# ECE 5424: Introduction to Machine Learning

Topics:

– Probability Review

**Readings: Barber 8.1, 8.2**

Stefan Lee

Virginia Tech

# Project

- Groups of 1-3
  - we prefer teams of 2

- Deliverables:
  - Project proposal (NIPS format): 2 page, due Sept 21
  - Midway presentations (in class)
  - Final report: webpage with results

# Administrative

- HW1
  - Due on Wed 09/14, 11:55pm
  - https://inclass.kaggle.com/c/vt-ece-introduction-to-machine-learning-hw-1

- Project Proposal
  - Due: Wed 09/21, 11:55 pm
  - <=2pages, NIPS format

# Proposal

- 2 Page (NIPS format)
  - https://nips.cc/Conferences/2015/PaperInformation/StyleFiles

- Necessary Information:
  - Project title
  - Project idea.
    - This should be approximately two paragraphs.
  - Data set details
    - Ideally existing dataset. No data-collection projects.
  - Software
    - Which libraries will you use?
    - What will you write?
  - Papers to read.
    - Include 1-3 relevant papers. You will probably want to read at least one of them before submitting your proposal.
  - Teammate
    - Will you have a teammate? If so, what's the break-down of labor? Maximum team size is 3 students.
  - Mid-semester Milestone
    - What will you complete by the project milestone due date? Experimental results of some kind are expected here.

# Project

- Rules
  - Must be about machine learning
  - Must involve real data
    - Use your own data or take from class website
  - Can apply ML to your own research.
    - Must be done this semester.
  - OK to combine with other class-projects
    - Must declare to both course instructors
    - Must have explicit permission from BOTH instructors
    - Must have a sufficient ML component
  - Using libraries
    - No need to implement all algorithms
    - OK to use standard SVM, MRF, Decision-Trees, etc libraries
    - More thought + effort => More credit

# Project

- Main categories
  - <span style="color:red">Application/Survey</span>
    - Compare a bunch of existing algorithms on a new application domain of your interest
  - <span style="color:red">Formulation/Development</span>
    - Formulate a new model or algorithm for a new or old problem
  - <span style="color:red">Theory</span>
    - Theoretically analyze an existing algorithm

- Support
  - List of ideas, pointers to dataset/algorithms/code
    - https://filebox.ece.vt.edu/~f16ece5424/project.html
    - We will mentor teams and give feedback.

# Procedural View

- Training Stage:
  - Raw Data $\to$ x                        (Feature Extraction)
  - Training Data { (x,y) } $\to$ f             (Learning)


- Testing Stage
  - Raw Data $\to$ x                        (Feature Extraction)
  - Test Data x $\to$ f(x)         (Apply function, Evaluate error)

# Statistical Estimation View

- Probabilities to rescue:
  - x and y are *random variables*
  - D = $(x_1,y_1)$, $(x_2,y_2)$, …, $(x_N,y_N)$        ~ P(X,Y)

- IID: Independent Identically Distributed
  - Both training & testing data sampled IID from P(X,Y)
  - Learn on training set
  - Have some hope of *generalizing* to test set

# Plan for Today

- **Review of Probability**
  - Discrete vs Continuous Random Variables
  - PMFs vs PDF
  - Joint vs Marginal vs Conditional Distributions
  - Bayes Rule and Prior
  - Expectation, Entropy, KL-Divergence

# Probability

- The world is a very uncertain place

- 30 years of Artificial Intelligence and Database research danced around this fact

- And then a few AI researchers decided to use some ideas from the eighteenth century

# Probability

- A is non-deterministic event
  - Can think of A as a boolean-valued variable

- Examples
  - A = your next patient has cancer
  - A = Donald Trump Wins the 2016 Presidential Election

# Interpreting Probabilities

- What does P(A) mean?

- Frequentist View
  - limit N$\rightarrow \infty$ #(A is true)/N
  - limiting frequency of a repeating non-deterministic event

- Bayesian View
  - P(A) is your "belief" about A

- Market Design View
  - P(A) tells you how much you would bet

**67.8%** Clinton ▼ -1.9% [CHARTS ⬍]

**28.8%** Trump ▲ 1.2%

**Presidential Winner Odds (% chance of winning)**

Legend: Clinton, Sanders, OMalley, Trump, Rubio, Cruz, Bush, Christie, Kasich, Carson, Paul, Fiorina, Huckabee

X-axis: Jan 2016, Apr 2016, Jul 2016
Y-axis: 0, 25, 50, 75, 100

The
Axioms
Of
Probabi
lity

# Axioms of Probability

- 0<= P(A) <= 1
- P(empty-set) = 0
- P(everything) = 1
- P(A or B) = P(A) + P(B) – P(A and B)

# Interpreting the Axioms

- 0<= P(A) <= 1
- P(empty-set) = 0
- P(everything) = 1
- P(A or B) = P(A) + P(B) – P(A and B)

Event space of all possible worlds

Its area is 1

Worlds in which A is true

Worlds in which A is False

P(A) = Area of reddish oval

# Interpreting the Axioms

- 0<= P(A) <= 1
- P(empty-set) = 0
- P(everything) = 1
- P(A or B) = P(A) + P(B) – P(A and B)

The area of A can't get any smaller than 0

And a zero area would mean no world could ever have A true

# Interpreting the Axioms

- 0<= P(A) <= 1
- P(empty-set) = 0
- P(everything) = 1
- P(A or B) = P(A) + P(B) – P(A and B)

The area of A can't get any bigger than 1

And an area of 1 would mean all worlds will have A true

# Interpreting the Axioms

- 0<= P(A) <= 1

- P(empty-set) = 0

- P(everything) = 1

- P(A or B) = P(A) + P(B) – P(A and B)



Simple addition and subtraction

# Concepts

- Sample Space
  - Space of events

- Random Variables
  - Mapping from events to numbers
  - Discrete vs Continuous

- Probability
  - Mass vs Density

# Discrete Random Variables

$X$ ⟶ discrete random variable

$\mathcal{X}$ or Val(X) ⟶ sample space of possible outcomes, which may be finite or countably infinite

$x \in \mathcal{X}$ ⟶ outcome of sample of discrete random variable

$p(X = x)$ ⟶ probability distribution (probability mass function)

$p(x)$ ⟶ shorthand used when no ambiguity

$$0 \leq p(x) \leq 1 \text{ for all } x \in \mathcal{X} \qquad \sum_{x \in \mathcal{X}} p(x) = 1$$



$\mathcal{X} = \{1, 2, 3, 4\}$

*uniform distribution*          *degenerate distribution*

# Continuous Random Variables

- On board

# Concepts

- Expectation

- Variance

# Most Important Concepts

- Marginal distributions / Marginalization

- Conditional distribution / Chain Rule
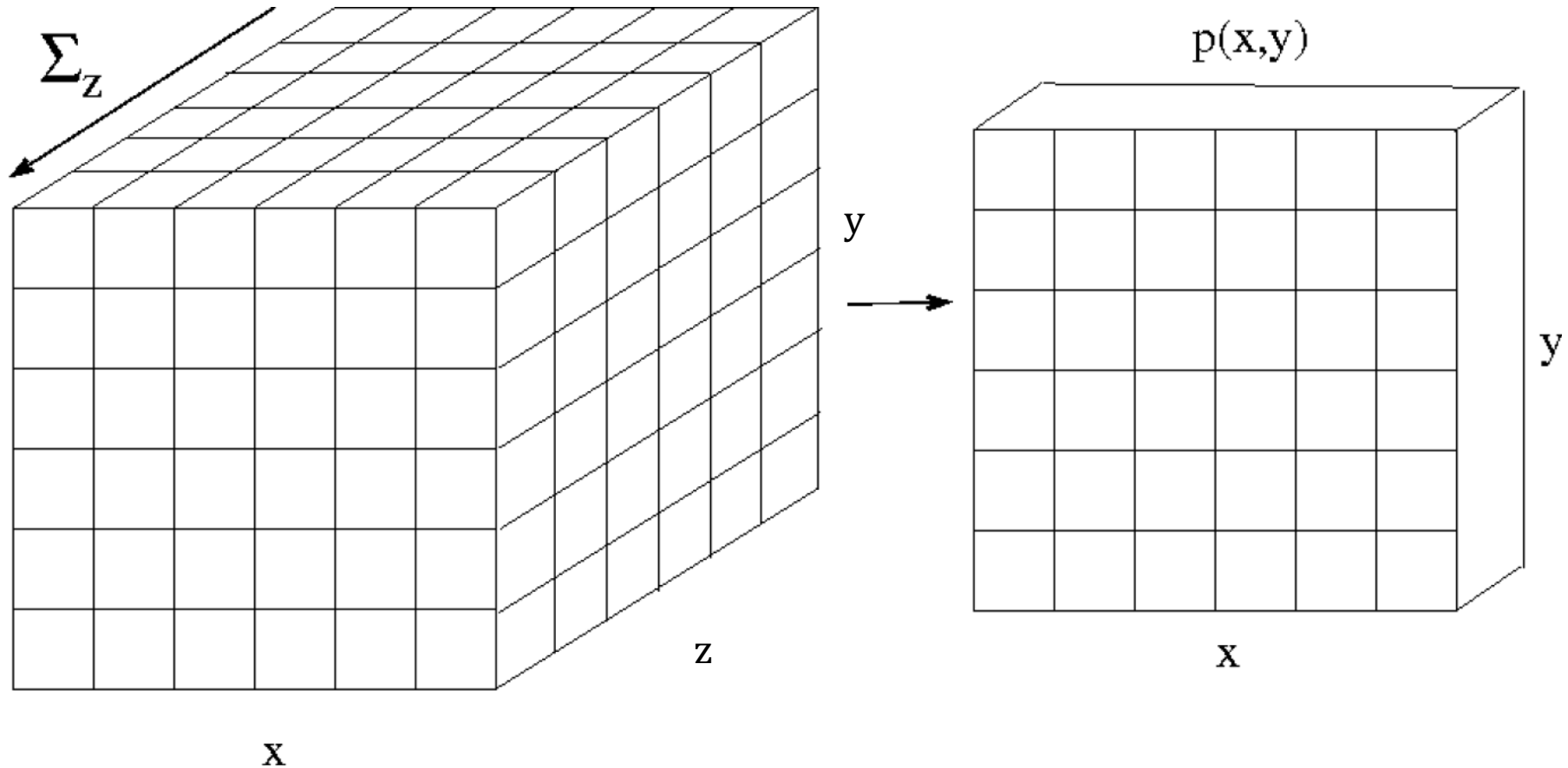
- Bayes Rule

# Joint Distribution



y

z

x

# Marginalization

- Marginalization
  - Events: P(A) = P(A and B) + P(A and not B)

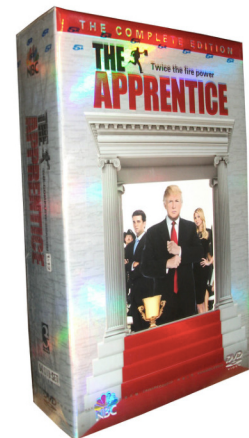  - Random variables $\quad P(X=x) = \sum_{y} P(X=x, Y=y)$

# Marginal Distributions



$$p(x,y) = \sum_{z \in \mathcal{Z}} p(x,y,z)$$

$$p(x) = \sum_{y \in \mathcal{Y}} p(x,y)$$

# Conditional Probabilities

- P(Y=y | X=x)

- What do you believe about Y=y, if I tell you X=x?

- P(Donald Trump Wins the 2016 Election)?

- What if I tell you:
  - He has the Republican nomination
  - His twitter history
  - The complete DVD set of The Apprentice

# Conditional Probabilities

- P(A | B) = In worlds that where B is true, fraction where A is true

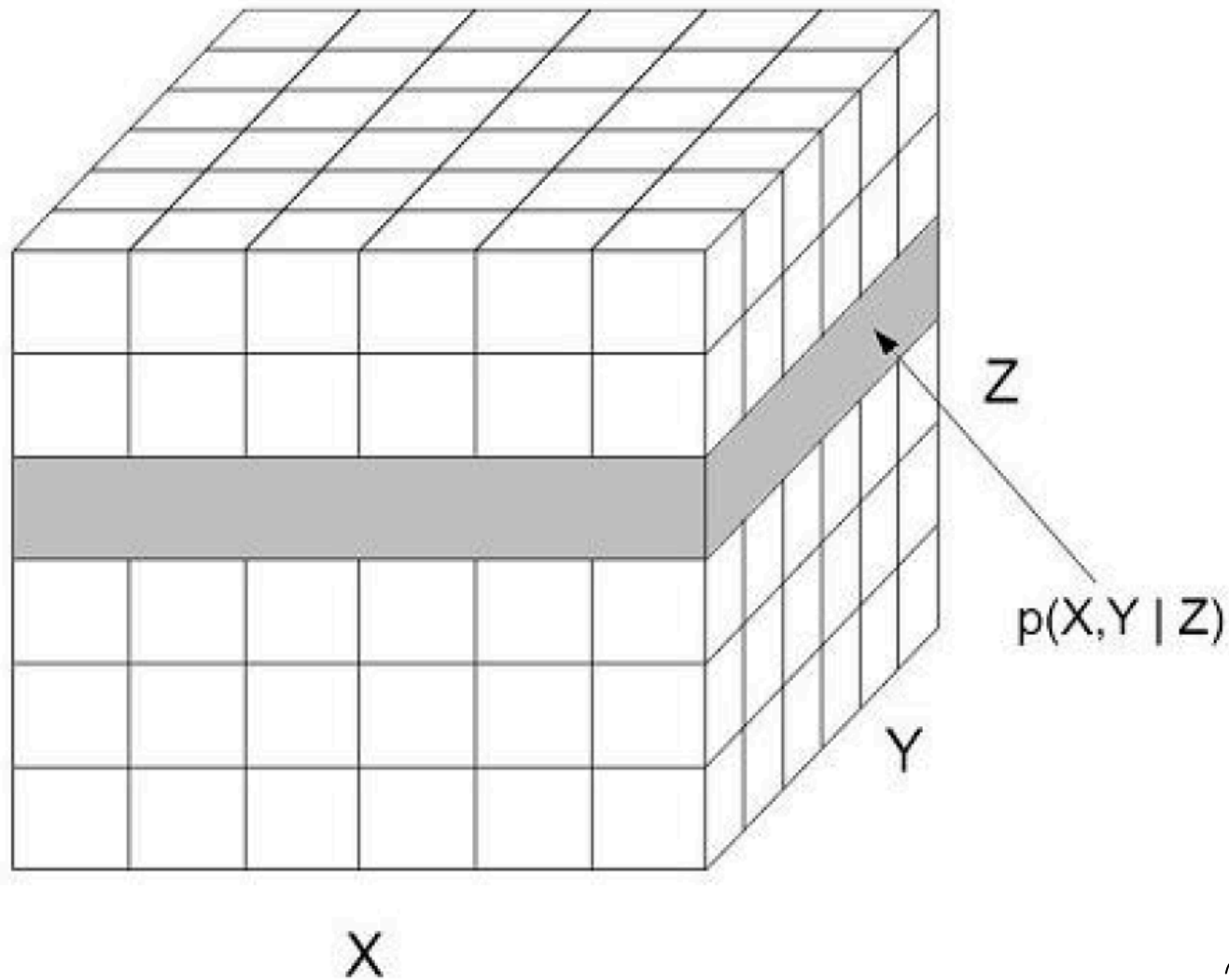- Example
  - H: "Have a headache"
  - F: "Coming down with Flu"



P(H) = 1/10
P(F) = 1/40
P(H|F) = 1/2

"Headaches are rare and flu is rarer, but if you're coming down with 'flu there's a 50-50 chance you'll have a headache."

# Conditional Distributions



$$p(x, y \mid Z = z) = \frac{p(x, y, z)}{p(z)}$$

# Conditional Probabilities

- Definition


- Corollary: Chain Rule

# Independent Random Variables

P(x,y)



$$X \perp Y$$

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

$$\text{for all } x \in \mathcal{X}, y \in \mathcal{Y}$$

# Marginal Independence

- **Sets** of variables **X**, **Y**

- **X** is independent of **Y**
  - Shorthand: $P \vdash (\mathbf{X} \perp \mathbf{Y})$

- **Proposition:** $P$ satisfies $(\mathbf{X} \perp \mathbf{Y})$ if and only if
  - P(**X=x**,**Y=y**) = P(**X=x**) P(**Y=y**), $\quad \forall\, x \in Val(X), \forall y \in Val(Y)$

# Conditional independence

- **Sets** of variables **X**, **Y**, **Z**

- **X** is independent of **Y** given **Z** if
  - Shorthand: $P \vdash (\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$
  - For $P \vdash (\mathbf{X} \perp \mathbf{Y} \mid \quad)$, write $P \vdash (\mathbf{X} \perp \mathbf{Y})$

- **Proposition:** $P$ satisfies $(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$ if and only if
  - P(**X**,**Y**|**Z**) = P(**X**|**Z**) P(**Y**|**Z**), $\quad \forall x \in Val(X), \forall y \in Val(Y), \forall z \in Val(Z)$

# Concept

- Bayes Rules
  - Simple yet fundamental

$$P(B|A) = \frac{P(A \wedge B)}{P(A)} = \frac{P(A|B)\,P(B)}{P(A)}$$

This is Bayes Rule

**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London,* **53:370-418**

# Bayes Rule

- **Simple yet profound**
  - Using Bayes Rules doesn't make your analysis Bayesian!

- **Concepts:**
  - Likelihood
    - How much does a certain hypothesis explain the data?
  - Prior
    - What do you believe before seeing any data?
  - Posterior
    - What do we believe after seeing the data?

# Entropy

- Measures the amount of ambiguity or uncertainty in a distribution:

$$H(p) = -\sum_x p(x) \log p(x)$$

- Expected value of $-\log p(x)$ (a function which depends on p(x)!).
- $H(p) > 0$ unless only one possible outcomein which case $H(p) = 0$.
- Maximal value when p is uniform.
- Tells you the expected "cost" if each event costs $-\log p$(event)
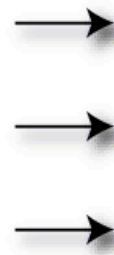
# KL-Divergence / Relative Entropy

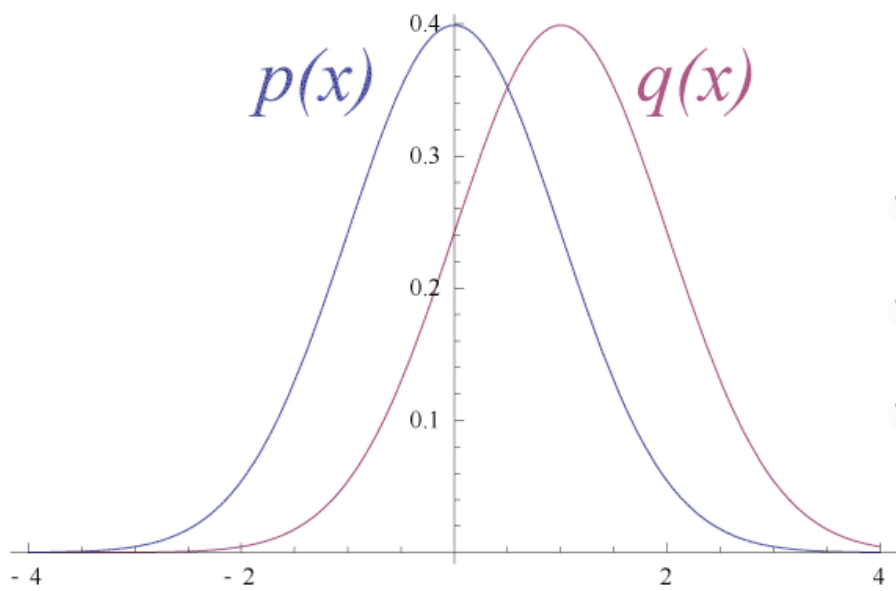- An assymetric measure of the distancebetween two distributions:

$$KL[p\|q] = \sum_x p(x)[\log p(x) - \log q(x)]$$

- $KL > 0$ unless $p = q$ then $KL = 0$

- Tells you the extra cost if events were generated by $p(x)$ but instead of charging under $p(x)$ you charged under $q(x)$.
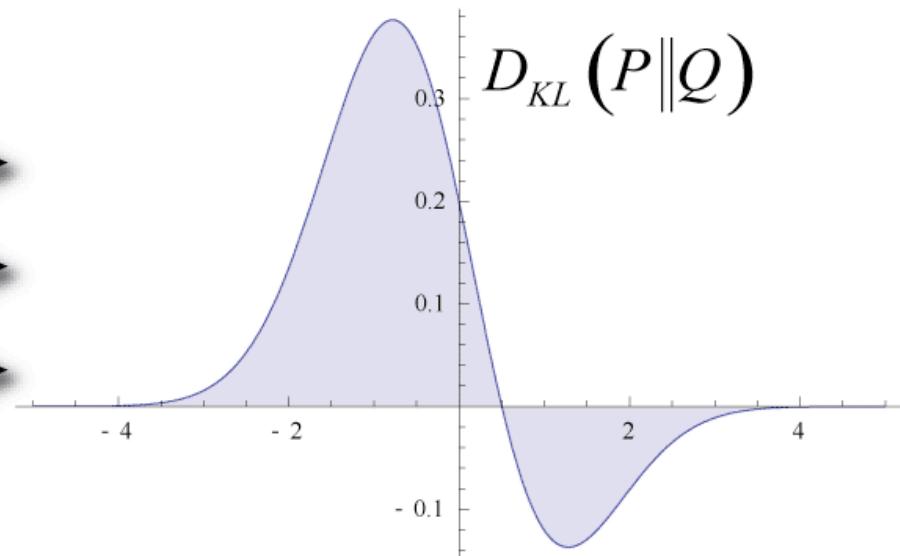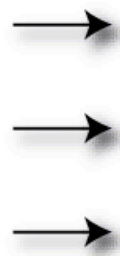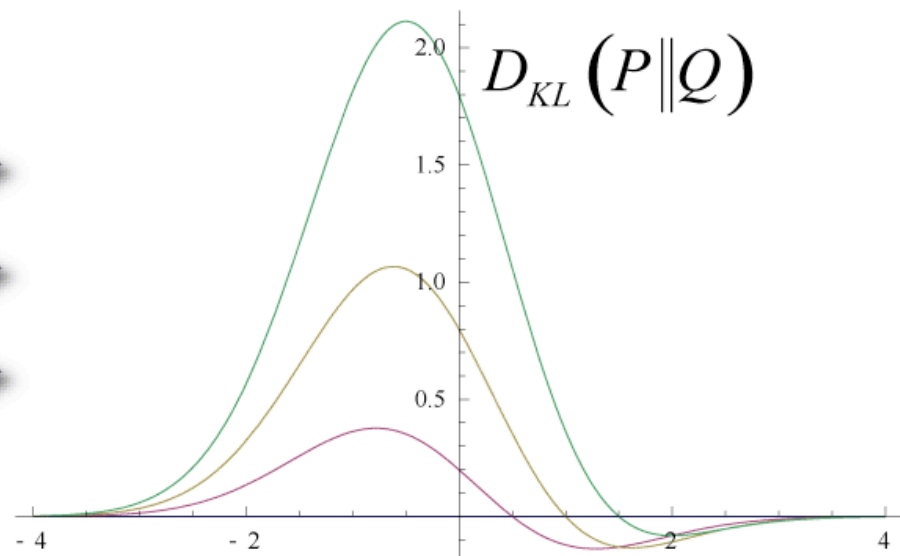
$p(x)$     $q(x)$

Original Gaussian PDF's

$D_{KL}\left(P\|Q\right)$

KL Area to be Integrated

$p(x)$  $q(x)$

Original Gaussian PDF's

$D_{KL}(P\|Q)$
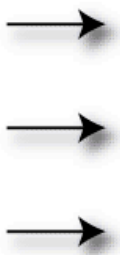
KL Area to be Integrated

$D_{KL}(P\|Q)$

- End of Prob. Review

- Start of Estimation