

# MLF

ML

→ Experience E,

Class of Problems T,

Performance P

(Learn)

(with respect to)

(improve)

Input → Model → Output

Human → Algorithm ← Data

\* Train Data

\* Test Data

\* Validation Data.

Classification :  $f = \text{sign}(w^T x + b)$  [Discrete]

$$L = \frac{1}{n} \sum_i \mathbb{1}(f(x^i) \neq y^i)$$
 [sup.]

Regression :  $f = w^T x + b$  [continuous]

$$L = \frac{1}{n} \sum_i (f(x^i) - y^i)^2$$
 [sup.]

Dimensionality Reduction :  $f(x), g(u)$  [Unsup.]

$$L = \frac{1}{n} \sum_i \|g(f(x^i)) - x^i\|^2$$

Density Estimation :  $P(x)$  [Unsup.]

$$L = \frac{1}{n} \sum_i -\log(P(x^i))$$

$\text{Loss}[f] < \text{Loss}[g]$   $\Rightarrow f$  is better.

Loss is aka reconstruction error.

Continuity :  $\lim_{x \rightarrow a^-} f(x) = f(a) = \lim_{x \rightarrow a^+} f(x)$

Differentiability :  $\lim_{x \rightarrow a^-} \frac{f(x) - f(a)}{x - a} = f'(a) = \lim_{x \rightarrow a^+} \frac{f(x) - f(a)}{x - a}$

Linear approximation :  $L(x) = f(a) + f'(a)(x-a)$

$n^{th}$  order approximation :

$$L(x) = f(a) + \frac{f'(a)(x-a)}{1!} + \frac{f''(a)(x-a)^2}{2!} + \dots$$

Approximations for multivariate functions :

a) Linear :  $L(x, y) = f(a, b) + \frac{\partial f}{\partial x} \Big|_{(a, b)} (x-a) + \frac{\partial f}{\partial y} \Big|_{(a, b)} (y-b)$

b) Higher :  $f(x) = f(v) + \nabla f(v)^T (x-v) + \frac{1}{2} (x-v)^T \nabla^2 f(v) (x-v)$

Gradient :  $\nabla f = \left[ \frac{\partial u}{\partial x}, \frac{\partial u}{\partial y} \right]$

Directional derivative :  $D_{\vec{u}} f(x, y) = \nabla f \cdot \vec{u}$

$$= u_1 \frac{\partial f}{\partial x} + u_2 \frac{\partial f}{\partial y}$$

Cauchy-Schwarz Inequality :

$$-\|a\| \cdot \|b\| \leq a^T b \leq \|a\| \cdot \|b\|$$

Matrix :  $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$

$\xleftarrow{x_1} \quad \xleftarrow{x_2} \quad \xleftarrow{x_m}$

↑      ↑      ↑

$u_1 \quad u_2 \quad u_n$

Columns.

Column Space :  $C(A) = \text{span}(u_1, u_2, \dots, u_n)$

Row Space :  $C(A^T) = \text{span}(x_1, x_2, \dots, x_n)$

Null Space :  $N(A) = \{x \mid Ax = 0\}$

Left Null Space :  $N(A^T) = \{y \mid A^T y = 0\}$

Length of vector :  $|v| = \sqrt{v_1^2 + v_2^2 + v_3^2}$

Dot / Inner Product :  $u \cdot v = u_1 v_1 + u_2 v_2 + u_3 v_3$

Rank-Nullity Theorem : Rank + Nullity = # Columns

Projections :

a) Proj. Matrix :  $P = \frac{aa^T}{a^T a}$

d) Error Matrix :

$$\boxed{e = b - P}$$

b) Proj of b onto a :  $\hat{b} = P \times b = \frac{aa^T}{a^T a} b$

c) Projection of a

along b :  $q = \frac{a \cdot b}{|b|} \times b$

Orthogonality :  $a \cdot b = 0$  [dot product = 0]

Least Squares Solution :  $Ax = b$

- (i) Find  $A^T A$  and  $A^T b$
- (ii) Find aug. matrix for  $A^T A x = A^T b$
- (iii) Row Reduction
- (iv) Solution  $\hat{x}$  is LSS.

Linear Regression : Goal  $\Rightarrow$  minimize  $L$

$$L(\theta) = \frac{1}{2} \sum (x_i^T \theta - y_i)^2$$

$$A = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

*feature matrix*

$$L(\theta) = \frac{1}{2} (A\theta - Y)^T (A\theta - Y)$$

$$\text{LSS} \Rightarrow (A^T A)\theta = A^T Y \quad \text{or} \quad \theta = (A^T A^{-1}) A^T Y$$

Maximum Likelihood :

$$L(\theta) = \prod_i \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2} (y_i - \theta^T x_i)^2\right)$$

Polynomial Regression :  $\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_m x^m$

$$\hat{y}(x) = \theta^T \phi(x), \quad \phi(x) = (1, x, x^2, \dots, x^m)$$

$$A = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\text{LSS} \rightarrow (A^T A) \theta = A^T y$$

Regularised Version (Ridge Regression) :

$$L(\theta) = \frac{1}{2} \sum_i (x_i^T \theta - y_i)^2 + \lambda \|\theta\|^2$$

$$\theta_{\text{reg}} = (A^T A + \lambda I)^{-1} A^T y$$

Eigen Values and Eigen Vectors :  $Ax = \lambda x$

$$|A - \lambda I| = 0 \leftarrow \text{Ch. eqn} \rightarrow \lambda \Rightarrow \text{eigen values}$$

$$(A - \lambda I)x = 0 \leftarrow \text{For each } \lambda, \text{ find } x \leftarrow \text{eigen vector}$$

Characteristic equations :

$$(i) \quad 2 \times 2 : \lambda^2 - \text{tr}(A)\lambda + \det(A) = 0$$

$$(ii) \quad 3 \times 3 : \lambda^3 - \text{tr}(A)\lambda^2 + (M_{11} + M_{22} + M_{33})\lambda - \det(A) = 0$$

## Properties of Eigenvalues & Eigenvectors :

- $\lambda_1 + \lambda_2 = \text{tr}(A)$
- $\lambda_1 \cdot \lambda_2 = |A|$  of A
- If  $\lambda_1, \dots, \lambda_m$  are EV of  $A^k$  then  $\lambda_1^k, \lambda_2^k, \dots, \lambda_m^k$  are EV of  $A^{k+1}$
- If  $\lambda_1, \dots, \lambda_m$  are EV of A then  $\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_m}$  are of  $A^{-1}$
- Symmetric matrices  $\Rightarrow$  Real eigenvalues
- Diagonal matrices  $\Rightarrow$  eigenvalues are diagonal elements
- UTM / LTM  $\Rightarrow$  eigenvalues are diagonal elements
- EV of A = EV of  $A^T$

Diagonalisation :  $\Lambda = S^{-1}AS$

- $\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$   $\lambda_1, \lambda_2, \lambda_3$  are eigenvalues

- $S = \begin{bmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ 1 & 1 & 1 \end{bmatrix}$   $x_1, x_2, x_3$  are eigenvectors

- $\Lambda^k = S^{-1}A^kS$

eigenvalues  
(of similar  
matrices are)  
same

$\Lambda$  and A  
is similar matrices

Fibonacci Series :  $0, 1, 1, 2, 3, 5, \dots \{F_k\}$

$$F_k = F_{k-1} + F_{k-2}$$

$$F_k = \frac{1}{\sqrt{5}} \left( \frac{1+\sqrt{5}}{2} \right)^k - \frac{1}{\sqrt{5}} \left( \frac{1-\sqrt{5}}{2} \right)^k$$

Orthogonally Diagonalizable Matrices :

- a) Find eigenvalues & eigenvectors
- b) Normalize eigenvectors
- c) Rest, follow as diagonalisation

$$A = Q \Lambda Q^{-1} \text{ or } A = Q \Lambda Q^T$$

$$\underbrace{Q^T Q = I}_{\text{Cloud}}$$

Complex Matrices :  $\mathbb{C}^n$  : Complex numbers

$$a+ib : a, b \in \mathbb{R}, i = \sqrt{-1} \Rightarrow i^2 = -1$$

a) Conjugate :  $(\bar{a}+ib) = (a-ib)$

$$x \cdot y = \overline{y \cdot x}$$

$$x(cy) = c(x \cdot y)$$

$$(cx) \cdot y = \overline{c}(x \cdot y)$$

b) Inner product :  $\bar{x}^T y = x \cdot y$

c) Length :  $\|x\|^2 = \bar{x}^T x$

$$(A^*)^* = A$$

d) Conjugate Transpose :  $A^* = (\bar{A})^T = \overline{(A^T)}$

$$(AB)^* = B^* A^*$$

Complex Matrices :a) Hermitian :  $A^* = A$ b) Skew-Hermitian :  $A^* = -A$ 

Properties : (i) All eigenvalues of HM are real  
 (ii) Eigenvectors corresponding to different eigenvalues are orthogonal.

(iii) All symmetric matrix are Hermitian

(iv) If EV are distinct, then matrix is diagonalizable

c) Unitary Matrices :  $U^* U = U U^* = I \Rightarrow U^* = U^{-1}$ 

Properties : (i)  $\|Ux\| = \|x\|$

(ii) eigenvalues  $\rightarrow \pm 1$ 

(iii) eigenvectors corresponding to different eigenvalues are orthogonal

Diagonalization :  $A = U \lambda U^*$ 

To show HM is unitarily diagonalizable we show that it is similar to an UTM.

Schur's Theorem : Any matrix  $A_{n \times n}$  is similar to UTM  
 i.e.  $A = U T U^*$

Spectral Theorem : A HM ' $A$ ' is unitarily diagonalizable  

$$\boxed{U^* A U = D}$$

For symmetric matrix,

$$Q^T A Q = D, \quad Q^T Q = I \quad \begin{matrix} \nearrow \text{Orthogonally} \\ \searrow \text{diagonalizable} \end{matrix}$$

• HM  $\iff$  Unitarily diagonalizable

Quadratic Form :  $f(x) = ax^2 + 2bxy + cy^2$

$$f(x, y) = [x \ y] \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

OR,

$$f(x, y) = v^T A v, \quad v = \begin{bmatrix} x \\ y \end{bmatrix}, \quad A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

Partial Derivatives :

a) 1<sup>st</sup> order :  $f_x = \frac{\partial f}{\partial x}, \quad f_y = \frac{\partial f}{\partial y}$

b) 2<sup>nd</sup> order :  $f_{xx} = \frac{\partial^2 f}{\partial x^2}, \quad f_{xy} = \frac{\partial^2 f}{\partial x \partial y}$

$$f_{yy} = \frac{\partial^2 f}{\partial y^2}, \quad f_{yx} = \frac{\partial^2 f}{\partial y \partial x}$$

Maxima & Minima : At  $(P, q)$  if 1<sup>st</sup> order derivative = 0, then its stationary pt.

$$D = \begin{vmatrix} f_{xx} & f_{xy} \\ f_{xy} & f_{yy} \end{vmatrix} = f_{xx}f_{yy} - f_{xy}^2$$

a) minima  $\Rightarrow f_{xx} > 0, D > 0$

b) maxima  $\Rightarrow f_{xx} < 0, D > 0$

c) saddle  $\Rightarrow D < 0$

d) inconclusive  $\Rightarrow D = 0$

Definiteness :  $f = v^T A v, A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}, v = \begin{bmatrix} x \\ y \end{bmatrix}$

+ve	$f > 0$	$a > 0, ac - b^2 > 0$	$EV > 0$
-ve	$f < 0$	$a < 0, ac - b^2 < 0$	$EV < 0$
+ve semi	$f \geq 0$	$a > 0, ac - b^2 = 0$	$EV \geq 0$
-ve semi	$f \leq 0$	$a < 0, ac - b^2 = 0$	$EV \leq 0$
indefinite	$f > 0, f < 0$	$ac - b^2 < 0$	$EV > 0, EV < 0$

## Singular Value Decomposition (SVD):

Given A,

a) Find  $\bar{A}^T A$

b) Find eigenvalues of  $\bar{A}^T A$  & sort them in descending order of their absolute values.  
Find singular values by taking square root of eigenvalues.

$$\sigma_i = \sqrt{\lambda_i}$$

c) Construct S by placing  $\sigma_i$ 's in descending order along its diagonal. Calculate  $S'$

$$S = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \ddots \end{bmatrix}$$

d) Use eigenvalues & compute eigenvectors of  $\bar{A}^T A$   
Place these along columns of V and find  $V^T$ .

$$V = \begin{bmatrix} 1 & 1 \\ x_1 & x_2 \\ 1 & 1 \end{bmatrix}$$

$x_1$ : eigenvector of  $\lambda_1$

$x_2$ : eigenvector of  $\lambda_2$

e) Find  $U = A V S^{-1}$

Result :  $A = U S V^T$

## Principal Component Analysis (PCA) :

a)  $X = [x_1 \dots x_n]$        $x_i = \begin{bmatrix} f_{e1} \\ \vdots \\ f_{em} \end{bmatrix}$

b)  $\bar{x} = \frac{1}{n} \sum x_i$  (mean vector)

c)  $X - \bar{x} = [x_1 - \bar{x}, \dots, x_n - \bar{x}]$   $\curvearrowright \text{symm}(m \times m)$

d)  $C = \frac{1}{n} \sum (x_i - \bar{x})(x_i - \bar{x})^T$  (cov. matrix)

e) Find eigenvalues & eigenvectors of C

$$\lambda = [\lambda_1, \lambda_2, \dots, \lambda_n], \quad U = [u_1, \dots, u_n]$$

f) Find transformed data points

$$\tilde{x}_i = \sum \alpha_j u_j, \quad \alpha_j = \tilde{x}_i^T u_j$$

g) Reconstruction Error :  $J = \frac{1}{n} \sum \|x_i - \tilde{x}_i\|^2$

h) Projected Variance :  $\lambda_1 + \lambda_2 + \dots + \lambda_k$

Gradient Descent : Given  $f(x)$ ,  $x_0$ ,  $n$ ,  $\epsilon$

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

$x_0$  = initial guess

$n$  = # iterations

$\epsilon$  = step.

For multivariate,  $f(x, y)$ ,  $(x_0, y_0)$ ,  $n$ ,  $\epsilon$

a)  $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$

b)  $x_n = x_{n-1} - \epsilon \nabla f$

Newton's Method : Given  $f(x)$ ,  $n$

$n$  = # iterations

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

Taylor's Series :  $f(x + \eta d) = f(x) + \eta d f'(x) + \frac{1}{2} \eta^2 d^2 f''(x) + \dots$

→ G.D :  $x_{t+1} = x_t + \eta \underset{t}{d}$ ,  $d = -f'(x_t)$

$$f(x + \eta d) \sim f(x) + \eta d f'(x) \quad (\text{for small } \eta)$$

$\Rightarrow df'(x) < 0 \quad \left[ \begin{matrix} \text{descent direction} \\ [d = -f'(x)] \end{matrix} \right]$

For higher dimensions :

Taylor Series :  $f(x + \eta d) \sim f(x) + \eta \underbrace{d^T \nabla f(x)}_{\text{small}} \quad d = -\nabla f(x)$

Location vector :  $(t+1)$  step along dir<sup>n</sup>,  $d$ ,  $\boxed{< 0}$

$$x_{t+1} = x_t + \eta d^T \nabla f(x_t), \quad d = -\nabla f(x_t)$$

Constrained Optimization :

a) Minimize  $f(x)$  s.t.  $g(x) \leq 0$

- (i) Descent direction :  $d^T \nabla f(x) < 0$
- (ii) Feasible direction :  $g(x) \leq 0$

b) If  $x^*$  is optimal then,

- (i)  $g(x^*) \leq 0$

     (ii) No descent direction is feasible.

Optimality : Minimize  $f(x)$  s.t.  $g(x) \leq 0$

a) necessary : ~~IT~~ at  $x^*$ ,  $d^T \nabla g(x^*) = 0$

b) optimal config. :  $\nabla f(x^*) = -\lambda \nabla g(x^*)$

Method of Lagrange Multipliers :

a) minimize :  $f(x)$  under  $g(x) = 0$

b) necessary : at  $x^*$ ,  $d^T \nabla g(x^*) = 0$

c) optimal config. :  $\nabla f(x^*) = \lambda \nabla g(x^*)$

$\nwarrow$   
 $L \cdot M$

Convexity :  $S \subseteq \mathbb{R}^d$  is convex if

$$\forall x_1, x_2 \in S, \lambda x_1 + (1-\lambda)x_2 \in S, \forall \lambda \in [0,1]$$

Property : If  $S_1, S_2 \subseteq \mathbb{R}^d$ , then  $S_{12} = S_1 \cap S_2 = \{x : x \in S_1, x \in S_2\}$

Intersection of convex sets is convex.

• Convex Hull :  $S = \{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}^d$

$$z = \lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n$$

$z$  is a convex combination of  $S$  if  $\exists \lambda_1, \lambda_2, \dots, \lambda_n$

$$\text{s.t. } \lambda_i \geq 0, \sum_i \lambda_i = 1$$

$$CH(S) = \left\{ z : z = \sum_i \lambda_i x_i, \exists x_i \in S, \lambda_i \geq 0, \sum_i \lambda_i = 1 \right\}$$

• Convex Hull is a convex set.

Convex Functions : a) A fcn  $f$  is convex if  $\text{epi}(f)$  is convex set.

$$\text{epi}(f) = \left\{ \begin{bmatrix} x \\ z \end{bmatrix} \in \mathbb{R}^d : z \geq f(x) \right\}$$

b)  $f$  is convex if  $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$  holds  $\forall \lambda \in [0,1]$

c)  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex iff for a point  $y$ :

$\underbrace{\text{differentiable}}_{\uparrow} \quad f(y) \geq f(x) + (y-x)^T \nabla f(x)$

d)  $f: \mathbb{R}^d \rightarrow \mathbb{R} \rightarrow \left( \begin{array}{l} \det(H) > 0 \\ \text{EV}(H) \geq 0 \end{array} \right), H \Rightarrow \text{Hessian matrix}$

$\left[ \begin{array}{l} \text{Positive / Positive} \\ \text{definite / semi-definite} \end{array} \right]$

e) For convex fcn, if  $x^*$  is local minima and  $z$  is global minima then,

$$f(z) = f(x^*) \leq f(x), \forall x \in [x^* - \delta, x^* + \delta], \delta > 0.$$

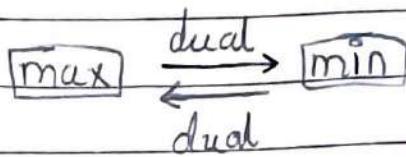
Linear Programming (LPP) : Constraints  $\rightarrow$  Linear  
Objective  $\rightarrow$  Linear

Karush-Kuhn-Tucker Conditions (KKT) :

- a) Stationarity :  $\nabla f(x) + \sum_i u_i \nabla g_i(x) + \sum_j v_j \nabla h_j(x) = 0$
- b) Complementary slackness :  $u_i g_i = 0 \quad \forall i$
- c) Primal feasibility :  $g_i(x) \leq 0 \quad \forall i$
- d) Dual feasibility :  $u_i \geq 0 \quad \forall i$

## Types of possible solutions for LPP :

- a) Infeasible
- b) Unbounded
- c) Feasible



Weak duality :  $g(\lambda^*) \leq f(x^*)$   
 Strong duality :  $g(\lambda^*) = f(x^*)$

Dual	Primal	Finite optimal	Unbounded	Feasible
Finite optimal		✓	✗	✗
Unbounded		✗	✗	✓
Feasible		✗	✓	✓

Sample Space :  $\Omega$  = set of outcomes

Events :  $E \subseteq \Omega$ ,  $F = P(\Omega)$

Probability :  $(\Omega, F, P)$ ,  $P : F \rightarrow [0, 1]$

Axioms :

a)  $P(\Omega) = 1$

b)  $P(A) \geq 0$

c)  $A_1, A_2, \dots, A_n$  are disjoint

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n)$$

Properties :

a)  $P(A^C) = 1 - P(A)$

b)  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

c)  $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B)$

$$- P(B \cap C) - P(A \cap C) + P(A \cap B \cap C)$$

Conditional Probability :  $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Total Probability :  $B_1, B_2, \dots, B_n$  are mut. exc & exh.

$$\forall i, j, \begin{cases} B_i \cap B_j = \emptyset \\ B_1 \cup B_2 \cup \dots \cup B_n = \Omega \end{cases}$$

$$\begin{aligned} P(A) &= P(A \cap B_1) + P(A \cap B_2) + \dots + P(A \cap B_n) \\ &= P(A|B_1) P(B_1) + \dots + P(A|B_n) P(B_n) \end{aligned}$$

Independent events :  $P(A \cap B) = P(A) \cdot P(B)$

Baye's Theorem :  $P(A|B) = \frac{P(B|A) P(A)}{P(B|A) P(A) + P(B|\bar{A}) P(\bar{A})}$

Probability Mass fn :  $X : \Omega \rightarrow \mathbb{R}$  (Random variable)  
 $f_X : \mathbb{R} \rightarrow [0, 1]$

$$f_X(x) = P(X=x) = P(\{w \in \Omega : X(w)=x\})$$

a)  $f_X(x) \in [0, 1]$

b)  $\sum_x f_X(x) = 1$

c)  $\sum_x f_X(x) = \sum_x P(\{w \in \Omega : X(w)=x\}) = P(\Omega) = 1$

Cumulative Distributive fn :  $F_X(x) = P(X \leq x)$

i)  $F_X(x) \in [0, 1]$ , ii)  $F_X(-\infty) = 0$ , iii)  $F_X(\infty) = 1$

Expectation :  $E[X] = \sum_x x f_x(x)$

Conditional Expectation :  $E[X|A] = \sum_x x f_{X|A}(x)$

Linearity : a)  $E[X+Y] = E[X] + E[Y]$   
 b)  $E[aX] = aE[X]$

Variance :  $\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - (E[X])^2$

Joint Distributions :  $f_{XY}(x, y) = P(X=x, Y=y)$

$f_{X,Y}(x, y) = P(\{w \in \Omega : X(w)=x\} \cap \{w \in \Omega : Y(w)=y\})$

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

## Marginal & Conditional Distributions :

$$f_x(x) = P(X=x) = \sum_y f_{xy}(x,y)$$

$$f_y(y) = P(Y=y) = \sum_x f_{xy}(x,y)$$

$$f_{x|y}(x|y) = P(X=x | Y=y) = \frac{f_{xy}(x,y)}{f_y(y)}$$

## Independent Random Variables :

~~$f_{xy}(x,y) = f_x(x) \cdot f_y(y)$~~



$\forall a, b \{X=a\} \& \{Y=b\}$  are independent



$$\forall g, h \quad E[g(x) \cdot h(y)] = E[g(x)] \cdot E[h(y)]$$

## Sum of Random Variables : $Z = X + Y$

~~$f_z(z) = P(Z=z) = \sum_x f_x(x) f_y(z-x)$~~

## Conditional Expectation : $X: \Omega \rightarrow \mathbb{R}$

$$A \subseteq \Omega$$

$E[X|A]$  defined

$$Y: \Omega \rightarrow \mathbb{R}$$

$E[X|Y]$  can be viewed as fxn of  $Y$

Covariance :

$$a) \text{cov}[x, y] = E[(x - EX)(y - EY)]$$

$$b) \text{cov}[x, x] = \text{var}(x)$$

c)  $\text{cov}[x, y] = 0 \Leftrightarrow x, y \text{ uncorrelated.}$

$$d) \text{cov}[x, y] = E[xy] - EX \cdot EY$$

If  $x, y$  independent,  $\text{cov}[x, y] = 0$

Bernoulli Distribution :  $f(1) = p$

$$f(0) = 1-p$$

$x \sim \text{Bernoulli}(p)$

$$E[x] = p$$

$$\text{Var}[x] = p(1-p)$$

Binomial Distribution :  $\text{Bin}(n, p)$

$x_1, x_2, \dots, x_n$  are <sup>independent</sup> Bern(p) RVs

$x \sim \text{Bin}(n, p)$

$$f_x(r) = P(X=r) = {}^n C_R p^r (1-p)^{n-r}$$

$$\sum_{R=0}^n f(r) = 1, \quad E[x] = np$$

$$\text{Var}(x) = npq, \quad q = 1-p$$

Poisson Distribution :  $x \sim \text{Poisson}(\lambda)$

$$f_x(r) = P(X=r) = \frac{e^{-\lambda} \lambda^r}{r!}$$

$$\sum_{R=0}^{\infty} f_x(r) = 1, \quad E[x] = \lambda$$

Geometric Distribution :  $X \sim \text{Geom}(P)$

$$f_x(r) = P(X=r) = (1-P)^{r-1}P$$

$$\sum_{r=1}^{\infty} (1-P)^{r-1}P = 1, E[X] = \frac{1}{P}$$

Variance Properties :  $\text{Var}[X] = EX^2 - (EX)^2$

$$a) \text{Var}[aX] = a^2 \text{Var}[X]$$

$$b) \text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{cov}[X,Y]$$

Continuous Random Variables :

$$a) \text{PDF} : f_x(x) = \frac{P(x \in [x, x+dx])}{dx}$$

$$b) \text{CDF} : F_x(x) = P(X \leq x)$$

$$\text{Properties} : (i) f_x(x) > 0$$

$$(ii) F_x(-\infty) = 0$$

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

$$(iv) F_x(\infty) = 1$$

(v)  $F_x$  is increasing.

$$c) \text{Conditional PDF} : f_{x|A}(x) = \frac{P(x \in [x, x+dx], A)}{dx}$$

Expectation & Variance :

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

$$(i) E[X+Y] = EX + EY$$

$$(ii) Y = g(X), E(Y) = \int_{-\infty}^{\infty} g(x) \cdot f_X(x) dx$$

$$b) \text{Var}[X] = E[(X-EX)^2] = EX^2 - (EX)^2$$

$$c) SD[X] = \sqrt{\text{Var}[X]}$$

$$(i) \text{Var}[X+Y] \neq \text{Var}[X] + \text{Var}[Y]$$

$$(ii) \text{Var}[aX] = a^2 \text{Var}[X]$$

$$(iii) \text{Var}[X] \geq 0$$

Conditional & Total Expectation :

$$a) E[X|A] = \int_{-\infty}^{\infty} x \cdot f_{X|A}(x) dx$$

$$b) E[X] = E[X|A] P(A) + E[X|A^C] P(A^C)$$

Density Functions :  $f_{X,Y} = \frac{P(X \in [x, x+dx], Y \in [y, y+dy])}{dx dy}$

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

$$(i) f_{X,Y}(x, y) \geq 0$$

$$(i) F_{X,Y}(-\infty, -\infty) = 0$$

$$(ii) \iint_{x,y} f_{X,Y}(x, y) dx dy = 1$$

$$(ii) F_{X,Y}(\infty, \infty) = 1$$

Marginals & Conditionals :

$$f_x(x) = \int_{-\infty}^{\infty} f_{xy}(x, y) dy$$

$$f_y(y) = \int_{-\infty}^{\infty} f_{xy}(x, y) dx \quad f_{x|y}(x|y) = \frac{f_{xy}(x, y)}{f_y(y)}$$

Sum of Random Variables :  $Z = X + Y$ 

$$f_z(z) = \int_{-\infty}^{\infty} f_x(x) f_y(z-x) dx$$

↓  
convolution  
operation

Maximum of Independent R.Vs :

$$Z = \max(X, Y)$$

$$F_Z(z) = P(Z \leq z) = F_X(z) \cdot F_Y(z)$$

Minimum of Independent RVs :

$$Z = \min(X, Y)$$

$$F_Z(z) = P(Z \leq z) = 1 - (1 - F_X(z)) \cdot (1 - F_Y(z))$$

## Covariance & Correlation :

$$a) \text{cov}[X, Y] = E(XY) - (EX)(EY)$$

$$b) \rho[X, Y] = \frac{\text{cov}(XY)}{\sqrt{\text{var}[X] \cdot \text{var}[Y]}}$$

\* Independent  $\Rightarrow$  uncorrelated

$$c) X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{cov}[X] = \begin{bmatrix} \text{var}[x_1] & \text{cov}[x_1, x_2] & \dots & \text{cov}[x_1, x_n] \\ \vdots & \ddots & & \vdots \\ \text{cov}[x_n, x_1] & \dots & \ddots & \text{var}[x_n] \end{bmatrix}$$

$$\text{cov}[X] = \begin{bmatrix} \text{var}[x_1] & \text{cov}[x_1, x_2] & \dots & \text{cov}[x_1, x_n] \\ \vdots & \vdots & & \vdots \\ \text{cov}[x_n, x_1] & \dots & \ddots & \text{var}[x_n] \end{bmatrix}$$

} Uniform Distribution :  $X \sim \text{Unif}(a, b)$

$$f_X(x) = \begin{cases} \frac{1}{b-a} & x \in (a, b) \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = \frac{a+b}{2}, \quad \text{Var}[X] = \frac{(b-a)^2}{12}$$

Exponential Distribution :  $X \sim \exp(\lambda)$

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$F_X(x) = P(X \leq x) = 1 - e^{-\lambda x}$$

$$E[X] = \frac{1}{\lambda}, \quad \text{Var}[X] = \frac{1}{\lambda^2}$$

Normal Distribution :  $Z \sim N(0,1)$

$$f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

$$A = \int_{-\infty}^{\infty} e^{-x^2/2} dx \Rightarrow A^2 = \int_0^{\infty} \int_0^{2\pi} e^{-r^2/2} r dr d\theta$$

$$A^2 = 2\pi \Rightarrow A = \sqrt{2\pi}$$

Gaussian Distribution :  $Z \sim N(0,1)$

$$X = \sigma Z + \mu, \quad Z = \frac{X - \mu}{\sigma}$$

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

$$E[X] = \mu, \quad \text{Var}[X] = \sigma^2$$

## Standard Normal Vector :

$$z_1 \sim N(0,1), z_2 \sim N(0,1), \dots, z_d \sim N(0,1)$$

$$\underline{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix}$$

$$f_z(\underline{z}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z_i^2\right) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|\underline{z}\|^2\right)$$

$$\underline{x} = \begin{bmatrix} 1 & 0 \\ p & \sqrt{1-p^2} \end{bmatrix} \underline{z} \quad x_1 = z_1$$

A →

$$x_2 = p z_1 + \sqrt{1-p^2} z_2$$

$$\underline{z} = \begin{bmatrix} 1 & 0 \\ -p & \frac{1}{\sqrt{1-p^2}} \end{bmatrix} \underline{x}$$

$\underline{A}^{-1} \rightarrow$

$$\det(A) = \sqrt{1-p^2}, \quad \det(A^{-1}) = \frac{1}{\sqrt{1-p^2}}$$

$$E[x_1] = E[x_2] = 0$$

$$\text{cov}[x_1, x_2] = p$$

$$\text{var}[x_1] = 1, \quad \text{var}[x_2] = 1, \quad E[\underline{x}] = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{cov}(\underline{x}) = \begin{bmatrix} 1 & p \\ p & 1 \end{bmatrix} = \underline{A} \underline{A}^T, \quad \det(\Sigma) = 1-p^2$$

$$[\text{cov}(\underline{x})]^{-1} = \frac{1}{1-p^2} \begin{bmatrix} 1 & -p \\ -p & 1 \end{bmatrix}, \quad \det(\Sigma^{-1}) = \frac{1}{1-p^2}$$

$$f_x(x) = f_z(A^{-1}x) |\det(A^{-1})|$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x_1^2\right) \frac{1}{\sqrt{2\pi}\sqrt{1-p^2}} \exp\left(-\frac{1}{2(1-p^2)}(x_2 - px_1)^2\right)$$

$$= f_{x_1}(x_1) \cdot f_{x_2|x_1}(x_2|x_1)$$

$$x_1 \sim N(0, 1)$$

$$x_2|x_1=x_1 \sim N(px_1, 1-p^2)$$

Bivariate Normal :  $x = Az$

$$f_x(x) = f_z(A^{-1}x) |\det(A^{-1})|$$

$$\Sigma = AAT = E(xx^T) = \text{cov}(x)$$

$$\det(A^{-1}) = \frac{1}{\sqrt{\det(\Sigma)}} \rightarrow f_x(x) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}x^T\Sigma^{-1}x\right)$$

$$\Sigma = \begin{bmatrix} a^2 & ab \\ ab & b^2 \end{bmatrix}$$

$$x = Az + \mu$$

$$z = A^{-1}(x - \mu)$$

$$f_x(x) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x-\mu)^T\Sigma^{-1}(x-\mu)\right)$$

## Multivariate Normal :

193 Bern

$$\mathbf{X} = \mathbf{A}\mathbf{Z} + \boldsymbol{\mu}$$

$$\mathbf{Z} = \mathbf{A}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

$$\Sigma = \mathbf{A}\mathbf{A}^T$$

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$$

Properties : a)  $\mathbf{y} = \mathbf{a}^T \mathbf{x} \sim N(\mathbf{a}^T \boldsymbol{\mu}, \mathbf{a}^T \Sigma \mathbf{a})$

104 Unif

b)  $\mathbf{y} = \mathbf{A}\mathbf{x} \sim N(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\Sigma\mathbf{A}^T)$

c)  $x_i, x_j$  are independent  $\Leftrightarrow \Sigma_{i,j} = 0$

## Parameter Estimation & Max. Likelihood :

$$\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$$

$x_1, x_2, \dots, x_n$  are iid from  $P_{\theta}$  for  $\theta \in \Theta$

$$L(\theta) = P(x_1 = x_1, \dots, x_n = x_n | \theta)$$

$$= \prod_{i=1}^n f_{x_i}(\mathbf{x}_i | \theta) = \prod_{i=1}^n P_{\theta}(x_i)$$

$$\log(L(\theta)) = \sum_{i=1}^n \log(P_{\theta}(x_i))$$

$$R(\theta) = -\log(L(\theta))$$

Bernoulli Bias :  $P = \{ \text{Bern}(\theta) : \theta \in [0,1] \}$

$$x_1, x_2, \dots, x_n \in \{0, 1\}$$

$$P_\theta(x) = \begin{cases} \theta & x=1 \\ 1-\theta & x=0 \end{cases} = \theta^x (1-\theta)^{1-x}$$

$$R(\theta) = - \sum_{i=1}^n \log(P_\theta(x_i)) = a \log \frac{1}{\theta} + (n-a) \log \frac{1}{1-\theta}$$

$\rightarrow \sum x_i$

$$\hat{\theta}_{ML} = \frac{a}{n} = \frac{\sum x_i^o}{n}$$

Uniform Limits :  $P = \{ \text{Unif}(a,b) : a, b \in \mathbb{R} \}$

$$P_\theta(x) = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{otherwise} \end{cases} = \frac{1}{b-a} 1(x \in [a,b])$$

$$R(\theta) = - \sum_{i=1}^n \log(P_\theta(x_i^o))$$

$$= \sum_{i=1}^n -\log\left(\frac{1}{b-a}\right) - \log(1(x_i^o \in [a,b]))$$

If  $a < \min(x_i^o)$  &  $b > \max(x_i^o)$

$$R(\theta) = n \log(b-a)$$

↑      ↑  
max  $x_i^o$    min  $x_i^o$

Normal Mean & Variance : a)  $P = \{N(\mu, I) : \mu \in \mathbb{R}\}$

$$P_0(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-\mu)^2\right)$$

$$R(\theta) = \sum_{i=1}^n \frac{1}{2} (x_i - \mu)^2 + c$$

$$\mu = \frac{1}{n} \sum x_i$$

b)  $P = \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+\}$

$$P_0(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2} \cdot \frac{(x_i - \mu)^2}{\sigma^2}\right)$$

$$R(\theta) = \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}$$

$$\mu = \frac{1}{n} \sum x_i$$

$$\sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

Multivariate Normal :  $P = \{N(\mu, \Sigma) : \mu \in \mathbb{R}^d, \Sigma \in S_d^+\}$

$$\mu = \frac{1}{N} \sum x_i$$

$$\Sigma = \frac{1}{N} \sum (x_i - \mu)(x_i - \mu)^T$$

## Linear Regression with Gaussian Noise :

$$x \in \mathbb{R}^d, y \in \mathbb{R}$$

$$y = w^T x + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

$$\cdot P_{Y|x} = \{ N(w^T x, \sigma^2) : w \in \mathbb{R}^d \}$$

$$P(Y_1 = y_1, \dots, Y_n = y_n | X_1 = x_1, \dots, X_n = x_n)$$

$$= \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2} (y_i - w^T x_i)^2\right)$$

$$R(w) = \sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 + c$$

## Gaussian Mixture Model :

$$f_x(x) = \sum_{R=1}^k \pi_R \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu_R)^T \Sigma_R^{-1} (x - \mu_R)\right)$$

$$P(x) = \sum_{R=1}^k \pi_R N(x | \mu_R, \Sigma_R)$$

## Max. Likelihood estimation :

$$P(\text{Data}) = \prod_{i=1}^N \left( \sum_{R=1}^k \pi_R N(x_i | \mu_R, \Sigma_R) \right)$$

## Inequalities and Laws :

a) Markov Inequality :  $P(X \geq t) \leq \mu/t$

$$EX = \int_0^{\infty} xf_X(x) dx \geq t \int_t^{\infty} f_X(x) dx$$

b) Chebyshev Inequality :  $P(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2}$

$$\leq \frac{E[(X - \mu)^2]}{t^2}$$

$$\bullet = \frac{\text{var}[X]}{t^2}$$

c) Hoeffding Inequality :  $a \leq x_i \leq b$   
 $\bar{X}_n = \frac{1}{n} \sum x_i$ ,  $\text{var}[\bar{X}_n] = \frac{1}{n^2} n \text{var}[x_i] = \frac{\sigma^2}{n}$

$$P(|\bar{X}_n - \mu| \geq t) \leq \frac{\text{var}[\bar{X}_n]}{t^2} = \frac{\sigma^2}{nt^2}$$

$$P(|\bar{X}_n - \mu| \geq t) \leq 2 \exp\left(-\frac{2nt^2}{(b-a)^2}\right)$$

d) Convergence in Probability :

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq t) = 0 \quad \forall t > 0$$

e) Convergence in Distribution :

$$\lim_{n \rightarrow \infty} |F_{X_n}(x) - F_X(x)| = 0 \quad \forall x$$

f) Law of Large Numbers :

$X_1, X_2, \dots, X_n$  iid from D

$$EX_i = \mu, \bar{X}_n = \frac{1}{n} \sum X_i$$

$$\bar{X}_n \xrightarrow{P} \mu$$

$$P(|\bar{X}_n - \mu| \geq t) \leq \frac{\sigma^2}{nt^2}$$

g) Central Limit Theorem :

$X_1, X_2, \dots, X_n$  iid from D

$$EX_i = \mu, \text{var}[X_i] = \sigma^2$$

$$Y_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu), Y_n \xrightarrow{D} N(0, \sigma^2)$$