

# **Mathematics of Big Data, I**

## **Lecture 2: Effective Optimization and Computation, Logistic Regression, and Generalized Linear Models**

**Weiqing Gu**  
Professor of Mathematics  
Director of the Mathematics Clinic

Harvey Mudd College  
Summer 2018

**<https://math189su18.github.io/>**

# Recall last time we covered following

- First: Big data introduction (answer first two questions)
  - Big Data Introduction
    - *Where does big data come from?*
    - *Different ways to describe big data*
- Second: Use linear regression as an example to give an overview of big data analytics

## Modeling Approaches:

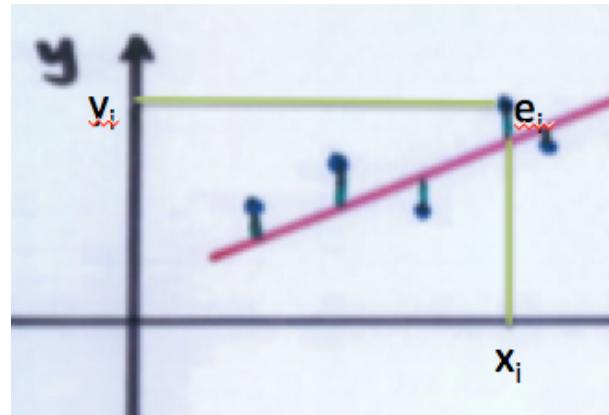
- *Statistical calculus*
- *Geometric analytic*
- *Probabilistic*

Each has its own merit

# Let's Recap

We had shown the following all three approaches give the same solution.

- Statistical calculus
- Geometric analytic
- Probabilistic



$$\begin{aligned}y &= \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_m x_m \\&\quad \uparrow \theta_0 = 1 \\&= (\theta_0, \theta_1, \dots, \theta_m) \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{pmatrix} \\&= \theta^T x \triangleq h_{\theta}(x)\end{aligned}$$

$$\theta = (X^T X)^{-1} X^T \vec{y}.$$

$$\text{Let } A = X^T X$$

Q1: What if A is not invertible? Want to achieve some perturbing of A. This is equivalent (your hw) to minimize:  $\|Ax - b\|_2^2 + \|\Gamma x\|_2^2$ . (This is Called ridge regression.)

Q2: For big data, is it really effective to compute  
? NO!  $(X^T X)^{-1}$

# Today' Topics

- Review or intro Probability Theory
- Logistic Regression
- Generalized Linear Models
- Effective Optimization and Computation  
*(only if time permits)*

# When we deal with big data, we must study **Effective optimization Techniques and Fast Computation**

- For this course, we will focus on
  - Gradient Descent
    - **Batch gradient descent**
    - **Stochastic gradient descent**
  - Newton's method
  - Various matrix decompositions, for examples
    - LU decomposition
    - Cholesky decomposition

For example: **We can use LU or Cholesky decomposition to solve the normal eqn:**

$$X^T X \theta = X^T \vec{y}$$

Recall: For linear regression, we want  
want to choose  $\theta$  to minimize  $J(\theta)$ .

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = \frac{1}{2}(X\theta - \vec{y})^T(X\theta - \vec{y})$$

**Key:  $J$  is quadratic on  $\theta$ ; Exists Unique Minimum!**

$$h_\theta(x^{(i)}) = (x^{(i)})^T \theta$$

**Note:  $h$  is linear on  $\theta$ !**

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$X\theta - \vec{y} = \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(m)})^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

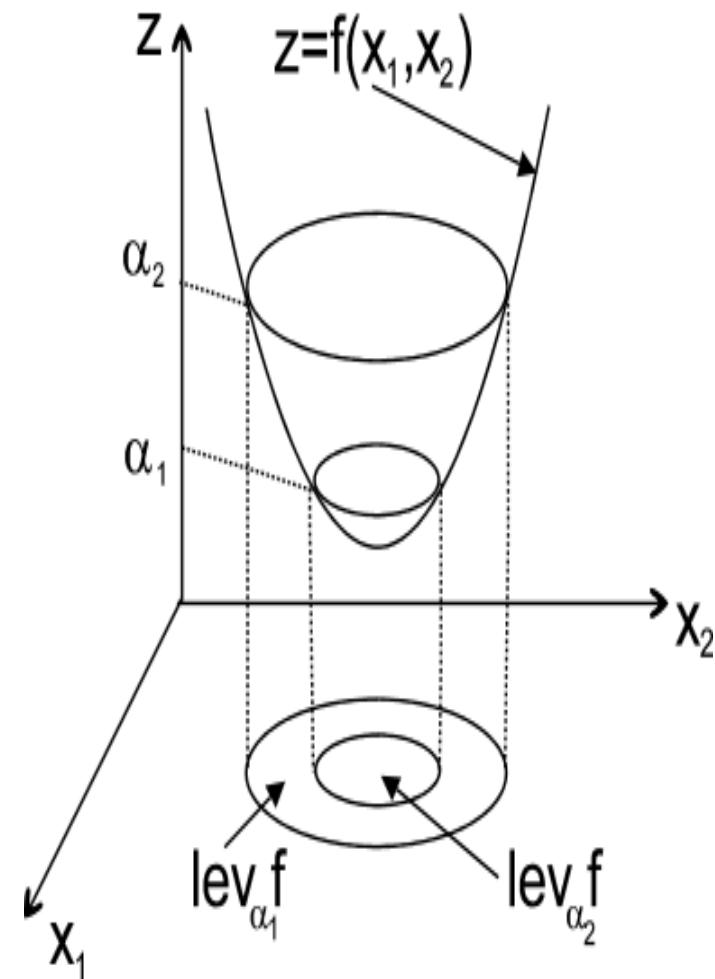
$$= \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}.$$

$$X = \begin{bmatrix} \quad (x^{(1)})^T \quad \\ \quad (x^{(2)})^T \quad \\ \vdots \\ \quad (x^{(m)})^T \quad \end{bmatrix}$$

# Why does $J(\theta)$ have a unique minimum?

(Exercise: Use two different ways to prove it—Hints below.)

- Since  $X^T X$  is positive definite when  $X^T X$  is invertible.
- (Hint for proof:  $v^T (X^T X) v = (Xv)^T (Xv) = \|Xv\|^2 \geq 0$  and the equality holds if and only if  $v$  is the zero vector. Use the rank of  $X$ .)
- In one variable,  $f(x) = ax^2 + bx + c$ , if  $a > 0$ , how does the graph of  $f$  look like?
- Another way: use geometric approach to get the normal equation and write down the unique solution.



# (Least Mean Square) LMS Algorithm

Q: Given a training set, how do we pick/learn, the parameters  $\theta$ ?

A: Find the gradient of  $J(\theta)$ .

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2.$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\ &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h_\theta(x) - y) x_j\end{aligned}$$

Note: Here it real should be the transpose of it times itself. But when you take derivative, you think it is a square.

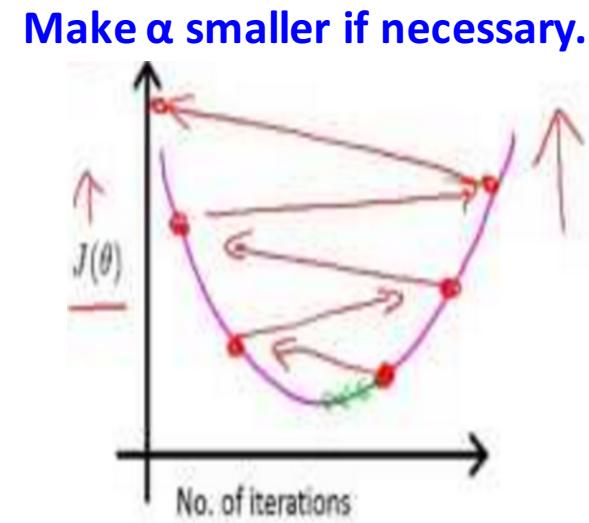
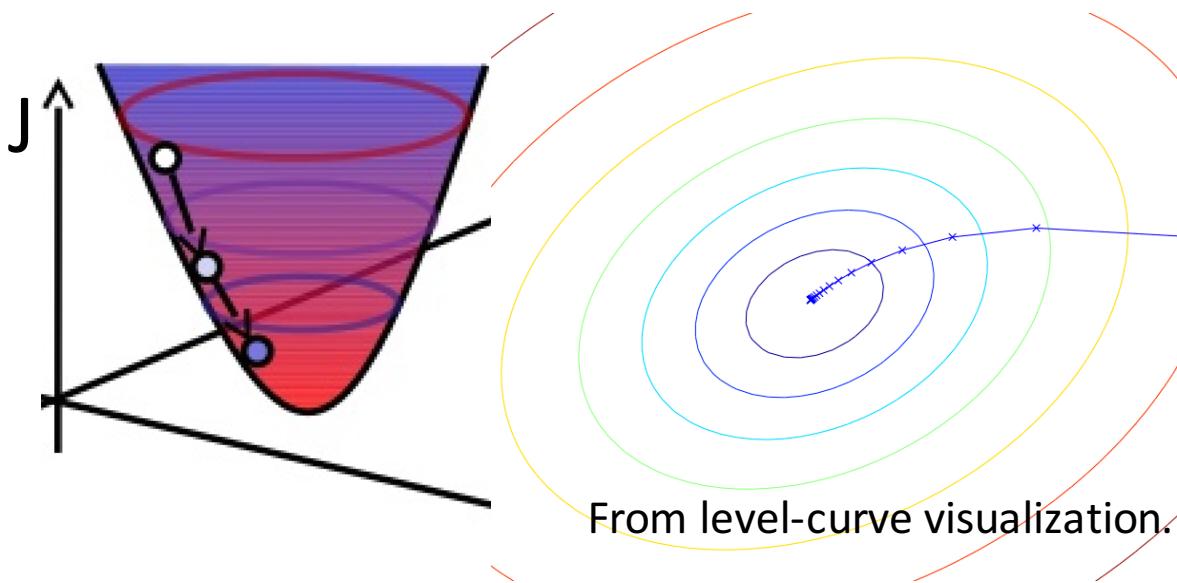


$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}.$$

This rule is called the LMS update rule (or Widrow-Hoff learning rule).

# Use the gradient descent algorithm

- Which starts with some initial  $\theta$ , and repeatedly performs the update.
- Here  $\alpha$  is called the learning rate.
- Geometrically, it repeatedly takes a step in the direction of steepest decrease of  $J$ .



# Batch Gradient Descent (BGD)

Repeat until convergence {

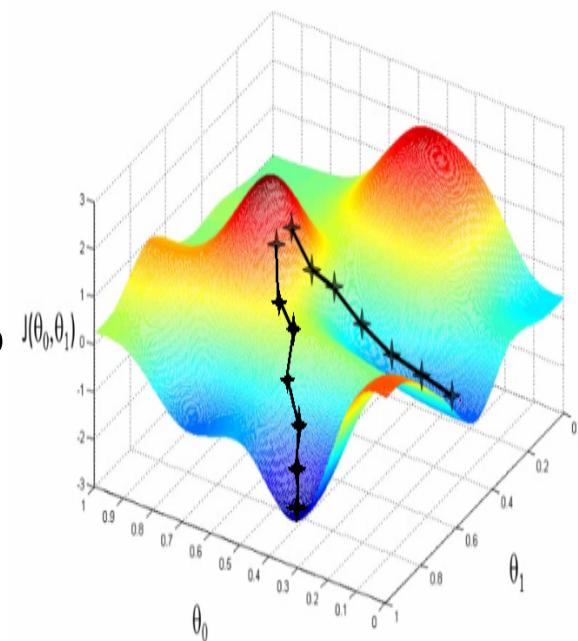
$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} \quad (\text{for every } j).$$

$\underbrace{\hspace{10em}}$   
 $- \partial J(\theta) / \partial \theta_j$

This is simply gradient descent on the original cost function  $J$ .

Remarks:

- 1) **This method looks at every example in the entire training set on every step**, and is called BGD.
- 2) It is well known that gradient descent can be susceptible to local minima in general (see the figure on right), **the optimization problem we have** posed here for linear regression **has only one global**, and no other local, **optima**; thus gradient descent always converges (assuming the learning rate  $\alpha$  is not too large) to the global minimum.
- 3) **The key is that our  $J$  is a convex quadratic function.**



# Stochastic Gradient Descent (SGD)

Loop {

    for i=1 to m, {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j).$$

    }

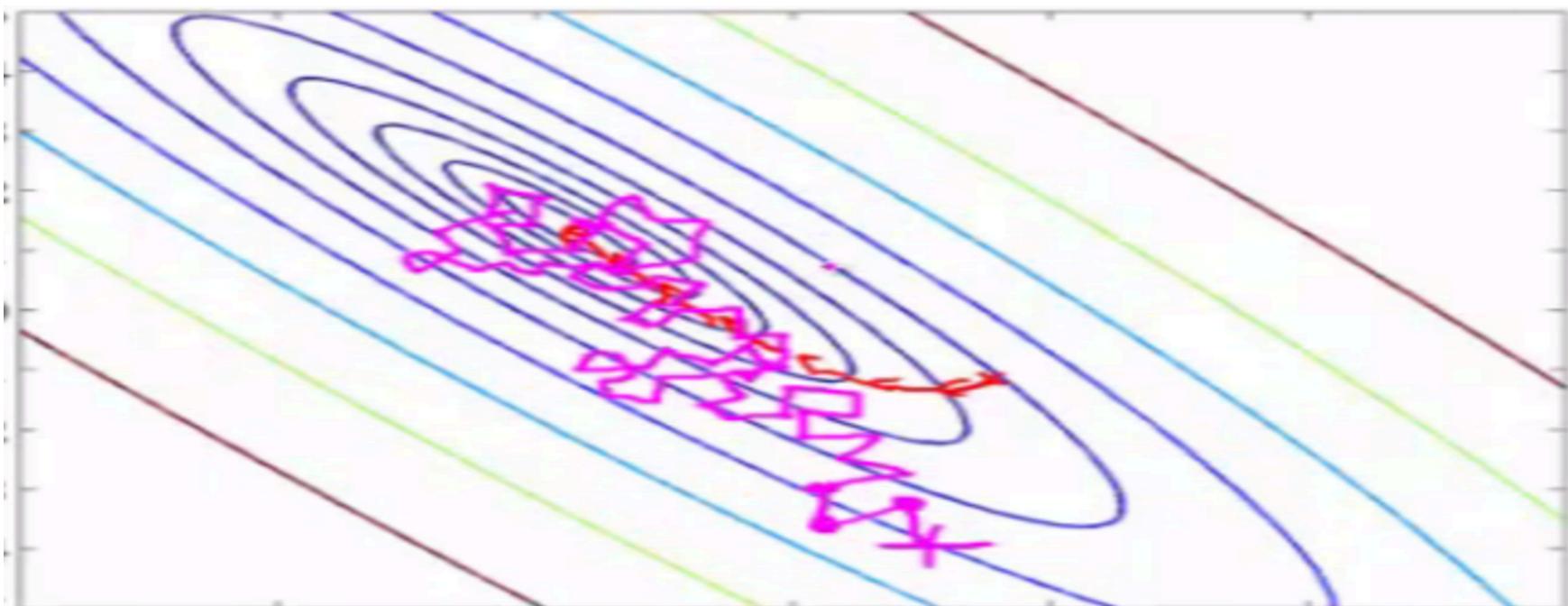
}

Remarks:

- 1) SGD **repeatedly run through the training set, and each time it encounters a training example, it updates the parameters** according to the gradient of the error with respect to that single training example only.
- 2) SGD **may never “converge” to the unique minimum**, and the parameters  $\theta$  will keep oscillating around the minimum of  $J(\theta)$ ; but **in practice** most of the values near the minimum will be **reasonably good approximations** to the true minimum.

# Comparing Batch gradient descent with Stochastic gradient descent

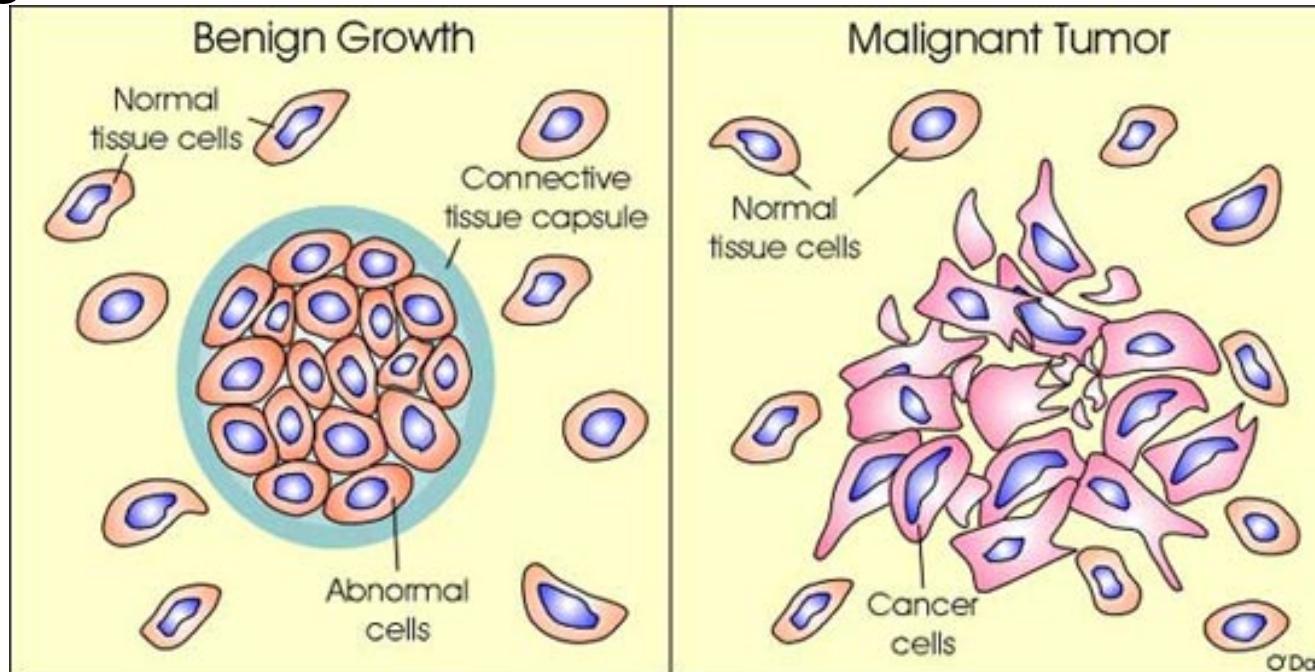
- *For big data*, often the training set is large, *people prefer use stochastic gradient descent* instead of batch gradient descent.
- Since *BGD has to scan thru the entire training set before taking a single step*—a costly operation if  $m$  is large—*SGD can start making progress right away*, & continues to make progress with each example it looks at.
- *SGD can run on dynamical data sets*. As data coming, it updates the parameters.
- Often, **SGD gets  $\theta$  “close” to the minimum much faster than BGD**.
- But SGD gets only approximation solution of  $\theta$ . This is a **trade off** when dealing with big data.



Now we switch gear:

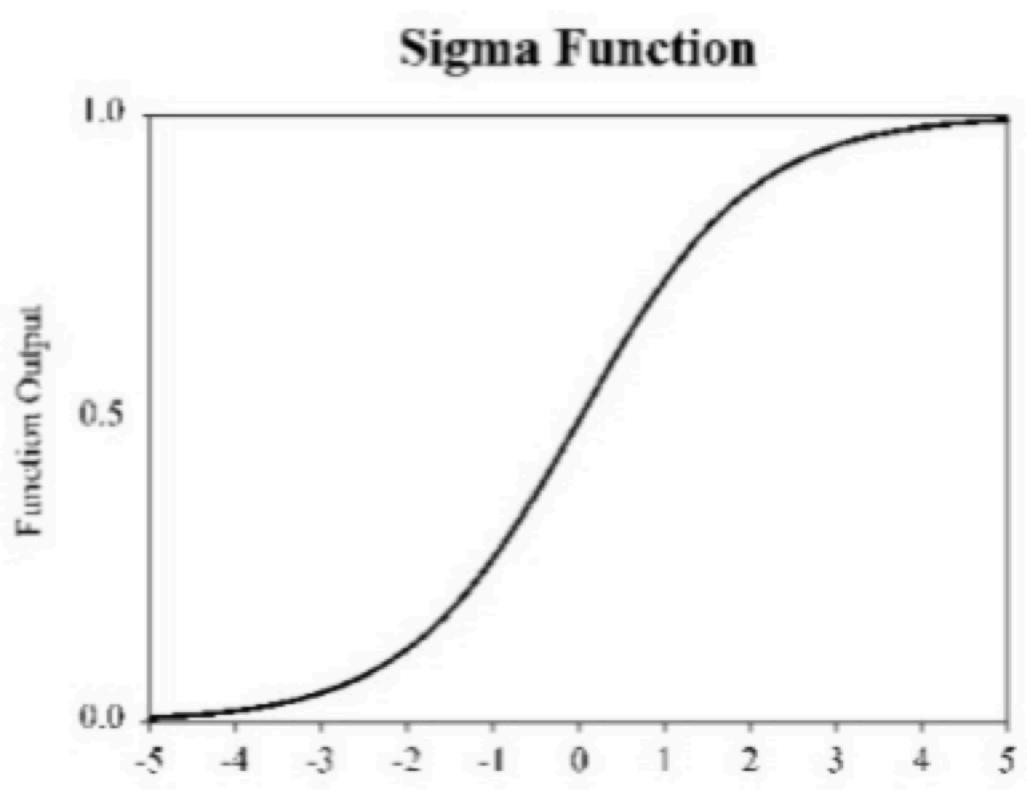
# Logistic Regression and Newton's Method

- Key: Logistic Regression is for Classification Problems.
- For example: distinguish between benign tumors and malignant tumors.



**Key idea:** try to utilize the linear regression techniques by transform a discrete problem to a smooth problem passing thru a sigma so that we can take gradient for optimization.

Logistic Regression maps the fitting straight line/hyperplane in linear regression to a monotone increasing curve, often a ***sigmoid function***.



*Work out the details  
of  
Logistic regression  
with  
the students on the  
board.*

Your homework: Problem 7

(a) Let  $\sigma(x) = \frac{1}{1+e^{-x}}$  be the sigmoid function. Show that

$$\sigma'(x) = \sigma(x) [1 - \sigma(x)].$$

# It is easier to maximize the log likelihood:

$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)})) \\ \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= (y - h_\theta(x)) x_j\end{aligned}$$

This gives us **the stochastic gradient ascent rule**:

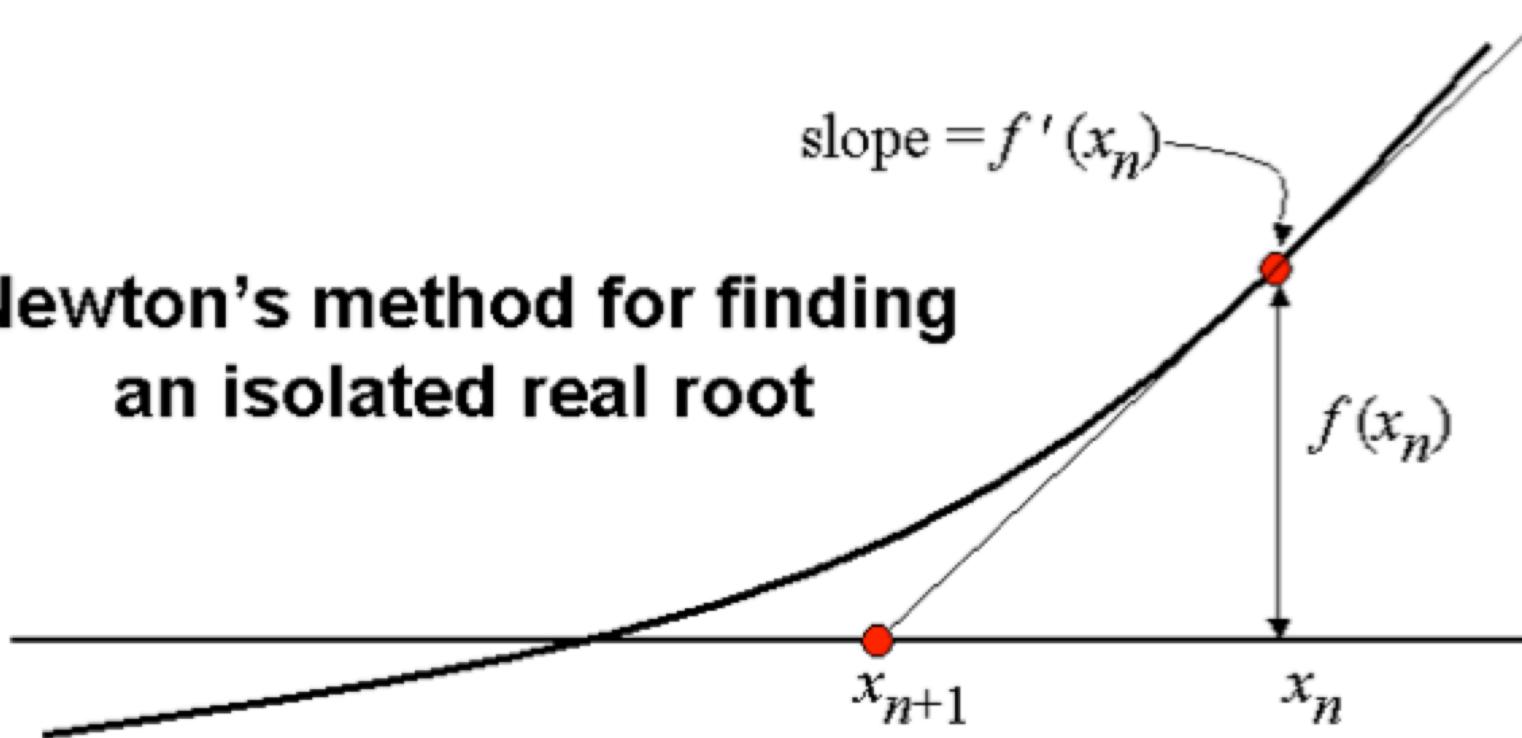
$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

*Note : There is another method running even fast than this one, called **Newton's method**.*

# Newton's method for fast computation

In the case of line, we just use the definition of the slope of  $f$ .

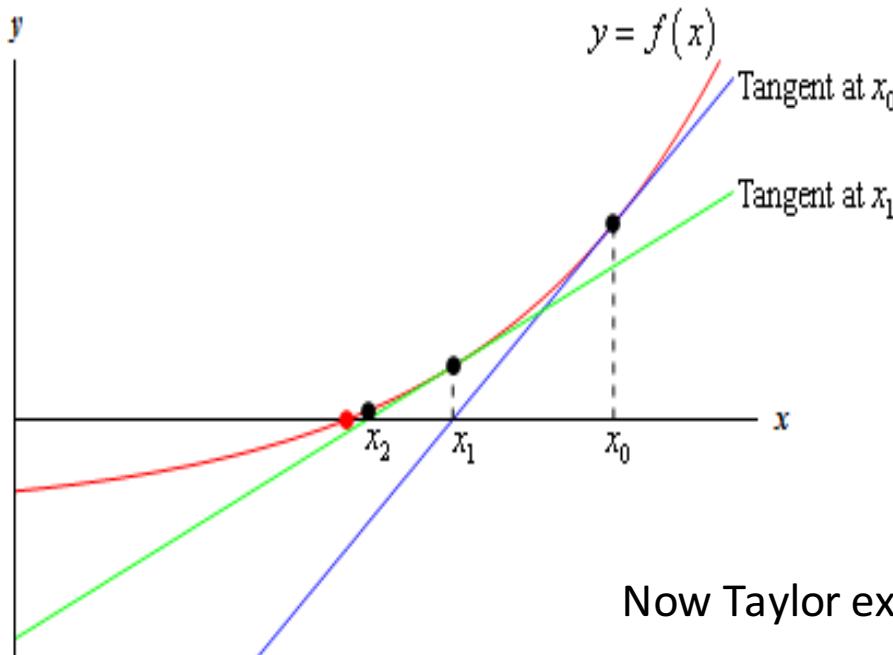
**Newton's method for finding  
an isolated real root**



$$x_{n+1} = x_n - \frac{f'(x_n)}{f(x_n)}$$

# Newton's method for fast computation

- Case 1 Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  (here, we just use the definition of the slope of  $f$ .)
- Newton's method for finding an isolated real root
- **Key:** In general using Taylor expansion at  $x = x_0$
- Take the linear best approximation & plug in  $x = x_1$ .



Note: if  $x_1$  is a root, then  $f(x_1) = 0$ , expand at  $x_0$  and plug  $x = x_1$ .

$$0 = f(x_0) + f'(x_0)(x_1 - x_0)$$

$$x_1 - x_0 = -\frac{f(x_0)}{f'(x_0)}$$

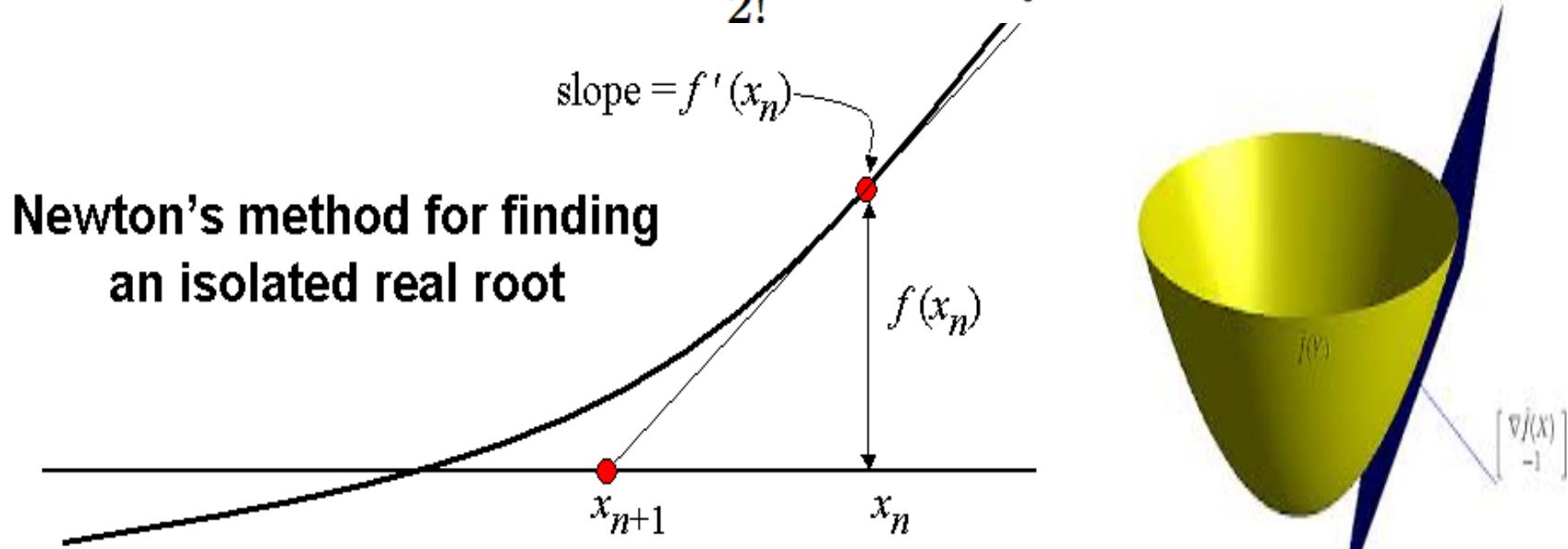
$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

Now Taylor expanding of  $f$  at  $x = x_1$  and similarly finding  $x_2$ .

*Iteratively:* Taylor expanding of  $f$  at  $x = x_n$ , plugging in  $x = x_{n+1}$ , and solving  $x_{n+1}$ .

Case 2 Multivariable:

$$f(\vec{x}) = f(\vec{a}) + (\vec{x} - \vec{a})^T \nabla f(\vec{a}) + \frac{1}{2!} (\vec{x} - \vec{a})^T H_f(\vec{a})(\vec{x} - \vec{a}) + \dots$$



$$x_{n+1} = x_n - \frac{f'(x_n)}{f(x_n)}$$

Newton's method is for finding a root of a function.

Keys: Taylor expansion, plug into linear part, solve, then iterate.

Now we switch gear again:

# Generalized Linear Models (GLMs)

- This topic includes: *exponential family* & *Softmax Regression*.
- *What is an exponential family?* A class of distributions is in the exponential family if

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

- $\eta$  = the natural parameter (or the canonical parameter) of the distribution
- $T(y)$  = the sufficient statistic ( often  $T(y) = y$ )
- $a(\eta)$  is the log partition function.

The quantity  $e^{-a(\eta)}$  essentially plays the role of a normalization constant, that makes sure the distribution  $p(y; \eta)$  sums/integrates over  $y$  to 1.

Let  $T$ ,  $a$  and  $b$  fixed and let the parameter  $\eta$  vary, then it defines a family of distribution. i.e. We get different distributions within this family.

Let's first show

**Bernoulli distributions are exponential family distribution.**

- Work out details with the students on the board.

Let's first show

## Gaussian distributions are exponential family distribution.

$$\begin{aligned} p(y; \mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp\left(\mu y - \frac{1}{2}\mu^2\right) \end{aligned}$$

Compare:

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

We get:

$$\eta = \mu$$

$$T(y) = y$$

$$a(\eta) = \mu^2/2$$

$$= \eta^2/2$$

$$b(y) = (1/\sqrt{2\pi}) \exp(-y^2/2).$$

# Constructing GLMs

Note: you need to know which distribution models what kind of problems  
(Reading assignment)

- Suppose you want to build a model to estimate the number ( $y$ ) of customers arriving in your store in any given hour, based on certain features  $x$  such as store promotions, recent advertising, weather, day-of-week, etc.
- We know that the Poisson distribution usually gives a good model for numbers of visitors.
- Knowing this, how can we come up with a model for this problem?
- Fortunately, the Poisson is an exponential family distribution, so we can apply a Generalized Linear Model (GLM). *(Homework or exam problem?)*
- Lots of known distributions are exponential families.
- Here, we will describe a method for constructing GLM models for problems such as these.

# Assumptions for Generalized Linear Models

- In generally, consider a classification or regression problem where we would like to predict the value of some random variable  $y$  as a function of  $x$ .
- To derive a GLM for this problem, we will make the following three assumptions about the conditional distribution of  $y$  given  $x$  and about our model:
- **1.  $y | x; \theta \sim \text{Exponential Family}(\eta)$ .** I.e., given  $x$  and  $\theta$ , the distribution of  $y$  follows some exponential family distribution, with parameter  $\eta$ .
- **2.** Given  $x$ , our goal is to predict the expected value of  $T(y)$  given  $x$ . Since often  $T(y) = y$ , so this means we would like the prediction  **$h(x)$  output by our learned hypothesis  $h$  to satisfy  $h(x) = E[y|x]$** . (Note that this assumption is satisfied in the choices for  $h_\theta(x)$  for both logistic regression and linear regression. For instance, in logistic regression, we had
$$h_\theta(x) = p(y = 1|x; \theta) = 0 \cdot p(y = 0|x; \theta) + 1 \cdot p(y = 1|x; \theta) = E[y|x; \theta].$$
)
- **3. The natural parameter  $\eta$  and the inputs  $x$  are related linearly:  $\eta = \theta^T x$ .** (Or, if  $\eta$  is vector-valued, then  $\eta_i = \theta_i^T x$ .)

# Examples: Least square and Logistic regression are GLM family of models

$$\begin{aligned} h_{\theta}(x) &= E[y|x; \theta] \\ &= \mu \\ &= \eta \\ &= \theta^T x. \end{aligned}$$

$$\begin{aligned} h_{\theta}(x) &= E[y|x; \theta] \\ &= \phi \\ &= 1/(1 + e^{-\eta}) \\ &= 1/(1 + e^{-\theta^T x}) \end{aligned}$$

Given that  $y$  is binary-valued, it therefore seems natural to choose the Bernoulli family of distributions to model the conditional distribution of  $y$  given  $x$ . In our formulation of the Bernoulli distribution as an exponential family distribution, we had  $\phi = 1/(1 + e^{-\eta})$ . Furthermore, note that if  $y|x; \theta \sim \text{Bernoulli}(\phi)$ , then  $E[y|x; \theta] = \phi$ .

# Softmax Regression

- Let's look at another example of a GLM. Consider a classification problem in which the response variable  $y \in \{1, 2, \dots, k\}$ .
- For example, rather than classifying email into the two classes spam or not-spam—which would have been a binary classification problem—this time we want to classify it into four classes, such as spam, family-mail, friends-mail, and work-related mail. The response variable is still discrete, but can now take on more than two values. We will thus model it as distributed according to a multinomial distribution.

Let's Derive

## A GLM using

# Multinomial distributions as exponential family distribution.

- **What are Multinomial distributions?**
- **For example:** If a 6 sided die has
  - 3 faces painted red
  - 2 faces painted white
  - 1 faces painted blueAnd rolled 100 times.  
Find  $P(60 \text{ red, } 30 \text{ white, and } 10 \text{ blue})$ .

*Work out details with the students on the board.*

**Generally an experiment with  $m$  outcomes with respective probabilities  $p_1, p_2, \dots, p_m$  is performed  $n$  times independently.**

**Let  $x_i = \# \text{ of times outcome } i \text{ appears, } i=1,2,\dots,m$**

**Then  $P(x_1=k_1, x_2=k_2, \dots, x_m = k_m) = ?$**

- Work out details with the students on the board.

# Details of Softmax Regression

- Work out details with the students on the board.

# Back Up Slides

- Review Probability Theory for those have taken the probability course.
- Introduction to Probability Theory for those have not taken the probability course.

A probability function is a special function which must satisfy:

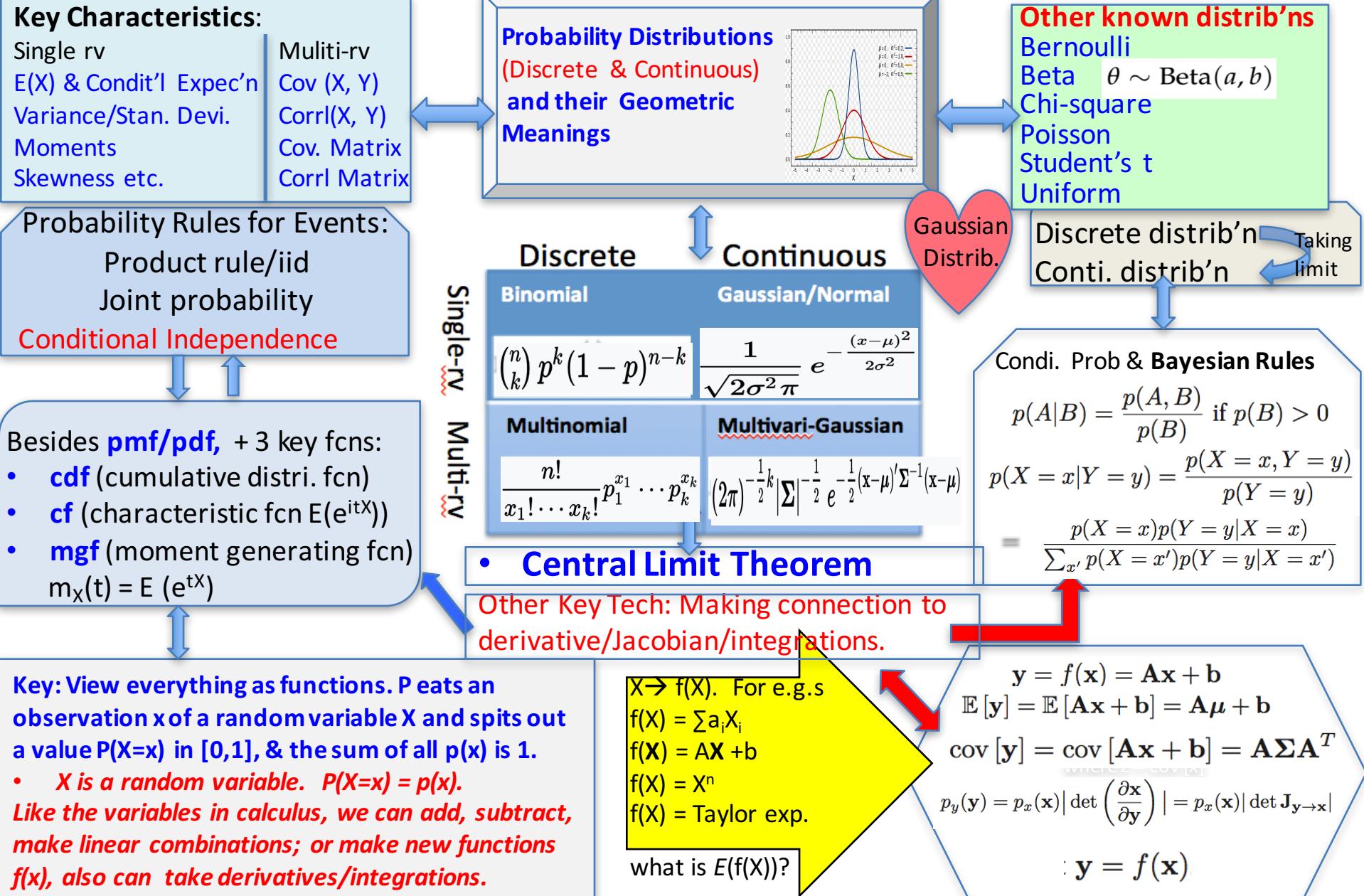
$$0 \leq P(X) \leq 1$$

$$\sum P(X) = 1$$

# A Big Picture of Probability Theory

$$0 \leq P(X) \leq 1$$

$$\sum P(X) = 1$$



# Bernoulli Distribution

If  $X$  is a random variable with this distribution, we have:

$$\Pr(X = 1) = p = 1 - \Pr(X = 0) = 1 - q.$$

The **probability mass function**  $f$  of this distribution, over possible outcomes  $k$ , is

$$f(k; p) = \begin{cases} p & \text{if } k = 1, \\ 1 - p & \text{if } k = 0. \end{cases}$$

This can also be expressed as

$$f(k; p) = p^k (1 - p)^{1-k} \quad \text{for } k \in \{0, 1\}$$

or as

$$f(k; p) = pk + (1 - p)(1 - k) \quad \text{for } k \in \{0, 1\}.$$

The Bernoulli distribution is a special case of the **binomial distribution** with  $n = 1$ .

# Binomial Distribution

- Probability Mass Function

In general, if the random variable  $X$  follows the binomial distribution with parameters  $n \in \mathbb{N}$  and  $p \in [0, 1]$ , we write  $X \sim B(n, p)$ . The probability of getting exactly  $k$  successes in  $n$  trials is given by the **probability mass function**:

$$Pr(k; n, p) = \Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

for  $k = 0, 1, 2, \dots, n$ , where

$$\binom{n}{k} = \frac{n!}{k!(n - k)!}$$

# **Two different ways to generalize Binomial distribution**

- From Binomial distribution to Poisson distribution
- From Binomial distribution to Multinomial Distribution

# Poisson's distribution

An event can occur 0, 1, 2, ... times in an interval. The average number of events in an interval is designated  $\lambda$  (lambda). Lambda is the event rate, also called the rate parameter. The probability of observing  $k$  events in an interval is given by the equation

$$P(k \text{ events in interval}) = e^{-\lambda} \frac{\lambda^k}{k!}$$

- Relations between Binomial Distribution and Poisson's distribution.
- Details on board.
- This gives an example of moving from discrete probability to continuous probability.

# Multivariate Distribution

- Details on board.
- **Example: Multinomial Distribution**

- **Recall:** What are Multinomial distributions?
- **For example:** If a 6 sided die has
  - 3 faces painted red
  - 2 faces painted white
  - 1 faces painted blueAnd rolled 100 times.  
Find  $P(60 \text{ red}, 30 \text{ white}, \text{ and } 10 \text{ blue})$ .

*Work out details with the students on the board.*

***Generally an experiment with  $m$  outcomes with respective probabilities  $p_1, p_2, \dots, p_m$  is performed  $n$  times independently.***

***Let  $x_i = \# \text{ of times outcome } i \text{ appears, } i=1,2,\dots,m$***

***Then  $P(x_1=k_1, x_2=k_2, \dots, x_m=k_m) = ?$***

**Claim: Multinomial distributions as exponential family distribution.**

**Claim: Multinomial distributions are exponential family distributions.**

- *Work out details with the students on the board.*

# correlation coefficient & correlation matrix

- The (Pearson) **correlation coefficient** between two rvs  $X$  and  $Y$  is defined as

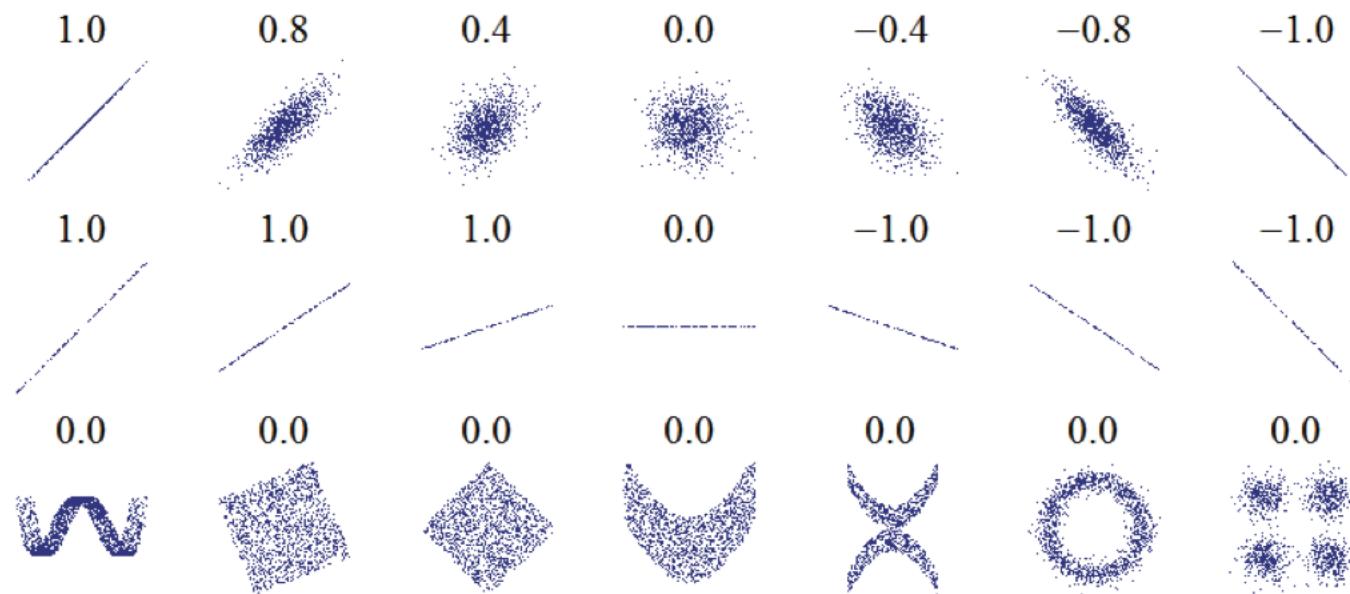
$$\text{corr} [X, Y] \triangleq \frac{\text{cov} [X, Y]}{\sqrt{\text{var} [X] \text{var} [Y]}}$$

- If  $X$  and  $Y$  are indep., then  $\text{cov} [X, Y] = 0$ ; say  $X$  and  $Y$  are uncorrelated.
- A **correlation matrix** of a random vector has the form:

$$\mathbf{R} = \begin{pmatrix} \text{corr} [X_1, X_1] & \text{corr} [X_1, X_2] & \cdots & \text{corr} [X_1, X_d] \\ \vdots & \vdots & \ddots & \vdots \\ \text{corr} [X_d, X_1] & \text{corr} [X_d, X_2] & \cdots & \text{corr} [X_d, X_d] \end{pmatrix}$$

Exercise: show that  $-1 \leq \text{corr} [X, Y] \leq 1$  and  
Show that  $\text{corr}[X, Y] = 1$  iff  $Y = aX + b$  for some parameters  $a$  and  $b$ .

# Example of Correlation Coefficients



**Figure 2.12** Several sets of  $(x, y)$  points, with the correlation coefficient of  $x$  and  $y$  for each set. Note that the correlation reflects the noisiness and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (bottom). N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of  $Y$  is zero. Source: [http://en.wikipedia.org/wiki/File:Correlation\\_examples.png](http://en.wikipedia.org/wiki/File:Correlation_examples.png)

# Conditional Probability

The **conditional probability** of event A,  
given that event B is true:

$$p(A|B) = \frac{p(A, B)}{p(B)} \text{ if } p(B) > 0$$

**Bayes rule:**

$$p(X = x|Y = y) = \frac{p(X = x, Y = y)}{p(Y = y)} = \frac{p(X = x)p(Y = y|X = x)}{\sum_{x'} p(X = x')p(Y = y|X = x')}$$

# Recall: Probability of an Event

- $p(A)$  denotes the probability that the event A is true.
- For example:
- A = a logical expression “it will rain tomorrow”  
We require that  $0 \leq p(A) \leq 1$ .

$p(A) = 0$  means the event definitely will not happen

$p(A) = 1$  means the event definitely will happen

$p(\bar{A})$  denotes the probability of the event not A

$$p(\bar{A}) = 1 - p(A)$$

We also write:

A=1 to mean the event A is true.

A=0 to mean the event A is false.

# Recall: Fundamental Rules

$$\begin{aligned} p(A \vee B) &= p(A) + p(B) - p(A \wedge B) \\ &= p(A) + p(B) \text{ if } A \text{ and } B \text{ are mutually exclusive} \end{aligned}$$

$$p(A, B) = p(A \wedge B) = p(A|B)p(B)$$

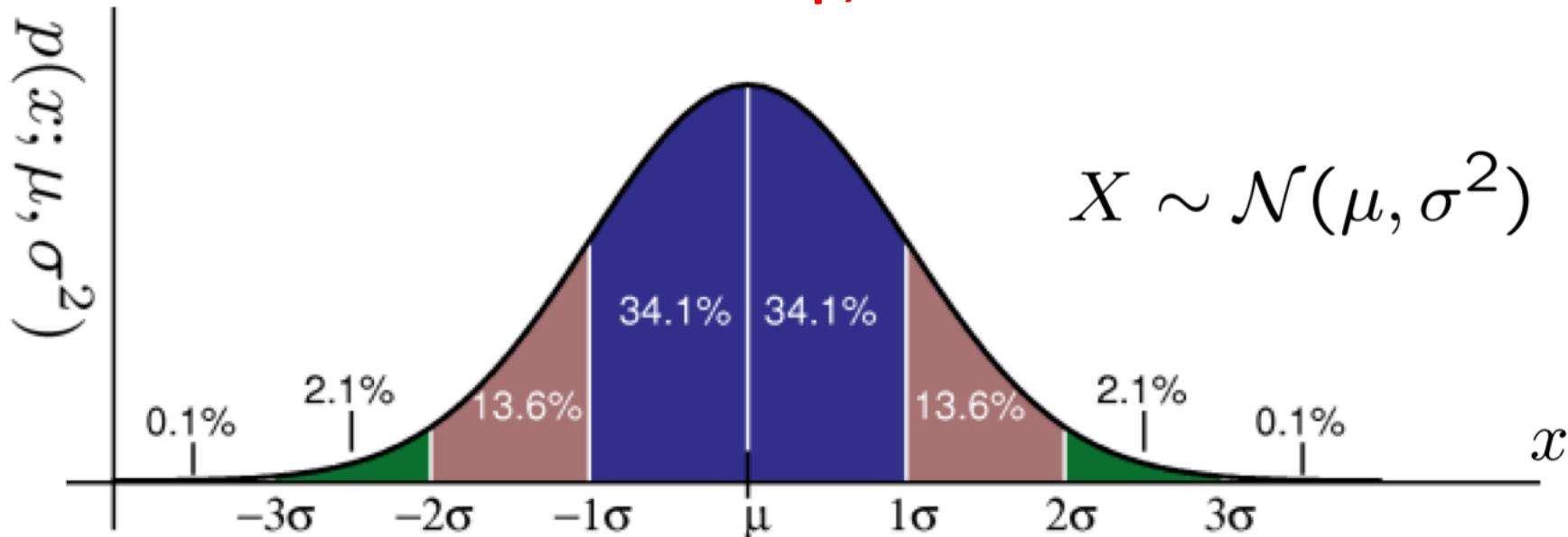
$$p(A) = \sum_b p(A, B) = \sum_b p(A|B = b)p(B = b)$$

$$p(X_{1:D}) = p(X_1)p(X_2|X_1)p(X_3|X_2, X_1)p(X_4|X_1, X_2, X_3)\dots p(X_D|X_{1:D-1})$$

Changing gear:

# Recall: Gaussian with one variable (called *Univariate Gaussian*)

Gaussian distribution with mean  $\mu$ , and standard deviation  $\sigma$ .



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

When  $\mu = 0$  and  $\sigma = 1$ , it is called the standard normal distribution.

# Different ways to find expected values

$$E[X] = \int_{-\infty}^{\infty} xf(x)dx$$

Where  $f(x)$  is the probability density function of  $X$ .

**Example:** Let  $f(x)$  be the density of the standard normal distribution.

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}}$$

$$E[X] = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} dx$$

**Method 1:** Since  $xe^{-x^2/2}$  is an odd function and the limits of the integral are symmetric, so we get  $E[X] = 0$ .

**Method 2:** Directly integrate.

**Method 3:** Using the moment generating function.

## Method 2

$$\begin{aligned} E[X] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} xe^{-\frac{x^2}{2}} dx \\ &= -\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} d(-\frac{x^2}{2}) \\ &= -\frac{1}{\sqrt{2\pi}} e^{-x^2/2} \Big|_{-\infty}^{\infty} \\ &= 0 \end{aligned}$$

# Method 3

- The moment generating function is defined as

$$\phi(t) = E[e^{tX}].$$

$$\phi(t) = C \int_{\mathbb{R}} e^{tx} e^{-x^2/2} dx = C \int_{\mathbb{R}} e^{-x^2/2+tx} dx = e^{t^2/2} C \int_{\mathbb{R}} e^{-(x-t)^2/2} dx.$$

$$t^2/2 - (x - t)^2/2 = t^2/2 + (-x^2/2 + tx - t^2/2) = -x^2/2 + tx$$

1

$$\phi(t) = e^{t^2/2} = 1 + (t^2/2) + \frac{1}{2}(t^2/2)^2 + \dots + \frac{1}{k!}(t^2/2)^k + \dots$$

$$E[e^{tX}] = E \left[ 1 + tX + \frac{1}{2}(tX)^2 + \dots + \frac{1}{n!}(tX)^n + \dots \right]$$

$$= 1 + E[X]t + \frac{1}{2}E[X^2]t^2 + \dots + \frac{1}{n!}E[X^n]t^n + \dots \rightarrow E[X] = 0$$

When  $k=1$ ,  
 $E[X^2]=1$ .  
Variance = 1.

Compare:

$$\frac{1}{(2k)!}E[X^{2k}]t^{2k} = \frac{1}{k!}(t^2/2)^k = \frac{1}{2^k k!}t^{2k},$$



$$E[X^{2k}] = \frac{(2k)!}{2^k k!}, k = 0, 1, 2, \dots$$

# Properties of Gaussians

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

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

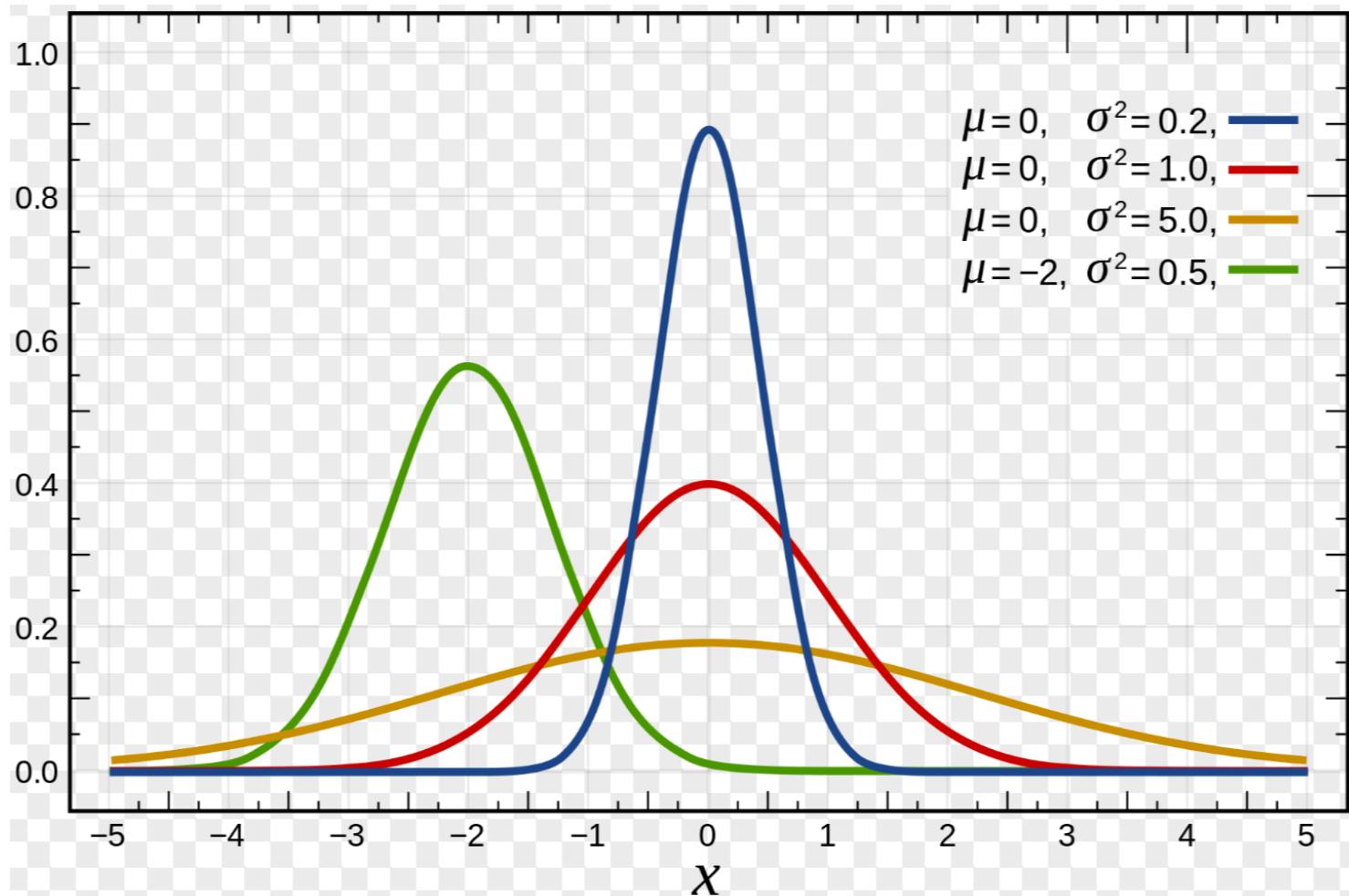
- Integration of the densities equals to 1.

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

- Mean:  $\begin{aligned} \mathbb{E}_X[X] &= \int_{-\infty}^{\infty} xp(x; \mu, \sigma^2) dx \\ &= \int_{-\infty}^{\infty} x \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx \\ &= \mu \end{aligned}$
- Variance:

$$\begin{aligned} \mathbb{E}_X[(X - \mu)^2] &= \int_{-\infty}^{\infty} (x - \mu)^2 p(x; \mu, \sigma^2) dx \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx \\ &= \sigma^2 \end{aligned}$$

In general, do translation and scale;  
i.e. change of variables when try to  
find those key characteristic values



# Covariance, and Covariance Matrix

- The **covariance** between two rv's X and Y measures the degree to which X and Y are (linearly) related; defined as

$$\text{cov}[X, Y] \triangleq \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

Exercise

$$= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

If  $\mathbf{x}$  is a d-dimensional random vector, its **covariance matrix** is defined to be the following symmetric, positive definite matrix:

Often denoted by  $\Sigma$

$$\text{cov}[\mathbf{x}] \triangleq \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T]$$
$$= \begin{pmatrix} \text{var}[X_1] & \text{cov}[X_1, X_2] & \cdots & \text{cov}[X_1, X_d] \\ \text{cov}[X_2, X_1] & \text{var}[X_2] & \cdots & \text{cov}[X_2, X_d] \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}[X_d, X_1] & \text{cov}[X_d, X_2] & \cdots & \text{var}[X_d] \end{pmatrix}$$

# correlation coefficient & correlation matrix

- The (Pearson) **correlation coefficient** between two rvs  $X$  and  $Y$  is defined as

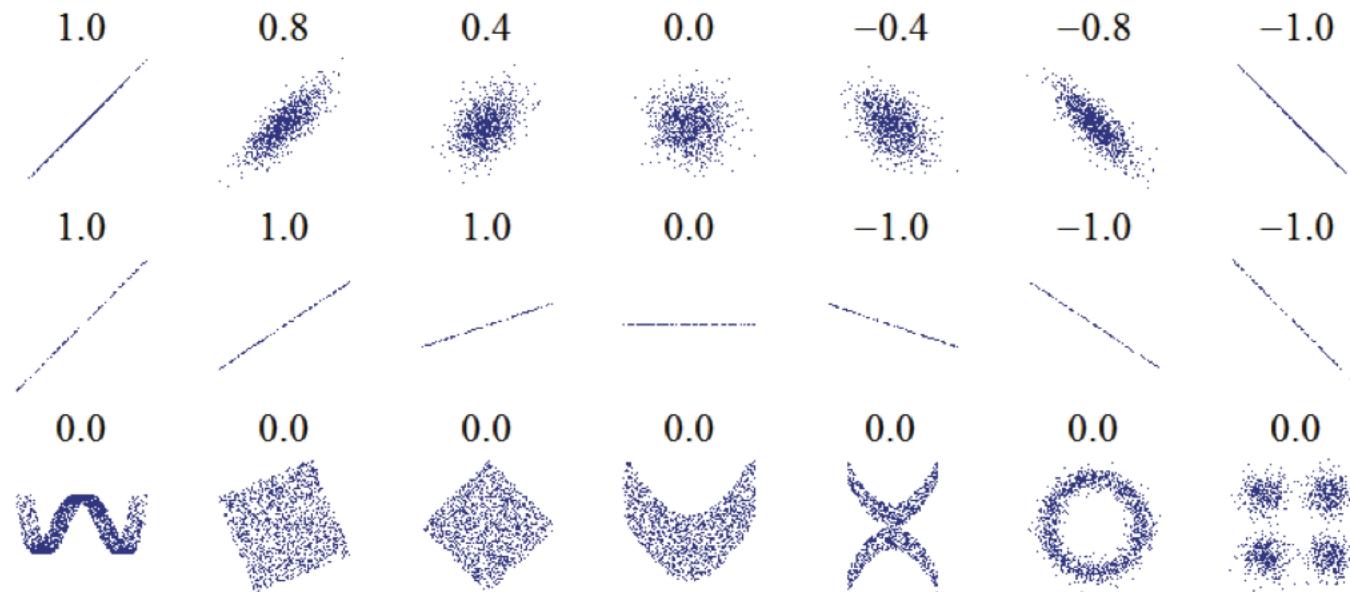
$$\text{corr} [X, Y] \triangleq \frac{\text{cov} [X, Y]}{\sqrt{\text{var} [X] \text{var} [Y]}}$$

- If  $X$  and  $Y$  are indep., then  $\text{cov} [X, Y] = 0$ ; say  $X$  and  $Y$  are uncorrelated.
- A **correlation matrix** of a random vector has the form:

$$\mathbf{R} = \begin{pmatrix} \text{corr} [X_1, X_1] & \text{corr} [X_1, X_2] & \cdots & \text{corr} [X_1, X_d] \\ \vdots & \vdots & \ddots & \vdots \\ \text{corr} [X_d, X_1] & \text{corr} [X_d, X_2] & \cdots & \text{corr} [X_d, X_d] \end{pmatrix}$$

Exercise: show that  $-1 \leq \text{corr} [X, Y] \leq 1$  and  
Show that  $\text{corr}[X, Y] = 1$  iff  $Y = aX + b$  for some parameters  $a$  and  $b$ .

# Example of Correlation Coefficients



**Figure 2.12** Several sets of  $(x, y)$  points, with the correlation coefficient of  $x$  and  $y$  for each set. Note that the correlation reflects the noisiness and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (bottom). N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of  $Y$  is zero. Source: [http://en.wikipedia.org/wiki/File:Correlation\\_examples.png](http://en.wikipedia.org/wiki/File:Correlation_examples.png)

# The multivariate Gaussian (distribution) or multivariate normal (MVN)

(The most widely used joint probability density function for continuous variables)

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

determinant

where  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] \in \mathbb{R}^D$  and  $\boldsymbol{\Sigma} = \text{cov}[\mathbf{x}]$

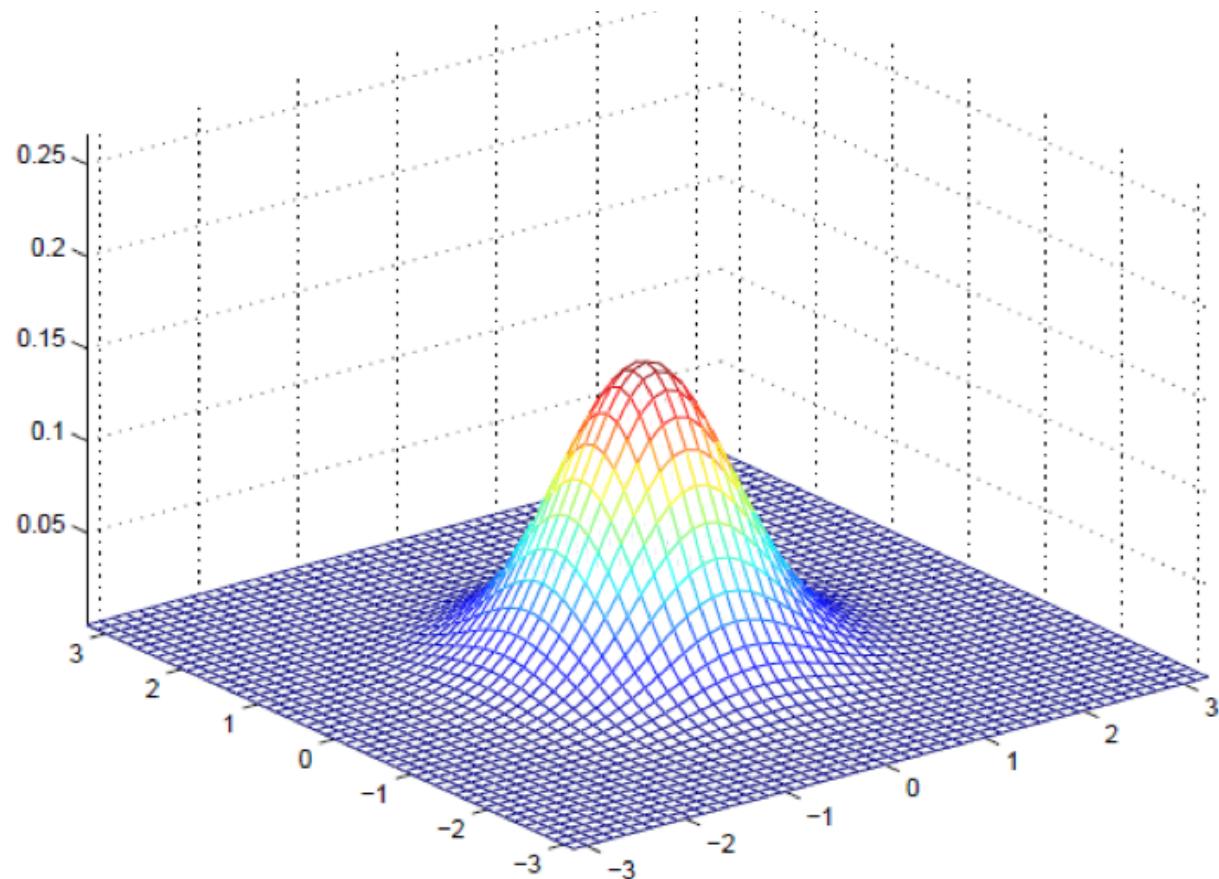
Note: the precision matrix or concentration matrix is just

the inverse covariance matrix,  $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$

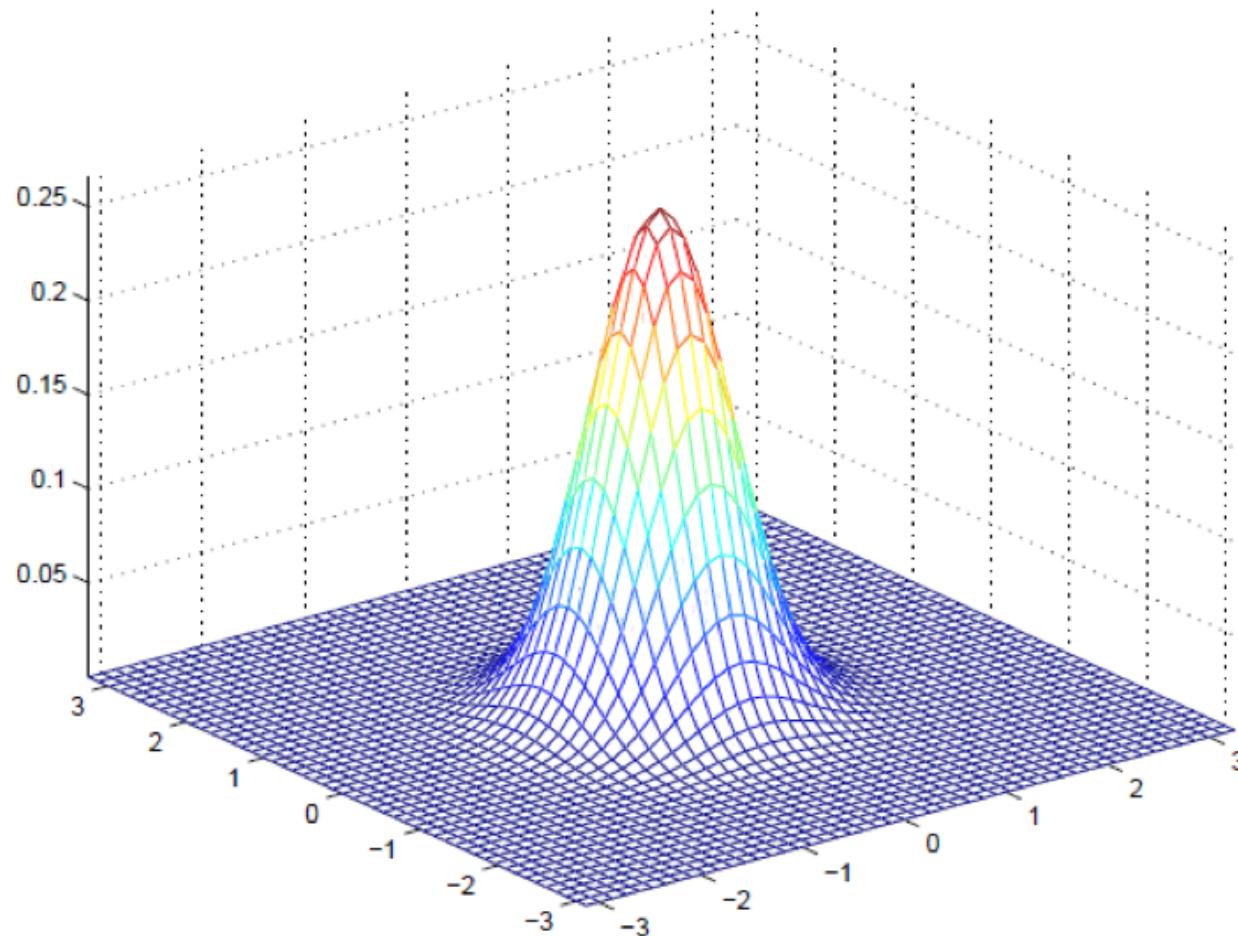
A spherical or isotropic covariance  $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_D$ ,  
has one free parameter.

$$\mu = [0; 0]$$

$$\Sigma = [I \ 0; 0 \ I]$$

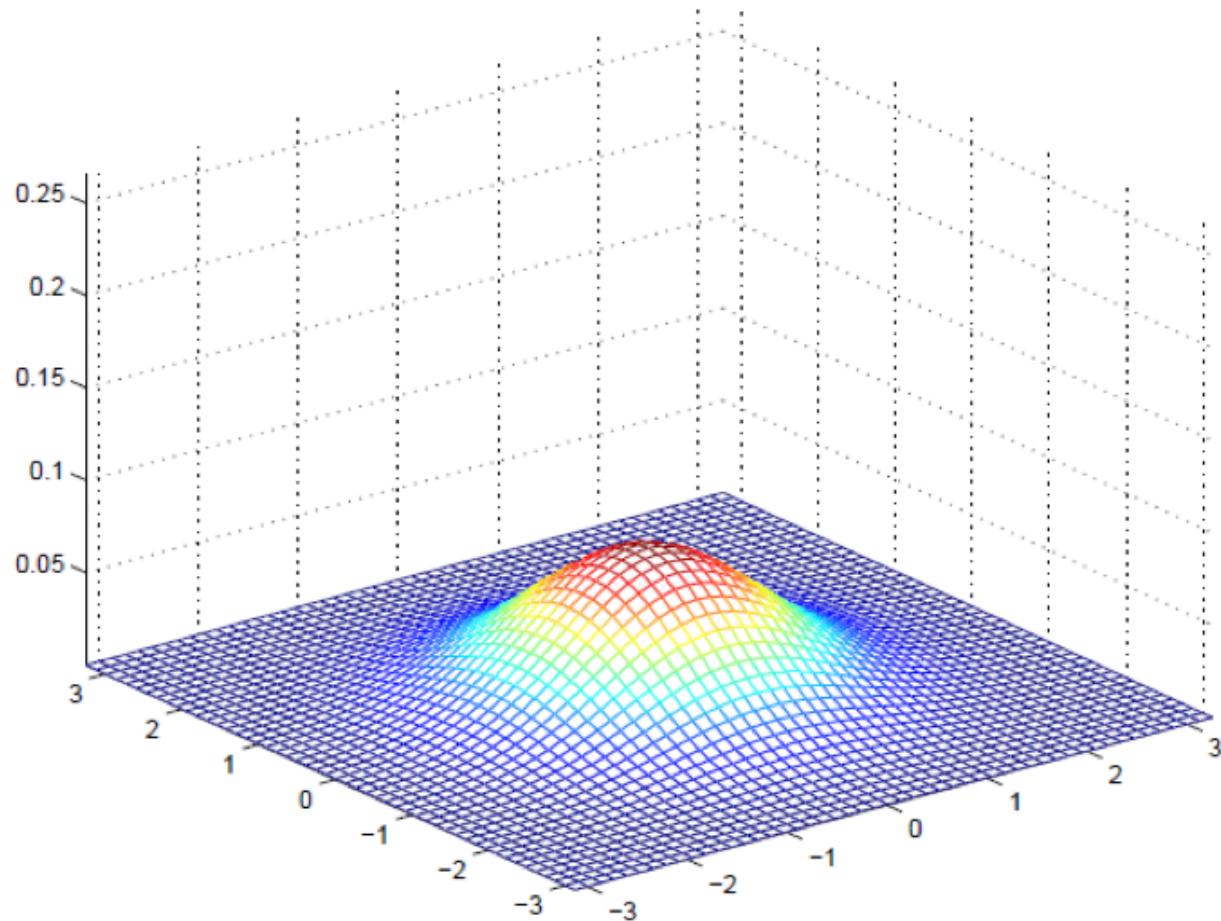


$$\mu = [0; 0]$$
$$\Sigma = [.6 \ 0 ; 0 .6]$$



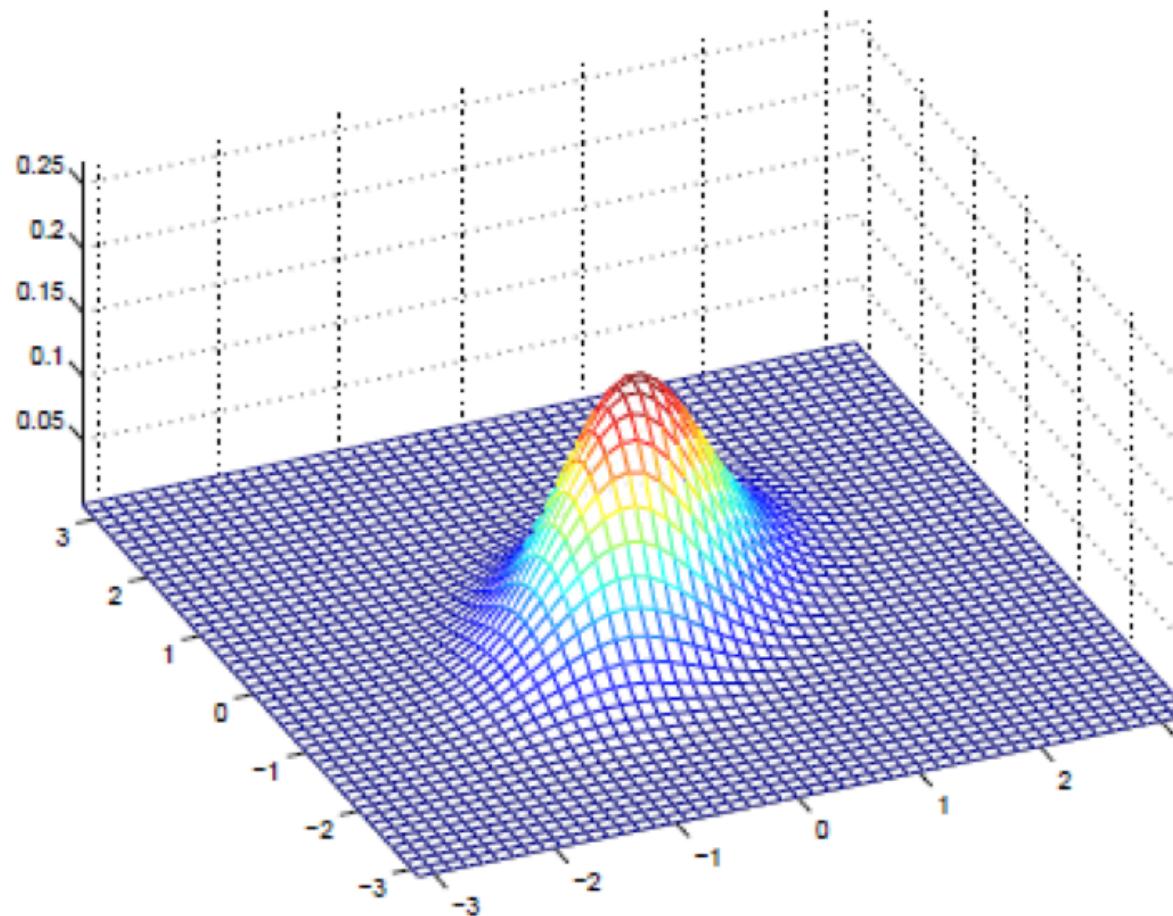
$$\mu = [0; 0]$$

$$\Sigma = [2 \ 0 ; 0 \ 2]$$



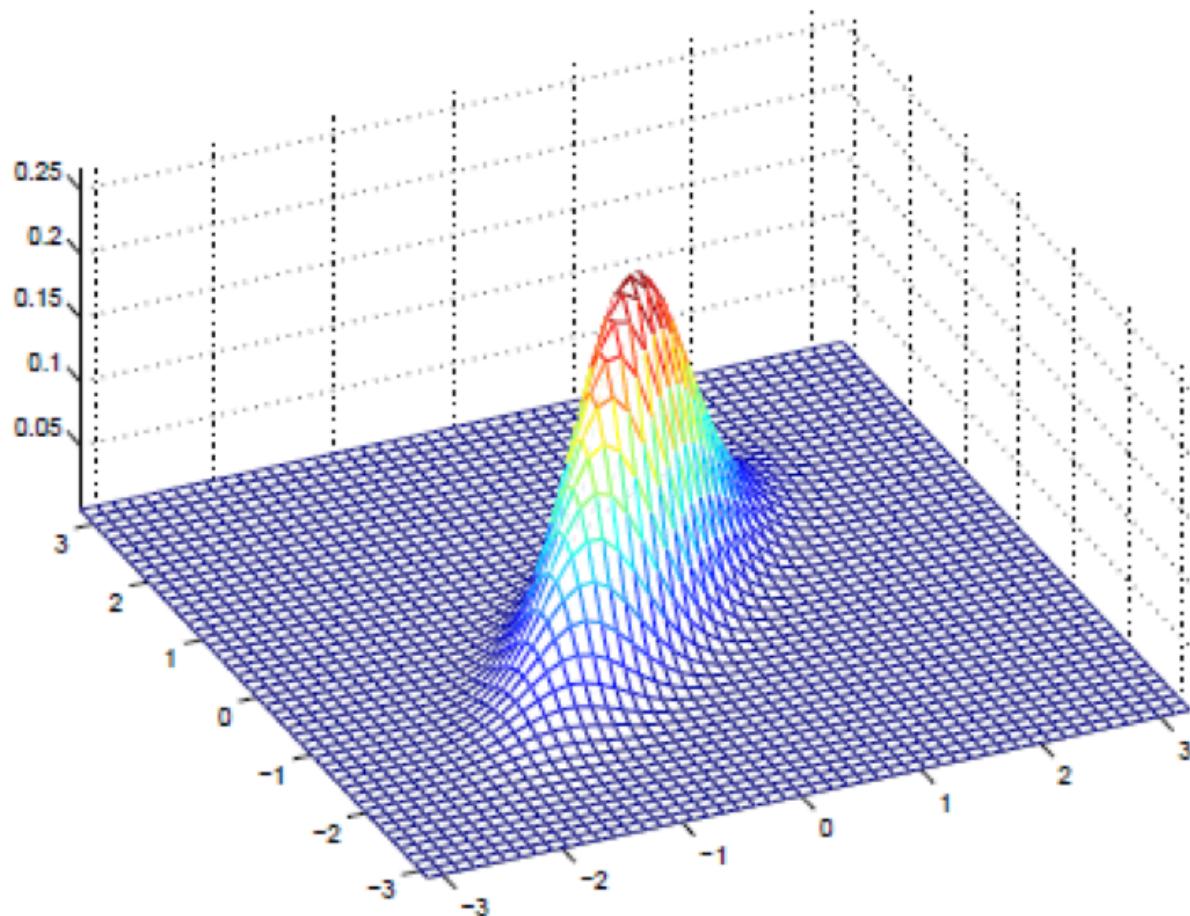
$$\mu = [0; 0]$$

$$\Sigma = [1 \ 0.5; 0.5 \ 1]$$



$$\mu = [0; 0]$$

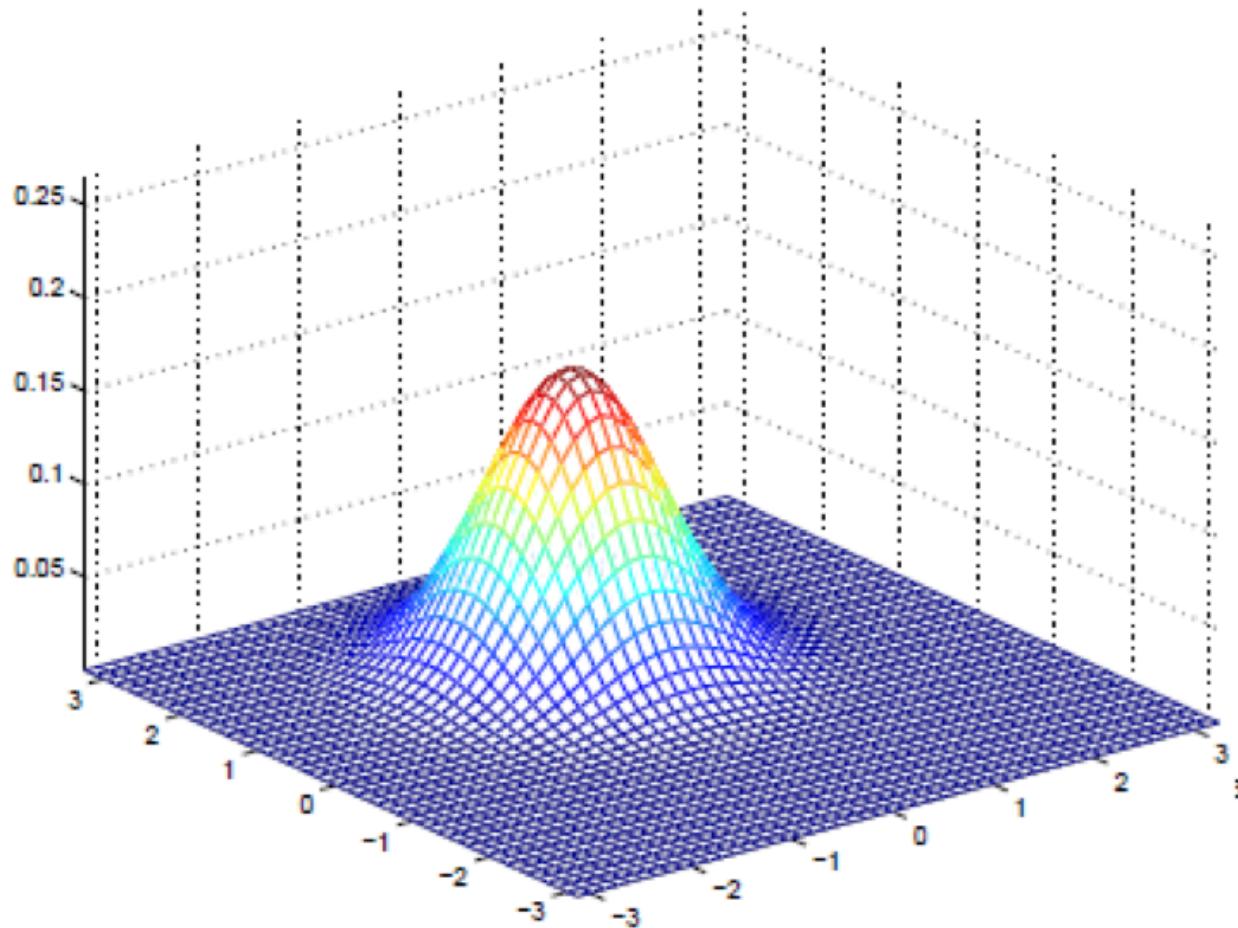
$$\Sigma = [1 \ 0.8; 0.8 \ 1]$$



# Now let's visualize as $\mu$ changes

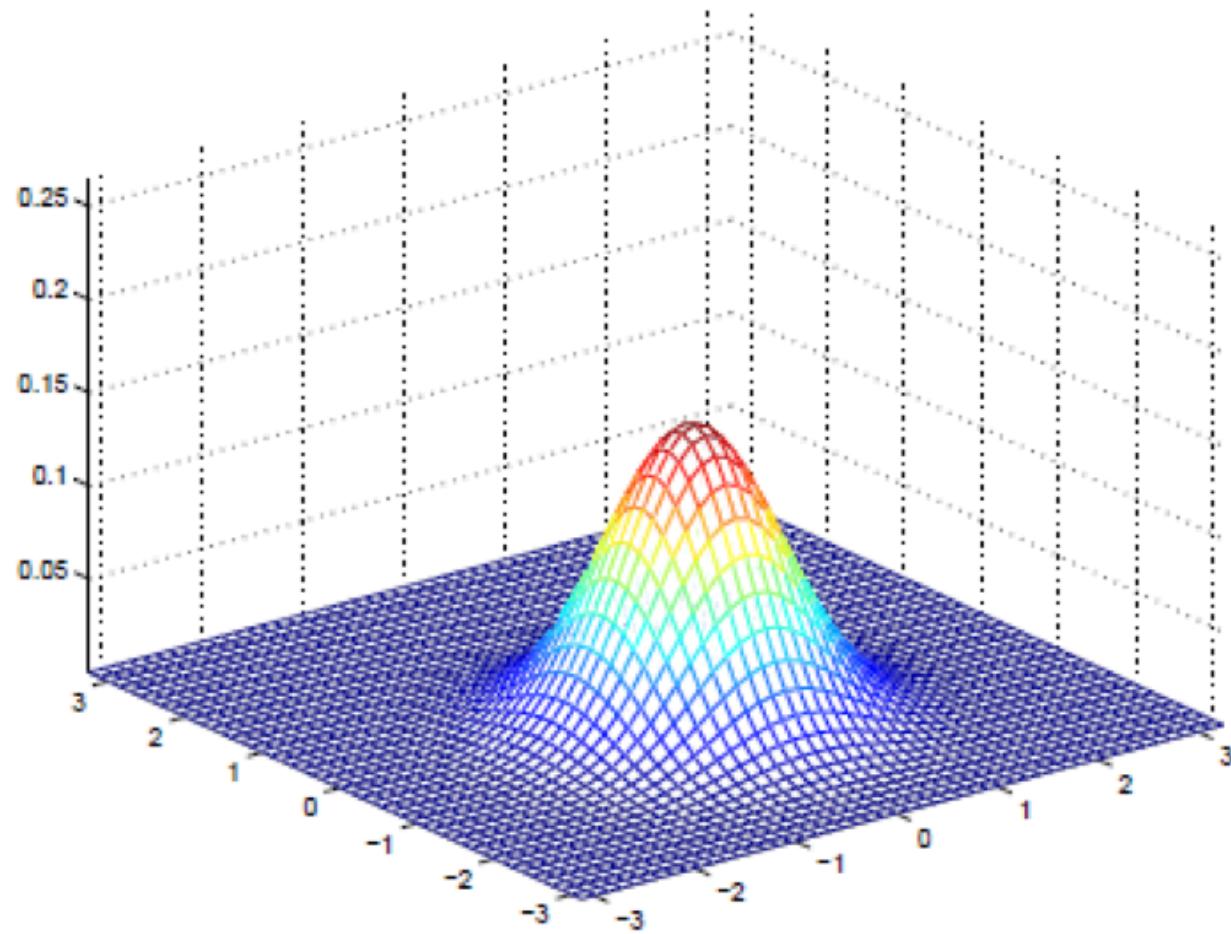
$$\mu = [1; 0]$$

$$\Sigma = [1 \ 0; 0 \ 1]$$



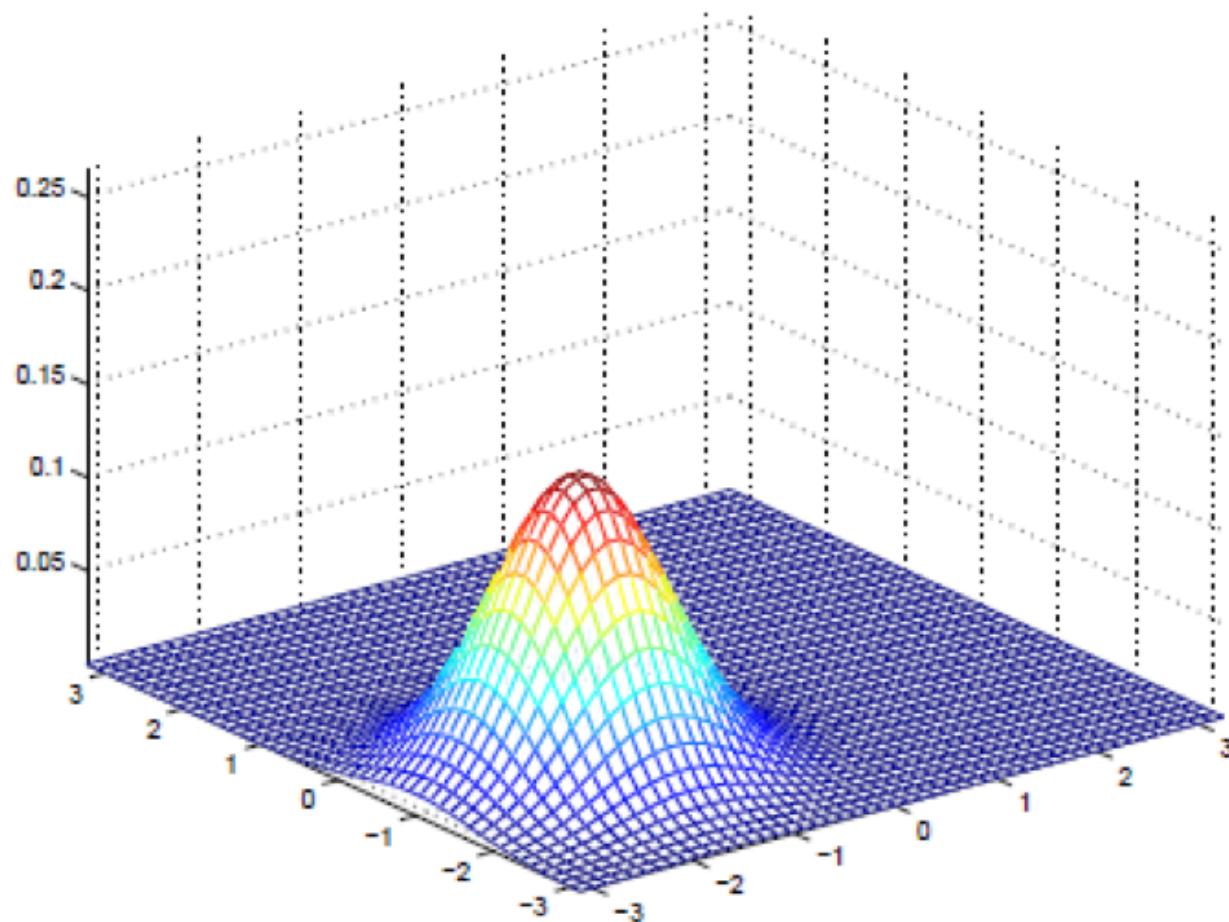
$$\mu = [-.5; 0]$$

$$\Sigma = [1 \ 0; 0 \ 1]$$



$$\mu = [-1; -1.5]$$

$$\Sigma = [1 \ 0; 0 \ 1]$$



# Level sets visualization

$$\mu = [0; 0]$$

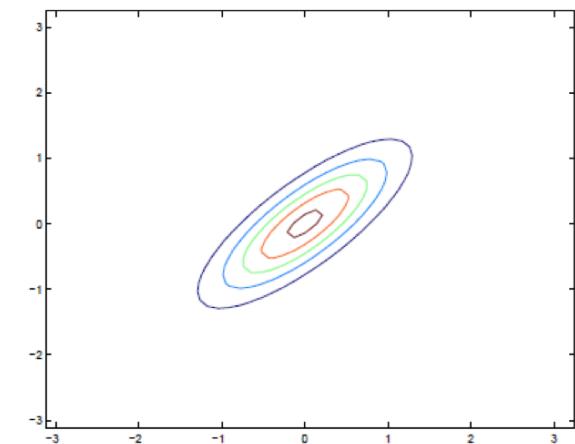
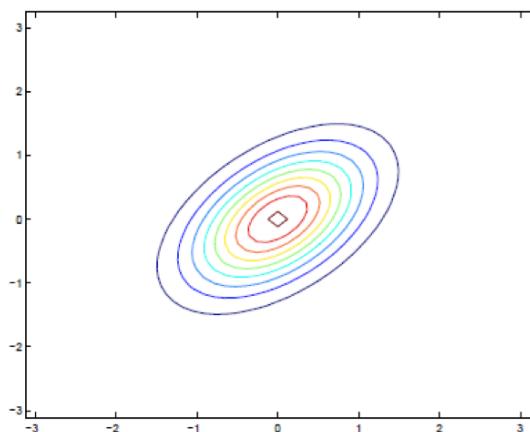
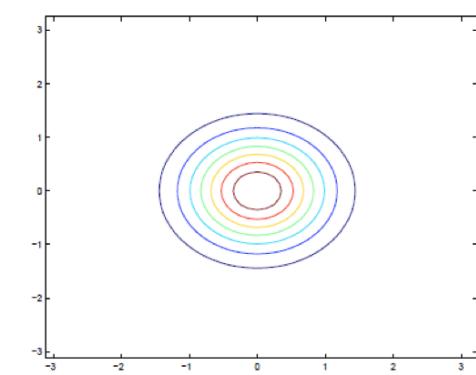
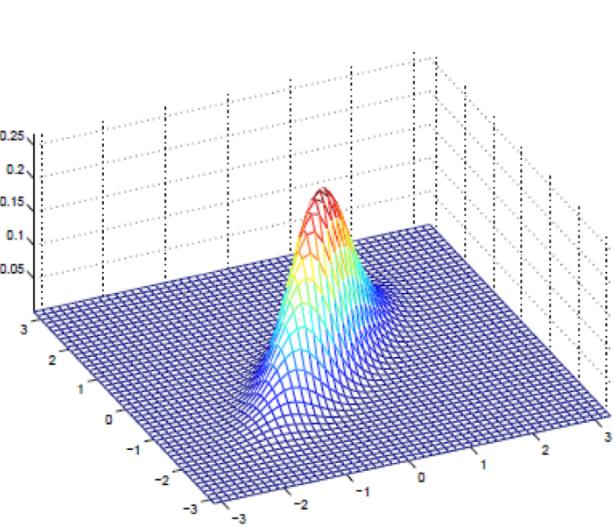
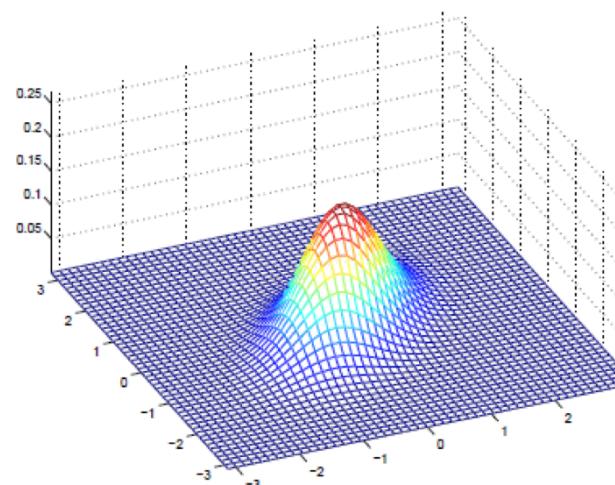
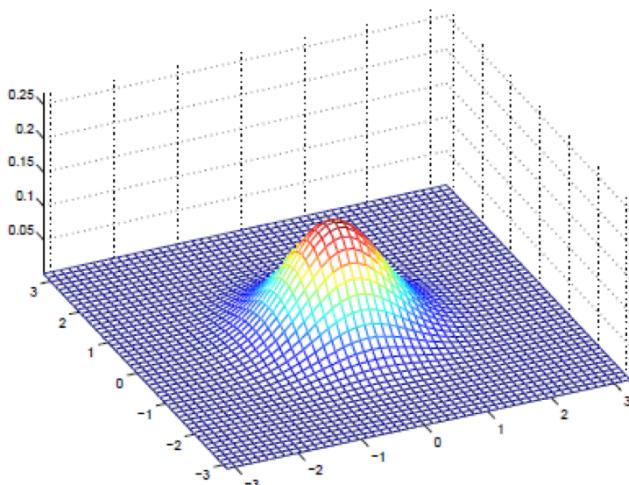
$$\Sigma = [1 \ 0; 0 \ 1]$$

$$\mu = [0; 0]$$

$$\Sigma = [1 \ 0.5; 0.5 \ 1]$$

$$\mu = [0; 0]$$

$$\Sigma = [1 \ 0.8; 0.8 \ 1]$$

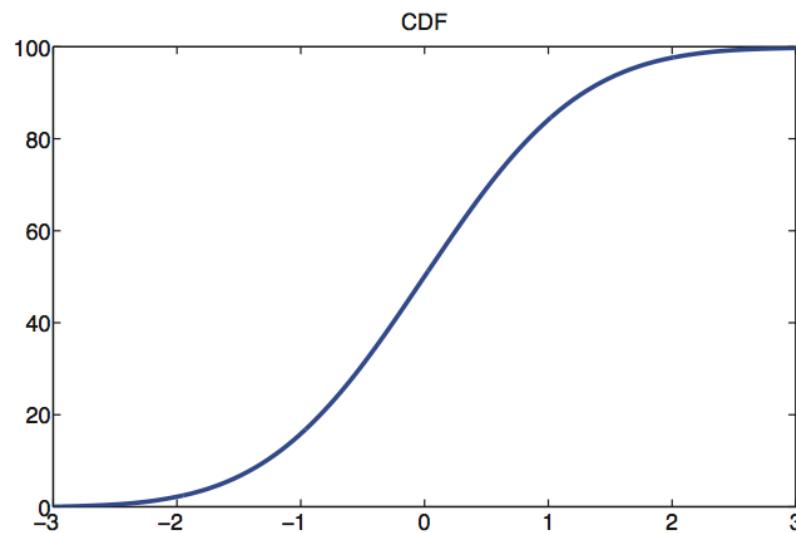


# The cumulative distribution function (cdf)

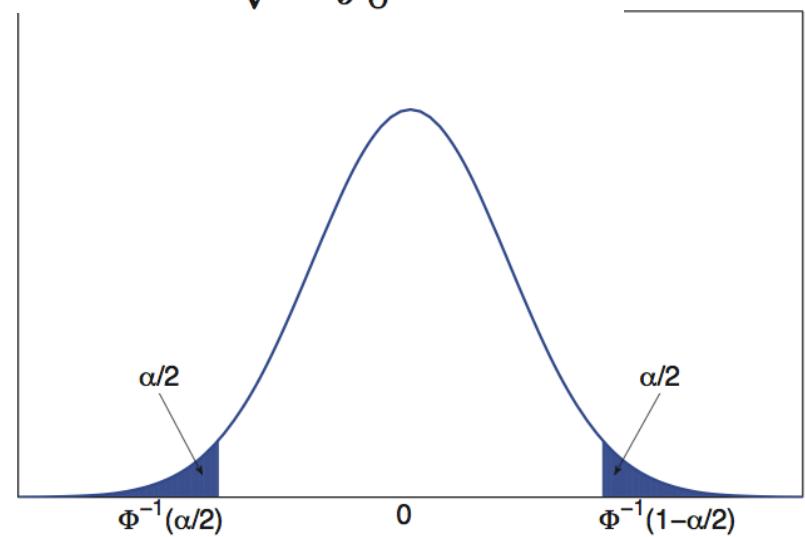
- For Gaussian distribution:  $\Phi(x; \mu, \sigma^2) \triangleq \int_{-\infty}^x \mathcal{N}(z|\mu, \sigma^2) dz$
- This integral has no closed form expression, but is built in to most software packages.

$$\Phi(x; \mu, \sigma) = \frac{1}{2}[1 + \text{erf}(z/\sqrt{2})] \quad \text{where } z = (x - \mu)/\sigma \text{ and}$$

$$\text{erf}(x) \triangleq \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$



(a) Plot of the cdf for the standard normal,  $\mathcal{N}(0, 1)$ .



(b) Corresponding pdf.

# About your homework...

## Beta Distribution

Study it in detail - Homework

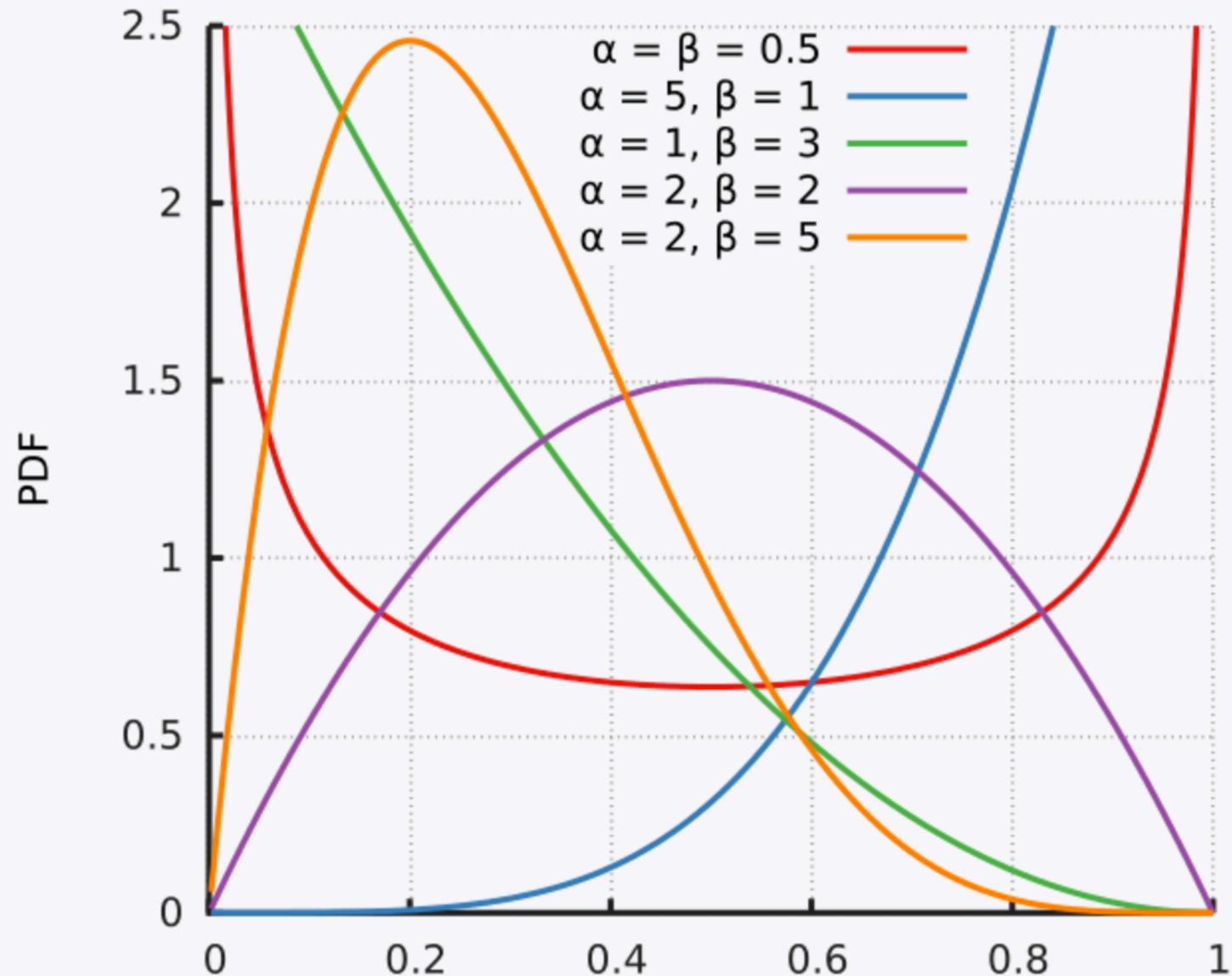
**PDF**

$$\frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$$

$$\text{where } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

# Beta

## Probability density function



# Review: Probability of an Event

- $p(A)$  denotes the probability that the event A is true.
- For example:
- A = a logical expression “it will rain tomorrow”

We require that  $0 \leq p(A) \leq 1$ .

$p(A) = 0$  means the event definitely will not happen

$p(A) = 1$  means the event definitely will happen

$p(\bar{A})$  denotes the probability of the event not A

$$p(\bar{A}) = 1 - p(A)$$

We also write:

A=1 to mean the event A is true.

A=0 to mean the event A is false.

# Review: Fundamental Rules

$$\begin{aligned} p(A \vee B) &= p(A) + p(B) - p(A \wedge B) \\ &= p(A) + p(B) \text{ if } A \text{ and } B \text{ are mutually exclusive} \end{aligned}$$

$$p(A, B) = p(A \wedge B) = p(A|B)p(B)$$

$$p(A) = \sum_b p(A, B) = \sum_b p(A|B = b)p(B = b)$$

$$p(X_{1:D}) = p(X_1)p(X_2|X_1)p(X_3|X_2, X_1)p(X_4|X_1, X_2, X_3)\dots p(X_D|X_{1:D-1})$$

- Independence (or unconditionally independent or marginally independent)  
denoted  $X \perp Y$ :

$$X \perp Y \iff p(X, Y) = p(X)p(Y)$$

- Conditional Independence

$$X \perp Y|Z \iff p(X, Y|Z) = p(X|Z)p(Y|Z)$$

Theorem:  $X \perp Y|Z$  iff there exist function  $g$  and  $h$  such that

$$p(x, y|z) = g(x, z)h(y, z)$$

for all  $x, y, z$  such that  $p(z) > 0$ .

The **conditional probability** of event A,  
given that event B is true:

$$p(A|B) = \frac{p(A, B)}{p(B)} \text{ if } p(B) > 0$$

**Bayes rule:**

$$p(X = x|Y = y) = \frac{p(X = x, Y = y)}{p(Y = y)} = \frac{p(X = x)p(Y = y|X = x)}{\sum_{x'} p(X = x')p(Y = y|X = x')}$$