

Stanford CS229 : Machine Learning

Complete Handwritten Notes

Devansh Chaudhary

Date	Title	Page No.	Teacher's Sign / Remarks
Lec 1	Intro to Linear Algebra		
Lec 2	Matrix Calculus & Probability		
Lec 3	Probability and Statistics		
Lec 4	Linear Regression		
Lee 5	Perception & logistic Regrs.		
Lee 6	Exponential family & GLM		
Lee 7	GDA, Naive Bayes & Laplace Smooth		
Lee 8	Kernel methods & SVM		
Lee 9	Bayesian Methods - Param & Non		
Lee 10	Deep Learning - I		
Lec 11	Deep Learning - II		
Lec 12	Bias and Variance & Regularisation		
Lec 13	Statistical learning Uniform Conv.		
Lec 14	Reinforcement Learning - I		
Lec 15	Reinforcement Learning - II		
Lec 16	K-means, GMM, EM		
Lec 17	Factor Analysis & ELBO		
Lee 18	Principal & Independent CA		
Lee 19	Maximum Entropy & Calibration		
Lee 20	Variational Autoencoder		
Lee 21	Evaluation Matrices		
Lee 22	Practical tips and Course Recap		
Lee 23	Course Recap and Wrap Up		

Stanford CS229Machine Learning | Lee I - Intro. and Linear Algeb.

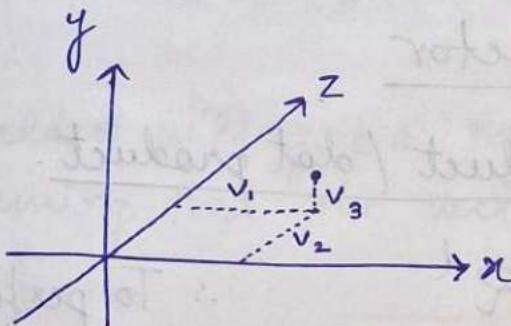
Machine learning is the study of computer algorithms that improve automatically through experience - Tom Mitchell, prof. CMU

Notation

Vector $v \in \mathbb{R}^d$: v is a vector and it is an element of this set and this set is a d -dimensional real space.

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

$$v^T = [v_1, \dots, v_d]$$



$$\mathbb{R}^d \rightarrow \mathbb{R}^3$$

- * Generally for vectors we use small letters notation and for Matrices we use Capital letters notation.

Matrix

$$A \in \mathbb{R}^{m \times n}$$

$$A = \begin{array}{c} \xleftarrow{n} \xrightarrow{n} \\ \downarrow m \quad \uparrow m \\ \begin{bmatrix} a_{11} & \cdots & \cdots & a_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ a_{m1} & \cdots & \cdots & a_{mn} \end{bmatrix} \end{array}$$

$m \rightarrow \text{rows}$
 $n \rightarrow \text{columns}$

Identity Matrix $(I) = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$

Diagonal Matrix $= \begin{bmatrix} \textcircled{1} & & & \\ & \textcircled{2} & & \\ & & \ddots & \\ & & & \textcircled{n} \end{bmatrix}$

Symmetric Matrix $= A = A^T$

$$A_{ij} = A_{ji}^T$$

Trace of a Matrix

$$\sum_i A_{ii}$$

Basic operations

- Vector - Vector

Inner product / dot product.

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

\therefore To perform inner product the vectors have to be of the same dimension

$$\sum_{i=1}^d x_i y_i ; x^T y$$

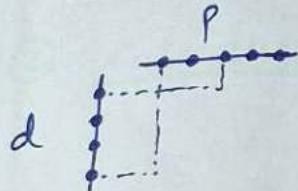
$$\overline{\quad} | d = \cdot \text{ (scalar)}$$

Outer product

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

(The two vectors need not to be of the same dimension)

xy^T , yx^T (The two are different).



$$= \begin{bmatrix} & a_{13} \\ a_{11} & \dots \end{bmatrix}$$

Rank-1 Matrix

- It is made of one pair of row and column.

Taking i^{th} element from x and j^{th} element from y , then multiply it and it gives ij element of the outer product

same dimension p

$$\begin{bmatrix} & + & \end{bmatrix} = \begin{bmatrix} & \end{bmatrix}$$

same dimension d

$$\rightarrow \text{Rank-2}$$

Two Rank-1 matrix added will give a Rank-2 Matrix provided the vectors are linearly independent.

→ LI vectors are a set of vectors that do not lie on the same plane or same line and cannot be expressed as a linear combination of each other.

$$\begin{bmatrix} & + & \end{bmatrix} + \begin{bmatrix} & \end{bmatrix} + \begin{bmatrix} & \end{bmatrix} = \begin{bmatrix} & \end{bmatrix} \leq \min(d, g, k)$$

Matrix - Vector

$$m \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad | \quad n \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$A \in \mathbb{R}^{m \times n}$ $x \in \mathbb{R}^n$

$$Ax = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$Ax \in \mathbb{R}^m$

$$m \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \quad | \quad n \begin{bmatrix} * \\ * \\ * \end{bmatrix}$$

A x

$$Ax = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

same result

Matrix - Matrix

$$m \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad k \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \end{bmatrix}$$

or can be thought of as

$$m \underbrace{\begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix}}_k \quad k \underbrace{\begin{bmatrix} n \\ \vdots \\ n \end{bmatrix}}_r$$

$$c_1 \left| \frac{x_1}{\dots} \right. + c_2 \left| \frac{x_2}{\dots} \right. + \dots = \begin{bmatrix} \vdots \end{bmatrix}$$

Think of it as
K column matrix
(where each col.
have m rows)

think of it as
K row matrix
(where each
row have n
col.).

Rank
 $\leq \min(d, p, k)$ Why Linear Algebra① Represent data

both Kitten -

$$X = \begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix} \quad Y = \begin{bmatrix} \text{cost} \\ \vdots \end{bmatrix}$$

② Covariance Matrices

$$S \in \mathbb{R}^{d \times d}$$

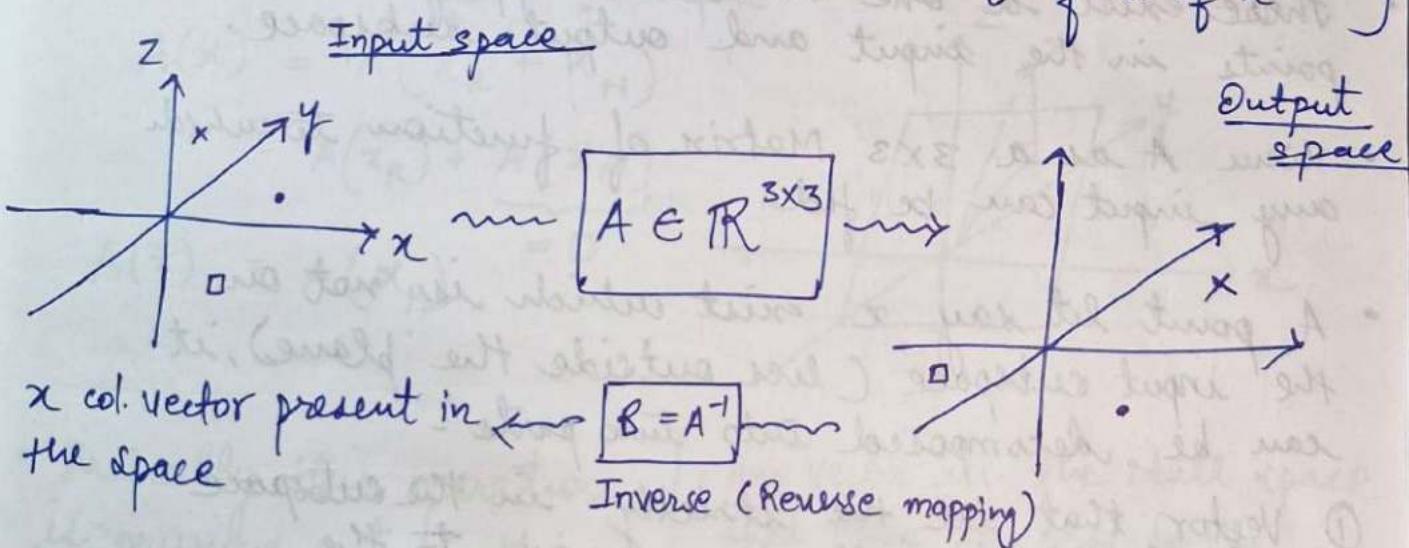
- ③ Calculus → Gradient - vector
 Hessians - Matrix (symmetric)
 Jacobians - Matrix

④ Kernal Method

Geometrical Interpretation

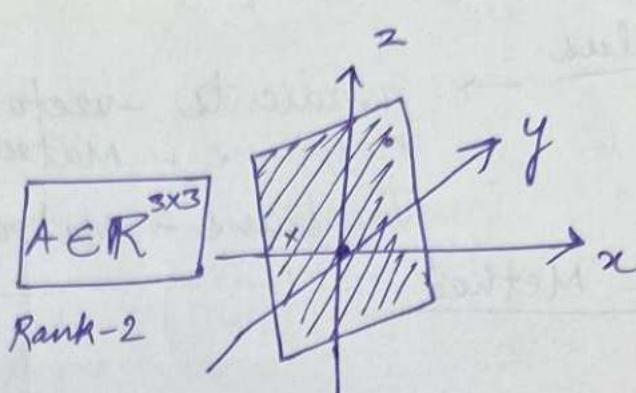
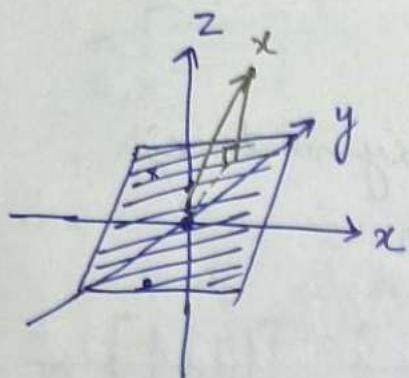
$$A \in \mathbb{R}^{m \times n}$$

$$x \in \mathbb{R}^n \Rightarrow A(x) \in \mathbb{R}^m \quad \left\{ \begin{array}{l} \text{think of } Ax \\ \text{as } A(x) - A \bar{x} \\ \text{a func. of } x \end{array} \right.$$



If A is full rank (i.e. Rank 3) there is going to be unique mapping of every point in the input to every point in the output.

FULL RANK

RANK DEFICIENT

- There exist a 2D subspace ~~in~~ that exists in the 3D space that is specific to this matrix. Also a corresponding 2D subspace in the output
- There exist a one-to-one mapping between the points in the input and output subspace.
- View A as a 3×3 Matrix of function to which any input can be fed
- A point let say x exist which is not on the input subspace (lies outside the plane), it can be decomposed into two parts -
 - ① Vector that lie the strictly in the subspace and it extends all the way to the nearest point to x .
 - ② Point (Vector) ~~to~~ to the subspace which originates from the nearest point to the point x .

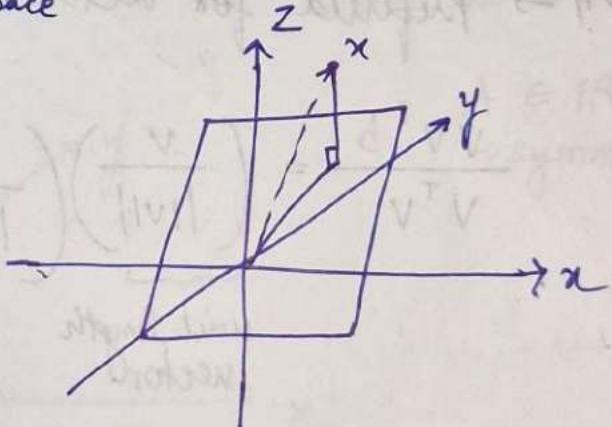
Basically splitting the vector into its components-

$$\textcircled{1} \quad \bar{x} = \text{Proj}(x; \text{Rowspace}) + \text{Proj}(x; \text{Nullspace}).$$

~~\bar{x}~~

- The input is called the rowspace because it is made precisely of those points which can be represented as the linear combinations of the rows of A .
- Similarly for the output space, just the change is that it is the linear combination of the columns of A . (Can be called as the range)
- Basically the input is precisely that subspace for which a bijection exists between the output and input.
- Subspace always passes through the origin by definition

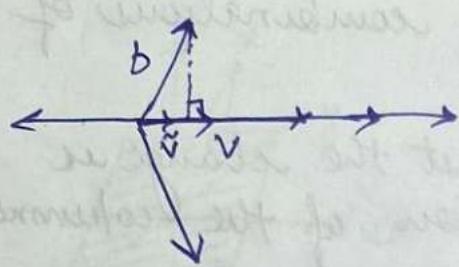
$$\begin{aligned} A(x) &= A(\overset{\text{Row}}{\underset{x_R}{\uparrow}} + \overset{\text{Null space}}{\underset{x_N}{\uparrow}}) \\ &= A(x_R) + \underbrace{A(x_N)}_{=0} \\ A(x) &= A(x_R) \end{aligned}$$



Basically this operation is pruning out anything that is outside the row space from the input and project on the output.

Any vector in the Null space will have its proj on the origin through A .
 $\Rightarrow A(x_N) = 0$

Projection



Projection Matrix

$$(v) = \left[\frac{vv^T}{v^Tv} \right] b$$

$$\frac{vv^T b}{v^Tv} = \left(\frac{v}{\|v\|} \right) \left(\frac{v}{\|v\|} \right)^T b$$

$\|v\| \rightarrow$ Norm of a vector
generalise to length or magnitude

$\|v\| \rightarrow$ Mag. (Preferred for scalars)

$\|v\| \rightarrow$ Preferred for vectors.

$$\frac{vv^T b}{v^Tv} = \underbrace{\left(\frac{v}{\|v\|} \right)}_{\text{unit length vector}} \left(\frac{v}{\|v\|} \right)^T b$$

$$= [\tilde{v} \ \tilde{v}^T] b$$

$$(v) = \underbrace{\tilde{v}}_{\text{scalar}} \underbrace{(\tilde{v}^T b)}_{\text{length}}$$

- If we want to project B onto subspace spanned by multiple vectors. Then in place of v we will use some matrix $x [1|1]$ \rightarrow complete subspace

$$\begin{bmatrix} VV^T \\ V^TV \end{bmatrix} b = V(V^T V)^{-1} V^T b$$

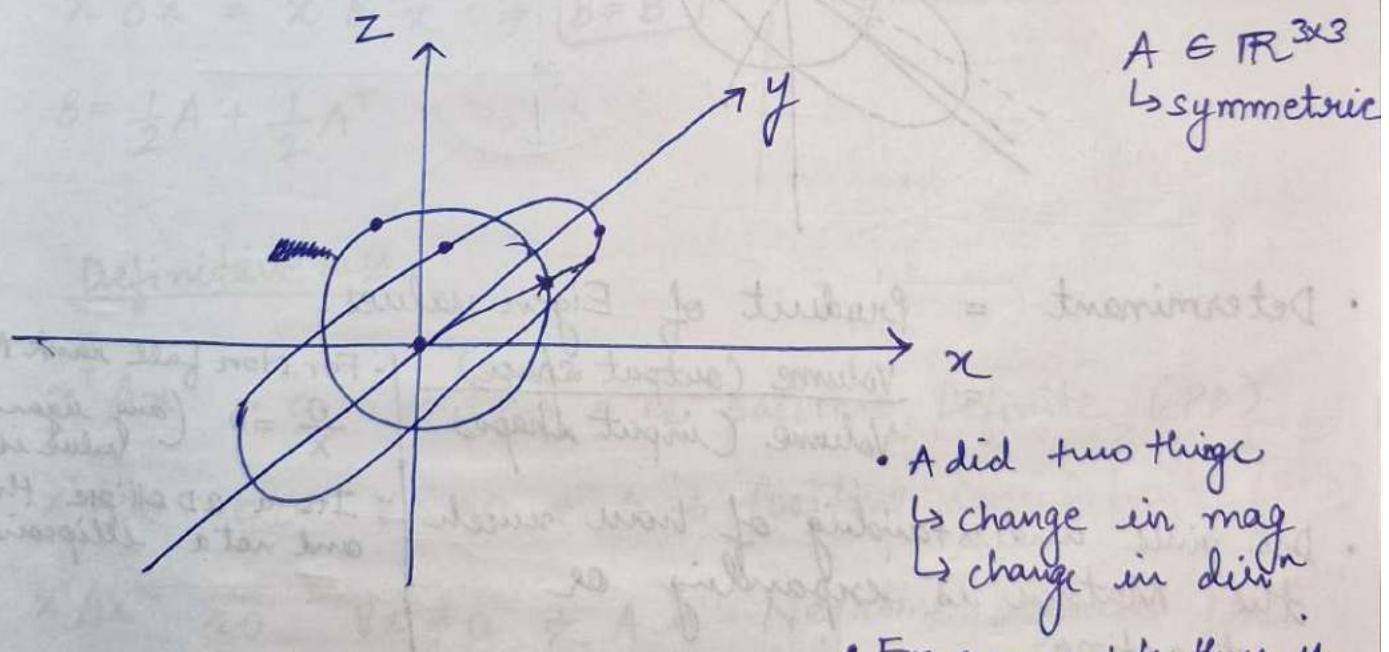
$$\Rightarrow X(X^T X)^{-1} X^T$$

$(m \times n) \quad (n \times n) \quad (n \times m) \rightarrow \text{dimension}$

- For no point in the input space can we reach a point that lies outside the column space or output subspace.

Square Matrices

↳ Input and output space have the same dimension.



On $\rightarrow A \rightsquigarrow$
sphere ellipsoid

- Did two things
 - change in mag
 - change in dir

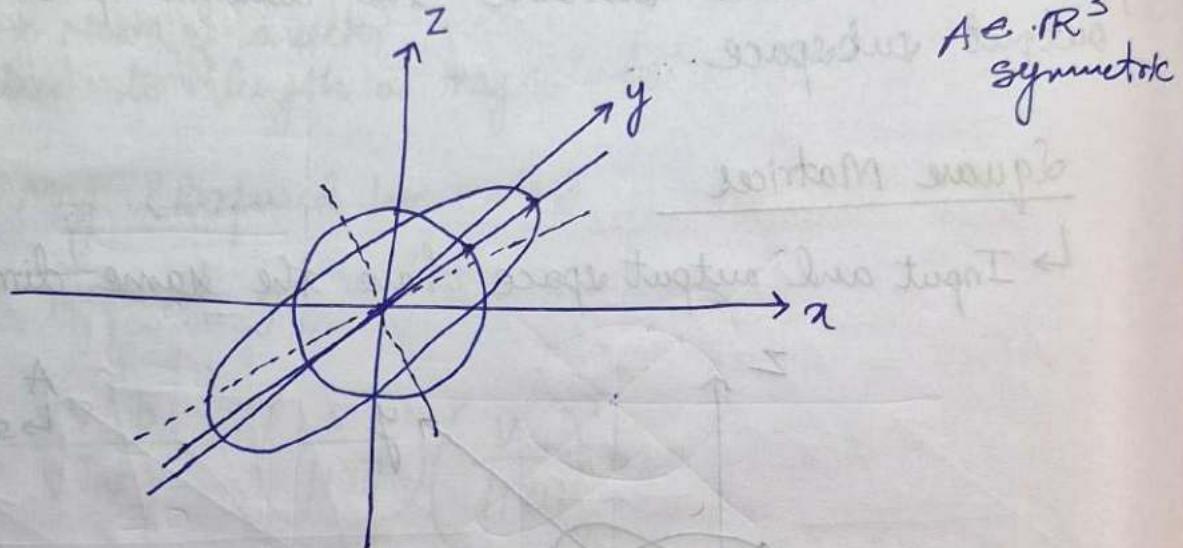
- For some points there is going to be only change in mag. and not direction
⇒ Called eigen vectors.
- It will have 3 eigen vectors.

Rank of a Matrix is precisely the number of non-zero eigen values

eigen value = $\frac{\text{length of output vector}}{\text{length of input vector}}$

Lecture - 2 : Matrix Calculus and Probability theory

- Eigen vectors are those special vectors for square matrices where the operation of the matrix does not change the angle, it only rescales it by some amount.
- The eigen value is the ratio by which the vector gets rescaled, so each eigen vector comes with corresponding eigen value.



- Determinant = Product of Eigen values.
 $= \frac{\text{Volume (output shape)}}{\text{Volume (input shape)}}$
 - For Non full rank Matrix $\frac{0}{x} = 0$ (any eigen value is 0)
 - It's a 2D ellipse then and not a ellipsoid
- Det. gives understanding of how much the Matrix is expanding or contracting.
- If any eigen value is 0. Then there doesn't exist Matrix inverse (2D ellipse cannot be mapped back to a 3D sphere)

Spectrum - Collection of Eigen values

Spectral theorem - For every square matrix $A \in \mathbb{R}^{d \times d}$ which is also symmetric $A = A^T$ has

- Real valued Eigen values
- Orthonormal Eigen vectors.

Quadratic form

$$A \in \mathbb{R}^{d \times d} \quad x \in \mathbb{R}^d$$

$x^T A x \rightarrow$ Quadratic form

$$x^T A x = x^T B x$$

$$x^T B x = x^T B^T x \Rightarrow B = B^T$$

$$B = \frac{1}{2} A + \frac{1}{2} A^T$$

It's good to assume that Matrix A is symmetric too, as it can always be expressed as in the form of other matrix B which is symmetric

Definiteness

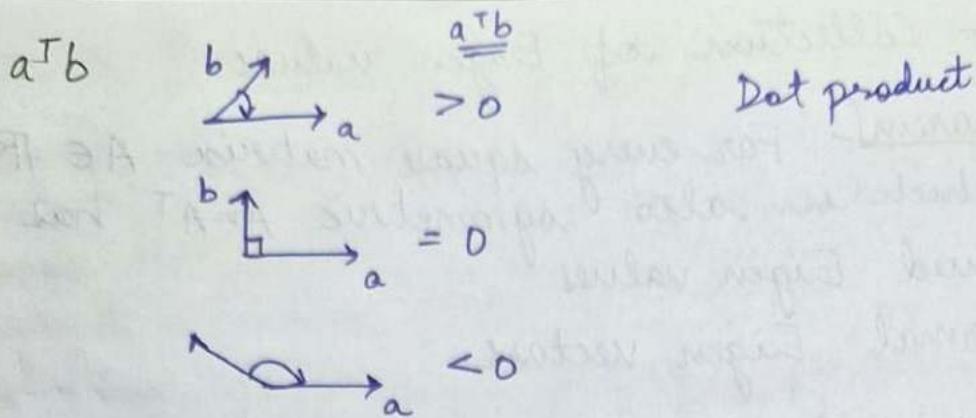
$x^T A x > 0 \quad \forall x \neq 0 \Rightarrow A$ is positive Definite (PD)

$x^T A x \geq 0 \quad \forall x \neq 0 \Rightarrow A$ is Positive Semi Definite (PSD)

$x^T A x < 0 \quad \forall x \neq 0 \Rightarrow A$ is Negative Definite (ND)

$x^T A x \leq 0 \quad \forall x \neq 0 \Rightarrow A$ is Negative Semi Definite (NSD)

$x^T A x >, \neq, <$ $\Rightarrow A$ is indefinite.



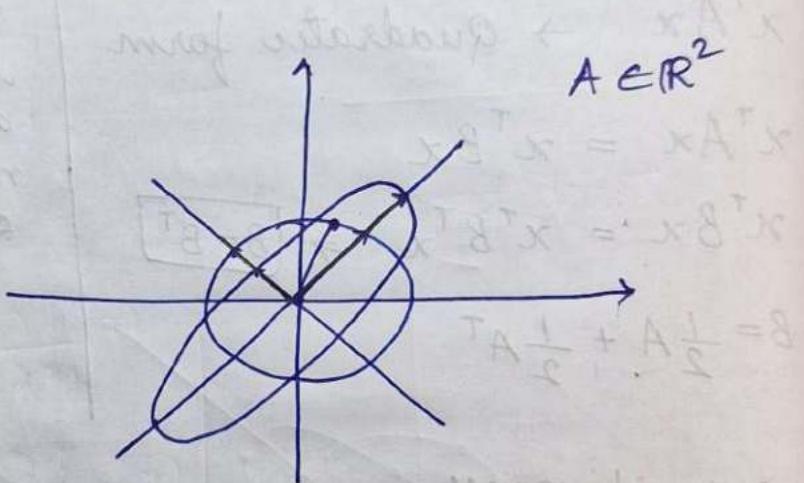
$$x^T A x \Rightarrow \underbrace{(x)}_{a}^T \underbrace{(Ax)}_{b}$$

- So

eigenvalues

PD	> 0
PSD	≥ 0
ND	< 0
NSD	≤ 0
ID	$>, <$

- The definitiveness of a matrix which is defined by the quadratic equation form has a one-to-one relation which the spectrum of the matrix.



- The eigen vectors kind of act as pivot
- It means vectors inside one quadrant remain in the same quadrant in the output space.

Decomposition

- Singular value decomposition
- Eigen value decomposition

- Cholesky decomposition

	A	Decomposition	
SVD	Any	$A = USV^T$ $\begin{bmatrix} \text{[I]} \\ \text{[I]} \end{bmatrix} \begin{bmatrix} \text{[..]} \\ \text{[..]} \end{bmatrix} \begin{bmatrix} \text{[=]} \\ \text{[=]} \end{bmatrix}$	$A(x) = U(S(V^T(x)))$ Nested pipeline
EVD	Square	$A = U D U^{-1}$ $\begin{bmatrix} \text{[I]} \\ \text{[I]} \end{bmatrix} \begin{bmatrix} \text{[..]} \\ \text{[..]} \end{bmatrix}$	$A(x) = U(D(U^{-1}(x)))$

A	Step 1	Step 2	Step 3	
SVD	V^T Rotation	Scaling along axis (Real valued)	U Rotation	<ul style="list-style-type: none"> • A arbitrary SVD
EVD	U^{-1} Rotation	Scaling along axis (Rotation) EV is complex	U inverse of step 1 (Rotation)	<ul style="list-style-type: none"> • A square <ul style="list-style-type: none"> ↳ can do both but the decomposed components may not be identical. • A square symm. ↳ $SVD = EVD$

- Eigen value decomposition defined only for square matrices while SVD is defined for any matrices.

$Ax = \lambda x \Rightarrow$ Eigen value is basically take a matrix A , multiply by x . we get an output vector x which is scaled by eigen value λ .
 $\Rightarrow A$ has to be square

MATRIX CALCULUS

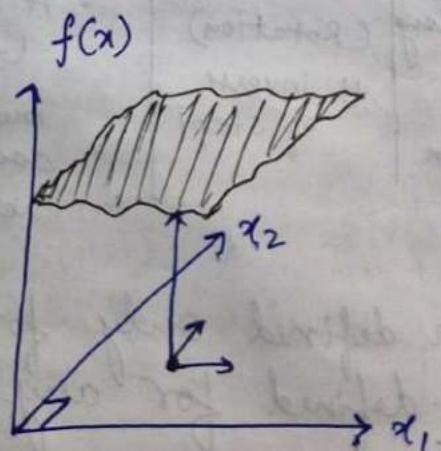
<u>function</u>	<u>Eg</u>	<u>Value</u>	<u>First deriv</u>	<u>Second deriv</u>
$f: \mathbb{R} \rightarrow \mathbb{R}$	x^2	\mathbb{R}	\mathbb{R}	\mathbb{R}
$f: \mathbb{R}^d \rightarrow \mathbb{R}$ (vec) (sc)	loss func.	\mathbb{R}	\mathbb{R}^d [Gradient]	$\mathbb{R}^{d \times d}$ [Hessian] S^d
$f: \mathbb{R}^d \rightarrow \mathbb{R}^p$ vec vec	NN layer	\mathbb{R}^p	$\mathbb{R}^{d \times p}$ [Jacobion]	$\mathbb{R}^{d \times p \times p}$ [High order tensor]

Gradient

It is the direction of the steepest ascent

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_d} \end{bmatrix}$$

$\nabla_x f(x_1, x_2, \dots, x_d)$



The gradient tells the direction in which we want to move to increase the value of the function the most.

Here we are taking vector as input and giving scalar as output.

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$A \rightarrow \text{Matrix}$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial a_1}(A) & \dots \\ \vdots & \ddots \\ \frac{\partial f}{\partial a_{mn}}(A) & \dots \end{bmatrix}$$

$\# f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

f is a function that takes Matrix A as input and gives scalar as output.

Hessian

$$\nabla_x^2 f(x) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\nabla_x^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_d \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_d \partial x_d} \end{bmatrix}$$

Example

$$\nabla_x b^T x = \begin{bmatrix} \vdots \\ \frac{\partial}{\partial x_i} (b^T x) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_i} (b_1 x_1 + b_2 x_2 - b_3 x_3 + \dots + b_d x_d) \end{bmatrix}$$

$$\nabla_x b^T x = \begin{bmatrix} \vdots \\ b_i \end{bmatrix} = b$$

Product Rule

$$\nabla_x \underbrace{x^T Ax}_{f(x) g(n)} \Rightarrow x^T \nabla_x (Ax) + Ax \nabla_x (x^T)$$

$$x^T A + A \mathbb{I} (I)$$

$$x^T (A + A^T) \quad \because A \text{ is symmetric}$$

$$\boxed{\nabla_x (x^T Ax) = 2Ax} \quad A = A^T$$

$$\boxed{\nabla_A \log |A| = A^{-1}}$$

$$\frac{d}{dx} \log(x) = x^{-1}$$

Review of Probability theory

- Set of all possible outcomes is called the sample space. $\Omega = \{HH, HT, TH, TT\}$
- An event A is subset of the sample space.
 $A \subseteq \Omega = \{HH, HT\}$
- The set of all possible subsets of sample space (F)

Probability Measure $P: F \rightarrow R$ (Here we are assigning prob. to events not outcomes)

$$P(A) \geq 0 \quad \forall A \in F$$

$$P(\Omega) = 1$$

- If A_1, A_2, \dots disjoint events then $P(\cup A_i) = \sum P(A_i)$

Conditional probability

Let $P(B) \neq 0$

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

$A \perp B$ only when $P(A \cap B) = P(A) \cdot P(B)$

$A \perp B$ only when $P(A|B) = \frac{P(A) P(B)}{P(B)} = P(A)$.

$\perp \rightarrow$ Independent of.

Random Variable (RV)

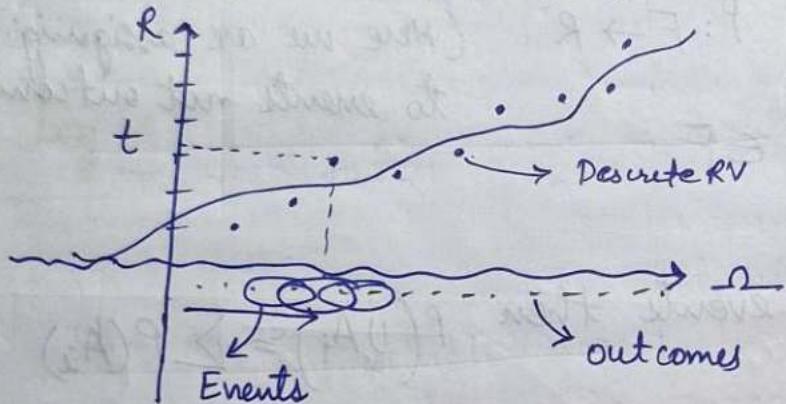
A random variable is a mapping from outcomes to real values.

$$\omega_0 = HHHHTTHTTTHTTTT$$

A RV is $X: \Omega \rightarrow \mathbb{R}$

of heads $X(\omega_0) = 5$

of tails $X(\omega_0) = 8$



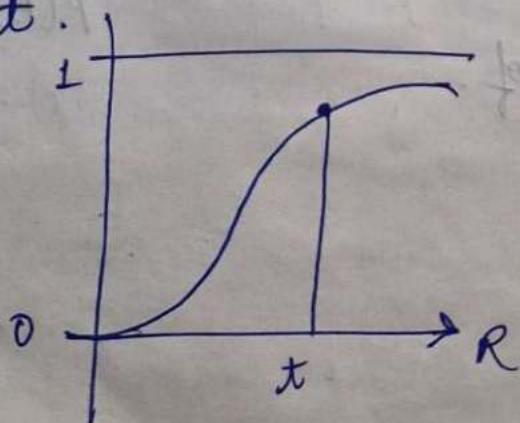
$$\text{Val}(X) = X(\Omega)$$

Cumulative Prob. distribution

$$F_X(x) = P(X < x)$$

$$P[\{\omega : X(\omega) < x\}]$$

Basically it means to get the set of all ω , such that $X(\omega) < x$ (lies below x). Get A set of those events and then measure of Prob. on that event.

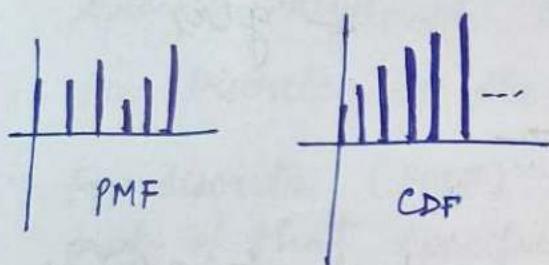


Discrete vs Continuous RV

Cumulative distribution func.
(CDF)

- ① Probability mass function
Discrete (PMF)
- ② Probability density function
Continuous. (PDF)

PMF



Discrete RV

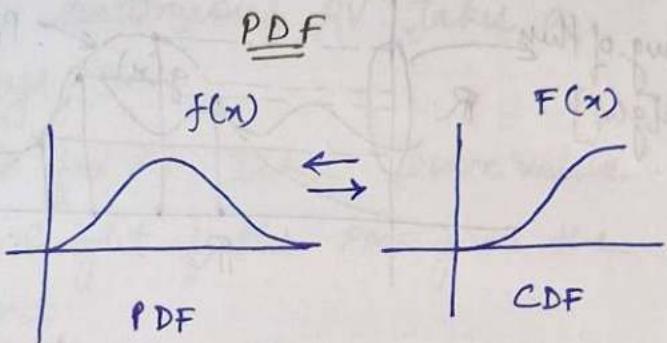
Val(X) countable

$$P(X=k) = P(\{w \mid X(w) = k\})$$

PMF

$$p(x) : \text{Val}(X) \rightarrow [0, 1]$$

$$\sum_{x \in \text{Val}(x)} p_x(x) = 1$$



$$\frac{d}{dx} F(x) = f(x)$$

or

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

Continuous RV

Val(X) uncountable

$$P(a \leq X \leq b) = P(\{w \mid a \leq X(w) \leq b\})$$

PDF

$$f_x : \mathbb{R} \rightarrow \mathbb{R}$$

$$f_x(x) = \frac{d}{dx} F_x(x)$$

$$f_x(x) \neq P(X=x)$$

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

$$P(x \leq X \leq x + dn)$$

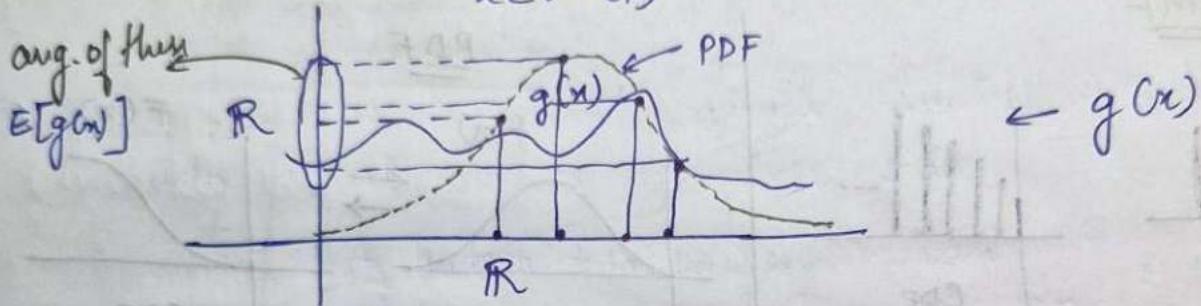
Expected value and Variance

$$g: \mathbb{R} \rightarrow \mathbb{R}$$

Exp value

Let X be a discrete RV with PMF p_X

$$E[g(X)] = \sum_{x \in \text{Val}(X)} g(x) p_X(x). \quad (\text{discrete})$$



$E[g(x)]$: If the input that we feed to $g(x)$ is a RV which means the actual input fed is random the expectation of $g(x)$ is sum over every possible x that can be fed as input to g . Multiplied by the prob. that x can be fed.

For continuous just take integral.

- Basically what is the avg. value / Exp. value of $g(x)$ we get after feeding the input.

for continuous RV with PDF f_x

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

$$= \int g(x) f_x(x) dx$$

$$\approx \int g(x) f_x(x) dx$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N g(x^{(i)}) \rightarrow \int_{-\infty}^{\infty} g(x) p(x) dx.$$

Law of Large Numbers

Monte Carlo Estimate

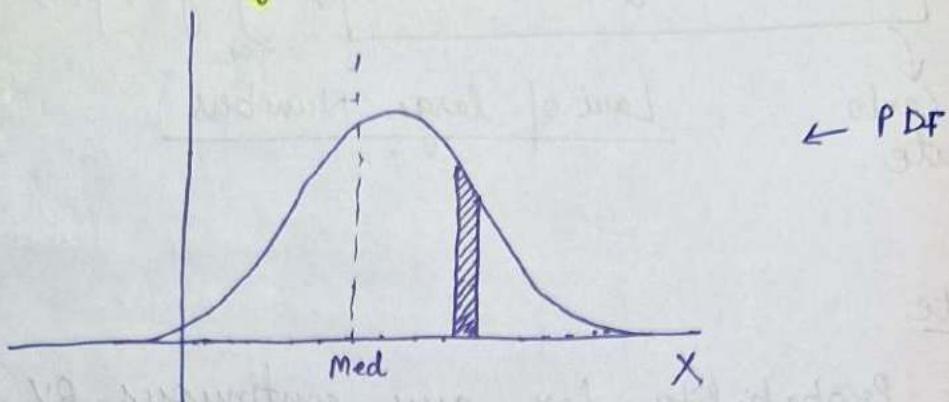
Note

- The probability for any continuous RV takes a specific value is always 0.
- For Discrete RV the Prob. for RV takes some value.
- For discrete (PMF) the height of the PMF gives the prob. of that specific input.
- But PDF i.e prob. density function is different as the height is not the prob. we can only find the probability in the intervals and not at specific point.
- Expectation is basically what this value RV takes the value on average.
- Variance is the measure of how spread apart is the distribution in general or how spread apart is this RV around its mean.

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

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

Lecture 3 - Probability and Statistics



- Median - It is the point which divides the mass into equal masses on either side
- For symm. distribution $E(X) = \text{Median}$
 $\text{expectation} = \text{Median}$

Two Random Variables

Binariate CDF

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

Binariate PMF

$$P_{XY}(x, y) = P(X=x, Y=y)$$

Joint Probability

- Joint Prob. Mass function -

$$\text{Joint(PMF)} \quad f: \mathbb{R}^n \rightarrow [0, 1]$$

$$f(x_1, \dots, x_n) = P(X_1=x_1, \dots, X_n=x_n)$$

for $n=2$

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

$$\text{CDF : } F: \mathbb{R}^n \rightarrow [0, 1]$$

$$F(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n)$$

- Marginal Probability mass function

$$X, Y \quad f(x, y) = P(X=x, Y=y)$$

$$\{X=x, Y=y\}$$

$$\{X=x\} \rightarrow P(X=x) = f_X(x)$$

(Marginal Prob.)

- # X and Y are two discrete RV with JP mass function $f(x, y)$. The prob. mass func. of ~~is~~ only X is $f_X(x)$ and is called marginal prob. mass func. of X . It basically doesn't consider what is Y but find prob for X only.
- # Similarly for Y : $f_Y(y) \rightarrow$ MPMF of Y .

$$\text{Range}(X) = \{x_1, \dots, x_n\} \text{ and } \text{Range}(Y) = \{y_1, \dots, y_m\}$$

Marg Prob : $f_X(x_i) = \sum_{j=1}^m f(x_i, y_j) \text{ for } i=1, \dots, n$

$$f_Y(y_j) = \sum_{i=1}^n f(x_i, y_j) \text{ for } j=1, \dots, m$$

- Joint Probability density function

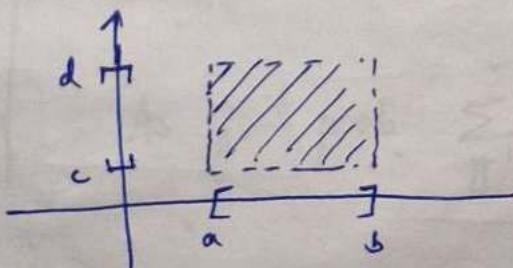
$$P(a \leq X \leq b) = \int_a^b f(x) dx \quad P(X=x) = 0$$

X, Y

$$P(X=x, Y=y) = 0$$

$$P(a \leq X \leq b, c \leq Y \leq d)$$

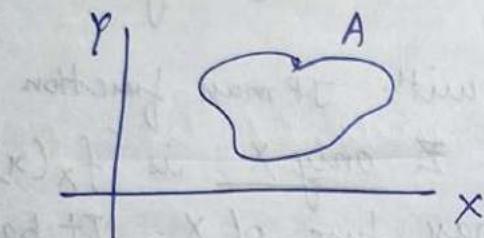
$$\{a \leq X \leq b, c \leq Y \leq d\}$$



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

$$\Rightarrow P(a \leq x \leq b, c \leq y \leq d) = \int_a^b \int_c^d f(x, y) dy dx$$

Generalization



Event: $\{(x, y) \in A\}$

$$P((x, y) \in A) = \int_A f(u, y) du dy$$

or

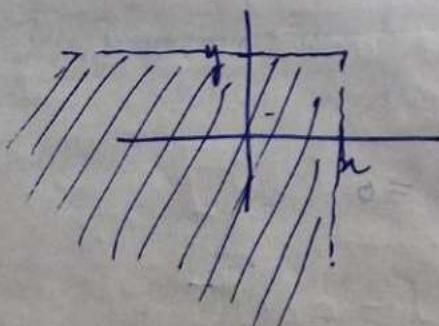
$$\iint_A f(x, y) dx dy$$

for all $A \in \mathbb{R}^2$

CDF

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(u) du \rightarrow f(x) = \frac{dF}{dx}$$

CDF: $F(x \leq x, y \leq y) = \int_{-\infty}^x \int_{-\infty}^y f(u, v) du dv$



$$\frac{\partial^2 F}{\partial x \partial y} = f(x, y)$$

Joint CDF

Joint PDF



- Marginal Probability density function

$$f(x, y) \quad \{a \leq x \leq b, c \leq y \leq d\} \quad \iint_a^b f(x, y) dy dx$$

$$\{a \leq x \leq b\} \quad \int_a^b f_x(x) dx$$

Now marginal of X and Y

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

$$f_y(y) = \int_{-\infty}^{\infty} f(x, y) dx$$

Bayes's Theorem

$$p(x, y) = p(x) \cdot p(y|x) \quad \text{--- ①}$$

Joint Marginal Conditional

$$p(x, y) = p(y) \cdot p(x|y) \quad \text{--- ②}$$

Equate ① and ②

$$p(y|x) = \frac{p(x|y) \cdot p(y)}{p(x)}$$

→ Given Cond. prob. $p(x|y)$
 → we find reverse for
 cond. prob. $p(y|x)$.
 using Bayes's theorem

$$p(y|x) = \frac{p(x|y) \cdot p(y)}{\sum_{y'} p(y') \cdot p(x|y')}$$

As joint prob. $\sum_{y'} p(x, y') = p(x)$

for continuous X and Y

$$f(y|x) = \frac{f(x|y) f(y)}{f(x)}$$

where: $f(x) = \int_{y' \in \text{val}(y)} f(x|y') f(y') dy'$

Independence

$P(A \cap B) = P(A) \cdot P(B)$

$$\underbrace{F(a,b)}_{\text{joint CDF}} = \underbrace{F_x(a)}_{\text{marginal CDFs}} \cdot \underbrace{F_y(b)}_{\text{marginal CDFs}}$$

$$f(x,y) = f_x(x) \cdot f_y(y)$$

Expectation $\underline{\text{eg}}$

$X, Y \rightarrow$ continuous RV

$g, R^2 \rightarrow R$

$$E(g(x,y)) = \int_{x \in \text{val}(x)} \int_{y \in \text{val}(y)} g(x,y) f_{XY}(x,y) dx dy$$

$$\begin{aligned} & g(x,y) = 3x \\ & f_{XY} = 4xy \\ & 0 < x < 1 \\ & 0 < y < 1 \end{aligned}$$

$$E = \int_0^1 \int_0^1 3x \cdot 4xy dy dx$$

$$12 \int_0^1 \int_0^1 x^2 y dy dx$$

$$12 \int_0^1 \int_0^1 \frac{x^3}{3} y dy dx \Rightarrow \int_0^1 4y dy = 2y^2 \Big|_0^1 \Rightarrow 2 \quad \text{②}$$

Covariance of two Random variable X and Y

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

$$\boxed{\text{Cov}(X, Y) = E[XY] - E[X]E[Y]}$$

$$\text{Cov}(X, X) = \text{Var}(X)$$

If X and Y independent then

$$E(XY) = E(X) \cdot E(Y) \implies \text{Cov}[X, Y] = 0$$

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

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$$
simple proof

Multivariate Gaussian (Normal) distribution

$$x \in \mathbb{R}^n$$

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

Σ : Covariance matrix
and it is PSD (Positive semi definite). Its Full Rank

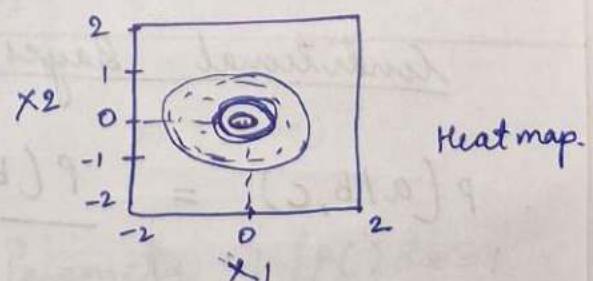
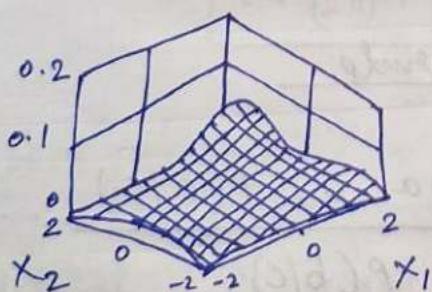
$|\Sigma|$: Determinant

The exp form is of the quadratic form: $V^T A V$

Ex

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0 \ 0]^T$$



$\mu[0 \ 0]^T \rightarrow$ mean the center of distribution is at 0,0

$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow$ Covariance b/w RV, $\text{Cov}=0$ mean no strong rel between x_1, x_2 (correlation)

$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow$ Variance of the marginal (Variation)

Conditional expectation

$E[X]$ = constant

$E[X|Y]$ = random variable

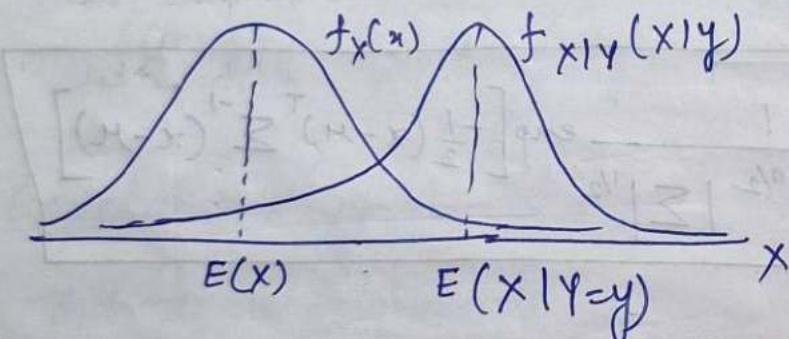
$E[X|Y=y]$ = function of y

X, Y

$Y=y$ (partial info).

$$f_X(x) \rightarrow f_{X|Y}(x|y) = \frac{f(x, y)}{f_Y(y)} \quad \begin{cases} \text{- Pdf change} \\ \text{- Exp will also change} \end{cases}$$

$$E(X) = \int_{-\infty}^{\infty} x \cdot f_X(x) dy \rightarrow E(X|Y=y) = \int_{-\infty}^{\infty} x \cdot f_{X|Y}(x|y) dx$$



Hence

Law of total expectation

$$E[X] = E[E[X|Y]]$$

Conditional Bayes rule

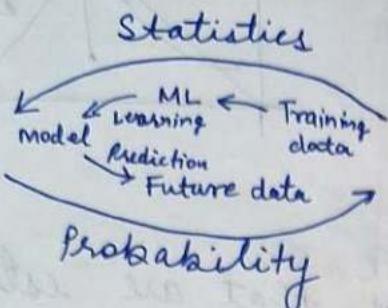
$$P(a|b,c) = \frac{P(b|a,c) P(a|c)}{P(b|c)}$$

Statistics in ML

Parameters

$$\mu, \Sigma$$

- Method of Least squares
- Maximum Likelihood estimation



Observations / Data

$$x \in \mathbb{R}^n$$

Training data
 (x, y)

Maximum likelihood estimation (MLE)

$$x^{(1)}, \dots, x^{(n)}$$

$x^{(i)} \rightarrow i^{\text{th}}$ example of $x \in \mathbb{R}^d$

Gaussian Data

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

{ sampled independently & identically distributed }

Probability density

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right]$$

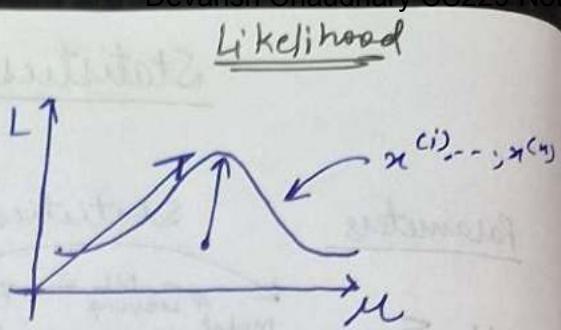
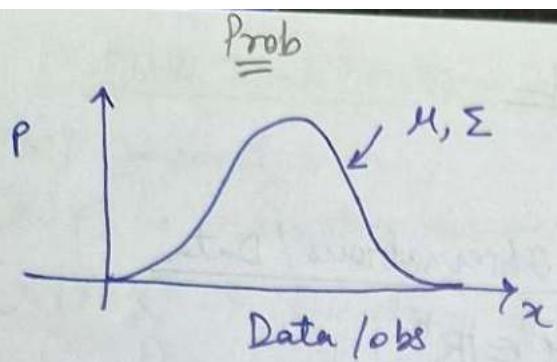
$$P(x^{(1)}, \dots, x^{(n)}) = \prod_{i=1}^n P(x^{(i)})$$

$$P(x^{(1)}, \dots, x^{(n)}; \mu, \Sigma) = \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)\right]$$

Likelihood: $L(\mu, \Sigma; x^{(1)}, \dots, x^{(n)}) = \dots$ (same)

$P \rightarrow$ Probability of data given Parameters $\left| \int P(x) dx = 1 \right.$

$L \rightarrow$ Likelihood of parameters given data $\left| \int L(\mu, \Sigma) = \text{nothing} \right.$



- We use subscript $\hat{\theta}$ over things that are estimated as output of some estimation process

$$L(\theta; x) = \prod_{i=1}^n L(\theta; x^{(i)})$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \prod_{i=1}^n L(\theta; x^{(i)})$$

$$= \arg \max_{\theta} \log \prod_{i=1}^n L(\theta; x^{(i)})$$

$$= \arg \max_{\theta} \sum_{i=1}^n l(\theta; x^{(i)}) \quad \left\{ l(\theta) = \log L(\theta) \right\}$$

For Gaussian

$$\hat{\mu}, \hat{\Sigma} = \arg \max_{\mu, \Sigma} \sum_{i=1}^n \log \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

$$= \arg \max_{\mu, \Sigma} \left[\sum_{i=1}^n \left(-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu) \right) \right]$$

To find arg max take diff.

$$\nabla_{\mu} \sum_{i=1}^n \kappa - \frac{1}{2} \log |\Sigma| \equiv -\frac{1}{2} (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}) \Sigma^{-1} (\boldsymbol{x}^{(i)} - \boldsymbol{\mu})$$

$$\nabla_{\mu} \sum_{i=1}^n \frac{-1}{2} \left[\boldsymbol{x}^{(i)\top} \Sigma^{-1} \boldsymbol{x}^{(i)} - \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{x}^{(i)} + \boldsymbol{\mu}^\top \Sigma \boldsymbol{\mu} \right]$$

$$\nabla_{\mu} \sum_{i=1}^n -\frac{1}{2} \left[-2(\Sigma^{-1} \boldsymbol{x}^{(i)})^\top \boldsymbol{\mu} + \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu} \right]$$

$$\sum_{i=1}^n (\Sigma^{-1} \boldsymbol{x}^{(i)} - \Sigma^{-1} \boldsymbol{\mu}) = 0$$

$$n \Sigma^{-1} \boldsymbol{\mu} = \sum_{i=1}^n \Sigma^{-1} \boldsymbol{x}^{(i)}$$

$$\Sigma^{-1} \boldsymbol{\mu} = \Sigma^{-1} \left[\frac{1}{n} \sum_{i=1}^n \boldsymbol{x}^{(i)} \right]$$

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}^{(i)}$$

$$\Sigma \equiv$$

$$\Sigma = \nabla_{\mu}^{-1}, \quad \nabla_{\mu} \ell = 0$$

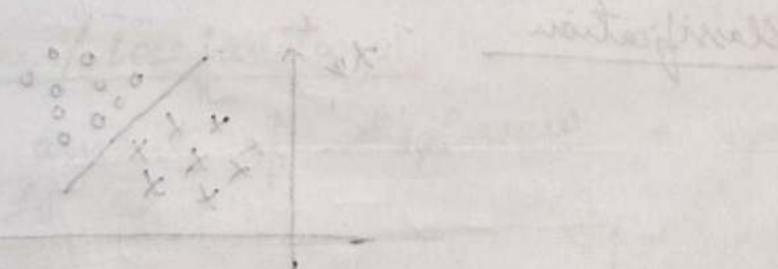
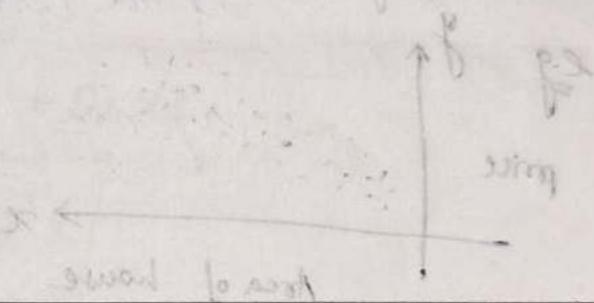
$$\Leftrightarrow \nabla_{\mu} \ell = 0$$

$$\nabla_{\mu} \sum_{i=1}^n \frac{1}{2} \log |\Sigma| - \frac{1}{2} ((\boldsymbol{x}^{(i)} - \boldsymbol{\mu})^\top \Sigma (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}))$$

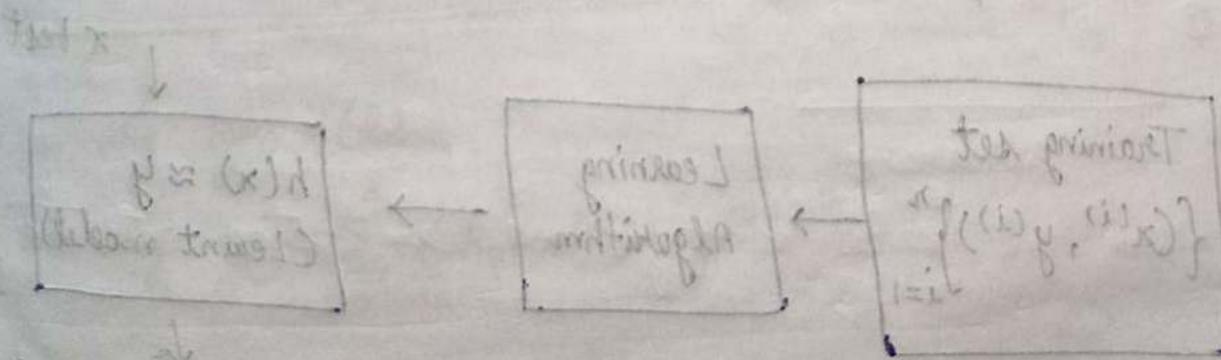
$$\frac{1}{2} \left[n\Sigma - \sum_{i=1}^n (\boldsymbol{x}^{(i)} - \boldsymbol{\mu})(\boldsymbol{x}^{(i)} - \boldsymbol{\mu})^\top \right] = 0$$

$$\Sigma^{-1} = \frac{1}{n} \sum_{i=1}^n (\boldsymbol{x}^{(i)} - \boldsymbol{\mu})(\boldsymbol{x}^{(i)} - \boldsymbol{\mu})^\top$$

$$\text{similar } \approx E[(x - E[x])^2]$$



using PCA for dimensionality reduction



Lecture 4 - Linear Regression

Supervised learning

$X \rightarrow Y$ mapping
 (Input) (Output)

n - No. of (x, y) examples in training set

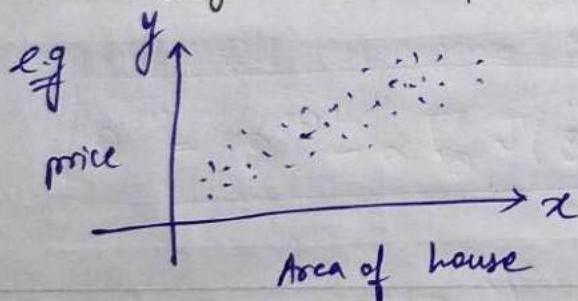
d - $x \in \mathbb{R}^d$ dimension of input

$x^{(i)}$ - i th example input

$y^{(i)}$ - i th example output

$(x^{(i)}, y^{(i)})$ - i th example

Regression (prediction)

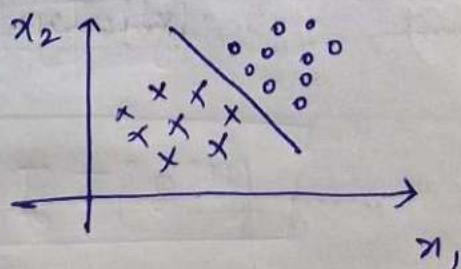


Learn

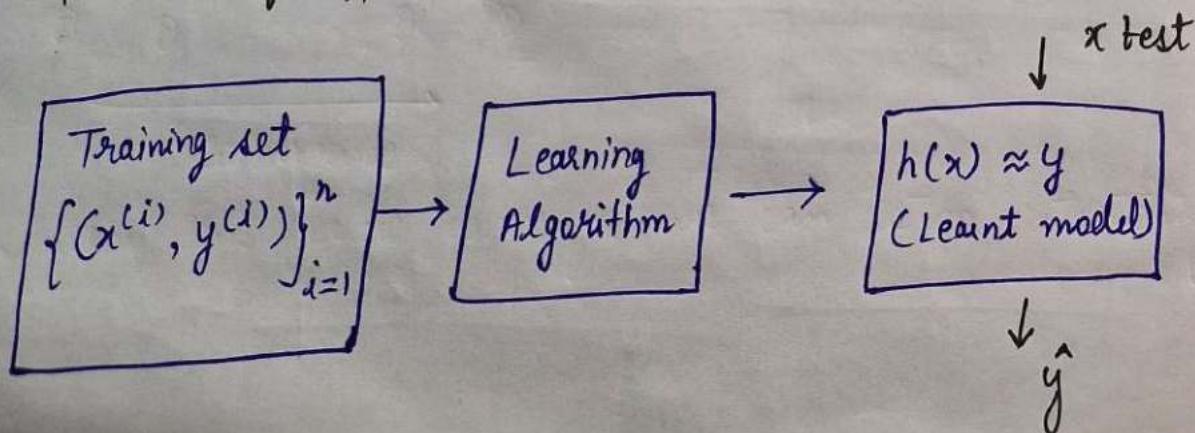
Hypothesis

$$h(x) \approx y$$

Classification



Flowchart of Approach



Linear Regression

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

The hypothesis that we want to learn has some parameter θ . The goal of linear regression is to learn a suitable set of parameters θ , that makes the y value as close as possible to $h_\theta(x)$.

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_d x_d$$

$$\theta \in \mathbb{R}^{d+1}$$

$$h_\theta(x) = \left(\sum_{i=1}^d \theta_i x_i \right) + \theta_0$$

$x_0 = 1$

- intercept term

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

$$h_\theta(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x$$

Cost function / Loss function

It is simply the amount of displeasure a specific hypothesis causes to us.

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

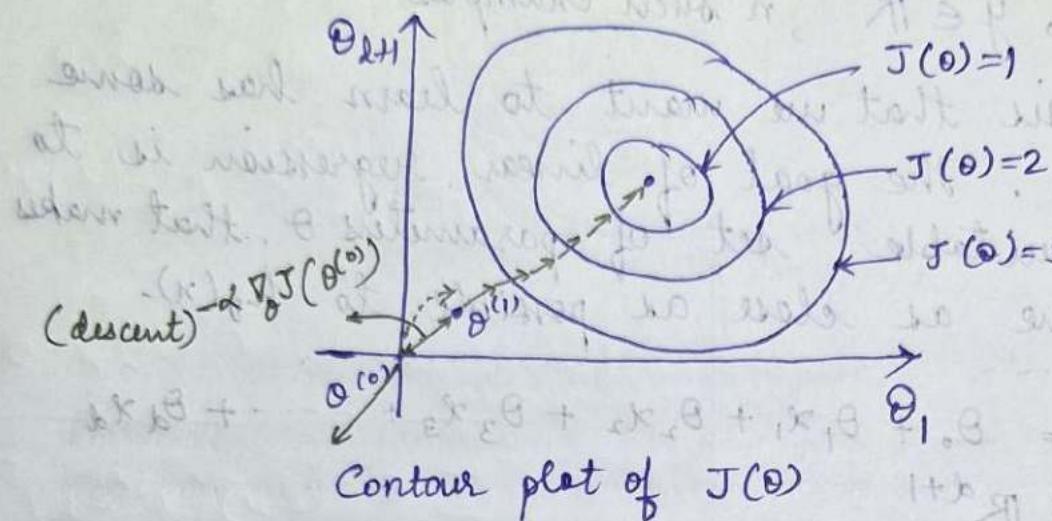
squared error

$$\hat{\theta} = \arg \min_{\theta} J(\theta)$$

$$= \arg \min_{\theta} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$

line with marks

Algorithm : Gradient Descent



$\theta^{(0)}$ = Initialization

$$\theta_j^{(1)} = \theta_j^{(0)} - \alpha \frac{\partial}{\partial \theta_j} J(\theta^{(0)})$$

Repeat till convergence:

$$\theta^{(1)} = \theta^{(0)} - \alpha \nabla_{\theta} J(\theta^{(0)})$$

α - Learning rate

$$J(\theta^{(1)}) < J(\theta^{(0)})$$

(hopefully)

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta^{(t)})$$

$$\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(t)}, \theta^{(t+1)}$$

Convergence

- when this series converges, we can stop iterating

check for convergence

$$\|\theta^{(t)} - \theta^{(t-1)}\|, \quad \|\mathcal{J}(\theta^{(t)}) - \mathcal{J}(\theta^{(t-1)})\|$$

If the norm of θ is too small

If the loss has going down

$$\|\nabla_{\theta} \mathcal{J}(\theta^{(t)})\|$$

If the norm of the gradient is too small.

Gradient descent on linear regression.

$$\theta^{(0)} = \text{---} \quad (\text{init})$$

repeat until convergence

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} \mathcal{J}(\theta^{(t)})$$

$$= \theta^{(t)} - \alpha \nabla_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right]$$

$$= \theta^{(t)} - \alpha \nabla_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 \right]$$

$$= \theta^{(t)} - \alpha \left[\frac{1}{2} \sum_{i=1}^n 2(\theta^T x^{(i)} - y^{(i)}) x^{(i)} \right]$$

$$\boxed{\theta^{(t+1)} = \theta^{(t)} - \alpha \left[\sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)}) x^{(i)} \right]}$$

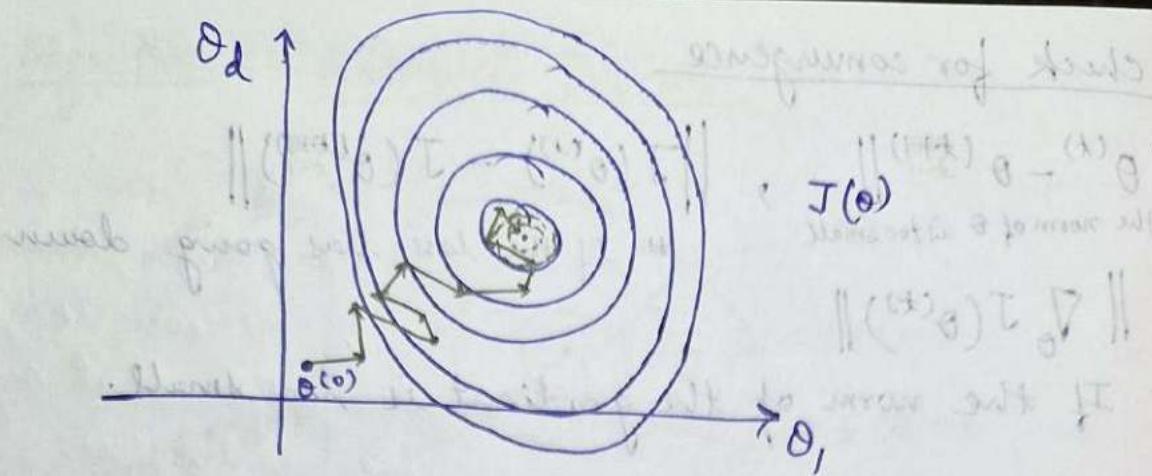
Stochastic Gradient Descent (SGD)

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} \tilde{\mathcal{J}}(\theta)$$

where,

$$\tilde{\mathcal{J}}(\theta) = \frac{1}{2} (\theta^T x^{(k)} - y^{(k)})^2$$

: K uniformly
at random
sampled from
training set.



- # In GD we perform descent in a proper manner i.e. approaching towards minima with each step using the loss function which is considering every example. Meanwhile in SGD we randomly pick an example and follow gradient descent using that. That is why it is noisy but saves a lot of expense/memory as we are not considering all examples.
- # Eventually using SGD we reach a region near minima where it is not a lot noisy.
- # Computational cost of SGD << GD. Hence we can trade off No. of steps to reach minima with computation cost.
- # GD and SGD are numerical algorithms (iterative algs) gives a value of θ as output.

cost function - $J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2$

Design Matrix $\xleftarrow{\text{d+1 dimension}}$

$$X = \begin{bmatrix} -x^{(1)} \\ \vdots \\ -x^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times d+1} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_{d+1} \end{bmatrix} \in \mathbb{R}^{d+1 \times 1}$$

$$X\theta - Y$$

$$\mathbb{R}^{n \times (d+1)} \rightarrow \mathbb{R}^{(d+1)}$$

$$\mathbb{R}^n - \mathbb{R}^n \rightarrow \boxed{\mathbb{R}^n}$$

$$X\theta - Y = \begin{bmatrix} x^{(1)\top}\theta - y^{(1)} \\ x^{(2)\top}\theta - y^{(2)} \\ \vdots \\ x^{(n)\top}\theta - y^{(n)} \end{bmatrix}$$

$$\Rightarrow J(\theta) = \frac{1}{2} \sum_{i=1}^n (x^{(i)\top}\theta - y^{(i)})^2 = \boxed{\frac{1}{2} (X\theta - Y)^T (X\theta - Y)}$$

$$\Rightarrow \nabla_{\theta} J(\theta) = 0 \quad \text{vector form}$$

$$\nabla_{\theta} \frac{1}{2} (X\theta - Y)^T (X\theta - Y)$$

$$\nabla_{\theta} \frac{1}{2} [(X\theta)^T X\theta - (X\theta)^T Y - Y^T (X\theta) + Y^T Y]$$

$$\nabla_{\theta} \frac{1}{2} \left[\underbrace{\theta^T (X^T X) \theta}_{\text{Quadratic form}} - \underbrace{2\theta^T (X^T Y)}_{\text{Vec} \times \text{Vec} = \text{scalar}} + \underbrace{Y^T Y}_{\text{Scalar}} \right]$$

$$\nabla_{\theta} \frac{1}{2} [-] = \frac{1}{2} [2(X^T X)\theta - 2X^T Y] = 0$$

$$(X^T X)\theta = X^T Y$$

Normal equation

Exact
solution
for cal. θ / min cost
time.

$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

(only if $X^T X$ invertible)

Probabilistic Interpretation

In this we make an assumption that $y^{(i)}$ is generated using this approach.

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

↓
 Unknown
Vector
that we
are
calculating

↓
 Eg.
Features
(price of
house)

Random Gaussian noise
distributed acc. to guass
variable.
 $\epsilon^{(i)} \sim N(0, \sigma^2)$

$$\epsilon^{(i)} = y^{(i)} - \theta^T x^{(i)} \sim N(0, \sigma^2)$$

Property of Gaussian noise

$$\hookrightarrow y^{(i)} - \theta^T x^{(i)} \sim N(0, \sigma^2)$$

$$y^{(i)} \sim N(\theta^T x^{(i)}, \sigma^2)$$

$$\Rightarrow p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi \sigma^2}} \exp\left[-\frac{1}{2} \frac{(y^{(i)} - \theta^T x^{(i)})^2}{\sigma^2}\right]$$

Data
 x, y

Parameters

$$\mu, \sigma^2$$

$$\theta, \sigma^2$$

$$\log L(\theta) = \log \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

$$Y^T X^T (X^T X)^{-1} \hat{\theta}$$

$$l(\theta) = \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(y^{(i)} - \theta^T x^{(i)})^2}{\sigma^2} \right] \right]$$

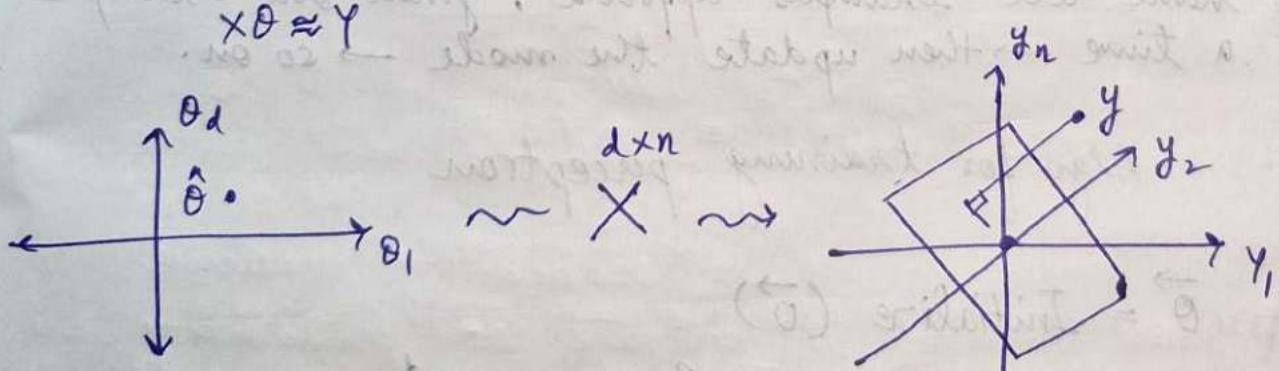
$$= K - \frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2$$

$$l(\theta) = K - \frac{1}{\sigma^2} \left[\sum_{i=1}^n \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2 \right]$$

$$\arg \max_{\theta} l(\theta) = \arg \min_{\theta} J(\theta).$$

- * The log likelihood function by making a probabilistic assumption about the noise is gonna give us the -ve of the scaled version of the original cost function
- * which means by performing max. likelihood, we are minimizing the squared error.

Another view : Projection Interpretation



$$X(X^T X)^{-1} X^T \quad (\text{Projection matrix})$$

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

- * We want $X\theta \approx \vec{y}$ but it's never the case, so we use projection

$$X\theta \approx \vec{y}$$

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

Lecture 5 - Perceptron and Logistic Regression

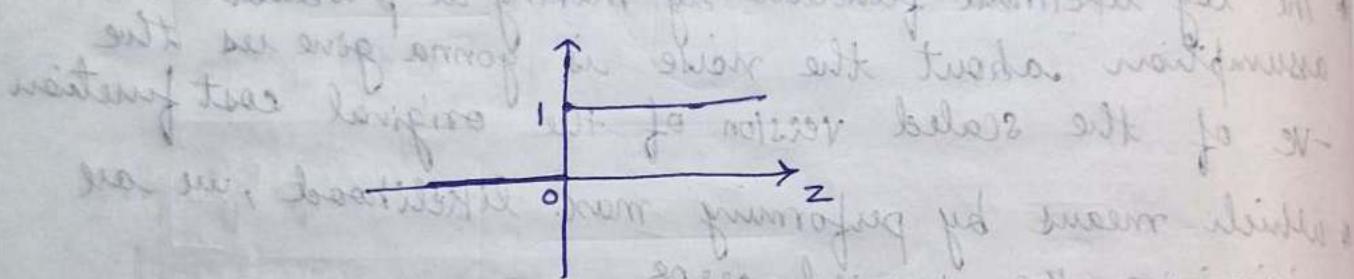
Perceptron algorithm

Simple classification algorithm

$$\vec{x}^{(i)} \in \mathbb{R}^d, y^{(i)} \in \{0, 1\}$$

Hypothesis $h_{\theta}(\vec{x}) = g(\vec{\theta}^T \vec{x})$

$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$



- # It is also called a streaming algorithm because it is designed to work in a setting where you are encountering examples one after another. We do not have all examples upfront, given one example at a time \rightarrow then update the model \rightarrow so on.

Algo for training perceptron

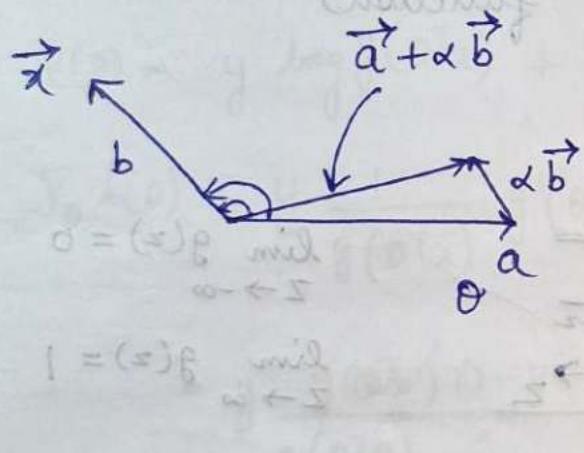
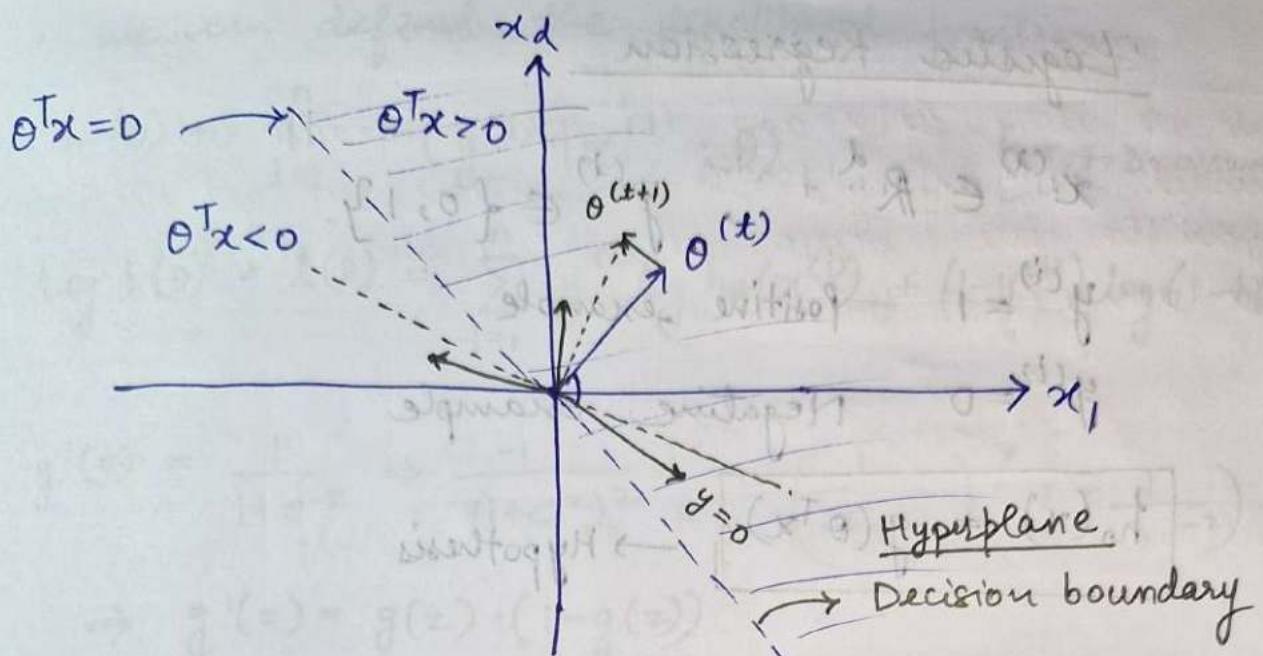
$$\vec{\theta} = \text{Initialize } (\vec{\theta})$$

For i in $1, 2, \dots, R$

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

(update rule)

- * This looks similar to update rule of lin. regression but is actually diff. since hypothesis is diff. It is $g(\vec{\theta}^T \vec{x})$ and not $\vec{\theta}^T \vec{x}$.



$a^T b < 0$ at present
 \Rightarrow so how do we make it positive

$$a^T(b + \alpha a) = a^T b + \alpha a^T a$$

always +ve

$$\Rightarrow a^T b + \alpha a^T a \geq a^T b$$

- Similarly, we can do for -ve.
- If $y^{(i)} = h_0(x^{(i)})$ means classification is done correctly, therefore no need to update θ .
- If you are given a dataset where you are given that the classes are linearly separable then this perceptron algorithm will find the separating hyperplane.

$$\begin{bmatrix} 1 & [x]_{art=1} \end{bmatrix} \cdot \begin{bmatrix} [x]_{art} \end{bmatrix} = (\theta : x) p$$

Logistic Regression

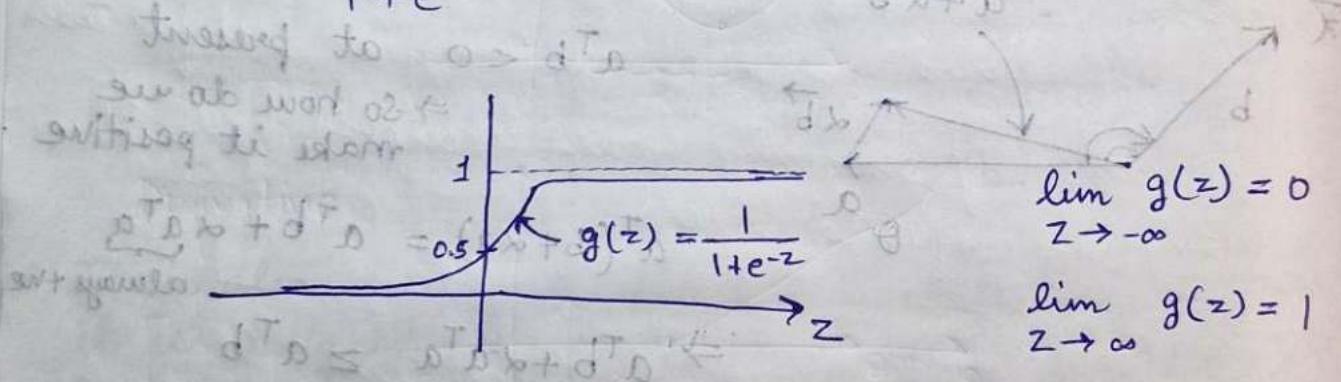
$$x^{(i)} \in \mathbb{R}^d \quad , \quad y^{(i)} \in \{0, 1\}$$

$y^{(i)} = 1$ Positive example

$y^{(i)} = 0$ Negative example

$$h_{\theta}(x) = g(\theta^T x) \rightarrow \text{Hypothesis}$$

$$\Rightarrow g(z) = \frac{1}{1 + e^{-z}} \quad (\text{logistic function})$$



$$P(y^{(i)} = 1 | x^{(i)}; \theta) = h_{\theta}(x)$$

$$P(y^{(i)} = 0 | x^{(i)}; \theta) = 1 - h_{\theta}(x)$$

- # Logistic regression is a probability machine. For any given example it outputs a probability and then it is upto us to choose a threshold and then compare the probability that is outputted to that and make decision.

$$P(y|x; \theta) = [h_{\theta}(x)]^y [1 - h_{\theta}(x)]^{1-y}$$

Similar to Bernoulli distribution

we can define the likelihood function.

$$L(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) \quad [\text{I.I.D assumption}]$$

$$\log L(\theta) = l(\theta) = \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)}))$$

$$g'(z) = \frac{1}{1+e^{-z}} \Rightarrow \frac{-1}{(1+e^{-z})^2} \cdot e^{-z} = \left(\frac{1}{1+e^{-z}}\right) \left(1 - \frac{1}{1+e^{-z}}\right)$$

$$\Rightarrow g'(z) = g(z) \cdot (1-g(z))$$

$$l(\theta) = y \log g(\theta^T x) + (1-y) \log (1-g(\theta^T x))$$

$$\nabla_\theta l(\theta) = y \cdot \frac{1}{g(\theta^T x)} \cdot g'(\theta^T x) \cdot x + (1-y) \frac{1}{1-g(\theta^T x)} \cdot (-1) g'(\theta^T x) \cdot x$$

$$= y \cdot \frac{g(\theta^T x) (1-g(\theta^T x)) \cdot x}{g(\theta^T x)} + (1-y) \frac{(-1)(g(\theta^T x)) \cdot (1-g(\theta^T x)) \cdot x}{(1-g(\theta^T x))}$$

$$\nabla_\theta l(\theta) = y(1-g(\theta^T x)) \cdot x + (1-y)(-1)g(\theta^T x) \cdot x$$

$$= x [y(1-g(\theta^T x)) - (1-y)g(\theta^T x)]$$

$$= [y - yg(\theta^T x) - g(\theta^T x) + yg(\theta^T x)] \cdot x$$

$$\nabla_\theta l(\theta) = [y - h_\theta(x)] \cdot x$$

$$\boxed{\theta = \theta + \alpha [y - h_\theta(x)] x}$$

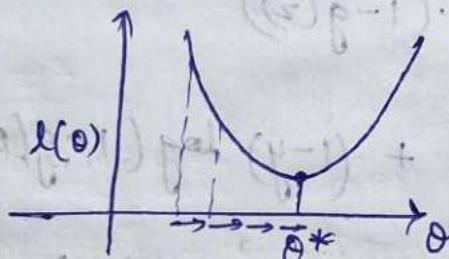
Note

- We are using $(y|x)$ and not (y,x) because when this model is deployed to a system we get the input from the user hence we use y given x . But we will be using (y,x) too later.

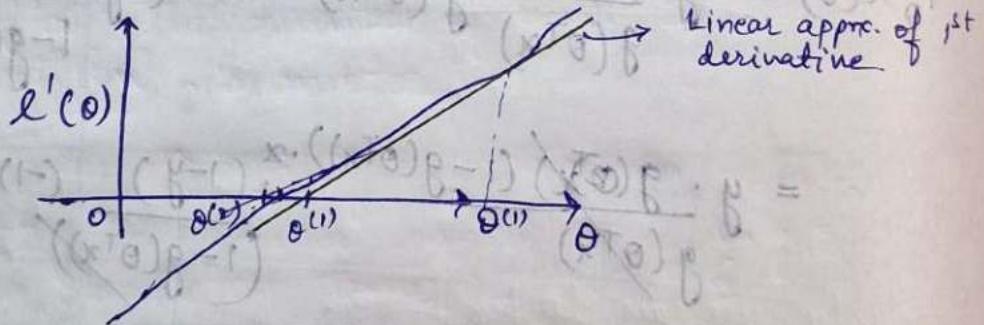
Newton's Method

- Alternative to gradient descent (GD).

Loss function



Derivative

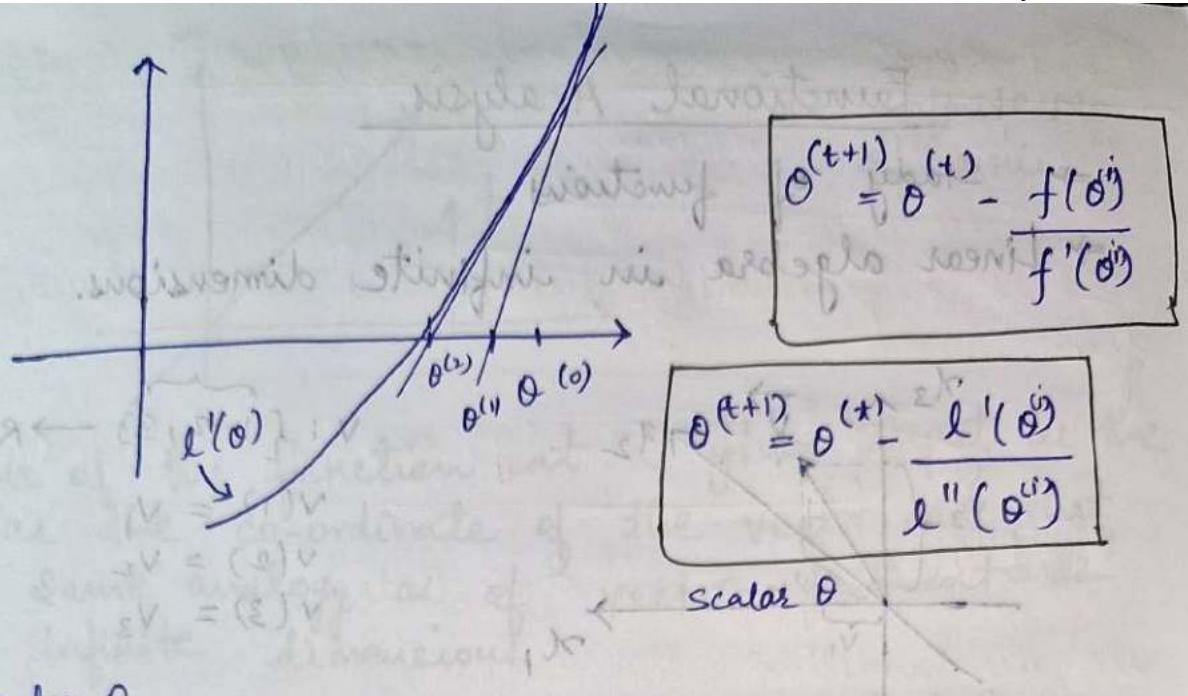


- Newton's Method is a way in which for any given function $f(x)$, it finds the point where $f(x)=0$.

It is basically a root finding method. Finds the input where $f(x)=0$.

$$\theta^{(0)}, \theta^{(1)}, \theta^{(2)}, \dots, \theta^*$$

The rate at which Newton's method converges is much faster than the regular gradient descent.

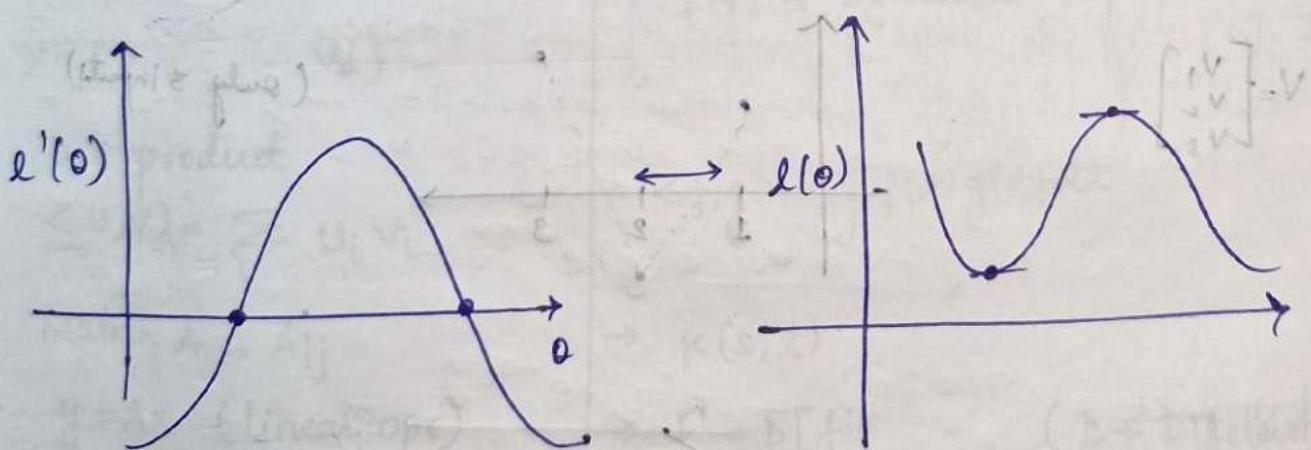


For vector θ

$$\theta^{(t+1)} = \theta^{(t)} - \alpha H^{-1} \nabla_{\theta} l(\theta^{(t)})$$

H = Hessian of loss $l(\theta)$

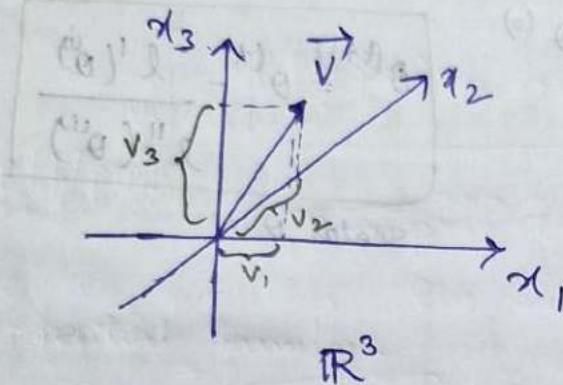
Newton Raphson method



- Newton's method will take you to the nearest stationary point of that function.
- It works well with a convex function than concave function

Functional Analysis

- ⇒ Study of functions
- ⇒ Linear algebra in infinite dimensions.



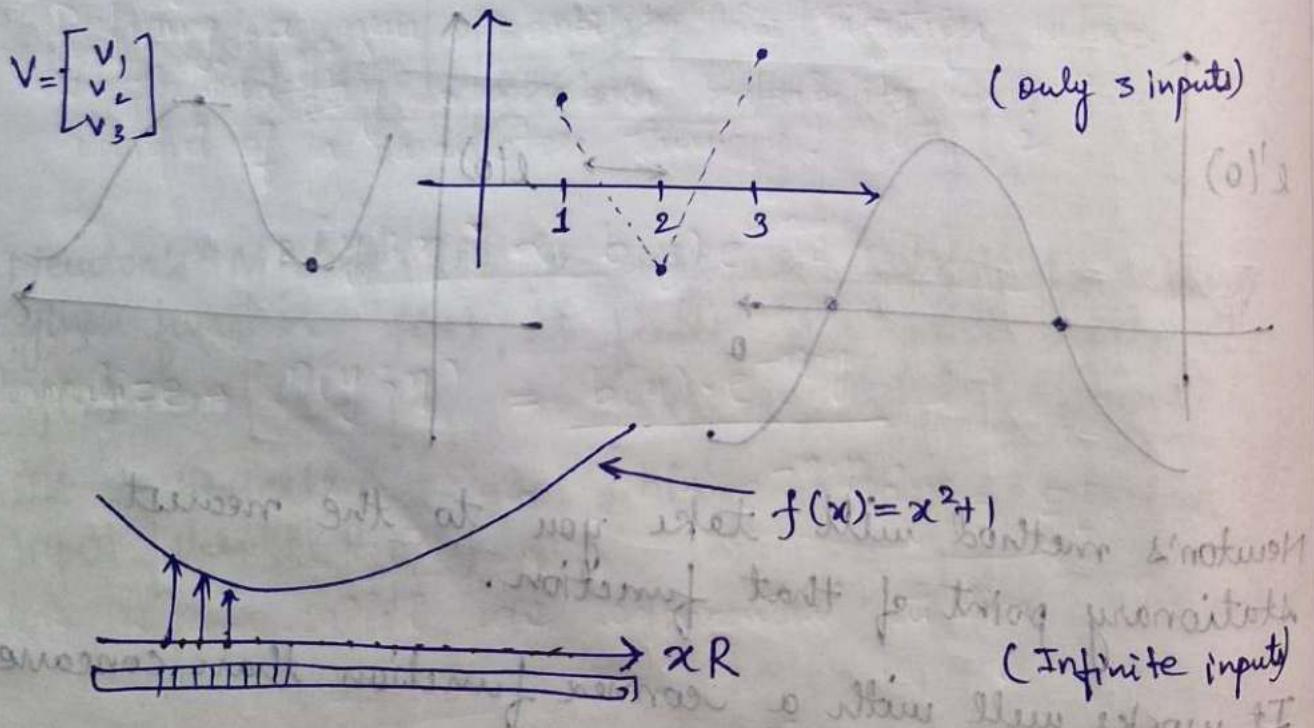
$$v : \{1, 2, 3\} \rightarrow \mathbb{R}$$

$$v(1) = v_1$$

$$v(2) = v_2$$

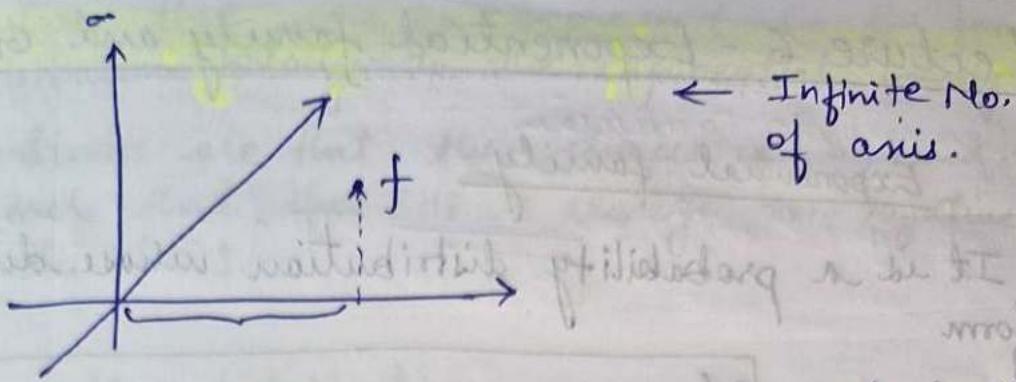
$$v(3) = v_3$$

- Think of functions and vectors as the same.
- Functions are infinite dimensional vector.
- In the vector, the number of axis is the domain of the function.



$$f = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix}$$

*// Basically functions are points in infinite dimensional space.



- The value of the function at a given input is the same as the co-ordinate of the vector at that axis. Same analogy as of vector, just extend it to infinite dimensions.

Finite

- Vector \vec{v}
- Index $\{1, \dots, d\}$
- Components
- Explicit RED
- $v = [v_1, \dots, v_d]$
- Dot product
 $\langle u, v \rangle = \sum_i u_i v_i$
- Matrix $A \rightarrow A_{ij}$
- $y = Ax$ (Linear opo)
- Eigen Vector A
 $Ax = \lambda x$
- $\vec{u} = A\vec{v}$
- $u_i = \sum_j A_{ij} v_j$

Infinite

- function $f(t)$
- Domain \mathbb{R}
- Values
- $f(t) = \text{Symbolic Representation}$
- $\langle f, g \rangle = \int f(t) g(t) dt$
- $K(s, t)$
- $f' = D[f]$ ($D \rightarrow \text{Differentiation}$)
- $f' = D[f] = \lambda f$
- $D[C^{kt}] = kC^{kt}$
- $g(s) = \int K(s, t) f(t) dt$

Lecture 6 - Exponential family and GLM

Exponential family

It is a probability distribution whose density has the form

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

pdf/pmf

$\eta \rightarrow$ parameter (eta)

We are modelling the output of our model with exponential family.

$b(y)$: Base measure

$T(y)$: Sufficient statistic

$\{ T(y) = y \text{ for most cases} \}$

$a(\eta)$: log-partition function

η : Natural Parameter

- # Defining a new probability distribution, whose density is defined by some base probability in y times the exponent of a parameters times y .

$$p(y; \eta) \propto b(y) e^{-\eta \cdot y}$$

Normalize $\Rightarrow p(y; \eta) = \frac{b(y) \cdot e^{-\eta \cdot y}}{\int b(y) e^{-\eta \cdot y} dy} = E[e^{-\eta y}]$

This is also called as exponential tilting

$$p(y; \eta) = \frac{b(y) e^{-\eta^T y}}{E[e^{-\eta^T y}]} = \frac{b(y) \cdot e^{\eta^T y - a(\eta)}}{a(\eta)}$$

A(η)

Normalizing constt.

- Any distribution that can be expressed in the form of this distribution, belongs to the exponential family.
- The only constraints are that the choices of b , T and a should be such that, the RHS is always non-negative and integrates to 1.

Bernoulli distribution

$$\begin{aligned}
 p(y; \phi) &= \phi^y (1-\phi)^{1-y} \\
 &= \exp \left[\log [\phi^y (1-\phi)^{1-y}] \right] \\
 &= \exp \left[y \log \phi + (1-y) \log (1-\phi) \right] \\
 &= \exp \left[\log \left(\frac{\phi}{1-\phi} \right) \cdot y + \log (1-\phi) \right]
 \end{aligned}$$

Compare this with general expression

$$p(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}$$

$$\Rightarrow \eta = \log \left(\frac{\phi}{1-\phi} \right) \Rightarrow \boxed{\phi = \frac{1}{1+e^{-\eta}}} \text{ logistic}$$

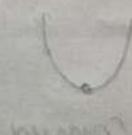
$$\Rightarrow T(y) = y$$

$$\Rightarrow a(\eta) = -\log (1-\phi) \Rightarrow \boxed{\log (1+e^{-\eta}) = a(\eta)}$$

$$\Rightarrow b(y) = 1$$

$$(r) \circ \frac{f}{g} = [r:f] \exists$$

$$(r) \circ \frac{f}{g} = [r:f] \forall$$



Gaussian distribution

Univariate gaussian

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

Assuming $\sigma^2 = 1$

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

$$P(y; \eta) = b(y) \exp\left\{\eta^T T(y) - a(\eta)\right\}$$

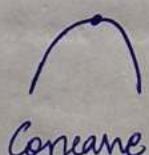
Matching pattern

$$b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) \quad \left| \begin{array}{c} \eta = \mu \\ \mu = \eta \end{array} \right| \quad \begin{array}{l} T(y) = y \\ a(\eta) = \mu^2/2 = \eta^2/2 \end{array}$$

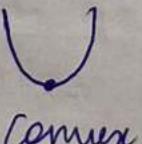
Properties of Exp families

- ① $\log p(y; \eta)$ is concave in η . MLE is concave in η , which means NLL is convex in η .
- ② $E[y; \eta] = \frac{\partial}{\partial \eta} a(\eta)$
- ③ $\text{Var}[y; \eta] = \frac{\partial^2}{\partial \eta^2} a(\eta)$

Global maxima



Global minima



Generalised linear Models (GLM)

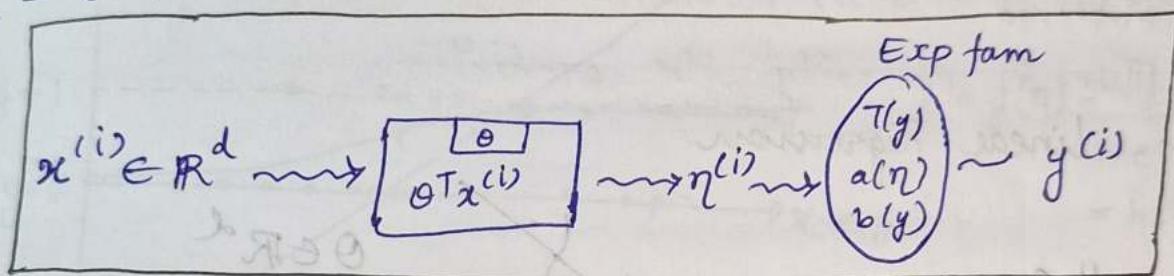
In GLM we use the exponential families to introduce a relation between a set of variables x to the variable y .

$$\textcircled{1} \quad y|x; \theta \sim \text{Exp Family } (\eta)$$

Given x and some θ , y belongs to exp family with parameter η .

$$\textcircled{2} \quad h_\theta(x) = E[y|x; \theta]$$

$$\textcircled{3} \quad \eta = \theta^T x$$



Flowchart of data generation

$\theta \rightarrow$ global (only one)

$\eta^{(i)} \rightarrow$ diff for each example

Each example will give its own η .

Ordinary least square (OLS) : Linear Regression

Exp Family in Gaussian

$$h_\theta(x) = E[y|x; \theta]$$

$$= \mu$$

$$= \eta$$

$$= \theta^T x$$

$$h_\theta(x) = \theta^T x$$

Step at 0 it is better to, (i) x enter twice since ROT +
two OLS with most with minimum σ so $R = 0$. Then line
 $y = \eta^{(i)} + \theta^T x$ with after further for the slopes β

Logistic Regression

Exp family is Bernoulli

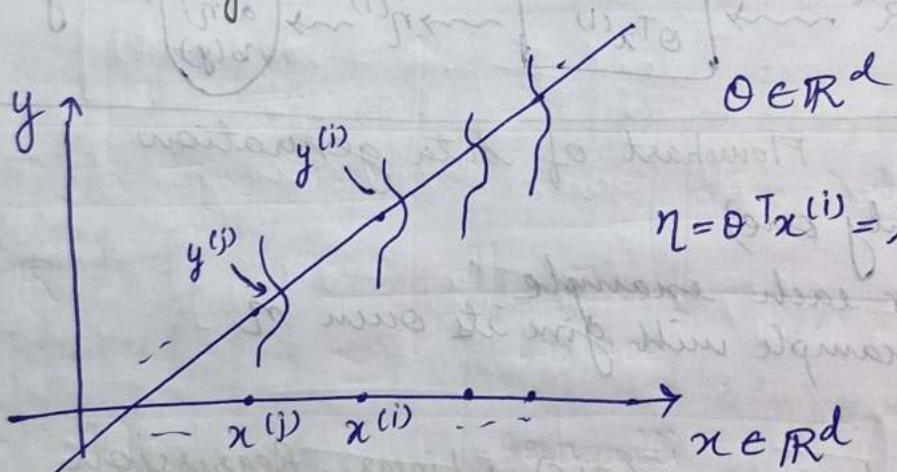
$$h_{\theta}(x) = E[y|x; \theta]$$

$$= \phi$$

$$= \frac{1}{1+e^{-z}}$$

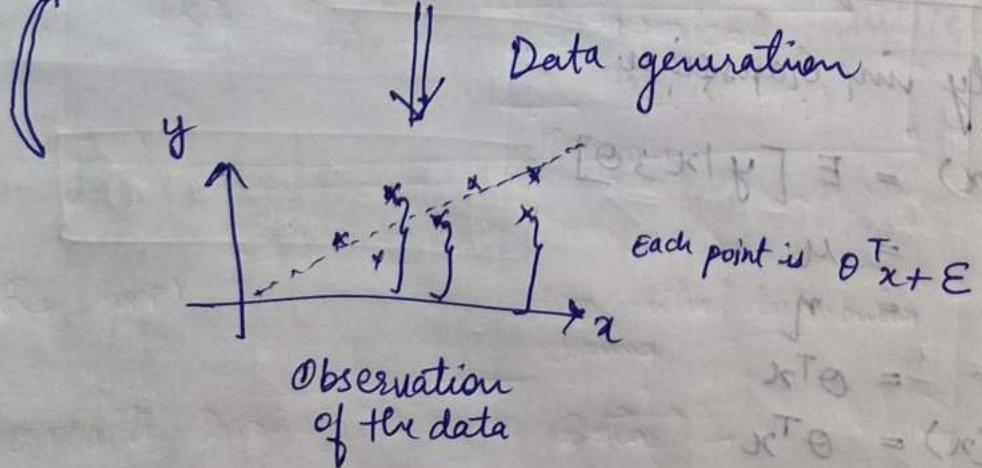
$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} = g(\theta^T x) ; g(z) = \frac{1}{1+e^{-z}}$$

linear Regression



θ ?

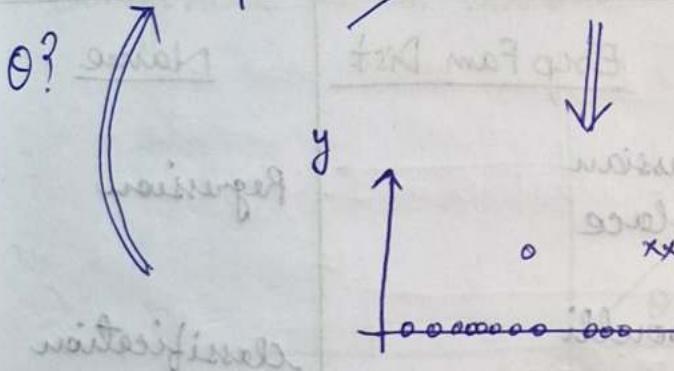
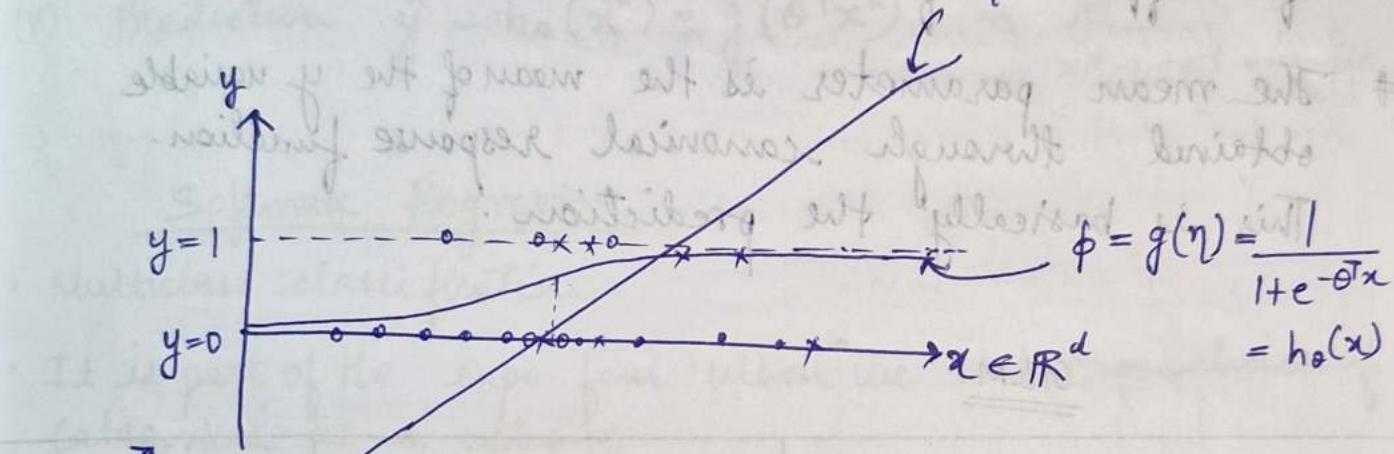
Data generation



For every input value $x^{(i)}$, we started with θ to get η and that $\eta = \mu$ of a Gaussian dist. From this we sample our y . That makes the $(x^{(i)}, y^{(i)})$ pair

Basically the visualization is like there is some θ which acts as the mean for each $y^{(i)}$ that we gonna sample from. The sampled $y^{(i)}$ will have some noise, we will start with these noisy observations and target to extract away the signal from that.

Logistic regression.



Model Parameters

$$\theta \in R^d$$

$$x^{(i)} \in R^d$$

Natural Parameter

$$\eta^{(i)}$$

$$g$$

Mean Parameter

$$\mu: \text{Gaussian} \quad \mu = g(\eta) = \eta$$

$$\phi: \text{Bernoulli} \quad \phi = g(\eta) = \frac{1}{1+e^{-\eta}}$$

$$\lambda: \text{Poisson}$$

$$h_\theta(x) = E[y|x; \theta] = g(\theta^T x)$$

y : Response Variable

g : Canonical response function.

g^{-1} : Canonical link function.

- # The model parameters: When we perform gradient descent, MLE on generalized linear model, these are the parameters that get trained. These are the parameters we are adjusting in each step of gradient descent or stochastic descent.
- # The Natural parameter is just the output of $\theta^T x^{(i)}$ for diff value of x
- # The mean parameter is the mean of the y variable obtained through canonical response function. This is basically the prediction.

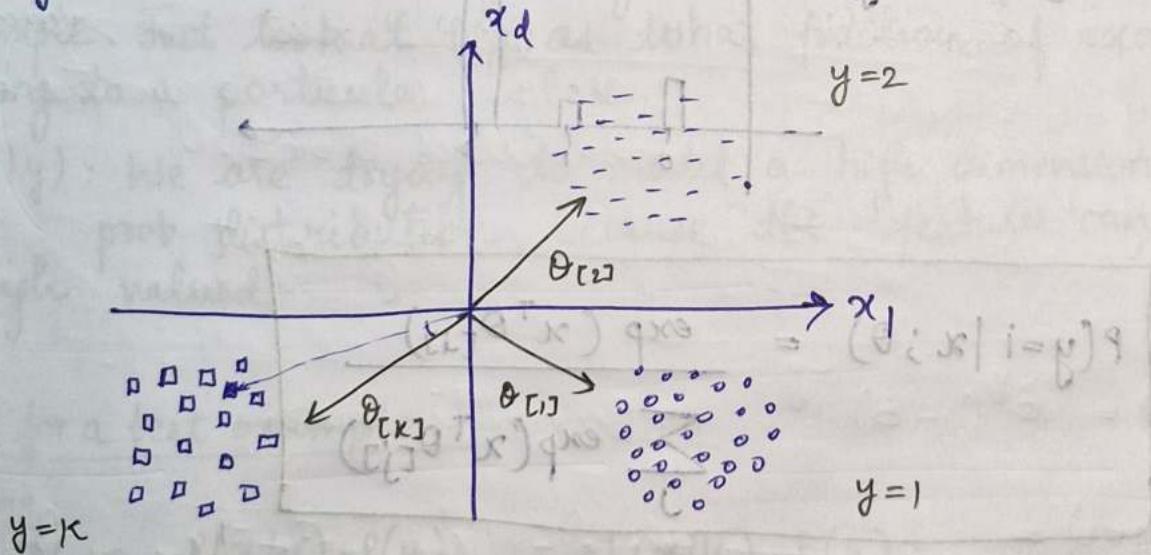
GLM

<u>Data type (y)</u>	<u>Exp Fam Dist</u>	<u>Name</u>
① \mathbb{R}	Gaussian Laplace	Regression
② $\{0, 1\}$	Bernoulli	classification
③ $\{1, \dots, K\}$	Categorical	Multiclass classification (softmax)
④ \mathbb{N}_+	Poisson	Count reg. Poisson reg.
⑤ $\mathbb{R}_+ (\text{time})$	Exponential distribution Gamma	Survival analysis
⑥ Bernoulli Distri	Beta	

- (i) make choice of Distribution according to Data type
 - (ii) Express in Exponential family. - $a(\eta)$, $b(y)$, $T(y)$
 - (iii) Hypothesis $h_\theta(x) = g(\theta^T x)$
 - (iv) $\theta = \theta + \alpha (y^{(i)} - h_\theta(x^{(i)})) \cdot x^{(i)}$
 - (v) Prediction $\hat{y} = h_\theta(x^*) = g(\theta^T x^*)$
- no. need to find
log, Now we can
directly find it
just use $h_\theta(x)$.
 x^* → test example
(unseen)

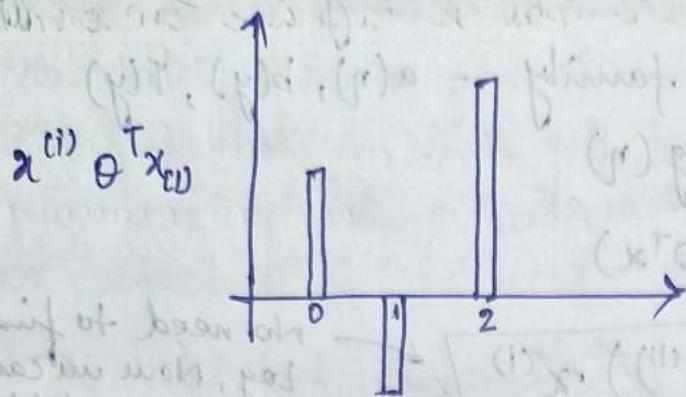
Softmax Regression

- Multiclass classification
- It is part of the expo. fam. When we make ~~a choice of~~ categorical as a choice.

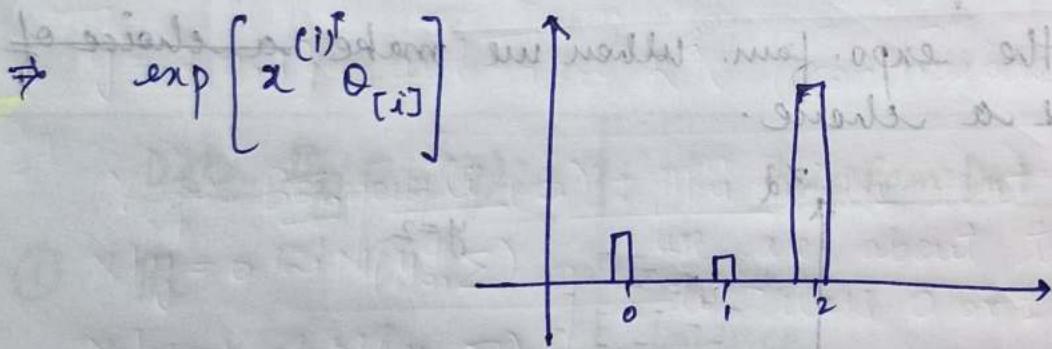


$$\left. \begin{array}{l} \theta_{[1]}^T x = R \\ \theta_{[2]}^T x = R \\ \vdots \\ \theta_{[K]}^T x = R \end{array} \right\} \text{arg max}$$

$$\boxed{\theta^T x = 0} \rightarrow \theta = \frac{1}{K} \begin{bmatrix} -\theta_{[1]}^T \\ -\theta_{[2]}^T \\ \vdots \\ -\theta_{[K]}^T \end{bmatrix} \begin{bmatrix} x \end{bmatrix}$$



- # Softmax takes the set of these scalar valued dot prod and converts them into a probability distribution.
- # So for prob. distribution
 - All the values should be non neg.
 - sum up to 1.



$$P(y=i | x; \theta) = \frac{\exp(x^T \theta_{[i]})}{\sum_j \exp(x^T \theta_{[j]})}$$

Softmax Function

- # $\vec{x}^T \text{Softmax}(\vec{x}) \approx \max(\vec{x})$

Lecture 7: GDA, Naive Bayes & Laplace Smoothing

Focus on $x \rightarrow$ Model $P(x, y) = \underbrace{P(x|y)}_{\text{High dimensional}} \cdot \underbrace{P(y)}_{\text{class prior}}$

So far we have studied Discriminative Algorithms where we directly model $P(y|x)$ but now our focus is on x .

$$\underbrace{P(x, y)}_{\text{model}} = \underbrace{P(x|y)}_{\text{High dimensional}} \cdot \underbrace{P(y)}_{\text{class prior}}$$

- We will be focusing on where y is discrete value
- GDA and Naive Bayes are best suited for classification problem.

- # y indicates the class of a given example. Just look at $P(y)$ as what fraction of examples belong to a particular class.
- * $P(x|y)$: We are trying to model a high dimensional prob. distribution, because the features can be high valued.

Let for a test example

$$\underbrace{P(y=1|x)}_{\substack{\text{posterior} \\ \text{Distribution}}} = \frac{\underbrace{P(x|y) \cdot P(y)}_{\substack{\text{prior}}}}{P(x)} = \frac{P(x|y) \cdot P(y)}{P(x|y=0) \cdot P(y=0) + P(x|y=1) \cdot P(y=1)}$$

Here y is assumed binary. {0, 1}.

$$\hat{y} = \arg \max_y p(y|x) = \arg \max_y \frac{p(x|y) \cdot P(y)}{P(x)} = \arg \max_y \frac{p(x|y) P(y)}{p(x)}$$

When we make prediction as a class rather than a probability. So $P(y=1 \text{ or } y=0 | x)$ gives a prob. for that class if we evaluate denominator too. Otherwise Numerator is enough to decide the class.

Model \Leftrightarrow Data Generating Process
 (Joint $p(x, y)$) \Downarrow Hierarchy of steps
 \Downarrow Factorize our joint probability

Gaussian Discriminant Analysis (GDA)

$y = \{0, 1\}$ and $x \in \mathbb{R}^d$ (continuous)

Step $y \sim \text{Bernoulli}(\phi)$: This decides that the example we are about to generate belongs to class 0 or 1.
 $x|y=0 \sim N(\mu_0, \Sigma)$
 $x|y=1 \sim N(\mu_1, \Sigma)$ \rightarrow Generating input

$$P(y) = \phi^y (1-\phi)^{1-y} \quad : \text{Bernoulli distri.}$$

$$P(x|y=0) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) \right\}$$

$$P(x|y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right\}$$

$$\therefore y \in \{0, 1\}$$

$$(x, p(x)) \in \mathbb{R}^d$$

In discriminative algorithms all we do is to discriminate if an example belongs to this class or the other. which can actually limit ourselves to look for certain queues in the input and ignore the rest of the example.

Whereas, in GDA we generate an entire example. we come up with full descriptive of what makes up $y=1$ or $y=0$. In simple terms Generative modelling is harder than Discriminative modelling.

$$\underbrace{P(x, y)}_{\text{Generative model}} = \underbrace{P(y|x) \cdot P(x)}_{\text{Discriminative mode}}$$

So in order to learn a gen. model, we have to learn a discriminative model and also something else.

Max Likelihood to learn parameters.

- * After generating the data, we observe the data. We can use the data to fit our model and learn these parameters $(\phi, \mu_0, \mu_1, \Sigma)$ using MLE.
- * A likelihood is a function over parameters

Log likelihood

$$\begin{aligned} l(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n P(x^{(i)}, y^{(i)}) \\ &= \log \prod_{i=1}^n P(x^{(i)}|y^{(i)}) P(y^{(i)}) \end{aligned}$$

$$l(\phi, \mu_0, \mu_1, \Sigma) = \log L(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^m P(x^{(i)}, y^{(i)})$$

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^m \log p(x^{(i)}, y^{(i)})$$

$$= \sum_{i=1}^m (\log p(x^{(i)} | y^{(i)}) + \log p(y^{(i)}))$$

* MLE: $\nabla \ell(\phi, \mu_0, \mu_1, \Sigma) = 0$

① Gradient with respect to ϕ :

$$P(y^{(i)}) = \phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}}$$

$$\sum_{i=1}^m \log p(y^{(i)}) = \sum_{i=1}^m [y^{(i)} \log \phi + (1-y^{(i)}) \log (1-\phi)]$$

$$\frac{\partial \ell}{\partial \phi} = \sum_{i=1}^m \left[\frac{y^{(i)}}{\phi} - \frac{1-y^{(i)}}{1-\phi} \right] = 0$$

$$\hat{\phi} = \frac{1}{m} \sum_{i=1}^m y^{(i)}$$

This is MLE for ϕ , which is simply the fraction of data points where $y^{(i)} = 1$.

② Gradient with respect to μ_0 :

$$\sum_{i: y^{(i)}=0} \log p(x^{(i)} | y^{(i)}=0)$$

$$\text{Ans} = \frac{-1}{2} ((x^{(i)} - \mu_0))^T \sum_{i: y^{(i)}=1} (x^{(i)} - \mu_0) + \log |\Sigma| + d \log(2\pi)$$

$$\frac{\partial l}{\partial \mu_0} = \sum_{i:y^{(i)}=0} \sum^{-1} (x^{(i)} - \mu_0)$$

$$\hat{\mu}_0 = \frac{\sum_{i:y^{(i)}=0} x^{(i)}}{\sum_{i:y^{(i)}=0} 1}$$

③ Gradient with respect to μ_1 ,

$$\sum_{i:y^{(i)}=1} \log p(x^{(i)} | y^{(i)}=1)$$

$$\frac{\partial l}{\partial \mu_1} = \sum_{i:y^{(i)}=1} \sum^{-1} (x^{(i)} - \mu_1)$$

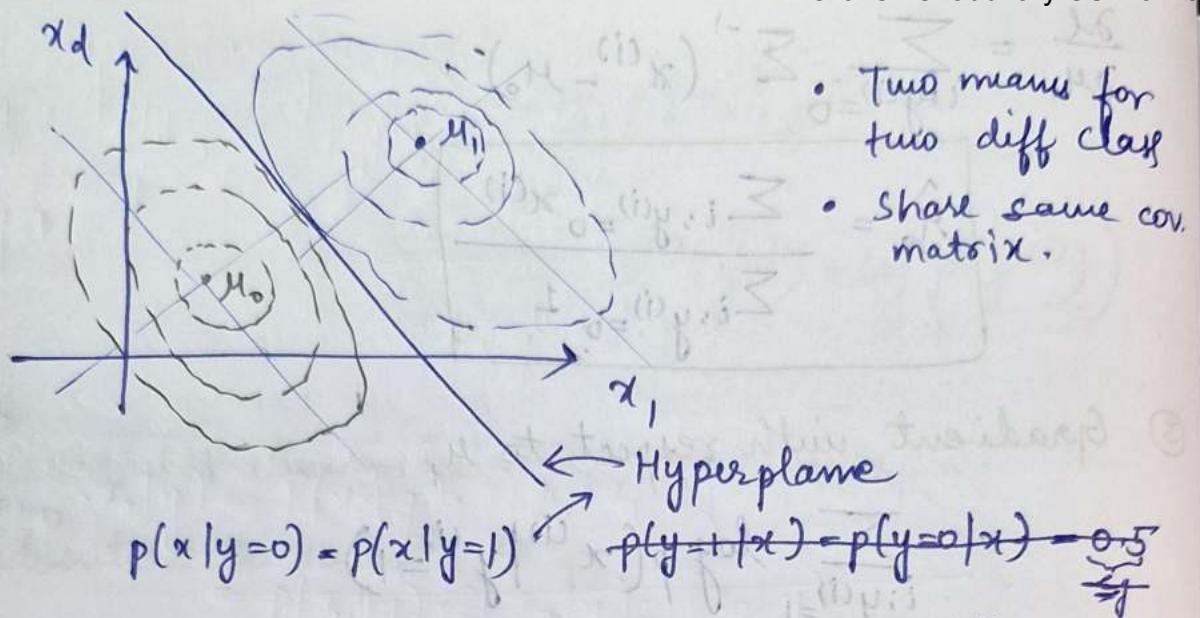
$$\hat{\mu}_1 = \frac{\sum_{i:y^{(i)}=1} x^{(i)}}{\sum_{i:y^{(i)}=1} 1}$$

④ Gradient w.r.t Σ

$$\sum_i \log p(x^{(i)} | y^{(i)})$$

$$\frac{\partial l}{\partial \Sigma} = -\frac{1}{2} \sum_{i=1}^m \left[\Sigma^{-1} - \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} \right]$$

$$\hat{\Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T$$



- Two means for two diff class
- Share same cov. matrix.

For any GDA models that share the same Σ , the posterior distribution of $(y|x)$ can be represented as a Logistic regression model, where the θ only depends on the parameters of the model ($\mu_0, \mu_1, \phi, \Sigma$)

$$\text{GDA} \Rightarrow \text{Logistic Regression : } P(y=1|x) = \frac{1}{1+e^{-\theta^T x}}$$

$$\left[\sum_{i=1}^m (y_i \phi(\mu_1 - \mu_0)^T x_i) - \sum_{i=1}^m \log \frac{1}{1+e^{-\theta^T x_i}} \right] \sum_{i=1}^m \frac{1}{1+e^{-\theta^T x_i}} = \frac{16}{28}$$

$$\sum_{i=1}^m (y_i \phi(\mu_1 - \mu_0)^T x_i) \sum_{i=1}^m \frac{1}{1+e^{-\theta^T x_i}}$$

Naive Bayes

x - Discrete valued

Used in for eg. text classification (spam filters)

Conditional Independence

$$P(x_j | x_k) = p(x_j) \quad [\text{independent}]$$

$$y: P(x_j | x_k, y) = P(x_j | y) \quad [\text{conditional independence}]$$

Bernoulli Event Model

"Buy our lottery" = $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$

aardvark
 aardwolf
 buy
 lottery
 our
 zygromorph

Vocabulary - d

$x \in \{0, 1\}^d$

$x_j \in \{0, 1\}$

Model

$$P(y=1) = \text{Bernoulli } (\phi_y) \quad - 1 \text{ parameter}$$

$$P(x_j=1 | y=0) = \text{Bernoulli } (\phi_j | y=0) \quad - d \text{ parameters}$$

$$P(x_j=1 | y=1) = \text{Bernoulli } (\phi_j | y=1) \quad - d \text{ parameters}$$

The intuition that eg $P(x_j=1 | y=1)$ what's the probability that this word corresponding to index j will show up in an email that is spam.

$$\lambda(\underbrace{\phi_y}_{\text{d}}, \underbrace{\phi_{j|y=0}}_{\text{d}}, \underbrace{\phi_{j|y=1}}_{\text{d}}) = \log \prod_{i=1}^n P(x^{(i)}, y^{(i)}; \phi)$$

$$= \log \prod_{i=1}^n P(y^{(i)}; \phi_y) \left(\frac{1}{n} \prod_{j=1}^d P(x_j^{(i)} | y^{(i)}; \phi) \right)$$

Property

$$P(x_1, x_2, \dots, x_n | y) = P(x_1 | y) \cdot P(x_2 | x_1, y) \cdot P(x_3 | x_1, x_2, y) \dots$$

$$= P(x_1 | y) \cdot P(x_2 | y) \cdot P(x_3 | y) \dots$$

MLE estimatesLemma three illustrated

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n I\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n I\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n I\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n I\{y^{(i)} = 0\}}$$

$$\phi_y = \frac{\sum_{i=1}^n I\{y^{(i)} = 1\}}{n}$$

- Now if a test word comes
 • calculate the probability and compare.

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)} = \frac{P(x|y) P(y)}{P(x|y=0) P(y=0) + P(x|y=1) P(y=1)}$$

$$P(y=1|x) = \frac{\prod_{j=1}^d P(x_j^{(i)}|y=1) P(y=1)}{\text{denom.}}$$

$$P(y=0|x) = \frac{\prod_{j=1}^d P(x_j^{(i)}|y=0) P(y=0)}{\text{denom}}$$

} Compare them
to decide
spam or not

Now what if at test time we encounter a word that was never seen in the training set.

$$P(y|x) = \frac{0}{0+0}$$

With this method the problem is that, in this case we cannot tell what our prediction should be. To overcome this a technique is used called Laplace smoothing.

Eg Suppose a coin toss, and on all 10 trials H turns up.
 therefore $x \in \{0, 1\}$ but to encounter such

$\phi = P(H) = 10/10 = 1$
 problem in Laplace smoothing we initialize count of every example by 1. (Basically assuming that we have seen that example at least once).

$$\leftarrow \phi \text{ preliminary}$$

After Laplace smoothing

$$\phi_{j|y=1} = \frac{1 + \sum_{i=1}^n \mathbb{1}(x_j^{(i)} = 1 \wedge y^{(i)} = 1)}{2 + \sum_{i=1}^n \mathbb{1}(y^{(i)} = 1)}$$

$$\phi_{j|y=0} = \frac{1 + \sum_{i=1}^n \mathbb{1}(x_j^{(i)} = 1 \wedge y^{(i)} = 0)}{2 + \sum_{i=1}^n \mathbb{1}(y^{(i)} = 0)}$$

Multinomial Event Model

$$y \sim \text{Bernoulli } (\phi) \quad (1 \text{ dimension})$$

$$x|y=0 \sim \text{Multinomial Categorical } (\phi) \quad (|\mathcal{V}| - 1 \text{ dimension})$$

MLE

$$\phi_{k|y=1} = \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} \mathbb{1}\{x_j^{(i)} = k \wedge y^{(i)} = 1\}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\} d_i}$$

- # k is the word in the vocabulary is given by \rightarrow summing over each message, within each message we are summing over every word and counting only those words in those messages which happen to be k and belong to the class we are interested in. Divided by total no. of words across all messages in that class.

Similarly $\phi_{k|y=0}$.

$$\phi_{k|y=0} = \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} I\{x_j^{(i)} = k \wedge y^{(i)} = 0\}}{\sum_{i=1}^n I\{y^{(i)} = 0\} d_i}$$

"just for zero prob"
"where zero prob"

Laplace smoothing

$$\# \phi_{k|y=1} = \frac{1 + N_{k|y=1}}{|V| + D_{k|y=1}}$$

• N - Numerator
from previous proof.

$$\# \phi_{k|y=0} = \frac{1 + N_{k|y=0}}{|V| + N_{k|y=0}}$$

"base message draw a
illustrated in several
ways all the way to
have all the ways in
the same way"

• D - Denominator
from the previous
proof.

Bernoulli method

	a	and work	buy	exam	
"buy our lottery"	0	0	0	1	0 0 1 -
"buy this watch, this watch only"	0	0	0	1 0	0 0 1 0 1 ...
"when is your exam"	$\phi_{j y=1}$	0	0	0	0 1 0 1 0
"When is your homework due"		0	0	0 1	0 0 1 0 0
					$\frac{1}{2}$
					$\phi_{j y=0}$

- * We put 1 only for the no. of message in which the word appear, doesn't matter how many times the word appears in the same message the count will be 1 only.

- * $\phi \rightarrow \frac{\text{Messages in which the word appears}}{\text{Total no. of messages.}}$

Multinomial model

"buy our lottery"

a a work -- buy -- exam --

0 0 1 - - 1 0

"buy this watch, this watch only"

0 0 1 - - 1 2

4 2 3 2 2

$$\Sigma = 4+2+3+2+2$$

"When is your exam"

"When is your homework due"

rearranged - C.

writing out word

Here we consider the number of times a word appear across all examples in your training set, whereas in Bernoulli it was the number of messages in which the word appear.

In Bernoulli each word has its own bernoulli variable, so we normalized it locally. But in Multinomial mode we want a distribution over words, so we normalize the entire thing.

0 1 0 1 0 0 0 0 0

"many ways to write"

0 0 1 0 0 1 0 0 0

"substitution way is will"

draw out string in question for am off road plus I took ski
 draw out with fewer word return first abcd, abcd
 - after feed New things with question easier off we change
 change draw off harder in question + phi +

- different for all later

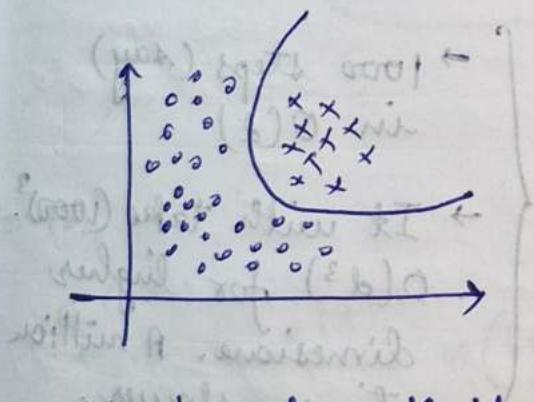
Lecture 8 - Kernel Methods & SVM

Kernel Methods

So far we learned the linear and logistic regression are linear in nature (linear models). They provide linear hypothesis.

But we can construct non-linear hypothesis using linear regression and similarly non-linear classifiers using logistic regression by including higher order features.

$$\underbrace{x}_{\text{attribute}} \rightarrow \underbrace{\phi(x)}_{\text{Features}} = \begin{bmatrix} 1 \\ x^1 \\ x^2 \\ \vdots \\ x^n \end{bmatrix} \rightarrow \begin{array}{c} y \\ \text{non-linear function} \end{array}$$



$$\text{linear function } \phi(x) = \begin{bmatrix} 1 \\ x^1 \\ x^2 \\ \vdots \\ x^n \end{bmatrix}$$

- * With kernel methods we can map an example to potentially infinite dimensional feature vector and perform the learning algorithm in the infinite dimensional space.

$$(i) x \phi^{(k)} \cdot \sum_{j=1}^k \theta_j = 0$$

Linear regression with gradient descent.

update rule is

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) \cdot x^{(i)} \quad [\theta \in \mathbb{R}^d]$$

With feature map.

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

$$\Rightarrow \phi : \mathbb{R}^d \rightarrow \mathbb{R}^p \quad [\theta \in \mathbb{R}^p]$$

The ϕ is taking from d to some high dimensional space, potentially ∞ .

Eg

$$\phi(x) = \begin{bmatrix} x_1 \\ x_1^2 \\ x_1^3 \\ x_1^4 \\ \vdots \\ x_1^2 x_2 \end{bmatrix}$$

monomial terms of order ≤ 3

$p \approx O(d^3)$

- $\rightarrow 1000$ steps (say) in $O(d)$
- \rightarrow It will take $(1000)^3$ $O(d^3)$ for higher dimensions. A million times slower.

- # The claim is that the $\theta^{(t)}$ that we encounter at any stage of our gradient descent can always be represented as a linear combination of the features.

$$\boxed{\theta^{(t)} = \sum_{i=1}^n \beta_i^{(t)} \phi(x^{(i)})}$$

$$\theta^{(t)} = \sum_{i=1}^n \underbrace{\alpha y^{(i)} \cdot \phi(x^{(i)})}_{\beta_i^{(t)}}$$

Inductive approach to prove.

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \sum_{i=1}^n (y^{(i)} - \theta^{(t)} \cdot \phi(x^{(i)})) \phi(x^{(i)})$$

$$= \sum_{i=1}^n \beta_i^{(t)} \phi(x^{(i)}) +$$

$$\alpha \sum_{i=1}^n (y^{(i)} - \left(\sum_{j=1}^n \beta_j^{(t)} \phi(x^{(j)}) \right)^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

$$\theta^{(t+1)} = \sum_{i=1}^n \left[\beta_i^{(t)} + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} \phi(x^{(j)}) \right)^T \phi(x^{(i)}) \right] \phi(x^{(i)})$$

↓

For $i=1, \dots, n$

$$\beta_i^{(t+1)} = \beta_i^{(t)} + \alpha \left(y^{(i)} - \sum_{j=1}^n (\beta_j^{(t)})^T \phi(x^{(j)}) \phi(x^{(i)}) \right)$$

Kernel

$$\triangleq K: X \times X \rightarrow \mathbb{R}$$

$x \in \mathbb{R}^d \quad X = \mathbb{R}^d$

$$K(x, z) = \langle \phi(x), \phi(z) \rangle$$

$$K(x, z) = \phi(x)^T \phi(z)$$

* A kernel is a function which takes two elements of that space and returns a real valued scalar. It also

satisfies the property that the evaluated value can be expressed as the inner product between two feature maps (feature map of input 1 and dot product of feature map of input 2).

Eq.

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$x \in \mathbb{R}^d + \underbrace{\dots}_{\approx O(d^3)}$$

$$K(x, z) = 1 + \underbrace{\langle x, z \rangle}_{O(d)} + \underbrace{\langle x, z \rangle^2}_{O(d)} + \underbrace{\langle x, z \rangle^3}_{O(d)} \quad \boxed{O(d^3)}$$

$$K(x, z) = \phi(x)^T \phi(z)$$

$\phi(d^3)$ - otherwise
but this approach
Reduce it to $O(d)$

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \cdot \begin{bmatrix} 1 \\ z_1 \\ z_2 \\ \vdots \\ z_d \end{bmatrix} = \underbrace{\frac{1}{d} x_1 z_1 + x_2 z_2 + \dots + x_d z_d}_{O(d^3)}$$

* Basically an order $O(d^2)$ operation is reduced to $O(d)$.

- # For a feature ϕ can have a corresponding kernel K where the kernel function has a more compact representation which is equivalent to performing a dot product in the high dimensional space.

$$\phi : X \rightarrow \mathbb{R}^d$$

for kernels and dot product between vectors x in length A and z in length B is $O(A \cdot B)$

Linear Regression (Kernelized)

① Precompute:

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

$$K_{ij} = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

We use same K notation for both matrix K function here K_{ij} is the square symm. matrix and $K(x^{(i)}, x^{(j)})$ is function.

② Loop: $\forall i \in \{1, 2, \dots, n\}$, $\beta^{(t)} \in \mathbb{R}^n$ (vector)

$$\beta_i^{(t+1)} = \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} K_{ij})$$

$$\beta^{(t+1)} = \beta^{(t)} + \alpha (\vec{y} - K\beta^{(t)})$$

Prediction

$$h_\theta(x) = \theta^T \phi(x)$$

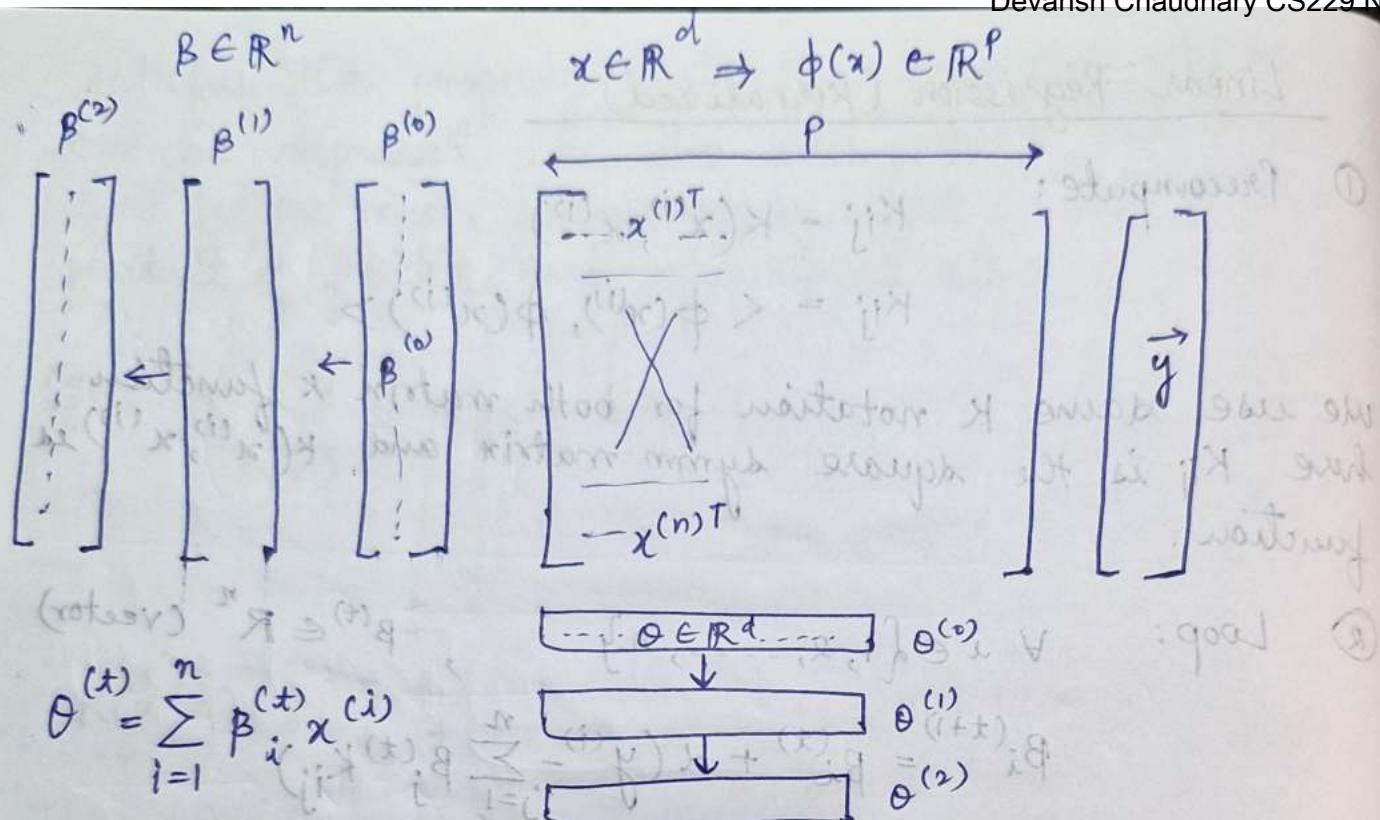
$$= \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x)$$

$$= \sum_{i=1}^n \beta_i K(x^{(i)}, x)$$

Observations

① Train = $\beta = \beta + \alpha (\vec{y} - K\beta)$ $\vec{\phi}(x)$ does not appear
 Test = $\hat{y} = \sum_{i=1}^n \beta_i K(x^{(i)}, x)$ If we had $\phi(x)$ we have to compute it and it could be as dimensional

- ② For prediction we need training examples to be stored in memory.
- ③ As can be inferred from ①, we need to remember all our training examples for testing



- # θ at any descent can be represented as a linear combination of vectors.
- # θ has one component per feature.
- # β has one component per example
- # Every β vector has a corresponding θ vector
- # This helps us to have infinite feature wide, and we achieve scaling.
- # Kernelization is a general technique that can be applied to a lot of places - Gen or discrim, super or unsup classification or regression etc.

Kernel examples

eg

$$A \quad K(x, z) = \langle x, z \rangle^2 \quad x \in \mathbb{R}^d$$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_d x_d \end{bmatrix}$$

$$B \quad K(x, z) = (x^T z + c)^2$$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ \vdots \\ x_d^2 \end{bmatrix} \quad K(x, z) = \phi(x)^T \phi(z)$$

$K(x, z) \uparrow$ (for similar x, z)

$K(x, z) \downarrow$ (for non similar x, z)

$$K(x, z) = \phi(x)^T \phi(z)$$

$$K(x, z) = e^{(-\|x - z\|^2 / 2\sigma^2)}$$

Necessary condition for K to be a kernel

- K should be symmetric

$$K(x, z) = K(z, x)$$

- $\{x^{(1)}, \dots, x^{(m)}\}$

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

K is symmetric & PSD

$$z^T K z = \sum_i \sum_j z_i K_{ij} z_j$$

$$= \sum_k \left(\sum_i z_i \phi_k(x^{(i)}) \right)^2 \geq 0$$

Mercer's theorem

Let $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be given.

For K to be a kernel, it is necessary and sufficient for any $\{x^{(1)}, \dots, x^{(m)}\}$ the corresponding kernel matrix $K_{ij} = K(x^{(i)}, x^{(j)})$ is PSD and symmetric.

Vector \longleftrightarrow function $\begin{bmatrix} f(x) \\ f'(x) \\ \vdots \\ f^{(k)}(x) \end{bmatrix} \in (\mathbb{R})^k$

Matrix \longleftrightarrow Operator / function (x, y)

$$K(x, z) = \begin{bmatrix} \cdot & \cdot & \cdots & \cdot & \cdot \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ \cdot & \cdot & \ddots & \cdot & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdots & \cdot & \cdot \end{bmatrix}_{x^{(1)} \quad \vdots \quad x^{(m)}} \quad \begin{matrix} \text{(\leftarrow x eliminated)} \\ \text{(\leftarrow z eliminated)} \end{matrix} \quad \begin{matrix} z^{(1)} \\ \vdots \\ z^{(n)} \end{matrix} \quad K \begin{bmatrix} \cdot & \cdot & \cdots & \cdot & \cdot \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ \cdot & \cdot & \ddots & \cdot & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdots & \cdot & \cdot \end{bmatrix}_{x^{(1)} \quad \vdots \quad x^{(m)}}$$

To prove K is kernel

① Construct ϕ

such that $K(\cdot, \cdot) = \phi^T \phi$

② Mercer's theorem

$$\{x^{(1)}, \dots, x^{(m)}\} \xrightarrow{\text{let } x_i^{(i)} = \phi_i} \{x_1, \dots, x_m\}$$

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

is PSD

③ $\int \int f(x) K(x, x') f(x') dx dx' \geq 0$

$$0 \leq \left(\int x^{(i)} \phi_i \right)^2 = \sum_i \phi_i^2 = \phi^T \phi$$

Support Vector Machines (SVM)

- Discriminative
- Classification

Notations

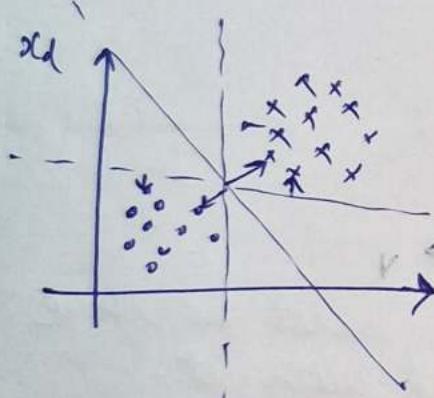
$$y^{(i)} \in \{+1, -1\}$$

Parameters w, b

$$w \in \mathbb{R}^d, b \in \mathbb{R}$$

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

$$\Rightarrow w^T x + b$$



→ What should be the ideal separating hyperplane that we should make
→ For this we use Margin.

$$\text{Margin} = y^{(i)} (w^T x^{(i)} + b) > 0$$

$$w^T x + b > 0 \quad \text{for } y = +1$$

$$w^T x + b < 0 \quad \text{for } y = -1$$

Desire: Margin to be large

The idea of SVM is to calculate the margin for w.r.t all the examples and maximize the smallest margin.

SVM is an algorithm that tries to maximize the geometric margin (i.e. the actual geometric distance b/w these separating hyperplane).

Problem with $y^{(i)}(w^T x + b) > 0$ is that it can be scaled just by increasing coeff. i.e. $y^{(i)}(2w^T x + 2b)$

SVM addresses this problem

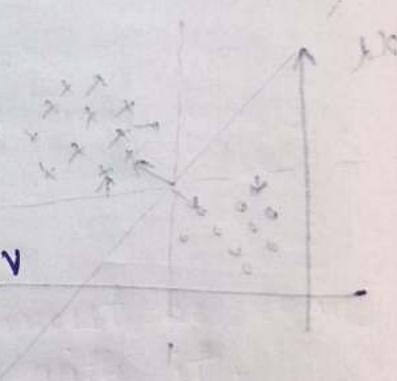
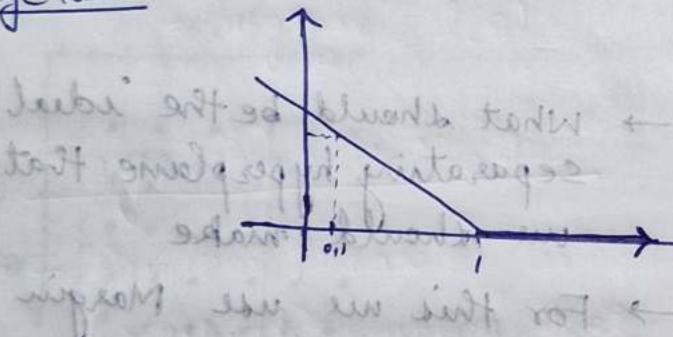
$$\underset{w, b}{\text{Min}} \left(\sum_{i=1}^n \max \left[0, 1 - y^{(i)}(w^T x^{(i)} + b) \right] \right) + \frac{1}{C} \|w\|^2$$

Margin = $y^{(i)}$

Safeguard for scaling

Hinge loss / SVM loss.

Hinge loss



$$\Rightarrow \underset{\xi, w, b}{\text{Min}} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

such that,

$$y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i \quad \forall i \in \{1, \dots, n\}$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$

Primal
Convex
Problem

$$\Rightarrow \text{Max} = \sum_{i,j} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y^{(i)} y^{(j)} \quad \text{such that } \alpha_i \alpha_j \leq x^{(i)}, x^{(j)} \geq 0$$

such that $0 \leq \alpha_i \leq C$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0$$

Dual
Convex
problem

A property of SVM is that once we solve the set of α_i , as a result the set of α_i will be sparse, which means only a small number of α_i will be non zero, but most of the $\alpha_i = 0$.

Those set of example for which $\alpha_i \neq 0$ are called the support vectors. Intuitively those examples will be nearest to the separating hyperplane.

$$(0)l - \min_{\alpha} g(\alpha) = 0$$

$$(0)l + \max_{\alpha} g(\alpha) = 0$$

homework problem 8

homework problems solved method - 8

and then with test sample which \leftarrow

the new feature vector being

methodology noise in a

$$(0|x)_q = (\alpha_0) x, Y$$

blue version:

$$\frac{(0)_q \cdot (\alpha|x)_q}{(x)_q} = (x|0)_q$$

$$\frac{(0)_q \cdot (\alpha|x)_q}{(x)_q} =$$

$$\frac{ab(f|x)_q}{ab(\alpha)_q \cdot (\alpha|x)_q}$$

in writing up in the homework missed out # in LA. this is a stimulus task, ELM set final

Lecture 9 - Bayesian Methods - Parametric & Non

Frequentist Methods

In the methods that we have studied previously we take θ as some constant

Unknown constant θ

$$l(\theta) = \log p(\text{data} ; \theta)$$

$-l(\theta)$ = loss function

$$\hat{\theta} = \arg \min_{\theta} -l(\theta)$$

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

Bayesian Method

θ - Random Variable unobserved

↳ which implies that there must be some prior distribution associated with it.

$\theta \sim \text{Prior Distribution}$

$$Y, X \text{ (Data)} \sim P(X|\theta)$$

$$\underbrace{P(\theta|X)}_{\text{Posterior Distr.}} = \frac{P(X|\theta) \cdot P(\theta)}{P(X)} \quad \therefore \text{ Bayesian rule}$$

$$= \frac{P(x|\theta) P(\theta)}{\int P(x|\theta) \cdot P(\theta) d\theta} = \frac{1}{\int P(x,y) d\theta}$$

In the Bayesian method there is no loss function, we don't take MLE, don't calculate gradients. All we

do is apply Bayes rule and get update our belief about θ given the data.

In Supervise ML

$$\theta \sim \text{Prior}, \quad \theta \perp x \quad (\text{Independent})$$

$$Y \sim P(y|x, \theta)$$

Posterior: $P(\theta|x, Y) = \frac{P(y|x, \theta) \cdot P(\theta)}{P(y|x)}$

Posterior predictive distribution

$$P(Y_*|X, Y, X_*) = \int P(Y_*|X_*, \theta) \cdot P(\theta|x, Y) d\theta$$

Interpretation

- * With the Posterior distribution we are coming up with a probability estimate of how likely a given θ is the correct θ .
- * $\int P(Y_*|X_*, \theta) \cdot P(\theta|x, Y) d\theta \Rightarrow$ Making a prediction for Y for some given value of θ . Then we are calculating the weighted avg. of all the predictions where the weights are decided by Posterior distribution.



Parametric

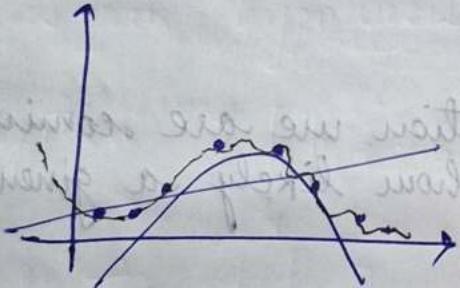
This approach we discussed will be called parametric because $p(y|x, \theta)$ has some functional form for which θ is a parameter

$$\Rightarrow Y|X; \theta = \frac{1}{1 + e^{-\theta^T x}}$$

- * The degree of freedom here we have is to just change θ and that in a way limits the set of all possible models that we can consider

$$[\theta^T (x, x | \theta) \cdot (\theta, x | \gamma)] = (x, x | \gamma)^T$$

- May be linear, quadratic or some higher degree polynomial



- * By limiting ourselves to certain kinds of parametric function, we are essentially saying that no matter how much data is given to us we are not flexible to fit all the variance.
- # That is the limitation of Parametric function that the functional form of the hypothesis function limits us from being flexible even in the presence of lot of data which suggests that there is some other pattern.

Bayesian Linear Regression

$$\{x^{(i)}, y^{(i)}\}_{i=1}^n$$

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$$\epsilon^{(i)} \sim N(0, \sigma^2)$$

$$\theta \sim N(\vec{\theta}, \gamma I)$$

Frequentist approach

$$l(\theta) = \log P(y|x; \theta)$$

$$= \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(y - \theta^T x)^2}{\sigma^2} \right\}$$

Bayes

$$l = \log (Z, H(x))$$

posterior

$$\text{Posterior: } \theta | s \sim N \left(\frac{1}{\sigma^2} A^{-1} x^T \vec{y}, A^{-1} \right) \text{ where } A = \frac{1}{\sigma^2} x^T x + \frac{1}{\gamma^2} I$$

$$\theta | s \sim N \left(\underbrace{\left(x^T x + \frac{\sigma^2}{\gamma^2} I \right)^{-1} x^T y}_{\text{similar to } (x^T x)^{-1} x^T y}, \left(\frac{1}{\sigma^2} x^T x + \frac{1}{\gamma^2} I \right)^{-1} \right)$$

Posterior
Predictive:

Property of Gaussians

$$\theta \sim N(\mu, \Sigma)$$

$$y_* = x^T \theta \sim N(\mu^T x, x^T \Sigma x)$$

$$Y_* | X_*, S \sim N \left(\frac{1}{\sigma^2} x_*^T A^{-1} x^T \vec{y}, x_*^T A^{-1} x_* + \sigma^2 \right)$$

- # In general, Bayesian methods tends to be heavy utilizers of probability theory because all we do is taking probability distributions and conditioning them to apply Bayes rule.

Gaussian Processes

Vector : functions :: MV Gaussians :: Gaussian Process

\downarrow
vectors. \downarrow
functions

Properties

The function over which we will define our gaussian process is our hypothesis function $h(x) = y$.

① Normalization

$$\int_{\mathcal{X}} P(x; \mu, \Sigma) dx = 1$$

② Marginalization

$$x = \begin{bmatrix} x_A \\ x_B \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix} \right)$$

$$P(x_A) = \int_{\mathcal{X}_B} P(x; \mu, \Sigma) dx_B$$

$$= N(\mu_A, \Sigma_{AA})$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1d} \\ \Sigma_{21} & \Sigma_{22} & \dots & \Sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{d1} & \Sigma_{d2} & \dots & \Sigma_{dd} \end{bmatrix} \right)$$

③ Conditioning

$$x_A | x_B \sim N\left(\mu_A + \sum_{AB} \sum_{BB}^{-1} (x_B - \mu_B), \sum_{AA} - \sum_{AB} \sum_{BB}^{-1} \sum_{BA}\right)$$

* Covariance = $\rho \cdot \text{std}(A) \cdot \text{std}(B)$

$$\begin{bmatrix} a \\ b \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \sigma_a^2 & \rho \sigma_a \sigma_b \\ \rho \sigma_a \sigma_b & \sigma_b^2 \end{bmatrix}\right)$$

$$a|b \sim N\left(\mu_a + \frac{\rho \sigma_a \sigma_b}{\sigma_b^2} (b - \mu_b), \sigma_a^2 - \rho \sigma_a \sigma_b \frac{\rho \sigma_a \sigma_b}{\sigma_b^2}\right)$$

$$\sim N\left(\mu_a + \sigma_a \cdot \rho \cdot \left(\frac{b - \mu_b}{\sigma_b}\right), \sigma_a^2 (1 - \rho)\right)$$

④ Summation

$$x \sim N(\mu_1, \Sigma_1)$$

$$y \sim N(\mu_2, \Sigma_2)$$

$$x+y \sim N(\mu_1 + \mu_2, \Sigma_1 + \Sigma_2)$$

(current E 2023, 15M) 429

$$\left(\begin{bmatrix} (x, x) & (x, y) \\ (y, x) & (y, y) \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right)_H = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

MV GaussianGaussian Process

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix}, \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \right) \Leftrightarrow \begin{bmatrix} f \\ f \\ \vdots \\ f \end{bmatrix} \sim GP \left(\begin{bmatrix} m \\ m \\ \vdots \\ m \end{bmatrix}, \begin{bmatrix} K & & & \\ & K & & \\ & & K & & \\ & & & K & \end{bmatrix} \right)$$

Vector Mean Vector Covariance Matrix functions Mean func. Kernel
 ↓ ↓ ↓

Marginalize out all irrelevant Examples i.e eg not in training or testing set.

$$\begin{bmatrix} f(x^{(1)}) \\ \vdots \\ f(x^{(n)}) \\ f(x_*^{(1)}) \\ \vdots \\ f(x_*^{(n*)}) \end{bmatrix} = N \left(\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} K(x^{(1)}, x^{(1)}) & \cdots & K(x^{(1)}, x_*^{(n*)}) \\ \vdots & \ddots & \vdots \\ K(x_*^{(n)}, x^{(1)}) & \cdots & K(x_*^{(n)}, x_*^{(n*)}) \end{bmatrix} \right)$$

nonterm?

(S+K)H = p+q

↓ PSD (Mercer Theorem)

$$\begin{bmatrix} \vec{f} \\ \vec{f}_* \end{bmatrix} = N \left(0, \begin{bmatrix} K(x, x) & K(x, x_*) \\ K(x_*, x) & K(x_*, x_*) \end{bmatrix} \right)$$

$$y = f(x) + \epsilon \quad \vec{\epsilon} \sim N(0, \sigma^2 I).$$

$$\begin{bmatrix} Y \\ Y_* \end{bmatrix} = \begin{bmatrix} f \\ f_* \end{bmatrix} + \begin{bmatrix} \epsilon \\ \epsilon_* \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K(X, X) + \sigma^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) + \sigma^2 I \end{bmatrix} \right)$$

Posterior Predictive distribution

$$Y_* | Y, X, X_* \sim N(\mu_*, \Sigma_*)$$

$$\mu_* = K(X_*, X) [K(X, X) + \sigma^2 I]^{-1} Y.$$

$$\Sigma_* = K(X_*, X_*) + \sigma^2 I - K(X_*, X) [K(X, X) + \sigma^2 I]^{-1} K(X, X_*)$$

Basically Kernel functions and Covariance functions are essentially the same.

$$\{x_i^\top\} \quad d + \left(w_i^\top \sum_{i=1}^n \right) \leftarrow$$

$$d + w^\top x = \Sigma \leftarrow$$

Algorithm for maximum likelihood estimation
using gradient descent method

Efficiency of the estimator depends on the

Lecture 10 : Deep Learning - I

- # In Deep learning the neural networks gives a way in which features themselves are learned automatically.
- # Models we have seen before given features we construct linear models. With neural networks not only are we learning that model given the features but we are also learning what the right set of features are.

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ \sin(x) \end{bmatrix}$$

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

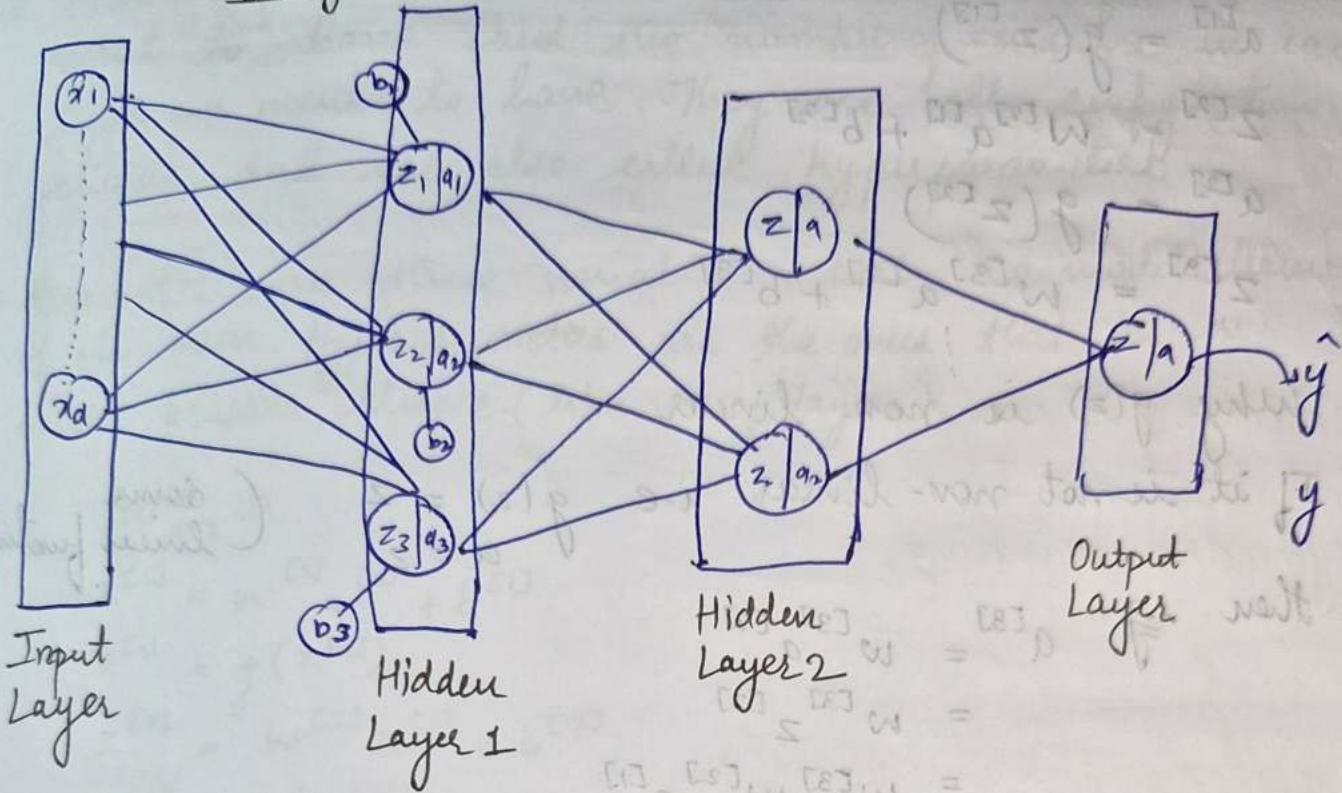
$$a = g(z)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\Rightarrow \left(\sum_{i=1}^d x_i w_i \right) + b \quad \begin{matrix} \{\theta^T x\} \\ \downarrow w, b \end{matrix}$$

$$\Rightarrow z = x^T w + b$$

- # This is just logistic regression in the example, we will just find \hat{y}, y then loss and perform gradient descent (table)
- # But in Neural networks we do it differently.

Fully connected Network

$w_{ij}^{[l]}$ - connection (It is a weight matrix)

$b_i^{[l]}$ - bias

$z_i^{[l]} = w \cdot c + b$ (linear combination)

$a_i^{[l]} = g(z)$

$x = a^{[0]}$

$$z_1^{[0]} = \sum_j w_{1j}^{[0]} a_j^{[0]} + b_1$$

No. of neuron in $l-1^{\text{th}}$ layer

$$z_2^{[0]} = \sum_j w_{2j}^{[0]} a_j^{[0]} + b_2$$

No. of neuron in l^{th} layer

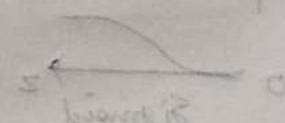
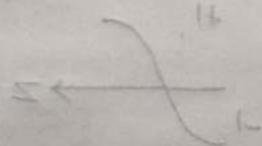
$$a_1^{[1]} = g(z_1^{[0]})$$

In vector notation

$$z^{[0]} = \underbrace{w^{[0]} a^{[0]}}_{3 \times 3} + \underbrace{b^{[0]}}_3$$

$= \text{dot product} = f$

$$\frac{1}{x+1} = \beta$$



$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

Why $g(z)$ is non linear

If it is not non-linear i.e. $g(z) = z$

then e.g.

$$a^{[3]} = w^{[3]} a^{[2]}$$

$$= w^{[3]} z^{[2]}$$

$$= w^{[3]} w^{[2]} a^{[1]}$$

$$= w^{[3]} w^{[2]} z^{[1]}$$

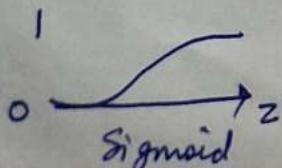
$$a^{[3]} = w^{[3]} w^{[2]} \cdot w^{[1]} \cdot x = \tilde{w}x$$

If g is a linear function then the entire network can now be represented as a single matrix, which means we haven't gone beyond the scope of linear models.

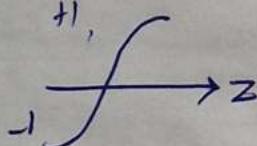
So g being non-linear is essential to have depth be meaningful. Otherwise no matter how many levels deep the network is, it can always be collapsed into a single matrix and it will effectively be as expressive as a single layer network.

Commonly used $g(z)$

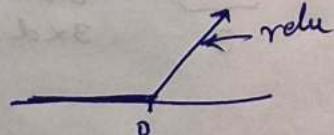
$$g = \frac{1}{1+e^{-z}}$$



$$g = \tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



$$g = \text{ReLU} = \max(z, 0)$$



From the fully connected network - The no. of layers we want to have and the number of neurons in each layer we want to have. They are all configuration choices and are also called hyperparameters.

Parameters are those variables which the model learns of its own. Hyperparameters are the ones that we as model builders choose. (like no. of layers).

$$\# \quad a^{[0]} = x^{(i)}$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

$$\vdots \quad \vdots$$

$$a^{[L]} = g(z^{[L]})$$

$$\hat{y}^{(i)} = a^{[L]}$$

$$\Rightarrow L(y^{(i)}, \hat{y}^{(i)}) = -[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

For l in $1, 2, \dots, L$

$$\left. \begin{aligned} W^{[l]} &= W^{[l]} - \alpha \frac{\partial L}{\partial W^{[l]}} \\ b^{[l]} &= b^{[l]} - \alpha \frac{\partial L}{\partial b^{[l]}} \end{aligned} \right\} \text{Update}$$

Backpropagation mathematically is just a chain rule. Algorithmically, we calculate it in a way such that it is memory efficient and we reuse a lot of the computation when calculating.

Wouldn't all the neurons learn the same thing in the first layer? Aren't we learning multiple copies of same logistic regression?

- # If we perform an initialization where entire network is initialized to 0, Then it will happen (Multiple copies).
- # In order to avoid this problem we perform Random Initialization. It is a necessary step also called symmetry breaking.

Xavier/ne
Initialization

$$w_{ij}^{[l]} \sim N(0, \sqrt{\frac{2}{n^{[l]} + n^{[l-1]}}})$$

OR

$$w_{ij}^{[l]} \sim \text{Unif} [-0.1, 0.1]$$

} Initialize w, b

$$w_{ij}^{[l]} = w_{ij}^{[l]} - \alpha \frac{\partial L}{\partial w_{ij}^{[l]}}$$

For t in 1, ...

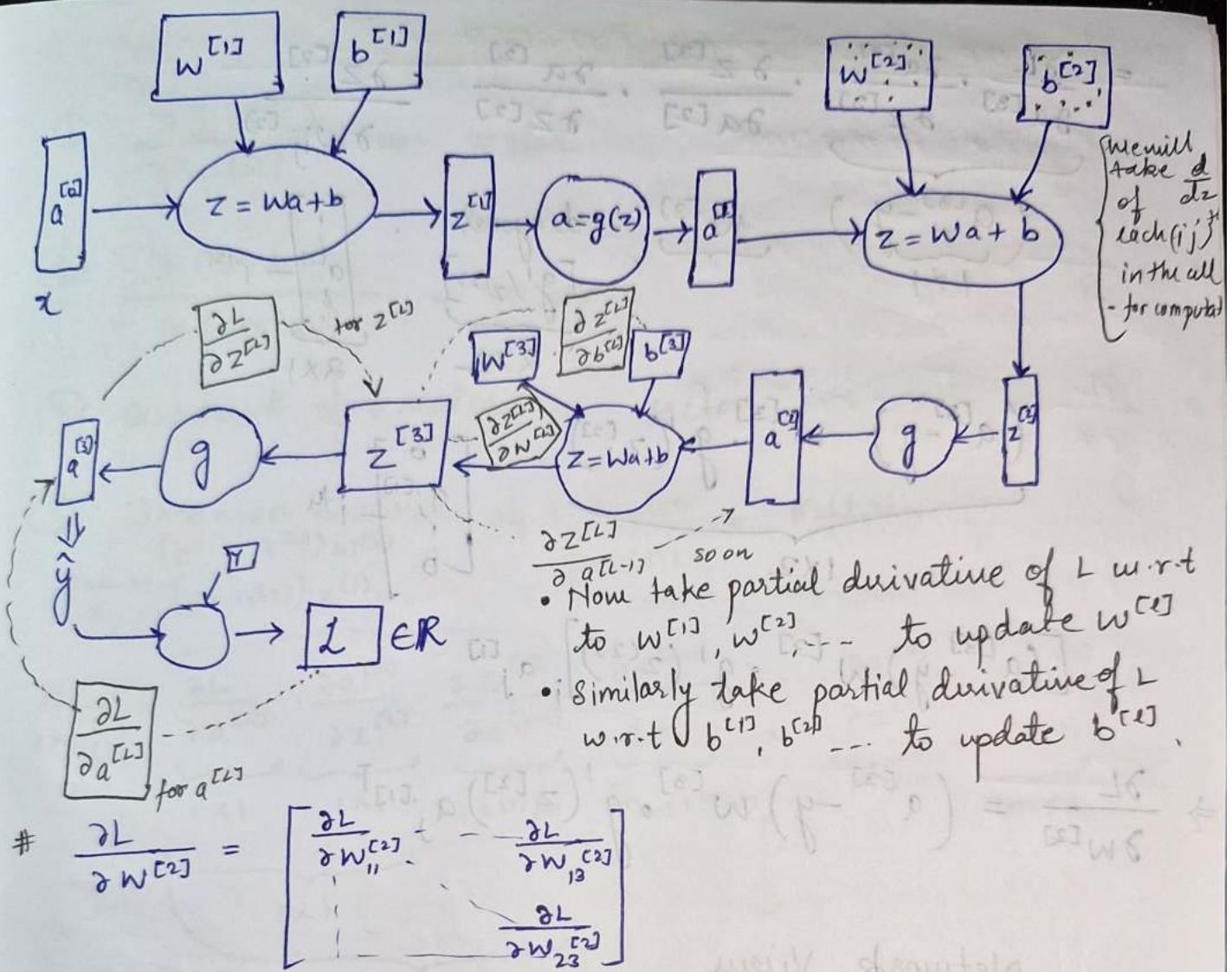
For lin 1, 2, ..., L

$$w^{[l]} = w^{[l]} - \alpha \frac{\partial L}{\partial w^{[l]}}$$

above needs a tiny alteration waiting for check

$$b^{[l]} = \underbrace{\dots}_{\text{we need to calculate this}} \text{ waiting for check}$$

we need to take a above two functions parallelly



$$\begin{aligned} \frac{\partial L}{\partial z^{[3]}} &= \frac{\partial}{\partial z^{[3]}} \left[-y \log \hat{y} - (1-y) \log (1-\hat{y}) \right] \\ &= \frac{\partial}{\partial z^{[3]}} \left[-y \log \sigma(z^{[3]}) - (1-y) \log (1-\sigma(z^{[3]})) \right] \end{aligned}$$

$$\begin{aligned} * \quad \frac{\partial L}{\partial w_{ij}^{[2]}} &= \frac{\partial L}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial w_{ij}^{[2]}} \\ &= \frac{\partial L}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial w_{ij}^{[2]}} \end{aligned}$$

$$\begin{aligned}
 &= \underbrace{\frac{\partial L}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[3]}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}}} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial w_{ij}^{[2]}}}_{\text{diagonal, } [g'/z^{[2]}]} \\
 &\quad \left(a^{[3]} - y \right) \underbrace{w^{[3]}}_{1 \times 1} \quad \left(a^{[3]} - y \right) \underbrace{w^{[3]} \cdot g'(z^{[2]})}_{1 \times 2} \\
 &\quad \left[(a^{[3]} - y) w^{[3]} \cdot g'(z^{[2]}) \right] \cdot a_j^{[1]} \\
 \Rightarrow \frac{\partial L}{\partial w^{[2]}} &= \left(a^{[3]} - y \right) w^{[3]} \cdot g'(z^{[2]}) a_j^{[1] T}
 \end{aligned}$$

Network View

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad \left[(\hat{y}_1 - y_1) \cdot \frac{\partial f}{\partial x} (\hat{y}_1) - (\hat{y}_2 - y_2) \cdot \frac{\partial f}{\partial x} (\hat{y}_2) \right] \frac{6}{56} = \frac{16}{56}$$

$$\frac{\partial f}{\partial x} \leftarrow \mathbb{R}^{m \times n} \quad (\text{Jacobian})$$

The chain rule of calculus tell us that because all these functions are compositions where the output of one becomes the input of the other and so on. We can take the local derivatives of each of these and multiply them out. This is the M.V chain rule of calculus.

$$\frac{16}{56} \cdot \frac{16}{56} \cdot \frac{16}{56} =$$

Cases

$$\textcircled{1} \quad \frac{\partial(\text{scalar})}{\partial(\text{scalar})} = \text{scalar}$$

$$\textcircled{2} \quad \frac{\partial(\text{vector})}{\partial(\text{scalar})} = \text{vector}$$

$$\textcircled{3} \quad \text{Gradient of a scalar} = \text{vector}$$

$$\textcircled{4} \quad \begin{matrix} \text{Jacobion matrix of a vector} \\ (y^{(i)} - a^{[L]}) w^{[L]} \\ (y - h_{\theta}(x)) \cdot x^{(i)} \end{matrix} = \text{matrix function}$$

$$\frac{\partial L}{\partial w_{ij}^{[L]}} = \underbrace{\frac{\partial L}{\partial a^{[L]}}}_{1 \times 1} \cdot \underbrace{\frac{\partial a^{[L]}}{\partial z^{[L]}}}_{1 \times 1} \cdot \underbrace{\frac{\partial z^{[L]}}{\partial a^{[L-1]}}}_{1 \times m_L} \cdot \underbrace{\frac{\partial a^{[L-1]}}{\partial z^{[L-1]}}}_{m_{L-1} \times m_{L-2}} \cdots \underbrace{\frac{\partial a^{[2]}}{\partial z^{[2]}}}_{m \times m} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial w_{ij}^{[2]}}}_{m \times 1}$$

1×1
 $1 \times m_{L-1}$

The chain of Jacobians will always condense to a scalar.

Backpropagation is an algorithm that tells us, what is the most optimal way to compute these gradients by reusing the intermediate computations as much as possible.

$$\boxed{id + \rho \beta W_2^T = 15}$$

Lecture 11 : Deep Learning - II

Logistic regression

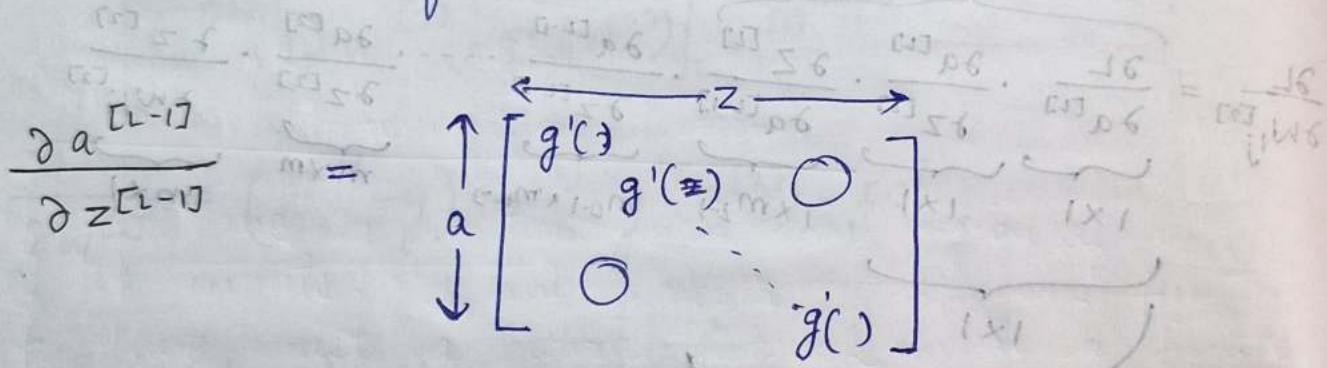
(value) 6
(value) 8

$$\# \frac{\partial L}{\partial \theta} = [y - h_{\theta}(x)] \cdot x$$

$$\text{value} = \frac{(\text{value}) 6}{(\text{value}) 8}$$

$$\# \frac{\partial L}{\partial x} = [y - h_{\theta}(x)] \theta \quad \text{value for bias term}$$

Therefore, derivative of loss w.r.t a will give us
and derivative of loss w.r.t w then a will show up.



$$q_i = g(z_i)$$

$$\frac{\partial q_j}{\partial z_i} = 0 \quad \text{if and only if } i = j$$

$$\# \frac{\partial z^{[l]}}{\partial a^{[l-1]}} = w^{[l]}$$

$$\Leftrightarrow \frac{\partial z^{[l]}}{\partial w_{ij}^{[l]}} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ q_j \\ 0 \end{bmatrix}_{i^{\text{th}}}$$

$$\text{Never do } \left[\begin{array}{c} z \\ \vdots \\ z \end{array} \right] = \left[\begin{array}{c} w \\ \vdots \\ w \end{array} \right] a$$

$$\boxed{z_i = \sum_j w_{ij} a_j + b_i}$$

$z_k \rightarrow \text{The rest is 0.}$

$$\left[\begin{array}{c} b_1 \\ \vdots \\ b_K \end{array} \right] \left[\begin{array}{c} 0 \\ \vdots \\ q_j \\ 0 \end{array} \right]_{i^{\text{th}} \text{ position}} \Rightarrow b_i q_j$$

$$\frac{\partial L}{\partial w_{ij}^{[2]}} = \delta^{[2]} \cdot \begin{bmatrix} 0 \\ a_j \\ 0 \end{bmatrix}$$

$$\Rightarrow \frac{\partial L}{\partial z^{[2]}} \begin{bmatrix} 0 \\ a_j \\ 0 \end{bmatrix} = \left[\frac{\partial L}{\partial z^{[2]}} \right]_i a_j$$

$$\Rightarrow \frac{\partial L}{\partial w^{[2]}} = \frac{\partial L}{\partial z^{[2]}} (a^{[2]})^T$$

$$\Rightarrow w^{[2]} = w^{[2]} - \alpha \frac{\partial L}{\partial w^{[2]}}$$

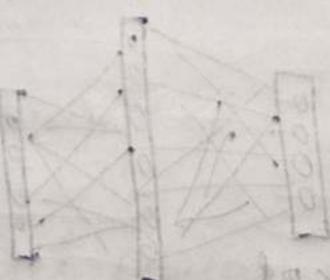
Training

- We don't always do stochastic gradient descent. We perform mini batch gradient descent. In this instead of sampling one example at random, we sample a batch of examples at random.

$$J(a, b) = \sum_{i=1}^B L(y^{(i)}, \hat{y}^{(i)})$$

B is the size of minibatch

$$\underbrace{z^{[1]}}_m = \underbrace{w^{[1]}}_{m \times d} \underbrace{x^{(i)}}_{d \times 1} + \underbrace{b^{[1]}}_{m \times 1}$$



$$= \underbrace{w^{[1]}}_{m \times d} \underbrace{[x^{(1)}, x^{(2)}, \dots, x^{(B)}]}_{d \times B} + \underbrace{b^{[1]}}_{m \times 1}$$

$$\Rightarrow m \times B + m \times 1 \rightarrow m \times B$$

- # We are trying to add two matrix of different dimensions ($m \times B$ and $m \times 1$). But this addition is not possible. The softwares that we use for training ML (tensorflow, Numpy, Pytorch etc). All these work in this scenario of adding $(m \times B$ to $m \times 1)$ with the help of Broadcasting.

$$(m \times B + (m \times 1) \rightarrow m \times B)$$

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 16 \\ 56 \end{pmatrix} = \begin{pmatrix} 16 \\ 56 \end{pmatrix}$$

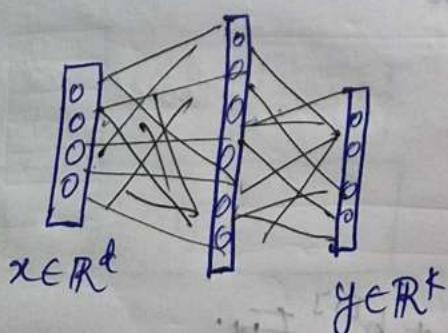
- # Every Neural network comes explicitly with a Kernel $\Rightarrow \text{Ker}(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(j)})$

Universal Approximation theorem.

If you have a function with given region of interest of the inputs. For any smooth / continuous function f , there exist a neural network with a finite number of hidden layers that can approximate f to an arbitrary degree of precision.

$$y = f(x), x \in \mathbb{R}^d, y \in \mathbb{R}^k$$

Exists a neural network with 1 hidden layer

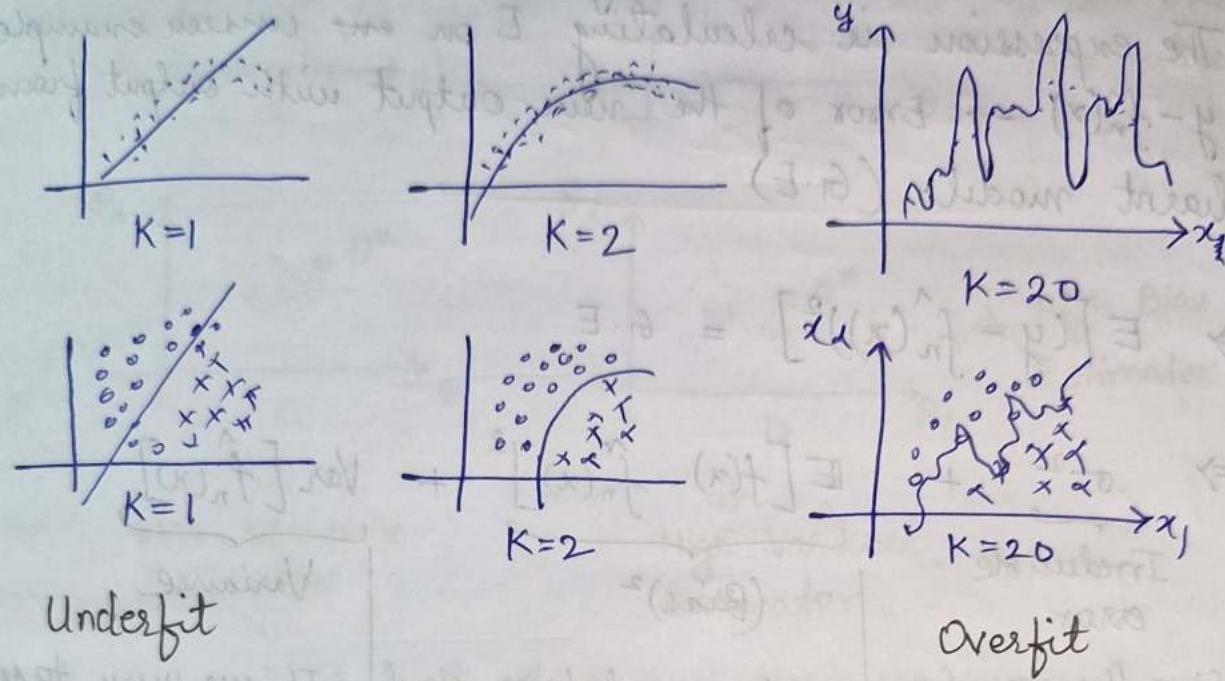


$$\text{eg } \epsilon = 10^{-6}$$

$$|\hat{f}(x) - f(x)| < \epsilon$$

- # However it doesn't tell you that how to recover that function.

Lecture 12: Bias, Variance & Regularization



Generalization Error (G.E.)

- # Bias: Component of G.E due to "Expressivity handicap".
- # Variance: Component of G.E due to finite sample of training set.

Squared error

$$y^{(i)} = f(x^{(i)}) + \varepsilon^{(i)}$$

$$E[\varepsilon] = 0$$

$$V[\varepsilon] = \sigma^2$$

$$\text{Test error [G.E.]} = E[(y - \hat{f}_n(x))^2]$$

- # here x, y are pair from the test set (unseen example).
- $\hat{f}_n(x)$ → It is the model we get by fitting on n training examples. n training examples were obtained using $g^{(i)} = f(x^{(i)}) + \varepsilon^{(i)}$.
- Thus it is a Random variable.

* $\mathbb{E}[(y - \hat{f}_n(x))^2]$: Expectation is over all x in training set and test examples.

* The expression is calculating E on one unseen example.
 $y - \hat{f}_n(x) \rightarrow$ Error of the new output with output from learnt 'model'. ($G \cdot E$)

$$\Rightarrow \mathbb{E}[(y - \hat{f}_n(x))^2] = G \cdot E$$

$$\Rightarrow \underbrace{\sigma^2}_{\text{Irreducible error}} + \underbrace{\mathbb{E}[f(x) - \hat{f}_n(x)]^2}_{(\text{Bias})^2} + \underbrace{\text{Var}[\hat{f}_n(x)]}_{\text{Variance}}$$

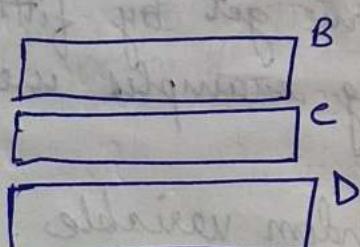
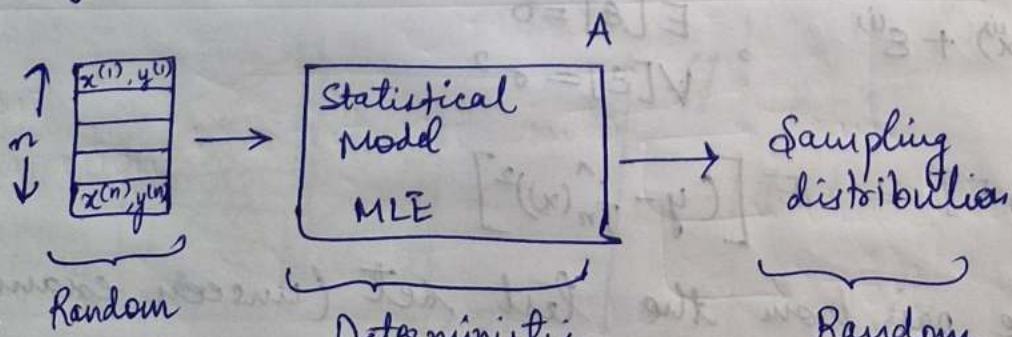
- Since the examples are noisy already

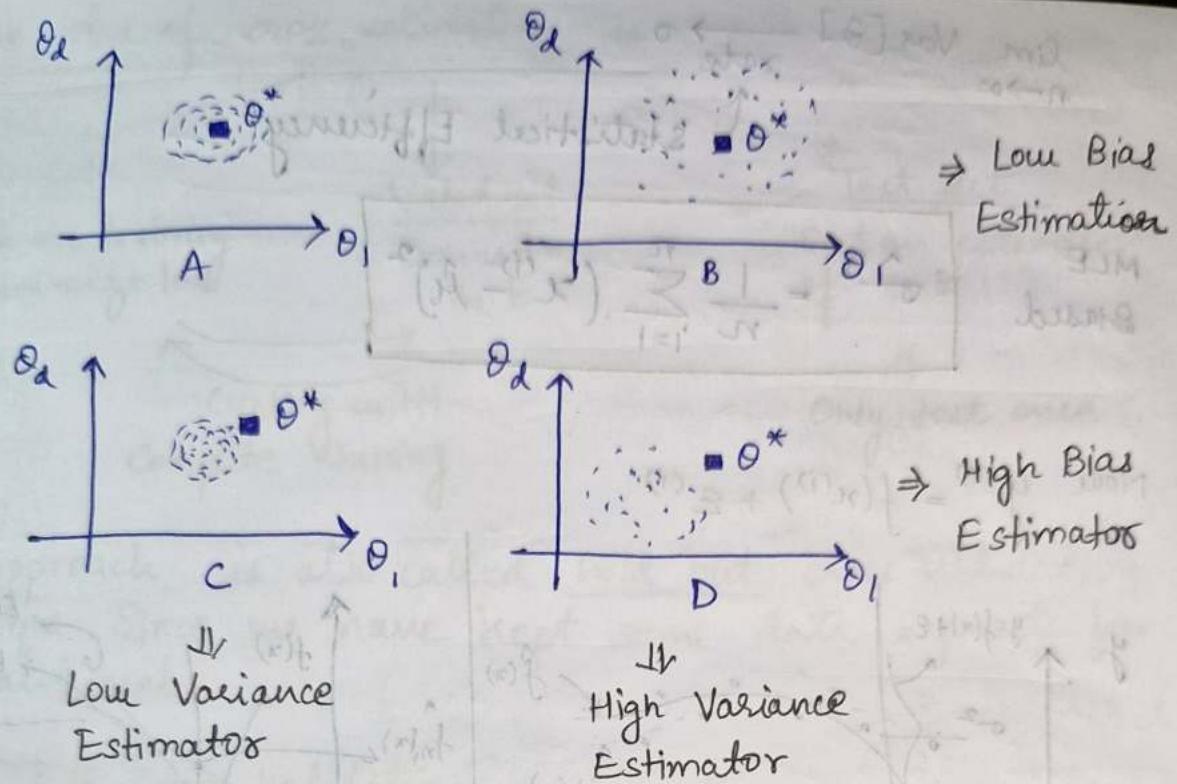
- A true underlying signal that we don't have access to. $f(x)$ will be different from the prediction made by our model
- Basically how systematically wrong we are

- If we were to repeat this example with new set of n examples. $\hat{f}_n(x) \rightarrow$ diff value
- Sensitive to noise.

Bias - Variance

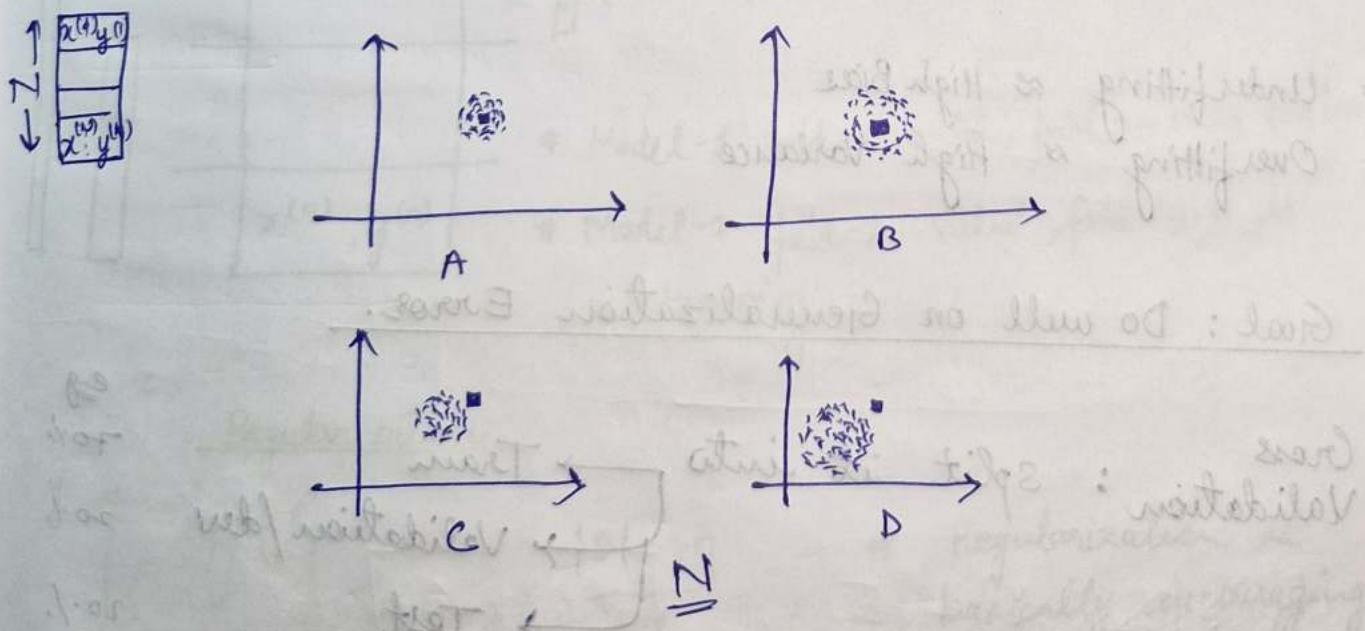
$$(x, y) \sim \text{Dist}(\theta^*)$$





This clearly suggest that Bias and variance are not correlated, knowing one doesn't give info about other.

Now with $n \rightarrow$ very high.



$$E[\hat{\theta} - \theta^*] = \text{Bias}$$

$$\text{Var}[\hat{\theta}] = \text{Variance}$$

\Rightarrow If $\text{Bias} \rightarrow 0$ as $n \rightarrow \infty$ \Rightarrow Consistent Estimation

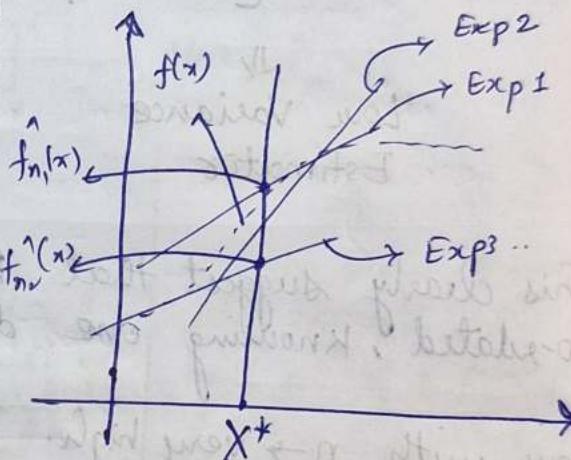
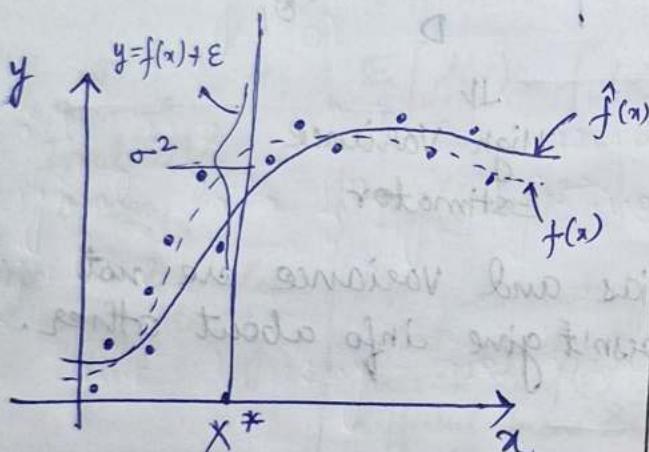
$$\lim_{n \rightarrow \infty} \text{Var}[\hat{\theta}] \xrightarrow{\text{rate}} 0$$

statistical Efficiency

- # MLE Biased

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\hat{x}^{(i)} - \bar{x})^2$$

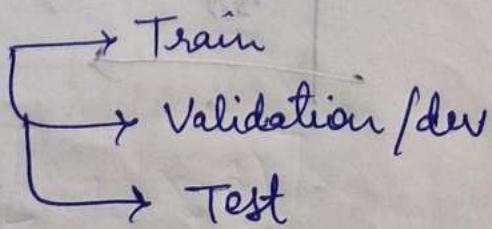
$$\# \text{ Now } y^{(i)} = f(x^{(i)}) + \epsilon^{(i)}$$



- # Underfitting \approx High Bias
- # Overfitting \approx High Variance

Goal : Do well on Generalization Error.

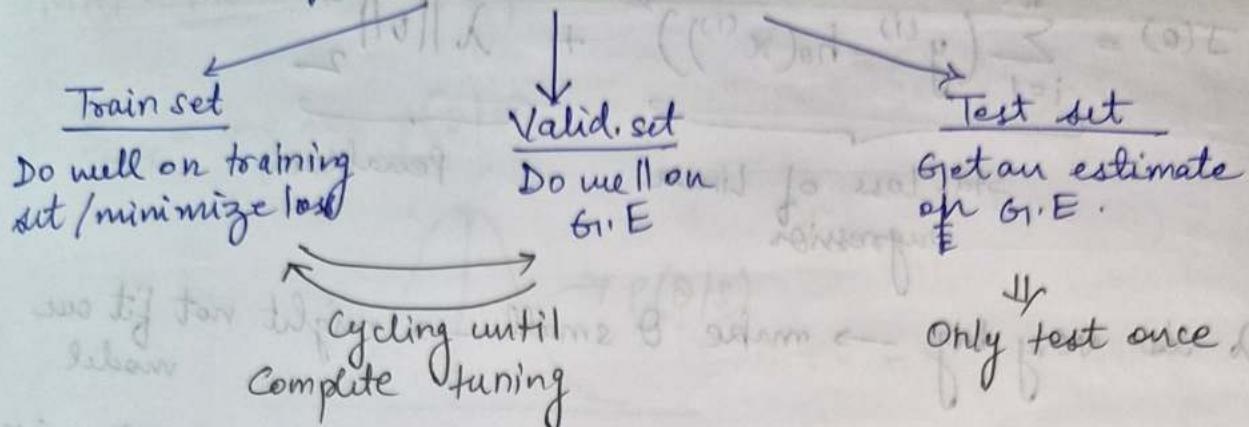
Gross Validation : split it into



eg
70%
20%
20%

- * No matter what your hyperparameter decision is, the process of cycling between training and validation set is essential. Universal ML concept.

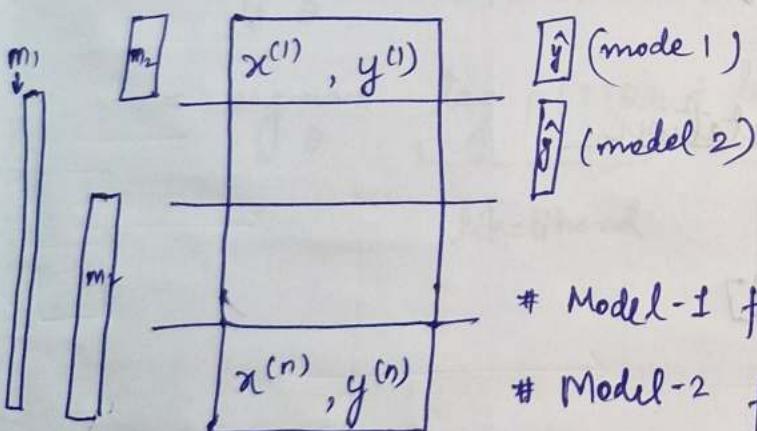
The idea of cross validation is



This approach is also called hold out cross validation technique. Since we have kept some data separate for test and val.

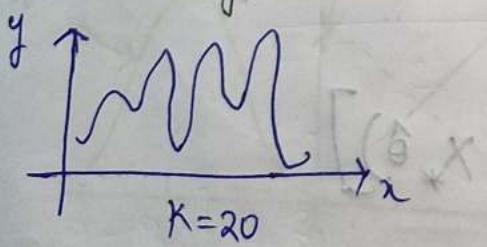
K-fold Cross Validation

- Small Data sets.



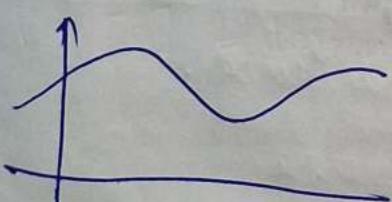
- # Model-1 fold-1 valid, fold 2-n is training
- # Model-2 fold-2 valid, fold 3-n is training

Regularization



$$\underbrace{|\theta| > 0}_{\text{Bad}}$$

- # Regularization is basically encouraging small $\|\theta\|$
- # Smoother curve



$$J(\theta) = \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)}))^2 + \lambda \|\theta\|_2^2$$

+ Std. loss of Linear regression Penalize

- # λ is very big \rightarrow make θ small \rightarrow might not fit our model
- # λ is very small \rightarrow make θ big \rightarrow model fit data without any regularization.
- # λ acts as a relative weight factor between the original data fitting objective and the regularization objective

We tune λ using cross-validation.

Bayesian Interpretation

$\theta \sim p(\theta)$ [Prior]

$s|\theta \sim$ [likelihood]

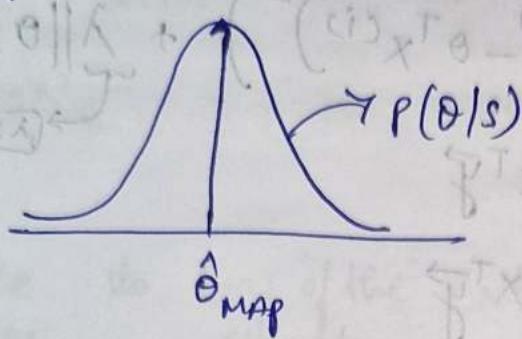
$\theta|s$ [Posterior]

$$P(\theta|s) = \frac{P(s|\theta) \cdot P(\theta)}{P(s)}$$

$$P(Y_* | X_*, X, Y) = E_{\theta \sim P(\theta|s)} [P(y_* | X_*, \theta)]$$

MAP - Maximum-a-posteriori Parameter Estimate

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | s)$$



$$\hat{\theta} = \arg \max_{\theta} p(\theta | s)$$

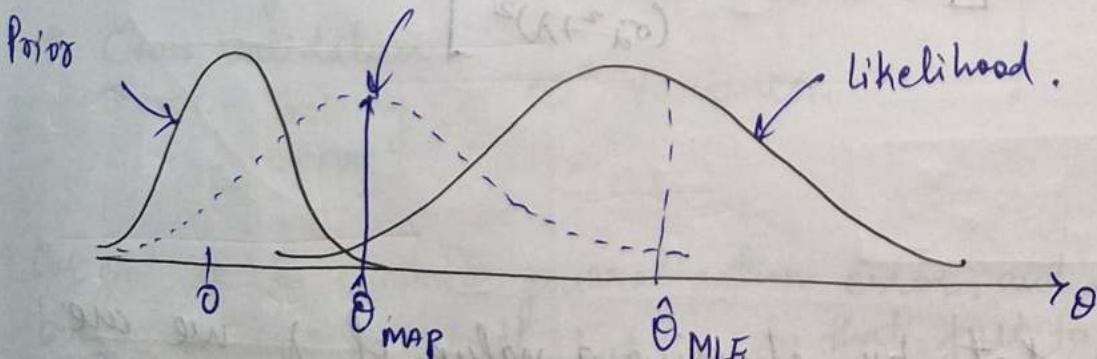
$$p(\theta) \sim N(\sigma^2) = \arg \max_{\theta} \frac{p(s|\theta) \cdot p(\theta)}{p(s)}$$

$$= \arg \max_{\theta} P(s|\theta) \cdot P(\theta)$$

$$= \arg \max_{\theta} \underbrace{\log p(s|\theta)}_{\text{likelihood}} + \underbrace{\log p(\theta)}_{\lambda \|\theta\|^2}$$

$$p(\theta) = \frac{1}{\sqrt{2\pi \cdot \sigma^2}} \exp \left[-\frac{1}{2} \frac{(\theta - \mu)^2}{\sigma^2} \right]$$

Prior x Likelihood



Lecture 13: Statistical learning Uniform Convergence

L_2 - Regularized Linear Regression

$$J(\theta) = \left(\sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 \right) + \lambda \|\theta\|_2^2$$

we know $\hat{\theta} = (x^T x)^{-1} x^T \vec{y}$

$$\hat{\theta}_n = (x^T x + \lambda I)^{-1} x^T \vec{y}$$

$$(x^T x + \lambda I)^{-1} = U \begin{bmatrix} (\sigma_1^2 + \lambda)^{-1} & & \\ & \ddots & \\ & & (\sigma_n^2 + \lambda)^{-1} \end{bmatrix} U^T$$

$$E[\hat{\theta}_n] = \left[U \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_n^2}{\sigma_n^2 + \lambda} \end{bmatrix} U^T \right] \theta^*$$

For Variance

$$\text{Cov}(\hat{\theta}_n) = U \begin{bmatrix} \frac{\sigma_1^2}{(\sigma_1^2 + \lambda)^2} & & \\ & \ddots & \\ & & \frac{\sigma_n^2}{(\sigma_n^2 + \lambda)^2} \end{bmatrix} U^T$$

$$\epsilon \sim N(0, \sigma^2)$$

$$y = \theta^T x + \epsilon$$

We can see that by choosing a value of λ we are doing trade off between bias and variance.

$\lambda \uparrow \rightarrow \text{Var} \downarrow, \text{Bias} \uparrow$ and vice versa

Bias Variance decomposition

$$\text{MSE}[\hat{f}_n] = \underbrace{\gamma^2}_{\text{Irreducible error}} + \underbrace{E[\hat{f}_n(x^*) - f(x^*)]^2}_{\text{Bias}^2} + \underbrace{V[\hat{f}_n(x^*)]}_{\text{Variance}}$$

$$f(x) = \theta^{*T} x$$

Now

Bias of the predictor \rightarrow Bias of the estimator

$$\begin{aligned}\text{Bias}(\hat{f}_n(x^*)) &= E[\hat{f}_n(x^*) - f(x^*)] \\ &= E[\hat{\theta}_n^T x^* - \theta^{*T} x^*] \\ &= E[\hat{\theta}_n - \theta^*]^T x^*\end{aligned}$$

$$\boxed{\text{Bias}(\hat{f}_n(x^*)) = \text{Bias}(\hat{\theta}_n)^T x^*}$$

$$\boxed{\text{Var}[\hat{f}_n] = x^T \text{Cov}[\hat{\theta}_n] x}$$

Heuristics for Bias and Variance

# Training Error	\approx Bias	$\left. \begin{array}{l} \approx \text{Appox} \\ \neq \text{Not equal} \end{array} \right\}$
# Cross Validation Error - Training error	\approx Variance.	

→ Our goal is to minimize generalization error, and not to go after bias or variance by itself but we do that just to minimize G.E. In order to do that, we have rough heuristic of bias and variance. Compare the two and go after the one which is larger

→ We do not want to take any arbitrary step like get more data, add more regularization. Therefore we construct heuristic estimation.

To fight Bias

- Add more layers to Neural network (Make model larger)
- Reduce regularization
- Add more features.

To fight Variance

- Increase regularization
- Collect more data
- Simpler model

Statistical Learning theory

- ① Train and Test data \sim Same distribution (Assumption).
- ② Examples are sampled IID

Risk of hypothesis

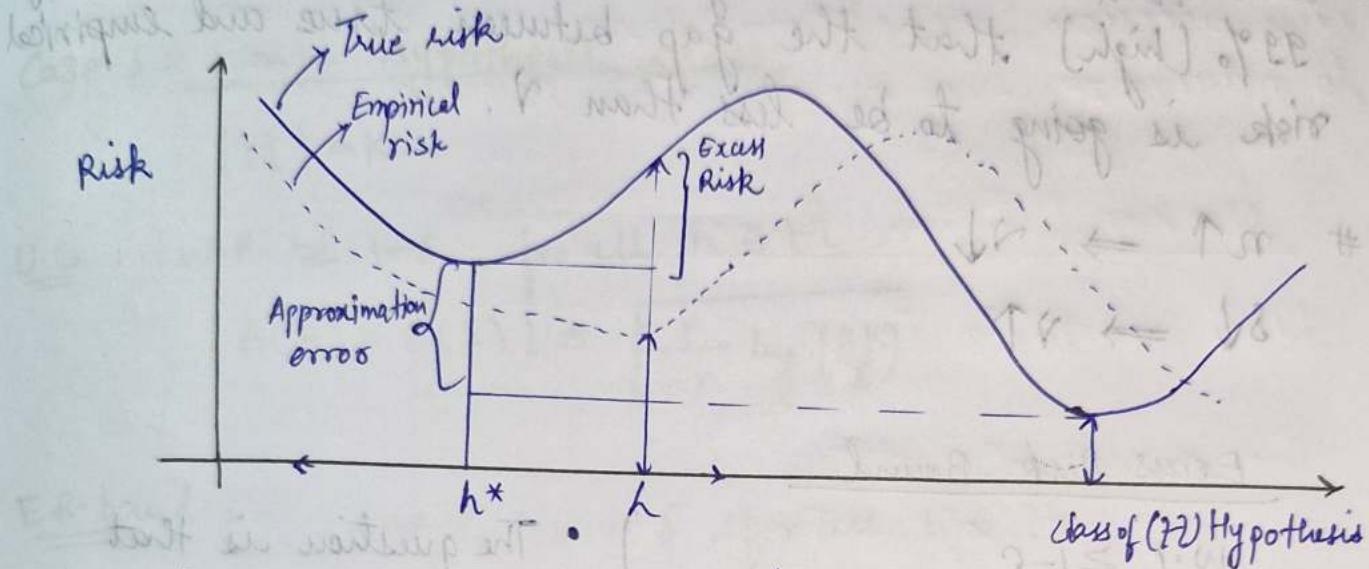
$$E(h) = \mathbb{E}_{(x,y) \sim D} [\text{loss}(y, h(x))]$$

$$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$$

Empirical Risk

$$\hat{E}(h) = \frac{1}{|S|} \sum_{(x,y) \in S} \text{loss}(y, h(x))$$

- Q) For a given h , what is the relation between $\hat{\epsilon}(h)$ and $\epsilon(h)$.
- Q) How does the generalization error compare to the best possible generalization error.



- Hypothesis - Set of all hypothesis for all possible models.
- h^* - best in class hypothesis.

Step 1: Uniform Convergence.

Step 2: Excess Risk Bound.

Uniform Convergence : How the Empirical risk function as we increase the number of data points we have will probabilistically converge to true risk function uniformly across all points.

We want to minimize the gap uniformly across all hypothesis.

Statement : With probability (W.P) \Rightarrow W.P = 1 - δ
for all $h \in \mathcal{H}$

$$|\epsilon(h) - \hat{\epsilon}(h)| \leq \underbrace{\gamma}_{\text{Term}(n, \delta, \mathcal{H})}$$

For all hypothesis on the H axis the gap of $|E(h) - \hat{E}(h)|$ is less than γ . If we are given a training set of size n , degree of confidence (δ) and γ then we will get some γ with probability 99% (high) that the gap between true and empirical risk is going to be less than γ .

$$\# n \uparrow \Rightarrow \gamma \downarrow$$

$$\delta \downarrow \Rightarrow \gamma \uparrow$$

Excess Risk Bound

$$W.P \geq 1-\delta$$

$$E(\hat{h}) \leq E(h^*) + 2\gamma$$

Proof

$$E(\hat{h}) \leq \hat{E}(\hat{h}) + \gamma, (\text{at } h)$$

$$\leq \hat{E}(h^*) + \gamma, [\text{By definition of E.P.M}]$$

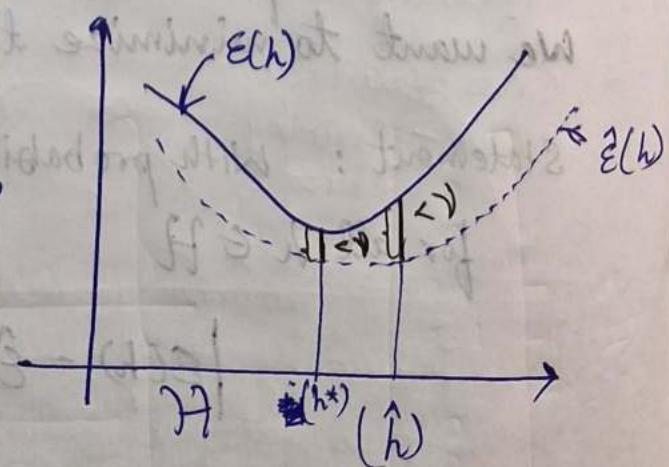
$$\leq [E(h^*) + \gamma] + \gamma$$

$$, [\text{U.C at } h^*]$$

$$\leq E(h^*) + 2\gamma$$

$$\boxed{E(\hat{h}) \leq E(h^*) + 2\gamma}$$

- The question is that uniform convergence gives the difference between $E(h)$ and $\hat{E}(h)$ i.e thick line (-) and dotted line (----).
- But in Risk bound we need to find difference between both the points on $E(h)$ i.e thick line (-)



This excess risk bound tells us how well our model is gonna perform in the real world relative to the best possible model that could have performed in the real world.

Case I: Finite Hypothesis class

$$|\mathcal{H}| = K$$

V.C : w.p. $\geq 1 - \delta$, for all $h \in \mathcal{H}$

$$|\hat{\epsilon}(h) - \hat{\epsilon}(h)| \leq \sqrt{\frac{1}{2n} \log\left(\frac{2K}{\delta}\right)}$$

E.R. bound :

w.p. $\geq 1 - \delta$, for all $h \in \mathcal{H}$

$$\hat{\epsilon}(h) \leq \underbrace{\hat{\epsilon}(h^*)}_{\min_{h \in \mathcal{H}} \hat{\epsilon}(h)} + 2 \sqrt{\frac{1}{2n} \log\left(\frac{2K}{\delta}\right)}$$

1+	2+	3+	4+	5+	6+
1-	-	-	-	-	-
-	-	-	-	-	-
-	-	-	-	-	-
-	-	-	-	-	-

	0	1	2	3	4	5
$(h, 1)$ -	0	1.0	1.0	1.0	1.0	1.0
$(S, 1)$ -	1.0	0	1.0	1.0	1.0	1.0
$(S, 2)$ -	1.0	1.0	0	1.0	1.0	1.0

$$\{(1, 1), (2, 1), (1, 2)\} = 2$$

$$|H| = |S|$$

$$\{3, 4, 2, 1\} = A$$

$$\text{by } \frac{1}{2}$$

Lecture 14: Reinforcement Learning - I

In reinforcement learning we are making multiple decisions, one decision at a time and with each decision we obtain some reward. Our goal is to maximize this accumulated reward into the future.

Markov Decision Process (MDP)

It is a tuple $(S, A, \{P_{sa}\}, \gamma, R)$

S - a set of states

A - a set of actions

P_{sa} - Transition probabilities

$\gamma \in [0, 1]$ - Discount factor

$R: S \times A \rightarrow R$

$R: S \xrightarrow{\text{or}} R$ - Reward function

Ex

	-0.02	-0.02	-0.02	+1
3	-	-	-	-1
2	-	-	-	-1
1	-	-	-	-
	1	2	3	4

$$S = \{(1,1), (1,2), \dots, \}$$

$$|S| = 11$$

$$A = \{N, S, W, E\}$$

$$P_{sa} \text{ and } \gamma$$

$$P_{(1,3)N} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.1 \\ 0.1 \\ 0.8 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{array}{l} (1,4) \\ (1,2) \\ (2,3) \end{array}$$

The reward function can be represented as a function of state

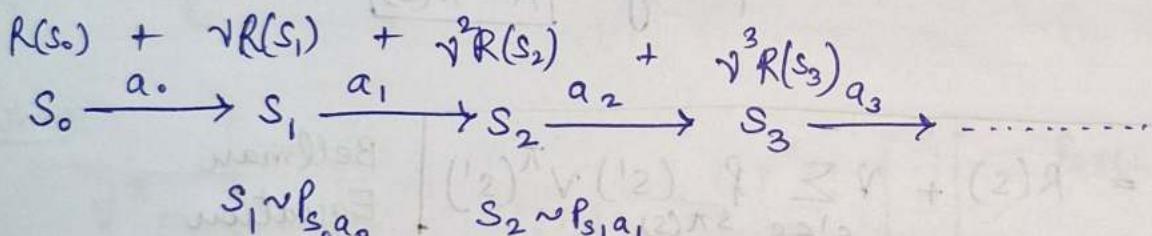
$$R(s) \rightarrow R$$

or it can be modelled as function $R(s, a)$

$$R(s, a) \rightarrow R$$

→ The reason is because taking different actions we have different costs, so that is why both s, a becomes important.

Trial / Episode / Trajectory



We want to get larger rewards sooner and negative rewards later.

Policy

Policy $\pi : S \rightarrow A$

Value $V^\pi : S \rightarrow R$

$$V^\pi(s) = \mathbb{E} [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s, \pi]$$

* Value can be defined as the long term reward we gain following a policy.

$s_0 \rightarrow$ not random

$s_1, s_2, \dots \rightarrow$ Random

$$V^\pi(s) = R(s) + \gamma E \left[R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots \right]$$

(2) \downarrow

$$V^\pi(s)$$

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

→ Bellman Equation

↓
Immediate reward

↓
Discount factor times expectation of future value of each particular state times the probability that we end up in that particular state by taking the action.

The action itself is based on the policy. $A = \pi(s)$

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Bellman Equation

Given Policy, how do we find Value

Policy
 $s \rightarrow A$

Value
 $s \rightarrow R$

① $\pi \rightarrow V^\pi$

② $\pi \leftarrow [R = \gamma \left(\dots + \gamma V^\pi \right)] \quad \boxed{V^\pi = \gamma R + \gamma V^\pi}$

Case I: $\pi \rightarrow V^\pi$ (Policy to Value)

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)} V^\pi(s')$$

↓ Known ↓ Solving for

$$V^\pi(s_1) = R(s_1) + \gamma \sum_{s' \in S} P_{s,\pi(s_1)}(s') V^\pi(s')$$

$$V^\pi(s_2) = R(s_2) + \gamma \sum_{s' \in S} P_{s,\pi(s_2)}(s') V^\pi(s')$$

$$V^\pi(s_{|S|}) = \dots$$

$$\frac{V^\pi}{|S|} = \frac{R}{|S|} + \frac{\gamma P^\pi}{|S| \times |S|} V^\pi$$

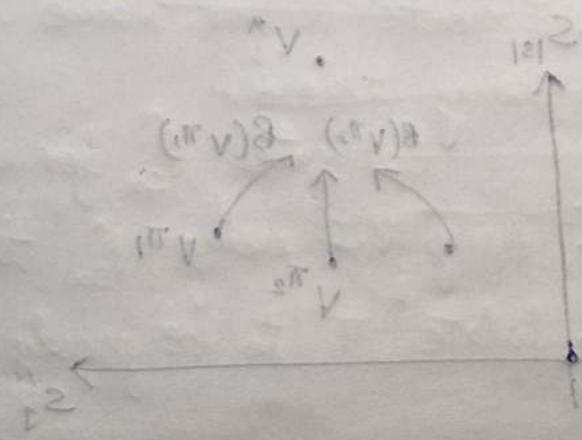
$$\Rightarrow V^\pi = [I - \gamma P^\pi]^{-1} R$$

Where

$$V^\pi = \begin{bmatrix} & s_1 & V^\pi(s_1) \\ & s_2 & V^\pi(s_2) \\ \vdots & \vdots & \vdots \\ & s_i & V^\pi(s_i) \end{bmatrix}$$

$$P^\pi = \begin{bmatrix} s_1 & P_{s_1,\pi(s_1)} & \dots \\ s_i & P_{s_i,\pi(s_i)} & \dots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

- # This approach gives us a way of starting from a policy π and calculating the long term values for each of the states if we follow the policy at all times.



Optimal Value function

$$V^*(s) = \max_{\pi} V^\pi(s)$$

$$V^*(s) = R(s) + \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s').$$

optimal policy

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

Algorithm used for computation

Value iteration algorithm:

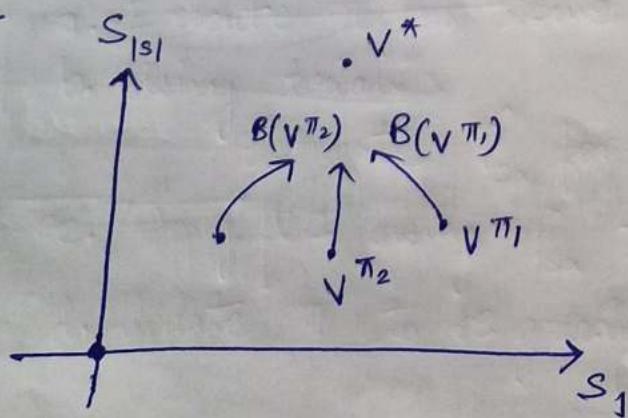
1. For each step/state s , initialize $V(s) = 0$
2. Repeat until convergence:

$O \rightarrow o \text{order}$ { For every state s , update }

$$O(|S|^2 |A|) \quad V(s) = R(s) + \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$$

(Bellman Backup operator)

Intuition



- # Bellman operator is called a contraction mapping because for any two inputs to the operator, the corresponding outputs will be closer to each other than the corresponding two inputs.
- # The point to which they all are converging to (V^*) is called a fixed point.

Policy iteration algorithm (Greedy Policy)

1. Initialize π randomly
2. Repeat until convergence

$$(a) \text{ Set } V = V^\pi \leftarrow O(|S|^3)$$

(b) For each state s , set

$$O(|A| \cdot |S|) \left\{ \begin{array}{l} \pi(s) = \arg \max_{a \in A} \sum_{s'} p_{sa}(s') V(s') \end{array} \right.$$

- # For large state spaces, policy iteration can be much more expensive than value iteration. However, with policy iteration after a finite no. of steps, we get the exact optimal value. Whereas with value iteration we get closer and closer to V^* (optimal value) but we may never reach such it.

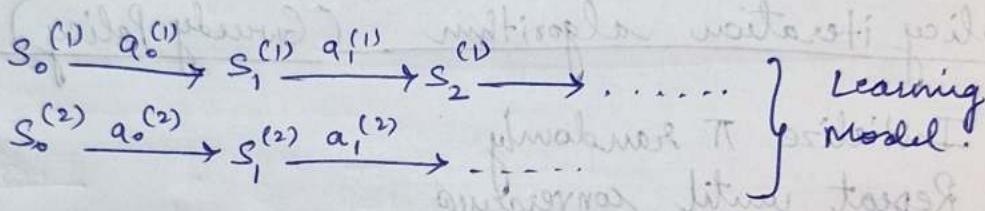
Lecture 15: Reinforcement Learning - II

P_{sa} is unknown

i. Model Based vs Model-free

$$\text{Model} = P_{sa}$$

Trials



$\hat{P}_{sa}(s')$ = No. of times we took action a at state s and got to state s'

No. of times we took action a at state s ,

$$= \frac{0}{0} \Rightarrow \text{uniform}$$

Algorithm (with Model learning)

1. Initialize π randomly

2. Repeat of

a. Execute π for some time

b. Using trials

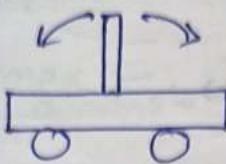
Estimate \hat{P}_{sa} for each s, a

c. Apply V.I using \hat{P}_{sa} to get updated $V^{(t)}$

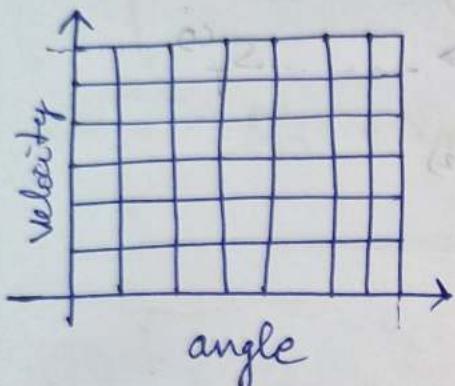
d. Update π to be greedy w.r.t $V^{(t)}$

Continuous States

$s \rightarrow \text{Space}$



: Balancing a stick
on a moving cart



- Discretization help us to solve continuous problems.
- But has limitations.

This is discrete now

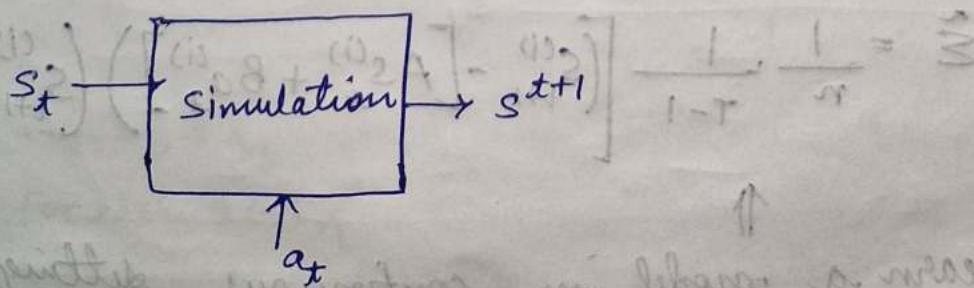
Curse of Dimensionality

12 parameters = $x, y, z, \theta, \phi, \psi, \dot{x}, \dot{y}, \dot{z}, \dot{\theta}, \dot{\phi}, \dot{\psi}$
 coordinates angles velocities angular velocities

- Consider this as an example for flying a helicopter
 - If we discretize it in each param. to 10 values. then
- Total parameters = 10^{12} (Massive state size)
- This is a shortcoming of discretization. (Limitation)

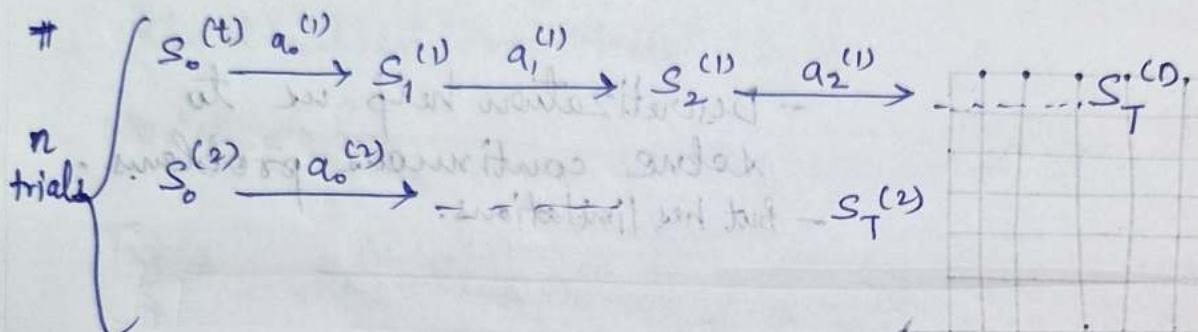
Value function Approximation

s is continuous and we don't discretize.



$$s^{(t+1)} = f(s^{(t)}, a^{(t)})$$

$$= As^{(t)} + Ba^{(t)}$$



Loss function

$$J(A, B) = \sum_{i=1}^n \sum_{t=0}^{T-1} \|s_{t+1}^{(i)} - (As_t^{(i)} + Ba_t^{(i)})\|_2^2$$

Now if we minimize it, then it will give our learnt model.

$$\hat{A}, \hat{B} = \arg \min_{A, B} J(A, B)$$

$$\Rightarrow s^{(t+1)} = \hat{A}s^{(t)} + \hat{B}a^{(t)} \rightarrow \text{Deterministic model}$$

$$s^{(t+1)} = \hat{A}s^{(t)} + \hat{B}a^{(t)} + \epsilon^{(t)} \rightarrow \text{Stochastic}$$

$$\boxed{\epsilon \sim N(0, \Sigma)}$$

where,

$$\hat{\Sigma} = \frac{1}{n} \cdot \frac{1}{T-1} \left[(s_{t+1}^{(i)} - [As_t^{(i)} + Ba_t^{(i)}]) (s_{t+1}^{(i)} - [As_t^{(i)} + Ba_t^{(i)}])^T \right]$$



Learn a model in continuous setting

Now making the value function also continuous.

$$V(s) = R(s) + \gamma \max_a E_{s' \sim P_{sa}} [V(s')]$$

$$P_{sa} = N(\hat{A}_s + \hat{B}_a; \Sigma)$$

$$V(s) = R(s) + \gamma \max_a \int P_{sa}(s') \cdot V(s') ds'$$

Gaussian
Pdf

Value function approximation

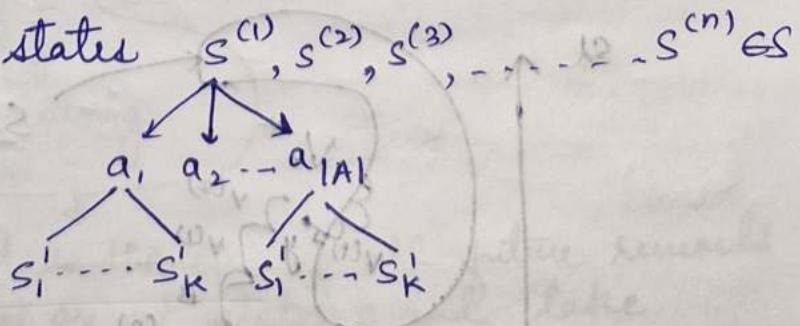
$$V(s) = \theta^T \phi(s)$$

Algorithm for Continuous state model value function approximation.

$$s^{(t+1)} = As^{(t)} + Ba^{(t)} + \epsilon$$

$$V(s) = \theta^T \phi(s)$$

1. Randomly samples n states



2. Initialize $\theta = 0$

3. Repeat {
a}

For $i=1, \dots, n$ {
b}

For each $a \in A$ {
c}

sample $s'_1, \dots, s'_K \sim P_{s^{(i)}}[s']$ $\stackrel{\text{e.g.}}{\sim} N(A s^{(i)} + B a, \Sigma)$

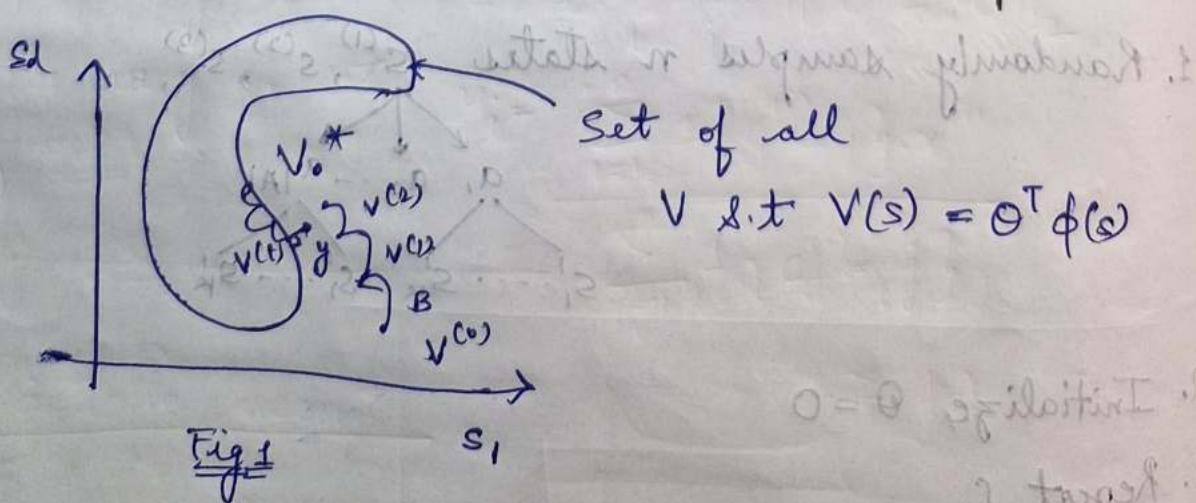
Set $q(a) = \frac{1}{K} \sum_{j=1}^K R(s^{(i)}) + \gamma V(s'_j)$ $[V(s) = \theta^T \phi(s)]$
 $\} _c$

Set $y^{(i)} = \max_a q(a)$ [label $V(s^{(i)})$]

$\} _b$ Set $\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^n (\theta^T \phi(s^{(i)}) - y^{(i)})^2$
 $\} _a$ $(\phi^T \theta) = V$

$q(s^{(i)}, a) \approx R(s^{(i)}) + \gamma \underset{s' \sim P_{sa}}{E}[V(s')]$

$y^{(i)} \approx R(s^{(i)}) + \max_a \gamma \underset{s' \sim P_{sa}}{E}[V(s')] \quad \begin{matrix} \text{Bellman} \\ \text{Backup} \\ \text{operator} \end{matrix}$



Value function based - fit function

- # In this setting we are learning the dynamics of A and B and using those dynamics we are performing value iteration using Bellman backup operator. Performing in such a way we use these fitted value function.
- # But now as can be seen in Fig 1. we are limiting ourselves in the space of value functions. Which means this whole algorithm is limited to some set of all V.

We don't know env.
but learn it

Value function based

Policy Based

Discrete finite

V.I

P.I

Infinite state
space

Fitted V.I

Model free

Q-function =

$Q(s, a) \approx \text{Value}$

↳ Here we do not care how the environment works. and also we don't care about it

$V^*(s)$

Q-learning

What is the sum of all future rewards if we are at state s and take action a.

Expected returns

$$\mathbb{E}[R_t | R_{t-1} = x] = \sum_{i=0}^{\infty} r_i p_i = (\pi, \gamma)^T \mathbf{v}_{\pi}$$

Lecture 16 - Unsupervised learning

K-means, GMM and EM.

In unsupervised learning, we are given the $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ but not the corresponding y labels.

$$x^{(i)} \in \mathbb{R}^d$$

K-Means Algorithm

$$S = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\} \quad x^{(i)} \in \mathbb{R}^d$$

K-clusters (given to us)

Steps I. V

1. Initialize cluster centroids

$$\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d \text{ randomly}$$

2. Repeat until convergence:

cluster identity - For every i , set $c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|_2^2$

Recalculate - For every j , set $\mu_j = \frac{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = j\} \cdot x^{(i)}}{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = j\}}$

(Centroid)

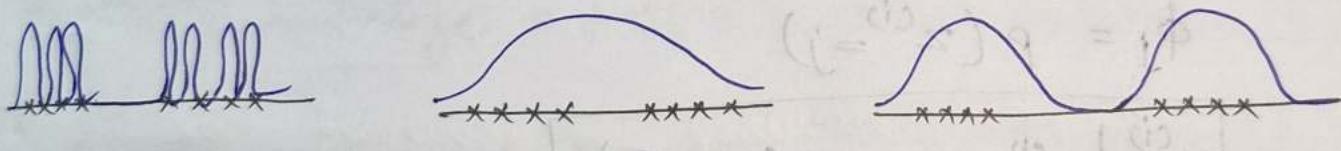
Co-ordinate descent

$$J(c, \mu) = \sum_{i=1}^n \|x^{(i)} - \mu_{c^{(i)}}\|_2^2$$

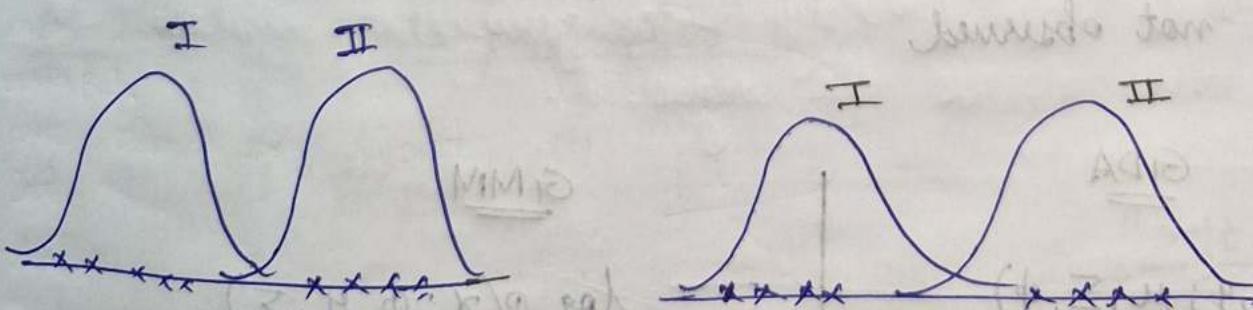
(Distortion function)

- # The step 1 is to minimize this function by holding μ fixed and optimizing c . and step 2 is to then minimizing j again by holding the c fixed and optimizing it with respect to μ .
- # K-means is co-ordinate descent of the distortion function J .
- # The μ and c we end up with can be different at different time, as it depends upon the random initialization. (Non convex)

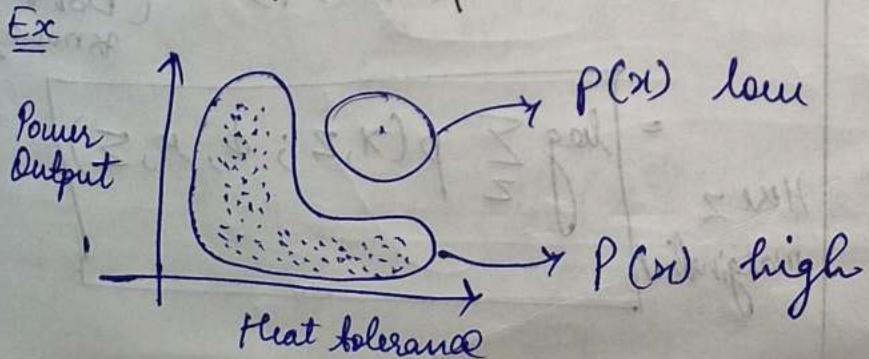
Density Estimation



Gaussian Mixture Model



To come up with $p(x)$.



Basically our objective is to construct an algorithm for mixture of gaussians and it is based purely on intuition. Then to describe this framework (EM) to derive (GMM) and check if we end up with the same algorithm.

GMM

$$S = \{x^{(1)}, \dots, x^{(m)}\} \rightarrow \text{Given training set}$$

$z^{(i)} \sim \text{Multinomial distribution } (\phi)$.

$$\phi_j \geq 0 \quad \sum_{j=1}^k \phi_j = 1$$

$$\phi_j = P(z^{(i)} = j)$$

$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \Sigma_j)$$

$z^{(i)}$ - latent variable : Random variable that we have not observed

GDA

$$\log p(x, y; \mu, \Sigma, \phi)$$

$$= l(\mu, \Sigma, \phi)$$

GMM

$$\log p(x; \phi, \mu, \Sigma)$$

$$= l(\mu, \Sigma, \phi)$$

Here z

Marginalized

$$= \log \sum_z p(x, z; \phi, \mu, \Sigma)$$

(Don't know z)

Notation for our study throughout from now.

$z \rightarrow$ Latent (unobserved)

$p(z) \rightarrow$ class prior / prior

$p(x, z) \rightarrow$ model

$p(z|x) \rightarrow$ posterior

$p(x) \rightarrow$ Evidence.

Our goal is to maximize the likelihood using the evidence.

[Inspired by K-means]

Repeat until convergence.

Randomly initialize μ, ϕ, Σ

E-step For each i, j set:

$$w_j^{(i)} = p(z^{(i)}=j | x^{(i)}; \phi, \mu, \Sigma) \quad \left\{ \begin{array}{l} \text{probability estimation} \\ p(z|x) = \frac{p(x|z) \cdot p(z)}{\sum p(x|z) \cdot p(z)} \end{array} \right.$$

M-step Update parameters

$$\phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$$

$$\Sigma_j = \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$

$$\Sigma_j = \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$

e.g. $K=3$

$$P(z^{(i)} = j | x^{(i)})$$

$$\begin{bmatrix} 0.1 \\ 0.7 \\ 0.2 \end{bmatrix} \begin{array}{l} K=1 \\ K=2 \\ K=3 \end{array} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

EM Algorithm (Expectation Maximization)

It is an algorithm where we perform MLE in the presence of latent (unknown) variables.

True model : $p(x, z; \theta)$

when z is unobserved

Maximization $l(\theta) = \log p(x; \theta)$ {
 $\underbrace{\qquad\qquad\qquad}_{z \text{ is marginalised}}$ }

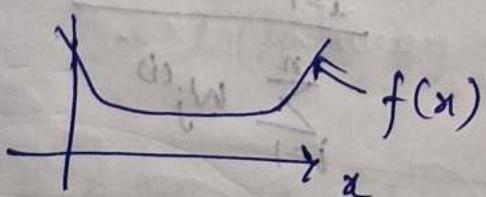
The EM gives us a framework for achieving this.
 (Indirect way)

EM algorithm is one of the key algorithm for modern Deep learning / Deep generative models. It is used across all algos.

Jensen's Inequality

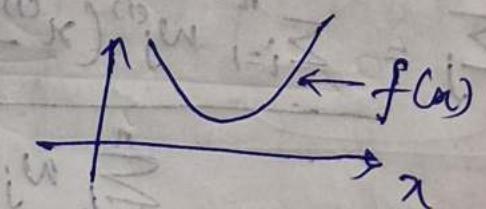
f to be convex

$$f''(x) \geq 0 \text{ for all } x$$



f is strictly convex

$$f''(x) > 0 \text{ for all } x$$

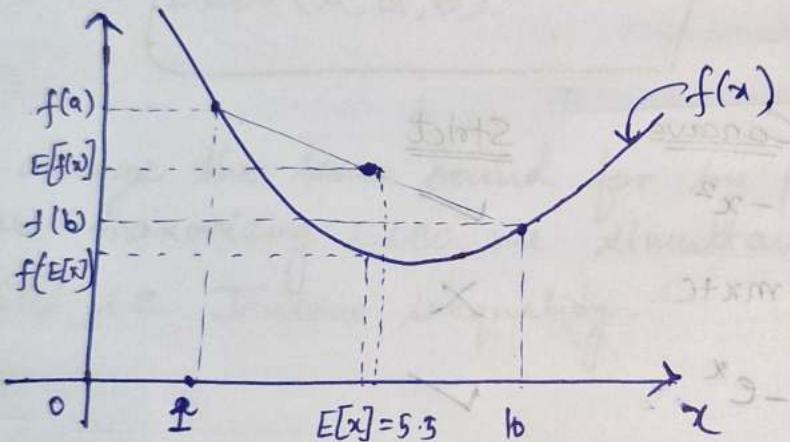


I $E[f(x)] \geq f(E[x])$ where f is convex
 $x \rightarrow$ any random variable

If f is strictly convex, then

II If $E[f(x)] = f(E[x])$
then $x = E[x]$ w.p. 1. → Means x is constant

Intuition for some disjoint news $[x]_3 \rightarrow (x)_3$



$E[x] = \int x \cdot p(x) dx$

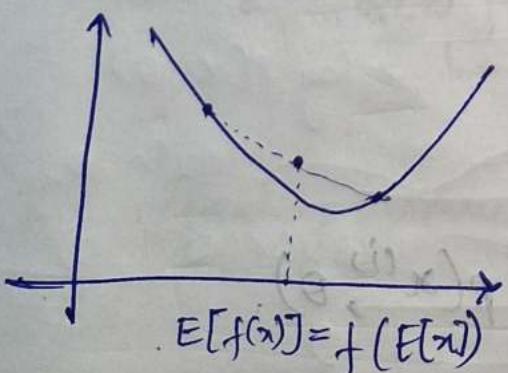
$E[f(x)] = \int f(x) \cdot p(x) dx$. because if we take

Strictly convex

means if f

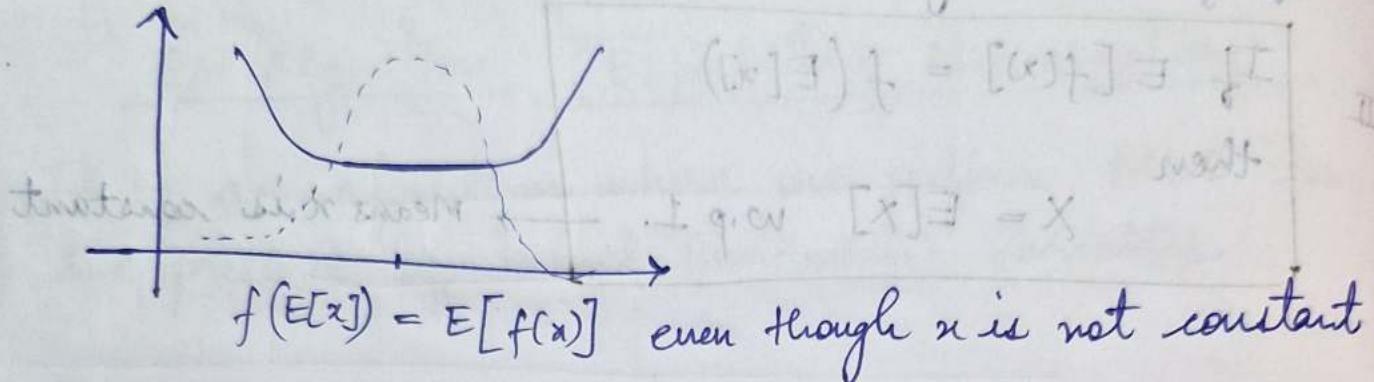
$(x)_{\text{gal}} = (x)_{\text{f}}$, f

$[x]_{\text{gal}} \geq [x]_{\text{f}}$



$E[f(x)] = f(E[x])$

The reason why we require f to be strictly convex is because if f were not strictly convex, we could have a case where $f(x)$ is flat and

Ex

$$\underline{\text{Convex}} \quad f(x) = x^2$$

$$\underline{\text{Concave}} \quad -x^2$$

$$f(x) = mx + c$$

$$mx + c$$

$$e^x$$

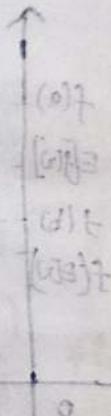
$$-e^x$$

$$-\log(x)$$

$$\log(x)$$

Strict

$$z = [x]$$



Jensen's inequality for concave function.
just the signs gets reversed.

If f is concave

$$\text{eg } f(x) = \log(x)$$

$$E[\log(x)] \leq \log E[x]$$

Goal is to maximize $\sum_{i=1}^n \log p(x^{(i)}, \theta)$.

Basically Maximize $\log p(x; \theta)$

$$\Rightarrow \log p(x; \theta) = \log \sum_z p(x, z; \theta)$$

$$= \log \sum_z Q(z) \cdot \frac{P(x, z; \theta)}{Q(z)}$$

$\left\{ \begin{array}{l} Q(z) > 0 \\ \text{for all } z \\ Q(z) \rightarrow \text{some} \\ \text{arbitrarily prob.} \\ \text{func. on } z \end{array} \right.$

$$\Rightarrow \log p(x; \theta) = \log E_{z \sim Q} \left[\frac{P(x, z; \theta)}{Q(z)} \right]$$

$$\geq E_{z \sim Q} \left[\log \frac{P(x, z; \theta)}{Q(z)} \right]$$

$$\log p(x; \theta) \geq \boxed{\text{ELBO}(x; Q, \theta)}$$

ELBO is always the lower bound for $\log p(x; \theta)$.
 If we are maximizing ELBO, we simultaneously minimize $\log p(x; \theta)$ i.e Jensen's inequality.

Are there cases when

$$\log p(x; \theta) = \text{ELBO}(x; \theta, \theta) \rightarrow \text{Yes}$$

If f is strictly convex and going back to II part of Jensen's inequality.

$$\Rightarrow E_{z \sim Q} \left[\log \frac{P(x, z; \theta)}{Q(z)} \right]$$

↑
Constant

The inequality becomes equality if $\frac{P(x, z; \theta)}{Q(z)} = \text{const.}$

↓
Constant

$$Q(z) = \frac{1}{c} p(x, z; \theta)$$

$$Q(z) \propto p(x, z; \theta)$$

$$Q(z) = \frac{p(x, z; \theta)}{\sum_z p(x, z; \theta)}$$

$$= \frac{p(x, z; \theta)}{P(x; \theta)}$$

$$Q(z) = p(z|x; \theta)$$

EM Algorithm:

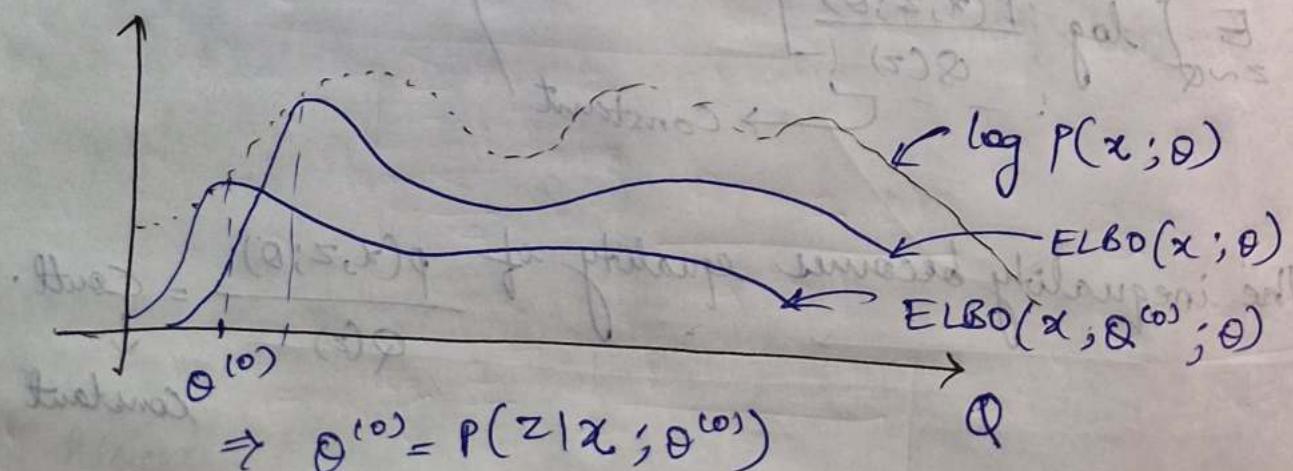
E-step

For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)}|x^{(i)}; \theta)$$

M-step

$$\theta = \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; \theta_i, \theta)$$



Lecture 17 - Factor Analysis & ELBO

Recap

① K-MEANS

Given: K , $\{x^{(i)}\}_{i=1}^n$

Randomly init $\boxed{\mu_j \in \mathbb{R}^d}$, $1 \leq j \leq K$

Loop till convergence

"E-Step": For each i

$$\text{Set } c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|_2^2$$

"M-Step": For each j

$$\mu_j = \frac{\sum_{i:c^{(i)}=j} x^{(i)}}{\sum_{i:c^{(i)}=j}}$$

$$J(c, \mu) = \sum_{i=1}^n \|x^{(i)} - \mu_{c^{(i)}}\|_2^2$$

② Gaussian Mixture Model [Soft K-means]

Given: $K, \{x^{(i)}\}_{i=1}^n$

Model: $z^{(i)} \sim \text{Multinomial}(\phi)$

$x^{(i)} | z^{(i)} \sim N(\mu_{z^{(i)}}, \Sigma_{z^{(i)}})$

Randomly Initialize ϕ, μ, Σ

Loop till convergence:

"E-step" For each :

$$\text{Set } w^{(i)} = P(z^{(i)} | x^{(i)}; \phi, \mu, \Sigma)$$

{ Use Bayes Rule }

"M-step" : For each j

$$\# \quad \phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)}$$

$$\# \quad \mu_j = \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$$

$$\# \quad \Sigma_j = \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$

③ E-M Algo $P(x, z; \theta)$ if z was observed
but z is not observed here

$$\log p(x; \theta) = \log \sum_z p(x, z; \theta)$$

$$\Downarrow = \log \sum_z Q(z) \frac{P(x, z; \theta)}{Q(z)}$$

$$(\log p(x; \theta)) = \log \mathbb{E}_{z \sim Q} \left[\frac{P(x, z; \theta)}{Q(z)} \right]$$

$$\begin{aligned} E_{z \sim Q}[g(z)] &= \sum_z Q(z) g(z) \\ g(z) &= \frac{P(x, z; \theta)}{Q(z)} \end{aligned}$$

$$\log P(\mathbf{x}; \theta) \geq \mathbb{E}_{z \sim Q} \left[\log \frac{P(\mathbf{x}, z; \theta)}{Q(z)} \right]$$

$$= \text{ELBO}(\mathbf{x}; Q, \theta)$$

Corollary

$$\log P(\mathbf{x}; \theta) = \text{ELBO}(\mathbf{x}; Q, \theta) \text{ iff } Q(z) = P(z|\mathbf{x}; \theta)$$

EM Algorithm [General form]Randomly init θ Loop till convergence

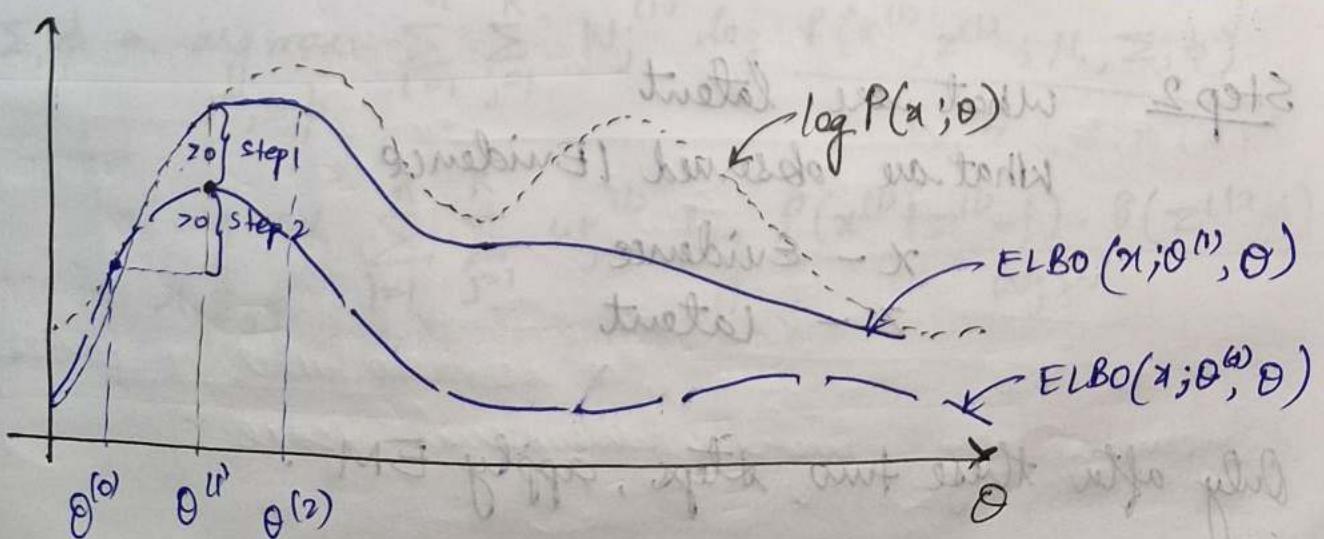
"E-Step": For each;

$$\text{set } Q_i(z) = P(z^{(i)} | \mathbf{x}^{(i)}; \theta)$$

:"M-step"

$$Q = \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(\mathbf{x}^{(i)}; \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^n \sum_z Q_i(z^{(i)}) \log \frac{P(\mathbf{x}, z; \theta)}{Q(z^{(i)})}$$



Proof of EM Convergence

For every t

$$l(\theta^{(t+1)}) \geq l(\theta^{(t)}) \quad \left\{ l(\theta) = \log P(x; \theta) \right\}$$

$$l(\theta^{(t+1)}) \geq \text{ELBO}(x; Q^{(t)}, \theta^{(t+1)}) \quad [\text{Jensen's Inequality}]$$

$$\geq \text{ELBO}(x; Q^{(t)}, \theta^{(t)}) \quad [\text{Due to M-step}]$$

$$l(\theta^{(t+1)}) = l(\theta^{(t)}) \quad [\text{corollary of Jensen's inequality}]$$

GMM via EM

Step 1 Data generating model information (Must be known)
Write out $p(x, z; \theta)$

$z \sim \text{Multinomial}$

$x|z \sim N(\mu_z, \Sigma_z)$

$$p(z) \cdot p(x|z) = p(x, z)$$

$$= p(x, z; \phi, \mu, \Sigma)$$

Step 2 What are latent

What are observed (Evidence)

x - Evidence

z - Latent

Only after these two steps, apply EM.

There are two types of unknown here

- ① Parameters $(\mu, \phi, \Sigma) \rightarrow$ global continuum (M-step)
- ② $x, z \rightarrow$ specific to each example (E-step)

E Step

For each i

$$\begin{aligned} Q_i(z^{(i)} = j) &= \frac{P(z^{(i)} | x^{(i)}; \mu, \Sigma, \phi)}{\sum_{j=1}^K P(z^{(i)} = j)} \\ &= \frac{P(x|z) \cdot P(z=j)}{\sum_{j=1}^K P(x|z=j) \cdot P(z=j)} \end{aligned}$$

$$\frac{\frac{1}{(2\pi)^{d/2}} |\Sigma_j|^{-1/2} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\} \phi_j}{\sum_j \left(\frac{1}{(2\pi)^{d/2}} |\Sigma_j|^{-1/2} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\} \phi_j \right)}$$

M-Step

$$\begin{aligned} \mu, \Sigma, \phi &= \arg \max_{\mu, \Sigma, \phi} \sum_{i=1}^n \sum_{j=1}^K w_j^{(i)} \log \frac{P(x^{(i)}, z^{(i)}; \mu, \Sigma, \phi)}{w_j^{(i)}} \\ &= \arg \max_{\mu, \Sigma, \phi} \sum_{i=1}^n \sum_{j=1}^K w_j^{(i)} \log \frac{P(x^{(i)} | z^{(i)} = j) \cdot P(z^{(i)} = j)}{w_j^{(i)}} \end{aligned}$$

$\hat{\phi} = \frac{1}{n} \sum_{i=1}^n w_j^{(i)}$ invariant for weight out size will

$\hat{\mu}_j = \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$ qot2

$\Sigma_j = \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j) (x^{(i)} - \mu_j)^T}{\sum_{i=1}^n (w_j^{(i)})}$ ideal fit

$$\frac{(i=s)q \cdot (i=s(x))q}{(x)q} =$$

Factor Analysis

$$x^{(i)} \in \mathbb{R}^d$$

$$S = \{x^{(i)}\}_{i=1}^n \quad \begin{matrix} (i=s)q \\ (i=s(x))q \end{matrix} =$$

dimensions $\gg n$ (examples)

$\left. \begin{array}{l} \text{usually examples} \\ \text{are much larger} \end{array} \right\} \text{the dimensions}$

$$\Sigma \in \mathbb{R}^{d \times d}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

$$\Sigma_{\text{rank}} = n \times 1 = n \text{ rank}$$

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

We are here considering the scenario where x comes from a very high dimensional space.

Now we assume $z^{(i)} \sim N(0, I)$, $I - K$ dimensional
 $z^{(i)} \in \mathbb{R}^K$ (K lower dimensional
subspace) $\rightarrow z$ resides in $\boxed{[k < n]}$

$\Rightarrow x^{(i)} | z^{(i)} \sim N(\mu + Lz, \Psi)$
 $L \in \mathbb{R}^{d \times K}$
 $\Psi \in \mathbb{R}^{d \times d}$ diagonal

+ There is this low dimensional (K) subspace where $K < n$, and the latent variable z that resides in the K dimensional subspace, which get mapped on to higher dimensional space (x). Using L (can be called an uplifting Matrix as it takes from $k \rightarrow d$), Lz : d dimensions, shifted by some offset μ and random diagonal noise Ψ .

Here z is continuous, not separate for each example. In factor analysis we are essentially trying to find the subspace of x . The finding low dimensional subspace (z) and assuming that x that we observe has corresponding latent variable z . The way x is generated from z is through the relation $N(\mu + Lz, \Psi)$.

Goal: $\log p(x; \mu, L, \Psi)$.
probabilistic model

- # We want to learn a model $p(x)$ so that we can make systems like for anomaly detection.
- # Using the set of assumptions, we want to maximize / learn $\log p(x; \mu, L, \Psi) \rightarrow$ This will be done by EM

Model

$$\boxed{z \sim N(0, I)}$$

$$x|z \sim N(\mu + Lz, \psi)$$
 $=$

$$\boxed{z \sim N(0, I)}$$

$$\epsilon \sim N(0, \psi)$$

$$x = \mu + Lz + \epsilon$$

$$\begin{bmatrix} z \\ x \end{bmatrix} \sim N\left(\mu_{zx}, \Sigma\right)$$

$$\mu_{zx} = \begin{bmatrix} \vec{\sigma} \\ \mu \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_z & \Sigma_{xz} \\ \Sigma_{zx}^T & \Sigma_x \end{bmatrix} = \begin{bmatrix} I & L^T \\ L & LL^T + \psi \end{bmatrix}$$

Model detail

$$P(x, z)$$

$$\Rightarrow \begin{bmatrix} z \\ x \end{bmatrix} \sim N\left(\begin{bmatrix} \vec{\sigma} \\ \mu \end{bmatrix}, \begin{bmatrix} I & L^T \\ L & LL^T + \psi \end{bmatrix}\right)$$

Parameters:

$$\mu, L, \psi$$

Observed - x

Latent - z

Plan of our approach

E step: $p(z|x)$

M step: $\mu, L, \psi = \arg \max \text{ELBO}$

Aside

$$\ell(\mu, L, \psi) = \log p(x)$$

$$x \sim N(\mu, LL^T + \psi)$$

$$\log p(x) = \frac{1}{(2\pi)^{d/2} |LL^T + \psi|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T [LL^T + \psi]^{-1} (x - \mu) \right\}$$

There is no close form solution for this in L and therefore we come up with factor analysis in Expectation maximization. E step, M step.

E StepCalculate $p(z|x)$

Using previous proofs i.e. $\begin{bmatrix} z \\ x \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & L^T \\ L & LL^T + \psi \end{bmatrix} \right)$

$$z|x \sim N \left[\underbrace{L^T(LL^T + \psi)^{-1}(x - \mu)}_{\text{Mean}}, \underbrace{I - L^T(LL^T + \psi)^{-1}L}_{\text{Covariance}} \right]$$

$$Q_i(z^{(i)}) = N(\text{Mean}, \text{Covariance})$$

M Step

$$\mu, L, \psi = \underset{\mu, L, \psi}{\arg \max} \sum_{i=1}^n \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}, \mu, L, \psi)}{Q_i(z^{(i)})}$$

$$\Rightarrow \underset{\mu, L, \psi}{\arg \max} \sum_{i=1}^n \underset{z^{(i)} \sim Q_i}{\mathbb{E}} \left[\log P(x^{(i)}|z^{(i)}) + \log P(z^{(i)}) - \log Q_i(z^{(i)}) \right]$$

Because it does not have $z^{(i)}$ which we are optimizing

$$\Rightarrow \underset{\mu, L, \psi}{\operatorname{argmax}} \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim Q_i} \left[\log \frac{1}{(2\pi)^{d/2} |\psi|^{1/2}} \exp \left\{ \frac{-1}{2} (x^{(i)} - \mu - L z^{(i)})^\top \psi^{-1} (x^{(i)} - \mu - L z^{(i)}) \right\} \right]$$

$$\sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim Q_i} \left[\left(\log \frac{1}{(2\pi)^{d/2}} \right) - \frac{1}{2} \log |\psi| - \frac{1}{2} (x^{(i)} - \mu - L z^{(i)})^\top \psi^{-1} (x^{(i)} - \mu - L z^{(i)}) \right]$$

$$\Rightarrow \hat{L} = \left(\sum_{i=1}^n (x^{(i)} - \mu) \mathbf{1}_{z^{(i)}}^\top \right) \left(\sum_{i=1}^n \mathbf{1}_{z^{(i)}} \mathbf{1}_{z^{(i)}}^\top + \sum_{z^{(i)}} \right)$$

$$\Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

$$\Rightarrow \hat{\Psi} = \frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top} = x^{(1)} \mathbf{1}_{z^{(1)}}^\top \mathbf{1}_{z^{(1)}}^\top L^\top - L \mathbf{1}_{z^{(1)}} x^{(1)\top}$$

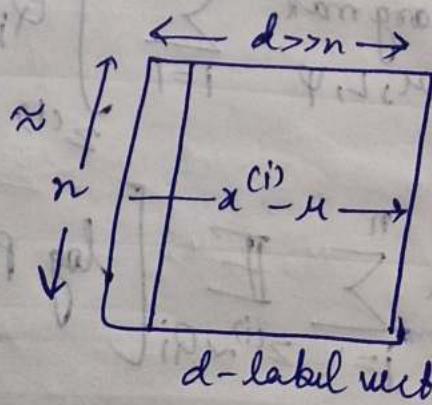
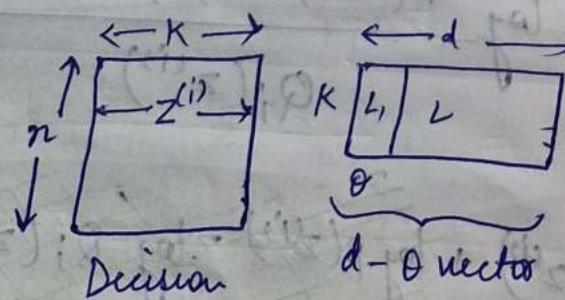
$$\Rightarrow V_{ij} = \hat{\Phi}_{ij}$$

Factor analysis sum up

Model: $z \sim N(0, I) \in \mathbb{R}^K$
 $x | z \sim N(\mu + Lz, \psi) \in \mathbb{R}^d$

$K < n < d$

Parameters: $\mu \in \mathbb{R}^d$, $L \in \mathbb{R}^{d \times K}$, $\psi \in \mathbb{R}^{d \times d}$



Lecture 18: Principal & Independent CA

On the basis of our previous derivations / proof we can make some best guesses about z and x .

$$\Rightarrow \frac{\mu_{z^{(i)}}}{x^{(i)}} \approx \hat{z}^{(i)} \quad (\text{Best estimate of } \hat{z}^{(i)})$$

$$\begin{aligned} \Rightarrow L \mu_{z^{(i)}} / x^{(i)} &= L \hat{z}^{(i)} \approx x^{(i)} - \mu \\ &= \hat{x}^{(i)} - \mu \end{aligned} \quad (\text{It can be seen as reconstruction})$$

Now writing L

$$L = \left[\sum_{i=1}^n y^{(i)} z^{(i)T} \right] \left[z^T z + \Sigma \right]^{-1}$$

$$L^T = (z^T z + \Sigma)^{-1} z^T y$$

For Ψ

$$\Psi_{ii} \approx \sum_{i=1}^n x^{(i)} x^{(i)T} - x^{(i)} \hat{x}^{(i)T} - \hat{x}^{(i)} x^{(i)T} + \hat{x}^{(i)} \hat{x}^{(i)T} + L \sum_{z^{(i)} / x^{(i)}} \Sigma$$

$$\Psi_{ii} \approx \left[\sum_{i=1}^n (x^{(i)} - \hat{x}^{(i)}) (x^{(i)} - \hat{x}^{(i)})^T + L \sum_{z^{(i)} / x^{(i)}} \Sigma \right]$$

Factor analysis can be thought of as a problem where we are given just the labels for d different problems and we are assuming a linear modelling of some hidden covariates. Don't know the design matrix, don't know the parameters, we are estimating both of these things given just the labels.

Just knowing labels $\xrightarrow{\text{Find}} x \text{ and parameters } (\mu, \Sigma, L)$

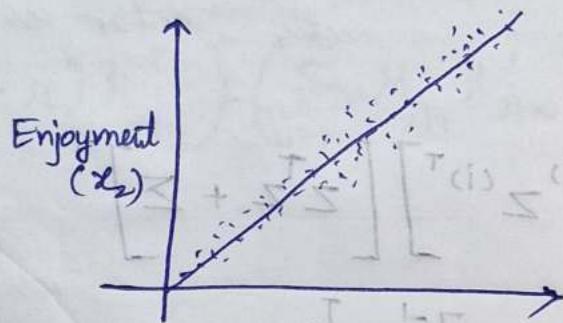
Principal Component Analysis (PCA)

$$\{x^{(i)}\}_{i=1}^n \in \mathbb{R}^d$$

Here we are back to our general case where n is much larger than d . Much more number of data points than the dimension of the space.

Ex

Flying helicopter →



Basically from the given feature relations can be correlated with each other, which will eventually help to reduce dimensions ($d \rightarrow k$)

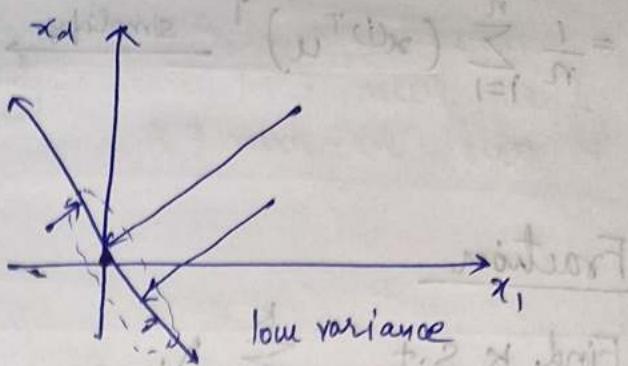
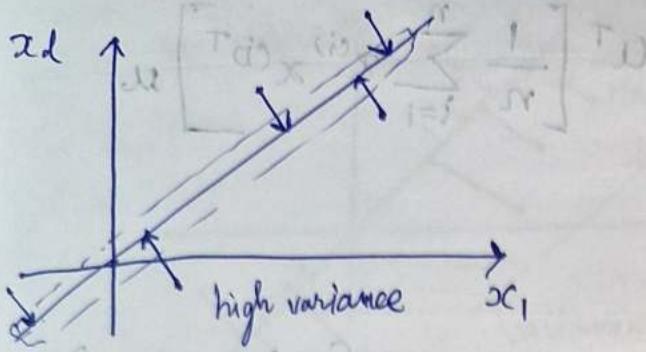
$$X: \begin{bmatrix} & \xleftarrow{\text{---}} d \xrightarrow{\text{---}} \\ \begin{matrix} n \\ \downarrow \\ \text{Skill} \end{matrix} & \left[\begin{matrix} | & | & | & | \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n)} \end{matrix} \right] + \begin{matrix} T \\ (x_1^{(1)} x_1^{(2)} \dots x_1^{(n)}) \end{matrix} \right] \approx \begin{bmatrix} & \xleftarrow{\text{---}} k \xrightarrow{\text{---}} \\ \left[\begin{matrix} | & | & | \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(n)} \end{matrix} \right] + \begin{matrix} T \\ (x_2^{(1)} x_2^{(2)} \dots x_2^{(n)}) \end{matrix} \right]$$

Step 1: Standardize

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j}$$

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n (x_j^{(i)} - \mu_j)^2$$



PCA says that the subspace onto which we project our data should be the one where the variance of the projected points is maximized.

$$\Rightarrow u \in \mathbb{R}^d \text{ unit length (Here } u \text{ is the basis vector)}$$

$$\Rightarrow \text{Proj}(u) \cdot \vec{x} = \frac{u u^T}{u^T u} \cdot \vec{x} = (\vec{x}^{(i)} \cdot u) \vec{u}$$

Our goal is to find the u , such that the variance across the projected points is maximum.

Basically find u such that $\frac{1}{n} \sum_{i=1}^n \|\text{Proj}(u) x_i\|^2$ is max.

$$\begin{aligned} u &= \underset{u}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \|\text{Proj}(u) x_i\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\vec{x}^{(i)} \cdot u)^2 \\ &= u^T \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \vec{x}^{(i)} \vec{x}^{(i)T} \right)}_{\text{Simple covariance Matrix}} u. \end{aligned}$$

$$u = \underset{u}{\operatorname{argmax}} u^T \left[\frac{1}{n} \sum_{i=1}^n \vec{x}^{(i)} \vec{x}^{(i)T} \right] u$$

Using prev. knowledge $\underset{u}{\operatorname{argmax}} u^T A u$

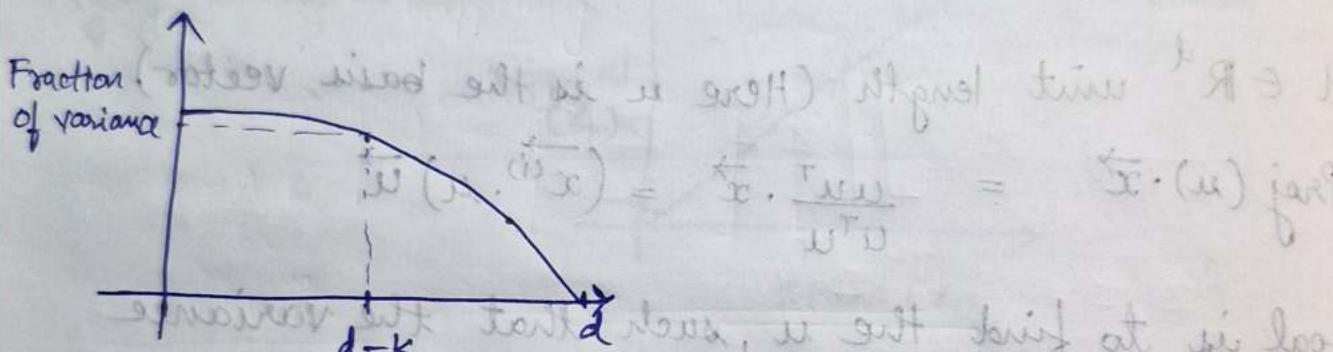
\Rightarrow Eigen value of A . (longer)

$$= \frac{1}{n} \sum_{i=1}^n (x^{(i)\top} u)^T \xrightarrow{\text{simplify}} u^\top \left[\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top} \right] u$$

FractionFind K s.t. $\sum_{i=1}^k \lambda_i$

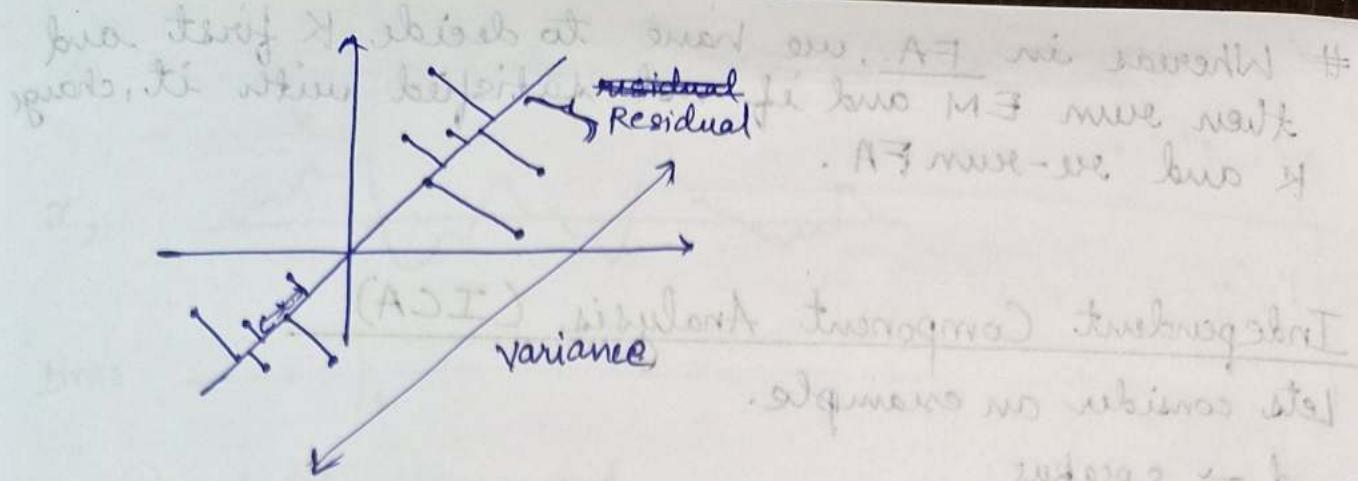
$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = 95\%$$

$\lambda \rightarrow \text{eigen}$
value

Value decomposition# $X \rightarrow$ Square and symmetricthen X has orthogonal eigen vectors & real eigen values# If X is PSD then it is obviously sq. & symm. Eigen values are positive.# We are using $X^\top X \rightarrow$ sq, symm and PSD

\Rightarrow Eigen decomposition on $X^\top X$ = SVD (singular value decomposition) on X

SVD is applied for any matrix.



If we rephrase the problem as finding the subspace where the projected points are as close as possible to the original points i.e. minimizing Residual.

Maximizing the variance and minimizing residual are equivalent

So far...

	Non Probabilistic	Probabilistic
Clustering	K-means	GMM
Subspace	PCA	Factor Analysis
	$U = [U_1 U_2]$ $n > d$ $XU^T = \text{Projecting } X \text{ on lower dimension}$ $U: x \rightarrow z$	

Loosely speaking, it can be thought of as the counterpart of classification

Regression

In PCA we first perform PCA to get the set of all eigenvalues & vectors and then decide what is the appropriate K for the problem.

Whereas in FA, we have to decide K first and then run EM and if not satisfied with it, change K and re-run FA.

Independent Component Analysis (ICA)

Let's consider an example.

$d \rightarrow$ speakers

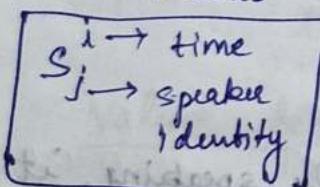
$d \rightarrow$ microphones placed in the room

$s \in \mathbb{R}^d \rightarrow$ One instantaneous recording of the speaker happening.

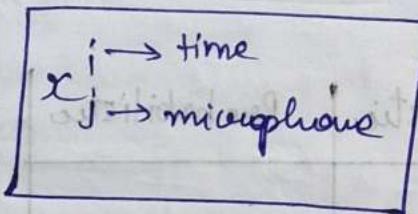
$x \in \mathbb{R}^d \rightarrow$ Recordings in microphone

$$\Rightarrow x = A \cdot s$$

\downarrow
Mixing
matrix



and

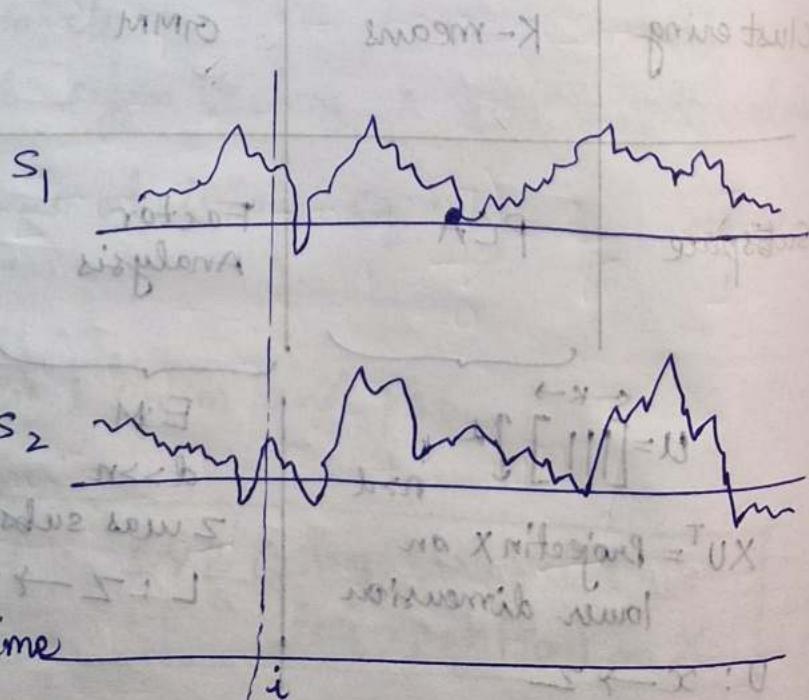


s_1

s_2

x_1

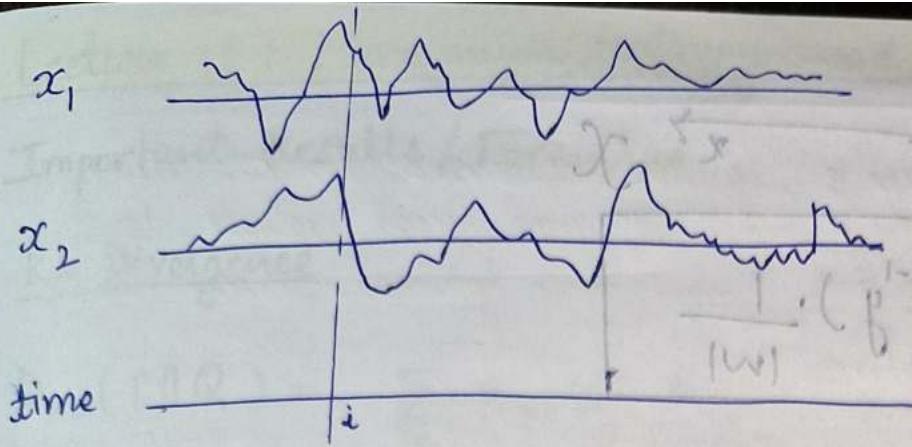
x_2



time

i

$$s^{(i)} = (s_1^{(i)}, s_2^{(i)}) \in \mathbb{R}^d$$



$$\frac{1}{|W|} \cdot (\mathbf{f}^T W) \mathbf{g} = (\mathbf{f}^T \mathbf{g}) \frac{1}{|W|}$$

$$\mathbf{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}) \in \mathbb{R}^d \quad d=2$$

$$\Rightarrow \mathbf{x}^{(i)} = \mathbf{A} \mathbf{s}^{(i)}$$

precise

$$\mathbf{w} = \mathbf{A}^{-1} \quad [\text{unmixing matrix}]$$

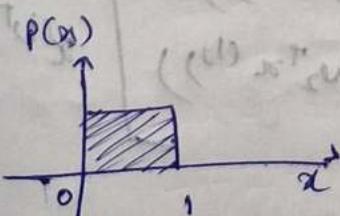
$$\Rightarrow \boxed{\mathbf{s}^{(i)} = \mathbf{w} \mathbf{x}^{(i)}} \rightarrow \text{extraction of source from mixed audio.}$$

Assumptions

- # No. of $s =$ No. of x .
- # Linear combination : $\mathbf{s} = \mathbf{w} \mathbf{x}$
- # The sources are independent : $s_j \perp s_k \quad j \neq k$
- # Gaussian; s_j is NOT Gaussian (Due to the problem of rotational invariance.)
(Therefore we use Laplace, logistic etc.)

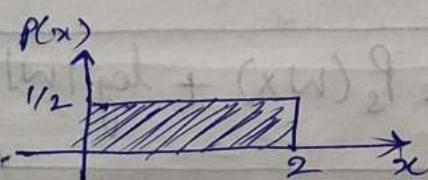
Proof

$$\mathbf{x} \sim \text{Unif}[0, 1]$$



$$y = 2x$$

$$P_y(y) = P_x\left(\frac{y}{2}\right) \cdot \frac{1}{2}$$



{ Area under
the PDF must
be 1 }

$$y = \underbrace{w^T x}_{R^d} + \underbrace{\epsilon}_{R^{d-d}} \quad \xrightarrow{x \in \mathcal{X}}$$

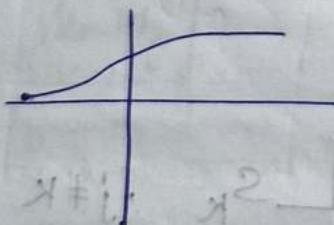
$$P_y(y) = P_x(w^T x) \cdot \frac{1}{|w|}$$

$$\Rightarrow P(x) = \prod_{j=1}^d P_s(w_j^T x) \cdot |w|$$

j^{th} source

$$[s_j] w = [-w_j -] [x] \quad \begin{matrix} \xrightarrow{\text{mixing principle}} \\ \xleftarrow{\text{unmixing matrix}} \end{matrix} \quad A = w$$

$P_s \sim \text{logistic distribution}$



$$\text{CDF of } \log F(x) = \frac{1}{1+e^{-x}} = \sigma(x)$$

$$\text{PDF of } f(x) = \sigma(x) \cdot (1 - \sigma(x))$$

$$\ell(w) = \sum_{i=1}^n \left[\left(\sum_{j=1}^d \log [\sigma(x^{(i)}) \cdot (1 - \sigma(x^{(i)}))] \right) + \log |w| \right]$$

$$w = w + \left[2 \begin{bmatrix} (1 - 2\sigma(w_1^T x^{(i)})) \\ (1 - 2\sigma(w_2^T x^{(i)})) \end{bmatrix} x^{(i)T} + (w^T)^{-1} \right]$$

$$S = Wx.$$

$$\hat{W} = \arg \max_w \log P_s(Wx) + \log |w|$$

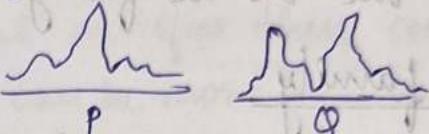
Lecture 19 : Maximum Entropy and Calibration

Important Results / Formulas.

KL Divergence

$$\# D_{KL}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = E_p \left[\log \frac{P}{Q} \right]$$

- $P, Q \rightarrow$ Probability distributions



- This is not symmetric. $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$

$$\# \text{ENTROPY } H(P) = - \sum_x P(x) \log P(x) = E_p \left[\log \frac{1}{P} \right]$$

$$\# \text{CROSS-ENTROPY } H(P, Q) = - \sum_x P(x) \log Q(x) = E_p \left[\log \frac{1}{Q} \right]$$

$$\# \text{Relative Entropy} \Rightarrow D_{KL}(P \parallel Q) = H(P, Q) - H(P).$$

= KL Div.

Exponential Family

$$P(y; \eta) = b(y) \exp \{ \eta \cdot T(y) - a(\eta) \}$$

$$\Rightarrow \text{Mean } (\mu) = a'(\eta)$$

$$\Rightarrow \eta = (a')^{-1}(\mu)$$

$$(\eta_j)' T \sum_{i=1}^n \frac{1}{\eta_i} = (\mu_j)' \alpha$$

$$\left((\eta_j)' T \sum_{i=1}^n \frac{1}{\eta_i} \right)' A_{ij} = \hat{\beta}_j$$

MLE of Exponential families

Previously exp. family, y was assumed to be always a scalar like in regression it was real values, logistic regression \rightarrow binary $(0, 1)$, some some generalization with softmax (some values).

Now we do it for generalized case ($y \rightarrow$ anything)

Exp family

$$P(y; \eta) = b(y) \exp\{n \cdot T(y) - a(n)\}$$

$$S = \{y^{(1)}, \dots, y^{(n)}\} \quad \text{i.i.d assumption}$$

$$\hat{\eta} = \boxed{\arg \max_{\eta} \log \prod_{i=1}^n P(y^{(i)}; \eta)}$$

$$= \arg \max_{\eta} \sum_{i=1}^n [\log b(y^{(i)}) + \eta \cdot T(y^{(i)}) - a(n)]$$

$g = a'$ = canonical and f^n

$g = (a')^{-1}$ = canonical link fn

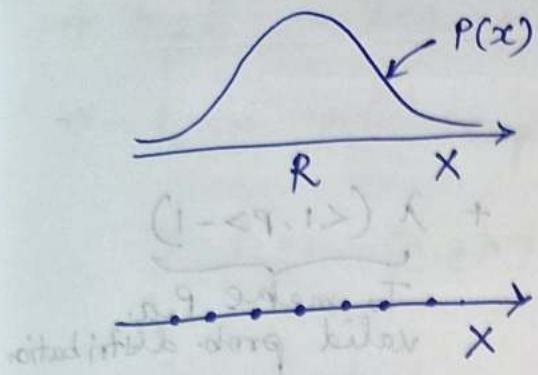
$$\nabla_n = \sum_{i=1}^n (0 + T(y^{(i)}) - a'(n)) = 0$$

$$\Rightarrow \underbrace{\sum_{i=1}^n a'(\eta)}_{\text{Total no. of n}} = \sum_{i=1}^n T(y^{(i)})$$

$$\Rightarrow a'(\eta) = \frac{1}{n} \sum_{i=1}^n T(y^{(i)})$$

General form

$$\boxed{\hat{\eta}_{MLE} = (a')^{-1} \left(\frac{1}{n} \sum_{i=1}^n T(y^{(i)}) \right)}$$



Constraints

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \hat{\mu})^2$$

- * Any probability distribution that satisfies these constraints is a feasible distribution. It can be more than one. Out of all these choose the one which has maximum entropy.

$$P^* = \arg \max_P H(P)$$

$$\text{s.t. } \sum_i T_j(p_i) = C_j$$

$$T_j : Y \rightarrow R, C_j$$

$$\sum_{j=1, \dots, m} T(y_i) p_i = C_j$$

for $j = 1, \dots, m$

$$T(y) = y$$

constraint for mean

$$\sum_{i=1}^N T(y) \cdot P(y) = C$$

$$= \left[\sum_{i=1}^N y_i P(y) \right] = E[y]$$

Support y is finite

N is the $|Y|$

M constraints (It means M such T functions)

$$C_j = \frac{1}{n} \sum_{i=1}^n T(y^{(i)})_j$$

$$\# \mathcal{L}(p, \eta, \lambda) = \underbrace{H(p)}_{\text{Entropy}} + \underbrace{\langle \eta, T_p - c \rangle}_{\substack{\text{constraints} \\ (\text{From data})}} + \underbrace{\lambda (\langle 1, p \rangle - 1)}_{\substack{\text{To make } p \text{ a} \\ \text{valid prob. distribution}}}$$

$$T_i \quad \begin{bmatrix} T_1(y_1) & T_1(y_2) & \dots & T_1(y_N) \\ \vdots & \vdots & \ddots & \vdots \\ T_m(y_1) & T_m(y_2) & \dots & T_m(y_N) \end{bmatrix} \quad \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_N \end{bmatrix} \quad N = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_N \end{bmatrix}$$

$$p_i \geq 0, \sum_{i=1}^N p_i = 1, p_2 = P(y_2)$$

$$\begin{bmatrix} E_p[T_1(y)] \\ E_p[T_2(y)] \\ \vdots \\ E_p[T_N(y)] \end{bmatrix} = \begin{bmatrix} c_1 \\ \vdots \\ c_N \end{bmatrix}$$

$$B = (B)^T$$

Solving Lagrangian (\mathcal{L})

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p_i} &= \frac{\partial}{\partial p_i} \left[H(p) + \langle \eta, T_p - c \rangle + \lambda (\langle 1, p \rangle - 1) \right] \\ &= \frac{\partial}{\partial p_i} \left[-p_i \log p_i + \langle \eta, T_p - c \rangle + \lambda (\langle 1, p \rangle - 1) \right] \\ &= \left[-(1/p_i + 1) + \langle \eta, T(y_i) \rangle + \lambda \right] = 0 \end{aligned}$$

$$\Rightarrow \log p_i = \lambda - 1 + \langle \eta, T(y_i) \rangle$$

$$\Rightarrow p_i = e^{\lambda - 1} \cdot \exp \{ \eta \cdot T(y_i) \}$$

$$\sum_{i=1}^n p_i = 1 \quad \sum_{i=1}^n e^{\lambda - 1} \cdot \exp \{ \eta \cdot T(y_i) \} = 1$$

$$e^{\lambda - 1} = \frac{1}{\sum \exp \{ \eta \cdot T(y_i) \}} \rightarrow z(\eta)$$

$$\Rightarrow P(y) = \frac{\exp \{ \eta \cdot T(y) \}}{z(\eta)}$$

$$a = \log z$$

$$\Rightarrow P(y|\eta) = \exp \{ \eta \cdot T(y) - a(\eta) \}$$

$$\Rightarrow \nabla_\eta \mathcal{L} = \nabla_\eta \left[- \sum_y P(y) \log P(y) + \eta \cdot T_p - \eta \cdot c + \lambda \right]$$

$$\Rightarrow \nabla_\eta \left[- \sum_y P(y) \{ \eta \cdot T(y) - a(\eta) \} + \eta \cdot E_p[T(y)] \cdot n_c + \lambda \right]$$

$$\Rightarrow \nabla_\eta \left[-\eta E_p[T(y)] + a(\eta) - \eta \cdot E_p[T(y)] - \eta \cdot c + \lambda \right]$$

$$\Rightarrow \nabla_\eta [a(\eta) - \eta \cdot c + \lambda]$$

$$\Rightarrow a'(\eta) - c = 0$$

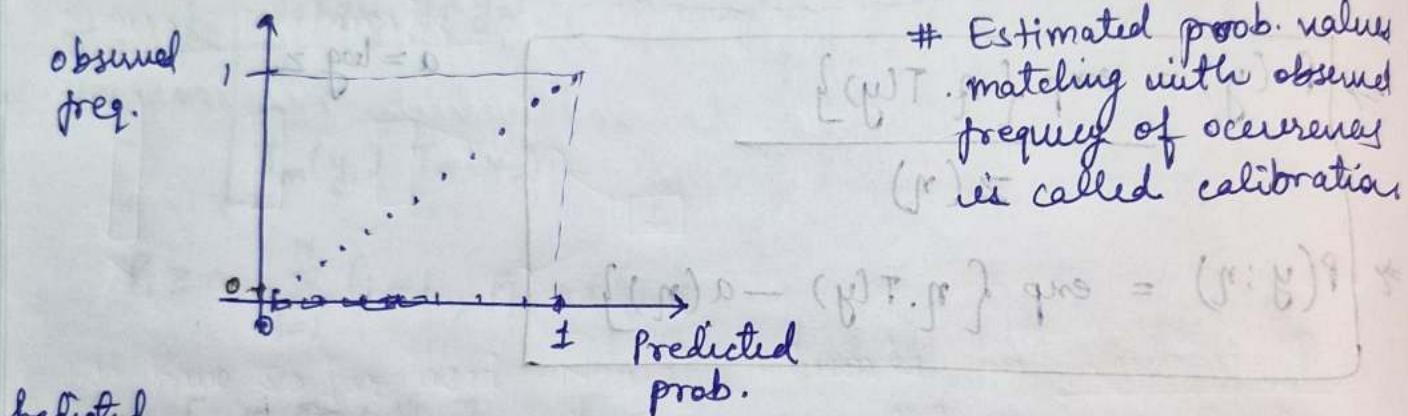
$$a'(\eta) = c$$

$$\hat{\eta} = (a')^{-1}(c) = (a')^{-1} \left[\frac{1}{n} \sum_{i=1}^n T(y^{(i)}) \right]$$

$$\hat{y} = \frac{1}{1 + \exp(-\theta^T x)} \in (0, 1)$$

Calibration

It tells us about quality of \hat{y} being an actual probability.

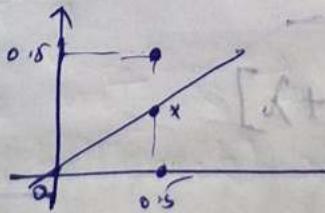


Predicted

$$\begin{array}{l} \hat{y}^{(i)} \\ (0,1) \end{array} \quad \begin{array}{l} y^{(i)} \\ \{0,1\} \end{array} \quad \begin{array}{l} \text{Ground truth} \\ \{0,1\} \end{array}$$

Calibration \Leftrightarrow Accuracy. (NO, both ways)

$$\begin{aligned} p(\text{correct}) &= 0.5 \\ p(\text{current}) &= 0.01 \end{aligned}$$



Entropy vs MLE

Max. Entropy \Rightarrow M.L.E

$$\text{loss}_{\theta} = -\log P_{\theta}(x)$$

- + Conventional way \rightarrow loss (y, \hat{y}) - 02 weeks
- + Here Scoring Rule $[P, y]$
 Prob distn over all possible outcomes \leftarrow The actual observed outcome.

Proper scoring rule, $f: (P, y) \rightarrow R$

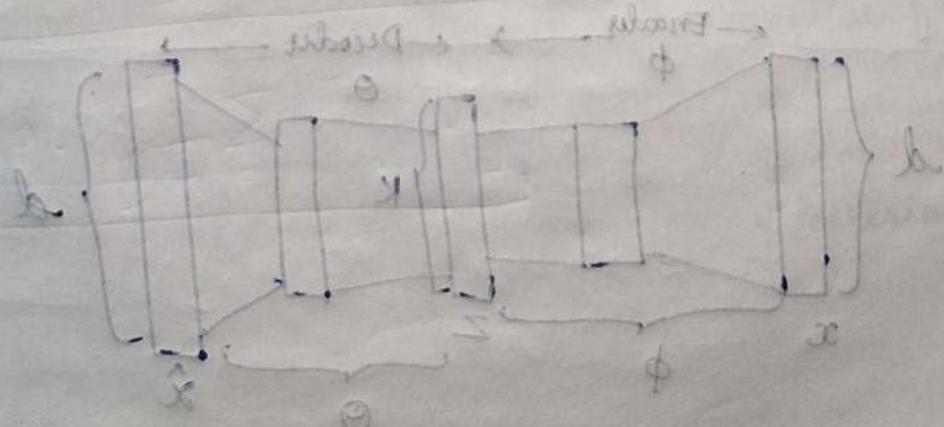
$$E_{y \sim Q} [f(Q, y)] \leq E_{y \sim Q} [f(P, y)] + P, Q$$

Strictly Proper scoring rule $\boxed{P = Q}$

\Rightarrow Basically f will assign lower values for better predictions.

$$\text{Loss} = E_{y \sim Q} [-\log \frac{1}{P}]$$

$$E_{y \sim Q} [-\log Q(y)] \leq E_{y \sim Q} [-\log P(y)] = 0 \leq E_{y \sim Q} [\log \frac{Q(y)}{P(y)}]$$

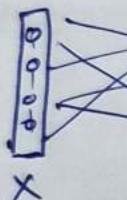


Lecture 20 - Variational Autoencoder

Variational autoencoders are the simplest Deep Generative models. VAE was an early method which made good progress in Deep gen. AI.

① Autoencoders

Neural Networks



$\theta \leftarrow (\theta, \beta) f_t$. other pristine signal

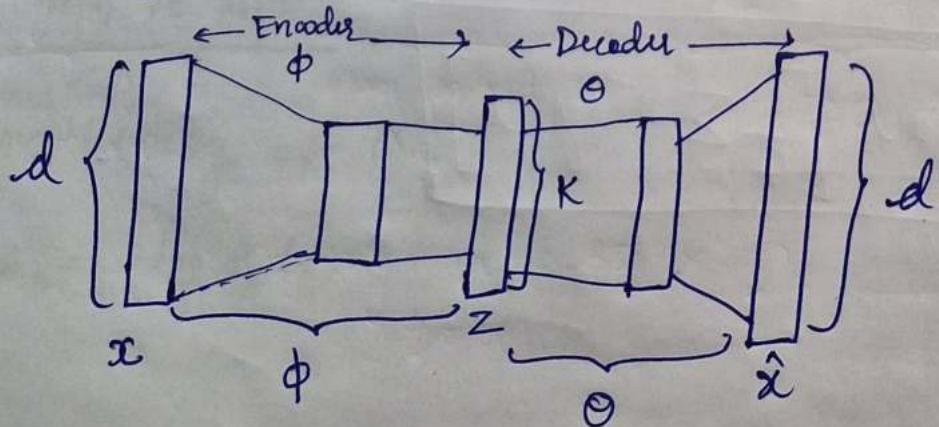
$$\nabla_{\theta} [(\theta, \beta) f_t] \Xi = [(\theta, \beta) f] \Xi$$

$$\nabla_{\beta} [(\theta, \beta) f_t] \Xi = [(\theta, \beta) f] \Xi$$

- # We have so far used Neural networks in supervised setting, where (X, Y) are given. And we use Backpropagation to update weights and biases.

Now for Autoencoders NN:

- # Now in unsupervised learning we are only given X . $S = \{x^{(1)}, \dots, x^{(n)}\}$ and not y .
- # The goal of autoencoder is that is to learn a way in which we obtain a z (bottleneck) and reconstruct the original data.



$$\# \text{ Loss} = \|\hat{x} - x\|_2^2$$

$$\text{Loss}_{(\theta, \phi)} = \sum_{i=1}^n \left\| x^{(i)} - \text{Dec}_{\theta} \left[\underbrace{\text{Enc}[x^{(i)}]}_{z^{(i)}} \right] \right\|_2^2$$

$\underbrace{\quad\quad\quad}_{\hat{x}^{(i)}}$

⇒ Here we take the data transformation from a bottleneck z ($d \rightarrow K$). It helps the model to learn some low dimensional representation of high dimensional input data (d). Then starting from the low.dim data we map it back to high dim. data itself.

- ⇒ Then minimize the loss ($\|\hat{x} - x\|_2^2$). Basically we have learned compression into a latent state / hidden state.
- ⇒ Weight & bias of Encoder are (ϕ) and of Decoder are (θ). These weight & biases (model parameters) are updated similarly that of backpropagation using loss calculated.

Extra

⇒ There is also a Denoising-autencoder where we learn how to denoise the input data. We provide the noisy input (\tilde{x}) and recover original (x) using this approach.

② MCMC : EM

E Step : For all i

$$\text{Set } Q_i^{(t)}(z) = P(z|x^{(i)}; \theta^{(t)}) \quad \{ \text{Posterior} \}$$

↳ Initialization → Iteration

M Step :

$$\theta^{(t+1)} = \arg \max_{\theta} \sum_{i=1}^n \sum_z Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z)}$$

$$= \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim Q_i^{(t)}} [\log P(x^{(i)}, z^{(i)}; \theta)]$$

$$\approx \arg \max_{\theta} \sum_{i=1}^n \frac{1}{M} \sum_{m=1}^M \log P(x^{(i)}, z_i^{(m)}; \theta)$$

Monte carlo estimate

$$z_i^{(m)} \sim Q_i^{(t)} \quad (\text{sampled})$$

In original EM we constructed the convergence proof where with every EM step the likelihood was only increasing. But here the guarantee does not hold any more because this is an approximation of lower bound and not the exact lower bound.

Variational Inference

$$\log P(x) \geq \text{ELBO}(x; Q)$$

$$\log P(x) = \text{ELBO}(x; Q) + ?$$

$$D_{KL}(Q || P_{z|x}) = \log P(x) - \text{ELBO}(x; Q)$$

$$\underbrace{\log P(x)}_{\text{Constant w.r.t } Q} = \underbrace{\text{ELBO}(x; Q)}_{\text{Max w.r.t } Q} + \underbrace{D_{KL}(Q || P_{z|x})}_{\geq 0}$$

[Variational]

$$P_{z|x} \approx \arg \max_{Q \in \mathcal{Q}} \text{ELBO}(x; Q) \quad \leftarrow \text{Variational Inference}$$

- # This approach where we maximize the lower bound (ELBO) w.r.t to Q and obtain a distribution which we approximate to $p(z)$ is variational inference.
 - # In the Monte Carlo technique, we do sampling while in variational inference we do optimization.
 - # In Monte Carlo we never know how good our estimate is all we know is that we eventually reach the exact solution. While in variational inference will converge and we know when to stop (The soln. is approximate).
- Q: In variational inference most of the times the question is how do we choose the family Q from which we want to perform optimization.

$$Q(z) = Q_1(z_1) \cdot Q_2(z_2) \cdot Q_3(z_3) \cdots Q_K(z_K)$$

$$z \in \mathbb{R}^K$$

This assumption is - Mean field assumption. VI

Variational Autoencoders

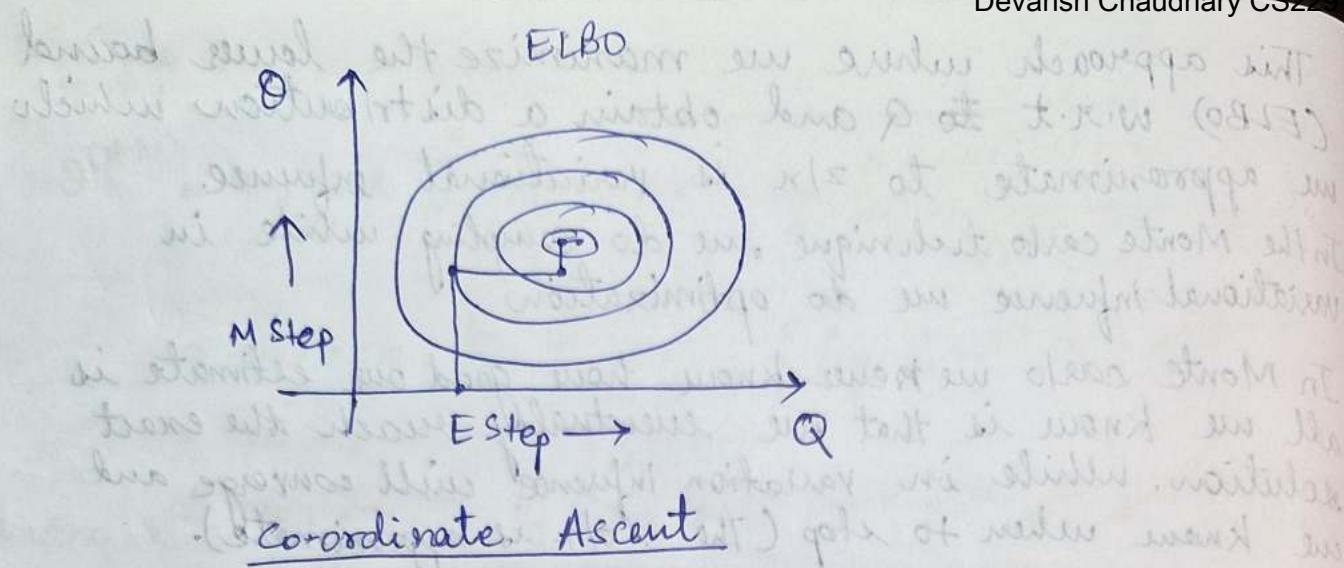
ELBO (Q, θ)

In EM we construct this ELBO, in E Step we find the best possible Q (posterior) & in the M Step we would update θ and while performing these steps we keep the other constant.

E Step : Calculate Q, θ fixed.

M Step : Calculate θ, Q fixed

This technique is also called Coordinate Ascent



Can we do gradient ascent?

In VAE we max. ELBO using gradient ascent as well

Prior $\rightarrow Z \sim N(\theta, I_{K \times K})$

Like Likelihood $\rightarrow X|Z \sim N(g(z; \theta), \sigma^2 I)$

g is a neural network with param θ .

Posterior $p(z|x)$

without A denoted by

Q family for V-I

$Q_i(z) = N\left(q(x^{(i)}; \phi), \text{diag}((V(x^{(i)}; \psi))^2)\right)$

Amortized Inference

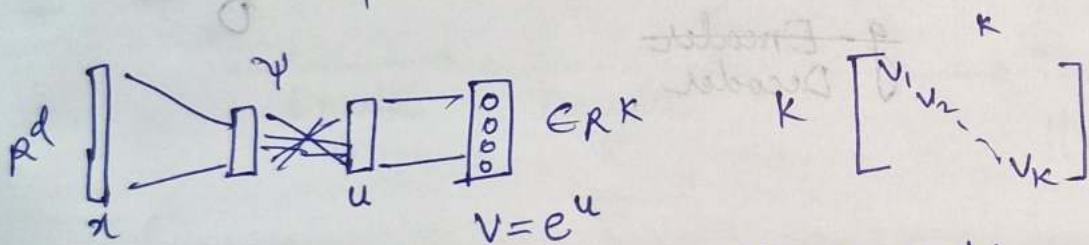
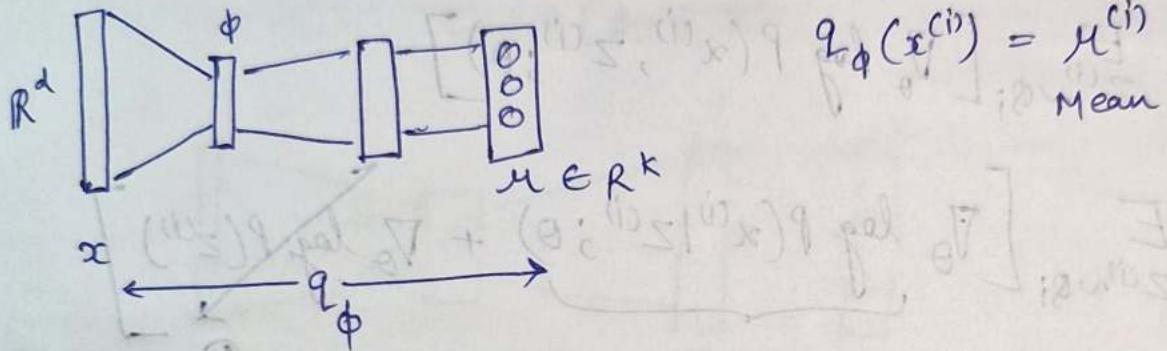
* The Difference b/w EM and VAE is that in EM we were separately calculating the parameters of the Q distribution for each example independently.

In Amortized Inference, we do not compute separately but we assume that all Q comes from Normal distribution N and the mean and variance of each Q_i is a function

of the x 's. These functions are Neural Networks.

+ Feed x as the input and the NN will output two scalars per input i.e Mean and Variance.

Note: This is not autoencoder (No. bottleneck).



We are making mean field assumption

Feed x in input $\xrightarrow{\text{Learned parameter}} \phi, \psi \xrightarrow{\text{Estimate Q from } \phi, \psi.}$ (E Step construction)

Basically

$$\text{ELBO}(\phi, \psi, \theta) = \sum_{i=1}^n \mathbb{E}_{\substack{z^{(i)} \sim Q_i}} \left[\log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

Where $Q_i = N(q(x^{(i)}; \phi), \text{diag}(v(x^{(i)}; \psi)))$

$\theta = \theta + \eta \nabla_\theta \text{ELBO}(\phi, \psi, \theta)$	→ Gradient updates until convergence.
$\phi = \phi + \eta \nabla_\phi (\text{ELBO})$	
$\psi = \psi + \eta \nabla_\psi (\text{ELBO})$	

For θ

$$\nabla_{\theta} \text{ELBO}(\phi, \psi, \theta)$$

$$\Rightarrow \nabla_{\theta} \sum_{i=1}^n E_{z^{(i)} \sim q_i} \left[\log \frac{P(x, z; \theta)}{Q(z^{(i)}; \theta)} \right] \quad \left. \begin{array}{l} \{\theta \text{ independent} \\ \text{of } Q \end{array} \right\}$$

$$\Rightarrow \sum_{i=1}^n E_{z^{(i)} \sim q_i} \left[\nabla_{\theta} \log P(x^{(i)}, z^{(i)}; \theta) \right]$$

$$\Rightarrow \sum_{i=1}^n E_{z^{(i)} \sim q_i} \left[\underbrace{\nabla_{\theta} \log P(x^{(i)} | z^{(i)}; \theta)}_{\text{Encoder}} + \nabla_{\theta} \log P(z^{(i)}) \right] \quad \left. \begin{array}{l} \text{Decoder} \\ \theta \end{array} \right\}$$

For ϕ

$$\nabla_{\phi} \text{ELBO}(\phi, \psi, \theta)$$

$$\Rightarrow \nabla_{\phi} \sum_{i=1}^n E_{z^{(i)} \sim q_i} \left[\log \frac{P(x, z; \theta)}{Q(z^{(i)}; \theta)} \right] \quad \left. \begin{array}{l} \{\phi \text{ is parameter} \\ \text{of } Q \end{array} \right\}$$

$$z \sim N(\mu, \sigma)$$

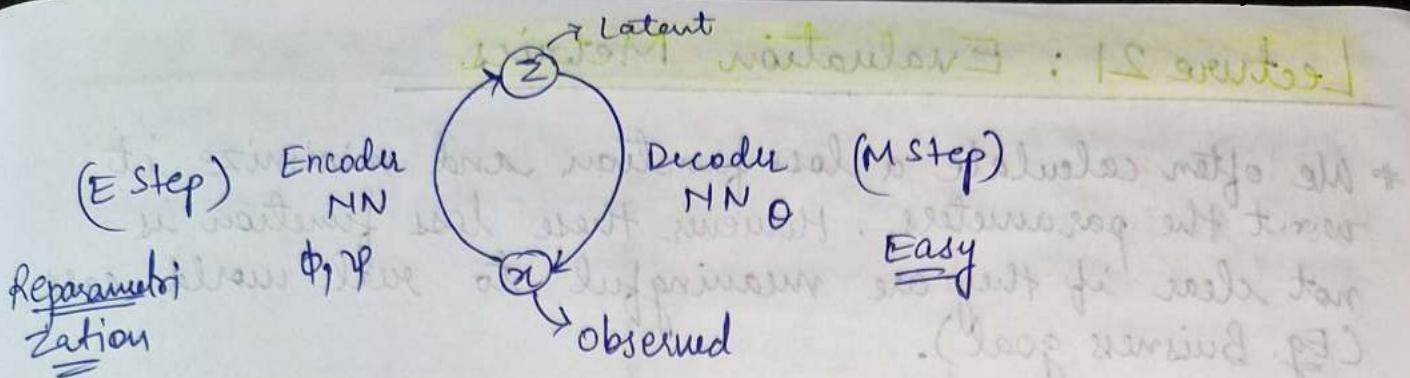
$$z = \epsilon \cdot \sigma + \mu \quad \epsilon \sim N(0, 1)$$

$$\nabla_{\phi} \sum_{i=1}^n E_{\epsilon^{(i)} \sim N(0, 1)} \left[\log \frac{P(x, \epsilon^{(i)} \cdot \sum^{(i)} + \mu^{(i)}; \theta)}{Q(\epsilon^{(i)} \cdot \sum^{(i)} + \mu^{(i)}; \theta)} \right]$$

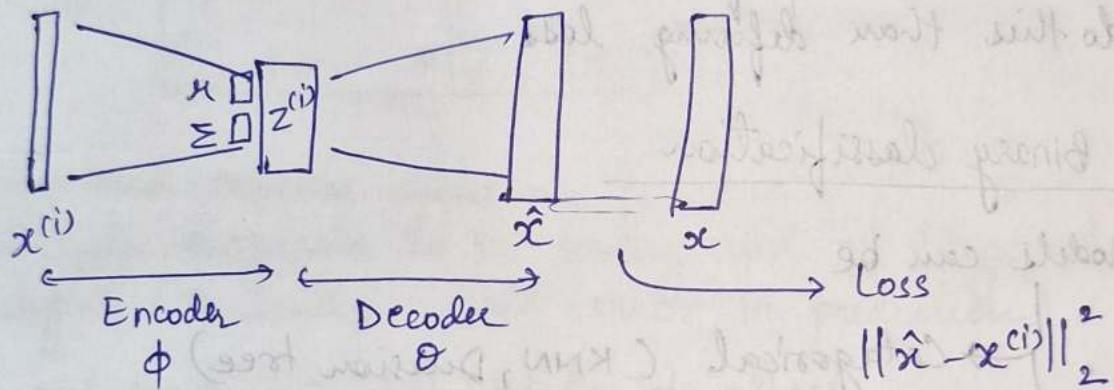
$$\text{where } \mu^{(i)} = q(x^{(i)}; \phi)$$

$$\sum^{(i)} = \text{diag}(v(x^{(i)}; \psi))$$

Now we can take gradient inside



#



(minimum integral, MV2) greater than last

TP	TN	FP	FN
3	3	0	0

where best was

signals labelled as + 0
signals labelled as - 0

$$I = 3/2$$

$$O = 8/2$$

$$\frac{3V + \#}{(3V - \#) + (3V + \#)} = \text{percentage} \#$$

Lecture 21 : Evaluation Metrics

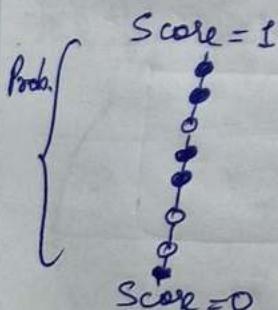
- # We often calculate a loss function and minimize it w.r.t the parameters. However these loss function is not clear if they are meaningful to real world case (Eg. Business goal).
- # Evaluation metrics help to quantify the gap between the desired performance and baseline. It is easier to do this than defining loss.

Binary classification

Models can be

- Categorical (KNN, Decision tree)
- Real value score (SVM, Logistic regression)

Score based models



- +ve labelled example
- -ve labelled example

$$\# \text{ Positive} = \frac{\# \text{ +ve}}{(\# \text{ +ve}) + (\# \text{ -ve})}$$

Point Metrics: Confusion Matrix

	Label +ve	Label -ve	
Perfect +ve	(9) TP • • • • • • • • •	FP (2) • •	Threshold = 0.5
Th = 0.5 Perfect -ve	TP + FN • • • • • • • • •	TN + FP • • • • • • • • •	
	(1) FN •	TN (8) • • • • • • • •	

- # Total sum and column sums are fixed.
- # We want the diagonals to be heavy and off diagonals to be light (Basically least error in prediction).
- # FP and FN have very diff. kinds of effect

Th	TP	TN	FP	FN	Acc.	Pr	Recall	Spec	Fi
0.5	9	8	2	1	0.85	0.81	0.90	0.8	0.85

Accuracy = $\frac{\text{Predicted correct}}{\text{Total}} = \frac{(9+8)}{(9+8+2+1)} = \frac{(TP+TN)}{(TP+TN+FP+FN)}$

→ Fraction of example we got right.

Precision = PPV (Positive predictive value)

→ Fraction of predicted examples actually positive

$$\Rightarrow \left(\frac{TP}{TP + FP} \right) = \frac{(9)}{(9+2)}$$

Recall (Sensitivity) or Positive recall:

Imagine if we were to use this classifier in deployment, what fraction of all the actual positives are we gonna recover. Nothing to deal with -ve.

$$\text{Recall} \quad (\text{Sensitivity}) = \frac{g}{g+1} = \left(\frac{\text{TP}}{\text{TP} + \text{FN}} \right)$$

Negative Recall (specificity)

What fraction of the actual examples (negatives only) did the model correctly classify as negative

$$\text{Neg. Recall} \quad (\text{Specificity}) = \frac{\text{TN}}{\text{TN} + \text{FP}} = \frac{8}{8+2}$$

F1 score

F1 score is the harmonic mean of precision and recall

$$\frac{1}{F_1} = \frac{1}{\text{Per}} + \frac{1}{R_c}$$

$$F_1 = \frac{\text{Pr} \cdot R_c}{\text{Pr} + R_c}$$

G1 score

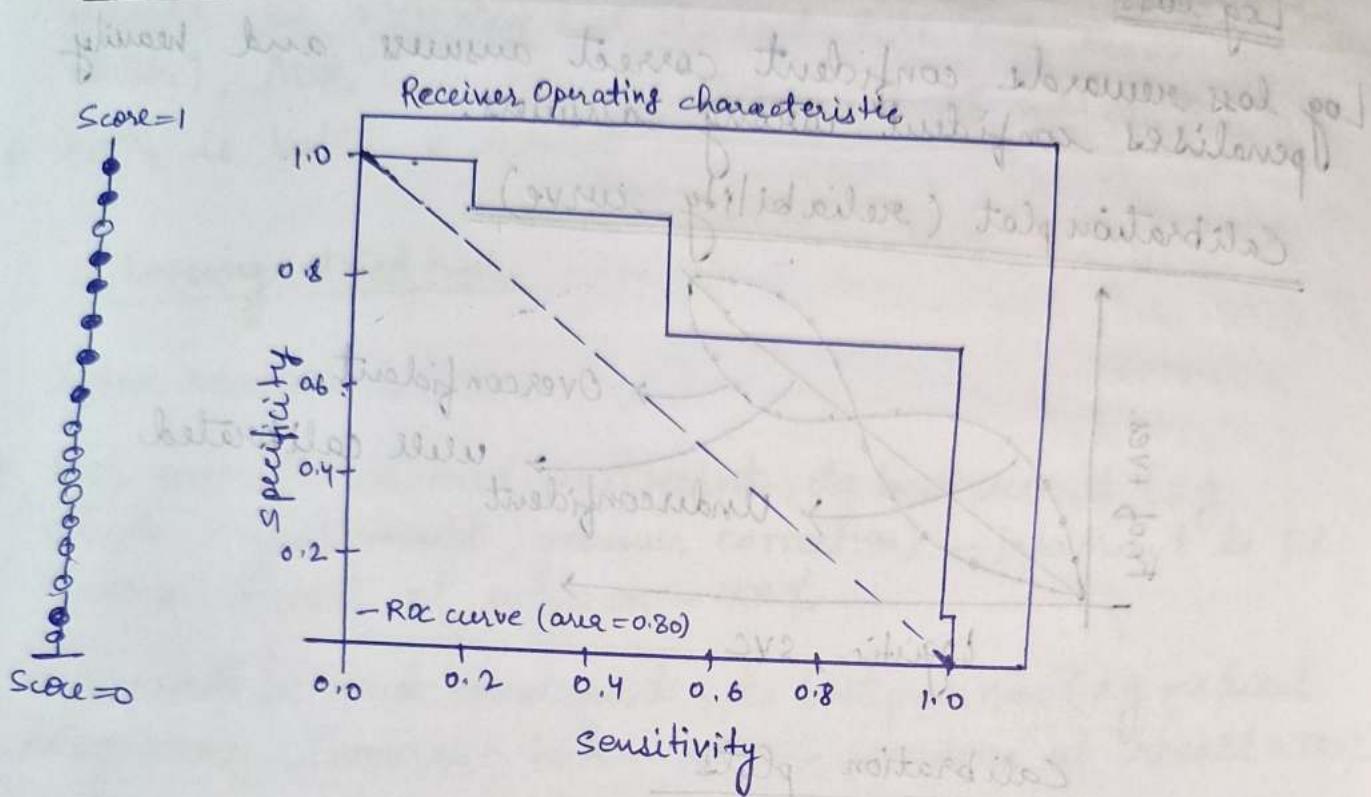
Geometric mean

$$G_1 = \sqrt{\text{Pr} \cdot R_c}$$

$$2 \log G_1 = \log \text{Pr} + \log R_c$$

As we change the threshold we get different values of all these param. Total thresholds = (1 + Total examples).

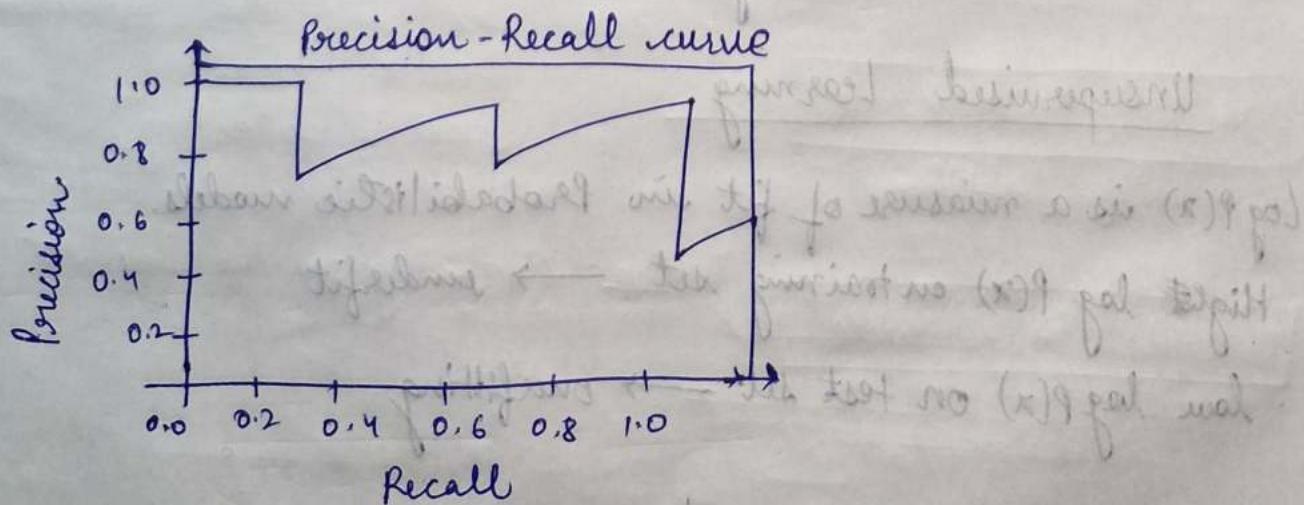
ROC curve



Area under the ROC curve will be high if most of the +ve examples are the ones we encounter first as we start scanning down.

Another intuition is that what is the probability that a random +ve example is ranked higher than a random -ve example. This probability is exactly equal to the area under the ROC curve.

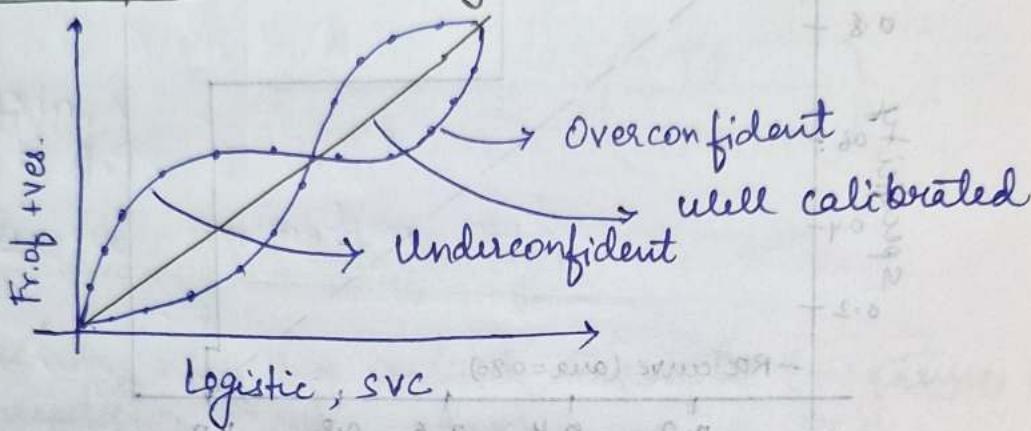
PRC curve



Log loss

Log loss rewards confident correct answers and heavily penalises confident wrong answers.

Calibration plot (reliability curve)



Calibration plots

Loss

$$\text{Booser} = \frac{1}{n} \sum_{i=1}^n (\hat{p}^{(i)} - y^{(i)})^2$$

$\hat{p} \rightarrow$ Predicted probability

Mean square error

- # Calibration: How meaningful the assigned probabilities are
- # Discrimination: How well the model can separate the +ve & -ve in a ranking sense.

Unsupervised Learning

$\log P(x)$ is a measure of fit in probabilistic models

High $\log P(x)$ on training set \rightarrow underfit

Low $\log P(x)$ on test set \rightarrow overfitting

- # In a class imbalance scenario one best preference should be accuracy (it doesn't tell any reasonable value), AUR_c can also be fooled/wrong intuition.
- # AVPR_c is better & robust

Choosing Metrics

Some common Patterns

- # High precision is hard constraint, do best recall (e.g. search engine result, grammar correction) -- intolerant to FP. Metric: Recall at precision = XX%.
- # High recall is hard constraint, do best precision (e.g. medical diagnosis). Intolerant to FN. Metric: Precision at recall = 100%.
- # Capacity constrained (by k). Metric: Precision in top k.

some other metric of interest -

generalized Jaccard measure

said at first : good

generalized precision

library ext for various : good

generalized rec - user

measured interesting thing : good

generalized test user

worst & user

[criticizes for test]

generalized full answer

test will ext no criticizes : good

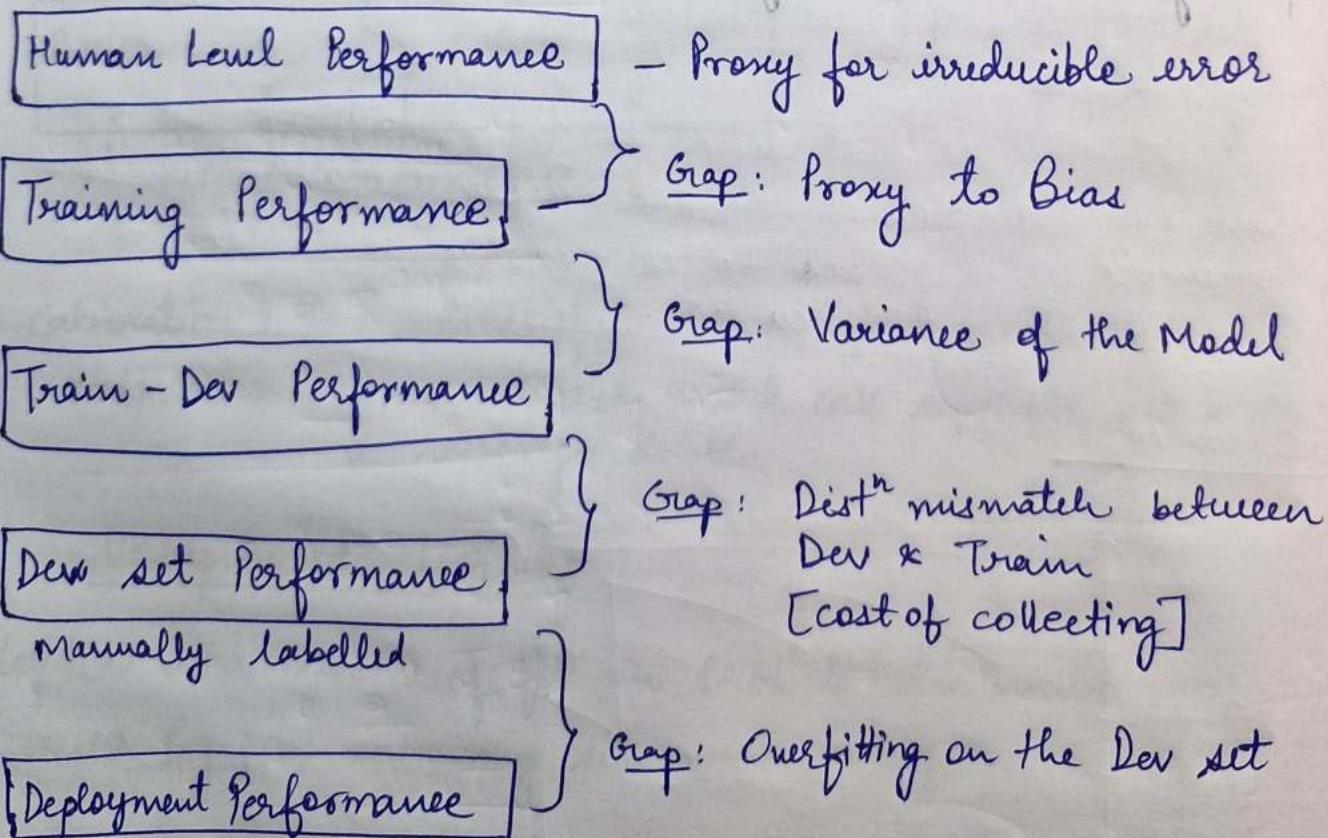
generalized true/false

Lecture 22 - Course Recap and Practical Tips

ML is Practice vs Research

- ① Start collecting a Dev set
- ② Dev set distribution match the production scenario.
- ③ Collect data & label
- ④ Define an evaluation metric
- ⑤ Training data
 - As close as possible to dev set
 - Automated / noisy labelling

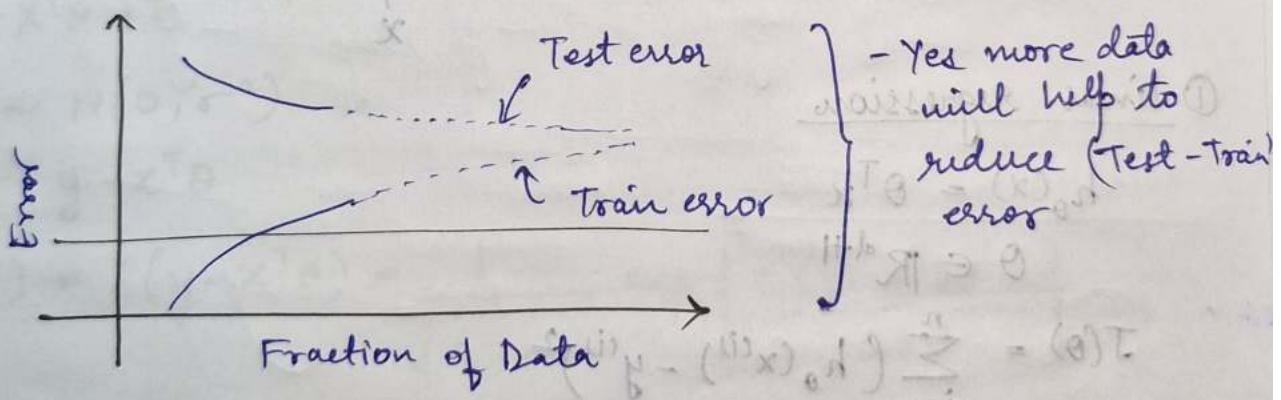
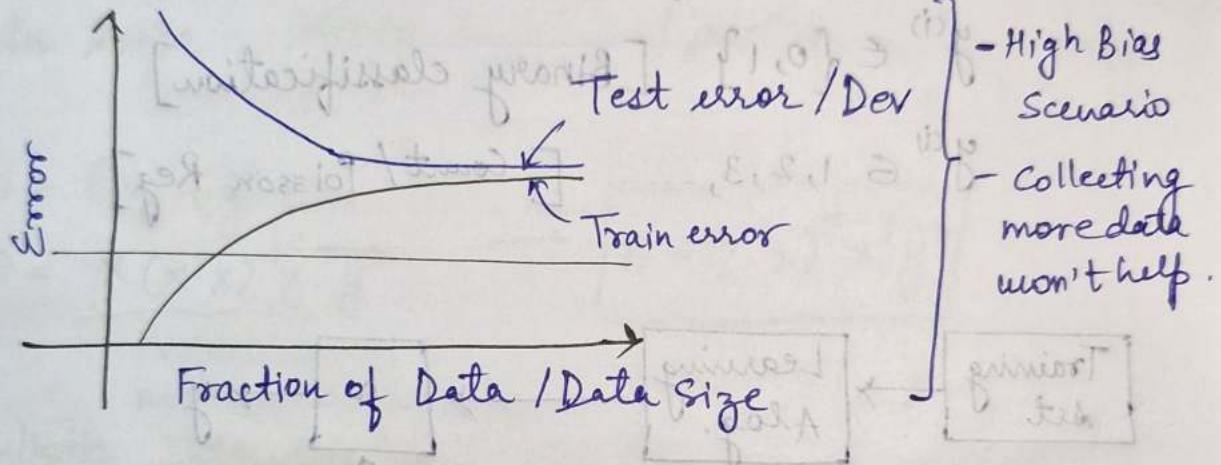
Breakdown of Performance



If the model has high bias, then these steps will help to fight bias

- add feature
- add depth
- Reduce regularization
- Complex kernels

Do the opposite of all these to fight variance
 \oplus Get more data



Our always focus should be on deployment performance

$$(x_i, (x_i)_\text{gt} - \hat{y}) \sum_{i=1}^n x_i + \theta = 0$$

target \hat{y}

$$(x_i, (x_i)_\text{gt} - \hat{y}) \sum_{i=1}^n x_i + \theta = 0$$

$$(x_i, (x_i)_\text{gt} - \hat{y}) \sum_{i=1}^n x_i + \theta = 0 \quad \text{target: } \hat{y}$$

SUPERVISED LEARNING

$$X \rightarrow Y$$

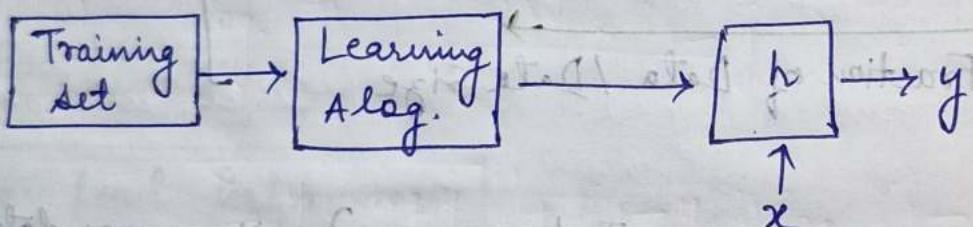
$$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$$

(input $x^{(i)} \in \mathbb{R}^d$)

(output $y^{(i)} \in \mathbb{R}$ [Regression])

($-y^{(i)} \in \{0, 1\}$ [Binary classification])

($y^{(i)} \in 1, 2, 3, \dots$ [Count / Poisson Reg])



① linear regression

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

$$\theta \in \mathbb{R}^{d+1}$$

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

$$= \|X\theta - \vec{y}\|_2^2$$

$$\hat{\theta} = \arg \min_{\theta} J(\theta)$$

② Gradient Descent

Repeat

$$\theta = \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x^{(i)}$$

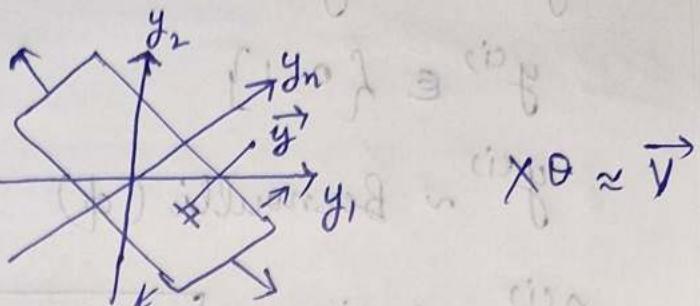
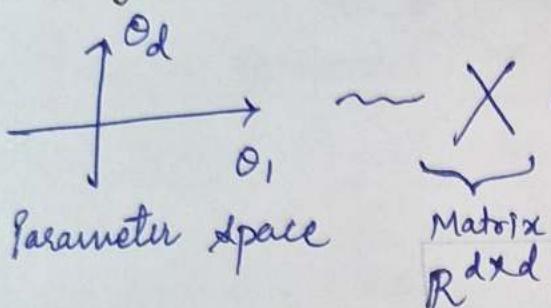
SGD: Repeat $\theta = \theta + \alpha (y^{(k)} - h_{\theta}(x^{(k)})) \cdot x^{(k)}$

, $k \sim \text{Unif}[1, n]$

③ Normal Equation

$$X^T X \theta = X^T y \Rightarrow \hat{\theta} = (X^T X)^{-1} X^T y$$

④ Projection



$$x\theta = \text{Proj } (\vec{y}; X)$$

$$x\theta = \underbrace{x(x^T x)^{-1} x^T}_{\text{Projection Matrix}} \vec{y}$$

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

⑤ Probabilistic Interpretation

$$y = x^T \theta + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

$$\epsilon = y - x^T \theta$$

$$P(\epsilon) = P(y - x^T \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y - x^T \theta)^2}{2\sigma^2} \right\}$$

How we do MLE which is to basically find the parameter θ so that it maximizes this likelihood term.

$$\ell(\theta) = \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(i)} - x^{(i)\top} \theta)^2}{2\sigma^2} \right\}$$

$$= \sum_{i=1}^n C + \left\{ -\frac{(y^{(i)} - x^{(i)\top} \theta)^2}{2\sigma^2} \right\}$$

$$= C' + \frac{1}{2\sigma^2} \left\{ \sum_{i=1}^n (y^{(i)} - x^{(i)\top} \theta)^2 \right\}$$

$$\arg \max_{\theta} l(\theta) = \arg \min_{\theta} \sum_{i=1}^n (y^{(i)} - x^{(i)T} \theta)^2$$

$B^T X = \theta X^T X$

Logistic Regression

$$y^{(i)} \in \{0, 1\}$$

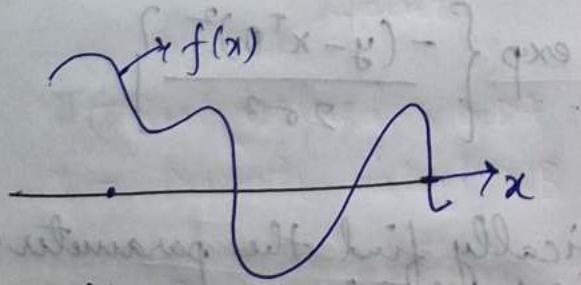
$$y^{(i)} \sim \text{Bernoulli}(\phi)$$

$$g^{(i)} = P(y^{(i)} = 1) = \frac{1}{1 + \exp\{-\theta^T x^{(i)}\}}$$

$$l(\theta) = \log \prod_{i=1}^n (\hat{y}^{(i)})^{y^{(i)}} \cdot (1 - \hat{y}^{(i)})^{(1-y^{(i)})}$$

$$= \sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})$$

Newton's Method (Scalar)



$$l(\theta) = (\theta^T x - b) q = (3)$$

$$f(x) = 0$$

$$f = l'$$

$$\theta = \theta + \frac{1}{l''(\theta)}$$

$$l'(\theta) =$$

$$l(\theta) = 0 \Leftrightarrow l'(\theta) = 0$$

Newton Raphson (Vector)

$$\theta = \theta + H^{-1}(\theta) l(\theta)$$

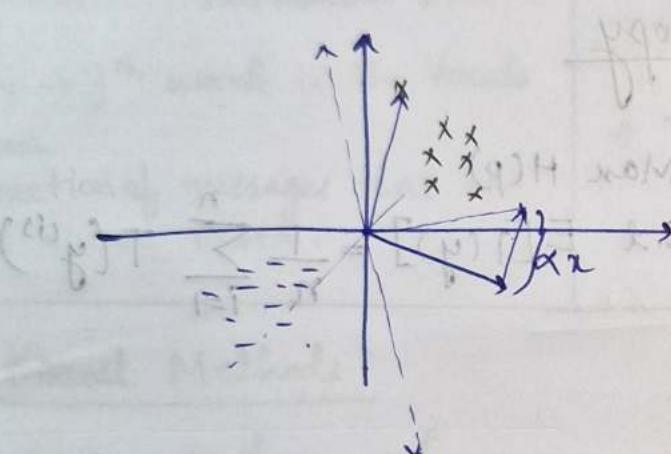
$$\text{where } H(\theta) = \nabla_{\theta}^2 l(\theta)$$

Lecture 23 - Course Recap & Wrap Up

Perceptron algorithm
Encounter one example at a time

$$\theta = \theta + \alpha [y - g(\theta^T x)] \cdot x$$

$$y(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$



$$(\theta + \alpha x)^T x \geq \theta^T x$$

If a separating hyperplane does exist, then no matter in what order you present the examples to the learning algo. It will eventually find some separating hyperplane.

Exp Family

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

y - variable

$T(y)$ - Sufficient statistics [$T(y) = y$]

η - Natural Parameter

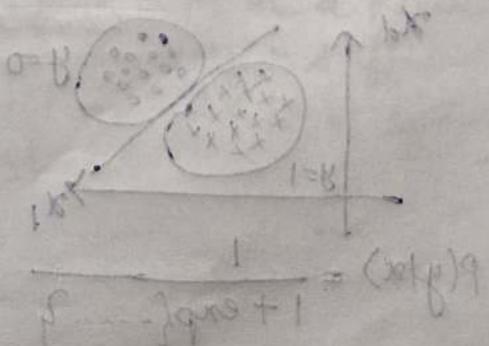
$b(y)$ - Base measure

$a(\eta)$ - log partition function

GLM

$$\eta = \theta^T x$$

Learnable parameter



Properties

$$\# E[T(y)] = E[y] = a'(\eta)$$

$$\# \text{Var}[y] = a''(\eta)$$

} Exponential family
suff. stat. and sufficient statistics

$$x[(x^T \theta) - \mu] \propto e^{\theta} = \theta$$

GLM

$$E[y|x; \theta] \cdot x = a'(\theta^T x)$$

$$a''(\eta) = \text{Var}[y|x; \theta] \cdot x^T x \geq 0 \quad \text{PSD}$$

Exp family & Max Entropy

$$\left\{ y^{(i)} \right\}_{i=1}^n \stackrel{\text{MLE}}{\leftarrow} x^{(i)}$$

Max H(p)

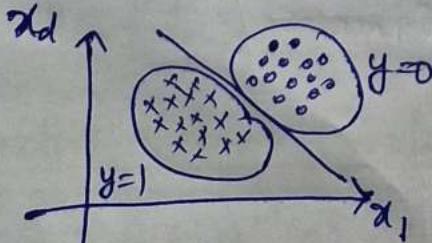
$$\text{s.t. } E[T(y)] = \frac{1}{n} \sum_{i=1}^n T(y^{(i)})$$

Generative ModelsDiscriminative - $P(y|x)$

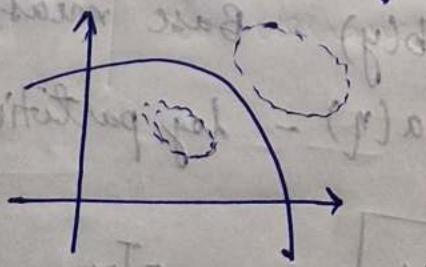
$$\underbrace{P(y, x)}_{\text{Joint}} = \underbrace{P(x|y) \cdot P(y)}_{\substack{\text{Likelihood} \\ \text{Prior}}} \quad [\text{Class Prior}]$$

 x - Real valued GDA x - Discrete valued NBGDA - $P(y) \sim \text{Bernoulli}(\phi)$

$$P(x|y) \sim N(\mu_y, \Sigma_y) \quad N(\mu_y, \Sigma_y)$$



$$P(y|x) = \frac{1}{1 + \exp\{-\dots\}}$$



$$x^T \theta = \mu$$

Naive Bayes

$$x_i \perp x_j | y \neq x_i \perp x_j$$

Bernoulli Event Model

$$P(y) = \phi_y \text{ class Prior}$$

$$P(x_j | y) = \phi_j | y$$

Bernoulli Dist.ⁿ

$x_j \rightarrow j^{\text{th}}$ word in the Vocab

Fraction of messages that has the spam word.

Multinomial Event Model

$$P(y) = \phi_y \text{ class Prior}$$

$$P(x_j | y) = \phi_j | y \text{ Multinomial Dist.}$$

$x_j \rightarrow j^{\text{th}}$ word in a message

Total no. of times the word appear in the set.

Kernel Methods

$$\phi(x) : \mathbb{R}^d \rightarrow \mathbb{R}^p$$

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

P - Dim of feature space p can be ∞

$$K(x, x') = \phi(x)^T \phi(x')$$

Properties

K is symmetric

$$\{x^{(1)}, \dots, x^{(m)}\}$$

$$\begin{bmatrix} K(x^{(1)}, x^{(1)}) & \dots & K(x^{(1)}, x^{(m)}) \\ \vdots & \ddots & \vdots \\ K(x^{(m)}, x^{(1)}) & \dots & K(x^{(m)}, x^{(m)}) \end{bmatrix} \rightarrow \text{PSD}$$

$$\underbrace{\text{Frobenius}}_{\text{Matrix Norm}} + \underbrace{\text{Determinant}}_{\text{Scalar}} = \underbrace{\text{SVD}}_{\text{Matrix}} \quad (\text{Fact})$$

[For func. K to be a kernel it is necessary & sufficient for the matrix to be symm and PSD]

Mercer's Theorem

This method allows us to have ~~as dimensional~~ feature maps.

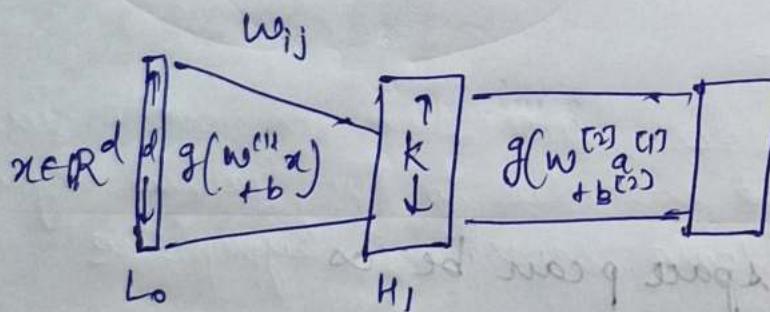
Gaussian Processes

Kernel method for regression

$$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \sim N\left[\mu, (\Sigma)\right]$$

posterior or we know $f_i \leftarrow x$
 \downarrow
 $f_i \sim GP \left(\mu, K \right)$

Fully connected Networks



$$w \in \mathbb{R}^{K \times d}$$

$$w_x \in \mathbb{R}^K$$

$$MSE_{(\text{test})} = \underbrace{\text{Irreducible error}}_{\text{Noise in Test example}} + \underbrace{\text{Bias}^2}_{\text{Inflexibility or Lin Capacity of Model class}} + \underbrace{\text{Variance}}_{\text{Training data noise contributes to variance}}$$

Irreducible error
 Inflexibility or Lin Capacity of Model class
 Training data noise contributes to variance

MSE & RMSE

$$\text{Loss} = \left(\sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})^2 \right) + \lambda \|\theta\|_2^2 - \text{MAP w/ Gaussian}$$

$$(h_{\theta}(x^{(i)}) + \lambda \|\theta\|_1) \|\theta\|_1 - \text{MAP w/ Laplace}$$

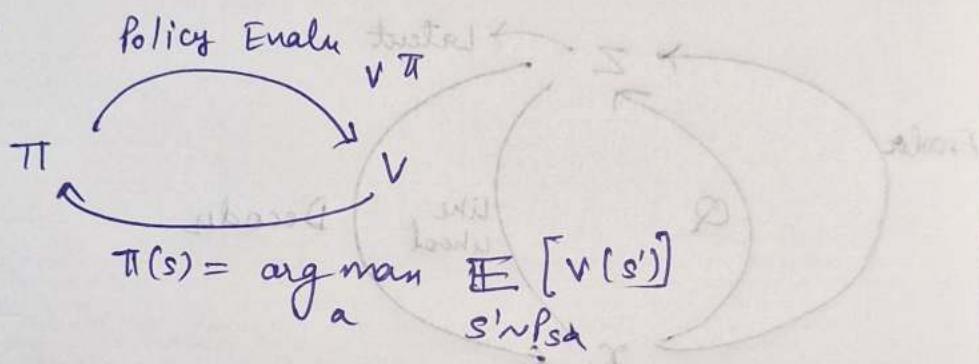
RL

$$\text{MDP} = \{S, A, \{P_{sa}\}, \mathcal{V}, R\}$$

$$\text{Value: } V^\pi(s)$$

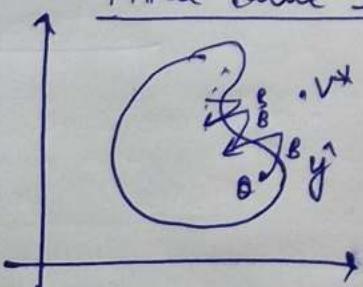
$$\text{Policy: } \pi(s) \rightarrow A$$

\rightarrow stochastic
 \Rightarrow total
 action

Policy Iteration# Bellman Equation

$$\text{Optimal } V^*(s) = \max_{\pi} V^\pi(s)$$

$$V(s) = R(s) + \gamma \max_a \mathbb{E}_{s' \sim P_{sa}} V(s')$$

Fitted value Iteration

EM

- # Model $P(x, z; \theta)$
- # Parameter θ
- Evidence x
- Latent variable z

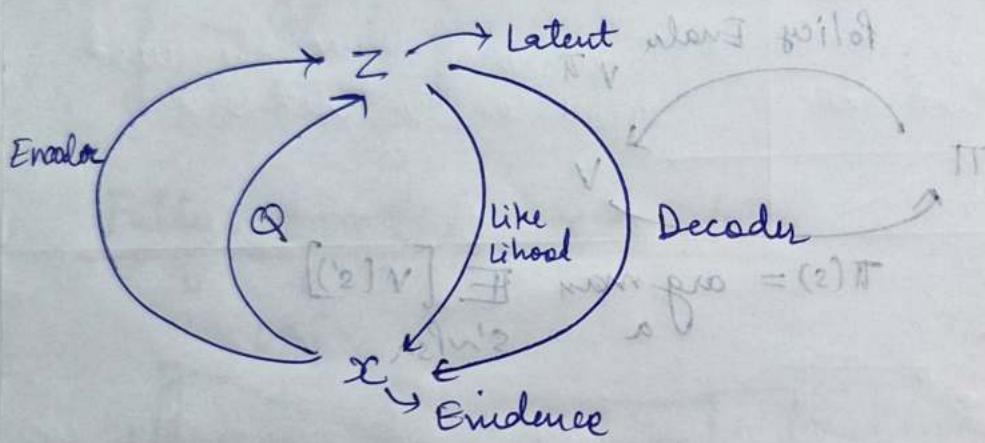
E-Step

$$Q(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

M Step

$$\theta = \underset{\theta}{\operatorname{argmax}} \sum_{z \sim Q_i} E \left[\log \frac{P(x, z; \theta)}{Q_i(z)} \right]$$

ELBO

VAE

$$(2)V \mathbb{E}[\text{mem}] + (2)I = (2)V$$

$\text{mem} = 0$

interests easy habit

