

# Final Review, Day 1 (Wednesday)

EECS 398: Practical Data Science, Winter 2025

practicaldsc.org • github.com/practicaldsc/wn25 •  See latest announcements [here on Ed](#)

# Agenda

- We'll work through (some of) the problems in this worksheet:  
<https://study.practicaldsc.org/fi-review-wednesday/index.html>
- I'll post these annotated slides after lecture.
- So that we focus our time on more challenging concepts, tell me which questions you want me to cover by voting here:

[practicaldsc.org/q](https://practicaldsc.org/q)



# Problem 1

Consider the vectors  $\vec{u}$  and  $\vec{v}$ , defined below.

$$\vec{u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \vec{v} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

We define  $X \in \mathbb{R}^{3 \times 2}$  to be the matrix whose first column is  $\vec{u}$  and whose second column is  $\vec{v}$ .

## Problem 1.1

In this part only, let  $\vec{y} = \begin{bmatrix} -1 \\ k \\ 252 \end{bmatrix}$ .

Find a scalar  $k$  such that  $\vec{y}$  is in  $\text{span}(\vec{u}, \vec{v})$ . Give your answer as a constant with no variables.

# Problem 1

Consider the vectors  $\vec{u}$  and  $\vec{v}$ , defined below.

$$\vec{u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \vec{v} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

We define  $X \in \mathbb{R}^{3 \times 2}$  to be the matrix whose first column is  $\vec{u}$  and whose second column is  $\vec{v}$ .

## Problem 1.2

Show that:

$$(X^T X)^{-1} X^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

*Hint: If  $A = \begin{bmatrix} a_1 & 0 \\ 0 & a_2 \end{bmatrix}$ , then  $A^{-1} = \begin{bmatrix} \frac{1}{a_1} & 0 \\ 0 & \frac{1}{a_2} \end{bmatrix}$ .*

# Problem 1

Consider the vectors  $\vec{u}$  and  $\vec{v}$ , defined below.

$$\vec{u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \vec{v} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

We define  $X \in \mathbb{R}^{3 \times 2}$  to be the matrix whose first column is  $\vec{u}$  and whose second column is  $\vec{v}$ .

## Problem 1.3

In parts 3 and 4 only, let  $\vec{y} = \begin{bmatrix} 4 \\ 2 \\ 8 \end{bmatrix}$ .

Find scalars  $a$  and  $b$  such that  $a\vec{u} + b\vec{v}$  is the vector in  $\text{span}(\vec{u}, \vec{v})$  that is as close to  $\vec{y}$  as possible. Give your answers as constants with no variables.

# Problem 1

Consider the vectors  $\vec{u}$  and  $\vec{v}$ , defined below.

$$\vec{u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \vec{v} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

We define  $X \in \mathbb{R}^{3 \times 2}$  to be the matrix whose first column is  $\vec{u}$  and whose second column is  $\vec{v}$ .

## Problem 1.4

Let  $\vec{e} = \vec{y} - (a\vec{u} + b\vec{v})$ , where  $a$  and  $b$  are the values you found in part (c).

What is  $\|\vec{e}\|$ ?

- ☐ 0
- ☐  $3\sqrt{2}$
- ☐  $4\sqrt{2}$
- ☐ 6
- ☐  $6\sqrt{2}$
- ☐  $2\sqrt{21}$

# Problem 1

Consider the vectors  $\vec{u}$  and  $\vec{v}$ , defined below.

$$\vec{u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \vec{v} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

We define  $X \in \mathbb{R}^{3 \times 2}$  to be the matrix whose first column is  $\vec{u}$  and whose second column is  $\vec{v}$ .

## Problem 1.5

Is it true that, for any vector  $\vec{y} \in \mathbb{R}^3$ , we can find scalars  $c$  and  $d$  such that the sum of the entries in the vector  $\vec{y} - (c\vec{u} + d\vec{v})$  is 0?

- ☐ Yes, because  $\vec{u}$  and  $\vec{v}$  are linearly independent.
- ☐ Yes, because  $\vec{u}$  and  $\vec{v}$  are orthogonal.
- ☐ Yes, but for a reason that isn't listed here.
- ☐ No, because  $\vec{y}$  is not necessarily in  $\text{span}(\vec{u}, \vec{v})$ .
- ☐ No, because neither  $\vec{u}$  nor  $\vec{v}$  is equal to the vector  $\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}^T$ .
- ☐ No, but for a reason that isn't listed here.





## Problem 1.6

Suppose that  $Q \in \mathbb{R}^{100 \times 12}$ ,  $\vec{s} \in \mathbb{R}^{100}$ , and  $\vec{f} \in \mathbb{R}^{12}$ . What are the dimensions of the following product?

$$\vec{s}^T Q \vec{f}$$

- ☐ scalar
- ☐  $12 \times 1$  vector
- ☐  $100 \times 1$  vector
- ☐  $100 \times 12$  matrix
- ☐  $12 \times 12$  matrix
- ☐  $12 \times 100$  matrix
- ☐ undefined

## Problem 2

One piece of information that may be useful as a feature is the proportion of SAT test takers who took the test in a particular month. The variable `lunch_props` contains 8 values, each of which are either "low", "medium", or "high". We will use the following Pipeline:

```
# Note: The FunctionTransformer is only needed to change the result  
# of the OneHotEncoder from a "sparse" matrix to a regular matrix  
# so that it can be used with StandardScaler;  
# it doesn't change anything mathematically.
```

```
pl = Pipeline([  
    ("ohe", OneHotEncoder(drop="first")),  
    ("ft", FunctionTransformer(lambda X: X.toarray())),  
    ("ss", StandardScaler())  
])
```

Given the above information, we can conclude that `lunch_props` has **(a)** value(s) equal to "low", **(b)** value(s) equal to "medium", and **(c)** value(s) equal to "high". (Note: Each of (a), (b), and (c) should be positive numbers, such that together, they add to 8.)

After calling `pl.fit(lunch_props)`, `pl.transform(lunch_props)` evaluates to the following

```
array([[ 1.29099445, -0.37796447],  
       [-0.77459667, -0.37796447],  
       [-0.77459667, -0.37796447],  
       [-0.77459667,  2.64575131],  
       [ 1.29099445, -0.37796447],  
       [ 1.29099445, -0.37796447],  
       [-0.77459667, -0.37796447],  
       [-0.77459667, -0.37796447]])
```

and `pl.named_steps["ohe"].get_feature_names()` evaluates to the following array:

```
array(["x0_low", "x0_med"], dtype=object)
```

## Problem 3

Suppose we have one qualitative variable that we convert to numerical values using one-hot encoding. We've shown the first four rows of the resulting design matrix below:

<b>a</b>	<b>b</b>	<b>c</b>
1	0	0
1	0	0
0	0	1
0	1	0

### Problem 3.1

Say we train a linear model  $m_1$  on these data. Then, we replace all of the 1 values in column **a** with 3's and all of the 1 values in column **b** with 2's and train a new linear model  $m_2$ . Neither  $m_1$  nor  $m_2$  have an intercept term. On the training data, the average squared loss for  $m_1$  will be \_\_\_\_\_ that of  $m_2$ .

- ☐ greater than
- ☐ less than
- ☐ equal to
- ☐ impossible to tell

## Problem 3

Suppose we have one qualitative variable that we convert to numerical values using one-hot encoding. We've shown the first four rows of the resulting design matrix below:

<b>a</b>	<b>b</b>	<b>c</b>
1	0	0
1	0	0
0	0	1
0	1	0

### Problem 3.2

To account for the intercept term, we add a column of all ones to our design matrix from part a. That is, the resulting design matrix has four columns: **a** with 3's instead of 1's, **b** with 2's instead of 1's, **c**, and a column of all ones. What is the rank of the new design matrix with these four columns?

- ☐ 1
- ☐ 2
- ☐ 3
- ☐ 4

## Problem 4

Suppose we build a binary classifier that uses a song's `"track_name"` and `"artist_names"` to predict whether its genre is `"Hip-Hop/Rap"` (1) or not (0).

For our classifier, we decide to use a brand-new model built into `sklearn` called the

`BillyClassifier`. A `BillyClassifier` instance has three hyperparameters that we'd like to tune. Below, we show a dictionary containing the values for each hyperparameter that we'd like to try:

```
hyp_grid = {  
    "radius": [0.1, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 100], # 12 total  
    "inflection": [-5, -4, -3, -2, -1, 0, 1, 2, 3, 4], # 10 total  
    "color": ["red", "yellow", "green", "blue", "purple"] # 5 total  
}
```

To find the best combination of hyperparameters for our `BillyClassifier`, we first conduct a train-test split, which we use to create a training set with 800 rows. We then use `GridSearchCV` to conduct  $k$ -fold cross-validation for each combination of hyperparameters in `hyp_grid`, with  $k = 4$ .

## Problem 4.1

When we call `GridSearchCV`, how many times is a `BillyClassifier` instance trained in total? Give your answer as an integer.

## Problem 4

Suppose we build a binary classifier that uses a song's "track\_name" and "artist\_names" to predict whether its genre is "Hip-Hop/Rap" (1) or not (0).

For our classifier, we decide to use a brand-new model built into `sklearn` called the

`BillyClassifier`. A `BillyClassifier` instance has three hyperparameters that we'd like to tune. Below, we show a dictionary containing the values for each hyperparameter that we'd like to try:

```
hyp_grid = {  
    "radius": [0.1, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 100], # 12 total  
    "inflection": [-5, -4, -3, -2, -1, 0, 1, 2, 3, 4], # 10 total  
    "color": ["red", "yellow", "green", "blue", "purple"] # 5 total  
}
```

To find the best combination of hyperparameters for our `BillyClassifier`, we first conduct a train-test split, which we use to create a training set with 800 rows. We then use `GridSearchCV` to conduct  $k$ -fold cross-validation for each combination of hyperparameters in `hyp_grid`, with  $k = 4$ .

## Problem 4.2

In each of the 4 folds of the data, how large is the training set, and how large is the validation set? Give your answers as integers.

size of training set =

size of validation set =

## Problem 4.3

Suppose that after fitting a `GridSearchCV` instance, its `best_params_` attribute is

```
{"radius": 8, "inflection": 4, "color": "blue"}
```

Select all true statements below.

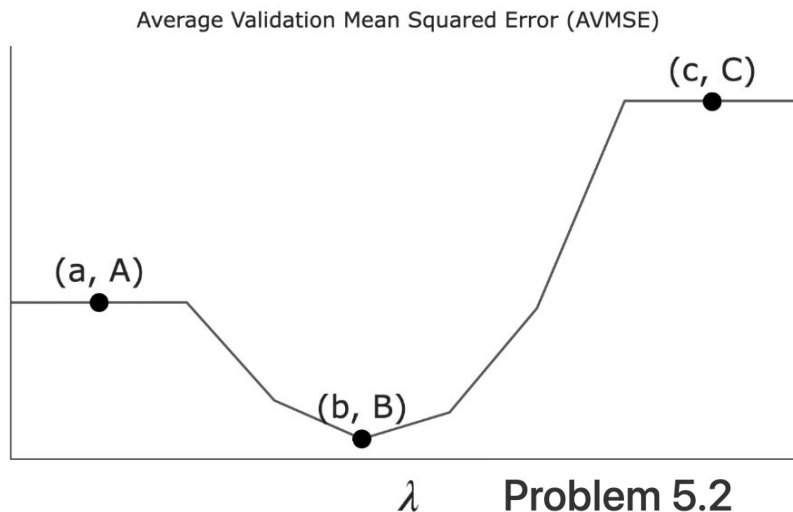
- ☐ The specific combination of hyperparameters in `best_params_` had the highest average training accuracy among all combinations of hyperparameters in `hyp_grid`.
- ☐ The specific combination of hyperparameters in `best_params_` had the highest average validation accuracy among all combinations of hyperparameters in `hyp_grid`.
- ☐ The specific combination of hyperparameters in `best_params_` had the highest training accuracy among all combinations of hyperparameters in `hyp_grid`, in each of the 4 folds of the training data.
- ☐ The specific combination of hyperparameters in `best_params_` had the highest validation accuracy among all combinations of hyperparameters in `hyp_grid`, in each of the 4 folds of the training data.
- ☐ A `BillyClassifier` that is fit using the specific combination of hyperparameters in `best_params_` is guaranteed to have the best accuracy on unseen testing data among all combinations of hyperparameters in `hyp_grid`.



## Problem 5

Suppose we want to use LASSO (i.e. minimize mean squared error with  $L_1$  regularization) to fit a linear model that predicts the number of ingredients in a product given its price and rating.

Let  $\lambda$  be a non-negative regularization hyperparameter. Using cross-validation, we determine the average validation mean squared error — which we'll refer to as AVMSE in this question — for several different choices of  $\lambda$ . The results are given below.



### Problem 5.1

As  $\lambda$  increases, what happens to model complexity and model variance?

### Problem 5.2

What does the value  $A$  on the graph above correspond to?

- ☐ The AVMSE of the  $\lambda$  we'd choose to use to train a model.
- ☐ The AVMSE of an unregularized multiple linear regression model.
- ☐ The AVMSE of the constant model.

## Problem 6

Suppose we want to minimize the function  $R(h) = e^{(h+1)^2}$

### Problem 6.1

Without using gradient descent or calculus, what is the value  $h^*$  that minimizes  $R(h)$ ?

# Problem 6

Suppose we want to minimize the function  $R(h) = e^{(h+1)^2}$

## Problem 6.2

Now, suppose we want to use gradient descent to minimize  $R(h)$ . Assume we use an initial guess of  $h^{(1)} = 0$ . What is  $h^{(1)}$ ? Give your answer in terms of a generic step size,  $\alpha$ , and other constants. ( $e$  is a constant.)

# Problem 6

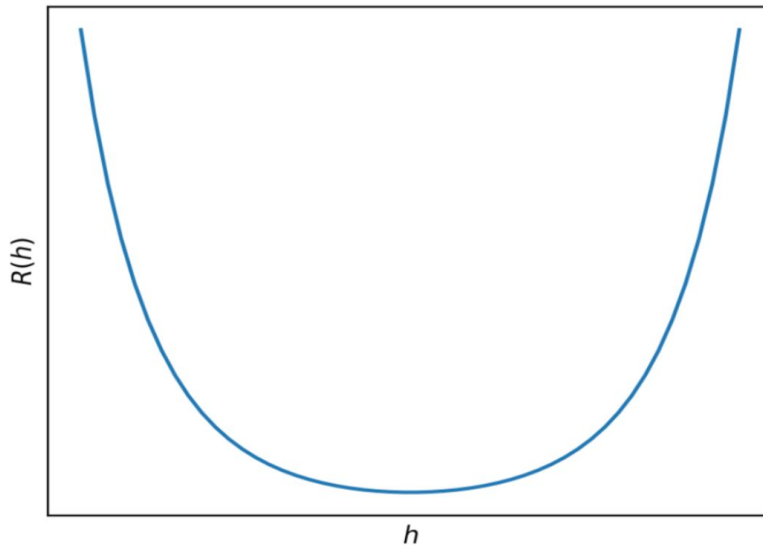
Suppose we want to minimize the function  $R(h) = e^{(h+1)^2}$

## Problem 6.3

Using your answers from the previous two parts, what should we set the value of  $\alpha$  to be if we want to ensure that gradient descent finds  $h^*$  after just one iteration?

## Problem 6.4

Below is a graph of  $R(h)$  with no axis labels.



True or False: Given an appropriate choice of step size,  $\alpha$ , gradient descent is guaranteed to find the minimizer of  $R(h)$ .

## Problem 7

We decide to build a classifier that takes in a state's demographic information and predicts whether, in a given year:

- The state's mean math score was greater than its mean verbal score (1), or
- the state's mean math score was less than or equal to its mean verbal score (0).

## Problem 7.2

Suppose we train a classifier, named Classifier 1, and it achieves an accuracy of  $\frac{5}{9}$  on our training set.

Typically, root mean squared error (RMSE) is used as a performance metric for regression models, but mathematically, nothing is stopping us from using it as a performance metric for classification models as well.

What is the RMSE of Classifier 1 on our training set? Give your answer as a **simplified fraction**.

## Problem 7.2

Suppose we train a classifier, named Classifier 1, and it achieves an accuracy of  $\frac{5}{9}$  on our training set.

Typically, root mean squared error (RMSE) is used as a performance metric for regression models, but mathematically, nothing is stopping us from using it as a performance metric for classification models as well.

What is the RMSE of Classifier 1 on our training set? Give your answer as a **simplified fraction**.

## Problem 7.3

While Classifier 1's accuracy on our training set is  $\frac{5}{9}$ , its accuracy on our test set is  $\frac{1}{4}$ . Which of the following scenarios is most likely?

- ☐ Classifier 1 overfit to our training set; we need to increase its complexity.
- ☐ Classifier 1 overfit to our training set; we need to decrease its complexity.
- ☐ Classifier 1 underfit to our training set; we need to increase its complexity.
- ☐ Classifier 1 underfit to our training set; we need to decrease its complexity.



For the remainder of this question, suppose we train another classifier, named Classifier 2, again on our training set. Its performance on the training set is described in the confusion matrix below. Note that the columns of the confusion matrix have been separately normalized so that each has a sum of 1.

	Actually 0	Actually 1
Predicted 0	0.9	0.4
Predicted 1	0.1	0.6

## Problem 7.4

Suppose `conf` is the DataFrame above. Which of the following evaluates to a Series of length 2 whose only unique value is the number `1`?

- ☐ `conf.sum(axis=0)`
- ☐ `conf.sum(axis=1)`

## Problem 7.5

Fill in the blank: the \_\_\_\_ of Classifier 2 is guaranteed to be 0.6.

- ☐ precision
- ☐ recall

## Problem 7.6

Suppose a fraction  $\alpha$  of the labels in the training set are actually 1 and the remaining  $1 - \alpha$  are actually 0. The accuracy of Classifier 2 is 0.65. What is the value of  $\alpha$ ?

Hint: If you're unsure on how to proceed, here are some guiding questions:

- Suppose the number of  $y$ -values that are actually 1 is  $A$  and that the number of  $y$ -values that are actually 0 is  $B$ . In terms of  $A$  and  $B$ , what is the accuracy of Classifier 2? Remember, you'll need to refer to the numbers in the confusion matrix above.
- What is the relationship between  $A$ ,  $B$ , and  $\alpha$ ? How does it simplify your calculation for the accuracy in the previous step?

	Actually 0	Actually 1
Predicted 0	0.9	0.4
Predicted 1	0.1	0.6

## Problem 8

Suppose we fit a logistic regression model that predicts whether a product is designed for sensitive skin, given its price,  $x^{(1)}$ , number of ingredients,  $x^{(2)}$ , and rating,  $x^{(3)}$ . After minimizing average cross-entropy loss, the optimal parameter vector is as follows:

$$\vec{w}^* = \begin{bmatrix} -1 \\ 1/5 \\ -3/5 \\ 0 \end{bmatrix}$$

In other words, the intercept term is  $-1$ , the coefficient on price is  $\frac{1}{5}$ , the coefficient on the number of ingredients is  $-\frac{3}{5}$ , and the coefficient on rating is 0.

Consider the following four products:

- **Wolfcare**: Costs \$15, made of 20 ingredients, 4.5 rating
- **Go Blue Glow**: Costs \$25, made of 5 ingredients, 4.9 rating
- **DataSPF**: Costs \$50, made of 15 ingredients, 3.6 rating
- **Maize Mist**: Free, made of 1 ingredient, 5.0 rating

Which of the following products have a predicted probability of being designed for sensitive skin of **at least 0.5 (50%)**? For each product, select Yes or No.



## Problem 9

Consider the following plot of data in  $d = 2$  dimensions. We'd like to use  $k$ -means clustering to cluster the data into  $k = 3$  clusters.

Suppose the crosses represent initial centroids, which are not themselves data points.

### Problem 9.1

Which of the following facts are true about the cluster assignment during the first iteration, as determined by these initial centroids?

- ☐ Exactly one cluster contains 11 data points.
- ☐ Exactly two clusters contain 11 data points.
- ☐ Exactly one cluster contains at least 12 data points.
- ☐ Exactly two clusters contain at least 12 data points.
- ☐ None of the above.



# Problem 9

Consider the following plot of data in  $d = 2$  dimensions. We'd like to use  $k$ -means clustering to cluster the data into  $k = 3$  clusters.

Suppose the crosses represent initial centroids, which are not themselves data points.

## Problem 9.2

The cross shapes in the plot above represent positions of the initial centroids before the first iteration. Now the algorithm is run for one iteration, after which the centroids have been adjusted.

We are now starting the second iteration. Which of the following facts are true about the cluster assignment during the **second** iteration? Select all that apply.

- ☐ Exactly one cluster contains 11 data points.
- ☐ Exactly two clusters contain 11 data points.
- ☐ Exactly one cluster contains at least 12 data points.
- ☐ Exactly two clusters contain at least 12 data points.
- ☐ None of the above.



## Problem 9

Consider the following plot of data in  $d = 2$  dimensions. We'd like to use  $k$ -means clustering to cluster the data into  $k = 3$  clusters.

Suppose the crosses represent initial centroids, which are not themselves data points.

### Problem 9.3

Compare the value of inertia after the end of the second iteration to the value of inertia at the end of the first iteration. Which of the following facts are true? Select all that apply.

- ☐ The inertia at the end of the second iteration is lower.
- ☐ The inertia doesn't decrease since there are actually 4 clusters in the data but we are using  $k$ -means with  $k = 3$ .
- ☐ The inertia doesn't decrease since there are actually 5 clusters in the data but we are using  $k$ -means with  $k = 3$ .
- ☐ The inertia doesn't decrease since there is an outlier that does not belong to any cluster.
- ☐ The inertia at the end of the second iteration is the same as at the end of the first iteration.



