

Tuan Tran

A20357888

CS 584- Summer 2020

* Numerical Solutions: (NS)

- 1, Numerical Solutions are necessary ^{and beneficial} when ~~the opt~~:
 - + The optimization problem ~~does~~ doesn't have an exact closed form solution (ex: non-linear gradient ~~at this point~~)
 - + Matrix is too large to invert \rightarrow Computationally expensive

2,

- Gradient Descent is the algo where we follow the gradient to decrease the ~~loss~~ curve and arrive at the point θ^* where our Objective Function is minimized.

This is possible since the property of gradient is that the gradient of a function points in the direction of max change

\rightarrow negative gradient is the direction of max decrease.

\hookrightarrow We can update our params:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla J(\theta^{(i)})$$

- Stop Condition is when the value of loss function doesn't change past some small threshold: $(J(\theta^{(i+1)}) - J(\theta^{(i)})) < \epsilon$

3,

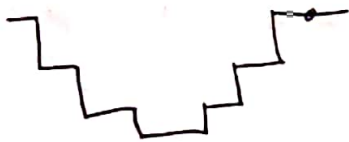
- learning rate is the size of the step we ~~take~~ take along direction of gradient
- learning rate is a hyperparameter \rightarrow we have to experiment to set it correctly

(2)

- When learning rate is small, we are likely to hit minimum point but ~~converge~~ ^{learning is} slower since we are taking small step size
- When ~~learn~~ Lr is large, we are taking larger step size and thus will descent faster but we might overshoot the minimum point

4,

- GD can't be applied to piecewise constant functions since it's easy to hit a point where grad is 0:



- GD is guaranteed to produce global min if objective function is convex
- We can attempt to avoid local minimum by ~~to~~ introducing momentum
- Idea is to use residual velocity from previous descent to move over local minimum.

5,

Loss func for Linear Regression:

$$J(\theta) = \frac{1}{2} (Z\theta - Y)^T (Z\theta - Y)$$

$$\rightarrow \nabla J(\theta) = Z^T (Z\theta - Y)$$

→ Update Equation:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta Z^T (Z\theta - Y)$$

→ Summation form:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \sum_{j=1}^m (z^{(j)T} \theta - y^{(j)}) z^{(j)}$$

(3)

6,

Basic GD computes gradient over all training examples

SGD computes gradient over a batch of examples only

↳ Basic GD:

+, Accurate, stable

+, Computationally expensive

SGD:

+, Converge faster but not as accurate and stable

+, less computationally expensive.

7,

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta (z^{(i)} \theta - y^{(i)}) z^{(i)}$$

$$\leftarrow \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} - \frac{1}{6} \cdot \left(\begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} - 2 \right) \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\leftarrow \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} - \frac{1}{6} (6) \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\leftarrow \cancel{\frac{2}{4}} \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\rightarrow \theta^{(i+1)} = \begin{bmatrix} 2 \\ 4 \\ 4 \end{bmatrix}$$

8,

Newton's method update equation:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - H^{-1} \nabla J(\theta^{(i)})$$

H^{-1} = inverse Hessian ~~is~~ is equivalent to the learning rate in GD

- Advantage of Newton's method:

+ Can compute learning rate instead of specifying it as hyper param.

+ We get different learning rate for different features

- Disadvantage:

+ Hessian computation is expensive and may be noisy

+ Sometimes cannot invert H .

9,

$$J(\theta) = \frac{1}{2} (Z\theta - Y)^T (Z\theta - Y)$$

$$\rightarrow \nabla J(\theta) = Z^T (Z\theta - Y)$$

$$\rightarrow H = \nabla (\nabla J(\theta)) = Z^T Z$$

$$\rightarrow \theta^{(i+1)} \leftarrow \theta^{(i)} - (Z^T Z)^{-1} Z^T (Z\theta - Y)$$

10,

The line search first finds the descent direction using GD and

then computes a stepsize that determines how ^{far} ~~much~~ to move along that direction and iteratively.

★ GDA:

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

- Σ gives the variance along the diagonal and covariance between features ~~at otherwise~~

(5)

- When features are uncorrelated, Σ becomes a diagonal matrix since covariance between 2 uncorrelated variables = 0

$$\rightarrow \Sigma = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{bmatrix}$$

2,

- ~~Euclidean~~ Euclidean assumes all directions are equally important

↳ Mahalanobis is basically Euclidean normalized by inverse of Σ

→ It gives different weights ($\frac{1}{\text{variance}}$) to different direction when computing the distance. It accounts for the variance (uncertainty) in different directions.

~~Euclidean~~

3,

$$x = \begin{bmatrix} 0 \\ 1 \end{bmatrix} : \mu = \begin{bmatrix} 2 \\ 3 \end{bmatrix} : \Sigma = \begin{bmatrix} 4 & 0 \\ 0 & 5 \end{bmatrix} \rightarrow \Sigma^{-1} = \begin{bmatrix} 1/4 & 0 \\ 0 & 1/5 \end{bmatrix}$$

$$\rightarrow \text{Mahalanobis dist} = \left(\begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 \\ 3 \end{bmatrix} \right)^T \cdot \begin{bmatrix} 1/4 & 0 \\ 0 & 1/5 \end{bmatrix} \cdot \left(\begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 \\ 3 \end{bmatrix} \right)$$

$$= \begin{bmatrix} -2 & -2 \end{bmatrix} \begin{bmatrix} 1/4 & 0 \\ 0 & 1/5 \end{bmatrix} \begin{bmatrix} -2 \\ -2 \end{bmatrix}$$

$$= \begin{bmatrix} -1/2 & -2/5 \end{bmatrix} \begin{bmatrix} -2 \\ -2 \end{bmatrix} = 1 + \frac{4}{5} = \frac{9}{5}$$

4,

In generative learning, we model the feature distribution in each class $P(X|Y)$ and class prior $P(Y)$

→ We can generate more examples given the class wanted ~~and posterior~~
~~posterior~~

(6)

~~but~~ We have posterior prob for each class:

$$p(y=k|x) = \frac{P(x|y=k) P(y=k)}{\sum_i P(x|y=i) P(y=i)}$$

→ Since when comparing $P(y=k|x)$ for different values of k , we only compare the numerator since as denominator is constant

→ We can define discriminant function as log of numerator.

5,

- Bayes rule helps ~~find~~^{describe} the poste prob of an event based on prior knowledge
 of conditions that may be ~~not~~ related to the event.

→ In our case:

$$P(y=j|x) = \frac{P(x|y=j) P(y=j)}{\sum_{i=1}^K P(x|y=i) P(y=i)}$$

where $P(y=j|x)$ = posterior

$P(x|y=j)$ = likelihood

$P(y=j)$ = prior.

- Evidence can be ignored since it is constant ~~for all~~^{For} $P(y=i|x)$ for all $i=1 \dots K$.

6,

For class k :

$$\mu_k = \frac{1}{m_k} \sum_{i=1}^{m_k} x^{(i)} 1(y^{(i)}=k)$$

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

$$p(y=k) = \frac{1}{m} \sum_{i=1}^m 1(y^{(i)}=k)$$

7,

- ~~the~~ behave for class k :

$$P(y=k|x) = \frac{P(x|y=k)P(y=k)}{\sum_{i \in K} P(x|y=i)P(y=i)}$$

-> denominator doesn't change

-> discriminant func: $g_k(x) = \log(P(x|y=k)P(y=k))$

- ~~Since~~ Since we choose gaussian

$$\rightarrow P(x|y=k) = \frac{1}{\sqrt{2\pi}\sigma_k} \cdot e^{-\frac{1}{2\sigma_k^2}(x-\mu_k)^2}$$

$$\begin{aligned} \rightarrow g_k(x) &= \log\left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2\sigma_k^2}(x-\mu_k)^2}\right) + \log(P(y=k)) \\ &= \log\left(\frac{1}{\sqrt{2\pi}}\right) - \log(\sigma_k) - \frac{1}{2\sigma_k^2}(x-\mu_k)^2 + \log(P(y=k)) \end{aligned}$$

\hookrightarrow constant

$$\rightarrow g_k(x) = -\log(\sigma_k) - \frac{1}{2\sigma_k^2}(x-\mu_k)^2 + \log(P(y=k))$$

- Equal priors and variances

$$\begin{aligned} \rightarrow \sigma_k &\rightarrow \sigma \\ P(y=k) &= \pi_k \end{aligned} \quad \left. \vphantom{\begin{aligned} \rightarrow \sigma_k &\rightarrow \sigma \\ P(y=k) &= \pi_k \end{aligned}} \right\} \rightarrow \text{constant}$$

$$\rightarrow g_k(x) = -\frac{1}{2\sigma^2}(x-\mu_k)^2$$

\hookrightarrow identical scale value for all classes

$$\rightarrow g_k(x) = (x-\mu_k)^2$$

\rightarrow Simplifies down to choosing class with closest mean

8,

- Model for class k :

$$P(X|y=k) = \frac{1}{(2\pi)^{n/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2} (x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)} ; d_k = P(y=k)$$

Similar to before, we have $g_k(x) = \log(P(X|y=k)) + \log d_k$.

$$\rightarrow g_k(x) = \underbrace{-\log(2\pi)^{n/2}}_{\text{Constant}} - \frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k) + \log d_k$$

$$\rightarrow g_k(x) = -\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k) + \log d_k$$

- Equal covariance and priors:

\rightarrow Similar to ~~before~~ ^{question 7}, we can simplify down to:

$$g_k(x) = \|x - \mu_k\|^2$$

\rightarrow euclidean distance to mean.

9,

- IID assumption means that ~~instead~~ instances are independent and follow the same distribution.

- ~~Model for k class:~~

$$p(x|y=k) =$$

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

$$\text{where } P(x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \frac{(x-\mu_k)^2}{\sigma^2}}$$

$$\rightarrow \ell(\theta) = \sum_{i=1}^m \log P(x^{(i)}; \theta)$$

$$\rightarrow \ell(\theta_k) = \sum_{i=1}^{m_k} \log \left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \frac{(x^{(i)} - \mu_k)^2}{\sigma_k^2}} \right)$$

$$\rightarrow \ell(\theta_k) = m \log \frac{1}{\sqrt{2\pi}} - m \log \sigma_k - \frac{1}{2\sigma_k^2} \sum_{i=1}^{m_k} (x^{(i)} - \mu_k)^2$$

$$\nabla \ell(\theta_k) = 0 \rightarrow \begin{bmatrix} \frac{\partial \ell}{\partial \mu_k} \\ \frac{\partial \ell}{\partial \sigma_k} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

+, For μ_k :

$$\frac{\partial \ell}{\partial \mu_k} = -\frac{1}{2\sigma_k^2} \sum_{i=1}^{m_k} 2(x^{(i)} - \mu_k)(-1)$$

$$= 0$$

$$\rightarrow \sum_{i=1}^{m_k} (x^{(i)} - \mu_k) = 0$$

$$\rightarrow \sum_{i=1}^{m_k} x^{(i)} = m_k \mu_k$$

$$\rightarrow \mu_k = \frac{1}{m_k} \sum_{i=1}^{m_k} x^{(i)}$$

+, For σ_k :

$$\frac{\partial \ell}{\partial \sigma_k} = -m \frac{1}{\sigma_k} - \frac{1}{2} (-2) \frac{1}{\sigma_k^3} \sum_{i=1}^{m_k} (x^{(i)} - \mu_k)^2 = 0$$

$$\rightarrow \sigma_k^2 = \frac{1}{m_k} \sum_{i=1}^{m_k} (x^{(i)} - \mu_k)^2$$

⊛ Naive Bayes:

1, Naive Bayes assumption is that we assume the features are independent

$$\rightarrow P(x) = P(x_1 \dots x_n) = \prod_{i=1}^n P(x_i)$$

This assumption is that such assumption leads to a simpler model since it reduces number of parameters to learn

2,

For Bernoulli NB:

$$P(x_j^{(i)} | y=1) = \alpha_{j|y=1}^{x_j^{(i)}} (1 - \alpha_{j|y=1})^{1-x_j^{(i)}}$$

$\alpha_{j|y=1}$ = prob that word j is in document i of class y=1

3,

For class y=1:

$$- \ell(\theta) = \log P(x^1 \dots x^{m_1}; \theta)$$

$$= \log \prod_{i=1}^{m_1} P(x^i; \theta)$$

$$= \log \prod_{i=1}^{m_1} P(x_1^i, x_2^i, \dots, x_n^i; \theta)$$

$$= \log \prod_{i=1}^{m_1} \prod_{j=1}^n P(x_j^{(i)}; \theta)$$

↳ NB assumption

$$\rightarrow \ell(\theta) = \sum_{i=1}^{m_1} \sum_{j=1}^n \log P(x_j^{(i)}; \theta)$$

$$= \sum_{i=1}^{m_1} \sum_{j=1}^n x_j^{(i)} \log(\alpha_{j|y=1}) + (1 - x_j^{(i)}) \log(1 - \alpha_{j|y=1})$$

$$- \frac{\partial \ell}{\partial \alpha_{j|y=1}} = \sum_{i=1}^{m_1} x_j^{(i)} \cdot \frac{1}{\alpha_{j|y=1}} + (1 - x_j^{(i)}) \cdot \frac{1}{1 - \alpha_{j|y=1}} \cdot (-1) = 0$$

$$\rightarrow \frac{1}{\alpha_{j|y=1}} \sum_{i=1}^{m_1} x_j^{(i)} = \frac{1}{1 - \alpha_{j|y=1}} \sum_{i=1}^{m_1} (1 - x_j^{(i)})$$

$$\rightarrow \boxed{\alpha_{j|y=1} = \frac{1}{m_1} \sum_{i=1}^{m_1} x_j^{(i)}}$$

4,

we define $g_k(x)$ for class $k = \prod_{j=1}^n \alpha_{j|y=k}^{x_j} (1 - \alpha_{j|y=k})^{1-x_j}$

$$\rightarrow \hat{y} = \arg \max_k g_k(x)$$

5,

Bernoulli doesn't give us useful info such as frequency of words.

→ With Binomial, we can get each feature as frequency x_j : how many times the word occurs in the document.

6,

~~First we~~

- ① Select binomial ~~model~~ model for $P(X|y)$
- ② Use MLE or MAP to find parameters of the binomial model
- ③ Derive discriminant function
- ④ $\hat{y} = \underset{k}{\operatorname{argmax}} g_k(x)$.

7,

- For class ~~$y = k$~~ we have following model log-likelihood:

~~$P(x|y)$~~

$$\begin{aligned} l(\theta) &= \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}) P(x^{(i)}; \theta) \\ &= \log \prod_{i=1}^m \prod_{j=1}^n P(x_j^{(i)}; \theta) \\ &= \sum_{i=1}^m \sum_{j=1}^n \log \binom{d^{(i)}}{x_j^{(i)}} \theta_{j|y=y^{(i)}}^{x_j^{(i)}} (1 - \theta_{j|y=y^{(i)}})^{(d^{(i)} - x_j^{(i)})} \end{aligned}$$

$$\rightarrow \frac{\partial l}{\partial \theta_{j|y=k}} = 0$$

$$\begin{aligned} \rightarrow \frac{\partial l}{\partial \theta_{j|y=k}} &= \frac{\sum_{i=1}^m 1(y^{(i)} = k) x_j^{(i)} + \epsilon}{\sum_{i=1}^m 1(y^{(i)} = k) d^{(i)} + k \epsilon} \end{aligned} \quad \left. \begin{array}{l} \text{Laplace smoothing} \end{array} \right\}$$

- Laplace smoothing is in place to prevent cases when a word j may not appear in any of the documents $\rightarrow P(x|y) \neq 0$

membership func used during classification for class k :

$$g_k(x) = \sum_{j=1}^n \left[\log \left(\frac{d_j}{x_j} \right) + x_j \log(d_{j|y=k}) + (1-x_j) \log(1-d_{j|y=k}) \right] + \log d_k$$

8,

- Due with gaussian feature vectors:

$$g_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log d_k.$$

$$\text{where } \mu_k = \frac{1}{m_k} \sum_{i=1}^m 1(y^{(i)}=k) \cdot x^{(i)}$$

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

- Due to NB assumption, Σ_k becomes a diagonal matrix

Since cov of 2 uncorrelated variables = 0.

→ instead of n^2 params in Σ_k , we only have n params.

9,

- Precision = $\frac{\text{True Positive}}{\text{TP} + \text{FP}}$ = Out of all positive prediction → the ^{percent (pct)} pct of actual positive

Recall = $\frac{\text{TP}}{\text{TP} + \text{FN}}$ = Out of all actual positive → the pct of ~~predict~~ true positive prediction.

Sensitivity = True positive rate

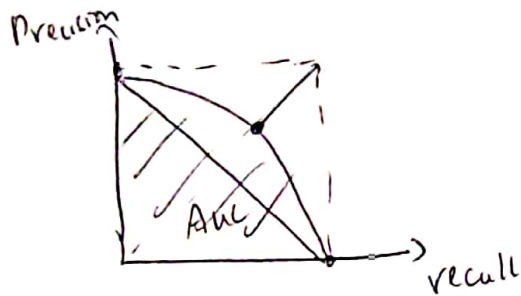
Specificity = False positive rate

Accuracy = ~~percentage~~ number of correct predictions over all instances

$$= \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

- The F-measure uses harmonic mean to combine and ~~avg~~ "average" precision and recall to account for their trade off

- The higher the area under the curve in precision-recall curve, the better



- k classes :

$$\text{Accuracy} = \frac{\sum C_{ii}}{\sum \sum C_{ij}}$$

$$\text{precision for class } i = \frac{C_{ii}}{\sum_j C_{ij}}$$

$$\text{recall} = \frac{C_{ii}}{\sum_i C_{ij}}$$

$$\rightarrow F\text{-measur} = 2 \frac{P \times R}{P + R}$$