# CMSC 473/673 Natural Language Processing

Instructor: Lara J. Martin (she/they)

TA: Duong Ta (he)

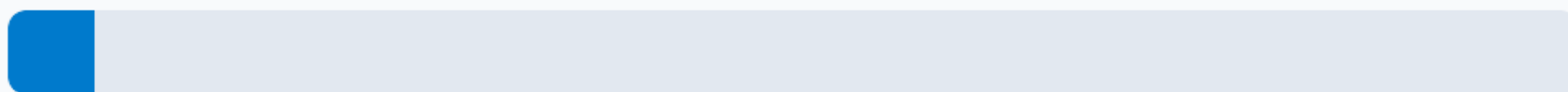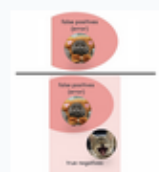*Slides modified from Dr. Frank Ferraro*

# Learning Objectives

Develop an intuition about precision & recall

Extend P/R to multi-class problems

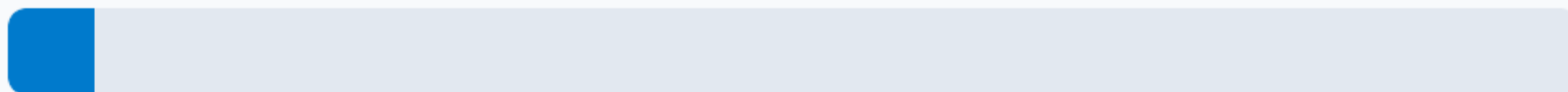Identify when you might want certain evaluation metrics over others

If you are classifying pictures of dogs, what would be the "equation" for *recall* (where the top of the image is the numerator and the bottom of the image is the denominator)?
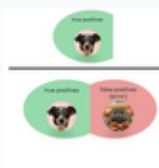
 — 3%

✓  — 56%

 — 3%

 — 38%

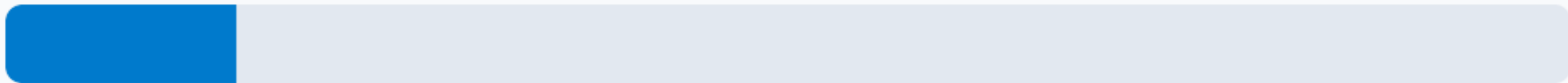# The difference between classification & regression is that a regression model will produce a continuous output.
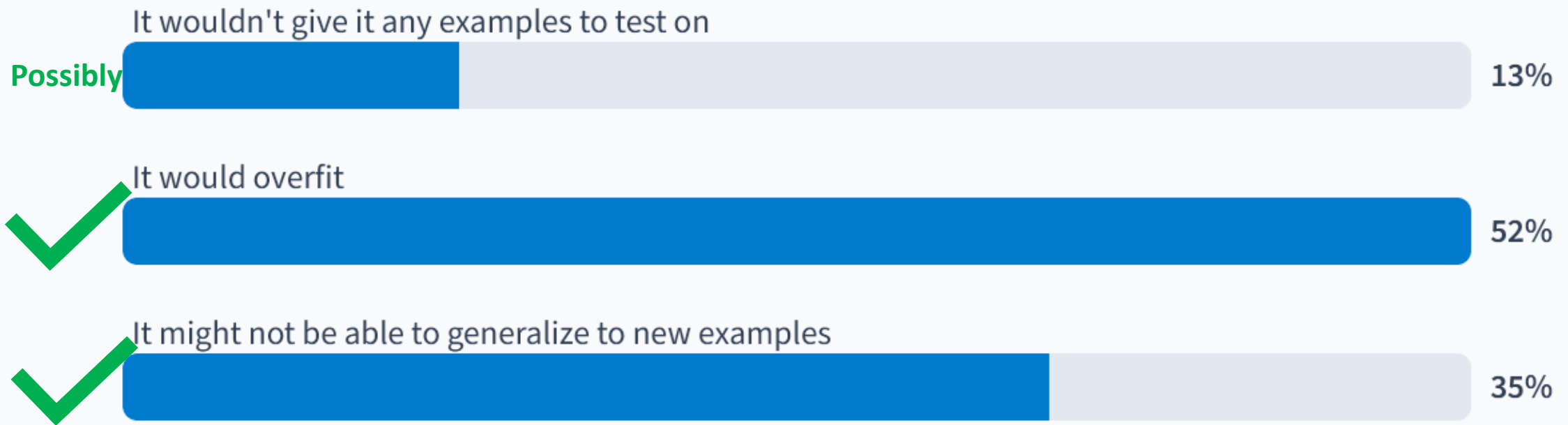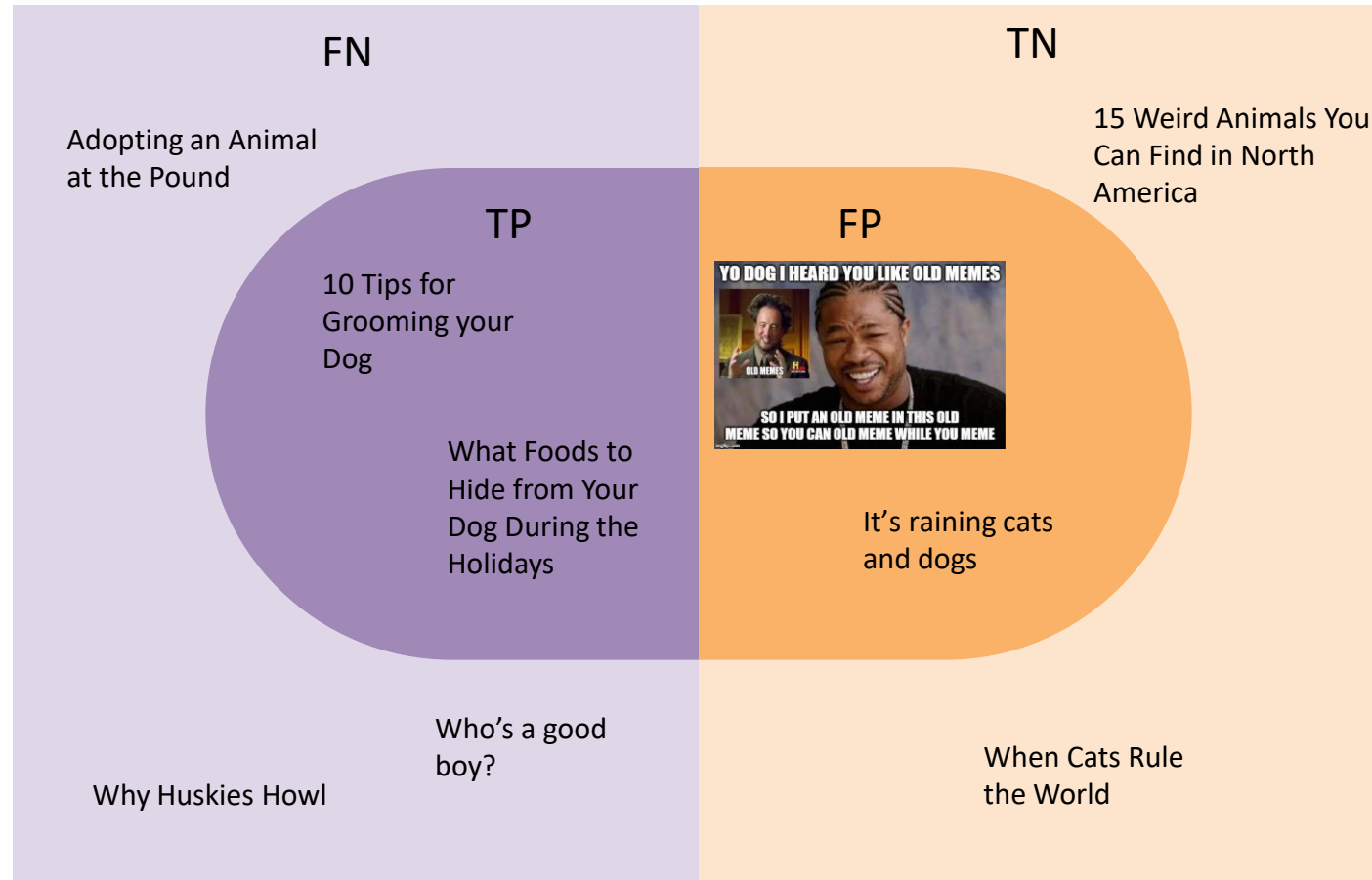
✓ True

87%

False

13%

# Why would you want to divide up your data (instead of training on it all)?

It wouldn't give it any examples to test on

**Possibly** 13%

It would overfit

✓ 52%

It might not be able to generalize to new examples

✓ 35%

# Contingency Table (out of table form)

Query:
Articles about dogs

**FN**

Adopting an Animal at the Pound

**TN**

15 Weird Animals You Can Find in North America

**TP**

10 Tips for Grooming your Dog

**FP**



What Foods to Hide from Your Dog During the Holidays

It's raining cats and dogs

Who's a good boy?

Why Huskies Howl

When Cats Rule the World

# Review: Steps



DO NOT ITERATE ON THE TESTING SET!!!

**Training**

Training Data

*perro*

*pato*

*...*

Training Labels

*dog*

*duck*

*...*

Word Features → Training → Learned model

Dev Set → Evaluate

**Testing**

Testing Data

*gato*

Image Features → Learned model → Prediction

# Review: Types of models



Classification

Regression

# Review: Classification Evaluation: the 2-by-2 contingency table

| *What label does our system predict? (↓)* | *What is the actual label?* | |
|---|---|---|
| | Actual Target Class ("●") | Not Target Class ("○") |
| **Selected/ Guessed ("●")** | True Positive (TP) ● *Actual* ● *Guessed* | False Positive (FP) ○ *Actual* ● *Guessed* |
| **Not selected/ not guessed ("○")** | False Negative (FN) ● *Actual* ○ *Guessed* | True Negative (TN) ○ *Actual* ○ *Guessed* |

# Precision and Recall Present a Tradeoff

1

precision

0

0                                    recall                                    1

# Precision and Recall Present a Tradeoff

precision

1

0

recall

0          1



Q: Where do you want your ideal model ?

Q: You have a model that always identifies correct instances. Where on this graph is it?

# Precision and Recall Present a Tradeoff



precision

recall

0, 1

# Precision and Recall Present a Tradeoff

# Precision and Recall Present a Tradeoff

precision

recall

1

0

0

1

For a given trained model, vary (certain) hyperparameters to adjust when your model makes a prediction

Q: Where do you want your ideal model ?

Q: You have a model that always identifies correct instances. Where on this graph is it?

Q: You have a model that only make correct predictions. Where on this graph is it?

Idea: measure the tradeoff between precision and recall

# Precision and Recall Present a Tradeoff



precision

Improve overall model: push the curve that way

0

0

recall

1

1

Q: Where do you want your ideal  model  ?

Q: You have a  model  that always identifies correct instances. Where on this graph is it?

Q: You have a  model  that only make correct predictions. Where on this graph is it?

Idea: measure the tradeoff between precision and recall

# Measure this Tradeoff:
# Area Under the Curve (AUC)

AUC measures the area under this tradeoff curve

Improve overall model: push the curve that way

precision

recall

0  1

0  1

Min AUC: 0 ☹

Max AUC: 1 😃

# Measure this Tradeoff:
# Area Under the Curve (AUC)



**AUC measures the area under this tradeoff curve**

1. **Computing the curve**

   You need true labels & predicted labels with some score/confidence estimate

   Threshold the scores and for each threshold compute precision and recall

Min AUC: 0 😞
Max AUC: 1 😀

# Measure this Tradeoff: Area Under the Curve (AUC)



1 precision recall 0

0 1

Improve overall model: push the curve that way

Min AUC: 0 😟
Max AUC: 1 😃

AUC measures the area under this tradeoff curve

1. Computing the curve

   You need true labels & predicted labels with some score/confidence estimate

   Threshold the scores and for each threshold compute precision and recall

2. Finding the area

   How to implement: trapezoidal rule (& others)

**In practice**: external library like the sklearn.metrics module

# A combined measure: F-score

Weighted (harmonic) average of **P**recision & **R**ecall

F1 measure: equal weighting between precision and recall

$$F_1 = \frac{2 * P * R}{P + R}$$

# A combined measure: F-score

Weighted (harmonic) average of **P**recision & **R**ecall

F1 measure: equal weighting between precision and recall

$$F_1 = \frac{2 * P * R}{P + R} = \frac{2 * TP}{2 * TP + FP + FN}$$

(useful when $P = R = 0$)

# Classification Evaluation: Accuracy, Precision, and Recall

**Accuracy**: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

$$F_1 = \frac{2*P*R}{P+R} = \frac{2*TP}{2*TP+FP+FN}$$

When would you want to use accuracy vs F1?

Accuracy works better if the dataset is <u>balanced</u>

Accuracy takes everything in consideration

F-Score is focused on TP

|  | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# Implementation: How To

1. scikit-learn: sklearn.metrics
   - very stable


2. huggingface evaluate module
   - community input
   - sometimes are based on sklearn


3. implement your own

# P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

*If we have more than one class, how do we combine multiple performance measures into one quantity?*

**Macroaveraging**: Compute performance for each class, then average.

**Microaveraging**: Collect decisions for all classes, compute contingency table, evaluate.

# P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

**Macroaveraging**: Compute performance for each class, then average.

$$\text{macroprecision} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c} = \frac{1}{C}\sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c} = \frac{1}{C}\sum_c \text{recall}_c$$

**Microaveraging**: Collect decisions for all classes, compute contingency table, evaluate.

$$\text{microprecision} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FP}_c} \qquad \text{microrecall} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FN}_c}$$

# Macro/Micro Example

# Macro-Average

Predicted "A"   Predicted "B"   Predicted "C"   Predicted "D"

**Class A**

Recall: **87%.**
Precision: 72%.

**Class B**

Recall: **33%.**
Precision: **20%.**

**Class C**

Recall: **90%.**
Precision: **90%.**

**Class D**

Recall: **93%.**
Precision: **100%.**

True "A"   True "B"   True "C"   True "D"

**Macro-average**

Recall = (0.87 + 0.33 + 0.9 + 0.93)/4 = **0.76**
Precision = (0.72+0.2+0.9+1)/4=**0.71**

https://www.evidentlyai.com/classification-metrics/multi-class-metrics

Each *instance* has equal weight

# Micro-Average

**All true positives**

All false negatives

All false positives

| Total TP | Total FP | Total FN |
|---|---|---|

$$\text{Precision}_{\text{Micro-average}} = \frac{13 + 1 + 9 + 14}{13 + 1 + 9 + 14 + (2 + 5 + 1 + 0)} = 0.82$$

$$\text{Recall}_{\text{Micro-average}} = \frac{13 + 1 + 9 + 14}{13 + 1 + 9 + 14 + (2 + 4 + 1 + 1)} = 0.82$$

https://www.evidentlyai.com/classification-metrics/multi-class-metrics

# Micro- vs Macro-Average

So when would we want to prefer micro-averaging vs macro-averaging?

$$\text{macroprecision} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c} = \frac{1}{C}\sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c} = \frac{1}{C}\sum_c \text{recall}_c$$

$$\text{microprecision} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FP}_c} \qquad \text{microrecall} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FN}_c}$$

# But how do we compute stats for multiple classes?

Either:

1. Compute "one-vs-all" 2x2 tables. OR

2. Generalize the 2x2 tables and compute per-class TP / FP / FN based on the diagonals and off-diagonals

# 1. Compute "one-vs-all" 2x2 tables

Predicted

Actual

| Look for ⬤ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

| Look for ◯ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

| Look for ▭ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

# 1. Compute "one-vs-all" 2x2 tables

Predicted

Actual

| Look for ● | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 2 | 1 |
| Not select/not guessed | 2 | 4 |

| Look for ○ | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 2 | 1 |
| Not select/not guessed | 1 | 5 |

| Look for ▭ | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 1 | 2 |
| Not select/not guessed | 1 | 5 |

ML EVALUATION

# 2. Generalizing the 2-by-2 contingency table

| | | Correct Value | | |
|---|---|---|---|---|
| | | 🟠 | ◯ | ▭ |
| **Guessed Value** | 🟠 | # | # | # |
| | ◯ | # | # | # |
| | ▭ | # | # | # |

This is also called a **Confusion Matrix**

# 2. Generalizing the 2-by-2 contingency table

Predicted    ● ○ ▭ ● ○ ▭ ● ○ ▭

Actual    ● ○ ○ ▭ ○ ▭ ● ● ●

| | | Correct Value | | |
|---|---|---|---|---|
| | | ● | ○ | ▭ |
| **Guessed Value** | ● | # | # | # |
| | ○ | # | # | # |
| | ▭ | # | # | # |

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

|  | | Correct Value | | |
|---|---|---|---|---|
| | | ● | ○ | ▭ |
| **Guessed Value** | ● | 2 | 0 | 1 |
| | ○ | 1 | 2 | 0 |
| | ▭ | 1 | 1 | 1 |

# 2. Generalizing the 2-by-2 contingency table

Predicted ● ○ ▭ ● ○ ▭ ● ○ ▭

Actual ● ○ ○ ▭ ○ ▭ ● ● ●

|  | **Correct Value** | | |
|---|---|---|---|
|  | ● | ○ | ▭ |
| **Guessed Value** ● | A 2 | B 0 | C 1 |
| ○ | D 1 | E 2 | F 0 |
| ▭ | G 1 | H 1 | I 1 |

How do you compute $TP$ ● ?

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

|  | | Correct Value | | |
|---|---|---|---|---|
|  | | ● | ○ | ▭ |
| **Guessed Value** | ● | A 2 | B 0 | C 1 |
|  | ○ | D 1 | E 2 | F 0 |
|  | ▭ | G 1 | H 1 | I 1 |

How do you compute $TP$ ● ?

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

| Correct Value | | |
|---|---|---|
| ● | ○ | ▭ |
| A 2 | B 0 | C 1 |
| D 1 | E 2 | F 0 |
| G 1 | H 1 | I 1 |

Guessed Value

How do you compute $FN_●$?

# 2. Generalizing the 2-by-2 contingency table

Predicted ● ○ ▭ ● ○ ▭ ● ○ ▭

Actual ● ○ ○ ▭ ○ ▭ ● ● ●

| | Correct Value | | |
|---|---|---|---|
| | ● | ○ | ▭ |
| **Guessed Value** ● | A 2 | B 0 | C 1 |
| ○ | D 1 | E 2 | F 0 |
| ▭ | G 1 | H 1 | I 1 |

How do you compute $FN_●$?

# 2. Generalizing the 2-by-2 contingency table

Predicted ● ○ ▭ ● ○ ▭ ● ○ ▭

Actual ● ○ ○ ▭ ○ ▭ ● ● ●

| | | Correct Value | | |
|---|---|---|---|---|
| | | ● | ○ | ▭ |
| **Guessed Value** | ● | A 2 | B 0 | C 1 |
| | ○ | D 1 | E 2 | F 0 |
| | ▭ | G 1 | H 1 | I 1 |

How do you compute $FP_{▭}$?

# 2. Generalizing the 2-by-2 contingency table

Predicted

Actual

| | | Correct Value | | |
|---|---|---|---|---|
| | | ● (circle) | ○ (ring) | ▭ (rounded rect) |
| **Guessed Value** | ● | A 2 | B 0 | C 1 |
| | ○ | D 1 | E 2 | F 0 |
| | ▭ | G 1 | H 1 | I 1 |

How do you compute $FP_{▭}$?

# Generalizing the 2-by-2 contingency table

| | | Correct Value | | |
|---|---|---|---|---|
| | | 🟠 (filled circle) | ⭕ (open circle) | 🔲 (rounded rectangle) |
| **Guessed Value** | 🟠 (filled circle) | 80 | 9 | 11 |
| | ⭕ (open circle) | 7 | 86 | 7 |
| | 🔲 (rounded rectangle) | 2 | 8 | 9 |

Q: Is this a good result?

# Generalizing the 2-by-2 contingency table

| | | Correct Value | | |
|---|---|---|---|---|
| **Q: Is this a good result?** | | 🟠 | ◯ | ▭ |
| **Guessed Value** | 🟠 | 30 | 40 | 30 |
| | ◯ | 25 | 30 | 50 |
| | ▭ | 30 | 35 | 35 |

# Generalizing the 2-by-2 contingency table

| | | Correct Value | | |
|---|---|---|---|---|
| | | ● | ○ | ▭ |
| **Guessed Value** | ● | 7 | 3 | 90 |
| | ○ | 4 | 8 | 88 |
| | ▭ | 3 | 7 | 90 |

Q: Is this a good result?