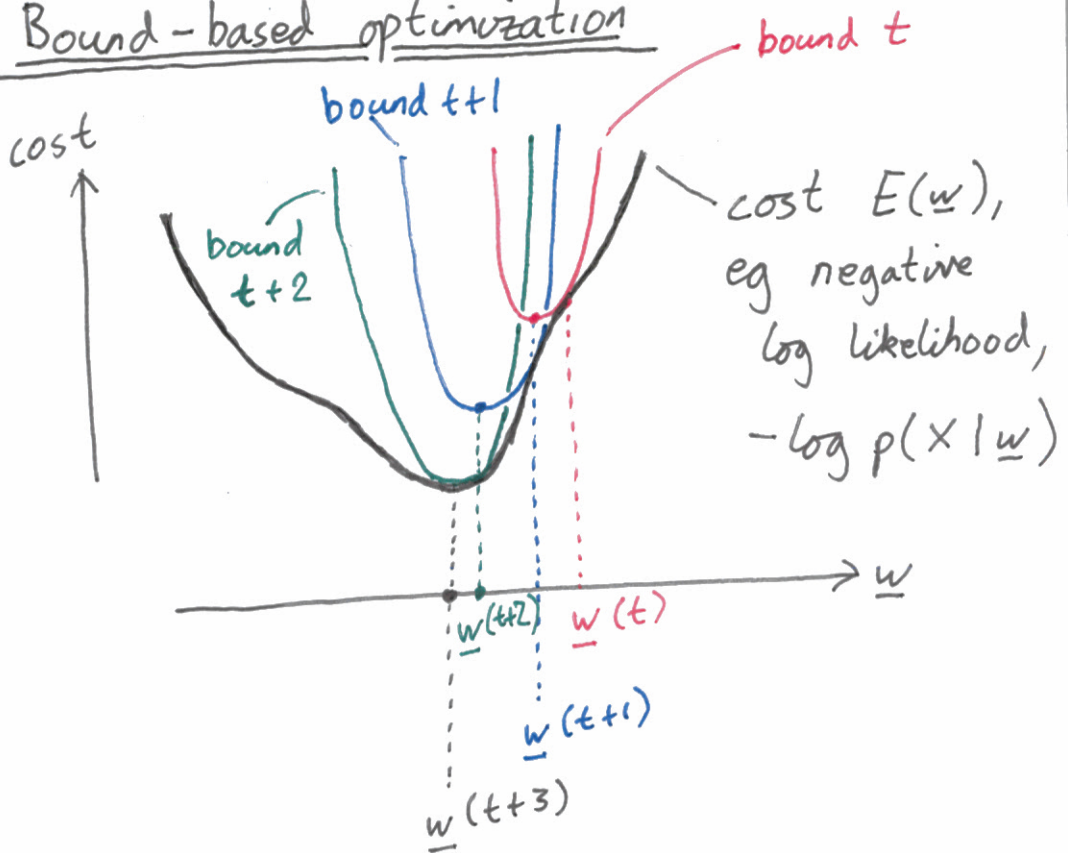


Bound-based optimization



For EM mix Gaussians

Bound minimized in closed form:

- 1) No step size / learning rate
- 2) Constraints on Σ , Π satisfied

Revision: Laplace Approximation

"Energy" $E(\underline{w})$

$$\text{Hessian, } H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j}$$



Approximate quadratic cost:

$$E(\underline{w}) \approx \frac{1}{2} (\underline{w} - \underline{w}^*)^T H (\underline{w} - \underline{w}^*) + \text{const.}$$

Newton's Method

- Initialize $\underline{w}^{(0)}$

$$- \underline{w}^{(t+1)} = \underline{w}^{(t)} - H^{-1} \underline{g}$$

eval at $\underline{w}^{(t)}$
 $\nabla_{\underline{w}} E(\underline{w})$

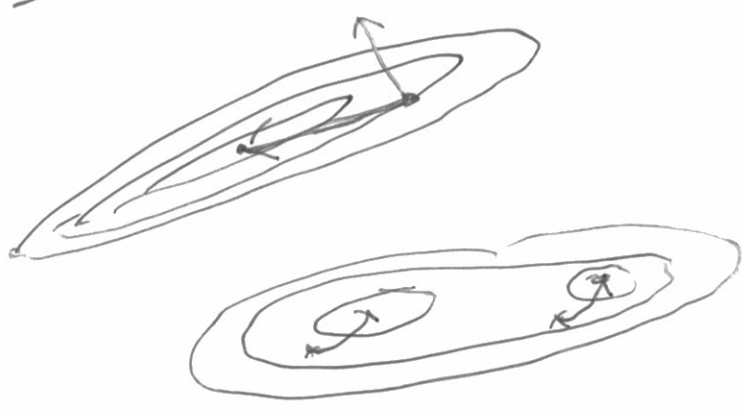
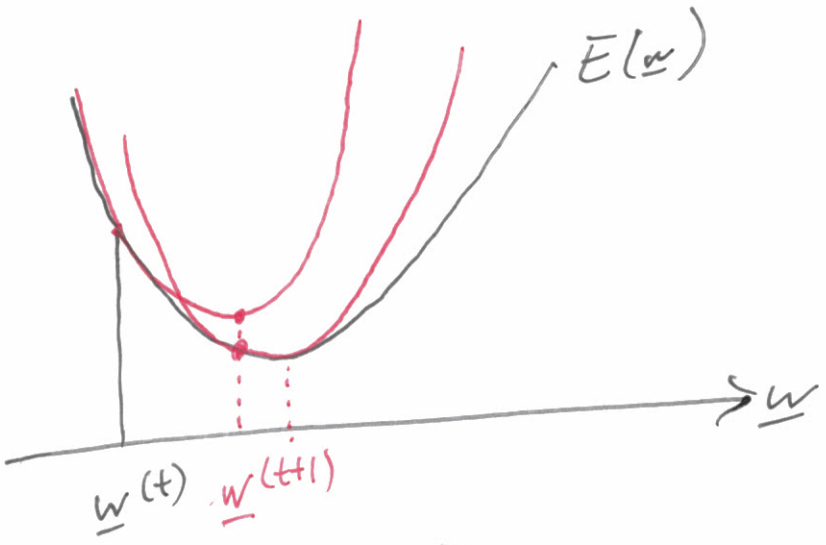
If cost is quadratic:

$$\underline{g} = H(\underline{w} - \underline{w}^*)$$

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - H^{-1} H(\underline{w}^{(t)} - \underline{w}^*)$$

Π if H invertible

$$= \underline{w}^*$$



Why use other optimizers?

- Convergence?
- Tuning?
- Constraints?
- SGD can't give "sparse" solutions
→ some $w_d = 0$

L1 Regularization

$$c(\underline{w}) = \underbrace{E(\underline{w})}_{\text{Training error}} + \lambda \underbrace{\sum_d |w_d|}_{\|\underline{w}\|_1}$$

The confection



m&m's
(185g)



Jelly Belly
(100g)



Chocolate Raisins
(200g)

Stuff Inf2b students wrote

Number M&Ms: ~~185~~ 204
 Number Jelly Belly: ~~146~~ 146
 Num. choc-raisin blobs: ~~87~~ 87

Number M&Ms: ~~185~~ 185
 Number Jelly Belly: ~~180~~ 180
 Num. choc-raisin blobs: ~~190~~ 190

Number M&Ms: ~~240~~ 240
 Number Jelly Belly: ~~150~~ 150
 Num. choc-raisin blobs: ~~130~~ 130

Number M&Ms: ~~247~~ 247
 Number Jelly Belly: ~~75~~ 75
 Num. choc-raisin blobs: ~~89~~ 89

Number M&Ms: ~~70~~ 70
 Number Jelly Belly: ~~83~~ 83
 Num. choc-raisin blobs: ~~100~~ 100

Number M&Ms: ~~150~~ 152 202 82
 Number Jelly Belly: ~~70~~ 72
 Num. choc-raisin blobs: ~~150~~ 132 102

Number M&Ms: ~~168~~ 168
 Number Jelly Belly: ~~98~~ 98
 Num. choc-raisin blobs: ~~139~~ 139

Number M&Ms: ~~84~~ 84
 Number Jelly Belly: ~~52~~ 52
 Num. choc-raisin blobs: ~~133~~ 133

F33) M3

Number M&Ms: 90
 Number Jelly Belly: 80
 Num. choc-raisin blobs: ~~80~~
 or more likely the average of all other guesses...
 Full name:
 (to award prize only)

Number M&Ms: 231.25
 Number Jelly Belly: 87.5
 Num. choc-raisin blobs: 133.34
 Full name: ANON
 (to award prize only)

$$\rho = 1 \frac{g}{cm^3}$$

. 5 cm³ each

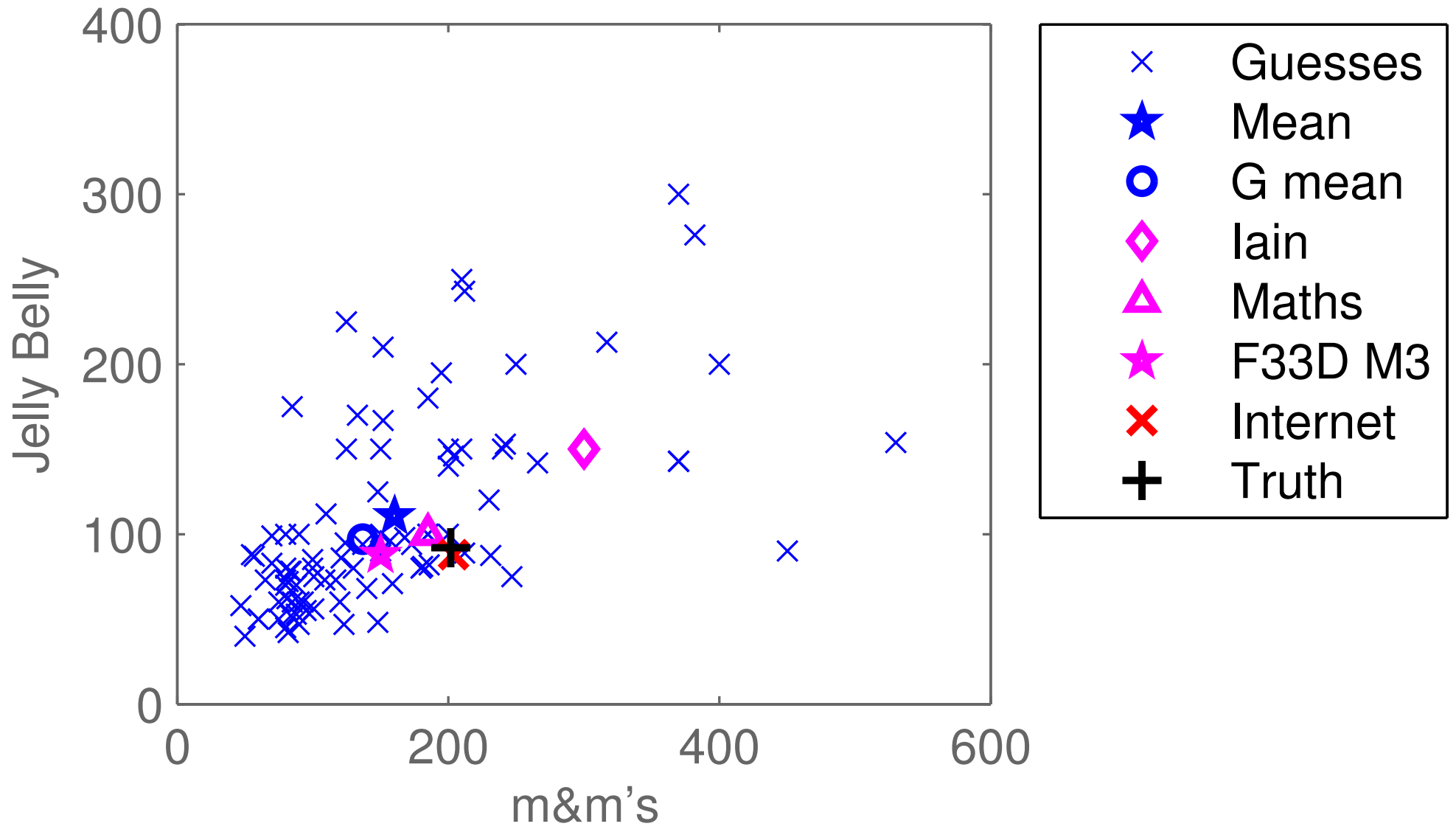
$$\rho = \frac{m}{V} \Rightarrow m = \rho V$$

$$\rho = 1.7 \frac{g}{cm^3}$$

$$m \frac{1}{cm^3} = \frac{1.7 \cdot 100}{1.5} = .5$$

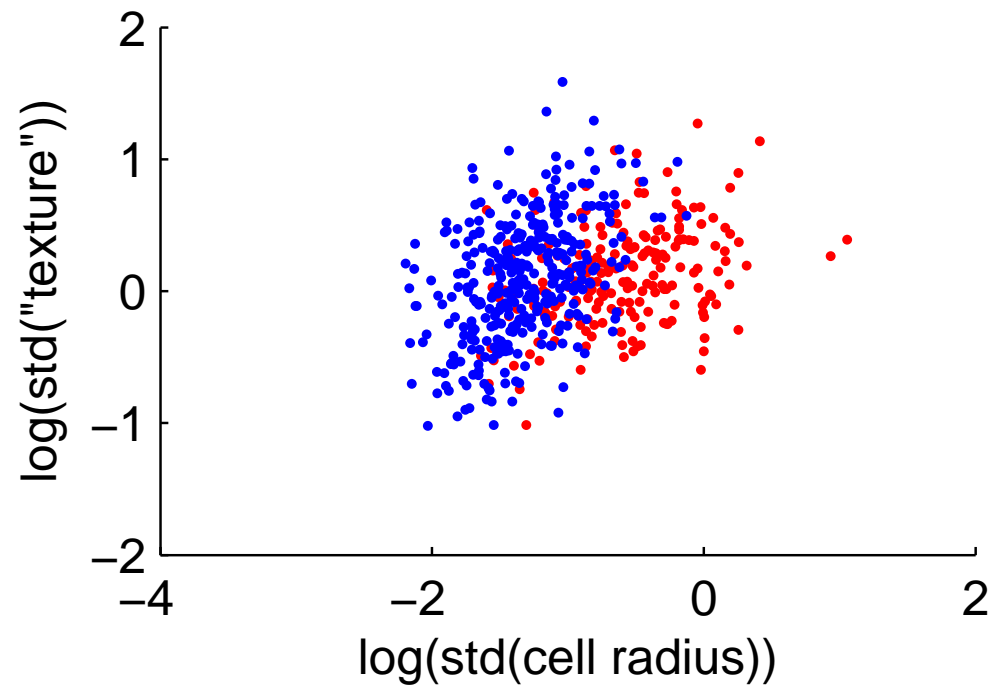
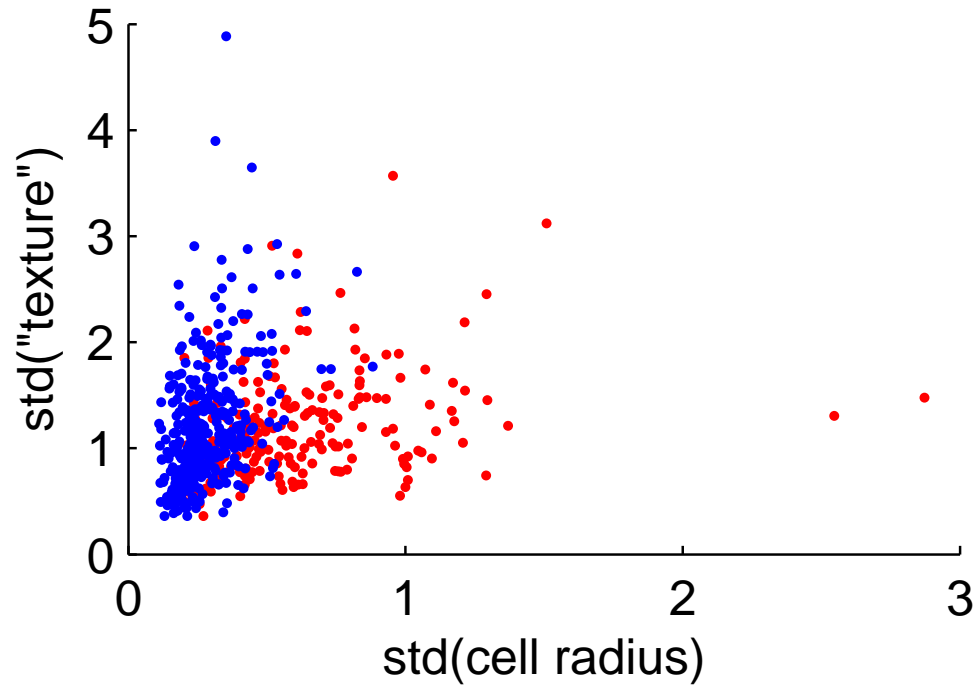
$$\frac{1.87}{-1}$$

A 2D space



For 3D and more, check out the code on the website.

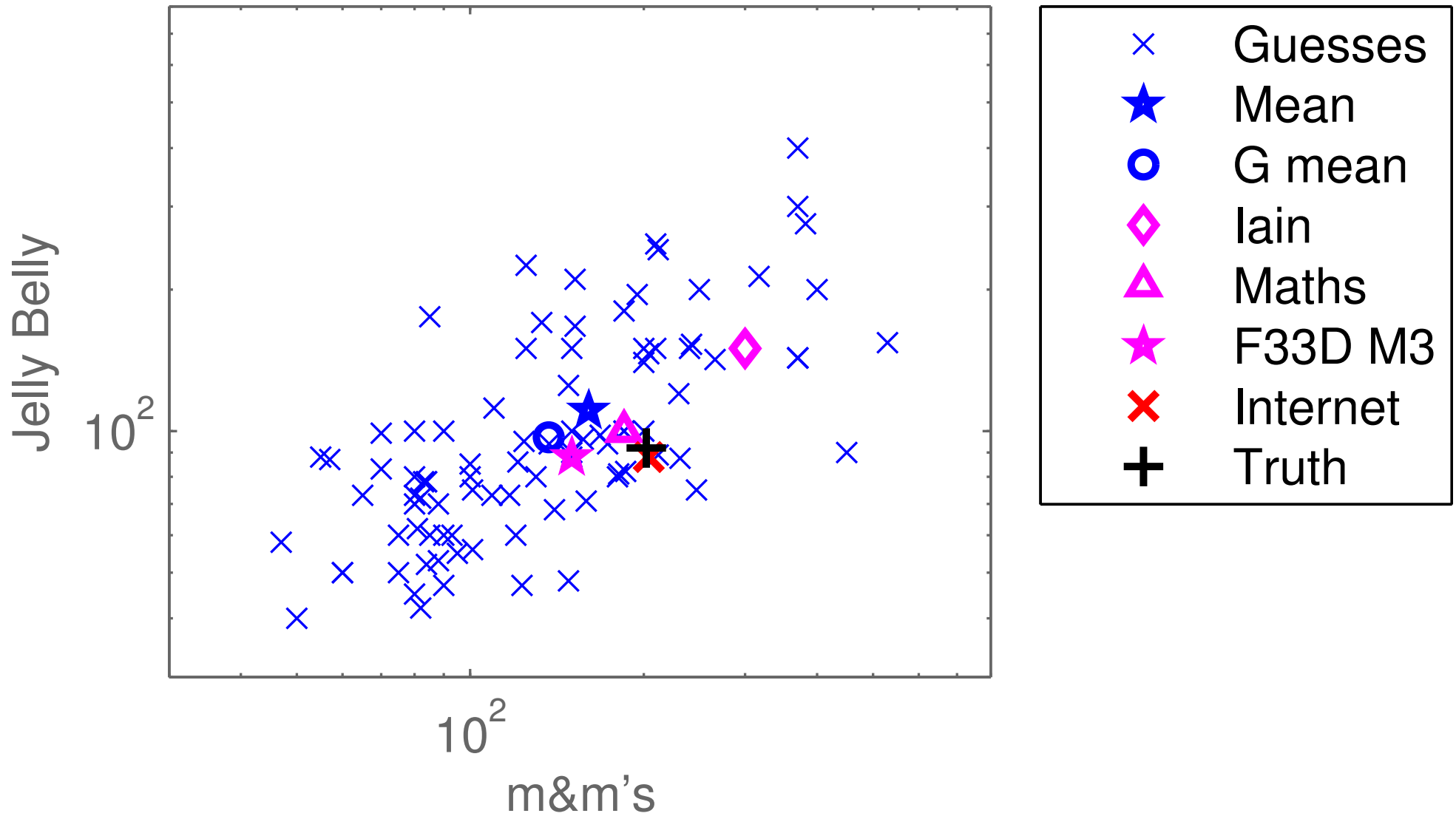
Often log-transform +ve data



Wisconsin breast cancer data

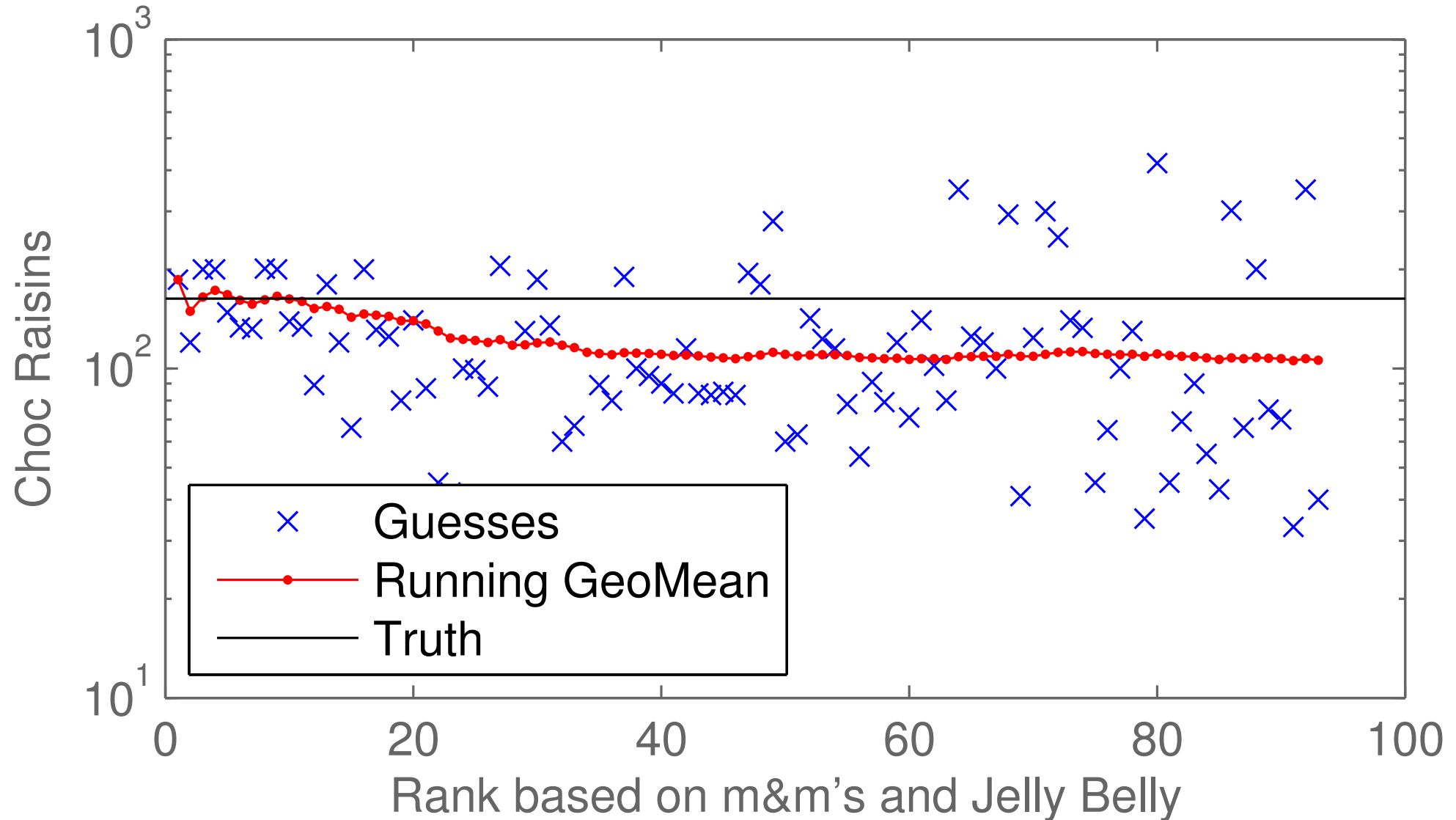
UCI ML repository

Count guesses on log-scale



Were some people just lucky?

Ranking by past performance



Ensemble of Models

Two motivations:

- ① Reduce over-fitting
- ② Reduce under-fitting

Example ①

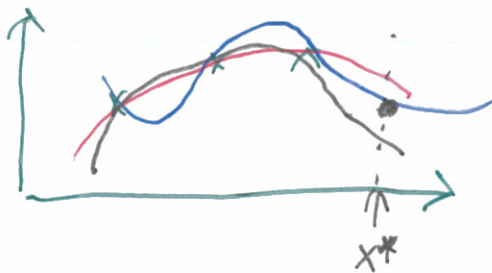
Bayesian model averaging:

$$p(y | \underline{x}, D) = \int p(y | \underline{x}, \underline{w}) p(\underline{w} | D) d\underline{w}$$

$$\approx \frac{1}{S} \sum_{s=1}^S p(y | \underline{x}, \underline{w}^{(s)})$$

$$\uparrow \quad \underline{w}^{(s)} \sim p(\underline{w} | D)$$

Ensemble of S predictors



Another similar ensemble Bagging

"Bootstrap aggregation"

N training examples

Training time:

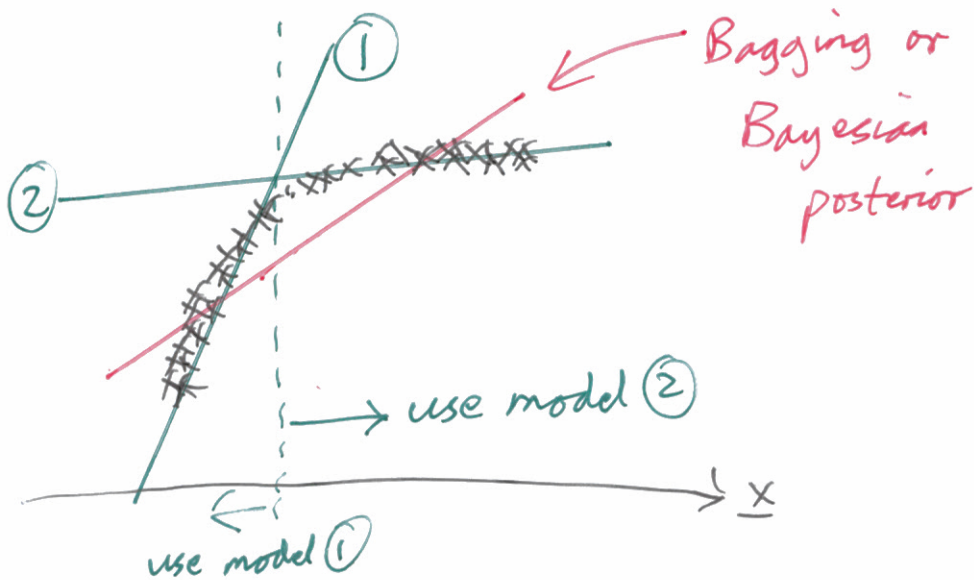
for $s = 1 \dots S$:

Bootstrap: create a new dataset,
sampling N datapoints from
training data with replacement

Fit model to dataset \rightarrow predictor s

Test time:

Average predictions
(or majority vote)



② Model combination

$$p(y | \underline{x}, \theta) = \sum_z \underbrace{p(y | \underline{x}, z, \theta)}_{\text{Any regression model}} \underbrace{p(z | \underline{x}, \theta)}_{\text{"Gating network"}}$$

"Mixture of experts"

choice expert

Any classifier

Fit θ , neg log likelihood,

Regularize fit, Bayesian, Bagging, ...

Another way to build complicated model
Boosting.