Hw1_writeup
Sunday, October 9, 2016   8:02 PM

1°
It is almost same as gradient descent but only takes small batches to estimate using a batch from the input to estimate the cost function.
taking the random batches from the Training set.

Pseudo code:

```
for i = 0 ..... N do:
        initialize  d_i randomly
        let t = 1
    for x in Training-batches:
            if (t ≤ T and stopping conditions not True):
                A_i ← xᵀx - Σ_{j=0}^{i-1} A_j d_j d_jᵀ  and
                y ← d_i - η ∇_{d_i}(-d_iᵀ A_i d_i)
                d_i ← y/‖y‖
                t ← t+1
        Else: return
    A_i ← d_iᵀ xᵀ x d_i
```

Here $A_i$ is estimated using the batches of training sets rather than the whole learning set.
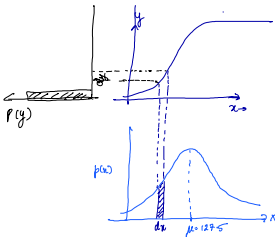
THEORETICAL:-

**2.d (i)**

$$p(x) = \frac{f(x)}{\int_0^{255} f(x)\,dx} \simeq f(x)$$

where $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$\sigma^2 = 1$ , $\mu = 127.5$



For any given mapping $y = g(x)$ between $x$ and $y$ we would have:

$$\boxed{p(y)\,dy = p(x)\,dx}$$

where - $0 \leq x \leq 255$

$$\int_0^{255} p(x)\,dx = 1$$

$$p(y) = \frac{1}{255}$$

and

$$y = g(x)$$
$$dy = d(g(x))$$

$\Rightarrow \quad d(g(x))_x = 255\, p(x)\,dx$

$\Rightarrow \quad \int_0^x d(g(x)) = 255 \int_0^x p(x)\,dx$

$\Rightarrow \quad g(x) = 255\, f(x)$    where $F(x)$ is Cummulative distⁿ fⁿ of X.

So $\boxed{y = g(x) = 255\, F(x)}$

**(ii)**

$$P(X=x, Y=y, Z=z) = \begin{cases} 8xyz & \text{for } x,y,z \in [0,1] \\ 0 & o/w \end{cases}$$

$$P(X=x) = \int_0^1 \int_0^1 P(X=x, Y=y, Z=z)\,dy\,dz$$

$$= \int_0^1 \int_0^1 8xyz \, dy\,dz$$

$$= \int_0^1 \left[\frac{8xyz^2}{2}\right]_0^1 dy$$

$$= \int_0^1 4xy \, dy$$

$$= \left[\frac{4xy^2}{2}\right]_0^1$$

$$= 2x$$

**lly** $P(Y=y) = \int_0^1 \int_0^1 P(X=x, Y=y, Z=z)\,dx\,dz$

$$= \int_0^\infty 4xy \, dx$$

$$= 2y$$

and $P(Z=z) = \int_0^1 \int_0^1 P(x=x, y=y, z=z) \, dy \, dx$

$$= \int_0^1 4zy \, dy$$

$$= 2z$$

$T = XYZ$

$T \in [0,1]$

$E(T) = \int_0^1 t \, P(T=t) \, dt$

$P(T=t) = \iint g(x=x, y=y, z=\frac{t}{xy}) \, dx \, dy$

$$= \int_0^1 \int_0^1 8xy \left(\frac{t}{xy}\right) dx \, dy$$

$$= 8t$$

$\Rightarrow \quad E[T] = \int_0^1 t(8t) \, dt$

$$= \left[ 8 \frac{t^3}{3} \right]_0^1$$

$$= \frac{8}{3}$$

$P(x=x, y=y \mid z=z_0) = \dfrac{P(x=x, y=y, z=z_0)}{P(z=z_0)}$   (applying Bayes rule).

$$= \frac{8 \overset{4}{x} y z_0}{2 z_0} = 4xy$$

$P(x=x \mid z=z_0) = \int_0^1 \dfrac{P(x=x, y=y, z=z_0) \, dy}{2 z_0}$

$$= \frac{8 x z_0 \int_0^1 y \, dy}{2 z_0}$$

$$= 2x$$

"y $P(y=y \mid z=z_0) \int_0^1 \dfrac{P(x=x, y=y, z=z_0) \, dx}{2 z_0}$

$$= 2y$$

$\Rightarrow \quad P(x=x, y=y \mid z=z_0) = P(x=x \mid z=z_0) \cdot P(y=y \mid z=z_0)$

$\therefore$ They are Conditionally independent.

$\underline{\underline{e}}$ (1) $x \sim N(\mu, \Sigma) = \dfrac{1}{(2\pi)^{\frac{n}{2}}} \dfrac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(\frac{-1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$

$b(\mu \mid \mu, \Sigma) = 1 \quad \frac{1}{\quad} \exp(-\frac{1}{2} (\mu-\mu)^T \Sigma^{-1}(\mu-\mu))$

$$p(\cdots | \mu_0, \Sigma_0) \quad (2\pi)^{n/2} \, |\Sigma_0|^{1/2} \quad \Gamma\left(-\tfrac{1}{2}(\mu-\mu_0)\Sigma_0^{-1}(\mu-\mu_0)\right)$$

$$\hat{\mu}_{MAP} = \arg\max_{\mu} f(\mu \mid X, \Sigma)$$

where $f(\mu \mid X, \Sigma)$ is posterior distribution p.d.f.

where $X = \{x^{i}\} \quad i = 1(1)m.$

$$f(\mu \mid X, \Sigma) = \frac{f(\mu, X \mid \Sigma)}{f(X \mid \Sigma)}$$

$$= \frac{f(X \mid \Sigma, \mu)\, p(\mu \mid \Sigma)}{f(X \mid \Sigma)}$$

$$= \frac{K \cdot \exp\left(-\tfrac{1}{2}\left[\sum_{i=1}^{m}(\mu-x^{(i)})^{T}\Sigma^{-1}(\mu-x^{(i)}) + (\mu-\mu_0)^{T}\Sigma_0^{-1}(\mu-\mu_0)\right]\right)}{\int_{\mu} K \exp\left(-\tfrac{1}{2}\left[\sum_{i}(\mu-x^{(i)})^{T}\Sigma^{-1}(\mu-x^{(i)}) + (\mu-\mu_0)^{T}\Sigma_0^{-1}(\mu-\mu_0)\right]\right)d\mu} \quad \text{—①}$$

where $K = \dfrac{1}{(2\pi)^{\frac{(m+1)}{2}n}} \dfrac{1}{|\Sigma_0|^{n/2}} \dfrac{1}{|\Sigma|^{\frac{mn}{2}}}$

writing expression ① as ②
with $\mu_n, \Sigma_n$ as $f^n$
of $m, \Sigma, \Sigma_0, \mu_0$

$$= \frac{K'(m,\Sigma_0,\Sigma,X,\mu_0)\dfrac{1}{(2\pi)^{n/2}}\dfrac{1}{|\Sigma_n|^{1/2}}\exp\left(-\tfrac{1}{2}(\mu-\mu_n)^{T}\Sigma_n^{-1}(\mu-\mu_n)\right)}{\int_{\mu} K'(m,\Sigma_0,\Sigma,X,\mu_0)\dfrac{1}{(2\pi)^{n/2}}\dfrac{1}{|\Sigma_n|^{1/2}}\exp\left(-\tfrac{1}{2}(\mu-\mu_n)^{T}\Sigma_n^{-1}(\mu-\mu_n)\right)d\mu} \quad \text{—②}$$

$$\underbrace{\text{multivariate normal distribution}}$$

$$= \frac{K' \dfrac{1}{(2\pi)^{n/2}}\dfrac{1}{|\Sigma_n|^{1/2}}\exp\left(-\tfrac{1}{2}(\mu-\mu_n)^{T}\Sigma_n^{-1}(\mu-\mu_n)\right)}{K'}$$

$$\Rightarrow \quad f(\mu \mid X, \Sigma) \sim N(\mu_n, \Sigma_n)$$

and the maximum value of $f(\mu \mid X, \Sigma)$

will be at $\mu = \mu_n.$

- $\mu_n$ and $\Sigma_n$ can be found by comparing the coefficients $\cdot$ and $\mu^T \Sigma_n \mu$ in following expressions

$$(\mu - \mu_n)^T \Sigma_n^{-1} (\mu - \mu_n), \quad \sum_{i=1}^{m} (\mu - x^{(i)})^T \Sigma^{-1} (\mu - x^{(i)}) + (\mu - \mu_0)^T \Sigma_0^{-1} (\mu - \mu_0)$$

- $$\mu^T \Sigma_n^{-1} \mu = \mu^T \left( m \Sigma^{-1} + \Sigma_0^{-1} \right) \mu$$

$$\Sigma_n = \left( m \Sigma^{-1} + \Sigma_0^{-1} \right)^{-1}$$

comparing coefficients of $\mu^T$

- $$\mu_n = \Sigma_n \left( m \Sigma^{-1} \bar{x} + \Sigma_0^{-1} \mu_0 \right)$$

$$\hat{\mu}_{MAP} = \left( m \Sigma^{-1} + \Sigma_0^{-1} \right)^{-1} \left( m \Sigma^{-1} \bar{x} + \Sigma_0^{-1} \mu_0 \right)$$

1) y.        assuming $\Sigma$ has no prior distribution.

$$f(\Sigma \mid x, \mu) = \frac{f(x, \Sigma \mid \mu)}{f(x \mid \mu)}$$

$$\hat{\Sigma}_{MAP} = \underset{\Sigma}{arg\,max} \left( f(\Sigma \mid x, \mu) \right)$$

where $f$ is p.d.f. of $\Sigma$.

$$f(\Sigma \mid x, \mu) = \frac{f(x \mid \Sigma, \mu)}{f(x \mid \mu)}$$

$$\Rightarrow f(\Sigma \mid x, \mu) = \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \frac{exp\left(-\frac{1}{2} \sum_i (x^{(i)} - \mu) \Sigma^{-1} (x^{(i)} - \mu)\right)}{f(x \mid \mu) \to \text{not a } f^n \text{ of } \Sigma.}$$

we can find the $\Sigma'$ s.t. $f(\Sigma' \mid x, \mu)$ is

max.    and    $\boxed{\hat{\Sigma}_{MAP} = \Sigma'}$

Since log is an monotonously increasing

function. $\therefore$ finding $\underset{\Sigma}{arg\,max} f(\Sigma \mid x, \mu)$

is same as    finding $\underset{\Sigma}{arg\,max} \log f(\Sigma \mid x, \mu)$

$\therefore$ MLE and MAP estimator for $\Sigma$

is Same.

and since $f(x \mid \mu)$ not a $f^n$ of $\Sigma$

then $\underset{\Sigma}{arg\,max} f(\Sigma \mid x, \mu) = \underset{\Sigma}{arg\,max}(\log f(\Sigma \mid x, \mu))$

$$= \underset{\Sigma}{arg\,max} \left[ \log f(x \mid \Sigma, \mu) \right]$$

$\therefore$ MLE and MAP of $\Sigma$ will be.

$$\hat{\Sigma}_{MAP} = \hat{\Sigma}_{MLE} = \underset{\Sigma}{arg\,max} \left( \log f(x \mid \Sigma, \mu) \right)$$

$$\hat{\Sigma}_{MAP} = \hat{\Sigma}_{MLE} = \underset{\Sigma}{arg\,max}\left(\log f(X|\Sigma,\mu)\right)$$

$f(X|\Sigma,\mu)$

$$= \frac{\prod\limits_{i=1}^{m} exp\left(-\frac{1}{2}(x^i-\mu)^T \Sigma^{-1}(x^i-\mu)\right)}{(2\pi)^{\frac{nm}{2}}|\Sigma|^{\frac{1}{2}}}$$

$$= \frac{1}{(2\pi)^{\frac{mn}{2}}|\Sigma|^{\frac{m}{2}}} exp\left(-\frac{1}{2}\sum_q (x^i-\mu)\Sigma^{-1}(x^i-\mu)\right) \quad -\textcircled{1}$$

now $\qquad A = \sum\limits_{i=1}^{m}(x^i-\mu)(x^i-\mu)^T$

$$\bar{x} = \frac{1}{m}\sum\limits_{i=1}^{m} x^i$$

as $\qquad \underbrace{\sum\limits_{i=1}^{m}(x^i-\mu)^T \Sigma^{-1}(x^i-\mu)}$

Scalar

$$= \sum\limits_{i=1}^{m}(x^i-\bar{x}+\bar{x}-\mu)^T \Sigma^{-1}(x^i-\bar{x}+\bar{x}-\mu)$$

$$= \underbrace{\left[\sum\limits_{i=1}^{m}(x^i-\bar{x})^T \Sigma^{-1}(x^i-\bar{x})\right]} + m(\bar{x}-\mu)^T \Sigma^{-1}(\bar{x}-\mu)$$

Scalar ( can take trace)

$$\Rightarrow tr\left(\Sigma^{-1}\sum\limits_{i=1}^{m}(x^i-\bar{x})(x^i-\bar{x})^T\right) + m(\bar{x}-\mu)^T \Sigma^{-1}(\bar{x}-\mu)$$

$$\left(tr(AB)=tr(BA)\right)$$

$\therefore$ ① becomes

$$f(X|\Sigma;\mu) = 2\pi^{-\frac{mn}{2}} |\Sigma^{-1}|^{\frac{m}{2}} \exp\left(-\frac{1}{2} tr(\Sigma^{-1}A) - \frac{m}{2}(\bar{x}-\mu)^T \Sigma^{-1}(\bar{x}-\mu\right.$$

$$\hat{\mu}_{MLE} = \bar{x}$$

solving for MLE of $\Sigma$ assuming that $\mu$ is $\hat{\mu}_{MLE}$

$$\log f(x|\Sigma,\mu))$$

$$= -\frac{mn}{2}\log 2\pi + \frac{m}{2}\log|\Sigma^{-1}| - \frac{1}{2}tr(\Sigma^{-1}A) - \underbrace{\frac{m}{2}(\bar{x}-\mu)^T \Sigma^{-1}(\bar{x}-\mu)}_{=0}$$

$$\Downarrow$$

$$\underbrace{\frac{m}{2}\log|\Sigma^{-1}A| - \frac{1}{2}tr\Sigma^{-1}A}_{②} \underbrace{- \frac{m}{2}\log|A|}_{\text{to maximise we can ignore this.}} -$$

$\partial_1, \partial_2 --- \partial_R$ be eigen values of $\Sigma^{-1}A$

So, ② $\Rightarrow$

$$\frac{m}{2}\log\left(\prod_{i=1}^{R}\partial_i\right) - \frac{1}{2}\left(\sum_{R=1}^{R}\partial_i\right) \qquad \left[\begin{array}{l}\text{This is maximised when} \\ \text{each } \partial_i = m.\end{array}\right]$$

matrix of eigenvectors

$$\therefore \Sigma^{-1}A = P(mI)P^T \longrightarrow P = [e_1\ e_2 --- e_k]$$

$$\Rightarrow \qquad \Sigma^{-1} = mA^{-1}$$

$$\Sigma = \frac{1}{m}A$$

$$\Rightarrow \qquad \hat{\Sigma}_{MAP} = \hat{\Sigma}_{MLE} = \frac{1}{m}\sum_{i=1}^{m}(x^{(i)}-\bar{x})(x^i-\bar{x})^T$$

Monday, October 10, 2016   6:47 PM

$$\hat{\mu}_{MAP} = (m\Sigma^{-1} + \Sigma_0^{-1})^{-1}(m\Sigma^{-1}\bar{x} + \Sigma_0^{-1}\mu_0)$$

$$\hat{\Sigma}_{MAP} = \frac{1}{m}\sum_{i=1}^{m}(x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T$$

## checking whether biased or not

(1) $\hat{\mu}_{MAP} = \left(m\Sigma^{-1} + \Sigma_0^{-1}\right)^{-1}(m\Sigma^{-1}\bar{x} + \Sigma_0^{-1}\mu_0)$

$$E(\hat{\mu}_{MAP}) = \left(m\Sigma^{-1} + \Sigma_0^{-1}\right)^{-1}E\left(m\Sigma^{-1}\bar{x} + \Sigma_0^{-1}\mu_0\right)$$

$$= \left(m\Sigma^{-1} + \Sigma_0^{-1}\right)^{-1}\left(m\Sigma^{-1}E(\bar{x}) + \Sigma_0^{-1}\mu_0\right)$$

as $E()$ is linear $\beta^m$ we can take inside bracket and $\Sigma_0^{-1}\mu_0$ is constant
∴ $E(\Sigma_0^{-1}\mu_0) = \Sigma_0^{-1}\mu_0$.

$$= \left(m\Sigma^{-1} + \Sigma_0^{-1}\right)^{-1}\left(m\Sigma^{-1}\mu + \Sigma_0^{-1}\mu_0\right)$$

∴ $E(\hat{\mu}_{MAP}) \neq \mu$
    estimate is biased

(2) $\hat{\Sigma}_{MAP} = \frac{1}{m}\sum_{i=1}^{m}\left[x^i x^{iT} - x^i\bar{x}^T - \bar{x}x^{iT} + \bar{x}\bar{x}^T\right]$

$$= \frac{1}{m}\sum_{i=1}^{m}\left[x^i x^{iT}\right] - \left(\frac{1}{m}\sum_{i=1}^{m}x^{(i)}\right)\bar{x}^T - \bar{x}\frac{1}{m}\sum_{i=1}^{m}x^{iT} + \bar{x}\bar{x}^T$$

$$= \frac{1}{m}\sum_{i=1}^{m}x^i x^{iT} - 2\bar{x}\bar{x}^T + \bar{x}\bar{x}^T$$

$$\hat{\Sigma}_{MAP} = \frac{1}{m}\sum_{i=1}^{m}x^i x^{iT} - \bar{x}\bar{x}^T \quad \text{———} \textcircled{1}$$

$$x^i \sim N(\mu, \Sigma)$$
$$\bar{x} \sim N\left(\mu, \frac{1}{m}\Sigma\right)$$

as $\bar{x} = \frac{1}{m}\sum x^i$

$$E(\bar{x}) = \frac{1}{m} \sum_{i=1} E(x^i) = \frac{m}{m}\mu = \mu$$

$$Cov(\bar{x}) = E(\bar{x}-\mu)(\bar{x}-\mu)^T$$

$$= E\left(\sum \frac{x^i}{m}-\mu\right)\left(\sum \frac{x^i}{m}-\mu\right)^T$$

$$= \frac{1}{m^2} E\left(\sum_{i=1}^{m} x^i - m\mu\right)\left(\sum_{i=1}^{m} x^i - m\mu\right)^T$$

$$= \frac{1}{m^2} E\left(\sum_{i=1}^{m}(x^i-\mu)\right)\left(\sum_{i=1}^{m}(x^i-\mu)\right)^T$$

$$= \frac{1}{m^2}\left[\sum_{i=1}^{m} E(x^i-\mu)(x^i-\mu)^T + \underbrace{\sum_{i}\sum_{j} E(x^i-\mu)(x^j-\mu)^T}_{=0}\right]$$

since $x^i$'s are independent

$$= \frac{1}{m^2} m \Sigma$$

$$= \frac{1}{m} \Sigma$$

now, from ①

$$E(\hat{\Sigma}_{MAP}) = E\left(\frac{1}{m}\sum_{i=1}^{m} x^i x^{iT} - \bar{x}\bar{x}^T\right)$$

$$= \frac{1}{m}\sum_{i=1}^{m} E(x^i x^{iT}) - E(\bar{x}\bar{x}^T)$$

(using linearity of $E()$ $f^n$)

$$= \frac{1}{m}\sum_{i=1}^{m}(\Sigma + \mu\mu^T) - \left(\frac{1}{m}\Sigma + \mu\mu^T\right)$$

$$= \Sigma + \mu\mu^T - \frac{1}{m}\Sigma - \mu\mu^T$$

$$= \frac{m-1}{m}\Sigma$$

∴ $\hat{\Sigma}_{MAP}$ is biased estimator

so is $\hat{\Sigma}_{MLE}$