

HW09- CS 189

1.

Who else did you work with on this homework? In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

I worked on this homework with Ehimare Okoyomon, Prashanth Ganeth, and Daniel Mockaitis. We worked by getting together throughout the week and communicating on facebook.

I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up}

Nicholas Lorio, 26089160

Question 6 Own Question.

How can we define entropy in the scope of machine learning?

Entropy is defined as the sum of probability of each label times the log probability of that same label. It is our measure of impurity in the decision tree. There is entropy in the probabilities of this tree which arises from the entropy/randomness of the probabilities which determine our classification labeling assignments. The outcome of these decisions are determined by tests performed on the attributes of our features.

At each node of the tree, a test is performed for every features. The feature that has the largest information gain, the difference in the entropy before and the entropy after the decision, is chosen for the decision split. ML favors features that produce decisions with low uncertainty/entropy. The process stops when the leaf node contains instances which all have the same class.

This methodology for ML helps to label for classification. We desire to minimize the incorrect labeling which in turns means we wish to minimize entropy in our decisions.

HW09

189 (2. a.) $L(F_{\text{opt}}(x), y) = \begin{cases} 0, & i=y \\ \lambda_c, & i \neq y \\ \lambda_d, & i=c+1 \end{cases}$

①

$$R(F(x)|x) = \sum_{i=1}^c L(F(x), i) P(Y=i|x)$$

$$= L(F(x), i) + \sum_{j \neq i} L(F(x), j) P(Y=j|x)$$

$$R(i|x) = 0 \cdot P(i|x) + \lambda_c \cdot \sum_{j \neq i} P(j|x) = \lambda_c \sum_{j \neq i} P(j|x)$$

$P(i|x)$ is the largest i : we know that choosing any other class assignment j would make our risk bigger as

$$\sum_{j \neq i} P(j|x) \leq \sum_{j \neq \text{any other } i \text{ classes}} P(j|x)$$

② Show that warrants the risk

$$R(i|x) \leq R(c+1|x)$$

$$\lambda_c \sum_{j \neq i} P(j|x) \leq \sum_i L(c+1|x) P(i|x)$$

$$\lambda_c \sum_{j \neq i} P(j|x) \leq \lambda_c \underbrace{\sum_{i=1}^c P(i|x)}$$

$$\lambda_c (1 - P(i|x)) \leq \lambda_c \quad \text{I encloses all probabilities}$$

$$1 - P(i|x) \leq \frac{\lambda_d}{\lambda_c}, \quad P(i|x) \geq 1 - \frac{\lambda_d}{\lambda_c} \quad \text{as required.}$$

as such the given $F_{\text{opt}}(x)$ obtains minimum risk.

HW09

$$(2.a. intuition) \min_{F_{opt}(x)} \sum_i P(Y=i|x) P(F(x), i) = \min_{F_{opt}(x)} E(F(x))$$

$\Rightarrow P(Y=i|x) \geq P(Y=j|x)$ For all j , class w/ largest prob given x .

② Choose class i : IF $P(Y=i|x) \geq 1 - \frac{\lambda_c}{\lambda_c} = \frac{\lambda_c - \lambda_d}{\lambda_c}$

③ choose Distr else

$F_{opt}(x) \rightarrow \text{Find } \max_i P(Y=i|x)$

①

return i IF $P(Y=i|x) \geq 1 - \frac{\lambda_d}{\lambda_c}$

②

③

return $C+1$ else

① we want to choose the class i which is most certain / has greatest probability of being associated w/ the given x .

② Assuming $\lambda_c < \lambda_d$: if our penalty for misclassification (λ_c) is much greater than our penalty for distract (λ_d) $\Rightarrow 1 - \frac{\lambda_d}{\lambda_c}$ becomes very close to 1. This means that we only choose i : given x IF we are very certain of our assignment. This is all respect to the values of λ_c , λ_d . If they are close to each other but $\lambda_c > \lambda_d$ then we are more likely to choose i .

③ we classify our x into $C+1$ IF this level of certainty is not reached. It is not worth the penalty risk to guess.

i. we capture our risk $F_{opt}(x)$

2.b. if $\lambda_d = 0$ then under our current decision rules we will always assign class 1 to x .

i) \Rightarrow Need to change decision rule ①

2* We should always select class $C+2$ since there is no longer any penalty associated with doubt.

ii) $\lambda_d > \lambda_C$ then $P(Y=1|X) \geq 1 - \frac{\lambda_d}{\lambda_C} = \text{negative}$

in this case we will always guess / assign class 1 to x regardless of how high the probability is. This because all probability are between 0 & 1 and will always be greater than a negative value.

- this is not intuitive as it defeats the point of trying to minimize risk associated with an incorrect classification assignment. (for $\lambda_d > \lambda_C$ case)

- the $\lambda_d = 0$ case is not intuitive as we effectively do not complete any classification whatsoever & are back to where we started if we classification assign every x to the $C+2$ (doubt class)

3. a. label 1 w/ $P = \pi_1$, $X \sim N(\mu_1, \Sigma)$
 $X \in \mathbb{R}^d$
 label 2 w/ $P = \pi_2$

$\mu_1, \mu_2 \in \mathbb{R}^d$
 $\Sigma \in \mathbb{R}^{d \times d}$
 i) MLE($\hat{\mu}_1, \hat{\mu}_2$) = arg max $P(\text{Data} | \text{true})$
 $P(\text{Data} | \text{label}) =$

(Saves Σ) • Using $Q_k(x)$, as defined in Note 18 we
 can determine the decision boundary for this case.

• MLE has no priors.

$$L(X|\theta) = \frac{1}{2} (x - \hat{\mu}_1)^T \sigma^{-2} I (x - \hat{\mu}_1) - \frac{1}{2} \ln(|\sigma^2 I|) = Q_1(x)$$

$$= \frac{1}{2} (x - \hat{\mu}_2)^T \sigma^{-2} I (x - \hat{\mu}_2) - \frac{1}{2} \ln(|\sigma^2 I|)$$

$$(x - \hat{\mu}_1)^T \sigma^{-2} I (x - \hat{\mu}_1) = (x - \hat{\mu}_2)^T \sigma^{-2} I (x - \hat{\mu}_2)$$

$$(x - \hat{\mu}_1)^T (x - \hat{\mu}_1) = (x - \hat{\mu}_2)^T (x - \hat{\mu}_2)$$

• Our Decision boundary is the set of points $x \in \mathcal{X}$ for
 when $\|x - \hat{\mu}_1\|_2^2 = \|x - \hat{\mu}_2\|_2^2$

$$\hat{L}_{\text{MLE}} = \begin{cases} 1, & \|x - \hat{\mu}_1\|_2^2 > \|x - \hat{\mu}_2\|_2^2 \\ 2, & \|x - \hat{\mu}_1\|_2^2 < \|x - \hat{\mu}_2\|_2^2 \\ 0.5, & \|x - \hat{\mu}_1\|_2^2 = \|x - \hat{\mu}_2\|_2^2 \end{cases}$$

ii) MAP w/ priors $P(1) \neq P(2)$

$$Q_1(x) = Q_2(x)$$

$$\left\{ \begin{array}{l} \text{let } P(1) = \pi, \neq P(2) = \pi_2 \\ \pi_2 = 1 - \pi_1 \end{array} \right.$$

$$\ln(P(1)) = \frac{1}{2}(x - \hat{\mu}_1)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_1) = \frac{1}{2} \ln(1/\hat{\pi}_1)$$

$$\ln(P(2)) = \frac{1}{2}(x - \hat{\mu}_2)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_2) = \frac{1}{2} \ln(1/\hat{\pi}_2)$$

$$2 \ln(P(1)) = (x - \hat{\mu}_1)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_1) \geq 2 \ln(P(2)) = (x - \hat{\mu}_2)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_2)$$

$$\Rightarrow (x^T \hat{\Sigma}^{-1} - \hat{\mu}_1^T \hat{\Sigma}^{-1})(x - \hat{\mu}_1) = x^T \hat{\Sigma}^{-1} x - x^T \hat{\Sigma}^{-1} \hat{\mu}_1 \\ - \hat{\mu}_1^T \hat{\Sigma}^{-1} x + \hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1$$

$$2 \ln P(1) = x^T \hat{\Sigma}^{-1} x + 2 x^T \hat{\Sigma}^{-1} \hat{\mu}_1 - \hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 = 2 \ln P(2) = x^T \hat{\Sigma}^{-1} x + 2 x^T \hat{\Sigma}^{-1} \hat{\mu}_2 - \hat{\mu}_2^T \hat{\Sigma}^{-1} \hat{\mu}_2$$

$$\nabla(\ln(P(1)) - \ln(P(2))) + x^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2) - (\hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 - \hat{\mu}_2^T \hat{\Sigma}^{-1} \hat{\mu}_2) = 0$$

$$\underbrace{x^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2)}_{w} + \left(\ln \left(\frac{P(1)}{P(2)} \right) - \left(\hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 - \hat{\mu}_2^T \hat{\Sigma}^{-1} \hat{\mu}_2 \right) \right) = 0$$

w : linear function, $f(x) = w^T x + b = x^T w + b$

$$\hat{\mu}_{MAP} = \begin{cases} \hat{\mu}_1, & > \\ \hat{\mu}_2, & < \end{cases}$$

" from below " "

$$\text{or } 2, 2 \ln(P(1)) - (x - \hat{\mu}_1)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_1) \geq 2 \ln(P(2)) - (x - \hat{\mu}_2)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_2)$$

$\ln(1) = 0$ so when the priors are equal then
MAP gives same decision boundary rules as MLE

$$F(x) = w^T x + \kappa \quad x \sim N(\mu_L, \Sigma) \quad w/ L=1, 2$$

$$3.b. F(x) = w^T (x - v), \text{ let } w^T = (\Sigma^{-1}(\hat{\mu}_1 - \hat{\mu}_2))^T$$

• let $y \in \mathbb{R}^2$ be one unit vector $y_i = e_k$ rather than $y_i = k$

• $\pi_1, \pi_2 = 1/2 \rightarrow$ Non-zero mean

• Compute $\Sigma_{xx}, \Sigma_{xy}, \Sigma_{yy}$ as function of μ_1, μ_2, Σ

• take property of expectation Hint

$\begin{pmatrix} 0 \\ \text{k-th spot } 1 \\ 0 \\ 0 \end{pmatrix}$ std Basis
each label
is a std basis
vector

$$\Sigma_{xx} = \sum$$

$$\Sigma_{xy} = E[(x - E(x))(y - E(y))^\top] = E[xy^\top] - E[x]E[y]^\top$$

$$E[x] = E[E(x|L=i)] = \frac{1}{2} E[x|L=1] + \frac{1}{2} E[x|L=2] = \frac{\mu_1 + \mu_2}{2} \text{ (d x 1 vector)}$$

$$E[y] = E[E(y|L=i)] = \frac{1}{2} \cdot E(y|L=1) + \frac{1}{2} E(y|L=2)$$

$$\frac{1}{2} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$\Sigma_{xy} = E[E(xy^\top | L=i)] = \frac{1}{2} E[xy^\top | L=1] + \frac{1}{2} E[xy^\top | L=2]$$

$$= \frac{1}{2} \left(\underset{d \times 1}{\mu_1} \cdot \underset{1 \times 2}{\begin{bmatrix} 1 & 0 \end{bmatrix}} + \underset{d \times 1}{\mu_2} \cdot \underset{1 \times 2}{\begin{bmatrix} 0 & 1 \end{bmatrix}} \right) = \frac{1}{2} [\vec{\mu}_1 \vec{0}] + \frac{1}{2} [\vec{0} \vec{\mu}_2] = \begin{bmatrix} \mu_1 & \mu_2 \end{bmatrix}$$

$$\Sigma_{xy} = \begin{bmatrix} \mu_1/2 & \mu_2/2 \end{bmatrix} - \frac{\mu_1 + \mu_2}{2} \begin{bmatrix} 1/2 & 1/2 \end{bmatrix}$$

$$= \begin{bmatrix} \mu_1/2 & \mu_2/2 \end{bmatrix} - \boxed{\begin{bmatrix} \frac{\mu_1 + \mu_2}{4} & \frac{\mu_1 + \mu_2}{4} \end{bmatrix}} = \boxed{\begin{bmatrix} \frac{\mu_1 - \mu_2}{4} & \frac{\mu_2 - \mu_1}{4} \end{bmatrix}}$$

3.b. cont.

$$\bullet \Sigma_{yy} = E[yy^T] - E[y]E[y]^T$$

w/ computers from earlier & the following

$$\begin{aligned} E[yy^T] &= E[E(yy^T | L=1)] \in \frac{1}{2} E[yy^T | L=1] + \frac{1}{2} E[yy^T | L=2] \\ &\stackrel{2 \times 1 \quad 1 \times 2}{=} \frac{1}{2} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} \\ &= \frac{1}{2} \left(\begin{bmatrix} 1 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 1 \end{bmatrix} \right) \\ &= \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \end{aligned}$$

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

$\begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix}$

3.1

3.c. Part 1

$$\rho(\mathbf{v}^T \mathbf{x}, \mathbf{v}^T \mathbf{y}) = \frac{\text{Cov}(\mathbf{v}^T \mathbf{x}, \mathbf{v}^T \mathbf{y})}{\sqrt{\text{Var}(\mathbf{v}^T \mathbf{x})} \sqrt{\text{Var}(\mathbf{v}^T \mathbf{y})}}^{1/2}$$

$$\text{Cov}(\mathbf{v}^T \mathbf{x}, \mathbf{v}^T \mathbf{y}) = E[(\mathbf{v}^T \mathbf{x} - E[\mathbf{v}^T \mathbf{x}])(\mathbf{v}^T \mathbf{y} - E[\mathbf{v}^T \mathbf{y}])] =$$

$$= E[(\mathbf{v}^T (\mathbf{x} - E[\mathbf{x}]) (\mathbf{y}^T \mathbf{v} - E[\mathbf{y}])^T \mathbf{v}]$$

$$= \mathbf{v}^T E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{y}^T - E[\mathbf{y}])^T] \mathbf{v}$$

By similar logic

$$\text{Var}(\mathbf{v}^T \mathbf{x}) = \mathbf{v}^T \text{Var}(\mathbf{x}) \mathbf{v} \quad \rho = \frac{\mathbf{v}^T \boldsymbol{\xi}_{xy} \mathbf{v}}{(\mathbf{v}^T \boldsymbol{\xi}_{xx} \mathbf{v} \mathbf{v}^T \boldsymbol{\xi}_{yy} \mathbf{v})^{1/2}}$$

$$\text{Var}(\mathbf{v}^T \mathbf{y}) = \mathbf{v}^T \text{Var}(\mathbf{y}) \mathbf{v}$$

lets Simplify this: $\rho = \frac{\mathbf{v}^T \frac{1}{n} [\mu_1 - \mu_2 \ \mu_2 - \mu_1] \mathbf{v}}{(\mathbf{v}^T \boldsymbol{\xi}_{xy} \mathbf{v} + \mathbf{v}^T \frac{1}{n} [\mu_1 - \mu_1] \mathbf{v})^{1/2}}$

$$\propto \frac{\mathbf{v}^T [\mu_1 \ \mu_2] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \mathbf{v}}{(\mathbf{v}^T \boldsymbol{\xi}_{xy} \mathbf{v} + \mathbf{v}^T \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \mathbf{v})^{1/2}}$$

• Square top & bottom
of the expression. (for ρ^2)

$$\propto \frac{(\mathbf{v}^T [\mu_1 \ \mu_2] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \mathbf{v})^2}{\mathbf{v}^T \boldsymbol{\xi}_{xy} \mathbf{v} + \mathbf{v}^T \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \mathbf{v}}$$

$$\cdot \delta/\delta \mathbf{v} = 2(\mathbf{v}^T [\mu_1 \ \mu_2]) \cdot [\mu_1 \ \mu_2] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \cdot (\mathbf{v}^T \boldsymbol{\xi}_{xy})^{-1}$$

$$+ (\mathbf{v}^T [\mu_1 \ \mu_2])^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \cdot (-1) \cdot (\mathbf{v}^T \boldsymbol{\xi}_{xy})^{-2} \cdot \mathbf{v}^T (\boldsymbol{\xi}_{xy} \boldsymbol{\xi}_{xy}^T)$$

• $\boldsymbol{\xi}$ is symmetric so
 $\boldsymbol{\xi}^T \boldsymbol{\xi}$

3.2

3.C.I. cont.

• Set each to zero, solve for U^*

$$\begin{aligned}
 & 2U^T[\mu_1, \mu_2][\mu_1, \mu_2] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \cdot (U^T \Sigma U)^{-1} = \\
 & (U^T[\mu_1, \mu_2])^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} (U^T \Sigma U)^{-2} 2U^T \Sigma \\
 & U^T[\mu_1, \mu_2]^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix (U^T \Sigma U)^{-1} = U^T \Sigma [\mu_1, \mu_2]^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix (U^T \Sigma U)^{-2} \\
 & \Sigma^{-1} U^T[\mu_1, \mu_2]^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix (U^T \Sigma U)^{-1} = U^T \Sigma [\mu_1, \mu_2]^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix (U^T \Sigma U)^{-2} \\
 & \Sigma^{-1} [\mu_1, \mu_2]^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix = U^T \Sigma [\mu_1, \mu_2]^2 \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix (U^T \Sigma U)^{-1} \\
 & \Sigma^{-1} [\mu_1, \mu_2] \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix (U^T \Sigma U) \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix^{-1} [\mu_1, \mu_2]^2 = U^T \Sigma^2 \\
 & U^* = \left(\left(\Sigma^{-1} [\mu_1, \mu_2] \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix U^T \Sigma U \begin{bmatrix} 1 & -1 \\ -1 & 1 \endbmatrix^{-1} [\mu_1, \mu_2]^2 \right)^T \right)^{1/2}
 \end{aligned}$$

thus U^* & $\Sigma^{-1}[\mu_1, \mu_2]$ as required.

3.3

3.C. part two

$$\mathbf{v}^* \propto \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$\& \mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

thus $\mathbf{v}^* \propto \mathbf{w}$

$$Y = F(\mathbf{x}) = \mathbf{w}^\top (\mathbf{x} - \mathbf{v}), f(\mathbf{x}) = \mathbf{v}^* (\mathbf{x} - \mathbf{v})$$

thus determining Decision boundary LDA & performing CCA are equivalent.

$$\text{given } \Sigma_{xx} \Sigma_{xy} \Sigma_{yy} \quad \begin{bmatrix} 1/\mu & -1/\mu \\ -1/\mu & \nu/\mu \end{bmatrix}^T = \begin{bmatrix} 1/\mu & -1/\mu \\ -1/\mu & 1/\mu \end{bmatrix} \quad 3, \text{ otherwise}$$

3.c.

$$\underset{U,V}{\text{Max}} \rho(U^T X, V^T Y) = \frac{\text{Cov}(U^T X, V^T Y)}{(\text{Var}(U^T X) \text{Var}(V^T Y))^{1/2}} = \frac{U^T \text{Cov}(X, Y) V}{(U^T \text{Var}(X) U V^T \text{Var}(Y) V)^{1/2}}$$

* From Note 11

$$= \frac{U^T \Sigma_{xy} V}{(U^T \Sigma_{xx} U V^T \Sigma_{yy} V)^{1/2}}$$

• let us perform a change of variables as follows:

$$a = \Sigma_{xx}^{-1/2} U, b = \Sigma_{yy}^{-1/2} V \Rightarrow U = \Sigma_{xx}^{-1/2} a, V = \Sigma_{yy}^{-1/2} b$$

$$\text{so } \rho = \frac{a^T \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} b}{(a^T a b^T b)^{1/2}} \quad * \Sigma_{xx} \text{ is DSD, Symmetric}$$

$$\text{so } \Sigma_{xx}^{-T} = \Sigma_{xx}$$

$$\text{, } \Sigma_{yy} \text{ symmetric: } \Sigma_{yy}^{-T} = \Sigma_{yy}$$

• By Cauchy Schwartz: $|\langle U, V \rangle|^2 \leq \langle U, U \rangle \langle V, V \rangle$

$$\begin{aligned} (a^T \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} b)^2 &\leq (a^T \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} (a^T \Sigma_{xx}^{1/2} \Sigma_{yy} \Sigma_{yy}^{-1/2})^T)^2 (b^T b)^{1/2} \\ &\leq (a^T \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} \Sigma_{yy}^{-T} \Sigma_{xy}^T \Sigma_{xx}^{-1/2} a)^{1/2} (b^T b)^{1/2} \end{aligned}$$

$$\text{thus } \rho \leq \frac{(a^T \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} \Sigma_{xy}^T \Sigma_{xx}^{-1/2} a)^{1/2}}{(a^T a)^{1/2}}$$

• this is the Rayleigh Quotient $R(x) = \frac{x^T A x}{x^T x}$ & is bounded by $\lambda_{\min}(A) \leq R(x) \leq \lambda_{\max}(A)$

$\Rightarrow a$ is an eigenvector of $(\Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} \Sigma_{xy}^T \Sigma_{xx}^{-1/2})$,

• Rayleigh Quotient is maximal by the eigenvector w/ the largest eigenvalue. Thus when we Max ρ then a is the eigenvector w/ the corresponding largest eigenvalue.

3.6. cont.

• changing the coordinate basis

$$a = \epsilon_{xx}^{1/2} u, u = \epsilon_{xx}^{-1/2} a$$

• u^* is an eigenvector of $\epsilon_{xx}^{-1} \epsilon_{xy} \epsilon_{yy}^{-1} \epsilon_{xy}^T$

$$\therefore \epsilon^{-1} \begin{bmatrix} M_1 & M_2 \\ M_2 & M_1 \end{bmatrix} \begin{bmatrix}$$

$$\epsilon_{yy}^{-1} = \begin{bmatrix} 1/4 & -1/4 \\ -1/4 & 1/4 \end{bmatrix}^{-1} = \underline{\underline{1}}$$

• ϵ_{yy} is not invertible therefore
this method will not work.

S, x_1, \dots, x_n iid RV from set $\{0, \dots, M\}$ $p_j = P(x_i=j)$ for $j = \{0, \dots, M\}$

w/ \hat{f}^n as empirical type is $\left(\hat{f}_0^{(n)}, \hat{f}_1^{(n)}, \dots, \hat{f}_M^{(n)}\right)$

'crude' Stirling approx $\Delta! \approx \left(\frac{\Delta}{e}\right)^{\Delta}$

\hat{f}^n is a vector w/ M entries where i th entry is $\hat{f}_i^{(n)}$

$$\text{Show: } \left(\hat{f}^n\right)^{\Delta!} = \left(\frac{\sum_{j=0}^M \hat{f}_j^{(n)}}{\Delta}, \frac{\hat{f}_1^{(n)}}{\Delta}, \dots, \frac{\hat{f}_M^{(n)}}{\Delta}\right) \approx \exp(\Delta H(\hat{f}^n / \Delta))$$

w/ $H(p) = \sum_{j=0}^M p_j \ln\left(\frac{1}{p_j}\right)$

$$\textcircled{1} \quad \Delta!$$

$$\frac{\hat{f}_0^{(n)}! \hat{f}_1^{(n)}! \dots \hat{f}_M^{(n)}!}{\left(\frac{\hat{f}_0^{(n)}}{\Delta}\right)^{\hat{f}_0^{(n)}} \dots \left(\frac{\hat{f}_M^{(n)}}{\Delta}\right)^{\hat{f}_M^{(n)}}} \approx$$

$$\approx \frac{\Delta^\Delta \cdot e^{\hat{f}_0^{(n)}} \dots e^{\hat{f}_M^{(n)}}}{e^\Delta \left(\frac{\hat{f}_0^{(n)}}{\Delta}\right)^{\hat{f}_0^{(n)}} \dots \left(\frac{\hat{f}_M^{(n)}}{\Delta}\right)^{\hat{f}_M^{(n)}}} = \frac{\Delta^\Delta \exp\left(\sum_{j=0}^M \hat{f}_j^{(n)}\right) e^{\Delta}}{\left(\frac{\hat{f}_0^{(n)}}{\Delta}\right)^{\hat{f}_0^{(n)}} \dots \left(\frac{\hat{f}_M^{(n)}}{\Delta}\right)^{\hat{f}_M^{(n)}}} = \frac{\Delta^\Delta e^{\Delta} e^{\sum_j \hat{f}_j^{(n)}}}{\prod_{j=0}^M \hat{f}_j^{(n)}}$$

$$\textcircled{2} \quad \exp(\Delta H(\hat{f}^n / \Delta))$$

$$= \exp\left(\Delta \sum_{j=0}^M \frac{\hat{f}_j^{(n)}}{\Delta} \cdot \ln\left(\frac{\Delta}{\hat{f}_j^{(n)}}\right)\right) = \exp\left(\sum_{j=0}^M \frac{\hat{f}_j^{(n)}}{\Delta} \ln\left(\frac{\Delta}{\hat{f}_j^{(n)}}\right)\right)$$

$$= \frac{\exp\left(\sum_j \hat{f}_j^{(n)} \ln\left(\frac{\Delta}{\hat{f}_j^{(n)}}\right)\right)}{\exp\left(\sum_j \frac{\hat{f}_j^{(n)}}{\Delta} \ln\left(\frac{\Delta}{\hat{f}_j^{(n)}}\right)\right)} \cdot \frac{\Delta^\Delta \exp\left(\sum_j \hat{f}_j^{(n)} / \Delta\right)}{\exp\left(\sum_j \ln\left(\hat{f}_j^{(n)} \Delta^{-1} \cdot \frac{\Delta}{\hat{f}_j^{(n)}}\right)\right)}$$

$$\begin{aligned}
 ② &= \frac{\Lambda^{\Lambda} e^{\sum f_j(\lambda) / \lambda}}{\sum \exp\left(\frac{f_0(\lambda) + f_1(\lambda)}{\lambda}\right)} = \frac{\Lambda^{\Lambda} e^{\sum f_j(\lambda) / \lambda}}{e^{f_0(\lambda)} \cdot e^{f_1(\lambda)}} \dots \\
 &= \frac{\Lambda^{\Lambda} e^{\sum f_j(\lambda) / \lambda}}{\left(\frac{e^{f_0(\lambda)}}{f_0}, \dots, \frac{e^{f_M(\lambda)}}{f_M}\right)} = \frac{\Lambda^{\Lambda} e^{\sum f_j(\lambda) / \lambda}}{\prod_{j=0}^M f_j(\lambda)^{f_j(\lambda)}} \\
 &= \frac{\Lambda^{\Lambda} e^{-\Lambda} e^{\sum f_j(\lambda)}}{\prod_{j=0}^M f_j(\lambda)^{f_j(\lambda)}}
 \end{aligned}$$

Showing approximate equivalence as round

S.b. Show $\lim_{n \rightarrow \infty} \log P(F^{(n)} = f^{(n)}) = -KL(F, p) = -\sum f_i \log \left(\frac{f_i}{p_i}\right)$

$$\lim_{n \rightarrow \infty} \log \left(\left(\frac{f_0^{(n)}, \dots, f_M^{(n)}}{\pi_0, \dots, \pi_M} \right) \prod_{j=1}^M p_j^{f_j^{(n)}} \right) \quad \text{large enough Stirling approximation}$$

$$\lim_{n \rightarrow \infty} \log \left(\exp \left(n \sum_{j=0}^M \left(\frac{f_j^{(n)}}{\pi_j} \right) \ln \left(\frac{n}{f_j^{(n)}} \right) \right) \prod_{j=1}^M p_j^{f_j^{(n)}} \right)$$

$$= \lim_{n \rightarrow \infty} n \sum_{j=0}^M \frac{f_j^{(n)}}{\pi_j} \ln \left(\frac{\pi_j}{f_j^{(n)}} \right) + \sum_{j=1}^M \frac{f_j^{(n)}}{\pi_j} \ln(p_j)$$

$$= \lim_{n \rightarrow \infty} \frac{f_0^{(n)}}{\pi_0} \ln \left(\frac{\pi_0}{f_0^{(n)}} \right) + \sum_{j=1}^M \left(\frac{f_j^{(n)}}{\pi_j} \ln \left(\frac{\pi_j}{f_j^{(n)}} \right) + \frac{f_j^{(n)}}{\pi_j} \ln(p_j) \right)$$

• $\lim_{n \rightarrow \infty} \frac{1}{n} F^{(n)} = f$ by definition

$$= \lim_{n \rightarrow \infty} \frac{f_0^{(n)}}{\pi_0} \ln \left(\frac{\pi_0}{f_0^{(n)}} \right) + \sum_{j=1}^M \left(\frac{f_j^{(n)}}{\pi_j} \left(\ln \left(\frac{\pi_j}{f_j^{(n)}} \right) + \ln(p_j) \right) \right) \\ + \sum_{j=1}^M \frac{f_j^{(n)}}{\pi_j} \ln \left(\frac{1 \cdot p_j}{f_j^{(n)}} \right)$$

$$= \lim_{n \rightarrow \infty} -\frac{f_0^{(n)}}{\pi_0} \ln \left(\frac{F_0^{(n)}}{\pi_0} \right) - \sum_{j=1}^M \frac{f_j^{(n)}}{\pi_j} \ln \left(\frac{f_j^{(n)}}{\pi_j \cdot p_j} \right)$$

$$= -F_0 \ln(F_0) - \sum_{j=1}^M f_j \ln \left(\frac{f_j}{p_j} \right) *$$

$$= -\sum_{j=0}^M f_j \ln \left(\frac{f_j}{p_j} \right) \quad \text{as required.}$$

S.C.

$$\text{Show: } \text{KL}(P, Q_\theta) = C - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log(Q_\theta(y|x))$$

$$Q_\theta(y|x) = \frac{Q_\theta(x,y)}{Q(x)} =$$

$$\begin{aligned} \text{KL}(P, Q_\theta) &= \sum_i p_i \log \left(\frac{p_i}{Q_{\theta,i}} \right) && \cdot \text{gm wrt } P \text{ wrt prob} \\ &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \left(\frac{p(x, y)}{Q_\theta(x, y)} \right) && \text{distrn } p(x, y) \neq Q_\theta(x, y) \\ &&& \text{wrt distrn } Q_\theta(x, y) \end{aligned}$$

$$\begin{aligned} \log \left(\frac{p(x, y)}{Q_\theta(x, y)} \right) &= \log(p(x, y)) - \log(Q_\theta(x, y)) \\ &= \log(p(x, y)) - \log(Q_\theta(y|x) Q(x)) \\ &= \log(p(x, y)) - \log(Q_\theta(y|x)) - \log(Q(x)) \end{aligned}$$

$$= \sum_{x \in X} \sum_{y \in Y} p(x, y) [\log(p(x, y)) - \log(Q_\theta(y|x)) - \log(Q(x))]$$

$$\begin{aligned} &= \underbrace{\sum_{x \in X} \sum_{y \in Y} p(x, y) \log(p(x, y))}_{C_2} - \underbrace{p(x, y) \log(Q_\theta(y|x))}_{C_1} - \underbrace{p(x, y) \log(Q(x))}_{C_1} \\ &= C - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log(Q_\theta(y|x)) \quad \text{as required} \end{aligned}$$

as C is some constant independent of θ .

logistic regression

S.D. $y_i = +1$ or $y_i = -1$ predicted labels from features x_i using

transition probabilities model: $q_\theta(y_i|x_i) = \frac{1}{1 + \exp(-y_i \theta^T x_i)}$

assume all x_i are distinct

vec, show that w/ $p(x,y) = \begin{cases} 1/n & \text{if } x=x_i \text{ & } y=y_i \\ 0 & \text{else} \end{cases} i=1,2,\dots,n$

we get $\min_{\theta} KL(p, q_\theta) = \min_{\theta} -\frac{1}{n} \sum_i \log(q_\theta(y_i|x_i))$

From c.

$$\begin{aligned} \min_{\theta} C &= \sum_x \sum_y p(x,y) \log q_\theta(y|x) \approx \min_{\theta} -\sum_{x,y} p(x,y) \log \left(\frac{1}{1 + \exp(-y \theta^T x)} \right) \\ &= \min_{\theta} -\sum_i \frac{1}{n} \log \left((1 + \exp(-y_i \theta^T x_i))^{-1} \right) + 0 \\ &= \min_{\theta} \theta - \frac{1}{n} \sum_i \log(p_q(y_i|x_i)) \quad \text{as required} \end{aligned}$$

or $\theta = 0$ if $p(x,y) = 0$ in the else scenario.

Question 4

```
In [17]: %matplotlib inline
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import scipy.spatial

# Gradient descent optimization
# The learning rate is specified by eta
class GDOptimizer(object):
    def __init__(self, eta):
        self.eta = eta

    def initialize(self, layers):
        pass

    # This function performs one gradient descent step
    # layers is a list of dense layers in the network
    # g is a list of gradients going into each layer before the nonlinear activation
    # a is a list of activations of each node in the previous layer going
    def update(self, layers, g, a):
        m = a[0].shape[1]
        for layer, curGrad, curA in zip(layers, g, a):
            layer.updateWeights(-self.eta / m * np.dot(curGrad, curA.T))
            layer.updateBias(-self.eta / m * np.sum(curGrad, 1).reshape(layer.b.shape))

    # Cost function used to compute prediction errors
class QuadraticCost(object):

    # Compute the squared error between the prediction yp and the observation y
    # This method should compute the cost per element such that the output is the
    # same shape as y and yp
    @staticmethod
    def fx(y, yp):
        return 0.5 * np.square(yp - y)

    # Derivative of the cost function with respect to yp
    @staticmethod
    def dx(y, yp):
        return y - yp

    # Sigmoid function fully implemented as an example
class SigmoidActivation(object):
    @staticmethod
    def fx(z):
        return 1 / (1 + np.exp(-z))

    @staticmethod
    def dx(z):
        return SigmoidActivation.fx(z) * (1 - SigmoidActivation.fx(z))

    # Hyperbolic tangent function
class TanhActivation(object):

    # Compute tanh for each element in the input z
    @staticmethod
    def fx(z):
        return np.tanh(z)
```

4.a

```
In [24]: import numpy as np
import scipy.spatial

#####
## Models used for predictions.
#####
def compute_update(single_obj_loc, sensor_loc, single_distance):
    """
    Compute the gradient of the log-likelihood function for part a.

    Input:
    single_obj_loc: 1 * d numpy array.
    Location of the single object.

    sensor_loc: k * d numpy array.
    Location of sensor.

    single_distance: k dimensional numpy array.
    Observed distance of the object.

    Output:
    grad: d-dimensional numpy array.

    """
    loc_difference = single_obj_loc - sensor_loc # k * d.
    phi = np.linalg.norm(loc_difference, axis=1) # k.
    grad = loc_difference / np.expand_dims(phi, 1) # k * 2.
    update = np.linalg.solve(grad.T.dot(grad), grad.T.dot(single_distance - phi))

    return update

def get_object_location(sensor_loc, single_distance, num_iters=20, num_repeats=10):
    """
    Compute the gradient of the log-likelihood function for part a.

    Input:
    sensor_loc: k * d numpy array. Location of sensor.

    single_distance: k dimensional numpy array.
    Observed distance of the object.

    Output:
    obj_loc: 1 * d numpy array. The mle for the location of the object.

    """
    obj_locs = np.zeros((num_repeats, 1, 2))
    distances = np.zeros(num_repeats)
    for i in range(num_repeats):
        obj_loc = np.random.randn(1, 2) * 100
        for t in range(num_iters):
            obj_loc += compute_update(obj_loc, sensor_loc, single_distance)

        distances[i] = np.sum((single_distance - np.linalg.norm(obj_loc - sensor_loc, axis=1))**2)
        obj_locs[i] = obj_loc

    obj_loc = obj_locs[np.argmin(distances)]

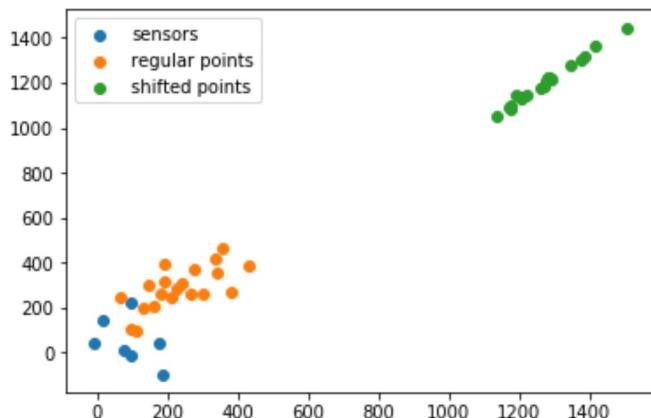
    return obj_loc[0]
```

4.b

```
In [19]: def main():
    np.random.seed(0)
    sensor_loc = generate_sensors()
    regular_loc, _ = generate_dataset(
        sensor_loc,
        num_sensors=sensor_loc.shape[0],
        spatial_dim=2,
        num_data=20,
        original_dist=True,
        noise=1)
    shifted_loc, _ = generate_dataset(
        sensor_loc,
        num_sensors=sensor_loc.shape[0],
        spatial_dim=2,
        num_data=20,
        original_dist=False,
        noise=1)

    plt.scatter(sensor_loc[:, 0], sensor_loc[:, 1], label="sensors")
    plt.scatter(regular_loc[:, 0], regular_loc[:, 1], label="regular points")
    plt.scatter(shifted_loc[:, 0], shifted_loc[:, 1], label="shifted points")
    plt.legend()
    plt.savefig("dataset.png")
    plt.show()

if __name__ == "__main__":
    main()
```



4.c

```
In [26]: def main():
    ##### PLOT PART 1 #####
    np.random.seed(0)

    ns = np.arange(10, 310, 20)
    replicates = 5
    num_methods = 6
    num_sets = 3
    mses = np.zeros((len(ns), replicates, num_methods, num_sets))

    def generate_data(sensor_loc, k=7, d=2, n=1, original_dist=True, noise=1):
        return generate_dataset(
            sensor_loc,
            num_sensors=k,
            spatial_dim=d,
            num_data=n,
            original_dist=original_dist,
            noise=noise)

    for s in range(replicates):
        sensor_loc = generate_sensors()
        X_test, Y_test = generate_data(sensor_loc, n=1000)
        X_test2, Y_test2 = generate_data(
            sensor_loc, n=1000, original_dist=False)
        for t, n in enumerate(ns):
            X, Y = generate_data(sensor_loc, n=n) # X [n * 2] Y [n * 7]
            Xs_test, Ys_test = [X, X_test, X_test2], [Y, Y_test, Y_test2]
            ### Linear regression:
            mse = linear_regression(X, Y, Xs_test, Ys_test)
            mses[t, s, 0] = mse

            ### Second-order Polynomial regression:
            mse = poly_regression_second(X, Y, Xs_test, Ys_test)
            mses[t, s, 1] = mse

            ### 3rd-order Polynomial regression:
            mse = poly_regression_cubic(X, Y, Xs_test, Ys_test)
            mses[t, s, 2] = mse

            ### Neural Network:
            mse = neural_network(X, Y, Xs_test, Ys_test)
            mses[t, s, 3] = mse

            ### Generative model:
            mse = generative_model(X, Y, Xs_test, Ys_test)
            mses[t, s, 4] = mse

            ### Oracle model:
            mse = oracle_model(X, Y, Xs_test, Ys_test, sensor_loc)
            mses[t, s, 5] = mse

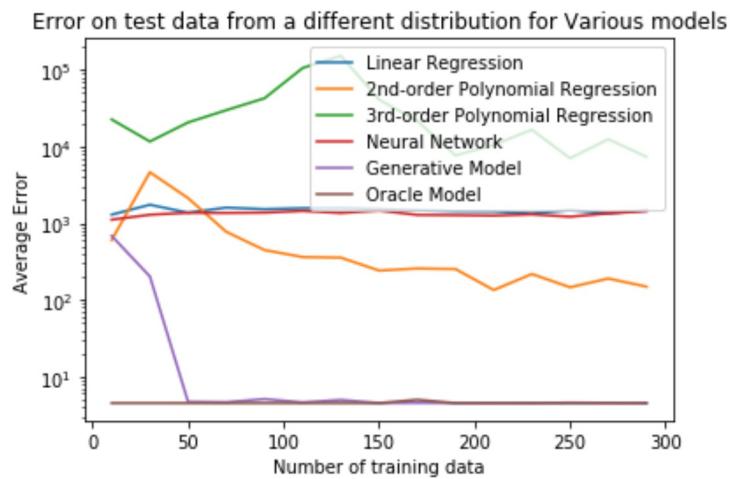
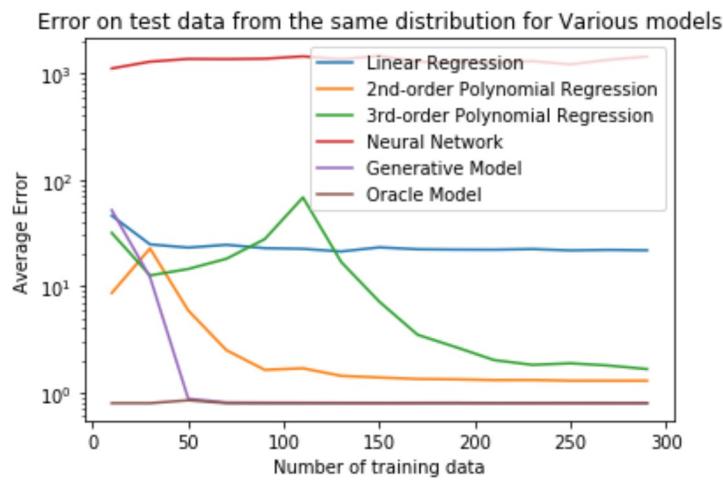
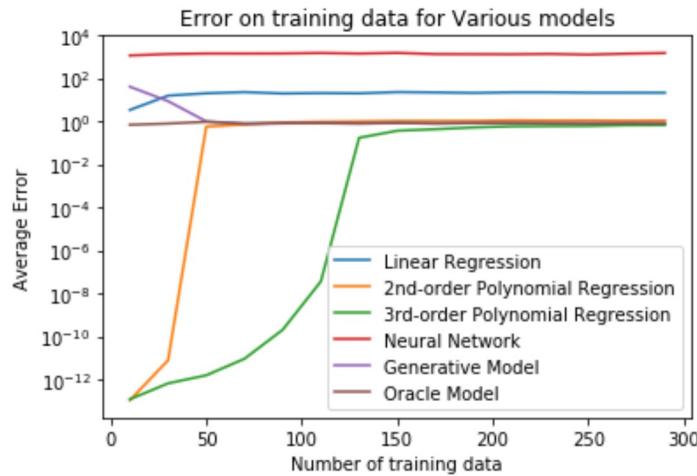
            print('{}th Experiment with {} samples done...'.format(s, n))

    ## Plot MSE for each model.
    plt.figure()
    regressors = [
        'Linear Regression', '2nd-order Polynomial Regression',
        '3rd-order Polynomial Regression', 'Neural Network',
        'Generative Model', 'Oracle Model'
    ]
    for a in range(6):
        plt.plot(ns, np.mean(mses[:, :, a, 0], axis=1), label=regressors[a])
```

C start time: 2018-03-23 20:32:13.431749

```
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:142: FutureWarning: `rcond` parameter will change to the default of machine precision times ``max(M, N)`` where M and N are the input matrix dimensions.
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:172: FutureWarning: `rcond` parameter will change to the default of machine precision times ``max(M, N)`` where M and N are the input matrix dimensions.
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:201: FutureWarning: `rcond` parameter will change to the default of machine precision times ``max(M, N)`` where M and N are the input matrix dimensions.
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.
```

0th Experiment with 10 samples done...
0th Experiment with 30 samples done...
0th Experiment with 50 samples done...
0th Experiment with 70 samples done...
0th Experiment with 90 samples done...
0th Experiment with 110 samples done...
0th Experiment with 130 samples done...
0th Experiment with 150 samples done...
0th Experiment with 170 samples done...
0th Experiment with 190 samples done...
0th Experiment with 210 samples done...
0th Experiment with 230 samples done...
0th Experiment with 250 samples done...
0th Experiment with 270 samples done...
0th Experiment with 290 samples done...
1th Experiment with 10 samples done...
1th Experiment with 30 samples done...
1th Experiment with 50 samples done...
1th Experiment with 70 samples done...
1th Experiment with 90 samples done...
1th Experiment with 110 samples done...
1th Experiment with 130 samples done...
1th Experiment with 150 samples done...
1th Experiment with 170 samples done...
1th Experiment with 190 samples done...
1th Experiment with 210 samples done...
1th Experiment with 230 samples done...
1th Experiment with 250 samples done...
1th Experiment with 270 samples done...
1th Experiment with 290 samples done...
2th Experiment with 10 samples done...
2th Experiment with 30 samples done...
2th Experiment with 50 samples done...
2th Experiment with 70 samples done...
2th Experiment with 90 samples done...
2th Experiment with 110 samples done...
2th Experiment with 130 samples done...
2th Experiment with 150 samples done...
2th Experiment with 170 samples done...
2th Experiment with 190 samples done...
2th Experiment with 210 samples done...
2th Experiment with 230 samples done...
2th Experiment with 250 samples done...
2th Experiment with 270 samples done...
2th Experiment with 290 samples done...
3th Experiment with 10 samples done...
3th Experiment with 30 samples done...
3th Experiment with 50 samples done...
3th Experiment with 70 samples done...
3th Experiment with 90 samples done...
3th Experiment with 110 samples done...
3th Experiment with 130 samples done...
3th Experiment with 150 samples done...
3th Experiment with 170 samples done...
3th Experiment with 190 samples done...
3th Experiment with 210 samples done...
3th Experiment with 230 samples done...
3th Experiment with 250 samples done...
3th Experiment with 270 samples done...
3th Experiment with 290 samples done...
4th Experiment with 10 samples done...
4th Experiment with 30 samples done...
4th Experiment with 50 samples done...
4th Experiment with 70 samples done...



C end time: 2018-03-23 21:01:20.650690

Briefly describe your observations. What are the strengths and weaknesses of each model?

4.d

```
In [20]: activations = dict(ReLU=ReLUActivation,
                           tanh=TanhActivation,
                           linear=LinearActivation)
lr = dict(ReLU=0.0000007,tanh=0.00001,linear=0.005)

def neural_network(X, Y, X_test, Y_test, num_neurons, activation):
    """
    This function performs neural network prediction.
    Input:
        X: independent variables in training data.
        Y: dependent variables in training data.
        X_test: independent variables in test data.
        Y_test: dependent variables in test data.
        num_neurons: number of neurons in each layer
        activation: type of activation, ReLU or tanh
    Output:
        mse: Mean square error on test data.
    """
    ## YOUR CODE HERE
    #####
    model = Model(X.shape[1])
    model.addLayer(DenseLayer(num_neurons, activations[activation]))
    model.addLayer(DenseLayer(num_neurons, activations[activation]))
    model.addLayer(DenseLayer(Y.shape[1], LinearActivation()))
    model.initialize(QuadraticCost())
    model.train(X,Y,500,GDOptimizer(eta=lr[activation]))
    Y_pred = model.predict(X_test)
    mse = np.mean(np.sqrt(np.sum((Y_pred - Y_test)**2, axis=1)))
    print(mse)
    return mse

#####
#####PLOT PART 2#####
#####

def generate_data(sensor_loc, k=7, d=2, n=1, original_dist=True, noise=1):
    return generate_dataset(
        sensor_loc,
        num_sensors=k,
        spatial_dim=d,
        num_data=n,
        original_dist=original_dist,
        noise=noise)

np.random.seed(0)
n = 200
num_neuronss = np.arange(100, 550, 50)
mses = np.zeros((len(num_neuronss), 2))

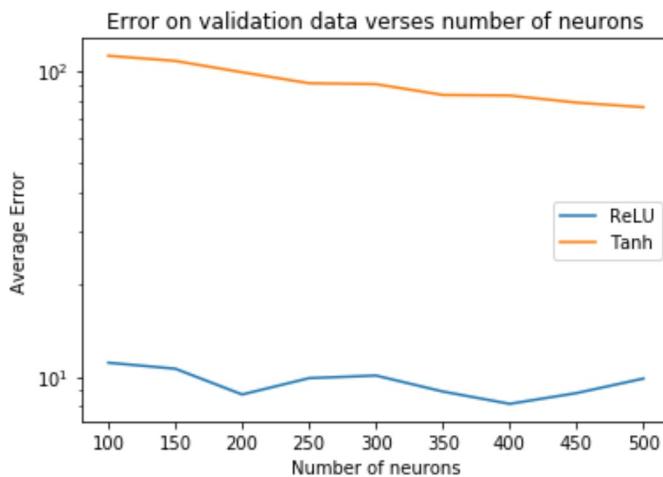
# for s in range(replicates):

    sensor_loc = generate_sensors()
    X, Y = generate_data(sensor_loc, n=n) # X [n * 2] Y [n * 7]
    X_test, Y_test = generate_data(sensor_loc, n=1000)
    for t, num_neurons in enumerate(num_neuronss):
        ### Neural Network:
        mse = neural_network(X, Y, X_test, Y_test, num_neurons, "ReLU")
        mses[t, 0] = mse

        mse = neural_network(X, Y, X_test, Y_test, num_neurons, "tanh")
        mses[t, 1] = mse

    print('Experiment with {} neurons done...'.format(num_neurons))
```

```
11.109407859943905
111.7986116835294
Experiment with 100 neurons done...
10.621046340776095
107.56930690987294
Experiment with 150 neurons done...
8.750118083604022
98.74964449460397
Experiment with 200 neurons done...
9.891959901212498
90.95010879210426
Experiment with 250 neurons done...
10.088487136454873
90.3092697911239
Experiment with 300 neurons done...
8.94582778492702
83.29509262423834
Experiment with 350 neurons done...
8.154480920424648
82.91806908671163
Experiment with 400 neurons done...
8.842727791373825
78.61146230571298
Experiment with 450 neurons done...
9.86359222938948
75.96352666746543
Experiment with 500 neurons done...
```



4.e

```
In [21]: activations = dict(ReLU=ReLUActivation,
                           tanh=TanhActivation,
                           linear=LinearActivation)
lr = dict(ReLU=0.0000007,tanh=0.00001,linear=0.005)

def neural_network(X, Y, X_test, Y_test, num_layers, activation):
    """
    This function performs neural network prediction.
    Input:
        X: independent variables in training data.
        Y: dependent variables in training data.
        X_test: independent variables in test data.
        Y_test: dependent variables in test data.
        num_layers: number of layers in neural network
        activation: type of activation, ReLU or tanh
    Output:
        mse: Mean square error on test data.
    """
    ## YOUR CODE HERE
    #####
    # Have K hidden layers and L neurons per layer.
    # Want total # weights to be 10000
    # X.shape[1]*L + (k-1)*(L*L) + L*Y.shape[1] = 10000
    # => (k-1)*L^2 + L*(X.shape + Y.shape[1]) - 10000 = 0
    # Let x = X.shape[1], y = Y.shape[1].
    # Quadratic formula:
    # L = (-x + y) +- sqrt((x+y)**2 - 4*(k-1)*(-10000)) / 2*(k-1)
    # Since expression on the right is positive and >= (x+1) for k >= 0,
    # only positive makes sense.
    k = num_layers
    x = X.shape[1]
    y = Y.shape[1]
    L = (-x + y) + np.sqrt((x + y)**2 + 40000*(k-1))/2*(k-1)
    L = int(L)

    model = Model(X.shape[1])
    for i in range(num_layers):
        model.addLayer(DenseLayer(L, activations[activation]))
    model.addLayer(DenseLayer(Y.shape[1], LinearActivation()))
    model.initialize(QuadraticCost())
    model.train(X,Y,500,GDOptimizer(eta=lr[activation])) # Learning rate of 0.02

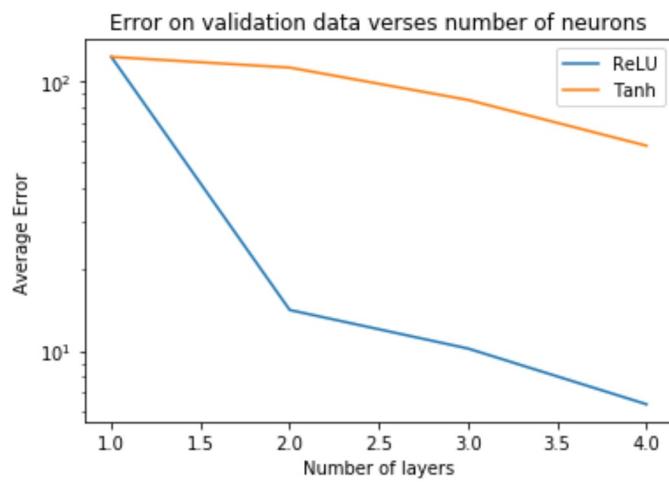
    Y_pred = model.predict(X_test)
    mse = np.mean(np.sqrt(np.sum((Y_pred - Y_test)**2, axis=1)))
    print(mse)
    return mse

#####
#####PLOT PART 2#####
#####

def generate_data(sensor_loc, k=7, d=2, n=1, original_dist=True, noise=1):
    return generate_dataset(
        sensor_loc,
        num_sensors=k,
        spatial_dim=d,
        num_data=n,
        original_dist=original_dist,
        noise=noise)

np.random.seed(0)
n = 200
num_layers = [1, 2, 3, 4]
```

```
123.34916925310644
123.37896279647254
Experiment with 1 layers done...
14.182293197206054
112.62096067500933
Experiment with 2 layers done...
10.202500813624377
85.39172088219182
Experiment with 3 layers done...
6.338764513589846
57.74814604922777
Experiment with 4 layers done...
```



4.f

```
In [27]: activations = dict(ReLU=ReLUActivation,
                           tanh=TanhActivation,
                           linear=LinearActivation)
lr = dict(ReLU=0.0000007,tanh=0.00001,linear=0.005)

def neural_network(X, Y, Xs_test, Ys_test):
    """
    This function performs neural network prediction.
    Input:
        X: independent variables in training data.
        Y: dependent variables in training data.
        X_test: independent variables in test data.
        Y_test: dependent variables in test data.
        num_layers: number of layers in neural network
        activation: type of activation, ReLU or tanh
    Output:
        mse: Mean square error on test data.
    """
    ## YOUR CODE HERE
    #####
    num_layers = 4
    activation = "ReLU"

    k = num_layers
    x = X.shape[1]
    y = Y.shape[1]
    L = (-(x + y) + np.sqrt((x + y)**2 + 40000*(k-1)))/2*(k-1)
    L = int(L)

    model = Model(X.shape[1])
    for i in range(num_layers):
        model.addLayer(DenseLayer(L, activations[activation]))
    model.addLayer(DenseLayer(Y.shape[1], LinearActivation()))
    model.initialize(QuadraticCost())
    model.train(X,Y,500,GDOptimizer(eta=lr[activation])) # Learning rate of 0.02

    Y_pred = model.predict(X_test)
    mse = np.mean(np.sqrt(np.sum((Y_pred - Y_test)**2, axis=1)))
    print(mse)
    return mse

def main():
    ##### PLOT PART 1#####
    np.random.seed(0)

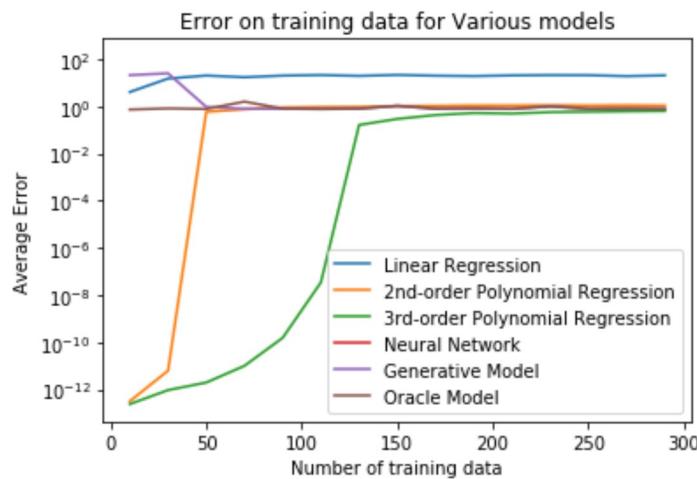
    ns = np.arange(10, 310, 20)
    replicates = 5
    num_methods = 6
    num_sets = 3
    mses = np.zeros((len(ns), replicates, num_methods, num_sets))

    def generate_data(sensor_loc, k=7, d=2, n=1, original_dist=True, noise=1):
        return generate_dataset(
            sensor_loc,
            num_sensors=k,
            spatial_dim=d,
            num_data=n,
            original_dist=original_dist,
            noise=noise)
```

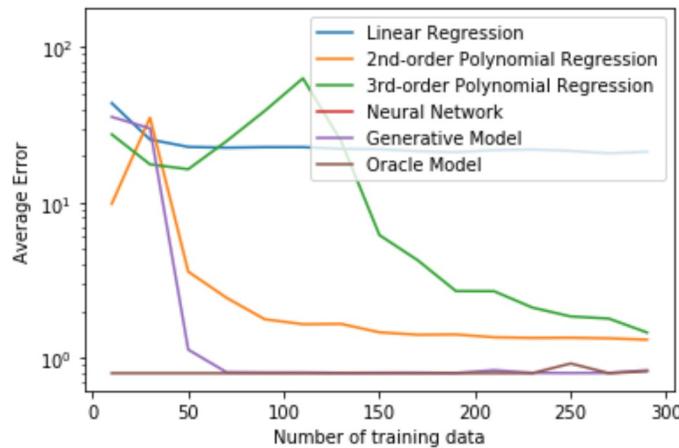
```
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:142: FutureWarning: `rcond` parameter will change to the default of machine precision time s ``max(M, N)`` where M and N are the input matrix dimensions.  
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.  
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:172: FutureWarning: `rcond` parameter will change to the default of machine precision time s ``max(M, N)`` where M and N are the input matrix dimensions.  
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.  
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:201: FutureWarning: `rcond` parameter will change to the default of machine precision time s ``max(M, N)`` where M and N are the input matrix dimensions.  
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.  
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:36: RuntimeWarning: overflow encountered in square  
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:224: RuntimeWarning: invalid value encountered in multiply  
/usr/local/lib/python3.5/dist-packages/numpy/core/_methods.py:32: RuntimeWarning: overflow encountered in reduce  
    return umr_sum(a, axis, dtype, out, keepdims)  
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:77: RuntimeWarning: invalid value encountered in greater
```

```
nan
0th Experiment with 10 samples done...
12.454824686649197
0th Experiment with 30 samples done...
nan
0th Experiment with 50 samples done...
nan
0th Experiment with 70 samples done...
20.41352178421555
0th Experiment with 90 samples done...
nan
0th Experiment with 110 samples done...
nan
0th Experiment with 130 samples done...
nan
0th Experiment with 150 samples done...
nan
0th Experiment with 170 samples done...
6.5905904523171195
0th Experiment with 190 samples done...
9.44426796016046
0th Experiment with 210 samples done...
8.066762207381036
0th Experiment with 230 samples done...
nan
0th Experiment with 250 samples done...
6.940402625030421
0th Experiment with 270 samples done...
7.2149912532420695
0th Experiment with 290 samples done...
187.06122990828078
1th Experiment with 10 samples done...
nan
1th Experiment with 30 samples done...
163.8230704844111
1th Experiment with 50 samples done...
nan
1th Experiment with 70 samples done...
145.32426742219718
1th Experiment with 90 samples done...
202.54499202115264
1th Experiment with 110 samples done...
nan
1th Experiment with 130 samples done...
175.0197396715583
1th Experiment with 150 samples done...
160.60171355672915
1th Experiment with 170 samples done...
nan
1th Experiment with 190 samples done...
175.11112759655146
1th Experiment with 210 samples done...
nan
1th Experiment with 230 samples done...
161.60692438275012
1th Experiment with 250 samples done...
nan
1th Experiment with 270 samples done...
195.2119073597044
1th Experiment with 290 samples done...
144.17268685961153
2th Experiment with 10 samples done...
nan
2th Experiment with 30 samples done...
```

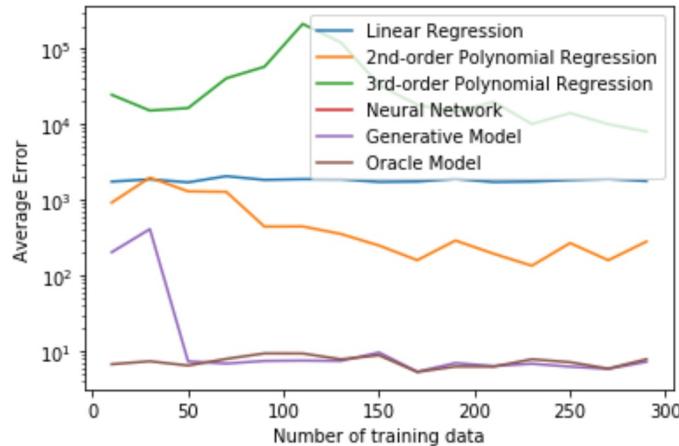
```
/usr/local/lib/python3.5/dist-packages/matplotlib/scale.py:114: RuntimeWarning:  
invalid value encountered in less_equal  
out[a <= 0] = -1000
```



Error on test data from the same distribution for Various models



Error on test data from a different distribution for Various models



HW09_Q4_Observations

4.C

Briefly describe your observations. What are the strengths and weaknesses of each Model?

Error on training data for various models:

Generative model has lowest average error, Oracle model quickly approaches this same error magnitude. Neural network performs the worst in terms of average error. The higher order polynomial regressions approach the same magnitude as generative and oracle.

With the exception of 2nd and 3rd order polynomial regression, average error stays relatively consistent as number of data training data increases.

Linear regression:

Error: Stable as number of training data increases
Same Dist: Stable as number of training data increases
Diff Dist: Stable as number of training data increases

Roughly same magnitude of error for all

2nd Order:

Error: Stable as number of training data increases
Same Dist: Peaks earlier than 3rd
Diff Dist: General trend towards decreasing error

3rd Order:

Error: Stable as number of training data increases
Same Dist: Error increases than decreases
Diff Dist: Multiple peaks in error, general trend towards decreases error

Neural Network:

Error: Stable as number of training data increases
Same Dist: Stable as number of training data increases
Diff Dist: Stable as number of training data increases

Generative Model:

Error: Lowest error magnitude

Same Dist: Lowest error magnitude

Diff Dist: Lowest error magnitude

Oracle Model:

Error: Lowest error magnitude

Same Dist: Lowest error magnitude

Diff Dist: Lowest error magnitude

4. D

Best choice for number of neurons is 200 or 400 for ReLu and is 500 for Tanh.

4.E

Best choice for number of layers is 2 or 4. For Tanh best is 4.