# Regression Tree
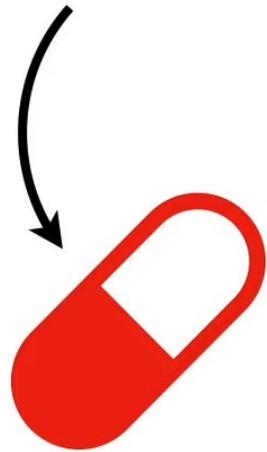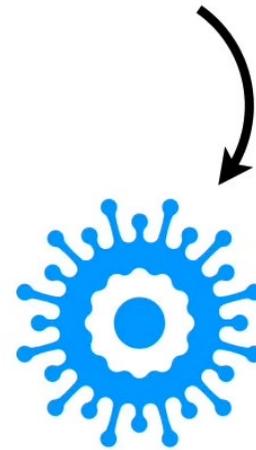
Imagine we developed a new drug… vs. …to cure the common cold.
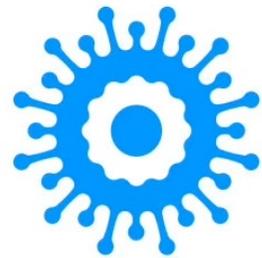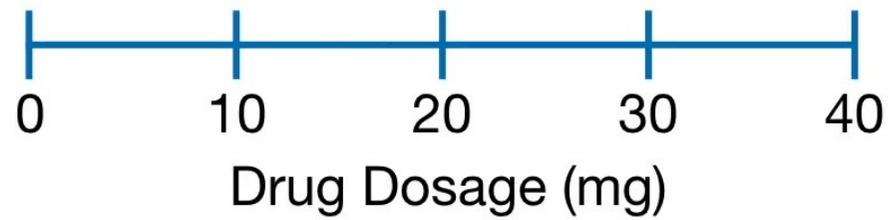
However, we don't know the optimal dosage to give to patients.

So we do a clinical trial with
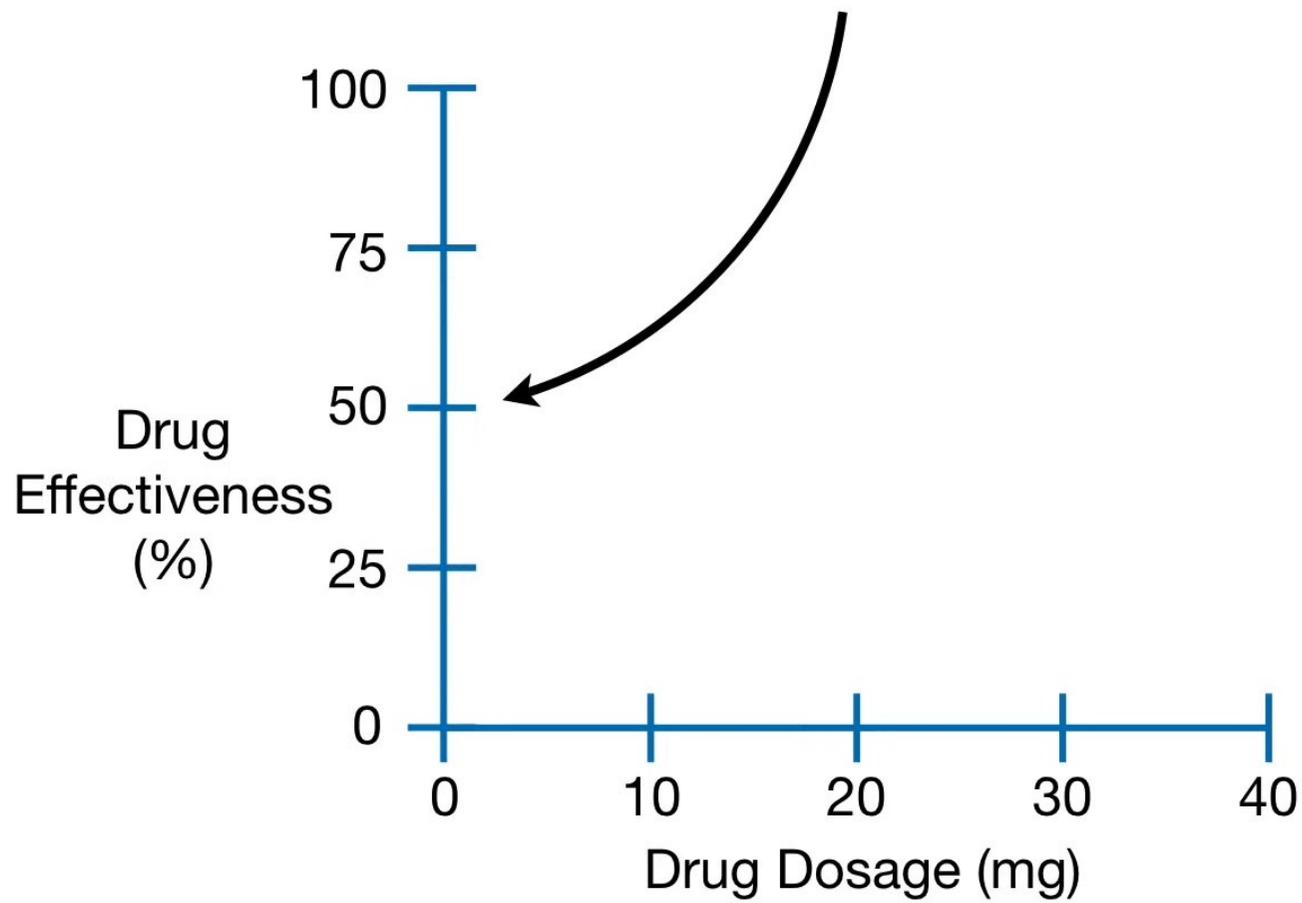different dosages…

0      10      20      30      40
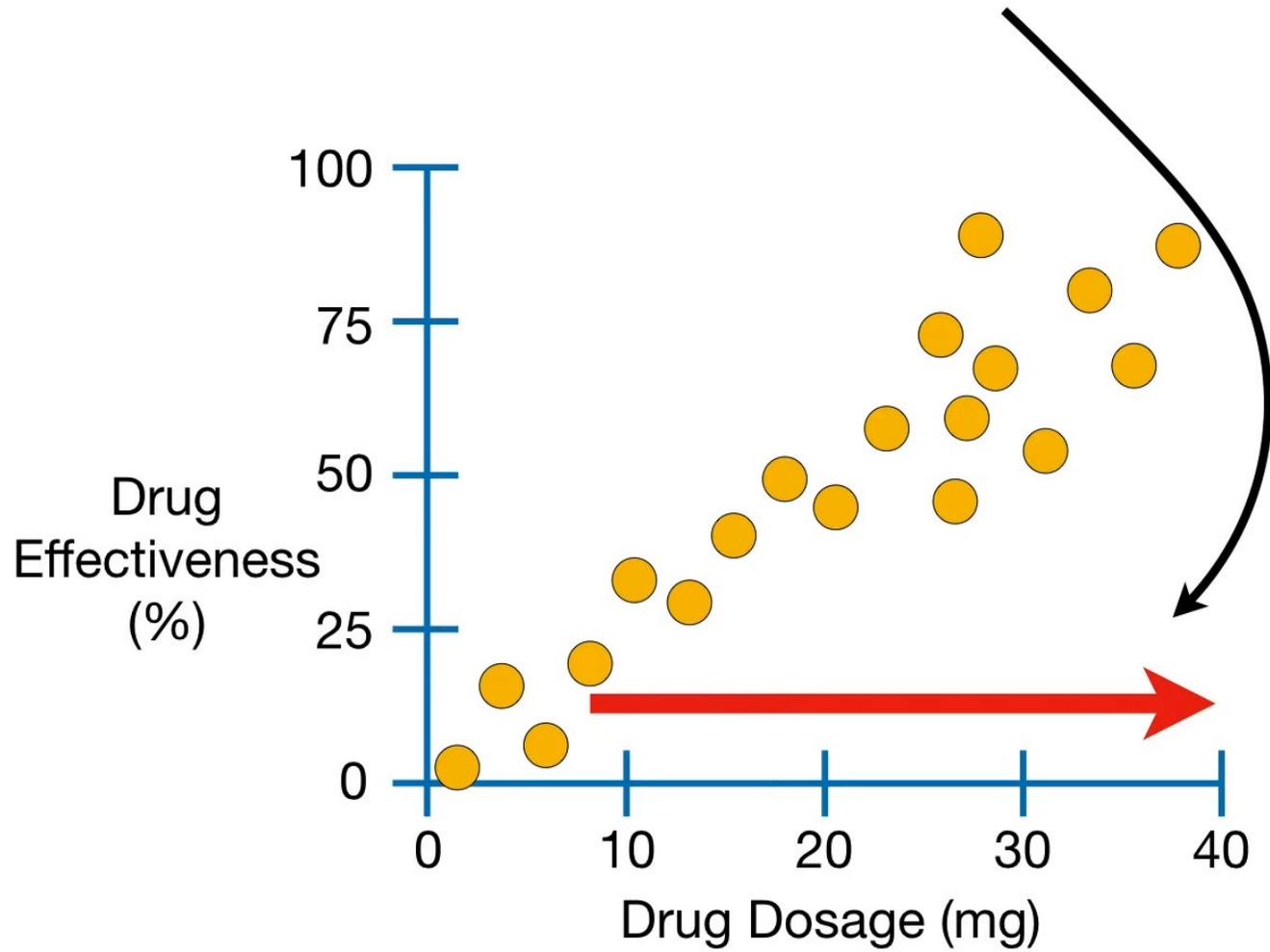
Drug Dosage (mg)

...and measure how effective each dosage is.

Drug Effectiveness (%)

100
75
50
25
0

Drug Dosage (mg)

0    10    20    30    40
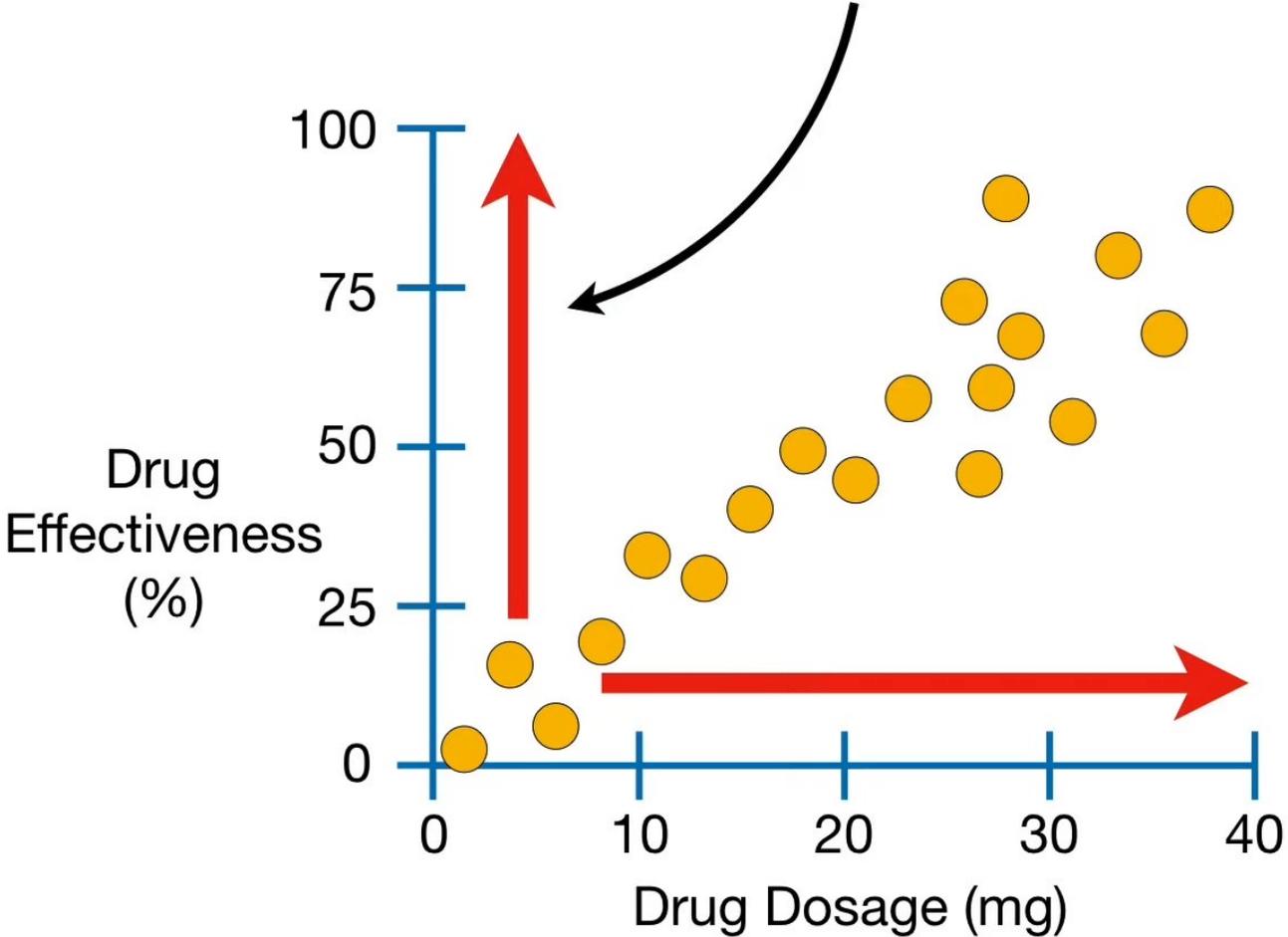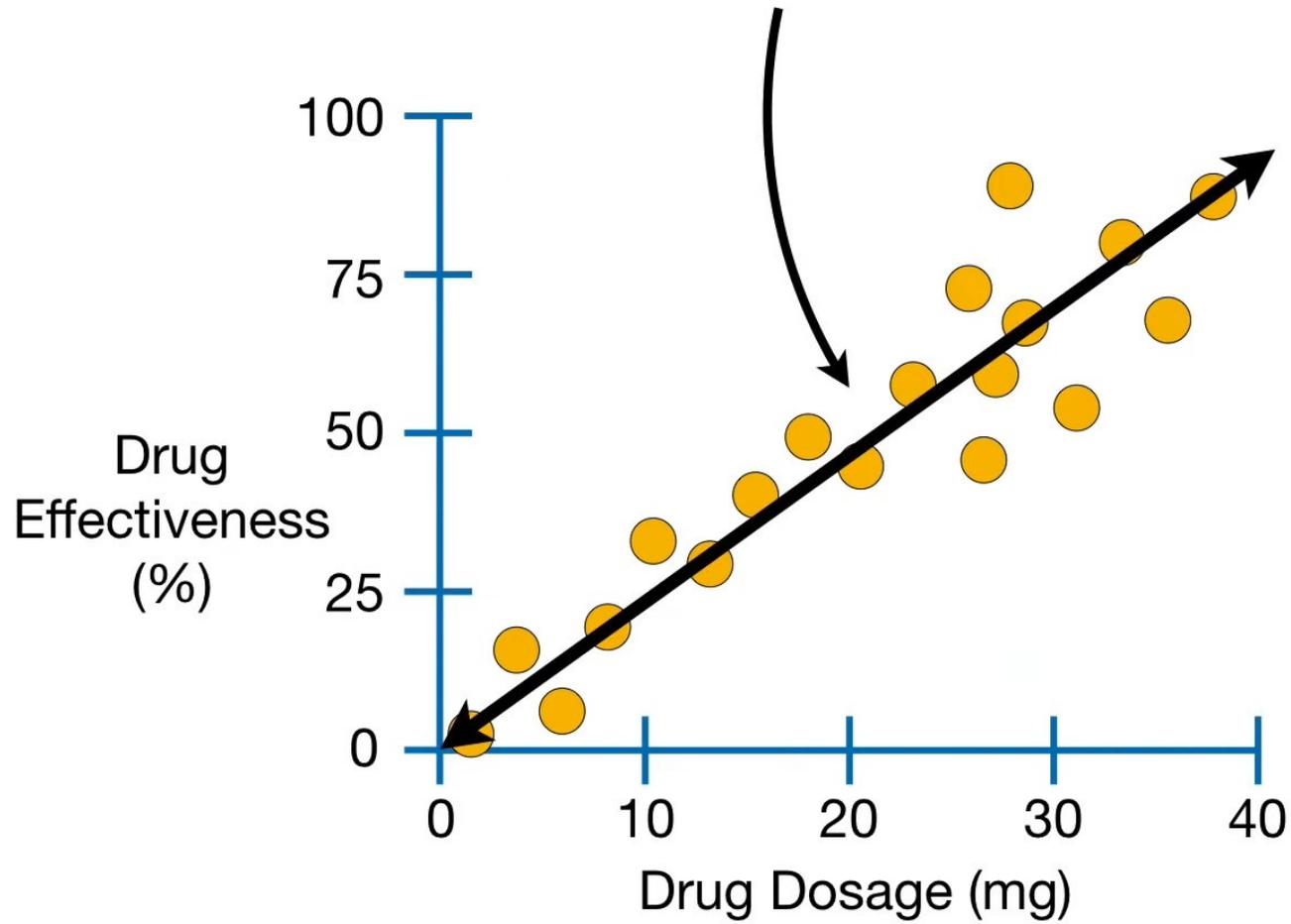
If the data looked like this…

...and, in general, the higher the dose,

…and, in general, the higher the dose, the more effective the drug…

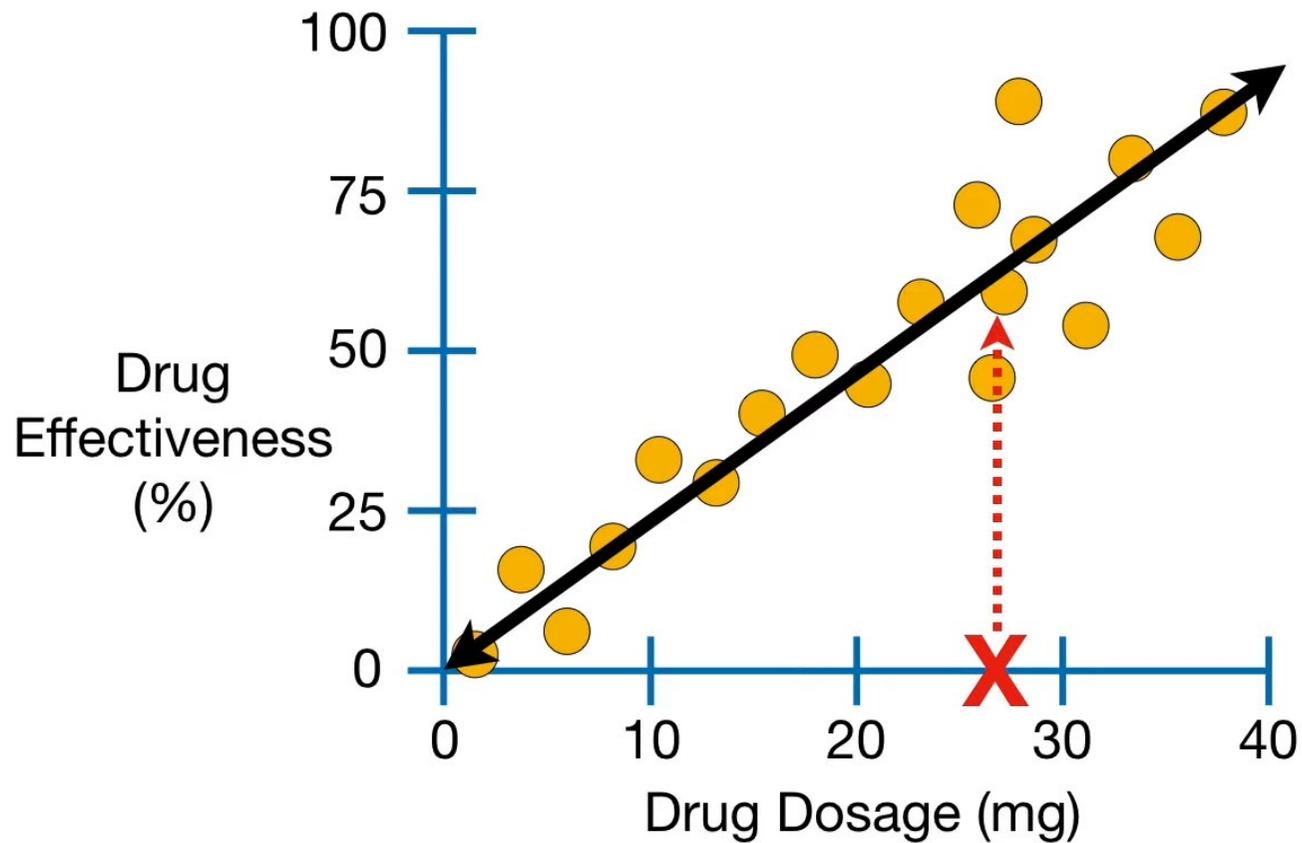…we could use the line to predict that a **27 mg Dose** should be **62% Effective**.
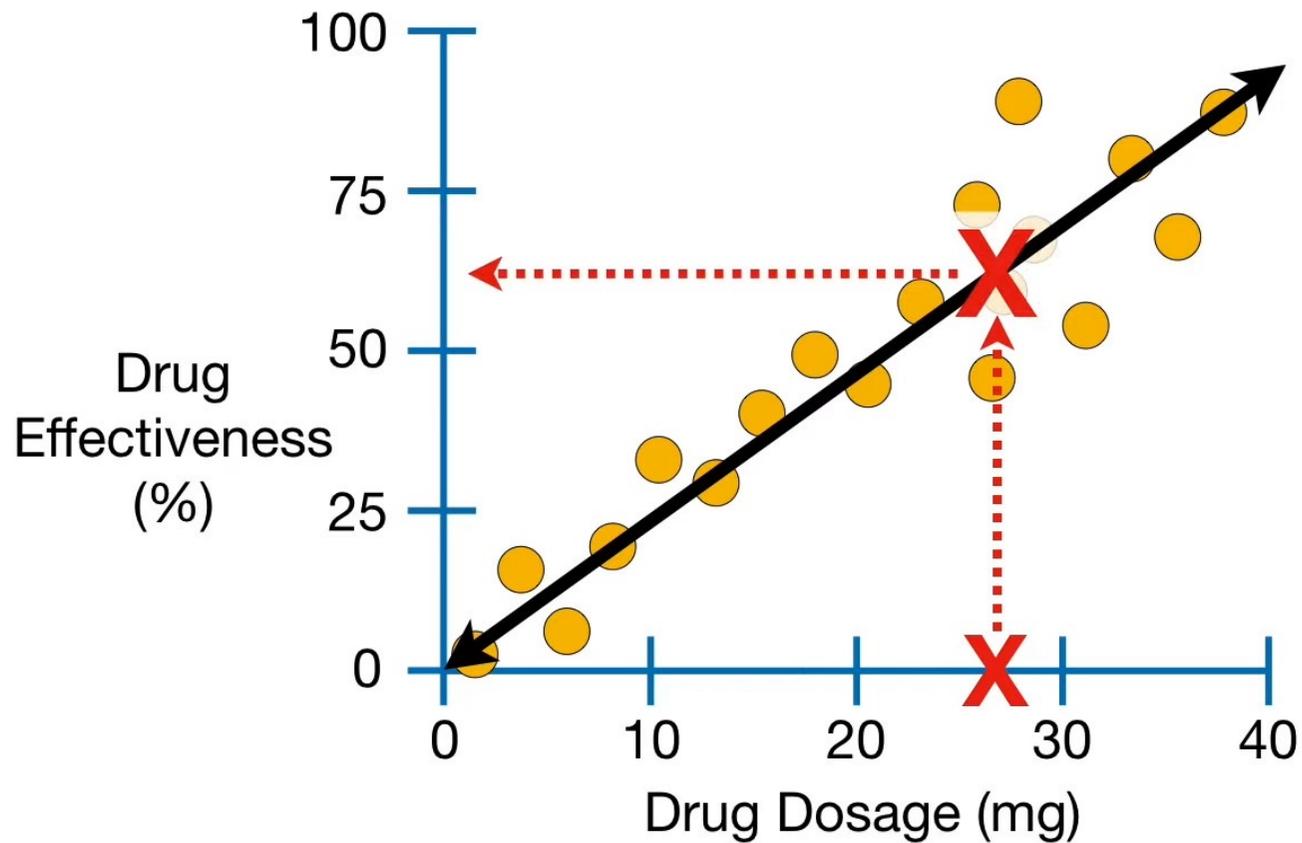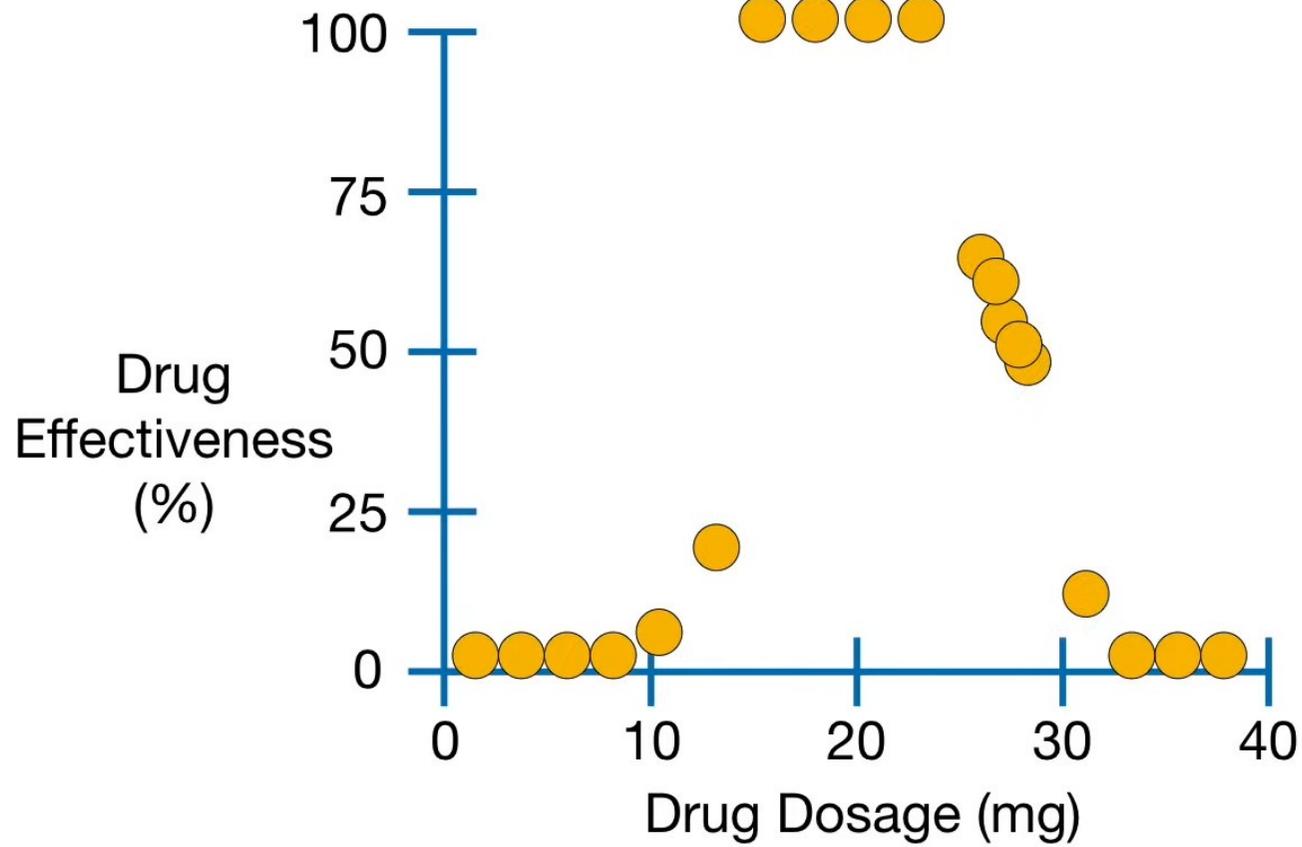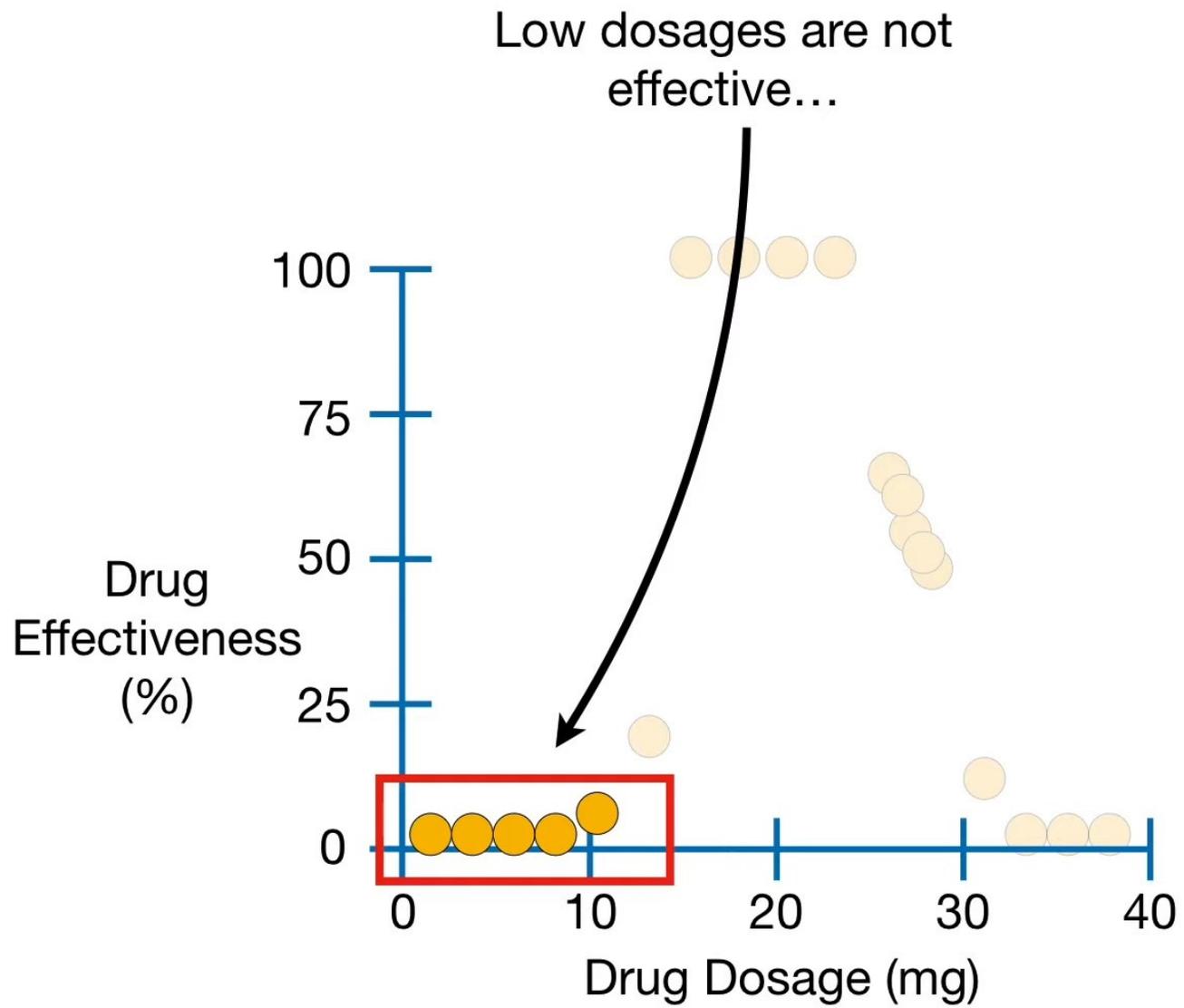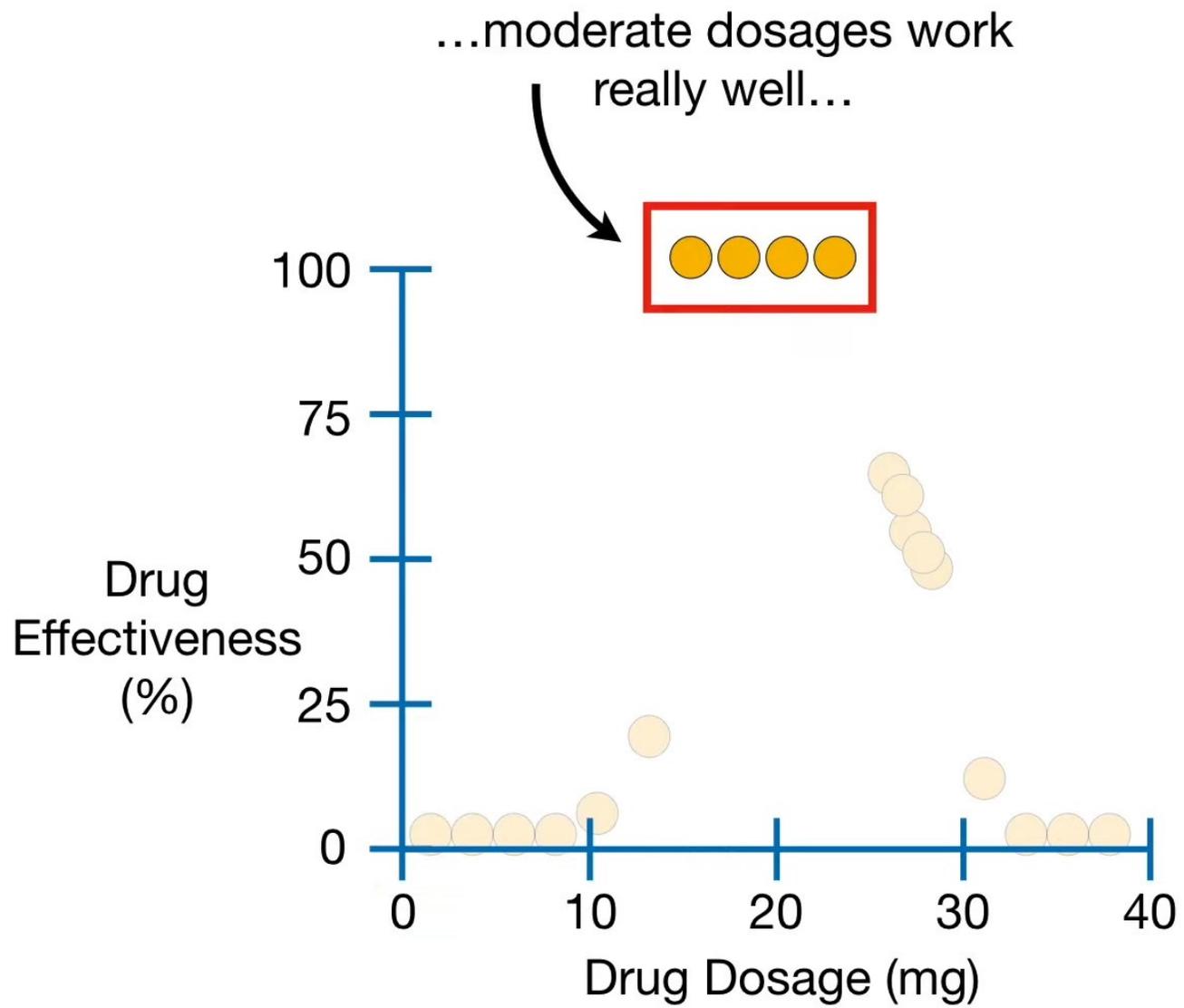
...we could use the line to predict that a **27 mg Dose** should be **62% Effective**.

However, what if the data looked like this?

Drug Effectiveness (%)

Drug Dosage (mg)

Low dosages are not effective…

...moderate dosages work really well...

Drug Effectiveness (%)

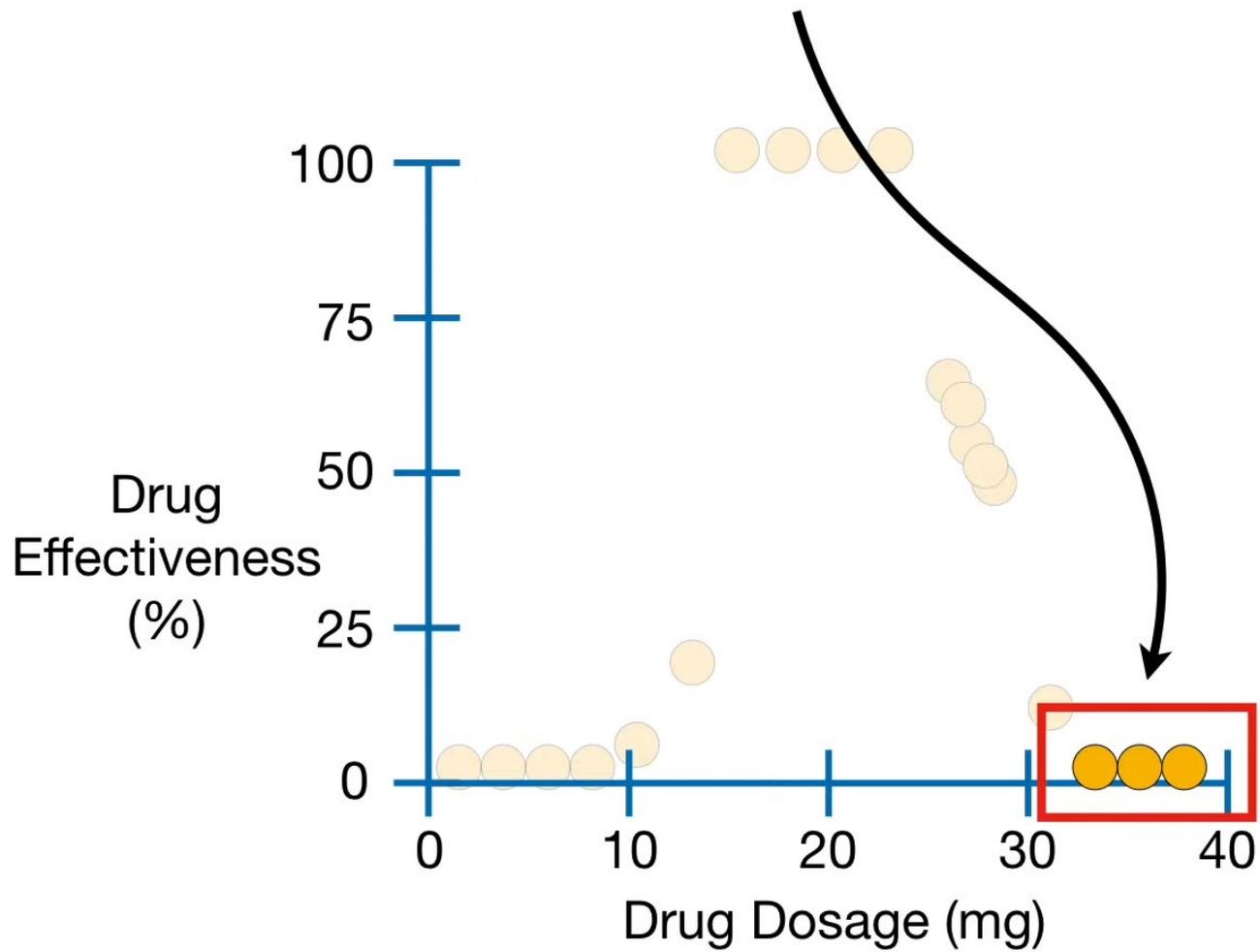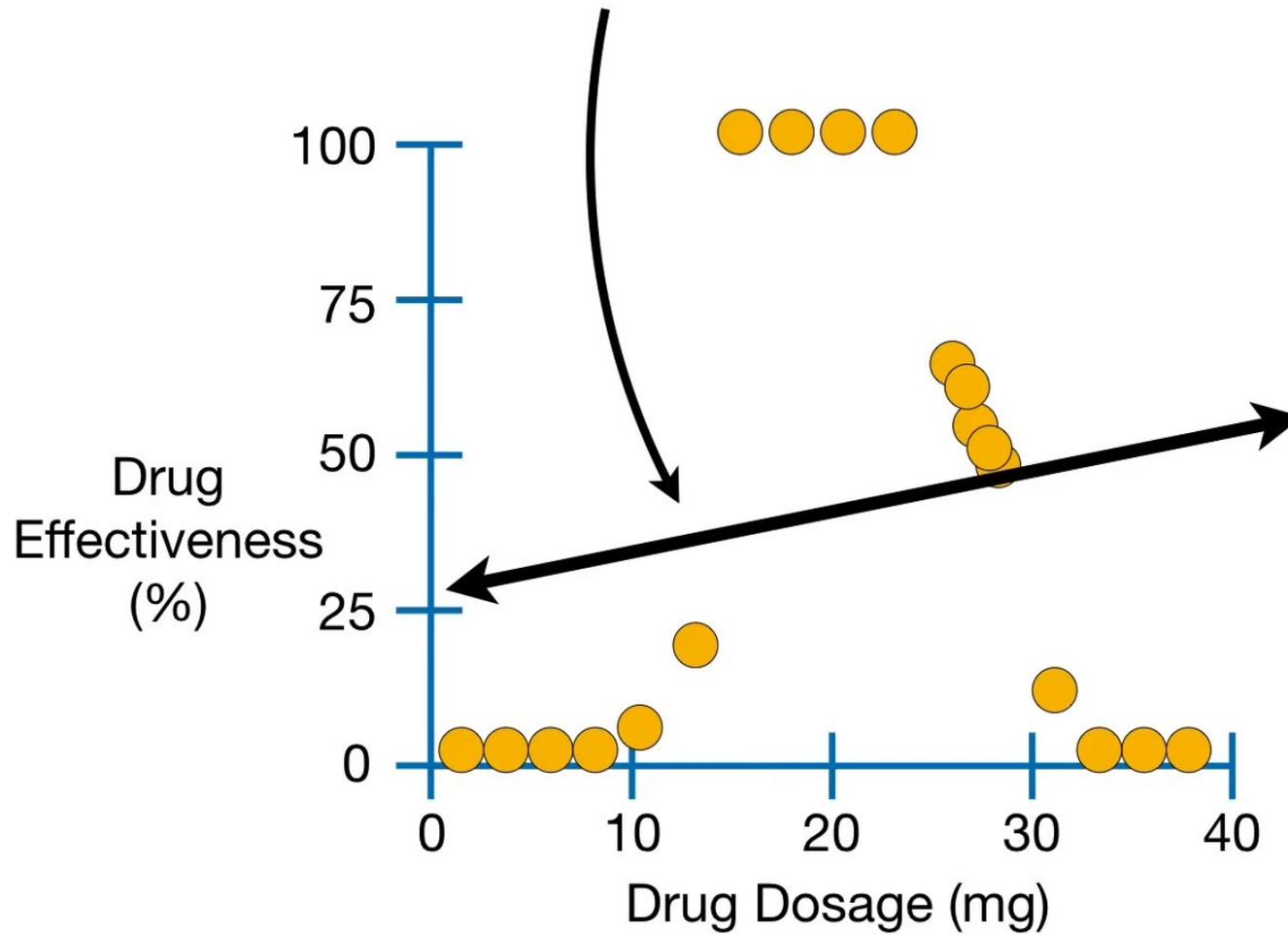Drug Dosage (mg)

...somewhat higher dosages work at about **50%** effectiveness...

...and high dosages are not effective at all.

Drug Effectiveness (%)

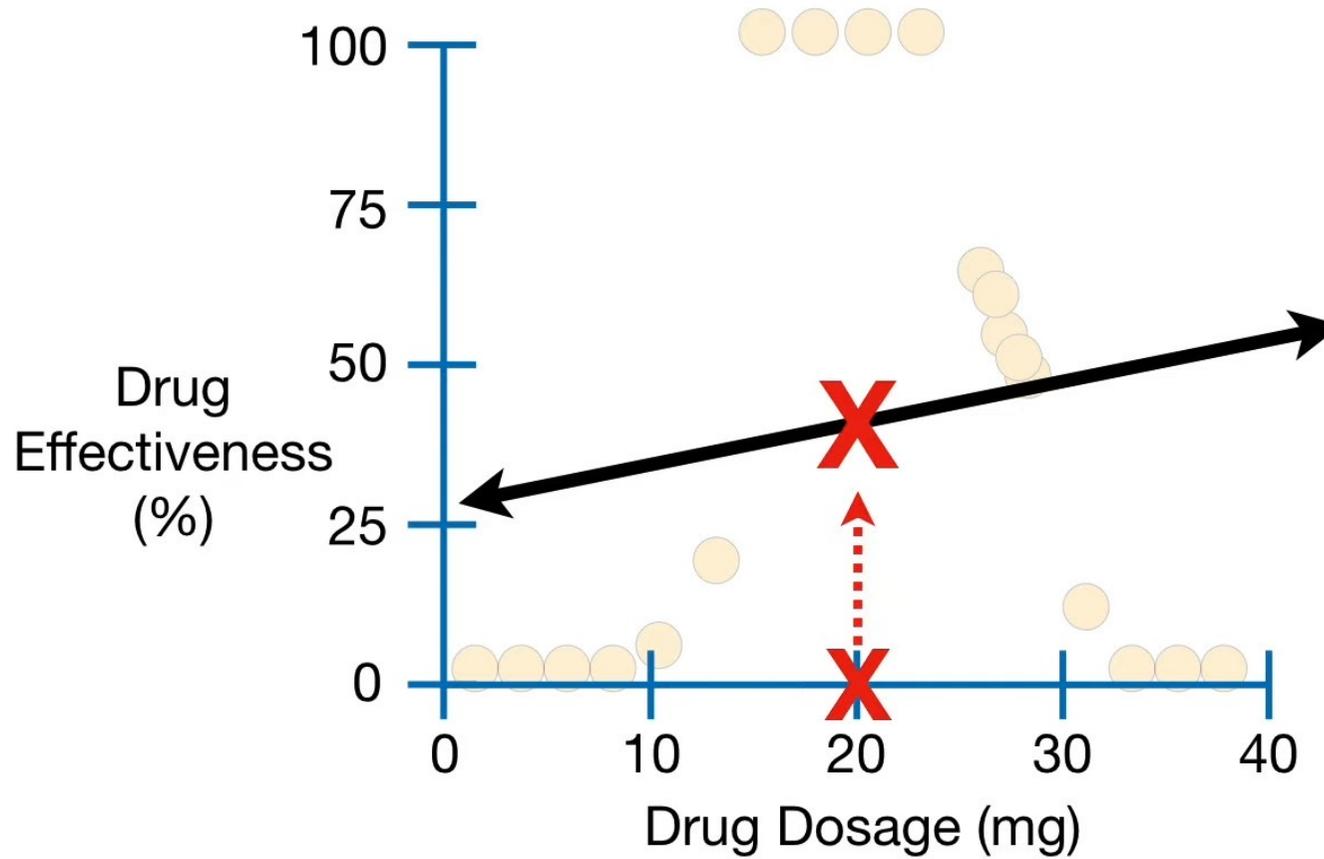Drug Dosage (mg)

In this case, fitting a straight line to the data will not be very useful.

For example, if someone told us they were taking a **20 mg Dose**…

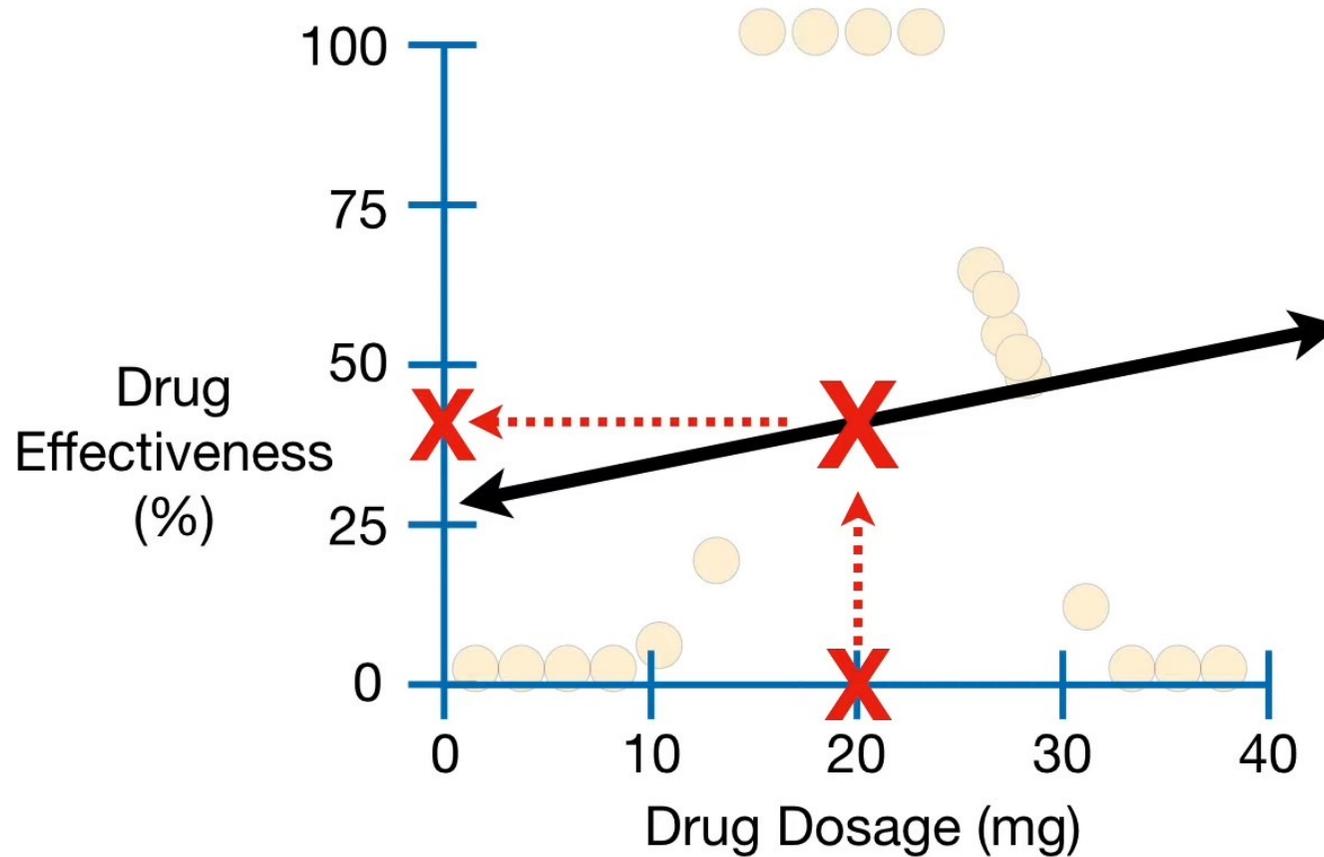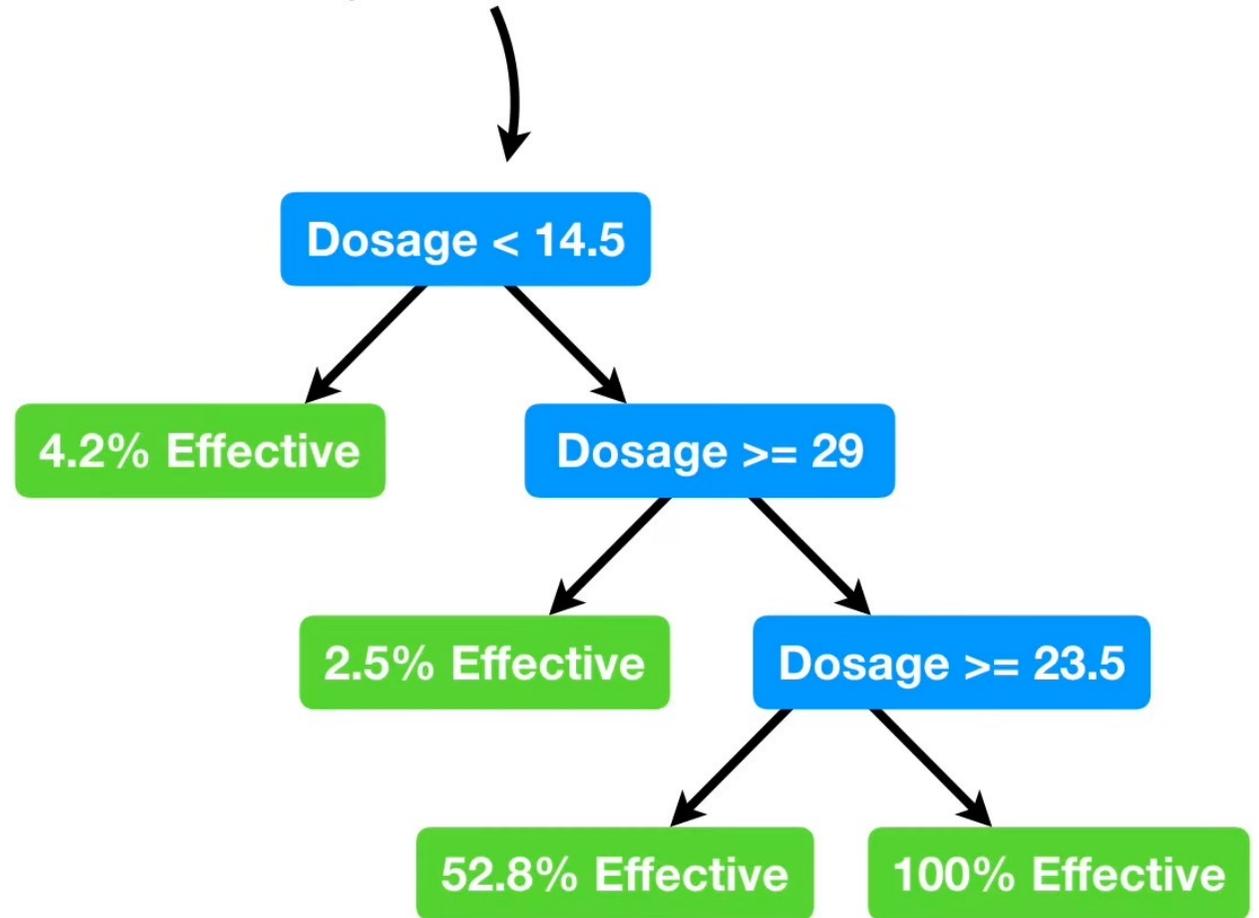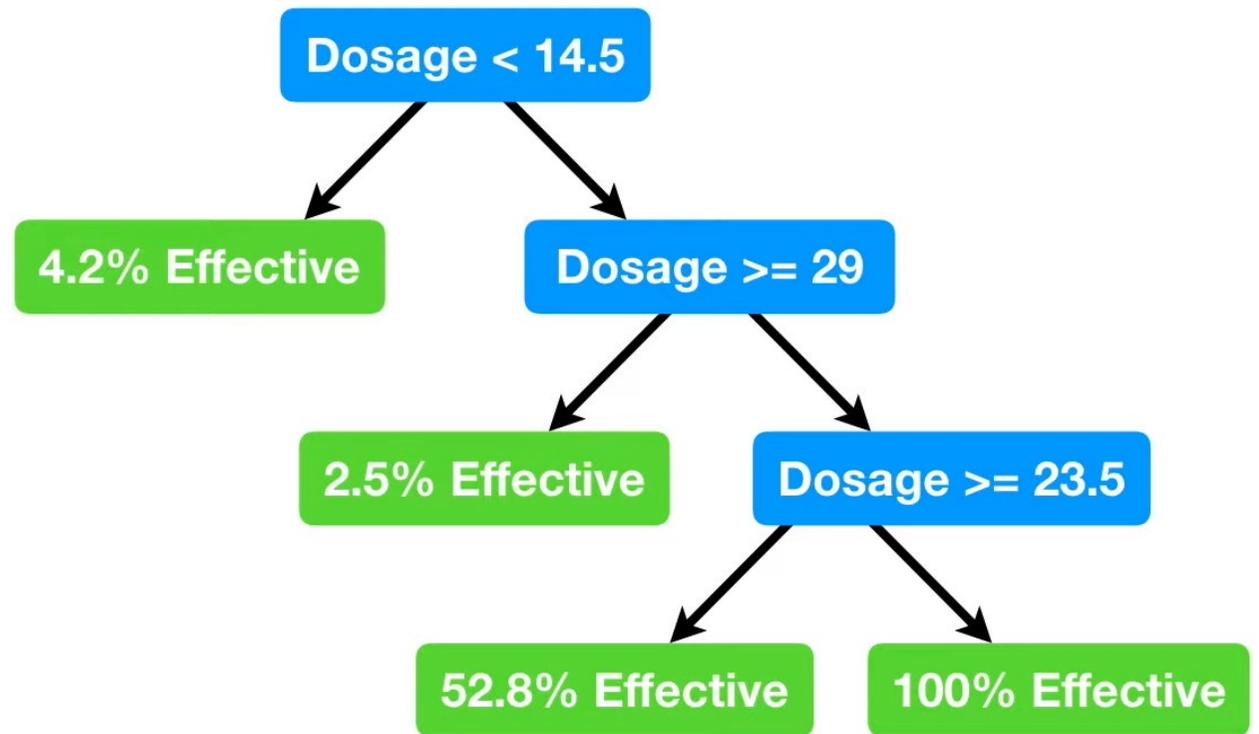…even though the observed data says that it should be **100% Effective**.

So we need to use something other
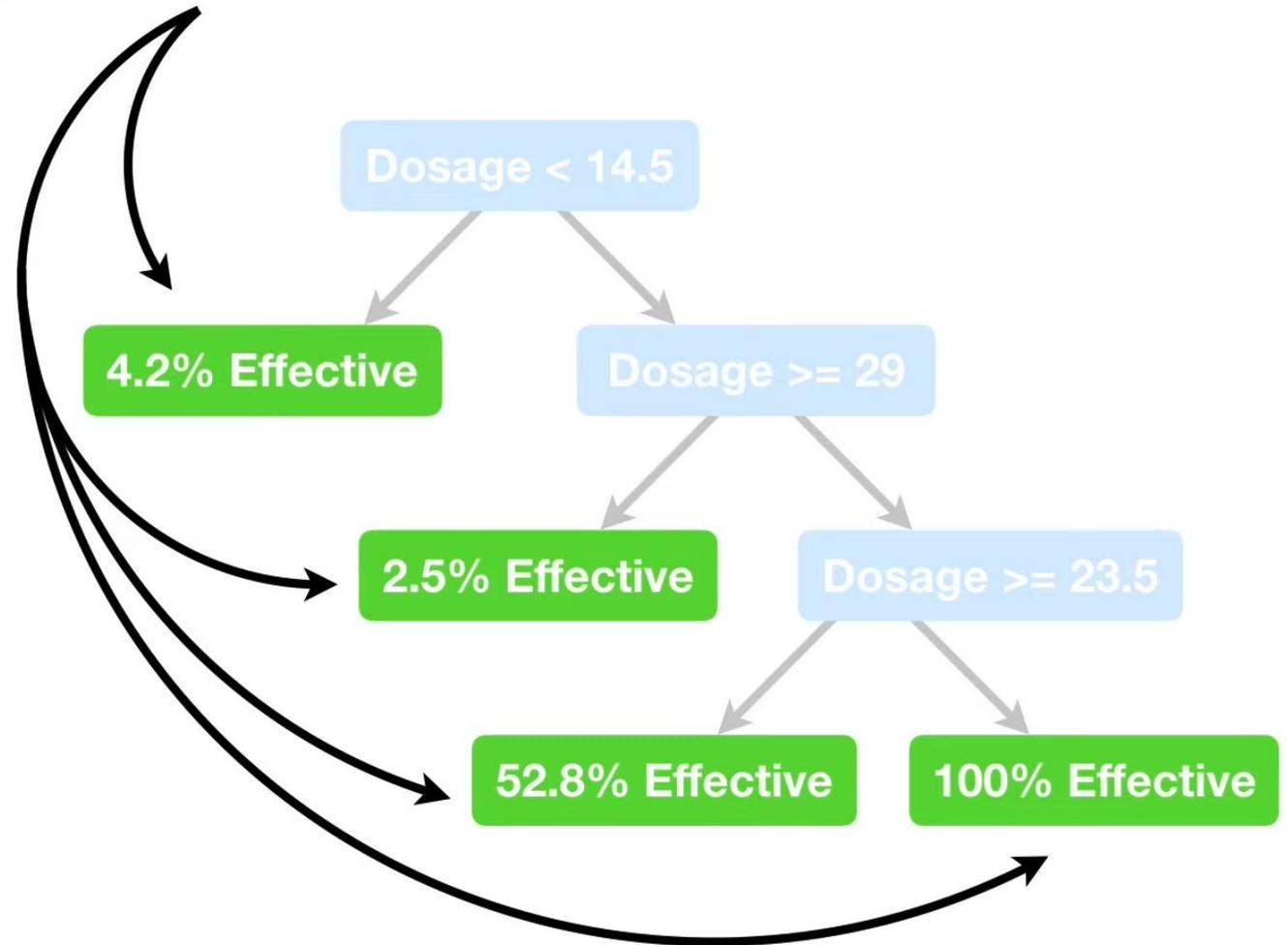than a straight line to make predictions.

One option is to use a **Regression Tree**.

**Regression Trees** are a type of **Decision Tree.**

In a **Regression Tree**, each leaf represents a numeric value.

Dosage < 14.5

4.2% Effective

Dosage >= 29

2.5% Effective

Dosage >= 23.5

52.8% Effective

100% Effective

**Has Hairy Toes**
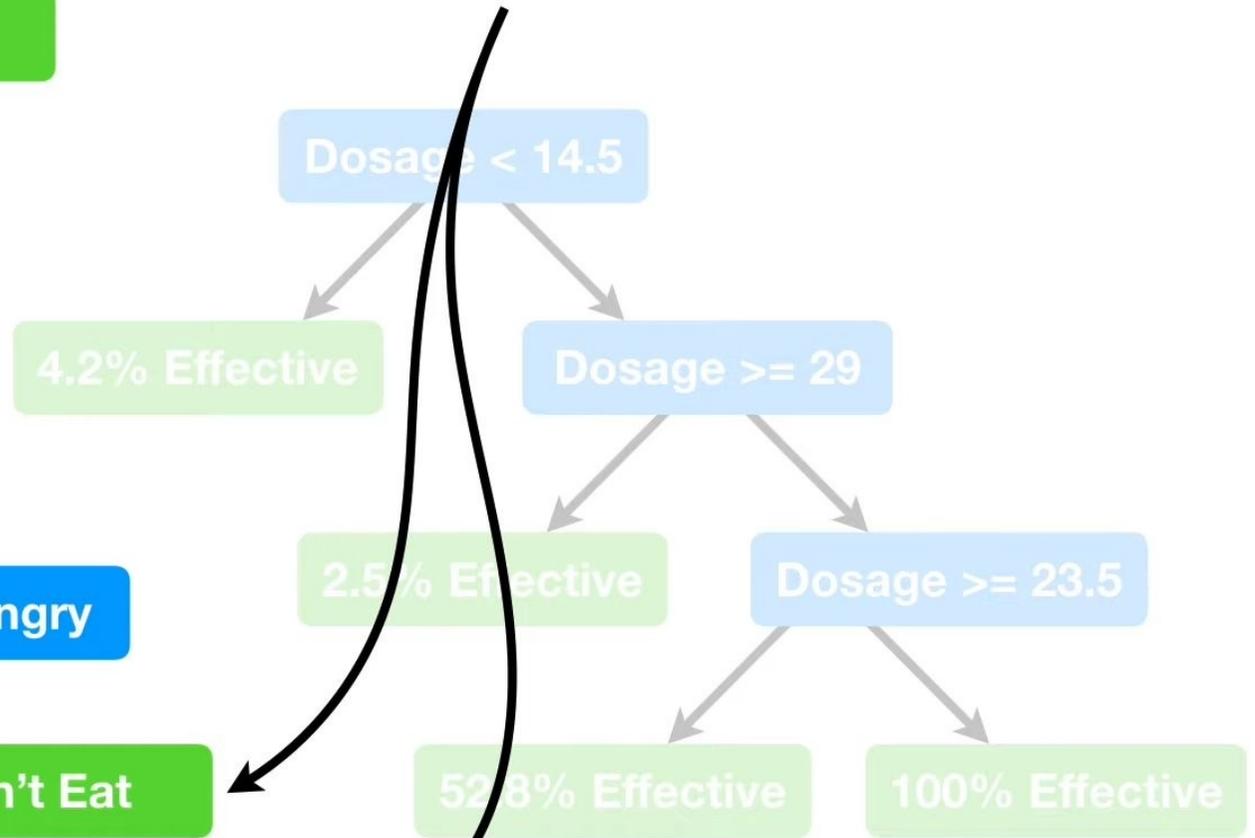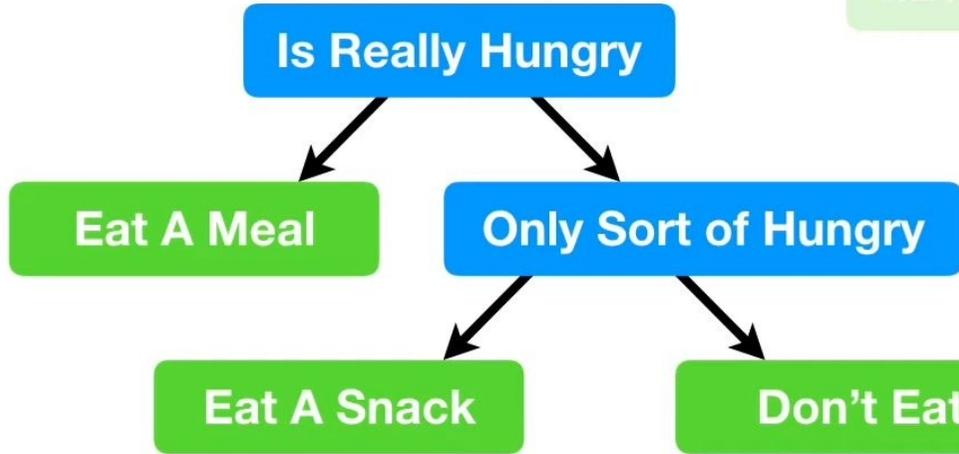
**True**  **False**

In contrast, **Classification Trees** have **True** or **False** in their leaves…

**Dosage < 14.5**

4.2% Effective  **Dosage >= 29**

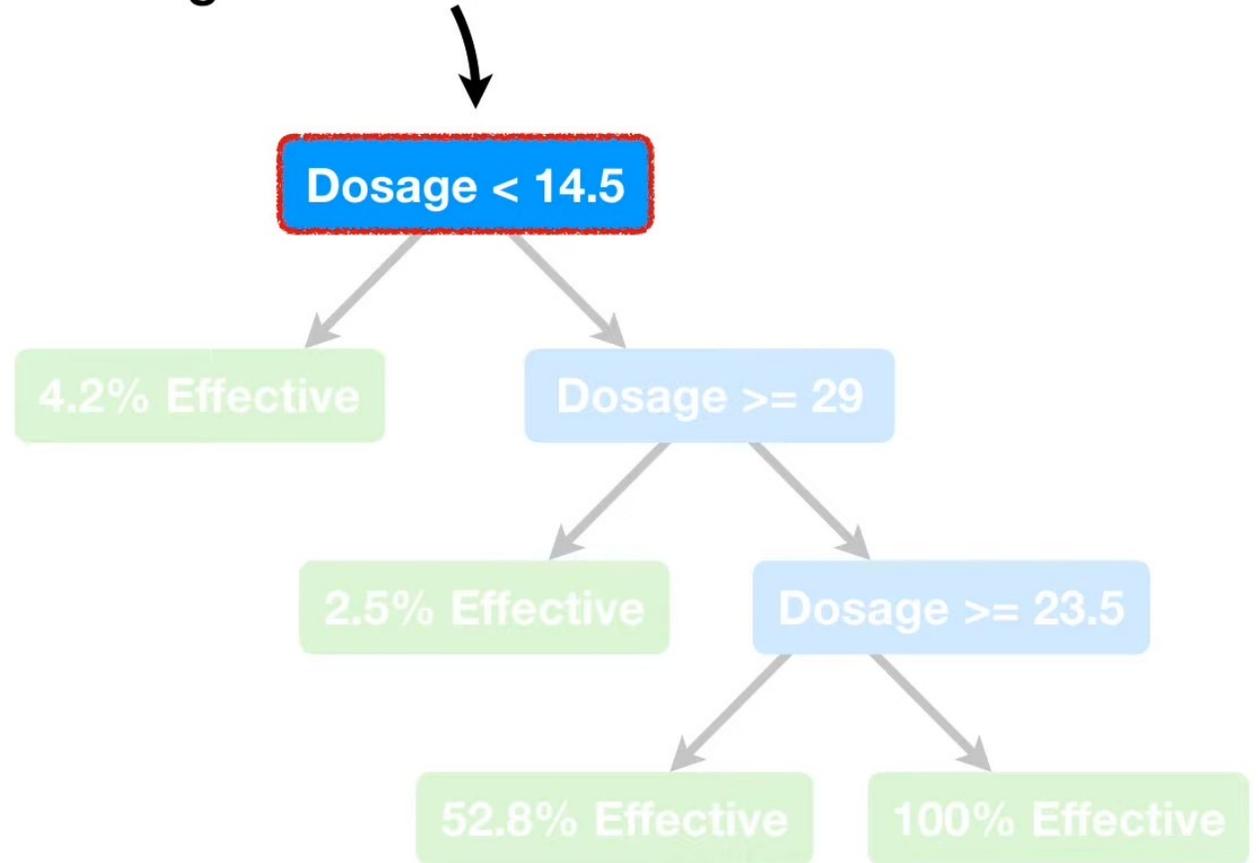2.5% Effective  **Dosage >= 23.5**

52.8% Effective  100% Effective

```
Has Hairy Toes
├── True
└── False
```

…or some other discrete category.

```
Is Really Hungry
├── Eat A Meal
└── Only Sort of Hungry
    ├── Eat A Snack
    └── Don't Eat
```

```
Dosage < 14.5
├── 4.2% Effective
└── Dosage >= 29
    ├── 2.5% Effective
    └── Dosage >= 23.5
        ├── 52.8% Effective
        └── 100% Effective
```
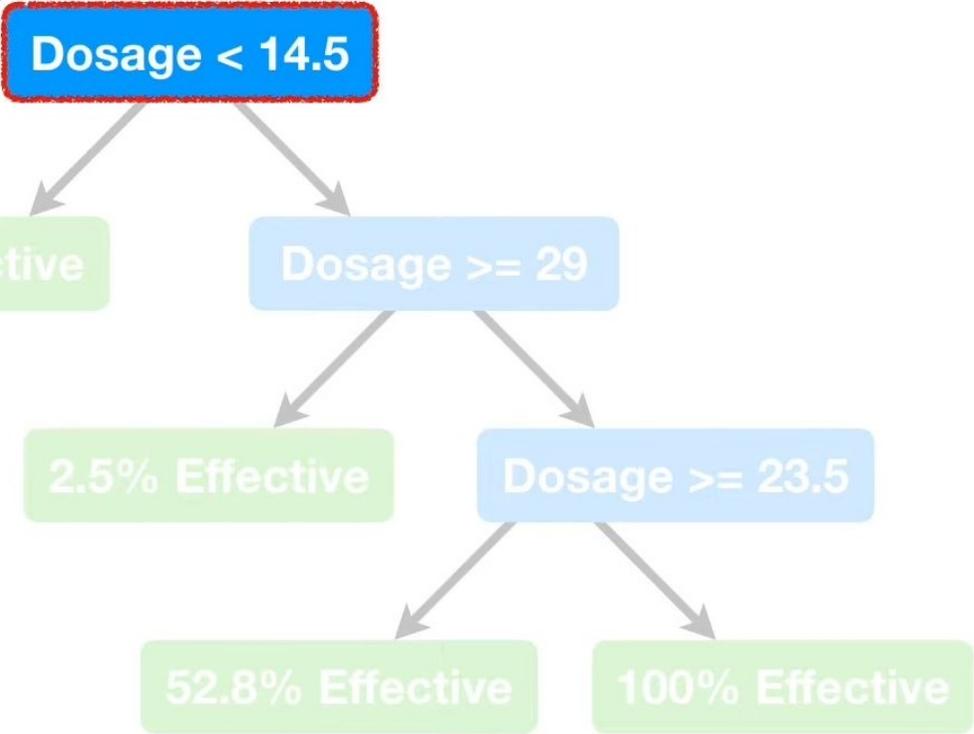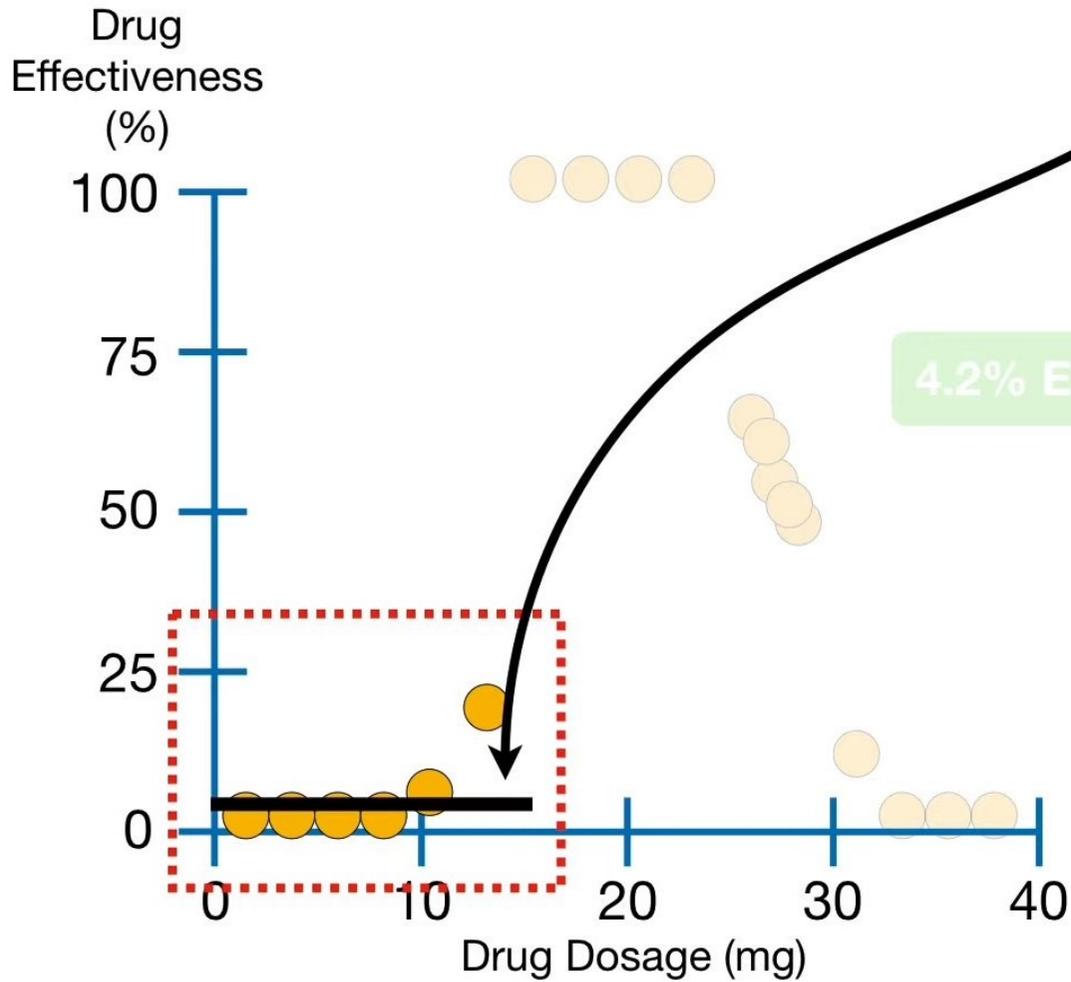
With *this* **Regression Tree**, we start by asking if the **Dosage** is less than **14.5**.

…and the average **Drug Effectiveness** for these **6** observations is **4.2%**…

...so the tree uses the average value, **4.2%**, as its prediction for people with **Dosages < 14.5**.

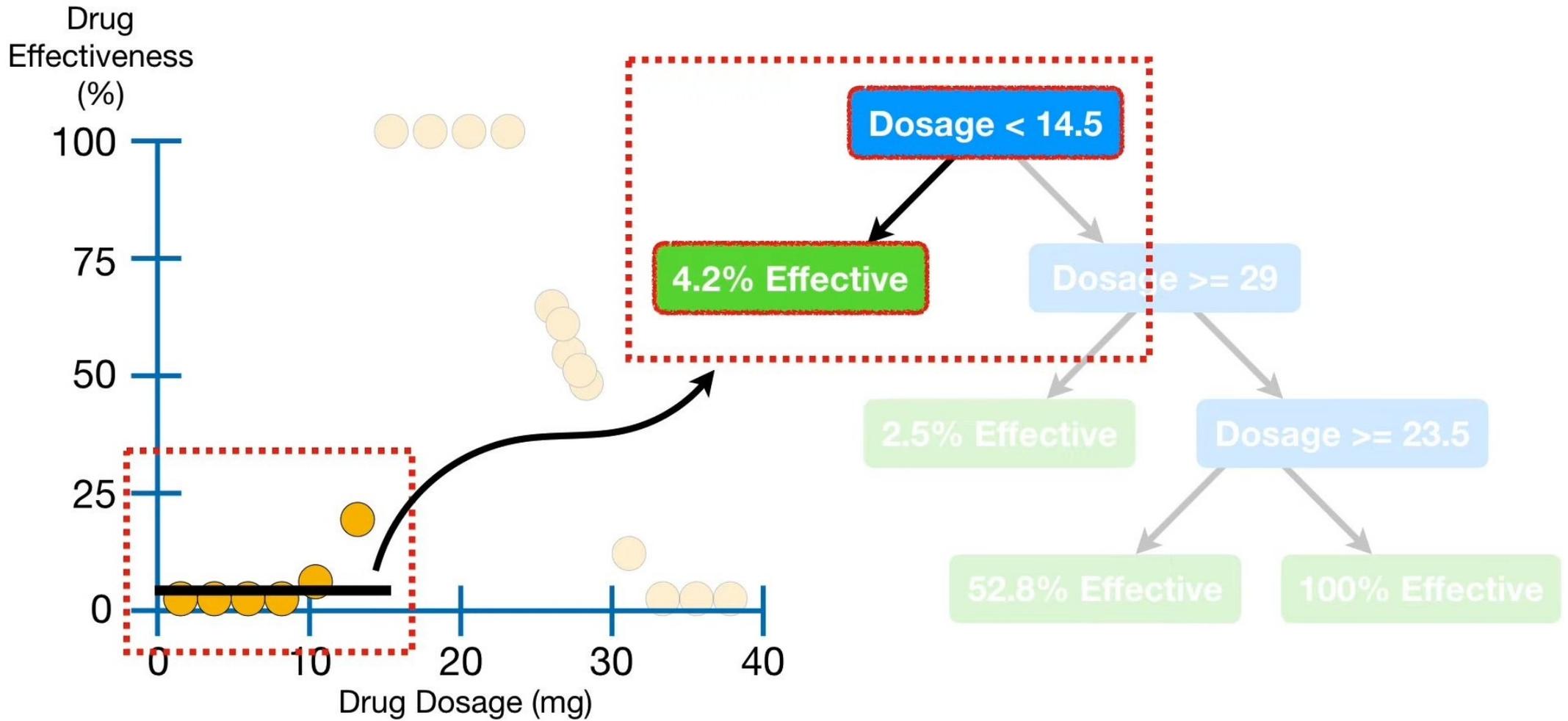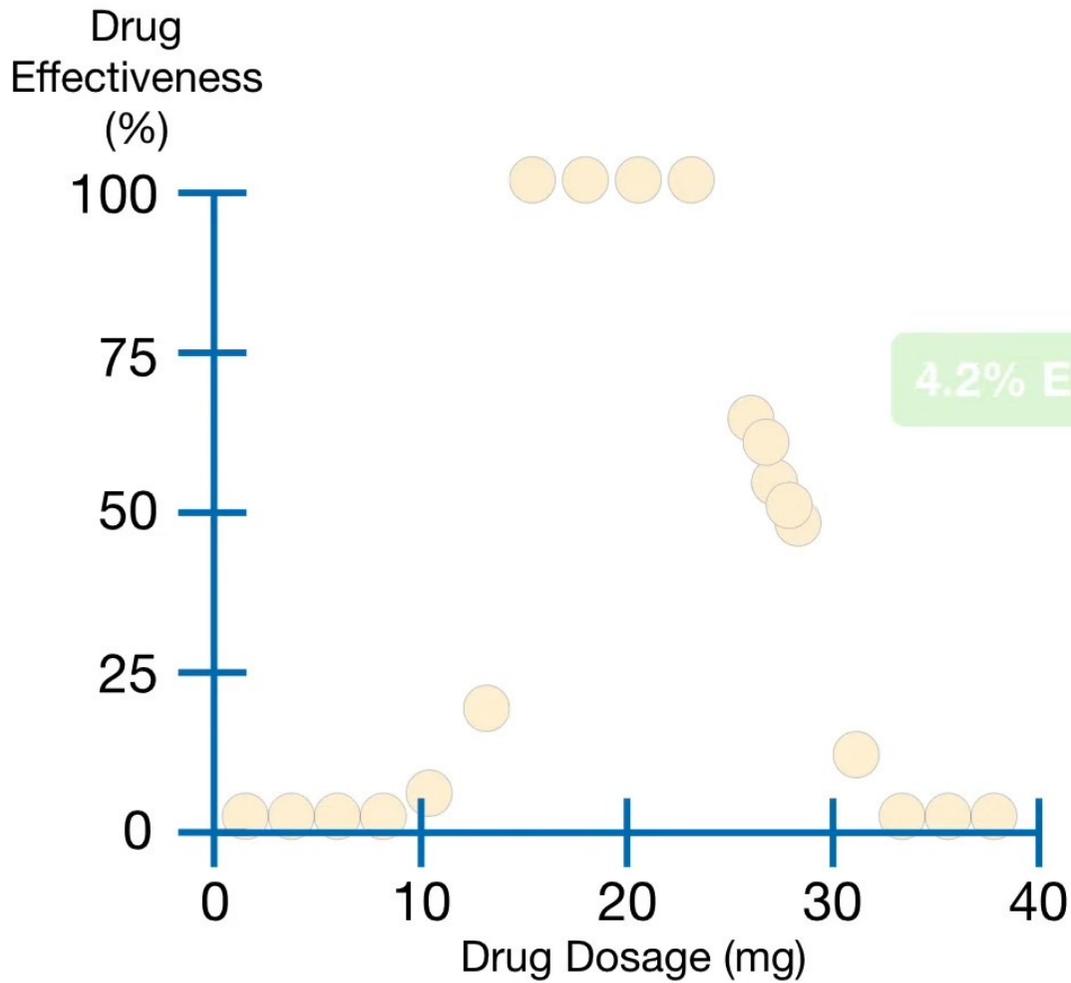…so the tree uses the average value, **2.5%**, as its prediction for people with **Dosages >= 29**.

...then we are talking about these **5** observations in the training dataset...

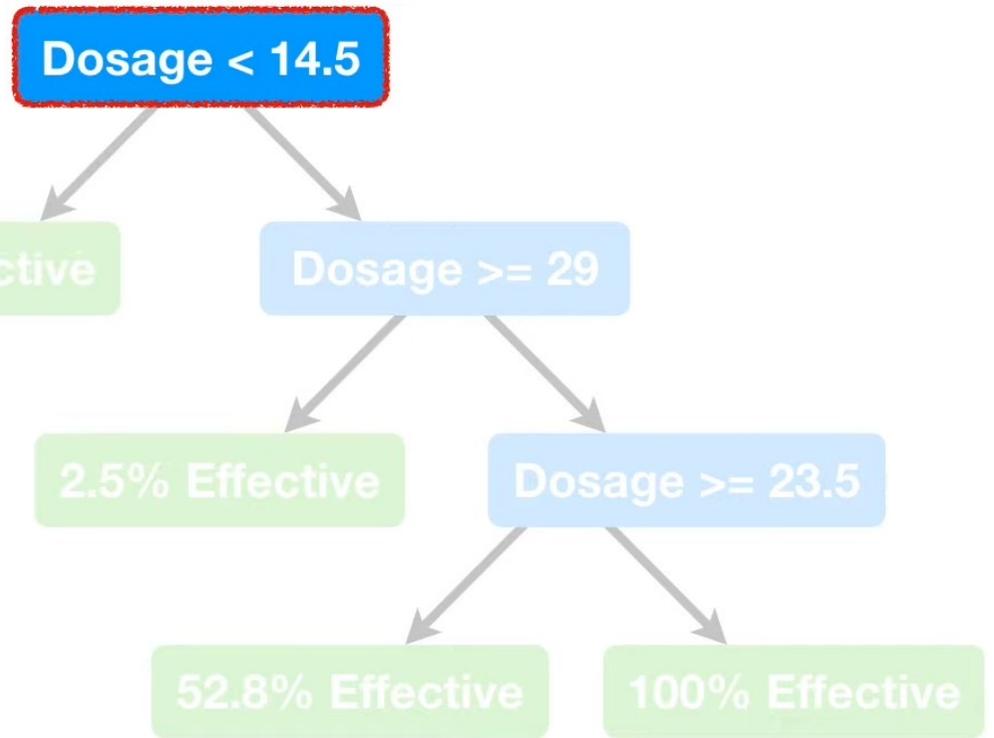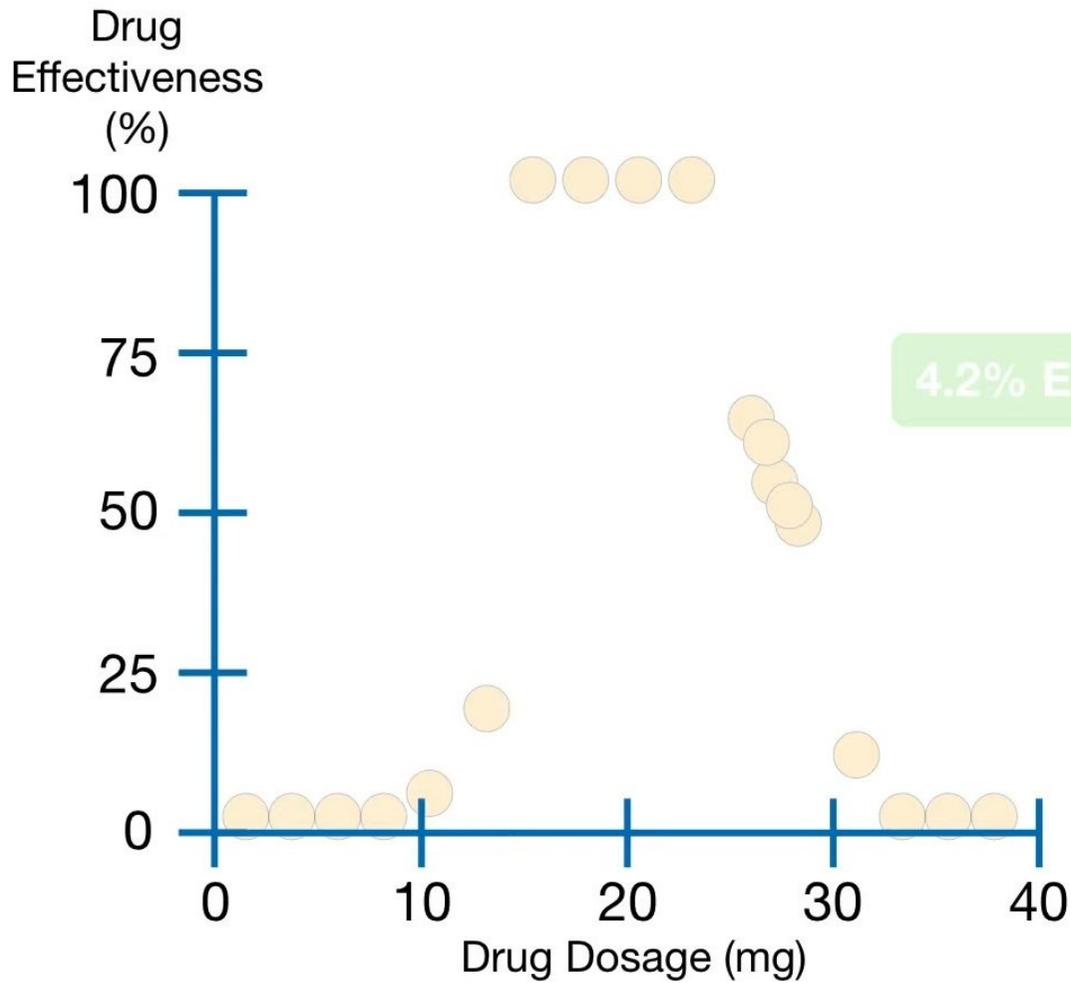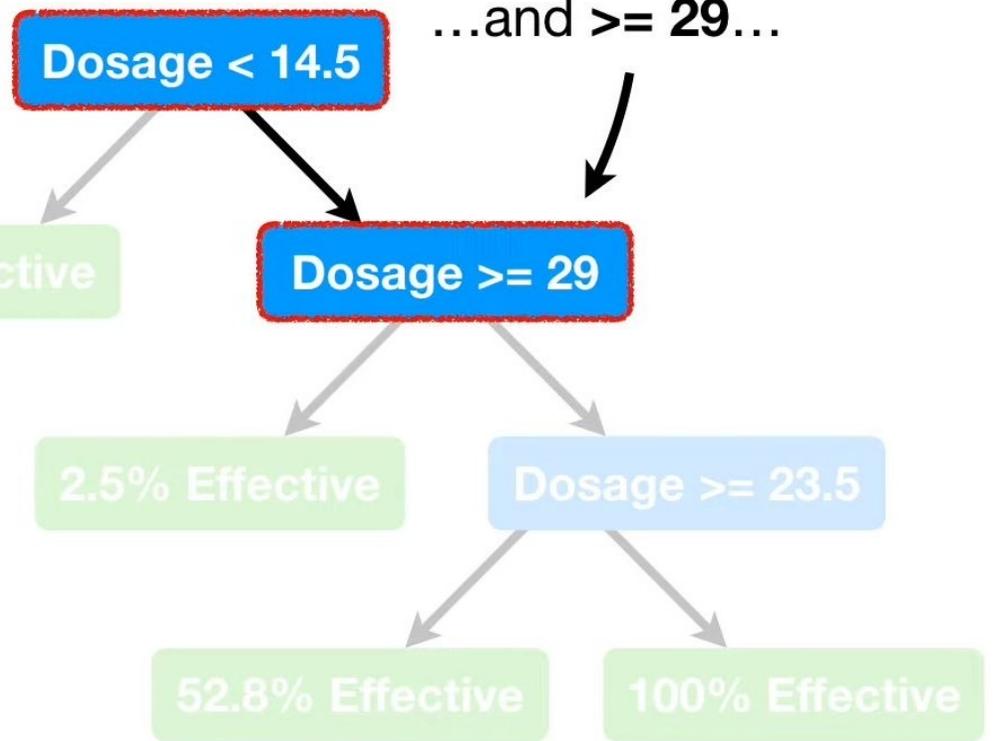…and the average **Drug Effectiveness** for these **5** observations is **52.8%**…

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 14.5

4.2% Effective

Dosage >= 29

2.5% Effective

Dosage >= 23.5

52.8% Effective

100% Effective

...so the tree uses the average value, **52.8%**, as its prediction for people with **Dosages** between **23.5** and **29**.
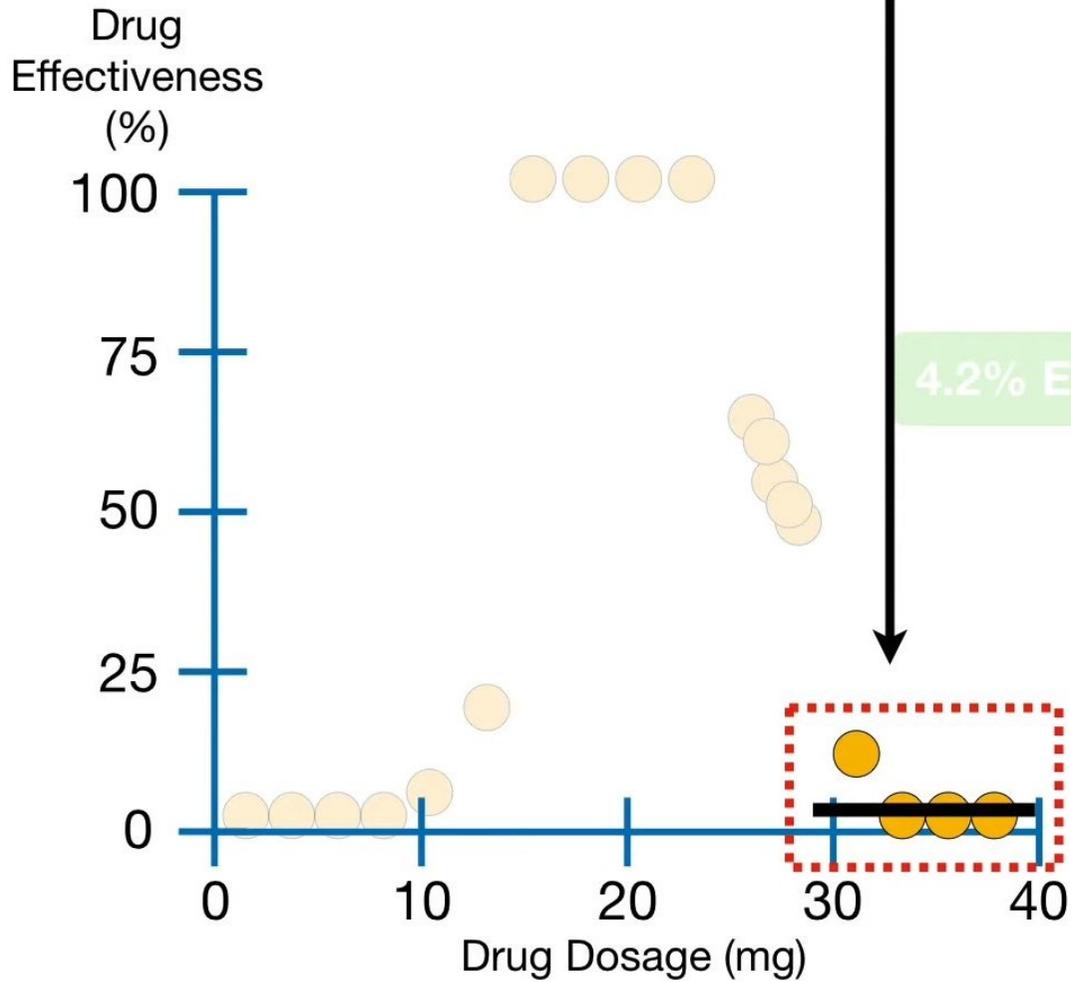
Lastly, if the **Dosage >= 14.5**…

Since each leaf corresponds to the average **Drug Effectiveness** in a different cluster of observations…

...the tree does a better job reflecting the data than the straight line.

When the data are super simple and we are only using one predictor, **Dosage**, to predict **Drug Effectiveness**, making predictions by eye isn't terrible.

Drug Effectiveness (%)

Drug Dosage (mg)

| Dosage | Drug Effect. |
|--------|--------------|
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc… | etc… |

But when we have **3** or more predictors, like **Dosage**, **Age** and **Sex**, to predict **Drug Effectiveness**, drawing a graph is very difficult, if not impossible.

| Dosage | Age | Sex | Etc. | Drug Effect. |
|--------|------|--------|------|--------------|
| 10 | 25 | Female | … | 98 |
| 20 | 73 | Male | … | 0 |
| 35 | 54 | Female | … | 100 |
| 5 | 12 | Male | … | 44 |
| etc… | etc… | etc… | etc… | etc… |

In contrast, a **Regression Tree** easily accommodates the additional predictors.



| Dosage | Age | Sex | Etc. | Drug Effect. |
|--------|-----|-----|------|--------------|
| 10 | 25 | Female | … | 98 |
| 20 | 73 | Male | … | 0 |
| 35 | 54 | Female | … | 100 |
| 5 | 12 | Male | … | 44 |
| etc… | etc… | etc… | etc… | etc… |

For example, if we wanted to predict the **Drug Effectiveness** for this patient…



| | | | | |
|---|---|---|---|---|
| **Age > 50** | | | | |

| | |
|---|---|
| **4% Effective** | **Dosage >= 29** |

| | |
|---|---|
| **20% Effective** | **Sex = Female** |

| | |
|---|---|
| **100% Effective** | **50% Effective** |

| Dosage | Age | Sex | Etc. | Drug Effect. |
|---|---|---|---|---|
| 10 | 25 | Female | … | 98 |
| 20 | 73 | Male | … | 0 |
| 35 | 54 | Female | … | 100 |
| 5 | 12 | Male | … | 44 |
| etc… | etc… | etc… | etc… | etc… |

...we start by asking if they are older than **50**...



| | Age > 50 | |
|---|---|---|
| 4% Effective | Dosage >= 29 | |
| | 20% Effective | Sex = Female |
| | 100% Effective | 50% Effective |

| Dosage | Age | Sex | Etc. | Drug Effect. |
|---|---|---|---|---|
| 10 | 25 | Female | ... | 98 |
| 20 | 73 | Male | ... | 0 |
| 35 | 54 | Female | ... | 100 |
| 5 | 12 | Male | ... | 44 |
| etc... | etc... | etc... | etc... | etc... |

…and since they *not* over **50**, we follow the branch
on the *right* and ask if their **Dosage >= 29**…



| Dosage | Age | Sex | Etc. | Drug Effect. |
|--------|-----|--------|------|--------------|
| 10 | 25 | Female | … | 98 |
| 20 | 73 | Male | … | 0 |
| 35 | 54 | Female | … | 100 |
| 5 | 12 | Male | … | 44 |
| etc… | etc… | etc… | etc… | etc… |

…and since their dosage is *not* **>= 29**, we follow the branch on the *right* and ask if they are **Female**…

Age > 50

4% Effective

Dosage >= 29

20% Effective

Sex = Female

100% Effective

50% Effective

| Dosage | Age | Sex | Etc. | Drug Effect. |
|--------|-----|-----|------|--------------|
| 10 | 25 | Female | … | 98 |
| 20 | 73 | Male | … | 0 |
| 35 | 54 | Female | … | 100 |
| 5 | 12 | Male | … | 44 |
| etc… | etc… | etc… | etc… | etc… |

…and since they *are* **Female**, we follow the branch on the *left* and predict that the dosage will be **100% Effective**…

…and that's not too far off from the truth, **98%**.
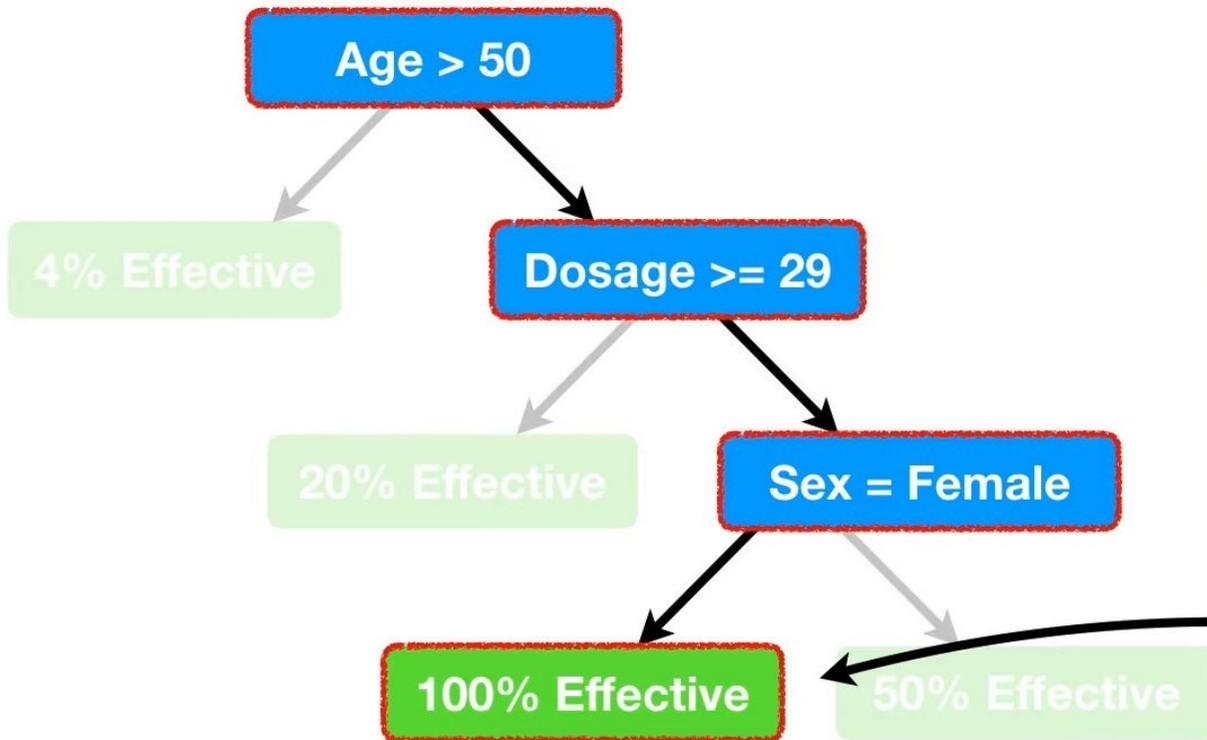


Decision tree:
- **Age > 50**
  - → 4% Effective
  - → **Dosage >= 29**
    - → 20% Effective
    - → **Sex = Female**
      - → **100% Effective**
      - → 50% Effective

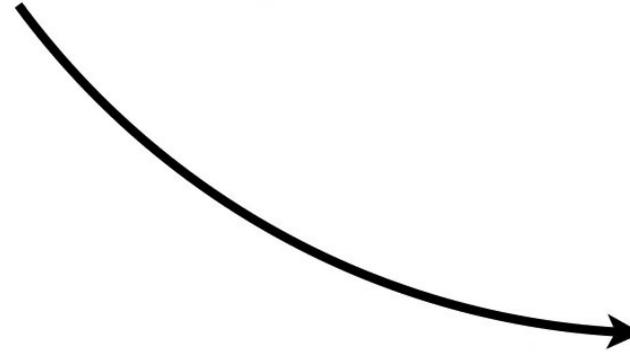| Dosage | Age | Sex | Etc. | Drug Effect. |
|--------|-----|--------|------|--------------|
| 10 | 25 | Female | … | 98 |
| 20 | 73 | Male | … | |
| 35 | 54 | Female | … | 100 |
| 5 | 12 | Male | … | 44 |
| etc… | etc… | etc… | etc… | etc… |

OK, now that we know that **Regression Trees**
can easily handle complicated data...



| Dosage | Age | Sex | Etc. | Drug Effect. |
|--------|-----|--------|------|--------------|
| 10 | 25 | Female | ... | 98 |
| 20 | 73 | Male | ... | |
| 35 | 54 | Female | ... | 100 |
| 5 | 12 | Male | ... | 44 |
| etc... | etc... | etc... | etc... | etc... |

…let's go back to the original data, with
just one predictor, **Dosage**…

| Dosage | Drug Effect. |
|--------|--------------|
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc… | etc… |

...and talk about how to build this
**Regression Tree** from scratch...



| Dosage | Drug Effect. |
| --- | --- |
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc... | etc... |

…and since **Regression Trees** are
built from the top down…

**Dosage < 14.5**

| Dosage | Drug Effect. |
|--------|--------------|
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc… | etc… |

...the first thing we do is figure out why
we start by asking if **Dosage < 14.5**.

Dosage < 14.5

4.2% Effective

Dosage >= 29

| Dosage | Drug Effect. |
|--------|--------------|
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc… | etc… |

Going back to the graph of the data…

| Dosage | Drug Effect. |
|--------|--------------|
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc… | etc… |

Their average **Dosage** is **3**, and that corresponds to this dotted **red line**.

Drug Effectiveness (%)

Drug Dosage (mg)

Now we can build a very simple tree that splits the observations into two groups based whether or not **Dosage < 3**.

Dosage < 3

The point on the far left is the only one with **Dosage < 3**…

Drug
Effectiveness
(%)

100

75

50

25

0

0    10    20    30    40

Drug Dosage (mg)

Dosage < 3

…so we put **0** in the leaf on the left side, for when **Dosage < 3**.

...and the average **Drug Effectiveness** for all of the points with **Dosages >= 3** is **38.8**, (the green line)...

The values in each leaf are the predictions that this simple tree will make for **Drug Effectiveness**.

Dosage < 3

Average=0

Average=38.8

The *prediction* for this point, **Drug Effectiveness = 0**, is pretty good since it is the same as the *observed* value.
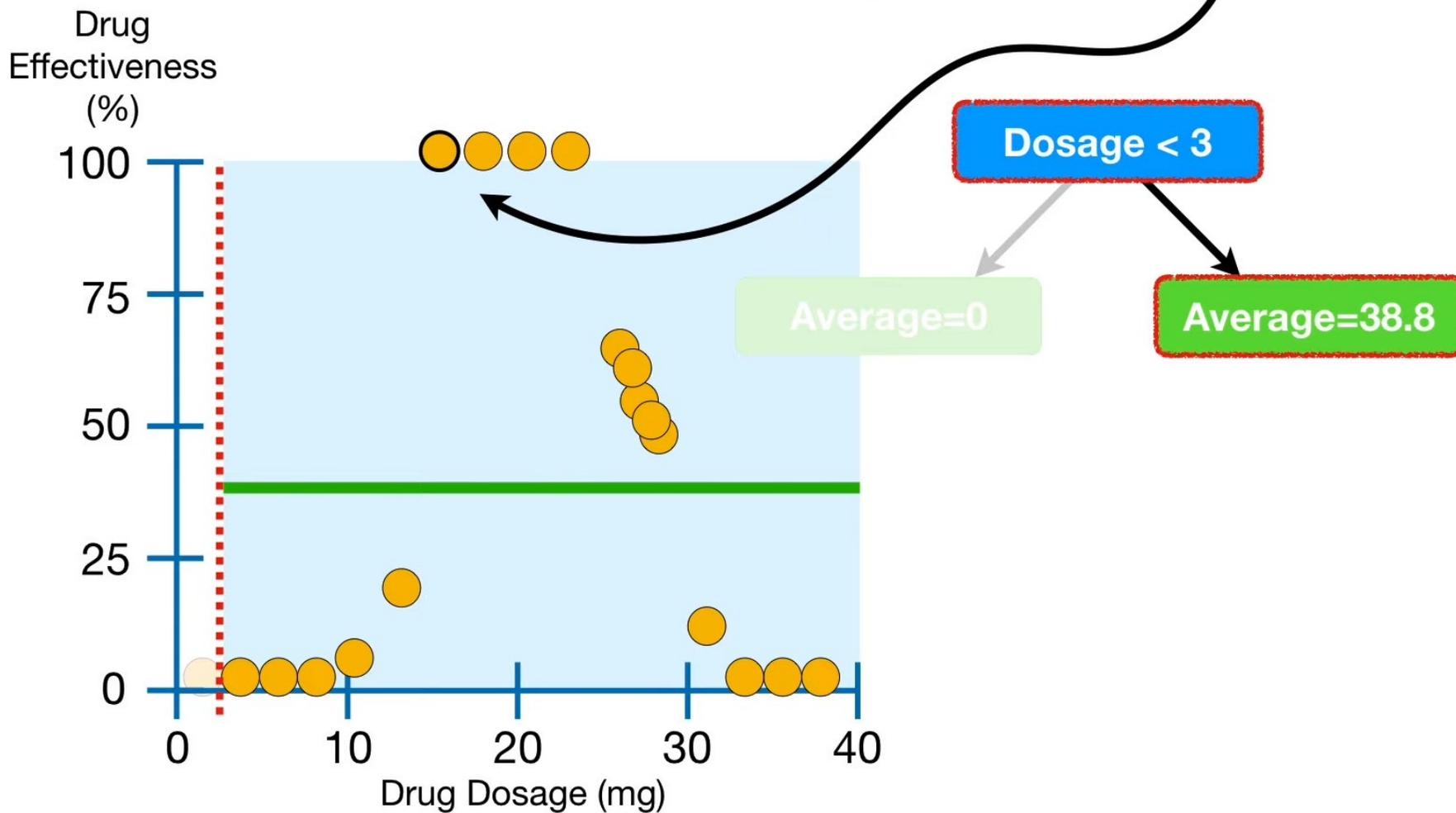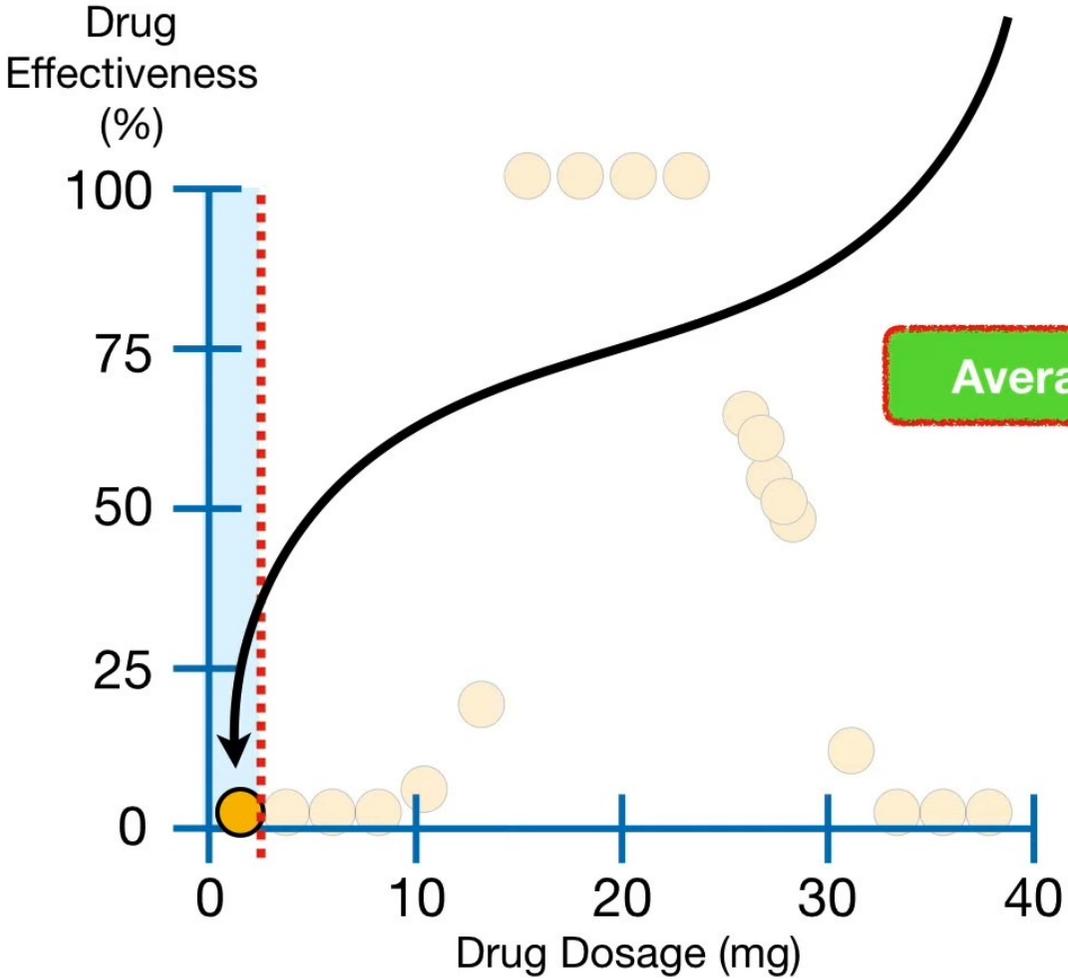
For each point in the data, we can draw its **residual**, the difference between the *observed* and *predicted* values…
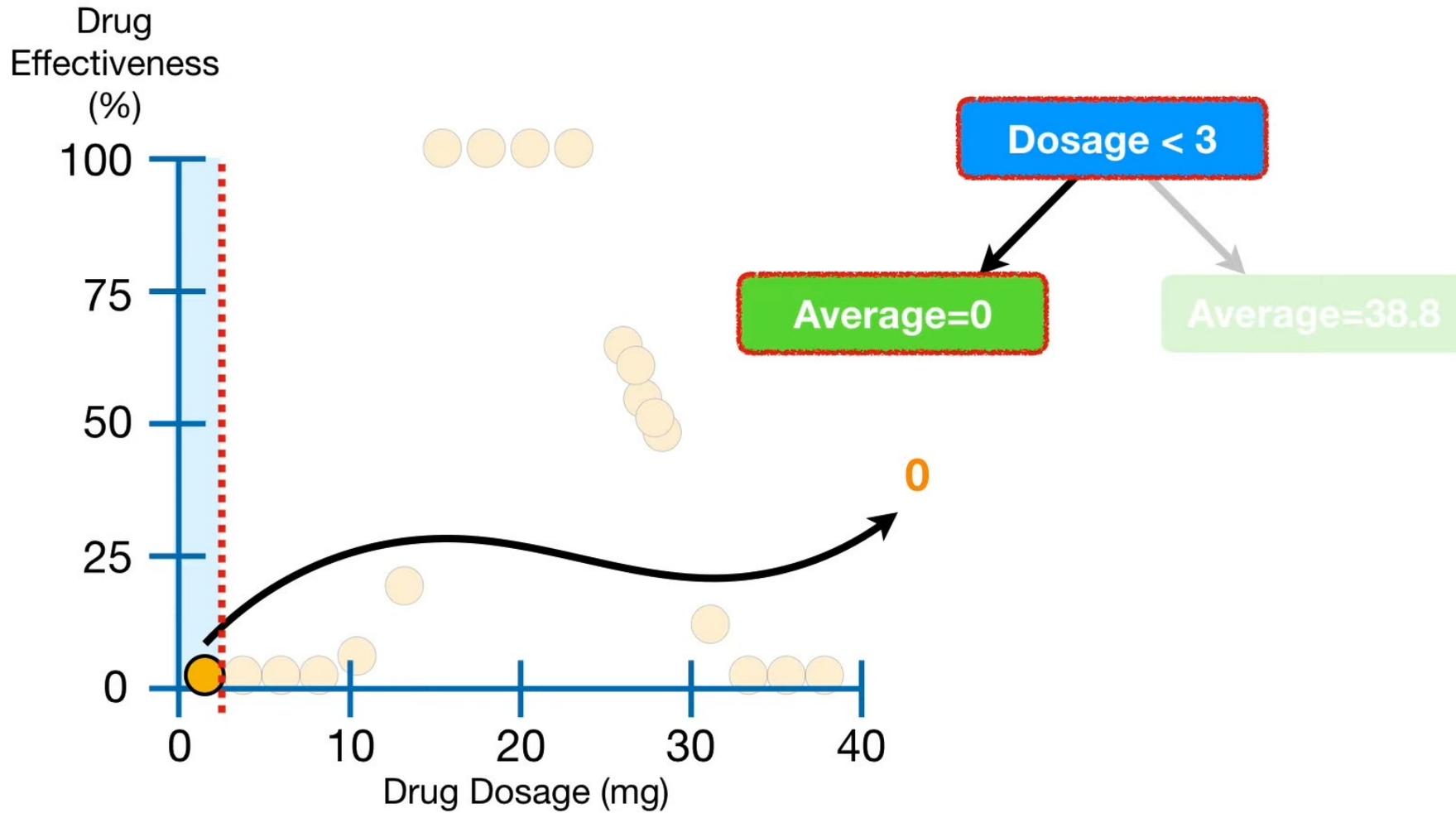
Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 3

Average=0

Average=38.8

…we calculate the difference between its *observed* **Drug Effectiveness**, **0**,…

...and then add it to the first term.

Drug Effectiveness (%)

100

75

50

25

0

0    10    20    30    40

Drug Dosage (mg)

Dosage < 3

Average=0

Average=38.8

$(0 - 0)^2 + (0 - 38.8)^2$

...and the rest of the points...

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 3

Average=0

Average=38.8

$(0 - 0)^2 + (0 - 38.8)^2 + (0 - 38.8)^2 + (0 - 38.8)^2$
$+ (5 - 38.8)^2 + (20 - 38.8)^2$
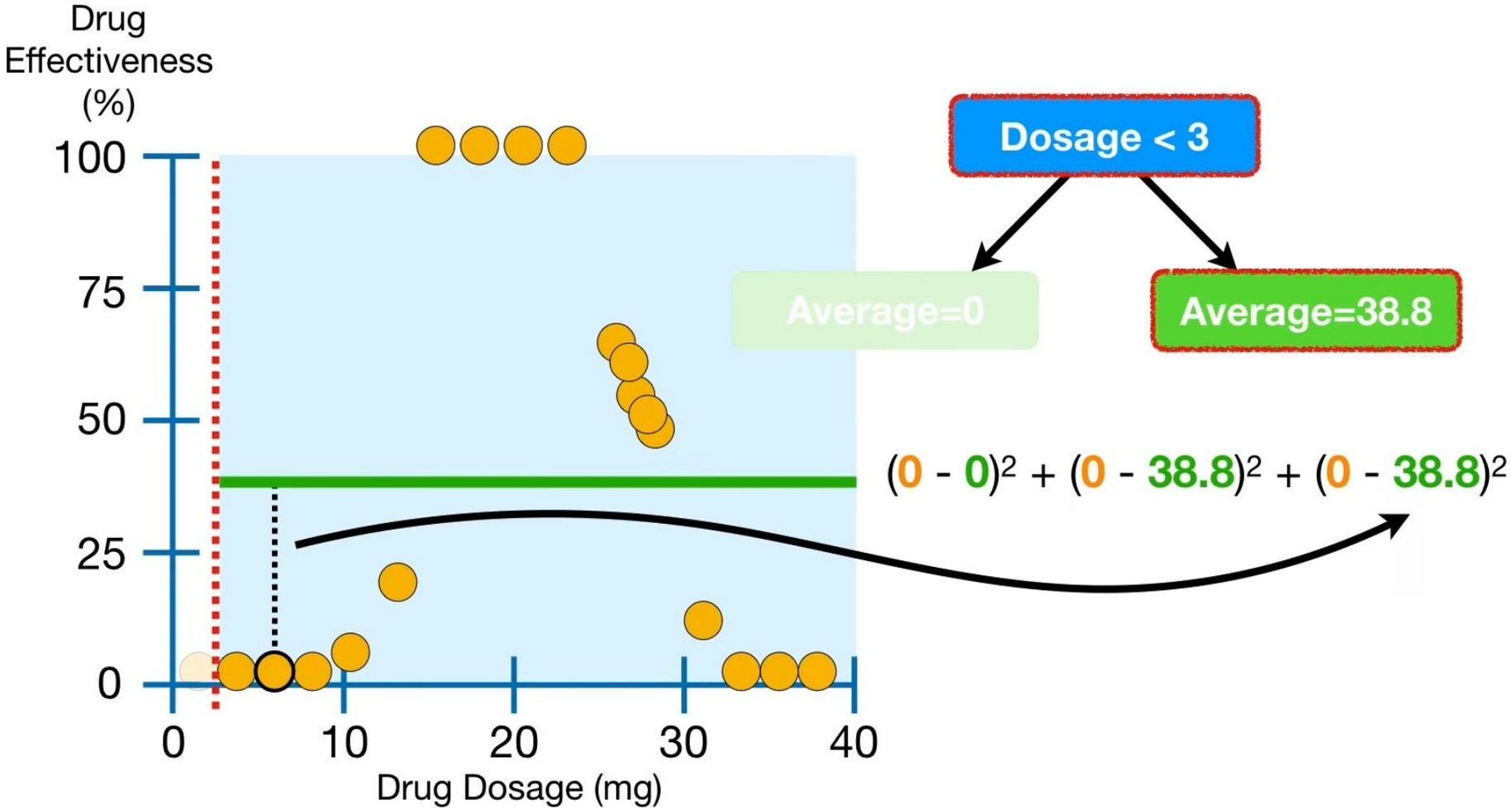
...until we have added squared residuals for every point.

Dosage < 3

Average=0          Average=38.8

$(0 - 0)^2 + (0 - 38.8)^2 + (0 - 38.8)^2 + (0 - 38.8)^2$

$+ (5 - 38.8)^2 + (20 - 38.8)^2 + (100 - 38.8)^2$

$+ (100 - 38.8)^2 + ... + (0 - 38.8)^2$

...we add up the squared residuals for every point...

Drug Effectiveness (%)

Dosage < 3

Average=0    Average=38.8
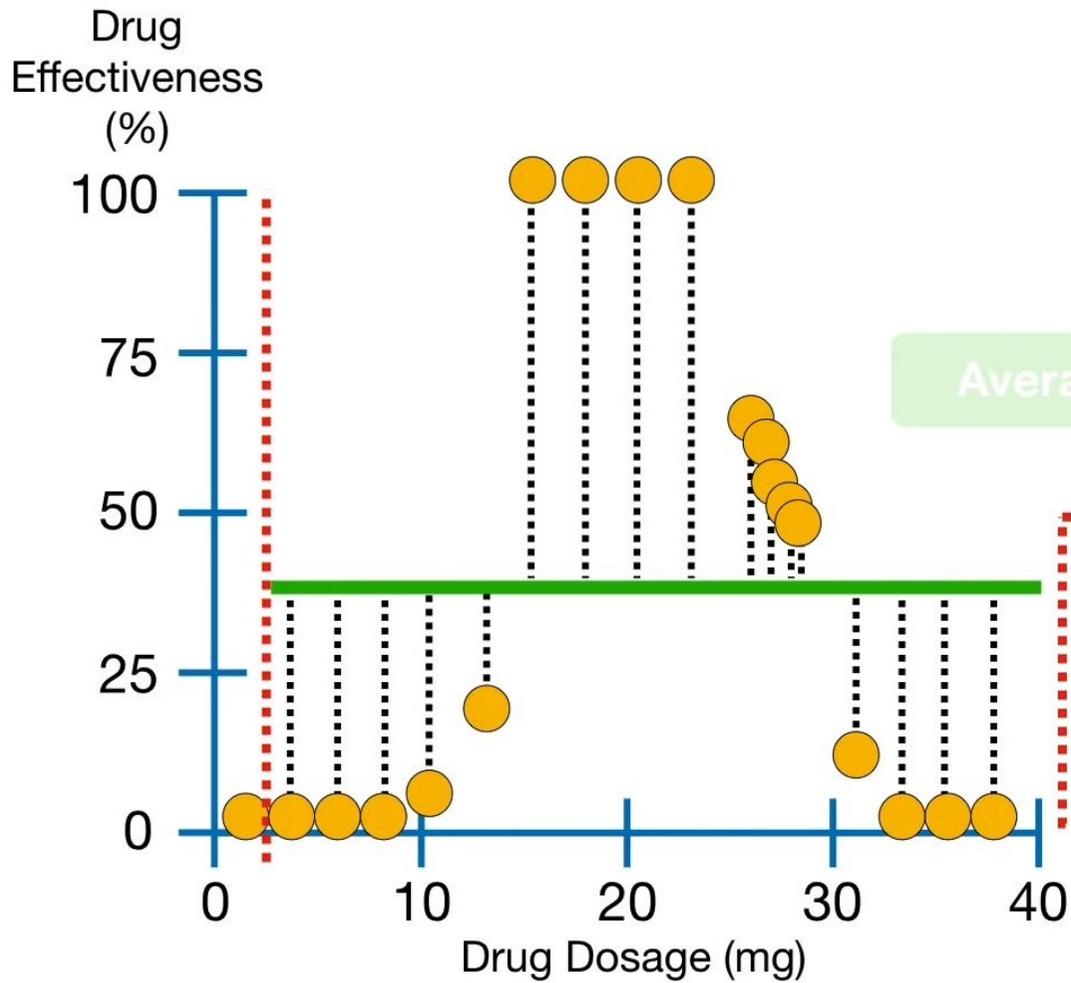
$(0 - 0)^2 + (0 - 38.8)^2 + (0 - 38.8)^2 + (0 - 38.8)^2$

$+ (5 - 38.8)^2 + (20 - 38.8)^2 + (100 - 38.8)^2$

$+ (100 - 38.8)^2 + \ldots + (0 - 38.8)^2$

Drug Dosage (mg)

...and get **27,468.5**.

Drug Effectiveness (%)

Dosage < 3

Average=0    Average=38.8

$(0 - 0)^2 + (0 - 38.8)^2 + (0 - 38.8)^2 + (0 - 38.8)^2$

$+ (5 - 38.8)^2 + (20 - 38.8)^2 + (100 - 38.8)^2$

$+ (100 - 38.8)^2 + \ldots + (0 - 38.8)^2$

$= 27,468.5$

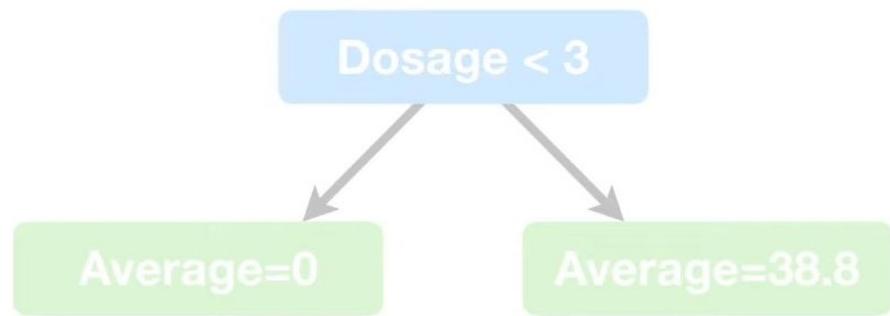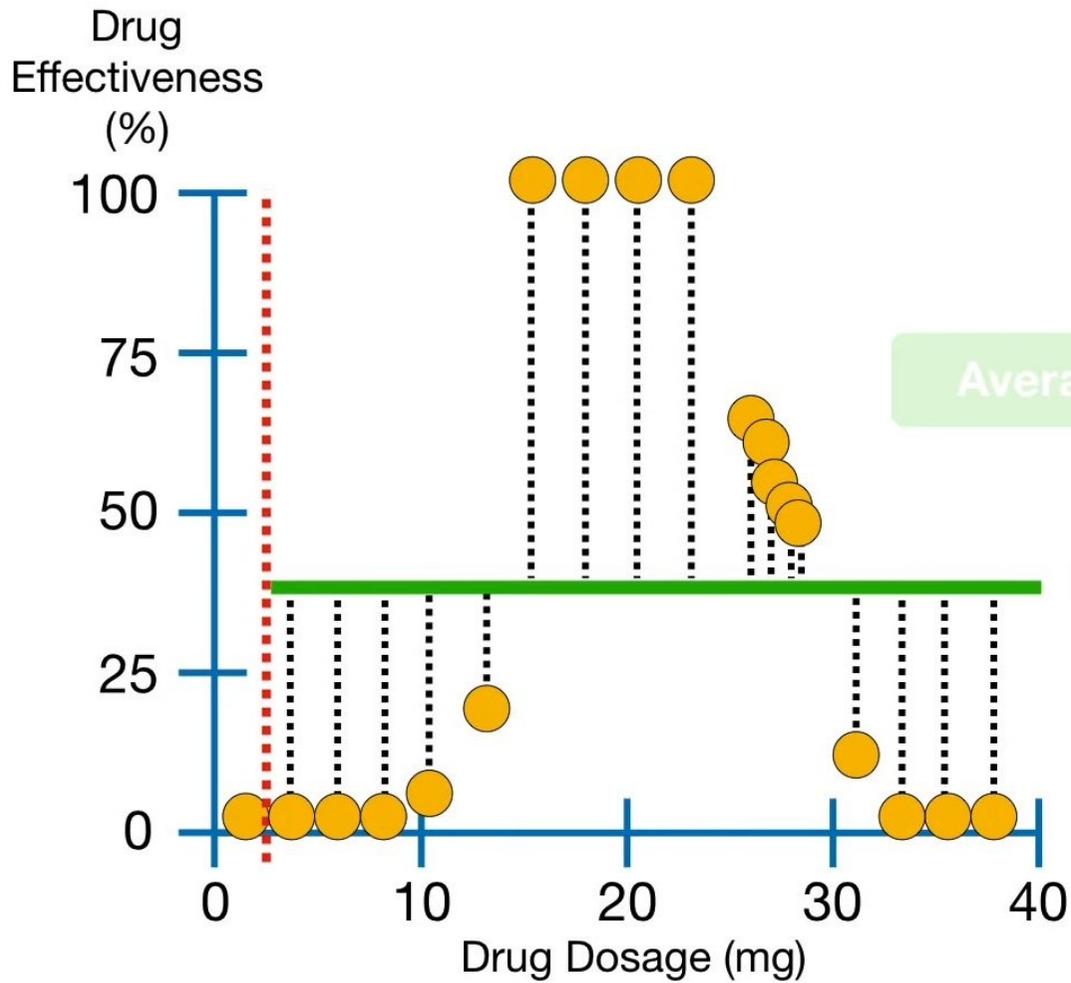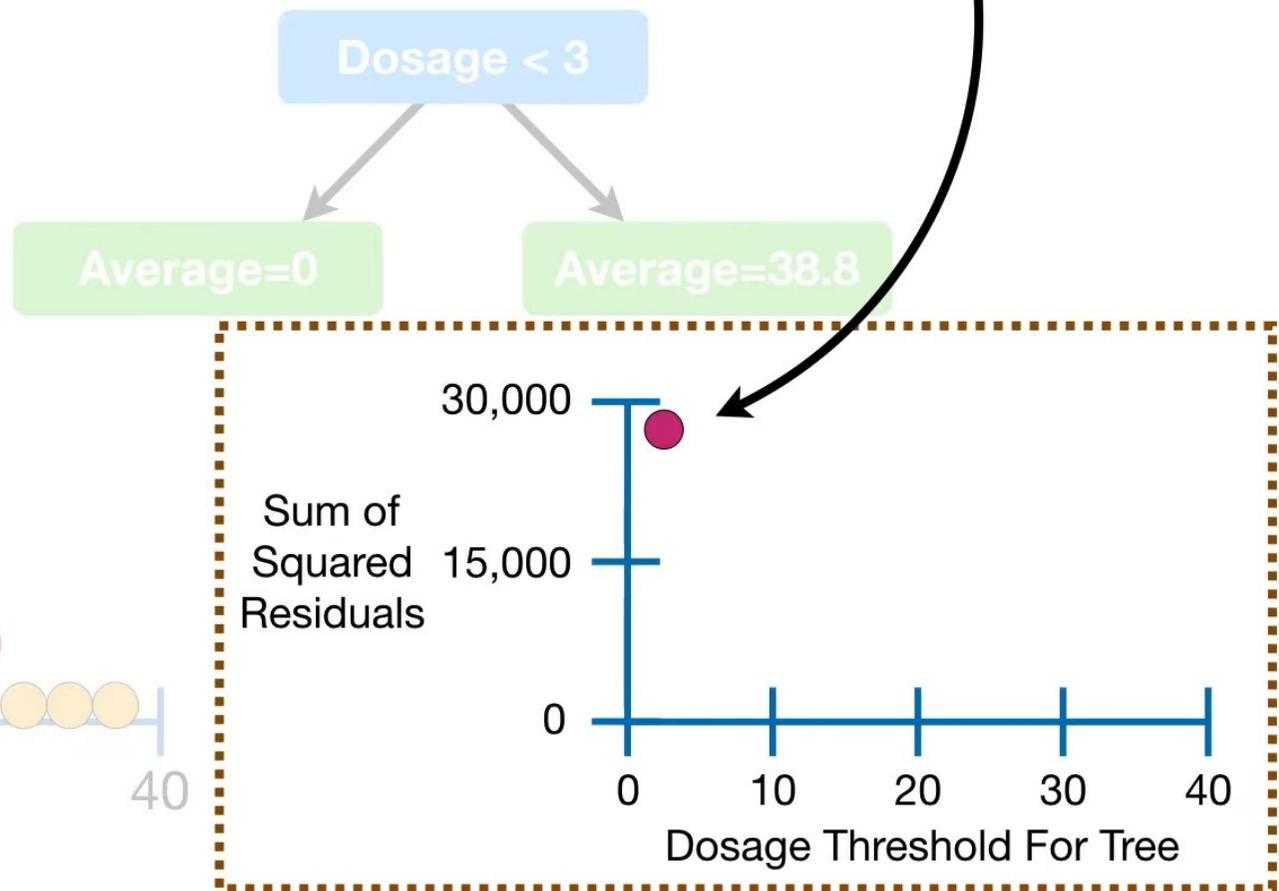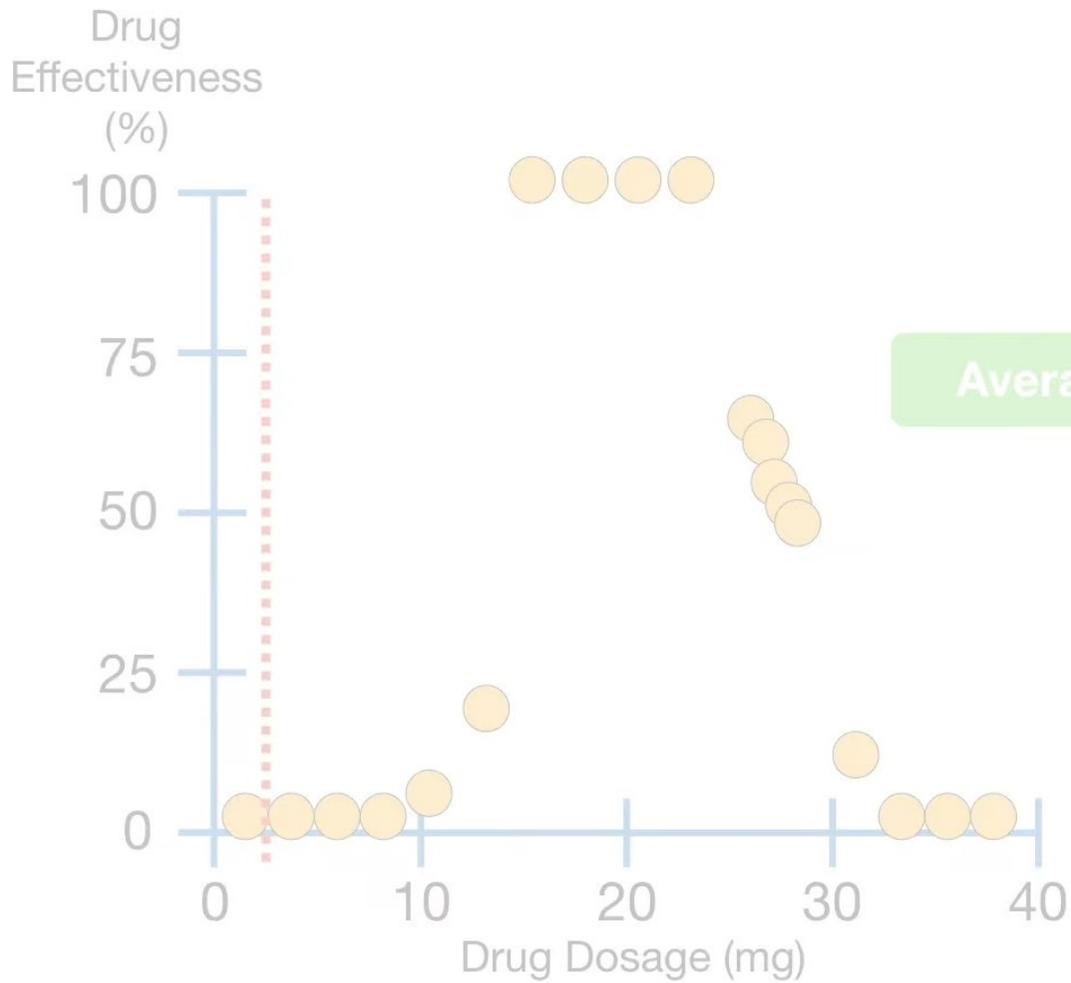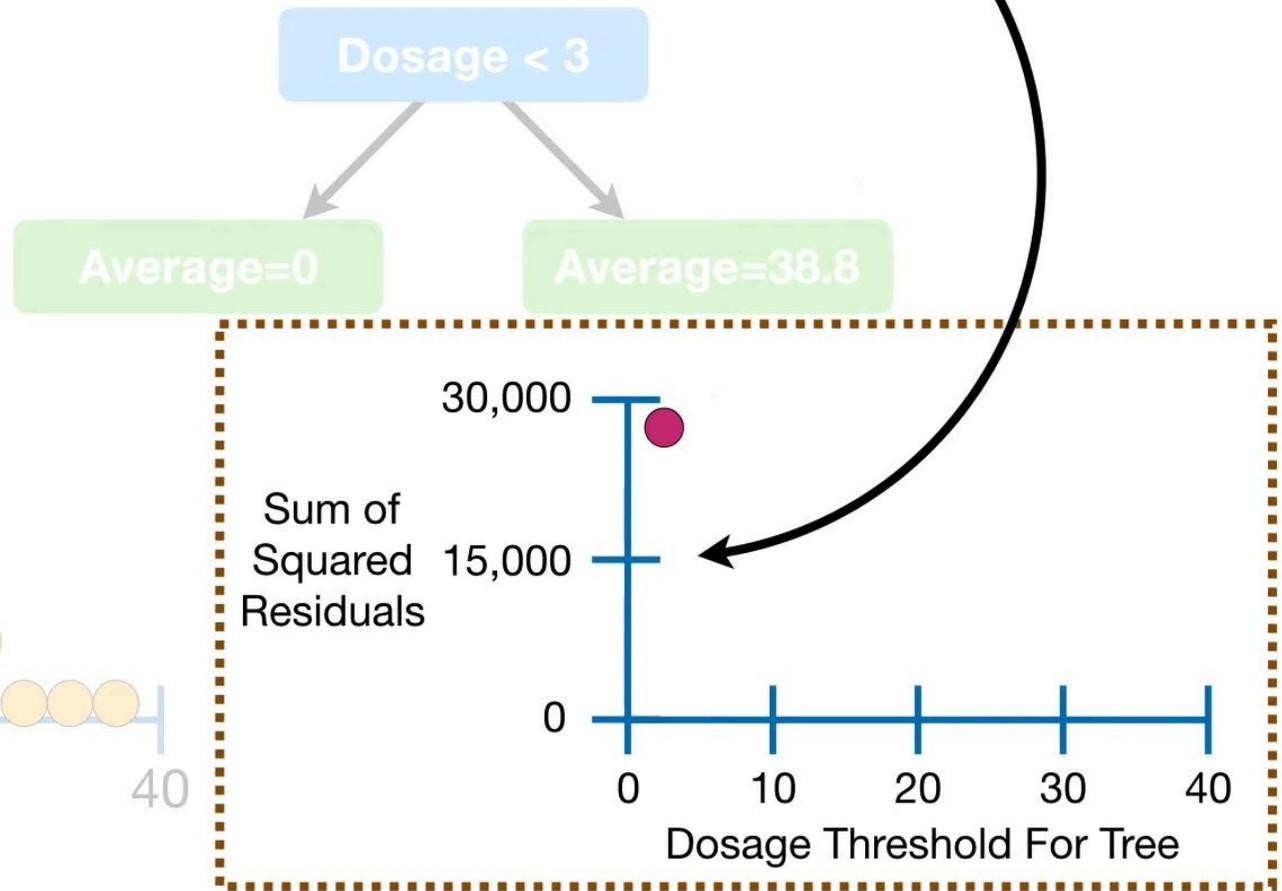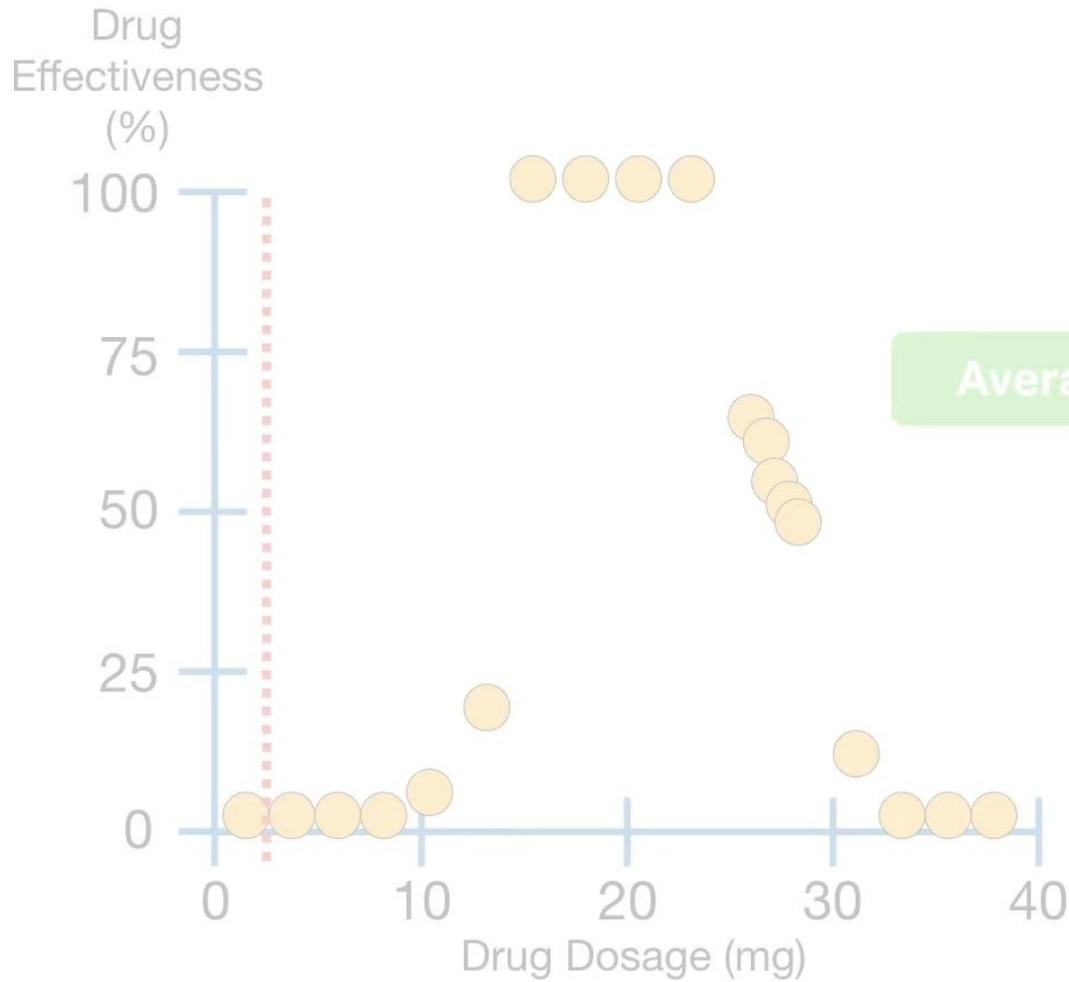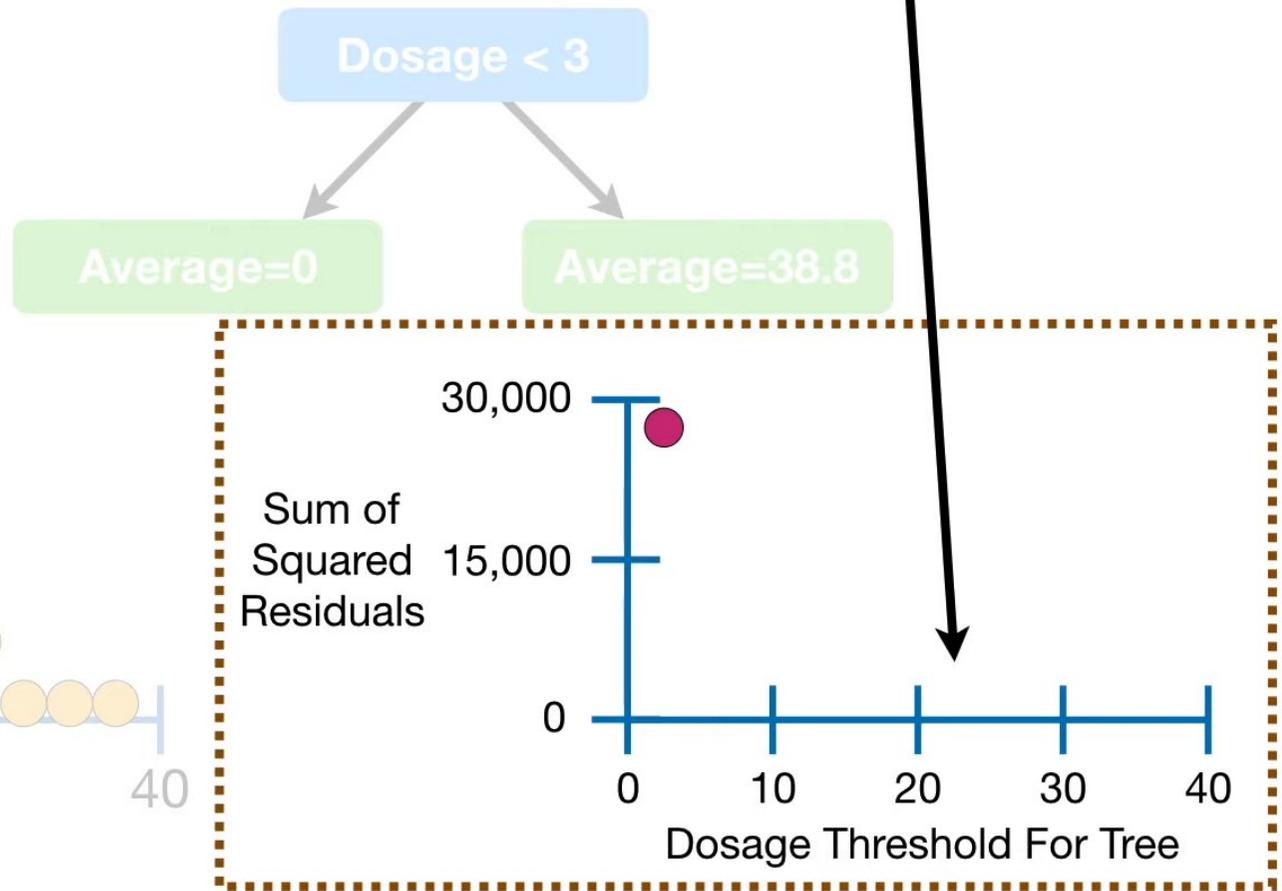Drug Dosage (mg)

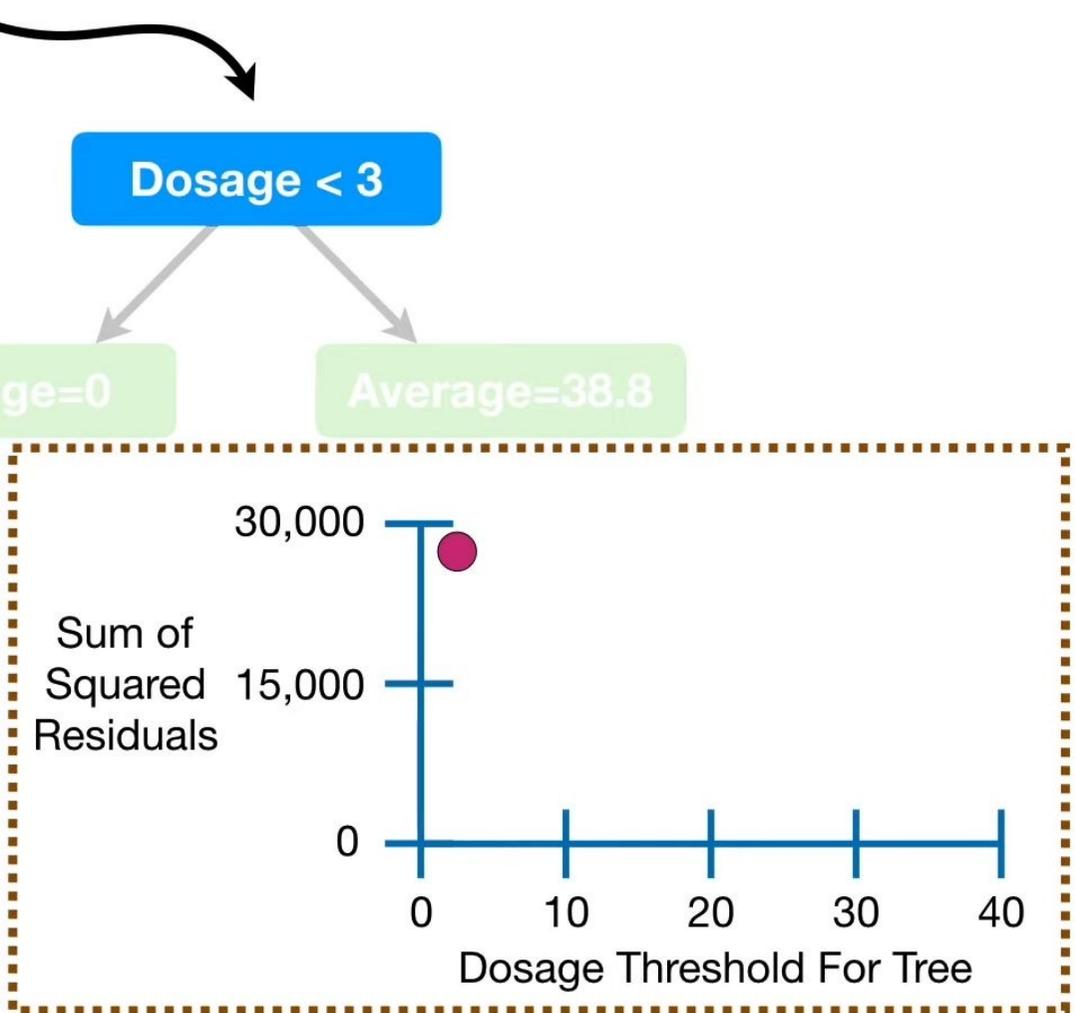NOTE: We can plot the sum of squared residuals on this graph.

Dosage < 3

Average=0          Average=38.8

The **y-axis** corresponds to the sum of squared residuals…

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 3

Average=0

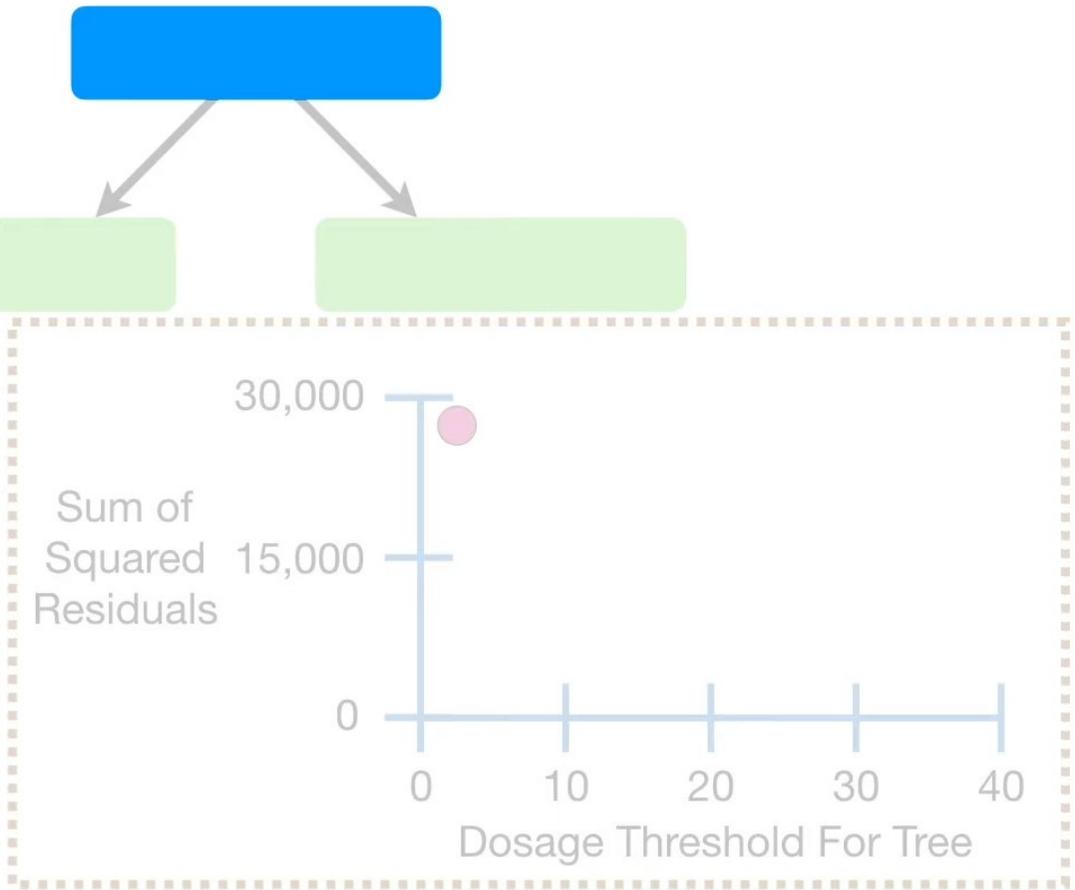Average=38.8

Sum of Squared Residuals

Dosage Threshold For Tree

...and the **x-axis** corresponds to **Dosage** thresholds.

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 3

Average=0

Average=38.8

Sum of Squared Residuals

30,000

15,000

0

Dosage Threshold For Tree

In this case, the **Dosage** threshold was **3**…

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 3

Average=0     Average=38.8

Sum of Squared Residuals

30,000

15,000

0

Dosage Threshold For Tree

...but if we focus on the next two points in the graph...

Drug Effectiveness (%)

100

75

50

25

0

0   10   20   30   40

Drug Dosage (mg)

Sum of Squared Residuals

30,000

15,000

0

0   10   20   30   40

Dosage Threshold For Tree

...and calculate their average **Dosage**, which is **5**...

Drug Effectiveness (%)

Drug Dosage (mg)

Sum of Squared Residuals

Dosage Threshold For Tree

…then we can use **Dosage < 5** as a new threshold.

Using **Dosage < 5** gives us new predictions…
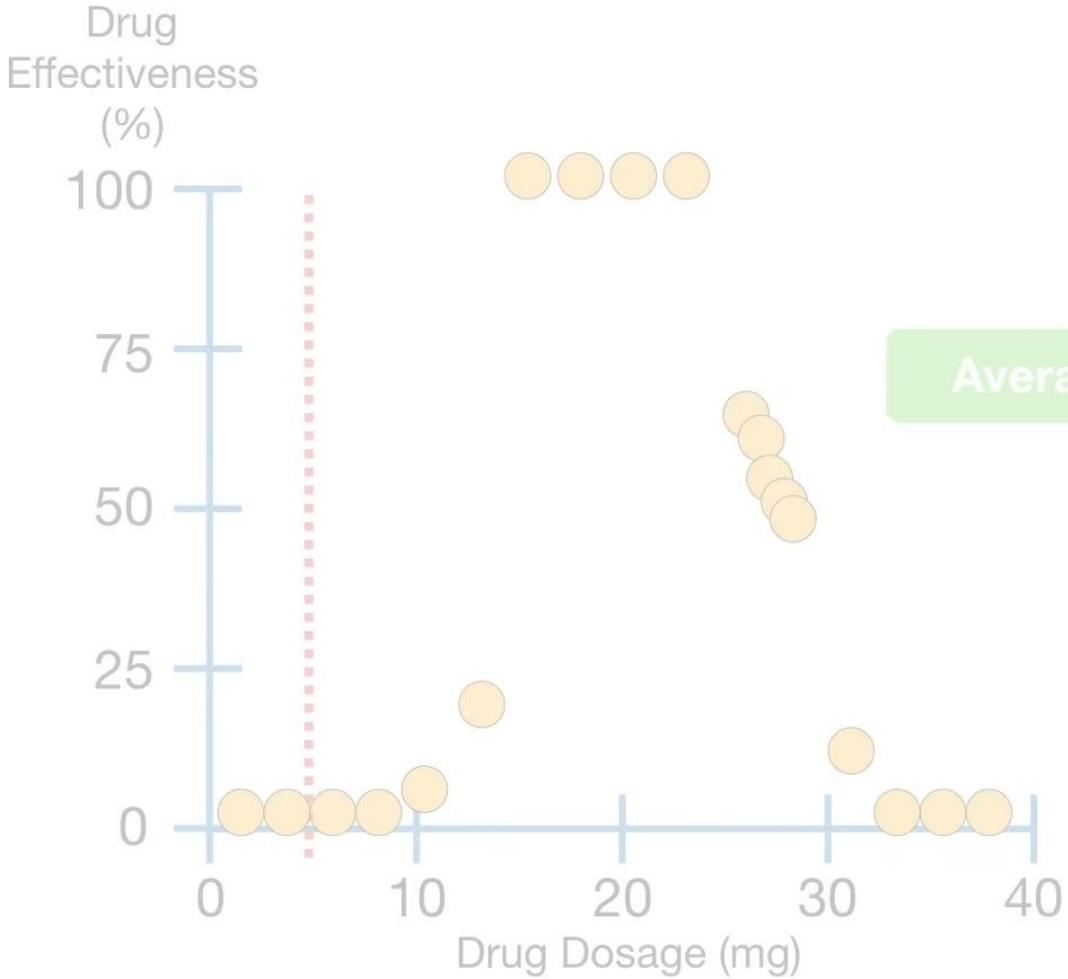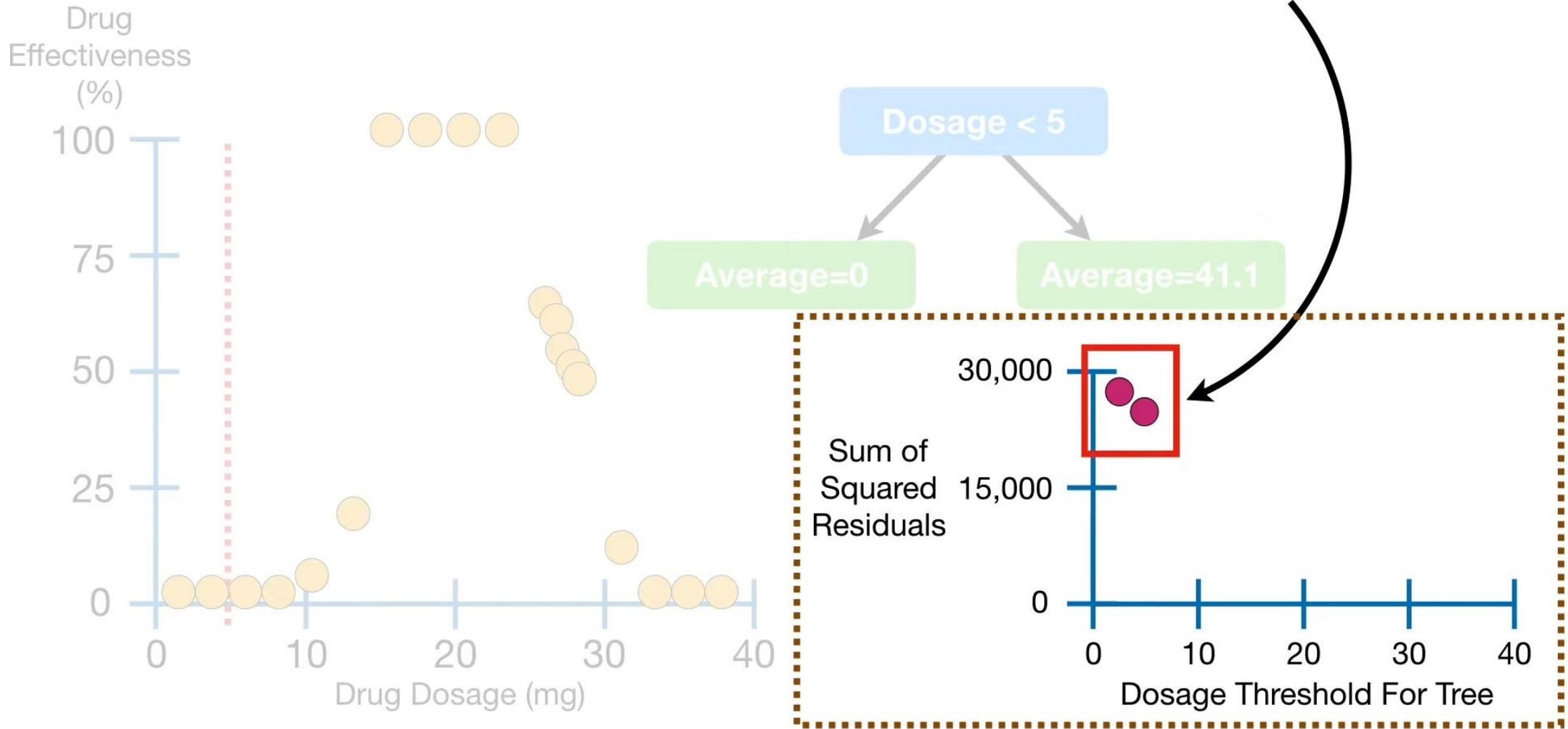
…and that means we can add a new sum of squared residuals to our graph.

Drug Effectiveness (%)

Dosage < 5

Average=0          Average=41.1
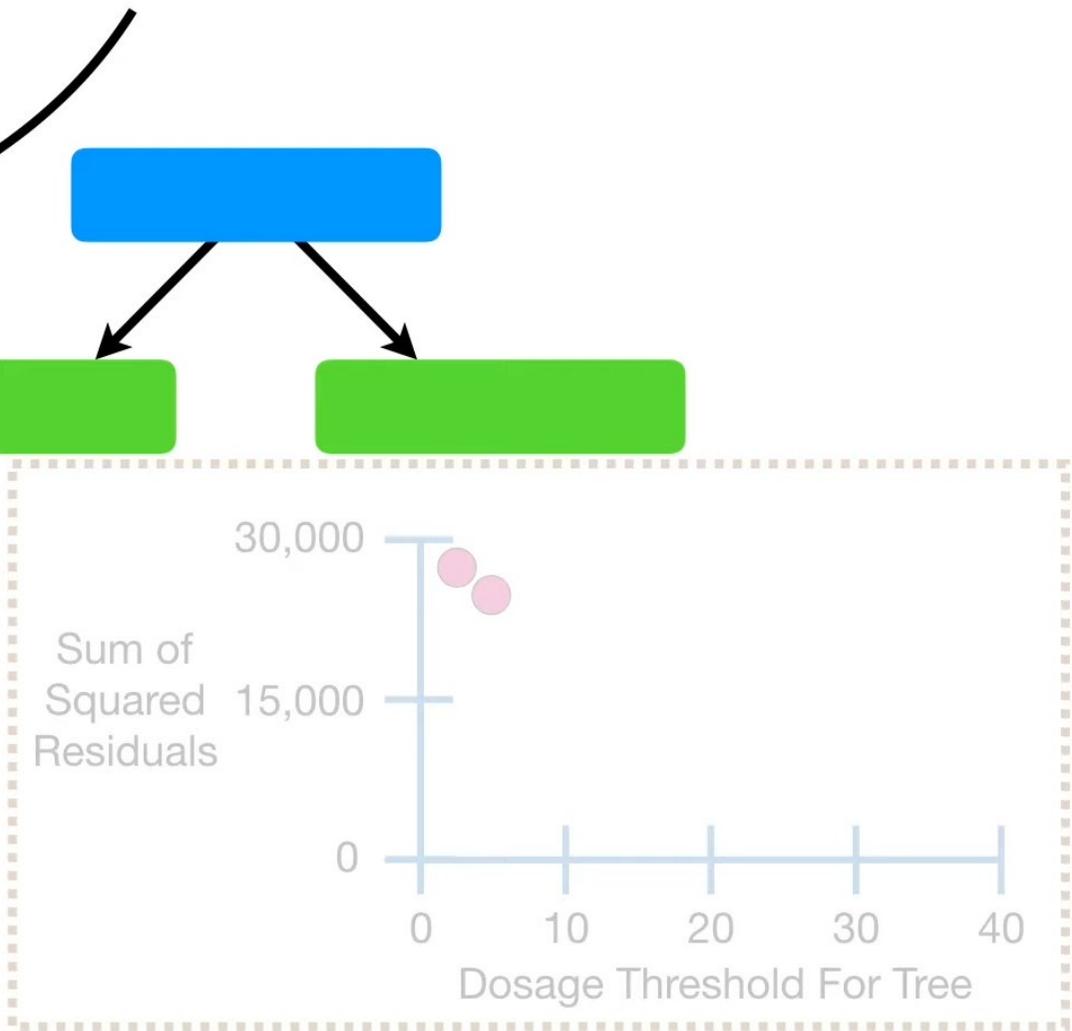
Sum of Squared Residuals

30,000

15,000

0

Dosage Threshold For Tree

Drug Dosage (mg)

In this case, the new threshold, **Dosage < 5**, results in a smaller sum of squared residuals…

...and that means using **Dosage < 5** as the threshold resulted in better predictions over all.

...calculate their average, which is **7**...
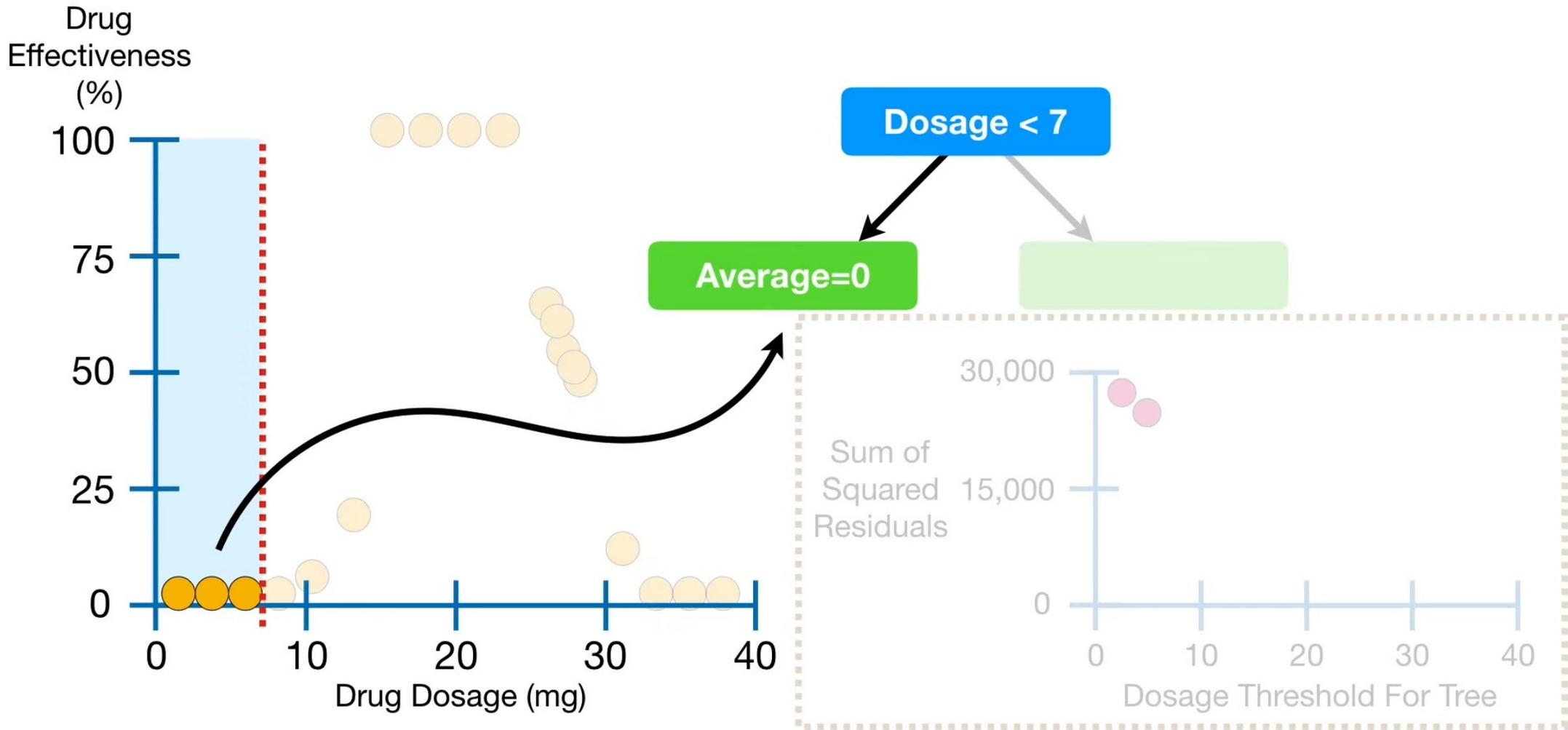
…and use **Dosage < 7** as a new threshold.

Again, the new threshold gives us new predictions…

Drug Effectiveness (%)

Dosage < 7

Average=0

Sum of Squared Residuals

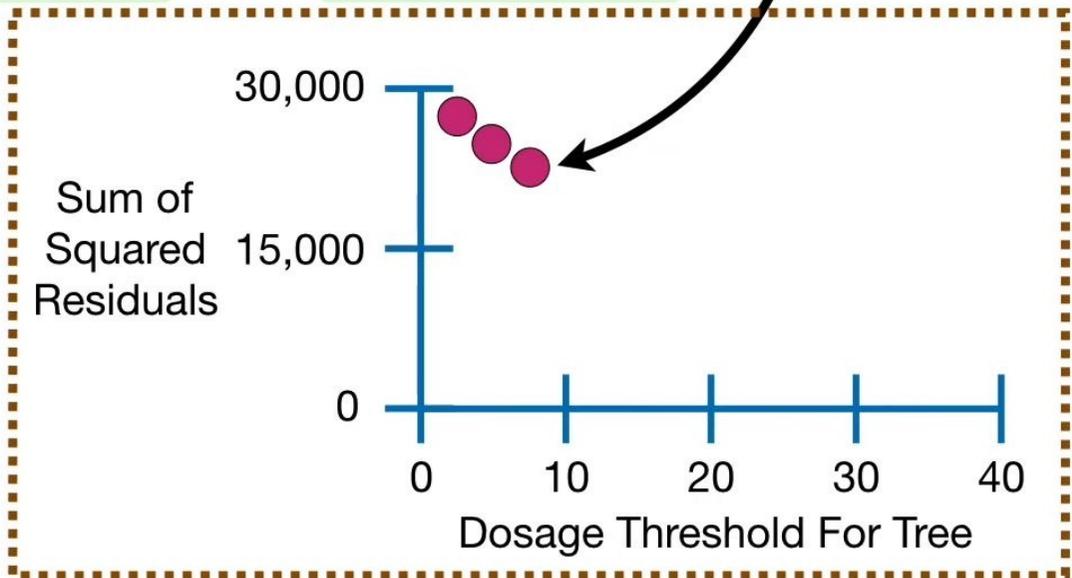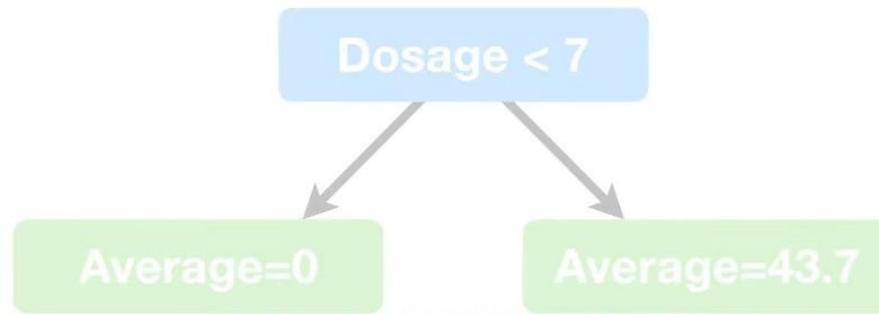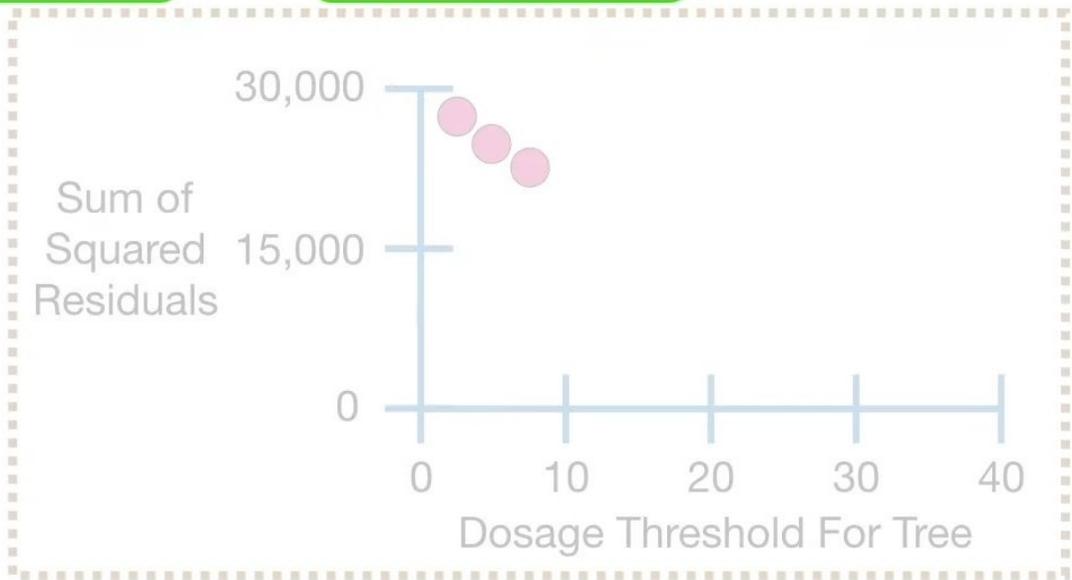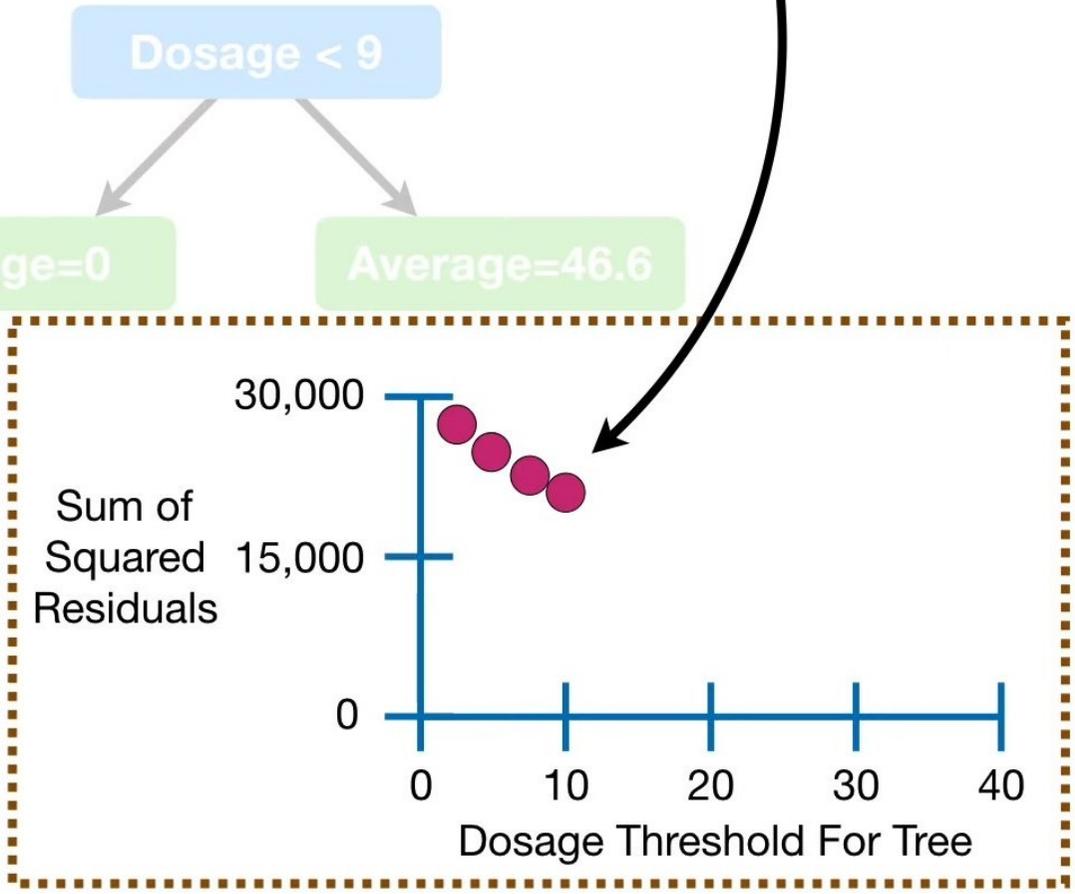Drug Dosage (mg)

Dosage Threshold For Tree

...and a new sum of squared residuals.

Dosage < 7

Average=0          Average=43.7

Drug Effectiveness (%)

100
75
50
25
0

0    10    20    30    40
Drug Dosage (mg)

Sum of Squared Residuals

30,000

15,000

0

0    10    20    30    40
Dosage Threshold For Tree

…and add the new sum of squared residuals to the graph.

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 9

Average=0

Average=46.6

Sum of Squared Residuals

Dosage Threshold For Tree

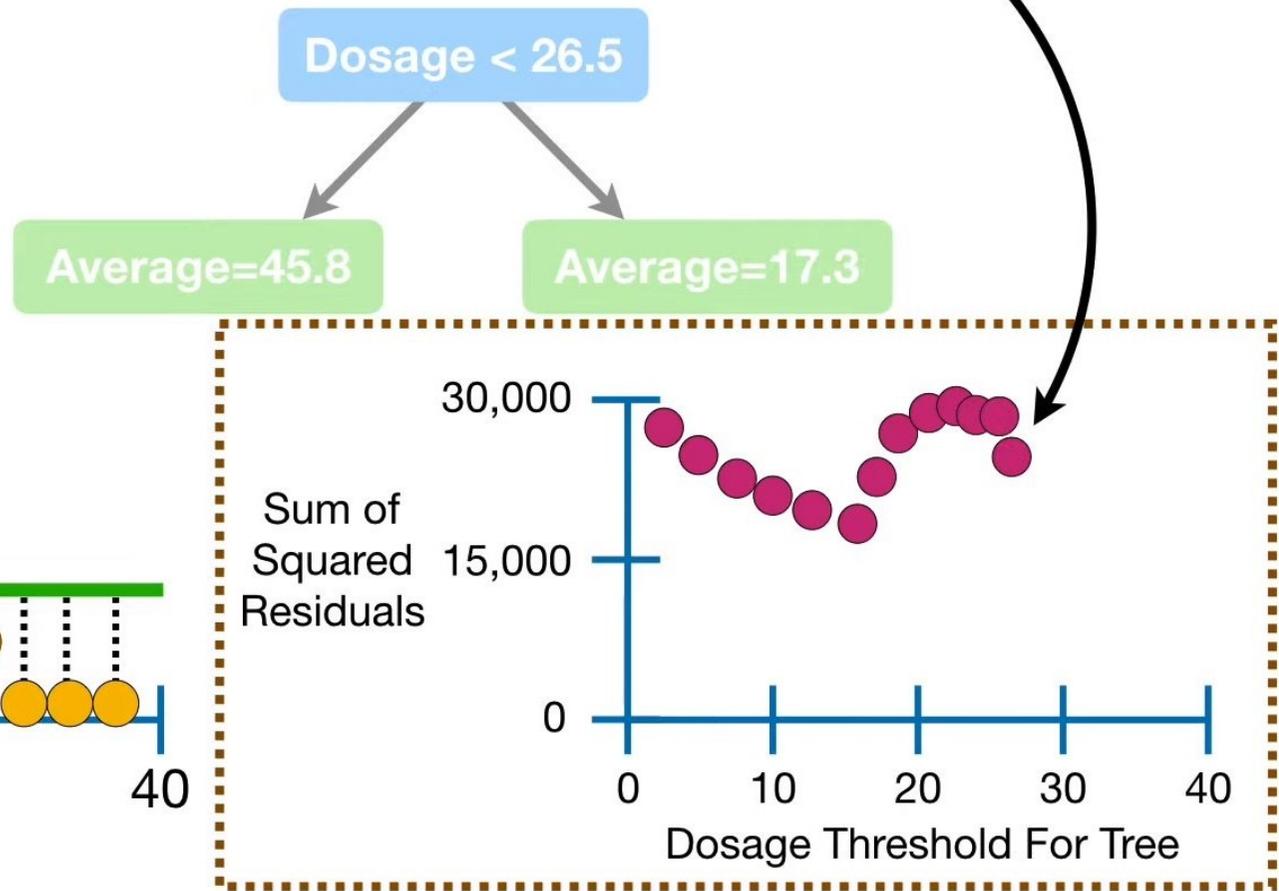And we repeat until we have calculated the sum of squared residuals for all of the remaining thresholds.

Drug Effectiveness (%)
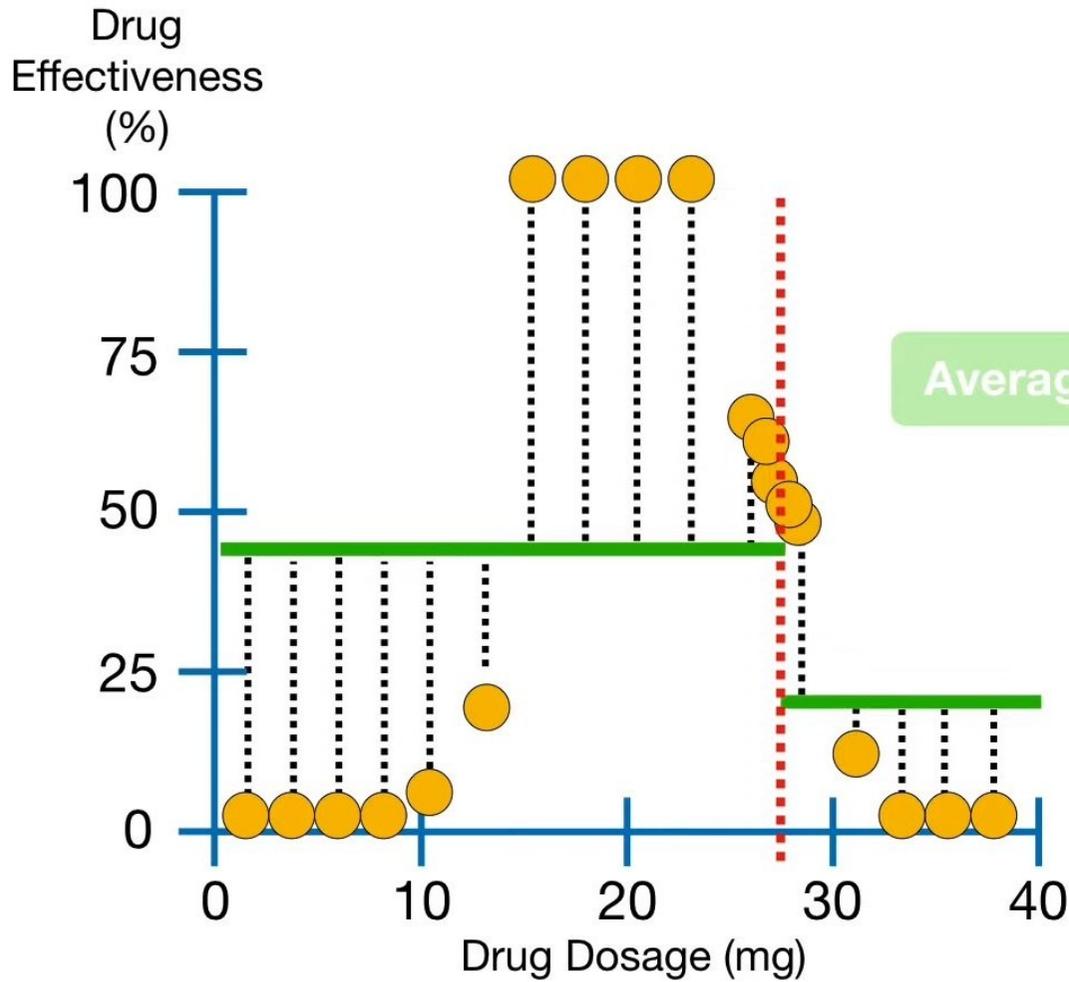
Drug Dosage (mg)

Dosage < 26

Average=44.1

Average=26.8

Sum of Squared Residuals
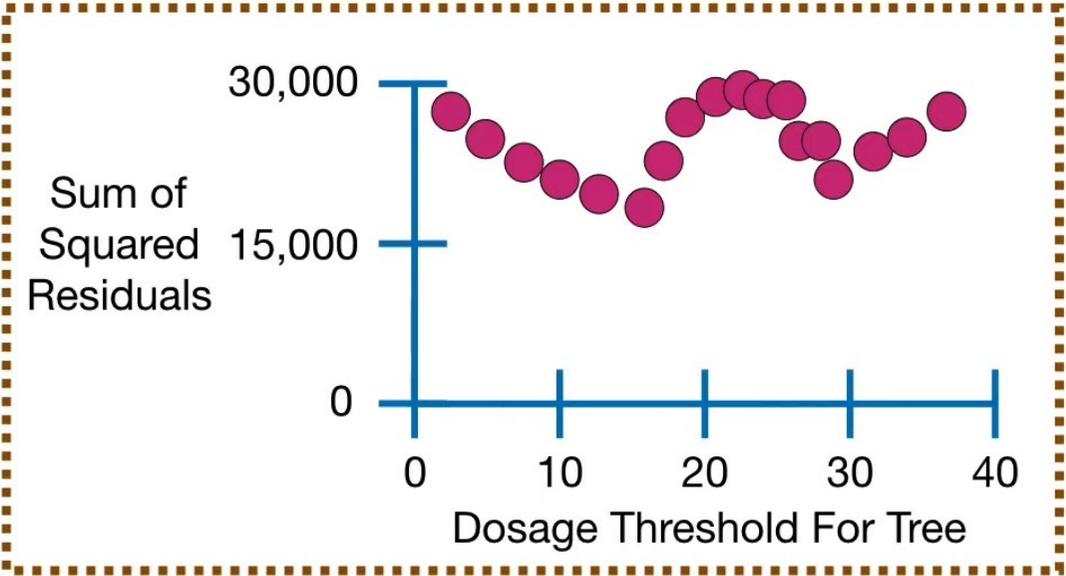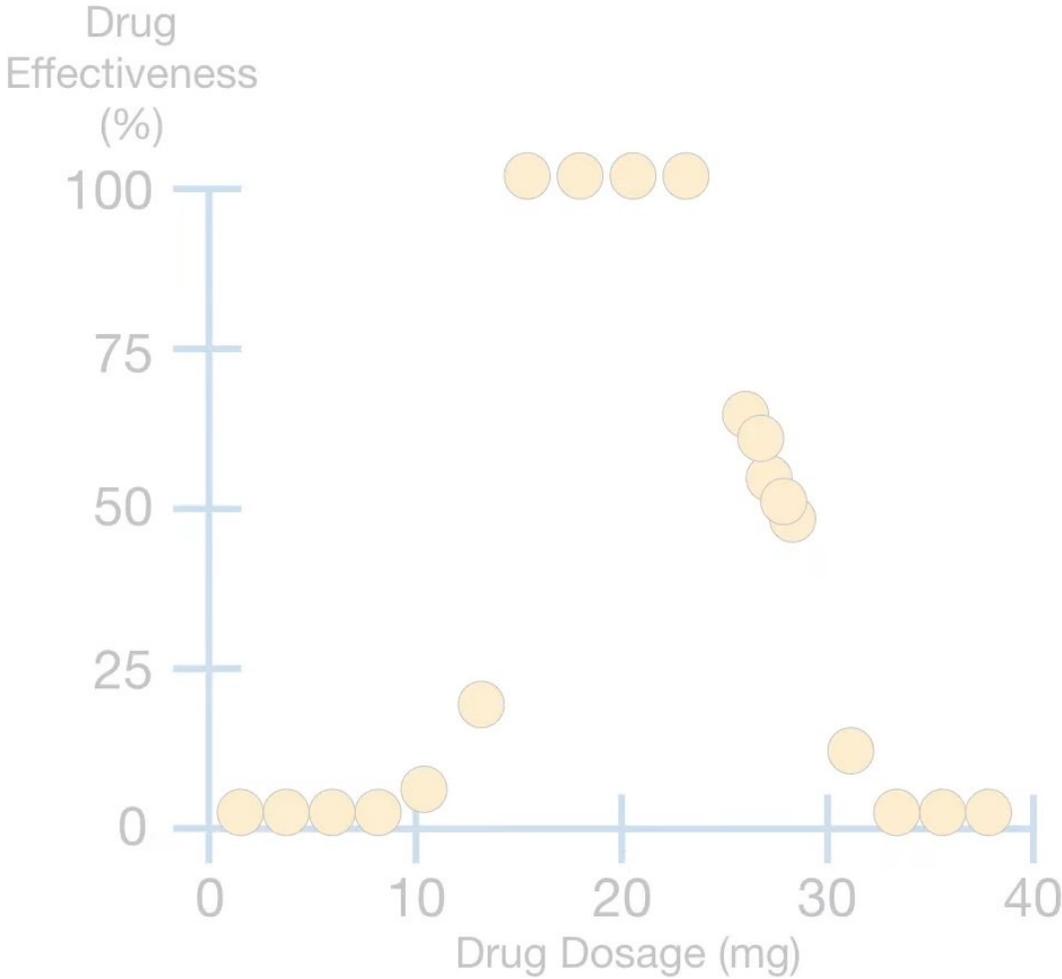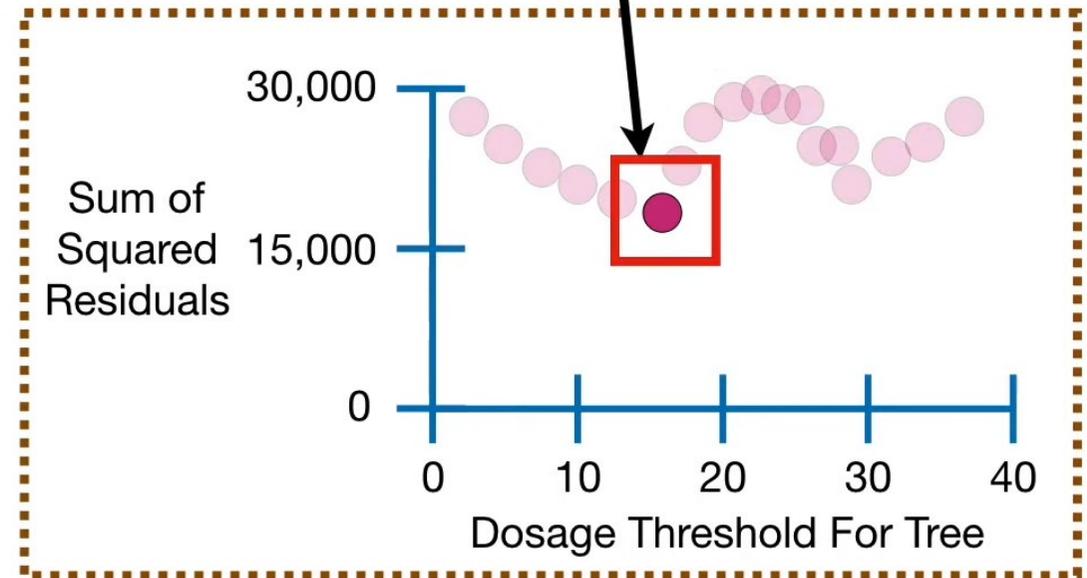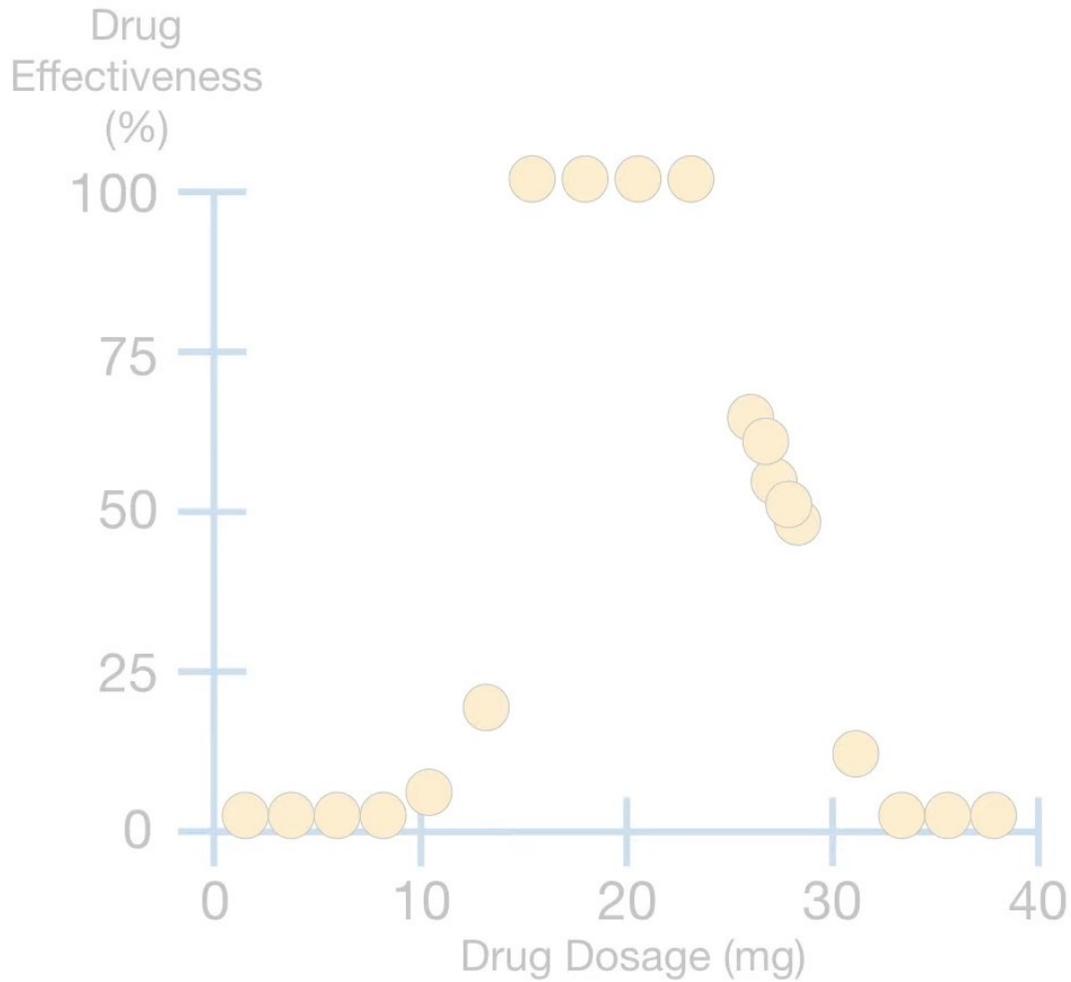
Dosage Threshold For Tree

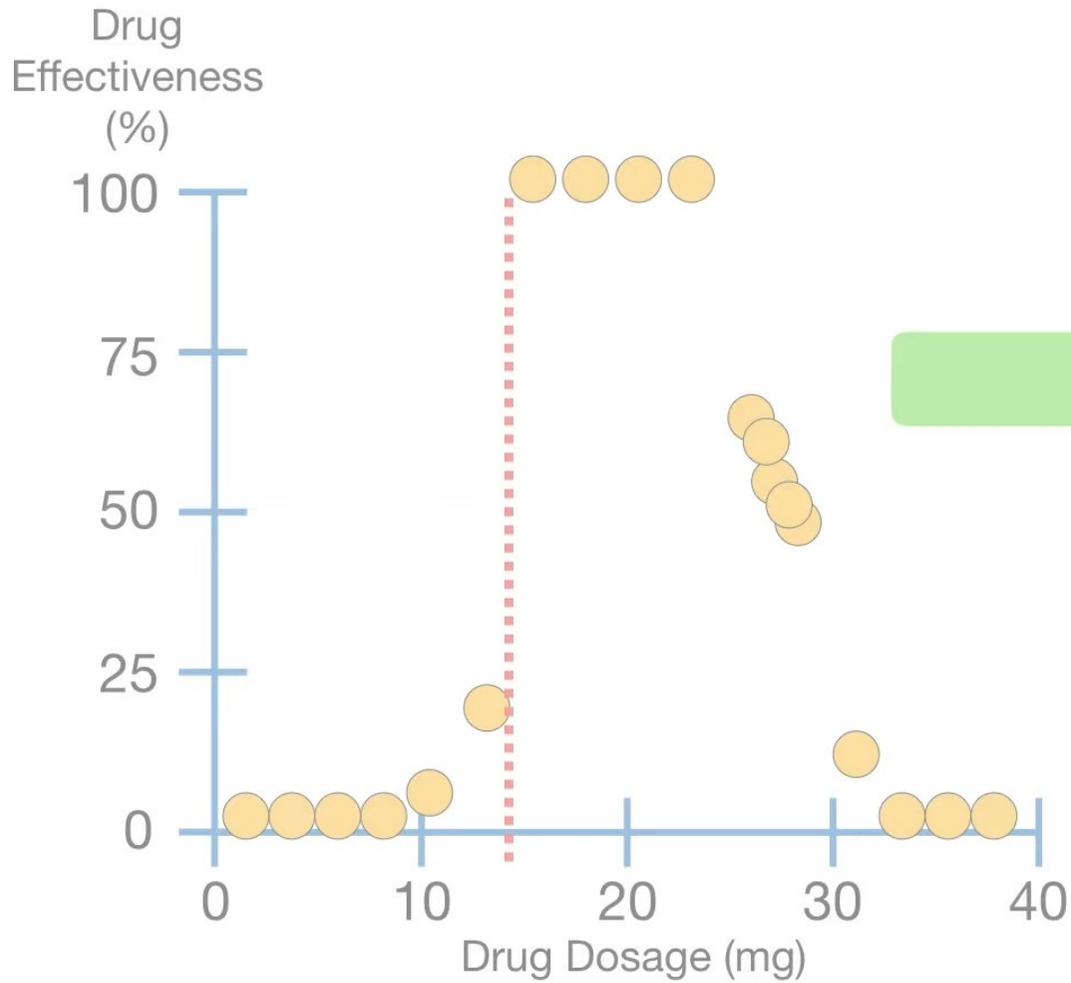And we repeat until we have calculated the sum of squared residuals for all of the remaining thresholds.

Now we can see the sum of squared residuals for all of the thresholds…

…and **Dosage < 14.5** had the smallest sum of squared residuals…
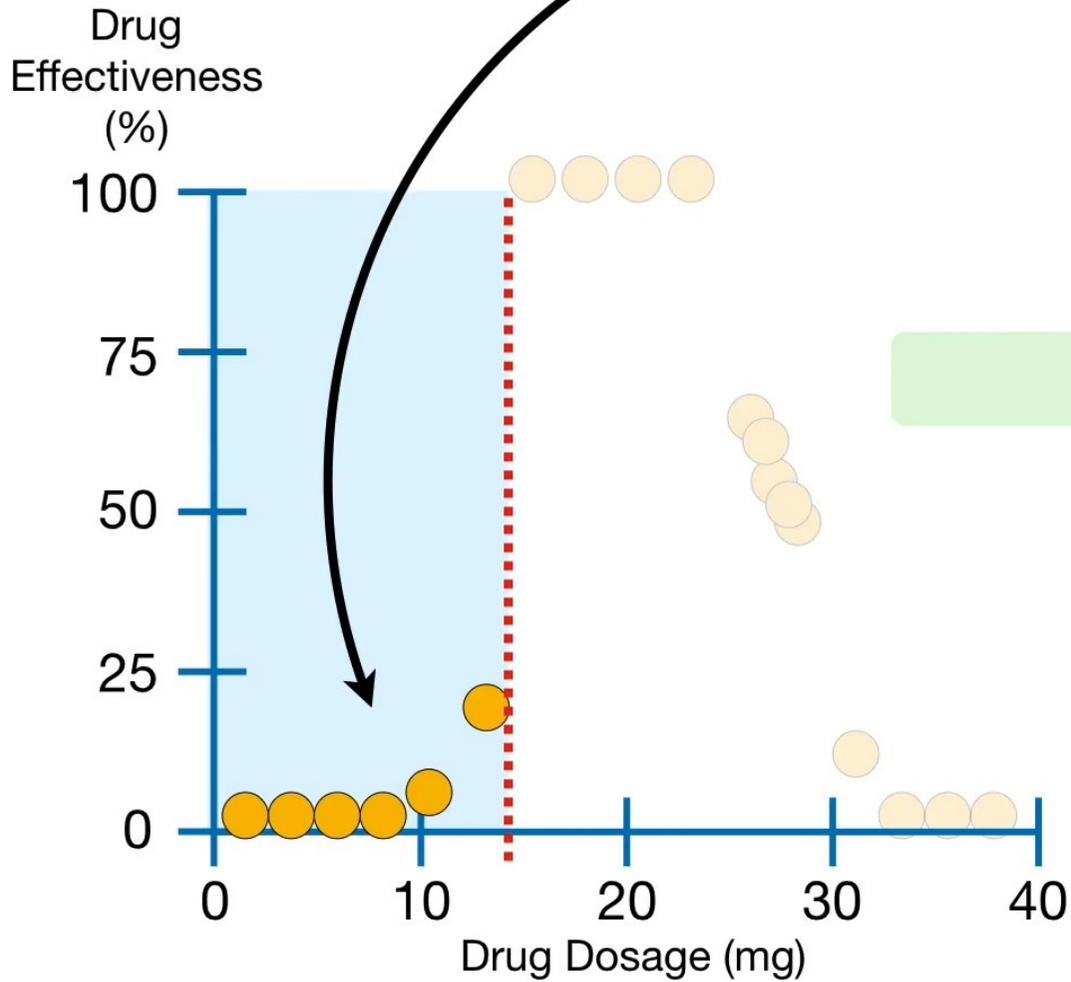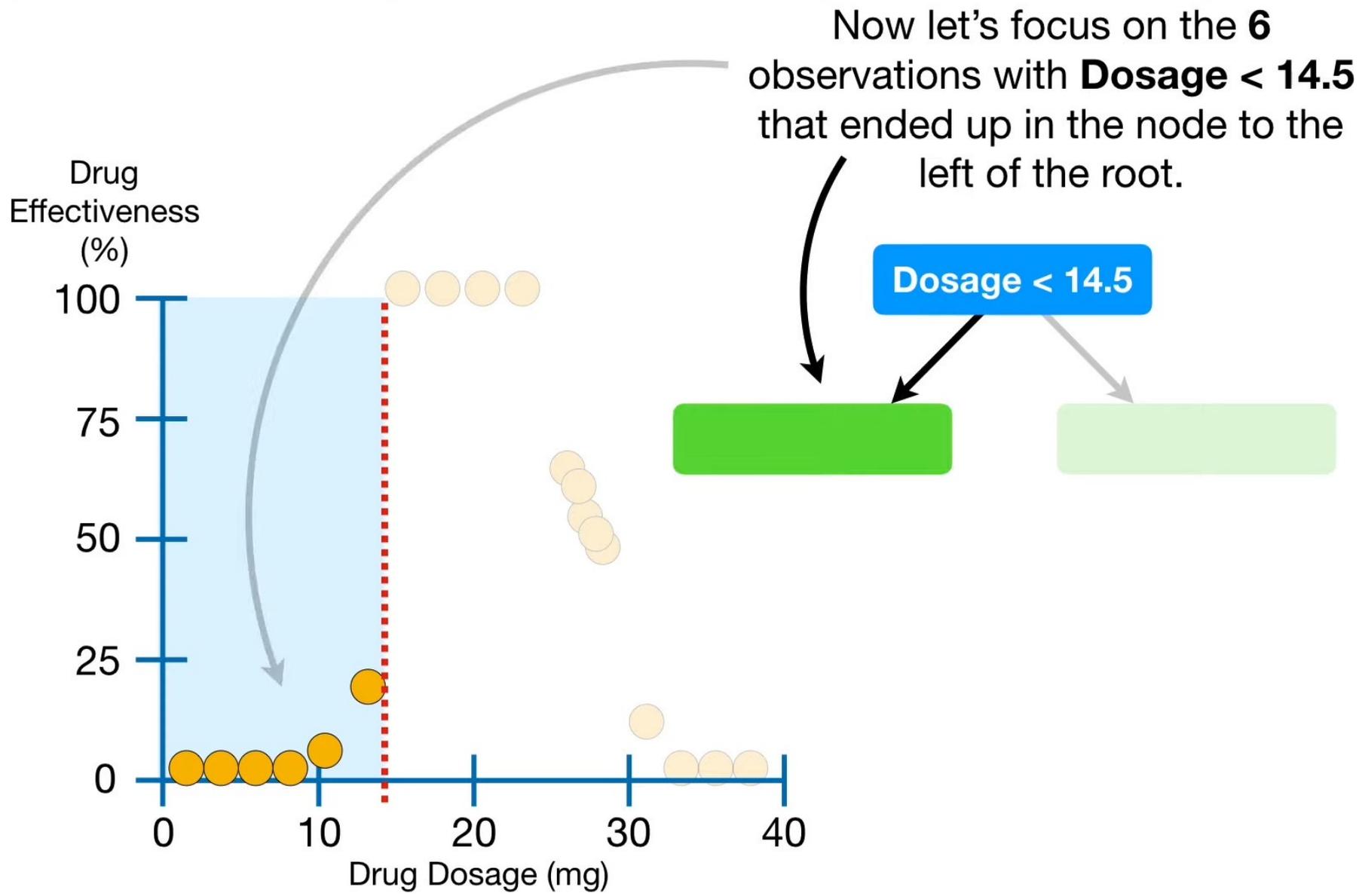
In summary, we split the data into two groups by finding the threshold that gave us the smallest sum of squared residuals.

Dosage < 14.5

Now let's focus on the **6** observations with **Dosage < 14.5** that ended up in the node to the left of the root.

Dosage < 14.5
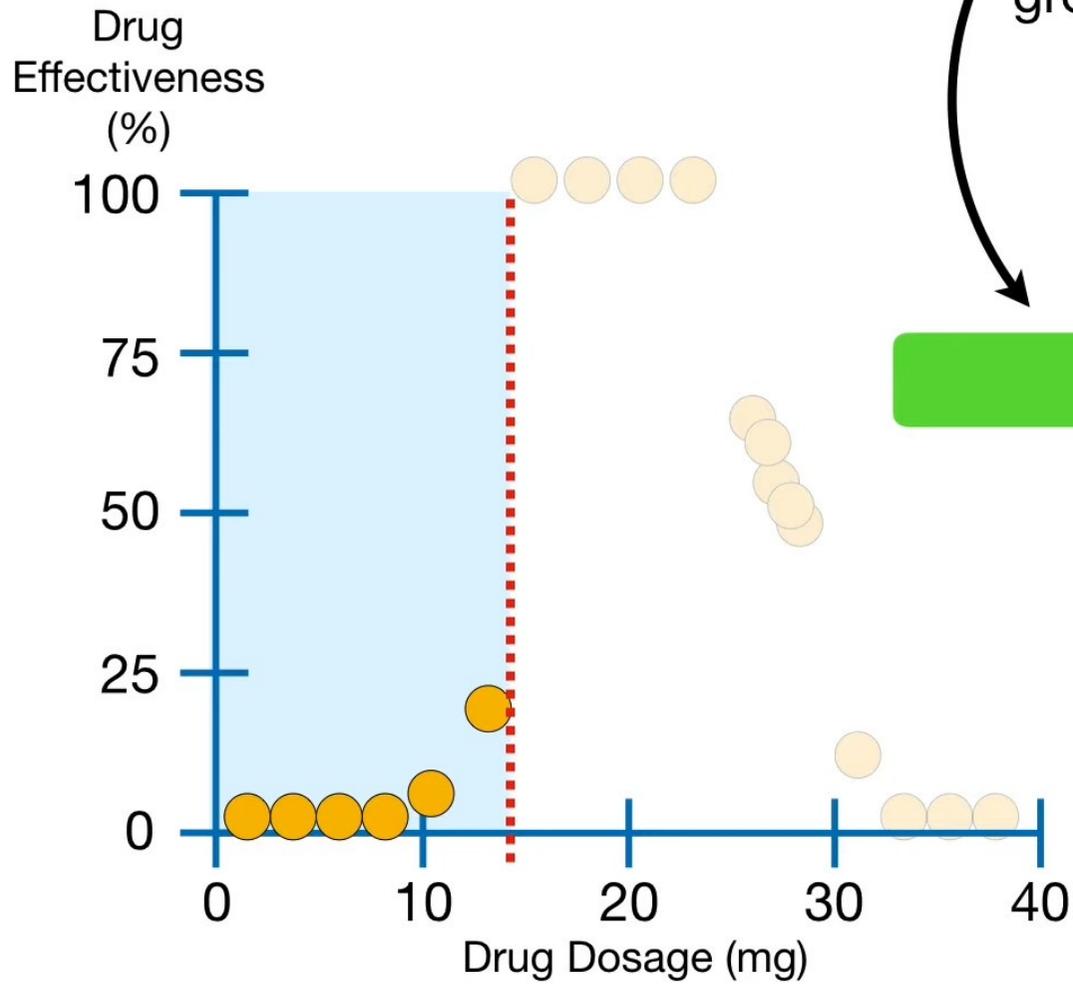
Now let's focus on the **6** observations with **Dosage < 14.5** that ended up in the node to the left of the root.

Dosage < 14.5

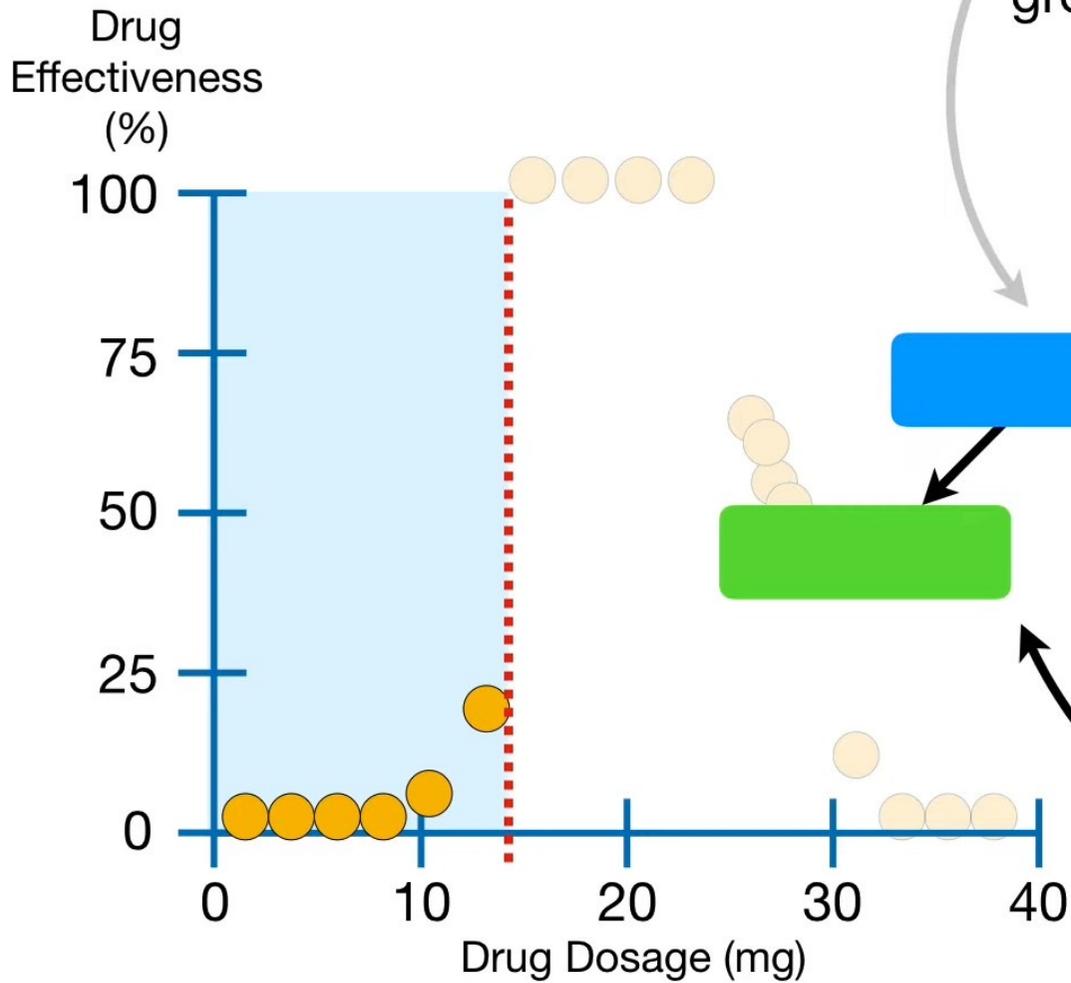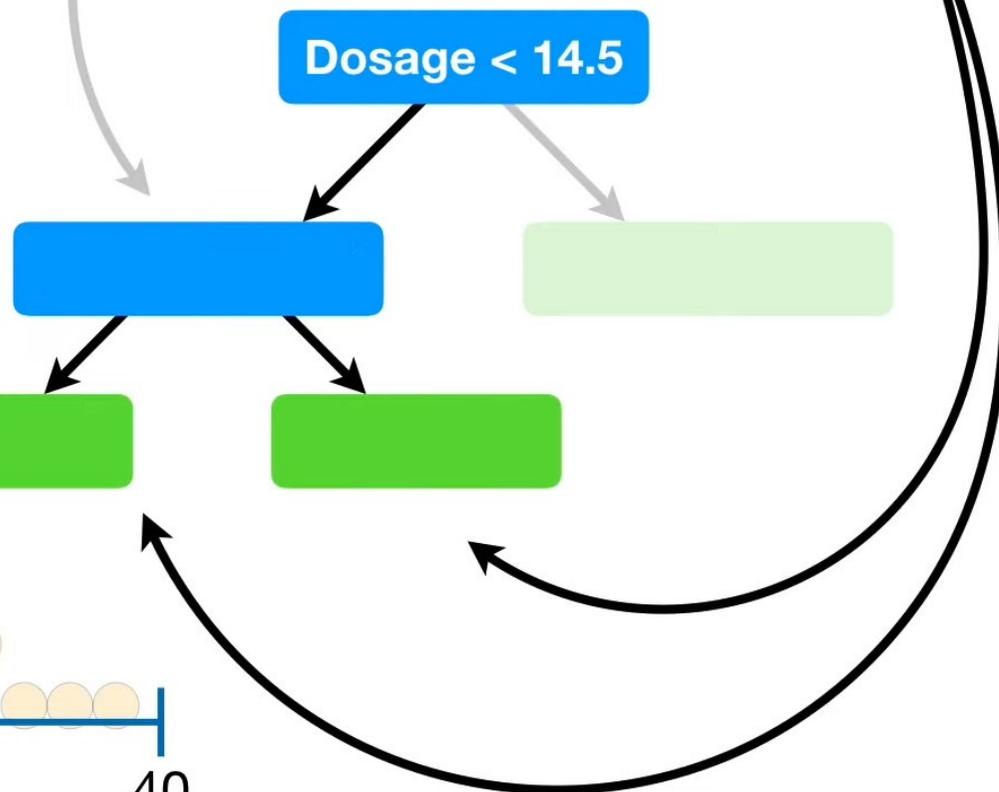...by calculating the sum of squared residuals for different thresholds...

Drug Effectiveness (%)

Dosage < 14.5

Sum of Squared Residuals

Drug Dosage (mg)

Dosage Threshold

...and choosing the threshold with the lowest sum of squared residuals.

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 14.5

Sum of Squared Residuals

Dosage Threshold

…has **Dosage < 14.5**…

...and since we can't split a single observation into two groups, we will call this node a leaf.
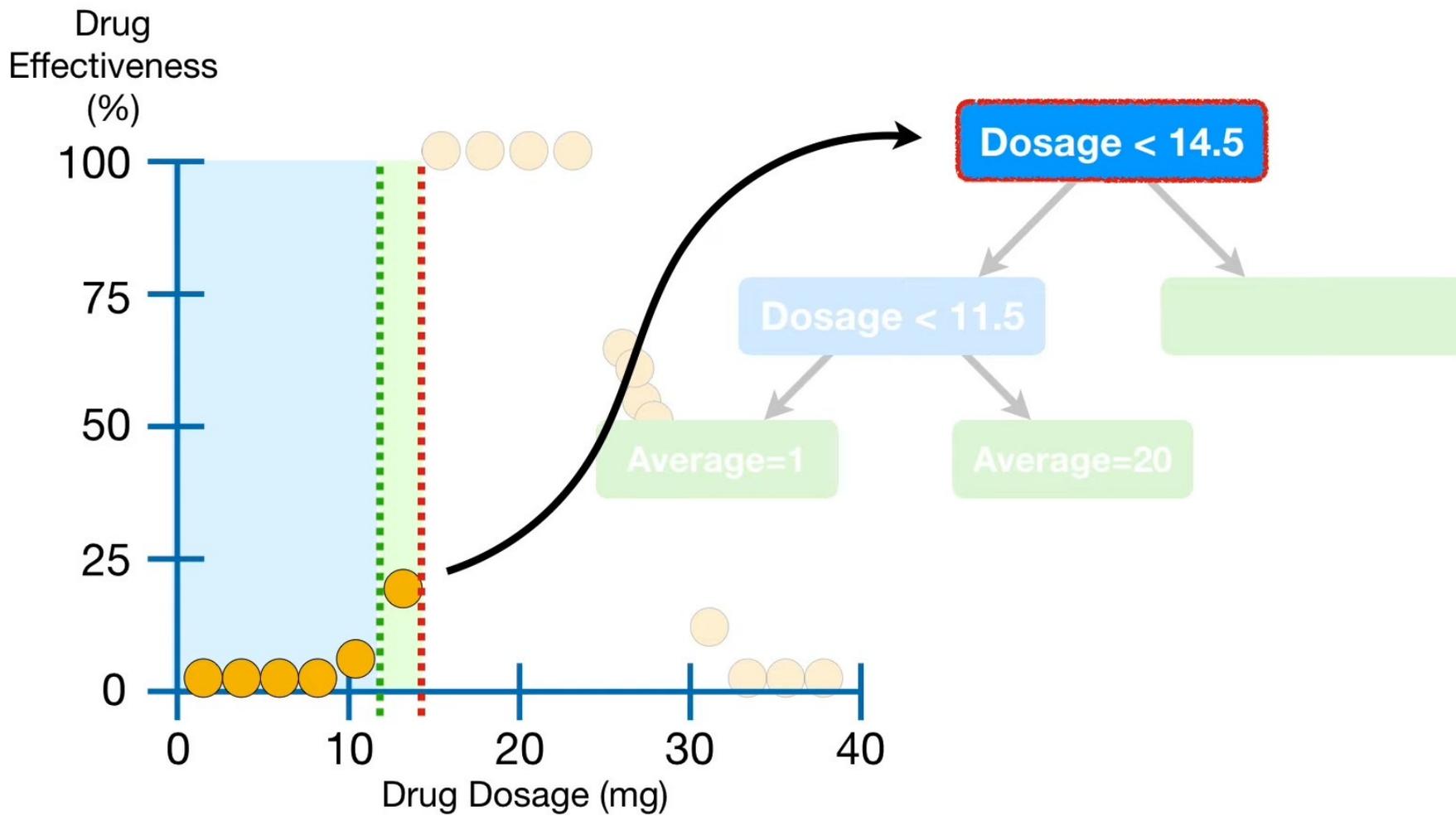
Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 14.5

Dosage < 11.5

Average=1

Average=20

Now we have divided the observations with **Dosage < 14.5** into **3** separate groups.

So we are done splitting the observations with **Dosage < 14.5** into smaller groups.

When a model fits the training data perfectly, it probably means it is overfit and will not perform well with new data.

In **Machine Learning Lingo**, the model has no *Bias*, but potentially large *Variance*.

Is there a way to prevent our tree from overfitting the training data?

The simplest is to only split observations when there are more than some minimum number.

Typically, the minimum number of observations to allow for a split is **20**.

However, since this example doesn't have many observations, I set the minimum to **7**.

Instead, this node will become a leaf...

Dosage < 14.5

...and the output will be the average **Drug Effectiveness** for the **6** observations with **Dosage < 14.5**, **4.2%**.

Dosage < 14.5

Since we have more than **7** observations on the right side (with **Dosage >= 14.5**), we can split them into two groups…

…thus, there are only **4** observations in this node…

…thus, we will make this a leaf because it contains fewer than **7** observations…

Dosage < 14.5

Dosage >= 29

4.2% Effective

...and the output will be average **Drug Effectiveness** for these **4** observations, **2.5%**.

Drug Effectiveness (%)

100

75

50

25

0

0    10    20    30    40

Drug Dosage (mg)

Dosage < 14.5

Dosage >= 29

4.2% Effective

2.5% Effective

Now we need to figure out what to do with the **9** observations with **Dosages** between **14.5** and **29**.

Dosage < 14.5

4.2% Effective

Dosage >= 29

2.5% Effective

...by finding the threshold that gives us the minimum sum of squared residuals.

Drug Effectiveness (%)

Drug Dosage (mg)

Dosage < 14.5

4.2% Effective

Dosage >= 29

2.5% Effective

Dosage >= 23.5

So we use the average **Drug Effectiveness** for observations with **Dosages** between **14.5** and **23.5**, **100%**, as the output for leaf on the right…

…and we use the average **Drug Effectiveness** for observations with **Dosages** between **23.5** and **29**, **52.8%**, as the output for leaf on the left.

…and each leaf corresponds to the average **Drug Effectiveness** from a different cluster of observations.

So far we have built a tree using a single predictor, **Dosage**, to predict **Drug Effectiveness**.

| Dosage | Drug Effect. |
|--------|--------------|
| 10 | 58 |
| 20 | 60 |
| 35 | 57 |
| 5 | 44 |
| etc… | etc… |

Now let's talk about how to build a tree to predict
**Drug Effectiveness** using a bunch of predictors.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Just like before, we will start by using
**Dosage** to predict **Drug Effectiveness**.

| Dosage | Age | Sex | Drug Effect. |
|--------|------|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Thus, just like before, we will try different thresholds for **Dosage** and calculate the sum of squared residuals at each step…

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

...and pick the threshold that gives us the minimum sum of squared residuals.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc... | etc... | etc... | etc... |

The best threshold becomes a *candidate* for the root.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Now we focus on using **Age** to predict **Drug Effectiveness**.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Just like with **Dosage**, we try different thresholds for **Age** and calculate the sum of squared residuals at each step…

| Dosage | Age | Sex | Drug Effect. |
|--------|------|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

...and pick the one that gives us the minimum sum of squared residuals.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

The best threshold becomes another *candidate* for the root.

Now we focus on using **Sex** to predict **Drug Effectiveness**.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

With **Sex**, there is only one threshold to try…

| Dosage | Age | Sex | Drug Effect. |
|---|---|---|---|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

...so we use that threshold to calculate the sum of squared residuals...

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc... | etc... | etc... | etc... |

...and that becomes another *candidate* for the root.

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc... | etc... | etc... | etc... |

Dosage < 14.5

Average=4.2    Average=51.8

Age > 50

Average=3    Average=52

Sex=Female

Average=12    Average=40

Dosage < 14.5

Average=4.2  Average=51.8

SSR = 19,564

Age > 50

Average=3  Average=52

SSR = 12,017

Sex=Female

Average=12  Average=40

SSR = 20,738

Now we compare the sum of squared residuals (SSRs) for each candidate…

Dosage < 14.5
SSR = 19,564
Average=4.2
Average=51.8

Age > 50
SSR = 12,017
Average=3
Average=52

...and pick the candidate with the lowest value.

Sex=Female
SSR = 20,738
Average=12
Average=40

Since **Age > 50** had the lowest sum of squared residuals, it becomes the root of the tree.



Age > 50

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Then we grow the tree just like before, except now
we compare the lowest sum of squared residuals
from each predictor.



| Dosage | Age | Sex | Drug Effect. |
|--------|------|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Then we grow the tree just like before, except now
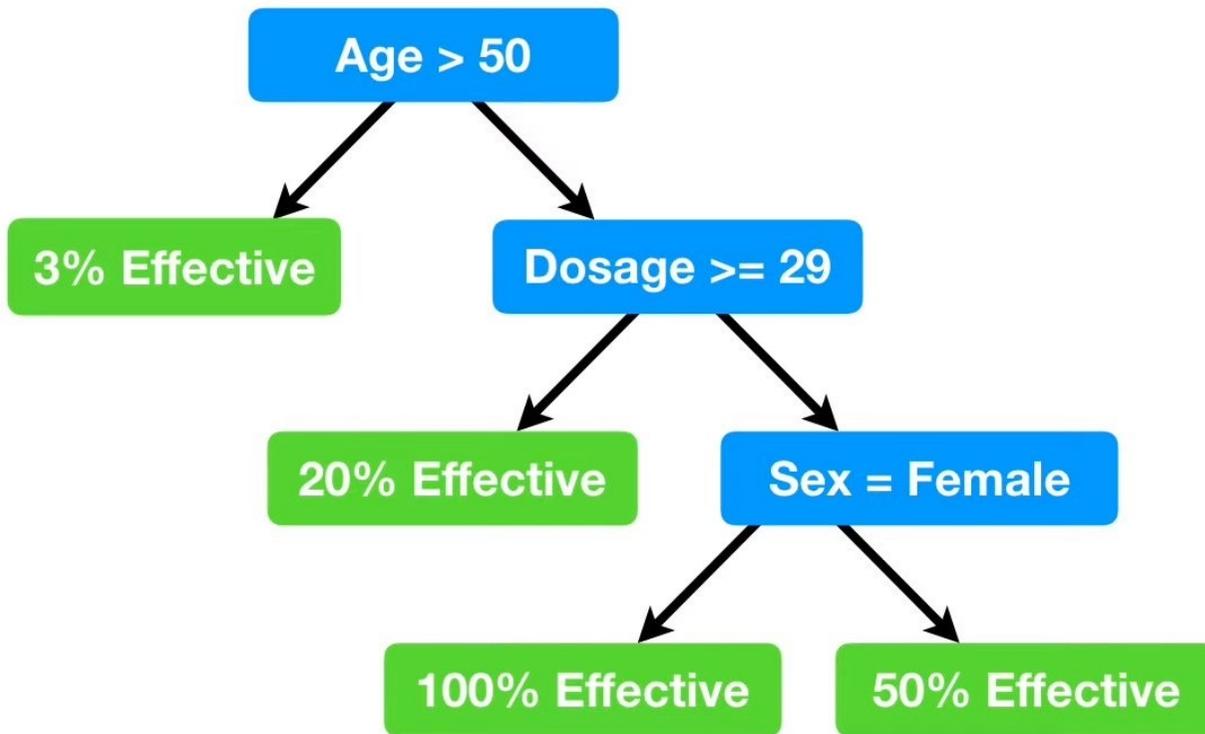we compare the lowest sum of squared residuals
from each predictor.



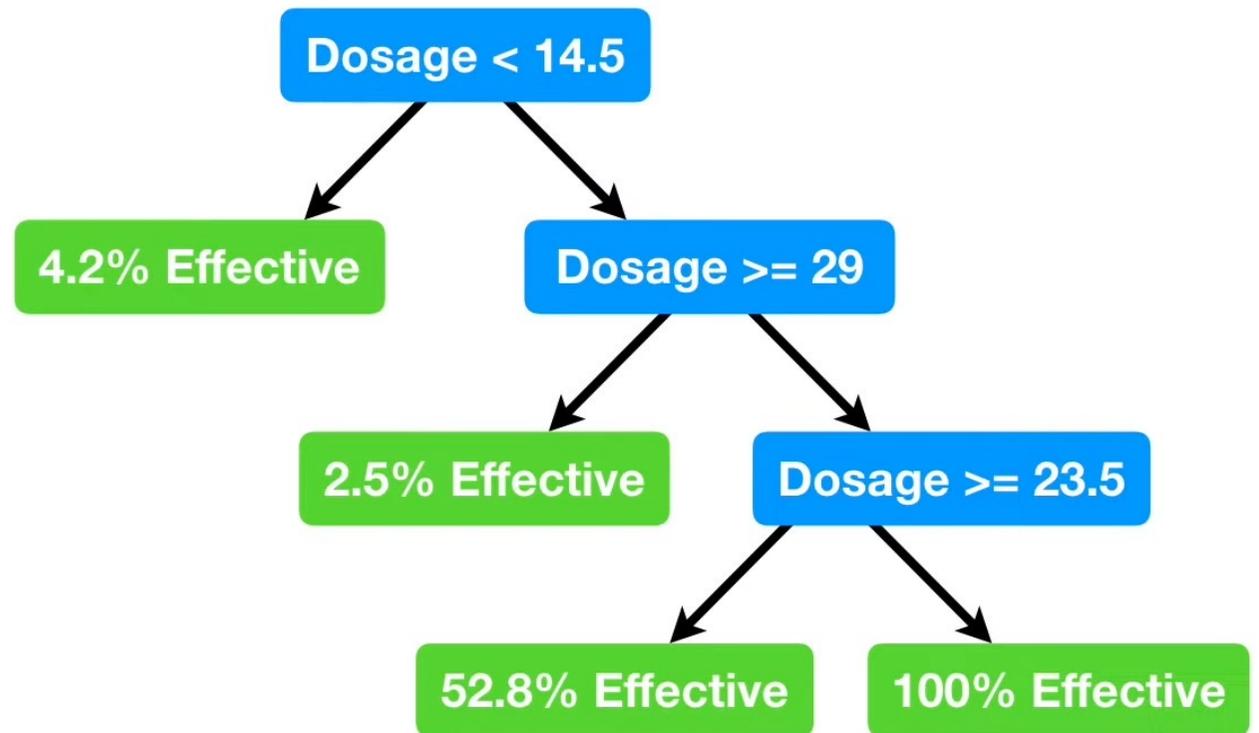| Dosage | Age | Sex | Drug Effect. |
|--------|------|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

And just like before, when a leaf has less than a
minimum number of observations, which is usually
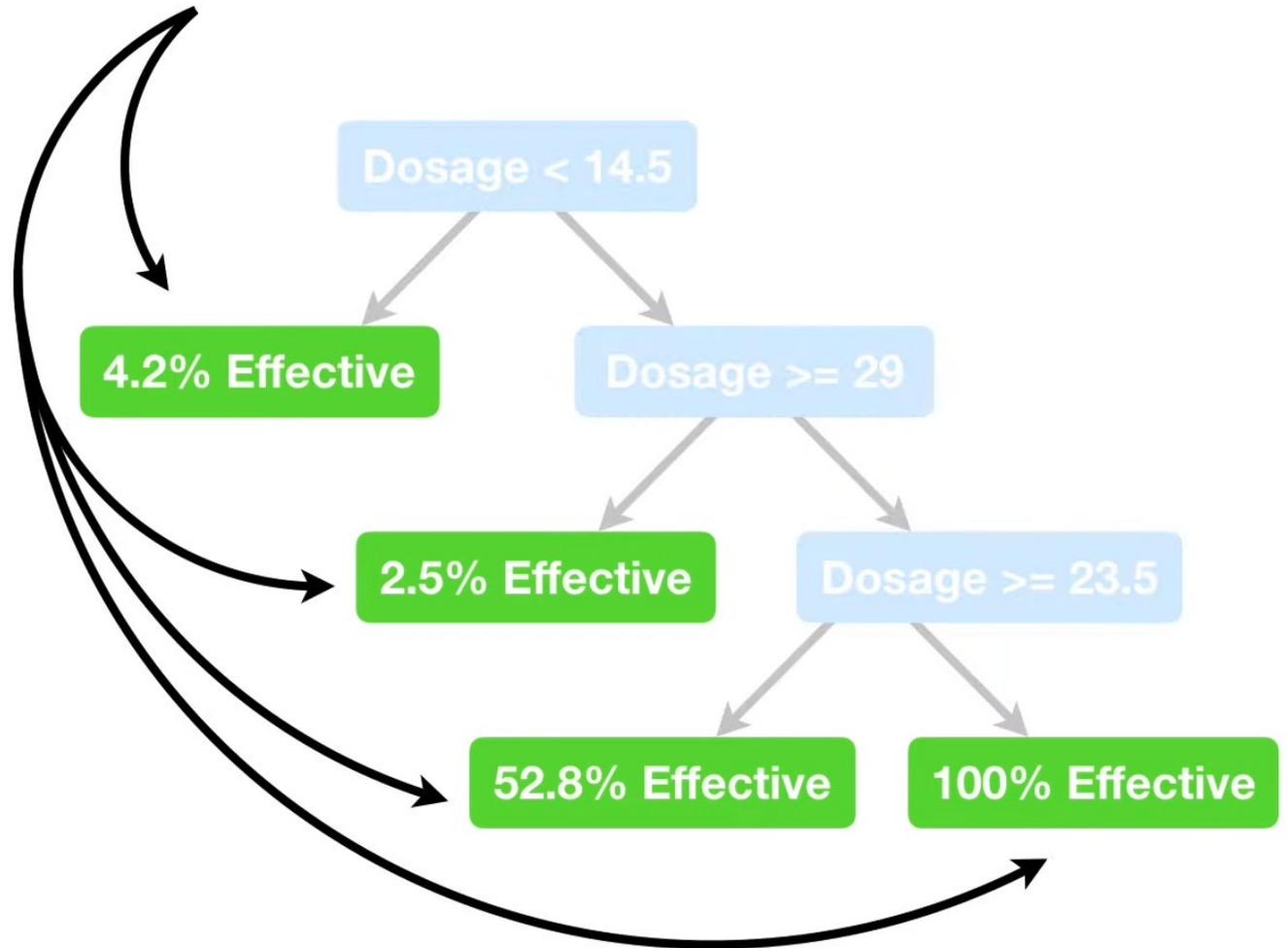**20**, but we are using **7**, we stop trying to divide them.



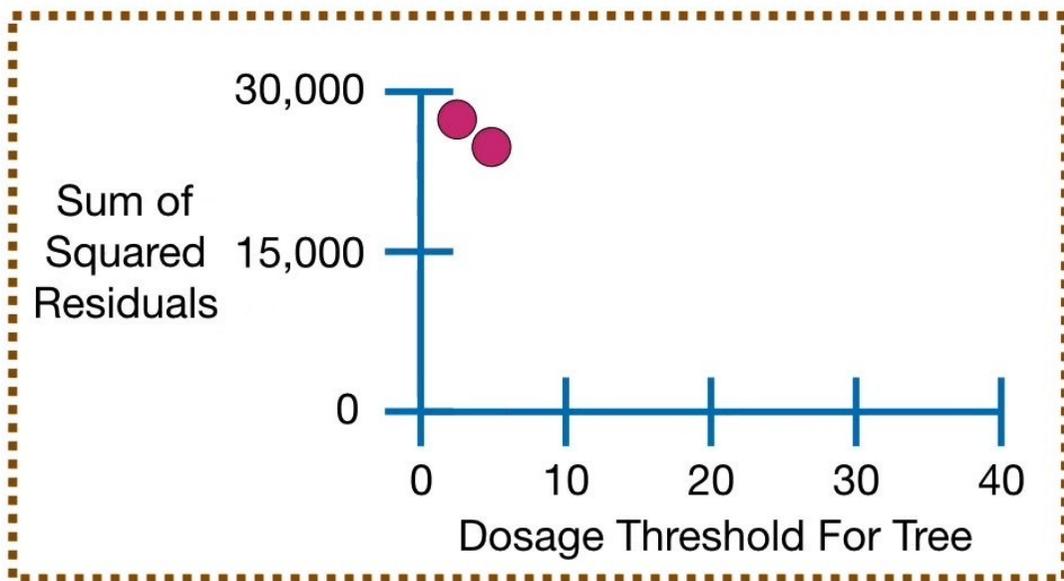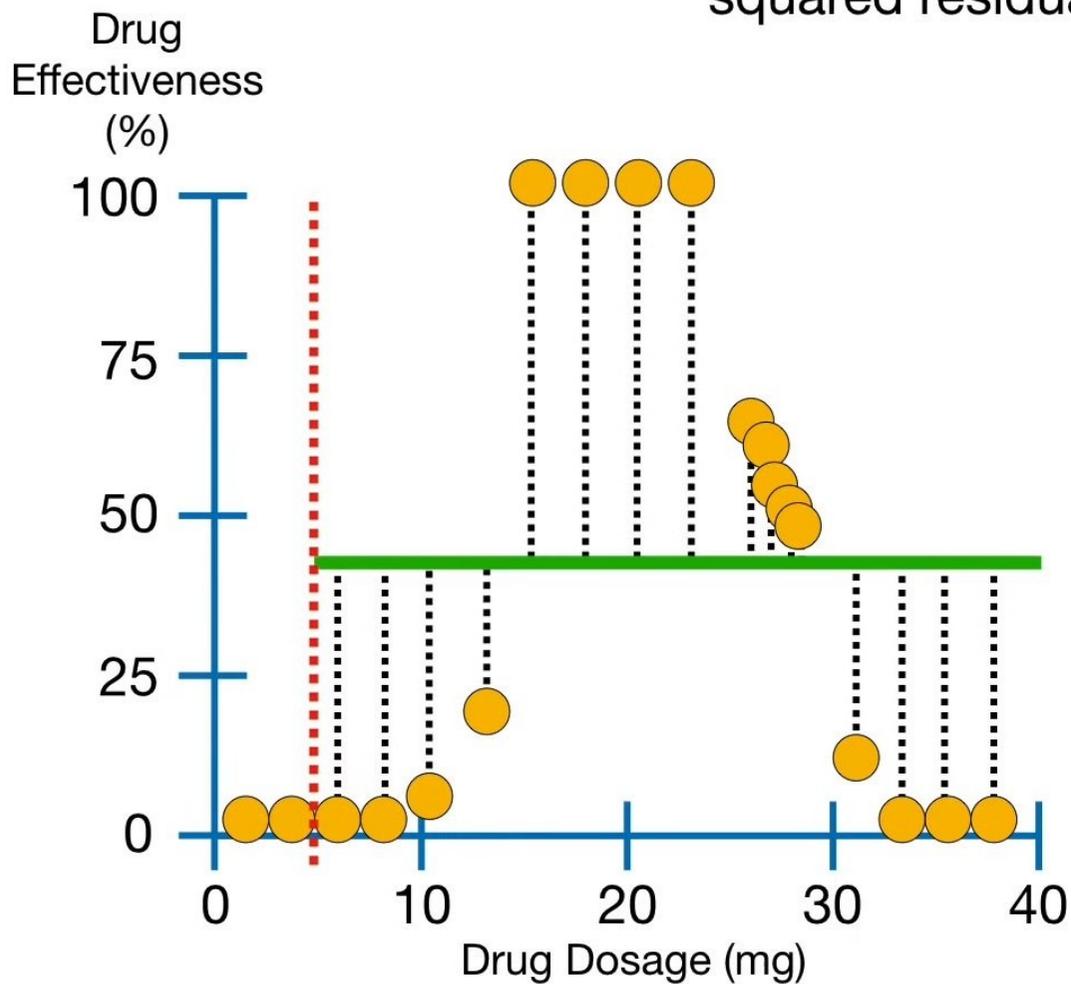| Dosage | Age | Sex | Drug Effect. |
|--------|-----|-----|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

In summary…
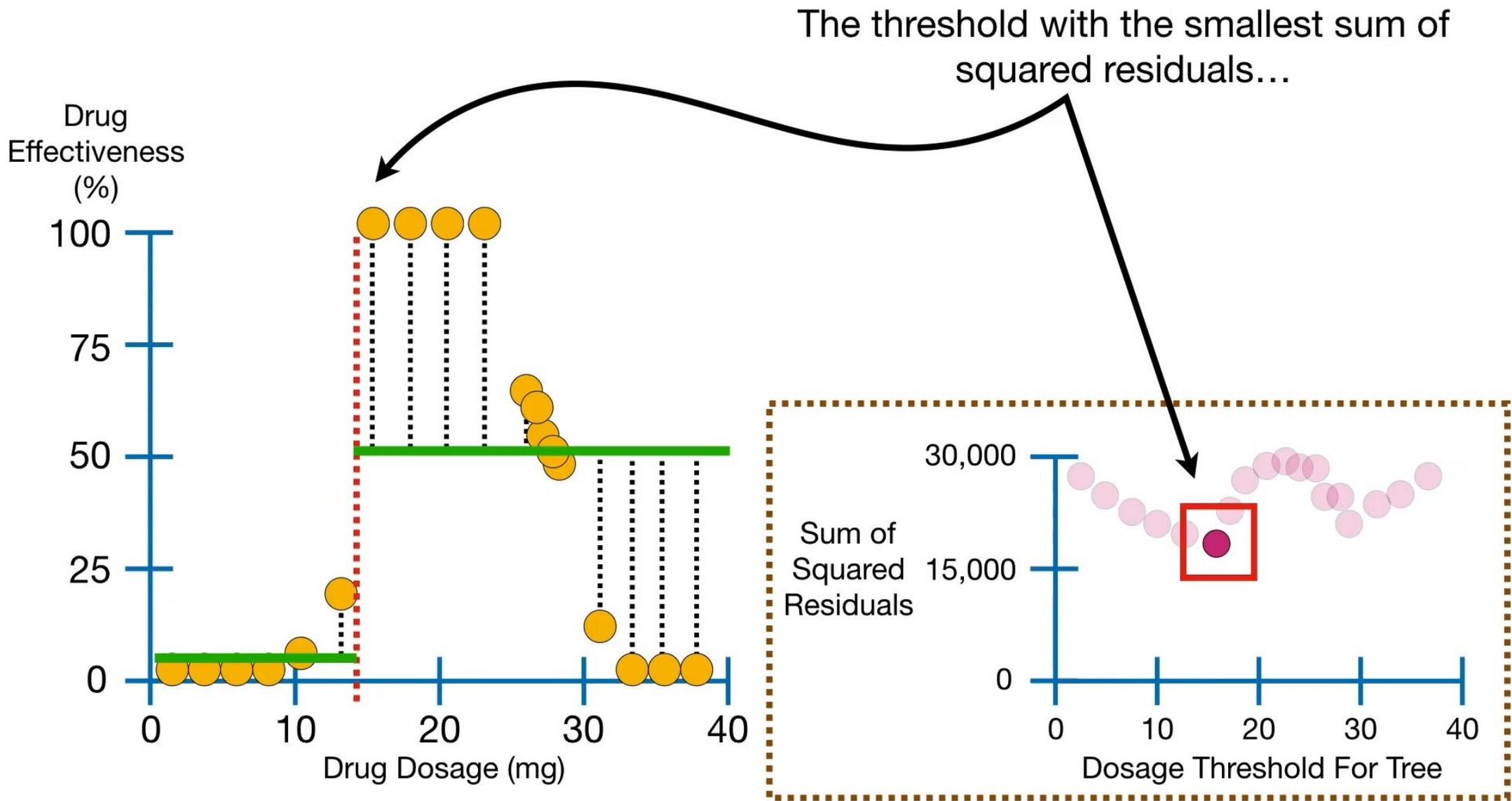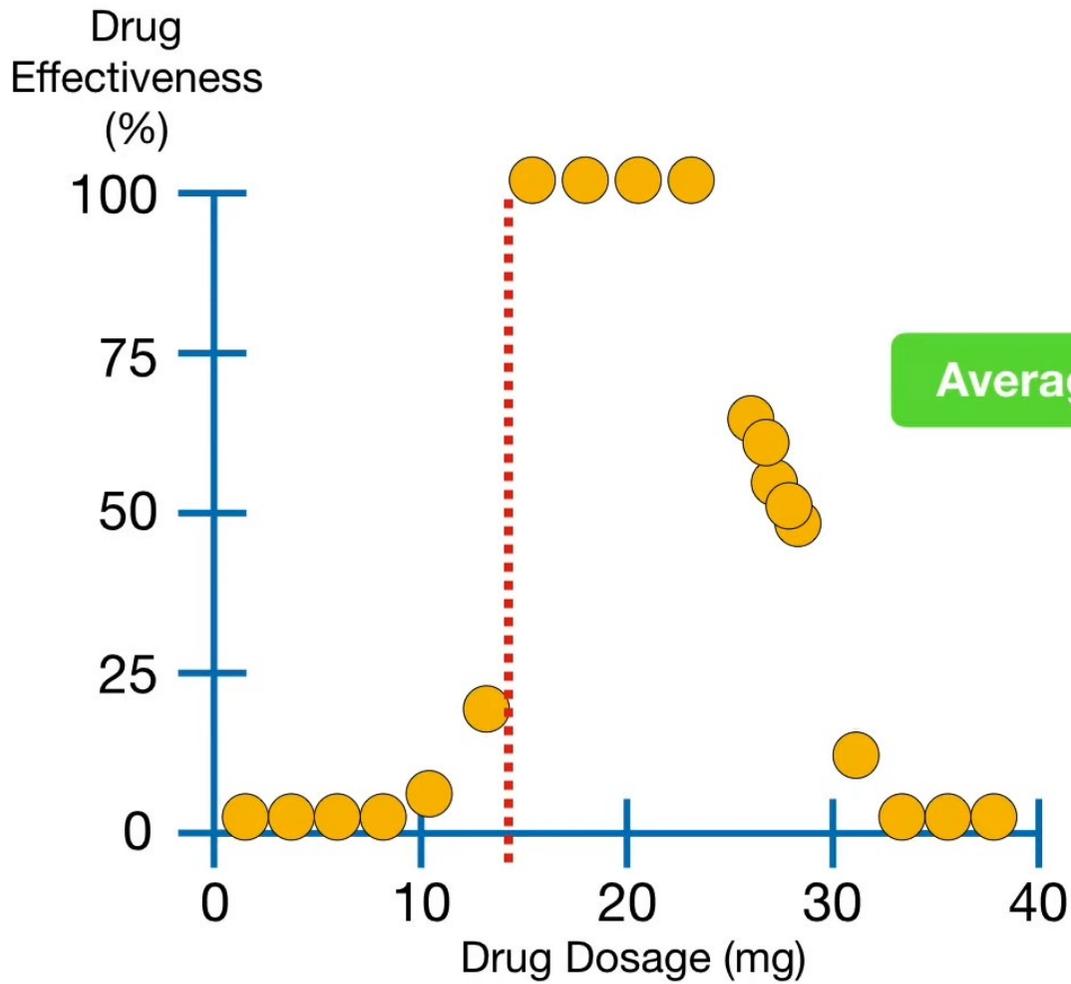
**Regression Trees** are a type of **Decision Tree.**

In a **Regression Tree**, each leaf represents a numeric value.

We determine how to divide the observations by trying different thresholds and calculating the sum of squared residuals at each step.
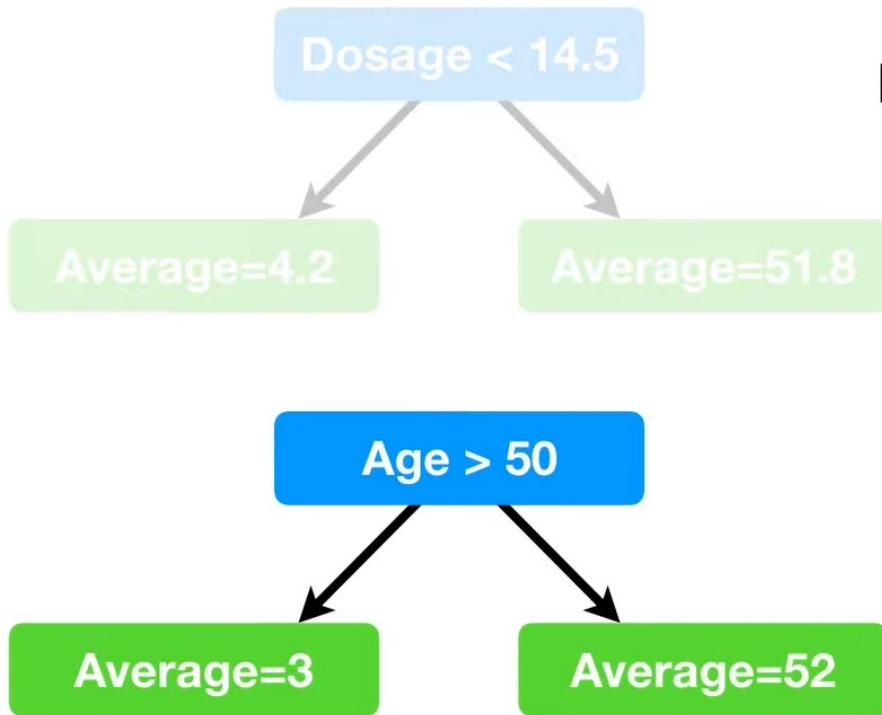
The threshold with the smallest sum of squared residuals…

Dosage < 14.5
→ Average=4.2
→ Average=51.8

Age > 50
→ Average=3
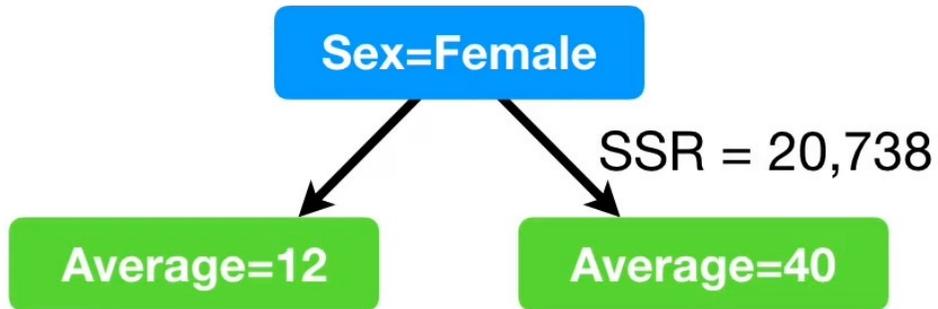→ Average=52

If we have more than one predictor, we find the optimal threshold for each one…

| Dosage | Age | Sex | Drug Effect. |
|--------|-----|--------|--------------|
| 10 | 25 | Female | 98 |
| 20 | 73 | Male | 0 |
| 35 | 54 | Female | 6 |
| 5 | 12 | Male | 44 |
| etc… | etc… | etc… | etc… |

Dosage < 14.5
Average=4.2    Average=51.8
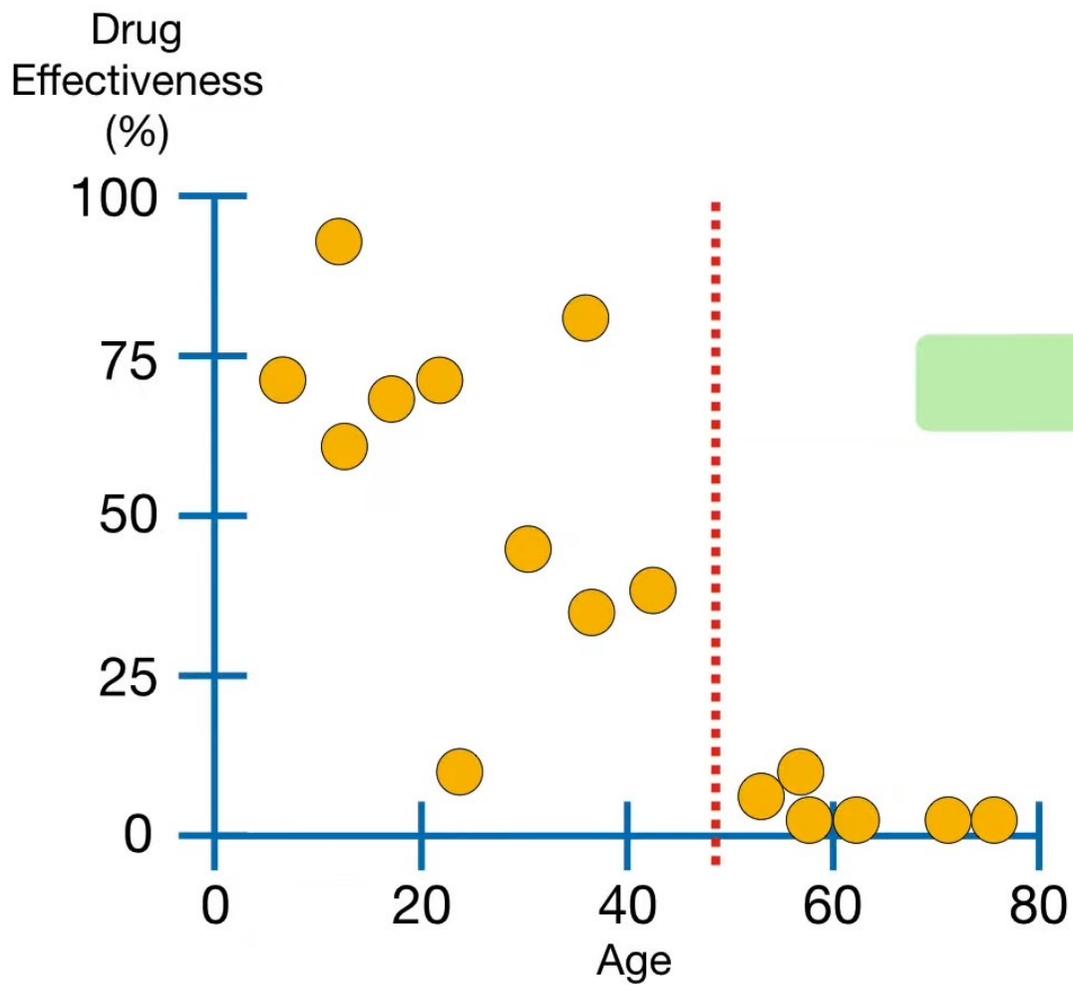SSR = 19,564

Age > 50
Average=3    Average=52
SSR = 12,017

Sex=Female
Average=12    Average=40
SSR = 20,738

…and we pick the candidate with the smallest sum of squared residuals…

…until we can no longer split the observations into smaller groups…

...and then we are done.