

**10-605/10-805**  
**Machine Learning with**  
**Large Datasets**

Fall 2020

# **Probabilities Recap**

# Outline

- Setup
- Random variables
- Distribution function
- Expectation
- Multivariate Distributions
- Independence
- ROC curve
- Probability in Hashing (birthday paradox)

# Setup

- **Sample Space**

- A set of all possible outcomes or realizations of some random trial.

- **Event**

- A subset of sample space

- **Probability Axioms**

- $P(A) \geq 0$  for every  $A$
- $P(\Omega)=1$ ;
- If  $A_1, A_2, \dots$  are disjoint, then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

# Random variables

- **Definition**

- A random variable is a function that maps from the sample space to the reals ( $X : \Omega \rightarrow \mathbb{R}$ ), i.e., it assigns a real number  $X(\omega)$  to each outcome  $\omega$ .

- **Example**

- $X$  returns 1 if a coin is heads and 0 if a coin is tails.  $Y$  returns the number of heads after 3 flips of a fair coin.
- Random variables can take on many values, and we are often interested in the distribution over the values of a random variable, e.g.,  $P(Y = 0)$

# Distribution function

- **Definition**
  - Suppose  $X$  is a random variable,  $x$  is a specific value that it can take,
  - Cumulative distribution function (CDF) is the function  $F : R \rightarrow [0, 1]$ , where  $F(x) = P(X \leq x)$ .
- **If  $X$  is discrete  $\Rightarrow$  probability mass function:  $f(x) = P(X = x)$ .**

## Distribution function (cont.)

- If  $X$  is continuous  $\Rightarrow$  probability density function for  $X$  if there exists a function  $f$  such that  $f(x) \geq 0$  for all  $x$ ,

$$\int_{-\infty}^{\infty} f(x)dx = 1$$

and for every  $a \leq b$ ,

$$P(a \leq X \leq b) = \int_a^b f(x)dx$$

If  $F(x)$  is differentiable everywhere,  $f(x) = F'(x)$ .

# Example of distributions

Discrete variable	Probability function	Mean	Variance
<b>Uniform</b> $X \sim U[1, \dots, N]$	$1/N$	$\frac{N+1}{2}$	
<b>Binomial</b> $X \sim \text{Bin}(n, p)$	$\binom{n}{x} p^x (1-p)^{(n-x)}$	$np$	
<b>Geometric</b> $X \sim \text{Geom}(p)$	$(1-p)^{x-1} p$	$1/p$	
<b>Poisson</b> $X \sim \text{Poisson}(\lambda)$	$\frac{e^{-\lambda} \lambda^x}{x!}$	$\lambda$	
Continuous variable	Probability density function	Mean	Variance
<b>Uniform</b> $X \sim U(a, b)$	$1/(b-a)$	$(a+b)/2$	
<b>Gaussian</b> $X \sim N(\mu, \sigma^2)$	$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$	$\mu$	
<b>Gamma</b> $X \sim \Gamma(\alpha, \beta) (x \geq 0)$	$\frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}$	$\alpha\beta$	
<b>Exponential</b> $X \sim \text{exponen}(\beta)$	$\frac{1}{\beta} e^{-x/\beta}$	$\beta$	

# Expectation

- **Expected Values**

- Discrete random variable  $X$

$$E[g(X)] = \sum_{x \in \chi} g(x)f(x)$$

- Continuous random variable  $X$

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



## Expectation (cont.)

- Mean and variance

$$\mu = E(X)$$

$$\text{var}[X] = E[(X - \mu)^2]$$

We also have

$$\text{var}[X] = E[X^2] - \mu^2$$

# Multivariate Distributions

- **Definition**

$$F_{X,Y}(x, y) := P(X \leq x, Y \leq y)$$

and

$$f_{X,Y}(x, y) := \frac{\partial^2 F_{X,Y}(x, y)}{\partial x \partial y}$$

- **Marginal Distribution of X (discrete case)**

$$f_X(x) = P(X = x) = \sum_y P(X = x, Y = y) = \sum_y f_{X,Y}(x, y)$$

**What about continuous variable?**

# Independence

- **Independent Variables**

- X and Y are independent if and only if

$$P(X = x, Y = y) = P(X = x)P(Y = y)$$

Or

$$f_{X,Y}(x, y) = f_X(x)f_Y(y)$$

# Independence (cont.)

- **IID variable**

- Independent and identically distributed (IID) random variables are drawn from the same distribution and are all mutually independent.

- **Linearity of Expectation**

- Even if the events are not independent, this property still holds

$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$$

# ROC curve

- Confusion matrix

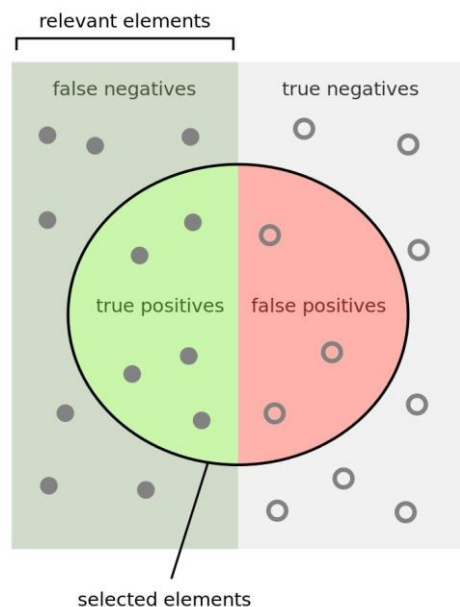
		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

$$\text{True Positive Rate (TPR)} = \frac{TP}{TP+FN}$$

$$\text{False Positive Rate (FPR)} = \frac{FP}{TN+FP}$$

# ROC curve

- **Statistics Computed from Confusion Matrix**
  - **Precision:** Out of all the predicted positive instances, how many were predicted correctly.
  - **Recall:** Out of all the positive classes how many instances were identified correctly (Same as TPR)

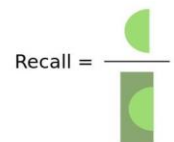


How many selected items are relevant?



$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

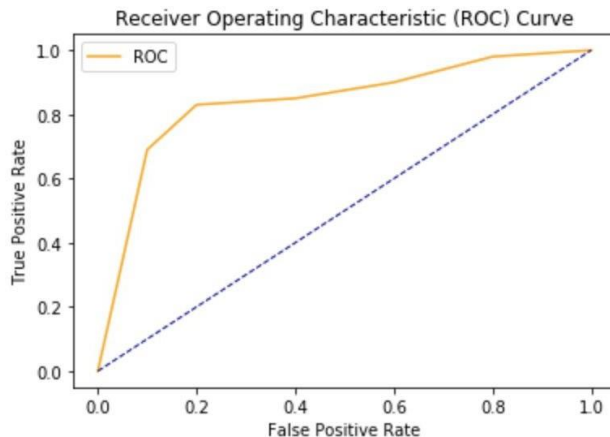
How many relevant items are selected?



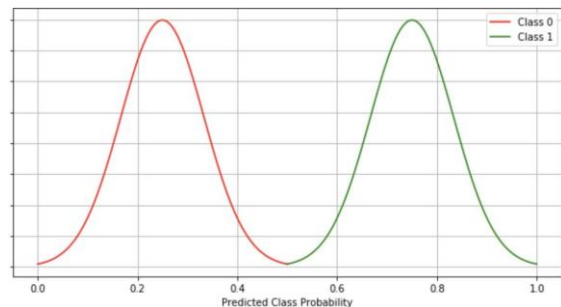
$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

# ROC curve

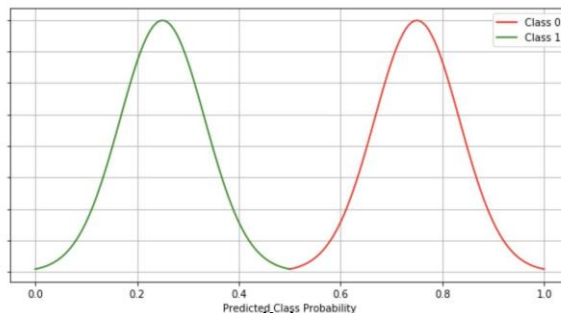
- **Introduction to AUC - ROC Curve**
  - how good the model is for distinguishing the given classes, in terms of the predicted probability



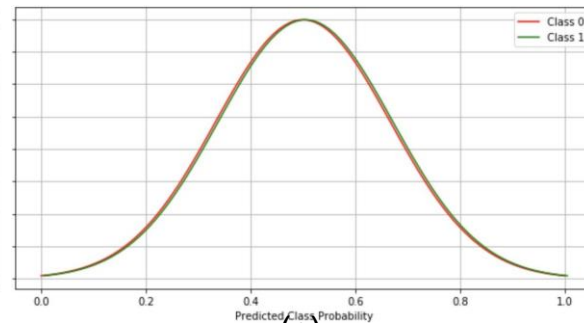
# ROC curve



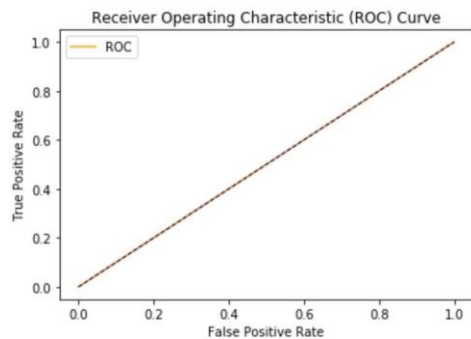
(a)



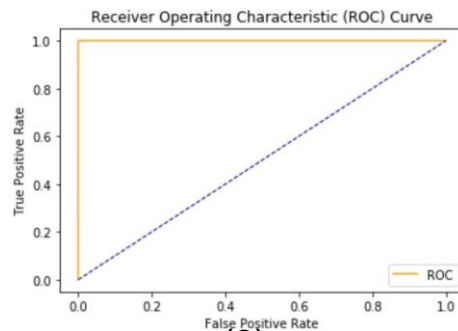
(b)



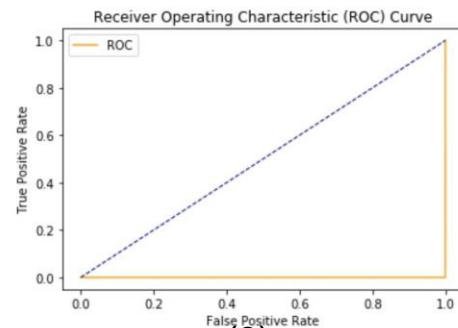
(c)



(1)



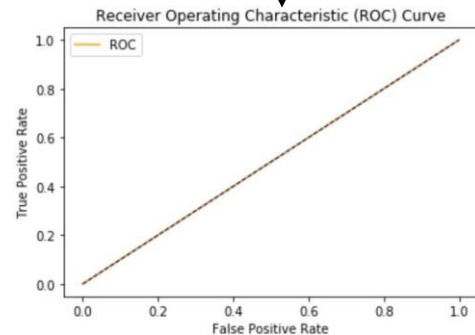
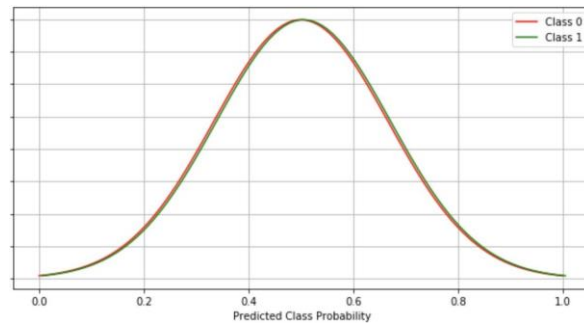
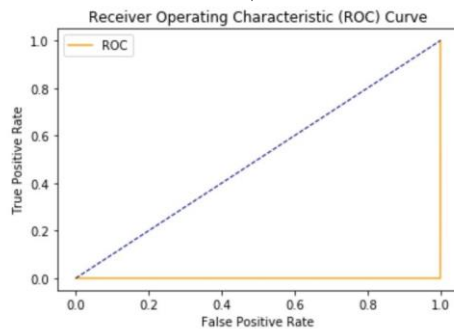
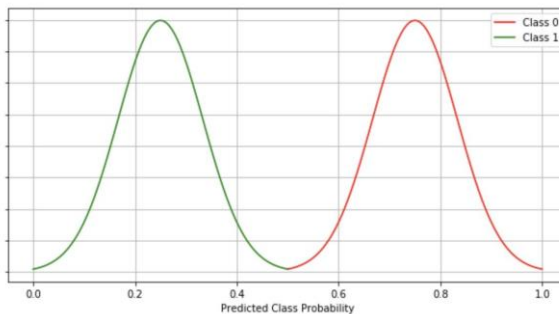
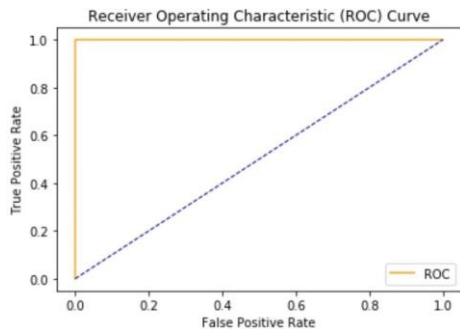
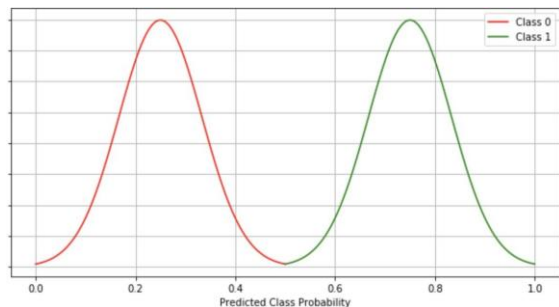
(2)



(3)



# ROC curve



# Probability in Hashing



# Probability in Hashing

- **Assumption**

- $n$ =number of people
- $k=365$
- $P(\text{ person } i \text{ is born on day } j) = 1/k$

We are interested in the event  $A$  that at least two people have the same birthday.

$$\begin{aligned} P(A) &= 1 - P(\bar{A}) \\ &= 1 - \frac{k}{k} \cdot \frac{k-1}{k} \cdot \dots \cdot \frac{k-n+1}{k} \\ &= 1 - \frac{k!}{(k-n)!k^n}. \end{aligned}$$

# Probability in Hashing

- **Hashing**
  - Similar to assignments of birthdays
  - $n$  items mapped into  $k$  slots
- **Hashing problems dealing with probabilities**
  - the expected number of items mapping to same slot
  - the expected number of empty slots
  - the expected number of collisions

# Probability in Hashing

- **Items per slot**

- Consider the following indicator random variable

$$X_i = \begin{cases} 1 & \text{if item } i \text{ is mapped to slot 1;} \\ 0 & \text{otherwise.} \end{cases}$$

- The number of items mapped to slot 1 is therefore

$$X = X_1 + X_2 + \dots + X_n$$

- The expected number of items mapped to slot 1 is

$$E(X) = \sum_{i=1}^n E(X_i) = \frac{n}{k}.$$

# Probability in Hashing

- **Empty slots**

- The probability that slot  $j$  remains empty after mapping all  $n$  items is

$$\left(1 - \frac{1}{k}\right)^n$$

- The expected number of empty slots is

$$E(X) = \sum_{j=1}^k E(X_j) = k \left(1 - \frac{1}{k}\right)^n.$$

- If  $k = n$ , we can get a max limitation of 0.367

# Probability in Hashing

- **Collisions**

- $X$  empty slots
- $(k-X)$  items hashed without collision
- $(n-k+X)$  collisions occur

$$\begin{aligned} E(Z) &= n - k + E(X) \\ &= n - k + k \left(1 - \frac{1}{k}\right)^n. \end{aligned}$$