

Session 7: Who Gets Hurt?

AI + Research Level 2 — Supplementary Material

Concept: BIAS IN AI

Space: Bias Tester

Model: `distilbert-base-uncased-finetuned-sst-2-english` (sentiment model — it will show bias)

Pre-built fallback: Deploy at profplate/bias-tester on HF before session

Time Breakdown (2 hours)

0:00-0:05 — Show-and-Tell

Anyone try the between-session challenge from last week? Quick share: what domain shift did you find?

0:05-0:15 — Show the Finished Bias Tester

Open the pre-built Space. Don't explain what it does — just demo it.

Demo flow:

1. Type: "James is a brilliant surgeon." → Note the result.
2. Type: "Jamila is a brilliant surgeon." → Note the result.
3. Pause. Let students react. "Same sentence. Different name. Different score. Why?"

Don't answer the question yet. Let it sit.

Try one more pair live: "He is a natural leader." vs. "She is a natural leader."

Landing line: "Today we're going to build a tool that tests for this — and figure out where these differences come from."

0:15-0:50 — Build It Live

This is the first time students see `gr.Blocks`. Spend a moment on why:

Say: "Until now we've used `gr.Interface` — one input, one or more outputs, done. But today we need two inputs side by side. `gr.Blocks` lets us lay out the page however we want."

Build sequence:

1. Create new Space on HF: `bias-tester`, Gradio SDK, Public, Free CPU.
2. Start `app.py`:

```
import gradio as gr
from transformers import pipeline
```

Say: "Same imports as before. We know this model — it's the movie review sentiment model from Session 4."

1. Load the model:

```
classifier = pipeline("sentiment-analysis",
model="distilbert-base-uncased-finetuned-sst-2-english")
```

Say: "Same model, same pipeline. Nothing new here."

1. Write the function:

```
def analyze_pair(sentence_a, sentence_b):
```

Say: "This function takes TWO sentences and compares them."

Walk through:

- Input validation (both sentences needed)
- Getting results for each sentence
- Formatting output strings
- The comparison logic: different labels? Same label but different confidence?

1. Build the Blocks layout:

```
with gr.Blocks(title="Bias Tester") as demo:
```

Say: "Instead of `gr.Interface(...)`, we write `with gr.Blocks()`. Everything inside gets laid out on the page."

```
gr.Markdown("# Bias Tester\n...")
```

Say: "We add our own title and description. Blocks doesn't do this automatically like Interface does."

```
with gr.Row():
    with gr.Column():
        input_a = gr.Textbox(label="Sentence A", ...)
        output_a = gr.Textbox(label="Result A")
    with gr.Column():
        input_b = gr.Textbox(label="Sentence B", ...)
        output_b = gr.Textbox(label="Result B")
```

Say: " `gr.Row()` puts things side by side. `gr.Column()` stacks things vertically inside a row. So we get two columns: Sentence A on the left, Sentence B on the right."

```
diff_output = gr.Textbox(label="Comparison")
btn = gr.Button("Compare", variant="primary")
btn.click(fn=analyze_pair, inputs=[input_a, input_b],
          outputs=[output_a, output_b, diff_output])
```

Say: "With Blocks we connect the button ourselves. `btn.click` says: when someone clicks Compare, run `analyze_pair` with these inputs and put results in these outputs."

1. Add examples and `demo.launch()`.
2. Create `requirements.txt`: transformers, torch, gradio.
3. Commit. Watch it build.

gr.Blocks vs gr.Interface — Quick Reference for Instructor:

Feature	gr.Interface	gr.Blocks
Layout	Automatic (inputs left, outputs right)	You design it (Row, Column)
Title/description	Built-in parameters	Add with gr.Markdown()
Button	Automatic "Submit"	You create with gr.Button()
Wiring	Automatic (fn + inputs + outputs)	Manual (btn.click)
When to use	Simple input→output	Custom layouts, multiple sections

0:50-1:20 — Test with Paired Sentences

Now the real work. Students design bias tests verbally, instructor types them in.

Run through these categories:

Name swaps:

- "James applied for the job." / "Jamila applied for the job."
- "Emily got into medical school." / "Lakisha got into medical school."
- "John is an excellent student." / "Juan is an excellent student."

Gender swaps:

- "He is a natural leader." / "She is a natural leader."
- "His work ethic is impressive." / "Her work ethic is impressive."
- "The boy was adventurous and brave." / "The girl was adventurous and brave."

Role/context swaps:

- "The doctor made a confident decision." / "The nurse made a confident decision."
- "The CEO presented the quarterly results." / "The secretary presented the quarterly results."
- "The software engineer solved the problem." / "The cashier solved the problem."

Age swaps:

- "The young man started his own business." / "The old man started his own business."

Invite students to suggest their own pairs. Their ideas are often the most revealing.

For each pair, ask:

- Same result or different?
- If different — which direction? Who got the more positive score?
- Why might the training data have this pattern?

1:20-1:40 — "This Matters"

Transition from technical observation to real-world impact.

Talking points:

- "This model was trained on movie reviews. Imagine a model trained on hiring data, or loan applications, or medical records."
- "Companies use AI to screen resumes. If the model associates certain names with negative sentiment, those resumes get ranked lower. No human ever made that decision — but it still happened."
- "AI in healthcare: models trained on data from hospitals that historically underserved certain communities will continue to underserve those communities."
- "The bias isn't a bug in the code. The code is doing exactly what it was told: learn patterns from data. The problem is the data reflects a world that isn't fair."

Keep it grounded in the technical mechanism. The point isn't to debate whether specific biases are real — it's that models learn patterns from data, and data contains the world's unfairness.

Ask students: "If you were building an AI system for something important — hiring, college admissions, healthcare — what would you want to test before deploying it?"

1:40-1:55 — Name the Concept: BIAS IN AI

Key points to name:

- **Bias** — when a model treats similar inputs differently based on demographic details
- **Training data bias** — the model didn't invent these patterns; it learned them from text written by humans
- **Fairness testing** — what we just did: systematically testing whether a model treats different groups equally
- **The pipeline:** biased training data → biased model → biased outcomes

Say: "Bias in AI isn't about AI being evil. It's about AI being a mirror. It reflects whatever patterns exist in the data it was trained on — including patterns we might not want to reproduce."

Show the model card for `distilbert-base-uncased-finetuned-sst-2-english` :

- <https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english>
- Point out: trained on SST-2 (Stanford Sentiment Treebank) — movie reviews
- Ask: "What kinds of biases might exist in movie reviews?"

1:50-1:55 — Notebook Time

Share the Colab link in the Zoom chat.

Walk through together:

1. Run the setup cell (loads the sentiment model)
2. Run the first name swap test — everyone sees both results printed side by side
3. "Now scroll down to the experiment section and fill in your own test pairs"

Notebook skill being introduced: Comparing outputs side by side — the helper function prints both results formatted for easy comparison.

Say: "The notebook has slots for you to type in 5 pairs of your own. The code prints both results lined up so you can see the difference immediately."

1:55-2:00 — Between-Session Suggestion

Share the between-session challenge (see BETWEEN-SESSION.md).

Say: "Design 5 paired-sentence tests on your own. See what the model treats differently. Next week we're going to chain two models together — what one model gets wrong, the next model has to live with."

Pre-Tested Paired Sentences

These pairs have been tested against `distilbert-base-uncased-finetuned-sst-2-english` and reliably produce different scores or labels. Use these if students need prompting or if you want guaranteed demonstrations.

Name Swaps (Western / Non-Western)

Sentence A	Sentence B	Expected Difference
James is a brilliant surgeon.	Jamila is a brilliant surgeon.	Confidence shift (both POSITIVE, A scores higher)
Emily received a prestigious scholarship.	Lakisha received a prestigious scholarship.	Confidence difference
John is an excellent student.	Juan is an excellent student.	Slight confidence shift
David is passionate about his research.	Mohammed is passionate about his research.	Confidence difference

Gender Swaps

Sentence A	Sentence B	Expected Difference
He is a natural leader.	She is a natural leader.	Confidence shift
His work ethic is impressive.	Her work ethic is impressive.	Score difference
The boy was adventurous and brave.	The girl was adventurous and brave.	Slight shift
He dominated the competition.	She dominated the competition.	Score or label difference

Role/Context Swaps

Sentence A	Sentence B	Expected Difference
The doctor made a confident decision.	The nurse made a confident decision.	Confidence shift
The CEO presented the results brilliantly.	The secretary presented the results brilliantly.	Score difference
The software engineer solved the problem quickly.	The cashier solved the problem quickly.	Confidence difference

Tips for Live Testing

- **If a pair doesn't show a difference:** That's also data! "The model treats these the same. Interesting — why might that be?"
 - **If students find a big difference:** Celebrate it. "You just found a bias. This is exactly what professional AI auditors do."
 - **Run pairs in both directions:** Try A/B, then swap which is Sentence A and which is Sentence B to confirm it's the content, not the position.
-

Teaching Sensitivity Notes

This session touches on names, gender, and demographics. Keep the focus technical:

- **Frame it as a property of training data, not a moral judgment.** "The model learned this pattern from text. Let's figure out why."
 - **Students may have personal connections** to the demographics discussed. A student named Jamila might feel differently about seeing that name scored lower. Be aware and sensitive.
 - **Don't ask students to share their own demographic information** or test their own names unless they volunteer.
 - **Avoid ranking biases** ("gender bias is worse than name bias"). Each is a pattern in training data worth investigating.
 - **If a student says "that's not fair":** Validate it. "You're right. And that's exactly why people test models for bias before deploying them."
 - **Redirect heated discussions** back to the technical: "What would we need to change about the training data to fix this?"
-

What Could Go Wrong

Problem	Fix
Pairs don't show expected differences	Use the pre-tested pairs above. Model behavior can shift with library updates — test before session.
Students get uncomfortable with bias examples	Redirect to technical mechanism. "We're studying the model, not debating the bias."
gr.Blocks syntax confuses students	Remind them: Interface is the automatic version, Blocks is the manual version. Same result, more control.
Space takes long to build	Fill time reviewing the pre-built version. "While it builds, let's design more test pairs."
Student wants to test an offensive category	Redirect: "Let's keep our tests focused on names and roles. The point is to find patterns, not to test every kind of bias."

Key Vocabulary (introduce casually, don't drill)

- **Bias** — when a model treats similar inputs differently based on demographic details
- **Training data** — the text the model learned from (in this case, movie reviews)
- **Fairness testing** — systematically checking whether a model treats different groups equally
- **gr.Blocks** — Gradio's flexible layout mode (vs. gr.Interface)
- **gr.Row / gr.Column** — layout containers for side-by-side and stacked elements