

[illegible]



Maximum Entropy Models

- An equivalent approach:
 - Lots of distributions out there, most of them very spiked, specific, overfit.
 - We want a distribution which is uniform except in specific ways we require.
 - Uniformity means **high entropy** – we can search for distributions which have properties we desire, but also have high entropy.

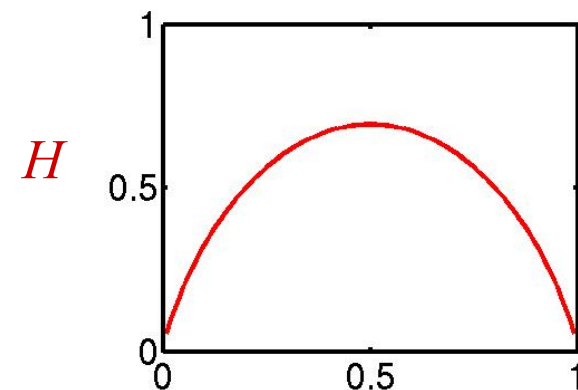
Ignorance is preferable to error and he is less remote from the truth who believes nothing than he who believes what is wrong – Thomas Jefferson (1781)



(Maximum) Entropy

- Entropy: the uncertainty of a distribution.
- Quantifying uncertainty (“surprise”):
 - Event x
 - Probability p_x
 - “Surprise” $\log(1/p_x)$
- Entropy: expected surprise (over p):

$$H(p) = E_p \left[\log_2 \frac{1}{p_x} \right] = - \sum_x p_x \log_2 p_x$$



A coin-flip is most uncertain for a fair coin.

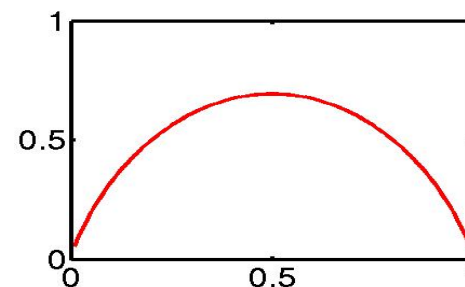


Maxent Examples I

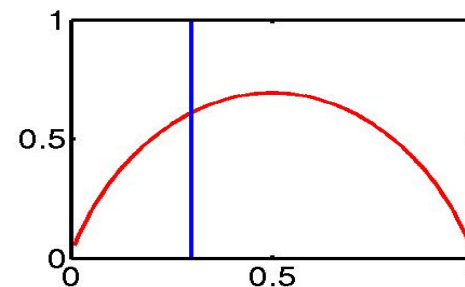
- What do we want from a distribution?
 - Minimize commitment = maximize entropy.
 - Resemble some reference distribution (data).
- Solution: maximize entropy H , subject to feature-based constraints:

$$E_p[f_i] = E_{\hat{p}}[f_i] \iff \sum_{x \in f_i} p_x = C_i$$

- Adding constraints (features):
 - Lowers maximum entropy
 - Raises maximum likelihood of data
 - Brings the distribution further from uniform
 - Brings the distribution closer to data



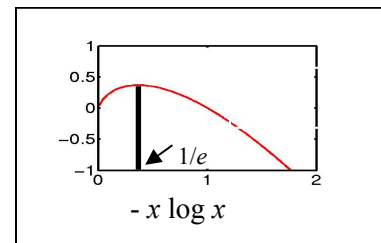
Unconstrained,
max at 0.5



Constraint that
 $p_{\text{HEADS}} = 0.3$



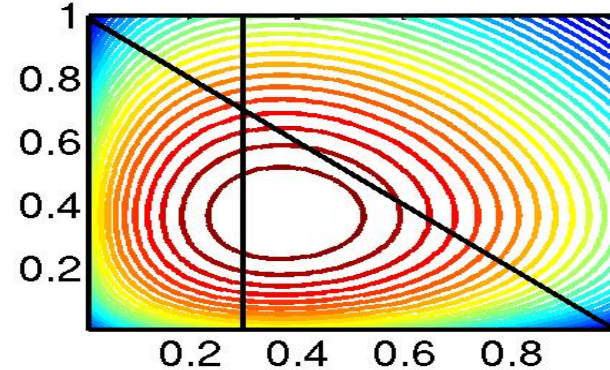
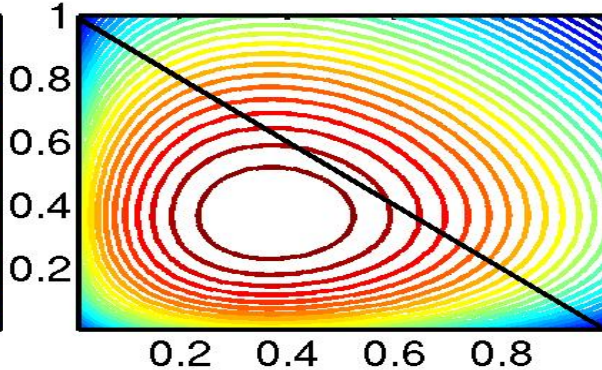
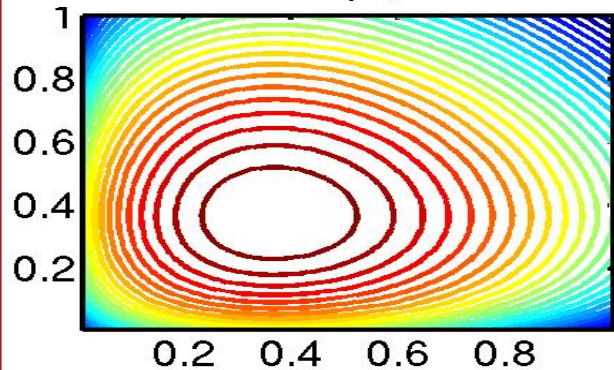
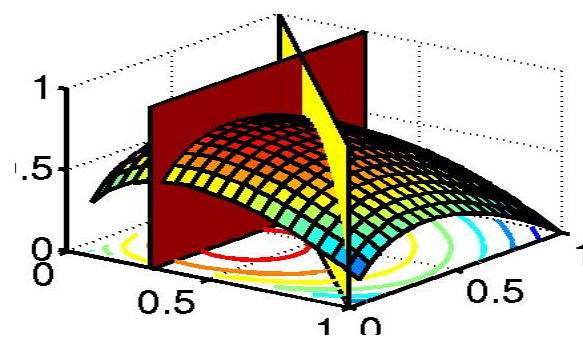
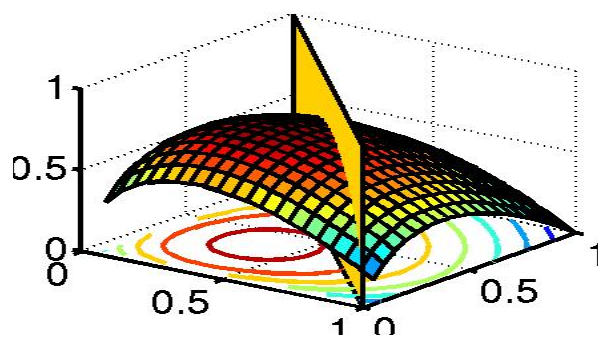
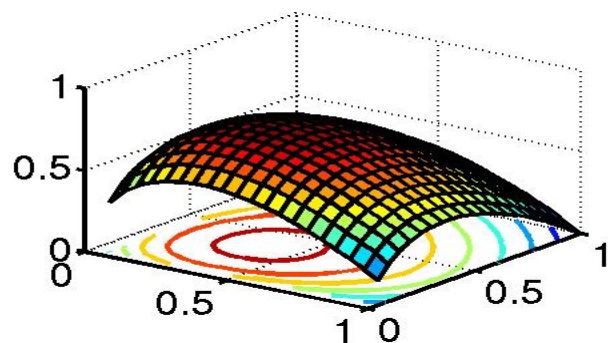
Maxent Examples II



$$H(p_H p_T)$$

$$p_H + p_T = 1$$

$$p_H = 0.3$$





Maxent Examples III

- Let's say we have the following event space:

NN	NNS	NNP	NNPS	VBZ	VBD
----	-----	-----	------	-----	-----

- ... and the following empirical data:

3	5	11	13	3	1
---	---	----	----	---	---

- Maximize H:

$1/e$	$1/e$	$1/e$	$1/e$	$1/e$	$1/e$
-------	-------	-------	-------	-------	-------

- ... want probabilities: $E[\text{NN}, \text{NNS}, \text{NNP}, \text{NNPS}, \text{VBZ}, \text{VBD}] = 1$

$1/6$	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$
-------	-------	-------	-------	-------	-------



Maxent Examples IV

- Too uniform!
- N^* are more common than V^* , so we add the feature $f_N = \{NN, NNS, NNP, NNPS\}$, with $E[f_N] = 32/36$

NN	NNS	NNP	NNPS	VBZ	VBD
8/36	8/36	8/36	8/36	2/36	2/36

- ... and proper nouns are more frequent than common nouns, so we add $f_p = \{NNP, NNPS\}$, with $E[f_p] = 24/36$

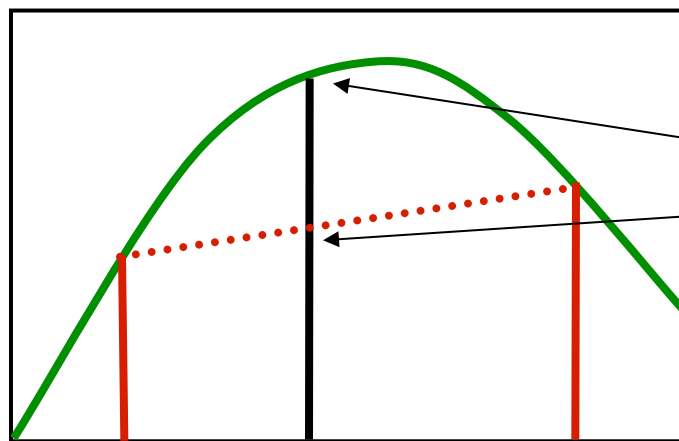
4/36	4/36	12/36	12/36	2/36	2/36
------	------	-------	-------	------	------

- ... we could keep refining the models, e.g., by adding a feature to distinguish singular vs. plural nouns, or verb types.

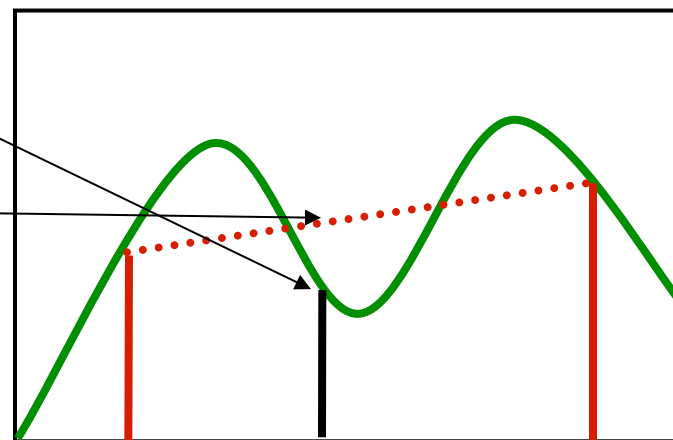


Convexity

$$f\left(\sum_i w_i x_i\right) \geq \sum_i w_i f(x_i) \quad \sum_i w_i = 1$$



Convex



Non-Convex

Convexity guarantees a single, global maximum because any higher points are greedily reachable.



Convexity II

- Constrained $H(p) = -\sum x \log x$ is convex:
 - $-x \log x$ is convex
 - $-\sum x \log x$ is convex (sum of convex functions is convex).
 - The feasible region of constrained H is a linear subspace (which is convex)
 - The constrained entropy surface is therefore convex.
- The maximum likelihood exponential model (dual) formulation is also convex.

