# Mathematical Basics

Govind Gopakumar

IIT Kanpur

## Announcements

- Pre-Course survey
- Programming assignments
- Project ideas and partners
- Installation of Jupyter / IPython notebook

# Overview

## Notations

**Dealing with data :**

- X : Data matrix (NxD)
- Y : Label matrix (Nx1)
- w : Model parameters
- $L(X,Y,w)$ : Loss of model w on X,Y

**Dealing with model:**

- $\lambda$ : Hyper parameters of a model
- $w^*$ : Optimal model (may or may not be unique)

## Mathematics in Machine Learning

- How do we describe and manipulate data?
- How do we "model" something?
- How do we analytically solve models?
- How do we mathematically "learn"?

# Probability

### Definitions

- Event : Some occurence that is desirable
- Sample space : All possible events
- $P(a) = \frac{\|a\|}{\|a\| + \|a'\|}$

### Terms

- $\prod p(a_i)$ - probability of multiple events
- Can also model likelihood of event
- Naturally leads to MLE (general technique, to be covered later)

## Random Variables

**What are they?**

- Map between events and some value
- Represented as a probability distribution function
- Discrete, continuous, categorical etc

**How do we use them?**

- Describe p(a) for a random variable
- Examples include normal, beta, poisson
- Integrate to 1

## Bayes theorem

**Invert the event!**

- Reverse the probability of events
- $P(a|b) = \frac{P(b|a)P(a)}{P(b)}$

**Terms in this expression**

- $P(a|b)$ - called the posterior
- $P(b|a)$ - called the likelihood
- $P(a)$ - called the prior

## Distributions - I

### Continuous

- Gaussian : Model any real number distribution
- Beta : Model number between [0,1]
- Dirichlet : Model a vector that sums to 1

### Discrete

- Bernoulli : Model number of heads in a coin toss
- Poisson : Model counts of a variable

These can be combined together (joint, marginal)

**Gaussian distribution :**

- $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-1}{2\sigma^2}(x-\mu)^2}$
- $\mu$ : Mean of the distribution
- $^2$ : Variance of the distribution

**Multivariate Gaussian :**

- $p(x) = \frac{1}{\sqrt{2\pi^k|\Sigma|}} e^{\frac{-1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$
- $\mu$ : Mean vector
- $\Sigma$ : Covariance matrix

# Statistics

## Statistics of a sample - I

**Mean of sample**

- $\mathbb{E}[X]$ - "average" of the distribution
- When can it be useless?
- When can it work as a representation?

**Variances and covariances**

- $\sigma^2$ - "spread" of the distribution
- Can be used to "normalize" data
- Can be used to see where data is useless

Generally, we do not come across other "moments" of the data in Machine Learning (skew, kurtosis etc).

## Statistics of a sample - II

**Of standard distributions**

- Gaussian : $\Sigma$
- Bernoulli : $p(1 - p)$

**Of a sample**

- Defined as "empirical" quantities
- Mean : $\mu$
- Variance / Covariance
- Used in "moment matching" techniques

# Linear Algebra

**Constituents :**

- Basis of the space
- Dot product or similarity measure

**Utility :**

- Our data "lives" in some space
- Our model describes "shapes" in that space
- Must deal with math of this space!

## Matrix Algebra

### Basics

- Matrix (NxD) : Can denote a set of points
- Vector (1xD) : Denotes a single point
- Usually denotes our data

### Properties

- Invertibility : $AA^{-1} = I$
- Definiteness : PD / PSD

## Other terms

### Eigenvalues

- $Av = \lambda v$ : $\lambda$ is an "eigenvalue"
- Denotes a direction in the space of the matrix

### Norm of vectors

- $\|x\|_p$ - denotes the p-norm
- Different norms have different interpretations

# Functions and Optimization

## Function shapes

### Convexity

- Convex (and concave) functions have single optima
- Easy to optimize over
- Follow the slope method
- Closed under summation (this is very very nice and important!)

### Smoothness and differentiability

- If a function is "smooth", it will be easy to find the slope.
- If it has kinks, slightly harder to find actual gradients!
- If it is discontinuous, no real way to find gradients!

## Optimization theory

**Basics :**

- Gradient descent : how to follow the slope
- Simple gradients for simple loss functions
- Combine gradients for sum of functions

**Examples of gradients :**

- $(w - x)^2 : 2(w - x)$
- $e^{-w} : -e^{-w}$

## Example of gradient descent

- For simple functions, easy to compute gradients
- General form of GD : $x^{t+1} = x^t - \eta g^t$
- Consider : $f(x) = (x + c)^2$
- Gradient : $g(x) = 2(x + c)$

Let's do gradient descent on this!

# Modelling

## Probabilistic modelling

**Coin tossing : model**

- What do we wish to model? : bias of coin (k)
- What data do we have? : H heads, T tails observed

**MLE modelling**

- p(H heads, T tails)?
- What can we do with this now?
- "Likelihood" can be our loss!
- What is the optimal choice here?
- Why could this fail?

**Conclusion**

## Takeaways

- How to write down probability of events
- What the mean and variance tell us about a random quantity
- Why matrices are used in Machine Learning, how we manipulate them
- What sort of loss functions should we consider? How do we actually use them?

## References

- Review lecture in CS771, IIT Kanpur
- Linear Algebra Overview
- Probability Overview
- Matrix Algebra Overview