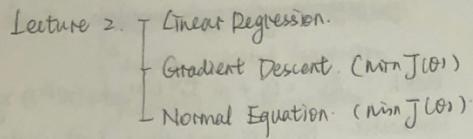


Cs229 (by Andrew Ng) – Notes by Xiaofan Wang

Name: XiaoFan Wang



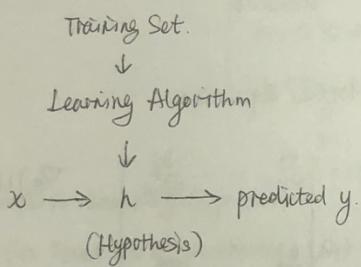
→ Supervised Learning.

Giving a training set: $\{(x^{(i)}, y^{(i)}) ; i=1, \dots, n\}$.

To learn a function h : $x \rightarrow y$.
(steering direction).

Regression: y is continuous.

Classification: y is discrete.



2. Linear Regression.

Living Area Size.	# Bedrooms	Price.
$x_1^{(1)} 2121$	3	400
$x_1^{(2)} 1600$	3	330
⋮	⋮	⋮

x : two-dimensional vectors in \mathbb{R}^2

x_1 : Size, x_2 : # bedrooms.

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$h(x) = \sum_{j=0}^2 \theta_j x_j$$

parameter $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$

$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$ Always!
Size
Bedrooms.

1. Linear Function: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ (θ : weights)

$$\text{Letting } x_0 = 1 : h(x) = \sum_{j=0}^2 \theta_j x_j = \theta^T x. (h_\theta(x) \approx y)$$

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



LMS Algorithm: choose θ to minimize $J(\theta)$.

(Least Mean Square).

→ Change θ to make $J(\theta)$ smaller.

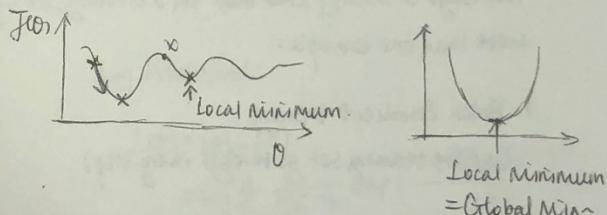
Until hopefully we converge to θ that $\min J(\theta)$.

that's consider the Gradient Descent Algorithm.

3. Gradient Descent Algorithm.

Start with θ . (Set $\theta = \vec{\theta}$).

Keep changing θ to reduce $J(\theta)$.



$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta) \quad (\alpha: \text{Learning Rate})$$

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \frac{1}{2} \sum_i (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\ &= (h_\theta(x) - y) \frac{\partial}{\partial \theta_j} \left(\sum_i \theta_j x_j - y \right) \\ &= (h_\theta(x) - y) \cdot x_j \end{aligned}$$

λ : training example #

x : Input (Features)

y : Output (Target Variable)

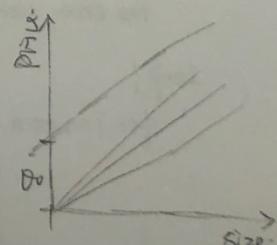
(x, y) : Training Examples

$(x^{(i)}, y^{(i)})$: i th training Examples

$x_1^{(i)}$: i runs from 1 to n .

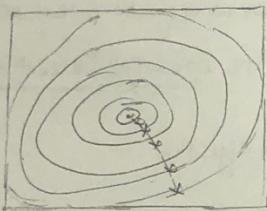
d : # features

$x^{(i)}$, $D(d+1)$ dimensional



LMS Rule: $\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$

$$\left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \cdot (h_\theta(x) - y) x_j \\ \theta_j := \theta_j - \alpha \cdot \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{array} \right. \quad \underbrace{\frac{\partial J(\theta)}{\partial \theta_j}}$$



$$\theta = \theta + \alpha \cdot \frac{1}{n} (y^{(i)} - h_\theta(x^{(i)})) x^{(i)}$$

Batch: $\theta_j := \theta_j + \alpha \cdot \frac{1}{n} (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$

Stochastic: $\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$

Comparison: Batch Gradient Descent has to scan through the entire training set before taking a single step \Rightarrow a costly operation if n is large.

Stochastic Gradient Descent can start making progress right away, and continues to make progress with each example it looks at.

IV. Normal Equation.

Two ways to modify LMS Rule for a training set of more than one example:

1. Batch Gradient Descent.

(entire training set go through every step).

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \cdot \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}, \quad (\text{for every } j)$$

}

$$\theta_j := \theta_j + \alpha \cdot \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}.$$

2. Stochastic Gradient Descent.

(update the parameters according to the gradient of the error with respect to that single training example only).

Loop?

for $i=1$ to n , {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}, \quad (\text{for every } j)$$

}

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} & \dots & \frac{\partial f}{\partial \theta_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial \theta_n} & \dots & \frac{\partial f}{\partial \theta_D} \end{bmatrix} \quad \nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_D} \end{bmatrix}$$

$$\nabla_{\theta} J(\theta) = \text{set } \vec{\theta}$$

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

since $h_\theta(x^{(i)}) = (x^{(i)})^T \theta$,

$$x\theta - \bar{y} = \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(n)})^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(n)}) - y^{(n)} \end{bmatrix}$$

using the fact. $\bar{z}^T z = \sum z_i^2$,

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (x\theta - \bar{y})^T (x\theta - \bar{y})$$

$$= \frac{1}{2} \nabla_{\theta} ((x\theta)^T x\theta - (x\theta)^T \bar{y} - \bar{y}^T (x\theta) + \bar{y}^T \bar{y})$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T (x^T x) \theta - 2(x^T \bar{y})^T \theta)$$

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

$$= x^T x\theta - x^T \bar{y} \quad (\text{set derivatives to 0})$$

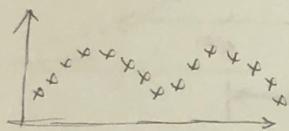
$$x^T x\theta = x^T \bar{y}$$

$$\Rightarrow \theta = (x^T x)^{-1} x^T \bar{y} \quad (\text{optimal value})$$

Name: XiaoFan Wang

- Lec 3. Linear Regression Model
- ↳ Locally Weighted Regression
 - ↳ Probabilistic Interpretation
 - ↳ Logistic Regression
 - ↳ Newton's Method.

- Locally Weighted Regression.



$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = \theta^T x$$

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

$$\begin{cases} \theta_0 + \theta_1 x_1 \\ \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 \log(x_3) \end{cases}$$

↑ ↑ ↑
 x_1 x_2 x_3

"Parametric" Learning Algorithm:

⇒ Fit Fixed set of parameters (θ_j) to data.

"Non-parametric" Learning Algorithm:

⇒ #parameters grow linearly with the size of data.

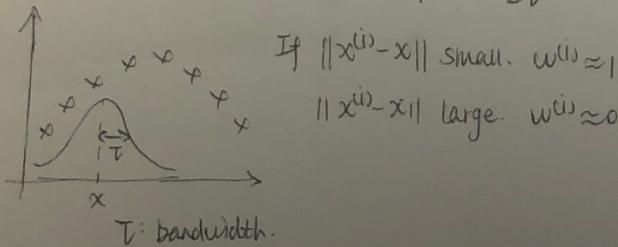
To evaluate h at certain x :

LR: Fit θ to minimize $\frac{1}{2} \sum (y^{(i)} - \theta^T x^{(i)})^2$. return $\theta^T x$

Locally Weighted Regression:

Fit θ to minimize $\sum_{i=1}^n w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$

$$w^{(i)} = \exp\left(-\frac{\|x^{(i)} - x\|^2}{2\sigma^2}\right)$$



= Probabilistic Interpretation.

$$\text{Assume: } y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$$

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

$$p(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

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

parametrized by

Random variable

mean

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

$$L(\theta) = p(\vec{y} | \vec{x}; \theta)$$

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

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

Log Likelihood:

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

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

$$= n \log \frac{1}{\sqrt{2\pi\sigma^2}} + \sum_{i=1}^n \frac{-(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

MLE: Maximum Likelihood Estimation.

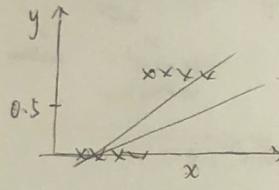
Choose θ to maximize $L(\theta)$

$$\text{minimize } \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = J(\theta)$$

III. Logistic Regression

Classification: ① Make assumption $p(y|x; \theta)$
 ② Compute θ by MLE

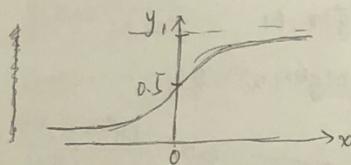
$y \in \{0, 1\}$. Binary Classification.



Want $h_{\theta}(x) \in [0, 1]$

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

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



Sigmoid function
 = logistic function

Linear regression: $h_{\theta}(x) = \theta^T x$

$$p(y=1|x; \theta) = h_{\theta}(x)$$

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

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

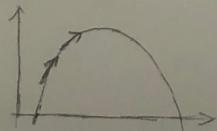
$$L(\theta) = p(\vec{y}|\vec{x}; \theta)$$

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

$$= \prod_{i=1}^n h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

$$l(\theta) = \log L(\theta)$$

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



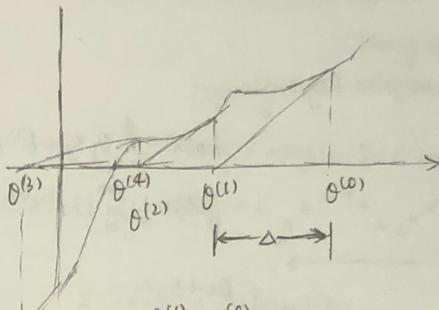
$$\theta_j := \theta_j + \alpha \cdot \sum_{i=1}^n \underbrace{(y^{(i)} - h_{\theta}(x^{(i)}))}_{\frac{\partial l(\theta)}{\partial \theta_j}} \cdot x_j^{(i)}$$

IV. Newton's Method

Want: maximize $L(\theta)$

$$l'(\theta) = 0$$

Max / Min \Leftrightarrow derivative = 0



$$\theta^{(1)} = \theta^{(0)} - \Delta$$

$$f'(\theta^{(0)}) = \frac{f(\theta^{(0)})}{\Delta} = \frac{\text{Height}}{\text{Base}} \quad (\Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})})$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})} \quad (f(\theta) = l'(\theta))$$

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

Quadratic Convergence: $0.1 \rightarrow 0.01 \rightarrow 0.0001$

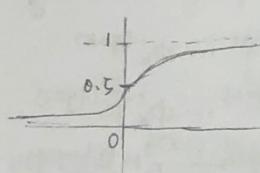
$$\overrightarrow{\theta} : \quad \theta^{(t+1)} = \theta^{(t)} + H^{-1} \nabla_{\theta} l \quad \begin{matrix} \leftarrow \\ \text{Vector } \mathbb{R}^{d+1} \end{matrix} \quad \text{IR}^{(d+1) \times (d+1)}$$

$$\text{Hessian } H = H_{ij} = \frac{\partial^2 l}{\partial \theta_i \cdot \partial \theta_j}$$

Name: XiaoFan Wang

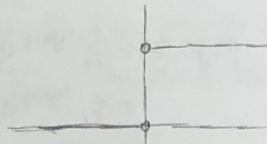
Leu 4. Perception
 Exponential Family
 Generalized Linear Models.
 Softmax Regression (Multi-class Classification)

1. Perception



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

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



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

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

Perception Learning Algorithm:

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

$$\begin{aligned} P(y; \phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp(y \log \phi + (1-y) \log (1-\phi)) \\ &= \exp((\log(\frac{\phi}{1-\phi})) y + \log(1-\phi)) \end{aligned}$$

Thus, $\eta = \log(\frac{\phi}{1-\phi})$, $\phi = \frac{1}{1+e^{-\eta}}$

(family sigmoid func.)

$$T(y) = y.$$

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

$$b(y) = 1$$

2. Gaussian Distribution

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

$$\text{Thus, } \eta = \mu$$

$$T(y) = y.$$

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

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

3. Exponential Family.

$$P(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

η : natural parameter

$T(y)$: sufficient statistic for Distributions.

$a(\eta)$: Log Partition Function

4. Bernoulli (ϕ) Distribution.

$$y \in \{0, 1\}. \quad \begin{cases} P(y=1; \phi) = \phi \\ P(y=0; \phi) = 1-\phi \end{cases}$$

5. Constructing GLM.

Assumption: 1. $y|x; \theta \sim \text{Exponential Family } (\eta)$

2. predict the expected value: $T(y) = y$
 prediction output: $h(x) = E[y|x]$

$$3. \eta = \theta^T x$$

$$\downarrow \text{vector valued } \eta_i = \theta_i^T x$$

1. Ordinary Least Squares.

Gaussian $N(\mu; \sigma^2)$

$y|x; \theta \sim N(\mu, \sigma^2)$

Gaussian as an Exponential Family Distribution: $\mu = \eta$.

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

$$= \mu$$

$$= \eta$$

$$= \theta^T x$$

$(Ty)_i$: the i th vector Ty

$$\{y_2 = 3\} = 0$$

$$\{y_3 = 5 - 2\} = 1$$

$$(Ty)_i = \{y_i = i\}$$

$$E[(Ty)_i] = p(y_i = i) = \phi_i$$

Multinomial as a member of Exponential Family:

$$p(y; \phi) = \phi_1^{y_1} \phi_2^{y_2} \cdots \phi_k^{y_k}$$

$$= \phi_1^{y_1} \phi_2^{y_2} \cdots \phi_k^{y_k} \prod_{i=1}^{k-1} \phi_i^{y_i}$$

$$= \phi_1^{(Ty)_1} \phi_2^{(Ty)_2} \cdots \phi_k^{(Ty)_k}$$

$$= \exp((Ty)_1 \log(\phi_1) + (Ty)_2 \log(\phi_2) + \cdots + (Ty)_k \log(\phi_k))$$

$$= \exp((Ty)_1 \log(\frac{\phi_1}{\phi_k}) + (Ty)_2 \log(\frac{\phi_2}{\phi_k}) + \cdots + (Ty)_k \log(\frac{\phi_k}{\phi_k}) + \log(\phi_k))$$

$$= b(y) \exp(\eta^T T(y) - a(\eta))$$

$$\eta = \begin{bmatrix} \log(\phi_1/\phi_k) \\ \log(\phi_2/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix}, \quad a(\eta) = -\log(\phi_k) \\ b(y) = 1$$

2. Logistic Regression.

Bernoulli Distribution as an Exponential Family:

$$\phi = \frac{1}{1+e^{-\eta}}$$

$$y|x; \theta \sim \text{Bernoulli}(\phi)$$

$$E[y|x; \theta] = \phi$$

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

$$= \phi$$

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

3. Softmax Regression. (Multinomial Distribution)

$\phi_1 \cdots \phi_k$ satisfy $\sum_{i=1}^k \phi_i = 1$

$$\phi_i = p(y=i; \phi)$$

$$p(y=k; \phi) = 1 - \sum_{i=1}^{k-1} \phi_i$$

$$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$$

Define $T(y) \in \mathbb{R}^{k-1}$:

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad T(3) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots, \quad T(k-1) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, \quad T(k) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

log-likelihood:

$$L(\theta) = \sum_{i=1}^n \log p(y^{(i)}|x^{(i)}, \theta) \\ = \sum_{i=1}^n \log \left(\frac{e^{\theta^T x^{(i)}}}{\sum_{j=1}^k e^{\theta^T x^{(i)}}} \right)^{I(y^{(i)}=j)}$$

$$= \sum_{i=1}^n \log \left(\frac{e^{\theta^T x^{(i)}}}{\sum_{j=1}^k e^{\theta^T x^{(i)}}} \right)^{I(y^{(i)}=j)}$$

$$\eta_i = \log \frac{\phi_i}{\phi_k}, \quad e^{\eta_i} = \frac{\phi_i}{\phi_k}, \quad \phi_k \cdot e^{\eta_i} = \phi_i$$

$$\phi_k \cdot \sum_{i=1}^k e^{\eta_i} = \sum_{i=1}^k \phi_i = 1, \quad \phi_k = \frac{1}{\sum_{i=1}^k e^{\eta_i}}$$

$$\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}, \quad p(y=i|x; \theta) = \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} = \frac{e^{\theta^T x}}{\sum_{j=1}^k e^{\theta^T x}}$$

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

$$= E \left[\begin{bmatrix} I(y=1) \\ \vdots \\ I(y=k-1) \end{bmatrix} | x; \theta \right] = \begin{bmatrix} \phi_1 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{e^{\theta^T x}}{\sum_{j=1}^k e^{\theta^T x}} \\ \vdots \\ \frac{e^{\theta^T x}}{\sum_{j=1}^k e^{\theta^T x}} \end{bmatrix}$$

Name: XiaoFan Wang

Lee 5 Generative Learning Algorithms.
Lee 6

Discriminative Learning Algorithms:

Input x mapping directly to label {0, 1}

Generative Learning Algorithms:

try to model $p(x|y)$, $p(y)$.

For instance: Dog(0), Elephant(1).

$$\begin{cases} p(x|y=0) : \text{dogs' features} \\ p(x|y=1) : \text{elephants' features} \end{cases}$$

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

$$p(x) = p(x|y=1)p(y=1) + p(x|y=0)p(y=0)$$

$$\begin{aligned} \arg \max_y p(y|x) &= \arg \max_y \frac{p(x|y)p(y)}{p(x)} \\ &= \arg \max_y p(x|y)p(y) \end{aligned}$$

1. Gaussian Discriminant Analysis (GDA)

1. Multivariate Normal Distribution.

(Multivariate Gaussian Distribution)

$$\left\{ \begin{array}{l} \text{mean vector } \mu \in \mathbb{R}^d \\ \text{covariance matrix } \Sigma \in \mathbb{R}^{d \times d} \quad (\Sigma \geq 0) \\ N(\mu, \Sigma) \end{array} \right.$$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu))$$

$$E[x] = \int_x x p(x; \mu, \Sigma) dx = \mu$$

$$\begin{aligned} \text{Cov}(z) &= E[(z - E[z])(z - E[z])^T] \\ &= E[zz^T] - (E[z])(E[z])^T \end{aligned}$$

If $x \sim N(\mu, \Sigma)$, then $\text{Cov}(x) = \Sigma$

2. Gaussian Discriminant Analysis Model

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

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

$$x|y=1 \sim N(\mu_1, \Sigma_1)$$

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

$$p(x|y=0) = \frac{1}{(2\pi)^{d/2} |\Sigma_0|^{1/2}} \exp(-\frac{1}{2}(x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0))$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma_1|^{1/2}} \exp(-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1))$$

Log-likelihood:

$$\begin{aligned} l(\phi, \mu_0, \mu_1, \Sigma_0) &= \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma_0) \\ &= \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma_0) p(y^{(i)}; \phi) \end{aligned}$$

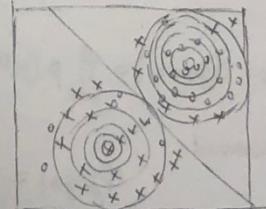
Maximize l :

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

$$\mu_0 = \frac{\sum_{i=1}^n \{y^{(i)}=0\} x^{(i)}}{\sum_{i=1}^n \{y^{(i)}=0\}}$$

$$\mu_1 = \frac{\sum_{i=1}^n \{y^{(i)}=1\} x^{(i)}}{\sum_{i=1}^n \{y^{(i)}=1\}}$$

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



3. GDA and Logistic Regression

$$p(y=1|x; \Phi, \Sigma_1, \mu_0, \mu_1) = \frac{1}{1 + \exp(-\Phi^T x)}$$

If $x_i | y=0 \sim \text{poisson}(\lambda_0)$
 $x_i | y=1 \sim \text{poisson}(\lambda_1)$

$$\begin{aligned} p(y=1|x) &= \frac{p(x|y=1)p(y=1)}{p(x)} \\ &= \frac{\prod_{j=1}^d p(x_j|y=1) \cdot p(y=1)}{\sum_{j=1}^d p(x_j|y=1)p(y=1) + \sum_{j=1}^d p(x_j|y=0)p(y=0)} \end{aligned}$$

Naive Bayes

$p(x|y) = x_i$ is conditionally independent given y .
(Naive Bayes Assumption).

Naive Bayes Classifier:

$$p(x_1 \dots x_{50000}|y)$$

$$= p(x_1|y)p(x_2|y, x_1)p(x_3|y, x_1, x_2) \dots p(x_{50000}|y, x_1 \dots x_{49999})$$

$$= p(x_1|y)p(x_2|y)p(x_3|y) \dots p(x_{50000}|y)$$

$$= \prod_{j=1}^d p(x_j|y)$$

$$\Phi_j|y=1 = p(x_j=1 | y=1)$$

$$\Phi_j|y=0 = p(x_j=1 | y=0)$$

$$\Phi_y = p(y=1)$$

$$L(\Phi_y, \Phi_j|y=1, \Phi_j|y=0) = \prod_{i=1}^n p(x_i^{(i)}, y^{(i)})$$

Maximum Likelihood:

$$\left\{ \begin{array}{l} \Phi_j|y=1 = \frac{\sum_{i=1}^n \{x_j^{(i)}=1 \wedge y^{(i)}=1\}}{\sum_{i=1}^n \{y^{(i)}=1\}} \\ \Phi_j|y=0 = \frac{\sum_{i=1}^n \{x_j^{(i)}=1 \wedge y^{(i)}=0\}}{\sum_{i=1}^n \{y^{(i)}=0\}} \\ \Phi_y = \frac{\sum_{i=1}^n \{y^{(i)}=1\}}{n} \end{array} \right.$$

1. Laplace Smoothing.

$$\Phi_{35000}|y=1 = \frac{\sum_{i=1}^n \{x_{35000}^{(i)}=1 \wedge y^{(i)}=1\}}{\sum_{i=1}^n \{y^{(i)}=1\}} = 0$$

$$\Phi_{35000}|y=0 = \frac{\sum_{i=1}^n \{x_{35000}^{(i)}=1 \wedge y^{(i)}=0\}}{\sum_{i=1}^n \{y^{(i)}=0\}} = 0$$

$$\text{Maximize Likelihood: } \Phi_j = \frac{\sum_{i=1}^n \{z^{(i)}=j\}}{n}$$

Φ_j might end up as zero \Rightarrow problem

$$\text{Laplace Smoothing: } \Phi_j = \frac{1 + \sum_{i=1}^n \{z^{(i)}=j\}}{k+n}$$

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

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

2. Event Models for Text Classification

$$p(y) \prod_{j=1}^d p(x_j|y)$$

$$\text{Likelihood: } L(\Phi_y, \Phi_k|y=0, \Phi_k|y=1)$$

$$= \prod_{i=1}^n p(x_i^{(i)}, y^{(i)})$$

$$= \prod_{i=1}^n \left(\prod_{j=1}^d p(x_j^{(i)}|y; \Phi_k|y=0, \Phi_k|y=1) \cdot p(y^{(i)}; \Phi_y) \right)$$

Maximizing Likelihood:

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

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

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

If applying Laplace smoothing:

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

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

While not necessarily the very best classification algorithm, the Naive Bayes classifier often work surprisingly well. It is often also a very good "first thing to try": given its simplicity and ease of implementation.

Name: Xiaofan Wang

Lec 7. Kernel Methods & SVM

→ Kernel Methods

1. feature maps

$$y = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0$$

let function $\phi: \mathbb{R} \rightarrow \mathbb{R}^d$ defined as :

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \in \mathbb{R}^d$$

let $\theta \in \mathbb{R}^d : \theta_0, \theta_1, \theta_2, \theta_3$.

$$\theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0 = \theta^T \phi(x)$$

ϕ : feature map (map the attributes to the features)

2. LMS with features

$$\begin{aligned} \text{Batch Gradient Descent: } \theta &:= \theta + \alpha \cdot \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \\ &= \theta + \alpha \cdot \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \end{aligned}$$

let $\phi: \mathbb{R}^d \Rightarrow \mathbb{R}^P$, map x to the feature $\phi(x)$ in \mathbb{R}^P

$$\left\{ \begin{array}{l} \theta := \theta + \alpha \cdot \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \end{array} \right.$$

$$\left. \begin{array}{l} \text{Stochastic Gradient Descent: } \theta := \theta + \alpha \cdot (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \end{array} \right.$$

3. LMS with the kernel trick.

$\phi(x)$ is high-dimensional (dimensional $\vec{\theta}$ by coefficients $\beta_1 \dots \beta_n$)

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_2 x_1 \\ x_3 \\ \vdots \end{bmatrix}$$

At initialization: $\theta = 0 = \sum_{i=1}^n 0 \cdot \phi(x^{(i)})$

At some point: $\theta = \sum_{i=1}^n \beta_i \phi(x^{(i)})$

$$\begin{aligned} \theta &:= \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \\ &= \sum_{i=1}^n \beta_i \phi(x^{(i)}) + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \\ &= \sum_{i=1}^n (\underbrace{\beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)}))}_{\text{New } \beta_i}) \cdot \phi(x^{(i)}) \end{aligned}$$

$$\beta_i := \beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)}))$$

Replace θ by $\theta = \sum_{j=1}^n \beta_j \phi(x^{(j)})$:

$$\forall i \in \{1 \dots n\}, \beta_i := \beta_i + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j \phi(x^{(j)})^T \phi(x^{(i)}))$$

$$\downarrow$$

$$\langle \phi(x^{(i)}), \phi(x^{(i)}) \rangle$$

all pairs: $O(p)$.

compute $\phi(x^{(i)})$:

$$\begin{aligned} \langle \phi(x), \phi(z) \rangle &= 1 + \sum_{i=1}^d x_i z_i + \sum_{i,j \in \{1 \dots d\}} x_i x_j z_i z_j + \sum_{i,k \in \{1 \dots d\}} x_i x_j x_k z_i z_j z_k \\ &= 1 + \sum_{i=1}^d x_i z_i + (\sum_{i=1}^d x_i z_i)^2 + (\sum_{i=1}^d x_i z_i)^2 \\ &= 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3 \end{aligned}$$

Kernel: $K(x, z) \triangleq \langle \phi(x), \phi(z) \rangle$

Final Algorithm:

compute all the values $K(x^{(i)}, x^{(j)}) \triangleq$
using equation for $i, j \in \{1 \dots n\}$. $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$
Set $\beta := 0$.

Loop: $\forall i \in \{1 \dots n\}$

$$\beta_i := \beta_i + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}))$$

let K be the $n \times n$ matrix
with $K_{ij} = K(x^{(i)}, x^{(j)})$

$$\beta := \beta + \alpha (\vec{y} - K\beta)$$

compute prediction:

$$\begin{aligned} \theta^T \phi(x) &= \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x) \\ &= \sum_{i=1}^n \beta_i K(x^{(i)}, x) \end{aligned}$$

4. Properties of kernels

$$\textcircled{1} \quad K(x, z) = \phi(x)^T \phi(z)$$

Suppose $x, z \in \mathbb{R}^d$, $K(x, z) = (x^T z)^2$

$$K(x, z) = \left(\sum_{i=1}^d x_i z_i \right) \left(\sum_{j=1}^d x_j z_j \right)$$

$$= \sum_{i=1}^d \sum_{j=1}^d x_i x_j z_i z_j$$

$$= \sum_{i,j=1}^d (x_i x_j)(z_i z_j)$$

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

feature mapping:

$$\begin{aligned} \Phi(x) &= \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \\ &\text{High-dimension.} \\ &\text{require } O(d^2) \text{ time} \\ K(x, z) &\text{take only } O(d) \text{ time.} \end{aligned}$$

$$\textcircled{2} \quad K(x, z) = (x^T z + c)^2$$

$$= \sum_{i,j=1}^d (x_i x_j)(z_i z_j) + \sum_{i=1}^d (\sqrt{c} x_i)(\sqrt{c} z_i) + c^2$$

\textcircled{3} Gaussian kernel (kernels as similarity metrics)

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

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

\textcircled{4} Necessary conditions for valid kernels.

$$\begin{aligned} \text{Kernel Matrix: } K_{ij} &= K(x^{(i)}, x^{(j)}) \\ &= \phi(x^{(i)})^T \phi(x^{(j)}) \\ &= \phi(x^{(j)})^T \phi(x^{(i)}) \\ &= K(x^{(j)}, x^{(i)}) = k_{ji} \end{aligned}$$

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

$$= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j$$

$$= \sum_i \sum_j z_i \sum_k \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j$$

$$= \sum_k \sum_i \sum_j z_i \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j$$

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

$$\sum_{i,j} a_i a_j = (\sum_i a_i)^2 \text{ for } a_i = z_i \phi_k(x^{(i)})$$

feature mapping:

$$\begin{aligned} \Phi(x) &= \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \\ \sqrt{c} x_1 \\ \sqrt{c} x_2 \\ \sqrt{c} x_3 \\ c \end{bmatrix} \end{aligned}$$

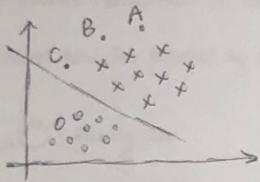
$K(x, z) = (x^T z + c)^2$ to a feature mapping $\binom{d+k}{k}$ feature space.

Despite working in $O(d^k)$ -dimensional space,

$K(x, z)$ takes only $O(d)$ time.

2. Support Vector Machines

1. Margins: Intuition



Training Example: Positive & Negative.

Separate hyperplane: A Decision Boundary (by the equation: $\mathbf{w}^T \mathbf{x} + b = 0$)

2. Notation

Class label: $y \in \{-1, 1\}$

classifier: $h_{\mathbf{w}, b}(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b)$

$\begin{cases} g(z) = 1 & \text{if } z \geq 0 \\ g(z) = -1 & \text{otherwise.} \end{cases}$

(separate from the other parameters)

intercept b

3. Functional and Geometric Margins

Given a training example $(\mathbf{x}^{(i)}, y^{(i)})$,

Define the functional margin of (\mathbf{w}, b) :

$$\hat{\gamma}^{(i)} = y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

If $y^{(i)}=1$, then for the functional margin to be large,
we need $\mathbf{w}^T \mathbf{x}^{(i)}+b$ to be a large positive number.

If $y^{(i)}=-1$, then for the functional margin to be large,
we need $\mathbf{w}^T \mathbf{x}^{(i)}+b$ to be a large negative number.

If $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)}+b) > 0$, then prediction is correct.

Hence, a large functional margin represents a confident
and a correct prediction.

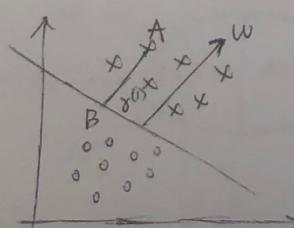
Given a training set $S = \{(\mathbf{x}^{(i)}, y^{(i)}) ; i=1 \dots n\}$,

$$\hat{\gamma} = \min_{i=1 \dots n} \hat{\gamma}^{(i)}$$

$A: \mathbf{x}^{(i)}$.

$B: \mathbf{x}^{(i)} - \hat{\gamma}^{(i)} \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}$ lies on decision boundary

$$\mathbf{w}^T (\mathbf{x}^{(i)} - \hat{\gamma}^{(i)} \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}) + b = 0$$



$$\gamma^{(i)} = \frac{\mathbf{w}^T \mathbf{x}^{(i)} + b}{\|\mathbf{w}\|} = \left(\frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x}^{(i)} + \frac{b}{\|\mathbf{w}\|}$$

$$\gamma^{(i)} = y^{(i)} \left(\left(\frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x}^{(i)} + \frac{b}{\|\mathbf{w}\|} \right)$$

$$\gamma = \min_{i=1 \dots n} \gamma^{(i)}$$

4. Optimal Margin Classifier

$$\max_{\gamma, w, b} \gamma$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \gamma, \quad (i=1 \dots n)$$

$$\text{If } \|w\|=1, \quad \gamma = \frac{\hat{\gamma}}{\|w\|}$$

$$\max_{\gamma, w, b} \frac{\hat{\gamma}}{\|w\|}$$

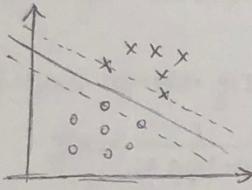
$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, \quad (i=1 \dots n)$$

$$\text{If } \hat{\gamma}=1, \quad \gamma = \frac{\hat{\gamma}}{\|w\|} \sim \frac{1}{\|w\|} \sim \|w\|^2$$

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad (i=1 \dots n)$$

$$\text{constraint: } g_i(w) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0$$



Lagrangian for optimization:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1]$$

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

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

$$\text{Derivative to } b: \frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

$$L(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} - b \cdot \sum_{i=1}^n \alpha_i y^{(i)}$$

$$L(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} \quad (\alpha_i \geq 0)$$

$$\max_{\alpha} \quad L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } \alpha_i \geq 0, \quad (i=1 \dots n)$$

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

$$b^* = - \frac{\max_{i:y^{(i)}=1} w^* T x^{(i)} + \min_{i:y^{(i)}=-1} w^* T x^{(i)}}{2}$$

$$w^T x + b = (\sum_{i=1}^n \alpha_i y^{(i)} x^{(i)})^T x + b$$

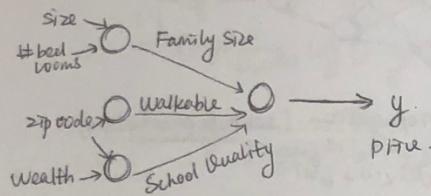
$$= \sum_{i=1}^n \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b.$$

Name: Xiaofan Wang

Lee 8. Lee's Neural Networks.

-> Neural Networks (end-to-end learning)

Small NN for predicting housing prices:



Input Layer: Four Features x_1, x_2, x_3, x_4
 Hidden Layer: Three Internal Neurons (Black Box)
 Output Layer: One Output Neuron.

If logistic regression $g(z)$ as a single neuron:

$$g(x) = \frac{1}{1 + \exp(-w^T x)}$$

Break $g(x)$ into two computations:

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

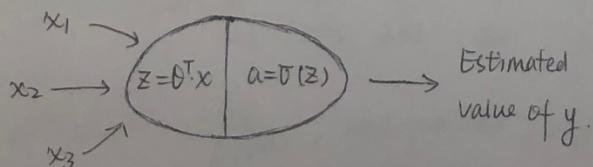
$$a = \sigma(z) \cdot (\sigma(z) = \frac{1}{1+e^{-z}})$$

$a = g(z) \rightarrow$ activation function:

$$\begin{cases} g(z) = \frac{1}{1+e^{-z}} \text{ (Sigmoid)} \\ g(z) = \max(z, 0) \text{ (ReLU)} \\ g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \text{ (tanh)} \end{cases}$$

In general, $g(z)$ is non-linear function.

For instance, Logistic Regression as a single neuron:



In the 1st layer:

$$1^{\text{st}} \text{ Hidden Unit: } z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = g(z_1^{[1]})$$

$$2^{\text{nd}} \text{ Hidden Unit: } z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = g(z_2^{[1]})$$

$$3^{\text{rd}} \text{ Hidden Unit: } z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = g(z_3^{[1]})$$

1st Layer Activations:

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

$$2^{\text{nd}} \text{ Layer Activation: } a_1^{[2]} = g(z_1^{[2]})$$

↓
final output prediction.

2. Vectorization

1. Vectorizing the Output Computation

$$\begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{bmatrix} = \begin{bmatrix} \cdots & w_1^{[1]T} & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & w_4^{[1]T} & \cdots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}$$

$$z^{[1]} \in \mathbb{R}^{4 \times 1} \quad w^{[1]} \in \mathbb{R}^{4 \times 3} \quad x \in \mathbb{R}^{3 \times 1} \quad b^{[1]} \in \mathbb{R}^{4 \times 1}$$

$$z^{[1]} = w^{[1]} x + b^{[1]}, \quad a^{[1]} = g(z^{[1]})$$

$$\text{Define } g(z) \text{ as: } g(z) = \frac{1}{1+e^{-z}} \quad (z \in \mathbb{R})$$

Matlab pseudocode: $g(z) = 1 ./ (1 + \exp(-z))$
where $z \in \mathbb{R}^d$

$$\begin{bmatrix} z^{[2]} \\ \vdots \\ z^{[4]} \end{bmatrix} = \begin{bmatrix} w^{[2]} \\ \vdots \\ w^{[4]} \end{bmatrix} \begin{bmatrix} a^{[1]} \\ \vdots \\ a^{[1]} \end{bmatrix} + \begin{bmatrix} b^{[2]} \\ \vdots \\ b^{[4]} \end{bmatrix}, \quad a^{[2]} = g(z^{[2]})$$

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

$$= w^{[2]} g(z^{[1]}) \quad \text{By definition}$$

$$= w^{[2]} z^{[1]} \quad \text{since } g(z) = z$$

$$= \tilde{w} x \quad \text{where } \tilde{w} = w^{[2]} w^{[1]}$$

2. Vectorization Over Training Examples

Activations for each examples :

$$\begin{cases} z^{1} = w^{[1]} x^{(1)} + b^{[1]} \\ z^{[1](2)} = w^{[1]} x^{(2)} + b^{[1]} \\ z^{[1](3)} = w^{[1]} x^{(3)} + b^{[1]} \end{cases}$$

$$\vec{x} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix}$$

$$z^{[1]} = \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} \\ | & | & | \end{bmatrix} = w^{[1]} \vec{x} + b^{[1]}$$

$$\tilde{b}^{[1]} = \begin{bmatrix} | \\ b^{[1]} \\ | \\ b^{[1]} \\ | \\ b^{[1]} \end{bmatrix} \quad z^{[1]} = w^{[1]} \vec{x} + \tilde{b}^{[1]}$$

① Initialize the parameters $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$.

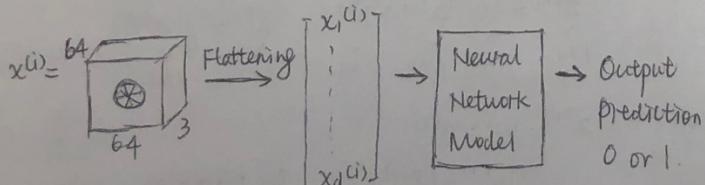
to small random numbers. For each example, we compute the output "probability" from the sigmoid function $a^{[2](i)}$.

② Using the logistic regression log likelihood:

$$\sum_{i=1}^n [y^{(i)} \log a^{[2](i)} + (1-y^{(i)}) \log (1-a^{[2](i)})]$$

③ Maximize the function using Gradient Descent. The Maximization procedure corresponds to training the Neural Network.

3. Backpropagation



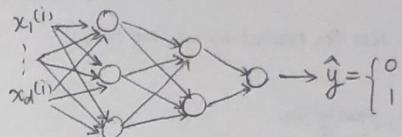
NN Model :

(i) Network Architecture : how many layers, neurons.
how layers are connected.

(ii) parameters : initialization, optimization.
Analyzing these parameters.

1. Parameter Initialization

Fully-connected Layer : All inputs are connected to (1st Hidden Layer), all neurons in the next layer.



① Forward propagation (compute parameters) :

$$z^{[1]} = w^{[1]} \vec{x} + b^{[1]}$$

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

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

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

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

$$\hat{y}^{(i)} = a^{[3]} = g(z^{[3]})$$

② Update parameters (compute Loss L) :

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

$$w^{[l]} = w^{[l]} - \alpha \frac{\partial L}{\partial w^{[l]}} \quad (\text{learning rate} \times \text{Gradient})$$

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

$$a^{[l+1]} = g(z^{[l]}) = g(w^{[l]} x^{(i)} + b^{[l]}).$$

Each Neuron will receive the exact same gradient update value.

$$\text{Initialization: } w^{[l]} \sim N(0, \sqrt{\frac{2}{n^{[l]} + n^{[l+1]}}})$$

↓
the Number of neuron in layer l.

Variance $\sigma^{(in)} \approx \sigma^{(out)}$.

2. optimization (using Stochastic Gradient Descent)

$w^{[3]}$ on the loss is more complex than that of $w^{[2]}$.

$w^{[2]}$ is "closer" to the output \hat{y} .

$$\begin{aligned}
 \frac{\partial L}{\partial w^{[3]}} &= \frac{-\partial}{\partial w^{[3]}} ((1-y) \log(1-\hat{y}) + y \log(\hat{y})) \\
 &= -(1-y) \frac{\partial}{\partial w^{[3]}} \log(1-g(w^{[2]}a^{[2]} + b^{[2]})) - y \frac{\partial}{\partial w^{[3]}} \log(g(w^{[2]}a^{[2]} + b^{[2]})) \\
 &= -(1-y) \frac{1}{1-g(w^{[2]}a^{[2]} + b^{[2]})} (-1) g'(w^{[2]}a^{[2]} + b^{[2]}) a^{[2]T} - y \cdot \frac{1}{g(w^{[2]}a^{[2]} + b^{[2]})} \cdot g'(w^{[2]}a^{[2]} + b^{[2]}) a^{[2]T} \\
 &= (1-y) \sigma(w^{[2]}a^{[2]} + b^{[2]}) a^{[2]T} - y(1-\sigma(w^{[2]}a^{[2]} + b^{[2]})) a^{[2]T} \\
 &= (a^{[2]} - y) \cdot a^{[2]T} \\
 &= (a^{[2]} - y) \cdot a^{[2]T}
 \end{aligned}$$

The derivative sigmoid function: $g = \sigma' = \sigma(1-\sigma)$

$$a^{[2]} = \sigma(w^{[2]}a^{[2]} + b^{[2]})$$

L depend on $\hat{y} = a^{[2]}$

$$\begin{aligned}
 \frac{\partial L}{\partial w^{[2]}} &= \frac{\partial L}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial w^{[2]}} \cancel{\frac{\partial L}{\partial a^{[3]}}} \\
 &= \frac{\partial L}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial w^{[2]}} \\
 &= \underbrace{\frac{\partial L}{\partial a^{[3]}}}_{\frac{\partial L}{\partial z^{[3]}}} \cdot \frac{\partial a^{[3]}}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial w^{[2]}} \\
 &= \underbrace{\frac{\partial L}{\partial a^{[3]}}}_{\frac{\partial L}{\partial z^{[3]}}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial z^{[3]}}}_{\frac{\partial a^{[3]}}{\partial w^{[2]}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}}}_{\frac{\partial z^{[3]}}{\partial w^{[2]}}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial w^{[2]}}}_{(a^{[2]} - y) \cdot w^{[2]}} = (a^{[2]} - y) \cdot w^{[2]} \cdot g'(z^{[2]}) \cdot a^{[2]T}
 \end{aligned}$$

$$\frac{\partial L}{\partial w^{[2]}} = \underbrace{(a^{[2]} - y)}_{2 \times 3} \underbrace{w^{[2]}}_{1 \times 1} \underbrace{g'(z^{[2]})}_{1 \times 2} \underbrace{a^{[2]T}}_{2 \times 1} = w^{[2]T} \circ g'(z^{[2]}) (a^{[2]} - y) a^{[2]T}$$

(mini-batch stochastic GD)

Another Method: Momentum

$$\begin{aligned}
 w^{[2]} &= w^{[2]} - \alpha \cdot \frac{\partial J}{\partial w^{[2]}} , \text{ cost function } J = \frac{1}{n} \cdot \sum_{i=1}^n L^{(i)} \\
 &\quad (\text{mini-batch GD}) \\
 J_{mb} &= \frac{1}{B} \cdot \sum_{i=1}^B L^{(i)}
 \end{aligned}$$

$$\begin{cases} v_{dw^{[2]}} = \beta \cdot v_{dw^{[2]}} + (1-\beta) \cdot \frac{\partial J}{\partial w^{[2]}} \\ w^{[2]} = w^{[2]} - \alpha \cdot v_{dw^{[2]}} \end{cases}$$

3. Analyzing the Parameters.

Training Set : 96% Accuracy.

Testing Set : 64% Accuracy.

Solution :

- Collecting more data
- Employing regularization
- Making the model shallower

① L₂ Regularization.

$$\begin{aligned} J_{L_2} &= J + \frac{\lambda}{2} \|w\|^2 \\ &= J + \frac{\lambda}{2} \sum_{ij} |w_{ij}|^2 \\ &= J + \frac{\lambda}{2} w^T w. \end{aligned}$$

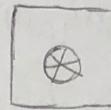
Update rule with L₂ Regularization :

$$\begin{aligned} w &= w - \alpha \cdot \frac{\partial J}{\partial w} - \alpha \cdot \frac{\lambda}{2} \cdot \frac{\partial w^T w}{\partial w} \\ &= (I - \alpha \lambda) w - \alpha \cdot \frac{\partial J}{\partial w} \end{aligned}$$

This penalization increases the cost J ,

which encourages individual parameters to be small in magnitude, which is a way to reduce overfitting.

② Parameter Sharing.



A ball may appear in any region of image and not always the center. The logistic regression will likely predict no ball.

Solution : Convolutional Neural Networks.

Suppose θ not a vector but a Matrix

$$a = \sum_{i=1}^4 \sum_{j=1}^4 \theta_{ij} x_{ij}$$

the parameter θ have "seen" all pixels of the image.

Name: Xiaofan Wang

Lee 10. Bias-Variance & Error Analysis

1. Bias-Variance Tradeoff.

$$E(x, y) \sim \text{test set} \mid \hat{f}(x) - y \mid^2$$

this MSE is too high:

(Mean Squared Error)

Overfitting: the model is too closely to training set
doesn't generalize well to other examples

Underfitting: the model didn't gather enough information
from training set. doesn't capture the
link between the feature x and target y .

The data is simply noisy: Neither overfitting nor underfitting.

The high MSE is due to the amount
of noise in the dataset.

$$y_i = f(x_i) + \varepsilon_i, \quad \varepsilon_i \sim E(\varepsilon_i) = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2.$$

$\hat{f}(x)$ is random, \Rightarrow Bias $E(\hat{f}(x) - f(x))$

compute MSE on the test set:

$$\begin{aligned}\text{Test MSE} &= E((\hat{y} - \hat{f}(x))^2) \\ &= E((\varepsilon + \hat{f}(x) - \hat{f}(x))^2) \\ &= E(\varepsilon^2) + E((\hat{f}(x) - \hat{f}(x))^2) \\ &= \sigma^2 + (E(\hat{f}(x) - \hat{f}(x)))^2 + \text{Var}(\hat{f}(x) - \hat{f}(x)) \\ &= \sigma^2 + (\text{Bias } \hat{f}(x))^2 + \text{Var } (\hat{f}(x))\end{aligned}$$

High Bias: Underfitting

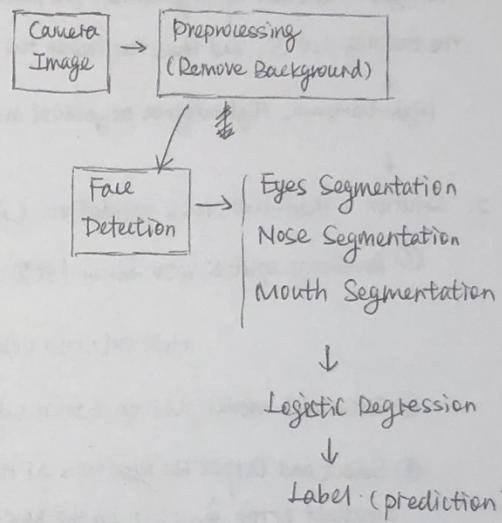
High Variance: Overfitting.

Large σ^2 : Noisy Data

Analyzing the performance of Machine Learning Algorithm:

reduce Bias without increasing the variance.
reduce Variance without increasing the Bias.
Bias-Variance Tradeoff.

2. Error Analysis



Name: Xiaofan Wang

Lee 11 Regularization Model Selection

- Cross Validation

1. Using empirical risk minimization for model selection:

- ① Train each model M_i on S_i to get some hypothesis h_i .
- ② Pick the hypothesis with the smallest training error.

The higher the order of polynomial, the better it will fit the training set S , and thus the lower the training error.

↓
High-Variance, High-degree Polynomial model.

↓

2. Solution: Hold-out cross validation (Algorithm).

- ① Randomly split S into Train (70%) and $\frac{S}{7}$ (30%)

Hold-out cross validation set.

- ② Train each model M_i on Train only, to get h_i

- ③ Select and output the hypothesis h_i that had the smallest error $\hat{\epsilon}_{cv}(h_i)$ on the hold-out cross validation set.

3. K-fold cross Validation (Hold Out less Data Each time)

- ① Randomly Split S into K disjoint subsets of $\frac{m}{K}$ training examples each. ($S_1 \dots S_K$).

- ② For each model M_i ,

For $j=1 \dots K$,

Train the Model M_i on $S_1 \cup \dots \cup S_{j-1} \cup S_{j+1} \cup \dots \cup S_K$
(train on all the data except S_j) to get some hypothesis h_{ij}

Test the hypothesis h_{ij} on S_j , to get $\hat{\epsilon}_j(h_{ij})$

- ③ pick the model M_i with the lowest estimated generation error, and retain that model on the entire training set S .

4. leave-one-out cross Validation.

Repeatedly train on all but one of the training examples $\setminus S_i$, and test on that held-out example.

The resulting $m=k$ errors are then ~~are~~ averaged together to obtain our estimate of the generalization error of a model.

- Feature Selection

1. Algorithm: Forward Search (Wrapper)

- ① Initialize $F = \emptyset$

- ② Repeat ?

- (a) For $i=1 \dots n$ if $i \notin F$, let $F_i = F \cup \{i\}$

use cross Validation to evaluate features F_i

- (b) Set F to be the best feature subset found on step(a)

}

- ③ Select and Output the best feature subset that was evaluated during the entire search procedure.

2. Algorithm: Backward Search (Filter)

Pick K features with the largest scores S_{ij}

Define S_{ij} to be the correlation between $x_i \sim y_j$

$$MI(x_i; y) = \sum_{x_i \in \{0,1\}} \sum_{y_j \in \{0,1\}} p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)}$$

S_{ij} : Mutual Information.

KL divergence:

$$MI(x_i; y) = KL(p(x_i, y) || p(x_i)p(y))$$

3. Bayesian Statistics and Regularization.

Parameter Fitting using Maximum Likelihood (ML) :

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

Bayesian View : θ random variable

prior distribution $p(\theta)$

Given a training set $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m$

Compute the posterior Distribution on the parameters:

$$\begin{aligned} p(\theta | S) &= \frac{p(S|\theta) \cdot p(\theta)}{p(S)} \\ &= \frac{\left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right) p(\theta)}{\int_{\theta} \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right) p(\theta) d\theta} \end{aligned}$$

Bayesian Logistic Regression :

$$p(y^{(i)} | x^{(i)}, \theta) = h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

$$h_{\theta}(x^{(i)}) = \frac{1}{1 + \exp(-\theta^T x^{(i)})}$$

Compute the posterior Distribution on θ :

$$p(\theta | x, S) = \int_{\theta} p(y|x, \theta) \cdot p(\theta | S) d\theta$$

$$E[y|x, S] = \int_y y p(y|x, S) dy$$

The MAP (maximum a posteriori) :

$$\theta_{MAP} = \arg \max_{\theta} \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) p(\theta)$$

$$\text{prior } p(\theta) = \theta \sim N(0, \tau^2 I)$$

Bayesian Logistic Regression turns out to be an effective algorithm for text classification, even though in text classification we usually have $n \gg m$.

Name: Xiaofan Wang.

Lee 12. K-means Clustering Algorithm.
(Unsupervised Learning problem).

Given a training set $\{x^{(1)}, \dots, x^{(n)}\}$.

Group the data into a few "clusters".

$x^{(i)} \in \mathbb{R}^d$, no labels $y^{(i)}$.

The K-means Clustering Algorithm:

1. Initializing cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d$ randomly.
2. Repeat until convergence:

For every i , set $c^{(i)} := \arg \min_j \|x^{(i)} - \mu_j\|^2$

For each j , set $\mu_j := \frac{\sum_{i=1}^n \mathbf{1}_{\{c^{(i)} = j\}} x^{(i)}}{\sum_{i=1}^n \mathbf{1}_{\{c^{(i)} = j\}}}$

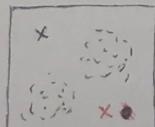
}

Inner Loop repeatedly carries out two steps:

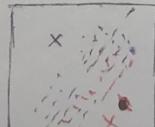
- (i) Assigning each training example $x^{(i)}$ to the ~~nearest~~ closest cluster centroid μ_j .
- (ii) Moving each cluster centroid μ_j to the mean of the points assigned to it.



(a)



(b)



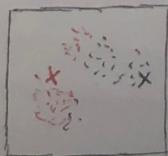
(c)

⇒

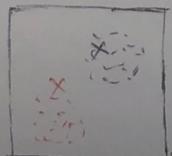
(a) Original dataset

(b) Random initial cluster centroids.

(c-f). running two iterations of K-means.



(d)



(e)



(f)

Define Distortion Function:

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

K-means is exactly coordinate descent on J .

The inner-loop of K-means repeatedly min- J with c , fixed μ . Then minimize J with μ , fixed c
⇒ J must monotonically decrease,
the value of J must converge.

Distortion Function J : non-convex function

coordinate descent on J is not guaranteed
to converge to the global minimum.

(K-means ~~are~~ ~susceptible to local optima.)

Very often K-means will work fine and
come up with very good clusterings ~

(To avoid getting stuck in bad Local optima,
running K-means many times.)

Then, out of all the different clusterings found,
pick the one that gives the lowest distortion

$$J(c, \mu)$$

Name: Xiaofan Wang

Lee 13. Mixtures of Gaussians and EM algorithm.

A joint Distribution: $p(x^{(i)}, z^{(i)}) = p(x^{(i)}|z^{(i)}) p(z^{(i)})$

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

$x^{(i)}|z^{(i)} = j \sim N(\mu_j, \Sigma_j)$

↓

$$\begin{aligned} l(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{j=1}^k p(x^{(i)}|z^{(i)}; \mu_j, \Sigma_j) p(z^{(i)}; \phi) \end{aligned}$$

$$l(\phi, \mu, \Sigma) = \sum_{i=1}^n \log p(x^{(i)}|z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi)$$

Maximizing this with respect to ϕ, μ, Σ gives the parameters =

$$\phi_j = \frac{1}{n} \cdot \sum_{i=1}^n I\{z^{(i)} = j\}$$

$$\mu_j = \frac{\sum_{i=1}^n I\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^n I\{z^{(i)} = j\}}$$

$$\Sigma_j = \frac{\sum_{i=1}^n I\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n I\{z^{(i)} = j\}}$$

In the E-step, calculate the posterior probability of our parameters the $z^{(i)}$'s.

Using Bayes Rule :

$$p(z^{(i)}_j | x^{(i)}; \phi, \mu, \Sigma) = \frac{p(x^{(i)}|z^{(i)}=j; \mu_j, \Sigma_j) p(z^{(i)}_j; \phi)}{\sum_{l=1}^k p(x^{(i)}|z^{(i)}=l; \mu_l, \Sigma_l) p(z^{(i)}_l; \phi)}$$

In the M-step. $z^{(i)}$: indicating from which Gaussian each datapoint had come, $\sim w_j^{(i)}$.

The EM Algorithm is also reminiscent of the k-means clustering algorithm, except that instead of the "hard" cluster assignments $C^{(i)}$, we instead have the "soft" assignments $w_j^{(i)}$. Similar to K-means, it is also susceptible to local optima, so reinitializing at several different initial parameters may be a good idea.

=. EM Algorithm. (Iterative algorithm. ~ two main steps)

1. In the E-step : try to guess the value of the $z^{(i)}$'s.
2. In the M-step : update the parameters of our model.

Repeat until convergence: }

(E-step) For each i, j , set

$$w_j^{(i)} := p(z^{(i)}_j | x^{(i)}; \phi, \mu, \Sigma)$$

(M-step) Update the parameters:

$$\phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$$

$$\Sigma_j = \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$