- Our first classification example will involve predicting whether or not a patient has diabetes, given other information about their health.

```
In [2]:  1 diabetes = pd.read_csv('data/diabetes.csv')
         2 display_df(diabetes, cols=9)
```

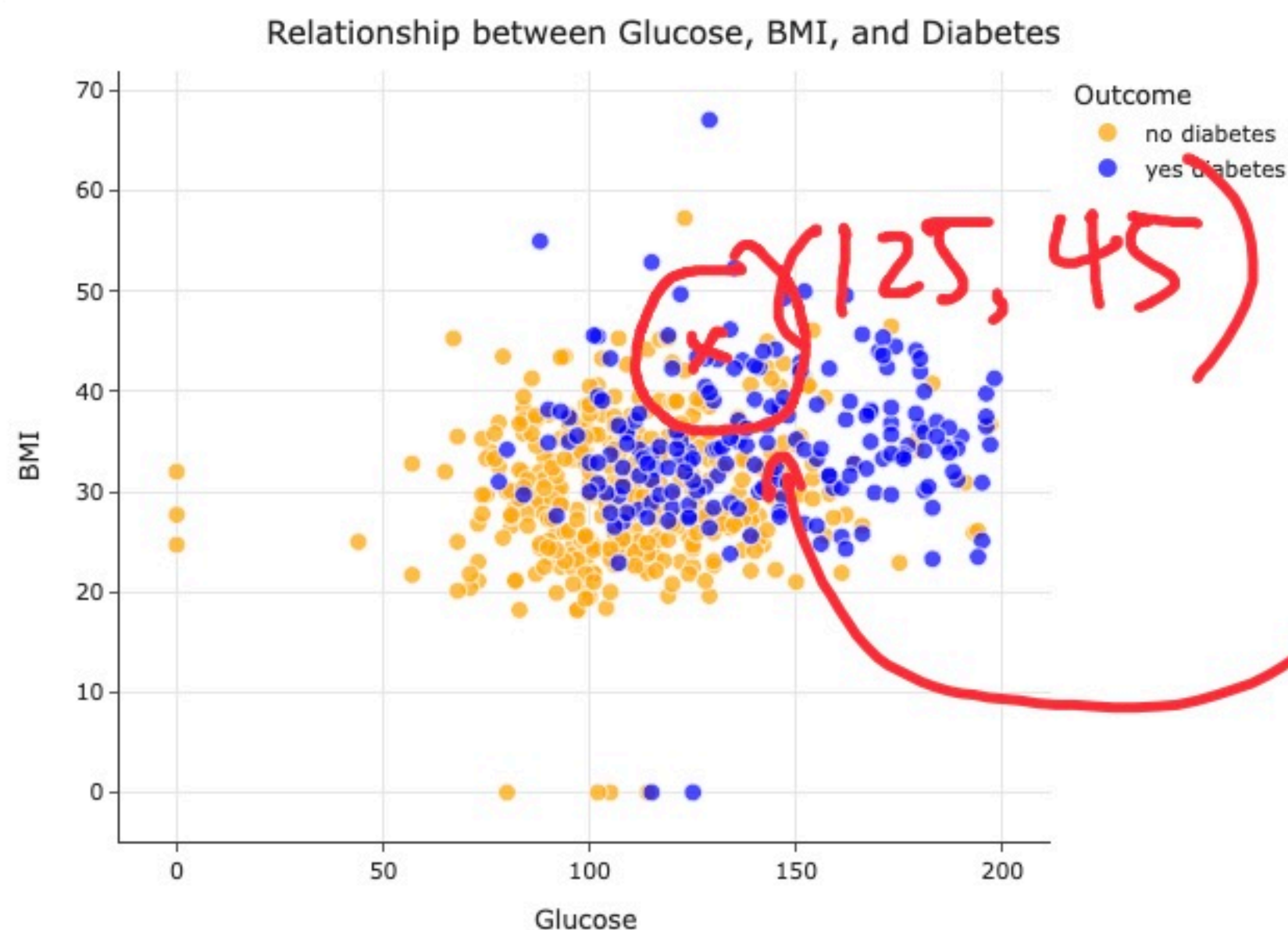| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.63 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.35 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.67 | 32 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.24 | 30 | 0 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.35 | 47 | 1 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.32 | 23 | 0 |

$y$

768 rows × 9 columns

```
In [3]:  1 # 0 means no diabetes, 1 means yes diabetes.
         2 diabetes['Outcome'].value_counts()
```

```
Out[3]: Outcome
        0    500
        1    268
        Name: count, dtype: int64
```

- Since there are only two possible `'Outcome'`s, we can draw a 2D scatter plot of `'BMI'` vs. `'Glucose'` and color each point by `'Outcome'`. Below, **class 0 (orange) is "no diabetes"** and **class 1 (blue) is "diabetes"**.

```
In [6]:  1  fig = util.create_base_scatter(X_train, y_train)
         2  fig
```



**Relationship between Glucose, BMI, and Diabetes**

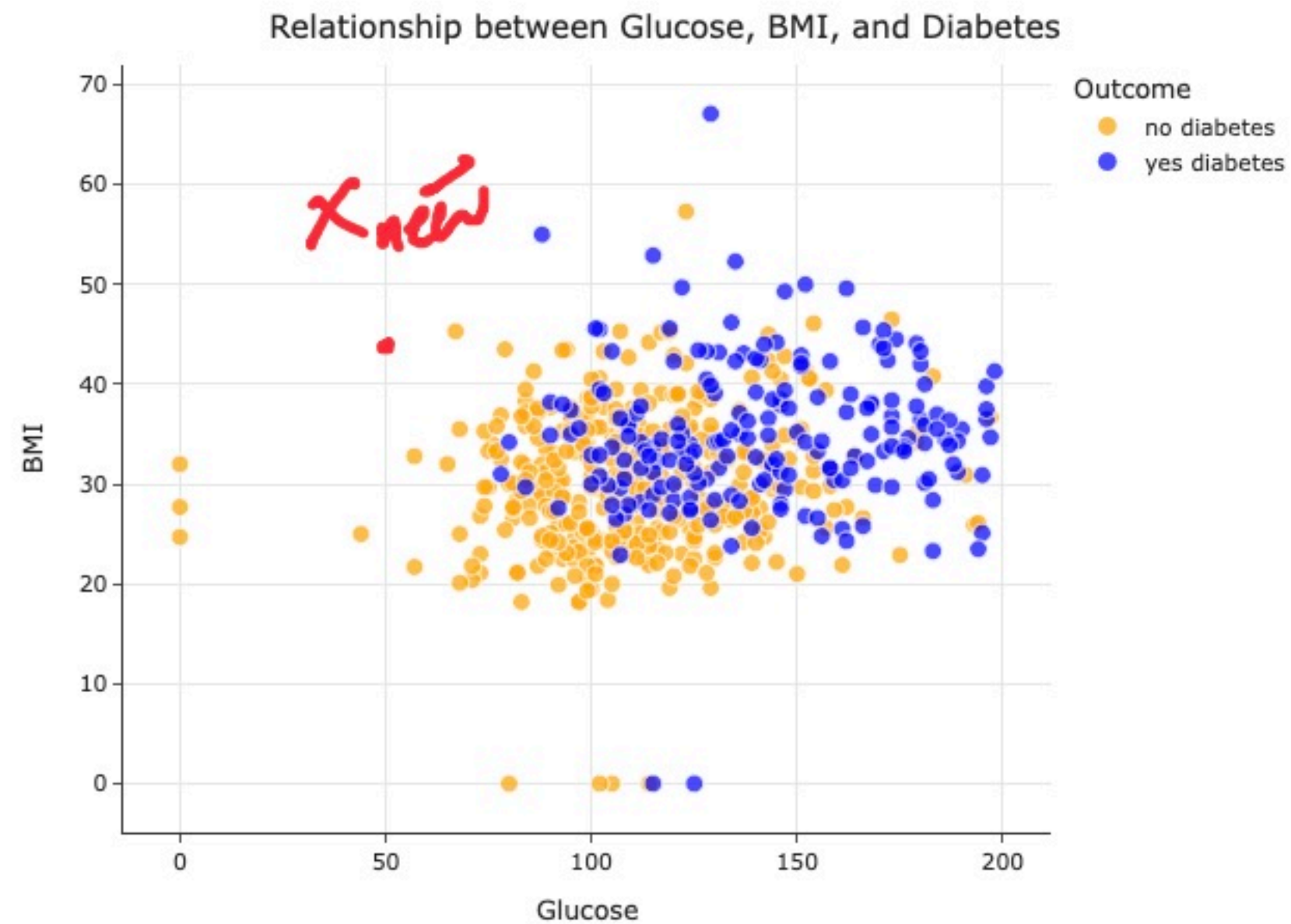(handwritten annotations: "x (125, 45)" circled; "look at people most similar to new patient")

- Using this dataset, how can we classify whether someone new (not already in the dataset) has diabetes, given their `'Glucose'` and `'BMI'`?

- **Intuition**: If a new person's feature vector is **close to the blue points**, we'll predict **blue (diabetes)**; if they're **close to the orange points**, we'll predict **orange (no diabetes)**.

9 . 1

2. Predicting that $\vec{x}_{\text{new}}$ belongs to the **most common class** among those $k$ closest points.

In [7]:
```
1 fig
```
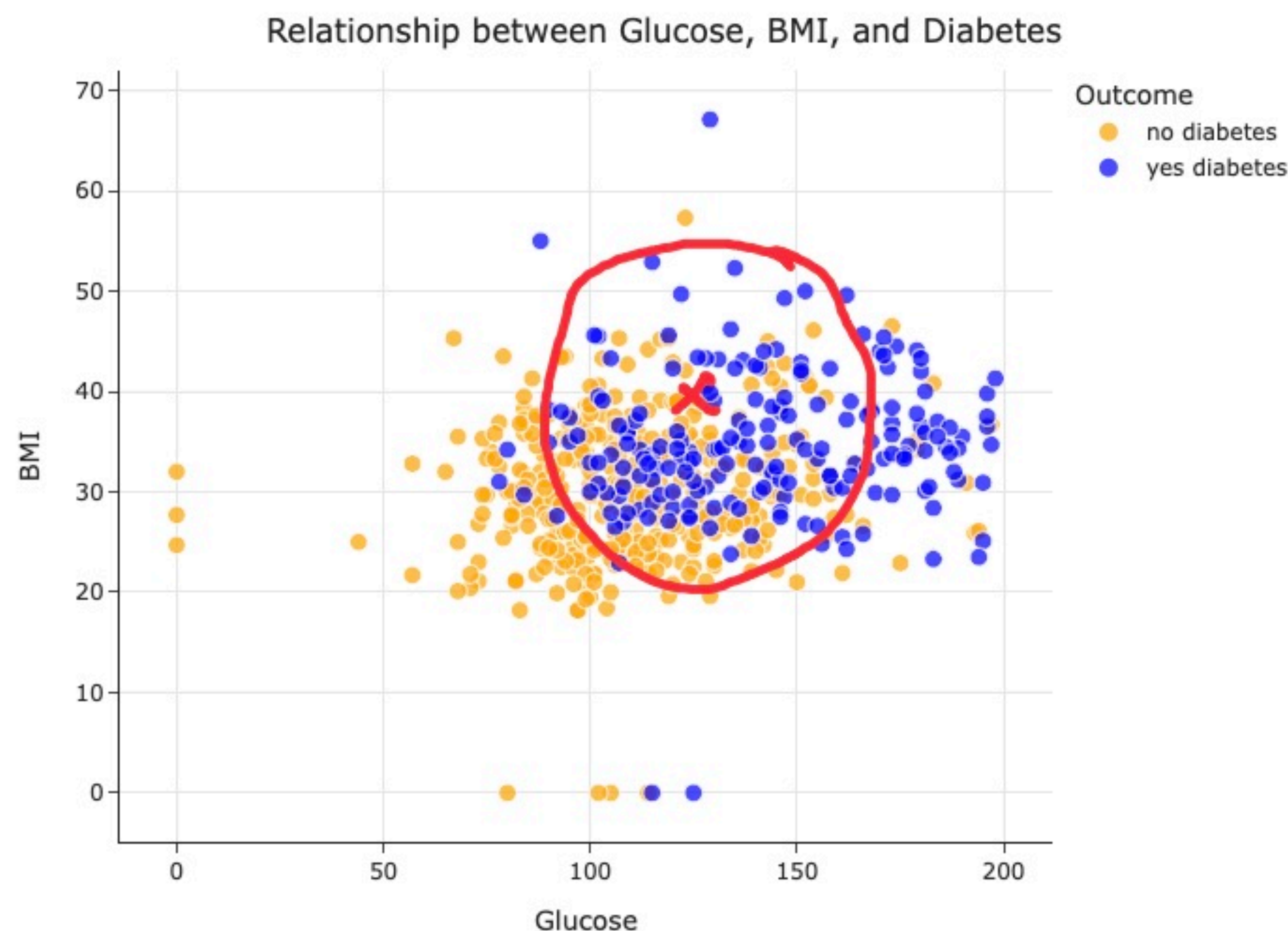


Relationship between Glucose, BMI, and Diabetes

- Example: Suppose $k = 6$. If, among the 6 closest points to $\vec{x}_{\text{new}}$, there are **4 blue** and **2 orange** points, we'd predict **blue (diabetes)**.

  What if there are ties? Read here.

```
In [11]:  1  # To know what reasonable values for 'Glucose' and 'BMI' might be, let's look at the plot again.
          2  fig
```

**Relationship between Glucose, BMI, and Diabetes**



Among the 28 closest points to (125, 40), the majority did not have diabetes.

```
In [12]:  1  model_knn.predict(pd.DataFrame([{
          2      'Glucose': 125,
          3      'BMI': 40
          4  }]))
```

```
Out[12]:  array([0])
```

- What does the resulting model **look like** 👀? Can we visualize it?

- The most common evaluation metric in classification is **accuracy**:

$$\text{accuracy} = \frac{\#\text{ points classified correctly}}{\#\text{ points}}$$

Accuracy ranges from 0 to 1, i.e. 0% to 100%. **Higher** values indicate **better** model performance.

*if model is a regression model, model.score returns $R^2$.*

```
In [28]:  1 # Equivalent to 75%.
          2 (model_knn.predict(X_test) == y_test).mean()
```

Out[28]: 0.75

- This is the default metric that the `score` method of a classifier computes, too.

```
In [29]:  1 model_knn.score(X_test, y_test)
```

Out[29]: 0.75

```
In [ ]:  1 # For future reference.
         2 test_scores = pd.Series()
         3 ...
         4 test_scores
```
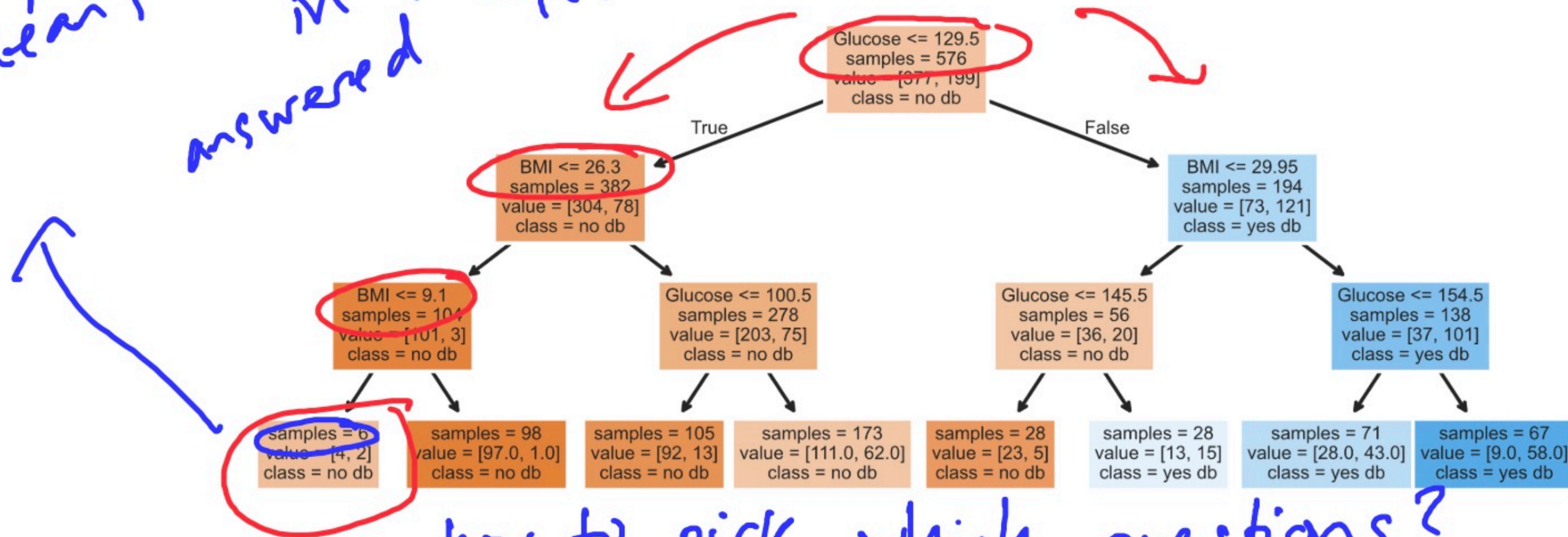
k like:

samples = 6

6 people training set
means answered "yes", then "yes",
in the "yes", then "yes"



Glucose <= 129.5
samples = 576
value = [377, 199]
class = no db

True

False

BMI <= 26.3
samples = 382
value = [304, 78]
class = no db

BMI <= 29.95
samples = 194
value = [73, 121]
class = yes db

BMI <= 9.1
samples = 104
value = [101, 3]
class = no db

Glucose <= 100.5
samples = 278
value = [203, 75]
class = no db

Glucose <= 145.5
samples = 56
value = [36, 20]
class = no db

Glucose <= 154.5
samples = 138
value = [37, 101]
class = yes db

samples = 6
value = [4, 2]
class = no db

samples = 98
value = [97.0, 1.0]
class = no db

samples = 105
value = [92, 13]
class = no db

samples = 173
value = [111.0, 62.0]
class = no db

samples = 28
value = [23, 5]
class = no db

samples = 28
value = [13, 15]
class = yes db

samples = 71
value = [28.0, 43.0]
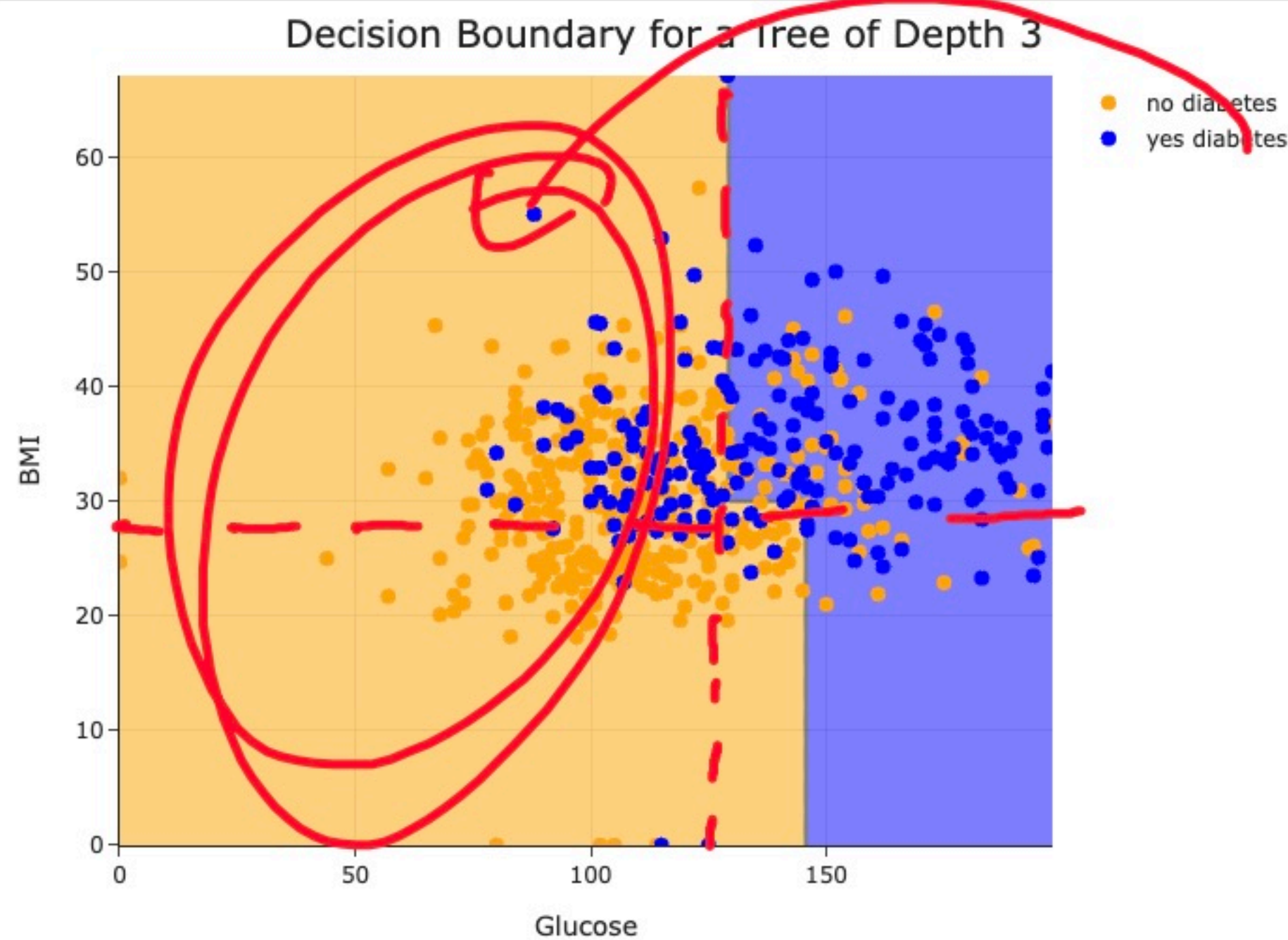class = yes db

samples = 67
value = [9.0, 58.0]
class = yes db

how to pick which questions?
Gini impurity, entropy → probability

# Decision boundaries for a decision tree classifier

```
In [51]:  1 util.show_decision_boundary(model_tree, X_train, y_train, title='Decision Boundary for a Tree of Depth 3')
```



Decision Boundary for a Tree of Depth 3

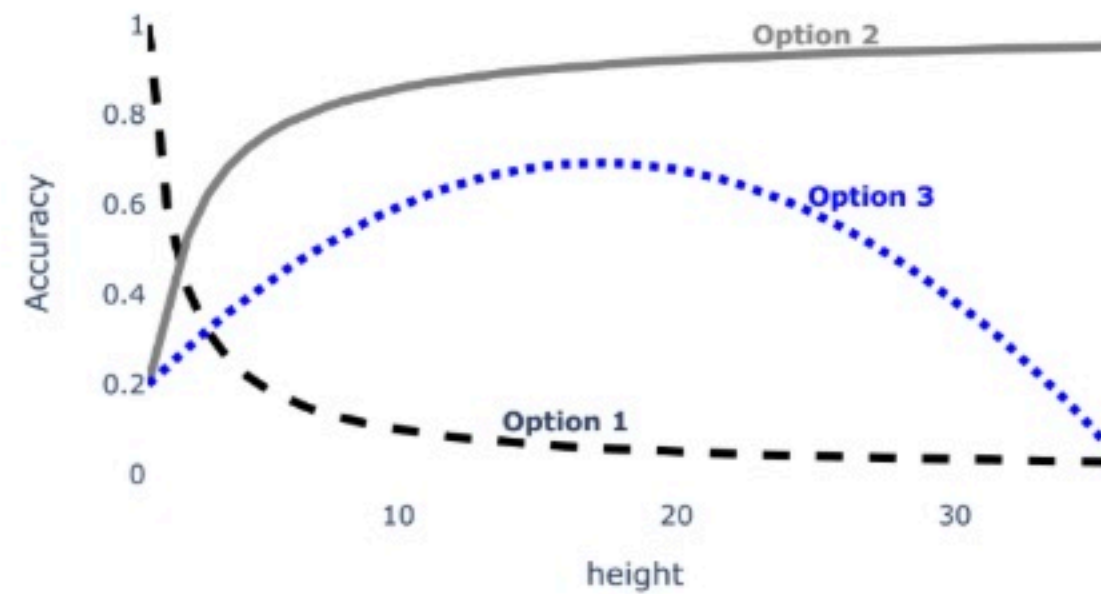*Handwritten annotation:* answered "Yes" to  Glucose <= 129.5

- Observe that the decision boundaries – at least when we set `max_depth` to 3 – look less "jagged" than with the $k$-NN classifier.

  Decision trees partition the feature space into rectangles.

ChickenClassifiers have many hyperparameters, one of which is `height`. As we increase the value of `height`, the model variance of the resulting `ChickenClassifier` also increases.

First, we consider the training and testing accuracy of a `ChickenClassifier` trained using various values of `height`. Consider the plot below.



Which of the following depicts **training accuracy vs. `height`**?

○ Option 1

○ Option 2

○ Option 3

Which of the following depicts **testing accuracy vs. `height`**?

○ Option 1

○ Option 2

○ Option 3

*(Handwritten annotations:)*

accuracy high = good model!

as height ↑, complexity ↑, tendency to overfit ↑

- When performing **binary** classification, there are four possible outcomes.

  Note: A "positive prediction" is a prediction of 1, and a "negative prediction" is a prediction of 0.

| Outcome of Prediction | Definition | True Class |
|---|---|---|
| **True** positive (TP) ✅ | The predictor **correctly** predicts the positive class. | P |
| False negative (FN) ❌ | The predictor incorrectly predicts the negative class. | P |
| **True** negative (TN) ✅ | The predictor **correctly** predicts the negative class. | N |
| False positive (FP) ❌ | The predictor incorrectly predicts the positive class. | N |

- We typically organize the four quantities above into a **confusion matrix**.

| | Predicted Negative | Predicted Positive |
|---|---|---|
| **Actually Negative** | TN ✅ | FP ❌ |
| **Actually Positive** | FN ❌ | TP ✅ |

- Note that in the four acronyms – TP, FN, TN, FP – the **first letter** is whether the prediction is correct, and the **second letter** is what the prediction is.

*[Handwritten annotations: "FP" with arrow pointing to FP, "first letter tells you if prediction is correct"]*

# Example: Accuracy of COVID tests

- The results of 100 Michigan Medicine COVID tests are given below.

|  | Predicted Negative | Predicted Positive |
|---|---|---|
| **Actually Negative** | TN = 90 ✅ | FP = 1 ❌ |
| **Actually Positive** | FN = 8 ❌ | TP = 1 ✅ |

*Michigan Medicine test results*

- 🤔 **Question:** What is the accuracy of the test?

$$\text{accuracy} = \frac{\text{\# points classified correctly}}{\text{\# points}}$$

$$\frac{91}{91+9}$$

$$= \frac{91}{100}$$

$$= 91\% .$$

# Example: Accuracy of COVID tests

- The results of 100 Michigan Medicine COVID tests are given below.

|  | Predicted Negative | Predicted Positive |
|---|---|---|
| **Actually Negative** | TN = 90 ✅ | FP = 1 ❌ |
| **Actually Positive** | FN = 8 ❌ | TP = 1 ✅ |

*Michigan Medicine test results*

→ 91 people don't have COVID

only 9 people actually do

- 🤔 **Question:** What is the accuracy of the test?

$$\text{accuracy} = \frac{\text{\# points classified correctly}}{\text{\# points}}$$

- 🙋 **Answer:**

$$\text{accuracy} = \frac{TP + TN}{TP + FP + FN + TN} = \frac{1 + 90}{100} = 0.91$$

- **Followup:** At first, the test seems good. But, suppose we build a classifier that predicts that **nobody has COVID**. What would its accuracy be?

- **Answer to followup:** Also 0.91! There is severe **class imbalance** in the dataset, meaning that most of the data points are in the

# Precision

| | Predicted Negative | Predicted Positive |
|---|---|---|
| **Actually Negative** | TN = 0 ✅ | FP = 91 ❌ |
| **Actually Positive** | FN = 0 ❌ | TP = 9 ✅ |

*everyone-has-COVID classifier*

- The **precision** of a binary classifier is the proportion of **predicted positive instances** that are correctly classified. We'd like this number to be as close to 1 (100%) as possible.

$$\text{precision} = \frac{TP}{\#\text{ predicted positive}} = \frac{TP}{TP + FP}$$

- To compute precision, look at the **right (positive) column** of the above confusion matrix.

  **Tip:** A good way to remember the difference between precision and recall is that in the denominator for 🅿recision, both terms have 🅿 in them (TP and FP).