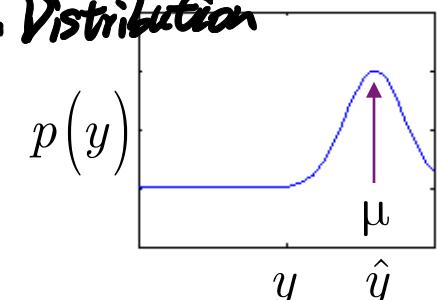
$$\hat{y} \Leftrightarrow p(y) = \delta(y - \hat{y})$$



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
 $p(y) \rightarrow p(y | \Theta)$ *Normal Distribution = Gaussian Distribution*
- E.g.: 1-dim Gaussian distribution
 'given' 'mean' parameter μ :

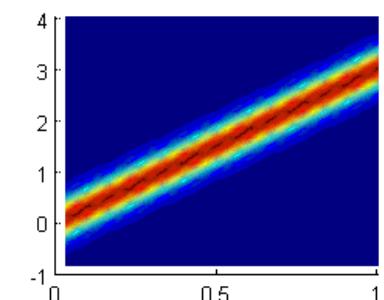
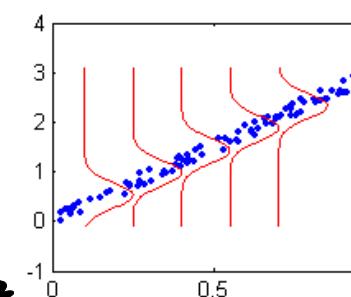
$$p(y | \mu) = N(y | \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\mu)^2}$$



- Want mean centered on $f(x)$'s value $p(y) = N(y | f(x))$
- Now, linear regression is:

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

$p(y|w) = N(y | u, \sigma^2) \Rightarrow f(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp(-\frac{(x-u)^2}{2\sigma^2}) = \text{pdf}$



$$N(\mu, 1) \Rightarrow f(x) = \frac{1}{\sqrt{2\pi}} \cdot \exp(-\frac{x^2}{2}) = \text{pdf}$$

$$N(\mu, 1) \Rightarrow \text{pdf} = p(y|x) = \frac{1}{\sqrt{2\pi}} \cdot \exp(-\frac{(y-f(x))^2}{2}) ; \text{when } u=f(x)=\theta^T x + b$$

Tony Jebara, Columbia University

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!

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

$$\stackrel{\text{pdf: } N(\mu, 1)}{=} -N \log(\sqrt{2\pi}) - \sum_{i=1}^N \frac{1}{2} (y_i - f(x_i))^2$$

$$L(\Theta) = \prod_{i=1}^N p(y_i | x_i) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \cdot \exp(-\frac{(y_i - f(x_i))^2}{2}) \quad / \log L(\Theta) = \log \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \cdot \exp(-\frac{(y_i - f(x_i))^2}{2})$$

$$= \frac{-N}{2} \cdot \log 2\pi + \sum_{i=1}^N -\frac{(y_i - f(\theta))^2}{2} = -\frac{N}{2} \log 2\pi - \sum_{i=1}^N \frac{(y_i - f(\theta))^2}{2}$$

Tony Jebara, Columbia University

Probability Models

- Can extend probability model to **2 bumps**:

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

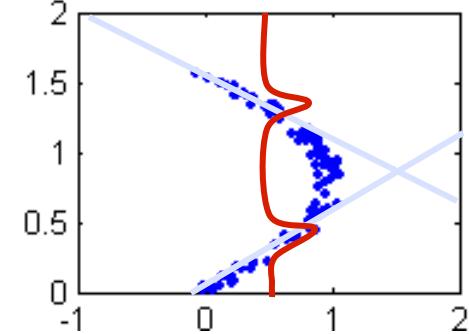
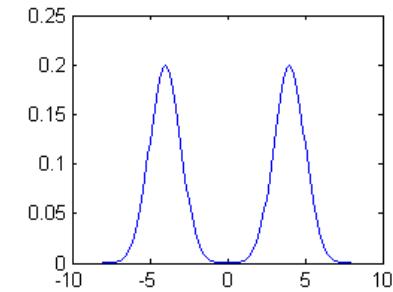
- Each mean can be a linear regression fn.

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

- 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



- 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 α

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

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

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

$\text{L}(y)$

$$\prod_{i=1}^N p(y_i | x_i) = \prod_{i=1}^N f(x_i)^{y_i} (1 - f(x_i))^{1-y_i}$$

$$\begin{aligned}
 & \log \prod_{i=1}^N p(y_i | x_i) = \sum_{i=1}^N y_i \cdot \log f(x_i) + (1-y_i) \cdot \log (1-f(x_i)) \quad ; \quad f(x_i) = \frac{1}{1+e^{-\theta^T x_i}} \quad \frac{\partial f(x_i)}{\partial \theta} = f(x_i)(1-f(x_i))(x_i). \\
 & \frac{\partial \log \prod_{i=1}^N p(y_i | x_i)}{\partial \theta} = \sum_{i=1}^N y_i \cdot \frac{f(x_i) \cdot (1-f(x_i))}{f(x_i)} \cdot x_i + (1-y_i) \cdot \frac{f(x_i) \cdot (1-f(x_i))}{1-f(x_i)} \cdot x_i \\
 & = \sum_{i=1}^N [y_i(1-f(x_i)) + (1-y_i)f(x_i)] \cdot x_i \quad \Rightarrow \text{Gradient Descent}
 \end{aligned}$$

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):

$$\begin{aligned}
 \sum_{i=1}^N \log p(y_i | x_i) &= \sum_{i=1}^N \log f(x_i)^{y_i} (1 - f(x_i))^{1-y_i} \\
 &= \sum_{i=1}^N y_i \log f(x_i) + (1 - y_i) \log (1 - f(x_i)) \\
 &= \sum_{i \in \text{class } 1} \log f(x_i) + \sum_{i \in \text{class } 0} \log (1 - f(x_i))
 \end{aligned}$$

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

- Use sigmoid or logistic again...

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

- Called logistic regression \rightarrow new loss function

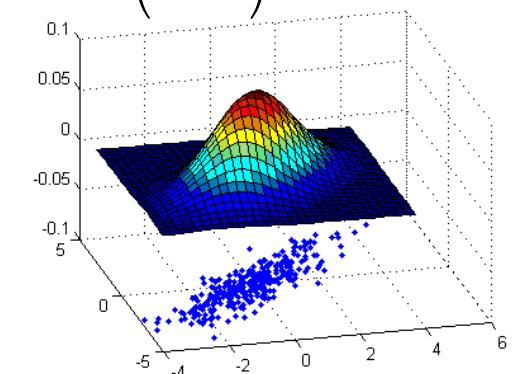
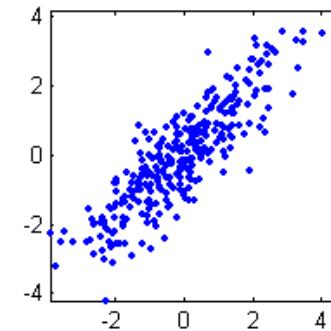
- Do gradient descent, similar to logistic output neural net!

- Can also handle multi-layer $f(x)$ and do backprop again!

Generative Probability Models

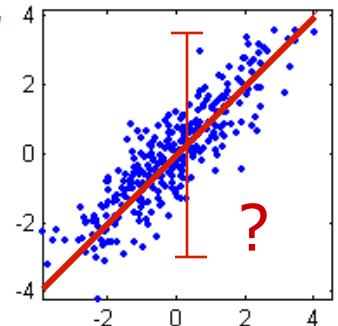
- Idea: Extend probability to describe *both* X and Y
- Find probability density function over both: $p(x, y)$

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:

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

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

- 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:

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

$$\sum_y \sum_x p(x, y) = 1$$

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 X being data and θ being a model

posterior
Likelihood

$$\text{posterior} \xrightarrow{\text{likelihood}} p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})} \xleftarrow{\text{prior}} \text{evidence}$$



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

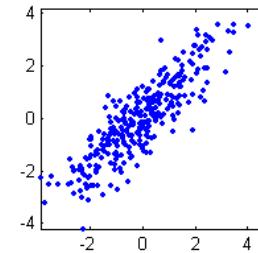
$$\begin{aligned} Var\{x\} &= E\left\{(x - E\{x\})^2\right\} = E\left\{x^2 - 2xE\{x\} + E\{x\}^2\right\} \\ &= E\{x^2\} - 2E\{x\}E\{x\} + E\{x\}^2 = E\{x^2\} - E\{x\}^2 \end{aligned}$$

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

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

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

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

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

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

More Properties of PDFs

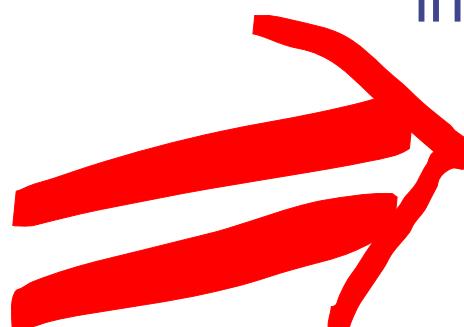
- Independence: probabilities of independent variables multiply. Denote with the following notation:

$$\begin{aligned} x \perp\!\!\!\perp y &\rightarrow p(x, y) = p(x)p(y) \\ x \perp\!\!\!\perp y &\rightarrow p(x | y) = p(x) \end{aligned}$$

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



$$\begin{aligned} x \perp\!\!\!\perp y | z &\rightarrow p(x | y, z) = p(x | z) \\ x \perp\!\!\!\perp y | z &\rightarrow p(x | y) \neq p(x) \end{aligned}$$

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 Θ) the probability of our data multiplies:
 - $p(x_1, \dots, x_N | \Theta) = \prod_{i=1}^N p_i(x_i | \Theta)$
- "Identically distributed" means that each marginal probability is the same for each data point
 - $p(x_1, \dots, x_N | \Theta) = \prod_{i=1}^N p_i(x_i | \Theta) = \prod_{i=1}^N p(x_i | \Theta)$

Θ : model

X : data

Tony Jebara, Columbia University

The IID Assumption

- Bayes rule says likelihood is probability of data given model

$$p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})}$$

likelihood → $p(\mathcal{X} | \theta)$ ← prior
 posterior → $p(\theta)$ ← evidence

- 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 $p(\mathcal{X} | \Theta)$ by maximum likelihood:

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

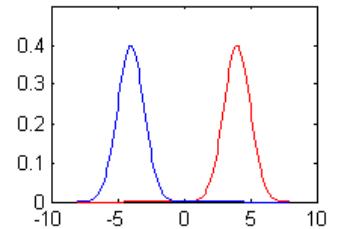
- Learn conditional $p(y | x, \Theta)$ by max conditional likelihood:

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

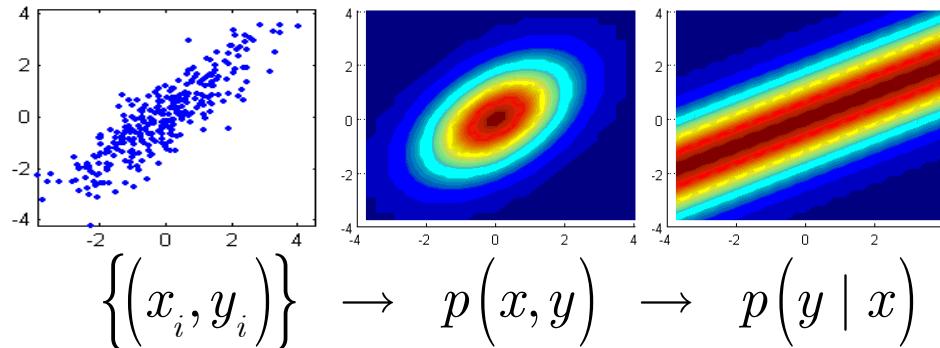
Uses of PDFs

- **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

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

