# Feature Overlap/ Feature Interaction

How overlapping features work in maxent models
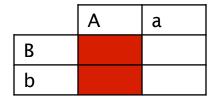
# Feature Overlap

- Maxent models handle overlapping features well.
- Unlike a NB model, there is no double counting!

**Empirical**

|   | A | a |
|---|---|---|
| B | 2 | 1 |
| b | 2 | 1 |

|   | A | a |
|---|---|---|
| B | | |
| b | | |

|   | A | a |
|---|---|---|
| B | | |
| b | | |

|   | A | a |
|---|---|---|
| B | | |
| b | | |

All = 1

|   | A | a |
|---|---|---|
| B | 1/4 | 1/4 |
| b | 1/4 | 1/4 |

A = 2/3

|   | A | a |
|---|---|---|
| B | 1/3 | 1/6 |
| b | 1/3 | 1/6 |

A = 2/3

|   | A | a |
|---|---|---|
| B | 1/3 | 1/6 |
| b | 1/3 | 1/6 |

|   | A | a |
|---|---|---|
| B | | |
| b | | |

|   | A | a |
|---|---|---|
| B | $\lambda_A$ | |
| b | $\lambda_A$ | |

|   | A | a |
|---|---|---|
| B | $\lambda'_A + \lambda''_A$ | |
| b | $\lambda'_A + \lambda''_A$ | |

# Example: Named Entity Feature Overlap

Grace is correlated with PERSON, but does not add much evidence on top of already knowing prefix features.

## Local Context

|       | Prev  | Cur   | Next  |
|-------|-------|-------|-------|
| State | Other | ???   | ???   |
| Word  | at    | Grace | Road  |
| Tag   | IN    | NNP   | NNP   |
| Sig   | x     | Xx    | Xx    |

## Feature Weights

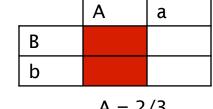| Feature Type         | Feature  | PERS  | LOC   |
|----------------------|----------|-------|-------|
| Previous word        | at       | -0.73 | 0.94  |
| Current word         | Grace    | 0.03  | 0.00  |
| Beginning bigram     | <G       | 0.45  | -0.04 |
| Current POS tag      | NNP      | 0.47  | 0.45  |
| Prev and cur tags    | IN NNP   | -0.10 | 0.14  |
| Previous state       | Other    | -0.70 | -0.92 |
| Current signature    | Xx       | 0.80  | 0.46  |
| Prev state, cur sig  | O-Xx     | 0.68  | 0.37  |
| Prev-cur-next sig     | x-Xx-Xx  | -0.69 | 0.37  |
| P. state - p-cur sig | O-x-Xx   | -0.20 | 0.82  |
| …                    |          |       |       |
| Total:               |          | -0.58 | 2.68  |

# Feature Interaction

- Maxent models handle overlapping features well, but do not automatically model feature interactions.

**Empirical**

|   | A | a |
|---|---|---|
| B | 1 | 1 |
| b | 1 | 0 |

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

**All = 1**

|   | A | a |
|---|---|---|
| B | 1/4 | 1/4 |
| b | 1/4 | 1/4 |

**A = 2/3**

|   | A | a |
|---|---|---|
| B | 1/3 | 1/6 |
| b | 1/3 | 1/6 |

**B = 2/3**

|   | A | a |
|---|---|---|
| B | 4/9 | 2/9 |
| b | 2/9 | 1/9 |

|   | A | a |
|---|---|---|
| B | 0 | 0 |
| b | 0 | 0 |

|   | A | a |
|---|---|---|
| B | $\lambda_A$ |   |
| b | $\lambda_A$ |   |

|   | A | a |
|---|---|---|
| B | $\lambda_A + \lambda_B$ | $\lambda_B$ |
| b | $\lambda_A$ |   |

# Feature Interaction

- If you want interaction terms, you have to add them:

Empirical

|   | A | a |
|---|---|---|
| B | 1 | 1 |
| b | 1 | 0 |

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

A = 2/3

|   | A | a |
|---|---|---|
| B | 1/3 | 1/6 |
| b | 1/3 | 1/6 |

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

B = 2/3

|   | A | a |
|---|---|---|
| B | 4/9 | 2/9 |
| b | 2/9 | 1/9 |

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

AB = 1/3

|   | A | a |
|---|---|---|
| B | 1/3 | 1/3 |
| b | 1/3 | 0 |

- A disjunctive feature would also have done it (alone):

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

|   | A | a |
|---|---|---|
| B | 1/3 | 1/3 |
| b | 1/3 | 0 |

# Quiz Question

- Suppose we have a 1 feature maxent model built over observed data as shown.
- What is the constructed model's probability distribution over the four possible outcomes?

Empirical

|   | A | a |
|---|---|---|
| B | 2 | 1 |
| b | 2 | 1 |

Features

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

Expectations

Probabilities

|   | A | a |
|---|---|---|
| B |   |   |
| b |   |   |

# Feature Interaction

- For loglinear/logistic regression models in statistics, it is standard to do a greedy stepwise search over the space of all possible interaction terms.

- This combinatorial space is exponential in size, but that's okay as most statistics models only have 4–8 features.

- In NLP, our models commonly use hundreds of thousands of features, so that's not okay.

- Commonly, interaction terms are added by hand based on linguistic intuitions.

# Example: NER Interaction

Previous-state and current-signature have interactions, e.g. P=PERS-C=Xx indicates C=PERS much more strongly than C=Xx and P=PERS independently.

This feature type allows the model to capture this interaction.

## Feature Weights

| Feature Type | Feature | PERS | LOC |
|---|---|---|---|
| Previous word | *at* | -0.73 | 0.94 |
| Current word | *Grace* | 0.03 | 0.00 |
| Beginning bigram | *<G* | 0.45 | -0.04 |
| Current POS tag | NNP | 0.47 | 0.45 |
| Prev and cur tags | IN NNP | -0.10 | 0.14 |
| Previous state | Other | -0.70 | -0.92 |
| Current signature | Xx | 0.80 | 0.46 |
| Prev state, cur sig | O-Xx | 0.68 | 0.37 |
| Prev-cur-next sig | x-Xx-Xx | -0.69 | 0.37 |
| P. state - p-cur sig | O-x-Xx | -0.20 | 0.82 |
| … | | | |
| **Total:** | | **-0.58** | **2.68** |

## Local Context

| | Prev | Cur | Next |
|---|---|---|---|
| State | Other | ??? | ??? |
| Word | at | Grace | Road |
| Tag | IN | NNP | NNP |
| Sig | x | Xx | Xx |

# Feature Overlap/ Feature Interaction

How overlapping features work in maxent models