# CSCI 5521 (002) Final Exam

Justine John Serdoncillo

TOTAL POINTS

## 98 / 100

QUESTION 1

## Open-ended 25 pts

*1.1* (a) **3 / 5**

✓ **- 0 pts** *Correct*

**- 5 pts** Incorrect solution

**- 2.5 pts** Incorrect solution for classification

**- 2.5 pts** Partially incorrect solution

**- 2.5 pts** Incorrect solution for regression

✓ **- 2 pts** *The explanation for each metric is not enough*

**- 2 pts** No explanation of the differences between metrics

**- 2.5 pts** Missing explanation for classification

*1.2* (b) **4 / 4**

✓ **- 0 pts** *Correct*

**- 2 pts** Explanation is not enough

*1.3* (c) **4 / 4**

✓ **- 0 pts** *Correct*

**- 2 pts** Incorrect explanation

*1.4* (d) **4 / 4**

✓ **- 0 pts** *Correct*

**- 2 pts** Explanation is not enough

**- 4 pts** Explanation is incorrect

*1.5* (e) **4 / 4**

✓ **- 0 pts** *Correct*

**- 2 pts** Explanation is not enough

**- 4 pts** Answer is missing

*1.6* (f) **4 / 4**

✓ **- 0 pts** *Correct*

QUESTION 2

## Gaussian Distribution 25 pts

*2.1* (a) **10 / 10**

✓ **- 0 pts** *Correct*

**- 10 pts** Missing

**- 3 pts** Minor mistake

**- 5 pts** Major mistake

*2.2* (b) **7 / 7**

✓ **- 0 pts** *Correct*

*2.3* (c) **8 / 8**

✓ **- 0 pts** *Correct*

**- 2 pts** Miner mistake

**- 4 pts** Major mistake

QUESTION 3

## Random Forest 25 pts

*3.1* (a) **10 / 10**

✓ **- 0 pts** *Correct*

**- 2 pts** Miner mistake

*3.2* (b) **8 / 8**

✓ **- 0 pts** *Correct*

*3.3* (c) **7 / 7**

✓ **- 0 pts** *Correct*

# SVM 25 pts

*4.1* (a) **10 / 10**

✓ **- 0 pts** *Correct*

**- 10 pts** Missing

**- 5 pts** Major mistake

**- 3 pts** Miner mistake

*4.2* (b) **7 / 7**

✓ **- 0 pts** *Correct*

**- 7 pts** Missing

**- 4 pts** Major mistake

**- 2 pts** Miner mistake

*4.3* (c) **8 / 8**

✓ **- 0 pts** *Correct*

**- 8 pts** Missing

ıl gradescope

Justine John A. Serdoncillo
"JJ"

Final Exam

CSCI 5521:
Machine Learning
Fundamentals

December 19, 2022

1)

a)

| Classification | regression |
|---|---|
| – Accuracy | – mean absolute error |
| – top k accuracy | – mean squared error |
| – f₁ score | – R² |
| – log loss | – Mean absolute percent error |
| – precission | |
| – recall | |

○ Classification is based on the accuracy of the modeling based on the true and predicted classes which are discreet.

○ regression is based on the closeness of the predicted with the actual value. These are not as discreet but are commonly continuous. Both cases can be scaled to the size of the data set.

b) Ensemble methods use sometimes multiple classifiers. They can also introduce randomness which can reduce variance. Additionally by combining diverse different classifiers it can also reduce bias by being uncorrelated to each other. All of these methods can improve the performance of the model.

c) An estimator is always composed of a variance and bias. Some methods can tweak these two. Depending on the performance of the model the bias and variance can be traded to improve performance. For a model that is underfitting, methods can be done like increasing complexity to reduce bias & increase variance. On the other hand, overfitting can use regularization to increase bias but decrease variance. These methods can improve accuracy.

d) Cross-validation can be used by cross-validating the data set to itself. The common one is k-fold where it is split into sets that can be the training and test sets. By doing so, the model can be tested how well it works with unseen data and see its consistency.

e) Machine learning is highly dependent on the model and data. The data quantity and quality should be good enough for the Machine to learn important patterns. A trash model will also produce trash outputs which means that achieving near perfection needs a near perfect model as well. ML is not that easy to understand the inner workings and is hard to generalize to multiple purposes. These limitations can be overcome by proper education and better methods for gathering quality data.

f) Data ethics is really necessary because it keeps ML in check with the purposes and the implications it shows. The model can be used for unethical reasons. At the same time, without further guidance, biases in the world can seep into the model and result to unethical concerns and a closed feedback loop. All these can be ensured by consistent monitoring, diversity and repeated questions of the effects of the ML algorithms created.

P_1

*1.1* (a) **3 / 5**

✓ **- 0 pts** *Correct*

   **- 5 pts** Incorrect solution

   **- 2.5 pts** Incorrect solution for classification

   **- 2.5 pts** Partially incorrect solution

   **- 2.5 pts** Incorrect solution for regression

✓ **- 2 pts** *The explanation for each metric is not enough*

   **- 2 pts** No explanation of the differences between metrics

   **- 2.5 pts** Missing explanation for classification

Justine John A. Serdoncillo
"JJ"

Final Exam

CSCI 5521:
Machine Learning
Fundamentals

December 19, 2022

1)

a)

| Classification | regression |
|---|---|
| – Accuracy | – mean absolute error |
| – top k accuracy | – mean squared error |
| – f1 score | – $R^2$ |
| – log loss | – Mean absolute percent error |
| – precision | |
| – recall | |

○ Classification is based on the accuracy of the modeling based on the true and predicted classes which are discreet.

○ regression is based on the closeness of the predicted with the actual value. These are not as discreet but are commonly continuos. Both cases can be scaled to the size of the data set.

b) Ensemble methods use sometimes multiple classifiers. They can also introduce randomness which can reduce variance. Additionally by combining diverse different classifiers it can also reduce bias by being uncorrelated to each other. All of these methods can improve the performance of the model.

c) An estimator is always composed of a variance and bias. Some methods can tweak these two. Depending on the performance of the model the bias and variance can be traded to improve performance. For a model that is underfitting, methods can be done like increasing complexity to reduce bias & increase variance. On the other hand, overfitting can use regularization to increase bias but decrease variance. These methods can improve accuracy.

d) Cross-validation can be used by cross-validating the data set to itself. The common one is k-fold where it is split into sets that can be the training and test sets. By doing so, the model can be tested how well it works with unseen data and see its consistency.

e) Machine learning is highly dependent on the model and data. The data quantity and quality should be good enough for the Machine to learn important patterns. A trash model will also produce trash outputs which means that achieving near perfection needs a near perfect model as well. ML is not that easy to understand the inner workings and is hard to generalize to multiple purposes. These limitations can be overcome by proper education and better methods for gathering quality data.

f) Data ethics is really necessary because it keeps ML in check with the purposes and the implications it shows. The model can be used for unethical reasons. At the same time, without further guidance, biases in the world can seep into the model and result to unethical concerns and a closed feedback loop. All these can be ensured by consistent monitoring, diversity and repeated questions of the effects of the ML algorithms created.

P 1

*1.2* (b) **4 / 4**

   ✓ **- 0 pts** *Correct*

      **- 2 pts** Explanation is not enough

ᵢᵢᵢ gradescope

Justine John A. Serdoncillo
"JJ"
Final Exam

CSCI 5521:
Machine Learning
Fundamentals

December 19, 2022

1)

a)

| Classification | regression |
|---|---|
| – Accuracy | – mean absolute error |
| – top k accuracy | – mean squared error |
| – f₁ score | – $R^2$ |
| – log loss | – Mean absolute percent error |
| – precision | |
| – recall | |

o Classification is based on the accuracy of the modeling based on the true and predicted classes which are discreet.

o regression is based on the closeness of the predicted with the actual value. These are not as discreet but are commonly continuous. Both cases can be scaled to the size of the data set.

b) Ensemble methods use sometimes multiple classifiers. They can also introduce randomness which can reduce variance. Additionally by combining diverse different classifiers it can also reduce bias by being uncorrelated to each other. All of these methods can improve the performance of the model.

c) An estimator is always composed of a variance and bias. Some methods can tweak these two. Depending on the performance of the model the bias and variance can be traded to improve performance. For a model that is underfitting, methods can be done like increasing complexity to reduce bias & increase variance. On the other hand, overfitting can use regularization to increase bias but decrease variance. These methods can improve accuracy.

d) Cross-validation can be used by cross-validating the data set to itself. The common one is k-fold where it is split into sets that can be the training and test sets. By doing so, the model can be tested how well it works with unseen data and see its consistency.

e) Machine learning is highly dependent on the model and data. The data quantity and quality should be good enough for the Machine to learn important patterns. A trash model will also produce trash outputs which means that achieving near perfection needs a near perfect model as well. ML is not that easy to understand the inner workings and is hard to generalize to multiple purposes. These limitations can be overcome by proper education and better methods for gathering quality data.

f) Data ethics is really necessary because it keeps ML in check with the purposes and the implications it shows. The model can be used for unethical reasons. At the same time, without further guidance, biases in the world can seep into the model and result to unethical concerns and a closed feedback loop. All these can be ensured by consistent monitoring, diversity and repeated question of the effects of the ML algorithms created.

P 1

*1.3* (c) **4 / 4**

✓ **- 0 pts** *Correct*

**- 2 pts** Incorrect explanation

Justine John A. Serdoncillo
"JJ"
Final Exam

CSCI 5521:
Machine Learning
Fundamentals

December 19, 2022

1)

a)

| Classification | regression |
|---|---|
| – Accuracy | – mean absolute error |
| – top k accuracy | – mean squared error |
| – f1 score | – $R^2$ |
| – log loss | – Much absolute percent error |
| – precision | |
| – recall | |

o Classification is based on the accuracy of the modeling based on the true and predicted classes which are discreet.

o regression is based on the closeness of the predicted with the actual value. These are not as discreet but are commonly continuous. Both cases can be scaled to the size of the data set.

b) Ensemble methods use sometimes multiple classifiers. They can also introduce randomness which can reduce variance. Additionally by combining diverse different classifiers it can also reduce bias by being uncorrelated to each other. All of these methods can improve the performance of the model.

c) An estimator is always composed of a variance and bias. Some methods can tweak these two. Depending on the performance of the model the bias and variance can be traded to improve performance. For a model that is underfitting, methods can be done like increasing complexity to reduce bias & increase variance. On the other hand, overfitting can use regularization to increase bias but decrease variance. These methods can improve accuracy.

d) Cross-validation can be used by cross-validating the data set to itself. The common one is k-fold where it is split into sets that can be the training and test sets. By doing so, the model can be tested how well it works with unseen data and see its consistency.

e) Machine learning is highly dependent on the model and data. The data quantity and quality should be good enough for the Machine to learn important patterns. A trash model will also produce trash outputs which means that achieving near perfection needs a near perfect model as well. ML is not that easy to understand the inner workings and is hard to generalize to multiple purposes. These limitations can be overcome by proper education and better methods for gathering quality data.

f) Data ethics is really necessary because it keeps ML in check with the purposes and the implications it shows. The model can be used for unethical reasons. At the same time, without further guidance, biases in the world can seep into the model and result to unethical concerns and a closed feedback loop. All these can be ensured by consistent monitoring, diversity and repeated questions of the effects of the ML algorithms created.

P 1

*1.4* (d) **4 / 4**

  ✓ **- 0 pts** *Correct*

    **- 2 pts** Explanation is not enough

    **- 4 pts** Explanation is incorrect

ıl gradescope

Justine John A. Serdoncillo
"JJ"

CSCI 5521:
Machine Learning
Fundamentals

December 17, 2022

Final Exam

1)

a)

| Classification | regression |
|---|---|
| – Accuracy | – mean absolute error |
| – top k accuracy | – mean squared error |
| – f1 score | – $R^2$ |
| – log loss | – Mean absolute percent error |
| – precision | |
| – recall | |

○ Classification is based on the accuracy of the modeling based on the true and predicted classes which are discreet.

○ regression is based on the closeness of the predicted with the actual value. These are not as discreet but are commonly continuos. Both cases can be scaled to the size of the data set.

b) Ensemble methods use sometimes multiple classifiers. They can also introduce randomness which can reduce variance. Additionally by combining diverse different classifiers it can also reduce bias by being uncorrelated to each other. All of these methods can improve the performance of the model.

c) An estimator is always composed of a variance and bias. Some methods can tweak these two. Depending on the performance of the model the bias and variance can be traded to improve performance. For a model that is underfitting, methods can be done like increasing complexity to reduce bias & increase variance. On the other hand, overfitting can use regularization to increase bias but decrease variance. These methods can improve accuracy.

d) Cross-validation can be used by cross-validating the data set to itself. The common one is k-fold where it is split into sets that can be the training and test sets. By doing so, the model can be tested how well it works with unseen data and see its consistency.

e) Machine learning is highly dependent on the model and data. The data quantity and quality should be good enough for the Machine to learn important patterns. A trash model will also produce trash outputs which means that achieving near perfection needs a near perfect model as well. ML is not that easy to understand the inner workings and is hard to generalize to multiple purposes. These limitations can be overcome by proper education and better methods for gathering quality data.

f) Data ethics is really necessary because it keeps ML in check with the purposes and the implications it shows. The model can be used for unethical reasons. At the same time, without further guidance, biases in the world can seep into the model and result to unethical concerns and a closed feedback loop. All these can be ensured by consistent monitoring, diversity and repeated questions of the effects of the ML algorithms created.

*1.5* (e) **4 / 4**

&check; **- 0 pts** *Correct*

   **- 2 pts** Explanation is not enough

   **- 4 pts** Answer is missing

gradescope

Justine John A. Serdoncillo
"JJ"

Final Exam

1)

a)

| Classification | regression |
|---|---|
| – Accuracy | – mean absolute error |
| – top k accuracy | – mean squared error |
| – f1 score | – $R^2$ |
| – log loss | – Much absolute percent error |
| – precission | |
| – recall | |

○ Classification is based on the accuracy of the modeling based on the true and predicted classes which are discreet.

○ regression is based on the closeness of the predicted with the actual value. These are not as discreet but are commonly continuos. Both cases can be scaled to the size of the data set.

b) Ensemble methods use sometimes multiple classifiers. They can also introduce randomness which can reduce variance. Additionally by combining diverse different classifiers it can also reduce bias by being uncorrelated to each other. All of these methods can improve the performance of the model.

c) An estimator is always composed of a variance and bias. Some methods can tweak these two. Depending on the performance of the model the bias and variance can be traded to improve performance. For a model that is underfitting, methods can be done like increasing complexity to reduce bias & increase variance. On the other hand, overfitting can use regularization to increase bias but decrease variance. These methods can improve accuracy.

d) Cross-validation can be used by cross-validating the data set to itself. The common one is k-fold where it is split into sets that can be the training and test sets. By doing so, the model can be tested how well it works with unseen data and see its consistency.

e) Machine learning is highly dependent on the model and data. The data quantity and quality should be good enough for the Machine to learn important patterns. A trash model will also produce trash outputs which means that achieving near perfection needs a near perfect model as well. ML is not that easy to understand the inner workings and is hard to generalize to multiple purposes. These limitations can be overcome by proper education and better methods for gathering quality data.

f) Data ethics is really necessary because it keeps ML in check with the purposes and the implications it shows. The model can be used for unethical reasons. At the same time, without further guidance, biases in the world can seep into the model and result to unethical concerns and a closed feedback loop. All these can be ensured by consistent monitoring, diversity and repeated questions of the effects of the ML algorithms created.

P1

*1.6* (f) **4 / 4**

  ✓ **- 0 pts** *Correct*

2) $\underline{X} = \{x_1, x_2 \ldots x_n\}$ $\quad x_i \in \mathbb{R}$ $\qquad M \in \mathbb{R}$ $\quad \sigma^2 \in \mathbb{R}_{++}$

a) $P(x) = \dfrac{1}{\sqrt{2\pi}\,\sigma}\exp\left[-\dfrac{(x-M)^2}{2\sigma^2}\right]$

$\ell(M,\sigma^2 \mid x) = \displaystyle\prod_{t=1}^{N} P(x \mid M,\sigma^2) = \prod_{t=1}^{N}\dfrac{1}{\sqrt{2\pi}\,\sigma}\exp\left[-\dfrac{\sum_{i=1}^{N}(x_i-M)^2}{2\sigma^2}\right]$

$\mathcal{I}(M,\sigma^2 \mid x) = \log \ell(M,\sigma^2 \mid x) = -\dfrac{N}{2}(\log 2\pi) - N\log\sigma - \dfrac{\sum_{i=1}^{n}(x_i-M)^2}{2\sigma^2}$

$\dfrac{\partial \mathcal{I}(M,\sigma^2 \mid x)}{\partial M} = -\left(\dfrac{2\sum_{i=1}^{n}x_i - 2MN}{2\sigma^2}\right) = 0$

$\dfrac{\partial \mathcal{I}(M,\sigma^2 \mid x)}{\partial \sigma} = \dfrac{-N}{\sigma} + \dfrac{2\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{2\sigma^3} = 0$

$2\sum_{i=1}^{n}x_i - 2MN = 0$

$\sum_{i=1}^{n}x_i = MN$

$M = \dfrac{\sum_{i=1}^{n}x_i}{N}$

$\boxed{\hat{M}_n = \dfrac{\sum_{i=1}^{n}x_i}{N}}$

$\dfrac{N}{\sigma} = \dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{\sigma^{-8}}$

$\sigma^2 = \dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{N}$

$\boxed{\hat{\sigma}_n^2 = \dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{N}}$

b) $E[\hat{M}_n] = E\left[\dfrac{\sum_{i=1}^{n}x_i}{N}\right] = \dfrac{1}{N}E\left[\sum_{i=1}^{n}x_i\right]$

$= \dfrac{1}{N}\sum_{i=1}^{n}E[x_i] = \dfrac{1}{N}NM = M$

$\boxed{E[\hat{M}_n] = M}$

it is true because for an unbiased $M$ it equals $M$

c) $E[\hat{\sigma}_n^2] = E\left[\dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{N}\right]$

$= \dfrac{1}{N}\left(E\left[\sum_{i=1}^{n}(x_i)^2 - N\hat{M}_n^2\right]\right)$

$= \dfrac{1}{N}\left(\underbrace{\sum_{i=1}^{n}E[(x_i)^2]}_{A} - \underbrace{N\,E[\hat{M}_n^2]}_{B}\right)$

by definition
$Var[x] = E[x^2] - E[x]^2$

$\hookrightarrow E[x^2] = Var[x] + E[x]^2$

(A) $E[(x_i)^2] = Var[x_i] + E[x_i]^2$

$E[(x_i)^2] = \sigma^2 + M^2$

*2.1* (a) **10 / 10**

✓ **- 0 pts** *Correct*

**- 10 pts** Missing

**- 3 pts** Minor mistake

**- 5 pts** Major mistake

2) $\underline{X} = \{x_1, x_2 \dots x_n\}$  $x_i \in \mathbb{R}$   $M \in \mathbb{R}$   $\sigma^2 \in \mathbb{R}_{++}$

a) $P(x) = \dfrac{1}{\sqrt{2\pi}\,\sigma} \exp\left[-\dfrac{(x-M)^2}{2\sigma^2}\right]$

$l(M, \sigma^2 \mid x) = \displaystyle\prod_{t=1}^{N} P(x \mid M, \sigma^2) = \prod_{t=1}^{N} \dfrac{1}{\sqrt{2\pi}\,\sigma} \exp\left[-\dfrac{\sum_{i=1}^{N}(x_i-M)^2}{2\sigma^2}\right]$

$\mathcal{L}(M, \sigma^2 \mid x) = \log l(M, \sigma^2 \mid x) = -\dfrac{N}{2}(\log 2\pi) - N\log\sigma - \dfrac{\sum_{i=1}^{n}(x_i-M)^2}{2\sigma^2}$

$\dfrac{\partial \mathcal{L}(M, \sigma^2 \mid x)}{\partial M} = -\left(\dfrac{2\sum_{i=1}^{n} x_i - 2MN}{2\sigma^2}\right) = 0$

$\dfrac{\partial \mathcal{L}(M, \sigma^2 \mid x)}{\partial \sigma} = -\dfrac{N}{\sigma} + \dfrac{2\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{2\sigma^3} = 0$

$2\sum_{i=1}^{n} x_i - 2MN = 0$

$\sum_{i=1}^{n} x_i = MN$

$M = \dfrac{\sum_{i=1}^{n} x_i}{N}$

$\boxed{\hat{M}_n = \dfrac{\sum_{i=1}^{n} x_i}{N}}$

$\dfrac{N}{\sigma} = \dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{\sigma^3}$

$\sigma^2 = \dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{N}$

$\boxed{\hat{\sigma}_n^2 = \dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{N}}$

b) $E[\hat{M}_n] = E\left[\dfrac{\sum_{i=1}^{n} x_i}{N}\right] = \dfrac{1}{N} E\left[\sum_{i=1}^{n} x_i\right]$

$= \dfrac{1}{N} \sum_{i=1}^{n} E[x_i] = \dfrac{1}{N} NM = M$

$\boxed{E[\hat{M}_n] = M}$   it is true because for an unbiased $M$ it equals $M$

c) $E[\hat{\sigma}_n^2] = E\left[\dfrac{\sum_{i=1}^{n}(x_i-\hat{M}_n)^2}{N}\right]$

$= \dfrac{1}{N}\left(E\left[\sum_{i=1}^{n}(x_i)^2 - N\hat{M}_n^2\right]\right)$

$= \dfrac{1}{N}\left(\underbrace{\sum_{i=1}^{n} E[(x_i)^2]}_{\text{(A)}} - \underbrace{N E[\hat{M}_n^2]}_{\text{(B)}}\right)$

by definition

$Var[x] = E[x^2] - E[x]^2$

$\hookrightarrow E[x^2] = Var[x] + E[x]^2$

(A) $E[(x_i)^2] = Var[x_i] + E[x_i]^2$

$E[(x_i)^2] = \sigma^2 + M^2$

*2.2* (b) **7 / 7**

✓ **- 0 pts** *Correct*

2) $\underline{X} = \{x_1, x_2 \ldots x_n\}$   $x_i \in \mathbb{R}$    $\mu \in \mathbb{R}$   $\sigma^2 \in \mathbb{R}_{++}$

a) $P(x) = \dfrac{1}{\sqrt{2\pi}\,\sigma} \exp\left[-\dfrac{(x-\mu)^2}{2\sigma^2}\right]$

$\ell(\mu, \sigma^2 \mid x) = \prod_{t=1}^{N} P(x \mid \mu, \sigma^2) = \prod_{t=1}^{N} \dfrac{1}{\sqrt{2\pi}\,\sigma} \exp\left[-\dfrac{\sum_{i=1}^{N}(x_i-\mu)^2}{2\sigma^2_j}\right]$

$\mathcal{L}(\mu, \sigma^2 \mid x) = \log \ell(\mu, \sigma^2 \mid x) = -\dfrac{N}{2}(\log 2\pi) - N\log\sigma - \dfrac{\sum_{i=1}^{n}(x_i-\mu)^2}{2\sigma^2}$

$\dfrac{\partial \mathcal{L}(\mu, \sigma^2 \mid x)}{\partial \mu} = -\left(\dfrac{2\sum_{i=1}^{n}x_i - 2\mu N}{2\sigma^2}\right) = 0$     $\dfrac{\partial \mathcal{L}(\mu, \sigma^2 \mid x)}{\partial \sigma} = \dfrac{-N}{\sigma} + \dfrac{2\sum_{i=1}^{n}|x_i - \hat{\mu}_n|^2}{2\sigma^3} = 0$

$2\sum_{i=1}^{n} x_i - 2\mu N = 0$       $\dfrac{N}{\sigma} = \dfrac{\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2}{\sigma^{-8=}}$

$\sum_{i=1}^{n} x_i = \mu N$

$\mu = \dfrac{\sum_{i=1}^{n} x_i}{N}$    $\boxed{\hat{\mu}_n = \dfrac{\sum_{i=1}^{n} x_i}{N}}$    $\sigma^2 = \dfrac{\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2}{N}$

$\boxed{\hat{\sigma}_n^2 = \dfrac{\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2}{N}}$

b) $E[\hat{\mu}_n] = E\left[\dfrac{\sum_{i=1}^{n} x_i}{N}\right] = \dfrac{1}{N} E\left[\sum_{i=1}^{n} x_i\right]$

$= \dfrac{1}{N}\sum_{i=1}^{n} E[x_i] = \dfrac{1}{N} N\mu = \mu$

$\boxed{E[\hat{\mu}_n] = \mu}$    it is true because for an unbiased $\mu$ it equals $\mu$

c) $E[\hat{\sigma}_n^2] = E\left[\dfrac{\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2}{N}\right]$

$= \dfrac{1}{N}\left(E\left[\sum_{i=1}^{n}(x_i)^2 - N\hat{\mu}_n^2\right]\right)$

$= \dfrac{1}{N}\left(\underbrace{\sum_{i=1}^{n} E[(x_i)^2]}_{\text{(A)}} - N\underbrace{E[\hat{\mu}_n^2]}_{\text{(B)}}\right)$

by definition
$\text{Var}[x] = E[x^2] - E[x]^2$
$\hookrightarrow E[x^2] = \text{Var}[x] + E[x]^2$

(A) $E[(x_i)^2] = \text{Var}[x_i] + E[x_i]^2$
$E[(x_i)^2] = \sigma^2 + \mu^2$

(B) $E[\hat{\mu}_n^2] = Var[\hat{\mu}_n] + E[\hat{\mu}_n]^2$

$Var[\hat{\mu}_n] = Var\left(\frac{\sum_{i=1}^{N} x_i}{N}\right)$

$= \frac{1}{N^2} \sum_{i=1}^{n} Var(x_i) = \frac{1}{N^2} N\sigma^2 = \frac{\sigma^2}{N}$

$E[\hat{\mu}_n^2] = \frac{\sigma^2}{N} + \mu^2$

$E[\hat{\sigma}_n^2] = \frac{1}{N}\left(N(\sigma^2 + \mu^2) - N\left(\frac{\sigma^2}{N} + \mu^2\right)\right)$

$= \sigma^2 + \mu^2 - \frac{\sigma^2}{N} - \mu^2 = \sigma^2\left(1 - \frac{1}{N}\right)$

$$\boxed{E[\hat{\mu}_n^2] = \sigma^2\left(1 - \frac{1}{N}\right) \neq \sigma^2}$$

the two are not equal to each other

- - - - - - - - - - - -

3) a) Random Forests are trained by training multiple decision trees that can be grown by only a few samples of training data or few attributes. The multiple trees are then compared to each other (and combined) to do a prediction done a test point. The prediction is done using the majority vote of features from the different Decision Trees. The two key hyper-parameters for RFs are

① Number of decision trees/samples

② Number of sample m of attributes to train on each node

b)

Pros of Random Forests:

- easy to code w/ fairly high accuracy
- lower chance of overfitting
- can handle weird data that can mess up training
- still has the advantages of normal decision trees, like setting key features

Cons of Random Forests:

- hard to understand the inner workings or interpret
- can be computationally expensive if generating lots of trees and takes a long time
- extra hyperparameters to tweak

c) If the number of features is a lot more than m=1, the decision tree may capture unimportant features which can be a problem because of possible pruning - This makes a lot of useless trees. If there is no pruning, this can still be a problem because in the end it is compiled using a majority vote which won't make sense because m=1 which means no other feature to compare with and is purely random.

*2.3* (c) **8 / 8**

    ✓ **- 0 pts** *Correct*

      **- 2 pts** Miner mistake

      **- 4 pts** Major mistake

$(B)$ $E[\hat{\mu}_n^2] = Var[\hat{\mu}_n] + E[\hat{\mu}_n]^2$

$Var[\hat{\mu}_n] = Var\left(\dfrac{\sum_{i=1}^{N} x_i}{N}\right)$

$= \dfrac{1}{N^2} \sum_{i=1}^{n} Var(x_i) = \dfrac{1}{N^2} N\sigma^2 = \dfrac{\sigma^2}{N}$

$E[\hat{\mu}_n^2] = \dfrac{\sigma^2}{N} + \mu^2$

$E[\hat{\sigma}_n^2] = \dfrac{1}{N}\left(N(\sigma^2 + \mu^2) - N\left(\dfrac{\sigma^2}{N} + \mu^2\right)\right)$

$= \sigma^2 + \mu^2 - \dfrac{\sigma^2}{N} - \mu^2 = \sigma^2\left(1 - \dfrac{1}{N}\right)$

$$\boxed{E[\hat{\mu}_n^2] = \sigma^2\left(1 - \dfrac{1}{N}\right) \neq \sigma^2}$$

the two are not equal to each other

- - - - - - - - - - -

3) a) Random Forests are trained by training multiple decision trees that can be grown by only a few samples of training data or few attributes. The multiple trees are then and combined compared to each other to do a prediction done a test point. The prediction is done using the majority vote of features from the different Decision Trees. The two key hyper-parameters for RFs are

① Number of decision trees/samples

② Number of sample m of attributes to train on each node

b)

Pros of Random Forests:

- easy to code w/ fairly high accuracy
- lower chance of overfitting
- can handle weird data that can mess up training
- still has the advantages of normal decision trees, like setting key features

Cons of Random Forests:

- hard to understand the inner workings or interpret
- can be computationally expensive if generating lots of trees and takes a long time
- extra hyperparameters to tweak

c) If the number of features is a lot more than $m=1$, the decision tree may capture unimportant features which can be a problem because of possible pruning - This makes a lot of useless trees. If there is no pruning, this can still be a problem because in the end it is compiled using a majority vote which won't make sense because $m=1$ which means no other feature to compare with and is purely random.

*3.1* (a) **10 / 10**

✓ **- 0 pts** *Correct*

**- 2 pts** Miner mistake

(B) $E[\hat{\mu}_n^2] = Var[\hat{\mu}_n] + E[\hat{\mu}_n]^2$

$Var[\hat{\mu}_n] = Var\left(\dfrac{\sum_{i=1}^{N} x_i}{N}\right)$

$= \dfrac{1}{N^2} \sum_{i=1}^{n} Var(x_i) = \dfrac{1}{N^2} N\sigma^2 = \dfrac{\sigma^2}{N}$

$E[\hat{\mu}_n^2] = \dfrac{\sigma^2}{N} + \mu^2$

$E[\hat{\sigma}_n^2] = \dfrac{1}{N}\left( N(\sigma^2 + \mu^2) - N\left(\dfrac{\sigma^2}{N} + \mu^2\right)\right)$

$= \sigma^2 + \mu^2 - \dfrac{\sigma^2}{N} - \mu^2 = \sigma^2\left(1 - \dfrac{1}{N}\right)$

$$\boxed{E[\hat{\mu}_n^2] = \sigma^2\left(1 - \dfrac{1}{N}\right) \neq \sigma^2}$$

the two are not equal to each other

---

3) a) Random Forests are trained by training multiple decision trees that can be grown by only a few samples of training data or few attributes. The multiple trees are then and combined compared to each other to do a prediction Vdone a test point. The prediction is done using the majority vote of features from the different Decision Trees. The two key hyper-parameters for RFs are

① Number of decision trees/samples

② Number of sample m of attributes to train on each node

b)

Pros of Random Forests:

- easy to code w/ fairly high accuracy
- lower chance of overfitting
- Can handle weird data that can mess up training
- still has the advantages of normal decision trees, like setting key features

Cons of Random Forests:

- hard to understand the inner workings or interpret
- can be computationally expensive if generating lots of trees and takes a long time
- extra hyperparameters to tweak

c) If the number of features is a lot more than m=1, the decision tree may capture unimportant features which can be a problem because of possible pruning - This makes a lot of useless trees. If there is no pruning, this can still be a problem because in the end it is compiled using a majority vote which won't make sense because m=1 which means no other feature to compare with and is purely random.

*3.2* (b) **8 / 8**

✓ **- 0 pts** *Correct*

(B) $E[\hat{\mu}_n^2] = Var[\hat{\mu}_n] + E[\hat{\mu}_n]^2$

$Var[\hat{\mu}_n] = Var\left(\dfrac{\overset{N}{\underset{i=1}{\Sigma}} x_i}{N}\right)$

$= \dfrac{1}{N^2} \overset{n}{\underset{i=1}{\Sigma}} Var(x_i) = \dfrac{1}{N^2} N \sigma^2 = \dfrac{\sigma^2}{N}$

$E[\hat{\mu}_n^2] = \dfrac{\sigma^2}{N} + \mu^2$

$E[\hat{\sigma}_n^2] = \dfrac{1}{N}\left(N(\sigma^2 + \mu^2) - N\left(\dfrac{\sigma^2}{N} + \mu^2\right)\right)$

$= \sigma^2 + \mu^2 - \dfrac{\sigma^2}{N} - \mu^2 = \sigma^2\left(1 - \dfrac{1}{N}\right)$

$$\boxed{E[\hat{\mu}_n^2] = \sigma^2\left(1 - \dfrac{1}{N}\right) \neq \sigma^2}$$

the two are not equal to each other

- - - - - - - - - - - -

3) a) Random Forests are trained by training multiple decision trees that can be grown by only a few samples of training data or few attributes. The multiple trees are then and combined compared to each other to do a prediction done a test point. The prediction is done using the majority vote of features from the different Decision Trees. The two key hyper-parameters for RFs are

① Number of decision trees/samples

② Number of sample m of attributes to train on each node

b)

Pros of Random Forests:

- easy to code w/ fairly high accuracy
- lower chance of overfitting
- can handle weird data that can mess up training
- still has the advantages of normal decision trees, like setting key features

Cons of Random Forests:

- hard to understand the inner workings or interpret
- can be computationally expensive if generating lots of trees and takes a long time
- extra hyperparameters to tweak

c) If the number of features is a lot more than m=1, the decision tree may capture unimportant features which can be a problem because of possible pruning - This makes a lot of useless trees. If there is no pruning, this can still be a problem because in the end it is compiled using a majority vote which won't make sense because m=1 which means no other feature to compare with and is purely random.

*3.3* (c) **7 / 7**

    ✓ **- 0 pts** *Correct*

4) $\underline{Z} = \{(\underline{x}^1, r^1), \ldots (\underline{x}^n, r^n)\}$  $\underline{x}^i \in \mathbb{R}^d$ and $r^i \in \{-1, 1\}$

Linear support vector machines (SVMs)

$$\min_{\underline{w}, w_0 \, \xi^i \, i=1 \ldots n} \frac{1}{2} \|\underline{w}\|^2 + C \sum_{i=1}^n \xi^i \text{ such that } r^i(\underline{w}^T \underline{x}^i + w_0) \geq 1 - \xi^i$$

$$\xi^i \geq 0 \quad i = 1 \ldots n$$

a) minimize the hinge loss on $(\underline{w}, w_0)$

$$\min_{\underline{w}, w_0} \sum_{i=1}^n \max\left(0, 1 - r^i[\underline{w}^T \underline{x}_i + w_0]\right) + \frac{\lambda}{2} \|\underline{w}\|^2$$

let $\xi^i = \max\left(0, 1 - r^i(\underline{w}^T \underline{x}_i + w_0)\right)$

$$\Rightarrow \min_{\underline{w}, w_0, \xi^i} \sum_{i=1}^n \xi^i + \frac{\lambda}{2} \|\underline{w}\|^2 = \frac{1}{\lambda} \sum_{i=1}^n \xi^i + \frac{1}{2} \|\underline{w}\|^2$$

let $C = \frac{1}{\lambda}$

$$\Rightarrow \min_{\underline{w}, w_0, \xi_i} \frac{1}{2} \|\underline{w}\|^2 + C \sum_{i=1}^n \xi^i$$

$\xi^i \geq 1 - r^i[\underline{w}^T \underline{x}_i + w_0]$

$\xi^i - 1 \geq -r^i[\underline{w}^T \underline{x}_i + w_0]$

$r^i[\underline{w}^T \underline{x}_i + w_0] \geq 1 - \xi^i$

how this is the same as above
but going backwards

w/ $C = \frac{1}{\lambda}$  &  $\xi^i = \max(0, 1 - r^i(\underline{w}^T \underline{x}_i + w_0))$

to produce a hinge loss optimization
problem & a
regularization term

b) Pseudo code

① initialize parameters to be used:
  - batch size, learning rate, #of epochs
  or loss tolerance when stopping

② train the SVM per mini batch

I thought mini-batch & Stochastic Gradient
Descent is different tho

run all epochs

  run all mini batches

    calculate $y = ?(\underline{w}^T \underline{x}^i + w_0)$

    if $y \neq r$

      keep w

    else

      calculate gradient
      $w = w - (\text{gradient})(\text{learning rate})$

done per batch

P4

*4.1* (a) **10 / 10**

   ✓ **- 0 pts** *Correct*

      **- 10 pts** Missing

      **- 5 pts** Major mistake

      **- 3 pts** Miner mistake

ili gradescope

4) $\underline{Z} = \{(\underline{x}', r'), \dots (\underline{x}^n, r^n)\}$ $\underline{x}^i \in \mathbb{R}^d$ and $r^i \in \{-1, 1\}$

Linear support vector machines (SVMs)

$$\min_{\underline{w}, w_0 \, \xi^i i = 1 \dots n} \frac{1}{2} \|\underline{w}\|^2 + C \sum_{i=1}^{n} \xi^i \text{ such that } r^i(\underline{w}^T \underline{x}^i + w_0) \geq 1 - \xi^i$$

$$\xi^i \geq 0 \quad i = 1 \dots n$$

a) minimize the hinge loss on $(\underline{w}, w_0)$

$$\min_{\underline{w}, w_0} \sum_{i=1}^{n} \max\left(0, 1 - r^i[\underline{w}^T \underline{x}_i + w_0]\right) + \frac{\lambda}{2} \|\underline{w}\|^2$$

let $\xi^i = \max\left(0, 1 - r^i(\underline{w}^T \xi_i + w_0)\right)$

$$\Rightarrow \min_{\underline{w}, w_0, \xi^i} \sum_{i=1}^{n} \xi^i + \frac{\lambda}{2} \|w\|^2 = \frac{1}{\lambda} \sum_{i=1}^{n} \xi^i + \frac{1}{2} \|w\|^2$$

let $C = \frac{1}{\lambda}$

$$\Rightarrow \min_{\underline{w}, w_0, \xi_i} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{n} \xi^i$$

$$\xi^i \geq 1 - r^i[\underline{w}^T \xi_i + w_0]$$

$$\xi^i - 1 \geq -r^i[\underline{w}^T \underline{x}_i + w_0]$$

$$r^i[\underline{w}^T \underline{x}_i + w_0] \geq 1 - \xi^i$$

how this is the same as above
but going backwards

w/ $C = \frac{1}{\lambda}$ & $\xi^i = \max(0, 1 - r^i(\underline{w}^T \underline{x}_i + w_0))$

to produce a hinge loss optimization
problem & a
regularization term

b) Pseudo code

① initialize parameters to be used:
   - batch size, learning rate, #of epochs
     or loss tolerance when stopping

② train the SVM per mini batch

I thought mini-batch & Stochastic Gradient
                                    Descent is
                                    different
                                    tho

run all epochs

  run all mini batches

   calculate $y = ?(\underline{w}^T \underline{x}^i + \underline{w}_0)$

done    if $y \neq r$
per batch     keep w

   else
    calculate gradient
    $w = W - (\text{gradient})(\text{learning rate})$

*4.2* (b) **7 / 7**

   ✓ **- 0 pts** *Correct*

     **- 7 pts** Missing

     **- 4 pts** Major mistake

     **- 2 pts** Miner mistake

ıll gradescope

c) $d=2$ $\underline{x}^i = (x_i, x_2^i)$

$$r^i = \begin{cases} +1 & \text{if } (x_i + 2x_2^i) \leq 1 \\ -1 & \text{otherwise} \end{cases}$$

$$r^i \begin{cases} +1 & 0 \leq 1 - (x_i + 2x_2^i)^2 \\ -1 & \text{otherwise} \end{cases}$$

$$r^i = \text{sign}\left[ 1 - (x_i + 2x_2^i)^2 \right]$$

map $\underline{x}^i$ to $\phi(\underline{x}^i)$ and train the linear SVM

Linear SVM

$$\phi(\underline{x}^i) = \begin{bmatrix} 1 & x_i^i & x_2^i & x_i^i x_2^i & (x_i^i)^2 & (x_2^i)^2 \end{bmatrix}^T$$
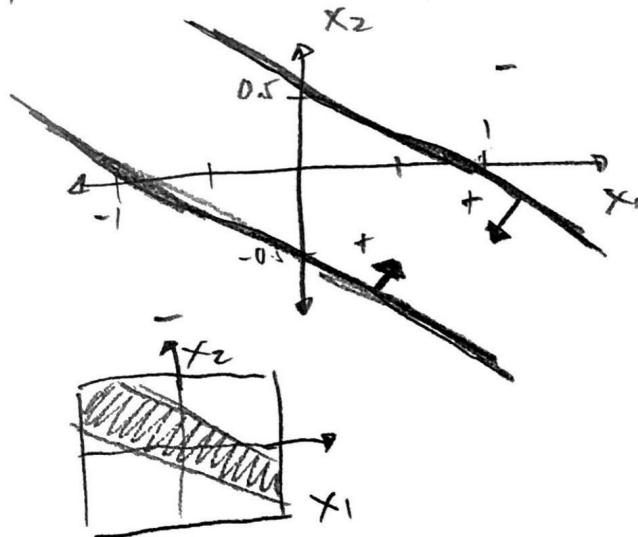
$$\min \frac{1}{2}\|\underline{w}\|^2 + C \sum_t \xi^t$$

$$r^t(\underline{w}^T \underline{x}^t + \underline{w}_0) \geq 1 - \xi^i$$

SVM solution

$$\sum_t \alpha^t r^t = 0 \quad 0 \leq \alpha^t \leq \frac{1}{N}$$

$$w = \sum_t \alpha^t r^t \phi(\underline{x}^t)$$

$$(x_i + 2x_2^i)^2 \leq 1$$

$$g(x) = w^T \phi(x) = \sum_t \alpha^t r^t \phi(\underline{x}^t)^T \phi(x)$$

$$-1 \leq x_i + 2x_2^i \leq 1 \quad \left.\right\} \text{equal to this}$$

$$\phi(\underline{x}^t)^T \phi(x) = K(\underline{x}^t, x)$$



$\phi(x)$ can be a highly accurate predictor
if the kernel trick works

$$\phi(\underline{x}^t)^T \phi(x)$$

$$= 1 + x_i^t x_1 + x_2^t x_2 + x_i^t x_2^t x_1 x_2$$

$$+ x_i^{t^2} x_1^2 + x_2^{t^2} x_2^2$$



this looks close to the vectorial kernel

with $q=2$ $\quad K(\underline{x}^t, x) = (x^{t^T} x + 1)^2$

given the polynomial kernel of degree 2 and the given boundary of the labels, it is possible that using this mapping can be a highly accurate predictor because it can split it into an outside and inside boundary

yes

*4.3* (c) **8 / 8**

✓ **- 0 pts** *Correct*

**- 8 pts** Missing