

Recitation: Final Exam Solutions

Shu-Hao Yu (Q1, Q8)

William Wang (Q2, Q4)

Avi Dubey (Q3)

Pengtao Xie (Q5, Q9)

Guanyu Wang (Q6, Q7, Q10)

CMU School of Computer Science

**Carnegie
Mellon
University**

Agenda

- final exam solutions
- final exam statistics

Final Exam Solutions

True/False

1. [2pts] Linear regression can only deal with linear relationships.

False, by using non-linear basis functions, we can deal with non-linear data.

2. [2pts] Theoretically the normal equation method and the least-mean-square (LMS) method should generate the same results for linear regression problem.

True, LMS can converge to the same solution as the normal equation method.

3. [2pts] The dual form of SVM solves the exact same problem as primary using a different representation.

True, they are different formulations of the same problem.

4. [2pts] The Dirichlet distribution is a conjugate prior for the multinomial distribution.

True, if the prior of the multinomial is Dirichlet, then the posterior is also Dirichlet.

True/False

5. [2pts] It is impossible to have infinite VC dimension for any set of functions.

False, $\sin(x)$ and $\cos(x)$ have infinite VC dimensions.

6. [2pts] When a dataset is not linearly separable, no SVM classifiers will find a hyperplane that separates the data points perfectly.

False, we can use kernel trick to deal with non-linear data.

7. [2pts] Logistic Regression is a linear classifier.

True, logistic regression also models $\langle W, X \rangle$.

8. [2pts] For a given fixed set of data points and a fixed k , k-means always converges to the same clustering of the data.

False, the initialization have impacts on the final results.

Multiple Choice

2.1 In the k-nearest neighbor classifier, which of the following statements are true? [2pts]

- A kNN is a supervised classifier. 😊
- B The hyperparameter k in kNN is typically set to an odd number. 😊
- C When k is set to an extremely large number, it is more likely that the classifier will overfit than underfit.
- D Both kNN and k-means are unsupervised learning techniques.

supervised

underfit

Multiple Choice

SGD will not converge if the objective function is convex but non-differentiable

2.2 Which of the following statements about stochastic gradient descent (SGD) are true? [2pts]

- A If an objective function has only one global optimum, then SGD is guaranteed to converge to the global optimum.
- B The standard SGD (excluding the mini-batch variant) is an online algorithm that updates the weights by looking at one example at a time.
- C In practice, when we set an inappropriate learning rate in SGD, it might take longer time for SGD to converge.
- D SGD can also be used to train a conditional random field (CRF) model.



Multiple Choice

2.3 Which of the following statements about the hidden Markov model (HMM) are true? [2pts]

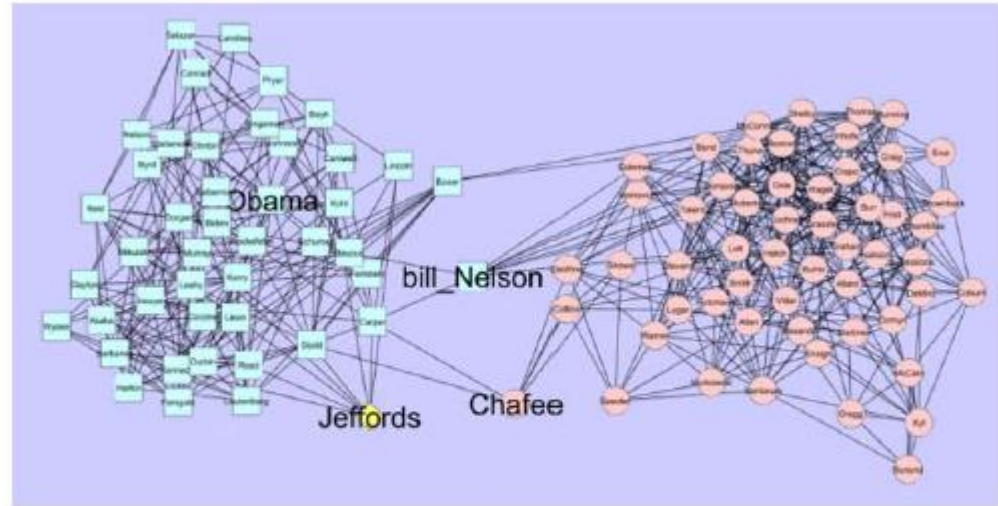
- A The forward-backward algorithm is a standard inference algorithm for HMM. 😊
- B Viterbi algorithm is a sum product algorithm for finding the most likely sequence of hidden states in a HMM.
- C HMM is a generative sequence modeling technique. 😊

max



Multiple Choice

2.4 Which of the following models would be possible for use with the data shown in figure 1 given that you have no labels for the nodes? [2pts]



The data shows the voting record for 100 senator in march 2005.

K-means doesn't model the links

A k-means

No labels

B Naive Bayes

C Stochastic block models (SBM) with 100 blocks, one for each senator.

D mixed membership SBM with 100 blocks, one for each senator.

E SBM with 2 blocks, one for each party

F mixed-membership SBM with 2 blocks, one for each party.

Figure 1: Network data.

Will generate 100 clusters

Short Answers

3 Short answers [32pts]

1. [2pts] Given samples from a uni-variate Gaussian distribution with mean 0 and variance 1, how could you create samples from a bi-variate Gaussian with means μ_1, μ_2 and co-variance

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$$

Assume Y is drawn from $N(0,1)$, to generate X from bi-variate Gaussian:

- Take $\mathbf{X} = \mathbf{S}\mathbf{Y} + \boldsymbol{\mu}_X$ where $\mathbf{C}_X = \mathbf{S}\mathbf{S}^T$ (Cholesky factorization)

$$\mathbf{C}_X = E\{(\mathbf{X} - \boldsymbol{\mu}_X)(\mathbf{X} - \boldsymbol{\mu}_X)^T\} = E\{(\mathbf{S}\mathbf{Y})(\mathbf{S}\mathbf{Y})^T\} = E\{\mathbf{S}\mathbf{Y}\mathbf{Y}^T\mathbf{S}^T\} = \mathbf{S}E\{\mathbf{Y}\mathbf{Y}^T\}\mathbf{S}^T = \mathbf{S}\mathbf{S}^T$$

$$\begin{array}{l} \mu_1 + \sqrt{\sigma_1} N(0,1) \\ \mu_2 + \sqrt{\sigma_2} N(0,1) \end{array}$$

Short Answers

2. [2pts] Suppose that in answering a question in a multiple choice test, an examinee either knows the answer, with probability p , or he guesses the answer with probability $1 - p$. Assume that the probability of answering a question correctly is $1 - \delta$ for an examinee who knows the answer and $1/m$ for an examinee who guesses, where m is the number of multiple choice alternatives. What is the probability that an examinee knew the answer to a question given that he correctly answered it.

$$\frac{(1 - \delta)p}{(1 - \delta)p + \frac{1}{m}(1 - p)}$$

Short Answers

3. [2pts] Why do we need smoothing in the naive Bayes classifier?

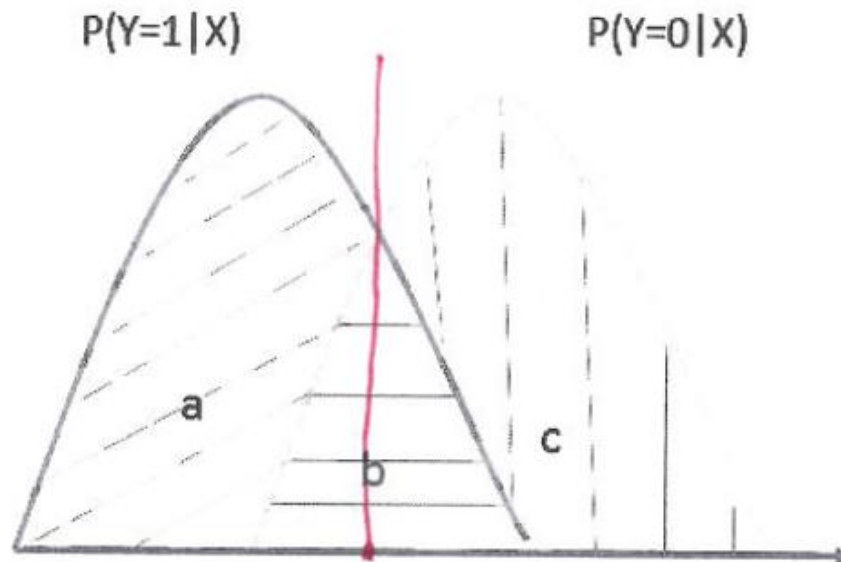
Avoid overfitting; avoid zero-count features causing problems in test.

4. [2pts] What is the main difference between semi-supervised learning and unsupervised learning?

Unsupervised learning only learns from unlabeled training data, while semi-supervised learning combines labeled and unlabeled data to train.

Short Answers

5. [2pts] Figure 2 shows the true distribution of the two classes. If your classifier discovers the true distribution then draw the decision boundary for your classifier. Also state whether this best classifier will have zero error or not.



Will not have
zero error.

Figure 2: Probability density of the two classes.

Short Answers

6. [2pts] For the data given in figure 3 draw the 1-NN decision boundary. Different shape represents Different classes.

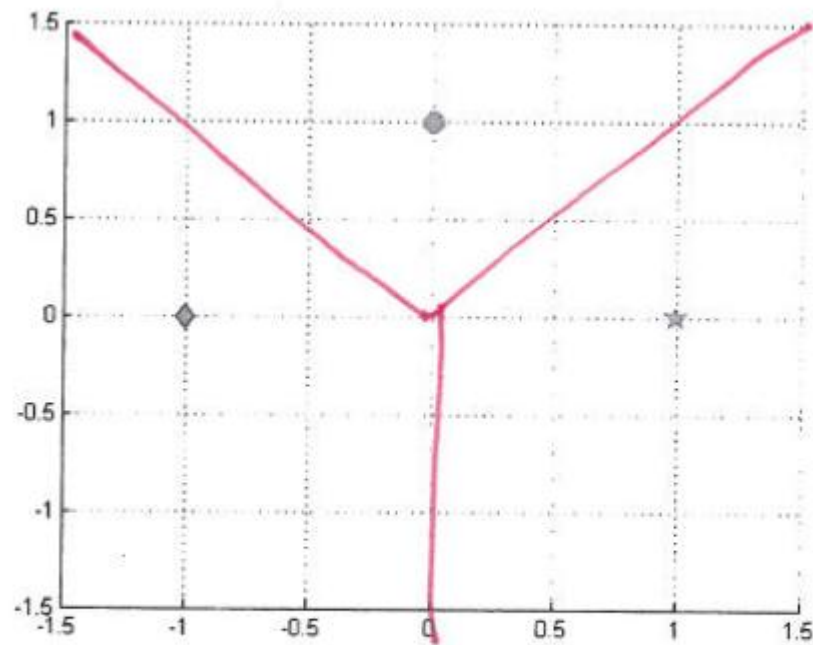


Figure 3: Points on two dimensional plane.

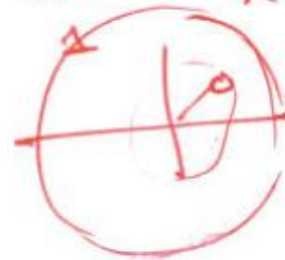
Short Answers

7. [2pts] What is the VC-dimension of the class of circles centered at the origin origin in \mathcal{R}^2 , or in other words the set of functions defined as

$$f(x_1, x_2) = \begin{cases} 1 & \text{if } (x_1)^2 + (x_2)^2 \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

1

For two points if you label the point ~~far~~ furthest from the origin as 1 and the one closest, you will not be able to shatter them.



Short Answers

8. [2pts] Suppose you have regression data generated from a polynomial of degree 5. Characterize the bias-variance of the estimates of the following models on the data with respect to the true model by circling appropriate entries

	Bias	Variance
Linear Regression	low/ high	low /high
Polynomial of degree 5	low /high	low /high
Polynomial of degree 20	low /high	low/ high

Short Answers

9. [2pts] When optimizing SVM with slack variables, we have the following objective function:

$$\arg \min_{w,z} \frac{1}{2} \|\vec{w}\|^2 + C \sum_i z_i \quad (1)$$

$$s.t. \forall i, y_i \vec{w} \cdot \vec{x}_i + z_i > 1, z_i \geq 0 \quad (2)$$

When we change C, how does the size of the margin change?

$C \uparrow$ margin \downarrow

Short Answers

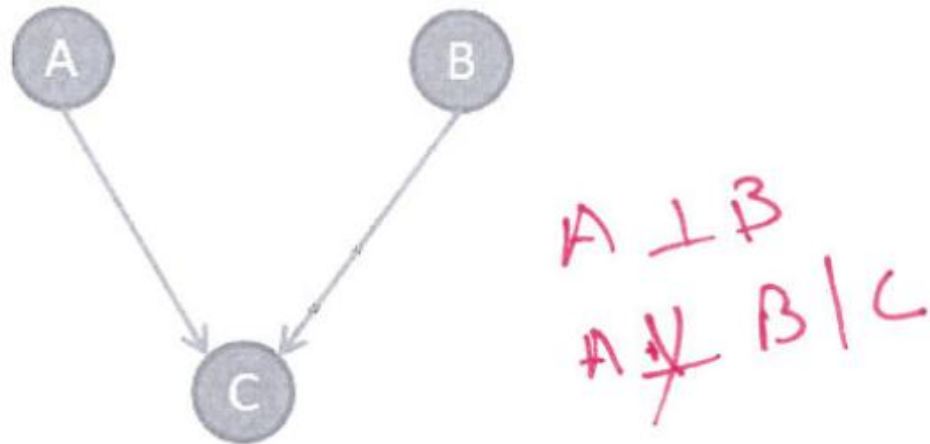


Figure 4: directed graphical model

10. [2pts] What are the independencies denoted by the graph in figure 4? Can the independencies represented by this directed graph be represented by an undirected graphical model?

No, not possible to represent using undirected graphical models.

Short Answers

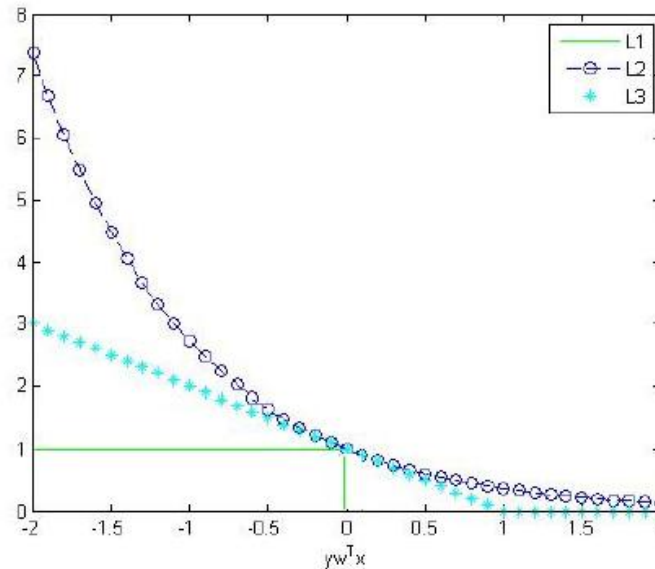


Figure 5: Loss function

11. [2pts] Label $L1$, $L2$ and $L3$ in figure 5 as one of “square loss” (used in regression), “exponential loss” (example used in adaboost), “zero one error” and “hinge loss” (example used in SVM).

L1: zero one loss.

L2: exponential loss.

L3: hinge loss.

Short Answers

12. [2pts] Name one method to enable linear regression to be used for modeling non-linear relationships.

Kernel trick.

13. [2pts] What is the reason for using regularization methods in machine learning problems?

Avoid overfitting.

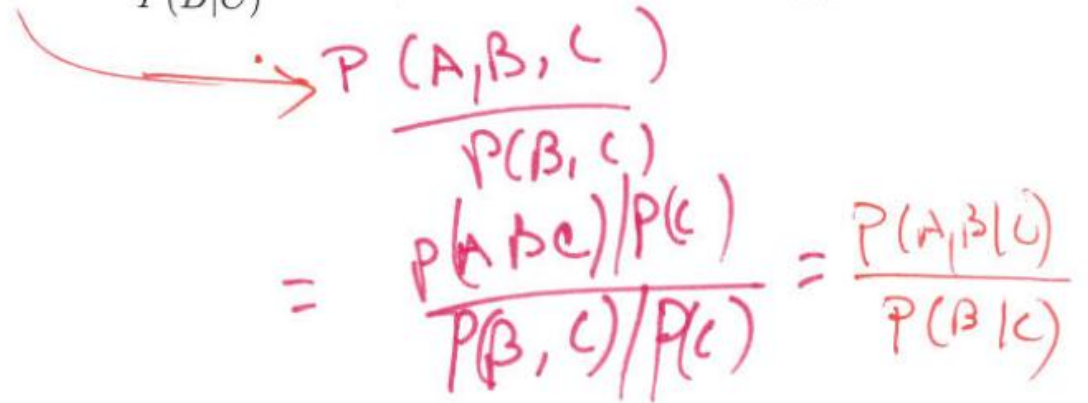
14. [2pts] What property does the l_1 regularization lead to?

Sparsity.

Short Answers

15. [2pts] Show that

$$p(A|B, C) = \frac{P(A, B|C)}{P(B|C)} \quad (3)$$



A handwritten derivation in red ink showing the steps to prove the formula. A red arrow points from the fraction $\frac{P(A, B|C)}{P(B|C)}$ in the printed equation to the first term of the handwritten expression. The handwritten expression is:

$$= \frac{\frac{P(A, B, C)}{P(C)}}{\frac{P(B, C)}{P(C)}} = \frac{P(A, B|C)}{P(B|C)}$$

Short Answers

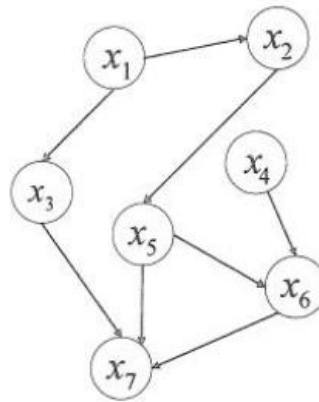


Figure 6: Graphical Model

16. [2pts] We have seven random variables x_1, x_2, \dots, x_7 and define a Bayesian network over them drawn in figure 6. From the Bayesian network, write down the joint distribution $p(x_1, x_2, x_3, x_4, x_5, x_6, x_7)$ of the seven variables.

$$p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_5|x_2) p(x_6|x_5, x_4) p(x_7|x_5, x_6, x_3)$$

Interpreting Data

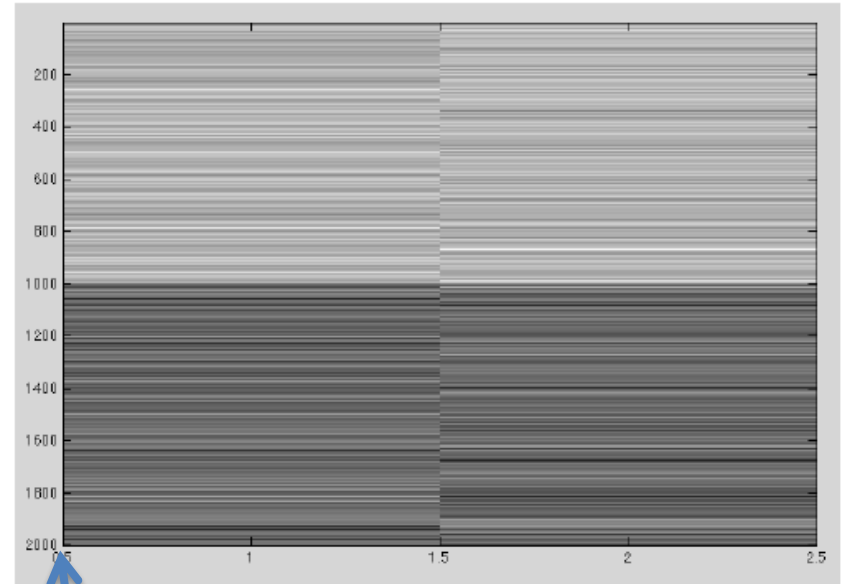
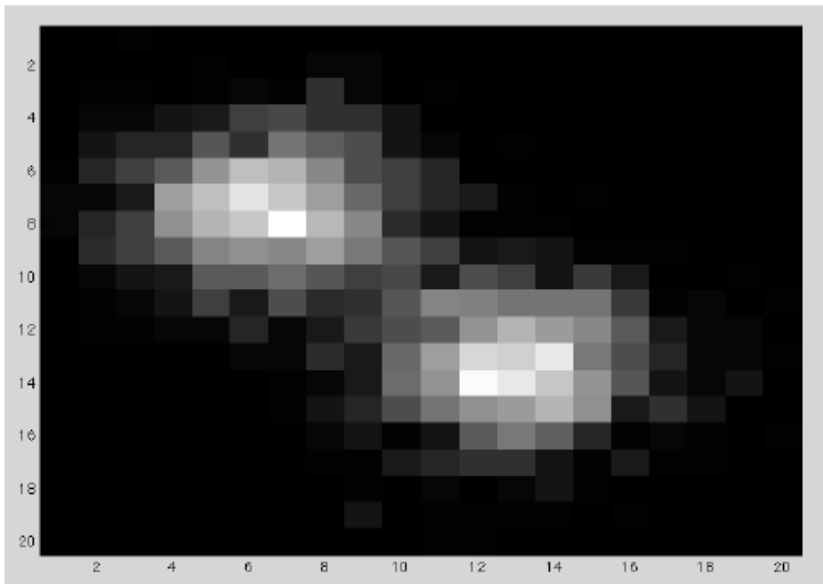


Figure 7 and figure 8 shows some Matlab images discussed in a lecture. These images summarize a data sample S where each instance is a point $\mathbf{x} = (x_1, x_2)$ in two-dimensional space. Figure 7 is a histogram where each dimension is placed into one of 20 equal-sized bins (white is a high count, black is a low count) and Figure 8 is a scaled grey image `imagesc` of the actual values of the data.

1. [2pts] About how many instances are in S ?

2000.

Interpreting Data

2. [2pts] Which of the following methods would be *reasonable* to use with this data?

(a) k -means ☺

(b) Gaussian naive Bayes

(c) k -nearest-neighbor

(d) mixtures of Gaussians ☺

(e) logistic regression

(f) principal component analysis ☺

No labels

Interpreting Data

3. [2pts] Which of the methods above would be *most appropriate* to use with this data, and why?

Any reasonable explanation of (a), (e), (f) receives full grade.
e.g. We can use PCA because we can analyze the dimension
With largest variance. Or, we can use GMM because data looks
like a bivariate Gaussian.

Interpreting Data

4. [2pts] Let X_1 and X_2 be random variables corresponding to the two dimensions of the data. From this sample, do you think that X_1 and X_2 are independent? Why or why not?

No, because from Figure 8, the values of X_1 and X_2 are highly correlated for all examples.

5. [2pts] Suppose S is a sample of the *positive* examples of a labeled data: i.e., all the instances in S have label $+1$. Kevin believes that the best classifier learner to use for this task is Gaussian naive Bayes. Do you agree? why or why not?

We accept two answers (any of them receives full grade):

1. GNB has strong independence assumption of its features, so we cannot use it here.
2. We only have positive examples, but not negative examples, so we cannot use it.

Linear Classifiers

1. [2pts] Kevin wants to improve the accuracy of his implementation of a multinomial naive Bayes text classifier by using a new smoothing method. How would he be most likely to improve performance, and why?
 - (a) By smoothing the estimates for the class priors.
 - ✓ (b) By smoothing the estimates for the conditional probabilities of features given a class.

Conditional probabilities of features are more likely to be zero than class priors.

Linear Classifiers

2. [2pts] Consider the perceptron applied to a series of examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, where k is the number of mistakes, \mathbf{v}_j is the j -th classifier produced by the perceptron, and R^2 is an upper bound on $\mathbf{x}_i \cdot \mathbf{x}_i$ for all instances \mathbf{x}_i . Suppose also that the data is separable with margin γ , as defined in the lecture.

Which of the following must be true?

- (a) There is some vector \mathbf{u} such that for all i , $(\mathbf{u} \cdot \mathbf{x}_i)y_i > \gamma$ and \mathbf{u} is orthogonal to \mathbf{x}_i .
- (b) The final classifier \mathbf{v}_k produced by the perceptron will have an error rate of no more than k/n on the training examples.
- (c) $k < \frac{R^2}{\gamma^2}$
- (d) $k < \frac{\gamma^2}{R^2}$
- (e) $\gamma \geq 1$
- (f) $\gamma \geq 0$

Linear Classifiers

3. [2pts] Suppose multinomial naive Bayes gets zero *training* error on a binary dataset D . Does that imply that D is linearly separable? Justify your answer.

true

multinomial naive Bayes is a linear classifier.

Linear Classifiers

4. A common regularization term for logistic regression is

$$\mu \sum_{j=1}^d (w_j)^2$$

where d is the number of dimensions of the data, and w_j is the weight for the j -th parameter.

- (a) [2pts] In Kevin's experiments with logistic regression, he used cross-validation to pick the best value of μ among these values: -1, -0.5, -0.1, 0, 0.1, 0.5, 1. He asks your advice about this selection—which set of values will you ask him to use among his suggested values and why?

use positive values.

Tradeoff parameters are required to be positive.

- (b) [2pts] Kevin wants to explore using a regularization term of the form $\mu \sum_{j=1}^d (w_j)^k$ for $k = 1, 2, 3, 4$. He asks your advice on this proposal—what do you say?

use even ~~numbers~~ k

odd k will make the regularization term to be negative infinity.

Decision Trees

Assume there are 4 binary attributes (value can only be T or F): Attribute#1, Attribute#2, Attribute#3, Attribute#4, two kinds of labels: 0 and 1. Now you are given the following 8 instances

Table 1: Instances

Instance	Label	Attribute#1	Attribute#2	Attribute#3	Attribute#4
1	1	T	T	T	F
1	1	T	T	T	F
1	1	F	T	T	F
1	1	F	T	F	F
0	0	T	T	F	F
0	0	T	T	F	F
0	0	F	T	F	F
0	0	F	T	T	F

- 1 [2pts] Compute the information gain for selecting each attribute as the separating attribute in the first (root) node. Which attribute should be selected?

$$\begin{aligned}
 IG(\text{Attribute \#1}) &= IG(\text{Attribute \#2}) = IG(\text{Attribute \#4}) = 0 \\
 IG(\text{Attribute \#3}) &= \left(-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}\right) - \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}\right) \times \frac{1}{2} \times 2 = 1 + \left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}\right) \\
 &\approx 0.189
 \end{aligned}$$

should choose Attribute #3

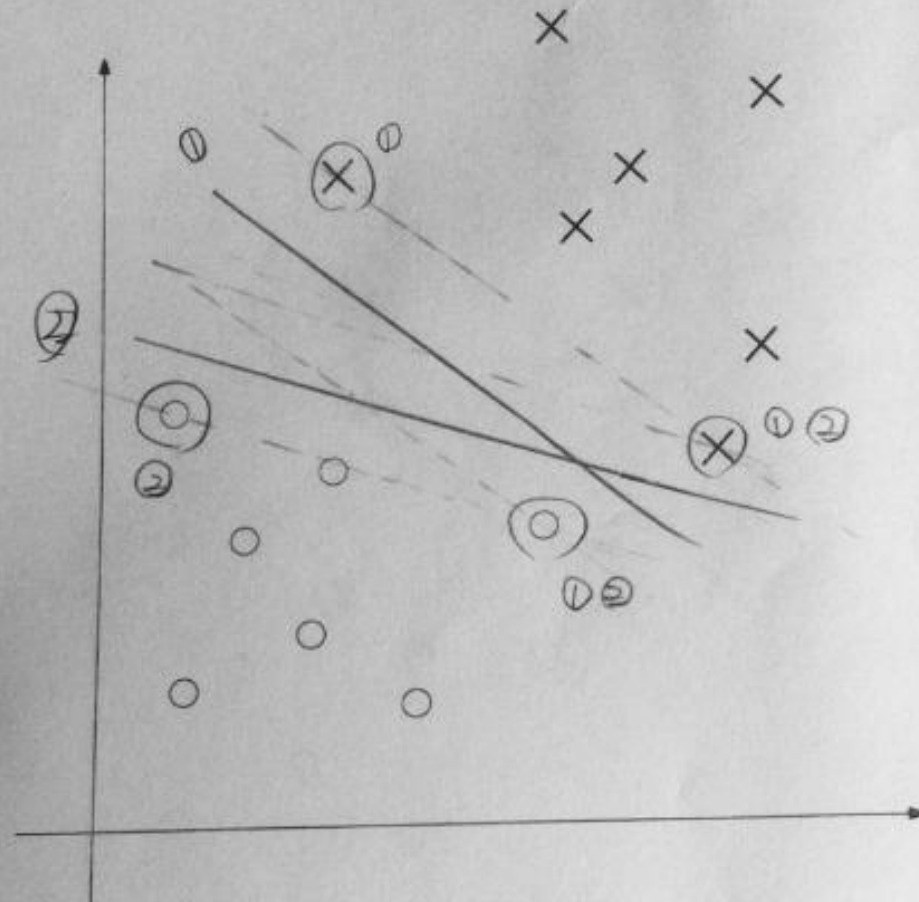
Decision Trees

2 [2pts] Does there exist a decision tree which can classify the given instances perfectly? If yes, draw that decision tree, otherwise give a simple explanation.

No, because two instances have the same features but different labels.

SVM

1. [2pts] Given the data points shown in the following figure with two different labels. Draw the linear SVM decision boundary for this binary classification problem. Also circle the support vectors (data points) for the boundary you find out



① ②

any one
would be
correct

SVM

2. [2pts] If you are asked to use the Perceptron algorithm to classify these data points, what will the decision boundary look like? Explain the similarities and differences between Perceptron and SVM.

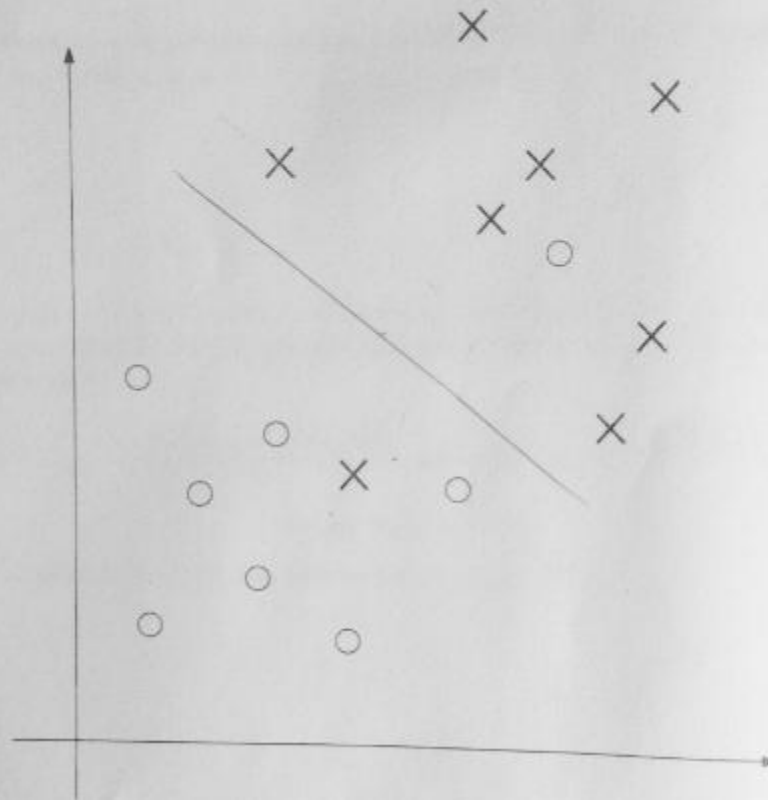
They both have linear decision boundary which can separate the dataset.

SVM provides the decision boundary which maximizes the margin. ¹⁵Perceptron does not

SVM

3. [2pts] Draw the SVM decision boundary for another set of data points on the following figure, and also explain the reason for achieving such a boundary. Please simply describe any other additional conditions or constraints you want to use, for example slack variable.

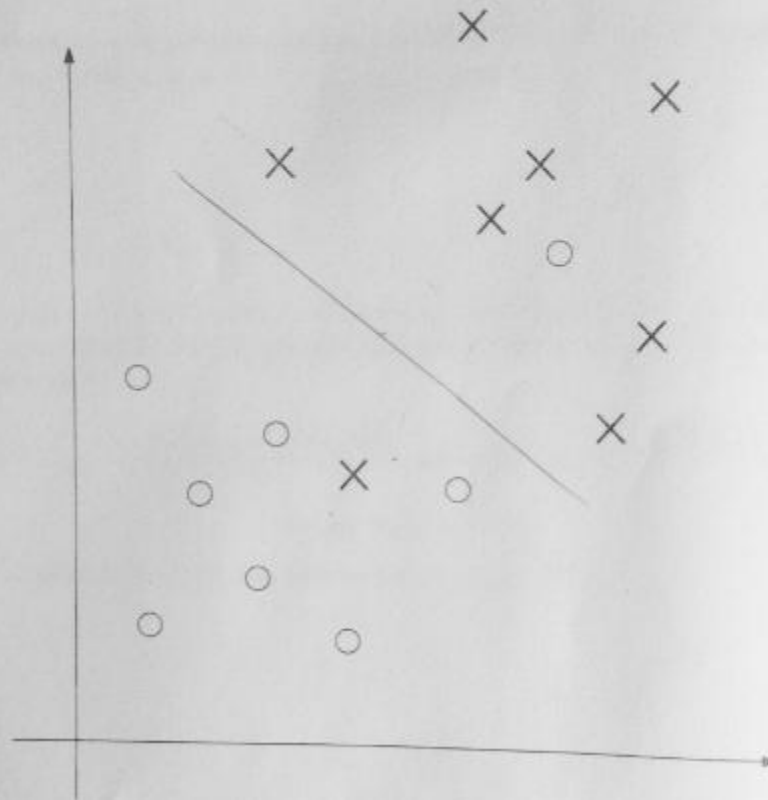
Need to use slack variables (soft margin)



SVM

3. [2pts] Draw the SVM decision boundary for another set of data points on the following figure, and also explain the reason for achieving such a boundary. Please simply describe any other additional conditions or constraints you want to use, for example slack variable.

Need to use slack variables (soft margin)



Comparing Classifiers

1. [2pts] When calculating confidence interval for a classifier, when we would like to use 10-fold cross-validation instead of holding out half of data?

When the data size is small. We want to fully utilize data for training.

2. [2pts] Given the sample mean and variance, calculating the p-value. Is the one-tailed p-value larger or the two-tailed p-value?

The two-tailed p-value is twice as large as the one-tailed p-value.

Comparing Classifiers

1. [2pts] When calculating confidence interval for a classifier, when we would like to use 10-fold cross-validation instead of holding out half of data?

When the data size is small. We want to fully utilize data for training.

2. [2pts] Given the sample mean and variance, calculating the p-value. Is the one-tailed p-value larger or the two-tailed p-value?

The two-tailed p-value is twice as large as the one-tailed p-value.

Comparing Classifiers

3. [2pts] Kevin proposed a hypothesis with accuracy 80% on a binary classification problem. While when you tested on 100 data points, you get accuracy for 70%. You can compute 95% confidence interval by

$$[Accu_S(h) - Z_N \sqrt{\frac{Accu_S(h)(1 - Accu_S(h))}{n}}, Accu_S(h) + Z_N \sqrt{\frac{Accu_S(h)(1 - Accu_S(h))}{n}}]$$

where $Z_{0.95} = 1.96$

Under 95% confidence level, do you believe Taiti's claim? Why?

Plug into formula.

$$80\% \notin \left[70\% - 1.96 \sqrt{\frac{0.7(1-0.7)}{100}}, 70\% + 1.96 \sqrt{\frac{0.7(1-0.7)}{100}} \right]$$

Thus, under 95% confidence level,
we don't believe Kevin's claim.

Ensemble Classifiers and SSL

1. [2pts] The instances in Figure 7 and Figure 8 are actually user profiles on an on-line gaming site, where for a user, x_1 is the number of hours spent per week logged on, and x_2 is the number of unique games played. You would like to find out if you can predict, from this information, some demographical information, namely whether the user is male or female. About 5% of the users have created profiles including this information. What methods would be *plausible* to use in this case?
- (a) PCA
 - (b) k -means
 - ☒ (c) seeded k -means
 - ☒ (d) seeded mixtures of Gaussians
 - (e) seeded k -nearest-neighbor

Ensemble Classifiers and SSL

2. [2pts] You suggest adding using additional features from a user's profile to make this task easier: specifically his/her age, how long he/she has had an account, and the specific games he/she have played. Kevin says that age is unnecessary to use as a feature, because you are trying to predict sex, and age and sex are independent. Is he right? Why or why not?

No.

They can be conditionally dependent.

3. [2pts] Would it be more appropriate to use a transductive semi-supervised learner, or an inductive semi-supervised learner? Why?

inductive

to apply to unseen data

Ensemble Classifiers and SSL

4. [2pts] Which method is easier to parallelize: boosting, or bagging? Why?

bagging

boosting is sequential

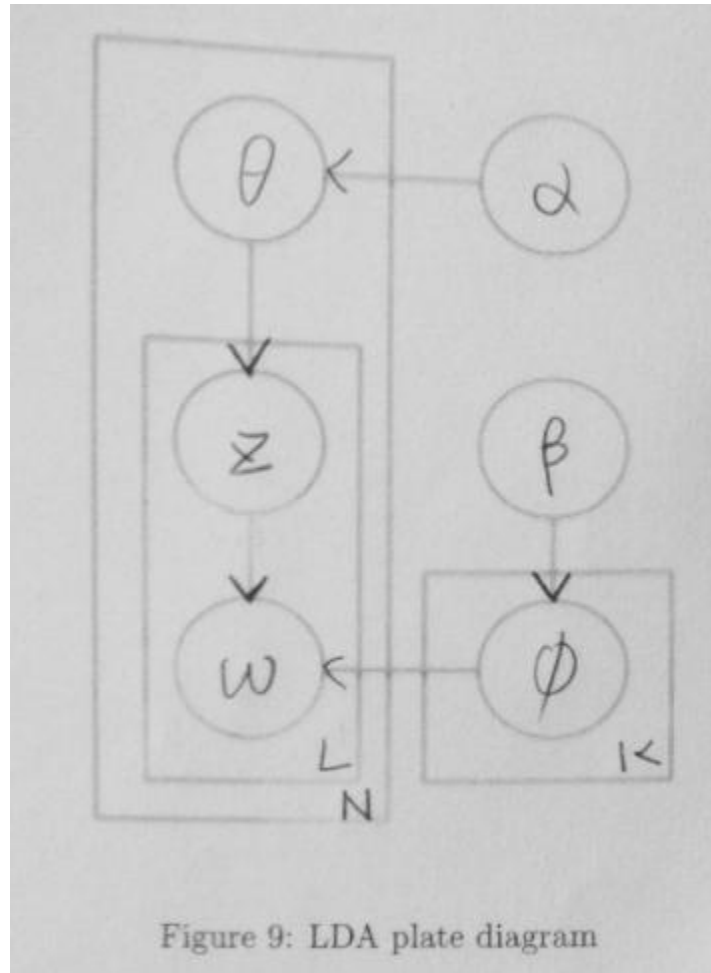
5. [2pts] You have implemented stacking using 5 base classifiers and 10-fold cross-validation. One of the five base classifiers is used as the final meta-classifier (the “top” of the stack). If each of your classifiers takes about the same time T to train on your dataset, and time t to evaluate a single example, how long does it take to train the stacked learner? How long does it take to evaluate an example with the learned classifier? Assume nothing has been parallelized.

5T

5t

Latent Dirichlet Allocation

1. [6pts] Draw the counts for each plate, the direction for each arrow and the variable for each circle.



Latent Dirichlet Allocation

2. [2pts] List the variables that are known in a typical use of lda

$W \quad (K, L, N)$

3. [2pts] List the variables that have a discrete domain.

$Z, w \quad (K, L, N)$

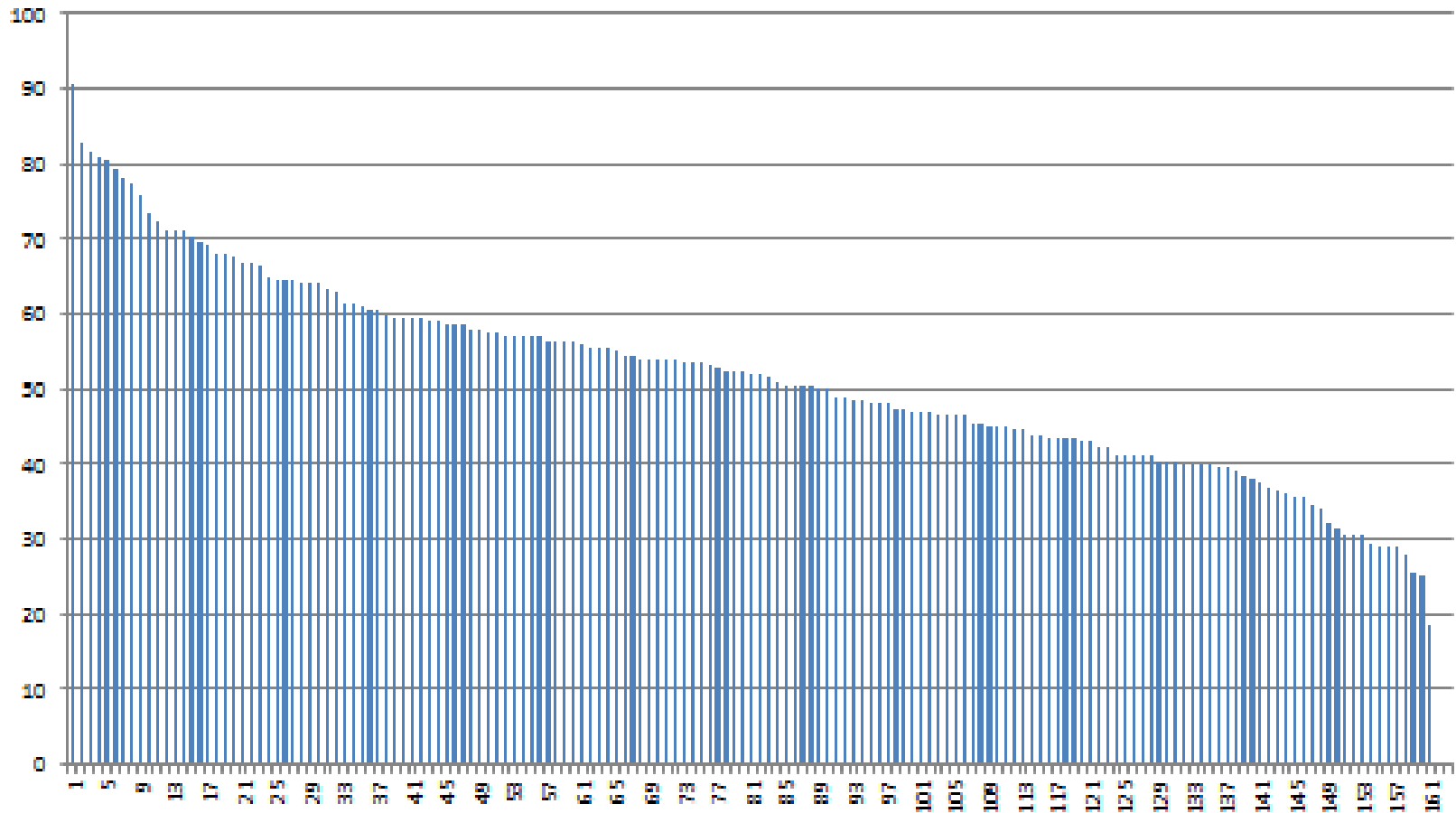
Latent Dirichlet Allocation

4. [2pts] Name one inference method that can be used for lda,

(variational) EM. Gibbs Sampling

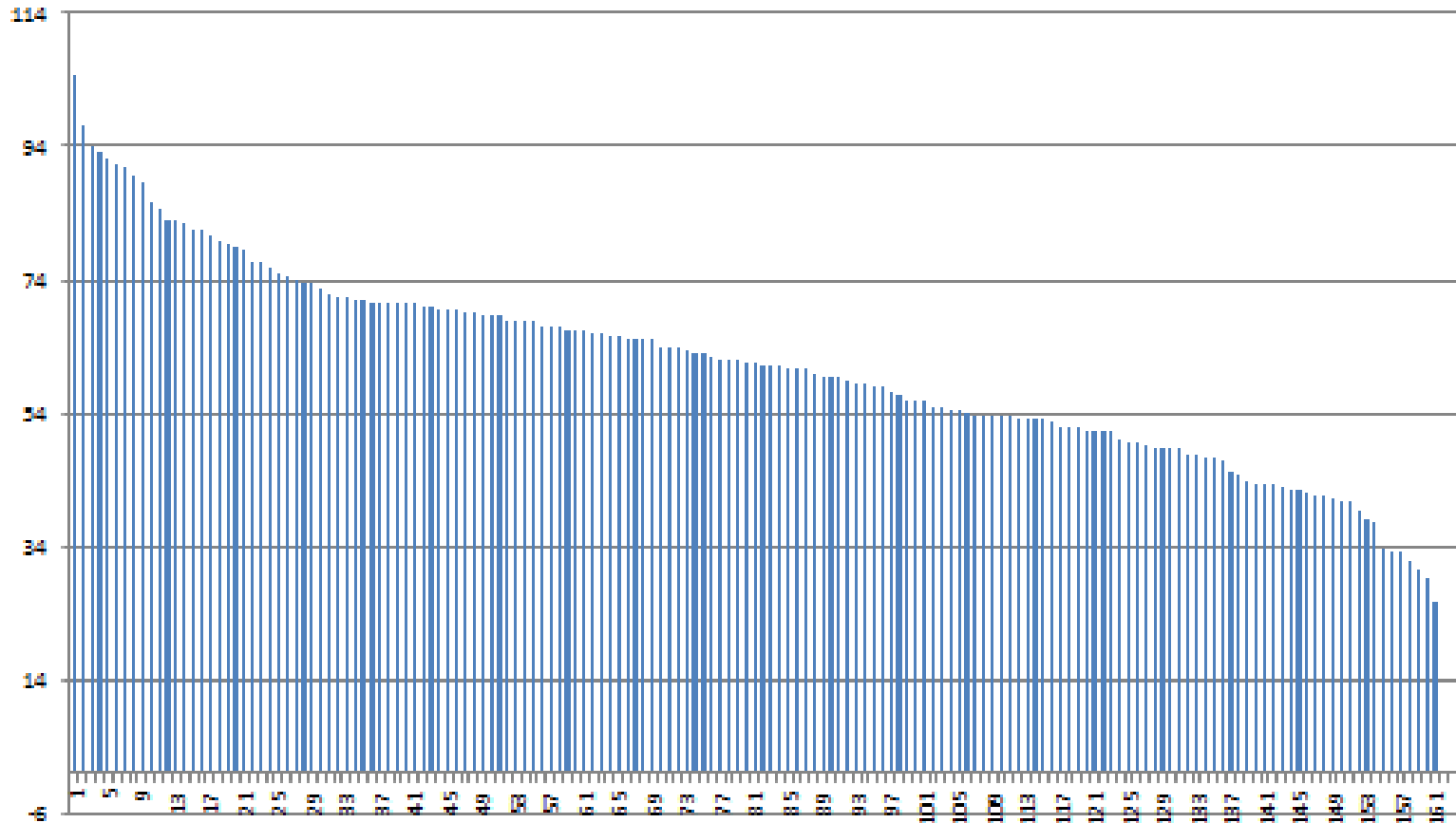
Final Exam Statistics

Final Exam w/o Extra Credits



max = 90.5, mean = 52, stdv = 13.

Final Exam w Extra Credits



max = 104.5, mean = 61, stdv = 15.