

COL774 - ML

parags

- 1) Basic Concepts
- 2) Learning Specific algorithms
- 3) Hands on experience
- 4) Art of applying ML

} What to
Expect

Group - 5, S.No 112
TA - Vishal

What is ML?

- ↳ identifying patterns
- predict future behaviour

Why ML now?

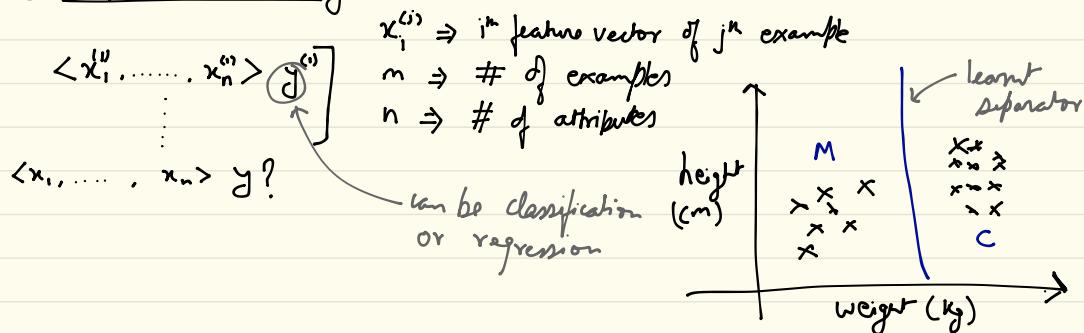
- ↳ Large amounts of data.
- Inc. in computation power
- Well founded models.

Types of ML algorithms -

- ↳ Supervised Learning

05/01/18

① Supervised Learning



We are looking for a hypothesis ($h: X \rightarrow Y$) $\in H$

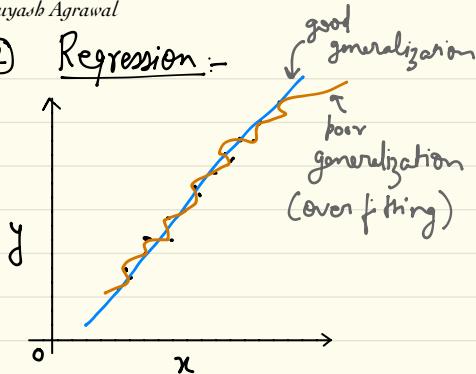
all possible input combination $\xrightarrow{\text{hypothesis class}}$

- 1) Which class of hypothesis?
- 2) Metric for being optimal
- 3) How to find the fⁿ?
- 4) Dealing with outliers.

Algos

- SVM
- Logistic Regression
- GDA

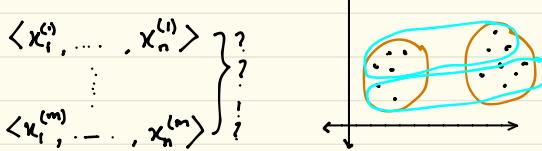
② Regression :-



$$h: X \rightarrow Y$$

Area	Prob.
2140	2
1100	1.5
.	.
.	.

③ Unsupervised Learning :-



- 1) How many clusters?
- 2) Metric for quality of clustering

Algo:-

• K-Means

④ Semi Supervised Learning :-

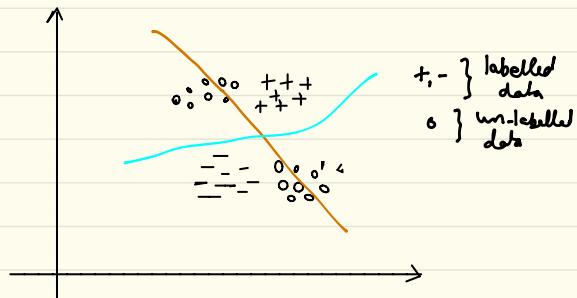
$$\langle x_1^{(1)}, \dots, x_n^{(1)} \rangle \quad y^{(1)}$$

\vdots

$$\langle x_1^{(i)}, \dots, x_n^{(i)} \rangle ?$$

\vdots

$$\langle x_1^{(m)}, \dots, x_n^{(m)} \rangle \quad y_m$$

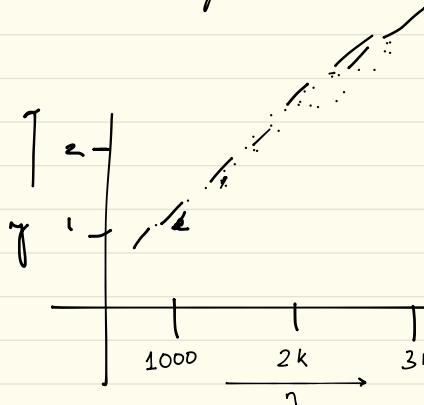


⑤ Reinforcement Learning

S- states, R-reward "Explore" / "Exploit"
 A- actions

Video Show

Linear Regression



$$h: X \rightarrow Y$$

given

$$\text{Aim: } h(x) \approx y$$

- 1) Sq error is differentiable. } Why choose sq error?
 2) hypothesis class

↳ set of all linear functions

$$h(x) = \theta_1 x + \theta_0$$

$$\langle x^{(i)}, \dots, x_n^{(i)} \rangle \quad y^{(i)}$$

$\{x^{(i)}, y^{(i)}\}_{i=1}^m$ ↳ one example
 ↳ training examples

$h: X \rightarrow Y$ | $h(x)$ should be as close to y as possible

$h \in H$ ↳ hypothesis class

$$h_{\theta}(x) = \theta_1 x + \theta_0 = \theta_2 x_2 + \theta_1 x_1 + \theta_0 \quad (2D)$$

$$\begin{aligned} \theta &= nx/|X| & \theta_0 &= 1/x \\ x &= nx/|X| & & \\ & & & \end{aligned}$$

$$\begin{aligned} &= \sum_{j=0}^n \theta_j x_j + \theta_0 \quad \text{dummy var. } x_0 = 1 \\ &= \sum_{j=0}^n \theta_j x_j \end{aligned}$$

We try to minimize :- Error = $\frac{1}{2} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2$

• How to minimize the error? $\left\{ \underset{\theta}{\operatorname{arg\,min}} J(\theta) \right\}$

Linear Regression:

- 1) Hypothesis Space
- 2) Loss function : $J(\Theta) = \frac{1}{2} \sum_{i=1}^m (h_{\Theta}(x^{(i)}) - y^{(i)})^2$
- 3) Optimization Algo.

) Vector Representation

- Default is column vector.

$$\Theta, x \in \mathbb{R}^{n+1} \quad x = \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix} \quad \Theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \quad h_{\Theta}(x) = \Theta^T x$$

$$J_{\Theta} = \frac{1}{2} \sum_{i=1}^m [\Theta^T x^{(i)} - y^{(i)}]^2 \quad \text{objective: } \underset{\Theta}{\operatorname{argmin}} J(\Theta) \quad \text{squared error.}$$

) Norm of a vector:-

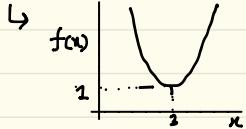
$$x \in \mathbb{R}^n, \quad \underbrace{\|x\|_p}_{\text{p-norm}} \text{ of } x \text{ is defined as} = \left[\sum_{j=1}^n |x_j|^p \right]^{1/p} \quad \|x\|_2 \Rightarrow 2 \text{ norm (L2)}$$

$$\|x\|_1 \Rightarrow 1 \text{ norm (L1)}$$

$L_0 \Rightarrow$ number of non zero entries

$L_{\infty} \Rightarrow \max_j |x_j|$

• How does one minimize a function?



$$f(x) = (x-3)^2 + 1$$

$$\underset{x}{\operatorname{argmin}} f(x) = 3$$

Methods:-

① find $\frac{\partial f(x)}{\partial x}$ & find zero point

then use $\frac{\partial^2 f(x)}{\partial x^2} > 0$ to check minima

$$\text{if } x \in \mathbb{R}^n : \quad \nabla_x f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}, \quad \text{at min/max} \quad \nabla_x f(x) = 0$$

Checking Min/Max :

Hessian (H) =

gives an idea of how
gradient is changing

✓ tells about curvature.

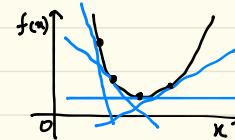
$$\begin{bmatrix} & & & \\ & \frac{\partial^2 f}{\partial x_1^2} & \dots & \\ & \vdots & & \end{bmatrix}$$

$$f(x) = x_1^2 x_2 + x_2$$

$$H = \begin{bmatrix} 2x_2 & 2x_1 \\ 2x_1 & 0 \end{bmatrix}$$

Gradient Descent

- It is not easy to find the roots of the gradient of function.
- Move opposite to gradient.



$$x^{t+1} = x^t - \eta \left. \frac{\partial f(x)}{\partial x} \right|_{x^t}$$

↑ Learning rate

- As we move towards minima the step size decreases & progress slows down.

Multivariate :-

$$x^{t+1} = x^{(t)} - \eta \cdot \nabla_x f(x)$$

Is η a vector/scalar or why?

12/01/18

• How to find $\arg \min_{\theta} J(\theta)$? $\hookrightarrow \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{2} \sum_{i=1}^m 2 [y^{(i)} - h_{\theta}(x^{(i)})] \underbrace{\left(-\frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)}) \right)}_{x_j^{(i)} \text{ (in linear regression)}}$

$$\therefore \theta_j^{t+1} = \theta_j^t + \eta \cdot \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] \cdot x_j^{(i)}$$

} Linear Regression update rule.

Vector: $\theta^{t+1} = \theta^t + \eta \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] x^{(i)}$

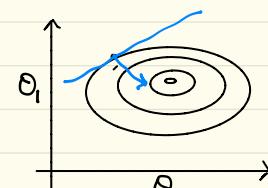
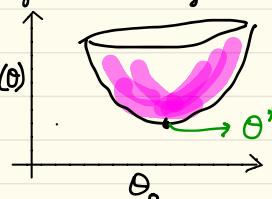
- Initialization
 - Learning rate
 - Convergence
- } Engineering Questions

Pseudo code for Gradient descent
use:- $\theta^{t+1} = \theta^t - \eta \nabla_{\theta} J(\theta)$

Understanding $J(\theta)$: (In case of linear regression)

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

Quadratic in θ_j



Contours of $J(\theta)$

- Initialization :-
• Randomly initializing close to zero works well (Practically).
- Learning Rate (η) :-
• Start slow & gradually increase until divergence signs
• Start large & gradually slow down.

- Converge :-

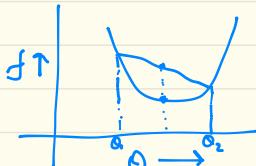
- * $|J(\theta^{t+1}) - J(\theta^t)| < \epsilon$
- * $|\theta^{t+1} - \theta^t| < \epsilon$ } Check L_∞ or L_2 norm.
- * $|\nabla_{\theta} J(\theta)| < \epsilon''$

If we choose ϵ to be very small then may take long time.

Lecture 07 - 16.01.18

- Stochastic Gradient Descent (SGD)

- Convexity: (Convex functions)



$f: R \rightarrow R$, f is convex iff:

$$\forall \theta_1, \theta_2 \in R \quad f(\alpha \theta_1 + (1-\alpha) \theta_2) \leq \alpha f(\theta_1) + (1-\alpha) f(\theta_2)$$

In strictly convex we have <

Convex functions have unique global minima. {at least one value}

If f is cont. & differentiable : convexity $\Leftrightarrow f''(x) \geq 0$
concave $\Leftrightarrow f''(x) \leq 0$

- For higher dimensions:

$f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is convex iff:

$$f(\alpha \theta_1 + (1-\alpha) \theta_2) \leq \alpha f(\theta_1) + (1-\alpha) f(\theta_2) \quad (0 \leq \alpha \leq 1) \quad \forall \theta_1, \theta_2 \in \mathbb{R}^n$$

2nd order condition:

$$\text{Hessian } H = \begin{bmatrix} & & & \\ & \ddots & & \\ & & \frac{\partial^2 f}{\partial \theta_i \partial \theta_j} & \\ & & & \ddots \end{bmatrix}_{n \times n}$$

Convex $\Leftrightarrow H \geq 0$ [H should be positive semi-definite]

Positive Semi-definite :-

Let $A \in \mathbb{R}^{n \times n}$ symmetric A is +ve semi-definite iff
 $\forall z \in \mathbb{R}^n \quad z^T A z \geq 0$

Eg. $I = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$

How to check if a matrix is positive semi-definite?
 ↳ Eigen values ≥ 0 ?

We can prove, $J_\theta = \frac{1}{2} \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})]^2$, is convex by calculating Hessian & showing it is +ve semi-definite.

H.W. calculate Hessian of $J(\theta) \Rightarrow X^T X \quad \{ \nabla_\theta J(\theta) = X^T x - X^T y \}$

Gradient Descent

$\theta = \text{init}()$

$t = 0$

do {

for $j = 1$ to n {

$$\theta_j^{t+1} = \theta_j^t + \eta \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] \cdot x_j^{(i)}$$

}

$t += 1$

} while (! converged)

m: examples, n: attributes

$O(mn)$ time expensive!

Stochastic Grad. Descent

- Examples are iid.

- We divide training set in batches

Batch size = r

$\theta = \text{init}()$

$t = 0; c = m/r$

do {

for $b = 1$ to c {

for $j = 1$ to n {

$$\theta_j^{t+1} = \theta_j^t + \eta \sum_{i=(b-1)r+1}^{br} [y^{(i)} - h_\theta(x^{(i)})] \cdot x_j^{(i)}$$

}

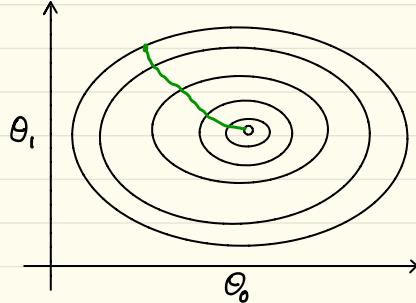
$t += 1;$

}

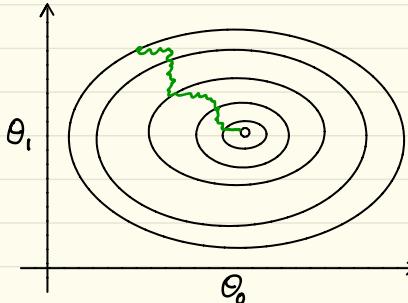
} while (! converged)

$O(rn)$ time

1 epoch: Seeing whole training data one time



Gradient Descent
 $\mathcal{O}(mnT)$



S.G.D.
 $\mathcal{O}(\tau n T')$

- $T \leq T'$ but iterations are faster & thus faster convergence
- Take smaller steps.
- Prove not converge to same minima.

Lecture 08 - 18.01.18

Heuristic: $\eta \propto \frac{1}{t}$ or $\eta \propto \frac{1}{\sqrt{t}}$ η : learning rate
 t: iterations done

Probabilistic Interpretation of Least Squares:

$$J_{\theta} = \frac{1}{2} \sum_{i=1}^m [\theta^T x^{(i)} - y^{(i)}]^2 \quad ? \text{ why? } \quad \text{Also, why vertical dist \not= dist from line?}$$

Generating the data :-

Assume underlying parameters θ

Input: $x^{(i)}$

① Apply Linear model: $\theta^T x^{(i)}$

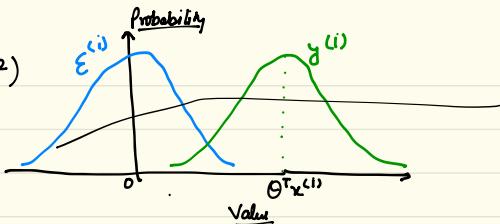
② $\epsilon^{(i)} \sim N(0, \sigma^2)$ (i.i.d noise)

③ $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$

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

$$P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\epsilon^{(i)} - \mu)^2}{2\sigma^2}}$$

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

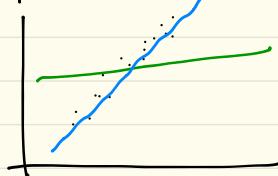


Maximum Likelihood Estimation - (ML Estimation)

Data = $\{x^{(i)}, y^{(i)}\}_{i=1}^m$, Estimate a set of parameters : θ

$$\boxed{\theta_{ML} = \underset{\theta}{\operatorname{argmax}} P(D; \theta)}$$

max. likelihood



Blue line has more probability of generating the data than Green line

$$P(y^{(1)}, \dots, y^{(m)} | x^{(1)}, \dots, x^{(m)}, \theta) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

$L(\theta)$ = Likelihood

$$\underset{\theta}{\operatorname{argmax}} L(\theta) = \underset{\theta}{\operatorname{argmax}} \text{LL}(\theta)$$

To log Likelihood.

$$\text{LL}(\theta) = \sum_{i=1}^m \left(\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2} \left(\frac{y^{(i)} - \theta^T x^{(i)}}{\sigma} \right)^2 \right)$$

$$\therefore \underset{\theta}{\operatorname{argmax}} \text{LL}(\theta) = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

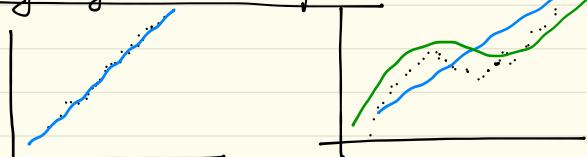
{ Removing the constants & taking }
argmin because of -

ignore, because its a +ve constant

$$= \underset{\theta}{\operatorname{argmin}} J(\theta)$$

\therefore Least squared minimization is equivalent to max. likelihood estimation under the assumption of i.i.d gaussian noise.

Locally weighted Linear Regression :



$$h_{\phi}(x) = \sum_{r=0}^k \theta_r x^r$$

Forwards as linear

$\phi : x \rightarrow \langle 1, x^2, x^3, \dots, x^k \rangle$
Kernel Trick?

We can also estimate data as being piecewise linear.

$$J(\theta) = \sum_{i=1}^m w^{(i)} [y^{(i)} - \theta^T x^{(i)}]^2, \quad w^{(i)} = f(x^{(i)}, x) = e^{-\frac{(x - x^{(i)})^2}{2\sigma^2}}$$

Locally weighted Linear Regression :

① Use Kernels : $x \rightarrow \mathcal{O}(n)$

$$\textcircled{2} \quad J(\theta) = \sum_{i=1}^m w^{(i)} [y^{(i)} - \theta^T x^{(i)}]^2$$

$$w^{(i)} = e^{-\frac{[x-x^{(i)}]^2}{2\sigma^2}}$$

point of interest

$\sigma \rightarrow \infty$ } regular linear regression

$\sigma \rightarrow 0$ } only x would have weight 1.

} ways to solve

- Underfitting when your hypothesis space is not able to fit the data.

- Lazy Learners** : no work done during training { this, K-nn }

- For every point we create a new loss function $J(\theta)$

- This is a **non-parametric method**.

- As new point comes, we create a new $J(\theta)$ for this of minimize this to get θ & then we predict y i.e. for every point we learn a new θ .

$$\text{for } x \in \mathbb{R}^n: \quad w^{(i)} = e^{-\frac{[(x^{(i)} - x)^T \Sigma^{-1} (x^{(i)} - x)]}{2}}$$

Linear Regression:

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

Analytical Solution for θ parameter:

$$\nabla_{\theta} J(\theta) = 0$$

$$\text{Design matrix: } X = \mathbb{R}^{m \times n} = \begin{bmatrix} \dots & x^{(1)\top} & \dots \\ \dots & x^{(2)\top} & \dots \\ \dots & x^{(m)\top} & \dots \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}, \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\therefore J(\theta) = \frac{1}{2} [(x\theta - Y)^T (x\theta - Y)] = \frac{1}{2} (\|x\theta - Y\|_2)^2$$

$$\textcircled{1} \quad \nabla_{\theta} \alpha^T \theta = \alpha$$

$$\nabla_{\theta} \theta^T \alpha = \alpha$$

number $\rightarrow \mathbb{R}^n$

$$\textcircled{2} \quad \text{a) } \nabla_{\theta} \theta^T A \theta = 2A\theta, \quad A \in \mathbb{R}^{n \times n} \text{ Symmetric}$$

$$\text{b) } \nabla_{\theta}^2 (\theta^T A \theta) = 2A$$

Show.

matrix of 2nd order derivatives, (Hessian)

What if not symmetric?

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

derive

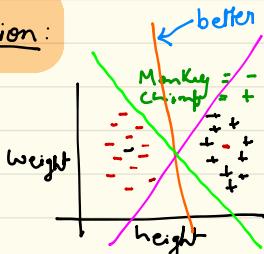
$$\theta = [(X^T X)^{-1} X^T] Y$$

pseudo inverse

Lecture 10 - 23.01.18

Logistic Regression:

Binary classification:



Linear Classifier { Decision boundary is linear }

Metric for good line ?

↳ Separation should be large in classification as large separation means more confidence

$g(z)$ (confidence for chimpanzee)

0.5

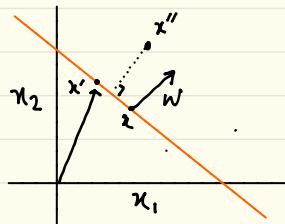
Logistic function

(also called)
sigmoid

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

z (signed distance)

Equation of line: $\theta^T X = 0$



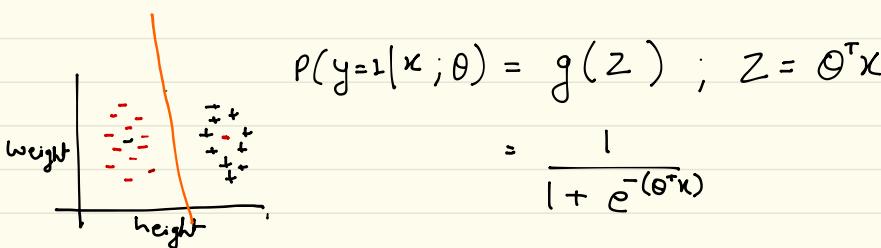
$$(w, b) \in \Theta \quad | \quad w_2 x_2 + w_1 x_1 + b \equiv \theta_2 x_2 + \theta_1 x_1 + \theta_0 x_0 = 0$$

$$\begin{aligned} \text{eqn: } & w^T(x - x') = 0 \\ & = w^T x - \underbrace{w^T x'}_b = 0 \\ & \rightarrow w^T x + b = 0 \end{aligned}$$

$$\begin{aligned} w^T(x'' - x') &= \text{un-normalized dist of } x'' \text{ from line} \\ &= w^T x'' + b \end{aligned}$$

\therefore To calculate the dist of point from line $w^T x + b = 0$, just put it in the equation gives un-normalized signed distance.

Process:



Generative Process:

1) Compute $z = \theta^T x$

2) Compute $P(y=1|x; \theta) = \frac{1}{1 + e^{-\theta^T x}}$

3) Sample $y \sim \text{Bernoulli}(\phi)$

\hookrightarrow Generate random number in $[0, 1]$ & if $< \phi$ then 1 else 0

How to sample from gaussian dist?

$$\left. \begin{array}{l} \text{Bernoulli } (\phi) \\ \therefore \phi = \frac{1}{1 + e^{-\theta^T x}} \end{array} \right\}$$

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

$$\therefore \Theta_{ML} = \underset{\theta}{\operatorname{argmax}} (L(\theta)) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} (\log(L(\theta))) = \sum_{i=1}^m \log(P(y^{(i)} | x^{(i)}; \theta))$$

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

Logistic Regression (Contd.):

Lecture 11 - 25.01.18

$$\frac{1}{1+e^{-z}} = g(z) \quad \boxed{g'(z) = g(z) \cdot (1-g(z))}$$

$$\underset{\theta}{\operatorname{argmax}} (\log(L(\theta))) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log(P(y^{(i)} | x^{(i)}; \theta)) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m y^{(i)} \cdot \log(g(x^{(i)}; \theta)) + (1-y^{(i)}) \cdot \log(1-g(x^{(i)}; \theta))$$

$L(\theta)$ is a concave function in θ .

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m y^{(i)} \cdot \log(g(\theta^T x^{(i)})) + (1-y^{(i)}) \cdot \log(1-g(\theta^T x^{(i)}))$$

$$\therefore \frac{\partial L(\theta)}{\partial \theta_j} = \sum_{i=1}^m \frac{y^{(i)} \cdot g(\theta^T x^{(i)}) \cdot (1-g(\theta^T x^{(i)})) \cdot x_j^{(i)}}{g(\theta^T x^{(i)})} + \frac{(1-y^{(i)}) \cdot g(\theta^T x^{(i)}) \cdot (1-g(\theta^T x^{(i)})) \cdot x_j^{(i)}}{(1-g(\theta^T x^{(i)}))}$$

$$= \sum_{i=1}^m [y^{(i)} \cdot (1-g(\theta^T x^{(i)})) - (1-y^{(i)}) \cdot g(\theta^T x^{(i)})] x_j^{(i)}$$

$$= \sum_{i=1}^m [y^{(i)} - g(\theta^T x^{(i)})] x_j^{(i)}$$

Here, we define $h_{\theta}(x) = P(y=1 | x; \theta) = g(\theta^T x)$

$$\therefore \frac{\partial L(\theta)}{\partial \theta_j} = \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] x_j^{(i)} \quad \left. \right\} \text{Similar to linear regression}$$

In fact we see that in both cases $h_{\theta}(x^{(i)})$ turns out to be the expected value of dist. of $y^{(i)}$

E.g. in linear reg: $\frac{\theta^T x^{(i)}}{E(y)}$ for normal \rightarrow , in logistic: $g(\theta^T x^{(i)}) \rightarrow E[y]$ for Bernoulli

$$\nabla_{\theta} LL(\theta) = \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x^{(i)}$$

- Softmax Regression : Multiclass classification

- Generalized Linear Models (GLMs):

This is a class of models

→ Linear Regression
→ Logistic Regression
↓

Training data: $\{y^{(i)}, x^{(i)}\}_{i=1}^m$

GLM :- Models: $P(y|x; \theta) = P(y; \eta)$ natural parameter

Conditions:

- ① $P(y; \eta) \in$ Exponential Family dist.
- ② $\eta = \theta^T x$ (η is linearly dependent on $\theta^T x$)
- ③ $h_{\theta}(x) = E[y; \eta]$

- Exponential Family distributions:

$$P(y; \eta) = b(y) \cdot \exp(\eta y - a(\eta)) = \frac{b(y) \cdot \exp(\eta y)}{\exp(a(\eta))}$$

$$\cdot y \sim N(\mu, 1) \Rightarrow P(y) = \frac{1}{\sqrt{2\pi}} \cdot e^{\frac{(y-\mu)^2}{2}} = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{[y^2 + \mu^2 - 2y\mu]}{2}}$$

$$\therefore \eta = \mu$$

$$b(y) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{y^2}{2}}$$

$$a(\eta) = \frac{\mu^2}{2}$$

Normal dist
also \in Exp. family dist

- Similarly, we can show that Bernoulli dist also \in Exp. family dist.

- Try with σ parameter in $N(\mu, \sigma^2)$

Numerical Algo: SVD

16/01/18

$$A: m \times n, \quad A: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\text{Nullspace}(A) = \{x : Ax = 0\}, \quad \text{Range}(A) = \{Ax : x \in \mathbb{R}^n\}$$

$$\dim(\text{Nullspace}(A)) + \dim(\text{Range}(A)) = n$$

$$\begin{array}{l} \downarrow \\ e_1, \dots, e_n \xrightarrow{\text{extend}} e_{n+1}, \dots, e_n \\ \downarrow \\ Ae_{n+1}, \dots, Ae_n \\ \text{basis of Range}(A) \end{array}$$

$$\dim(\text{Range}(A)) = \text{col. rank}(A)$$

$$[A^1 \cdots A^n] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = x_1 A^1 + \cdots + x_n A^n$$

2 norm of diag matrix:

$$A = \begin{bmatrix} d_1 & & 0 \\ & d_2 & \\ 0 & & \ddots & d_n \end{bmatrix} \quad \|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 = (d_1^2 x_1^2 + d_2^2 x_2^2 + \cdots + d_n^2 x_n^2)^{1/2} = \max\{d_1, \dots, d_n\}$$

① Suppose U is an invertible $n \times n$ matrix. Then, $\text{rank}(A \cdot U) = \text{rank}(A)$

$$\begin{array}{l} \text{range}(A) = \text{range}(AU) \\ \downarrow \\ y = Ax \\ = AU \cdot (U^{-1}x) \\ \Rightarrow y \in \text{range}(AU) \end{array} \quad \begin{array}{l} \uparrow \\ y = AUx \\ = A(Ux) \\ \downarrow \\ y \in \text{range}(A) \end{array}$$

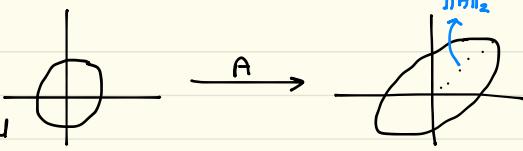
② V : $m \times m$ inv. matrix : $\text{rank}(VA) = \text{rank}(A)$

$$\begin{array}{l} \text{nullspace}(A) = \text{nullspace}(VA) \\ \downarrow \\ Ax = 0 \\ VAx = 0 \\ VAx = 0 \quad (\text{multiply by } V^{-1}) \\ \therefore x \in \text{nullspace}(VA) \end{array} \quad \begin{array}{l} \downarrow \\ Ax = 0 \\ An = 0 \\ \therefore n \in \text{nullspace}(A). \end{array}$$

③ Suppose U is unitary ($U^T U = I$)

$$\|AU\|_2 = \|A\|_2$$

Multiplying unitary = changing orthonormal basis.



$$\|AU\|_2 = \max_{\|x\|_2=1} \|AUx\|_2$$

$$\text{Let } y = Ux, \quad \|y\|_2 = 1$$

$$= \max_{\|y\|=1} \|Ay\|_2$$

$$= \|A\|_2$$

$$\|AU\|_1 = \|A\| ? \quad \times$$

$$\|AU\|_F = \|A\|_F \quad \checkmark \quad (\text{Frobenius norm})$$

④ V : $m \times n$ Unitary $\|VA\|_2 = \|A\|_2$ (as in rotation of ellipse the length of points remain same)

$$\|V^T y\|_2 = \|y\|_2$$

$$\therefore \|VA\|_2 = \max_{\|x\|_2=1} \|VAn\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 = \|A\|_2$$

$$\text{Also, } \|VA\|_F = \|A\|_F$$

$$\|A\|_F^2 = \|A^1\|_2^2 + \|A^2\|_2^2 + \dots + \|A^n\|_2^2$$

$$\|VA\|_F^2 = \|VA^1\|_2^2 + \dots = \underbrace{\|A^1\|_2^2}_{\|A\|_F^2}$$

SVD Theorem: A : $m \times n$ matrix

Then \exists unitary matrices U : $m \times m$

V : $n \times n$

$$U = [U^1 \ U^2 \ \dots \ U^m]$$

right singular vectors

$$V = [V^1 \ V^2 \ \dots \ V^n]$$

left singular vectors

s.t. $A = U \Sigma V^T$

$$\Sigma: \text{diag matrix} : \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \\ & & & \sigma_{\min(m,n)} \end{bmatrix}$$

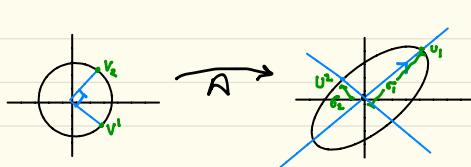
$\sigma_1 \geq \sigma_2 \dots \geq \sigma_{\min(m,n)} \geq 0$
(singular values)

and are uniquely determined by A

$\Rightarrow AV = U\Sigma$, Now $\|A\|_2 = \sigma_1$ {as V & U are unitary}

$$AV' = \sigma_1 U'$$

$$AV^2 = \sigma_2 U^2$$



$$\|x\|_2 = 1, X = \lambda_1 V^1 + \lambda_2 V^2, \lambda_1^2 + \lambda_2^2 = 1$$

$$AX = \lambda_1 \sigma_1 U^1 + \lambda_2 \sigma_2 U^2$$

$$\|AX\|_2 = (\lambda_1^2 \sigma_1^2 + \lambda_2^2 \sigma_2^2)^{1/2} \in [\sigma_2, \sigma_1]$$

- If A is singular then $\sigma_2 = 0$ & ellipse will be a straight line.
- Rank (A) = rank (Σ) = number of non zero singular values.

Proof: Define $\sigma_i = \|Av_i\|_2$

Let v' be the unit vector s.t. $\|Av'\|_2 = \sigma_1$

Let U' be the unit vector s.t. $AV' = \sigma_1 U'$

Let L : (set of vector \perp to v') = $\{x : \langle x, v' \rangle = 0\}$

M: $\{x | \langle x, v' \rangle = 0\}$

Key claim: If $x \in L$, $AX \in M$

Suppose the claim is false:

$\exists x \in L$, s.t.

$AX = \alpha U^1 + \beta U^2, U \in M, \alpha \neq 0, x$ is unit vector

Consider the vector $x' = x + \epsilon v'$

We will show that $\frac{\|Ax'\|}{\|x'\|} > \sigma_1$

$$Ax' = \alpha U^1 + \beta U^2 + \epsilon \sigma_1 U'$$

$$= (\alpha + \epsilon \sigma_1) U^1 + \beta U^2$$

$$\frac{\|Ax'\|}{\|x'\|} = \frac{(\alpha + \epsilon \sigma_1)^2 + \beta^2}{1 + \epsilon^2} > \frac{\sigma_1^2 \epsilon^2 + 2\alpha \epsilon \sigma_1}{1 + \epsilon^2}$$

$$= \sigma_1^2$$

if we take $\epsilon = \frac{\sigma_1}{2\alpha}$

GLM:

$P(y; \eta) \xrightarrow{f_n(\theta, x)}$ Exponential Family

① $\eta = \theta^T x$ } linearity

② $h_\theta(x) = E[y | x; \eta]$

• Exponential dist.:

$$P(y; n) = b(y) \cdot \exp(ny - a(n))$$

① $\hookrightarrow y \sim N(\mu, 1)$ then $P(y) \in$ Exponential family.

$$\eta = \mu$$

② $\hookrightarrow y \sim \text{Bernoulli}(\phi)$, then $P(y) \in$ Exp. family

Hohoh $\phi = \frac{1}{1 + e^{-\eta}}, \eta = \theta^T x$

all of
these is
 $y|x$

• Gradient update rule for GLMs :-

$$LL(\theta) = \sum_{i=1}^m \log(P(y^i | x^i; \theta))$$

$$= \sum_{i=1}^m \log(P(y^i; n^{(i)}))$$

$$= \sum_{i=1}^m \log(b(y^{(i)}) \cdot [\exp(z^{(i)}y^{(i)} - a(n^{(i)})])]$$

$$= \sum_{i=1}^m \log(b(y^{(i)})) + \sum_{i=1}^m [n^{(i)} \cdot y^{(i)} - a(n^{(i)})]$$

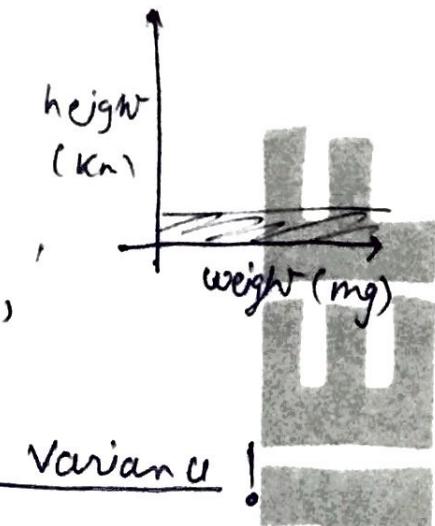
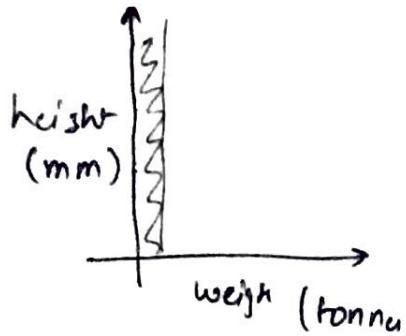
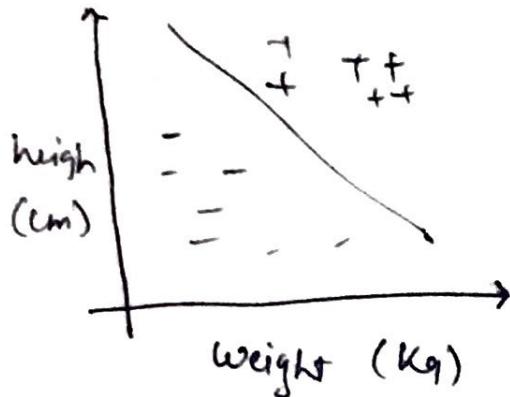
$$\frac{\partial LL(\theta)}{\partial \theta_j} = 0 + \sum_{i=1}^m (y^{(i)} \cdot x_j^{(i)} - a'(n^{(i)}) \cdot x_j^{(i)})$$

$$= \sum_{i=1}^m [y^{(i)} - a'(n^{(i)})] \cdot x_j^{(i)}$$

linear or logistic
for $a'(n^{(i)})$ becomes $h(\theta^T x^{(i)})$

• Convexity?

• Normalization: $\{x^{(i)}, y^{(i)}\}_{i=1}^m$



Make each attribute 0 mean & unit variance!

for $j = 1 \text{ to } n \}$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^i$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^i - \mu_j)^2$$

{ m or $m-1$? [Doesn't matter]

}

for $j = 1 \text{ to } n \}$

$$\text{# i: } x_j^i = \left[\frac{x_j^i - \mu_j}{\sigma_j} \right]$$

{

$$\text{Var}(x) = \sigma^2$$

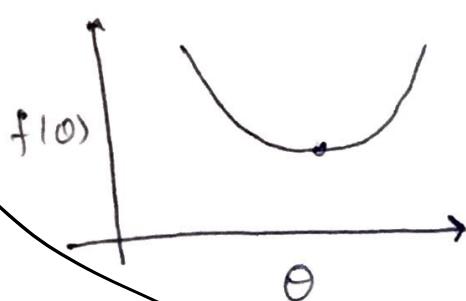
$$\text{Var}(x_k) = \frac{\sigma^2}{K^2}$$

* We almost always should normalize the data we are working on [even in closed form sol''].

• Should we also normalize y ? [Not really needed].

- When not normalize? → When data comes from same metric, we don't do normalization. E.g. Naive Bayes classification we don't normalize data.

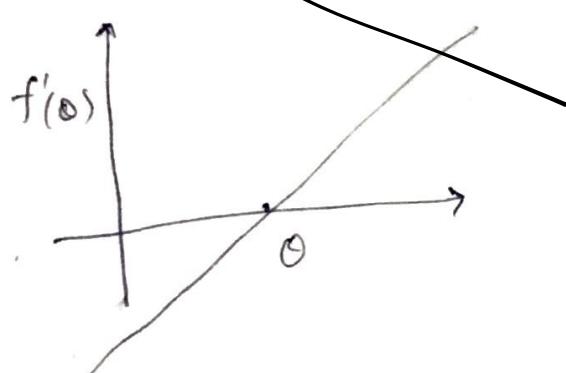
* Newton's Method.



$$f(\theta) = (\theta - 3)^2 + 1$$

$$f'(\theta) = 2(\theta - 3)$$

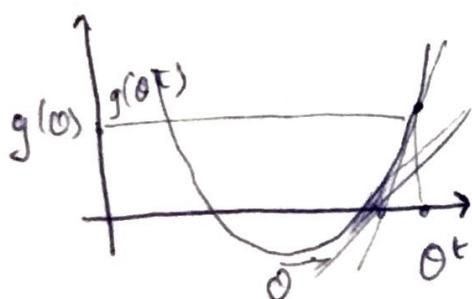
- ① Analytic
- ② Gradient Descent
- ③ ~~Newton's meth.~~



To optimize we find zero of $f'(\theta)$

~~θ⁰ θ¹ θ² θ³ θ⁴ θ⁵ θ⁶~~

* Finding zero of a function:



iterative process

$$\frac{\theta - g(\theta^t)}{\theta^{t+1} - \theta^t} = g'(\theta^t)$$

$$\Rightarrow \theta^{t+1} = \theta^t - \frac{g(\theta^t)}{g'(\theta^t)}$$

Newton's update

• BFGS
• LBFGS

for optimization: (G)

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f'(\theta^t)}{f''(\theta^t)} \quad \boxed{\theta \in \mathbb{R}}$$

if $\theta \in \mathbb{R}^n$:

$$\underline{\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla_{\theta} f(\theta^{(t)})}$$

(↳ concrete Proof? (See piazzza))

• Hessian ~~non~~ singular?

•

CODE

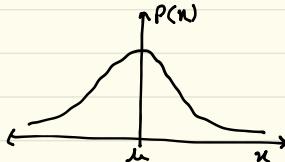
Last class: • GLM

- Normalization
- Newton's Method

GDA: Gaussian Discriminant analysis {Classification Algorithm}Normal Distribution:

$$x \in \mathbb{R}, \quad x \sim N(\mu, \sigma^2)$$

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



- For $x \in \mathbb{R}^n$:

$$x = (x_1, \dots, x_n)$$

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

$\begin{matrix} \text{n sized} \\ \text{vector} \\ \in \mathbb{R}^n \end{matrix}$ covariance matrix ($n \times n$)

$$P(x=x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \cdot e^{-\frac{1}{2} ((x-\mu)^T \Sigma^{-1} (x-\mu))}$$

- If x is a R.V. : $E[x] = \int P(x=x) \cdot x \, dx$
- $(x \in \mathbb{R}) \quad \text{Var}[x] = E[(x - E(x))^2] = E(x^2) - E(x)^2$
- x_1, x_2 is R.V. $\Rightarrow \text{Cov}(x_1, x_2) = E[(x_1 - E[x_1])(x_2 - E[x_2])] = E[x_1 \cdot x_2] - E[x_1] \cdot E[x_2]$

- Independent $\Rightarrow 0$ Covariance

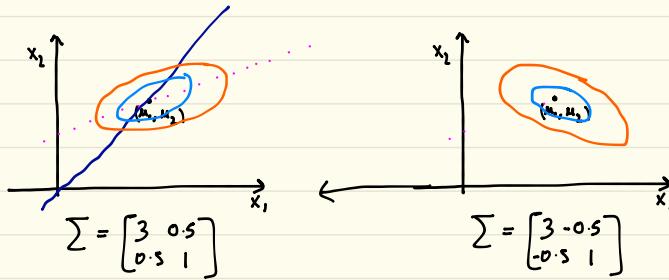
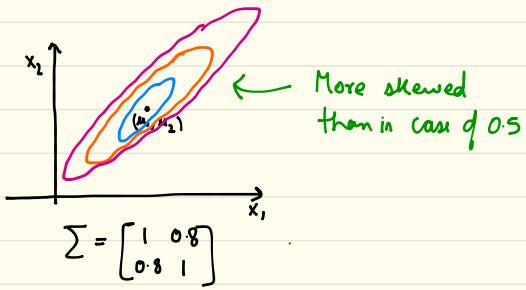
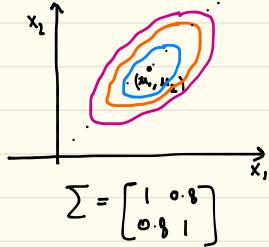
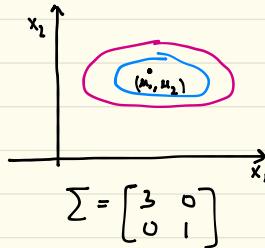
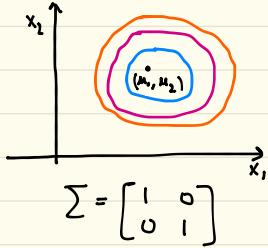
- If $x = (x_1, \dots, x_n)$, then $\Sigma = \text{Covariance matrix}$

$$= \begin{bmatrix} & & & i \\ & & \vdots & \\ & \dots & & \Sigma_{ii} \\ j & \dots & & \end{bmatrix} \quad \Sigma_{ij} = \text{Cov}(x_i, x_j)$$

Properties of Σ :

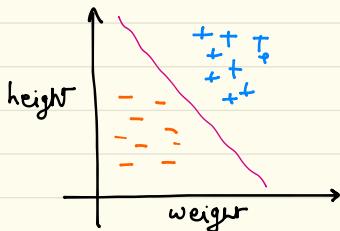
- Symmetric
- Positive semi definite {Prove}
- $\Sigma_{ii} = \text{Var}(X_i)$
- If Σ is diagonal, then all features are un-correlated.

Eg: $X = (x_1, x_2)$, $\mu = (\mu_1, \mu_2)$



See the gda pdf [CS229N] for better visualization.

Gaussian Discriminant Analysis:



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

Logistic:

Given x_i we model $P(y^i|x_i; \theta)$

$y^{(i)} \sim \text{Bernoulli}(\phi)$

$$P(y=1|x; \theta) = \frac{1}{1 + e^{-\theta^T x}}$$

- In GDA we try to generate the entire data { constraint to Logistic where we generated y_i given x_i }
- This is an example of Generative Model.

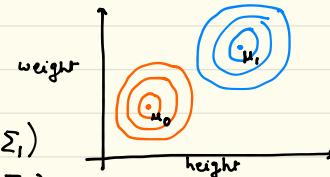
$P(x, y; \theta)$

↳ a) $y \sim \text{Bernoulli}(\phi)$

b) Model $P(x|y=1)$ & $P(x|y=0)$

Here we make the assumption: $P(x|y=1) \sim N(\mu_1, \Sigma_1)$

$$P(x|y=0) \sim N(\mu_0, \Sigma_0)$$

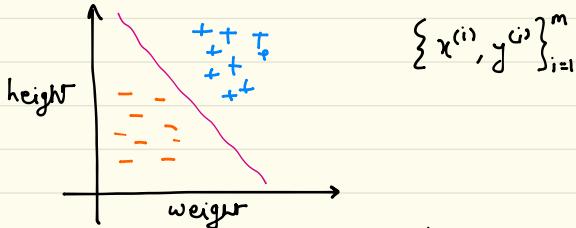


$$LL(\theta) = \sum_{i=1}^m \log(P(x^{(i)}, y^{(i)}; \theta)) = \sum_{i=1}^m \log(P(y^{(i)}; \theta) \cdot P(x^{(i)}|y^{(i)}; \theta))$$

We try to find $\operatorname{argmax}_\theta (LL(\theta))$

MINOR: 2:25pm - 3:35pm , 1:10 hour, Closed Notes

• Gaussian Discriminant Analysis:



We want to model $P(y, x; \theta)$.

$$\cdot LL(\theta) = \sum_{i=1}^m \log(P(x^{(i)}, y^{(i)}; \theta)) = \sum_{i=1}^m \log(P(y^{(i)}; \theta) \cdot P(x^{(i)}|y^{(i)}; \theta))$$

$$\Theta = (\phi, \mu_0, \Sigma_0, \mu_1, \Sigma_1)$$

$$= \sum_{i=1}^m \log(P(y^i; \phi)) + \log(P(x^i|y^i; \Theta))$$

$$= \sum_{i=1}^m y^i \cdot [\log(P(y^i=1; \phi)) + \log(P(x^i|y^i=1; \mu_1, \Sigma_1))] \\ + (1-y^i) \cdot [\log(P(y^i=0; \phi)) + \log(P(x^i|y^i=0; \mu_0, \Sigma_0))]$$

To maximize $LL(\theta)$: $\nabla_\phi LL(\theta) = 0, \nabla_{\mu_1} LL(\theta) = 0$

$$\nabla_{\mu_0} LL(\theta) = 0, \nabla_{\Sigma_0} LL(\theta) = 0, \nabla_{\Sigma_1} LL(\theta) = 0$$

Optimal parameters: $\phi = \frac{\sum_{i=1}^m 1\{y^i=1\}}{m}$ $(1\{<\text{bool}>\} = 1 \text{ if bool is T})$ 0 otherwise

$$\mu_1 = \frac{\sum_{i=1}^m 1\{y^i=1\} \cdot x^{(i)}}{\sum_{i=1}^m 1\{y^i=1\}}, \quad \mu_0 = \frac{\sum_{i=1}^m 1\{y^i=0\} \cdot x^{(i)}}{\sum_{i=1}^m 1\{y^i=0\}}$$

$$\Sigma_1 = \frac{\sum_{i=1}^m 1\{y^i=1\} (x^{(i)} - \mu_1)(x^{(i)} - \mu_1)^T}{\sum_{i=1}^m 1\{y^i=1\}}$$

$$\Sigma_0 = \frac{\sum_{i=1}^m 1\{y^i=0\} (x^{(i)} - \mu_0)(x^{(i)} - \mu_0)^T}{\sum_{i=1}^m 1\{y^i=0\}}$$

How to do classification using this?

↪ We want to calculate $P(y|x; \theta)$

$$= P(x|y; \theta) \cdot P(y; \theta)$$

$$\frac{P(x; \theta)}{\rightarrow \text{same for } y=0} \cdot P(y=0)$$

$$+ P(x|y=1; \theta) \cdot P(y=1)$$

height



Decision boundary

• Decision Boundary:

$$1 = \frac{P(y=1|x; \theta)}{P(y=0|x; \theta)} = \frac{P(x|y=1) \cdot P(y=1)}{P(x|y=0) \cdot P(y=0)} = \frac{\frac{\phi}{(2\pi)^{N_x/2} |\Sigma_1|^{1/2}} \cdot e^{-\frac{(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)}{2}}}{\frac{(1-\phi)}{(2\pi)^{N_x/2} |\Sigma_0|^{1/2}} \cdot e^{-\frac{(x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0)}{2}}}$$

Take denominator to other side & take log.

$$\log(1-\phi) - \log((1|\Sigma_0|^{\frac{1}{2}}) - \frac{(x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0)}{2}) = \log(\phi) - \log((1|\Sigma_1|^{\frac{1}{2}}) - \frac{(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)}{2})$$

$$\frac{1}{2} \left[(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0) \right] = \underbrace{\log\left(\frac{\phi}{1-\phi}\right)}_{\text{constant (indep. of } x)} - \underbrace{\log\left(\frac{1|\Sigma_1|^{\frac{1}{2}}}{1|\Sigma_0|^{\frac{1}{2}}}\right)}_{\text{Quadratic in } x}$$

Quadratic Expression: $x^T \Sigma_1^{-1} x - x^T \Sigma_0^{-1} x$

∴ If $\Sigma_1 = \Sigma_0$ then quad. term = 0 & thus boundary is linear.

H.W.: Compute $P(y=1|x; \theta)$ with $\Sigma_0 \neq \Sigma_1$ as same

↪ Turns out to be $\frac{1}{1 + e^{-\theta^T x}}$

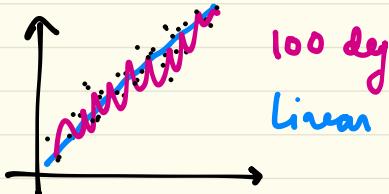
- Minor 1 Copy distribution : 1:00-2:30, SIT006. Grp 1,2 : 1-1:30
3,4 : 1:30-2 , 09.02.18
5,6 : 2-2:30

• Question Discussion:

Q.3) Checking convergence of S-G-D.

- Take difference in cost of consecutive iterations Too jittery
- Take difference of moving average over last r examples? Smoother, but may still be jittery
- Have a band of ϵ when if the last β batches have error b/w this then we say that we have converged.

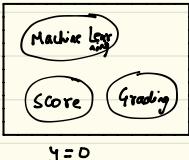
Q4) b) • If we have less data then 1) classes are really distributed as gaussian then GDA would give us much better generalization. E.g. Linear VS 100 deg poly fit in a linearly spread data



MINOR-2

• Naive Bayes:

Eg. Spam Classification:



i^{th} Example: $\langle x_1^{(i)}, x_2^{(i)}, \dots, x_{N_i}^{(i)} \rangle, y^{(i)}$
 $x_j^{(i)} \in \{0, 1\} \Rightarrow j^{\text{th}}$ word present in i^{th} doc or not
 $y^{(i)} \in \{0, 1\} \Rightarrow i^{\text{th}}$ document is spam or not.

Training Examples: $\{x^{(i)}, y^{(i)}\}_{i=1}^m$

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{P(x^{(i)}, y^{(i)}; \theta)}{P(x^{(i)}; \theta)} = \frac{P(x^{(i)}|y^{(i)}; \theta) \cdot P(y^{(i)}; \theta)}{P(x^{(i)}; \theta)}$$

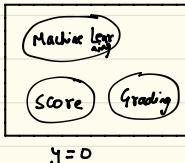
$y \sim \text{Bernoulli}(\phi)$

$\underbrace{x^{(i)}}_{2^{|V|}},$ since $2^{|V|}$ is humongous, we make a **very strong assumption**: $x_j \perp x_k \mid y$
↳ Naive Bayes Assumption

$$\therefore P(x^{(i)}|y^{(i)}; \theta) = \prod_{j=1}^{|V|} P(x_j^{(i)}|y^{(i)}; \theta) \quad \} \text{ due to assumption}$$

Last class:

- Minor 1 Discussion
- Naive Bayes Start

Naive Bayes:

ith Example: $\langle x_1^{(i)}, x_2^{(i)}, \dots, x_{IV}^{(i)} \rangle, y^{(i)}$
 $x_j^{(i)} \in \{0, 1\} \Rightarrow j^{\text{th}}$ word present in ith doc or not
 $y^{(i)} \in \{0, 1\} \Rightarrow$ ith document is spam or not

Bernoulli Event Model.

We are making a generative model.

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y=0 = (x_1, \dots, x_{IV}) | y=0$$

$$\therefore P(x|y) = P(x_1, \dots, x_{IV}|y) = \prod_{j=1}^{IV} P(x_j|y)$$

Assumption: x_1, \dots, x_{IV} are independent | y
Using assumption

$$\text{Now, } \forall_j x_j | y=1 \sim \text{Bernoulli}(\theta_j | y=1)$$

$$\forall_j x_j | y=0 \sim \text{Bernoulli}(\theta_j | y=0)$$

$$\# \text{params} = 2IV + 1$$

Note: Here we are giving up on the order of occurrence of words & frequency of occurring.

Log Likelihood:

$$\begin{aligned} LL(\theta) &= \sum_{i=1}^m \log P(y^{(i)}, x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log (P(y^{(i)}; \theta)) + \sum_{i=1}^m \log (P(x^{(i)} | y^{(i)}; \theta)) \\ &= \sum_{i=1}^m y^{(i)} \log \phi + (1-y^{(i)}) \log (1-\phi) \\ &\quad + \sum_{i=1}^m y^{(i)} \log \left(\underbrace{\sum_{j=1}^n P(x_j^{(i)} | y^{(i)}=1; \theta_j | y=1)}_{x_j^{(i)}, \theta_j | y=1} \right) + (1-y^{(i)}) \log \left(\sum_{j=1}^n P(x_j^{(i)} | y^{(i)}=0; \theta_j | y=0) \right) \end{aligned}$$

$$\underset{\theta}{\operatorname{argmax}} LL(\theta) \rightarrow \nabla_{\theta} LL(\theta)=0, \quad \forall_j \nabla_{\theta_j | y=1} LL(\theta)=0, \quad \forall_j \nabla_{\theta_j | y=0} LL(\theta)=0$$

$$\text{Parameters: } \theta = \phi, \{ \theta_j | y=0 \}_{j=1}^n, \{ \theta_j | y=1 \}_{j=1}^n$$

$$\text{After solving: } \phi = \frac{\sum_{i=1}^m 1 \{ y^{(i)}=1 \}}{m}, \quad \theta_j | y=1 = \frac{\sum_{i=1}^m 1 \{ y^{(i)}=1 \wedge x_j^{(i)}=1 \}}{\sum_{i=1}^m 1 \{ y^{(i)}=1 \}}, \quad \theta_j | y=0 = \frac{\sum_{i=1}^m 1 \{ y^{(i)}=0 \wedge x_j^{(i)}=1 \}}{\sum_{i=1}^m 1 \{ y^{(i)}=0 \}}$$

Prediction time:

x : document

$$\underset{y}{\operatorname{argmax}} P(y|x) = \underset{y}{\operatorname{argmax}} \frac{P(x|y) \cdot P(y)}{P(x)} = \underset{y}{\operatorname{argmax}} P(y) \cdot P(x|y) = \underset{y}{\operatorname{argmax}} \left[P(y) \cdot \prod_{j=1}^n P(x_j|y) \right]$$

↑ same thing

Warning: Never multiply.
Always take log.
to avoid underflows.

Note: For computing probabilities Naive Bayes may not be good, but for classification it works reasonably.

Note: Let us suppose that we have an exotic word w_i (say "archi lecture") & it never occurred in our training data. $\therefore P(w_i | y=0) = 0$

Now if a new document contains this word w_i then $P(d|y=0)$ will turn out to be 0!

To Fix this: { Laplace Smoothing }

$$\Theta_j|_{y=1} = \frac{\sum_{i=1}^m 1 \{ y^{(i)} = 1 \wedge x_j^{(i)} = 1 \} + c}{\sum_{i=1}^m 1 \{ y^{(i)} = 1 \} + 2c}, \quad \Theta_j|_{y=0} = \frac{\sum_{i=1}^m 1 \{ y^{(i)} = 0 \wedge x_j^{(i)} = 1 \} + c}{\sum_{i=1}^m 1 \{ y^{(i)} = 0 \} + 2c}$$

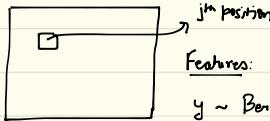
Intuition: Kind of adding two documents: one will all the words & other with no words

c is large \Rightarrow

c is small \Rightarrow

- In Bernoulli event model we ignored the count of the words.

- Alternative:



Features: X_j : For position j in the document which word occurs $\in \{1, 2, \dots, |V|\}$

$$y \sim \text{Bernoulli}(\theta)$$

$$X_{j,y} \sim \text{Multinomial}(\{\theta_{j,y}^k\}_{k=1}^{|V|}), \# \text{param} = 2(|V|-1) + 1$$

Since these are very large, we assume that all positions are independent

Assume:

$$\{\theta_{j,y}^k\} = \{\phi_{j,y}^k\} + j, j' \quad | \quad \sum \phi_{j,y}^k = 1, \# \text{param} = 2(|V|-1) + 1$$

Write expression for LLB & solve.

Optimal Params:

$$\phi = \sum_{i=1}^m \frac{1\{j^i=1\}}{m}$$

Bag of Words Model

$$\theta_{|y|=1}^k = \frac{\sum_{i=1}^m 1\{j^i=1\} n_k^{(i)}}{\sum_{i=1}^m 1\{j^i=1\} n^{(i)}} \quad \begin{matrix} \text{\# of times kth} \\ \text{word occurs} \\ \text{in i-th document} \end{matrix}, \quad \theta_{|y|=0}^k = \frac{\sum_{i=1}^m 1\{j^i=0\} n_k^{(i)}}{\sum_{i=1}^m 1\{j^i=0\} n^{(i)}}$$

We just need two arrays in order to implement this (each of $|V|$ size)

$$\text{Laplace Smoothing: } \theta_{|y|=0}^k = \frac{\sum_{i=1}^m 1\{j^i=0\} n_k^{(i)} + c}{\sum_{i=1}^m 1\{j^i=0\} n^{(i)} + c|V|}, \quad \theta_{|y|=1}^k = \frac{\sum_{i=1}^m 1\{j^i=1\} n_k^{(i)} + c}{\sum_{i=1}^m 1\{j^i=1\} n^{(i)} + c|V|}$$

\Downarrow adding c document will add words in both classes

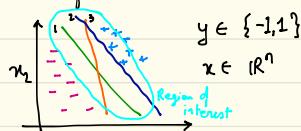
① Linear Regression: $y \in \mathbb{R}$, $x \in \mathbb{R}^n$

② Classification:

Logistic Regression	GDA	Naive Bayes
1) $y \in \{0, 1\}$ if y is discrete: Softmax	$y \in \{0, 1\}$ $y \in \{0, \dots, r\}$	$y \in \{0, 1\}$ $y \in \{0, \dots, r\}$
2) $x \in \mathbb{R}^n$ Logistic cannot give us $P(x)$!	$x \in \mathbb{R}^n$ <small>GDA subsumes Naive Bayes if $x \in \mathbb{R}^n$ is naive bayes & we use normal this is called gaussian naive bayes</small>	$\forall x_j \in \{0, \dots, V \}$ $x_j \in \mathbb{R} \mid P(x_j y)$ as some dist. [e.g. normal]

Support Vector Machines:

Classification:



① is better than ② & ③

Methods seen till now:

- Logistic Regression {Discriminative}, $P(y=1|x) = \frac{1}{1+e^{-\theta^T x}}$
- GDA {Generative}: $P(y,x) = P(y) \cdot P(x|y)$
Bernoulli
 $\begin{aligned} & N(\mu_0, \Sigma_0) \\ & N(\mu_1, \Sigma_1) \end{aligned}$
- Naive Bayes: $P(x,y) = P(y) \cdot \frac{P(x|y)}{\prod_i P(x_i|y)}$

Similarity: All of these are probabilistic models {define dist of data}. That's why we try to maximize log likelihood $LL(\theta)$.

• SVMs are more discriminative than Linear Models

• The hyperplane of interest is: $w^T x + b = 0$, $w \in \mathbb{R}^n$, $x \in \mathbb{R}^n$
instead of $\theta^T x = 0$, $\theta \in \mathbb{R}^{n+1}$, $x \in \mathbb{R}^n$

• Signed Distance of $x^{(i)}$: $\frac{(w^T x^{(i)} + b) \cdot y^{(i)}}{\|w\|} = \hat{\gamma}^{(i)}$, $y^{(i)} \in \{-1, 1\}$
geometric margin
functional margin

with no outliers

• We currently are assuming the data is linearly separable, but will relax this constraint later

• We try to: $\max_{w,b} (\min_i \gamma^{(i)})$

• Now: $\max_{w,b} (\min_i \gamma^{(i)}) \equiv \max_{w,b} \gamma, \gamma \leq \gamma^{(i)} + i$
 $\gamma \leq y^{(i)} \cdot \frac{(w^T x^{(i)} + b)}{\|w\|}$

$$\therefore \max_{w,b} \frac{\gamma}{\|w\|}, \quad \gamma \leq y^{(i)} (w^T x^{(i)} + b) \quad \text{[defined } \hat{\gamma} = \|w\| \cdot \gamma \text{]}$$

If $\hat{\gamma}^*, w^*, b^*$ is a soln then $k \hat{\gamma}^*, k w^*, k b^*$ is also a soln.

Thus we take $\hat{\gamma}^* = 1$ {WLOG}

Hence, our optimization becomes:

$$\max_{w,b} \frac{1}{\|w\|}, \quad 1 \leq y^{(i)} (w^T x^{(i)} + b)$$

Changing max to min:

$$\min_{w,b} \frac{1}{\|w\|}, \quad 1 \leq y^{(i)} (w^T x^{(i)} + b)$$

\downarrow
 $\|w\|^2$
 $w^T w$

∴ Final optimization:

$$\min_{w,b} \left[\frac{1}{2} w^T w \right], \quad y^{(i)} (w^T x^{(i)} + b) \geq 1$$

Convex Optimization Problem

Stephen Boyd {Book}

Convex Optimization Problems : (Constrained)

General form:

$$\min_w f(w) \text{ , subject to: } g_k(w) \leq 0, h_l(w) = 0$$

convex \Downarrow affine \Downarrow $\{w^T x + b\}$

$$k \in \{1, \dots, r\}$$

$$l \in \{1, \dots, s\}$$

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

If form is this then feasible set is convex set. C is convex if $w_1, w_2 \in C$ then $\alpha w_1 + (1-\alpha) w_2 \in C$ If f is convex then this is convex optimization.

SVM:

$$\text{SVM Objective: } \min_{w,b} \left[\frac{1}{2} w^T w \right], \quad y^{(i)} (w^T x^{(i)} + b) \geq 1$$

Is SVM optimization a convex constrained optimization problem? Yes

$$\hookrightarrow f(w) = \frac{1}{2} w^T w \text{ is convex}$$

We have no equality constraints

Also, $y^{(i)} (w^T x^{(i)} + b)$ is linear & hence convex

} \therefore This is a convex optimization problem.

- We will go from primal space to dual space. {where it is easier to solve}
- Think on why can't we use SGD to optimize this?

Theory of Lagrangians:

$$\underbrace{L(w, \alpha, \beta)}_{\text{Lagrangian}} = f(w) + \sum_{k=1}^m \alpha_k g_k(w) + \sum_{l=1}^n \beta_l h_l(w)$$

α, β = Lagrangian multipliers

$$\begin{aligned} \text{Optimization Problem: } & \min_w f(w) \\ & \forall k \quad g_k(w) \leq 0 \\ & \forall l \quad h_l(w) = 0 \end{aligned}$$

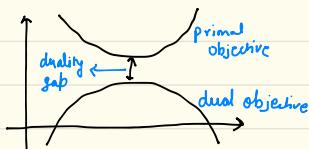
$$\text{Primal objective: } \max_{\substack{\alpha, \beta \\ \alpha \geq 0}} L(w, \alpha, \beta) = \Theta_p(w)$$

$$\text{form of } \Theta_p(w) = \begin{cases} f(w) & \text{if } \text{--- is feasible} \quad \{ h_l = 0 \text{ & } g_k \leq 0 \therefore \alpha_k = 0 \} \\ \infty & \text{otherwise} \quad \{ \text{make the corresponding } \alpha_k \text{ or } \beta_l \text{ large} \} \end{cases}$$

Now, $\min_w \Theta_p(w)$ is exactly same as our original optimization

$$\min_w \underbrace{\Theta_p(w)}_{\text{primal objective}} = \min_w \left[\max_{\alpha, \beta} (L(w, \alpha, \beta)) \right]$$

$$\max_{\substack{\alpha, \beta \\ \alpha \geq 0}} \underbrace{\Theta_D(\alpha, \beta)}_{\text{dual objective}} = \max_{\alpha, \beta} \left[\min_w L(w, \alpha, \beta) \right] \quad \text{dual problem}$$



$$\Theta_p(w) \geq \Theta_D(\alpha, \beta)$$

$$\therefore \min_w \Theta_p(w) \geq \max_{\alpha, \beta} \Theta_D(\alpha, \beta)$$

• Theorem: $\max_{\gamma} \min_{\omega} f(\omega, \gamma) \leq \min_{\omega} \max_{\gamma} f(\omega, \gamma)$

\hookrightarrow Let $g(\gamma) = \min_{\omega} f(\omega, \gamma)$

and γ^* be s.t. $\gamma^* = \arg \max_{\gamma} g(\gamma)$

$$\therefore g(\gamma^*) = \min_{\omega} f(\omega, \gamma^*)$$

$$\text{Now, } f(\omega, \gamma^*) \leq \max_{\gamma} f(\omega, \gamma) \quad \forall \omega$$

$$\begin{aligned} \text{by def of } g(\gamma) \quad \min_{\omega} f(\omega, \gamma^*) &\leq \min_{\omega} \max_{\gamma} f(\omega, \gamma) \\ \max_{\gamma} \min_{\omega} f(\omega, \gamma) &\leq \min_{\omega} \max_{\gamma} f(\omega, \gamma) \end{aligned}$$

Hence Proved

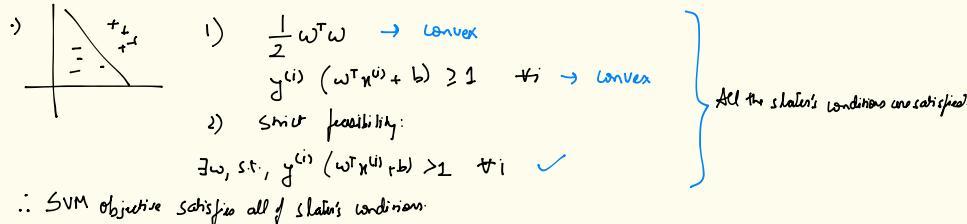
• Thus, $\max_{\substack{\alpha, \beta \\ \alpha \geq 0}} \Theta_D(\alpha, \beta) = \max_{\substack{\alpha, \beta \\ \alpha \geq 0}} \left[\min_{\omega} L(\omega, \alpha, \beta) \right] \leq \min_{\omega} \left[\max_{\substack{\alpha, \beta \\ \alpha \geq 0}} L(\omega, \alpha, \beta) \right] = \min_{\omega} \Theta_P(\omega)$

$$\therefore d^* \leq p^*$$

We will prove that in convex case: $d^* = p^*$

• Slater's conditions:

- ① ④ $L(\omega)$ is convex
- ⑤ $g_k(\omega)$'s are convex $\forall k$
- ⑥ $h_\lambda(\omega)$'s are affine $\forall \lambda$
- ② $\exists \omega$, $g_k(\omega)$'s are strictly feasible $\forall k$
 $\exists \omega$, s.t. $g_k(\omega) < 0 \quad \forall k$
- ⑦ $h_\lambda(\omega) = 0 \quad \forall \lambda$



∴ SVM objective satisfies all of slater's conditions.

Now, Slater's conditions satisfied \Rightarrow

- ① $P^* = D^*$ primal optimal equals dual optimal
- ② $\exists \omega^*, \alpha^*, \beta^*$

ω^* is primal optimal point
 α^*, β^* are dual optimal point

$$P^* = L(\omega^*, \alpha^*, \beta^*) = D^*$$

• KKT conditions (Karush Kuhn Tucker):

Consider a convex optimization problem

if $\exists \omega^*, \alpha^*, \beta^*$ s.t. $P^* = L(\omega^*, \alpha^*, \beta^*) = D^*$ \iff

- ① ④ $\nabla_{\omega} L(\omega, \alpha^*, \beta^*) \Big|_{\omega^*} = 0$ {gradient variables}
- ② $\nabla_{\beta} L(\omega^*, \alpha^*, \beta) \Big|_{\beta^*} = 0$ {dual variables}
- ③ $g_k(\omega^*) \leq 0 \quad \forall k$ } initial feasibility
- $h_\lambda(\omega^*) = 0 \quad \forall \lambda$
- $\delta \alpha_k^* \geq 0 \quad \forall k$ } dual feasibility
- ④ $\alpha_k^* g_k(\omega^*) = 0 \quad \forall k$ \Rightarrow $\begin{cases} g_k(\omega) < 0 \\ \alpha_k = 0 \end{cases}$ } complementary slackness

SVM Objective: $\min_{\omega, b} \frac{1}{2} \omega^T \omega , y^{(i)} (\omega^T x^{(i)} + b) \geq 1$

$$\therefore L(\omega, b, \alpha) = \frac{1}{2} \omega^T \omega + \sum_{i=1}^m \alpha_i (1 - y^{(i)} (\omega^T x^{(i)} + b))$$

$$\Theta_P(\omega, b) = \min_{\omega, b} \left[\max_{\alpha \geq 0} L(\omega, b, \alpha) \right]$$

$$\Theta_D(\alpha) = \max_{\omega, b} \left[\min_{\alpha \geq 0} L(\omega, b, \alpha) \right]$$

$$\text{Now, } \nabla_{\omega} L(\omega, b, \alpha) = \omega - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\nabla_{\omega} L(\omega, b, \alpha) \Big|_{\omega^*} = 0 \Rightarrow \omega^* = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\text{Similarly, } \nabla_b L(\omega, b, \alpha) = -\sum_{i=1}^m \alpha_i y^{(i)}$$

$$\therefore \nabla_b L(\omega, b, \alpha) \Big|_{b^*} = 0 \Rightarrow \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

$$\text{At, } \omega^*, b^* \Rightarrow L(\alpha, \omega^*, b^*) = \frac{1}{2} \omega^{*T} \omega + \sum_{i=1}^m \alpha_i (1 - y^{(i)} (\omega^{*T} x^{(i)} + b))$$

$$\omega^* = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad \& \quad \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

Substitution:

$$L(\alpha, \omega^*, b^*) = \frac{1}{2} \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T \left(\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right) + \sum_{i=1}^m \alpha_i - \left[\sum_{i=1}^m \alpha_i y^{(i)} \left(\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right)^T x^{(i)} \right. \\ \left. + b \sum_{i=1}^m \alpha_i y^{(i)} \right]$$

$$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)} x^{(j)}$$

Constraints: ① $\alpha_i \geq 0 \forall i$

$$\textcircled{2} \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

Dual problem: $\max_{\alpha \geq 0} L(\alpha, \omega^*, b^*)$

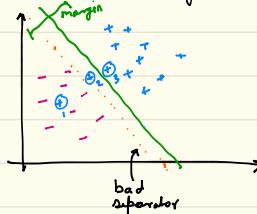
Primal Objective: $\min_{w,b} \frac{1}{2} w^T w, \quad y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad \forall i$

\downarrow Dual problem

Dual: $\max_{\alpha} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \cdot \alpha_j y^{(i)} y^{(j)} x^{(i)^T} x^{(j)} \right), \quad \alpha_i \geq 0 \quad \forall i, \quad \sum_{i=1}^m \alpha_i y^{(i)} = 0, \quad w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$

SMO: Sequential Minimal Optimization, LIBSVM

Handling non-linearly separable data:



$$\min_{w,b} \frac{1}{2} w^T w, \quad y^{(i)}(w^T x^{(i)} + b) \geq 1$$

How to handle cases like ①, ② & ③?

$$\hookrightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0$$

But how to control ϵ_i 's so that we get some sensible soln? [else $w=0$ is a soln]

$$\hookrightarrow \min_{w,b,\epsilon} \frac{1}{2} w^T w + C \sum_{i=1}^m \epsilon_i$$

sum loss
controls the effect of the penalty

Note: $C \rightarrow \infty$ in linearly separable case.

In this case too Slater's conditions are satisfied. Primal soln = Dual soln.

$$L(w,b,\epsilon,\gamma) = \frac{1}{2} w^T w + C \sum \epsilon_i + \sum_{i=1}^m \alpha_i [1 - \epsilon_i - y^{(i)}(w^T x^{(i)} + b)] + \sum_{i=1}^m \gamma_i (-\epsilon_i)$$

Dual objective: $\min_{w,b,\epsilon} L(w,b,\epsilon,\gamma)$

$$\nabla_w L = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\nabla_b L = 0 \Rightarrow \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

$$\nabla_\epsilon L = 0 \Rightarrow C - \alpha_i = \gamma_i$$

$$\max_{\alpha} \left[\sum \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)^T} x^{(j)} \right], \quad \alpha_i \geq 0, \quad \sum_{i=1}^m \alpha_i y^{(i)} = 0, \quad 0 \leq \alpha_i \leq C$$

$x \rightarrow \phi(x) \Rightarrow \text{Kernel}$

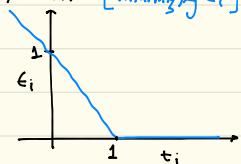
• Solving Primal Problem:

$$\min_{w,b} \frac{1}{2} w^T w + c \sum \epsilon_i, \quad y^{(i)} (w^T x^{(i)} + b) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0$$

$$\Leftrightarrow \epsilon_i \geq 1 - y^{(i)} (w^T x^{(i)} + b)$$

$$\therefore \min_{w,b,t} \frac{1}{2} w^T w + c \sum \epsilon_i, \quad \begin{array}{l} \epsilon_i \geq 1 - t_i, \\ t_i = y^{(i)} (w^T x^{(i)} + b) \end{array} \rightarrow \epsilon_i \geq \max(0, 1 - t_i)$$

$$t_i = \max(0, 1 - t_i) \quad \text{[as we are minimizing } \epsilon_i \text{]}$$



∴ We can transform to unconstraint problem:-

$$\min_{w,b} \frac{1}{2} w^T w + c \sum_{i=1}^m \max(0, 1 - t_i)$$

where $t_i = y^{(i)} (w^T x^{(i)} + b)$

But we can't use gradient descent because this is not differentiable.

Some possible soln:-

$$\epsilon_i = \begin{cases} 0 & t_i > 1 \\ 1 - t_i & t_i \leq 1 \end{cases} \quad \nabla_{t_i} L = \begin{cases} 0 & t_i > 1 \\ -1 & t_i \leq 1 \end{cases}$$

we are using left gradient.

This is called Subgradient.

$$\begin{aligned} \nabla_w L &= w + c \sum_{i=1}^m \nabla_w [\max(0, 1 - t_i)] \\ &= w + c \sum_{i=1}^m \nabla_t [\max(0, 1 - t_i)] \cdot \nabla_w t_i \\ &= w + c \sum_{i=1}^m \left\{ \begin{array}{ll} 0 & t_i > 1 \\ -1 & t_i \leq 1 \end{array} \right\} * (y^{(i)} x^{(i)}) \end{aligned}$$

SVM:

$$\min_{w,b} \frac{1}{2} w^T w + C \sum_{i=1}^m \epsilon_i$$

Optimization methods

→ Primal space → Pegasos

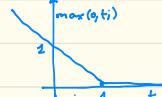
→ Dual space → SMO

$$\min_{w,b} \frac{1}{2} w^T w + \sum_{i=1}^m \max(0, 1 - t_i)$$

regularization

where $t_i = y^{(i)}(w^T x^{(i)} + b)$

Hinge loss



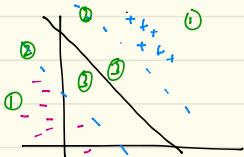
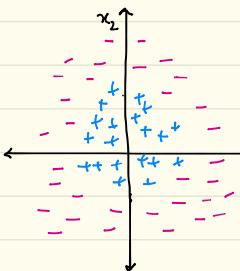
Note: the gradient is not zero at minima!

Dual Formulation:

$$\frac{1}{2} w^T w + C \sum_{i=1}^m \epsilon_i ; y^{(i)}(w^T x^{(i)} + b) \geq 1 - \epsilon_i ; \epsilon_i \geq 0$$

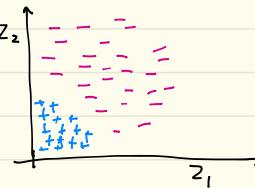
↓

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} [x^{(i)}]^T x^{(j)} ; 0 \leq \alpha_i \leq C ; \sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad \text{solve for } \alpha_i$$

Analysis of ϵ_i , α_i , y_i & support vectorsIn region ① $\Rightarrow \alpha_i = 0, \gamma_i = C$ ③ $\Rightarrow \alpha_i = C, \gamma_i = 0, \epsilon_i > 0$ ② $\Rightarrow 0 < \alpha_i < C, \gamma_i = 0, \epsilon_i = 0$ Kernels:

How to classify here? Can we use a linear separator?

$$Z_1 = x_1^2 \\ Z_2 = x_2^2$$



ϕ is kernel $(x_1, x_2) \xrightarrow{\phi} (x_1^2, x_2^2)$

$$(x_1, x_2) \xrightarrow{\phi} \begin{matrix} x_1^2, x_2^2, x_1 x_2 \\ x_1, x_2 \end{matrix} \in \mathbb{R}^4$$

$\phi(x) \in \mathbb{R}^N$

Kernels:

$$(x_1, x_2) \xrightarrow{\phi} \underbrace{\phi(x_1, x_2)}_{\{x_1^1, x_1^2, x_1^3, x_1^4, x_1, x_2\}}$$

$$(x_1, x_2, \dots, x_n) \xrightarrow[\substack{\text{polynomial of} \\ \text{size } d}]{} \text{terms of degree } \leq d \quad \# = {}^{n+1}C_d$$

$$\in \mathbb{R}^N, N = O(n^d)$$

If we could compute: $K(x, z) = \phi(x)^T \phi(z)$ efficiently then we are ok. $\Theta \{ \text{Kernel Trick}$

E.g.:

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

Goal: $K(x, z)$ represents $\phi(x)^T \phi(z)$ for some $\phi(x)$ (feature transformation)

$$\begin{aligned} & \Rightarrow \left(\sum_{i=1}^n x_i z_i \right)^2 + 2c \sum x_i z_i + c^2 \\ & = \left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right) + \sum_{i=1}^n (\sqrt{2c} x_i) (\sqrt{2c} z_i) + c^2 \\ & = \sum_{j=1}^n \sum_{i=1}^n x_i x_j z_i z_j + \sum_{i=1}^n (\sqrt{2c} x_i) (\sqrt{2c} z_i) + c^2 \end{aligned}$$

$$\phi(x) = \begin{bmatrix} x_1^1, x_1^2 \\ x_1^2, x_1^3 \\ x_1^3, x_1^4 \\ \vdots \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_1 \\ c \end{bmatrix} \quad \phi(z) = \begin{bmatrix} z_1^1, z_1^2 \\ z_1^2, z_1^3 \\ z_1^3, z_1^4 \\ \vdots \\ \sqrt{2c} z_1 \\ \sqrt{2c} z_1 \\ c \end{bmatrix}$$

Similarly: $K(x, z) = \underbrace{(x^T z + c)^d}_{O(n + \log d)}$, then $\phi(x)$ = degree d poly. on x $\Rightarrow O(n^d)$

$$\text{SVM Objective: } \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{[x^{(i)}]^T x^{(j)}]}_{\phi(x^{(i)})^T \phi(x^{(j)})}; \quad 0 \leq \alpha_i \leq C$$

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad \{ \text{How to calculate } w \text{ if } \phi(x) \text{ is exp?} \}$$

\hookrightarrow We don't! we only need $w^T x$ at classification which is computable using kernel trick.

$$f_{w_b}(x) = \begin{cases} 1 & [\text{if } w^T \phi(x) \geq 0] \rightarrow (\sum_{i=1}^m \alpha_i y^{(i)} \underbrace{\phi(x^{(i)})^T}_{\phi(x)}) \\ 0 & \text{otherwise} \end{cases}$$

Similarly for b $\longleftarrow \{ \text{Check!} \}$

Kernel Function :

$$\text{u}, \text{z} \in \mathbb{R}^n, \quad K(\text{u}, \text{z}) : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$$

We say that $K(\text{u}, \text{z})$ is a valid kernel function iff \exists a feature transformation ϕ s.t. $K(\text{u}, \text{z}) = \phi(\text{u})^\top \phi(\text{z})$

Another example of a Kernel function:

$$K(\text{u}, \text{z}) = e^{-\frac{\|\text{u}-\text{z}\|^2}{2\sigma^2}} \Rightarrow \text{Gaussian / RBF Kernel}$$

↳ Similarity b/w two points.

(Q) How to prove that $K(\text{u}, \text{z})$ is a valid kernel?

→ Find $\phi(\text{u})$ s.t. $K(\text{u}, \text{z}) = \phi(\text{u})^\top \phi(\text{z})$

→ Mercer's Theorem

Merger's Theorem :

$K(\text{u}, \text{z}) : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ on a finite dimensional space then for K to be valid it is necessary & sufficient that:

Let $\{\text{x}^{(1)}, \dots, \text{x}^{(m)}\}$ be a finite set of points in the original space of kernel matrix:

$$K^m = \begin{bmatrix} & & & \\ & \ddots & & \\ & & \vdots & \\ i & \cdots & K_{ij}^m & = \\ & & & K(\text{x}^{(i)}, \text{x}^{(j)}) \end{bmatrix} \quad \text{then:}$$

1) K^m is symmetric

2) K^m is positive semi-definite

Proof that this is necessary:

$$\begin{aligned} \text{① } K^m \text{ is symmetric: } K_{ij}^m &= \phi(\text{x}^{(i)})^\top \phi(\text{x}^{(j)}) \\ &= \phi(\text{x}^{(j)})^\top \phi(\text{x}^{(i)}) \\ &= K_{ji}^m \end{aligned}$$

② K^m is true semi-definite:

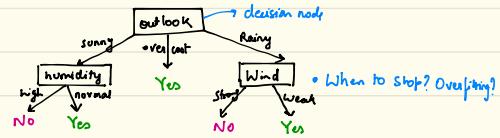
if $\text{u}^\top K^m \text{u} \geq 0 \forall \text{u}$ then true

$$\begin{aligned} \text{u}^\top K^m \text{u} &= \sum_{i,j=1}^m u_i K_{ij}^m u_j \\ &= \sum_{i,j=1}^m u_i (\sum_k \phi(\text{x}^{(i)})_k \phi(\text{x}^{(j)})_k) u_j \\ &= \sum_{i,j=1}^m \sum_k u_i \phi(\text{x}^{(i)})_k \phi(\text{x}^{(j)})_k u_j \\ &= \sum_k \left(\sum_i u_i \phi(\text{x}^{(i)})_k \right) \left(\sum_j u_j \phi(\text{x}^{(j)})_k \right) \\ &= \sum_k \left[\sum_i u_i \phi(\text{x}^{(i)})_k \right]^2 \geq 0 ; \quad \text{Hence Proved} \end{aligned}$$

Decision Trees:

Example: Learning to play Tennis:

x_1, x_2, x_3, y Outlook ∈ {sunny, overcast, rainy}
 outlook humidity wind ∈ {high, normal}



* When to Stop? Overfitting!

Example Tree

D Internal Nodes → Decision Nodes

2) Leaves → Labels

3) Test single attribute at each node {Training difficult with multiple attributes even though they may have more advantage due to look ahead}

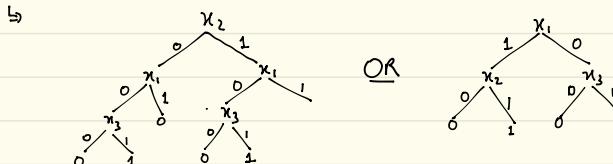
Note: Multi attribute has same representation power as single attribute decision node.

Claim: A decision tree can represent any Boolean function over n variables

↳ One way is to expand the tree fully!

But we can also do better {shortest tree}

Eg. $(x_1 \wedge x_2) \vee (\neg x_1 \wedge x_3)$, built shortest possible tree



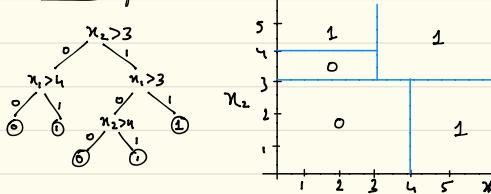
Not shortest

• How to handle continuous data?

→ Binning
→ Learning the Splitter

• Variations: testing same attribute again down the path {in case of continuous / multivalued (binary tree)}

• Decision surface:



For oblique we need to test two attributes at a time.

Func growTree(D) {

if ($y^{(i)} = 0$) ++i;

return Leaf(D, 0);

if ($y^{(i)} = 1$) ++i;

return Leaf(D, 1);

X_j = choose Best Attr(D);

$D_0 = \{x^{(i)}, y^{(i)}\}_{i=1}^m$ s.t. $y^{(i)} = 0$

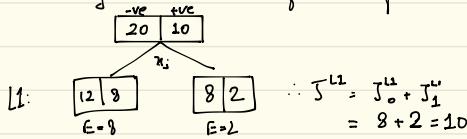
$D_1 = \{x^{(i)}, y^{(i)}\}_{i=1}^m$ s.t. $y^{(i)} = 1$

return new Node(D, X_j , growTree(D_0), growTree(D_1));

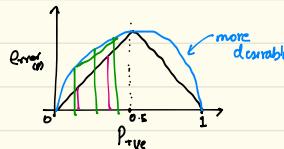
}

Training error:

Let us say their error is # of misclassified labels.



But this error does not ensure that data labels are segregated!



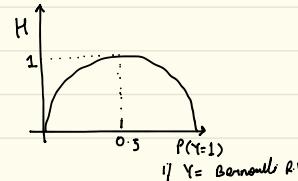
Entropy:

Y: discrete valued R.V.

$$H(Y) = \sum_y P(Y=y) [-\log_2(P(Y=y))]$$

entropy (Expected surprise)

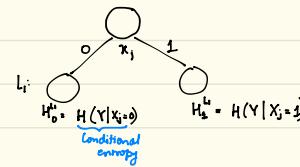
surprise



Now, let us use entropy as our error metric.

* Split on the feature which maximizes $H(Y) - H(Y|X_j)$

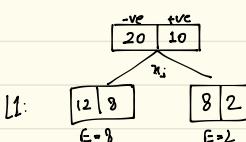
$$H(Y|X_j) = \sum_{x_j} P(X_j=x_j) \cdot H(Y|X_j=x_j)$$



$$I(Y, X_j) = H(Y) - H(Y|X_j)$$

Information Gain

$I(X_j, Y)$



$$H(Y) = \frac{2}{3} \left[-\log \left(\frac{2}{3} \right) \right] + \frac{1}{3} \left(\log \left(\frac{1}{3} \right) \right) = 0.9183$$

$$H(Y|X_j=0) = \frac{3}{5} \left(-\log \frac{3}{5} \right) + \frac{2}{5} \left(\log \frac{2}{5} \right) = 0.9710$$

$$H(Y|X_j=1) = \frac{4}{5} \left(-\log \frac{4}{5} \right) + \frac{1}{5} \left(\log \frac{1}{5} \right) = 0.71$$

$$H(Y|X_j) = \frac{2}{3} \cdot H(Y|X_j=0) + \frac{1}{3} \cdot H(Y|X_j=1) = 0.8978$$

- Slides Today!
- Pathological case for decision trees \rightarrow Parity with noise.
- Definition of overfitting

Random Forests:

Decision trees are unstable learners \rightarrow High variance

Based on small changes in data, the model can drastically change

Soln: Bagging



Bagging:

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

function generate Bootstrap Samples (D_k) {

For $k=1$ to r {

$$D_k = \{\}$$

For $i=1$ to m {

$$c \sim \text{Uniformly}(\{1 \dots m\})$$

$$D_k \leftarrow \{x^{(c)}, y^{(c)}\} \cup D_k$$

}

}

Sampling with replacement

For prediction $\begin{cases} \text{Classification: Voting} \\ \text{Regression : Average} \end{cases}$

Random Forests:

Bag of Decision Trees

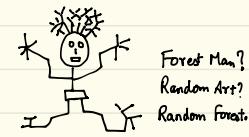
1) # estimators (# of trees)

2) Max-features \rightarrow Select best feature from a random sample of all features (of size maxfeat)

3) Boot strap

4) Out of bag error. \exists H.W.

Trade Off: Accuracy of individual trees vs Low variance of models
(bias-variance) high ns vs low ns

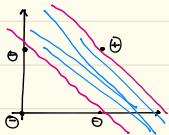


Forest Man?
Random Art?
Random forest.

Basic unit: Perceptron

\therefore This perceptron learns linear boundary $[\theta^T x \geq 0]$.

$$\text{Ex: } f(x_1, x_2) = x_1 \wedge x_2 \quad x_1, x_2 \in \{0, 1\}$$



$$h_\theta(x) = I\{\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0\}$$

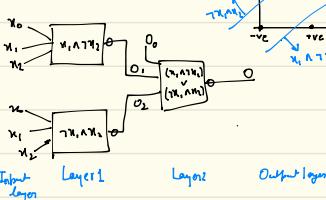
Optimal parameter (if slope is constrained to -1): $\begin{cases} \theta_0 = 1 \\ \theta_1 = 1 \\ \theta_2 = -1 \end{cases}$

• if $f(x_1, x_2) = x_1 \vee x_2 \Rightarrow$ if $\theta_2 = 1, \theta_1 = 1$ then $\theta_0 \in (-1, 0)$

• if $f(x) = \gamma x \Rightarrow$ if $\theta_1 = -1$ then $\theta_0 \in (0, 1)$

• if $f(x) = K \Rightarrow$ if $\theta_1 = 1$ then $\theta_0 \in (-1, 0)$

• if $f(x_1, x_2) = x_1 \oplus x_2 \Rightarrow$ Cannot Learn this!!! But $x_1 \oplus x_2 = (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$ of this can be represented by composition of perceptrons.



• How do we learn a perceptron?

$$\text{Input: } \begin{array}{c} x_0 \\ x_1 \\ x_2 \end{array}, \quad J(\theta) = \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)}))^2 \quad \left[\begin{array}{l} \text{Use Gradient descent. But problem:-} \\ \Rightarrow h_\theta(x) \text{ is non-differentiable! [Also subgradient without wort Why?]} \\ \Rightarrow \text{Not convex, but we don't care} \oplus \end{array} \right]$$

Delta Rule (ad hoc heuristic algo for training):

init (θ);

do {

$x^{(i)} \leftarrow$ wrongly classified example

for all $j \{$

$$\Delta \theta_j = \eta [y^{(i)} - h_\theta(x^{(i)})] x_j^{(i)}$$

$$\theta_j = \theta_j + \Delta \theta_j;$$

} (until convergence)

Note: This actually converges if our data is linearly separable

for a sufficiently small η .

Let us replace our step function to sigmoid function. (to make loss differentiable)

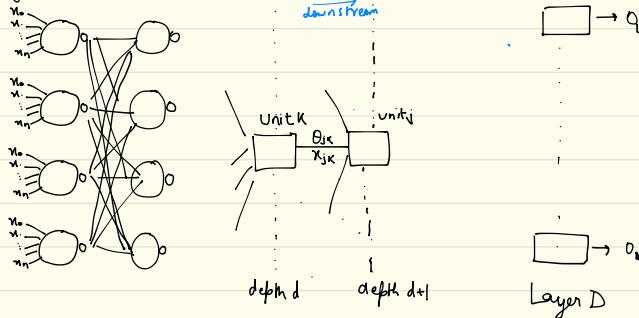
$$\text{Input: } \begin{array}{c} x_0 \\ x_1 \\ x_2 \end{array}, \quad f_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

Training of a perceptron:

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)}))^2 \\ \frac{\partial J(\theta)}{\partial \theta_j} &= -\frac{1}{2} \sum_{i=1}^m 2(y^{(i)} - h_\theta(x^{(i)})) \cdot \frac{\partial (h_\theta(x^{(i)}))}{\partial \theta_j} \\ &= -\sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] (h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))) x_j^{(i)} \end{aligned}$$

[Check sign] [not convex in general]
extra term as compared to logistic

Training a general network:



Note: hidden layer one layers
hidden from the output

$$x_{jk} = \text{output from unit } k \text{ to unit } j$$

$$\theta_{jk} = \text{parameter associated with input } x_{jk}$$

$$\text{net}_j = \text{net linear input associated with unit } j = \sum_{k \in \text{upnbr}(j)} \theta_{jk} x_{jk}$$

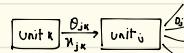
$$o_j = \text{output of unit } j = g(\text{net}_j) = \frac{1}{1 + e^{-\text{net}_j}}$$

↳ upneighbours of unit j

$$J(\theta) = \sum_{i=1}^m \sum_{j=1}^n [y_j^{(i)} - o_j^{(i)}]^2$$

↳ indirectly dependent
on param. in hidden layer

$$\textcircled{1} \quad \frac{\partial J(\theta)}{\partial \theta_{jk}} = \frac{\partial J(\theta)}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial \theta_{jk}} = \frac{\partial J(\theta)}{\partial \text{net}_j} x_{jk}$$



$$\frac{\partial J(\theta)}{\partial \theta_{jk}} = \frac{\partial J(\theta)}{\partial \text{net}_j} x_{jk} = -\delta_j x_{jk}$$

↳ we are defining this

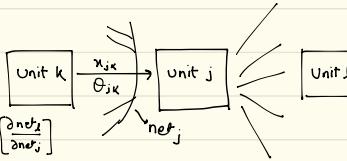
$$\textcircled{2} \quad \frac{\partial J(\theta)}{\partial \text{net}_j} = o_j(1 - o_j), \quad \forall k \in \text{upnbr}(j) \quad x_{jk} = o_k$$

Note: We will assume: $m = 1$

(A) Unit $j \in$ output layer

$$\frac{\partial J(\sigma)}{\partial \text{net}_j} = \frac{\partial}{\partial \text{net}_j} \left[\frac{1}{2} \sum_{c=1}^C (y_c - o_j)^2 \right] = -1 \cdot (y_j - o_j) o_j (1-o_j)$$

(B) unit $j \in$ hidden layer



$$\frac{\partial J(\sigma)}{\partial \text{net}_j} = \sum_{l \in \text{domain}(j)} \left[\frac{\partial J(\sigma)}{\partial \text{net}_l} \right] \left[\frac{\partial \text{net}_l}{\partial \text{net}_j} \right]$$

$$= \sum_{l \in \text{domain}(j)} -d_l \cdot \frac{\partial \text{net}_l}{\partial \text{net}_j} o_j = \sum_{l \in \text{domain}(j)} -d_l o_j O_j (1-O_j)$$

$$\therefore d_j = \sum_{l \in \text{domain}(j)} -d_l O_j O_j (1-O_j)$$

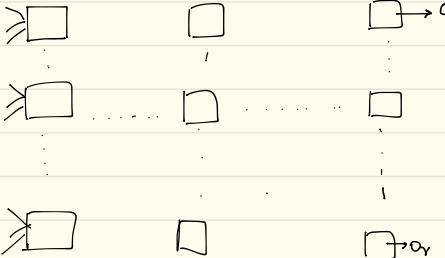
Note:

$f(y_1, \dots, y_n)$, where each y_c is a fn of x

$$\frac{\partial f(y_1, \dots, y_n)}{\partial x} = \sum_{c=1}^n \frac{\partial f(y_1, \dots, y_n)}{\partial y_c} \frac{\partial y_c}{\partial x}$$

Lecture 29 - 23.03.18

- Layer by layer training:



Training:

- ① Forward Propagation
- ② Back prop
- ③ Weight update

- Slides

- Auto encoder

Deep Learning:

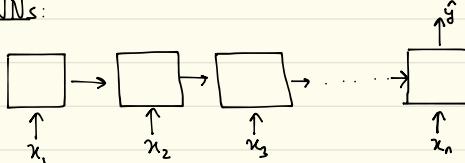
"Deep" neural network

Architectures:

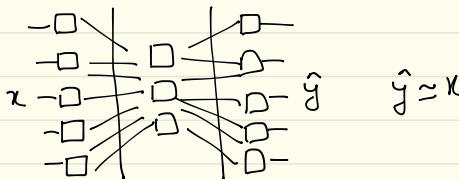
1) CNNs

Lecture 31 - 05.04.18

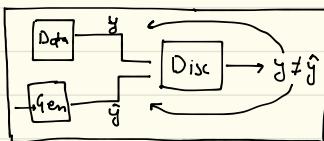
• RNNs:



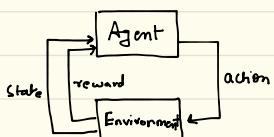
• Auto encoders:



• GANs: (Generative Adversarial Networks)

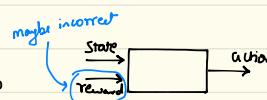


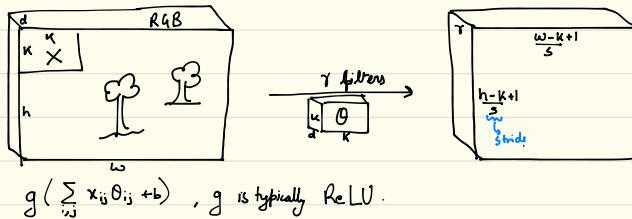
• Deep Q-learning: { Reinforcement learning }



Example: Alpha Go

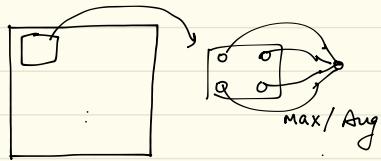
Prof: Balramenon Ravindran {IITM}



• CNNs:

We typically pad our images to make input & output dimensions same.

Lecture 32 - 06.04.18

• Pooling:

\therefore for 2×2 pooling $\Rightarrow 500 \times 500 \xrightarrow{\text{pool}} 250 \times 250$

We can use either avg or max pooling

• Slides

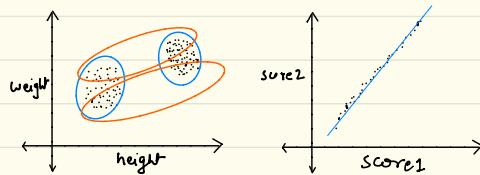
• Unsupervised Learning :

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

1) Clustering

2) Density estimation: $p(w) = \sum_z p(x|z) p(z)$

3) Dimensionality Reduction : PCA {Principal Component Analysis}



• K-Means :

Good clustering :

- Points in the same cluster are close by
- Points in different clusters are far apart

Let us assume that we know the number of clusters = k

Parameters:

- $\mu = (\mu_1, \dots, \mu_k)$, $i \in \{1 \dots k\}$
- $c = (c^{(1)}, c^{(2)}, \dots, c^{(m)})$, $i \in \{1 \dots m\}$

Objective:

$$\min_{\mu, c} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

• Given cluster assignment we can find μ

• Given μ we can find cluster assignments.

Algorithm:

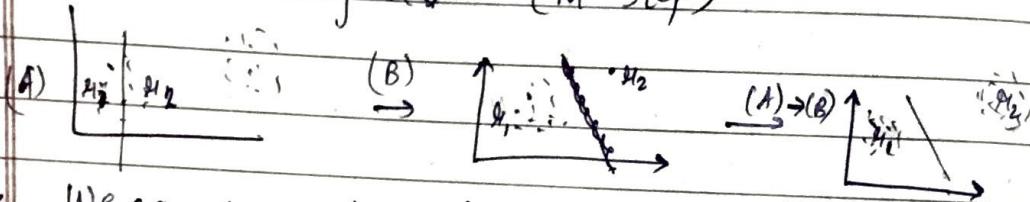
- Randomly initialize μ
- While not converged:
 - Given μ find C {E-step}
 - Given C find μ {M-step}

Points in each cluster should be as close to their representatives as possible. This automatically implies that representatives are far apart.

* If we know the assignments, we can compute the representatives. We can do this by minimising sum of distances of representatives from mean of each clusters. Given

K-Means Algorithm:-

- 1) Given μ_i 's compute cluster assignments.
 - 2) Given c 's compute cluster representatives.
~~This is any~~ a fixed number. Our algorithm should first compute c and then find out the clustering.
- (A) Given μ_i 's find c 's (E-step)
 → (B) Given c 's, find μ_i 's (M-step)



* We can show that K-Means finds local minima of $\underset{c, \mu}{\operatorname{argmin}} J(c, \mu)$

12/04/2018

→ InitParams(); $\{ \mu_1, \dots, \mu_c \} \rightarrow \text{parameters?}$

Repeat [(A) → AssignPoint(); → E step
 until convergence (B) → EstimateParams() → M steps]

$$\star J(c, \mu) = \underset{i=1}{\operatorname{argmin}} \sum_{i=1}^m \|x^{(i)} - \mu^{(i)}\|^2 \quad c = (c^{(1)}, \dots, c^{(m)}) \\ \mu = (\mu^{(1)}, \dots, \mu^{(m)})$$

Here C is discrete, so we cannot apply gradient descent.
 $\mu_k \in \mathbb{R}^n$

Coordinate Descent :-

$\underset{x_1, x_2}{\operatorname{argmin}} f(x_1, x_2) \rightarrow \text{fix everything except one direction}$

Step 1: $\min_{x_2} f(x_1, x_2) \rightarrow$ stops at tangent to contour
 \downarrow
 fixed

Step 2: then fix x_2 and move along x_1 .

$$\min_{x_1} f(x_1, x_2)$$

go to step 1.

\rightarrow gradient descent

\rightarrow SGD

\rightarrow coordinate descent



* Convergence criteria - when there is no direction to move (no improvement in any direction). This has similar guarantees as gradient descent. When to use it?
Single variable optimization can be easier (analytical solution available). Apply coordinate descent where gradient descent is difficult to compute while it may be possible to solve entire coordinate descent analytically.

* Block Coordinate Descent:- fix some coordinates and descend on others. (some maybe discrete while others continuous)

* Q: For SVM, we use block coordinate descent. Why can't we use coordinate descent in SMO?

Because of similar guarantees, block coordinate descent also converges to local optima.

$$J(C, \alpha) = \underset{C, \alpha}{\operatorname{argmin}} \sum_{i=1}^m [\|\alpha^{(i)} - \mu_{c(i)}\|^2]$$

① Fix $\alpha^{(i)}$'s

$$\underset{\alpha^{(1)}, \dots, \alpha^{(m)}}{\operatorname{argmin}} \sum_{i=1}^m \|\alpha^{(i)} - \mu_{c(i)}\|^2$$

$$= \underset{\alpha^{(i)}}{\operatorname{argmin}} \sum_{c} \|\alpha^{(i)} - \mu_{c(i)}\|^2 \rightarrow E \text{ step}$$

② Fix $c^{(i)}$'s

$$\underset{\mu_1, \dots, \mu_K}{\operatorname{argmin}} \sum_{i=1}^m \|\alpha^{(i)} - \mu_{c(i)}\|^2 = \underset{\mu_1, \dots, \mu_K}{\operatorname{argmin}} \sum_{i=1}^m \sum_{k=1}^K I_{\{c(i)=k\}} \|\alpha^{(i)} - \mu_k\|^2$$

$$= \underset{\mu_1, \dots, \mu_K}{\operatorname{argmin}} \sum_{i=1}^m I_{\{c(i)=k\}} \|\alpha^{(i)} - \mu_k\|^2$$

(K-Medians chooses one of the points as representatives)

* $\mu_k \in \mathbb{R}^n$

$$\Rightarrow \mu_k = \frac{1}{m} \sum_{i=1}^m$$

We can apply SGD or CGD.

$$I\{x^{(i)} = k\} x^{(i)}$$

$$\frac{\partial}{\partial \mu_k} I\{x^{(i)} = k\} x^{(i)}$$

∇R

$$\sum_{i=1}^m I\{x^{(i)} = k\} x^{(i)}$$

$\hookrightarrow \text{Mean}$

* K-Means: ① Start with some initial cluster means.

② Find c 's given μ 's.

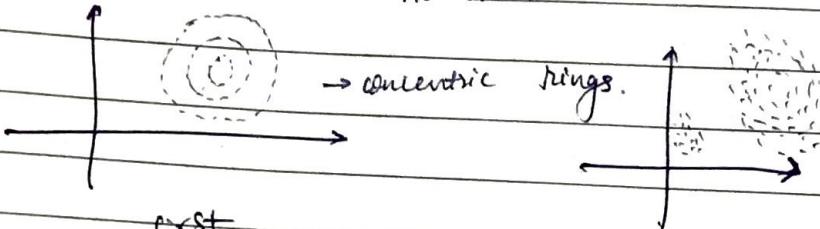
③ Find μ 's given c 's.

repeat
until
converges

→ Convergence: No change in c or μ in consecutive steps.

This converges to the local optima.

Q: When does K-Means not work?



* In K-Means first

pick a centroid among one of the points.

* Clusters must be of similar size and well separated.

K-Means assumes "K". Some algorithms figure out K and then proceed.

13/04/2018

- Some other clustering algorithm apart from K-Means:
 - DBSCAN
 - Chinese Restaurant Process

Gaussian Mixture Models (GMM):

$Z^{(i)} \sim \text{Multinomial}(\phi)$

$x^{(i)} | Z^{(i)} \sim \text{Normal}(\mu_k, \Sigma_k)$ (given $Z^{(i)} = k$)

$$\Theta_{\text{ML}} = \arg \max_{\theta} \prod_{i=1}^m P(x^{(i)}; \theta)$$

$$= \arg \max_{\theta} \prod_{i=1}^m \left(\sum_{j=1}^k P(x^{(i)}, j; \theta) \right)$$

$$P(x^{(i)}, j; \theta) \cdot P(j; \theta)$$

} Since there is a Σ here we can't have closed form solution. Also gradient descent is bad because of Σ inside a log. Thus we use iterative style optimization (EM)

Algorithm:

① Init Params () ~ $(\phi, \{\mu_k, \Sigma_k\}_{k=1}^K)$

② A) Estimate $Z^{(i)}$'s given the parameters.

$$\hookrightarrow P(Z^{(i)} | x^{(i)}; \theta) = \frac{P(x^{(i)} | Z^{(i)}; \theta) \cdot P(Z^{(i)}; \theta)}{\sum_{Z^{(i)}} P(x^{(i)} | Z^{(i)}; \theta) \cdot P(Z^{(i)}; \theta)}$$

} E step

repeat

③ B) Update the parameters:

$$\hookrightarrow \phi_k = \frac{\sum_{i=1}^m 1\{Z^{(i)}=k\}}{m}$$

$$\mu_k = \frac{\sum_{i=1}^m 1\{Z^{(i)}=k\} x^{(i)}}{\sum_{i=1}^m 1\{Z^{(i)}=k\}}$$

$P(Z^{(i)}=k) \Rightarrow \text{Soft EM}$

Since otherwise we throw away info & impose a hard bound.

$$\Sigma_k = \frac{\sum_{i=1}^m 1\{Z^{(i)}=k\} (x^{(i)} - \mu_k)(x^{(i)} - \mu_k)^T}{\sum_{i=1}^m 1\{Z^{(i)}=k\}}$$

} M step

④ ϕ : uniform

⑤ Σ : identity

⑥ Do hard assignment

• K-Means vs GMM:

K-Means

① Init Params (); // set of centroids

Hard ② E-step: Assign cluster to every point

③ M-step: Compute (μ_1, \dots, μ_r) given $(c^{(1)}, \dots, c^{(m)})$

GMM

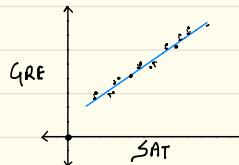
① Init params () // set $\phi, \{\mu_k, \Sigma_k\}$

② E-step: Compute $P(Z^{(i)} | x^{(i)}; \theta)$ Soft

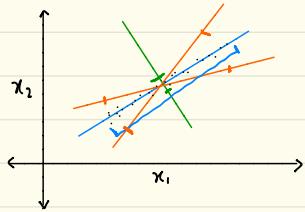
③ M-step: Compute $\phi, \{\mu_k, \Sigma_k\}$ given $P(Z^{(i)} | x^{(i)}; \theta)$

Principal Component Analysis:

- Dimensionality Reduction
- Difference with Autoencoders?
- Eigenface: $n \times n$ image to ~ 50 dimensional vector.



- Advantages:
- Discern Pattern
 - Reduce computation
 - Reducing noise



We try to maximize variability while projecting.

- Given: $\{x^{(i)}\}_{i=1}^m$, $x^{(i)} \in \mathbb{R}^n$
 - To find: $\{u_1, \dots, u_k\}$ where $u_i \in \mathbb{R}^n$, $k \ll n$
- $U_i \perp U_j \quad \forall i, j \quad \& \quad \|u_i\|^2 = 1$
- $$x^{(i)} \approx \sum_{j=1}^k \lambda_{ij} u_j$$

- Normalizing Data:

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \quad \forall j: x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

We may not want to preprocess if all attributes are from same dimension. {E.g. not in image pixels}

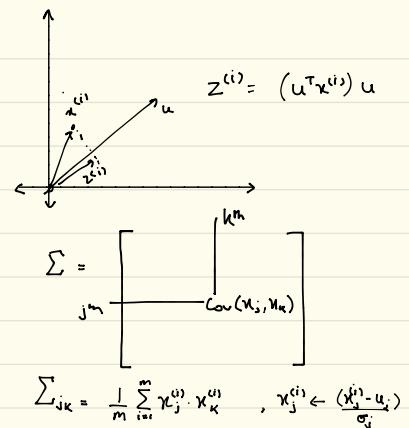
- We can either try to maximize variance or minimize dist. projected

$$\frac{1}{m} \sum (z^{(i)})^2 = \frac{1}{m} \sum (\chi^{(i)\top} u)^2 = \frac{1}{m} \sum (\chi^{(i)\top} u) (\chi^{(i)\top} u)$$

$$= u^\top \left[\frac{1}{m} \sum_{i=1}^m \chi^{(i)} \chi^{(i)\top} \right] u$$

$\underbrace{\Sigma}_{\text{empirical cov matrix}}$

PCA (single component) $\Rightarrow \arg \max_u [u^\top \Sigma u]$



Find $\{u_1, \dots, u_m\}$ s.t. $u_1 \perp u_2, \|u_k\|^2 = 1$

\hookrightarrow Projection will be: $z_1^{(i)}, \dots, z_K^{(i)}$

$$\therefore (u_1, \dots, u_K) = \arg \max_{u_1, \dots, u_K} \sum_{k=1}^K [u_k^\top \Sigma u_k], \text{ where } u_1 \perp u_2 \text{ and } \|u_k\|^2 = 1$$

Solved using SVD. Eigenvalues based computation

Lecture 37 - 19.04.18

Class skipped due to Imagine Cup finals.

19/04/2018

Exam ★ Maximising variance is equivalent to minimising sum of squared distances from projected line.

$u_1, \dots, u_K \rightarrow$ orthonormal basis

$$\underset{u_1, \dots, u_K}{\operatorname{argmax}} \sum_{l=1}^K (u_l^T \Sigma u_l) \quad \{u_1, \dots, u_K\} \rightarrow \text{orthonormal basis}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \quad (\text{an } n \times 1 \times 1 \times n \equiv n \times n) \quad \|u_l\|^2 = 1$$

$$\text{cov}(x_j, x_k) = E[(x_j - E[x_j])(x_k - E[x_k])]$$

$X = \text{Design Matrix } (n \times m)$

$$\Rightarrow \Sigma = [X^T X]$$

$$\Rightarrow \underset{u_1, \dots, u_K}{\operatorname{argmax}} \sum_{l=1}^K u_l^T [X^T X] u_l$$

★ $A \in \mathbb{R}^{n \times n}$

λ : - eigenvalue of A if $\exists x$ s.t. $Ax = \lambda x$ $\xrightarrow{\lambda \in \mathbb{R} + \cup \{0\}}$

$\lambda \neq 0$ \nearrow eigenvector

$$(A - \lambda I)x = 0$$

$$\Rightarrow \det(A - \lambda I) = 0$$

↳ polynomial in λ (degree n)

★ Use eigenvalue decomposition:

$$A = Q \Lambda Q^{-1} \quad (\Lambda \rightarrow \text{diagonal matrix})$$

$Q^{-1} = Q^T$, Q : rows & columns are \perp^T to each other (norm of each row & col. is 1)

(Q : orthonormal)

$O(n^3)$

↳ +ve semi-definite

$$A = X^T X$$

$$\Rightarrow \underset{u_1, \dots, u_K}{\operatorname{argmax}} \sum_{l=1}^K u_l^T (X^T X) u_l$$

u_1, \dots, u_K are principal eigenvectors of $X^T X$ (takes those eigenvectors whose eigenvalues are maximum)

Here, pick K principal eigenvectors.

* $X \in \mathbb{R}^{m \times n}$ $O(n^2 m)$ in $X^T X$ computation

* SVD:- Singular Value Decomposition

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

$$X = U D V^T$$

m n

($n-m$ entries 0 and initial m entries along diagonal non-zero)

$$m [U] \quad m [D] \quad n [V]$$

Columns of U are orthonormal eigenvectors of $X X^T$.

Columns of V are eigenvectors of $X^T X$.

$O(m^2 n)$ time

* If n is smaller do SVD on X^T else on X .

exam question

\rightarrow D:- entries of D are square roots of eigenvalues of $X^T X$ or $X X^T$.

$$X^T X = (U D V^T) (V D U^T)$$

$$= U (U D V^T V D U^T)$$

$$= U D^2 U^T$$

$$X^T X = \cancel{U D^2 U^T} \rightarrow \text{eigenvalues}$$

20/04/2018

If $m < n \rightarrow$ fat matrix

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

* SVD:- $X = U D V^T$ $D \in \mathbb{R}^{m \times n}$

$$m [U] \quad \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ 0 & 0 & 0 \end{bmatrix} \quad n [V]$$

U, V :- orthonormal (all rows & columns \perp to each other and norm = 1)

* Columns of U are eigenvectors of $X X^T$

columns of V are eigenvectors of $X^T X$.

Entries of $D \rightarrow$ sqrt of eigenvalues of $X X^T$.

Proof \rightarrow Assume $X^T X$ is full rank.

$$X^T X = V D^T U^T U D V^T = \cancel{V D^T U^T} V (D^T D) V^T \rightarrow \text{Eigen value decm. of } X^T X$$

SVD:

$$X = UDV^T \in \mathbb{R}^{m \times n} \quad m < n$$

$$= \begin{matrix} m \\ \vdots \\ m \end{matrix} \left[\begin{matrix} x_1 & x_2 & \dots & x_m \\ 0 & 0 & \dots & 0 \end{matrix} \right] \begin{matrix} n \\ \vdots \\ n \end{matrix}$$

U, V are orthonormal.

Columns of U are eigen vectors of XX^T

$$\| \cdot \|_V \quad \| \cdot \|_U \quad \| \cdot \|_{X^T X}$$

Assume XX^T is full rank.

$$X^T X = V D^T U^T U D V^T = V (D^T D) V^T \quad \left. \begin{matrix} \text{eigen value decomposition} \\ \dots \end{matrix} \right\}$$

$$XX^T = U D V^T D^T U^T = U (D D^T) V^T$$

$$D^T D = R^{n \times n} = \text{diag}(\Delta_{11}, \Delta_{22}, \dots, \Delta_{mm}, 0, 0, \dots, 0)$$

EM (Expectation Maximization):

Given $\{x^{(i)}\}_{i=1}^m$

Suppose: We have model for $\{x^{(i)}, y^{(i)}\}_{i=1}^m$

$P(x^{(i)}, y^{(i)}; \theta)$: Model

If we had labelled data :

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

$$\theta_{ML} = \arg \max_{\theta} \prod_{i=1}^m P(x^{(i)}, y^{(i)}; \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log(P(x^{(i)}, y^{(i)}; \theta))$$

Our case: (EM)

$$\{x^{(i)}\}_{i=1}^m$$

$$\theta_{ML} = \arg \max_{\theta} \prod_{i=1}^m P(x^{(i)}; \theta)$$

$$= \arg \max_{\theta} \prod_{i=1}^m \left[\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) \right]$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log \left(\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) \right)$$

EM:

L(A) : if we knew θ 's estimate $z^{(i)}$'s

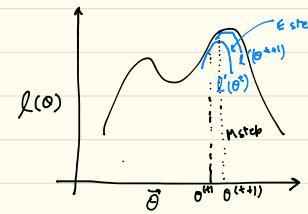
L(B) : if we knew $z^{(i)}$'s estimate θ 's

EM (Expectation Maximization) :

$$\{x^{(i)}\}_{i=1}^m, \text{ hidden: } \{z^{(i)}\}_{i=1}^m$$

$$\begin{aligned}\theta_{ML} &= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P(x^{(i)}; \theta) \\ &= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m \left[\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) \right] \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left(\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) \right)\end{aligned}$$

$\ell(\theta) \geq \ell'(\theta)$

① Jensen's Inequality :

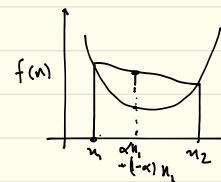
$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f(\alpha n_1 + (1-\alpha)n_2) \leq \alpha f(n_1) + (1-\alpha)f(n_2)$$

$P(x)$: over X

f is convex:

$$f(E_p(x)) \leq E_p[f(x)]$$



f is concave:

$$f(E_p(x)) \geq E_p(f(x))$$



f is strictly concave:

$$f(E_p(x)) > E_p(f(x)) \quad \text{unless } \exists c \text{ s.t. } P(X=c) = 1$$

Also,

$$E[f(x)] = f(E(x)) \text{ iff } \exists c, P(X=c) = 1$$

$$\ell(\theta) = \sum_{i=1}^m \left[\log \left(\sum_{z^{(i)}} \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \cdot Q_i(z^{(i)}) \right) \right]$$

where $Q_i(z^{(i)})$ is some distribution over $Z^{(i)}$.

$$\Downarrow \Downarrow \Downarrow \Downarrow$$

$$\geq \sum_{i=1}^m \left[\sum_{z^{(i)}} \left(\log \left(\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) Q_i(z^{(i)}) \right) \right] = \ell'(\theta)$$

If the quantity inside the log becomes constant then inequality becomes equality.

$$\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = c$$

$$\Rightarrow \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) = c \cdot \sum_{z^{(i)}} Q_i(z^{(i)})$$

$$P(x^{(i)}; \theta) = c \cdot 1$$

$$\therefore c = P(x^{(i)}; \theta)$$

$$Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{P(x^{(i)}; \theta)} = \frac{P(z^{(i)} | x^{(i)}; \theta)}{P(x^{(i)}; \theta)}$$

$Q_i(z^{(i)})$ is filling in missing values (by computing probability for $z^{(i)}$'s)

EM:

Intuition:

- ① If we knew θ 's, then we can compute
 $z^{(i)}$'s :- $Q_i(z^{(i)}) \leftarrow P(z^{(i)} | x^{(i)}; \theta)$

- ② If we knew $z^{(i)}$'s | $Q_i(z^{(i)})$, then we
can compute $\arg\max_{\theta} \sum_{z^{(i)}} \log \left(\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) Q_i(z^{(i)})$
 $\equiv \arg\max_{\theta} l(\theta)$

Formal:

$$t \leftarrow 0$$

$$\theta^{(t)} \leftarrow \text{initParams}()$$

do {

E step

$$t \leftarrow t + 1$$

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

M step

$$\theta^{(t+1)} = \arg\max_{\theta} \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \times \log \left(\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right)$$

constant

}

(until convergence)

Learning Theory:

$(x^{(i)}, y^{(i)}) \sim D \rightarrow$ underlying distribution.

$y^{(i)} \in \{0, 1\}$, $h \in H \rightarrow$ hypothesis space

$$\varepsilon^{(i)}(h) = 1\{h(x^{(i)}) \neq y^{(i)}\}$$

$$\hat{\varepsilon}(h) = \frac{1}{m} \sum_{i=1}^m 1\{h(x^{(i)}) \neq y^{(i)}\} \Rightarrow \text{Training Error}$$

$$\varepsilon(h) = \underset{(x,y) \sim D}{E} [1\{h(x) \neq y\}] \Rightarrow \text{generalization error}$$

Goal: Minimize $\varepsilon(h)$ \longrightarrow ① Minimize $\hat{\varepsilon}(h)$
 ② Minimize $|\hat{\varepsilon}(h) - \varepsilon(h)|$

$$P(|\varepsilon(h) - \hat{\varepsilon}(h)| > \gamma) \leq \delta$$

PAC-L \equiv Probably Approximately Correct

- ① $\gamma \rightarrow$ error
- ② $\delta \rightarrow$ probability of error
- ③ m
- ④ $|H|$:- important

Setting: ① Binary Classification
 ② H is finite {We will generalize later}

Coin Toss:

ϕ : Probability of Head (H)

$$\therefore P(Z=1) = \phi \quad \{1 \equiv \text{Head}\}$$

$Z \sim \text{Bernoulli}(\phi)$

$$E(Z) = \phi \cdot 1 + (1-\phi) \cdot 0 = \phi$$

$$\text{Var}(Z) = E(Z^2) - E^2(Z) = \phi - \phi^2 = \phi(1-\phi)$$

$$\text{Trial 1: } 51H \quad 49T \Rightarrow \hat{\phi} = 0.51$$

$$\text{Trial 2: } 45H \quad 55T \Rightarrow \hat{\phi} = 0.45$$

↓ γ times

$\hat{\phi}$: Random Variable ! {How?}

$\hat{\phi} = \frac{1}{m} \sum_{i=1}^m Z^{(i)}$ } Estimator of ϕ ,

$$E(\hat{\phi}) = E\left(\frac{1}{m} \sum_{i=1}^m Z^{(i)}\right) = \frac{1}{m} \sum_{i=1}^m E(Z^{(i)}) = \frac{\phi m}{m} = \phi \quad \{ \text{Unbiased (since } E(\hat{\phi}) = \phi\text{)}$$

To we don't need $Z^{(i)}$'s to be independent

$$\text{Var}(\hat{\phi}) = \text{Var}\left(\frac{1}{m} \sum_{i=1}^m Z^{(i)}\right) = \frac{1}{m^2} \sum_{i=1}^m \text{Var}(Z^{(i)}) = \frac{m\phi(1-\phi)}{m^2} = \frac{\phi(1-\phi)}{m} \quad \{ \text{As } m \uparrow \text{ Var}(\hat{\phi}) \downarrow \}$$

Pac Learning: $Z \sim \text{Bernoulli}(\phi)$

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

$E[\hat{\phi}] = \phi, \text{Var}(\hat{\phi}) = \frac{\phi(1-\phi)}{m}$

$P(|\phi - \hat{\phi}| > \gamma) \leq 2e^{-2\gamma^2 m} \quad \xrightarrow{\text{Chernoff Bound / Hoeffding inequality}}$

$\mathcal{E}(h) = \underset{(x,y) \sim D}{E}[1\{h(x) \neq y\}]$

$Z = 1\{h(x) \neq y\} \in \{0,1\}$

$E_Z[Z] = \mathcal{E}(h) = \phi$

$\hat{\phi} = \frac{1}{m} \sum_{i=1}^m Z^{(i)} = \hat{\mathcal{E}}(h)$

$\therefore P_Z(|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \geq \gamma) \leq \delta = 2e^{-2\gamma^2 m} \quad \text{independent of } H!$

$P_T(\exists h \in H, |\mathcal{E}(h) - \hat{\mathcal{E}}(h)| > \gamma) \leq \delta \equiv P_{\text{true}}(|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \leq \gamma) \geq 1 - \delta$

Define : $A_\delta := |\mathcal{E}(h_\delta) - \hat{\mathcal{E}}(h_\delta)| > \gamma$

$\forall \delta = P_T(A_\delta) \leq \delta$

$P_T(\bigcup A_\delta) \leq \sum_{\delta=1}^k P_T(A_\delta)$

$\therefore P_T(\exists h \in H, |\mathcal{E}(h) - \hat{\mathcal{E}}(h)| > \gamma) \leq \delta = k \cdot 2e^{-2\gamma^2 m} \quad \left. \right\} \text{Uniform Convergence Bound}$

$\text{As, } \delta \geq 2k e^{-2\gamma^2 m}$

$\Rightarrow m \geq \frac{1}{2\gamma^2} \log\left(\frac{2k}{\delta}\right) \quad \text{or} \quad \gamma \geq \sqrt{\frac{1}{2m} \log\left(\frac{2k}{\delta}\right)}$

What about best hypothesis?

 $\hat{h}^* := \text{Best hypothesis on the training set}$ $\hat{h}^* = \underset{h \in H}{\operatorname{argmin}} \hat{\mathcal{E}}(h) \quad \left. \right\} \text{Best in practice}$ $\hat{h}^* = \underset{h \in H}{\operatorname{argmin}} \mathcal{E}(h) \quad \left. \right\} \text{Best in principle}$ since \hat{h}^* was best on training

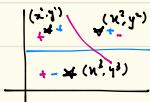
$$\mathcal{E}(\hat{h}^*) \leq \hat{\mathcal{E}}(\hat{h}^*) + \gamma \leq \hat{\mathcal{E}}(h^*) + \gamma \leq \mathcal{E}(h^*) + 2\gamma \quad \left. \right\} \text{with probability atleast } 1 - \delta.$$

due to Chernoff convergence bound

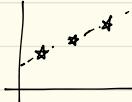
using uniform convergence bound

VC Dimension (Vapnik Chervonenkis Dimension):

$y \in \{0, 1\}$, H set of all hypothesis



These points are shattered



Are not shattered by linear separator. (E.g. config + - +)

by the space of linear separator

VC dimension: "capacity" of a hypothesis space to classify a set of points.

Size of the largest set that can be shattered



$h \in$ linear classifier

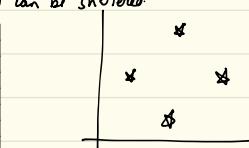
$|H|$ is ∞

\Rightarrow These points can't be shattered

by a linear separator

Infact, there is no set of 4 points

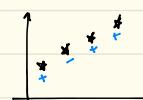
which can be shattered by linear sep.



$h \in$ quadratic classifier

$|H|$ is infinite

These set of points is shattered by a quad. classifier. Even 5 points can be done.



This can't be shattered by quad. classifier

The game is to find a set of K points S for any possible labelling, we are able to shatter them using our hypothesis.

VC dimension of Linear separator is 3, while that of quad. classifier is at least 5 (in 2-D space)

Definition:

Let H be a hypothesis class.

Let $S = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ where $x^{(i)} \in X \rightarrow$ instance space

Further, let $h \in H$ be boolean valued

$h: X \rightarrow \{0, 1\}$, $h(x) \in \{0, 1\}$

S is shattered by H iff $\exists h \in H$ st. h can realize all possible labelling over S .

Let $L = \{l^{(1)}, \dots, l^{(m)}\}$ be a labelling of S

$\nexists h \in H$ st. $h(x^{(i)}) = l^{(i)}$ $\forall i = 1 \dots m$

VC Dimension: Given a hypothesis space H of instance space X , it is the size of largest set S that is shattered by H .

Example:

$H = \text{set of all possible axis } || \text{ rectangles}$

$$X = \mathbb{R}^2$$

$$h_R \in H, h_R(x) = \begin{cases} 1 & \text{if } x \text{ lies inside } R \\ 0 & \text{o.w.} \end{cases}$$



4 points can be shattered



5 points can't be shattered

[using + to topmost, leftmost, rightmost, bottom most of the last on -].

$$\therefore \text{VC dim} = 4$$

• VC dimension of degree k polynomial over $x \in \mathbb{R} \rightarrow k+1$?

PAC bounds:

Finite $|H|$, $K = |H|$

$$\textcircled{1} \quad m \geq \frac{1}{2\gamma^2} \log \left(\frac{2K}{\delta} \right)$$

$$\textcircled{2} \quad \gamma \geq \sqrt{\frac{1}{2m} \log \left(\frac{2K}{\delta} \right)}$$

$$\textcircled{3} \quad \delta \geq 2K e^{-2\gamma^2 m}$$

Uniform convergence bound.

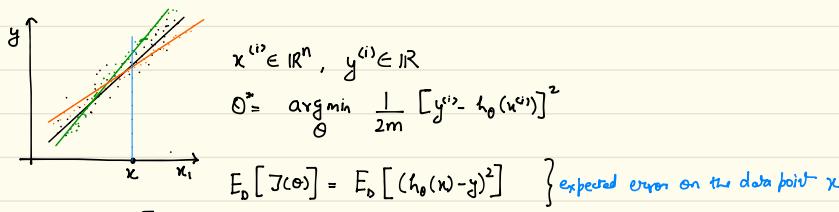
Infinite $|H|$, $d = \text{VC}(H)$

$$m \geq \frac{1}{\gamma} \left(4 \log \left(\frac{2}{\delta} \right) + 5d \log \left(\frac{12}{\delta} \right) \right)$$

$$\therefore m = O_{\epsilon, \gamma}(d) \Rightarrow \text{Linear in } d$$

$$\gamma \geq O(d m \log \frac{2}{\delta} + \frac{1}{m} \log \frac{1}{\delta})$$

- Bias Variance Tradeoff :



$$h_{\theta}(x) = E_D [h_{\theta}(x)]$$

$$[y - h_{\theta}(x)] \Rightarrow \text{Bias}$$

$$E_D [(h_{\theta}(x) - h_{\theta}(x))^2] \Rightarrow \text{Variance}$$

$$E_D [J(\theta)] = (\text{Bias})^2 + \text{Variance} \quad] \Rightarrow \text{Proof?}$$

inherent capacity
to fit training

Variation caused by different training samples

Now, high degree poly will have high variance {as slight change in training data can change pred quite a bit}

But, they also tend to have low bias since they are able to learn well.

High bias \Rightarrow underfitting

High variance \Rightarrow overfitting

- Summary:

$$\{x^{(i)}, y^{(i)}\}_{i=1}^m \Rightarrow \text{Supervised}$$

- Classification

Logistic	GDA / Naive Bayes	SVM	Decision Trees	Neural Nets	Deep Learning
----------	-------------------	-----	----------------	-------------	---------------

- ① Hypothesis Space
- ② Optimization metric / Loss fn
- ③ Optimization fn